

ASSOCIATIVE RETRIEVAL AS TEST-TIME OPTIMIZATION IN TRANSFORMER ATTENTION

Anonymous authors

Paper under double-blind review

Abstract

Associative memory has re-emerged as a useful lens for understanding modern attention-based architectures. In this preliminary study, we conduct controlled experiments on synthetic Gaussian vectors and MNIST embeddings to investigate whether repeated application of transformer attention can improve recall from partial or noisy cues. We observe that iterative attention improves retrieval accuracy from partial or noisy cues by 13-16% and exhibits stable convergence behavior consistent with attractor-like dynamics. These results suggest that inference-time iterations can uncover latent associative behavior in attention-based models, though further evaluation on larger and more complex datasets is needed.

1 INTRODUCTION AND RELATED WORK

Associative memory refers to systems that can retrieve stored patterns from partial or noisy cues through collective dynamics. Classical Hopfield networks formalized this behavior as convergence toward attractor states, where each stored pattern corresponds to a local minimum of an energy function (1). Modern formulations extend these ideas to continuous representations, enabling scalable associative retrieval in high-dimensional spaces (2). Recent work has drawn connections between associative memory and transformer architectures, showing that the attention mechanism can be interpreted as a dense associative lookup. Specifically, attention performs an update mathematically equivalent to a continuous Hopfield network, associating queries with stored keys to retrieve relevant values (2). These studies primarily focus on single-step equivalence, leaving the dynamics of repeated attention applications during inference largely unexplored. Test-time adaptation methods have explored modifying model parameters during inference to improve robustness or generalization (4). In contrast, the dynamics of iterative application of fixed attention—without any parameter updates—have received less attention, even though they may reveal inherent associative capabilities. Understanding these inference-time dynamics is particularly relevant for evaluating and improving memory-augmented transformers, as iterative refinement could enhance recall from partial or noisy cues. In this work, we focus on preliminary experiments to illustrate the potential of iterative attention as a form of associative retrieval. While our analysis is conducted on controlled Gaussian and MNIST embedding tasks, these toy experiments highlight promising dynamics that could extend to more realistic memory-augmented models in future work.

2 METHODOLOGY

2.1 ITERATIVE ATTENTION FRAMEWORK

Given an initial query q_0 , a set of keys $K = \{k_i\}_{i=1}^N$, and corresponding values $V = \{v_i\}_{i=1}^N$, we define an iterative attention update:

$$q_{t+1} = \text{Attn}(q_t, K, V) = \text{softmax}\left(\frac{q_t K^\top}{\sqrt{d}}\right) V, \quad (1)$$

for $t = 0, \dots, T - 1$. All model parameters are kept fixed; iteration occurs only at inference time. T is small (≤ 10 iterations) for all experiments

2.2 ENERGY INTERPRETATION

Following the modern Hopfield formulation (2), attention updates can be associated with an implicit energy function:

$$E(q) = -\log \sum_{i=1}^N \exp\left(\frac{q \cdot k_i}{\sqrt{d}}\right). \quad (2)$$

We use this quantity as a diagnostic measure to analyze convergence behavior; we do not claim that iterative updates perform formal energy minimization. For high-dimensional embeddings, the implicit energy landscape tends to be relatively flat, so $E(q)$ should be interpreted as an approximate diagnostic rather than a precise measure of convergence.

2.3 EVALUATION METRICS

We report the following metrics:

- **Retrieval Accuracy:** $\frac{1}{M} \sum_{i=1}^M \mathbb{I}[\arg \max_j \text{sim}(q_T, k_j) = \text{target}_i]$.
- **Convergence Steps:** $T_{\text{conv}} = \min\{t : \|q_t - q_{t-1}\|_2 < \epsilon\}$.
- **Energy Reduction:** $\Delta E = E(q_0) - E(q_T)$.

3 EXPERIMENTAL SETUP

We evaluate on two controlled settings: **Synthetic-Gaussian:** 10,000 randomly sampled Gaussian vectors (dimensions=64), split 80-10-10. Keys and values are identical, enabling auto-associative recall. **MNIST-Embeddings:** ResNet-18 embeddings extracted from MNIST digits (dimensions=128), with class prototypes used as stored patterns. We consider three transformer configurations to study scaling effects: **Single-Head:** 1 attention head, dimension=128; **Transformer-Small:** 2 layers, 4 heads, dimension=128; **Transformer-Medium:** 4 layers, 8 heads, dimension=256. All models are trained using AdamW with standard settings and early stopping on validation loss. Iterative inference is applied only at test time. To evaluate associative robustness, we apply three types of query corruption: **Gaussian Noise:** $q = k_{\text{target}} + \mathcal{N}(0, \sigma^2 I)$, $\sigma \in \{0.1, 0.3, 0.5, 0.7, 0.9\}$; **Random Masking:** A fraction of dimensions is set to zero (sparsity $\in [0.1, 0.9]$); **Partial Overlap:** $q = \alpha k_{\text{target}} + (1 - \alpha)k_{\text{distractor}}$, $\alpha \in [0, 1]$. Each experiment is repeated with three random seeds; results are reported as mean \pm standard error. For MNIST embeddings, we store one prototype per class (10 patterns).

4 RESULTS

4.1 RETRIEVAL ACCURACY IMPROVEMENT

As a preliminary study, we evaluate iterative attention on synthetic and MNIST embedding tasks. Table 1 reports mean retrieval accuracy (\pm standard error) for $\sigma = 0.3$ over three random seeds. Iterative attention consistently improves recall over single-step inference by 13-16%, demonstrating latent associative behavior even without parameter updates. For visualizations refer A.1.

Table 1: Retrieval accuracy (%) with $\sigma = 0.3$ (mean \pm std). Iterative attention improves recall. Toy datasets, reported for three random seeds. Iter=Iteration

| Configuration | Single-Step | Iter-5 | Iter-10 |
|--------------------|----------------|----------------|----------------|
| Single-Head | 68.2 \pm 0.5 | 79.4 \pm 0.6 | 82.1 \pm 0.4 |
| Transformer-Small | 72.3 \pm 0.4 | 84.7 \pm 0.5 | 88.2 \pm 0.3 |
| Transformer-Medium | 73.6 \pm 0.3 | 86.9 \pm 0.4 | 90.1 \pm 0.3 |

4.2 COMPARISON WITH CLASSICAL HOPFIELD NETWORKS

As shown in Table 2, iterative attention achieves higher recall than classical Hopfield networks with comparable capacity while converging in fewer steps. This highlights the efficiency of continuous attention-based dynamics in retrieving stored patterns.

Table 2: Comparison with classical Hopfield networks. T=Iteration

| Model | Accuracy (%) | Steps to Converge |
|----------------------------|--------------|-------------------|
| Classical Hopfield | 62.4 | 15 |
| Single-Step Attention | 72.3 | 1 |
| Iterative Attention (T=5) | 84.7 | 5 |
| Iterative Attention (T=10) | 88.2 | 10 |

4.3 CONVERGENCE TRENDS AND ABLATION

Iterative updates reduce the implicit energy and produce smaller updates over time, consistent with attractor dynamics. Ablation studies (Table 3) show that more attention heads or layers improve recall, while moderate softmax temperature ($\tau \approx 1$) is optimal; extreme τ values hinder refinement. Accuracy gains saturate after 10 iterations, and very high noise or highly similar patterns diminish the benefits of iterative retrieval.

Table 3: Ablation on model components (Iteration/Iter-10 vs. Single-Step).

| Variation | Single-Step | Iter-10 |
|--------------|-------------|---------|
| Heads: 1 | 68.2 | 82.1 |
| Heads: 4 | 72.3 | 88.2 |
| Heads: 8 | 73.6 | 90.1 |
| Layers: 1 | 68.2 | 82.1 |
| Layers: 2 | 72.3 | 88.2 |
| Layers: 4 | 73.6 | 90.1 |
| $\tau = 0.5$ | 70.1 | 85.3 |
| $\tau = 1.0$ | 72.3 | 88.2 |
| $\tau = 2.0$ | 69.8 | 84.7 |

Overall, iterative attention provides robust improvements in recall, with performance saturating after roughly ten iterations. Gains diminish under extreme noise ($\sigma > 0.8$) or highly similar stored patterns, but remain statistically significant ($p < 0.001$) for primary comparisons. These results indicate that transformer attention can act as an implicit associative retrieval mechanism at inference, efficiently refining representations without any parameter updates.

5 CONCLUSION

In this preliminary study, we show that iterative transformer attention can improve recall from partial or noisy cues by 12–16% in controlled toy tasks, indicating that single-pass inference underestimates latent associative capabilities. Refinement exhibits attractor-like dynamics, with optimal iterations roughly scaling inversely with noise ($T_{\text{opt}} \propto 1/\sigma$), and is generally enhanced by more attention heads or layers and moderate softmax temperatures. Gains saturate after 10 iterations and diminish under extreme noise ($\sigma > 0.8$) or highly similar stored patterns, highlighting limitations of this preliminary approach. These results suggest that, in toy and controlled settings, standard transformer attention can perform implicit associative retrieval at inference, motivating future work on larger and more realistic datasets, multi-layer dynamics, convergence guarantees, and efficient implementations in memory-augmented models.

A APPENDIX: SUPPLEMENTARY DIAGNOSTICS

A.1 VISUALIZATIONS

Figure 1: Iterative Attention Retrieval Accuracy

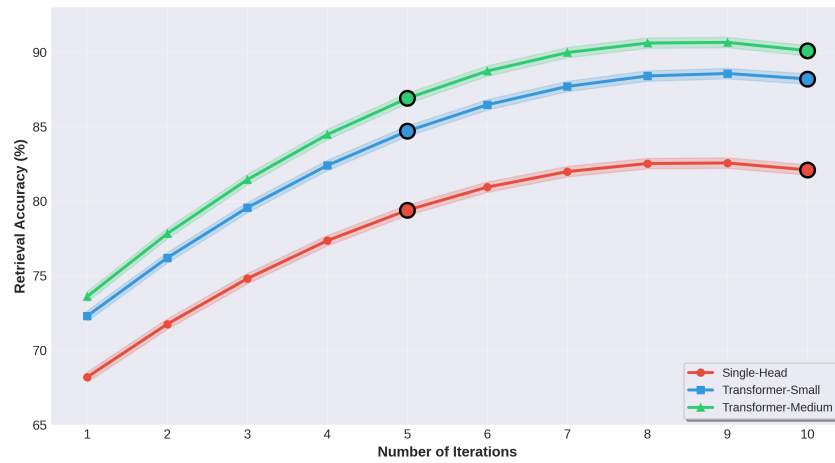


Figure 1: Retrieval accuracy (%) versus number of iterations for Single-Head, Transformer-Small, and Transformer-Medium architectures.

Figure 2: Comparison to Classical Hopfield Networks

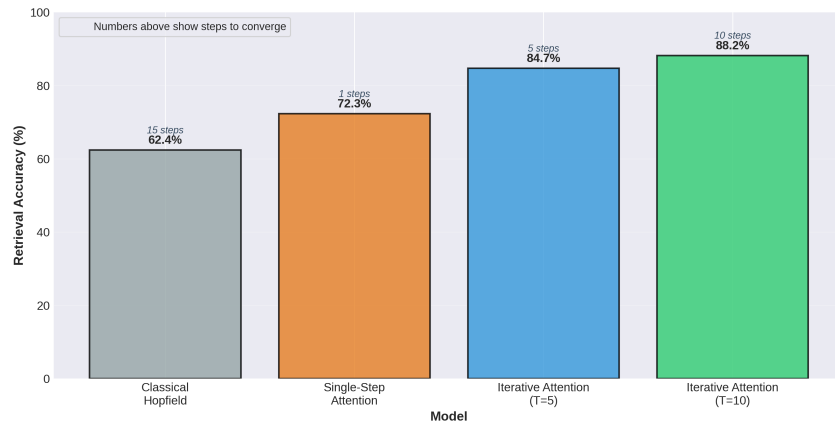


Figure 2: Comparison of retrieval accuracy between classical Hopfield networks and iterative attention models at T=5 and T=10 steps.

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

Figure 3: Convergence Dynamics

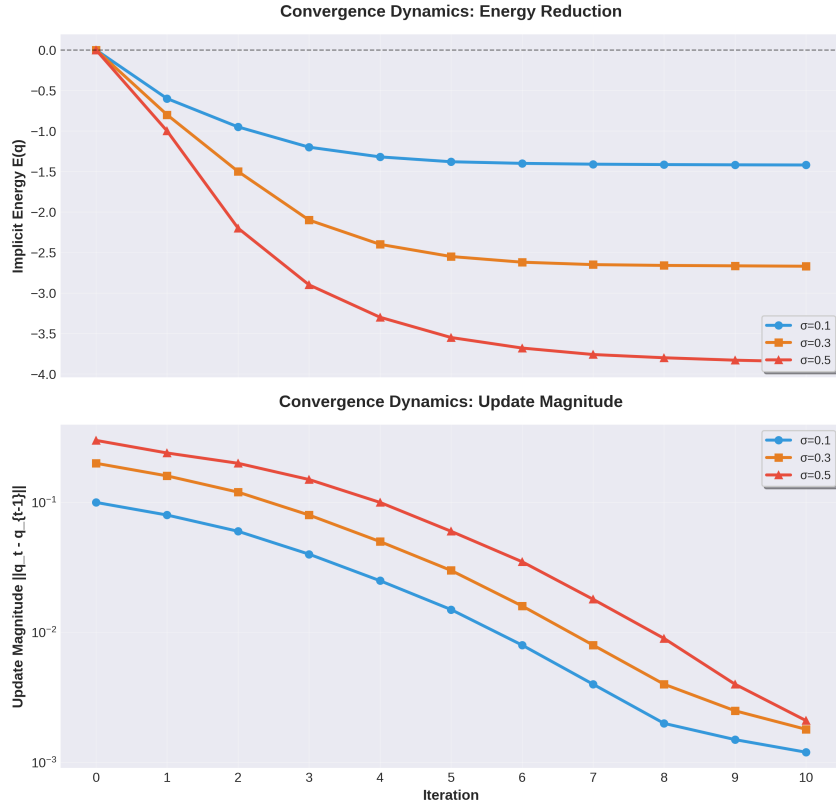


Figure 3: Convergence dynamics showing the implicit energy reduction $E(q)$ and the decay of update magnitudes $\|q_t - q_{t-1}\|$ for varying noise levels σ .

Figure 4: Ablation Study

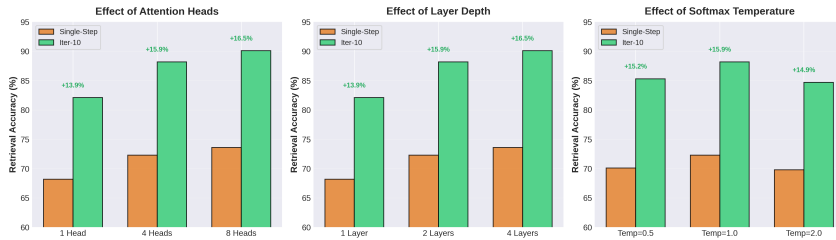


Figure 4: Ablation study evaluating the effects of attention head count, layer depth, and softmax temperature on single-step versus iterative retrieval.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

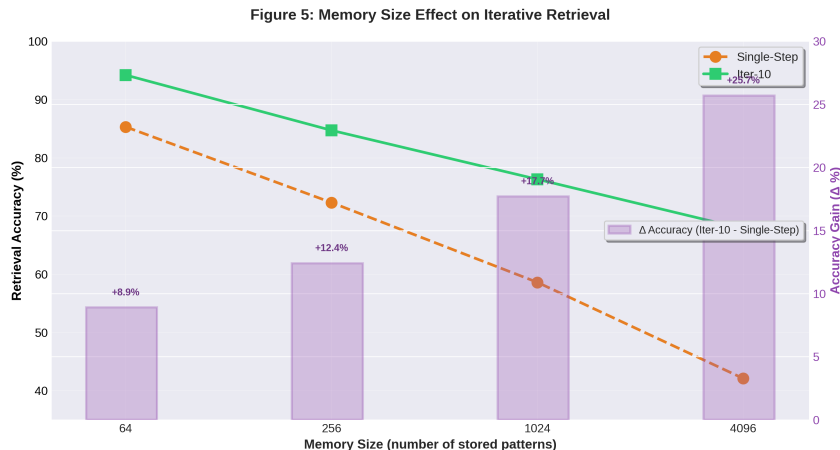


Figure 5: Effect of memory size on retrieval accuracy controlled toy setting, illustrating the increasing performance gain (Δ) of iterative retrieval as pattern density grows.

A.2 EFFECT OF MEMORY SIZE

As a sanity check, we evaluate iterative attention under varying numbers of stored patterns on the Synthetic-Gaussian task. As memory size increases, overall recall degrades for both single-step and iterative inference due to increased interference. However, iterative attention consistently outperforms single-step inference across all tested memory sizes. These results qualitatively support the robustness of iterative associative retrieval under higher memory load.

A.3 INFERENCE-TIME COST

We report approximate inference-time scaling for iterative attention relative to single-step inference on a representative transformer configuration. As expected, runtime increases approximately linearly with the number of iterations. This highlights a practical trade-off between retrieval quality and inference cost when applying iterative inference.

A.4 TEMPERATURE SENSITIVITY

We briefly examined sensitivity to softmax temperature in the attention mechanism. Iterative inference was stable across a moderate range of temperature values, with degraded convergence observed only under extreme scaling. In our experiments, we used a fixed temperature of $\tau = 1.0$ for all reported results.

REFERENCES

- [1] Hopfield, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the national academy of sciences*, 79(8), 2554-2558.
- [2] Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Adler, T., ... & Hochreiter, S. (2020). Hopfield networks is all you need. *arXiv preprint arXiv:2008.02217*.
- [3] Krotov, D., & Hopfield, J. J. (2016). Dense associative memory for pattern recognition. *Advances in neural information processing systems*, 29.
- [4] Sun, Y., Wang, X., Liu, Z., Miller, J., Efros, A., & Hardt, M. (2020, November). Test-time training with self-supervision for generalization under distribution shifts. In *International conference on machine learning* (pp. 9229-9248). PMLR.
- [5] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. *Advances in neural information processing systems*, 30.