

Yu Xia<sup>1</sup> Junda Wu<sup>1</sup> Sungchul Kim<sup>2</sup> Tong Yu<sup>2</sup> Ryan A. Rossi<sup>2</sup> Haoliang Wang<sup>2</sup> Julian McAuley<sup>1</sup><sup>1</sup>University of California San Diego <sup>2</sup>Adobe Research

# Knowledge-Aware Query Expansion with Large Language Models for Textual and Relational Retrieval

## Motivation

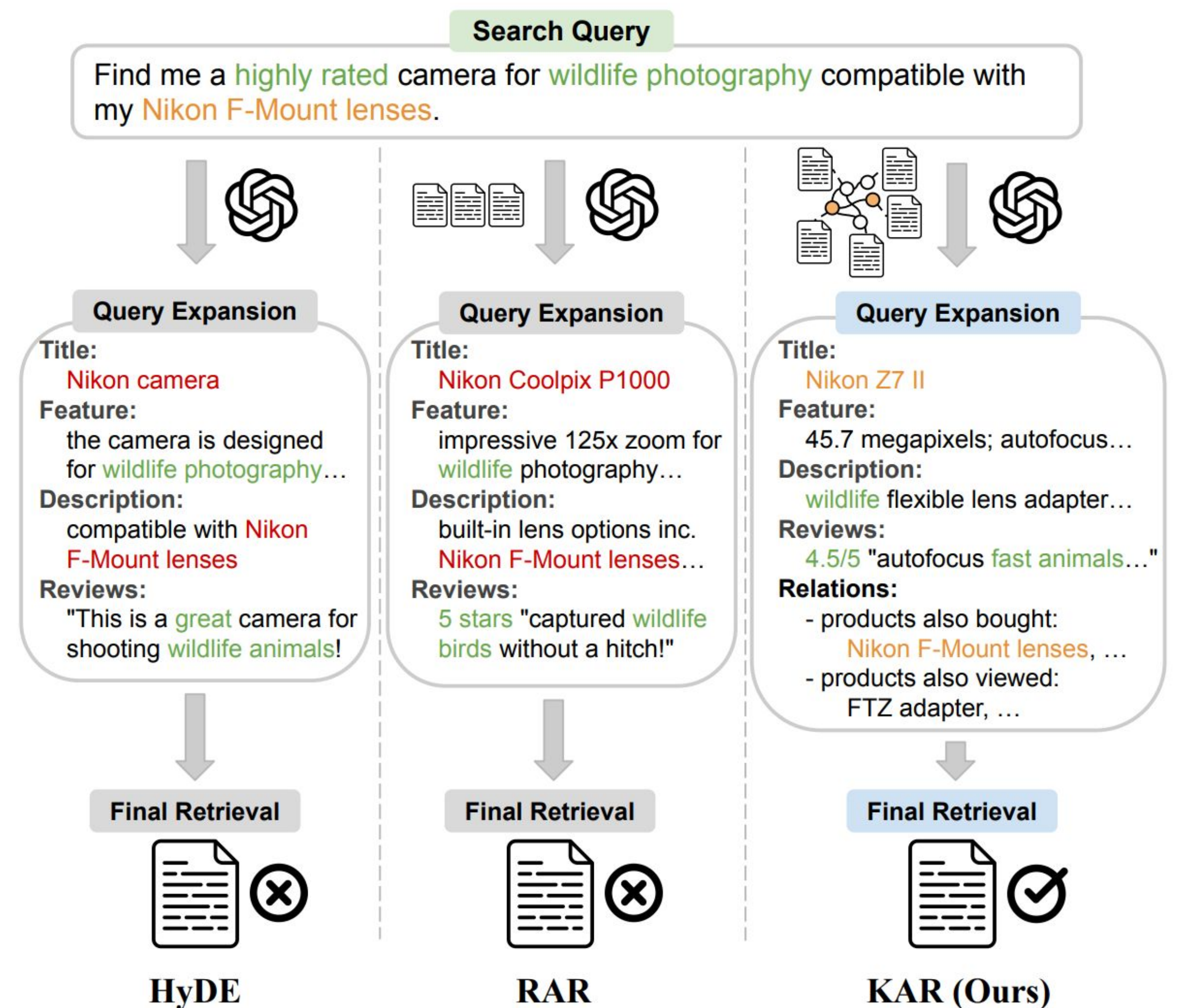
LLMs have been used to generate **query expansions** augmenting information retrieval.

Existing methods:

- focus on **textual similarities** between queries and documents
- overlook **relations** between documents.
- often fail to handle complex queries with both **textual** and **relational** requirements

Our method:

- augment LLMs with structured document relations from Knowledge Graph (KG)
- use both **textual** and **relational** information for query expansion



## Methodology

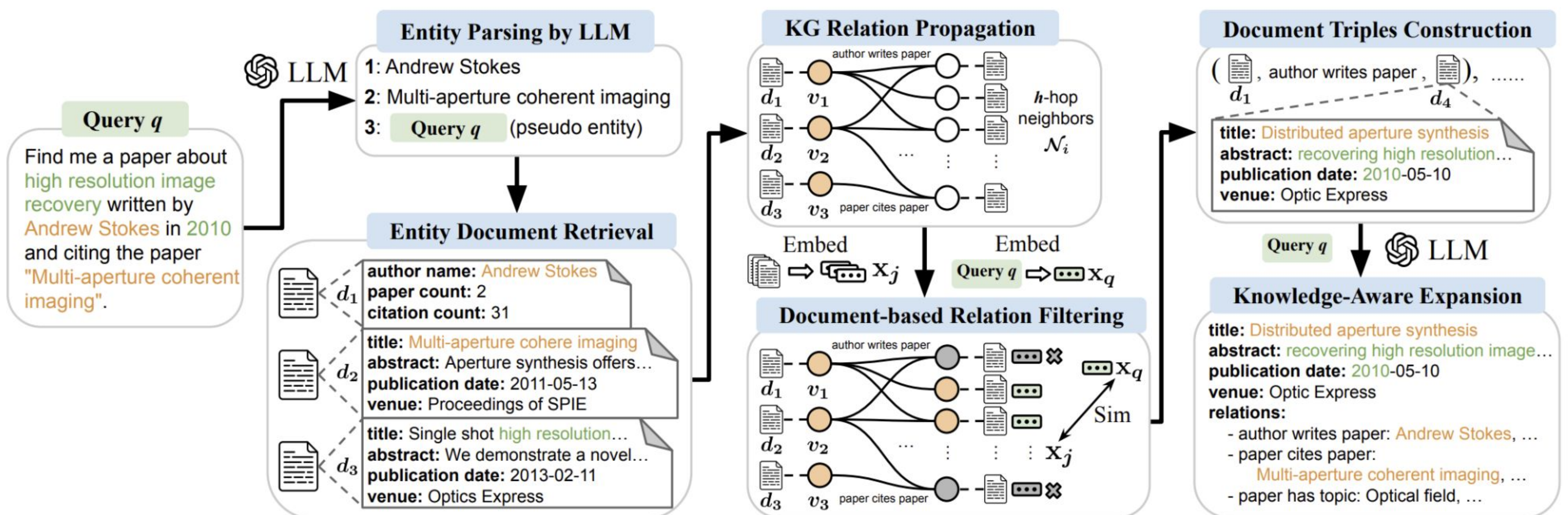


Figure 2: Overview of our knowledge-aware query expansion framework illustrated with an example academic paper search query with **textual** and **relational** requirements.

## Experiments

- Three semi-structured retrieval datasets from STaRK benchmark.

- Table 2 is the results with

- GPT4o as LLM
- text-embedding-ada-002 as embedding.

- More evaluations and results with other LLMs and embeddings please refer to our paper.

Method	AMAZON				MAG				PRIME			
	Hit@1	Hit@5	R@20	MRR	Hit@1	Hit@5	R@20	MRR	Hit@1	Hit@5	R@20	MRR
<i>Supervised Settings</i>												
DPR	15.29	47.93	44.49	30.20	10.51	35.23	42.11	21.34	4.46	21.85	30.13	12.38
QAGNN	26.56	50.01	52.05	37.75	12.88	39.01	46.97	29.12	8.85	21.35	29.63	14.73
AvaTaR	49.87	<b>69.16</b>	<b>60.57</b>	58.70	44.36	59.66	50.63	51.15	18.44	36.73	39.31	26.73
<i>Zero-Shot Settings</i>												
Base	39.16	62.73	53.29	50.35	29.08	49.61	48.36	38.62	12.63	31.49	36.00	21.41
PRF	40.07	60.66	51.24	49.79	29.04	47.65	46.69	37.90	12.46	28.63	33.04	20.06
HyDE	40.31	64.43	53.71	51.42	29.98	50.10	50.02	39.58	16.85	37.59	43.55	26.56
RAR	<u>51.52</u>	66.63	54.63	<u>58.73</u>	39.02	52.87	50.87	45.74	22.53	40.84	44.50	30.93
AGR	49.82	62.97	53.38	56.77	39.29	53.66	51.89	46.20	<u>25.85</u>	44.41	46.63	35.04
KAR <sub>w/o</sub> KG	43.54	60.29	51.83	51.80	31.14	46.75	46.86	38.88	18.03	36.27	42.00	26.84
KAR <sub>w/o</sub> DRF	47.99	67.54	56.91	57.14	<u>45.44</u>	<u>63.83</u>	<u>58.67</u>	<u>53.85</u>	<u>25.85</u>	<u>46.52</u>	<u>48.10</u>	<u>35.52</u>
KAR	<b>54.20</b>	<u>68.70</u>	<u>57.24</u>	<b>61.29</b>	<b>50.47</b>	<b>65.37</b>	<b>60.28</b>	<b>57.51</b>	<b>30.35</b>	<b>49.30</b>	<b>50.81</b>	<b>39.22</b>

Table 2: Retrieval results on test sets of synthetic search queries.

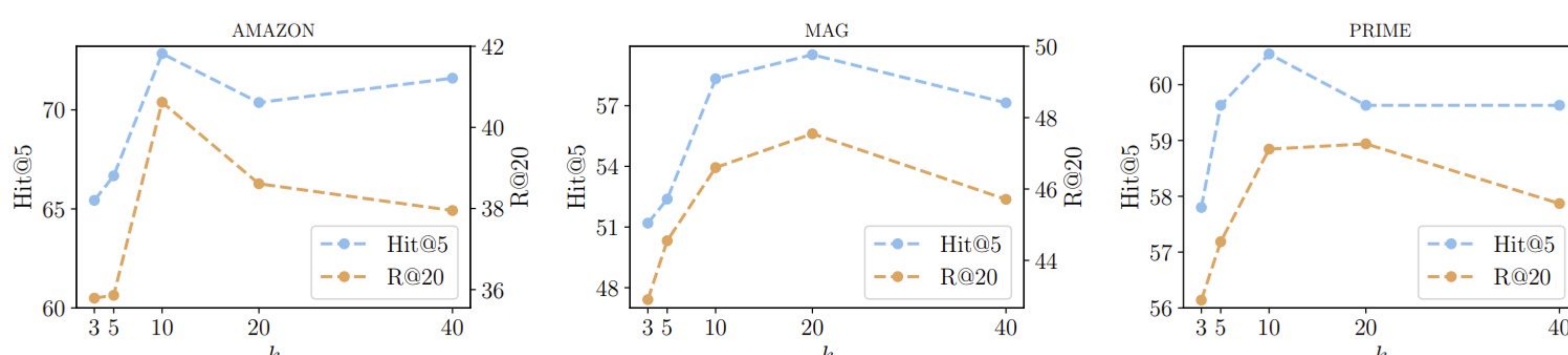


Figure 3: Influence of different values of  $k$  for filtered top- $k$  neighbors in KAR.

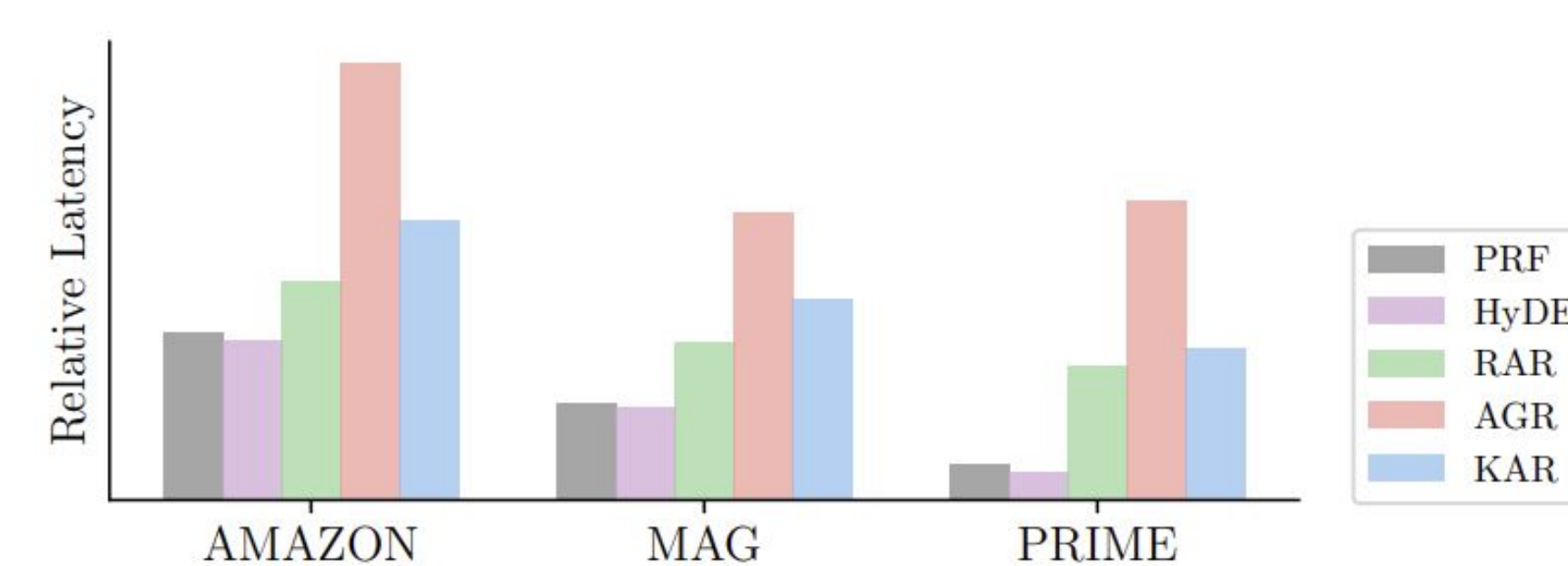


Figure 5: Latency comparison of query expansions.



Paper