

# Bounded-Compute Multimodal Regression for Product-Rating Prediction

William Leach   Ru He   Sizhuo Ma   Yizhen Jia   Min Cao   Jian Wang   Rick Cao  
Snap Inc.

wleach@snapchat.com

## Abstract

*Vision-language models (VLMs) are increasingly attractive for multimodal quality assessment, but their default reliance on autoregressive text generation and dynamic visual processing is poorly matched to scalar regression under strict latency budgets. We present a bounded-compute adaptation of SmolVLM2-256M-Video-Instruct [24] for product-rating prediction in the LoViF 2026 Efficient VLM challenge. Motivated by recent multimodal engagement-prediction results showing that feature-based regression can outperform token-based score generation [26], we replace the language-modeling head with a lightweight two-layer MLP fed by pooled decoder states, and we enforce deterministic inputs through fixed  $384 \times 384$  images and truncated metadata. Across controlled ablations, static global image processing slightly outperforms dynamic tiling, and scaling from 100K to 16M training examples substantially improves validation correlation. Under the official held-out evaluation, our 228M-parameter model achieves 0.39 PLCC and 0.40 CES, providing a strong and reproducible baseline for resource-constrained multimodal regression.*

## 1. Introduction

Vision-language models (VLMs) have become a standard foundation for multimodal reasoning, with strong performance on visual question answering, captioning, document understanding, and image-text dialogue [1, 16, 20, 21]. However, most high-performing VLMs are optimized for open-ended text generation rather than bounded-cost scalar prediction. This mismatch is especially problematic in deployment settings that require predictable latency, fixed memory usage, and high throughput.

We address this bounded-compute regression challenge through the lens of *automated product-rating prediction* using product images and structured metadata. Although this task only requires a single scalar output, a standard generative VLM pipeline still incurs autoregressive decoding overhead and often relies on dynamic image tiling or resizing to preserve fine-grained visual detail [4, 17]. These

design choices are sensible for open-ended recognition and reasoning, but they produce variable sequence lengths and input-dependent compute, which is undesirable for large-scale or real-time regression systems.

This problem is formalized by the LoViF 2026 Challenge on Efficient VLM for Multimodal Creative Quality Scoring, which evaluates product-rating prediction under joint accuracy and efficiency constraints [7]. Recent work on short-video engagement prediction [18, 19] provides an important clue for this setting: Sun et al. [26] show that feature-based regression on top of multimodal hidden states outperforms token-based score generation for Qwen2.5-VL, suggesting that hidden-state regression is often a better fit than autoregressive decoding when the target is a continuous score.

Motivated by this observation, we adapt SmolVLM2-256M-Video-Instruct into a deterministic multimodal regressor [24]. We replace the language-modeling head with a lightweight two-layer MLP, pool the decoder hidden states into a fixed representation, and enforce bounded compute by resizing all images to a fixed resolution and aggressively truncating metadata. Our backbone choice is deliberate: SmolVLM2 is a compact VLM family built for efficient multimodal inference, and the 256M Video-Instruct checkpoint combines a SigLIP vision encoder with a SmolLM2 decoder in an Idefics3-style design [14, 16, 24].

Our contributions are threefold. First, we present a simple but effective recipe for converting a compact generative VLM into a bounded-compute multimodal regressor. Second, through controlled ablations, we show that static global image processing can slightly outperform dynamic tiling for this global-context scoring task, while short metadata truncation provides a favorable efficiency–accuracy trade-off. Third, we demonstrate that a 256M-class VLM continues to benefit from large-scale supervision, scaling effectively from 100K to 16M Amazon Reviews’23 examples while maintaining strong performance under the official LoViF evaluation protocol.

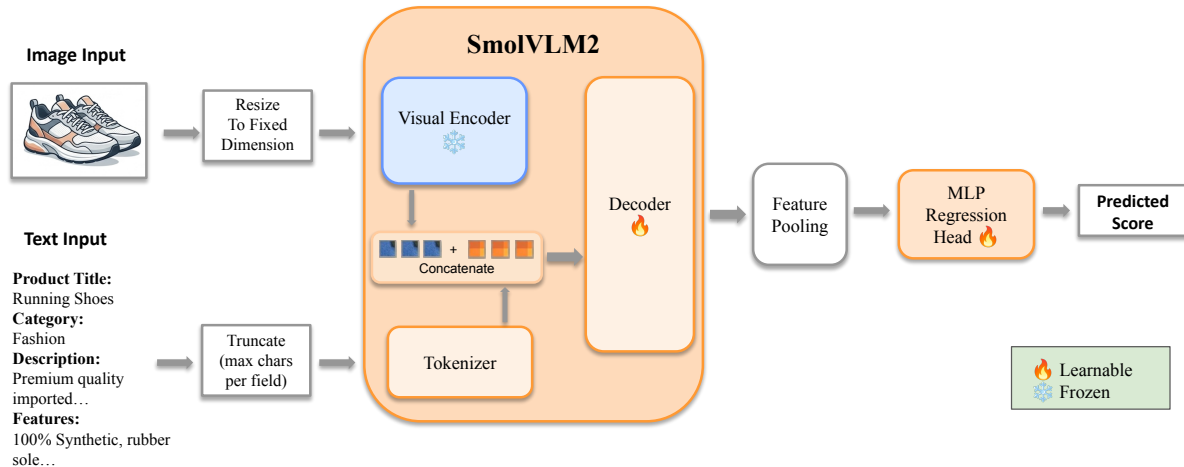


Figure 1. Overview of our bounded-compute regression pipeline. A fixed-resolution image and truncated text metadata are encoded into a multimodal token sequence, processed by SmolVLM2, pooled into a single representation, and mapped to a scalar rating by a lightweight MLP head.

## 2. Related Work

### 2.1. Vision-Language Models

Representative modern vision-language model (VLM) families include Flamingo [1], BLIP-2 [20], LLaVA-1.5 [21], and Idefics3-style architectures [16]. The Qwen family has significantly influenced the field: Qwen-VL established multilingual grounding and OCR [3], while Qwen2-VL and Qwen2.5-VL introduced dynamic resolution and enhanced document/video parsing [2, 4]. More recently, the evolution of the backbone family from Qwen2 to Qwen3 has further pushed the boundaries of multimodal performance [2, 28].

However, these models primarily focus on open-ended text generation, which introduces variable latency and input-dependent compute [4, 17]. A parallel line of research focuses on compact and efficient VLMs for constrained deployment, such as MobileVLM [5], MobileVLM V2 [6], MiniCPM-V [29], PaliGemma 2 [25], and DeepSeek-VL2 [27]. Marafioti et al. [24] introduced SmolVLM, aimed at resource-efficient inference, by systematically studying the factors that affect compact VLMs, including the balance between vision encoder and language-model decoder capacity, tokenization, positional encoding, text prompt engineering, and training-data composition. They subsequently released SmolVLM2, a family of VLMs (ranging from 256M to 2.2B parameters) with even better visual understanding capabilities [14].

While our work aligns with this compact-VLM trajectory, we differentiate ourselves by focusing on **bounded-compute scalar regression** rather than generative tasks.

### 2.2. Multimodal Regression and Quality Scoring

Recent research has increasingly focused on adapting large multimodal models for scalar prediction and perceptual assessment. For instance, DepictQA explores multimodal language models for image quality assessment [30], while Aes-Bench provides a benchmark for evaluating image aesthetics perception in large models [13]. In the domain of short-video engagement, Li et al. [18] introduce a large-scale benchmark dataset comprising real-world user-generated content (UGC) short videos, and proposed a predictive model that uses a cross-modal attention mechanism to incorporate rich multi-modal features.

Critically, Sun et al. [26] evaluated both feature-based regression and token-based score generation for engagement prediction on UGC using large multimodal models. Their findings indicate that regression from hidden features via a lightweight MLP consistently outperforms autoregressive token generation when the target is a continuous score [26]. This empirical evidence directly motivates our decision to repurpose a compact VLM into a deterministic regressor, replacing the language-modeling head with a regression-specific architecture to ensure both higher accuracy and stable computational overhead [24, 26].

### 2.3. Multimodal E-commerce Prediction and Recommendation

Our task is closely related to multimodal recommendation and e-commerce modeling, where product images and metadata provide complementary signals regarding item quality and consumer preference. Early approaches like VBPR [11] pioneered the integration of visual item features into recommendation frameworks. More recently, multimodal review analysis has explored predicting review help-

fulness by combining product, review, image, and text signals through multi-perspective coherent reasoning [22] or selective-attention contrastive learning [10].

As noted in recent surveys, multimodal item representations are increasingly critical for predictive tasks across online platforms [23]. Our setting differs from these lines in its focus: we predict a product-level scalar rating directly from a primary image and structured metadata under strict parameter and FLOP budgets.

### 3. Method Details

Our goal is to adapt a compact generative VLM into a deterministic regressor for product-rating prediction. The input to our model is a single product image paired with structured metadata fields (title, description, features, and main category), and the target is the product’s average star rating in [1, 5].

#### 3.1. Backbone Choice

We build upon SmolVLM2-256M-Video-Instruct, a compact instruction-tuned VLM that combines a SigLIP vision encoder with a SmolLM2 decoder in an Idefics3-style multimodal architecture [14, 16, 24]. This model is particularly well-suited for the LoViF efficiency challenge because it retains strong multimodal grounding capabilities while operating at a much smaller scale than mainstream VLMs. Our adapted regression model contains 228M parameters in total, with 135M trainable during fine-tuning.

#### 3.2. Deterministic Multimodal Processing

A key design goal of our model is *bounded compute*. Standard VLM pipelines often use dynamic image tiling or native dynamic resolution to preserve fine detail [4, 17]. While beneficial for OCR-heavy or fine-grained recognition tasks, those mechanisms introduce input-dependent visual token counts and therefore variable latency and memory usage.

To avoid this variability, we disable dynamic resizing and image splitting, and resize every image to a fixed  $384 \times 384$  resolution using bilinear interpolation. On the text side, we concatenate the product’s metadata fields into a structured natural language prompt. To guarantee a strictly bounded language context, each formatted key–value pair is individually truncated to a maximum character budget  $L$  (e.g.,  $L = 100$ ) before concatenation. The final multimodal prompt  $P$  follows this fixed template:

```
<image> The average user rating
for this product. Text metadata:
Title: <title>[:L], Description:
<description>[:L], Features:
<features>[:L], Main Category:
<category>[:L]
```

The entire prompt  $P$  is tokenized and processed by the SmolVLM2 decoder. Unlike generative architectures that output numerical scores via autoregressive next-token generation (e.g., Qwen2.5-VL), our model aggregates the hidden states of this deterministic sequence and projects them directly into a continuous scalar rating. This yields strictly bounded multimodal sequence lengths, stable FLOPs, and highly reproducible runtime characteristics across all samples.

#### 3.3. Regression Adaptation

While SmolVLM2 processes the visual and textual tokens with an autoregressive decoder, our task is constrained to a single scalar output. Following recent evidence that feature-based regression can outperform token-based score generation for multimodal engagement prediction [26], we remove the language-modeling head and attach a lightweight two-layer MLP regression head.

Given an RGB image  $x$  and tokenized metadata, the SigLIP vision encoder first produces patch-level visual embeddings. These embeddings are compressed by the model’s pixel-shuffle connector and projected into the decoder embedding space ( $d = 576$ ), then concatenated with text tokens and processed by the SmolLM2 decoder. Let  $T$  denote the maximum sequence length in a batch, and  $h_i \in \mathbb{R}^d$  the final hidden state at position  $i$ . To aggregate the variable-length decoder sequence into a fixed-dimensional representation, we employ mask-aware mean pooling:

$$h_{\text{pool}} = \frac{\sum_{i=1}^T m_i h_i}{\sum_{i=1}^T m_i}. \quad (1)$$

Here,  $m_i \in \{0, 1\}$  represents the standard binary attention mask generated natively by the multimodal processor during batching,  $m_i = 1$  if the token at position  $i$  is a valid visual or textual token, and  $m_i = 0$  if the token at position  $i$  is a padding token.

This formulation ensures that only meaningful multimodal features contribute to the final pooled state, explicitly preventing padding tokens from diluting the representation.

The pooled vector is then passed through a two-layer MLP head,  $\text{Linear}(576, 288) \rightarrow \text{ReLU} \rightarrow \text{Linear}(288, 1)$ , which outputs a scalar logit  $x$ . To restrict predictions to the valid rating range, we apply a scaled sigmoid:

$$\hat{y} = 1 + 4 \sigma(x), \quad (2)$$

where  $\sigma(\cdot)$  is the logistic sigmoid function. This bounded parameterization explicitly constrains the output to the valid [1, 5] interval, ensuring that the model avoids implausible out-of-range predictions by construction.

### 3.4. Training Objective

We optimize the model with mean squared error between predicted and ground-truth ratings:

$$L = \frac{1}{N} \sum_{i=1}^N (y_{\text{pred}}^{(i)} - y_{\text{true}}^{(i)})^2. \quad (3)$$

During fine-tuning, the SigLIP vision encoder and pixel-shuffle connector are frozen, while the decoder and regression head are updated. This choice preserves the compact model’s pretrained visual representation while focusing optimization on the multimodal decoder and the task-specific scalar head.

## 4. Experiments

### 4.1. Dataset and Evaluation Metrics

We train our model on Amazon Reviews’23 [12, 15], using item-level metadata paired with a representative product image. Each sample contains one primary product image together with four metadata fields: title, description, features, and main category. The regression target is the product’s average star rating.

**Category balancing.** Amazon Reviews’23 is highly imbalanced across product categories: large groups such as `Unknown`, `Clothing_Shoes_and_Jewelry`, and `Books` contain far more items than long-tail categories such as `Subscription_Boxes`, `Gift_Cards`, and `Magazine_Subscriptions` [15]. To reduce category dominance, we apply popularity-stratified sampling within each category: we select the top-1M and bottom-1M items ranked by `rating_number`, and retain all items for categories with fewer than 2M candidates.

**Filtering and preprocessing.** We further filter the sampled items to those with (i) at least 10 reviews and (ii) an available `MAIN` image, preferring high-resolution URLs when available. This yields 16,455,671 training samples, with 10,000 held out for validation. All images are resized to  $384 \times 384$  using bilinear interpolation; dynamic resizing and image splitting are disabled. Metadata fields are formatted as key–value text and truncated to 100 characters per field unless otherwise noted.

To evaluate predictive accuracy and ranking quality, we report:

- **RMSE (Root Mean Square Error):** average magnitude of the prediction error; lower is better.
- **PLCC (Pearson Linear Correlation Coefficient):** linear correlation between predicted and target ratings; higher is better.

- **SRCC (Spearman’s Rank Correlation Coefficient):** rank correlation between predictions and targets; higher is better.

### 4.2. Training and Implementation Details

Unless otherwise noted, all experiments use SmolVLM2-256M-Video-Instruct with the deterministic  $384 \times 384$  preprocessing pipeline described in Sec. 3. The SigLIP vision encoder and pixel-shuffle connector are frozen, while the decoder and regression head (with 135M parameters) are fine-tuned with MSE loss. Training is performed with `DistributedDataParallel` on 4 NVIDIA A100 GPUs using a global batch size of 64. We use 8-bit AdamW optimizer states [9], a peak learning rate of  $4 \times 10^{-4}$ , linear decay, 3% warmup, and 0.1 dropout in the regression head. We train for up to 5 epochs, employing early stopping to select the optimal model checkpoint based on peak validation performance.

At inference time, the adapted model contains 228M active parameters. On a single NVIDIA A100 with batch size 64, the final model runs at 0.0084 seconds per image (119.3 images/s). We also use FlashAttention-2 [8] and `bfloat16` automatic mixed precision.

### 4.3. Ablation Studies

We perform controlled ablations to characterize the accuracy–efficiency trade-offs of compact multimodal regression. Unless otherwise stated, these ablations are conducted on a 1M-sample subset of the training data.

**Model capacity.** We first compare the 256M and 500M SmolVLM2 variants under the original dynamic  $512 \times 512$  preprocessing pipeline. Table 1 illustrates the predictive performance of both backbones.

Table 1. Effect of model capacity (1M samples).

Backbone	RMSE ↓	PLCC ↑	SRCC ↑
SmolVLM2-256M	0.333	0.683	0.652
SmolVLM2-500M	<b>0.329</b>	<b>0.694</b>	<b>0.664</b>

While the 500M model yields better predictive correlation, the gain comes at a substantial resource cost. Specifically, upgrading to the 500M variant roughly doubles the active parameter count (460M vs. 228M) and raises the theoretical maximum compute from 113 GFLOPs to 165 GFLOPs. Given the strict efficiency constraints of our target application, the +0.011 gain in PLCC does not mathematically justify the increased memory footprint and computational overhead. We therefore retain the 256M model for all subsequent experiments to prioritize inference efficiency.

**Static vs. dynamic visual preprocessing.** The LoViF benchmark’s canonical inference protocol strictly requires deterministic compute. Modern dynamic tiling generates a variable number of visual tokens depending on the input image’s aspect ratio, resulting in heavily fluctuating, image-dependent FLOPs that cannot be objectively evaluated under the challenge’s fixed-cost scoring system. To meet the benchmark requirements, we ablated the default dynamic tiling strategy against a forced static global resize at the same nominal  $512 \times 512$  resolution. Because dynamic tiling compute varies per sample, we do not report a static FLOP comparison for this ablation.

Table 2. Dynamic vs. static image preprocessing at  $512 \times 512$ .

Preprocessing	RMSE ↓	PLCC ↑	SRCC ↑
Dynamic tiling	0.333	0.683	0.652
Static global resize	<b>0.331</b>	<b>0.689</b>	<b>0.657</b>

Surprisingly, enforcing this strict operational constraint actually improved predictive performance. Despite using a theoretically less adaptive visual pipeline, static preprocessing proved slightly better in this setting. One possible explanation is that global product-rating prediction depends more on overall object presentation and coarse semantics than on localized high-resolution crops, making dynamic tiling unnecessary or even mildly distracting.

**Input resolution.** Having selected static preprocessing, we ablate image resolution. Table 3 illustrates how scaling the image size impacts predictive performance.

Table 3. Effect of static image resolution (1M samples).

Resolution	RMSE ↓	PLCC ↑	SRCC ↑
$512 \times 512$	<b>0.331</b>	<b>0.689</b>	<b>0.657</b>
$384 \times 384$	0.335	0.679	0.646

While the higher  $512 \times 512$  resolution yields a slight improvement in raw correlation, it also substantially increases the computational cost. Specifically, increasing the resolution from 384 to 512 raises the theoretical maximum compute from 72 GFLOPs to 113 GFLOPs. Given the strict operational constraints of the target benchmark, we find that this  $\sim 57\%$  increase in FLOPs outweighs the  $+0.010$  gain in PLCC. Consequently, we adopt  $384 \times 384$  as our final operating point to maintain a highly efficient compute profile.

**Text truncation.** To bound language-side compute, we vary the metadata truncation limit. Table 4 illustrates the impact of different character limits on predictive performance.

Table 4. Effect of metadata truncation at  $384 \times 384$ .

Char limit	RMSE ↓	PLCC ↑	SRCC ↑
50	0.340	0.666	0.628
100	0.335	0.679	0.646
200	<b>0.333</b>	<b>0.685</b>	<b>0.648</b>

While longer metadata naturally improves correlation, it introduces a substantial computational overhead. Specifically, expanding the truncation limit from 50 to 100 and 200 characters increases the theoretical maximum compute from 65 GFLOPs to 72 GFLOPs and 86 GFLOPs, respectively. Although the 200-character limit yields the highest overall PLCC, the additional  $+0.006$  gain over the 100-character limit requires a disproportionate increase of 14 GFLOPs per forward pass. Balancing this trade-off, we select 100 characters per field as our highly efficient default operating point.

#### 4.4. Data Scaling and Final Results

Using the final efficient configuration (SmolVLM2-256M, static  $384 \times 384$  preprocessing, 100-character truncation), we evaluate how performance changes with training-set scale.

Table 5. Effect of training data scale on validation performance.

Training data	RMSE ↓	PLCC ↑	SRCC ↑
$\sim 100\text{K}$ samples	0.363	0.605	0.558
$\sim 1\text{M}$ samples	0.335	0.679	0.646
$\sim 16\text{M}$ samples	<b>0.326</b>	<b>0.700</b>	<b>0.664</b>

As shown in Table 5, the compact backbone continues to benefit from additional supervision: scaling from 100K to 16M samples improves PLCC by  $+0.095$  and SRCC by  $+0.106$  on the same validation split. This is an important result for the challenge setting, since it shows that a small VLM can convert large external data into measurable quality gains without changing the deployment footprint.

**Official challenge results.** Following the LoViF challenge protocol [7], the canonical evaluation heavily emphasizes operational cost via the Comprehensive Efficiency Score (CES):

$$\text{CES} = \text{PLCC}^+ \times \mathcal{E}(\mathcal{C})$$

Here,  $\text{PLCC}^+ = \max(0, \text{PLCC}(y, \hat{y}))$  clips negatively correlated predictions to zero to avoid misaligned incentives, and  $\mathcal{E}(\mathcal{C})$  is an asymmetric efficiency multiplier based on the geometric resource cost,  $\mathcal{C}$ . The cost  $\mathcal{C}$  balances the

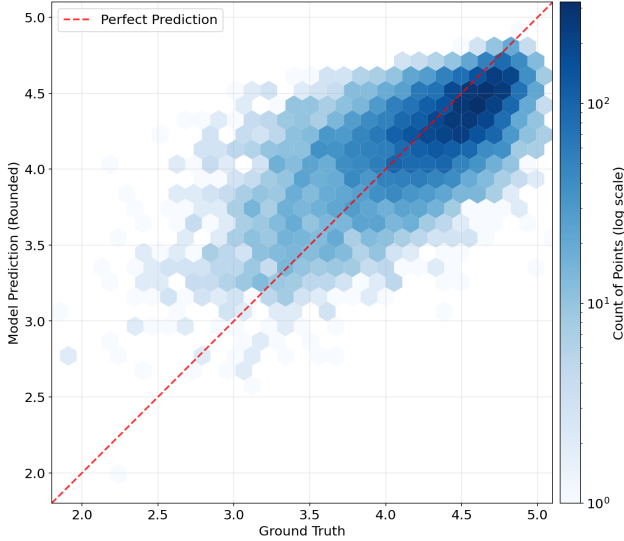


Figure 2. Hexbin density of rounded predictions versus rounded ground-truth ratings on the validation set. Predictions are rounded to one decimal place for visualization to match the label granularity. The dashed red line denotes perfect prediction.

model’s parameters and compute against baseline targets ( $P_{\text{tgt}} = 1000\text{M}$  parameters,  $F_{\text{tgt}} = 20\text{G}$  FLOPs):

$$C = \left( \frac{\text{Params}}{P_{\text{tgt}}} \right)^{0.5} \cdot \left( \frac{\text{FLOPs}}{F_{\text{tgt}}} \right)^{0.5}$$

To encourage lightweight design, the efficiency factor  $\mathcal{E}(C)$  applies a bounded bonus for efficient models ( $C \leq 1$ ) and a steep logarithmic penalty for over-budget models ( $C > 1$ ):

$$\mathcal{E}(C) = \begin{cases} \min(1 + 0.05 \ln(1/C), 1.10) & \text{if } C \leq 1 \\ \frac{1}{1 + 2.0 \ln(C)} & \text{if } C > 1 \end{cases}$$

By calculating the average text length of the metadata fields on the final test set, we estimate an operational compute cost of 68 GFLOPs. Given our 228M parameter footprint, this yields a resource cost of  $C \approx 0.881$ . Because this remains under the reference budget, we achieve a positive efficiency multiplier of  $\mathcal{E}(C) \approx 1.006$ . Under the official held-out test evaluation, our final submission achieves **0.39 PLCC** and **0.40 CES**, **ranked 3rd on the leaderboard**. Because the official benchmark uses a separate hidden split and challenge-specific efficiency scoring, these results should be interpreted separately from the validation ablations above.

#### 4.5. Prediction–Target Alignment

Figure 2 shows a strong concentration of mass near the diagonal, indicating that the model tracks the overall rating

trend well. Errors are concentrated primarily in the mid-range ratings, where many products occupy visually similar quality bands and the target distribution is denser. The bounded sigmoid head also keeps all predictions within the valid  $[1, 5]$  interval by construction.

## 5. Conclusion

We presented a bounded-compute adaptation of SmolVLM2 for multimodal product-rating prediction. By replacing autoregressive score generation with a lightweight regression head, and by enforcing fixed-resolution visual inputs and aggressively truncated metadata, we obtain a deterministic inference pipeline with stable latency and memory usage.

Our experiments show that strong scalar prediction does not require the full dynamic visual processing pipeline commonly used in open-ended VLMs. In particular, static global image processing slightly outperforms dynamic tiling in our setting, while compact text truncation preserves most of the accuracy at a much lower compute cost. We also find that the 256M backbone continues to benefit from large-scale supervision, improving steadily up to 16M training examples and achieving 0.39 PLCC / 0.40 CES on the official held-out evaluation.

These results position compact VLMs as a practical foundation for resource-constrained multimodal regression. At the same time, our design deliberately favors deterministic global processing over adaptivity, which may limit performance on products whose rating depends on fine-grained local details or long-form metadata. Future work could explore conditional compute allocation, stronger multimodal pretraining tailored to scoring tasks, or explicit uncertainty estimation for ambiguous examples.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. *arXiv (Cornell University)*, 2022. 1, 2
- [2] Yang An, Baosong Yang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Zhou, Cheng-Peng Li, Chengyuan Li, Dayiheng Liu, Fei Huang, Guanglong Dong, Haoran Wei, Huan Lin, Jialong Tang, Jialin Wang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Ma, Jianxin Yang, Xu Jin, Jingren Zhou, Jinze Bai, Jinzheng He, Junyang Lin, Kai Dang, Keming Lu, Keqin Chen, Kexin Yang, Mei Li, Mingfeng Xue, Na Ni, Pei Zhang, Peng Wang, Ru Peng, Rui Men, Ruize Gao, Runji Lin, Shijie Wang, Shuai Bai, Sinan Tan, Tianhang Zhu, Tian-

- hao Li, Tianyu Liu, Wenbin Ge, Xiaodong Deng, Xiaohuan Zhou, Xingzhang Ren, Xinyu Zhang, Xipin Wei, Xuancheng Ren, Xuejing Liu, Fan Yang, Yao Yang, Yichang Zhang, Yu Wan, Yunfei Chu, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, Zhifang Guo, and Zhihao Fan. Qwen2 technical report. *arXiv (Cornell University)*, 2024. 2
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. 2
- [4] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *ArXiv*, abs/2502.13923, 2025. 1, 2, 3
- [5] Xiangxiang Chu, Limeng Qiao, Xinyang Lin, Shuang Xu, Yang Yang, Yiming Hu, Fei Wei, Xinyu Zhang, Bo Zhang, Xiaolin Wei, and Chunhua Shen. Mobilevlm : A fast, strong and open vision language assistant for mobile devices. *arXiv (Cornell University)*, 2023. 2
- [6] Xiangxiang Chu, Limeng Qiao, Xinyu Zhang, Shuang Xu, Fei Wei, Yang Yang, Xiaofei Sun, Yi-Ming Hu, Xinyang Lin, Bo Zhang, and Chunhua Shen. Mobilevlm v2: Faster and stronger baseline for vision language model. *arXiv (Cornell University)*, 2024. 2
- [7] CodaBench and LoViF Organizers. Lovif @ cvpr 2026: Challenge on efficient vlm for multimodal creative quality scoring. <https://www.codabench.org/competitions/13463/>, 2026. Accessed 2026-03-19. 1, 5
- [8] Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691*, 2023. 4
- [9] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. *arXiv preprint arXiv:2110.02861*, 2021. 4
- [10] Wei Han, Hui Chen, Zhen Hai, Soujanya Poria, and Lidong Bing. Sancl: Multimodal review helpfulness prediction with selective attention and natural contrastive learning. *arXiv preprint arXiv:2209.05040*, 2022. 3
- [11] Ruining He and Julian McAuley. Vbpr: Visual bayesian personalized ranking from implicit feedback. *Proceedings of the AAAI Conference on Artificial Intelligence*, 2016. 2
- [12] Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiuxi Chen, and Julian McAuley. Bridging language and items for retrieval and recommendation. *arXiv preprint arXiv:2403.03952*, 2024. 4
- [13] Yipo Huang, Quan Yuan, Xiangfei Sheng, Zhichao Yang, Haoning Wu, Pengfei Chen, Yuzhe Yang, Leida Li, and Weisi Lin. Aesbench: An expert benchmark for multimodal large language models on image aesthetics perception. *arXiv (Cornell University)*, 2024. 2
- [14] Hugging Face. Huggingface/tb/smolvlm2-256m-video-instruct model card. <https://huggingface.co/HuggingFaceTB/SmolVLM2-256M-Video-Instruct>, 2025. Accessed 2026-03-19. 1, 2, 3
- [15] Hugging Face Datasets. Mcauley-lab/amazon-reviews-2023 dataset card. <https://huggingface.co/datasets/McAuley-Lab/Amazon-Reviews-2023>, 2024. Accessed 2026-03-19. 4
- [16] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding vision-language models: insights and future directions. *arXiv (Cornell University)*, 2024. 1, 2, 3
- [17] Hugo Laurençon, Léo Tronchon, Matthieu Cord, and Victor Sanh. What matters when building vision-language models? *arXiv (Cornell University)*, 2024. 1, 2, 3
- [18] Dasong Li, Wenjie Li, Baili Lu, Hongsheng Li, Sizhuo Ma, Gurunandan Krishnan, and Jian Wang. Delving deep into engagement prediction of short videos. In *European Conference on Computer Vision*, pages 289–306. Springer, 2024. 1, 2
- [19] Dasong Li, Sizhuo Ma, Hang Hua, Wenjie Li, Jian Wang, Chris Wei Zhou, Fengbin Guan, Xin Li, Zihao Yu, Yiting Lu, et al. Vqula 2025 challenge on engagement prediction for short videos: Methods and results. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pages 3391–3401, 2025. 1
- [20] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 1, 2
- [21] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 26296–26306, 2024. 1, 2
- [22] Junhao Liu, Zhen Hai, Min Yang, and Lidong Bing. Multi-perspective coherent reasoning for helpfulness prediction of multimodal reviews. 2021. 3
- [23] Qidong Liu, Jiayi Hu, Yutian Xiao, Xiangyu Zhao, Jingtong Gao, Wanyu Wang, Qing Li, and Jiliang Tang. Multimodal recommender systems: A survey. *ACM Computing Surveys*, 2024. 3
- [24] Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvlm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025. 1, 2, 3
- [25] Andreas Steiner, André Susano Pinto, Michael Tschanen, Daniel Keysers, Xiao Wang, Yonatan Bitton, Alexey A. Gritsenko, Matthias Minderer, Anthony Sherbondy, Sichang Long, Siyang Qin, Reeve Ingle, Emanuele Bugliarello, Sahar Kazemzadeh, Thomas Mesnard, Ibrahim Alabdulmohsin, Lucas Beyer, and Xiaohua Zhai. Paligemma 2: A family of versatile vlms for transfer. *arXiv (Cornell University)*, 2024. 2
- [26] Wei Sun, Linhan Cao, Yuqin Cao, Weixia Zhang, Wen Wen, Kaiwei Zhang, Zijian Chen, Fangfang Lu, Xiongkuo Min, and Guangtao Zhai. Engagement prediction of short videos with large multimodal models. In *Proceedings of*

*the IEEE/CVF International Conference on Computer Vision Workshops*, pages 3402–3411, 2025. [1](#), [2](#), [3](#)

- [27] Zhiyu Wu, Xiaokang Chen, Zizheng Pan, Xingchao Liu, Wen Liu, Damai Dai, Huazuo Gao, Yiyang Ma, Chengyue Wu, Bingxuan Wang, Zhenda Xie, Yu Wu, Kai Hu, Jiawei Wang, Yaofeng Sun, Yukun Li, Yishi Piao, Kang Guan, Aixin Liu, Xin Xie, Yuxiang You, Kai Dong, Xingkai Yu, Haowei Zhang, Liang Zhao, Yisong Wang, and Chong Ruan. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multimodal understanding. *arXiv (Cornell University)*, 2024. [2](#)
- [28] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, B. X. Yu, Chang Gao, C. Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, H Feng, Hao Ge, Haoran Wei, Lin Huan, Jialong Tang, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiayi Yang, Jing Zhou, Jingren Zhou, Junyang Lin, Kai Dang, Keqin Bao, Kexin Yang, Le Yu, Linyu Deng, Mei Li, Mingfeng Xue, Mingze Li, Pei Zhang, Peng Wang, Qin Zhu, Rui Men, Ruize Gao, Shixuan Liu, Shuang Luo, Tianhao Li, Tianyi Tang, Wenbiao Yin, Xingzhang Ren, Xinyu Wang, Xinyu Zhang, Xuancheng Ren, Fan Yang, Su Yang, Yichang Zhang, Yinger Zhang, Yu Wan, Yuqiong Liu, Zekun Wang, Zeyu Cui, Zhenru Zhang, Zhipeng Zhou, and Zihan Qiu. Qwen3 technical report. *ArXiv.org*, 2025. [2](#)
- [29] Yuan Yao, Tianyou Yu, Ao Zhang, Chongyi Wang, Junbo Cui, Hanlong Zhu, Tianchi Cai, Haoyu Li, Wei Zhao, Zhihui He, Qianyu Chen, Huarong Zhou, Zhixiang Zou, Haoye Zhang, Shengding Hu, Zhi Zheng, Jie Zhou, Jie Cai, Xu Han, Guoyang Zeng, Dahai Li, Zhiyuan Liu, and Maosong Sun. Minicpm-v: A gpt-4v level mllm on your phone. *arXiv (Cornell University)*, 2024. [2](#)
- [30] Zhiyuan You, Zhaodonghui Li, Jinjin Gu, Zhenfei Yin, Tianfan Xue, and Chao Dong. Depicting beyond scores: Advancing image quality assessment through multi-modal language models. *Lecture notes in computer science*, 2024. [2](#)