

$G(q)$  behavior:  $\overline{IDF}(q) \ll \overline{\pi}_m$  (common words in query),  $G(q) \approx 0$  and  $R(q^{\text{mic}}, d^{\text{mic}})$  doesn't contribute  
 $\overline{IDF}(q) \gg \overline{\pi}_m$  (rare queries),  $G(q) \approx 1$  and  $R(q^{\text{mic}}, d^{\text{mic}})$  is weighted by  $w_{\text{mic}}$  (0.12)

The double log structure for term weight: additional matched terms do not contribute additively, but provide diminishing returns

$B_{\text{cov}}$ : A document match all query terms gets a  $1 + \gamma_{\text{cov}}$  multiplier (x1.25)  
 match none gets  $\times 1$  multiplier, but in that case,  $E$  is also 0

$\text{PMI}(t, d)$ : kind of like  $\log\left(\frac{P(t|d)}{P(t|c)}\right)$        $P(t|d) = \frac{tf(t, d)}{|d|_{\text{eff}}}$        $P(t|c) = \frac{df(t)}{N}$

$B_{\text{spec}}$ : how much information the observation "t is in d" carries beyond what the corpus prior would predict

$B_{\text{anc}}$ : if no matched term has  $IDF(t) > \overline{\pi}(4.2)$  then  $A=0$  and  $B_{\text{anc}}=1$

Optimal Multiplier implemented by GPT 5.2 but was not used ( $\alpha_{\text{gini}} = 0$ )

$$B_{\text{gini}} = 1 + \alpha_{\text{gini}} \cdot \text{Conc}(q, d)$$

$$\text{Conc}(q, d) = \text{clip}\left(\frac{H - \frac{1}{|M|}}{1 - \frac{1}{|M|}}, 0, 1\right)$$

Conc high  
If document's evidence is concentrated on a few terms,  
(topical focus?)  
Conc low if evidence spreads evenly

$$H = \sum_{t \in M} (w(t) \cdot tf(t, d))^2$$

$$\left( \sum_{t \in M} (w(t) \cdot tf(t, d)) \right)^2$$

Bigram: [neural, network, model]  $\rightarrow$  [neural-network, network-model]

Tokenization Channels:

base: Lucene tokenizer

prefix: neural  $\rightarrow$  neura, netwo

(only keep first  $L_p$  characters  $L_p=5$ )

Micro: [neural]  $\rightarrow$  [neu, eur, ura, ral]

(all overlapping  $L_m$ -grams,  $L_m=3$ )