# Critical Analysis of Perception Tokens Enhance Visual Reasoning in Multimodal Language Models

Part II: Strengths, Limitations, and Research Implications

Kai-Yu Lu

2025/10/16

## 1 Methodological Strengths

**Experimental design choices**

- **Single-model inference with explicit visual intermediates**. The method augments an instruction-tuned multimodal language model with discrete perception tokens so that intermediate depth codes or bounding-box coordinates are generated and then consulted during reasoning. This isolates the effect of in-model perceptual reasoning and avoids tool routing variability.

- **Progressive curriculum**. Training mixes three regimes per task, namely atomic token generation, chain-of-thought with perception tokens, and direct labeling. A temperature-controlled schedule formalizes an easy-to-hard progression and mitigates catastrophic forgetting.

- **Evaluation without options**. Multiple-choice options are removed from depth and counting benchmarks, enforcing free-form prediction and suppressing option-prior biases.

**Novel evaluation protocols or metrics**

- **Hard depth splits**. Relative-depth sets with three to five markers placed near mid height reduce trivial height cues and provide a controlled difficulty ladder.

- **Programmatic evidence checks**. Decoded depth-token maps are compared at marked points to verify that intermediate evidence is consistent with final answers, adding an interpretable consistency signal.

**Data collection strategies**

- **Aligned tokenization**. Depth is vector-quantized into a fixed $10 \times 10$ code grid using a compact codebook with start and end delimiters. Counting uses discrete coordinate tokens after resizing to a canonical resolution. Supervision thus matches the discrete structures expected at inference.

- **Task-specific mixers**. Depth uses a large set for token learning together with smaller curated sets for chain-of-thought and direct labeling. Counting follows an analogous design. This separation helps disentangle token acquisition from downstream reasoning.

1

### Systematic comparisons and ablations

- **Reasoning-step ablations**. Removing either the coordinate-identification step or the depth-token step degrades relative-depth accuracy, clarifying their complementary roles.

- **Token-type ablations**. Discrete coordinate tokens outperform plain text numerals for localization-driven counting, indicating advantages from structured token spaces.

- **Reconstruction objective**. Adding a decoder-based reconstruction penalty offers small but interpretable gains, which helps characterize its cost–benefit profile at the reported scale.

### Transparency and reproducibility

- **Implementation specifics**. The paper states backbone choice, frozen components, LoRA usage, vocabulary expansion, decoding constraints, training epochs, and hardware, which supports reproducibility.

- **Cross-task checks**. Depth-trained models are evaluated on related depth benchmarks without architectural changes, providing evidence of transfer beyond the construction set.

## 2 Key Limitations

### Dataset scale, diversity, and supervision

- **Small supervised chain-of-thought sets**. The curated chain-of-thought and direct-labeling splits are modest, limiting analysis of long-tail phenomena and cross-domain robustness.

- **Pseudo-label dependence**. Depth supervision is derived from an estimator rather than ground truth, which can propagate estimator biases into token learning.

### Modeling scope

- **Limited perception families**. Experiments instantiate depth tokens and box tokens for counting. Other mid-level signals such as surface normals, optical flow, and keypoints are discussed conceptually but not validated empirically.

### Evaluation breadth and error analysis

- **Benchmark concentration**. Results emphasize curated depth and counting suites. Open-world scenes without markers, heavy occlusions, and extreme scale variation are underexplored.

- **Failure-mode taxonomy**. Qualitative examples are present, yet a large-scale labeled taxonomy of errors is not reported, which limits targeted remediation.

### System cost reporting

- **Parameter and latency deltas**. Vocabulary growth enlarges embeddings and the language modeling head. Detailed parameter increments and throughput or latency impacts per token family are not fully quantified.

# 3 Technical Bottlenecks

## Core architectural and algorithmic constraints

- **Discrete codebook compression**. A fixed-size codebook with a fixed grid compresses continuous geometry into coarse codes, which can underfit fine depth gradients and thin structures.

- **Rigid constrained decoding**. Enforcing a fixed-length depth-token block simplifies validation but prevents adaptive spatial granularity and variable-length evidence chains.

## Information bottlenecks and integration

- **Coupling between chain-of-thought and tokens**. Best performance requires both coordinate extraction and depth-token reasoning, which implies sensitivity to prompt templates and step ordering and can reduce robustness under prompt variations.

## Trade-offs

- **Reconstruction objective versus compute**. Decoder-based reconstruction improves interpretability with limited quantitative gains at the reported scale, creating tension between accuracy and added compute.

# 4 Research Implications

## Capabilities versus requirements

- **Mid-level structure as a catalyst for reasoning**. Gains on relative depth and counting indicate that explicit visual abstractions used as intermediate steps are beneficial when tasks depend on perception rather than language priors.

## Benchmark–deployment gap

- **Marked versus unmarked scenes**. Relative depth benchmarks rely on marked points, whereas real applications require saliency selection and occlusion reasoning without markers, indicating a protocol gap.

## Connections to broader challenges

- **Unified token spaces for auditable decisions**. Requiring models to use generated visual tokens for answers aligns with needs in robotics, medical triage, and embodied agents where intermediate states must be auditable.

# 5 Potential Research Directions

## Representations and architectures

- Extend token families to surface normals, keypoints, and instance masks. Adopt hierarchical or multi-scale codebooks that allow variable-length evidence instead of a fixed grid.

- Replace rigid fixed-length constraints with grammar-guided constrained decoding that preserves structural validity while permitting variable token counts.

**Evaluation methodologies**

- Construct unmarked relative-depth tests with automatic saliency selection and controlled occlusions. Report per-sample evidence–answer consistency rates using deterministic validators over decoded tokens.

- Provide parameter, latency, and throughput deltas attributable to each token family. Include ablations on decoding constraints to quantify accuracy–efficiency trade-offs.

**Integration strategies**

- Share projection heads across perception vocabularies and apply mixture-of-experts routing to amortize embedding and output growth while preserving specialization.

- Jointly train on marked and unmarked variants to reduce prompt brittleness at the interface between chain-of-thought and perception tokens.

**Robustness and reliability**

- Calibrate confidence over token sequences and add abstention rules when decoded evidence conflicts with answers, together with entropy-based early-stop policies.

- Quantify pseudo-label bias by training with multiple depth estimators and measuring variance in downstream accuracy and evidence–answer agreement.

**Personalization and adaptation**

- Introduce lightweight domain adapters that specialize token vocabularies for aerial, endoscopic, or industrial imagery while retaining a shared global vocabulary for portability.

# 6   Conclusion

The study demonstrates that perception tokens combined with a progressive curriculum improve perception-heavy visual reasoning under single-model inference. Strengths include principled curriculum design, interpretable evidence checks, option removal, and clear ablations. The most significant limitations involve modest supervised chain-of-thought scale, reliance on pseudo labels, fixed token budgets, and limited failure-mode taxonomy. The most promising directions are multi-scale tokenization with grammar-guided decoding, evidence–answer auditing, explicit cost reporting, and shared heads with expert routing to bound parameter growth.