

# Dynamic Gatekeeper

A Batch-Relative Filtering Algorithm for RAG Judge Scores

Technical Report — February 2026

## 1. Problem Statement and Motivation

Retrieval-Augmented Generation (RAG) systems in high-stakes domains such as legal precedent analysis retrieve a fixed-size candidate set — typically the top- $k$  documents by vector similarity — and pass them to a Judge Model that assigns relevance scores in  $[0, 1]$ . The core difficulty is *contextual calibration drift*: the Judge’s scores are relative to the difficulty of the current batch, not on a universal absolute scale.

A fixed threshold  $\tau$  fails in both directions. A strict threshold (e.g.,  $\tau = 0.8$ ) returns zero results when the entire batch is mediocre — a catastrophic failure for a production system. A lenient threshold (e.g.,  $\tau = 0.4$ ) admits irrelevant noise when the batch is uniformly strong. The problem therefore demands an algorithm that is *simultaneously* capable of noise reduction in easy batches and signal salvaging in hard batches — two directly opposing goals satisfied by a single adaptive mechanism.

## 2. Approach and Methodology

The proposed algorithm, **Dynamic Gatekeeper**, computes a final threshold by blending three complementary components and enforcing a hard keep-ratio safety net. Each component targets a distinct failure mode.

### 2.1 Component 1 — Statistical Threshold

Inspired by the MAIN-RAG adaptive thresholding framework, the statistical threshold anchors the decision boundary to the batch’s own distribution:

$$\tau_{\text{stat}} = \text{clip}(\bar{x} - \alpha \sigma, 0.10, 0.95)$$

where  $\bar{x}$  is the batch mean,  $\sigma$  the population standard deviation, and  $\alpha = 0.6$  a tunable multiplier. When  $\bar{x}$  is high (easy batch),  $\tau_{\text{stat}}$  is high, enforcing strictness. When  $\bar{x}$  is low (hard batch),  $\tau_{\text{stat}}$  falls, enabling leniency. Critically,  $\alpha$  is a fixed constant, so  $\sigma$  retains genuine influence on the threshold throughout — avoiding the algebraic cancellation flaw present in the formula  $n = (0.5 - \bar{x})/\sigma$ , where the standard deviation terms cancel in the unclamped regime.

### 2.2 Component 2 — Natural Gap Threshold

Otsu’s method, while standard in image processing, requires histogram-level sample sizes ( $N \gg 100$ ) to estimate between-class variance reliably. With only  $k = 10$  retrieved documents, it produces numerically arbitrary results. The natural gap detector is a valid replacement for small  $n$ :

$$\tau_{\text{gap}} = \frac{s_{(j)} + s_{(j+1)}}{2}, \quad j = \arg \max_i (s_{(i)} - s_{(i+1)})$$

where  $s_{(1)} \geq s_{(2)} \geq \dots \geq s_{(k)}$  are the sorted scores in descending order. The threshold is placed at the midpoint of the single largest consecutive drop. For bimodal (mixed) batches this cleanly

isolates the high-quality cluster. For unimodal batches the largest gap is small, so the component contributes modestly and does not distort the blend.

### 2.3 Component 3 — IQR Floor Guard

The interquartile-range lower bound prevents the blended threshold from falling below the region of non-outlier scores:

$$\tau_{\text{IQR}} = \max(0, Q_1 - 1.5 \times \text{IQR})$$

Unlike prior approaches where this quantity was computed but never applied to the threshold logic, here it acts as a *hard floor*: the final threshold cannot drop below  $\tau_{\text{IQR}}$ , preventing statistically anomalous low-scoring documents from being retained even in lenient batches.

### 2.4 Blending and Safety Net

The two primary components are combined with weights  $w_s = 0.65$  and  $w_g = 0.35$ , then the IQR floor is applied:

$$\tau_{\text{blend}} = w_s \tau_{\text{stat}} + w_g \tau_{\text{gap}}, \quad \tau_{\text{final}} = \text{clip}(\max(\tau_{\text{blend}}, \tau_{\text{IQR}}), 0.05, 0.99)$$

The statistical component is weighted higher because it is always well-defined; the gap component carries more influence only when a natural separation is detectable. After applying  $\tau_{\text{final}}$ , a keep-ratio safety net enforces  $r \in [r_{\min}, r_{\max}] = [0.20, 0.80]$ , overriding the threshold if the resulting fraction of kept documents would otherwise be too small (catastrophic return of nothing) or too large (uncritical acceptance of noise).

## 3. Assumptions

The algorithm rests on three explicit assumptions. First, Judge Model scores are internally consistent within a single batch, even if they drift across batches — this is precisely the calibration drift problem being addressed. Second, the top- $k$  retrieved documents contain at least one genuinely useful document in the majority of query cases; the algorithm is not designed for queries with zero relevant documents in the corpus. Third, the batch size is small and fixed ( $k \approx 10$ ), which motivates the gap-based approach over histogram methods and the use of population (rather than sample) standard deviation for stability.

## 4. Empirical Findings

The algorithm was validated against seven controlled scenarios covering the full spectrum of batch types. Table 1 summarises the key results.

Three findings are of particular note. In the **easy batch**, the high batch mean pushes  $\tau_{\text{stat}}$  to 0.888 and the gap component detects the small drop around 0.905, producing a blended threshold of 0.908 that retains only the top five documents — demonstrating noise reduction. In the **mixed batch**, the gap component ( $\tau_{\text{gap}} = 0.68$ ) successfully identifies the natural break between the high-quality cluster ( $\approx 0.88\text{--}0.92$ ) and the noise cluster ( $\approx 0.41\text{--}0.48$ ), and the blend correctly keeps exactly four documents. In the **extreme low batch**, the IQR floor and the clipped  $\tau_{\text{stat}}$  (floored at 0.10) combine with the min-keep safety net to salvage the three best documents rather than returning an empty set.

<b>Scenario</b>	<b>Mean</b>	$\tau_{\text{final}}$	<b>Kept</b>	<b>Ratio</b>
Easy Batch (all high)	0.905	0.908	5 / 10	0.50
Hard Batch (all mediocre)	0.445	0.448	5 / 10	0.50
Mixed Batch (bimodal)	0.627	0.558	4 / 10	0.40
Near-Identical Scores	0.608	0.608	6 / 10	0.60
Extreme Low Batch	0.097	0.109	3 / 10	0.30
Single Document	0.730	0.730	1 / 1	1.00
Empty Input	—	—	0 / 0	—

Table 1: Test scenario results.  $\tau_{\text{final}}$  adapts to each batch; no scenario returns zero results.

## 5. Conclusions

The Dynamic Gatekeeper resolves the calibration drift problem through a principled, three-component adaptive mechanism. By replacing Otsu’s method with a gap detector suited to small  $n$ , using a direct  $\bar{x} - \alpha\sigma$  formula that preserves the influence of standard deviation, applying the IQR bound as an active floor rather than passive metadata, and enforcing keep-ratio constraints as a final guarantee, the algorithm satisfies both the noise reduction and signal salvaging requirements simultaneously across all tested batch types. The design is transparent, auditable via the returned metadata dictionary, and tunable via three well-defined hyperparameters ( $\alpha$ ,  $r_{\min}$ ,  $r_{\max}$ ) that can be calibrated against domain-specific Judge Model behaviour.

---

*Implementation verified against 7 test scenarios. All referenced statistical techniques (MAIN-RAG adaptive thresholding, IQR outlier detection, gap-based natural break detection) are standard methods from the distributional analysis literature.*