

CARES: Context-Aware Resolution Selector for VLMs

Moshe Kimhi^{1,2*} Nimrod Shabtay^{2,3*} Raja Giryes³ Chaim Baskin^{4†} Eli Schwartz^{2†}

¹Technion ²IBM Research

³Tel-Aviv University ⁴Ben-Gurion University

Project Page: <https://mkimhi.github.io/CARES/>

Abstract

Large vision–language models (VLMs) commonly process images at native or high resolution to remain effective across tasks. This inflates visual tokens up to to 99% of total tokens of the prefill stage, resulting in high compute and latency, even when low-resolution images would suffice. We introduce *CARES*—a Context-Aware Resolution Selector, a lightweight preprocessing module that, given an image–query pair, predicts the *minimal* sufficient input resolution. CARES uses a compact VLM (350M) to extract features and predict when a target pretrained VLM’s response converges to its peak ability to answer correctly. Though trained as a discrete classifier over a set of optional resolutions, CARES interpolates continuous resolutions at inference for fine-grained control. Across nine multimodal benchmarks spanning documents and natural images, as well as diverse target VLMs, CARES preserves task performance while reducing compute by up to 78% on average across 9 benchmarks.

1 Introduction

Large vision–language models (VLMs) are increasingly used as general-purpose systems that solve a broad variety of visual tasks using a single model. Since the complexity and nature of each task are not known in advance, these models typically process images at very high resolutions to preserve the visual detail necessary for any potential query. This leads to a sharp increase in the number of visual tokens, as modern architectures map higher resolutions to proportionally more tokens. Strategies like AnyRes and tiling further increase token counts in order to capture fine-grained information (Liu et al., 2024a; Wang et al., 2024). In practical settings, visual tokens make up to 99% of all tokens

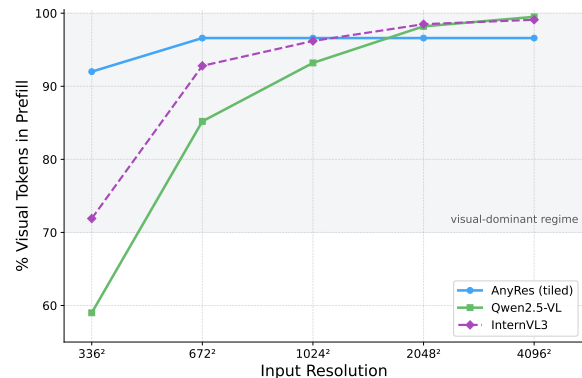


Figure 1: **Visual token dominance across resolutions.** Fraction of visual tokens relative to a fixed 100-token text prompt. As resolution increases, visual tokens quickly dominate the context window, particularly in dynamic-resolution models where scaling is quadratic. AnyRes refer to tiling of multiple views. More details in Appx A.1.

processed per request, which significantly impacts latency and memory consumption (Fig 1), even when the actual query may only require a coarse understanding of the scene.

A key observation is that *not all queries require the same visual granularity*. Coarse queries (e.g., “What is the breed of the dog?”) are typically answerable from a small image; fine-grained queries (e.g., “What is the name on the collar?”) benefit from higher resolution. Existing efficiency methods typically operate *after* tokenization, on the output of the vision encoder -pruning, pooling, merging, or compressing with Q-former style architecture (Arif et al., 2025; Zhang et al., 2025c; Xing et al., 2025; Lin et al., 2025; Rao et al., 2021; Liang et al., 2022; Bolya et al., 2023; Hu et al., 2025; Cai et al., 2025). While complementary, these methods typically operate on the output of the visual encoder alone and are unaware of the text input or the current query. Yet a more fundamental lever remains untouched: *Can we choose the input granularity as a pre-processing step?*

* Equal contribution

† Equal supervision

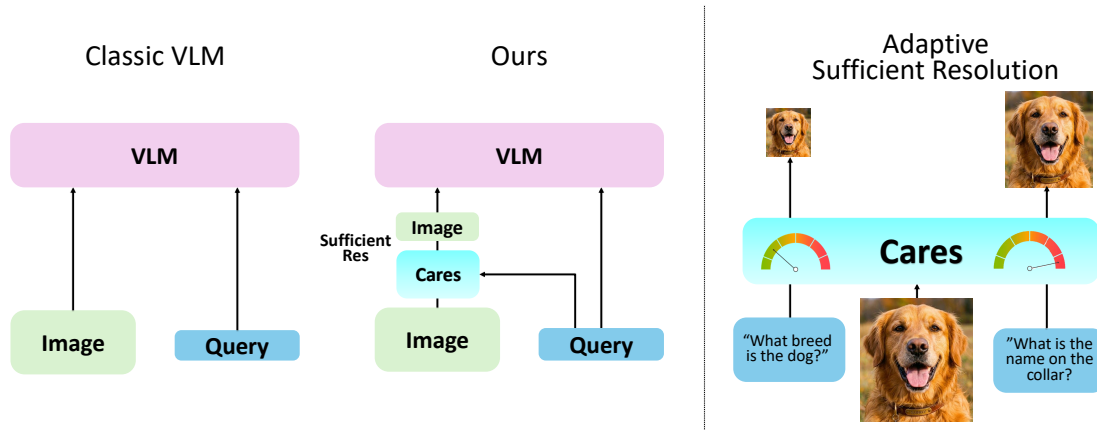


Figure 2: Overview of **CARES**. On the left, we compare the traditional pipeline of a use of VLM vs the pipeline using CARES. Given an image and its query, CARES predicts the minimal sufficient input resolution. The image is resized accordingly and, together with the query, passed to a downstream VLM. Coarse queries are routed to lower resolution; fine-grained queries that require more detail trigger higher resolution, which yields more visual tokens in the VLM.

We propose a *Context-Aware Resolution Selector* (**CARES**), a lightweight model that, for a given image-query pair, selects the *minimal* sufficient resolution to answer the query (Fig. 2). CARES is model-agnostic, placed *in front of* an arbitrary VLM. While our main instantiation uses a compact frozen VLM with a lightweight discriminative classifier, the CARES formulation is not tied to a specific predictor architecture. We also study a closely related autoregressive instantiation based on Granite-Docling, fine-tuned with LoRA, and report it separately on document-centric benchmarks.

It operates in three steps:

- A cheap low-resolution pass (e.g., $\leq 384^2$) extracts a joint image–query representation using a small proxy VLM.
- Given this representation, a lightweight classifier predicts the minimal resolution required for the task.
- The image is resized to the predicted resolution and passed to the target VLM. No changes to the VLM’s architecture, weights, or training are required.

A central challenge is supervision: what resolution is *truly* sufficient for each example? We introduce a simple labeling procedure based on a discrete set of resolutions \mathcal{R} and a task performance metric. For each image, query, and GT response, we evaluate a pretrained VLM with increasingly higher resolution up to convergence in

terms of the task metric (or reaching the native resolution). The lowest resolution at which the convergence occurs is selected as the ground-truth optimal resolution for training CARES. Using a discrete resolution set avoids the cost of exhaustively searching over continuous values. Since the labels are discrete, the model is trained as a classifier. At inference time, however, we interpolate between the predicted class probabilities to recover a continuous resolution estimate.

Across 9 multimodal benchmarks, varying from natural images to document understanding (Section 4) and different open and api-based model, CARES reduces average visual tokens and GFLOPS by 70-80%, with minimal to no accuracy drop compared to always using the highest (native) resolution.

Our contributions are as follows:

1. We define the task of *query- and image-conditioned resolution selection* for vision-language models, aimed at reducing input size without sacrificing accuracy.
2. We propose a simple yet effective supervision strategy based on multi-resolution rollouts and a convergence rule, yielding per-example sufficient resolution ground-truth, enabling training and evaluation.
3. We introduce CARES, a lightweight, model-agnostic module that selects resolution as a pre-processing step, requiring no changes to the target VLM.

4. We demonstrate that many visual tokens are unnecessary: CARES preserves performance across tasks while reducing visual compute by up to 78% on average across 9 benchmarks, and is orthogonal with post-tokenization token compression.

2 Related Work

Visual-token sparsification at inference A growing line of work trims visual tokens *after* tokenization inside the VLM stack. HiRED uses [CLS] attention to allocate a per-partition token budget and drop the least-informative vision tokens under a fixed budget, yielding large speedups on high-resolution inputs without retraining (Arif et al., 2025). SparseVLM proposes a training-free, text-guided strategy: self-attention matrices rank visual tokens with an adaptive layer-wise sparsification ratio and a token-recycling mechanism to preserve information (Zhang et al., 2025c). PyramidDrop stages the model and progressively reduces tokens at stage boundaries, motivated by the observation that redundancy increases with depth; it accelerates both training and inference and can also be used in a plug-and-play inference mode (Xing et al., 2025). Complementary to these, Visual Tokens Withdrawal (VTW) argues that visual information migrates to text tokens in early layers and thus withdraws vision tokens beyond a learned layer, cutting compute while maintaining quality (Lin et al., 2025). In contrast, CARES decides *before* tokenization which input resolution to use and leaves all VLM’s components frozen.

Training for flexible token budgets TokenFLEX trains VLMs to operate across a range of visual–token counts by stochastically modulating tokens during training and adding a lightweight projector with adaptive pooling (Hu et al., 2025). *Matryoshka Multimodal Models* (MMM) further pursue elastic compute, training nested representations that remain useful under progressively smaller token/feature budgets (Cai et al., 2025). *LLaVA-Mini* pushes efficiency to the extreme by compressing visual information into (nearly) a single vision token while retaining competitive performance for both images and videos (Zhang et al., 2025b). CARES targets the complementary axis of *adaptive pixel allocation* before tokenization: it selects the minimal input resolution needed for a target utility and can front-end TokenFLEX/Matryoshka/LLaVA-Mini-style models to reduce pixels (and thus to-

kens) further.

Any-resolution inputs and tiling Many modern ViTs (Dehghani et al., 2023; Beyer et al., 2023) and VLMs boost fine-grained perception with AnyRes/dynamic-high-resolution tiling (e.g., LLaVA-NeXT) or native dynamic resolution that maps larger images to more tokens (e.g., Qwen2-VL) (Liu et al., 2024a; Wang et al., 2024). While effective, these strategies often increase visual tokens substantially. CARES explicitly *avoids* unnecessary tiling by routing easy cases to low resolutions and only escalating when the query and low-res cues predict a benefit.

Dynamic computation Vision-only methods reduce computation via token pruning/merging inside ViTs-e.g., DynamicViT prunes tokens hierarchically with learned importance (Rao et al., 2021), EViT reorganizes/discards inattentive tokens (Liang et al., 2022), and ToMe merges similar tokens on the fly (Bolya et al., 2023). WAVE-CLIP replaces patch tokenization with a multi-level wavelet tokenizer and performs coarse-to-fine inference in a single ViT (Kimhi et al., 2025). For VLMs, SGL routes easy cases via a small ‘stitch’ model and defers hard ones to a larger counterpart, akin to early-exit routing (Zhao et al., 2024). These operate *within* the encoder after tokenization; CARES is complementary, deciding how many pixels to tokenize in the first place.

Adaptive input resolution selection Outside VLMs, dynamic-resolution networks learn a per-image resolution predictor that trades accuracy for cost in classification (Zhu et al., 2021). CARES brings this idea to multimodal QA, conditions the policy on the query text, and supervises it with *per-example* multi-resolution rollouts of the target VLM using a sufficiency rule, which yields unambiguous labels at deployment resolutions.

Extreme compression and design insights Recent analyses argue that, under fixed inference budgets, compute-optimal VLMs may prefer very few visual tokens and a larger LLM (Li et al., 2024). Such results support approaches that minimize visual tokens when possible; methods like *LLaVA-Mini* instantiate the “one-token vision” regime in practice (Zhang et al., 2025b). CARES provides a query-conditioned mechanism to reduce pixels upstream, complementing these token-minimal designs.

3 CARES

This section outlines the problem addressed by CARES (3.1), followed by a description of the dataset generation procedure (3.2). We then detail the architecture and the training details of CARES (3.3). Finally we outline our continuous resolution approach (3.4).

3.1 Problem Definition

Given an image x and query q , let $\mathcal{R} = [r_{\min}, r_{\max}] \subset \mathbb{R}^+$ denote the range of valid input resolutions and let F be a fixed VLM. For any resolution $r \in \mathcal{R}$, we denote by $x^{(r)}$ the image x resized such that its largest dimension equals r . Feeding $x^{(r)}$ and q into F yields an output $y = F(x^{(r)}, q)$. The VLM forms $T(r)$ visual tokens at resolution r (including AnyRes/tiling effects). Our goal is to learn a *selector* f_θ that predicts, from a single inexpensive low-resolution pass at r_{\min} , the minimal *sufficient* resolution $r_s \in \mathcal{R}$ for accurately answering the query q given image x .

3.2 Labeling Strategy for Training CARES

Since searching for the optimal $r^* \in \mathcal{R}$ is prohibitively expensive, we chose to use a small, discrete set of valid resolutions for the annotation $\mathcal{R}_d = \{r_1, \dots, r_K\} \subset \mathcal{R}$. For each sample, we render the image at the fixed resolutions, \mathcal{R}_d , and use a pretrained VLM to generate predictions at each resolution. The predictions are evaluated against the ground-truth annotations using the ANLS metric. The supervision label is assigned as the lowest resolution whose ANLS score exceeds a threshold, without significant improvement at higher resolutions. The procedure yields a *discrete* sufficiency label $r^* \in \mathcal{R}_d$ per example. We emphasize that discretization is only used for supervision efficiency; at inference, we deploy a *continuous* finer-grained selector (§3.4). Algorithm 1 outlines the data generation process, and Table. 1 visualizes the concept.

Formally, we compute the ANLS score for each resolution:

$$u_k = \text{ANLS}\left(F(x^{(r_k)}, q), \text{gt}\right) \in [0, 1] \quad (1)$$

and select the minimal sufficient resolution as:

$$r^* = \min \left\{ r_k \mid u_k \geq \tau, \max_{\ell > k} (u_\ell - u_k) \leq \delta \right\} \quad (2)$$

where we default to r_K if no resolution satisfies the condition. We set $\tau = 0.85$ and use a small

Algorithm 1: Labeling via multi-resolution sufficiency rollouts.

Input: (x, q) ; resolutions \mathcal{R} ; VLM F ; utility U ; threshold τ ; margin δ

Output: Label $r^* \in \mathcal{R}$

for $k \leftarrow 1$ **to** K **do**

$y_k \leftarrow F(x^{(r_k)}, q)$; $u_k \leftarrow U(y_k, \text{gt})$

for $k \leftarrow 1$ **to** K **do**

if $u_k \geq \tau$ **and** $\max_{\ell > k} (u_\ell - u_k) \leq \delta$

then

return $r^* \leftarrow r_k$

return $r^* \leftarrow r_K$

margin δ (e.g., 0.1) to prevent rewarding negligible performance improvements. We define the full resolution range as $\mathcal{R} = [384, 1024]$, and use a discrete set $\mathcal{R}_d = \{384, 768, 1024\}$ for annotation.

3.3 Model Instantiations

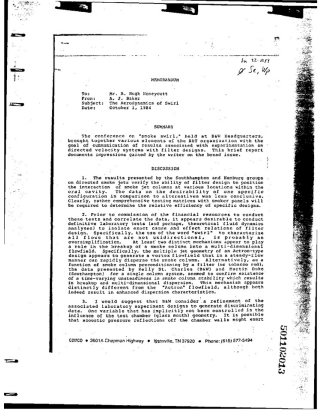
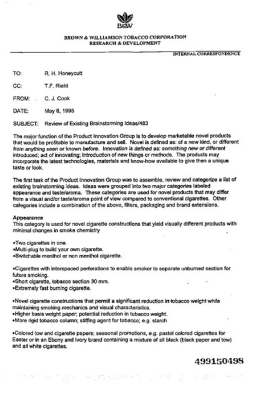
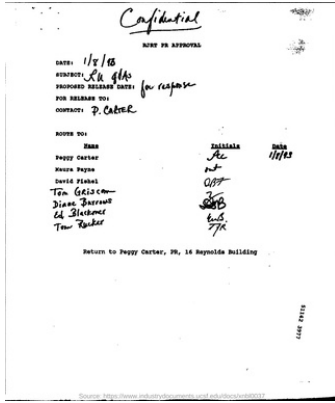
Unless otherwise stated, all main experiments in this paper use the following discriminative instantiation of CARES.

We design CARES as a lightweight resolution selector that can be deployed in front of any vision–language model (VLM) to improve efficiency. Its behavior is governed by three core principles:

1. **Compactness:** minimal overhead in computation and memory.
2. **Preprocessing role:** determines resolution directly from raw inputs before invoking the VLM.
3. **VLM-agnosticism:** works with any VLM, whether run locally or accessed via API, with no architecture changes or retraining required.

To implement these principles, we use a compact frozen VLM backbone as a joint vision–text feature extractor, followed by a lightweight classifier head.

Specifically, we adopt the pretrained SmolVLM-500M model (Marafioti et al., 2025), with layers 17–32 removed, as the backbone. Given an image at resolution r_{\min} and a text query, we feed both into the model and extract the hidden state of the final token at layer 16. This representation encodes the joint image–query context and is passed to a classifier that outputs a soft distribution over target resolutions. This design is motivated by recent findings showing that intermediate layer activations



Query	what is the contact person name mentioned in letter?	Who is in cc in this letter?	One variable that has implicitly not been controlled?
GT	P. Carter	T.F. Riehl	influence of the test chamber (glass mouth) geometry.
Resp@384 ANLS	P. Carter 1.0	T.F. Rosel 0.7	concentration of the final product 0.0
Resp@768 ANLS		T.F. Riehl 1.0	the influence of the test chamber (i.e. ash seath) geometry on the flow 0.65
Resp@1024 ANLS			the influence of the test chamber (glass mouth) geometry. 0.93
Sufficient Resolution	384	768	1024

Table 1: Data generation pipeline for training CARES. We process each input through a pretrained VLM (Granite-Vision) at three fixed resolutions and select the smallest resolution that produces a sufficient answer quality according to the ANLS metric.

in LLMs and VLMs encode rich perceptual and semantic information that may not be surfaced at the output layer (Orgad et al., 2024; Zhang et al., 2025a). In addition to being more informative, as also evidenced by the performance gap in Table 3 where using intermediate features outperforms last-layer features by about 1%, this choice substantially reduces computation since only roughly half of the LLM is used for feature extraction.

The resulting CARES module has approximately 350M parameters and is trained with supervision over discrete resolution labels (see §3.2).

Autoregressive document-specialized instantiation. In addition to the discriminative selector above, we also instantiate CARES using an autoregressive vision-language model. Concretely, we start from Granite-Docling-258M (Auer et al., 2024) and fine-tune it with LoRA (rank 8) on the same resolution-selection training set. Given the low-resolution image and the query, the model is prompted to predict one resolution label from the discrete set $\mathcal{R}_d = \{384, 768, 1024\}$. To avoid tokenization ambiguity, we map these labels to dedi-

cated tokens $\langle 1 \rangle$, $\langle 2 \rangle$ and $\langle 3 \rangle$.

At inference time, we read the first-step logits over the resolution tokens, apply a softmax to obtain class probabilities, and use the same expectation-based interpolation described in Eq. 3 to produce a continuous resolution. This preserves the deployment rule of CARES while replacing the discriminative head with an autoregressive predictor.

3.4 From Discrete Supervision to a Continuous Resolution

Although CARES is trained as a K -way classifier over a discrete set of resolutions $\mathcal{R}_d = \{r_1 < \dots < r_K\}$, we deploy it as a *continuous* selector over $\mathcal{R} = [r_{min}, r_{max}]$. Given features z from the low-resolution image and query, compute logits $\ell(z) \in \mathbb{R}^K$ and class probabilities

$$p = \text{softmax}(\ell),$$

We use the probability-weighted expectation

over \mathcal{R}_d :

$$\tilde{r} = \sum_{k=1}^{|\mathcal{R}_d|} p_k r_k, \quad (3)$$

This yields a *continuous* resolution that varies smoothly with confidence and is insensitive to the specific discretization used for labeling. In practice, \tilde{r} preserves the routing behavior of the classifier while allowing finer control.

Algorithm 2: Continuous resolution selection.

Input: (x, q) ; low-res r_1 ; logits ℓ .

Output: Continuous resolution

$$\tilde{r} \in [r_1, r_K].$$

$z \leftarrow$ features from proxy VLM at r_1

$p \leftarrow \text{softmax}(\ell(z))$

$\tilde{r} \leftarrow \sum_{k=1}^K p_k r_k$

return \tilde{r}

Continuous inference algorithm.

Deployment. The target VLM receives x with the largest dimension resized to \tilde{r} (or to the nearest supported side length to avoid under-allocation). For backbones that only accept a discrete set of input sizes, we round *up* to the next supported size.

4 Results & Analysis

This section presents the experimental evaluation of CARES. We begin by describing the benchmarks and evaluation metrics (4.1), followed by the main results (4.2), and finally a comprehensive ablation study (4.4).

4.1 Experimental Setup

Training Data To train the resolution selector, we construct a dataset of images and queries (x, q) we automatically annotated with the minimal sufficient resolution r^* . We construct an 80K-sample training set by randomly sampling 20K instances from each of four datasets: TextVQA (Singh et al., 2019), ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), and LLaVA-Multi (Jiang et al., 2024), covering documents and natural images domains.

Training details We train CARES on the curated data described in 3.2 for 6 epochs using a learning rate of $1e - 3$ and a batch size of 32. We optimize the standard cross-entropy loss over the fixed resolution labels:

$$\mathcal{L}(\theta) = \text{CE}(f_\theta(z), r^*).$$

Where $f_\theta(z)$ is CARES composed of a frozen VLM and the lightweight classifier. In addition, we apply label smoothing of 0.05 to support continuous resolutions at inference time.

VLM variant training details. For the autoregressive (AR) Granite-Docling instantiation, we use the same training set and the same discrete supervision labels. The model is fine-tuned with LoRA of rank 8, while the base model remains frozen. Training is performed with next-token supervision over the resolution tokens, and for efficiency, generation length is set to 1. Learning rate is set to $1e - 5$ and a batch size of 64 for 3 epochs.

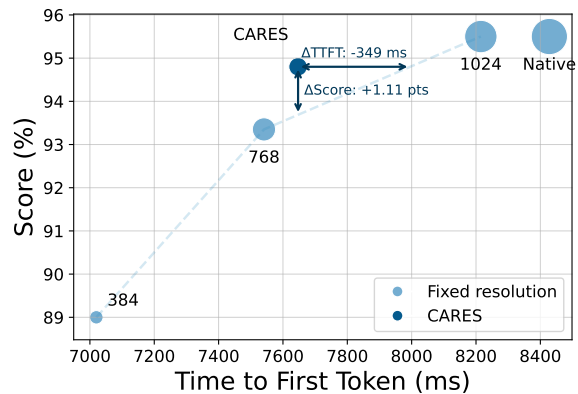


Figure 3: Accuracy vs. TTFT for DocVQA with Qwen2.5-VL-72B across native and fixed-resolution settings versus CARES. Bubble size indicates the number of pixels processed by the model.

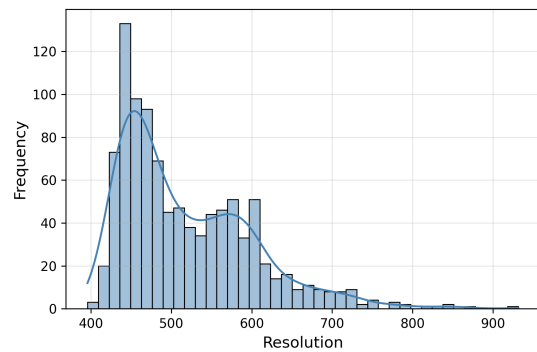


Figure 4: Histogram of the predicted resolutions \tilde{r} by CARES for OCRBench.

Evaluation We evaluate on nine public benchmarks varying from documents to natural images: Ai2D (Kembhavi et al., 2016), ChartQA (Masry et al., 2022), DocVQA (Mathew et al., 2021), OCRBench (Liu et al., 2024b), and SeedBench-2 (Li et al., 2023), MMMU (Yue et al., 2024), RealWorldQA (xAI, 2024), InfoVQA (Mathew et al.,

Model	Ai2D		ChartQA		DocVQA		OCRBench		SeedBench-2		MMMU		RealWorldQA		InfoVQA		MathVista		Average	
	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost	Score	Cost
Granite-Vision-2B	0.74		0.86		0.90		0.80		0.72		0.29		0.17		0.35		0.48		0.59	
+ CARES	0.73	-67%	0.87	-69%	0.90	-68%	0.80	-68%	0.72	-44%	0.29	-85%	0.19	-72%	0.40	-72%	0.48	-22%	0.60	-63%
+ CARES-AR	0.71	-81%	0.84	-81%	0.88	-82%	0.77	-75%	0.72	-10%	0.30	-84%	0.15	-82%	0.39	-81%	0.44	-25%	0.58	-67%
InternVL3-8B	0.84		0.86		0.92		0.85		0.79		0.56		0.68		0.72		0.69		0.77	
+ CARES	0.84	-66%	0.86	-68%	0.92	-69%	0.85	-70%	0.79	-44%	0.56	-86%	0.68	-82%	0.74	-72%	0.69	-22%	0.77	-64%
+ CARES-AR	0.84	-86%	0.86	-81%	0.92	-80%	0.85	-78%	0.72	-84%	0.55	-85%	0.68	-82%	0.74	-81%	0.68	-31%	0.76	-76%
Qwen2.5-VL-72B	0.87		0.87		0.96		0.75		0.81		0.62		0.77		0.73		0.74		0.79	
+ CARES	0.87	-85%	0.84	-77%	0.95	-84%	0.76	-64%	0.79	-77%	0.62	-86%	0.79	-82%	0.84	-72%	0.74	-7%	0.80	-70%
GPT-4o	0.78		0.56		0.80		0.77		0.76		0.57		0.61		0.75		0.64		0.69	
+ CARES	0.78	-60%	0.56	-60%	0.80	-36%	0.75	-33%	0.75	-47%	0.56	-85%	0.61	-84%	0.73	-76%	0.61	-17%	0.68	-55%
+ CARES-AR	0.74	-85%	0.52	-85%	0.78	-88%	0.73	-84%	0.71	-82%	0.56	-85%	0.62	-84%	0.71	-82%	0.58	-28%	0.66	-78%

Table 2: **Benchmark performance** and estimated prefill-stage savings for **Cost** (measured in FLOPS for local models or \$ for API models). Reporting CARES-AR for auto-regressive prediction. CARES was trained on document data, where other domain datasets shows similar performance with less aggressive Cost saved.

2022) and MathVista (Lu et al., 2024). For Ai2D, ChartQA, and SeedBench-2 we report exact-match accuracy. For DocVQA and OCRBench we report Average Normalized Levenshtein Similarity (ANLS). All evaluations were performed with the standard Imms-eval (Zhang et al., 2024) setup. We also report a macro-averaged Performance (%) across all datasets.

4.2 Main results

We evaluate CARES across **Granite-Vision 3.3-2B** (Team et al., 2025), **InternVL3-8B** (Zhu et al., 2025), **Qwen2.5-VL-72B** (Bai et al., 2025), and **GPT-4o** (Achiam et al., 2023). We also report prefill-stage FLOPS savings for locally run models, and estimated dollar savings in API usage for GPT-4o. As summarized in Table 2, CARES maintains accuracy while cutting prefill compute: averaged over models and datasets, prefill FLOPs drop by **65–85%** with at most a sub-point change in macro performance relative to always using the highest/native resolution. The effect is consistent for compact (Granite-Vision 3.3-2B) and large (Qwen2.5-VL-72B) backbones, and holds for GPT-4o accessed via API (accuracy parity at comparable quality).

Fig. 3 shows the accuracy–latency frontier: CARES matches near-native accuracy while using far fewer TFLOPs (e.g., 2.58 vs. 7.5) and achieving ~ 1 second lower time-to-first-token (TTFT); static high-res inputs (e.g., 1024^2) incur substantial compute with limited TTFT gains, whereas fixed low-res (384^2) improves TTFT at the cost of quality. The query-aware routing yields a superior Pareto point.

Finally, the distribution of predicted continuous resolutions \tilde{r} (Fig. 4) and the comparison in Table 5 indicate that continuous routing adapts per instance, matches or slightly improves accuracy

over a discrete menu, and saves additional compute without quality loss.

4.3 Cross-Teacher Agreement for Resolution Labels

Because our supervision is generated by rolling out a pretrained VLM at multiple resolutions, one natural question is whether the resulting labels depend strongly on the specific annotating model. To test this, we compare labels generated by two substantially different teachers: **Granite-Vision-2B** and **Qwen3-VL-235B**, on a shared subset of 1000 examples.

We find a high degree of agreement between the two annotators. The two teachers predict the same sufficient resolution for more than **95%** of examples, with Pearson correlation **0.908** and mutual information **1.116** between their predicted sufficiency levels. The confusion matrix is shown in Table 5. These results suggest that the notion of sufficient resolution is largely shared across architectures and scales, and is not tied to a single model family.

This result complements the downstream transfer results in the main paper, where a selector trained using labels derived from one setup transfers well across multiple target VLMs. Together, these findings support the view that CARES captures a broadly shared notion of *resolution adequacy*, rather than overfitting to one teacher’s idiosyncrasies.

4.4 Ablation study

We conduct a series of ablations to isolate the effect of key training design choices on resolution selection accuracy and downstream benchmark performance.

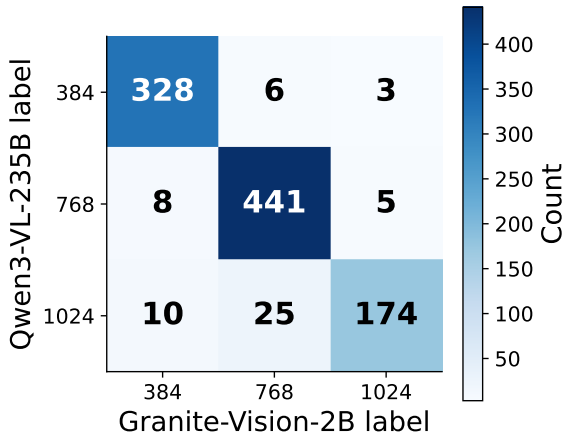


Figure 5: Confusion matrix between sufficient-resolution labels generated by Granite-Vision-2B and Qwen3-VL-235B on a shared subset. Most mass lies on the diagonal, indicating strong agreement across teachers.

Feature extractor. We ablate several frozen backbones used for feature extraction in CARES, varying both model type and layer depth. As shown in Table 3, both Qwen2.5-3B and SmolVLM achieve higher accuracy when using intermediate-layer features, outperforming their own final-layer variants. This aligns with prior findings suggesting that intermediate representations in VLMs often encode richer signals than final outputs.

Qwen2.5-3B and SmolVLM both process the image and query jointly within a unified transformer, in contrast to SigLIP v2’s dual-encoder architecture, where vision and language are encoded separately. For SigLIP, we follow the original design by pooling the outputs of each tower, concatenating them, and passing the result to the classifier head. While this setup is architecturally simple, it underperforms joint encoding by a considerable margin (56.1% accuracy), and it requires more parameters than the lightweight SmolVLM.

Although Qwen2.5-3B achieves the best overall accuracy, we adopt SmolVLM as our default backbone due to its favorable trade-off between performance, size, and efficiency, making it a more practical choice for real-world pre-processing.

Resolution menu size. We compare training with binary $\mathcal{R}_d = \{384, 1024\}$ ($|\mathcal{R}_d| = 2$) vs. ternary $\mathcal{R}_d = \{384, 768, 1024\}$ ($|\mathcal{R}_d| = 3$) resolution choices. Table 4 reports both the classification accuracy and the downstream performance of Granite Vision, averaged over 5 benchmarks. As expected, the two-way classification yields higher

Model	Layer	Params	Accuracy
SigLIP v2	-	0.8B	56.1%
SmolVLM	Mid	0.35B	63.3%
SmolVLM	Last	0.5B	62.3%
Qwen2.5-3B	Mid	2.3B	67.2%
Qwen2.5-3B	Last	3.75B	66.2%

Table 3: **Feature extractor.** Validation accuracy and parameter count for different frozen feature extractors used in CARES. All models are trained to classify among three resolution choices. For SmolVLM and Qwen2.5-3B, we compare features extracted from intermediate (MID) and final (LAST) layers. For SigLIP, the pooled outputs from the vision and language towers are concatenated and passed to the classifier head. Qwen2.5-3B provides the best performance, while SmolVLM offers strong accuracy with minimal size.

validation accuracy in the resolution classification task compared to the more challenging three-way classification. But the ternary setup leads to better downstream benchmark performance due to the finer-grained control.

$ \mathcal{R}_d $	Resolution Accuracy	Downstream Accuracy
2	96.2%	0.76
3	67.2%	0.80

Table 4: **Binary vs. Ternary Resolution Classification.** We compare binary ($|\mathcal{R}_d| = 2$, using $\{384, 1024\}$) and ternary ($|\mathcal{R}_d| = 3$, using $\{384, 768, 1024\}$) resolution selection setups. The binary classifier achieves higher accuracy on the resolution prediction task due to its reduced complexity, while the ternary classifier improves downstream performance by enabling finer control over resolution. Reported downstream accuracy is averaged over 5 vision-language benchmarks using Granite Vision.

Discrete vs. continuous. CARES is trained as a discrete resolution classifier, but at inference time, it can produce either discrete predictions or a continuous estimate via interpolation. In Table 5, we compare the impact of discrete versus continuous inference across three VLM backbones. All scores and FLOPS deltas are averaged over nine benchmarks. We find that continuous resolution selection achieves comparable accuracy to both discrete and native strategies, while significantly reducing compute. For example, with Granite-Vision 3.3-2B and InternVL3-8B, FLOPS are reduced by 63% us-

ing continuous prediction, compared to 46% with discrete. These results suggest that continuous inference allows finer control over input resolution and leads to more efficient inference without compromising performance.

Model	Resolution	Score	FLOPS
Granite-Vision 3.3-2B	Native	0.803	
	Discrete	0.801	-46%
	Continuous	0.804	-63%
InternVL3-8B	Native	0.851	
	Discrete	0.851	-46%
	Continuous	0.851	-63%
Qwen2.5-VL-72B	Native	0.851	
	Discrete	0.852	-74%
	Continuous	0.839	-80%

Table 5: **Discrete vs. Continuous Resolution Selector.** The overall score and relative FLOPS delta per resolution strategy are averaged over 5 benchmarks. Using continuous resolutions allows finer control of the resolution, resulting in a lower resolution and computation with no drop in accuracy.

Label smoothing. To bridge the mismatch between *discrete* supervision and our *continuous* inference policy, we apply label smoothing when training the classifier over \mathcal{R}_d . Smoothing softens class boundaries and discourages over-confident logits, yielding better-calibrated probability distributions p that are subsequently mapped to a scalar resolution via expectation (Eq. 3). This improves the stability of the continuous selector, reduces spurious hard escalations near decision thresholds, and translates to higher downstream utility at similar—or lower—compute. Empirically, Table 6 shows that adding label smoothing improves OCR-Bench performance for Qwen2.5-VL-7B (0.821 vs. 0.811) while slightly *reducing* expected FLOPS, supporting its role as a simple but effective regularizer for continuous-resolution deployment.

Setting	Score	FLOPS
Native resolution	0.824	
CARES Without label-smoothing	0.811	-60.5%
CARES With label-smoothing	0.821	-63.8%

Table 6: **Label smoothing effect.** Evaluated on OCR-Bench with Qwen2.5-VL-7B. Comparison of native resolution and training with or without label smoothing. FLOPs indicate relative change.

5 Discussion and Conclusion

Inference efficiency has become a critical concern for modern vision-language systems. Most user queries do not require high-resolution inputs, yet current deployments often process all images at native or tiled resolutions by default. This leads to bloated token counts, slower response times, and higher costs. CARES addresses this challenge with a lightweight, model-agnostic approach that dynamically selects input resolution based on the query. By acting before tokenization, it provides a clean and practical lever for controlling inference cost while maintaining output quality.

Key Takeaways

- CARES reduces compute and latency across a wide range of models and benchmarks, with minimal to no loss in task accuracy.
- It requires no changes to the vision-language model and works as a plug-in component, making it easy to integrate into real-world pipelines.
- CARES adapts resolution based on the specific query, using a single low-cost pass to determine how much visual detail is needed.
- The design is compact and efficient, enabling wide applicability without adding large overhead to the main model.

Overall, CARES highlights the value of adaptive pixel allocation as a simple yet powerful strategy for efficient multimodal inference. It complements existing techniques for token-level compression and opens up a new path for practical deployment of vision-language models at scale.

Limitations

CARES depends on a frozen proxy VLM for low-resolution features; domains requiring extremely fine cues (e.g., dense OCR, medical imagery) may be under-allocated. Our supervision uses multi-resolution rollouts of a target VLM and thus inherits that model’s biases and limited language support. We evaluate single-image, single-turn inputs only; multi-image, video, streaming, and joint resolution-tiling selection are left to future work. We do not study safety, robustness to adversarial prompts, or detailed cost-latency trade-offs across hardware.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Kazi Hasan Ibn Arif, JinYi Yoon, Dimitrios S Nikolopoulos, Hans Vandierendonck, Deepu John, and Bo Ji. 2025. Hired: Attention-guided token dropping for efficient inference of high-resolution vision-language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 1773–1781.
- Christoph Auer, Maksym Lysak, Ahmed Nassar, Michele Dolfi, Nikolaos Livathinos, Panos Vagenas, Cesar Berrospi Ramis, Matteo Omenetti, Fabian Lindlbauer, Kasper Dinkla, Lokesh Mishra, Yusik Kim, Shubham Gupta, Rafael Teixeira de Lima, Valery Weber, Lucas Morin, Ingmar Meijer, Viktor Kuropiatnyk, and Peter W. J. Staar. 2024. [Docling technical report](#). *Preprint*, arXiv:2408.09869.
- 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, and 8 others. 2025. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*.
- Lucas Beyer, Pavel Izmailov, Alexander Kolesnikov, Mathilde Caron, Simon Kornblith, Xiaohua Zhai, Matthias Minderer, Michael Tschannen, Ibrahim Alabdulmohsin, and Filip Pavetic. 2023. [Flexivit: One model for all patch sizes](#). *Preprint*, arXiv:2212.08013.
- Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. 2023. Token merging: Your ViT but faster. In *International Conference on Learning Representations*.
- Mu Cai, Jianwei Yang, Jianfeng Gao, and Yong Jae Lee. 2025. Matryoshka multimodal models. *Proceedings of the International Conference on Learning Representation*.
- Mostafa Dehghani, Basil Mustafa, Josip Djolonga, Jonathan Heek, Matthias Minderer, Mathilde Caron, Andreas Peter Steiner, Joan Puigcerver, Robert Geirhos, Ibrahim Alabdulmohsin, Avital Oliver, Piotr Padlewski, Alexey A. Gritsenko, Mario Lucic, and Neil Houlsby. 2023. Patch n’ pack: Navit, a vision transformer for any aspect ratio and resolution. In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Junshan Hu, Jialiang Mao, Zhikang Liu, Zhongpu Xia, Peng Jia, and Xianpeng Lang. 2025. [Tokenflex: Unified vlm training for flexible visual tokens inference](#). *Preprint*, arXiv:2504.03154.
- Dongfu Jiang, Xuan He, Huaye Zeng, Cong Wei, Max Ku, Qian Liu, and Wenhua Chen. 2024. Mantis: Interleaved multi-image instruction tuning. *arXiv preprint arXiv:2405.01483*.
- Aniruddha Kembhavi, Mike Salvato, Eric Kolve, Minjoon Seo, Hannaneh Hajishirzi, and Ali Farhadi. 2016. [A diagram is worth a dozen images](#). *Preprint*, arXiv:1603.07396.
- Moshe Kimhi, Erez Koifman, Ehud Rivlin, Eli Schwartz, and Chaim Baskin. 2025. [Waveclip: Wavelet tokenization for adaptive-resolution clip](#). *Preprint*, arXiv:2509.21153.
- Bohao Li, Yuying Ge, Yixiao Ge, Guangzhi Wang, Rui Wang, Ruimao Zhang, and Ying Shan. 2023. Seed-bench-2: Benchmarking multimodal large language models. *arXiv preprint arXiv:2311.17092*.
- Kevin Y. Li, Sachin Goyal, Joao D. Semedo, and J. Zico Kolter. 2024. [Inference optimal vlms need only one visual token but larger models](#). *Preprint*, arXiv:2411.03312.
- Youwei Liang, Chongjian Ge, Zhan Tong, Yibing Song, Jue Wang, and Pengtao Xie. 2022. [Not all patches are what you need: Expediting vision transformers via token reorganizations](#). *Preprint*, arXiv:2202.07800.
- Zhihang Lin, Mingbao Lin, Luxi Lin, and Rongrong Ji. 2025. Boosting multimodal large language models with visual tokens withdrawal for rapid inference. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 5334–5342.
- Haotian Liu, Chunyuan Li, Yuheng Li, Bo Li, Yuanhan Zhang, Sheng Shen, and Yong Jae Lee. 2024a. [Llava-next: Improved reasoning, ocr, and world knowledge](#).
- Yuliang Liu, Zhang Li, Mingxin Huang, Biao Yang, Wenwen Yu, Chunyuan Li, Xu-Cheng Yin, Cheng-Lin Liu, Lianwen Jin, and Xiang Bai. 2024b. [Ocr-bench: on the hidden mystery of ocr in large multimodal models](#). *Science China Information Sciences*, 67(12).
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. 2024. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts. In *International Conference on Learning Representations (ICLR)*.
- Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, Vaibhav Srivastav, Joshua Lochner, Hugo Larcher, Mathieu Morlon, Lewis Tunstall, Leandro von Werra, and Thomas Wolf. 2025. Smolvlm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*.

- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. Chartqa: A benchmark for question answering about charts with visual and logical reasoning. *arXiv preprint arXiv:2203.10244*.
- Minesh Mathew, Viraj Bagal, Rubèn Pérez Tito, Dimosthenis Karatzas, Ernest Valveny, and C.V. Jawahar. 2022. Infographicvqa. In *2022 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 2582–2591.
- Minesh Mathew, Dimosthenis Karatzas, and CV Jawahar. 2021. Docvqa: A dataset for vqa on document images. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2200–2209.
- Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Reichart, Idan Szepkator, Hadas Kotek, and Yonatan Belinkov. 2024. Llm know more than they show: On the intrinsic representation of llm hallucinations. *arXiv preprint arXiv:2410.02707*.
- Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. 2021. Dynamicvit: Efficient vision transformers with dynamic token sparsification. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Dhruv Batra, Devi Parikh, and Marcus Rohrbach. 2019. Towards vqa models that can read. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8317–8326.
- Granite Vision Team, Leonid Karlinsky, Assaf Arbelle, Abraham Daniels, Ahmed Nassar, Amit Alfassi, Bo Wu, Eli Schwartz, Dhiraj Joshi, Jovana Kondic, and 1 others. 2025. Granite vision: a lightweight, open-source multimodal model for enterprise intelligence. *arXiv preprint arXiv:2502.09927*.
- IBM Granite Vision Team. 2026. [Granite 4.0 vision](#).
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.
- xAI. 2024. Realworldqa. <https://x.ai/news/grok-1.5v>.
- Long Xing, Qidong Huang, Xiaoyi Dong, Jiajie Lu, Pan Zhang, Yuhang Zang, Yuhang Cao, Conghui He, Jiaqi Wang, Feng Wu, and Dahua Lin. 2025. Pyramidrop: Accelerating your large vision-language models via pyramid visual redundancy reduction. *Preprint*, arXiv:2410.17247.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, Cong Wei, Botao Yu, Ruibin Yuan, Renliang Sun, Ming Yin, Boyuan Zheng, Zhenzhu Yang, Yibo Liu, Wenhao Huang, and 3 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of CVPR*.
- Jiarui Zhang, Mahyar Khayatkhoei, Prateek Chhikara, and Filip Ilievski. 2025a. [MLLMs know where to look: Training-free perception of small visual details with multimodal LLMs](#). In *The Thirteenth International Conference on Learning Representations*.
- Kaichen Zhang, Bo Li, Peiyuan Zhang, Fanyi Pu, Joshua Adrian Cahyono, Kairui Hu, Shuai Liu, Yuanhan Zhang, Jingkang Yang, Chunyuan Li, and Ziwei Liu. 2024. [Lmms-eval: Reality check on the evaluation of large multimodal models](#). *Preprint*, arXiv:2407.12772.
- Shaolei Zhang, Qingkai Fang, Zhe Yang, and Yang Feng. 2025b. Llava-mini: Efficient image and video large multimodal models with one vision token. In *International Conference on Learning Representations (ICLR)*.
- Yuan Zhang, Chun-Kai Fan, Junpeng Ma, Wenzhao Zheng, Tao Huang, Kuan Cheng, Denis Gudovskiy, Tomoyuki Okuno, Yohei Nakata, Kurt Keutzer, and 1 others. 2025c. Sparsevlm: Visual token sparsification for efficient vision-language model inference. In *International Conference on Machine Learning*.
- Wangbo Zhao, Yizeng Han, Jiasheng Tang, Zhikai Li, Yibing Song, Kai Wang, Zhangyang Wang, and Yang You. 2024. A stitch in time saves nine: Small vlm is a precise guidance for accelerating large vlms. *arXiv preprint arXiv:2412.03324*.
- Jinguo Zhu, Weiyun Wang, Zhe Chen, Zhaoyang Liu, Shenglong Ye, Lixin Gu, Hao Tian, Yuchen Duan, Weijie Su, Jie Shao, and 1 others. 2025. Internvl3: Exploring advanced training and test-time recipes for open-source multimodal models. *arXiv preprint arXiv:2504.10479*.
- Mingjian Zhu, Kai Han, Enhua Wu, Qiulin Zhang, Ying Nie, Zhenzhong Lan, and Yunhe Wang. 2021. Dynamic resolution network. In *Advances in Neural Information Processing Systems (NeurIPS)*, Red Hook, NY, USA.

A Additional Analysis and Results

This appendix provides additional qualitative and quantitative analysis of CARES.

A.1 Extended token count evaluation

A.1.1 Textual token statistics

We analyze the number of textual tokens across the evaluated benchmarks using the Qwen2.5-VL tokenizer. For each dataset, we compute the average number of tokens in the full input prompt, including the question, answer choices (when applicable), and the instruction suffix (e.g., “*Answer with the option’s letter...*”). The results are summarized in Table 8.

We observe that most benchmarks contain relatively short textual inputs, typically in the range of 20–50 tokens. The longest prompts appear in MMMU, with an average of approximately 100 tokens, due to its multi-choice and instruction-heavy format.

Based on this analysis, we adopt a fixed text length of $T = 100$ tokens in our token composition study (Section A.1). This choice reflects a conservative upper bound over the evaluated benchmarks and ensures that our analysis does not underestimate the contribution of textual tokens. In practice, this assumption is favorable to text, as most datasets contain substantially fewer tokens, further increasing the relative dominance of visual tokens in real-world settings.

For a fixed text prompt of $T = 100$ tokens, the visual-token fraction is

$$P(V; T) = 100 \cdot \frac{V}{V + T}.$$

For Qwen2.5-VL, we approximate the number of LLM-side visual tokens by

$$V_{\text{Qwen}}(H, W) = \left\lceil \frac{H}{28} \right\rceil \left\lceil \frac{W}{28} \right\rceil,$$

reflecting patch size 14 with spatial merge size 2.

For AnyRes tiled models (e.g., LLaVA-NeXT / Granite-Vision (Team, 2026) processing), a square-input abstraction is

$$V_{\text{AnyRes}}(s) = 576 \left(1 + \min \left\{ \left\lceil \frac{s}{336} \right\rceil^2, 4 \right\} \right),$$

where 576 is the per-image embedding length and the default square grid saturates at a 2×2 local tiling plus a global view.

Resolution	Model	Visual	Total	% Visual
336×336	AnyRes (Tiled)	1152	1252	92.0%
	Qwen2.5-VL	144	244	59.0%
	InternVL3	256	356	71.9%
672×672	AnyRes (Tiled)	2880	2980	96.6%
	Qwen2.5-VL	576	676	85.2%
	InternVL3	1280	1380	92.8%
1024×1024	AnyRes (Tiled)	2880	2980	96.6%
	Qwen2.5-VL	1369	1469	93.2%
	InternVL3	2560	2660	96.2%
2048×2048	AnyRes (Tiled)	2880	2980	96.6%
	Qwen2.5-VL	5476	5576	98.2%
	InternVL3	6656	6756	98.5%
4096×4096	AnyRes (Tiled)	2880	2980	96.6%
	Qwen2.5-VL	21609	21709	99.5%
	InternVL3	10496	10596	99.1%

Table 7: Visual vs. text token composition (assuming 100 text tokens) across representative VLM architectures. Qwen2.5-VL scales approximately quadratically with input size, AnyRes-style tiled models saturate due to a fixed square tiling budget, and InternVL3 follows a dynamic tile-based scaling with 448×448 tiles and an additional thumbnail view.

Dataset	Avg. Tokens
AI2D	44.3
ChartQA	21.6
DocVQA	21.2
OCRBench	22.5
SeedBench-2	45.8
MMMU	100.4
RealWorldQA	37.5
InfoVQA	24.0
MathVista	67.3

Table 8: Average number of textual tokens per benchmark. Token counts are computed using the Qwen2.5-VL tokenizer and include the full input prompt (question, answer choices when applicable, and instruction suffix such as “*Answer with the option’s letter...*”).

For InternVL3, we use a tile-based approximation consistent with its dynamic-resolution preprocessing:

$$V_{\text{Intern}}(s) = 256 \left(n(s) + \mathbf{1}[n(s) > 1] \right),$$

$$n(s) = \min \left(\left\lceil \frac{s}{448} \right\rceil^2, 40 \right),$$

where each 448×448 tile contributes 256 LLM-side tokens after pixel unshuffle, and an additional thumbnail is used when more than one tile is present.

These expressions explain the different regimes in Table 7: Qwen2.5-VL grows approximately quadratically with resolution, AnyRes saturates early due to a fixed tile budget, and InternVL3 follows a stepwise tile-based growth pattern.

Setting	Accuracy
Native	95.50
1024	95.50
768	93.35
384	89.00
CARES	94.80

Table 9: Comparison between fixed-resolution inference and CARES. While naive downscaling reduces accuracy, CARES recovers most of the native performance at much lower average compute.

A.2 Label generation pipeline

Figure 6 illustrates the supervision pipeline used to train CARES. For each image-query pair, we evaluate a pretrained VLM at several fixed resolutions and compare its prediction to the ground-truth answer. The smallest resolution whose score satisfies the sufficiency criterion is used as the training target. This process transforms downstream task behavior into per-example supervision for resolution selection, enabling CARES to learn when higher visual detail is genuinely needed.

A.3 Adaptive Selection vs. Fixed-Resolution Baselines

To disentangle the benefit of adaptive resolution selection from the general robustness of VLMs to downscaling, we compare CARES against fixed-resolution baselines. Table 9 reports an example comparison using [MODEL / benchmark setting], where all inputs are processed at a single fixed resolution.

These results clarify that the gains are not simply due to the model tolerating smaller images. Rather, the adaptive policy selectively preserves high resolution for the examples that need it, while routing easier cases to much smaller inputs. This explains why CARES achieves a better accuracy–efficiency trade-off than any single fixed-resolution operating

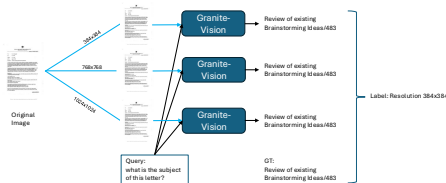


Figure 6: **Label generation pipeline for training CARES.** For each image-query pair, we evaluate a pretrained VLM at multiple fixed resolutions and assign the smallest resolution that satisfies the sufficiency criterion as the supervision label.

Table 10: **Time to First Token (TTFT, ms)** measured on H100 with batch size 1, averaged over 100 DocVQA examples. *Native* denotes the model’s default input pipeline. CARES reduces TTFT substantially compared to native and fixed high-resolution settings while preserving strong downstream accuracy.

Model	Native	1024 ²	768 ²	384 ²	CARES
Qwen2.5-VL-7B	435.7	433.8	220	76.12	270.1
Granite-Vision 3.3-2B	228.6	201.3	140.1	96.1	108.9

point.

A.4 Time-to-first-token analysis

Table 10 reports time-to-first-token (TTFT) on DocVQA for representative downstream VLMs. The results mirror the FLOPS trends in the main paper: lower resolutions substantially reduce latency, while CARES achieves a favorable trade-off by approaching the latency of low-resolution inference without incurring the accuracy loss of always using a small input. In particular, CARES significantly improves TTFT relative to native or fixed high-resolution processing, confirming that adaptive resolution selection translates into practical end-to-end inference gains.

A.5 Robustness to Proxy–Target Feature Mismatch

Since CARES uses a proxy VLM to extract low-resolution image-query features, one may worry that mismatch between the proxy representation and the downstream target model could introduce bias. To directly test this, we evaluate CARES on **Qwen2.5-VL-3B** using two feature extractors: (i) features from the same target-family model, and (ii) features from **SmolVLM**, a substantially smaller proxy.

Table 11 shows that the two variants perform similarly across all tested benchmarks, with differences that are small relative to the benchmark scale. This indicates that the resolution-selection decision depends primarily on coarse visual-textual cues that are preserved across different VLMs, rather than requiring tight alignment between proxy and target feature spaces.

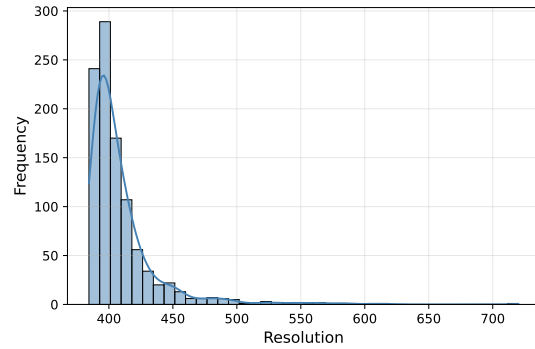
A.6 Predicted resolution distributions

Figure 7 shows the distribution of continuous resolutions predicted by CARES across different benchmarks. The histograms highlight that the selector adapts its behavior to the underlying task:

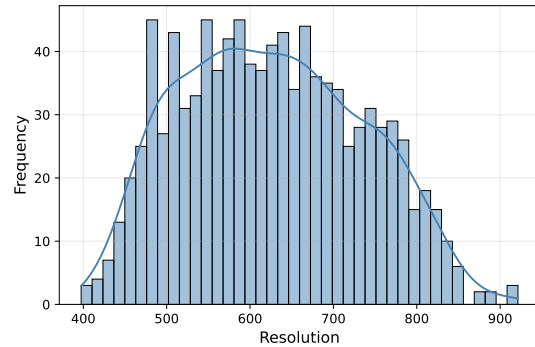
Task	Qwen features	SmolVLM features
Ai2D	0.7830 ± 0.0074	0.7824 ± 0.0074
ChartQA	0.8164 ± 0.0077	0.8080 ± 0.0079
DocVQA	0.8814 ± 0.0044	0.8640 ± 0.0040
OCRBench	0.7500 ± 0.0001	0.7600 ± 0.0003

Table 11: Downstream performance of CARES when using same-family versus proxy-family features. The small differences suggest that CARES is robust to moderate proxy–target representation mismatch.

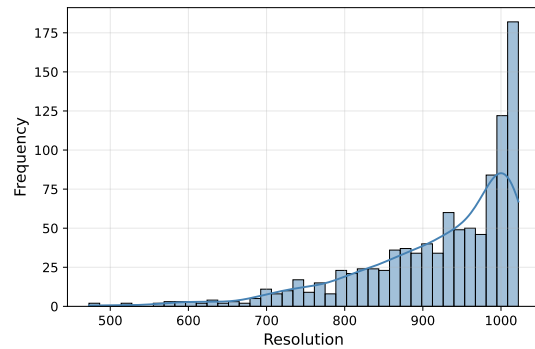
Ai2D is dominated by lower-resolution predictions, suggesting that many diagram-understanding questions require only coarse visual information; SeedBench-2 shifts toward higher resolutions, reflecting the need for finer-grained visual recognition; and DocVQA and OCRBench exhibit broader distributions, indicating a mixture of easy and detail-sensitive examples. This behavior is consistent with the intended design of CARES, which escalates resolution only when the image-query pair appears to demand additional visual detail.



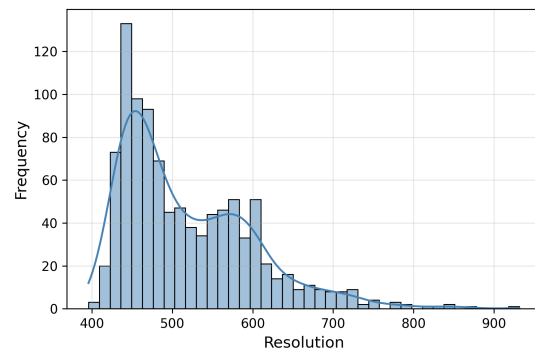
(a) Ai2D.



(b) DocVQA.



(c) SeedBench-2.



(d) OCRBench.

Figure 7: **Histograms of the predicted continuous resolutions \tilde{r} by CARES.** CARES routes many Ai2D examples to lower resolutions, while SeedBench-2 shifts toward higher resolutions. DocVQA and OCRBench show broader distributions, reflecting their mixture of coarse and fine-grained queries, including dense text and complex layouts.