

Urban-MAS: Human-Centered Urban Prediction with LLM-Based Multi-Agent System

Shangyu Lou

shangyulou@ucsb.edu, slou4820@sdsu.edu

University of California, Santa Barbara & San Diego State University
California, USA

Abstract

Urban Artificial Intelligence (Urban AI) has advanced human-centered urban tasks such as perception prediction and human dynamics. Large Language Models (LLMs) can integrate multimodal inputs to address heterogeneous data in complex urban systems but often underperform on domain-specific tasks. Urban-MAS, an LLM-based Multi-Agent System (MAS) framework, is introduced for human-centered urban prediction under zero-shot settings. It includes three agent types: Predictive Factor Guidance Agents, which prioritize key predictive factors to guide knowledge extraction and enhance the effectiveness of compressed urban knowledge in LLMs; Reliable UrbanInfo Extraction Agents, which improve robustness by comparing multiple outputs, validating consistency, and re-extracting when conflicts occur; and Multi-UrbanInfo Inference Agents, which integrate extracted multi-source information across dimensions for prediction. Experiments on running-amount prediction and urban perception across Tokyo, Milan, and Seattle demonstrate that Urban-MAS substantially reduces errors compared to single-LLM baselines. Ablation studies indicate that Predictive Factor Guidance Agents are most critical for enhancing predictive performance, positioning Urban-MAS as a scalable paradigm for human-centered urban AI prediction. Code is available on the project website.¹

CCS Concepts

• **Computing methodologies** → **Spatial and physical reasoning**.

Keywords

Large Language Models, Urban AI, Multi-Agent System, Human-centered Urban Prediction

ACM Reference Format:

Shangyu Lou. 2025. Urban-MAS: Human-Centered Urban Prediction with LLM-Based Multi-Agent System. In *The 3rd ACM SIGSPATIAL International Workshop on Advances in Urban-AI (UrbanAI '25)*, November 3–6, 2025, Minneapolis, MN, USA. <https://doi.org/10.1145/3764926.3771951>

1 Introduction

Urban AI has played an increasingly important role in addressing diverse urban challenges. Human-centered studies—such as

urban perception [2, 16, 21] and human dynamics [3], deepening understanding of urban systems and support evidence-based policymaking to improve quality of life [3, 16, 21]. However, the complexity of urban systems makes selecting and representing predictive features difficult, limiting accuracy [4]. Large Language Models (LLMs)² have shown strong potential in integrating heterogeneous modalities like text and imagery [4]. However, it remains limited in domain-specific applications [7–10, 22]. LLMs often compress vast knowledge, making it difficult to model complex domain problems [13, 17], and single-LLM approaches struggle to handle the specialized and multifaceted requirements of urban tasks [5], leading to biased or incomplete outputs that may misinform policymaking. To overcome these limitations, researchers are increasingly adopting MAS [6, 11], where multiple LLM-based agents collaborate. Compared with single LLMs, MAS offers stronger specialization and fault tolerance, mitigating common issues such as hallucinations and insufficient domain expertise [5]. Through division of labor and collaborative reasoning, MAS demonstrates improved scalability and reasoning ability for complex urban tasks [6]. However, to the best of knowledge, no prior work has applied MAS approaches to enhance human-centered urban predictions, motivating this work to examine: To what extent do LLM-based multi-agent systems enhance human-centered urban prediction?

Existing approaches to enhance single-LLM performance mainly rely on fine-tuning [4, 17, 20], which demands extensive data and computation, or on chain-of-thought (CoT) prompting [14], which requires manually designed reasoning chains that are time-consuming and labor-intensive. In human-centered urban prediction such as urban perception prediction, much attention has been given to identifying the most important features that contribute to task accuracy [15]. yet none explicitly extract the most predictive influential factors to guide LLMs. Given the complexity of cities, determinants span multiple dimensions (social and built environments) and scales (macro and street levels) [3], making manual identification highly demanding. By contrast, MAS enables parallelized and specialized exploration[1] across dimensions, providing a more automatic, comprehensive, and efficient way to uncover task-relevant influential factors. While deep-research MAS has shown promise in general AI [1], it has not been applied to systematically identify such factors in human-centered urban prediction.

In addition, Some studies attempt to leverage task-related knowledge compressed within a single LLM [14, 17]. However, they lack mechanisms to verify and enhance the reliability of extracted information. Biased or inconsistent information, once incorporated into prediction, can cause errors, thereby reduce result reliability. MAS offers the advantages through role allocation and task

¹<https://github.com/THETUREHOOHA/UrbanMAS>



This work is licensed under a Creative Commons Attribution 4.0 International License. *UrbanAI '25, Minneapolis, MN, USA*

© 2025 Copyright held by the owner/author(s).

ACM ISBN 979-8-4007-2189-2/2025/11

<https://doi.org/10.1145/3764926.3771951>

²Also referred to as MLLMs; used interchangeably in this paper.

specialization[1, 6], it could enable simultaneous cross-validation and refinement during knowledge extraction, substantially improving credibility [5]. Nevertheless, this capability has yet to be applied to human-centered urban tasks.

To address these limitations and explore MAS potential in urban prediction, this study introduces Urban-MAS, a multi-agent system for human-centered urban tasks under zero-shot conditions. Urban-MAS integrates three classes of agents: (1) Predictive Factor Guidance Agents, which collaborate to prioritize influential factors and guide knowledge extraction; (2) Reliable UrbanInfo Extraction Agents, which improve robustness through multi-output comparison, consistency validation, and re-extraction when conflicts occur; and (3) Multi-UrbanInfo Inference Agents, which integrate multi-source information for prediction. This framework overcomes the constraints of single-LLM methods and advances LLM performance in human-centered urban prediction. The main contributions are:

- Presents the first MAS framework for human-centered urban prediction. Urban-MAS integrates three agent layers to enable prioritized factor selection, reliable urban information extraction, and integrated multi-source prediction under zero-shot conditions, outperforming single-LLM baselines across metrics.
- Introduces a MAS-based mechanism to improve the reliability of urban knowledge extraction. Reliable UrbanInfo Extraction Agents use output comparison, consistency checks, and selective re-extraction to mitigate bias and enhance predictive robustness.
- Evaluates performance on urban perception and human dynamics prediction across cities on three continents. Results demonstrate that Urban-MAS achieves efficient, low-cost, and significant zero-shot gains, advancing human-centered Urban AI research.

2 Urban-MAS

To enhance LLM performance in human-centered urban research, the Urban-MAS framework comprises three agent layers (Figure 1). Predictive Factor Guidance Agents identify the most relevant predictive factors, guiding extraction and improving the utility of compressed urban knowledge. Reliable UrbanInfo Extraction Agents enhance stability by generating multiple outputs, checking consistency, and re-extracting when needed to ensure trustworthy information. Multi-UrbanInfo Inference Agents integrate these refined multi-source signals across dimensions and scales to deliver robust, task-specific urban predictions.

Predictive Factor Guidance Agents. Urban prediction requires focusing on the most influential factors for each task and data source, as general prompts often yield noisy cues. To address this, the Predictive Factor Guidance Agents layer employs targeted *Deep-research* subagents based on the *opendeepresearch* framework from LangChain[12] to generate research-level reports on the most influential factors of the task. Then, the *Summary* subagents summarize the findings into concise predictive factors, organized by social and environmental dimensions and macro and street levels. Given an urban prediction task τ , the input is the task description. Dimensions are defined as $D = \text{Social, Built Environmental}$ and levels as $R = \text{Macro, Street}$. For each (d, r) pair under the task, deep-research

subagents produce a brief report $t_{d,r}$ containing six key predictive factors, which summary subagents compress into a predictive factor set $P_{d,r}$ as output, providing dimension- and level-specific guidance for the UrbanInfo Extraction Agents.

Reliable UrbanInfo Extraction Agents. Urban information from a single LLM call is often noisy or inconsistent, weakening inference. To enhance stability, each extraction agent generates two output variants via an *Extractor* subagent and compares them with an *Evaluator* subagent. Conflicting fields are selectively re-extracted by a *Refiner* subagent. For each location ℓ , four extraction agents are aligned to $P_{d,r}$: social-macro, social-street, environment-macro, and environment-street. Candidate outputs are compared by similarity metrics; if agreement exceeds a threshold, one is accepted, otherwise only conflicting fields are regenerated, yielding reliable outputs $U_{d,r}$. This dual-variant and conflict-repair mechanism ensures consistent UrbanInfo across dimensions and levels, strengthening multi-source integration. Specifically, for each dimension-scale pair, the *Evaluator* subagent compares two urban information independently generated variants (A and B) from *Extractor*. The two generated variants are normalized (lowercasing, punctuation removal, whitespace collapsing), and field-level hybrid soft similarity is computed as

$$\text{soft_sim}(a, b) = 0.4 \times \text{Jaccard}(a, b) + 0.6 \times \text{SequenceMatcher}(a, b), \quad (1)$$

where *Jaccard*[19] measures token overlap and *SequenceMatcher*[18] captures phrase-level alignment. A stability threshold of 0.72 is used: if similarity ≥ 0.72 , Variant A is accepted; otherwise, the *Refiner* regenerates only differing fields.

Multi-UrbanInfo Inference Agents. Urban prediction requires reasoning over complementary UrbanInfo across social and built dimensions and macro- and street-level scales. Even when each source is reliable, isolated reasoning risks imbalance. To address this, the inference layer employs an *LLM-based Inference Agent* that jointly processes the four reliable inputs— $U_{\text{social,macro}}^*$, $U_{\text{social,street}}^*$, $U_{\text{environment,macro}}^*$, and $U_{\text{environment,street}}^*$ —to infer task-specific outputs such as running amount or perception scores. The Inference Agent receives the four structured JSON inputs, enforces schema constraints (e.g., $\{\text{"running_amount"}: 0.0\}$). As the final Urban-MAS stage, this layer integrates reliable extraction to deliver robust, coherent predictions.

3 Experiments

Tasks. We evaluate Urban-MAS on two representative human-centered urban prediction tasks: (i) urban perception prediction, focusing on one positive (lively) and one negative (boring) dimension of perception, and (ii) human dynamics prediction via running amount estimation. These tasks capture complementary aspects of human–environment interaction by linking perceptual and behavioral responses to urban form.

Datasets. Experiments are conducted on 300 samples across Tokyo, Milan, and Seattle, covering three continents to assess cross-regional generalizability. For each sampled location, we process the raw geographic coordinates into location text via OpenStreetMap’s Nominatim API for reverse geocoding (convert coordinates to address), and further enrich the inputs by querying nearby points of interest using the Overpass API. In addition, street-view imagery is retrieved through the Google Maps API. These inputs are used as

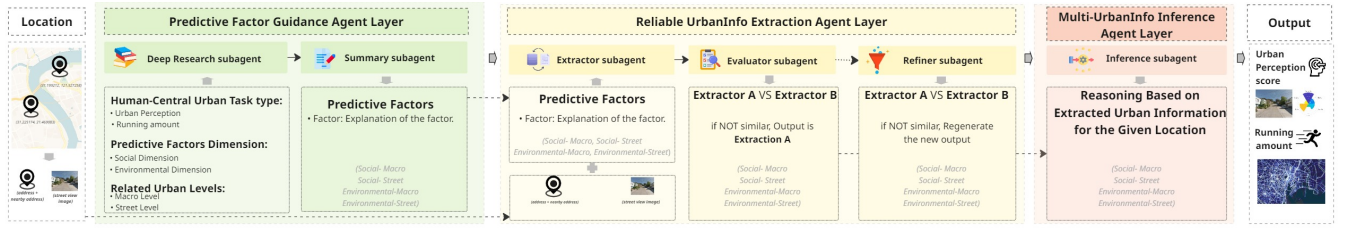


Figure 1: Urban-MAS comprises three agent layers: (Left) Predictive Factor Guidance Agents prioritize factor selection; (Middle) Reliable UrbanInfo Extraction Agents ensure reliable urban information extraction through consistency checks; (Right) Multi-UrbanInfo Inference Agents integrate multi-source urban information for robust urban prediction.

multi-source data for Urban-MAS across different human-centered tasks. Although the validation dataset is relatively small, its scale is constrained by the cost of collecting and verifying multi-source urban data (street-view imagery, POI queries), yet remains sufficient for evaluating the quality of the LLM-based solution using several hundred representative examples[5].

For running amount prediction, we use Strava heatmaps, a widely recognized source of physical activity data. Following Le’s agent-based modeling approach, high-resolution Strava heatmaps were collected for each study area, and raster values were extracted in QGIS at each sampling point. Brightness values, representing running intensity, were rescaled to [0,10]. For urban perception, we adopt Place Pulse 2.0 [21], the largest benchmark of human perceptions of urban environments. Pairwise Google Street View comparisons were aggregated using the TrueSkill algorithm, yielding continuous perception scores that we rescaled to [0,10].

Models. The experimental setup was designed to evaluate the performance and capabilities of LLM-based multi-agent systems in human-centered urban perception and prediction. All experiments were conducted using GPT-5, a closed-source model representing the state of the art in its class. The exclusive use of LLM-based configurations aims to examine how multi-agent interaction improves efficiency, reasoning, and overall performance compared with a single-LLM baseline, as well as to assess the contribution of each agent component within the MAS framework. The proposed Urban-MAS framework comprises three layers of agents: Predictive Factor Guidance Agents, Reliable UrbanInfo Extraction Agents, and Multi-UrbanInfo Inference Agents. The Predictive Factor Guidance layer leverages the opendeepresearch framework from LangChain—ranked sixth on the Deep Research Bench Leaderboard with GPT-4o—to ensure consistency and relevance of the extracted predictive factors. The Reliable UrbanInfo Extraction and Multi-UrbanInfo Inference layers are implemented on GPT-5 in JSON mode using zero-shot prompting, with a single-LLM GPT-5 baseline included for comparison.

Evaluation. Performance is assessed using MAE, MSE, and RMSE for continuous outcomes (running amount and perception scores). We also report ablation studies to isolate the contributions of factor guidance and reliability enhancement.

3.1 Main Results

Table 1 summarizes the performance improvements achieved by the integrated Urban-MAS compared to the baseline single-LLM

Table 1: Performance gains with Urban-MAS with single LLM across multiple tasks (all metrics: lower is better).

Method	MAE	MSE	RMSE
<i>People Running Amount</i>			
Single LLM	2.99	13.73	3.70
Urban-MAS (Ours)	2.97 (↓0.73%)	13.20 (↓3.82%)	3.63 (↓1.93%)
<i>People Urban Perception (Boringness)</i>			
Single LLM	2.83	9.95	3.15
Urban-MAS (Ours)	2.05 (↓27.37%)	5.84 (↓41.33%)	2.42 (↓23.40%)
<i>People Urban Perception (Liveliness)</i>			
Single LLM	2.69	9.10	3.02
Urban-MAS (Ours)	1.73 (↓35.81%)	4.40 (↓51.67%)	2.10 (↓30.48%)

model (GPT-5). Relative to GPT-5, Urban-MAS achieves substantial reductions across all error metrics for every task, demonstrating its effectiveness in enhancing prediction performance for human-centered urban applications. A closer look at the baseline results further reveals variations in the degree of improvement across task types: in particular, error reductions are more pronounced for the safety perception task than for the running amount prediction task.

3.2 Ablation Study

Table 2: Ablation study on Predictive Factor Guidance and Reliable UrbanInfo Extraction (all metrics: lower is better).

Method Variant	MAE	MSE	RMSE
<i>People Running Amount</i>			
Urban-MAS	2.97	13.20	3.63
- PredictiveFactors	4.53 (↑52.84%)	26.80 (↑102.98%)	5.18 (↑42.47%)
- ReliabilityBoost	2.98 (↑0.30%)	13.39 (↑1.46%)	3.66 (↑0.73%)
<i>People Urban Perception (Boringness)</i>			
Urban-MAS	2.05	5.84	2.42
- PredictiveFactors	2.39 (↑16.37%)	7.52 (↑28.69%)	2.74 (↑13.44%)
- ReliabilityBoost	2.29 (↑11.52%)	6.98 (↑19.45%)	2.64 (↑9.29%)
<i>People Urban Perception (Liveliness)</i>			
Urban-MAS	1.73	4.40	2.10
- PredictiveFactors	2.54 (↑46.89%)	8.09 (↑83.79%)	2.84 (↑35.57%)
- ReliabilityBoost	2.21 (↑28.06%)	6.47 (↑47.02%)	2.54 (↑21.25%)

We conducted an ablation study to evaluate two strategies for improving LLM performance on human-centered urban tasks: (1)

Predictive Factor Guidance Agents, which prioritize influential predictive factors, and (2) Reliable UrbanInfo Extraction Agents, which enhance extraction reliability. As shown in Table 2, the full Urban-MAS configuration achieves the lowest errors across all tasks. Removing the Predictive Factor Guidance module causes a larger error increase than removing the Reliability module, indicating that prioritizing predictive factors is more critical. Removing the Predictive Factor Guidance module causes a larger error increase than removing the Reliability module, indicating that prioritizing predictive factors is more critical. Disabling factor prioritization (“–Predictive-Factors”) consistently raises errors, especially in running amount prediction, showing its importance for modeling activity patterns. Conversely, disabling reliability enhancement (“–ReliabilityBoost”) increases errors across all tasks, particularly in perception-related dimensions such as liveliness and boringness, highlighting the necessity of reliable urban information extraction.

3.3 Case Study

Predictive Factors Agent. Predictive factors were derived through an LLM-guided process across two dimensions (social and environmental) and two spatial scales (macro and street level). These structured factor sets serve as measurable descriptors guiding information extraction and inference within Urban-MAS. The complete list of predictors are available on the project website.

UrbanInfo Extraction Agents. The Reliable UrbanInfo Extraction Agents Layers extract consistent urban information for each location. When both variants generated from Extractor subagent, the Evaluator subagent confirms stability and retains one version directly. However, when factual discrepancies arise, the Refiner subagent performs conflict-only reconciliation—regenerating only inconsistent fields while preserving verified ones. This selective correction ensures factual reliability without redundant regeneration. Detailed examples are available on the project website.

4 Conclusion

This paper introduced Urban-MAS, a novel LLM-based multi-agent system for human-centered urban tasks with multi-source data inputs. Urban-MAS integrates automated prioritization of the most predictive factors to guide subsequent knowledge extraction, enhanced reliability in extracting urban information from multi-source data, and predictor agents that collaborate under a zero-shot setting to improve performance on human-centered urban tasks. Experimental results across multiple cities and tasks demonstrate that, compared with single-LLM baselines, Urban-MAS significantly reduces prediction error. Urban-MAS also provides methodological insights: prioritizing the most predictive factors is crucial for enhancing human-centered urban prediction tasks, and improving the reliability of extracting urban knowledge from LLMs is also essential. Future work will extend this framework by incorporating MAS-based automatic optimization of urban prediction performance and applying Urban-MAS to a broader range of urban tasks and larger test samples.

References

- [1] Anthropic. 2024. How we built our multi-agent research system. <https://www.anthropic.com/engineering/multi-agent-research-system>.

- [2] Ziqi Cui and Shangyu Lou. 2025. Syncperception: A Real-Time Urban Perception Prediction Tool Based on Graph Neural Networks. In *2025 Annual Modeling and Simulation Conference (ANNSIM)*.
- [3] Lin Dong, Hongchao Jiang, Wenjing Li, Bing Qiu, Hao Wang, and Waishan Qiu. 2023. Assessing impacts of objective features and subjective perceptions of street environment on running amount: A case study of Boston. *Landscape and Urban Planning* 235 (2023), 104756.
- [4] Jie Feng, Shengyuan Wang, Tianhui Liu, Yanxin Xi, and Yong Li. 2024. Urban-LLaVA: A Multi-modal Large Language Model for Urban Intelligence with Spatial Reasoning and Understanding. *arXiv preprint* (2024). arXiv:2406.05294
- [5] Anna Kalyuzhnaya, Sergey Mityagin, Elizaveta Lutsenko, Andrey Getmanov, Yaroslav Aksentkin, Kamil Fatkhiev, Kirill Fedorin, Nikolay O. Nikitin, Natalia Chichkova, Vladimir Vorona, and Alexander Boukhanovsky. 2025. LLM Agents for Smart City Management: Enhancing Decision Support Through Multi-Agent AI Systems. *Smart Cities* 8, 1 (2025), 19.
- [6] Zixuan Ke, Fangkai Jiao, Yifei Ming, Xuan-Phi Nguyen, Austin Xu, Do Xuan Long, Minzhi Li, Chengwei Qin, Peifeng Wang, Silvio Savarese, Caiming Xiong, and Shafiq Joty. 2025. A Survey of Frontiers in LLM Reasoning: Inference Scaling, Learning to Reason, and Agentic Systems. *arXiv preprint arXiv:2504.09037* (2025).
- [7] Zixuan Ke and Bing Liu. 2023. Continual Learning of Natural Language Processing Tasks: A Survey. arXiv:2211.12701
- [8] Zixuan Ke, Yifei Ming, and Shafiq Joty. 2025. NAACL2025 Tutorial: Adaptation of Large Language Models. arXiv:2504.03931 [cs.CL]
- [9] Zixuan Ke, Yifei Ming, Xuan-Phi Nguyen, Caiming Xiong, and Shafiq Joty. 2025. Demystifying Domain-adaptive Post-training for Financial LLMs. arXiv:2501.04961
- [10] Zixuan Ke, Yijia Shao, Haowei Lin, Tatsuya Konishi, Gyuhaek Kim, and Bing Liu. 2023. Continual Pre-training of Language Models. arXiv:2302.03241
- [11] Zixuan Ke, Austin Xu, Yifei Ming, Xuan-Phi Nguyen, Caiming Xiong, and Shafiq Joty. 2025. MAS-ZERO: Designing Multi-Agent Systems with Zero Supervision. *arXiv preprint arXiv:2505.14996* (2025).
- [12] LangChain. 2024. Open Deep Research. <https://blog.langchain.com/open-deep-research/>.
- [13] Xinjin Li, Yu Ma, Yangchen Huang, Xingqi Wang, Yuzhen Lin, and Chenxi Zhang. 2024. Synergized Data Efficiency and Compression (SEC) Optimization for Large Language Models. In *2024 4th International Conference on Electronic Information Engineering and Computer Science (EIECS)*.
- [14] Zongrong Li, Junhao Xu, Siqin Wang, Yifan Wu, and Haiyang Li. 2024. StreetViewLLM: Extracting Geographic Information Using a Chain-of-Thought Multimodal Large Language Model. *arXiv preprint arXiv:2411.14476* (2024).
- [15] Yunzhe Liu, Meixu Chen, Meihui Wang, Jing Huang, Fisher Thomas, Kazem Rahimi, and Mohammad Mamouei. 2023. An interpretable machine learning framework for measuring urban perceptions from panoramic street view images. *PLOS Computational Biology* 19, 3 (2023), e1010911.
- [16] Shangyu Lou, Gabriele Stancato, and Barbara EA Piga. 2024. Assessing in-motion urban visual perception: analyzing urban features, design qualities, and people's perception. In *Advances in Representation: New AI-and XR-Driven Transdisciplinarity*. Springer, 691–706.
- [17] Rohin Manvi, Samar Khanna, Gengchen Mai, Marshall Burke, David B. Lobell, and Stefano Ermon. 2023. GeoLLM: Extracting Geospatial Knowledge from Large Language Models. *arXiv preprint* (2023). arXiv:2310.06213
- [18] Python Library. [n. d.]. DiffLib Python Libaray. <https://docs.python.org/3/library/difflib.html>.
- [19] Gonzalo Travieso, Alexandre Benatti, and Luciano da F. Costa. 2024. An Analytical Approach to the Jaccard Similarity Index. *arXiv preprint arXiv:2410.16436* (2024).
- [20] Yiming Zeng, Wanhao Yu, Zexin Li, Tao Ren, Yu Ma, Jinghan Cao, Xiyan Chen, and Tingting Yu. 2025. Bridging the Editing Gap in LLMs: FineEdit for Precise and Targeted Text Modifications. arXiv:2502.13358
- [21] Fan Zhang, Bolei Zhou, Liu Liu, Yu Liu, Helene H. Fung, Hui Lin, and Carlo Ratti. 2018. Measuring human perceptions of a large-scale urban region using machine learning. *Landscape and Urban Planning* 180 (2018), 148–160.
- [22] Sanqiang Zhang, Xinyu Liu, et al. 2023. A Survey on Large Language Models: Applications, Challenges, and Opportunities. *arXiv preprint arXiv:2305.18703* (2023).