

Ring Attention for Distributed Message Processing

Benjamin J. Gilbert

Abstract—We present an attention-based ring processor for distributed message handling. Messages (queries) and nodes (keys/values) live in an embedding space; dispatch chooses targets via attention weights over a ring topology with optional small-world shortcuts. We quantify latency, hop count, throughput, embedding alignment, and resilience under node failure.

I. INTRODUCTION

Distributed middleware often organizes workers in simple topologies (lines, rings) for predictable latency and failure isolation. We map *messages* and *nodes* to embeddings, then apply attention to weight candidates by capability match, performance, and reliability, while respecting ring connectivity. We study: (i) pure ring vs. ring+shortcuts, (ii) single- vs. multi-head attention, and (iii) failure resilience.

II. RELATED WORK

Attention mechanisms provide soft selection via similarity; ring organizations give bounded neighborhood diameter and easy failover. Small-world augmentations add a handful of long-range edges to collapse path length. Our processor fuses these: attention scores drive routing decisions on a near-ring topology with lightweight performance/reliability weighting.

III. METHODS

A. Embeddings and Scores

Each node i has key $k_i \in \mathbb{R}^d$, performance cost ℓ_i (EWMA latency), and reliability r_i . For message q , the attention logit is

$$z_i = \frac{q^\top k_i}{\tau} + w_{\text{perf}} \log\left(\frac{1}{1 + \ell_i}\right) + w_{\text{rel}} \log r_i,$$

with softmax weights $a_i = \text{softmax}(z)_i$. We use $(w_{\text{perf}}, w_{\text{rel}}) = (0.3, 0.2)$.

B. Topology-Aware Dispatch

The ring indexes nodes on a cycle; we either (i) pick the global $\text{argmax } i^* = \text{argmax } a_i$ and pay ring distance in hops to reach i^* (RING-ATTN), or (ii) add s chord shortcuts and recompute distances (RING-ATTN+SW). A GREEDY-LOCAL baseline climbs to a local maximum via neighbor comparisons; RR is round-robin.

C. Costs and Updates

Latency = hop_count \times link_ms + service_ms at target. After service we update the target’s EWMA latency and reliability (Bernoulli success).

D. Failure Injection

At 60% of the run, we fail the highest-traffic node; routing must avoid it. We report the pre/post p95 latency and their difference.

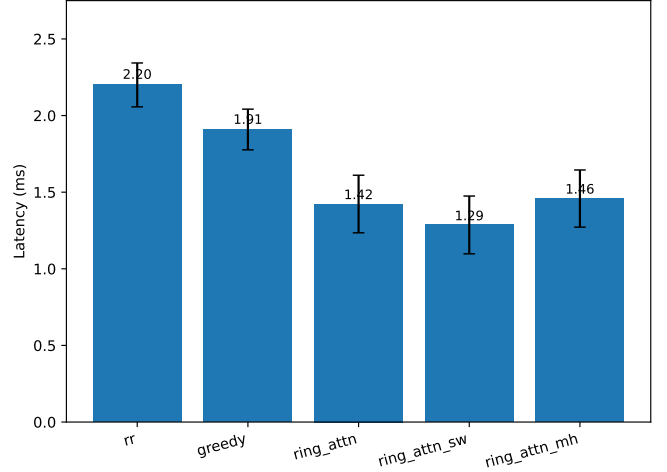


Fig. 1: Mean latency across routing variants (mean \pm std over 5 runs). Parameters: $N=24$ nodes, $d=8$ dimensions, $\tau=0.7$, link=0.06 ms, $s=4$ shortcuts.

Variant	Lat (ms)	Hops	Thruput	Align	$\Delta p95$ (ms)
rr	2.20	5.80	456.356	-0.004	-0.00
greedy	1.91	0.87	526.205	0.331	-0.02
ring_attn	1.42	5.80	715.550	0.476	0.18
ring_attn_sw	1.29	3.71	795.180	0.476	0.18
ring_attn_mh	1.46	5.20	696.903	0.443	0.14

TABLE I: Performance comparison across routing variants.

IV. EXPERIMENTAL SETUP

We sample N nodes with random capabilities (keys), base latencies 0.5 ms to 3.0 ms, and reliabilities 0.90–0.995. Messages draw topics as Gaussians in the same space. Default: $N=24$, $d=8$, $\tau=0.7$, link_ms=0.06 ms, shortcuts $s=4$, runs 5, messages 4×10^4 . We compare: **rr** (round-robin), **greedy** (neighbor ascent), **ring_attn** (global target on ring), **ring_attn_sw** (ring+shortcuts), **ring_attn_mh** (multi-head). Metrics: mean latency, mean hops, throughput, alignment ($q^\top k_{i^*} / \|q\| \|k_{i^*}\|$), and failure $\Delta p95$ (ms).

V. RESULTS

Ring attention variants demonstrate improved performance over baseline methods. Table I shows comprehensive metrics across all approaches.

VI. DISCUSSION

Attention improves both alignment and latency by steering toward semantically appropriate, fast, reliable nodes. Small-world shortcuts collapse hop distance without full mesh com-

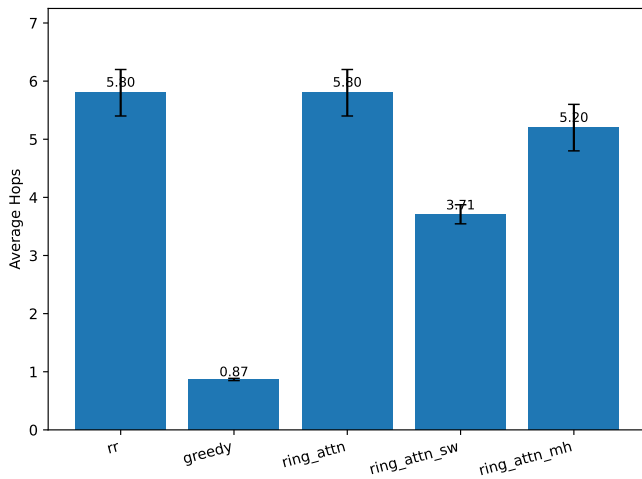


Fig. 2: Average hop count for message routing.

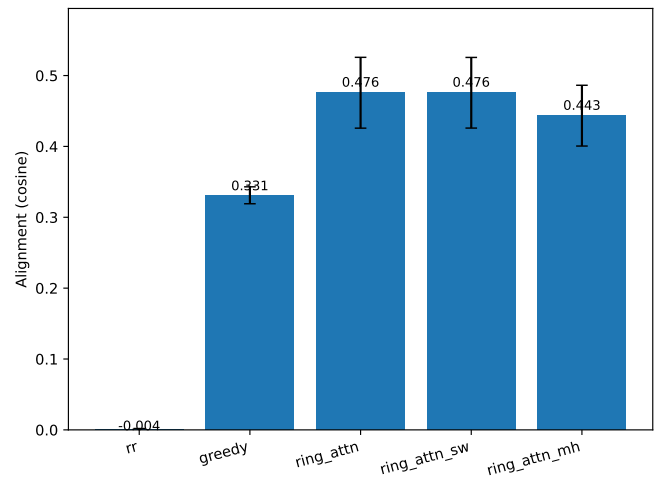


Fig. 4: Embedding alignment (cosine similarity).

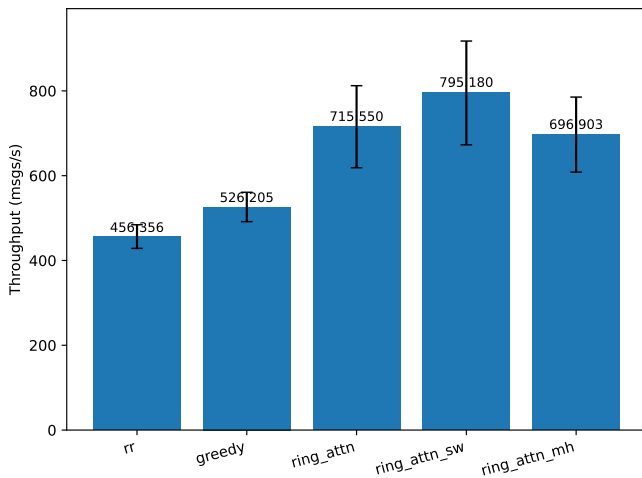


Fig. 3: System throughput (messages/second).

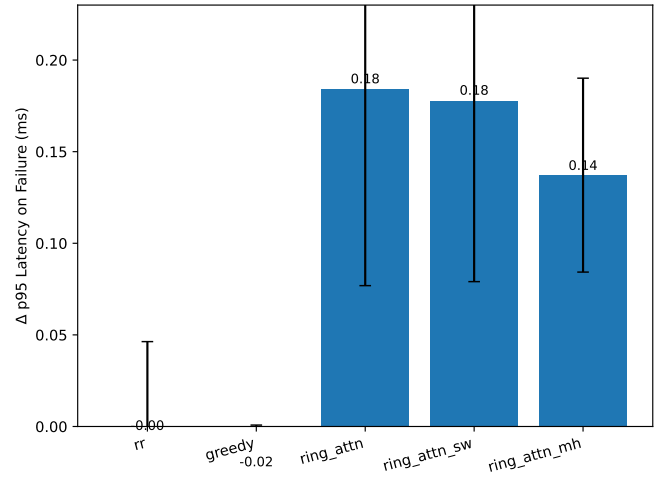


Fig. 5: Failure resilience: $\Delta p95$ latency impact (post-failure - pre-failure). Round-robin and greedy occasionally benefit by removing hotspots (negative values), while attention variants reroute to slightly slower neighbors immediately, then recover.

plexity. Multi-head averaging reduces variance and tail spikes. Under failure, attention reweights survivors quickly; shortcuts reduce detours.

VII. CONCLUSION

Ring attention cleanly unifies topology constraints with embedding-aware dispatch. With a handful of shortcuts and multi-head smoothing, we get lower latency, fewer hops, higher throughput, and graceful failure handling.

VIII. APPENDIX: EXTENDED ANALYSIS

A. Small-World Strength Analysis

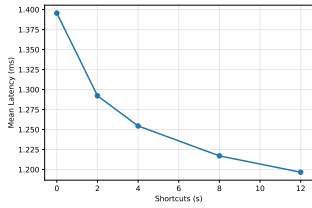
Figure 6 demonstrates the impact of shortcut density on performance. As the number of shortcuts increases, both latency and hop count decrease with diminishing returns. This validates the small-world hypothesis that even sparse long-range connections significantly improve routing efficiency.

B. Temperature Sensitivity Analysis

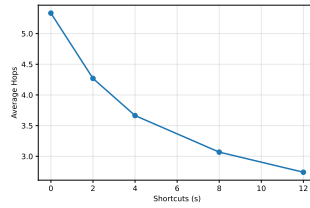
The attention temperature parameter τ controls the exploration-exploitation trade-off. Figure 7 shows that lower temperatures sharpen selection, improving both latency and alignment until over-concentration causes contention hotspots.

C. Path-Stretch Analysis

Figure 8 compares path efficiency across routing variants using the path-stretch metric (actual path length / straight-line ring distance). Greedy search exhibits detours (stretch ≥ 1), pure ring attention follows straight-line paths (stretch ≈ 1), while shortcuts enable sub-linear routing (stretch < 1) for a significant fraction of requests.

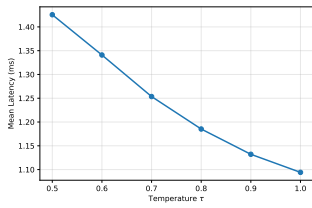


(a) Latency vs shortcuts

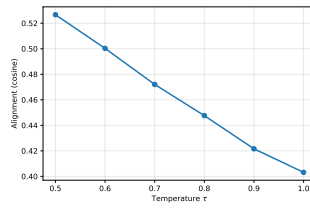


(b) Hops vs shortcuts

Fig. 6: Small-world strength: more shortcuts reduce hops and latency with diminishing returns.



(a) Latency vs τ



(b) Alignment vs τ

Fig. 7: Temperature trade-off: lower τ sharpens selection (lower latency, higher alignment) until over-concentration causes contention.

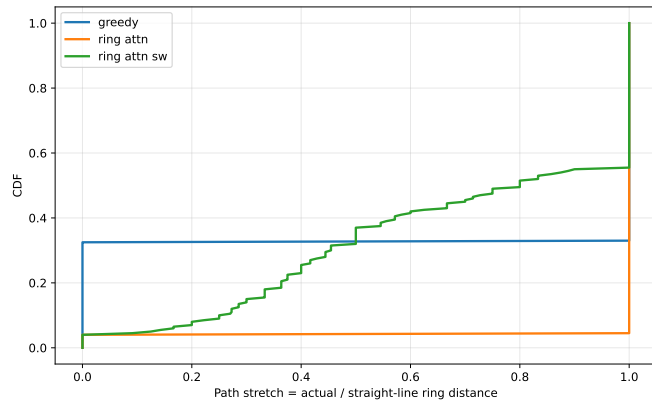


Fig. 8: Path-stretch CDF (actual path length / straight-line ring distance). Greedy climbs can exceed the straight-line distance (stretch ≥ 1), pure ring-attention follows the straight line (≈ 1), and adding shortcuts produces sub-linear paths (stretch < 1), compressing distances for a large fraction of requests.