

Beyond Sight: Towards Cognitive Alignment in LVLm via Enriched Visual Knowledge

Yaqi Zhao^{1*}, Yuanyang Yin^{2*}, Lin Li¹, Mangan Lin³, Victor Shea-Jay Huang¹
 Siwei Chen¹, Weipeng Chen³, Baoqun Yin², Zenan Zhou^{3†}, Wentao Zhang^{1†}
¹Peking University ²University of Science and Technology of China ³Baichuan Inc.

Abstract

Does seeing always mean knowing? Large Vision-Language Models (LVLms) integrate separately pre-trained vision and language components, often using CLIP-ViT as vision backbone. However, these models frequently encounter a core issue of “cognitive misalignment” between the vision encoder (VE) and the large language model (LLM). Specifically, the VE’s representation of visual information may not fully align with LLM’s cognitive framework, leading to a mismatch where visual features exceed the language model’s interpretive range. To address this, we investigate how variations in VE representations influence LVLm comprehension, especially when the LLM faces VE-Unknown data—images whose ambiguous visual representations challenge the VE’s interpretive precision. Accordingly, we construct a multi-granularity landmark dataset and systematically examine the impact of VE-Known and VE-Unknown data on interpretive abilities. Our results show that VE-Unknown data limits LVLm’s capacity for accurate understanding, while VE-Known data, rich in distinctive features, helps reduce cognitive misalignment. Building on these insights, we propose Entity-Enhanced Cognitive Alignment (EECA), a method that employs multi-granularity supervision to generate visually enriched, well-aligned tokens that not only integrate within the LLM’s embedding space but also align with the LLM’s cognitive framework. This alignment markedly enhances LVLm performance in landmark recognition. Our findings underscore the challenges posed by VE-Unknown data and highlight the essential role of cognitive alignment in advancing multimodal systems.

1. Introduction

Large Vision Language Models (LVLms) [6, 9, 15–17, 19, 22, 27, 41, 45] have recently achieved significant advance-

ments. By mapping visual inputs into the embedding space of large language models, LVLms harness the powerful interpretative capabilities of language models [1, 26, 36] to address complex tasks like visual question answering, grounding, counting, etc.

Despite the impressive advancements in LVLms, these models still struggle with fundamental recognition tasks. As illustrated in Figure 1, even state-of-the-art models like GPT-4o and Qwen2-VL fail to recognize iconic landmarks from images, although they can describe these landmarks accurately when prompted with text alone. This raises critical questions: *Why do these challenges persist? To what extent do these models truly understand what they perceive?* Drawing inspiration from [8], we attribute this issue to a broader problem that we call **cognitive misalignment**—a fundamental disconnect between the representations generated by the vision encoder (VE) and the interpretive framework of the large language model (LLM).

Currently, most LVLms [3, 6, 19] integrate separately pretrained vision [29, 34, 44] and language models [1, 28, 36] through different adapters [2, 16, 19] to achieve multimodal comprehension. However, this simple connection often results in fundamental misalignment, as recent studies suggest [37, 39, 43], which limits the potential of LVLms. In this paper, we start by examining the cognitive framework of the vision encoder (VE) in LVLms (Section 3). To systematically investigate this, we define evaluation metrics based on CLIP’s training paradigm to assess VE knowledge. Using these metrics, we categorize data into VE-Known and VE-Unknown. Our study empirically demonstrates the following insights:

- **Enhanced alignment with VE-Known data:** VE-Known data, comprising images with rich and distinctive features, provides a strong, discriminative foundation that enhances the utilization of visual knowledge in LVLms. This data enables smoother alignment between the vision encoder (VE) and the language model (LLM), allowing for more effective cognitive integration.
- **Challenges with VE-Unknown data:** In contrast, VE-Unknown data, characterized by weak alignment with

*Equal contribution.

†Corresponding Author.

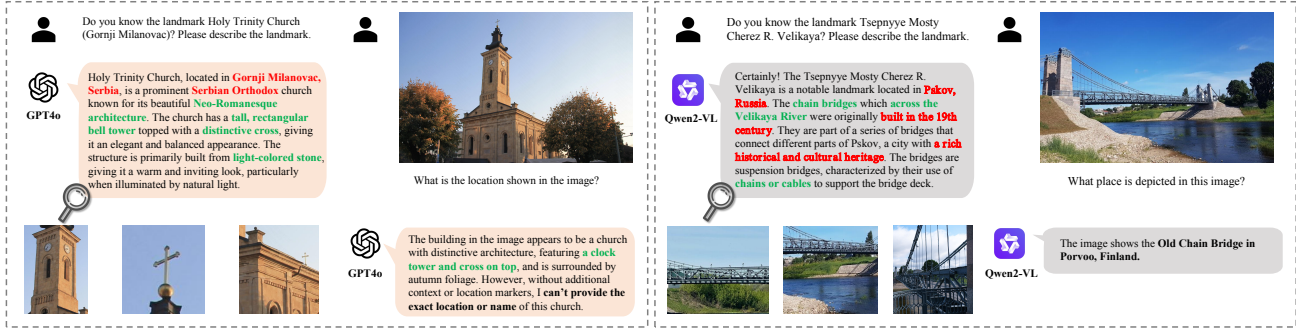


Figure 1. Instances of *Cognitive misalignment* are systematically identified, even in advanced models like GPT-4o and Qwen2-VL. Although the image closely aligns with the description generated from the text-only prompt, both models fail to recognize the landmark when presented with the image. Text highlighted in **green** emphasizes details that are particularly relevant to the image.

textual labels, poses significant alignment challenges. These representations hinder the VE’s ability to align with the LLM, leading to reduced interpretability and performance degradation on downstream tasks.

- **Importance of data quality over quantity:** The contrast between VE-Known and VE-Unknown data underscores that data quality has a greater impact than sheer volume. VE-Unknown samples tend to hinder cognitive integration, limiting the LVLM’s ability to effectively process complex visual inputs. This finding highlights the need to prioritize high-quality, VE-Known data to maximize LVLM performance.

Having identified the cognitive framework of the vision encoder in LVLM, we shift our focus toward enriching the visual knowledge in a way that aligns with the LLM’s cognitive framework, effectively “opening the eyes” of the LLM and achieving cognitive alignment (Section 4). To this end, we construct **Multi-granularity Landmark Dataset (MGLD)** that ensures consistency between the visual input and the language model’s output. Building on this, we propose the Entity-Enhanced Cognitive Alignment (EECA), a structured approach that enhances the richness and discriminative power of visual tokens through supervised multi-granularity learning. Through this multi-faceted approach, EECA not only minimizes information loss in the adapter’s transformation but also enables the transformed visual tokens to mimic VE-Known representations.

Our contributions are as follows. First, we provide an in-depth analysis of the impact of VE representations on the LLM’s understanding, shedding light on the challenges of cognitive alignment. Building on these insights, we propose EECA, a novel method that introduces discriminative features from the LLM’s recognition space into visual tokens through supervised learning, enhancing their effectiveness for cognitive alignment. Finally, extensive experiments and analyses show that our method significantly enhances landmark recognition performance, with notable improvements

observed on both VE-Known and VE-Unknown data, consistently outperforming baseline models.

2. Preliminary

Notation. In typical Large Vision Language Models (LVLMs) [6, 7, 15, 18–20, 22, 38, 41, 42], an **adapter** is used to connect VE and LLM seamlessly. This adapter, denoted by g_θ , can be a simple linear layer or a more complex attention module. For clarity, we define the features output by the vision encoder as *visual patches* and the features after passing through the adapter as *visual tokens*.

In our framework, the vision encoder f_v and text encoder f_t are from CLIP’s dual modules. The vision encoder f_v processes an image I_i to produce $V_i = f_v(I_i) \in \mathbb{R}^{P \times d}$, where d is the feature dimension and P is the number of visual patches. The text encoder f_t maps text T_i to $t_i = f_t(T_i) \in \mathbb{R}^d$. During LVLM training, the LLM input is constructed as follows:

$$X_v = g_\theta(f_v(I_i)), \quad X_t = \phi(T_i), \quad (1)$$

where g_θ transforms the vision encoder’s output to visual tokens $X_v \in \mathbb{R}^{N_v \times C}$, and ϕ maps text T_i to text tokens $X_t \in \mathbb{R}^{N_t \times C}$ in the LLM’s embedding space. The visual tokens X_v and text tokens X_t are then concatenated and fed into the LLM to generate a response.

Cognitive misalignment: issues and challenges. We identified a critical issue in LVLMs, illustrated through the example in Figure 1, which reveals a phenomenon we term *cognitive misalignment*. This issue arises from discrepancies between Vision Encoder’s representations and the language model’s cognitive framework. For instance, when asked, “Do you know the landmark Holy Trinity Church (Gornji Milanovac)? Please describe the landmark,” the model provides an accurate description based on its textual knowledge. However, when shown an image that corresponds to this description and asked, “What is the location

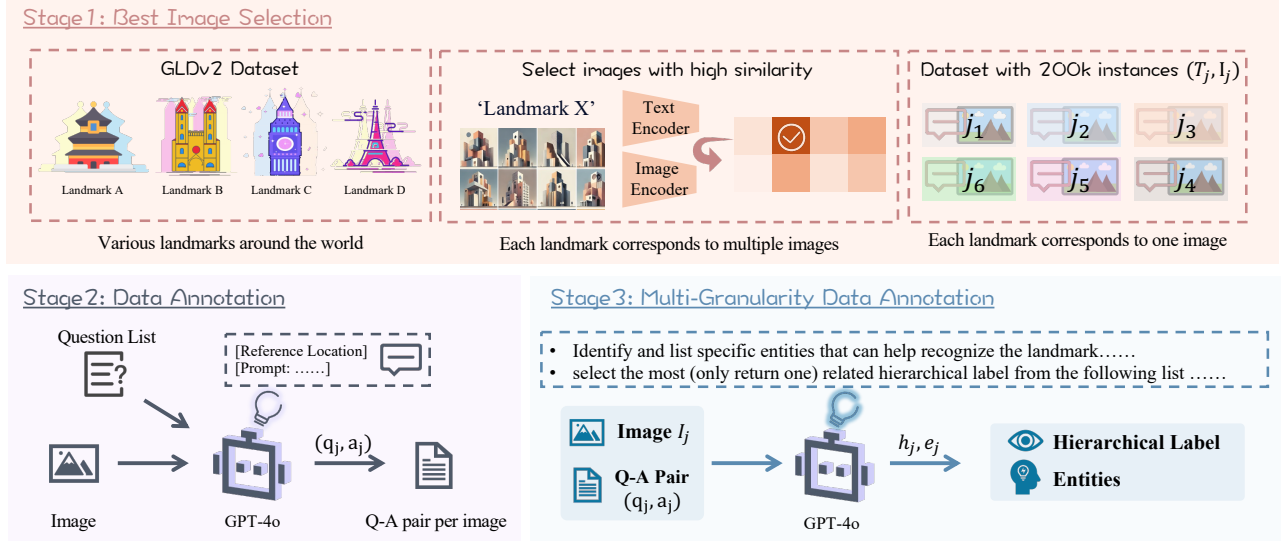


Figure 2. Illustration of the dataset construction process for the Multi-granularity Landmark Dataset (MGLD), showing three stages: best image selection using CLIP similarity, data annotation with Q-A pairs, and multi-granularity data annotation.

shown in the image?”, it fails to recognize it. This misalignment highlights a core challenge: although the language model can describe landmarks accurately in text, the VE’s visual outputs do not align closely enough with the language model’s cognitive framework to enable reliable recognition from images. This insight drives our exploration of different types of visual knowledge—*VE-Known* and *VE-Unknown*—which we elaborate on in later sections.

3. From Sight to Insight

In this section, we explore the cognitive alignment between the VE and the LLM within LVLMs. By constructing a fine-tuned dataset and evaluating visual knowledge across different levels, we aim to systematically analyze the quality of visual outputs and their impact on achieving cognitive alignment with the LLM’s understanding.

3.1. Study setup

Dataset construction. To explore cognitive alignment between Vision Encoders (VE) and large language models (LLMs) in LVLMs, we use an entity-related dataset where both the visual and language models have their own distinct representations and understandings. This allows us to systematically examine how well visual outputs align with the LLM’s understanding space. In light of this, we construct a fine-tuning dataset, termed the **Multi-Granularity Landmark Dataset (MGLD)**, comprising approximately 200k samples from the GLDv2 [40] dataset, following the process illustrated in Figure 2 (detailed in Appendix B). In this section, we focus solely on image-text pairs (I_i, T_i) in Stage 1 and Q-A pairs (q_i, a_i) in Stage 2.

Specifically, for Stage 1, each landmark in the original dataset corresponds to multiple images, and we pair the landmark name T_i with the image I_i that has the highest CLIP similarity to its corresponding landmark names. This selection process ensures that each chosen image aligns closely with the landmark identity. In Stage 2, we generate Question-Answer (Q-A) pairs to support landmark recognition. We randomly select a question from a predefined set and provide it to GPT-4o along with the selected image and corresponding landmark name. The model then responds with answers describing the landmark. Then, the dataset used in this section can be denoted as $D = \{(I_i, T_i, q_i, a_i)\}_{i=1}^N$, where I_i represents the image, q_i is a question related to the landmark, and a_i is the corresponding answer, including both image descriptions and location-related information. T_i is the ground-truth landmark name (e.g., “Eiffel Tower”).

Measuring knowledge of vision encoder. We calculate the CLIP similarity between the visual representation of each image I_i and all the landmark names T_j (e.g., “Eiffel Tower”) in the dataset, using the CLIP vision and text encoders. The similarity, denoted as $\text{Sim}_{\text{CLIP}}(I_i, T_j)$, is computed as follows:

$$\text{Sim}_{\text{CLIP}}(I_i, T_j) = \frac{\langle f_v(I_i), f_t(T_j) \rangle}{\|f_v(I_i)\| \|f_t(T_j)\|} \quad (2)$$

where $f_v(I_i)$ and $f_t(T_j)$ represent the visual and textual embeddings of the image and the landmark name, respectively.

We record the **Similarity Score** between the image I_i and its corresponding ground-truth label T_i , denoted as

$\text{Sim}_{\text{CLIP}}^i$. We also determine the **Relative Similarity Rank (RSR)** of the ground-truth T_i among all the landmark candidates T_j for the image I_i , based on the similarity scores, denoted as RSR^i , indicating the discriminative power of the visual representation.

We apply various data selection methods to the full dataset based on $\text{Sim}_{\text{CLIP}}^i$ and RSR^i , each aimed at capturing specific characteristics that indicate whether the Vision Encoder is in a Known or Unknown state. For consistency, all subsets are of equal size.

- **High Discrimination Selection (HDS):** This method selects images with very high RSR^i values, capturing instances where the model effectively distinguishes the ground-truth T_i from other candidates. HDS emphasizes strong visual discrimination in the VE’s representations, aligning with “*VE-Known*” characteristics.
- **High Similarity Selection (HSS):** This method focuses on images with high $\text{Sim}_{\text{CLIP}}^i$ values, highlighting the VE’s capacity for capturing rich visual features. HSS prioritizes visual richness over discriminative power, aligning with general “*VE-Known*” characteristics.
- **Low Clarity Selection (LCS):** This method selects images with both low $\text{Sim}_{\text{CLIP}}^i$ and low RSR^i values, targeting visually ambiguous cases where the model struggles with feature extraction and differentiation. These instances reflect “*VE-Unknown*” characteristics, indicating limited VE understanding.
- **Balanced Reference Selection (BRS):** This method randomly selects images to serve as a comparative baseline.

Evaluation metric. To evaluate LVLm performance, we generate five inference outputs for each question in the test dataset. These responses are assessed across four recognition levels based on their accuracy and detail relative to the correct answer by GPT-4o, with specific examples provided in Appendix C.

- **Strongly Known.** The model consistently delivers detailed, accurate information that closely aligns with the ground truth across multiple responses.
- **Known.** At least one response is correct and includes reasonable explanations.
- **Weakly Unknown.** The correct entity is not identified, but responses provide relevant hints, such as geographic, architectural, or contextual cues, suggesting an association with the target landmark.
- **Unknown.** None of the responses accurately identify the target. Information is either unrelated or overly generic, lacking specific clues about the target’s identity.

The evaluation metrics are defined as follows: each recognition level is quantified by its proportion within the test dataset. **Accuracy** is defined as the proportion of recognized responses (*i.e.*, the sum of Strongly Known and

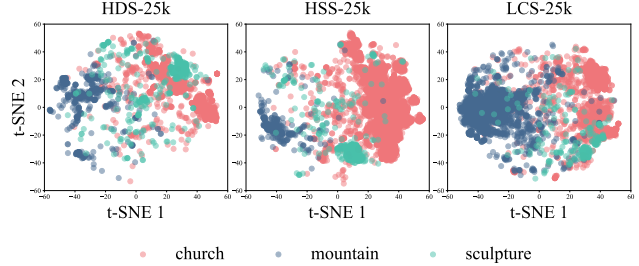


Figure 3. t-SNE visualization of image features. **Left:** The HDS subset shows more dispersed representations for categories (*e.g.*, “church”). **Middle:** The HSS subset shows distinct inter-class separations. **Right:** The LCS subset shows reduced intra-class variability and less distinct inter-class separations.

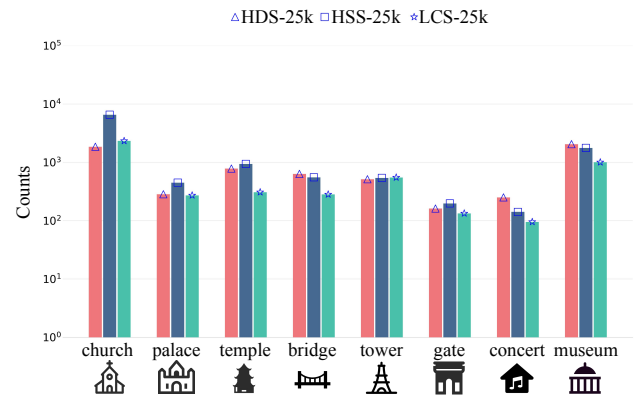


Figure 4. Category counts across subsets.

Known) in the test dataset, representing the model’s overall recognition capability.

3.2. Visual patterns of different knowledge

Firstly, we analyze the characteristics of different types of visual knowledge, focusing on how variations in visual representation affect cognitive alignment between VE and LLM. As shown in Figure 3, we used t-SNE to visualize the visual knowledge patterns extracted by the VE across distinct subsets. Specifically, we constructed three datasets: HDS-25k (top 25k samples with the highest RSR^i), HSS-25k (top 25k samples with the highest $\text{Sim}_{\text{CLIP}}^i$), and LCS-25k (25k samples with both low RSR^i and low $\text{Sim}_{\text{CLIP}}^i$). Additionally, Figure 4 shows the distribution of major categories in each subset, offering insights into their visual class composition.

The t-SNE results reveal distinct patterns across these subsets. The HDS subset shows dispersed visual representations within categories, indicating that the visual encoder captures fine-grained intra-class diversity, which enhances subtle landmark recognition. Conversely, the HSS subset forms compact, well-separated clusters, suggesting a focus on broad, shared patterns that aid in distinguish-

ing distinct categories like “mountain” and “sculpture”. The HSS subset also shows the most variation in category counts (see Figure 4), supporting its strong inter-class separability. The LCS subset displays reduced intra-class variability and less distinct inter-class boundaries, suggesting challenges in category differentiation.

Category	Selection Method	Data Size	
		25k	50k
Reference	BRS	40.1	56.2
VE-Unknown	LCS	23.0 (-17.1)	28.1 (-28.1)
VE-Known	HSS	40.1	60.4 (+4.2)
VE-Known	HDS	49.3 (+9.2)	64.1 (+7.9)

Table 1. The baseline accuracy for non-extra data is 8.68%. The values in the table represent the percentage increase in test accuracy over this baseline, with values in parentheses indicating changes relative to BRS. VE-Unknown shows reduced performance compared to BRS, while VE-Known methods improve accuracy, with HDS achieving the highest increase.

3.3. Uncovering the impact of visual knowledge

After analyzing the visual patterns of different knowledge extracted by the VE, we investigate a key question: *how do the strength and quality of VE representations influence cognitive alignment with the LLM’s understanding?* To answer this, we evaluated models trained with various data selection strategies, starting from the LLaVA [19] initialization and optimizing both the LLM and the adapter. Strong alignment—indicated by a higher proportion of recognized responses (sum of Strongly Known and Known)—suggests that visual features are both robust and supportive of cross-modal cognitive integration.

To measure improvement, we compared the proportion of recognized responses to the LLaVA baseline, focusing on the percentage increase in accurate landmark identifications. Our findings, presented in Table 1 and Figure 5, reveal three insights:

- VE-Known data mitigates cognitive misalignment.** As shown in Table 1, VE-Known subsets (HDS and HSS) consistently outperform BRS, especially with smaller datasets. The t-SNE analysis indicates that these subsets have richer, more discriminative visual features, facilitating smoother knowledge transfer from the VE to the LLM.
- VE-Unknown data exacerbates cognitive misalignment.** The Low Clarity Set (LCS), which includes VE-Unknown instances with ambiguous representations, shows the lowest performance (see Table 1). While VE-Unknown data may contribute partially to LVLM learning, its low confidence and information loss hinder effective

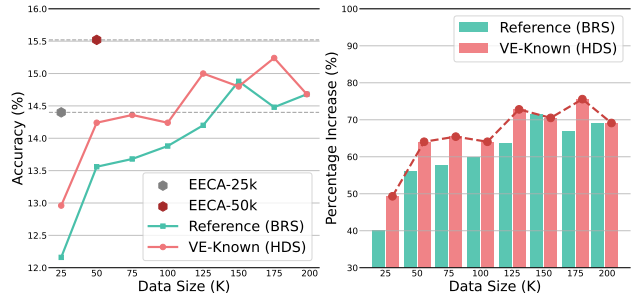


Figure 5. Comparative performance of HDS and BRS selection methods across different dataset sizes. **Left:** Accuracy (%) vs. Data Size for BRS and HDS with two point of EECA at 25k and 50k(best). **Right:** Percentage increase over baseline. VE-Known data outperforms Reference, demonstrating its effectiveness.

tive cognitive integration, underscoring the need to address VE-Unknown data for better model performance.

- Diminishing returns from increased data volume with VE-Unknown samples.** As shown in Figure 5, performance gains plateau with larger datasets, especially when VE-Unknown samples from LCS are included, indicating that data quality has a greater impact than data volume beyond a certain threshold.

Our experiment demonstrates that models trained on quality-driven selections (*e.g.*, HDS), can achieve competitive or superior performance with reduced data sizes. This finding highlights a crucial insight from the perspective of the VE’s cognitive framework: to strengthen cognitive alignment between the VE and the LLM, the information passed from the adapter must be both discriminative and richly informative, helping the LLM distinguish landmarks with greater accuracy. Building on this insight, the following section introduce a method that leverages entity-aware supervision to mitigate cognitive integration challenges in LVLMs.

4. Open the Eyes of LVLM

After examining the cognitive framework of the VE, a natural question arises: *how can we fully “open the eyes” of the LVLM to achieve cognitive alignment between visual and language components?* In this section, we take initial steps toward answering this question. First, we design a data annotation pipeline to ensure consistency between the visual input and the language model’s output (see Section 4.1). Building on this, we introduce the Entity-Enhanced Cognitive Alignment (EECA) framework, which supervises the adapter’s visual tokens to retain richness and discriminative power, minimizing information loss (see Section 4.2). This approach encourages transformed visual tokens to “mimic” VE-Known representations, facilitating alignment with the LLM’s cognitive framework.

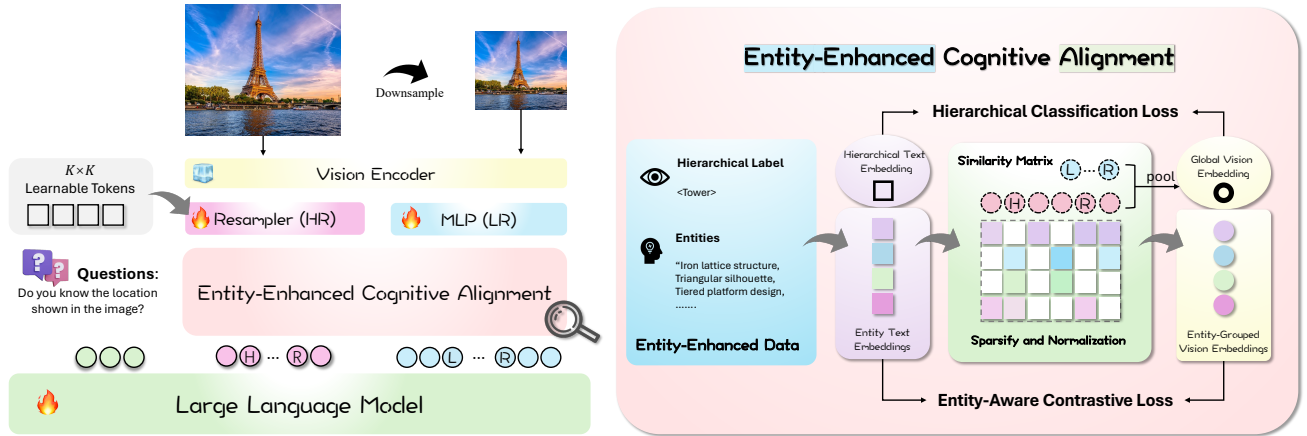


Figure 6. Overview of our model’s framework. EECA approach combines high- and low-resolution visual features through a dual-branch architecture, supporting alignment with the language model’s cognitive framework. The framework includes a hierarchical classification loss and an entity-aware contrastive loss to encourage rich, discriminative representations aligned with entity-specific information.

4.1. Multi-granularity data annotation pipeline

To construct a landmark instruction dataset that enhances the model’s recognition and differentiation capabilities, we designed a 3-stage data pipeline as shown in Figure 2. Stage 1 and 2 have been introduced in Section 3.1. Here, we focus on providing a detailed description of Stage 3.

The final stage, Multi-granularity data annotation, enriches each landmark with hierarchical labels and unique entities, providing a layered structure that captures both broad classifications and fine-grained details. Firstly, each landmark is assigned a hierarchical label (e.g., “mountain,” “lake,” “church,”), placing it within a general category. In contrast, entities provides a finer level of detail, capturing landmark-specific attributes. For instance, the Eiffel Tower might be associated with entities like “iron lattice structure” and “triangular silhouette”, which highlight its distinctive physical attributes. These entities are annotated based on both image content and QA pairs, making them directly relevant to the LLM’s answer outputs.

The resulting dataset, which we denote as $D = \{(I_j, q_j, a_j, h_j, e_j)\}_{j=1}^N$, is structured such that each entry contains an image I_j , a question q_j , an answer a_j incorporating both descriptive and location-specific information, a hierarchical label h_j representing the broader category, and unique entities e_j detailing fine-grained, landmark-specific attributes. This dataset construction enables EECA to capture the unique identity of each landmark through visual tokens, facilitating cognitive alignment between the vision encoder and the language model.

4.2. Entity-Enhanced cognitive alignment

Architecture overview. Within this enriched dataset, we train our model using EECA, as illustrated in Figure 6.

EECA promotes cognitive alignment by supervising the transformation from visual patches to tokens, ensuring that the tokens retain rich, discriminative information linked to specific entities in MGLD. These entities are annotated to align with the LLM’s cognitive framework, supporting cross-modal understanding.

Inspired by recent work [17, 32], we employ a dual-branch visual architecture to leverage high-resolution (HR) information, enhancing the richness of visual representations. Starting with a high-resolution image I^H , we generate a low-resolution version I^L via bilinear interpolation. I^H is divided into four sub-images, which, along with I^L , are fed to a shared vision encoder, producing low-resolution and high-resolution visual features, V^L and V^H .

In the low-resolution branch, a 2-layer MLP adapter outputs visual tokens $X_v^L \in \mathbb{R}^{N_{vL} \times C}$. In the high-resolution branch, a shared perceiver resampler [2] generates high-resolution visual tokens $X_v^H \in \mathbb{R}^{N_{vH} \times C}$. To make HR information efficient, EECA supervises the compression of HR features, aligning visual tokens with the LLM’s embedding space. The LLM then integrates textual and visual information for comprehensive understanding. The EECA framework is optimized with two loss functions: Entity-Aware Contrastive Loss and Hierarchical Classification Loss, detailed below.

Entity-Aware contrastive loss. Given the concatenated visual tokens $X_v = [X_v^L; X_v^H]$, where $X_v^H = (X_{v_1}^H, \dots, X_{v_{N_{vH}}}^H) \in \mathbb{R}^{N_{vH} \times C}$ represents high-resolution tokens and X_v^L represents low-resolution tokens, this loss function encourages the model to incorporate fine-grained, entity-specific details from X_v^H into the primary representation in X_v^L .

Motivated by the idea that different queries can capture diverse aspects of visual information [12], we aim to learn a combination of query embeddings that corresponds to each entity’s token. Specifically, we construct an entity-grouped vision embedding for each entity token $\phi(e_j) = X_{e_j} \in \mathbb{R}^C$ as a weighted combination of high-resolution tokens X_v^H , based on a sparsified and normalized similarity matrix, inspired by [4]. This approach preserves entity-relevant characteristics over general visual features.

Following [43], we propose an Entity-Aware Contrastive Loss to align high-resolution tokens with entity-specific information in the LLM’s embedding space. Formally, let $W_{i,j}$ be the weight matrix based on the similarity between high-resolution tokens $X_{i,v}^H$ and entity embeddings $e_{i,j}$, where i is the i -th image and j is the j -th entity within that image. The entity-grouped visual embedding $\tilde{X}_{e_{i,j}}$ is defined as:

$$\tilde{X}_{e_{i,j}} = \sum_{v=1}^{N_{vH}} W_{i,j} X_{i,v}^H \quad (3)$$

The Entity-Aware Contrastive Loss \mathcal{L}_{ec} is then calculated as:

$$\begin{aligned} \mathcal{L}_e = & -\frac{1}{2B} \sum_{i=1}^B \sum_{j=1}^{E_i} \left(\log \frac{\exp(S(X_{e_{i,j}}, \tilde{X}_{e_{i,j}})/\tau)}{\sum_{k=1}^{E_i} \exp(S(X_{e_{i,j}}, \tilde{X}_{e_{i,k}})/\tau)} \right. \\ & \left. + \log \frac{\exp(S(\tilde{X}_{e_{i,j}}, X_{e_{i,j}})/\tau)}{\sum_{k=1}^{E_i} \exp(S(\tilde{X}_{e_{i,j}}, X_{e_{i,k}})/\tau)} \right) \end{aligned} \quad (4)$$

where B is the batch size, E_i is the number of entities for the i -th image, $S(a, b) = \frac{a \cdot b}{\|a\| \|b\|}$, and τ is the temperature.

Hierarchical classification loss. To enhance the model’s understanding across multiple levels of granularity, we introduce a hierarchical classification task. Using the hierarchical labels in MGLD (e.g., “church,” “lake,”), this task builds between-category distinctions on top of the within-category distinctions captured through entities.

For each image, classification features are generated by concatenating high-resolution and low-resolution visual tokens, $X_{i,v}^H$ and $X_{i,v}^L$, followed by average pooling to obtain a comprehensive representation \mathbf{h}_i . This representation is then used to compute the hierarchical classification loss, defined as:

$$\mathcal{L}_h = -\frac{1}{B} \sum_{i=1}^B \log P(y_i | \mathbf{h}_i) \quad (5)$$

where $P(y_i | \mathbf{h}_i)$ is the predicted probability of the correct class y_i for the i -th image.

Overall objective. The total loss function used to train the model combines three components: a standard language

modeling loss (\mathcal{L}_g), a hierarchical classification loss (\mathcal{L}_h), and an Entity-Aware Contrastive Loss (\mathcal{L}_e), which together support EECA learning.

The language modeling loss \mathcal{L}_g encourages accurate text generation based on the visual input and is defined as:

$$\mathcal{L}_g = -\frac{1}{B} \sum_{i=1}^B \sum_{t=1}^{T_i} \log P(w_{i,t} | w_{i,<t}, X_i) \quad (6)$$

Overall combined objective can be expressed as $\mathcal{L} = \lambda \mathcal{L}_g + \mu_e \mathcal{L}_e + \mu_h \mathcal{L}_h$ where λ , μ_e , and μ_h are balancing coefficients for each loss component.

This layered learning process supports cognitive alignment between modalities, as enriched and distinct visual tokens better align with the LLM’s cognitive framework. Optimizing this combined loss enables multi-granularity learning.

5. Experiments

In this section, we analyze the effectiveness of EECA on landmark recognition tasks. To evaluate the impact of our approach, we conduct experiments on different variations and data subsets, comparing performance across multiple configurations.

5.1. Experimental setup

Our experiments are conducted within the LLaVA-1.5 [19], using a dual-branch visual architecture for enhanced feature representation. We jointly optimize both the LLM and the two adapters within the EECA framework. Further experimental details, including hyperparameter configurations, are provided in Appendix A. For data preparation, we construct the dataset according to the process described in Figure 2, creating subsets with varying data selection methods (See Section 3.1).

5.2. Results and analysis

Effectiveness of EECA. We evaluate EECA against three configurations. *Baseline* serves as a reference without entity prompts. *Entity Prompt (Inference Only)* includes entities only during inference, confirming that entity prompts do not negatively impact performance. *Entity Prompt (Training + Inference)* incorporates entities during both training and inference, yielding significantly higher accuracy and indicating alignment with the LLM’s cognitive framework. As shown in Figure 5, with only 25k data, EECA reaches the performance of the 125k reference dataset, and with 50k data, it achieves the best results. EECA thus provides balanced improvements without extensive reliance on entity prompts during inference, demonstrating effectiveness in enhancing performance without overfitting. Overall, EECA offers a more efficient, focused approach, achieving strong alignment with the LLM’s cognitive framework.

Method	Strongly Known	Known
Baseline	4.12	4.56
Entity Prompt (Inference Only)	4.56	3.24
Entity Prompt (Training + Inference)	19.52	9.32
EECA	8.52	7.0

Table 2. Comparison of different methods.

Method	Strongly Known	Known	Accuracy
Baseline	4.12	4.56	8.68
+ <i>HSS-50k</i>	7.48 (+3.36)	6.44 (+1.88)	13.92 (+5.24)
+ <i>HR Branch</i>	7.92 (+0.44)	5.96 (-0.48)	13.88 (-0.04)
+ \mathcal{L}_e	8.48 (+0.56)	5.92 (-0.04)	14.4 (+0.52)
+ \mathcal{L}_h	8.52 (+0.04)	7.00 (+1.08)	15.52 (+1.12)

Table 3. Ablation Study Results.

Ablation study. Table 3 demonstrates that the main accuracy gains come from the proposed entity-aware contrastive loss (\mathcal{L}_e) and hierarchical classification loss (\mathcal{L}_h). In contrast, adding the high-resolution (HR) branch alone does not improve performance without targeted supervision. These results highlight that EECA’s supervision allows the HR branch to deliver richer, more discriminative visual information, aligning the visual encoder’s outputs with the language model’s cognitive framework and enhancing interpretability in landmark recognition.

5.3. Generalizability of EECA

To validate EECA’s generalizability, we applied it to three distinct data subsets representing different levels of visual knowledge. Results in Table 4 show that EECA consistently improves performance across all subsets, demonstrating adaptability even in challenging VE-Unknown scenarios.

For LCS-25k (VE-Unknown data), adding the HR branch significantly boosts performance, indicating that enriched visual information in the HR branch effectively reduces visual ambiguity by enhancing feature richness and separability. In contrast, HDS-25k and HSS-25k (VE-Known data) achieve the largest gains with the entity-aware contrastive loss (\mathcal{L}_e) and hierarchical classification loss (\mathcal{L}_h), showing that VE-Known data benefit from targeted supervision to extract discriminative features.

Overall, these results underscore EECA’s robustness and adaptability, addressing visual ambiguity in VE-Unknown while enhancing feature discrimination in VE-Known.

6. Related Work

Large vision language models. The integration of vision and language models has advanced Large Vision Language

Method	HDS-25k	HSS-25k	LCS-25k
Baseline	8.68	8.68	8.68
+ <i>25k Data</i>	13.00 (+4.32)	13.00 (+4.32)	10.68 (+2.00)
+ <i>HR Branch</i>	13.60 (+4.92)	13.84 (+5.16)	12.08 (+3.40)
+ \mathcal{L}_e	14.00 (+5.32)	14.40 (+5.72)	12.32 (+3.64)
+ \mathcal{L}_h	14.40 (+5.72)	13.84 (+5.16)	12.32 (+3.64)

Table 4. Performance comparison across VE-Known (HDS, HSS) and VE-Unknown (LCS) subsets.

Models for enhanced cross-modal understanding. Foundational models like CLIP [29] demonstrated the potential of cross-modal applications, while recent models such as Flamingo [2], BLIP-2 [16], and LLaVA [19, 20] refine alignment through efficient strategies and large-scale data. Specialized models like Shikra [5], which outputs spatial coordinates in natural language, and Mini-Gemini [17], which employs an additional visual encoder for high-resolution enhancement, further expand functionality. Qwen2-VL [36] introduces a dynamic resolution mechanism for adaptive tokenization across varying image resolutions. Despite these advancements, achieving cognitive alignment between VE and LLM remains challenging.

Visual challenges in LVLMs. Despite significant progress, many visual challenges remain due to inherent limitations in vision encoders like CLIP. While widely adopted, CLIP often fails to capture comprehensive visual information, hindering LVLm performance. Research [39] shows that CLIP struggles with various visual tasks, reducing the effectiveness of CLIP-based LVLms. Another study [35] suggests that CLIP may treat visual inputs as a “bag of concepts,” missing higher-level structures like part-whole relationships. Furthermore, [37] found that LVLms often fail on simple questions due to CLIP’s pre-training limitations, which overlook critical visual details and struggle to prioritize significant patterns. Additionally, the vision encoder’s performance is constrained by resolution and data quality [36, 45]. These studies underscore the need for techniques to fully unlock the potential of LVLms.

7. Conclusion and Discussion

This study revisits the question: does seeing always mean knowing? In Large Vision Language Models (LVLms), we find that this is often not the case. Our investigation reveals a critical cognitive misalignment between the vision encoder (VE) and the large language model (LLM), where VE’s visual representations do not fully align with the LLM’s cognitive framework. Our results show that LVLm performance is closely tied to the knowledge within the VE: VE-Known data alleviate cognitive misalignment,

while VE-Unknown data exacerbate it. To address this gap, we propose Entity-Enhanced Cognitive Alignment (EECA), which supervises the adapter’s visual tokens to retain richness and discriminative power through multigranular representation, enabling them to “mimic” VE-Known behavior. Our results demonstrate that EECA is effective, robust, and generalizable, improving performance across both VE-Known and VE-Unknown data.

Despite these advancements, EECA has certain limitations. Its effectiveness relies on labeled data to generate supervision signals, limiting its applicability to tasks with available annotations. Additionally, EECA’s design focuses on entity recognition and does not provide a general solution for cognitive misalignment beyond this scope. We hope our findings inspire future work toward broader solutions for cognitive alignment in multimodal systems.

References

- [1] AI@Meta. Llama 3 model card. 2024. 1
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. 1, 6, 8
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A frontier large vision-language model with versatile abilities. *arXiv preprint arXiv:2308.12966*, 2023. 1
- [4] Ioana Bica, Anastasija Ilić, Matthias Bauer, Goker Erdogan, Matko Bošnjak, Christos Kaplanis, Alexey A Gritsenko, Matthias Minderer, Charles Blundell, Razvan Pascanu, et al. Improving fine-grained understanding in image-text pre-training. *arXiv preprint arXiv:2401.09865*, 2024. 7
- [5] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal llm’s referential dialogue magic. *arXiv preprint arXiv:2306.15195*, 2023. 8
- [6] Wenliang Dai, Junnan Li, D Li, AMH Tiong, J Zhao, W Wang, B Li, P Fung, and S Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv 2023. arXiv preprint arXiv:2305.06500*, 2, 2023. 1, 2
- [7] Xiaoyi Dong, Pan Zhang, Yuhang Zang, Yuhang Cao, Bin Wang, Linke Ouyang, Songyang Zhang, Haodong Duan, Wenwei Zhang, Yining Li, et al. Internlm-xcomposer2-4khd: A pioneering large vision-language model handling resolutions from 336 pixels to 4k hd. *arXiv preprint arXiv:2404.06512*, 2024. 2
- [8] Zorik Gekhman, Gal Yona, Roei Aharoni, Matan Eyal, Amir Feder, Roi Reichart, and Jonathan Herzig. Does fine-tuning llms on new knowledge encourage hallucinations? *arXiv preprint arXiv:2405.05904*, 2024. 1
- [9] Google. Bard, 2023. 1
- [10] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the V in VQA matter: Evaluating the role of image understanding in visual question answering. In *CVPR*, 2017. 1
- [11] Drew A Hudson and Christopher D Manning. GQA: A new dataset for real-world visual reasoning and compositional question answering. In *CVPR*, 2019. 1
- [12] Oğuzhan Fatih Kar, Alessio Tonioni, Petra Poklukar, Achin Kulshrestha, Amir Zamir, and Federico Tombari. Brave: Broadening the visual encoding of vision-language models. *arXiv preprint arXiv:2404.07204*, 2024. 7
- [13] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *EMNLP*, 2014. 1
- [14] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *IJCV*, 2017. 1
- [15] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv preprint arXiv:2405.02246*, 2024. 1, 2
- [16] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv:2301.12597*, 2023. 1, 8
- [17] Yanwei Li, Yuechen Zhang, Chengyao Wang, Zhisheng Zhong, Yixin Chen, Ruihang Chu, Shaoteng Liu, and Jiaya Jia. Mini-gemini: Mining the potential of multi-modality vision language models. *arXiv preprint arXiv:2403.18814*, 2024. 1, 6, 8
- [18] Dongyang Liu, Renrui Zhang, Longtian Qiu, Siyuan Huang, Weifeng Lin, Shitian Zhao, Shijie Geng, Ziyi Lin, Peng Jin, Kaipeng Zhang, et al. Sphinx-x: Scaling data and parameters for a family of multi-modal large language models. *arXiv preprint arXiv:2402.05935*, 2024. 2
- [19] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *arXiv preprint arXiv:2310.03744*, 2023. 1, 5, 7, 8
- [20] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. 2023. 2, 8, 1
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2017. 1
- [22] Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhuoshu Li, Hao Yang, et al. Deepseek-vl: towards real-world vision-language understanding. *arXiv preprint arXiv:2403.05525*, 2024. 1, 2
- [23] Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. Generation and comprehension of unambiguous object descriptions. In *CVPR*, 2016. 1
- [24] Kenneth Marino, Mohammad Rastegari, Ali Farhadi, and Roozbeh Mottaghi. OK-VQA: A visual question answering benchmark requiring external knowledge. In *CVPR*, 2019.
- [25] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. OCR-VQA: Visual question answering by reading text in images. In *ICDAR*, 2019. 1
- [26] OpenAI. Gpt-4 technical report, 2023. 1

- [27] OpenAI. GPT-4o System Card, 2024. 1
- [28] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*, 2023. 1
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 1, 8
- [30] Dustin Schwenk, Apoorv Khandelwal, Christopher Clark, Kenneth Marino, and Roozbeh Mottaghi. A-OKVQA: A benchmark for visual question answering using world knowledge. In *ECCV*, 2022. 1
- [31] Piyush Sharma, Nan Ding, Sebastian Goodman, and Radu Soricut. Conceptual captions: A cleaned, hypernymed, image alt-text dataset for automatic image captioning. In *ACL*, 2018. 1
- [32] Baifeng Shi, Ziyang Wu, Maolin Mao, Xin Wang, and Trevor Darrell. When do we not need larger vision models? In *European Conference on Computer Vision*, pages 444–462. Springer, 2025. 6
- [33] Oleksii Sidorov, Ronghang Hu, Marcus Rohrbach, and Amanpreet Singh. Textcaps: a dataset for image captioning with reading comprehension. In *ECCV*, 2020. 1
- [34] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. Eva-clip: Improved training techniques for clip at scale. *arXiv preprint arXiv:2303.15389*, 2023. 1
- [35] Yingtian Tang, Yutaro Yamada, Yoyo Zhang, and Ilker Yildirim. When are lemons purple? the concept association bias of vision-language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14333–14348, 2023. 8
- [36] Qwen Team. Qwen2.5: A party of foundation models, 2024. 1, 8
- [37] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann LeCun, and Saining Xie. Eyes wide shut? exploring the visual shortcomings of multimodal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9568–9578, 2024. 1, 8
- [38] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 2
- [39] Haoran Wei, Lingyu Kong, Jinyue Chen, Liang Zhao, Zheng Ge, Jinrong Yang, Jianjian Sun, Chunrui Han, and Xiangyu Zhang. Vary: Scaling up the vision vocabulary for large vision-language model. In *European Conference on Computer Vision*, pages 408–424. Springer, 2025. 1, 8
- [40] T. Weyand, A. Araujo, B. Cao, and J. Sim. Google Landmarks Dataset v2 - A Large-Scale Benchmark for Instance-Level Recognition and Retrieval. In *Proc. CVPR*, 2020. 3, 1
- [41] Le Xue, Manli Shu, Anas Awadalla, Jun Wang, An Yan, Senthil Purushwalkam, Honglu Zhou, Viraj Prabhu, Yutong Dai, Michael S Ryoo, et al. xgen-mm (blip-3): A family of open large multimodal models. *arXiv preprint arXiv:2408.08872*, 2024. 1, 2
- [42] Jiabo Ye, Haiyang Xu, Haowei Liu, Anwen Hu, Ming Yan, Qi Qian, Ji Zhang, Fei Huang, and Jingren Zhou. mplug-owl3: Towards long image-sequence understanding in multi-modal large language models. *arXiv preprint arXiv:2408.04840*, 2024. 2
- [43] Yuanyang Yin, Yaqi Zhao, Yajie Zhang, Ke Lin, Jiahao Wang, Xin Tao, Pengfei Wan, Di Zhang, Baoqun Yin, and Wentao Zhang. Sea: Supervised embedding alignment for token-level visual-textual integration in mllms. *arXiv preprint arXiv:2408.11813*, 2024. 1, 7
- [44] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11975–11986, 2023. 1
- [45] Pan Zhang, Xiaoyi Dong, Yuhang Zang, Yuhang Cao, Rui Qian, Lin Chen, Qipeng Guo, Haodong Duan, Bin Wang, Linke Ouyang, et al. Internlm-xcomposer-2.5: A versatile large vision language model supporting long-contextual input and output. *arXiv preprint arXiv:2407.03320*, 2024. 1, 8

Beyond Sight: Towards Cognitive Alignment in LVLm via Enriched Visual Knowledge

Supplementary Material

A. Experiment Details

Architecture. We utilize CLIP-ViT-L-14 [29] as the vision encoder, with a default resolution of 336×336 , and Meta-Llama-3-8B [1] as the language model. The low-resolution (LR) branch employs a 2-layer MLP as the adapter, while the high-resolution (HR) branch compresses visual tokens using a shared Perceiver resampler layer [2]. In the HR branch, each high-resolution image is divided into four sub-patches. The resampler processes the visual tokens from these sub-images, compressing them from 2,880 down to 128 tokens via cross-attention with query vectors. These 128 tokens are then concatenated with the 576 tokens from the low-resolution overview image and fed into the LLM.

Pretrain Datasets. We use the same dataset for LLaVA-1.5 experiments. Specifically, stage 1 uses CC595k [31] and stage 2 uses DataMix 665k [10, 11, 13, 14, 20, 23–25, 30, 33] proposed in [19].

Hyperparameters. In this work, we adopt the same set of hyperparameters as LLaVA-1.5 [19]. We show the training hyperparameters for LLaVA-1.5 experiments in Table 5. All experiments are conducted using a maximum of 8 Nvidia H800 GPUs. We set the value of $\mu_e = 7.32$ and $\mu_h = 4.38$, respectively, with a sparsification threshold $\theta = 0.5$ applied to selectively filter out lower-relevance tokens. Additionally, the temperature parameter, initialized at zero, is set as learnable to dynamically adjust throughout training.

Robustness of EECA. We evaluate the robustness of EECA across different hyperparameter settings. The model maintains comparable accuracy within specific ranges of key parameters, such as the balancing coefficient μ_e , sparsification threshold θ , and high-resolution visual tokens N_{vH} (see Figure 7), demonstrating stable performance despite small variations.

B. Datasets Construction for MGLD

This section expands Section 3 with additional details about our data preprocessing steps.

B.1. Details of google landmarks dataset v2

Overview. The Google Landmarks Dataset v2 (GLDv2) [40] is the largest benchmark for fine-grained

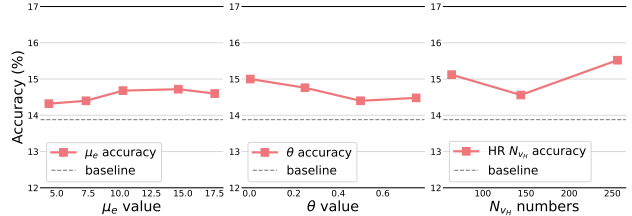


Figure 7. The robustness of the different hyperparameters.

Hyperparameter	LLaVA-1.5		EECA
	Stage 1	Stage 2	Stage 3
batch size	256	128	128
lr	2e-3	2e-5	2e-5
lr schedule decay	cosine	cosine	cosine
lr warmup ratio	0.03	0.03	0.03
weight decay	0	0	0
epoch	1	1	1*
optimizer	AdamW [21]		
DeepSpeed stage	2	2	2

Table 5. Hyperparameters for EECA training on LLaVA-1.5. By default, EECA is trained for 1 epoch, denoted by *. Unless otherwise stated, the results presented in Table 4 are based on 2 epochs of training.

instance recognition and image retrieval, comprising over 5 million images with 200,000 distinct instance labels. It sourced from Wikimedia Commons, and is characterized by real-world challenges such as imbalanced class distribution and high intra-class variability.

Data usage in this work. In this study, we leverage the GLDv2 dataset to construct a fine-tuning dataset, utilizing all data from the training set (train.csv, train_label_to_category.csv, and train_label_to_hierarchical.csv). This dataset comprises 4.1 million images spanning 203,000 landmarks.

- **train.csv:** Contains fields `id`, `url`, and `landmark_id`. Here, `id` is a 16-character string, `url` is a string representing the image’s URL, and `landmark_id` is an integer identifier for the landmark.
- **train_label_to_category.csv:** Includes `landmark_id` and `category` fields. `landmark_id` is an integer, while `category` is a Wikimedia URL linking to the class definition of the landmark.



Figure 8. The structure of the GLDv2 train set

• **train_label_to_hierarchical.csv:**

Contains fields `landmark_id`, `category`, `supercategory`, `hierarchical_label`, and `natural_or_human_made`. `Supercategory` refers to the type of landmark (e.g., `natural` or `human-made`), mined from Wikimedia. `Hierarchical_label` corresponds to the landmark’s hierarchical classification, and `natural_or_human_made` indicates whether the landmark is naturally occurring or man-made.

The structure of the GLDv2 training dataset is depicted in Figure 8. Each `hierarchical_label` encompasses multiple categories, and each category consists of a varying number of images, reflecting the diverse and hierarchical nature of the dataset. The category distribution in GLDv2 training dataset is highly imbalanced, as illustrated in ???. Approximately 57% of the categories contain at most 10 images, and 38% have 5 or fewer images. This makes the dataset diverse, covering a wide range of landmarks, from globally renowned sites to more obscure, local landmarks.

B.2. Image selection methodology

In the first stage (Figure 2), we select representative images from the GLDv2 dataset. Since each landmark name corresponds to multiple images, typically sourced from Wikipedia entries related to the landmark, we use the landmark’s simple name (e.g. “Eiffel Tower”) as the sole reference. Using CLIP-based similarity measures, we select the image that best matches this name, filtering out low-quality or ambiguous photos and ensuring a high-quality visual representation aligned with the landmark’s identity. From the top three images with the highest similarity scores, we conduct weighted sampling based on their similarity to select a unique image corresponding to each landmark. Figure 10 presents a specific example of image selection.

B.3. Prompt design for data annotation

Q-A pair. Once we have a refined set of images, we generate Question-Answer (Q-A) pairs to facilitate landmark recognition. The prompt used for image Q-A pair annota-

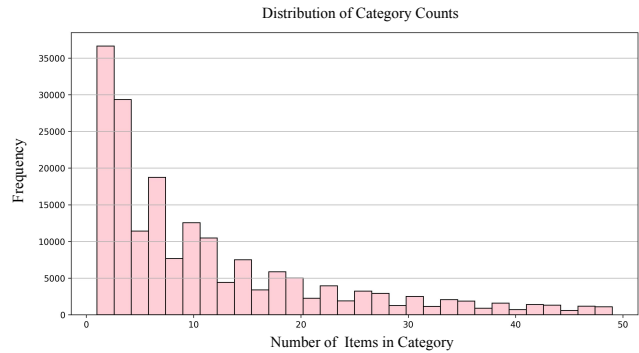


Figure 9. Frequency of the counts of images per category



Figure 10. Example for the image selection methodology (the Castle of Pardo de Cela). **Red:** The top3 images with the highest similarity scores. **Green:** The final image obtained through weighted sampling.

tion is shown in Figure 12. For each selected image, we random select one question from the questions set and add the landmark name as a reference to ensure the accuracy of the annotation. While answering these questions, the model is encouraged to provide descriptive details about the landmark, drawing on both the visual features and contextual information. This approach aims to broaden the model’s understanding of each landmark’s visual and con-

You are tasked with labeling a batch of images depicting landmarks from around the world. You will be provided with the name and description of the landmark in the image.

[Location]: {Location}
 [Description]: {Description}

Using both the image and the description provided, your task is to generate rich content that includes:

1. Entities Identification:
 Identify and list specific entities that can help recognize the landmark. These entities should include both visual features (e.g., structural elements, colors, shapes) and conceptual features understood by a language model (e.g., architectural style, historical significance, cultural context, environment) but avoid explicitly using the landmark name.
2. Given the original hierarchical label: '{Hierarchical_label}' (may be unknown and incorrect), return the most (only return one) related new hierarchical label from the following list:
 "mountain", "volcano", "lake", "waterfall", "river", "wetland", "ocean area", "beach", "cliff", "cave", "island", "tree", "botanical garden", "parks", "trail", "agricultural land", "stone", "canyon", "desert", "fjord", "glacier", "peninsula", "biosphere reserve", "well", "salt flat", "church", "temple", "castle / fort", "palace", "monastery", "tower", "memorial", "cemetery", "ruins", "pyramid", "winery", "sports venue", "road", "museum", "theatre", "library", "zoo", "shopping", "house", "square", "hotel", "restaurant", "school", "hospital", "prison", "embassy", "concert hall", "town hall", "cinema", "fountain", "post office", "bank", "gate", "bridge", "lighthouse", "harbor", "canal", "mine", "dam", "factory", "power plant", "air transportation", "rail transportation", "cable transportation", "aqueduct", "tunnel", "ship", "market", "sculpture", "artwork", "bath", "swimming pool", "amusement park", "windmill", "stairs", "observatory", "skyscraper", "conference center", "casino", "festival", "aquarium"]

Example Output Structure for Eiffel Tower:

```
{
  "entities": [
    "wrought iron lattice structure",
    "Parisian skyline",
    "Gustave Eiffel",
    "19th-century engineering",
    "tourist attraction",
    "World's Fair 1889"
  ],
  "hierarchical_label": "tower",
}
```

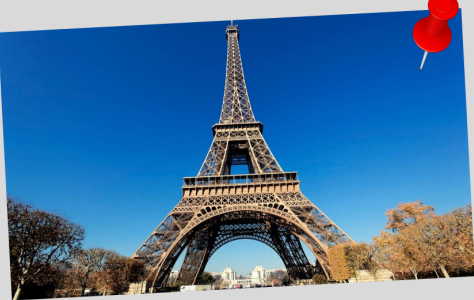


Figure 11. Multi-Granularity Data Generation Prompt. The description is the Q-A pair.

textual uniqueness, laying a foundation for aligning with VE’s cognitive framework.

Multi-granularity data annotation. Following the Q-A pair annotation for each image, we generate multi-granularity data using the multi-granularity data generation prompt shown in Figure 11. In the original dataset, some hierarchical labels were already provided. To enhance accuracy, we employed GPT-4o to refine and expand the annotations, using the original labels as a reference. This process provides us with entities and an updated hierarchical label.

B.4. MGLD overview.

We structure the data as illustrated in Figure 13. Each image is associated with a Q-A conversation, its landmark name, entities that capture both visual and conceptual features, and a hierarchical label representing its general category.

From the final dataset of approximately 203k samples, we set aside 5k samples as the test dataset.

C. Evaluation Detail

C.1. Prompt design

We provide GPT4o with 5 inference runs of the VLLM and the ground-truth answer. The GPT4o is asked to evaluate the overall recognition of the landmark based on all 5 re-

Questions

1. Where was this photo taken?
2. Identify the location where this photo was taken.
3. What is the location shown in the image?
4. Tell me where this photo was taken.
5. Where might this photo have been taken?
6. What place is depicted in this image?

Prompt

PROMPT_QA = "{QUESTION}
 [Reference Location]: {LOCATION}
 [Note: The reference location may be incorrect, so you need to rely on your own ability to answer.]
 [Your response must adhere to the following requirements]:
 1. The response must be in English;
 2. Besides identifying the location where the photo was taken, you should also describe the photo and share some knowledge related to the location;
 3. Do not mention the reference location or this note in your response."

Figure 12. Q-A pair Prompt.

sponses together, and classify the level of recognition into one of the four levels: Strongly Known, Known, Partially Known, or Unknown. The classification criteria is clearly defined in the prompt (see Figure 14).

C.2. Rating criteria examples

This section provides an example (“Kinderdijk Windmills“) evaluated by GPT-4o, where the answer across different models is assessed at four different levels—*Strongly*

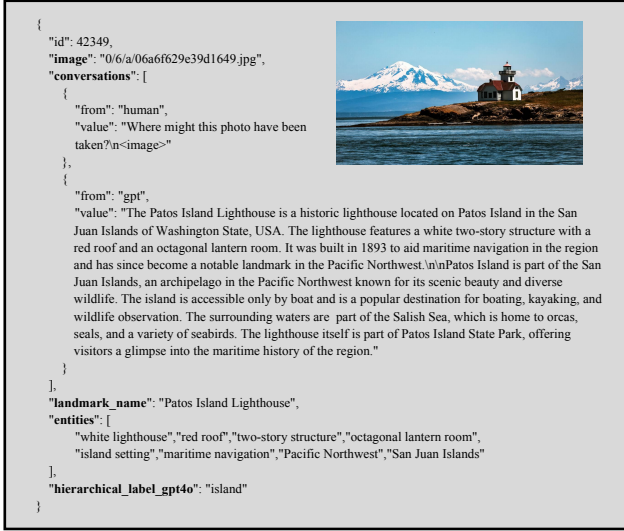


Figure 13. Example of the MGLD datasets. The conversation is the Q-A pair, and the hierarchical_label_gpt4o is the new hierarchical label annotated by GPT-4o.

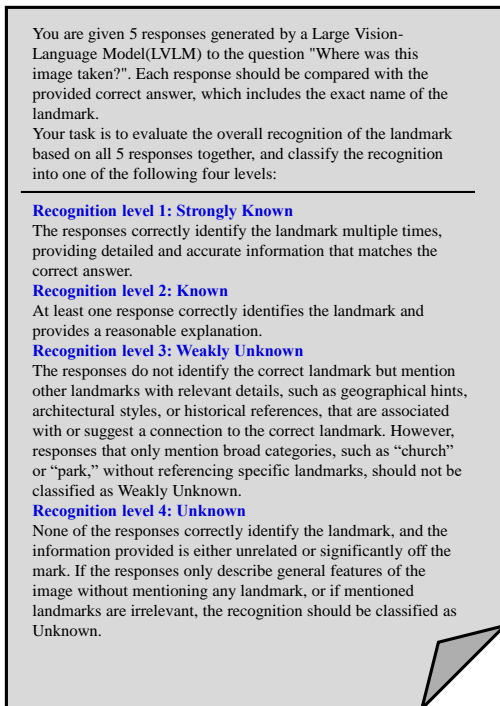


Figure 14. Evaluation Prompt

Known, Known, Partially Known, and Unknown. Figure 15 shows the specific answers from each model, followed by GPT-4o’s explanations for the corresponding evaluations.

- Explanation of strongly known : **The responses correctly identify the Kinderdijk Windmills multiple times**, specifically in Responses 2 and 3, providing de-

tailed and accurate information that matches the correct answer. These responses include precise details about the location, historical significance, and characteristics of the Kinderdijk Windmills, demonstrating a strong recognition of the landmark.

- Explanation of known: Response 1 **correctly identifies the Kinderdijk Windmills** and provides detailed and accurate information about their location, history, and significance. This response alone is sufficient to classify the recognition as ‘Known’. Although Responses 2 and 3 incorrectly identify the location as Zaanse Schans, they do **provide relevant details about Dutch windmills**, which are related to the correct landmark. Responses 4 and 5 describe general features of Dutch windmills and mention regions like Kinderdijk, but do not specifically identify the landmark. Overall, the presence of one accurate identification with detailed information justifies the ‘Known’ classification.
- Explanation of weakly unknown: None of the responses correctly **identify the ‘Kinderdijk Windmills’ specifically**. However, several responses **mention relevant Dutch windmills and regions associated with windmills**, such as Kinderdijk and Zaanse Schans, providing geographical hints and cultural context related to Dutch windmills. This suggests a connection to the correct landmark, even though it is not explicitly named.
- Explanation of unknown: **None of the responses correctly identify the ‘Kinderdijk Windmills’**. The responses mention other windmills and locations such as ‘Weerschipmolen’, ‘Zaanse Schans’, and ‘Huis op de Hoop’, but these are not related to the correct landmark. The information provided is either unrelated or significantly off the mark, as none of the responses provide any specific details or hints that connect to the Kinderdijk Windmills.

D. Statistical Analysis for Partial Knowledge Parts

D.1. Detail results for ablation study

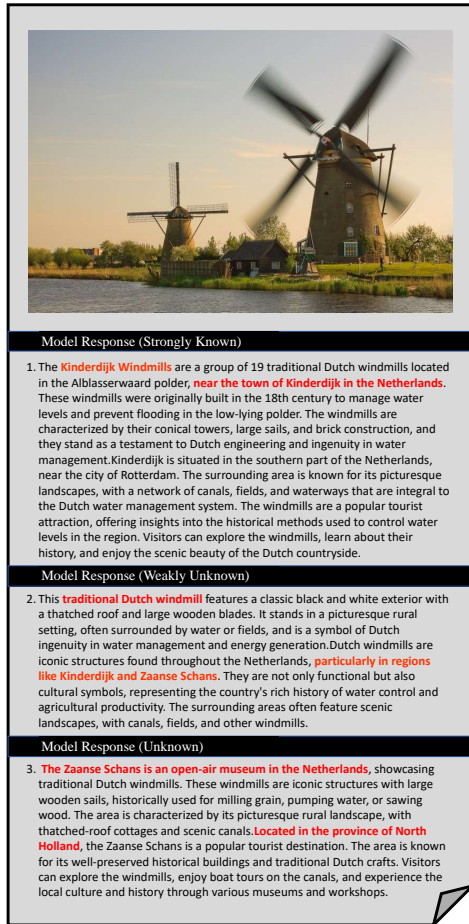
We reports the results of the ablation study in terms of proportions relative to the entire dataset (see Table 3). Here, in Table 6, we provide the corresponding absolute counts for each knowledge level (Strongly Known, Known, Weakly Unknown, Unknown).

D.2. Detailed results for generalizability experiment

We summarizes the results of the generalizability experiments using accuracy (see Table 4). To complement this, Table 7 presents the absolute counts for each knowledge level (Strongly Known, Known, Weakly UnKnown, and Unknown).

Method	Strongly Known	Known	Weakly Unknown	Unknown
<i>Baseline</i>	103	114	145	2138
+ <i>HSS-50k</i>	187	161	145	2007
+ <i>HR-Branch</i>	198	149	152	2001
+ \mathcal{L}_e	212	148	163	1977
+ \mathcal{L}_h	213	175	159	1953

Table 6. Detailed ablation study results. Absolute counts for each level of knowledge.



Method	Strongly Known	Known	Weakly Unknown	Unknown
Baseline	103	114	145	2138
For HDS-25k				
+Data	182	142	151	2025
+HR Branch	218	122	185	1975
+ \mathcal{L}_h	233	127	175	1965
+ \mathcal{L}_e	229	121	199	1951
For HSS-25k				
+Data	169	135	201	1995
+HR Branch	202	144	177	1977
+ \mathcal{L}_h	205	141	193	1961
+ \mathcal{L}_e	226	134	153	1987
For LCS-25k				
+Data	179	88	253	1980
+HR Branch	203	99	243	1955
+ \mathcal{L}_h	192	116	254	1938
+ \mathcal{L}_e	198	110	266	1926

Table 7. Generalizability experiment results. Absolute counts for each knowledge level across datasets and methods.

Figure 15. Model responses from different recognition level. There is a significant gap in recognition ability among “Strongly Known,” “Weakly Unknown,” and “Unknown”. The difference between “Strongly Known” and “Known” is the number of times the model correctly identifies the landmark.