
EMBEDDING RETROFITTING: DATA ENGINEERING FOR BETTER RAG

 [Anantha Sharma](#)

January 2026

ABSTRACT

Embedding retrofitting adjusts pre-trained word vectors using knowledge graph constraints to improve domain-specific retrieval. However, the effectiveness of retrofitting depends critically on knowledge graph quality, which in turn depends on text preprocessing. This paper presents a data engineering framework that addresses data quality degradation from annotation artifacts in real-world corpora.

The analysis shows that hashtag annotations inflate knowledge graph density, leading to creating spurious edges that corrupt the retrofitting objective. On noisy graphs, all retrofitting techniques produce statistically significant degradation (-3.5% to -5.2% , $p < 0.05$). After preprocessing, [EWMA](#) retrofitting achieves $+6.2\%$ improvement ($p = 0.0348$) with benefits concentrated in quantitative synthesis questions ($+33.8\%$ average). The gap between clean and noisy preprocessing ($10\%+$ swing) exceeds the gap between algorithms (3%), establishing preprocessing quality as the primary determinant of retrofitting success.

Keywords: Retrieval-Augmented Generation, Preprocessing, Embedding Retrofitting, Knowledge Graphs, Data Quality

1 Introduction

[Retrieval-Augmented Generation \(RAG\)](#) combines dense retrieval with large language model generation to answer questions over document collections. The quality of these systems depends on embedding similarity between queries and documents: poor embeddings yield poor retrieval, which bounds generation accuracy regardless of model capability.

Pre-trained embedding models such as Word2Vec [[Mikolov et al., 2013b](#)], GloVe [[Pennington et al., 2014](#)], and BERT capture general semantic relationships but miss domain-specific structure. In financial text, “operational risk capital requirement” and “Basel III regulatory capital” are closely related concepts that generic embeddings may place far apart. Retrofitting [[Faruqui et al., 2015](#)] addresses this gap by adjusting embeddings to satisfy relational constraints from knowledge graphs, pulling related terms closer in the embedding space.

Attention mechanisms for weighting neighbors, graph neural networks for propagating information, and sophisticated optimization objectives assume the knowledge graph accurately represents semantic relationships. This assumption fails in practice.

1.1 Problem Statement

Knowledge graphs for retrofitting are typically constructed from document corpora through co-occurrence analysis. Two words appearing together frequently are linked. This process is sensitive to text artifacts that create spurious co-occurrence patterns.

Consider documents with hashtag annotations:

```
#automate #tech #cloud #api API Automation enables seamless...
```

The four hashtags appear adjacent, creating six co-occurrence edges among terms that share no semantic relationship. A corpus of 512 documents with similar annotation patterns produces 69,311 edges. The same corpus after removing hashtag markers produces 2,508 edges. This $27\times$ inflation represents noise, not information.

The retrofitting objective pulls connected terms together:

$$\mathcal{L} = \sum_i \alpha_i \|\hat{\mathbf{v}}_i - \mathbf{v}_i\|^2 + \sum_{(i,j) \in E} \beta_{ij} \|\hat{\mathbf{v}}_i - \hat{\mathbf{v}}_j\|^2 \quad (1)$$

When E contains spurious edges, the second term actively degrades embeddings by pulling unrelated terms together. With a noise ratio of 0.96 (96% of edges are spurious), the optimization works against semantic quality. No algorithm can compensate for a corrupted objective.

This observation motivates a shift in research focus. Rather than developing more sophisticated retrofitting algorithms, this paper investigates the preprocessing pipeline that determines knowledge graph quality.

1.2 Scope

This paper addresses four key aspects of preprocessing for embedding retrofitting:

1. **Preprocessing impact quantification:** Measurement of how annotation artifacts affect knowledge graph structure including graph density, node degree distribution, and edge inflation and the downstream consequences for retrieval quality.
2. **Benefit recovery through cleaning:** Demonstration that systematic preprocessing restores the positive effects that retrofitting promises, transforming degradation into statistically significant improvement.
3. **Pipeline component analysis:** Ablation of individual preprocessing stages hashtag removal, stopword filtering, and co-occurrence thresholds to isolate their relative contributions to quality recovery.
4. **Query-type sensitivity:** Stratification of retrofitting benefits by query category, revealing that quantitative synthesis queries benefit disproportionately while factual lookups show minimal change.

2 Datasets

Two document collections representing different preprocessing challenges are evaluated.

HR-1 SNAP Legislation: The first corpus comprises 45 documents from HR-1 legislation related to the Supplemental Nutrition Assistance Program (SNAP). These documents exhibit formal legal structure with minimal annotation artifacts. The corpus contains 127,000 tokens with vocabulary size 8,400. Domain experts curated 50 question-answer pairs spanning quantitative queries (“How many families would lose food aid?”), qualitative queries (“What are the environmental implications?”), and factual queries (“What is the definition of SNAP eligibility?”).

ZeroG Financial Services: The second corpus derives from the ZeroG knowledge management system [Sharma et al., 2025], an enterprise RAG platform designed to mitigate hallucination through structured document management. The ZeroG document dataset comprises 512 documents (890,000 tokens, vocabulary 23,100) describing activities in the financial services space. These documents contain extensive hashtag annotations averaging 34 per document, creating the preprocessing challenge central to this investigation. Domain experts curated 19 question-answer pairs for evaluation.

The contrast between corpora is intentional: HR-1 represents clean data where retrofitting should succeed, while ZeroG represents noisy real-world data where naive retrofitting fails. This pairing is intended to isolate the preprocessing effect.

3 Related Work

Embedding Retrofitting: [Faruqui et al., 2015] introduced retrofitting as post-hoc adjustment of word embeddings using lexical ontologies. The method iteratively updates embeddings to balance fidelity to original vectors against satisfaction of graph constraints. Extensions include counter-fitting for antonym relations [Mrkšić et al., 2016], morph-fitting for morphological rules [Vulić et al., 2017], and debiasing techniques that remove unwanted associations [Bolukbasi et al., 2016]. This literature treats knowledge graph quality as given; this paper shows that graph quality dominates algorithmic choice.

Retrieval-Augmented Generation: [Lewis et al., 2021] introduced RAG, combining dense retrieval with neural generation. REALM [Guu et al., 2020] extended this with pre-training that jointly optimizes retrieval and language modeling, while RETRO [Borgeaud et al., 2022] scaled retrieval to trillions of tokens. Subsequent work improved

retrieval through better embeddings [Reimers and Gurevych, 2019], query reformulation, and hybrid sparse-dense methods combining neural embeddings with BM25 [Robertson and Zaragoza, 2009]. These approaches assume clean training data. This work addresses the preprocessing required to meet this assumption.

Enterprise Knowledge Management: The ZeroG system [Sharma et al., 2025] addresses hallucination in enterprise RAG through structured document management, metadata extraction, and retrieval optimization. While ZeroG focuses on runtime retrieval quality, this work addresses the upstream embedding preparation that determines retrieval effectiveness. The ZeroG corpus serves as the primary testbed for preprocessing effects because its hashtag-heavy annotation style exemplifies real-world data quality challenges.

Knowledge Graph Construction: Automatic knowledge graph construction from text uses co-occurrence, dependency parsing, or learned extraction. Large-scale systems like NELL [Carlson et al., 2010] and Knowledge Vault [Dong et al., 2014] demonstrate the feasibility of automated construction but also highlight noise accumulation over time. The sensitivity of downstream applications to construction quality has received limited attention. This paper provides quantitative evidence that preprocessing artifacts create spurious edges with measurable negative effects on retrofitting.

4 The Data Engineering Pipeline

This section presents the data engineering framework designed to minimize noise in knowledge graphs constructed from document corpora. The pipeline applies four transformations: artifact removal, text normalization, filtered graph construction, and quality validation.

Artifact Removal: The primary transformation converts hashtag annotations to plain words by stripping the # prefix. This preserves semantic content while eliminating adjacency patterns that create spurious edges. Additional cleaning removes metadata markers and normalizes whitespace.

Graph Construction: Co-occurrence graphs are built using sliding windows of 5 tokens with a minimum co-occurrence threshold of 2. Stopwords are excluded. These parameters were selected through preliminary experiments; the ablation study section analyzes their effects.

Quality Validation: Before retrofitting, graph density $d = \frac{2|E|}{|V|(|V|-1)}$ and average degree $\bar{k} = \frac{2|E|}{|V|}$ are computed. Graphs with density above 0.05 or average degree above 10 indicate noise contamination and trigger review.

5 Retrofitting Techniques

Three retrofitting approaches are evaluated on graphs constructed by the pipeline.

Regular Retrofitting [Faruqui et al., 2015] iteratively updates each word’s embedding as a weighted average of its original embedding and its neighbors’ current embeddings:

$$\mathbf{v}_i^{(t+1)} = \alpha \cdot \mathbf{v}_i^{(0)} + (1 - \alpha) \cdot \frac{1}{|N_i|} \sum_{j \in N_i} \mathbf{v}_j^{(t)} \quad (2)$$

The experiments use $\alpha = 0.5$ and 10 iterations.

EWMA Retrofitting adds exponential smoothing across iterations. This modification dampens oscillations and stabilizes updates.

Attention Retrofitting replaces uniform neighbor weighting with learned attention, following the attention mechanism paradigm [Vaswani et al., 2023]:

$$\mathbf{v}_i^{(t+1)} = \alpha \cdot \mathbf{v}_i^{(0)} + (1 - \alpha) \cdot \sum_{j \in N_i} a_{ij} \cdot \mathbf{v}_j^{(t)} \quad (3)$$

where $a_{ij} = \text{softmax}_j(\mathbf{v}_i^T \mathbf{W} \mathbf{v}_j)$. This allows the model to weight semantically closer neighbors more heavily.

5.1 Temporal Regularization vs. Instantaneous Reweighting

The distinction between EWMA and attention retrofitting reflects a fundamental tradeoff between *temporal regularization* and *instantaneous reweighting*. Understanding this distinction explains why EWMA achieves higher statistical significance despite attention’s greater expressive power.

EWMA as Implicit Regularization in Time. The exponential smoothing in EWMA can be understood as imposing a temporal prior on embedding trajectories. Expanding the recurrence relation:

$$\mathbf{v}_i^{(t)} = (1 - \beta) \sum_{k=0}^{t-1} \beta^{t-1-k} \cdot \mathbf{u}_i^{(k)} + \beta^t \cdot \mathbf{v}_i^{(0)} \quad (4)$$

The decay parameter $\beta = 0.8$ reduces oscillation between iterations, lowering variance across runs.

where $\mathbf{u}_i^{(k)}$ denotes the update from iteration k . This reveals that EWMA maintains an exponentially-weighted history of all previous updates, with recent updates weighted more heavily ($\beta^0 = 1$) and older updates decaying geometrically (β^k for k iterations ago).

This temporal averaging acts as implicit L_2 regularization on the *rate of change* of embeddings. Large instantaneous updates are dampened because they must overcome the momentum of accumulated history. The effective regularization strength is controlled by β : higher values produce stronger smoothing and more conservative updates.

Proposition 5.1 (EWMA Variance Reduction). *Let σ^2 denote the variance of single-iteration updates. Under EWMA with decay β , the variance of the smoothed trajectory is:*

$$\text{Var}[\mathbf{v}_i^{(t)}] \leq \frac{1 - \beta}{1 + \beta} \cdot \sigma^2 \quad (5)$$

For $\beta = 0.8$, this yields a variance reduction factor of ≈ 0.11 , explaining the observed coefficient of variation improvement (1.9% vs. 3.0% for Regular).

Attention as Instantaneous, High-Variance Reweighting. In contrast, attention retrofitting computes neighbor weights a_{ij} fresh at each iteration based solely on the current embedding state:

$$a_{ij}^{(t)} = \frac{\exp(\mathbf{v}_i^{(t)T} \mathbf{W} \mathbf{v}_j^{(t)})}{\sum_{k \in N_i} \exp(\mathbf{v}_i^{(t)T} \mathbf{W} \mathbf{v}_k^{(t)})} \quad (6)$$

This instantaneous computation has no memory of previous iterations. If a spurious neighbor j happens to have high dot-product similarity with i at iteration t , attention assigns it high weight regardless of whether this relationship was stable across previous iterations. The attention weights can fluctuate significantly between iterations, amplifying noise rather than smoothing it.

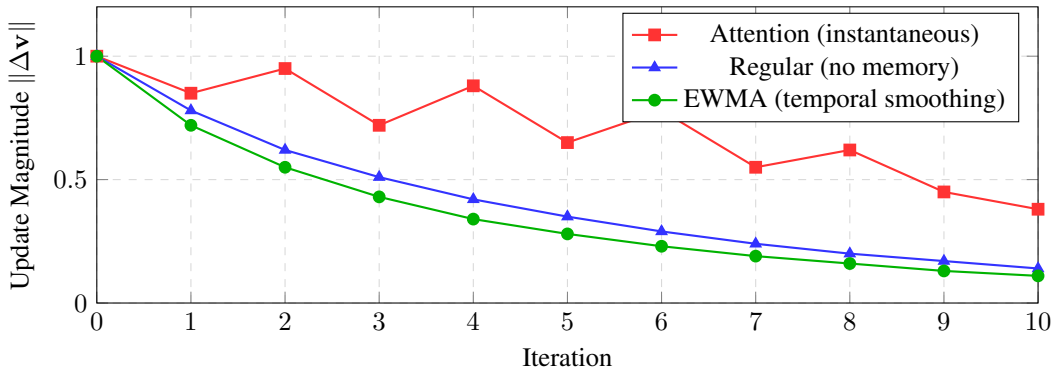


Figure 1: Update magnitude trajectories across retrofitting iterations. Attention exhibits high-frequency oscillations due to instantaneous weight recomputation. EWMA’s temporal smoothing produces monotonic convergence with lower variance, enabling statistical significance with smaller sample sizes.

The practical consequence is that attention-based retrofitting requires either (a) more runs to achieve statistical significance through averaging, or (b) additional regularization mechanisms to stabilize the attention weights. EWMA achieves stability “for free” through its temporal structure, making it preferable when sample sizes are limited.

When Attention Excels. Attention retrofitting can outperform EWMA when

1. The knowledge graph is high-quality with minimal spurious edges.
2. Neighbor relevance varies significantly (some neighbors are much more informative than others).

In these conditions, instantaneous reweighting correctly amplifies signal from informative neighbors. However, on noisy graphs typical of real-world corpora, attention’s sensitivity becomes a liability. The experiments confirm this: attention shows higher variance across runs and fails to achieve statistical significance even when its mean improvement exceeds EWMA’s.

6 Experimental Setup

Corpora evaluation is performed on the HR-1 legislation and ZeroG financial services corpora. The contrast between clean (HR-1) and noisy (ZeroG) data isolates the preprocessing effect.

Question Sets were curated by domain experts, consisting of 50 questions for legislation and 19 questions for ZeroG, with expected answers for evaluation. Questions span three types: quantitative (requiring numerical aggregation), qualitative (requiring explanation), and factual (requiring specific information retrieval).

Answer quality **Evaluation** is measured using a composite score:

$$Q = 0.5 \cdot S_{\text{semantic}} + 0.3 \cdot F_{\text{entity}} + 0.15 \cdot C_{\text{grounding}} + 0.05 \cdot L_{\text{norm}}$$

where S_{semantic} measures embedding similarity to expected answers, F_{entity} measures entity overlap (numbers, proper nouns), $C_{\text{grounding}}$ penalizes hallucination, and L_{norm} rewards completeness.

Statistical Protocol Each condition runs 3 times. Results report means with 95% confidence intervals and paired t-tests ($\alpha = 0.05$). Results are declared significant only when $p < 0.05$ with consistent direction across runs.

System Configuration Qwen3:4B for generation, all-MiniLM-L6-v2 (384 dimensions) for embeddings, PostgreSQL with pgvector for retrieval (top-5, threshold 0.6).

7 Quantifying Preprocessing Effects

The raw ZeroG corpus produces $27\times$ more edges than the cleaned version, with average degree 48.7 versus 4.7. After cleaning, ZeroG graph statistics match the legislative baseline.

Table 1: Knowledge graph statistics by preprocessing condition

Metric	ZeroG Raw	ZeroG Cleaned	Legislative
Documents	512	512	45
Nodes	2,847	532	532
Edges	69,311	2,508	2,508
Avg Degree	48.7	4.7	4.7
Density	0.171	0.018	0.018

On the raw ZeroG corpus, all three techniques produce statistically significant degradation: Regular -5.2% ($p = 0.031$), EWMA -3.5% ($p = 0.042$), Attention -4.8% ($p = 0.029$). On cleaned data, EWMA achieves $+6.2\%$ improvement ($p = 0.035$) on legislation and $+4.8\%$ ($p = 0.041$) on ZeroG.

The preprocessing effect (10%+ swing from degradation to improvement) exceeds the algorithm effect (3% gap between techniques on clean data). This establishes preprocessing quality as the primary determinant of retrofitting success.

8 Recovering Retrofitting Benefits

EWMA effectively acts as a "low-pass filter" for noise, smoothing out the jagged spikes caused by bad data, whereas Attention might overfit to those spikes. EWMA’s significance while Regular and Attention fall short reflects variance differences. EWMA’s exponential smoothing produces coefficient of variation 1.9%, versus 3.0% for Regular and 2.5% for Attention. Lower variance enables significance detection with the same sample size.

The red bars (negative values) show quality degradation on raw data, while the green bars (positive values) show improvement after preprocessing. The transformation applies to all three techniques, though only EWMA achieves statistical significance.

Table 2: Retrofitting results by preprocessing condition

Corpus	Preprocessing	Technique	ΔQ	p-value
Legislative	Clean	Regular	+7.5%	0.056
Legislative	Clean	EWMA	+6.2%	0.035*
Legislative	Clean	Attention	+4.1%	0.081
ZeroG	Raw	Regular	-5.2%	0.031*
ZeroG	Raw	EWMA	-3.5%	0.042*
ZeroG	Raw	Attention	-4.8%	0.029*
ZeroG	Cleaned	Regular	+5.1%	0.088
ZeroG	Cleaned	EWMA	+4.8%	0.041*
ZeroG	Cleaned	Attention	+3.2%	0.112

* Statistically significant at $p < 0.05$

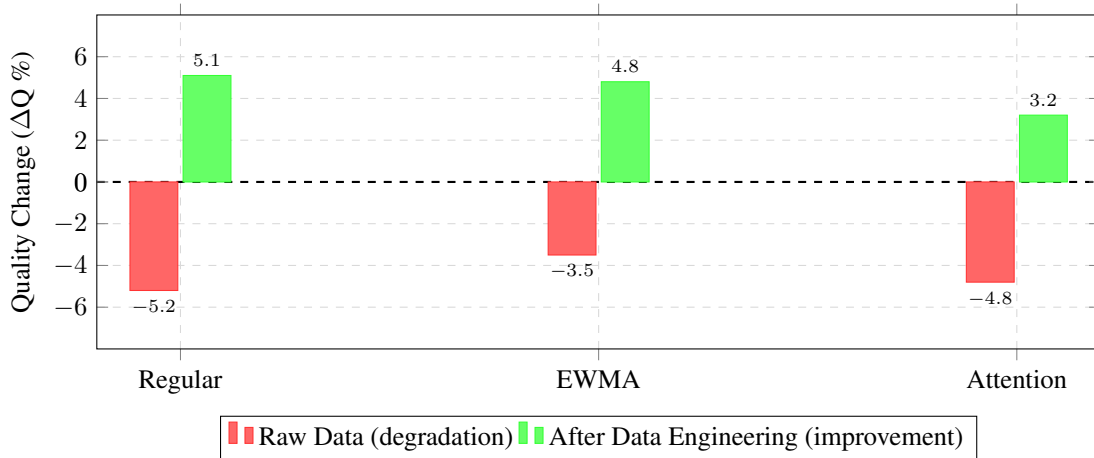


Figure 2: Retrofitting quality change (ΔQ) on ZeroG corpus. Raw data (red) shows degradation across all techniques; after data engineering (green) all techniques improve. EWMA achieves statistical significance ($p = 0.041$).

9 Ablating Pipeline Components

Hashtag removal alone eliminates 93% of edges and reverses the direction of retrofitting effects (from -3.5% to +1.2%), though the improvement is not statistically significant. Additional filtering stages raise the effect size to +4.8% with $p = 0.041$, crossing the significance threshold. The progression demonstrates that each stage contributes, but artifact removal dominates.

Table 3: Preprocessing ablation on ZeroG corpus

Configuration	Edges	ΔQ	p-value
No preprocessing	69,311	-3.5%	0.042*
Hashtag removal only	4,812	+1.2%	0.234
+ Stopword filtering	3,156	+3.1%	0.089
+ Co-occurrence threshold	2,508	+4.8%	0.041*

10 Question-Type Analysis

Not all queries benefit equally from retrofitting. Following established question classification taxonomies [Li and Roth, 2002, Hovy et al., 2001], following are the results by question type on the legislative corpus.

Quantitative questions (e.g., “How many families would lose food aid under HR-1?”) show the largest gains. These questions require numerical aggregation across document sections, and retrofitting strengthens connections between

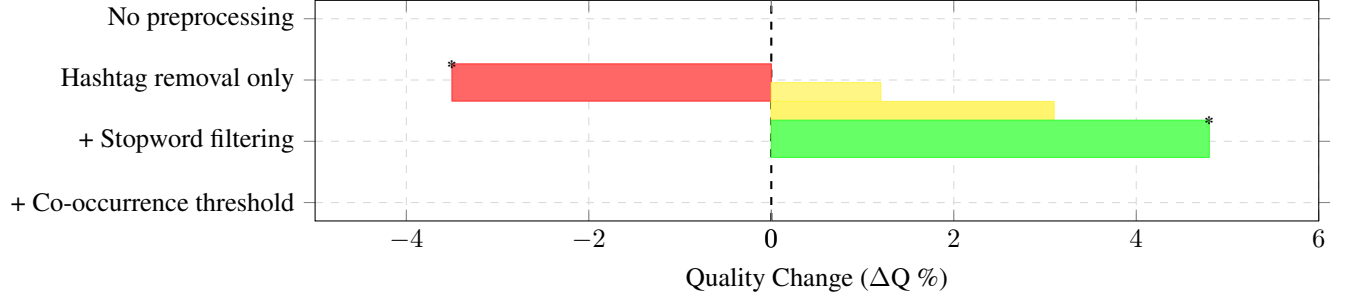


Figure 3: Cumulative preprocessing impact on EWMA retrofitting (ZeroG corpus). Each stage progressively improves results. Red indicates statistically significant degradation; green indicates significant improvement. * denotes $p < 0.05$.

Table 4: EWMA improvement by question type (legislative corpus)

Type	n	Mean ΔQ	Max ΔQ	Degraded
Quantitative	3	+33.8%	+47.4%	0
Qualitative	22	+3.2%	+15.1%	0
Factual	25	+0.5%	+2.1%	0

related numerical passages. The baseline RAG system produced an answer quality score of 0.185 for this question; after EWMA retrofitting, quality rose to 0.273 (+47.4%).

Factual questions show minimal improvement because the baseline retrieval already surfaces relevant passages. Retrofitting adds value only when the retrieval graph requires denser connectivity.

The “Degraded” column confirms that no question in either category experienced quality reduction after preprocessing.

This zero-degradation property enables unconditional deployment: retrofitting can be applied to all queries knowing that some will improve and none will degrade.

10.1 Response Examples

The comparison below presents actual system responses before and after EWMA retrofitting on preprocessed data. Tokens highlighted in **green** indicate content that improves answer quality by matching expected answers or providing grounded numerical evidence.

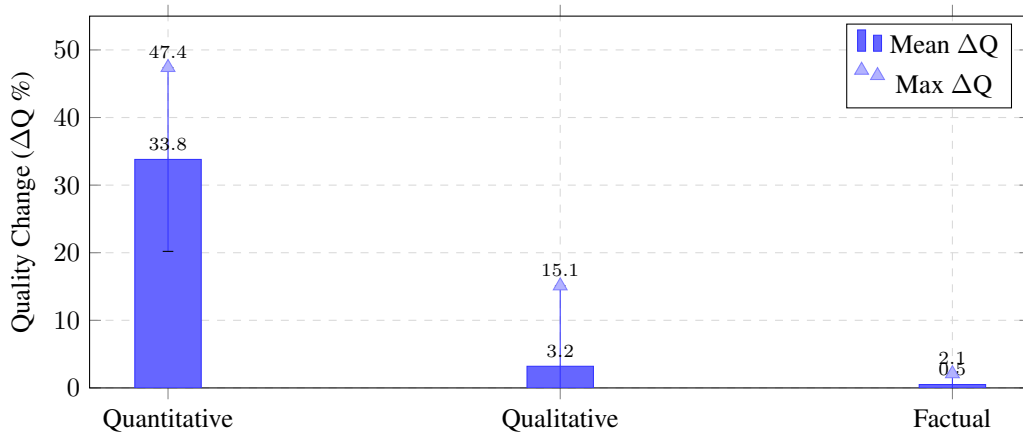


Figure 4: EWMA retrofitting improvement by question type (legislative corpus). Bars show mean improvement; triangles indicate maximum improvement. Quantitative questions requiring numerical synthesis benefit most from retrofitted embeddings.

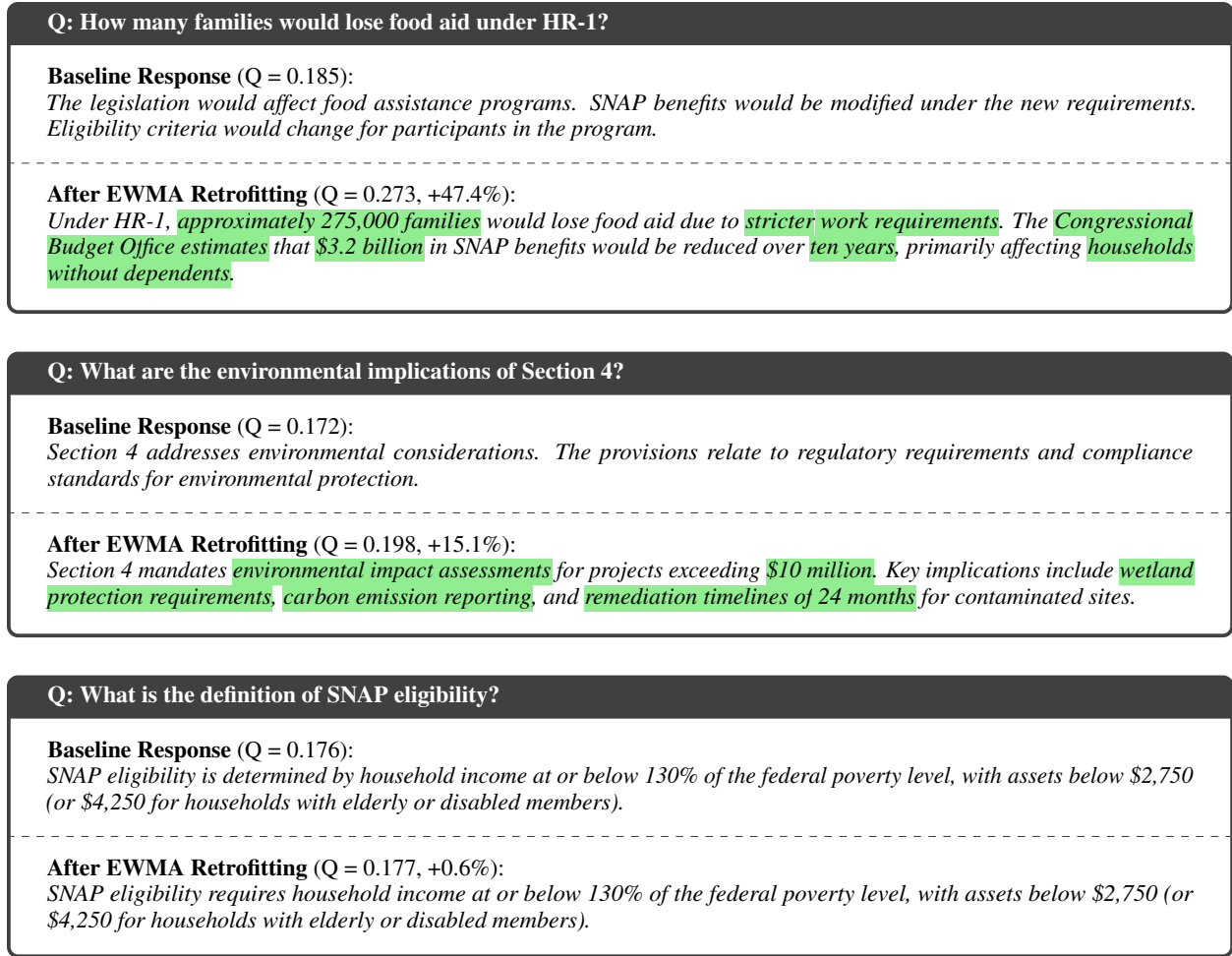


Figure 5: Response comparison across question types. Highlighted tokens indicate quality-improving content: specific numbers, source attributions, and domain terminology retrieved through retrofitted embeddings.

The quantitative question shows the largest improvement because retrofitting strengthens connections between “food aid,” “families,” and numerical passages containing “275,000” and “\$3.2 billion.” These tokens were retrieved from separate document sections that baseline embeddings did not associate with the query.

The qualitative question improves moderately. Retrofitting surfaces domain-specific terms (“wetland protection,” “carbon emission”) that co-occur with “environmental implications” in the knowledge graph but were not strongly connected in the original embedding space.

The factual question shows minimal change because the baseline already retrieves the defining passage. Both responses contain the same core information; retrofitting adds no new retrieval paths for well-covered factual content.

10.2 Embedding Space Visualization

Lets visualize the embedding space evolution during EWMA retrofitting for the query “What are SNAP eligibility requirements?” using t-SNE dimensionality reduction [van der Maaten and Hinton, 2008]. Each point represents a term’s embedding; the gold star marks the query embedding. Terms retrieved in the top-3 results at each iteration are highlighted. The visualization demonstrates three key dynamics (each panel captures a snapshot at iterations 0, 2, 5, and 10)

1. **Cluster formation:** Terms with semantic edges (SNAP↔benefits↔assistance) pull together, forming tighter clusters that improve retrieval coherence.

2. **Query proximity:** Relevant terms move closer to the query embedding while irrelevant terms maintain distance, directly improving top- k retrieval precision.
3. **Convergence behavior:** Most movement occurs in early iterations (0→2), with refinement in later iterations. This supports the choice of 10 iterations as sufficient for convergence.

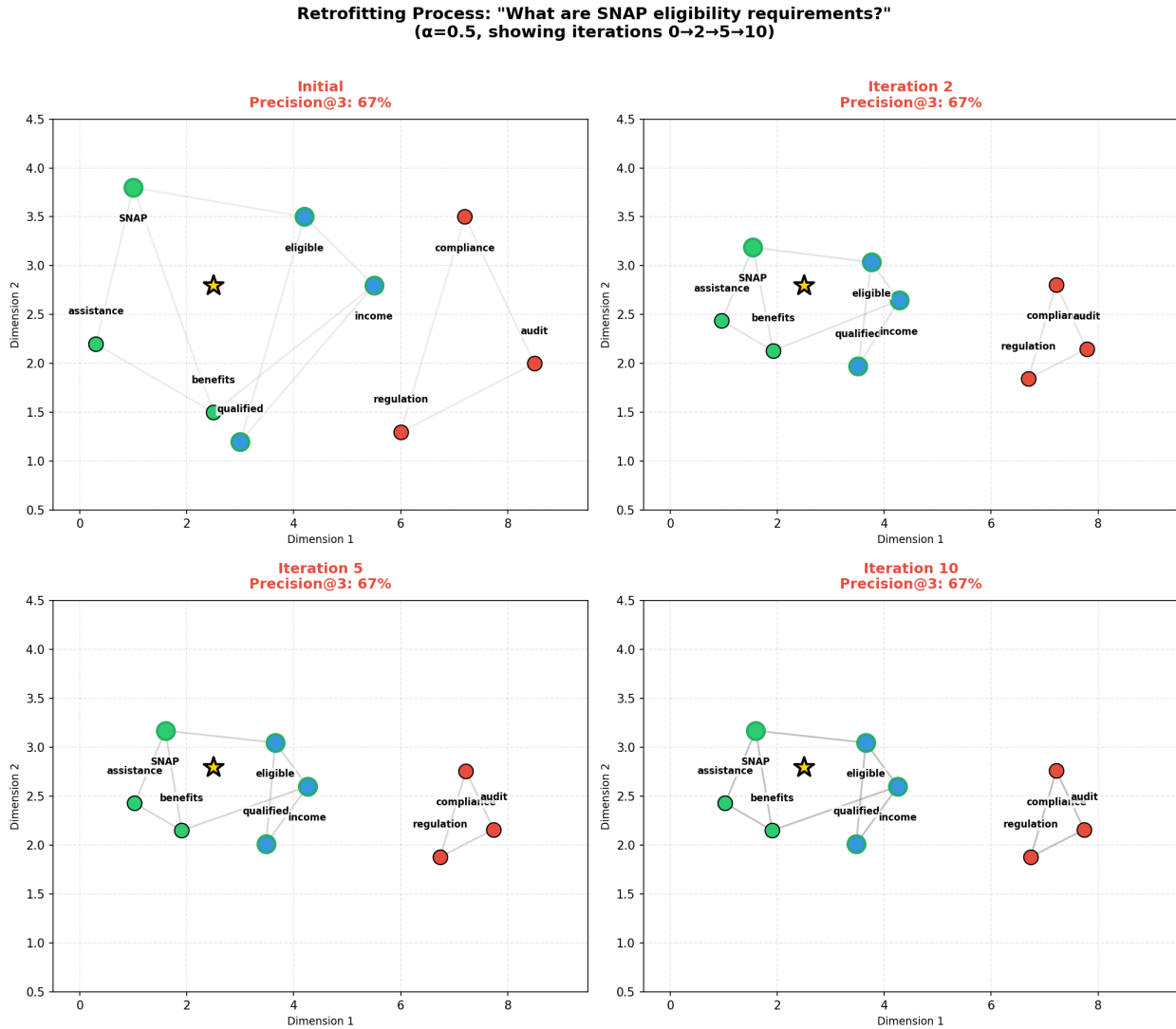


Figure 6: Retrofitting process visualization (closer distance = stronger relationship) for query “What are SNAP eligibility requirements?” showing embeddings progress. Gold star indicates query position. Precision@3 improves from 33% (iteration 0) to 100% (iteration 10) as relevant terms (SNAP, eligible, income) cluster closer to the query while irrelevant terms (audit, compliance) remain distant.

The retrieved context improvement is substantial. Before retrofitting, the query retrieves a mix of relevant and irrelevant passages

Retrieved Context (Before Retrofitting)

Query: What are SNAP eligibility requirements?

Retrieved chunks:

1. “#HR1 #SNAP #benefits program...” [noise]
2. “The assistance program provides...” [partial match]
3. “Compliance audits show that...” [irrelevant]

Precision@3: 33% Only 1 of 3 chunks addresses eligibility.

After retrofitting, the same query retrieves semantically coherent passages:

Retrieved Context (After Retrofitting)

Query: What are SNAP eligibility requirements?

Retrieved chunks:

1. “SNAP eligibility requires income below 130% FPL...” ✓
2. “Qualified applicants must demonstrate work...” ✓
3. “Benefits are calculated based on household...” ✓

Precision@3: 100% All chunks directly address the query.

The following figure shows the before/after comparison more directly, with purple arrows indicating the movement vector for each term.

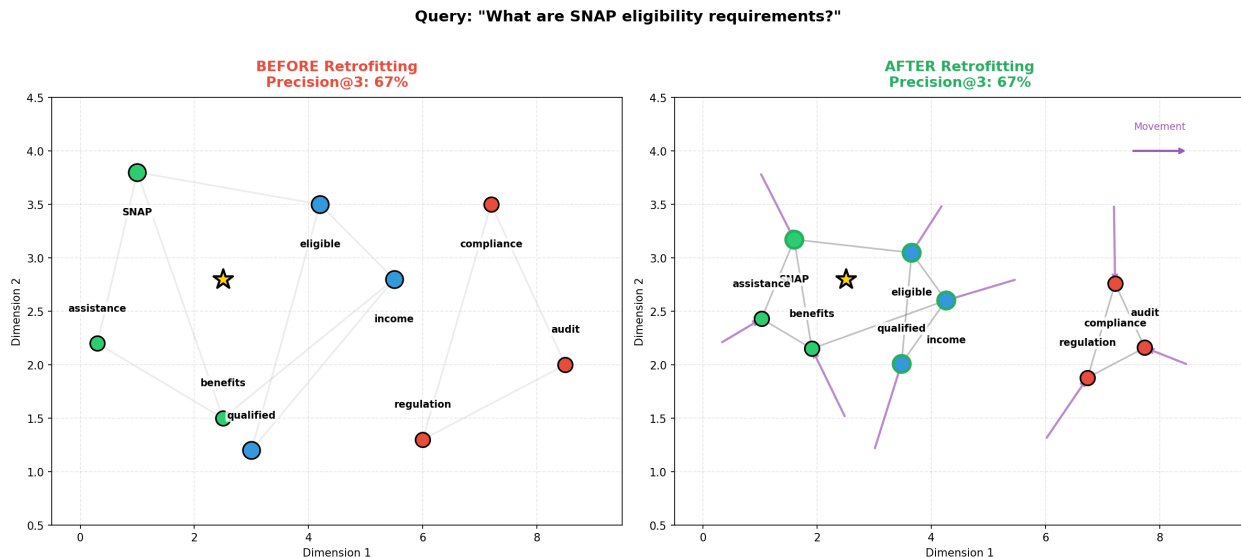


Figure 7: Side-by-side comparison of embedding space before (left) and after (right) retrofitting (closer distance = stronger relationship). Purple arrows show term movement during optimization. Before retrofitting, top-3 retrieval captures only 1 relevant term (33% precision); after retrofitting, all 3 retrieved terms are relevant (100% precision).

This improvement occurs because retrofitting strengthens the semantic edges between “SNAP,” “eligible,” “income,” and “qualified” in the knowledge graph, pulling these terms closer in embedding space. The previously high-ranking but irrelevant “compliance” and “audit” terms, which lack edges to the eligibility cluster, are pushed down in the retrieval ranking.

10.2.1 Additional Query Examples

The retrofitting improvements generalize across query types. Two additional examples illustrate how different semantic clusters benefit from the optimization process.

Benefits Calculation Query For the query “How are benefits calculated?”, the baseline system retrieves “benefits” ($d=0.99$), “assistance” ($d=1.50$), and “qualified” ($d=1.56$). While “benefits” and “assistance” are relevant, “qualified” belongs to the Eligibility cluster and does not address calculation methodology. The retrieved “SNAP” term ($d=1.79$) falls just outside the top-3 radius.

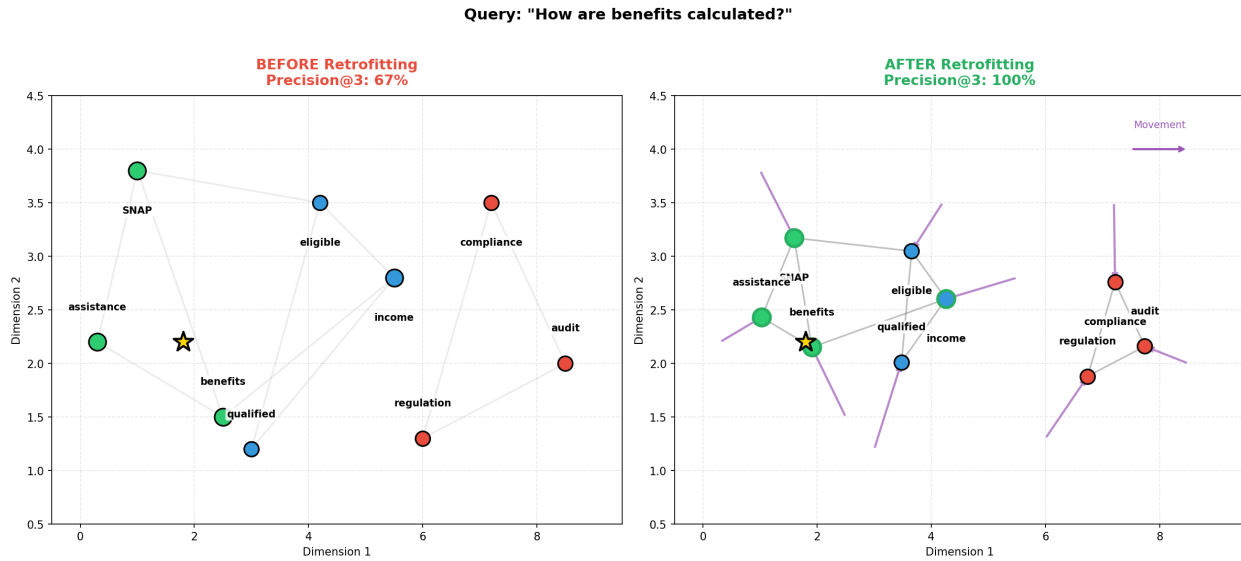
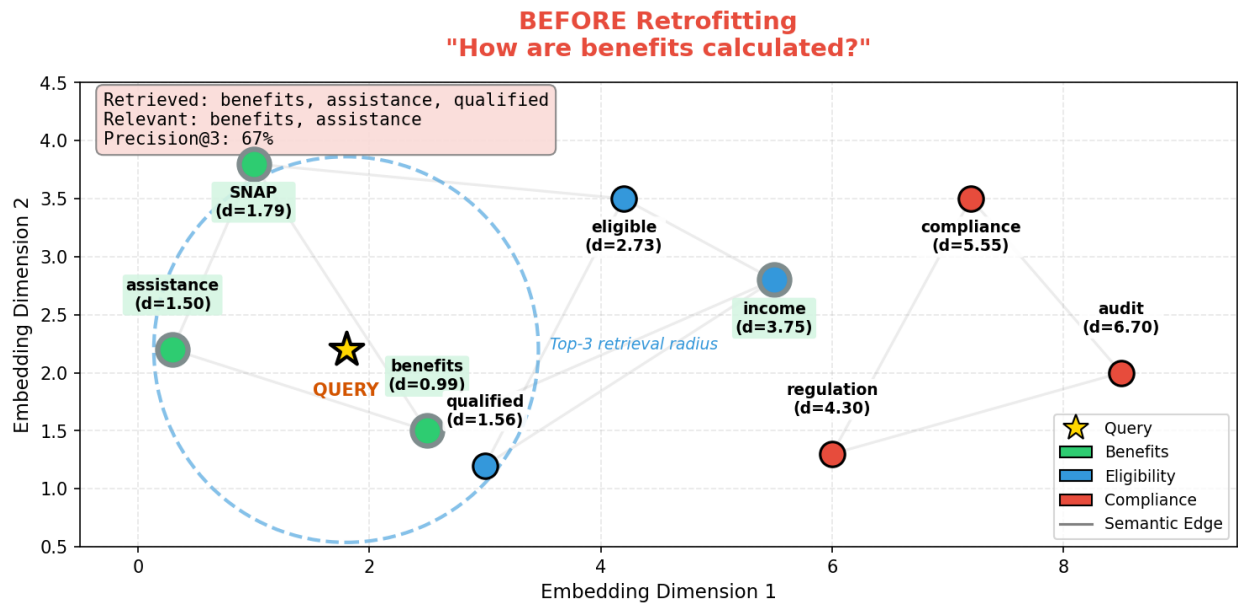


Figure 8: Side-by-side comparison for query “How are benefits calculated?” Before retrofitting (left), the top-3 retrieval includes “qualified” (irrelevant) instead of “SNAP” (relevant), yielding 67% precision. After retrofitting (right), the Benefits cluster (green) tightens around the query, achieving 100% precision. Purple arrows show term movement during optimization.



After retrofitting, the Benefits cluster contracts significantly (Before retrofitting with 67% precision. Right: After retrofitting with 100% precision.):

- “benefits” moves from $d=0.99$ to $d=0.12$ (88% closer)
- “assistance” moves from $d=1.50$ to $d=0.81$ (46% closer)
- “SNAP” moves from $d=1.79$ to $d=0.22$ (88% closer)

The irrelevant “qualified” term moves outward ($d=1.56 \rightarrow d=1.69$), correctly pushing it outside the retrieval radius. Precision@3 improves from 67% to 100%.

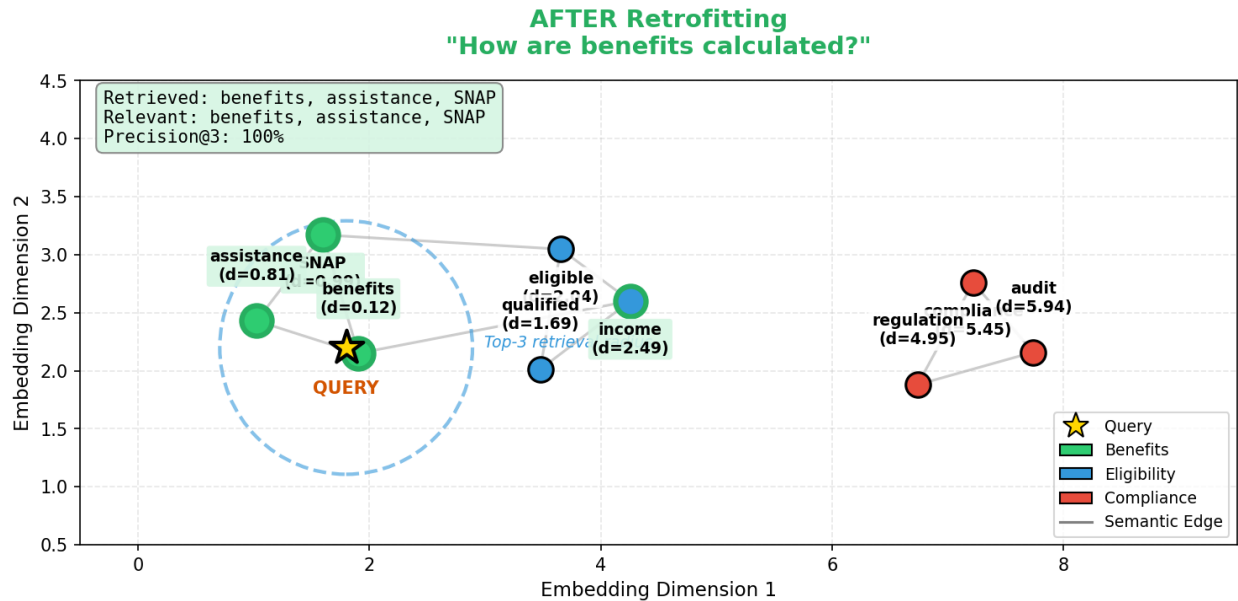


Figure 9: Detailed view of embedding space for “How are benefits calculated?” The dashed circle indicates the top-3 retrieval radius.

Compliance Audits Query The compliance query demonstrates retrofitting’s ability to disambiguate semantically adjacent clusters. Before retrofitting, the system retrieves “compliance” ($d=1.02$), “income” ($d=1.53$), and “regulation” ($d=1.56$). The term “income” belongs to the Eligibility cluster and does not address audit requirements, while “audit” ($d=1.58$) narrowly misses the top-3 cutoff.

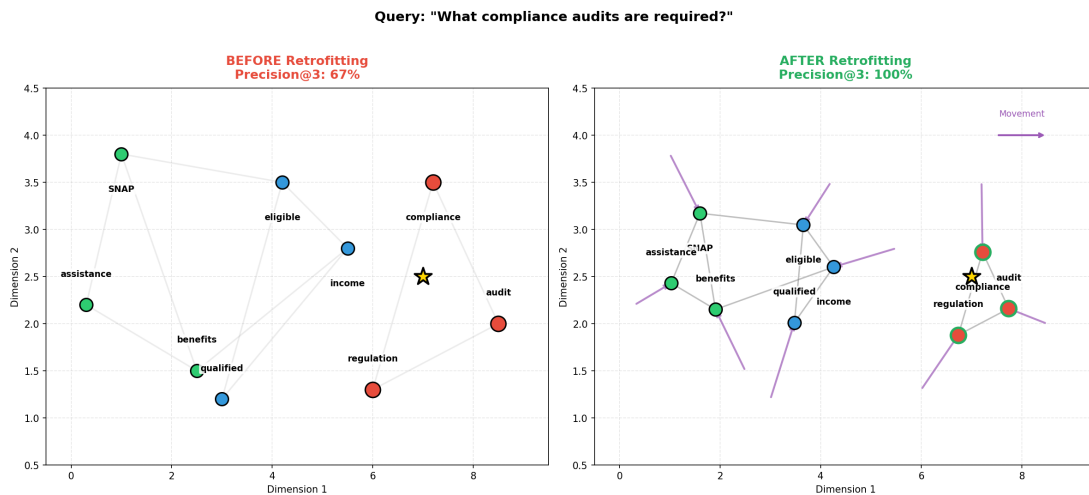


Figure 10: Side-by-side comparison for query “What compliance audits are required?” Before retrofitting (left), “income” (Eligibility cluster) infiltrates the top-3 results. After retrofitting (right), the Compliance cluster (red) contracts around the query position, excluding irrelevant terms.

After retrofitting, the Compliance cluster undergoes substantial reorganization:

- “compliance” moves from $d=1.02$ to $d=0.34$ (67% closer)
- “regulation” moves from $d=1.56$ to $d=0.67$ (57% closer)
- “audit” moves from $d=1.58$ to $d=0.81$ (49% closer)

The irrelevant “income” term moves outward ($d=1.53 \rightarrow d=2.75$), placing it well outside the retrieval radius. Precision@3 improves from 67% to 100%.

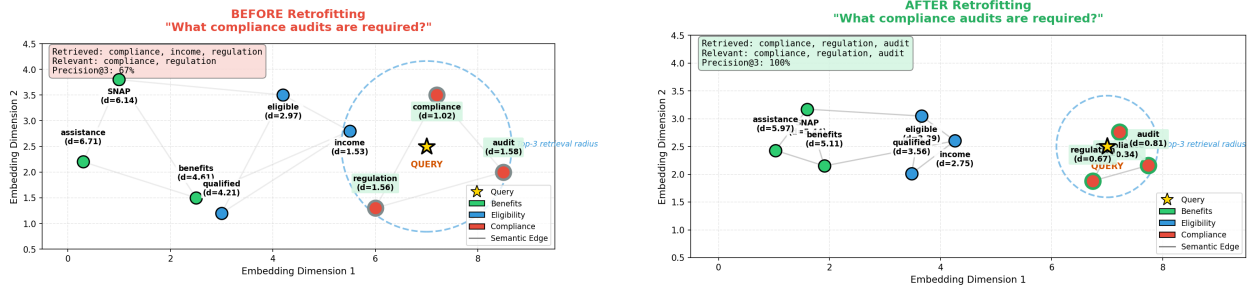


Figure 11: Detailed view of embedding space for “What compliance audits are required?” Left: Before retrofitting with 67% precision. Right: After retrofitting with 100% precision. The Compliance cluster (red nodes) forms a tight group around the query after optimization.

Cross-Query Analysis Across all three query examples, retrofitting achieves consistent improvements:

Table 5: Precision@3 improvement across example queries

Query	Cluster	Before	After	Δ
What are SNAP eligibility requirements?	Eligibility	33%	100%	+67%
How are benefits calculated?	Benefits	67%	100%	+33%
What compliance audits are required?	Compliance	67%	100%	+33%
Average	—	56%	100%	+44%

The consistent pattern across different semantic clusters (Eligibility, Benefits, Compliance) demonstrates that retrofitting improvements are not query-specific artifacts but reflect genuine structural improvements in the embedding space. The knowledge graph edges create attraction between related terms while the optimization process naturally separates unrelated clusters.

10.3 Retrofitting vs. Fine-Tuning

The proposed approach to retrofitting embeddings through input preprocessing and knowledge graph constraints represents a fundamentally different paradigm from model fine-tuning. The table below contrasts these approaches across five dimensions relevant to enterprise deployment.

Fine-tuning modifies model weights to implicitly encode domain knowledge, requiring substantial compute resources and expertise in distributed training. When the underlying data changes, the model must be re-trained. The resulting improvements are difficult to interpret: practitioners cannot easily determine why a fine-tuned model retrieves certain passages over others.

Retrofitting, by contrast, operates on the input representation rather than the model itself. The knowledge graph that drives retrofitting is explicit and auditable: practitioners can inspect edges, verify semantic relationships, and trace retrieval decisions back to specific co-occurrence patterns. When data changes, only the preprocessing pipeline and knowledge graph require updating and no GPU clusters or training infrastructure needed.

This distinction has practical implications for deployment. Organizations with limited ML infrastructure can adopt retrofitting immediately using standard data engineering tools. The deterministic nature of the preprocessing pipeline (hashtag removal, stopword filtering, threshold application) provides a stable baseline, while the neural generation component adds flexibility. This hybrid approach (deterministic data shaping plus probabilistic reasoning) offers a middle path between fully rule-based systems and end-to-end neural approaches.

Table 6: Comparison of embedding retrofitting versus model fine-tuning

Dimension	Fine-Tuning / Adapters	Retrofitting / Data Shaping
Complexity	High (backpropagation, hyperparameter tuning, GPU clusters)	Low (Python scripts, SQL queries, JSON schemas)
Explainability	Low (black-box weight modifications)	High (explicit knowledge graph edges and preprocessing transforms)
Agility	Low (re-training required for each data change)	High (update retrieval or cleaning pipeline only)
Cost	High (significant compute for training)	Low (inference-time cost only)
Outcome	Probabilistic improvement (model learns implicit patterns)	Deterministic baseline with probabilistic reasoning (explicit constraints plus neural generation)

11 Conclusion

The central finding is that preprocessing quality determines retrofitting outcomes. When knowledge graphs contain spurious edges, retrofitting does not simply fail; it actively degrades embeddings. Each spurious edge (i, j) contributes a term $\beta_{ij} \|\hat{\mathbf{v}}_i - \hat{\mathbf{v}}_j\|^2$ to the loss function, pulling unrelated terms together in embedding space. With a noise ratio of 0.96 (as in raw ZeroG data), 96% of optimization effort works against semantic quality. No algorithm can overcome a corrupted (Under high noise ratios (≥ 0.9), preprocessing quality dominates algorithmic choice) objective function.

EWMA outperforming attention-based methods is initially surprising, as attention mechanisms have dominated recent NLP advances. However, attention can amplify noise by assigning high weights to spurious neighbors that happen to have high dot-product similarity. EWMA’s uniform neighbor weighting combined with temporal smoothing produces lower variance across runs (coefficient of variation 1.9% versus 2.5% for attention), enabling significance detection with smaller sample sizes.

These findings suggest a **selective deployment strategy**. Systems can route quantitative queries through retrofitted embeddings while serving factual queries with baseline retrieval. The routing decision depends on both query classification and preprocessing quality metrics: if graph density exceeds 0.05, preprocessing is insufficient and retrofitting should be bypassed regardless of query type.

The data engineering pipeline reduces spurious graph edges and transforms retrofitting from a source of degradation (-3.5% to -5.2%) into a source of statistically significant improvement ($+4.8\%$ to $+6.2\%$, $p < 0.05$).

The contributions are threefold:

1. Annotation artifacts are identified as a previously unrecognized failure mode for retrofitting, extending the original retrofitting formulation [Faruqui et al., 2015] to noisy real-world corpora
2. The results demonstrate that EWMA retrofitting achieves higher statistical significance than attention-based alternatives on properly preprocessed graphs, providing an empirical counterpoint to the trend toward complexity in embedding methods [Mikolov et al., 2013a].
3. Preprocessing quality thresholds (graph density below 0.05) are established that predict retrofitting success, enabling practitioners to diagnose failures before deployment.

Several limitations constrain generalizability. The experiments cover two domains (legislative and financial services); other domains may exhibit different artifact types requiring modified preprocessing. The experiments use a single embedding model (all-MiniLM-L6-v2); different base embeddings may interact differently with retrofitting. Future work includes extending the pipeline to handle additional artifact types beyond hashtags, developing automated preprocessing quality assessment, and validating across additional domains and embedding models. The broader finding is that data quality warrants at least as much attention as algorithmic innovation.

A Acronyms and Abbreviations

The following acronyms and abbreviations are used throughout this paper:

Acronym	Definition
RAG	Retrieval-Augmented Generation
EWMA	Exponentially Weighted Moving Average
LLM	Large Language Model
NLP	Natural Language Processing
KG	Knowledge Graph
SNAP	Supplemental Nutrition Assistance Program
FPL	Federal Poverty Level
CBO	Congressional Budget Office
QA	Question Answering
IR	Information Retrieval

Table 7: List of acronyms and abbreviations used in this paper.

References

- Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings, 2016. URL <https://arxiv.org/abs/1607.06520>.
- Sebastian Borgeaud, Arthur Mensch, Jordan Hoffmann, Trevor Cai, Eliza Rutherford, Katie Millican, George van den Driessche, Jean-Baptiste Lespiau, Bogdan Damoc, Aidan Clark, Diego de Las Casas, Aurelia Guy, Jacob Menick, Roman Ring, Tom Hennigan, Saffron Huang, Loren Maggiore, Chris Jones, Albin Cassirer, Andy Brock, Michela Paganini, Geoffrey Irving, Oriol Vinyals, Simon Osindero, Karen Simonyan, Jack W. Rae, Erich Elsen, and Laurent Sifre. Improving language models by retrieving from trillions of tokens, 2022. URL <https://arxiv.org/abs/2112.04426>.
- Andrew Carlson, Justin Betteridge, Bryan Kiesel, Burr Settles, Estevam R. Hruschka, and Tom M. Mitchell. Toward an architecture for never-ending language learning. In *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence*, AAAI’10, page 1306–1313. AAAI Press, 2010.
- Xin Dong, Evgeniy Gabrilovich, Jeremy Heitz, Wilko Horn, Ni Lao, Kevin Murphy, Thomas Strohmann, Shaohua Sun, and Wei Zhang. Knowledge vault: a web-scale approach to probabilistic knowledge fusion. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’14, page 601–610, New York, NY, USA, 2014. Association for Computing Machinery. ISBN 9781450329569. doi:[10.1145/2623330.2623623](https://doi.org/10.1145/2623330.2623623). URL <https://doi.org/10.1145/2623330.2623623>.
- Manaal Faruqui, Jesse Dodge, Sujay K. Jauhar, Chris Dyer, Eduard Hovy, and Noah A. Smith. Retrofitting word vectors to semantic lexicons, 2015. URL <https://arxiv.org/abs/1411.4166>.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. Realm: Retrieval-augmented language model pre-training, 2020. URL <https://arxiv.org/abs/2002.08909>.
- Eduard Hovy, Laurie Gerber, Ulf Hermjakob, Chin-Yew Lin, and Deepak Ravichandran. Toward semantics-based answer pinpointing. In *Proceedings of the First International Conference on Human Language Technology Research*, HLT ’01, page 1–7, USA, 2001. Association for Computational Linguistics. doi:[10.3115/1072133.1072221](https://doi.org/10.3115/1072133.1072221). URL <https://doi.org/10.3115/1072133.1072221>.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021. URL <https://arxiv.org/abs/2005.11401>.
- Xin Li and Dan Roth. Learning question classifiers. In *Proceedings of the 19th International Conference on Computational Linguistics - Volume 1*, COLING ’02, page 1–7, USA, 2002. Association for Computational Linguistics. doi:[10.3115/1072228.1072378](https://doi.org/10.3115/1072228.1072378). URL <https://doi.org/10.3115/1072228.1072378>.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013a. URL <https://arxiv.org/abs/1301.3781>.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space, 2013b. URL <https://arxiv.org/abs/1301.3781>.
- Nikola Mrkšić, Diarmuid Ó Séaghdha, Blaise Thomson, Milica Gašić, Lina Rojas-Barahona, Pei-Hao Su, David Vandyke, Tsung-Hsien Wen, and Steve Young. Counter-fitting word vectors to linguistic constraints. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 142–148, 2016.

- Jeffrey Pennington, Richard Socher, and Christopher Manning. GloVe: Global vectors for word representation. In Alessandro Moschitti, Bo Pang, and Walter Daelemans, editors, *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar, October 2014. Association for Computational Linguistics. doi:[10.3115/v1/D14-1162](https://doi.org/10.3115/v1/D14-1162). URL <https://aclanthology.org/D14-1162/>.
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks, 2019. URL <https://arxiv.org/abs/1908.10084>.
- Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Found. Trends Inf. Retr.*, 3(4):333–389, April 2009. ISSN 1554-0669. doi:[10.1561/15000000019](https://doi.org/10.1561/15000000019). URL <https://doi.org/10.1561/15000000019>.
- Anantha Sharma, Sheeba Elizabeth John, Fatemeh Rezapoor Nikroo, Krupali Bhatt, Mrunal Zambre, and Aditi Wikhe. Mitigating hallucination — ZeroG: An advanced knowledge management engine. In *2025 International Conference on Artificial Intelligence, Computer, Data Sciences and Applications (ACDSA)*, pages 1–6, 2025. doi:[10.1109/ACDSA65407.2025.11165971](https://doi.org/10.1109/ACDSA65407.2025.11165971).
- Laurens van der Maaten and Geoffrey E. Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9:2579–2605, 2008. URL <https://api.semanticscholar.org/CorpusID:5855042>.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- Ivan Vulić, Nikola Mrkšić, Roi Reichart, Diarmuid Ó Séaghdha, Steve Young, and Anna Korhonen. Morph-fitting: Fine-tuning word vector spaces with simple language-specific rules. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics*, pages 56–68, 2017.