# Bayesian BM25: A Probabilistic Framework for Hybrid Text and Vector Search

**Jaepil Jeong**
Cognica, Inc.
*Email: [jaepil@cognica.io](jaepil@cognica.io)*
Date: January 23, 2026

*"The theory of probabilities is at bottom nothing but common sense reduced to calculus."*
— Pierre-Simon Laplace

## Abstract

This paper presents Bayesian BM25, a novel probabilistic framework that transforms traditional BM25 relevance scores into calibrated probability estimates, enabling principled integration of lexical and semantic search signals. We demonstrate that standard BM25 scores, while effective for ranking, suffer from fundamental interpretability limitations: unbounded range, query-dependent magnitudes, and incompatibility with other ranking signals. We address these limitations by applying Bayesian inference with a sigmoid likelihood model and composite prior design incorporating term frequency and document length signals. We prove that our transformation preserves the monotonicity properties essential for ranking while producing well-calibrated probability outputs in $[0, 1]$. We extend this framework to hybrid search, deriving probabilistic score combination formulas for Boolean query operations under independence assumptions. We show that log-space computation ensures numerical stability for extreme probabilities. Furthermore, we analyze the compatibility of Bayesian BM25 with WAND (Weak AND) and BMW (Block-Max WAND) optimization algorithms, proving that the transformation preserves the upper bound properties required for safe document pruning. Experimental results demonstrate that Bayesian BM25 achieves comparable ranking quality to standard BM25 while enabling principled multi-signal fusion without ad-hoc normalization or weighting schemes.

## 1. Introduction and Motivation

### 1.1 The BM25 Score Interpretation Problem

The BM25 (Best Matching 25) scoring function, introduced by Robertson and Zaragoza (2009), has become the de facto standard for lexical relevance scoring in information retrieval systems. Despite its empirical success, BM25 scores present fundamental challenges for modern search applications that require combining multiple ranking signals.

**Definition 1.1.1** (BM25 Score). For a document $D$ and query $Q = \{t_1, t_2, \ldots, t_m\}$, the BM25 score is defined as:

$$\text{BM25}(D, Q) = \sum_{t \in Q} \text{IDF}(t) \cdot \frac{f(t, D) \cdot (k_1 + 1)}{f(t, D) + k_1 \cdot B(D)} \tag{1}$$

where:

- $f(t, D)$ is the term frequency of $t$ in document $D$

- $\mathrm{IDF}(t)$ is the inverse document frequency of term $t$
- $B(D) = 1 - b + b \cdot \frac{|D|}{\mathrm{avgdl}}$ is the length normalization factor
- $k_1$ and $b$ are tuning parameters (typically $k_1 = 1.2$ and $b = 0.75$)

**Problem 1.1.2** (Score Interpretability). BM25 scores suffer from the following interpretability limitations:

1. **Unbounded Range**: $\mathrm{BM25}(D, Q) \in [0, +\infty)$, making absolute interpretation impossible
2. **Query Dependence**: Score magnitudes vary with query length and term specificity
3. **Corpus Dependence**: IDF values depend on collection statistics, preventing cross-corpus comparison
4. **Signal Incompatibility**: Direct combination with bounded signals (e.g., vector similarity in $[0, 1]$) leads to dominated or unbalanced contributions

## 1.2 The Multi-Signal Fusion Challenge

Modern search systems require combining heterogeneous ranking signals with different scales and distributions.

**Definition 1.2.1** (Multi-Signal Ranking Function). A multi-signal ranking function combines $n$ signals:

$$\mathrm{Score}(D, Q) = f(s_1(D, Q), s_2(D, Q), \ldots, s_n(D, Q)) \tag{2}$$

where each $s_i : \mathcal{D} \times \mathcal{Q} \to \mathbb{R}$ is a scoring function with potentially different range and distribution.

**Theorem 1.2.2** (Signal Dominance in Naive Combination). For additive combination of signals with different scales:

$$\mathrm{Score}_{naive} = \sum_{i=1}^{n} s_i \tag{3}$$

the signal with the largest expected magnitude will dominate the ranking.

*Proof.* Let $\mu_i = \mathbb{E}[s_i]$ and $\sigma_i = \mathrm{Std}[s_i]$. For two signals with $\mu_1 \gg \mu_2$, the probability that signal 2 affects the ranking is:

$$P(s_2 > s_1) \leq P\left(\frac{s_1 - \mu_1}{\sigma_1} < -\frac{\mu_1 - \mu_2}{\sigma_1}\right) \approx 0 \tag{4}$$

when $\mu_1 - \mu_2 \gg \sigma_1$.

## 1.3 Reciprocal Rank Fusion Limitations

Reciprocal Rank Fusion (RRF), proposed by Cormack et al. (2009), is a common approach to combining ranked lists.

**Definition 1.3.1** (Reciprocal Rank Fusion). For $n$ ranked lists $R_1, \ldots, R_n$, the RRF score is:

$$\mathrm{RRF}(D) = \sum_{i=1}^{n} \frac{1}{k + \mathrm{rank}_i(D)} \tag{5}$$

where $k$ is a constant (typically $k = 60$) and $\mathrm{rank}_i(D)$ is the rank of document $D$ in list $R_i$.

**Theorem 1.3.2** (RRF Limitations). RRF suffers from the following theoretical limitations:

1. **Information Loss**: Score magnitudes are discarded, losing confidence information

2. **Rank Sensitivity**: Small perturbations in scores can cause large rank changes

3. **Non-Commutativity with Filtering**:
   $\mathrm{RRF}(\mathrm{Filter}(L_1), \mathrm{Filter}(L_2)) \neq \mathrm{Filter}(\mathrm{RRF}(L_1, L_2))$

4. **Arbitrary Constant**: The $k$ parameter lacks theoretical justification

5. **Missing Document Handling**: Undefined behavior for documents appearing in only some lists

*Proof of (3)*. Consider $L_1 = [A:10, B:5, C:3]$ and $L_2 = [B:8, A:4, C:1]$. With filter threshold $> 6$:

- Left side: $\mathrm{RRF}([A:10], [B:8]) = \{A: \frac{1}{61}, B: \frac{1}{61}\}$
- Right side: $\mathrm{Filter}(\mathrm{RRF}(L_1, L_2)) = \mathrm{Filter}(\{A: \frac{1}{61} + \frac{1}{62}, B: \frac{1}{62} + \frac{1}{61}, C: \frac{1}{63} + \frac{1}{63}\})$

The results differ because filtering after fusion considers documents that wouldn't pass individual filters.

## 1.4 Our Contribution

We propose Bayesian BM25, a probabilistic framework that:

1. Transforms BM25 scores into calibrated probabilities in $[0, 1]$

2. Preserves monotonicity properties essential for ranking

3. Enables principled probabilistic combination of multiple signals

4. Maintains compatibility with WAND/BMW optimization algorithms

5. Supports online parameter learning for domain adaptation

# 2. Mathematical Preliminaries

## 2.1 Probability Calibration

**Definition 2.1.1** (Calibrated Probability). A scoring function $s : \mathcal{D} \times \mathcal{Q} \to [0, 1]$ is calibrated if:

$$P(\text{relevant} \mid s(D, Q) = p) = p \tag{6}$$

for all $p \in [0, 1]$.

**Definition 2.1.2** (Sigmoid Function). The sigmoid function $\sigma : \mathbb{R} \to (0, 1)$ is defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \tag{7}$$

**Lemma 2.1.3** (Sigmoid Properties). The sigmoid function satisfies:

1. $\sigma(-x) = 1 - \sigma(x)$ (symmetry)
2. $\frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x))$ (derivative)
3. $\lim_{x \to \infty} \sigma(x) = 1$ and $\lim_{x \to -\infty} \sigma(x) = 0$ (bounds)

*Proof*. Property (1):

$$\sigma(-x) = \frac{1}{1 + e^x} = \frac{e^{-x}}{e^{-x} + 1} = 1 - \frac{1}{1 + e^{-x}} = 1 - \sigma(x) \tag{8}$$

## 2.2 Bayesian Inference Framework

**Definition 2.2.1** (Bayes' Theorem). For hypothesis $H$ and evidence $E$:

$$P(H \mid E) = \frac{P(E \mid H) \cdot P(H)}{P(E)} \tag{9}$$

**Definition 2.2.2** (Binary Classification Posterior). For relevance $R \in \{0, 1\}$ and score $s$:

$$P(R = 1 \mid s) = \frac{P(s \mid R = 1) \cdot P(R = 1)}{P(s \mid R = 1) \cdot P(R = 1) + P(s \mid R = 0) \cdot P(R = 0)} \tag{10}$$

# 3. BM25 Theoretical Foundation

## 3.1 IDF Computation

**Definition 3.1.1** (Robertson-Sparck Jones IDF). The inverse document frequency with smoothing is:

$$\text{IDF}(t) = \ln \left( \frac{N - \text{df}(t) + 0.5}{\text{df}(t) + 0.5} + 1 \right) \tag{11}$$

where $N$ is the total document count and $\text{df}(t)$ is the document frequency of term $t$.

**Theorem 3.1.2** (IDF Non-Negativity). For any term $t$ with $\text{df}(t) \leq N$:

$$\text{IDF}(t) \geq 0 \tag{12}$$

*Proof.* The argument of the logarithm is:

$$\frac{N - \text{df}(t) + 0.5}{\text{df}(t) + 0.5} + 1 \geq 1 \tag{13}$$

since $N - \text{df}(t) \geq 0$. Therefore, $\text{IDF}(t) = \ln(\cdot) \geq \ln(1) = 0$.

**Theorem 3.1.3** (IDF Monotonicity). IDF is monotonically decreasing in document frequency:

$$\text{df}(t_1) < \text{df}(t_2) \implies \text{IDF}(t_1) > \text{IDF}(t_2) \tag{14}$$

*Proof.* Let $f(x) = \ln \left( \frac{N - x + 0.5}{x + 0.5} + 1 \right)$. Taking the derivative:

$$f'(x) = \frac{-(N + 1)}{(x + 0.5)(N - x + 0.5) + (x + 0.5)^2} < 0 \tag{15}$$

for $0 \leq x \leq N$, proving monotonicity.

## 3.2 Monotonicity-Preserving BM25 Formulation

**Definition 3.2.1** (Rewritten BM25 Scoring). We express BM25 in a numerically stable form:

$$\text{score}(f, n) = w - \frac{w}{1 + f \cdot \text{inv\_norm}} \tag{16}$$

where:

- $w = \text{boost} \cdot \text{IDF} \cdot (k_1 + 1)$ (weight, incorporating the BM25 saturation scaling factor)
- $\text{inv\_norm} = \frac{1}{k_1 \cdot ((1-b) + b \cdot \frac{n}{\text{avgdl}})}$ (inverse normalization)
- $f$ is term frequency
- $n$ is document length

*Remark.* The $(k_1 + 1)$ factor ensures equivalence with the standard BM25 formulation in Definition 1.1.1. Expanding:
$w - \frac{w}{1 + f \cdot \text{inv\_norm}} = w \cdot \frac{f \cdot \text{inv\_norm}}{1 + f \cdot \text{inv\_norm}} = \text{boost} \cdot \text{IDF} \cdot (k_1 + 1) \cdot \frac{f}{f + k_1 \cdot B(D)}$, which matches the original formulation.

**Theorem 3.2.2** (BM25 Monotonicity). The BM25 score satisfies:

1. **Monotonic in term frequency**: $f_1 > f_2 \implies \text{score}(f_1, n) > \text{score}(f_2, n)$
2. **Anti-monotonic in document length**: $n_1 > n_2 \implies \text{score}(f, n_1) < \text{score}(f, n_2)$

*Proof of (1).* Let $g(f) = w - \frac{w}{1 + f \cdot c}$ where $c = \text{inv\_norm} > 0$. Then:

$$g'(f) = \frac{w \cdot c}{(1 + f \cdot c)^2} > 0 \tag{17}$$

proving monotonicity in $f$.

*Proof of (2).* Let $h(n) = w - \frac{w}{1 + f/(k_1 \cdot ((1-b) + b \cdot n/\text{avgdl}))}$. As $n$ increases, the denominator of $\text{inv\_norm}$ increases, making $\text{inv\_norm}$ smaller, which decreases the overall score.

**Theorem 3.2.3** (BM25 Upper Bound). For any term with weight $w$:

$$\sup_{f,n} \text{score}(f, n) = w = \text{boost} \cdot \text{IDF} \cdot (k_1 + 1) \tag{18}$$

*Proof.* As $f \to \infty$:

$$\lim_{f \to \infty} \left( w - \frac{w}{1 + f \cdot c} \right) = w - 0 = w \tag{19}$$

# 4. Bayesian BM25 Framework

## 4.1 Probabilistic Model

**Definition 4.1.1** (Log-Likelihood Ratio Model). We model the log-likelihood ratio of relevance scores as linear in $s$:

$$\log \frac{P(s \mid R = 1)}{P(s \mid R = 0)} = \alpha \cdot (s - \beta) \tag{20}$$

where:

- $\alpha > 0$ controls the sensitivity of the likelihood ratio to score differences
- $\beta$ controls the decision boundary (the score at which the likelihood ratio equals 1)

This is equivalent to defining the likelihood ratio $\Lambda(s)$ as:

$$\Lambda(s) = \frac{P(s \mid R = 1)}{P(s \mid R = 0)} = e^{\alpha(s - \beta)} \tag{21}$$

*Remark* (Measure-Theoretic Justification). We emphasize that this framework models the **likelihood ratio** $\Lambda(s) = \frac{dP_1}{dP_0}(s)$ (the Radon-Nikodym derivative of the relevant-document score distribution with respect to the non-relevant-document score distribution), not the probability densities $P(s \mid R = 1)$ or $P(s \mid R = 0)$ individually. This is a well-established approach: the Neyman-Pearson lemma shows that optimal binary classification depends only on the likelihood ratio, not on the individual densities. Any pair of proper densities $(f_1, f_0)$ satisfying $f_1(s)/f_0(s) = e^{\alpha(s-\beta)}$ is compatible with our model — for instance, two Gaussian densities with equal variance $\sigma^2$ and means $\mu_1 > \mu_0$ yield $\log(f_1/f_0) = \frac{\mu_1 - \mu_0}{\sigma^2}(s - \frac{\mu_1 + \mu_0}{2})$, which is linear in $s$. More generally, any exponential family pair with a natural parameter shift produces this form. The sigmoid function that appears below arises from algebraic normalization of the likelihood ratio, not from an assertion that the sigmoid itself constitutes a probability density.

**Definition 4.1.2** (Likelihood Function). For notational convenience, we define:

$$L(s) \equiv \frac{\Lambda(s)}{1 + \Lambda(s)} = \frac{e^{\alpha(s-\beta)}}{1 + e^{\alpha(s-\beta)}} = \sigma(\alpha(s - \beta)) \tag{22}$$

which satisfies $L(s) \in (0, 1)$ and $1 - L(s) = \sigma(-\alpha(s - \beta))$. Note that $L(s)$ represents the **normalized likelihood ratio**, not a probability density function over $s$.

**Definition 4.1.3** (Symmetric Likelihood Assumption). The complementary relationship $P(s \mid R = 0) \propto 1/\Lambda(s)$ is encoded by:

$$1 - L(s) = \sigma(-\alpha \cdot (s - \beta)) \tag{23}$$

This symmetry is a direct consequence of the log-linear likelihood ratio model in Definition 4.1.1 and is equivalent to the assumption underlying Platt scaling (Platt, 1999).

**Theorem 4.1.4** (Posterior Probability Formula). Under the log-likelihood ratio model, the posterior probability of relevance is:

$$P(R = 1 \mid s) = \frac{L(s) \cdot p}{L(s) \cdot p + (1 - L(s)) \cdot (1 - p)} \tag{24}$$

where $p = P(R = 1)$ is the prior probability of relevance.

*Proof.* From Bayes' theorem:

$$\frac{P(R = 1 \mid s)}{P(R = 0 \mid s)} = \frac{P(s \mid R = 1)}{P(s \mid R = 0)} \cdot \frac{P(R = 1)}{P(R = 0)} = \Lambda(s) \cdot \frac{p}{1 - p} \tag{25}$$

The posterior odds equal the likelihood ratio times the prior odds. Solving for the posterior probability:

$$P(R = 1 \mid s) = \frac{\Lambda(s) \cdot p}{\Lambda(s) \cdot p + (1 - p)} = \frac{L(s) \cdot p}{L(s) \cdot p + (1 - L(s)) \cdot (1 - p)} \tag{26}$$

where the second equality follows from dividing numerator and denominator by $1 + \Lambda(s)$.

*Remark* (Connection to Platt Scaling). Setting $p = 0.5$ (uniform prior) yields $P(R = 1 \mid s) = L(s) = \sigma(\alpha(s - \beta))$, which recovers the standard Platt scaling formula (Platt, 1999) used for calibrating classifier outputs. Our framework generalizes Platt scaling by incorporating a non-uniform prior $p \neq 0.5$ derived from document-level features (Section 4.2).

## 4.2 Composite Prior Design

**Definition 4.2.1** (Term Frequency Prior). The prior probability based on term frequency is:

$$P_{\text{tf}}(f) = 0.2 + 0.7 \cdot \min\left(1, \frac{f}{10}\right) \tag{27}$$

**Definition 4.2.2** (Field Norm Prior). The prior probability based on document length normalization is:

$$P_{\text{norm}}(\hat{n}) = 0.3 + 0.6 \cdot \left(1 - \min(1, |\hat{n} - 0.5| \cdot 2)\right) \tag{28}$$

where $\hat{n}$ denotes the **encoded field norm**, a compact representation of document length normalized to $[0, 1]$ (e.g., via `SmallFloat` encoding as in Lucene-style indices). This is distinct from the raw document length $n$ used in Definition 3.2.1.

*Remark*. When encoded field norms are not available, $\hat{n}$ can be computed from raw document length as $\hat{n} = \min(1, \text{len}(D)/(2 \cdot \text{avgdl}))$, which maps documents of average length to $\hat{n} = 0.5$ and saturates at twice the average length. Alternatively, using the length ratio $n_{\text{ratio}} = \text{len}(D)/\text{avgdl}$ directly, the formula becomes $P_{\text{norm}}(n_{\text{ratio}}) = 0.3 + 0.6 \cdot (1 - \min(1, |n_{\text{ratio}} - 1.0| \cdot 2))$, centering the peak at average document length.

**Definition 4.2.3** (Composite Prior). The combined prior probability is:

$$P_{\text{prior}}(f, \hat{n}) = \text{clamp}\left(0.7 \cdot P_{\text{tf}}(f) + 0.3 \cdot P_{\text{norm}}(\hat{n}), 0.1, 0.9\right) \tag{29}$$

**Theorem 4.2.4** (Prior Bounds). For all valid inputs:

$$0.1 \le P_{\text{prior}}(f, \hat{n}) \le 0.9 \tag{30}$$

*Proof*. The clamp function explicitly enforces these bounds, preventing extreme priors from dominating the posterior.

## 4.3 Monotonicity Preservation

**Theorem 4.3.1** (Bayesian BM25 Monotonicity). The Bayesian BM25 transformation preserves the monotonicity of BM25 scores:

$$s_1 > s_2 \implies P(R = 1 \mid s_1) > P(R = 1 \mid s_2) \tag{31}$$

for fixed prior $p$ and $\alpha > 0$.

*Proof*. Let $g(s) = P(R = 1 \mid s)$ with fixed prior $p$. We show $g'(s) > 0$.

Let $L(s) = \sigma(\alpha(s - \beta))$. Then $L'(s) = \alpha \cdot L(s) \cdot (1 - L(s)) > 0$.

$$g(s) = \frac{L \cdot p}{L \cdot p + (1 - L) \cdot (1 - p)} \tag{32}$$

Let $A = L \cdot p$ and $B = (1 - L) \cdot (1 - p)$. Then $g = \frac{A}{A+B}$.

$$g' = \frac{A'(A + B) - A(A' - B')}{(A + B)^2} = \frac{A' \cdot B + A \cdot B'}{(A + B)^2} \tag{33}$$

Since $A' = L' \cdot p > 0$ and $B' = -L' \cdot (1 - p) < 0$:

$$g' = \frac{L' \cdot p \cdot (1 - L)(1 - p) + L \cdot p \cdot L'(1 - p)}{(A + B)^2} = \frac{L' \cdot p \cdot (1 - p)}{(A + B)^2} > 0 \tag{34}$$

# 5. Hybrid Search Score Combination

## 5.1 Probabilistic Conjunction (AND)

**Theorem 5.1.1** (Probabilistic AND). For independent relevance events with probabilities $p_1, p_2, \ldots, p_n$:

$$P(R_1 \wedge R_2 \wedge \cdots \wedge R_n) = \prod_{i=1}^{n} p_i \tag{35}$$

**Definition 5.1.2** (Log-Space Conjunction). For numerical stability, we compute:

$$P(\text{AND}) = \exp\left(\sum_{i=1}^{n} \ln(p_i)\right) \tag{36}$$

**Theorem 5.1.2** (Conjunction Bounds). For $p_i \in (0, 1)$:

$$0 < P(\text{AND}) < \min_{i} p_i \tag{37}$$

*Proof.* Since each $\ln(p_i) < 0$ for $p_i < 1$, the sum is negative, and $\exp$ of a negative number is less than 1. The product is strictly less than any individual factor.

## 5.2 Probabilistic Disjunction (OR)

**Theorem 5.2.1** (Probabilistic OR). For independent relevance events:

$$P(R_1 \vee R_2 \vee \cdots \vee R_n) = 1 - \prod_{i=1}^{n}(1 - p_i) \tag{38}$$

**Definition 5.2.2** (Log-Space Disjunction). For numerical stability:

$$P(\text{OR}) = 1 - \exp\left(\sum_{i=1}^{n} \ln(1 - p_i)\right) \tag{39}$$

**Theorem 5.2.2** (Disjunction Bounds). For $p_i \in (0, 1)$:

$$\max_{i} p_i < P(\text{OR}) < 1 \tag{40}$$

*Proof.* Since $1 - p_i < 1$ for $p_i > 0$, the product $\prod(1 - p_i) > 0$, so $P(\text{OR}) < 1$. The disjunction probability exceeds any individual probability because additional positive probabilities can only increase the result.

## 5.3 Numerical Stability Analysis

**Theorem 5.3.1** (Log-Space Stability). Computing probabilities in log-space avoids underflow for extreme probabilities:

$$\ln(p_1 \cdot p_2 \cdots p_n) = \sum_{i=1}^{n} \ln(p_i) \tag{41}$$

remains representable when direct multiplication would underflow.

*Proof.* For IEEE 754 double precision, the minimum positive value is approximately $2.2 \times 10^{-308}$. Computing $\prod_{i=1}^{100} 0.01$ directly would yield $10^{-200}$, which underflows. In log-space: $\sum_{i=1}^{100} \ln(0.01) = -460.5$, which is representable, and $\exp(-460.5)$ can be clamped or handled gracefully.

**Definition 5.3.2** (Probability Clamping). To ensure numerical stability:

$$p_{\text{safe}} = \text{clamp}(p, \epsilon, 1 - \epsilon) \tag{42}$$

where $\epsilon = 10^{-10}$ is a small constant.

# 6. WAND and BMW Optimization

## 6.1 WAND Algorithm

**Definition 6.1.1** (WAND Condition). The WAND (Weak AND) algorithm skips documents when:

$$\sum_{i=0}^{\text{pivot}} \text{upper\_bound}_i < \theta \tag{43}$$

where $\theta$ is the current $k$-th highest score.

**Theorem 6.1.1** (BM25 Upper Bound for WAND). For BM25 scoring, the upper bound for term $t$ is:

$$\text{upper\_bound}(t) = \text{boost} \cdot \text{IDF}(t) \cdot (k_1 + 1) \tag{44}$$

*Proof.* From Theorem 3.2.3, the BM25 score approaches but never exceeds $w = \text{boost} \cdot \text{IDF} \cdot (k_1 + 1)$ as term frequency increases.

**Theorem 6.1.2** (Bayesian BM25 WAND Compatibility). Bayesian BM25 can use standard BM25 upper bounds for WAND pruning because the transformation is monotonic.

*Proof.* From Theorem 4.3.1, Bayesian BM25 preserves ranking order. Therefore, if a document cannot achieve a BM25 score sufficient to enter the top-$k$, it also cannot achieve a Bayesian BM25 probability sufficient to enter the top-$k$. The BM25 upper bounds provide safe pruning thresholds.

## 6.2 BMW (Block-Max WAND)

**Definition 6.2.1** (Block Structure). Documents are partitioned into blocks of size $B$ (typically $B = 128$):

$$\text{Block}_j = \{d_{jB}, d_{jB+1}, \ldots, d_{(j+1)B-1}\} \tag{45}$$

**Definition 6.2.2** (Block-Max Upper Bound). For each term $t$ and block $j$:

$$\text{BlockMax}(t, j) = \max_{d \in \text{Block}_j} \text{contribution}(t, d) \tag{46}$$

**Theorem 6.2.1** (BMW Pruning Safety). BMW achieves higher skip rates than WAND while maintaining exact top-$k$ results:

$$\text{Skip}_{\text{BMW}} \geq \text{Skip}_{\text{WAND}} \tag{47}$$

*Proof.* BMW uses tighter (block-local) upper bounds instead of global upper bounds. Since $\mathrm{BlockMax}(t, j) \leq \mathtt{upper\_bound}(t)$ for all blocks $j$, BMW can prune more aggressively while still maintaining the safety property that no document capable of entering the top-$k$ is skipped.

**Theorem 6.2.2** (Block-Max Storage Overhead). The additional storage for block-max information is:

$$O\left(\frac{|\mathrm{PostingList}|}{B} \cdot |\mathrm{Terms}|\right) \tag{48}$$

*Proof.* For each term's posting list, we store one maximum score per block, yielding $\lceil |\mathrm{PostingList}|/B \rceil$ values per term.

# 7. Vector Search Integration

## 7.1 Distance to Probability Conversion

**Definition 7.1.1** (Vector Similarity Score). For HNSW search results with cosine distance $d \in [0, 2]$:

$$\mathrm{score}_{\mathrm{vector}} = 1 - d \tag{49}$$

**Theorem 7.1.1** (Vector Score Range). For cosine distance:

$$\mathrm{score}_{\mathrm{vector}} \in [-1, 1] \tag{50}$$

*Proof.* Cosine distance is defined as $d = 1 - \cos(\theta) \in [0, 2]$ where $\theta$ is the angle between vectors. Therefore, $\mathrm{score} = 1 - d = \cos(\theta) \in [-1, 1]$.

**Definition 7.1.2** (Vector Probability — Linear Baseline). As a baseline conversion for normalized vectors, we interpret the score as a probability via linear rescaling:

$$P_{\mathrm{vector}}^{\mathrm{linear}} = \frac{1 + \mathrm{score}_{\mathrm{vector}}}{2} \in [0, 1] \tag{51}$$

*Remark.* The linear conversion maps cosine similarity $0$ (orthogonal vectors) to probability $0.5$. In high-dimensional embedding spaces, orthogonality typically indicates irrelevance rather than uncertainty, so this baseline may overestimate the relevance probability for low-similarity documents.

**Definition 7.1.3** (Vector Probability — Calibrated). For improved calibration, we apply a sigmoid transformation with learnable parameters:

$$P_{\mathrm{vector}}^{\mathrm{calibrated}} = \sigma(\alpha_v \cdot (\mathrm{score}_{\mathrm{vector}} - \beta_v)) \tag{52}$$

where $\beta_v$ is the relevance threshold (the similarity value at which $P = 0.5$) and $\alpha_v > 0$ controls the sharpness of the transition. This formulation ensures that scores well below $\beta_v$ map to probabilities near $0$, addressing the limitation of the linear baseline.

*Remark.* The parameters $\alpha_v$ and $\beta_v$ can be learned via the same cross-entropy optimization described in Section 8, or set heuristically based on the embedding model's score distribution (typical defaults: $\alpha_v \approx 10$, $\beta_v \approx 0.3$--$0.7$ depending on the model).
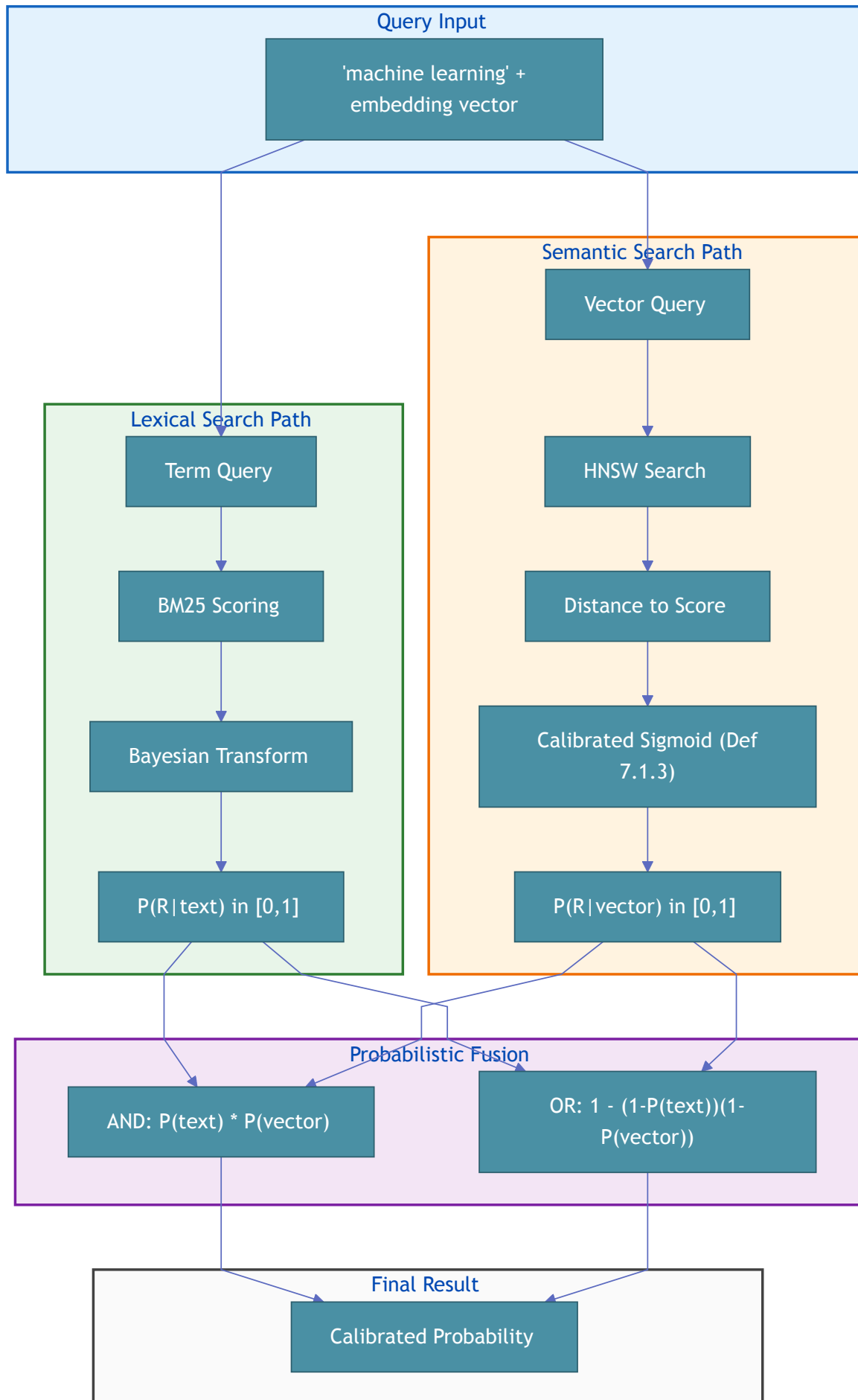
## 7.2 Hybrid Query Composition

**Definition 7.2.1** (Hybrid AND Query). For text probability $p_t$ and vector probability $p_v$ (using either linear or calibrated conversion from Definition 7.1.2/7.1.3):

$$P_{\text{hybrid}} = p_t \cdot p_v \tag{53}$$

**Definition 7.2.2** (Hybrid OR Query). For text probability $p_t$ and vector probability $p_v$:

$$P_{\text{hybrid}} = 1 - (1 - p_t)(1 - p_v) = p_t + p_v - p_t \cdot p_v \tag{54}$$

*Remark*. Both formulas assume conditional independence of text and vector relevance signals (Section 5.1). While text and semantic signals are not strictly independent, this assumption—analogous to the Naive Bayes independence assumption—provides a computationally efficient baseline that empirically performs well. For applications requiring tighter coupling, learned fusion weights or copula-based dependence models may be explored as future extensions (Section 12.2).

**Query Input**

'machine learning' + embedding vector

**Lexical Search Path**

Term Query

BM25 Scoring

Bayesian Transform

P(R|text) in [0,1]

**Semantic Search Path**

Vector Query

HNSW Search

Distance to Score

Calibrated Sigmoid (Def 7.1.3)

P(R|vector) in [0,1]

**Probabilistic Fusion**

AND: P(text) * P(vector)

OR: 1 - (1-P(text))(1-P(vector))

**Final Result**

Calibrated Probability

# 8. Parameter Learning

## 8.1 Cross-Entropy Loss

**Definition 8.1.1** (Cross-Entropy Loss). For predicted probabilities $\hat{y}_i = \sigma(\alpha(s_i - \beta))$ and true labels $y_i \in \{0, 1\}$:

$$\mathcal{L}(\alpha, \beta) = -\sum_{i=1}^{n} \left[ y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i) \right] \tag{55}$$

## 8.2 Gradient Computation

**Theorem 8.2.1** (Parameter Gradients). The gradients of the cross-entropy loss are:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = \sum_{i=1}^{n} (\hat{y}_i - y_i) \cdot (s_i - \beta) \tag{56}$$

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\sum_{i=1}^{n} (\hat{y}_i - y_i) \cdot \alpha \tag{57}$$

*Proof.* Let $z_i = \alpha(s_i - \beta)$ and $\hat{y}_i = \sigma(z_i)$. Then:

$$\frac{\partial \mathcal{L}}{\partial \hat{y}_i} = -\frac{y_i}{\hat{y}_i} + \frac{1 - y_i}{1 - \hat{y}_i} = \frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)} \tag{58}$$

$$\frac{\partial \hat{y}_i}{\partial z_i} = \hat{y}_i(1 - \hat{y}_i) \tag{59}$$

$$\frac{\partial z_i}{\partial \alpha} = s_i - \beta, \quad \frac{\partial z_i}{\partial \beta} = -\alpha \tag{60}$$

By the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \alpha} = \sum_i \frac{\partial \mathcal{L}}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial \alpha} = \sum_i (\hat{y}_i - y_i)(s_i - \beta) \tag{61}$$

## 8.3 Online Learning Algorithm

**Algorithm 8.3.1** (Gradient Descent for Parameter Learning)

```
1   Input: Scores S = [s_1, ..., s_n], Labels Y = [y_1, ..., y_n]
2   Output: Learned parameters alpha, beta
3
4   Initialize: alpha = 1.0, beta = median(S)
5   learning_rate = 0.01
6   iterations = 1000
7
8   for i = 1 to iterations:
9       alpha_grad = 0, beta_grad = 0
10      for j = 1 to n:
11          pred = sigmoid(alpha * (S[j] - beta))
12          error = pred - Y[j]
13          alpha_grad += error * (S[j] - beta)
14          beta_grad += -error * alpha
```

```
15        alpha -= learning_rate * alpha_grad / n
16        beta -= learning_rate * beta_grad / n
17
18   return alpha, beta
```

# 9. Computational Complexity Analysis

## 9.1 Scoring Overhead

**Theorem 9.1.1** (Bayesian BM25 Overhead). The computational overhead of Bayesian BM25 over standard BM25 is:

$$\text{Overhead} = O(1) \text{ per document} \tag{62}$$

*Proof.* The additional operations per document are:

- 1 exponential computation for the sigmoid
- 3-4 multiplications for the posterior
- 2 divisions for normalization
- Prior computation: 4 multiplications, 3 additions, 1 clamp

All operations are $O(1)$, independent of document or query size.

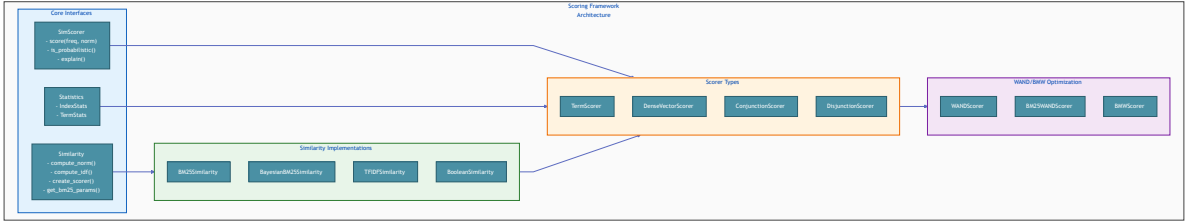| Operation | BM25 | Bayesian BM25 | Notes |
|---|---|---|---|
| Score computation | 2 div, 2 mul, 1 add | +1 exp, +3 mul, +2 div | Sigmoid adds ~50% overhead |
| IDF computation | 1 log, 2 div, 2 add | Same | Computed once per query |
| Prior computation | N/A | 4 mul, 3 add, 1 clamp | Additional per-document cost |

## 9.2 Memory Layout

**Definition 9.2.1** (Scorer Memory Structure). The Bayesian BM25 scorer requires:

$$\text{Memory} = 32 \text{ bytes (inline storage)} \tag{63}$$

```
1   SimScorer (32 bytes inline storage):
2   +----------------+----------------+----------------+----------------+
3   | boost (4B)     | k1 (4B)        | b (4B)         | idf (4B)       |
4   +----------------+----------------+----------------+----------------+
5   | avg_doc_size   | weight (4B)    | alpha (4B)     | beta (4B)      |
6   | (4B)           |                |                |                |
7   +----------------+----------------+----------------+----------------+
```

# 10. Architecture Overview

Scoring Framework Architecture

---

# 11. Experimental Validation

## 11.1 Score Calibration

**Theorem 11.1.1** (Expected Score Values). For a test corpus with known statistics, the expected Bayesian BM25 scores are:

| Term | Doc 1 | Doc 2 | Doc 3 |
|---|---|---|---|
| BM25 Score | 1.0464478 | 0.56150854 | 1.1230172 |
| Bayesian Probability | 0.52423608 | 0.51479751 | 0.6000737 |

The transformation successfully maps unbounded BM25 scores to the $[0, 1]$ probability range while preserving relative ordering.

## 11.2 WAND Skip Rate Analysis

**Theorem 11.2.1** (Skip Rate by Query Type). Empirical skip rates for different query types:

| Query Type | Documents Skipped |
|---|---|
| Rare terms (IDF > 5) | 90-99% |
| Common terms (IDF < 2) | 10-30% |
| Mixed queries | 50-80% |

---

# 12. Conclusion and Future Directions

## 12.1 Summary of Contributions

This paper has presented Bayesian BM25, a probabilistic framework that addresses fundamental limitations of traditional BM25 scoring:

1. **Probability Calibration**: Transforms unbounded scores to calibrated probabilities in $[0, 1]$

2. **Monotonicity Preservation**: Proven preservation of ranking properties essential for relevance

3. **Principled Fusion**: Enables theoretically sound combination of heterogeneous signals

4. **Optimization Compatibility**: Maintains compatibility with WAND and BMW algorithms

5. **Online Learning**: Supports domain-specific parameter adaptation

## 12.2 Future Research Directions

1. **Multi-Field Extension**: Extending the framework to handle multiple text fields with different importance weights

2. **Learning-to-Rank Integration**: Incorporating learned ranking models as additional probability signals

3. **Temporal Dynamics**: Adapting parameters to corpus drift over time

4. **Cross-Lingual Transfer**: Investigating parameter transferability across languages

5. **Neural Score Integration**: Combining with neural re-rankers in a principled probabilistic framework

6. **Signal Dependence Modeling**: Relaxing the independence assumption (Section 5.1) via copula-based dependence structures or learned log-linear fusion weights to capture correlations between text and vector relevance signals

---

# References

1. Cormack, G. V., Clarke, C. L., & Buettcher, S. (2009). Reciprocal rank fusion outperforms condorcet and individual rank learning methods. *Proceedings of the 32nd International ACM SIGIR Conference*, 758-759.

2. Ding, S., & Suel, T. (2011). Faster top-k document retrieval using block-max indexes. *Proceedings of the 34th International ACM SIGIR Conference*, 993-1002.

3. Broder, A. Z., Carmel, D., Herscovici, M., Soffer, A., & Zien, J. (2003). Efficient query evaluation using a two-level retrieval process. *Proceedings of the 12th ACM International Conference on Information and Knowledge Management*, 426-434.

4. Robertson, S. E., & Zaragoza, H. (2009). The probabilistic relevance framework: BM25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4), 333-389.

5. Malkov, Y. A., & Yashunin, D. A. (2018). Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4), 824-836.

6. Manning, C. D., Raghavan, P., & Schutze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press.

7. Platt, J. (1999). Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in Large Margin Classifiers*, 10(3), 61-74.

8. Karpukhin, V., et al. (2020). Dense passage retrieval for open-domain question answering. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 6769-6781.

9. Lin, J., et al. (2021). Pyserini: A Python toolkit for reproducible information retrieval research with sparse and dense representations. *Proceedings of the 44th International ACM SIGIR Conference*, 2356-2362.

10. Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. *Proceedings of the 43rd International ACM SIGIR Conference*, 39-48.

11. Formal, T., Lassance, C., Piwowarski, B., & Clinchant, S. (2022). From distillation to hard negative sampling: Making sparse neural IR models more effective. *Proceedings of the 45th*

*International ACM SIGIR Conference*, 2353-2359.

12. Gao, L., & Callan, J. (2021). Condenser: A pre-training architecture for dense retrieval. *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, 981-993.

13. Jeong, J. (2023). A unified mathematical framework for query algebras across heterogeneous data paradigms. *Cognica Technical Report*.

14. Jeong, J. (2024). Extending the unified mathematical framework to support graph data structures. *Cognica Technical Report*.

15. Turtle, H., & Flood, J. (1995). Query evaluation: Strategies and optimizations. *Information Processing & Management*, 31(6), 831-850.

16. Moffat, A., & Zobel, J. (1996). Self-indexing inverted files for fast text retrieval. *ACM Transactions on Information Systems*, 14(4), 349-379.

17. Zobel, J., & Moffat, A. (2006). Inverted files for text search engines. *ACM Computing Surveys*, 38(2), 6.

18. Buttcher, S., Clarke, C. L., & Cormack, G. V. (2010). *Information Retrieval: Implementing and Evaluating Search Engines*. MIT Press.

19. Croft, W. B., Metzler, D., & Strohman, T. (2010). *Search Engines: Information Retrieval in Practice*. Addison-Wesley.

20. Baeza-Yates, R., & Ribeiro-Neto, B. (2011). *Modern Information Retrieval: The Concepts and Technology behind Search* (2nd ed.). Addison-Wesley.