

博士論文

Study on Statistical Inference for Unknown Sources of
Atmospheric Pollutants in Urban Environment
(都市環境における未知空気汚染発生源の
確率的推定に関する研究)

賈 鴻源

Abstract of dissertation

This dissertation studies the statistical inference method for an unknown source of atmospheric pollutants in the complicated urban environment based on Bayesian inference. To enable the inference method to handle the characteristics of the urban area, the main works of the dissertation include: extending the inference's ability to estimate the geometry of the unknown source; increasing the estimation accuracy by introducing a sophisticated dispersion model into the inference; proposing a sensor configuration optimization method to improve the quality of measurements and guide the sensor network design.

Nowadays, a high level of urbanization results in a considerable amount of people gathering in the urban area, where a safe and healthy atmospheric environment has never been more important. However, unexpected atmospheric pollution emitted from unknown sources sometimes occurred because of terrorism, nuclear accidents, illegal industrial emission, and other emergencies, which is a serious threat to humankind and the earth's environment. Therefore, it is important to identify the unknown source as soon as possible after these dispersion emergencies happened. Until now, intensive research has proposed various methods to realize source term estimation (STE). The basic framework of STE is using the estimation algorithm to find the true source based on the measurements obtained during the emergency and the source-receptor relationship (S-RR), which is the concentration prediction modeled in advance. However, the existing methods are still not capable enough of handling complex scenarios in the urban environment because of its unique features as follows.

First of all, possible sources of pollutants are diverse in the urban environment. Although the previous research has considered mobile sources and multiple sources, their estimation algorithm always assumed that the unknown source is an ideal point without geometry. Actually, some sources have non-negligible shapes or volumes, which may fail the estimation based on the point assumption. It is important to extend the STE method for geometry estimation.

Besides, the dispersion mechanism of atmospheric pollutants is extremely complicated in the urban environment. The transportation of pollutants involves complex turbulence caused by buildings, equipment, moving objects, and heat discontinuity. However, the existing simulation techniques for S-RR are still limited in steady models and cannot accurately predict turbulent dispersion. Since STE relies on S-RR as the prior prediction of pollution dispersion, it is necessary to introduce a precise modeling technique for the S-RR in order to promise the accuracy of STE.

What's more, ideal concentration measurements are difficult to acquire in cities because of the dense building distributions and land property limitations. Sensor configurations used in real life are nearly random or empirical, but their effectiveness to all possible sources has no promise. Because the performance of STE is highly dependent on the quality of measurements, it is meaningful to develop a sensor configuration optimization method to guide the sensor deployment and provide high-quality measurements.

Based on this research background, this dissertation aims to develop a statistical inference method for STE in the urban environment by dealing with 3 unsolved problems mentioned above. Because the complexity of real atmospheric pollution is out of the range of a single dissertation, the research subject is limited to one single source with constant emission strength and fixed location in a statistically steady turbulent flow field in a neighborhood-scaled urban area. The main research contents are summarized below.

As a beginning, a new method was proposed to estimate the geometry of unknown sources based on the super-Gaussian function. The coefficients of this function can control its distribution to approximate common shapes: line, rectangular, and ellipse. These coefficients are added into Bayesian inference to realize the geometry estimation. The applicability of the proposed method was first confirmed using a numerical experiment of an ideal boundary layer. The method successfully inferred that the source is line-like without any prior knowledge. Based on this case, the effects of different sensor configurations on the line source estimation were discussed. Because the line source contained more geometric information than point sources, the conventional sensor configuration for the point source may fail in the line source estimation. It was found that the requirements on the sensor configuration become higher. Both the sensors near the source and null-measurement sensors are indispensable.

To examine the robustness of the proposed method against measurement and modeling errors, the second case of a simplified urban square with wind tunnel experiment measurements was conducted. The line source was successfully identified by the proposed method again. By comparing to the conventional STE method with ideal point assumption, it is confirmed that the proposed method can not only provide the geometry estimation but also reduce the inference errors caused by the point source assumption. Hence, it is important to include the geometry estimation when the geometry of the source has unignorable effects on the dispersion.

Then, to improve the accuracy of STE in complex urban applications, large-eddy simulation (LES) was introduced to model the S-RR by unsteady adjoint equations. The LES of adjoint equations has rarely been conducted in the literature because the adjoint equation describes an inverse dispersion process. The time-series flow field data of the entire domain must be produced by forward simulation and stored in advance, thus the volume of data simulated with LES is too large for practical applications. This research proposed to use the wavelet-based compression method to mitigate the storage pressure. The LES flow field can be compressed into a portable database to make the simulation of unsteady adjoint equations easier.

As the first step, to evaluate the accuracy of compression and usefulness of the compressed database, a turbulent flow field in a block-arrayed building group model was simulated by LES and compressed into a database by the wavelet-based compression method. The influence of compression on the quality of the data was checked from the perspectives of a single snapshot and time series. It was confirmed that about 100 times compression can still satisfy the requirement of flow field visualization and afterward simulation. Large-scaled turbulent structures were well preserved after compression, and the dispersion simulation can be reliably reproduced with compression data. Therefore, it is reasonable to expect that the unsteady simulation of adjoint equations can be realized based on the compression database.

Afterward, the compression database above was used in the LES of adjoint equations to model S-RR, which was combined with Bayesian inference as a new STE method. The concentration measurements obtained from a wind tunnel experiment were applied to testify the performance of the proposed method. As a comparison, another STE was also conducted with a conventional method, where steady adjoint equations were simulated with the Reynolds-averaged Navier-Stokes (RANS) model. The results showed that the modeling of S-RR and the

accuracy of STE were significantly improved by the LES of the adjoint equation. The complicated turbulent flows caused by buildings destroyed the reliability of the conventional RANS model. Although the proposed method needs more computational resources, to effectively perform STE in the complicated urban environment, it is valuable to apply LES for adjoint equation simulation.

At last, a sensor configuration optimization method for STE was proposed by the design of an objective function and application of the simulated annealing method. The objective function was set as the information entropy of the spatial distribution of the adjoint concentration field. Its ability to represent the measurement ability of sensor configurations was proved from the views of mathematics and physical meanings. Simulated annealing was applied to find the optimal configuration which owns the largest value of the objective function. The proposed method was utilized to design an optimal sensor configuration for the block-arrayed building group model. The performance of the optimal configuration in STE was compared to uniform and random configurations through estimations for 25 unknown sources. The results revealed that the accuracy of STE is related to the entropy contained in the adjoint concentration of the configuration such that the design of the objective function is reliable. The optimal configuration outperforms the other two in STEs. It is valuable to use the proposed method to guide the configuration design in real applications.

Table of contents

ABSTRACT OF DISSERTATION	I
ABBREVIATIONS.....	1
CHAPTER 1 INTRODUCTION	3
Abstract	4
1.1 Research Background.....	5
1.1.1 Atmospheric pollution in the urban environment	5
1.1.2 Unexpected dispersion of atmospheric pollutants	6
1.1.3 Source term estimation	7
1.1.4 Elements in source term estimation.....	7
1.1.5 Difficulty of source estimation in the urban environment.....	9
1.2 Research objective.....	10
1.3 Research problem.....	11
1.4 Structure of the thesis	12
CHAPTER 2 BASIC METHODOLOGY FOR SOURCE TERM ESTIMATION	14
Abstract	15
2.1 Problem definition.....	16
2.2 Measurement	17
2.3 Estimation algorithm.....	18
2.3.1 Optimization framework	18

2.3.2 Statistical framework.....	21
2.4 Source-receptor relationship	24
2.4.1 Gaussian puff model.....	25
2.4.2 Markov chain model.....	25
2.4.3 Adjoint equation method	26
2.5 Conclusion.....	31
Symbols.....	32
CHAPTER 3 LINE SOURCE ESTIMATION USING SUPER- GAUSSIAN FUNCTION	35
Abstract	36
3.1 Introduction	37
3.2 Line source model	37
3.2.1 Ordinary Gaussian function.....	38
3.2.2 Super-Gaussian function.....	40
3.3 Case I: numerical experiment of urban boundary layer	41
3.3.1 Numerical simulation	41
3.3.2 Bayesian inference settings	46
3.3.3 Estimation results	46
3.4 Discussion about sensor configuration.....	48
3.4.1 Uniform sensor configuration.....	49
3.4.2 Importance of null measurements.....	52
3.4.3 Necessary sensors for line source estimation	56
3.5 Case II: wind tunnel experiment of the urban square model.....	59
3.5.1 Wind tunnel experiment for measurements	60

3.5.2 Numerical simulation	61
3.5.3 Bayesian inference settings	63
3.5.4 Estimation results	63
3.5.5 Comparison to conventional method with point assumption	67
3.6 Conclusion.....	68
Symbols.....	71
CHAPTER 4 CONSTRUCTION OF URBAN TURBULENT FLOW DATABASE BY LARGE-EDDY SIMULATION AND WAVELET-BASED COMPRESSION METHOD	74
Abstract	75
4.1 Introduction	76
4.2 Compression methodology.....	79
4.2.1 Wavelet decomposition.....	80
4.2.2 Quantization	83
4.2.3 Entropy encoding.....	84
4.3 Case description	85
4.3.1 Production of raw data with large-eddy simulation.....	85
4.3.2 Method to verify the compression database	87
4.4 Results and discussion.....	89
4.4.1 Validation of large-eddy simulation.....	89
4.4.2 Compression ability.....	92
4.4.3 Compression effects on a single snapshot	95
4.4.4 Compression effects on time-series data	101
4.4.5 Dispersion simulation with compression database	112

4.5 Conclusions and discussions	115
4.5.1 Conclusions	115
4.5.2 Method limitation and future research.....	117
Symbols.....	119
CHAPTER 5 SOURCE TERM ESTIMATION WITH UNSTEADY ADJOINT EQUATIONS MODELED BY LARGE-EDDY SIMULATION AND COMPRESSION DATABASE.....	122
Abstract	123
5.1 Introduction	124
5.2 Simulations for adjoint equation	125
5.2.1 Reynolds averaged Navier-Stokes simulation of the adjoint equation	125
5.2.2 Large eddy simulation of the adjoint equation	126
5.3 Case description	127
5.3.1 Wind tunnel experiment	127
5.3.2 Simulation settings	128
5.3.3 Bayesian inference settings	129
5.4 Results and discussions	129
5.4.1 Flow fields of forward simulations.....	129
5.4.2 Comparison of adjoint concentration	131
5.4.3 Source term estimation results.....	134
5.5 Conclusions	138
Symbols.....	140
CHAPTER 6 SENSOR CONFIGURATION OPTIMIZATION FOR SOURCE TERM ESTIMATION BASED ON THE ENTROPY OF ADJOINT EQUATION.....	143

Abstract	144
6.1 Introduction	145
6.2 Entropy-based configuration optimization.....	146
6.2.1 Objective function	146
6.2.2 Simulated Annealing	152
6.2.3 Calculation of $H(R I)$	154
6.2.4 Limitation of the method	155
6.3 Case study	155
6.3.1 Numerical simulation	156
6.3.2 Unknown sources	156
6.3.3 Sensor candidates	157
6.3.4 Comparison configurations	161
6.4 Results and Discussion.....	161
6.4.1 Optimum sensor configurations	161
6.4.2 Estimation results of one source.....	164
6.4.3 Estimation results of all sources	167
6.5 Conclusions	169
Symbols.....	171
CHAPTER 7 CONCLUSION AND FUTURE RESEARCH ...	174
7.1 Conclusions	175
7.2 Future research	177
REFERENCE	180
APPENDIX.....	192
Appendix A. Numerical simulation for turbulence and dispersion of pollution	

.....	193
A.1 Fluid and turbulence	193
A.2 Governing equations for fluid and dispersion	194
A.3 Large eddy simulation	196
A.4 Reynolds averaged Navier-Stokes model	201
A.5 Comparisons of two turbulence model	204
Appendix B. Publications related to dissertation	206
B.1 Journal publications	206
B.2 Conference publications	206
ACKNOWLEDGMENT	208

Abbreviations

CDF9/7	: Cohen–Daubechies–Feauveau 9/7
CFD	: computational fluid dynamics
DNS	: direct numerical simulation
GDH	: gradient diffusion hypothesis
LES	: large-eddy simulation
MCMC	: Monte-Carlo Markov chain
MHMC	: Metropolis-Hastings within Gibbs algorithm
PDF	: probability density function
PSD	: power spectrum density
RANS	: Reynolds-averaged Navier-Stokes
RHS	: right hand side
SA	: simulated annealing
SCO	: sensor configuration optimization
SGF	: super-Gaussian function
SGS	: sub-grid scale

S-RR : source-receptor relationship

STE : source term estimation

UAV : unmanned aerial vehicle

WCM : wavelet-based compression method

WD : wavelet decomposition

WTE : wind tunnel experiment

Chapter 1

Introduction

Abstract

This dissertation studies about statistical estimation method for an unknown source of atmospheric pollutants in the complex urban environment based on Bayesian inference. This chapter is an introduction of the dissertation including the research background, objective, and its structure.

1.1 Research Background

1.1.1 Atmospheric pollution in the urban environment

With the quick development of urbanization, the population is shifting from rural to urban areas. Until 2020, more than half of people in the world, over 4 billion, are living in urban areas (Fig. 1.1). Under this circumstance, it is important to keep a healthy atmospheric environment for such a massive population in urban areas. Meanwhile, the situation that massive population gathers in limited urban lands with extremely complicated infrastructures just causes frequent air pollution.

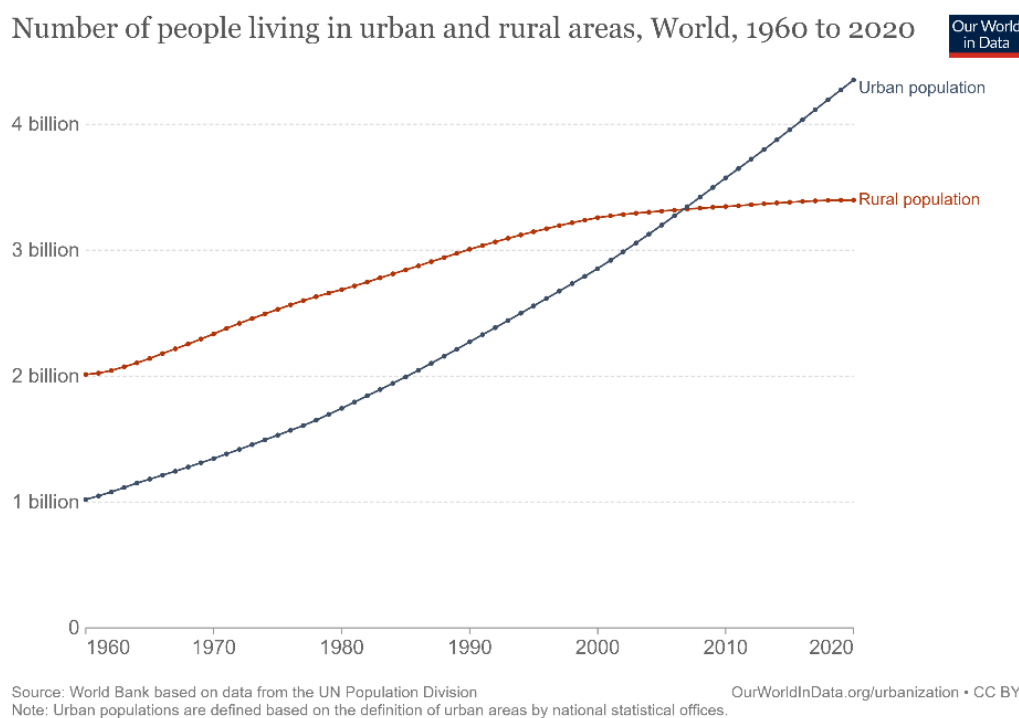


Figure 1.1 The trend of the number of people living in urban and rural areas from 1960 to 2020. (Data source: <https://ourworldindata.org/urbanization>)

Some air pollutants are well-known by people due to media promotion. For example, carbon dioxide is a famous greenhouse pollutant for its leading effects on global warming. Sulfur oxides produced by various industrial processes often cause consideration environmental damages like acid rain. With the popularity of auto-mobile, nitrogen oxides get the public's attention and become one of the main factors to the atmospheric pollution. The sources of these pollutants are known but are numerous or even infinite. They exist all around people in daily life. In recent years, with international consensus, people are dealing with these pollutants by

technological innovation, legal restriction, and pollution treatment. It has been reported that most pollutions are improving due to these actions.

In contrast, another kind of pollution has different characteristics: (1) the source is single or finite but unknown; (2) the source does not commonly exist in the atmosphere but suddenly appears due to accidents and emergencies. This pollution is caused by the unexpected dispersion of atmospheric pollutants. Although the occurrence is relatively low, people still need a comprehensive countermeasure system to handle it owing to its high level of danger.

1.1.2 Unexpected dispersion of atmospheric pollutants

Unexpected dispersion of atmospheric pollutants jeopardizes the safety of humankind and the earth's environment. These kinds of emergencies can be caused by terrorism. On 20 March 1995, the most severe terrorism attack in Japan, the Tokyo subway sarin incident happened in the central districts of the city. The toxic gas killed 14 people and injured over 6000 people to varying degrees. Nuclear accidents also caused harmful unexpected dispersion. From 11 March 2011 to 15 March 2011, the explosion of 4 reaction units in the Fukushima Daiichi Nuclear Power Plant released a large amount of radioactive contamination into the atmosphere and ocean. The harm caused by this dispersion accident is immeasurable and is still ongoing. Besides, illegal industrial emission also needs special attention. In the spring of 2019, Seoul suffered from its worst-ever atmospheric pollution that the concentration of PM_{2.5} reached almost 4 times the South Korean standard. The reason was found to be the illegal industrial emission of LG Chem and Hanwha Chemical, which has continued for over 4 years until the investigation. Apart from these remarkable emergencies, relatively small-scaled accidents like gas leaks and fire in the common life may also result in severe damage and require attention.

One of the most dangerous properties of unexpected dispersion is that the source of pollutants is unknown at the beginning. People cannot evaluate the emergency or take effective countermeasures until they gain enough information about the source, which often costs considerable time and allows the exacerbation of the emergency. Therefore, it is important to identify the unknown source as soon as possible after the dispersion emergencies happened.

For the realization of this target, the available knowledge generally includes discrete concentration measurements provided by limited sensors, meteorological conditions like wind speed and direction, and geographic information of the incident site. The problem that

identifying the information of the unknown source by utilizing this information is called source term estimation (STE), which is an ill-posed inverse problem with high nonlinearity, strong dependence on input data, and non-single solutions.

1.1.3 Source term estimation

In recent years, intensive research efforts have been devoted to proposing solutions for STE (Hutchinson et al., 2017). Existed methods can be divided into 2 directions by the type of sensors they rely on: mobile sensors and stationary sensors. The development of robotics and unmanned aerial vehicles (UAVs) attributes to the appearance of STE with mobile sensors. People can deploy several sensors after the emergency, and sensors can automatically track down the boundaries or sources of the contaminants. This concept was inspired by biological behaviors like swarms of beetles. The mobile sensors can communicate and cooperate to realize the route planning during STE (Zarzhitsky and Spears, 2005). However, the main problem of this direction is the immature techniques of robots and UAVs, which still have a long way to go before real applications. Besides, its applicability is also doubtful considering the strict regulation of UAVs in urban areas.

The other direction, stationary sensor, is more mature and practical. Many cities already deployed measurement stations at different positions as civil infrastructures. They can provide the basic measurement: the time-averaged concentration of a period. Some stations can also provide advanced information like the time series of concentration and eddy covariance (concentration flux). As a result, more STE methods are based on stationary sensors than mobile ones, and this dissertation also focuses on the STE with stationary sensors as below.

1.1.4 Elements in source term estimation

There are three basic elements in the STE: measurement, source-receptor relationship (S-RR), and estimation algorithm. Measurement is the foundation of STE and is provided by sensors. It dominantly affects the accuracy of STE and is dependent on the sensing strategies, equipment, and techniques, which means sensing in the real application is a complicated task and inevitably involved uncertainties. In the previous research, uniform and random sensors configuration are commonly used to evaluate the performance of STE methods.

S-RR describes the relationship between a source and a sensor. It predicts the concentration measurement of one sensor when a source with certain strength appears

somewhere before the emergency by computational methods. Popular methods include the Gaussian plume model, adjoint equation model, and Markov chain model. Because knowledge of sources before the emergency is unavailable, as a prior database, S-RR should be able to count in a large number of possible sources, which is a high requirement for the calculation cost of methods. Meanwhile, since accurate S-RR is critical to a successful STE, the prediction accuracy of computational methods should be promised. It should be able to capture the influence of meteorological conditions and terrain characteristics on the dispersion.

Gaussian plume model and Markov chain model satisfy the first requirement. By assuming an ideal dispersion in a homogeneous flow field, the Gaussian plume model predicts the concentration of a sensor with an algebraic equation, which is easy and fast to calculate. However, its reliability is vulnerable to real dispersion fields like buildings immersed into a flow. Markov chain model assumes that movements of pollutants among different zones follow a probability matrix, which can be calculated by computational fluid dynamics (CFD) simulation. Although CFD can bring some inhomogeneous information of flow field, to keep the calculation speed, the probability matrix is set as unchanged, which means the turbulent dispersion cannot be properly predicted. Therefore, neither the Gaussian plume model nor the Markov chain model can maintain the prediction accuracy when it comes to complex urban applications.

Considering the main research interest of the thesis is STE in the urban environment, the adjoint equation model is used here to calculate S-RR. This model indeed needs more calculation resources than the other two models, but its prediction accuracy is much higher because the dispersion process is described as partial differential equations, which can theoretically deal with turbulent dispersion. With the help of CFD techniques, adjoint equations can be simulated to predict S-RR for complicated urban flow fields.

The estimation algorithm finds the most credible source as the estimation result based on the measurements and S-RR. Because the estimation algorithm is executed right after the emergency, both the speed and accuracy are important. Existing methods can be classified into two frameworks: optimization and statistics.

The optimization framework applies the objective function to evaluate each candidate, find a candidate that minimizes the value of the objective function and regards it as a unique solution. Different designs of the objective function and plenty of algorithms for minimization

have been developed. Optimization methods are often less computationally expensive than the statistical framework, and their convergence speed is faster. Despite that, sometimes their performance is highly dependent on the initial guess, which may cause strong unsteadiness. More importantly, the estimation is limited to a single solution. In real applications, the inevitable measurements errors in the sensing and numerical errors of S-RR modeling degrade the reliability of the single solution. The optimization may be misled by errors, yield a wrong estimation and the truth is ignored.

In contrast, the statistical framework calculates the probability of each candidate to produce a probability distribution, and consequently, the noise of measurements and simulations can be reflected. In this case, more information like confidence intervals will be provided, and the truth is more likely to be noticed. Most of the statistical methods are based on Bayesian inference, where the posterior probability can be estimated based on the prior information and measurements. Because Bayesian inference can reflect the influence of noise on the estimation results, it is probable to apply it to the STE in the complex urban environment. Therefore, the Bayesian inference is selected as the basic estimation algorithm of STE in the thesis.

1.1.5 Difficulty of source estimation in the urban environment

Although plenty of STE methods, which combine different choices for the above three elements, have been proposed, it is still difficult to claim that the existed methods can handle complex reality problems in the urban environment. In fact, there is still large room for further development of STE because of the features of the urban environment as follows.

First, possible sources of pollutants are diverse in the urban environment. They can be a point or shaped, fixed or moving, single or multiple. In the beginning, the estimation algorithm can only deal with a point source at a fixed position. Afterward, advanced algorithms were proposed to manage moving sources and multiple sources. However, little research is concerned about the estimation of geometry information of the source, which is also valuable for risk management and evaluation of seriousness.

Furthermore, the complicated turbulence field and dispersion phenomenon in the urban environment brought a challenge to the accuracy of S-RR simulation. The turbulent complexity of the atmospheric boundary layer is increased further by buildings, equipment, transportation

systems, and heat discontinuity after it reached the urban area. Therefore, modeling the turbulence flow field and pollutants dispersion in the urban environment is troublesome. In recent years, with the development of computers and numerical simulation techniques, sophisticated CFD models like the large eddy simulation are proposed to reproduce unsteady turbulence fields. However, models for S-RR are still limited that the most accurate solution till now is the steady Reynolds-averaged Navier-Stokes, which may fail to accurately predict the dynamic properties of dispersion behavior. Thus, it is necessary to upgrade the simulation method of S-RR and improve the STE performance.

Last but not least, ideal concentration measurements are difficult to acquire in the urban environment. Even though the existed STE methods perform well with uniform sensor configurations in the previous research, it is impractical to deploy a uniform sensor network in the urban area due to the irregular building distributions and land property. Actually, measurement stations in reality are closer to a random configuration, which is a consequence caused by check and balance of deployment difficulty and sensing efficiency. However, considering that unknown sources may appear anywhere in the city, a random configuration may cause considerable errors in the measurements of certain sources. It is meaningful to develop a sensor configuration optimization method to guide the sensor deployment. Thus, no matter where the unknown source of emergencies appears in the target area, the optimum sensor configuration can provide high-quality concentration measurements that help STE effectively identify the source, and consequently is an important part of the atmospheric pollution detection system.

1.2 Research objective

According to the introduction above, in order to protect the atmospheric quality and strengthen the defense ability of urban areas against dispersion emergencies caused by unknown sources, it is necessary to improve the performance of the current STE for the urban environment. It should be noted that this target cannot be accomplished only by one research given the extreme complexity of the urban environment. As one step of this process, this research mainly tried to realize three terms as follows:

- a) Extend the estimation algorithm's ability to estimate the geometry information of common

sources in the urban area: point, line, rectangular, and ellipse.

- b) Apply large eddy simulation in the S-RR simulation and evaluate its effectiveness in STE.
- c) Develop a sensor configuration optimization method independent of the prior knowledge of unknown sources.

Concentrating on the statistical STE problem, the above three improvements are made based on three elements as shown in **Fig. 1.2**.

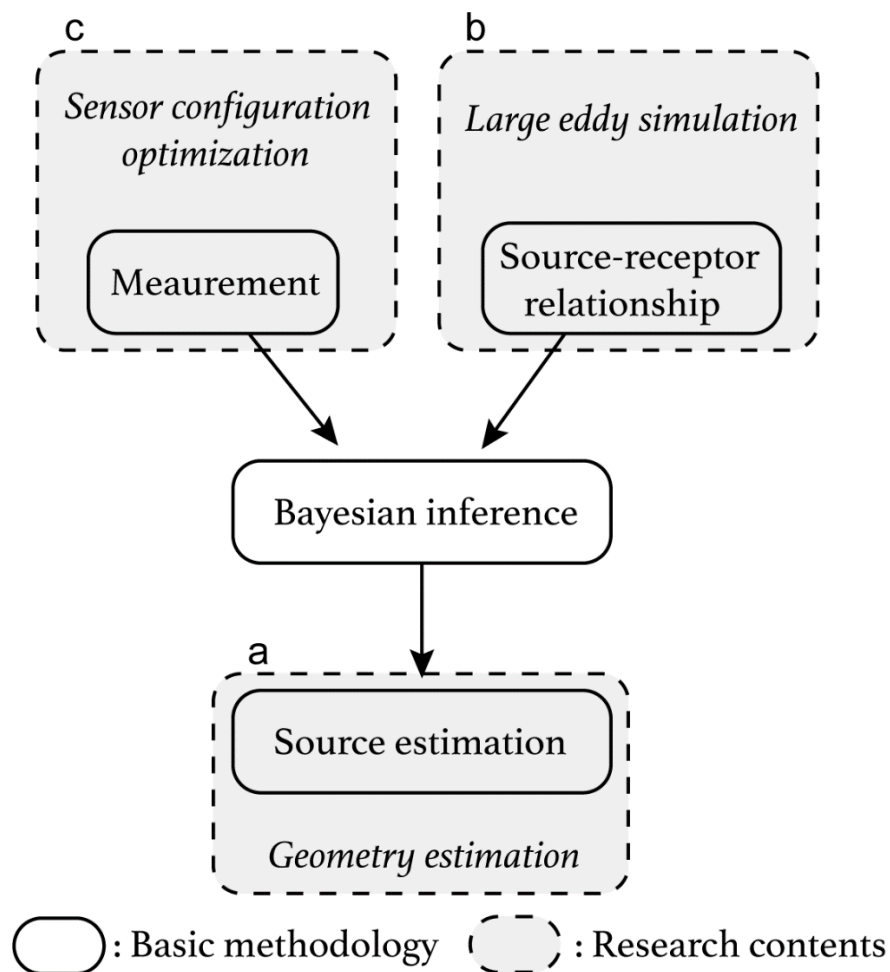


Figure 1.2. The relationship between basic statistical STE methodology and the research contents of the dissertation.

1.3 Research problem

The dispersion phenomena of pollutants are affected by many factors, including but not

limited to unsteady meteorological conditions, coupling of multiple sources, sources' movement, fluctuating emission strength, and chemical reactions. Thereby, a comprehensive STE for the urban environment needs the efforts of a wide range of research. Considering that it is impossible to solve all these issues by single research, to focus on the current target and keep the conciseness, the author needs to limit the range of this thesis by some basic assumptions. The research problem is assumed to be the identification of one single source with constant emission strength and fixed location in a statistically turbulent flow field in a neighborhood-scaled urban area. The emitted pollutants are regarded as passive scalars that will not react with other substances. These assumptions also have been widely used in the previous STE research.

1.4 Structure of the thesis

The thesis consists of 7 chapters as shown in **Fig. 1.3**. The content of each chapter is as follows:

- Chapter 1 introduces the research background, objective, and structure of the thesis.
- Chapter 2 gives a short review of previous STE research and explains the basic framework of the Bayesian inference STE method.
- Chapter 3-6 are the main content of this research. Chapter 3 proposes a new method to estimate the geometry information of the unknown source in STE.
- Chapters 4 & 5 apply large eddy simulation to model S-RR and compare its performance with conventional methods.
- Chapter 6 develops an entropy-based sensor configuration optimization method and evaluates its effectiveness in an ideal urban model.
- Chapter 7 is the concluding remark for this research. The future research plan based on the current progress is also discussed.

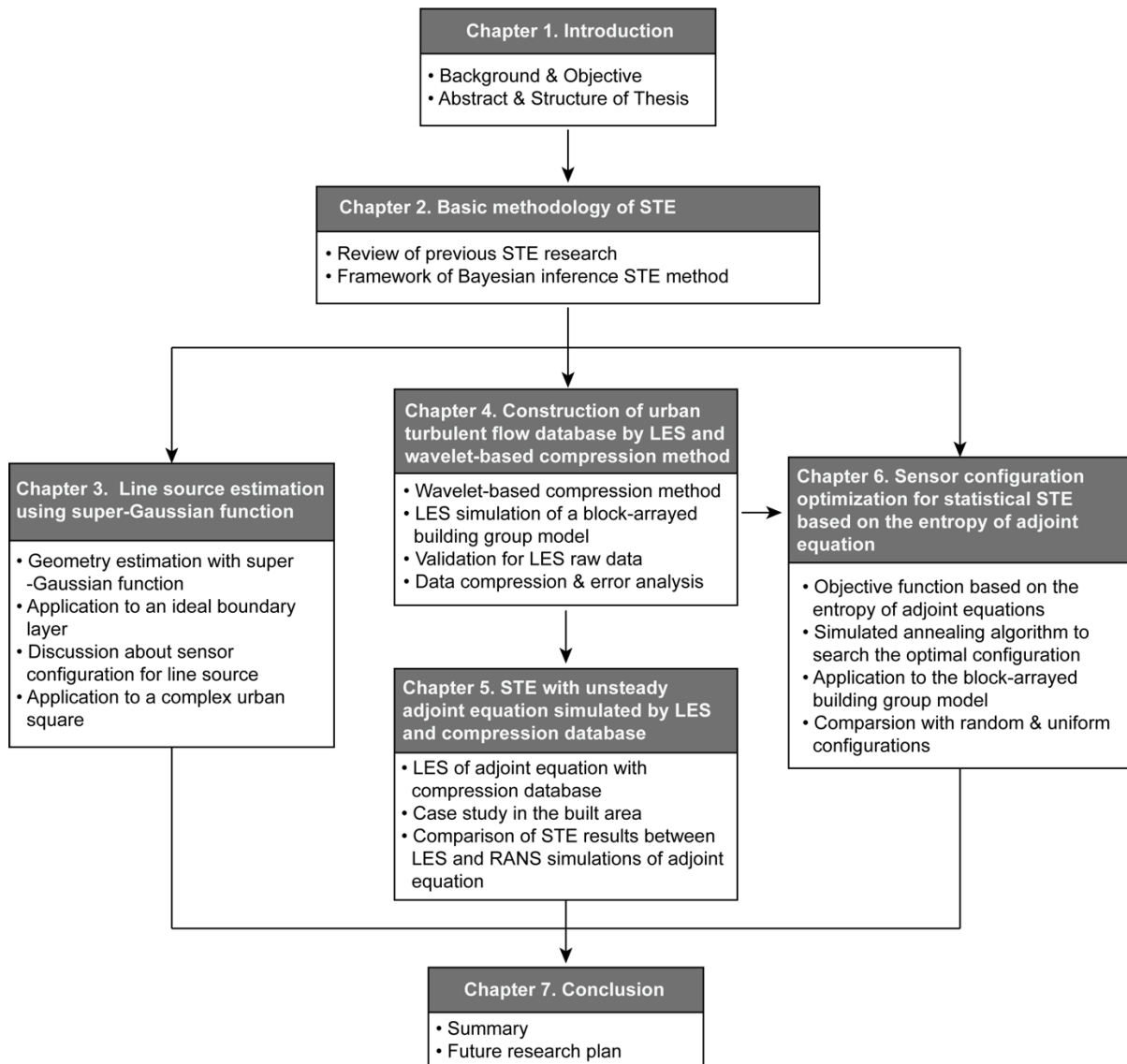


Figure 1.3. The structure of the dissertation.

Chapter 2

Basic methodology for source term estimation

Abstract

In this chapter, the basic methodology for source term estimation is introduced from the view of three elements: measurements, estimation algorithm, and source-receptor relationship. In each part, a short literature review about recent progress will be provided first. Then, for complicated urban applications, appropriate methods for each element are selected and described in detail.

2.1 Problem definition

The source term estimation (STE) problem is defined to identify unknown sources of pollutants in the urban environment. The focus is specified as the atmospheric pollutants here. Even so, the range of STE is still wide in the reality. It may involve different types of sources: single or multiple, fixed or mobile, point sources, or sources with geometry (volume). The emission patterns are also diverse. Some sources release the pollutants with constant strength, while some sources with varying strength. Besides, the properties of pollutants matter in STE. The inert gas can be regarded as the passive scalar which is easy to deal with, but some pollutants may be chemically active that they will react with other gases during dispersion. What's more, the pollution particles could experience buoyance, condensation, collision, and other phenomena during the dispersion. In addition, the meteorological condition is important to STE. In real applications, the incoming wind in the urban areas is unsteady. The wind speed and directions may fluctuate or even suddenly change during the dispersion emergency. It is inevitable to consider this possibility.

Therefore, STE is an inter-disciplinary research problem including numerous targets and complex dispersion situations. Each possible STE target mentioned above needs a series of research to solve, and even though intense research efforts have been devoted to STE, some issues remain unresolved. For instance, the forward modeling of pollution dispersion with chemical reaction is still in the early stage, let alone the STE for such pollution, which needs inverse simulation of such dispersion. Obviously, it is reasonable to infer that plenty of research is still necessary to comprehensively realize general STE in the future.

Since it is almost impossible to cover the total range of STE in a single research, and the current thesis is part of a series of studies to realize STE in the urban environment, it is necessary to limit the range of STE here by some assumption in order to focus on the main research target.

Considering that STE research is still rudimentary in the urban environment, estimations for multiple sources and mobile sources are too early. The STE target is limited to a single source \mathbf{s} located at a fixed position \mathbf{x}_s . Here, \mathbf{s} is a vector describing the information of the unknown source. \mathbf{x}_s is the time-independent coordinate of the source in the domain. Because the source can have a geometry to estimate (**Chapter 3**), \mathbf{x}_s includes all the coordinates that the source has.

As for the emission state, it is assumed that the pollutants are constantly released from the source with unchanged strength q_s . The temporal fluctuations of emission strength will not be considered here. After the emission, the pollutant is regarded as the passive scalar in the dispersion modeling. As mentioned before, the precise modeling techniques incorporating chemical reactions and other physical mechanisms are immature, which hinders STE for such pollution.

The meteorological condition also needs some limitations. Since the wind flow in the real-life cannot get rid of turbulence, turbulent inflow is taken into consideration. Meanwhile, it is expected that STE should be realized in a short time after the dispersion emergency happened. During this short time, it is appropriate to assume that the meteorological condition is statistically steady, which means that the wind direction will not change and the incoming turbulence owns a constant statistical property that can be described by several profiles.

According to these assumptions, the target is limited to a single source located at a fixed position in a statistically stationary flow field with a constant release strength, which can be described by a vector \mathbf{s} :

$$\mathbf{s} = (\mathbf{x}_s, q_s) \quad (2.1)$$

2.2 Measurement

After the pollutants were emitted from the source \mathbf{s} , the sensor network can monitor the information about the concentration and provide the measurement data. One of the most basic and common data types is the mean concentration over a certain period, which can be provided by low-cost sensors. In real applications, the time scale of averaging operation is influential to the measurements. At the beginning of the dispersion, because the spatial concentration distribution is developing, mean concentration may strongly fluctuate with the averaging time scale and start point. To avoid this, it is effective to start the averaging after the dispersion reached to statistically steady state and make sure the time scale is long enough.

With the development of sensing techniques, more information can be measured by sensors like time-series data and eddy covariance. Time series data can reflect the temporal change of concentration caused by multiple sources, moving sources, or unsteady emission

strength. Therefore, it is valuable in these scenarios. For instance, Wang et al. (2021) proposed a method for multi-point source identification based on the correlation of time-series measurements. Apart from that, eddy covariance measurement is also becoming a monitoring exercise rather than a purely scientific activity (Aubinet et al., 2012). The concentration flux can reveal the coupling effects of source and flow field on the sensor, which is valuable complementary information for STE. The flux has been well used in the modeling of the sensor's footprint (Hellsten et al., 2015), whose meaning is similar to S-RR which represents the measurement ability of the sensor and relationships between possible sources and a certain sensor. Thus, it is reasonable to deduce that the accuracy of STE could be improved further by adding flux measurements. However, it is necessary to predict flux measurements in advance in the source-receptor relationship (S-RR), which needs complicated simulation models. It still needs further research in the future.

Because the analysis of advanced measurement data is out of the interest of the thesis, the measurement data here is the time-averaged concentration data, which can be obtained by most sensors. The measurements provided by n sensors can be represented by a n elements vector D .

2.3 Estimation algorithm

In the next step, it is necessary to estimate the source terms based on the measurements D . The estimation algorithms based on static sensors can be divided into two frameworks: optimization and statistics.

2.3.1 Optimization framework

The optimization framework relies on a kernel objective function, which is constructed based on the difference between D and modeled concentration R . The most common form is the residual sum square between D and R . The true source is believed to be the one that minimizes the objective function in the vector space of s . To effectively find the target, different types of methods are proposed for searching.

a. Gradient-based algorithm

The gradient-based methods find the optimal solution by decreasing the value of the

objective function. The descending direction is determined by the gradient of the objective function, which is approximated by the partial derivatives of the objective function over the source's parameters. Li and Niu (2005) used the gradient method to estimate the emission of volatile organic compounds in dry building materials. Sharan et al. (2012) applied the least square method to identify single and multiple point sources based on synthetic measurements and real measurements.

An extension of this method is called re-normalization (Issartel, 2003), which is a linear approximation for source terms based on measurements. It used the inverse linear relationship to estimate the source. To avoid artifacts generation in the inversion process, a renormalized weight function was introduced to limit the search space and can be updated in each iteration (Issartel, 2005a). The effectiveness of renormalization in STE was demonstrated in Kumar et al. (2015), where a point source in an urban-like environment was successfully reconstructed.

Overall, one drawback that needs to be noted in the gradient-based method is that the performance is vulnerable to the initial guess, and the optimization process may be stuck in a local minima without global searching.

b. Heuristics algorithm

Another searching strategy is the heuristics algorithms. Unlike gradient-based strategy, heuristics methods directly start searching without prior mathematical processing of objective function. They sacrifice precision or completeness for speed when there is no traditional way to efficiently find an accurate solution. The common heuristics searching methods in STE are pattern search, simulated annealing, and genetic algorithm.

Pattern search firstly defines the initial values for source parameters. Then, the algorithm changes each parameter by specified step length, which is called axis exploration, and the objective function value is calculated once again. The new value and original value are compared to decide whether a new estimation is accepted. If the objective value keeps unchanged in all directions, the step size is adjusted for new exploration, which is called pattern move. Zheng and Chen (2010) applied this method to locate a point source in a field experiment and found that it is limited to the local search, which needs to be improved.

Simulated annealing analogies the optimization process to the cooling process of a material, where the objective function is regarded as the thermodynamic energy of the system.

It also starts from an initial guess and generates new estimations by random disturbances. The special feature is that the cooling process analogy brings the concept of the temperature, which is used to adjust the acceptance probability of the new estimation. This feature enables the algorithm to jump out from the local minima and conduct a global search. More details can be found in **Chapter 6**. Thomson et al. (2007) used this method to locate an unknown point source in a desert environment.

Genetic algorithm is another popular global search method in engineering. At the beginning of the genetic algorithm, not a single guess but a random population of candidates are initialized. The most important feature is that the genetic algorithm focuses on the population rather than the parameter. It is inspired by the natural genetic evolution and finishes the optimization by selection, crossover, and mutation. Selection is to pick up the high-quality candidates to construct a gene pool by their values of the objective function. Crossover is to change part of the parameters of candidates in the pool to create the next generation. Mutation means that parameters may have random disturbances during the crossover to realize global search. The applications of the genetic algorithm in STE can be found in Allen et al. (2007) and Wang et al. (2018).

It is worth mentioning that the convergence of heuristic algorithms cannot be promised owing to their direct search property. It is difficult to judge whether the solution found by the heuristic is the global optimal or just a local peak. Many modifications have been proposed to improve its credibility like hybrid algorithms and combination with a good initial guess, but the effectiveness still needs more validations of real applications.

c. Limitation of optimization methods

Optimization methods frequently appear in the STE research, which demonstrates their effectiveness to the inverse problem. However, no matter in gradient-based algorithms or heuristic algorithms, only a single estimate, the final result, is provided to users. The complex, unavoidable noise in the real world cannot be separated from the result for further analysis. The substantial noise may make the optimization result totally deviate from the truth, thus users fail to identify the source.

2.3.2 Statistical framework

A better way to evaluate the noise in real applications is solving STE in a statistical framework, which estimates the probability of all candidates rather than a single result. The most popular statistical STE method is Bayesian inference. It treats each source term as an independent random variable. The posterior probabilities of these variables can be calculated by updating prior probabilities with sensor observations. Then stochastic sampling methods are used to attain the probability density function (PDF) of posterior probability for each parameter. As a result, users can obtain more information from PDF like the most probable estimate, the mean estimate, and credible interval when compared to optimization results. The uncertainty or noise during STE is reflected in PDFs. Keats et al. (2007) successfully identified a point source in a complex urban environment using Bayesian inference. Since their study, this method has been extended to estimate multiple point sources (Wade and Senocak, 2013a; Yee, 2012), mobile sources (Kopka et al., 2016).

Seeing that complex noise is inevitable in STE for the urban area, the Bayesian inference is selected as the main methodology in the thesis. Here a brief introduction is given.

a. Bayes' theorem and problem formulation

The definition of Bayes' theorem is:

$$p(A|B) = \frac{p(B|A)p(A)}{p(B)} \quad (2.2)$$

where A and B are events and $p(B) \neq 0$. $p(A|B)$ is a conditional probability called posterior probability, which is the probability of event A occurring given that B is true. $p(B|A)$ is called likelihood function, which has the similar probability meaning with $p(A|B)$. Following the Bayes' theorem, the posterior probability that \mathbf{s} is the true source given measurements \mathbf{D} can be written as:

$$p(\mathbf{s}|\mathbf{D}, I) = \frac{p(\mathbf{D}|\mathbf{s}, I)p(\mathbf{s}|I)}{p(\mathbf{D}|I)} \quad (2.3)$$

Here, I is the background information, including meteorological and geographical data. $p(\mathbf{D}|\mathbf{s}, I)$ is the likelihood function. $p(\mathbf{s}|I)$ is the prior probability. $p(\mathbf{D}|I)$ is called the evidence. According to the definition of conditional probability, it is the probability of \mathbf{D} can be measured given the background information. In this case, $p(\mathbf{D}|I)$ is independent of sources

and can be assumed as a constant when a certain I was imposed. It acts as a normalizing factor in Eq. (2.3), which does not affect the posterior probability distribution. Therefore, the posterior distribution can be estimated using a likelihood function and prior probability:

$$p(\mathbf{s}|\mathbf{D}, I) \propto p(\mathbf{D}|\mathbf{s}, I)p(\mathbf{s}|I) \quad (2.4)$$

b. Likelihood function

A likelihood function represents the probability that measurements, \mathbf{D} , can be obtained when the source is \mathbf{s} . This probability is evaluated based on the difference between the real measurements, \mathbf{D} , and the modeled concentration, \mathbf{R} , resulting from source \mathbf{s} . Considering the inevitable errors caused by measurement processes and modeling, the true concentration \mathbf{D}_{true} , \mathbf{D} , and \mathbf{R} exhibit the following relationship:

$$\begin{aligned} \mathbf{D} &= \mathbf{D}_{true} + \mathbf{e}^d \\ \mathbf{R} &= \mathbf{D}_{true} + \mathbf{e}^m \end{aligned} \quad (2.5)$$

Here, \mathbf{e}^d is the measurement error and \mathbf{e}^m is the modeling error. In real-world applications, these errors play an important role in STE, though they are difficult to quantify (Yee et al., 2014). In the previous research (Keats et al., 2007a; Xue et al., 2017, 2018a), it is common to assume that the errors of each sensor follow Gaussian distributions with means of zero and variances of $\sigma_{d,i}^2$ and $\sigma_{m,i}^2$, where i is the index for the sensors. In this case, the likelihood can be calculated as:

$$p(\mathbf{D}|\mathbf{D}_{true}, I) \propto \exp\left[-\frac{1}{2}\sum\frac{(\mathbf{D} - \mathbf{D}_{true})^2}{\sigma_{d,i}^2}\right] \quad (2.6)$$

$$p(\mathbf{D}_{true}|\mathbf{s}, I) \propto \exp\left[-\frac{1}{2}\sum\frac{(\mathbf{D}_{true} - \mathbf{R}(\mathbf{s}))^2}{\sigma_{m,i}^2}\right] \quad (2.7)$$

$$\begin{aligned} p(\mathbf{D}|\mathbf{s}, I) &= \int p(\mathbf{D}|\mathbf{D}_{true}, I)p(\mathbf{D}_{true}|\mathbf{s}, I)d\mathbf{D}_{true} \\ &\propto \exp\left[-\frac{1}{2}\sum\frac{(\mathbf{D} - \mathbf{R}(\mathbf{s}))^2}{\sigma_{d,i}^2 + \sigma_{m,i}^2}\right] \end{aligned} \quad (2.8)$$

c. Prior probability

A prior probability $p(\mathbf{s}|I)$ sets the probability distribution of source parameters before

the inference. The prior knowledge of an unknown source can be contained to make the inference faster or more certain. However, it is usually difficult to acquire credible information just after the emergency. It is usually assumed that no information except for the measurements was available before the inference. Therefore, the source can appear anywhere within the target domain and with any strength. Additionally, the source parameters are independent of each other, meaning that a uniform distribution was imposed on the prior probability.

$$p(\mathbf{s}|I) = \text{constant} \quad (2.9)$$

d. Posterior probability

The posterior probability can be obtained by combining Eq. (2.4), (2.8), and (2.9), as follows:

$$p(\mathbf{s}|D, I) \propto p(\mathbf{D}|\mathbf{s}, I)p(\mathbf{s}|I) \propto \exp \left[-\frac{1}{2} \sum \frac{(\mathbf{D} - \mathbf{R}(\mathbf{s}))^2}{\sigma_{d,i}^2 + \sigma_{m,i}^2} \right] \quad (2.10)$$

The only unknown information in Eq. (2.10) is the modeling concentration, \mathbf{R} , of source \mathbf{s} .

e. Sampling method

Although the target distribution can be explicitly formulated as Eq. (2.10), it is still difficult to determine its shape because the PDF lies in a multidimensional space composed of the elements of \mathbf{s} . Many solutions have been proposed to sample posterior distributions in Bayesian inference, among which the most commonly employed is the Monte-Carlo Markov chain (MCMC). However, the conventional MCMC often suffers a large rejection ratio in the high dimensional sampling, which diminishes the efficiency. In recent years, a hybrid MCMC method, the Metropolis–Hastings–within–Gibbs algorithm (MHMC) (Gilks et al., 1995; Hastings, 1970), was widely applied during the sampling process. MHMC introduces a statistical acceptance ratio adjustment based on the current sampling to mitigate this problem. This algorithm includes the following steps:

Step 1 : Propose an initial guess of the source: \mathbf{m}_1

For $i=1:n$

Step 2 : Generate a new estimate $\tilde{\mathbf{m}}$ by a proposal distribution $q(\cdot)$: $\tilde{\mathbf{m}} \sim q(\tilde{\mathbf{m}}|\mathbf{m}_i)$

Step 3 : Evaluate the MH acceptance probability:

$$\text{Set } \alpha = \min \left[1, \frac{P(\tilde{\mathbf{m}}|D, I)q(\tilde{\mathbf{m}}|\mathbf{m}_i)}{P(\mathbf{m}_i|D, I)q(\tilde{\mathbf{m}}|\mathbf{m}_i)} \right]$$

Step 4 : Update the Markov Chain according to acceptance or rejection:

$$\mathbf{m}_{i+1} = \begin{cases} \tilde{\mathbf{m}} & \text{if } \alpha \geq N[0,1]; \\ \mathbf{m}_i & \text{otherwise;} \end{cases}$$

End For

Here $N[0,1]$ represents the uniform distribution bound between 0 and 1, and n is the sampling number.

On the other hand, some methods also have been proposed to accelerate MCMC with more complicated techniques, such as the application of variational inference (de Freitas et al., 2013), resampling (Gelfand and Sahu, 1994), and subspace construction (Constantine et al., 2016). However, their efficiency and robustness still need testification by practical applications. Heavy tails and local peaks may appear in the complex sampling process (Vrugt et al., 2009). Therefore, this thesis employed MHMC for the balance between speed and steadiness. The creditability of this method has been proved in the previous research of STE (Keats et al., 2007a; Xue et al., 2018a).

For each case in the dissertation, the Markov chain starts from an initial guess that is far enough from the true source. The total sampling number is set to 5.0×10^6 , in which the first 5.0×10^5 samplings are discarded for the burn-in process. The rest 4.5×10^6 samplings are used to estimate the PDF of the posterior probability.

2.4 Source-receptor relationship

As shown in Eq. (2.10), the modeling concentration, \mathbf{R} , of each sensor corresponding to each possible source, \mathbf{s} , is needed to evaluate the posterior probability distribution. Apparently, it is necessary to calculate \mathbf{R} for all possible sources and construct a database, which is exactly the S-RR. Because the number of possible combinations of parameters in vector \mathbf{s} is so large, special simulation techniques including the Gaussian puff model, Markov chain method, and adjoint equation method have been proposed.

2.4.1 Gaussian puff model

One popular forward calculation method is the Gaussian puff model. The concentration distribution for the continuous release is formulated as Zheng and Chen (2010):

$$R(x, y, z) = \frac{q_s}{2\pi u \sigma_y \sigma_z} \exp\left(-\frac{y^2}{2\sigma_y^2}\right) \times \left\{ \exp\left(-\frac{(z - H_e)^2}{2\sigma_z^2}\right) + \exp\left(-\frac{(z + H_e)^2}{2\sigma_z^2}\right) \right\} \quad (2.11)$$

Here, $R(x, y, z)$ is the modeling concentration at the location (x, y, z) in the downwind direction from the source following Cartesian coordinate. u is the wind speed. σ_y and σ_z are the dispersion coefficients for each direction. H_e is the effective height of the source.

It can be seen that Eq. (2.11) is an algebraic model. Given the source terms and sensor's position, the modeling simulation can be quickly calculated. In this case, even with a large number of possible sources, S-RR still can be easily constructed. However, the simplicity also brings with an important flaw that it can only describe an ideal dispersion in a homogenous flow field. Its accuracy cannot be promised when an obstacle appears in the domain. Until now, the literature using this model (Cui et al., 2019; Kormi et al., 2018; Ma et al., 2017; Wade and Senocak, 2013b; Wang et al., 2018; Zheng and Chen, 2010) also focus on the homogenous dispersion in the open field. As a result, the applicability of the Gaussian model in the urban environment is limited.

2.4.2 Markov chain model

Markov chain model is another forward calculation technique that was first proposed to solve indoor contaminant transport by Nicas (2000). The basic idea is to divide the dispersion domain into n zones or cells, and the concentration distribution of these zones is called states related to time. After the source appears at a certain state, the pollutant of each zone will move to an adjacent zone or remain still in the next step. Then, the transport behavior can be driven by a transition probability matrix:

$$\mathbf{P} = (p_{i,j})_{(n \times n)} = \begin{bmatrix} p_{1,1} & p_{1,2} & \cdots & p_{1,n} \\ p_{2,1} & p_{2,2} & \cdots & p_{2,n} \\ \vdots & \vdots & \cdots & \vdots \\ p_{n,1} & p_{n,2} & \cdots & p_{n,n} \end{bmatrix} \quad (2.12)$$

In this matrix, $p_{i,i}$ is the probability that the pollutant will stay in the original cell, and $p_{i,j}$ represents the probability that pollutant will move from cell i to cell j in the next step. According to the research of Chen et al. (2015), the probability matrix can be calculated by the computational fluid dynamics (CFD) simulation of airflow. The concentration distribution at any time step is the result of iterating multiple of \mathbf{P} to the initial state.

There are two drawbacks in this model. To decrease the computational burden, \mathbf{P} is often kept constant or owns limited patterns. This treatment is acceptable when the flow field and the turbulent diffusion are steady. If strongly fluctuating turbulence exists in the target domain, the statistical properties of the flow field are difficult to be described only by several patterns. Furthermore, the number of divided cells or elements in \mathbf{P} cannot be too large, which means the resolution of simulation and STE is restricted. As a result, the best application scenario of this method is the place where the flow field is almost unchanged and the low-resolution is allowable, such as a small indoor space with a fixed ventilation pattern (Li et al., 2020). Its applicability in the large-scaled outdoor space is in doubt and related literature is sparse.

2.4.3 Adjoint equation method

It can be stated that the accuracy of the two methods above is not enough in the complex urban area despite their computational simplicity. Speaking of accuracy, with the development of CFD techniques and high-performance computers, the dispersion behavior of atmospheric pollutants can be well predicted using turbulence models and transport equations. The related contents about turbulence modeling and CFD simulation of dispersion of atmospheric pollution are introduced in detail in **Appendix A**. Some representative literature about applications in the urban environment can be referred to (Branford et al., 2011a; Coceal et al., 2014a; Tominaga and Stathopoulos, 2011a, 2013).

However, the forward CFD dispersion method needs massive computation cost and long calculation time, it is almost impossible to directly predict the dispersion of all possible sources.

To overcome this problem, Pudykiewicz (1998) introduced the adjoint equation method into S-RR simulation to decrease the extreme computational costs. This approach enables users to construct S-RR from the perspective of sensors using inverse simulations, rather than from possible sources via forward simulation. Since the number of sensors is finite, the use of adjoint equations makes the construction of S-RR much more convenient. A brief introduction of the adjoint equation is given below.

a. Adjoint equation theory

For any source, \mathbf{s} , that releases a passive scalar value in a spatiotemporal domain with a strength of q_s and location of \mathbf{x}_s , the conservation equation of passive scalar transport can be written as (also Eq. (A.3) in **Appendix. A**):

$$\frac{\partial C}{\partial t} + \frac{\partial u_j C}{\partial x_j} - \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C}{\partial x_j} \right) = q_s \delta(\mathbf{x} - \mathbf{x}_s) \quad (2.13)$$

with the boundary conditions:

$$\begin{aligned} \nabla_n C &= 0 \text{ at } \partial\Omega \\ C(\mathbf{x}, t = 0) &= 0 \end{aligned} \quad (2.14)$$

Here, C is the concentration function of the passive scalar in a spatial domain of Ω and time of $[0, \mathbb{T}]$, D_m is the mass diffusivity, $\delta(\cdot)$ is the Dirac delta function, and ∇_n is a directional derivative normal to the boundary. The left-hand side (LHS) of Eq. (2.13) can be regarded as the resultant transformation of linear operator $L(\cdot)$ on function C in a Hilbert space, where $L(\cdot)$ is defined as:

$$L(\cdot) \equiv \frac{\partial(\cdot)}{\partial t} + \frac{\partial(u_j \cdot)}{\partial x_j} - \frac{\partial}{\partial x_j} \left(D_m \frac{\partial(\cdot)}{\partial x_j} \right) \quad (2.15)$$

and Eq. (2.13) turns into:

$$L(C) = q_s \delta(\mathbf{x} - \mathbf{x}_s) \quad (2.16)$$

Then, we introduce the conjugate concentration field C^* and consider the following spatiotemporal integration.

$$\langle C^* \mathbf{L}(C) \rangle = \int_T \int_{\Omega} (C^* \cdot \mathbf{L}(C)) dxdt \quad (2.17)$$

For convenience, the spatiotemporal integration is expressed as $\langle \cdot \rangle$. If we expand the LHS of Eq. (2.17), we have

$$\begin{aligned} \langle C^* \mathbf{L}(C) \rangle &= \langle C^* \left\{ \frac{\partial(C)}{\partial t} + \frac{\partial(u_j C)}{\partial x_j} - \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C}{\partial x_j} \right) \right\} \rangle \\ &= \langle C^* \frac{\partial(C)}{\partial t} + C^* \frac{\partial(u_j C)}{\partial x_j} - C^* \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C}{\partial x_j} \right) \rangle \\ &= \left\langle \frac{\partial(C^* C)}{\partial t} - C \frac{\partial(C^*)}{\partial t} + \frac{\partial(CC^* u_j)}{\partial x_j} - C u_j \frac{\partial(C^*)}{\partial x_j} \right. \\ &\quad \left. - \frac{\partial}{\partial x_j} \left(C^* D_m \frac{\partial C}{\partial x_j} \right) + D_m \frac{\partial C}{\partial x_j} \frac{\partial C^*}{\partial x_j} \right\rangle \\ &= \left\langle \frac{\partial(C^* C)}{\partial t} - C \frac{\partial(C^*)}{\partial t} + \frac{\partial(CC^* u_j)}{\partial x_j} - C u_j \frac{\partial(C^*)}{\partial x_j} \right. \\ &\quad \left. - \frac{\partial}{\partial x_j} \left(C^* D_m \frac{\partial C}{\partial x_j} \right) + \frac{\partial}{\partial x_j} \left(C D_m \frac{\partial C^*}{\partial x_j} \right) \right. \\ &\quad \left. - C \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C^*}{\partial x_j} \right) \right\rangle \\ &= \left\langle \frac{\partial(C^* C)}{\partial t} \right\rangle \\ &\quad + \left\langle -C \frac{\partial(C^*)}{\partial t} - C u_j \frac{\partial(C^*)}{\partial x_j} - C \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C^*}{\partial x_j} \right) \right\rangle \\ &\quad + \left\langle \frac{\partial(CC^* u_j)}{\partial x_j} - \frac{\partial}{\partial x_j} \left(C^* D_m \frac{\partial C}{\partial x_j} \right) + \frac{\partial}{\partial x_j} \left(C D_m \frac{\partial C^*}{\partial x_j} \right) \right\rangle \end{aligned} \quad (2.18)$$

If the integration happened at the time when C and C^* reach steady state or the integration was conducted following the periodicity of C and C^* , the first term is 0. According to the divergence theorem, the last term can be transferred into the surface integration of the space, which is called the boundary term. If the spatial integration is conducted in the space far away from the source, the boundary terms are almost 0. Therefore, Eq. (2.18) changes into

$$\langle C^* \mathbf{L}(C) \rangle = \left\langle -C \frac{\partial(C^*)}{\partial t} - C u_j \frac{\partial(C^*)}{\partial x_j} - C \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C^*}{\partial x_j} \right) \right\rangle = \langle \mathbf{L}^*(C^*) C \rangle \quad (2.19)$$

Here, \mathbf{L}^* is the adjoint operator of \mathbf{L} ,

$$\mathbf{L}^* \equiv -\frac{\partial(\cdot)}{\partial t} - u_j \frac{\partial(\cdot)}{\partial x_j} - \frac{\partial}{\partial x_j} \left(D_m \frac{\partial(\cdot)}{\partial x_j} \right) \quad (2.20)$$

Similarly, the transformation of $\mathbf{L}^*(\cdot)$ in C^* can be regarded as the governing equation of an adjoint tracer dispersion equation, which can be expressed as:

$$\mathbf{L}^*(C^*) = q_m \delta(\mathbf{x} - \mathbf{x}_m) \quad (2.21)$$

with the boundary conditions:

$$\begin{aligned} \nabla_{\mathbf{n}} C^* &= 0 \text{ at } \partial\Omega \\ C^*(\mathbf{x}, t = \mathbb{T}) &= 0 \end{aligned} \quad (2.22)$$

In contrast to the dispersion behavior of source \mathbf{s} , where the passive scalar is transported by $\mathbf{u}(\mathbf{x}, t)$ from time 0 to time \mathbb{T} , the physical meaning of Eq. (2.21) is the dispersion transported by $-\mathbf{u}(\mathbf{x}, t)$ from time \mathbb{T} to time 0, where the source is located at \mathbf{x}_m and has the strength of q_m . The importance of the adjoint equation can be revealed after Eq. (2.16) and Eq. (2.21) are substituted into Eq. (2.19):

$$q_s \langle C^*(\mathbf{x}_s) \rangle = q_m \langle C(\mathbf{x}_m) \rangle \quad (2.23)$$

If we impose the location of a sensor to \mathbf{x}_m and unit to q_m , the right-hand side (RHS) of Eq. (2.23) is the modeling concentration \mathbf{R} of the sensor, which is equal to the LHS, and the time-averaged concentration of tracers at source \mathbf{s} emitted from the sensor. Therefore, in this method, the number of dispersion equations that must be solved declines remarkably from the number of possible sources, \mathbf{s} , to the number of sensors applied. Owing to the convenience of this calculation, adjoint equations have been applied to model S-RR in several previous studies (Efthimiou et al., 2018a; Kumar et al., 2015a; Rajaona et al., 2015; Xue et al., 2017, 2018a) for outdoor STE.

b. Simulation of the adjoint equation

According to Eq. (2.23), the modeling concentration \mathbf{R} can be obtained after the adjoint concentration field C^* was calculated. Because Eq. (2.21) is in the form of a partial differential equation like the transport equation of Eq. (2.13), it can be solved by numerical simulation. In the previous research, it was commonly simulated by Reynolds averaged Navier-Stokes model based on the Reynolds-averaged form.

$$-\bar{\mathbf{u}} \frac{\partial(\overline{C^*})}{\partial \mathbf{x}} - \frac{\partial}{\partial \mathbf{x}} \left([D_t + D_m] \frac{\partial(\overline{C^*})}{\partial \mathbf{x}} \right) = q_m \delta(\mathbf{X} - \mathbf{X}_m) \quad (2.24)$$

where $\bar{\mathbf{u}}$ is the mean velocity field obtained from forward simulation and $\overline{C^*}$ is the simulated mean distribution of the adjoint concentration. D_t is the turbulent diffusivity used to model the turbulent scalar fluxes, which is an extra term produced by Reynolds averaging. The molecular diffusion coefficient D_m was set as $1.5 \times 10^{-5} \text{m}^2/\text{s}$ in this thesis, which corresponds to the diffusion of C_2H_4 , a common gas used in the dispersion experiment. It is true that different gases have different D_m values, but most of them lies in the range of 10^{-6} to $10^{-5} \text{m}^2/\text{s}$ (Wilke and Lee, 1955). When compared with the D_t , this small order difference would not cause too much effects on the dispersion behaviors.

Meanwhile, the wind speed determines the Reynolds number of the flow field in the domain. For the flow field with a high Reynolds number, which is common in the atmospheric boundary layer, the turbulence is well-developed, and its effects on the dispersion are large enough to make the molecular diffusion neglectable. Therefore, in the microscale or mesoscale of the atmospheric environment, where the ratio between the wind speed and diffusion coefficient is relatively large, the dispersion is dominated by turbulence rather than molecular diffusion.

Furthermore, for a certain dangerous pollutant with a special diffusion coefficient that needs attention, the adjoint equation can be specially simulated with this coefficient to build the S-RR, which would be applied later in the Bayesian inference when sensors monitor this pollutant.

Before the simulation of Eq. (2.24), a forward numerical simulation is necessary to obtain the mean velocity field $\bar{\mathbf{u}}$ of the target domain. One of the important factors in this forward simulation is the inflow boundary condition for $\bar{\mathbf{u}}$, which should correctly reflect the wind

characteristics of the dispersion emergency. It was assumed that the coming wind is statistically stationary during the dispersion process of pollutants. In fact, because adjoint concentration databases can be constructed long before emergencies, there is enough time to simulate them with different wind speeds and directions to cover the possible meteorological situations of the target domain. When the dispersion emergencies happen, as important meteorological information for STE, the properties of wind can be measured and included in the prior information. Then, the adjoint concentration corresponding to the wind properties would be selected for the Bayesian inference. In this case, because the adjoint equation was simulated under the same inflow as reality, the wind effects on the source-receptor relationship can be correctly built in the proposed method.

Because the simulation of the adjoint equation requires a preparatory forward simulation for the flow field information, theoretically, its calculation cost is a little bit higher than the Gaussian puff model and Markov chain model. Despite that, the complex turbulence dispersion is included in the adjoint equation to realize much higher accuracy, it is the most promising method for S-RR simulation in the outdoor STE among these three and consequently is applied in this thesis.

2.5 Conclusion

In this chapter, three core parts of STE: measurements, estimation algorithm, and S-RR were introduced. The commonly used methods in three parts were reviewed first. Then, considering the characteristics of STE in the urban environment, suitable methods for the thesis were selected to construct the basic structure of STE. The measurements are time-averaged concentrations measured by discrete sensors. These measurements are inputted into Bayesian inference to estimate posterior probability distributions of source parameters based on the comparison with modeled concentrations. The adjoint equations are used to predict the modeled concentrations before emergencies. In the following chapters, all the improvements and STEs will be conducted based on this basic structure.

Symbols

- C : the concentration distribution caused by a source
- C^* : adjoint concentration distribution
- \mathbf{D} : measurements vector
- D_m : the mass diffusivity of the pollutant ($1.5 \times 10^{-5} m^2/s$)
- D_t : turbulent diffusivity
- D_{true} : the true concentration at the sensor
- e^d : measurement error of the sensor
- e^m : modeling errors in the modeled concentration
- H_e : effective height of the source
- I : background information for Bayesian inference
- $L(\cdot)$: linear operator for the conservation equation of passive scalar transport
- $L^*(\cdot)$: the adjoint operator of $L(\cdot)$
- \mathbf{m}_i : the i th sampling in the MHMC
- $\min[a, b]$: the minimum value between a and b
- $N[a, b]$: the uniform distribution bound between a and b
- $p(A)$: the probability of event A
- $p(A|B)$: conditional probability of event A occurring given that B is true
- \mathbf{P} : transition probability matrix in the Markov chain dispersion model

-
- $p_{i,j}$: the probability that pollutant will move from cell i to cell j in the next step
- q_m : release strength of the sensor in the adjoint equation
- q_s : release strength of the source
- \mathbf{R} : modeled concentration vector
- $R(x, y, z)$: the modeling concentration at the location (x, y, z) in the downwind direction from the source following Cartesian coordinate in the Gaussian puff model
- \mathbf{s} : a vector representing the unknown source
- t : time from 0 to \mathbb{T}
- \mathbf{u} : wind velocity
- u_j : wind velocity in j direction
- \mathbf{x}_m : coordinates of the sensor
- x_j : coordinate in j direction
- \mathbf{x}_s : coordinates of the source
- α : acceptance probability for a new sampling in MHMC
- $\delta(\cdot)$: Dirac delta function
- $\sigma_{d,i}^2$: the variance of error in the measurement of the sensor with index i
- $\sigma_{m,i}^2$: the variance of error in the modeling concentration for the sensor with index i

σ_y : dispersion coefficients for y direction in the Gaussian puff model

σ_z : dispersion coefficients for z direction in the Gaussian puff model

∇_n : a directional derivative normal to the boundary

Ω : a spatial domain for dispersion and adjoint dispersion

$\langle \cdot \rangle$: spatiotemporal integration for domain Ω and time t

$-$: Reynolds average operator

1

2

3

4

5

6

7

8

9 Chapter 3

10 Line source estimation using

11 super-Gaussian function

12

13

1

2

3

4

5

6

Abstract

7

8

9

10

11

12

13

14

15

16

17

18

19

This chapter tries to extend the standard statistical source term estimation method to be capable of line source estimation, which includes more geometric information. Firstly, the proposed solution based on the super-Gaussian function is introduced. To justify its effectiveness, it is applied to the estimation of a line source in an ideal urban boundary layer with simulated measurements. Based on the estimation results, the effects of different sensor configurations on results and special requirements of line source on the configuration are discussed. After that, the performance of the proposed method is confirmed further by identifying a line source in a complex urban square with measurements in a wind tunnel experiment. By comparing with the conventional point source estimation method, the necessity and effectiveness of the proposed method are demonstrated.

3.1 Introduction

According to **Chapter 1&2**, most prior research aimed to estimate the parameters of a single point source, assuming that sources can be represented by a point without geometry; however, not all sources can be regarded as a point and most have shapes or volumes that cannot be neglected. The point assumption may result in noisy or even wrong estimation of these sources. Additionally, the source's geometric information is important for risk management and evaluation of seriousness. The recent public controversy concerning the Amazon forest fire (BBC, 2019) is a typical result of a lack of geometric information. Therefore, extending the source term estimation (STE) method for the estimation of sources with geometry is meaningful.

In addition, the line is a common and elementary geometry for pollution sources in the atmospheric environment (Gromke et al., 2008; Meroney et al., 1996; Salim et al., 2011a, 2011b). Traffic pollution in a street canyon has been extensively studied as a typical line source. Ideally, a method of source estimation handling any geometry should be possible; however, owing to its complexity and the numerous different kinds of geometry, this would require a series of research beyond the scope of one thesis. Correctly estimating the line source terms, which are not only position and strength but also length, width, and inclined angle is a good starting point for research.

The objective of this chapter is to extend the current STE method from point source estimation to line source estimation, which includes more geometric information. The basic structure is a combination of the super-Gaussian function for the source geometry and the Bayesian inference method. The method has the potential to estimate several common shapes: point, line, rectangular, and ellipse. The effectiveness of the method is evaluated by estimating a line source in two cases: an ideal urban boundary layer with simulated measurements and a complex urban square with discrete measurements, obtained via a wind tunnel experiment. The results of the new method are also compared with the conventional method's performance with a point assumption in the second case to demonstrate the necessity of the proposed method and geometry estimation.

3.2 Line source model

Eq. (2.23) has been widely used in previous research for point source estimation due to

1 convenience. During the Bayesian inference calculation, the modeled concentration \mathbf{R} for any
 2 possible point sources $\langle C^*(\mathbf{x}_s) \rangle$ can be picked up simply by their coordinates. However, the
 3 situation changes when the sources have various shapes or volumes instead of being a point.
 4 Considering the speed and efficiency of the sampling algorithm, the shape estimation of sources
 5 should be realized not by combining points randomly, but by some controllable functions that
 6 can be adjusted directly with several coefficients. For line source estimation, this chapter
 7 proposes that SGF is a good choice. This section introduces the super-Gaussian function as
 8 below.

9 3.2.1 Ordinary Gaussian function

10 A general format of the multivariate Gaussian (or normal) function is

$$f_{\mathbf{x}}(x_1, \dots, x_m) = \frac{1}{\sqrt{(2\pi)^m |\boldsymbol{\Sigma}|}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right] \quad (3.1)$$

11 where $\boldsymbol{\Sigma}$ is the covariance matrix of variable \mathbf{x} , and $\boldsymbol{\mu}$ is the mean value matrix of variable
 12 \mathbf{x} . When $m = 3$, this function corresponds to any ellipsoid or sphere in a given 3-dimensional
 13 space by changing the correlation coefficients in $\boldsymbol{\Sigma}$ and mean value in $\boldsymbol{\mu}$. Without loss of
 14 generality, the bivariate case can be transformed to a more intuitive format:

$$f(x, y) = A \exp[-(a(x - x_0)^2 + 2b(x - x_0)(y - y_0) + c(y - y_0)^2)] \quad (3.2)$$

15 where

$$\begin{aligned} a &= \frac{\cos^2 \theta}{2\sigma_X^2} + \frac{\sin^2 \theta}{2\sigma_Y^2} \\ b &= -\frac{\sin 2\theta}{4\sigma_X^2} + \frac{\sin 2\theta}{4\sigma_Y^2} \\ c &= \frac{\sin^2 \theta}{2\sigma_X^2} + \frac{\cos^2 \theta}{2\sigma_Y^2} \end{aligned} \quad (3.3)$$

16 A is a normalized factor. In this case, the Gaussian function can be adjusted directly through
 17 five coefficients: x_0 , y_0 , θ , σ_X , σ_Y , which determine the geometric shape independently. If
 18 we let the ratio σ_X / σ_Y become very large or very small, the shape is very close to a line as
 19 shown by **Fig. 3.1(a)**. Meanwhile, (x_0, y_0) is the coordinate of the middle point of the line, θ
 20 represents the inclined angle, $2\sigma_X$ and $2\sigma_Y$ are the width and length of the line.

1 More importantly, this method is not limited to the line source only. It also has the potential
 2 to estimate sources with shapes of ideal point, ellipse, and rectangular. During the STE process,
 3 the Bayesian inference will tune the coefficients automatically according to the measurement
 4 information to find the most appropriate shape: very small σ_X and σ_Y when the source is a
 5 point, a large ratio σ_X / σ_Y when the source is a line and normal σ_X, σ_Y when the source is
 6 ellipsoid or rectangular. In other words, it is not necessary to know the shape of the source
 7 beforehand, which is impossible in practical application. This model can estimate the
 8 approximate shape using measurement information only, as long as the source's shape is close
 9 to a sphere, ellipse, and line. If the source is an ideal point like the previous research, this
 10 method will regress into the conventional process automatically and yield the same result.

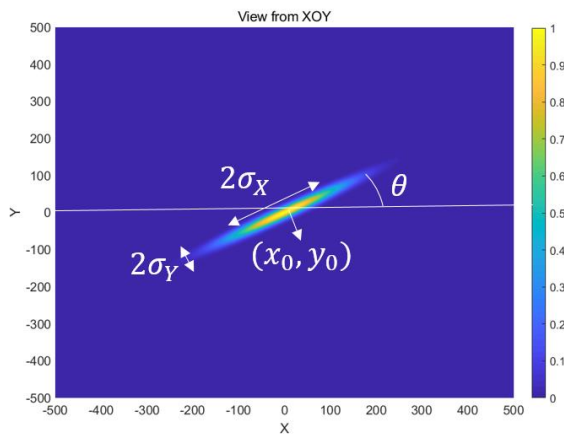
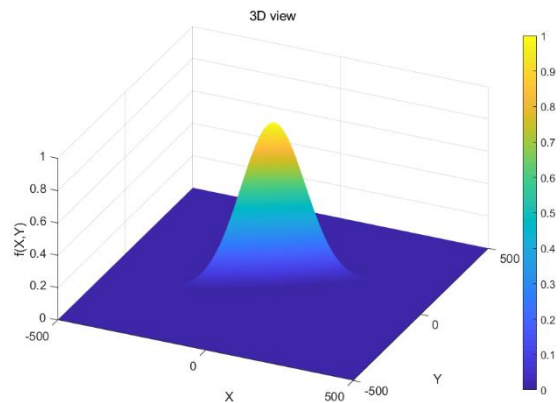
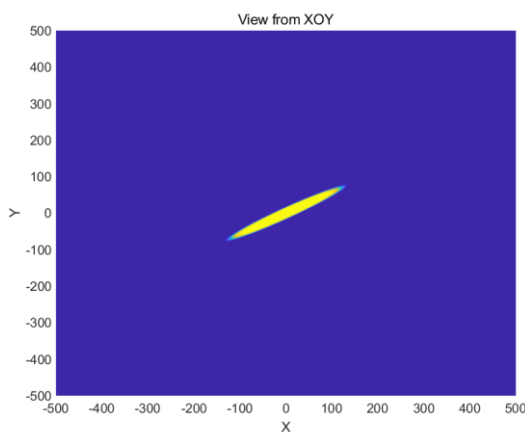
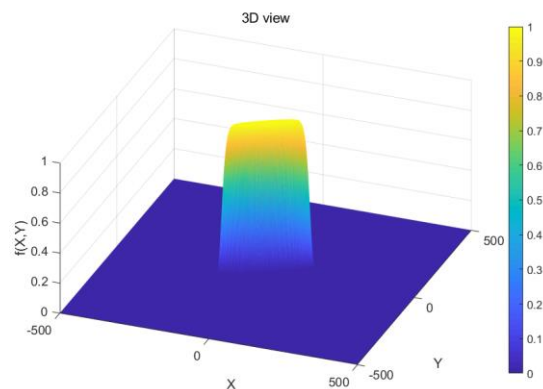
(a) Horizontal plane (ordinary, $\lambda = 1$)(b) 3-D view (ordinary, $\lambda = 1$)(c) Horizontal plane (Super, $\lambda = 8$)(d) 3-D view (Super, $\lambda = 8$)

Figure 3.1. Schematic of ordinary and super Gaussian function

1

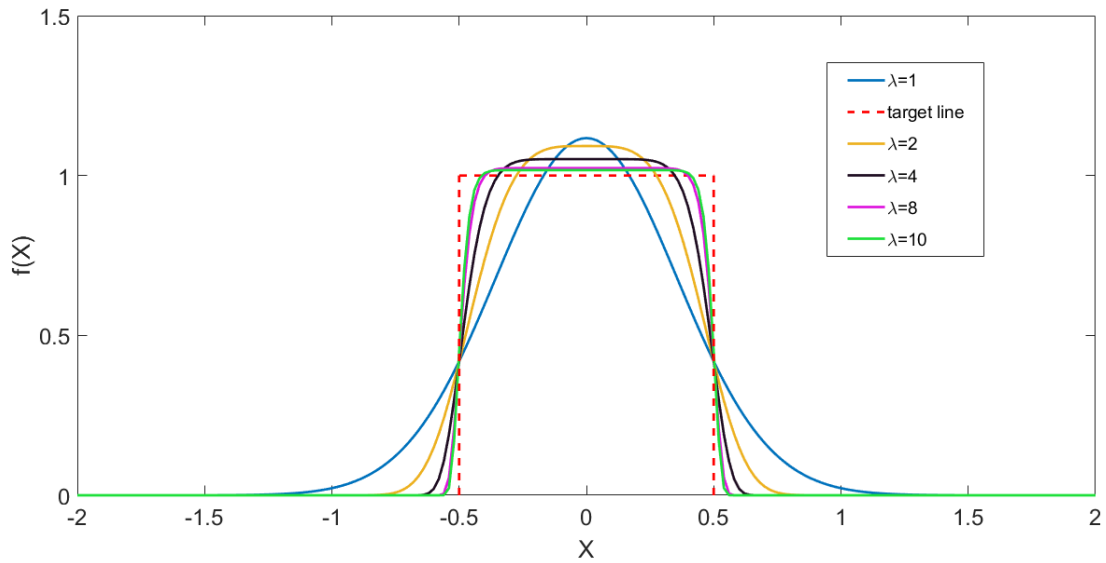


Figure 3.2. The shape of the super-Gaussian function with different powers

2

3 3.2.2 Super-Gaussian function

4 Nonetheless, the ordinary Gaussian function may not be the most appropriate one because
 5 of nonuniformity and tail error. Simply, if we use a one-dimensional normal distribution with
 6 the variance σ_x^2 to approximate the target line segment, it can be obviously noticed that the
 7 shape deviates from the target since the value is not uniform in the middle area (in **Fig. 3.2**).
 8 Besides, there are residual errors in the two tails where the value should be 0. Thereby, this
 9 model will bring extra error, unwanted in the inference. To minimize this error, this paper
 10 applied the super-Gaussian function (Parent et al., 1992) instead of the ordinary one.

$$f_{\mathbf{x}}(x_1, \dots, x_m) = A \exp \left[- \left(\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)^\lambda \right] \quad (3.4)$$

11 This equation modifies the shape with a flat-top and Gaussian fall-off by raising the
 12 content of the exponent to a power λ . It can be observed in **Fig. 3.2** that when λ increases, the
 13 slope becomes precipitous and the flat top becomes longer, which makes the shape closer to the
 14 target shape. In this research, λ was set to eight to balance between error control and
 15 calculation simplicity. The two-dimensional super-Gaussian function is

$$f(x, y) = A \exp[-(a(x - x_0)^2 + 2b(x - x_0)(y - y_0) + c(y - y_0)^2)^\lambda] \quad (3.5)$$

1 and the shape is shown in **Fig. 1(c)(d)**.

2 In real urban areas, most line sources can be regarded as a line in the horizontal plane,
 3 which is close to the surface. Moreover, the right-hand side of Eq. (2.23) involves huge matrixes
 4 product in the numerical calculation, if the position of the line source has no restraint. Therefore,
 5 this paper only deals with a line source in a horizontal plane with a fixed height above the
 6 surface. The height of the line is included in the background information I and does not need
 7 to be estimated, due to which Eq. (3.5) is adequate for this research. In spite of that, it is worth
 8 mentioning that this does not mean that the super-Gaussian function is incapable of estimating
 9 the height or vertical inclined angle of other lines. In fact, that kind of estimation can be
 10 acquired quickly by increasing the dimension k to three and adding two more correlation
 11 coefficients in Σ . Hence, the parameters vector of the line sources can be written as

$$\mathbf{s} = (x_s, y_s, \theta, \sigma_X, \sigma_Y, q_s) \quad (3.6)$$

12 In the application, these six parameters will be estimated.

13

14 **3.3 Case I: numerical experiment of urban boundary layer**

15 To evaluate the applicability of the proposed method, this case uses it to estimate the
 16 parameters of a line source in a simple urban area, based on computational fluid dynamics (CFD)
 17 simulation data. One of the most elementary dispersion cases, a steady-state urban boundary
 18 layer without obstacles on the ground, is analyzed. This setting is appropriate to explore the
 19 basic properties of the proposed method.

20

21 **3.3.1 Numerical simulation**

22 It is necessary to run a forward CFD simulation first to synthesize the measurements, then
 23 adjoint equation simulation is run based on the forward flow field.

24 The calculation domain of the forward simulation, with a size of $1500 \text{ m}(x) \times$
 25 $1000 \text{ m}(y) \times 500 \text{ m}(z)$, is shown in **Fig. 3.3**. In the x and y directions, the cubic grid is

1 uniform, with 10 m length. To ensure that the line source has a sufficiently high resolution in
 2 the vertical direction and to simulate the dispersion behavior precisely near the bottom surface,
 3 the finest mesh is generated from the bottom, with a 0.2-m length. The mesh becomes rougher
 4 along with the height at a 1.02 ratio. To generate the effects of the urban boundary layer, a half
 5 channel flow is produced by mapping the flow properties of the outlet to the inlet. After each
 6 step, the pressure gradient is adjusted to keep the mass flow rate constant, at $\bar{U} = 4$ m/s. At
 7 the bottom wall, the wall function proposed by Blocken et al. (2007) is applied. k_s is set as
 8 45 m, which corresponds to a Davenport roughness of 1.5 m, a common case for urban areas.
 9 The turbulence is modeled via a standard $k - \varepsilon$ model. The details of other configuration are
 10 listed in **Table 3.1**.

Table 3.1. Numerical schemes and boundary conditions for Case I

Time marching	Steady state, Semi-Implicit Method for Pressure-Linked Equations (SIMPLE) method.
Spatial discretization	Advection term: total variation diminishing (TVD) scheme;
Inlet	U, k, ε, D_t : mapped Outlet; C : constant (= 0); C^* : zero-gradient
Outlet	Flow: zero-gradient; C : zero-gradient; C^* : constant (= 0)
Top wall	Slip;
Bottom wall	Flow: generalized logarithmic law; Roughness height: 1.5 m Davenport Roughness; C, C^* : zero-gradient
Source	C, C^* : Constant injection rate (= 1 1/s)

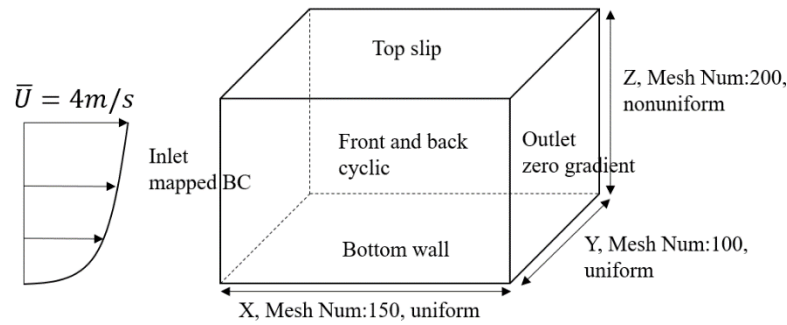


Figure 3.3. The calculation domain and boundary condition of the numerical simulation in Case I

1 The $100\sqrt{2}$ -m-long line source is settled on the horizontal plane, 10 m above the bottom
 2 wall. The coordinate of the middle point is (750 m, 500 m, 10 m). The inclined angle is $\pi/4$
 3 from the positive x - and z -axes directions. In the numerical simulation, this line source is
 4 resolved by discrete cells that are connected along the vertical edge. With the increase of
 5 iterations, the boundary layer gradually became steady. The wind profile along the vertical edge
 6 in the middle of the outlet settled down eventually after 90000 iterations as shown in **Fig. 3.4**.
 7 In this research, the data of iteration 130000 is used to make sure the flow is in a steady state.
 8 The steady dispersion field of the line source is presented in **Fig. 3.5**. It can be confirmed from
 9 the concentration profile that the influence of the dispersion reaches up to 200 m height in the
 10 boundary layer.

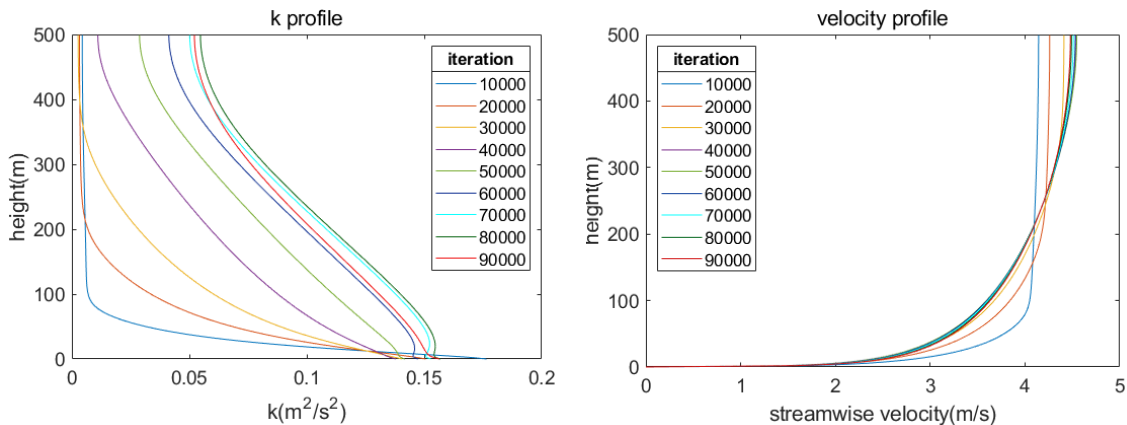
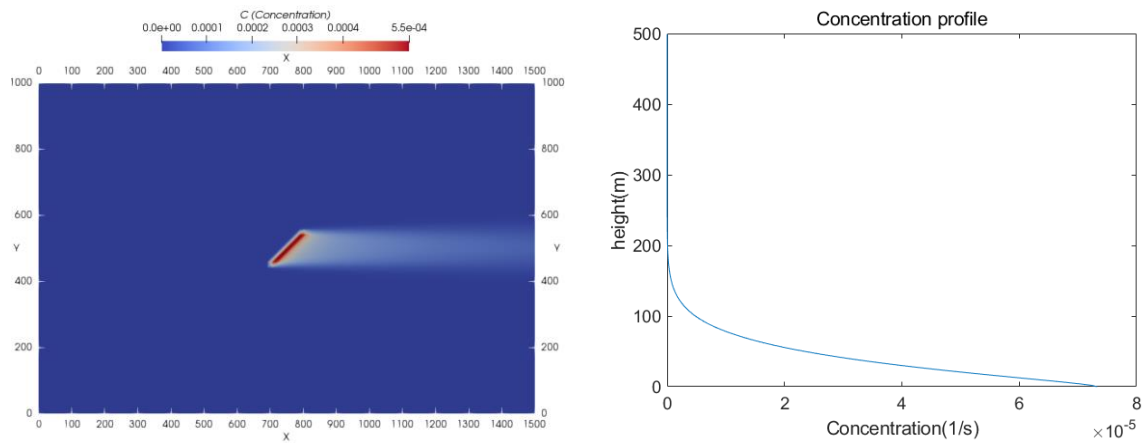


Figure 3.4. Profile along the vertical line in the middle of the outlet
 (k: turbulent energy)

1



The dispersion field in the horizontal plane (height = 10 m)

The concentration profile along the vertical line in the middle of the outlet

Figure 3.5. Steady dispersion field of the line source

2 For the sensor configuration, a common setting in STE research is the uniform distribution
 3 in the domain (Keats et al., 2007b; Xue et al., 2018b). However, the distribution of stationary
 4 sensors cannot be uniform among complicated urban terrain in real life. Therefore, a random
 5 walk was added to the location of each sensor with a uniform configuration (**Fig. 3.6**); this
 6 could also test the robustness of the proposed method. The coordinate was changed from
 7 (x_m, y_m) to $(x_m + 10 \times [N[-10, 10]], y_m + 10 \times [N[-5, 5]])$. $N[a, b]$ represents the
 8 uniform distribution bound between a and b . The sensor heights are also different to simulate
 9 the real application. The adjoint concentration field was calculated in the same domain with a
 10 reverse velocity field for each sensor. One particular adjoint transport field from a sensor is
 11 shown in **Fig. 3.7**.

12 The simulated concentration at each sensor is regarded as measurement \mathbf{D} in this case.
 13 As these concentrations are synthesized via the simulation, random noises $N[-0.5, 0.5] \times D_i$
 14 are added to them to ensure that this case study is closer to reality.

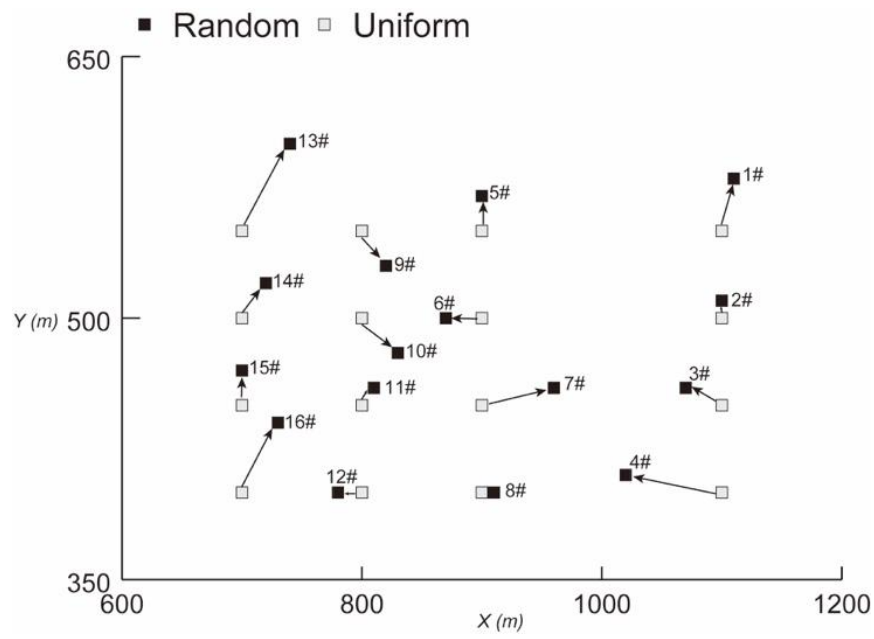


Figure 3.6. The sensor configuration in Case I

1

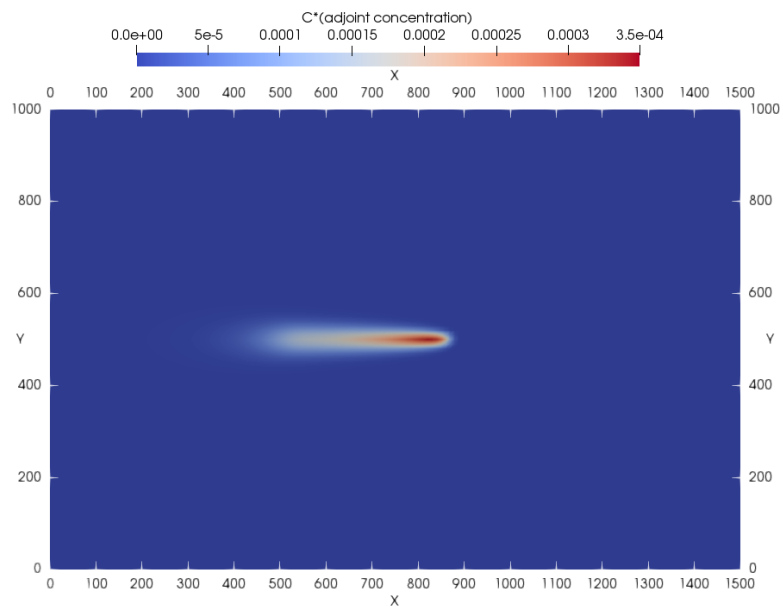


Figure 3.7. Adjoint transport field from a sensor in a horizontal plane (height = 10 m)

2

3.3.2 Bayesian inference settings

Bayesian inference was conducted after the simulation, as described in **Chapter 2**. In this process, one of the most critical factors is the error variance $\sigma_{d,i}^2 + \sigma_{m,i}^2$ in Eq. (2.8). The value of this factor subtly affects the effectiveness and accuracy of the estimation procedure. Conversely, the optimal value of this factor is difficult to determine and still in controversy. Keats et al. (2007b) assigned the values manually for each sensor, according to the observed trends. These values lie between 10% and 270% of the mean measurements. Xue et al. (2018b) applied the variance of the measurements of each sensor. Since the determination of this factor still requires further study and is not the core part of this paper, this research assigns values in the same way as Keats et al. (2007b); the values lie between 10% and 250% of the mean forward measurements. Meanwhile, because a minor variance could cause unsteadiness in the calculation of Eq. (2.8), the variances of the sensors that could barely measure concentration were assigned a fixed value. The Bayesian inference of STE was executed with MATLAB on a personal computer with Intel® Core™ i7-6700 CPU @ 3.4GHz and 16GB of RAM. The averaged computational time for estimating one source is about 42 min.

3.3.3 Estimation results

Fig. 3.8 shows the estimation results as the marginal probability distribution of each term. When the sensor measurements are combined with the super-Gaussian function, Bayesian inference can provide an ideal estimation for the line source, regardless of the initial guess provided to the Markov chain. All the parameters are correctly estimated. The coordinates, length, and strength estimations agree well with the true values. The peak value of width is no more than 20 m from the true value. Considering that the resolution of the simulation data is 10 m rough, the proposed method is considerably accurate. The discrepancy for the angle is also smaller than 10 degrees. Considering that the posterior probability density functions (PDFs) in the results are usually highly skewed, like **Fig. 3.8(d)**, only one estimator may not be able to objectively represent the distribution of estimations, we select both the 50th percentile and standard mean values of PDFs as the parameter estimators, the comparison between the true values and estimations are summarized in **Table 3.2**.

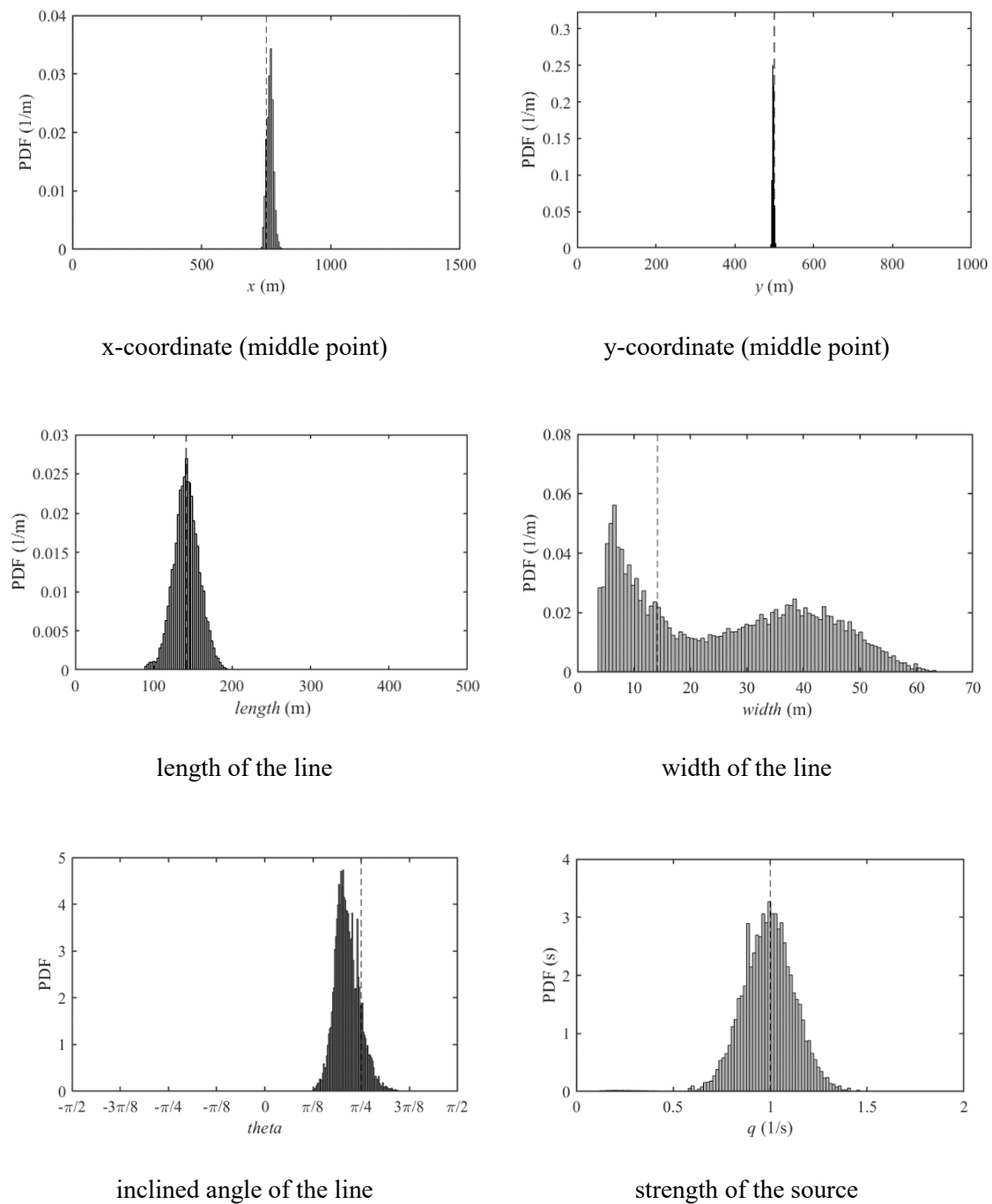


Figure 3.8. Posterior PDFs of source parameters in Case I (dotted line: true value)

- 1 The expectation of line source geometry $\int p(\mathbf{s}|\mathbf{D}, I) \times f(x, y|\mathbf{s})d\mathbf{s}$ is shown in **Fig. 3.9**.
- 2 It can be confirmed that the proposed method restricts the estimation results to a small area
- 3 around the true line. The shape of the contour is also an ellipse, with a similar line source angle.
- 4 Therefore, it can be concluded that under ideal conditions, without measurement and modeling
- 5 errors, the proposed method successfully identifies the line source information.

1

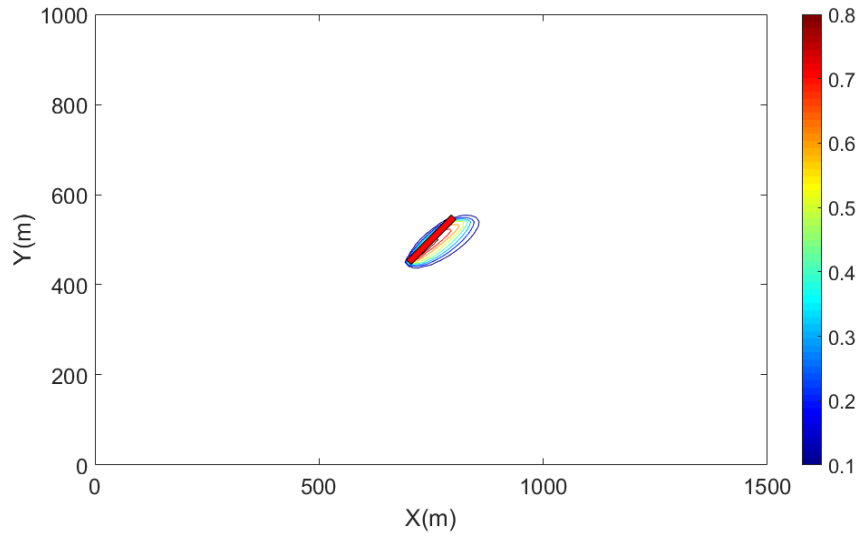


Figure 3.9. The expectation of line source geometry $\int p(\mathbf{s}|\mathbf{D}, I) \times f(x, y|\mathbf{s}) d\mathbf{s}$ estimated by the proposed method in Case I. (red patch: true line source)

2

Table 3.2. Summarized estimation results of the proposed method

Method	x_s (m)	y_s (m)	length(m)	width(m)	angle	strength(1/s)
True value	750	500	$100\sqrt{2}$	$10\sqrt{2}$	$\pi/4$	1
Estimations (50 th percentile)	765	496	142	25	$\pi/4.7$	0.99
Estimations (standard mean)	764	496	142	26	$\pi/4.7$	0.99

3

4 3.4 Discussion about sensor configuration

5 Sensor configuration is crucial for all STE methods as sensor measurements are essentially
6 the only information on which the STE method is based. Considering that there are limited
7 discussions concerning the sensor configuration for line source estimation in the literature, it is

1 necessary to study the basic properties of the line source's dispersion. Furthermore, the author
 2 wants to analyze the success and errors of line source estimation using different sensor
 3 configurations based on Case I. Notably, the contents in this section apply to all the line source
 4 identification methods based on stationary sensors, rather than to a specific design for the
 5 current method.

6

7 3.4.1 Uniform sensor configuration

8 In previous research regarding STE methods for point sources, a regular sensor network
 9 downstream of the source has commonly been used, yielding successful estimations. However,
 10 the situation has entirely changed in estimating the line source as it is more concerned with
 11 geometric information, which makes STE more difficult with a conventional sensor
 12 configuration.

13 In the first configuration, the sensor network is like the conventional ones used for the
 14 point source estimation in previous research studies (Keats et al., 2007b; Kumar et al., 2015b;
 15 Xue et al., 2018b). Only one sensor is in front of the source and all the others are at the
 16 downstream with regular positions. The schematic diagram is in **Fig. 3.10**.

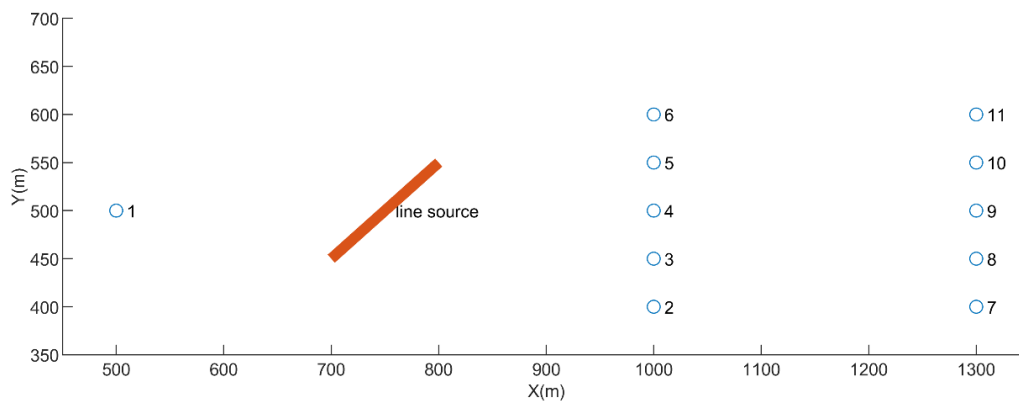


Figure 3.10. Schematic of uniform sensor configuration. (red patch: true line source)

17 **Fig. 3.11** shows the marginal parameter distributions of the posterior probability. It can be
 18 noted that the estimation of the middle point coordinate, strength, and length is proper with
 19 small discrepancies from the true value. There is only one peak in each histogram. On the other
 20 hand, the angle and width results are totally distorted. It seems that the sampling algorithm

- 1 randomly travels in the width space, losing the ability to identify the true value. The accuracy
 2 is even worse in the angle estimation, in which the result is totally opposite to the true value. In
 3 consequence, the proposed method, coupled with the conventional sensor set, fails to estimate
 4 all the parameters of the line source, especially in estimating the angle and width.

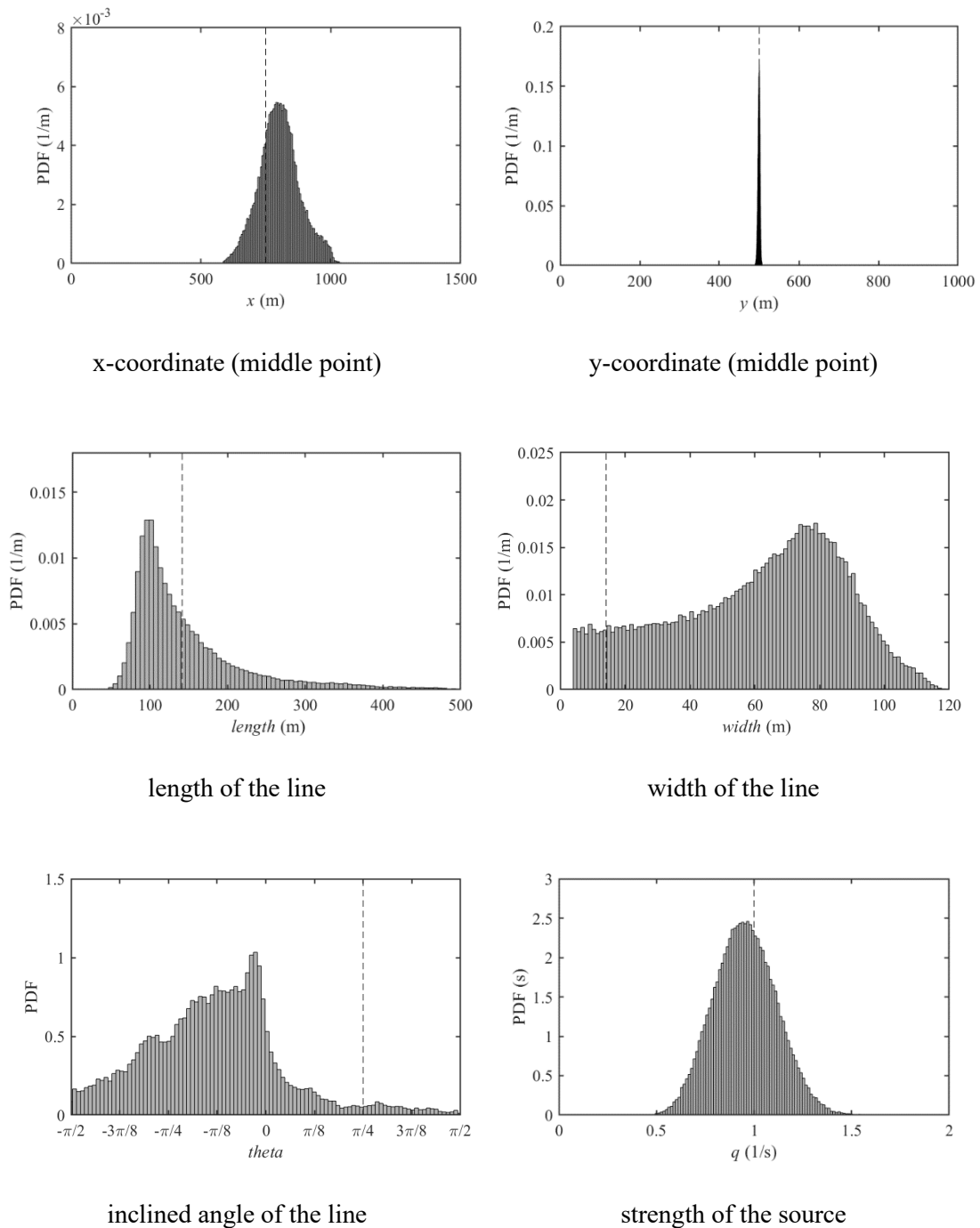


Figure 3.11. Posterior PDF of source parameters with uniform sensor configuration.

(dotted line: true value)

1 Since the conventional configuration has been proved effective enough by previous
 2 research studies, and the proposed method has successfully identified the line source in **Section**
 3 **3.3**, the problem should be in the different properties between the point source and line source.
 4 Actually, if we check the results of the conventional configuration in **Fig. 3.11**, even though the
 5 angle and width estimations are completely wrong, the coordinates of the middle point are still
 6 correctly estimated. This fact implies that the conventional configuration is capable of locating
 7 the position of the source, which is enough for a point source estimation. If geometric
 8 information is needed, like the inclined angle of the line, the conventional set cannot be relied
 9 on anymore. This can be illustrated by the three different line sources located at the same
 10 position as shown in **Fig. 3.12**. Although these three line sources are distinct from each other,
 11 their dispersion fields, caused by the same flow and the influence on the sensor downstream,
 12 are similar to a high degree. The only difference that can be captured by the downstream sensors
 13 is the concentration discrepancy in the span-wise direction, which is vulnerable under the noise
 14 of the model and measurement. Therefore, the inclination angle of the line source cannot be
 15 estimated through the measurement of the regular sensor configuration downstream of the
 16 source. For the same reason, the width is also free to change and is thus impossible to estimate.
 17 Accordingly, the estimation of the line source requires more geometric knowledge, which can
 18 be provided only by the sensors close to the circular area, whose diameter is the same as the
 19 length of the line source.

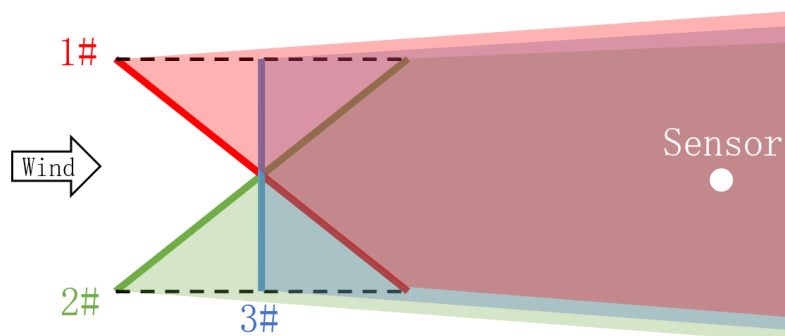


Figure 3.12. Three different line sources with similar concentration fields (areas with different colors)

20 The inclination angle and width were successfully estimated in Case I. It is reasonable to
 21 believe that the sensors around the line source (No. 9, 14, 15, 16) are the key factor. During the
 22 inference, the information provided by these sensors prevents the generated source from over-

1 rotating or expanding in width. Despite of that, the sources generated during sampling still
 2 fluctuated to some degree because the density of sensors is not high enough. It can be observed
 3 in **Fig. 3.8** that second peak appeared in the posterior PDF of width, and the estimation of angle
 4 also has a deviation. It means that a wider line source with smaller inclined angle would have
 5 similar impacts on the current sensor configuration, which is intuitively shown in **Fig. 3.9**.

6

7 **3.4.2 Importance of null measurements**

8 Apart from these indispensable sensors near the source, the effects of the “region of
 9 influence” should also be considered. This concept was proposed by Keats et al. (2007b),
 10 emphasizing the importance of sensors with null measurements to the point source estimation.
 11 The estimation of the line source has the same requirements. Considering the case shown in
 12 **Fig. 3.13**, whereby 9 sensors are used to estimate the red line source, the incorrect estimation,
 13 i.e., the green line, will destroy the inference because the red and green lines have almost equal
 14 influence on the sensors. The failure is essentially due to the lack of sensors in the upper-right
 15 area; although sensors in this area will measure none of the true red line, only the null
 16 measurement informs the inference of the length of the source and excludes the possibility that
 17 the green line is the true source. According to Keats et al. (2007b), this kind of sensor should
 18 also be classified in the ‘region of influence’.

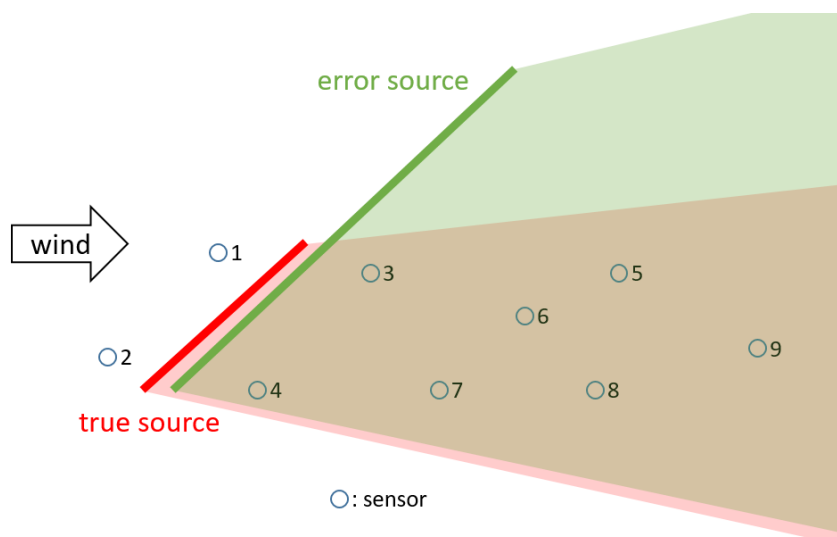


Figure 3.13. Diagram of typical estimation error caused by the lack of null measurement

19 This characteristic can be better demonstrated by the following trial with Case I. This

1 sensor configuration is totally produced by randomness. The horizontal distribution of sensors
 2 for the trail is shown in **Fig. 3.14**. The configuration has several sensors around the line source,
 3 which should improve the estimation accuracy when compared to the uniform configuration.

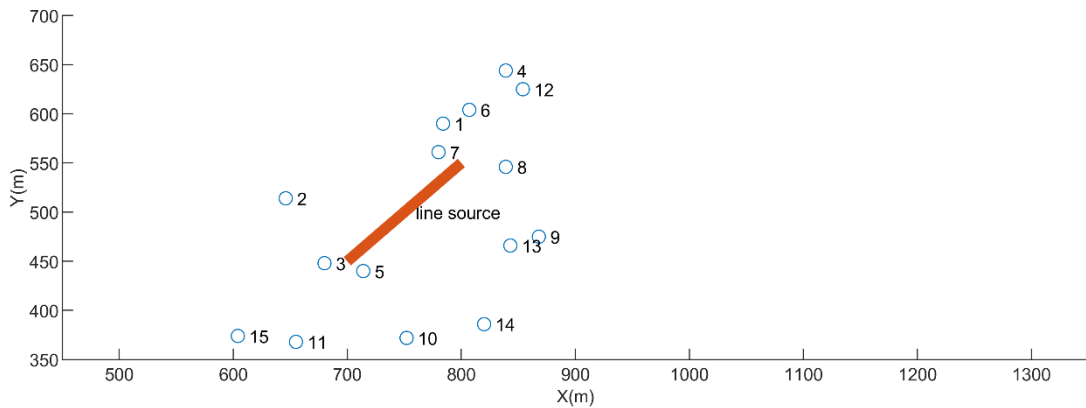


Figure 3.14. Schematic of sensor configuration (totally random)

4 The estimation result is shown in **Fig. 3.15**. From the first look, all the parameters are
 5 almost correctly estimated. The accuracy of the coordinates, length, and strength is the same or
 6 even better than that of conventional configuration. More importantly, the inference
 7 successfully narrows the width scale to under 25 m. Similarly, the estimation of the angle is
 8 improved significantly using the measurements from this configuration. The inference
 9 eliminates the possibilities of other angles and gathers only around the true value. It is
 10 reasonable to believe that the sensors around the line source (No. 3, 5, 7, 8) cause the
 11 improvement. During the inference, the information provided by these sensors prevents the
 12 generated source from rotating too much or expanding the width.

13 Despite this advancement over the conventional configuration, the Bayesian inference will
 14 be found unsteady and sensitive to the initial guess if different Markov chains with different
 15 start points are tested. Although for most of the time, the estimations are similar to the one in
 16 **Fig. 3.15**, one particular error in **Fig. 3.16** will appear repeatedly from time to time. In this
 17 particular error, wrong peaks show up in the x coordinate as well as the length and dominate
 18 over the true peak, which is unacceptable in practical use. Hence, the measurements provided
 19 by this configuration are still not enough to reach a steady and accurate inference.

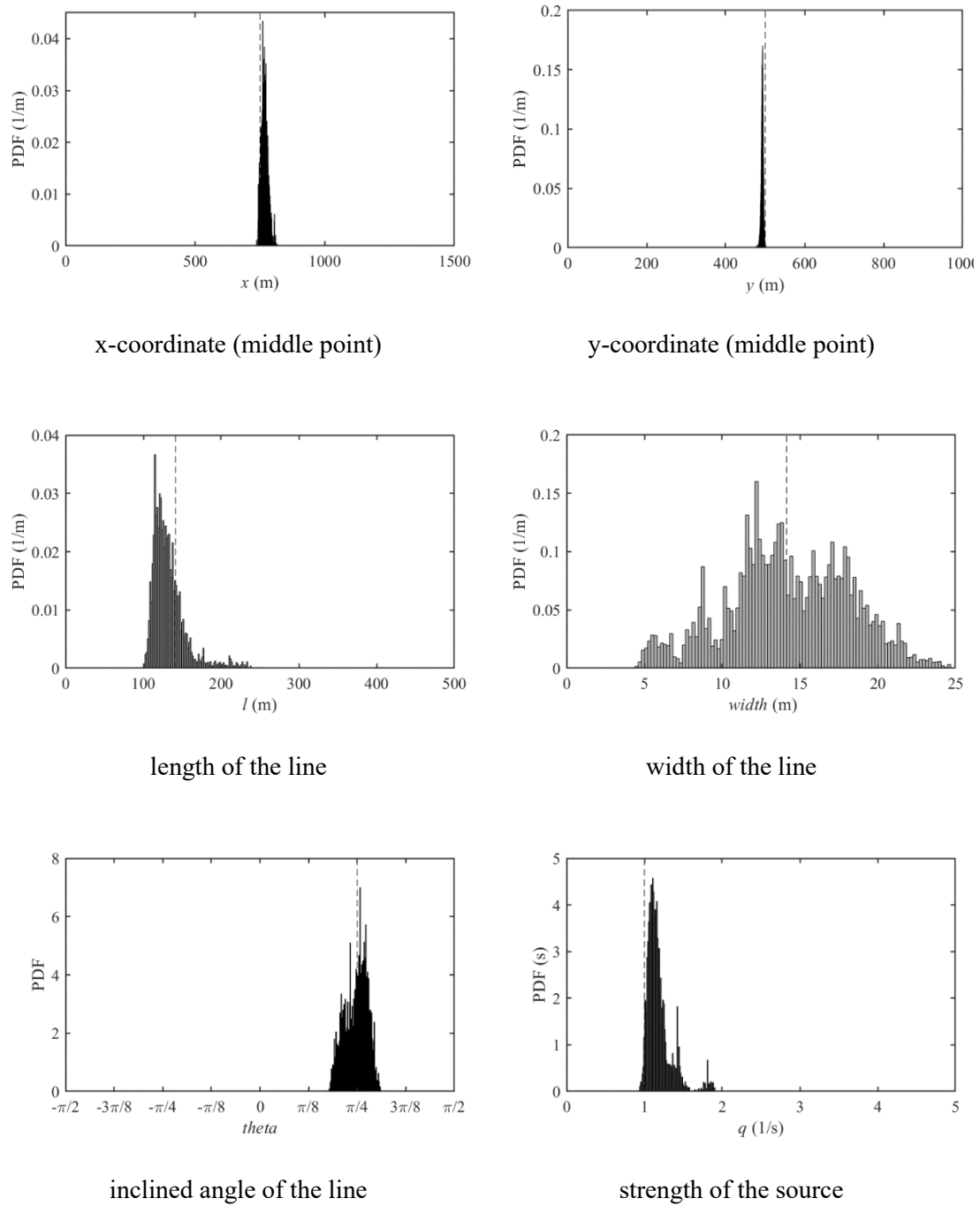


Figure 3.15. Posterior PDF of source parameters with sensor configuration of **Fig. 3.14**.

(dotted line: true value)

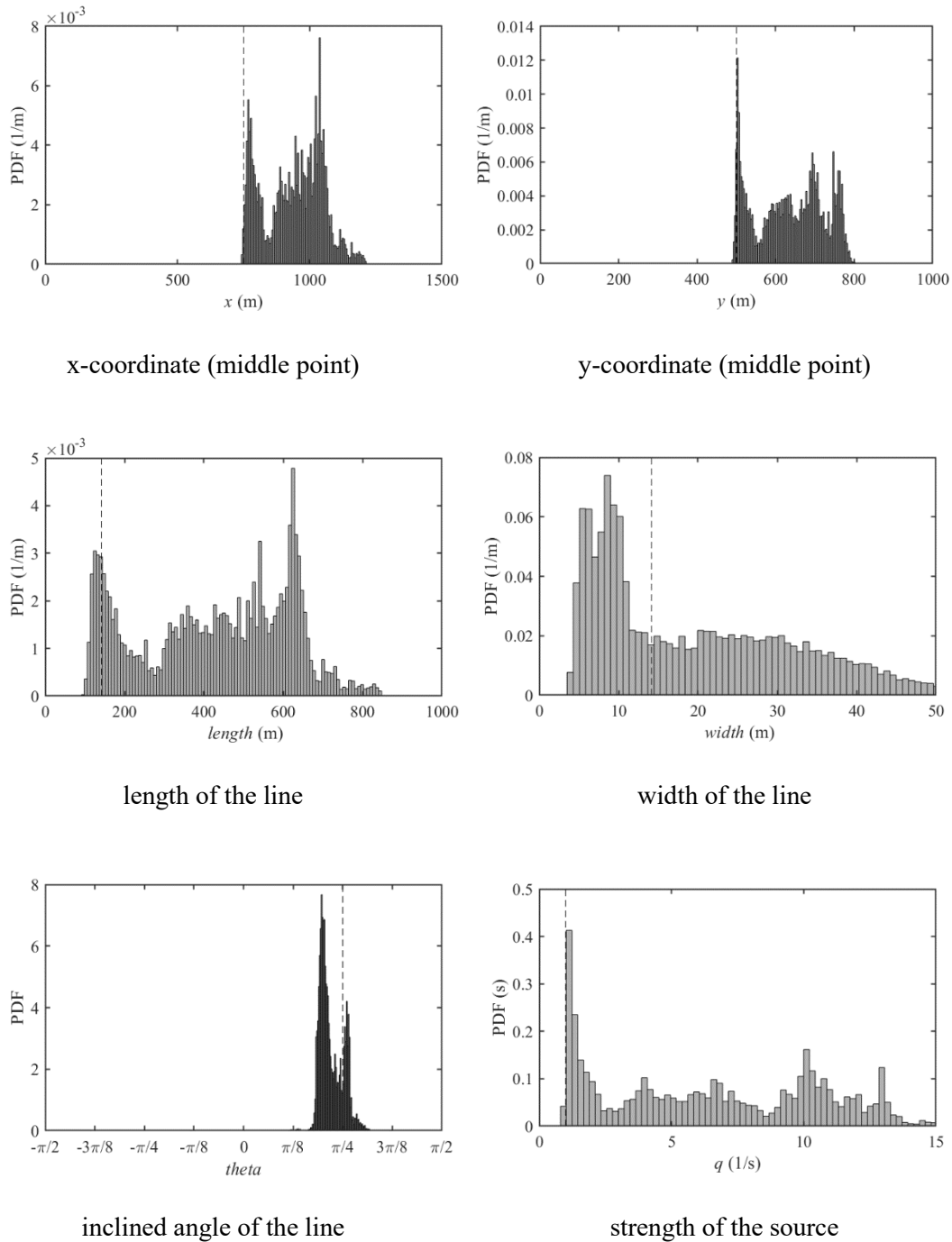


Figure 3.16. Posterior PDF of source parameters with sensor configuration of **Fig. 3.14**.

(a particular wrong example) (dotted line: true value)

- 1 The failure of the length estimation is mainly caused by the lack of sensors at the upper
- 2 right corner downstream to the line source. In this case, the line source can extend freely in that
- 3 direction without significantly changing the influence on each sensor. One particular error is
- 4 shown in **Fig. 3.17**. Consequently, apart from the sensors near the line source, we still need

1 some null-measurement sensors to prevent the error extension.

2 In **Section 3.3**, the length was well reconstructed. The zero measurements provided by
 3 sensors No. 1, 4, 5, 8, and 12 prevent error extension. Otherwise, the difference between the
 4 measured and modeled values would become significant, and the sampling algorithm (MHMC
 5 in **Chapter 2**) would abandon the sampling.

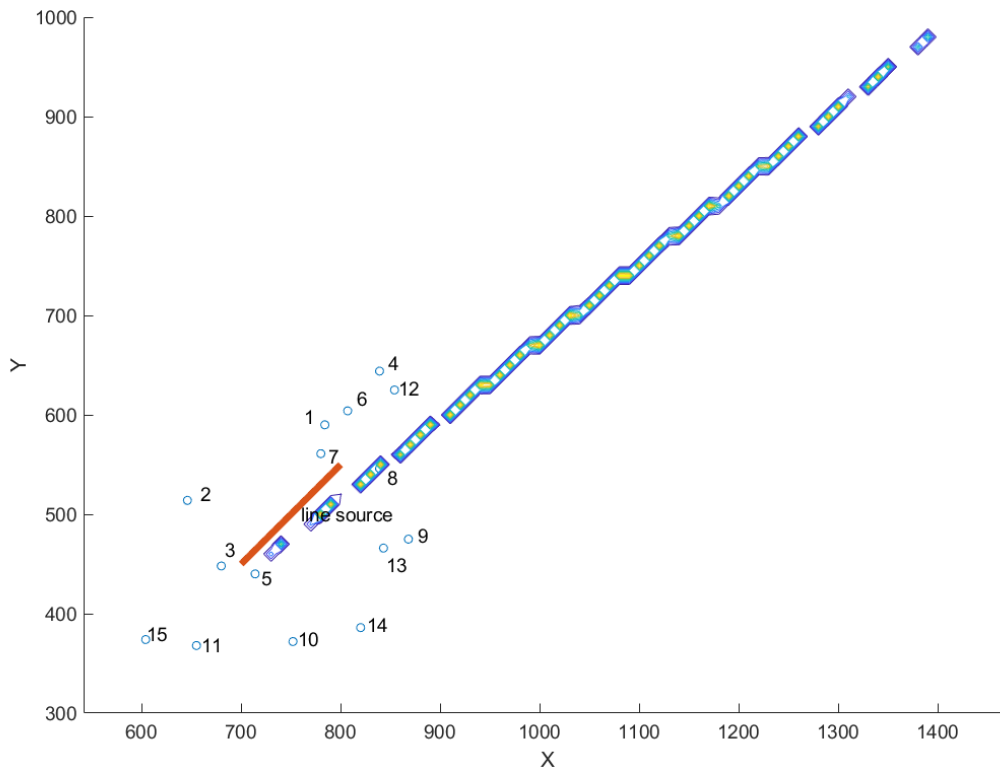


Figure 3.17. Wrong estimation with sensor configuration of **Fig. 3.14**.

(red patch: line source; blue line: wrong estimation)

6

7 **3.4.3 Necessary sensors for line source estimation**

8 It is natural to ask how can we know a set of sensors is enough for the estimation or not?
 9 Or, what is the minimum number of the necessary sensors for the line source estimation? Based
 10 on the analysis above, it is proposed here that at least six sensors located appropriately are
 11 enough for the purpose. The schematic diagram is in **Fig. 3.18**. Among them, three sensors
 12 around the source are mainly for the angle and width estimation. The upper one and lower one
 13 in the downstream can restrain the length. It does not mean that each sensor is only in charge

1 of the estimation of one parameter. In fact, all the sensors are helpful to the estimation of several
 2 parameters, but certain sensors are more important to the estimation of certain parameters.

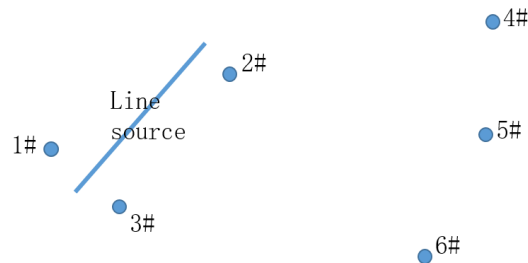


Figure 3.18. Schematic of the necessary sensors for the line source estimation

3 In order to prove the correctness of the proposal above, six sensors (No. 1, 3, 4, 9, 15, 16)
 4 are picked out from the configuration of Case I in **Section 3.3** according to **Fig. 3.18** to conduct
 5 a new STE. The results are summarized in **Fig. 3.19**. Apparently, the accuracy is almost
 6 consistent with that of Case I in **Fig. 3.8**. All the parameters are successfully estimated with
 7 tolerable error. Without a doubt, the samplings in **Fig. 3.8** gathered tightly around the true value
 8 while in **Fig. 3.19** more residual probability appears at wrong values. The redundant 10 sensors
 9 in Case I infuse more information to the inference so that the true value is captured effectively.
 10 This fact is also well illustrated by the cumulative probability density of six sensors in **Fig. 3.20**.
 11 Thereby, we can conclude that, in order to estimate the line source term, both for position and
 12 geometry, six sensors with proper location are indispensable. Any sensor network including
 13 these six sensors can perform well in line estimation cooperated with the super-Gaussian
 14 function method. Apart from these six sensors, the more sensors we have, the more effective
 15 and accurate will the inference be.

16

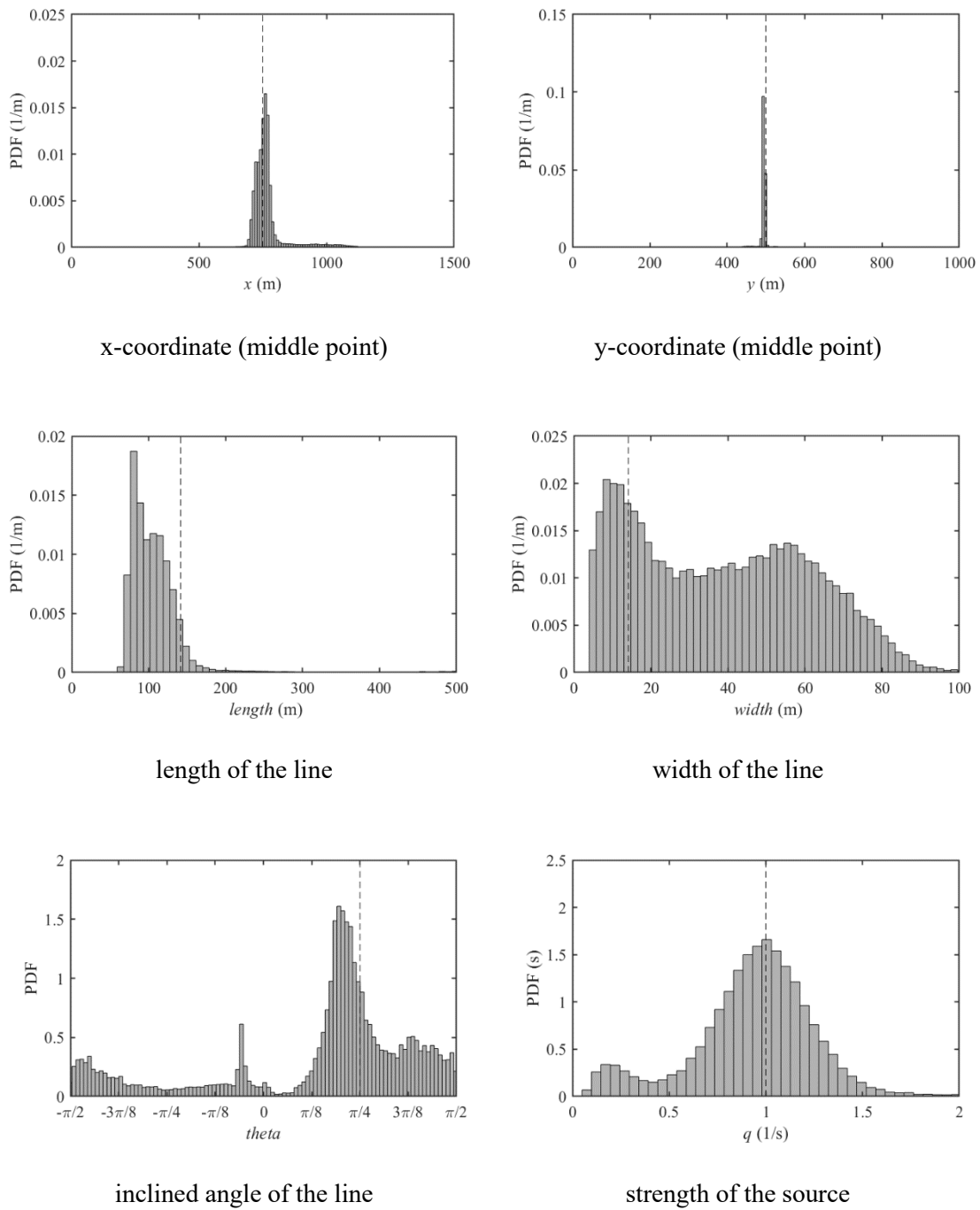


Figure 3.19. Posterior PDF of source parameters (selected 6 sensors from Case I).

(dotted line: true value)

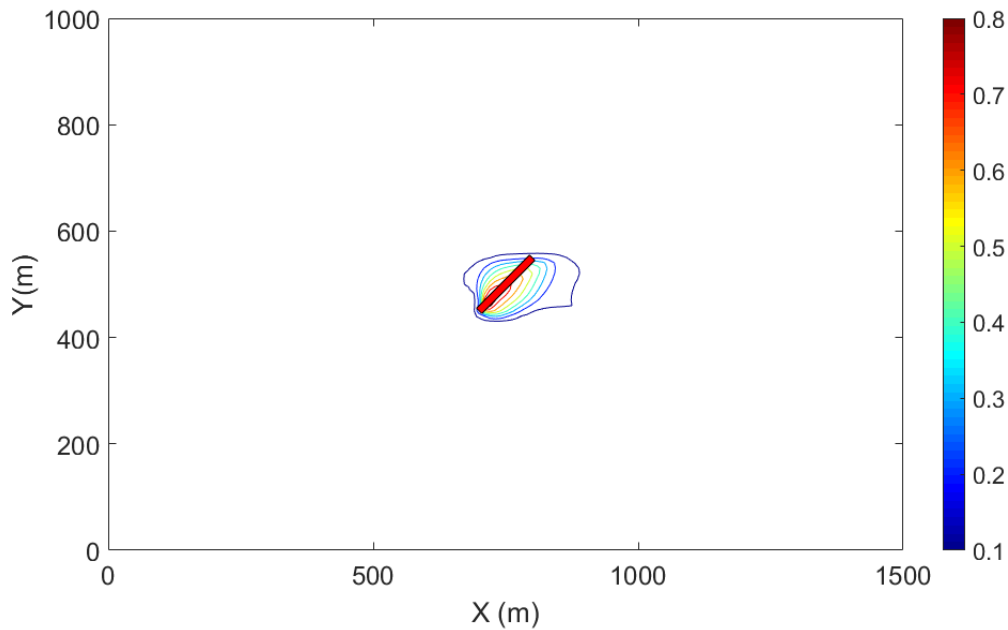


Figure 3.20. The joint posterior probability density distributions $P(x_s, y_s, \theta, \sigma_X, \sigma_Y | D, I)$ of selected 6 sensors. (red patch: true line source)

1

2 **3.5 Case II: wind tunnel experiment of the urban square model**

3 After the success with the elementary case presented above, it is necessary to examine the
 4 robustness of the proposed method under the complicated measurement and modeling errors
 5 within the practical case. Although the line source's dispersion has been frequently studied in
 6 the literature, there are very few experimental databases accessible online for validating the
 7 inversion methodology. Balczó and Lajos (2015) measured the dispersion of a line source in a
 8 simplified urban square, as shown in **Fig. 3.21(a)**. Their case includes buildings of different
 9 sizes, and the configuration is representative of real urban neighborhoods. In addition, they
 10 performed a comprehensive measurement in the space between buildings, which makes the
 11 STE possible. Hence, we demonstrate the feasibility and necessity of the proposed method
 12 based on this practical case.

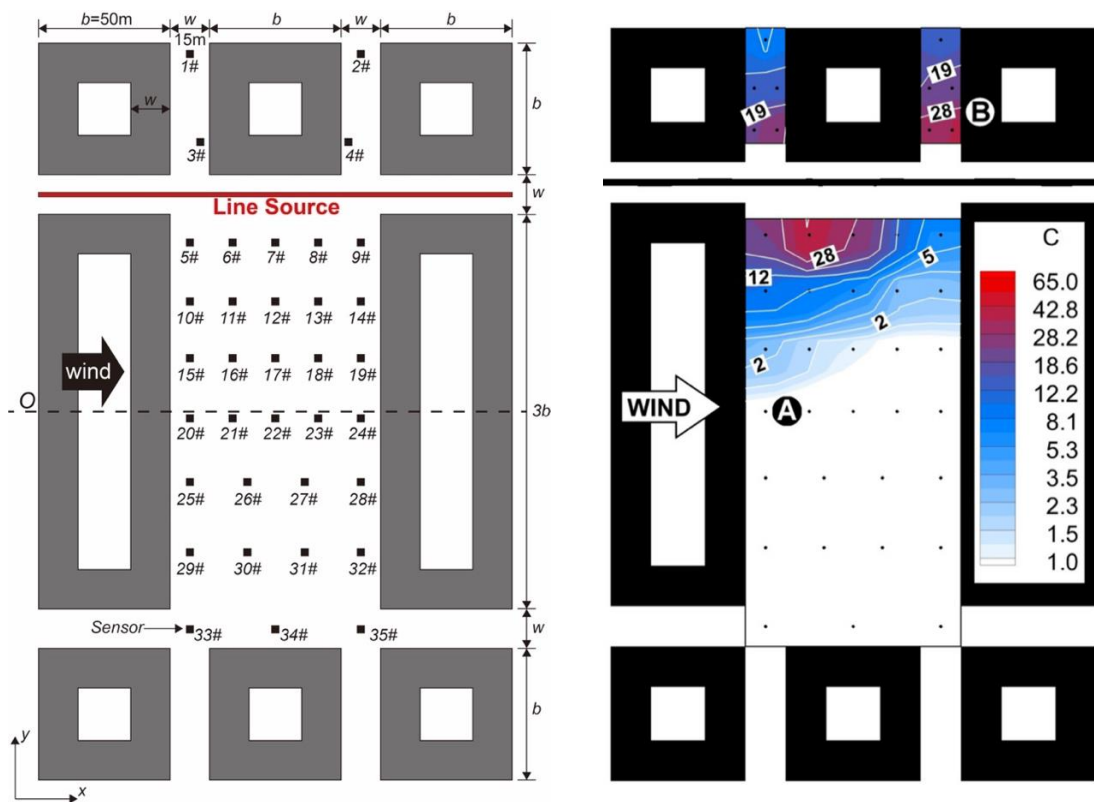
13 Meanwhile, to illustrate the necessity of geometry estimation, we also conducted the STE
 14 via the conventional method, with the ideal-point assumption. The results are compared with
 15 those obtained using the proposed method.

16

1 3.5.1 Wind tunnel experiment for measurements

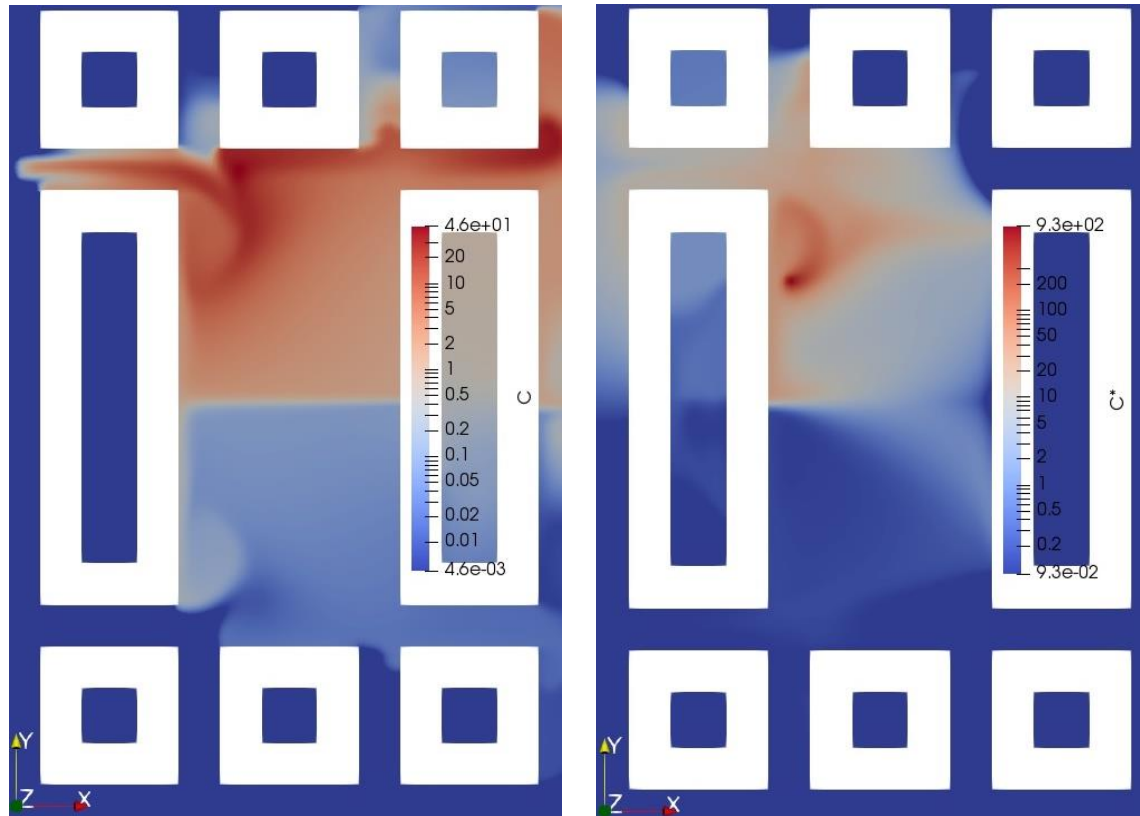
2 In this case, measurements D were obtained via the wind tunnel experiment in Balczó
 3 and Lajos (2015), conducted in an Eiffel-type wind tunnel with a closed test section of
 4 $0.5\text{ m} \times 0.5\text{ m}$. All buildings had the same height. The model had a scale of 1:650, which means
 5 the building height was $H = 30\text{ m}$ at full scale and $h = 46\text{ mm}$ at the model scale. The
 6 turbulent boundary inflow generated by roughness moved from left to right (**Fig. 3.22**). The
 7 reference velocity U_r was measured at the building height.

8 The line source was located on the bottom wall in the upper part of the building area, as
 9 shown in **Fig. 3.21(a)**. The gas emitted from the line source was measured using slow flame
 10 ionization detector sensors located at 35 positions in a horizontal plane with a height of $H/10$.
 11 The measured concentration was non-dimensionalized via the reference concentration $\bar{C} =$
 12 $\frac{Q}{U_r h^2}$, where $Q[\text{g/s}]$ is the source strength. The non-dimensionalized concentration distribution
 13 in the horizontal plane of sensors is shown in **Fig. 3.21(b)**.



The basic settings for the simplified urban square model

The measured distribution of concentration emitted from the line source in the WTE



The simulated distribution of concentration emitted from the line source

The simulated adjoint tracer field emitted from the No. 15 sensor

Figure 3.21. Schematic of the wind tunnel experiment and simulation in Case II (horizontal plane with $z = H/10$): (b) is referred to Balczó and Lajos (2015).

1

2 3.5.2 Numerical simulation

3 A CFD simulation was conducted to calculate the flow field and adjoint tracer field. The
 4 calculation domain was $26H(x) \times 10.67H(y) \times 6H(z)$. The building area was $10H$
 5 from the inlet and outlet and $5H$ from the top boundary. The experimental model was close to
 6 the side walls in the small wind tunnel and thus the distance between the building area and
 7 lateral boundaries was also set small, as $0.67H$, to simulate the effects of side walls. In the
 8 building area, the cubic mesh was uniform, with a $H/18$ (1.67 m) length in all directions, then
 9 expanded along with the coordinates at a 1.08 ratio.

10 The boundary conditions were set according to the AIJ basic guidelines (Tominaga et al.,
 11 2008b). The streamwise velocity and turbulence intensity profiles were prescribed on the inlet

1 based on the experimental data, as shown in **Fig. 3.22**. At the bottom and side walls, a no-slip
 2 condition with the Spalding wall function was applied. The turbulence was modeled using the
 3 renormalization group $k - \varepsilon$ model (Yakhot et al., 1992). The details of other configurations
 4 are listed in **Table 3.3**. The adjoint transport field was calculated in the simulation domain, with
 5 a reverse velocity field for each sensor in the wind tunnel experiment. One adjoint transport
 6 field from the No. 15 sensor is shown in **Fig. 3.21 (d)**.

Table 3.3. Numerical schemes and boundary conditions for case II

Time marching	Steady state (SIMPLE method)
Spatial discretization	Advection term: TVD scheme;
Inlet	\bar{U}, k, ε : experimental profile; Pressure: zero-gradient C : constant (= 0); C^* : zero-gradient
Outlet	Flow: zero-gradient; Pressure: fixed value 0; C : zero-gradient; C^* : constant (= 0)
Top wall	Slip;
Bottom & side walls	Flow: Spalding wall function; Pressure: zero-gradient; C, C^* : zero-gradient
Source	C, C^* : Constant injection rate

7 In the numerical simulation, apart from the hypothetical tracers, the dispersion from the
 8 line source was also predicted to reveal the model error. A comparison between the
 9 measurements and simulation can be observed in **Fig. 3.21(b) & (c)**. The general distribution
 10 was predicted accurately via simulation; the pollutant mainly accumulated in the upper part of
 11 the square, while the concentration in the lower part was very low. The high pollutant
 12 concentration downstream of the building area was not reflected in the measurement, owing to
 13 the limited number of sensors. One evident simulation error occurred in the wake region behind

1 the buildings in the upper part of the domain. The simulation failed to predict the convection of
 2 pollutants in this area, which may lead to noticeable modeling errors during Bayesian inference.
 3 This concentration underestimation was also noted by Balczó and Lajos (2015), suggesting that
 4 the steady simulation based on Reynolds averaged Navier-Stokes equations still needs
 5 improvement to accurately predict dispersion in complex urban areas.

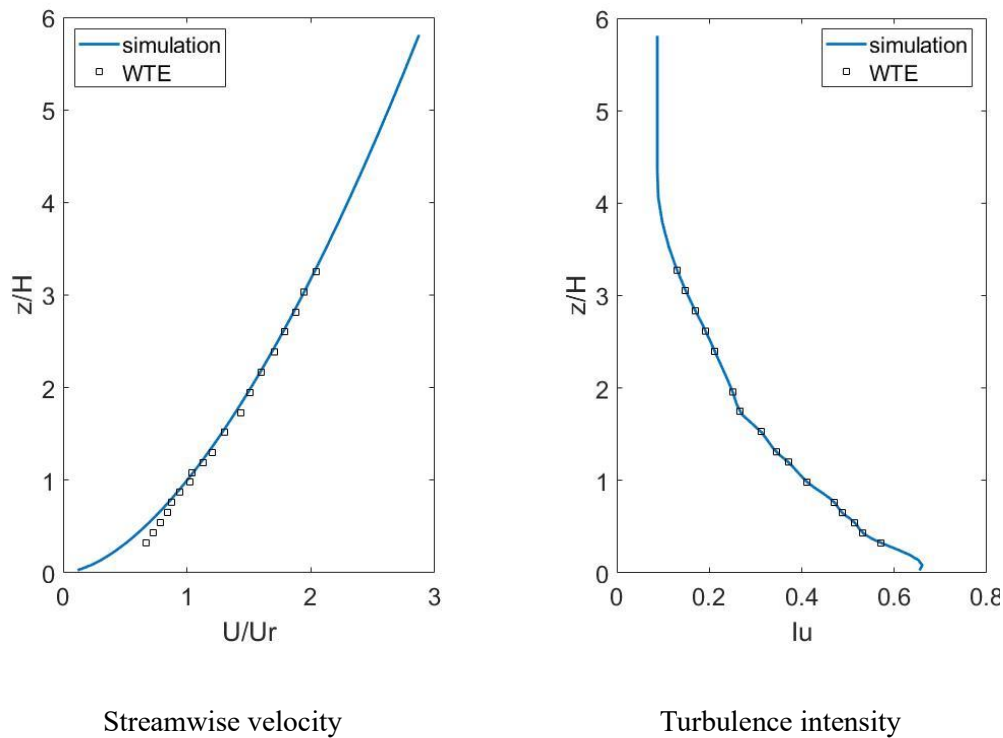


Figure 3.22. Vertical profiles of streamwise velocity and turbulence intensity for inflow in Case II

6 3.5.3 Bayesian inference settings

7 As in Keats et al. (2007b), the values of $\sigma_{d,i}^2 + \sigma_{m,i}^2$ lie between 10% and 250% of the
 8 mean sensor measurements in the wind tunnel experiment. The variance of the sensors with
 9 minimum concentrations was set as a small fixed value to maintain steady Bayesian inference.
 10 For prior information, this research applied a common assumption in STE that the source is in
 11 the main area $x[-10m, 190m] \times y[-150m, 150m]$ and will not appear in the buildings.

12

13 3.5.4 Estimation results

14 The estimation results of the proposed method for a single term are shown in Fig. 3.23. In

1 general, for each term, except length, the peak probability distribution value is close to the true
2 value. The probability distribution for coordinate x_0 of the middle point demonstrates that the
3 reconstructed source moves unsteadily along the x direction. The result for y_0 appears better
4 as the bulk of the probability mass concentrates in several grid cells around the true source. The
5 zero measurements provided by sensors No. 17 ~ 35 limit the sampling of y_0 to move
6 downward; if not, the difference in the measured and modeled value would be too large to
7 invalidate the sampling. Samplings in θ tightly surround the true value, with small
8 discrepancies. Besides, the peak value of the width is no more than 2 m away from the true
9 value. Again, the estimation is considerably accurate since the resolution of the simulation data
10 is approximately 1.67 m. The sensors around the line source (No. 3~9) help the algorithm
11 successfully estimate the angle and width.

12 There is a deviation in the estimation of source strength. It should be mentioned that the
13 scale of the horizontal axis was set small in the results of width and strength to clearly show
14 the distributions. Hence, the uncertainty seems to be amplified. The estimation for length has a
15 wide distribution range, and the estimated peak value has a relatively large bias compared with
16 the true value. Combining the fact that x_0 estimation also moves along the x -direction, it can
17 be concluded that the reconstructed source stretches and moves unsteadily in the x -direction
18 during inference to find the most probable sampling. The reason for this behavior could be the
19 lack of sensors along with the x -axis upstream and downstream of the source, which is exactly
20 the null-measurements discussed in **Section 3.5.2**. These sensors are sensitive to the change in
21 length; thus, their measurements can provide accurate information on length for the estimation
22 process.

23 Combining the estimation performance in two cases, it can be concluded that line source
24 estimation has a stricter requirement for sensor configuration than point source estimation, to
25 obtain the geometric information. Sensors near the source and sensors in the ‘region of
26 influence’ with zero measurements are all necessary. As a result, in practical applications, the
27 density of the sensors in urban areas should be high enough to capture all possible unknown
28 line sources, which could appear anywhere.

29 Admittedly, the model error of steady CFD simulation for the source-receptor
30 relationship, especially the concentration difference in the wake region behind buildings
31 (sensors No. 1~4), affects the accuracy of estimation. Hence, it is necessary to improve the

- 1 adjoint concentration simulation to ensure that the STE method is accurate enough in
- 2 complicated real applications. Even with model error, the proposed method restricts the
- 3 sampling between 0.5 and 2 times the true value, which is acceptable.

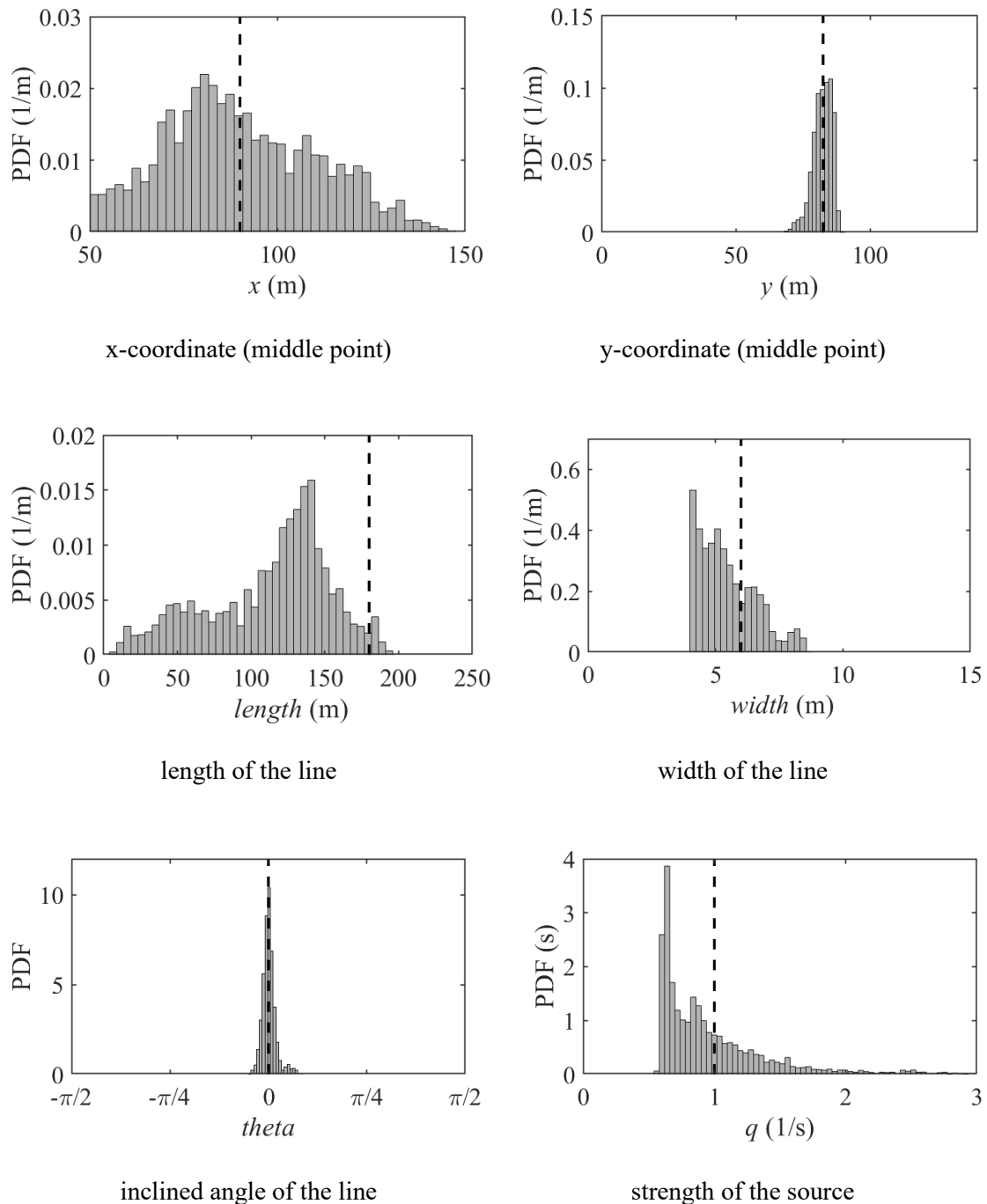
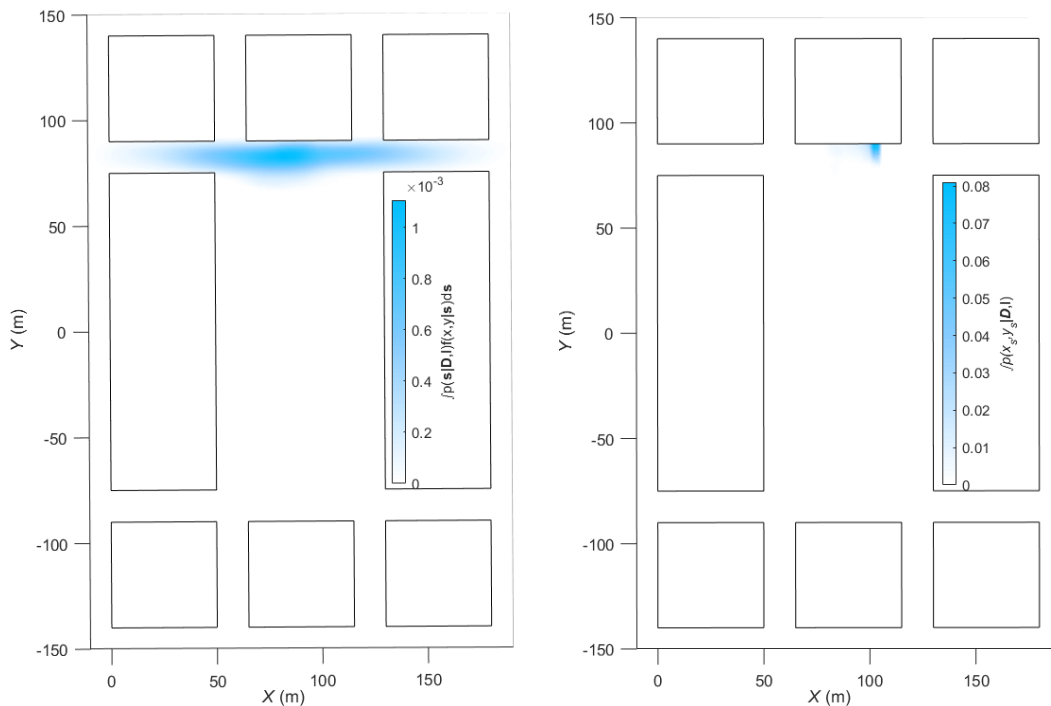


Figure 3.23. Posterior PDF of source parameters estimated by the proposed method in Case II.

(dotted line: true value)

1 The power of the proposed method can be better illustrated by the estimated line source
 2 geometry expectation $\int p(\mathbf{s}|\mathbf{D}, I) \times f(x, y|\mathbf{s})d\mathbf{s}$, as shown in **Fig. 3.24 (a)**. Although the
 3 estimations for a single term are not perfect, the joint probability distribution clearly reflects
 4 that the target is a line-shaped source, with a geometry that is very similar to the actual geometry.
 5 The proposed method constricts the distribution into an area around the true line while the
 6 possibility of other areas is eliminated by inference. Therefore, it can be confirmed that the
 7 proposed method effectively estimates the geometry information of the source in this practical
 8 case.

9



$\int p(\mathbf{s}|\mathbf{D}, I) \times f(x, y|\mathbf{s})d\mathbf{s}$ obtained from
the proposed method

$p(x_s, y_s|\mathbf{D}, I)$ obtained from conventional
method with point assumption

Figure 3.24. The estimation results of the proposed method and conventional method in Case II.

1 3.5.5 Comparison to conventional method with point assumption

2 To prove the necessity of geometry estimation in the proposed method, we also conducted
 3 STE using the conventional method introduced in **Chapter 2**, in which the source is always
 4 assumed to be an ideal point without any geometry. All the settings for Bayesian inference and
 5 measurement data are the same as those used above.

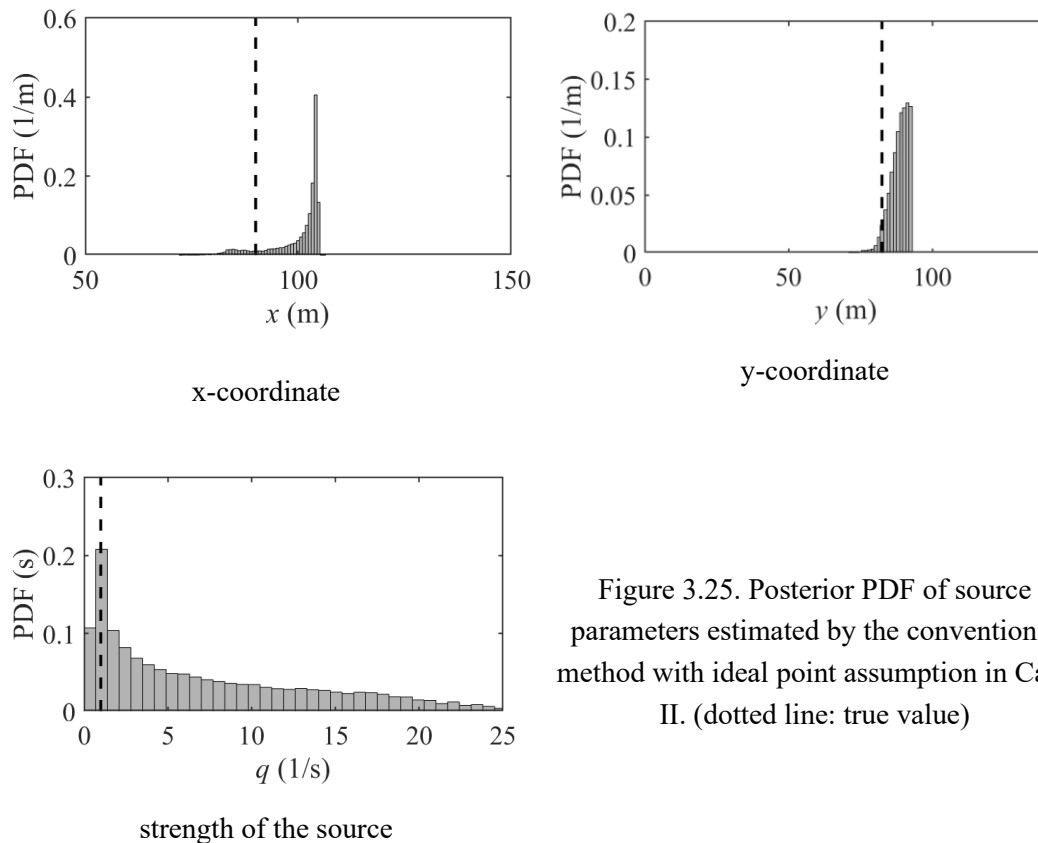


Figure 3.25. Posterior PDF of source parameters estimated by the conventional method with ideal point assumption in Case II. (dotted line: true value)

6 The results are summarized in **Fig. 3.25**. Since the conventional method does not include
 7 the geometry estimation function, it can only provide the samplings for coordinate and strength.
 8 The coordinate estimation diverges from the true value of the middle point of the line source.
 9 Although the peak value of strength PDF coincides with the true value, it is notable that the
 10 sampling range spreads up to 25 times the true value while the proposed method constrains this
 11 range by a factor of less than 2. In actual applications, this wide range may confuse the risk
 12 management process. It is also possible that the complicated dispersion in real urban may
 13 produce a wider range or a wrong peak. Again, regarding the 50th percentile and standard mean
 14 values of PDFs as the parameter estimators, the comparison between the true values,
 15 estimations of the proposed method, and estimations of conventional methods are summarized

1 in **Table 3.4**. The proposed method can not only provide valuable information about the
 2 geometry of the unknown source but also perform better than the conventional method in the
 3 aspects of locations and strength.

4 The joint posterior PDF $p(x_s, y_s | D, I)$ is shown in **Fig. 3.24 (b)**. Compared with the
 5 results of the proposed method, the entire distribution is limited to a small area. The main body
 6 of the line source was excluded in the probability distribution. The result cannot provide enough
 7 information about the location or geometry of the target source, suggesting that the STE with
 8 the conventional method fails in estimating a source that is not an ideal point.

Table 3.4. Summarized estimation results of the proposed method and conventional method

Method		x_s (m)	y_s (m)	length(m)	width(m)	angle	strength(1/s)
True value		90	82.5	140	6	0	1
Estimations (proposed method)	50 th percentile	87	83	123	5	0	1.7
	standard mean	88	82	113	5.5	0	1.1
Estimations (point assumption)	50 th percentile	103	89	/	/	/	5.7
	standard mean	100	88	/	/	/	7.5

9

10 3.6 Conclusion

11 In this Chapter, a method was proposed for the identification of a line source with the
 12 application of the super-Gaussian function. The main conclusions are as below.

13 The applicability of the method was first evaluated through Case I: a numerical experiment
 14 involving an ideal urban boundary layer. The numerical measurements were simulated using
 15 the Reynolds averaged Navier-Stokes model. According to the estimation results, all the
 16 parameters were correctly identified with almost none or acceptable errors under ideal
 17 conditions, without measurements and modeling errors. The proposed method successfully
 18 inferred that the source's shape is line-like, by automatically adjusting the coefficients in the
 19 super-Gaussian function, without any prior knowledge of the source geometry.

1 After that, the effects of different sensor configurations on the estimation results based on
2 Case I were discussed to reveal the measurement requirements of a line source. Because the
3 line source contained more geometric information (such as length, angle, and width) compared
4 with the point source, it cannot be identified correctly by the conventional sensor network, in
5 which all the sensors are placed regularly downstream. This chapter summarizes that, in order
6 to collect enough information about the line source, at least six sensors are indispensable, which
7 includes the sensors near the source and null-measurement sensors. In the real application, it
8 was noted that the density of sensors should be sufficiently high to cover the “region of
9 influence” and handle any possible unknown sources.

10 Next, to examine the robustness of the proposed method under the complicated errors
11 caused by measurements and numerical modeling, a published practical case of a simplified
12 urban square was utilized. In this case, the measurements were obtained from discrete sensors
13 in wind tunnel experiments rather than through simulation. The estimation results showed that,
14 in the practical situation, the proposed method could still estimate the line source parameters
15 efficiently, with few and simple coefficients that were tractable for Bayesian inference. More
16 importantly, the estimation was accomplished without any prior geometric information.

17 Furthermore, with the experimental measurements in Case II, the performance of the
18 proposed and conventional methods with ideal points was compared. The results implied that
19 the ideal point assumption makes the conventional method incapable of providing any
20 information about the source geometry. This simplified assumption also resulted in significant
21 errors while estimating location and strength. Consequentially, the ideal point assumption failed
22 in estimating the unknown source with geometry, which is a common scenario in real life. The
23 proposed method is an important improvement as it can estimate the geometry of the source
24 and reduce the error caused by the point assumption.

25 Admittedly, the current model may still bring a tail error into the length and width
26 estimation as the super-Gaussian function cannot be the same as a line no matter what the value
27 of λ is, though it has demonstrated that the error is acceptable.

28 It should also be noted that this is an early trail to estimate an unknown source with
29 geometry. Few related research, either the validation database or comparable algorithm, is
30 available in the literature. Most of previous research still concentrates on the point source
31 estimation. Therefore, in the following chapters, the focus will change from the geometry

-
- 1 estimation to point source estimation to utilize the existing database and benchmarks, and the
 - 2 super-Gaussian function will not be used.
 - 3

1 Symbols

- A : normalization factor of Gaussian function
- a, b, c : coefficients of Gaussian function
- C : the concentration distribution caused by a source
- C^* : adjoint concentration distribution
- D : measurements vector
- D_i : the measurement of the sensor with index i
- D_t : turbulent diffusivity
- f_x : the Gaussian function for variables x
- $f(x, y)$: the bivariate Gaussian function
- H : the building height in the real life of Case II
- h : the height of the building model in the experiment of Case II
- I : background information for Bayesian inference
- I_u : the turbulence intensity in the streamwise direction
- k : turbulent kinematic energy
- k_s : roughness coefficient for the wall function applied in Case I
- m : the dimension of the Gaussian function
- $N[a, b]$: the uniform distribution bound between a and b
- $p(A|B)$: conditional probability of event A occurring given that B is true

- Q : the source strength in the experiment of Case II
- q_s : release strength of the source
- \mathbf{R} : modeled concentration vector
- \bar{U} : the spatially averaged wind velocity of the inlet in Case I (4 m/s)
- U_r : the reference velocity measured at the building height in the experiment of Case II
- x_0 : x coordinate of the middle point of Gaussian function
- x_m : x coordinate of the sensor
- x_s : x coordinate of the middle point of line source
- y_0 : y coordinate of the middle point of Gaussian function
- y_m : y coordinate of the sensor
- y_s : y coordinate of the middle point of line source
- Σ : the covariance matrix of variable \mathbf{x} in the Gaussian function
- $\boldsymbol{\mu}$: the mean value matrix of variable \mathbf{x} in the Gaussian function
- θ : inclined angle of the Gaussian function or the line source
- $\sigma_{d,i}^2$: the variance of error in the measurement of the sensor with index i
- $\sigma_{m,i}^2$: the variance of error in the modeling concentration for the sensor with index i
- σ_X : the covariance of Gaussian function in the x direction

σ_y : the covariance of Gaussian function in the y direction

λ : power coefficient of the super-Gaussian function

ε : turbulent energy dissipation

1

2

1

2

3

4

5

6

7

8

9 Chapter 4

10 Construction of urban turbulent
11 flow database by large-eddy
12 simulation and wavelet-based
13 compression method

14

15

Abstract

According to **Chapter 2**, the prediction of the source-receptor relationship relies on the simulation of the adjoint equation, in which the hypothetical tracer emitted from each sensor is transported in the inverse spatiotemporal flow field. As shown in **Chapter 3**, an adjoint equation is a partial differential equation that is similar in form to a dispersion equation and can be simulated by computational fluid dynamics (CFD) model like Reynolds averaged Navier-Stokes (RANS). However, the prediction accuracy of RANS models for the time-averaged flow fields around buildings has been shown to be insufficient when compared with the large eddy simulation (LES) model (Tominaga et al., 2008a). Besides, the turbulence diffusion is approximated based on the mean field in the RANS model. These limitations undermine the accuracy of adjoint equation simulations and source term estimation (STE) with RANS. It is believed that LES may improve the prediction accuracy of adjoint equation simulation and enhance the reliability of statistical STE in urban applications.

Until now, the LES of the adjoint equation has been regarded as impractical because the time-series flow field data of the entire domain must be produced by forward simulation and stored in advance to realize the inverse simulation. One important challenge here is that the volume of data acquired by forward LES is too large for practical application. This research proposed to use the wavelet-based compression method to mitigate the storage pressure. The LES flow field can be compressed into a portable database for later usages like adjoint equation simulation, new dispersion simulation, and validation for new models.

Before the detailed introduction of LES for adjoint equation and STE in **Chapter 5**, this chapter constructed a compression turbulent flow database for a block-arrayed building group model simulated via LES. The objective has three aspects: (1) Evaluate the accuracy of compression and the applicability of compressed database; (2) Construct a compressed database for the adjoint equation simulation in **Chapter 5**; (3) Prepare a raw dispersion field for the STE in **Chapter 6**.

More details of the last two aspects can be found in their own chapters. This chapter focuses on the compression method and database construction. The compression performance was analyzed from two viewpoints: single snapshot and time-series data.

4.1 Introduction

In environmental flow simulations, pollutant dispersions in complicated urban areas are often simulated together with the flow fields by computational fluid dynamics (CFD). To simulate complex flow fields accurately, sophisticated numerical models, such as large eddy simulation (LES) and direct numerical simulation (DNS), have been developed, and refined grids are involved. With the help of supercomputers or clusters of workstations, most of the details needed in different scales can be simulated by proven techniques with sufficient calculation time. For example, Kikumoto and Ooka (2012) simulated air pollutant dispersion with bimolecular chemical reactions in a street canyon by LES. Coceal et al. (2007), Jacob and Sagaut (2018) explored a dynamic wind environment in a real and an idealized urban area with LES and DNS. These cases often require more than a hundred calculation hours, even if the simulations are conducted in parallel by dozens of cores. However, when new sources and questions arise, these massive simulations have to be repeated many times. A more economical way is to construct a flow field database with accurate simulation. It can be used to simulate new dispersion cases and validate new simulation methods. One example is the Johns Hopkins Turbulence Database (<http://turbulence.pha.jhu.edu>) for DNS simulations of turbulence. This database has been used, fully or partially, in several research projects, contributing to more than 40 publications (Kanov et al., 2015). For instance, Graham et al. (2016) used this database to validate a new wall model for LES.

However, another issue has been brought to the attention of researchers by this trend: the resultant data of a highly resolved flow simulation by LES or DNS has an enormous volume (gigabytes or more) per time step. The handling of such big data has challenged the limitations of the storage space and I/O speed of computers. To alleviate the pressure on storage and promote sharing, new calculation, and postprocessing of databases, an appropriate data compression method is necessary. In addition, considering that there are strong correlations in turbulent structures, temporally and spatially, the information of the flow could be represented by fewer data given that some data can be reconstructed by others with these correlations. Consequently, the compression of CFD data is theoretically realizable.

Generally, there are two kinds of compression methods: lossless and lossy. Lossless compression methods are developed for situations in which any difference between the decompressed and original data cannot be tolerated. One important application is text

1 compression, where even a different character from the original text may change the meaning
2 completely. Compression of confidential data, such as financial records, also requires strict
3 accuracy. In this case, because all the data need to be preserved rigorously, lossless compression
4 methods, such as the Lempel–Ziv–Markov chain algorithm, can reduce the file size to less than
5 20% (Hadjidoukas and Wermelinger, 2019).

6 However, data loss is allowed in lossy compression, so the compression ratio is
7 correspondingly higher. It happens that most CFD data have robustness, to some extent, against
8 proper loss. It is a common practice to save the data with lower accuracy than full precision
9 considering that numerical error makes saving with full precision meaningless. Moreover, as
10 discussed above, some data can be deleted in the compression and reconstructed later by
11 physical correlations. Hence, lossy compression methods with a high compression ratio seem
12 feasible. Still, it is worth mentioning that, since the compression is based on the correlated
13 turbulence structures, for small turbulence that is highly random with minimum correlation, the
14 effective application of compression still needs confirmation. Despite that small-scale
15 turbulence seems vulnerable to even small errors, it is reasonable to believe that the compressed
16 data is applicable when this small turbulence is not the dominant part.

17 One popular direction in the field of fluid mechanics is mode decomposition, in which the
18 nonlinear, complicated flow field is decomposed into different modes and can be approximated
19 by the reconstruction of a limited number of dominant modes. Several methods, such as proper
20 orthogonal decomposition (Berkooz et al., 1993) and dynamic mode decomposition (Schmid,
21 2010), have been invented and have achieved success in application. These methods can be
22 naturally utilized as lossy compression methods where the original time-series data are
23 represented by the dominant modes and corresponding mode-temporal coefficients, whose
24 volume is much smaller. Meanwhile, a loss is caused by neglecting the nondominant modes.
25 The main drawback of these methods is that the decomposition process always involves
26 convoluted eigen-calculation of a huge matrix. If a database involves with long time series and
27 large number of grids, the matrix will occupy a massive memory space that is difficult to
28 process for a common computer. The requirements of calculation ability and the time history
29 of the flow confine their general application as compression tools. However, its merit is the
30 convenience of decompression. Once the modes and coefficients are successfully decomposed,
31 the reconstruction only needs simple algebraic operation, which is fast.

1 In contrast, another decomposition method, wavelet decomposition (WD), can produce a
2 set of bases for the fluid with a much lighter calculation burden. Furthermore, it can be applied
3 to one single snapshot without time history. The only demand for achieving a high compression
4 ratio is that the data must change smoothly in the space (Schmalzl, 2003), which is guaranteed
5 by the diffusion term in the governing equations of the flow field. Therefore, wavelet-based
6 compression is a probable choice for CFD data. Zubair et al. (1992) introduced this idea first to
7 decompose and compress turbulent signals. Wilson (2002) applied wavelets to compress the
8 turbulent simulation data, and the compression ratio reached 256 for the maximum. Schmalzl
9 (2003) pointed out that the existing image wavelet compression method provides an available,
10 easily adopted tool for CFD data. The efficiency of different compression schemes, such as
11 JPEG and MPEG, were evaluated. Sakai et al. (2013) combined the WD with quantization and
12 entropy encoding to construct the basic structure of the wavelet-based compression method
13 (WCM) and increased the compression ratio further. They verified the WCM by simulation data
14 in a structured Cartesian mesh, which is a common situation in CFD. Following that method,
15 Kolomenskiy et al. (2018) proposed an empirical equation for error control, which makes the
16 WCM more controllable and user-friendly. They also implemented the method in different
17 simulation cases with different scales and resolutions. The wavelet-based method has shown a
18 high compression ratio and accurate decompression in general. Therefore, it is a promising
19 solution to construct a small-sized CFD database by the WCM.

20 However, to the authors' knowledge, the WCM has not been applied to construct a CFD
21 database including both the spatial distribution and time-dependent dynamics. In previous
22 research, the verification of the WCM was limited to the compression of a single snapshot of
23 the flow field. The major concern was the effects of compression of a snapshot on the
24 postprocessing or restarting simulation from that snapshot, but the cumulative effects resulting
25 from compression of time series data have not been reported. It is not clear whether the
26 dispersion simulation can proceed successfully with such a compressed database. Moreover,
27 studies from the perspective of the limitation of the compression ratio and source of errors have
28 been few. There is little research in which the feasibility of compression in complicated flow
29 field simulation, such as building clusters in urban areas, has been addressed. These aspects
30 should be investigated further before the construction of the CFD database with WCM and
31 applying it for inverse simulation later.

32 In this chapter, a small urban flow database was constructed by using the WCM, and the

1 applicability was verified. The original data for the database were obtained from an LES
 2 simulation of a block-array urban area. The LES data were compressed with different specified
 3 errors by the WCM to construct the database. The relationship between compression error and
 4 the compression ratio, as well as the limitation of the compression, are demonstrated.
 5 Furthermore, the effects of compression on a single snapshot and time-series data are discussed.
 6 To confirm the usability of the compressed database, it was used to re-simulate the dispersion
 7 of a passive scalar, which is compared with the dispersion field of the original LES and wind
 8 tunnel experiment.

9

10 4.2 Compression methodology

11 First of all, the compression method is introduced briefly. The basic idea of using WD to
 12 realize data compression for the CFD simulation was proposed by (Zubair et al., 1992) for
 13 turbulent signals. Subsequent research (Sakai et al., 2013; Schmalzl, 2003; Wilson, 2002)
 14 enriched the idea by adjusting the algorithm for general cases and combining it with other data
 15 compression concepts. Kolomenskiy et al. (2018) proposed an empirical formula for error
 16 control of the WCM and verified its accuracy in several possible application scenarios.

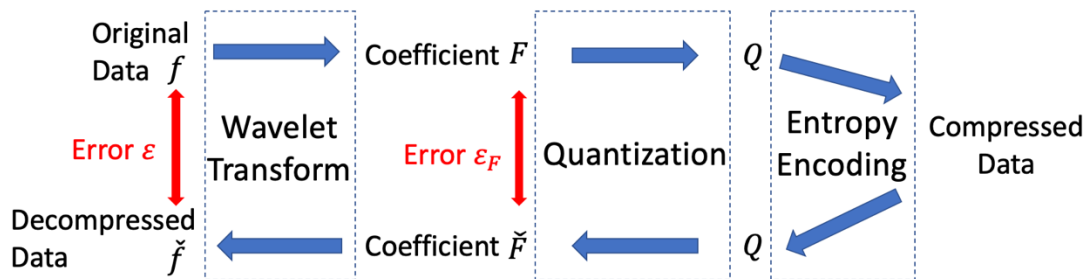


Figure 4.1. The structure of the compression method

17 Here, the method proposed by Kolomenskiy et al. (2018) was applied exactly because the
 18 structure of this method is relatively complete and, more importantly, the error control is
 19 explicit. There are three steps in this compression method: WD, quantization, and entropy
 20 encoding (**Fig. 4.1**). Decomposition can be simply realized through the inverse of the process.
 21 In the following, these three steps are introduced. The original data are assumed to constitute a
 22 three-dimensional scalar field obtained from the Cartesian indexing mesh grid $\{f_{i,j,k}\}$. The

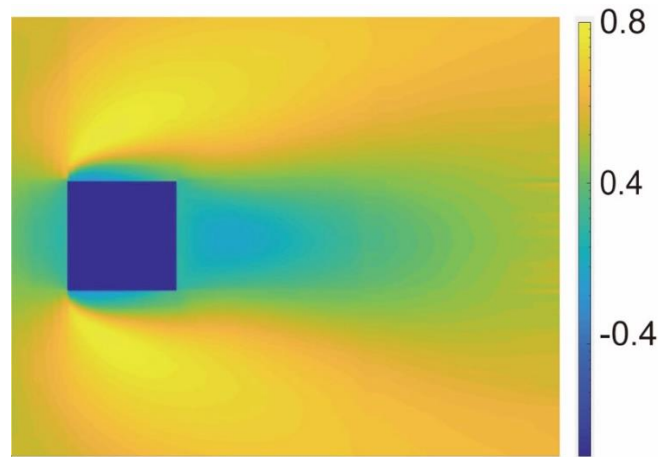
1 index i, j, k indicates the coordinate of a grid in the x, y, z directions.

2

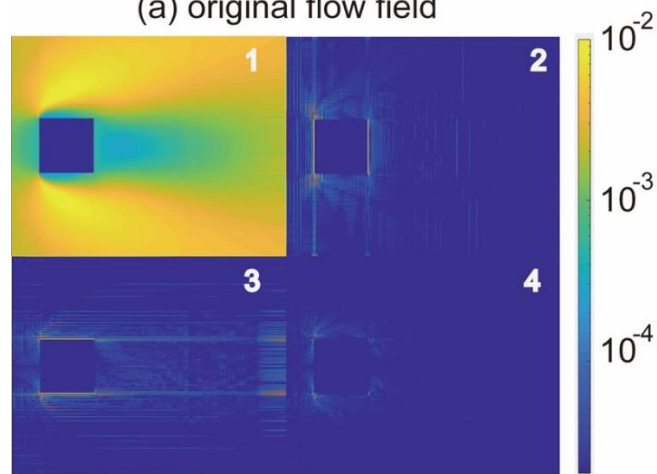
3 **4.2.1 Wavelet decomposition**

4 WD is popular in the engineering field and has shown its power in signal processing (Rioul
5 and Vetterli, 1991), noise control (Qiu et al., 2016), and image compression (Usevitch, 2001).
6 The essence of the WD is to decompose the original data into the approximation part with low
7 frequency and the detailed part with high frequency. On the one hand, for the flow field, most
8 of the energy is contained in the low-frequency structures, which means the decomposition
9 coefficients are large for the approximation part and small for the detailed part. On the other
10 hand, the approximation part only occupies part of the original saving space. Because the WD
11 can be implemented on the decomposed one repeatedly, the resulting approximation part uses
12 quite a small space to represent most of the information.

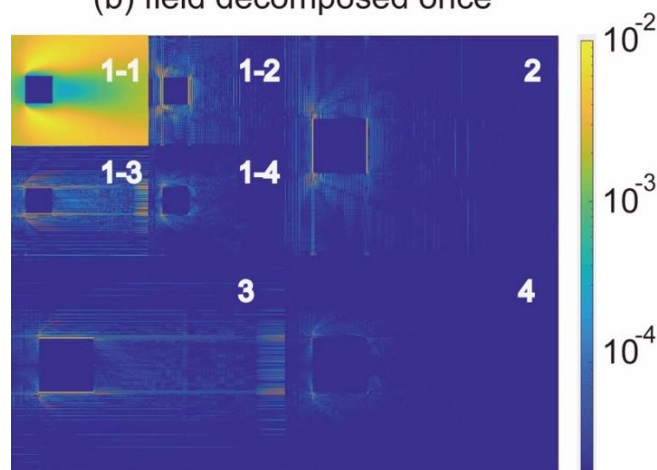
13 A two-dimensional WD was used as an example. **Fig. 4.2a** is an averaged streamwise
14 velocity field around an isolated building on a horizontal plane. This flow field was decomposed
15 once to yield **Fig. 4.2b** using the Cohen–Daubechies–Feauveau 9/7 (CDF9/7) (Daubechies and
16 Sweldens, 1998) wavelet. It indicates that the decomposed flow field consists of four parts,
17 which correspond to the low-low frequency part (1), high-low frequency part (2), low-high
18 frequency part (3), and high-high frequency part (4), respectively, transformed by the wavelet
19 in the horizontal and vertical directions. In the decomposition, wavelets work as a filter to
20 convert the flow field to the wavelet coefficients. Large filters can extract large structures and
21 small filters can extract small structures. The filtering process is conducted in each axis
22 separately. In this 2-dimensional example, the horizontal decomposition puts the coefficients
23 obtained from the large filters to the left and coefficients of small filters to the right. Similarly,
24 the vertical decomposition puts the coefficients of large filters to the up and coefficients of
25 small filters below. Therefore, the coefficients at the top-left are the ones corresponding to the
26 large filters in both horizontal and vertical directions, representing the large-scale structures
27 with low frequencies in both directions, and labeled with “low-low”. The other parts are labeled
28 in the same way.



(a) original flow field



(b) field decomposed once



(c) field decomposed twice

Figure 4.2. A 2-dimensional example of wavelet decomposition of flow field: (a) is the original mean streamwise velocity field around a building; (b) is the magnitude of wavelet coefficients obtained from one-time decomposition of (a); (c) is the magnitude of wavelet coefficients obtained from two-times decomposition of (a). (parts 1 and 1-1 are the approximation parts and use the same color bar in (a), and others are detailed parts using the logarithm color bar to show the coefficients clearly)

1 Among these four parts, only the coefficients in part (1) have the same order of magnitude
2 as the original field, which means this is the approximation part with the most energy. In
3 contrast with part (1), the order of coefficients in the other three detailed parts is extremely
4 small. These coefficients can be approximated to reduce the volume of the data considerably if
5 the minor energy is not in the scope of interest. Another advantage of WD is that it can be
6 repeated several times on the decomposed field. **Fig. 4.2c** shows the further decomposed field
7 from **Fig. 4.2b**. In the second transform, parts (2)–(4) are kept unchanged while part (1) is
8 further decomposed into four smaller parts. Similar to the first transform, these four parts
9 represent different frequencies. Part (1-1) is the approximation part with major energy, while
10 parts (1-2), (1-3), and (1-4) are the detailed parts with a small order of magnitude. After WD is
11 done several times, the approximation part containing the dominant information occupies a
12 limited storage volume in the upper left corner, while the rest of the detailed parts can be
13 approximated further.

14
15 In this research, a three-dimensional CDF9/7 WD was performed four times on the original
16 data. It was confirmed that CDF9/7 can reach a higher compression ratio than other low-order
17 wavelets, such as Haar, Daubechies, and Symlets (Kolomenskiy et al., 2018). The basic wavelet
18 like Haar is the low-order orthogonal type, i.e. the basis functions of wavelet are orthogonal to
19 each other and have zero inner product. Although it is simple to use, it may produce too many
20 wavelet coefficients, which are unfavorable to the compression. Therefore, the image
21 processing community suggests the use of biorthogonal wavelets, which can transform the
22 image into fewer coefficients because there is more flexibility to design and optimize the basis
23 functions (Li, 2015). Among the biorthogonal wavelets, CDF 9/7 is confirmed to be the best
24 one for image processing (Villasenor et al., 1995) and employed as the default wavelet in the
25 JPEG2000 (Skodras et al., 2001). However, it has also been pointed out that the input
26 coefficients of CDF 9/7 need more hardware resources and processing time because they are
27 irrational and need high precision representation in the computer (Martina and Masera, 2005;
28 Naik and Holambe, 2014). Several methods have been proposed by them to improve this point
29 and their effects on turbulence compression still need further tests. Meanwhile, according to
30 the author's knowledge, research about the relationship between turbulence and wavelet
31 compression is very limited. Systematic research concerning the best wavelets for turbulence
32 decomposition is still in need. Therefore, in the current research, the most popular wavelet CDF

1 9/7 is utilized. The wavelet coefficients $\{F_{i,j,k}\}$ are processed further in the next two steps.

2

3 4.2.2 Quantization

4 Considering that most of the wavelet coefficients $\{F_{i,j,k}\}$ have a small order of magnitude,
 5 one brute-force compression algorithm changes all the detailed parts with 0 and only keeps the
 6 approximation part. Although this algorithm can ensure a high compression ratio, the resultant
 7 accuracy may be too poor because there are still relatively large coefficients in the detailed parts
 8 that are removed directly. Moreover, there is no selection leeway for users to control the error
 9 margin of the compression. Only one choice remains because all the details are abandoned. A
 10 better algorithm is quantization, in which the user can specify how much of the detail should
 11 be removed to satisfy the application requirements and error control.

12 In the computer system, eight 1-byte memory spaces are commonly used to represent the
 13 double-precision floating number $\{F_{i,j,k}\}$ in a lossless manner, in which one bit represents the
 14 sign, 11 bits represent exponent and 52 bits represent the fraction. To achieve a high
 15 compression ratio, one must quantize $\{F_{i,j,k}\}$ with fewer 1-byte memory spaces and accept the
 16 inevitable truncation errors ε_F . In the quantization system, the bits do not function in the
 17 conventional way but only create positions to quantize $\{F_{i,j,k}\}$. The algorithm will scan all
 18 $\{F_{i,j,k}\}$ to pick up the maximum and minimum values, and assign the minimum value to the
 19 first bit array and maximum value to the last bit array. The positions in the middle represent the
 20 values with an even difference between the maximum and minimum. In this case, if one
 21 allocates N 1-byte spaces for the quantization, there are $2^{8N} = 256^N$ positions, and the data
 22 $\{F_{i,j,k}\}$ should be assigned to the corresponding positions. The true value is replaced by the
 23 rounded one nearest it, and the differences lying in the gap of the truncation error $\varepsilon_F =$
 24 $\frac{\max(F) - \min(F)}{256^N}$ are removed. In this case, the algorithm can control the width of the gap ε_F to
 25 decide how much of the detailed part in the wavelet coefficients is contaminated. Furthermore,
 26 the value with a large order of magnitude is added to the truncation error instead of being
 27 dumped directly. In the algorithm, the number of 1-byte memory spaces N was calculated by
 28 the truncation error ε_F .

$$N = \log_{256} \left(\frac{\max(F) - \min(F)}{\varepsilon_F} \right) \quad (4.1)$$

1 Meanwhile, it is obvious that

$$\varepsilon_F \approx \max|F - \check{F}| \quad (4.2)$$

2 where \check{F} is the quantized wavelet coefficients.

3 However, this truncation error is not directly correlated with the original data $\{f_{i,j,k}\}$ but
 4 rather with the wavelet coefficient $\{F_{i,j,k}\}$. It is necessary to figure out the effects of the
 5 truncation error on the original data to control the compression process effectively.
 6 Kolomenskiy et al. (2018) proposed that ε_F and the user-specified error ε for the original
 7 data have following relationship:

$$\varepsilon_F = \frac{\varepsilon \times \max|f|}{\eta} = \frac{\max|f - \check{f}|}{\eta} \quad (4.3)$$

8 Here, η is an empirical coefficient. Hence, in the compression process, the user can specify ε
 9 according to the requirement. The algorithm calculates the truncation error ε_F and necessary
 10 memory spaces N through Eqs. (4.1) and (4.3) for the quantization, which can ensure that the
 11 real compression error metric $\varepsilon_r = \frac{\max|f - \check{f}|}{\max|f|}$ is almost the same as ε . Moreover, η is
 12 proposed to be approximately 1.75 to realize the relationship according to several numerical
 13 cases. This value was examined also in this study.

14 According to the definition of ε_r and $\varepsilon \approx \varepsilon_r$, if ε is larger than 1, the largest
 15 compression error would be even larger than the largest flow field data, which makes the
 16 compression data too inaccurate. Meanwhile, if ε is too small, the compression ratio is so
 17 small that the compression is meaningless. Therefore, ε should be appropriately specified
 18 between 0 and 1. The attained coefficients \check{F} are input in the entropy encoding.

19

20 4.2.3 Entropy encoding

21 Entropy encoding is a lossless compression technique to reduce the storage volume by
 22 representing subsequences of data with different symbols according to the times of their

1 appearance (Song, 2008). Although quantization does compress the data significantly, there is
2 still redundancy in the data because the output of quantization does not occur with equal
3 probability. This redundancy can be optimally removed by entropy encoding (Ruttimann and
4 Pipberger, 1979). Specifically, in the algorithm, the appearance times of each value in the input
5 data arrays are counted first. Then, values with large appearance times are represented by short
6 symbols, while ones with small times are represented by long symbols, which reduces the
7 volume of the input data further. In this research, the entropy encoding method called “range
8 coding” developed by Martin (1979) was employed. The quantized wavelet coefficients are
9 represented by different symbols and saved eventually.

10

11 **4.3 Case description**

12 In this part, the WCM is used to construct a compression database based on LES raw data.
13 The details of the study case are as below.

14

15 **4.3.1 Production of raw data with large-eddy simulation**

16 The time series raw data for the database was first produced by CFD. To evaluate the
17 compression ability of the WCM, the simulation case should have complicated flow fields and
18 turbulence structures. Meanwhile, as a database for later source identification, the case should
19 represent the common situation for dispersion in the urban area. In recent years, as a
20 theoretically representative form of a real urban area, the block-arrayed idealized urban model
21 has been widely used to investigate the wind field (Abd Razak et al., 2013; Ikegaya et al., 2017),
22 pollutant dispersion (Branford et al., 2011b; Coceal et al., 2014b), and thermal effects (Uehara
23 et al., 2000). Hence, the simulation case was set as regular building arrays.

24 LES was applied to simulate the turbulence flow because it can resolve most of the details
25 of the flow. Its reliability and accuracy of simulating both the time-dependent dynamics and
26 spatial distribution in wind engineering have been confirmed by multiple research (Blocken,
27 2018). The flow was simulated with a finite volume method implemented in OpenFOAM v1906.
28 The standard Smagorinsky model was adjusted with the Smagorinsky constant $C_s = 0.12$,
29 which has been confirmed in previous work (Tominaga et al., 2008a) for the flow past buildings.

1 The geometry of the calculation domain is shown in **Fig. 4.3**. The settings are the same as
2 the wind tunnel experiment conducted at Tokyo Polytechnic University. The experimental data
3 are open in the AIJ database (https://www.aij.or.jp/jpn/publish/cfdguide/index_e.htm). All
4 blocks are the same cubes with edges $H = 60$ mm. They were arranged uniformly with a
5 distance H between each other. The domain size was set as $16H(x) \times 14H(y) \times 4H(z)$. For
6 the simulation of block-arrayed urban flow, Coceal et al. (2006) has explored the influence of
7 domain height on the simulation with $4H$, $6H$ and $8H$. It was found that although $4H$ is too
8 small to capture the largest turbulence structures, like the long streaky one, the differences in
9 mean and turbulence statistics were negligible except in the vicinity of the top. Considering the
10 calculation burden and the fact that flow around blocks is the main target, several research (Abd
11 Razak et al., 2013; Claus et al., 2012; Ikegaya et al., 2017; Xie and Castro, 2006) applied $4H$
12 as domain height. Therefore, we set the same height following these research. Uniform
13 structured meshes with a medium size ($H/20$) were applied. In previous research (Ikegaya et
14 al., 2019; Xie and Castro, 2006), it was reported that the medium-sized mesh is sufficiently
15 reliable to reproduce the mean flow distribution, turbulent kinematic energy, and large-scale
16 flow dynamics. The total number of mesh cells was approximately 6.9 million. The bottom side
17 of the domain and the surfaces of all blocks were defined as a nonslip boundary with a Spalding
18 wall function. At the top of the domain, a free slip condition was imposed for velocity, and the
19 zero-gradient Neumann condition was imposed for pressure. The four sides of the domain were
20 all periodic boundaries. The turbulent inflow was generated by coupled inlet and outlet and was
21 driven by a pressure gradient, which was adjusted at each step to ensure a fixed flux through
22 the outlet.

23 The free stream speed u_r was measured at $(x, y, z) = (0, 7H, 3.33H)$ in both simulation
24 and the experiment. The inflow moves along the positive x axis. The Reynolds number based
25 on u_r and H is 1.68×10^4 . The friction speed u_* was calculated by $0.07u_r$, which was
26 proposed in previous work (Cheng and Castro, 2002). Approximately a $200T$ ($T = H/u_*$)
27 calculation time was spent for initialization to ensure that the simulation converges to a
28 statistically steady state. Then, the time series data of $240T$ at all grid points were compressed
29 and saved to construct the database. Each time step was 0.00075 s (approximately $T/240$).

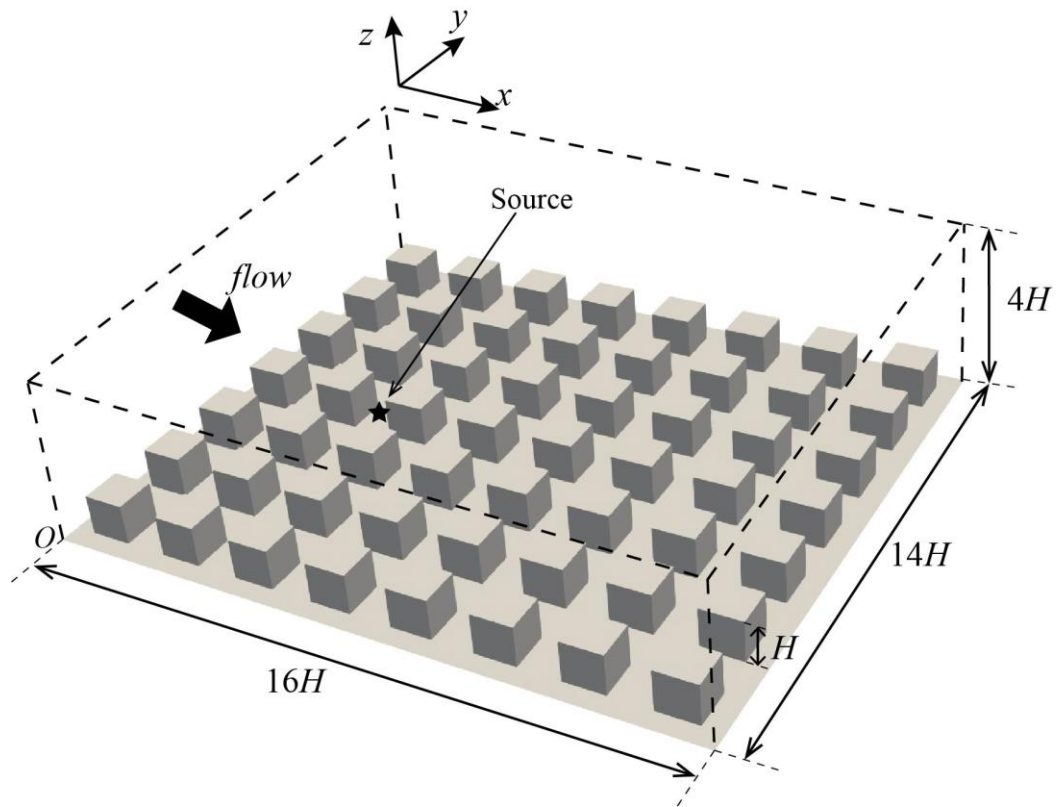


Figure 4.3. Schematic of simulation domain and the block-arrayed model.

1 4.3.2 Method to verify the compression database

2 Commonly, databases should satisfy usages in two time scales. The single snapshot is used
 3 to visualize the instantaneous flow fields and turbulent structures. The time-series data can
 4 provide statistical information and be applied later for inverse simulation or dispersion
 5 simulation of new sources. The compression should not affect these applications to protect the
 6 quality of databases.

7 Correspondingly, in this research, the instant effects on the single snapshot and cumulative
 8 effects on the time-series data of the compression process are confirmed. The spatial
 9 distribution of the flow field in one snapshot and statistics at several points in the compression
 10 database are analyzed in detail. To confirm that whether the compression database can be used
 11 for the dispersion simulation, the dispersion of the passive scalar emitted by the same source
 12 was simulated with both compressed databases and raw LES data. The original concentration
 13 distribution field of passive scalar was simulated online together with the raw LES while the
 14 concentration field of the same source transported by the compressed flow field was simulated

1 offline after the compression databases were constructed. The source was arranged at (x, y, z)
 2 $= (4H, 7H, 0)$ with the same injection speed in the wind tunnel experiment. The simulated
 3 concentrations in different cases are compared.

4 The governing equation for the transportation of the passive scalar in the LES can be
 5 expressed as

$$\frac{\partial C}{\partial t} + \frac{\partial UC}{\partial x} = S + \frac{\partial}{\partial x} \left(D_e \frac{\partial C}{\partial x} \right) \quad (4.4)$$

6 where C is the concentration in the grid scale, \mathbf{U} is the velocity in the grid scale, S is the
 7 source term, and D_e is the effective diffusion coefficient including the molecular diffusivity
 8 D_m and the sub-grid scale (SGS) turbulent diffusivity D_{sgs} . Here, D_m is a parameter
 9 dependent on temperature and does not change in this simulation. Its value was given in
 10 **Chapter 2**. D_{sgs} is modeled as $D_{sgs} = \nu_{sgs}/Sc_{sgs}$. Here, the SGS turbulent Schmidt number
 11 Sc_{sgs} is assigned as 0.7, and ν_{sgs} is the eddy viscosity coefficient at the SGS. During the
 12 simulation, \mathbf{U} and D_{sgs} are time-dependent. The database compressed these two fields at
 13 each time step with different error control ε : 10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} , and 10^0 .
 14 Then, the dispersion fields were repeatedly simulated with the decompressed $\tilde{\mathbf{U}}$ and $\widetilde{D_{sgs}}$
 15 data. The calculations restarted from the same time step in which the simulation reached
 16 statistical steadiness. A snapshot of each component of velocity field or effective diffusion
 17 coefficient field is about 100MB. The original data size of $240T$ sampling time is about 24 TB.

18 In this research, the implementation of WCM can be divided into the following procedures:

19

- 20 ● LES: the raw data of \mathbf{U} and D_{sgs} were produced at each timestep using OpenFOAM,
 21 and the original concentration field of passive scalar was calculated simultaneously.
- 22 ● Compression: At each timestep, the raw data of \mathbf{U} and D_{sgs} were inputted into
 23 WCM to be compressed and stored to construct the database.
- 24 ● Decompression: The database was decompressed by WCM into $\tilde{\mathbf{U}}$ and $\widetilde{D_{sgs}}$ data for
 25 each timestep.
- 26 ● Dispersion analysis: OpenFOAM read $\tilde{\mathbf{U}}$ and $\widetilde{D_{sgs}}$ to calculate the concentration
 27 field of passive scalar for each case.
- 28 ● Post-processing and Comparison.

29

1 The compression source code is exactly the same one which is opened in Kolomenskiy et
2 al. (2018). A script was made by the author to enable it to process the output of OpenFOAM in
3 each processor in parallel automatically.

4

5 **4.4 Results and discussion**

6 The quality of the compressed database was examined by analyzing the effects of
7 compression on the spatial distribution of the flow field in one snapshot and time-series
8 statistics at several points. The compression ability of WCM was also examined in terms of the
9 compression ratio and compression limitation. All the velocity results were nondimensionalized
10 by u_r , and the concentration results were nondimensionalized by the standard concentration
11 $C_r = C_{\text{gas}}q/(u_r h_r^2)$. Here, C_{gas} is the emission strength of the source in the simulation and
12 concentration of the source in the experiment, q is the gas flow rate, and h_r is the reference
13 height $3.33H$.

14

15 **4.4.1 Validation of large-eddy simulation**

16 First of all, the raw data of LES results are validated in **Fig. 4.4 & 4.5**. Because the flow
17 field is regular in the spanwise direction, we validated two rows in the streamwise direction: the
18 profiles in the wake region of buildings ($y=7H$) and profiles in the open street region ($y=6H$).
19 Samplings for $240T$ are used for the validation. It can be confirmed that the simulated flow field
20 and concentration agree well with the experimental measurements. The mean, standard deviation of
21 streamwise velocity, and the turbulent kinetic energy agreed well in most places.

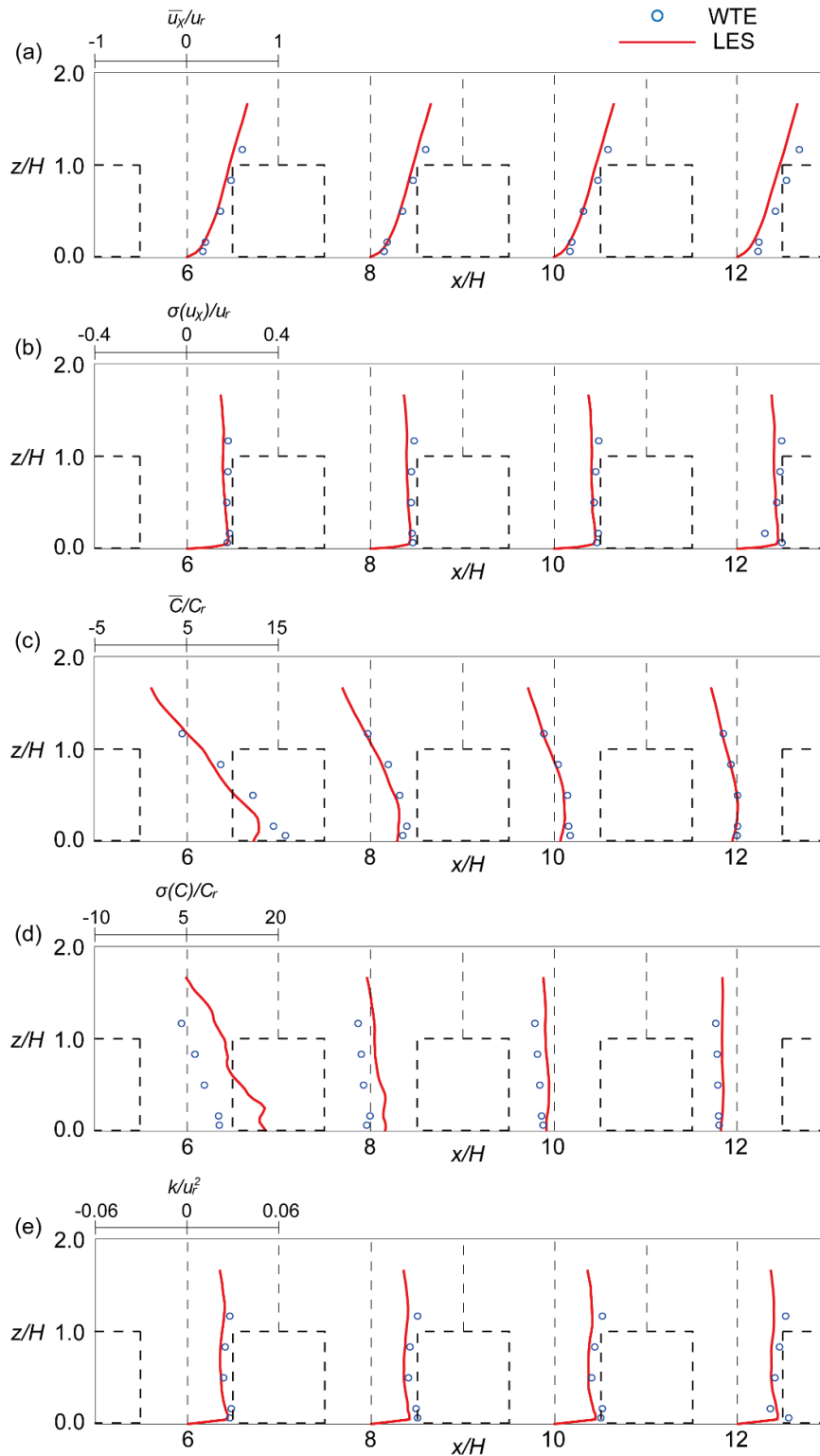


Figure 4.4. Validation of the velocity and concentration fields predicted by LES based on the experimental measurements in the open street region ($y = 6H$). (a) Time-averaged streamwise velocity; (b) standard deviation of streamwise velocity; (c) time-averaged concentration; (d) standard deviation of concentration; and (e) turbulent kinetic energy.

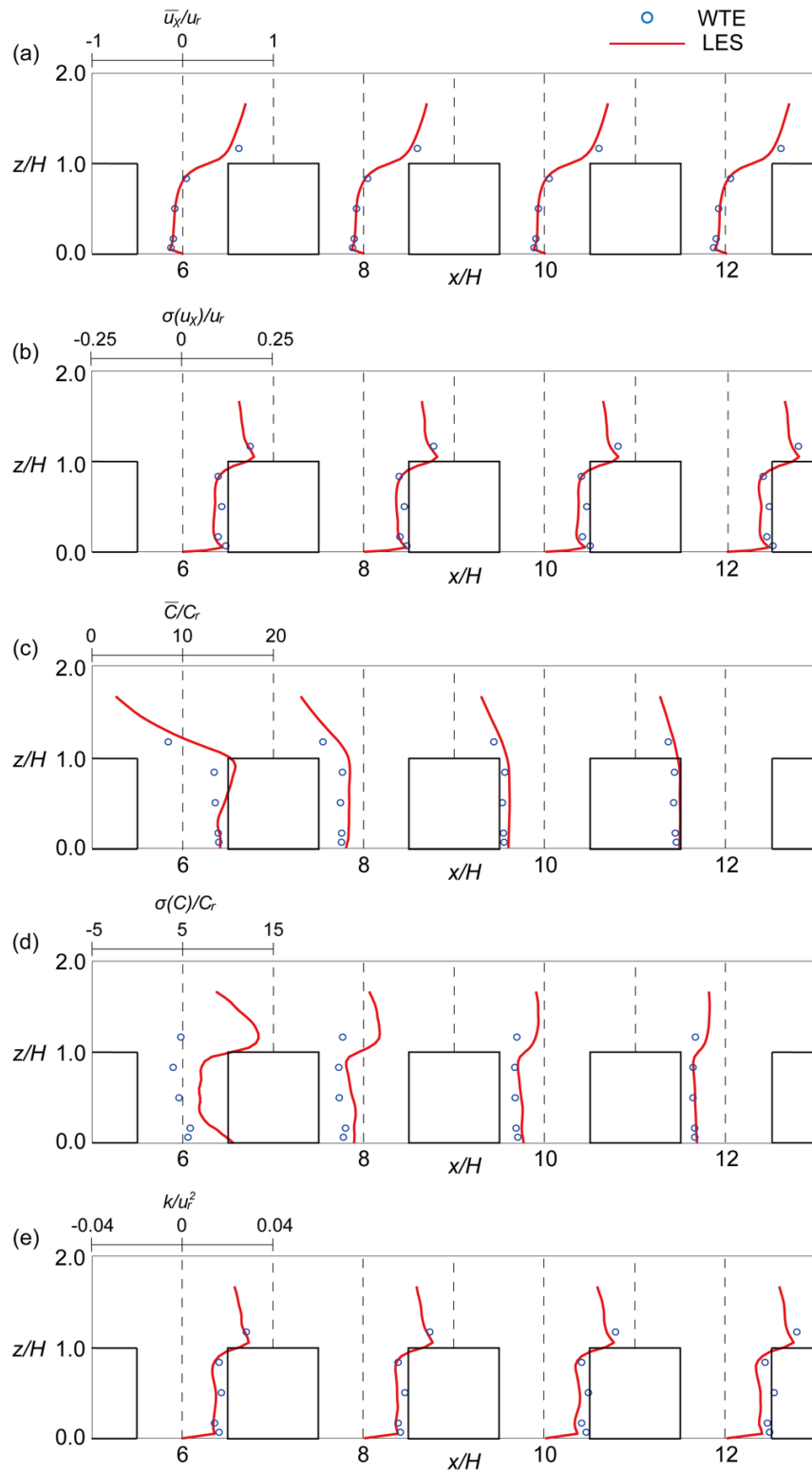


Figure 4.5. Validation of the velocity and concentration fields predicted by LES based on the experimental measurements in the wake region ($y = 7H$). (a) Time-averaged streamwise velocity; (b) standard deviation of streamwise velocity; (c) time-averaged concentration; (d) standard deviation of concentration; and (e) turbulent kinetic energy.

1 4.4.2 Compression ability

2 The compression ability of the WCM was explored. All the following discussions in this
 3 section are based on the compression result of the effective diffusion coefficient D_e in a single
 4 snapshot. The velocity field compression has similar results and is discussed later.

5 First, it is necessary to determine whether $\eta = 1.75$ in Eq. (4.3), proposed by
 6 Kolomenskiy et al. (2018), is an appropriate empirical coefficient for the error control. It is
 7 indicated in **Fig. 4.6** that the prescribed error control ε and resultant real error $\varepsilon_r = \frac{\max|f-\tilde{f}|}{\max|f|}$
 8 are almost the same in the compression method using $\eta = 1.75$. Only a small discrepancy
 9 appears after ε is larger than 10^{-1} . The reason why the error control fails afterward is the
 10 limitation of the compression algorithm, which is addressed specifically later in this section. As
 11 a consequence, the error control with the empirical coefficient is accurate for complicated urban
 12 flow.

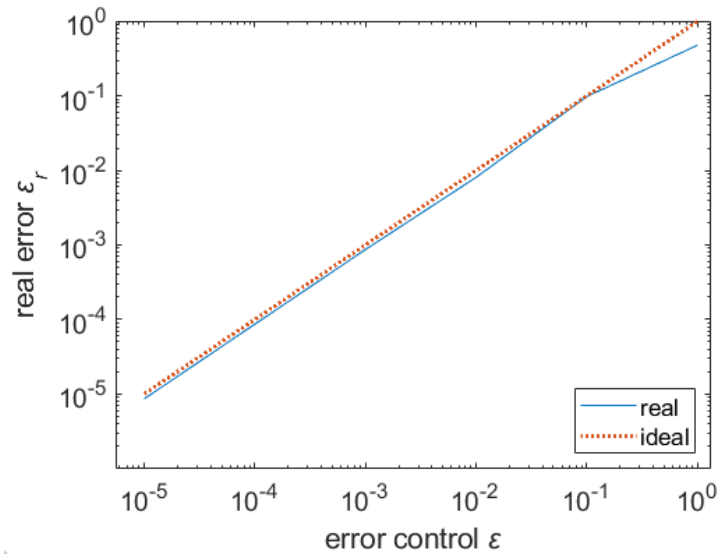


Figure 4.6. Relationship between the error control ε and real error ε_r : the red dotted line is the ideal one where $\varepsilon = \varepsilon_r$, and the blue line is the real situation (compression of effective diffusion coefficients)

13 **Fig. 4.7** shows the relationship between the compression ratio $r = V_o/V_c$ (where V_o is
 14 the storage volume of the original data, and V_c is the storage volume of the compressed data)
 15 and the error control ε . When ε is set small, r increases with ε exponentially. After ε

1 reaches approximately $10^{-0.8}$, the increasing tendency starts to change. Large oscillations of
 2 the compression ratio appear in the red area, which indicates the compression limitation of the
 3 WCM. It is shown that r reaches its largest value when ε is larger than 1, where the largest
 4 compression error is almost the same as the largest value in the original field. The dotted line
 5 denotes the compression ratio realized by a conventional compression method, in which the
 6 flow data was represented by single-precision numbers and then compressed with the standard
 7 ZIP algorithm. It can be confirmed that the current method outperforms the conventional one
 8 in the compression ratio and provides error control to users.

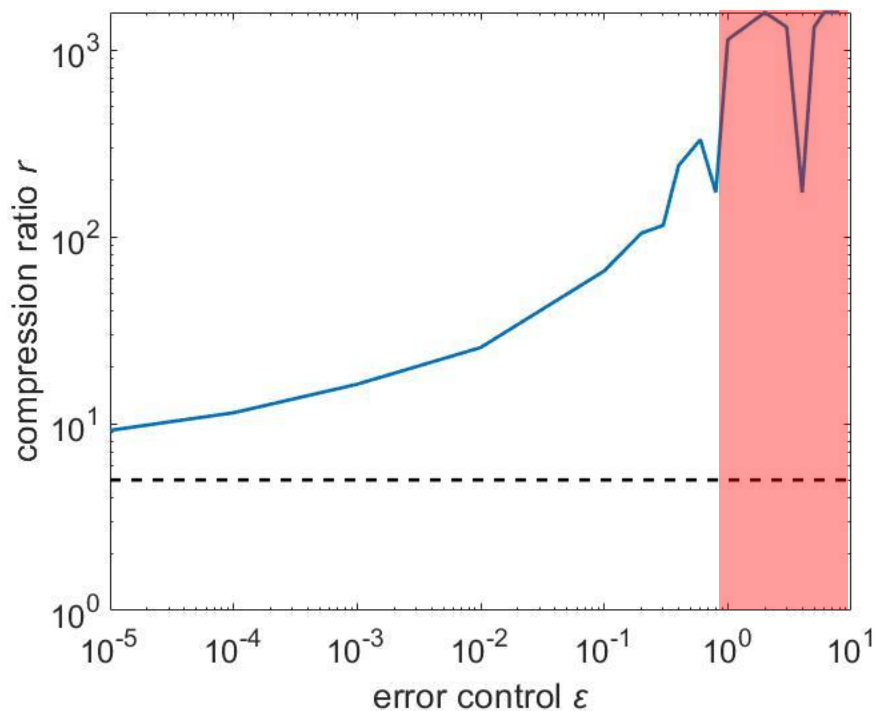


Figure 4.7. Relationship between the error control ε and the compression ratio r (compression of effective diffusion coefficients). The dotted line marked the compression ratio realized by a conventional method: single-precision representation + standard ZIP.

9 The reason for the appearance of the limitation was studied here. Because the compression
 10 method in this research retains all the wavelet coefficients for the quantization and the entropy
 11 encoding causes no error, the main error results from the truncation in the quantization process.
 12 **Fig. 4.8** shows the probability density function of the wavelet coefficients $\{F_{i,j,k}\}$ of the D_e
 13 field in one snapshot compression. Although the coefficients ranged to more than 3×10^{-3} ,
 14 most of the information lies in the area between -1×10^{-4} and 1×10^{-4} . Because the

1 coefficient difference that is smaller than the truncation error ε_F is discarded in the
 2 quantization process, if ε_F is larger than 2×10^{-4} in this case, the main information is lost,
 3 and the whole structure of the wavelet coefficient is destroyed. In Eq. (4.3), if ε_F equals
 4 2×10^{-4} and the true value of $\max|f|$ is substituted, the limitation of the user-specified error
 5 ε can be calculated as

$$\varepsilon \approx \frac{0.0002 \times 1.75}{0.00029} = 1.2 \quad (4.5)$$

6 This value is close to the one where the compression ratio confronted the limitation.
 7 Consequently, the main reason for the compression limitation is that the truncation error is too
 8 large so that most wavelet coefficients are removed by the quantization. Although the error
 9 control ε that reaches the limitation of the compression method can change for each case, it is
 10 approximately 1 in the condition of this study. There is no universal value suitable for all cases,
 11 but the limitation can be calculated by the range of the wavelet coefficient to ensure that the
 12 core information is not damaged.

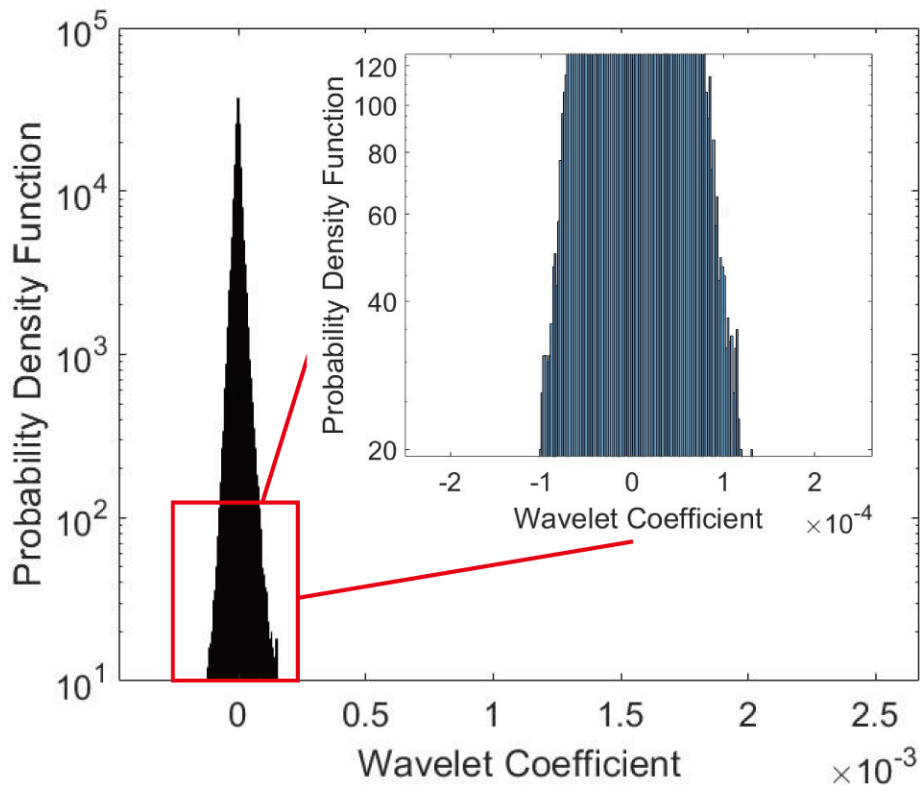


Figure 4.8. Distribution of wavelet coefficients of effective diffusion coefficients in one snapshot

4.4.3 Compression effects on a single snapshot

In this section, the focus is on the compression effects on a single snapshot, which is important for the postprocessing and flow visualization with the database. Because the WCM compressed each timestep independently, a snapshot was arbitrarily selected and analyzed.

Fig. 4.9 compares the instantaneous streamwise velocity field U_x in the vertical plane in the middle of the y direction between the original LES and compressed databases. The difference between these fields can be checked in **Fig. 4.10**. When the error control ε is smaller than 10^{-1} , the compression effects are so small that the visualization cannot identify the difference. When ε reaches 1 (10^0), the situation changes so that discrete large values appear and mottle the velocity field.

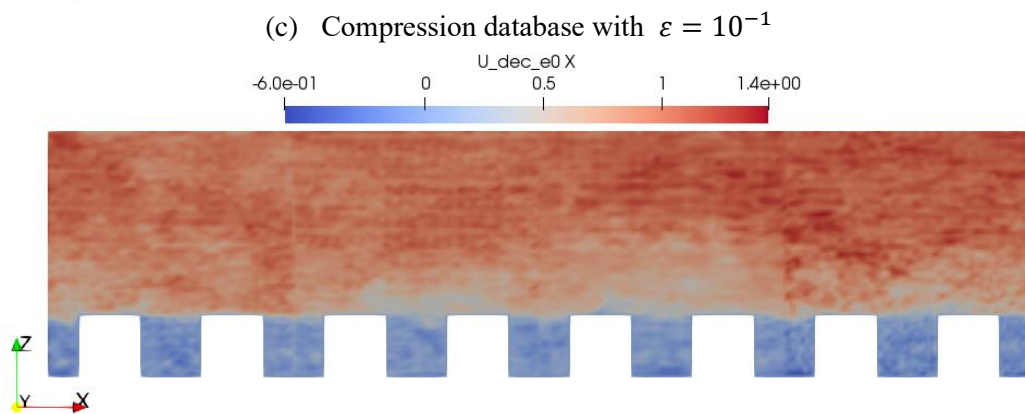
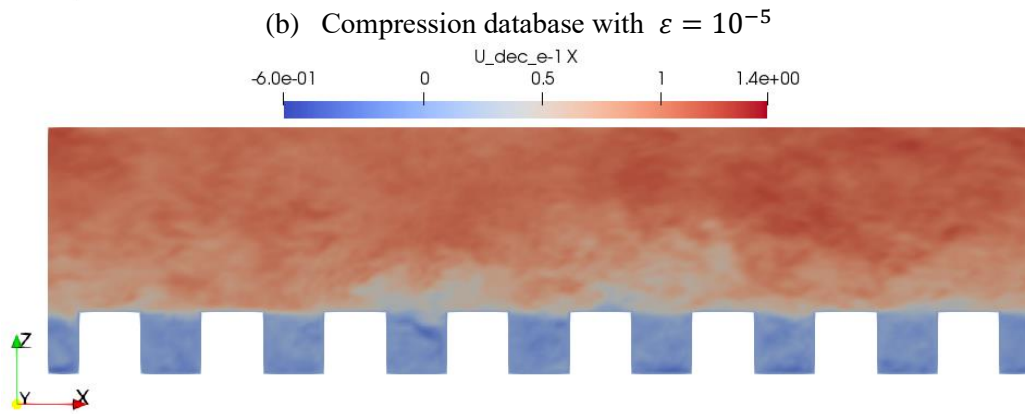
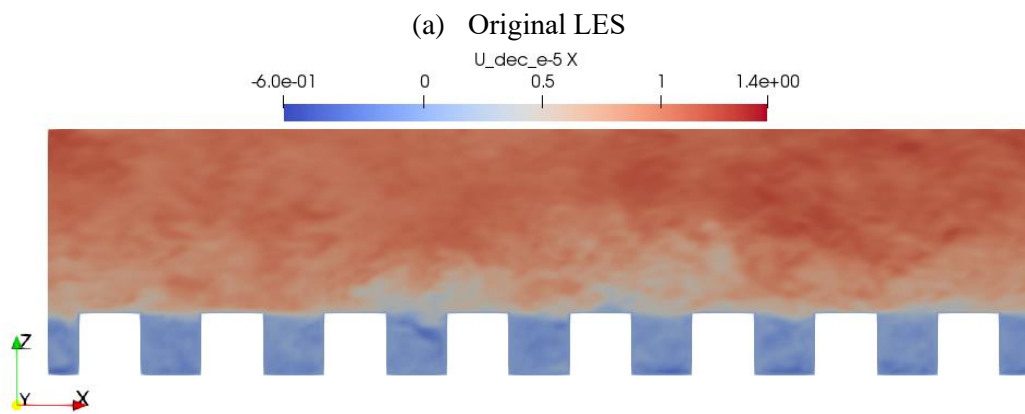
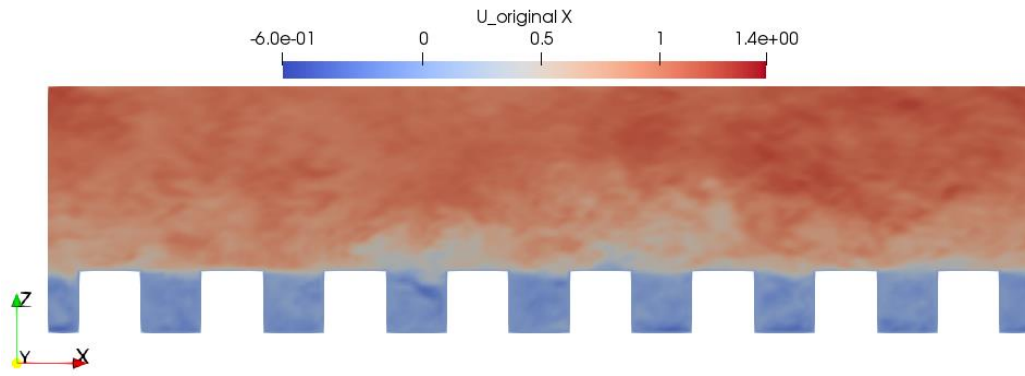
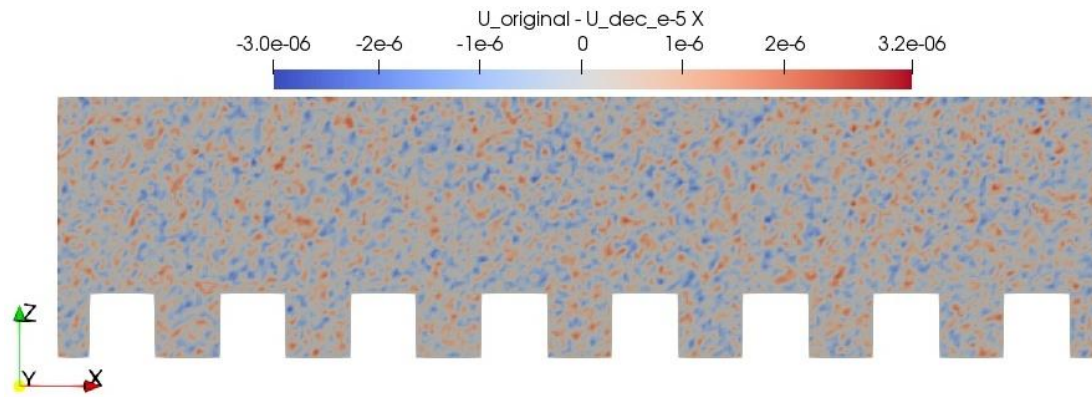
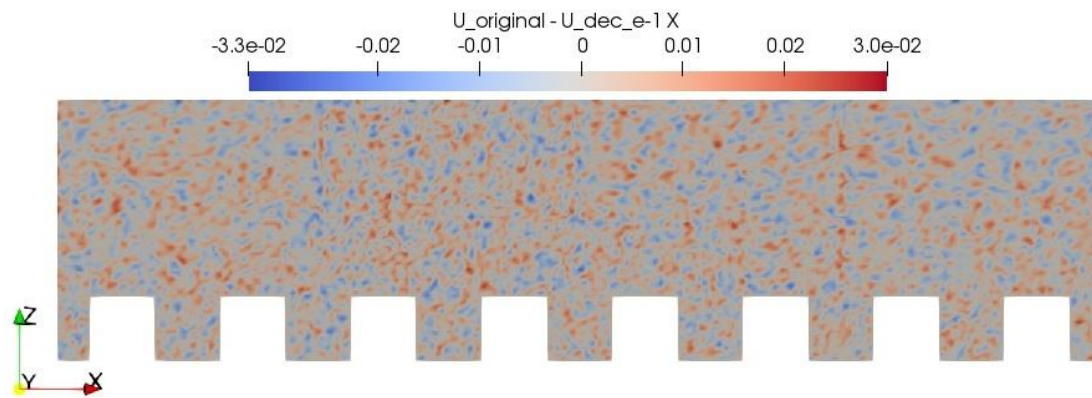


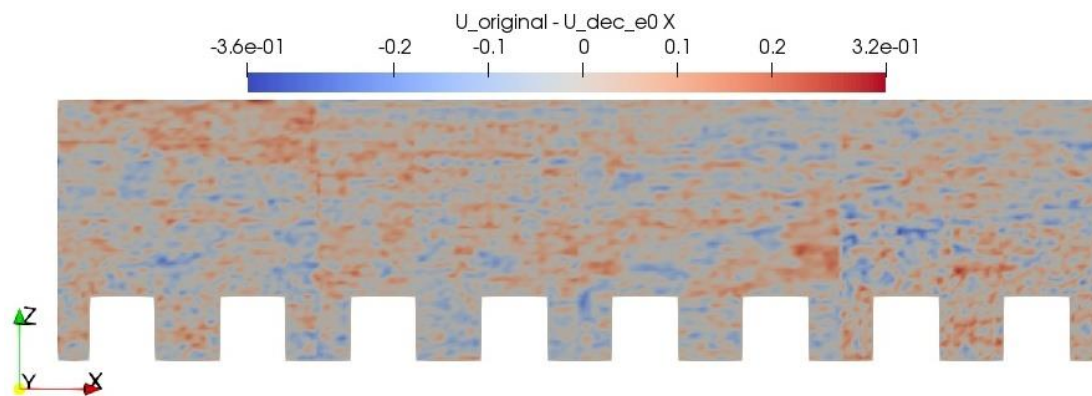
Figure 4.9. Streamwise velocity field in a vertical plane ($y = 7H$)



(a) Difference between original LES and compression database with $\varepsilon = 10^{-5}$



(b) Difference between original LES and compression database with $\varepsilon = 10^{-1}$



(c) Difference between original LES and compression database with $\varepsilon = 10^0$

Figure 4.10. Streamwise velocity difference field between original LES and compression databases in a vertical plane ($y = 7H$)

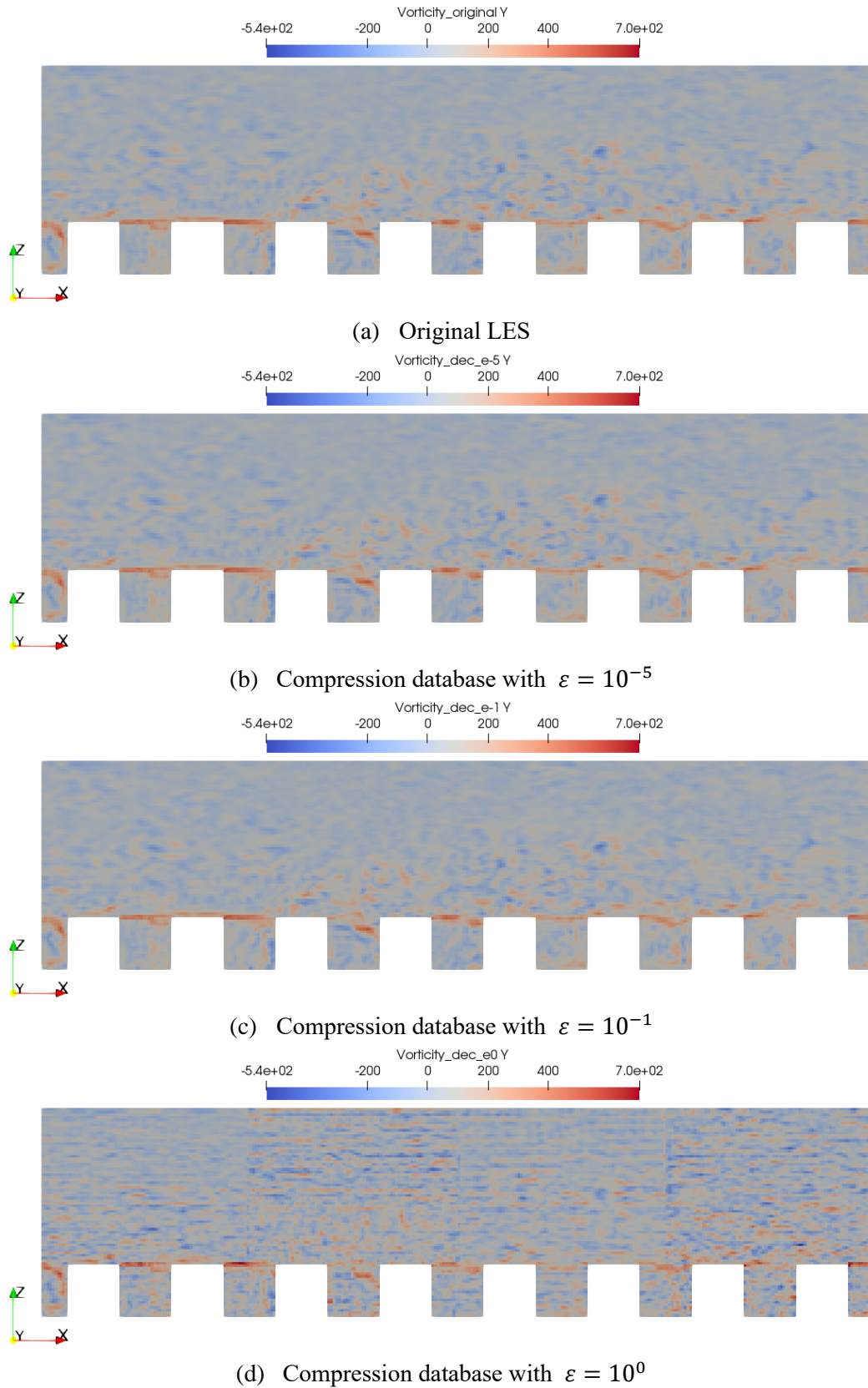
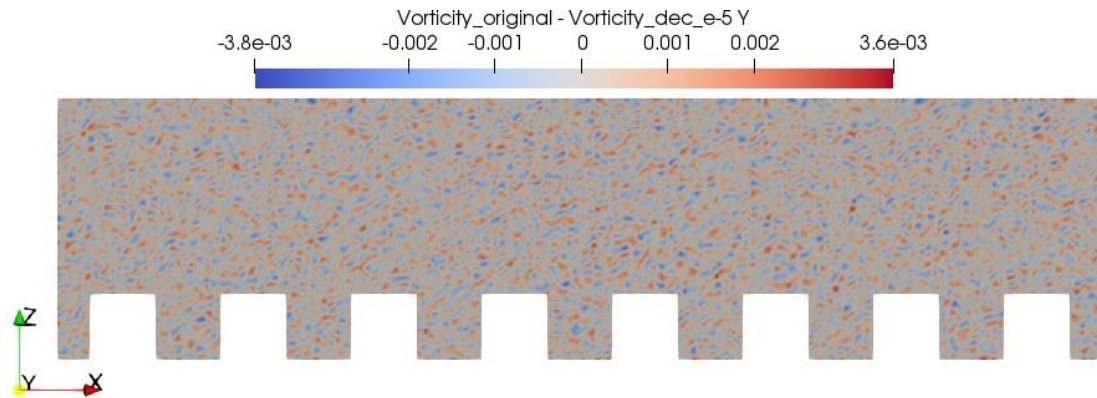
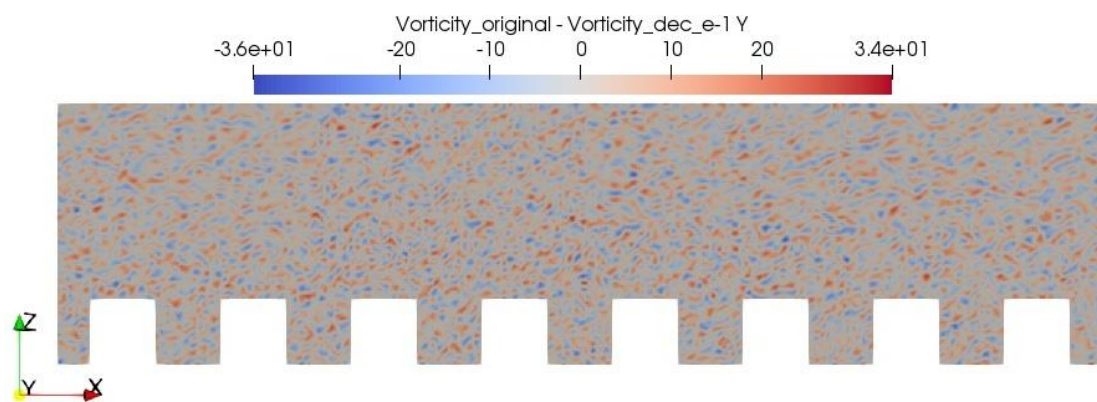


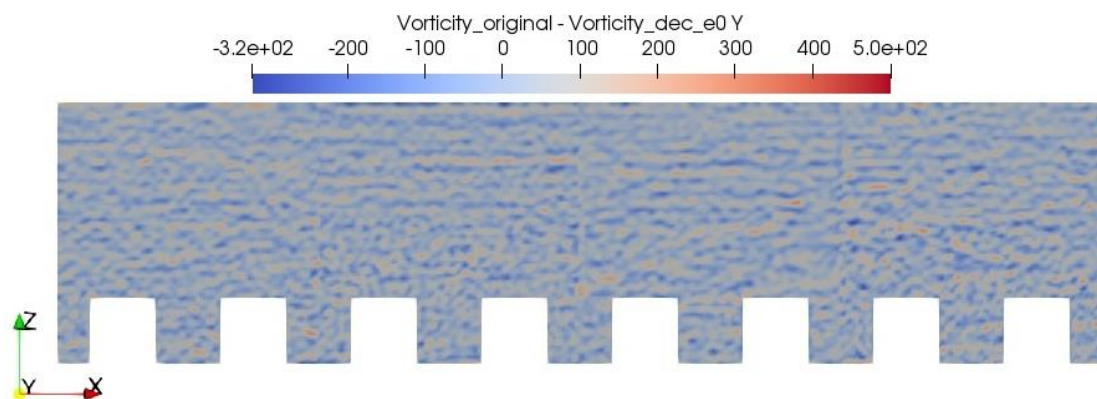
Figure 4.11. Vorticity field of spanwise direction in a vertical plane ($y = 7H$)



(a) Difference between original LES and compression database with $\varepsilon = 10^{-5}$



(b) Difference between original LES and compression database with $\varepsilon = 10^{-1}$



(c) Difference between original LES and compression database with $\varepsilon = 10^0$

Figure 4.12. Spanwise vorticity difference field between original LES and compression databases in a vertical plane ($y = 7H$)

1 As for the vortex structures, the vorticity field in the spanwise direction of the same plane
2 in **Fig. 4.11** shows the same results. The difference fields between the original LES and
3 compression databases can be checked in **Fig. 4.12** in the supplementary material. When ε is
4 smaller than 1, the vorticity field is well maintained in the compressed database. The separation
5 region and wake vortex are clearly seen. While $\varepsilon = 1$, many small vorticities are produced in
6 the whole plane by the compression error, contaminating the original field.

7 A comparison of the velocity and vorticity field reveals that the compression database with
8 $\varepsilon = 10^{-1}$ can be applied for the postprocessing and visualization of a single snapshot because
9 the compression error effects are relatively small. In this case, the compression ratio can reach
10 approximately 100 times according to **Fig. 4.7**. However, when $\varepsilon = 1$, the compression
11 reaches the threshold where significant errors appear. This threshold where most information
12 in the wavelet coefficients has been removed by the quantization process was confirmed in
13 **Section 4.4.2**. In this case, the decompression cannot reconstruct the flow field properly.

14 Considering that $\varepsilon \approx \frac{\max|f-\tilde{f}|}{\max|f|}$, the largest compression error is almost the same as the largest
15 value at this threshold. Hence, there are many discontinuous large values in the velocity field,
16 further resulting in unphysical, small vortexes, as shown in the vorticity field.

17 Another important issue about the velocity field is continuity, which is one of the most
18 fundamental properties in fluid mechanics. Since the simulation case here is incompressible
19 flow, the divergence of the velocity field in one snapshot is a suitable metric. In the ideal status,
20 the divergence should be 0 everywhere in the flow when the continuity is promised, but the
21 value will not be 0 in one snapshot of the numerical simulation due to the differential process
22 and truncation error. **Fig. 4.13** compares the probability density functions of divergence of the
23 velocity of all cells in a snapshot between original LES and three compression cases. Because
24 the curves are too close to the original LES when $\varepsilon < 10^{-2}$, they are not shown here. The
25 interpolation scheme was set as Gauss Linear in the OpenFOAM. As a reference, the average
26 of the absolute value of velocity divergence in the original LES case is about $10(I/s)$.

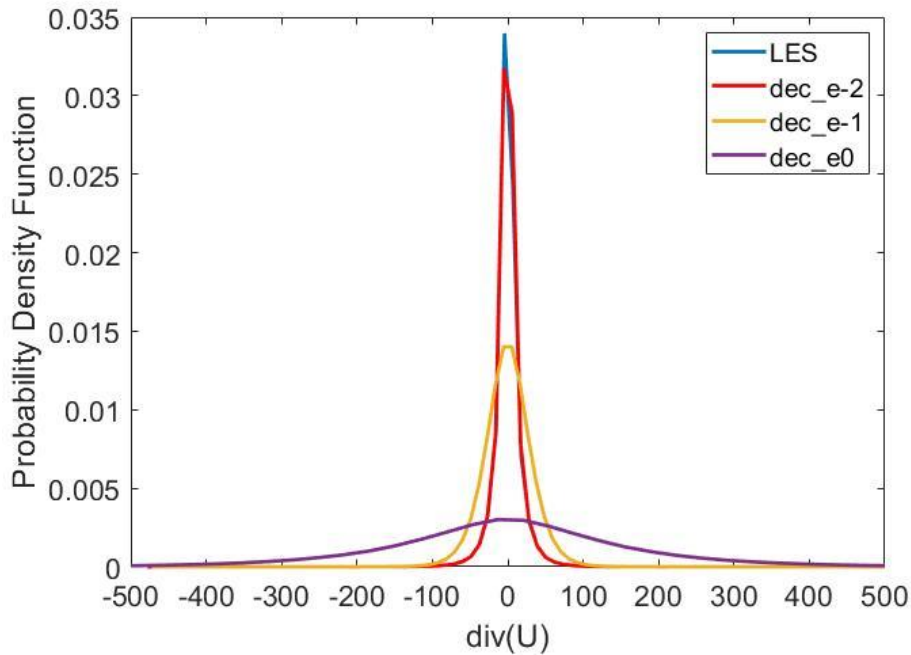


Figure 4.13. Comparison of probability density functions of velocity divergence in one snapshot between original LES and several compression databases: $\varepsilon = 10^{-2}, 10^{-1}, 10^0$

1 From **Fig. 4.13** it can be observed that the curve of the compression database with $\varepsilon =$
 2 10^{-2} is similar to that of the original LES, where the divergence concentrates near 0. The
 3 continuity only suffers small damage due to the compression. As ε becomes larger, the curve
 4 becomes wider and lower, which means the continuity in more areas was disturbed. When the
 5 compression error ε is imposed to the value 10^{-1} , the standard deviation of the velocity
 6 divergence changes from 17 (1/s) to 31 (1/s) holding the same mean value compared with the
 7 original LES. It is important to note that the curve changes significantly when ε turns to 10^0 .
 8 The flattening curve with a large range shows that the continuity of the whole field has been
 9 corrupted. This result agrees well with the above one that the compression reaches its threshold
 10 when $\varepsilon = 1$.

11

12 4.4.4 Compression effects on time-series data

13 The cumulative effects of compression on the database were investigated from the
 14 viewpoint of the velocity field and concentration of the simulated passive scalar.

15 First, the power spectrum densities (PSDs) of velocities at two points located at the canyon
 16 region behind the buildings ($x = 6H, y = 7H, z = 0.5H$) and the open street region between the

1 rows of blocks ($x = 6H, y = 6H, z = 0.5H$) were compared for the original LES, compression
2 database with $\varepsilon = 10^{-5}$, and compression database with $\varepsilon = 10^{-1}$. The results in **Fig. 4.14**
3 show that the PSD curves of the three cases are almost the same when the frequency is relatively
4 small, indicating that large-scale vortex structures are well preserved in the compression. For
5 the high-frequency area, the larger ε , the larger the difference between the compressed and
6 original data. This phenomenon can be explained by the principles of the WCM. WD
7 decomposes the flow field into an approximation part with low frequency and a detailed part
8 with high frequency. Because the truncation error in the quantization mainly deals with the
9 detailed part whose wavelet coefficients are much smaller than those of the approximation part,
10 after the compression, the high-frequency structures are highly coupled with truncation errors,
11 which are manifested as white noise in the PSD. For the same reason, the compression data
12 with $\varepsilon = 10^{-5}$ coincide with the original data longer and more closely than the data with $\varepsilon =$
13 10^{-1} because the truncation error is set smaller. Additionally, for the point in the open street
14 region, these three curves are closer than that in the wake region. The reason is that the velocity
15 value is higher in this area, so the effects of compression error are weaker.

16 The probability density functions of the velocity at the same two points are summarized
17 in **Fig. 4.15**. At both points, the velocity distribution of case $\varepsilon = 10^{-5}$ coincides well with the
18 original data. As for $\varepsilon = 10^{-1}$, the distribution deviates at the place where the velocity is near
19 0 for the point in the canyon region. These small velocities are mainly contributed by a high-
20 frequency small vortex, which suffers from the damage from the compression process to some
21 degree. It is reasonable to conclude that the compression barely disturbs the occurrence
22 probability distribution of the velocity field.

23

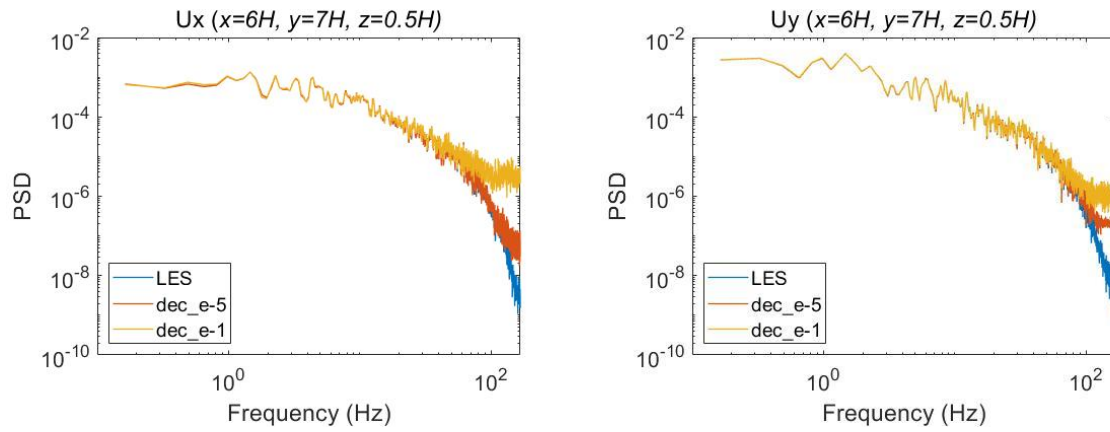
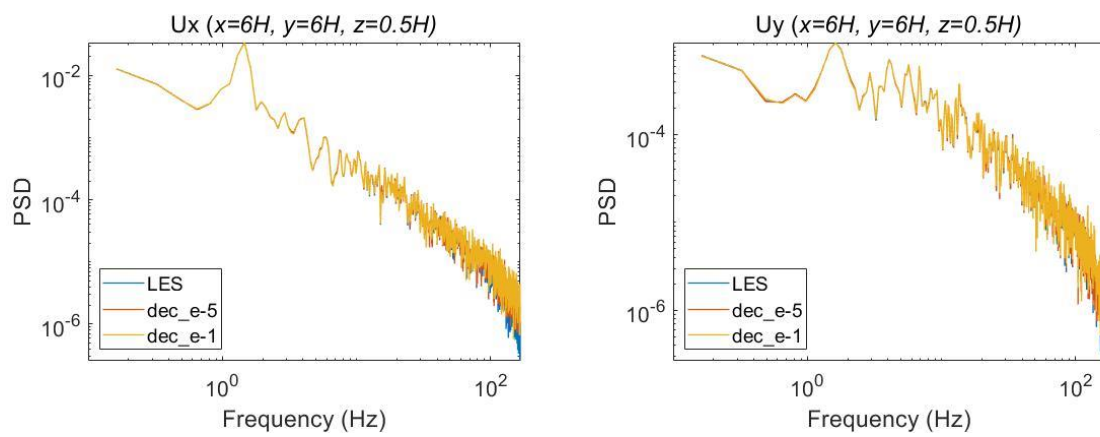
(a) Point at canyon region ($x = 6H, y = 7H, z = 0.5H$)(b) Point at open street region ($x = 6H, y = 6H, z = 0.5H$)

Figure 4.14. Comparison of the PSD of velocity (streamwise and spanwise) at two points in three cases: original LES, compression database with $\varepsilon = 10^{-5}$ (dec_e-5), and compression database with $\varepsilon = 10^{-1}$ (dec_e-1)

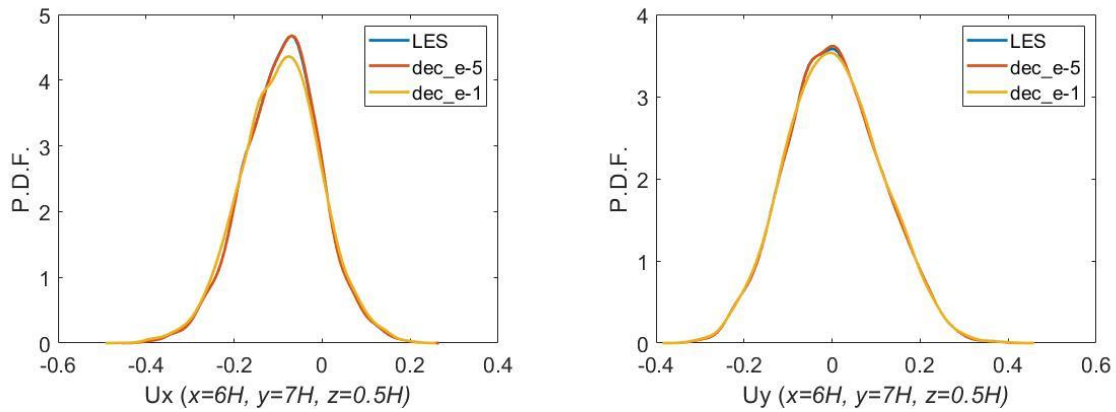
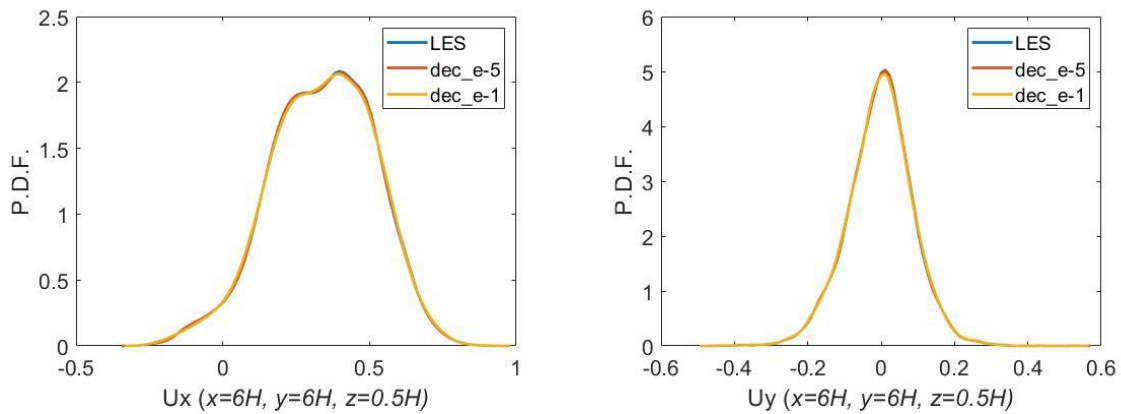
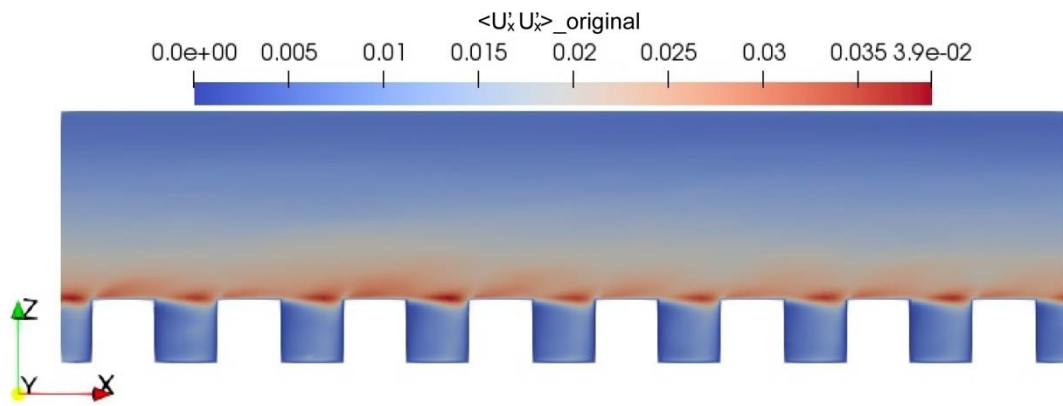
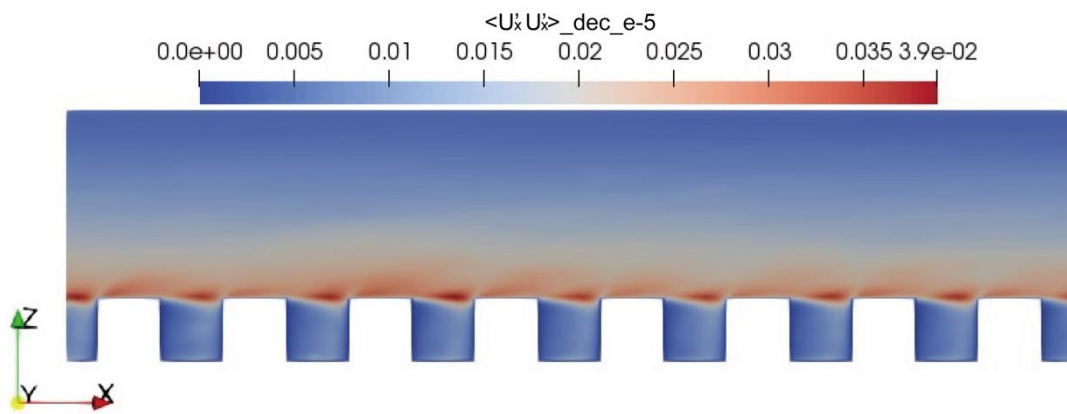
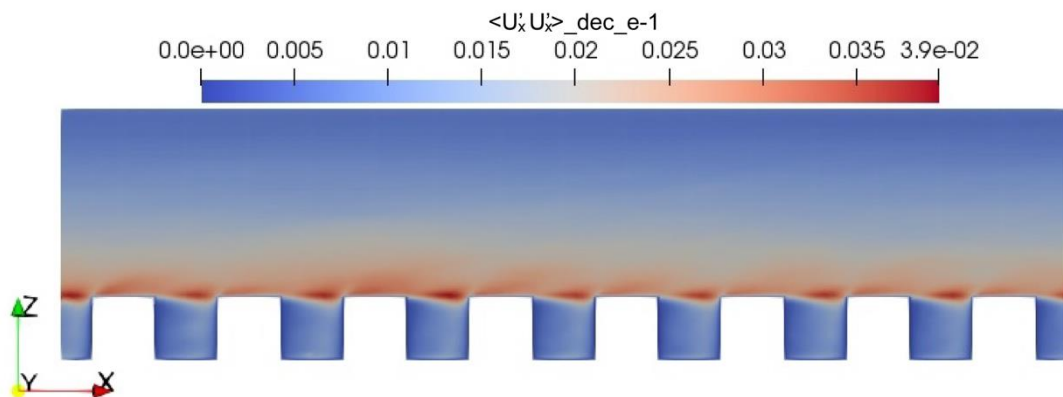
(a) Point at canyon region ($x = 6H, y = 7H, z = 0.5H$)(b) Point at open street region ($x = 6H, y = 6H, z = 0.5H$)

Figure 4.15. Comparison of probability density functions of velocity at two points in three cases: original LES, compression database with $\varepsilon = 10^{-5}$ (dec_e-5), and compression database with $\varepsilon = 10^{-1}$ (dec_e-1)

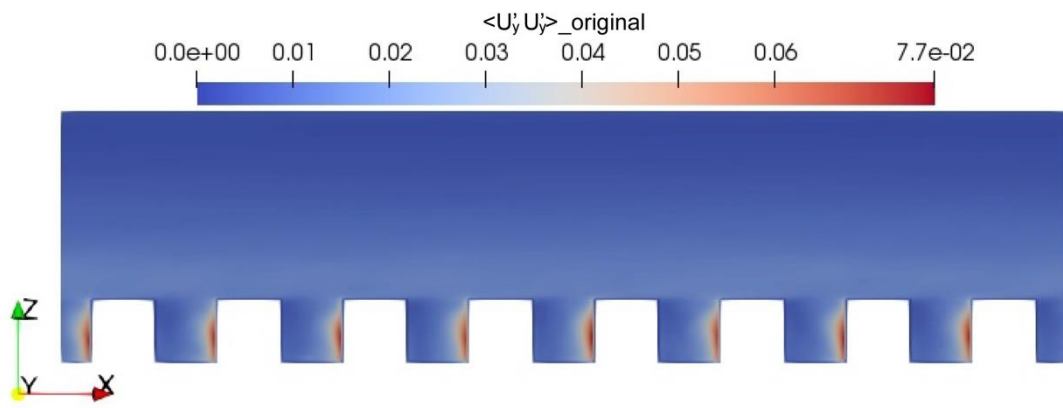
1 The effects of compression on the higher order statistics (variance and covariance of
 2 velocity fluctuations) are also checked. The distributions of variance $\langle U'_x U'_x \rangle$, $\langle U'_y U'_y \rangle$,
 3 $\langle U'_z U'_z \rangle$, $\langle U'_x U'_z \rangle$, $\langle U'_x U'_y \rangle$ and $\langle U'_y U'_z \rangle$ in the middle vertical plane are presented
 4 in **Fig. 4.16 ~ 4.21**. Here, x, y, z denote the three components of the velocity field, $'$ means
 5 the fluctuations, and $\langle \bullet \rangle$ is the time-averaged operator. It is shown that the 100 times
 6 compression preserves the second order statistics as well as the instantaneous flow field. Both
 7 the magnitude and shapes are almost the same as the original LES. The dominant co-variance
 8 structures caused by the separation and wake of the buildings are completely maintained.



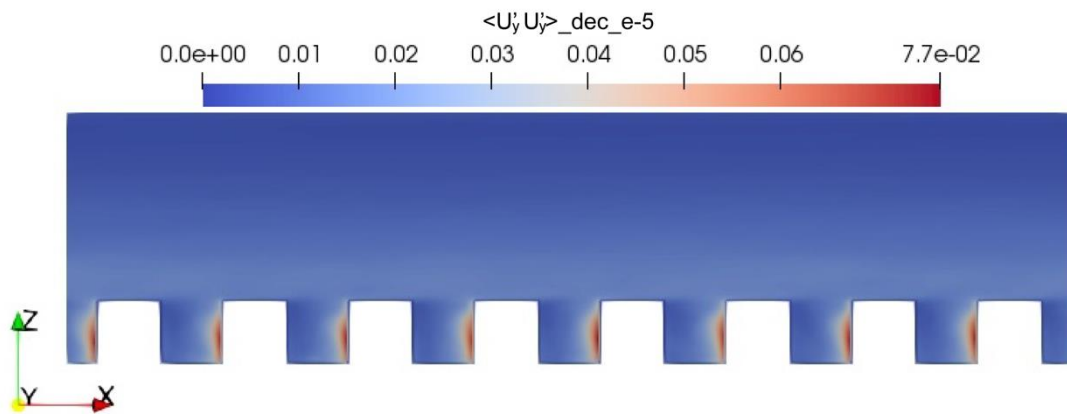
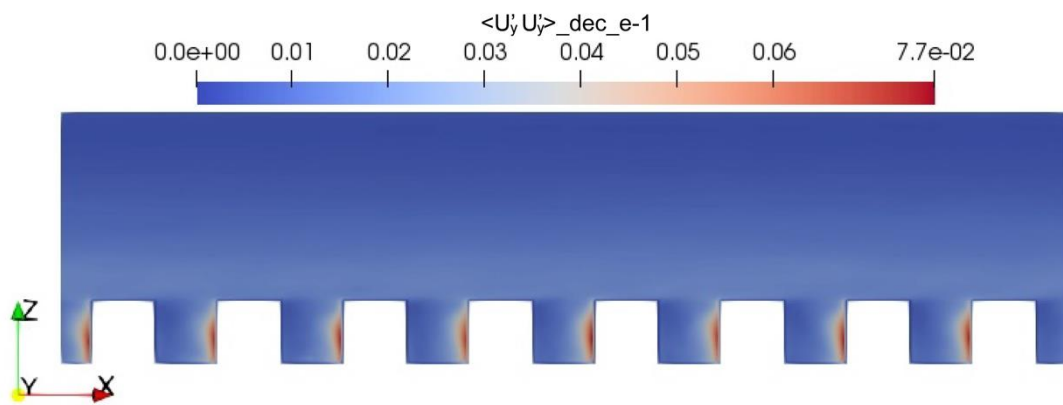
(a) Original LES

(b) Compression database with $\varepsilon = 10^{-5}$ (c) Compression database with $\varepsilon = 10^{-1}$ Figure 4.16. The distribution of $\langle U'_x U'_x \rangle$ of different cases in a vertical plane ($y = 7H$)

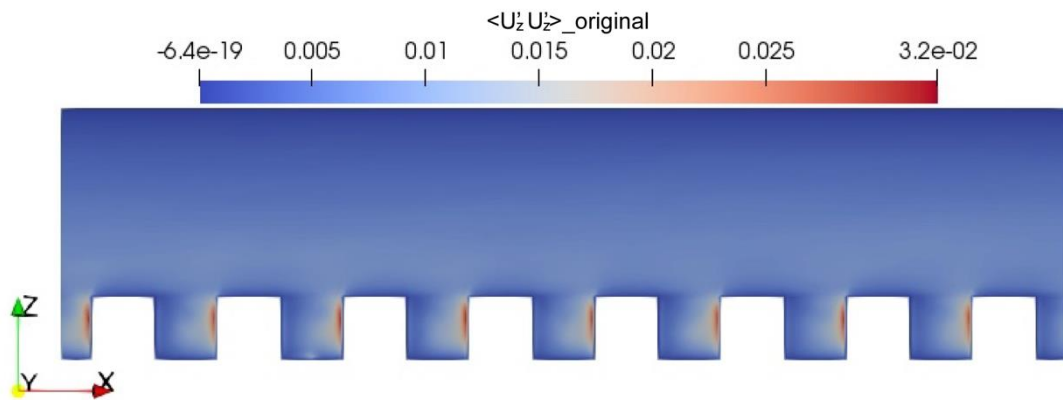
1



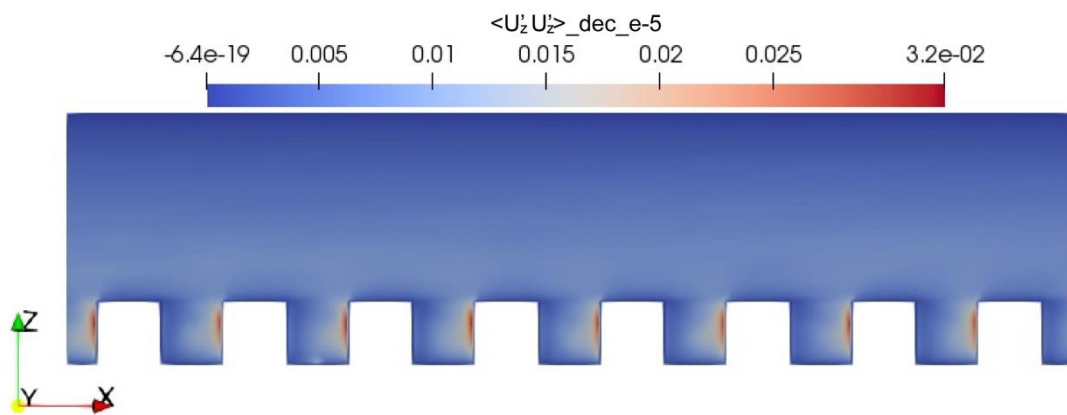
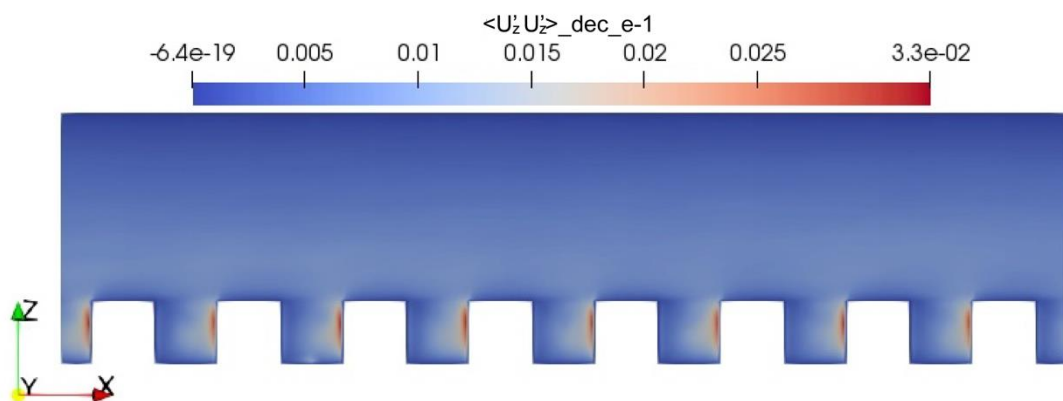
(a) Original LES

(b) Compression database with $\varepsilon = 10^{-5}$ (c) Compression database with $\varepsilon = 10^{-1}$ Figure 4.17. The distribution of $\langle U'_y U'_y \rangle$ of different cases in a vertical plane ($y = 7H$)

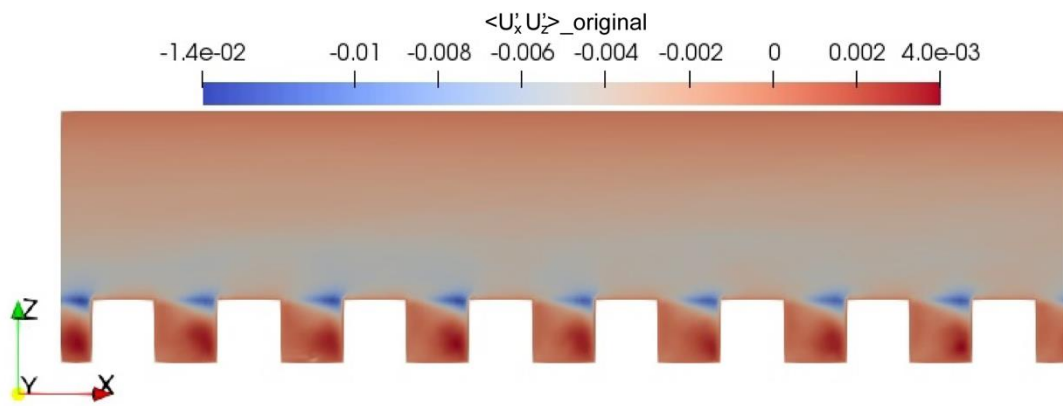
2



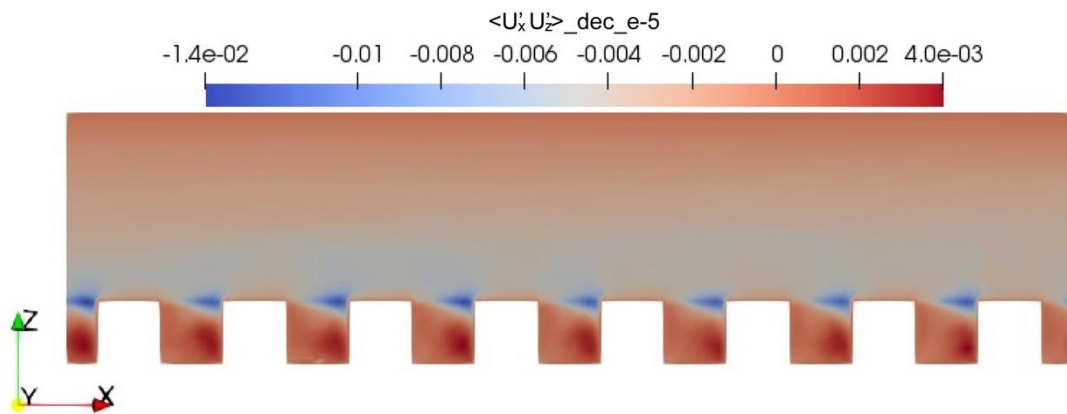
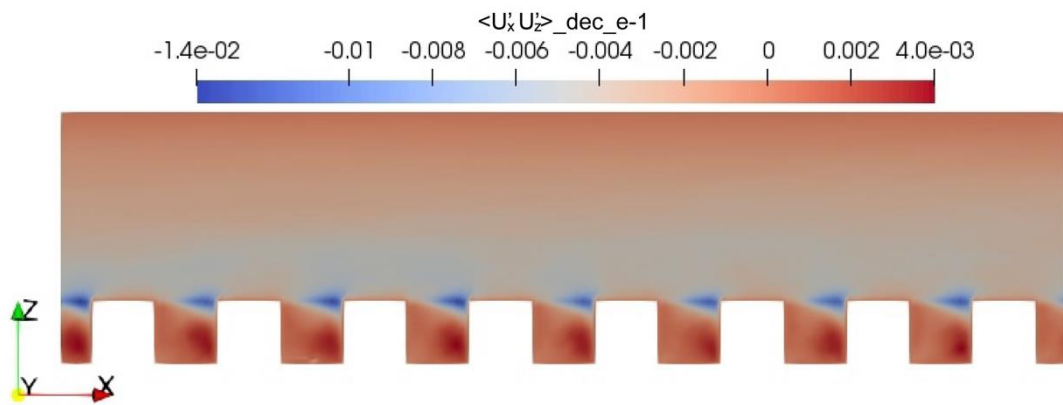
(a) Original LES

(b) Compression database with $\varepsilon = 10^{-5}$ (c) Compression database with $\varepsilon = 10^{-1}$ Figure 4.18. The distribution of $\langle U'_z U'_z \rangle$ of different cases in a vertical plane ($y = 7H$)

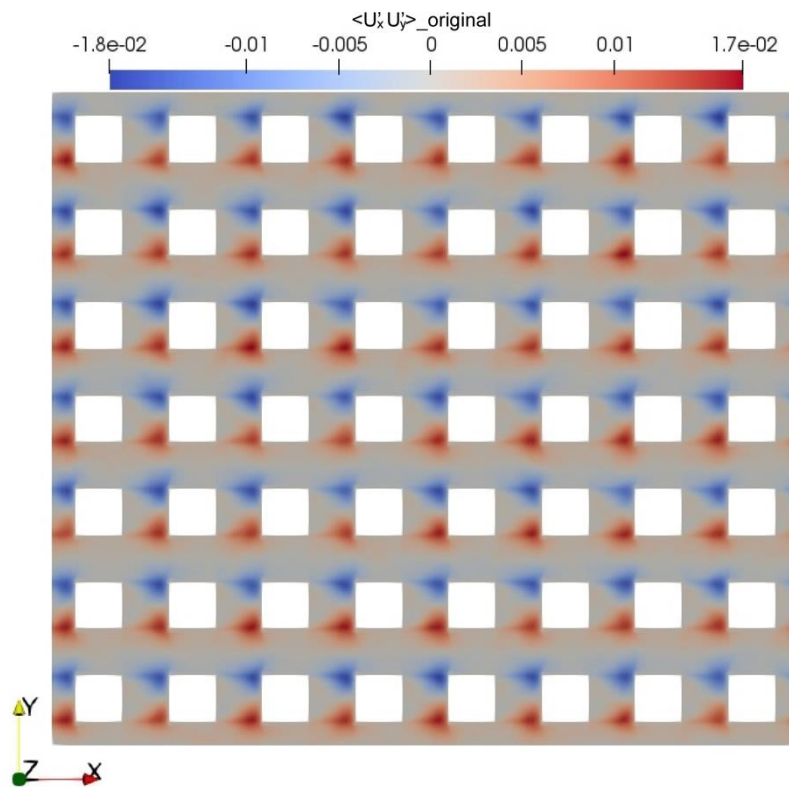
1



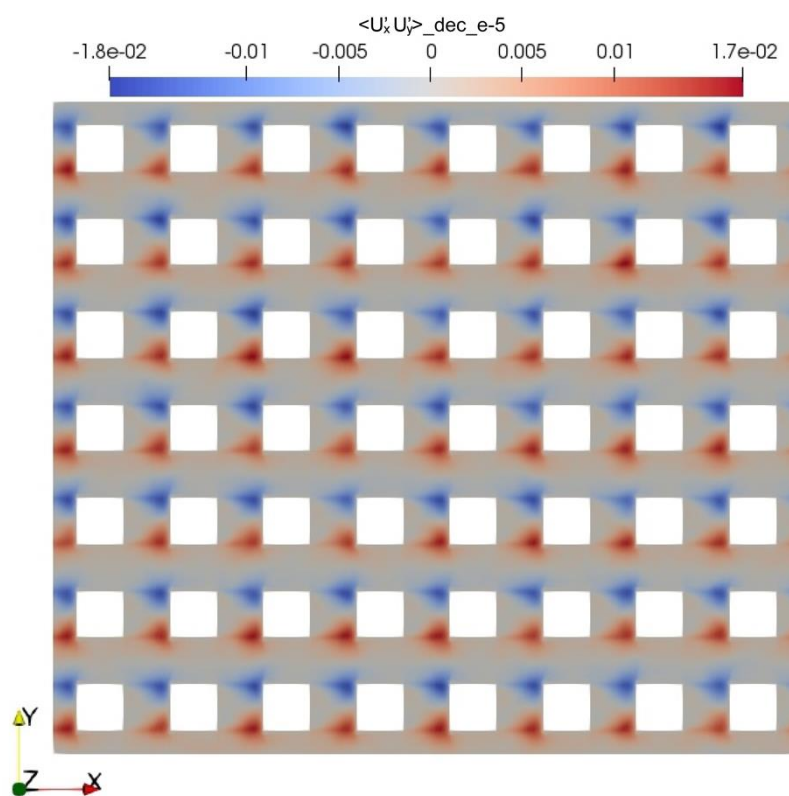
(a) Original LES

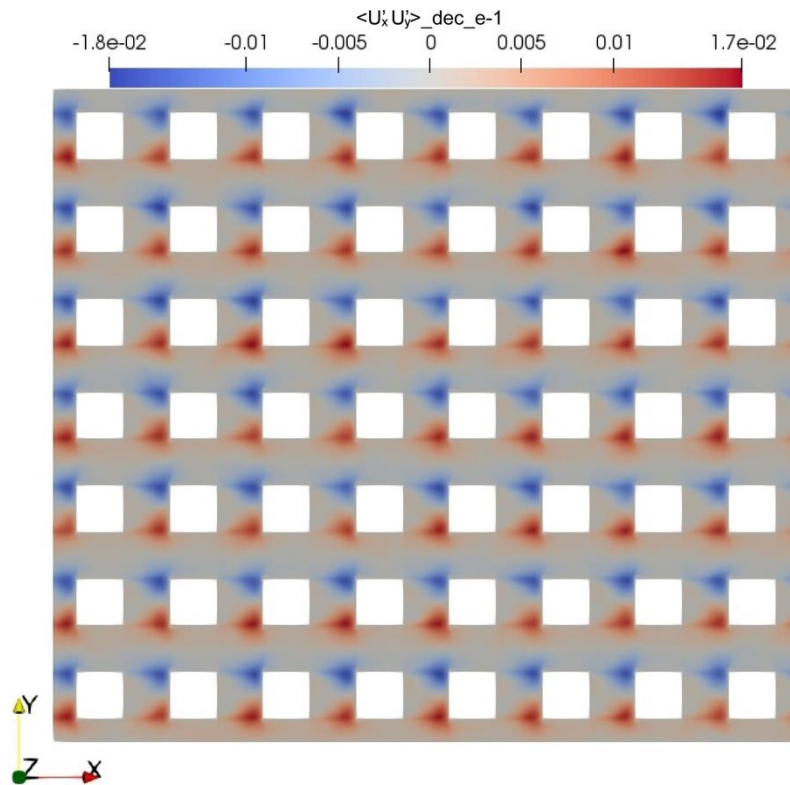
(b) Compression database with $\epsilon = 10^{-5}$ (c) Compression database with $\epsilon = 10^{-1}$ Figure 4.19. The distribution of $\langle U'_x U'_z \rangle$ of different cases in a vertical plane ($y = 7H$)

2



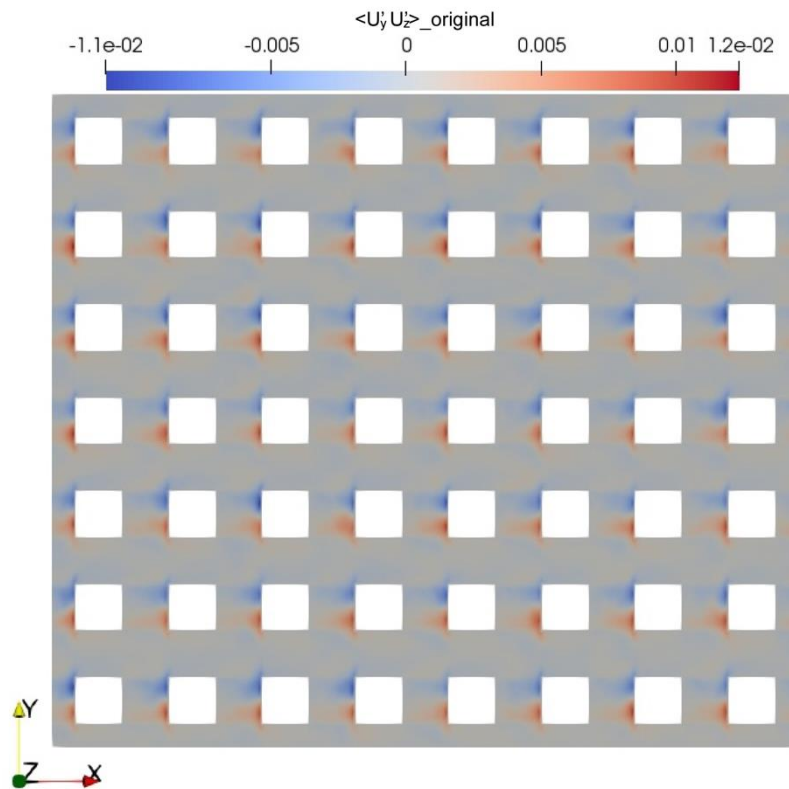
(a) Original LES

(b) Compression database with $\varepsilon = 10^{-5}$

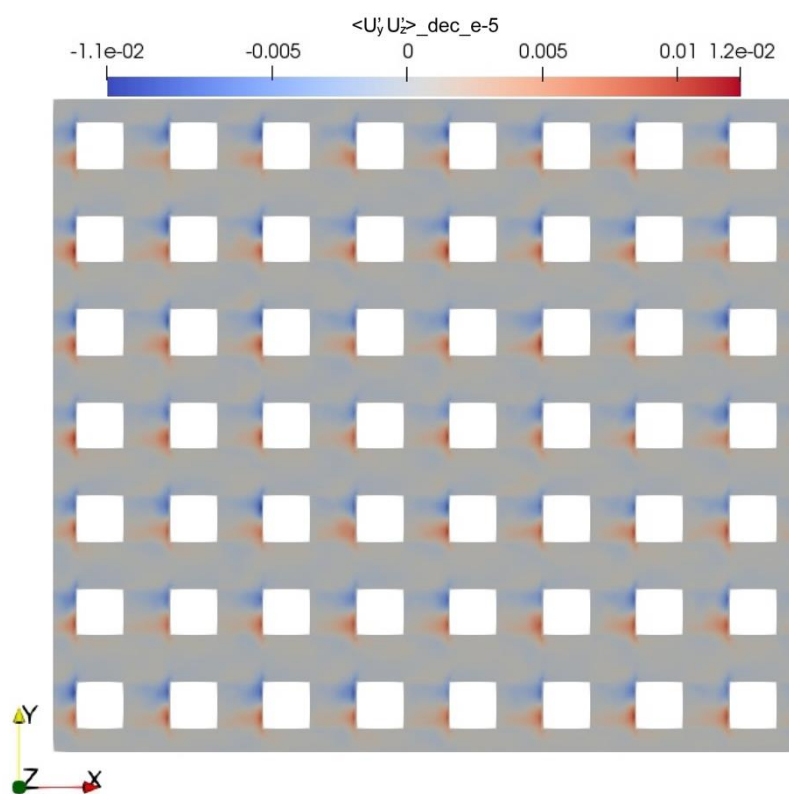


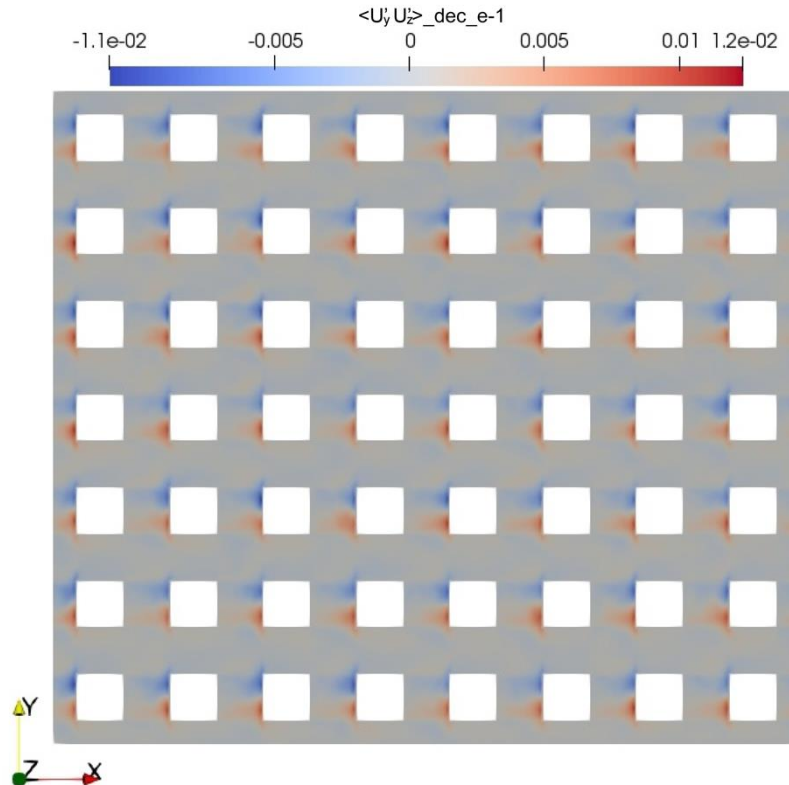
(c) Compression database with $\varepsilon = 10^{-1}$

Figure 4.20. The distribution of $\langle U'_x U'_y \rangle$ of different cases in a vertical plane ($y = 7H$)



(a) Original LES

(b) Compression database with $\varepsilon = 10^{-5}$

(c) Compression database with $\varepsilon = 10^{-1}$ Figure 4.21. The distribution of $\langle U'_y U'_z \rangle$ of different cases in a vertical plane ($y = 7H$)

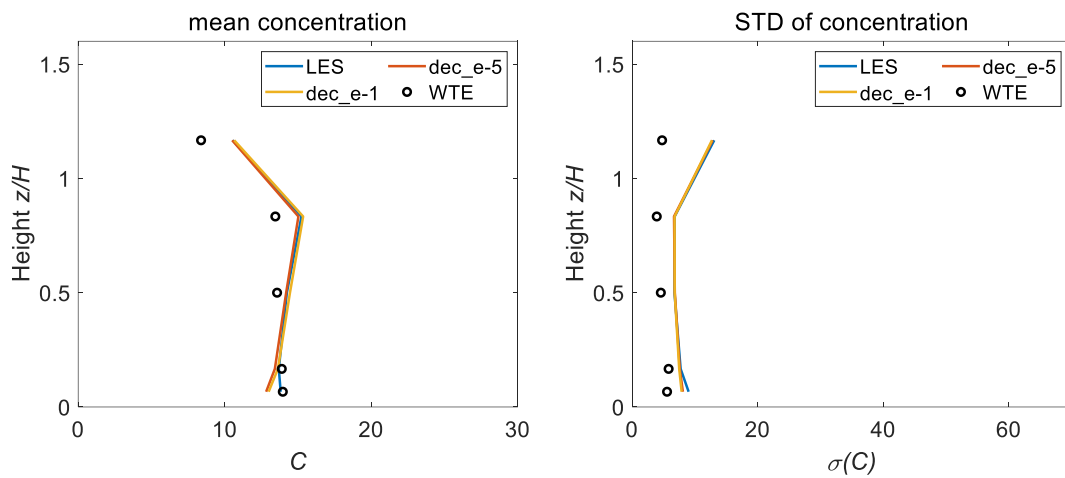
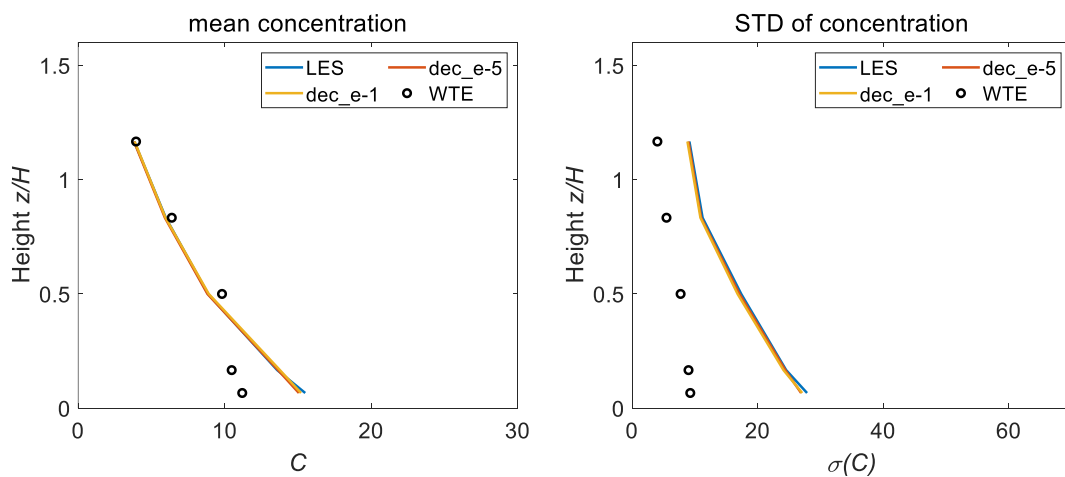
1

2 **4.4.5 Dispersion simulation with compression database**

3 Then, the re-simulation of the passive scalar dispersion was evaluated with the compressed
 4 database. Comparisons between the wind tunnel experiment, original LES, and compressed
 5 databases at four representative locations — the canyon region near the source ($x = 6H, y =$
 6 $7H$), street region near the source ($x = 6H, y = 8H$), canyon region away from the source ($x =$
 7 $10H, y = 7H$), and street region away from the source ($x = 10H, y = 8H$) — can be observed in
 8 **Fig. 4.22**. The original LES basically predicted the concentration distribution in accordance
 9 with the experiment. At the location in the open street region near the source, LES slightly
 10 overestimated the turbulent scalar flux. Because a method to simulate accurately the scalar
 11 dispersion in the complicated building blocks is still being researched (Tominaga and
 12 Stathopoulos, 2012) and was not the concern of this study, the overestimation was not analyzed
 13 in detail.

14 Meantime, the concentration of the original LES was successfully reproduced by the

1 compressed data, among which the smallest data volume was only approximately 1% of the
 2 original one (compression database with $\varepsilon = 10^{-1}$). At the point in the canyon region near the
 3 source, the concentration of the compressed database had a small deviation near the ground.
 4 The velocity at this height was somewhat small, indicating that it is vulnerable to compression
 5 error. Given that the concentration flux is rather strong near the source, a slightly different
 6 velocity can result in a different concentration distribution and cause this kind of deviation. At
 7 the places that were relatively far from the source, the effects of compression became almost
 8 unnoticeable. Generally, not only the mean concentrations but also the standard deviation of all
 9 the simulations of the compressed database were close to that of the original LES results.

(a) $x = 6H$ $y = 7H$ (b) $x = 6H$ $y = 8H$

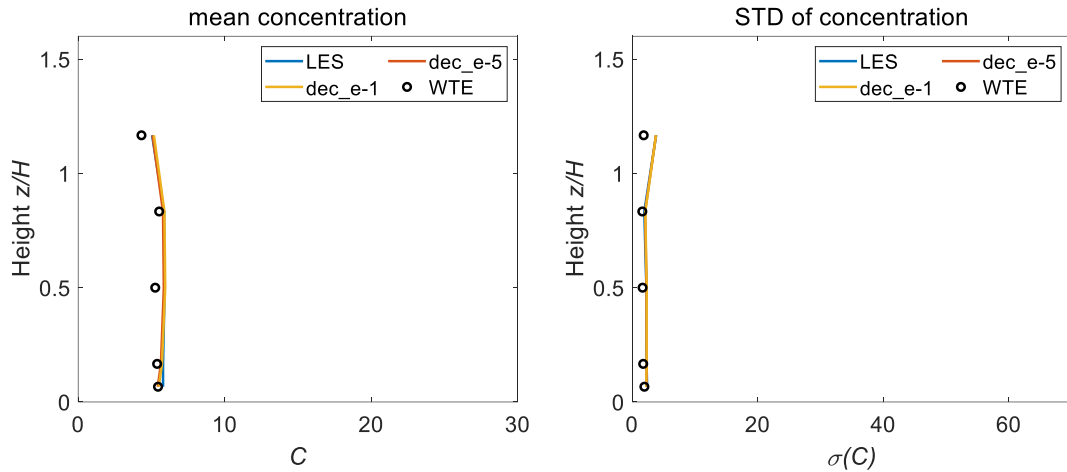
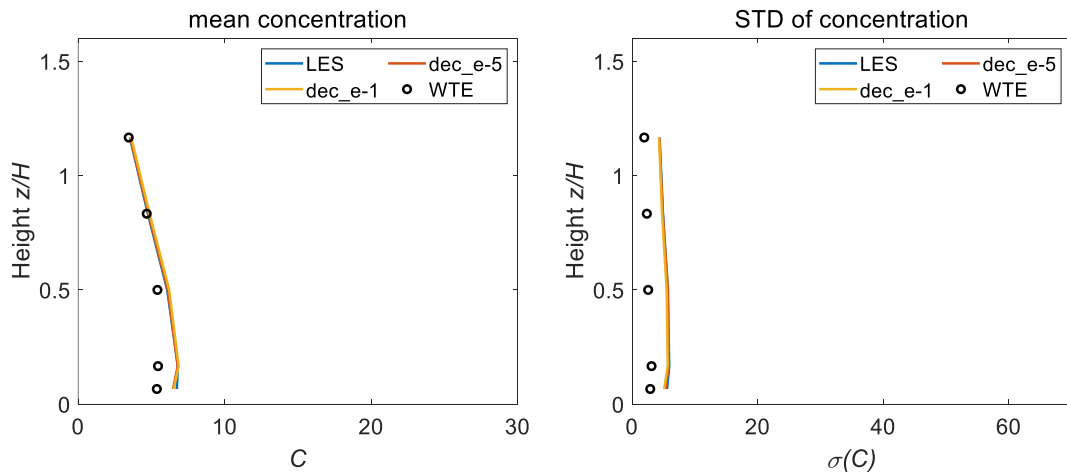
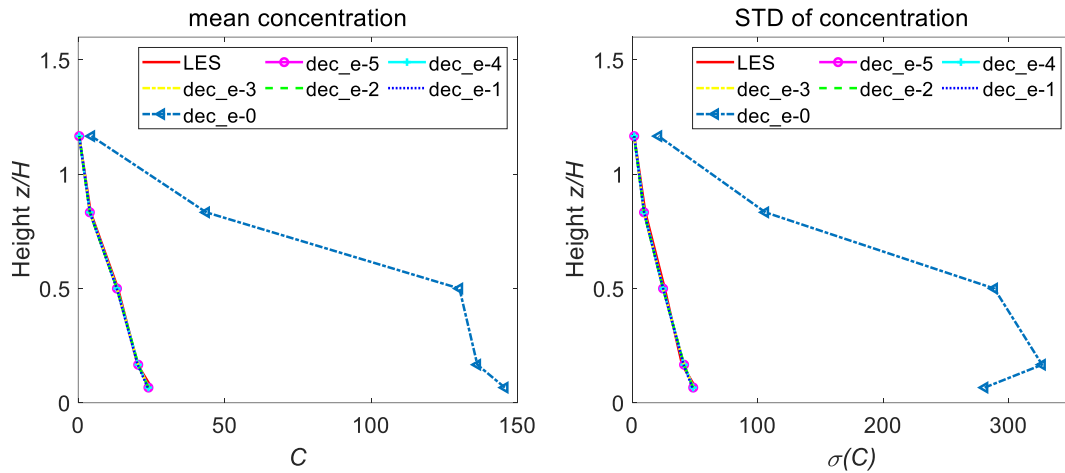
(c) $x = 10H$ $y = 7H$ (d) $x = 10H$ $y = 8H$

Figure 4.22. Comparison of the concentration at different locations in four cases: wind tunnel experiment, original LES, compression database with $\varepsilon = 10^{-5}$ (dec_e-5), and compression database with $\varepsilon = 10^{-1}$ (dec_e-1) — left: time mean concentration; right: standard deviation of concentration

1 When the compressed database with $\varepsilon = 10^0$ was used, the statistical result of
 2 concentration totally collapsed within only 1s (approximately $5T$) simulation time, while the
 3 other cases remained stable (**Fig. 4.23**). As a consequence, the concentration simulation was
 4 insensitive to the compression error under the threshold. When ε was loosened to 1, an
 5 enormous error appeared. Once again, the reason was that the truncation error in the
 6 compression was so large that the main information in the wavelet coefficients was removed;

- 1 therefore, the flow structures were damaged, and the reconstructed flux could no longer
 2 transport the passive scalars correctly.



$$x = 4.5H \quad y = 6H$$

Figure 4.23. Comparison of the concentration sampled in a short time in seven cases: original LES and all the compression cases — left: time mean concentration; right: standard deviation of the concentration

- 3 To sum up, the dispersion of passive scalars can be obtained later, as in the original LES
 4 simulation with a compressed database with $\varepsilon = 10^{-1}$, whose volume was approximately 1%
 5 of the original one according to **Fig. 4.7**.

6

7 4.5 Conclusions and discussions

8 4.5.1 Conclusions

- 9 The feasibility of constructing a small urban flow database by compressing the original
 10 CFD data with the WCM was investigated. The original data were flow fields in a block-arrayed
 11 urban area simulated by the LES model. Several compressed databases were constructed with
 12 different error controls ε : 10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} , and 10^0 .

- 13 The compression method was based on WD, quantization, and entropy encoding. The
 14 original data were first transformed by WD to the approximate part with a small volume and
 15 detailed parts with a large volume. The quantization process decided the proportion of detailed
 16 parts that should be retained according to the user-specified error control ε and caused the

1 truncation error ε_F . Finally, entropy encoding was implemented to compress the data further.

2 The compression ability of the WCM was studied. The empirical formula for error control
3 guarantees that the prescribed error is exactly realized in the compression. The compression
4 ratio increases with ε exponentially until it reaches the compression limitation, approximately
5 1000 times in this study, where $\varepsilon = 10^0$ is so large that most of the wavelet coefficients are
6 deeply stained. It was found that, although the error control that can reach this compression
7 limitation is found according to each case, it can be calculated from the distribution of the
8 wavelet coefficient in advance.

9 The influence of compression on the quality of the database was checked from two
10 viewpoints: single snapshot and time-series data. On one hand, for a single snapshot used in
11 postprocessing and visualization, the velocity and vorticity fields are almost the same as long
12 as ε is less than the threshold of 10^0 , which means that the compression database with
13 approximately 1% of the volume of the original data is sufficient for visualization. When
14 compression reaches the threshold, a considerable number of unphysical large values in velocity
15 and small vorticities appear, and the flow structures are highly contaminated. On the other hand,
16 the cumulative effects were inspected by comparing the velocity field and different dispersion
17 results re-simulated with the compressed database to the original LES results. According to the
18 PSD velocity results, the large-scale vortex structures with low frequencies are well preserved
19 in the compression. The compression errors function as white noise superimposed on the high-
20 frequency structures. Moreover, the dispersion re-simulation is insensitive to the compression
21 error below the compression limitation. The dispersion process can be reproduced closely, even
22 by rough compressed data with approximately 1% of the original volume. A considerable
23 discrepancy appears near the compression limitation where the flow field can no longer be
24 reconstructed properly.

25 Overall, the results in this chapter show that WCM is a powerful tool to construct a
26 portable flow database. The possible applications of a database, such as visualization of a single
27 snapshot and re-simulation with time-series data, can be appropriately realized, even under
28 approximately 100-times compression. Therefore, in the next chapter, the unsteady adjoint
29 equation will be simulated by the compression database constructed here, and the performance
30 in source identifications will be checked.

31

4.5.2 Method limitation and future research

It has to be admitted that this chapter is one of the early attempts of a compressed turbulence database focusing on a complicated urban environmental flow, more situations in other areas should be examined before the general application of WCM. The current method still has some limitations which can be improved in future research. Since the method consists of three steps: WD, quantization, and entropy encoding, a brief discussion about method limitation and possible future research plans is conducted in these three aspects.

The first one is WD. In the current method, all the wavelet coefficients are passed exactly to quantization, which means it is lossless and unrelated to compression error control. The improvement strategy should focus on the calculation speed and coefficients production. The core part of WD is the wavelet kernel function. The kernel function that requires less calculation resource and yields fewer coefficients is favorable to the compression process. However, it seems that these two properties are contradictory to each other. Although the current kernel CDF9/7 performs well in the coefficients production, it is one of the most complicated kernels because of its irrational decomposition coefficients. It is worthwhile to testify the performance of other improved kernels based on CDF9/7 in the future. Meanwhile, even though CDF9/7 is highly recommended in image compression, it cannot be simply asserted to be the best one for turbulence compression. After all, image compression focuses on the visual effects to humankind, in which the low-frequency domain is the most important part, while turbulence focuses on the physical information, which is contained in different scales structures with a large range of frequency. Deeper research is necessary to decide the most appropriate kernel for turbulence decomposition.

The only step concerning error control is quantization. The main limitation of the method is the way to control compression error. According to the definition of the error metric $\varepsilon_r = \frac{\max|f-\tilde{f}|}{\max|f|}$, the error control is closely related to the maximum value of the flow field data and maximum compression error rather than the spatial mean value or the range of all values. However, the place where the maximum compression error happens is not promised to be the place where the flow field takes the maximum value. In other words, the situation of all places in the field cannot be independently determined by the current error metric so that someplace may suffer unexpected errors resulting from the maximum field data. For example, if the inflow

1 speed is increased in the case of this paper, the maximum velocity in the high place will increase
2 correspondingly. It is possible that the compression error of the low velocity area near walls
3 will be increased even if the compression algorithm is set with the same ε and yields a close
4 compression ratio. Therefore, it has to be very cautious to apply this method to compress the
5 flow field with a wide range of values. In the common urban flow simulation, the range of
6 physical values is mild and high Reynolds number flow is the main concern, author is optimistic
7 about the application of the current method. However, special attention needs to be paid to when
8 both the low and high Reynolds number turbulence matters.

9 Entropy encoding is a lossless compression algorithm that is also included in other
10 compression methods like JPEG2000 as a standard process. It is reasonable to say that the
11 effectiveness of this step would not be changed by different application scenarios of turbulence
12 compression.

13

14

1 Symbols

C	: the concentration distribution
C_r	: standard (reference) concentration
C_s	: Smagorinsky constant for LES model (=0.12)
C_{gas}	: the emission strength of the source
D_e	: effective diffusion ($D_m + D_{sgs}$)
D_m	: molecular diffusivity
D_{sgs}	: sub-grid scaled turbulent diffusivity
\widetilde{D}_{sgs}	: decompressed sub-grid scaled turbulent diffusivity
$f_{i,j,k}$: the original data of a three-dimensional scalar field in the Cartesian indexing mesh grid
\check{f}	: the decompressed flow field data
$F_{i,j,k}$: the wavelet coefficients obtained from the wavelet decomposition of raw data $f_{i,j,k}$
\check{F}	: the wavelet coefficients reconstructed from quantization results Q
H	: the edge length of blocks in the simulation (=60 mm)
h_r	: reference height
N	: the number of 1-byte memory spaces
q	: the gas flow rate at the source

Q	: the results produced by quantization operation
r	: the compression ratio
S	: source term in the transport equation
Sc_{sgs}	: sub-grid scaled turbulent Schmidt number
T	: standard time-scale of LES simulation
u_r	: reference speed
u_*	: friction speed
U_x, U_y, U_z	: streamwise, spanwise, and vertical velocity
\mathbf{U}	: velocity field
$\check{\mathbf{U}}$: decompressed velocity field
V_o	: storage volume of raw data
V_c	: storage volume of compressed data
ε_F	: truncation errors between F and \check{F} caused by quantization operation
ε	: user-specified allowable error for the WCM
ε_r	: real compression error metric between f and \check{f}
η	: empirical coefficient for error control of WCM (=1.75)
ν_{sgs}	: eddy viscosity coefficient at the sub-grid scale
$\sigma(\cdot)$: standard deviation

$(\cdot)'$: fluctuation value with time

$\langle \cdot \rangle$: time-averaged operator

1

2

1

2

3

4

5

6

7

8

9

Chapter 5

10 Source term estimation with
11 unsteady adjoint equations
12 modeled by large-eddy simulation
13 and compression database

14

1

2

3

4

5

6

Abstract

7

8

9

This chapter tries to develop a new source term estimation (STE) method in which the unsteady simulation of adjoint equations is embedded via large eddy simulation (LES) into a Bayesian inference framework. The LES of adjoint equations is based on the compression database constructed in **Chapter 4**. The performance of the proposed model is validated using a point source dispersion case in a regular, block-arrayed, urban model wind tunnel experiment. To clarify the effects of different computational fluid dynamics models on the simulation of adjoint equations and the accuracy of STE, an existing (Xue et al., 2018b) Reynolds averaged Navier-Stokes model using a time-averaged LES flow field is selected for comparison.

17

5.1 Introduction

An adjoint equation is a partial differential equation that is similar in form to a dispersion equation and can be simulated by computational fluid dynamics (CFD) models. Among them, two popular methods are the Reynolds averaged Navier-Stokes (RANS) model and large eddy simulation (LES). The main difference between these models is the way in which turbulent diffusion is simulated. The RANS model is based on an ensemble-averaged governing equation and approximates fluctuations in the turbulent flux using a mean field according to the gradient diffusion hypothesis (GDH), which assumes that turbulent diffusion is proportional to the gradient of the mean concentration field. Meanwhile, an LES explicitly resolves most turbulent effects using a fine mesh. The RANS model has been frequently used to simulate adjoint equations (Efthimiou et al., 2018b; Keats et al., 2007b; Kumar et al., 2015b) and is capable of providing more accurate dispersion predictions than a Gaussian dispersion model, with a mild increase in computational resources. However, the prediction accuracy of RANS models for the time-averaged flow fields around buildings has been shown to be insufficient when compared with LES (Tominaga et al., 2008a), which undermines the accuracy of their adjoint equation simulations and source term estimation (STE). Xue et al. (2018b) noticed this defect and suggested that the coupling of a RANS-like simulation of adjoint equations with the time-averaged flow of LES can improve the accuracy of STEs. In these RANS and RANS-like simulations, turbulent diffusion is always modeled according to the GDH, regardless of whether the mean flow field was produced by a RANS model or LES. However, the validity of the GDH on dispersion remains unclear, especially across complex urban terrains.

It has been confirmed that the dispersion prediction accuracy of RANS models operating according to the GDH is limited when compared with that of explicit models with LES (Tominaga and Stathopoulos, 2012, 2011b). Until now, the LES of adjoint equations has been regarded as impractical because the embedded adjoint equation represents an inverse dispersion process, such that the time-series flow field data of the entire domain must be produced by forward simulation and stored in advance to realize the unsteady simulation. One important challenge is that the volume of data acquired by LES seems too large for practical application. Nevertheless, this problem can be solved by the application of the compression method. It has been shown in **Chapter 4** that the wavelet-based compression method can now compress data to ~100 times their original size, while simultaneously conserving their accuracy.

Another limitation is the heavy calculation burden brought by LES of highly-complicated real urban areas and data compression. Considering that the source-receptor relationship can be constructed as databases long before the emergent injection of hazardous materials into the atmosphere, the requirement for more calculations caused by this method is not a problem. Indeed, it enables the real-world application of the LES of adjoint equations. Therefore, this chapter evaluates the applicability of STE with LES of the adjoint equation in an urban-like experiment.

5.2 Simulations for adjoint equation

Before the Bayesian inference, it is necessary to clarify the difference between RANS and LES of adjoint equation simulation.

5.2.1 Reynolds averaged Navier-Stokes simulation of the adjoint equation

The RANS simulation for the adjoint equation has been introduced in **Chapter 2**. It can be noticed that Eq. (2.24), RANS form of the adjoint equation, only explicitly resolves the mean flow advection, so the effects of turbulent diffusion, $\overline{\mathbf{u}'C^{*'}}'$ (where $'$ denotes the temporal fluctuations of variables) must be modeled by the mean field. One simple way of achieving this is through the GDH (Combest et al., 2011):

$$\overline{\mathbf{u}'C^{*'}} = -D_t \frac{\partial(\overline{C^*})}{\partial \mathbf{x}} = -\frac{\nu_t}{Sc_t} \frac{\partial(\overline{C^*})}{\partial \mathbf{x}} \quad (5.1)$$

where D_t is the turbulent diffusivity, ν_t is the eddy viscosity, which can be estimated based on the mean velocity and Reynolds stresses, and Sc_t is the turbulent Schmidt number. However, there are at least three problems in this GDH approximation. First, it assumes that the turbulent scalar flux is aligned with the mean scalar gradient, which is invalid under anisotropic turbulence. Secondly, ν_t is estimated by the mean velocity and Reynolds stresses predicted by the RANS model, which has been shown to inaccurately predict reality (Tominaga et al., 2008a). Finally, Sc_t is usually set as a global constant in the domain for simplicity, but it actually spans a large range of spatial differences (Combest et al., 2011; Tominaga and Stathopoulos, 2012). Thus, there remains no satisfactory way to model the spatial distribution of Sc_t .

Following Tominaga and Stathopoulos (2012), who clarified that accurate predictions of mean flow fields can improve the estimation of ν_t and further improve RANS predictions of dispersion fields, Xue et al. (2018) conducted LES-based forward simulations to obtain a time-averaged

1 velocity field and used it in a RANS-like (hereafter, RANS) simulation for adjoint equations. In this
2 case, the eddy viscosity was calculated as:

$$v_t = -\frac{\sum_{i,j} \overline{u'_i u'_j} S_{ij}}{\sum_{i,j} 2S_{ij}^2} \quad (5.2)$$

3 where S_{ij} is the main strain rate:

$$S_{ij} = \frac{1}{2} \left(\frac{\partial \bar{u}_i}{\partial x_j} + \frac{\partial \bar{u}_j}{\partial x_i} \right) \quad (5.3)$$

4 The mean velocity, Reynolds stresses, and strain rate were obtained by the forward LES. The
5 accuracy of STE was confirmed to have improved because LES can predict the mean flow field in
6 urban areas more accurately than RANS models. Despite this, the RANS simulation of adjoint
7 equations is still imperfect due to the existence of the other two problems related to the GDH. In
8 real applications, this simplification may cause significant errors in the modeling concentration, \mathbf{R} ,
9 especially in complex urban areas with many buildings.

10

11 5.2.2 Large eddy simulation of the adjoint equation

12 In contrast, LES appears to be a promising choice since it can explicitly resolve most of the
13 turbulent diffusion with the following adjoint equation.

$$-\frac{\partial(\widetilde{C}^*)}{\partial t} - \widetilde{\mathbf{u}} \frac{\partial(\widetilde{C}^*)}{\partial \mathbf{x}} - \frac{\partial}{\partial \mathbf{x}} \left([D_{sgs} + D_m] \frac{\partial(\widetilde{C}^*)}{\partial \mathbf{x}} \right) = \delta(\mathbf{x} - \mathbf{x}_m) \quad (5.4)$$

14 Here \widetilde{C}^* represents the grid-scale value of the variables. The sub-grid scale (SGS) turbulent
15 diffusion is modeled by:

$$\widetilde{\mathbf{u}' C^{*'}} = -D_{sgs} \frac{\partial(\widetilde{C}^*)}{\partial \mathbf{x}} = -\frac{v_{sgs}}{Sc_{sgs}} \frac{\partial(\widetilde{C}^*)}{\partial \mathbf{x}} \quad (5.5)$$

16 where v_{sgs} is the SGS eddy viscosity and Sc_{sgs} is the SGS turbulent Schmidt number. When the
17 grid is fine enough, most of the turbulent diffusion is directly calculated and the modeled SGS is
18 very small. Compared with RANS models, LES can more accurately simulate adjoint equations.

19 The main limitation to the practical application of LES is that since the adjoint tracer dispersion
20 field is a reverse simulation with $-\widetilde{\mathbf{u}}(\mathbf{x}, t)$, simulating the unsteady dispersion field requires saving
21 all of the $\widetilde{\mathbf{u}}(\mathbf{x}, t)$ and D_{sgs} data for the entire domain and across all timesteps in advance.

1 Furthermore, LES always requires a fine mesh and short timesteps to ensure the accuracy and
2 convergence of simulations, meaning that the data saved for the LES of adjoint equations require
3 such massive storage volumes that it renders the method impractical. As a proposed solution, the
4 compression database in **Chapter 4** is used for an LES of adjoint equations.

6 **5.3 Case description**

7 As for the study case, the block-arrayed urban model in **Chapter 4** is applied to represent the
8 geometry of a real urban neighborhood. This model has been frequently used in previous studies on
9 airflow (Cheng and Castro, 2002; Uehara et al., 2000; Xie and Castro, 2006) and the dispersion
10 characteristics of atmospheric pollutants (Branford et al., 2011b; Coceal et al., 2014b; Tominaga
11 and Stathopoulos, 2012) in urban areas. It includes complicated unsteady turbulent structures that
12 challenge the robustness of STE methods. Moreover, there is an open WTE database with the same
13 model that can be used as the measurements of concentration *D*.

14 A forward simulation with an LES model was utilized to predict the unsteady flow field, which
15 was compressed and stored in each time step. The source–receptor relationship was then obtained
16 by the unsteady simulation of adjoint equations using compressed flow data. To evaluate the
17 improvement made to the estimation accuracy of the proposed method, the RANS simulation of
18 adjoint equations was also performed with the time-averaged flow field of the LES. Bayesian
19 inference of the source was conducted twice based on the RANS and LES models.

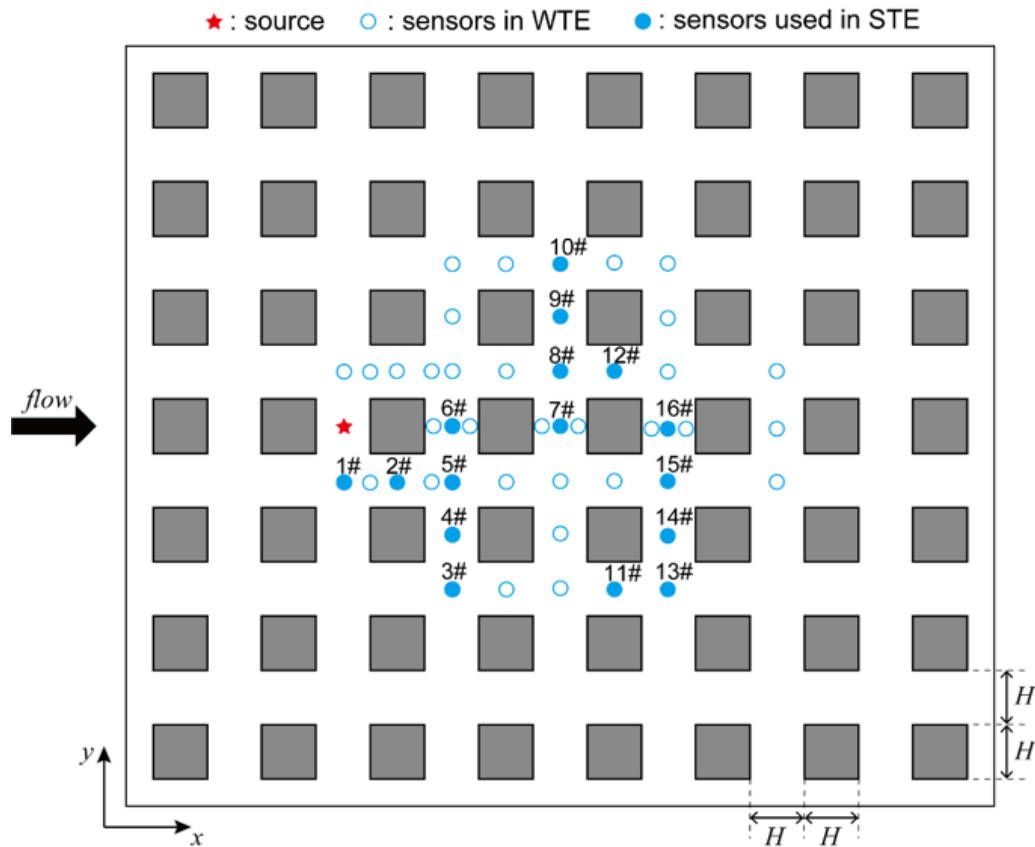
21 **5.3.1 Wind tunnel experiment**

22 In the wind tunnel experiment, a continuously releasing point source (a hole with a diameter
23 of 2 mm) was placed at the ground with coordinates of $(4H, 7H, 0)$. Pure ethylene gas (C_2H_4) was
24 released at a flow rate of $Q=0.216$ L/min from the source. Because the gas was released in the
25 wake region of one block and transported downstream via complex turbulence, the setting was
26 identified as a difficult scenario for the STE.

27 Velocities and concentrations were measured at different heights ($H/15, H/6, H/2, 5H/6, 7H/6$)
28 following the same symmetrical horizontal sensor configuration shown in **Fig. 5.1**. The wind
29 velocity was measured with split-film probes and the concentration was measured with a rapid-

1 response flame ionization detector. The measurements were sampled every 120 s in the frequency
 2 of 1000 Hz. Here, considering the symmetrical characteristics of the sensor configuration and the
 3 calculation burden of simulating adjoint equations, the time-averaged measurements of only 16
 4 sensors in the horizontal plane with a height of $0.5H$ were used to construct the vector \mathbf{D} (solid
 5 cycles in **Fig. 5.1**).

6



7

8 Figure 5.1. Schematic of sensor configuration (horizontal plane with $z = H/2$).

9

10 **5.3.2 Simulation settings**

11 In order to evaluate the improvements on STE accuracy brought by LES of the adjoint equation,
 12 two inferences based on RANS and LES of the adjoint equation are conducted and compared. One
 13 of the databases in **Chapter 4** for which the data volume was compressed 10 times to ensure the
 14 error $\frac{\max|f-\tilde{f}|}{\max|f|} < 10^{-5}$ is selected. The original volume of $240T$ data was ~ 24 TB. The compression
 15 reduced this size significantly to ~ 2.4 TB, which made it much easier to handle, even with a portable

1 hard disk.

2 The RANS simulation of adjoint equations followed Eqs. (2.24, 5.1-5.3) with a global turbulent
 3 Schmidt number $Sc_t = 0.7$. All time-averaged values were obtained from the 240T original LES.
 4 The simulated ensemble-averaged concentration of the adjoint tracer was regarded as the modeling
 5 concentration \mathbf{R} . The unsteady adjoint concentration distribution was predicted by LES according
 6 to Eq. (5.4-5.5), where $\tilde{\mathbf{u}}$ and D_{sgs} were the instantaneous values of the compressed LES flow
 7 field at each time step and Sc_{sgs} was set as 0.7 for the entire domain. Of the total 240T of
 8 compressed flow field data, the first 120T were used to initialize the simulation of the adjoint
 9 equation. The time-averaged C^* distribution of the second 120T was then sampled as the modeling
 10 concentration \mathbf{R} .

11

12 5.3.3 Bayesian inference settings

13 Bayesian inference was implemented twice based on the different modeling concentrations, \mathbf{R} ,
 14 obtained from the RANS and LES models of adjoint equations. The measurements, \mathbf{D} , were
 15 collected in the wind tunnel experiment. The variance of errors is set as $\sigma_{d,i}^2 + \sigma_{m,i}^2 = r * D_i$, where
 16 $r = 0.2$. Other settings are the same as inference introduced in **Chapter 2**.

17

18 5.4 Results and discussions

19 Here, all the results were nondimensionalized. The coordinate and length values were
 20 nondimensionalized by H . The velocity results were nondimensionalized by U_r . The concentration
 21 results were nondimensionalized by the standard concentration $C_r = C_{gas}Q/(U_r h_r^2)$. Here, C_{gas} is
 22 the concentration of gas in the injected flow from the source in the wind tunnel experiment, Q is
 23 the injected flow rate, and h_r is the reference height, which was set as $3.33H$.

24

25 5.4.1 Flow fields of forward simulations

26 In the LES flow fields of the forward simulations, the velocity results were nondimensionalized
 27 by the reference velocity, U_r . The mean velocity fields in several planes are shown in **Fig. 5.2**. The
 28 horizontal velocity distribution shows the symmetry caused by the regular block-arrayed

1 configuration. The most dominant velocity is the streamwise velocity in open street areas and the
 2 separation between the windward sides of buildings and circular vortices in wake regions can be
 3 distinguished through the spanwise velocity distribution. Apart from this, the spanwise velocity in
 4 most areas is nearly zero. In the vertical plane of open street regions ($y = 6H$, **Fig. 5.2(c, d)**), the
 5 mean velocity is layered such that the streamwise velocity at the same height does not change with
 6 the x -coordinate and the vertical velocity is nearly zero. There is essentially no momentum exchange
 7 in the vertical direction. This averaged flow field was used in the RANS simulation of adjoint
 8 equations for all sensors as the conventional method.

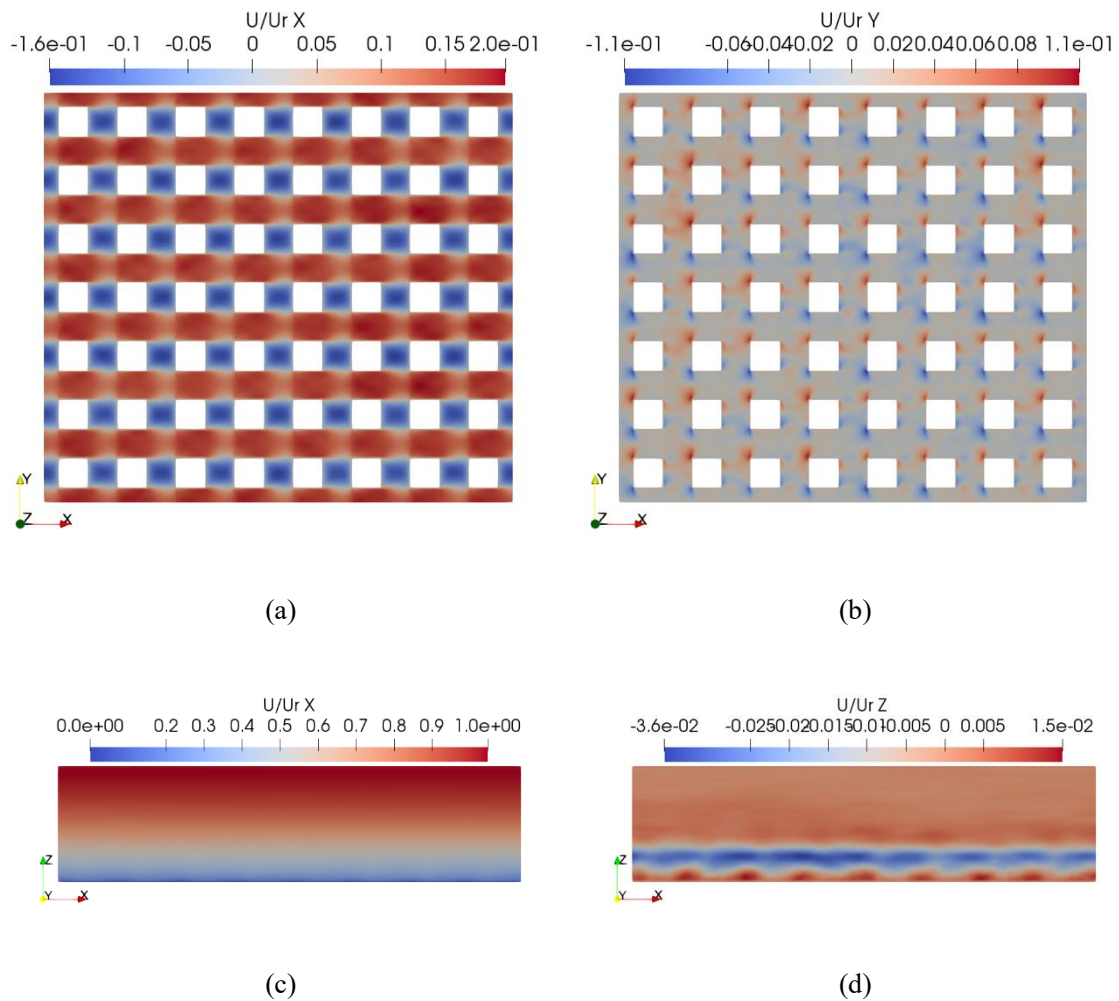


Figure 5.2. Mean flow fields produced by the forward LES. (a) streamwise velocity in a horizontal plane; (b) spanwise velocity in a horizontal plane; (c) streamwise velocity in a vertical plane; (d) vertical velocity in a vertical plane. Horizontal: $z = H/20$; vertical: $y = 6H$.

9 The validation profiles of the streamwise velocity can be checked in **Fig. 4.4 & 4.5**. Since the
 10 compression error was very limited in this study, it is reasonable to assume that the velocity field in

1 the LES of adjoint equations does not contribute greatly to the overall modeling error.

2

3 5.4.2 Comparison of adjoint concentration

4 The adjoint concentration fields simulated by RANS and LES models are compared. **Fig. 5.3**
 5 shows the mean distribution of the adjoint tracer released from sensor No. 1, which was located in
 6 the open street region between buildings ($x = 4H, y = 6H, z = 0.5H$). The horizontal plane is shown
 7 as the lowest plane in the domain, which contains the true source and can demonstrate its adjoint
 8 concentration. The vertical plane contains the sensor.

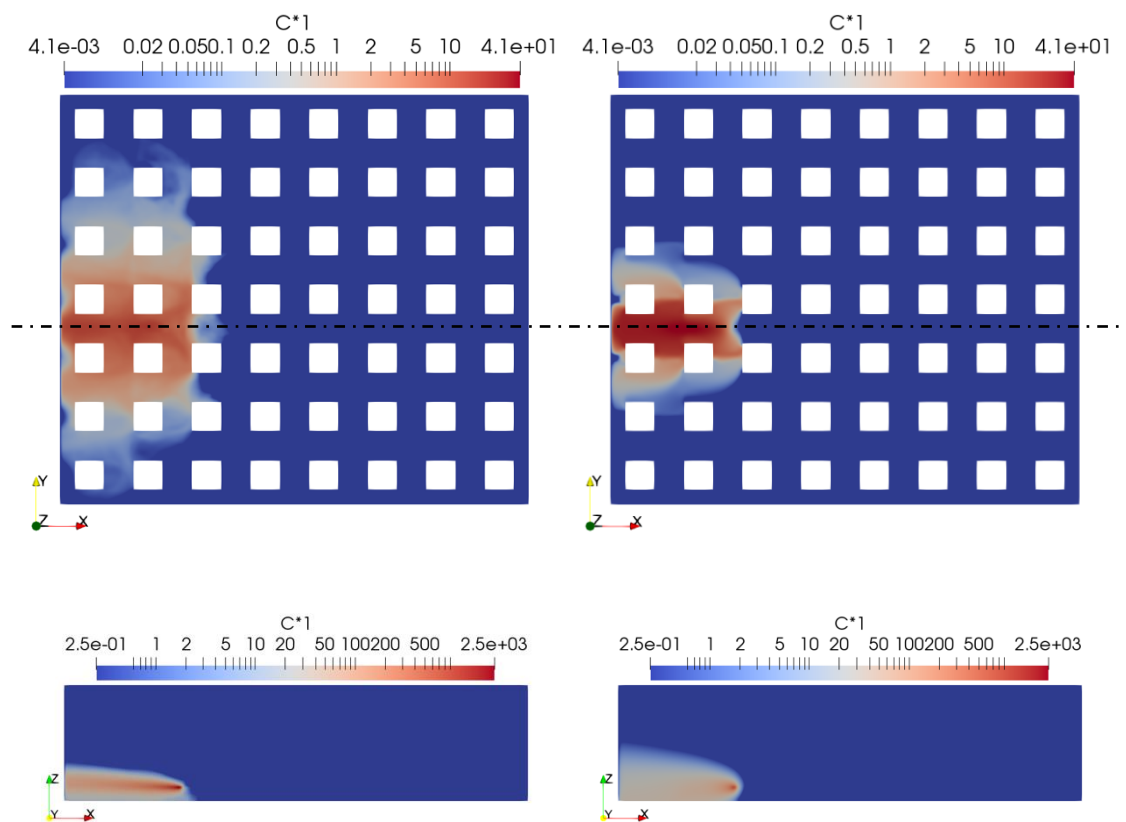


Figure 5.3. Mean adjoint concentration field of sensor No. 1 located in the open street region simulated by different methods: (left) LES and (right) RANS. (Top) Distribution in the horizontal plane ($z = H/20$); (bottom) distribution in the vertical plane ($y = 6H$).

9 In the RANS simulation of adjoint equations, it was confirmed that the dispersion area in the
 10 horizontal plane was limited within $5H$ in width, which is less than that of the LES results. Most of
 11 the adjoint tracers were clustered along the open street where the sensor is located. The GDH in the
 12 RANS model makes the turbulent diffusion proportional to the gradient of the mean concentration,

1 which is only true for small turbulent structures that can be simplified as isotropic. However,
2 according to the mean flow shown in **Fig. 5.2**, the advection of the tracers was dominated by the
3 streamwise velocity and separation flow, which constructed large-scale flow structures stretched in
4 the streamwise direction. The strong anisotropy of the flow fields invalidated the efficiency of the
5 GDH, making the turbulent diffusion also follow this large structure and causing simulation errors
6 in the adjoint equations.

7 As most tracers are transported along the streamwise direction, dispersion in the spanwise
8 direction is very weak. Even though the separation flow caused by the front walls of blocks can suck
9 tracers into the wake region in the reverse flow, its strength is limited when compared with the
10 streamwise velocity. The asymmetry of wake flow structures prevents the dispersion from
11 expanding into the other half of the wake. Moreover, the mean velocity field is layered in the vertical
12 direction, so the transport of the adjoint tracer in this direction can only rely on the GDH and
13 molecular diffusion. Hence, in the vertical plane, the concentration distribution of the RANS model
14 demonstrated a higher level of isotropy than the LES, where the distribution was stretched more in
15 the streamwise direction and squeezed downward. Consequently, fewer tracers were moved to the
16 bottom boundary in the RANS simulation. When the source was far away from a sensor in the
17 spanwise direction or they were at different heights, it would not be notably impacted from the tracer
18 and the adjoint relationship would be incorrectly constructed. In contrast, the situation changed in
19 the LES of adjoint equations. Because most of the turbulent diffusion was explicitly resolved, the
20 dispersion of adjoint tracers was wide throughout the domain in both the horizontal and vertical
21 directions. Instantaneous turbulent flows carried the tracers emitted from the sensor across several
22 wake regions.

23 **Fig. 5.4** shows the mean distribution of the adjoint tracer released from sensor No. 4 in the
24 wake region behind a building ($x = 6H, y = 5H, z = 0.5H$). In the RANS model, compared with the
25 sensor in the open street area, the dispersion area widened and the concentration in the lowest
26 horizontal plane increased. The time-averaged flow structures in the wake region were more
27 effective for dispersion than the layered flow in the open street region. Nevertheless, the dispersion
28 was not able to expand across two wake regions in the spanwise direction; instead, it clustered
29 around the blocks in the same row as the sensor. Thus, the prevention of dispersion caused by
30 asymmetric wake flows remained. As before, in the LES, the adjoint tracers were transported over
31 a larger area by the explicit turbulent diffusion.

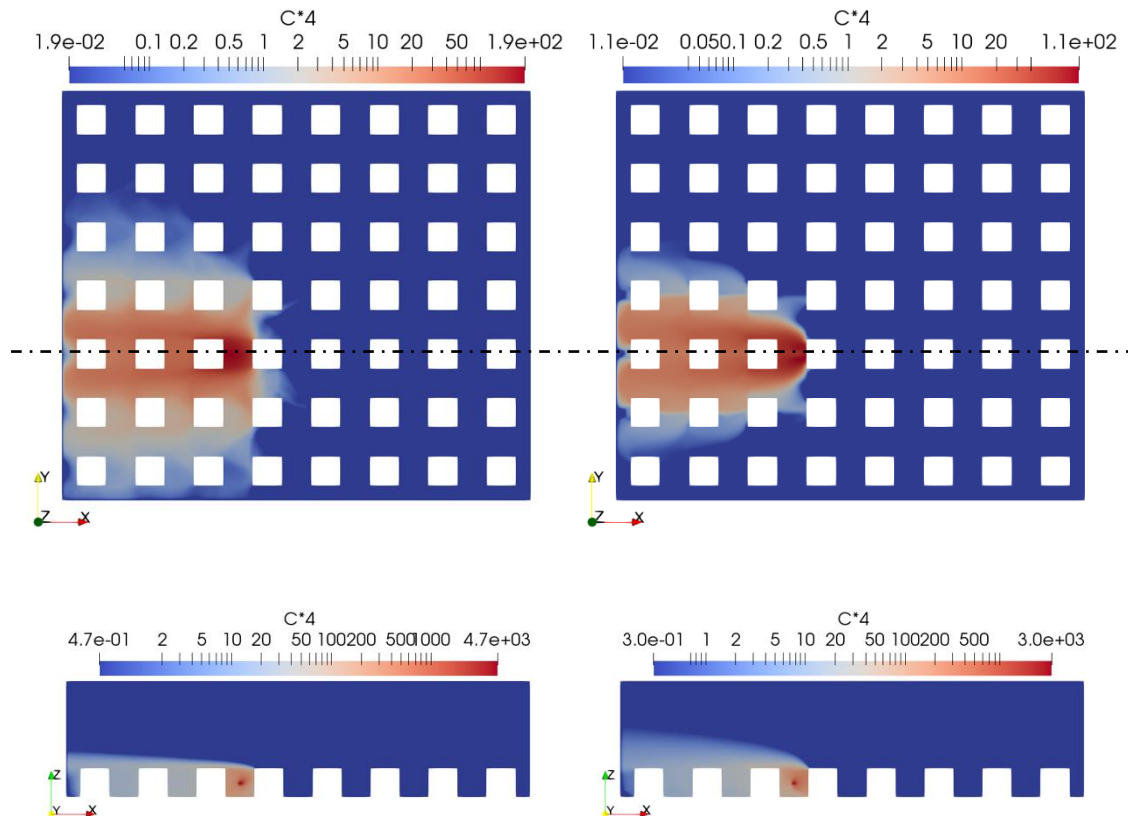
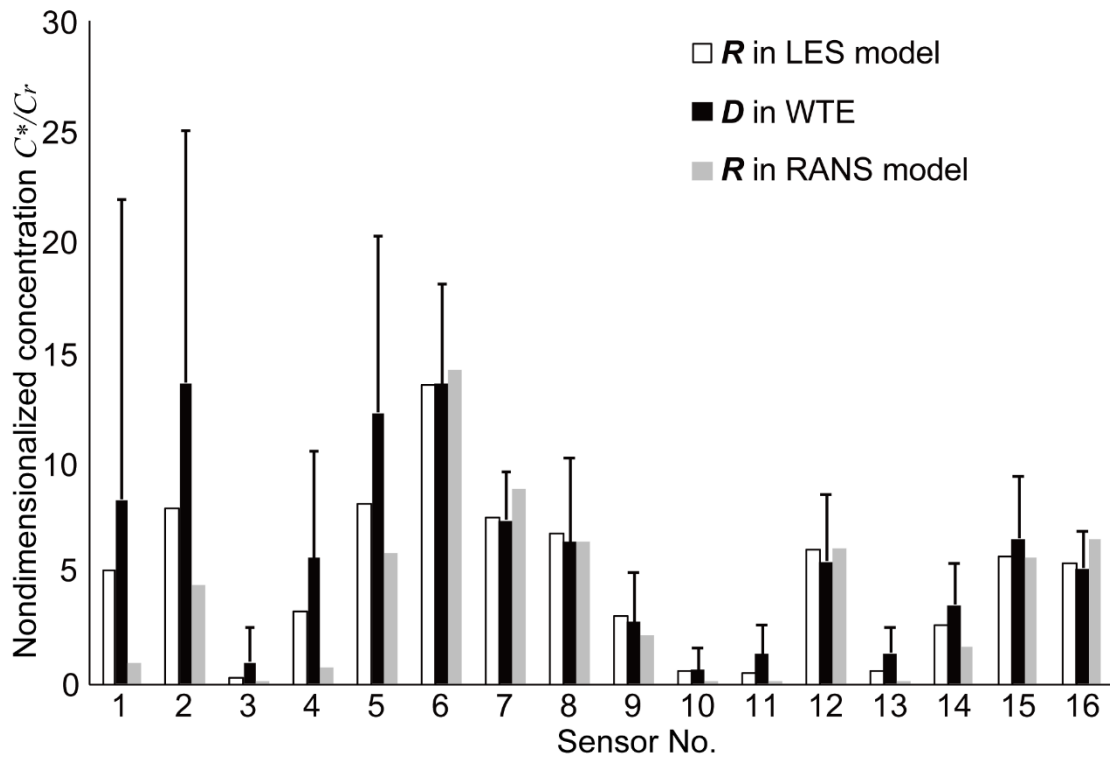


Figure 5.4. Mean adjoint concentration field of sensor No. 4 located behind a building simulated using different methods: (left) LES and (right) RANS. (Top) Distribution in the horizontal plane ($z = H/20$); (bottom) distribution in the vertical plane ($y = 5H$).

1 The measurements, \mathbf{D} , from the wind tunnel experiment are compared with the modeling
2 concentration, \mathbf{R} , of the RANS and LES models in **Fig. 5.5**. According to the adjoint relationship
3 shown in Eq. (2.23), the \mathbf{R} of each sensor in **Fig. 5.5** is the nondimensionalized concentration of
4 adjoint tracers emitted by that sensor at the true source's location. We confirmed that the LES
5 modeling concentration was closer to the experimental data than that in the RANS model because
6 most of the adjoint tracers were transported along the streamwise direction and the spanwise
7 diffusion was insufficient in the latter model. For the sensors placed at large distances from the true
8 source in the y -direction, the RANS model underestimated the concentration measurements.
9 Meanwhile, for sensors placed in the same row as the true source, concentrations were
10 overestimated due to the over-concentration of the adjoint tracers. Therefore, regardless of whether
11 the sensor was in an open street or wake region, simulating turbulent diffusion is critical to
12 dispersion.



1

2 Figure 5.5. Comparison between the measurements, \mathbf{D} , in the wind tunnel experiment (WTE) and
 3 modeling concentration, \mathbf{R} , of the RANS and LES models. Error bars denote the standard
 4 deviation of \mathbf{D} in the WTE.

5 The numerical model employed to simulate adjoint equations requires special attention for
 6 complex urban areas, as in the case study. In general, the LES performed better than the RANS
 7 model for most sensors even though there were also discrepancies when compared to the
 8 experimental data. One possible reason for this result could be due to the fact that the RANS Sc_t
 9 was still set as a global constant and ν_t was modeled under the assumption of isotropic turbulence
 10 in the current case. Thus, deeper research is still needed to advance the modeling of Sc_t and ν_t .

11

12 5.4.3 Source term estimation results

13 We compared the estimation results of Bayesian inference using the modeling concentrations,
 14 \mathbf{R} , of the RANS and LES models of adjoint equations. **Fig. 5.6** presents the marginal probability
 15 distribution of each source parameter. In the results of the proposed method based on the LES of
 16 adjoint equations, the estimation of each parameter was accurate. The peak values of the x - and y -
 17 coordinates were close to the true ones and the deviations were smaller than $0.5H$. The strength of
 18 the source was also well-defined since the peak value was almost the same as the true one, and the

- 1 probability density function (PDF) was concentrated around the true value, similar to a normal
- 2 distribution with a small standard deviation.

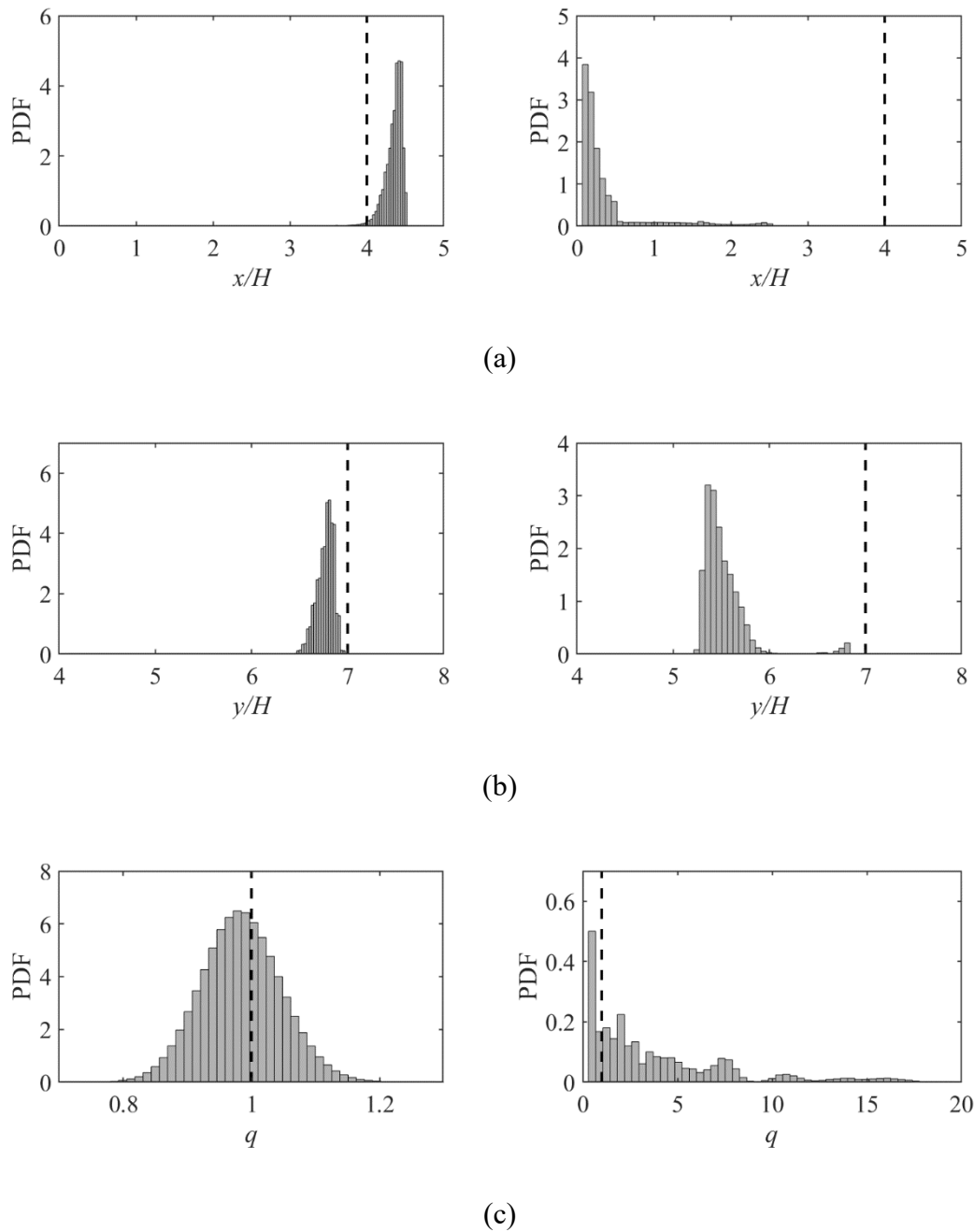


Figure 5.6. Marginal probability distributions of single-source parameters (a: x -coordinate; b: y -coordinate; c: strength) estimated by Bayesian inference using different adjoint concentration fields: (left) LES and (right) RANS. The dashed line represents the true value.

- 3 The estimation performance with RANS modeling of the adjoint equations was worse than that
- 4 of the LES. The PDFs of the x - and y -coordinates greatly deviated from true values. The zero

1 probability of estimating the true value indicates that the Markov chain did not explore this area or
2 did not credit it after exploration; both situations mean that the probability of the true value was too
3 low because the difference was too great between the modeling concentration and measured value.
4 The estimated PDF of the strength revealed even more limitations of the conventional method. In
5 spite of the fact that the peak value was close to the true strength, the distribution was so wide that
6 it rendered the estimate invalid. There was also a second peak in the distribution far from the true
7 value. In a real-world application, with more complicated measurement noise, it is possible that the
8 width of the distribution would increase and the wrong peak would be identified as the dominant
9 one, which may hinder subsequent risk management.

10 The main reason for the estimation failure of the RANS model was attributed to the inaccurate
11 prediction of adjoint concentration fields, where the turbulent diffusion was modeled according to
12 the GDH. The difference between the measurements and modeling concentrations made it difficult
13 to identify the true parameters via Bayesian inference. It is worth noting that, although the RANS
14 modeling concentration errors may seem tolerable when compared with the variance within the
15 measured values, the resulting estimation error is considerably larger. The reason for this is that in
16 sampling the posterior distribution, the modeling concentration is coupled with the inference input
17 coefficient, $\sigma_{d,i}^2 + \sigma_{m,i}^2$. As the optimal selection of this variance remains unknown, this
18 amplification could be inefficient in practical applications. For modeling concentrations that are not
19 accurate enough, the inference is vulnerable to the selection of the input coefficient. When
20 confronting complex measurements in reality, establishing an accurate modeling concentration is
21 key to effective STE before systematic research into the selection of the input coefficient is
22 completed.

23 The improvements made to the accuracy of the STE yielded by the proposed model is better
24 demonstrated by the joint probability distribution, $p((x, y)|D, I)$, of the two methods (**Fig. 5.7**). We
25 confirmed that the estimated source locations based on the LES model were very close to that of the
26 true source. The inference narrowed the probability into the same wake region and the distance
27 between the estimated location and the true one was smaller than $0.5H$. In contrast, the estimation
28 based on the RANS model was inaccurate. The joint probability clustered around a small area that
29 was approximately two columns away from the true location, which suggests the limited credibility
30 of the conventional RANS-based method.

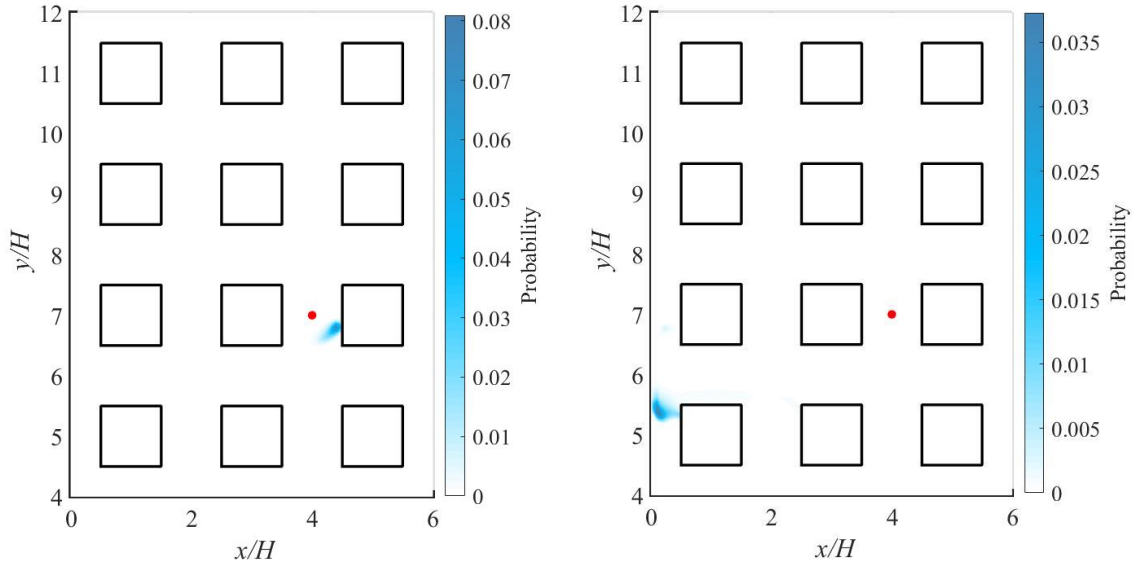


Figure 5.7. Joint probability distribution of source locations estimated by Bayesian inference using different adjoint concentration fields: (left) LES and (right) RANS. The red point represents the true source.

1 To quantify the improvement to STE achieved by the LES modeling of adjoint equations, we
 2 selected the 50th percentile values of the PDFs shown in **Fig. 5.6** as the estimated results of the two
 3 methods and defined the following two indices.

4 The location error E_d is defined as the expectation of the Euclidian distance between the
 5 estimated location and the true location, which was calculated using the following equation:

$$E_d = \int \int p(x, y | \mathbf{D}, \mathbf{I}) \sqrt{(x - x_s)^2 + (y - y_s)^2} dx dy \quad (5.6)$$

6 where (x_s, y_s) is the true location of the source.

7 The strength error E_q is the expectation of the normalized difference between the estimated
 8 strength and true strength, which is expressed as follows:

$$E_q = \int p(q | \mathbf{D}, \mathbf{I}) \frac{|q - q_s|}{q_s} dq \quad (5.7)$$

9 where q_s is the true strength of the source.

10 These indices are summarized in **Table 5.1**. We found that the LES model reduced the errors
 11 of location and strength estimates by 89% and 99%, respectively, when compared with the RANS
 12 model.

Table 5.1. Summarized estimation results of the RANS and LES models.

Method	Location		Strength	
	(x_s, y_s)	E_d	q_s	E_q
True value	$(4H, 7H)$		1	
Steady adjoint equation (RANS)	$(0.21H, 5.46H)$	$4.09H$	2.5	1.5
Unsteady adjoint equation (LES)	$(4.38H, 6.78H)$	$0.44H$	0.98	-0.02

1

2 5.5 Conclusions

3 In this chapter, to improve the accuracy of STE in complex urban applications, a new method
4 was developed based on Bayesian inference coupled with unsteady adjoint equation modeling via
5 LES. The performance of the proposed method was evaluated through a regular, block-arrayed
6 urban model with the continuous dispersion of a point source in the wind tunnel experiment. The
7 LES approach was applied to predict the flow fields of the urban model and the full spatial
8 distribution and time-dependent dynamics were saved for the simulation of adjoint equations. To
9 relieve the pressure imposed by storing such a large amount of data, the wavelet-based compression
10 method was employed to compress the original data to 10% of its original volume. The source-
11 receptor relationship obtained from the LES model of the adjoint equations and the measurements
12 from 16 sensors in the wind tunnel experiment were used as inputs for Bayesian inference to
13 calculate the posterior distribution of the source term. To clarify the improvements made to the
14 accuracy of STE using the proposed method, an existing RANS-based method of simulating adjoint
15 equations with a time-averaged LES flow field was conducted for comparison.

16 The results showed that the modeling of the adjoint equation was significantly improved by
17 the new LES model. The proposed method reflected the crucial effects of an explicit resolution of
18 turbulent diffusion with time-series data. The modeling concentration, \mathbf{R} , of the LES model was
19 closer to the measurements, \mathbf{D} , than that of the RANS model. Since turbulent diffusion was
20 modeled by the GDH in the RANS method, the adjoint concentration predicted with the mean flow
21 field did not reflect the true source-receptor relationship well. The different degrees of accuracy of
22 the modeling concentrations between the two methods further resulted in different estimated

1 posterior probabilities of the source parameters. For the RANS model, neither the coordinates nor
2 the strength of the point source was successfully identified. In contrast, the LES model provided
3 much more accurate estimates, which were robust to variations in noise, and the estimated location
4 was in the same wake region of the true source. The errors in estimates of location and strength were
5 reduced by 89% and 99%, respectively, using the LES model of adjoint equations when compared
6 with the RANS approach.

7 It must be noted that the regular, block-arrayed model employed in this chapter is an idealized
8 scenario of real urban geometry. Even in this simple case, the difference between the adjoint tracer
9 concentration predicted by the RANS and LES models was large enough to significantly affect the
10 STE results. The critical factor was the simulation of dispersion in the adjoint equations. Even
11 though the RANS model imposed lower computational costs, it was demonstrated that the GDH
12 could not satisfy the accuracy requirements in strongly anisotropic flow fields. In real urban areas,
13 the situation is much more complicated because of the diverse geometries and configurations of
14 buildings and equipment. Hence, to effectively perform STE, it is beneficial to spend more time
15 explicitly resolving turbulent diffusion for the source–receptor relationship via LES until more
16 accurate numerical models for turbulent diffusion are proposed.

17 Even though unsteady LES of adjoint equation was realized in the current research, the time-
18 averaged adjoint relationship was used in the inference instead of time-series data. The reason is
19 that only time-averaged measured concentration data is available in the case study. Besides,
20 advanced measurement equipment that can record concentration in a time-series form is still
21 unpopular due to high cost. Existing research in the literature is still mainly based on mean
22 concentration measurement. Hence, for the convenience of comparison, this research also used time-
23 averaged simulated results. Although it seems to be a waste of calculation resources to only use
24 mean value of unsteady simulation, it was confirmed that mean results of the unsteady simulation
25 outperform that of the steady simulation.

26 However, it must be admitted that time-series data contains much more useful information to
27 STE than time-averaged data. The performance of STE will be improved if the time-series data can
28 be effectively utilized, which is also helpful to estimations considering multiple source or response
29 time. Therefore, it is necessary to develop appropriate ways to fully use unsteady adjoint
30 relationships.

31

1 Symbols

C	: the concentration field
C^*	: the adjoint concentration field
C_r	: standard (reference) concentration
C_{gas}	: the emission strength of the source
D	: the measurements vector
D_i	: the measurement of the sensor with index i
D_m	: molecular diffusivity
D_{sgs}	: sub-grid scaled turbulent diffusivity
D_t	: turbulent diffusivity
E_d	: the expectation of the Euclidian distance between the estimated location and the true location
E_q	: the expectation of the normalized difference between the estimated strength and true strength
f	: the original data of a three-dimensional scalar field in the Cartesian indexing mesh grid
\check{f}	: the decompressed flow field data
H	: the edgy length of blocks in the simulation (=60 mm)
h_r	: reference height
I	: the background information for Bayesian inference

- $p(A|B)$: conditional probability of event A occurring given that B is true
- Q : the gas flow rate at the source
- q : strength samplings produced in MCMC for the point source
- q_s : the true strength of the point source
- r : the ratio between the specified error covariance in the Bayesian inference and the measurements
- S_{ij} : mean strain rate of i and j directions
- Sc_{sgs} : sub-grid scaled turbulent Schmidt number
- Sc_t : turbulent Schmidt number
- T : standard time-scale of LES simulation
- U_r : reference speed
- \mathbf{u} : velocity field
- \mathbf{x}_m : coordinates of sensors
- (x, y) : location samplings produced in MCMC for the point source
- (x_s, y_s) : the true location of the point source
- $\delta(\cdot)$: Dirac delta function
- ν_{sgs} : eddy viscosity coefficient at the sub-grid scale
- ν_t : eddy viscosity
- $\sigma_{d,i}^2$: the variance of error in the measurement of the sensor with index i

$\sigma_{m,i}^2$: the variance of error in the modeling concentration for the sensor with index i

$\sigma(\cdot)$: standard deviation

$(\cdot)'$: fluctuation value with time

$\bar{\quad}$: Reynolds average operator

$\tilde{\quad}$: filtering operator in LES

1

2

1

2

3

4

5

6

7

8

9

Chapter 6

10 Sensor configuration optimization
11 for source term estimation based
12 on the entropy of adjoint equation

13

1

2

3

4

5

6

Abstract

7

8

9 Until now, few studies have developed sensor configuration optimization methods aiming to
10 ensure good source term estimation (STE) performance in the monitoring area, such that the
11 posterior probability could aggregate around the truths while addressing most sources. This chapter
12 proposes a method by designing an objective function and applying a simulated annealing algorithm.
13 The objective function is set as the information joint entropy of the adjoint concentration. The
14 performance of the proposed method was assessed by Bayesian inference STE for 25 unknown
15 sources based on the obtained optimal configuration in a regular block-arrayed building group
16 model. The STE results were compared with those of uniform and random configurations.

17

6.1 Introduction

Source term estimation (STE) has three critical factors: measurements, source-receptor relationship, and estimating algorithm. In the previous research, most efforts have been devoted to the source-receptor relationship modeling and estimating algorithms. Many methods have been proposed to accurately simulate the predicted concentration and accelerate the estimation algorithm, while measurements have been neglected for a long time, and the related research is sparse. However, measurements are the foundation of STE. Theoretically, all STE methods would be useless if sensors could not efficiently measure the concentration information (Keats et al., 2010). In **Chapter 3**, it has been confirmed that some configurations may not provide sufficient measurement information for the STE of the line source, and the resultant estimations are inaccurate. In addition, the uniform and random sensor configurations used in wind tunnel experiments or field tests in previous research are impractical because of the irregular building distribution and the limitation of deployment cost in real applications. Random configuration may also cause considerable errors in the measurements of certain sources, which may be dominant over the truth.

Most existing research on sensor configuration optimization (SCO) design is still insufficient because of the special requirements of STE applications. Several studies used the difference between measurements and simulated concentrations to evaluate the sensor configuration. For example, a sensor network has been designed to monitor nuclear power stations in France (Abida et al., 2008; Saunier et al., 2009). Kouichi et al., (2019) proposed an optimization method to reduce the existing sensors for an urban model area. However, the simulated concentration in these methods needs prior knowledge of sources in advance, which is unavailable in STE applications. Other published SCO methods aim to monitor specific atmospheric pollutants such as ozone, rather than a specific source (Araki et al., 2015; Fuentes et al., 2007; Wu and Bocquet, 2011). Sources emitting ozone are too numerous to effectively identify, so their target is to estimate the global distribution using several point measurements. Keats et al. (2010) proposed a design method related to STE, but the research problem is to place an additional sensor as a supplementary for a fixed configuration after the dispersion emergency. The method cannot be used in the preparatory sensor configuration design before the emergency. Until now, one of the most applicable SCO methods for STE was proposed by Ngae et al. (2019), in which the sensor configuration is evaluated based on the entropic criterion of measurement information provided by sensors (Issartel, 2005b) without the prior knowledge of sources. However, this method has a deep relationship with the renormalization inversion theory, which belongs to the deterministic STE category. It is still necessary to develop other SCO methods

1 from the view of stochastic STE theory.

2 Therefore, in this chapter, an SCO method for the stochastic STE is proposed. This method has
3 the same basic idea as Ngae et al. (2019) that evaluates the quality of sensor configuration based on
4 the information entropy contained in sensors' measurements. Despite that, these two methods use
5 totally different ways to calculate the information entropy. Ngae et al. (2019) obtained the entropy
6 by the determinant of the weighted Gram matrix (Issartel, 2005b), which is an important element
7 transformed from adjoint concentration field in the renormalization inversion theory, while our
8 method calculated the entropy with the spatial distribution of adjoint concentration, which is
9 expected to be more intuitive and easy to understand. Besides, because entropy is an important
10 concept in stochastic theory, its physical meaning is closely related to stochastic STE and is
11 interpreted in this chapter.

12 The application scenario is assumed to be that a point source with constant release strength
13 may appear anywhere in the target domain with a statistically steady flow field. To monitor the
14 target area, the proposed method finds a user-specified number of sensors among candidates to
15 construct an optimal configuration. Regarding stochastic STE, the Bayesian inference method in
16 **Chapter 2** is applied. In an accurate Bayesian STE, the probability mass in the posterior probability
17 distribution should concentrate around the true value with a narrow distribution width. The optimal
18 configuration and resultant estimations are expected to be superior to others such as uniform,
19 random, and experience-based configurations.

20

21 **6.2 Entropy-based configuration optimization**

22 In this section, the proposed SCO method is introduced. It consists of two parts: the objective
23 function to evaluate the rank of any configuration and the simulated annealing (SA) algorithm to
24 efficiently identify the optimal configuration.

25

26 **6.2.1 Objective function**

27 The objective function is used to evaluate the measurement ability and steadiness of each
28 sensor configuration. It is proposed that all configurations can be ranked by the joint entropy of the
29 spatial probability distribution of adjoint concentration fields of their sensors. The optimal

1 configuration should have the largest entropy. The author first proves this mathematically and then
2 explains its physical meaning.

3 **a. Mathematical explanation**

4 Here is the definition of the information entropy of a probability distribution, $p(x)$ (Cover
5 and Thomas, 2006).

$$H(x) = - \int p(x) \log p(x) dx \quad (6.1)$$

6 In the statistics, entropy represents the information or uncertainties of $p(x)$. A distribution with a
7 narrow spread and a sharp peak possesses small entropy, meaning that the uncertainty of the
8 objective or information contained in the distribution is small. Similarly, a distribution with a wide
9 spread and slight slope has a large entropy. Naturally, a uniform distribution has the largest entropy,
10 information, and uncertainty (Fuentes et al., 2007).

11 For any unknown source \mathbf{s} , if we only have background information I without any sensors,
12 the prior probability distribution of the unknown source parameter is $p(\mathbf{s}|I)$. When we have n
13 sensors to monitor an area with the measurement vector $\mathbf{D} = (D_1, D_2, D_3, \dots, D_n)$, after the Bayesian
14 inference, the posterior probability distribution of the unknown source parameter becomes
15 $p(\mathbf{s}|\mathbf{D}, I)$. The information about unknown sources provided by these sensors can be quantified by
16 mutual information I_m , which is defined as

$$I_m(\mathbf{s}; \mathbf{D}|I) = H(\mathbf{s}|I) - H(\mathbf{s}|\mathbf{D}, I) \quad (6.2)$$

17 where $H(\mathbf{s}|I) = - \int p(\mathbf{s}|I) \log p(\mathbf{s}|I) d\mathbf{s}$ and $H(\mathbf{s}|\mathbf{D}, I) = - \int p(\mathbf{s}|\mathbf{D}, I) \log p(\mathbf{s}|\mathbf{D}, I) d\mathbf{s}$ are
18 the entropies of the corresponding probability distributions. The optimum sensor configuration
19 should provide the largest information $I_m(\mathbf{s}; \mathbf{D}|I)$, which means it should minimize $H(\mathbf{s}|\mathbf{D}, I)$
20 because $H(\mathbf{s}|I)$ is constant before and after measurement. According to the properties of the
21 entropy definition, $H(\mathbf{s}|\mathbf{D}, I)$ becomes smaller when $p(\mathbf{s}|\mathbf{D}, I)$ has a narrower distribution.
22 Therefore, the minimum $H(\mathbf{s}|\mathbf{D}, I)$ denotes the posterior probability density function (PDF) of
23 unknown source $p(\mathbf{s}|\mathbf{D}, I)$ concentrates around a certain value, which is expected to be close to
24 the true value when the Bayesian inference is sufficiently accurate.

25 To find \mathbf{D} , which minimizes $H(\mathbf{s}|\mathbf{D}, I)$, let's start with a simple scenario. Imagine the situation
26 that we determine (randomly or by experience) first as $n - 1$ sensors, fix them, and find the

1 optimal position for the next sensor e . The measurement vector can be divided into $\mathbf{D} =$
 2 (\mathbf{D}_{n-1}, D_n) . The probability that sensor e can provide the measurement D_n is denoted as

$$\begin{aligned} p(D_n|\mathbf{D}_{n-1}, I) &= \int p(D_n|\mathbf{s}, \mathbf{D}_{n-1}, I)p(\mathbf{s}|\mathbf{D}_{n-1}, I)d\mathbf{s} \\ &= \int p(D_n|\mathbf{s}, I)p(\mathbf{s}|\mathbf{D}_{n-1}, I)d\mathbf{s} \end{aligned} \quad (6.3)$$

3 When we assume that the measurement of each sensor is only determined by source \mathbf{s} and would
 4 not be affected by other sensors, \mathbf{D}_{n-1} can be removed from $p(D_n|\mathbf{s}, \mathbf{D}_{n-1}, I)$ on the right-hand
 5 side (RHS). This probability contains the real measurement error and the posterior guess of the
 6 source using the known data \mathbf{D}_{n-1} . The optimal selection of e should minimize the entropy
 7 $H(\mathbf{s}|\mathbf{D}, I)$ as follows:

$$H(\mathbf{s}|\mathbf{D}, I) = H(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I) = - \int p(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I) \log p(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I) d\mathbf{s} \quad (6.4)$$

8 Because D_n also has uncertainty, as shown in Eq. (6.3), the expectation of $H(\mathbf{s}|\mathbf{D}, I)$ can be
 9 evaluated by

$$E[H(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I)] = \int p(D_n|\mathbf{D}_{n-1}, I)H(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I)dD_n \quad (6.5)$$

10 The joint entropy has a chain rule (Cover and Thomas, 2006):

$$H(X, Y|Z) = H(X|Z) + H(Y|X, Z) = H(X|Z) + \int p(x)H(Y|x, Z)dx \quad (6.6)$$

11 If we let $X = \mathbf{s}$, $Y = D_n$, $Z = (\mathbf{D}_{n-1}, I)$ and notice that neglecting Z in the equation for a
 12 moment for simplicity would not change the final result because it appears in each term, we have:

$$H(\mathbf{s}, D_n) = H(\mathbf{s}) + \int p(\mathbf{s})H(D_n|\mathbf{s})d\mathbf{s} \quad (6.7)$$

13 By switching the position of \mathbf{s} and D_n , we can further obtain:

$$H(\mathbf{s}, D_n) = H(D_n) + \int p(D_n)H(\mathbf{s}|D_n)dD_n \quad (6.8)$$

14 The second term on the RHS of Eq. (6.8) is the same as the RHS in Eq. (6.5), the expectation of
 15 $H(\mathbf{s}|\mathbf{D}, I)$, and RHS of Eq. (6.8) also equals the RHS of Eq. (6.7). Then, we return $Z = (\mathbf{D}_{n-1}, I)$
 16 to obtain

$$\begin{aligned}
E[H(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I)] &= \int p(D_n|\mathbf{D}_{n-1}, I)H(\mathbf{s}|\mathbf{D}_{n-1}, D_n, I)dD_n \\
&= H(\mathbf{s}|\mathbf{D}_{n-1}, I) + \int p(\mathbf{s}|\mathbf{D}_{n-1}, I)H(D_n|\mathbf{s}, \mathbf{D}_{n-1}, I)d\mathbf{s} \\
&\quad - H(D_n|\mathbf{D}_{n-1}, I)
\end{aligned} \tag{6.9}$$

1 In the RHS of Eq. (6.9), the first term is the entropy of the posterior PDF of the STE by \mathbf{D}_{n-1} .
2 Because these sensors are fixed, the value is constant and independent of D_n . The second term
3 describes the influence of measurement noise on the data. $H(D_n|\mathbf{s}, \mathbf{D}_{n-1}, I)$ means the uncertainty
4 contained in the measurement signals when a source \mathbf{s} and a No. n sensor are selected. This
5 uncertainty is usually affected only by the measurement noise since both source and sensor are fixed.
6 If we assume that the noise is independent of the measurement location, the second term is also
7 constant (Loredo, 2004). Therefore, to minimize the left-hand side of Eq. (6.9), the third term must
8 be maximized.

9 Before further derivation, we first check the meaning of the third term $H(D_n|\mathbf{D}_{n-1}, I)$, which
10 is the entropy of probability $p(D_n|\mathbf{D}_{n-1}, I)$. According to Eq. (6.3), $p(D_n|\mathbf{D}_{n-1}, I) =$
11 $\int p(D_n|\mathbf{s}, I)p(\mathbf{s}|\mathbf{D}_{n-1}, I)d\mathbf{s}$. $p(\mathbf{s}|\mathbf{D}_{n-1}, I)$ indicates that measurements \mathbf{D}_{n-1} constrain the
12 estimation of an unknown source in a limited parameter space \mathbb{V} with certainty. When the unknown
13 source is located at different places in \mathbb{V} , it may cause different concentration measurements D_n
14 at sensor e , which is the meaning of the first term $p(D_n|\mathbf{s}, I)$.

15 Assuming that the adjoint equation is accurately simulated, according to the adjoint
16 relationship between the source and sensor, as shown in Eq. (2.24), we obtain:

$$D_n = \overline{C_s(e)} = q_s \overline{C_e^*(\mathbf{s})} \tag{6.10}$$

17 Therefore, $p(D_n|\mathbf{s}, I) = p(\overline{C_s(e)}|\mathbf{s}, I) = p(q_s \overline{C_e^*(\mathbf{s})}|\mathbf{s}, I)$.

18 To maximize $H(D_n|\mathbf{D}_{n-1}, I)$, $p(D_n|\mathbf{D}_{n-1}, I)$ should have as wide a distribution as possible;
19 consequently, $p(D_n|\mathbf{s}, I)$ or $p(q_s \overline{C_e^*(\mathbf{s})}|\mathbf{s}, I)$ should be wide, since $p(\mathbf{s}|\mathbf{D}_{n-1}, I)$ has been fixed
20 in Eq. (6.3). In the Bayesian inference for STE, q_s is sampled by the MCMC with a uniform prior
21 distribution. Meanwhile, in the process of location optimization for sensor e , changing the location
22 of e would only affect $\overline{C_e^*(\mathbf{s})}$, the spatial distribution of adjoint concentration corresponding to \mathbb{V} .
23 The wider this spatial distribution is, the wider $p(q_s \overline{C_e^*(\mathbf{s})}|\mathbf{s}, I)$ becomes. In other words, the
24 uncertainty that the unknown source causes different measurements in sensor e can be transferred
25 to the uncertainty of the spatial adjoint concentration resulting from e .

1 Returning to the optimal process, we use the chain rule of joint entropy (Eq. (6.6)) once more
 2 to obtain

$$H(\mathbf{D}|I) = H(\mathbf{D}_{n-1}, D_n|I) = H(\mathbf{D}_{n-1}|I) + H(D_n|\mathbf{D}_{n-1}, I) \quad (6.11)$$

3 Because the first term on the RHS of Eq. (6.11) is fixed and our target is to maximize the
 4 second term in the RHS, $H(\mathbf{D}|I)$ must be maximized. It is the joint entropy of probability of all
 5 measurements at fixed sensors (if e has been decided) caused by different sources \mathbf{s} . According
 6 to the discussion on the adjoint relationship above, this joint probability can be transferred to the
 7 joint spatial probability of adjoint concentrations at different places $p(\mathbf{R}|I)$ resulting from each
 8 sensor. Therefore, we find a sensor location e that makes the spatial distribution $p(\mathbf{R}|I)$ as wide
 9 as possible, and the corresponding entropy $H(\mathbf{R}|I)$ or $H(\mathbf{D}|I)$ is maximized.

10 In the above content, we proved the theorem that when $n - 1$ sensors are fixed, the best
 11 location of the remaining sensor should be the one that makes the entropy $H(\mathbf{R}|I)$ larger than the
 12 other places. Now we can say that the configuration which owns the maximum $H(\mathbf{R}|I)$ is the
 13 optimized one. If it is not, we have to adjust the sensors' locations to improve it further. However,
 14 for each sensor, no matter where we move it, $H(\mathbf{R}|I)$ will decrease since the original one has the
 15 maximum value. According to the above theorem, the new configuration is worse than the original
 16 one. There is no more space for further improvement. Therefore, it is reasonable to set our objective
 17 function as $H(\mathbf{R}|I)$ to evaluate the efficiency of the sensor configuration. The optimal
 18 configuration should have the largest $H(\mathbf{R}|I)$ values. In the designing process, we simultaneously
 19 determine the locations for all sensors to maximize $H(\mathbf{R}|I)$.

20 **b. Physical explanation**

21 The reason the optimal configuration should have the largest $H(\mathbf{R}|I)$ can be intuitively
 22 explained by **Fig. 6.1**. An unknown source, the red star, appeared in the target domain. To facilitate
 23 understanding, it is assumed in this example that the strength of the source q_s is already limited in
 24 a small interval and the STE is operated in a two-dimensional space, which indicate that we only
 25 need to estimate its position in a plane here.

26 If we do not have any sensors, nothing is known. After we deployed one sensor, it can provide
 27 measurements D_1 , and the corresponding adjoint concentration field R_1 can also be simulated, as
 28 shown in **Fig. 6.1**. According to the adjoint relationship in Eq. (2.24), it is reasonable to believe that

1 R_1 at the source should be close to D_1 . The Bayesian inference evaluates the probability of each
 2 possible source at different locations using Eq. (2.10). based on the difference between the
 3 measurement \mathbf{D} and the adjoint concentration \mathbf{R} . The probability is large when the difference is
 4 small. According to the adjoint concentration distribution of sensor No. 1 in **Fig. 6.1**, values of R_1
 5 in the red potential area is similar to that of the true source because they have the same color bar. In
 6 this case, the posterior probability of the potential area is similar to that of the true source. The
 7 estimation result will show that the source may locate at anywhere in the potential area. Therefore,
 8 the estimation results based on one sensor are not accurate. It is still impossible to infer the true
 9 parameters of an unknown source.

10 Then, we can supplement another sensor to improve the estimation. Naturally, we can obtain a
 11 measurement vector $\mathbf{D} = (D_1, D_2)$, and each place has an adjoint concentration vector $\mathbf{R} =$
 12 (R_1, R_2) . Owing to the information provided by sensor No. 2, the potential area estimated by one
 13 sensor rapidly shrinks because the adjoint concentration R_2 of most parts does not agree with the
 14 R_2 of the source or measurement D_2 . After removing these areas, the results of Bayesian inference
 15 become much better than before that the potential area tightly concentrated around the true source.

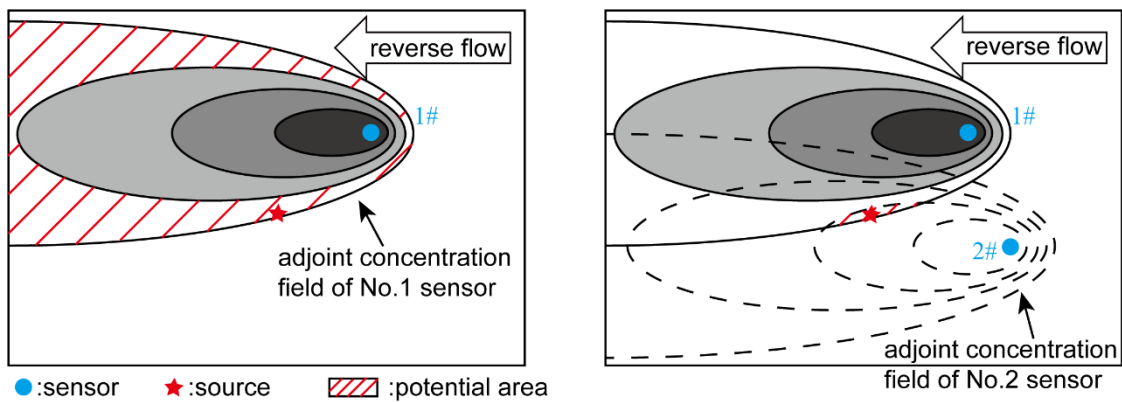


Figure 6.1. Diagram of the proposed objective function $H(\mathbf{R}|\mathbf{I})$. Left: Adjoint concentration field of sensor No. 1; Right: Adjoint concentration fields of sensor No. 1 & 2.

16 A statistical perspective highlights that two sensors collect much more information about the
 17 dispersion of the unknown source than one sensor. Hence, the resultant estimations contain
 18 considerably fewer uncertainties. The increase in the measured information can be represented by
 19 the entropy of the spatial probability distribution of the adjoint concentration in the target domain.
 20 When compared with one sensor, two sensors complicate the adjoint concentration field by adding

1 one dimension in \mathbf{R} , making the spatial probability distribution wider, and increasing the
2 corresponding entropy. For this reason, there are significantly fewer areas owning a similar adjoint
3 concentration vector with the true source, which can then be successfully identified.

4 If we want to re-deploy the second sensor, the location of sensor No. 1 should not be considered.
5 Under such circumstances, the adjoint concentration field of sensor No. 2 is the same as that of
6 sensor No. 1, which means that sensor No. 2 will fail to measure any new information regarding
7 pollutant dispersion. Equally, the spatial probability distribution, entropy of the adjoint
8 concentration field, and potential area remain unchanged. The estimation of the STE would not be
9 improved. As a result, it is always desirable for the new sensor to increase the entropy of adjoint
10 concentration; in other words, fewer points have the same adjoint concentration vector as the true
11 source. Moreover, because the unknown source can appear anywhere in the target area, the ideal
12 state is that each place can own a unique adjoint concentration vector in which the spatial probability
13 distribution is close to a uniform distribution and the entropy reaches the maximum value. In other
14 words, the physical meaning of the maximum entropy of spatial distribution of adjoint concentration
15 is that the sensor configuration has unique relationship with all possible sources. When any source
16 appears, the measured concentration vector and the adjoint concentration vector are special enough
17 to eliminate all other options and make sure the STE can identify the true source. Hence, when the
18 number of sensors is fixed and limited, we can find the sensor configuration with the largest entropy
19 of adjoint concentration, which is expressed by the objective function.

21 6.2.2 Simulated Annealing

22 After determining the objective function to rank each sensor configuration, it is necessary to
23 quickly finish the optimization process and determine the best configuration with the largest
24 $H(\mathbf{R}|I)$. However, the optimization design for a large number of candidates is troublesome because
25 of the heavy calculation burden. If we want to select j sensors out of n candidates to construct the
26 network, the number of possible combinations is $C_n^j = \frac{n!}{(n-j)!j!}$. When j is 8 and n is 100, as a
27 common example, the total number of combinations reaches over one thousand billion. It is
28 impossible to directly calculate the objective function for each combination and select the optimum
29 function. Therefore, this research applied the SA algorithm (Kirkpatrick et al., 1983), which is a
30 heuristic method inspired by the process of cooling a liquid to the lowest possible energy state. The

1 detailed steps of the SA are as follows:

2

3 *Initialization:* Randomly select j sensors to form the first sensor network \mathbf{m}_1 and calculate its joint
4 entropy of adjoint concentration $H(\mathbf{R}_1|I)$

5 *For* $k = 2$ to N :

6 Randomly change a sensor in \mathbf{m}_{k-1} to generate a new network \mathbf{m}_k . Calculate $H(\mathbf{R}_k|I)$.

7 *If* $H(\mathbf{R}_k|I) > H(\mathbf{R}_{k-1}|I)$

8 *Accepted*

9 *Else*

10 *If* $N[0,1] < \exp\left(-\frac{H(\mathbf{R}_{k-1}|I) - H(\mathbf{R}_k|I)}{T_k}\right)$

11 *Accepted*

12 *Else*

13 *Rejected;* $\mathbf{m}_k = \mathbf{m}_{k-1}$, $H(\mathbf{R}_k|I) = H(\mathbf{R}_{k-1}|I)$

14 *End if*

15 *End if*

16 *End For*

17

18 In this algorithm, N is the specified loop number for the stop criterion. $N[0,1]$ is a random
19 variable with a uniform distribution between 0 and 1, which means that the new configuration has
20 a certain possibility to be accepted even if its entropy is smaller. This setting enables the
21 optimization process to jump out from the local maximum. The critical factor T_k is a virtual
22 temperature defined by $T_k = T_0 \times a^k$, where $0 < a < 1$ is the cooling coefficient. In the
23 beginning, when T_k is large, the probability of jumping to an inferior configuration is higher, and
24 the algorithm is flexible for exploring the entire parameter space. As k increases, T_k tends to 0,
25 making it almost impossible to move to an inferior configuration even if the difference is quite small.

1 T_k adjusts the balance between local and global exploration, thereby affecting the efficiency of the
 2 optimization process. Because T_k is closely related to the value of the objective function case by
 3 case, there is no general way to set the value of T_k . In this research, T_0 was imposed as 1 and a
 4 as 0.9, which is similar to previous research (Fuentes et al., 2007). To remove the influence of the
 5 initialization point, eight SA optimization chains were run in parallel with $N = 5000$. The sensor
 6 configuration with the largest entropy among the eight chains was regarded as the optimal
 7 configuration.

8

9 **6.2.3 Calculation of $H(\mathbf{R}|I)$**

10 During SA optimization, $H(\mathbf{R}|I)$ needs to be calculated for each configuration. It is the joint
 11 entropy of probability $p(\mathbf{R}|I)$, where $\mathbf{R} = (R_{i,j})$, $R_{i,j} = q_{s,i} \overline{C_j^*}(\mathbf{x}_{s,i})$, i is the index for all
 12 possible sources, and j is the index for sensors. In this definition, $\overline{C_j^*}$ is a three-dimensional
 13 adjoint concentration function of sensor No. j . Assume that there are m possible sources in total and
 14 we have n sensors to form a configuration, \mathbf{R} can be expressed as:

$$\mathbf{R} = \begin{bmatrix} R_{1,1} & R_{1,1} & \cdots & R_{1,n} \\ R_{2,1} & \ddots & & \\ \vdots & & R_{i,j} & \vdots \\ R_{m,1} & \cdots & \ddots & R_{m,n} \end{bmatrix} \quad (6.12)$$

15 Here $m \gg n$. Each column represents the adjoint concentration of a sensor No. j
 16 corresponding to all possible sources. Each row represents the adjoint concentrations of all sensors
 17 for a single source No. i . Because the sensor configuration is designed for all possible sources \mathbf{s} in
 18 the target domain, all rows are included in $p(\mathbf{R}|I)$. In this case, we regard each row as a vector, and
 19 calculate the probability distribution of these vectors, which is $p(\mathbf{R}|I)$.

20 Each $R_{i,j}$ is a multiple of source strength and adjoint concentration of a spatial point. However,
 21 the source strength is unavailable when the sensor configuration is designed. It could have a wide
 22 range, which makes m in Eq. (6.12) huge and calculation cost of $p(\mathbf{R}|I)$ heavy. To decrease this
 23 cost, we notice that in the Bayesian inference for STE, q_s is sampled by the MCMC with a uniform
 24 prior distribution. Meanwhile, in the process of location optimization for sensors, changing the
 25 location of e would only affect $\overline{C_j^*}$, so the entropy $H(\mathbf{R}|I)$ mainly depends on the $\overline{C_j^*}$. Hence, in
 26 the case study of Section 6.3, we calculate $H(\overline{\mathbf{C}^*}|I)$ instead of $H(\mathbf{R}|I)$ and assume that larger

1 $H(\overline{\mathbf{C}^*|I})$ is equivalent to larger $H(\mathbf{R}|I)$. In the following content, $H(\mathbf{R}|I)$ means $H(\overline{\mathbf{C}^*|I})$. In this
2 case, m in Eq. (6.12) equals to the number of spatial grids of the adjoint concentration field \overline{C}_j^* .

3 Future research is still necessary to clarify the effects of this calculation method on the final
4 result. It is still recommended that a wide range of discrete q_s should be multiplied with the adjoint
5 concentration if enough calculation resources is available.

6

7 **6.2.4 Limitation of the method**

8 Because this method was proposed based on several assumptions, there are some limitations
9 that need to be noted. First of all, this method corresponds with the statistical STE problem of a
10 point source in the target area with stationary meteorological situations. Its effectiveness cannot be
11 ensured when applied to other environmental monitoring. More research is still necessary to
12 evaluate its performance in real urban applications with complex meteorology. However, because
13 numerical simulations of adjoint concentration include the information of meteorology, it is
14 reasonable to say that the proposed method has the potential ability in such scenarios.

15 Besides, during the mathematical proof, one requirement is that the adjoint concentration
16 should be accurately simulated to ensure that Eq. (2.24) is satisfied. As a result, the performance of
17 this proposed method may depend on the quality of adjoint concentration simulation. What's more,
18 it is useful to note that there is no guarantee that SA would converge at the global optimal, however,
19 it is likely that a 'near optimal' can be obtained (Altinel et al., 2008). Considering the complicated
20 errors in real applications, this imperfection is considered to be acceptable here.

21

22 **6.3 Case study**

23 To verify the performance of the proposed method, it is used to design an optimum sensor
24 configuration for an ideal urban area: a regular block-arrayed building group model, just the one in
25 **Chapter 4**.

26 To comprehensively evaluate the quality of the sensor configurations, it is necessary to identify
27 different sources based on the configuration and check the corresponding STE performance.
28 However, it is troublesome to measure the concentrations of many sources in the wind tunnel
29 experiment, let alone that measurements of several sensor configurations are needed, and it is almost

1 impractical to employ the experiment or field test in this evaluation. Moreover, LES is increasingly
2 gaining credit in the prediction of urban flow fields and dispersion of pollutants (Kikumoto and
3 Ooka, 2012). Some research also relies on the simulation data for sensor network design (Abida et
4 al., 2008) or STE method verification (Mons et al., 2017). As a result, the well-validated LES
5 database produced in **Chapter 4** is applied for the evaluation of sensor configurations.

6

7 **6.3.1 Numerical simulation**

8 The statistical data of the LES are validated in **Chapter 4** in **Fig. 4.4 & 4.5**. It can be confirmed
9 that the simulated flow field and concentration agree well with the wind tunnel experiment
10 measurements. Therefore, using the LES results as concentration measurements would not damage
11 the credibility of the research.

12

13 **6.3.2 Unknown sources**

14 A red area is set as the target monitoring domain where an unknown source may appear, as
15 shown in **Fig. 6.2**. The target area is smaller than the entire calculation domain because the
16 measurement information around the outer bounds could be valuable, especially downstream
17 measurements. If the target domain is the same as the calculation domain, the source is difficult to
18 estimate with most sensor configurations when it is on the streamwise boundary.

19 Then, some sources are set in the target domain to evaluate the quality of each sensor
20 configuration through the STE. It is difficult to obtain the concentration fields of many sources,
21 regardless of experiments or simulations. Because the block configuration is regular in the current
22 case, the concentration field of one point source located at $(4H, 7H, 0)$ simulated in **Chapter 4** was
23 used and copied for the other 24 sources, as shown in **Fig. 6.2**. All these sources are located in the
24 middle of the wake region, where the unsteady turbulent flow is strong, posing a challenging STE
25 problem for sensor configurations. Note that the concentration simulation was well-validated, and
26 the creditability of the concentration database was maintained.

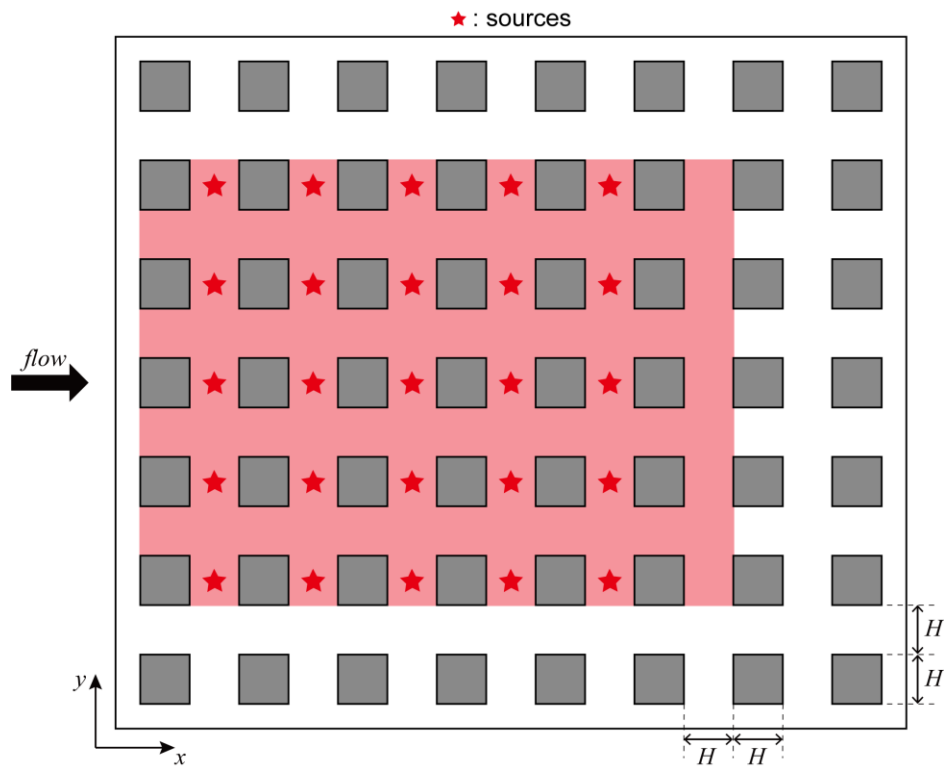


Figure 6.2. Target monitoring domain (red area) and distribution of unknown sources

1

2 6.3.3 Sensor candidates

3 In this part, a number of sensor candidates are prepared, from which the optimization process
 4 can select part of sensors to construct the sensor configuration. In real applications, considering the
 5 deployment limitation of the stationary sensors and complex terrain in the urban areas, most places
 6 cannot be used for reasons like space is not enough or the land property problems. It is more practical
 7 to conduct optimization among the candidates than from all positions. Meanwhile, in the proposed
 8 method, the adjoint concentration fields of all sensor candidates have to be simulated in advance. If
 9 there are too many candidates, the computational cost is very large, and a huge database has to be
 10 saved for the entropy calculation later. Therefore, it is difficult to regard all the places as the location
 11 candidates. In the previous research, the design is usually selected from the sensor candidates
 12 prepared in advance (Kouichi et al., 2019; Ngae et al., 2019). The distribution of the sensor
 13 candidates in this case study is illustrated in **Fig. 6.3**. A total of 113 sensors were located on the
 14 horizontal plane with $z = H/2$. The objective of the optimization is to select the eight best sensor

1 combinations among these candidates.

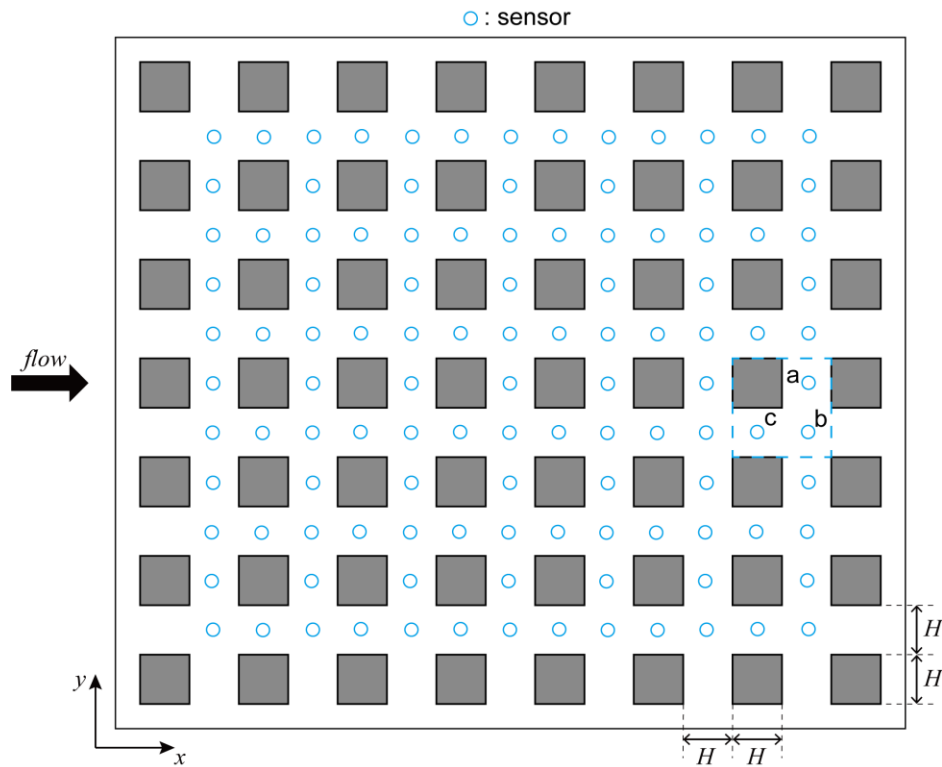


Figure 6.3. Distribution of sensor candidates

2 Because the objective function is based on the entropy of the adjoint concentration, the heaviest
 3 calculation burden is the adjoint equation simulation of all sensors. Once again, the advantage of a
 4 regular flow field is used and the adjoint concentration fields for only three sensors (a, b, and c in
 5 **Fig. 6.3**) are simulated. The fields of the other sensors are obtained by copying these three fields.

6 In **Chapter 5**, the performance of Reynolds averaged Navier-Stokes (RANS) and large-eddy
 7 simulation (LES) of adjoint equations has been compared and discussed. Although LES has been
 8 confirmed to be more accurate, to balance accuracy and calculation costs, the approach proposed by
 9 Xue et al. (2018b) is still adopted here, in which the adjoint equation was simulated by a RANS-
 10 like model based on the mean velocity field predicted by LES in the forward simulation. It should
 11 be noted that this method will bring some modeling errors into the adjoint concentration field and
 12 SCO. The obtained optimal configuration may not be the theoretically best one. Meanwhile, the
 13 simulated adjoint concentrations are also used in the following STE for all configurations, so it is
 14 believed that the modeling errors of adjoint equations simulation will not affect the ranking of
 15 different configurations.

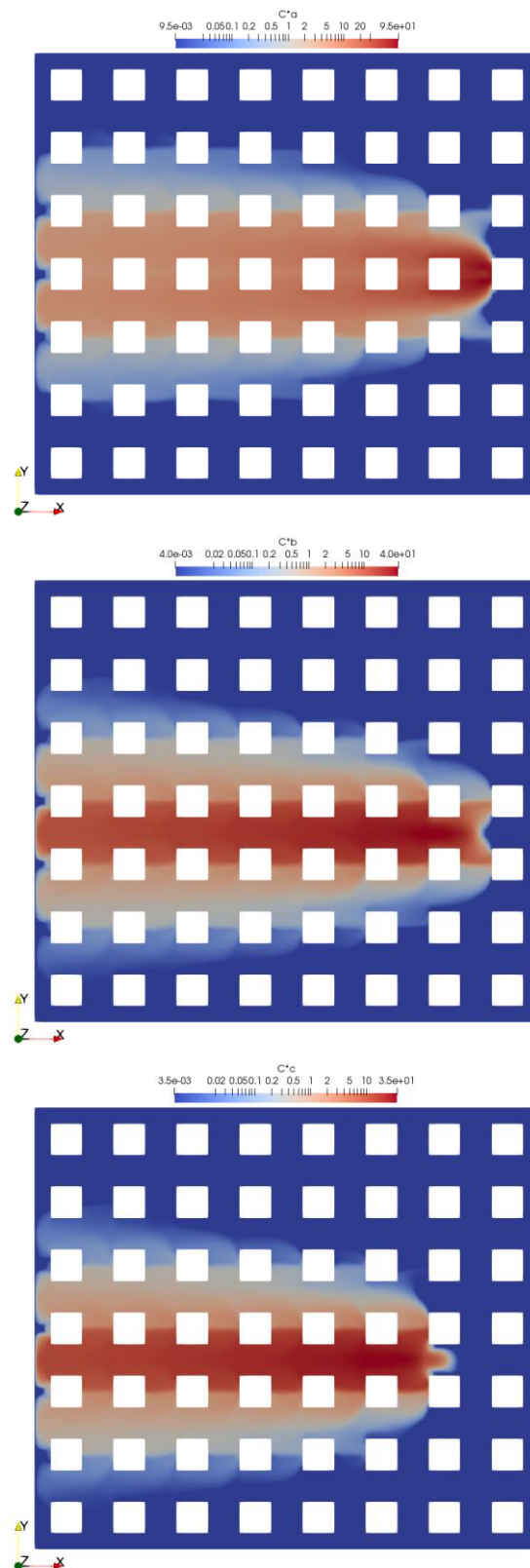


Figure 6.4. Simulated adjoint concentration distribution of sensors a, b, and c in Fig. 6.3. (Horizontal plane with $z = 0$; From top to bottom: sensor a, b, c)

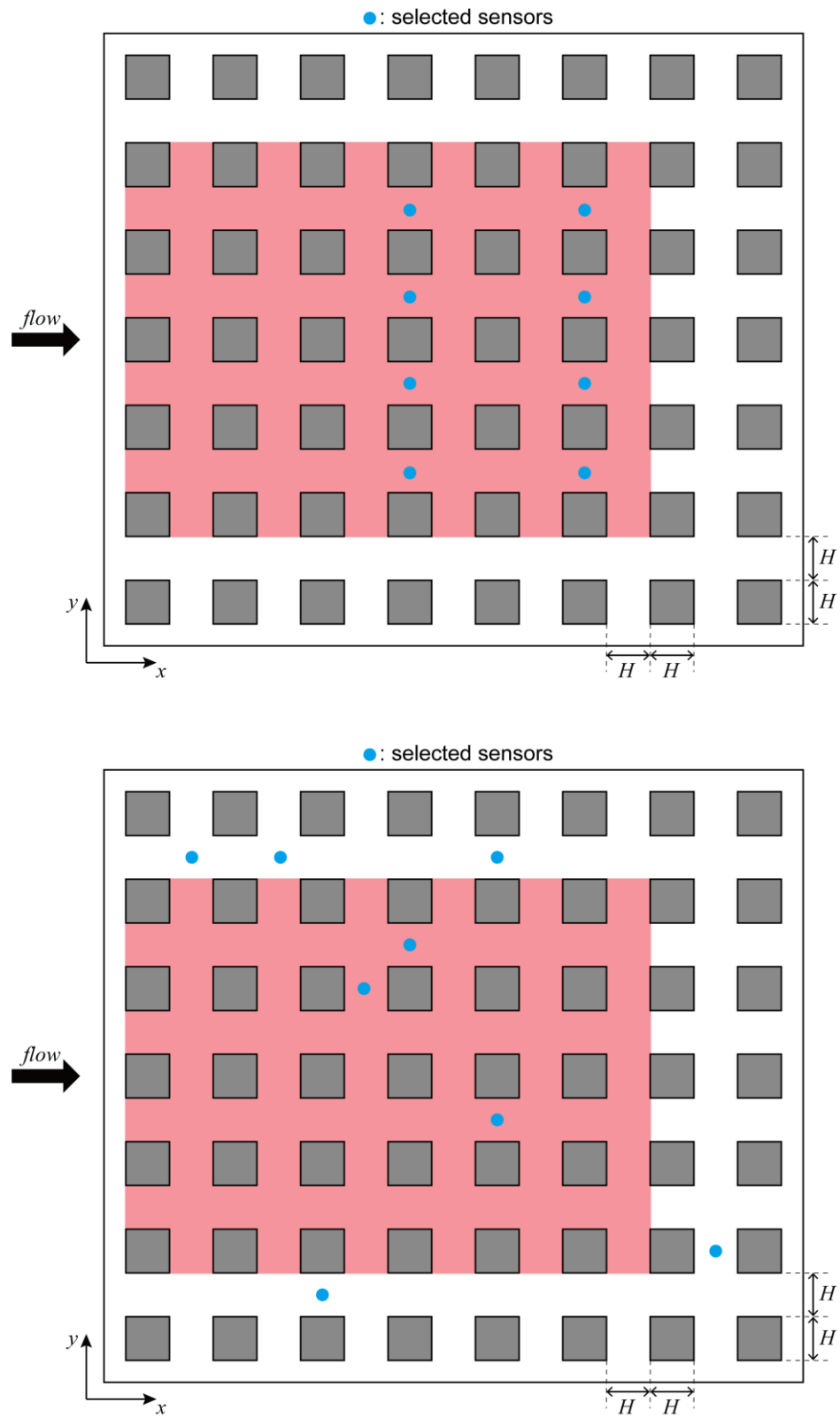


Figure 6.5. Comparison configurations: uniform (top) and random (below)

1 The adjoint concentration fields of three sensors are shown in **Fig. 6.4**.

2

3 **6.3.4 Comparison configurations**

4 To demonstrate the advantage of the optimum sensor configuration, two comparison
5 configurations that were commonly used in previous STE research are selected, namely, the uniform
6 configuration and random configuration (**Fig. 6.5**). The STE results of these two configurations are
7 compared with those of the optimal configuration.

8

9 **6.4 Results and Discussion**

10 **6.4.1 Optimum sensor configurations**

11 The proposed method is applied to determine the optimum sensor configuration. The value of
12 the objective function of each of the eight sensor combinations $H(\mathbf{R}|I)$ was calculated, and the SA
13 quickly adjusted the configuration to reach a larger value. **Fig. 6.6** illustrates the changing process
14 of the objective function value during the eight SA searching chains. Seven chains are
15 transparentized to emphasize the one whose final objective function value is the maximum. The
16 logarithmic horizontal axis was used to clearly illustrate the search process. The SA optimization
17 algorithm was run in MATLAB on a personal computer with Intel® Core™ i7-6700 CPU @ 3.4GHz
18 and 32GB of RAM. The averaged computational time for one SA chain with 5000 optimization
19 steps is about 180s.

20 The changing patterns of the eight chains were similar. Initially, because the SA chain started
21 from a random configuration, the entropy was small. Before approximately 100 steps, SA moved to
22 a smaller entropy rather than a larger one several times. At this moment, the virtual temperature T_k
23 still has a large value, making the SA flexible to jump into the inferior configuration for global
24 exploration. After, the inferior jumping became strict, and the chain only moved to a larger entropy.
25 Meanwhile, the rate of increase also decreased because it became increasingly difficult to find a
26 better configuration. Finally, the value of the objective function remained unchanged for
27 approximately 1500 steps, yielding the largest value of 11.21 and the corresponding configuration.

28 As a comparison, the $H(\mathbf{R}|I)$ values of the random and uniform configurations are 6.35 and
29 9.06, respectively. It can be confirmed that 6.35 is almost the smallest value in the chain, which

1 means that the random configuration is ineffective for measurements. This may challenge the
 2 robustness of STE methods and, consequently, should not be used in real applications. The uniform
 3 configuration is located at the intermediate level of the chain, with a value of 9.06. It should be
 4 noted that the regular block distribution provides some merits to the uniform configuration. Whether
 5 the performance can be maintained in a real urban area with a complicated terrain remains unclear.

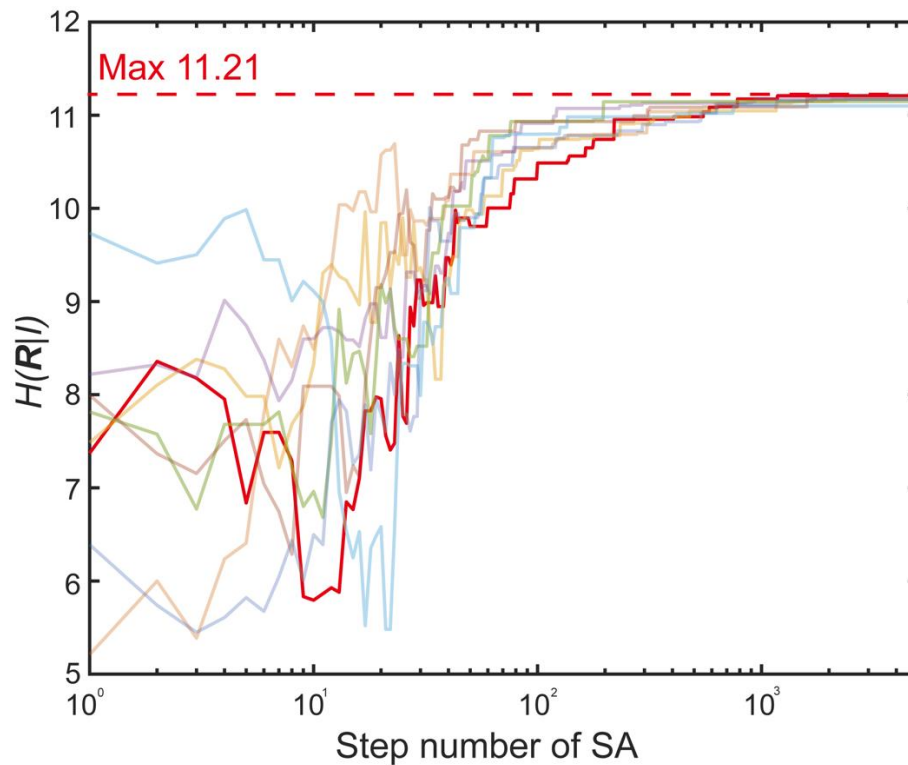


Figure 6.6. Variation in $H(R|I)$ among eight SA chains. The red line indicates the chain whose final objective function value is the maximum. The other chains are shown with faded colors.

6 The obtained optimum sensor configuration is shown in **Fig. 6.7**. This configuration has several
 7 interesting characteristics. First, the configuration is not symmetric as the calculation domain. This
 8 could result from the simulation errors of the adjoint equation because the simulated adjoint
 9 concentration fields in **Fig. 6.4** were also slightly skewed. It is probable that the optimization result
 10 is sensitive to the adjoint concentration field, so the accurate simulation of the adjoint equation is
 11 important to the proposed method. From a good perspective, the proposed optimization method
 12 captures the information contained in the adjoint concentration field well, which further indicates
 13 that the source-receptor relationship is well considered.

1 Second, most sensors were located in the downstream area of the target domain. This
 2 distribution coincides with the intuitive decision that people usually tend to arrange all sensors
 3 downstream of sources based on empiricism. If the wind direction is constant, the adjoint
 4 concentration fields of these downstream sensors cover a sizable part of the target area, indicating
 5 a strong measurement capability. Although the logic seems obvious, it is difficult to quantize it for
 6 decision-making in the optimization algorithm. The produced optimal configuration implies that the
 7 entropy-based objective function in the proposed method is an appropriate solution.

8 Moreover, there is one characteristic beyond empiricism that not all sensors gather downstream.
 9 The algorithm detected that the downstream area was already crowded with six sensors. It is
 10 ineffective to add more sensors there in terms of cost performance because the adjoint concentration
 11 field of the new sensor would nearly overlap with existing ones, and the new source–receptor
 12 information is limited. Consequently, two sensors were arranged in the middle of the target domain
 13 to distinguish the front and back parts of the domain. The importance of these two sensors is easy
 14 to ignore when the sensor configuration is designed based on the human experience. The proposed
 15 method can deepen the insights from the entropy of the adjoint concentration.

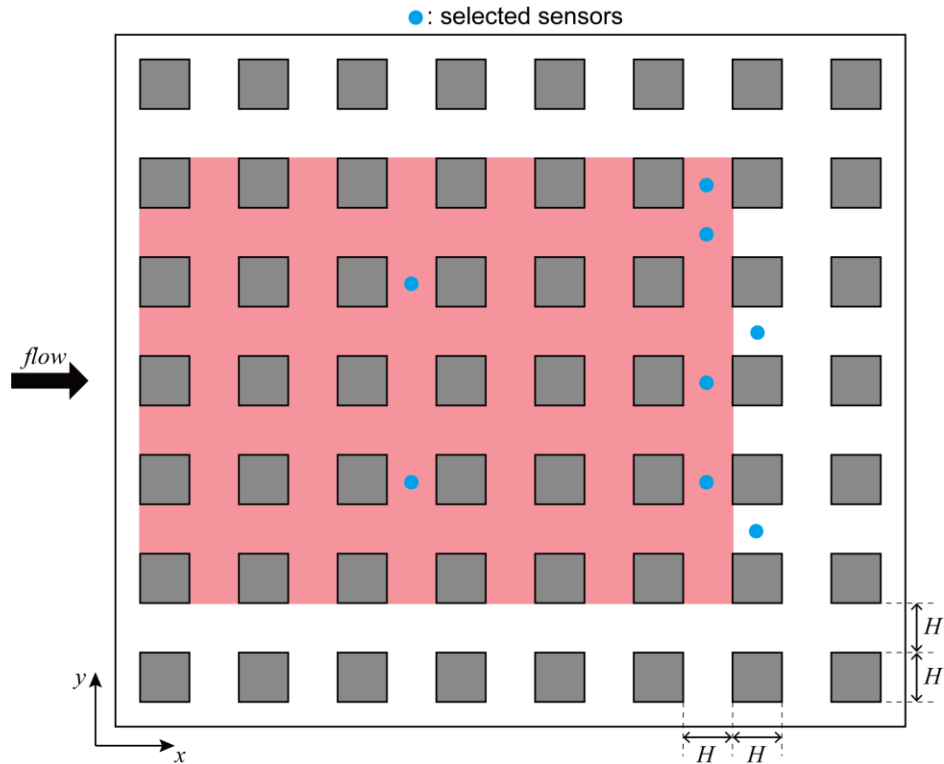


Figure 6.7. Optimum sensor configuration produced by the proposed method

1 According to the discussion above, the chosen optimum sensor configuration is reasonably
2 designed and seems capable of satisfying the STE performance. We employed this configuration as
3 well as random and uniform configurations to estimate 25 unknown sources. In the Bayesian
4 inference, the variance of errors is set as $\sigma_{d,i}^2 + \sigma_{m,i}^2 = r * D_i$, where $r = 0.3$. The results are as
5 follows.

7 6.4.2 Estimation results of one source

8 First, we present the STE results of one source located at $(2H, 7H, 0)$ based on the
9 measurements of the three sensor configurations. The Bayesian inference of STE was executed with
10 MATLAB on a personal computer with Intel® Core™ i7-6700 CPU @ 3.4GHz and 32GB of RAM.
11 The averaged computational time for estimating one source is about 11s. **Fig. 6.8** summarizes the
12 posterior joint probability distribution $p(x, y|\mathbf{D}, I)$, which represents the location estimation of the
13 true source. **Fig. 6.9** shows the posterior marginal probability distribution for the strength estimation.

14 The optimal configuration produced the best estimations among the three. In **Fig. 6.8(a)**, most
15 of the probability mass is concentrated around the true source. Only a fraction deviates from the
16 truth, and the distance is constrained to approximately $2H$. The accuracy is better than the STE in
17 **Chapter 5** with the same adjoint equation simulation method. Possible reasons could be that the
18 optimal sensor configuration is more informative. Besides, it can be noticed that the downstream
19 distance between source and sensors is larger than that in **Chapter 5**. According to the previous
20 research (Xue et al., 2017), the simulated adjoint concentration distribution tends to be a narrower
21 plume compared with that in the real scenario. Therefore, a smaller distance may result in a larger
22 difference between \mathbf{D} and \mathbf{R} .

23 In addition, the emission strength was also well identified, as shown in **Fig. 6.9(a)**. All
24 probability masses lie between 0 and 2. Although the peak value does not perfectly agree with the
25 truth, the error is acceptable. This estimation error mainly results from the simulation method of the
26 adjoint equation. The LES model could effectively mitigate this problem and further improve the
27 STE accuracy. Because this is not relevant in the current study, the details are not expanded upon
28 here.

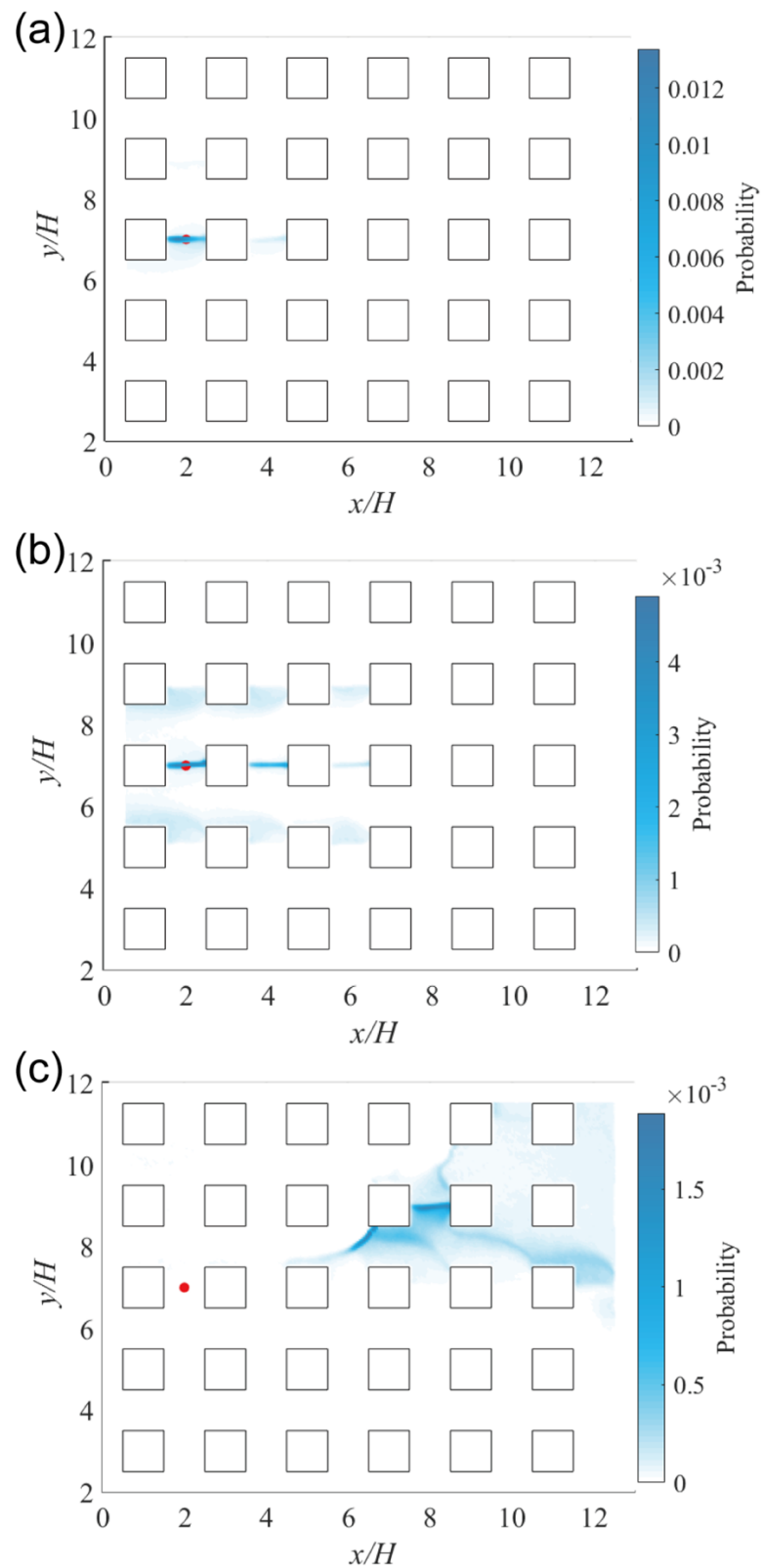


Figure 6.8. Posterior joint probability distribution $p(x, y | \mathbf{D}, \mathbf{I})$ obtained by three sensor configurations: (a) optimum; (b) uniform; and (c) random. Red point: true source.

1

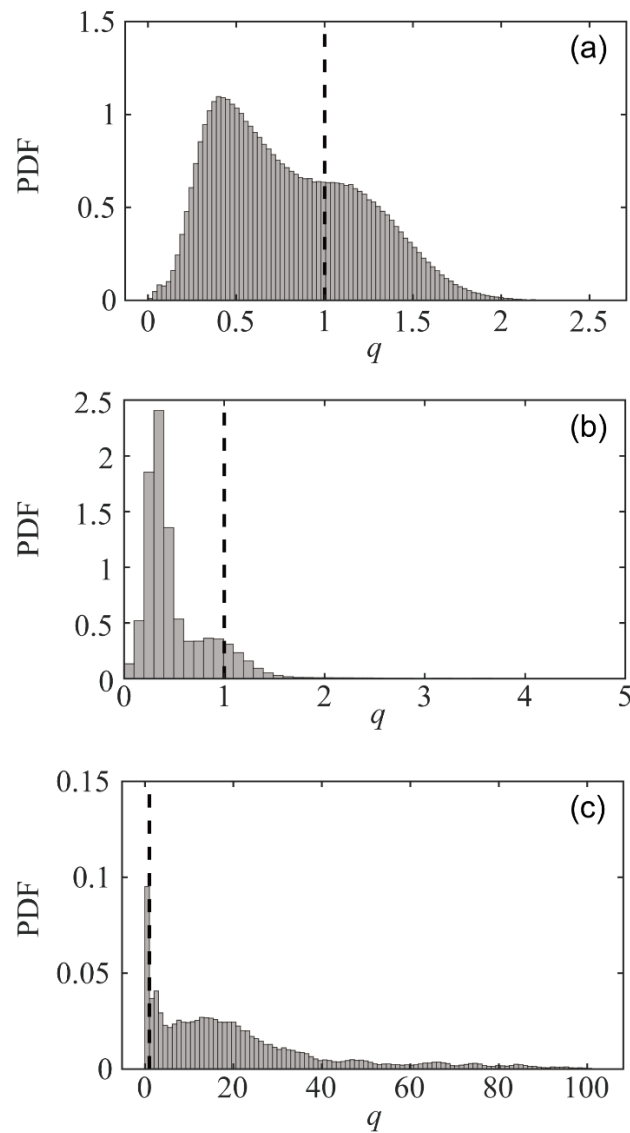


Figure 6.9. Posterior probability distribution $p(q|\mathbf{D}, \mathbf{I})$ obtained by three sensor configurations: (a) optimum; (b) uniform; and (c) random. Dotted line: true value

2 The uniform configuration failed to provide estimations with comparable accuracy to that of
 3 the optimal configuration. Although the true location has a distribution peak, many probability
 4 masses diverge broadly. This divergence represents redundant uncertainties caused by low-quality
 5 measurements, which is insufficient for STE sampling to firmly adhere to the truth. The same
 6 situation can also be confirmed in the strength estimation shown in **Fig. 6.9(b)**. Although the width
 7 and the general shape are similar to those of the optimal configuration, more probability masses
 8 move from the truth to the wrong peak. In short, the STE estimation contains more noisy

1 uncertainties than that of the optimal configuration because the measurements of this source
2 provided by the uniform configuration are less informative.

3 The performance of the random configuration was the worst. All probability mass in **Fig. 6.8(c)**
4 departs from the true source, which reveals that the STE broke down. The information measured by
5 the random configuration cannot support the STE in identifying the true source. Similarly, severe
6 unsteadiness was observed in the strength estimation. Although the peak is close to the truth, the
7 width of the distribution spreads over 100 times, as shown in **Fig. 6.9(c)**. In this case, the estimations
8 have little reliability for risk management.

9 In general, the performance of these three configurations has the same order as their entropy
10 values, which proves that the objective function design in the proposed method is reasonable.

11

12 **6.4.3 Estimation results of all sources**

13 To avoid the contingency in STE for only one source, we conducted STE for all 25 sources
14 based on the three configurations. Considering the length of the thesis, it is difficult to present all
15 posterior PDFs here. Two indices $E_{d,m}$ and $E_{q,m}$ introduced by Eq. (5.6 & 5.7) are used to
16 quantify the estimation accuracy of STE for all unknown sources based on the three sensor
17 configurations. m is the index for source *No. m*

18 Accurate STE always has smaller $E_{d,m}$ and $E_{q,m}$. Because the posterior PDFs are often
19 highly skewed, $E_{d,m}$ and $E_{q,m}$ are first-order statistical moments and cannot comprehensively
20 reflect real situations. The two PDFs may be different, even if their expectations are the same. In
21 the Bayesian inference STE, it is expected that the probability mass is concentrated around the true
22 value with limited distribution rather than a uniform distribution with too many uncertainties. As a
23 result, we introduced the third index to evaluate the concentration degree of the posterior PDF: the
24 joint entropy $H(x, y, q)$ of $p(x, y, q|\mathbf{D}, \mathbf{I})$. A PDF with a small $H(x, y, q)$ and small $E_{d,m}$ &
25 $E_{q,m}$ indicates a good estimation that most probability is close to the truth.

26 The values of the three indices for STE based on the three configurations are summarized in
27 **Table 6.1**. All values were averaged for 25 unknown sources. The optimal configuration has the
28 smallest values for all three indices. The location error of $2.83H$ and strength error of 1.5 are at the
29 normal level for STE with a steady simulation of the adjoint equation, which has been confirmed in
30 **Chapter 5** with 16 sensors. When the sensor configuration becomes uniform, the estimation

1 accuracy becomes inaccurate according to the three indices. The location error increased by
 2 approximately 65%, and the strength error significantly increased by approximately 25 times when
 3 compared with that of the optimal configuration. The limited amplification of $H(x, y, q)$
 4 demonstrated that considerable deviation was caused in $E_{q,m}$ by ineffective measurements of
 5 uniform configuration while the concentration rate of PDF was less damaged. However, the random
 6 configuration thoroughly eroded the estimation accuracy in terms of deviation and concentration
 7 rate in the PDF. The location error, strength error, and $H(x, y, q)$ increased by more than 80%,
 8 7300%, and 47%, respectively, compared with the optimal configuration. It should be mentioned
 9 that most sensors in the random configuration were located in the upper-left area. When the source
 10 appears at the lower right corner, only one or two sensors can measure the concentration changes,
 11 which is obviously too insufficient and results in STE failures.

12 Consequently, the efficiency of the optimum sensor configuration in STE was confirmed. The
 13 performance of the configurations is positively proportional to the joint entropy of adjoint
 14 concentrations $H(\mathbf{R}|I)$. Among these three indices, the strength error is the most sensitive to the
 15 measurement quality of the sensor configuration. During the Bayesian inference, the exploration
 16 space of the location is limited in the current case (approximately $12H$ for x and $10H$ for y), while
 17 the sampling range is infinite for strength. When the measurements are not sufficient for source
 18 identification, Bayesian inference would prefer to adjust the strength rather than the location
 19 parameter to meet the gap between the measurements and simulated concentrations. Once a
 20 probable parameter combination was found, most sampling would occur in that combination. As a
 21 consequence, E_{di} and $H(x, y, q)$ were not promoted comparably to E_{qi} in the random
 22 configuration.

Table 6.1. Summarized estimation indices of three sensor configurations

	Mean $E_{d,m}$	Mean $E_{q,m}$	Mean $H(x, y, q)$
Optimal configuration	$2.83H$	1.54	12.23
Uniform configuration	$4.67H$	39.54	14.17
Random configuration	$5.17H$	112.53	18.01

1

2 **6.5 Conclusions**

3 This chapter proposes an SCO method for identifying an optimum sensor configuration that
4 can provide informative concentration measurements for STE, regardless of where the unknown
5 source appears in the target domain. When compared with other configurations, it is expected that
6 Bayesian inference STE can produce better estimations based on measurements of the optimal
7 configuration, in which the probability mass in the posterior PDF gathers around the true value with
8 a narrow distribution.

9 The proposed method is composed of an objective function and an SA algorithm. The objective
10 function was set as the joint entropy of the spatial probability distribution of the adjoint
11 concentration of a configuration to evaluate its measurement ability. The reasonability of this setting
12 was explained from mathematical and physical perspectives. The value of the objective function
13 represents the uncertainty about the source that can be measured by the sensor configuration. If
14 more uncertainty has been measured, less probability divergence would appear in the posterior PDF,
15 and the probability mass could concentrate around the truth. Therefore, the optimum sensor
16 configuration should have the largest objective function value. SA was applied to quickly find the
17 optimal configuration using the objective function value among countless possible combinations.

18 To evaluate the performance of the proposed method, 25 unknown point sources in the regular
19 block-arrayed building group model were estimated based on measurements of optimum, random,
20 and uniform configurations. According to the results, it was found that the main calculation cost in
21 the proposed method results from the adjoint equation simulations for all sensor candidates. By
22 comparing the STE performance of the three configurations, it was proved that the estimation
23 accuracy is positively correlated with the objective function value such that the optimal
24 configuration is the best, the uniform configuration is less accurate, and the random one is the least
25 accurate. Hence, the objective function is appropriately selected. In the presented posterior PDFs
26 for the STE of one source, the probability mass settled around the truth with a limited spread in the
27 optimal configuration, while the estimations totally deviated in the random configuration. This
28 finding was confirmed again when the STE results of 25 unknown sources were quantized using
29 three indices. The average estimation errors of the optimal configuration were limited. The location
30 and strength estimation errors increased by approximately 80% and 7300%, respectively, when the
31 configuration changed from optimum to random.

1 The measurements in the case study were synthesized from a well-validated LES simulation.
2 Future research with complicated urban terrain and real measurements is necessary to verify the
3 robustness of the proposed method. Besides, the meteorological condition was set stationary by
4 constant inflow boundary in adjoint equation simulation. The optimal sensor configuration may
5 change with different wind directions/speeds. However, it has been presented that the adjoint
6 concentration simulation is able to include the effects of terrain geometry, meteorological situations,
7 and building structures. It is reasonable to expect that the proposed method can be extended to real
8 urban applications with variable conditions.

9

10

1 Symbols

C	: the concentration field
C_s	: the concentration field of source s
C_e^*	: the adjoint concentration field of sensor e
C_n^j	: the number of possible j combinations out of a set with n elements
D	: the measurements vector
D_i	: the measurement of the sensor with index i
e	: the selected sensor for further optimization when other sensors are fixed
$E_{d,m}$: the expectation of the Euclidian distance between the estimated location and the true location of source $No. m$
$E_{q,m}$: the expectation of the normalized difference between the estimated strength and true strength of source $No. m$
H	: the edgy length of blocks in the simulation (=60 mm)
$H(x)$: the information entropy of a probability distribution $p(x)$
I	: the background information for Bayesian inference
$I_m(\mathbf{s}; \mathbf{D} I)$: the mutual information representing the information about the source \mathbf{s} provided by sensors with measurements \mathbf{D}
m_i	: the i th searching in the Simulated Annealing
N	: the specified number of searching steps in Simulated Annealing

-
- $N[a, b]$: the uniform distribution bound between a and b
- $p(x)$: the probability of event x occurring
- $p(A|B)$: conditional probability of event A occurring given that B is true
- Q : the gas flow rate at the source
- q : strength samplings produced in MCMC for the point source
- q_s : the true strength of the point source
- r : the ratio between the specified error covariance in the Bayesian inference and the measurements
- \mathbf{R} : adjoint concentration fields
- \mathbf{s} : a vector representing the unknown source
- T_0 : the initial virtual temperature for Simulated Annealing
- T_k : the virtual temperature of the system at the k th step in Simulated Annealing
- \mathbb{V} : a parameter space where the unknown source is possibly located
- (x, y) : location samplings produced in MHMC for the point source
- (x_s, y_s) : the true location of the point source
- a : the cooling coefficient in the Simulated Annealing
- $\sigma_{d,i}^2$: the variance of error in the measurement of the sensor with index i
- $\sigma_{m,i}^2$: the variance of error in the modeling concentration for the sensor with index i

— : Reynolds average operator

1

2

1

2

3

4

5

6

7

8

9

Chapter 7

10

Conclusion and Future research

11

12

13

14

15

16

17

18

19

1 **7.1 Conclusions**

2 This thesis develops a statistical estimation method for an unknown source of atmospheric
3 pollutants in the complicated urban environment based on Bayesian inference. The main
4 conclusions of each chapter are as follows.

5 Chapter 1 introduces the research background, objective, and structure of the thesis.

6 Dispersion emergencies caused by unknown sources occurred in the urban area from time
7 to time and may cause considerable damage to people and the environment. It is necessary to
8 estimate the source parameters as soon as possible after the emergencies happen. Dealing with
9 the characteristics of source term estimation (STE) applications in the urban environment, this
10 thesis proposed three improvements for the statistical STE to realize better estimation
11 performance, and evaluate their effectiveness by a series of case studies.

12 Chapter 2 provided a short review of the recent progress in STE of atmospheric pollutants,
13 and the basic methodology used in the thesis is introduced.

14 STE consists of three basic elements: measurements, estimation algorithm, and source-
15 receptor relationship. Different methods for these elements have been proposed in the previous
16 research. Considering the applications in the urban environment, in this thesis, the time-
17 averaged concentrations which can be measured by most sensors are the measurements; the
18 Bayesian inference is selected as the estimation algorithm to assess the noise in STE; the adjoint
19 equation method is applied for the source-receptor relationship simulation to deal with the
20 complicated dispersion phenomenon in the built area.

21 Chapter 3 proposed a line source estimation method by combining the basic methodology
22 with the super-Gaussian function.

23 The coefficients of the super-Gaussian function were added into Bayesian inference to
24 realize the geometry estimation. The applicability was first confirmed by a numerical
25 experiment of an ideal urban boundary layer. The proposed method successfully inferred that
26 the source is line-like without any prior knowledge. Based on this case, the effects of different
27 sensor configurations on the line source estimation were discussed. Because the line source
28 contained more geometric information than point sources, the requirements on the sensor
29 configuration become higher that both sensors near the source and null-measurement sensors
30 are indispensable. The conventional sensor configuration may fail in the line source estimation.

1 To examine the robustness of the proposed method against measurement and modeling
2 errors, the second case of a simplified urban square with wind tunnel experiment measurements
3 was conducted. The line source was successfully identified by the proposed method again. By
4 comparing to the conventional method with ideal point assumption, it is confirmed that the
5 proposed method can not only provide the geometry estimation but also reduce the inference
6 errors caused by the point source assumption. It is meaningful to include the geometry
7 estimation into STE.

8 Chapter 4 constructed a compression database by wavelet-based compression method for
9 the turbulent flow field of a block-arrayed building group model based on large eddy simulation
10 (LES) raw data. It is a prerequisite content for Chapters 5 & 6.

11 The compression ability and error control of the wavelet-based compression method was
12 analyzed. The influence of compression on the quality of the data was checked from a single
13 snapshot and time-series perspectives. In the case study, it was found that about 100 times
14 compression can satisfy the requirement of flow field visualization, large-scale turbulent
15 structure preservation, and afterward dispersion simulation. Therefore, it is reasonable to say
16 that the wavelet-based compression method is a powerful tool to construct a portable flow
17 database. The unsteady simulation of the adjoint equation can be realized based on a
18 compressed flow field.

19 In Chapter 5, to improve the accuracy of STE in complex urban applications, a new method
20 was developed based on Bayesian inference coupled with unsteady adjoint equation modeling
21 via LES.

22 The performance of the proposed method was evaluated in the block-arrayed urban model
23 with the wind tunnel experiment measurements of continuous dispersion of a point source. The
24 LES of the adjoint equation for the source-receptor relationship was realized based on the
25 compressed flow field constructed in Chapter 4. As a comparison, another STE was also
26 conducted with a conventional method, where steady adjoint equations were simulated with
27 Reynolds averaged Navier-Stokes (RANS) model. The results showed that the modeling of the
28 adjoint equation and STE were significantly improved by the LES. The complicated turbulent
29 flows caused by buildings destroyed the reliability of the gradient diffusion hypothesis in the
30 conventional RANS simulation of the adjoint equation. Although the proposed method needs
31 more computational resources, to effectively perform STE in the complicated urban

1 environment, it is valuable to apply LES to adjoint equation simulation.

2 Chapter 6 proposed a sensor configuration optimization method for STE by the design of
3 an objective function and application of the simulated annealing algorithm.

4 The objective function was set as the information entropy of the spatial distribution of the
5 adjoint concentration field. Its ability to represent the measurement ability of sensor
6 configurations was proved from the views of mathematics and physics. Simulated annealing
7 was applied to find the optimal configuration which owns the largest value of the objective
8 function.

9 The proposed method was utilized to design an optimal sensor configuration for the block-
10 arrayed urban model in Chapter 4. The performance of the optimal configuration in STE was
11 compared to uniform and random configurations through estimations for 25 unknown sources.
12 The results revealed that the accuracy of STE is related to the entropy contained in the adjoint
13 concentration of the configuration such that the design of the objective function is reliable. The
14 optimal configuration outperforms the other two in STEs. It is valuable to use the proposed
15 method to guide the configuration design in real applications.

16

17 **7.2 Future research**

18 Although some progress has been achieved in this thesis, there are still several limitations
19 in the statistical STE method which need to be noticed. To ultimately realize effective STE for
20 the urban environment, here are future research summarized based on the author's knowledge.

21 Firstly, more advanced measurements should be utilized in STE. In all STEs in the thesis
22 as well as a considerable amount of previous research, only the time-averaged concentration of
23 pollutants was used as measurements. Although this information is the simplest and most basic
24 data to measure for sensors, if other advanced measurements can be accounted in, the accuracy
25 of STE could be improved further following the measured information entropy in Chapter 6.
26 For instance, adding the turbulent flux into measurements is likely to better distinguish the
27 wrong estimations and reduce errors in STE further. Research using time-series measurements
28 in STE is still sparse until now. It is believable that much more information can be dug out from
29 time-series data than average value, thus more efforts should be devoted to how to use it in STE.

1 More importantly, the main limitation of time-averaged concentration is that it should be
2 measured after the dispersion has reached into the statistically steady state by a large enough
3 averaging time scale. It means that plenty of time is cost by measurement before the STE, which
4 potentially brings more risk due to the continuous dispersion of unknown sources during this
5 time. It is meaningful to develop a technique which can identify the unknown source faster
6 based on the sudden change of concentration or time lags between sensors.

7 Secondly, the estimation target is assumed to be a single source. It means that the proposed
8 methods may suffer from identifying multiple sources that the accuracy cannot be promised
9 any more. Since dispersion emergency with multiple sources is a possible scenario in the real
10 applications, it is necessary to extend the proposed methods to handle with them in the future.

11 Thirdly, the simulation method for adjoint equations can be improved further. Although
12 the proposed LES of adjoint equations realized accurate prediction of source-receptor
13 relationship, its calculation cost is still large for real applications. Methods with better cost-
14 performance ratios should be developed in the future. According to the results in Chapter 4, the
15 decompressed flow fields only persevering large-scaled turbulent structures can yield almost
16 the same prediction accuracy for the dispersion simulation. It is possible that the unsteady
17 simulation of adjoint equations could be successfully driven only by large-scale turbulent
18 structures. Under such circumstances, bridging models between RANS and LES like unsteady
19 RANS, partially-averaged Navier-Stokes, or detached eddy simulation have the opportunities
20 to simulate adjoint equations. About the data storage for the inverse simulation, apart from the
21 wavelet-compression method in the thesis, the large-scaled turbulent structures can also be
22 stored by low-dimensional modes like proper orthogonal decomposition or dynamic mode
23 decomposition.

24 Furthermore, the case study in the thesis is a regular block-arrayed building group model,
25 which is an ideal situation for dispersion in the urban environment. Dense buildings with
26 diverse geometries and irregular distribution may cause an unexpected challenge to STE. The
27 moving objects like automobiles, trains, and people will make the source-receptor relationship
28 more complicated to predict in the neighborhood scale. Their influence on the accuracy of STE
29 should be addressed in the future, and if the influence was unignorable, it is also necessary to
30 study how to model them in the source-receptor relationship.

31 Another critical factor to source-receptor relationship modeling is the meteorological

1 condition, which was assumed to be stationary during the dispersion emergency in the thesis.
2 However, in the reality, meteorological conditions like wind speed and directions are unsteady
3 at different times. It is likely that the inflow condition would suddenly change during the
4 dispersion emergency. Hence, it is valuable to extend the current method to handle unsteady
5 dispersions.

6 Meanwhile, the mechanisms of pollutants' dispersion are extremely complicated. The
7 particles may go through processes like collision, sediment, condensation, and so on. Some
8 toxic gases may also react with other gases during the dispersion. For these pollutants, the
9 current method may totally fail because of the passive scalars assumption. Even the model for
10 forward simulation is still on research to predict this kind of complicated dispersion. It is a huge
11 challenge for STE to cover a comprehensive dispersion mechanism.

12 Last but not least, the proposed methods in the thesis have not been testified by field test.
13 In fact, available databases of field experiments are limited in the literature. Most STE research
14 used the data from wind tunnel experiments or numerical experiments, which are relatively
15 ideal compared to real scenarios. One of the most famous field tests is the Mock Urban Setting
16 Test, where the containers are regularly placed in an open area to design an ideal urban model.
17 However, the field experiment with dispersion in the real urban area is far too sparse, even
18 though it is beneficial to future research like new method development, validation, and noise
19 analysis for STE. It is expected that such a field test database could be built.

20 The above is a brief description of possible research contents to develop a more effective
21 STE method for the urban environment in the future.

22

23

Reference

- Abd Razak, A., Hagishima, A., Ikegaya, N., Tanimoto, J., 2013. Analysis of airflow over building arrays for assessment of urban wind environment. *Building and Environment* 59, 56–65. <https://doi.org/10.1016/J.BUILDENV.2012.08.007>
- Abida, R., Bocquet, M., Vercauteren, N., Isnard, O., 2008. Design of a monitoring network over France in case of a radiological accidental release. *Atmospheric Environment* 42, 5205–5219. <https://doi.org/10.1016/j.atmosenv.2008.02.065>
- Allen, C.T., Haupt, S.E., Young, G.S., 2007. Source Characterization with a Genetic Algorithm–Coupled Dispersion–Backward Model Incorporating SCIPUFF. *Journal of Applied Meteorology and Climatology* 46, 273–287. <https://doi.org/10.1175/JAM2459.1>
- Altinel, I.K., Aras, N., Güney, E., Ersoy, C., 2008. Binary integer programming formulation and heuristics for differentiated coverage in heterogeneous sensor networks. *Computer Networks* 52, 2419–2431. <https://doi.org/10.1016/j.comnet.2008.05.002>
- Araki, S., Iwahashi, K., Shimadera, H., Yamamoto, K., Kondo, A., 2015. Optimization of air monitoring networks using chemical transport model and search algorithm. *Atmospheric Environment* 122, 22–30. <https://doi.org/10.1016/j.atmosenv.2015.09.030>
- Aubinet, M., Vesala, T., Papale, D. (Eds.), 2012. *Eddy Covariance: A Practical Guide to Measurement and Data Analysis*. Springer Netherlands, Dordrecht. <https://doi.org/10.1007/978-94-007-2351-1>
- Balczó, M., Lajos, T., 2015. Flow and dispersion phenomena in a simplified urban square. *Periodica Polytechnica Civil Engineering* 59, 347–360. <https://doi.org/10.3311/PPci.7852>
- BBC, 2019. The Amazon in Brazil is on fire - how bad is it? [WWW Document]. URL <https://www.bbc.com/news/world-latin-america-49433767>
- Berkooz, G., Holmes, P., Lumley, L.J., 1993. The Proper Orthogonal Decomposition in the Analysis of Turbulent Flows. *Annual Review of Fluid Mechanics* 25, 539–575. <https://doi.org/10.1146/annurev.fluid.25.1.539>

Blocken, B., 2018. LES over RANS in building simulation for outdoor and indoor applications: A foregone conclusion?, *Building Simulation*. <https://doi.org/10.1007/s12273-018-0459-3>

Blocken, B., Stathopoulos, T., Carmeliet, J., 2007. CFD simulation of the atmospheric boundary layer: wall function problems. *Atmospheric Environment* 41, 238–252. <https://doi.org/10.1016/j.atmosenv.2006.08.019>

Branford, S., Coceal, O., Thomas, T.G., Belcher, S.E., 2011a. Dispersion of a Point-Source Release of a Passive Scalar Through an Urban-Like Array for Different Wind Directions. *Boundary-Layer Meteorology* 139, 367–394. <https://doi.org/10.1007/s10546-011-9589-1>

Branford, S., Coceal, O., Thomas, T.G., Belcher, S.E., 2011b. Dispersion of a Point-Source Release of a Passive Scalar Through an Urban-Like Array for Different Wind Directions. *Boundary-Layer Meteorology* 139, 367–394. <https://doi.org/10.1007/s10546-011-9589-1>

Chen, C., Liu, W., Lin, C.-H., Chen, Q., 2015. A Markov chain model for predicting transient particle transport in enclosed environments. *Building and Environment* 90, 30–36. <https://doi.org/10.1016/j.buildenv.2015.03.024>

Cheng, H., Castro, I.P., 2002. Near wall flow over urban-like roughness. *Boundary-Layer Meteorology* 104, 229–259. <https://doi.org/10.1023/A:1016060103448>

Claus, J., Coceal, O., Thomas, T.G., Branford, S., Belcher, S.E., Castro, I.P., 2012. Wind-Direction Effects on Urban-Type Flows. *Boundary-Layer Meteorology* 142, 265–287. <https://doi.org/10.1007/s10546-011-9667-4>

Coceal, O., Dobre, A., Thomas, T.G., 2007. Unsteady dynamics and organized structures from DNS over an idealized building canopy. *International Journal of Climatology* 27, 1943–1953. <https://doi.org/10.1002/joc.1549>

Coceal, O., Goulart, E.V., Branford, S., Glyn Thomas, T., Belcher, S.E., 2014a. Flow structure and near-field dispersion in arrays of building-like obstacles. *Journal of Wind Engineering and Industrial Aerodynamics* 125, 52–68. <https://doi.org/10.1016/J.JWEIA.2013.11.013>

- Coceal, O., Goulart, E. V., Branford, S., Glyn Thomas, T., Belcher, S.E., 2014b. Flow structure and near-field dispersion in arrays of building-like obstacles. *Journal of Wind Engineering and Industrial Aerodynamics* 125, 52–68.
<https://doi.org/10.1016/J.JWEIA.2013.11.013>
- Coceal, O., Thomas, T.G., Castro, I.P., Belcher, S.E., 2006. Mean flow and turbulence statistics over groups of urban-like cubical obstacles. *Boundary-Layer Meteorology* 121, 491–519. <https://doi.org/10.1007/s10546-006-9076-2>
- Combest, D.P., Ramachandran, P.A., Dudukovic, M.P., 2011. On the Gradient Diffusion Hypothesis and Passive Scalar Transport in Turbulent Flows. *Industrial & Engineering Chemistry Research* 50, 8817–8823. <https://doi.org/10.1021/ie200055s>
- Constantine, P.G., Kent, C., Bui-Thanh, T., 2016. Accelerating Markov Chain Monte Carlo with Active Subspaces. *SIAM Journal on Scientific Computing* 38, A2779–A2805. <https://doi.org/10.1137/15M1042127>
- Cover, T.M., Thomas, J.A., 2006. *Elements of Information Theory*. Wiley-Interscience, USA.
- Cui, J., Lang, J., Chen, T., Cheng, S., Shen, Z., Mao, S., 2019. Investigating the impacts of atmospheric diffusion conditions on source parameter identification based on an optimized inverse modelling method. *Atmospheric Environment* 205, 19–29.
<https://doi.org/10.1016/j.atmosenv.2019.02.035>
- Daubechies, I., Sweldens, W., 1998. Factoring wavelet transforms into lifting steps. *The Journal of Fourier Analysis and Applications* 4, 247–269.
<https://doi.org/10.1007/BF02476026>
- de Freitas, N., Hojen-Sorensen, P., Jordan, M.I., Russell, S., 2013. Variational MCMC.
- Efthimiou, G.C., Kovalets, I.V., Argyropoulos, C.D., Venetsanos, A., Andronopoulos, S., Kakosimos, K.E., 2018a. Evaluation of an inverse modelling methodology for the prediction of a stationary point pollutant source in complex urban environments. *Building and Environment* 143, 107–119. <https://doi.org/10.1016/J.BUILDENV.2018.07.003>
- Efthimiou, G.C., Kovalets, I. V., Argyropoulos, C.D., Venetsanos, A., Andronopoulos, S., Kakosimos, K.E., 2018b. Evaluation of an inverse modelling methodology for the

prediction of a stationary point pollutant source in complex urban environments. *Building and Environment* 143, 107–119. <https://doi.org/10.1016/J.BUILDENV.2018.07.003>

Fuentes, M., Chaudhuri, A., Holland, D.M., 2007. Bayesian entropy for spatial sampling design of environmental data. *Environmental and Ecological Statistics* 14, 323–340. <https://doi.org/10.1007/s10651-007-0017-0>

Gelfand, A.E., Sahu, S.K., 1994. On markov chain monte carlo acceleration. *Journal of Computational and Graphical Statistics* 3, 261–276. <https://doi.org/10.1080/10618600.1994.10474644>

Gilks, W.R., Richardson, S., Spiegelhalter, D., 1995. *Markov Chain Monte Carlo in Practice*, Chapman & Hall/CRC Interdisciplinary Statistics. Taylor & Francis.

Graham, J., Kanov, K., Yang, X.I.A., Lee, M., Malaya, N., Lalescu, C.C., Burns, R., Eyink, G., Szalay, A., Moser, R.D., Meneveau, C., 2016. A web services accessible database of turbulent channel flow and its use for testing a new integral wall model for LES. *Journal of Turbulence* 17, 181–215. <https://doi.org/10.1080/14685248.2015.1088656>

Gromke, C., Buccolieri, R., Di Sabatino, S., Ruck, B., 2008. Dispersion study in a street canyon with tree planting by means of wind tunnel and numerical investigations – Evaluation of CFD data with experimental data. *Atmospheric Environment* 42, 8640–8650. <https://doi.org/10.1016/J.ATMOSENV.2008.08.019>

Hadjidoukas, P., Wermelinger, F., 2019. A Parallel Data Compression Framework for Large Scale 3D Scientific Data.

Hastings, W.K., 1970. Monte carlo sampling methods using Markov chains and their applications. *Biometrika* 57, 97–109. <https://doi.org/10.1093/biomet/57.1.97>

Hellsten, A., Luukkonen, S.M., Steinfeld, G., Kanani-Sühring, F., Markkanen, T., Järvi, L., Lento, J., Vesala, T., Raasch, S., 2015. Footprint Evaluation for Flux and Concentration Measurements for an Urban-Like Canopy with Coupled Lagrangian Stochastic and Large-Eddy Simulation Models. *Boundary-Layer Meteorology* 157, 191–217. <https://doi.org/10.1007/s10546-015-0062-4>

Hutchinson, M., Oh, H., Chen, W.-H., 2017. A review of source term estimation methods for atmospheric dispersion events using static or mobile sensors. *Information Fusion*

36, 130–148. <https://doi.org/10.1016/j.inffus.2016.11.010>

Ikegaya, N., Ikeda, Y., Hagishima, A., Razak, A.A., Tanimoto, J., 2017. A prediction model for wind speed ratios at pedestrian level with simplified urban canopies. *Theoretical and Applied Climatology* 127, 655–665. <https://doi.org/10.1007/s00704-015-1655-z>

Ikegaya, N., Okaze, T., Kikumoto, H., Imano, M., Ono, H., Tominaga, Y., 2019. Effect of the numerical viscosity on reproduction of mean and turbulent flow fields in the case of a 1:1:2 single block model. *Journal of Wind Engineering and Industrial Aerodynamics* 191, 279–296. <https://doi.org/10.1016/J.JWEIA.2019.06.013>

Issartel, J.-P., 2005a. Emergence of a tracer source from air concentration measurements, a new strategy for linear assimilation. *Atmos. Chem. Phys.* 5, 249–273. <https://doi.org/10.5194/acp-5-249-2005>

Issartel, J.-P., 2005b. Emergence of a tracer source from air concentration measurements, a new strategy for linear assimilation. *Atmospheric Chemistry and Physics* 5, 249–273. <https://doi.org/10.5194/acp-5-249-2005>

Issartel, J.-P., 2003. Rebuilding sources of linear tracers after atmospheric concentration measurements. *Atmos. Chem. Phys.* 3, 2111–2125. <https://doi.org/10.5194/acp-3-2111-2003>

Jacob, J., Sagaut, P., 2018. Wind comfort assessment by means of large eddy simulation with lattice Boltzmann method in full scale city area. *Building and Environment* 139, 110–124. <https://doi.org/10.1016/j.buildenv.2018.05.015>

Kanov, K., Burns, R., Lalescu, C., Eyink, G., 2015. The Johns Hopkins Turbulence Databases: An Open Simulation Laboratory for Turbulence Research. *Computing in Science and Engineering* 17, 10–17. <https://doi.org/10.1109/MCSE.2015.103>

Keats, A., Yee, E., Lien, F.-S., 2010. Information-driven receptor placement for contaminant source determination. *Environmental Modelling & Software* 25, 1000–1013. <https://doi.org/10.1016/j.envsoft.2010.01.006>

Keats, A., Yee, E., Lien, F.-S., 2007a. Bayesian inference for source determination with applications to a complex urban environment. *Atmospheric Environment* 41, 465–479. <https://doi.org/10.1016/j.atmosenv.2006.08.044>

Keats, A., Yee, E., Lien, F.-S., 2007b. Bayesian inference for source determination with

applications to a complex urban environment. *Atmospheric Environment* 41, 465–479.
<https://doi.org/10.1016/j.atmosenv.2006.08.044>

Kikumoto, H., Ooka, R., 2012. A numerical study of air pollutant dispersion with bimolecular chemical reactions in an urban street canyon using large-eddy simulation. *Atmospheric Environment* 54, 456–464. <https://doi.org/10.1016/j.atmosenv.2012.02.039>

Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P., 1983. Optimization by Simulated Annealing. *Science* 220, 671–680. <https://doi.org/10.1126/science.220.4598.671>

Kolomenskiy, D., Onishi, R., Uehara, H., 2018. Data Compression for Environmental Flow Simulations.

Kopka, P., Wawrzynczak, A., Borysiewicz, M., 2016. Application of the Approximate Bayesian Computation methods in the stochastic estimation of atmospheric contamination parameters for mobile sources. *Atmospheric Environment* 145, 201–212.
<https://doi.org/10.1016/J.ATMOSENV.2016.09.029>

Kormi, T., Mhadhebi, S., Bel Hadj Ali, N., Abichou, T., Green, R., 2018. Estimation of fugitive landfill methane emissions using surface emission monitoring and Genetic Algorithms optimization. *Waste Management* 72, 313–328.
<https://doi.org/10.1016/j.wasman.2016.11.024>

Kouichi, H., Ngae, P., Kumar, P., Feiz, A.-A., Bekka, N., 2019. An optimization for reducing the size of an existing urban-like monitoring network for retrieving an unknown point source emission. *Geoscientific Model Development* 12, 3687–3705.
<https://doi.org/10.5194/gmd-12-3687-2019>

Kumar, P., Feiz, A.-A., Singh, S.K., Ngae, P., Turbelin, G., 2015a. Reconstruction of an atmospheric tracer source in an urban-like environment. *Journal of Geophysical Research: Atmospheres* 120, 12589–12604. <https://doi.org/10.1002/2015JD024110>

Kumar, P., Feiz, A.-A., Singh, S.K., Ngae, P., Turbelin, G., 2015b. Reconstruction of an atmospheric tracer source in an urban-like environment. *Journal of Geophysical Research: Atmospheres* 120, 12589–12604. <https://doi.org/10.1002/2015JD024110>

Li, F., Liu, X., Liu, J., Cai, H., Wang, H., Zhang, K., Dai, C., 2020. Solutions to mitigate the impact of measurement noise on the air pollution source strength estimation in a multi-

zone building. *Building Simulation* 13, 1329–1337. <https://doi.org/10.1007/s12273-020-0635-0>

Li, F., Niu, J., 2005. An inverse approach for estimating the initial distribution of volatile organic compounds in dry building material. *Atmospheric Environment* 39, 1447–1455. <https://doi.org/10.1016/j.atmosenv.2004.11.021>

Li, S., 2015. Evaluating the Efficacy of Wavelet Compression for Turbulent-Flow Data Visualization.

Loredo, T.J., 2004. Bayesian Adaptive Exploration. *arXiv* 707, 330–346. <https://doi.org/10.1063/1.1751377>

Ma, D., Tan, W., Zhang, Z., Hu, J., 2017. Parameter identification for continuous point emission source based on Tikhonov regularization method coupled with particle swarm optimization algorithm. *Journal of Hazardous Materials* 325, 239–250. <https://doi.org/10.1016/j.jhazmat.2016.11.071>

Martin, G.N., 1979. Range encoding: an algorithm for removing redundancy from a digitised message, in: *Video & Data Recording Conference*. Southampton, UK.

Martina, M., Masera, G., 2005. Low-complexity, efficient 9/7 wavelet filters implementation, in: *IEEE International Conference on Image Processing 2005*. IEEE, p. III–1000. <https://doi.org/10.1109/ICIP.2005.1530563>

Meroney, R.N., Pavageau, M., Rafailidis, S., Schatzmann, M., 1996. Study of line source characteristics for 2-D physical modelling of pollutant dispersion in street canyons. *Journal of Wind Engineering and Industrial Aerodynamics* 62, 37–56. [https://doi.org/10.1016/S0167-6105\(96\)00057-8](https://doi.org/10.1016/S0167-6105(96)00057-8)

Mons, V., Margheri, L., Chassaing, J.-C., Sagaut, P., 2017. Data assimilation-based reconstruction of urban pollutant release characteristics. *Journal of Wind Engineering and Industrial Aerodynamics* 169, 232–250. <https://doi.org/10.1016/j.jweia.2017.07.007>

Naik, A.K., Holambe, R.S., 2014. New Approach to the Design of Low Complexity 9/7 Tap Wavelet Filters With Maximum Vanishing Moments. *IEEE Transactions on Image Processing* 23, 5722–5732. <https://doi.org/10.1109/TIP.2014.2363733>

Ngae, P., Kouichi, H., Kumar, P., Feiz, A.A., Chpoun, A., 2019. Optimization of an

urban monitoring network for emergency response applications: An approach for characterizing the source of hazardous releases. *Quarterly Journal of the Royal Meteorological Society* 145, 967–981. <https://doi.org/10.1002/qj.3471>

Nicas, M., 2000. Markov Modeling of Contaminant Concentrations in Indoor Air. *AIHAJ - American Industrial Hygiene Association* 61, 484–491. <https://doi.org/10.1080/15298660008984559>

Parent, A., Morin, M., Lavigne, P., 1992. Propagation of super-Gaussian field distributions. *Optical and Quantum Electronics* 24, S1071–S1079. <https://doi.org/10.1007/BF01588606>

Pudykiewicz, J.A., 1998. Application of adjoint tracer transport equations for evaluating source parameters. *Atmospheric Environment* 32, 3039–3050. [https://doi.org/10.1016/S1352-2310\(97\)00480-9](https://doi.org/10.1016/S1352-2310(97)00480-9)

Qiu, Z., Lee, C.M., Xu, Z.H., Sui, L.N., 2016. A multi-resolution filtered-x LMS algorithm based on discrete wavelet transform for active noise control. *Mechanical Systems and Signal Processing* 66–67, 458–469. <https://doi.org/10.1016/j.ymssp.2015.05.024>

Rajaona, H., Septier, F., Armand, P., Delignon, Y., Olry, C., Albergel, A., Moussafir, J., 2015. An adaptive Bayesian inference algorithm to estimate the parameters of a hazardous atmospheric release. *Atmospheric Environment* 122, 748–762. <https://doi.org/10.1016/j.atmosenv.2015.10.026>

Rioul, O., Vetterli, M., 1991. Wavelets and signal processing. *IEEE Signal Processing Magazine* 8, 14–38. <https://doi.org/10.1109/79.91217>

Ruttimann, U.E., Pipberger, H. V., 1979. Compression of the ECG by Prediction or Interpolation and Entropy Encoding. *IEEE Transactions on Biomedical Engineering BME-26*, 613–623. <https://doi.org/10.1109/TBME.1979.326543>

Sakai, R., Sasaki, D., Obayashi, S., Nakahashi, K., 2013. Wavelet-based data compression for flow simulation on block-structured Cartesian mesh. *International Journal for Numerical Methods in Fluids* 73, 462–476. <https://doi.org/10.1002/flid.3808>

Salim, S.M., Buccolieri, R., Chan, A., Di Sabatino, S., 2011a. Numerical simulation of atmospheric pollutant dispersion in an urban street canyon: Comparison between RANS and

LES. *Journal of Wind Engineering and Industrial Aerodynamics* 99, 103–113.

<https://doi.org/10.1016/J.JWEIA.2010.12.002>

Salim, S.M., Cheah, S.C., Chan, A., 2011b. Numerical simulation of dispersion in urban street canyons with avenue-like tree plantings: Comparison between RANS and LES.

Building and Environment 46, 1735–1746.

<https://doi.org/10.1016/J.BUILDENV.2011.01.032>

Saunier, O., Bocquet, M., Mathieu, A., Isnard, O., 2009. Model reduction via principal component truncation for the optimal design of atmospheric monitoring networks.

Atmospheric Environment 43, 4940–4950. <https://doi.org/10.1016/j.atmosenv.2009.07.011>

Schmalzl, J., 2003. Using standard image compression algorithms to store data from computational fluid dynamics. *Computers and Geosciences* 29, 1021–1031.

[https://doi.org/10.1016/S0098-3004\(03\)00098-0](https://doi.org/10.1016/S0098-3004(03)00098-0)

Schmid, P.J., 2010. Dynamic mode decomposition of numerical and experimental data.

Journal of Fluid Mechanics 656, 5–28. <https://doi.org/10.1017/S0022112010001217>

Sharan, M., Singh, S.K., Issartel, J.P., 2012. Least Square Data Assimilation for Identification of the Point Source Emissions. *Pure Appl. Geophys.* 169, 483–497.

<https://doi.org/10.1007/s00024-011-0382-3>

Skodras, A., Christopoulos, C., Ebrahimi, T., 2001. The JPEG 2000 still image compression standard. *IEEE Signal Processing Magazine* 18, 36–58.

<https://doi.org/10.1109/79.952804>

Song, M.-S., 2008. Entropy Encoding in Wavelet Image Compression, in: Jorgensen, P.E.T., Merrill, K.D., Packer, J.A. (Eds.), *Representations, Wavelets, and Frames*. Birkhäuser Boston, Boston, MA, pp. 293–311. https://doi.org/10.1007/978-0-8176-4683-7_14

Thomson, L.C., Hirst, B., Gibson, G., Gillespie, S., Jonathan, P., Skeldon, K.D., Padgett, M.J., 2007. An improved algorithm for locating a gas source using inverse methods.

Atmospheric Environment 41, 1128–1134. <https://doi.org/10.1016/j.atmosenv.2006.10.003>

Tominaga, Y., Mochida, A., Murakami, S., Sawaki, S., 2008a. Comparison of various revised k- ϵ models and LES applied to flow around a high-rise building model with 1:1:2 shape placed within the surface boundary layer. *Journal of Wind Engineering and Industrial*

Aerodynamics 96, 389–411. <https://doi.org/10.1016/j.jweia.2008.01.004>

Tominaga, Y., Mochida, A., Yoshie, R., Kataoka, H., Nozu, T., Yoshikawa, M., Shirasawa, T., 2008b. AIJ guidelines for practical applications of CFD to pedestrian wind environment around buildings. *Journal of Wind Engineering and Industrial Aerodynamics* 96, 1749–1761. <https://doi.org/10.1016/j.jweia.2008.02.058>

Tominaga, Y., Stathopoulos, T., 2013. CFD simulation of near-field pollutant dispersion in the urban environment: A review of current modeling techniques. *Atmospheric Environment* 79, 716–730. <https://doi.org/10.1016/j.atmosenv.2013.07.028>

Tominaga, Y., Stathopoulos, T., 2012. CFD Modeling of Pollution Dispersion in Building Array: Evaluation of turbulent scalar flux modeling in RANS model using LES results. *Journal of Wind Engineering and Industrial Aerodynamics* 104–106, 484–491. <https://doi.org/10.1016/J.JWEIA.2012.02.004>

Tominaga, Y., Stathopoulos, T., 2011a. CFD modeling of pollution dispersion in a street canyon: Comparison between LES and RANS. *Journal of Wind Engineering and Industrial Aerodynamics* 99, 340–348. <https://doi.org/10.1016/J.JWEIA.2010.12.005>

Tominaga, Y., Stathopoulos, T., 2011b. CFD modeling of pollution dispersion in a street canyon: Comparison between LES and RANS. *Journal of Wind Engineering and Industrial Aerodynamics* 99, 340–348. <https://doi.org/10.1016/J.JWEIA.2010.12.005>

Uehara, K., Murakami, S., Oikawa, S., Wakamatsu, S., 2000. Wind tunnel experiments on how thermal stratification affects flow in and above urban street canyons. *Atmospheric Environment* 34, 1553–1562. [https://doi.org/10.1016/S1352-2310\(99\)00410-0](https://doi.org/10.1016/S1352-2310(99)00410-0)

Usevitch, B.E., 2001. A tutorial on modern lossy wavelet image compression: Foundations of JPEG 2000. *IEEE Signal Processing Magazine* 18, 22–35. <https://doi.org/10.1109/79.952803>

Villasenor, J.D., Belzer, B., Liao, J., 1995. Wavelet Filter Evaluation for Image Compression. *IEEE Transactions on Image Processing* 4, 1053–1060. <https://doi.org/10.1109/83.403412>

Vrugt, J.A., ter Braak, C.J.F., Diks, C.G.H., Robinson, B.A., Hyman, J.M., Higdón, D., 2009. Accelerating Markov Chain Monte Carlo Simulation by Differential Evolution with

Self-Adaptive Randomized Subspace Sampling. *International Journal of Nonlinear Sciences and Numerical Simulation* 10, 1–4. <https://doi.org/10.1515/IJNSNS.2009.10.3.273>

Wade, D., Senocak, I., 2013a. Stochastic reconstruction of multiple source atmospheric contaminant dispersion events. *Atmospheric Environment* 74, 45–51.
<https://doi.org/10.1016/J.ATMOSENV.2013.02.051>

Wade, D., Senocak, I., 2013b. Stochastic reconstruction of multiple source atmospheric contaminant dispersion events. *Atmospheric Environment* 74, 45–51.
<https://doi.org/10.1016/J.ATMOSENV.2013.02.051>

Wang, J., Wang, B., Liu, J., Cheng, W., Zhang, J., 2021. An inverse method to estimate the source term of atmospheric pollutant releases. *Atmospheric Environment* 118554.
<https://doi.org/10.1016/j.atmosenv.2021.118554>

Wang, Y., Huang, H., Huang, L., Zhang, X., 2018. Source term estimation of hazardous material releases using hybrid genetic algorithm with composite cost functions. *Engineering Applications of Artificial Intelligence* 75, 102–113.
<https://doi.org/10.1016/J.ENGAPPAL.2018.08.005>

Wilke, C.R., Lee, C.Y., 1955. Estimation of Diffusion Coefficients for Gases and Vapors. *Industrial & Engineering Chemistry* 47, 1253–1257.
<https://doi.org/10.1021/ie50546a056>

Wilson, J.P., 2002. Wavelet-based lossy compression of barotropic turbulence simulation data, in: *Proceedings DCC 2002. Data Compression Conference*. pp. 479–.
<https://doi.org/10.1109/DCC.2002.1000022>

Wu, L., Bocquet, M., 2011. Optimal redistribution of the background ozone monitoring stations over France. *Atmospheric Environment* 45, 772–783.
<https://doi.org/10.1016/j.atmosenv.2010.08.038>

Xie, Z., Castro, I.P., 2006. LES and RANS for turbulent flow over arrays of wall-mounted obstacles. *Flow, Turbulence and Combustion* 76, 291–312.
<https://doi.org/10.1007/s10494-006-9018-6>

Xue, F., Kikumoto, H., Li, X., Ooka, R., 2018a. Bayesian source term estimation of atmospheric releases in urban areas using LES approach. *Journal of Hazardous Materials* 349,

68–78. <https://doi.org/10.1016/j.jhazmat.2018.01.050>

Xue, F., Kikumoto, H., Li, X., Ooka, R., 2018b. Bayesian source term estimation of atmospheric releases in urban areas using LES approach. *Journal of Hazardous Materials* 349, 68–78. <https://doi.org/10.1016/j.jhazmat.2018.01.050>

Xue, F., Li, X., Ooka, R., Kikumoto, H., Zhang, W., 2017. Turbulent Schmidt number for source term estimation using Bayesian inference. *Building and Environment* 125, 414–422. <https://doi.org/10.1016/j.buildenv.2017.09.012>

Yakhot, V., Orszag, S.A., Thangam, S., Gatski, T.B., Speziale, C.G., 1992. Development of turbulence models for shear flows by a double expansion technique. *Physics of Fluids A* 4, 1510–1520. <https://doi.org/10.1063/1.858424>

Yee, E., 2012. Probability theory as logic: Data assimilation for multiple source reconstruction. *Pure and Applied Geophysics* 169, 499–517. <https://doi.org/10.1007/s00024-011-0384-1>

Yee, E., Hoffman, I., Ungar, K., 2014. Bayesian Inference for Source Reconstruction: A Real-World Application. *International Scholarly Research Notices* 2014, 1–12. <https://doi.org/10.1155/2014/507634>

Zarzhitsky, D., Spears, D.F., 2005. Swarm approach to chemical source localization, in: *Conference Proceedings - IEEE International Conference on Systems, Man and Cybernetics*.

Zheng, X., Chen, Z., 2010. Back-calculation of the strength and location of hazardous materials releases using the pattern search method. *Journal of Hazardous Materials* 183, 474–481. <https://doi.org/10.1016/j.jhazmat.2010.07.048>

Zubair, L., Sreenivasan, K.R., Wickerhauser, M.V., 1992. Characterization and Compression of Turbulent Signals and Images Using Wavelet-Packets, in: Gatski, T.B., Speziale, C.G., Sarkar, S. (Eds.), *Studies in Turbulence*. Springer New York, New York, NY, pp. 489–513. https://doi.org/10.1007/978-1-4612-2792-2_37

Appendix

Appendix A. Numerical simulation for turbulence and dispersion of pollution

A.1 Fluid and turbulence

Atmospheric pollution is caused by the dispersion of pollutants in the air, which is a mass transfer behavior in fluids. Fluids like gases and liquids are defined as substances that cannot permanently keep their shapes under certain stress. Fluids also do not have a definite shape. They can be divided into Newtonian or non-Newtonian based on the relationship between the shear stress and shear rate. Most gases are Newtonian whose relationship is linear with the slope of molecular viscosity. Hence, only the Newtonian fluids are discussed in what follows.

Fluid flow has two states: laminar and turbulent. According to the fluid experiments, a key factor to decide the state of the fluid is Reynolds number $Re = \frac{UL}{\nu}$. Here ν is the dynamic viscosity coefficient. U is the standard velocity of the flow. L is the standard length of the flow. Hence, Re measures the relative power of inertia forces and viscous force on the state of fluids. In fact, for each kind of fluid, there exists a critical Reynolds number Re_c . When Re is smaller than Re_c , the flow is smooth and adjacent parts slide past each other in order. If the boundary condition of the flow does not change with time, the flow will keep the steady state everywhere in the domain and it is called laminar flow. If Re gradually increases after it reaches Re_c , the laminar flow will go through an important state called transition, where complicated events happen to change the flow into turbulence. In this state, turbulent structures with different sizes appear spatially and temporally. They include the information of nonlinear correlation among a certain space of the flow. This nonlinearity is mathematically represented by the advection term (see Eq. (A.2)). However, if a single point was monitored, the flow velocity there is totally random and chaotic. It is difficult to build a physical system for such a flow and find the mathematical description of it.

From a microcosmic point of view, fluid is a system consisting of numerous small molecular even though it is a continuum in the human's eyes. In order to describe the system in

mathematics, Eulerian and Lagrangian approaches were used. In the Lagrangian approach, the fluid system is divided into small fluid parcels. When the system is changing through space and time, the location and velocity of each parcel are followed to describe the system. In contrast, the Eulerian approach focuses on the properties of specific locations in the domain. Its view will not follow any fluid parcels, but it cares about the flow velocity somewhere as time changes. Although the two approaches are based on totally different concepts, they are both effective in the mathematical description of fluids. In this thesis, fluids are processed from the Eulerian view, and the corresponding theories are briefly introduced here.

A.2 Governing equations for fluid and dispersion

In this study, the dispersion of pollution is assumed to occur in an incompressible, isothermal fluid flow. Besides, the pollutants are assumed as the passive scalar which will not affect the behavior of fluids. Under this circumstance, the Eulerian descriptions of the fluids and the dispersion are the following three equations: continuity equation, momentum equation (or Navier-Stokes equation), and transport equation.

$$\frac{\partial u_i}{\partial x_i} = 0 \quad (\text{A.1})$$

$$\frac{\partial u_i}{\partial t} + \frac{\partial u_j u_i}{\partial x_j} = -\frac{1}{\rho} \frac{\partial p}{\partial x_i} + \frac{\partial}{\partial x_j} \left\{ \nu \left(\frac{\partial u_i}{\partial x_j} + \frac{\partial u_j}{\partial x_i} \right) \right\} \quad (\text{A.2})$$

$$\frac{\partial C}{\partial t} + \frac{\partial u_j C}{\partial x_j} = \frac{\partial}{\partial x_j} \left(D_m \frac{\partial C}{\partial x_j} \right) + S \quad (\text{A.3})$$

Here, u_i is the instantaneous velocity in the x_i direction. ρ is the constant density of fluids. p is the instantaneous pressure. ν is the dynamic viscosity coefficient. C is the concentration. D_m is the molecular diffusion coefficient. S is the source term of the passive scalar. Eq. (A.1) follows the conservation rule of mass. Eq. (A.2) follows the conservation rule of momentum. Eq. (A.3) describes the transport process of passive scalar released from the source.

Although the governing equations for fluids and dispersion have been established, and the number of unknown variables equals that of available equations, it is still one of the most challenging mathematical problems to find the analytical solution for them. The reason is that Navier-Stokes equation is a non-linear partial differential equation. The second term in the left hand side of Eq. (A.2), which is called as advection term, represents the nonlinear coupling of velocities in two directions.

In this case, a compromise approach is to solve it by a numerical method like the finite volume method. In this method, the calculation domain is subdivided into finite small, discrete volumes by meshing, then these partial differential equations are transferred into a huge set of algebraic equations by integration for each volume. The numerical solutions can be calculated by simultaneously solving this set of equations for a short time step with time marching. It is straightforward to notice that the success of numerical simulation needs a considerable amount of computational resources for the iterative calculation of algebraic equations.

In recent years, the development of high performance computers makes numerical simulation realizable. The computational fluid dynamics (CFD) technique has become a powerful tool in industrial applications and research of fluid mechanisms. Therefore, this thesis models the source-receptor relationship and pollution dispersion by CFD.

In theory, CFD can provide a nearly perfect solution for Eqs. (A.1-A.3), and such technique does exist with a name of direct numerical simulation (DNS). However, the reality is not so easy.

On the one hand, in the turbulent flow field of the atmospheric environment, the range of sizes of turbulence structures is considerably wide. During the movement of the flow field, the large turbulence will break into smaller turbulence. The kinetic energy is also transferred into smaller turbulence, which is called the energy cascade. This process will continue to pass the energy to smaller turbulence until it reaches the smallest size and eventually dissipates in the form of thermal energy caused by the molecular viscosity. Therefore, in order to completely simulate the turbulence flow, it is necessary to reproduce all the structures from the largest ones to the smallest ones.

On the other hand, in CFD, the partial differential equations are integrated into algebraic equations for each mesh (volume). In this case, the turbulence structures smaller than the mesh cannot be simulated. The mesh should be made extremely fine to capture all the structures. According to the previous research, the mesh number needs to have the order of $Re^{9/4}$. A common case in wind engineering can be checked here. If the flow field around a building with $L = 15m$ height caused by a coming flow with $U = 1m/s$ is simulated by DNS, with the dynamic viscosity coefficient of air $\nu = 1.5 \times 10^{-6}m^2/s$, Re is 10^6 and the number of mesh astonishingly becomes about 3.1×10^{13} . The resultant calculation burden is unaffordable even with the most advanced computers right now. Therefore, DNS has only been

used to simulate the flow field with a small Reynolds number. Most engineering simulations are constrained.

Since a complete simulation of all turbulence structures with CFD is almost impossible for high Reynolds number, in order to realize the CFD for engineering applications, alternative approaches have been developed by researchers. One common method is that rather than all turbulence, only the turbulence of interest with relatively large sizes are explicitly simulated, while the small turbulence are modeled. This is so-called turbulence modeling. In the past decades, among plenty of modeling methods, two popular turbulence models are the Reynolds averaged Navier-Stokes (RANS) model and the large-eddy simulation (LES) model. These two models are also used in this thesis and will be introduced in the following contents.

A.3 Large eddy simulation

A.3.1 Filtering process

The principal idea behind LES is to reduce the computational cost by modeling the small turbulence structures, which are the most computationally expensive to resolve, via low-passing filtering of the Navier-Stokes equations. In other words, for any wanted physical variables $f(x, t)$, LES will decompose it into the grid scale (GS) component $\tilde{f}(x, t)$ which can be explicitly resolved by the mesh, and the sub-grid scale (SGS) component $f''(x, t)$ which is filtered out and implicitly modeled.

$$f(x, t) = \tilde{f}(x, t) + f''(x, t) \quad (\text{A.4})$$

The separation for GS and SGS components is operated by filter function $G(\xi)$, which is defined by:

$$\tilde{f}(x, t) = \int_{-\infty}^{\infty} G(\xi) f(x - \xi) d\xi \quad (\text{A.5})$$

The filter function has to satisfy the following requirements:

$$\lim_{\xi \rightarrow \pm\infty} G(\xi) = 0 \quad (\text{A.6})$$

$$\int_{-\infty}^{\infty} G(\xi) d\xi = 1 \quad (\text{A.7})$$

Common filter functions in LES include top hat filter, Gaussian filter, and sharp cut-off filter. After the filtering, the filtered momentum equation becomes:

$$\frac{\partial \tilde{u}_i}{\partial t} + \frac{\partial \tilde{u}_i \tilde{u}_j}{\partial x_j} = -\frac{1}{\rho} \frac{\partial \tilde{p}}{\partial x_i} + \frac{1}{\nu} \frac{\partial^2 \tilde{u}_i}{\partial x_j \partial x_j} \quad (\text{A.8})$$

A.3.2 Sub-grid scaled stress modeling

The form of Eq. (A.8) is different from that of Eq. (A.2) because the filtered product $\tilde{u}_i \tilde{u}_j$ is different from the product of the filtered velocity $\tilde{u}_i \tilde{u}_j$. The difference is the residual-stress tensor defined by

$$T_{ij} = \tilde{u}_i \tilde{u}_j - \tilde{u}_i \tilde{u}_j \quad (\text{A.9})$$

This tensor is also called the SGS stress tensor. In this research, it is modeled by the standard Smagorinsky SGS model, according to which, the SGS tensor is split into an isotropic part $\frac{1}{3} T_{kk} \delta_{ij}$ and an anisotropic part $T_{ij} - \frac{1}{3} T_{kk} \delta_{ij}$, where δ_{ij} is the Kronecker delta.

$$\begin{aligned} T_{ij} &= \tilde{u}_i \tilde{u}_j - \tilde{u}_i \tilde{u}_j \\ &= \frac{1}{3} T_{kk} \delta_{ij} + \left(T_{ij} - \frac{1}{3} T_{kk} \delta_{ij} \right) \\ &\approx \frac{1}{3} T_{kk} \delta_{ij} - 2\nu_{sgs} dev(\tilde{D})_{ij} \\ &= \frac{2}{3} k_{sgs} \delta_{ij} - 2\nu_{sgs} dev(\tilde{D})_{ij} \end{aligned} \quad (\text{A.10})$$

Here ν_{sgs} is the SGS eddy viscosity. The resolved-scale strain rate tensor \tilde{D}_{ij} is defined as

$$\tilde{D}_{ij} = \frac{1}{2} \left(\frac{\partial \tilde{u}_i}{\partial x_j} + \frac{\partial \tilde{u}_j}{\partial x_i} \right) \quad (\text{A.11})$$

The SGS kinetic energy k_{sgs} is

$$k_{sgs} = \frac{1}{2} T_{kk} = \frac{1}{2} (\tilde{u}_i \tilde{u}_j - \tilde{u}_i \tilde{u}_j) \quad (\text{A.12})$$

The Smagorinsky SGS model is based on two elementary assumptions: eddy viscosity approximation and local equilibrium. In the first assumption, the anisotropic part in equation (A.10) is approximated by relating it to the resolved rate of the strain tensor \tilde{D}_{ij} .

$$T_{ij} - \frac{1}{3}T_{kk}\delta_{ij} \approx -2\nu_{sgs}dev(\tilde{D})_{ij} \quad (A.13)$$

The SGS scale viscosity is computed as:

$$\nu_{sgs} = C_k\Delta\sqrt{k_{sgs}} \quad (A.14)$$

Here C_k is a model constant with a default value of 0.094. Δ is the cubic root of the cell volume that defines the sub-grid length scale.

In the second assumption about local equilibrium, there is a balance between the sub-grid scale energy production and dissipation. The SGS kinetic energy k_{sgs} is calculated by this assumption.

$$\begin{aligned} \tilde{D}:T_{ij} + C_e\frac{k_{sgs}^{1.5}}{\Delta} &= 0 \\ \tilde{D}:\left(\frac{2}{3}k_{sgs}I - 2\nu_{sgs}dev(\tilde{D})\right) + C_e\frac{k_{sgs}^{1.5}}{\Delta} &= 0 \\ \tilde{D}:\left(\frac{2}{3}k_{sgs}I - 2C_k\Delta\sqrt{k_{sgs}}dev(\tilde{D})\right) + C_e\frac{k_{sgs}^{1.5}}{\Delta} &= 0 \\ \sqrt{k_{sgs}}\left(\frac{C_e}{\Delta}k_{sgs} + \frac{2}{3}tr(\tilde{D})\sqrt{k_{sgs}} - 2C_k\Delta(dev(\tilde{D}):\tilde{D})\right) &= 0 \\ ak_{sgs} + b\sqrt{k_{sgs}} - c &= 0 \\ k_{sgs} &= \left(\frac{-b + \sqrt{b^2 + 4ac}}{2a}\right)^2 \end{aligned} \quad (A.15)$$

The operator $:$ is a double inner product of two second-rank tensors which is the sum of the 9 products of the tensor components. I is the identity matrix. In Eq. (A.15),

$$\begin{aligned}
 a &= \frac{C_e}{\Delta} \\
 b &= \frac{2}{3} \text{tr}(\tilde{D}) \\
 c &= 2C_k \Delta (\text{dev}(\tilde{D}) : \tilde{D})
 \end{aligned} \tag{A.16}$$

In the case of incompressible flow, which is the same as this research, the Equation (A.16) reduces to

$$\begin{aligned}
 b &= \frac{2}{3} \text{tr}(\tilde{D}) = 0 \\
 c &= 2C_k \Delta (\text{dev}(\tilde{D}) : \tilde{D}) = C_k \Delta |\tilde{D}|^2
 \end{aligned} \tag{A.17}$$

where

$$|\tilde{D}| = \sqrt{2\tilde{D} : \tilde{D}} \tag{A.18}$$

By substituting the Eq. (A.17) into Eq. (A.15), we have:

$$k_{sgs} = \frac{c}{a} = \frac{C_k \Delta^2 |\tilde{D}|^2}{C_e} \tag{A.19}$$

The following expression can be obtained for the SGS eddy viscosity in the case of incompressible flows by substituting the Equation (A.19) into the Equation (A.14).

$$\nu_{sgs} = C_k \sqrt{\frac{C_k}{C_e}} \Delta^2 |\tilde{D}| \tag{A.20}$$

In the literature, the SGS eddy viscosity is commonly expressed as below.

$$\nu_{sgs} = (C_s \Delta)^2 |\tilde{D}| \tag{A.21}$$

We can get the relation between the Smagorinsky constant C_s and other coefficients.

$$C_s^2 = C_k \sqrt{\frac{C_k}{C_e}} \tag{A.22}$$

The default value of the Smagorinsky constant is 0.173. For the flow field around the solid structure, this constant is commonly set to 0.10~0.15. In this research, it was set as 0.12.

A.3.3 Van driest damping function

In the simulation, if the boundary condition is the non-slip wall, there should be no turbulence at the boundary. However, in the standard SGS model, the eddy viscosity is nonzero, which is contrary to reality. To fix this problem, one way is to add a Van Driest damping function into the length scale. The SGS eddy viscosity changes to:

$$\nu_{sgs} = (C_s f_s \Delta)^2 |\bar{D}| \quad (\text{A.23})$$

Here, the f_s is the Van Driest damping function.

$$f_s = 1 - \exp\left(\frac{-y^+}{A^+}\right) \quad (\text{A.24})$$

Here, A^+ is the Van Driest constant, which is set to 26 in this research. $y^+ = yu_\tau/\nu$ is the non-dimensional wall unit. u_τ is the shear velocity. In other words, the characteristic spatial length of the filter in the turbulent boundary layer is not necessarily related to the mesh size, but the minimum value between Δ and the one obtained from the damping function in equation (A.24) is locally adapted in space and time.

A.3.4 Dispersion equation for passive scalar

After the filtering and modeling for the momentum equation, the transport equation for passive scalar in LES can be dealt with in a similar way. By filtering operation, Eq. (3) becomes:

$$\frac{\partial \tilde{C}}{\partial t} + \frac{\partial \tilde{u}_j \tilde{C}}{\partial x_j} = \frac{\partial}{\partial x_j} \left(D_m \frac{\partial \tilde{C}}{\partial x_j} \right) + S \quad (\text{A.25})$$

The SGS turbulent diffusion term caused by filtering will also appear as:

$$\widetilde{u_i'' C''} = \widetilde{u_i C} - \tilde{u}_i \tilde{C} \quad (\text{A.26})$$

One of the most common ways to model this term is relating it to the gradient of GS concentration:

$$\widetilde{u_i'' C''} \approx -D_{sgs} \frac{\partial(\tilde{C})}{\partial x_j} \approx -\frac{\nu_{sgs}}{Sc_{sgs}} \frac{\partial(\tilde{C})}{\partial x_j} \quad (\text{A.27})$$

where Sc_{sgs} is the SGS turbulent Schmidt number. When the grid is fine enough, most of the turbulent diffusion is directly calculated and the modeled SGS is limited.

A.4 Reynolds averaged Navier-Stokes model

A.4.1 Reynolds average

Another turbulence model RANS can also be regarded as a filtering approach, but unlike LES where the filtering will change with the mesh size, the filtering in RANS is the ensemble average, which is constant in the space. The physical variable of flow f can be divided into ensemble-averaged component \bar{f} and fluctuation component f' .

$$f(x, t) = \bar{f}(x, t) + f'(x, t) \quad (\text{A.28})$$

Meanwhile, for the convenience of mathematical processing, it is required that

$$\bar{\bar{f}} = \bar{f}, \quad \bar{f'} = 0, \quad \overline{f'f} = 0 \quad (\text{A.29})$$

The ensemble average that satisfies the above relationships is called as Reynolds average.

If Reynolds average is operated on the continuity equation and momentum equation, they change to

$$\frac{\partial \bar{u}_i}{\partial x_i} = 0 \quad (\text{A.30})$$

$$\frac{\partial \bar{u}_i}{\partial t} + \frac{\partial \bar{u}_j \bar{u}_i}{\partial x_j} = -\frac{1}{\rho} \frac{\partial \bar{p}}{\partial x_i} + \frac{\partial}{\partial x_j} \left\{ \nu \left(\frac{\partial \bar{u}_i}{\partial x_j} + \frac{\partial \bar{u}_j}{\partial x_i} \right) \right\} - \frac{\partial \tau_{ij}}{\partial x_j} \quad (\text{A.31})$$

The extra term τ_{ij} is the Reynolds stress resulting from the Reynolds average.

$$\tau_{ij} = \overline{u'_i u'_j} = \bar{u}_i \bar{u}_j - \bar{u}_j \bar{u}_i \quad (\text{A.32})$$

It can be noticed that in order to simulate Eq. (A.30 & A.31) without closure problem, the Reynolds stress needs to be modeled, which is similar to the SGS stress tensor. Several approaches have been developed to model Reynolds stress under the RANS framework, as an instance, here one of the most used models standard $k - \varepsilon$ model is introduced.

A.4.2 Standard $k - \varepsilon$ model

The Standard $k - \varepsilon$ model is based on the assumption that there is an analogy between viscous stress and Reynold stress. In 1877, Boussinesq proposed that the Reynolds stress could be expressed as stress that is proportional to the mean rate of deformation.

$$\tau_{ij} = \nu_t \left(\frac{\partial \bar{u}_i}{\partial x_j} + \frac{\partial \bar{u}_j}{\partial x_i} \right) - \frac{2}{3} k \delta_{ij} \quad (\text{A.33})$$

where ν_t is the dynamic turbulent viscosity coefficient. k is the turbulent kinetic energy defined as $\frac{1}{2} (\overline{u_x'^2} + \overline{u_y'^2} + \overline{u_z'^2})$ or $\frac{1}{2} \tau_{ii}$. In the next step, it is necessary to provide the modeling expression for ν_t and k .

In fact, the governing equation for k can be derived. First, we multiple the Navier-Stokes equation Eq. (A.2) with the fluctuation velocity components and add them together. Then, the same process can be repeated again on the averaged equations Eq. (A.31). At last, we subtract the resultant two summed equations. Because the process is too long, only the final result is presented here:

$$\frac{\partial k}{\partial t} + \frac{\partial \bar{u}_j k}{\partial x_j} = P_{ij} - \varepsilon_{ij} + \mathbb{T}_{ijk}^p + \mathbb{T}_{ijk}^v + \mathbb{T}_{ijk}^R \quad (\text{A.34a})$$

$$P_{ij} = \tau_{ij} \cdot S_{ij} \quad (\text{A.34b})$$

$$\varepsilon_{ij} = \overline{\nu S_{ij}' \cdot S_{ij}'} \quad (\text{A.34c})$$

$$\mathbb{T}_{ijk}^p = -\frac{1}{\rho} (\overline{p' u_i' \delta_{jk}} + \overline{p' u_j' \delta_{ik}}) \quad (\text{A.34d})$$

$$\mathbb{T}_{ijk}^v = \nu \frac{\partial \tau_{ij}}{\partial x_k} \quad (\text{A.34e})$$

$$\mathbb{T}_{ijk}^R = -\overline{u_i' u_j' u_k'} \quad (\text{A.34f})$$

$$S_{ij} = \frac{\partial \bar{u}_i}{\partial x_j} + \frac{\partial \bar{u}_j}{\partial x_i}, \quad s_{ij} = \frac{\partial u_i'}{\partial x_j} + \frac{\partial u_j'}{\partial x_i} \quad (\text{A.34g})$$

Here, P_{ij} is the rate of production of k . ε_{ij} is the rate of dissipation of k . \mathbb{T}_{ijk}^p is the transport of k by pressure. \mathbb{T}_{ijk}^v is the transport of k by viscous stresses. \mathbb{T}_{ijk}^R is the transport of k by Reynolds stresses.

As for ν_t , it was modeled based on the dimensional analysis.

$$\nu_t = C_\mu \frac{k^2}{\varepsilon} \quad (\text{A.35})$$

where C_μ is a dimensionless constant.

Noticing that $k = \frac{1}{2} \tau_{ii}$ and using Eq. (A.33), the production term for k can be expressed as:

$$P_{ij} = 2\nu_t S_{ij} \cdot S_{ij} \quad (\text{A.36})$$

Besides, the transport effects of pressure and Reynolds stresses can be approximated together by

$$\mathbb{T}_{ijk}^p + \mathbb{T}_{ijk}^R = \frac{\nu_t}{\sigma_k} \frac{\partial k}{\partial x_j} \quad (\text{A.37})$$

Substituting Eq. (A.36 & A.37) into Eq. (A.34a), the governing equation for k changes to

$$\frac{\partial k}{\partial t} + \frac{\partial \bar{u}_j k}{\partial x_j} = P_{ij} - \varepsilon_{ij} + \frac{\partial}{\partial x_j} \left\{ \left(\nu + \frac{\nu_t}{\sigma_k} \right) \frac{\partial k}{\partial x_j} \right\} \quad (\text{A.38})$$

The dissipation term ε in the above equation represents the energy-flow rate from the large turbulence to small turbulence in the cascade, which means it is mainly determined in the large-scale motion. The governing equation of ε can also be derived for numerical simulation. However, the derivation process is based on the dissipative range which represents the smallest scale of motion in the flow field. As a result, it is recommended to check this equation from the empirical view that it follows a similar transport form of Eq. (A.38).

$$\frac{\partial \varepsilon}{\partial t} + \frac{\partial \bar{u}_j \varepsilon}{\partial x_j} = (C_{\varepsilon 1} P_{ij} - C_{\varepsilon 2} \varepsilon) \frac{\varepsilon}{k} + \frac{\partial}{\partial x_j} \left\{ \left(\nu + \frac{\nu_t}{\sigma_\varepsilon} \right) \frac{\partial \varepsilon}{\partial x_j} \right\} \quad (\text{A.39})$$

In fewer words, the RANS model uses Eq. (A.35, A.38, A.39) to calculate the turbulent viscosity coefficient, and substitute it into Eq. (A.33) to model the Reynolds stress, which satisfies the closure requirement of Eq. (A.30 & A.31). In these equations, in total five dimensionless coefficients exist to tune the performance of RANS. Their standard values are set by Launder and Sharma (1974) by data fitting for different turbulent flows.

$$C_\mu = 0.09, C_{\varepsilon 1} = 1.44, C_{\varepsilon 2} = 1.92, \sigma_k = 1.0, \sigma_\varepsilon = 1.3 \quad (\text{A.40})$$

A.4.3 Dispersion equation for passive scalar

The Reynolds average can also be operated to passive scalar transport equation, Eq. (A.3) changes to

$$\frac{\partial \bar{C}}{\partial t} + \frac{\partial \bar{u}_j \bar{C}}{\partial x_j} = \frac{\partial}{\partial x_j} \left(D_m \frac{\partial \bar{C}}{\partial x_j} \right) + S - \frac{\partial \overline{u_j' C'}}{\partial x_j} \quad (\text{A.41})$$

The extra term $\overline{u_j' C'}$ is the Reynolds average of turbulent diffusion term which needs modeling. The most common way in RANS is the gradient diffusion hypothesis, which assumes that the turbulent diffusion is proportional to the gradient of the averaged concentration field.

$$\overline{u_j' C'} = -D_t \frac{\partial \bar{C}}{\partial x_j} = -\frac{\nu_t}{Sc_t} \frac{\partial \bar{C}}{\partial x_j} \quad (\text{A.42})$$

where D_t is the turbulent diffusivity, and Sc_t is the turbulent Schmidt number.

A.5 Comparisons of two turbulence model

As two of the most popular approaches for turbulence modeling, RANS and LES have distinguishing features.

RANS focuses on the mean flow due to the Reynolds averaging processing. The unsteady effects of turbulent flow are also removed out from the averaging and approximated by the mean physical quantities. Therefore, the computational cost of RANS is mild and makes RANS widely used in engineering applications for the past decades. Despite that, its disadvantage is also obvious that the unsteady turbulence is unavailable. More importantly, since the effects of turbulence on the mean flow are also implicitly modeled, the accuracy of mean flow is not good enough in the flow fields like strong separation from the corner of the building.

In contrast, LES focuses on the relatively large-scaled turbulent structures with filtering techniques. It only models the SGS behavior of turbulence flow. As a result, when the mesh is fine enough, most dominant turbulence structures can be accurately and explicitly captured, which brings a more reliable simulation of the flow field than RANS. The corresponding cost is the high requirement on the computational resources. In order to promise the performance of LES, the mesh should be carefully designed with fine resolution, which needs additional resources apart from that for unsteady simulation caused by time marching. Meanwhile, due to a large amount of mesh, the storage pressure of LES results is also much heavier than RANS.

Nowadays, with the development of high performance computers and semiconductor products, the calculation burden of LES becomes lighter and the importance of accuracy becomes more desired. This is the reason that more and more engineering applications start to switch from RANS to LES, and the application of LES in the adjoint equation simulation is also part of the main work of this dissertation.

Appendix B. Publications related to dissertation

B.1 Journal publications

(with peer review)

1.1) Hongyuan Jia, Hideki Kikumoto, “Construction of urban turbulent flow database with wavelet-based compression: A study with large-eddy simulation of flow and dispersion in block-arrayed building group model”, *Journal of Wind Engineering and Industrial Aerodynamics*, vol. 208, 104433, 2021.

1.2) Hongyuan Jia, Hideki Kikumoto, “Line source estimation of environmental pollutants using super-Gaussian geometry model and Bayesian inference”, *Environmental Research*, vol. 194, 110706, 2021.

1.3) Hongyuan Jia, Hideki Kikumoto, “Source term estimation in complex urban environments based on Bayesian inference and unsteady adjoint equations simulated via large eddy simulation”, *Building and Environment*, vol. 193, 107669, 2021.

1.4) Hongyuan Jia, Hideki Kikumoto, “Sensor configuration optimization based on the entropy of adjoint concentration distribution for stochastic source term estimation in built environment”, *Sustainable Cities and Society*, accepted.

(without peer review)

1.5) 賈 鴻源, 菊本 英紀, 「スーパーガウス関数とベイズ推定を用いた環境汚染物質線形発生源の同定」, 生産研究, 2020, 72巻, 1号, pp.49-55.

1.6) 賈 鴻源, 菊本 英紀, 「LESによる随伴濃度解析を用いた市環境汚染物質発生源のベイズ推定」, 生産研究, 2021, 73巻, 1号, pp.71-76.

B.2 Conference publications

(with peer review)

2.1) Hongyuan Jia, Hideki Kikumoto, “Line source estimation of environmental pollutants using Bayesian inference coupled with super-Gaussian geometry model”, 15th RoomVent Conference, Torino, Italy, (2021).

(without peer review)

2.2) 賈 鴻源, 菊本 英紀: 「スーパーガウス幾何モデルを用いた環境汚染物質線形発生源のベイズ推定」, 第 33 回数値流体力学シンポジウム, 北海道, 2019 年 11 月.

2.3) 賈 鴻源, 菊本 英紀: 「CFD を用いた市街地気流解析におけるデータ圧縮に関する研究 (その 1) ウェーブレット変換に基づいた圧縮手法の概要」, 日本建築学会大会, 千葉, 2020 年 9 月.

2.4) 胡 書媛, 賈 鴻源, 菊本 英紀: 「CFD を用いた市街地気流解析におけるデータ圧縮に関する研究 (その 2) 圧縮した LES 流れ場データを用いた汚染物質拡散の再解析」, 日本建築学会大会, 千葉, 2020 年 9 月.

2.5) 賈 鴻源, 菊本 英紀: 「ウェーブレット変換に基づく圧縮手法を用いた市街地気流の LES データベースの構築」, 日本流体力学会年会, 山口, 2020 年 9 月.

2.6) 賈 鴻源, 菊本 英紀: 「確率的発生源同定におけるセンサー配置最適化手法に関する研究 (その 1) 随伴濃度分布のエントロピーに基づく最適化アルゴリズム」, 日本建築学会大会, 名古屋, 2021 年 9 月.

2.7) 李 栄茂, 賈 鴻源, 菊本 英紀: 「確率的発生源同定におけるセンサー配置最適化手法に関する研究 (その 2) 室内モデル空間における随伴濃度分布とセンサー配置の最適化」, 日本建築学会大会, 名古屋, 2021 年 9 月.

Acknowledgment

This dissertation is a summary of my research during three years as a doctoral candidate in the graduate school of engineering, the University of Tokyo. In this unforgettable period, I received a great deal of support and help from people around me, to whom I would like to express my sincere gratitude.

First and foremost, I would like to thank my supervisor Associate Professor Hideki Kikumoto. I feel extremely lucky to finish my doctoral degree under your supervision. I missed our long-lasting chatting in the office where I learned about how to start my career as a researcher. You gave me a lot of inspiring suggestions when I encountered difficulties. Any research idea was warmly welcomed by the professor even if it sounds immature. The joy of interest-driven research and insightful discussions provided in the Kikumoto Laboratory are the energetic fuels to this dissertation.

I would also like to express my deepest appreciation to my deputy supervisor Professor Ryoza Ooka. You gave me valuable advice on the general direction of this research. The late lighting of your office always reminds me of the importance of hard-working to success.

I must give my sincere gratitude to my doctoral advising committee member, Professor Shinichi Sakamoto, Associate Professor Yosuke Hasegawa, and Associate Professor Masayuki Mae. Their careful reviewing and helpful suggestions improve this dissertation with no doubt. Encouragement from Prof. Sakamoto gave me confidence to continue my research. Prof. Masayuki guided me from the structure of the dissertation. Prof. Hasegawa made this research more rigorous in mathematics and expressions.

I am also grateful to Mrs. Iizumi and Mrs. Honma, the secretaries of the laboratory. They handled the paperwork of my administrative procedures, which makes me focus on my research at ease.

Special thanks go to Associate Professor Mengtao Han in Huazhong University of Science and Technology. As an OB of the laboratory, he taught me the knowledge about Linux servers and OpenFOAM, which helped me smoothly start the research.

I would like to thank Mr. Chao Lin, who is not only a colleague but also a good friend in daily life. He gave me continuous support starting from my preparation for the entrance

examination. I always enjoyed chatting with him about research or other trivial things around.

With a special mention to Dr. Wonjun Choi, Dr. Wonseok Oh, Dr. Doyun Lee, Dr. Qi Zhou, Dr. Bingchao Zhang, Dr. Shan Gao, Dr. Mingzhe Liu, Dr. Huaiyu Zhong, Dr. Xiaochen Liu, Yuki Matsuda, Ke Wen, Christopher O'Malley, Chenghao Wei, Wenchao Wang, Rongmao Li, Hong Hu, Hiroyuki Ichikawa, Yunchen Bu, Chaoyi Hu, Xiang Wang, Shuyuan Hu, Ken Takahashi, Mayuka Nakai, Sorana Ozaki, Toshiki Nakamura, Eiki Sakamoto, Tomoya Oosaki, Ran Chen, Keisuke Ogasawara, Dun Zhu, Qianwen Guo in general. It was wonderful to have the opportunity to work with such a big group, which exposes me to a wide range of research topics.

I would like to acknowledge my best friend Dr. Zhipeng Zhang for his encouragement to my career pursuit and for accompanying me in the academic journey for over 13 years.

I must express my profound gratitude to my parents for providing me with financial support and understanding about my life choice. This dissertation cannot be finished without their love and care.

Finally, my appreciation goes to my wife, who studies abroad with me, always believes me to be able to accomplish this dissertation and comfort me at depressed moments. Over two of the past three years, I had to work from home because of COVID19, but your love made it the happiest time for both research and life. This dissertation will not be possible without you.