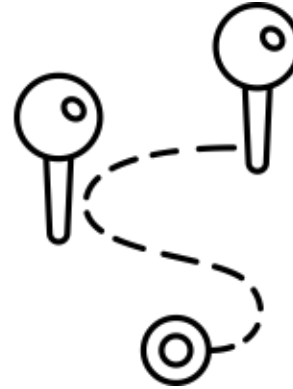


# Trajectory-User Linking using Higher- order Mobility Flow Representations

Supervisor: Manos Papagelis

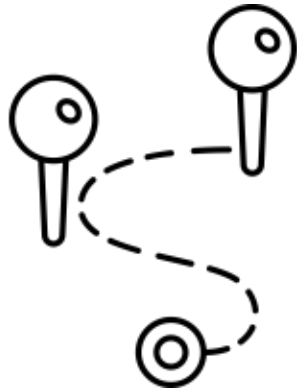
M.SC. THESIS OF MAHMOUD ALSAEED





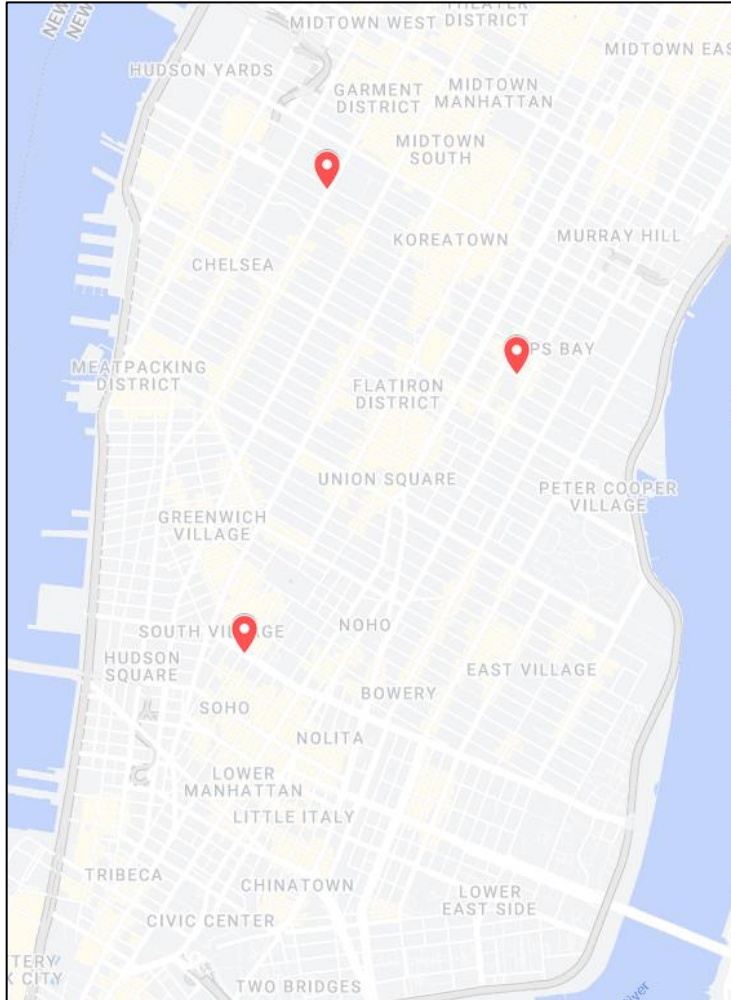
trajectories are  
**unique** to each person

# Trajectory-user Linking (TUL)

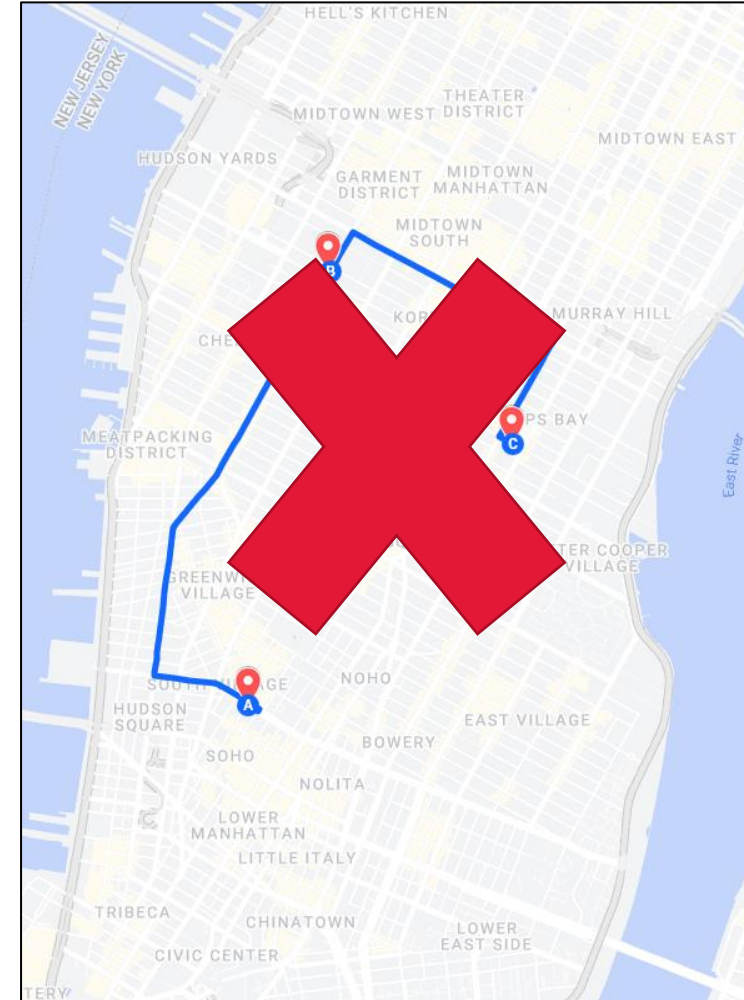


trajectory-user linking **aims to link** anonymous trajectories to users who generate them

# Data for Trajectory-user Linking (TUL)



Check-ins Trajectory

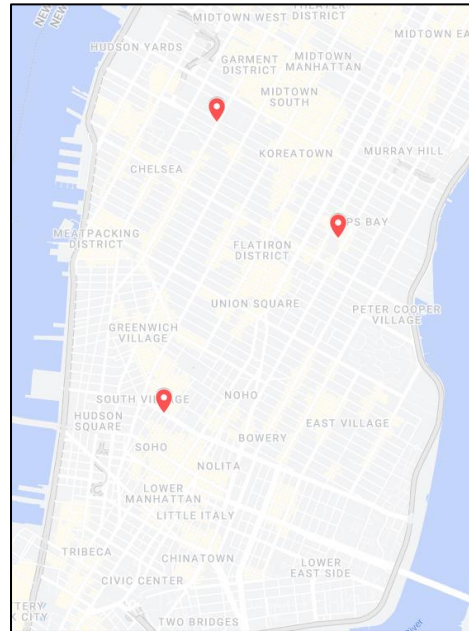


Mobility flow

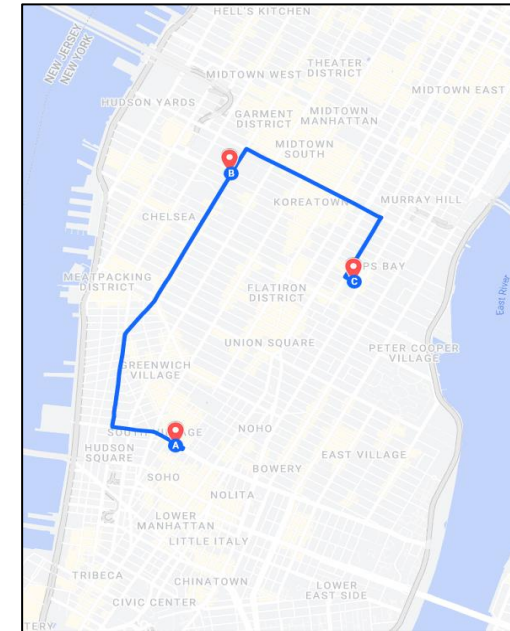
# Why do the limitations exist ?

➤ Check-ins data

➤ Active data collection



Check-ins trajectory



Mobility flow

## Limitations of the current approaches

- Data Quality
  - low accuracy and completeness
- Data sparsity
  - limited data
- Imbalanced Data
  - 80% of the data is generated by 20% of the users



# Problem Definition

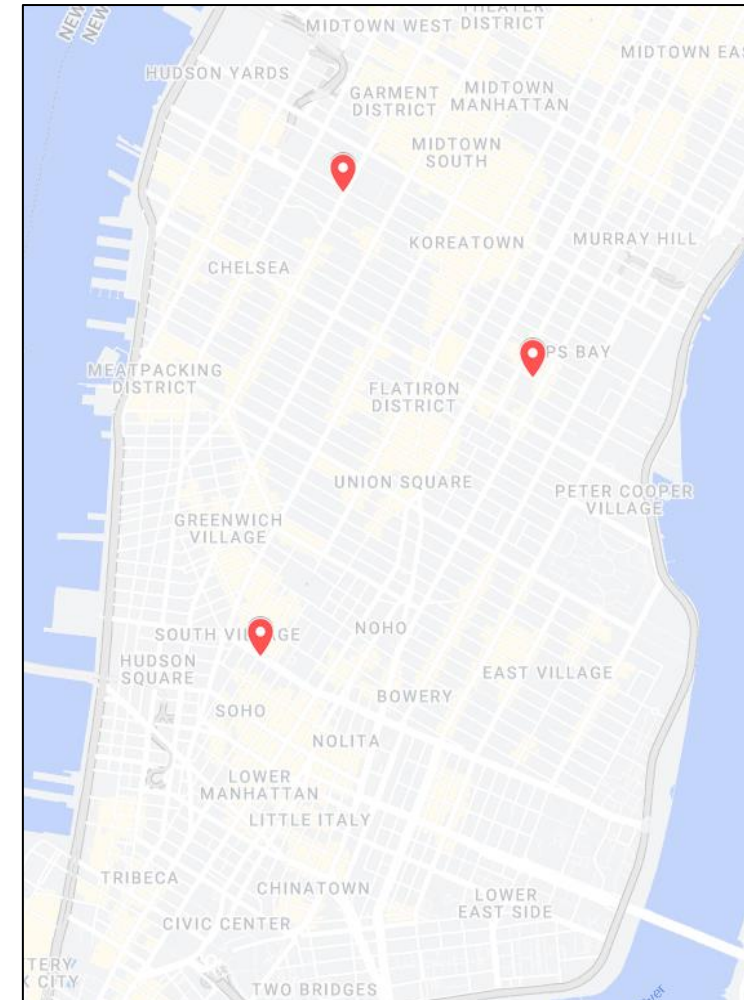
# What is a check-in trajectory ?

## > Check-in record/visit

- $r = (u, p, t, \langle x, y \rangle)$

## > Check-in trajectory

- $Tr = \{r_1, r_2, \dots, r_m\}$



Check-ins Trajectory



# Problem Definition

- Trajectory-user linking
  - aims to link anonymous **trajectories** to **users** who generate them
- $\mathcal{T} = \{Tr_1, Tr_2, \dots, Tr_n\}$  - trajectories
- $\mathcal{U} = \{u_1, u_2, u_3, \dots, u_c\}$  - users

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(y_i - y'_i)$$

$\mathcal{F}$  is the set of all classifiers in the hypothesis space

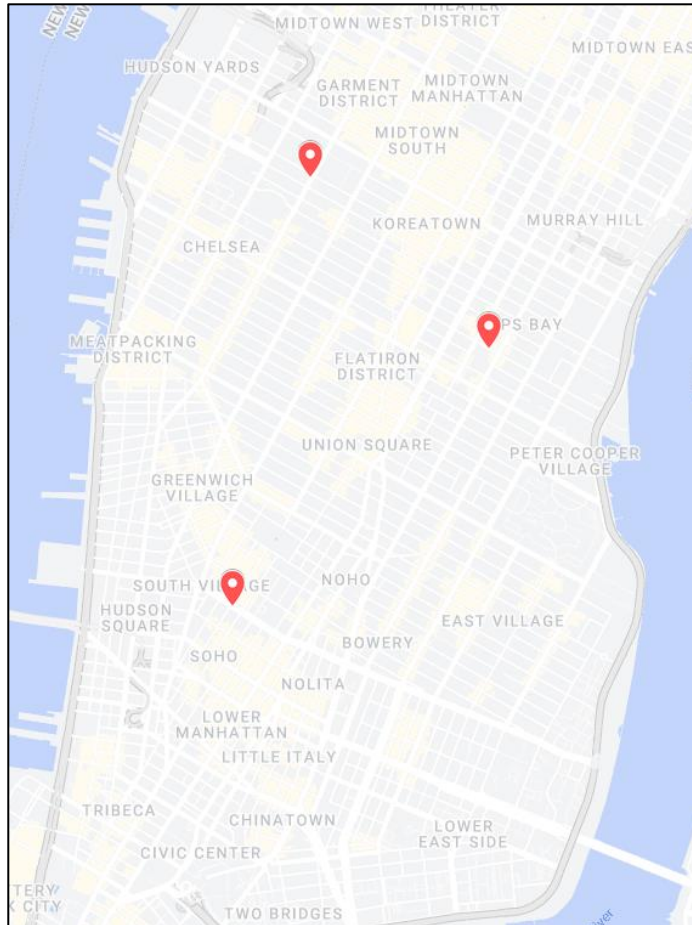
$\mathcal{L}(\cdot)$  is the loss between the predicted label  $y'_i$  and the true label  $y_i$

# Methodology

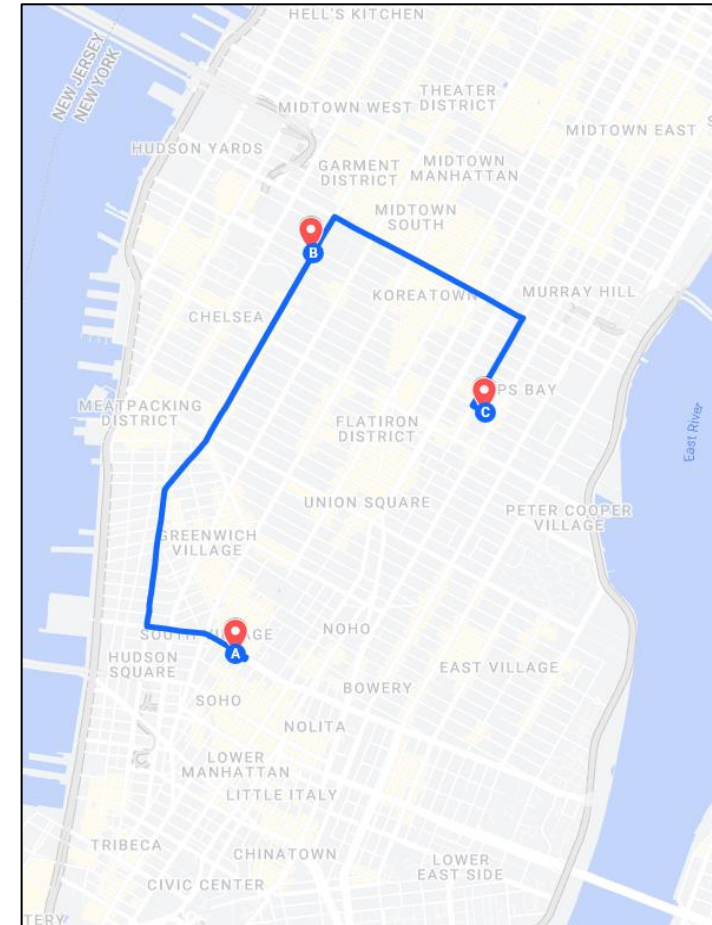
## Overview

- Generating higher-order mobility flow representations
  - generating mobility flow data from check-ins
  - generating higher-order mobility flow and check-ins
  
- Modeling trajectory-User linking

# Generating Mobility flow data

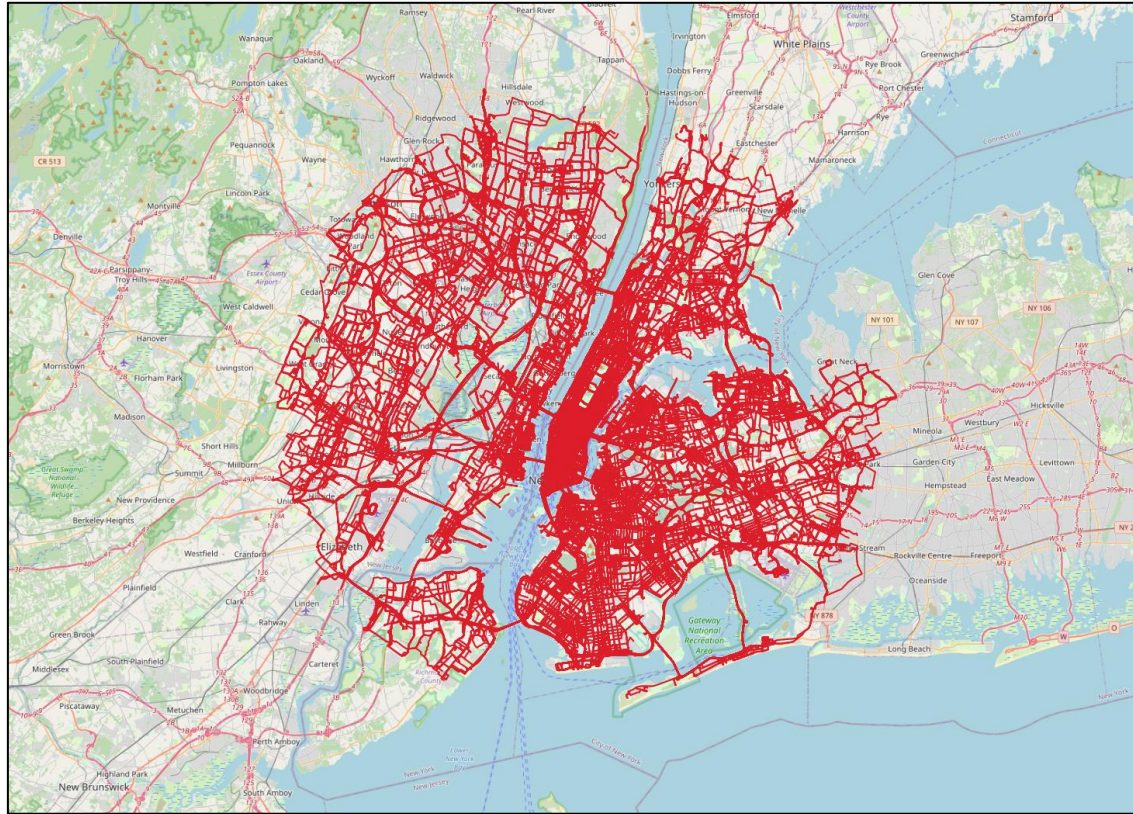


Check-ins Trajectory

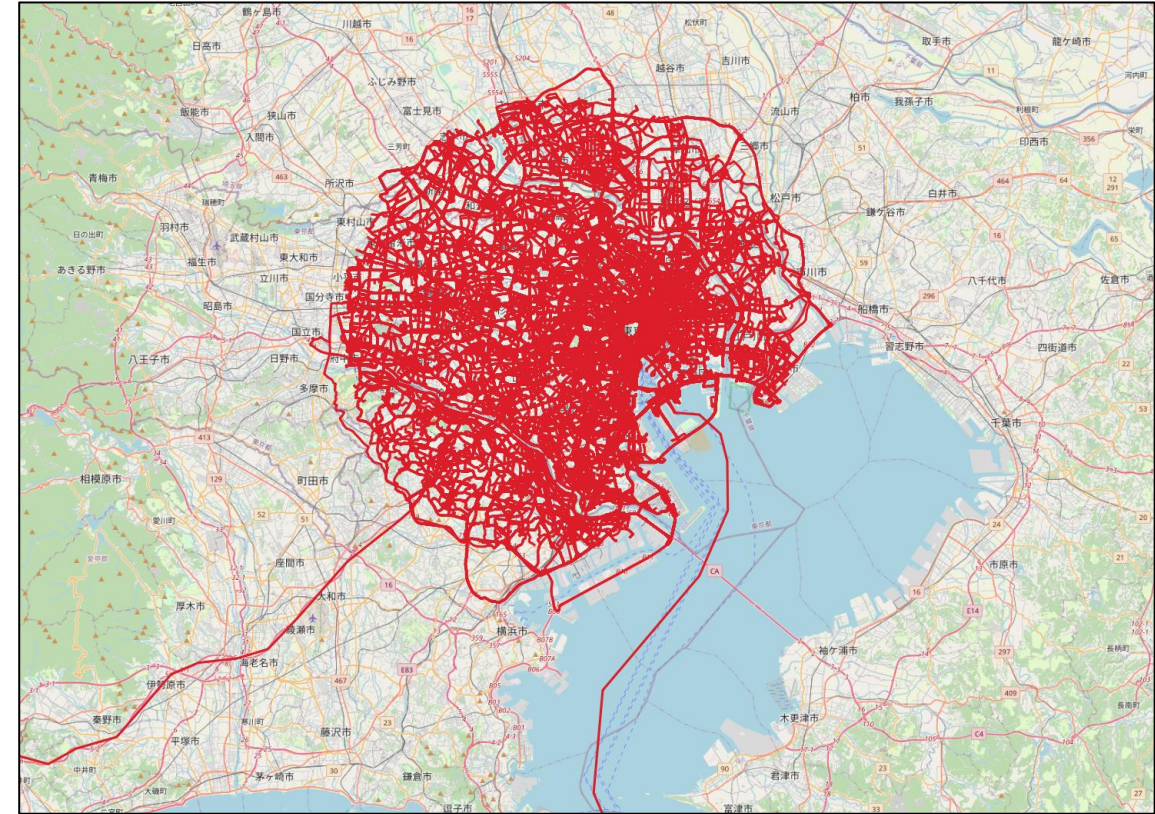


Mobility flow

# Mobility flow of NYC and TKY

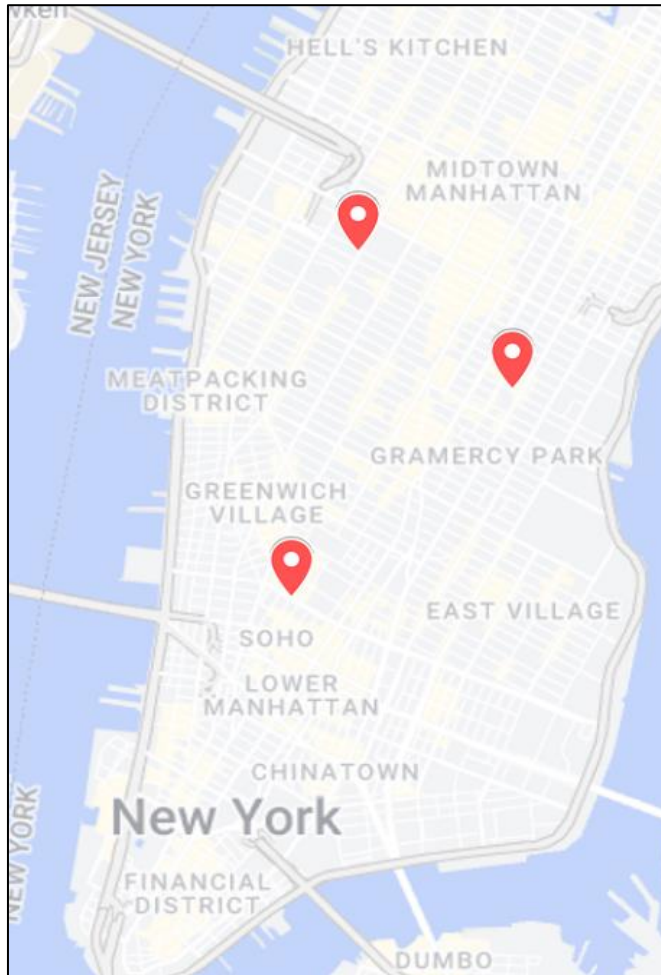


NYC

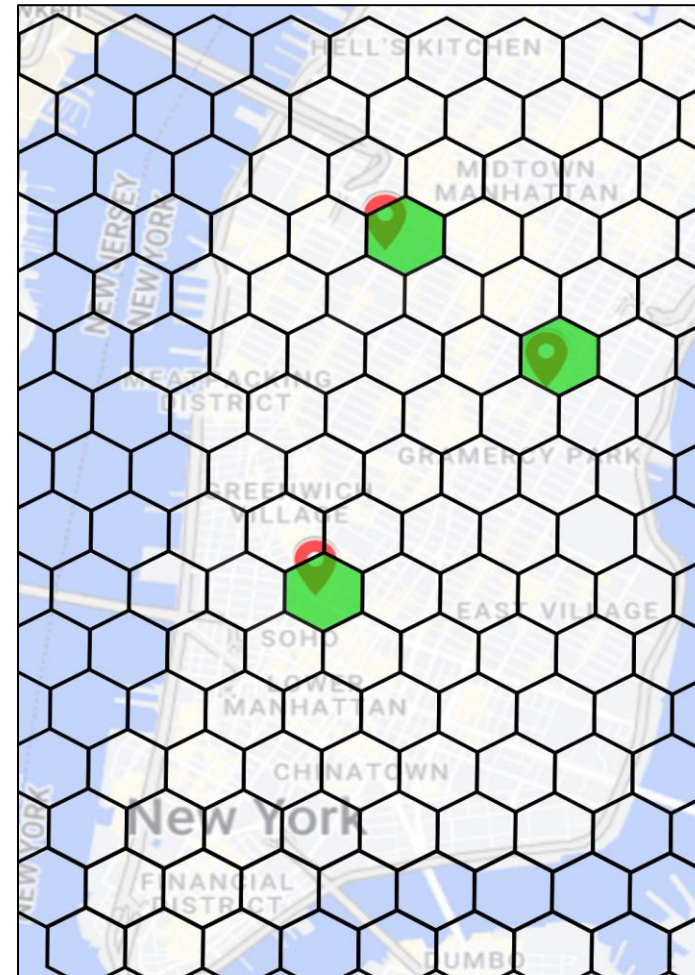


TKY

# Generating higher-order check-ins



Check-ins Trajectory



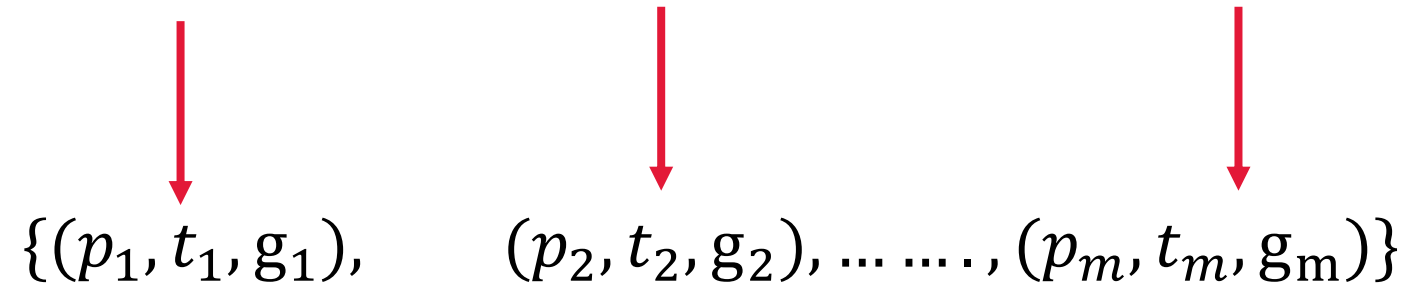
Higher-order  
check-ins

# Translate check-ins to Higher-order

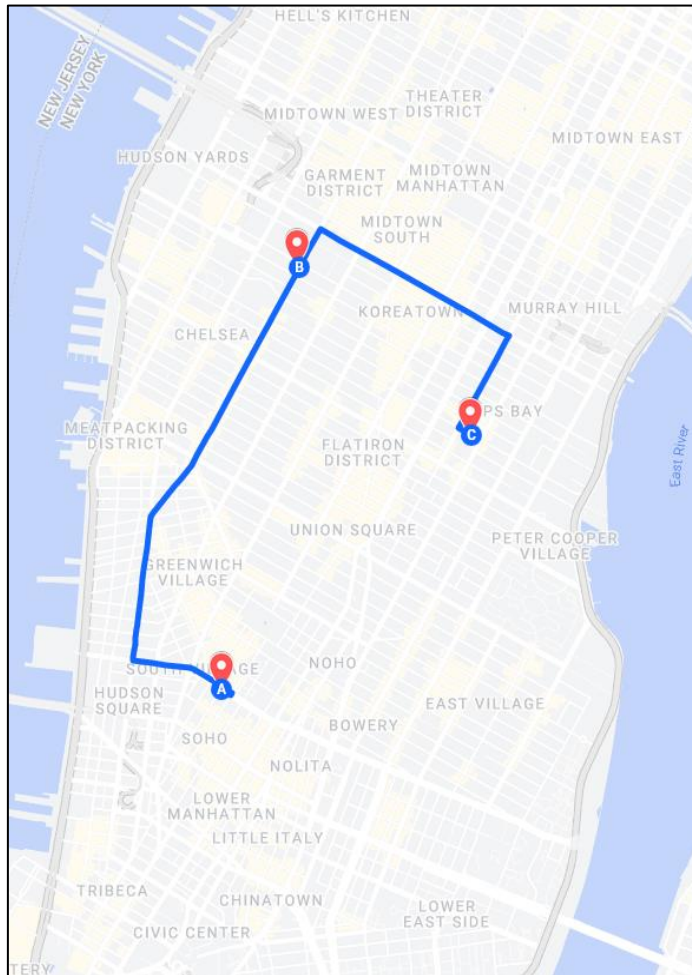
Check-ins

$$Tr = \{r_1, r_2, \dots, r_m\} = \{(p_1, t_1, \langle x_1, y_1 \rangle), (p_2, t_2, \langle x_2, y_2 \rangle), \dots, (p_m, t_m, \langle x_m, y_m \rangle)\}$$

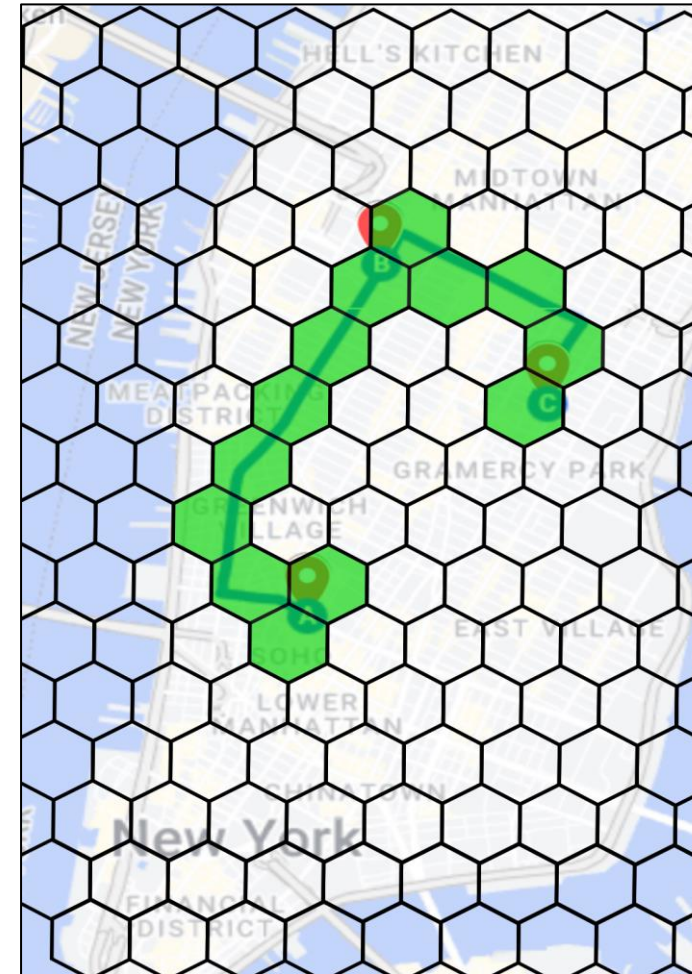
Higher-order


$$\{(p_1, t_1, g_1), (p_2, t_2, g_2), \dots, (p_m, t_m, g_m)\}$$

# Generating Higher-order Mobility flow



Mobility flow



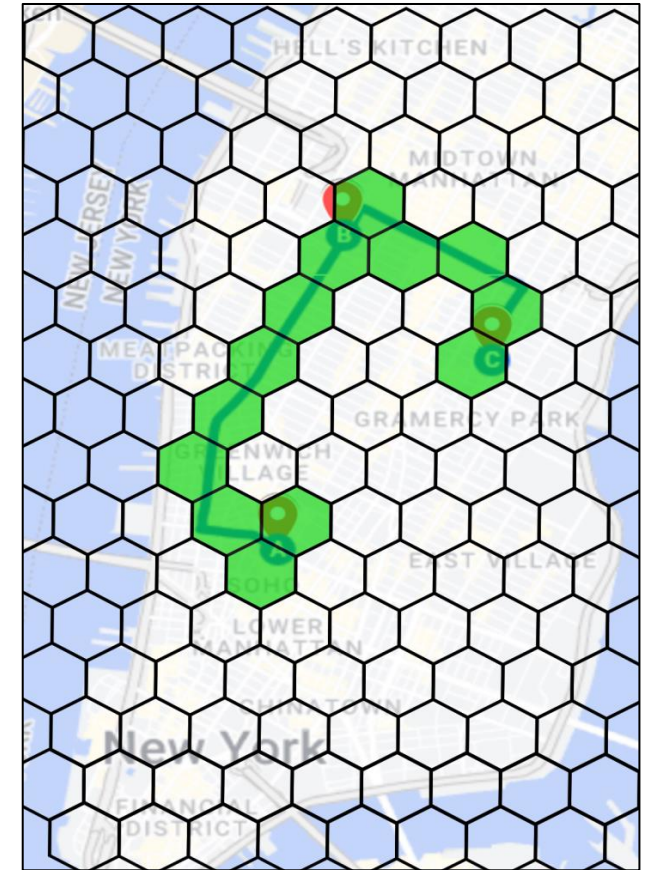
Higher-order  
Mobility flow



## Higher-order Mobility flow data

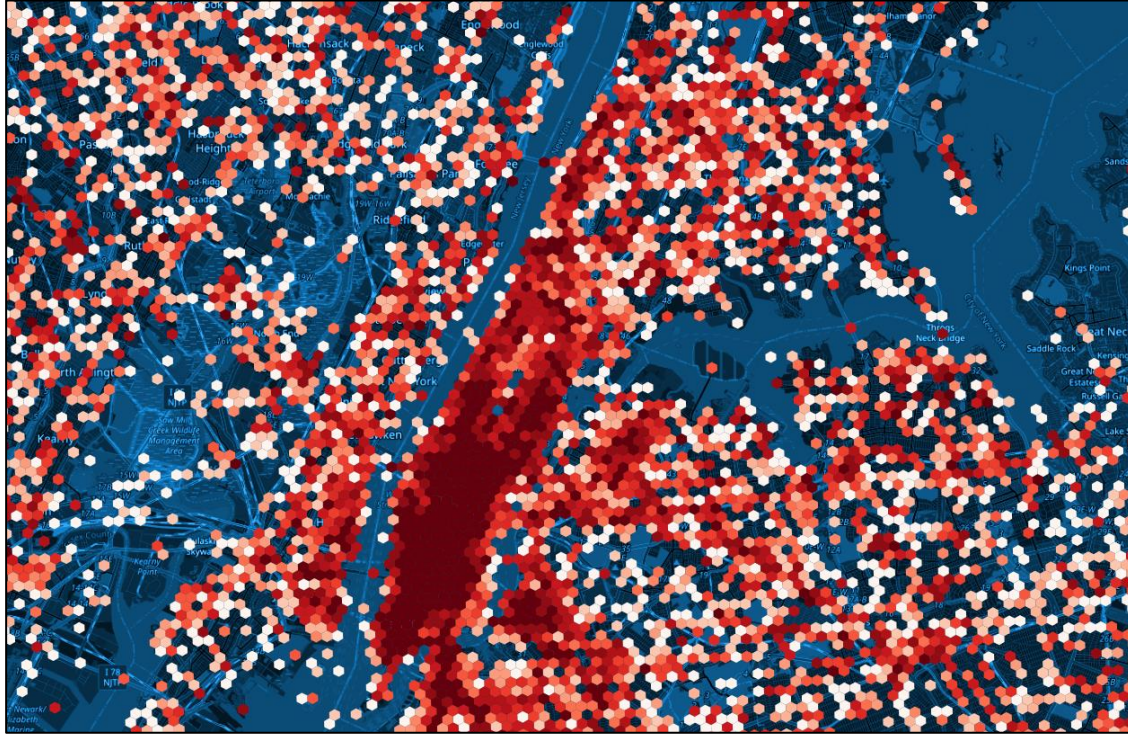
➤ Represented as a sequence of grid cells

➤  $\{g_1, g_2, \dots\}$

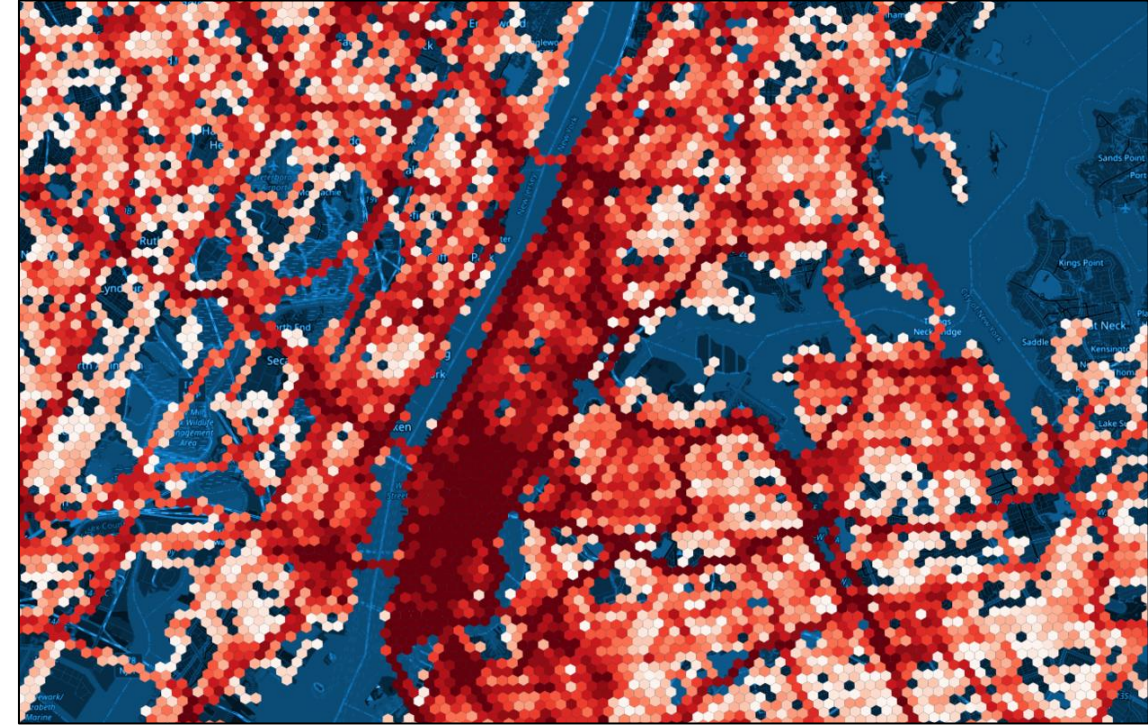


Higher-order  
Mobility flow

# FOURSQUARE-NYC Heatmap



Higher-order check-ins



Higher-order Mobility flow

## How to calculate Sparsity ?



Alex  
{1,3}



Eve  
{2}



Bob  
{1,2,3}

	$p_1$	$p_2$	$p_3$
Alex	1	0	1
Eve	0	1	0
Bob	1	1	1

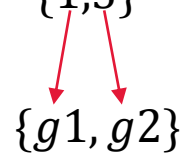
Sparsity = % of **zeros** in **User-POI** matrix

$$\frac{3}{9} = 30\%$$

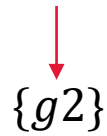
# Higher-order Sparsity



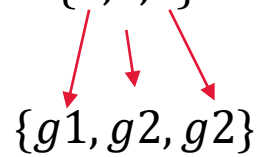
Alex  
{1,3}



Eve  
{2}



Bob  
{1,2,3}

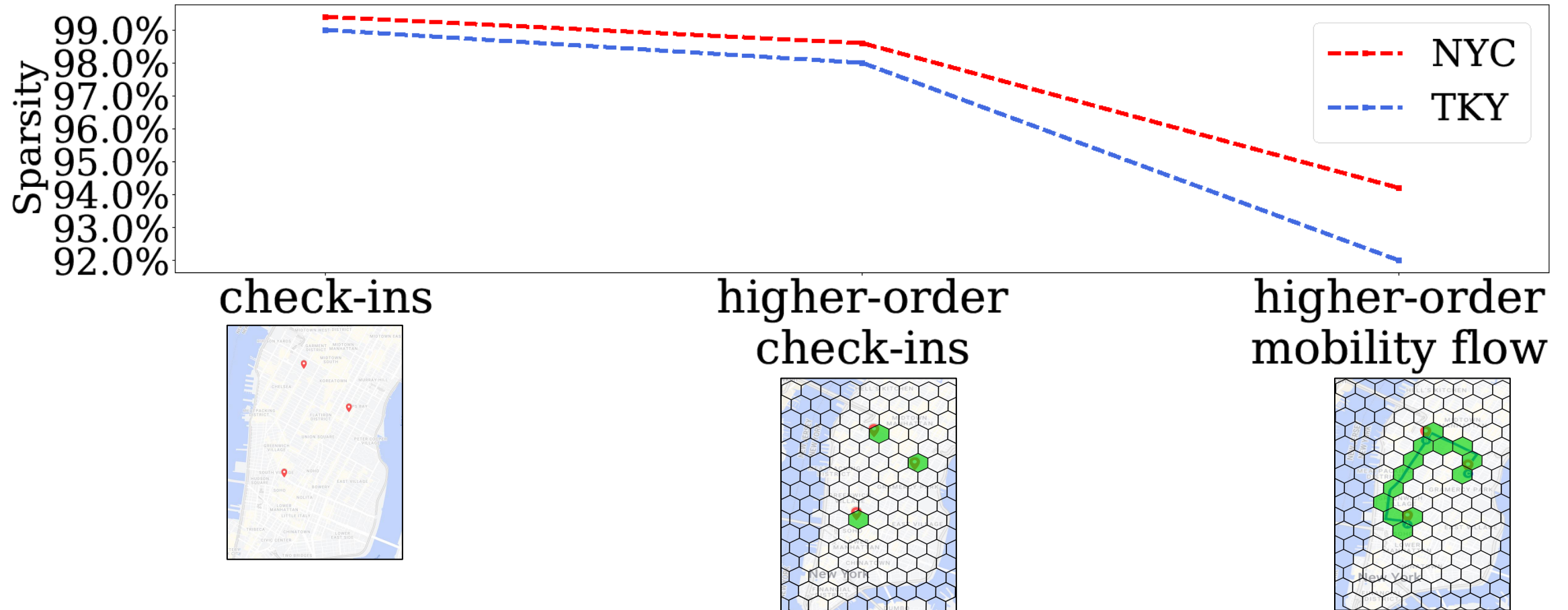


	$g_1$	$g_2$
Alex	1	1
Eve	0	1
Bob	1	2

$$\frac{1}{6} = 16\%$$

Check-in Sparsity  $\geq$  Higher-order Sparsity

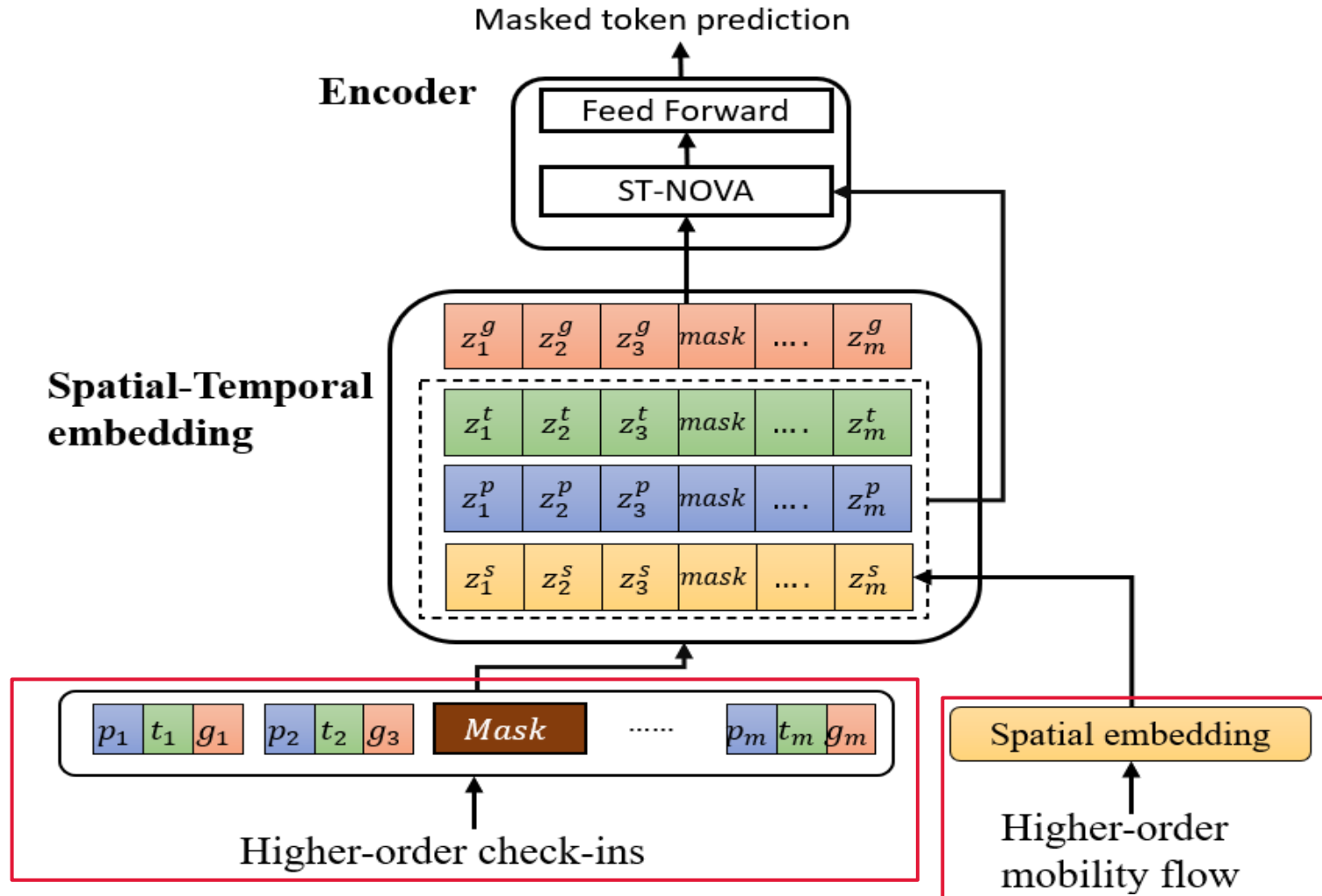
# Impact of higher-order abstraction on sparsity



# Overview

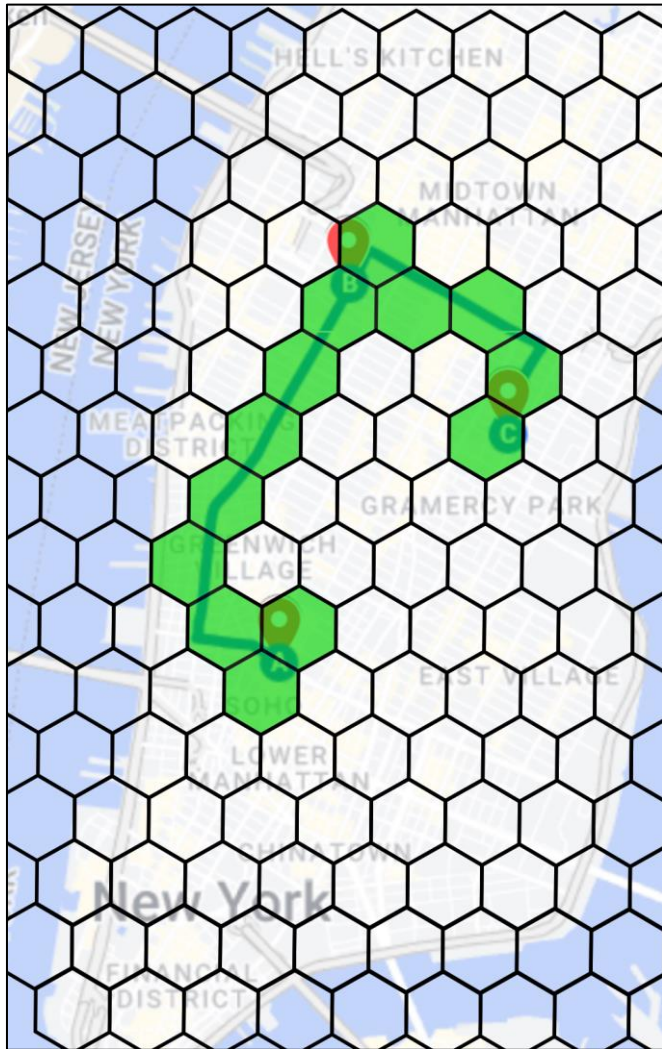
- Generating higher-order mobility flow representations
  - generating mobility flow data from check-ins
  - generating higher-order mobility flow and check-ins
  
- Modeling trajectory-User linking

# TULHOR (trajectory-user linking using higher-order representations )

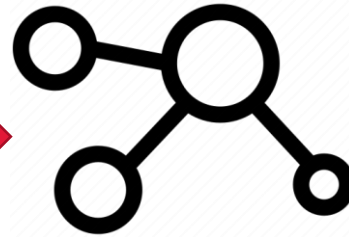




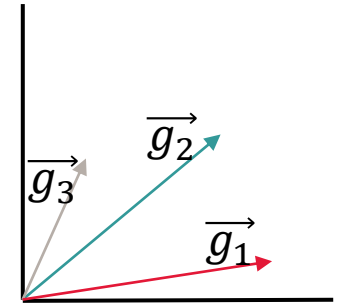
# Spatial embedding



Higher-order  
Mobility flow

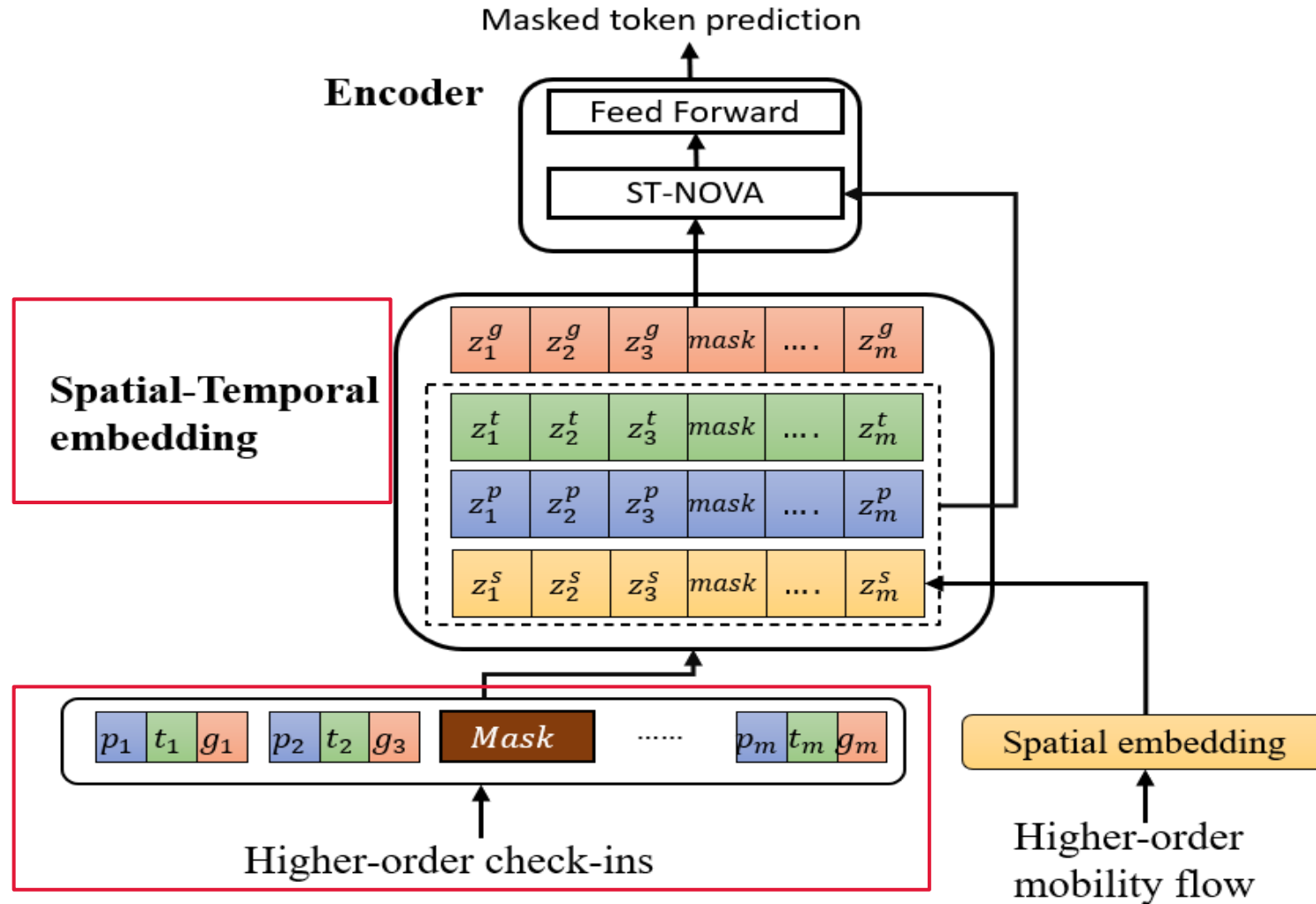


Mobility flow graph



Spatial  
Embedding

# TULHOR (trajectory-user linking using higher-order representations )



# Spatial-temporal embedding

Higher-order  
check-ins  $\{(p_1, t_1, g_1), (p_2, t_2, g_2), \dots, (p_m, t_m, g_m)\}$

$$\begin{aligned} z_i^g &= \phi_g(g_i, W_g) \\ z_i^p &= \phi_p(p_i, W_p) \\ z_i^s &= \phi_s(g_i, W_s) \end{aligned}$$

# Temporal-aware positional encoding

$$\{(p_1, t_1, g_1), (p_2, t_2, g_2), \dots \dots, (p_m, t_m, g_m)\}$$

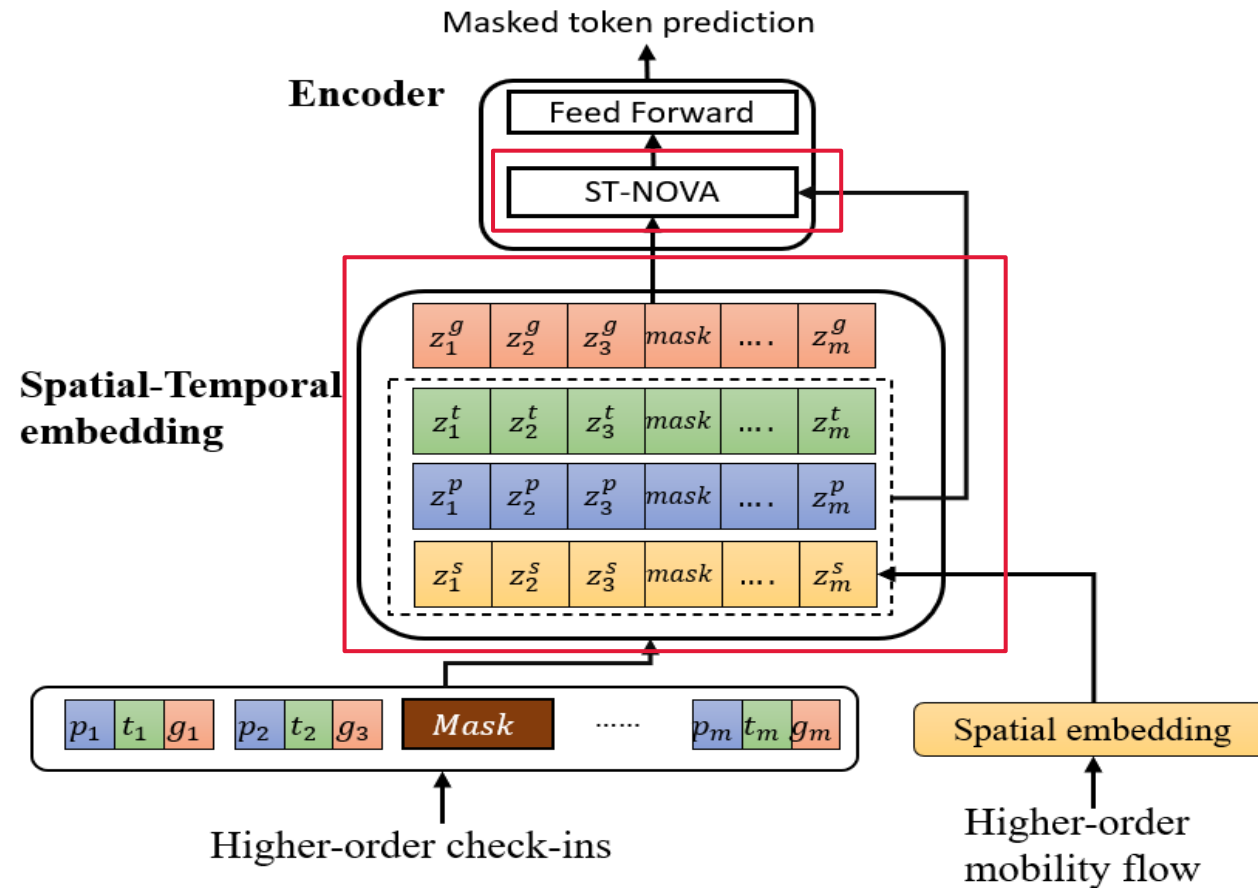
$$[z_i^t]_j = \begin{cases} \sin(w_j t_i), & \text{if } j \text{ is odd} \\ \cos(w_j t_i), & \text{if } j \text{ is even} \end{cases}$$

## Output of the Spatial-temporal embedding

$$R^{(id)} = z_1^g, z_2^g, \dots, z_m^g$$

$$R = (\{z_1^s, z_2^s, \dots, z_m^s\}, \{z_1^p, z_2^p, \dots, z_m^p\}, \{z_1^t, z_2^t, \dots, z_m^t\})$$

# TULHOR



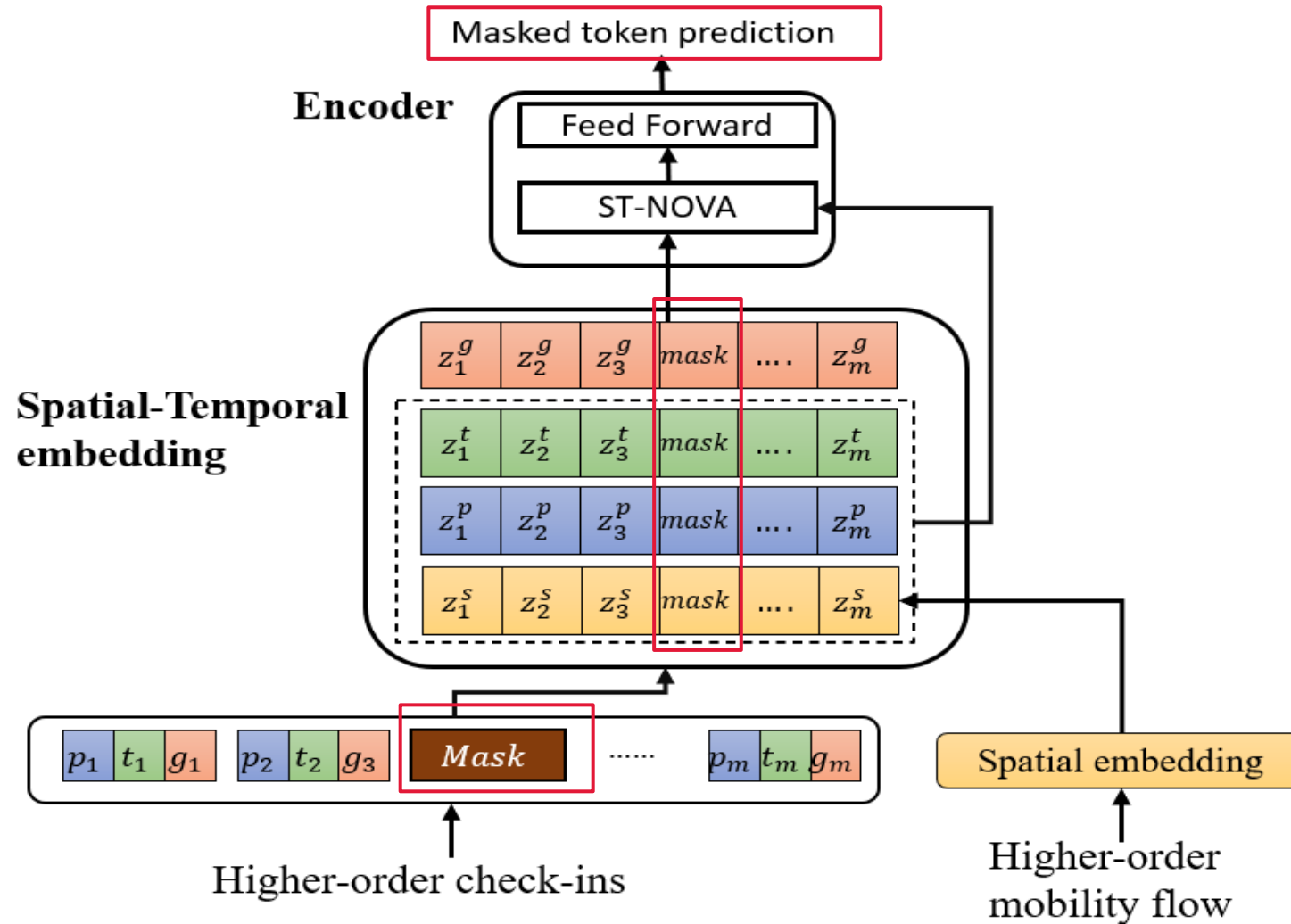
## ST-NOVA

$$NOVA(R^{(id)}, R) = \sigma\left(\frac{QK^T}{\sqrt{d_L}}\right)V$$

$$Q = R^{id} \times W_Q, K = F \times W_k, V = F \times W_V$$

$$F = MLP(R^{(id)} || R)$$

# TULHOR





## Two stages

### ➤ Pre-training TULHOR

- Input : higher-order check-ins + Masking , higher-order mobility flow
- Output: predicting masked token

### ➤ Fine-tuning TULHOR

- Input : higher-order check-ins
- Output: user who generated the higher-order check-ins

Experiment

# Overview

## > Datasets

- Foursquare NYC and TKY

## > Experiments

- TULHOR accuracy performance (vs SOTA and baselines)
- TULHOR Ablation study
- Tessellation granularity (grid size) effect

<b>DATASET</b>	$ \mathcal{U} $	$ \mathcal{T} $
<b>FOURSQUARE-NYC</b>	108	6795
	209	9,637
	234	10,133
<b>FOURSQUARE-TKY</b>	108	9343
	209	14,151
	451	20,964

# Baselines

## ➤ Conventional ML:

- Decision Tree
- Linear Discriminant Analysis (LDA)
- Linear Support Vector Machine (SVM)

## ➤ TULER:

- RNN
- LSTM
- GRU

## ➤ DeepTUL

- RNN (DeepTUL)
- LSTM (Attn-LSTM)
- GRU (Attn-GRU)

# TULHOR performance

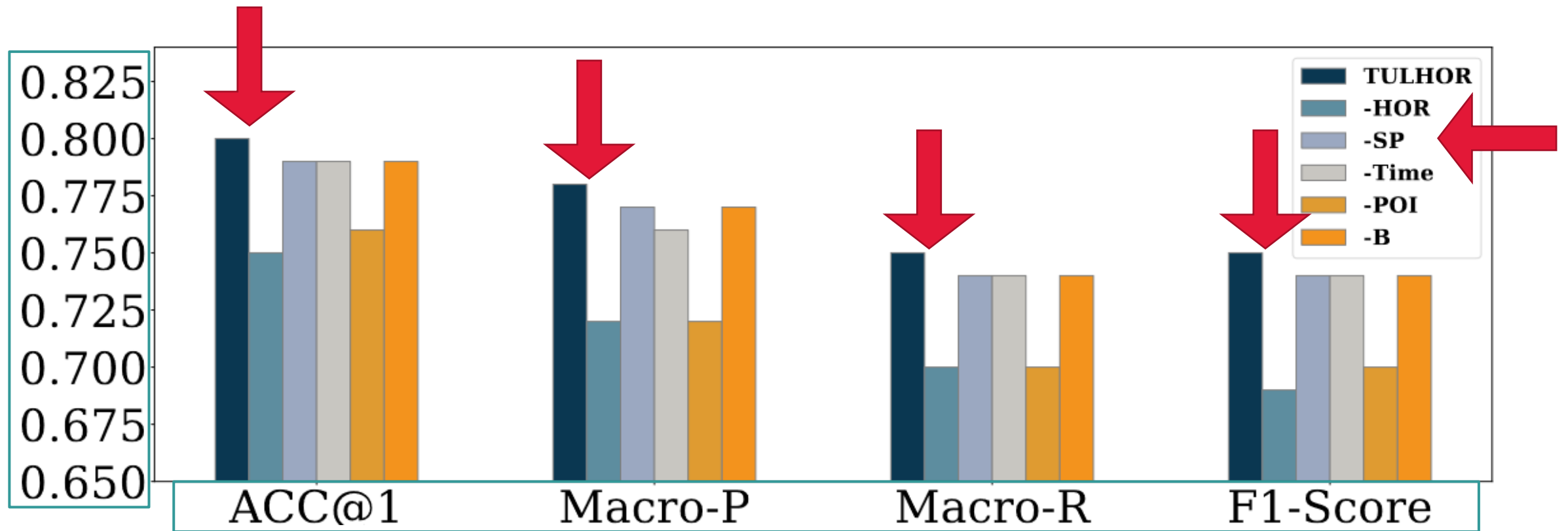
MODEL	FOURSQUARE-TKY														
	$ \mathcal{U}  = 108$					$ \mathcal{U}  = 209$					$ \mathcal{U}  = 451$				
	Acc@1	Acc@5	P	R	F1	Acc@1	Acc@5	P	R	F1	Acc@1	Acc@5	P	R	F1
DT	0.789	0.793	0.785	0.777	0.775	0.658	0.664	0.629	0.615	0.613	0.522	0.525	0.446	0.437	0.431
LDA	0.853	0.912	0.927	0.847	0.874	0.722	0.808	0.778	0.692	0.713	0.574	0.720	0.553	0.501	0.495
LINEAR-SVM	0.890	0.948	0.923	0.886	0.898	0.769	0.878	0.794	0.736	0.748	0.609	0.761	0.610	0.539	0.550
TULER	0.870	0.933	0.871	0.860	0.860	0.768	0.864	0.762	0.735	0.736	0.637	0.74	0.588	0.554	0.548
TULER-L	0.905	0.952	0.904	0.898	0.897	0.848	0.911	0.837	0.825	0.824	0.739	0.827	0.708	0.675	0.675
TULER-G	0.915	0.954	0.916	0.910	0.909	0.851	0.911	0.842	0.824	0.825	0.738	0.823	0.701	0.672	0.671
ATT-LSTM	0.908	0.966	0.916	0.901	0.908	0.752	0.871	0.795	0.729	0.760	0.407	0.584	0.362	0.326	0.343
ATT-GRU	0.933	0.975	0.932	0.928	0.930	0.869	0.937	0.872	0.856	0.864	0.742	0.821	0.715	0.689	0.695
DEEPTUL	0.922	0.966	0.927	0.913	0.920	0.773	0.904	0.820	0.747	0.782	0.660	0.790	0.631	0.587	0.608
<b>TULHOR</b>	<b>0.939</b>	<b>0.973</b>	<b>0.937</b>	<b>0.934</b>	<b>0.933</b>	<b>0.893</b>	<b>0.953</b>	<b>0.883</b>	<b>0.877</b>	<b>0.875</b>	<b>0.801</b>	<b>0.888</b>	<b>0.783</b>	<b>0.755</b>	<b>0.752</b>
<b>Improvement</b>	<b>0.58%</b>	<b>-0.26%</b>	<b>0.59%</b>	<b>0.71%</b>	<b>0.37%</b>	<b>2.7%</b>	<b>1.77%</b>	<b>1.33%</b>	<b>2.53%</b>	<b>1.30%</b>	<b>7.86%</b>	<b>7.47%</b>	<b>9.52%</b>	<b>9.53%</b>	<b>8.11%</b>

TULHOR **outperforms** every baseline

TULHOR has better **scalability**

MODEL	FOURSQUARE-NYC														
	$ \mathcal{U}  = 108$					$ \mathcal{U}  = 209$					$ \mathcal{U}  = 234$				
	Acc@1	Acc@5	P	R	F1	Acc@1	Acc@5	P	R	F1	Acc@1	Acc@5	P	R	F1
<b>DT</b>	0.884	0.892	0.878	0.867	0.868	0.785	0.788	0.753	0.728	0.730	0.778	0.782	0.722	0.712	0.705
<b>LDA</b>	0.822	0.851	<u>0.962</u>	0.810	0.868	0.746	0.781	0.791	0.687	0.718	0.696	0.752	0.724	0.615	0.650
<b>LINEAR-SVM</b>	0.873	0.929	<b>0.966</b>	0.878	<u>0.909</u>	0.776	0.839	0.785	0.702	0.727	0.731	0.798	0.724	0.628	0.657
<b>TULER</b>	0.870	0.929	0.869	0.851	<u>0.852</u>	0.776	0.853	0.749	0.722	0.718	0.768	0.844	0.733	0.707	0.703
<b>TULER-L</b>	0.903	0.942	0.904	0.890	0.890	0.847	<u>0.898</u>	0.828	0.803	0.807	0.845	0.889	<u>0.821</u>	<u>0.806</u>	<u>0.803</u>
<b>TULER-G</b>	<u>0.909</u>	<u>0.949</u>	0.914	<u>0.897</u>	0.898	<u>0.854</u>	<u>0.892</u>	<u>0.835</u>	<u>0.811</u>	<u>0.812</u>	0.846	0.891	0.821	0.805	<u>0.803</u>
<b>ATT-LSTM</b>	0.823	<u>0.896</u>	0.715	0.703	0.709	0.716	0.832	0.554	0.559	0.556	0.712	0.830	0.569	0.557	0.563
<b>ATT-GRU</b>	0.886	0.933	0.779	0.779	0.791	0.835	0.891	0.663	0.680	0.671	<u>0.889</u>	<b>0.936</b>	0.741	0.738	0.740
<b>DEEPTUL</b>	0.853	0.923	0.765	0.738	0.751	0.733	0.840	0.614	0.597	0.606	0.789	0.891	0.607	0.617	0.612
<b>TULHOR</b>	<b>0.940</b>	<b>0.966</b>	0.938	<b>0.931</b>	<b>0.932</b>	<b>0.903</b>	<b>0.943</b>	<b>0.890</b>	<b>0.877</b>	<b>0.876</b>	<b>0.892</b>	<u>0.932</u>	<b>0.876</b>	<b>0.864</b>	<b>0.860</b>
<b>Improvement</b>	3.42%	1.85%	-2.89%	3.85%	2.53%	5.82%	5.07%	6.58%	7.83%	7.87%	0.35%	-0.49%	6.61%	7.13%	7.19%

# Ablation study



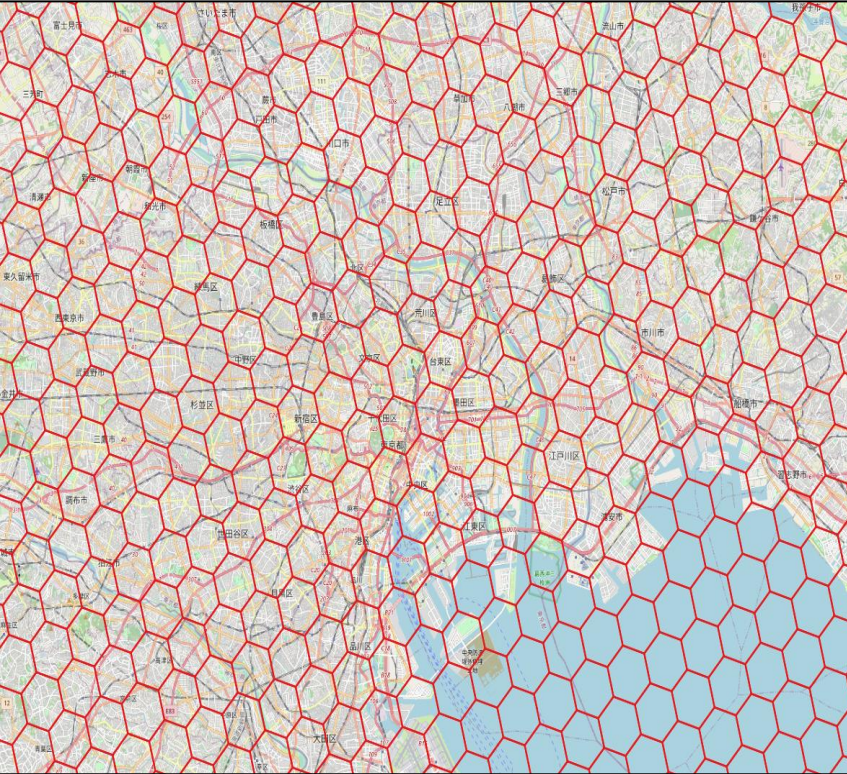
Removing **Higher-order significantly reduces** the performance



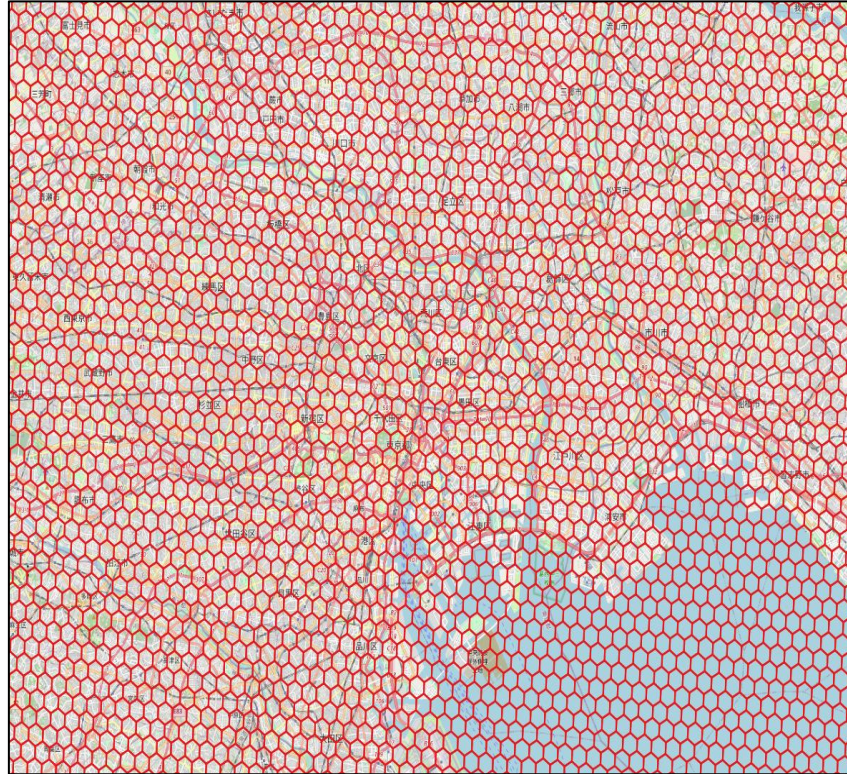
## Tessellation granularity (grid size) effect

<b>RESOLUTION</b>	<b># OF CELLS</b>	<b>CELL SIZE (<math>km^2</math>)</b>
<b>HEX@7</b>	334	5.160
<b>HEX@8</b>	2,003	0.730
<b>HEX@9</b>	11,036	0.015

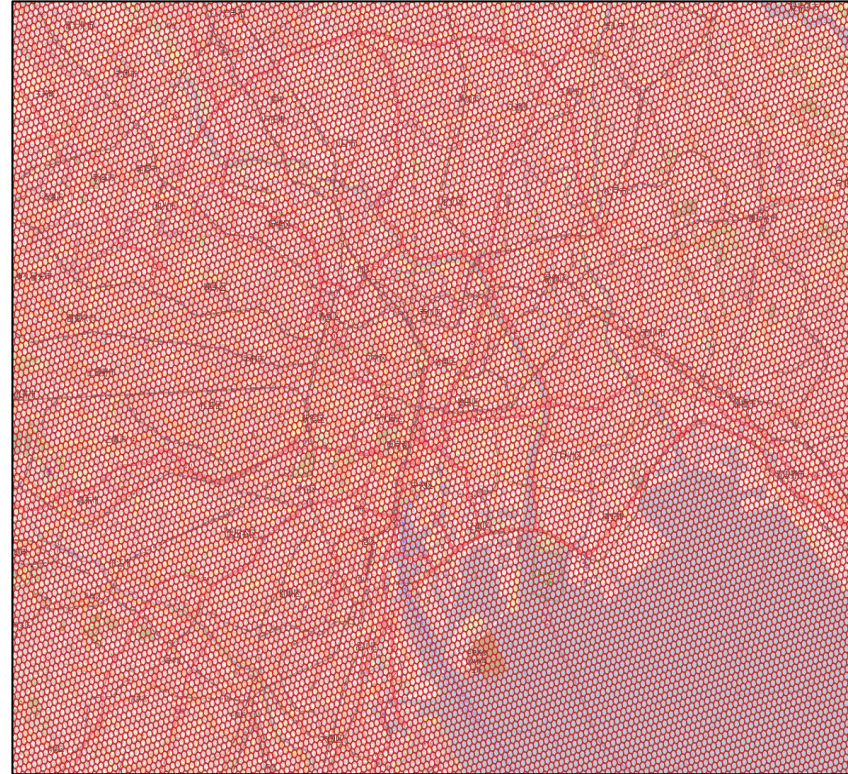
# Tessellations of Tokyo



Hex@7



Hex@8



Hex@9

# Results of grid size study

METHOD	FOURSQUARE-TKY														
	#USERS = 108					#USERS = 209					#USERS = 451				
	ACC@1	ACC@5	P	R	F1	ACC@1	ACC@5	P	R	F1	ACC@1	ACC@5	P	R	F1
<b>HEX@7</b>	0.923	0.971	0.920	0.911	0.913	0.868	0.943	0.832	0.817	0.815	0.711	0.883	0.734	0.734	0.711
<b>HEX@8</b>	0.926	<b>0.977</b>	0.925	0.917	0.917	0.868	0.940	0.862	0.849	0.849	0.790	0.884	0.753	0.740	0.733
<b>HEX@9</b>	<b>0.939</b>	0.973	<b>0.937</b>	<b>0.934</b>	<b>0.933</b>	<b>0.893</b>	<b>0.953</b>	<b>0.883</b>	<b>0.877</b>	<b>0.875</b>	<b>0.801</b>	<b>0.888</b>	<b>0.783</b>	<b>0.755</b>	<b>0.752</b>

Hex@9 **outperforms** other **sizes** as the number of users increases

The **smaller the cells** are, the **better the scalability**

## Computing Environment

- Implemented in Python
- Run in a GPU cluster (lion)

Conclusion

## Summary

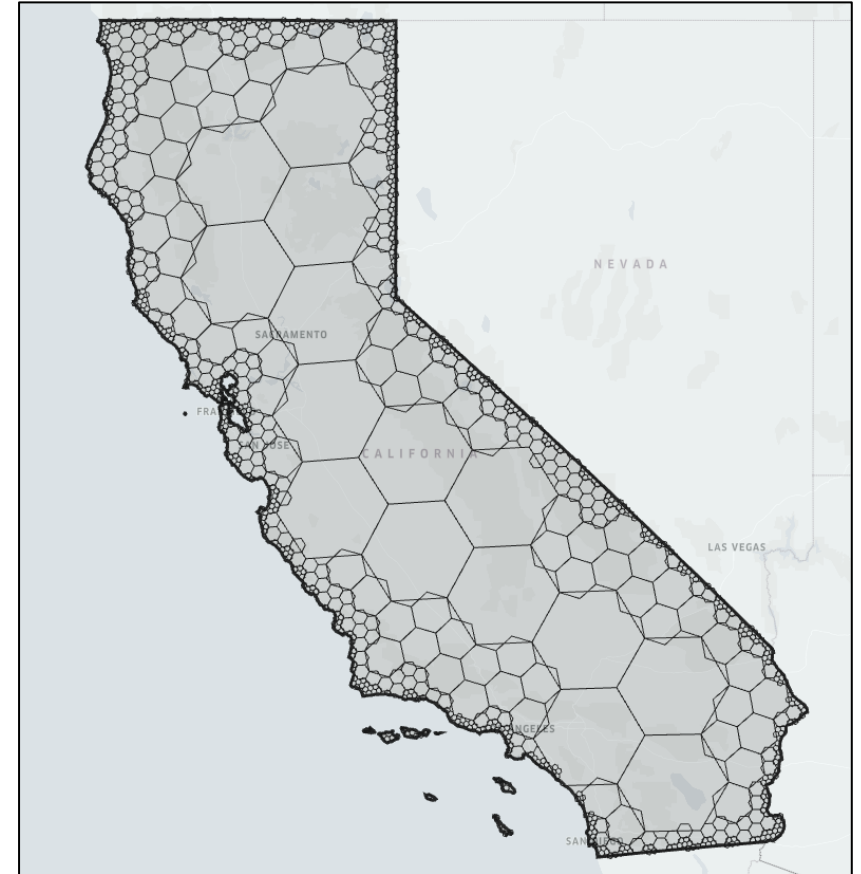
- Method for generating higher-order mobility flow & check-ins
- Novel deep learning framework – TULHOR (trajectory-user linking using higher-order representations )
- Extensive set of experiments

## Contributions

- Model agnostic method for dealing with sparsity and low data quality
- Implementation of multiple TUL models
- Publicly available dataset
- Accepted to **The 24<sup>th</sup> IEEE International Conference on Mobile Data Management (MDM)**

## Future works

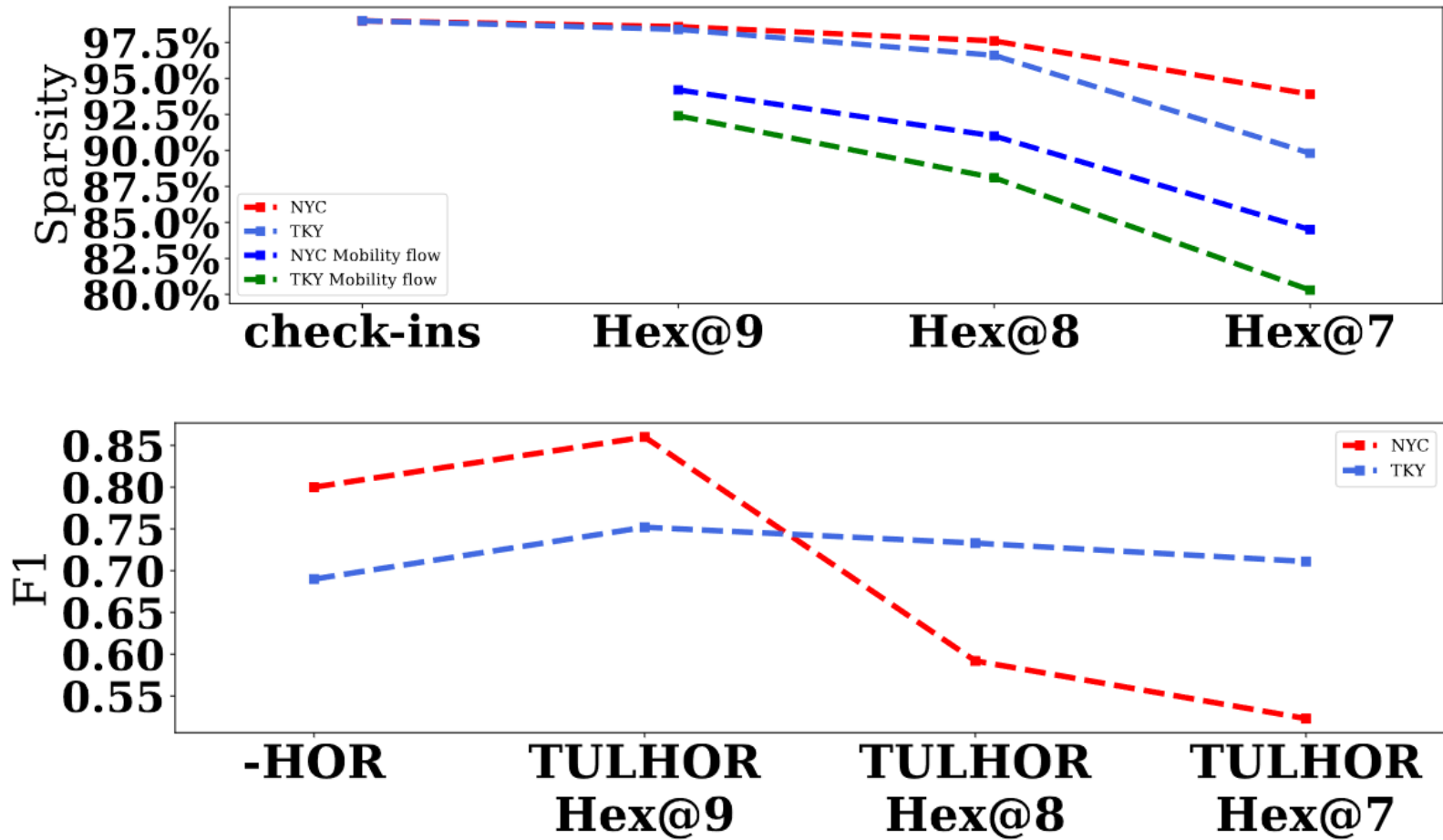
- Multi-trajectory User Linking
- POI recommendation
- Mixed hierarchical representation







# Impact of Sparsity



## Balanced loss

$$\mathcal{L}(Tr, u_i) = \frac{1 - \beta}{1 - \beta^{n_{u_i}}} \log(\sigma(y'))$$