# CityRiSE: Reasoning Urban Socio-Economic Status in Vision-Language Models via Reinforcement Learning

**Tianhui Liu[1], Hetian Pang[2], Xin Zhang[2], Jie Feng[2], Yong Li[2], Pan Hui[1]**
[1]Information Hub, The Hong Kong University of Science and Technology (Guangzhou)
[2]Department of Electronic Engineering, BNRist, Tsinghua University

## Abstract

Harnessing publicly available, large-scale web data, such as street view and satellite imagery, urban socio-economic sensing is of paramount importance for achieving global sustainable development goals. With the emergence of Large Vision-Language Models (LVLMs), new opportunities have arisen to solve this task by treating it as a multi-modal perception and understanding problem. However, recent studies reveal that LVLMs still struggle with accurate and interpretable socio-economic predictions from visual data. To address these limitations and maximize the potential of LVLMs, we introduce **CityRiSE**, a novel framework for **R**eason**i**ng urban **S**ocio-**E**conomic status in LVLMs through pure reinforcement learning (RL). With carefully curated multi-modal data and verifiable reward design, our approach guides the LVLM to focus on semantically meaningful visual cues, enabling structured and goal-oriented reasoning for generalist socio-economic status prediction. Experiments demonstrate that CityRiSE with emergent reasoning process significantly outperforms existing baselines, improving both prediction accuracy and generalization across diverse urban contexts, particularly for prediction on unseen cities and unseen indicators. This work highlights the promise of combining RL and LVLMs for interpretable and generalist urban socio-economic sensing.

## 1 Introduction

Socio-economic indicators, such as GDP, population density, education level, and healthcare access, are critical for assessing the development, equity, and sustainability of urban regions. These metrics play a foundational role in guiding urban planning, policy-making, and international initiatives such as the United Nations Sustainable Development Goals (UN SDGs), which emphasize the importance of data-driven insights for fostering sustainable and inclusive growth. Traditionally, socio-economic data have been collected through large-scale surveys and censuses. While reliable, these methods are time-consuming, labor-intensive, and often lack the temporal granularity required for dynamic decision-making.

In recent years, the emergence of *web platforms* and open geospatial data sources has enabled new forms of urban perception. In particular, *web data* sources such as satellite imagery and street view images provide rich, frequently updated visual information that reflects the physical and infrastructural characteristics of urban environments. This has led to a growing body of research exploring the use of visual data to estimate socio-economic status, forming a new paradigm of visual perception for urban analysis.

Despite this progress, significant challenges remain. The urban perception involves complex multi-modal signals, requiring models to extract and integrate useful visual features from diverse views and reason about their relationship to abstract socio-economic targets. More importantly, for these models to be practically useful, they are expected to generalize not only across geographic regions but also to

Table 1: Comparison of Methods on Universality, City Transferability, Indicator Generalization, and Reasoning Ability.

| Method | Universal Model | City Transferability | Indicator Generalization | Reasoning Ability |
|---|---|---|---|---|
| UrbanKG [21] | ✗ | ✓ | ✗ | ✗ |
| SVF [6] | ✗ | ✓ | ✗ | ✗ |
| UrbanCLIP [38] | ✗ | ✓ | ✗ | ✗ |
| GeoLLM [24] | ✓ | ✓ | ✗ | ✗ |
| GeoSEE [11] | ✓ | ✓ | ✗ | ✗ |
| UrbanMLLM [41] | ✓ | ✗ | ✗ | ✗ |
| CityRiSE | ✓ | ✓ | ✓ | ✓ |

novel indicators. However, such generalization remains underexplored and technically demanding due to the heterogeneous nature of the input data and the abstraction level of the prediction targets.

A growing body of work has explored socio-economic indicator prediction from urban imagery, spanning a range of modeling paradigms. Early approaches, such as SVF [6], use a two-stage pipeline where visual features are extracted from street view images using a computer vision model and then fed into regression model for socio-economic prediction. More recent efforts explore more integrated approaches based on knowledge graphs and contrastive learning frameworks to better model complex urban patterns [45, 21, 2]. With the emergence of large vision-language models, these powerful models have also started to be applied to socio-economic indicator prediction tasks [38, 13, 41]. However, as summarized in Table 1, existing approaches face several limitations: First, prior to the advent of LVLMs, most models are designed for specific city. This city-specific training severely limits their scalability and generalization. Second, while many models achieve strong performance on indicators or cities seen during training, their transferability to unseen regions or targets remains limited. Some methods explore cross-city prediction, but with poor performance, and few address the more challenging task of cross-indicator generalization. Third, most existing models lack interpretability. They output numerical predictions without providing reasoning steps or explanations, making it difficult to analyze or trust the model's decision-making process.

To address these limitations, we propose *CityRiSE*, a reinforcement learning framework built on top of the large vision-language model for generalist urban socio-economic status prediction. CityRiSE achieves strong generalization across both geographic regions and socio-economic indicators, and notably is the first to demonstrate cross-indicator generalization in this domain. In addition, it enables interpretable, step-by-step reasoning without requiring costly manual annotations. In CityRiSE, we first adopt Group Relative Policy Optimization (GRPO) as our training strategy, enabling effective reward-based learning across groups of responses without relying on ground-truth labels. To support transferability to unseen cities, we further construct two auxiliary datasets that target complementary reasoning skills. The Perceptual Urban Reasoning Data focuses on intermediate tasks grounded in urban perception, such as city identification that mirror key structures of the socio-economic prediction problem. In contrast, the General Visual Reasoning Data, on the other hand, emphasizes logic-driven tasks like object counting, which cultivate abstract reasoning abilities. In parallel, we design a verifiable reward mechanism to support the emergence of interpretable reasoning. This includes a regression reward that enforces numerical correctness, and a keyword reward that implicitly guides the model to produce coherent, goal-directed reasoning chains. The main contributions of this paper are summarized as follows:

- We are the first, to our knowledge, to use reinforcement learning with LVLMs to enable emergent visual reasoning for high-level urban socio-economic status prediction.
- Through verifiable reward design, CityRiSE induces emergent, coherent, and interpretable reasoning chains for indicator prediction, offering transparency and explanatory insights beyond black-box outputs.
- By constructing auxiliary datasets that target transferable perceptual and reasoning skills, CityRiSE achieves strong generalization across both cities and indicators, with cross-indicator generalization demonstrated for the first time in this domain.
- Extensive experiments demonstrate that CityRiSE outperforms existing SOTA baselines in 11 urban socio-economic status prediction tasks, particularly for prediction in unseen cities and

unseen indicators. Our codes and datasets are open-sourced via `https://github.com/tsinghua-fib-lab/CityRiSE`.
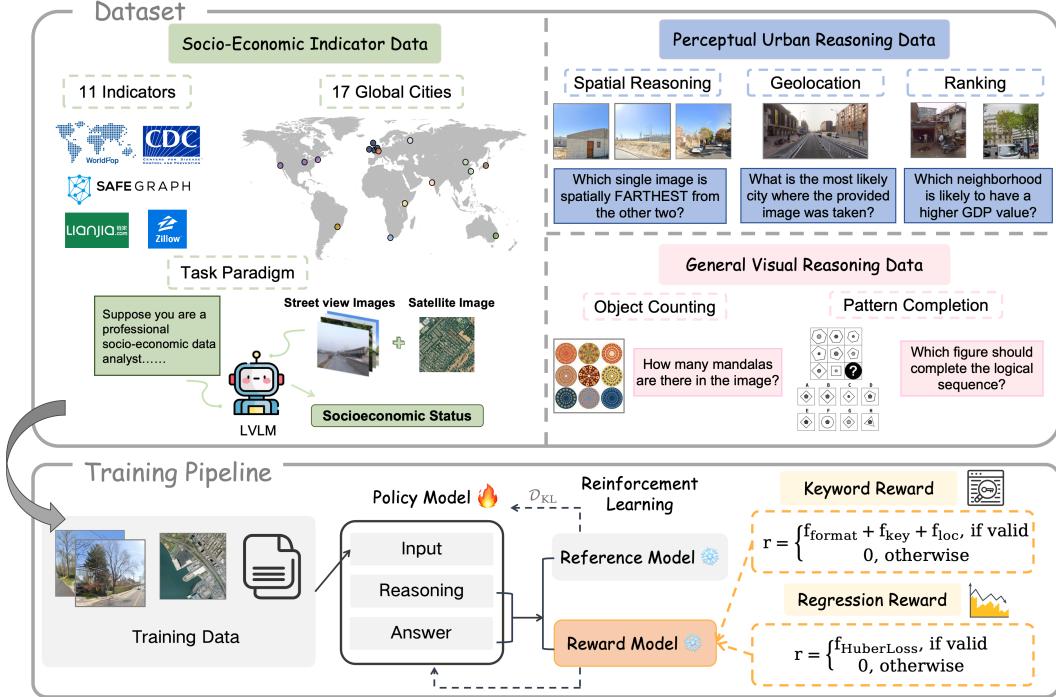
## 2 Methods

### 2.1 Transferable Datasets



Figure 1: Framework of CityRiSE, which integrates three types of training data with a GRPO training pipeline. The GRPO algorithm incorporates both Keyword Reward and Regression Reward to guide learning.

#### 2.1.1 Socio-Economic Indicator Data

We utilize CityLens [20] dataset as the source of our socio-economic indicators. CityLens is a multi-modal dataset covering 17 globally distributed cities, and provides a structured set of 11 socio-economic indicators that comprehensively reflect 6 major urban domains: economy, education, crime, transport, health, and environment. Specifically, the 11 indicators include: Gross Domestic Product (GDP), Population, House Price, Public Transport, Drive Ratio, Violent Crime, Mental Health, Accessibility to Health, Life Expectancy, Building Height, and Bachelor Ratio. Each spatial unit in the dataset is associated with its corresponding ground-truth indicator values, and is paired with one satellite image and ten street view images, enabling multi-modal learning and visual reasoning.

Following the method in GeoLLM [24], we discretize each socio-economic indicator into 10 bins and normalize the values to a range of 1 to 10. This binning formulation enables the model to handle heterogeneous indicator scales in a unified manner, facilitating stable optimization and allowing a consistent output space across different prediction targets. To evaluate city transferability, we split the 17 cities in CityLens into two groups. 10 cities are used for training—Beijing, New York, Cape Town, London, Mumbai, Moscow, Sydney, Paris, Tokyo, and Chicago—while the remaining 7 cities—Shanghai, San Francisco, Sao Paulo, Nairobi, Leeds, Liverpool, and Birmingham—are held out entirely and used for testing. This split enables the evaluation of cross-city transferability, where the model is required to predict indicators for previously unseen urban regions. In terms of prediction targets, we use 5 indicators—GDP, Population, Public Transport, Mental Health, and Bachelor Ratio—which are seen during both training and evaluation. The remaining 6 are reserved exclusively for testing, enabling the evaluation of the model's ability to generalize to unseen indicator

types. The task format in Socio-Economic Indicator Data-Train is consistent with that of the test set: each input consists of one satellite image and ten street view images from a given region, and the model is required to predict a specific socio-economic indicator associated with that region.

### 2.1.2 Perceptual Urban Reasoning Data

To enhance the model's ability to understand and reason about urban environments from street view images, we design three perceptual tasks that capture spatial, geographic, and socio-economic aspects of visual reasoning for training.

**Spatial Reasoning**   In this task, each input sample consists of three street view images, and the model is asked to determine which image is spatially furthest from the other two. This encourages the model to infer spatial proximity using only visual cues. We introduce two settings: (1) two images are sampled from the same city, while the third comes from a different city; (2) two images are taken from the same neighborhood or block, while the third comes from a different neighborhood within the same city. These settings require the model to develop a fine-grained perception of urban layouts, building density, architectural styles, and environmental continuity.

**Geolocation Task**   This task draws on a well-established tradition in urban geospatial research, where inferring a city's identity from visual data has long served as a benchmark for evaluating spatial understanding [16, 17]. Given a single street view image, the model is required to predict its city of origin. This setting not only tests the model's ability to extract high-level geographic and architectural signals, like signage, vegetation, or road infrastructure, but also provides valuable priors for socio-economic status prediction. By grounding visual inputs in specific geographic contexts, the model can leverage its internally stored world knowledge, such as the average socio-economic profile of the predicted city, to make more informed and context-aware predictions.

**Socio-Economic Ranking**   Inspired by prior work in image quality assessment that formulates visual comparison as a learning-to-rank problem [35, 40], we design a perceptual reasoning task in which the model compares two street view images and determines which region has a higher value for a given socio-economic indicator, relying solely on visual cues. This pairwise comparison not only helps the model learn to associate visual patterns with latent urban conditions, but also mirrors the underlying structure of the final prediction objective, which involves distinguishing subtle differences in socio-economic status across regions.

### 2.1.3 General Visual Reasoning Data

To support the model's ability in urban socio-economic sensing and abstraction, we incorporate two types of general visual reasoning tasks to build fundamental perceptual skills essential for urban understanding. The first is an object counting task [27], where the model must determine the number of instances of a given object type within an image. This requires the model to localize, distinguish, and quantify multiple objects under diverse spatial arrangements. The second is a pattern completion task, inspired by classic geometric reasoning problems [15]. In this setting, the model is presented with a visual sequence or matrix with a missing element and is required to infer the underlying rule to select the correct option. These tasks encourage the model to recognize abstract visual patterns, including regularity, symmetry, analogy, and transformation, which are transferable to urban reasoning challenges, such as identifying socio-spatial structures or comparing urban forms across cities.

Table 2: Reward design and instance statistics in the CityRiSE training data.

| Data Type | Format Reward | Accuracy Reward | Instances |
|---|---|---|---|
| Socio–Economic Indicator | Keyword | Regression | 2,828 |
| Spatial Reasoning | Standard | Standard | 632 |
| Geolocation Task | Standard | Standard | 350 |
| Socio-Economic Ranking | Standard | Standard | 699 |
| Object Counting | Standard | Regression | 300 |
| Pattern Completion | Standard | Standard | 300 |

## 2.2 Verifiable Reward

We train the LVLM using the GRPO algorithm, and design a verifiable reward mechanism to guide the optimization process. Table 2 summarizes the reward types used for different training data. In particular, we introduce two novel reward functions — Keyword Reward and Regression Reward — tailored specifically for the socio-economic indicator data. These rewards are designed to guide the LVLM toward more goal-directed perception and generation, encouraging outputs that align with specific socio-economic targets and can be objectively evaluated.

### 2.2.1 Keyword Reward

To guide the model toward producing responses that reflect meaningful urban visual features, we design a keyword-based reward component. This module is integrated into the format reward of GRPO, serving as an auxiliary signal that encourages attention to socio-economically relevant concepts during training. Concretely, we define a list of six urban-perceptual keywords: "person", "vehicle", "greenery", "road infrastructure", "street furniture", and "building", drawn from prior literature on urban visual intelligence [6]. These concepts represent key visual cues that correlate with a wide range of socio-economic indicators. In addition to these perceptual tags, we assign a special reward to the mention of "location", which indicates that the model attempts to anchor the visual input to a specific city. This grounding process enables the model to retrieve its internally stored prior knowledge about the city, such as its average socio-economic levels, and integrate it into its reasoning. This step is crucial for models that aim to generalize across diverse urban environments.

$$R_{\text{keyword}}(r) = \lambda_{\text{base}} \cdot \mathbf{1}_{\text{format}(r)} + \sum_{k \in \mathcal{K}} \lambda_k \cdot \mathbf{1}_{k \in r} + \lambda_{\text{loc}} \cdot \mathbf{1}_{\text{location} \in r}. \tag{1}$$

Here, $\mathbf{1}_{\text{condition}}$ is the indicator function, which equals 1 if the condition holds, and 0 otherwise.

### 2.2.2 Regression Reward

To better reflect the nature of our task as a regression problem, we design a reward function that is sensitive to the magnitude of prediction errors. In socio-economic indicator prediction, the model outputs a continuous score, and small deviations from the ground truth should not be penalized as harshly as large ones. For example, if the ground truth is 8, predicting 7 should be rewarded more than predicting 1—something that binary correctness-based metrics fail to capture.

To address this, we introduce a regression reward based on the Huber loss:

$$\text{error} = y_{\text{pred}} - y_{\text{true}}, \quad |\text{error}| = |y_{\text{pred}} - y_{\text{true}}|, \tag{2}$$

$$L_{\text{Huber}}(\text{error}) = \begin{cases} \frac{1}{2} \text{error}^2, & \text{if } |\text{error}| \leq \delta \\ \delta \left( |\text{error}| - \frac{1}{2}\delta \right), & \text{otherwise} \end{cases}, \tag{3}$$

$$R_{\text{regression}} = \exp\left( -\alpha \, L_{\text{Huber}}(\text{error}) \right). \tag{4}$$

This formulation yields a reward that decays smoothly and non-linearly with prediction error. The use of Huber loss ensures stability near the correct value due to its quadratic form, while reducing sensitivity to outliers through its linear behavior for larger errors. The exponential transformation ensures that predictions closer to the true value receive rewards close to 1, while distant predictions are penalized more strongly. Overall, this reward better aligns with the continuous, fine-grained nature of the prediction task, and provides a more informative learning signal during training. A similar formulation is also applied to object counting data, where the model produces continuous-valued outputs and is likewise trained using the regression reward.

## 2.3 Training Pipeline

We fine-tune the LVLM using GRPO on the dataset described in Section 2.1. As Algorithm 1 shows, GRPO operates by generating multiple candidate responses for each data point and calculating the reward for each candidate based on the reward design outlined in Section 2.2. To encourage the model to favor higher-reward completions within each group, GRPO computes a group-normalized advantage for each response. This advantage reflects how well each response performs relative to

other candidates in the same group. Finally, the model is updated by minimizing a loss function that incorporates both the advantage and the KL divergence, ensuring the model learns to improve its performance on the dataset.

---

**Algorithm 1** Training Algorithm for CityRiSE using GRPO (Single Data Point)

---

**Input:** Pretrained LVLM policy $\pi_\theta$, reference policy $\pi_{\text{ref}}$, reward model $R$, training data $\mathcal{D} = \{(x_i, y_i)\}$, KL coefficient $\beta$, clipping coefficient $\epsilon$.

**Output:** Fine-tuned policy $\pi_\theta$.

Initialize model parameters $\theta$

**for** each $(x_i, y_i) \in \mathcal{D}$ **do**

 Generate N responses: $\{o_i^{(j)}\}_{j=1}^N \sim \pi_\theta(\cdot \mid x_i)$

 Compute reward: $r_i^{(j)} = R(o_i^{(j)})$

 Compute advantage: $A_i^{(j)} = r_i^{(j)} - \frac{1}{N} \sum_{l=1}^N r_i^{(l)}$

 Compute ratio: $s_i^{(j)} = \frac{\pi_\theta(o_i^{(j)}|x_i)}{\pi_{\theta_{\text{old}}}(o_i^{(j)}|x_i)}$

 Compute loss for the data point:

$$\mathcal{L}_i = \frac{1}{N} \sum_{l=1}^N \left[ \min\left( s_i^{(l)} \cdot A_i^{(l)}, \text{clip}(s_i^{(l)}, 1 - \epsilon, 1 + \epsilon) \cdot A_i^{(l)} \right) - \beta D_{\text{KL}}[\pi_\theta \parallel \pi_{\text{ref}}] \right]$$

 Update $\theta$ via gradient descent on $\mathcal{L}_i$

**end for**

**return** trained model $\pi_\theta$

---

# 3 Experiments

## 3.1 Settings

### 3.1.1 Baselines

To comprehensively evaluate our approach, we compare it against 8 baselines falling into 3 categories: (1) classic methods that do not involve LVLMs, (2) general large vision-language models, and (3) LVLMs trained on domain-specific data.

**Classic methods**  **Random**: For each input, a value is randomly sampled from a uniform distribution over 1–10. Despite its simplicity, this baseline serves as a critical reference point to understand performance gains. **SVF-GT** and **SVF-Rank**: These methods are adapted from [6]. They first extract 13 specific visual features from each street view image using `ResNet18dilated + PPM_deepsup` [43]. Then, a Least Absolute Shrinkage and Selection Operator (LASSO) regression model is trained to predict the socio-economic indicators. The GT variant uses the ground-truth indicators as targets, while the Rank variant aligns with the normalization setup used by all other models.

**General-LVLMs**  **Qwen2.5-VL-7B**: This model is a large vision-language model with strong visual understanding capabilities. It serves as our foundational LVLM, upon which both supervised fine-tuning (SFT) and reinforcement learning are later applied. **Qwen3-VL-235B A22B Thinking**, **GPT-4.1-Mini**, and **Gemini 2.5-Flash**: These are state-of-the-art general-purpose LVLMs that exhibit strong capabilities across a wide range of tasks. We include them to benchmark how CityRiSE compares with powerful generalist models.

**SFT-LVLMs**  **UrbanMLLM**: This model is fine-tuned from Qwen2.5-VL-7B using four epochs of supervised learning on the Socio-Economic Indicator Data-Train dataset, following the paradigm proposed in [41]. **UI-CoT**: Inspired by [42], we construct a Chain-of-Thought (CoT) prompting dataset from Socio-Economic Indicator Data-Train (example provided in Appendix 5), and perform four epochs of SFT on Qwen2.5-VL-7B. This baseline evaluates the impact of incorporating structured reasoning during supervised fine-tuning.

### 3.1.2 Metrics

To evaluate model performance, we adopt the coefficient of determination ($R^2$), which is one of the most widely used metrics in socio-economic prediction. The $R^2$ score measures the proportion of variance in the ground-truth values that is captured by the model's predictions. An $R^2$ score of 1 indicates a perfect fit, while lower values represent worse performance. Notably, a negative $R^2$ indicates that the model performs worse than simply predicting the mean of the ground-truth values for all instances.

## 3.2 Overall Performance

Table 3: Main results on Socio-Economic Indicator Data-Test. The values in the table represent $R^2$ scores. In each row, bold indicates the best result, and underline denotes the second-best.

| Domain | In-domain | | | | | Unseen Cities | | | | | Unseen Indicators | | | | | | Overall |
| Tasks | GDP | Pop. | MH | PT | BR | GDP | Pop. | MH | PT | BR | DR | BH | VC | AH | LE | HP | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Random | -1.293 | -1.220 | -0.798 | -0.562 | -1.191 | -4.426 | -0.874 | -2.902 | -11.543 | -2.085 | -0.916 | -0.734 | -1.106 | -0.902 | -0.941 | -0.951 | -0.909 |
| SVF-GT | 0.329 | 0.268 | 0.087 | 0.417 | 0.210 | -1.847 | 0.091 | -1.964 | -4.404 | -0.517 | / | / | / | / | / | / | / |
| SVF-Rank | 0.366 | 0.267 | _0.162_ | 0.358 | 0.287 | -0.830 | 0.172 | -2.154 | -5.351 | -0.720 | / | / | / | / | / | / | / |
| Qwen2.5-VL-7B | -0.280 | -0.330 | -0.545 | -0.013 | -0.094 | -1.772 | -0.209 | -2.172 | -1.459 | -0.918 | -0.130 | 0.178 | -0.251 | _-0.311_ | _0.083_ | -0.026 | -0.096 |
| Qwen3-VL-235B | -0.752 | -0.944 | -1.288 | -0.621 | 0.173 | -1.573 | -0.405 | -5.924 | **-0.165** | -0.169 | -1.534 | 0.290 | -1.150 | -1.069 | -0.874 | 0.056 | -0.517 |
| GPT-4.1-Mini | -0.323 | -0.222 | -0.989 | 0.342 | 0.239 | -0.798 | -0.183 | -4.943 | _-0.354_ | -1.386 | _0.252_ | 0.285 | -0.908 | -0.897 | -0.023 | _0.111_ | -0.216 |
| Gemini2.5-Flash | -0.825 | -0.918 | -1.078 | 0.529 | -0.220 | -2.676 | -0.529 | -6.709 | -0.984 | -1.801 | -0.903 | _0.395_ | -1.435 | -1.717 | -0.525 | _-0.211_ | -0.674 |
| UrbanMLLM (SFT) | **0.691** | **0.574** | -0.005 | **0.607** | 0.534 | **0.045** | _0.274_ | **0.329** | -0.470 | -0.946 | -2.249 | 0.314 | -0.087 | -2.008 | -1.400 | -0.209 | -0.057 |
| UI-CoT (SFT) | _0.546_ | _0.499_ | 0.053 | 0.446 | **0.616** | -0.952 | -0.055 | 0.264 | -0.951 | _0.063_ | -0.939 | **0.398** | _0.103_ | -1.780 | -0.493 | -0.105 | _0.099_ |
| CityRiSE (only RL) | 0.385 | 0.334 | **0.399** | _0.594_ | _0.603_ | _-0.255_ | **0.365** | 0.259 | -4.387 | **0.286** | **0.289** | 0.366 | **0.305** | **-0.196** | **0.200** | **0.218** | **0.361** |

Table 3 presents the main results on the Socio-Economic Indicator Data-Test set. The evaluation tasks are broadly categorized into three types based on their relationship to the training distribution:

- In-domain. Both the cities and indicators in these tasks are consistent with those seen during training. Importantly, the test set is constructed with entirely disjoint samples to prevent any data leakage and ensure a fair evaluation.
- Unseen Cities. The indicators remain the same as in training, but the evaluation is conducted on cities that were not present in the training set.
- Unseen Indicators. These tasks involve indicators that do not appear in the training set. Except for the Life Expectancy task, which involves both unseen cities and unseen indicators, all other tasks in this group share cities with the training set but introduce new indicators.

This categorization allows us to assess the model's generalization ability across different dimensions: geographical (new cities) and semantic (new indicators).

### 3.2.1 In-domain Analysis

Under the in-domain setting, CityRiSE significantly outperforms its base model Qwen2.5-VL-7B across all five indicators, with $R^2$ scores improving from $-0.330$ to $0.334$ on Population and from $-0.013$ to $0.594$ on Public Transport, ultimately reaching up to $0.603$ on Bachelor Ratio. These results demonstrate the strong predictive capability of CityRiSE. Qwen3-VL-235B A22B Thinking, GPT-4.1-Mini, and Gemini2.5-Flash are all strong general large vision-language models, while their zero-shot performance on socio-economic status prediction remains suboptimal, further highlighting the challenge of this task. Interestingly, SFT-based models, including UrbanMLLM and UI-CoT, achieve the highest performance on some indicators, particularly GDP and Population. This aligns with the insight proposed in [4] that "SFT Memorizes, RL Generalizes". Supervised fine-tuning is highly effective for in-domain memorization, whereas reinforcement learning tends to improve generalization beyond training data. We also observe that SVF-GT and SVF-Rank achieve very similar performance across all in-domain indicators. This suggests that normalizing the ground-truth values in the SVF-Rank setup does not substantially degrade predictive accuracy compared to directly regressing the raw values. While this does not prove full equivalence between the two formats, it provides empirical support for the practicality of our normalizing strategy. Importantly, CityRiSE consistently outperforms both SVF baselines, highlighting the advantage of leveraging large vision-language models and reinforcement learning in this task.

### 3.2.2 Transferability to Unseen Cities

In the unseen cities setting, CityRiSE demonstrates strong spatial transferability across all indicators. Notably, it achieves the highest performance on Bachelor Ratio with an $R^2$ of 0.286, outperforming SFT-based models, general-LVLMs and classic methods. Despite being trained on a different set of cities, the model effectively transfers learned representations to novel urban contexts. While CityRiSE and UrbanMLLM show comparable overall performance, CityRiSE holds a distinct advantage that it autonomously generates step-by-step reasoning during prediction. This allows the model to provide interpretable justifications for its socio-economic estimates, offering greater transparency and accountability. Such structured reasoning is not naturally produced by UrbanMLLM, underscoring the benefit of reinforcement learning in producing both accurate and explainable predictions.

### 3.2.3 Generalization to Unseen Indicators

In the unseen indicators setting, CityRiSE achieves the best overall performance, outperforming all baselines across most of the novel socio-economic indicators. This highlights its strong ability to generalize not only to new visual contexts, but also to entirely unseen prediction targets. The model maintains positive $R^2$ scores on all unseen indicators except Accessibility to Health, including challenging ones such as Life Expectancy (0.200) and Housing Price (0.218), where other models struggle or fail to generalize. The relatively low performance on AH is understandable, as it represents the 'average walking-only travel time to healthcare facilities', an indicator that is inherently difficult to estimate from limited visual input alone, due to its dependence on infrastructural and service distribution information that may not be directly observable. Notably, traditional regression-based baselines such as SVF-GT and SVF-Rank are absent from this evaluation. This is because such models require a separate predictor for each target variable, making them inherently limited in handling unseen indicators. This limitation underscores a key advantage of large vision-language models: their ability to serve as universal predictors for diverse tasks through language-based prompts.

CityRiSE exemplifies this capability—as a universal model, it can reason about a wide range of socio-economic indicators across cities, even without explicit training on each target. Through guided optimization on the training set, CityRiSE learns a transferable perception-to-prediction paradigm that it discovers how to interpret satellite and street view imagery to estimate socio-economic conditions. This learned paradigm enables the model to adapt to novel indicators and urban contexts, highlighting its potential as a general framework for socio-economic status prediction tasks.

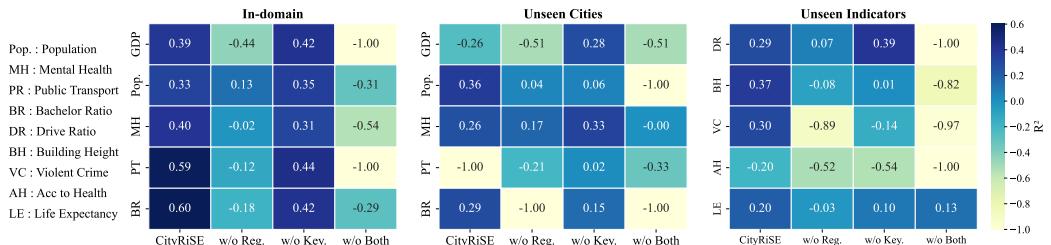### 3.3 Ablation Study

### 3.3.1 Reward Ablation



Figure 2: Results of Reward Ablation, where 'Reg.' refers to the Regression Reward, 'Key.' refers to the Keyword Reward, and 'w/o Both' indicates the setting where both rewards are removed.

Figure 2 presents the results of a reward ablation study, comparing the performance of the full CityRiSE model with three variants: one trained without the Regression Reward, one without the Keyword Reward, and one with both rewards removed. When either reward is removed, the model is optimized using only the standard reward. To enhance visualization clarity, we clip all $R^2$ values less than $-1$, which are considered non-informative for evaluation, to $-1$ in the figure. The complete, unprocessed results are provided in Appendix A.4.

We observe that removing either reward consistently degrades performance, and removing both leads to severe performance collapse across most indicators. For instance, in the in-domain setting,

removing both rewards causes the model's $R^2$ score on GDP to drop from 0.39 to below $-1.00$, and on Public Transport from 0.59 to below $-1.00$, indicating that the model fails to receive meaningful learning signals during training and is thus unable to learn how to predict socio-economic indicators effectively. Interestingly, the Keyword Reward appears particularly critical in certain indicators. For example, removing only the Keyword Reward leads to a sharp decline on Violent Crime from 0.30 to $-0.14$ and Life Expectancy from 0.20 to 0.10, suggesting its role in enhancing semantic grounding. And the Regression Reward contributes strongly to numeric precision, as evidenced by the drop in Population on unseen cities from 0.36 to $-0.04$ when it is removed. Overall, the impact of the Regression Reward is more significant than that of the Keyword Reward, which we attribute to the fact that relevant indicator keywords are already explicitly present in the prompts. This explicit guidance may reduce the relative contribution of the Keyword Reward during training. These results confirm that both rewards contribute complementary supervision signals. The Regression Reward aligns outputs with quantitative targets, while the Keyword Reward strengthens interpretability and semantic relevance. Together, they enable CityRiSE to make more accurate and generalizable predictions across diverse socio-economic indicators.

### 3.3.2 Data Ablation

We incorporate the Perceptual Urban Reasoning Data and General Visual Reasoning Data with the goal of enhancing the model's reasoning capabilities that are relevant to socio-economic status prediction. In particular, we aim to improve the model's generalization when transferring to out-of-domain settings. To evaluate the impact of these datasets, Figure 3 presents the performance of models trained with different data ablations on the unseen cities and unseen indicators subsets. To enable clearer comparison across models, we exclude tasks where all model configurations, including the full model and its ablated variants, achieve negative $R^2$ scores, as they offer limited value for analysis. The complete results, including all tasks, are provided in Appendix A.5.
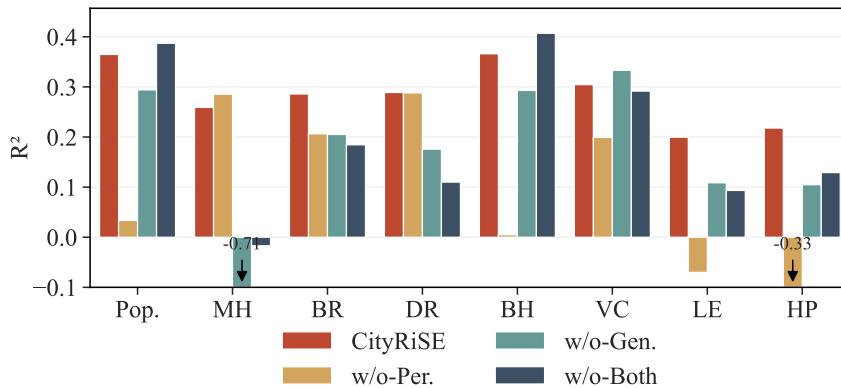


Figure 3: Results of Data Ablation, where 'Per.' refers to the Perceptual Urban Reasoning Data, 'Gen.' refers to the General Visual Reasoning Data.

As shown in Figure 3, removing either of the auxiliary datasets consistently leads to performance drops across most indicators, confirming the effectiveness of both Perceptual and General Visual Reasoning data. The General Visual Reasoning Data exhibits a broad impact that its removal leads to noticeable declines on nearly all indicators, most notably on Mental Health, where the $R^2$ score drops from 0.26 to $-0.71$. In contrast, the Perceptual Urban Reasoning Data plays a more targeted yet critical role. For example, its removal substantially affects Life Expectancy and House Price, suggesting that learning from spatial comparisons and urban identity reasoning contributes to accurate estimation for indicators grounded in regional characteristics. When both datasets are removed, the model suffers from severe degradation across almost all tasks, including a complete collapse on Mental Health, demonstrating the complementary nature of the two datasets in supporting robust visual reasoning and generalization.
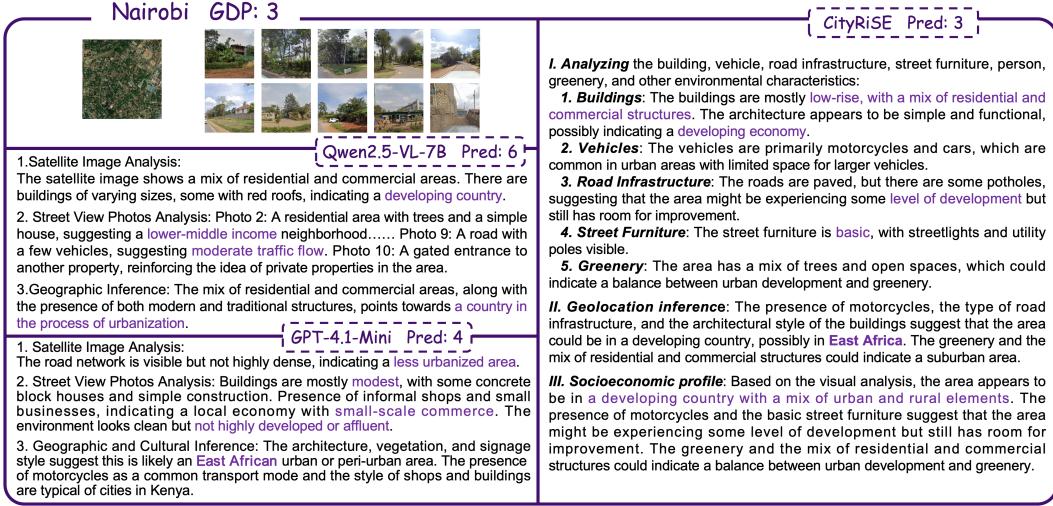
Figure 4: Qualitative comparison of reasoning patterns produced by different models. CityRiSE exhibits more interpretable and goal-directed reasoning behavior, while baseline models often generate less informative responses.

## 3.4 Reasoning Comparison

### 3.4.1 Reasoning Case

Figure 4 presents a qualitative comparison of partial reasoning chains generated by three different models when responding to the same socio-economic prediction task. Notably, although CityRiSE is trained via reinforcement learning without access to answer labels, it still develops a clear and consistent reasoning structure. This behavior stems from the combined effect of the Keyword Reward and prompt design, the former offering implicit guidance, and the latter providing explicit instruction, both steering the model toward goal-directed reasoning. Compared to its base model (Qwen2.5-VL-7B), CityRiSE demonstrates a more holistic and higher-level perception across multiple images, rather than producing low-information summaries of each image in isolation. It effectively integrates visual cues to form richer inferences. In particular, CityRiSE is able to more accurately perceive and reason about the street view images. It successfully infers the likely location as "East Africa", aligning with the output of GPT-4.1-Mini, a much larger proprietary model. This geographic inference provides valuable prior knowledge for socio-economic status prediction, as recognizing the region as East Africa naturally helps the model adjust its expectations for GDP, which tends to fall within a lower range. This demonstrates not only robust visual grounding, but also a natural and uninterrupted reasoning process, emerging organically through RL training.

### 3.4.2 Reasoning Induction

Table 4 presents a comparison of models trained under different strategies on three representative tasks spanning in-domain, unseen cities, and unseen indicators. Full results are available in Appendix A.6. The three models differ in how reasoning is introduced during training. UI-CoT refers to the supervised fine-tuning model trained on manually constructed CoT annotations. UI-CoT + RL further fine-tunes the UI-CoT model using reinforcement learning. In contrast, CityRiSE is trained purely via RL, without access to human-crafted reasoning chains, relying instead on reward design to induce emergent reasoning behavior. Interestingly, the UI-CoT model tends to produce slightly rigid and formulaic reasoning structures that perform well on in-domain tasks. However, this fixed reasoning template fails to generalize to tasks outside the training distribution. As shown in the table, its performance drops significantly on unseen tasks. Consistently, applying RL to the UI-CoT model leads to further gains in in-domain tasks, but has limited impact on out-of-domain performance. We attribute this to the rigidity of its reasoning structure, which is inherited from its supervised training and remains largely unchanged during RL fine-tuning. In contrast, CityRiSE develops reasoning abilities organically during RL training, guided by the verifiable reward design. This leads to flexible, transferable inference behaviors that extend beyond training tasks, with strong performance on unseen

cities and unseen indicators. These results support the hypothesis that emergent reasoning, induced rather than explicitly provided, is more adaptable to novel contexts. We further experiment with applying RL on top of UrbanMLLM, a SFT-only model. Even when the prompt explicitly instructs the model to "Please reason step-by-step and write your reasoning inside the <think> </think> tags" during RL, the model consistently ignores the instruction and produce only the final answer. This highlights that once reasoning behavior is not introduced during early supervised fine-tuning, it is difficult to induce it later through RL alone.

Table 4: Performance of models trained with different methods on three types of tasks.

| Model | In-domain | Unseen Cities | Unseen Indicators |
|---|---|---|---|
| | Public Transport | Bachelor Ratio | Violent Crime |
| UI-CoT (SFT-CoT) | 0.446 | 0.063 | 0.103 |
| UI-CoT + RL | 0.658 | 0.013 | 0.107 |
| CityRiSE (RL) | 0.594 | 0.286 | 0.305 |

## 4 Related Work

**Urban Socio-Economic Status Prediction**   It often referred to as urban sensing [14], involves assessing and forecasting various social and economic statuses of cities using multi-modal urban information like street view and satellite data, covering predictions for economic levels [25, 19], environmental status [32, 22, 31], and health levels [26, 12]. Researchers have developed various modeling approaches that leverage the diverse, multi-modal urban data, ranging from classic statistical methods to modern deep learning and representation learning techniques. The primary methods used for this task include: Classical Regression Methods [6] that extract manually defined features from visual data using classic computer vision models [43] and apply standard regression analysis; Knowledge Graph Methods [21, 45] which structure various urban data into a well-defined urban knowledge graph to facilitate automated feature extraction before inputting them into deep neural networks for prediction; and Representation Learning Approaches [2, 39] which often through modern contrastive learning, automatically learn robust, universal representations of various urban data by capturing relationships across different data modalities, subsequently using these representations to obtain better prediction results.

**LLM for Socio-Economic Status Prediction**   With the advent of LLMs, a promising new direction involves leveraging their embedded commonsense and world knowledge to enrich urban visual data with contextual information such as background and culture [18, 8, 20]. The typical paradigms for this integration fall into two categories: using LLMs as a supplementary tool for extra information to assist the aforementioned classical methods [38, 13], and directly employing LLMs as the predictor by leveraging their rich internal knowledge and modeling capabilities. In the first paradigm, a representative example is UrbanCLIP, which uses an LLM to annotate existing data, obtaining new auxiliary features that enhance the effectiveness of subsequent contrastive learning. The second approach, which more tightly integrates the LLM's inherent knowledge with indicator prediction, has garnered more attention  [24, 23, 41, 7]. Furthermore, GeoSEE [11] utilizes the LLM's In-Context Learning capability to achieve rapid adaptation and accurate prediction. In addition, Zhou et al. [44] use LLMs as an Agent to automatically construct and parse urban knowledge graphs for socio-economic indicator prediction, Deng et al. [5] endow multi-modal LLMs (MLLMs) with the ability to analyze massive street view data for predicting socio-economic status changes through a cleverly designed mechanism.

**Generalist Reasoning in Large Vision-Language Models**   Prior to leveraging Reinforcement Learning to encourage autonomous visual reasoning, existing methods primarily guided MLLMs through fixed prompt design for visual inference [10, 3, 1, 34]. More recently, autonomous visual reasoning in MLLMs based purely on Reinforcement Learning has become a popular topic. The core goal is to train MLLMs via RL, to independently learn to analyze and reason, thereby boosting performance and crucially generalization ability in typical visual tasks [33, 30]. Current work in this space has mainly concentrated on the understanding and reasoning of low-level, intuitive visual elements, such as mathematical geometry problems and natural image comprehension. By utilizing

datasets annotated with explicit reasoning processes [37, 29] and carefully designed reward functions [28, 36], similar to the approach of DeepSeek-R1 [9], researchers guide MLLMs toward more general visual reasoning, leading to higher performance and better generalization. Distinct from these efforts focusing on intuitive, low-level visual elements, our work is the first to apply the RL paradigm to implicit, high-level visual reasoning. We achieve the association and mining of intuitive elements with hidden knowledge, thus enabling the generalized prediction of unseen socio-economic indicators for the first time while simultaneously enhancing the generalized prediction capability for unseen geographic regions.

## 5 Conclusion

In this paper, we propose CityRiSE, the first framework that uses reinforcement learning to elicit the emergent reasoning capabilities of large vision-language models, thereby achieving generalist urban socio-economic status prediction across diverse urban contexts. Within CityRiSE, we lay the groundwork for generalist prediction by utilizing carefully curated multi-source, multi-modal data with readily obtainable labels. Crucially, a verified multi-component reward design is then employed, activating the LVLMs' emergent reasoning and subsequently strengthening their generalized prediction ability while also providing interpretability. Extensive experimental results demonstrate the superiority of our proposed framework over state-of-the-art baselines for 11 urban socio-economic status predictions. Furthermore, CityRiSE provides clear, realizable reasoning steps, a distinct advantage when compared against powerful closed-source models like GPT-4.1-Mini and Gemini2.5-Flash.

## References

[1] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465, 2024.

[2] Meng Chen, Zechen Li, Weiming Huang, Yongshun Gong, and Yilong Yin. Profiling urban streets: A semi-supervised prediction model based on street view imagery and spatial topology. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 319–328, 2024.

[3] An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang, and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision-language models. *Advances in Neural Information Processing Systems*, 37:135062–135093, 2024.

[4] Tianzhe Chu, Yuexiang Zhai, Jihan Yang, Shengbang Tong, Saining Xie, Dale Schuurmans, Quoc V Le, Sergey Levine, and Yi Ma. Sft memorizes, rl generalizes: A comparative study of foundation model post-training. *arXiv preprint arXiv:2501.17161*, 2025.

[5] Boyang Deng, Songyou Peng, Kyle Genova, Gordon Wetzstein, Noah Snavely, Leonidas Guibas, and Thomas Funkhouser. Visual chronicles: Using multimodal llms to analyze massive collections of images. *arXiv preprint arXiv:2504.08727*, 2025.

[6] Zhuangyuan Fan, Fan Zhang, Becky P. Y. Loo, and Carlo Ratti. Urban visual intelligence: Uncovering hidden city profiles with street view images. *Proceedings of the National Academy of Sciences*, 120(27):e2220417120, 2023.

[7] Jie Feng, Shengyuan Wang, Tianhui Liu, Yanxin Xi, and Yong Li. Urbanllava: A multi-modal large language model for urban intelligence. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6209–6219, 2025.

[8] Jie Feng, Jun Zhang, Tianhui Liu, Xin Zhang, Tianjian Ouyang, Junbo Yan, Yuwei Du, Siqi Guo, and Yong Li. Citybench: Evaluating the capabilities of large language models for urban tasks. In *Proceedings of the 31st ACM SIGKDD Conference on Knowledge Discovery and Data Mining V. 2*, pages 5413–5424, 2025.

[9] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Peiyi Wang, Qihao Zhu, Runxin Xu, Ruoyu Zhang, Shirong Ma, Xiao Bi, et al. Deepseek-r1 incentivizes reasoning in llms through reinforcement learning. *Nature*, 645(8081):633–638, 2025.

[10] Qiushan Guo, Shalini De Mello, Hongxu Yin, Wonmin Byeon, Ka Chun Cheung, Yizhou Yu, Ping Luo, and Sifei Liu. Regiongpt: Towards region understanding vision language model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13796–13806, 2024.

[11] Sungwon Han, Donghyun Ahn, Seungeon Lee, Minhyuk Song, Sungwon Park, Sangyoon Park, Jihee Kim, and Meeyoung Cha. Geosee: Regional socio-economic estimation with a large language model. *arXiv preprint arXiv:2406.09799*, 2024.

[12] Zhenyu Han, Tong Xia, Yanxin Xi, and Yong Li. Healthy cities, a comprehensive dataset for environmental determinants of health in england cities. *Scientific data*, 10(1):165, 2023.

[13] Xixuan Hao, Wei Chen, Yibo Yan, Siru Zhong, Kun Wang, Qingsong Wen, and Yuxuan Liang. Urbanvlp: Multi-granularity vision-language pretraining for urban socioeconomic indicator prediction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 28061–28069, 2025.

[14] Yuhao Kang, Fan Zhang, Song Gao, Hui Lin, and Yu Liu. A review of urban physical environment sensing using street view imagery in public health studies. *Annals of GIS*, 26(3):261–275, 2020.

[15] Sicong Leng, Jing Wang, Jiaxi Li, Hao Zhang, Zhiqiang Hu, Boqiang Zhang, Yuming Jiang, Hang Zhang, Xin Li, Lidong Bing, Deli Zhao, Wei Lu, Yu Rong, Aixin Sun, and Shijian Lu. Mmr1: Enhancing multimodal reasoning with variance-aware sampling and open resources, 2025.

[16] Ling Li, Yu Ye, Bingchuan Jiang, and Wei Zeng. Georeasoner: Geo-localization with reasoning in street views using a large vision-language model. In *International Conference on Machine Learning (ICML)*, 2024.

[17] Ling Li, Yao Zhou, Yuxuan Liang, Fugee Tsung, and Jiaheng Wei. Recognition through reasoning: Reinforcing image geo-localization with large vision-language models. *arXiv preprint arXiv:2506.14674*, 2025.

[18] Zhuoheng Li, Yaochen Wang, Zhixue Song, Yuqi Huang, Rui Bao, Guanjie Zheng, and Zhenhui Jessie Li. What can llm tell us about cities? *arXiv preprint arXiv:2411.16791*, 2024.

[19] Yuming Lin, Xin Zhang, Yu Liu, Zhenyu Han, Qingmin Liao, and Yong Li. Long-term detection and monitory of chinese urban village using satellite imagery. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 7349–7357, 2024.

[20] Tianhui Liu, Jie Feng, Hetian Pang, Xin Zhang, Tianjian Ouyang, Zhiyuan Zhang, and Yong Li. Citylens: Benchmarking large language-vision models for urban socioeconomic sensing. *arXiv preprint arXiv:2506.00530*, 2025.

[21] Yu Liu, Xin Zhang, Jingtao Ding, Yanxin Xi, and Yong Li. Knowledge-infused contrastive learning for urban imagery-based socioeconomic prediction. In *Proceedings of the ACM web conference 2023*, pages 4150–4160, 2023.

[22] Ying Long and Liu Liu. How green are the streets? an analysis for central areas of chinese cities using tencent street view. *PloS one*, 12(2):e0171110, 2017.

[23] Rohin Manvi, Samar Khanna, Marshall Burke, David Lobell, and Stefano Ermon. Large language models are geographically biased. *arXiv preprint arXiv:2402.02680*, 2024.

[24] Rohin Manvi, Samar Khanna, Gengchen Mai, Marshall Burke, David B. Lobell, and Stefano Ermon. Geollm: Extracting geospatial knowledge from large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

[25] Nikhil Naik, Scott Duke Kominers, Ramesh Raskar, Edward L Glaeser, and César A Hidalgo. Computer vision uncovers predictors of physical urban change. *Proceedings of the National Academy of Sciences*, 114(29):7571–7576, 2017.

[26] Tianjian Ouyang, Xin Zhang, Zhenyu Han, Yu Shang, and Yong Li. Health clip: Depression rate prediction using health related features in satellite and street view images. In *Companion Proceedings of the ACM Web Conference 2024*, pages 1142–1145, 2024.

[27] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten, 2023.

[28] Haozhan Shen, Peng Liu, Jingcheng Li, Chunxin Fang, Yibo Ma, Jiajia Liao, Qiaoli Shen, Zilun Zhang, Kangjia Zhao, Qianqian Zhang, et al. Vlm-r1: A stable and generalizable r1-style large vision-language model. *arXiv preprint arXiv:2504.07615*, 2025.

[29] Yifan Shen, Yuanzhe Liu, Jingyuan Zhu, Xu Cao, Xiaofeng Zhang, Yixiao He, Wenming Ye, James Matthew Rehg, and Ismini Lourentzou. Fine-grained preference optimization improves spatial reasoning in vlms. *arXiv preprint arXiv:2506.21656*, 2025.

[30] Zhaochen Su, Peng Xia, Hangyu Guo, Zhenhua Liu, Yan Ma, Xiaoye Qu, Jiaqi Liu, Yanshu Li, Kaide Zeng, Zhengyuan Yang, et al. Thinking with images for multimodal reasoning: Foundations, methods, and future frontiers. *arXiv preprint arXiv:2506.23918*, 2025.

[31] Esra Suel, John W Polak, James E Bennett, and Majid Ezzati. Measuring social, environmental and health inequalities using deep learning and street imagery. *Scientific reports*, 9(1):6229, 2019.

[32] Jingxian Tang and Ying Long. Measuring visual quality of street space and its temporal variation: Methodology and its application in the hutong area in beijing. *Landscape and Urban Planning*, 191:103436, 2019.

[33] Yaoting Wang, Shengqiong Wu, Yuecheng Zhang, Shuicheng Yan, Ziwei Liu, Jiebo Luo, and Hao Fei. Multimodal chain-of-thought reasoning: A comprehensive survey. *arXiv preprint arXiv:2503.12605*, 2025.

[34] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13084–13094, 2024.

[35] Tianhe Wu, Jian Zou, Jie Liang, Lei Zhang, and Kede Ma. VisualQuality-R1: Reasoning-induced image quality assessment via reinforcement learning to rank. *arXiv preprint arXiv:2505.14460*, 2025.

[36] Tong Xiao, Xin Xu, Zhenya Huang, Hongyu Gao, Quan Liu, Qi Liu, and Enhong Chen. Advancing multimodal reasoning capabilities of multimodal large language models via visual perception reward. *arXiv preprint arXiv:2506.07218*, 2025.

[37] Guowei Xu, Peng Jin, Ziang Wu, Hao Li, Yibing Song, Lichao Sun, and Li Yuan. Llava-cot: Let vision language models reason step-by-step. *arXiv preprint arXiv:2411.10440*, 2024.

[38] Yibo Yan, Haomin Wen, Siru Zhong, Wei Chen, Haodong Chen, Qingsong Wen, Roger Zimmermann, and Yuxuan Liang. Urbanclip: Learning text-enhanced urban region profiling with contrastive language-image pretraining from the web. In *Proceedings of the ACM Web Conference 2024*, pages 4006–4017, 2024.

[39] Xixian Yong and Xiao Zhou. Musecl: predicting urban socioeconomic indicators via multi-semantic contrastive learning. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence*, pages 7536–7544, 2024.

[40] Weixia Zhang, Kede Ma, Guangtao Zhai, and Xiaokang Yang. Uncertainty-aware blind image quality assessment in the laboratory and wild. *IEEE Transactions on Image Processing*, 30:3474–3486, Mar. 2021.

[41] Xin Zhang, Tianjian Ouyang, Yu Shang, Qingmin Liao, and Yong Li. UrbanMLLM: Joint learning of cross-view imagery for urban understanding, 2025.

[42] Yunke Zhang, Ruolong Ma, Xin Zhang, and Yong Li. Perceiving urban inequality from imagery using visual language models with chain-of-thought reasoning. In *Proceedings of the ACM on Web Conference 2025*, page 5342–5351. Association for Computing Machinery, 2025.

[43] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017.

[44] Zhilun Zhou, Jingyang Fan, Yu Liu, Fengli Xu, Depeng Jin, and Yong Li. Synergizing llm agents and knowledge graph for socioeconomic prediction in lbsn. *arXiv preprint arXiv:2411.00028*, 2024.

[45] Zhilun Zhou, Yu Liu, Jingtao Ding, Depeng Jin, and Yong Li. Hierarchical knowledge graph learning enabled socioeconomic indicator prediction in location-based social network. In *Proceedings of the ACM web conference 2023*, pages 122–132, 2023.

# A   Appendix

## A.1   Limitation

While CityRiSE demonstrates strong performance across a range of indicators, it still struggles with more abstract targets such as Accessibility to Health, where relevant cues are difficult to infer directly from imagery. Improving the model's capability on such challenging yet socially critical indicators remains an important direction for future work. In addition, we observe some instability during GRPO training, which occasionally leads to unexpected behaviors. Enhancing training robustness and sample efficiency presents another avenue for improvement.

## A.2   Training Details

In the CityRiSE setting, we use Qwen2.5-VL-7B as the base model. As described in Section 2, we apply the GRPO algorithm with verifiable reward to fine-tune the model on three datasets: Socio-Economic Indicator Data-Train, Perceptual Urban Reasoning Data, and General Visual Reasoning Data. The model is trained for 4 epochs using 4 A800 GPUs. We summarize the key hyper-parameters in training CityRiSE in Table 5.

Table 5: The hyper-parameter settings of CityRiSE.

| Hyper-parameters | Value |
|---|---|
| Global Batch Size | 8 |
| Learning Rate | 1e-6 |
| Weight Decay | 1e-2 |
| Optimizer Strategy | AdamW |
| Rollout N | 5 |
| Max Model Length | 120000 |
| Max Pixel | 100352 |
| Min Pixel | 50176 |

## A.3   Example of CoT Prompt and Answer

Figure 5 presents an example of UI-CoT training data for the Bachelor Ratio estimation task. This CoT format illustrates a step-by-step reasoning process, analyzing visual features such as buildings, vehicles, road infrastructure, street furniture, greenery, and other environmental cues, as well as geographic inference. The example concludes with a final prediction in the <answer>NUMBER</answer> format.

## A.4   Complete Results of Reward Ablation

As shown in Table 6, the Reward Ablation results on the Socio-Economic Indicator Data-Test are evaluated using the $R^2$ metric. The table compares the full CityRiSE model with its variants where different reward components are removed, across all status prediction tasks. The results show that the complete CityRiSE model achieves the best overall performance (Overall = 0.361), indicating

Figure 5: CoT prompt and answer example for the Bachelor Ratio task.

that both the keyword and regression rewards contribute significantly to improving the model's effectiveness.

Table 6: Reward Ablation results on Socio-Economic Indicator Data-Test. The values in the table represent $R^2$ scores. In each row, bold indicates the best result, and underline denotes the second-best.

| Domain | In-domain | | | | | Unseen Cities | | | | | | Unseen Indicators | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Tasks | GDP | Pop. | MH | PT | BR | GDP | Pop. | MH | PT | BR | DR | BH | VC | AH | LE | HP | |
| CityRiSE | 0.385 | 0.334 | **0.399** | **0.594** | **0.603** | -0.255 | **0.365** | 0.259 | -4.387 | **0.286** | 0.289 | **0.366** | **0.305** | -0.196 | **0.200** | **0.218** | **0.361** |
| CityRiSE w/o Keyword | -0.438 | 0.133 | -0.017 | -0.118 | -0.182 | -0.513 | 0.044 | 0.175 | -0.210 | -3.308 | 0.065 | -0.080 | -0.888 | -0.521 | -0.033 | -0.233 | -0.099 |
| CityRiSE w/o Regression | **0.417** | **0.349** | 0.310 | 0.442 | 0.416 | 0.279 | 0.055 | **0.333** | 0.016 | 0.146 | **0.391** | 0.014 | -0.142 | -0.541 | 0.097 | -0.177 | 0.303 |
| CityRiSE w/o Key. & Reg. | -1.197 | -0.313 | -0.543 | -1.483 | -0.289 | -0.508 | -1.001 | -0.003 | -0.326 | -2.883 | -1.243 | -0.825 | -0.966 | -1.387 | 0.129 | -1.067 | -0.489 |

## A.5 Complete Results of Data Ablation

Table 7 reports the results of our data ablation study on the Socio-Economic Indicator Data-Test set. We evaluate the contribution of different data modalities by progressively removing Perceptual Urban Reasoning Data and General Visual Reasoning Data from CityRiSE. The full model consistently achieves the best or second-best performance across most indicators, especially in the In-domain and Unseen Indicators settings, demonstrating the strong complementarity between Perceptual Urban Reasoning and General Visual Reasoning data sources. Removing Perceptual Urban Reasoning Data causes a notable performance drop in both In-domain and Unseen Cities scenarios, while removing General Visual Reasoning Data mainly affects cross-indicator generalization. The results confirm that multi-source integration plays a critical role in supporting robust and transferable urban socio-economic prediction.

Table 7: Data Ablation results on Socio-Economic Indicator Data-Test. The values in the table represent $R^2$ scores. In each row, bold indicates the best result, and underline denotes the second-best.

| Domain | In-domain | | | | | Unseen Cities | | | | | Unseen Indicators | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Tasks | GDP | Pop. | MH | PT | BR | GDP | Pop. | MH | PT | BR | DR | BH | VC | AH | LE | HP | |
| CityRiSE | 0.385 | 0.334 | 0.399 | 0.594 | **0.603** | -0.255 | 0.365 | 0.259 | -4.387 | **0.286** | **0.289** | 0.366 | 0.305 | -0.196 | **0.200** | **0.218** | **0.361** |
| CityRiSE w/o Per. | 0.459 | 0.400 | 0.309 | 0.428 | 0.147 | -0.118 | 0.034 | **0.285** | -12.750 | 0.206 | 0.288 | 0.005 | 0.199 | -1.014 | -0.070 | -0.334 | 0.034 |
| CityRiSE w/o Gen. | 0.453 | **0.512** | -0.152 | **0.613** | 0.597 | -0.381 | 0.294 | -0.715 | -4.080 | 0.205 | 0.176 | 0.293 | **0.333** | -0.148 | 0.109 | 0.105 | 0.295 |
| CityRiSE w/o Per. & Gen. | **0.485** | 0.489 | **0.402** | 0.531 | 0.511 | **-0.027** | **0.387** | -0.017 | **-2.474** | 0.184 | 0.110 | **0.407** | 0.292 | -0.298 | 0.093 | 0.129 | 0.356 |

Table 8: Reasoning Induction results on Socio-Economic Indicator Data-Test. The values in the table represent $R^2$ scores. In each row, bold indicates the best result, and underline denotes the second-best.

| Domain | In-domain | | | | | Unseen Cities | | | | | Unseen Indicators | | | | | | Overall |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Tasks | GDP | Pop. | MH | PT | BR | GDP | Pop. | MH | PT | BR | DR | BH | VC | AH | LE | HP | |
| UI-CoT (SFT-CoT) | 0.546 | 0.499 | 0.053 | 0.446 | 0.616 | -0.952 | -0.055 | 0.264 | -0.951 | 0.063 | -0.939 | **0.398** | 0.103 | -1.780 | -0.493 | -0.105 | 0.099 |
| UI-CoT + RL | **0.598** | **0.510** | **0.460** | **0.658** | **0.631** | -1.323 | 0.318 | **0.380** | **-0.199** | 0.013 | -1.154 | 0.339 | 0.107 | -1.644 | -1.439 | -0.014 | 0.091 |
| CityRiSE (RL) | 0.385 | 0.334 | 0.399 | 0.594 | 0.603 | **-0.255** | **0.365** | 0.259 | -4.387 | **0.286** | **0.289** | 0.366 | **0.305** | **-0.196** | **0.200** | **0.218** | **0.361** |

## A.6 Complete Results of Reasoning Induction

Table 8 presents the results of the Reasoning Induction experiments on the Socio-Economic Indicator Data-Test set. The table compares models trained with different reasoning induction strategies, including UI-CoT, UI-CoT + RL, and CityRiSE model. The results show that CityRiSE consistently achieves superior and more balanced performance across both in-domain and out-of-domain settings, suggesting that an RL-from-scratch training regime with reward-driven objectives encourages more flexible and transferable reasoning. By contrast, although UI-CoT + RL undergoes RL fine-tuning, it largely retains the rigid reasoning structures inherited from supervised training, which remain mostly unchanged during optimization and constrain its ability to generalize to unseen cities and indicators.

## A.7 Evaluation Data Statistics

We summarize here the composition of the evaluation datasets used across different generalization settings. For the in-domain evaluation, each of the GDP and Population prediction tasks contains 200 cases, while the Mental Health, Public Transport, and Bachelor Ratio tasks each include 50 cases. In the out-of-domain setting, all five indicator tasks—GDP, Population, Mental Health, Public Transport, and Bachelor Ratio—contain 200 cases each. Finally, the out-indicator evaluation covers six novel tasks that are unseen during training: Drive Ratio (200), Building Height (200), Violent Crime (200), Accessibility to Health (200), Life Expectancy (193), and House Price (200). Together, these datasets provide a comprehensive evaluation dataset spanning cross-city and cross-indicator generalization.