APPENDIX

A ARCHITECTURAL GENERALIZATION

A.1 EXTENDED RESULTS ON DIFFERENT ARCHITECTURES

Table 6: RAH-LoRA performance across different MLLM architectures. Results show consistent improvements except for architectures with bottleneck designs.

Model	Size	VQAv2	TextVQA	GQA	POPE	Avg Δ	
Standard Multi-Head Attention							
LLaVA-1.5	7B	+1.3	+1.4	+1.1	+1.2	+1.25	
LLaVA-1.5	13B	+1.1	+1.2	+0.9	+1.3	+1.13	
Qwen-VL-Chat	7B	+1.2	+1.3	+1.0	+1.1	+1.15	
VILA	7B	+1.0	+1.1	+0.8	+1.0	+0.98	

B DOMAIN MISMATCH ANALYSIS

B.1 Cross-Domain Calibration

Table 7: Performance degradation with domain mismatch between calibration and evaluation data.

Eval Calib	VQA	TextVQA	GQA	SciQA
VQA	1.35	1.12	1.08	0.82
TextVQA	1.05	1.42	0.95	0.73
GQA	1.10	0.98	1.31	0.85
SciQA	0.75	0.68	0.71	1.48
CC3M	0.88	0.82	0.85	0.90

C SPECTRAL ANALYSIS OF WEIGHT UPDATES

C.1 SINGULAR VALUE ANALYSIS

Table 8: Singular value decay and variance captured across layers.

Layer Group	Decay Rate	Var @ r=4	Var @ r=8	Var @ r=16	Effective Rank
Early (0-7)	0.68	68%	82%	91%	5.2
Middle (8-15)	0.71	71%	85%	93%	6.1
Late (16-23)	0.73	73%	87%	94%	6.8
Deep (24-31)	0.75	75%	89%	95%	7.3

C.2 ANALYSIS

Exponential decay rates (0.68-0.75) confirm the low-rank structure of beneficial updates. Deeper layers show slightly higher effective rank, suggesting more complex cross-modal patterns. The 87% variance captured at r=8 (late layers) validates our default rank choice, balancing expressiveness and regularization.

Table 9: Concentration of improvements across calibrated heads.

Head Percentile	# Heads	Cum. Gain	Avg α	Avg $ \Delta W _F$
Top 20%	8	65%	0.124	0.218
20-40%	8	82%	0.091	0.156
40-60%	8	91%	0.073	0.112
60-80%	8	97%	0.058	0.089
Bottom 20%	8	100%	0.042	0.065

D PER-HEAD CONTRIBUTION ANALYSIS

D.1 CONTRIBUTION STATISTICS

D.2 ANALYSIS

Pareto principle confirmed: top 20% of heads contribute 65% of gains with $3\times$ larger updates ($\alpha=0.124~{\rm vs}~0.042$). These high-impact heads cluster in layers 14-20, corresponding to peak cross-modal pattern formation. The correlation between update magnitude and contribution (r=0.78) suggests our selection criteria effectively identify bottlenecks.

E FAILURE CASE ANALYSIS

E.1 QUALITATIVE FAILURE EXAMPLES

Figure 4 shows representative success and failure cases. Success cases (left) demonstrate improved focus on question-relevant regions, particularly for OCR and spatial tasks. Failure cases (right) typically involve counting or pure language reasoning where dispersed attention may be beneficial.

F IMPLEMENTATION DETAILS

F.1 CALIBRATION DATA SELECTION

- Random sampling from training split with fixed seed (42) - Balanced sampling across question types when available - Context length: maximum 256 tokens for questions - Image resolution: standard model preprocessing



Figure 4

F.2 ALGORITHM PSEUDOCODE

Algorithm 1 RAH-LoRA: Representative Anchor Head Low-Rank Adaptation

Require: Model \mathcal{M} , unlabeled calibration data \mathcal{D}_{cal}

Ensure: Calibrated model \mathcal{M}'

- 1: // Step 1: Profile attention patterns
- 2: Compute I-SAL scores for all heads using \mathcal{D}_{cal}
- 3: Compute CAF scores via gradient-based importance
- 4: // Step 2: Select calibration targets
- 5: **for** each layer *l* **do**
- 6: Identify heads with low I-SAL (bottom 10-15%)
- 7: Filter by CAF to exclude critical auxiliary heads
- 8: Add surviving heads to target set TH
- 9: end for

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

- 10: // Step 3: Calibrate each target head
- 11: **for** each target head $(l, h) \in \mathcal{TH}$ **do**
- 12: Find high-performing anchor heads in layer *l*
- 13: Construct RAH via weighted aggregation of anchors
- 14: Compute low-rank approximation of difference
- 15: Apply update with trust-region bounded step size
- 16: **end for**
- 17: return \mathcal{M}'

G ADDITIONAL VALIDATION EXPERIMENTS

H IMPACT OF LAYER RANGE SELECTION

H.1 LAYER RANGE ANALYSIS

Table 10: Calibration statistics and performance across different layer ranges (LLaVA-1.5-7B).

Layer Range	Layers	Targets	Calibrated	Avg Δ	Time
Early (0-7) Middle (8-15) Late (16-23) Deep (24-31)	8 8 8	42 65 71 38	8 18 23 12	+0.65 +0.92 +0.78 +0.41	1.2m 1.8m 2.1m 1.5m
Optimal (0-15) Default (12-23) Full (0-31)	16 12 32	107 124 284	26 38 72	+1.48 +1.35 +1.51	3.0m 3.2m 8.5m

H.2 KEY FINDINGS

The optimal range (layers 0-15) achieves the best performance (+1.48%) by capturing both early visual feature extraction and initial cross-modal integration patterns. This range provides:

- Early layers (0-7): Foundation visual processing that benefits from alignment
- Middle layers (8-15): Critical cross-modal pattern formation where most coordination failures occur
- Efficiency gain: 98% of full-model performance with only 50% of layers

While the full model achieves marginally higher gains (+1.51%), the 0-15 range offers the best efficiency-performance trade-off with 2.8× faster calibration. For deployment, we recommend layers 0-15 as the optimal configuration.

H.3 KEY FINDINGS

The default range (layers 12-23) captures 73% of problematic heads while achieving 95% of full-model gains. This sweet spot emerges because: - Early layers (0-7): Primarily unimodal processing, few targets - Middle-late layers (12-23): Peak cross-modal pattern formation with coordination failures - Deep layers (24-31): Most heads already well-aligned

Full model calibration yields marginal gains (+0.07%) for 2.6× computation, confirming diminishing returns. For deployment, we recommend the default 12-23 range as optimal balance between coverage and efficiency.