

# Temporal-Aware Neural Retrieval: Exponential Decay Reranking for Time-Sensitive Document Search

Phantastic Labs Research<sup>1</sup>

<sup>1</sup>*Phantastic Labs*

(Dated: January 28, 2026)

We investigate methods for incorporating temporal relevance into neural embedding-based document retrieval. Using a corpus of 730 dated documents (real memory notes and synthetic entries spanning 6 months) and 60 queries (30 temporal-sensitive, 30 neutral), we compare four retrieval strategies: pure cosine similarity baseline, exponential time-decay reranking, temporal string injection into embeddings, and hybrid temporal pre-filtering. All methods use `all-MiniLM-L6-v2` (384-dimensional) sentence embeddings. We find that exponential time-decay with  $\lambda = 0.005$  yields the best overall performance (nDCG@10 = 0.822), improving temporal query nDCG@10 by 14.2% over the cosine baseline while maintaining neutral query performance ( $\Delta < 0.001$ ). Temporal injection into embedding text surprisingly *degrades* temporal retrieval. We validate embedding quality on BEIR SciFact (nDCG@10 = 0.645), confirming alignment with published benchmarks.

## INTRODUCTION

Information retrieval systems increasingly serve as memory layers for AI agents and personal knowledge management tools. In these settings, documents are naturally timestamped and queries often carry implicit or explicit temporal intent (“What happened this week?”, “Recent infrastructure changes”). Standard dense retrieval treats all documents equally regardless of recency, which is suboptimal for temporal queries.

We systematically evaluate four approaches to temporal-aware retrieval, measuring their effect on both temporal and non-temporal queries to understand the precision-recency tradeoff.

## METHODS

### Embedding Model

We use `all-MiniLM-L6-v2` from Sentence-Transformers [1], a 384-dimensional model producing L2-normalized embeddings. We validate on BEIR SciFact, obtaining nDCG@10 = 0.645, consistent with published results.

### Corpus

Our corpus contains 730 documents: 30 sections from real dated memory files (January 2026) and 700 synthetic documents spanning August 2025–January 2026 across five topic categories (infrastructure, ML research, product, team, personal).

## Queries and Ground Truth

We construct 60 queries: 30 temporal-sensitive (requiring recency or specific time-period awareness) and 30 neutral (topic-only). Ground truth relevance is determined by keyword matching combined with date-range constraints.

## Retrieval Methods

**Cosine Baseline.**  $s(q, d) = \cos(\mathbf{e}_q, \mathbf{e}_d)$

**Time Decay.**  $s(q, d) = \cos(\mathbf{e}_q, \mathbf{e}_d) \cdot \exp(-\lambda \cdot \text{age}(d))$  where  $\text{age}(d)$  is document age in days from the reference date. We sweep  $\lambda \in \{0, 0.001, 0.005, 0.01, 0.02, 0.05, 0.1, 0.2, 0.5\}$ .

**Temporal Injection.** Append “[Date: YYYY-MM-DD]” to document text before embedding, then rank by cosine similarity with the unmodified query.

**Hybrid Pre-filter.** Zero out scores for documents outside a time window  $W$ , then rank by cosine. We test  $W \in \{30, 60, 90, 180\}$  days.

## Metrics

We report nDCG@10, Recall@10, Recall@100, and MRR, computed separately for temporal, neutral, and all queries.

## RESULTS

### Overall Performance

Time decay with  $\lambda = 0.005$  achieves the best overall nDCG@10 (0.822), a 5.3% relative improvement over the cosine baseline.

TABLE I. Key method comparison across all 60 queries.

Method	nDCG@10	R@10	R@100	MRR
Cosine	0.781	0.447	0.752	0.754
Decay $\lambda=0.005$	<b>0.822</b>	<b>0.453</b>	<b>0.810</b>	<b>0.809</b>
Injection	0.763	0.443	0.757	0.728
Hybrid 60d	0.810	0.295	0.434	0.800

### Temporal vs. Neutral Split

TABLE II. Performance split by query type.

Method	Temp. nDCG@10	Neut. nDCG@10
Cosine	0.580	0.982
Decay $\lambda=0.005$	<b>0.662</b>	0.981
Injection	0.539	<b>0.986</b>
Hybrid 60d	0.637	0.983

The time-decay method improves temporal nDCG@10 by 14.2% relative ( $0.580 \rightarrow 0.662$ ) with negligible degradation on neutral queries ( $-0.001$ ). Notably, temporal injection *hurts* temporal retrieval ( $-7.0\%$ ), likely because date strings in embeddings create spurious similarity between temporally proximate but semantically unrelated documents.

### $\lambda$ Parameter Sweep

TABLE III. Effect of  $\lambda$  on nDCG@10 by query type.

$\lambda$	All	Temporal	Neutral
0.0	0.781	0.580	0.982
0.001	0.811	0.642	0.981
<b>0.005</b>	<b>0.822</b>	<b>0.662</b>	0.981
0.01	0.805	0.644	0.967
0.02	0.784	0.634	0.935
0.05	0.729	0.602	0.855
0.1	0.624	0.539	0.710
0.2	0.476	0.460	0.491
0.5	0.331	0.356	0.306

The optimal  $\lambda = 0.005$  corresponds to a half-life of  $\ln 2 / 0.005 \approx 139$  days. For  $\lambda > 0.01$ , neutral query performance degrades significantly, indicating over-penalization of older but topically relevant documents. The crossover point where temporal queries outperform neutral occurs at  $\lambda \approx 0.2$ , far past the useful range.

### Hybrid Pre-filter Windows

Window size 90 days achieves the best temporal nDCG@10 (0.641) but severely limits Recall@100 (0.317 for 60d window vs. 0.752 for cosine). The hard cutoff nature of pre-filtering makes it unsuitable when temporal relevance is uncertain.

### Latency

All methods achieve sub-millisecond mean latency (0.06–0.48ms) for scoring 730 documents, confirming that temporal reranking adds negligible overhead to dense retrieval.

### BEIR SciFact Validation

Our all-MiniLM-L6-v2 implementation achieves nDCG@10 = 0.645 on BEIR SciFact (5,183 documents, 300 test queries), consistent with published MTEB leaderboard results, validating our embedding and evaluation pipeline.

### DISCUSSION

The exponential time-decay method emerges as the clear winner for temporal-aware retrieval. Its key advantages are:

- 1. Simplicity:** A single multiplicative factor with one tunable parameter.
- 2. Graceful degradation:** At  $\lambda = 0$ , it reduces to cosine baseline. Small  $\lambda$  provides temporal boost without harming topical relevance.
- 3. No re-embedding:** Unlike temporal injection, it reuses existing document embeddings.

The failure of temporal injection is instructive: embedding date strings creates a “temporal proximity” signal that interferes with semantic similarity. Documents from the same week become artificially similar regardless of topic.

For production systems, we recommend  $\lambda \in [0.001, 0.01]$  depending on the desired recency bias, with  $\lambda = 0.005$  as a strong default. Query-dependent  $\lambda$  selection (detecting temporal intent) is a promising direction for future work.

### CONCLUSION

Exponential time-decay reranking with  $\lambda = 0.005$  provides a 14% improvement in temporal query nDCG@10

over pure cosine similarity with negligible impact on non-temporal queries. This simple, parameter-efficient method outperforms both temporal string injection and hard time-window pre-filtering, making it the recommended approach for temporally-aware neural retrieval.

- 
- [1] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,” in *Proc. EMNLP*, 2019.