

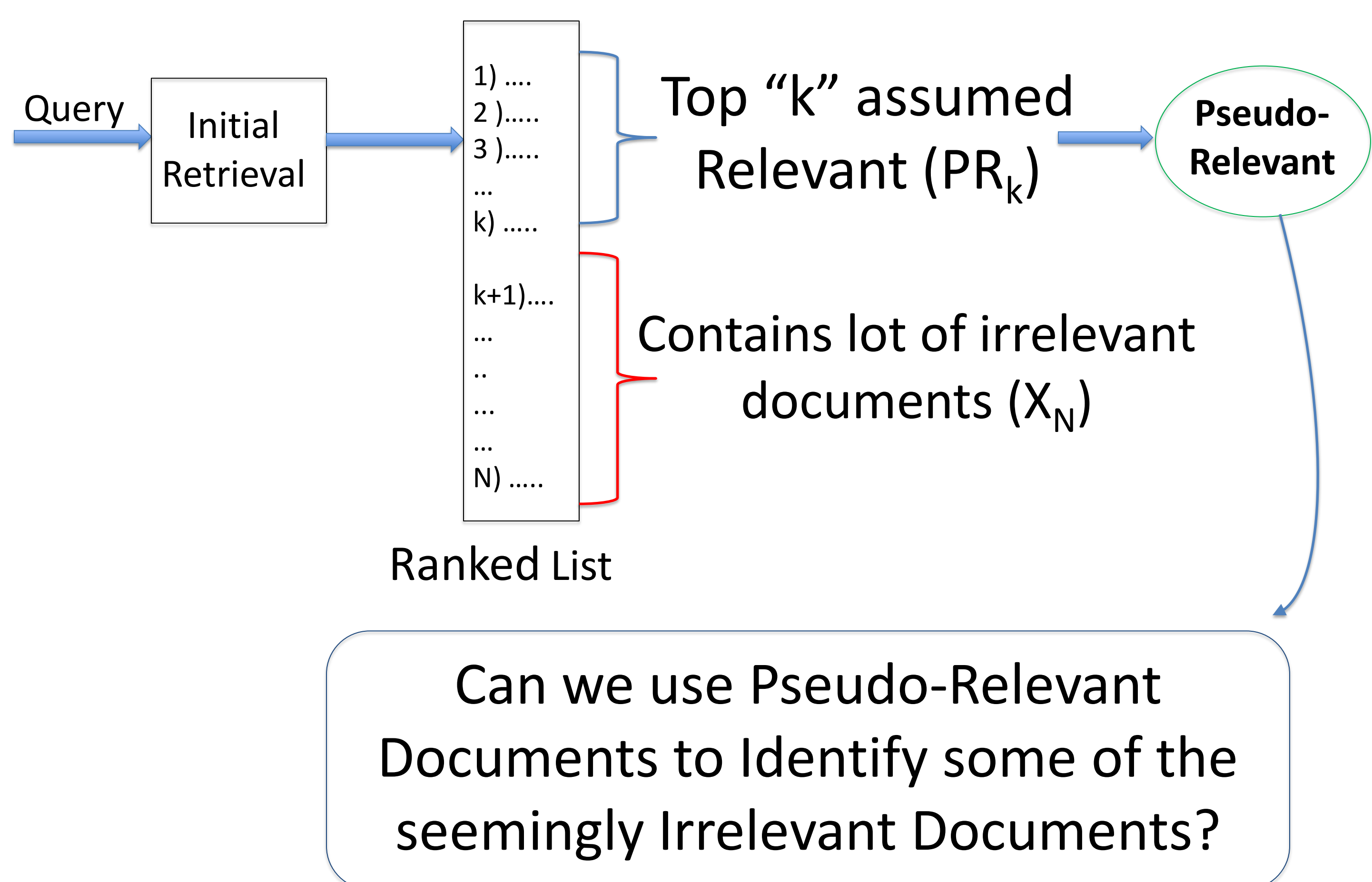
Improving Pseudo-Relevance Feedback Using Pseudo-Irrelevant Documents

Karthik Raman¹, Raghavendra Udupa² and Pushpak Bhattacharya¹

¹Indian Institute of Technology Bombay

²Microsoft Research India

PRF and Motivation behind Pseudo-Irrelevance



"Pseudo-Irrelevant Documents" & their Identification

Y_{PR} : Documents similar to PR_k

$$Y_{PR} = \bigcup_D Q_D \text{ for } \forall D \in PR_k$$

where Q_D : Top Documents retrieved if D was a query

PI: Pseudo – Irrelevant Documents

$$PI = X_N - (X_N \cap Y_{PR})$$

High-Scoring Documents in X_N (Ranked $k+1$ to N) which are dissimilar to pseudo-relevant docs.

Using Pseudo-Irrelevant Documents

- Use Rocchio Algorithm with PI as –ve Feedback
- Zhang et. al 's Distribution Separation Model
- Other Negative Feedback Algorithms

Query Expansion Using "Discriminative Terms"

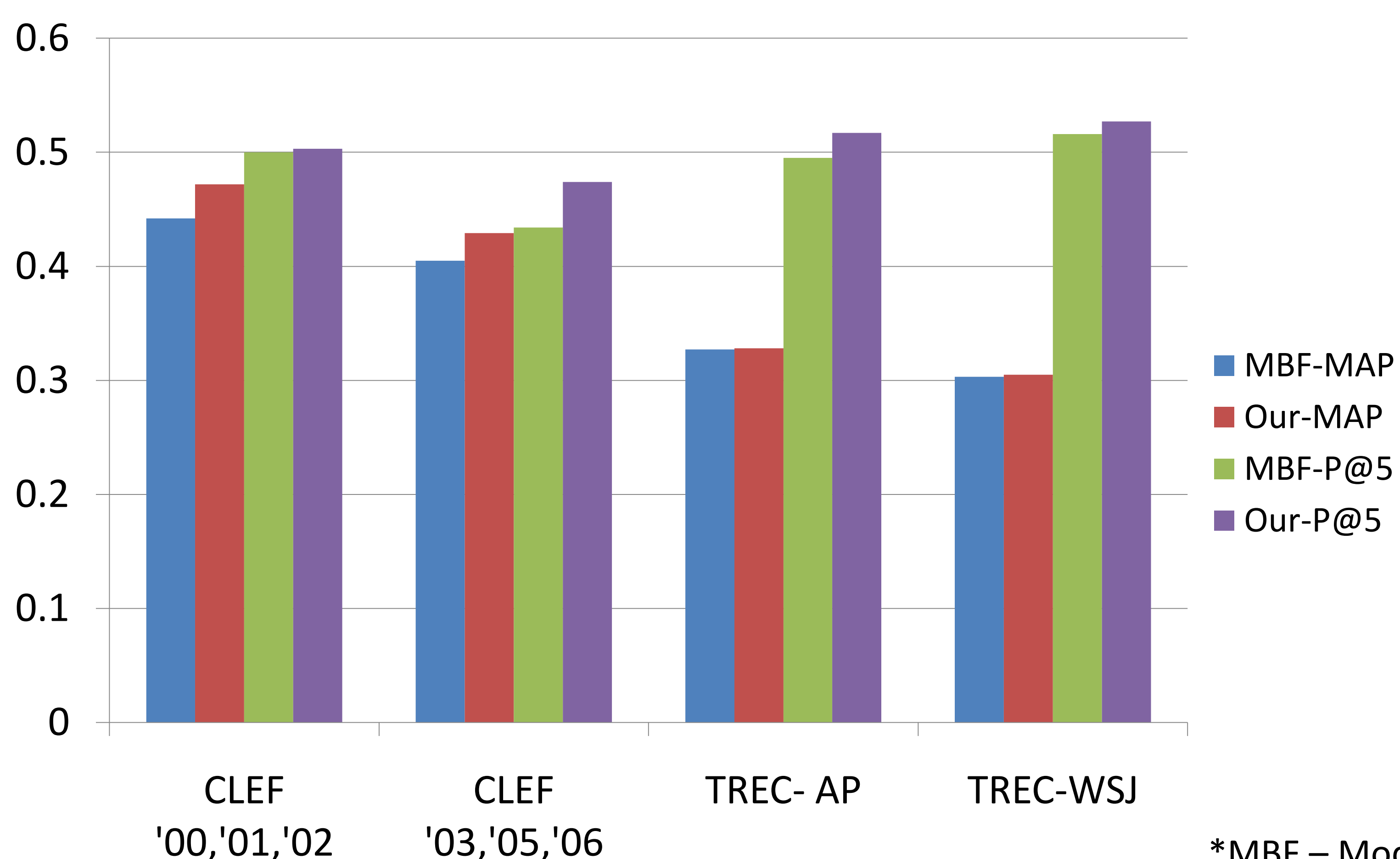
Terms which distinguish Documents in PR_k from those in PI

Run Logistic Classifier with:

- PR_k :+ve ; PI :–ve
- Term TF-IDF Values as Features

Expand Query with Most discriminative terms i.e. highest +ve feature weights,

Preliminary Results?



What Next?

- Can Use Pseudo-Irrelevant Documents to Identify Irrelevant Documents in PR_k .
- Pseudo-Irrelevant Documents found to be closer to Irrelevant than Relevant Documents