**Modelling Functional Spatial Structure in Mega-city Region Based on Human Mobility Pattern**

Zhang, Bowen

*Awarding institution:*
King's College London

# Modelling Functional Spatial Structure in Mega-city Region Based on Human Mobility Pattern

**Bowen Zhang**

Student ID: 19029097

A dissertation submitted in partial fulfilment of the requirements for the degree of

**Doctor of Philosophy**

of

**Geography**

King's College London

Supervised by:  Dr Zahratu Shabrina,

Dr Chen Zhong and Dr James Millington

26th August 2024

# Declaration

I hereby declare that except for cited instances explicitly acknowledging the contributions of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this or any other university. This dissertation represents my individual effort and does not include any material produced through collaborative efforts except as delineated in the text and acknowledgements section.

<div align="right">

Bowen Zhang

August 2024

</div>

# Acknowledgements

Embarking on my PhD journey on the eve of the pandemic, I navigated through an era marked by unprecedented challenges and uncertainties. Completing this doctorate represents a long journey that could not have been achieved without those wonderful people's selfless help and support.

First of all, I would like to express my profound gratitude to my supervisors, Dr Zahratu Shabrina, Dr Chen Zhong, and Dr James Millington. I would appreciate Zara for her patience and time spent discussing my work and thesis. Her encouragement and insistence on high standards propelled me to refine my work to new heights. I would especially thank Chen for her consistent support throughout the whole PhD journey. Her expertise and insightful guidance have been instrumental in shaping my research trajectory. I want to thank James for his kind support during the special time when switching supervisors due to unexpected circumstances.

I want to thank my colleagues and friends at King's Geography and CASA UCL for spending an unforgettable PhD life with me. Sharing this path with you has been both a privilege and a source of endless inspiration. I have greatly enjoyed and appreciated the support and friendship.

A special tribute is owed to my family, whose unwavering support and love have been the bedrock of my resilience throughout this endeavour. To my parents, thank you for believing in me and encouraging me to pursue my dreams. A heartfelt acknowledgement to Mingzhu Cai, whose understanding, love, and unwavering support have been my stronghold.

This thesis stands as a milestone in my academic journey, and I am grateful for the opportunity to grow and learn under the support of so many. Thank you!

# Abstract

Mega-city regions have emerged from urban growth and improved inter-city connectivity. Due to the complexity of its spatial arrangement, quantitatively describing and predicting the functional spatial structure within mega-city regions has become a new challenge. To address this challenge, human mobility has become the hot spot of research, as it enables the exchange of ideas, goods, services, and cultural interactions that shape the dynamics of urban spaces. Human mobility is integral to the functioning of society. It could explain the relationship between micro-level individual behaviour and macro-level urban phenomenon. Therefore, this research proposes to develop an analytical framework for modelling the functional spatial structure in mega-city regions through the lens of human mobility, predicting the dynamic shift in the urban spatial structure. To achieve the main research aim of predicting the urban spatial structure, this research set a series of research objectives as follows: (1) To identify urban functional zones within mega-city regions by examining travel behaviour and differentiating intra-city from inter-city trips; (2) To develop a novel spatial interaction model that enhances travel flow predictions by incorporating residents' socio-economic characteristics; and (3) To predict the impact of urban interventions and policies on travel patterns and the mega-city region's functional spatial structure through localised changes.

This study proposes several novel algorithms based on spatial-interaction models to achieve its research objectives and then tests them with case studies. The study first designed a regionalisation algorithm for delineating urban functional zones, utilising cell phone signalling data in the Great Bay Area in China. Secondly, this research proposes a novel variation of the

spatial interaction model, combined with the k-means clustering algorithm, to predict the travel flows of residents using census data in the Greater London Area in the United Kingdom. Furthermore, this research integrated the tools to simulate how urban interventions and policies affect the functional spatial structure in the Great Bay Area in China from the perspective of human mobility patterns.

The primary research outcome of this study suggests that the distance decay in the spatial interaction model exhibits significant spatial heterogeneity, and this parameter could be used to represent the functional urban spatial structure. This distance decay parameter could also be associated with various factors, including spatial arrangement and non-spatial factors, such as socio-economic factors. By predicting the local variation of the distance decay with socio-economic characteristics and travel flows, we can forecast the dynamics of the urban spatial structure in the mega-city region for future scenarios. This simulation model could help governments and urban planners make informed decisions by forecasting the impacts of urban interventions on the spatial structure of mega-city regions.

Furthermore, this thesis advances the discussion on long-standing issues in spatial interaction models using human mobility big data research, such as localisation, calibration methods, and spatial heterogeneity, which contribute to solving these long-standing issues through novel approaches.

# List of publications based on the thesis

Parts of the materials included in this thesis have been published or are under consideration for publication in the form of journal articles or book chapters as follows:

**Journal Papers**

Zhang, B., Zhong, C., Gao, Q., Shabrina, Z., & Tu, W. (2022). Delineating urban functional zones using mobile phone data: A case study of cross-boundary integration in Shenzhen-Dongguan-Huizhou area. *Computers, Environment and Urban Systems*, *98*, 101872.

Zhang, B., Zhong, C., Gao, Q., & Shabrina, Z., Exploring the Associations of Socioeconomic Characteristics and Distance Decay effects in Spatial Interaction (Preprint available at SSRN 4733461, Under review by *Sustainable Cities and Society*)

**Selected Conference Proceedings**

Zhang, B., Zhong, C., Gao, Q., & Shabrina, Z., Exploring the Distance-decay Effect in Commuting Behaviour at the Local-level with a Localised Spatial Interaction Model. *GISRUK 2023*, Glasgow, United Kingdom

Zhang, B., Zhong, C., Gao, Q., Shabrina, Z., & Tu, W., Delineating urban functional zones using mobile phone data: A case study of cross-boundary integration in Shenzhen-Dongguan-Huizhou area. *The 2021 European Colloquium on Theoretical and Quantitative Geography (ECTQG)*, Manchester, United Kingdom

# Table of Content

# List of Figures

# List of Tables

# 1 Introduction

## 1.1 Research Background

During (the second half of) the 20<sup>th</sup> century, a significant change in urban spatial structure occurred. The typical monocentric organisation of cities transformed into the decentralised polycentric urban land-use form in many areas (Anas et al., 1998; Smith, 2009). The most recent phase is that city regions and super mega-city regions have been formed as a consequence of urban growth and a vast improvement in inter-cities connectivity (Hall & Pain, 2006). The mega-city regions are characterised by a network of closely linked urban areas centred around one or more core cities, which gained vast attention from the public and scholars because of the massive population and its economic agglomeration effects (Scott, 2019). With the spatial structure of an urban environment/space becoming more complex, our knowledge about the mechanism of how cities grow. However, even after a century of work, the understanding of how cities evolve is still insufficient (Batty, 2008). The concept of urban functional zones (UFZs) varies based on the research objectives and the data utilised, including different factors such as land-use patterns, human activities, and regional planning principles (Chen & Yeh, 2022; Niu & Jin, 2020; Zhong et al., 2014). In the context of mega-city regions, daily activities tied to urban functions—like work, living, and leisure—spill over established administrative borders, happening across various cities. This expansion gives rise to the new concept of urban functional zones (UFZs), where these activities are no longer contained within

a single city's limits but spread across a larger metropolitan region. Therefore, in this PhD thesis, the UFZs primarily refer to the geographical extent covered by each city's functional area within a mega-city region. In urban analytics, it is a challenge to quantitatively describe and predict the complex urban spatial structure within mega-city regions. Researchers have addressed this challenge and developed urban models as simplified representations of reality.

A model is an approximation of truth that mediates between theory and the physical world, allowing scientific experimentation and a means of analytical testing (Morgan & Morrison, 1999). Urban models originated in the late 1950s as simulations designed to represent the patterns of land use and the transport flows within cities, then further developed as the Land Use-Transport Interaction (LUTI) model during the 1960s in the United States. Computational models have long been used to support planners' decisions to eliminate complexity and codify a straightforward and concise understanding of many aspects of urban structure and transport. In the last decades, various types of urban models have been created to explain the mechanism of urban space forms and changes and how urban space will influence residents' daily activities. The famous examples include the gravity model, fractal models, cellular automata, and Agent-based models (Batty & Longley, 1994; Matthews et al., 2007; Santé et al., 2010; Wilson, 1971). However, traditional urban models primarily utilise conventional data sources, emphasising location and infrastructure while often neglecting individual social characteristics and preferences.

Additionally, computational limitations have historically restricted the scale of urban systems these models can accurately represent. Recently, the evolution of computing power and data collection/storage techniques enabled the possibility of using spatial interaction models to build larger urban models with granular spatial resolutions, and this trend is known as "the Renaissance of large-scale modelling" (Batty & Milton, 2021). More research suggests that

cities are complex systems that grow from the bottom up (Batty, 1976a), and there is optimism that big-data techniques will address the lack of individual dimensions in research (Geurs & Van Wee, 2004).

The emerging data sources gathered from sources like smartphones, GPS devices, and smart cards provide a high-resolution image of how people move across and interact with urban spaces (González et al., 2008), and the amount of these urban mobility data has grown exponentially in recent years (Zhong et al., 2016). Human mobility is integral to the functioning of society, as it enables the exchange of ideas, goods, services, and cultural interactions that shape the dynamics of urban spaces (Yuan et al., 2012). Especially with the concept of big data, human mobility patterns have become the hot spot for explaining the relationship between micro-level individual behaviour and macro-level urban phenomenon (Anejionu et al., 2019; Huang et al., 2019; Shelton et al., 2015). This granular understanding enables urban geographers and analysts to unravel daily movement patterns, pinpoint critical areas of activity, and model the complex urban spatial structure. The unprecedented speed of urban growth and the change in human mobility patterns in a data-rich era would provide a valuable perspective to observe how accessibility and geospatial planning policy reshape the urban spatial structure.

## 1.2  Research Question and Objectives.

The thesis explores methods for modelling the mega-city region's functional spatial structure based on human mobility patterns. The doctoral research attempts to answer the question, "**How can we predict the changes of functional urban spatial structure in mega-city regions from a human mobility perspective**?" The hypothesis is that residents' travel behaviour could reflect the functional urban spatial structure. In addition, we also hypothesise

that local socioeconomic demographics and travel behaviour correlate, thus enabling improvements in travel flow prediction. Therefore, these research questions can be answered through three specific research objectives as follows:

1. To detect the urban functional zones of emerging mega-city regions reflected by the human mobility pattern and the difference between intra-city and inter-city trips.

2. To establish a novel spatial interaction model to predict travel flows more accurately by highlighting the socioeconomic characteristics of local residents.

3. To predict how specific urban interventions and policy can influence human mobility patterns via changing localised characteristics, further affecting the functional spatial structure within the mega-city region.

## 1.3 Significance of the study

This doctoral research aims to provide an in-depth analytic framework for understanding the mechanism of functional spatial structure reshaping in mega-city regions through the lens of human mobility. In his book *Overview of Land Use Transport Models*, Wegener, M. (2004) defined the urban spatial processes and temporal scales. Fast processes include information flows (which may only take seconds via the internet) and everyday urban transport loops. Regarding land use and urban function, the mechanisms shaping communities' physical configuration are typically medium-term temporal relations (10 years or less), mainly involving industries, residents' location decisions, and urban growth cycles. Meanwhile, these medium-term processes are correlated with faster complex processes like travel patterns and slower urban functions like structural economic/major demographic change or techno-economic paradigm shift (Wegener, 2004). One of the main challenges of urban analytics is building a

framework that connects different spatiotemporal scales. Thus, this doctoral research would use figures of fast processes, such as human mobility patterns, to assess mid-term and long-term processes, such as urban spatial structure, contributing to a better understanding of how to incorporate the interrelationships between urban spatial structure at various levels. Methodologically, this study aims to make several contributions:

Firstly, the research contributes to improving the current spatial interaction models. Most current spatial interaction models and other flow-predicting models at an aggregated level assume that the interior space of the modelling region is spatially isogenous, meaning that the distribution of trips only obeys a general law associated with distance between locations (De Vries et al., 2009; Fotheringham & O'Kelly, 1989; Simini et al., 2021). Researchers found that a global approach to spatial analysis may not be suitable for the local area within the sub-case study area due to spatial heterogeneity (Fotheringham & Sachdeva, 2022). Thus, adding local socioeconomic characteristics to improve the accuracy of spatial interaction models is a long-standing research topic. Thus, rethinking and building a more localised model is not only the first stage needed for this study but also a contribution to the methodology in urban analytics. This research seeks to fill a knowledge gap by answering how variations in urban spatial structure and social groups could be reflected in the spatial interaction of travel behaviours. This helps to establish a link between socioeconomic characteristics, urban spatial configuration, and spatial interaction.

Secondly, this research aims to provide a new definition of the functional urban structure of cities/city regions by human mobility patterns. An urban system often does not coincide with the administrative boundaries, which may cause distortion and lead to planning failure (Calafati & Veneri, 2013). Due to the vagueness of the word 'urban' and the uncertainty of which feature of urbanity produces the efficiency premium, there is a lack of theoretical formulation on

agglomeration economies' urban boundary (Bretagnolle et al., 2002). Compared to traditional monocentric cities, boundaries in the form of city-region are even more challenging to define. To discuss the urban spatial structures of the mega-city regions, it is necessary to determine the boundaries that fit the specific urban context. The urban space is formed by residential and industrial activities (Lynch, 1960). Modern big data techniques enable the possibility to observe how individuals' daily activities agglomerate to create urban spaces, indicating the in-fact boundary of cities/city-regions. The view from human mobility would help measure the functional urban boundaries and explain their geographical meaning.

Lastly, this research also establishes a simulation-based approach to understanding how specific policies and interventions affect the human mobility pattern and ultimately influence the spatial structure of cities in relatively medium- or long-term periods. Urban systems are becoming even broader and more complex with the development of economies and social and transport structures. Computational models have long been used to support planning decisions to eliminate complexity and codify a straightforward and concise understanding of many aspects of urban structure and transport. Previous LUTI models determine the land-use patterns and then predict the transport interactions. Besides, it mostly focuses on separated urban systems and rarely considers socioeconomic characteristics. Therefore, this research would establish a model for a mega city region that predicts the spatial structure change based on the human mobility pattern and considering socioeconomic characteristics. This simulation framework will support decision-making by the government and planners in predicting future scenarios with/without urban interventions.

## 1.4 Methodology Summary

In this thesis, we would like to advocate an urban analytic methodology framework for modelling functional spatial structure about human mobility patterns in the context of mega-city regions. The methodological framework is shown in Figure 1.1, the specific research and modelling method is introduced below:



*Figure 1.1 Flowchart of the methodology framework*

## 1.4.1  Spatial interaction models

The spatial interaction model is a framework for predicting and understanding the flow of goods, services, people, and information between different locations based on their distance, economic size, and other intervening factors. The spatial interaction model (SIM) addresses research inquiries concerning how specific urban interventions can impact urban spatial structures by influencing human mobility. The traditional, or global, spatial interaction model assumes that all trips adhere to a general law, typically characterised by a negative power or exponential function with uniform distance decay parameters. However, our assumptions suggest that multiple human mobility patterns may coexist within the same region. Therefore, we propose a novel variant of the spatial interaction model for predicting regional human mobility flows. This research refined existing spatial interaction models to better elucidate local variations in travel flow distributions. The fitted localisation parameters (especially for the distance decay) and updated predicted movement flow could serve to measure social indicators

and urban spatial structures through the analytical framework proposed. To emphasise the local socioeconomic characteristics in the SIM model, we utilise the k-means clustering algorithm to partition our modelling areas into k-clusters, facilitating the prediction of local distance decay parameters while limiting the computing complexity. Additionally, the goodness of fit in spatial interaction models could be employed as an indicator to assess the validity of delineating the boundaries of sub-models.

## 1.4.2  Regionalisation algorithm

The regionalisation algorithm aims to delineate N larger regions from the aggregation of M smaller regions, where M exceeds N (Duque et al., 2007; Shortt, 2009). Research Objective 2 seeks to identify urban functional zones within the mega-city region, aligning with the regionalisation algorithm's fundamental concept. Consequently, the regionalisation algorithm serves as a core research method in this study to delineate the urban functional zones from human mobility. Specifically, the study examines the spatial interaction volumes between areas with two urban cores, which typically follow a distance decay law. If the spatial interaction suggests a continuous urban space between these cores, it implies integration within the urban space. Conversely, significant discrepancies between predicted and observed commuting flows between cores may indicate the presence of an invisible barrier, suggesting two separate urban spaces in spatial structure terms. We have designed and applied a novel regionalisation algorithm based on the spatial interaction model (SIM) to operationalise this approach. This algorithm aims to identify the optimal partitioning scheme that best fits the observed variation in SIM. The iteration-based algorithm iteratively refines the partitioning scheme until the best partition, characterised by the highest goodness of fit, is achieved.

### 1.4.3  Urban simulation model

Urban simulation models are powerful tools in this effort, providing insights into the complex interactions between land use patterns, transportation networks, and population trends (Harris & Batty, 1993). Urban spaces are changing dynamically, with population growth, economic development, and technological advances reshaping the structure of mega-city regions. Quantitative analysis of the physical urban environment and the study of the relationship between form and function are the future trends of urban science (Wu et al., 2024). In this research, we built an urban simulation model integrated with the other methods and models mentioned above to predict how specific urban interventions and policies can influence the functional spatial structure within the mega-city region as human mobility patterns change.

## 1.5  Case Study Area

In the global north context, the formation of mega-city regions usually followed the pathway of cities - city regions - mega-city regions (Scott et al., 2001). It finally grew as a super mega-city region with a huge population and complex spatial structure. However, recent research pointed out that the formation of super mega-city regions in some emerging markets countries, such as China and Mexico, didn't follow this typical pathway (Scott, 2019; Yeh & Chen, 2019). In global northern countries, especially in the global mega city-region like London, data availability and reliability are significant advantages compared to developing countries because of their publicly available advanced facilities and mature data platforms. However, the formation of urban spatial structures in the majority of mega-cities was completed before the 1980s in the global north countries (Brenner, 2002; Scott, 2019), which means we cannot obtain

the benefits of big data techniques to understand the spatial transformation in the mega-city regions. Fortunately, global southern countries like China have been experiencing urban growth and spatial structure change in recent decades along with their economic growth (Li, 2020; Wang et al., 2016). Meanwhile, the inequality in mobilities and accessible opportunities is still remarkable across regions and social groups in global southern countries. Thus, this research will take the Great London Area in the UK and the Great Bay Area in China as two case study areas, exploring the interplay between socioeconomic status, human mobility patterns and urban spatial structure.

## 1.5.1 Shenzhen-Dongguan-Huizhou (SDH) metropolitan area

This doctoral research takes the Shenzhen-Dongguan-Huizhou (SDH) metropolitan area, which is one of the sub-regions in the Pearl River Delta Great Bay Area (GBA), as the first case study area. The Greater Bay Area comprises the two Special Administrative Regions of Hong Kong and Macao and the nine Pearl River Delta (PRD) cities of Guangzhou, Shenzhen, Zhuhai, Foshan, Huizhou, Dongguan, Zhongshan, Jiangmen and Zhaoqing in Guangdong Province. The GBA area has a total area of 56,000 square kilometres and a total population of 86 million at the end of 2022, according to the yearbook of Guangdong Province (2022). Shenzhen-Dongguan-Huizhou area as a geographical concept formally appeared in the 2004 PRD Urban Cluster Coordinated Development Plan. It consists of one vice-provincial-level municipality, Shenzhen (subordinate to the central government and the Guangdong Provincial Government) and two prefecture-level municipalities, Dongguan and Huizhou (under the Guangdong Provincial Government). SDH area covers a total area of 15,800 square kilometres, with a resident population of 34.15 million and a total GDP of RMB 4.9 trillion in 2022. This area has been experiencing rapid urban growth and change in urban spatial structure since the 1980s and has become one of China's most open and economically vibrant regions.

*Figure 1.2 Great Bay Area (GBA)and Shenzhen-Dongguan-Huizhou (SDH) areas, China*

Since 1978, reform and opening-up policies have led to exponential population growth and industrial prosperity, especially concentrated in super mega-city regions. Because of their crucial role in China's urbanisation and economic growth, mega-city regions have recently been put at the forefront of policy. The Chinese government has issued a series of policies since 2004 to encourage cities within one city-region growth and integration as one city, which further promoted the urban spatial structure change and urban space integrations (Li et al., 2015; Wu, 2016). The national and international spotlight has increasingly been on this region, particularly following the proposition of the Guangdong-Hong Kong-Macao Greater Bay Area in 2015.

## 1.5.2  Greater London Area (GLA)

Another case study area is the Greater London Area, which is substantial both in size and population, making it a significant urban zone in the United Kingdom and a vital player in the global economy. GLA covers approximately 1,572 square kilometres (about 607 square miles). This area encompasses the City of London, the historic and financial heart of the metropolis, and 32 boroughs. Greater London is home to nearly 9 million people, making it the most populous municipality in the United Kingdom (Census, 2021). It is the engine of the UK's economy, contributing 23% of the country's GDP (ONS, 2021). The city's economic activity benefits from its status as a global transport hub, served by extensive underground and rail networks, and major airports like Heathrow and Gatwick. As a global mega-city region, it presents a rich cultural, economic, and environmental amalgamation, rendering it a prime subject for a broad spectrum of geographical studies. London could provide the best data availability in multi-dimensions including geographical data, flow data, and social characteristics data. In addition, numerous previous literatures taking London as a case study area could be referenced to help us understand the urban phenomenon, and better design our modelling method.

*Figure 1.3 The Great London Area (GLA), United Kingdom*

## 1.6  Structure of this report

The thesis is divided into seven chapters:

**Chapter 1** introduces the research background, focusing on how to model urban spatial structure in mega-city regions, and outlines the research questions and objectives aimed at addressing the identified gaps in the literature. It then delineates the rationale for undertaking this research by highlighting the significance of this study. Finally, the chapter presents an overview of the report's structure, guiding readers through the subsequent text for clear navigation.

**Chapter 2** is a literature review chapter. It first introduces previous research about the travel flow prediction method, particularly focusing on spatial interaction models. Then the spatial

transformation of the city region and the urban boundaries, as the key element of the spatial structure, were examined by reviewing the related research. In later sections, the current large-scale urban simulation model and its trend of the 'Renaissance' are reviewed. Lastly, the previous research about mobility big data utilisation in urban analytics will be summarised. The specific research gaps are identified in this chapter, contributing to establishing the methodologies of this study.

**Chapter 3** introduces a regionalisation algorithm for delineating urban functional zones using human mobility data, aiming to achieve the objectives of detecting the functional spatial organisation. By analysing the mobile phone data in the Great Bay Area in China as a case study, this chapter provides insights into cross-boundary city integration and related policy implications. This chapter is a "*thesis incorporating publications*" chapter, as the content has been published in the Journal of *Computers, Environment and Urban Systems* (Zhang et al., 2022).

**Chapter 4** proposes a novel travel prediction tool with consideration of localised characteristics. In this chapter, a two-step spatial interaction model emphasises this variance in the local distance-decay effect by reflecting the socioeconomic characteristics of residents. As a case study in the Great London Area in the UK, this model predicts commuting behaviours using census 2021 data.

**Chapter 5** establishes a simulation model to predict how the specific urban intervention could influence the urban spatial structure. Based on the two-step spatial interaction model in Chapter 4, specific policy assumptions (e.g. population growth & migration, transport facilities development, and changes in socioeconomic characteristics) are tested in this simulation model

to predict travel behaviour change. The possible change in functional spatial organisation will also be examined in relation to the regionalisation algorithm proposed in Chapter 3.

**Chapter 6** is a discussion chapter based on the previous chapters. The chapter first discusses the relationship between distance decay and urban spatial structure in Mega-city regions through analytical results. This thesis concludes that the distance decay in the spatial interaction model exhibits significant spatial heterogeneity, and this parameter could be used to represent the functional urban spatial structure. Furthermore, this chapter discusses some long-standing issues in spatial interaction models using human mobility big data research, such as localisation, calibration method, and spatial heterogeneity. It illustrates how this doctoral research could contribute to solving these long-standing issues.

**Chapter 7** summarises the significant contributions of this doctoral study, including understanding the transformation of urban spatial structures through the lens of human mobility, methodological contribution, and support for urban planning and policymaking via an urban simulation model. Furthermore, this chapter suggests directions for future research in related fields of urban analytics.

# 2 Literature review

The interplay between human mobility and urban spatial structures has emerged as a key area of research within geography and urban studies, particularly in emphasising the dynamic of urban space in the 21st century. This literature review chapter aims to introduce the historical development of these concepts and explore their relevance and application in the modern context.

This chapter first reviews the spatial interaction models by comprehensively examining the evolution, theories, and methodologies underpinning this dynamic field. The literature review delves into the challenges and opportunities presented by applying spatial interaction models within urban systems, highlighting the need for localised models to accommodate the granular spatial resolutions in the big data era. Then, this discussion around the concept of urban spatial structure examines cities' morphological and functional aspects and the transition from monocentric to polycentric forms as urban areas expand and evolve. The spatial transformation of the city region and the urban boundaries, as the key element of the spatial structure, will be examined by reviewing the related research. Lastly, the previous research about mobility big data utilisation in urban analytics is summarised.

By identifying existing gaps in the literature, this review chapter proposes directions for future research, focusing on integrating new data sources and localisation of the models in studying human mobility and urban spatial structures.

## 2.1 Human mobility and spatial interaction model

The foundational concept of human mobility patterns was first introduced during the 19th century to elucidate the frequency of travel between neighbouring cities, considering factors such as population size and distance. Subsequent decades witnessed the development of migration laws and Zipf's formulation (Zipf, 1946), which laid the groundwork for the emergence of the widely recognised gravity law, serving as the cornerstone for contemporary analyses of human mobility patterns. Quantitative investigations into human mobility commenced in the United States metropolitan areas, driven by the conceptualisation of 'geography as spatial interaction' in the 1950s, as articulated by Haynes and Fotheringham (1985).

As data collection efforts improved, particularly with the advent of Information and Communication Technology (ICT) data, the concept of time geography gained prominence. This approach considers both temporal and spatial constraints and has been extensively utilised to measure, comprehend, and forecast spatiotemporal trajectories at the individual level. Researchers have examined various travel scales, ranging from daily commuting to international travel, to unravel their associations with practical real-world applications, such as traffic flow forecasting (Lopane et al., 2023), urban planning (Bokányi et al., 2019), and epidemic modelling (Spooner et al., 2021; Zhou et al., 2020), risk management (Song et al., 2014). In the realm of urban studies, significant indicators such as accessibility and employment density can exert considerable influence on human mobility patterns within cities (Bocarejo S & Oviedo H, 2012; Preston & Rajé, 2007). Correspondingly, the human mobility model, along with its outcomes and parameters, can serve as valuable indicators for unveiling the distinctive characteristics of cities (Zhong et al., 2017).

Human mobility models can be employed to replicate both individual mobility patterns and population movement patterns (Alessandretti et al., 2020). In both scenarios, it is imperative to account for the distinctive geographical and temporal dimensions inherent to the mobility process, encompassing distances ranging from hundreds of meters to thousands of kilometres and time spans varying from hours to years (Wegener, 2004). At the individual level, a range of models, such as Brownian motion and Lévy flight, have been developed based on the principles of random walk theory. These models predict the likelihood of an individual's location or travel range (Barbosa et al., 2018; Rhee et al., 2011). Conversely, at the population level, models are designed to represent collective mobility behaviours and aim to replicate Origin-Destination (OD) matrices by estimating the average number of individuals travelling between any two geographical zones over a given unit of time (Willumsen, 2001). The spatial interaction model is the prevailing model for forecasting travel flows in this context.

## 2.1.1 Classic gravity model theory and its development

Spatial interaction studies have employed a variety of methods, with the gravity model emerging as the most frequently utilised approach (Haynes & Fotheringham, 1985). In 1946, George K. Zipf introduced an equation to forecast movement flows, positing that the volume of movement flow is directly proportional to the product of the population sizes of any two communities, denoted as $P_i \times P_j$, and inversely proportional to the shortest transport distance, represented as $d_{ij}$ (Zipf, 1946).

$$T_{ij} \propto \frac{P_i \times P_j}{d_{ij}} \tag{1}$$

While the gravity model has gained popularity due to its formal simplicity and successful application in modelling empirical flows and movements, its theoretical underpinnings have

39

remained a subject of debate. In 1970, a significant breakthrough occurred when Wilson introduced the entropy-maximizing methodology into the gravity model, bridging the gap between its theoretical foundation and empirical utility (Wilson, 1970). The gravity model's popularity can be attributed to its concise formulation, and its simplest version, the unconstrained gravity model, can be written down as the equation (2).

$$T_{ij} = \sum_i \sum_j K \frac{O_i^\alpha D_j^\gamma}{f(d_{ij})} \tag{2}$$

In the gravity model, each variable and parameter have a straightforward geographical interpretation, whereas $O_i^\alpha$ and $D_j^\gamma$ represents the production form origins and attraction of the destinations, respectively. The distance decay effect is represented as $f(d_{ij})$, which is the core of the gravity models.

Wilson's contributions extended beyond this methodological innovation, as he also delineated four general types of spatial interaction models (as presented in Table 2.1 below): unconstrained (where both Oi and Dj are unknown), production-constrained (where Oi is known but Dj is unknown), attraction-constrained (where Oi is unknown, but Dj is known), and doubly constrained spatial interaction models (where both Oi and Dj are known) (Wilson, 1971).

*Table 2.1 Family of gravity models*

| Model Forms | Formula |
|---|---|
| Unconstrained | $T_{ij} = A_i B_j f(d_{ij})$ |
| Production-constrained | $T_{ij} = A_i O_i f(d_{ij})$ |
| Attraction-constrained | $T_{ij} = B_j D_j f(d_{ij})$ |
| Doubly constrained spatial interaction model | $T_{ij} = A_i B_j O_i D_j f(d_{ij})$ |

Depending on the availability of data and research objectives, researchers can select the appropriate form and approach to calibrate the model to approximate its parameters. These fundamental forms of the gravity model gained further credibility in 1971 when Alan Wilson demonstrated their reliability by integrating classical transportation theory and entropy maximization theory into the framework.

After determining the basic formulation of the gravity model, the function form of the relationship $f(d_{ij})$ between travel cost with the travel flow needs to be chosen. The most common formats of $f(d_{ij})$ are for populations of origin and destination and exponential laws $(\exp(\beta d_{ij}))$ or laws of power $(d_{ij}^{-\beta})$ for dependency on distances (Dennett, 2012). These specific functional forms are chosen to allow for quick and precise calibration of the model. This makes it possible for researchers to use linear regression approaches to predict the value of parameters by functions of the population and distance (Barbosa et al., 2018).

During the application of the SI model, the parameters calibration process is seen as the critical step for fitting the prediction model to match the real-world flows. One of the mainstream method to calibrate parameters is applying general linear regression after exponential transformation of equations (Flowerdew & Aitkin, 1982). After the log formed transformation, equation (2) can be written as equation (3)

$$t_{ij} = \exp\left(K + \alpha O_i + \gamma D_j - \beta ln d_{ij}\right) \tag{3}$$

The flows addressed by spatial interaction models, such as migration or commuting, pertain to non-negative integer counts. Consequently, the probability of migration or commuting is not delineated by a continuous (normal) probability distribution, which typically underlies the error distribution in standard linear regression models. Instead, it is characterised by a discrete probability distribution, such as the Poisson distribution or the negative binomial distribution, with the Poisson distribution being a special case of the latter (Dennett, 2018; Flowerdew & Aitkin, 1982). The equation (3), therefore, could apply the Poisson regression to estimate the value of parameters $\alpha$, $\gamma$, and $\beta$ in the programming platform or software.

Another method to calibrate the parameters is the iteration-based calibration algorithm based on Maximum Likelihood Estimation. (Batty & Mackie, 1972) proposed calibration procedure which picks the distance decay parameters $\beta$ by continually executing the iterations of standard non-linear optimised until the difference between the predicted mean trip cost C and the observed mean trip cost $C^{obs}$ is less than the pre-set threshold $\varepsilon$. Technically, this method usually performs better but can only apply to single (distance decay) parameter models.

$$| C^{pre} - C^{obs} | < \varepsilon \tag{4}$$

$$C = \frac{\sum_i \sum_j T_{ij} d_{ij}}{\sum_i \sum_j T_{ij}} \qquad (5)$$

Another commonly used iteration-based calibration method is designed specifically for doubly constrained gravity models. This method typically involves establishing two balancing factors, each associated with origins and destinations, and allowing them to mutually influence each other. Subsequently, an initial value is set, and iterations are performed until a specified convergence criterion is achieved (Dennett, 2012; Plane, 1984).

## 2.1.2   Localised gravity models

Most previous spatial interaction models have operated under the assumption that the internal space of the modelling area exhibits spatial isogeneity, meaning that the distribution of trips follows a single overarching law associated with $f(d_{ij})$. Previous research suggested that spatial heterogeneity widely exists in the spatial interaction model and may reflect uneven trip distribution within urban space due to factors like system boundaries, transport accessibility, and other intricate urban contexts (Fotheringham, 1981; Nakaya, 2001; Zhang et al., 2022). The spatial heterogeneity effect, often referred to as the Modifiable Areal Unit Problem (MAUP), stemming from the configuration of zoning systems, has been extensively discussed. Spatial heterogeneity, which can lead to inconsistent results in spatial interaction models, has been perceived as a challenge, prompting efforts to identify optimal zoning systems or technical solutions to mitigate its impact (Arbia & Petrarca, 2011; Marceau, 1999; Openshaw, 1977). Some researchers have explored hierarchical structures to address MAUP issues during interaction estimation (Masser & Brown, 1975). Building upon the MAUP concept, others have proposed that implementing a hierarchical structure within spatial interaction models could

diminish spatial heterogeneity among sub-systems, thereby enhancing predictive accuracy (Fotheringham et al., 2001; Nazara et al., 2006; Qian et al., 2020). The hierarchical spatial interaction model distinguishes trips between intra and inter-subsystems for flow estimation, effectively reducing spatial heterogeneity at sub-system borders.

Since the 1970s, scholars have been aware that the spatial factor could affect the SI significantly, and correctly representing the spatial effect in the model could improve the performance of SI models (Masser & Brown, 1975; Openshaw, 1977; Oshan, 2020). Thus, another attempt is to modify the gravity models by highlighting the local characters in spatial interaction models. The Origin-Specific SI model is the most successful branch, which improves the fitting performance by separating the flows by subsets of each origin and calibrating parameters in each sub-model (Fotheringham, 1981; Fotheringham & Brunsdon, 1999; Gould, 1975). Subsequently, a set of parameter estimates for each model term is derived for each origin, which can be mapped to explore potential spatial variation (Oshan, 2020). Following the classic unconstrained (or sum-constrained) gravity model written as equation (2) above, this study adopted a disaggregated spatial interaction model referenced by the previous research (Fotheringham & Brunsdon, 1999), which divides the flows by origins and then fits the flows with separate models in the formatting of the classic unconstrained gravity model (6). For giving a specific origin, the $O_i$ is part of the constant (7). Each sub-gravity model has its own distance-decay parameters calibrated by the general linear regression model.

$$t_i = \sum_j t_{ij} = O_i{}^{\alpha_i} \sum_j K \frac{D_j^{\gamma_i}}{d_{ij}^{\beta_i}} \tag{6}$$

$$T = \sum_i t_i \qquad (7)$$

Based on this method, Fotheringham (1983) developed and extended the Origin-Specific that the distance-decay parameter can be an accessibility indicator to show the competing destinations. Nakaya (2001) proposed a method to model the immigration flow within Japan by using Geographically Weighted Regression (GWR) to calibrate the spatial interaction model with local parameters. By differing the bandwidth of GWR, different local parameters for both origin-specific and destination-specific distance decay parameters are generated.

On the contrary, some scholars have argued that spatial heterogeneity may have a positive aspect in detecting agglomeration effects (Menon, 2012). Given that the functional space of cities relies on how residents perceive their activity domains and interact with urban environments (Lynch, 1960), certain researchers have recognised the connection between border effects and spatial heterogeneity within spatial interaction models. They have endeavoured to quantify the border effect between zones using spatial interaction models (Engel & Rogers, 1994; McCallum, 1995; Yin et al., 2017).

### 2.1.3 Other types of spatial interaction models

Unlike taking distance as the core factor of spatial interaction in the gravity model, another set of human mobility models is the intervening opportunities model. In 1940, Stouffer suggested that "the number of persons going a given distance is directly proportional to the number of opportunities at that distance and inversely proportional to the number of intervening opportunities". The definition of opportunities could vary depending on the social phenomena investigated. Specifically, the opportunities could be jobs, market size or public services. Since the distribution of opportunities can be extremely heterogeneous in space, distance, therefore,

has an indirect effect on the final assignment of trip destinations and, as a result, on the decay of overall flows. The concept can be written as the formula (8) and (9) below:

$$T_{ij} \propto \frac{1}{x} \frac{dx(r)}{dr} \tag{8}$$

and

$$V_{ij} = \frac{dx(r)}{dr} \tag{9}$$

Where $x(r)$ in the equation stands a cumulative number of intervening opportunities in given the travel distance $r$. An application of the intervening opportunities concept uses intervening opportunities $V_{ij}$ to replace the distance or deterrence in the gravity models like the formula below. The intervening opportunities model may be described as a specific form of the gravity model.

$$T_{ij} = A_i B_j O_i D_j f(V_{ij}) \tag{10}$$

Based on the statistical results of abundant previous research, the advantages of the intervening opportunities model have been proven to explain the mobility data at a broad range (Kang et al., 2015; Stouffer, 1940). However, the intervening opportunities model has lost popularity in recent years due to the lack of research effort into the implementation and calibration of the model (Barbosa et al., 2018).

A new variant of the intervening opportunities model is the so-called radiation model, which established on Schneider's hypothesis: "The probability that a trip ends in a given location is equal to the probability that this location offers an acceptable opportunity times the probability that an acceptable opportunity in another location closer to the origin of the trips has not been

chosen" (Schneider, 1959). It assumes travellers would make the best possible choice for statistical treatment simplification (Simini et al., 2012). The density of opportunities is related to the population. There are two steps for how a traveller chooses a destination based on the assumption of the radiation model:

1. The quality of the traveller's opportunity in every location is represented by a number, z, which is allocated by some distribution p(z).

2. The traveller would pick the closest locations with opportunity quality higher than the traveller's threshold, which is another random number extracted from the fitness distribution p(z).

$$T_{ij} = O_i \frac{m_i m_j}{(m_i + s_{ij})(m_i + m_j + s_{ij})} \tag{11}$$

Where:

$O_i$: Total number of trips departing from location i

$m_i$ and $m_j$: Total opportunities at location i and location j.

$s_{ij}$: Total population in the circle of radius rij centred at i (excluding the source and

destination population)

The most important feature of the radiation model, also the reason why it has been favoured by practitioners, is the simple form and parameter-free property (Kang et al., 2015). On the other hand, one shortcoming of radiation the elements with substantial importance like spatial scale and heterogeneity, are overlooked in the model due to the parameter-free property, limiting the simulation abilities at the city level (Masucci et al., 2013).

Some other spatial interaction models have been proposed recently, including the population-weighted opportunities model (Yan et al., 2014) and the rank-based gravity model (Noulas et al., 2012), which are still based on the two basic frameworks of gravity and intervening opportunities theory. Pieces of research have been conducted to compare the performance of the gravity model framework and the intervening opportunities model framework in the decades since both frameworks developed (Lenormand et al., 2016; Pyers, 1966) Still, no conclusion has been drawn for which model is better than the other (Barbosa et al., 2018). Thus, some scholars attempted to combine the characteristics of these two frameworks with improving the goodness of fitting, Kang et al. (2015) proposed a 'generalised radiation model' which introduced gravity model-liked parameters into the radiation model system to fit more variety of mobility systems, improved the goodness of fitting of origin radiation model. In addition, updating the calibrating methods related to specific research topics by various data sources is the current research hotspot for spatial interaction models. Mobile phone data (Gao et al., 2013), social media check-in data(Liu et al., 2014), and smart card data (Zhong et al., 2015) have been used in the spatial interaction model under urban research topics.

### 2.1.4 Combining socioeconomic characteristics and travel behaviour

An existing criticism states that many studies explore the connections between land use and travel patterns neglecting to incorporate the socioeconomic dimension, and the omission causes oversimplified perspectives on the relationship between land use and travel (Stead, 2001). Various studies have supported the idea that the spatial relationship between workplaces and residences cannot be the sole explanation for observed commuting patterns (Kitamura et al., 1997; van de Coevering & Schwanen, 2006). Even in areas with abundant transport options, marginalised groups might find it challenging to utilise urban services due to other major obstacles (e.g., health conditions, poverty) hindering their participation in these opportunities

(Bradshaw et al., 2004). Thus, individual-level socioeconomic characteristics, individual preferences, and attitudes also play a role in influencing people's commuting behaviours (Lin et al., 2015). In 1982, Hanson posited that socio-demographic factors often have a greater impact on various travel behaviours than the aspects of urban spatial layout. In a study of the Boston metropolitan area, Shen (2000) elucidated that both the configuration of urban spatial structures and specific socioeconomic attributes play a pivotal role in determining the patterns of commuting durations.

Various socioeconomic characteristics could reflect the variation in travel behaviours. It has been widely recognised that age, income, and life stage significantly and interactively influence travel behaviour, impacting how individuals move and commute (Kattiyapornpong & Miller, 2009; Lin et al., 2015; Srinivasan & Rogers, 2005). With the trend of globalisation, migration and ethnic groups have also been identified as important factors that could affect travel behaviours (Hu, 2017; Mattioli & Scheiner, 2022). As an effective indicator to reflect individuals' socioeconomic characteristics, housing conditions have been drawn attention to their relationship with travel behaviours and the correlation has been confirmed (Jain & Tiwari, 2019; Scheiner, 2006). In addition, multiple research projects found that private car ownership could lead to residents having a larger coverage area by travelling longer distances and enjoying higher travel efficiency(Gao et al., 2022; Haque et al., 2019) .

In previous research, the relationship usually be investigated by establishing a straightforward regression model between socioeconomic characteristics and typical indicators of travel behaviours (i.e., average travel distance/time, travel frequency) (Hu, 2017; Mattioli & Scheiner, 2022; Srinivasan & Rogers, 2005). Thus, the majority of research about socioeconomic chrematistics and travel behaviour tends to draw analytic conclusions, but previous research has not focused on using socioeconomic characteristics in its model to predict travel flows.

## 2.1.5 Research gap for applying spatial interaction models in small zones within urban system

Recently, the evolution of computing power and data collection/storage techniques enabled the possibility of using spatial interaction models to build larger urban models with granular spatial resolutions, and this trend is known as "the renaissance of large-scale modelling" (Batty & Milton, 2021). However, a research gap for applying the localised spatial interaction model is that most localised spatial interaction models stay at the relatively macro level (e.g., province/state/ regional level), and do not go further into a finer spatial resolution within the urban systems (Dennett & Wilson, 2013). This is due to some associated issues that sometimes prevent the localised spatial model from being utilised in predicting flow within the urban system. As the spatial resolution becomes finer, the number of origin-destination pairs grows exponentially. This proliferation of data points increases the model's complexity and computational demands. Each sub-model may have its own set of parameters, necessitating separate fitting processes. The need to fit and validate these sub-models further adds to the computational burden.

Another issue is the local calibration in the origin-specific gravity model may be invalid in low-flow volume areas within the urban system. In granular spatial resolutions, some areas may have low flow volumes recorded. Local calibration requires enough data to make reliable estimates, in the case of grouping the flow by its origin area, those areas with a lot of zeros flow could lead to imprecise and unstable parameter estimates (Fotheringham & O'Kelly, 1989). Practically, Poisson regression sometimes cannot provide reliable results due to uncontrolled model error and the non-linear iteration method may not be able to converge in some situations.

## 2.2  Urban spatial structure for the mega-city region

The study of the spatial structure within cities can be traced back to von Thünen's *The Isolated State* (1966). The location theory focuses on how agricultural resources can be allocated to different distances from the market (urban centres), thereby improving the efficiency of agricultural operations. Urban spatial structure refers to the arrangement of urban space concerning the relationships arising out of urban form and its underlying interactions, composed of people, freight and materials, and information (Rodrigue, 2020). The urban spatial structure is one research topic in the urban geography field, that attracts much attention. Like politics, economic activities, topography, history, infrastructures, and policies, various factors interact with the urban structure and eventually form how city elements are located (Dadashpoor & Yousefi, 2018; Engelfriet & Koomen, 2018). Therefore, researchers have dedicated themselves to finding an 'optimal' spatial structure for encouraging the achievement of specific development goals of cities, such as economic performance (Wu & Yeh, 1999).

### 2.2.1  Describing the urban spatial structure

Although academics were aware of the importance of the urban spatial structure, the definition of urban spatial structure is arguable. Moreover, urban spatial structure is an interdisciplinary object of study, and it is difficult to form a common conceptual framework due to the different perspectives of various disciplines. The current prevailing interpretations may be categorised into **morphological form** and **functional form**, distinguished from the data sources and how urban structures are interpreted (Green, 2007). This debate is fuelled by evolving patterns in

location and transportation behaviour, as well as the influx of new data sources such as big data and crowd-sourced transport information (Thomas et al., 2018).

The urban spatial structures are traditionally described using morphological properties. The morphological method is based on traditional geographical data and survey data, such as population concentration and density, employment, and built-up area, to detect the CBDs and sub-centres (Zhang et al., 2021). This method is intuitional and easy to understand but has been criticised for failing to take the dynamic elements, such as human activities and interaction, into account. In contrast to the concept of morphology, functional structure emphasises the socioeconomic links between urban areas. Therefore, two distant areas can be integrated into a community because of the strong links of functional elements (Zhang et al., 2021). The previous research used various urban flows, i.e. commuting flow (Zhong et al., 2014), for detecting the functional structure of cities.

Another dimension of the urban spatial structure is **monocentric-polycentric**. Monocentric refers to the concentration of population, employment, and other factors within a certain area in a single centre. At the same time, polycentric describes the relatively balanced distribution of population and employment in multiple centres. Previous studies have focused on understanding polycentric structures in terms of both morphological and functional (Burger & Meijers, 2012). In the monocentric model, the gradient of the density function gradually decreases as the total population rises, incomes increase, land rents rise, and transport costs fall. As a result, the spatial structure within the city shows a decentralization trend. With rapid urbanisation processing, decentralised population and employment have led to the monocentric city model becoming less effective in explaining the spatial structure of the inner city. To obtain the economic effects of agglomeration, the decentralised population and employment regroup in the suburbs, creating new urban sub-centres independent or dependent on the central

business district (Anas et al., 1998). The urban space organisation can be more decentralised and complicated with a city expansion. Apart from the central business districts (CBD), the more urban hub has built up around a traditional urban centre (Zhong et al., 2014). Thus, cities began to shift towards a polycentric structure from their original monocentric forms. In terms of monocentric-polycentric measures, city polycentricity is often measured based on the size pattern of the centres within the city. Basically speaking, more equilibrium in the size of the centres means more polycentricity in the city (Burger & Meijers, 2012).

At the city level, spatial structure is usually described by the degree of **concentration** or **decentralisation**. In contrast, activity at the local level can be described as clustered or dispersed, depending on whether the distribution is a polycentric/monocentric pattern. Anas et al. (1998) suggested that centralisation depends on the extent to which urban activity is concentrated in the immediate area of a city's central business district (CBD).

Last but not least, the conceptions we introduced are not isolated. These conceptions can be employed at the same time to describe the urban structure of the same city. It may depend on the research topic and measuring method. For instance, Greater London has been identified as a morphologically monocentric region but a functionally more polycentric region (Hall & Pain, 2006). Taking the description by Zhang et al. (2021), the relationship between different conceptions of urban spatial structure is shown in Figure 2.1 below:

Metropolitan Boundary

Core Centre

Non-Centre settlements

Centralised

Monocentric — Polycentric

Dispersed

*Figure 2.1 The Conceptions in Urban Spatial Structure*

Based on the basic conceptions of urban spatial structures, scholars attempted to establish various models by empirical research to explain the urban structure and its evolutionary pathway. The negative exponential model is the most widely used in studying monocentric urban spatial structure among traditional monocentric models like the classic Alonso-Mills-Muth model (AMM model). Therefore, the negative exponential model's density gradient is commonly used to measure the concentration/decentralisation in urban structures with evolutionary characteristics (Mieszkowski & Smith, 1991). With more discussions about decentralisation and polycentric cities, various measuring methods about monocentric-polycentric, and centralised-decentralised were emerging. Helsley and Sullivan (1991) proposed a dynamic model of intra-urban spatial structure, in which the polycentric spatial structure is considered as a system formed by primary and secondary centres together, which

gradually goes through the stages of primary centre development, secondary centre development, and primary and secondary centre co-development in chronological order, with secondary centres emerging after the primary centres and their formation being influenced by both agglomeration economy and agglomeration diseconomy. Thus, the formation of secondary centres is influenced by both agglomeration and deagglomeration economies. Henderson and Mitra (1996) proposed the Edge City model, which incorporates land developers and historical factors into a model of the evolution of spatial structure within cities. They argued that land developers play an important role in forming urban sub-centres. The building plans of land developers in medium-sized cities can influence the locational choices of residents and businesses.

## 2.2.2 Detecting the urban spatial structure

Researchers have conducted studies on the urban spatial structure at two levels due to spatial scale differences: the intra-city level and the city-region level. For intra-city polycentric spatial structures, scholars have focused on the location, and morphological attributes of the newly emerged centres in the evolution of cities from monocentric to polycentric and then analysed the systemic characteristics and interrelationships between the internal centres. Meanwhile, in regional-level or country-level polycentric studies of spatial structure, studies usually take administrative cities as the centre of regional spatial structure rather than searching the urban centres by detecting method (Gao et al., 2017). Some studies (Gordon & Richardson, 1997; Richardson, 1969) have found that there is a significant relationship between metropolitan spatial structure and economic growth, depending on metropolitan size and its structural organisation.

In the study of intra-city spatial structure, individual cities are generally considered as a surface. Thus, the intra-city spatial structure reflects the interaction of elements within the urban territory. According to the study by McDonald's (1987), urban sub-centres are generally defined by these two perspectives: (i) they are areas of high regional population (employment) density; (ii) they concentrate enough population (employment) to have an impact on the surrounding area.

Based on the definition of sub-centres, scholars have developed various methods to identify morphological urban sub-centres. In the early research, the threshold method was the most common method to identify sub-centres. However, the threshold method is highly subjective, as different thresholds for the same city often give different results. Therefore, thresholds in different cities are not comparable (Anas et al., 1998). In addition, this method usually ignores the suburban area due to relatively lower density. The relative threshold approach was introduced to reduce the subjectivity of threshold setting and, more importantly, compare cities. For example, some studies have selected employment density above the mean (or above the mean plus one or two standard deviations) and employment above 1% of total urban employment as criteria for identifying urban centres (Garcia-López & Muñiz, 2010; Muñiz et al., 2008). Moreover, the relative threshold approach was applied not only to street and neighbourhood vector spatial data but also to population density raster data (Liu & Wang, 2016). Methods such as kernel density mapping are also used for identifying urban centres. Gordon et al. (1986) and Maoh and Kanaroglou (2007) searched for urban centres from the high-value areas in the raster map employing density mapping. Although the graphical method can visualise the internal spatial structure of the city, shortcomings such as the inability to identify the extent of urban centres limit the application of the graphical method. Besides, scholars have used spatial statistical analysis to identify urban areas with significantly higher employment

(population) densities than their neighbours as urban centres (Arribas-Bel & Sanz-Gracia, 2014; Asikhia & Nkeki, 2013; Hajrasouliha & Hamidi, 2017). For this approach, only the neighbourhood range and significance thresholds need to be selected based on the local Moran index and the local G-statistic.

Excluding the numerical methods introduced above, geometry methods like city fractals have also been employed to detect the urban morphological structure. This is because the geographic world contains many phenomena without characteristic scales, which cannot be effectively portrayed by traditional mathematical methods but can be described spatially and statistically using fractal geometry (Zhang, 2018). However, due to how land-use data is summarised and ignoring differences in processes operating at the micro and macro scales (Anas et al., 1998), the graphical methods are not taken as often as quantitative methods.

In parallel with morphological discussion, the functional structural is another strand for urban structure detection. Green (2007) firstly proposed a framework for measuring functional polycentricity drawing on the conceptions used in social network analysis. It defined the indicators of polycentricity by classic topology and network analysis methods, such as nodality, centrality, and network density.

Human mobility big data can reflect the interaction between regions, using human mobility big data to detect the functional urban structure has become a hotspot of urban studies. For example, trajectory data such as taxi and bus swipe cards, and mobile phone signalling can generate OD flow. These data can identify the centre areas in the network, further identifying the city's spatial structure (Jiang et al., 2017; Zhong et al., 2017). In 2014, Zhong et al. detected and depicted urban structure in Singapore by graph-based community detection algorithm, and it is one of the representative studies for functional urban structure detection. According to their

analysis, they suggested collective human mobility can shape geographic communities like social networks. The network method may explain the composition via structural shifts of transient sub-centres. For example, it might describe the increasing interaction between certain developing sub-centres (Zhang et al., 2021). Recently, Shen and Batty (2019) detected community structures in the London Metropolitan area based on disaggregated flow data, suggesting that the functional structure may vary for people with different occupations. Zhang et al. (2021) analysed multi-year transport smart card data in London, the results of network community detection show that Greater London can be clustered into five communities based on the travel pattern, but London moved towards a more polycentric and compact urban structure.

## 2.2.3 Research related to detecting cities' boundary

As we introduced in last section, the main task of detecting the Urban spatial structure at the intra-city level is mainly about detecting the urban (sub-)centres and describing the relationship between centres. But at the city-region or larger scale, studies usually take administrative cities as regional spatial centres. Therefore, the priority in city-region spatial structure shifts to detecting the city boundaries. Once the border is fixed, the size and the interaction relationship between centres can be used to discuss the city region spatial structure similar to the intra-city urban structure analysis. However, it is always a great challenge for urban geographers and planners to define city boundaries. Cities combine multiple systems from local to national on different spatial scales and varying temporal scales from day-to-day operations to those that run over decades. Therefore, the concept of system boundary is a significant issue. The word 'boundary' refers to something that indicates or fixes a limit or extent (Merriam-Webster dictionary). Owing to the constant interaction of urban sub-systems into dynamic wholes, attempting to separate systems' facets is difficult. Besides, a persistent source of error known

as the Modifiable Areal Unit Problem (MAUP) is one of the key reasons why cities' boundaries must be specified (Openshaw, 1984). They pointed out that the changing system's spatial boundaries in a zonal system may significantly affect the overall statistical properties. All zonal data could be influenced by the MAUP and compounded by the fact that zonal borders are often arbitrary or set for purposes incidental to the study intention (Openshaw, 1996). A prevalent example is that local governments often only consider administrative boundaries as zonal boundaries due to their job duty. In this case, the difference between administrative boundaries and the range of people's movements and industries' activity may distort simulation results. The situation could be worse for defining urban clusters since the urban term is a very vague one. The issue of which feature of urbanity produces the efficiency premium is also unclear due to the lack of theoretical formulation on the spatial side of agglomeration economies (Bretagnolle et al., 2002; Parr, 2007).

If we focus on previous research in terms of urban boundary or boundary detection, many previous researches dedicated to identifying the area of the built environment from remote sensing data by novel classification algorithms (Henderson et al., 2003). However, there is relatively little research that discusses the boundaries defined by invisible activities such as people's movement or economic agglomeration.

Calafati and Veneri (2013) highlighted that Italy's spatial polarisation and territorial integration processes since the 1950s have not been matched with necessary institutional changes. This mismatch has led to a significant gap between the territory's functional and political-administrative organization. They argued that in Italy, the central city, along with its first and second rings of municipalities, should be recognised as a "fact city" rather than considering only the central city or a broader area like the entire functional urban region (FUR). To support their viewpoint, they employed basic geographical and economic statistical data, drawing

comparisons based on distances among urban settlements, commuting-to-work patterns, changes in the spatial distribution of population and employment, and patterns of residential and employment density.

Arcaute et al. (2015) developed a framework to consistently define cities, using commuting to work and population density thresholds, and construct thousands of realizations of systems of cities with different boundaries for England and Wales. The research was based on the census data of the UK in 2001/02 given by ward level, along with demographical data like household income, and land use data including road facilities, paths and buildings. They employed the concept of ubiquitous scaling law:

$$Y(t) = Y_0(t)N(t)^\beta \qquad (12)$$

Where Y(t) and N(t) represent the urban indicator and the population size of a city at time t respectively, and Y0(t) is a time-dependent normalization constant, with the scaling exponent β represented defined by the nature of the urban observable. Based on the 10 years period simulation results and the calibration with the real data, they suggest continuous wards with a density of more than 14 persons per hectare could be an appropriate threshold for defining the urbanised space. In addition, they particularly mentioned that for a mega-city like London, its strong role as an information and economic hub suggests that the urban system is highly integrated and that it is difficult to partition the system into individual cities that capture these social interaction effects. The research scaling of mega-city should not be limited to regional or even national levels. Following Arcaute et al,' s work, Cottineau et al. (2019) attempted to build a comprehensive representation of where cities extend by a case study in France. They used census data to identify the 'night-time cities' and 'day-time cities' by using residential and workplace geographies. After predicting aggregated wages Y from the scaling equations

specified for population and density by the OLS-regression model, they confirmed that the urban aggregation economic effect would relate to jobs and residential density. Based on this result, the belief that the scale of analysis (local, metropolitan, regional) is therefore critical in detecting the agglomeration effect on wages or not, since it means that the entire or only parts of the urban system are counted in the calculation. They also suggested that mega-cities tend to be either wealthier or equally affluent as smaller cities, but never poorer. In addition, larger towns appear to be more or similarly unequal than smaller towns, but on average never more equal. In the last part of the paper, they admitted that there should be a way to define the urban cluster to discuss the effect of agglomeration economies, and the effect could be affected by policy and infrastructure such as roads, the topic should be a research gap which needs further investigation.

Besides, some researchers also attempted to define or detect the urban boundaries from morphological observation (Tannier et al., 2011b), and the transport network's density (Long, 2016; Long et al., 2018). Different data and analysis methods have been employed to achieve the targets, but these researchers observed cities from a macro insight. The statistical number or geographical could shape an image of cities from up to bottom. However, the spatial extents, in other words, the cities' boundaries, often overlap and agglomerate depending on how citizens perceive their activity space and interact with their urban environments (Lynch 1960). In recent years, defining cities' boundaries by new forms of data, especially spatial big data, enables detection of cities' boundaries from the daily activities of individuals living in the cities.

Yin et al.(2017) defined a method focused on human experiences with physical space inferred from social media to delineate urban boundaries. The hierarchical definition of non-administrative urban boundaries is from various movement activity ranges extracted from 1,153,891 users' collective mobility habits reflected by 69,847,497 tweets made geotagged on

Twitter (now X). They found a 92% possibility for collective displacements within the range (10m, 70km) and 10km is the distinct geographic distance for separating two main groups of Twitter users regarding the UK's spatial coverage. They established a spatial interaction model to depict urban structures by overserving the distance decay effects and drew the nonoverlapping boundaries of UK cities. The limitation of the boundary is that social media data like Twitter data cannot represent the complete real-world image, and the spatial sparseness of geo-located Twitter data could make the relatively small cities be ignored (Stefanidis et al., 2013). Therefore, it is worth to use more representative data to identify the boundaries of citizens' activities.

To conclude the previous research about urban boundaries, defining cities' boundaries is a meaningful research question since an urban system often does not coincide with the administrative boundaries, which may cause distortion and planning failure due to an effect similar to MAUP. However, current research for defining and measuring the urban system's boundary is minimal, and the measuring methods and indicators are usually limited to population density. Fortunately, defining urban boundaries from "bottom to up" becomes possible since the big data records individuals' daily movements. Still, more data dimensions and methods activities need to be considered in further research. Besides, none of the research considers the city's boundaries as dynamic processing means the relationship between city boundaries and specific urban interventions is still unclear. Thus, as one of the most commonly used indicators to modelling the functional urban structure, human mobility could be a good entry point to solve this unanswered question.

## 2.3 Big data in human mobility research

Acquiring proper data has been identified as the major challenge to using metrics in human mobility and comparative studies between cities (Boisjoly & El-Geneidy, 2017). The starting point of research on a relatively large spatial scale normally is census data, which can provide basic demographical information at the city, regional, and country levels. Currently, a census survey in most countries includes the name, ethnic group, age, gender, occupation, and marriage of each resident in a household (Tizzoni et al., 2014). One obvious advantage of census data is it almost covers all population in a country since it is mandatory for citizens. Census data in some countries also include income, education, workplace location, and household type (Statistics, 2021). Thus, census data is very reliable and informative for researchers. However, it is implausible to access the individual-level information of census data due to privacy issues: the individual socioeconomic data is highly sensitive because it can identify personal identity without pre-authorised, which may cause an ethical argument even risk of criminal in case of leaking. Thus, the census data has often been aggregated at a relatively rough level. Therefore, research at the finer level (spatial and/or temporal) usually requires additional data sources.

The travel survey data can be seen as another data source other than census data for measuring human mobility. The census is more about demographics, not designed for travel, but the travel survey is designed to provide data on personal travel. The most common travel survey data is at cities level, which mainly focuses on short-term trips with related information such as travel purpose, time, cost, and transport mode, enabling researchers to investigate the movement pattern and establish the transport model. Nevertheless, unlike the census data, most travel demand surveys take sampling data in a specific region. Therefore, sample selection bias is a

common issue for the survey data. Besides, when considering inter-city and intra-city trips simultaneously, survey data quality is not satisfying since the study area usually cannot be fully covered by a single trip survey. In addition, the National Travel Survey (NTS) is another widely used data source of travel data, which primarily designed to track long-term development of trends since 1960s (Morris et al., 2013). Excluding census and trip survey data, some traditional data sources can indirectly reflect human mobility, such as tax revenue data or currency bill data. These sources have not been regarded as widely used data sources in existing human mobility research (Barbosa et al., 2018).

In recent decades, new forms of data have driven a revolution in human mobility measuring. The traditional method of obtaining travel origins and destinations (OD) is to conduct travel surveys, which are very time-consuming and costly, and the accuracy of which needs to be improved. In the era of Big Data, various data sources have emerged to obtain travel OD without conducting field surveys, including mobile phone signalling data, GPS location data, Smart Card Data, etc.

*Table 2.2 Mainstream type of new form of mobility data*

| | Mobile Phone Signalling Data | GPS Data | Smart Card Data | Social Media Data | Mobile Application Data |
|---|---|---|---|---|---|
| **Coverage** | All mobile phone users in signal-covered areas | Individuals or vehicles with GPS tracker | Public Transport passengers who use smart card | Social media users who publish content with a geotag. | Users who use related mobile application |
| **Location/spatial precision** | Up to but usually provided by aggregated level due to privacy | 3-5m for outdoor | Locations of public transport nodes | Depending on the Geotag | Depending on the method of location information collected |
| **Data Availability** | Hard to obtain | Hard to obtain | Easy to obtain | Easy to obtain | Hard to obtain |
| **Pros** | -Huge data size<br><br>-Relatively high resolution both in spatial and temporal | -Highest spatial and temporal accuracy | -Huge data size<br><br>-Record individuals' daily movement activities | -Easy to obtain<br><br>-Textual message is associated with geoinformation | -Rich dimension in user's behaviour<br><br>-Large geographical extent |
| **Cons** | -Computational expense is high<br><br>-Data availability due to privacy issues | -The sample size is very limited<br><br>- Unavailable or lower precision when indoor | -Data is limited by public transport nodes, lacking other information. | -Data points usually are scattered in temporal and temporal distribution<br><br>-User group cannot cover the majority of the population | -User group cannot cover the majority of the population<br><br>- Data quality largely depends on the use frequency |

Mobile phone data, which has been seen as the most important 'game-change data' (Barbosa et al., 2018), has been widely used in related research. Mobile phone signalling data is generated by cell phone users in the event of calls, text messages or mobile location, captured by the operator's communication base station and recorded by the same user signalling trajectory, the sample of mobile phone signalling data can be found in Figure 2.2 below. After decryption, desensitization, expansion of the sample and other processing, it can be applied to research on human mobility and other urban applications such as town spatial layout. The spatial



*Figure 2.2 Data sample of Mobile phone Signalling data. (Source: (Song et al., 2010))*

resolutions of mobile phone signalling data usually are in the range of 50m to 5km, depending on the density of stations, while the temporal resolution can be accurate to seconds (Jiang et al., 2017). Mobile phone data has been widely used in discussing the human mobility pattern both at the individual level (De Domenico et al., 2015; González et al., 2008; Lu et al., 2013) and population-levels (Palmer et al., 2013; Phithakkitnukoon et al., 2012; Zhou et al., 2020). One significant advantage of mobile signalling data is its wide coverage - as long as the user switches on the phone, the data will be automatically captured (González et al., 2008). Besides,

because it records high-precision individual movement trajectory, mobile signalling data have been applied to research with different spatial scales from neighbourhoods to the country-wide level (even international travel and movement across borders). The shortcoming of mobile phone signalling data is the lack of data availability due to privacy issues because it may contain private information about the user, requiring applying desensitisation algorithms before being provided to researchers (Barbosa et al., 2018; Williams et al., 2015). Another disadvantage is since the data size of signalling data is huge, it may require an extremely high computational cost.

Another data source of human mobility analysis is the Global Positioning System (GPS) data, which records the highest precision trajectories of movement located by satellites at regular intervals (usually in seconds) (Zheng et al., 2009). GPS device-attached vehicles are the one of main data providers for GPS data worldwide (Pappalardo et al., 2013; Yuan et al., 2010). Human mobility topics related to road traffic benefited from these data for discussing issues such as congestion, travel cost, and taxi accessibility (Bazzani et al., 2011; Li et al., 2012; Pappalardo et al., 2013). Individual GPS trackers are another part of GPS data which requires pedestrians or cyclists to carry the GPS device with them. An example of research utilising individual-carrying GPS data is Geolife (Figure 2.3), which contains 17,621 trajectories recorded per 1–5 s or 5–10 m for 182 individuals in a period of over three years. Those data have been applied to mining the transport mode, point of interest (POI), and purpose of travel (Zheng et al., 2010; Zheng et al., 2009). Compared with mobile phone data, the drawback of GPS data is that the typical data size of GPS data is tiny, usually no more than several thousand individual users.

*Figure 2.3 Example of taxi GPS dataset in Beijing, source (Yao et al., 2021)*

Smart Card data arouses researchers' attention because Smart Card has been replacing paper tickets or single-use tokens in cities worldwide in recent decades. When cards become affiliated with specific individuals, it enables the possibility of recording each trip by individuals, capturing relatively precise spatial and temporal attributes such as the origin/destination stations and the stay durations, the data sample is shown in Figure 2.4 below. Thus, smart card data is an ideal data source for research related to human mobility patterns. Smart card data have also been widely applied to optimising public transport design and management systems for dealing with issues like delays, disruptions and congestion, improving the passengers' experience (Uniman et al., 2010). Transport planners and researchers benefit from Smart Card data for analysing travel patterns with heuristic and stochastic approaches at a disaggregated

level (Sari Aslam & Cheng, 2018). For instance, a series of research has been conducted utilising the Oyster Card in London, identifying the service reliability, deprived areas and temporal mobility patterns (Smith et al., 2013; Uniman et al., 2010; Zhong et al., 2016). The advantage of smart card data is it can reflect individuals' spatiotemporal movement patterns as a sequence of activity locations and durations daily. However, the spatial resolution of smart card data is limited by public transport nodes. Inferring secondary activities (other than work and residence) is difficult when supplementary information is lacking (Sari Aslam et al., 2021).



*Figure 2.4 Dataset of smart card data (source: Song et al. (2018))*

The last data format discussed in this review is the social media data. Since the era of smartphones, social media providers have been collecting valuable data, including social networks and geographical information. Whenever users publish content, social network providers such as Twitter, Facebook, and WeChat would collect geoinformation including the geographic coordinates, the time stamp, and additional information (such as POIs or contacts

nearby) (Barbosa et al., 2018). Various research methodologies have been applied to reveal human mobility patterns from massive social media data containing spatiotemporal information, including data mining, spatial statistical analysis and geo-visualisation (Gao et al. 2013, Zheng et al., 2010), while the typical applications of human mobility researches utilising social network data include recommendation POIs to the user based on the users' past trajectory, estimating local commuting patterns and visualisation the social interactions (Bao et al., 2015; Zheng, 2015). Data availability is one of the pros of social network data since the social media world is open to browsing for everyone, which means massive data can be simply acquired through free APIs.



*Figure 2.5 The social media checked-in data (source: Hu and Jin (2017)).*

Meanwhile, more contextual information can be collected associated with the geoinformation, contributing to classifying the data for application to the research under specific contexts. The cons of social media data mainly include the representativeness for the general population is questionable and massive invalid/fake information needs extensive data clean work. In addition, privacy concerns recently drove some countries to tighten laws to limit getting user data from social network providers, which may affect data availability in the future.

Analysing human mobility patterns through mobile application data offers many insights into how people move and interact within various environments. This analysis method has become increasingly popular with advancements in technology and the widespread use of smartphones. The most significant advantage of mobile application data is that it can provide rich dimensions of user behaviour depending on the purpose of the mobile application, particularly benefiting those research focuses on specific travel behaviour (e.g. shopping, tourism, dating). The spatial resolution of this type of data varies. Applications using GPS information (like navigation applications) provide highly accurate and precise location data, allowing for detailed analysis of movement patterns, but other application datasets can only provide the IP address as recorded to represent the geolocation information (Wang et al., 2021). Mobile application data are mainly collected by online mobile application operators, who treat their users' data as commercial assets(C. Hu et al., 2022). Thus, the open published dataset of mobile application data is rare. The main disadvantage of this data type is the data bias: most mobile applications serve specific use groups, and there are demographic differences between user and non-user groups. This can result in biased data not representing the entire population's mobility patterns. In addition, the quality of mobility data largely relies on the frequency of use: location data is collected when using those applications, but the location information during other times remains unknown, which makes it hard for low-use-frequency applications to generate reliable

trajectories. Some "Big Tech" companies (like Google, Baidu, Alibaba) have capability to combines data from different applications to overcome the disadvantages of data bias and low frequency (Quilty et al., 2020; Ruktanonchai et al., 2018), but this data collection behaviour itself is arguable due to the wide concern of privacy issues (Cohen & Mello, 2019).

There are also some other new form data sources such as point of interest (POIs) data, toll-fee data, and tickets data also be applied to human mobility research. However, these data sources often are used as supplementary data sources due to their limited number and application, insufficient sample size and lack of precision. A new trend is emerging that some researchers combine more than one data format to analyse the human mobility pattern, overcoming the drawback of a single data format. For example, the Smart card data and POI data were combined to infer secondary activities (Sari Aslam et al., 2021), and GPS Probe Data and social media data are incorporated to measure indicators of urban traffic congestion (Wang et al., 2017). Overall, big data and new data formats are changing the measuring method for human mobility, and it enables more possibilities for research. On the other hand, limitations persist in all other data sources. Keeping the balance between rich information and privacy protection is a significant issue for human mobility big data.

# 3 Delineating Urban Functional Zones using Mobile Phone Data

## 3.1 Introduction

The formation of city regions has been driven by urban growth and significant improvements in inter-city connectivity since the second half of the 20th century (Hall & Pain, 2006). Governments worldwide have encouraged regional cooperation to leverage efficiency benefits from agglomeration economies (Brenner, 2002). For instance, the Chinese government has implemented a series of policies since 2004 to foster integration among cities within one city-region. These policies aim to promote industrial cooperation and facilitate the sharing of urban functions (Li et al., 2015; Wu, 2016). As cities continue to expand and interact with one another, human daily activities related to urban functions, such as work, residence, and recreation, have extended beyond their original administrative boundaries and now occur across different cities. This phenomenon has given rise to the concept of urban functional zones (UFZs) (Gao et al., 2017; Yeh & Chen, 2019; Zhai et al., 2019; Zhong et al., 2014). The ambiguity surrounding UFZs poses new challenges for regional planning and management in response to the rapid development of mega-city regions.

Despite the widespread existence of this phenomenon, there are limited accurate quantitative methods to assess how UFZs have been integrated across cities. The emerging mobility data

provides an opportunity for a breakthrough to delineate UFZs of cross-city from human activity. The current mobility data mining techniques can track the daily movement flows of a huge population. Moreover, unlike the traditional surveys data conducted by local authorities, the new form of data (e.g., mobile signalling data and social media data) enables us to analyse a finer-grained networks beyond city/county boundaries. By taking the benefit of mobility big data, previous research mainly applied network analysis methods in many urban analytics applications, including detecting spatial structure and community detections (Jin et al., 2021; Shen & Batty, 2019; Wu et al., 2021; Zhong et al., 2014). However, the network-based analysis also show limitations as the distance decays effect often has not been appropriately reflected in the topology relationship (Liu et al., 2014; Yin et al., 2017). When detecting the communities across boundaries, this would cause the network-based interaction model to be insensitive for changes in cross-boundary flows (Liu et al., 2014). Thus, a novel method for detecting UFZs with more sensitivity for cross-boundary travel flow and distance decay effect is needed.

In this chapter, a critical hypothesis is that the boundary of a functional zone is highly associated with local distance decay. Zipf (1946) proposed that human mobility follows a spatial distribution with a distance decay from centres to the periphery. This concept has been accepted and applied in previous trip estimation models (Batty & Milton, 2021; Masucci et al., 2013; Wilson, 1971). Thus, the heterogeneity of trip distribution can be seen as an indicator to reveal the discontinuity of urban functional space or the "border effect" of urban functional zones (Brown et al., 2020; Jin et al., 2021). When crossing different urban functional zones, the border effect can be observed in human travel activity. Such border effects could be used as indicators for delineating urban functional zones (Jin et al., 2021; Rinzivillo et al., 2012; Shen & Batty, 2019). As one of the most widely used methods for predicting interaction flows, spatial interaction models predict the strength of spatial interaction based on the distance decay

effect. Previous studies confirmed the border effect can be represented by spatial heterogeneity in the spatial interaction model (Jin et al., 2021; McCallum, 1995), which provides a new method for delineating the UFZs and overcoming the limitations of network-based methods.

To achieve research objective 2 raised in Chapter 1.2 of delineating UFZs, this chapter will introduce a new method that first applies a two-level hierarchical spatial interaction model (HSIM) to generate the flow of spatial interaction between zones, then redraws non-overlap boundaries of urban functional zones by searching for the best partition with the best goodness of fitting in HSIM. By applying this algorithm to delineate the cities' functional regions within a specific mega-city region, the Shenzhen-Dongguan-Huizhou (SDH) area, in two different settings, empirical results prove that the goodness of fitting in HSIM can represent reasonable cities' boundaries. The results show that current UFZs almost coincide with administrative boundaries. Meanwhile, the results of long-term predictions remind policymakers to give more attention to the areas near the Dongguan-Huizhou boundary to promote industry cooperation and avoid serval mismatches between functional and administrative regions, providing implications for related regional planning policies.

## 3.2 Challenge in delineating the urban functional zones.

The urban spatial structure is topical research in urban geography. Various factors, such as politics, economic activities, topography, history, infrastructures, and policies, interact with the urban spatial structure and eventually form how city elements are geographically located (Dadashpoor & Yousefi, 2018; Engelfriet & Koomen, 2018). The current prevailing interpretations of urban spatial structure can be categorised into morphological structures and functional structures, distinguished based on the data sources and how urban structures are

interpreted (Green, 2007). In previous <u>literature review sections</u>, the basic concept about morphological vs functional urban spatial structure has been discussed. Some researchers also attempted to define or detect the urban boundaries from morphological observation (Tannier et al., 2011a), using data analytic framework such as the scaling law (Alvioli, 2020; Arcaute et al., 2015; Cottineau et al., 2019) and the transport network density (Long, 2016; Long et al., 2018). Compared with the morphological structure, the functional structure is more temporal and dynamic (Wu et al., 2021), which would better correspond to the rapid changes in the urban environment. Various urban flows like commuting and logistic flows within the city are being used to describe the urban spatial structure by spatial interactions (Burger & Meijers, 2012; Sohn, 2005; Zhong et al., 2014). Therefore, two distant areas can be integrated into a community because of the strong links of functional elements (Zhang et al., 2021). This feature would benefit understanding the integration of urban functional integration across cities.

Researchers also conducted studies on the urban spatial structure at two levels due to spatial scale differences:  the intra-city level and city-region level. For intra-city polycentric spatial structures, scholars have focused on the location, and morphological attributes of the newly emerged centres in the evolution of cities from monocentric to polycentric and then analysed the systemic characteristics and interrelationships between the internal centres. Meanwhile, in regional-level or country-level polycentric studies of spatial structure, studies usually take administrative cities as the centre of regional spatial structure rather than searching the urban centres by detecting method (Huang et al., 2015b; Gao et al., 2017). Therefore, delineating the urban functional zones between cities is essential for discussing the functional spatial structure in a city region or larger scope. There are some existing regionalisation algorithms are well-known for delineating regions based on indicators or objective functions. For example, P-regions and max-p is based on a defined objective function, meanwhile, REDCAP and

SKATER is based on hierarchical structure reflecting neighbourhood relationship (Duque et al., 2011; Guo, 2008; Helbich et al., 2013). These methods mainly use socioeconomic indicators (e.g., house price, income) to find the spatial clustering or non-spatial similarity rather than using flow data to evaluate the connection between areas.

### 3.2.1 Modifiable Areal Unit Problem (MAUP) and boundary effects

As reviewed in the section 2.1.2, most of the previous spatial interaction models assume that the inner space of the modelling area has spatial isogeneity, which means the distribution of trips only follows one general law related to $f(d_{ij})$. Meanwhile, section 2.1.2 also introduced the basic conception of the Modifiable Areal Unit Problem (MAUP). Since spatial heterogeneity may cause inconsistent results in spatial interaction models, previous research has regarded spatial heterogeneity as a "problem" and has attempted to find an optimal zoning system or technical solution to mitigate its effect (Arbia & Petrarca, 2011; Marceau, 1999; Openshaw, 1977). Besides, a few researchers attempted to adopt hierarchical structures to eliminate the MAUP issues during estimating interactions (Masser & Brown, 1975). Following the conception of MAUP, some researchers proposed that applying a hierarchical structure in the spatial interaction model may eliminate the spatial heterogeneity between each sub-system, improving the overall performance of prediction (Fotheringham et al., 2001; Nazara et al., 2006; Qian et al., 2020). The hierarchical spatial interaction model distinguished the trips between inner and inter sub-systems for estimating the flows respectively. By applying this framework, spatial heterogeneity on borders between sub-systems can be largely eliminated (Qian et al., 2020).

On the other hand, some scholars also argued that MAUP issues have the 'bright side' for detecting the agglomeration effect (Menon, 2012). Considering the cities' functional space

depends on how citizens perceive their activity space and interact with their urban environments (Lynch, 1960), some researchers were aware of the linkage between the border effect and spatial heterogeneity in the spatial interaction model and attempted to quantify the border effect between zones by spatial interaction model (Engel & Rogers, 1994; McCallum, 1995; Yin et al., 2017). However, how to use the variability of hierarchical boundary in the spatial interaction model as the indicator for delineating the urban functional zone or other types of communities has not been further discussed. Thus, it would be an interesting perspective to observe the boundary effect by trip distribution, providing a solid reference for delineating the boundaries of the urban functional zones.

### 3.2.2   Using human mobility data to understand urban functional zones

Origin-destination (OD) flow matrix generated from human mobility data (e.g. taxi and bus swipe cards, mobile phone signalling data) can be used as a proxy of the interaction between regions (González et al., 2008). Based on that, some studies have been developed to identify urban functional zones and urban spatial structures. Network-based methods is a commonly used approach based on the intensity of human interactions between different spatial units (Jiang & Miao, 2015; Louail et al., 2015; Zhang et al., 2020; Zhong et al., 2015). Each spatial unit is seen as a node, and human interactions are represented as edges between the two nodes.

In 2014, Zhong et al. detected and depicted urban structures in Singapore using a graph-based community detection algorithm, and it is one of the representative studies for urban functional zones detection. The network method may explain the composition via structural shifts of transient sub-centres. For example, it can describe the increasing interaction between certain developing sub-centres (Zhang et al., 2021). Shen and Batty (2019) detected community structures in the London Metropolitan area based on disaggregated flow data, suggesting that

the functional structure may vary for people with different occupations. Zhang et al. (2021) analysed several years of transport smart card data in London and the results of network community detection shows that Greater London can be clustered into five communities based on the travel pattern, but London moved towards a more polycentric and compact urban structure. However, the traditional network analysis and most community detection algorithms usually only consider the absolute value of flow volume (edges) for dividing the partitions regardless of the spatial factors such as distance decay or time consumption (Adam et al., 2018; Hong & Yao, 2019; Jin et al., 2021).

Some researchers have been aware of this and tried to apply spatial interaction to improve their method. Jin et al.(2021) identified the activity broad within Shenzhen city and discussed the boundary effect by using a modified spatial interaction model. Yin et al. (2017) proposed a method to delineate urban boundaries for Great Britain based on the physical space inferred from human activities of social media then verified the results by a gravity model. However, both still applied a network-based algorithm to identify the functional urban regions. Currently, due to the limitation of the data and computation, most of the previous studies on spatial structure from movement flow investigated one city only (Jin et al., 2021; Wu et al., 2021; Zhang et al., 2021; Zhong et al., 2014). This study would explore functional urban spatial structure at a larger scale by delineating the urban functional zones based on human mobility and movement flows. In addition, Because of the importance of distance factors and the absence of a method that can detect urban boundaries by distance-based trip distribution, we believe it is worth establishing a new method for depicting the form of urban boundaries.

## 3.3  Data collection

The case study area of this study is Shenzhen-Dongguan-Huizhou (SDH) area. SDH area covers a total area of 15,800 square kilometres, with a resident population of 26.25 million and a total GDP of RMB 3.7 trillion in 2019. This area has been experiencing rapid urban growth and change of urban spatial structure since the 1980s and became one of the most open and economically vibrant regions in China. SDH persistent attracts national and global focuses, especially after the Guangdong-Hong Kong-Macao Greater Bay Area was proposed in 2015.

This research uses the mobile phone data provided by one of the main mobile phone operators called China Unicom. The data contains Origin sub-district ID, Destination sub-district ID, the volume of travel flow and travel time. The spatial resolution of the original data is collected as 500m*500m but is provided as aggregated form into 172 sub-district level zones ("jiedao level" or "街道级" in China). The specific study units are shown in Figure 3.1 below. The mobile phone data detected 13,588,846 commuters (about 37% of the overall population), including both intra-city and inter-city. The observed period of the data is February 2019. To identify commuting trips, home and workplaces are first determined from one-month sequent locations of mobile phones.

Specifically, the site with the most prolonged stay during the observation period (9:00 pm-08:00 am) in a day is considered the candidate place of residence. When a candidate residence lasts for more than 15 days in a month, it is deemed to be valid. Similarly, the location with the most prolonged stay between 09:00 am and 5:00 pm is determined to be the workplace. Commuting is defined as a journey from one's home to the workplace. Individual commuting trips of mobile phone users are aggregated at the street scale, generating links between streets across the study area (SDH region), the data sample could be found in Table 3.1.

*Table 3.1 Data sample of GBA mobile phone signalling data*

| Origin Area Name | Origin Area No. | Destination Area Name | Destination Are No. | The volume of the flow | Avg time (by min) |
|---|---|---|---|---|---|
| Guiyuan Street | 440303001 | Guiyuan Street | 440303001 | 4024 | 11.87 |
| Guiyuan Street | 440303001 | Huangbei Street | 440303002 | 1513 | 10.08 |
| Guiyuan Street | 440303001 | Dongmen Street | 440303003 | 1609 | 9.43 |
| Guiyuan Street | 440303001 | Cuizhu Street | 440303004 | 1345 | 16.67 |

*Figure 3.1 Distribution of Cross-city flows within the SDH area*

The data used for our study is at the sub-district level. In total, there are 8,921 pairs of Origin-destinations (OD) summarised from commuting trips. For each pair of ODs, the data records the original street ID, the destination street ID, the number of commuters, average commuting time and distance. Figure 3.1 shows the distribution of inter-city flows within the SDH area.

## 3.4 Methodology

This study adopted disaggregated spatial interaction model for simulating the flow of spatial interaction between zones. Furthermore, a Hierarchical Spatial Interaction Model (HSIM) is applied to reflect the boundary effect between cities. For detecting the urban functional zones, this research proposed a novel regionalisation algorithm that redraws non-overlap boundaries of urban functional zones by searching for the best partition with the best goodness of fitting in HSIM.

### 3.4.1 Basic spatial interaction model

In this study, we established a set of spatial interaction models using the singly constrained gravity model, which assumes the distribution of trips roughly follows the format of the negative-power function for predicting the flow between zones. The core spatial interaction model can be represented as the following Equation (13):

$$T_{ij} = O_i^{obs} \frac{D_j^{obs} c_{ij}^{-\beta}}{\sum_k D_k^{obs} c_{ik}^{-\beta}} \tag{13}$$

Where $O_i^{obs}$ is observed origins totals from zone $i$ and $D_j^{obs}$ refers observed destinations totals to zone $j$, $c_{ij}$ is the main travel time between origins and destinations, $\beta$ is a parameter related to the travel cost.

The basic framework of this model is a form of classic gravity models (Wilson, 1971). The calibrating processing is a parameter-free since the model picks the distance decay parameters $\beta$ by continually executing the iterations of standard non-linear optimised (Batty, 1976b; Batty & Milton, 2021) until the difference between the predicted mean trip cost $C$ and the observed

mean trip cost $C^{obs}$ is less than the pre-set threshold $\varepsilon$ (default is 5% for balancing the calculation time and accuracy) (Equation (14)-(15)).

$$| C^{pre} - C^{obs} | < \varepsilon \qquad (14)$$

Where:

$$C = \frac{\sum_i \sum_j T_{ij} c_{ij}}{\sum_i \sum_j T_{ij}} \qquad (15)$$

## 3.4.2 Hierarchical spatial interaction model

Although spatial heterogeneity exists within a mega-city region, most traditional (or "global") spatial interaction models assume the inner space of the modelling area is spatial isogeneity, which means all trips flow one general law. Thus, we further adopted a two-level hierarchical spatial interaction model for estimating the travel flow between zones (Figure 3.2). By applying this framework, it can eliminate spatial heterogeneity because of the boundary between sub-systems. It divides the global spatial interaction model into some intra-city interaction models. For each sub-model, the form is the same as the basic spatial interaction model introduced in section 3.4.1. Since the distance decay parameters β is an auto-fitted value that is different in each sub-system, each sub-system describes a distinguished travel pattern. The model can be written as Equation (16)-(18).

$$T_{total} = \sum T_{intra(n)} + T_{inter} \qquad (16)$$

$$T_{intra(n)} = \sum_{i \in n} \sum_{j \in n} O_i^{obs} \frac{D_j^{obs} c_{ij}^{-\beta_n}}{\sum_{k \in n} D_k^{obs} c_{ik}^{-\beta_n}} \qquad (17)$$

$$T_{inter} = \sum_{i \in n} \sum_{j \in m} O_i^{obs} \frac{D_j^{obs} c_{ij}^{-\beta_{inter}}}{\sum_k D_k^{obs} c_{ik}^{-\beta_{inter}}} \quad (n \neq m) \qquad (18)$$

Where $n$ and $m$ represents set of zones in different cities, $O_i^{obs}$ is observed origins totals from zone $i$ and $D_j^{obs}$ and $D_k^{obs}$ refers to observed destinations totals to zone $j$ and zone $k$.



*Figure 3.2 The Hierarchical Spatial Interaction Model: The total predicting trips equals to city-level intra-city models plus one inter-city model, Equation (5) = Equation (6) + Equation (7)*

As a key condition, the difference of boundary of cities can affect the performance of this model because of the boundary effect and spatial heterogeneity. If the cities' boundaries in this model coincide with the boundary effect, the overall performance of this model will improve since the processing of splitting has eliminated the spatial heterogeneity between sub-models. Therefore, we believe the goodness of fitting of this model can be an indicator for assessing

the reasonableness when drawing the functional boundary of cities. The detailed proof of this hypothesis will be described in section 3.5.1. Since all trip flows can be allocated to one of the sub-models, the sum of the total trip and the constrained factor (volume origin trips in our model) would keep constant by applying the HSIM framework without any loss of information.

### 3.4.3   Regionalisation algorithm for delineating urban functional zones

Based on the Hierarchical Spatial Interaction Model, we propose a novel regionalisation algorithm for delineating urban functional zones by searching the best partition with the best goodness of fitting in the HSIM. After determining a predefined number of regions (in this case study is three because there are three cities-level governments within this area), our iteration-based algorithm will run several times until the best partition which has the highest $R^2$. This algorithm design takes the conception of the tabu search algorithm. As an evolutional method of local search, it inherits the basic concepts of greedy algorithms that continually choose the optimal choice at each step to find the optimal solution to reduce complexity and time consumption. At the same time, it can avoid being trapped in local optima by adopting the "tabu list". This design of this algorithm referenced previous works of by Openshaw and Rao (1995).

*Figure 3.3 Flowchart of the redrawing boundaries tabu search algorithm*

The basic workflow for each iteration is:

*Step one* Testing to find reassigning which zones will improve the goodness of fitting of HISM and then update the boundary.

*Step two* When the algorithm finds that reassigning any zone not in the tabu list cannot improve the result anymore, the algorithm will test if reassigning zones currently in the tabu list can further improve the result, called the "aspiration move".

*Step three* If no further improvement or aspiration move can be made, the algorithm would reassign the zones with the best result even current assigning improving, then back to step one for starting a new iteration.

Every reassigned zone will be recorded into the "'tabu list", which not be considered in the following iterations. In addition, a stopping criterion has been set to avoid endless iterations. The iteration will be terminated if the best partition has not been updated after N (N=20 in this study) times iterations. We have conducted a sensitivity analysis for these two parameters, which determines R=11, which can maximise the goodness of fitting. The results keep the same when R is within the range 1~10. Then the result slightly improved and then kept the same when R equals 11 or continually increased until reaching the length of the candidate zone list. In addition, the sensitivity analysis finds that the algorithm is not sensitive to the value of N. Whether the N increase to 50 or 100, the result would not change.

Although theoretically, the candidate zones could be any zone within the region, we could customise a set of prioritised zones to improve the algorithm's efficiency. For instance, in our experience, we set the scope of search space to all zones whose intercity commuters are more than 1%, which matches the average ratio of inter-city commuters in the case study area. A flow chart of this algorithm for reassigning the boundaries can be found in Figure 3.3.

To provide appropriate decision support, this algorithm should not only assess the current functional zones but also predict the long-term situation. Thus, we designed two different settings that have a minor difference when we execute the algorithm. The first setting is based on the situation that current cities' core functional regions can only spill over to zones close to the administrative boundary due to the local authority's current land-use planning and management scope. Therefore, the proportion of inter-city trips in each zone would not change

further by updating functional zones during each step of iterations. In other words, the inter-city flows for each zone is static according to the administrative boundary.

The second setting is the inter-city flows for each zone is dynamically updated according to the current boundary in iteration processing. That is, inter-city flows may be re-classified as intra-city flow after iterations of a boundary. This setting could be used for predicting long-term scenarios. in which the cities' core functional zones can spill over freely without restriction by the current administrative boundary, forecasting the potential functional boundary in the long term.

## 3.5   Results

### 3.5.1   Goodness of fitting for HSIM

For this case study, we split the global spatial interaction model of the whole SDH area into four sub-SIM models: three intra-city trips models for Shenzhen, Dongguan, Huizhou respectively based on its original administrative boundaries, plus one model for only predicting the inter-city trip.

To verify the hypothesis, we raised before that if the goodness of fitting can be an indicator for reflecting the boundary effect of cities, we introduced a controlled group Since the spatial heterogeneity and the boundary effect are often more significant around the boundaries between cities, this controlled group is set as it still has the same trips and zoning system (172 sub-district level zones) but with randomly urban boundaries. The boundaries applied in models for shown in Figure 3.4.

*Figure 3.4 (a) The GSIM model with one whole modelling area (left); Figure 3.4(b) the HSIM model with random boundaries (middle); Figure 3.4(c) the*

*HSIM model with administrative boundaries (right)*

There are some flow trips produced by the spatial interaction model and hierarchical spatial interaction model (HSIM) introduced in section 3.4.1 and 3.4.2 according to the different boundaries. As introduced before, the goodness of fitting can be an indicator for assessing the reasonableness when drawing the functional boundary of cities. The estimated distance decay parameters in sub-models have been attached as Table 3.2.

*Table 3.2 The distance decay parameters in sub-models*

|            | Global-SIM | HSIM with RB | HSIM with AB |
|------------|------------|--------------|--------------|
| Area 1     | 2.5425     | 3.2006       | 2.4759       |
| Area 2     | N/A        | 2.3841       | 3.4837       |
| Area 3     | N/A        | 2.6334       | 2.0937       |
| Inter-city | N/A        | 1.7279       | 1.6927       |

*RB means Random Boundaries, and AB means administrative boundaries.

For assessing the goodness of fitting, we calculated the mean-square error (MSE), Mean absolute error (MAE), Root Mean Square Error (RMSE) and R-square ($R^2$) compared with the observed flow, the results are represented in Table 3.3 below.

*Table 3.3 Goodness of fitting for GSIM and HSIMs*

| | Global-SIM | HSIM with RB | HSIM with AB |
|---|---|---|---|
| MSE | 14,945,798 | 9,082,981 | 5,013,151 |
| MAE | 335.73 | 278.188 | 239.90 |
| RMSE | 3865.80 | 3013.79 | 2239.01 |
| $R^2$ | 0.4531 | 0.6564 | 0.8165 |

*RB means Random Boundaries, and AB means administrative boundaries.

As the statistical measures are shown in Table 3.3, compared with the traditional GSIM model, the $R^2$ for the HSIM model with Random Boundaries and administrative boundary sharply rise to 0.6564 and 0.8165 from 0.45310. Meanwhile, MSE, MAE, and RMSE decreased significantly, which shows that the HSIM model largely shortened the difference between the estimated and actual values. Thus, all statistical measures indicators prove that the modular spatial interaction can significantly improve the goodness of fitting from the traditional Global methods in regional-scale scenarios. This result indicating the broader effect can be partly represented by this random boundary.

By comparing the result of HSIM with arbitrary boundaries and the HSIM with administrative boundaries, all statistical measures indicators reveal that applying appropriate boundary that reflects the spatial heterogeneity in HSIM would significantly improve the model's performance. This finding suggests that the travel behaviours of people who belong to the same functional city may yield better performance in fitting the specific distribution of trips. In other

words, in the case of the spatial resolution and number of sub-models keeping constant, the goodness of fitting by HSIM can be an indicator for assessing the reasonability of functional boundaries of cities. This finding provides a solid theoretical reference for the algorithm that we will introduce in the next section.

### 3.5.2 Result for detection of urban functional zones in SDH area

Table 3.4 reports the statistical result for the models with different boundaries. Although the performance of the basic scenario with the administrative boundary is already good enough, the models of both settings with new boundaries still slightly outperform the basic model. Comparing the minor improvement of statistical measurements, the new boundary itself is more meaningful for assessing the urban functional integration. The estimated distance decay parameters in sub-models have been attached in Table 3.5.

*Table 3.4 Goodness of fitting for HSIM in different scenarios*

|  | Base Scenario (AB) | Setting 1 (CB) | Setting 2 (LB) |
|---|---|---|---|
| MSE | 5,013,151 | 4,842,008 | 3,969,357 |
| MAE | 239.90 | 242.71 | 225.90 |
| RMSE | 2239.01 | 2200.456 | 1992.32 |
| $R^2$ | 0.8165 | 0.8228 | 0.8548 |

*Table 3.5 Distance decays of sub-models*

|  | Base Scenario (AB) | Setting 1 (CB) | Setting 2 (LB) |
|---|---|---|---|
| Shenzhen | 2.4759 | 2.4594 | 2.4594 |
| Dongguan | 3.4837 | 3.5455 | 3.7646 |
| Huizhou | 2.0937 | 2.0552 | 2.3514 |
| Inter City | 1.6927 | 1.7429 | 1.7746 |

* AB means administrative boundaries; CB means current boundaries; and LB means long-term boundaries.

**3.5.2.1   Setting 1-Current functional boundary**



*Figure 3.5 Setting 1- Current functional boundary within SDH area*

Figure 3.5 shows the result of setting 1 (statistic inter-city flow), indicating the current functional boundary within the SDH area. This result suggests that the current administrative boundary explains the boundary effect of trip distribution well. The statistical measurements in Table 3.2 support it. Compared with the current administrative boundary in the base scenario, the $R^2$ and other statistical indicators improved very slightly as only a few zones changed their belonging. For example, the functional core of Dongguan city is in the west of its administrative boundary because of its good transport connection with Shenzhen and Guangzhou. Therefore, the only zone in Dongguan that should be re-assigned to Shenzhen is Fenggang, as it has been known as the 'sleep city' for workers in Shenzhen, which is a typical example for cross-cities

functional integration. Moreover, a few zones near the Dongguan-Huizhou boundary will be re-allocated to Huizhou from Dongguan because these zones are away from the city centre and lack commuting connection with the city centre. It might be the main reason why trips in these zones would better fit the trip distribution in Huizhou rather than Dongguan. Similarly, a few zones in east Shenzhen that have been re-assigned to Huizhou.

### 3.5.2.2  Setting 2- Predicted functional boundary within SDH area in long-term



*Figure 3.6 Setting 2- Predicted functional boundary within SDH area in long-term*

As for the dynamic inter-city flow setting, the result (shown in figure 3.6) predicts that the functional areas will have more reassigning between Huizhou and Dongguan. The re-assigned zones in Dongguan are mainly from the 'East industrial park', Songshan Lake, and Dalang.

Historically, East industrial park areas are the cluster of manufacturing industries but lack commuting connection with the city centre. The algorithm also finds the Songshan Lake area in the middle of Dongguan is reassigned. This region has been assigned because of its strong linkage with the 'East industrial park' area. Local governments have recently emphasised such connections in their planning report. Besides, the Dalang area may also be reassigned for long-term prediction. Unlike the Fenggang area, this area lacks road linkage directly with Shenzhen though it is physically close to Shenzhen. Thus, this area has been re-assigned to Huizhou following the reassigning of the Songshan Lake area, which is one of the main workplaces for residents in Dalang. These results show that there will be more potential interaction and functional integration opportunities between zones between Dongguan and Huizhou because of the chain reaction in long-term prediction.

### 3.5.3  Policy implications for city integration in SDH area

These empirical-based results can help the governments and planners to understand the spatial structure in mega city-region and support their urban integration policy. Previous studies have always focused more on the north-western part of Shenzhen and the south-western part of Dongguan since it has the most volume of the cross-boundary trip statistically. Our study argues that in the case of balanced bidirectional flows present, the functional boundary effect between the cities would not change obviously. Because trips in this area fit their original intra-city trip distributions, the high inter-city flow might be a natural consequence of the high population density and spatially relatively close to their original urban centres. In contrast, this study reveals that Fenggang and Shenzhen have a very high degree of functional integration, indicating that the urban function (e.g., housing or employment) are shared within these areas.

When considering such an integration between the two regions, the policymakers should pay more attention to amenities and public service for inter-city commuters.

For the long-term prediction, zones in mid-Dongguan should be given more attention. These areas are very 'sensitive' to any change of trips since fits in these areas do not fit the intra-city trip distribution of their original cities and are far away from city centres. Thus, our algorithm predicts a severe mismatch between functional zones and administrative boundaries could occur in these areas, even with tiny inter-city interactions. This result proves that transport linkages are vital for reshaping the urban functional zones in the long term because of the chain reaction of the previous reassigned zones. An example is the Dalang area. Though it is physically close to Shenzhen, it has been re-assigned to Huizhou because the road linkage with Songshan Lake is better than those with Shenzhen.

Overall, these empirical results imply that there will be more potential interaction and functional integration opportunities between zones between Dongguan and Huizhou in the future. Besides, policymakers should consider improving transport connectivity between the reassigned areas and Dongguan city centres to eliminate the boundary effect of city centres in trip distribution. Such measures would also avoid severe mismatches between functional zones and administrative regions, which may cause extra difficulty for management.

# 3.6 Methodological discussion

### 3.6.1 Spatial interaction methods vs Network-based methods

The network-based community detection method is the mainstreaming method employed for detecting the boundary of communities and functional spatial structure at the cities-level in previous studies. However, there is some limitation as well. The traditional network analysis and most community detection algorithms usually only consider the absolute value of flow volume for dividing the partitions but overlook the spatial factors like travel distance/cost (Liu et al., 2014; Yin et al., 2017).

Typical community detection algorithms (e.g., Louvain algorithm) are always trying to search a partition for maximising the ratio of intra-city flows in overall flows. However, because the percentage of inter-city trips is usually tiny (3% or less compared to intra trips) among all trips, thus when ignoring the spatial factors, the traditional community detection prefers to split space within the origin of administrative boundaries rather than break it. Similar to the phenomenon observed by previous research (Liu et al., 2014), the detected communities are precisely the same as the original administrative boundary when we use the Louvain algorithm and adjust the minimum resolution point to let the number of communities equals to the number of cities. According to the definition by OECD, the city or town whose 10% of the population exhibits cross-boundary commuting behaviour can be regarded as the satellite city of the mega-city. Thus, the traditional network analysis is not sensitive enough for cross-boundary commuting trips, which may fail to support planners and policymakers appropriately when discussing the cross-boundary integration of the functional region.

In contrast with the network analysis-based method, our proposed spatial interaction-based algorithm will more consider the distance decay effect when detecting the boundary effect reflecting spatial heterogeneity. Because zones close to the cities' boundary are usually spatially far from the city centre, our algorithm would be more sensitive to cross-boundary trips even with a relatively small volume.

Besides, another limitation of network-based methods is the difficulty of predicting the future situation. Almost all research applied community detection methods must base on the existing data of travel flow. If the data is unavailable, the estimation of flows would still rely on spatial interaction models (Wu et al., 2021). It will cause more deviation when switching between the multi-methods. Because of the strength for estimating the travel flow, the spatial interaction-based methods would have a special advantage for the prediction and simulation of future urban regions dynamically.

### 3.6.2   Methodological limitations

There are some limitations, and several directions can be further explored. First, different forms of spatial interaction models can be adopted for predicting the trip distribution. This method only employed the most widely used gravity model with a negative power functional form. Therefore, more conditions of the spatial interaction model, including the intervening opportunity and radiation models, can be discussed and employed in future work.

The second point is the "scaling issue". Additional experiments have been conducted to validate the model with more communities and different boundaries. One of the experiments attempted to extend the case study area to a border area, the Great Bay Area (GBA) in Pearl River Delta China, for nine cities with the same zoning system (sub-district level units). It

confirms algorithm still works appropriately for this extended area, but the algorithm would yield different results for local results in the SDH areas. The reason could be that the added areas and the additional trips will affect the existing results when applying the inter-city trip estimation models. The difference would be extended in long-term prediction due to the chain reaction. Thus, choosing the spatial extent needs to be associated with the specific research question and focus study area.

Lastly, our regionalised algorithm considers connections and flows between any pairs of units, not just neighbours. The spatial factors have mainly been reflected by travel time in this study. On the one hand, this is one of the advantages of emphasising mobility flows compared with other regionalisation algorithms. However, on the other hand, spatial adjacency is crucial in some cases (land-use planning, air pollution, etc.). Thus, spatial constraints on physical distance may need to be added to this algorithm to handle more situations.

## 3.7  Chapter conclusion

There are several contributions from this Chapter. First, this research confirms that the results of the hierarchical spatial interaction model (HSIM) can assess if the boundary of subsystems appropriately represents the inter-city boundary effect in trip distribution. Furthermore, this study proposes a novel method to delineate UFZs by searching for the best partitions in HSIM. By adopting the proposed model into a specific mega-city region, China, Shenzhen-Dongguan-Huizhou (SDH) area, this research confirmed the model's effectiveness in delineating UFZs based on spatial interaction from the perspective of human activity behaviour.

# 4 Exploring the Associations of Socioeconomic Characteristics and Distance Decay Effects in Spatial Interaction

## 4.1 Background

Understanding and predicting commuting behaviours are long-standing topics in urban analytics and numerous attempts have been made to use various quantitative models to achieve this (Barbosa et al., 2018; Lenormand et al., 2016; Schläpfer et al., 2021). The spatial Interaction (SI) model is one of the most powerful techniques for modelling and predicting flows. It forecasts the strength of spatial interaction based on the influence of distance decay, which means the interaction strength would be decreased along with the distance increasing. Most current spatial interaction models and other flow-predicting models at an aggregated-level assume that the interior space of the modelling region is spatially isogenous, meaning that the distribution of trips only obeys a general law associated with distance between locations (De Vries et al., 2009; Fotheringham & O'Kelly, 1989; Simini et al., 2021). However, previous research has evidenced that spatial heterogeneity widely exists in the spatial interaction model

103

and may reflect the border effect of trip distribution within urban space (Zhang et al., 2022). Researchers found that a global approach to spatial analysis may not be suitable for the local area within the sub-case study area due to spatial heterogeneity. (Fotheringham & Sachdeva, 2022). Thus, adding local characteristics to improve the accuracy of spatial interaction models is a long-standing research topic.

As introduced in Chapter 2.1, previous research considers the variation in the distance decays resulting from the spatial structure (Curry, 1972; Griffith & Jones, 1980; Oshan, 2020); thus, previous research has attempted to represent the local spatial attributes in spatial interaction models with different methods, such as the singly constrained gravity model, which is the most successful branch (Fotheringham, 1981; Nakaya, 2001; Oshan, 2016; Zhang & Li, 2024). However, some existing issues prevented the localised spatial interaction model from being applied to urban systems with granular spatial resolution, such as the excessive computing complexity and difficulty in calibrating parameters in low-flow or zero-flow areas due to data sparsity (Fotheringham & O'Kelly, 1989).

Recently, the evolution of computing power and data collection/storage techniques enabled the possibility of using spatial interaction models to build large-scale urban models with granular spatial resolutions, and this trend is known as "the renaissance of large-scale modelling" (Batty & Milton, 2021). However, a research gap for applying the localised spatial interaction model is that most localised spatial interaction models stay at the relatively macro level (e.g., province/state/ regional level), and do not go further into a finer spatial resolution within the urban systems (Dennett & Wilson, 2013). This is due to some associated issues that sometimes prevent the localised spatial model from being utilised in predicting flow within the urban system. As the spatial resolution becomes finer, the number of origin-destination pairs grows exponentially. This proliferation of data points increases the model's complexity and

computational demands. Each sub-model may have its own set of parameters, necessitating separate fitting processes. The need to fit and validate these sub-models further adds to the computational burden.

Another issue is the local calibration in the origin-specific gravity model may be invalid in low-flow volume areas within the urban system. In granular spatial resolutions, some areas may have low flow volumes recorded. Local calibration requires enough data to make reliable estimates. In the case of grouping the flow by its origin area, those areas with a lot of zeros flow could lead to imprecise and unstable parameter estimates (Fotheringham & O'Kelly, 1989). Practically, Poisson regression usually cannot provide reliable results and the non-linear iteration method could meet the issue of being unable to converge.

Besides, various studies have supported the idea that the spatial relationship between workplaces and residences is not the sole factor explaining the spatial heterogeneity observed in commuting patterns. Individual-level socioeconomic characteristics, individual preferences, and attitudes also play a role in influencing people's travel behaviours (Gao et al., 2024; S. Hu et al., 2022; Lin et al., 2015). While research on localised spatial interaction has predominantly focused on its relationship with spatial structure (Chen et al., 2019; Zhang et al., 2022), non-spatial factors like socioeconomic characteristics are normally not considered in the models. Reviewing previous studies, the research gap in integrating socioeconomic characteristics into spatial interaction models has been identified in Chapter 2.1.4. Therefore, this chapter will introduce a novel method to fill this research gap.

Thus, this chapter aims to introduce an algorithm by addressing the central question of how local socioeconomic characteristics can be integrated into a localised spatial interaction model for predicting commuting trips'. We proposed a two-step spatial interaction model framework

for achieving research objective 2 to tackle this. This framework quantitatively captures the variance of distance-decay effects in local commuting behaviours, substantiating the correlation between socioeconomic characteristics, urban spatial configuration, and spatial interaction. Furthermore, by integrating a clustering algorithm based on socioeconomic factors into the localised spatial interaction model, we demonstrate that the performance of the trip prediction model can be remarkably enhanced with only a limited parameter increase. This opens avenues to devise trip prediction models that more precisely align with the evolving socioeconomic scenarios of residents.

## 4.2  Methodology

### 4.2.1  Two-step localised distance-decay with origin-specific gravity model

This research applies an origin-specific gravity model (OSGM) to observe the variant of distance decay for commuting by different origins and destinations in London. Following the classic unconstrained (or sum-constrained) gravity model written as equation (2) above, this study adopted a disaggregated spatial interaction model referenced the previous research (Fotheringham & Brunsdon, 1999), which divides the flows by origins and then fits the flows with separate models in the formatting of the classic unconstrained gravity model (19). For giving a specific origin, the $O_i$ is part of the constant (20). Each sub-gravity model has its own distance-decay parameters calibrated by the general linear regression model.

$$t_i = \sum_j t_{ij} = O_i{}^{\alpha_i} \sum_j K \frac{D_j^{\gamma_i}}{d_{ij}^{\beta_i}} \tag{19}$$

$$T = \sum_i t_i \tag{20}$$

Where $t_{ij}$ is travel flow between zone $i$ and zone $j$, $O_i$ is observed origins totals from zone $i$ and $D_j$ refers to observed destinations totals to zone $j$, $d_{ij}$ is the main travel distance between origins and destinations, and $\beta$ is a parameter related to the distance decay.

To determine the distance decay parameters $\beta$, we employ the general linear model (Poisson regression) to calibrate the sum-constrained model (Dennett & Wilson, 2013) after the log formed transformation. (21)

$$t_{ij} = \exp\left(K + \gamma_i D_j - \beta_i ln d_{ij}\right) \tag{21}$$

The OSGM has a better prediction ability by aggregating the predicted flows as one predicting O-D matrix (6). It performs better in statistical measurements (e.g., R-square and Root Mean Squared Error) compared with a classic constrain gravity model fitted by general linear regression models, proving the reasonability for highlighting the local distance-decay parameters.

## 4.2.2  Grouping k-local groups and two-step flow prediction method

Equation 3 has a key parameter $\beta_i$ which controls the local distance decay effect in the flow prediction. From previous literature, predicting the local level of $\beta_i$ in the urban system has been identified as a significant research gap and, in this Chapter, we propose a novel approach in utilising socioeconomic characteristics to shape the local distance decay as a potential

solution, because similar social groups may have similar commuting behaviours. Thus, we design a two-step spatial interaction model introduced after the fitting process of the OSGM has determined the localised distance-decay parameter is described below:

Step One: Identifying areas in which specific social groups with a k-means clustering algorithm by applying residents' socioeconomic status.

Step Two: calibrating the localised distance-decay parameter in origin specific gravity model by the general linear model in equation (21) for predicting the travel flows.

Clustering, a versatile tool for identifying groups or clusters within multivariate datasets, has found extensive application across domains such as biology, psychology, and economics (Kodinariya & Makwana, 2013). This research uses the k-means method to partition our dataset into k distinct, non-overlapping clusters to predict local distance decay parameters. k-means is a prevalent clustering technique that groups data points based on their similarity in specific features. The Sum of Squared Errors (SSE) represents the aggregate of these squared distances across all data points, and the primary objective of the k-means algorithm is its minimisation. It operates with the assumption that data within a given cluster is more akin to each other than to data in other clusters. This similarity is quantified by calculating the squared distance of each point to its assigned centroid.

Determining the value of k is a key step for the k-means algorithm. Thus, we utilised a method similar to the current determining k elbow method which finds the optimal value for k by finding the local optimum point. This technique is rooted in the principle that the ideal number of clusters, k, is identified at the point where the reduction in the sum of squared distances between data points and their corresponding centroids begins to plateau or 'level off'. However, unlike the current method, which merely determines k from the clustering performance, our

method adds the goodness of fitting for the OSGM model as a second indicator to determine the value of k. The Root Mean Squared Error (RMSE) estimates how well the model can predict the target and actual values. RMSE is employed to find the optimised number of k to improve the goodness of fitting of the OSGM model.

After identifying the number of k and social group within the modelling space area, we apply the OSGM in fitting the distance decay as the localised distance-decay parameter as the proxy of distance-decay parameters for the whole group to reduce the number of parameters and complexity of the OSGM model. Thus, the model in equation (21) will be further transformed as (22) to let the local parameters of equal zones be represented by a set of parameters calibrated with other zones belonging to the same social group.

$$t_{ij} = \exp(K + \alpha_k D_j - \beta_k lnd_{ij}) \mid Zone\ i\ \in Group\ k \qquad (22)$$

## 4.3  Case study

### 4.3.1  Data collection

The case study focuses on the Great London Area in the United Kingdom and utilises research based on UK census data 2011[1]. The data sources include Origin-Destination (O-D) data, which

---

[1] Link of Data source: https://www.nomisweb.co.uk/sources/census_2011

tracks the location of usual residence and place of work, as well as the socioeconomic characteristics of local residents.

Meanwhile, socioeconomic data encompasses various demographic and economic indicators. The research covers a wide range of socioeconomic variables, such as gender ratio, mean age, one-person household ratio, house ownership rate, minority ethnic group rate, higher education rate, economic active rate, car ownership rate, and marriage rate (figure 4.1). The selection of variables is referenced previous research introduced in the literature review and subjected to the availability of census dataset. These indicators offer valuable information about the social and economic dynamics of the local population.



*Figure 4.1 Socioeconomic characteristics in London*

The spatial resolution employed in the study is based on the Middle Layer Super Output Area (MSOA) level, with 983 areas in total for London. For comparison the possible change from

2001 to 2011, we also download the 2001 census data. The data sample is shown at Table 4.1 below:

*Table 4.1 Data sample of Census data*

| Category | Class | Example |
|----------|-------|---------|
| Location | Area Name | City of London 001 |
| | MSOA Code | E02000001 |
| Education | No qualifications (%) | 6.7 |
| | Highest level of qualification: Level 1 qualifications (%) | 4.3 |
| | Highest level of qualification: Level 2 qualifications (%) | 6.6 |
| | Highest level of qualification: Apprenticeship (%) | 0.7 |
| | Highest level of qualification: Level 3 qualifications (%) | 7.2 |
| | Highest level of qualification: Level 4 qualifications and above (%) | 68.4 |
| | Highest level of qualification: Other qualifications (%) | 6.2 |
| | Schoolchildren and full-time students: Ages 16 to 17 (%) | 1 |
| | Schoolchildren and full-time students: Ages 18 and over (%) | 6.2 |
| Housing | Owned (%) | 42.3 |

| | | |
|---|---|---|
| | Shared ownership (part owned and part rented) (%) | 0.3 |
| | Social rented (%) | 16.5 |
| | Private rented (%) | 35.9 |
| | Living rent-free (%) | 5.0 |
| Car Ownership | No cars or vans in household (%) | 69.4 |
| | 1 cars or vans in household (%) | 25.1 |

## 4.3.2 Exploring spatial heterogeneity from distance decays in commuting behaviour

Utilising the Origin-Specific Gravity Model (OSGM) demonstrated a notable enhancement in model fitting compared to classic gravity models. Specifically, the r-square value improved markedly from 0.522 when using the sum-constrained spatial interaction model fitted by the general linear regression model to 0.815 with the OSGM. Concurrently, there was a significant reduction in the RMSE from 12.15 in the sum-constrained spatial interaction approach to 7.009 in the OSGM. This enhancement in performance suggests that incorporating local distance decay can considerably mitigate errors stemming from spatial heterogeneity. Consequently, local distance decay parameters emerge as pivotal variables, warranting further exploration.

Figures 4.2 provide an analysis of the local distance-decay parameters in London in 2011. In this figure, the lighter colours indicate smaller absolute values of the distance-decay parameter, suggesting smoother distance-decay effects for the corresponding areas. Conversely, the border

areas depicted in darker colours represent sharper distance-decay effects within these zones. The distribution of distance decay remained relatively stable during this relatively long-term period by comparing with the 2001 distance decay parameters, and the comparison results will be further discussed later.

The clustering of highlight values in the central area indicates that the distance-decay effects in these regions are more gradual and less influenced by distance. This smoother decay suggests that commuting patterns or other spatial interactions within these central areas are relatively consistent, regardless of the distance between locations. On the other hand, the border areas displaying darker colours signify stronger distance-decay effects. In these zones, the impact of distance on commuting or other spatial interactions is more pronounced, implying that the ease of interaction decreases significantly as one moves away from the central areas.



*Figure 4.2 Distance decay parameters in London in 2011*

To identify the spatial clustering pattern, Moran's I was utilised to measure spatial autocorrelation for the distance decay parameters (beta) for 2011's results. The Global Moran's I index yielded a value of 0.801, indicating a strong positive spatial autocorrelation. The associated p-value, which is less than 0.01, signifies that this result is statistically significant.

### 4.3.3 Paradox between local distance decay and spatial structure

Typically, the sharper distance-decay effect means the commuters living in these areas are more sensitive about travel distance and vice versa. Therefore, those people who are more sensitive to travel distances should live in the central areas. However, our results denied that statement, a proof of this is the distance decay effect shows a negative linear relationship with the travel distances. Figure 4.3 clearly illustrates the distance-decay effect is more significant as the absolute value of beta expanded along with the average travel distance increasing. Please note that this distance decay mainly indicates to what extent the travel frequency between two areas is decreased along with the distance increasing, however, the absolute value of travel frequency and average travel distance cannot be represented by the distance decay. The subplots in Figure 4.3 show the histogram for the number of travellers vs travel distance in two areas with similar average travel distances but different distance decay. The above subplot is in the relatively central area, Newham, which has a high travel frequency, but the distribution of travel frequency does not significantly decrease with the commuting distance before 15000m. The bottom subplot is from an area in Hillingdon, a border Borough in London. It could be spotted that the distance decay effect here is sharper.

*Figure 4.3 Relationship between average travel distance and the distance-decay effect*

A possible explanation is most people are commuting from their residence towards city centres, but not in the opposite direction. Still, most commuters would not go beyond the central area to another side of the city, that is why the previous research pointed out that travel frequency has sharply decayed along with the increasing of travel distance. This edge effect could be determined by the functional spatial organisation in London's case, as most job opportunities are mainly distributed in the central area rather than the border area. In contrast, commuters living in the central areas would freely choose their workplace without significant limitations of specific directions or areas. Thus, the distance-decay effect for those people shows smooth trends. Based on this explanation, socioeconomic characteristics, such as housing prices and affordability, could play a significant role in shaping distance-decay patterns in commuting behaviour. In other words, Socioeconomic factors can influence residential and employment

decisions and in turn, impact commuting behaviours and the distribution of distance-decay effects.



*Figure 4.4 Difference of local distance decay parameter between 2001-2011.*

Another Interesting discussion is about the temporal variation in the distance decay effect. We also use the same commuting dataset from the 2001 census. The areas within the Great London Area are developed, so the spatial structure generally remained stable within these 10 years. Even though the overall trend of distance decay is stable because the spatial structure is stable, the local change still exists. Figure 4.4 illustrates those certain areas in East London, such as Canning Town, Greenwich, and Woolwich, have experienced a notable reduction in the distance decay effect. These East London areas have also been identified as gentrification zones

since the 21st century. The reduction of distance decay clustering in these areas indicates changes in the commuting behaviours of residents as the results of gentrification.

Conversely, some areas in West London and central regions have experienced a slight increase in the distance decay effect. This could indicate that commuting patterns between these areas and other locations have become more influenced by distance over time. One possible explanation for this trend is the emergence of employment opportunities in East London, which could draw workers from the city's western and central parts. But compared with the distance-decay reduced area, these increased areas do not show spatial clustering in some specific areas. The two-year comparison result confirms that the local residents' socioeconomic characteristics will significantly affect the distance decay, and this effect is beyond the explanation range of pure spatial factors.

## 4.3.4 Grouping areas based on socioeconomic characteristics associations with localised distance decay

From the last section, the distance decay effect significantly correlated with spatial structure, which socioeconomic factors could drive. Based on this hypothesis, we designed an algorithm to predict the variation of the local distance decay effect with functional spatial structure reflected by residents' socioeconomic characteristics, predicting the more accurate travel flows

*Figure 4.5 The fitting figures to demine the K (both R-square reports and RMSE reports goodness of*

*fitting)*

For our case study data, we found k=2 and k=6 could be the optimum value of k (Figure 4.5). k=2 is a local optimum point since it is the peak for both indicators when k<5, and k=5 has been identified as another local optimum point because it then increasing slowly after k=6 (over-fitting region).

*Figure 4.6 K-means cluster results for the socioeconomic characteristics when K=2*

Figure 4.6 delineates the results from k-means clustering. Notably, even though the algorithm does not incorporate spatial data, the classification based on socioeconomic characteristics still reveals a discernible spatial clustering pattern. The results depict roughly two concentric circles from the centre outwards, evincing the London area's mono-centric spatial configuration. This result confirms that the clustering of socioeconomic characteristics can represent the spatial autocorrelation in the distance decay effect, though the dataset is non-spatial.

We refer to Group 1 as the "Inner London" group, indicated by the light-yellow colour. The areas boasting a mean distance decay value of -1.57, and primarily inhabit central London. They tend to exhibit a high education level and are economically active, with a higher level of one-person household and ethnic minority rate with lower car and house ownership rates.

Group 2 is the "Outer London" group, which manifests a more pronounced distance decay effect, marked by a beta value of -1.95. Distinctively, these zones have elevated percentages of car and house ownership, along with the oldest mean age and the fewest one-person households. These characteristics imply that Group 2 offers the most stability compared with the inner London Group.

*Table 4.2 The distance decay and socioeconomic characteristics for groups of areas when K=2*

| Group | Distance decay ($\beta$) | Car Ownership (%) | Higher Education (%) | Economic Active (%) | Ethnic Minority (%) | Housing Ownership (%) | One-person Household (%) |
|---|---|---|---|---|---|---|---|
| Inner London | -1.13 | 44.93 | 43.86 | 72.23 | 40.62 | 33.76 | 16.16 |
| Outer London | -1.71 | 72.82 | 31.9 | 70.96 | 38.39 | 62.58 | 10.01 |

When K=6, the grouping map (Figure 4.7) shows a different trend compared with the results when K=2. On the one hand, the grouping results are still clustered in the spatial aspect, reflecting that the grouping results still can reflect spatial factors. On the other hand, this clustering result can go beyond the continuous spatial restriction of the concept of 'boundary', the detailed socioeconomic characteristics and distance decay could be found in Table 4.3.



*Figure 4.7 K-means cluster results for the socioeconomic characteristics with K=6*

*Table 4.3 The distance decay and socioeconomic characteristics for groups of areas when K=6*

| Group | Name | Car Ownership (%) | Higher Education (%) | Economic Active (%) | Ethnic Minority (%) | Housing Ownership (%) | One-person Household (%) | Distance decay ($\beta$) |
|---|---|---|---|---|---|---|---|---|
| 1 | Inner-city professionals | 40.17 | 51.28 | 72.68 | 32.39 | 32.92 | 22.21 | -0.96 |
| 2 | Mixed-ethnicities city renters | 39.93 | 35.15 | 69.59 | 52.11 | 26.13 | 14.04 | -1.26 |
| 3 | Settled city achievers | 60.34 | 55.81 | 77.08 | 24.82 | 50.79 | 14.90 | -1.31 |
| 4 | Ethnic-minorities family | 66.27 | 30.54 | 68.45 | 69.74 | 55.27 | 7.22 | -1.68 |
| 5 | Working-class suburbs | 63.01 | 26.84 | 69.59 | 40.50 | 49.01 | 11.19 | -1.70 |
| 6 | City Fringe Homeowners | 81.01 | 32.38 | 72.47 | 20.57 | 76.42 | 10.53 | -1.85 |

**Inner-city professionals** (Group 1) areas are primarily located in the most central regions, with a few situated in sub-centres of employment, such as Canary Wharf and Central Croydon. The distance decay of travel patterns for these areas is the most gradual, with a value of -0.96. This group's profile displays relatively low percentages for car and housing ownership but has higher education percentages and moderate economic activity. Residents in these areas likely consist of younger professionals.

**Suburban renters** (Group 2) areas envelop the central regions, extending towards the northwest (for example, Harlesden), northeast (such as Hackney), east (like Whitechapel), and southeast (e.g., Old Kent Road). This group's beta value is -1.26, signifying a more pronounced distance decay effect than in Group 1 areas, although it remains moderate. This group appears to have a relatively higher percentage of ethnic minorities, coupled with the lowest rates of car and housing ownership and education. These areas might predominantly be inhabited by a diverse working-class population, possibly renting in pre-gentrification zones.

**Suburban achievers** (Group 3) areas span affluent southwest regions, like Fulham and Richmond, and northern areas like Hampstead. The distance decay effect here is slightly sharper than in Group 2 areas, with a beta value of -1.31. Boasting the highest Economic Active Percentage, these areas also have relatively high car and housing ownership rates. However, they have the second lowest ethnic minority percentage among all six groups. It's plausible that residents in Group 3 areas consist mainly of affluent locals and middle-aged professionals.

**Ethnic family life** (Group 4) and **blue-collar communities** (Group 5) exhibit similar distance decay levels, with beta values of -1.68 and -1.71, respectively. Areas such as Southall, Wembley, and Finsbury Park typify Group 4. These regions have the highest minority proportions and, concurrently, relatively high car and homeownership rates. Notably, they also have the lowest

rates of one-person households. This suggests that immigrant families with larger household sizes and stable economic backgrounds primarily inhabit Group 4 areas. On the other hand, Group 5 areas, represented by places like Dagenham and Uxbridge, display moderate percentages in most metrics. The notably low higher education percentage might indicate that a majority of the residents in these zones belong to the working-class segment rather than the professional class.

Lastly, **Outer London white families** (Group 6) areas are predominantly found on the outskirts of the Greater London Areas. These areas showcase the highest car and housing ownership percentages, coupled with the lowest figures for ethnic minorities and one-person households. This composition implies a dominant presence of stable, local white residents.

### 4.3.5 Goodness of fitting comparison

The next step is using this cluster results as the proxy of distance decay in the OSGM model. Table 4.4 compares the number of parameters and the performance among the four models. The Global gravity model encompasses just a single distance decay parameter. However, its performance leaves much to be desired, lagging in effectiveness. While the OSGM model stands out with superior performance to the other models, it demands a hefty 983 parameters to make its predictions. When K=2, this adaptation of the OSGM with clustering algorithm uses only two parameters as proxies for the localised distance decay parameters within OSGM. Even with very limited degrees of freedom, this algorithm can still significantly improve flow prediction accuracy by effectively representing London's monocentric spatial structure. When K=6, this algorithm uses six sets of parameters to represent the variation in distance decay in areas with each clustering group, promoting the goodness of fitting further and approaching the OSGM without clustering

*Table 4.4 Comparison results for different gravity models*

|  | **Global gravity model** | **OSGM with clustering (K=2)** | **OSGM with clustering (K=6)** | **OSGM** |
|---|---|---|---|---|
| **Number of parameters set** | 1 | 2 | 6 | 983 |
| **R-Square** | 0.522 | 0.655 | 0.722 | 0.815 |
| **RMSE** | 12.150 | 9.536 | 8.532 | 7.009 |

These results confirm that utilising socioeconomic characteristics to detect social groups within the urban system can effectively represent the variation of distance decay at the local level, and this representation could significantly improve the accuracy of prediction with spatial interaction.

## 4.4 Discussion for localising distance decay within urban space

After reviewing the previous research about SI and spatial factors in past research, this research proposed that socioeconomic factors should be considered to improve traditional gravity models, which consider merely spatial arrangement. Our results in traditional OSGM confirm that distance decay effects vary within the city space and long-term period, significantly affecting the flow prediction results and considering localised distance decay parameters could significantly improve the performance of the spatial interaction model. At the same time, this variation is beyond the traditional concept of 'boundaries' because the space shared similar distance decay is spatially discontinued. By applying the two-step gravity model, the localised

distance decay effect could be well-represented with limited parameters. One obvious advantage of this method is reducing the complexity of computing in the calibration process. In addition, the calibration issue in the low-flow area is mitigated because the subset flow has not been separated too granularly. In our case study, all sub-models are significant in the Poisson regression models, which means the calibration of local parameters is reliable.

Another significant pro of this algorithm is it can capture the heterogeneity in distance decay not only due to spatial factors but also non-spatial factors. Even though the socioeconomic dataset is non-spatial, the clustering results still show clustering in the spatial aspect, reflecting that spatial factors are still one of the main factors affecting distance decay variations. Meanwhile, this clustering result can go beyond the continuous spatial restriction of the concept of 'boundary.' In some areas, the different grouping areas are staggered. These results confirm that non-spatial factors, especially socioeconomic characteristics, influence distance decay variation. By appropriately representing spatial and non-spatial factors simultaneously, this model significantly improves the performance of the gravity models within the urban system.

In current spatial interaction models, there is a widespread assumption that the distance decay parameter remains consistent throughout the entire urban system and that its value will persist unchanged into the future (Batty & Milton, 2021; Lopane et al., 2023). One primary reason for adhering to this assumption has been the absence of effective methodologies to predict shifts in the distance decay parameter, especially at a localised level. This research bridges this gap by establishing a connection between socioeconomic status and localised distance decay effects. As a result, we can quantitatively evaluate how the distance decay effect could be predicted or policy scenarios might influence the distance decay effect, laying the groundwork for more robust and reliable predictive models in the future. Using socioeconomic status to group the area could help to solve the issues of the localised spatial interaction model. This method could

quickly respond to the scenario when the potential socioeconomic changes need to be considered (e.g., gentrification), providing more accurate prediction tools for travel flows in future scenarios, which will be introduced in next Chapter.

Last but not least, the optimum k number could vary depending on urban form and spatial structure. London has been widely recognised as a morphologically monocentric but more functionally polycentric region. The two local optimum values of k, k=2 and k=6, can be seen as reflecting the difference between morphological and functional in London, and this could also be observed from the visualisation of the k-means results. Thus, the factors that affect the optimum value K and an automatic algorithm to determine could be interesting topics in further research.

## 4.5  Chapter conclusion

To tackle the current research gap for lack of a granular-level spatial interaction framework to predict travel flows and consider behaviour variation due to local socioeconomic characteristics, this Chapter adopted a localised spatial interaction model to quantitatively assess how local commuting behaviours decayed with distance and how this is related to socioeconomic status and urban spatial structure. This Chapter also proposed a two-step algorithm for improving trip prediction models by incorporating a clustering algorithm based on socioeconomic factors to reflect the socioeconomic contexts of inhabitants better. With successfully defined social groups in two different settings, this prediction accuracy using spatial interaction is significantly improved with a limited increase in computing complexity. We also discussed some long-standing issues in tradition and disaggregated spatial interaction models in the later discussion part and proposed further research direction.

# 5 Application: Predicting the Impact of Changes in Transport Infrastructure on Urban Integration

During the last three decades, the Pearl River Delta (PRD) has been shaped by a variety of integration policies, including regional plans, industrial policies, and regional agendas. This approach emphasises the formation of city regions, which represent a novel spatial scale comprising clusters of cities and their surrounding hinterlands. City regions have emerged as significant geographical entities facilitating population growth, industrial advancement, and urbanisation (M. Chen et al., 2016; Zhang & Sun, 2019). These initiatives have aimed to restructure economic, social, and institutional intercity linkages within the region. Recently, the emergence of the Great Bay Area (GBA) of the Pearl River Delta emphasises the integration of cross-city space toward a super mega city-region as a key national development strategy.

As mentioned in Chapter 1.5.1, the GBA has a huge population size and geographical area. Under the master regional plan of the super mega-city region GBA, the government also defines some sub-mega-city regions to indicate the direction of urban integration in the relatively short-term future. Shenzhen Dongguan-Huizhou (SDH) area is one mega-city region around the core city of Shenzhen and includes two other prefecture-level municipalities, Dongguan and Huizhou. In 2012, the three cities' governments published the SDH Regional

The "Coordinated Development Master Plan (2012-2020)" promises to launch a series of urban integration policies, including highway and rail transit facilities, education and public health services, and encourage residents of the three cities to live as one city. However, local governments have different goals for their plan: the Shenzhen government hopes that sharing the urban functions with Dongguan and Huizhou can alleviate land resource shortages and "control the population's excessive growth" via promoting population immigration to other cities' directions. The governments of Dongguan and Huizhou hope to share in Shenzhen's economic growth by taking over Shenzhen's industrial and population transfer.

Notable, these regional plans are not static, the specific delineation of mega-city regions has been changed multiple times since the 1990s (Zhang et al., 2018). Zhongshan, as one of the 9 prefecture-level municipalities within the GBA, was a part of the SDH area in the master plan because of its location on the western coast of the bay separated by the Pearl River Estuary. Chinese government believes that large-scale infrastructure projects can effectively promote inter-city cooperation in the GBA area, and some previous qualitative research supported this from statement (Li et al., 2014; Xu & Yeh, 2013). The *Shenzhen-Zhongshan Bridge* is a key infrastructure development plan that connects Shenzhen and Zhongshan by crossing the Pearl River Estuary and is planned to open by the end of 2024 (Figure 5.1). The 24 km long bridge equipped with a motorway stand would be able to cut commuting time from Shenzhen to Zhongshan from the current 2 hours to about 30 minutes.    Based on this infrastructure construction, the local governments raised the "Shenzhen-Zhongshan same city" integration policy in 2019 (Government, 2019).

*Figure 5.1 Map of SDH area and Zhongshan City*

Predicting future urban structure and traffic flows in cities, especially by visually demonstrating the impact of urban intervention measures is critical for supporting relevant policy decisions about urban planning, infrastructure development and promoting sustainable growth. Urban simulation models offer several strengths, including their ability to simulate complex urban processes, explore alternative scenarios, and support evidence-based decision-making. In this chapter, a simulation model is established to forecast the potential impacts of specific urban interventions on urban spatial structure, utilising the urban analytic framework proposed in this doctoral research.

## 5.1 Data collection

Boeing (2017) introduced OSMnx, a Python package designed for downloading, analysing, and visualising street networks sourced from OpenStreetMap (OSM). With OSMnx, users can generate network graphs and execute network-based computation upon the generated graph. In this study, the road network is downloaded from the OpenStreetMap via the package OSMnx, and the download parameters have been set as "drive" to load the drivable road network open until Jan 2024 (Figure 5.2). The network does not include the proposed cross-sea bridge between Shenzhen and Zhongshan. Thus, the edge has been manually added to the network according to the polished road planning in order to calculate the future scenarios' travel costs.



*Figure 5.2 Road network downloaded from OSMnx*

The second dataset originates from Baidu, a prominent internet search and navigation service provider in China. Baidu provides a range of widely used internet services and mobile applications. Their Location-Based Service (LBS) automatically collects and combines human movement data from different applications with fine granularity, encompassing all populations. Within the SDH region, Baidu's data covers approximately 38,270,000 residents distributed across 15,970,000 working and residential sites. The dataset used in this chapter is collected and analysed based on the one-month period in November 2020. This dataset provides a more comprehensive source for investigating human mobility within the mega-city region compared to traditional surveys.

Moreover, Baidu's dataset includes user portrait data associated with travel information, which is collected by area simultaneously. This additional information allows for the analysis of potential connections between residents' mobility patterns and social characteristics. Although they can only provide aggregated data by the percentages of each attribute due to the privacy issue with mobility data, similar to the London census data we introduced in Chapter 4.3.1, it would still fulfil the data requirement for establishing the simulation model. Table 5.1 below lists the columns of the user portrait data and the data sample.

*Table 5.1 Data sample of user portrait data*

| Category | Class | Example |
|---|---|---|
| Location | Area Name | 松岗街道 (Songgang Street) |
| | Area Number | 440303001 |
| Population | Sum_all | 16547 |
| | male | 12177 |
| | female | 4370 |
| age | <18 | 562 |
| | 18-24 | 4601 |
| | 25-34 | 7338 |
| | 35-44 | 3525 |
| | 45-54 | 517 |
| | 55-64 | 4 |
| | >=65 | 0 |
| Education | High school or below | 10987 |
| | College | 3996 |
| | Bachelor or above | 595 |

| Category | Class | Example |
|---|---|---|
| Income | <=2499 | 527 |
| | 2500-3999 | 7378 |
| | 4000-7999 | 4379 |
| | >=8000 | 2211 |
| Car Ownership (Private Car) | >=1 | 5578 |
| | 0 | 10969 |

## 5.2  Design of the simulation framework

In this chapter, a simulation workflow is established to forecast the potential impacts of specific urban interventions on urban spatial structure, utilizing the urban analytic framework proposed in this research. Building upon the two-step spatial interaction model introduced in Chapter 3, the simulation model tests various policy assumptions, including population growth and migration, development of transport facilities, and changes in socioeconomic characteristics. These assumptions are incorporated into the simulation model to predict changes in travel behaviour within the urban area. Furthermore, the simulation model evaluates the potential alterations in functional spatial organization, drawing on the regionalization algorithm proposed in Chapter 4. This algorithm plays a crucial role in delineating urban functional zones and enables a comprehensive examination of how specific intervention may influence the urban spatial structure within the SDH regions.  The simulation model is designed as Figure 5.3 below:

*Figure 5.3 Design of integrated simulation model*

- Step1: Generating the travel cost based on the current road map.

Based on the road network download by OSMnx, we applied the Dijkstra's shortest path search algorithm to search the shortest distance for each pair of zones (Hagberg & Conway, 2020), the parameter of weight has been set as the travel time, which could highlight the difference in travel cost after the new Shenzhen-Zhongshan bridge opened.

- Step 2 Using socioeconomic data to cluster the Shenzhen-Dongguan Huizhou three cities:

Like the algorithm introduced in chapter 4, we employ the k-means method to partition Shenzhen-Dongguan-Huizhou (note the new area, Zhongshan, is not included) into k distinct, non-overlapping clusters. Determining the value of k depends on the goodness of fitting for the OSGM model, which calculates the Root Mean Squared Error (RMSE) of the delineated OSGM model to find the optimised number of k to improve the goodness of fitting, the detailed description of this method could refer to the section 4.2.2.

- Step3 predict the localised distance decay in the Shenzhen-Dongguan-Huizhou (SDH)

After determining the optimal number of clusters (k) and identifying the social groups within the modelling space, the Origin-Specific Gravity Model (OSGM) is utilised to fit the distance decay. This involves using the localised distance-decay parameter as a proxy for the entire group's distance-decay parameters.

- Step 4: Using same classifier to group areas in Zhongshan to predict the localised distance-decay parameters.

In this step, an assumption has been made here is once Zhongshan joins the SDH area, the residents with similar socioeconomic characteristics will share similar travel behaviour with existing residents in the SDH area. Thus, we use the same classifier trained in the last step to fit areas in Zhongshan to the existing social group, estimating the localised distance-decay parameters.

- Step 5: Updating the population, road network and travel cost.

For updating the travel cost after applying the intervention, the new Shenzhen-Zhongshan bridge is added to the existing road network segmentally based on the transport plan (LTD, 2014). The travel time has been calculated according to its designed speed (100km/h). After adding the new edges, the Dijkstra's shortest search re-runs to derive the updated travel cost.

- Step 6: Generating the "synthetic" travel flows.

With the localised distance-decay parameters in SDHZ areas and updated travel flow, we could generate the forecast travel flow after the intervention applied. The synthetic travel flows are generated from the travel flows are predicted based on the origin-specific gravity model introduced in section 4.2.1.

$$t_i = \sum_j t_{ij} = O_i{}^{\alpha_i} \sum_j K \frac{D_j^{\gamma_i}}{d_{ij}^{\beta_i}} \tag{22}$$

$$T = \sum_i t_i \tag{23}$$

Where $t_{ij}$ is travel flow between zone $i$ and zone $j$, $O_i$ is observed origins totals from zone $i$ and $D_j$ refers to observed destinations totals to zone $j$, $d_{ij}$ is the main travel distance between origins and destinations, and $\beta$ is a parameter related to the distance decay.

To The Oi and Dj could be further adjust based on the specific policy and scenarios assumptions.

- Step 7: Regionalisation Algorithm for Delineating the Urban Functional zones (UFZs) for the Shenzhen-Dongguan-Huizhou-Zhongshan area.

Once we get the "synthetic" travel flows from the last step, the final step of this urban simulation model is applying the novel regionalisation algorithm introduced in section 3.4.3 for delineating urban functional zones by searching for the best partition with the best goodness of fitting the spatial interaction models.

In this chapter, we adopt the second setting which allows the inter-city flows for each zone are dynamically updating according to the current boundaries during iteration processing. This

means that inter-city flows may be reclassified as intra-city flows following the iterations of boundary adjustments.

The tabu search algorithm starts from the current administrative boundary of the four cities and the output results could indicate that future UFZs within this mega-city region consist of the four cities.

## 5.3  Simulation results

In this section, we will first present the urban simulation model baseline for the year 2018 in the SDHZ area. This baseline model will provide an overview of the current situation in the area before any urban interventions are applied. The baseline model primarily involves steps 1-4 of the simulation model framework.

Following the presentation of the baseline model, we will introduce the updated simulation results after interventions are applied. These results will demonstrate the impact of the interventions on the urban spatial structure and will reflect changes in travel behaviour, socioeconomic characteristics, and functional spatial organization within the SDHZ area.

## 5.3.1 Baseline in 2020

Figure 5.4 provides an analysis of the local distance-decay parameters in the SDH area for the year 2020. The areas shaded in lighter blue represent smaller absolute values of the distance-decay parameter, suggesting a gentler decline in interaction as distance increases within those regions. Conversely, the regions depicted in darker colours indicate a steeper distance-decay effect, demonstrating a more rapid decrease in interaction as distance from these zones increases. This pattern mirrors the spatial distribution observed in London (section 4.3.2), where highlighted values tend to cluster in city centre areas. Regions displaying darker colours are primarily located at the borders or are distant from city centres, signifying stronger distance-decay effects.



*Figure 5.4 The local distance-decay parameters in the SDH area*

In Chapter 4, the methodology and results revealed a significant correlation between the distance decay effect and the spatial structure, which socioeconomic factors could influence. We adapted and applied a similar model to the SDH area as one of the key steps in the simulation model. Specifically, in our case study of the SDH region, we identified that k=6 is the optimal value for maximising the fit of the OSGM. This optimum value effectively illustrates the urban spatial structure, particularly when linked with the local socioeconomic characteristics. The detailed results and explanations are below, alongside Figure 5.5 and Table 5.2.



*Figure 5.5 The k-means cluster results for the socioeconomic characteristics in SDH areas*

*Table 5.2 The distance decay and socioeconomic characteristics for groups of areas in SDH areas*

| Group | Group Name | Median Age | Average Income (RMB/Month) | Higher Education Rate (%) | Car Ownership (%) | Travel time (seconds) | Beta |
|---|---|---|---|---|---|---|---|
| 1 | Central Business District elites | 36.78 | 12535.40 | 0.48 | 0.19 | 1309 | -3.24 |
| 2 | Urban core areas achievers | 34.88 | 10695.79 | 0.41 | 0.15 | 1381 | -3.46 |
| 3 | Shenzhen commuters | 32.35 | 9008.55 | 0.31 | 0.10 | 1529 | -3.75 |
| 4 | Suburban industrial workers | 34.70 | 7456.71 | 0.26 | 0.12 | 1596 | -3.39 |
| 5 | Outer suburbs residents | 32.64 | 6081.39 | 0.15 | 0.09 | 1774 | -4.12 |
| 6 | Non-urban rural areas | 34.16 | 4902.50 | 0.12 | 0.07 | 2377 | -3.94 |

**Central Business District elites** (group 1) are only located in the most central regions in Shenzhen. They have the highest median age, average income, and higher education rate. They also have the highest rate of car ownership and the shortest travel time. The distance decay factor is lower compared to other groups, indicating less sensitivity to distance or better accessibility to destinations (likely due to higher income and car ownership enabling more efficient commuting options). This suggests that wealthier, more educated individuals have

better access to resources and can afford to live closer to work or have more efficient commuting mean.

**Urban core areas achievers (**group2**)** are the second oldest group which has a slightly lower income, car ownership, and higher education rate compared to Group 1 but much higher than other groups. This area is located in Shenzhen's relatively central areas, as well as the most central single area for Dongguan and Huizhou. The distance decay is slightly higher compared to Group 1, this group's travel time is still short, maintaining relatively efficient travel but begins to show signs of reduced mobility compared to the wealthiest group.

Group 3 is mainly for the **Shenzhen commuters** who live in the area surrounding the city's core areas. The only exception is Songshan Lake district in Dongguan, which is widely considered Shenzhen's satellite area. This group is the youngest, with significantly lower income and education levels than the first two. Their car ownership is significantly lower, and their travel time has significantly increased. The distance decay is also higher than the two previous groups, indicating a much greater sensitivity to distance, possibly due to the inability to afford living closer to central areas or workplaces. However, the income of this group is still higher than any other group below, indicating the unique privileges of Shenzhen's economics within the SDH area.

**Suburban industrial workers** (group 4) see a slight increase in median age but continues the trend of lower income and education rates compared to Shenzhen commuters. These areas, mainly in Dongguan and Huizhou, surround the cities' central areas. Their car ownership slightly increases compared to Shenzhen commuters, but travel time also increases. This fact may mean they could have a mix of travel modes, making their sensitivity to travel costs moderate.

**Outer suburbs residents (**group 5) are younger, with the second-to-lowest income and education. They also have low car ownership and the second-highest travel times. This group likely lives in Dongguan and Huizhou areas, which are far from employment centres (or referring to the word "exurban"). The highest distance decay value indicates the most considerable distance sensitivity, reflecting this group's substantial barriers to efficient mobility.

**Non-urban rural areas** (Group 6) only appear in Huizhou. Unlike Shenzhen and Dongguan, whose urbanisation rates are almost 100%, Huizhou has a lot of rural and mountain areas. The groups 6 are similar in age to Groups 2 and 4 but has the lowest income and education rates of all groups. They have the lowest car ownership, the longest travel times and the second-highest distance decay, indicating their mobility challenges. Considering the regional background, and more than 80% of travel flows are local flows within the same area, these areas are not a part of the city-region areas but are more like rural areas.

*Figure 5.6 The k-means cluster results for the socioeconomic characteristics in SDH plus Zhongshan*

*areas*

The subsequent phase in the simulation model applies the same classifier to categorise areas within Zhongshan to predict localized distance-decay parameters. In this model, the city centres of Zhongshan are classified into Group 2 (Urban core areas achievers), aligning with similar categorizations in Dongguan and Huizhou. Surrounding the city centre, the areas are segmented into Suburban industrial workers (Group 4) and Outer suburbs residents (Group 5); specifically, the Northern and Western areas fall into Group 4, while the Southern and Eastern areas are categorized into Group 5. Additionally, the northern region adjacent to Foshan is incorporated into the same group. The outer suburbs and coastal regions of Zhongshan are designated as non-urban rural areas.

## 5.3.2 Predilect the future of mega-city region in 2035

As discussed in Section 5.1, Shenzhen is aiming to regulate its population growth, leading to a redistribution of migration to neighbouring cities such as Huizhou and Zhongshan. Dongguan tends to promote a mild increase in their population as well. In alignment with the official spatial master plan covering the period from 2020 to 2035, the projected data indicates that Huizhou and Zhongshan will undergo considerable population increases by the year 2035, with each city expected to see growth rates exceeding 40% as the Table 5.3 listed. Conversely, Shenzhen and Dongguan are anticipated to experience more moderate increases in population, with projected growth rates of 8.200% and 11.975% respectively.

*Table 5.3 The population growth targets set for the four cities*

| City | Official Documents | Current Residents in 2020 (Million) | Planned residence in 2035 (Million) | Increase rate (%) |
|------|-------------------|------------------------------------|-------------------------------------|-------------------|
| Shenzhen | Outline of Territorial Spatial Planning, Shenzhen (2020-2035) | 17.560 | 19.000 | 8.200 |
| Dongguan | Dongguan City Population Development Plan (2020-2035) | 9.645 | 10.800 | 11.975 |
| Huizhou | Huizhou City Land and Space Master Plan (2021-2035) | 6.042 | 8.500 | 40.682 |
| Zhongshan | Zhongshan Population Development Plan (2020-2035) | 4.418 | 6.200 | 40.335 |

Following the adjustment based on Zhongshan City's master plan for developing the Cuiheng City New Centre, significant changes are anticipated for the Nanlang area. This development aims to transform the region into a hub for high-end manufacturing and a base for cultural and technological cooperation. The area, expected to house over 200,000 'high-end talent' residents, will undergo considerable demographic and socioeconomic changes. As a result of these planned changes, Nanlang's population and socioeconomic status will be significantly revised upwards, transitioning from Non-urban rural areas to Shenzhen commuters to reflect the improvement in the area's socioeconomic characteristics.

Thus, we adjust the origins and destinations (denoted as Oi and Dj) for each zone in accordance with the rate of population increase in each respective city. This adjustment is based on the assumption that the number of job opportunities, which is more closely related to Dj, will grow at the same rate as the population increases.

### 5.3.2.1    Applying the urban intervention

The edge has been manually incorporated into the network by referencing the construction plan of the Shenzhen-Zhongshan Bridge to facilitate the calculation of travel costs for future scenarios. As Figure 5.7 shows, At Zhongshan side, the first exit of the road is at Ma'an Island, the second exit is the Licun and the third exit is Bo'ai Road. At Shenzhen side, the only exit is Bao'an Airport.



*Figure 5.7 The exits of Shenzhen-Zhongshan bridges (image source: Baidu map)*

The parameters for the new edges have been set according to the master plan, with the travel time calculated based on the designed speed of 120 km/h, the detailed number could be found in the Table 5.4.

*Table 5.4 Parameters of new edge*

| Road | Entry | Exit | Distance (km) | Travel time (seconds) |
|------|-------|------|---------------|------------------------|
| 1 | Bo'ai Road | Licun | 10.48 | 377 |
| 2 | Licun | Ma'an Island | 3.42 | 123 |
| 3 | Ma'an Island | Bao'an Airport. | 24.13 | 868 |

After incorporating the new edge, the nearest path has been recalculated using Dijkstra's shortest path search algorithm to update the travel cost. This updated travel cost demonstrates that the distance and travel time from Zhongshan to Shenzhen have significantly decreased. Figure 5.8 illustrates an example of the shortest road from Nanlang to Bao'an, comparing the scenarios before and after the road opening.

Distance=89.14km
Travel time=4512s

Distance=35.57km
Travel time=2018s

*Figure 5.8 The shortest road from Nanlang to Bao'an, comparing the scenarios before (left))and*

*after(right) the road opening*

### 5.3.2.2  Generating the synthetic travel flows

The synthetic travel flows are designed to capture the dynamics of human mobility patterns following the application of an intervention. These travel flows are predicted based on the origin-specific gravity model (OSGM) with updated travel costs. Given that the OSGM operates on the principle of being specific to the origin, it can be considered analogous to an origin-constrained gravity model. Therefore, the description and analysis of the travel flows will be discussed based on their origin.

*Figure 5.9 Change of average travel time by origin in SDHZ.*

Figure 5.9 depicts the changes in average travel time originating from different areas within the four cities. In this figure, areas marked in red indicate an increase in average travel time from those locations, while areas in blue denote a decrease in average travel time. Specifically, the city centre areas around Shenzhen, Huizhou, and Zhongshan exhibit reduced average travel times. In contrast, most areas in Dongguan and the peripheral regions of Shenzhen, Zhongshan show increasing trends in travel time. Notably, Huizhou shows a different trend for their edge areas; the north-eastern areas of Huizhou exhibit a significant decrease in travel time, but other peripheral areas show an increase in travel time.

### 5.3.2.3   Output of the regionalisation algorithm

After getting the updated travel flow, the last output of the simulation model is outputting the results of delineating the urban functional zones in the SDHZ area. For comparison impact

between before/after the urban intervention applied, the regionalisation algorithm runs two time based on the different travel flows. The results are list below:

Figure 5.10 illustrates the results from setting 1, focusing on statistical inter-city flow, which reveals the current functional boundaries within the SDH area. Since the synthetic travel flows are generated based on distance decay laws, the original statistics for goodness of fit are already quite high, with an R-squared value around 0.90 for the administrative boundaries. However, some zones still have noticeable shifts that will boost the new R-square to 0.916, indicating noncoincidence between administrative and functional boundaries.

The most significant observation from the figure is the expansion of Huizhou's functional zones, which extend into several zones in the northwest of Shenzhen. Additionally, some areas in the eastern part of Dongguan have been reassigned to Huizhou. This reassignment likely reflects the impact of population growth control policies. While Shenzhen and Dongguan are set to minimise their population growth until 2035 nearly, Huizhou is expected to see a substantial increase in its population, enhancing its role within the SDH area. Additionally, some areas located near the junction of the administrative boundaries of the three cities, which originally belonged to Dongguan and Huizhou, will be reassigned to Shenzhen. These areas are far from the city centre of Dongguan/Huizhou but have convenient motorway access to Shenzhen. Furthermore, the northern area of Shenzhen has been reallocated to Dongguan, possibly due to deeper integration within the urban space.

Despite expectations for significant population growth in Zhongshan by 2035, the Zhongshan areas remain separate from the SDH (figure 5.10a) area due to the lack of a cross-sea bridge linking it to Shenzhen and the travel time from west coast to east coast is still very high.

After the urban intervention of the new road applied, we re-executed the flow generating algorithm and regionalisation algorithm for the four cities involved. The outcomes, depicted in Figure 5.10b, show a significant alteration in the spatial dynamics of the region. Notably, this modification has shifted the affiliation of the Nanlang area from being part of Zhongshan to now being associated with Shenzhen, highlighting the impact of the new interventions on regional boundaries of urban functional zones.

*Figure 5.10(a) and Figure 5.10 (b)The delineation of unfunctional zones in the SDHZ area before intervention was applying (left); after the intervention*

*applying (right)*

## 5.4 Discussion and policy implication

### 5.4.1 Driving factor of the urban functional zones change



*Figure 5.11 Simulation results without population growth*

The first driving factor of the simulation results is population change. Population growth within the SDH area is unsynchronised for policy and planning reasons. Figure 5.11 shows the simulation results without the assumed population growth. The functional zones of Shenzhen are extended to both the west wing and east wing of Dongguan. Meanwhile, Huizhou's functional zones are generally the same as its administrative boundaries. Comparing Figure

5.11 and Figure 5.10 (a), in which Huizhou extended its urban functional zones to Dongguan and Shenzhen, it confirms that the population growth control policies could be one of the most significant driving factors for the simulation results.

Besides, the transport infrastructure is another driving factor for the simulation model. As illustrated in Chapter 5.4.3, applying urban intervention through the new road will promote the change of regional boundaries of urban functional zones. However, unlike population change, new transport infrastructure only has local influence on those areas that are directly linked to it.

## 5.4.2  The debate of good commuting pattern

The simulation results present an intriguing scenario where the average travel time increases after the opening of the proposed new road. This outcome appears paradoxical, as it is commonly believed that new transportation infrastructure would enhance the mobility of residents, thereby reducing travel time. However, our simulation contradicts this assumption, with Figure 10 clearly demonstrating that travel times have actually increased, particularly for areas that were expected to benefit significantly from these transport facilities. The rationale for this unexpected result is that the new bridge brings new job opportunities within the accessible travel cost to the Zhongshan zones, encouraging residents in Zhongshan to commute long distances.

*Figure 5.12 The increased average travel distance after the new bridge open*

The unexpected increase in travel time could be attributed to the fact that the new road infrastructure, such as a bridge, may lead to new job opportunities becoming accessible within a reasonable travel cost for Zhongshan zones. This, in turn, encourages residents to commute longer distances, possibly due to better employment prospects or higher wages, leading to an overall increase in average travel times.

This result could link to a long-term discussion of "what constitutes good commuting." It raises questions about whether the primary goal should be to minimise travel time or increase accessible job opportunities. Practically, this debate could relate to the specific goals that planners want to achieve. On the one hand, trends to reduce commuting time and create walkable commuting, such as the concept of a "15-minute city," have been listed as priorities for some governments. On the other hand, the concept of the "(super) mega city-region" and

"urban integration" encourages people to share industry cooperation and job opportunities in the larger area, boosting economic growth.

The case of the new bridge in Zhongshan, as per the master plan, illustrates this tension between facilitating shorter commutes and fostering economic growth through improved regional connectivity. Finding the right balance between reducing travel time and expanding job opportunities becomes a central challenge for urban planners. This balance will significantly impact how transport infrastructure development priorities are reassessed, highlighting the need for a holistic approach that considers the diverse impacts of new transportation projects on urban areas.

## 5.5  Chapter conclusion

In this chapter, a novel simulation model has been established to understand how specific interventions affect the human mobility pattern and ultimately influence the spatial structure of the Shenzhen-Dongguan-Huizhou (SDH) area after the Zhongshan joined this city region in the future. This model echoes the last research objective of this PhD research, which connects the human mobility pattern, socioeconomic characteristics, and spatial structure, supporting the related decisions by the government and planners to predict future scenarios with specific urban interventions.

# 6 Discussion

## 6.1 Improving the modelling framework for understanding urban spatial structure in mega city-regions

The relationship between spatial interaction and spatial structure has been debated for over 50 years without a clear conclusion (Griffith, 2007; Griffith & Jones, 1980; Oshan, 2020). A key step of the model framework proposed by this doctoral research is the investigation of localised variations in distance decay within spatial interactions through the calibration of an origin-specific gravity model. By employing two case study areas, the Greater London Area (GLA) and the Shenzhen-Dongguan-Huizhou (SDH) regions, this research elucidates the difference in describing the urban spatial structure of single-city metropolitan areas and mega-city regions from the perspective of distance decays in spatial interaction. This thesis suggests significant spatial variability in the distance decay component of the spatial interaction model, and such distance decay effect can reveal the functional urban spatial structure.

## 6.1.1 Incorporating socioeconomic factors when modelling travel pattern and urban spatial structure

We examine the spatial distribution of the localised distance decay effect concerning commuting behaviours within the GLA in Chapter 4. The associated Figure 6.1(a) illustrates that in the central regions of London, the commuting intensity decays with distance—a phenomenon known as the 'distance decay effect'— which occurs more gradually compared to the more pronounced decay observed in the peripheral areas. One hypothesis attributes this pattern to the significant edge effect influencing distance decay distribution due to the distance from the city centre. This suggests that the observed spatial patterns in local distance-decay parameter estimates arise from varying distributions of distances between each origin and all destinations. In other words, the distance from the city centre plays a crucial role in shaping the decay of spatial interactions, with areas closer to the centre exhibiting different interaction patterns than those farther away (Clark, 1951). This framework has garnered acknowledgement in early studies, particularly those focusing on morphological spatial structures (Johnston, 1973). The conventional approach to modelling this phenomenon involves quantifying the decay of commuting flows as a function of the distance from the 'city centre' (Halás et al., 2014).

*Figure 6.1 (a) and 6.1 (b) The localised distance decay in GLA area (left) and SDH area (right)*

This theory appears to appropriately explain dynamics within a single-city area, as the functional centre is typically located within the central areas inside the boundary—this is corroborated by our Greater London Area (GLA) case study. However, when the scale expanded to encompass a polycentric mega-city region, the application of this theory became more complex. Multiple functional centres exist in such regions, each with distinct functional roles contributing to shaping the urban structure, and the functional centres are not located in the geometric centre of the mega city-region.

The Shenzhen-Dongguan-Huizhou (SDH) region is an illustrative case where Shenzhen's functional centre, which also is the most crucial functional centre in the SDH area, is positioned at the southwest corner of the SDH area. In this instance, the most pronounced distance decay effects are not observed near the geometric centre of the area as they are more gradual in regions distant from this geometric midpoint yet nearer to the functional centres at the southern end of the entire region. This pattern is depicted in Figure 6.1(b), emphasising the importance of

functional centres over geometric centrality in influencing spatial dynamics and distance decay within the mega-city region. The comparative findings suggest that while distance contributes to the distance decay effect, it is not the sole determinant influencing the intensity of interactions within urban systems. This observation aligns with conclusions drawn from recent research, indicating that other non-spatial elements also play critical roles (Park et al., 2021; Šveda & Madajová, 2023). Thus, when modelling human mobility in the mega-city area, the functional urban spatial structure we should consider more complicated effects interplayed rather than single distance-related effects.

In this thesis, we propose a novel model framework that uses residents' socioeconomics characteristics to capture the local variations of the distance decay. Previous research has confirmed the link between travel patterns, especially commuting patterns and socioeconomic characteristics (Gao et al., 2024; Shen & Batty, 2019). Park et al. (2021) raised a hypothesis that the localised distance decay effect in travel behaviour could be associated with the socioeconomic characteristics of residents after calibrating distance decay parameters of travel flow in London. Our research validates the hypothesis that socioeconomic characteristics can effectively represent local variations in distance decay. The improved goodness of fit provides evidence in our models, which indicates a more accurate reflection of real-world travel pattern behaviours. Another advantage of using residents' socioeconomic characteristics is that they do not separate from spatial factors. As in both case studies shown in Figure 6.1, the clustering results represent significant spatial autocorrelation, even though our k-means algorithm does not have any spatial factors involved. One possible explanation is that the decision of locations, especially for residents' areas, is usually related to spatial factors, e.g., accessibility. This result implies the interaction between the socioeconomic and spatial elements, reinforcing the

validity of our approach in capturing variation in travel behaviour and distance decay within varied urban spaces.

### 6.1.2 Modelling mega city region using hierarchical vs unified system

Another critical discussion is whether to treat the polycentric mega-city region as a unified system in modelling human mobility within mega-city regions. Recent research increasingly suggests that the concept of city boundaries is diminishing in importance as urban spaces begin to fuse and traditional boundaries become indistinct (Batty, 2023; Dong et al., 2024). This trend aligns with the concept of a 'functional urban area' or 'metropolitan region', which encourages applying a holistic approach to understanding and modelling mobility patterns encompassing the entire mega-city region and even larger scales regions (Batty & Milton, 2021; Lomax et al., 2022; Lopane et al., 2023). However, our results suggest that there are still boundary effects in the human mobility pattern that could be observed across the urban system within the SDH area. The results from section 4.5.1 indicate that hierarchical modelling the area with a suitable boundary for sub-system could enhance the model's efficacy. Additionally, modelling the travel behaviours of individuals within the same functional city might improve the accuracy of representing the patterns of trips. Besides, the distance decays in Figure 6.1 (b) also show some discontinuity around the administrative boundary of Shenzhen-Huizhou and Dongguan-Huizhou. This discontinuity also implies the boundary effect may exist between the functional zones of different urban systems. Localised distance decay patterns typically emerge within each administrative boundary of these sprawling urban areas, transitioning from denser, granular centres to areas of more pronounced decay in the peripheries. Therefore, modelling the mega-city region, such as the SDH areas, using a hierarchical system may offer a more

nuanced and effective approach than a unified model. A hierarchical system can accommodate the varying levels of influence exerted by different centres within the region, reflecting the multi-layered nature of urban interactions and mobility patterns indicate that implementing a suitable boundary could enhance the model's efficacy. Additionally, modelling the travel behaviours of individuals within the same functional city might improve the accuracy of representing the patterns of trips. Besides, the distance decays in Figure 6.1 (b) also show some discontinuity around the administrative boundary of Shenzhen-Huizhou and Dongguan-Huizhou. This discontinuity also implies the boundary effect may exist between the functional zones of different urban systems. Localised distance decay patterns typically emerge within each administrative boundary of these sprawling urban areas, transitioning from denser, granular centres to areas of more pronounced decay in the peripheries. Therefore, modelling the mega-city region, such as the SDH areas, using a hierarchical system may offer a more nuanced and effective approach than a unified model. A hierarchical system can accommodate the varying levels of influence exerted by different centres within the region, reflecting the multi-layered nature of urban interactions and mobility patterns.

Political practices in China are another reason to promote this thesis to treat SDH areas as a hierarchical system rather than a continuous space. Local governments, particularly those at the prefecture level, have more considerable authority over urban planning and the development of public transportation infrastructure than Western countries' local governments (Chen & Yeh, 2023). This substantial local autonomy means that urban development and transportation networks can significantly differ from one jurisdiction to another, contributing to the heterogeneity observed within the mega-city regions and leading to distinct urban characteristics and functional zones. Another realistic issue is that the urban administrative

boundaries, to some extent, reflect the power struggle over territories among local governments, which makes it challenging for cities to accurately and consistently represent the true spatial extent of urban areas (Cartier, 2022; Chen & Yeh, 2023). An interesting observation is that Chinese scholars are more interested in identifying urban boundaries than their colleagues in Europe and America (Chen & Yeh, 2022; Gu et al., 2023; Li et al., 2020). It may be associated with the country-specific administrative and planning framework, particularly the prefecture-level city-led planning approach. Consequently, the spatial structure and mobility patterns in regions like SDH are needed to delineate boundaries of hierarchical systems and to better model and understand the urban dynamics in such polycentric areas. Meanwhile, delineating the urban functional area could better support the planning practice in China.

## 6.2  Calibration spatial interaction modelling with fine granularity

Initially, SIMs were developed as an aggregate method designed to forecast flows between zones using broad data indicators such as population, GDP, etc. The era of big data allows for the development of models based on granular spatial resolution, leading to a more accurate representation of real-world interactions on a large-scale (Batty & Milton, 2021). However, this trend also brings new challenges in calibrating the spatial interaction model.

In section 2.1.5, some research gaps have been identified for applying spatial Interaction models in small zones within urban systems from previous studies. This subchapter will explore various considerations and methodologies for calibrating spatial interaction models, drawing upon empirical evidence gathered during this doctoral research. The discussion will

include choosing an estimation method, and trade-offs between granular resolution and prediction reliability.

## 6.2.1 Choosing estimation methods

Selecting the appropriate estimation method for calibrating parameters in a Spatial Interaction Model (SIM) is indeed crucial, as it greatly affects the model's accuracy and interpretability. The calibration of SIMs primarily employs two methodologies: the regression-based method and the iteration-based maximum likelihood estimation (MLE) calibration algorithm, as introduced in Section 2.1.1. While both methods are theoretically grounded in MLE theory, the choice between them depends on the specific characteristics of the model and the nature of the data at hand which may cause a significant difference in goodness-of-fitting (Fotheringham & O'Kelly, 1989).

In this doctoral research, both calibration methods were applied to fit different research contexts: Chapter 3 utilises the iteration-based method for calibrating parameters of the Hierarchical Spatial Interaction Model (HSIM). In Chapters 4 and 5, the Origin-Specific Gravity Model (OSGM) is calibrated using the Poisson regression method.

The calibration process for each sub-model in the Hierarchical Spatial Interaction Model (HISM) generally follows the calibration method utilised in the QUANT model, as outlined by (Batty & Milton, 2021). This iteration-based method is a more flexible framework to set the iteration conditions (e.g. average travel cost or difference in flows). In practical applications, it has been observed that this iteration-based framework typically yields a slightly better goodness-of-fit for predicting flows, possibly because of its good totality, compared to the same

constraint-condition models calibrated using the standard Poisson regression method. Table 1 compares fitting performance between attractive-constrained iteration-based and attractive-constrained Poisson regression-based method results in London with 2001 census data. However, some issues prevent us from adopting the iteration-based calibration method applied for all models.

*Table 6.1 Comparison of fitting performance between attractive-constrained gravity model using iteration-based method regression-based method.*

|  | Iteration-based Method | Regression-based Method |
|---|---|---|
| MAE | 2.132 | 2.17 |
| RMSE | 12.288 | 12.625 |
| $R^2$ | 0.688 | 0.495 |

The first issue concerning the use of iteration-based methods revolves around the demand for computational resources. As the research shifts focus more localised distance decay effects in Chapter 4 and Chapter 5, the necessity of running an iterative process for each relatively small area becomes impractical, particularly in extensive regions like the GLA, where the number of sub-models could approach 1000. This high number of sub-models significantly increases the computational burden, making the approach less feasible.

The second issue relates to the challenges of achieving convergence with iteration-based methods, which are particularly sensitive to the distance function, iteration criteria, threshold settings, and starting values. Typically, iteration-based methods perform well when applied to the entire area under study. However, issues arise when the analysis narrows down to subsets of travel flows, such as those originating from the same area or within the same social group. In these instances, it has been observed that some iterative processes fail to converge. One possible assumption is that the flows in the sub-areas do not follow the distance decay functions we set or that the increase/decrease of average travel cost does not follow the same ratio of the distance decay parameters ($\beta$).

One possible assumption for this lack of convergence could be that the travel flows in these subsets do not adhere to the predetermined distance decay functions. Alternatively, the variations in average travel costs might not align proportionally with the expected values derived from parameter b, which typically modulates the impact of distance on travel behaviour in spatial interaction models. This assumption gains further credence when considering that other researchers have identified the coexistence of different distance decay functions within a single urban system (Šveda & Madajová, 2023). Furthermore, the patterns illustrated in Figure 6.2 indicate that some areas may exhibit a weak or non-existent distance decay effect which may cause the failure of the convergence.

*Figure 6.2  Different distance decay formats could exist within one city (source: Šveda and Madajová,,2023)*

As a result, the local difference of distance decay might be hidden if only one set of general parameters was calibrated for the whole modelling system. This discrepancy may arise because sub-areas can exhibit unique travel patterns not adequately captured by universal distance decay parameters. This convergence issue also highlights the importance of localising the gravity model to reflect travel behaviours within different segments of the urban area.

## 6.2.2  Predicting zero interaction

The issue of zero interactions between origin-destination (O-D) pairs is indeed not a new problem in modelling spatial interaction. Previous research raised this issue mainly because zero interactions would cause difficulty during the calibration process, particularly because logarithmic transformations are undefined at zero (Flowerdew & Aitkin, 1982; Fotheringham & O'Kelly, 1989; Sen & Sööt, 1981). Historically, this problem was not considered a major concern: when flows were aggregated over large areas, fewer instances of zero interactions could appear in the dataset. However, the shift towards more granular modelling units has exacerbated the zero-interaction issue. In finer resolution models, where the spatial units are smaller and the analysis more detailed, zero flows between O-D pairs become much more common due to the greater number of O-D pairs being considered, which naturally includes more pairs with no observed interactions, especially in less populated or less connected areas.

Fotheringham and O'Kelly (1989) listed several potential solutions for this issue: One approach is to eliminate all zero interactions from the analysis, and another is to remove all origins and destinations related to zero interactions. The third approach, which is also the most widely used method for handling zero interactions, is adding a constant to the elements of the interaction matrix.

However, all these solutions are imperfect when we apply them to our gravity models. Excluding zero flows from the dataset can induce biased parameter estimates, as this approach neglects the genuine absence or minimal interactions between specific origin-destination (O-D) pairs. Such omissions can markedly skew regression outcomes and interpretations, particularly within a framework based on regression analysis. Furthermore, eliminating origins

and destinations linked to zero interactions is impractical in granular models, as this would lead to removing the whole dataset. In detailed models, it is common for nearly all areas to exhibit zero interactions with at least one other area. While adding a constant value to the O-D matrix might seem like a viable strategy to retain maximum information, this approach becomes problematic at the granular modelling level. For example, 63.5% of the O-D pairs recorded zero interactions in our GLA commuting flow data set, and the mean flow value across all observations was merely 2.88. Consequently, even the addition of a small constant can disproportionately influence the calibration outcomes, altering the distribution and magnitude of flow values significantly.

Some techniques are suggested to reduce the influence of zero interactions based on empirical experience to establish the gravity model more granularly. Segmenting the data to isolate areas with high zero interactions and treating them as special cases within the model can also be a way to manage this issue. For instance, the HSIM model and OSGM applied in this doctoral study could significantly reduce the zero interactions by dividing the whole dataset into some sub-dataset. The zero-inflated gravity model has been widely used in economic gravity models to handle datasets with excessive zeros in the Poisson regression (Burger et al., 2009). Compared with the current method we used to predict the possibilities of travel in a single function, it first uses a model that predicts the occurrence of zero interaction, then uses another model for positive flows, offering a structured way to manage zero values (Martin & Pham, 2020).

In practical applications, it has been observed that the iteration-based calibration method exhibits greater resilience to zero interactions compared to regression-based methods,

primarily because it does not depend on logarithmic transformations. Nevertheless, a significant challenge with the iteration-based approach based on the O-D matrix is its tendency to overlook the underlying factors leading to the absence of flows between certain origins and destinations. This overlooking could result in the generation of unrealistic flows between areas where actual interactions should not occur. Moreover, excessive zero interactions can lead to convergence issues within the iteration-based method. This problem arises because a large presence of zeros challenges the foundational assumptions of the distance decay laws. One feasible solution is adjusting the value in the travel cost matrix to let the travel cost between the areas that are unlikely to have interactions (e.g., not accessible routes or policy prohibits) become extremely high, which could effectively prevent generating unrealistic flows.

### 6.2.3 Errors in converting predicted flows to integers

Low interaction presents a challenge that, while related to zero interaction, is distinct in its implications for spatial modelling. Given that travel flow data inherently count data, the values represent actual numbers of movements or interactions and, as such, should naturally be integers. Therefore, any flow predictions derived from iteration- or regression-based methods must be adjusted to integer values, as fractional values do not match real-world scenarios where individuals move between locations, not fractions of individuals. As discussed earlier, employing small geographical units with high spatial resolution in modelling tends to yield very low absolute flow values between most areas. This high granularity makes the data and subsequent operations extremely sensitive to processing or manipulation, potentially introducing significant errors in the analysis and modelling processes. Thus, this adjustment

must be conducted carefully to avoid distorting the underlying data patterns and to ensure that the model remains as accurate and representative of the observed phenomena as possible.

Thus, we import an additional iteration process to check if the sum flows at the origin match the observed value. If not, it should continue fitting with the observed values until the discrepancies are reduced to a small fraction. This method could be performed better in the modelling practice than simply rounding up or rounding down, as the change of differences in the goodness-of-fitting indicators is usually less after executing this process.

## 6.3  Uncertainty in urban modelling based on human mobility

As George Box's renowned aphorism states, "all models are wrong"(Box, 1976). The topic of error and uncertainty has been a longstanding subject of discussion among scholars, especially within the context of urban modelling, where the complexity and dynamic nature of cities make it less feasible to establish a model with perfect accuracy. Understanding and addressing these uncertainties is crucial for establishing urban simulation models and supporting planning and policy decision-making. Previous research has mentioned several factors, such as the model structure (Casman et al., 1999), input parameters (Ševčíková et al., 2007), and data transaction and processing (Yeh & Li, 2006).

During the modelling process of this doctoral study, we identified the primary source of uncertainty as stemming from the human mobility of big data. In our case study of the Shenzhen-Dongguan-Huizhou (SDH) area, we depended heavily on mobile phone data to

analyse and model human mobility within this region because there is no reliable open dataset such as census data published by the government. Shifting from traditional survey data to emerging mobile big data, such as mobile phone data, is undoubtedly huge progress because new data mobile methods greatly reduce data collection costs and expand the sample size. Unlike survey data or other traditional datasets collected explicitly for research purposes, mobile phone data is considered 'passive data'. This term refers to data generated for purposes other than research and is subsequently utilised in studies without having been actively solicited or collected through direct inquiries (C. Chen et al., 2016).

In the Shenzhen-Dongguan-Huizhou (SDH) area, our research employed two distinct mobile phone datasets provided by different operators, as detailed in Chapters 3 and 5, respectively. Both datasets were compiled in 2018, before the COVID-19 pandemic, and are comparable in size. Additionally, each data provider asserts that their dataset accurately reflects the general human mobility trends within the mega-city region. However, despite these similarities, notable discrepancies in travel behaviour patterns were observed between the two datasets.



*Figure 6.3 Difference of the urban functional zones when using different dataset*

As the model's complexity increases, the dataset's discrepancies could lead to significant variations in the simulation results when delineating urban functional zones. Figure 6.3 shows the difference between the two datasets using the same regionalisation algorithm. In the worst case, this uncertainty could influence stakeholders to make incorrect judgments based on the simulation results. Verifying the reliability of mobility big data is particularly concerning for urban simulation modelling. The lack of supplementary datasets for comparison or validation makes it challenging to ascertain the accuracy and representativeness of the big data being used. This situation highlights the need to add the validation process to control the uncertainty in data, e.g. running the model with additional data, which could enhance the reliability of the simulation results.

# 7 Conclusion

## 7.1 Contributions of this study

This study contributes to the field of urban analytics by proposing a comprehensive analytical framework to understand urban spatial structure transformation in mega-city regions from the human mobility perspective. Here are the key contributions:

**Improvement of Spatial Interaction Models**: This research contributes in addressing the limitations of current spatial interaction models, which often assume uniform spatial distribution and neglect local variations by emphasising the importance of incorporating local characteristics to enhance model accuracy. This study improves spatial interaction models by localising the model and including local socioeconomic and spatial characteristics. In addition, this study also discussed some long-standing issues in applying the spatial interaction models and providing distinctive insight with empirical evidence of establishing gravity models in a more granular level.

**Delineation the Urban Functional Zones**: The research proposes a novel approach to defining and delineating urban functional zones based on distance decay in human mobility patterns, moving beyond the traditional administrative boundaries. This method aims to better reflect the functional boundaries of urban areas, addressing the challenges of traditional models that fail to capture the extent and dynamics of urban spaces.

**Analytical Framework for Urban spatial structure in Mega-city region**: The research offers a framework to analyse urban spatial processes by integrating human mobility patterns, socioeconomic characteristics, and spatial structures. By establishing an integrated simulation model framework and testing it on the Great Bay Area, this doctoral study aims to help governments and planners to make informed decisions by predicting the outcomes of urban interventions on the spatial structure of mega-city regions over the medium to long term period.

## 7.2 Future research direction

Firstly, future research could focus on exploring distance decay variations in spatial interactions in more detail. As discussed in the discussion chapter, more than one distance decay law could exist within the same urban system. Therefore, developing more sophisticated models, for example, including multi-format distance decay functions, would better capture differences in spatial interaction within urban space. Meanwhile, the mechanism of how socioeconomic characteristics influence distance decay in spatial interactions is still unclear. Thus, more investigation is needed to explore the local variations in distance decay within different urban contexts and socioeconomic groups.

Secondly, future research may address the challenges of modelling human mobility and spatial interactions within polycentric mega-city regions and delve into hierarchical modelling approaches that reflect the complex nature of these urban areas. The current framework is generally based on the assumption that multiple functional centres and their spatial influence

do not overlap, but the realistic situation could be more complex. Future studies could explore how the overlay influence of the polycentric functional centres could be reflected in human mobility and functional zones.

Thirdly, tackling the specific technical challenges identified in establishing gravity models, especially at a granular level. These include issues related to zero interactions by applying the zero-inflated gravity model to predict the zeros interaction and a more detailed discussion of the potential reasons preventing convergence. Tackling these technical challenges effort would be beneficial in establishing a more reliable gravity model based on granular spatial units.

Finally, more effort could be made to explore uncertainties in urban modelling to better understand the associated challenges, particularly those arising from using big data for human mobility analysis. This could involve validating models with additional datasets and developing methods to reconcile discrepancies between data sources. Future research could focus on methodologies for assessing and mitigating uncertainties in spatial data and model design, thereby improving the robustness of urban analytical frameworks.

# References

Adam, A., Delvenne, J.-C., & Thomas, I. (2018). Detecting communities with the multi-scale Louvain method: robustness test on the metropolitan area of Brussels. *Journal of Geographical Systems*, *20*(4), 363-386.

Alessandretti, L., Aslak, U., & Lehmann, S. (2020). The scales of human mobility. *Nature*, *587*(7834), 402-407. https://doi.org/10.1038/s41586-020-2909-1

Alvioli, M. (2020). Administrative boundaries and urban areas in Italy: A perspective from scaling laws. *Landscape and Urban Planning*, *204*, 103906. https://doi.org/https://doi.org/10.1016/j.landurbplan.2020.103906

Anas, A., Arnott, R., & Small, K. A. (1998). Urban Spatial Structure. *Journal of Economic Literature*, *36*(3), 1426-1464. http://www.jstor.org/stable/2564805

Anejionu, O. C. D., Thakuriah, P., McHugh, A., Sun, Y., McArthur, D., Mason, P., & Walpole, R. (2019). Spatial urban data system: A cloud-enabled big data infrastructure for social and economic urban analytics. *Future Generation Computer Systems*, *98*, 456-473. https://doi.org/10.1016/j.future.2019.03.052

Arbia, G., & Petrarca, F. (2011). Effects of MAUP on spatial econometric models. *Letters in Spatial and Resource Sciences*, *4*(3), 173.

Arcaute, E., Hatna, E., Ferguson, P., Youn, H., Johansson, A., & Batty, M. (2015). Constructing cities, deconstructing scaling laws. *Journal of the royal society interface*, *12*(102), 20140745.

Arribas-Bel, D., & Sanz-Gracia, F. (2014). The validity of the monocentric city model in a polycentric age: US metropolitan areas in 1990, 2000 and 2010. *Urban Geography*, *35*(7), 980-997.

Asikhia, M. O., & Nkeki, N. F. (2013). Polycentric employment growth and the commuting behaviour in Benin Metropolitan Region, Nigeria. *Journal of Geography and Geology*, *5*(2), 1-17.

Barbosa, H., Barthelemy, M., Ghoshal, G., James, C. R., Lenormand, M., Louail, T., Menezes, R., Ramasco, J. J., Simini, F., & Tomasini, M. (2018). Human mobility: Models and applications. *Physics Reports*, *734*, 1-74. https://doi.org/https://doi.org/10.1016/j.physrep.2018.01.001

Batty, M. (1976a). *Urban modelling*. Cambridge University Press Cambridge.

Batty, M. (1976b). *Urban Modelling. Algorithms, Calibrations, Predictions*.

Batty, M. (2008). The Size, Scale, and Shape of Cities. *Science*, *319*(5864), 769-771. https://doi.org/10.1126/science.1151419

Batty, M. (2023). The boundary problem. *Environment and Planning B: Urban Analytics and City Science*, *50*(7), 1707-1710. https://doi.org/10.1177/23998083231202903

Batty, M., & Longley, P. A. (1994). *Fractal cities: a geometry of form and function*. Academic press.

Batty, M., & Mackie, S. (1972). The calibration of gravity, entropy, and related models of spatial interaction. *Environment and planning A*, *4*(2), 205-233.

Batty, M., & Milton, R. (2021). A new framework for very large-scale urban modelling. *Urban Studies*, *58*(15), 3071-3094. https://doi.org/10.1177/0042098020982252

Bazzani, A., Giorgini, B., Gallotti, R., Giovannini, L., Marchioni, M., & Rambaldi, S. (2011). Towards congestion detection in transportation networks using GPS data. 2011 IEEE Third International Conference on Privacy, Security, Risk and Trust and 2011 IEEE Third International Conference on Social Computing,

Bocarejo S, J. P., & Oviedo H, D. R. (2012). Transport accessibility and social inequities: a tool for identification of mobility needs and evaluation of transport investments. *Journal of Transport Geography*, *24*, 142-154. https://doi.org/https://doi.org/10.1016/j.jtrangeo.2011.12.004

Boeing, G. (2017). OSMnx: New methods for acquiring, constructing, analyzing, and visualizing complex street networks. *Computers, Environment and Urban Systems*, *65*, 126-139.

Boisjoly, G., & El-Geneidy, A. M. (2017). The insider: A planners' perspective on accessibility. *Journal of Transport Geography*, *64*, 33-43. https://doi.org/10.1016/j.jtrangeo.2017.08.006

Bokányi, E., Kallus, Z., & Gódor, I. (2019). Collective sensing of evolving urban structures: From activity-based to content-aware social monitoring. *Environment and Planning B: Urban Analytics and City Science*, *48*(1), 115-131. https://doi.org/10.1177/2399808319848760

Box, G. E. P. (1976). Science and Statistics. *Journal of the American Statistical Association*, *71*(356), 791-799. https://doi.org/10.1080/01621459.1976.10480949

Bradshaw, J., Kemp, P., Baldwin, S., & Rowe, A. (2004). The drivers of social exclusion. *London: Social Exclusion Unit*.

Brenner, N. (2002). Decoding the Newest "Metropolitan Regionalism" in the USA: A Critical Overview. *Cities*, *19*(1), 3-21. https://doi.org/10.1016/S0264-2751(01)00042-7

Bretagnolle, A., Paulus, F., & Pumain, D. (2002). Time and space scales for measuring urban growth. *Cybergeo: European Journal of Geography*.

Brown, W. M., Dar-Brodeur, A., & Tweedle, J. (2020). Firm networks, borders, and regional economic integration. *Journal of Regional Science*, *60*(2), 374-395.

Burger, M., & Meijers, E. (2012). Form Follows Function? Linking Morphological and Functional Polycentricity. *Urban Studies*, *49*(5), 1127-1149. https://doi.org/10.1177/0042098011407095

Burger, M., van Oort, F., & Linders, G.-J. (2009). On the Specification of the Gravity Model of Trade: Zeros, Excess Zeros and Zero-inflated Estimation. *Spatial Economic Analysis*, *4*(2), 167-190. https://doi.org/10.1080/17421770902834327

Calafati, A. G., & Veneri, P. (2013). Re-defining the boundaries of major Italian cities. *Regional Studies*, *47*(5), 789-802.

Cartier, C. (2022). "There are no cities in China" and the paradox of urban theory. *Eurasian Geography and Economics*, 1-18. https://doi.org/10.1080/15387216.2022.2047750

Casman, E. A., Morgan, M. G., & Dowlatabadi, H. (1999). Mixed levels of uncertainty in complex policy models. *Risk Analysis*, *19*(1), 33-42.

Chen, C., Ma, J., Susilo, Y., Liu, Y., & Wang, M. (2016). The promises of big data and small data for travel behavior (aka human mobility) analysis. *Transportation Research Part C: Emerging Technologies*, *68*, 285-299. https://doi.org/https://doi.org/10.1016/j.trc.2016.04.005

Chen, M., Liu, W., & Lu, D. (2016). Challenges and the way forward in China's new-type urbanization. *Land Use Policy*, *55*, 334-339.

Chen, T., Hui, E. C. M., Wu, J., Lang, W., & Li, X. (2019). Identifying urban spatial structure and urban vibrancy in highly dense cities using georeferenced social media data. *Habitat International*, *89*, 102005. https://doi.org/https://doi.org/10.1016/j.habitatint.2019.102005

Chen, Z., & Yeh, A. G.-O. (2022). Delineating functional urban areas in Chinese mega city regions using fine-grained population data and cellphone location data: A case of Pearl River Delta. *Computers, Environment and Urban Systems*, *93*, 101771. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2022.101771

Chen, Z., & Yeh, A. G.-O. (2023). Is prefecture-level city a "city" in China: a critical review. *Eurasian Geography and Economics*, 1-26. https://doi.org/10.1080/15387216.2023.2267064

Clark, C. (1951). Urban population densities. *Journal of the Royal Statistical Society. Series A (General)*, *114*(4), 490-496.

Cohen, I. G., & Mello, M. M. (2019). Big data, big tech, and protecting patient privacy. *Jama*, *322*(12), 1141-1142.

Cottineau, C., Finance, O., Hatna, E., Arcaute, E., & Batty, M. (2019). Defining urban clusters to detect agglomeration economies. *Environment and Planning B: Urban Analytics and City Science*, *46*(9), 1611-1626.

Curry, L. (1972). A spatial analysis of gravity flows. *Regional Studies*, *6*(2), 131-147.

Dadashpoor, H., & Yousefi, Z. (2018). Centralization or decentralization? A review on the effects of information and communication technology on urban spatial structure. *Cities*, *78*, 194-205. https://doi.org/https://doi.org/10.1016/j.cities.2018.02.013

De Domenico, M., Lima, A., González, M. C., & Arenas, A. (2015). Personalized routing for multitudes in smart cities. *EPJ Data Science*, *4*(1), 1. https://doi.org/10.1140/epjds/s13688-015-0038-0

De Vries, J. J., Nijkamp, P., & Rietveld, P. (2009). Exponential or power distance-decay for commuting? An alternative specification. *Environment and planning A*, *41*(2), 461-480.

Dennett, A. (2012). Estimating flows between geographical locations:'get me started in'spatial interaction modelling. *UCL working paper series* (181), 1-24.

Dennett, A. (2018). Modelling population flows using spatial interaction models. *Australian Population Studies*, *2*(2), 33-58.

Dennett, A., & Wilson, A. (2013). A Multilevel Spatial Interaction Modelling Framework for Estimating Interregional Migration in Europe. *Environment and Planning A: Economy and Space*, *45*(6), 1491-1507. https://doi.org/10.1068/a45398

Dong, L., Duarte, F., Duranton, G., Santi, P., Barthelemy, M., Batty, M., Bettencourt, L., Goodchild, M., Hack, G., Liu, Y., Pumain, D., Shi, W., Verbavatz, V., West, G. B.,

Yeh, A. G. O., & Ratti, C. (2024). Defining a city — delineating urban areas using cell-phone data. *Nature Cities*, *1*(2), 117-125. https://doi.org/10.1038/s44284-023-00019-z

Duque, J. C., Church, R. L., & Middleton, R. S. (2011). The p-Regions Problem [https://doi.org/10.1111/j.1538-4632.2010.00810.x]. *Geographical Analysis*, *43*(1), 104-126. https://doi.org/https://doi.org/10.1111/j.1538-4632.2010.00810.x

Duque, J. C., Ramos, R., & Suriñach, J. (2007). Supervised regionalization methods: A survey. *International Regional Science Review*, *30*(3), 195-220.

Engel, C., & Rogers, J. (1994). How wide is the border? In: National Bureau of Economic Research Cambridge, Mass., USA.

Engelfriet, L., & Koomen, E. (2018). The impact of urban form on commuting in large Chinese cities. *Transportation*, *45*(5), 1269-1295. https://doi.org/10.1007/s11116-017-9762-6

Flowerdew, R., & Aitkin, M. (1982). A method of fitting the gravity model based on the Poisson distribution. *Journal of Regional Science*, *22*(2), 191-202.

Fotheringham, A. S. (1981). Spatial Structure and Distance-Decay Parameters. *Annals of the Association of American Geographers*, *71*(3), 425-436. http://www.jstor.org/stable/2562901

Fotheringham, A. S. (1983). A new set of spatial-interaction models: the theory of competing destinations. *Environment and Planning A: Economy and Space*, *15*(1), 15-36.

Fotheringham, A. S., & Brunsdon, C. (1999). Local Forms of Spatial Analysis. *Geographical Analysis*, *31*(4), 340-358. https://doi.org/https://doi.org/10.1111/j.1538-4632.1999.tb00989.x

Fotheringham, A. S., Nakaya, T., Yano, K., Openshaw, S., & Ishikawa, Y. (2001). Hierarchical Destination Choice and Spatial Interaction Modelling: A Simulation Experiment. *Environment and Planning A: Economy and Space*, *33*(5), 901-920. https://doi.org/10.1068/a33136

Fotheringham, A. S., & O'Kelly, M. E. (1989). *Spatial interaction models: formulations and applications* (Vol. 1). Kluwer Academic Publishers Dordrecht.

Fotheringham, A. S., & Sachdeva, M. (2022). On the importance of thinking locally for statistics and society. *Spatial Statistics*, *50*, 100601. https://doi.org/https://doi.org/10.1016/j.spasta.2022.100601

Gao, Q.-L., Yue, Y., Zhong, C., Cao, J., Tu, W., & Li, Q.-Q. (2022). Revealing transport inequality from an activity space perspective: A study based on human mobility data. *Cities*, *131*, 104036. https://doi.org/https://doi.org/10.1016/j.cities.2022.104036

Gao, Q.-L., Zhong, C., Yue, Y., Cao, R., & Zhang, B. (2024). Income estimation based on human mobility patterns and machine learning models. *Applied Geography*, *163*, 103179. https://doi.org/https://doi.org/10.1016/j.apgeog.2023.103179

Gao, S., Janowicz, K., & Couclelis, H. (2017). Extracting urban functional regions from points of interest and human activities on location-based social networks. *Transactions in GIS*, *21*(3), 446-467.

Gao, S., Liu, Y., Wang, Y., & Ma, X. (2013). Discovering Spatial Interaction Communities from Mobile Phone Data. *Transactions in GIS*, *17*(3), 463-481. https://doi.org/10.1111/tgis.12042

Garcia-López, M.-À., & Muñiz, I. (2010). Employment Decentralisation: Polycentricity or Scatteration? The Case of Barcelona. *Urban Studies*, *47*(14), 3035-3056. https://doi.org/10.1177/0042098009360229

Geurs, K. T., & Van Wee, B. (2004). Accessibility evaluation of land-use and transport strategies: review and research directions. *Journal of Transport Geography*, *12*(2), 127-140. https://doi.org/10.1016/j.jtrangeo.2003.10.005

González, M. C., Hidalgo, C. A., & Barabási, A.-L. (2008). Understanding individual human mobility patterns. *Nature*, *453*(7196), 779-782. https://doi.org/10.1038/nature06958

Gordon, P., & Richardson, H. W. (1997). Are compact cities a desirable planning goal? *Journal of the American Planning Association*, *63*(1), 95-106.

Gordon, P., Richardson, H. W., & Wong, H. L. (1986). The Distribution of Population and Employment in a Polycentric City: The Case of Los Angeles. *Environment and Planning A: Economy and Space*, *18*(2), 161-173. https://doi.org/10.1068/a180161

Gould, P. (1975). Acquiring spatial information. *Economic Geography*, *51*(2), 87-99.

Government, C. C. (2019). *Guangdong-Hong Kong-Macao Greater Bay Area - Outline Development Plan*. Retrieved from https://www.bayarea.gov.hk/en/outline/plan.html

Green, N. (2007). Functional Polycentricity: A Formal Definition in Terms of Social Network Analysis. *Urban Studies*, *44*(11), 2077-2103. https://doi.org/10.1080/00420980701518941

Griffith, D. A. (2007). Spatial structure and spatial interaction: 25 years later. *Review of Regional Studies*, *37*(1), 28-38.

Griffith, D. A., & Jones, K. G. (1980). Explorations into the relationship between spatial structure and spatial interaction. *Environment and planning A*, *12*(2), 187-201.

Gu, Y., Shi, R., Zhuang, Y., Li, Q., & Yue, Y. (2023). How to determine city hierarchies and spatial structure of a megaregion? *Geo-spatial Information Science*, 1-13.

Guo, D. (2008). Regionalization with dynamically constrained agglomerative clustering and partitioning (REDCAP). *International Journal of Geographical Information Science*, *22*(7), 801-823.

Hagberg, A., & Conway, D. (2020). Networkx: Network analysis with python. *URL: https://networkx. github. io*.

Hajrasouliha, A. H., & Hamidi, S. (2017). The typology of the American metropolis: monocentricity, polycentricity, or generalized dispersion? *Urban Geography*, *38*(3), 420-444.

Halás, M., Klapka, P., & Kladivo, P. (2014). Distance-decay functions for daily travel-to-work flows. *Journal of Transport Geography*, *35*, 107-119. https://doi.org/https://doi.org/10.1016/j.jtrangeo.2014.02.001

Hall, P. G., & Pain, K. (2006). *The polycentric metropolis: learning from mega-city regions in Europe*. Routledge.

Hanson, S. (1982). The determinants of daily travel-activity patterns: relative location and sociodemographic factors. *Urban Geography*, *3*(3), 179-202.

Haque, M. B., Choudhury, C., Hess, S., & dit Sourd, R. C. (2019). Modelling residential mobility decision and its impact on car ownership and travel mode. *Travel Behaviour and Society*, *17*, 104-119.

Harris, B., & Batty, M. (1993). Locational Models, Geographic Information and Planning Support Systems. *Journal of Planning Education and Research*, *12*(3), 184-198. https://doi.org/10.1177/0739456X9301200302

Haynes, K. E., & Fotheringham, A. S. (1985). Gravity and spatial interaction models.

Helbich, M., Brunauer, W., Hagenauer, J., & Leitner, M. (2013). Data-driven regionalization of housing markets. *Annals of the Association of American Geographers*, *103*(4), 871-889.

Helsley, R. W., & Sullivan, A. M. (1991). Urban subcenter formation. *Regional Science and Urban Economics*, *21*(2), 255-275. https://doi.org/10.1016/0166-0462(91)90036-m

Henderson, M., Yeh, E. T., Gong, P., Elvidge, C., & Baugh, K. (2003). Validation of urban boundaries derived from global night-time satellite imagery. *International Journal of Remote Sensing*, *24*(3), 595-609. https://doi.org/10.1080/01431160304982

Henderson, V., & Mitra, A. (1996). The new urban landscape: Developers and edge cities. *Regional Science and Urban Economics*, *26*(6), 613-643. https://doi.org/10.1016/s0166-0462(96)02136-9

Hong, Y., & Yao, Y. (2019). Hierarchical community detection and functional area identification with OSM roads and complex graph theory. *International Journal of Geographical Information Science*, *33*(8), 1569-1587.

Hu, C., Li, Y., & Zheng, X. (2022). Data assets, information uses, and operational efficiency. *Applied Economics*, *54*(60), 6887-6900.

Hu, L. (2017). Changing travel behavior of Asian immigrants in the US. *Transportation research part A: policy and practice*, *106*, 248-260.

Hu, S., Xiong, C., Younes, H., Yang, M., Darzi, A., & Jin, Z. C. (2022). Examining spatiotemporal evolution of racial/ethnic disparities in human mobility and COVID-19 health outcomes: Evidence from the contiguous United States. *Sustainable Cities and Society*, *76*, 103506. https://doi.org/https://doi.org/10.1016/j.scs.2021.103506

Hu, W., & Jin, P. J. (2017). An adaptive hawkes process formulation for estimating time-of-day zonal trip arrivals with location-based social networking check-in data. *Transportation Research Part C: Emerging Technologies*, *79*, 136-155. https://doi.org/https://doi.org/10.1016/j.trc.2017.02.002

Huang, B., Zhou, Y., Li, Z., Song, Y., Cai, J., & Tu, W. (2019). Evaluating and characterizing urban vibrancy using spatial big data: Shanghai as a case study. *Environment and Planning B: Urban Analytics and City Science*, 239980831982873. https://doi.org/10.1177/2399808319828730

Jain, D., & Tiwari, G. (2019). Explaining travel behaviour with limited socio-economic data: Case study of Vishakhapatnam, India. *Travel Behaviour and Society*, *15*, 44-53. https://doi.org/https://doi.org/10.1016/j.tbs.2018.12.001

Jiang, B., & Miao, Y. (2015). The Evolution of Natural Cities from the Perspective of Location-Based Social Media. *The Professional Geographer*, *67*(2), 295-306. https://doi.org/10.1080/00330124.2014.968886

Jiang, S., Ferreira, J., & Gonzalez, M. C. (2017). Activity-Based Human Mobility Patterns Inferred from Mobile Phone Data: A Case Study of Singapore. *IEEE Transactions on Big Data*, *3*(2), 208-219. https://doi.org/10.1109/tbdata.2016.2631141

Jin, M., Gong, L., Cao, Y., Zhang, P., Gong, Y., & Liu, Y. (2021). Identifying borders of activity spaces and quantifying border effects on intra-urban travel through spatial

interaction network. *Computers, Environment and Urban Systems*, *87*, 101625. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2021.101625

Johnston, R. J. (1973). On frictions of distance and regression coefficients. *Area*, 187-191.

Kang, C., Liu, Y., Guo, D., & Qin, K. (2015). A Generalized Radiation Model for Human Mobility: Spatial Scale, Searching Direction and Trip Constraint. *PLOS ONE*, *10*(11), e0143500. https://doi.org/10.1371/journal.pone.0143500

Kattiyapornpong, U., & Miller, K. E. (2009). Socio-demographic constraints to travel behavior. *International Journal of Culture, Tourism and Hospitality Research*, *3*(1), 81-94.

Kitamura, R., Mokhtarian, P. L., & Laidet, L. (1997). A micro-analysis of land use and travel in five neighborhoods in the San Francisco Bay Area. *Transportation*, *24*, 125-158.

Kodinariya, T. M., & Makwana, P. R. (2013). Review on determining number of Cluster in K-Means Clustering. *International Journal*, *1*(6), 90-95.

Lenormand, M., Bassolas, A., & Ramasco, J. J. (2016). Systematic comparison of trip distribution laws and models. *Journal of Transport Geography*, *51*, 158-169.

Li, X., Gong, P., Zhou, Y., Wang, J., Bai, Y., Chen, B., Hu, T., Xiao, Y., Xu, B., & Yang, J. (2020). Mapping global urban boundaries from the global artificial impervious area (GAIA) data. *Environmental Research Letters*, *15*(9), 094044.

Li, X., Pan, G., Wu, Z., Qi, G., Li, S., Zhang, D., Zhang, W., & Wang, Z. (2012). Prediction of urban human mobility using large-scale taxi traces and its applications. *Frontiers of Computer Science*, *6*(1), 111-121.

Li, Y. (2020). Towards concentration and decentralization: The evolution of urban spatial structure of Chinese cities, 2001–2016. *Computers, Environment and Urban Systems*, *80*, 101425.

Li, Y., Wu, F., & Hay, I. (2015). City-region integration policies and their incongruous outcomes: The case of Shantou-Chaozhou-Jieyang city-region in east Guangdong Province, China. *Habitat International*, *46*, 214-222. https://doi.org/10.1016/j.habitatint.2014.12.006

Li, Z., Xu, J., & Yeh, A. G. O. (2014). State Rescaling and the Making of City-Regions in the Pearl River Delta, China. *Environment and Planning C: Government and Policy*, *32*(1), 129-143. https://doi.org/10.1068/c11328

Lin, D., Allan, A., & Cui, J. (2015). The impacts of urban spatial structure and socio-economic factors on patterns of commuting: a review. *International Journal of Urban Sciences*, *19*(2), 238-255.

Liu, X., & Wang, M. (2016). How polycentric is urban China and why? A case study of 318 cities. *Landscape and Urban Planning*, *151*, 10-20. https://doi.org/10.1016/j.landurbplan.2016.03.007

Liu, Y., Sui, Z., Kang, C., & Gao, Y. (2014). Uncovering Patterns of Inter-Urban Trip and Spatial Interaction from Social Media Check-In Data. *PLOS ONE*, *9*(1), e86026. https://doi.org/10.1371/journal.pone.0086026

Lomax, N., Smith, A. P., Archer, L., Ford, A., & Virgo, J. (2022). An Open-Source Model for Projecting Small Area Demographic and Land-Use Change. *Geographical Analysis*, *54*(3), 599-622.

Long, Y. (2016). Redefining Chinese city system with emerging new data. *Applied Geography*, *75*, 36-48. https://doi.org/https://doi.org/10.1016/j.apgeog.2016.08.002

Long, Y., Zhai, W., Shen, Y., & Ye, X. (2018). Understanding uneven urban expansion with natural cities using open data. *Landscape and Urban Planning*, *177*, 281-293. https://doi.org/https://doi.org/10.1016/j.landurbplan.2017.05.008

Lopane, F. D., Kalantzi, E., Milton, R., & Batty, M. (2023). A land-use transport-interaction framework for large scale strategic urban modeling. *Computers, Environment and Urban Systems*, *104*, 102007. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2023.102007

Louail, T., Lenormand, M., Picornell, M., García Cantú, O., Herranz, R., Frias-Martinez, E., Ramasco, J. J., & Barthelemy, M. (2015). Uncovering the spatial structure of mobility networks. *Nature Communications*, *6*(1), 6007. https://doi.org/10.1038/ncomms7007

LTD, C. O. E. S. (2014). *Shenzhen to Zhongshan Cross-River Corridor Project Environmental Impact Report*. http://www.sz.gov.cn/attachment/0/85/85124/2123722.pdf

Lu, X., Wetter, E., Bharti, N., Tatem, A. J., & Bengtsson, L. (2013). Approaching the Limit of Predictability in Human Mobility. *Scientific Reports*, *3*(1). https://doi.org/10.1038/srep02923

Lynch, K. (1960). *The image of the city* (Vol. 11). MIT press.

Maoh, H., & Kanaroglou, P. (2007). Geographic clustering of firms and urban form: a multivariate analysis. *Journal of Geographical Systems*, *9*(1), 29-52. https://doi.org/10.1007/s10109-006-0029-6

Marceau, D. J. (1999). The scale issue in the social and natural sciences. *Canadian Journal of Remote Sensing*, *25*(4), 347-356.

Martin, W., & Pham, C. S. (2020). Estimating the gravity model when zero trade flows are frequent and economically determined. *Applied Economics*, *52*(26), 2766-2779.

Masser, I., & Brown, P. J. (1975). Hierarchical aggregation procedures for interaction data. *Environment and planning A*, *7*(5), 509-523.

Masucci, A. P., Serras, J., Johansson, A., & Batty, M. (2013). Gravity versus radiation models: On the importance of scale and heterogeneity in commuting flows. *Physical Review E*, *88*(2). https://doi.org/10.1103/physreve.88.022812

Matthews, R. B., Gilbert, N. G., Roach, A., Polhill, J. G., & Gotts, N. M. (2007). Agent-based land-use models: a review of applications. *Landscape Ecology*, *22*, 1447-1459.

Mattioli, G., & Scheiner, J. (2022). The impact of migration background, ethnicity and social network dispersion on air and car travel in the UK. *Travel Behaviour and Society*, *27*, 65-78. https://doi.org/https://doi.org/10.1016/j.tbs.2021.12.001

McCallum, J. (1995). National Borders Matter: Canada-U.S. Regional Trade Patterns. *The American Economic Review*, *85*(3), 615-623. http://www.jstor.org/stable/2118191

McDonald, J. F. (1987). The identification of urban employment subcenters. *Journal of Urban Economics*, *21*(2), 242-258. https://doi.org/10.1016/0094-1190(87)90017-9

Menon, C. (2012). The bright side of MAUP: Defining new measures of industrial agglomeration. *Papers in Regional Science*, *91*(1), 3-28.

Mieszkowski, P., & Smith, B. (1991). Analyzing urban decentralization. *Regional Science and Urban Economics*, *21*(2), 183-199. https://doi.org/10.1016/0166-0462(91)90033-j

Morgan, M. S., & Morrison, M. (1999). *Models as mediators*. Cambridge University Press Cambridge.

Morris, S., Humphrey, A., Pickering, A., Tipping, S., Templeton, I., & Hurn, J. (2013). National travel survey 2013. *The Department for Transport, NatCen Social Research, London, UK*.

Municipality, D. a. R. B. o. D. (2019). *Dongguan City Population Development Plan (2020-2035)*. Development and Reform Bureau of Dongguan Municipality

Municipality, D. a. R. B. o. D. (2022, 2022). *Zhongshan Population Development Plan (2020-2035)*. Development and Reform Bureau of Dongguan Municipality. http://www.zs.gov.cn/zwgk/zfgb/zfgb202204/zcjd/content/mpost_2088026.html

Municipality, D. a. R. B. o. H. (2023). *Huizhou City Land and Space Master Plan (2021-2035)*. https://www.gd.gov.cn/zwgk/wjk/qbwj/yfh/content/post_4243751.html

Municipality, N. R. B. o. S. (2021). *Outline of Territorial Spatial Planning, Shenzhen (2020-2035) (in Chinese)*. Natural Resources Bureau of Shenzhen Municipality Retrieved from http://pnr.sz.gov.cn/attachment/0/794/794784/8858879.pdf

https://pnr.sz.gov.cn/xxgk/gggs/content/post_8858879.html

Muñiz, I., Garcia-López, M. À., & Galindo, A. (2008). The Effect of Employment Sub-centres on Population Density in Barcelona. *Urban Studies*, *45*(3), 627-649. https://doi.org/10.1177/0042098007087338

Nakaya, T. (2001). Local spatial interaction modelling based on the geographically weighted regression approach. *GeoJournal*, *53*(4), 347-358. http://www.jstor.org/stable/41147623

Nazara, S., Hewings, G. J., & Sonis, M. (2006). An exploratory analysis of hierarchical spatial interaction: the case of regional income shares in Indonesia. *Journal of Geographical Systems*, *8*(3), 253-268.

Niu, N., & Jin, H. (2020). Integrating multiple data to identify building functions in China's urban villages. *Environment and Planning B: Urban Analytics and City Science*, *48*(6), 1527-1542. https://doi.org/10.1177/2399808320938796

Noulas, A., Scellato, S., Lambiotte, R., Pontil, M., & Mascolo, C. (2012). A Tale of Many Cities: Universal Patterns in Human Urban Mobility. *PLOS ONE*, *7*(5), e37027. https://doi.org/10.1371/journal.pone.0037027

Openshaw, S. (1977). Optimal zoning systems for spatial interaction models. *Environment and planning A*, *9*(2), 169-184.

Openshaw, S. (1984). Ecological fallacies and the analysis of areal census data. *Environment and planning A*, *16*(1), 17-31.

Openshaw, S. (1996). Developing GIS-relevant zone-based spatial analysis methods. *Spatial analysis: modelling in a GIS environment*, 55-73.

Openshaw, S., & Rao, L. (1995). Algorithms for Reengineering 1991 Census Geography. *Environment and Planning A: Economy and Space*, *27*(3), 425-446. https://doi.org/10.1068/a270425

Oshan, T. M. (2016). A primer for working with the Spatial Interaction modeling (SpInt) module in the python spatial analysis library (PySAL). *Region*, *3*(2), R11-R23.

Oshan, T. M. (2020). The spatial structure debate in spatial interaction modeling: 50 years on. *Progress in Human Geography*, *45*(5), 925-950. https://doi.org/10.1177/0309132520968134

Palmer, J. R. B., Espenshade, T. J., Bartumeus, F., Chung, C. Y., Ozgencil, N. E., & Li, K. (2013). New Approaches to Human Mobility: Using Mobile Phones for Demographic Research. *Demography*, *50*(3), 1105-1128. https://doi.org/10.1007/s13524-012-0175-z

Pappalardo, L., Rinzivillo, S., Qu, Z., Pedreschi, D., & Giannotti, F. (2013). Understanding the patterns of car travel. *The European Physical Journal Special Topics*, *215*(1), 61-73.

Park, S., Oshan, T. M., El Ali, A., & Finamore, A. (2021). Are we breaking bubbles as we move? Using a large sample to explore the relationship between urban mobility and segregation. *Computers, Environment and Urban Systems*, *86*, 101585. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2020.101585

Parr, H. (2007). Mental health, nature work, and social inclusion. *Environment and Planning D: Society and Space*, *25*(3), 537-561.

Phithakkitnukoon, S., Smoreda, Z., & Olivier, P. (2012). Socio-Geography of Human Mobility: A Study Using Longitudinal Mobile Phone Data. *PLOS ONE*, *7*(6), e39253. https://doi.org/10.1371/journal.pone.0039253

Plane, D. A. (1984). Migration space: Doubly constrained gravity model mapping of relative interstate separation. *Annals of the Association of American Geographers*, *74*(2), 244-256.

Preston, J., & Rajé, F. (2007). Accessibility, mobility and transport-related social exclusion. *Journal of Transport Geography*, *15*(3), 151-160.

Pyers, C. E. (1966). Evaluation of intervening opportunities trip distribution model. *Highway Research Record*(114).

Qian, Y., Zhou, W., Pickett, S. T. A., Yu, W., Xiong, D., Wang, W., & Jing, C. (2020). Integrating structure and function: mapping the hierarchical spatial heterogeneity of urban landscapes. *Ecological Processes*, *9*(1), 59. https://doi.org/10.1186/s13717-020-00266-1

Quilty, B. J., Diamond, C., Liu, Y., Gibbs, H., Russell, T. W., Jarvis, C. I., Prem, K., Pearson, C. A., Clifford, S., & Flasche, S. (2020). The effect of travel restrictions on the geographical spread of COVID-19 between large cities in China: a modelling study. *BMC medicine*, *18*(1), 259.

Rhee, I., Shin, M., Hong, S., Lee, K., Kim, S. J., & Chong, S. (2011). On the levy-walk nature of human mobility. *IEEE/ACM transactions on networking*, *19*(3), 630-643.

Richardson, H. W. (1969). Regional Economics. Location theory, urban structure and regional change. *Regional economics. Location theory, urban structure and regional change.*

Rinzivillo, S., Mainardi, S., Pezzoni, F., Coscia, M., Pedreschi, D., & Giannotti, F. (2012). Discovering the geographical borders of human mobility. *KI-Künstliche Intelligenz*, *26*(3), 253-260.

Rodrigue, J.-P. (2020). *The geography of transport systems*. Routledge.

Ruktanonchai, N. W., Ruktanonchai, C. W., Floyd, J. R., & Tatem, A. J. (2018). Using Google Location History data to quantify fine-scale human mobility. *International journal of health geographics*, *17*, 1-13.

Santé, I., García, A. M., Miranda, D., & Crecente, R. (2010). Cellular automata models for the simulation of real-world urban processes: A review and analysis. *Landscape and Urban Planning*, *96*(2), 108-122.

Sari Aslam, N., & Cheng, T. (2018). Smart Card Data and Human Mobility. In (pp. 111). https://doi.org/10.14324/111.9781787353886.

Sari Aslam, N., Zhu, D., Cheng, T., Ibrahim, M. R., & Zhang, Y. (2021). Semantic enrichment of secondary activities using smart card data and point of interests: a case study in London. *Annals of GIS*, *27*(1), 29-41. https://doi.org/10.1080/19475683.2020.1783359

Scheiner, J. (2006). Housing mobility and travel behaviour: A process-oriented approach to spatial mobility: Evidence from a new research field in Germany. *Journal of Transport Geography*, *14*(4), 287-298.

Schläpfer, M., Dong, L., O'Keeffe, K., Santi, P., Szell, M., Salat, H., Anklesaria, S., Vazifeh, M., Ratti, C., & West, G. B. (2021). The universal visitation law of human mobility. *Nature*, *593*(7860), 522-527. https://doi.org/10.1038/s41586-021-03480-9

Schneider, M. (1959). Gravity models and trip distribution theory. *Papers in Regional Science*, *5*(1), 51-56.

Scott, A. J. (2019). City-regions reconsidered. *Environment and Planning A: Economy and Space*, *51*(3), 554-580. https://doi.org/10.1177/0308518X19831591

Scott, A. J., Agnew, J., Soja, E. W., & Storper, M. (2001). Global city-regions: an overview. *Global City Regions, Oxford University Press, Oxford*.

Sen, A., & Sööt, S. (1981). Selected procedures for calibrating the generalized gravity model. *Papers in Regional Science*, *48*(1), 165-176.

Ševčíková, H., Raftery, A. E., & Waddell, P. A. (2007). Assessing uncertainty in urban simulations using Bayesian melding. *Transportation Research Part B: Methodological*, *41*(6), 652-669. https://doi.org/https://doi.org/10.1016/j.trb.2006.11.001

Shelton, T., Poorthuis, A., & Zook, M. (2015). Social media and the city: Rethinking urban socio-spatial inequality using user-generated geographic information. *Landscape and Urban Planning*, *142*, 198-211. https://doi.org/https://doi.org/10.1016/j.landurbplan.2015.02.020

Shen, Q. (2000). Spatial and Social Dimensions of Commuting. *Journal of the American Planning Association*, *66*(1), 68-82. https://doi.org/10.1080/01944360008976085

Shen, Y., & Batty, M. (2019). Delineating the perceived functional regions of London from commuting flows. *Environment and Planning A: Economy and Space*, *51*(3), 547-550. https://doi.org/10.1177/0308518x18786253

Shortt, N. K. (2009). Regionalization/Zoning Systems. In R. Kitchin & N. Thrift (Eds.), *International Encyclopedia of Human Geography* (pp. 298-301). Elsevier. https://doi.org/https://doi.org/10.1016/B978-008044910-4.00506-X

Simini, F., Barlacchi, G., Luca, M., & Pappalardo, L. (2021). A Deep Gravity model for mobility flows generation. *Nature Communications*, *12*(1), 6576. https://doi.org/10.1038/s41467-021-26752-4

Simini, F., González, M. C., Maritan, A., & Barabási, A.-L. (2012). A universal model for mobility and migration patterns. *Nature*, *484*(7392), 96-100. https://doi.org/10.1038/nature10856

Smith, C., Quercia, D., & Capra, L. (2013). Finger on the pulse: identifying deprivation using transit flow analysis. Proceedings of the 2013 conference on Computer supported cooperative work,

Smith, D. (2009). Polycentric Cities and Sustainable Development. *Regional Science. Available online: http://www. slideshare. net/DuncanSmith/polycentric-cities-and-sustainable-development (accessed on 4 September 2009)*.

Sohn, J. (2005). Are commuting patterns a good indicator of urban spatial structure? *Journal of Transport Geography*, *13*(4), 306-317.

Song, C., Qu, Z., Blumm, N., & Barabási, A.-L. (2010). Limits of Predictability in Human Mobility. *Science*, *327*(5968), 1018-1021. https://doi.org/10.1126/science.1177170

Song, X., Zhang, Q., Sekimoto, Y., & Shibasaki, R. (2014). Prediction of human emergency behavior and their mobility following large-scale disaster. Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining,

Song, Y., Fan, Y., Li, X., & Ji, Y. (2018). Multidimensional visualization of transit smartcard data using space–time plots and data cubes. *Transportation*, *45*, 311-333.

Spooner, F., Abrams, J. F., Morrissey, K., Shaddick, G., Batty, M., Milton, R., Dennett, A., Lomax, N., Malleson, N., Nelissen, N., Coleman, A., Nur, J., Jin, Y., Greig, R., Shenton, C., & Birkin, M. (2021). A dynamic microsimulation model for epidemics. *Social Science & Medicine*, *291*, 114461. https://doi.org/https://doi.org/10.1016/j.socscimed.2021.114461

Srinivasan, S., & Rogers, P. (2005). Travel behavior of low-income residents: studying two contrasting locations in the city of Chennai, India. *Journal of Transport Geography*, *13*(3), 265-274. https://doi.org/https://doi.org/10.1016/j.jtrangeo.2004.07.008

Statistics, O. f. N. (2021). *Census 2021*. https://census.gov.uk/

Stead, D. (2001). Relationships between Land Use, Socioeconomic Factors, and Travel Patterns in Britain. *Environment and Planning B: Planning and Design*, *28*(4), 499-528. https://doi.org/10.1068/b2677

Stefanidis, A., Cotnoir, A., Croitoru, A., Crooks, A., Rice, M., & Radzikowski, J. (2013). Demarcating new boundaries: mapping virtual polycentric communities through social media content. *Cartography and Geographic Information Science*, *40*(2), 116-129. https://doi.org/10.1080/15230406.2013.776211

Stouffer, S. A. (1940). Intervening Opportunities: A Theory Relating Mobility and Distance. *American sociological review*, *5*(6), 845. https://doi.org/10.2307/2084520

Šveda, M., & Madajová, M. S. (2023). Estimating distance decay of intra-urban trips using mobile phone data: The case of Bratislava, Slovakia. *Journal of Transport Geography*, *107*, 103552. https://doi.org/https://doi.org/10.1016/j.jtrangeo.2023.103552

Tannier, C., Thomas, I., Vuidel, G., & Frankhauser, P. (2011a). A Fractal Approach to Identifying Urban Boundaries. *Geographical Analysis*, *43*(2), 211-227. https://doi.org/https://doi.org/10.1111/j.1538-4632.2011.00814.x

Tannier, C., Thomas, I., Vuidel, G., & Frankhauser, P. (2011b). A Fractal Approach to Identifying Urban Boundaries. 城市边界识别的分形方法. *Geographical Analysis*, *43*(2), 211-227. https://doi.org/https://doi.org/10.1111/j.1538-4632.2011.00814.x

Thomas, I., Jones, J., Caruso, G., & Gerber, P. (2018). City delineation in European applications of LUTI models: review and tests. *Transport Reviews*, *38*(1), 6-32. https://doi.org/10.1080/01441647.2017.1295112

Tizzoni, M., Bajardi, P., Decuyper, A., Kon Kam King, G., Schneider, C. M., Blondel, V., Smoreda, Z., González, M. C., & Colizza, V. (2014). On the Use of Human Mobility Proxies for Modeling Epidemics. *PLOS Computational Biology*, *10*(7), e1003716. https://doi.org/10.1371/journal.pcbi.1003716

Uniman, D. L., Attanucci, J., Mishalani, R. G., & Wilson, N. H. (2010). Service reliability measurement using automated fare card data: application to the London underground. *Transportation research record*, *2143*(1), 92-99.

van de Coevering, P., & Schwanen, T. (2006). Re-evaluating the impact of urban form on travel patternsin Europe and North-America. *Transport policy*, *13*(3), 229-239.

Wang, H., Zeng, S., Li, Y., & Jin, D. (2021). Predictability and Prediction of Human Mobility Based on Application-Collected Location Data. *IEEE Transactions on Mobile Computing*, *20*(7), 2457-2472. https://doi.org/10.1109/TMC.2020.2981441

Wang, L., Wong, C., & Duan, X. (2016). Urban growth and spatial restructuring patterns: The case of Yangtze River Delta Region, China. *Environment and Planning B: Planning and Design*, *43*(3), 515-539.

Wang, S., Zhang, X., Cao, J., He, L., Stenneth, L., Yu, P. S., Li, Z., & Huang, Z. (2017). Computing Urban Traffic Congestions by Incorporating Sparse GPS Probe Data and Social Media Data. *ACM Trans. Inf. Syst.*, *35*(4), Article 40. https://doi.org/10.1145/3057281

Wegener, M. (2004). Overview of land use transport models. In *Handbook of transport geography and spatial systems*. Emerald Group Publishing Limited.

Williams, N. E., Thomas, T. A., Dunbar, M., Eagle, N., & Dobra, A. (2015). Measures of Human Mobility Using Mobile Phone Records Enhanced with GIS Data. *PLOS ONE*, *10*(7), e0133630. https://doi.org/10.1371/journal.pone.0133630

Willumsen, L. G. (2001). *Modelling transport*. Wiley-Blackwell.

Wilson, A. (1970). *Entropy in urban and regional modelling* (Vol. 1). Routledge.

Wilson, A. G. (1971). A Family of Spatial Interaction Models, and Associated Developments. *Environment and Planning A: Economy and Space*, *3*(1), 1-32. https://doi.org/10.1068/a030001

Wu, C., Smith, D., & Wang, M. (2021). Simulating the urban spatial structure with spatial interaction: A case study of urban polycentricity under different scenarios. *Computers, Environment and Urban Systems*, *89*, 101677. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2021.101677

Wu, C., Wang, J., Wang, M., & Kraak, M.-J. (2024). Machine learning-based characterisation of urban morphology with the street pattern. *Computers, Environment and Urban Systems*, 102078. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2024.102078

Wu, F. (2016). China's Emergent City-Region Governance: A New Form of State Spatial Selectivity through State-orchestrated Rescaling. *International Journal of Urban and Regional Research*, *40*(6), 1134-1151. https://doi.org/10.1111/1468-2427.12437

Wu, F., & Yeh, A. G.-O. (1999). Urban Spatial Structure in a Transitional Economy. *Journal of the American Planning Association*, *65*(4), 377-394. https://doi.org/10.1080/01944369908976069

Xu, J., & Yeh, A. G. O. (2013). Interjurisdictional Cooperation through Bargaining: The Case of the Guangzhou–Zhuhai Railway in the Pearl River Delta, China. *The China Quarterly*, *213*, 130-151. https://doi.org/10.1017/S0305741013000283

Yan, X.-Y., Zhao, C., Fan, Y., Di, Z., & Wang, W.-X. (2014). Universal predictability of mobility patterns in cities. *Journal of the royal society interface*, *11*(100), 20140834. https://doi.org/10.1098/rsif.2014.0834

Yao, X., Gao, Y., Zhu, D., Manley, E., Wang, J., & Liu, Y. (2021). Spatial Origin-Destination Flow Imputation Using Graph Convolutional Networks. *IEEE Transactions on Intelligent Transportation Systems*, *22*(12), 7474-7484. https://doi.org/10.1109/TITS.2020.3003310

Yeh, A. G.-O., & Chen, Z. (2019). From cities to super mega city regions in China in a new wave of urbanisation and economic transition: Issues and challenges. *Urban Studies*, 0042098019879566. https://doi.org/10.1177/0042098019879566

Yeh, A. G.-O., & Li, X. (2006). Errors and uncertainties in urban cellular automata. *Computers, Environment and Urban Systems*, *30*(1), 10-28. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2004.05.007

Yin, J., Soliman, A., Yin, D., & Wang, S. (2017). Depicting urban boundaries from a mobility network of spatial interactions: a case study of Great Britain with geo-located Twitter data. *International Journal of Geographical Information Science*, *31*(7), 1293-1313. https://doi.org/10.1080/13658816.2017.1282615

Yuan, J., Zheng, Y., & Xie, X. (2012). Discovering regions of different functions in a city using human mobility and POIs. Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining,

Yuan, J., Zheng, Y., Zhang, C., Xie, W., Xie, X., Sun, G., & Huang, Y. (2010). *T-drive: driving directions based on taxi trajectories* Proceedings of the 18th SIGSPATIAL International Conference on Advances in Geographic Information Systems, San Jose, California. https://doi.org/10.1145/1869790.1869807

Zhai, W., Bai, X., Shi, Y., Han, Y., Peng, Z.-R., & Gu, C. (2019). Beyond Word2vec: An approach for urban functional region extraction and identification by combining Place2vec and POIs. *Computers, Environment and Urban Systems*, *74*, 1-12. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2018.11.008

Zhang, B., Zhong, C., Gao, Q., Shabrina, Z., & Tu, W. (2022). Delineating urban functional zones using mobile phone data: A case study of cross-boundary integration in Shenzhen-Dongguan-Huizhou area. *Computers, Environment and Urban Systems*, *98*, 101872. https://doi.org/https://doi.org/10.1016/j.compenvurbsys.2022.101872

Zhang, S. (2018). Computing Local Fractal Dimension Using Geographical Weighting Scheme.

Zhang, W., Fang, C., Zhou, L., & Zhu, J. (2020). Measuring megaregional structure in the Pearl River Delta by mobile phone signaling data: A complex network approach. *Cities*, *104*, 102809. https://doi.org/10.1016/j.cities.2020.102809

Zhang, X., Guo, Q.-e., Cheung, D. M.-w., & Zhang, T. (2018). Evaluating the institutional performance of the Pearl River Delta integration policy through intercity cooperation network analysis. *Cities*, *81*, 131-144. https://doi.org/https://doi.org/10.1016/j.cities.2018.04.002

Zhang, X., & Li, N. (2024). An activity space-based gravity model for intracity human mobility flows. *Sustainable Cities and Society*, *101*, 105073. https://doi.org/https://doi.org/10.1016/j.scs.2023.105073

Zhang, X., & Sun, Y. (2019). Investigating institutional integration in the contexts of Chinese city-regionalization: Evidence from Shenzhen–Dongguan–Huizhou. *Land Use Policy*, *88*, 104170. https://doi.org/https://doi.org/10.1016/j.landusepol.2019.104170

Zhang, Y., Marshall, S., Cao, M., Manley, E., & Chen, H. (2021). Discovering the evolution of urban structure using smart card data: The case of London. *Cities*, *112*, 103157. https://doi.org/10.1016/j.cities.2021.103157

Zheng, Y., Xie, X., & Ma, W.-Y. (2010). Geolife: A collaborative social networking service among user, location and trajectory. *IEEE Data Eng. Bull.*, *33*(2), 32-39.

Zheng, Y., Zhang, L., Xie, X., & Ma, W.-Y. (2009). *Mining interesting locations and travel sequences from GPS trajectories* Proceedings of the 18th international conference on World wide web, Madrid, Spain. https://doi.org/10.1145/1526709.1526816

Zhong, C., Arisona, S. M., Huang, X., Batty, M., & Schmitt, G. (2014). Detecting the dynamics of urban structure through spatial network analysis. *International Journal of Geographical Information Science*, *28*(11), 2178-2199. https://doi.org/10.1080/13658816.2014.914521

Zhong, C., Batty, M., Manley, E., Wang, J., Wang, Z., Chen, F., & Schmitt, G. (2016). Variability in Regularity: Mining Temporal Mobility Patterns in London, Singapore and Beijing Using Smart-Card Data. *PLOS ONE*, *11*(2), e0149222. https://doi.org/10.1371/journal.pone.0149222

Zhong, C., Manley, E., Müller Arisona, S., Batty, M., & Schmitt, G. (2015). Measuring variability of mobility patterns from multiday smart-card data. *Journal of Computational Science*, *9*, 125-130. https://doi.org/https://doi.org/10.1016/j.jocs.2015.04.021

Zhong, C., Schläpfer, M., Müller Arisona, S., Batty, M., Ratti, C., & Schmitt, G. (2017). Revealing centrality in the spatial structure of cities from human activity patterns. *Urban Studies*, *54*(2), 437-455. https://doi.org/10.1177/0042098015601599

Zhou, Y., Xu, R., Hu, D., Yue, Y., Li, Q., & Xia, J. (2020). Effects of human mobility restrictions on the spread of COVID-19 in Shenzhen, China: a modelling study using mobile phone data. *The Lancet Digital Health*, *2*(8), e417-e424. https://doi.org/10.1016/s2589-7500(20)30165-5

Zipf, G. K. (1946). The P 1 P 2/D hypothesis: on the intercity movement of persons. *American sociological review*, *11*(6), 677-686.

# Appendix A: List of Abbreviations

**GLA**          Great London area

**HSIM**         Hierarchical spatial interaction model

**LUTI**         Land use/Transport interaction

**MAE**          Mean absolute error

**MAUP**         Modifiable areal unit problem

**MSE**          Mean squared error

**O-D**          Origin-Destination

**OSGM**         Origin-specific gravity model

**POI**          Point of interest

**RMSE**         Root-mean-square deviation

**SDH**          Shenzhen-Dongguan-Huizhou

**SDHZ**         Shenzhen-Dongguan-Huizhou-Zhongshan

**SI**           Spatial interaction

**SIM**          Spatial interaction model

**UFZ(s)**       Urban functional zone(s)

# Appendix B: Data Inventory

| Datasets | Year | Provider | Data categories | Ethnics clarence needed* |
|---|---|---|---|---|
| UK Census data | 2001, 2011 | Office for National Statistics (ONS) | Public dataset | No |
| Great Bay Area mobile phone data | 2018 | China Unicom | Completely anonymous aggerated data | No |
| Great Bay Area road network | 2023 | OpenStreet Map | Public dataset | No |
| Great Bay Area residents/work location data and socioeconomic data | 2018, 2020 | Baidu | Completely anonymous aggerated data | No |

*The requirement of ethics clarence is based on the research ethics guideline for the "use of pre-existing data" published by King's College London in October 2019