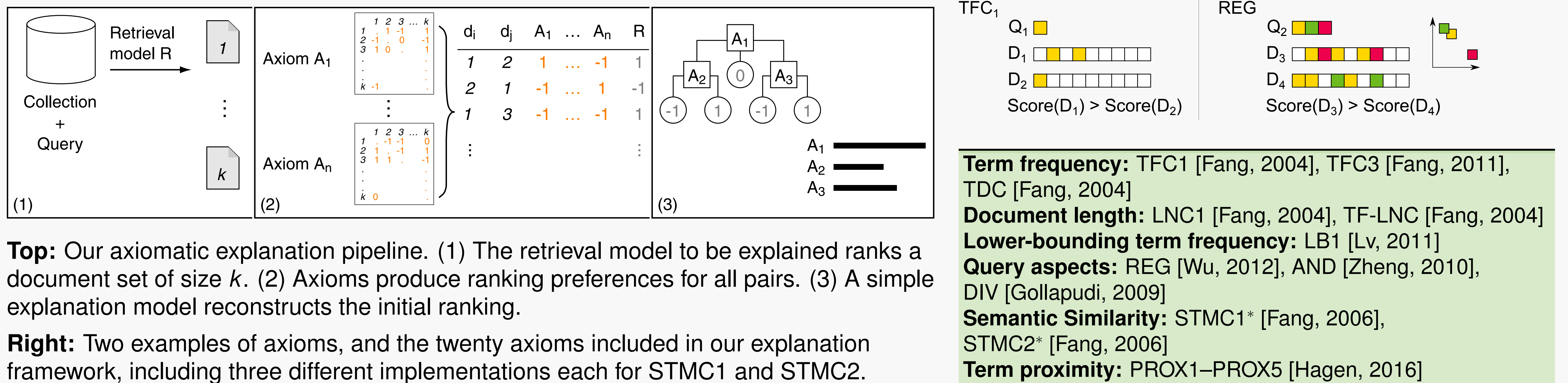


# Towards Axiomatic Explanations for Neural Ranking Models

## Our Axiomatic Explanation Framework Explains Rankings by Reconstruction

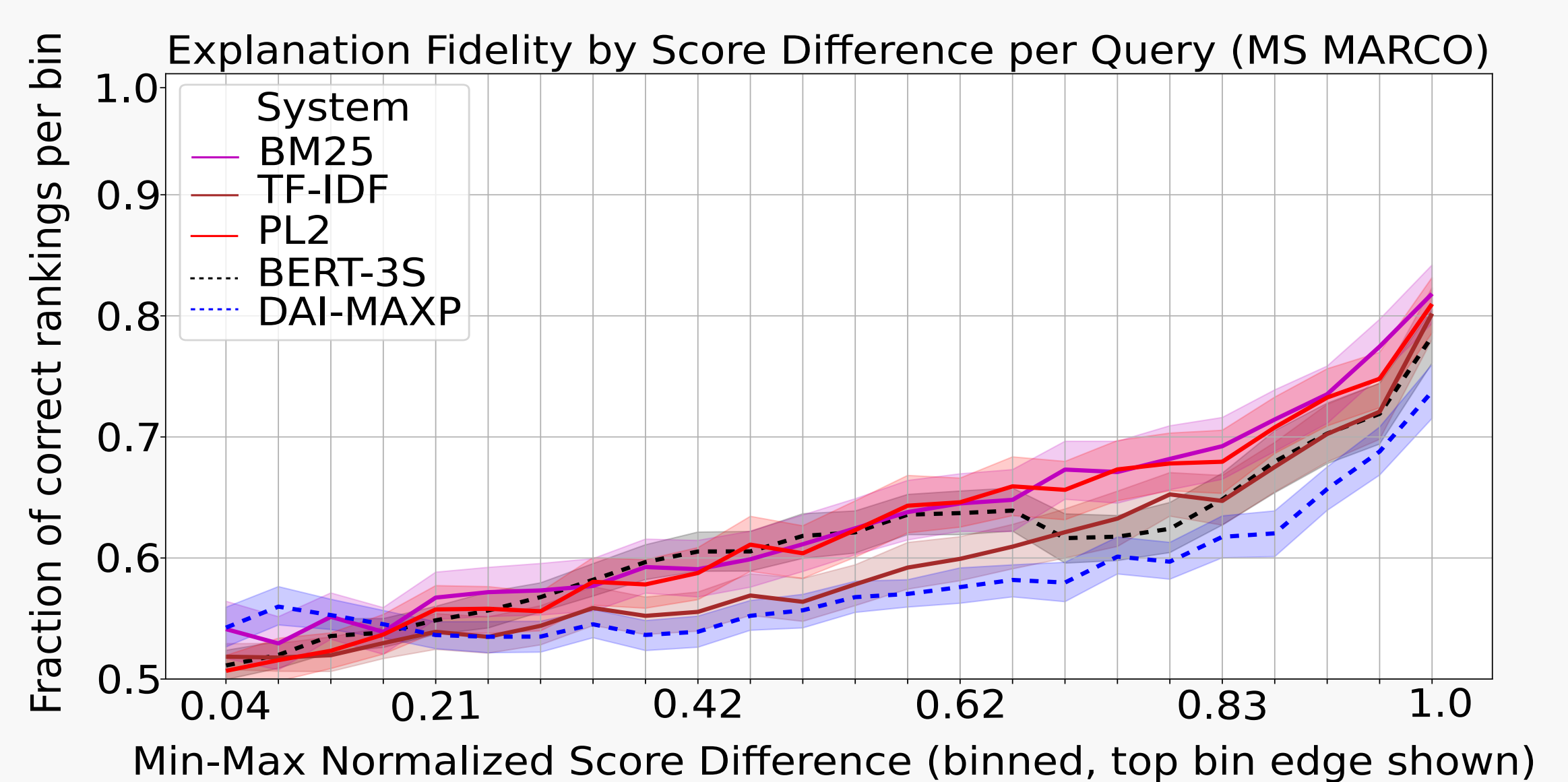
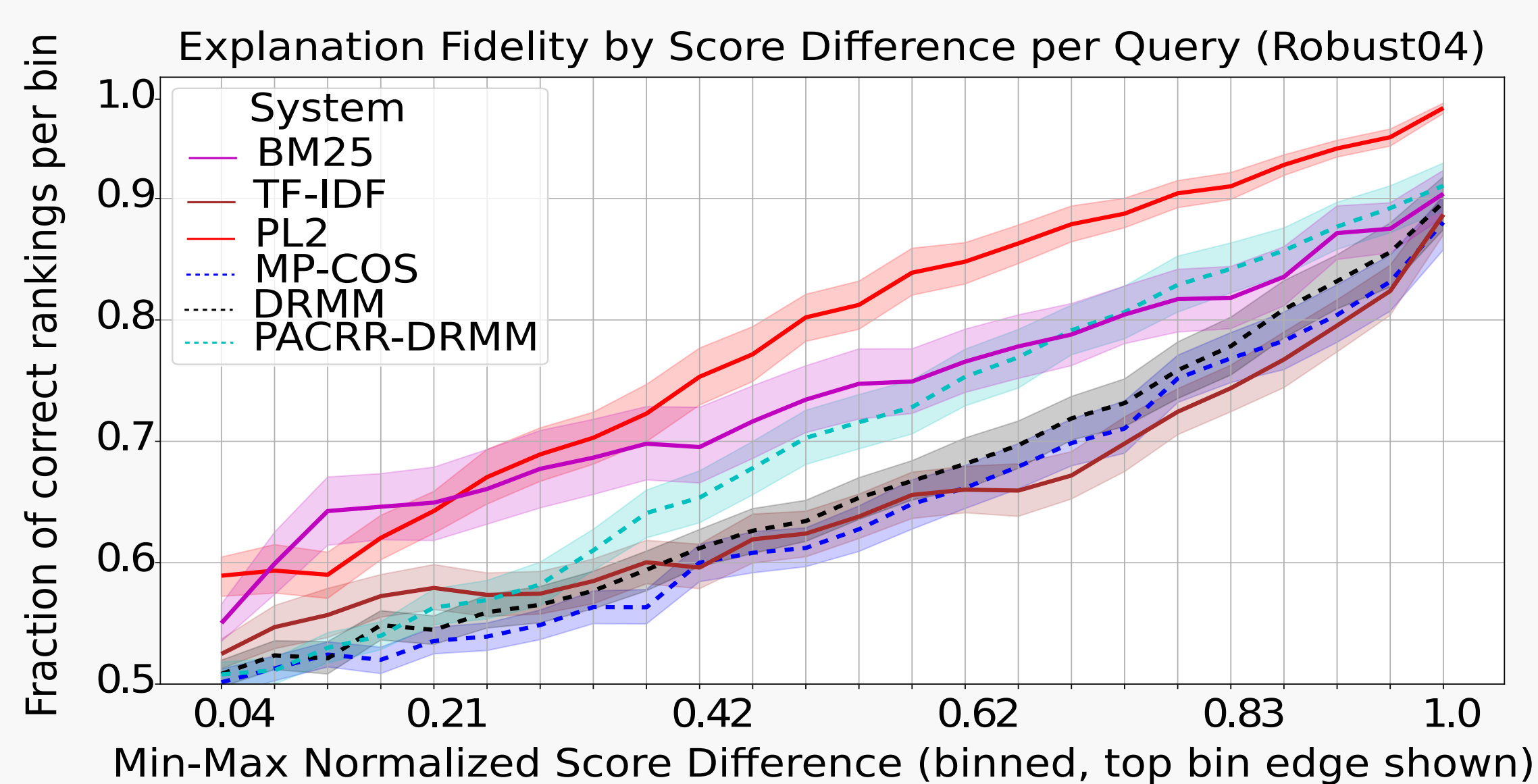


## With our Current Axiom Set, Explanation Fidelity is Limited Except for Distant Pairs

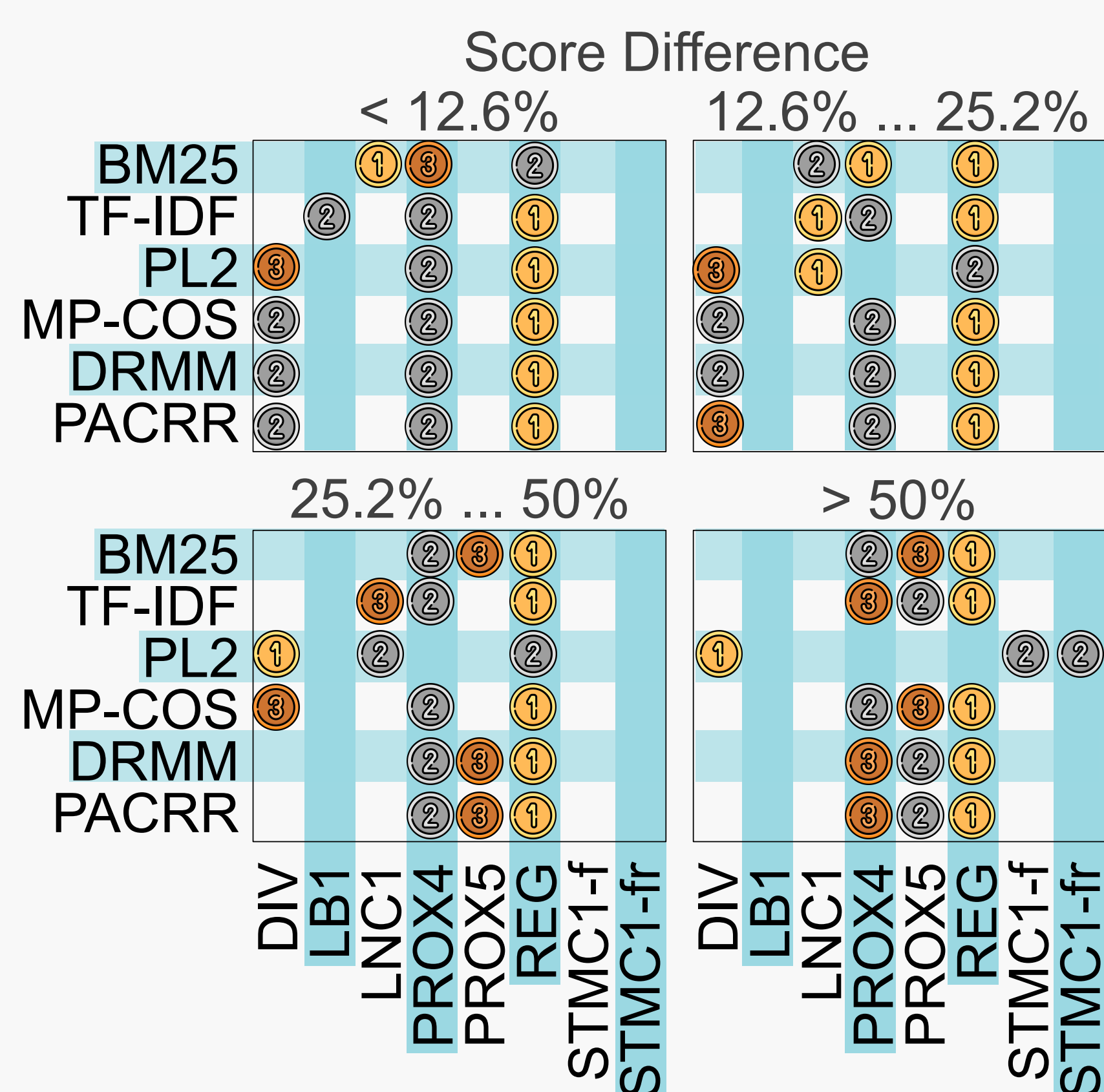
Explanation Models Scope	Per Retr. Model	Explanation Fidelity (Fraction Correct Pairs)					
		Classical Retrieval Models			Neural Retrieval Models		
<i>Robust04</i>		<b>BM25</b>	<b>TF-IDF</b>	<b>PL2</b>	<b>MP-COS</b>	<b>DRMM</b>	<b>PACRR-DRMM</b>
query	100	<b>0.75</b>	<b>0.66</b>	0.78	<b>0.67</b>	<b>0.68</b>	<b>0.72</b>
rank-diff bin	24	0.71	0.63	0.77	0.59	0.61	0.67
score-diff bin	24	0.72	0.64	0.78	0.59	0.61	0.68
query, rank-diff bin	2,368	0.73	0.64	0.77	0.65	0.66	0.70
query, score-diff bin	2,394	0.74	0.65	<b>0.79</b>	0.64	0.66	0.70
<i>MS MARCO</i>		<b>BM25</b>	<b>TF-IDF</b>	<b>PL2</b>	<b>BERT-3S</b>	<b>DAI-MAXP</b>	
query	100	<b>0.64</b>	<b>0.60</b>	<b>0.63</b>	<b>0.61</b>	<b>0.59</b>	
rank-diff bin	24	0.60	0.56	0.59	0.57	0.54	
score-diff bin	24	0.61	0.56	0.59	0.59	0.55	
query, rank-diff bin	2,400	0.62	0.58	0.61	0.60	0.57	
query, score-diff bin	2,376	0.63	0.60	0.62	0.61	0.58	

**Left:** Explanation models for 100 queries each from Robust04 and MS MARCO. For the score- and rank-difference-scope models, document pair samples are divided into 24 bins based on the min-max normalized difference in retrieval scores. Explanation fidelity is computed for each bin separately and is further macro-averaged.

**Below:** The increase in explainability as the score difference grows is more pronounced for Robust04, whereas the explanations on both datasets perform very similarly at the low-score difference end (0.5 is no better than random guessing).



## Top Axioms Overlap Between Models, and “Easier” Queries are Easier to Explain



Query	Explanation fidelity			Topic title	nDCG		
	DRMM	MP-COS	PACRF		DRMM	MP-COS	PACRR
344	0.57	0.55	0.60	Abuses of E-Mail	0.26	0.27	0.33
352	0.64	0.56	0.61	British Chunnel impact	0.15	0.17	0.15
...							
618	0.82	0.81	0.83	Ayatollah Khomeini death	0.41	0.52	0.44
684	0.67	0.55	0.66	Part-time benefits	0.41	0.26	0.39

**Top:** Explanation fidelity and nDCG are weakly positively correlated (Pearson  $\approx 0.1$ ).

**Left:** Top-3 axioms per retrieval model by relative score difference (Robust04).

**Conclusions:**

- Twenty well understood IR axioms to explain black-box neural rankers.
- Explanation fidelity on the smaller, more genre-focused Robust04 with its shorter queries is superior to that on MS MARCO.
- Well-grounded axiomatic constraints capturing other retrieval aspects seem to be needed to further improve explanation fidelity.