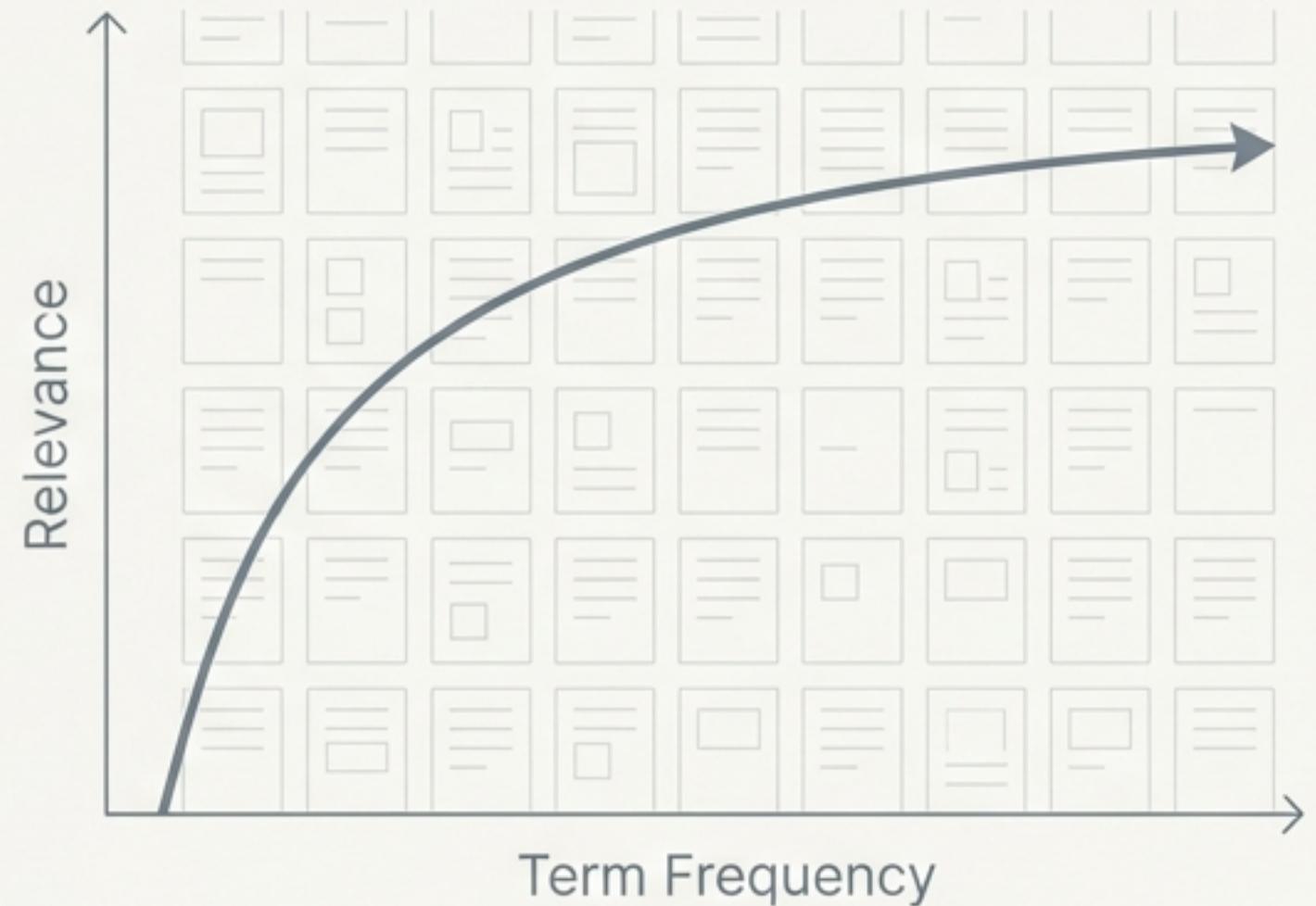


Deconstructing Okapi BM25

*Probabilistic Foundations,
Mathematical Mechanics, and
the Future of Hybrid Search
Search.*

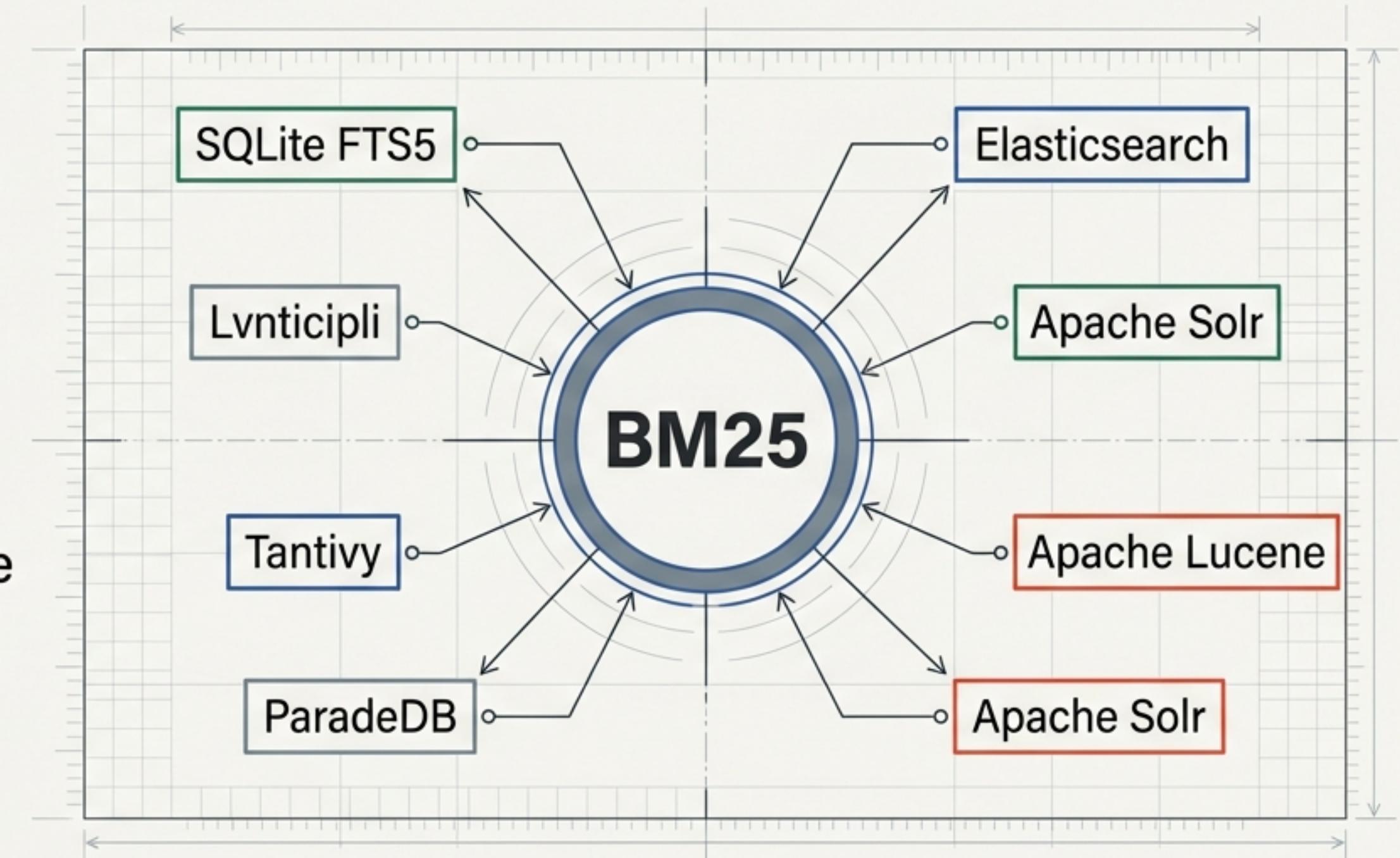


Best Matching version 25 • Developed at City University, London (1980s-90s)

The Enduring Engine of Information Retrieval

BM25 is not a relic; it is the default industry standard. It powers the search infrastructure for the world's largest applications, balancing precision with computational efficiency.

It outperforms older heuristics like TF-IDF by introducing saturation and length normalization.

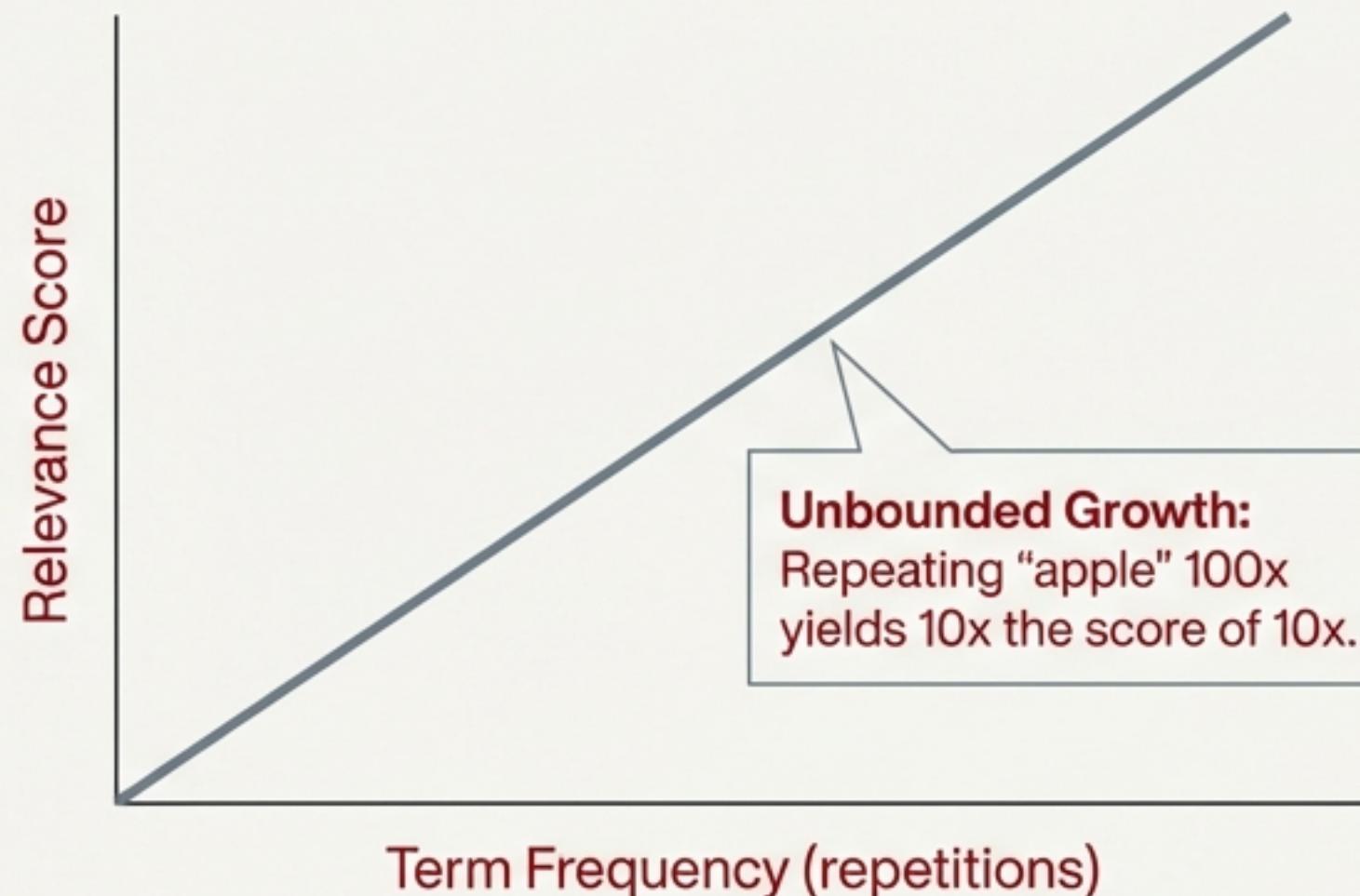


Key Insight: The “Best Matching” function acts as a probabilistic proxy for human relevance judgment.

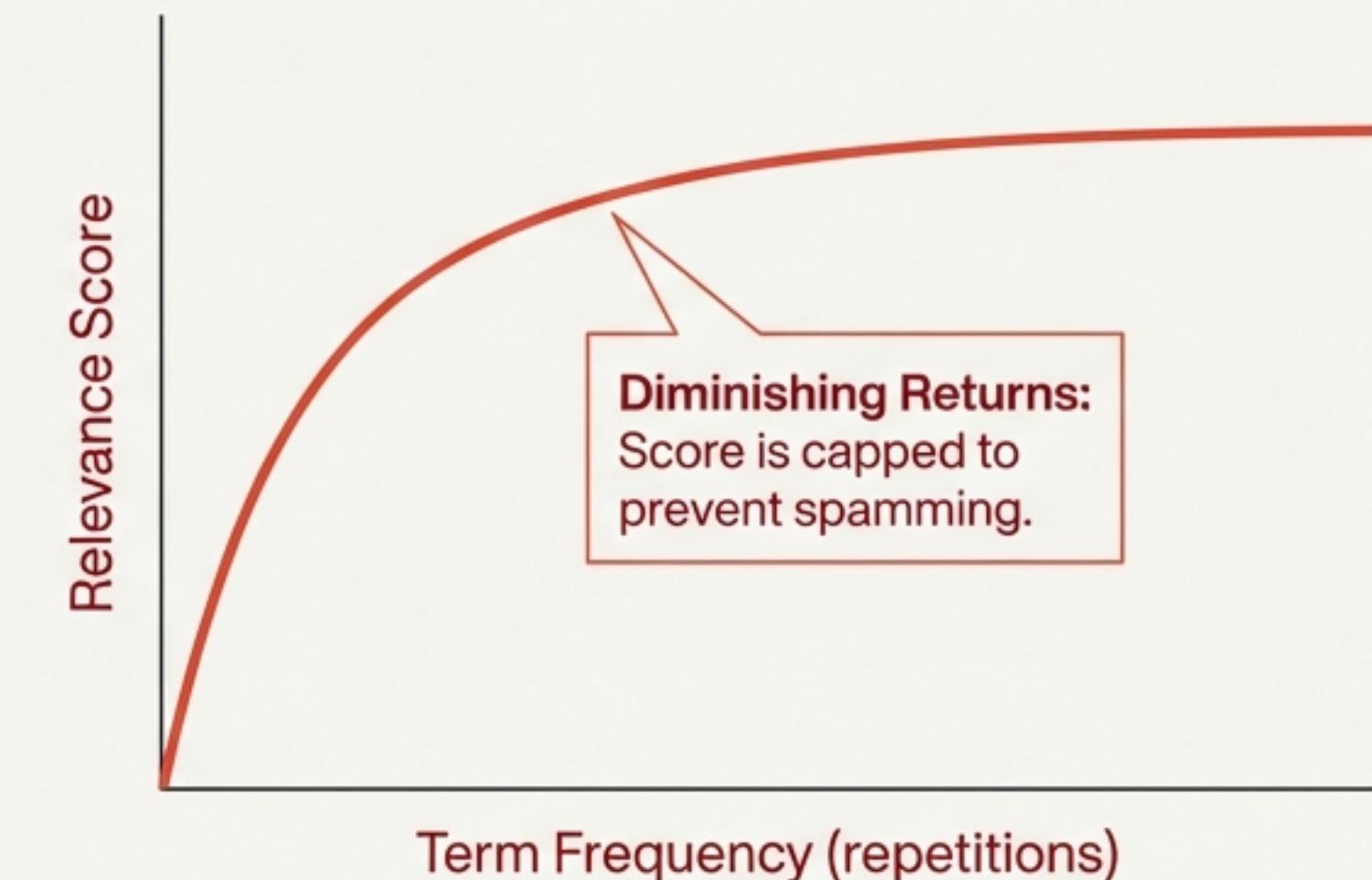
The Predecessor: Why TF-IDF Wasn't Enough

The limitations of heuristic scoring in a complex document landscape.

TF-IDF: Linear Bias



BM25: Saturation



Problem: Raw TF-IDF favors long documents unfairly (Length Bias) and allows keyword stuffing.

The Probabilistic Shift

The Probability Ranking Principle

“Documents should be ranked in order of their probability of relevance to the user’s information need.” — Robertson & Spärck Jones

BM25 estimates $P(R|D, Q)$: The probability that Document D is Relevant given Query Q.

This moves us from simple “counting” to “evidence weighting”.



The Bag-of-Words Assumption: Structure is ignored; only token evidence matters.

The Anatomy of the Equation

$$score(D, Q) = \sum_i IDF(q_i) \cdot \left[\frac{f(q_i, D) \cdot (k_1 + 1)}{f(q_i, D) + k_1 \cdot \left(1 - b + b \cdot \frac{|D|}{avgdl}\right)} \right]$$

Annotations:

- 1. The Rarity Signal (IDF) - Points to the $IDF(q_i)$ term.
- 2. The Saturation Signal (Term Frequency) - Points to the $f(q_i, D)$ term.
- 2. The Saturation Signal (Term Frequency) - Points to the $f(q_i, D)$ term.
- 3. The Fairness Signal (Length Normalization) - Points to the denominator term $\left(1 - b + b \cdot \frac{|D|}{avgdl}\right)$.

Three mechanical levers that translate text into probability.



Signal 1: IDF (The Rarity Engine)

Not all matches are created equal.

Inverse Document Frequency (IDF) measures the discriminative power of a term.

Common words (the, is, and) appear everywhere -> Score approaches 0 (or negative).

Rare words (Okapi, Hydroplane) appear rarely -> High Score.

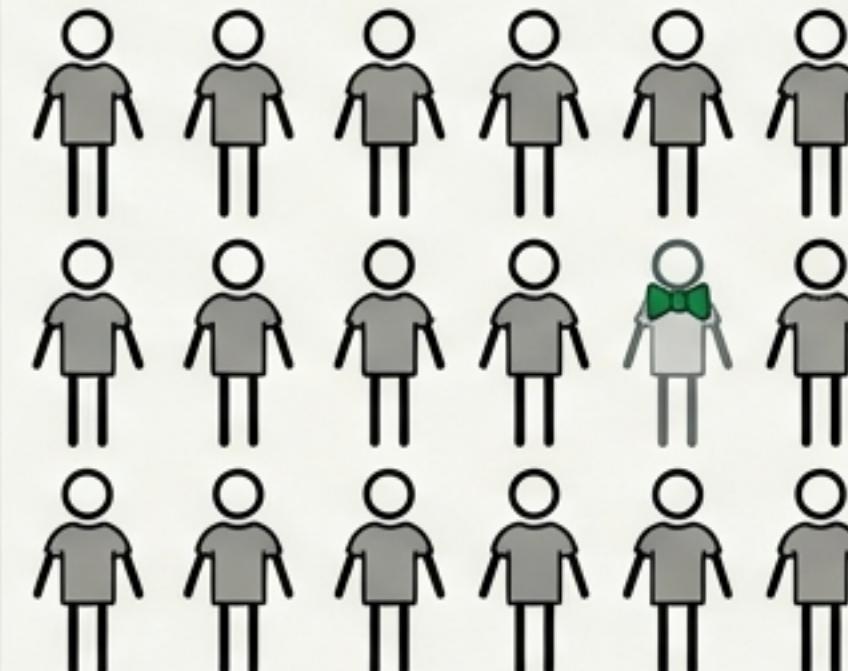
$$\text{IDF}(q) = \ln \left(\frac{N - n(q) + 0.5}{n(q) + 0.5} + 1 \right)$$

The Classroom Search Analogy



Signal Red

Query: "Person in shirt".
Result: Matches Everyone.
(Low Information)

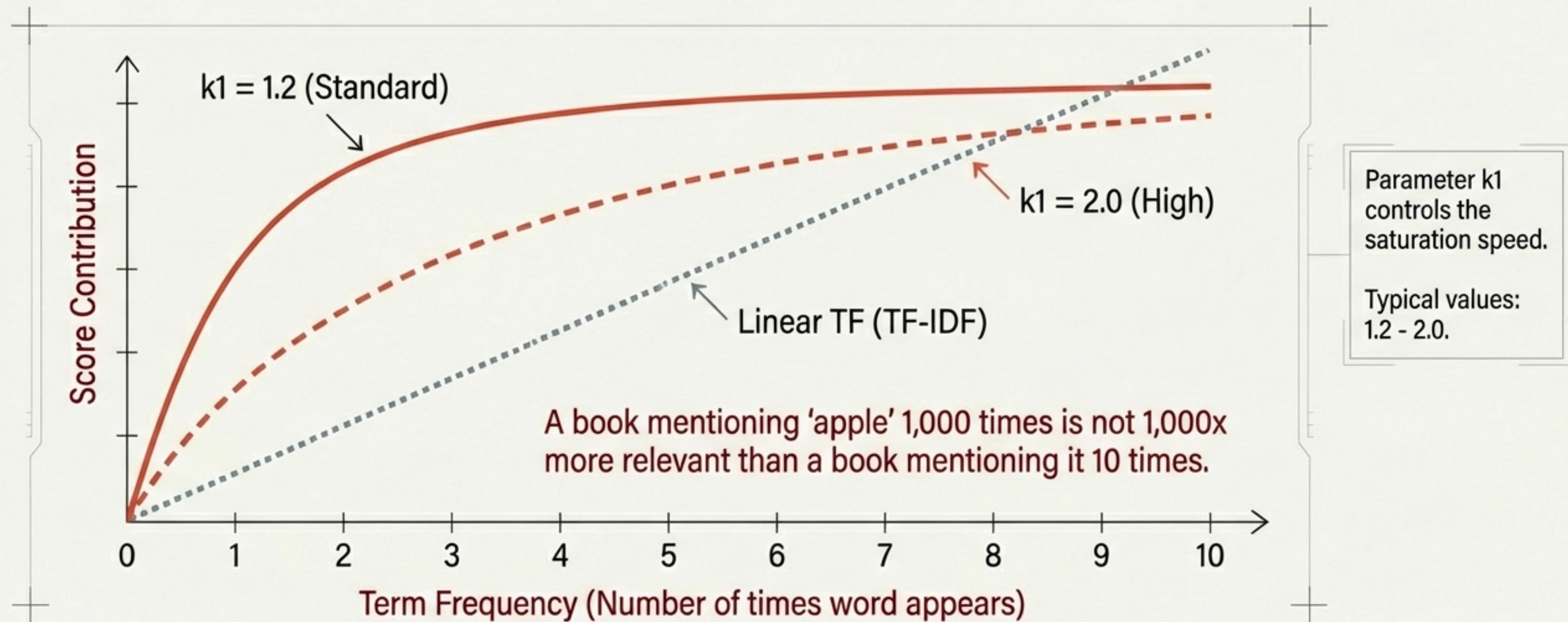


Forest Green

Query: "Neon Green Bowtie".
Result: One Specific Match.
(High Information)

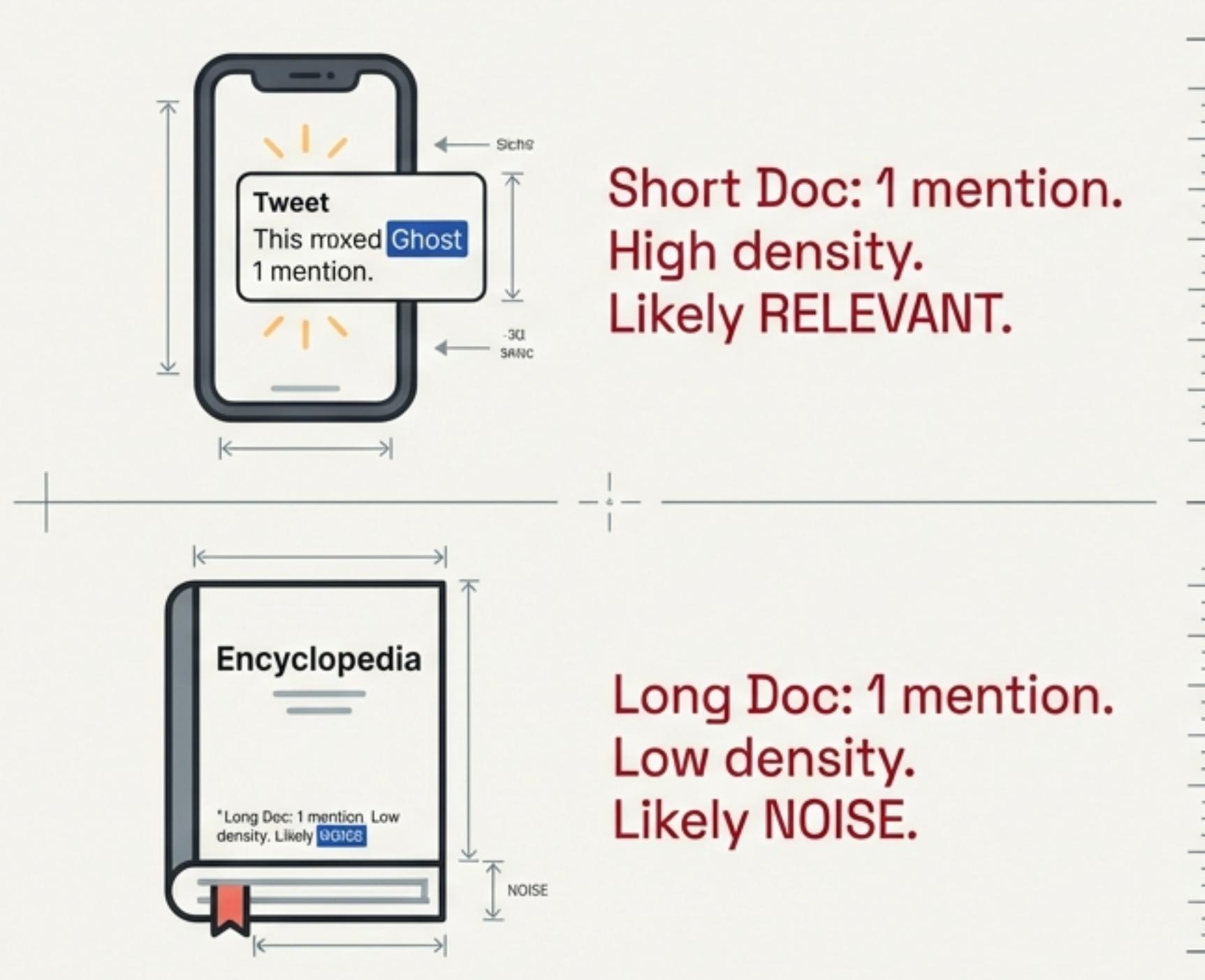
Signal 2: Term Frequency Saturation

The Law of Diminishing Returns, controlled by parameter k_1 .



Signal 3: Length Normalization

Ensuring fairness across document sizes, controlled by parameter b.



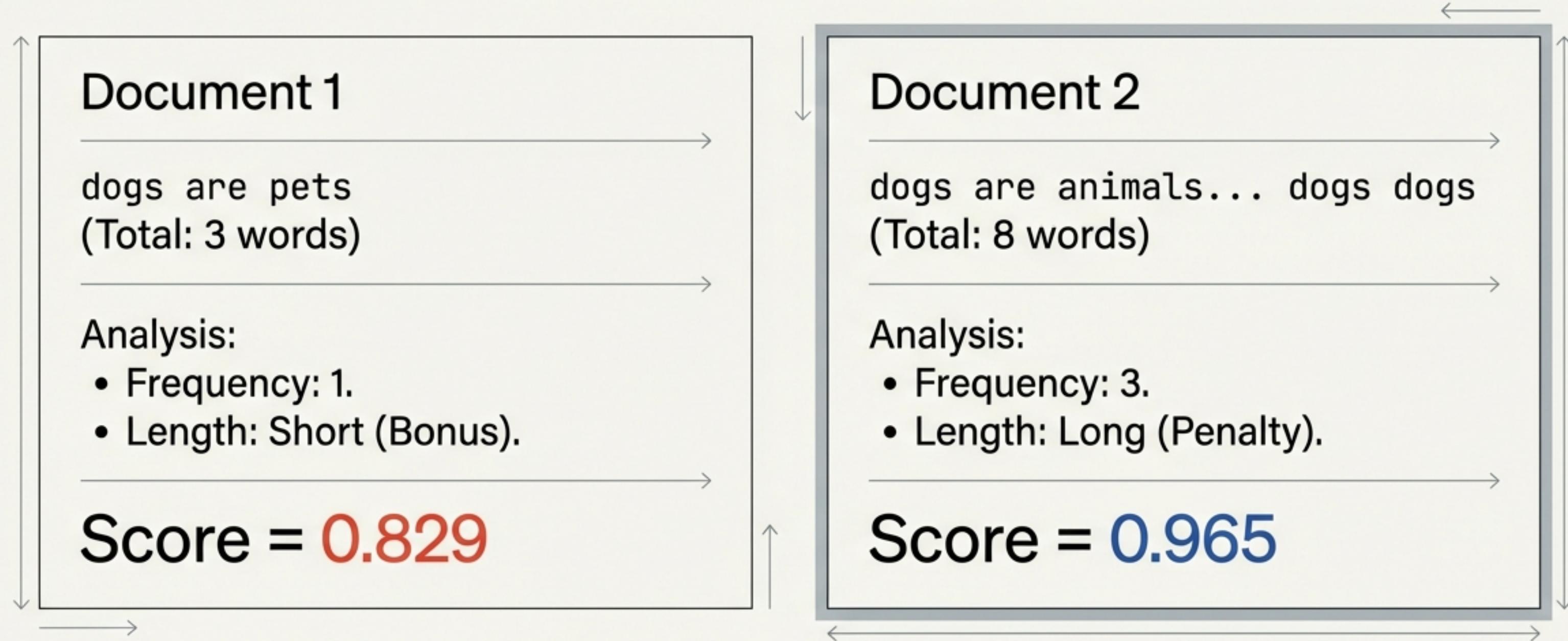
- Parameter b (0 to 1) controls the penalty.
- $b = 1.0$: Full Penalty (Strictly normalizes by length).
- $b = 0.0$: No Penalty (Ignores length, behaves like TF-IDF).
- $b = 0.75$: Industry Standard.

$$\text{Normalization Factor} = 1 - b + b \cdot \frac{|D|}{\text{avgdl}}$$

$|D|$ = document length
 avgdl = average document length

The Math in Action: A Numerical Walkthrough

Query: "dogs". Average Corpus Length: 5 words.



Result: Document 2 wins. The higher frequency (3x) overcomes the length penalty, but saturation keeps the score from being 3x higher (only ~16% higher).

The Showdown: BM25 vs. TF-IDF

→ TF-IDF (Legacy)

- Foundation: Heuristic / Statistical
- Term Frequency: Unbounded (Linear/Log)
- Length Bias: Biased toward long docs
- Tuning: None (Parameter-free)
- Best For: Small, uniform datasets

Okapi BM25 (Standard) ←

- Foundation: Probabilistic Relevance
- Term Frequency: Saturated (Asymptotic limit)
- Length Bias: Normalized (Fair)
- Tuning: Highly Tunable (k_1, b)
- Best For: Web scale, Enterprise, Varied lengths

← Previous slide's design language

← Previous slide's design language

Tuning the Engine: An Industry Guide

Optimizing parameters k_1 and b for specific domains.

Legal / Medical



High k_1 (2.0),
High b (0.8)

Crimson:
Precision is key.
Documents are long.
Repetition is significant.

E-Commerce



Low k_1 (1.0),
Low b (0.4)

Rationale
Descriptions are short.
Prevent keyword
spamming.

General Web



Default k_1 (1.2),
Default b (0.75)

Rationale
Balanced approach for
diverse content types.

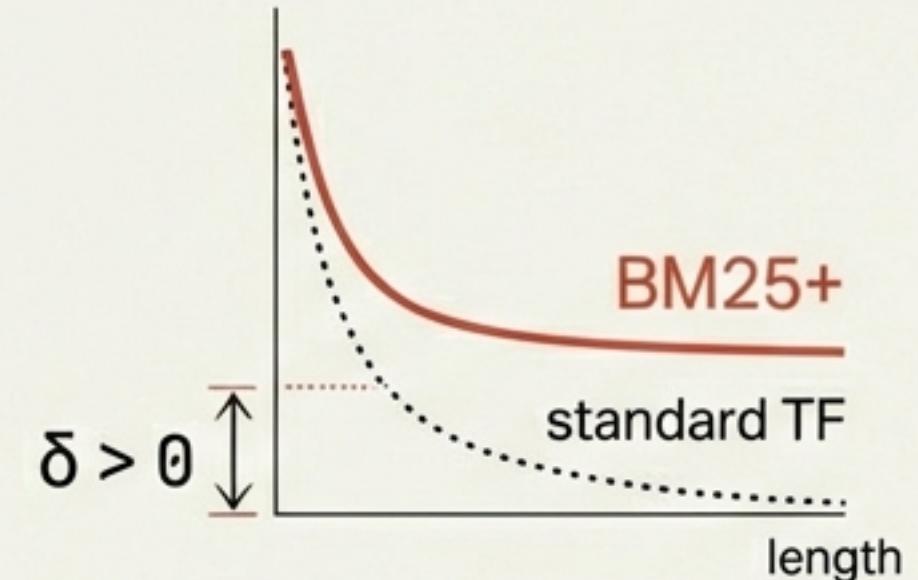
Beyond the Basics: Advanced Variants



BM25+ (The Long Doc Fix)

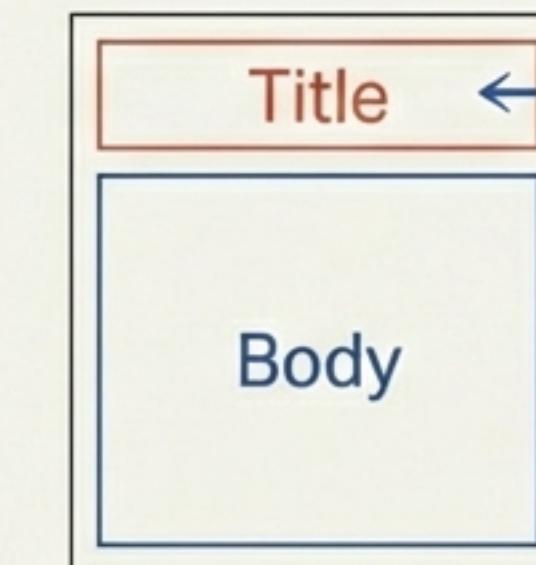
Standard BM25 can over-penalize extremely long documents, driving scores to zero.

- >Adds a lower-bound constant (δ) to the Term Frequency.
Ensures a match always contributes **something**.



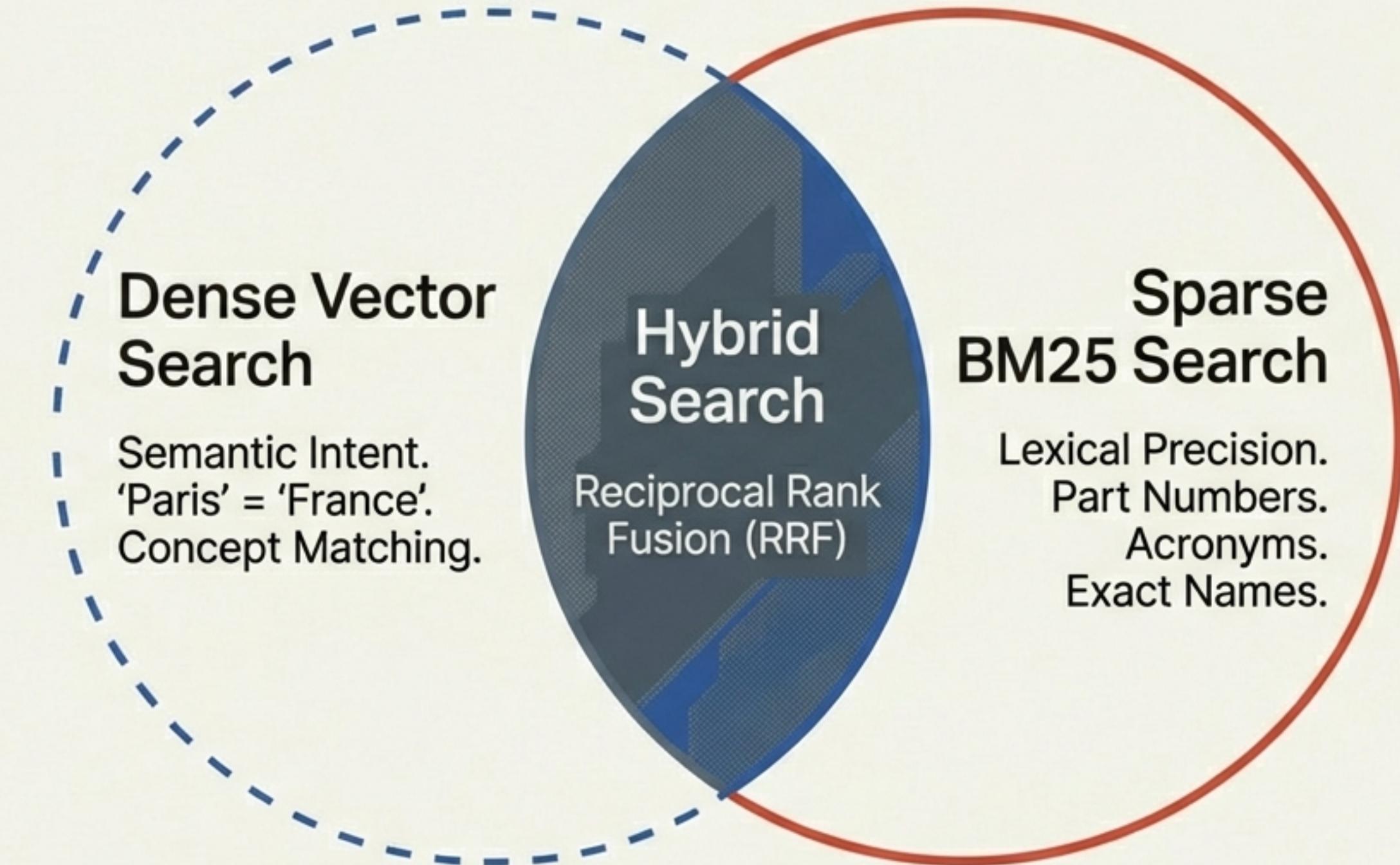
BM25F (Fielded Search)

Documents have structure
(Title, Body, Anchor).



Calculates a weighted average of frequencies across fields before saturation.

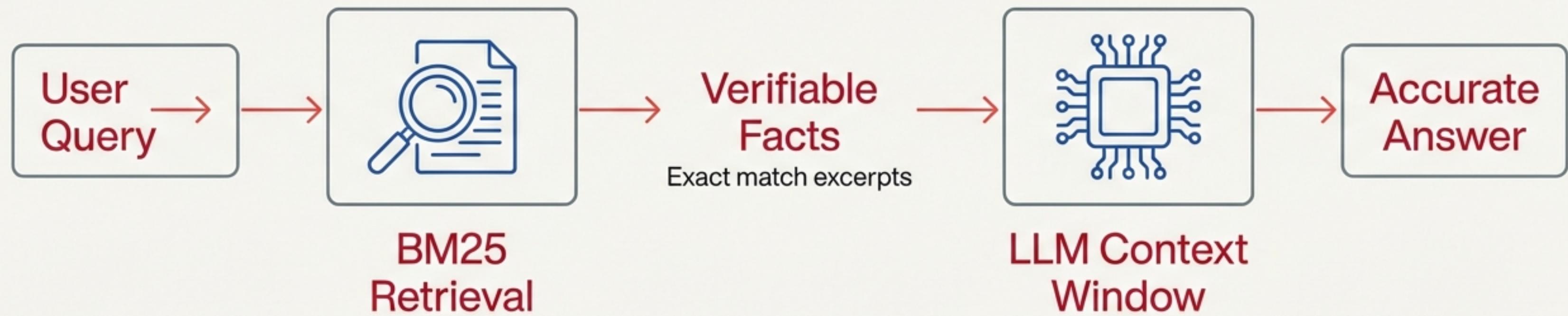
BM25 in the Age of AI: Hybrid Search



“Vectors understand intent. BM25 guarantees precision.”

Grounding AI: Preventing Hallucinations

How BM25 supports Retrieval-Augmented Generation (RAG).



Anthropic's Contextual Retrieval: Combining BM25 with AI-generated context chunks significantly improves retrieval over embeddings alone.

The Verdict: Augmented, Not Replaced

BM25 remains the undisputed champion of explainable, efficient, keyword-based retrieval.

- 1. SATURATION:** Handles repetition intelligently.
- 2. NORMALIZATION:** Ensures fairness across document lengths.
- 3. EFFICIENCY:** Powers the largest search engines and modern RAG pipelines.

“In the era of complex Neural Networks, the mathematically elegant ‘Bag of Words’ remains essential.”