

Technical Summary of An Improved Attention for Visual Question Answering

Part I: Problem Formulation, Methods, and Evidence

Kai-Yu Lu

2025/09/20

1 Research Problem and Motivation

Task. Visual Question Answering (VQA) requires producing an accurate natural-language answer given an image and a free-form question. The task demands joint reasoning over visual content and text, with answers spanning categories such as Yes/No, Number, and Other.

Motivation. Attention mechanisms are central to strong VQA systems. Dense co-attention improves cross-modal interactions between question tokens and image regions yet always returns a weighted combination of values, even when relevant context is absent. This may inject noise into downstream reasoning. The paper augments attention with a mechanism that can gate out irrelevant outputs.

2 Related Work

Prior VQA architectures stack *self-attention* for intra-modal reasoning and *guided attention* for cross-modal reasoning, with MCAN as a representative design. Attention-on-Attention (AoA) has been explored in captioning to assess compatibility between the raw attention output and the query. Multimodal fusion spans linear fusion and bilinear pooling such as MUTAN. The present work integrates these ideas into a gated, modular co-attention-on-attention architecture with principled fusion.

3 Dataset Construction

Benchmark. VQA-v2 combines MS COCO images with human-authored questions and ten human answers per question. Splits include train, validation, and test, where the test split is partitioned into *test-dev* and *test-std* for online evaluation.

4 Query Protocol and Task Definitions

Input and output. The input is an image and a tokenized question. The output is a single answer from a fixed answer vocabulary. Training adopts a Binary Cross-Entropy objective over multi-label answer scores. Evaluation reports overall accuracy and per-type accuracy for Yes/No, Number, and Other according to the official VQA-v2 protocol.

5 Modeling Approach

Terminology and Abbreviations

Self Attention on Attention (SAoA). Intra-modal self-attention augmented with an AoA gate.

Guided Attention on Attention (GAoA). Cross-modal attention augmented with an AoA gate where one modality guides the other.

Modular Co-Attention on Attention (MCAoA). A layer that composes SAoA for both modalities with GAoA. Stacking yields the overall model.

MUTAN. A bilinear multimodal fusion operator.

Backbone Feature Extraction

Visual features. Object-centric region features are extracted using Faster R-CNN with a ResNet backbone. The features of detected regions form the image feature matrix X .

Text features. Question tokens are embedded with GloVe vectors and encoded by an LSTM to produce a sequence feature matrix Y .

Core Attention and Attention-on-Attention

Scaled dot-product attention. For a single head,

$$f_{\text{att}}(Q, K, V) = \text{Softmax}\left(\frac{QK}{\sqrt{d}}\right) V \quad (1)$$

Symbols. Q is the query matrix, K the key matrix, V the value matrix, and d the key dimension.

Intuition. The operator computes similarity between queries and keys, normalizes to attention weights, and aggregates values accordingly.

Role. This is the base attention used in both self- and guided-attention units before gating.

Attention on Attention (AoA) gate. Let $V' = f_{\text{att}}(Q, K, V)$ denote the attention output. AoA produces an information vector I and an attention gate G by

$$I = W_Q Q + W_{V'} V' + b_I, \quad (2)$$

$$G = \sigma(W_G Q + W_{G'} V' + b_G). \quad (3)$$

Symbols. $W_Q, W_{V'}, W_G, W_{G'} \in \mathbb{R}^{d \times d}$ are learned projections; $b_I, b_G \in \mathbb{R}^d$ are biases; $\sigma(\cdot)$ is the logistic sigmoid; d is the feature dimension.

Intuition. The gate G estimates compatibility between the attention output and the query context. The elementwise product $I \odot G$ preserves information aligned with the query and suppresses irrelevant components.

Role. This addresses the dense-attention failure mode where a head is forced to return a misleading weighted average when no relevant context exists.

SAoA and GAoA blocks. Each block applies multi-head attention, then AoA gating, followed by a positionwise feed-forward network $\text{FC}(4d) \rightarrow \text{ReLU} \rightarrow \text{Dropout}(0.1) \rightarrow \text{FC}(d)$ with residual connections and LayerNorm. SAoA uses self-attention within a modality; GAoA uses cross-attention where one modality guides the other.

MCAoA layer and stacking. A single MCAoA layer applies SAoA to both modalities, then GAoA in which question features guide image features. Stacking L layers yields refined outputs Y^L and X^L . Empirically, $L = 6$ provides the strongest validation performance.

Multimodal Fusion and Answer Head

From $X^L = [x_1, \dots, x_m]$ and $Y^L = [y_1, \dots, y_n]$, attended global vectors are computed via learned token and region weights:

$$X' = \sum_{i=1}^m \text{Softmax}(\text{MLP}(X_L)) x_i, \quad (4)$$

$$Y' = \sum_{i=1}^n \text{Softmax}(\text{MLP}(Y_L)) y_i. \quad (5)$$

Symbols. MLP is a two-layer perceptron $\text{FC}(d) \rightarrow \text{ReLU} \rightarrow \text{Dropout}(0.1) \rightarrow \text{FC}(1)$ producing scalar logits per token or region; Softmax normalizes them to weights.

Intuition. The equations compute modality-specific global summaries by softly selecting informative regions and tokens.

Role. These pooled vectors are fused for answer prediction. Two fusion variants are considered:

- **Multi-modal Attention Fusion.** Concatenate X' and Y' , estimate modality weights with a small classifier, form a weighted combination, and predict answers with a final classifier.
- **Multi-modal MUTAN Fusion.** Replace concatenation with MUTAN bilinear pooling, keeping the rest identical.

Hyperparameters and Training Protocol

Eight attention heads with feature dimension $d = 512$ are used, giving head dimension 64. Optimization uses Adam with $\beta_1 = 0.9$ and $\beta_2 = 0.98$, batch size 64, and approximately 13 epochs. A learning-rate schedule warms up to 10^{-4} and decays thereafter. The answer vocabulary size is approximately 3129. Training adopts Binary Cross-Entropy.

6 Empirical Results

Ablations

Depth L . Validation accuracy increases from $L = 2$ to $L = 6$ and saturates or drops at $L = 8$. The best validation accuracy occurs at $L = 6$.

Table 1: Validation accuracy on VQA-v2 by MCAoA layer depth L .

L	All	Other	Yes/No	Number
2	81.88	74.47	96.11	69.00
4	83.34	76.48	96.65	71.00
6	83.45	76.45	96.83	71.44
8	82.20	75.42	95.87	68.53

Table 2: Validation accuracy: baseline vs proposed components.

Method	All	Other	Yes/No	Number
MCAN	81.20	73.73	95.86	67.30
MCAoAN	82.91	75.92	96.47	70.38
MCAoAN + MUTAN Fusion	83.00	76.13	96.36	70.42
MCAoAN + Attention Fusion	83.25	76.51	96.58	70.40

Test-Set Comparisons

On *test-dev* and *test-std*, the proposed model outperforms competitive baselines including Bottom-Up, BAN, MuRel, and MCAN, improving the All metric while maintaining strong Yes/No and Number accuracy.

Table 3: Test-dev accuracy comparison.

Method	All	Other	Yes/No	Number
Bottom-Up	65.32	56.05	81.82	44.21
MFH	68.76	59.89	84.27	49.56
BAN	69.52	60.26	85.31	50.93
BAN + Counter	70.04	60.52	85.42	54.04
MuRel	68.03	57.85	84.77	49.84
MCAN	70.63	60.72	86.82	53.26
MCAoA (proposed)	70.90	60.97	87.05	53.81

Table 4: Test-std accuracy comparison.

Method	All	Other	Yes/No	Number
Bottom-Up	65.67	56.26	82.20	43.90
BAN + Counter	70.35	sub-metrics not reported		
MuRel	68.41	sub-metrics not reported		
MCAN	70.90	sub-metrics not reported		
MCAoA (proposed)	71.14	61.18	87.25	53.36

7 Summary

Contributions. The work introduces AoA-gated self and guided attention blocks that explicitly evaluate compatibility between attention outputs and queries, mitigating noise from irrelevant contexts. Stacked MCAoA layers paired with an effective fusion head deliver consistent gains on VQA-v2 over a strong MCAN baseline.

Evidence. Improvements hold across validation and both test splits, with ablations demonstrating benefits from depth and fusion.

Methodological Strengths

- **Principled gating of attention outputs.** AoA introduces an explicit compatibility gate to reject uninformative attention outputs.
- **Clear modular design.** SAoA and GAOA decompose intra- and inter-modal reasoning, enabling interpretable stacking in MCAoA layers.
- **Comprehensive ablations.** Layer-depth sweeps and component-wise comparisons quantify contributions and support design choices.
- **Dual fusion heads.** Both attention-based and MUTAN-based fusion are evaluated, improving robustness to modality weighting.
- **Use of standard features and splits.** Faster R-CNN features and VQA-v2 evaluation facilitate reproducibility and comparability.

Key Limitations and Future Directions

- **Task scope.** Evaluation focuses on VQA-v2. Transfer to other VQA datasets or broader vision–language tasks is not reported.
- **Computation and efficiency.** Hardware details and efficiency analyses are limited. Scaling behavior with respect to region counts and sequence length is not studied.
- **Answer space.** A fixed answer vocabulary is used. Handling free-form generative answers is not considered.
- **Error modes.** The paper includes qualitative failure cases but lacks a systematic bias analysis. Future work could relate AoA gating scores to uncertainty estimation and abstention.

Structured Technical Details

- **Datasets.** VQA-v2 with ten human answers per question; standard train, validation, test-dev, and test-std splits. Tokenization with truncation for long questions.
- **Model architecture.** Faster R-CNN features for regions; GloVe-embedded LSTM for questions; stacked MCAoA layers; fusion via attention or MUTAN; final sigmoid classifier over an answer vocabulary of about 3129.
- **Hyperparameters.** Adam ($\beta_1 = 0.9, \beta_2 = 0.98$), batch size 64, approximately 13 epochs, warmup then decay schedule, $N_{\text{heads}} = 8, d = 512$.
- **Evaluation metrics.** Overall and per-type accuracy for Yes/No, Number, Other following the VQA-v2 protocol.
- **Baselines.** Bottom-Up, MFH, BAN, BAN+Counter, MuRel, MCAN; comparisons on validation, test-dev, and test-std.
- **Environment.** Implementation based on common PyTorch stacks. GPU type and exact software versions are not specified.