Forecasting Fine-Grained Urban Flows Via Spatio-Temporal Contrastive Self-Supervision

Hao Qu, Yongshun Gong[®], *Member, IEEE*, Meng Chen[®], *Member, IEEE*, Junbo Zhang[®], *Member, IEEE*, Yu Zheng[®], *Fellow, IEEE*, and Yilong Yin[®]

Abstract—As a critical task of the urban traffic services, fine-grained urban flow inference (FUFI) benefits in many fields including intelligent transportation management, urban planning, public safety. FUFI is a technique that focuses on inferring fine-grained urban flows depending solely on observed coarse-grained data. However, existing methods always require massive learnable parameters and the complex network structures. To reduce these defects, we formulate a contrastive self-supervision method to predict fine-grained urban flows taking into account all correlated spatial and temporal contrastive patterns. Through several well-designed self-supervised tasks, uncomplicated networks have a strong ability to capture high-level representations from flow data. Then, a fine-tuning network combining with three pre-training encoder networks is proposed. We conduct experiments to evaluate our model and compare with other state-of-the-art methods by using two real-world datasets. All the empirical results not only show the superiority of our model against other comparative models, but also demonstrate its effectiveness in the resource-limited environment.

Index Terms—Contrastive self-supervision, fine-grained urban flow inference, spatio-temporal data

1 INTRODUCTION

WITH the developing trend of urbanization, intelligent transportation system has become one of the crucial components in the realm of smart cities [1], [2]. A critical requirement from urban planners and administrators is to monitor fine-grained urban flows, along with warnings in case of traffic congestion, public risk, etc [2], [3], [4], [5]. For example, streamed people caused a chaotic crowd stampede at the Falls Festival in Shanghai, leaved up to 36 people died and 80 people injured in a catastrophic stampede [6]. Urban Managers can locate high-risk regions and prevent people from such real tragedies by utilizing emergency mechanisms based on the fine-grained crowd warning and prediction model. Furthermore, with the telecommunication construction from 4 G to 5 G, the distance between base stations gradually decreases [7]. Fine-grained inference tasks, such as

(Corresponding authors: Yongshun Gong and Yilong Yin.) Recommended for acceptance by R. C.-W. Wong. Digital Object Identifier no. 10.1109/TKDE.2022.3200734 fine-grained urban flow prediction can provide a more accurate guidance to set 5 G base stations from the human mobility aspect.

However, forecasting fine-granularity urban flows signify that large numbers of monitoring equipment (e.g., mobile devices, surveillance cameras and piezoelectric sensors) have to be developed over the city [8], [9], [10]. Despite thousands of sensing devices bring convenience to the public, they consume huge amounts of power resources. For example, authorities get costly in operating ubiquitous monitoring equipment in terms of the procurement, manpower and maintenance fees, which increases the financial pressure of the government [11], [12].

To address such problems, fine-grained urban flow inference (FUFI) is proposed recently, which focuses on estimating fine-grained flows depending solely on observed coarse-grained data [3], [4], [5]. Fig. 1 gives an example of this process. Fig. 1a and 1b illustrate the same city area but with two different division scales, the left sub-figure is the coarse-granularity map (32×32) and the right one represents the fine-granularity map (64×64). The goal of FUFI is to make an accurate prediction for the fine-grained flow map from the coarse flow data. Intuitively, FUFI is also recognized as a variant of image super-resolution but has its unique *structural constraint*, i.e., the sum of the flow volumes in fine-grained regions strictly equals that of the corresponding super-region.

Despite achievements in FUFI problem [3], [4], [5], most of them require a complex neural network architecture, a huge number of parameters, and a long-term training period. To present, contrastive self-supervision is an effective method to handle such issues, which has been well-performed in the field of computer vision [13], [14], [15], [16] and natural language process [17], [18]. These models have shown a strong representation learning ability in a large

<sup>Hao Qu, Yongshun Gong, Meng Chen, and Yilong Yin are with the School of Software, Shandong University, Jinan 250012, China. E-mail: hao. qu@mail.sdu.edu.cn, {ysgong, mchen, ylyin}@sdu.edu.cn.
Junbo Zhang and Yu Zheng are with JD Intelligent Cities Research, Bei-</sup>

Junbo Zhang and Yu Zheng are with JD Intelligent Cities Research, Beijing 101111, China. E-mail: {msjunbozhang, msyuzheng}@outlook.com.

Manuscript received 16 June 2021; revised 17 July 2022; accepted 6 August 2022. Date of publication 22 August 2022; date of current version 21 June 2023.

This work was supported in part by Shandong Excellent Young Scientists Fund (Oversea) under Grant 2022HWYQ-044, in part by the Natural Science Foundation of Shandong Province under Grant ZR2021QF034, in part by the Open Fund of Key Laboratory of Urban Land Resources Monitoring and Simulation, Ministry of Natural Resources, in part by the Shandong Provincial Natural Science Foundation for Distinguished Young Scholars under Grant ZR2021JQ26, in part by the Major Basic Research Project of Natural Science Foundation of Shandong Province under Grant ZR2021ZD15, in part by the Young Scholars Program of Shandong University, and in part by the Fundamental Research Funds of Shandong University.

^{1041-4347 © 2022} IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.



(a) Coarse-grained crowd flows (32×32)

(b) Fine-grained crowd flows (64×64)

Fig. 1. Traffic flow of two different granularities in Beijing, where each grid denotes a region. Fine-grained urban flow inference aims to infer from (a) Coarse-grained crowd flows to (b) Fine-grained crowd flows.

amount of unlabeled data or fewer labeled data. To the best our knowledge, existing contrastive self-supervised learning strategies cannot be utilized in the FUFI problem directly. In reality, this work faces several specific challenges when we formulate the problem:

• Spatial Contrastive Self-Supervision. Essentially, the flows of a region are mainly affected by the surrounding regions. However, two regions can have similar flows when they fall into the same functional area (e.g., business center, residential area, and tourist area) [19], [20]. As shown in Fig. 1a, even there is a long distance between regions A and B, they have a similar flow property. Previous FUFI studies focus mainly on neighboring correlations, while ignoring semantic similarities. Moreover, existing contrastive self-supervision usually uses the entire flow maps to set a comparison pair, but neglects the comparisons at the regional level. Therefore, how to devise an effective spatial contrastive learning method is a principal challenge that needs to be resolved.

• *Temporal Contrastive Self-Supervision*. Existing FUFI methods aim to predict one fine-grained flow map from a snapshot of the coarse-grained flow map at the current moment. This one-to-one approach does not make effective use of temporal information from the self-supervision perspective. The prediction of fine-grained urban flow is not only inferred from the current timestamp but also affected by previous conditions. Besides, the overall traffic flow changes in an area have strong periodic characteristics, which indicates that both sequential neighborhoods and semantic similarity points contribute to the flow inference. Failure to use of this information will lead to a poor performance.

• *External Factors*. External factors also play a crucial role in FUFI [4], [5]. For example, during peak hours of commuting traffic, the traffic flow of arterial roads is greater than other time periods. When severe weather occurs, people tend to be indoors rather than outdoors. Various external factors have different effects on the real-world fine-grained flow inference.

To address all challenges well, we propose a spatio-temporal contrastive self-supervision method for the FUFI, named as UrbanSTC. UrbanSTC contains three self-supervised pretext tasks: regional contrast, spatial super-resolution inference and temporal contrast. Regional contrast focuses on exploring similarities among regional-level flows based on the intrinsic spatial characteristics. Spatial super-resolution inference is an inference network that learns the spatial and upscaling patterns in the super-resolution process. Given the triplet sets of flow maps, the temporal contrast task bridges the distances between all positive pairs, while requires all negative pairs far away from each other. Finally, a fine-tuning network combining with three pre-training encoder networks is devised to make the fine-grained flow prediction. Differing from UrbanFM [3] and FODE [4], the proposed UrbanSTC achieves a significant performance improvement with a light architecture. The main contributions and innovations of this paper are summarized as follows:

• We propose a general framework of spatio-temporal contrastive self-supervision for the FUFI problem. We design two pretext tasks from the spatial aspect, i.e., the regional contrast and the spatial super-resolution inference. These two pretext tasks can identify the spatial underlying relationships among regions in terms of the surrounding property and semantic similarity.

• Two kinds of temporal contrast sampling methods, hard sampling and weight sampling, are proposed in this paper. The former method selects the most confident examples as the positive and negative pairs, and the latter leverages an adaptive weight strategy to rebuild positive and negative pairs of the anchor example.

• We incorporate external factors in the fine-tuning UrbanSTC network. Experimental results prove that the external influences benefit the final results because they have drawn useful information from events and weather conditions.

• We perform a collection of experiments on two types of dense and sparse real-world datasets to prove the effectiveness of our method compared with other state-of-the-art models. All evaluation results show that the proposed method UrbanSTC yields the best performance. Specifically, when the training data reduces, our model shows an outperformed prediction performance, which demonstrates that UrbanSTC has its own advantages in the absence of training data resources.

The rest of this paper is organized as follows: Section 2 includes a literature review. Section 3 formally defines our problem. The proposed method is shown in Section 4. All experimental results are shown in Section 5. Finally, conclusions are drawn in Section 6.

2 RELATED WORK

In this section, we first review the current studies on the fine-grained urban flow inference (FUFI) and self-supervised learning methods. Since the FUFI problem [3] can be treated as a variant of single image super-resolution (SISR), we then introduce the SISR problem and reveal the difference between them.

2.1 Fine-Grained Urban Flow Inference

FUFI aims at inferring fine-grained crowed flows in a city based on the coarse-grained observations, which is a variant of SISR in the traffic prediction field [21], [22]. Liang et al. [3] first propose a neural network named UrbanFM to address the FUFI problem, which mainly leverages the SRResNet [23] under the *structural constraint*. UrbanFM devises an M^2 -Normalization layer, which outputs a *distributions* across every patch of *M*-by-*M* subregions of an associated superregion. Shen et al. design a weather-affected FUFI Predictor (WFRFP) model based on the super-resolution scheme [24]. WFRFP explores the relationship between the weather conditions and flow distributions, and reduces the scope of the predicting area based on the corresponding coarse-grained flow map. However, the proposed architecture heavily relies on empirically stacking deep neural networks. To solve this problem, Chen et al. introduces the deep neural network model from the perspective of the combination of differential equations and neural networks [25]. They regard the training and prediction of neural networks as the ordinary differential equation problems. Since the neural ordinary differential equations (NODE) is proposed, Zhou et al. find that NODE can be used as a core module to solve the FUFI problem, which proposes a more general neural ODE architecture called FODE [4]. FODE can address the numerical instability problem of the previous method without causing additional memory costs. The key idea of FODE is to incorporate an affine coupling layer in each ODE block to avoid the inaccurate gradient issue. The difference between FODE and UrbanFM is that FODE utilizes ODE block instead of the ResNet block. Despite the success of the above models, existing techniques rely on massive parameters and complex neural network architectures.

2.2 Self-Supervised Learning

Self-supervised learning has gained popularity because it can avoid the cost of annotating large-scale datasets. It mainly uses auxiliary tasks (pretext) to mine some specific supervised information from the large-scale unsupervised data, and trains the network through this constructed supervised information, in order to learn valuable representations for downstream tasks. According to the manifestation of self-supervision tasks, self-supervision is divided into the following three types: Context-based, Temporalbased and Contrastive-based approaches.

Early Context-based self-supervised technique focuses on common rules to generate labels, such as Jigsaw puzzle [13], Image restoration [14], Color transformation [15] and Image rotations [16]. The above mentioned methods are applied in the field of computer vision. Besides, in the field of natural language processing, Word2vec [17] is a popular model to use sequence of sentences to construct auxiliary tasks for predicting words. Large-scale pre-training model Bert [18] uses MASK word method to construct auxiliary task. They have achieved remarkable results in many fields. Most of the methods introduced above are based on the samples' meta information but with specific constraints between samples. One of the Temporal-based methods uses the concept of similar features in the video frame [26], [27]. The assumption is that the features of adjacent frames in the video are similar, while the video frames are far apart that are dissimilar. Self-supervised constraints are performed by constructing such similar (positive) and dissimilar (negative) samples. Another temporal-based method constructs positive and negative example features by tracking different frames of an object [28]. Recently, Contrastive-based has become a dominant component in self-supervised learning, which builds representations by encoding dissimilar or similar properties [29], [30].

While self-upervised learning shines in the field of computer vision, natural language processing, video processing, etc, there is limited study focusing on the urban flow forecasting, especially in the FUFI problem. We will explore a spatio-temporal contrastive self-supervision method to predict fine-grained urban flows.

2.3 Image Super-Resolution

Single image super-resolution (SISR) refers to the reconstruction of a high-resolution image with only one low-resolution observation image, combining with some prior knowledge of the target image. It is one of the basic issues related to the image processing, and has a wide range of practical needs and application scenarios, e.g., applied in the digital imaging technology [31], video coding communication technology [32] and fine-grained crowdsourcing [33]. To date, there are three mainstream algorithms of SISR: interpolationbased, reconstruction-based and learning-based methods. In the interpolation-based method, early techniques focused on bicubic interpolation [34] and Lanczos resampling [35], which is fast but not accurate. Reconstruction-based SR methods [36], [37], [38] adopt sophisticated prior knowledge to solve Single image super-resolution with flexible and sharp details. However, as the scale factor increases, the performance of many reconstruction-based methods declines rapidly and usually time-consuming. Learning-based SISR methods utilize machine learning algorithms to analyze statistical relationships between the low-resolution (LR) and its corresponding high-resolution (HR) counterpart from a large quantity training dataset. Change et al. [39] proposed the neighbor embedding method that used the similar local geometry between LR and HR to restore HR image blocks. Meanwhile, many researchers focus on combining the advantages of reconstruction-based with learning-based methods [40], [41], [42].

With the rapid development of deep learning in recent years, many studies have achieved great success since they do not require many human-engineered features. An end-toend mapping method represented as CNNs between the LR and HR images is first proposed by Dong et al. [43]. Inspired by the superior performance of CNN, various models for CNN began to be applied for SR. Among them, Shi et al. [32] proposed an efficient sub-pixel convolutional layer to recover HR images with little additional computational cost compared with the deconvolutional layer. Due to the great progress of VGG-net in image classification [44], a deep CNN was applied for SISR in [45]. However, the deep network is prone to model degradation in the training phase. Kim et al. [45] proposed a residual structure that makes the training of deeper convolutional neural networks possible, which has greatly promoted the development of SISR.

However, there is a great disparity between the FUFI and image super-resolution task, i.e., the unique *structural constraint* in FUFI. *structural constraint* requires mining changes within the data from a coarse-grained view, while single image super-resolution on natural images is more inclined to recover the lost high-frequency information.

3 PROBLEM STATEMENT

Before clarifying our method, we first introduce some basic notations and then formulate the problem of FUFI. The

| TABLE 1 |
|--------------------|
| Symbol Description |

| Symbols | Descriptions |
|--|--|
| $\mathbf{X} = [\mathbf{X}_1^c, \mathbf{X}_2^c, \dots, \mathbf{X}_T^c]$ | The flow map contains <i>T</i> moments |
| I; J | The granularity of division of latitude and longitude |
| M | The upscaling factor |
| $x_{i,j}$ | The small region in flow |
| H; W | The length and width of feature maps in \mathbf{H}^{reg} and \mathbf{H}^{tcs} |
| C | The channel number of convolution kernel |
| \mathbf{X}^{mc} ; \mathbf{X}^{c} ; \mathbf{X}^{f} | The down-scaling coarse-grained flows map; The coarse-grained flows map; The fine-grained flows map; |
| \mathbf{H}^{reg} ; \mathbf{H}^{inf} ; \mathbf{H}^{tcs} | The low-level hidden feature maps for regional contrast, inference net, and temporal contrast |
| $\mathbf{Z}^{reg}; \mathbf{D}^{tcs}$ | The high-level semantic features in regional contrast and temporal contrast |
| \mathbf{U}^{f} | The flow inference high-level semantic representation |
| \mathbf{U}_{a}^{f} | The flow inference distribution map of the hidden state |
| W ^{<i>f</i>} | The weight matrix of flow inference |

main sysmbols used in this paper are summarized in Table 1.

- **Definition 1 (Grid Flow Maps).** Given a timestamp t, assume that $\mathbf{X} \in \mathbb{R}^{I \times J}_+$ is an urban flow map partitioned evenly into a $I \times J$ grid map at t, where a grid denotes a region as shown in Fig. 1. Each entry $x_{i,j} \in \mathbb{R}_+$ denotes the volume of the observed flow.
- **Definition 2 (Superregion & Subregion).** Figs. 1a and 1b illustrate the same city area but with two different division scales, the left sub-figure is the coarse-grained flows map (32×32) and the right one represents the fine-grained flows map (64×64) . M denotes the scaling factor controlling the resolution changes between the coarse- and fine-grained maps. Fig. 1 represents an example when M = 2. We use supperregion and subregions to define the larger grid and its constituent smaller regions respectively [3], [5].
- **Definition 3 (Structural Constraint).** The sum of the flow volumes in fine-grained subregions $x_{i',j'}^f$ strictly equals that of the corresponding superregion $x_{i,j}^c$

$$x_{i,j}^c = \sum_{i',j'} x_{i',j'}^f \quad \text{s.t.} \quad i = \left\lfloor \frac{i'}{M} \right\rfloor, j = \left\lfloor \frac{j'}{M} \right\rfloor, \tag{1}$$

where i = 1, 2, ..., I and j = 1, 2, ..., J.

Fine-Grained Urban Flow Inference. Given a coarse-grained map $\mathbf{X}^c \in \mathbb{R}^{I \times J}_+$ and the upscaling factor $M \in \mathbb{Z}_+$, the goal of this paper is to infer the fine-grained flow map $\mathbf{X}^f \in \mathbb{R}^{MI \times MJ}_+$ under the *structural constraint*.

4 THE PROPOSED METHOD

Fig. 2 illustrates the flowchart of UrbanSTC. Our model is pre-trained by spatial self-supervision and temporal selfsupervision, and then the pre-trained encoders are copied to the final network for fine-tuning. We propose three kinds of pretext strategies separately for the spatial and temporal self-supervision methods.

4.1 Spatial Self-Supervision

Urban flow data has typical spatial characteristics. Inspired by self-supervised learning, we provide two types of selfsupervision tasks on the spatial perspective: regional contrast and spatial super-resolution inference network.

4.1.1 Regional-Level Contrast Pre-Training

Regional contrast self-supervision is dedicated to mining flow relationships at the regional level. At any timestamp t, there are many regions having similar or dissimilar flow conditions in the coarse-grained flow map X^c . In Fig. 2, the light blue block (Reg) depicts an example for the regionallevel contrastive learning. Assume the black rectangle is an anchor region x^q . The regions with red and blue rectangles can be treated as positive and negative samples respectively via a semantic distance with x^q , as expressed in Equations (2) and 3

$$dist^{s}(\mathbf{x}^{q}, \mathbf{x}_{i,j}) = \sqrt{\left(\mathbf{x}^{q} - \mathbf{x}_{i,j}\right)^{2}},$$
(2)

where $\mathbf{x}_{i,j}$ is a candidate area in the flow map.

$$\mathbf{x}_{i,j} \in \begin{cases} \text{positive,} & dist^s(\mathbf{x}^q, \mathbf{x}_{i,j}) \leq \lambda \\ \text{negative,} & dist^s(\mathbf{x}^q, \mathbf{x}_{i,j}) > \lambda \end{cases}$$
(3)

in which λ is a threshold for distinguishing between positive and negative samples. Because of the different semantic distances among regions, we hope to remain such properties in their high-level representations, i.e, the representation distances between \mathbf{x}^q and positive regional samples $\{\mathbf{x}_{k_1}^+\}_{k_1=1}^{K_1}$ are closer enough, while all negative representations $\{\mathbf{x}_{k_2}^-\}_{k_2=1}^{K_2}$ are moving away from \mathbf{x}^q , where K_1 and K_2 are the numbers of selected positive and negative regional samples.

For a coarse-grained flow map \mathbf{X}^c , we first project it into a low-level hidden feature map $\mathbf{H}^{reg} \in \mathbb{R}^{H \times W \times \hat{C}}$ by utilizing a non-linear encoder. This component of our network is named regional level encoder $\mathbf{Enc}_{reg}(\cdot)$ which will be used in the fine-tuning process. Thereafter, H^{reg} is normalized by a batch normalization method [46] and reshaped to $\mathbf{S}^{reg} \in$ $\mathbb{R}^{HW \times C}$. At last, a fully connected layer with C hidden units produces high-level semantic features $\mathbf{Z}^{reg} \in \mathbb{R}^{HW \times C}$ for the coarse-grained flow map \mathbf{X}^c . Unlike some previous contrastive loss functions, such as InfoNCE contrastive loss, only select one example as the positive example strictly [30], [47], our method considers a set of regions from **Z**^{*reg*} as positive samples, and put all the rest as negative samples, which is similar as the strategy in [48]. Given a coarse-grained flow map \mathbf{X}^{c} , we can obtain its dense representation \mathbf{Z}^{reg} . For each \mathbf{X}^{c} , we will randomly select the regional anchor point, and



Fig. 2. The framework of UrbanSTC. There are there major parts: spatial self-supervision, temporal self-supervision and fine-tuning stage. Reg (light blue block) represents the Regional-level contrast; Inf (Light pink block) represents the spatial super-resolution inference; TCS (Light yellow block) represents the temporal contrast. Dec indicates a decoder that can convert the embedding vectors generated by the spatio-temporal self-supervision into the output fine-grained maps. Urban STC includes pre-training and fine-tuning stages. Among, Reg and Inf belong to the spatial self-supervision pre-training. TCS belongs to the temporal self-supervision pre-training. We first learn encoders through a spatio-temporal pre-training, and finally complete the network in the fine-tuning stage.

distinguish positive and negative regional samples by calculating the euclidean distances based on a pre-defined threshold λ . Then our contrastive loss function is expressed as

$$\mathcal{L}_{reg} = -\log \frac{\sum_{k_1=1}^{K_1} \exp sim(\mathbf{z}^q, \mathbf{z}_{k_1}^+)}{\sum_{k_1=1}^{K_1} \exp sim(\mathbf{z}^q, \mathbf{z}_{k_1}^+) + \sum_{k_2=1}^{K_2} \exp sim(\mathbf{z}^q, \mathbf{z}_{k_2}^-)},$$
(4)

where $\mathbf{z}^{q}, \mathbf{z}_{k_{1}}^{+}, \mathbf{z}_{k_{2}}^{-} \in \mathbf{Z}^{reg}$ and $sim(\mathbf{u}, \mathbf{v})$ is similarity function between two representations (e.g., inner product).

Through this method, positive regional samples should make similar representations close to each other rather than negative types of samples.

4.1.2 Spatial Super-Resolution Inference Network Pre-Training

Given a coarse-grained map $\mathbf{X}^c \in \mathbb{R}^{I \times J}_+$ and upscaling factor $M \in \mathbb{Z}_+$, FUFI aims to learn a super-resolution model to infer the fine-grained flow map $\mathbf{X}^f \in \mathbb{R}^{MI \times MJ}_+$ under the *structural*

constraint. The most important learning mechanism is how to split a coarse region x_{ij}^c to its M^2 fine-grained cells, which can be represented as $\mathbb{I} \in \mathbb{R}^{1 \times 1}_+ \to \mathbb{M} \in \mathbb{R}^{M \times M}_+$. To simulate this process, we design a spatial super-resolution inference network in our pre-training.

Our intention is to use a coarser granularity map to infer the pattern $\mathbb{I} \to \mathbb{M}$ with a pretext-task. In detail, we first get a down-scaling coarser granularity map $\mathbf{X}^{mc} \in \mathbb{R}^{\lfloor \frac{I}{M} \rfloor \times \lfloor \frac{J}{M} \rfloor}_{+}$ based on the coarse-grained map \mathbf{X}^c and M, where each entry of \mathbf{X}^{mc} equals to the sum of corresponding M^2 flow volumes in \mathbf{X}^c . Then we can construct a spatial super-resolution network inferring \mathbf{X}^c from \mathbf{X}^{mc} . This pre-text task is able to capture the $\mathbb{I} \to \mathbb{M}$ pattern in advance, and could be benefit for improving the inference capability of surrounding flows.

For a X^{mc} , we first encode it by two convolutional layers with *C* channels and 3×3 kernel size, each layer followed by Relu nonlinearity as shown in Fig. 3. The two convolutional layers are taken as a feature learning network to map \mathbf{X}^{mc} to the low-level hidden feature maps $\mathbf{H}^{inf} \in \mathbb{R}^{\underline{H} \times \underline{W} \times C}$. This component of our network is named spatial super-resolution encoder $\mathbf{Enc}_{inf}(\cdot)$ which is used later in the fine-tuning process. Then we can leverage the prior FUFI methods distributional upsampling at the end of their networks [3], [4], [5]. We also adopt M^2 -Normalization¹ to impose the *structural constraint* on the network. The final loss is computed by the pixel-wise Mean Square Error (MSE)

$$\mathcal{L}_{inf} = \frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{X}_{t}^{c} - \mathcal{F}_{inf} \left(\mathbf{X}_{t}^{mc}; \theta \right) \right\|^{2},$$
(5)

where θ represents all learnable parameters in the inference network.

This inference structure and function \mathcal{F}_{inf} are similar to our final fine-tuning UrbanSTC, please refer to Section 4.4 for details.

4.2 Temporal Self-Supervision

Existing FUFI studies focus on inferring the fine-grained flow map based solely on its coarse-grained one, ignoring that similar flow conditions at different moments will also contribute to the inference. Here we devise a temporal-contrastive self-supervision network (TCS) to extract the similarity information in the temporal dimension. For any timestamp *t*, we can get an anchor point \mathbf{X}_{t}^{c} , and then collect its positive ($\{\mathbf{X}_{t,k_3}^+\}_{k_3=1}^{K_3}$) and negative samples ($\{\mathbf{X}_{t,k_4}^-\}_{k_4=1}^{K_4}$) by identifying the similarities among samples, where K_3 and K_4 are the numbers of selected positive and negative temporal samples.

TCS constructs a self-supervised auxiliary task that narrows encoder features between the anchor example and positive samples, and keeps the negative samples far away. The TCS encoder $\text{Enc}_{tcs}(\cdot)$ has a similar structure to the spatial super-resolution inference network. It projects coarse-grained map \mathbf{X}_{t}^{c} to the low-level hidden feature map $\mathbf{H}_{t}^{tcs} \in \mathbb{R}^{H \times W \times C}$. Then we adopt a batch normalization layers and the global average pooling layer. Finally, Multilayer Perceptron (MLP) with Relu activation function is used to make nonlinearity, converting the encoder feature map \mathbf{H}_{t}^{tcs} to the high-level semantic features $\mathbf{D}_{t}^{tcs} \in \mathbb{R}^{C}$. As shown in Fig. 2 the light yellow block (TCS), there are three kinds of samples: anchor point, positive and negative samples. Next, we will introduce how to select them.

4.2.1 Hard Sampling

We first use a straightforward way to pick the closest and the farthest samples of anchor point as its positive and negative pair. The distances between samples are calculated by the euclidean distance method

$$dist^{t}(\mathbf{X}^{c}, \mathbf{X}^{k}) = \sqrt{\frac{1}{HW} \sum_{i=1}^{H} \sum_{j=1}^{W} (\mathbf{X}_{i,j}^{c} - \mathbf{X}_{i,j}^{k})^{2}},$$
(6)

where X^c is the current coarse-grained flow map and X^k is the flow map at other times. As shown in Fig. 2, the module of Temporal self-supervision, there are three types of samples indicating by green, red and blue points. They represent the



Fig. 3. Spatial Super-resolution Inference Network Pre-training. We get a down-scaling coarser granularity map (more coarse-grained) X^{mc} based on the coarse-grained map X^c and upscaling factor M. Spatial super-resolution inference network simplifies the difficulty of the task and imitates the process of inferring.

anchor point, positive samples and negative samples respectively. Hard sampling method aims to select the closest (positive) sample with the current anchor, and find the farthest one as the negative sample. Note that, time-contrastive approaches are widely used in the video processing, such as [26], [27], which only picks the positive samples within a time window, and put all the rest into the negative pool. It is because the natural analogies between adjacent frames of video data. However, the previous and next traffic snapshots are probably not the closest semantic samples of the anchor point due to the high periodicity in traffic flow prediction problems [49]. Thus we choose to calculate distances between the anchor point and all training samples.

4.2.2 Weight Sampling

Considering that the hard sampling cannot fully use the correlations among all temporal samples $\{\mathbf{X}_t^c\}_{t=1}^T$, we further propose a weight sampling method in this section. In detail, we select Top-*K* positive and negative samples with a weighted combination approach

$$\mathbf{X}_{t}^{+} = \sum_{k=1}^{K} \frac{1/dist_{k}^{t}}{\sum_{j=1}^{K} 1/dist_{j}^{t}} \mathbf{X}_{t,k}^{+},$$
(7)

$$\mathbf{X}_{t}^{-} = \sum_{k=1}^{K} \frac{dist_{k}^{t}}{\sum_{j=1}^{K} dist_{j}^{t}} \mathbf{X}_{t,k}^{-},$$

$$\tag{8}$$

where $dist_k^t$ denotes the euclidean distance between the anchor point and *k*th selected sample.

Algorithm 1 shows the detailed procedure of the weight sampling method. The results affected by these two sampling methods have been presented in Section 5.2.3.

| Algorithm 1. Weight Sampling |
|--|
| Input : original coarse data $\{\mathbf{X}^c\}$. |
| Output : complete data $\{\mathbf{X}^c, \mathbf{X}^+, \mathbf{X}^-\}$. |
| 1: for $x \in \{\mathbf{X}_1^c, \dots, \mathbf{X}_T^c\}$ do |
| 2: Build Max-heap and Min-heap. |
| 3: for $y \in \{\mathbf{X}_1^c, \dots, \mathbf{X}_T^c\}$ do |
| 4: if $x \neq y$ then |
| 5: Calculate the euclidean distance dis^t between x and y |
| 6: Adjust Max-head and Min-heap. |
| 7: Select Top- <i>k</i> positive and negative samples respectively. |
| 8: get \mathbf{X}_t^+ by Equation (7) |
| 9: get \mathbf{X}_t^- by Equation (8) |
| |

TCS uses a triplet loss [50] to optimize the pre-trained model. Given a triplet constraint $\mathcal{I} = \langle \mathbf{X}^c, \mathbf{X}^+, \mathbf{X}^- \rangle$. The triplet loss ensures that a pair of co-occuring \mathbf{X}_t^c (*anchor*) and



Fig. 4. External Factors Fusion. External Factors are separated into continuous features (blue block) and categorical features (yellow block).

 \mathbf{X}_{t}^{+} (*positive*) are closer to each other in the embedding space while moving away from \mathbf{X}_{t}^{-} (*negative*). We define the score of this triplet as

$$d_f(\mathcal{I}) = \left\| f(\mathbf{X}_t^c) - f(\mathbf{X}_t^+) \right\|_2^2 - \left\| f(\mathbf{X}_t^c) - f(\mathbf{X}_t^-) \right\|_2^2 + \alpha, \quad (9)$$

$$\mathcal{L}_{TCS} = \frac{1}{T} \sum_{t=1}^{T} (max \{ d_f(\mathcal{I}), 0 \}),$$
(10)

where f(.) is a non-linear mapping function that needs to be learned, and α is a positive margin parameter. Notably, the triplet constraint is more flexible to adapt to different levels of intra-class variances [51], [52], which guarantees the differences between various timestamps.

4.3 External Factor Fusion

External factors (e.g., temperature, wind speed, weather and holidays) affect the flow distribution over the subregions. For example, people are more inclined to walk out of the office area on holidays. And when bad weather comes, people prefer to stay indoors instead of outdoors. Therefore, we should take such external factors into consideration.

We initialize the external factors into continuous features and categorical features. Among them, continuous features including temperature and wind speed are directly concatenated to form a vector \mathbf{e}_{con} . Categorical features include timestamps, days, holidays and weather conditions (e.g., windy, rainy). We use the method in UrbanFM [3] to initializes external information. The categorical features are transformed into low-dimensional vectors by feeding into separate embedding layers, and then use concatenate operation to construct the categorical vector \mathbf{e}_{cat} . Then, we splice the two vectors \mathbf{e}_{con} and \mathbf{e}_{cat} to the final external embedding ($\mathbf{e} = [\mathbf{e}_{con}; \mathbf{e}_{cat}]$).

As shown in Fig. 4. we use two layers of multi-layer perception with nonlinear transformation to feed external embedding **e**. By using nonlinear transformation, different external factors are converged into a hidden state $\mathbf{X}^e \in \mathbb{R}^{I \times J}_+$. We regard it as a bias of flow graph. In the previous sections, we only used coarse-grained views without external information for pre-training. Finally, we use the tensor addition operation $\mathbf{X}^e + \mathbf{X}^e$ as the input of the model in the fine-tuning stage.

4.4 Fine-Tuning UrbanSTC

We derive three encoders when completing the above pretraining tasks, i.e., regional constrastive encoder $\mathbf{Enc}_{reg}(\cdot)$, spatial super-resolution inference encoder $\mathbf{Enc}_{inf}(\cdot)$ and TCS encoder $\mathbf{Enc}_{tcs}(\cdot)$. As illustrated in Fig. 2, three encoders are used for fine-tuning the downstream task. First, we combine three low-level hidden feature maps by encoders. This step can be described as

$$\mathbf{H}^{reg} = \mathbf{Enc}_{reg}(\mathbf{X}^c), \tag{11}$$

$$\mathbf{H}^{inf} = \mathbf{Enc}_{inf}(\mathbf{X}^c),\tag{12}$$

$$\mathbf{H}^{tcs} = \mathbf{Enc}_{tcs}(\mathbf{X}^c),\tag{13}$$

$$\mathbf{H}^{a} = \mathbf{Concat}(\mathbf{H}^{reg}, \mathbf{H}^{inf}, \mathbf{H}^{tcs}), \tag{14}$$

where Concat is the tensor concatenate operation. Then **Decoder** has a convolutional layer $(3 \times 3, C)$ with ReLU nonlinearity, which is used to decode three low-level hidden features. Besides, we adopt another convolutional layer $(3 \times 3, C \times M^2)$ and PixelShuffle layers, which rearranges features and increases sizes by the upscaling factor M. At the end of PixelSuffle, we use a ReLU activation function. After the above operations, a feature $\mathbf{U}^{f} \in \mathbb{R}^{MH \times MW \times C}$ is obtained where the first two dimensions have been increased M times. Next, we use a 3×3 convolution with the 1-size channel to get a fine-grained flow distribution map of the hidden state $\mathbf{U}_{o}^{f} \in \mathbb{R}^{MH \times MW \times 1}$. Due to the structural constraint of FUFI problem, the MSE loss cannot be used directly. Refer to the distributional upsampling in UrbanFM [3] and FODE [4], we choose a M^2 -Normalization that makes the sum of subregions equal to their corresponding superregion, which is described as

$$W^{f}_{(i,j)} = \frac{U^{f}_{o(i,j)}}{\sum_{i' \in \left(\lfloor \frac{i}{M} \rfloor M, (\lfloor \frac{i}{M} \rfloor + 1)M\right)]} U^{f}_{o(i'j')}},$$

$$(15)$$

where $U_{o(i,j)}^{f}$ is the *i*th row and *j*th column cell in \mathbf{U}_{o}^{f} and $W_{(i,j)}^{f} \in [0,1]$ represents probability.

 M^2 -Normalization aims to learn the probability mapping from a coarse-grained view to a fine-grained view. Finally, we infer the fine-grained crowds map by $\hat{\mathbf{X}}^f = \mathbf{X}^c \odot \mathbf{W}^f$. Mean Square Error (MSE) is used as the loss function

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^{T} \left\| \mathbf{X}_{t}^{f} - \mathcal{F} \left(\mathbf{X}_{t}^{c}; \theta \right) \right\|^{2},$$
(16)

where \mathcal{F} represents the UrbanSTC model and θ represents all learnable parameters used in this model.

5 EXPERIMENTS

In this chapter, we have conducted comprehensive experiments to demonstrate the effectiveness of our method. The source code has been released at https://github.com/ HaoQu59/UrbanSTC.

5.1 Experimental Settings

5.1.1 Datasets

We evaluate the performance of our model as well as baselines on two real-world urban flow datasets. The dataset statistics are shown in Table 2. In the experiments, we partition the data into non-overlapping training, validation and test data by a ratio of 2:1:1 respectively.

• *TaxiBJ* [3], [5] This dataset is collected from Beijing taxi flows, including four different periods: P1 to P4. The time interval is 30 minutes.

| Dataset | TaxiBJ | BikeNYC |
|---------------------|-------------------------------------|-----------------------------------|
| Time span | P1: 7/1/2013-10/31/2013 | |
| 1 | P2: 2/1/2014-6/30/2014 | 1/1/2019- |
| | P3: 3/1/2015-6/30/2015 | 31/3/2019 |
| | P4: 11/1/2015-3/31/2016 | |
| Time interval | 30 minutes | 1 hour |
| Coarse-grained size | 32×32 | 40×20 |
| Fine-grained size | 128×128 | $80{\times}40$ |
| Upscaling factor(M) | 4 | 2 |
| Latitude range | $39.82^{\circ}N - 39.99^{\circ}N$ | $40.65^{\circ}N - 40.81^{\circ}N$ |
| Longitude range | $116.26^{\circ}E - 116.49^{\circ}E$ | $73.93^{\circ}W - 74.01^{\circ}W$ |
| 0 | External Factors (meterology, | |
| | time and event) in TaxiBJ dataset | |
| Temperature /°C | [-24.6, 41.0] | λ. |
| Wind speed/mph | [0, 48.6] | λ. |
| Weather conditions | 16 types (e.g., Sunny) | λ. |
| Holidays | 18 | λ. |

TABLE 2 Statistics of Datasets

• *BikeNYC*² This dataset is collected from an open website that contains bike flow data in New York City from Jan 1 to Mar 31, 2019. We partition the city area into 40×20 grids as the coarse-grained map, and define the fine-granularity map with 80×40 .

5.1.2 Baselines

We compare the proposed method UrbanSTC with the following 13 baselines, including three types of methods, Heuristic, state-of-the-art image super-resolution and FUFI methods. All parameters of the proposed method and baselines adopt M^2 -Normalization to obey the *structural constraint* of FUFI.

Heuristic Methods.

• *Mean Partition (Mean)*: We evenly distribute coarsegrained maps into fine-grained maps according to the scaling factor.

• *Historical Average (HA)*: Predict the fine-grained subregions by the historical average of its corresponding superregion, and distribute flows into sub-regions based on historical split proportions.

Image Super-Resolution Methods.

• *SRCNN* [43]: It is the first method to introduce convolutional neural networks (CNNs) into image super-resolution problems. SRCNN first uses bicubic interpolation to enlarge the low-resolution image to the target size, then fits the non-linear mapping through a three-layer convolutional network, and finally outputs the high-resolution image result.

• *ESPCN* [32]: ESPCN proposes a sub-pixel convolution method to extract features directly from low-resolution image size, and calculate an efficient method to obtain high-resolution images.

• *VDSR* [45]: It is different from the three-stage architecture of SRCNN and ESPCN. VDSR is based on the idea of residual structure and uses a resolution method of deep neural networks with a depth of up to 20.

• SRResNet [23]: SRResNet uses perceptual loss and adversarial loss to improve the realism of the restored

picture. Perceptual loss is the feature extracted by the convolutional neural network.

• *DeepSD* [53]: DeepSD is the state-of-the-art method on statistical upscaling (i.e., super-resolution) for meteorological data. It uses a stacked strategy to use multiple SRCNNs for intermediate-level downscaling, and performs further upsampling by simply stacking these SRCNNs.

• *LapSRN* [54]: LapSRN is divided into two parts: feature extraction and image reconstruction. It uses low-resolution images directly as input to the network, and through stepby-step amplification, while reducing the amount of calculation, it also effectively improves the accuracy. And between the levels of each pyramid and within each level, parameter sharing is carried out through recursive.

• *IMDN* [55]: IMDN is a lightweight network architecture which contains distillation and selective fusion parts to address issues that excessive convolutions will limit the application of super-resolution technology in low computing power devices. They first use the distillation module to extract the hierarchical structure, and then use the contrast-based channel attention to fuse the features.

• *SCN* [56]: It is proved that modeling the scale invariance into the neural network can significantly improve the image restoration performance. Inspired by the spatial convolution of shift-invariance, "scale-wise convolution" is proposed to convolve across multiple scales for scale invariance.

FUFI Methods.

• *UrbanFM* [3]: UrbanFM first proposes Fine-grained urban flow super-resolution. Its difficulty is that the sum of the flow of multiple fine-grained areas is equal to the flow of a coarse-grained area and the mutual influence between adjacent areas. UrbanFM designs stacking ResNet-based neural networks and M^2 -Normalization layer to overcome.

• *UrbanPy* [5]: A progressive method of UrbanFM which uses a cascading model for forecasting fine-grained urban flows by decomposing the original tasks into multiple subtasks.

• *FODE* [4]: FODE is the state-of-the-art method in finegrained Urban Flow Super-Resolution. Inspired by the Neural Ordinary Differential Equations (NODE) [25]. They propose FODE block replaces ResNet as the backbone. IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 35, NO. 8, AUGUST 2023

TABLE 3 The Average RMSE, MAE and MAPE on TaxiBJ Dataset (P1) With Different Proportions of Training Data

| Methods | | P1(20%) |) | | P1(40% |) | | P1(60% |) | | P1(80% |) | | P1(100% | 76) |
|----------|--------|---------|---------|--------|--------|---------|--------|--------|---------|--------|--------|---------|--------|---------|---------|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE |
| MEAN | 20.918 | 12.019 | 4.469 | 20.918 | 12.019 | 4.469 | 20.918 | 12.019 | 4.469 | 20.918 | 12.019 | 4.469 | 20.918 | 12.019 | 4.469 |
| HA | 4.794 | 2.269 | 0.339 | 4.802 | 2.263 | 0.338 | 4.793 | 2.258 | 0.338 | 4.785 | 2.256 | 0.337 | 4.772 | 2.251 | 0.336 |
| SRCNN | 4.737 | 2.767 | 0.804 | 4.498 | 2.578 | 0.706 | 4.506 | 2.587 | 0.712 | 4.290 | 2.425 | 0.631 | 4.275 | 2.430 | 0.642 |
| ESPCN | 4.552 | 2.583 | 0.682 | 4.493 | 2.540 | 0.657 | 4.264 | 2.346 | 0.558 | 4.216 | 2.316 | 0.544 | 4.208 | 2.318 | 0.546 |
| DeepSD | 4.532 | 2.535 | 0.652 | 4.346 | 2.373 | 0.566 | 4.883 | 2.834 | 0.805 | 4.287 | 2.349 | 0.556 | 4.128 | 2.248 | 0.516 |
| VDŜR | 4.546 | 2.556 | 0.669 | 4.299 | 2.354 | 0.562 | 4.198 | 2.279 | 0.527 | 4.119 | 2.229 | 0.503 | 4.054 | 2.186 | 0.485 |
| SRResNet | 4.734 | 2.800 | 0.844 | 4.383 | 2.520 | 0.696 | 4.276 | 2.437 | 0.654 | 4.179 | 2.366 | 0.618 | 4.079 | 2.291 | 0.580 |
| LapSRN | 4.676 | 2.738 | 0.801 | 4.642 | 2.715 | 0.789 | 4.309 | 2.432 | 0.635 | 4.153 | 2.305 | 0.567 | 4.083 | 2.255 | 0.542 |
| IMDN | 4.696 | 2.748 | 0.794 | 4.388 | 2.464 | 0.635 | 4.251 | 2.376 | 0.601 | 4.159 | 2.295 | 0.554 | 4.085 | 2.253 | 0.538 |
| SCN | 4.395 | 2.491 | 0.661 | 4.219 | 2.351 | 0.588 | 4.096 | 2.250 | 0.536 | 4.028 | 2.203 | 0.515 | 3.965 | 2.162 | 0.494 |
| UrbanFM | 4.560 | 2.343 | 0.398 | 4.321 | 2.213 | 0.369 | 4.195 | 2.140 | 0.350 | 4.108 | 2.095 | 0.340 | 4.042 | 2.062 | 0.337 |
| UrbanPy | 4.665 | 2.471 | 0.547 | 4.363 | 2.233 | 0.415 | 4.112 | 2.077 | 0.349 | 4.033 | 2.041 | 0.343 | 3.944 | 1.998 | 0.333 |
| FODE | 4.476 | 2.304 | 0.391 | 4.260 | 2.170 | 0.349 | 4.161 | 2.116 | 0.344 | 4.084 | 2.078 | 0.338 | 4.002 | 2.044 | 0.336 |
| UrbanSTC | 4.083 | 2.022 | 0.302 | 3.988 | 1.983 | 0.302 | 3.941 | 1.962 | 0.301 | 3.900 | 1.942 | 0.298 | 3.845 | 1.922 | 0.298 |
| Δ | +7.10% | +12.24% | +10.91% | +5.48% | +8.62% | +10.65% | +3.78% | +5.54% | +10.95% | +3.18% | +4.85% | +11.57% | +2.51% | +3.80%- | +10.51% |

The best results are bold and the second best are underlined.

5.1.3 Evaluation Metrics

We evaluate different methods with three widely used metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE) and Mean Absolute Percentage Error (MAPE)

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^{N} \left(\mathbf{X}_{i} - \hat{\mathbf{X}}_{i}\right)^{2}}$$
$$MAE = \frac{1}{N} \sum_{i=1}^{N} |\mathbf{X}_{i} - \hat{\mathbf{X}}_{i}|$$
$$MAPE = \frac{1}{N} \sum_{i=1}^{N} \left|\frac{\mathbf{X}_{i} - \hat{\mathbf{X}}_{i}}{\mathbf{X}_{i}}\right|$$

where $\hat{\mathbf{X}}_i$ is a prediction for fine-grained flow, and \mathbf{X}_i is the ground truth; *N* is the number of prediction values.

5.1.4 Training Details & Hyperparameters

Our model and baselines are completely implemented by PyTorch 1.60 with a RTX 2080 GPU. The network is trained using Adam with the first and second moment estimates equaling to 0.9 and 0.999, respectively [57]. The initial learning rate is set to be 1e-3, and is divided by 2 after 50 epochs, which allows smoother search near the convergence point. The minibatch size is 16, and the number of base channels is 128.

5.2 Results on TaxiBJ

We first assess the performances of our model and baselines on TaxiBJ with a varying ratio of training data. Tables 3, 4, 5, and 6 report the prediction results. Note that, the variances of the results are almost in the range of 0.000 - 0.002, thus we omit the variances. We summarize the tables with several key observations:

(1) UrbanSTC outperforms all competitive methods across the entire time spans (P1-P4). By comparing to

TABLE 4 The Average RMSE, MAE and MAPE on TaxiBJ Dataset (P2) With Different Proportions of Training Data

| Methods | P2(20%) | | | | P2(40% |) | | P2(60% |) | | P2(80% |) | P2(100%) | | |
|----------|---------|---------|---------|--------|--------|---------|--------|--------|---------|--------|--------|---------|----------|--------|--------|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE |
| MEAN | 26.729 | 15.350 | 5.364 | 26.729 | 15.350 | 5.364 | 26.729 | 15.350 | 5.364 | 26.729 | 15.350 | 5.364 | 26.729 | 15.350 | 5.364 |
| HA | 6.568 | 2.889 | 0.358 | 5.875 | 2.679 | 0.342 | 5.669 | 2.620 | 0.338 | 5.544 | 2.587 | 0.335 | 5.512 | 2.576 | 0.334 |
| SRCNN | 5.613 | 3.201 | 0.837 | 4.994 | 2.855 | 0.706 | 5.172 | 3.036 | 0.801 | 4.924 | 2.839 | 0.713 | 4.978 | 2.896 | 0.748 |
| ESPCN | 5.461 | 3.062 | 0.738 | 5.186 | 2.987 | 0.740 | 4.934 | 2.779 | 0.637 | 4.554 | 2.473 | 0.482 | 5.072 | 2.957 | 0.749 |
| DeepSD | 5.412 | 2.991 | 0.704 | 5.608 | 3.290 | 0.892 | 4.716 | 2.585 | 0.546 | 5.018 | 2.816 | 0.659 | 4.909 | 2.738 | 0.625 |
| VDŜR | 5.449 | 3.024 | 0.727 | 4.753 | 2.608 | 0.561 | 4.954 | 2.795 | 0.660 | 4.494 | 2.444 | 0.492 | 4.429 | 2.402 | 0.475 |
| SRResNet | 5.801 | 3.420 | 0.992 | 4.946 | 2.878 | 0.749 | 4.702 | 2.760 | 0.653 | 4.572 | 2.600 | 0.614 | 4.548 | 2.573 | 0.605 |
| LapSRN | 5.717 | 3.343 | 0.931 | 4.844 | 2.751 | 0.664 | 4.818 | 2.753 | 0.673 | 4.535 | 2.525 | 0.554 | 4.555 | 2.556 | 0.569 |
| IMDN | 5.790 | 3.547 | 1.123 | 4.927 | 2.971 | 0.853 | 4.710 | 2.792 | 0.755 | 4.573 | 2.688 | 0.703 | 4.476 | 2.608 | 0.661 |
| SCN | 5.222 | 2.932 | 0.721 | 4.640 | 2.567 | 0.579 | 4.487 | 2.475 | 0.528 | 4.402 | 2.422 | 0.505 | 4.336 | 2.381 | 0.490 |
| UrbanFM | 5.546 | 2.855 | 0.433 | 4.805 | 2.469 | 0.353 | 4.588 | 2.365 | 0.336 | 4.489 | 2.309 | 0.324 | 4.414 | 2.272 | 0.318 |
| UrbanPy | 5.528 | 2.803 | 0.485 | 4.728 | 2.412 | 0.370 | 4.464 | 2.276 | 0.334 | 4.446 | 2.279 | 0.341 | 4.315 | 2.210 | 0.323 |
| FODE | 5.362 | 2.734 | 0.395 | 4.704 | 2.416 | 0.337 | 4.538 | 2.331 | 0.323 | 4.434 | 2.285 | 0.325 | 4.366 | 2.248 | 0.317 |
| UrbanSTC | 4.975 | 2.424 | 0.297 | 4.454 | 2.231 | 0.294 | 4.347 | 2.185 | 0.288 | 4.274 | 2.157 | 0.288 | 4.225 | 2.136 | 0.288 |
| Δ | +4.73% | +11.34% | +17.04% | +4.01% | +7.50% | +12.76% | +2.62% | +4.00% | +10.84% | +2.91% | +5.35% | +11.11% | +2.09% | +3.35% | +9.15% |

The best results are bold and the second best are underlined.

TABLE 5 The Average RMSE, MAE and MAPE on TaxiBJ Dataset (P3) With Different Proportions of Training Data

| Methods | | P3(20%) |) | | P3(40%) | | | P3(60% |) | | P3(80% |) | P3(100%) | | |
|----------|--------|---------|---------|--------|---------|---------|--------|--------|---------|--------|--------|---------|----------|--------|--------|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE |
| MEAN | 27.442 | 16.029 | 5.612 | 27.442 | 16.029 | 5.612 | 27.442 | 16.029 | 5.612 | 27.442 | 16.029 | 5.612 | 27.442 | 16.029 | 5.612 |
| HA | 5.833 | 2.741 | 0.337 | 5.746 | 2.713 | 0.333 | 5.731 | 2.707 | 0.331 | 5.692 | 2.695 | 0.330 | 5.675 | 2.670 | 0.328 |
| SRCNN | 5.581 | 3.317 | 0.906 | 5.150 | 2.962 | 0.728 | 5.082 | 2.936 | 0.718 | 4.923 | 2.821 | 0.666 | 4.891 | 2.817 | 0.673 |
| ESPCN | 5.273 | 3.013 | 0.717 | 5.043 | 2.848 | 0.638 | 5.091 | 2.888 | 0.656 | 4.796 | 2.668 | 0.556 | 4.853 | 2.716 | 0.579 |
| DeepSD | 5.257 | 2.935 | 0.666 | 5.048 | 2.796 | 0.606 | 4.960 | 2.749 | 0.583 | 4.878 | 2.690 | 0.559 | 4.720 | 2.580 | 0.510 |
| VDŜR | 5.285 | 2.982 | 0.699 | 4.963 | 2.748 | 0.591 | 4.786 | 2.626 | 0.536 | 4.695 | 2.568 | 0.512 | 4.616 | 2.522 | 0.495 |
| SRResNet | 5.578 | 3.352 | 0.945 | 5.120 | 2.998 | 0.776 | 4.934 | 2.857 | 0.705 | 4.773 | 2.734 | 0.643 | 4.658 | 2.648 | 0.602 |
| LapSRN | 5.832 | 3.535 | 1.019 | 5.135 | 2.970 | 0.740 | 5.041 | 2.920 | 0.721 | 4.923 | 2.828 | 0.675 | 4.641 | 2.589 | 0.550 |
| IMDN | 5.635 | 3.493 | 1.077 | 5.143 | 3.107 | 0.876 | 4.908 | 2.930 | 0.788 | 4.745 | 2.794 | 0.715 | 4.690 | 2.765 | 0.704 |
| SCN | 5.090 | 2.899 | 0.694 | 4.826 | 2.702 | 0.601 | 4.670 | 2.593 | 0.549 | 4.575 | 2.531 | 0.522 | 4.514 | 2.494 | 0.506 |
| UrbanFM | 5.299 | 2.738 | 0.379 | 4.951 | 2.558 | 0.350 | 4.761 | 2.456 | 0.336 | 4.656 | 2.408 | 0.330 | 4.578 | 2.356 | 0.314 |
| UrbanPy | 5.342 | 2.827 | 0.529 | 4.946 | 2.532 | 0.382 | 4.743 | 2.443 | 0.362 | 4.578 | 2.346 | 0.332 | 4.436 | 2.272 | 0.318 |
| FODE | 5.165 | 2.686 | 0.380 | 4.875 | 2.521 | 0.347 | 4.712 | 2.434 | 0.331 | 4.616 | 2.387 | 0.327 | 4.536 | 2.345 | 0.319 |
| UrbanSTC | 4.781 | 2.383 | 0.287 | 4.607 | 2.309 | 0.292 | 4.512 | 2.271 | 0.288 | 4.439 | 2.240 | 0.288 | 4.382 | 2.215 | 0.285 |
| Δ | +6.07% | +11.28% | +14.84% | +4.54% | +8.41% | +12.31% | +3.38% | +6.70% | +12.99% | +2.97% | +4.52% | +11.93% | +1.22% | +2.51% | +9.24% |

The best results are bold and the second best are underlined.

current state-of-the-art methods, UrbanSTC has improved 2.51%, 3.80% and 10.51% for RMSE, MAE and MAPE on average on TaxiBJ-P1 with 100.00% training data.

(2) It is apparent that UrbanSTC can achieve the best results when training data decreases. Taking TaxiBJ-P1 (20% training data) for example, UrbanSTC yields 7.10%, 12.24% and 10.91% relative improvements in terms of RMSE, MAE and MAPE, respectively.

The above results show that UrbanSTC has its own advantages in the absence of training data resources. This is consistent with our motivation that spatio-temporal contrastive self-supervision can better learn flow feature representations and improve FUFI performance. Image super-resolution method SCN [56] performers better than other baselines with metric RMSE on 20% - 80% TaxiBJ datasets, while shows deteriorate scores on MAE and MAPE. It is mainly because SCN is a state-of-the-art image super-resolution method with the root mean square loss function. However, most image super-resolution methods are not adapt to the FUFI problem since they do not consider the *structural constraint* when designing models. Compared with UrbanFM, UrbanPy, and FODE, spatio-temporal contrastive learning method UrbanSTC can provide the better latent representations that performs most through all experiments.

5.2.1 Ablation Analysis

To analyses the contribution of each component of UrbanSTC, we analyze the ablation study in this section. We only report the evaluation metrics on TaxiBJ dataset (average result of P1 to P4) because the experimental results on BikeNYC can make similar conclusions. All the results are shown in Table 7. The term "Reg" means the regional-level contrast pre-training; "Inf" illustrates the spatial super-resolution inference network; "TCS" indicates the temporal contrast is used or not.

| TABL | E 6 |
|------|-----|
|------|-----|

The Average RMSE, MAE and MAPE on TaxiBJ Dataset (P4) With Different Proportions of Training Data

| Methods | | P4(20%) |) | | P4(40%) | | | P4(60%) | | | P4(80%) |) | P4(100%) | | |
|----------|--------|---------|---------|--------|---------|---------|--------|---------|---------|--------|---------|--------|----------|--------|--------|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE |
| MEAN | 19.049 | 11.070 | 4.192 | 19.049 | 11.070 | 4.192 | 19.049 | 11.070 | 4.192 | 19.049 | 11.070 | 4.192 | 19.049 | 11.070 | 4.192 |
| HA | 4.306 | 2.067 | 0.319 | 4.238 | 2.052 | 0.319 | 4.209 | 2.043 | 0.319 | 4.223 | 2.045 | 0.320 | 4.201 | 2.039 | 0.320 |
| SRCNN | 4.048 | 2.369 | 0.668 | 4.065 | 2.381 | 0.660 | 3.799 | 2.182 | 0.569 | 3.944 | 2.277 | 0.613 | 3.813 | 2.188 | 0.571 |
| ESPCN | 3.983 | 2.290 | 0.600 | 3.865 | 2.187 | 0.542 | 4.112 | 2.430 | 0.684 | 3.853 | 2.194 | 0.552 | 3.914 | 2.277 | 0.607 |
| DeepSD | 3.980 | 2.240 | 0.562 | 3.910 | 2.181 | 0.527 | 3.924 | 2.215 | 0.552 | 3.806 | 2.121 | 0.511 | 3.662 | 2.030 | 0.472 |
| VDŜR | 3.952 | 2.239 | 0.573 | 3.741 | 2.075 | 0.489 | 3.655 | 2.015 | 0.462 | 3.644 | 2.007 | 0.457 | 3.555 | 1.948 | 0.431 |
| SRResNet | 4.118 | 2.463 | 0.738 | 4.053 | 2.431 | 0.729 | 3.761 | 2.184 | 0.591 | 3.710 | 2.102 | 0.508 | 3.630 | 2.067 | 0.523 |
| LapSRN | 4.467 | 2.753 | 0.884 | 4.150 | 2.489 | 0.745 | 3.705 | 2.103 | 0.530 | 3.673 | 2.080 | 0.520 | 3.679 | 2.118 | 0.544 |
| IMDN | 4.100 | 2.530 | 0.818 | 3.828 | 2.301 | 0.686 | 3.703 | 2.203 | 0.635 | 3.619 | 2.119 | 0.580 | 3.848 | 2.340 | 0.720 |
| SCN | 3.798 | 2.154 | 0.550 | 3.660 | 2.048 | 0.496 | 3.573 | 1.987 | 0.467 | 3.524 | 1.952 | 0.450 | 3.486 | 1.927 | 0.439 |
| UrbanFM | 4.054 | 2.126 | 0.373 | 3.794 | 1.969 | 0.330 | 3.677 | 1.908 | 0.323 | 3.601 | 1.865 | 0.315 | 3.559 | 1.841 | 0.305 |
| UrbanPy | 3.959 | 2.088 | 0.413 | 3.740 | 1.936 | 0.342 | 3.644 | 1.889 | 0.332 | 3.606 | 1.868 | 0.325 | 3.470 | 1.801 | 0.313 |
| FODE | 3.912 | 2.042 | 0.350 | 3.725 | 1.930 | 0.321 | 3.627 | 1.879 | 0.314 | 3.565 | 1.846 | 0.308 | 3.529 | 1.828 | 0.304 |
| UrbanSTC | 3.640 | 1.837 | 0.278 | 3.542 | 1.796 | 0.282 | 3.474 | 1.769 | 0.282 | 3.454 | 1.759 | 0.283 | 3.416 | 1.742 | 0.278 |
| Δ | +4.16% | +10.04% | +12.85% | +3.22% | +6.94% | +11.60% | +2.77% | +5.85% | +10.19% | +1.99% | +4.71% | +8.12% | +1.56% | +3.28% | +8.55% |

The best results are bold and the second best are underlined.

| Regional | Spatial | Temporal | | TaxiBJ | | | | | | |
|----------|----------------------|----------|-------|--------|-------|--|--|--|--|--|
| contrast | super- resolution | contrast | RMSE | MAE | MAPE | | | | | |
| 1 | | | 4.118 | 2.100 | 0.311 | | | | | |
| | 1 | | 4.019 | 2.040 | 0.297 | | | | | |
| | | 1 | 4.008 | 2.027 | 0.290 | | | | | |
| 1 | 1 | | 3.970 | 2.009 | 0.289 | | | | | |
| 1 | | 1 | 3.983 | 2.009 | 0.287 | | | | | |
| | 1 | 1 | 3.975 | 2.008 | 0.288 | | | | | |
| 1 | 1 | 1 | 3.967 | 2.004 | 0.287 | | | | | |

TABLE 7 Ablation Studies

We report the strategies used in different models on TaxiBJ dataset's average results.

We can clearly see that the combination of any two components is better than the single one, which proves the effectiveness of our proposed components. When only considering one strategy, temporal contrast performs better than regional-level contrast and the spatial super-resolution inference network. Spatial contrast contains two components, "Reg" and "Inf". We find that the effect of the spatial super-resolution network (Inf) is better than the regionallevel contrast (Reg). It is mainly because the kernel of the Reg encoder is 1×1 , while that of in Inf encoder is 3×3 , where the larger convolution kernel size helps to capture more information in the encoder. The results of combination of "Reg" + "TCS" and "Inf" + "TCS" are slightly worse than the final model, indicating that such prior knowledge considered both spatial and temporal information is significant for the fine-grained urban flow inference.

To better present the ablation results, we draw some comparable images in Fig. 5. Fig. 5 shows the inference errors $\|\mathbf{X}^f - \hat{\mathbf{X}}^f\|_{1,1}$ generated by UrbanSTC and other ablation parts, where a brighter pixel indicates a large error. A (West TuCheng Road) and B (Sanyuan bridge) are the main traffic arteries in Beijing. It is apparent that UrbanSTC achieves better results than other ablation experiments, which proves that the final structure of the proposed model can better capture the spatio-temporal characteristics of flow data.



Fig. 5. Visualization of the Ablation Study.

5.2.2 End-to-End and Two-Stage Comparison

To verify the effectiveness of the two-stage training process and end-to-end training process, we conduct experiments in TaxiBJ (average result of P1 to P4) and BikeNYC datasets. The end-to-end model integrates three proposed modules, i.e., the coarse-grained flow map is introduced to the spatial self-supervision, temporal self-supervision and external factor learning simultaneously, and optimize these three loss functions integrally. As shown in Table 8, we can clearly find that the two-stage experimental results are better than the end-to-end training process. The end-to-end training method needs to adjust the balance factors between each loss function, while the two-stage training method is not required to adjust the balances among pretexts. The advantage of the self-supervised learning lies in two-stage training. The pretexts help the model in learning the internal characteristics of the data in advance, and the fine-tuning stage then learns the corresponding label information [30], [47], [58], [59].

5.2.3 Temporal Contrastive Sampling Analysis

To evaluate the effect of hard sampling and weight sampling methods, we report the experimental results on TaxiBJ-P1. The tests drawn in Fig. 6 demonstrate that the weight sampling is better than the hard sampling when the proportion of used training data is lower than 60%. This is because the weight sampling method can comprehensively use top K related samples, while the hard sampling only uses the most similar or dissimilar data. With the amount of training data increases, hard sampling begins to show a better performance than weight sampling. When the training dataset is small, we can hardly to pick up the global most similar sample, but use top-K similar samples instead.

TABLE 8 End-to-End and Two-Stage Comparison

| Mehtods | | TaxiBJ | | BikeNYC | | | | | |
|------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--------------------------------|--|--|--|
| | RMSE | MAE | MAPE | RMSE | MAE | MAPE | | | |
| End-to-End Two-stage Δ | 3.980 3.958 0.55% | 2.053 1.998 2.68% | 0.294 0.284 3.40% | 1.120 1.093 2.41% | 0.245 0.236 3.67% | 0.077 0.073 5.19% | | | |



Fig. 6. Performance comparison between Hard Sampling and Weight Sampling on TaxiBJ-P1 dataset.

| P1 | | | | P2 | | | P3 | | P4 | | | |
|-----------------------|--|--|---|--|---|--|---|--|---|--|--|--|
| RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | RMSE | MAE | MAPE | |
| 3.970 | 2.023 | 0.334 | 4.355 | 2.239 | 0.317 | 4.530 | 2.335 | 0.321 | 3.528 | 1.824 | 0.303 | |
| 3.909 | 1.981 | 0.330 | 4.353 | 2.230 | 0.327 | 4.466 | 2.294 | 0.323 | 3.498 | 1.817 | 0.317 | |
| 3.915 | 1.996 | 0.332 | 4.348 | 2.235 | 0.316 | 4.505 | 2.329 | 0.314 | 3.505 | 1.821 | 0.311 | |
| 3.845 3.841 | 1.922 1.917 | 0.298 0.292 | 4.225 4.209 | 2.136 2.125 | 0.288 0.284 | 4.382 4.376 | 2.215 2.210 | 0.285 0.283 | 3.416 3.404 | 1.742 1.738 | 0.278 0.275 | |
| | RMSE 3.970 3.909 3.915 3.845 3.841 | P1 RMSE MAE 3.970 2.023 3.909 1.981 3.915 1.996 3.845 1.922 3.841 1.917 | P1 RMSE MAE MAPE 3.970 2.023 0.334 3.909 1.981 0.330 3.915 1.996 0.332 3.845 1.922 0.298 3.841 1.917 0.292 | P1 RMSE MAE MAPE RMSE 3.970 2.023 0.334 4.355 3.909 1.981 0.330 4.353 3.915 1.996 0.332 4.348 3.845 1.922 0.298 4.225 3.841 1.917 0.292 4.209 | P1 P2 RMSE MAE MAPE RMSE MAE 3.970 2.023 0.334 4.355 2.239 3.909 1.981 0.330 4.353 2.230 3.915 1.996 0.332 4.348 2.235 3.845 1.922 0.298 4.225 2.136 3.841 1.917 0.292 4.209 2.125 | P1 P2 RMSE MAE MAPE RMSE MAE MAPE 3.970 2.023 0.334 4.355 2.239 0.317 3.909 1.981 0.330 4.353 2.230 0.327 3.915 1.996 0.332 4.348 2.235 0.316 3.845 1.922 0.298 4.225 2.136 0.288 3.841 1.917 0.292 4.209 2.125 0.284 | P1 P2 RMSE MAE MAPE RMSE MAE MAPE RMSE 3.970 2.023 0.334 4.355 2.239 0.317 4.530 3.909 1.981 0.330 4.353 2.230 0.327 4.466 3.915 1.996 0.332 4.348 2.235 0.316 4.505 3.845 1.922 0.298 4.225 2.136 0.288 4.382 3.841 1.917 0.292 4.209 2.125 0.284 4.376 | P1 P2 P3 RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAPE RMSE MAE 1.30 2.335 2.335 2.335 2.335 3.909 1.981 0.330 4.353 2.230 0.327 4.466 2.294 3.915 1.996 0.332 4.348 2.235 0.316 4.505 2.329 3.845 1.922 0.298 4.225 2.136 0.288 4.382 2.215 3.841 1.917 0.292 4.209 2.125 0.284 4.376 2.210 | P1 P2 P3 RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAPE 3.970 2.023 0.334 4.355 2.239 0.317 4.530 2.335 0.321 3.909 1.981 0.330 4.353 2.230 0.327 4.466 2.294 0.323 3.915 1.996 0.332 4.348 2.235 0.316 4.505 2.329 0.314 3.845 1.922 0.298 4.225 2.136 0.288 4.382 2.215 0.285 3.841 1.917 0.292 4.209 2.125 0.284 4.376 2.210 0.283 | P1 P2 P3 RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAPE RMSE RMSE MAE RMSE MAE MAPE RMSE RMSE MAE RMSE MAE MAPE RMSE 3.528 3.528 3.528 3.528 3.528 3.498 3.505 3.498 3.505 3.845 1.922 0.298 4.225 2.136 0.288 4.382 2.215 0.285 3.416 3.841 1.917 0.292 4.209 2.125 0.284 4.376 2.210 0.283 3.404 | P1 P2 P3 P4 RMSE MAE MAPE RMSE MAE MAPE RMSE MAE MAE | |

TABLE 9 The Average RMSE. MAE and MAPE on TaxiBJ Dataset With External Factors

Note that "+E" represents a model that incorporates external factors. The best results are bold.

Otherwise, if the most similar sample is found with the training data increasing, the hard sampling method can achieve the better result. Therefore, a combination of two methods can be adopted in different training scenarios.

5.2.4 Study on External Factor Fusion

In reality, there are complicated external factors in the FUFI problem. In order to verify the effectiveness of the external information in our method, we introduce external factors and conduct experiments on TaxiBJ datasets with different time spans (P1-P4). We only compare our method with available baselines. As test shown in Table 9, we clearly see that UrbanSTC+E performs better than other models across all time spans, which reveals that the combination of our UrbanSTC and external factors can improve the model performance. Note that, even some compared FUFI methods have the well-designed external information fusion module, our proposed method UrbanSTC can leverage external information with a simple network.

5.2.5 Configurations and Parameters Analysis

In this section, we try to explore the learning abilities of our method in various setting environments. Compared with different channels (32, 64, 128), we can get the results shown in Fig. 8. Fig. 8a illustrates that the larger number of channels, the better performance of UrbanSTC. Besides, Figs. 8b, 8c and 8d show that a larger number of channels can improve the efficiency of the learning convergence.

We analyze the influence of λ in the regional-level contrastive learning. Fig. 7a shows the different performances with a varying setting of λ . The regional-level contrast judges which regions are positive and negative samples based on the threshold λ . The experimental result shows that the best result is achieved when the threshold is 1e-4 on the taxi dataset. Fig. 7b represents that $\lambda = 5e-5$ yields the best performance.



Fig. 7. Effect of λ . We explore the influence of λ in the spatial contrastive learning.

For the parameter analysis, Fig. 9a represents that UrbanSTC can get better results than other models with different training data fractions. Fig. 9b indicates the traditional image super-resolution methods, e.g., IMDN, VDSR and SRResNet are not suitable for the FUFI problem due to the inherent difference. Although SRResNet and UrbanFM have similar structures, the M^2 -Normalization layer in UrbanFM contributes to the FUFI problem. UrbanPy uses a cascading model for forecasting fine-grained urban flows by decomposing the original task into multiple subtasks, which leads to the increase of computing complexity. FODE utilizes ODE module to replace the ResNet strucure in the UrbanFM. Because the above modules can be viewed as a discretization of a continuous ODE operator, which greatly improves the convergence speed and reduces the number of parameters.



Fig. 8. Study on Configurations. The convergence rate loss error of the self-supervised module under different channel dimensions.



Fig. 9. Study on Parameters. Experiments on the P1 dataset with different training data fractions and the comparison of parameter cost.

TABLE 10 Efficiency Evaluated on the P1 Dataset

| Method | Params | Training Time | Inference Time | Total Time | RMSE | |
|-------------|---------|---------------|----------------|------------|-------|--|
| VDSR 4.79 M | | 4.37 s | 0.76 s | 9.10mins | 4.054 | |
| SRResNet | 5.79 M | 11.4 s | 1.57 s | 33.25mins | 4.079 | |
| IMDN | 2.63 M | 4.21 s | 0.69 s | 13.33mins | 4.085 | |
| SCN | 18.55 M | 23.86 s | 5.21 s | 59.65mins | 3.965 | |
| UrbanFM | 5.94 M | 12.28 s | 1.80 s | 10.23mins | 4.042 | |
| UrbanPv | 11.28 M | 27.39 s | 11.89 s | 79.89mins | 3.944 | |
| FODE | 4.23 M | 14.05 s | 1.91 s | 17.56mins | 4.002 | |
| Reg | 0.03 M | 3.89 s | - | 6.48mins | - | |
| Inf | 1.41 M | 0.80 s | - | 1.60mins | - | |
| TCS | 0.16 M | 0.99 s | - | 1.65mins | - | |
| Fine-tuning | 1.98 M | 3.34 s | 0.66 s | 2.78mins | 3.845 | |
| UrbanSTC | 3.58 M | 9.02 s | 0.66 s | 12.51mins | 3.845 | |

The entire model UrbanSTC and its four components have been tested separately.

For our model UrbanSTC, we design several self-supervised pretext tasks to make encoders rich in spatio-temporal information. As shown in Table 10, "Reg" indicates the regionallevel contrast pre-training; "Inf" illustrates the spatial superresolution inference network; "TCS" denotes the temporal contrast. UrbanSTC consists of three self-supervised modules and a fine-tuning stage. The parameter cost of UrbanSTC is slightly higher than IMDN because the latter is a lightweight image super-resolution method designed in mobile devices. Based on our well-designed self-supervised tasks, UrbanSTC can capture spatio-temporal knowledge in advance and perform better than other baselines with a relatively small amount of parameters.

We further conduct a comparison between UrbanSTC and baselines in terms of the training time and inference time. We reported the training time of each epoch and the total training time until model convergence in P1 dataset of TaxiBJ, which contains 1530 training snapshots and 765 test snapshots respectively. Even our model contains two stages, the training time of each epoch is less than all previous FUFI models (UrbanFM, UrbanPy and FODE) as shown in Table 10. It is mainly because the structure of proposed encoders is simple while well-designed that can capture rich spatio-temporal characteristics in advance. Two image super-resolution methods, VDSR and IMDN are efficient in the training process, yet their performances are far more worse than our method.

As shown in Fig. 10, UrbanSTC can efficiently converge with a small number of epochs. Even UrbanSTC spent slightly more total training time than UrbanFM and VDSR, it is efficient with the best results achieved. In summary, extensive experiments demonstrate that UrbanSTC can achieve the best results efficiently by using a small amount of parameters.

5.3 Results on BikeNYC

Table 11 presents the comparison results on the BikeNYC dataset. Since we cannot get the external factors of this dataset, we will do not add such information in the experiments. In this experiment, the baseline DeepSD will be the same as SRCNN when M is $2 \times$, therefore we remove the DeepSD.

BikNYC dataset is more sparse than TaxiBJ dataset. Nonetheless, UrbanSTC still yields 2.93% and 5.60% improvements on average in terms of RMSE and MAE, respectively. Note that due to the extremely sparsity of BikeNYC dataset, the metric MAPE is not available. It is apparent that the experimental results lead to similar conclusions to the test on TaxiBJ. The proposed model outperforms other baseline methods on both sparse and dense datasets, which has a good robustness.

5.4 Visualization of Fine-Grained Flow Prediction

Fig. 11 gives an intuitive presentation of the fine-grained urban flow prediction in BikeNYC data. Fig. 11a represents the coarse-grained crowd flows and 11b is the ground-truth of the fine-grained flow map from (a), and (c) is our prediction result. This visualization illustrates the effectiveness of our model.

Fig. 12 shows the inference errors $\|\mathbf{X}^f - \hat{\mathbf{X}}^f\|_{1,1}$ generated by UrbanSTC and the other three baselines for a sample at the 4× task, where a brighter pixel indicates a large error. Overall, UrbanSTC has obtained more detailed inference effects and less global error. To better visualize the quality of inference, we select four busy subregions (A, B, C and D) where the UrbanSTC performs better than other methods



Fig. 10. Convergence rate of different methods.

| | | 0 <i>i</i> | | | | | • | | 0 | | |
|----------|--------------|-------------------|--------------|--------|--------|--------------|--------|--------------|--------|---------------|--|
| Methods | BikeNYC(20%) | | BikeNYC(40%) | | BikeNY | BikeNYC(60%) | | BikeNYC(80%) | | BikeNYC(100%) | |
| | RMSE | MAE | RMSE | MAE | RMSE | MAE | RMSE | MAE | RMSE | MAE | |
| MEAN | 3.776 | 1.281 | 3.776 | 1.281 | 3.776 | 1.281 | 3.776 | 1.281 | 3.776 | 1.281 | |
| HA | 1.498 | 0.359 | 1.476 | 0.355 | 1.511 | 0.365 | 1.506 | 0.364 | 1.502 | 0.364 | |
| SRCNN | 1.419 | 0.452 | 1.306 | 0.421 | 1.228 | 0.373 | 1.201 | 0.364 | 1.262 | 0.413 | |
| ESPCN | 1.458 | 0.489 | 1.432 | 0.495 | 1.302 | 0.402 | 1.322 | 0.451 | 1.295 | 0.411 | |
| VDSR | 1.888 | 0.838 | 1.740 | 0.758 | 1.616 | 0.700 | 1.531 | 0.665 | 1.476 | 0.626 | |
| SRResNet | 1.843 | 0.891 | 1.713 | 0.781 | 1.607 | 0.712 | 1.488 | 0.643 | 1.443 | 0.600 | |
| LapSRN | 1.582 | 0.635 | 1.448 | 0.550 | 1.392 | 0.516 | 1.339 | 0.492 | 1.320 | 0.464 | |
| IMDN | 1.407 | 0.521 | 1.345 | 0.456 | 1.292 | 0.447 | 1.241 | 0.422 | 1.220 | 0.402 | |
| SCN | 1.331 | 0.424 | 1.276 | 0.404 | 1.191 | 0.356 | 1.200 | 0.362 | 1.162 | 0.332 | |
| UrbanFM | 1.405 | 0.316 | 1.302 | 0.309 | 1.215 | 0.283 | 1.215 | 0.265 | 1.172 | 0.263 | |
| UrbanPy | 1.381 | 0.315 | 1.310 | 0.301 | 1.271 | 0.286 | 1.200 | 0.273 | 1.126 | 0.250 | |
| FODE | 1.293 | 0.302 | 1.214 | 0.279 | 1.167 | 0.265 | 1.146 | 0.258 | 1.134 | 0.253 | |
| UrbanSTC | 1.267 | 0.276 | 1.191 | 0.257 | 1.146 | 0.246 | 1.107 | 0.239 | 1.093 | 0.236 | |
| Δ | +2.01% | +8.61% | +1.89% | +7.89% | +1.80% | +7.17% | +3.40% | +7.36% | +2.93% | +5.60% | |

TABLE 11 The Average RMSE, MAE and MAPE on BikeNYC Dataset With Different Proportions of Training Data

The best results are bold and the second best are underlined.



(a) Coarse-grained Crowd Flows

(b) Fine-grained Crowd Flows

(c) UrbanSTC Fine-grained Inference

Fig. 11. Visualization of crowd flows in BikeNYC.

obviously. Area A is the Sanyuan bridge (the main entrance to downtown); area B is the Beijing zoo (a large number of tourists); areas C and D cover the Beijing and Beijing west railway stations. Compared with existing FUFI methods, we observe that UrbanSTC has made great improvements in the above areas. Besides, UrbanSTC shows a darker tone than other methods from the heat map, which corresponds to the quantitive results from Table 3.



Fig. 12. Visualization for inference errors among different methods on P1 dataset. Best view in color.

6 CONCLUSION

In this paper, we propose a spatio-temporal contrastive selfsupervision method named UrbanSTC for the fine-grained urban flow inference problem. Our model can extract rich spatial and temporal characteristics from urban flows. In detail, we establish self-supervision pretext tasks from two aspects, that are spatial and temporal correlations. For the spatial correlation, regional contrast and spatial super-resolution inference network make great contributions to capture similarities among regional-level flows and upscaling patterns. Moreover, we devise two sampling strategies based on temporal attributes. The overall architecture of our model obeys the self-supervised training mode: pre-training & fine-tuning. Through welldesigned self-supervised tasks, uncomplicated networks have a strong ability to learn high-level representations from urban flows. We conduct intensive experiments on two real-world datasets to compare the performances between UrbanSTC and other state-of-the-art approaches. The results not only show that our approach outperforms all other methods, but also represent a high performance when the training data decrease.

ACKNOWLEDGMENTS

Hao Qu and Yongshun Gong contributed equally to this work.

REFERENCES

 Z. Li, J. Zhang, Q. Wu, Y. Gong, J. Yi, and C. Kirsch, "Sample adaptive multiple kernel learning for failure prediction of railway points," in *Proc. 25th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2019, pp. 2848–2856.

- Y. Gong, Z. Li, J. Zhang, W. Liu, and Y. Zheng, "Online spatio-[2] temporal crowd flow distribution prediction for complex metro system," IEEE Trans. Knowl. Data Eng., vol. 34, no. 2, pp. 865-880, Feb. 2020.
- Y. Liang et al., "UrbanFM: Inferring fine-grained urban flows," in [3] Proc. 25th ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2019, pp. 3132-3142.
- F. Zhou, L. Li, T. Zhong, G. Trajcevski, K. Zhang, and J. Wang, [4] "Enhancing urban flow maps via neural odes," in Proc. 29th Int. Joint Conf. Artif. Intell., 2020, pp. 1295-1302.
- [5] K. Ouyang et al., "Fine-grained urban flow inference," IEEE Trans. Knowl. Data Eng., vol. 34, no. 6, pp. 2755–2770, Jun. 2022
- [6] J. Zhang, Y. Zheng, D. Qi, R. Li, X. Yi, and T. Li, "Predicting citywide crowd flows using deep spatio-temporal residual networks," Artif. Intell., vol. 259, pp. 147-166, 2018.
- W. S. H. M. W. Ahmad et al., "5G technology: Towards dynamic [7] spectrum sharing using cognitive radio networks," IEEE Access, vol. 8, pp. 14460-14488, 2020.
- Y. Zheng, L. Capra, O. Wolfson, and H. Yang, "Urban computing: Concepts, methodologies, and applications," ACM Trans. Intell. Syst. Technol., vol. 5, no. 3, pp. 1–55, 2014.
- Y. Gong, Z. Li, J. Zhang, W. Liu, and J. Yi, "Potential passenger [9] flow prediction: A novel study for urban transportation development," in Proc. AAAI Conf. Artif. Intell., 2020, pp. 4020-4027.
- [10] S. Sarkar et al., "Effective urban structure inference from traffic flow dynamics," IEEE Trans. Big Data, vol. 3, no. 2, pp. 181-193, Jun. 2017.
- [11] J. Gutiérrez Bayo, "International case studies of smart cities: Santander, Spain," Inter-Amer. Develop. Bank, pp. 1-52, 2016.
- C. Schreiner, "International case studies of smart cities: Rio de janeiro, Brazil," *Inter-Amer. Develop. Bank*, pp. 1–72, 2016. [12]
- C. Doersch, A. Gupta, and A. A. Efros, "Unsupervised visual [13] representation learning by context prediction," in Proc. IEEE Int. Conf. Comput. Vis., 2015, pp. 1422–1430. [14] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros,
- "Context encoders: Feature learning by inpainting," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2016, pp. 2536-2544.
- R. Y. Zhang et al., "Real-time user-guided image colorization with learned deep priors," ACM Trans. Graph., vol. 36, no. 4, 2017, [15] Art. no. 119.
- [16] S. Gidaris, P. Singh, and N. Komodakis, "Unsupervised representation learning by predicting image rotations," in Int. Conf. Learn. Representations, 2018, pp. 1–16.
- [17] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," 2013, arXiv:1301.3781.
- [18] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pretraining of deep bidirectional transformers for language understanding," 2018, arXiv:1810.04805.
- [19] Y. Gong, Z. Li, J. Zhang, W. Liu, Y. Yin, and Y. Zheng, "Missing value imputation for multi-view urban statistical data via spatial correlation learning," IEEE Trans. Knowl. Data Eng., early access, Apr. 13, 2021, doi: 10.1109/TKDE.2021.3072642
- [20] Y. Gong, Z. Li, J. Zhang, W. Liu, B. Chen, and X. Dong, "A spatial missing value imputation method for multi-view urban statistical data," in Proc. 29th Int. Conf. Int. Joint Conf. Artif. Intell., 2021, pp. 1310–1316.
- [21] J. Cai, S. Gu, R. Timofte, and L. Zhang, "NTIRE 2019 challenge on real image super-resolution: Methods and results," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops, 2019, pp. 2211-2223.
- [22] Z. Wang, J. Chen, and S. C. Hoi, "Deep learning for image super-resolution: A survey," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 10, pp. 3365–3387, Oct. 2020.
 [23] C. Ledig et al., "Photo-realistic single image super-resolution
- using a generative adversarial network," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 105–114. [24] R. Shen, J. Xu, Q. Bao, W. Li, H. Yuan, and M. Xu, "Fine-grained
- urban flow prediction via a spatio-temporal super-resolution scheme," in Proc. Asia-Pacific Web Web-Age Inf. Manage. Joint Int. Conf. Web Big Data, 2020, pp. 360-375.
- [25] R. T. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, "Neural ordinary differential equations," in Proc. 32nd Int. Conf. Neural Inf. Process. Syst., 2018, pp. 6572–6583. P. Sermanet et al., "Time-contrastive networks: Self-supervised
- [26] learning from video," in Proc. IEEE Int. Conf. Robot. Automat., 2018, pp. 1134-1141.

- [27] P. Sermanet, C. Lynch, J. Hsu, and S. Levine, "Time-contrastive networks: Self-supervised learning from multi-view observation," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops, 2017, pp. 486-487.
- [28] X. Wang and A. Gupta, "Unsupervised learning of visual representations using videos," in Proc. IEEE Int. Conf. Comput. Vis., 2015, pp. 2794-2802.
- [29] R. D. Hjelm et al., "Learning deep representations by mutual information estimation and maximization," in Proc. Int. Conf. Learn. Representations, 2018, pp. 1-24.
- [30] A. v. d. Oord, Y. Li, and O. Vinyals, "Representation learning with contrastive predictive coding," 2018, arXiv:1807.03748.
- [31] K. Nazeri, H. Thasarathan, and M. Ebrahimi, "Edge-informed single image super-resolution," in Proc. IEEE/CVF Int. Conf. Comput. Vis. Workshops, 2019, pp. 3275-3284.
- [32] W. Shi et al., "Real-time single image and video super-resolution using an efficient sub-pixel convolutional neural network," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2016, pp. 1874-1883.
- [33] M. W. Thornton, P. M. Atkinson, and D. Holland, "Sub-pixel mapping of rural land cover objects from fine spatial resolution satellite sensor imagery using super-resolution pixel-swapping," Int. J. Remote Sens., vol. 27, no. 3, pp. 473–491, 2006.
- [34] R. Keys, "Cubic convolution interpolation for digital image processing," IEEE Trans. Acoust., Speech, Signal Process., vol. 29, no. 6, pp. 1153–1160, Dec. 1981. C. E. Duchon, "Lanczos filtering in one and two dimensions,"
- [35] J. Appl. Meteorol. Climatol., vol. 18, no. 8, pp. 1016–1022, 1979.
- [36] S. Dai, M. Han, W. Xu, Y. Wu, Y. Gong, and A. K. Katsaggelos, "SoftCuts: A soft edge smoothness prior for color image superresolution," IEEE Trans. Image Process., vol. 18, no. 5, pp. 969-981, May 2009.
- [37] J. Sun, Z. Xu, and H.-Y. Shum, "Image super-resolution using gradient profile prior," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2008, pp. 1-8.
- [38] Q. Yan, Y. Xu, X. Yang, and T. Q. Nguyen, "Single image superresolution based on gradient profile sharpness," IEEE Trans. Image Process., vol. 24, no. 10, pp. 3187-3202, Oct. 2015.
- H. Chang, D.-Y. Yeung, and Y. Xiong, "Super-resolution through neighbor embedding," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, 2004, pp. I–I. [39]
- [40] K. Zhang, D. Tao, X. Gao, X. Li, and J. Li, "Coarse-to-fine learning for single-image super-resolution," IEEE Trans. Neural Netw. Learn. Syst, vol. 28, no. 5, pp. 1109-1122, May 2017.
- [41] C. Deng, J. Xu, K. Zhang, D. Tao, X. Gao, and X. Li, "Similarity constraints-based structured output regression machine: An approach to image super-resolution," IEEE Trans. Neural Netw. Learn. Syst., vol. 27, no. 12, pp. 2472-2485, Dec. 2015.
- [42] W. Yang, Y. Tian, F. Zhou, Q. Liao, H. Chen, and C. Zheng, "Consistent coding scheme for single-image super-resolution via independent dictionaries," IEEE Trans. Multimedia, vol. 18, no. 3, pp. 313–325, Mar. 2016. [43] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution
- using deep convolutional networks," IEEE Trans. Pattern Anal. Mach. Intell., vol. 38, no. 2, pp. 295-307, Feb. 2016.
- [44] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, arXiv:1409.1556.
- [45] J. Kim, J. K. Lee, and K. M. Lee, "Accurate image super-resolution using very deep convolutional networks," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2016, pp. 1646–1654.
- [46] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in Proc. Int. Conf. Mach. Learn., 2015, pp. 448–456. [47] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple
- framework for contrastive learning of visual representations," in Proc. Int. Conf. Mach. Learn., 2020, pp. 1597-1607.
- [48] A. Miech, J.-B. Alayrac, L. Smaira, I. Laptev, J. Sivic, and A. Zisserman, "End-to-end learning of visual representations from uncurated instructional videos," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., 2020, pp. 9879-9889.
- [49] Y. Gong, Z. Li, J. Zhang, W. Liu, Y. Zheng, and C. Kirsch, "Network-wide crowd flow prediction of sydney trains via customized online non-negative matrix factorization," in Proc. 27th ACM Int. Conf. Inf. Knowl. Manage., 2018, pp. 1243–1252.
- [50] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A unified embedding for face recognition and clustering," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2015, pp. 815-823.

- [51] C.-Y. Wu, R. Manmatha, A. J. Smola, and P. Krahenbuhl, "Sampling matters in deep embedding learning," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 2859–2867.
- [52] Y. Gong, J. Yi, D.-D. Chen, J. Zhang, J. Zhou, and Z. Zhou, "Inferring the importance of product appearance with semisupervised multi-modal enhancement: A step towards the screenless retailing," in *Proc. 29th ACM Int. Conf. Multimedia*, 2021, pp. 1120–1128.
- [53] T. Vandal, E. Kodra, S. Ganguly, A. Michaelis, R. Nemani, and A. R. Ganguly, "Deepsd: Generating high resolution climate change projections through single image super-resolution," in *Proc. 23rd ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining*, 2017, pp. 1663–1672.
- pp. 1663–1672.
 [54] W.-S. Lai, J.-B. Huang, N. Ahuja, and M.-H. Yang, "Deep laplacian pyramid networks for fast and accurate super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 5835–5843.
- [55] Z. Hui, X. Gao, Y. Yang, and X. Wang, "Lightweight image superresolution with information multi-distillation network," in *Proc.* 27th ACM Int. Conf. Multimedia, 2019, pp. 2024–2032.
- [56] Y. Fan, J. Yu, D. Liu, and T. S. Huang, "Scale-wise convolution for image restoration," in *Proc. AAAI Conf. Artif. Intell.*, 2020, pp. 10770–10777.
- [57] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, arXiv:1412.6980.
- [58] O. Sumer, T. Dencker, and B. Ommer, "Self-supervised learning of pose embeddings from spatiotemporal relations in videos," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 4308–4317.
- [59] J. Yang, J. M. Alvarez, and M. Liu, "Self-supervised learning of depth inference for multi-view stereo," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 7522–7530.



Hao Qu is currently working toward the postgraduate degree with the School of Software, Shandong University, China. His Principal research interest covers the data science and machine learning, in particular, the following areas: traffic analysis; spatio-temporal data mining; natural language processing. He has published 2 papers in *Computer Science journal.*



Yongshun Gong (Member, IEEE) received the PhD degree from the University of Technology Sydney, in 2021. He is an associate professor with the School of Software, Shandong University, China. His Principal research interest covers the data science and machine learning, in particular, the following areas: spatio-temporal data mining and traffic prediction. He has published more than 30 papers in top journals and refereed conference proceedings, including the *IEEE Transactions on Knowledge and Data Engineering*,

IEEE Transactions on Neural Networks and Learning Systems, IEEE Transactions on Cybernetics, NeurIPS, KDD, CIKM, AAAI and IJCAI.



Meng Chen (Member, IEEE) received the PhD degree in computer science and technology from Shandong University, China, in 2016. He worked as a postdoctoral fellow from 2016 to 2018 with the School of Information Technology, York University, Canada. He is currently an assistant professor with the School of Software, Shandong University, China. His research interest is in the area of trajectory data mining and traffic management.



Junbo Zhang (Member, IEEE) is a Senior Researcher of JD Intelligent Cities Research. He is leading the Urban AI Product Department of JD iCity at JD Technology, as well as AI Lab of JD Intelligent Cities Research. His research interests include Spatio-Temporal Data Mining and AI, Urban Computing, Deep Learning, Federated Learning. He has published more than 50 research papers (e.g., *AI Journal, IEEE Transactions on Knowledge and Data Engineering*, KDD, AAAI, IJCAI, WWW, ACL, Ubi-Comp) in refereed journals and conferences. He

received the ACM Chengdu Doctoral Dissertation Award, in 2016, the Chinese Association for Artificial Intelligence (CAAI) Excellent Doctoral Dissertation Nomination Award, in 2016, the Si Shi Yang Hua Medal of SWJTU, in 2012, and the Outstanding PhD Graduate of Sichuan Province, in 2013.



Yu Zheng (Fellow, IEEE) is the cice president of JD. COM and the chief data scientist with JD Digits, passionate about using Big Data and AI technology to tackle urban challenges. Before Joining JD.COM, he was a senior research manager with Microsoft Research. He is also a chair professor with Shanghai Jiao Tong University. He currently serves as the editor-in-chief of ACM Transactions on Intelligent Systems and Technology and has served as chair on more than 10 prestigious international conferences. He is also a keynote speaker of AAAI 2019,

KDD 2019 Plenary Keynote Panel and IJCAI 2019 Industrial Days. His monograph, entitled Urban Computing, has been used as the first text book in this field. In 2013, he was named one of the Top Innovators under 35 by MIT Technology Review (TR35) and featured by Time Magazine for his research on urban computing. In 2017, he is honored as an ACM distinguished scientist.



Yilong Yin received the PhD degree from Jilin University, Changchun, China, in 2000. From 2000 to 2002, he was a postdoctoral fellow with the Department of Electronic Science and Engineering, Nanjing University, Nanjing, China. He is the director of the Machine Learning and Applications Group and a professor with Shandong University, Jinan, China. His research interests include machine learning, data mining, computational medicine, and biometrics.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/csdl.