Lexically-Accelerated Dense Retrieval

Hrishikesh Kulkarni¹, Sean MacAvaney², Nazli Goharian¹, Ophir Frieder¹

¹Georgetown University ²University of Glasgow

ACM SIGIR 2023



Problem

1. Lexical retrieval:

Inverted index - high efficiency

Vocabulary based - low effectiveness

2. Dense retrieval:

Exhaustive - low efficiency

Semantic - high effectiveness



Existing Solutions

1. Approximation methods: Methods which aim to reduce the computational overhead by approximating results of a full dense retriever e.g. IVF, ScaNN, HNSW, GAR etc

Good at approximating top results - nDCG@10

Suffer in terms of Recall -Recall@1000



Where do the Approximation Methods Lack?

- 1. Seed points influenced by graph building algorithms and don't take query characteristics into account
- 2. Random exploration of search space is expensive
- 3. Approaches like reranking only explore documents having lexical overlap with the query
- 4. Termination at local minimas



How we address this gap?

- 1. LADR: improves efficiency without compromising on retrieval effectiveness
- 2. Informed exploration by identifying right seed points
- 3. Exploration: Pre-computed document proximity graph

Seed Points: Results from Efficient lexical model



Iteratively explores neighbors of top c results until they converge









BM25











Proactive LADR

- 1. Seed with lexical results
- 2. Add neighbors
- 3. Score seeds and their neighbors

Vector comparison: O(kn) Storage: O(kD)

D is the corpus size, k is no. of neighbors

Adaptive LADR

- 1. Seed with lexical results
- 2. Add neighbors
- 3. Score seeds and their neighbors
- 4. Repeat 2 and 3 until no new neighbors can be identified

Vector comparison: O(D)* Storage: O(kD)

*O(D) is worst case - does not happen in practice due to clustering hypothesis

RQ1: How does LADR compare to other approximation techniques in terms of effectiveness and efficiency?

- 1. Operational points: 4ms/query, 8ms/query
- 2. Dense retrieval models: TAS-B, RetroMAE, TCT-ColBERT-HNP, ANCE
- 3. Baselines: IVF, ScaNN, HNSW, GAR, Re-ranking
- 4. Datasets: DL 2019, DL 2020, Dev (small)

RQ1: nDCG and Recall@1k

Statistically significant improvements observed across exhaustive evaluations! Effectiveness-efficiency Pareto-frontier created by LADR among approx k nearest neighbor techniques!



RQ2: Is LADR applicable to a variety of single-representation dense retrieval models?

- 1. Dense retrieval models: TAS-B, RetroMAE, TCT-ColBERT-HNP, ANCE
- 2. Statistically significant results observed on all four dense retrieval models.
- 3. Hence, LADR can be used on top of wide range of dense retrieval models.

Proactive LADR DL 2019 ~ 4 ms

Dense Retrieval Model	nDCG@10	Recall@1k
TAS-B	0.690	0.771
RetroMAE	0.691	0.765
TCT-Colbert-HNP	0.680	0.747
ANCE	0.645	0.751

Dense Retrieval Model	nDCG@10	Recall@1k
TAS-B	0.738	0.872
RetroMAE	0.740	0.866
TCT-Colbert-HNP	0.729	0.848
ANCE	0.665	0.820

Adaptive LADR DL 2019 ~ 8ms

RQ3: What are the computational overheads of LADR?

RQ4:

How do the parameters introduced by LADR affect the effectiveness and efficiency of LADR?

RQ3 and RQ4: Overhead and Parameters

- 1. Latency increases with increase in no. of neighbors, seed set size and exploration depth but so does RBO.
- 2. Best nDCG and Recall is obtained with higher no. of neighbors but lower seed set size.
- 3. Adaptive LADR should be preferred over Proactive for larger time budgets.
- 4. k, c, n parameters introduced in LADR can be tuned deliver best performance in given time budget.
- 5. Adaptive LADR: 0.98 RBO in 27.7ms/q
- 6. Proactive LADR: 0.92 RBO in 69.8ms/q

		Proactive LADR				Adaptive LADR						
		4 -	.62	.66	.70	.72	-	.63	.64	.66	.67	.69
	IS	8 -	.65	.69	.72	.72	-	.64	.66	.68	.70	.71
	bbd	16 -	.69	.71	.73	.73	-	.66	.68	.70	.71	.72
	nD(eig	32 -	.72	.74	.73	.73	-	.68	.69	.71	.72	.72
	k n	64 -	.73	.74	.73	.72	-	.70	.71	.72	.72	.72
		128 -	.74	.73	.73	.72	-	.74	.74	.74	.73	.72
			1	1	1	1		1	1	1	1	1
		4 -	.65	.71	.78	.82	-	.65	.68	.71	.74	.77
	ors	8 -	.71	.76	.81	.84	-	.68	.72	.76	.79	.81
	@1k ghb	16 -	.77	.81	.84	.85	-	.72	.75	.79	.83	.85
	R@ <i>k</i> neiç	32 -	.81	.85	.85	.85	-	.75	.78	.82	.85	.85
		64 -	.85	.86	.85	.85	-	.79	.81	.85	.85	.84
		128 -	.86	.86	.85	.85	-	.85	.86	.87	.86	.85
		4	65	72	70	01		60	70	74		70
	10	4-	.05	.72	.70	.01		.00	.70	.74	.//	.79
	Sor	0- 16	.71	.70	.01	.04	-	.12	.75	.01	CO.	.07
	ght	10 -	.70	.79	.85	.00	-	.//	.80	.00	.89	.91
	hei	32 -	.80	.83	.80	.89	-	.80	.85	.90	.92	.93
	×	64 -	.84	.87	.89	.90	-	.83	.88	.93	.95	.96
		128 -	.87	.89	.91	.92	1	.88	.92	.95	.97	.98
		4 -	3.8	5.1	8.3	12.0	_	4.6	4.7	4.9	5.4	6.6
	s'q) s'	8 -	4.2	5.9	9.4	13.6	_	4.6	4.8	5.2	6.0	7.8
	sm)	16 -	4.8	6.8	11.0	16.9	_	4.5	4.7	5.3	6.4	8.8
	eigh	32 -	5.8	8.2	14.4	26.6	_	4.8	5.2	6.3	8.1	12.2
	k ne	64 -	7.6	11.1	20.4	42.2	_	5.3	5.9	7.8	10.6	17.8
	La	128 -	10.2	15.9	34.9	69.8	_	5.9	7.0	9.9	15.1	27.7
			100	200	500	1000		10	20	50	100	200
						TO	20	ratio		200		

RQ5: Is an exact nearest neighbor graph needed for LADR to be effective?

Document Proximity Sources -Statistically equivalent performance observed across:

- a. Exact
- b. HNSW
- c. BM25

Method	nDCG@10	Recall@1k
Exact	0.730	0.850
Approx	0.731	0.845
BM25	0.732	0.835

RQ6: What are the trade-offs between proactive and adaptive LADR in terms of precision, recall and latency?

- 1. Adaptive LADR performs better than Proactive LADR under same time budget.
- 2. Proactive LADR can deliver good results in very low time budgets.
- Adaptive LADR has a O(D) worst case time complexity (even though unlikely). While Proactive LADR has O(kn).

Lexically-Accelerated Dense Retrieval

- <u>Reduces computational overhead</u> of dense retrieval maintaining high efficacy
- Delivers a proactive and an adaptive strategy for <u>optimal use of time budget</u>
- Establishes a <u>new Pareto frontier</u> for low latency approximate dense retrieval
- <u>Robust</u> to parameters and alternate sources of document proximity

Hrishikesh Kulkarni¹, Sean MacAvaney², Nazli Goharian¹, Ophir Frieder¹

¹Georgetown University ²University of Glasgow

