

Urban Data Scarcity: from Sparse Crowdsensing to Urban Transfer Learning

Leye Wang, Peking University

leyewang@pku.edu.cn, <https://sites.google.com/site/wangleye/>

2019.08.17



"DATA IS THE NEW OIL"

Since the beginning of recorded time until 1990, we created **5 exabytes** of data.

In 2011 the same amount was created every two days.

By 2020, it's expected that the time will shrink to 10 minutes.

7 billion DVDs.

That's the size of what's lost in a single day of a disaster.

Google's 2009 data center footprint is larger than any other company's. The company plans to spend \$10 billion on data centers in 2011.

English is the dominant language of the world. But by 2050, it will be Chinese.



267 billion EMAILS

are sent every day. (That's the daily equivalent of 100 million letters.)

4 million websites have been created since 2000.

As of August 2010, there were just over 100 million websites.

There are **111 million** BLOGS on the web.

80% of all humans use a mobile phone or some form of mobile device.

7 billion and counting. In Singapore, 80% of residents are using mobile phones.

100% of all humans use a mobile phone or some form of mobile device.

In 2010, 100,000 new websites were created every day.

High-frequency traders use the edge of computer algorithms, not big time to make money, and to outpace the market.

Just an increase of activity on Twitter game accounts, Twitter accounts, and Facebook accounts, meaning it's about the economy, not just about the market.

These specialized algorithms create new money, but it's not as much as you think. They're not as much as you think.

1 million new websites are created every day. (That's the daily equivalent of 100 million letters.)

400 new fiber-optic cables are laid every day. (That's the daily equivalent of 100 million letters.)

They're not as much as you think. They're not as much as you think.

How they save 3 million new websites are created every day. (That's the daily equivalent of 100 million letters.)

The new cables will be an order of magnitude faster than the current cables. (That's the daily equivalent of 100 million letters.)

The new cables will be an order of magnitude faster than the current cables.

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

50% of 5-year old data in the U.S. are gone. (That's the daily equivalent of 100 million letters.)

But data is not always available

Urban Data Scarcity

Collect New Data

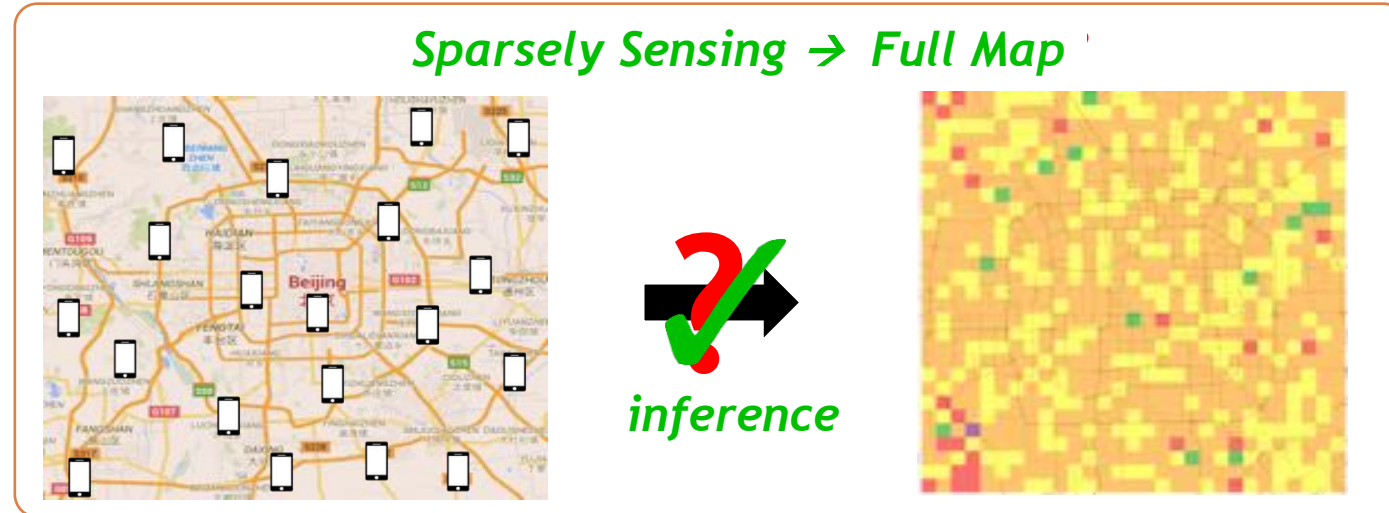
Exploit Existing Data

Sparse Crowdsensing

Urban Transfer Learning

Smart City

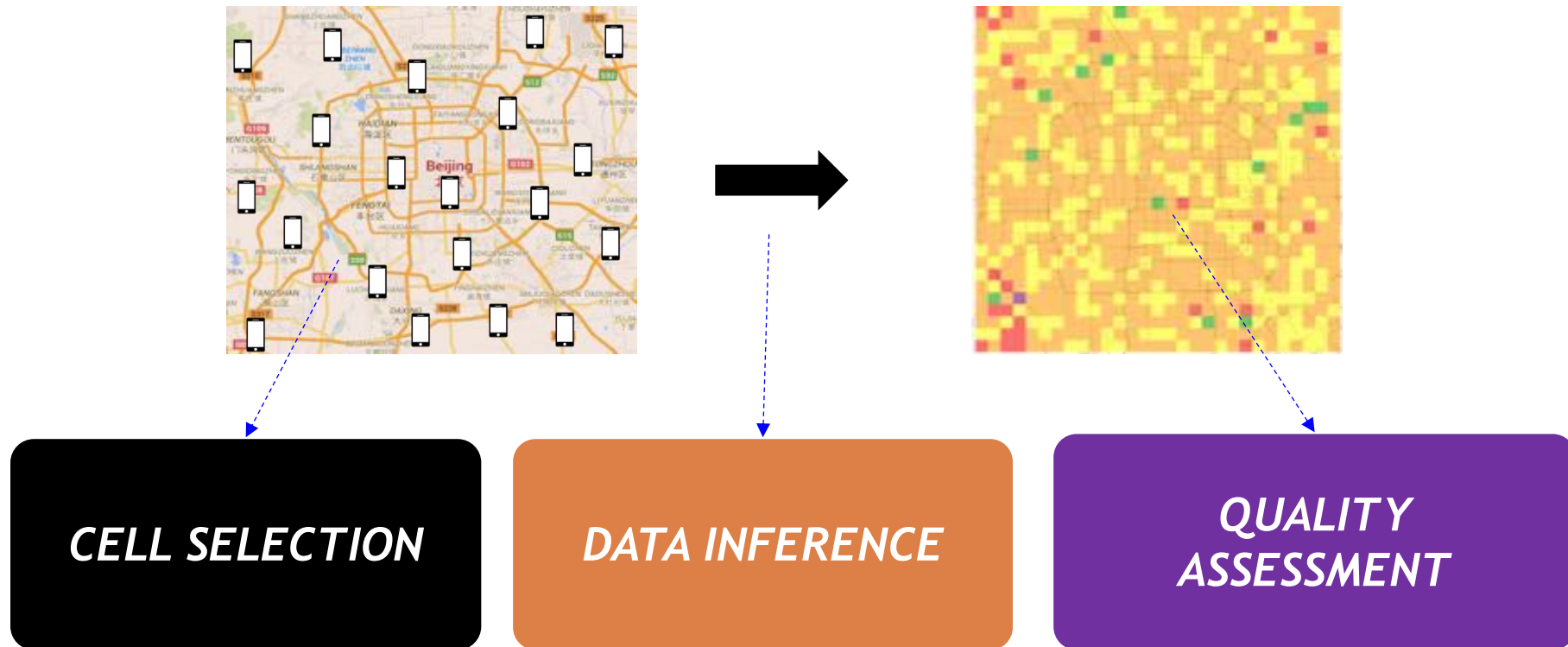
Sparse Crowdsensing



Sense Less, Infer More

	Traditional MCS	Sparse MCS
High Quality	high/full coverage	high inference accuracy (probably sparse coverage)

Key Issues in Sparse MCS



Three-step framework

Area: fixed-size cells (e.g. 100m*100m); *Duration*: fixed-length cycle

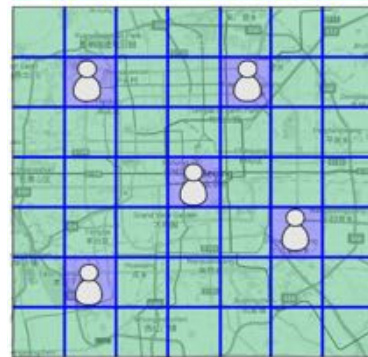
Step 3. Quality Assessment

Use collected data to estimate inference error of unsensed cells.



Step 1: Cell Selection

after five iterations



Step 2: Data Inference



Step 3: Quality Assessment



Data Inference

- Feature of Environment Data
 - Spatial and temporal correlations*
 - Low-rank property in sensing matrix*
- Spatio-Temporal Compressive Sensing (STCS)*** considers *spatial and temporal correlations*, and *low-rank property* all together.

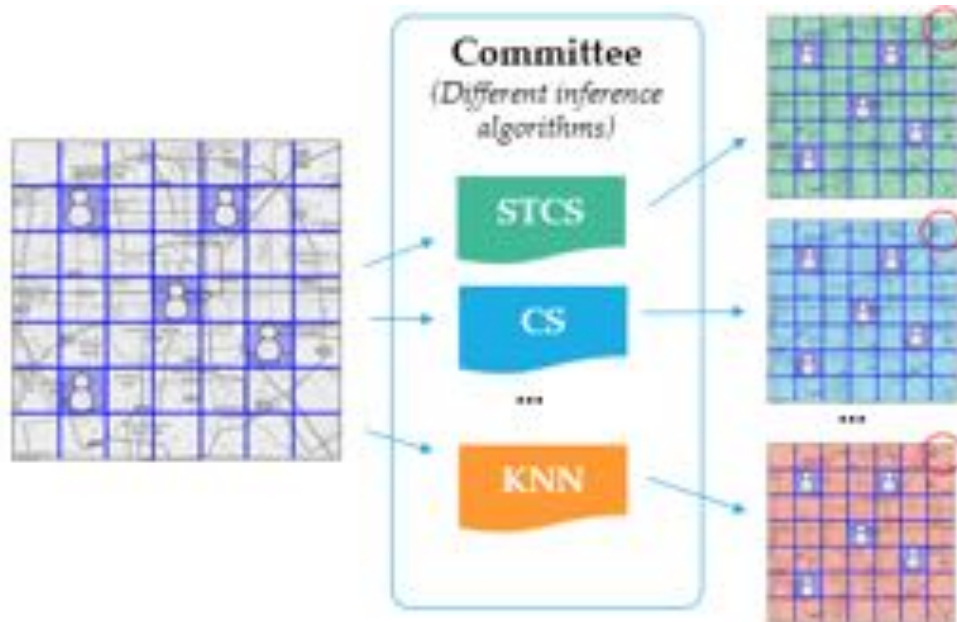
sensing matrix: entry $[i, j]$ means the sensing data of cell i in cycle j .

$$\begin{array}{c} \text{cycles} \\ \text{cells} \end{array} \begin{array}{c} \text{sensing} \\ \text{matrix} \end{array} = \begin{array}{c} L \end{array} \begin{array}{c} R^T \end{array}$$

$$\begin{aligned} \min \quad & \lambda_r (\|L\|_F^2 + \|R\|_F^2) + \|LR^T \circ S - C\|_F^2 \\ & + \lambda_s \|\mathbb{S}(LR^T)\|_F^2 + \lambda_t \|(LR^T)\mathbb{T}^T\|_F^2 \end{aligned}$$

Cell Selection

- Allocate Task to the cell that is hard to be inferred. Then which cell?
 - *Intuition:* If different inference algorithms get significantly different inferred values for a cell, then this cell is said to be hard to be inferred.
- Called *Query by Committee (QBC)*



Selected Cell:
Its ... vary most significantly

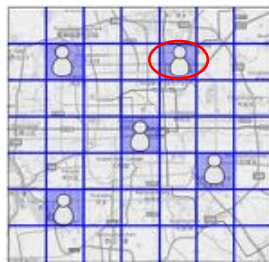
Quality Assessment

- **LOO-BI**

- **Step 1: Leave-one-out Resampling**
- **Step 2: Bayesian Inference**
 - probability distribution of error



Example

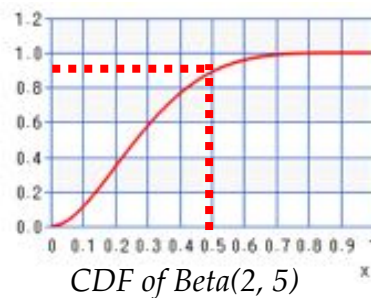


Leave-one-out Result:

One cell out of the five (1/5) is inferred with wrong air quality category (good/medium/bad)

Bayesian Inference:

Classification error $\sim \text{Beta}(2, 5)$



$$P(\text{classification error} \leq 0.5) = 0.9$$



Current quality: In 90% of cycles, the classification error is $\leq 50\%$

Experiment setting

- Datasets

- Temperature: 57 cells (30m*50m), 30-min cycle; **mean absolute error**
- PM2.5: 36 cells (1km*1km), 1-hour cycle; **classification error** (six AQI categories)

- Our Method

- CCS-TA

- Baselines

- RAND-TA: randomly selecting next sensing cell.
- FIX-TA-k: keep same number of sensing cells (k) in each cycle.

Q. How many sensed cells are necessary for sparse crowdsensing and baselines for the same quality requirement?

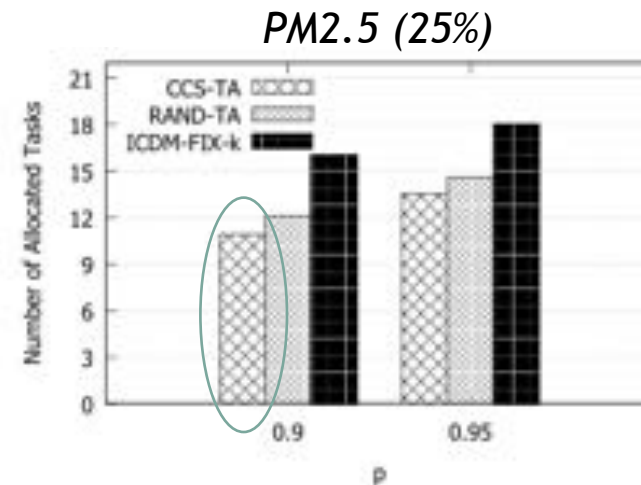
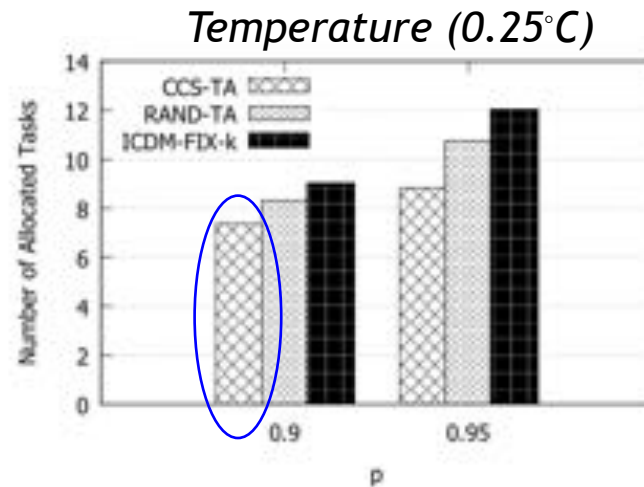
temperature



PM2.5



Number of Cells Required for CCS-TA



*Sense 13% of cells \rightarrow
Mean absolute error of temperature $\leq 0.25^{\circ}\text{C}$ in 90% of cycles*

*Sensing 29% of cells \rightarrow
Classification error of PM2.5 $\leq 25\%$ in 90% of cycles*

Initial Proposal of the idea of Sparse Crowdsensing:

- **L. Wang**, D. Zhang, A. Pathak, C. Chen, H. Xiong, D. Yang, Y. Wang. CCS-TA: Quality-Guaranteed Online Task Allocation in Compressive Crowdsensing. UbiComp 2015
- **L. Wang**, D. Zhang, Y. Wang, C. Chen, X. Han, A. Mhamed. Sparse Mobile Crowdsensing: Challenges and Opportunities. IEEE Communications Magazine, vol. 54, no. 7, pp. 161-167, 2016

From Single Task to Multi Tasks:

- **L. Wang**, D. Zhang, D. Yang, A. Pathak, C. Chen, X. Han, H. Xiong, Y. Wang. SPACE-TA: Cost-Effective Task Allocation Exploiting Intradata and Interdata Correlations in Sparse Crowdsensing. ACM Transactions on Intelligent Systems and Technology, vol. 9, no. 2, pp. 20:1-20:28, 2018

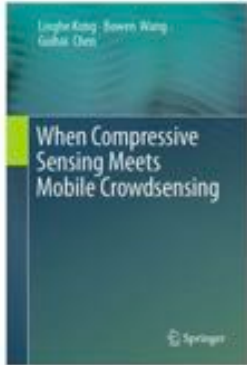
Privacy Protection:

- **L. Wang**, D. Zhang, D. Yang, B. Y. Lim, X. Ma. Differential Location Privacy in Sparse Mobile Crowdsensing. ICDM 2016: 1257-1262.
- T. Zhou, Z. Cai, B. Xiao, **L. Wang**, M. Xu, Y. Chen. Location Privacy-Preserving Data Recovery for Mobile Crowdsensing. UbiComp 2018.

Reinforcement Learning for Cell Selection:

- **L. Wang**, W. Liu, D. Zhang, Y. Wang, E. Wang, Y. Yang. Cell Selection with Deep Reinforcement Learning in Sparse Mobile Crowdsensing. ICDCS 2018
- W. Liu, Y. Yang, E. Wang, **L. Wang**, D. Zeghlache, D. Zhang. Multi-Dimensional Urban Sensing in Sparse Mobile Crowdsensing. IEEE Access vol. 7, pp. 82066-82079, 2019

Related Research from Other Teams



When Compressive Sensing Meets Mobile Crowdsensing

Authors ([view affiliations](#))

Linghe Kong, Bowen Wang, Guihai Chen

Book

Cost-aware compressive sensing for networked sensing systems

2015, *Information Processing in Sensor Networks*, pp 130–141

Liwen Xu (Tsinghua University), Xiaohong Hao (Tsinghua University), Nicholas D. Lane (Microsoft), Xin Liu (University of California, Davis), Thomas Moscibroda (Microsoft)

Density-aware compressive crowdsensing

2017, *Information Processing in Sensor Networks*, pp 29–39

Xiaohong Hao (Tsinghua University), Liwen Xu (Tsinghua University), Nicholas D. Lane (University College London), Xin Liu (University of California, Davis), Thomas Moscibroda (Microsoft)

Active Sparse Mobile Crowd Sensing Based on Matrix Completion

Full Text: PDF Get this Article

Authors: [Kun Xie](#) [Hunan University, Changsha, China](#)
[Xiaocan Li](#) [Hunan University, Changsha, China](#)
[Xin Wang](#) [Stony Brook University, Stony Brook, NY, USA](#)
[Gaogang Xie](#) [Institute of Computing Technology & University of Chinese Academy of Sciences, Beijing, China](#)
[Jigang Wen](#) [Institute of Computing Technology & Chinese Academy of Sciences, Beijing, China](#)
[Dafang Zhang](#) [Hunan University, Changsha, China](#)



Bibliometrics

- Citation Count: 0
- Downloads (cumulative): 107
- Downloads (12 Months): 107
- Downloads (6 Weeks): 61

Published in:



· Proceeding
[SIGMOD '19](#) Proceedings of the 2019 International Conference on Management of Data
Pages 195-210

Amsterdam, Netherlands — June 30 - July 05, 2019

[ACM](#) New York, NY, USA ©2019

[table of contents](#) ISBN: 978-1-4503-5643-5 doi>[10.1145/3299869.3319856](#)

More with less: lowering user burden in mobile crowdsourcing through compressive sensing

2015, *Ubiquitous Computing*, pp 659–670

Liwen Xu (Tsinghua University), Xiaohong Hao (Tsinghua University), Nicholas D. Lane (Bell Labs), Xin Liu (University of California, Davis), Thomas Moscibroda (Microsoft)

Urban Data Scarcity

Collect New Data

Exploit Existing Data

Sparse Crowdsensing

Urban Transfer Learning

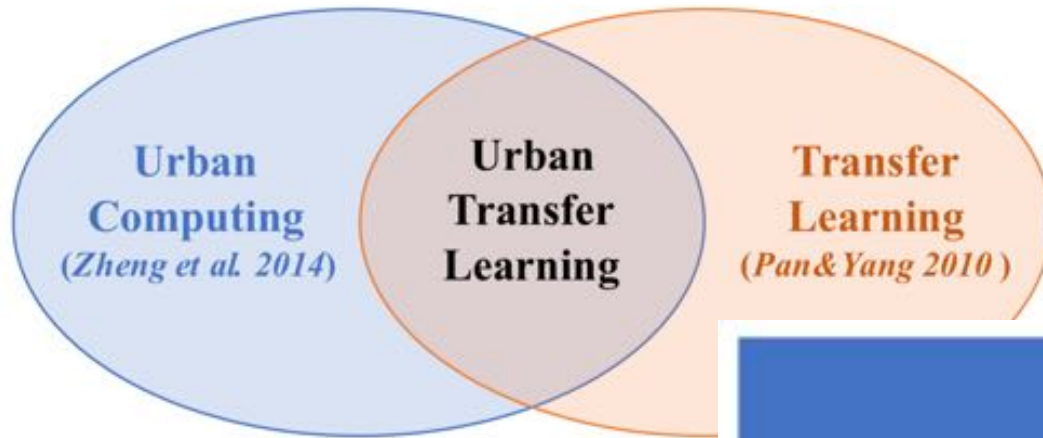
Smart City

Transfer Learning



Instance Transfer, Feature Transfer, Model Transfer ...

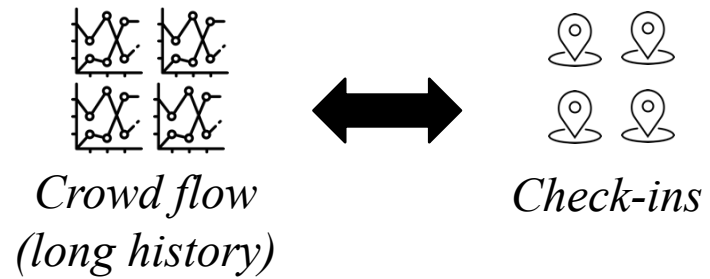
Urban Transfer Learning



	Transfer Learning on Recommendation [5]	Urban Transfer Learning on Traffic Prediction [17]
<i>Source Domain</i>	Movie Rating	City 1's Traffic
<i>Target Domain</i>	Book Rating	City 2's Traffic
<i>Remark</i>	Single data modality (user-item ratings)	<ol style="list-style-type: none">1. Heterogeneous data modalities (daytime, weather, check-ins)2. Spatio-temporal patterns

Basic Idea in Urban Transfer

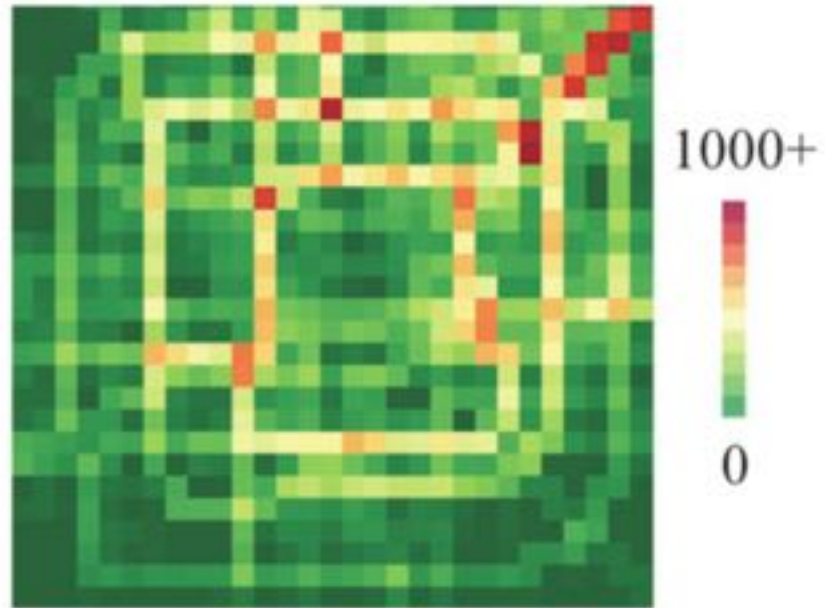
Cross-Modality (heterogeneous data)



Cross-City (spatial transfer)



Crowd Flow Prediction



Predict inflow/outflow of a region in next 30 minutes

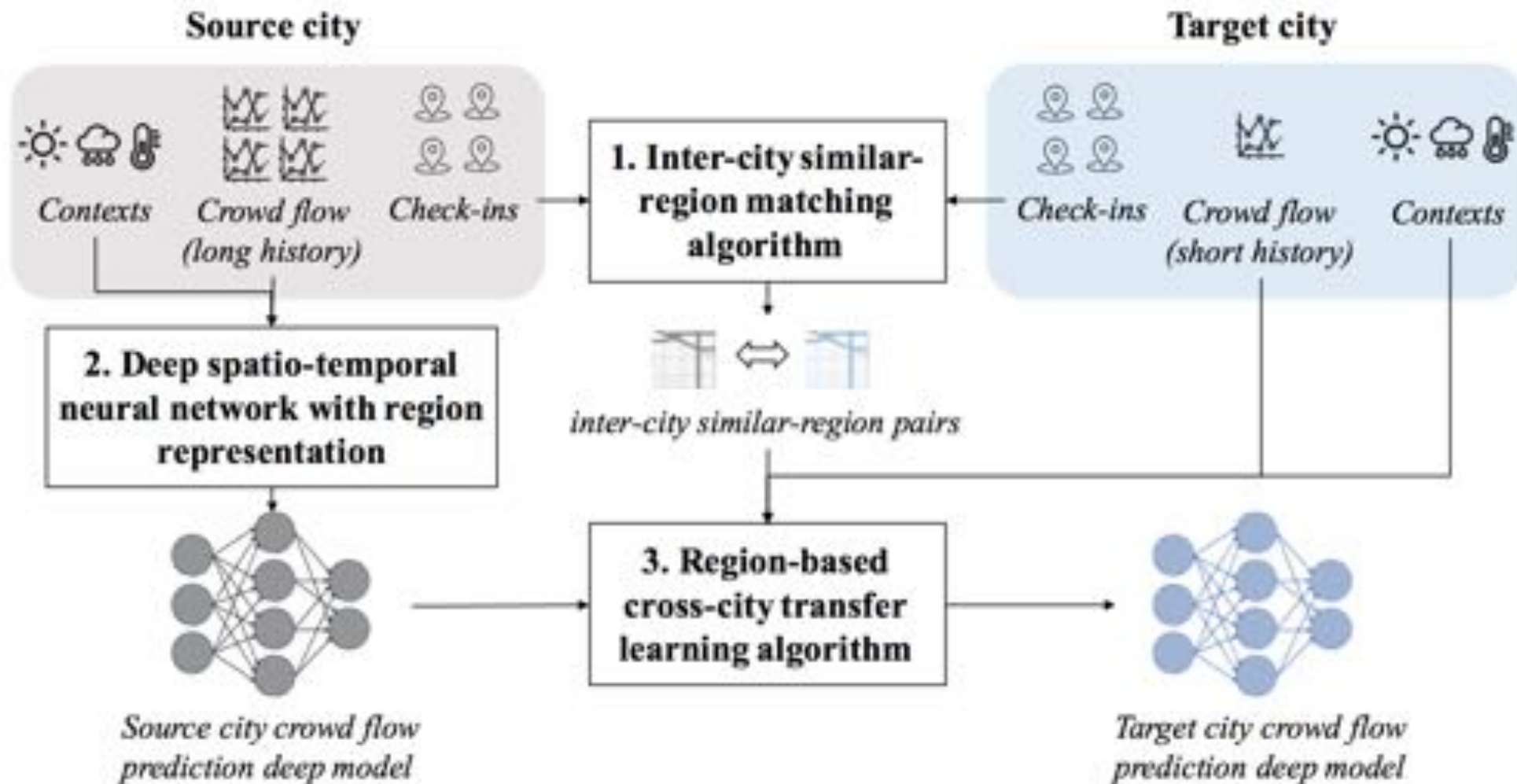
- Deep learning has shown its advantage

Zhang J. et al. "Deep Spatio-Temporal Residual Networks for Citywide Crowd Flows Prediction." AAAI. 2017.

For a new city without so many data?

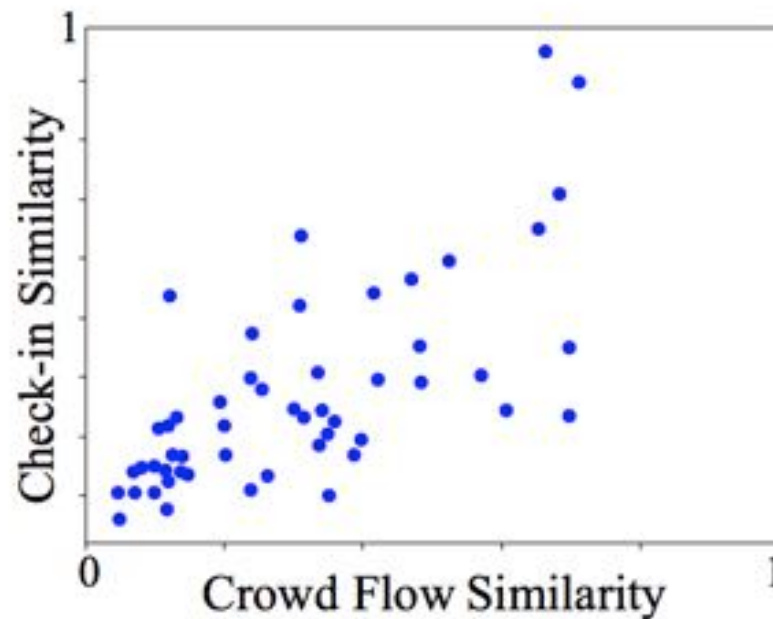
Cross-City Transfer Learning for Deep Spatio-Temporal Prediction. IJCAI, 2019.

RegionTrans

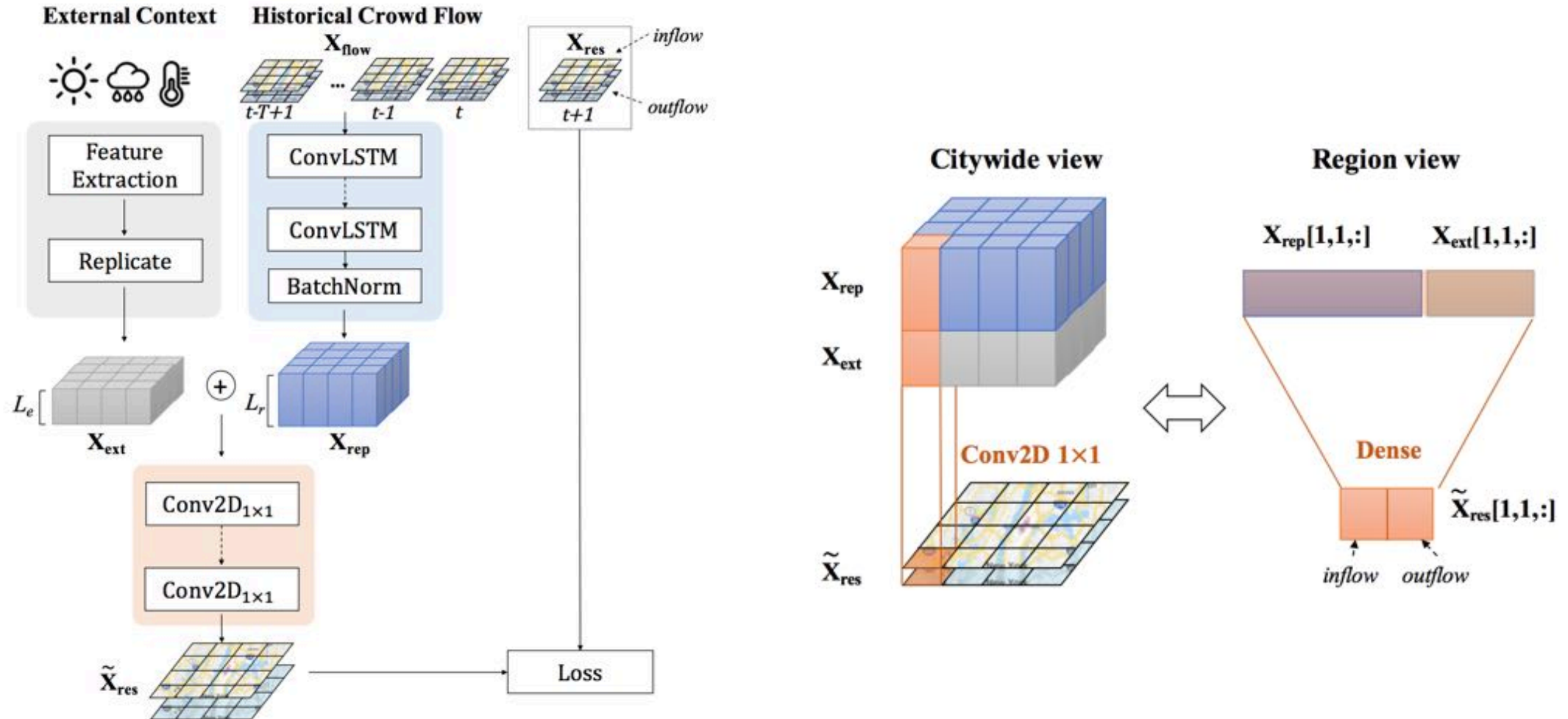


Inter-city Similar Region Matching

- Target city has very little crowd flow data
- For each target region, use **social network check-in** data to find *top-k* similar source regions



DNN with region representation



Cross-city Transfer Learning

- Learn a model in source city, use the parameter as the start of the target city model
- Use the little data in the target city to optimize the parameter
- Optimization objective

$$\arg \min_{\theta} \left[w \left(\frac{1}{k} \sum_{1 \leq i \leq k} \|\rho_i \circ (\mathbf{X}_{rep} - \hat{\mathbf{X}}_{rep}^i)\|_2^2 \right) + (1 - w) \|\tilde{\mathbf{X}}_{res} - \mathbf{X}_{res}\|_2^2 \right]$$

Minimize representation discrepancy

Minimize prediction error

Experiment: Bike Sharing

	Washington D.C.	Chicago
<i>#Trip records</i>	6,519,741	6,690,351
<i>Time span</i>	2015.1.1 - 2016.12.31	
<i>Time interval</i>	30 minutes	
<i>Region size</i>	$1km \times 1km$	
<i>City map size</i>	20×20	

RegionTrans can outperform state-of-the-art deep learning with fine-tune by reducing up to 10% error (RMSE)

	D.C.→Chicago			Chicago→D.C.		
	1-day	3-day	7-day	1-day	3-day	7-day
Target Data Only						
<i>ARIMA</i>	0.740	0.694	0.679	0.707	0.661	0.647
<i>DeepST</i>	0.771	0.711	0.636	1.075	0.767	0.691
<i>ST-ResNet</i>	0.914	0.703	1.053	0.869	0.738	1.054
Source & Target Data						
<i>DeepST (FT)</i>	0.652	0.611	0.566	0.672	0.619	0.586
<i>ST-ResNet (FT)</i>	0.667	0.615	0.613	0.695	0.623	0.608
<i>RegionTrans</i>	0.587	0.576	0.553	0.600	0.581	0.573

Robust Against Negative Transfer

	D.C.→NYC		NYC→D.C.	
	1-day	3-day	1-day	3-day
Target Data Only				
<i>ARIMA</i>	0.360	0.341	0.707	0.661
<i>DeepST</i>	0.350	0.359	1.075	0.767
<i>ST-ResNet</i>	0.376	0.349	0.869	0.738
Source & Target Data				
<i>DeepST (FT)</i>	0.363	0.369	0.713	0.711
<i>ST-ResNet (FT)</i>	0.385	0.349	0.696	0.691
<i>RegionTrans</i>	0.328	0.305	0.665	0.593

The performance of existing deep models gets worse when using fine-tune for D.C. → NYC, but *RegionTrans* still makes an effective transfer.

Future Opportunities

Multi-Source Transfer



Privacy-Preserving Transfer



How is Uber Movement preserving the privacy of Uber riders and drivers?

Preserving rider and driver privacy is our #1 priority. All data is anonymized and aggregated to ensure no personally identifiable information or user behavior can be surfaced through the Movement tool. All data shared through Movement adheres to Uber's privacy policy, and at no point will Movement provide a means for partners to access individual driver or rider details in any way, shape or form. Further details of how we're ensuring this are available in our FAQ.

Uber Movement 共享数据平台的隐私相关声明



Leye Wang, Key Lab of High Confidence Software Technologies, Peking University

Bin Guo, Northwestern Polytechnical University

Qiang Yang, Hong Kong University of Science and Technology

L. Wang, B. Guo, Q. Yang, “*Smart City Development with Urban Transfer Learning*” . IEEE Computer 51(12): 32-41 (2018).

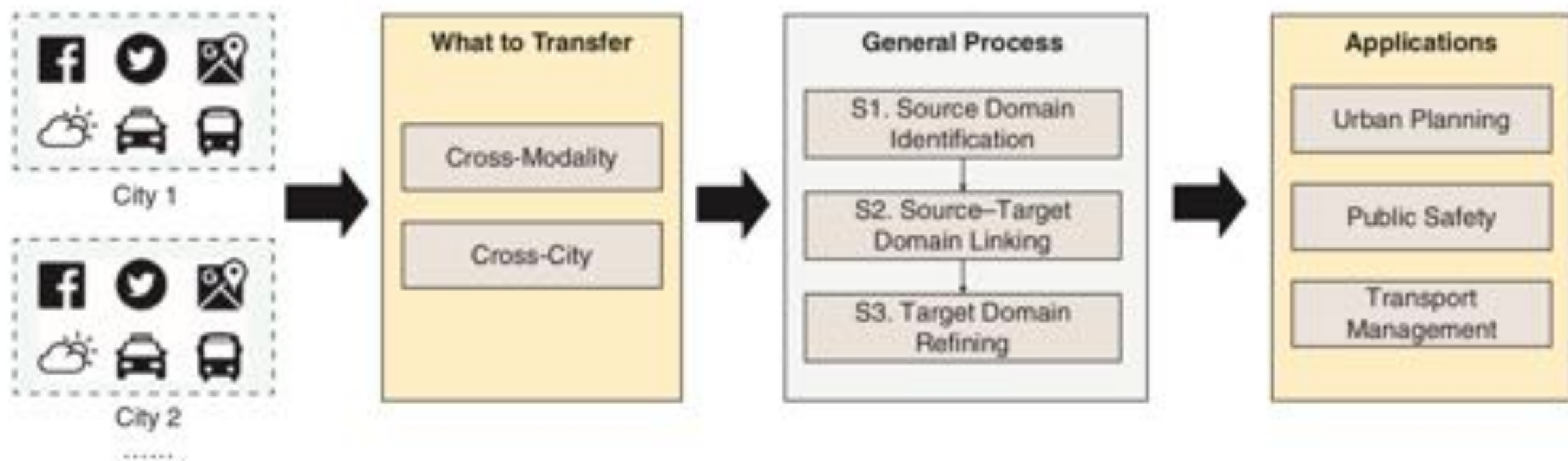
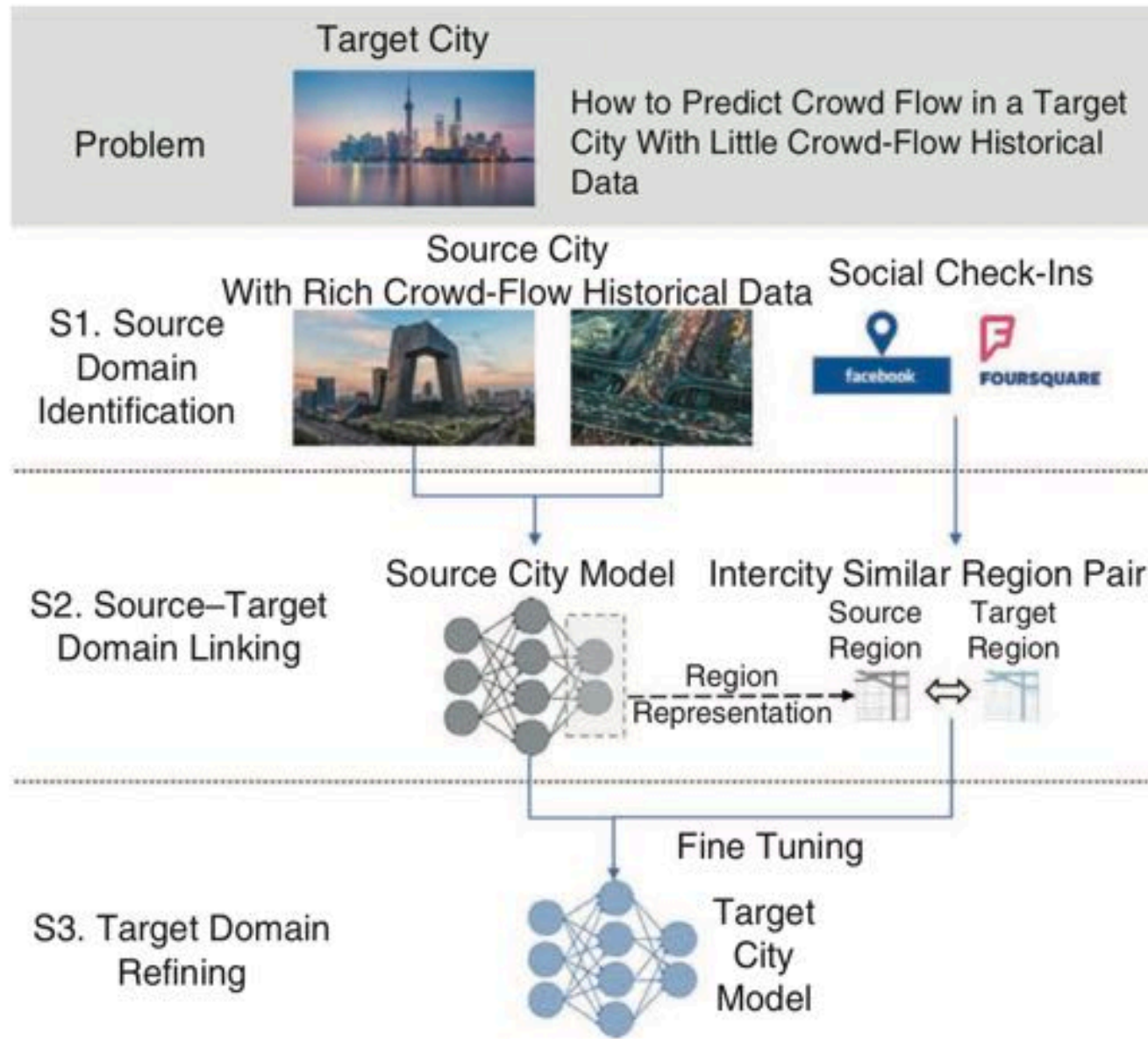
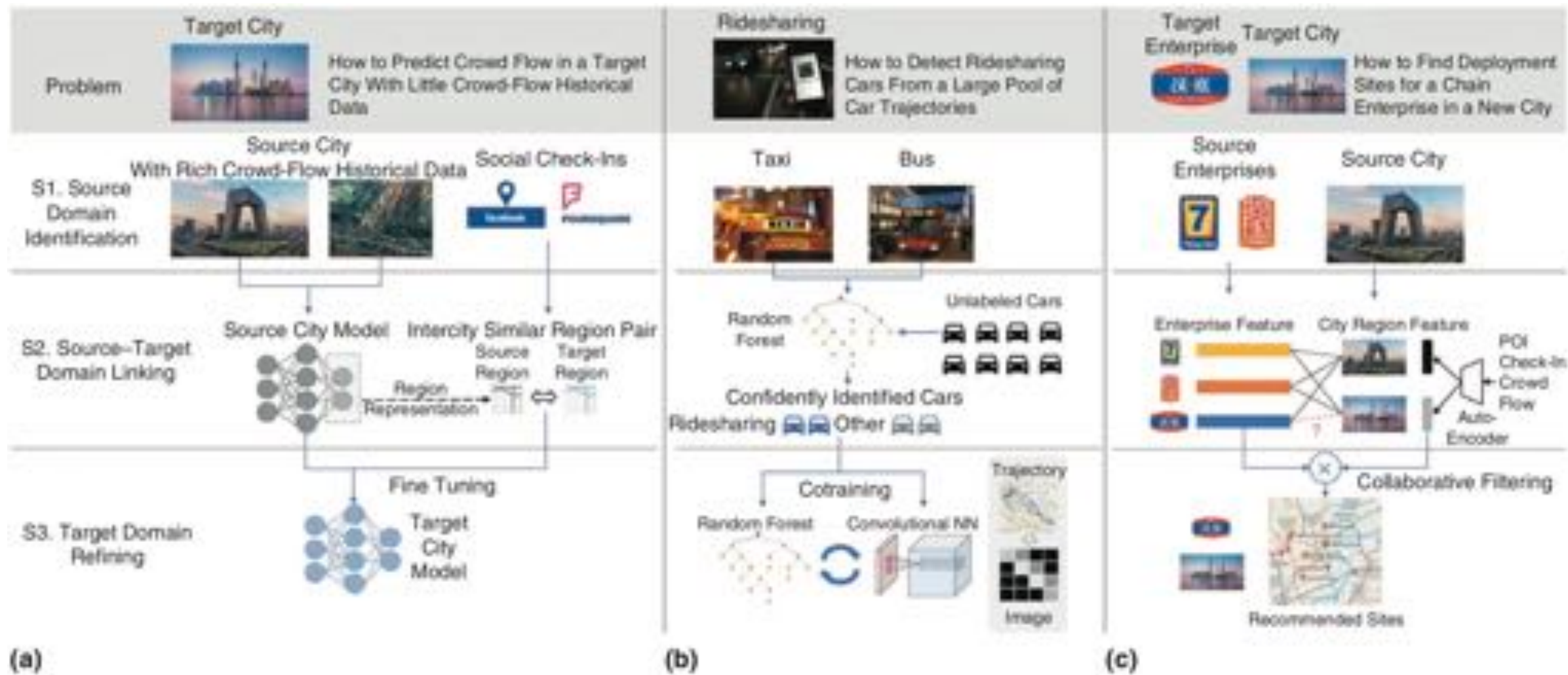


FIGURE 2. The framework of urban transfer learning.





(a)

(b)

(c)

Related Papers

- Wang, L., Geng, X., Ma, X., Liu, F., & Yang, Q. (2019). Cross-city transfer learning for deep spatio-temporal prediction. *IJCAI*.
- Wang, L., Geng, X., Ma, X., Zhang, D., & Yang, Q. (2019). Ridesharing Car Detection by Transfer Learning. *Artificial Intelligence*.
- Wang, L., Guo, B., & Yang, Q. (2018). Smart City Development with Urban Transfer Learning. *IEEE Computer*.
- Yao, H., Liu, Y., Wei, Y., Tang, X., & Li, Z. (2019). Learning from Multiple Cities: A Meta-Learning Approach for Spatial-Temporal Prediction. *In The World Wide Web Conference* (pp. 2181-2191).
- Guo, B., Li, J., Zheng, V. W., Wang, Z., & Yu, Z. (2018). CityTransfer: Transferring Inter-and Intra-City Knowledge for Chain Store Site Recommendation based on Multi-Source Urban Data. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies (UbiComp)*, 1(4), 135.
- Wei, Y., Zheng, Y., & Yang, Q. (2016). Transfer knowledge between cities. *In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (pp. 1905-1914).
- Jean, N., Burke, M., Xie, M., Davis, W. M., Lobell, D. B., & Ermon, S. (2016). Combining satellite imagery and machine learning to predict poverty. *Science*, 353(6301), 790-794.
- Zheng, Y., Capra, L., Wolfson, O., & Yang, H. (2014). Urban computing: concepts, methodologies, and applications. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 5(3), 38.
- Pan, S. J., & Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10), 1345-1359.



Thanks !

leyewang@pku.edu.cn

<https://sites.google.com/site/wangleye/>