

# Progressive Knowledge-Guided Distillation for Multimodal Reasoning Models

Prasanth Yadla

Independent Researcher

pyadla2@alumni.ncsu.edu

## Abstract

Contemporary Multimodal Large Language Models (MLLMs) demonstrate exceptional capabilities in synthesizing visual and linguistic information with external knowledge repositories for sophisticated reasoning applications. Nevertheless, their substantial computational requirements present significant obstacles for implementation in resource-constrained settings. This research presents a knowledge-guided distillation methodology that facilitates the transfer of reasoning capabilities from large, knowledge-enriched teacher networks to streamlined student frameworks. Our technique preserves 87.3% of the teacher model’s performance while achieving a  $1.4\times$  acceleration in inference speed and a 49% reduction in parameter count. Evaluations on knowledge-enhanced visual question answering datasets demonstrate that our distillation approach surpasses conventional distillation methods by 0.4 percentage points while maintaining comparable factual accuracy. These findings establish a viable pathway for developing efficient MLLMs optimized for knowledge-intensive applications demanding real-time processing capabilities.

## 1 Introduction

Multimodal Large Language Models have demonstrated unprecedented performance in visual-textual comprehension tasks through the integration of external knowledge repositories, including structured knowledge graphs [Lu et al., 2022; Li et al., 2019; Tan and Bansal, 2019]. However, the massive parametric complexity of these models, frequently exceeding one billion parameters, introduces substantial deployment constraints, particularly in edge computing environments where computational resources remain severely limited.

Knowledge distillation presents a theoretically grounded approach for model compression, enabling the transfer of learned representations from computationally intensive teacher models to architecturally efficient student networks [Hinton, Vinyals, and Dean, 2015; Romero et al., 2014]. Nevertheless, conventional distillation methodologies

encounter significant challenges when applied to knowledge-grounded multimodal architectures, frequently failing to preserve the intricate reasoning patterns necessary for effective external knowledge integration.

This work introduces a knowledge-guided distillation framework that systematically transfers knowledge-grounded reasoning capabilities from teacher to student models. In contrast to traditional distillation approaches that prioritize output mimicry, the proposed methodology incorporates external knowledge graphs to guide the distillation process, ensuring that student models acquire the capacity to effectively utilize external knowledge for complex reasoning tasks.

The primary contributions of this research encompass: first, a multi-level knowledge distillation algorithm that transfers factual knowledge, reasoning patterns, and cross-modal attention mechanisms; second, a progressive training strategy that incrementally introduces knowledge complexity during the distillation process; third, comprehensive experimental validation demonstrating improvements over standard distillation baselines while maintaining deployment efficiency.

## 2 Related Work

### 2.1 Multimodal Knowledge Integration

Recent developments in multimodal architectures, including LXMERT, VL-BERT, and UNITER, have established that incorporating structured knowledge substantially enhances performance on reasoning-intensive tasks [Su et al., 2019; Chen et al., 2020; Tan and Bansal, 2019]. These architectures typically employ knowledge retrieval mechanisms from knowledge graphs during inference or integrate knowledge during pre-training phases. However, the computational overhead associated with knowledge retrieval and processing compounds the deployment challenges inherent in large-scale models.

### 2.2 Knowledge Distillation

Knowledge distillation facilitates the transfer of learned representations from large teacher models to compact student architectures [Hinton, Vinyals, and Dean, 2015]. Recent methodological advances include feature-level distillation [Romero et al., 2014], attention transfer mechanisms [Zagoruyko and Komodakis, 2016], and multi-teacher frameworks [You, Xu, and Tao, 2017]. However, limited

research addresses knowledge-grounded reasoning in multimodal contexts, where external knowledge integration introduces additional complexity to the distillation process.

### 2.3 Multimodal Model Compression

Contemporary efforts in multimodal compression have investigated various techniques including network pruning [Michel, Levy, and Neubig, 2019], quantization methods [Zafirir et al., 2019], and architectural simplification. However, these approaches frequently neglect the preservation of knowledge-grounded reasoning capabilities, resulting in disproportionate performance degradation on knowledge-intensive tasks.

## 3 Methodology

### 3.1 Problem Formulation

Given a large-scale pre-trained MLLM teacher model  $\mathcal{T}$  with parameters  $\theta_T$  and a compact student model  $\mathcal{S}$  with parameters  $\theta_S$ , along with an external multimodal knowledge graph  $\mathcal{KG}$ , the objective is to train  $\mathcal{S}$  to preserve  $\mathcal{T}$ ’s knowledge-grounded reasoning capabilities while achieving substantial computational efficiency.

### 3.2 Knowledge-Guided Distillation Algorithm

The proposed algorithm incorporates multiple distillation objectives to transfer distinct aspects of knowledge-grounded reasoning:

$$\mathcal{L}_{total} = \alpha\mathcal{L}_{task} + \beta\mathcal{L}_{output} + \gamma\mathcal{L}_{feature} + \delta\mathcal{L}_{knowledge} + \epsilon\mathcal{L}_{attention} \quad (1)$$

where each component targets specific aspects of knowledge transfer.

#### Output-Level Distillation

The standard knowledge distillation loss utilizing temperature-scaled softmax distribution:

$$\mathcal{L}_{output} = \text{KL} \left( \sigma \left( \frac{z_S}{T} \right) \parallel \sigma \left( \frac{z_T}{T} \right) \right) \quad (2)$$

where  $z_S$  and  $z_T$  represent student and teacher logits respectively, and  $T$  denotes the temperature parameter.

#### Feature-Level Knowledge Transfer

Intermediate representation alignment between teacher and student models, with emphasis on knowledge-aware feature mappings:

$$\mathcal{L}_{feature} = \frac{1}{N} \sum_{i=1}^N \|\phi_S(x_i) - W \cdot \phi_T(x_i)\|_2^2 \quad (3)$$

where  $\phi_S$  and  $\phi_T$  denote student and teacher feature extractors respectively, and  $W$  represents a learned projection matrix to accommodate dimensionality differences.

### Knowledge Consistency Loss

A knowledge-specific loss function that encourages alignment with external knowledge facts:

$$\mathcal{L}_{knowledge} = \frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} \text{CE}(p_S(k|x), p_T(k|x)) \quad (4)$$

where  $\mathcal{K}$  represents the set of retrieved knowledge facts, and  $p_S(k|x)$  and  $p_T(k|x)$  denote student and teacher probabilities for fact  $k$  conditioned on input  $x$ .

### Cross-Modal Attention Transfer

To preserve cross-modal reasoning patterns, attention mechanisms connecting visual and textual information are distilled:

$$\mathcal{L}_{attention} = \frac{1}{H} \sum_{h=1}^H \text{MSE}(A_S^{(h)}, A_T^{(h)}) \quad (5)$$

where  $A_S^{(h)}$  and  $A_T^{(h)}$  represent attention weights for head  $h$  in student and teacher models respectively.

### 3.3 Progressive Knowledge Distillation

A curriculum learning approach is adopted to gradually increase the complexity of knowledge presented to the student model. In the first stage, the model focuses on **fundamental multimodal understanding**, where it is trained on visual-textual alignment tasks without relying on external knowledge, thereby learning essential cross-modal representations. The second stage introduces **elementary knowledge integration**, in which single-hop knowledge facts are incorporated, enabling the student to utilize fundamental external information. Finally, the third stage targets **complex reasoning patterns**, where multi-hop reasoning examples are employed to transfer the teacher model’s ability to chain multiple knowledge facts for sophisticated inference.

## 4 Experiments

### 4.1 Experimental Setup

The experimental setup involves both teacher and student models, benchmark datasets, external knowledge sources, baseline methods, and multiple evaluation metrics. The teacher model is LXMERT-Base with 183M parameters, augmented with knowledge graph integration modules. The student models are compact architectures with 93M parameters, employing DistilBERT as the language backbone and ResNet-34 for visual encoding. For evaluation, we use OK-VQA [Marino et al., 2019] (14,031 questions) and FVQA [Wang et al., 2017] (5,826 questions), which are widely adopted benchmarks for knowledge-grounded visual question answering. External knowledge is provided by ConceptNet 5.7 [Speer, Chin, and Havasi, 2017] and Visual Genome [Krishna et al., 2017], covering both commonsense and visual concepts. The proposed approach is compared against several baselines, including standard knowledge distillation [Hinton, Vinyals, and Dean, 2015], feature-level distillation [Romero et al., 2014], attention distillation [Zagoruyko and Komodakis, 2016], and multi-teacher distillation frameworks. Model performance is assessed using accuracy, knowledge preservation score, inference speedup, and parameter efficiency as evaluation metrics.

## 4.2 Main Results

Table 1 presents comprehensive comparative results of the proposed knowledge-guided distillation against standard baseline methodologies.

Table 1: Performance comparison of distillation methods on OK-VQA

Method	Params (M)	Accuracy (%)	Knowledge Preserv. (%)	Speedup (x)
Teacher (LXMERT-Base)	183	42.7	100.0	1.0x
Standard Distillation	93	35.1	82.2	1.3x
Feature Distillation	93	36.4	85.3	1.3x
Attention Distillation	93	36.8	86.2	1.3x
Multi-Teacher Distillation	93	36.9	86.5	1.2x
<b>Ours (Single-Stage)</b>	<b>93</b>	<b>37.0</b>	<b>86.7</b>	<b>1.4x</b>
<b>Ours (Progressive)</b>	<b>93</b>	<b>37.3</b>	<b>87.3</b>	<b>1.4x</b>

The proposed knowledge-guided distillation achieves 37.3% accuracy with progressive training, representing 87.3% retention of teacher performance (42.7%). The methodology maintains 87.3% of the teacher’s factual reasoning capabilities while achieving 1.4x speedup and 49% parameter reduction.

## 4.3 Progressive Training Analysis

Table 2 demonstrates the effectiveness of the proposed progressive curriculum approach.

Table 2: Progressive training stage analysis

Training Stage	Accuracy (%)	Knowledge Preserv. (%)	Training Time (hrs)
Stage 1: Basic MM	31.8	74.5	12
Stage 2: Simple KG	35.2	82.5	18
Stage 3: Complex Reason.	37.3	87.3	24
End-to-End Training	36.7	85.9	28

Progressive training methodology achieves superior final performance while providing better knowledge retention and more stable convergence patterns compared to end-to-end training approaches.

## 4.4 Ablation Studies

Table 3 quantifies the contribution of each distillation component.

Table 3: Ablation study on loss components

Configuration	Accuracy (%)	Knowledge Preserv. (%)
Full Framework	<b>37.3</b>	<b>87.3</b>
w/o Knowledge Loss	36.1	84.6
w/o Attention Transfer	36.8	86.2
w/o Feature Distillation	36.9	86.4
w/o Progressive Training	36.7	85.9
Output Distillation Only	35.1	82.2

Table 4: Cross-dataset evaluation results

Method	OK-VQA	FVQA	GQA
Teacher Model	42.7	56.8	39.2
Standard Distillation	35.1	46.3	31.8
<b>Ours (Progressive)</b>	<b>37.3</b>	<b>49.6</b>	<b>33.9</b>
<b>Retention Rate</b>	<b>87.3%</b>	<b>87.3%</b>	<b>86.5%</b>

## 4.5 Cross-Dataset Evaluation

Generalization capability assessment across heterogeneous VQA datasets:

The proposed approach maintains consistent performance across domains, achieving 86.5-87.3% retention rates compared to 81.2-82.2% for standard distillation methodologies.

## 4.6 Computational Efficiency Analysis

The proposed knowledge-guided distillation algorithm introduces a moderate training overhead when compared to standard distillation due to the additional processes of knowledge retrieval and multi-level loss computation. Despite this increase in training cost, the distilled models offer meaningful efficiency gains at inference time, achieving a 1.4x speedup over teacher models while retaining 87.3% of their accuracy. Furthermore, memory consumption is reduced from 1.2 GB to 0.61 GB, demonstrating the practicality of the approach for deployment in resource-constrained environments.

## 5 Analysis and Discussion

### 5.1 Knowledge Transfer Effectiveness

An analysis of the proposed framework highlights its effectiveness in transferring distinct reasoning patterns. For **single-hop reasoning**, which involves direct fact lookup, the transfer effectiveness reaches 78.1%. In the case of **multi-hop reasoning**, where complex inference chains are required, the effectiveness is lower at 71.3%. Finally, **commonsense reasoning**, which relies on the implicit application of external knowledge, achieves 74.7% transfer effectiveness. These results demonstrate that the framework successfully preserves diverse reasoning capabilities, with particularly strong performance in direct fact retrieval tasks.

### 5.2 Limitations

The proposed approach exhibits certain limitations in specific scenarios. First, in the case of **complex visual reasoning**, performance degradation becomes more evident for questions that demand fine-grained visual analysis in combination with extensive background knowledge, showing a 8-12% additional drop compared to simpler visual tasks. Second, with respect to **domain-specific knowledge**, transfer effectiveness diminishes when applied to highly specialized areas that are insufficiently represented in general-purpose knowledge graphs. Finally, the framework shows a notable **dependency on knowledge graph coverage**, with results indicating a 6-9% reduction in accuracy when coverage falls below 75%. These findings suggest that future improvements may require enhanced visual reasoning modules, domain-adaptive knowledge sources, and robustness against incomplete knowledge coverage.

## 6 Conclusion

This research presents a knowledge-guided distillation framework for compressing multimodal language models while preserving knowledge-grounded reasoning capabilities. The proposed approach leverages external knowledge graphs to guide the teacher-student transfer process, achieving modest but consistent improvements compared to standard distillation techniques.

Experimental results demonstrate practical efficiency gains: progressive distillation achieves 87.3% accuracy retention with 1.4× inference speedup and 49% parameter reduction. Knowledge preservation scores of 87.3% confirm that the proposed approach successfully transfers reasoning capabilities to compact student models, though with some expected degradation in complex reasoning tasks.

Future research directions should explore adaptive knowledge selection during distillation and investigate domain-specific knowledge transfer techniques. The broader impact encompasses enabling deployment of multimodal reasoning models in resource-constrained environments while acknowledging the trade-offs between efficiency and performance.

**Limitations:** The proposed approach requires high-quality knowledge graphs and introduces moderate training complexity. Performance gains are most pronounced for knowledge-intensive tasks, with diminishing benefits for purely perceptual reasoning. The method shows sensitivity to knowledge graph coverage and may require domain-specific adaptations for specialized applications.

## References

- [Chen et al., 2020] Chen, Y.-C.; Li, L.; Yu, L.; Kholy, A. E.; Ahmed, F.; Gan, Z.; Cheng, Y.; and Liu, J. 2020. Uniter: Universal image-text representation learning. In *ECCV*.
- [Hinton, Vinyals, and Dean, 2015] Hinton, G.; Vinyals, O.; and Dean, J. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
- [Krishna et al., 2017] Krishna, R.; Zhu, Y.; Groth, O.; Johnson, J.; Hata, K.; Kravitz, J.; Chen, S.; Kalantidis, Y.; Li, L.-J.; Shamma, D. A.; et al. 2017. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *International journal of computer vision* 123(1):32–73.
- [Li et al., 2019] Li, L. H.; Yatskar, M.; Yin, D.; Hsieh, C.-J.; and Chang, K.-W. 2019. Visualbert: A simple and performant baseline for vision and language. *arXiv preprint arXiv:1908.03557*.
- [Lu et al., 2022] Lu, J.; Batra, D.; Parikh, D.; and Lee, S. 2022. Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *NeurIPS*.
- [Marino et al., 2019] Marino, K.; Rastegari, M.; Farhadi, A.; and Mottaghi, R. 2019. Ok-vqa: A visual question answering benchmark requiring external knowledge. In *CVPR*.
- [Michel, Levy, and Neubig, 2019] Michel, P.; Levy, O.; and Neubig, G. 2019. Are sixteen heads really better than one? In *NeurIPS*.

- [Romero et al., 2014] Romero, A.; Ballas, N.; Kahou, S. E.; Chassang, A.; Gatta, C.; and Bengio, Y. 2014. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*.
- [Speer, Chin, and Havasi, 2017] Speer, R.; Chin, J.; and Havasi, C. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *AAAI*.
- [Su et al., 2019] Su, W.; Zhu, X.; Cao, Y.; Li, B.; Lu, L.; Wei, F.; and Dai, J. 2019. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.02530*.
- [Tan and Bansal, 2019] Tan, H. and Bansal, M. 2019. Lxmert: Learning cross-modality encoder representations from transformers. In *EMNLP*.
- [Wang et al., 2017] Wang, P.; Wu, Q.; Shen, C.; Dick, A.; and van den Hengel, A. 2017. Fvqa: Fact-based visual question answering. *IEEE transactions on pattern analysis and machine intelligence* 40(10):2413–2427.
- [You, Xu, and Tao, 2017] You, S.; Xu, C.; and Tao, D. 2017. Learning from multiple teacher networks. In *KDD*.
- [Zafir et al., 2019] Zafir, O.; Boudoukh, G.; Izsak, P.; and Wasserblat, M. 2019. Q8bert: Quantized 8bit bert. In *5th Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition*.
- [Zagoruyko and Komodakis, 2016] Zagoruyko, S. and Komodakis, N. 2016. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*.