

Downscaling Intelligence: Exploring Perception and Reasoning Bottlenecks in Small Multimodal Models

Supplementary Material

A1. Additional LLM Downscaling Results and Details

Per-tasks results. We present plots showing the performance dropoff from LLM downscaling across all evaluated tasks in Figure A1. As described in the main text, most tasks exhibit minimal performance decline when downscaling the LLM, except for a handful of vision-centric tasks that exhibit substantially larger drops (e.g., Grounding, NIGHTS, PieAPP).

Full decoupled results. We plot the performance dropoff from LLM downscaling of the perception and reasoning modules in Figure A2. We find that LLM downscaling of either module leads to performance degradation across a wide range of tasks. Notably, downscaling the perception module has a large effect on both tasks assessing perception (e.g., OCR-VQA, Fine-grained Perception) and reasoning (e.g., Logical Reasoning). One exception is Math, where LLM downscaling of the perception module has little impact. We expect this is because mathematical ability is limited primarily by the downstream process of operating on visual information (reasoning) rather than by the foundational perception ability.

soning modules in Figure A2. We find that LLM downscaling of either module leads to performance degradation across a wide range of tasks. Notably, downscaling the perception module has a large effect on both tasks assessing perception (e.g., OCR-VQA, Fine-grained Perception) and reasoning (e.g., Logical Reasoning). One exception is Math, where LLM downscaling of the perception module has little impact. We expect this is because mathematical ability is limited primarily by the downstream process of operating on visual information (reasoning) rather than by the foundational perception ability.

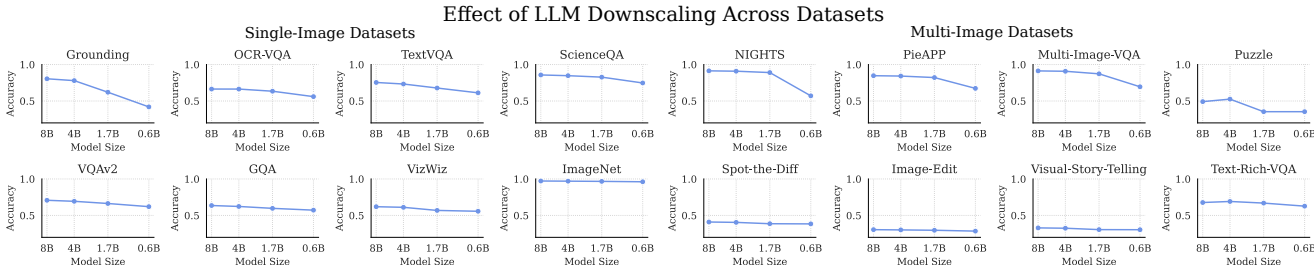


Figure A1. Performance dropoff from downscaling LLM across all datasets.

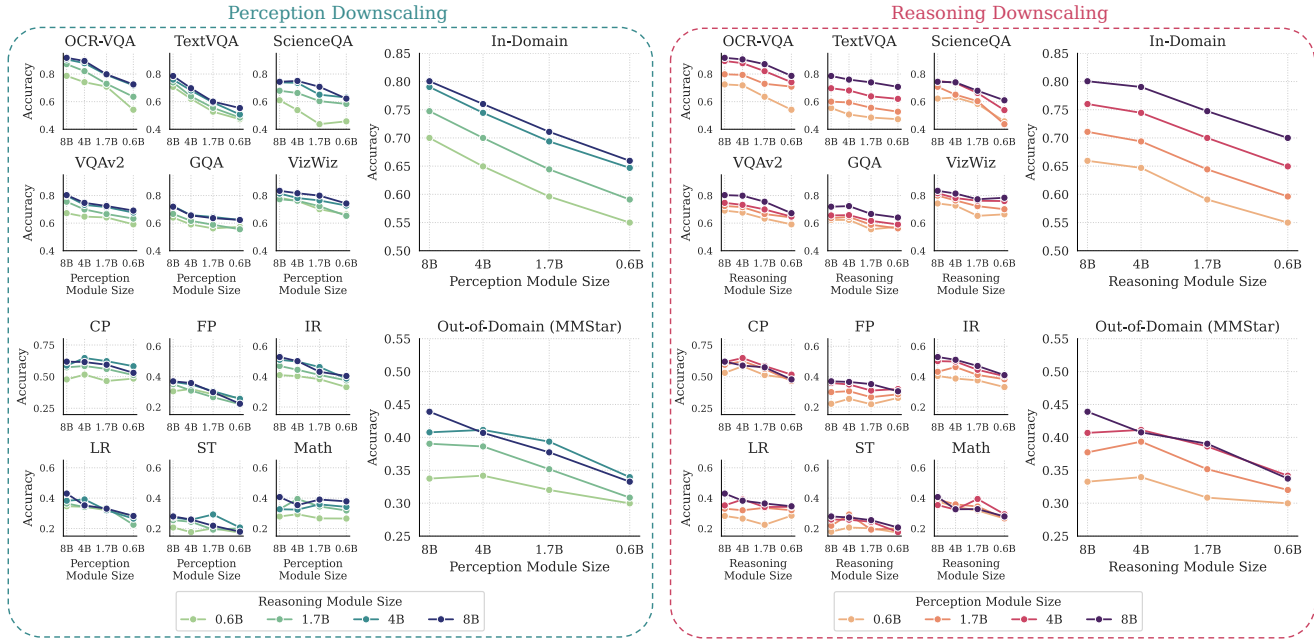


Figure A2. Full decoupled results. CP=Coarse Perception, FP=Fine-grained Perception, IR=Instance Reasoning, LR=Logical Reasoning, ST=Science & Technology.

LLaVA-OneVision as the perception module. While our main analysis used models trained from scratch for a controlled study, here we also experiment with using LLaVA-OneVision (\in 0.5B, 7B) as the perception module in the decoupled framework. We first present decoupled results using the same reasoning module as in our experiments (Qwen3). As shown in Figure A3, this configuration produces results that are largely consistent with those obtained using our controlled model as the perception module, where LLM downscaling of either module hinders performance. We do observe, however, that downscaling the perception module has a smaller effect than in our controlled study for

in-domain data. This is likely because LLaVA-OneVision includes extensive training on captioning, which we demonstrate alleviates the perception bottleneck.

Additionally, we experiment with using Qwen2 as the reasoning model in this setup (the LLM used in LLaVA-OneVision). As shown in Figure A4, relative to the earlier results using Qwen3 as the reasoning module, we see a larger impact from downscaling the reasoning module on the in-domain tasks, and overall performance is weaker than when using Qwen3. This outcome is not surprising, as Qwen3 has demonstrated stronger performance than Qwen2 on textual tasks, particularly for smaller model variants.

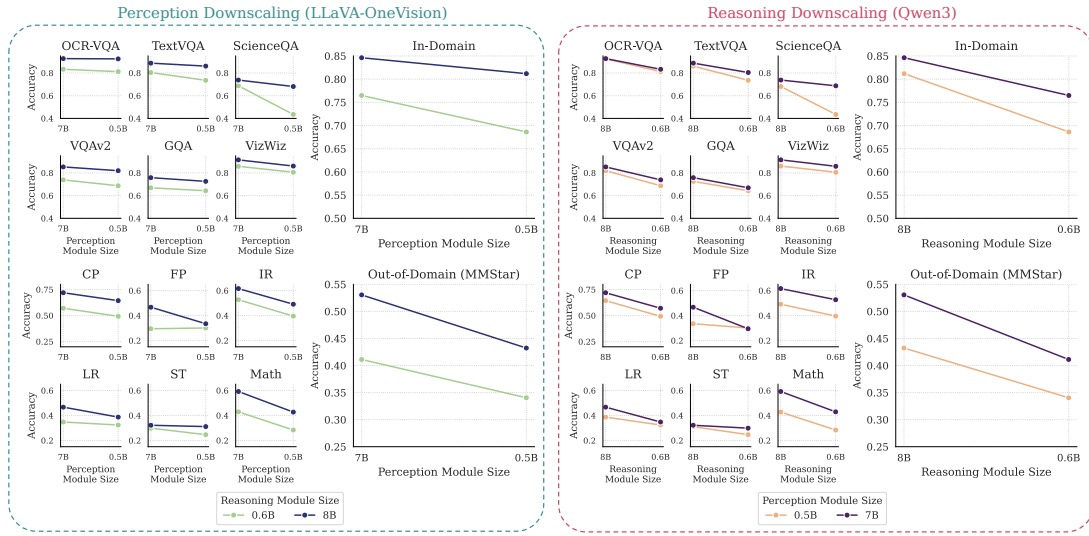


Figure A3. Decoupled analysis using LLaVA-OneVision as the perception module and Qwen3 as the reasoning module. CP=Coarse Perception, FP=Fine-grained Perception, IR=Instance Reasoning, LR=Logical Reasoning, ST=Science & Technology.

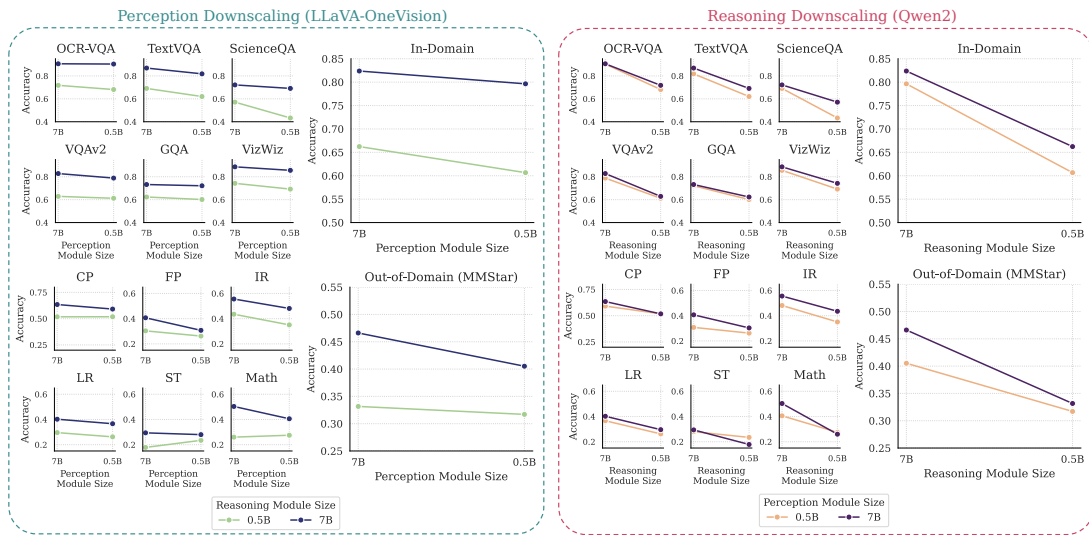


Figure A4. Decoupled analysis using LLaVA-OneVision as the perception module and Qwen2 as the reasoning module. CP=Coarse Perception, FP=Fine-grained Perception, IR=Instance Reasoning, LR=Logical Reasoning, ST=Science & Technology.

Prompts for decoupled analysis. We list prompt templates for our decoupled perception and reasoning analysis in Figure A5. These prompts follow the Prism framework, except that the initial instruction for obtaining question-specific information is run offline using the same model throughout, ensuring that deriving question-specific instructions does not influence our analysis of perception and reasoning downscaling. Thus, the question-specific information inserted into the perception module prompt is consistent across all model setups.

Perception bottleneck across model families. We additionally analyze Gemma3 (12B → 4B) and InternVL2.5 (8B → 2B), which offer multiple LLM sizes with a fixed vision encoder, enabling a controlled analysis. As shown in

Table A1, evaluating on MMStar using the same procedure as §3.3, we observe consistent trends of decreasing performance when downscaling the perception module’s LLM.

P Module	R Module	P LLM Size	Avg. Acc. Drop
Ours (§3)	Qwen3	8B → 1.7B	3.30
InternVL2.5	InternLM2	8B → 2B	4.33
Gemma3	Gemma3	12B → 4B	5.26

Table A1. Effect of downscaling the perception module on MM-Star accuracy. For each model family, we vary the perception module’s LM size while holding the reasoning module fixed, and report the accuracy drop averaged across both reasoning module sizes. P=Perception, R=Reasoning.

Question-specific Instruction Prompt

Your task is to give a concise instruction about what basic elements are needed to be described based on the given question. Ensure that your instructions do not cover the raw question, options or thought process of answering the question.

Examples:

Question: In which period the number of full time employees is the maximum?

Contents to observe: the number of full time employees

Question: What is the value of the smallest bar?

Contents to observe: the heights of all bars and their values

Question: What is the main subject of the image?

Contents to observe: the central theme or object

Question: What is the position of the catcher relative to the home plate?

Contents to observe: the spatial arrangement of the objects

Question: What is the expected ratio of offspring with white spots to offspring with solid coloring? Choose the most likely ratio.

Contents to observe: the genetic information

Now, perform the task, and format your answer as "Contents to observe:"

Question: <question>

Perception Module Prompt

Describe the fine-grained content of the image, including scenes, objects, relationships, instance location, and any text present. Especially, pay attention to <question-specific info>

Reasoning Module Prompt (w/o thinking)

You are an excellent text-based reasoning expert. You are required to answer the question based on the detailed description of the image.

Description: <description>

Question: <question>

Answer directly with the option’s letter in the format of "Answer:". Do not add anything other than the letter answer after "Answer:".

Reasoning Module Prompt (w/ thinking)

You are an excellent text-based reasoning expert. You are required to answer the question based on the detailed description of the image.

Description: <description>

Question: <question>

Please reason step by step, and give the final answer on the last line by itself in the format of "Answer:". Do not add anything other than the letter answer after "Answer:".

Figure A5. Prompt templates for the decoupled perception / reasoning analysis.

A2. Additional Visual Extraction Tuning Details

Visual extraction tuning data generation pipeline. We present the prompt templates for our pipeline generating visual extraction tuning data in Figure A6. In the first stage, each question–answer pair in a visual instruction example is converted into a declarative statement. These statement(s) are then inserted into a prompt that instructs the model to describe fine-grained visual details, with explicit emphasis on the information contained in the visual instruction example.

Visual extraction tuning data examples. We provide examples produced by our visual extraction tuning data generation pipeline in Figure A7. We see that the generated examples preserve the original information required to answer the questions from the visual instruction while also offering a more detailed and comprehensive depiction of the specific image content.

Question-answer Pair to Declarative Statement Prompt

Your task is to convert each question–answer pair about an image into a concise, fully self-contained declarative statement. The resulting statements should be understandable on their own, without requiring the reader to refer to the original question.

```
% for each QA pair in conv include:
```

```
    Question: <question>
```

```
    Answer: <answer>
```

```
% if len(conv) > 1 include:
```

```
    As there are <len(conv)> questions, you should respond with <len(conv)> statements. Include each statement on its own line
```

```
    Declarative Statement(s):
```

Visual Extraction Response Generation Prompt

Your task is to describe the fine-grained content of the image, including scenes, objects, relationships, instance location, and any text present.

As part of your description, you should incorporate the following information about the image.

```
<declarative statements>
```

```
Description:
```

Figure A6. Prompts used for visual extraction tuning data generation pipeline.

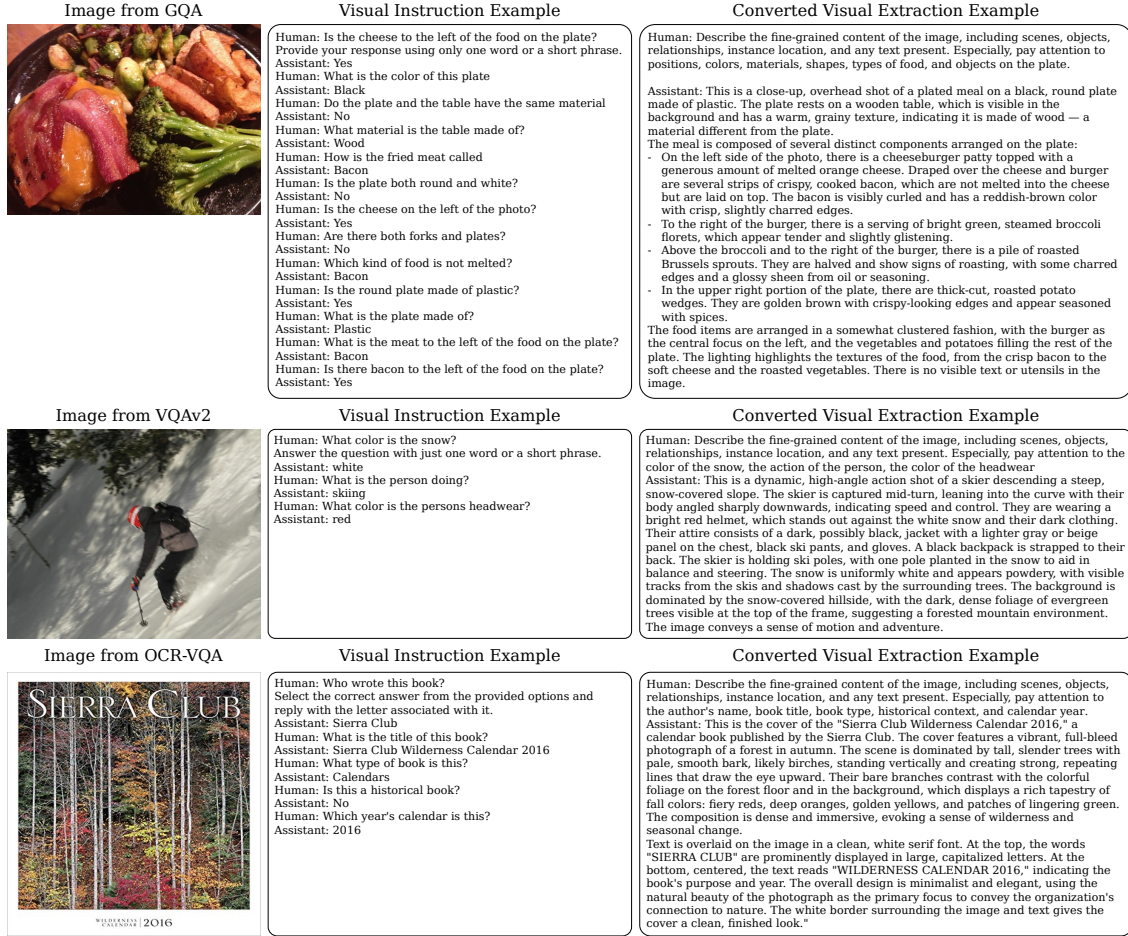


Figure A7. Visual extraction tuning data examples.

A3. Additional Step-by-step Reasoning Details and Results

NoWait Setup. To reduce overthinking of Qwen3 with thinking mode enabled, we use a logits processor that suppresses self-reflection tokens. Namely, we mask the logits of any token that contains one of the following keywords: {wait, alternatively, hmm, but, however, alternative, another, check, double-check, oh, maybe, verify, other,

again, now, ah, anyway, anyhow}, while manually excluding words that only contain a keyword as a substring but are not reflexive (e.g., waiter).

Full results. We present results from performing step-by-step visual reasoning across all tasks in Figure A8. Expectedly, we find that Math heavily benefits from CoT reasoning (consistent with findings in text-only Math tasks).

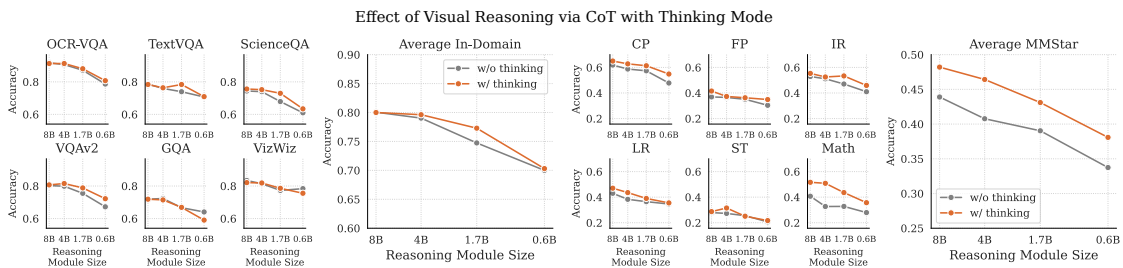


Figure A8. Full results showing impact of step-by-step reasoning on in-domain and out-of-domain (MMStar) performance.

A4. Additional EXTRACT+THINK Analyses

Inference latency analysis. In Figure A9, we characterize the tradeoff between our approach’s parameter/data efficiency and its inference latency. We plot latency across end-to-end baselines, PrismCaptioner (1.8/70B), and our EXTRACT+THINK method (1.7/4B w/ and w/o CoT). While end-to-end approaches incur lower latency from generating far fewer tokens, our method achieves 27.4% improvement on MMStar and 8.8% on VMCBench tasks. Additionally, vLLM greatly increases generation throughput for longer output sequences, narrowing the gap, particularly at a higher batch size.

Comparing against PrismCaptioner underscores the benefit of visual extraction tuning: even w/o CoT, our method improves performance by 9.80% on MMStar and 10.9% on VMCBench while using a smaller reasoning module and reducing latency. This indicates that increasing the token budget via a two-stage design alone does not drive performance; rather, it is the visual extraction tuning that drives gains within this framework. Applying test-time scaling through CoT further improves performance, particularly on MMStar.

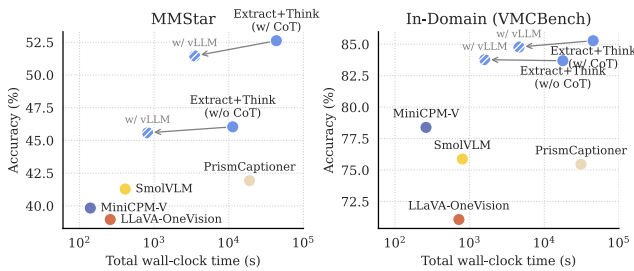
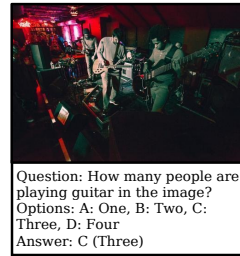


Figure A9. Inference latency on 1 L40s using Transformers (BS=1) or vLLM (BS=32), evaluated on MMStar and In-Domain (VMCBench).

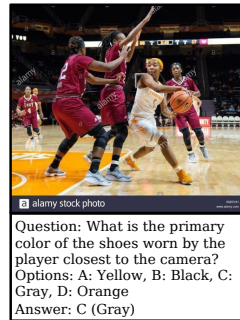
Information bottleneck in two-stage design. While the two-stage approach offers a modular framework for independently studying and improving perception and reasoning and has been shown to enable strong performance, a drawback is that its design requires all relevant visual information for answering the question to be captured and adequately described in the initial stage. This stage acts as an intermediate representation: the reasoning module bases its final response on this description rather than the image itself. In Figure A10, we show examples where this leads to loss of critical visual information, ultimately leading the reasoning module astray. For instance, a musician may be described as playing an instrument without specifying which one, or players may be described without indicating their exact, relative spatial relationship to the camera, rendering certain questions unanswerable. Nonetheless, the proposed visual extraction tuning approach is specifically



Extraction output: “In the foreground, a guitarist is prominently featured, holding a black electric guitar and playing it with a focused expression... Behind him, another guitarist is visible, also playing an electric guitar, though slightly out of focus. A third musician, partially obscured, is visible in the background, also playing an instrument...”
Final prediction: B (Two)



Extraction output: “The player on the right, wearing a white pinstripe jersey with red sleeves and blue pants, is in a defensive stance, reaching out with her right hand to block or intercept the ball. Her jersey has the word “Raiders” in red script across the chest, and the number “4” is visible below it. She is wearing a brown leather catcher’s mitt on her left hand, which is holding a bright yellow softball...”
Final prediction: C (Yellow)



Extraction output: “The central focus is on a player in a white jersey with orange trim and the word “TENNESSEE” across the chest, who is dribbling a basketball with both hands. She is wearing bright orange sneakers and is in motion, leaning forward with her body low to the ground, suggesting she is maneuvering around defenders... Surrounding her are three opposing players in maroon jerseys... **The players are positioned close to the camera**, creating a sense of immediacy and intensity in the moment.”
Final prediction: D (Orange)

Figure A10. Examples showing failure modes of two-stage framework for visual question answering (using EXTRACT+THINK w/ 1.7B perception module).

designed to mitigate such issues by targeting the extraction of question-relevant visual details and has been shown to effectively improve performance across evaluated tasks.