

Seed1.5-VL Technical Report

ByteDance Seed

See [Contributions and Acknowledgments](#) section for a full author list.

Abstract

We present Seed1.5-VL, a vision-language foundation model designed to advance general-purpose multimodal understanding and reasoning. Seed1.5-VL is composed with a 532M-parameter vision encoder and a Mixture-of-Experts (MoE) LLM of 20B active parameters. Despite its relatively compact architecture, it delivers strong performance across a wide spectrum of public VLM benchmarks and internal evaluation suites, achieving the state-of-the-art performance on 38 out of 60 public benchmarks. Moreover, in agent-centric tasks such as GUI control and gameplay, Seed1.5-VL outperforms leading multimodal systems, including OpenAI CUA and Claude 3.7. Beyond visual and video understanding, it also demonstrates strong reasoning abilities, making it particularly effective for multimodal reasoning challenges such as visual puzzles. We believe these capabilities will empower broader applications across diverse tasks. In this report, we mainly provide a comprehensive review of our experiences in building Seed1.5-VL across model design, data construction, and training at various stages, hoping that this report can inspire further research. Seed1.5-VL is now accessible on [Volcano Engine^a](#).

Date: May 11, 2025

Correspondence: shiguang.sg@bytedance.com

^aModel ID: doubao-1-5-thinking-vision-pro-250428

Contents

1	Introduction	4
2	Architecture	5
2.1	Vision Encoder	5
2.1.1	Architecture	6
2.1.2	ViT Pre-training Stage	6
2.2	Video Encoding	7
3	Pre-training	8
3.1	Pre-training Data	8
3.1.1	Generic Image-Text Pairs & Knowledge Data	8
3.1.2	Optical Character Recognition (OCR)	9
3.1.3	Visual Grounding & Counting	10
3.1.4	3D Spatial Understanding	11
3.1.5	Video	11
3.1.6	Science, Technology, Engineering, and Mathematics (STEM)	12
3.1.7	Graphical User Interface (GUI)	12
3.2	Training Recipe	13
3.3	Scaling Laws	14
4	Post-training	15
4.1	Supervised Fine-tuning	16
4.1.1	SFT Data Construction	16
4.1.2	Training Recipe	16
4.2	Reinforcement Learning from Human Feedback	17
4.2.1	Preference Data	17
4.2.2	VLM as a Reward Model	17
4.2.3	Data Curation for Reinforcement Learning	18
4.3	Reinforcement Learning with Verifiable Rewards	18
4.3.1	Visual STEM	18
4.3.2	Visual Perception and Reasoning	18
4.4	Hybrid Reinforcement Learning	19
4.5	Iterative Update by Rejection Sampling Fine-tuning	20
5	Training Infrastructure	21
5.1	Large-Scale Pre-training	21
5.1.1	Hybrid Parallelism	21
5.1.2	Workload Balancing	21
5.1.3	Parallelism-Aware Data Loading	21
5.1.4	Fault Tolerance	21
5.2	Post-Training Framework	22
6	Evaluation	22
6.1	Public Benchmarks	22
6.1.1	Vision Encoder as a Zero-shot Classifier	22
6.1.2	Vision Task Evaluation	23
6.1.3	Video Task Evaluation	25
6.2	Multimodal Agent	25
6.3	Internal Benchmarks	28
6.3.1	Motivation and Design Principles	28
6.3.2	Comparison with State-of-the-arts	29
6.3.3	Out-of-distribution Generalization	30

6.4	Limitations	30
7	Conclusion and Next Steps	32
8	Contributions and Acknowledgments	44
A	Qualitative examples	47
A.1	Reasoning Cases: Visual Reasoning	48
A.2	Reasoning Cases: Geolocation Prediction	49
A.3	Visual Reasoning: Solving Rebus Puzzles	50
A.4	Visual Reasoning: Emoji Quiz	51
A.5	Visual Reasoning: Word Game I	52
A.6	Visual Reasoning: Word Game II	53
A.7	Visual Reasoning: Visual Pattern Recognition	54
A.8	Visual Puzzles: Find the Differences	55
A.9	Geometry	56
A.10	Counting in a complex scene	57
A.11	Spatial Understanding: Depth Sorting	58
A.12	Video Temporal Grounding	58
A.13	OCR Parsing and Document Understanding	59
A.14	Multilingual OCR Parsing	60
A.15	Generate Code for a Diagram of Novel Format	61
A.16	Image-conditioned Creative Writing	62
A.17	Failure Cases: 3D Spatial Imagination	63
A.18	Failure Cases: Hallucination (Knowledge Prior)	64
A.19	Failure Cases: Combinatorial Search I	65
A.20	Failure Cases: Combinatorial Search II	66
B	Evaluation Details	67
B.1	Internal Benchmark Structure	67
B.2	Comprehensive Comparisons on internal benchmarks	69
B.3	Capabilities and Benchmark Tasks	70
B.4	Evaluation Prompts	71

1 Introduction

Vision-language models (VLMs) have emerged as a foundational paradigm for enabling general-purpose AI to perceive, reason, and act in open-ended virtual and physical environments. By aligning visual and textual modalities within a unified model, VLMs have rapidly advanced research frontiers in areas, such as multimodal reasoning [96, 129, 141], image editing [35, 97], GUI agents [5, 98, 105], autonomous driving [103, 131, 157], and robotics [31, 55, 63], while also powering real-world applications across education, healthcare, chatbots, and wearable devices.

However, despite substantial progress, current VLMs still fall short of human-level generality, particularly in tasks requiring 3D spatial understanding, object counting, imaginative visual inference, and interactive game play. These limitations highlight the inherent challenges in VLM development. Unlike large language models (LLMs), which benefit from abundant, high-quality textual corpora that capture a wide spectrum of human knowledge, VLMs lack access to equally rich and diverse vision-language annotations, especially for concepts grounded in low-level perceptual phenomena. Moreover, the heterogeneous nature of multimodal data introduces additional complexity in both training and inference, complicating data pipeline design, parallel training strategies, and evaluation protocols.

In this report, we share the efforts during the development of Seed1.5-VL, our latest multimodal foundation model for vision-language understanding. To address the scarcity of high-quality annotations, we developed a suite of diversified data synthesis pipelines targeting key capabilities, including optical character recognition (OCR), visual grounding, counting, video understanding, and long-tail knowledge during pre-training, as well as visual puzzles and games during post-training. Seed1.5-VL is pre-trained on trillions of multimodal tokens spanning diverse modalities, *i.e.*, images, videos, text, and human-computer interaction data, to acquire broad visual knowledge and master core visual competencies. We also share the scaling behavior in the pre-training stage. In the post-training phase, we incorporate both human feedback and verifiable reward signals to further strengthen its general reasoning abilities.

We also address the challenge of efficiently training large-scale multimodal models with asymmetrical architecture, especially the imbalance between the vision encoder and the language model. Our contributions include (1) a novel *hybrid parallelism* scheme optimized for this asymmetry and (2) a *vision token redistribution strategy* to balance GPU workloads. In addition, we implement a customized data loader that minimizes I/O bottlenecks under 3D parallelism. These innovations, combined with standard system-level optimizations (*e.g.*, kernel fusion, selective activation checkpointing, offloading), collectively enhance overall training throughput.

To establish a comprehensive understanding of the current landscape of VLM capabilities, thereby informing future research directions towards model improvements, we evaluate Seed1.5-VL on an extensive suite of public and internal benchmarks, covering a wide range of tasks including visual reasoning, grounding, counting, video understanding, and computer usage. Specifically, we report results on 60 public benchmarks, where Seed1.5-VL achieves state-of-the-art performance on 38 of them, including 21 out of 34 in vision-language benchmarks, 14 out of 19 in the video benchmarks, and 3 out of 7 in GUI agent tasks. Beyond benchmark performance, we also deploy Seed1.5-VL within an internal chatbot system to monitor its real-world and out-of-distribution (OOD) performance in dynamic, interactive environments.

Despite its strong capabilities, Seed1.5-VL maintains a compact and efficient architecture, featuring a 532-million-parameter vision encoder and a language model with 20 billion active parameters. This streamlined design reduces inference costs and computational demands, making the model well-suited for interactive applications. The efficiency of Seed1.5-VL enhances accessibility for a broader user base via API services and contributes to a smoother user experience within the Doubao chatbot. Access to Seed1.5-VL will soon be available on the Volcano Engine API platform¹.

The remainder of this report is organized as follows. We begin by presenting an overview of the model architecture and detailing the image and video encoding methods (section 2). Section 3 describes the data curation strategies and the pre-training procedure, including initial findings on multimodal model scaling laws and metric prediction—a relatively underexplored area. Section 4 details the data and techniques

¹<https://www.volcengine.com>

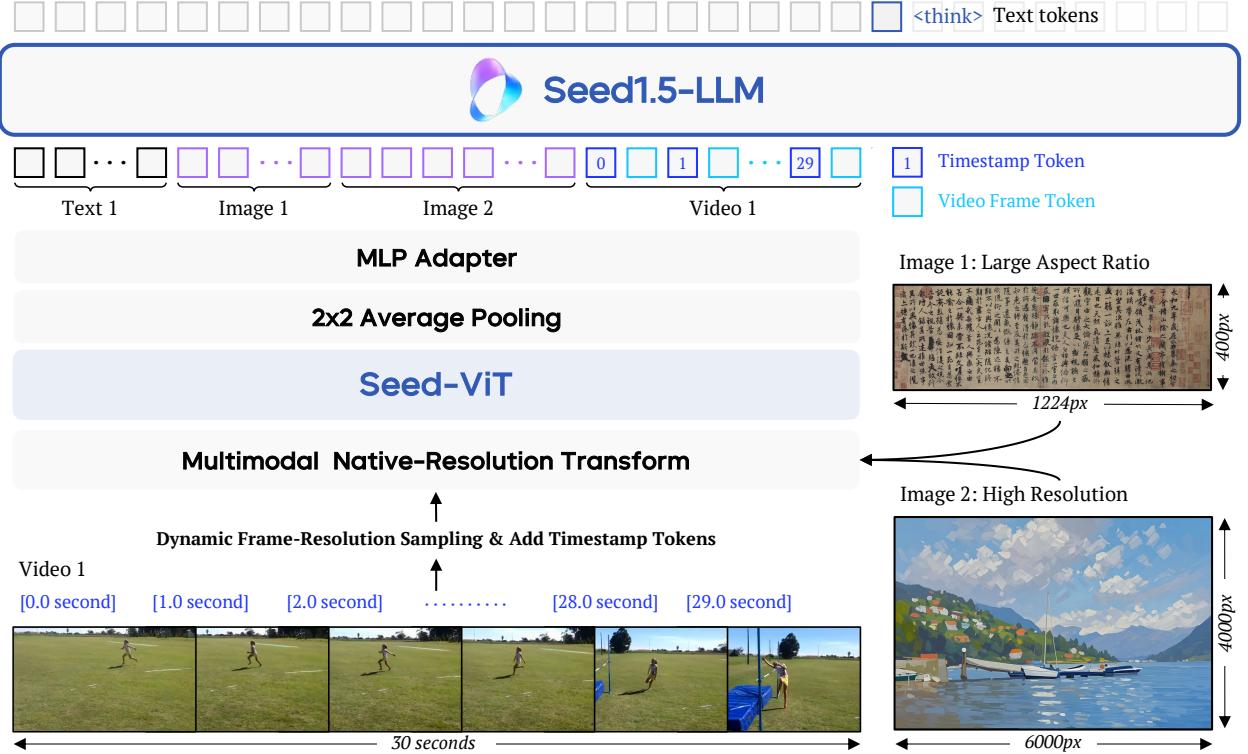


Figure 1 The architecture of Seed1.5-VL. The proposed Seed1.5-VL comprises three main components: (1) SeedViT to encode images and videos, (2) an MLP adapter to project visual features into multimodal tokens, and (3) a Large Language Model to process multimodal inputs. Seed1.5-VL accepts images at various resolutions and processes them using a native-resolution transform to preserve maximum image detail. For video inputs, we propose the dynamic frame-resolution sampling strategy, which dynamically adjusts the sampling frame rate and resolution. Additionally, a timestamp token is added before each frame to enhance the model’s temporal awareness.

employed during the post-training phase to enhance alignment with human preferences and improve reasoning capabilities. Section 5 elaborates on the necessary infrastructure innovations developed to enable scalable pre-training and post-training. Finally, section 6 presents comprehensive evaluation results on public benchmarks, showcases model capabilities via qualitative examples, discusses limitations of current multimodal models, and proposes directions for future research.

2 Architecture

The architecture of Seed1.5-VL consists of three components: a vision encoder, an MLP adapter, and a large language model (LLM). The vision encoder natively supports dynamic image resolutions and employs 2D RoPE [126] for positional encoding, enabling flexible adaptation to images of arbitrary dimensions. To enhance computational efficiency, the architecture applies average pooling over adjacent 2×2 feature patches; a two-layer MLP subsequently processes these pooled features before being input to the LLM. Encoder-free architectures [1, 23, 127] are not considered, as the vision encoder provides efficient image compression, enabling high-resolution image representation with fewer tokens. The overall architecture is shown in figure 1.

2.1 Vision Encoder

Many contemporary Vision-Language Models (VLMs) [2, 5, 7, 16, 37, 54, 71, 78, 104, 128, 141] commonly integrate pre-trained vision encoders designed for a fixed input resolution, typically square images. While this approach simplifies model architecture, it can inadvertently discard fine-grained visual information when processing high-resolution images, videos, or handling tasks requiring intricate detail such as OCR.

Recent efforts, such as those in Qwen2-VL [141] and InternVL-2.5 [16], have explored fine-tuning pre-trained vision encoders to accommodate dynamic-resolution inputs, offering a partial alleviation of this limitation. Nevertheless, these methods still largely depend on adapting existing fixed-resolution architectures and necessitate adjustments to position encodings (e.g., transitioning from 1D flatten position embedding to 2D RoPE [16, 141] or interpolation of 1D position embeddings to various shapes [99, 135]), which may not fully retain visual details and precision post-adaptation. Furthermore, we incorporate video data into the pretraining phase to enable the model to learn not only spatial features from images but also spatial-temporal dynamics, thereby enhancing its capacity to process dynamic scenes and complex visual content.

Addressing the challenges posed by fixed-resolution processing, we developed Seed-ViT, a vision encoder specifically designed for native-resolution feature extraction. Based on the well-established Vision Transformer (ViT) architecture [26], Seed-ViT consists of 532 million parameters. It demonstrates strong capabilities in general visual perception across diverse domains. Notably, on zero-shot classification benchmarks, Seed-ViT attains performance comparable to models with substantially more parameters, such as InternVL-C (6 billion parameters), highlighting its efficiency. Further architectural details and our pretraining approach for Seed-ViT are provided in sections 2.1.1 and 2.1.2, respectively.

2.1.1 Architecture

The architectural hyper-parameters of Seed-ViT can be found in table 1.

Patch size	Pos embed	Head dim	Num heads	Embed dim	MLP ratio	Depth
14	2D RoPE	64	20	1280	4.0	27

Table 1 The architectural hyperparameters of Seed-ViT.

Our vision encoder is designed to accommodate input images of varying dimensions. Initially, input images undergo a pre-processing step involving bilinear interpolation to adjust their resolutions to the nearest multiple of 28×28 pixels. Subsequently, each image is segmented into a sequence of non-overlapping patches, each of 14×14 pixels. Following the approach outlined in NaViT [20], we concatenate patch sequences from multiple input images into a unified sequence. These raw patch sequences are then projected into tokens in the embedding space via a linear patch embedding layer, which are then fed into the transformer blocks. To ensure that tokens belonging to one image do not attend to tokens from other images within the batched sequence, we employ appropriate attention masks during the self-attention computations within the transformer blocks. Finally, a 2×2 average pooling operation is applied to the output patch embeddings before they are passed to the subsequent MLP adaptor and the LLM, as described above.

2.1.2 ViT Pre-training Stage

Categories	Unlabeled image	Image-text pairs	Video-audio-text tuples
Training samples	2.2B	4.8B	65M
Token percentages	4.0%	91.2%	4.8%
Batch sizes	55,296	32,768	1,024
LR warm up steps	1,692	2,000	12,800
Maximum LR	7.06×10^{-3}	1.0×10^{-4}	5.0×10^{-5}
Minimum LR	1.05×10^{-5}	1.2×10^{-6}	2.02×10^{-7}

Table 2 Training setup and hyperparameters used in the three ViT pre-training stages.

Our vision transformer, Seed-ViT, undergoes a dedicated pre-training pipeline before integration with the LLM. Guided by empirical evidence, we establish three key guidelines for our pre-training methodology:

- **Better Training Efficiency with ViT-pretraining.** Most successful VLMs [7, 16, 141] follow the setup of having a vision encoder (e.g., CLIP or SigLIP [171]) and a few work [1, 24] have attempted to remove vision encoder entirely and directly pass image patches in decoder-only LLMs but with mixed results. Beyer et al. [10] also concluded that encoder-free VLMs may be a promising future direction but still suffer in training efficiency.
- **Early Integration of Native-Resolution Modeling.** We prioritize the early introduction of native-resolution modeling within the pre-training pipeline. The architecture of Seed-ViT is maintained consistently throughout both the ViT pre-training and VLM stages. This ensures the prevention of performance degradation stemming from architectural modifications and eliminates the need for extensive fine-tuning to compensate for such discrepancies.
- **Comprehensive Data Utilization.** The pre-training stage leverages the full spectrum of data intended for VLM training, encompassing unlabeled images, image-text pairs, and videos accompanied by visual and audio captions.

Based on the above guidelines, the ViT pre-training pipeline is divided into three stages: (i) Masked Image Modeling (MIM) [145] with 2D RoPE, (ii) Native-Resolution Contrastive Learning, and (iii) Omni-modal Pre-training. Below, we provide more details of each stage.

MIM with 2D RoPE. In the first stage, our goal is to enhance the visual perception ability on visual geometry and structure awareness by MIM. We leverage the EVA02-CLIP-E [29] as the teacher model, and the student model is randomly initialized following the architecture defined in [table 1](#). During training, we randomly mask out 75% image patches and the corresponding RoPE embeddings and use the CLIP [107] features produced by the teacher as reconstruction targets. This process is optimized by a simple cosine similarity loss between masked-out patches in the student’s and teacher’s outputs. We find that the discrepancy in visual position embeddings between student and teacher models does not harm the performance, as the teacher employs learnable positional embeddings while the student uses 2D RoPE. Instead, 2D RoPE empowers the student with robust native dynamic-resolution recognition. As we scale up this MIM process, the abilities of VLMs on chart/document understanding and OCR are significantly improved.

Native-Resolution Contrastive Learning. In the contrastive learning stage, the vision encoder is initialized with our MIM-trained student model, while the text encoder is initialized using the text encoder from EVA-02-CLIP-E. For each given image-text pair, we aggregate the extracted patch features from the vision encoder into a single 1280-dimensional image embedding using attention pooling. Alignment between the image and text embeddings is then achieved by jointly optimizing the SigLIP loss [171] and the SuperClass loss [52].

Omni-modal Pre-training. This stage adopts the MiCo framework [174], constructing aligned tuples consisting of video frames, audio, visual captions, and audio captions from video data. The ViT encodes both video frames and audio, while a separate text encoder processes captions. Through alignment of these embeddings, the ViT learns unified omni-modal representations. Despite consuming only 4.8% of the token budget allocated for the entire ViT pre-training process, this stage significantly enhances the ViT’s performance on image and video understanding tasks.

[Table 2](#) summarizes the training setup and hyperparameters used in each stage.

2.2 Video Encoding

Effectively encoding video, beyond static image representation, remains a core challenge. A model’s ability to interpret temporal sequences, adapt to varying frame rates, and perceive absolute time is critical for understanding dynamic visual content. Seed1.5-VL addresses these challenges by introducing **Dynamic Frame-Resolution Sampling**, a novel strategy that jointly optimizes sampling across both the temporal (frame) and spatial (resolution) dimensions to balance semantic richness and computational efficiency.

Under this Dynamic Frame-Resolution Sampling strategy, videos are processed as sequences of image frames. The temporal dimension is managed through dynamic frame sampling. Instead of a uniform rate, Seed1.5-VL adjusts the frame sampling frequency based on content complexity and task requirements. The default

sampling rate is set at 1 frame per second (FPS), suitable for capturing a general understanding of video content. For tasks [73, 139] requiring detailed temporal information, the frame sampling rate is increased to 2 FPS. For tasks such as video counting [27] or motion tracking [48], the rate is increased to 5 FPS. To explicitly ground each frame within the video’s timeline, we prepend timestamp tokens (i.e., [1.5 second]) to each frame. This explicit timing annotation substantially enhances the model’s temporal awareness and enables it to handle variable frame rates common in real-world scenarios effectively.

Considering computational constraints inherent in processing long video sequences, the spatial dimension of the sampling is governed by dynamically adjusting the resolution allocated to each selected frame, managed within a maximum budget of 81,920 tokens per video. The model dynamically adjusts spatial resolutions, assigning tokens per frame through a hierarchical allocation system offering six predefined levels: {640, 512, 384, 256, 160, 128}. This allows for a flexible trade-off, i.e., using higher resolution for fewer frames or lower resolution to accommodate more frames from longer videos. In cases where a video is exceptionally long and exceeds the maximum encoding length even when using the lowest token allocation (128 tokens per frame), a fallback mechanism is triggered. The model then reduces the total frame count through uniform sampling across the video. While this reduces temporal density, it ensures that the entire video is represented, balancing processing efficiency with the preservation of significant temporal information.

This flexible strategy allows Seed1.5-VL to efficiently and accurately process varying video lengths and frame rates, maintaining essential temporal details crucial for diverse video understanding tasks.

3 Pre-training

This section describes the data curation process (section 3.1) and training recipe (section 3.2) used in the pre-training stage of Seed1.5-VL. In section 3.3, we present the scaling behavior of our model.

3.1 Pre-training Data

The Seed1.5-VL pre-training corpus contains 3 trillion diverse, high-quality source tokens. This data is categorized based on target capabilities, with the curation process for each category detailed in the following subsections.

3.1.1 Generic Image-Text Pairs & Knowledge Data

Web-sourced image-text pair data, including alt text, image captions, and surrounding text, is available at an unprecedented scale (billions of instances) and exhibits high diversity in both visual and textual concepts. However, this data is inherently noisy (e.g., irrelevant or inaccurate text) and often exhibits class imbalance.

To mitigate these challenges, we first employ a series of filtering techniques, including image-text similarity scoring (e.g., CLIP-score thresholding), image-based criteria (e.g., removal of undersized images or those with extreme aspect ratios), text-based criteria (e.g., filtering of excessively short or long text), deduplication strategies (e.g., exact and near-duplicate image removal), and URL/domain-based filtering.

Furthermore, the distribution of visual concepts within the raw image-text pairs adheres to a long-tail pattern. To empirically test this observation, we conduct a sandbox experiment using Biotrove [159], a large-scale dataset for species classification containing 161.9 million images spanning 366,600 species. We train a 1.1 billion-active-parameter variant of our VLM using three distinct data distributions:

- **Random-46M.** 46 million samples randomly selected from the training set.
- **Max1k-46M.** 46 million samples selected with a maximum of 1,000 samples per species, ensuring inclusion of rare species.
- **Max100-15M.** 15 million samples with a maximum of 100 samples per species, providing greater relative exposure to rare species.

We evaluate the models on two specially filtered test sets derived from the original dataset: Balanced10k (sampled from BioTrove-Balanced representing common species) and Rare2k (sampled from BioTrove-Unseen

representing rare species). Our experiment shown in [table 3](#) indicates that the Random-46M configuration performs poorly on rare species recognition. In contrast, limiting the maximum samples per common species (Max1k-46M) significantly improves performance on rare species. Further restricting common species' representation (Max100-15M) enhances memorization of rare species but adversely affects common species recognition. Thus, effectively capturing visual knowledge requires maintaining diverse examples of common visual concepts while ensuring sufficient training iterations for rare visual concepts.

	Training tokens	Balanced10k	Rare2k	Average
Random-46M (1 epoch)	12B	78.92	10.46	44.69
Max1k-46M (1 epoch)	12B	79.17	44.85	62.01
Max100-15M (3 epochs)	12B	60.31	89.41	74.86

Table 3 Performance comparison on Balanced10k and Rare2k under three training data distributions, Random-46M, Max1k-46M, and Max100-15M. Evaluation was conducted using an open-ended Question Answering (QA) task, with responses automatically scored by a LLM judge. All models were trained with a fixed budget of 12 billion tokens.

To address the imbalance between common and rare visual knowledge acquisition from image-alt-text pairs, we propose a targeted pre-processing framework. Initially, this framework utilizes a precursor version of our VLM to automatically annotate the data with pertinent semantic domains (e.g., landmarks, food, commodities, biology) and associated named entities (e.g., product brands, species names). Named entities exhibiting low corpus frequency are identified as instances of rare visual knowledge. To mitigate data sparsity, we identify domains whose representation constitutes less than 50% of the average domain frequency. Alt-texts corresponding to these underrepresented domains are subsequently duplicated. By merging this augmented subset, enriched with samples from less frequent domains, back into the original corpus, we achieve a more balanced distribution of visual concepts. This re-balancing is designed to enhance the visual knowledge learning component, crucial to our pre-training methodology.

3.1.2 Optical Character Recognition (OCR)

To enhance the Optical Character Recognition (OCR) capabilities of the VLM, particularly for multilingual text, special symbols, and the analysis of structurally complex documents, as shown in [figure 2](#), we adopt large volumes of both annotated and synthetic data to train Seed1.5-VL.

We build an in-house OCR training dataset containing over 1 billion samples, covering documents, scene text, tables, charts, and flowcharts. For document data, we collected a large volume of pages from various sources and applied our internal tools to extract content and layout information. Furthermore, we curated a diverse set of fonts, including artistic, handwritten, and non-Latin scripts, and subsequently synthesized over 200 million text-intensive images utilizing tools such as SynthDog [62] and LaTeX (see [figure 2\(a\)](#) for an example). To improve the model's robustness in understanding textual content within images, we apply various data augmentation techniques to the synthetic data, including blurring, the addition of moiré patterns, and image distortion. [Figure 2\(c\)](#) illustrates an example of a document image after applying distortion-based augmentation.

Our chart dataset combines existing open-source datasets (e.g., FigureQA [58]) with newly generated synthetic data. Synthetic charts were generated using both conventional tools (ECharts [70], Matplotlib [53]) and a novel LLM-based pipeline. In our pipeline, an LLM generates textual chart components (titles, legends, etc.), which are then transformed by an LLM into LaTeX or Python code for rendering ([figure 2\(b\)](#)). Chart images were obtained via execution of this code. This multi-pronged approach resulted in a large-scale dataset exceeding 100 million chart examples.

For table data, we extract text in HTML, LaTeX, and Markdown formats from various sources, including web page HTML, GitHub README files, and LaTeX files from arXiv. Using this text, we render over 50 million table images, creating a comprehensive dataset for table parsing. This dataset enables our model to efficiently convert tables into formats such as HTML, LaTeX, and Markdown.

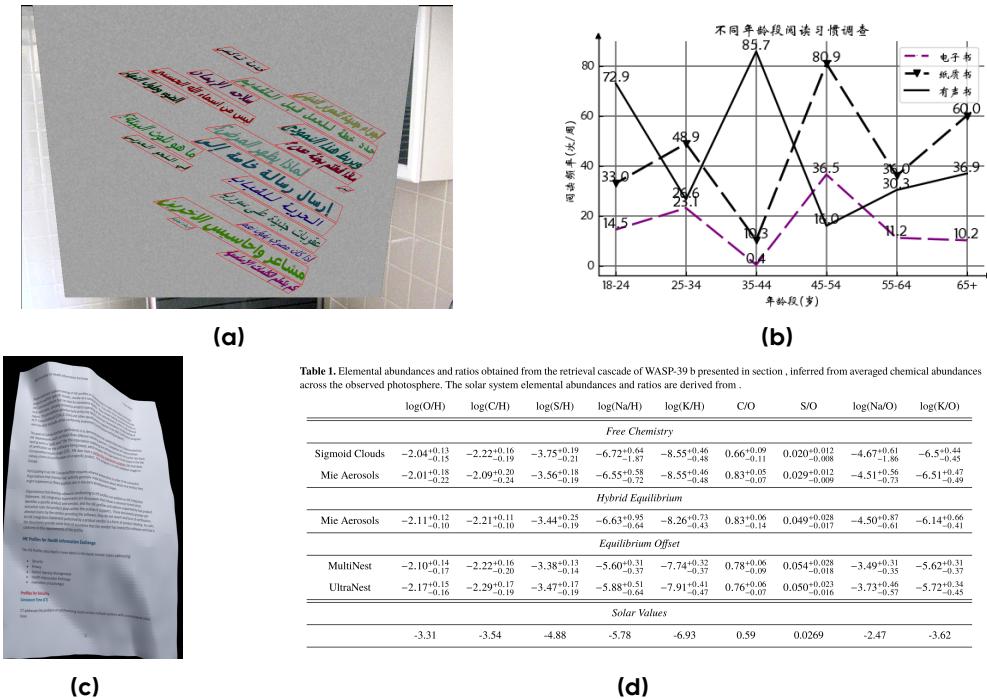


Figure 2 (a) An image generated by SynthDog and the corresponding textual annotations are organized in the following format: `<text>...</text><polygon>...</polygon>`; (b) The synthesized chart data includes two types of annotations: chart-to-text parsing and QA pairs; (c) The original document image undergoes transformations to simulate real-world distortions, such as perspective shifts, bends, and wrinkles. These augmentations enhance the model’s robustness and improve its ability to recognize texts under diverse and challenging conditions; (d) An example of a QA pair generated for the above synthesized table image: *Question: What is the value of $\log(\text{C/H})$ for Sigmoid Clouds? Give analytical steps.* *Answer: We look for the row labeled “Sigmoid Clouds” and the column labeled “ $\log(\text{C/H})$ ”. The value in that cell is $-2.22^{+0.16}_{-0.19}$.*

To further enhance the model’s comprehension of textual content within images, we constructed a visual question answering (VQA) dataset to complement the structured image-text representations. Specifically, we employed a previous version of our VLM to generate question-answer pairs by conditioning on OCR outputs, chart content, table text, and the images themselves, utilizing a few-shot prompting approach. Figure 2(d) gives an example of an input table image and the corresponding generated QA pair. Subsequently, we applied an internal LLM to filter the generated question-answer pairs, removing instances exhibiting low semantic relevance between the question and the answer. Our experiments indicate that the inclusion of this VQA dataset significantly improved the model’s ability to understand textual information present in images.

3.1.3 Visual Grounding & Counting

Object grounding, a fundamental capability for multimodal models, involves interpreting user instructions to identify and locate specific object regions within images. In this work, we employ two primary grounding representations for Seed1.5-VL: bounding boxes and center points. Building upon this localization foundation, we extend Seed1.5-VL’s capabilities to include object counting. Accordingly, our training strategy primarily utilizes three data types: bounding box annotations, point annotations, and counting data.

Bounding Box Data. Firstly, we adopt widely-used open-source datasets for generic object grounding, including Objects365 [118], OpenImages [66], and RefCOCO/+g [60, 92, 164]. Rather than directly incorporating those datasets for training, we filter low-quality samples of the open-source datasets and construct diverse grounding tasks. Specifically, we render all object bounding boxes for each category onto the images and adopt the previous version of our VLM to perform data inspection, which allows us to filter out samples with incorrect annotations, missing labels, or redundant annotations. Furthermore, we use these open-source

datasets to construct diverse multi-task training data, including: (1) generic 2D grounding, (2) question answering about spatial relationships, and (3) question answering with visual prompts, which results in about 48 million samples and 41 billion tokens. Considering the limitations in the diversity of open-source grounding datasets in terms of both data domains and categories, we develop an efficient automatic annotation pipeline for generic multi-object grounding with large-scale image-text pairs. Specifically, we follow previous work [17] and extract noun phrases and entities from captions, and then adopt Grounding DINO [14, 80] to annotate diverse open-vocabulary objects in web images. We filter out low-quality annotations with CLIP [106] and heuristic metrics, e.g., non-maximum suppression. The automatic annotation pipeline brings about 200 million samples and 200 billion tokens.

Point Data. Initially, we utilized the public data provided by PixMo-Points [21]. Recognizing limitations in the diversity and quantity of the available PixMo data, we developed a dedicated pipeline for generating additional pointing data. This pipeline employs Molmo [21] and CountGD [3] to annotate the center points of objects within a large collection of web images. Notably, CountGD proved particularly effective in annotating objects in dense image scenarios. Following annotation, low-quality data samples were filtered out, resulting in a final dataset comprising approximately 170 million instructions and 110 billion tokens.

Counting Data. We further sample from the aforementioned bounding box and point data to construct a counting dataset, containing approximately 8 million samples and 13 billion tokens. Specifically, we developed two variants: box-based counting and point-based counting, following a two-stage pipeline of 1) detection or pointing, then 2) generating counting results based on the numbers of the bounding boxes or points.

During training, we employ relative coordinates and normalize all coordinate values such that the output bounding boxes and points fall within the range [0, 999], which enables Seed1.5-VL to accurately predict corresponding bounding boxes and points irrespective of the input image resolution. We apply this normalization strategy to all data related, including Optical Character Recognition (OCR) and Graphical User Interfaces (GUI).

3.1.4 3D Spatial Understanding

To enable the model’s 3D spatial understanding ability from a single image, we construct data targeting the following three tasks: relative depth sorting, absolute depth estimation, and 3D grounding. To generate the **relative depth sorting** data, we employed DepthAnything V2 [160] to infer depth relationships among objects sampled from 2 million internet images. This process yielded a dataset component comprising 3.2 billion tokens associated with this task. In particular, we select the average depth of objects with a relative depth gap beyond 20%.

Data for **absolute depth estimation** was derived from publicly available datasets. For each entity identified by its semantic mask, we determined its absolute depth using the corresponding annotated depth map. This procedure resulted in 18 million instruction pairs (e.g., query/depth value) and contributed 28 billion tokens to our pre-training corpus.

For **3D grounding** data, we utilized publicly available datasets from the internet. These datasets were then processed and reformulated into question-answering (QA) pairs. Specifically, our reformulation involved prompting for the 3D locations of objects belonging to a particular category. This process yielded a dataset of 770K instruction-following pairs, comprising 1.3 billion tokens.

3.1.5 Video

This part of data is used to improve the model’s understanding of multi-frame time-series images in video. It comprises three primary categories. Firstly, general video understanding data, this portion encompasses a variety of tasks, including video captioning, video question answering, action recognition, action grounding, and multi-image understanding. Data are sourced from public datasets and internally collected video-caption pairs. Secondly, we include several publicly available datasets for video temporal grounding and moment retrieval to enhance the model’s temporal awareness. Specifically, Seed1.5-VL directly predicts the start and end timestamps based on user prompts, with the default seconds format. Temporal grounding capability benefits complex reasoning tasks in videos. Lastly, video streaming data is crucial for understanding dynamic

and continuous video content. The data is drawn from various sources and structured into three main components:

- **Interleaved Caption/QA Data.** First, we construct interleaved video text sequences either by directly captioning segmented video clips or by constructing multi-turn question-answer pairs in chronological order. These captions and QA pairs are inserted at the corresponding timestamps within the video to enhance real-time video understanding.
- **Proactive Reasoning Data.** Second, we reconstruct grounded video question answering and dense caption data into a frame-by-frame response format. This data requires the model to continuously monitor the video stream and proactively determine the appropriate timestamps to produce responses.
- **Realtime Commentary Data.** Third, we leverage naturally temporally synchronized video commentary data to provide fine-grained interleaving and alignment of video frames and texts. This formation enables the model to handle interruptions and dynamically update responses in real-time according to the video stream.

Together, these datasets form a comprehensive foundation for effective video training.

3.1.6 Science, Technology, Engineering, and Mathematics (STEM)

To enhance the model’s reasoning capabilities during pre-training, we incorporated a diverse collection of problem-solving data across various STEM domains, obtained through both crawling and manual annotation. This effort culminated in the creation of comprehensive STEM datasets, structured around two primary components: **image comprehension data** and **problem-solving data**.

The **image comprehension data** comprises several subsets. We collected 3.2 million high-quality educational grounding samples across 300 categories within mathematics, physics, chemistry, and biology. Additionally, we synthesized 10 million structured tables with diverse formats, generated 4.5 million chemical structural diagrams, and produced 1.5 million synthetic coordinate system diagrams, including function plots and positional graphs. A specific subset, K12 Caption data, includes 100,000 human-annotated captions for educational images, 1 million visual question-answering (VQA) pairs, 1 million machine-generated captions using an automated pipeline, and hundreds of thousands of geometry-specific captions.

For the **problem-solving data** component, we processed over 100 million K12-level exercises through a rigorous cleaning and reformulation process. This was complemented by tens of millions of curated Chinese adult education problems and several million English-language image-associated questions.

The construction of these datasets employed hybrid acquisition strategies, integrating manual annotation, automated synthesis, and stringent quality control measures. This approach ensures multimodal coverage encompassing textual, visual, and diagrammatic representations across core STEM domains such as mathematics, physics, and chemistry.

3.1.7 Graphical User Interface (GUI)

For GUI data, we mainly include data curated from UI-TARS [105, 116]. Specifically, to support robust GUI perception, grounding, and reasoning, we curated a large-scale dataset across web, app, and desktop environments. Each screenshot is paired with structured metadata—element type, bounding box, text, and depth—collected via automated parsing and human-assisted exploration. For **perception**, we constructed tasks including element description, dense captioning, and state transition captioning. These tasks teach the model to identify small UI components, understand overall layouts, and detect subtle visual changes across frames. Visual markers (Set-of-Mark) are also overlaid to strengthen spatial correspondence. For **grounding**, we train the model to predict element coordinates from textual descriptions. Bounding boxes are normalized across resolutions. For **reasoning**, we collect multi-step task trajectories, each annotated with observations, intermediate thoughts, and actions. This data, combining in-house and standardized open-source traces, enables the model to learn step-by-step planning, correction, and reflection.

Stages	Stage 0	Stage 1	Stage 2
Training budget (tokens)	16B	3T	240B
Sequence length	32,768	32,768	131,072
Trainable components	MLP adaptor	all	all
Batch sizes (tokens)	8.4M	71M	71M
LR warmup steps	100	500	0
Maximum LR	2.52×10^{-4}	5.22×10^{-5}	5.22×10^{-6}
Minimum LR	4.50×10^{-5}	5.22×10^{-6}	5.22×10^{-6}

Table 4 Training setup and hyperparameters in three pre-training stages.

3.2 Training Recipe

Large multimodal models are typically trained either through joint multimodal learning from the start [54, 128], or via post-hoc adaptation after language model pre-training [16, 141]. Seed1.5-VL currently adopts the latter for flexible ablation and fast iterative development.

As delineated in section 2, our proposed model comprises three primary modules: a vision encoder, an MLP adapter, and a language model. Prior to the VLM pre-training phase, the vision encoder undergoes an independent training procedure as detailed in section 2.1. The language model is initialized from an internal pre-trained model with approximately 20 billion active parameters. This language model employs a decoder-only Mixture-of-Experts (MoE) architecture [119] and has been trained on a large-scale corpus consisting of trillions of high-quality text-only tokens. Our VLM pre-training methodology is structured into three distinct stages, as summarized in table 4:

1. In stage 0, we align the vision encoder with the language model by only training the MLP adapter while keeping the vision encoder and the language model frozen. Omitting this stage yields a slightly higher loss and worse performance.
2. In stage 1, all model parameters are trainable. This stage focuses on knowledge accumulation and mastering visual grounding and OCR capabilities of the model by training on a multimodal corpus of 3 trillion tokens, mainly composed of captions, interleaved image-text, visual grounding, and OCR data. Empirically, we found that adding a small amount of text-only tokens (e.g., 5%) can maintain the model’s language-only capabilities. Also, adding a small amount of instruction following data results in more reliable evaluation results, which allows us to decouple pre-training development from post-training’s.
3. In stage 2, we create a more balanced data mixture across different tasks, as well as adding data from new domains, such as video understanding, coding, and 3D spatial understanding. In addition, we increase the sequence length from 32,768 to 131,072, which better accommodates modeling long dependencies in videos and complex reasoning problems. Same as in stage 1, all model parameters are trainable.

We also experimented with an alternative training strategy, similar to approaches employed by [16, 141], where in stage-0 both the MLP adaptor and the vision encoder are trained while the language model remains frozen. Empirical evaluation, however, demonstrated that our training recipe yields superior performance. We hypothesize that this difference may stem from the vision encoder attempting to compensate for potential inabilities within the frozen LLM, which could consequently compromise its perceptual capabilities.

We employ the AdamW optimizer [64] in all three stages’ training with $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a weight decay of 0.1. The bias and normalization parameters are omitted from the weight decay, and other training hyperparameters can be found in table 4. Stage-0 and stage-1 training follow a full cosine decay learning rate schedule, while the starting learning rate in stage 2 is equal to the ending learning rate from stage 1 and is kept constant throughout the training. In stage 2, we load the optimizer states from stage 1, so no learning rate warmup is used.

3.3 Scaling Laws

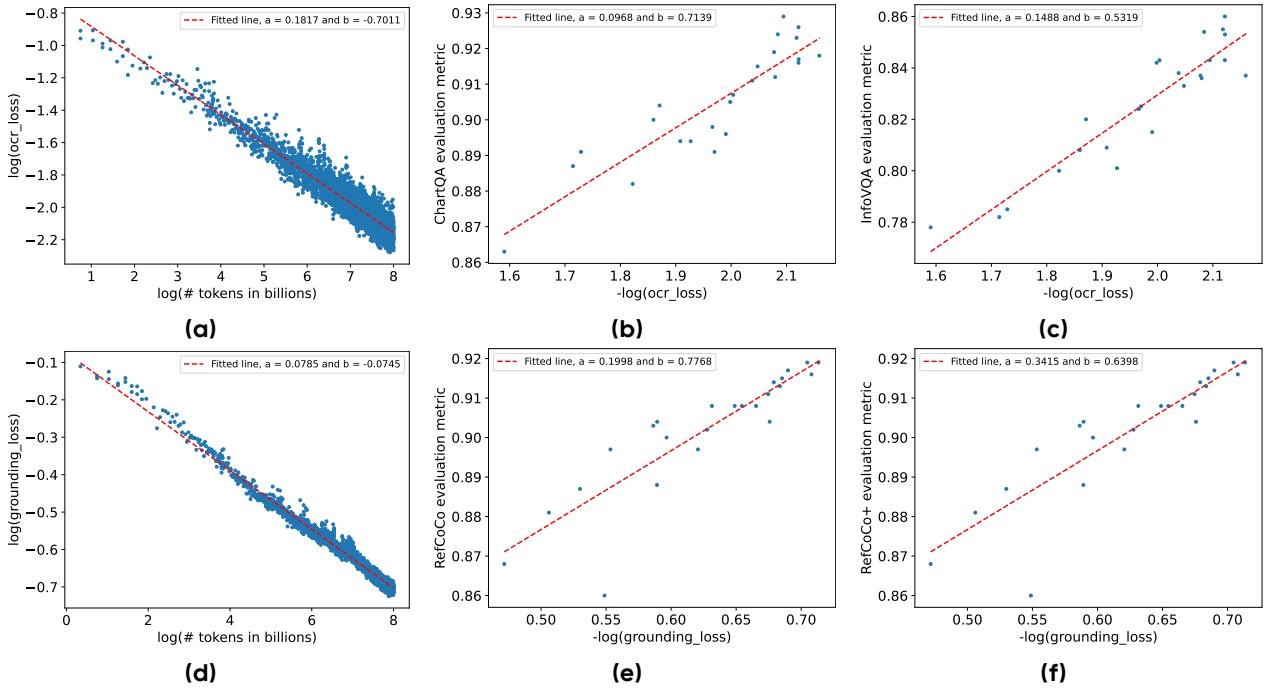


Figure 3 The relationship between the training loss of most sub-categories and training tokens obeys the power law [46]. Also, the relationship between the training loss of a sub-category and the corresponding downstream evaluation metric appears to be log-linear (e.g., metric $\sim \log(\text{loss})$) within a local neighborhood. (a) The training loss of OCR related dataset as a function of training tokens; (b) Top-1 accuracy on ChartQA [88] as a function of the training loss; (c) Top-1 accuracy on InfographicVQA [90] as a function of the training loss; (d) The training loss of grounding related dataset as a function of training tokens; (e) Precision@IoU=0.5 on RefCOCO [60, 164] as a function of the training loss; (f) Precision@IoU=0.5 on RefCOCO+ [60, 164] as a function of the training loss. Note that the evaluation metrics displayed in this figure represent performance after pre-training and are therefore not directly comparable to the final results, which are achieved following reinforcement learning (RL) as detailed in Section 6.

The pre-training of Vision-Language Models (VLMs) like Seed1.5-VL differs fundamentally from the standard practice for Large Language Models (LLMs), which typically involves random initialization of all model parameters. In contrast, Seed1.5-VL is built upon pre-trained components, including a vision encoder, an MLP adaptor, and a language model. This section focuses on understanding the scaling behavior of Seed1.5-VL during the stage-1 phase of pre-training. Based on prior work on LLM scaling laws [45, 46, 59], the average negative log-likelihood loss L is modelled as a function of model parameters N and training tokens D :

$$\hat{L} \sim \frac{A}{N^\alpha} + \frac{B}{D^\beta}. \quad (1)$$

Given that our model architecture and thus the number of parameters are fixed during this stage, equation (1) simplifies to a dependency primarily on the scale of the training data:

$$\hat{L} \sim \frac{B}{D^\beta}. \quad (2)$$

To facilitate analysis, we examine this relationship in log-log space by taking the logarithm of both sides:

$$\log(\hat{L}) \sim \log(B) - \beta \log(D) = -a \log(D) + b. \quad (3)$$

We organized our pre-training dataset into distinct categories corresponding to specific capabilities (as detailed in section 3.1). We observed that the training loss for the majority of these data sub-categories exhibits a

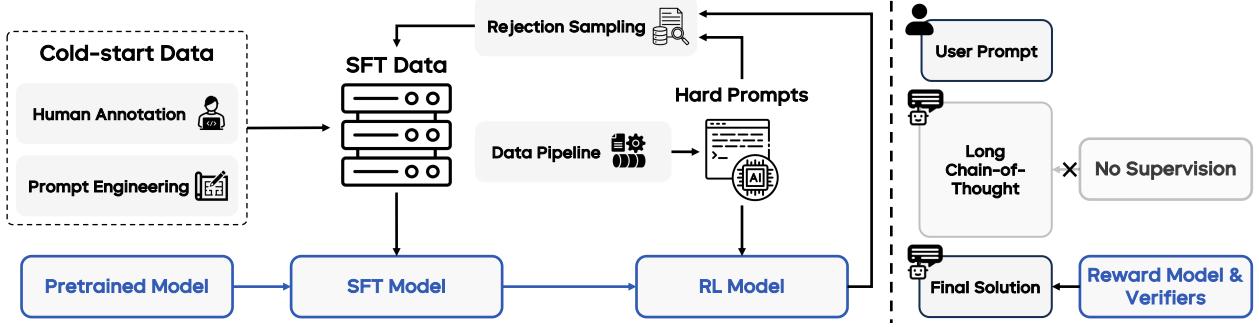


Figure 4 The overview of post-training for Seed1.5-VL. The post-training for Seed1.5-VL includes an iterative update combining rejection sampling and online reinforcement learning. We build a data pipeline including collection and curation of hard prompts for augmenting post-training data. A key aspect of our reinforcement learning implementation is that supervision, mediated by reward models and rule verifiers, is applied solely to the final generated output. We intentionally refrain from supervising the detailed chain-of-thought reasoning itself, a distinction highlighted in the illustration’s right section.

clear adherence to the scaling relationship defined by equation (3). As shown in figure 3 (a) and (d), the training losses for OCR and grounding related datasets can be modeled as follows:

$$\log(\hat{L}_{\text{ocr}}) \approx -0.1817 \log(D) - 0.7011$$

$$\log(\hat{L}_{\text{grounding}}) \approx -0.0785 \log(D) - 0.0745.$$

Beyond the scaling laws of training loss, our analysis reveals that the training loss achieved on specific data sub-categories can serve as a predictor for performance on related downstream tasks. We find that the relationship between a sub-category’s training loss and its corresponding downstream metric is approximately log-linear. However, it is important to note that such a log-linear relationship is likely sustainable only within a local neighborhood of performance values, as the range of typical evaluation metrics (e.g., accuracy, F1 score) is inherently bounded, usually between 0 and 1. As demonstrated in figure 3 (b) and (c), the top-1 accuracies on the ChartQA and InfographicVQA datasets show a clear correlation with the logarithm of the OCR training loss ($\log(\text{loss}_{\text{ocr}})$), as captured by the following approximate linear models:

$$\text{Acc}_{\text{ChartQA}} \approx -0.0968 \log(\text{loss}_{\text{ocr}}) + 0.7139$$

$$\text{Acc}_{\text{InfoVQA}} \approx -0.1488 \log(\text{loss}_{\text{ocr}}) + 0.5319$$

Analogously, figure 3 (e) and (f) detail the estimated relationship between the model’s grounding loss during training and its performance on the RefCOCO evaluation benchmark. Performance prediction remains an active research area, and prior works have used a sigmoid function to model the relationship between LLM performance and loss [37, 151] or compute [101].

4 Post-training

The post-training stage equips Seed1.5-VL with robust instruction-following and reasoning abilities through a combination of Supervised Fine-tuning (SFT) and Reinforcement Learning (RL). Depicted in figure 4, this begins with an SFT model trained on curated cold-start data. A crucial component is our data pipeline, continuously gathering hard and diverse prompts that feed into RL and improve SFT data via rejection sampling. Post-training proceeds iteratively: the SFT model is progressively enhanced by distilling the RL model’s learnings on diverse prompts. This iterative refinement continues until the prompt pool is exhausted and performance metrics converge. Ultimately, this process yields Seed1.5-VL, capable of generating both swift, succinct replies and in-depth responses featuring long Chain-of-Thought (LongCoT) reasoning [56]. We discuss details of each component in the following subsections.

4.1 Supervised Fine-tuning

The Supervised Fine-tuning (SFT) stage is integral to equipping Seed1.5-VL with foundational instruction-following and reasoning capabilities prior to reinforcement learning. Our SFT dataset comprises two primary components targeting distinct capabilities. The first component, General Instruction data, trains Seed1.5-VL on diverse, complex instructions, emphasizing the generation of concise and accurate responses. The second, Long Chain-of-Thought (LongCoT) data, focuses on generating detailed, step-by-step reasoning. This data is generated via prompt engineering and rejection sampling (inspired by [134]), mainly using high-quality outputs from Seed1.5-VL; specifics are detailed in [section 4.5](#). Besides, each data type is associated with a distinct system prompt, which allows users to dynamically toggle LongCoT reasoning during inference. The construction methodology for the SFT dataset and the specifics of Seed1.5-VL’s SFT training regimen are further elaborated in [sections 4.1.1](#) and [4.1.2](#), respectively.

4.1.1 SFT Data Construction

In the initial phase of SFT data construction, we aimed to equip the model with the ability to address a broad spectrum of application scenarios. To this end, we developed a model capability taxonomy informed by the classification of traditional visual tasks and the empirical application requirements of vision-language models. Guided by this taxonomy, we utilized crowdsourcing to collect images from the internet and generate approximately 13,000 high-quality instruction-tuning data, each comprising a prompt and a corresponding response. These initial responses were designed to exhibit strong alignment with human preferences.

To further enhance the model’s performance, we incorporated an additional 30,000 high-quality data samples sourced from the research community. These samples were curated from our carefully collected open-source repository containing approximately 1.5 million entries. Initially, we utilized a proprietary image-text embedding model to cluster the image-text pairs into task-specific categories. This clustering enabled targeted downsampling, ensuring the dataset preserved a high degree of diversity across various tasks. Subsequently, we leveraged our trained SFT model, aligned with human preferences, to perform multiple roll-outs on this sampled subset. The generated responses were filtered by LLM-as-a-judge [177], which justifies the correctness of the model’s generated responses with the original ground truth as reference. On this basis, we further adopted the Reward Model ([section 4.2.2](#)) to screen out the responses that are most aligned with human preferences from the retained results, thus obtaining the final rejection sampling fine-tuning data [134]. Eventually, we compressed the amount of open-source data in the SFT data from 1.5 million to approximately 30,000 high-quality data. The other open-source data was used in the pre-training stage in advance.

Building upon the enhanced capabilities acquired during pre-training, including complex chart understanding, STEM-related reasoning, grounding, and 3D perception, and video analysis, we iteratively increased the complexity of our fine-tuning data and instructions. This involved reducing the proportion of simple prompts readily solvable with individual capabilities and introducing more challenging questions that previously exposed limitations in the pre-trained model. Leveraging a self-instruct methodology [143], we synthesized novel complex prompts and their corresponding model responses by combining multiple simpler prompts according to various logical structures. Responses generated through self-instruct and rejection sampling underwent a manual secondary verification process to identify and rectify errors. Compared to direct human annotation, this approach of refining model-generated responses significantly improves human annotation efficiency. Moreover, it enables the exclusion of data exceeding the model’s current capacity, thereby mitigating the risk of hallucinations.

4.1.2 Training Recipe

For the SFT stage, we assembled a concise and high-quality dataset comprising approximately 50,000 samples. This multimodal SFT data was integrated with an in-house text-only SFT dataset. Together with the Long Chain-of-Thought (LongCoT) SFT data, as described in [section 4.5](#), this combined corpus was used for training over two epochs. During SFT, the vision encoder’s parameters were frozen, while all other model parameters remained trainable. The training was conducted with a sequence length of 131,072 tokens and a batch size equivalent to 16 times the sequence length. We utilized the AdamW optimizer [64] for training, with hyperparameters set to $\beta_1 = 0.9$, $\beta_2 = 0.95$, and a weight decay of 0.1. The training process included a

warm-up phase spanning 10% of the total steps, after which the learning rate decayed from a peak value of 2×10^{-5} to 2×10^{-6} following a cosine decay schedule.

4.2 Reinforcement Learning from Human Feedback

To further boost both human evaluation performance and multimodal understanding capabilities, we conduct reinforcement learning from human feedback (RLHF) [180], which involves preference data collection, reward model training, and optimization with reinforcement algorithms.

4.2.1 Preference Data

To train the reward model, we collect list-wise multimodal preference datasets for reward modeling through human annotation and heuristic synthesis.

Human annotations. The human-annotated preference data involves comparing several candidate model responses using a 5-scale rating system. The prompts for generating preference data cover all general visual understanding abilities, and maintain a balanced scale across all abilities. We utilize the current top-performing in-house models to randomly sample responses through nucleus sampling [47]. To ensure the diversity of responses, we apply filtering techniques—such as editing distance, semantic similarity, and length-balancing strategies—prior to selecting responses for human annotation. Beyond ranking the responses by quality, we instruct human annotators to select one model response that requires minimal editing to correct or improve its quality, which further compensates for the lack of diversity in the limited response sampling. Annotators are also tasked with identifying and highlighting issues within the responses—such as hallucinations, helpfulness, informativeness, etc.—and providing detailed explanations for these issues. To further enhance the efficiency of the annotation process, we employ the latest reward models to pre-annotate the rankings, offering initial guidance for human annotators. This approach not only streamlines the annotation workflow but also ensures more consistent and objective evaluations.

Synthetic data. While some recent approaches [172, 179] have used deliberate error introduction to synthesize preference pairs, multiple studies [4, 75, 162] demonstrate that such synthetic data often fails to generalize effectively, as the reward model tends to learn the inherent patterns between edited and original responses. Instead, we aggregate a diverse set of multimodal prompts with clear ground-truths, while implementing format constraints such as “Final Answer.”. For each prompt, we generate model responses K times and use existing vision-language models to evaluate their correctness and adherence to format based on the ground-truth. Consequently, we establish list-wise preferences with clear rankings: correct responses with well-defined formats rank highest, followed by incorrect responses with well-defined formats, and lastly, incorrect responses that do not follow the format. Additionally, we follow FeedQuill [162] to generate image captioning preference pairs, which helps in reducing hallucinations. All the synthetic preference data is refined using preference strength following [137].

4.2.2 VLM as a Reward Model

We initialize the reward model with an instruction-tuned VLM. Then, following [86, 120] we prompt the model π_ϕ to act as a generative classifier that directly outputs answer indicator token \hat{I} regarding the preference between two responses, y_1 and y_2 , given the prompt x . This process can be formulated as $\hat{I} \sim \pi_\phi(I|x, y_1, y_2)$.

We find that this approach yields a more robust and superior reward model compared to traditional Bradley-Terry reward modeling [100] due to its direct handling of token probabilities and response comparisons. To mitigate the potential positional bias inherent in vision-language models [176], we compute the probabilities for both possible orderings of the responses, i.e., both (x, y_1, y_2) and (x, y_2, y_1) . This ensures that the model’s preference judgment is fair and not affected by the order in which responses are presented. Additionally, during training, we apply an iterative learning strategy to maintain the consistency of annotation principles as standards evolve. This strategy involves continuously updating the training data and annotation guidelines to reflect the most current and accurate criteria. By doing so, we ensure that the reward model remains reliable and adaptable to changing requirements. This approach helps in improving the generalization capability of the model and maintaining high-quality performance over time.

4.2.3 Data Curation for Reinforcement Learning

Our online reinforcement learning implementation employs a variant of the Proximal Policy Optimization (PPO) algorithm [155]. In this approach, the reward signal is derived from the probability assigned by a reward model to the generated answer tokens. In addition, the ground truth response or the best-of-N responses from an SFT model are given as the reference answer to the reward model during PPO training.

Prompts utilized for RL training were derived from the preference dataset. It was observed that the coverage of the prompt distribution critically influences RL performance. Consequently, our data collection strategy aimed to mirror the distribution of the preference data. However, the collected prompts demonstrated significant heterogeneity in quality, characterized by highly skewed distributions across both task difficulty and ability categories. To address these issues, a multi-stage data refinement pipeline was implemented. Initially, a tagging model was trained to assign capability category labels to prompts, followed by stratified sampling to ensure a balanced representation across different ability categories. Subsequently, for each prompt, K responses were generated using state-of-the-art internal models and evaluated using the most recent iteration of our reward model. A filtering criterion was applied based on the reward score variance: prompts where the difference between the maximum and mean reward across the K responses fell below a predefined threshold were excluded. This step ensures the retention of prompts for which the reward model exhibits significant discriminative capability. Finally, during the initial phases of RL training, prompts exhibiting rapid concurrent increase in both reward and KL divergence, indicative of lower task difficulty, were subject to downsampling.

4.3 Reinforcement Learning with Verifiable Rewards

In addition to human feedback, Reinforcement Learning with Verifiable Rewards (RLVR) [68] emerges as an efficient training method for various tasks [39, 69], such as mathematical reasoning and coding where we simply use answer matching or constraint verification to train the model, instead of leveraging model-based reward estimation. In this section, we design several visual tasks whose final solutions can be precisely verified by rules or external executors, which will later be incorporated into the RLVR training.

4.3.1 Visual STEM

STEM (science, technology, engineering, and mathematics) questions usually have unique and verifiable answers, which are suitable for RLVR. We collect over one million problems with images in STEM fields, mostly on mathematics, from both open-sourced resources [85] and internal K-12 education collections.

To prepare the training data, multiple-choice questions were initially transformed into an open-ended format by removing the choices, thus forcing the model to generate the correct answer’s content and preventing random guessing. Subsequently, difficult questions were selected via rejection sampling based on the performance of the SFT model. We carefully remove questions that can be answered by text only or text and captions, ensuring shortcuts on text or superficial visual elements will not be reinforced in RL. Specifically, 16 responses were generated per question, and questions achieving either 0% or greater than 75% accuracy with the SFT model were discarded. This filtering isolates challenging prompts ($0\% < \text{accuracy} \leq 75\%$) appropriate for RLVR exploration while removing potentially erroneous or trivial questions. Lastly, a preamble instruction was prepended to prompts, instructing the model to format the final answer using designated LaTeX identifiers (e.g., `\boxed{answer}`) to enable straightforward automated extraction.

Our STEM verifier transforms the predicted answers into a sympy expression and matches it with ground truths. To ensure the accuracy of our verifier, we also remove prompts that contain multiple questions or whose ground truths are complex phrases.

4.3.2 Visual Perception and Reasoning

Verifier feedback can also be collected from various visual tasks to enhance the perception and reasoning capabilities of VLMs. Here we present some early explorations on grounding, visual puzzles, and perception-related games.

Grounding. The grounding task aims to evaluate a model’s ability to accurately associate (“ground”) textual descriptions with corresponding visual elements within an input image. For easier answer extraction, we add an instruction in the prompt to encourage the model to output the predicted bounding boxes enclosed between `<bbox>` and `</bbox>` tokens. The reward is computed as the intersection over union (IoU) between the predicted bounding box and the ground-truth one. We also optimize for pointing capability in a similar way and put the object’ center point position between `<point>` and `</point>`.

Visual Instruction Following. Instruction-following capabilities can be improved with synthetic data and rule-based verifiers [25, 161]. Following this idea, we synthesize diverse visual instructions whose outcomes can be verified by corresponding regular expressions to further enhance visual instruction-following capabilities.

Visual Puzzles & Games. Visual puzzles are tasks that require the model to gather information from a visual scene and apply reasoning techniques such as abstract reasoning, inductive reasoning, and deductive reasoning. Similar to [18, 132], we synthesize over 20k visual puzzles and their corresponding solutions for RLVR. We carefully decontaminate our synthetic training data with existing visual puzzle benchmarks, such as PuzzleVQA [18]. We also involve puzzles in graph reasoning [146] and pattern identification. Similar to the STEM verifier, we prompt models to enclose final answers of puzzles in `\boxed{answer}` and verify the prediction through a string matching algorithm.

Beyond generating natural language responses, we are exploring VLM output formats that enable direct interaction with or manipulation of image content, aiming to facilitate broader VLM applications through more intuitive and engaging interactions. Imagine, for example, AI-enhanced glasses overlaying a navigation route directly onto the user’s view, rather than relying solely on text or speech—a potentially more intuitive approach. As an initial step towards developing these interactive capabilities, we focus on visual games, which are suitable testbeds because they require strong perceptual skills and have clearly verifiable outcomes indicating success. Specifically, we target the “Spot the Differences” game, tasking the model with identifying discrepancies between two images. Crucially, the model must not only explain these differences using natural language but also output bounding boxes that precisely localize the differing regions directly on the image. We train this capability using synthetically generated data employing two methods: (1) We take images from open-sourced datasets, randomly mask segments, use a diffusion model for inpainting (see figure 5 for an example), and then filter out pairs where the inpainted content is too similar to the original; (2) To ensure the model perceives subtle differences like line width or object size, we generate additional image pairs by systematically modifying SVG properties from open-sourced datasets.



Figure 5 An example of a synthesized image pair used for training the “Spot the Differences” game, with the differences highlighted by red boxes in the left image.

4.4 Hybrid Reinforcement Learning

The Seed1.5-VL model is trained utilizing a hybrid RL framework derived from a variant of the PPO algorithm. This framework incorporates a generative RM, as detailed in [156], and integrates several advancements and

exploration techniques from recent RL research [121, 165, 167, 170]. Specifically, our training is a combination of RLHF and RLVR. We present more detailed implementations as follows:

Format reward. We predefine a response format of `<think>{thought}</think>{solution}` to ensure models provide comprehensive thoughts before giving the final solution. We set rewards to zero if the model’s responses do not comply with this format. We also apply penalties if responses fail to follow format requirements for different verifiers in various tasks.

Hybrid reward. Our training prompts are categorized into general and verifiable prompts based on tasks, rewarded with RM and the verifier, respectively. Prompts are randomly shuffled in each epoch. So, general and verifiable prompts are mixed in each batch. We truncate the thought and only keep the solution in response to the reward model. Therefore, RM will ignore the CoT thought and only focus on providing rewards for the final solution. Such modification can ease constraints on thoughts and encourage models to explore more effective CoT thoughts.

Shared critic. A single critic model architecture is employed to estimate the value function corresponding to both reward sources (i.e., the reward model and verifiers). This unified approach is viable due to both reward signals operating within the same normalized range of $[0, 1]$. Specifically, the reward model inherently generates outputs within this interval, while the outcomes derived from all verifiers are explicitly scaled to conform to the same $[0, 1]$ range. The critic model’s parameters are initialized using the weights of the pre-trained reward model. Subsequently, the critic undergoes an initial warm-up phase consisting of 100 training steps, utilizing trajectory data (rollouts) generated by the SFT model.

KL coefficients. We employ distinct KL divergence coefficients for general and verifiable prompts. Specifically, a coefficient of 1×10^{-5} is applied to general prompts, while a coefficient of 0 is used for verifiable prompts. The application of a small KL coefficient for general prompts serves to mitigate potential reward hacking. Conversely, training verifiable tasks without a KL divergence term facilitates greater exploratory capacity for the model.

Training recipe. The context length and max output length of hybrid RL training are 8,192 and 16,384, respectively. We sample 4,096 roll-outs in each episode. For training updates, we use a mini-batch size of 512 samples, performing 8 gradient steps per episode. PPO clip range for the training is 0.2. Learning rates for the actor and critic are 6×10^{-7} and 7.5×10^{-7} , respectively. The number of roll-outs is different for each prompt, as harder prompts need more comprehensive exploration. We only sample once for each prompt rewarded by the reward model, while sampling 4 or 8 times for the counterpart rewarded by verifiers. Noticeably, although we only train Seed1.5-VL with LongCoT responses in the RL stage, we still witness a significant improvement in regular responses without extended reasoning.

4.5 Iterative Update by Rejection Sampling Fine-tuning

In this work, we employ an iterative training strategy to enhance Seed1.5-VL during the RL stage. The process commences with a cold-start SFT model for LongCoT, initially trained on a limited number of low-quality LongCoT samples generated via in-context prompting of the base model with a small set of hand-annotated examples. Observing that a stronger cold-start SFT naturally leads to a stronger final model after LongCoT RL, we adopt a rejection sampling fine-tuning approach to obtain an improved starting point. Specifically, following the release of each iteration of the LongCoT RL model, we gather additional challenging prompts through our data pipeline and evaluate the latest RL model on these prompts. Correctly answered responses are then collected, in the vein of rejection sampling, and incorporated into the data for the subsequent SFT release. The same verifiers used in the RL phase are utilized to confirm the correctness of these responses. Furthermore, we implement manually crafted regular expression-based filters to remove undesirable patterns such as infinite repetition, overthinking, and other linguistic artifacts. The current iteration of Seed1.5-VL has undergone four such rounds of iteration, demonstrating consistent improvements, and this iterative refinement is expected to further enhance its performance.

5 Training Infrastructure

5.1 Large-Scale Pre-training

To accelerate and stabilize pretraining, we have developed a number of training optimizations, including hybrid parallelism, workload balancing, parallelism-aware data loading and robust training. We also apply high-performance attention kernels for context parallelism, selective activation checkpointing and offloading, kernel fusion, and fine-grained communication overlapping [13, 173]. The pretraining phase consumes 1.3 million GPU hours in total².

5.1.1 Hybrid Parallelism

Training a VLM model faces unique challenges due to the heterogeneity of both the data, which consists of visual data and natural language data, and the model, which consists of a small vision encoder and a significantly larger language model. Existing training frameworks are primarily designed for sequential unimodal tasks and fall short in VLM training. They either treat the encoder as preprocessing for the LLM’s data, or completely disaggregate the encoder from the LLM, leading to imbalanced workloads, prolonged device stalls and poor scalability. To tackle these challenges, we develop a hybrid parallelism approach [30] that parallelizes the vision encoder and the language model differently. For the vision encoder and the MLP adaptor, we leverage ZeRO data parallelism [109], while for the language model, we use standard 4-D parallelism, which combines expert parallelism [65, 123], interleaved pipeline parallelism [50, 93, 94], ZeRO-1 data parallelism [109] and context parallelism [77] for context extension. We separate the parallelism strategies for the encoder/adaptor and the LLM for efficiency and simplicity—it is challenging to integrate the encoder and the adaptor into 4-D parallelism without introducing pipeline-level imbalance. Our hybrid parallelism is simple and efficient, significantly accelerating training with minimal changes to model code.

5.1.2 Workload Balancing

Vision samples contain a varying number of images, causing computation imbalance among GPUs. We adopt a classical greedy algorithm to redistribute the vision data to achieve load balancing for the vision encoder and adaptor. Firstly, we sort the images in descending order according to their computation intensity, which is defined as the number of floating-point operations (FLOPS) needed to process each image. Secondly, we scan these images in the sorted order, and assign each image to the GPU with the lowest total computation intensity. Additionally, we leverage group-wise balancing to reduce data redistribution overhead. Instead of balancing vision data across all GPUs, we divide them into evenly sized groups and only balance vision data within each group only. Empirically, we set the group size to 128-256 GPUs.

5.1.3 Parallelism-Aware Data Loading

To reduce multimodal data IO overhead, we have also built a parallelism-aware data loader. For example, GPUs within non-data-parallel groups are expected to consume the same set of training samples. Redundantly reading the same data from the distributed file system can significantly amplify data read and preprocessing overhead, slowing down microbatch readiness. We address this problem using a parallelism-aware data loader. For example, only one GPU within a PP group loads the data while the other PP ranks receive the necessary metadata from it via broadcast. Additionally, since we use pure data parallelism for the vision encoder, each GPU only processes a portion of the loaded image data. We filter out unnecessary images before moving training batches to the GPU, reducing PCIe traffic. To hide these data broadcast and transfer costs, we use a prefetcher to ensure IO and computation fully overlap.

5.1.4 Fault Tolerance

To handle various hardware and software faults during training, we use the robust training framework MegaScale [57] to achieve fault tolerance. Once the robust training framework detects a fault, it triggers the

²For consistency, all computational costs mentioned in this report are normalized to GPU hours based on the H800.

recovery process and resumes training from the last successful checkpoint. We leverage ByteCheckpoint [136] for efficient checkpoint saving and resuming.

5.2 Post-Training Framework

We conduct hybrid reinforcement learning with both human feedback (RLHF) and verifier feedback (RLVF) of Seed1.5-VL on a verl-based [122] framework. It combines a single-controller for managing inter-RL-role dataflow and multi-controllers for managing intra-RL-role data and model parallelism. Verifiers are deployed in process-based services to isolate potential verifier faults. This design greatly simplifies deployment and development for various experiments. We use the same training system and optimization techniques as in the pretraining phase for efficient actor and critic updates, and vLLM [67] for autoregressive generation of rollouts. Specifically, actor and critic training employs 3-D parallelism [50, 93, 109, 123]; rollout generation and reward/reference model inference use replicas, each configured with tensor parallelism [115]. The RL phase of Seed1.5-VL costs 60k GPU hours. The reward model is trained using the same framework as the Seed1.5-VL pretraining phase, requiring 24k GPU hours. Post-training phases also leverage ByteCheckpoint [136] for efficient checkpoint saving and resuming.

6 Evaluation

This section is structured as follows. Quantitative results on public benchmarks are presented in [section 6.1](#), followed by an assessment of performance on agentic tasks in [section 6.2](#). The design of our internal benchmark and a comparison of our model against industry-leading models are subsequently detailed in [section 6.3](#). Model limitations are discussed in [section 6.4](#). Qualitative examples are provided in [appendix A](#), and comprehensive evaluation settings are described in [appendix B](#).

6.1 Public Benchmarks

6.1.1 Vision Encoder as a Zero-shot Classifier

We evaluate Seed-ViT using zero-shot image classification benchmarks, including ImageNet-1K [22], ImageNet-V2 [112], ImageNet-A [44], ImageNet-R [43], ImageNet-S [138], and ObjectNet [8]. As detailed in [table 5](#), Seed-ViT achieves an average zero-shot accuracy of 82.5 across these datasets, which is comparable to that of InternVL-C-6B [16], despite the fact that the number of parameters of Seed-ViT is only 9% of that of InternVL-C-6B. Impressively, compared to EVA-CLIP-18B, which has 30 \times more parameters, Seed-ViT achieves comparable accuracies on most of the ImageNet variants. Furthermore, compared to DFN-5B-CLIP-H/14++ [28], Seed-ViT demonstrates superior performance on ObjectNet (which contains images with challenging backgrounds, rotations, and viewpoints) and ImageNet-A (which contains natural adversarial examples), suggesting greater robustness of Seed-ViT to real-world variations.

Models	Seed-ViT	OpenCLIP-G/14	DFN-5B-CLIP-H/14++	InternVL-C	EVA-CLIP-18B
#Param	532M	1.8B	632M	6B	17.5B
ImageNet-1K	83.6	80.4	84.3	83.2	83.8
ImageNet-V2	77.6	73.6	78.3	77.3	77.9
ImageNet-A	85.5	69.3	79.6	83.8	87.3
ImageNet-R	95.2	92.8	94.9	95.7	95.7
ImageNet-S	74.1	69.9	73.6	74.3	74.7
ObjectNet	79.2	73.0	78.0	80.6	82.2
<i>Avg.</i>	82.5	76.5	81.4	82.5	83.6

Table 5 Comparisons of pre-trained Seed-ViT (before integration with the LLM) and existing competitors with more parameters on the common zero-shot benchmarks.

6.1.2 Vision Task Evaluation

We evaluated the performance of Seed1.5-VL on a comprehensive suite of public image benchmarks, comparing it against several state-of-the-art multimodal models including Gemini 2.5 Pro (0325 version), OpenAI o1, Claude 3.7 Sonnet, OpenAI GPT-4o, and Qwen 2.5-VL 72B. We compare Seed1.5-VL with Gemini 2.5 Pro (Preview 03-25) instead of Gemini 2.5 Pro (Preview 05-06) as Gemini 2.5 Pro (Preview 03-25) shows stronger capabilities in open visual-language benchmarks (81.7_{Preview 03-25} v.s. 79.6_{Preview 05-06} in MMMU)³. The evaluation covers capabilities ranging from multimodal reasoning and general visual question answering to document understanding, grounding, and spatial reasoning. Table 6 presents the detailed results, highlighting the highest score in bold and the second highest score underlined for each benchmark, except for FSC-147 and NYU-Depth V2 where lower is better. We report results for Seed1.5-VL in both its standard ‘non-thinking’ mode and an enhanced ‘thinking’ mode, incorporating long chain-of-thought to improve reasoning.

Multimodal Reasoning. In complex multimodal reasoning tasks, Seed1.5-VL demonstrates strong capabilities in both thinking and non-thinking modes. Notably, it achieves state-of-the-art (SOTA) performance on MathVista (85.6 thinking), V* (89.5 non-thinking), VLM are Blind (92.1 thinking), ZeroBench (sub) (30.8 thinking), and VisuLogic (35.0 thinking). On MathVista and VLM are Blind, Seed1.5-VL significantly outperforms all listed counterparts. While Gemini 2.5 Pro leads on benchmarks like MMMU (81.7 vs. 77.9 for the thinking mode in Seed1.5-VL), MMMU-Pro (68.8 vs. 67.6), MathVision (73.3 vs. 68.7), and OlympiadBench (69.8 vs. 65.0), Seed1.5-VL remains competitive, securing the second position. For ZeroBench (main), Seed1.5-VL in the thinking mode solves 2 cases, ranking second alongside OpenAI o1, behind Gemini 2.5 Pro and Claude 3.7 Sonnet. Seed1.5-VL in the non-thinking mode also significantly excels in all multimodal reasoning compared with its non-thinking counterparts.

We observed that the model naturally exhibited diverse vision-centric strategies during our first round of LongCoT RL training, such as “let me look at the image again” and “analyze details before recognizing a location”, as shown in figure 9 and figure 10, even though we had not labeled related SFT data at that time.

General Visual Question Answering. For general visual question answering benchmarks, Seed1.5-VL shows robust performance. It achieves SOTA results on RealWorldQA (78.4 thinking) and SimpleVQA (63.4 thinking). On MMStar, Seed1.5-VL (77.8 thinking) also achieves the highest score among the compared models. Similarly, on MMBench-en (89.9 thinking) and MMBench-cn (89.1 thinking), Seed1.5-VL scores are near the top performers like Gemini 2.5 Pro and Qwen 2.5-VL 72B. On HallusionBench, Seed1.5-VL (60.3 thinking) secures the second-best score, slightly behind Gemini 2.5 Pro (63.7).

Document and Chart Understanding. Seed1.5-VL excels in document and chart understanding tasks. It sets new SOTA benchmarks on TextVQA (84.2 non-thinking), InfographicVQA (91.2 thinking), and DocVQA (96.9 non-thinking), surpassing strong models like Qwen 2.5-VL 72B and Gemini 2.5 Pro in these areas. On ChartQA, Seed1.5-VL (89.1 thinking) achieves the second-highest score, only behind Qwen 2.5-VL 72B (89.5). It also delivers strong performance on AI2D (88.5 non-thinking) and OCRBench (881 non-thinking), ranking competitively behind Qwen 2.5-VL 72B and Gemini 2.5 Pro. For CharXiv (DQ), Seed1.5-VL (92.6 thinking and non-thinking) ranks second to Gemini 2.5 Pro (94.4). However, on CharXiv (RQ), its performance (60.2 thinking) lags behind the leaders Gemini 2.5 Pro (69.9) and Claude 3.7 Sonnet (68.9).

Grounding and Counting. This category highlights a significant strength of Seed1.5-VL. It achieves SOTA performance across *all* listed grounding and counting benchmarks. Specifically, Seed1.5-VL leads on BLINK (72.1 thinking), LVIS-MG (73.8 non-thinking), VisualWebBench (87.8 non-thinking), RefCOCO-avg (91.6 non-thinking), CountBench (93.7 thinking), and FSC-147 (17.9 thinking, lower is better). Notably, Seed1.5-VL achieves better performance on LVIS-MG against traditional detectors, i.e., Grounding DINO-L [14, 80], which obtains 54.4 F1-score, demonstrating the strong capability of Seed1.5-VL in terms of multi-object grounding. The consistent top performance across these diverse tasks underscores Seed1.5-VL’s superior capabilities in object localization, fine-grained visual understanding, and counting.

3D Spatial Understanding. We select depth estimation, 3D object detection, and multi-view reasoning as the three tasks to evaluate Seed1.5-VL’s capability on 3D spatial understanding. In particular, for depth

³<https://deepmind.google/technologies/gemini/pro/>

Capability	Benchmark	Seed 1.5-VL	Seed 1.5-VL	Gemini 2.5 Pro	OpenAI o1	Claude 3.7 Sonnet	OpenAI GPT-4o	Qwen 2.5-VL 72B
		thinking	non-thinking	thinking	thinking	thinking	non-thinking	non-thinking
Multimodal reasoning	MMMU	<u>77.9</u>	73.6	81.7	77.6	75.2*	70.7*	70.2
	MMMU-Pro	<u>67.6</u>	59.9	68.8 *	66.4*	50.1*	54.5*	51.1
	MathVision	<u>68.7</u>	65.5	73.3 *	63.2*	58.6*	31.2*	38.1
	OlympiadBench	<u>65.0</u>	60.4	69.8 *	48.5*	54.2*	25.9*	35.9
	MathVista	85.6	<u>83.0</u>	82.7*	71.8	74.5*	63.8*	74.8
	V*	<u>89.0</u>	89.5	79.1*	69.7*	86.4*	73.9*	86.4
	VLM are Blind	92.1	<u>90.8</u>	84.3*	57.0*	69.0*	50.4*	69
	ZeroBench (main)	<u>2</u>	0	3 *	0*	3 *	0*	0
	ZeroBench (sub)	30.8	<u>29.0</u>	26.0*	20.2*	20.4*	19.6*	13.0
	VisuLogic	35.0	<u>33.0</u>	31.0*	29.0*	24.8*	26.3*	28.0
General visual question answering	RealWorldQA	78.4	77.0	<u>78.0</u> *	77.1*	67.8*	76.2*	75.7
	SimpleVQA	63.4	<u>63.1</u>	62.0*	58.8*	50.1*	52.4*	52.4
	MMStar	77.8	76.2	<u>77.5</u> *	67.5*	68.8*	65.1*	70.8
	MMBench-en	<u>89.9</u>	88.0	90.1 *	83.8*	82.0*	84.3*	88.6
	MMBench-cn	<u>89.1</u>	88.1	89.7 *	81.3*	82.7*	82.0*	87.9
	MMVP	<u>69.3</u>	70.7	70.7 *	—†	—†	70.7 *	66.7
	HallusionBench	<u>60.3</u>	60.0	63.7 *	55.6*	58.3*	56.2*	55.2
	TextVQA	81.8	84.2	76.8*	66.2*	62.4*	81.4*	83.5
Document and chart understanding	AI2D	87.3	<u>88.5</u>	88.4*	79.5 *	82.1*	84.9*	88.7
	ChartQA	<u>89.1</u>	87.4	83.3*	83.1*	56.5*	86.7*	89.5
	InfographicVQA	91.2	<u>89.3</u>	84.3*	65.4*	66.5*	79.2*	87.3
	DocVQA	96.9	<u>96.7</u>	94.0*	81.6*	87.4*	66.2*	96.4
	OCRBench	861	<u>881</u>	866*	750*	793*	806*	885
	CharXiv (RQ)	60.2	59.8	69.9 *	55.1*	<u>68.9</u> *	52.0*	49.7*
	CharXiv (DQ)	<u>92.6</u>	92.6	94.4 *	88.9*	92.0*	86.5*	87.4*
	BLINK	72.1	70.2	<u>70.6</u> *	66.1*	62.5*	65.9*	64.4
Grounding & counting	LVIS-MG	<u>72.5</u>	73.8	63.8*	—†	—†	—†	—†
	VisualWebBench	<u>87.3</u>	88.0	<u>87.3</u> *	80.9*	85.9*	80.2*	82.3*
	RefCOCO-avg	<u>91.3</u>	91.6	74.6*	—†	—†	—†	90.3
	CountBench	93.7	93.5	91.0*	86.6*	86.1*	85.7*	93.6
	FSC-147 ↓	17.9	<u>18.6</u>	24.5*	34.3*	33.4*	46.8*	28.6*
	DA-2K	<u>91.7</u>	91.9	73.0*	72.3*	40.1*	66.9*	69.6*
3D Spatial understanding	NYU-Depth V2 ↓	<u>13.6</u>	11.6	27.5*	82.1*	92.4*	73.8*	35.5*
	All-Angles Bench	<u>58.6</u>	59.0	53.4*	54.0*	50.0	49.1*	55.7

* Results self-collected via API in April 2025.

† Invalid results due to failures in following format requirements.

Table 6 Performance of Seed1.5-VL on public visual-language benchmarks (appendix B.3) compared to previous models. All benchmarks are evaluated with greedy decoding except for Claude-3.7 Sonnet where a default sampling mode is recommended. We report Pass@1 in these benchmarks. For FSC-147 and NYU-Depth V2, Mean Absolute Error (MAE) and Absolute Relative Error (AbsRel) are used as the metrics, respectively, so lower numbers are better. For all other benchmarks, higher numbers are better. The highest score in each benchmark is marked in **bold**, and the second is underlined.

estimation, we report results on two public benchmarks, DA-2K [160] and NYU-Depth V2 [95]. In DA-2K, we follow [160] and report the accuracy of relative depth estimation between two pixels (e.g., which pixel is closer). In NYU-Depth V2, we report the standard absolute relative error measured as $|\text{dist}_{\text{pred}} - \text{dist}_{\text{gt}}|/\text{dist}_{\text{gt}}$ where $\text{dist}_{\text{pred}}$ and dist_{gt} are the predicted and ground truth distances, respectively. As shown in [table 6](#), Seed1.5-VL-thinking scores 91.7 on DA-2K and 0.136 error rate on NYU Depth V2, which surpasses previous VLMs by a large margin. In non-thinking mode, Seed1.5-VL achieves 91.9 and 0.116 error rate on DA-2K and NYU-Depth V2, respectively. For 3D object detection, we report results on SUN-RGBD [125]. In non-thinking mode, our model scores 33.5 AP@15 on SUN-RGBD surpassing Gemini 2.0 Pro Experimental, which scores 32.5 AP@15 [129]. However, we observed a performance regression using thinking mode for this task. Namely, the result is decreased to 32.0 AP@15. For multi-view reasoning, we conduct evaluation on All-Angles Bench [163]. Seed1.5-VL attains 59.0 in non-thinking mode and 58.6 in thinking mode, which significantly surpasses previous models.

In summary, Seed1.5-VL exhibits state-of-the-art or highly competitive performance across a wide range of visual language benchmarks. It particularly excels in grounding, counting, 3D spatial understanding, document understanding (TextVQA, DocVQA, InfographicVQA), and certain reasoning tasks (MathVista, VLM are Blind, etc.), establishing itself as a powerful and versatile multimodal model.

6.1.3 Video Task Evaluation

We conduct an evaluation of Seed1.5-VL’s proficiency in video understanding, assessing its capabilities across five dimensions: short video, long video, streaming video, video reasoning, and video grounding. [Table 7](#) benchmarks Seed1.5-VL against state-of-the-art (SOTA) models. Due to API limitations (e.g., network timeouts, video processing errors), we cannot evaluate certain proprietary models such as Gemini 2.5 Pro across all benchmarks. Therefore, the table reports the highest score obtained, either sourced from public reports or self-collected via API.

For short video understanding, Seed1.5-VL achieves SOTA performance on MotionBench, TVBench, Dream1K, and TempCompass, demonstrating its exceptional proficiency in processing temporal dynamics and motion patterns characteristic of concise video segments. For long video understanding, it also attains strong results with a 128K token context (up to 640 frames). We recognize the importance of extended temporal understanding and plan future work focused on expanding this context window capacity to further enhance long-form video comprehension. Regarding streaming video understanding, we evaluate on OVBench [51], OVOBench [74], StreamBench [153], and the proactive sub-task of StreamingBench [76]. Seed1.5-VL achieves SOTA performance across all these benchmarks, indicating strong potential for real-time applications such as interactive video dialogue systems. In video reasoning (Video-MMMU [49], MMVU [175]), Seed1.5-VL scores 81.4 and 70.1, respectively, currently trailing top models such as Gemini 2.5 Pro. Furthermore, Seed1.5-VL excels in video grounding tasks, specifically designed to locate temporal segments within videos corresponding to textual descriptions. It achieves SOTA performance on Charades-STA [34] and TACoS [114], demonstrating precise localization capabilities.

6.2 Multimodal Agent

Multimodal agents are systems that perceive the world through visual inputs, understand instructions in natural language, and take actions to complete tasks. Two key scenarios for evaluating such agents are GUI interaction and gameplay, which test real-world usability and complex reasoning. GUI agents simulate human-computer interaction by perceiving and acting on screen interfaces across desktops, browsers, and mobile devices. These tasks require precise visual grounding and multi-step execution. Game agents operate in visually rich and interactive environments, requiring strategic planning, real-time decision-making, and commonsense reasoning. We benchmark Seed1.5-VL across both domains—GUI operation and gameplay—using a diverse set of evaluations. Results are shown in [tables 8](#) and [9](#), where we report Seed1.5-VL’s performance under the thinking mode.

GUI Grounding. GUI grounding refers to the model’s ability to understand and localize interface elements—a fundamental skill for vision-based agents. We evaluate this capability on ScreenSpot Pro [72], which focuses on expert-annotated tasks in professional settings, and ScreenSpot v2 [149], which covers grounding across

Capability	Benchmark	Seed1.5-VL thinking	Seed1.5-VL non-thinking	Prior SOTA
Short video	MotionBench [48]	68.4	68.4	62.8 GLM-4V 76.4 InternVL-2.5 46.9* Gemini 2.5 Pro
	MVBench [73]	74.4	74.3	62.6* Gemini 2.5 Pro
	TOMATO [117]	44.7	44.2	42.0 Tarsier2 75.8* Gemini 2.5 Pro
	TVBench [19]	63.6	61.5	
	Dream-1K [139]	43.9	42.6	
	TempCompass [82]	83.7	83.1	
Long video	LongVideoBench [147]	74.0	74.4	66.7 GPT-4o 69.2* Gemini 2.5 Pro
	LVBench [142]	64.6	64.0	81.2* Gemini 2.5 Pro
	MLVU [178]	82.1	81.8	87.0* Gemini 2.5 Pro
	VideoMME(w/o sub) [32]	77.9	77.6	
	TemporalBench [12]	79.8	78.9	73.3 GPT-4o
Streaming video	OVBench [51]	60.0	59.6	54.9 PMB [51]
	OVOBench [74]	72.3	72.0	67.7 Gemini1.5-Pro
	StreamBench [153]	72.8	71.2	68.7 GPT-4o
	StreamingBench(proactive) [76]	68.0	82.8	64.7 Claude 3.5 Sonnet
Video reasoning	Video-MMMU [49]	81.4	72.1	76.7 Kimi-K1.6
	MMVU [175]	70.1	70.1	75.8* Gemini 2.5 Pro
Video grounding [†]	Charades-STA [34]	64.0	64.7	60.7 SG-DETR [36]
	TACoS [114]	49.6	47.8	42.4 SG-DETR [36]

* Results self-collected via API in April 2025.

† We adopt mIoU as the main metric for video grounding tasks.

Table 7 Seed1.5-VL performance on public video benchmarks compared to previous models. For all benchmarks, higher numbers are better. The evaluation frame rates are 2 FPS for MotionBench, MVBench, TOMATO, and TVBench, 3 FPS for Dream-1K, and 1 FPS for all other datasets.

Capability	Benchmark	Seed 1.5-VL	OpenAI CUA [98]	Claude 3.7 Sonnet [6]	UI-TARS 1.5 [116]	Kimi VL-A3B [130]	Qwen 2.5 VL 72B [7]
GUI Grounding	ScreenSpot-V2 [149]	95.2	87.9	87.6	<u>94.2</u>	92.8	-
	ScreenSpot-Pro [72]	<u>60.9</u>	23.4	27.7	61.6	34.5	43.6
Computer Use	OSWorld [152]	36.7	<u>38.1</u>	28.0	42.5	8.2	8.8
	Windows Agent Arena [11]	39.6	-	38.9	42.1	10.4	-
Browser Use	WebVoyager [42]	87.2	<u>87.0</u>	84.1	84.8	-	-
	Online-Mind2Web [158]	76.4	71.0	62.9	<u>75.8</u>	-	-
Phone Use	Android World [111]	<u>62.1</u>	-	-	64.2	-	<u>35.0</u>

Table 8 Seed1.5-VL performance on public GUI online benchmarks compared to previous models.

Game	Seed1.5-VL	UI-TARS-1.5	OpenAI CUA	Claude 3.7 Sonnet
2048 (score)	870.6	721.3	611.2	800.0
Cubinko (level)	2.0	0.0	0.0	0.0
Energy (level)	2.3	1.8	0.8	1.0
Free-The-Key (level)	1.0	0.0	0.0	0.0
Gem-11 (score)	35.1	10.8	8.7	0.0
Hex-Frvr (score)	1414.0	1583.7	651.6	523.1
Infinity-Loop (level)	1.4	0.7	0.4	0.1
Laser-Maze-Puzzle (level)	2.6	2.2	1.4	1.4
Maze:Path-of-Light (level)	1.3	0.3	0.3	0.8
Shapes (level)	2.2	1.5	0.9	0.2
Snake-Solver (level)	1.3	0.2	0.2	0.2
Tiles-Master (level)	2.3	1.7	1.5	1.6
Wood-Blocks-3d (score)	864.0	213.3	18.1	0.0
Yarn-Untangle (level)	6.0	5.7	5.1	1.6

Table 9 Seed1.5-VL performance on 14 Poki games with scores or levels completed. Models are evaluated over multiple runs, allowing up to 100 steps. For all games, higher numbers are better.

desktop, mobile, and web interfaces. Seed1.5-VL demonstrates strong grounding performance, achieving 60.9 on ScreenSpot Pro and 95.2 on ScreenSpot v2, which outperforms both OpenAI CUA and Claude 3.7 Sonnet. As the foundation of multimodal interaction, GUI grounding enables agents to perceive actionable elements and bridge perception with control.

GUI Agent. For GUI agent capability evaluation, we compare Seed1.5-VL with strong baselines such as OpenAI CUA [98] and Claude 3.7 Sonnet [6] on different GUI scenarios covering computer use, browser use, and phone use. As illustrated in [table 8](#), Seed1.5-VL consistently outperforms previous models on several key benchmarks. For instance, on OSWorld [152] and Windows Agent Arena [11], Seed1.5-VL achieves 36.7% and 39.6%, respectively, surpassing Claude 3.7 Sonnet’s 28.0% and 38.9%. In browser use, Seed1.5-VL scores 87.2% on WebVoyager [42] and 76.4% on Online-Mind2Web [158], outperforming OpenAI CUA and Claude 3.7 Sonnet, setting new state-of-the-art results. On AndroidWorld [111], a challenging mobile interface task, Seed1.5-VL also achieves a high score of 62.1%. Overall, among all the foundation VLMs (i.e., Claude 3.7 Sonnet, Kimi VL-A3B, and Qwen 2.5-VL), Seed1.5-VL achieves significantly better performance in GUI agent tasks. These results underscore Seed1.5-VL’s exceptional capabilities in executing GUI tasks and its strong generalization across diverse environments and devices, firmly establishing it as a premier position in GUI domain.

Game Agent. Gameplay serves as a rigorous benchmark for multimodal models, combining visually rich

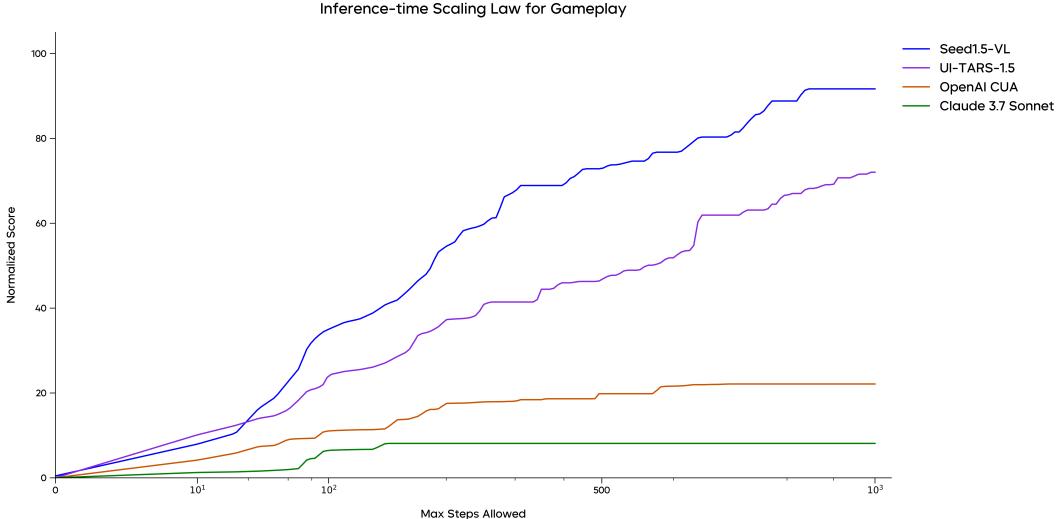


Figure 6 For each game, we compute a scaling curve per model using normalized reference scores, and averaged them to produce an overall inference-time scaling trend.

environments with complex logic that challenges models to handle intricate reasoning, sequential decision-making, and rapid adaptation. Success in gameplay depends on intuitive commonsense reasoning, long-term strategic planning, and the ability to adapt to dynamic challenges—making it an ideal testbed for showcasing the advanced cognitive capabilities of state-of-the-art multimodal agents.

We assemble a benchmark of 14 diverse games from Poki.com⁴, which assess Seed1.5-VL’s abilities in grounding, perception, and reasoning. As shown in [table 9](#), Seed1.5-VL outperforms previous models across multiple games. For example, Seed1.5-VL achieves 870.6 in 2048, surpassing OpenAI CUA (611.2) and Claude 3.7 Sonnet (800.0), and 1414.0 in Hex-Frvr, a considerable lead over OpenAI CUA (651.6) and Claude 3.7 Sonnet (523.1). These results highlight Seed1.5-VL’s exceptional performance in completing game levels and achieving high scores. In addition, the long-horizon nature of gameplay makes it particularly well-suited for evaluating inference-time scaling behaviors. As depicted in [figure 6](#), Seed1.5-VL demonstrates strong scalability, maintaining higher performance as interaction rounds increase. This showcases its robust design and advanced reasoning abilities, ensuring consistent improvement even as the complexity of tasks grows over time.

6.3 Internal Benchmarks

Besides public benchmarks, we also build internal benchmarks to comprehensively evaluate our models. We present motivation and design principles of our internal benchmarks in [section 6.3.1](#), show results in [section 6.3.2](#), and demonstrate model’s Out-of-distribution (OOD) generalization ability in [section 6.3.3](#).

6.3.1 Motivation and Design Principles

In addition to leveraging public benchmarks for exhaustive evaluation, we developed an internal benchmark suite to address several limitations inherent in existing resources. First, the predominance of English in public benchmarks necessitated the creation of comprehensive benchmarks to evaluate model performance specifically in Chinese, aligning with operational requirements. Second, the rapid pace of progress in multimodal research has resulted in saturation on many public benchmarks, reducing their sensitivity to incremental model improvements and hindering effective differentiation among leading models. Finally, limitations associated with the prevalent rule-based evaluation methods in public datasets, including challenges in answer parsing

⁴<https://poki.com>

and potential data quality issues like label errors, underscored the need for tailored internal benchmarks with potentially more robust evaluation protocols and curated data.

Consequently, we developed our in-house benchmarks guided by several core principles:

- **Focus on Core Capabilities over User Alignment:** The benchmarks prioritize assessing fundamental model abilities (e.g., perception, reasoning) rather than superficial alignment characteristics, such as preferences for response verbosity. This approach minimizes the confounding influence of alignment tuning on the evaluation of iterative model improvements.
- **Comprehensive Scope (Atomic and Integrated Capabilities):** The evaluation suite encompasses assessments of both specific, atomic capabilities (e.g., fine-grained visual recognition) and complex, integrated multimodal tasks spanning diverse application domains.
- **Evaluation Accuracy and Methodology:** We employ Large Language Models (LLMs) as judges, advancing beyond traditional rule-based metrics. The prompts and reference answers utilized by these “evaluator” models undergo continuous refinement to ensure high evaluation fidelity. Current evaluator accuracy averages above 95% for multiple-choice or simple-answer questions (e.g., single word/number responses) and exceeds 90% for open-ended questions (further details in [appendix B.1](#)).
- **Mitigation of Benchmark Overfitting:** To prevent inflated performance scores resulting from model overfitting to the benchmark data, we implement a rigorous data deduplication pipeline. Furthermore, task types and data sources within the benchmarks are periodically refreshed.
- **Task and Input Diversity:** Recognizing the critical role of diversity for VLMs, our benchmarks emphasize variety in both task types and input images. Image sourcing prioritizes non-publicly crawled data when feasible. We structure the benchmarks across numerous distinct dimensions, resulting in over 100 tasks and more than 12,000 samples from varied sources and domains. This includes a dedicated Out-of-Distribution (OOD) category featuring unconventional tasks designed to probe model generalization capabilities. A detailed taxonomy of targeted capabilities is provided in [appendix B.1](#).

6.3.2 Comparison with State-of-the-arts

Level-1 Capabilities	Level-2 Capabilities	Weight	Seed 1.5-VL thinking	Gemini 2.5 Pro thinking	OpenAI o1 thinking	OpenAI o4-mini w/o tool use	Claude 3.7 Sonnet thinking
Overall		1.0	<u>59.3</u>	61.6	54.0	55.4	48.6
Vision Capabilities	Perception	0.1	<u>63.0</u>	64.4	51.6	56.8	48.4
	Recognition	0.1	<u>72.4</u>	74.8	<u>74.5</u>	64.8	55.7
	OCR	0.1	<u>67.2</u>	70.7	55.7	64.4	57.1
	Caption & Counterfactual	0.05	<u>47.7</u>	54.9	43.6	27.6	34.1
Integrated Capabilities	OOD	0.15	44.1	<u>43.1</u>	42.3	38.4	35.9
	STEM	0.04	<u>63.3</u>	64.0	56.1	55.0	45.2
	Knowledge	0.06	64.9	73.6	<u>68.5</u>	57.8	50.8
	Reasoning	0.1	47.6	52.4	44.9	57.4	39.6
	Document & Diagram Understanding	0.1	<u>73.1</u>	75.5	66.3	70.9	64.7
	Agent	0.1	63.1	63.1	53.2	52.9	53.2
	Atomic Instruction Following	0.03	69.6	<u>69.2</u>	63.8	68.7	50.5
	Code	0.05	44.0	43.7	39.9	60.6	<u>54.6</u>
	ToB	0.02	<u>47.1</u>	54.7	30.2	39.8	29.1

Table 10 Evaluation results comparing Seed1.5-VL and state-of-the-art models on the internal benchmark. The overall score is calculated as a weighted average across performance in defined sub-categories. Data for other models was sourced via API access in April 2025. Weights for averaging are set for minimizing variance of evaluation and highlighting the importance of each category. The highest scores are marked in **bold** and the second is underlined.

We compare Seed1.5-VL with leading industry models (Gemini 2.5 Pro, OpenAI o1, OpenAI o4-mini, Claude 3.7) in [table 10](#) under *thinking mode*. The leading score of 61.6 (Gemini 2.5 Pro) highlights substantial room

for improvement on this benchmark, unlike many public benchmarks nearing saturation above 80 in [table 6](#). A more comprehensive comparison including *non-thinking* models can be found in [appendix B.2](#).

Seed1.5-VL achieves the second-highest overall score. It achieves state-of-the-art performance in OOD, Agent, Atomic Instruction Following categories, and shows strong capabilities in STEM and Document & Diagram Understanding. Its primary weaknesses relative to the top performer are observed in knowledge, reasoning, code, and captioning/counterfactual tasks. We attribute this gap partly to the scale of the current model, which utilizes a language model with approximately 20B active parameters. Evidence supporting potential gains from further scaling is presented in [figure 3](#), where the training loss shows no sign of saturation after 3 trillion tokens, and evaluation metrics correlate strongly with loss. Therefore, we expect the performance gap to diminish as we increase the model size and the training compute.

Grouping models strictly by parameter count is challenging due to the lack of public disclosure of specific parameter details for many models. Our model’s size is comparable to the recently released Llama 4 Maverick [91], which is reported to utilize 17 billion active parameters and employs a Mixture-of-Experts (MoE) architecture. Our evaluation demonstrates that Seed1.5-VL achieves significantly better performance than Llama 4 Maverick on this benchmark ([figure 29](#)).

6.3.3 Out-of-distribution Generalization

As shown in [table 10](#), our model demonstrates performance comparable to that of two leading industry models, Gemini Pro 2.5 and OpenAI o1, within the Out-of-Distribution (OOD) category of the benchmark. Complementary to standardized evaluations, an internal Chatbot platform was developed to assess the model’s ability to integrate multiple atomic capabilities in complex real-world scenarios. Three representative examples are highlighted. First, the model successfully solves a Rebus puzzle shown in [figure 7](#), leveraging its OCR, knowledge retrieval, and reasoning abilities. Additional examples, some of which may be challenging for humans, are provided in [appendix A.3](#). In the second example ([figure 8](#)), it processes code captured from a challenging visual input (a shiny whiteboard), analyzes and corrects errors, and validates the corrected code using user-provided examples. In the third example ([figure 23](#)), the model parses a diagram of an unseen format, performing “inverse graphics” to generate the underlying Mermaid code while adhering to a specific spatial constraint (“south region”).

Get it Get it Get it Get it	<p>Prompts: can you solve this Rebus puzzle?</p> <p>Model Responses: <think>...</think>The Rebus puzzle features the phrase “Get it” repeated four times. The number of repetitions (four) combined with the phrase “Get it” creates a homophonic play: “Four get it” sounds like “Forget it”. This wordplay uses the repetition count (four) to form the first part of the phrase, which phonetically matches “forget”. Thus, the solution is Forget it.</p>
--	--

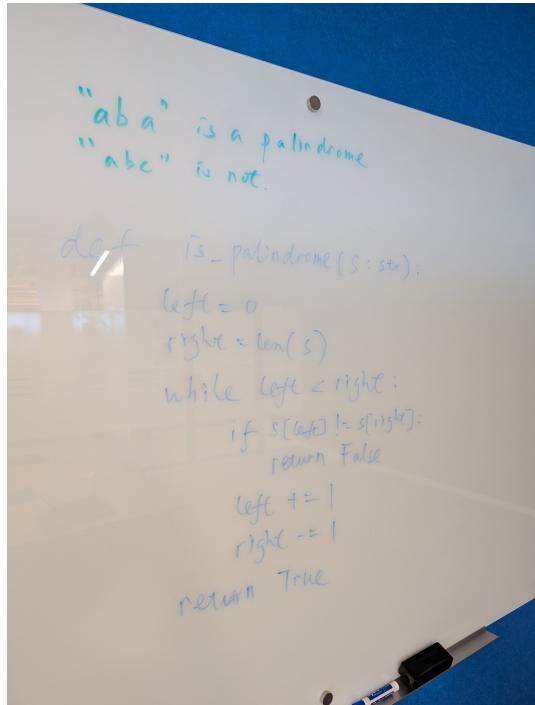
Figure 7 To solve this Rebus puzzle, our model has to correctly recognize characters in the image, count the number of repetitions, retrieve knowledge about Rebus puzzles, and perform reasoning to obtain the final answer. Thinking tokens are omitted for brevity.

A distinct benchmark, supplementing the internal evaluation, was employed to assess the new model’s user preference alignment capability. This evaluation involved augmenting the model with search tools and serves as a proxy for its utility to the Doubaot⁵ user base. The usefulness rate for Seed 1.5VL is 62.6%. Similarly, we have also tested Gemini 2.5 pro 0325, with the usefulness rates reaching 57.4%.

6.4 Limitations

Despite strong performance across many benchmarks, Seed1.5-VL exhibits certain limitations, particularly in fine-grained visual perception and complex reasoning.

⁵<https://www.doubaot.com/chat/>



Prompts: There is a bug in the code on the whiteboard. Can you spot it and return the correct code?

Model Responses: <think>...</think>The bug is in the initialization of the right pointer. It should be `len(s) - 1` (since Python uses zero-based indexing, the last character is at index `len(s)-1`, not `len(s)`). Here's the corrected code:

```
def is_palindrome(s: str) -> bool:
    left = 0
    right = len(s) - 1 # Fixed: use last index (len(s)-1)
    instead of len(s)
    while left < right:
        if s[left] != s[right]:
            return False
        left += 1
        right -= 1
    return True
```

Explanation:

1. **Right Pointer Initialization:** `right = len(s) - 1` ensures we start at the last character (correct index), avoiding an `IndexError`.
2. **Loop Logic:** Compares characters from the start (`left`) and end (`right`), moving pointers inward. If any pair mismatches, returns `False`. If the loop completes (all pairs match), returns `True`.

Testing:

- "aba": `left=0, right=2` → match ('a' == 'a'), then `left=1, right=1` → loop ends, return `True` (correct).
- "abc": `left=0, right=2` → 'a' != 'c' → return `False` (correct).

This fixes the index error and correctly implements the palindrome check using two pointers.

Figure 8 An example of Seed1.5-VL correcting code written on a whiteboard using its OCR and coding capabilities.

In complex visual perception tasks, Seed1.5-VL struggles with accurately counting objects when they are irregularly arranged, similar in color, or partially occluded. Identifying subtle differences between images also presents a challenge, sometimes leading to missed details or inaccurate descriptions. Furthermore, similar to some contemporaries (e.g., OpenAI GPT-4o and Gemini 2.5 Pro), difficulties can arise in precisely interpreting complex spatial relationships, especially with varying perspectives, and accurately responding to visual prompts, occasionally misidentifying content or quantity within specified regions or making localization errors.

Beyond perception, challenges also emerge in higher-level reasoning, as suggested by both open and internal benchmarks. Seed1.5-VL demonstrates suboptimal performance on reasoning tasks trivial for humans, such as solving Klotski puzzles or navigating simple mazes, suggesting a need for future exploration into techniques like visual Chain-of-Thought (CoT) where the model might augment its process with auxiliary visual aids (e.g., lines). Reasoning requiring combinatorial search poses a significant challenge for many existing VLMs. Figures 27 and 28 provide two illustrative examples of problems falling into this category. While challenging for current VLM architectures, combinatorial search tasks are often more readily addressed through programmatic or code-based approaches. Consequently, incorporating code-use and other external tools into VLM frameworks represents an important direction for future research aimed at enhancing such reasoning capabilities.

Limitations are observed in 3D spatial reasoning tasks for most VLMs. These challenges include, for example, tasks involving 3D object manipulation or reasoning about the projection of 3D objects. Specific instances illustrating such failure cases are provided in figure 25. A potential direction to address this limitation is the incorporation of image generation capabilities into the foundation model, which could further enable visual chain-of-thought mechanisms. This approach remains a subject for future research.

Additionally, VLMs sometimes still produce incorrect inferences, particularly in tasks such as visual puzzles.

These reasoning errors may stem from underlying perceptual mistakes (misinterpreting shapes or conditions) or from limitations in the logical deduction process itself. In tasks requiring planning or adherence to complex instructions, Seed1.5-VL may overlook specific conditions or introduce unfounded assumptions, which can lead to incomplete or invalid responses.

Our internal evaluations also revealed model’s deficiency in temporal reasoning capability, as the model faced difficulties in discerning the chronological sequence of continuous actions or inferring order from the before-and-after states of objects. The capacity for multi-image reasoning is limited, with performance degrading on tasks requiring the synthesis of clues across multiple images with strong logical interdependencies.

Finally, hallucination persists as a significant challenge for all VLMs. [Figure 26](#) illustrates a particularly notable case where the visual input conflicts with the prior knowledge acquired from the language model component. In such instances, models tend to prioritize this acquired knowledge, effectively overriding or conforming the perceived visual information to align with learned priors.

7 Conclusion and Next Steps

In this paper, we presented Seed1.5-VL, our latest multimodal foundation model demonstrating strong capabilities in reasoning, OCR, diagram understanding, visual grounding, 3D spatial understanding, and video understanding. Despite its relatively moderate size, Seed1.5-VL achieves state-of-the-art results on 38 out of 60 evaluated public benchmarks, including a score of 77.9 on the MMMU benchmark, widely regarded as a key indicator of multimodal reasoning ability.

Beyond benchmark performance, Seed1.5-VL exhibits significant integrated capabilities and generalization to tasks dissimilar to its training data. Examples include solving complex visual reasoning tasks such as Rebus puzzles, interpreting and correcting handwritten code from whiteboard images, and functioning as an agent for computer interaction and gameplay. Further exploration of these emergent abilities is warranted.

Our scaling analysis indicates that model performance shows no sign of saturation, suggesting that increasing model parameters and training compute represents a promising immediate direction. Through our evaluations, we also identified limitations common to contemporary VLMs, such as robust 3D spatial reasoning, hallucination mitigation, and complex combinatorial search. Addressing these challenges constitutes a core part of our ongoing research, which includes efforts towards unifying existing model capabilities with image generation (potentially enabling visual Chain-of-Thought) and incorporating robust tool-use mechanisms.

The advancements presented here build upon substantial prior work within the AI research community, leveraging foundational developments like the Transformer and Vision Transformer architectures. To contribute to future progress, we have detailed our model architecture, data synthesis pipeline, training methodology, training framework innovations, and internal evaluation design in this report.

References

- [1] Fuyu-8b: A multimodal architecture for ai agents. <https://www.adept.ai/blog/fuyu-8b>, 2023.
- [2] Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. [arXiv preprint arXiv:2404.14219](https://arxiv.org/abs/2404.14219), 2024.
- [3] Niki Amini-Naieni, Tengda Han, and Andrew Zisserman. Countgd: Multi-modal open-world counting. *Advances in Neural Information Processing Systems*, 37:48810–48837, 2024.
- [4] Elmira Amirloo, Jean-Philippe Fauconnier, Christoph Roesmann, Christian Kerl, Rinu Boney, Yusu Qian, Zirui Wang, Afshin Dehghan, Yinfei Yang, Zhe Gan, et al. Understanding alignment in multimodal llms: A comprehensive study. [arXiv preprint arXiv:2407.02477](https://arxiv.org/abs/2407.02477), 2024.
- [5] Anthropic. Claude 3.7 sonnet system card. 2025.
- [6] anthropic. Claude’s extended thinking, 2025. URL <https://www.anthropic.com/news/visible-extended-thinking>.
- [7] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. Qwen2. 5-vl technical report. [arXiv preprint arXiv:2502.13923](https://arxiv.org/abs/2502.13923), 2025.
- [8] Andrei Barbu, David Mayo, Julian Alverio, William Luo, Christopher Wang, Dan Gutfreund, Josh Tenenbaum, and Boris Katz. Objectnet: A large-scale bias-controlled dataset for pushing the limits of object recognition models. *Advances in neural information processing systems*, 32, 2019.
- [9] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, et al. Arkitscenes: A diverse real-world dataset for 3d indoor scene understanding using mobile rgb-d data. [arXiv preprint arXiv:2111.08897](https://arxiv.org/abs/2111.08897), 2021.
- [10] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschannen, Emanuele Bugliarello, et al. Paligemma: A versatile 3b vlm for transfer. [arXiv preprint arXiv:2407.07726](https://arxiv.org/abs/2407.07726), 2024.
- [11] Rogerio Bonatti, Dan Zhao, Francesco Bonacci, Dillon Dupont, Sara Abdali, Yinheng Li, Yadong Lu, Justin Wagle, Kazuhito Koishida, Arthur Bucker, et al. Windows agent arena: Evaluating multi-modal os agents at scale. [arXiv preprint arXiv:2409.08264](https://arxiv.org/abs/2409.08264), 2024.
- [12] Mu Cai, Reuben Tan, Jianrui Zhang, Bocheng Zou, Kai Zhang, Feng Yao, Fangrui Zhu, Jing Gu, Yiwu Zhong, Yuzhang Shang, et al. Temporalbench: Benchmarking fine-grained temporal understanding for multimodal video models. [arXiv preprint arXiv:2410.10818](https://arxiv.org/abs/2410.10818), 2024.
- [13] Li-Wen Chang, Wenlei Bao, Qi Hou, Chengquan Jiang, Ningxin Zheng, Yinmin Zhong, Xuanrun Zhang, Zuquan Song, Chengji Yao, Ziheng Jiang, et al. Flux: Fast software-based communication overlap on gpus through kernel fusion. [arXiv preprint arXiv:2406.06858](https://arxiv.org/abs/2406.06858), 2024.
- [14] Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, Zheng Zhang, Dazhi Cheng, Chenchen Zhu, Tianheng Cheng, Qijie Zhao, Buyu Li, Xin Lu, Rui Zhu, Yue Wu, Jifeng Dai, Jingdong Wang, Jianping Shi, Wanli Ouyang, Chen Change Loy, and Dahua Lin. MMDetection: Open mmlab detection toolbox and benchmark. [arXiv preprint arXiv:1906.07155](https://arxiv.org/abs/1906.07155), 2019.
- [15] Lin Chen, Jinsong Li, Xiaoyi Dong, Pan Zhang, Yuhang Zang, Zehui Chen, Haodong Duan, Jiaqi Wang, Yu Qiao, Dahua Lin, et al. Are we on the right way for evaluating large vision-language models? [arXiv preprint arXiv:2403.20330](https://arxiv.org/abs/2403.20330), 2024.
- [16] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 24185–24198, 2024.
- [17] Tianheng Cheng, Lin Song, Yixiao Ge, Wenyu Liu, Xinggang Wang, and Ying Shan. Yolo-world: Real-time open-vocabulary object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16901–16911, June 2024.

- [18] Yew Ken Chia, Vernon Toh Yan Han, Deepanway Ghosal, Lidong Bing, and Soujanya Poria. Puzzlevqa: Diagnosing multimodal reasoning challenges of language models with abstract visual patterns, 2024. URL <https://arxiv.org/abs/2403.13315>.
- [19] Daniel Cores, Michael Dorkenwald, Manuel Mucientes, Cees GM Snoek, and Yuki M Asano. Tvbench: Redesigning video-language evaluation. *arXiv preprint arXiv:2410.07752*, 2024.
- [20] Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim M Alabdulmohsin, et al. Patch n'pack: Navit, a vision transformer for any aspect ratio and resolution. *Advances in Neural Information Processing Systems*, 36: 2252–2274, 2023.
- [21] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muenmighoff, Kyle Lo, Luca Soldaini, et al. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*, 2024.
- [22] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009.
- [23] Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.
- [24] Haiwen Diao, Yufeng Cui, Xiaotong Li, Yueze Wang, Huchuan Lu, and Xinlong Wang. Unveiling encoder-free vision-language models. *arXiv preprint arXiv:2406.11832*, 2024.
- [25] Guanting Dong, Keming Lu, Chengpeng Li, Tingyu Xia, Bowen Yu, Chang Zhou, and Jingren Zhou. Self-play with execution feedback: Improving instruction-following capabilities of large language models, 2024. URL <https://arxiv.org/abs/2406.13542>.
- [26] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [27] Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman. Counting out time: Class agnostic video repetition counting in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10387–10396, 2020.
- [28] Alex Fang, Albin Madappally Jose, Amit Jain, Ludwig Schmidt, Alexander Toshev, and Vaishaal Shankar. Data filtering networks. *arXiv preprint arXiv:2309.17425*, 2023.
- [29] Yuxin Fang, Wen Wang, Binhuai Xie, Quan Sun, Ledell Wu, Xinggang Wang, Tiejun Huang, Xinlong Wang, and Yue Cao. Eva: Exploring the limits of masked visual representation learning at scale. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19358–19369, 2023.
- [30] Weiqi Feng, Yangrui Chen, Shaoyu Wang, Yanghua Peng, Haibin Lin, and Minlan Yu. Optimus: Accelerating large-scale multi-modal llm training by bubble exploitation. *arXiv preprint arXiv:2408.03505*, 2024.
- [31] Figure AI. Helix: A vision-language-action model for generalist humanoid control. <https://www.figure.ai/news/helix>, 2025. Accessed: 2025-04-23.
- [32] Chaoyou Fu, Yuhang Dai, Yongdong Luo, Lei Li, Shuhuai Ren, Renrui Zhang, Zihan Wang, Chenyu Zhou, Yunhang Shen, Mengdan Zhang, et al. Video-mme: The first-ever comprehensive evaluation benchmark of multi-modal llms in video analysis. *arXiv preprint arXiv:2405.21075*, 2024.
- [33] Xingyu Fu, Yushi Hu, Bangzheng Li, Yu Feng, Haoyu Wang, Xudong Lin, Dan Roth, Noah A Smith, Wei-Chiu Ma, and Ranjay Krishna. Blink: Multimodal large language models can see but not perceive. In *European Conference on Computer Vision*, pages 148–166. Springer, 2024.
- [34] Jiyang Gao, Chen Sun, Zhenheng Yang, and Ram Nevatia. Tall: Temporal activity localization via language query. In *Proceedings of the IEEE international conference on computer vision*, pages 5267–5275, 2017.
- [35] Google. Experiment with gemini 2.0 flash native image generation. <https://developers.googleblog.com/en/experiment-with-gemini-20-flash-native-image-generation>, 2025.
- [36] Aleksandr Gordeev, Vladimir Dokholyan, Irina Tolstykh, and Maksim Kuprashevich. Saliency-guided detr for moment retrieval and highlight detection. *arXiv preprint arXiv:2410.01615*, 2024.

- [37] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. [arXiv preprint arXiv:2407.21783](#), 2024.
- [38] Tianrui Guan, Fuxiao Liu, Xiyang Wu, Ruiqi Xian, Zongxia Li, Xiaoyu Liu, Xijun Wang, Lichang Chen, Furong Huang, Yaser Yacoob, et al. Hallusionbench: an advanced diagnostic suite for entangled language hallucination and visual illusion in large vision-language models. In [Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition](#), pages 14375–14385, 2024.
- [39] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. [arXiv preprint arXiv:2501.12948](#), 2025.
- [40] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In [Proceedings of the IEEE/CVF conference on computer vision and pattern recognition](#), pages 5356–5364, 2019.
- [41] Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, Jinyi Hu, Xu Han, Yujie Huang, Yuxiang Zhang, et al. Olympiadbench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. [arXiv preprint arXiv:2402.14008](#), 2024.
- [42] Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. Webvoyager: Building an end-to-end web agent with large multimodal models. [arXiv preprint arXiv:2401.13919](#), 2024.
- [43] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. The many faces of robustness: A critical analysis of out-of-distribution generalization. In [Proceedings of the IEEE/CVF international conference on computer vision](#), pages 8340–8349, 2021.
- [44] Dan Hendrycks, Kevin Zhao, Steven Basart, Jacob Steinhardt, and Dawn Song. Natural adversarial examples. In [Proceedings of the IEEE/CVF conference on computer vision and pattern recognition](#), pages 15262–15271, 2021.
- [45] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. [arXiv preprint arXiv:2010.14701](#), 2020.
- [46] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. [arXiv preprint arXiv:2203.15556](#), 2022.
- [47] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. [arXiv preprint arXiv:1904.09751](#), 2019.
- [48] Wenyi Hong, Yean Cheng, Zhuoyi Yang, Weihan Wang, Lefan Wang, Xiaotao Gu, Shiyu Huang, Yuxiao Dong, and Jie Tang. Motionbench: Benchmarking and improving fine-grained video motion understanding for vision language models. [arXiv preprint arXiv:2501.02955](#), 2025.
- [49] Kairui Hu, Penghao Wu, Fanyi Pu, Wang Xiao, Yuanhan Zhang, Xiang Yue, Bo Li, and Ziwei Liu. Video-mmmu: Evaluating knowledge acquisition from multi-discipline professional videos. [arXiv preprint arXiv:2501.13826](#), 2025.
- [50] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, HyoukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. [Advances in neural information processing systems](#), 32, 2019.
- [51] Zhenpeng Huang, Xinhao Li, Jiaqi Li, Jing Wang, Xiangyu Zeng, Cheng Liang, Tao Wu, Xi Chen, Liang Li, and Limin Wang. Online video understanding: A comprehensive benchmark and memory-augmented method. [arXiv preprint arXiv:2501.00584](#), 2024.
- [52] Zilong Huang, Qinghao Ye, Bingyi Kang, Jiashi Feng, and Haoqi Fan. Classification done right for vision-language pre-training. [Advances in Neural Information Processing Systems](#), 37:96483–96504, 2024.
- [53] J. D. Hunter. Matplotlib: A 2d graphics environment. [Computing in Science & Engineering](#), 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.

- [54] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. [arXiv preprint arXiv:2410.21276](#), 2024.
- [55] Physical Intelligence, Kevin Black, Noah Brown, James Darpinian, Karan Dhabalia, Danny Driess, Adnan Esmail, Michael Equi, Chelsea Finn, Niccolo Fusai, Manuel Y. Galliker, Dibya Ghosh, Lachy Groom, Karol Hausman, Brian Ichter, Szymon Jakubczak, Tim Jones, Liyiming Ke, Devin LeBlanc, Sergey Levine, Adrian Li-Bell, Mohith Mothukuri, Suraj Nair, Karl Pertsch, Allen Z. Ren, Lucy Xiaoyang Shi, Laura Smith, Jost Tobias Springenberg, Kyle Stachowicz, James Tanner, Quan Vuong, Homer Walke, Anna Walling, Haohuan Wang, Lili Yu, and Ury Zhilinsky. $\pi_{0.5}$: A vision-language-action model with open-world generalization. [arXiv preprint arXiv:2504.16054](#), 2025. URL <https://arxiv.org/abs/2504.16054>.
- [56] Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. [arXiv preprint arXiv:2412.16720](#), 2024.
- [57] Ziheng Jiang, Haibin Lin, Yinmin Zhong, Qi Huang, Yangrui Chen, Zhi Zhang, Yanghua Peng, Xiang Li, Cong Xie, Shibiao Nong, et al. {MegaScale}: Scaling large language model training to more than 10,000 {GPUs}. In *21st USENIX Symposium on Networked Systems Design and Implementation (NSDI 24)*, pages 745–760, 2024.
- [58] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. [arXiv preprint arXiv:1710.07300](#), 2017.
- [59] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. [arXiv preprint arXiv:2001.08361](#), 2020.
- [60] Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. Referitgame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 787–798, 2014.
- [61] Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. A diagram is worth a dozen images. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 235–251. Springer, 2016.
- [62] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang, Sangdoo Yun, Dongyo Han, and Seunghyun Park. Ocr-free document understanding transformer. In *European Conference on Computer Vision (ECCV)*, 2022.
- [63] Moo Jin Kim, Karl Pertsch, Siddharth Karamcheti, Ted Xiao, Ashwin Balakrishna, Suraj Nair, Rafael Rafailev, Ethan P Foster, Grace Lam, Pannag R Sanketi, Quan Vuong, Thomas Kollar, Benjamin Burchfiel, Russ Tedrake, Dorsa Sadigh, Sergey Levine, Percy Liang, and Chelsea Finn. Openvla: An open-source vision-language-action model. In Pukit Agrawal, Oliver Kroemer, and Wolfram Burgard, editors, *Proceedings of The 8th Conference on Robot Learning*, volume 270 of *Proceedings of Machine Learning Research*, pages 2679–2713. PMLR, 06–09 Nov 2025. URL <https://proceedings.mlr.press/v270/kim25c.html>.
- [64] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. [arXiv preprint arXiv:1412.6980](#), 2014.
- [65] Vijay Anand Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. Reducing activation recomputation in large transformer models. *Proceedings of Machine Learning and Systems*, 5:341–353, 2023.
- [66] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Matteo Malloci, Alexander Kolesnikov, Tom Duerig, and Vittorio Ferrari. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. *IJCV*, 2020.
- [67] Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention, 2023. URL <https://arxiv.org/abs/2309.06180>.
- [68] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. T\ ulu 3: Pushing frontiers in open language model post-training. [arXiv preprint arXiv:2411.15124](#), 2024.
- [69] Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang,

- Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025. URL <https://arxiv.org/abs/2411.15124>.
- [70] Deqing Li, Honghui Mei, Yi Shen, Shuang Su, Wenli Zhang, Junting Wang, Ming Zu, and Wei Chen. Echarts: a declarative framework for rapid construction of web-based visualization. *Visual Informatics*, 2(2):136–146, 2018.
- [71] Feng Li, Renrui Zhang, Hao Zhang, Yuanhan Zhang, Bo Li, Wei Li, Zejun Ma, and Chunyuan Li. Llava-next-interleave: Tackling multi-image, video, and 3d in large multimodal models. *arXiv preprint arXiv:2407.07895*, 2024.
- [72] Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. Screenspot-pro: Gui grounding for professional high-resolution computer use. *arXiv preprint arXiv:2504.07981*, 2025.
- [73] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Yi Liu, Zun Wang, Jilan Xu, Guo Chen, Ping Luo, et al. Mvbench: A comprehensive multi-modal video understanding benchmark. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22195–22206, 2024.
- [74] Yifei Li, Junbo Niu, Ziyang Miao, Chunjiang Ge, Yuanhang Zhou, Qihao He, Xiaoyi Dong, Haodong Duan, Shuanggrui Ding, Rui Qian, et al. Ovo-bench: How far is your video-llms from real-world online video understanding? *arXiv preprint arXiv:2501.05510*, 2025.
- [75] Zichao Li, Xueru Wen, Jie Lou, Yuqiu Ji, Yaojie Lu, Xianpei Han, Debing Zhang, and Le Sun. The devil is in the details: Tackling unimodal spurious correlations for generalizable multimodal reward models. *arXiv preprint arXiv:2503.03122*, 2025.
- [76] Junming Lin, Zheng Fang, Chi Chen, Zihao Wan, Fuwen Luo, Peng Li, Yang Liu, and Maosong Sun. Streaming-bench: Assessing the gap for mllms to achieve streaming video understanding. *arXiv preprint arXiv:2411.03628*, 2024.
- [77] Hao Liu, Matei Zaharia, and Pieter Abbeel. Ring attention with blockwise transformers for near-infinite context. *arXiv preprint arXiv:2310.01889*, 2023.
- [78] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023.
- [79] Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang Yue. Visual-webbench: How far have multimodal llms evolved in web page understanding and grounding? *arXiv preprint arXiv:2404.05955*, 2024.
- [80] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55. Springer, 2024.
- [81] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, et al. Mmbench: Is your multi-modal model an all-around player? In *European conference on computer vision*, pages 216–233. Springer, 2024.
- [82] Yuanxin Liu, Shicheng Li, Yi Liu, Yuxiang Wang, Shuhuai Ren, Lei Li, Sishuo Chen, Xu Sun, and Lu Hou. Tempcompass: Do video llms really understand videos? *arXiv preprint arXiv:2403.00476*, 2024.
- [83] Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. Ocrbench: on the hidden mystery of ocr in large multimodal models. *Science China Information Sciences*, 67(12):220102, 2024.
- [84] Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. *arXiv preprint arXiv:2310.02255*, 2023.
- [85] Ruilin Luo, Zhuofan Zheng, Yifan Wang, Yiyao Yu, Xinzhe Ni, Zicheng Lin, Jin Zeng, and Yujiu Yang. Ursu: Understanding and verifying chain-of-thought reasoning in multimodal mathematics. *arXiv preprint arXiv:2501.04686*, 2025.
- [86] Dakota Mahan, Duy Van Phung, Rafael Rafailov, Chase Blagden, Nathan Lile, Louis Castricato, Jan-Philipp Fränken, Chelsea Finn, and Alon Albalak. Generative reward models. *arXiv preprint arXiv:2410.12832*, 2024.

- [87] Karttikeya Mangalam, Raiymbek Akshulakov, and Jitendra Malik. Egoschema: A diagnostic benchmark for very long-form video language understanding. *Advances in Neural Information Processing Systems*, 36:46212–46244, 2023.
- [88] Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*, 2022.
- [89] Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209, 2021.
- [90] Minesh Mathew, Viraj Bagal, Rubén Tito, Dimosthenis Karatzas, Ernest Valveny, and CV Jawahar. Info-graphicvqa. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1697–1706, 2022.
- [91] AI Meta. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, 2025.
- [92] Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. Modeling context between objects for referring expression understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer, 2016.
- [93] Deepak Narayanan, Amar Phanishayee, Kaiyu Shi, Xie Chen, and Matei Zaharia. Memory-efficient pipeline-parallel dnn training. In *International Conference on Machine Learning*, pages 7937–7947. PMLR, 2021.
- [94] Deepak Narayanan, Mohammad Shoeybi, Jared Casper, Patrick LeGresley, Mostafa Patwary, Vijay Korthikanti, Dmitri Vainbrand, Prethvi Kashinkunti, Julie Bernauer, Bryan Catanzaro, et al. Efficient large-scale language model training on gpu clusters using megatron-lm. In *Proceedings of the international conference for high performance computing, networking, storage and analysis*, pages 1–15, 2021.
- [95] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgbd images. In *ECCV*, 2012.
- [96] OpenAI. Gpt-4v(ision) system card. <https://openai.com/index/gpt-4v-system-card/>, 2023. Accessed: 2025-04-23.
- [97] OpenAI. Addendum to gpt-4o system card: 4o image generation. <https://openai.com/index/gpt-4o-image-generation-system-card-addendum/>, 2025.
- [98] openai. Operator, 2025. URL <https://openai.com/index/introducing-operator/>.
- [99] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- [100] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- [101] David Owen. How predictable is language model benchmark performance? *arXiv preprint arXiv:2401.04757*, 2024.
- [102] Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel. Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3170–3180, 2023.
- [103] Chenbin Pan, Yujun Shen, Yujie Wang, Yujing Wang, Yifan Liu, Jiajun Shen, and Yiming Qian. Vlp: Vision language planning for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12345–12354, 2024.
- [104] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *arXiv preprint arXiv:2306.14824*, 2023.
- [105] Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, et al. Ui-tars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*, 2025.

- [106] 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. In *International Conference on Machine Learning*, 2021. URL <https://api.semanticscholar.org/CorpusID:231591445>.
- [107] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [108] Pooyan Rahmazadehgervi, Logan Bolton, Mohammad Reza Taesiri, and Anh Totti Nguyen. Vision language models are blind: Failing to translate detailed visual features into words, 2025. URL <https://arxiv.org/abs/2407.06581>.
- [109] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- [110] Viresh Ranjan, Udbhav Sharma, Thu Nguyen, and Minh Hoai. Learning to count everything. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3394–3403, 2021.
- [111] Christopher Rawles, Sarah Clinckemaillie, Yifan Chang, Jonathan Waltz, Gabrielle Lau, Marybeth Fair, Alice Li, William Bishop, Wei Li, Folawiyo Campbell-Ajala, et al. Androidworld: A dynamic benchmarking environment for autonomous agents. *arXiv preprint arXiv:2405.14573*, 2024.
- [112] Benjamin Recht, Rebecca Roelofs, Ludwig Schmidt, and Vaishaal Shankar. Do imagenet classifiers generalize to imagenet? In *International conference on machine learning*, pages 5389–5400. PMLR, 2019.
- [113] Jonathan Roberts, Mohammad Reza Taesiri, Ansh Sharma, Akash Gupta, Samuel Roberts, Ioana Croitoru, Simion-Vlad Bogolin, Jialu Tang, Florian Langer, Vyas Raina, Vatsal Raina, Hanyi Xiong, Vishaal Udandarao, Jingyi Lu, Shiyang Chen, Sam Purkis, Tianshuo Yan, Wenye Lin, Gyungin Shin, Qiaochu Yang, Anh Totti Nguyen, David I. Atkinson, Aaditya Baranwal, Alexandru Coca, Mikah Dang, Sebastian Dziadzio, Jakob D. Kunz, Kaiqu Liang, Alexander Lo, Brian Pulfer, Steven Walton, Charig Yang, Kai Han, and Samuel Albanie. Zerobench: An impossible visual benchmark for contemporary large multimodal models, 2025. URL <https://arxiv.org/abs/2502.09696>.
- [114] Anna Rohrbach, Marcus Rohrbach, Weijian Qiu, Annemarie Friedrich, Manfred Pinkal, and Bernt Schiele. Coherent multi-sentence video description with variable level of detail. In *German Conference on Pattern Recognition*, 2014.
- [115] ByteDance Seed. Seed-thinking-v1. 5: Advancing superb reasoning models with reinforcement learning. Technical report, Technical report, ByteDance, 2025. URL <https://github.com/ByteDance-Seed/>, 2025.
- [116] ByteDance Seed. Ui-tars-1.5. <https://seed-tars.com/1.5>, 2025.
- [117] Ziyao Shangguan, Chuhan Li, Yuxuan Ding, Yanan Zheng, Yilun Zhao, Tesca Fitzgerald, and Arman Cohan. Tomato: Assessing visual temporal reasoning capabilities in multimodal foundation models. *arXiv preprint arXiv:2410.23266*, 2024.
- [118] Shuai Shao, Zeming Li, Tianyuan Zhang, Chao Peng, Gang Yu, Xiangyu Zhang, Jing Li, and Jian Sun. Objects365: A large-scale, high-quality dataset for object detection. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [119] N Shazeer, A Mirhoseini, K Maziarz, A Davis, Q Le, G Hinton, and J Dean. The sparsely-gated mixture-of-experts layer. *Outrageously large neural networks*, 2017.
- [120] Wei Shen, Guanlin Liu, Zheng Wu, Ruofei Zhu, Qingping Yang, Chao Xin, Yu Yue, and Lin Yan. Exploring data scaling trends and effects in reinforcement learning from human feedback. 2025. URL <https://api.semanticscholar.org/CorpusID:277435161>.
- [121] Wei Shen, Guanlin Liu, Zheng Wu, Ruofei Zhu, Qingping Yang, Chao Xin, Yu Yue, and Lin Yan. Exploring data scaling trends and effects in reinforcement learning from human feedback, 2025. URL <https://arxiv.org/abs/2503.22230>.
- [122] Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems, EuroSys ’25*, page 1279–1297. ACM, March 2025. doi: 10.1145/3689031.3696075. URL <http://dx.doi.org/10.1145/3689031.3696075>.

- [123] Mohammad Shoeybi, Mostafa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-lm: Training multi-billion parameter language models using model parallelism. [arXiv preprint arXiv:1909.08053](#), 2019.
- [124] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In [Proceedings of the IEEE/CVF conference on computer vision and pattern recognition](#), pages 8317–8326, 2019.
- [125] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In [Proceedings of the IEEE conference on computer vision and pattern recognition](#), pages 567–576, 2015.
- [126] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. [Neurocomputing](#), 568:127063, 2024.
- [127] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. [arXiv preprint arXiv:2405.09818](#), 2024.
- [128] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. [arXiv preprint arXiv:2312.11805](#), 2023.
- [129] Gemini Robotics Team, Saminda Abeyruwan, Joshua Ainslie, Jean-Baptiste Alayrac, Montserrat Gonzalez Arenas, Travis Armstrong, Ashwin Balakrishna, Robert Baruch, Maria Bauza, Michiel Blokzijl, et al. Gemini robotics: Bringing ai into the physical world. [arXiv preprint arXiv:2503.20020](#), 2025.
- [130] Kimi Team, Angang Du, Bohong Yin, Bowei Xing, Bowen Qu, Bowen Wang, Cheng Chen, Chenlin Zhang, Chenzhuang Du, Chu Wei, et al. Kimi-vl technical report. [arXiv preprint arXiv:2504.07491](#), 2025.
- [131] Xiaoyu Tian, Junru Gu, Bailin Li, Yicheng Liu, Zhiyong Zhao, Yang Wang, Kun Zhan, Peng Jia, Xianpeng Lang, and Hang Zhao. Drivevlm: The convergence of autonomous driving and large vision-language models. [arXiv preprint arXiv:2402.12289](#), 2024.
- [132] Vernon Y. H. Toh, Yew Ken Chia, Deepanway Ghosal, and Soujanya Poria. The jumping reasoning curve? tracking the evolution of reasoning performance in gpt-[n] and o-[n] models on multimodal puzzles, 2025. URL <https://arxiv.org/abs/2502.01081>.
- [133] 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.
- [134] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikell, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Bin Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [135] Michael Tschanne, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. [arXiv preprint arXiv:2502.14786](#), 2025.
- [136] Borui Wan, Mingji Han, Yiyao Sheng, Yanghua Peng, Haibin Lin, Mofan Zhang, Zhichao Lai, Menghan Yu, Junda Zhang, Zuquan Song, et al. Bytecheckpoint: A unified checkpointing system for large foundation model development. [arXiv preprint arXiv:2407.20143](#), 2024.
- [137] Binghai Wang, Rui Zheng, Lu Chen, Yan Liu, Shihan Dou, Caishuang Huang, Wei Shen, Senjie Jin, Enyu Zhou, Chenyu Shi, et al. Secrets of rlhf in large language models part ii: Reward modeling. [arXiv preprint arXiv:2401.06080](#), 2024.

- [138] Haohan Wang, Songwei Ge, Zachary Lipton, and Eric P Xing. Learning robust global representations by penalizing local predictive power. *Advances in neural information processing systems*, 32, 2019.
- [139] Jiawei Wang, Liping Yuan, Yuchen Zhang, and Haomiao Sun. Tarsier: Recipes for training and evaluating large video description models, 2024. URL <https://arxiv.org/abs/2407.00634>.
- [140] Ke Wang, Junting Pan, Weikang Shi, Zimu Lu, Houxing Ren, Aojun Zhou, Mingjie Zhan, and Hongsheng Li. Measuring multimodal mathematical reasoning with math-vision dataset. *Advances in Neural Information Processing Systems*, 37:95095–95169, 2024.
- [141] 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.
- [142] Weihan Wang, Zehai He, Wenyi Hong, Yean Cheng, Xiaohan Zhang, Ji Qi, Xiaotao Gu, Shiyu Huang, Bin Xu, Yuxiao Dong, et al. Lvbench: An extreme long video understanding benchmark. *arXiv preprint arXiv:2406.08035*, 2024.
- [143] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language models with self-generated instructions. *arXiv preprint arXiv:2212.10560*, 2022.
- [144] Zirui Wang, Mengzhou Xia, Luxi He, Howard Chen, Yitao Liu, Richard Zhu, Kaiqu Liang, Xindi Wu, Haotian Liu, Sadhika Malladi, et al. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *Advances in Neural Information Processing Systems*, 37:113569–113697, 2024.
- [145] Chen Wei, Haoqi Fan, Saining Xie, Chao-Yuan Wu, Alan Yuille, and Christoph Feichtenhofer. Masked feature prediction for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14668–14678, 2022.
- [146] Yanbin Wei, Shuai Fu, Weisen Jiang, Zejian Zhang, Zhixiong Zeng, Qi Wu, James T. Kwok, and Yu Zhang. Gita: Graph to visual and textual integration for vision-language graph reasoning, 2024. URL <https://arxiv.org/abs/2402.02130>.
- [147] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. Longvideobench: A benchmark for long-context interleaved video-language understanding. *Advances in Neural Information Processing Systems*, 37:28828–28857, 2024.
- [148] Penghao Wu and Saining Xie. V*: Guided visual search as a core mechanism in multimodal llms, 2023. URL <https://arxiv.org/abs/2312.14135>.
- [149] Zhiyong Wu, Zhenyu Wu, Fangzhi Xu, Yian Wang, Qiushi Sun, Chengyou Jia, Kanzhi Cheng, Zichen Ding, Liheng Chen, Paul Pu Liang, et al. Os-atlas: A foundation action model for generalist gui agents. *arXiv preprint arXiv:2410.23218*, 2024.
- [150] xAI. Realworldqa: A benchmark for real-world spatial understanding. <https://huggingface.co/datasets/xai-org/RealworldQA>, 2024. Accessed: 2025-04-26.
- [151] Chaojun Xiao, Jie Cai, Weilin Zhao, Guoyang Zeng, Biyuan Lin, Jie Zhou, Zhi Zheng, Xu Han, Zhiyuan Liu, and Maosong Sun. Densing law of llms. *arXiv preprint arXiv:2412.04315*, 2024.
- [152] Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh J Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, et al. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *Advances in Neural Information Processing Systems*, 37:52040–52094, 2024.
- [153] Haomiao Xiong, Zongxin Yang, Jiazu Yu, Yunzhi Zhuge, Lu Zhang, Jiawen Zhu, and Huchuan Lu. Streaming video understanding and multi-round interaction with memory-enhanced knowledge. *arXiv preprint arXiv:2501.13468*, 2025.
- [154] Weiye Xu, Jiahao Wang, Weiyun Wang, Zhe Chen, Wengang Zhou, Aijun Yang, Lewei Lu, Houqiang Li, Xiaohua Wang, Xizhou Zhu, Wenhui Wang, Jifeng Dai, and Jinguo Zhu. Visulogic: A benchmark for evaluating visual reasoning in multi-modal large language models, 2025. URL <https://arxiv.org/abs/2504.15279>.
- [155] Wenyuan Xu, Xiaochen Zuo, Chao Xin, Yu Yue, Lin Yan, and Yonghui Wu. A unified pairwise framework for rlhf: Bridging generative reward modeling and policy optimization. *arXiv preprint arXiv:2504.04950*, 2025.

- [156] Wenyuan Xu, Xiaochen Zuo, Chao Xin, Yu Yue, Lin Yan, and Yonghui Wu. A unified pairwise framework for rlhf: Bridging generative reward modeling and policy optimization, 2025. URL <https://arxiv.org/abs/2504.04950>.
- [157] Zhenhua Xu, Yujia Zhang, Enze Xie, Zhen Zhao, Yong Guo, Kwan-Yee. K. Wong, Zhenguo Li, and Hengshuang Zhao. Drivegpt4: Interpretable end-to-end autonomous driving via large language model, 2024. URL <https://arxiv.org/abs/2310.01412>.
- [158] Tianci Xue, Weijian Qi, Tianneng Shi, Chan Hee Song, Boyu Gou, Dawn Song, Huan Sun, and Yu Su. An illusion of progress? assessing the current state of web agents. *arXiv preprint arXiv:2504.01382*, 2025.
- [159] Chih-Hsuan Yang, Benjamin Feuer, Talukder Jubery, Zi Deng, Andre Nakkab, Md Zahid Hasan, Shivani Chiranjevi, Kelly Marshall, Nirmal Baishnab, Asheesh Singh, et al. Biotrove: A large curated image dataset enabling ai for biodiversity. *Advances in Neural Information Processing Systems*, 37:102101–102120, 2024.
- [160] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2. *arXiv:2406.09414*, 2024.
- [161] Shunyu Yao, Howard Chen, Austin W. Hanjie, Runzhe Yang, and Karthik Narasimhan. Collie: Systematic construction of constrained text generation tasks, 2023. URL <https://arxiv.org/abs/2307.08689>.
- [162] Qinghao Ye, Xianhan Zeng, Fu Li, Chunyuan Li, and Haoqi Fan. Painting with words: Elevating detailed image captioning with benchmark and alignment learning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=636MOnNbPs>.
- [163] Chun-Hsiao Yeh, Chenyu Wang, Shengbang Tong, Ta-Ying Cheng, Rouyu Wang, Tianzhe Chu, Yuexiang Zhai, Yubei Chen, Shenghua Gao, and Yi Ma. Seeing from another perspective: Evaluating multi-view understanding in mllms. *arXiv preprint arXiv:2504.15280*, 2025.
- [164] Licheng Yu, Patrick Poirson, Shan Yang, Alexander C Berg, and Tamara L Berg. Modeling context in referring expressions. In *Computer Vision-ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part II 14*, pages 69–85. Springer, 2016.
- [165] Qiyi Yu, Zheng Zhang, Ruofei Zhu, Yufeng Yuan, Xiaochen Zuo, Yu Yue, Tiantian Fan, Gaohong Liu, Lingjun Liu, Xin Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Guangming Sheng, Yuxuan Tong, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Jinhua Zhu, Jiaze Chen, Jiangjie Chen, Chengyi Wang, Hongli Yu, Weinan Dai, Yuxuan Song, Xiangpeng Wei, Hao Zhou, Jingjing Liu, Wei-Ying Ma, Ya-Qin Zhang, Lin Yan, Mu Qiao, Yonghui Wu, and Mingxuan Wang. Dapo: An open-source llm reinforcement learning system at scale, 2025. URL <https://arxiv.org/abs/2503.14476>.
- [166] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *arXiv preprint arXiv:2308.02490*, 2023.
- [167] Yufeng Yuan, Yu Yue, Ruofei Zhu, Tiantian Fan, and Lin Yan. What’s behind ppo’s collapse in long-cot? value optimization holds the secret, 2025. URL <https://arxiv.org/abs/2503.01491>.
- [168] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556–9567, 2024.
- [169] Xiang Yue, Tianyu Zheng, Yuansheng Ni, Yubo Wang, Kai Zhang, Shengbang Tong, Yuxuan Sun, Botao Yu, Ge Zhang, Huan Sun, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. *arXiv preprint arXiv:2409.02813*, 2024.
- [170] Yu Yue, Yufeng Yuan, Qiyi Yu, Xiaochen Zuo, Ruofei Zhu, Wenyuan Xu, Jiaze Chen, Chengyi Wang, TianTian Fan, Zhengyin Du, Xiangpeng Wei, Xiangyu Yu, Gaohong Liu, Juncai Liu, Lingjun Liu, Haibin Lin, Zhiqi Lin, Bole Ma, Chi Zhang, Mofan Zhang, Wang Zhang, Hang Zhu, Ru Zhang, Xin Liu, Mingxuan Wang, Yonghui Wu, and Lin Yan. Vapo: Efficient and reliable reinforcement learning for advanced reasoning tasks, 2025. URL <https://arxiv.org/abs/2504.05118>.
- [171] 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.

- [172] Di Zhang, Junxian Li, Jingdi Lei, Xunzhi Wang, Yujie Liu, Zonglin Yang, Jiatong Li, Weida Wang, Suorong Yang, Jianbo Wu, et al. Critic-v: Vlm critics help catch vlm errors in multimodal reasoning. [arXiv preprint arXiv:2411.18203](#), 2024.
- [173] Shulai Zhang, Ningxin Zheng, Haibin Lin, Ziheng Jiang, Wenlei Bao, Chengquan Jiang, Qi Hou, Weihao Cui, Size Zheng, Li-Wen Chang, et al. Comet: Fine-grained computation-communication overlapping for mixture-of-experts. [arXiv preprint arXiv:2502.19811](#), 2025.
- [174] Yiyuan Zhang, Handong Li, Jing Liu, and Xiangyu Yue. Explore the limits of omni-modal pretraining at scale. [arXiv preprint arXiv:2406.09412](#), 2024.
- [175] Yilun Zhao, Lujing Xie, Haowei Zhang, Guo Gan, Yitao Long, Zhiyuan Hu, Tongyan Hu, Weiyuan Chen, Chuhan Li, Junyang Song, et al. Mmvu: Measuring expert-level multi-discipline video understanding. [arXiv preprint arXiv:2501.12380](#), 2025.
- [176] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. Judging llm-as-a-judge with mt-bench and chatbot arena. [Advances in Neural Information Processing Systems](#), 36:46595–46623, 2023.
- [177] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL <https://arxiv.org/abs/2306.05685>.
- [178] Junjie Zhou, Yan Shu, Bo Zhao, Boya Wu, Shitao Xiao, Xi Yang, Yongping Xiong, Bo Zhang, Tiejun Huang, and Zheng Liu. Mlvu: A comprehensive benchmark for multi-task long video understanding. [arXiv preprint arXiv:2406.04264](#), 2024.
- [179] Yiyang Zhou, Chenhang Cui, Rafael Rafailov, Chelsea Finn, and Huaxiu Yao. Aligning modalities in vision large language models via preference fine-tuning. [arXiv preprint arXiv:2402.11411](#), 2024.
- [180] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. [arXiv preprint arXiv:1909.08593](#), 2019.

8 Contributions and Acknowledgments

The authors are listed in alphabetical order by their first names. Some names refer to the authors' internal aliases at the company.

Core Contributors

Dong Guo
Faming Wu
Feida Zhu
Fuxing Leng
Guang Shi
Haobin Chen
Haoqi Fan
Jian Wang
Jianyu Jiang
Jiawei Wang
Jingji Chen
Jingjia Huang
Kang Lei
Liping Yuan
Lishu Luo
Pengfei Liu
Qinghao Ye
Rui Qian
Shen Yan
Shixiong Zhao
Shuai Peng
Shuangye Li
Sihang Yuan
Sijin Wu
Tianheng Cheng
Weiwei Liu
Wenqian Wang
Xianhan Zeng
Xiao Liu
Xiaobo Qin
Xiaohan Ding
Xiaojun Xiao
Xiaoying Zhang
Xuanwei Zhang
Xuehan Xiong
Yanghua Peng
Yangrui Chen
Yanwei Li
Yanxu Hu
Yi Lin
Yiyuan Hu
Yiyuan Zhang
Youbin Wu
Yu Li
Yudong Liu
Yue Ling
Yujia Qin

Zanbo Wang
Zhiwu He

Contributors

Aoxue Zhang
Bairen Yi
Bencheng Liao
Can Huang
Can Zhang
Chaorui Deng
Chaoyi Deng
Cheng Lin
Cheng Yuan
Chenggang Li
Chenhui Gou
Chenwei Lou
Chengzhi Wei
Chundian Liu
Chunyuan Li
Deyao Zhu
Donghong Zhong
Feng Li
Feng Zhang
Gang Wu
Guodong Li
Guohong Xiao
Haibin Lin
Haihua Yang
Haoming Wang
Heng Ji
Hongxiang Hao
Hui Shen
Huixia Li
Jiahao Li
Jialong Wu
Jianhua Zhu
Jianpeng Jiao
Jiashi Feng
Jiaze Chen
Jianhui Duan
Jihao Liu
Jin Zeng
Jingqun Tang
Jingyu Sun
Joya Chen
Jun Long
Junda Feng
Junfeng Zhan

Junjie Fang	Weihao Yu
Junting Lu	Wenhai Huang
Kai Hua	Wenjia Zhu
Kai Liu	Wenli Yang
Kai Shen	Wenzhi Wang
Kaiyuan Zhang	Xiang Long
Ke Shen	XiangRui Yin
Ke Wang	Xiao Li
Keyu Pan	Xiaolei Zhu
Kun Zhang	Xiaoying Jia
Kunchang Li	Xijin Zhang
Lanxin Li	Xin Liu
Lei Li	Xinchen Zhang
Lei Shi	Xinyu Yang
Li Han	Xiongcai Luo
Liang Xiang	Xiuli Chen
Liangqiang Chen	Xuantong Zhong
Lin Chen	Xuefeng Xiao
Lin Li	Xujing Li
Lin Yan	Yan Wu
Liying Chi	Yawei Wen
Longxiang Liu	Yifan Du
Mengfei Du	Yihao Zhang
Mingxuan Wang	Yining Ye
Ningxin Pan	Yonghui Wu
Peibin Chen	Yu Liu
Pengfei Chen	Yu Yue
Pengfei Wu	Yufeng Zhou
Qingqing Yuan	Yufeng Yuan
Qingyao Shuai	Yuhang Xu
Qiuyan Tao	Yuhong Yang
Renjie Zheng	Yun Zhang
Renrui Zhang	Yunhao Fang
Ru Zhang	Yuntao Li
Rui Wang	Yurui Ren
Rui Yang	Yuwen Xiong
Rui Zhao	Zehua Hong
Shaoqiang Xu	Zehua Wang
Shihao Liang	Zewei Sun
Shipeng Yan	Zeyu Wang
Shu Zhong	Zhao Cai
Shuaishuai Cao	Zhaoyue Zha
Shuangzhi Wu	Zhecheng An
Shufan Liu	Zhehui Zhao
Shuhan Chang	Zhengzhuo Xu
Songhua Cai	Zhipeng Chen
Tenglong Ao	Zhiyong Wu
Tianhao Yang	Zhuofan Zheng
Tingting Zhang	Zihao Wang
Wanjun Zhong	Zilong Huang
Wei Jia	Ziyu Zhu
Wei Weng	Zuquan Song

Acknowledgments

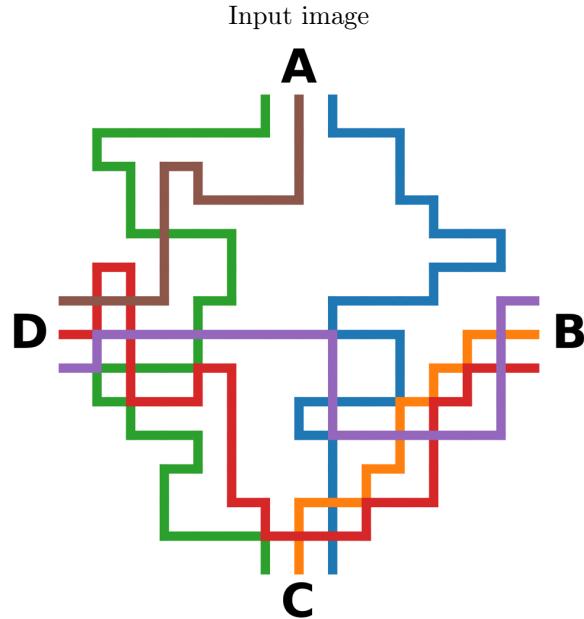
We would like to sincerely thank Allie Guo, Bingyi Kang, Borui Wan, Chaoran Guo, Chenyuan Wang, Chixiang Ma, Fei Xiong, Fu Li, Fuxiang Li, Gaohong Liu, Hongbin Ren, Hongyu Zhu, Jinxin Chi, Juncai Liu, Kaihua Jiang, Kayden, Lei Zuo, Lianke Qin, Lingjun Liu, Liyang Liu, Minchao Wang, Mingji Han, Mofan Zhang, Pengyuan Zhao, Qianli Ma, Qifan Yang, Qinlong Wang, Shibiao Nong, Tong Zhou, Weiqiang Lou, Xiangpeng Wei, Xiangrui Yin, Xiangtai Li, Xiaokai Li, Xudong Sun, Xun Wang, Yaowei Zheng, Yixin Wu, Yingping Zhang, Yun Zhang, Yuwen Tang, Zhe Nan, Zhelun Shi, Zheng Zhong, Zhenyuan Yang, Zhi Zhang, Zhongjia Wei, Zhuolin Zheng, Zilong Zhou, Ziqian Wei, Ziwen Xu, Zixin Chen, Ziyuan Feng, Zuquan Song for their insightful discussions and unwavering support. Their valuable input has been instrumental in advancing Seed1.5-VL's development, evaluation, defect analysis, and exploration of future research directions.

Appendix

A Qualitative examples

This section presents a selection of qualitative examples illustrating the capabilities of the Seed1.5-VL model through responses generated from various prompts. Examples demonstrating the model’s visual reasoning ability are provided in [figures 9 to 15](#). [Figure 16](#) shows a case of model solving visual puzzles, such as “Find-the-Differences”. The model’s proficiency in solving geometry problems is showcased in [figure 17](#). [Figure 18](#) illustrates model’s ability of accurately counting objects presented in a crowded scene. Document understanding and OCR capabilities are highlighted in [figures 21 and 22](#). Beyond standard image understanding, our model is also capable of 3D spatial understanding from a single image, exemplified by sorting objects based on their depth relative to the camera as shown in [figure 19](#). Extending beyond static images, [figure 20](#) presents an example of the model localizing events within a video based on user queries. [Figure 24](#) showcases an instance of image-conditioned creative writing generated by our model. Finally, failure cases illustrating the current limitations observed in VLMs, including issues related to spatial imagination, hallucination, and combinatorial search, are presented in [figures 25 to 28](#).

A.1 Reasoning Cases: Visual Reasoning



Prompts: Count the one-colored routes that go from C to A. Answer with a number in curly brackets e.g. {3}

Model Responses: <think>So, let's try to figure out how many one-colored routes go from C to A. First, I need to identify each color and check if there's a continuous path from C to A for each color.

Let's list the colors: green, brown (or maybe dark red/brown), blue, red, orange, purple. Wait, let's check each color's path. First, green: Let's see if there's a green path from C to A. Starting at C, follow green lines. Does it reach A? Let's trace. Green lines: from C, there's a green segment, then up, etc. Wait, need to check continuity.

Brown: The brown path. Does brown go from C to A? Let's check. Brown starts... Wait, A is at the top, so brown's path: does it connect C to A?

Blue: Blue path. From C, blue lines—do they go up to A?

Red: Red path. From C, red lines—do they reach A?

Orange: Orange path. From C, orange lines—do they reach A?

Purple: Purple path. From C, purple lines—do they reach A?

...

Let's look at the image again.

Wait, the colors are:

- Green: several segments forming a path from C to A.

- Blue: several segments forming a path from C to A.

Are there any other colors? Let's check each color's continuity.

Green: Yes, continuous from C to A.

Blue: Yes, continuous from C to A.

Brown: Starts at A, goes down, but does it reach C? No, because the brown path ...

So that's two. So the answer is 2?</think>

{2}

Figure 9 A case demonstrating one of the visual-centric thinking strategies that emerged during LongCoT RL.

A.2 Reasoning Cases: Geolocation Prediction

Input image



Prompts: You are participating in a geolocation challenge. Based on the provided image:

1. Carefully analyze the image for clues about its location (architecture, signage, vegetation, terrain, etc.)
2. Think step-by-step about what continent, country, and city this is likely to be in and why
3. Estimate the approximate latitude and longitude based on your analysis Take your time to reason through the evidence.

Your final answer MUST include these five lines somewhere in your response:

continent: [continent name]

country: [country name]

city: [city name]

lat: [latitude as a decimal number]

lng: [longitude as a decimal number]

You can provide additional reasoning or explanation, but these three specific lines MUST be included.

Model Responses: <think>So, let's try to figure out the location. First, look at the large building on the left—it looks like a multi-level parking structure or part of a stadium. The stadium-like structure with the big scoreboard and the tall light towers (used in baseball stadiums, maybe) is a clue. The number "2" on a sign might be a gate number. The vegetation: there are trees and hedges, typical of an urban area.

Now, the architecture style—modern stadium, maybe in the US. Let's think about major stadiums. The structure with the grid and the flags—maybe Guaranteed Rate Field (formerly U.S. Cellular Field), home of the Chicago White Sox. Let's check. The stadium's exterior, the parking structure, and the surrounding area. Chicago, Illinois, USA.

To confirm, the style of the stadium, the urban setting, and the type of infrastructure. Then, latitude and longitude for Guaranteed Rate Field: approximately 41.8316° N, 87.6348° W. ... </think>

...

continent: North America

country: United States

city: Chicago

lat: 41.8316

lng: -87.6348

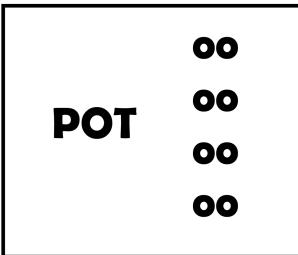
Figure 10 A case demonstrating the "analyze visual details before recognizing a location" strategy that emerged during LongCoT RL.

A.3 Visual Reasoning: Solving Rebus Puzzles



Prompts: can you solve this Rebus puzzle?

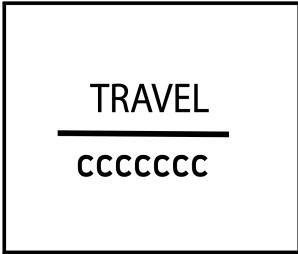
Model Responses: <think>...</think>The Rebus puzzle features the word "FAST" with cracked lines, suggesting it is being "broken." Combining "break" (from the visual of cracks, implying breaking) with "fast" results in the word **breakfast**



Prompts: can you solve this Rebus puzzle?

Model Responses: <think>...</think>The Rebus puzzle features the word "POT" on the left and eight circles (four pairs of two circles each, resembling "OO") on the right. Combining "POT" with "8 O's" (where "8 O's" phonetically sounds like "ato" when spoken quickly, as in "eight O's" → "ato") forms the word **Potato**. This wordplay uses the visual of eight circles (O's) to create the suffix "ato" after "POT".

Answer: Potato



Prompts: can you solve this Rebus puzzle?

Model Responses: <think>...</think>The rebus puzzle consists of "TRAVEL" above a line (representing a hyphen or connection) and seven "C"s. Each "C" sounds like "sea," so seven Cs are "seven seas." Combining these, the solution is "**Travel the seven seas**".

Figure 11 Three examples of Seed1.5-VL solving Rebus puzzles. Our model gave the correct reasoning and answers to the first two. For the third one, although the model gave an incorrect answer it still output logical reasoning. The correct answer is "travel overseas". We omit thinking tokens for brevity.

A.4 Visual Reasoning: Emoji Quiz



Prompts: Which movie do these emojis remind you of?

Model Responses: <think>...</think>The emojis—a tropical island (representing isolation/stranded setting) and a volleyball (referencing the character “Wilson” from the movie)—remind us of “**Cast Away**”. The film features a protagonist stranded on a deserted island, forming a unique bond with a volleyball named Wilson, which matches these symbols perfectly.



Prompts: Which movie do these emojis remind you of?

Model Responses: <think>...</think>The emojis—a Italian flag (representing Italian heritage), a horse (referencing “Stallion,” as in Rocky’s nickname “The Italian Stallion”), and a boxing glove (boxing, Rocky’s sport)—collectively evoke “**Rocky**”. Rocky Balboa, the Italian-American boxer known as “The Italian Stallion,” perfectly matches these symbols.



Prompts: Which movie do these emojis remind you of?

Model Responses: <think>...</think>The emojis include a UFO (alien space-craft), landmarks like the Statue of Liberty (New York) and Golden Gate Bridge (San Francisco), cityscapes, and explosions—all key elements from “**Independence Day**”, where aliens attack major cities with UFOs, resulting in widespread destruction. These symbols align perfectly with the movie’s plot and iconic scenes.

Figure 12 Three illustrative instances of Seed1.5-VL successfully identifying movie names from their corresponding emoji inputs. These examples highlight the model’s ability to recognize and semantically link emoji icons to movie plot information.

A.5 Visual Reasoning: Word Game I

LION JOB ENDING
TOOK MY ORE
CRIMSON ELBOW
RAPID DISARM

Prompts: These seemingly random phrases hide the names of eight capital cities around the world. Each phrase contains the name of two cities. Can you unscramble them and figure out what they are?

Model Responses: <think>...</think> The phrases contain the following capital cities, each pair from a line:

1. **LION JOB ENDING** → **London** (L-O-N-D-O-N) and **Beijing** (B-E-I-J-I-N-G).
 2. **TOOK MY ORE** → **Tokyo** (T-O-K-Y-O) and **Rome** (R-O-M-E).
 3. **CRIMSON ELBOW** → **Berlin** (B-E-R-L-I-N) and **Moscow** (M-O-S-C-O-W).
 4. **RAPID DISARM** → **Paris** (P-A-R-I-S) and **Madrid** (M-A-D-R-I-D).
- These are all capital cities (London, Beijing, Tokyo, Rome, Berlin, Moscow, Paris, Madrid), each hidden in the letter combinations of the phrases.

Figure 13 An example of Seed1.5-VL playing a word game. The model correctly recognize letters from the image, and retrieve capital cities that are spelled using those letters. This exemplifies the model’s integrated OCR, reasoning, and knowledge retrieval capabilities.

A.6 Visual Reasoning: Word Game II

Input image

Alice's Guess	Betty's Answer
HUNKY	2
JUNKS	0
WHISK	2
CENTS	2
STINK	0

Prompt:

Alice and Betty are playing “guess my word”. Betty thinks of a five-letter word (all distinct letters), and Alice guesses. For each guess, Betty tells Alice how many letter from her guess are right. What is Betty’s word?

Model Responses: <think>...</think>To determine Betty’s five-letter word (with all distinct letters), we analyze the clues from Alice’s guesses:

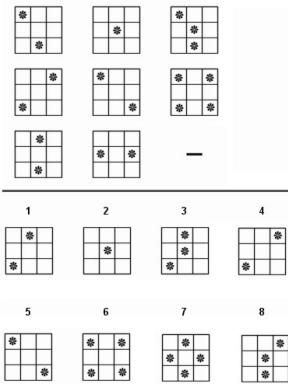
1. **Eliminate Letters Using Guesses with 0 Correct Letters**
 - When Alice guesses **JUNKS** and gets 0, it means the letters **J, U, N, K, S** are not in Betty’s word.
 - When Alice guesses **STINK** and gets 0, it means the letters **S, T, I, N, K** are not in Betty’s word.
 - By taking the union of the eliminated letters from these two guesses, we rule out **J, U, N, K, S, T, I**.
2. **Identify Valid Letters from Guesses with 2 Correct Letters**
 - For the guess **HUNKY** (which has 2 correct letters), after eliminating the invalid letters (**U, N, K**), the valid letters that could be correct are **H, Y**.
 - For the guess **WHISK** (also with 2 correct letters), after eliminating the invalid letters (**I, S, K**), the valid letters that could be correct are **W, H**.
 - For the guess **CENTS** (with 2 correct letters), after eliminating the invalid letters (**N, T, S**), the valid letters that could be correct are **C, E**.
 - Combining these valid letters, we get the set **H, Y, W, C, E**, which has exactly 5 distinct letters, satisfying the condition of a five-letter word.
3. **Form the Word**
 - We can form the word **CHEWY** using the letters **C, H, E, W, Y**.
 - Let’s verify this word against all the guesses:
 - **HUNKY**: The letters H and Y match, so there are 2 correct letters.
 - **JUNKS**: None of the letters J, U, N, K, S match, so there are 0 correct letters.
 - **WHISK**: The letters W and H match, so there are 2 correct letters.
 - **CENTS**: The letters C and E match, so there are 2 correct letters.
 - **STINK**: None of the letters S, T, I, N, K match, so there are 0 correct letters.

All the conditions are satisfied.

Answer: CHEWY

Figure 14 An example of solving a word game requiring multi-step reasoning. First, eliminate invalid letters. Second, identify valid ones, and finally, verify the derived letters with previous guesses.

A.7 Visual Reasoning: Visual Pattern Recognition



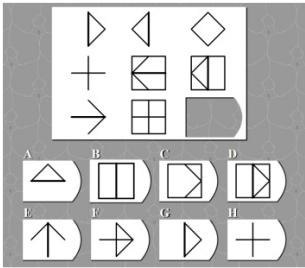
Prompts: here is a visual reasoning puzzle. what is the last pattern indicated by “-”? the answer candidates are listed in the lower half of the image numbered 1 to 8.

Model Responses: <think>...</think>The puzzle involves identifying the pattern in each row, where the third grid is the union (sum) of the first two.

- **First row:** $2^* + 1^* = 3^*$ (matches third grid).
- **Second row:** $2^* + 2^* = 4^*$ (matches third grid, four corners).
- **Third row:** First two grids each have 2^* , so the third must have 4^* .

Analyzing candidates, **Candidate 7** has four asterisks arranged in a column (top-center, middle-center, bottom-center, and another? No, wait—correcting: no, finally, recognizing the union of third row’s first two grids (2^* each) results in 4^* , and **Candidate 7** matches this pattern (four asterisks in a vertical column, aligning with the combined positions).

Answer: 7



Prompts: here is a visual reasoning puzzle. what is the missing pattern listed as the last element in the 3rd row? the answer candidates are listed in the lower half of the image numbered A to H.

Model Responses: <think>...</think>To solve the puzzle, analyze each row for patterns. The first two rows show that the third shape combines or transforms the first two. In the third row:

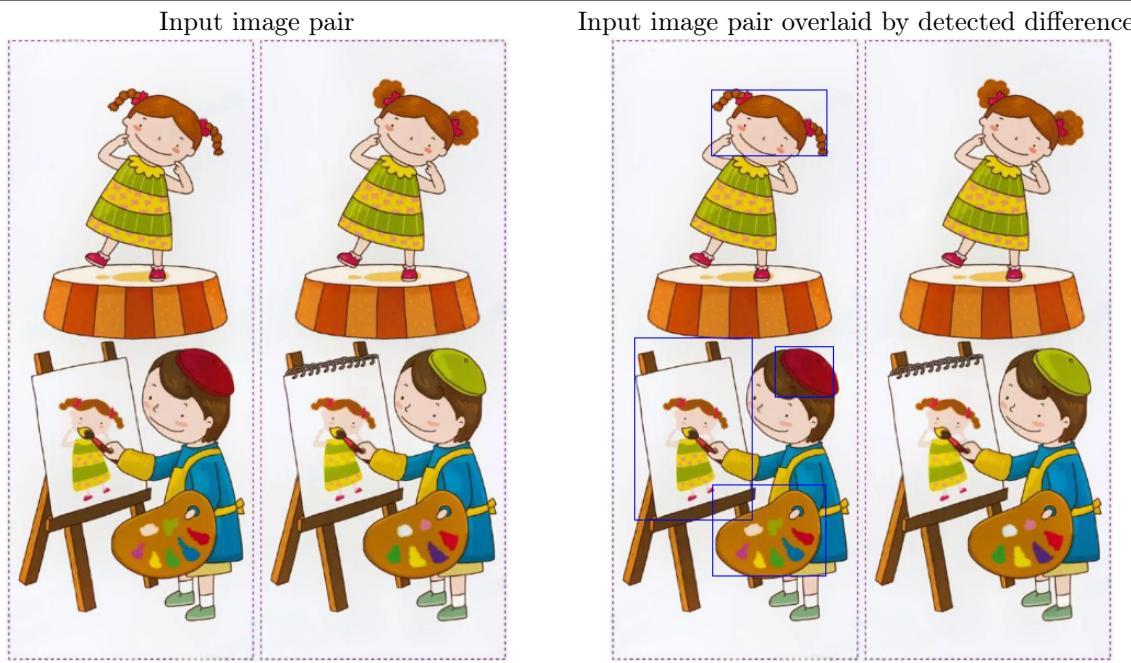
- **First shape:** Arrow (right, stem).
- **Second shape:** Grid (square divided into four).

The missing third shape (candidate D) matches the pattern: it integrates the arrow’s stem (horizontal line) with the grid’s structure, forming a square with a diagonal and stem, consistent with the visual logic of combining row elements.

Answer: D

Figure 15 Two examples of Seed1.5-VL using its visual inductive reasoning ability to correctly predict the missing patterns.

A.8 Visual Puzzles: Find the Differences



Prompt: 请仔细对比两张图，找出它们的不同之处。

针对每一处差异，请在第一张图中用边界框标记对应位置。每个差异区域使用**<bbox>**格式标注,格式为:

<bbox>左上角x坐标左上角y坐标右下角x坐标右下角y坐标</bbox>

要求:

- 所有坐标值为0-999的整数,表示像素位置

- 左上角坐标必须小于右下角坐标

Model response: <think>...</think><bbox>202 82 432 188</bbox> - 第一张图中女孩的双麻花辫

<bbox>329 494 445 576</bbox> - 第一张图中画家的红色帽子

<bbox>49 480 282 773</bbox> - 第一张图中无螺旋装订的画板

<bbox>204 717 429 863</bbox> - 第一张图中调色板的颜色 (包含蓝色、紫色等, 与第二张图不同)

Translated prompt: Please carefully compare the two images and identify the differences between them.

For each difference, mark the corresponding location with bounding boxes on the first image. Annotate each differing region using the **<bbox>** tag in the following format: **<bbox>top-left-x top-left-y bottom-right-x bottom-right-y</bbox>**.

Requirements:

- All coordinate values must be integers between 0 and 999, representing pixel positions.

- The top-left coordinates must be smaller than the bottom-right coordinates.

Translated model response: <think>...</think><bbox>202 82 432 188</bbox> - Girl's double braids in the first image

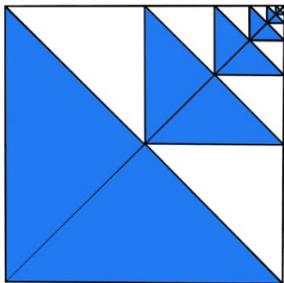
<bbox>329 494 445 576</bbox> - Painter's red hat in the first image

<bbox>49 480 282 773</bbox> - Regular sketchpad without spiral binding in the first image

<bbox>204 717 429 863</bbox> - Palette colors in the first image (containing blue, purple, etc., differing from the second image)

Figure 16 An example of model playing “Find-the-differences” game. The original prompt and model response are in Chinese. The English translation are shown in gray.

A.9 Geometry



Prompts: The square is filled in with a pattern of triangles that become smaller and smaller until they are infinitely small. what fraction of the square is shaded blue?

Model Responses: <think>...</think>The shaded area forms a geometric series:

1. **Largest segment:** The two large blue triangles in the bottom-left occupy $\frac{1}{2}$ of the square.

2. **Subsequent segments:** In the top-right quadrant, each blue triangle is $\frac{1}{4}$ the size of the previous shaded region (e.g., $\frac{1}{8}, \frac{1}{32}, \frac{1}{128}, \dots$).

This is a geometric series with:

- First term $a = \frac{1}{2}$
- Common ratio $r = \frac{1}{4}$

Using the formula for the sum of an infinite geometric series $S = \frac{a}{1-r}$:

$$S = \frac{\frac{1}{2}}{1-\frac{1}{4}} = \frac{\frac{1}{2}}{\frac{3}{4}} = \frac{2}{3}$$

Answer: $\frac{2}{3}$

Prompts: the area of the blue triangle is 12 and the area of the gray triangle is 9. what is the area of the orange rectangle?

Model Responses: <think>...</think>To find the area of the orange rectangle, we analyze the similar triangles (blue and gray) and use their areas:

1. **Blue Triangle:**

$$\text{Area} = \frac{1}{2} \times \text{base} \times \text{height} = 12 \\ \frac{1}{2} \times 3 \times H = 12 \Rightarrow H = 8 \text{ (height of blue triangle).}$$

2. **Gray Triangle:**

$$\text{Its height is } H - 2 = 8 - 2 = 6 \text{ (since the rectangle's height is 2).} \\ \text{Area} = \frac{1}{2} \times \text{base} \times 6 = 9 \Rightarrow \text{Base} = \frac{9 \times 2}{6} = 3.$$

3. **Orange Triangle:**

The rectangle has **length=3** (matching the gray triangle's base) and **width=2**.

$$\text{Area} = 3 \times 2 = 6$$

Answer: 6

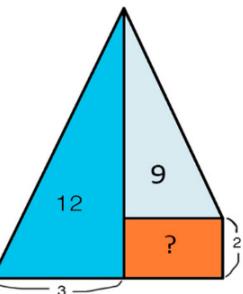
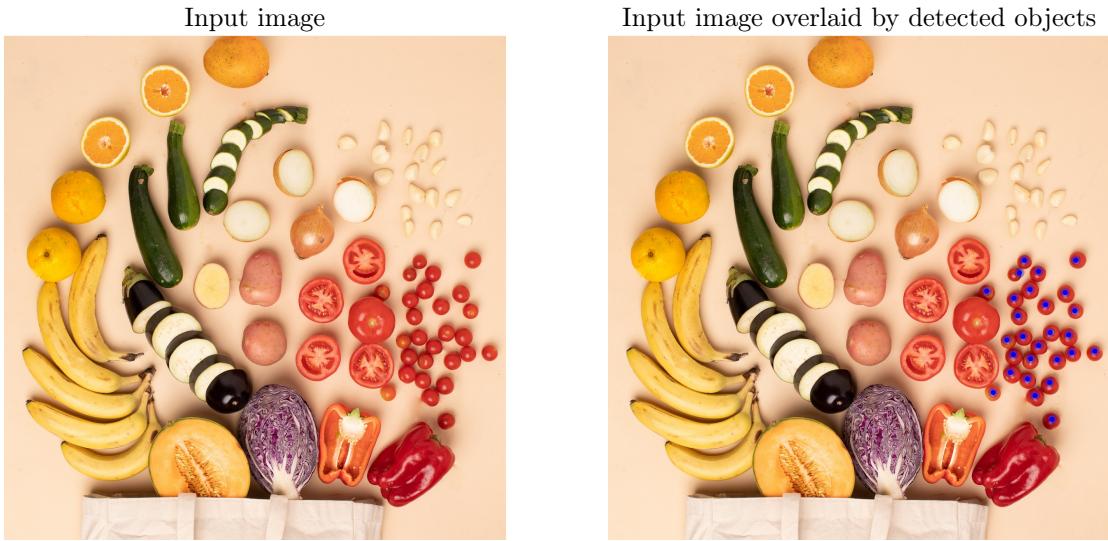


Figure 17 Two examples of Seed1.5-VL solving geometry problems, which requires multi-step reasoning.

A.10 Counting in a complex scene

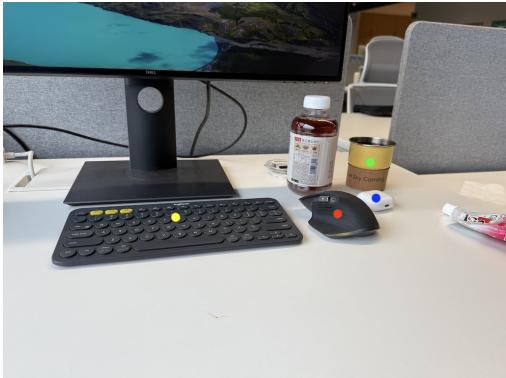


Prompt: How many cherry tomatoes are there in the picture? Point them out first in <point>x y</point> format.

Model response: <point>766 708</point><point>818 471</point><point>828 446</point> <point>856 468</point><point>839 504</point><point>807 521</point><point>815 556</point><point>870 534</point><point>909 510</point><point>930 446</point><point>928 546</point><point>753 507</point><point>793 604</point><point>825 597</point><point>879 589</point><point>916 594</point><point>806 633</point><point>840 641</point><point>856 615</point><point>893 643</point><point>922 629</point><point>968 626</point><point>800 668</point><point>834 681</point><point>878 688</point><point>849 714</point><point>883 761</point> There are 27 cherry tomatoes in the picture

Figure 18 An example of Seed1.5-VL recognizing queried objects and counting their occurrences. Seed1.5-VL first detects the objects based on the user's prompt and then answers the query using the number of detected objects.

A.11 Spatial Understanding: Depth Sorting

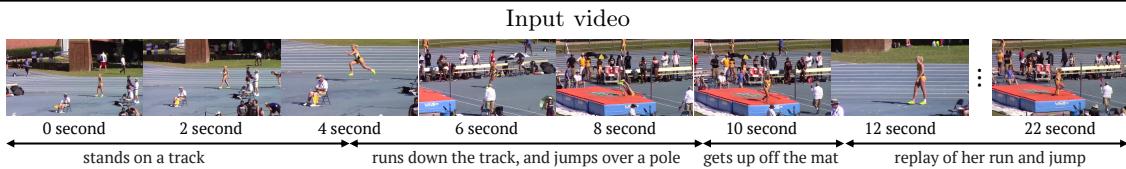


Prompts: There are several points with different colors in the image: point1: blue, point2: yellow, point3: green, point4: red. Each point represents an object. Please sort the object represented by these points in a close to far order and give their name.

Model Responses: point2: yellow (key-board), point4: red (mouse), point1: blue (earphone), point3: green (cup)

Figure 19 An example showing model's ability to estimate objects' depth ordering from a single image.
Source: photo taken by an author from the ByteDance Seed team.

A.12 Video Temporal Grounding



Prompt: Please watch this video and find out all key events in this video, and output the events along with the start/end timestamps.

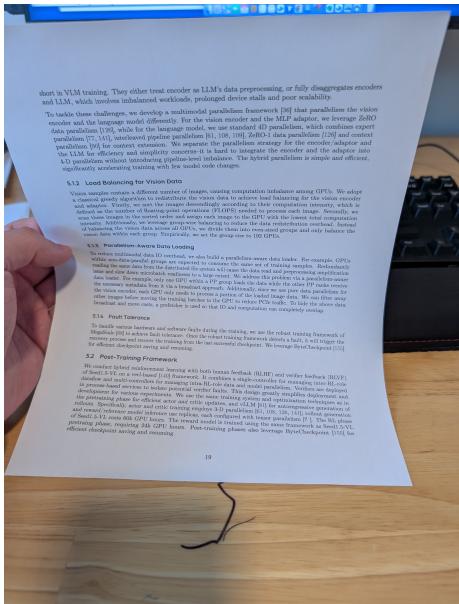
Model response: 0.0–3.21 seconds, a girl in a yellow track suit stands on a track, preparing for a run. 3.32–10.02 seconds, she runs down the track, and jumps over a pole, landing on a mat. 10.13–12.16 seconds, she gets up off the mat, and walks away. 12.27–22.36 seconds, a replay of her run and jump is shown.

Prompt: When does the woman jump? Please output the exact timestamps. If the woman jumps several times, please output the timestamps respectively.

Model response: 6–8 seconds, 16–22 seconds.

Figure 20 The first query shows the model's ability to summarize/segment key events in a video (e.g., dense video captioning), and the second query shows model's ability of localizing an event in a video (e.g., video grounding).

A.13 OCR Parsing and Document Understanding



Prompts: Based on the paper, what framework does the author use to achieve fault tolerance?

Model Responses: According to the paper, to handle various hardware and software faults during the training, the author uses the robust training framework of MegaScale to achieve fault tolerance.

Prompts: Based on the paper, which one takes longer to train, RL phase or reward model?

Model Responses:

1. First, find the relevant information in the paper:
 - The paper states: “The RL phase of Seed1.5-VL costs 60k GPU hours. The reward model is trained using the same framework as Seed1.5-VL pretraining phase, requiring 24k GPU hours.”
 2. Then, compare the GPU-hour requirements:
 - The RL phase requires 60k GPU hours, and the reward model requires 24k GPU hours.

So, the RL phase takes longer to train.

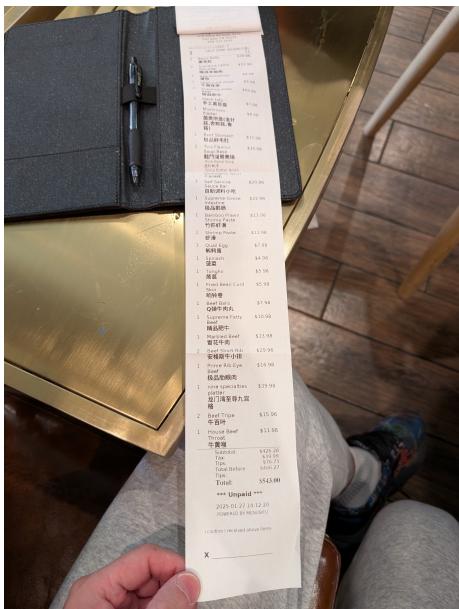
Prompts: Can you repeat section 5.1.4 but skip the first sentence?

Model Responses: Once the robust training framework detects a fault, it will trigger the recovery process and recover the training from the last successful checkpoint. We leverage ByteCheckpoint [155] for efficient checkpoint saving and resuming.

Figure 21 Demonstration of Seed1.5-VL’s ability of accurately extracting texts from an image, understanding the user’s query, and retrieving the answer from the document. The final example also highlights model’s instruction following capability (e.g., “skip the first sentence”).

Source: a printout of an earlier draft of the technical report.

A.14 Multilingual OCR Parsing



Prompts: can you extract the dish names, quantities, and their prices into a markdown table?

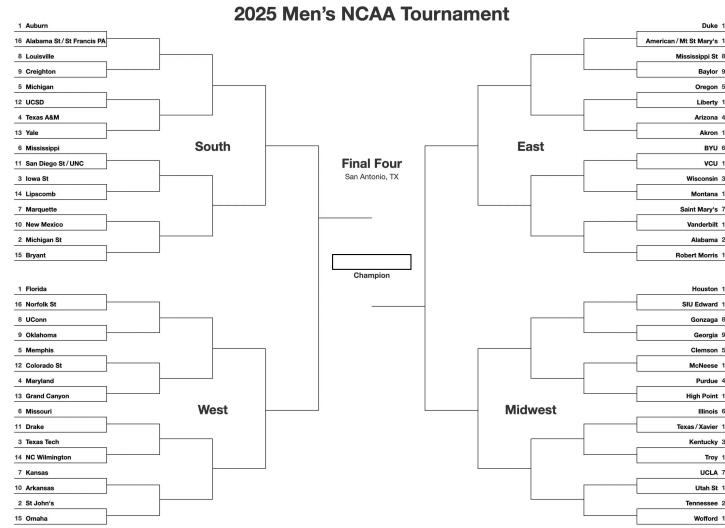
Model Responses:

Dish Name	Quantity	Price
Black Belly (黑毛肚)	2	\$39.96
Supreme Lamb Shoulder (精选羊肩肉)	2	\$23.96
Egg Dumplings (蛋饺)	1	\$8.98
Bean Curd Sheet (干煎豆皮)	1	\$5.98
Supreme Kobe Beef (极品和牛)	2	\$59.96
Black Tofu (手工黑豆腐)	1	\$7.98
Mushroom Platter (菌类拼盘(金针菇,杏鲍菇,香菇))	1	\$8.98
Beef Stomach (极品鲜毛肚)	1	\$11.98
Two Flavour Soup Base (龙门湾鸳鸯锅)	1	\$15.98
Self Service Sauce Bar (自助调料小吃)	7	\$20.86
Supreme Goose Intestine (极品鹅肠)	1	\$22.98
Bamboo Prawn Shrimp Paste (竹荪虾滑)	1	\$13.98
Shrimp Paste (虾滑)	1	\$11.98
Quail Egg (鹌鹑蛋)	1	\$7.98
Spinach (菠菜)	1	\$4.98
Tongho (茼蒿)	1	\$5.98
Fried Bean Curd Skin (响铃卷)	1	\$5.98
Beef Balls (Q弹牛肉丸)	1	\$7.98
Supreme Fatty Beef (精品肥牛)	1	\$10.98
Marbled Beef (雪花牛肉)	1	\$13.98
Beef Short Rib (安格斯牛小排)	2	\$29.96
Prime Rib Eye Beef (极品肋眼肉)	1	\$16.98
nine specialties platter (龙门湾至尊九宫格)	1	\$39.98
Beef Tripe (牛百叶)	2	\$15.96
House Beef Throat (牛黄喉)	1	\$11.98

Figure 22 An example of extracting multilingual texts from a long receipt.

Source: photo taken by an author from the ByteDance Seed team.

A.15 Generate Code for a Diagram of Novel Format



Prompts: can you generate the Mermaid code that produce the bracket for ONLY the South region?

Model Responses: <think>...</think>Here's the Mermaid code for the South Region bracket:

flowchart TD

```

S1[1\nAuburn] --> S1_16[1v16\nAuburn vs Alabama St/St Francis PA]
S16[16\nAlabama St/St Francis PA] --> S1_16
S8[8\nLouisville] --> S8_9[8v9\nLouisville vs Creighton]
S9[9\nCreighton] --> S8_9
S5[5\nMichigan] --> S5_12[5v12\nMichigan vs UCSD]
S12[12\nUCSD] --> S5_12
S4[4\nTexas A&M] --> S4_13[4v13\nTexas A&M vs Yale]
S13[13\nYale] --> S4_13
S6[6\nMississippi] --> S6_11[6v11\nMississippi vs San Diego St/UNC]
S11[11\nSan Diego St/UNC] --> S6_11
S3[3\nIowa St] --> S3_14[3v14\nIowa St vs Lipscomb]
S14[14\nLipscomb] --> S3_14
S7[7\nMarquette] --> S7_10[7v10\nMarquette vs New Mexico]
S10[10\nNew Mexico] --> S7_10
S2[2\nMichigan St] --> S2_15[2v15\nMichigan St vs Bryant]
S15[15\nBryant] --> S2_15

S1_16 --> S2nd1[2nd Round\nWinner 1v16 vs 8v9]
S8_9 --> S2nd1
S5_12 --> S2nd2[2nd Round\nWinner 5v12 vs 4v13]
S4_13 --> S2nd2
S6_11 --> S2nd3[2nd Round\nWinner 6v11 vs 3v14]
S3_14 --> S2nd3
S7_10 --> S2nd4[2nd Round\nWinner 7v10 vs 2v15]
S2_15 --> S2nd4

S2nd1 --> S3rd1[Elite 8\nWinner 2nd1 vs 2nd2]
S2nd2 --> S3rd1
S2nd3 --> S3rd2[Elite 8\nWinner 2nd3 vs 2nd4]
S2nd4 --> S3rd2

S3rd1 --> SWinner[South Winner\nWinner 3rd1 vs 3rd2]
S3rd2 --> SWinner

```

Figure 23 An example of model's several capabilities, such as OCR, coding, instruction following, working together to generate Mermaid code for a portion of the diagram. Summarization texts are omitted from the response. Paste the code into the [Mermaid Live Editor](#) to visualize the South region bracket.

A.16 Image-conditioned Creative Writing

Prompts: write a romantic poem in English inspired by this picture.

Model Responses: <think>...</think>

Fingers Entwined



Beneath the hum of a world that rushes by,
Our hands meet---soft, a silent lullaby.
Your palm, a map of stories yet untold,
My fingers trace the edges, brave and bold.

Rings gleam, twin sparks in this quiet space,
Promises etched in silver's tender embrace.
Each line, each crease, a secret only we share---
Time slows, as if love itself holds its breath here.

Remember when our shadows first aligned?
Now skin to skin, no distance left to find.
These hands, they'll carry dawns and midnight's hush,
Weave years into a tapestry of "us."

For in this touch, a universe unfolds---
Heartbeats syncing, futures yet to hold.
Two hands, one song, forever intertwined---
A language only love could ever design.

The poem mirrors the intimacy of clasped hands, weaving themes of connection, memory, and enduring commitment, with imagery of touch and time to echo the quiet romance in the photograph.

Figure 24 An example of image-conditioned creative text generation. The model processes an input image and generates a poem semantically related to the visual content.

Source: picture taken by an author from the ByteDance Seed team.

A.17 Failure Cases: 3D Spatial Imagination

Prompts: if we move the dice along the arrow directions, what will be the number on top?

Model Responses: <think>...</think>To determine the number on top after moving the dice along the arrow directions, we analyze the initial configuration and movements:

1. **Initial Dice Faces:**

- Front: 1 dot (opposite back: 6 dots, $1 + 6 = 7$).
- Right: 2 dots (opposite left: 5 dots, $2 + 5 = 7$).
- Top: 5 dots (opposite bottom: 2 dots, $5 + 2 = 7$).

2. **Movements:**

- **Down Arrow:** Rolling the die downward (front to bottom) does not affect the top face directly (since top/bottom and front/back are separate axes).
- **Right Arrow:** Rolling the die rightward (right to front) also does not change the top face (top remains independent of right/left movements in this sequence).

Thus, the top face (initially 5 dots) remains unchanged after the specified movements.

Answer: 5

Prompts: View the blocks on the left from the orange point and direction of the arrow. Which image on the right is the correct view?

Model Responses: <think>...</think>To determine the correct front view from the orange point and arrow direction, analyze the vertical height of each column in the 3D structure:

- **Leftmost column:** 1 cube (bottom layer only).
- **Second column:** 2 cubes (two layers).
- **Third column (central, tallest):** 3 cubes (three layers).
- **Fourth column:** 2 cubes (two layers).

Matching this height distribution (1, 2, 3, 2) to the options:

- **A:** Incorrect (leftmost column too tall).
- **B:** Incorrect (symmetric, not matching the structure).
- **C:** Incorrect (heights don't align).
- **D:** Correct (heights 1, 2, 3, 2 match the front view).

Answer: D

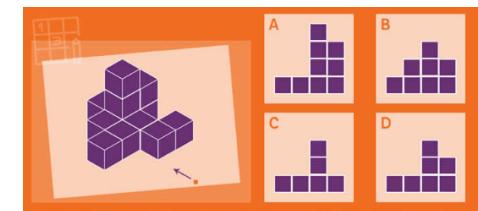


Figure 25 Two examples where Seed1.5-VL, Gemini 2.5 Pro, and OpenAI o4-mini all failed to produce the correct reasoning. In the Seed1.5-VL case, the reasoning is incorrect despite a correct final answer. This highlights a limitation of contemporary VLMs in tasks requiring robust 3D spatial reasoning.

A.18 Failure Cases: Hallucination (Knowledge Prior)

Prompts: can you generate the DOT code that renders a similar graph to this one?

Model Responses: <think>...</think>Here's the DOT code to render a graph similar to your Transformer layer diagram:

```
digraph TransformerLayer {
    rankdir=TB;
    node [fontname="Sans-Serif"];

    // Define nodes
    input [label="x_{l+1}", shape=rectangle];
    ln1 [label="LN", shape=rectangle];
    mha [label="MHA", shape=rectangle];
    add1 [label="+", shape=circle];
    ln2 [label="LN", shape=rectangle];
    ffn [label="FFN", shape=rectangle];
    add2 [label="+", shape=circle];
    output [label="x_{l+1}", shape=rectangle];

    // Define edges
    input -> ln1;
    ln1 -> mha;
    mha -> add1;
    input -> add1; // Residual connection (first skip)

    add1 -> ln2;
    ln2 -> ffn;
    ffn -> add2;
    add1 -> add2; // Residual connection (second skip)

    add2 -> output;

    // Optional styling (you can adjust colors/shapes)
    edge [arrowhead="vee", arrowsize=0.8];
    node [style=filled, fillcolor=white, penwidth=1];
    add1, add2 [fillcolor=lightblue];
}
```

...

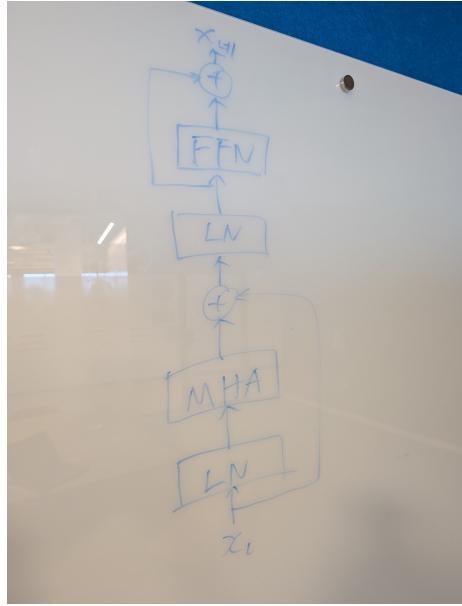
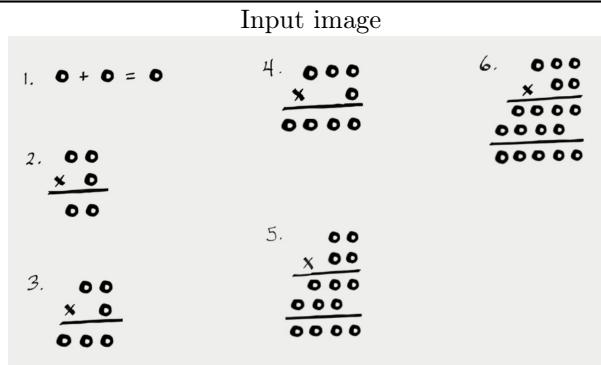


Figure 26 A case study illustrating the inability of Seed1.5-VL, Gemini 2.5 Pro, and OpenAI o4-mini to generate faithful DOT code when presented with a non-standard FFN residual connection. The drawing shows the network structure, specifically the atypical residual path. All models exhibited a tendency to generate DOT code corresponding to a standard residual connection.

A.19 Failure Cases: Combinatorial Search I



Prompts: In this puzzle, each circle is a prime number (2, 3, 5 or 7), solve puzzle #6 which only has one solution.

Model Responses: <think>...</think> To solve Puzzle #6 where each circle represents a prime number (2, 3, 5, or 7), we analyze the multiplication:

1. Prime Digit Constraints

Each digit in the multiplicand (three-digit number), multiplier (two-digit number), and the product must be among 2, 3, 5, or 7. We start by looking at the units place of the multiplication. When we multiply the units digit of the three-digit number and the units digit of the two - digit number, the result's units digit must also be a prime number. Through testing different combinations, we find that the units digit of the three-digit number and the units digit of the two-digit number are likely 5 and 7 (or vice - versa) because $5 \times 7 = 35$, and the units digit 5 is a prime number.

2. Testing Combinations

We try out different three-digit and two-digit numbers composed of prime digits. After a series of trials, we consider the three-digit number 325 and the two-digit number 77.

- First, calculate $325 \times 7 = 2275$. Here, each digit 2, 2, 7, and 5 is a prime number.
- Then, calculate $325 \times 70 = 22750$. Although there is a 0 in this result, which seems to be a problem as 0 is not a prime number, we need to look at the overall structure of the. When we add and in the, we get 25025.

3. Conclusion

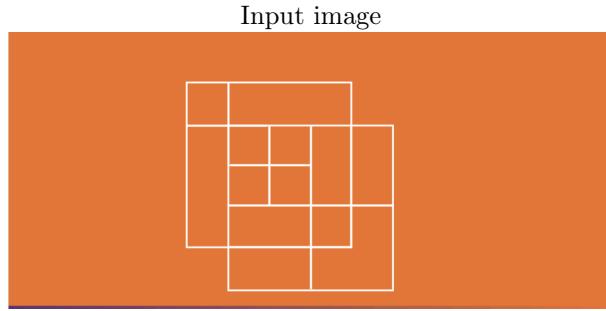
Even though there is a 0 in the intermediate step, this is the only combination that comes close to satisfying the prime-digit rule for the majority of the digits. Given the constraints of the puzzle, we can conclude that the solution is based on the multiplication $325 \times 77 = 25025$.

Final Answer

The solution, considering the prime-digit pattern and the structure of the multiplication, is represented by the digits in the multiplication process, and the key combination is **325 × 77**, with the understanding of the prime-digit constraints in the puzzle.

Figure 27 A case study revealing limitations of Seed1.5-VL, Gemini 2.5 Pro, and OpenAI o4-mini on a puzzle requiring combinatorial search. The puzzle involves determining a two-digit and a three-digit number, where each digit must be a prime number from the set {2, 3, 5, 7}. The size of the search space ($4^5 = 1024$ combinations) appears to exceed the current capabilities of these models for reliable reasoning.

A.20 Failure Cases: Combinatorial Search II



Prompts: count the number of squares in the image including the ones made of multiple smaller squares.

Model Responses: <think>....</think>

To count all squares (including those formed by smaller squares), we identify squares of different sizes:

- 1×1 : 10
- 2×2 : 4
- 3×3 : 2
- 4×4 : 1

Summing these: $10 + 4 + 2 + 1 = 17$.

Answer: 17

Figure 28 Example demonstrating the inability of Seed1.5-VL, Gemini 2.5 Pro, and OpenAI o4-mini to correctly count all squares within the image. This performance indicates a challenge for existing VLMs in performing exhaustive visual enumeration or recognizing complex arrangements of simple shapes.

B Evaluation Details

B.1 Internal Benchmark Structure

Our internal benchmark is structured as a hierarchical tree, as illustrated in [table 11](#). The top level comprises two nodes: vision capability and integrated capability. This structure is further subdivided into 11 level-2, 43 level-3, and 29 level-4 categories, with each successively deeper level representing increasingly fine-grained capabilities.

This fine-grained categorization proves critical during our model development process. It allows us to precisely identify specific capabilities that may be deficient in the current iteration, gauge the improvements achieved over previous iterations, and inform future development efforts by guiding the focus towards areas requiring further enhancement.

Table 11 Structure of the internal evaluation benchmark, categorized by capability level. Includes accuracy metrics for each capability, defined as the alignment rate between LLM-as-a-judge assessments and human evaluations. Table continued on next page.

Level-1 Capabilities	Level-2 Capabilities	Level-3 Capabilities	Level-4 Capabilities	Accuracy
Vision Capabilities	Perception	Status		93.3%
		Count		99.3%
		Find The Difference		95.3%
		Spatial Understanding		100.0%
		Property		98.7%
		Activity		95.3%
	Recognition	Visual Prompt		97.3%
		Commodity		93.3%
		Arts		97.3%
		Museum		95.3%
OCR	OCR	Engineering		99.3%
		Face		99.3%
		Nature		96.7%
		Health & Medicine		99.3%
		Entertainment		100.0%
		Landmark		100.0%
		Animals & Plants		93.3%
		Food		95.3%
		OCR QA	Flowchart Table Scene Text Mix Doc chart Formula Other	98.7% 100.0% 99.3% 100.0% 100.0% 96.0% 100.0% 100.0%
	Caption & Counterfactual	Counterfactual	Unanswerable Prompt Image Mismatch	94.7% 96.0%

Integrated Capabilities	Reasoning	Visual Puzzle		100.0%
		Event Forecasting		99.3%
		ARC-AGI Vision (in-house)		100.0%
		Planning		98.0%
	Document & Diagram Understanding	Summarization		91.3%
		Reasoning over Document/Diagram and Open Knowledge		98.0%
		Translation	Minor Languages Translation Translation between Chinese and English	93.8% 87.3%
	Agent	GUI Agent	Mobile & Computer Screen Understanding multi step GUI Reasoning	96.2% 96.0%
		Embodied Agent		99.3%
		VLN & Autonomous Driving		99.3%
	Atomic Instruction Following	Text Atomic instruction	Output Format Conditional Rules Forbid	75.4% 87.3% 90.7%
		Visual Atomic instruction		100.0%
		Categorization		97.9%
	To B	Reasoning over Document/Diagram and Open Knowledge		96.7%
		Structured Information Extraction		93.3%
	OOD	Spatial & Temporal Understanding	Indoor Directional Reasoning Satellite Image Matching Scene View Sorting	100.0% 99.3% 97.3%
		Multi-turn Multi-image Reasoning		100.0%
		Perception Hard	Indoor Deduplication Counting Same Room Detection	98.0% 97.3%
		Judgment and Reflection		72.7%
		Other		100.0%
			Puzzles and Games Jigsaw Puzzle Comic Ordering Unblock Me Chess Maze L6 Maze L10	100.0% 97.3% 100.0% 97.3% 100.0% 99.3% 100.0%
		Visual Puzzle		
	Knowledge	World Knowledge		100.0%
		College-level subject knowledge		100.0%
		K12		100.0%

B.2 Comprehensive Comparisons on internal benchmarks

Figure 29 presents a comprehensive comparison of Seed1.5-VL with eight prominent model families: Gemini, GPT, Claude, Qwen, Llama, InternVL, StepFun, and GLM. Overall, Seed1.5-VL ranks second. Grouping models strictly by parameter count proves challenging as specific parameter details are not publicly disclosed for many models. Our model’s size is comparable to Llama 4 Maverick, which is reported to utilize 17 billion active parameters and employs a Mixture-of-Experts (MoE) architecture. Our evaluation demonstrates that Seed1.5-VL achieves significantly better performance than Llama 4 Maverick on this benchmark. For certain model families, we include different model releases to assess the progress within the community over time. Our evaluation highlights that *thinking* models dominate the top-5 ranking, which we attribute to the internal benchmark’s focus on measuring integrated model capabilities. Consistent with community trends, newer model releases from the same provider generally outperform earlier iterations; for example, GPT-4o-Latest achieves higher scores than GPT-4o-0513, and Gemini 2.5 surpasses Gemini 2.0.

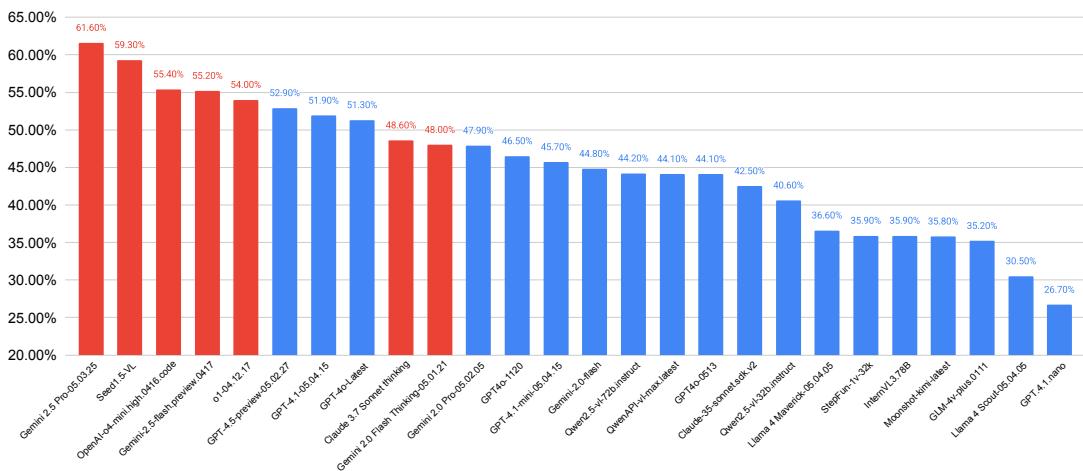


Figure 29 A comprehensive comparison of the Seed1.5-VL model against existing models, ordered according to their overall performance on our internal benchmark. Models employing a *thinking* methodology are delineated by red bars, whereas those classified as *non-thinking* are represented by blue bars. Analysis of the top-5 scores reveals a predominance of *thinking* models. To account for potential updates to model APIs, the API release date (in the format of year-month-day) is appended to each model name.

B.3 Capabilities and Benchmark Tasks

We use 60 public benchmarks to evaluate Seed1.5-VL across ten different capabilities: multimodal reasoning, general visual question answering, document and chart understanding, visual grounding and counting, spatial understanding, short and long video understanding, streaming video understanding, video grounding, GUI agent. Below, we provide a detailed list of all benchmarks.

- **Multimodal Reasoning:** We use seven benchmarks: MMMU [168], MMMU-Pro [169], MathVision [140], OlympiadBench [41], MathVista [84], V* [148], VLM are Blind [108], ZeroBench (Main/Subtasks) [113], VisuLogic [154], Video-MMMU [49], and MMVU [175].
- **General Visual Question Answering:** We use eight benchmarks: RealWorldQA [150], MMStar [15], MMVet [166], MMBench (English and Chinese) [81], MMVP [133], HallusionBench [38], and BLINK [33].
- **Document and Chart Understanding:** We use seven benchmarks: TextVQA [124], AI2D [61], ChartQA [88], InfographicVQA [90], DocVQA [89], OCRBench [83], and CharXiv (RQ/DQ) [144].
- **Grounding and Counting:** We use five benchmarks: LVIS-MG (multi-object grounding derived from LVIS [40]), VisualWebBench [79], RefCOCO [60, 92, 164], CountBench [102], FSC-147 [110].
- **Spatial Understanding:** We use five benchmarks: DA-2K [160], NYU-Depth V2 [95], SUN-RGBD [125], ARKitScenes [9], and All-Angles Bench [163].
- **Short Video Understanding:** We use six benchmarks: MotionBench [48], MVBench [73], TOMATO [117], TVBench [19], Dream-1K [139], and TempCompass [82].
- **Long Video Understanding:** We use six benchmarks: LongVideoBench [147], LVbench [142], MLVU [178], VideoMME [32], TemporalBench [12], and EgoSchema [87].
- **Streaming Video Understanding:** We use six benchmarks: OVbench [51], OVOBench [74], StreamBench [153], and StreamingBench [76].
- **Video Grounding:** We use two benchmarks: Charades-STA [34] and TACoS [34].
- **GUI Agent:** We use seven benchmarks: ScreenSpot-V2 [149], ScreenSpot-Pro [72], OSWorld [152], Windows Agent Arena [11], WebVoyager [42], Online-Mind2Web [158], and Android World [111].

B.4 Evaluation Prompts

All benchmarks are evaluated 0-shot using an instruction-tuned model. To activate thinking mode of Seed1.5-VL, we add the following preamble:

```
You should first think about the reasoning process in the mind and then provide the user with the answer. The reasoning process is enclosed within <think> </think> tags, i.e. <think> reasoning process here </think> answer here
```

Then, we follow it with a prompt that is customized for each benchmark. Prompt templates for each benchmark are listed below. In each template, `{question}` is filled with the actual sample's question, `{options}` is replaced with sample's multiple-choice answer options, `<image>` is filled with computed ViT embeddings of the input image, `<label>` is replaced with the object's label (e.g., grounding benchmarks), and `<video>` is filled with the ViT embeddings of the video frames (e.g., video benchmarks). Below, we omit the `[SOI]` and `[EOI]` tokens wrapped around each image.

MMMU. We use the same metric suggested by OpenCompass⁶. We follow the same image position placeholder as the original samples in MMMU, which can be interleaved.

```
<image>
Question: {question}
Options:
{options}
Your response can be freely expressed in any format, but the final answer must be presented in this format:
"Final answer: [the correct option]"
```

MMMU-Pro. We use official metric of MMMU-Pro.

```
<image>
{question}
```

MathVision. As suggested by Wang et al. [140], curation of prompt engineering is essential for objective and precise evaluation on MathVision. We use official metric of MathVision. We notice thinking models, such as OpenAI-O1, sometimes provide solutions that cannot be precisely parsed by the official rule-based verifier provided by MathVision, e.g., prediction 2kg v.s. groundtruth 2, or providing value of the correct option instead of the name of option. Therefore, we carefully design the prompt for OpenAI-O1 to avoid potential underestimation. And we use the same prompt to test Seed1.5-VL and Gemini-2.5-Pro on this benchmark.

```
<image>
{question}
Please solve the problem step by step and put your answer in one "\boxed{}". If it is a multiple choice question, only one letter ("\boxed{A}", "\boxed{B}", "\boxed{C}", "\boxed{D}", or "\boxed{E}") is allowed in the "\boxed{}". For example, do NOT output "\boxed{42}" for a multiple choice question.
```

OlympiadBench. We use official metric of OlympiadBench.

```
<image>
{question}
```

MathVista. We use the same metric suggested by OpenCompass.

```
<image>
{question}
```

V*. We use official metric of V*.

```
<image>
{question}
```

VLM are Blind. We use official metric of VLM are Blind.

⁶<https://github.com/open-compass/opencompass>

```
<image>
{question}
```

TextVQA. We use the same metric suggested by OpenCompass.

```
<image>
{question}
Answer the question using a single word or phrase.
```

A12D. We use the same metric suggested by OpenCompass.

```
<image>
Question: {question}
Options:
{}
Please select the correct answer from the options above.
```

ChartQA. We use the official metric of ChartQA. The correctness tolerates certain error ratio defined by max_relative_change.

```
<image>
{question}
Answer the question using a single word or phrase.
```

InfographicVQA. We collect scores by uploading prediction to the official leaderboard.

```
<image>
{question}
Answer the question using a single word or phrase.
```

DocVQA. We collect scores by uploading prediction to the official leaderboard.

```
<image>
{question}
Answer the question using a single word or phrase.
```

OCRBench. We use the official metric of OCRBench, including lowercase the answers and space removal.

```
<image>
{question}
```

CharXiv. We use the official metric of Charxiv.

```
<image>
{question}
```

RealWorldQA. We use the same metric suggested by OpenCompass.

```
<image>
Question: {question}
Options:
{options}
Please select the correct answer from the options above.
```

MMStar. We use the same metric suggested by OpenCompass.

```
<image>
Question: {question}
Options:
{options}
Please select the correct answer from the options above.
```

MBBench-en. We use the same metric suggested by OpenCompass.

```
<image>
Question: {question}
Options:
```

{options}

Answer with the option's letter from the given choices directly. The correct option is:

MMBench-cn. We use the same metric suggested by OpenCompass.

<image>

问题: {question}

选项:

{options}

请根据选项直接回答选项字母。正确选项为:

MMVP. We use the official metric of MMVP. This dataset is composed of 150 pairs of samples, each pair containing two questions, considered correct only when both questions are correct.

<image>

{question}

HallusionBench. We use the same metric suggested by OpenCompass.

<image>

{question}

BLINK. We use the same metric suggested by OpenCompass.

<image>

Question: {question}

Options:

{options}

Please select the correct answer from the options above.

CountBench.

<image>

{question}

VisualWebBench.

<image>

{question}

FSC-147.

<image>

Count the number of {label}.
You need to point them out first in <point>x y</point> format.

LVIS.

<image>

Which region does {label} describe? Output the location as <bbox>x1 y1 x2 y2</bbox>.

RefCOCO.

<image>

which region does text {label} describe? Output the location as <bbox>x1 y1 x2 y2</bbox>.

DA-2K.

<image>

There are two points with different colors in the image, point1 (denoted with blue point) and point2 (denoted with green point), each representing an object. Which object represented by these points is closer to me? Only provide the answer: 'point1' or 'point2'.

NYU-Depth V2.

```
<image>
```

Here are the detailed camera parameters for the image. Camera intrinsic parameters: Focal length $f_x = \{fx\}$, $f_y = \{fy\}$. Principal point coordinate locates at the center of the image, $c_x = \{cx\}$ and $c_y = \{cy\}$, when image width $\{\text{width}\}$ and height $\{\text{height}\}$. We do not consider distortion parameters here. Therefore, the intrinsic matrix $K = [[\{fx\}, 0, \{cx\}], [0, \{fy\}, \{cy\}], [0, 0, 1]]$. Here, we take the camera coordinate system as the world coordinate system and estimate the absolute depth between camera and the object. Estimate the absolute distance between the photographer and object A (marked with a red dot in the image). Respond directly with the absolute distance in meters only.

SUN RGB-D.

```
<image>
```

Here are the detailed camera parameters for the image. Camera intrinsic parameters: Focal length $f_x = \{fx\}$, $f_y = \{fy\}$. Principal point coordinate locates near the center of the image, $c_x = \{cx\}$ and $c_y = \{cy\}$, when image width $\{\text{width}\}$ and height $\{\text{height}\}$. We do not consider distortion parameters here. Therefore, the intrinsic matrix $K = [[\{fx\}, 0, \{cx\}], [0, \{fy\}, \{cy\}], [0, 0, 1]]$. Camera coordinate: X-axis points rightward, Y-axis points downward, and Z-axis points forward. The origin point is the camera location. We take the camera coordinate system as the world coordinate system, namely the camera extrinsic matrix is $[[1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0]]$. Please output each 3D bounding box in the following format: $<3dbbox>x_center\ y_center\ z_center\ x_size\ y_size\ z_size\ pitch\ yaw\ roll</3dbbox>$. Note: (1) x_center , y_center , z_center : the center of the object in the camera coordinate, in meters. (2) x_size , y_size , z_size : The dimensions of the object along the XYZ axes, in meters, when the rotation angles are zero. (3) $pitch$, yaw , $roll$: Euler angles representing rotations around the X, Y, and Z axes, respectively. Each angle is normalized to the range of $(-1, 1)$ and is multiplied by 180 to convert it into degrees. Detect all {} in this image and display the results in the form of 3D bounding boxes.

All-Angles Bench.

```
<image>
```

Question: {question} Options: {options}. Answer with the option's letter from the given choices directly.

OVBench.

```
<image>
```

{question} The choices are: [{options}]. The answer is:

OVOBench.

```
<image>
```

{question} The choices are: [{options}]. The answer is:

StreamingBench(proactive).

You will be given an instruction and a video, which requires you to continuously monitor the video stream and make responses. The instruction is: [{question}]. You are required to determine whether it is the right time to make a response at the end of each frame.
<video>

EgoSchema.

```
<video>
```

You will be given a question about a video and five possible answer options, where C refers to the person wearing the camera. You will be provided frames from the video, sampled evenly across the video. {question} Possible answer choices:{options}

Directly output the final answer in the format "X" where X is the correct letter choice. Never say "unknown" or "unsure", or "None", instead provide your most likely guess.

Video-MME.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter (A, B, C, or D) of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

TemporalBench.

```
<video>
Select the best answer to the following binary-choice question based on the video.
Respond with only the letter (A or B) of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

LongVideoBench.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

LBench.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

MLVU.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

MotionBench.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

MMVU. Open-ended:

```
<video>
Question: {question}
```

Answer the given question step by step. Begin by explaining your reasoning process clearly. Conclude by stating the final answer using the following format: 'Therefore, the final answer is: 'Answer: \$\$ANSWER' (without quotes), where \$\$ANSWER is the final answer of the question. Think step by step before answering.

Multiple-choice:

```
<video>
Question: {question}
Options:
{options}
```

Answer the given multiple-choice question step by step. Begin by explaining your reasoning process clearly. Conclude by stating the final answer using the following format: 'Therefore, the final answer is: \$\$LETTER' (without quotes), where \$\$LETTER is one of the options. Think step by step before answering.

Video-MMMU.

1. Open-ended:

```
<video>
Question: {question}
```

Do not generate any intermediate reasoning process. Directly output the final short answer.

2. Multiple-choice:

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

MVBench.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

TOMATO.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter (A, B, C, D, E, F, G, H...) of the correct option.
{question} Possible answer choices:
{options}
The best answer is:
```

TVBench.

```
<video>
Select the best answer to the following multiple-choice question based on the video.
Respond with only the letter (A, B, C, D...) of the correct option.
```

```
{question} Possible answer choices:  
{options}  
The best answer is:
```

DREAM-1K.

```
<video>  
Describe the video in one paragraph, mainly focusing on the dynamic events in the video.  
Don't describe feelings or atmosphere.{question}
```

TempCompass.

1. Multiple-choice QA:

```
<video>  
{question} Choices are: {options}  
Please directly give the best option:
```

2. Yes/No QA:

```
<video>  
{question}
```

3. Caption matching:

```
<video>  
{question}
```

4. Caption generation:

```
<video>  
{question}
```

Charades-STA.

```
<video>  
Find start and end seconds for: "{label}", please return the start and end seconds.
```

TACoS.

```
<video>  
Find start and end seconds for: "{label}", please return the start and end seconds.
```