

Data-Driven Similarity-based Worker Recruitment Towards Multi-task Data Inference for Sparse Mobile Crowdsensing



En Wang*, Zijie Tian*, Yongjian Yang*, Wenbin Liu*, Baoju Li*, Nan Jiang‡ and Jie Wu§

*Jinlin University, ‡East China Jiaotong University , § Temple University

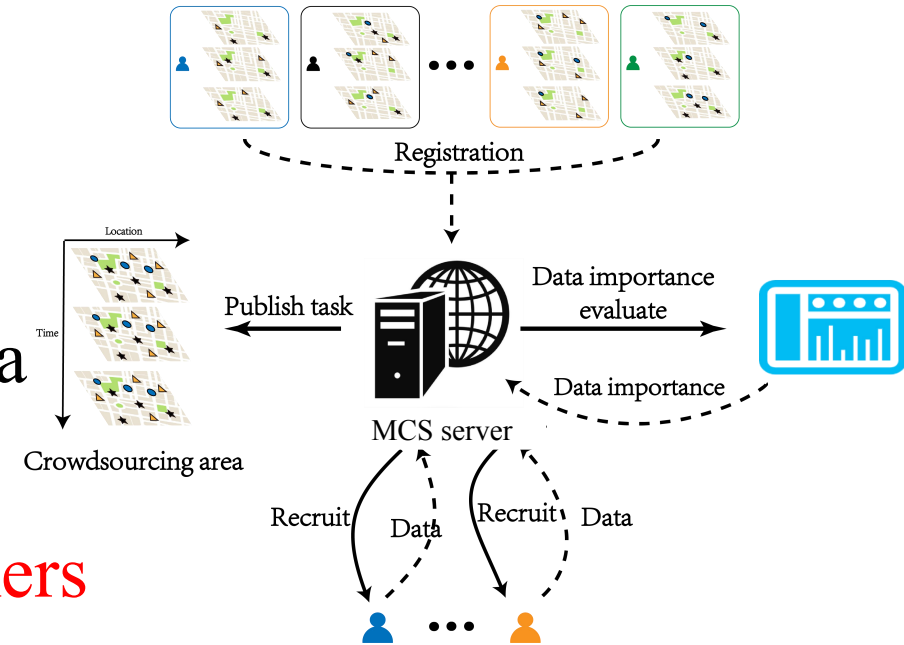
I. Background

II. Challenge

III. Method

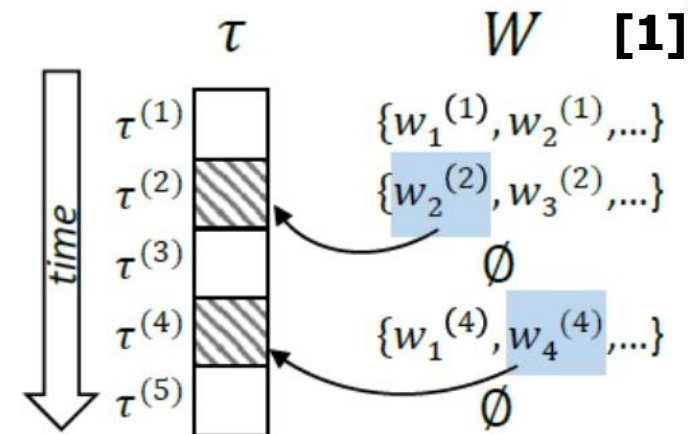
IV. Experiments

- Crowdsensing recruits workers to collect various data
- The limitation of budget & workers
- Inferring unsensed areas using a little sensed data
- Finding important data > Recruiting useful workers



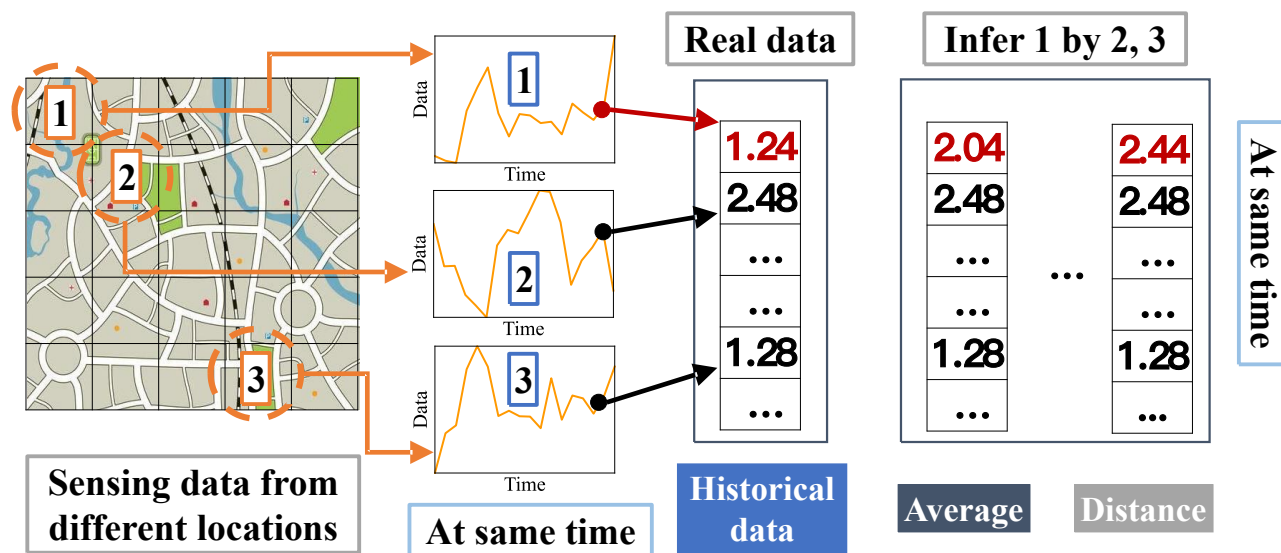
Existing work:

- Spatial-temporal distance between sensed data.
- Using entropy to evaluate importance and select workers.



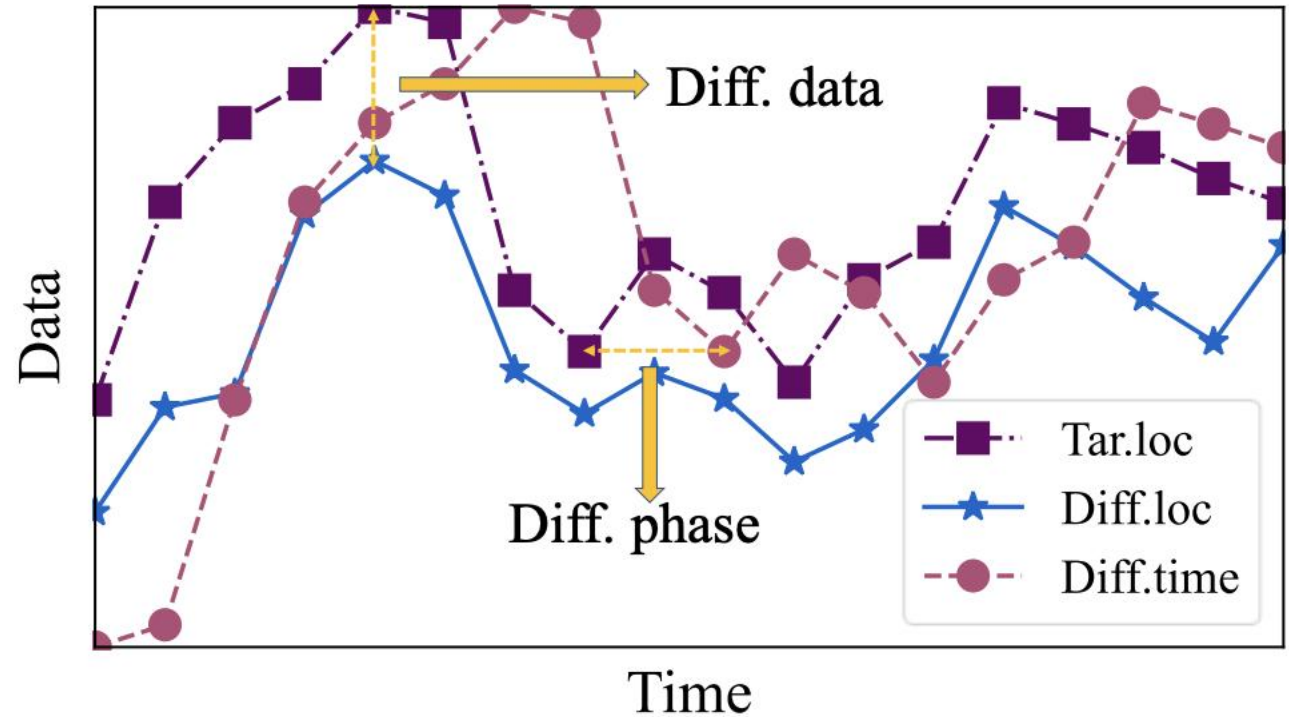
Disadvantage:

- Distance vs. Historical data knowledge



[1] On Efficient and Scalable Time-Continuous Spatial Crowdsourcing

- Diff. data & Diff. phase
- Data reliability

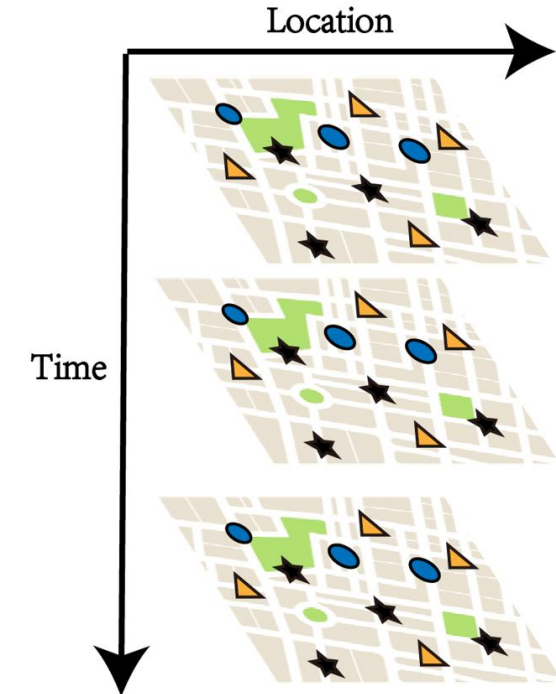


- How to calculate similarity according to the historical data?

- Multiple data-types in practice

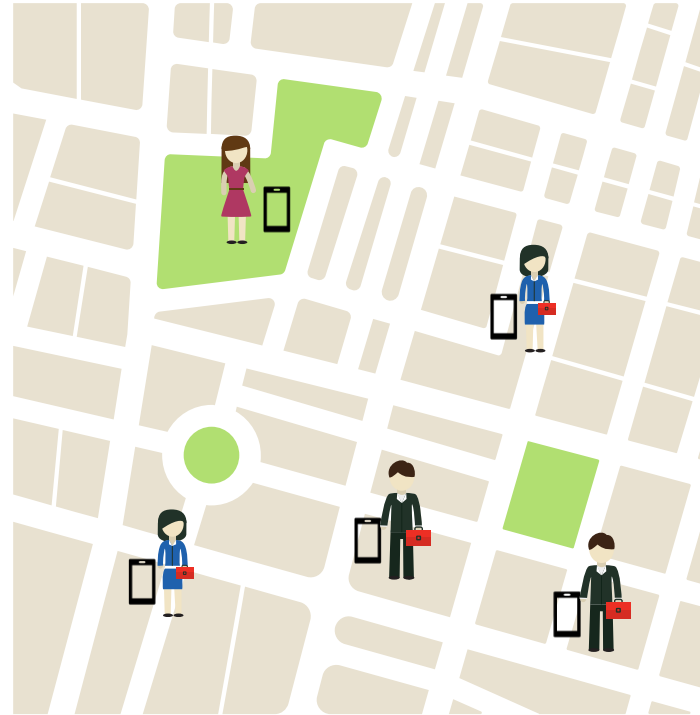
LVTH				
	Light	Voltage	Humidity	Temperature
Country/Region - City	Intel Berkeley Research Lab - California			
Subarea	33 subareas			
Time slots	384 time slots each with half an hour			
Mean \pm Std.	546.96 \pm 649.31 (<i>lux</i>)	2.62 \pm 0.08 (<i>v</i>)	35.74 \pm 7.06 (%)	22.68 \pm 7.08 ($^{\circ}$ C)

	PM		HT	
	PM2.5	PM10	Humidity	Temperature
Country/Region - City	Chinese mainland - Beijing		Switzerland - Lausanne	
Subarea	36 subareas		57 subareas	
Time slots	264 time slots each with an hour		336 time slots each with half an hour	
Mean \pm Std.	79.11 \pm 81.21 (μ g/m ³)	63.12 \pm 48.56 (μ g/m ³)	84.52 \pm 6.32 (%)	6.04 \pm 1.87 ($^{\circ}$ C)



- How to measure similarity and use it to evaluate data importance in three dimensional (spatial, temporal and data-type) ?

- Worker unreliability^[2]
- Multiple equipped sensors
- Platform budget limitation



Thermometer

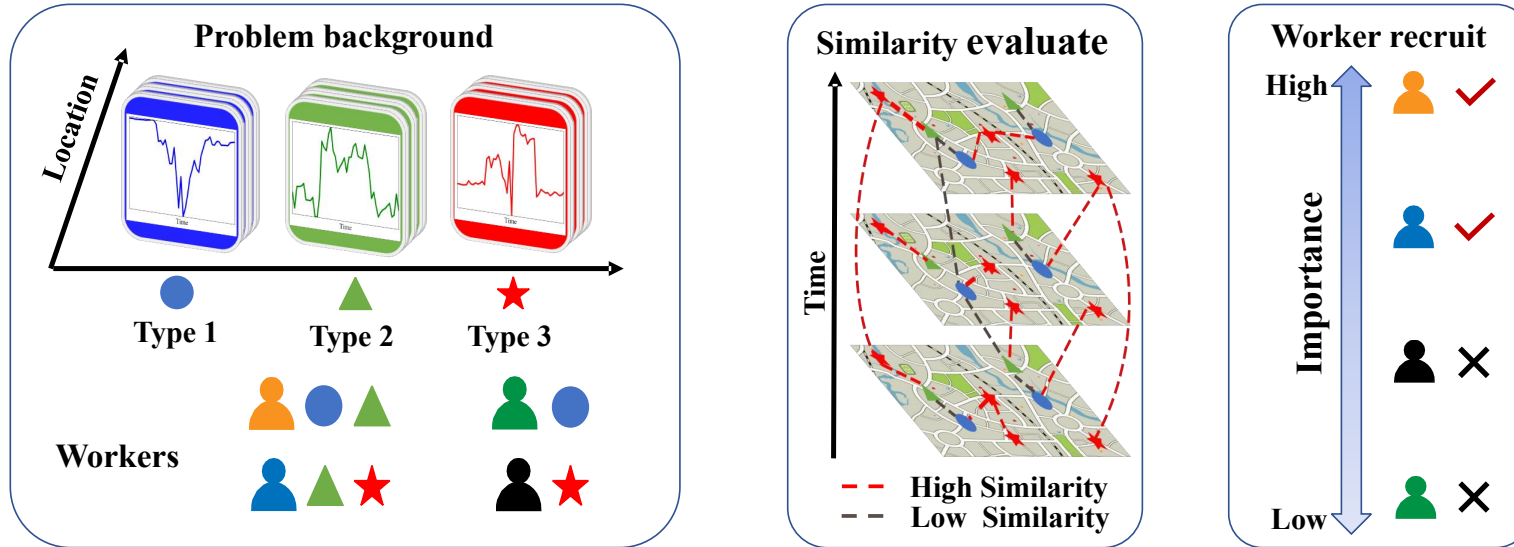


Air Quality Sensor

.....

- How to select a group of suitable workers?

[2] Reliable Diversity-Based Spatial Crowdsourcing by Moving Workers

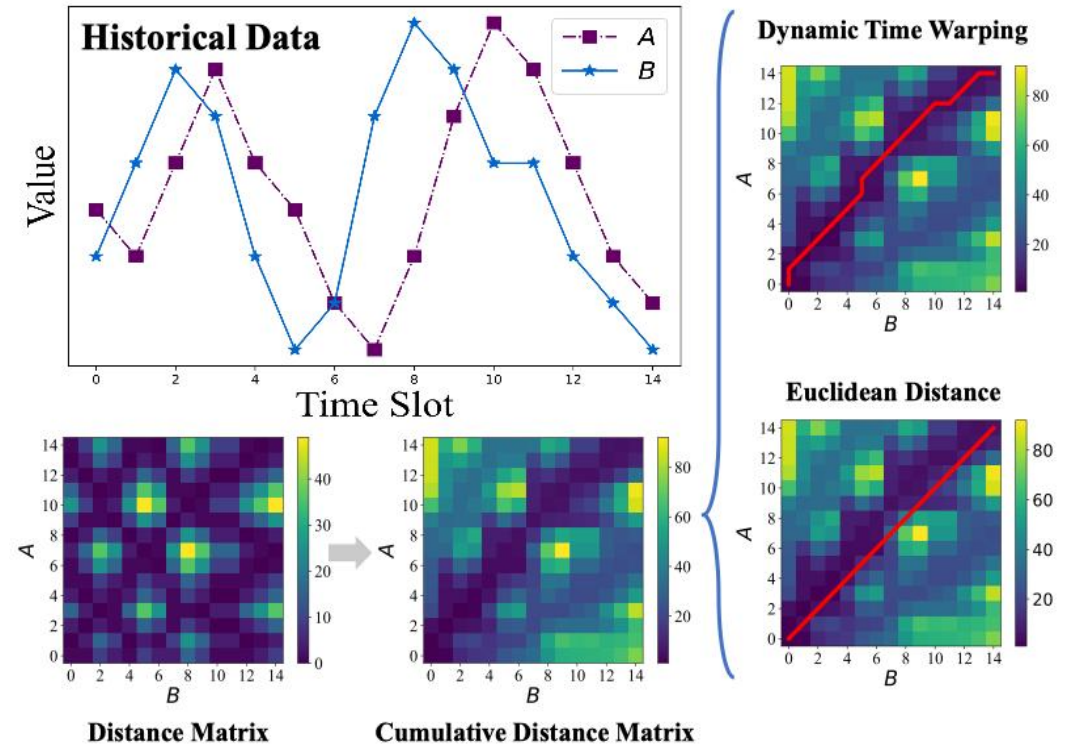


- SWDTW to evaluate time-series similarity
- Entropy-Weighted method to calculate similarity in three dimensional
- WRGSA to recruit workers according reliability and equipped sensors

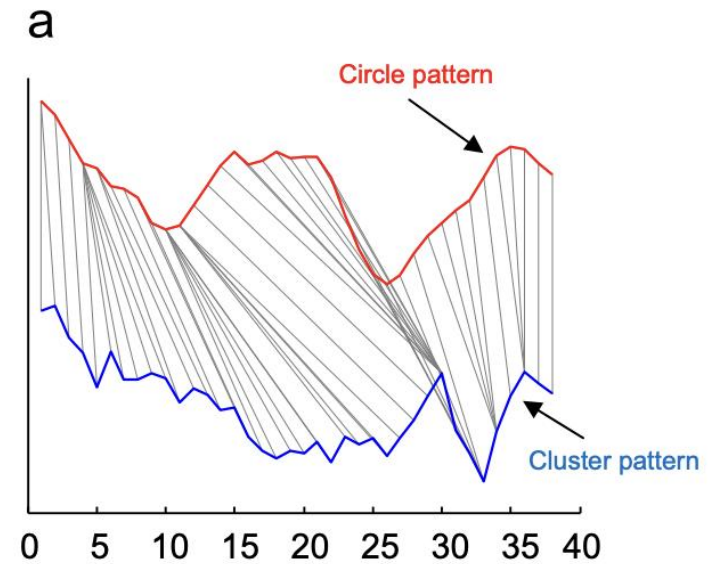
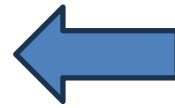
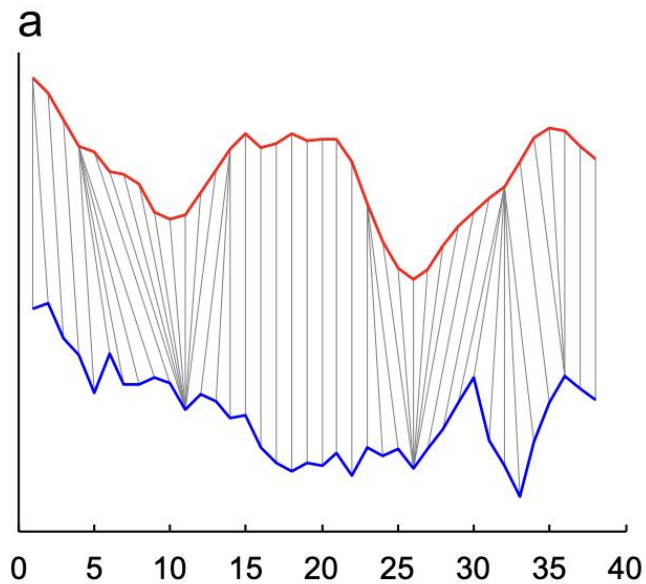
- Similarity-Weighted Dynamic Time Warping (SWDTW)
 - Dynamic Time Warping (DTW)

Euclidean Distance

$$\Upsilon_{n_1, n_2} = dis_{n_1, n_2} + \min\{\Upsilon_{m, n-1}, \Upsilon_{m-1, n}, \Upsilon_{m-1, n-1}\}$$



- Similarity-Weighted Dynamic Time Warping (SWDTW)
 - Weighted Dynamic Time Warping (WDTW)

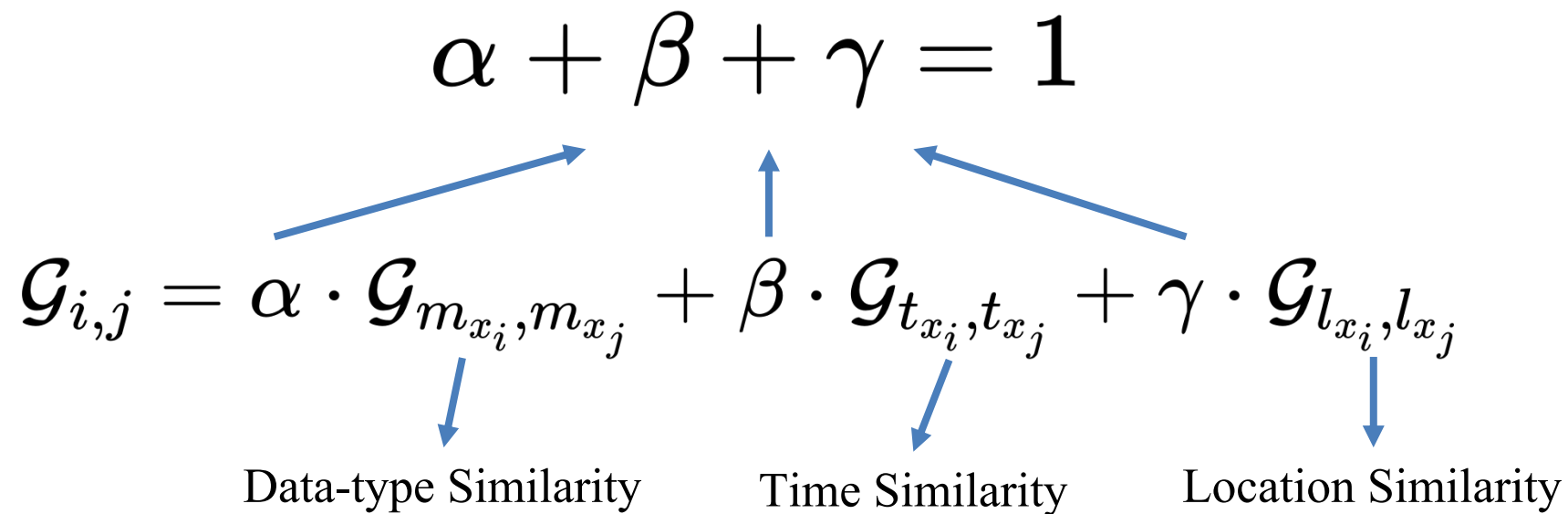


$$\text{dis}_{n_1, n_2} = \frac{\omega_{max}}{1 + \exp(-g(|n_1 - n_2| - m_c))} \times (a_{n_1} - b_{n_2})^2$$

- 三维空间相似性

$$\alpha + \beta + \gamma = 1$$
$$\mathcal{G}_{i,j} = \alpha \cdot \mathcal{G}_{m_{x_i}, m_{x_j}} + \beta \cdot \mathcal{G}_{t_{x_i}, t_{x_j}} + \gamma \cdot \mathcal{G}_{l_{x_i}, l_{x_j}}$$

Data-type Similarity Time Similarity Location Similarity



- 子任务选取

$$err(x_i) = \frac{\sum_{j \in S_k(x_i)} (1 - \mathcal{G}_{i,j})}{k}$$

$$kr_i = \frac{1}{|X|} (1 - err(x_i)) \quad \text{The knowledge of } x_i \text{ for inferring. } [0, \frac{1}{|X|}]$$

$$Q(X|\zeta) = - \sum_{\forall i \in X} kr_i \log_2 kr_i \quad \text{The data quality[1] of } X \text{ with selected datapoints[1].}$$

For $x \in X - \zeta$:

$$\Delta_x = Q(X|\zeta + x) - Q(X|\zeta)$$

[1] On Efficient and Scalable Time-Continuous Spatial Crowdsourcing

- 参与者招募

- 参与者集合: $w \in W = (p_w, c_w, \mathbf{S}_w)$

- 子任务可靠性: $rel(x_i, W_{x_i}) = 1 - \prod_{w_j \in W_{x_i}} (1 - p_{w_j})$

$$d'_{i,j} = (1 - G_{i,j}) \cdot (1 - rel(j, W_j))$$

$$err(x_i) = \frac{\sum_{j \in S'_k(x_i)} d'_{i,j}}{k}$$

$$kr_i = \frac{1}{|X|} \left(\frac{\sum_{j \in S'_k(x_i)} rel(j, W_j)}{k} - err(x_i) \right)$$

• Worker Recruitment

Algorithm 3 WRGSA

Input: $B, W, X, T, T_{max}, \alpha;$

Output: ε : the set of recruited workers

$R = \{rel(x_1, W_{x_1}), rel(x_2, W_{x_2}), \dots, rel(x_n, W_{x_n})\}$: the reliabilities of all datapoints

- 1: Initialize ε^* by algorithm 2
 - 2: **while** stop condition not met **do**
 - 3: generate a new solution ε' from ε^* by Algorithm 4
 - 4: **if** $Q(X|\varepsilon') > Q(X|\varepsilon^*)$ **then**
 - 5: $\varepsilon^* \leftarrow \varepsilon'$
 - 6: **else**
 - 7: $\varepsilon^* \leftarrow \varepsilon'$ with probability $\mathcal{J}_{\varepsilon^*, \varepsilon'}$
 - 8: **end if**
 - 9: **if** $Q(X|\varepsilon^*) > Q(X|\varepsilon)$ **then**
 - 10: $\varepsilon \leftarrow \varepsilon^*, T_b \leftarrow T$
 - 11: calculate R according Eq.(1)
 - 12: **end if**
 - 13: $T \leftarrow \alpha \times T$
 - 14: **if** $T < 0.01$ **then**
 - 15: $T_b \leftarrow 2 \times T_b, T \leftarrow \min \{T_b, T_{max}\}$
 - 16: **end if**
 - 17: **end while**
-

Algorithm 4 Generate New Solution

Input: $B, W, X, cost$

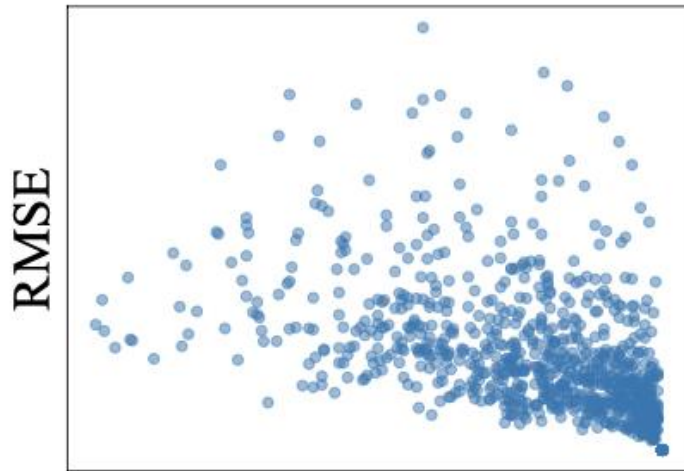
R : the reliabilities of all Data

ε^* : a set of already recruited workers

Output: ε' : the new set of recruited workers

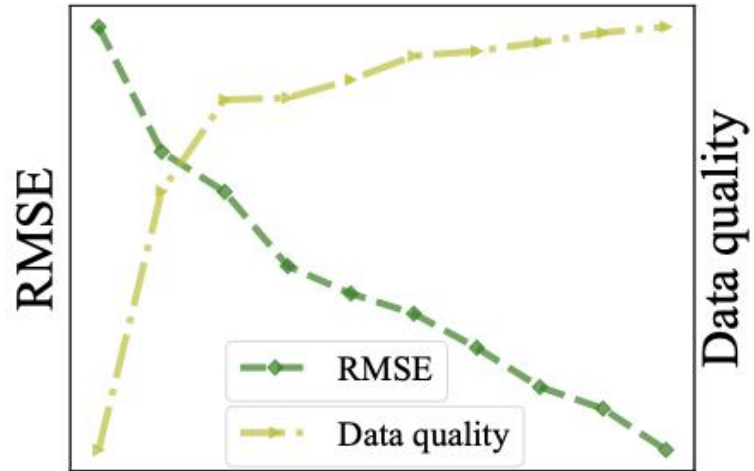
- 1: $\varepsilon' \leftarrow \varepsilon^*, S_{min} \leftarrow \emptyset, \mathcal{L} \leftarrow \emptyset$
 - 2: **for each** $\omega_i \in \varepsilon^*$ **do**
 - 3: compute $\frac{Q(X|\varepsilon' - \omega_i) - Q(X|\varepsilon')}{c_{\omega_i}}$ as Δ_{ω_i}
 - 4: $S_{min} \leftarrow S_{min} \cup (\omega_i, \Delta_{\omega_i})$
 - 5: **end for**
 - 6: sort S_{min} by ascending order of Δ_{ω}
 - 7: **for each** $\omega \in S_{min}$ with its rank i **do**
 - 8: **if** ω is removed with \mathcal{P}_i **then**
 - 9: $\varepsilon' \leftarrow \varepsilon' - \omega, cost \leftarrow cost + c_{\omega}$
 - 10: add ω in \mathcal{L}
 - 11: **end if**
 - 12: **end for**
 - 13: Select workers \mathcal{W} from $\subseteq W - \varepsilon' - \mathcal{L}$ by greedy
 - 14: Add \mathcal{W} in ε
-

- Datapoint Selection
 - Similarity based data quality vs. Inference error



Data quality

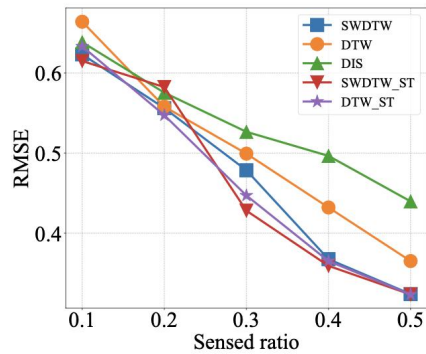
(a) Scatter



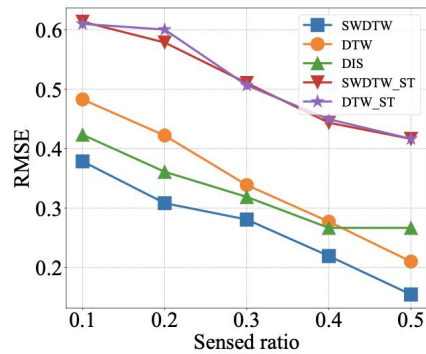
Sensed ratio

(b) Line chart

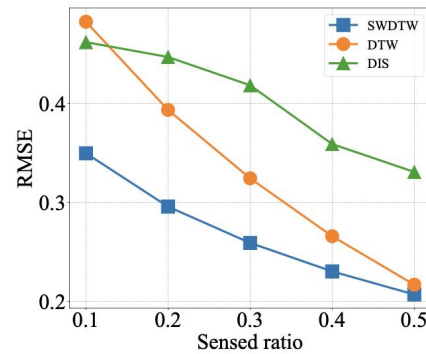
- Datapoint Selection
 - The performance in datasets



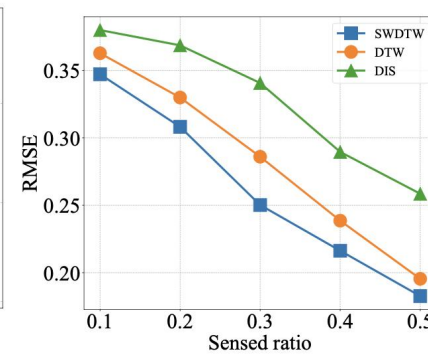
(a) PM10



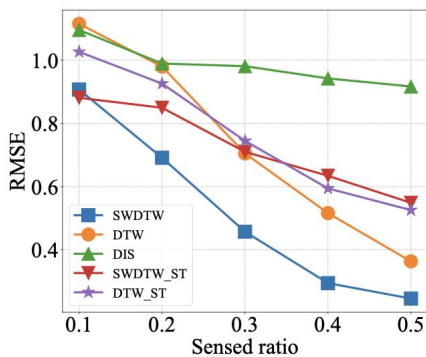
(b) PM2.5



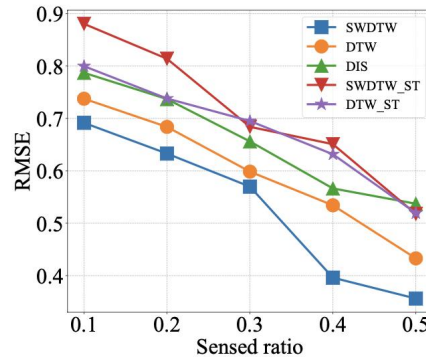
(c) Humidity



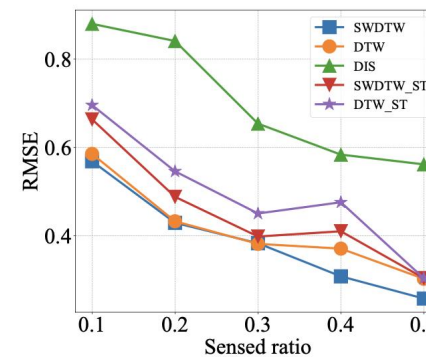
(d) Temperature



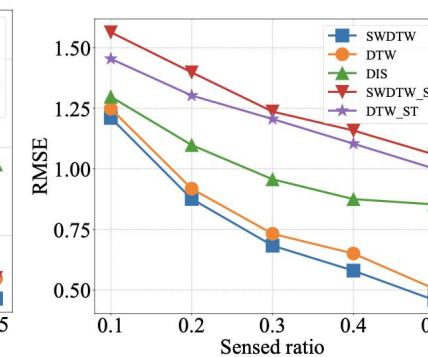
(a) Light



(b) Voltage

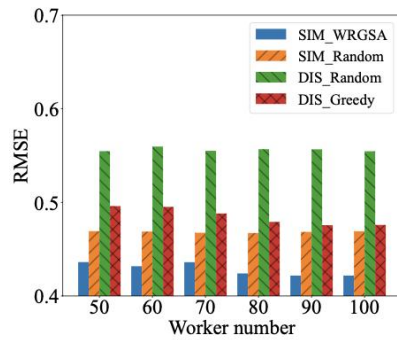


(c) Temperature

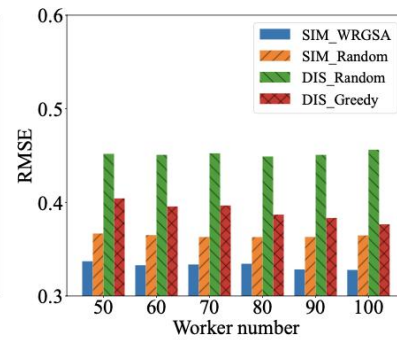


(d) Humidity

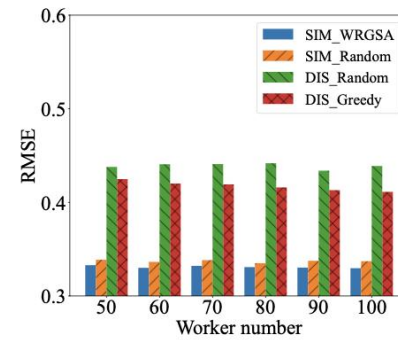
- Worker selection
 - Increasing worker number



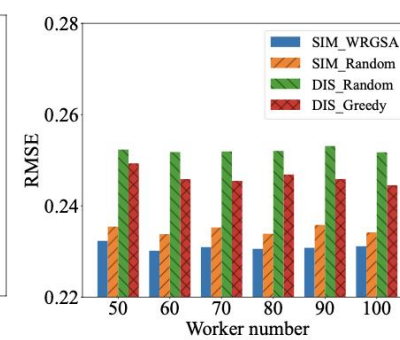
(a) PM2.5



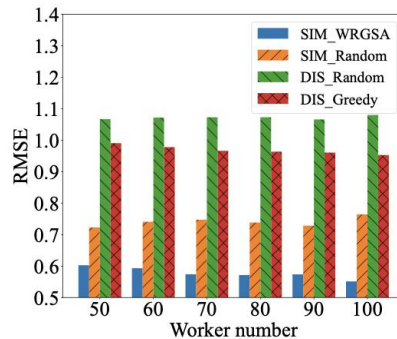
(b) PM10



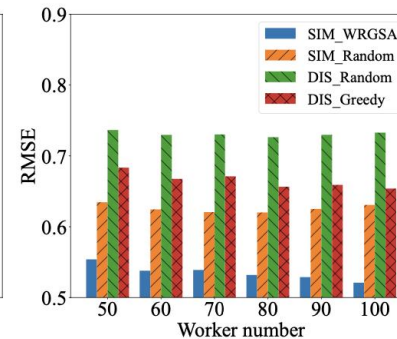
(c) Humidity



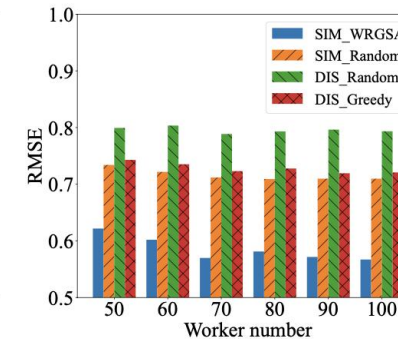
(d) Temperature



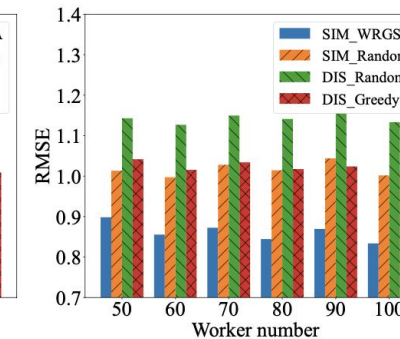
(a) Light



(b) Voltage

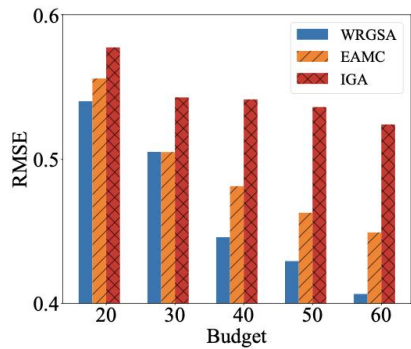


(c) Temperature

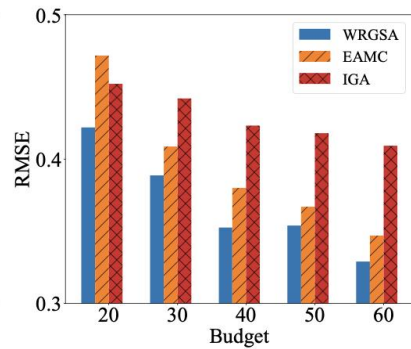


(d) Humidity

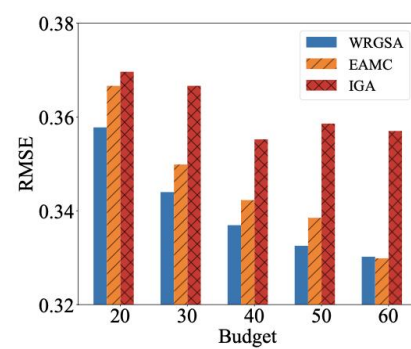
- Worker selection
 - Increasing budget



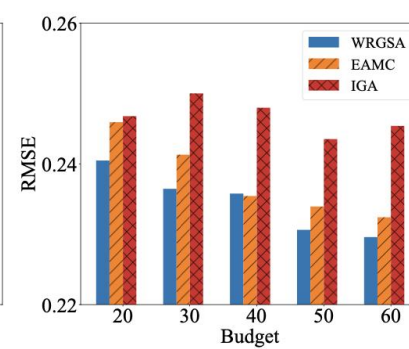
(a) PM2.5



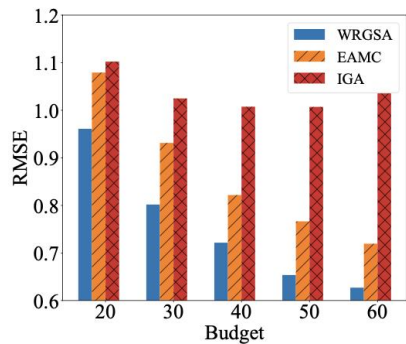
(b) PM10



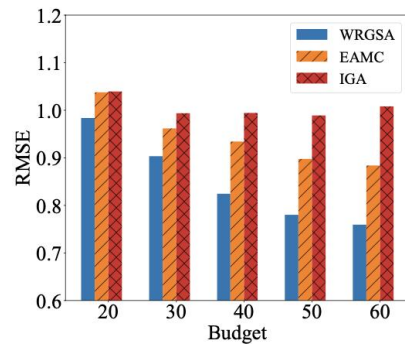
(c) Humidity



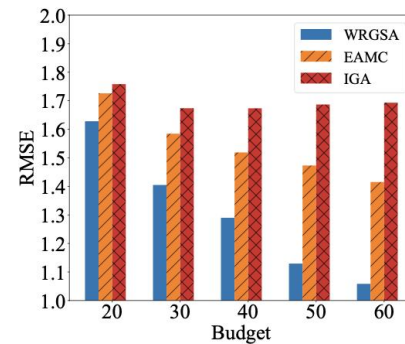
(d) Temperature



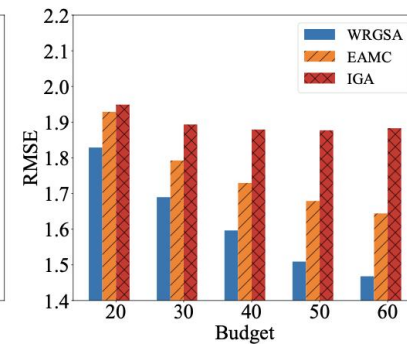
(a) Light



(b) Voltage



(c) Temperature



(d) Humidity



Thank you!
