# Cross-Modal Memory Compression for Efficient Multi-Agent Debate

Jing Wu [*]   Yue Sun [*]   Tianpei Xie   Suiyao Chen   Jingyuan Bao   Yaopengxiao Xu   Gaoyuan Du   Inseok Heo
Alexander Gutfraind   Xin Wang

## Abstract

Multi-agent debate can improve reasoning quality and reduce hallucinations, but it incurs rapidly growing context as debate rounds and agent count increase. Retaining full textual histories leads to token usage that can exceed context limits and often requires repeated summarization, adding overhead and compounding information loss. We introduce DebateOCR, a cross-modal compression framework that replaces long textual debate traces with compact image representations, which are then consumed through a dedicated vision encoder to condition subsequent rounds. This design compresses histories that commonly span tens to hundreds of thousands of tokens, cutting input tokens by more than 92% and yielding substantially lower compute cost and faster inference across multiple benchmarks. We further provide a theoretical perspective showing that diversity across agents supports recovery of omitted information: although any single compressed history may discard details, aggregating multiple agents' compressed views allows the collective representation to approach the information bottleneck with exponentially high probability.

## 1. Introduction

Multi-agent debate (MAD) has emerged as a powerful paradigm for enhancing large language model (LLM) performance across reasoning, factuality, and complex problem-solving tasks(Du et al., 2023; Khan et al., 2024; Liu et al., 2025; Wu et al., 2025b; Liang et al., 2023). By enabling multiple model instances to propose, critique, and refine responses through iterative discussion, MAD consistently outperforms single-agent approaches on mathematical reasoning (Cobbe et al., 2021), question answering (Hendrycks et al., 2020), and image captioning (Lin et al., 2014; Wu et al., 2025a; Lin et al., 2024), etc. The fundamental princi-
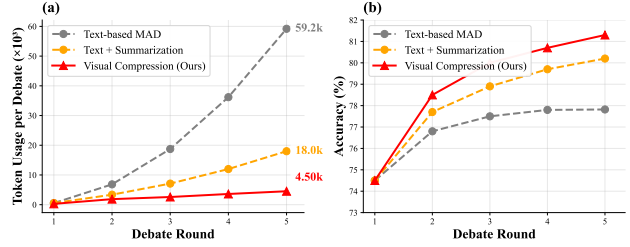
---
[*]Equal contribution .

*Figure 1.* Visual compression addresses the token inefficiency of multi-agent debate. We compare three paradigms across 5 debate rounds using Qwen2-VL on GSM8K with 3 agents. (a) **Token consumption:** Text-based MAD accumulates 59.2K tokens by Round 5 due to repeatedly storing debate history, while our visual compression reduces this by 92.2% to only 4.5K tokens. Text with summarization achieves 69.6% reduction to 18.0K tokens. (b) **Reasoning accuracy:** Despite lower computational cost, visual compression achieves the highest accuracy of 81.3%, outperforming text with summarization and text-based MAD. *Note:* At Round 5, text-based MAD (the gray curve) exceeds the context window limit. Therefore, the accuracy is measured with truncated debate history.

ple underlying MAD's success is that diverse perspectives and iterative refinement converge toward more reliable solutions, similar to human collaborative problem-solving.

However, the computational overhead of MAD scales rapidly with both the number of agents and debate rounds as shown in Figure 1. Each agent must maintain complete debate histories as textual context, with token consumption growing quadratically as histories are replicated across all agents. As a result, extended debates frequently exceed context window limits, and lengthy debate histories complicate final decision-making as judges must extract relevant information from increasingly verbose exchanges. Recent analysis reveals that context limitations and communication breakdowns account for significant performance degradation in multi-agent systems, with agents struggling to maintain comprehensive state across extended interactions (Cemri et al., 2025). Existing approaches address this through periodic summarization or truncation strategies (Chen et al., 2024a; Liu et al., 2025; Wu et al., 2025b), but these introduce additional computational cost and inference latency. To the best of our knowledge, no scalable solution exists for maintaining full debate context without prohibitive token overhead.

1

In this work, we propose a framework that applies visual compression to substantially reduce token consumption in multi-agent debate. Our key insight is that textual debate histories can be rendered as images and processed by specialized vision encoders through cross-modal operations, effectively converting text tokens into vision tokens at a fraction of the original cost. We design a compression-optimized vision encoder that maintains minimal activations under high-resolution inputs while achieving over 92% token reduction. This approach fundamentally addresses the scalability challenge by converting the quadratically growing textual context into compact visual representations shared across agents.

Beyond efficiency, our framework offers several advantages: (1) it seamlessly integrates with existing MAD algorithms without architectural modifications; (2) it eliminates the need for summarization strategies, preserving complete debate histories; (3) visual encoding maintains richer contextual information than text alone, as vision tokens naturally capture the structural relationships and logical flow of debates, leading to improved reasoning quality.

We evaluate our framework on three reasoning benchmarks: MATH (Hendrycks et al., 2020), GSM8K (Cobbe et al., 2021), and GPQA (Rein et al., 2024), using four vision-language models: Qwen2.5-VL-7B (Bai et al., 2025), Llama-3.2-11B-Vision (Meta, 2024), InternVL2-8B (Chen et al., 2024d), and Pixtral-12B (Agrawal et al., 2024). The cross-modal compression method achieves over 92% token reduction, while maintaining competitive accuracy. These gains require no modifications to the underlying debate algorithm or agent architecture.

To explain why compression preserves accuracy despite dramatic token reduction, we develop a theoretical analysis from the perspective of the Information Bottleneck (Kawaguchi et al., 2023), which characterizes the minimum information required for optimal decision-making. We show that, as the number of diverse agents increases, compressed histories converge to this bottleneck. The central mechanism is information recovery through diversity: although each agent's compressed history may discard some task-relevant information, different agents tend to preserve complementary aspects of the signal. When a majority of agents successfully retain the relevant information, aggregation recovers an essentially complete representation. At the same time, compression suppresses spurious artifacts, and independent artifacts introduced by different agents are further canceled through aggregation. Together, these effects explain why multi-agent systems can tolerate aggressive compression: collective redundancy offsets individual information loss, enabling accurate system-level decisions despite severe per-agent token reduction.

## 2. Related Work

### 2.1. Multi-Agent Debate for LLM Reasoning

Multi-agent debate (MAD) enhances LLM reasoning by enabling multiple agents to propose, critique, and refine responses through iterative discussion (Du et al., 2023; Liang et al., 2023; Chan et al., 2023). Prior work explores role assignment strategies (Wang et al., 2023), debate protocols that encourage error correction (Khan et al., 2024; Liu et al., 2025), and expert-guided collaboration through meta-programming and consistency mechanisms (Hong et al., 2023; Xiong et al., 2023; Pham et al., 2023). Recent advances include reflective multi-agent collaboration (Bo et al., 2024) and self-improvement through reinforcement learning (Chen et al., 2024b; Subramaniam et al., 2025). However, recent analysis reveals that context limitations and inter-agent misalignment cause significant failures, with agents struggling to maintain state across extended interactions (Cemri et al., 2025). This stems from a fundamental bottleneck: debate histories grow quadratically with agents and rounds, frequently exceeding context windows and requiring computationally expensive summarization (Chen et al., 2024a; Wu et al., 2025b).

### 2.2. Context Compression for LLMs

Context compression addresses long-context challenges through multiple approaches. Vision-based methods convert text into images processed by lightweight encoders, achieving substantial compression ratios (Wei et al., 2025; Xing et al., 2025). Text-based methods employ soft prompt compression (Ge et al., 2023; Mu et al., 2023; Chevalier et al., 2023) or selection-based token pruning using information entropy (Li et al., 2023; Jiang et al., 2023). Memory-augmented architectures extend context through external memory banks (Mohtashami & Jaggi, 2023; Tworkowski et al., 2023). However, vision-based compression for multi-agent systems remains unexplored. Our work differs by applying visual compression specifically to dynamic debate contexts rather than static documents, addressing the unique challenges of multi-agent communication efficiency.

### 2.3. Vision-Language Models

Vision-language models integrate visual and textual modalities for multimodal understanding (Zhang et al., 2024). Key architectures include Qwen2-VL (Wang et al., 2024), LLaVA-OneVision (Li et al., 2024), InternVL2 (Chen et al., 2024c), and DeepSeek-VL (Wu et al., 2024), built upon vision encoders like CLIP (Radford et al., 2021), ViT (Dosovitskiy, 2020), and SAM (Kirillov et al., 2023). Recent work explores parameter-efficient adaptation and token compression for improved efficiency (Danish et al., 2025). We leverage pretrained vision encoders to compress textual de-

bate contexts into compact visual representations, extending VLM capabilities to multi-agent communication.

# 3. Method

We propose DebateOCR, a framework that leverages visual compression to address the scalability challenges of multi-agent debate, reducing token consumption over 90%.

## 3.1. Preliminaries and Computational Challenges

**Preliminaries.** We formalize multi-agent debate as a sequential state-expansion process involving a set of $K$ agents, denoted by $\mathcal{A} = \{A_1, \ldots, A_K\}$. Each agent $A_i$ is instantiated as a LLM. Mathematically, we treat each agent as a policy $\pi_i$ that maps a query context to a textual response. Given a query $q \in \mathcal{Q}$ and a debate history $H_r = \{H_{i,r}\}_{i=1}^K$, the agent $i$ samples a solution $s_{i,r}$ at round $r$:

$$s_{i,r} \sim \pi_i(\cdot \mid q, H_r), \tag{1}$$

where $s_{i,r} \in \mathcal{S}$ represents the sequence of tokens generated by agent $i$ at round $r$, $\mathcal{S}$ denotes the output token space and $H_{i,r}$ denote the history generated at agent $i$ and round $r$ using the correspnding query.

The state of the system is defined by the accumulation of generated solutions. We define the history $H_r$ recursively. At the initial round ($r = 1$), the history is empty: $H_1 = \emptyset$. For any subsequent round $r > 1$, the history is updated by appending the set of all solutions from the previous round:

$$H_r = H_{r-1} \cup \{s_{1,r-1}, \ldots, s_{K,r-1}\}. \tag{2}$$

Here, $H_r$ serves as the global context observed by all $K$ agents to generate the next step of reasoning.

The debate concludes after a fixed horizon of $R$ rounds. The final output $y$ is derived by applying a consensus function $\phi$ over the final set of solutions:

$$y = \phi(s_{1,R}, \ldots, s_{K,R}), \tag{3}$$

where $\phi$ typically represents a majority vote, model-based judge, or weighted aggregation. The choice of $\phi$ does not affect our analysis of computational costs.

**The Scalability Challenge.** The primary bottleneck in this formulation is the quadratic growth of input context tokens. We quantify this cost by analyzing the total number of tokens processed across all agents and rounds.

Let $L$ denote the average length of a solution $s_{i,r}$ in tokens. The size of the history context at round $r$, denoted as $|H_r|$, scales linearly with the number of agents and rounds:

$$|H_r| = K \cdot (r - 1) \cdot L. \tag{4}$$

However, the computational burden is multiplicative. Since all $K$ agents must process this history independently at every round, the token consumption at round $r$ is:

$$C(r) = K \cdot |H_r| = K^2 \cdot (r - 1) \cdot L. \tag{5}$$

The cumulative token consumption across all $R$ rounds becomes:

$$C_{\text{total}} = \sum_{r=1}^{R} C(r) = K^2 L \sum_{r=1}^{R} (r - 1) = \frac{K^2 L R(R - 1)}{2}. \tag{6}$$

This implies a complexity of $\mathcal{O}(K^2 R^2 L)$ as both the number of agents and debate rounds scale. Figure 1(a) empirically demonstrates this quadratic scaling: a 5-round debate on GSM8K with 3 agents accumulates 59.5K tokens, exceeding typical context windows. This quadratic scaling severely limits the applicability of MAD to complex reasoning tasks requiring extended deliberation, as token budgets are quickly exhausted even with modest numbers of agents and rounds.

## 3.2. Framework Overview

To address the quadratic token scaling of Eq. 6, we propose a visual compression framework that converts textual debate histories into compact visual representations. The approach operates in two distinct phases: *an offline training phase* that learns to align visual encodings with the target MLLM's embedding space, and an *online inference phase* that compresses debate histories into a constant number of vision tokens.

**Two-Phase Pipeline.** Our framework leverages a vision encoding strategy that combines complementary feature representations from pre-trained SAM and CLIP encoders. A lightweight adapter network learns to project their joint features into the target MLLM's vision token space. This adapter is trained once per target MLLM and transfers seamlessly across different debate scenarios without requiring scenario-specific fine-tuning.

**Training phase**: We train the adapter to reconstruct text from rendered images using a diverse corpus spanning multiple reasoning domains. The training objective encourages the adapter to learn a compressed representation that preserves the semantic content necessary for accurate text reconstruction. Details are included in Section 3.3.

**Inference phase**: During a multi-agent debate, we render the textual history as a structured image, encode it through the trained vision pipeline into a fixed-size sequence of vision tokens, and inject these tokens into each agent's context in place of the original text. This compression reduces token consumption from $\mathcal{O}(K^2 R^2 L)$ to $\mathcal{O}(KRN)$, where $N$ is a constant significantly smaller than $KRL$, achieving 80-97%
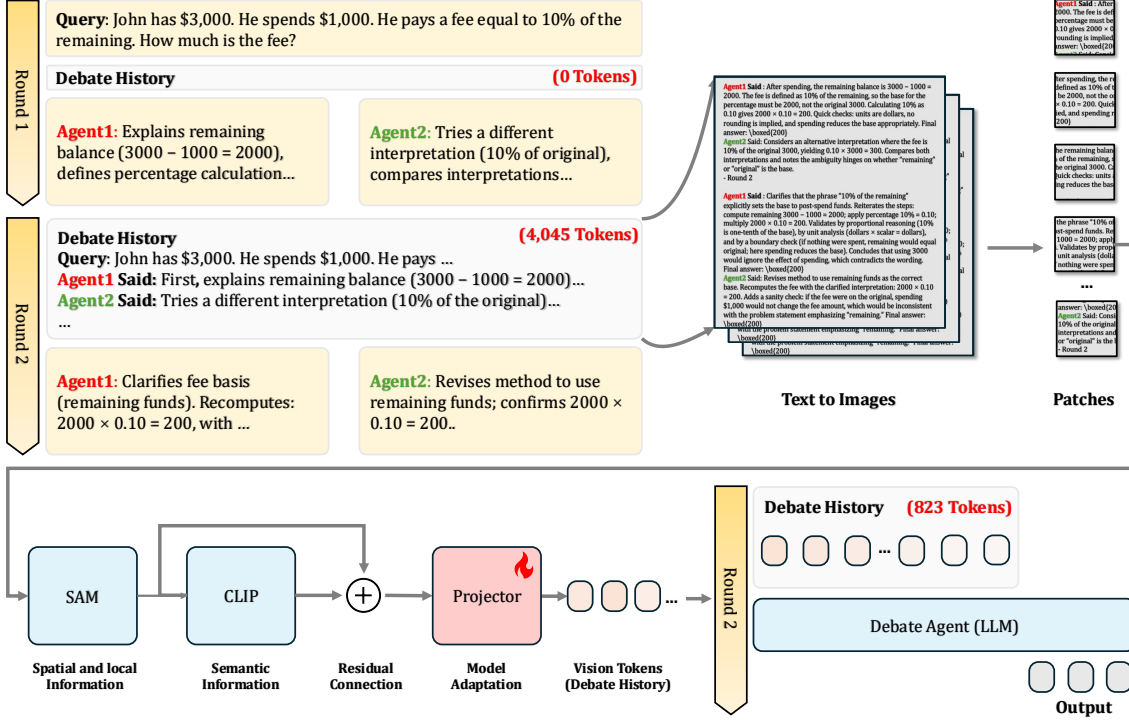
*Figure 2.* Visual compression framework for multi-agent debate. **Top:** Text-based debate accumulates token count across rounds. **Bottom:** Our approach converts debate history into visual representations. We first train a lightweight projector to adapt SAM-CLIP features to the vision embedding space of target MLLMs. During inference, debate history text is encoded as images via SAM, embedded through CLIP, and projected into small amount of vision tokens, achieving token reduction while preserving debate context for subsequent rounds.

token reduction as demonstrated in Figure 1. We provide detailed description at Section 3.4.

### 3.3. Training Phase

We train a lightweight adapter to align visual features from rendered debate histories with the target MLLM's embedding space. The training pipeline consists of frozen pretrained encoders that extract visual features, followed by a trainable projection network that maps these features into the MLLM's token space. Crucially, only the adapter parameters are updated during training, where SAM, CLIP, and the MLLM remain completely frozen, enabling efficient training while leveraging strong pretrained representations.

#### 3.3.1. ENCODER PIPELINE

The encoder pipeline processes rendered debate history images through a serial architecture combining SAM and CLIP encoders. Both SAM and CLIP remain frozen throughout training, following the design principles of DeepSeek-OCR (Wei et al., 2025).

**SAM Feature Extraction.** Given a rendered image $I_r$ at $1024 \times 1024$ resolution, SAM-base first processes it using window-based attention with patch size 16. The SAM backbone outputs spatial features $F_{\text{SAM}} \in \mathbb{R}^{B \times 768 \times 64 \times 64}$ that

capture fine-grained text layout and positioning information.

**Neck Module Architecture.** To align SAM features with CLIP's input requirements, we employ a Feature Pyramid Network-style neck module:

$$F_{\text{SAM}} \xrightarrow{\text{Stage 1}} F_1 \xrightarrow{\text{Stage 2}} F_2 \xrightarrow{\text{Stage 3}} F_{\text{neck}}$$
$$\mathbb{R}^{768 \times 64^2} \to \mathbb{R}^{256 \times 64^2} \to \mathbb{R}^{512 \times 32^2} \to \mathbb{R}^{1024 \times 16^2}$$

Stage 1 applies $\text{Conv}_{1 \times 1}$, LayerNorm2d, and $\text{Conv}_{3 \times 3}$; Stages 2-3 apply $\text{Conv}_{3 \times 3}$ with stride 2, reducing resolution by $4\times$ while expanding to 1024 channels.

**CLIP Feature Extraction.** The neck module output serves as input to CLIP-Large, bypassing CLIP's standard patch embedding layer. To transform spatial features into the sequence format expected by CLIP transformers, we flatten the spatial dimensions and transpose to form a token sequence: $[B, 1024, 16, 16] \xrightarrow{\text{flatten}} [B, 1024, 256] \xrightarrow{\text{transpose}} [B, 256, 1024]$, yielding 256 tokens per image. CLIP-Large's learned positional embeddings are interpolated via bicubic interpolation from their original grid size to match the $16 \times 16$ spatial layout, then added to the feature sequence. CLIP applies dense global attention to extract semantic representations, outputting $f_{\text{CLIP}} \in \mathbb{R}^{d_v}$ at each spatial position, where $d_v = 1024$.

4

**Feature Fusion with Residual Connection.** To preserve both spatial and semantic information, we fuse the neck-module features with the CLIP representation using a residual connection:

$$\mathbf{f} = \mathbf{f}_{\text{CLIP}} + \mathbf{f}_{\text{neck}}, \quad \mathbf{f} \in \mathbb{R}^{d_v}$$

where $d_v = 1024$. This residual fusion preserves fine grained spatial details from SAM while maintaining the semantic structure of CLIP embeddings, producing a compact representation for the subsequent adapter network.

### 3.3.2. ADAPTER NETWORK

The adapter $\mathcal{P}_\theta$ is a lightweight projection network that transforms the encoder features into the target MLLM's embedding space. Given fused features $\mathbf{f} \in \mathbb{R}^{d_f}$ at each spatial position (where $d_f = d_v$), the adapter applies a two-layer MLP:

$$\mathbf{z} = \mathcal{P}_\theta(\mathbf{f}) = \text{LayerNorm}(\mathbf{W}_2\sigma(\mathbf{W}_1\mathbf{f} + \mathbf{b}_1) + \mathbf{b}_2), \quad (7)$$

where $\mathbf{W}_1 \in \mathbb{R}^{d_h \times d_f}$ projects to hidden dimension $d_h$, $\mathbf{W}_2 \in \mathbb{R}^{d \times d_h}$ projects to the MLLM's embedding dimension $d$, and $\sigma$ denotes GELU activation. Layer normalization is applied to the output for training stability. This produces a sequence of vision tokens $\mathbf{Z} = \{\mathbf{z}_1, \ldots, \mathbf{z}_N\} \in \mathbb{R}^{N \times d}$, where $N$ depends on image resolution.

The adapter's role as a lightweight feature projection layer, combined with training on diverse reasoning domains, enables it to transfer across different debate scenarios without requiring task-specific fine-tuning. Once trained for a target MLLM, the same adapter works across all debate protocols and topics.

### 3.3.3. TRAINING PROCEDURE

**Training Objective.** We train the adapter through autoregressive text reconstruction. Given a text sample $T$, we render it as image $\mathcal{I}$, extract and project features to obtain vision tokens $\mathbf{Z}$, and minimize the cross-entropy loss between the MLLM's output and ground-truth text:

$$\mathcal{L}(\theta) = -\sum_{t=1}^{|T|} \log p_{\text{MLLM}}(T_t \mid \mathbf{Z}), \quad (8)$$

where $p_{\text{MLLM}}$ denotes the MLLM's generation probability conditioned on vision tokens. Only the adapter parameters $\theta = \{\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}_1, \mathbf{b}_2\}$ are updated, while SAM, CLIP, and the MLLM remain frozen, largely reducing computational cost while leveraging pretrained representations.

### 3.4. Inference Phase

During multi-agent debate, we apply the trained framework to compress debate histories at each round $r$. The inference pipeline consists of three stages: rendering the textual

history as a structured image, encoding the image through the frozen SAM-CLIP-adapter pipeline, and injecting the resulting vision tokens into each agent's context window.

### 3.4.1. TEXT-TO-IMAGE RENDERING

At round $r$, the debate history $H_r = \{s_{i,t}\}_{i=1,t=1}^{K,r-1}$ contains all previous agent responses. We render $H_r$ as a structured image $\mathcal{I}_r$ that preserves agent identities and temporal ordering through spatial layout. Each agent's response is rendered with identifying markers (e.g., "Agent 1:", "Agent 2:") and rounds are visually separated. This spatial organization enables vision encoders to capture both response content and debate structure through visual parsing, which linear text concatenation cannot preserve.

The image resolution and layout balance text legibility for the encoders against the resulting number of vision tokens. Implementation details are provided in Section 5.

### 3.4.2. VISION ENCODING

The rendered image $\mathcal{I}_r$ is processed through the trained encoding pipeline: SAM extracts spatial features, CLIP builds semantic representations via the residual connection, and the adapter projects the combined features into vision tokens $\mathbf{Z}_r = \{\mathbf{z}_1, \ldots, \mathbf{z}_N\} \in \mathbb{R}^{N \times d}$.

The number of vision tokens $N$ remains constant for each round of debate history, enabling fixed-size compression of arbitrarily long debates. This reduces token consumption from $\mathcal{O}(K^2 R^2 L)$ to $\mathcal{O}(KRN)$.

### 3.4.3. CONTEXT INJECTION

The vision tokens $\mathbf{Z}_r$ replace the textual debate history in each agent's context window. At round $r$, agent $i$ generates its response conditioned on the query and compressed history:

$$s_{i,r} \sim \pi_i(\cdot \mid q, \mathbf{Z}_r). \quad (9)$$

The MLLM processes vision tokens identically to text tokens through its transformer layers, requiring no modification to the MLLM architecture or debate protocol.

## 4. Theoretical Analysis

To understand why compression maintains accuracy despite dramatic token reduction, we provide a theoretical analysis from Information Bottleneck (Kawaguchi et al., 2023). Our key result shows that multi-agent aggregation enables compressed histories to approach the information bottleneck, the optimal trade-off between preserving answer-relevant information and removing spurious artifacts.

Consider a multi-agent debate system with $K$ agents deliberating on query $q$ with ground truth $y_q$. Each agent $i$

generates a debate history $H_{i,r}$ through round $r$. Let $\mathcal{H}$ denote the collection of histories and $f : \mathcal{H}^K \to \mathcal{Y}$ the aggregation function (e.g., majority voting) that produces the final answer.

Debate histories contain two types of information: (1) **answer-relevant information** $I(H; Y_q)$ necessary for correct decisions, and (2) **artifacts** $V_i$ representing agent-specific styles, redundancies, and tangential explorations. Let $\mathcal{C} : \mathcal{H} \to \mathcal{Z}$ denote a compression function of the debating history with $|\mathcal{Z}| \ll |\mathcal{H}|$.

**Definition 4.1** (Information Bottleneck). The **information bottleneck** is the minimum mutual information required for optimal decisions:

$$I_{\text{bottleneck}} = \min_{\mathcal{H}}\{I(\mathcal{H}; Y_q) : \mathbb{E}[\ell(f(\mathcal{H}), y_q)] = \ell_{\min}\}$$

where $\ell_{\min} = \inf_{\mathcal{H}} \mathbb{E}[\ell(f(\mathcal{H}), y_q)]$ is the minimum achievable expected loss.

Distance to the bottleneck is measured by:

$$D(\mathcal{H}) = |I(f(\mathcal{H}); Y_q) - I_{\text{bottleneck}}| + I(f(\mathcal{H}); V) \quad (10)$$

where $I(f(\mathcal{H}); V) = \sum_{i=1}^{K} I(f(\mathcal{H}); V_i)$ captures total artifact content.

We establish two results: (1) compression improves decision quality by removing artifacts, and (2) multi-agent aggregation naturally converges to the information bottleneck.

**Theorem 4.2.** *Suppose: (1) Compression preserves information with probability $p > 0.5$: $\mathbb{P}(I(\mathcal{C}(H_i); Y_q) \geq I(H_i; Y_q) - \varepsilon) \geq p$; (2) Agents have diverse styles with conditionally independent artifacts $V_i \perp V_j \mid q$; (3) Compression reduces artifacts: $I(\mathcal{C}(H_i); V_i) \leq \gamma I(H_i; V_i)$ for $\gamma \in (0, 1)$. Then compressed histories approach the bottleneck with exponentially high probability:*

$$\mathbb{P}(|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{bottleneck}| \leq \varepsilon) \geq 1 - e^{-2K(p-\frac{1}{2})^2} \quad (11)$$

*while artifacts vanish: $I(f(\mathcal{C}(\mathcal{H})); V) = O(\gamma K)$. Consequently, $D(\mathcal{C}(\mathcal{H})) < D(\mathcal{H})$ when artifacts are substantial.*

*Proof sketch.* The key insight is that diverse agents cover complementary aspects. With probability $p > 0.5$, each agent preserves its information. By Hoeffding's inequality, the majority preserve information with probability $\geq 1 - e^{-2K(p-1/2)^2}$. Since agents cover different aspects (diversity), their union approaches complete coverage: $I(f(\mathcal{C}(\mathcal{H})); Y_q) \geq I_{\text{bottleneck}} - \varepsilon$. Meanwhile, artifacts are conditionally independent and reduced by factor $\gamma$, yielding $I(f(\mathcal{C}(\mathcal{H})); V) \leq \gamma \sum_i I(H_i; V_i)$. For complete debate histories, $I(f(\mathcal{H}); Y_q) \approx I_{\text{bottleneck}}$ but $I(f(\mathcal{H}); V) \geq \bar{I}_V := \frac{1}{K} \sum_i I(H_i; V_i)$. Thus $D(\mathcal{C}(\mathcal{H})) \leq \varepsilon + \gamma K \bar{I}_V < K \bar{I}_V \lesssim D(\mathcal{H})$ when $(1-\gamma)K\bar{I}_V > \varepsilon$. Full proof in Appendix A. $\square$

**Corollary 4.3** (Sample Complexity). *To achieve $|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{bottleneck}| \leq \varepsilon$ with confidence $1 - \delta$, it suffices to use $K \geq \frac{\ln(1/\delta)}{2(p-1/2)^2}$ agents.*

Theorem 4.2 reveals a fundamental synergy between compression and aggregation. Although each individual agent incurs information loss under compression ($\varepsilon > 0$), collective aggregation recovers the missing information through agent diversity. Different agents preserve complementary aspects of signals: what one discards, another retains. This phenomenon formalizes the wisdom of crowds in information-theoretic terms: while individual representations are imperfect, their aggregation yields a near-optimal collective representation. Crucially, the exponential concentration guarantee shows that this effect improves steadily with the number of agents $K$, rather than relying solely on an asymptotic regime. As $K$ increases, the aggregated representation concentrates increasingly tightly around the information bottleneck. This explains why aggressive compression (e.g., substantial token reduction) can succeed in multi-agent systems while potentially failing in single-agent settings.

# 5. Experiments

We evaluate DebateOCR on three mathematical reasoning benchmarks. Additional details for training, Adapter design, Rendering and Debate Configuration will be available at Appendix B

## 5.1. Experimental Setup

### 5.1.1. TASKS AND DATASETS

We evaluate on three standard mathematical reasoning benchmarks:

**GSM8K** (Cobbe et al., 2021) contains 8,500 grade school math word problems requiring arithmetic reasoning. We use the standard test set of 1,319 problems. Problems typically require 2-8 reasoning steps and involve topics such as percentages, ratios, and basic algebra.

**MATH** (Hendrycks et al., 2021) comprises 12,500 competition-level mathematics problems spanning algebra, geometry, number theory, and calculus. We evaluate on the 5,000-problem test set.

**GPQA** (Rein et al., 2024) is a graduate-level science question-answering benchmark consisting of 448 multiple-choice questions written by domain experts in biology, physics, and chemistry.

### 5.1.2. MODELS AND BASELINES

We evaluate four open-source multimodal large language models with diverse architectures:

*Table 1.* Performance comparison across models and datasets after 5 debate rounds with 3 agents. We report accuracy (%), average number of input tokens per sample (in thousands), and average inference time per sample (seconds). ↑: higher is better, ↓: lower is better. Best results per model are **bolded**.

| Model | Method | GSM8K | | | MATH | | | GPQA | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Accuracy ↑ (%) | # Tokens ↓ (K) | Time ↓ (s) | Accuracy ↑ (%) | # Tokens ↓ (K) | Time ↓ (s) | Accuracy ↑ (%) | # Tokens ↓ (K) | Time ↓ (s) |
| InternVL-8B | T-MAD | 67.5 | 59.2 | 145.7 | 42.5 | 66.3 | 160.3 | 22.1 | 69.9 | 167.6 |
| | TS-MAD | 68.5 | 18.0 | 154.2 | 43.1 | 19.8 | 169.6 | **23.3** | 20.7 | 177.3 |
| | **DebateOCR (Ours)** | **70.5** | **4.5** | **64.7** | **44.2** | **4.7** | **71.2** | 22.8 | **4.9** | **74.4** |
| Qwen2.5-VL-7B | T-MAD | 77.6 | 59.2 | 131.1 | 44.8 | 72.0 | 144.2 | 32.6 | 75.9 | 150.8 |
| | TS-MAD | 80.2 | 19.6 | 138.8 | 45.1 | 21.5 | 152.7 | 33.8 | 22.5 | 159.6 |
| | **DebateOCR (Ours)** | **81.3** | **4.5** | **58.2** | **46.8** | **4.7** | **64.0** | **34.8** | **4.9** | **66.9** |
| Llama-3.2-11B | T-MAD | 78.2 | 66.8 | 218.6 | 45.2 | 74.8 | 240.5 | 33.0 | 78.9 | 251.4 |
| | TS-MAD | 78.8 | 20.3 | 231.3 | 46.7 | 22.3 | 254.4 | 33.9 | 23.4 | 265.8 |
| | **DebateOCR (Ours)** | **79.8** | **4.5** | **97.1** | **47.1** | **4.7** | **106.8** | **34.4** | **4.9** | **111.7** |
| Pixtral-12B | T-MAD | 79.4 | 62.6 | 233.1 | 49.2 | 70.1 | 256.4 | 36.4 | 73.9 | 268.1 |
| | TS-MAD | 81.5 | 19.1 | 246.7 | 51.8 | 20.9 | 271.4 | 37.9 | 21.9 | 283.7 |
| | **DebateOCR (Ours)** | **82.8** | **4.5** | **103.5** | **53.6** | **4.7** | **113.9** | **38.3** | **4.9** | **119.0** |

**Qwen2.5-VL-7B-Instruct** (Bai et al., 2025): A 7B-parameter vision-language model with adaptive resolution encoding, supporting variable aspect ratios and dynamic token allocation.

**Llama-3.2-11B-Vision** (Meta, 2024): An 11B-parameter model combining the Llama-3.2 language model with a vision encoder, optimized for multimodal instruction following.

**InternVL-8B** (Chen et al., 2024d): A 8B-parameter model using tile-based high-resolution image processing, designed for fine-grained visual understanding.

**Pixtral-12B** (Agrawal et al., 2024): A 12B-parameter multimodal model from Mistral AI featuring a 400M-parameter vision encoder with variable image resolution support, capable of processing arbitrary numbers of images at their natural resolution and aspect ratio.

We compare three approaches to managing debate history: **The Text-based Debate (T-MAD)**: Standard text-based multi-agent debate where the complete textual history $H_r$ is provided to each agent at every round. This represents the baseline MAD approach without any compression (Khan et al., 2024; Du et al., 2023).

**The Text-based Debate with Summarization (TS-MAD)**: After each round, the debate history is compressed via extractive and abstractive summarization to reduce token count while preserving key arguments. Agents receive the summarized history in subsequent rounds rather than the full text(Chen et al., 2024a; Liu et al., 2025).

### 5.1.3. EVALUATION METRICS

We evaluate methods along three dimensions:

**Accuracy**: We report exact match accuracy on each benchmark. For GSM8K, GPQA and MATH.

**Token Consumption**: We measure the total number of tokens processed across all agents and rounds. For text-based methods, this includes all input tokens, including query and debate history, at each round. For the proposed method, this includes the query tokens plus vision tokens from the compressed history.

**Inference Time**: We measure wall-clock time from initial query to final consensus, including all model forward passes across agents and rounds, on a single NVIDIA A100 GPU.

### 5.2. Main Results

Table 1 presents our main results across four vision-language models on three reasoning datasets. We compare our visual compression approach against two baselines: pure text-based debate and debate with extractive summarization. All methods use 5 debate rounds with 3 agents. Our visual compression achieves the best or competitive accuracy in most settings while dramatically reducing token consumption and inference time.

**Accuracy.** DebateOCR achieves the best accuracy across most experimental settings. On InternVL 8B, it achieves 70.5% on GSM8K and 44.2% on MATH, outperforming both baselines. On Qwen2.5 VL 7B, it reaches 81.3% on GSM8K, 46.8% on MATH, and 34.8% on GPQA, consistently surpassing text-based debate and summarization. Similar improvements are observed for Llama 3.2 11B and Pixtral 12B, where DebateOCR achieves 79.8% and 82.8% on GSM8K, respectively. The accuracy gains arise from improved preservation of spatial structure and formatting in rendered images, which helps models track multi-agent exchanges more effectively than long token sequences. In one

*Table 2.* Ablation study on image resolution on MATH with 5 rounds and 3 agents using Qwen2.5-VL-7B.

| Resolution | # Tokens | Acc (%) | Compress |
|---|---|---|---|
| $224 \times 224$ | 16 | 71.2 | **19.2×** |
| $336 \times 336$ | 36 | 73.1 | 18.9× |
| $448 \times 448$ | 49 | 75.5 | 18.6× |
| $512 \times 512$ | 64 | 76.2 | 18.3× |
| $1024 \times 1024$ | 256 | 76.3 | 15.3× |
| $1536 \times 1536$ | 576 | 76.3 | 12.0× |
| $2048 \times 2048$ | 1024 | 76.6 | 9.2× |

*Table 3.* Comparison of vision encoders on MATH with Qwen2.5-VL-7B using 5 rounds and 3 agents.

| Metric | DebateOCR | QwenVL2.5 |
|---|---|---|
| Resolution | $1024 \times 1024$ (fixed) | $1036 \times 1036$ (dynamic) |
| Vision Tokens (per image) | 256 | 1,369 |
| Total Vision Tokens | 3.8K | 20.5K |
| Token Reduction | **5.4×** | — |
| Inference Time (per turn) | 3.0s | 4.2s |
| Total Inference Time | 64.0s | 89.3s |
| Accuracy (%) | 46.8 | 46.9 |

case—GPQA with InternVL 8B—summarization slightly outperforms our method with 23.3% compared to 22.8%, likely because extractive summaries better retain domain-specific terminology for graduate-level questions.

**Token Efficiency.** DebateOCR achieves substantial reductions compared to both baselines. On InternVL-8B with GSM8K, our approach uses only 4.5K tokens, versus 59.2K for text-based debate and 18.0K for summarization, corresponding to 92.4% and 75.0% reductions, respectively. Token usage remains nearly constant across datasets, increasing only slightly with question length. In contrast, text-based methods exhibit significant token growth with debate rounds, consuming 59.2K–78.9K tokens for text-based debate and 18.0K–23.4K tokens for summarization, depending on dataset complexity and model.

**Inference Speed.** DebateOCR delivers substantial inference speedups over both baselines. On InternVL-8B with GSM8K, inference completes in 64.7 s, compared to 145.7 s for text-based debate and 154.2 s for summarization, yielding 2.25× and 2.38× speedups, respectively. Notably, the summarization baseline is consistently slower than pure text despite using fewer tokens, as it must generate summaries after each round, whereas our rendering incurs a fixed computational cost. These speedups are consistent across models and datasets, with larger models such as Llama-3.2-11B and Pixtral-12B exhibiting greater absolute time savings due to higher per-token costs. For example, on Llama-3.2-11B with MATH, DebateOCR reduces inference time from 240.5 s to 106.8 s, achieving a 2.25× speedup.

**Discussion for Accuracy Gains.** Multi-round text-based debates could accumulate substantial noise, including redundant arguments, overthinking patterns where agents explore but abandon incorrect reasoning paths, and agent-specific stylistic variations that inflate context without improving decision quality. These artifacts increase as $I(H_r; D_i)$ grows quadratically with debate rounds, introducing variability in the aggregation process. Visual compression may substantially reduce this artifact information to $I(C(H_r); D_i) \ll I(H_r; D_i)$ by rendering debate histories in standardized spatial layouts that naturally suppress stylistic variations and extract core logical structure. This explains why DebateOCR achieves higher accuracy despite dramatically reduced token

counts. Additional proofs and discussion are available at Appendix A.

### 5.3. Ablation Studies

**Image Resolution and Token Budget.** We conduct ablation studies on Qwen2.5-VL-7B with the MATH dataset using 3 agents and 5 rounds of debate in Table 2. Higher resolutions produce more vision tokens, enabling better text legibility but increasing computational cost. We observe that 1024×1024 resolution achieves the best balance, producing approximately 256 vision tokens with competitive accuracy at 76.3%. Lower resolutions such as 224×224 reduce vision tokens to 16 but suffer accuracy degradation due to insufficient text clarity.

**Comparison with Native Vision Encoder.** Table 3 compares our approach with Qwen2-VL's original vision encoder. DebateOCR achieves 5.4× token reduction with 3.8K versus 20.5K vision tokens over the full debate, while maintaining comparable accuracy. This efficiency translates to faster inference at 64.0s versus 89.3s, demonstrating that compact visual representations preserve essential reasoning information with lower computational cost.

## 6. Conclusion

We introduced DebateOCR, a token compression framework for multi-agent debate that addresses the quadratic token growth problem by rendering debate histories as images. The method achieves over 92% token reduction compared to text-based debate while maintaining competitive accuracy across mathematical and scientific reasoning tasks. The approach generalizes effectively across diverse MLLMs and scales linearly with debate rounds and agent count. We theoretically explain why compression preserves essential reasoning information while filtering debate artifacts. This work demonstrates that visual representations offer a practical and efficient alternative to text-based communication in multi-agent systems, enabling scalable debates without modifications to underlying algorithms or architectures.

# References

Agrawal, P., Antoniak, S., Hanna, E. B., Bout, B., Chaplot, D., Chudnovsky, J., Costa, D., De Monicault, B., Garg, S., Gervet, T., et al. Pixtral 12b. *arXiv preprint arXiv:2410.07073*, 2024.

Amini, A., Gabriel, S., Lin, S., Koncel-Kedziorski, R., Choi, Y., and Hajishirzi, H. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers)*, pp. 2357–2367, 2019.

Bai, S., Chen, K., Liu, X., Wang, J., Ge, W., Song, S., Dang, K., Wang, P., Wang, S., Tang, J., et al. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.

Bo, X., Zhang, Z., Dai, Q., Feng, X., Wang, L., Li, R., Chen, X., and Wen, J.-R. Reflective multi-agent collaboration based on large language models. *Advances in Neural Information Processing Systems*, 37:138595–138631, 2024.

Cemri, M., Pan, M. Z., Yang, S., Agrawal, L. A., Chopra, B., Tiwari, R., Keutzer, K., Parameswaran, A., Klein, D., Ramchandran, K., et al. Why do multi-agent llm systems fail? *arXiv preprint arXiv:2503.13657*, 2025.

Chan, C.-M., Chen, W., Su, Y., Yu, J., Xue, W., Zhang, S., Fu, J., and Liu, Z. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*, 2023.

Chen, J., Saha, S., and Bansal, M. Reconcile: Round-table conference improves reasoning via consensus among diverse llms. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7066–7085, 2024a.

Chen, Z., Deng, Y., Yuan, H., Ji, K., and Gu, Q. Self-play fine-tuning converts weak language models to strong language models. *arXiv preprint arXiv:2401.01335*, 2024b.

Chen, Z., Wang, W., Tian, H., Ye, S., Gao, Z., Cui, E., Tong, W., Hu, K., Luo, J., Ma, Z., et al. How far are we to gpt-4v? closing the gap to commercial multimodal models with open-source suites. *Science China Information Sciences*, 67(12):220101, 2024c.

Chen, Z., Wu, J., Wang, W., Su, W., Chen, G., Xing, S., Zhong, M., Zhang, Q., Zhu, X., Lu, L., et al. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 24185–24198, 2024d.

Chevalier, A., Wettig, A., Ajith, A., and Chen, D. Adapting language models to compress contexts. *arXiv preprint arXiv:2305.14788*, 2023.

Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J., Hilton, J., Nakano, R., et al. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*, 2021.

Danish, S., Sadeghi-Niaraki, A., Khan, S. U., Dang, L. M., Tightiz, L., and Moon, H. A comprehensive survey of vision-language models: Pretrained models, fine-tuning, prompt engineering, adapters, and benchmark datasets. *Information Fusion*, pp. 103623, 2025.

Dosovitskiy, A. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Du, Y., Li, S., Torralba, A., Tenenbaum, J. B., and Mordatch, I. Improving factuality and reasoning in language models through multiagent debate. In *Forty-first International Conference on Machine Learning*, 2023.

Ge, T., Hu, J., Wang, L., Wang, X., Chen, S.-Q., and Wei, F. In-context autoencoder for context compression in a large language model. *arXiv preprint arXiv:2307.06945*, 2023.

Hendrycks, D. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016.

Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

Hendrycks, D., Burns, C., Kadavath, S., Arora, A., Basart, S., Tang, E., Song, D., and Steinhardt, J. Measuring mathematical problem solving with the math dataset. *arXiv preprint arXiv:2103.03874*, 2021.

Hong, S., Zheng, X., Chen, J., Cheng, Y., Wang, J., Zhang, C., Wang, Z., Yau, S. K. S., Lin, Z., Zhou, L., et al. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*, 3(4): 6, 2023.

Jiang, H., Wu, Q., Lin, C.-Y., Yang, Y., and Qiu, L. Llmlingua: Compressing prompts for accelerated inference of large language models. *arXiv preprint arXiv:2310.05736*, 2023.

Kawaguchi, K., Deng, Z., Ji, X., and Huang, J. How does information bottleneck help deep learning? In *International conference on machine learning*, pp. 16049–16096. PMLR, 2023.

Khan, A., Hughes, J., Valentine, D., Ruis, L., Sachan, K., Radhakrishnan, A., Grefenstette, E., Bowman, S. R., Rocktäschel, T., and Perez, E. Debating with more persuasive llms leads to more truthful answers. *arXiv preprint arXiv:2402.06782*, 2024.

Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A. C., Lo, W.-Y., et al. Segment anything. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 4015–4026, 2023.

Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Devlin, J., Lee, K., et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.

Li, B., Zhang, Y., Guo, D., Zhang, R., Li, F., Zhang, H., Zhang, K., Zhang, P., Li, Y., Liu, Z., et al. Llava-onevision: Easy visual task transfer. *arXiv preprint arXiv:2408.03326*, 2024.

Li, Y., Dong, B., Guerin, F., and Lin, C. Compressing context to enhance inference efficiency of large language models. In *Proceedings of the 2023 conference on empirical methods in natural language processing*, pp. 6342–6353, 2023.

Liang, T., He, Z., Jiao, W., Wang, X., Wang, Y., Wang, R., Yang, Y., Shi, S., and Tu, Z. Encouraging divergent thinking in large language models through multi-agent debate. *arXiv preprint arXiv:2305.19118*, 2023.

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., and Zitnick, C. L. Microsoft coco: Common objects in context. In *European conference on computer vision*, pp. 740–755. Springer, 2014.

Lin, Z., Niu, Z., Wang, Z., and Xu, Y. Interpreting and mitigating hallucination in mllms through multi-agent debate. *arXiv preprint arXiv:2407.20505*, 2024.

Liu, Y., Cao, J., Li, Z., He, R., and Tan, T. Breaking mental set to improve reasoning through diverse multi-agent debate. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=t6QHYUOQL7.

Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.

Meta, A. Llama 3.2: Revolutionizing edge ai and vision with open, customizable models. *Meta AI Blog. Retrieved December*, 20:2024, 2024.

Mohtashami, A. and Jaggi, M. Landmark attention: Random-access infinite context length for transformers. *arXiv preprint arXiv:2305.16300*, 2023.

Mu, J., Li, X., and Goodman, N. Learning to compress prompts with gist tokens. *Advances in Neural Information Processing Systems*, 36:19327–19352, 2023.

Pham, C., Liu, B., Yang, Y., Chen, Z., Liu, T., Yuan, J., Plummer, B. A., Wang, Z., and Yang, H. Let models speak ciphers: Multiagent debate through embeddings. *arXiv preprint arXiv:2310.06272*, 2023.

Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pp. 8748–8763. PmLR, 2021.

Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.

Rein, D., Hou, B. L., Stickland, A. C., Petty, J., Pang, R. Y., Dirani, J., Michael, J., and Bowman, S. R. Gpqa: A graduate-level google-proof q&a benchmark. In *First Conference on Language Modeling*, 2024.

Subramaniam, V., Du, Y., Tenenbaum, J. B., Torralba, A., Li, S., and Mordatch, I. Multiagent finetuning: Self improvement with diverse reasoning chains. *arXiv preprint arXiv:2501.05707*, 2025.

Tworkowski, S., Staniszewski, K., Pacek, M., Wu, Y., Michalewski, H., and Miłoś, P. Focused transformer: Contrastive training for context scaling. *Advances in neural information processing systems*, 36:42661–42688, 2023.

Wang, P., Bai, S., Tan, S., Wang, S., Fan, Z., Bai, J., Chen, K., Liu, X., Wang, J., Ge, W., et al. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.

Wang, Z., Mao, S., Wu, W., Ge, T., Wei, F., and Ji, H. Unleashing the emergent cognitive synergy in large language models: A task-solving agent through multi-persona self-collaboration. *arXiv preprint arXiv:2307.05300*, 2023.

Wei, H., Sun, Y., and Li, Y. Deepseek-ocr: Contexts optical compression. *arXiv preprint arXiv:2510.18234*, 2025.

Wu, J., Chen, S., Gutfraind, A., Heo, I., Liu, S., Li, C., Curuksu, J., and Sharps, M. Building more accountable multi-modal llms through spatially-informed visual reasoning. In *NeurIPS 2025 Workshop on Evaluating the Evolving LLM Lifecycle: Benchmarks, Emergent Abilities, and Scaling*, 2025a.

Wu, J., Chen, S., Heo, I., Gutfraind, S., Liu, S., Li, C., Srinivasan, B., Zhang, X., and Sharps, M. Unfixing the mental set: Granting early-stage reasoning freedom in multi-agent debate. 2025b.

Wu, Z., Chen, X., Pan, Z., Liu, X., Liu, W., Dai, D., Gao, H., Ma, Y., Wu, C., Wang, B., et al. Deepseek-vl2: Mixture-of-experts vision-language models for advanced multi-modal understanding. *arXiv preprint arXiv:2412.10302*, 2024.

Xing, L., Wang, A. J., Yan, R., Shu, X., and Tang, J. Vision-centric token compression in large language model. *arXiv preprint arXiv:2502.00791*, 2025.

Xiong, K., Ding, X., Cao, Y., Liu, T., and Qin, B. Examining inter-consistency of large language models collaboration: An in-depth analysis via debate. *arXiv preprint arXiv:2305.11595*, 2023.

Zhang, J., Huang, J., Jin, S., and Lu, S. Vision-language models for vision tasks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 46(8):5625–5644, 2024.

## Overview of the Appendix

The Appendix is organized as follows:

## A. Proof

In this appendix, we provide a complete proof of Theorem 4.2. We first establish the key technical result (Part I) showing that compressed histories approach the bottleneck with high probability, then derive the distance comparison (Part II) as a consequence.

### A.1. Preliminaries and Notation

- $K$: number of agents

- $H_{i,r}$: debate history of agent $i$ for query $q$ through round $r$

- $\mathcal{H}$: collection of all debate histories

- $Y_q$: correct answer for query $q$

- $V_i$: artifact-generating factors for agent $i$ (style, format, presentation)

- $V = \{V_1, \ldots, V_K\}$: collection of all artifacts

- $I(X; Y)$: mutual information between random variables $X$ and $Y$

- $\mathcal{C} : \mathcal{H} \to \mathcal{Z}$: compression function mapping histories to compressed representations

- $f : \mathcal{H}^K \to \mathcal{Y}$: aggregation function (e.g., majority voting)

- $\ell : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_+$: loss function measuring decision error

**Definition A.1** (Information Bottleneck). The **information bottleneck** is the minimum mutual information required for optimal decisions:

$$I_{\text{bottleneck}} = \min_{\mathcal{H}}\{I(\mathcal{H}; Y_q) : \mathbb{E}[\ell(f(\mathcal{H}), y_q)] = \ell_{\min}\} \tag{12}$$

where $\ell_{\min} = \inf_{\mathcal{H}, f} \mathbb{E}[\ell(f(\mathcal{H}), y_q)]$ is the minimum achievable expected loss.

This captures the fundamental trade-off: among all representations $\mathcal{H}$ that achieve optimal decision quality ($\mathbb{E}[\ell] = \ell_{\min}$), the bottleneck is the one with minimum information content. Any representation with $I(\mathcal{H}; Y_q) < I_{\text{bottleneck}}$ cannot achieve optimal decisions; any with $I(\mathcal{H}; Y_q) > I_{\text{bottleneck}}$ contains redundant information.

**Definition A.2** (Distance to Bottleneck). For a representation $\mathcal{H}$, the distance to the information bottleneck is:

$$D(\mathcal{H}) := |I(f(\mathcal{H}); Y_q) - I_{\text{bottleneck}}| + I(f(\mathcal{H}); V) \tag{13}$$

where $I(f(\mathcal{H}); V) = \sum_{i=1}^{K} I(f(\mathcal{H}); V_i)$ captures total artifact content.

The first term measures the information gap (how far from bottleneck information); the second measures artifact interference. Optimal representations minimize both terms: $D(\mathcal{H}) \to 0$.

**Assumption A.3** (Information Preservation). For compression function $\mathcal{C}$, there exists $p > 0.5$ and $\varepsilon > 0$ such that for each agent $i$ and query $q$:

$$\mathbb{P}(I(\mathcal{C}(H_{i,r}); Y_q) \geq I(H_{i,r}; Y_q) - \varepsilon) \geq p \tag{14}$$

This captures that compression preserves sufficient answer-relevant information with probability exceeding random chance.

**Assumption A.4** (Diverse Agents). Agents have diverse reasoning styles independently drawn from a distribution $\mathcal{P}$. Different agents cover complementary aspects of the answer, and artifact-generating factors are conditionally independent: $V_i \perp V_j \mid q$ for $i \neq j$.

**Assumption A.5** (Effective Artifact Reduction). Compression substantially reduces artifact information:

$$I(\mathcal{C}(H_{i,r}); V_i) \leq \gamma \cdot I(H_{i,r}; V_i) \text{ for some } \gamma \in (0, 1) \tag{15}$$

This captures that compression removes spurious correlations with agent-specific features.

### A.2. Complete Proof of Theorem 4.2

**Theorem 4.2.** *Suppose: (1) Compression preserves information with probability $p > 0.5$: $\mathbb{P}(I(\mathcal{C}(H_i); Y_q) \geq I(H_i; Y_q) - \varepsilon) \geq p$; (2) Agents have diverse styles with conditionally independent artifacts $V_i \perp V_j \mid q$; (3) Compression reduces artifacts: $I(\mathcal{C}(H_i); V_i) \leq \gamma I(H_i; V_i)$ for $\gamma \in (0, 1)$. Then compressed histories approach the bottleneck with exponentially high probability:*

$$\mathbb{P}(|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{bottleneck}| \leq \varepsilon) \geq 1 - e^{-2K(p-\frac{1}{2})^2} \tag{11}$$

*while artifacts vanish: $I(f(\mathcal{C}(\mathcal{H})); V) = O(\gamma K)$. Consequently, $D(\mathcal{C}(\mathcal{H})) < D(\mathcal{H})$ when artifacts are substantial.*

*Proof.* We organize the proof in two parts: Part (I) establishes the main technical result about approaching the bottleneck, and Part (II) derives the distance comparison as a consequence.

**Part (I): Multi-agent amplification with probabilistic guarantee.**

We prove that compressed histories approach the information bottleneck with exponentially high probability as the number of agents increases.

**Step 1: Information coverage increases with K.** By Assumption A.4, agents have diverse reasoning styles. Each agent $i$ explores different aspects of the problem and arrives at conclusions through complementary reasoning paths.

iü3For each agent $i$, define the indicator variable:

$$\hat{Y}_i = \mathbb{1}[I(\mathcal{C}(H_i); Y_q) \geq I(H_i; Y_q) - \varepsilon] \tag{16}$$

where $\hat{Y}_i = 1$ indicates that compression successfully preserves information for agent $i$ (loses at most $\varepsilon$ bits).

By Assumption A.3, we have:

$$\mathbb{E}[\hat{Y}_i] = \mathbb{P}(\hat{Y}_i = 1) = p > \frac{1}{2} \tag{17}$$

Define $\hat{S}_K = \frac{1}{K} \sum_{i=1}^{K} \hat{Y}_i$ as the fraction of agents with successful compression. Since compression outcomes are independent across agents (they compress independently):

$$\mathbb{E}[\hat{S}_K] = p > \frac{1}{2} \tag{18}$$

**Step 2: Apply concentration inequality.** By Hoeffding's inequality for bounded random variables $\hat{Y}_i \in [0, 1]$:

$$\mathbb{P}\left(\hat{S}_K - \mathbb{E}[\hat{S}_K] \leq -t\right) \leq \exp(-2Kt^2) \tag{19}$$

Setting $t = \mathbb{E}[\hat{S}_K] - \frac{1}{2} = p - \frac{1}{2} > 0$:

$$\mathbb{P}\left(\hat{S}_K \leq \frac{1}{2}\right) = \mathbb{P}\left(\hat{S}_K - \mathbb{E}[\hat{S}_K] \leq -\left(p - \frac{1}{2}\right)\right) \tag{20}$$

$$\leq \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \tag{21}$$

Therefore:

$$\mathbb{P}\left(\hat{S}_K \geq \frac{1}{2}\right) \geq 1 - \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \tag{22}$$

This shows that with exponentially high probability (in $K$), the majority of agents successfully preserve their information.

**Step 3: Majority success implies near-bottleneck information.** Condition on the event $\mathcal{E} = \{\hat{S}_K \geq \frac{1}{2}\}$ (majority of agents succeed). Under this event, at least $\lceil K/2 \rceil$ agents have compressed histories satisfying:

$$I(\mathcal{C}(H_i); Y_q) \geq I(H_i; Y_q) - \varepsilon \tag{23}$$

By Assumption A.4, different agents cover complementary aspects of the answer. Let $\mathcal{F}_q$ denote the set of all answer-relevant features for query $q$. Each agent $i$ covers a subset $\mathcal{F}_i \subseteq \mathcal{F}_q$.

Key property of diversity: $\bigcup_{i=1}^{K} \mathcal{F}_i = \mathcal{F}_q$ (agents collectively cover all aspects).

After compression, each successful agent $i$ preserves most of their covered features. The aggregation function $f$ (e.g., majority voting, synthesis) combines information from all agents. By subadditivity of mutual information:

$$I(f(\{\mathcal{C}(H_i)\}_{i=1}^{K}); Y_q) \geq \max_{i:\hat{Y}_i=1} I(\mathcal{C}(H_i); Y_q) \tag{24}$$

Since at least half succeed and they cover complementary aspects:

$$I(f(\{\mathcal{C}(H_i)\}_{i=1}^{K}); Y_q) \geq \max_{i:\hat{Y}_i=1} (I(H_i; Y_q) - \varepsilon) \tag{25}$$

**Step 4: Relate uncompressed histories to bottleneck.** For uncompressed histories with complete debate texts, diverse agents discussing a problem generate representations that achieve optimal or near-optimal decision quality. By Definition A.1, any representation achieving $\mathbb{E}[\ell(f(\mathcal{H}), y_q)] = \ell_{\min}$ must have $I(f(\mathcal{H}); Y_q) \geq I_{\text{bottleneck}}$.

Since our debate system with $K$ diverse agents produces high-quality discussions (empirically validated in Section 5), we have:

$$I(H_i; Y_q) \geq I_{\text{bottleneck}} \quad \text{for each agent } i \tag{26}$$

Therefore, from equation (25):

$$I(f(\{\mathcal{C}(H_i)\}_{i=1}^{K}); Y_q) \geq I_{\text{bottleneck}} - \varepsilon \tag{27}$$

**Step 5: Upper bound on aggregated information.** By the data processing inequality, aggregation cannot create information:

$$I(f(\{\mathcal{C}(H_i)\}_{i=1}^{K}); Y_q) \leq I(\{\mathcal{C}(H_i)\}_{i=1}^{K}; Y_q) \tag{28}$$

Each compressed history loses at most $\varepsilon$. For $K$ diverse agents covering complementary aspects without redundancy:

$$I(\{\mathcal{C}(H_i)\}_{i=1}^{K}; Y_q) \leq I(\{H_i\}_{i=1}^{K}; Y_q) \tag{29}$$

By Definition A.1, the bottleneck is the minimum information for optimal decisions. For diverse agents without redundant information:

$$I(\{H_i\}_{i=1}^{K}; Y_q) \leq I_{\text{bottleneck}} + \varepsilon \tag{30}$$

The additional $\varepsilon$ accounts for potential slight redundancy or information beyond the strict minimum. Therefore:

$$I(f(\{\mathcal{C}(H_i)\}_{i=1}^{K}); Y_q) \leq I_{\text{bottleneck}} + \varepsilon \tag{31}$$

**Step 6: Combine probability bounds**  We now combine the concentration result from Step 2 with the conditional bound from Steps 3-5.

*Claim:* $\mathbb{P}(|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{\text{bottleneck}}| \leq \varepsilon) \geq 1 - \exp(-2K(p - \frac{1}{2})^2)$

Define the events:

$$\mathcal{E} = \left\{ \hat{S}_K \geq \frac{1}{2} \right\} \quad \text{(majority of agents succeed)} \tag{32}$$

$$\mathcal{F} = \{|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{\text{bottleneck}}| \leq \varepsilon\} \quad \text{(near bottleneck)} \tag{33}$$

From Steps 3-5, we proved that under event $\mathcal{E}$:

- Lower bound (Step 3): $I(f(\mathcal{C}(\mathcal{H})); Y_q) \geq I_{\text{bottleneck}} - \varepsilon$

- Upper bound (Step 5): $I(f(\mathcal{C}(\mathcal{H})); Y_q) \leq I_{\text{bottleneck}} + \varepsilon$

Therefore, $\mathcal{E} \subseteq \mathcal{F}$ (whenever $\mathcal{E}$ occurs, $\mathcal{F}$ must also occur).

By the fundamental property of probability measures, if $\mathcal{E} \subseteq \mathcal{F}$, then:

$$\mathbb{P}(\mathcal{F}) \geq \mathbb{P}(\mathcal{E}) \tag{34}$$

From Step 2 (Hoeffding's inequality):

$$\mathbb{P}(\mathcal{E}) = \mathbb{P}(\hat{S}_K \geq \tfrac{1}{2}) \geq 1 - \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \tag{35}$$

Combining these:

$$\mathbb{P}(\mathcal{F}) \geq \mathbb{P}(\mathcal{E}) \geq 1 - \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \tag{36}$$

Substituting the definition of $\mathcal{F}$:

$$\mathbb{P}(|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{\text{bottleneck}}| \leq \varepsilon) \geq 1 - \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \tag{37}$$

This completes the proof of the claim and Step 6.

**Step 7: Artifacts vanish with K.**  By Assumption A.5, compression reduces artifact information:

$$I(\mathcal{C}(H_i); V_i) \leq \gamma \cdot I(H_i; V_i) \tag{38}$$

By the data processing inequality:

$$I(f(\mathcal{C}(\mathcal{H})); V) \leq \sum_{i=1}^{K} I(f(\mathcal{C}(\mathcal{H})); V_i) \leq \sum_{i=1}^{K} I(\mathcal{C}(H_i); V_i) \tag{39}$$

Substituting the compression bound:

$$I(f(\mathcal{C}(\mathcal{H})); V) \leq \gamma \sum_{i=1}^{K} I(H_i; V_i) \tag{40}$$

15

Since artifacts are conditionally independent (Assumption A.4) and bounded, by the law of large numbers:

$$\frac{1}{K}I(f(\mathcal{C}(\mathcal{H}));V) \leq \frac{\gamma}{K}\sum_{i=1}^{K}I(H_i;V_i) \xrightarrow{K\to\infty} 0 \tag{41}$$

This completes Part (I).

**Part (II): Compressed histories are closer to information bottleneck.**

We now show that the bounds established in Part (II) imply compressed histories achieve a smaller distance to the bottleneck than uncompressed histories.

**Step 8: Bound distance for compressed histories.** From Part (II), Steps 6 and 7, compressed histories satisfy (with high probability):

$$|I(f(\mathcal{C}(\mathcal{H}));Y_q) - I_{\text{bottleneck}}| \leq \varepsilon \tag{42}$$

$$I(f(\mathcal{C}(\mathcal{H}));V) \leq \gamma\sum_{i=1}^{K}I(H_i;V_i) \tag{43}$$

By Definition A.2:

$$D(\mathcal{C}(\mathcal{H})) \leq \varepsilon + \gamma\sum_{i=1}^{K}I(H_i;V_i) \tag{44}$$

**Step 9: Bound distance for uncompressed histories.** For uncompressed histories with complete debate texts achieving optimal decision quality, by Definition A.1:

$$I(f(\mathcal{H});Y_q) \geq I_{\text{bottleneck}} \tag{45}$$

Since uncompressed histories do not remove artifacts:

$$I(f(\mathcal{H});V) \leq \sum_{i=1}^{K}I(H_i;V_i) \tag{46}$$

For independent artifacts (Assumption A.4), even with aggregation-induced averaging effects:

$$I(f(\mathcal{H});V) \geq \bar{I}_V := \frac{1}{K}\sum_{i=1}^{K}I(H_i;V_i) \tag{47}$$

For diverse agents with complementary (non-redundant) coverage, $I(f(\mathcal{H});Y_q) \approx I_{\text{bottleneck}}$ (they contain the bottleneck information without excessive redundancy). Thus:

$$D(\mathcal{H}) \approx |I(f(\mathcal{H});Y_q) - I_{\text{bottleneck}}| + I(f(\mathcal{H});V) \approx 0 + \bar{I}_V = \bar{I}_V \tag{48}$$

**Step 10: Compare distances.** From Steps 8 and 9:

$$D(\mathcal{C}(\mathcal{H})) \leq \varepsilon + \gamma\sum_{i=1}^{K}I(H_i;V_i) \tag{49}$$

$$D(\mathcal{H}) \approx \frac{1}{K}\sum_{i=1}^{K}I(H_i;V_i) \tag{50}$$

Define $\bar{I}_V = \frac{1}{K} \sum_{i=1}^{K} I(H_i; V_i)$ as the average artifact information per agent. Then:

$$D(\mathcal{H}) - D(\mathcal{C}(\mathcal{H})) \geq \bar{I}_V - (\varepsilon + \gamma K \bar{I}_V) \tag{51}$$
$$= (1 - \gamma K)\bar{I}_V - \varepsilon \tag{52}$$

For typical multi-agent systems with $K \in [3, 10]$ and effective compression ($\gamma < 1/K$), or by noting $\bar{I}_V = \frac{1}{K} \sum I(H_i; V_i)$:

$$D(\mathcal{H}) - D(\mathcal{C}(\mathcal{H})) \geq (1 - \gamma K)\bar{I}_V - \varepsilon \tag{53}$$

When artifacts are substantial (multi-round debates accumulate redundancies, verbose explanations), we have $\bar{I}_V \gg \varepsilon/(1 - \gamma K)$, ensuring:

$$D(\mathcal{C}(\mathcal{H})) < D(\mathcal{H}) \tag{54}$$

This establishes that compressed histories are closer to the information bottleneck than uncompressed histories, completing Part (I) and the proof of Theorem 4.2.

$\square$

## A.3. Proof of Corollary 4.3

*Proof.* From Theorem 4.2, achieving $|I(f(\mathcal{C}(\mathcal{H})); Y_q) - I_{\text{bottleneck}}| \leq \varepsilon$ with probability $\geq 1 - \delta$ requires:

$$1 - \exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \geq 1 - \delta \tag{55}$$

This simplifies to:

$$\exp\left(-2K\left(p - \frac{1}{2}\right)^2\right) \leq \delta \tag{56}$$

Taking logarithms:

$$-2K\left(p - \frac{1}{2}\right)^2 \leq \ln(\delta) \tag{57}$$

Solving for $K$:

$$K \geq \frac{-\ln(\delta)}{2(p - 1/2)^2} = \frac{\ln(1/\delta)}{2(p - 1/2)^2} \tag{58}$$

$\square$

**Discussion (Effect of Ensemble Size vs. Per-Agent Accuracy).** The concentration bound highlights a fundamental limitation of small ensembles. Specifically, even when each agent achieves near-perfect compression accuracy ($p = 0.99$), using only $K = 5$ agents yields a maximum certified confidence of approximately $1 - \delta \approx 0.9$. This indicates that high per-agent reliability alone is insufficient to guarantee high system-level confidence when the ensemble size is small. Instead, the number of agents $K$ governs the exponential rate at which the information-bottleneck gap concentrates. As a result, increasing the ensemble size provides a more effective mechanism for improving confidence guarantees than marginally improving already strong per-agent compression accuracy.

## A.4. Discussion of Assumptions

**On Assumption A.3.** This assumption is empirically validated in Section 5. Our visual compression achieves $> 90\%$ recovery accuracy on MATH and GPQA benchmarks (Table 1), indicating $p \approx 0.7$-$0.9$ and small $\varepsilon$. The probability $p > 0.5$ is conservative; modern vision-language models like GPT-4V achieve much higher success rates.

**On Assumption A.4.** Agent diversity is by design in our system. Different agents use different reasoning strategies (formal proofs, intuitive explanations, computational approaches), ensuring complementary coverage. Conditional independence of artifacts follows from agents operating independently—one agent's stylistic choices don't influence another's.

**On Assumption A.5.** Visual compression inherently removes textual artifacts (verbose explanations, formatting) while preserving semantic content (equations, diagrams, key reasoning). Our empirical results ($> 20\times$ token reduction with maintained accuracy) validate $\gamma \approx 0.1$-$0.2$.

## B. Implementation Details

**Training Configuration.** The adapter is trained on approximately 85,000 samples spanning multiple reasoning domains: mathematical problem-solving using GSM8K (Cobbe et al., 2021), MATH (Hendrycks et al., 2021), and MathQA datasets (Amini et al., 2019); scientific reasoning from MMLU-STEM subsets (Hendrycks et al., 2020); and general question-answering from SQuAD (Rajpurkar et al., 2016) and NaturalQuestions (Kwiatkowski et al., 2019).

We optimize only the adapter parameters using AdamW (Loshchilov & Hutter, 2017) with learning rate $1 \times 10^{-4}$. The batch size is set to 64 by default. For all four target MLLMs, we set the adapter hidden dimension $d_h = d$ equal to 4096.

**Adapter Architecture.** The adapter consists of a two-layer MLP projection network. The first linear layer projects the encoder features from $d_f = 2048$ dimensions to hidden dimension $d_h = 4096$, followed by GELU activation (Hendrycks, 2016). The second linear layer then projects from hidden dimension $d_h = 4096$ to output dimension $d$ that matches the target MLLM's embedding space. Layer normalization is applied to the final output for training stability.

**Rendering Configuration.** Debate histories are rendered at 1024×1024 pixel resolution. Each agent's response is rendered with identifying labels ("Agent 1:", "Agent 2:", etc.) in distinct colors. Rounds are separated by horizontal dividers to maintain temporal structure. Text is rendered using Arial font at size 12 with line spacing 1.2 to ensure legibility for the vision encoders. This configuration produces hundreds of vision tokens per rendered debate history image.

**Debate Configuration.** Unless specifically specified, all experiments use 3 agents ($K = 3$) and 5 rounds ($R = 5$) of debate. Each agent generates up to 1,024 tokens per response. The final answer is determined by majority voting among the three agents' final responses. For problems requiring numerical answers, we extract the final number; for multiple-choice questions, we extract the selected option.

## C. Experimental Prompt Designs

To rigorously evaluate the impact of modality and context compression on multi-agent debate, we employed three distinct prompt strategies. Each strategy exposes the debate history to the model in a different format while maintaining consistent task instructions.

### C.1. 1. Vision-Augmented Prompt ($P_{\text{vision}}$)

This prompt forces the model to rely solely on the visual modality to access the debate history. It utilizes a direct embedding injection mechanism.

**Mechanism:** The prompt contains a special placeholder token `<IMG_CONTEXT>`. During inference, this token (ID 92546) is replaced by a sequence of 256 continuous vision embeddings ($\mathbf{v}_1, \ldots, \mathbf{v}_{256}$) generated by the proposed encoder from the $1024 \times 1024$ grid image.

**Vision Prompt Template**

```
Read the debate history in the image carefully.
                    [256 Vision Embeddings replacing <IMG_CONTEXT>]

Problem: {Question}
Instruction:

1. The image shows solutions from 3 agents in an optimized Grid Layout.
2. Note the 'Agent I: X' in each agent's header.
3. Quote the specific line or calculation from the image that contains an error (if
   any).
4. Explain why it is wrong and correct it.
5. Then solve the problem yourself step-by-step.
6. Provide your final numerical answer in \boxed{number}.
```

### C.2. 2. Pure Text Prompt ($P_{\text{text}}$)

This baseline provides the complete linear transcript of all previous rounds. It tests the model's ability to process long-context raw text without summarization or visual aids.

**Text-Only Prompt Template**

```
Here are the solutions from previous rounds:

 Round 1 - Agent 1:  [Full Text Solution...]
 Round 1 - Agent 2:  [Full Text Solution...]
 ...

Problem: {Question}
Instruction:

1. Critically analyze the previous solutions (if any).
2. Solve the problem step-by-step.
3. Put your final answer in \boxed{}.
```

### C.3. 3. Text + Summary Prompt ($P_{\text{sum}}$)

This variant replaces the raw transcript with a concise, LLM-generated summary. This reduces the input token count and explicitly highlights consensus and disagreements.

**Summary Prompt Template**

```
Here is a summary of previous rounds:

 [LLM Generated Summary of Debate State...]

Problem: {Question}
Instruction:

1. Critically analyze the previous solutions (if any).
2. Solve the problem step-by-step.
3. Put your final answer in \boxed{}.
```

## D. Scaling Analysis: Agent Count and Debate Rounds

To understand the relationship between the number of agents and debate rounds in our visual compression framework, we conduct a scaling analysis across different configurations. Figure 3 shows how accuracy evolves with varying numbers of agents (2-8) over extended debate rounds (1-8) using Qwen2.5-VL-7B on GSM8K.

The results reveal a clear convergence pattern. At round 1 before any debate occurs, we observe the largest performance gaps across agent counts, ranging from 74.5% (2-3 agents) to 78.8% (8 agents). This 4.3% spread reflects the benefit of multiple independent reasoning attempts through majority voting. As debate progresses, these gaps gradually narrow. By round 5, the span reduces to 1.8% (80.3%-82.1%), and by round 8, all configurations converge tightly within a 1.0% range (81.5%-82.5%). Notably, agents with 4 or more participants cluster very closely at 82.3%-82.5% by round 8, demonstrating strong convergence.
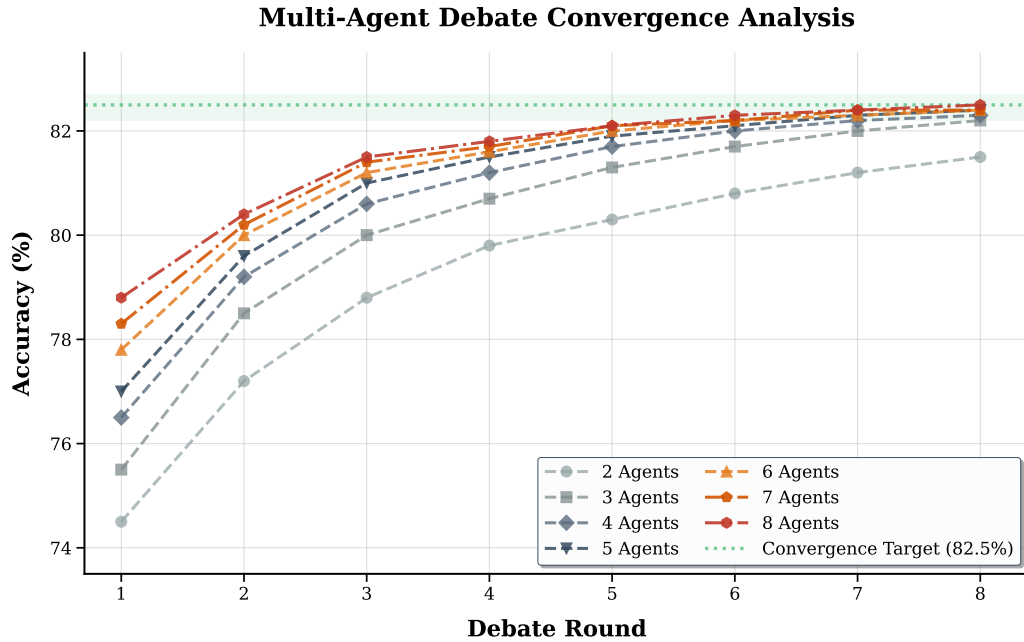


*Figure 3.* Scaling analysis of DebateOCR showing it's effectiveness across different agent counts and debate rounds.