

KG-FIT: Knowledge Graph Fine-Tuning Upon Open-World Knowledge

Pengcheng Jiang, Lang Cao, Cao Xiao*, Parminder Bhatia*, Jimeng Sun, and Jiawei Han
University of Illinois at Urbana Champaign GE HealthCare*

Patrick (Pengcheng) Jiang
CS Ph.D. at UIUC

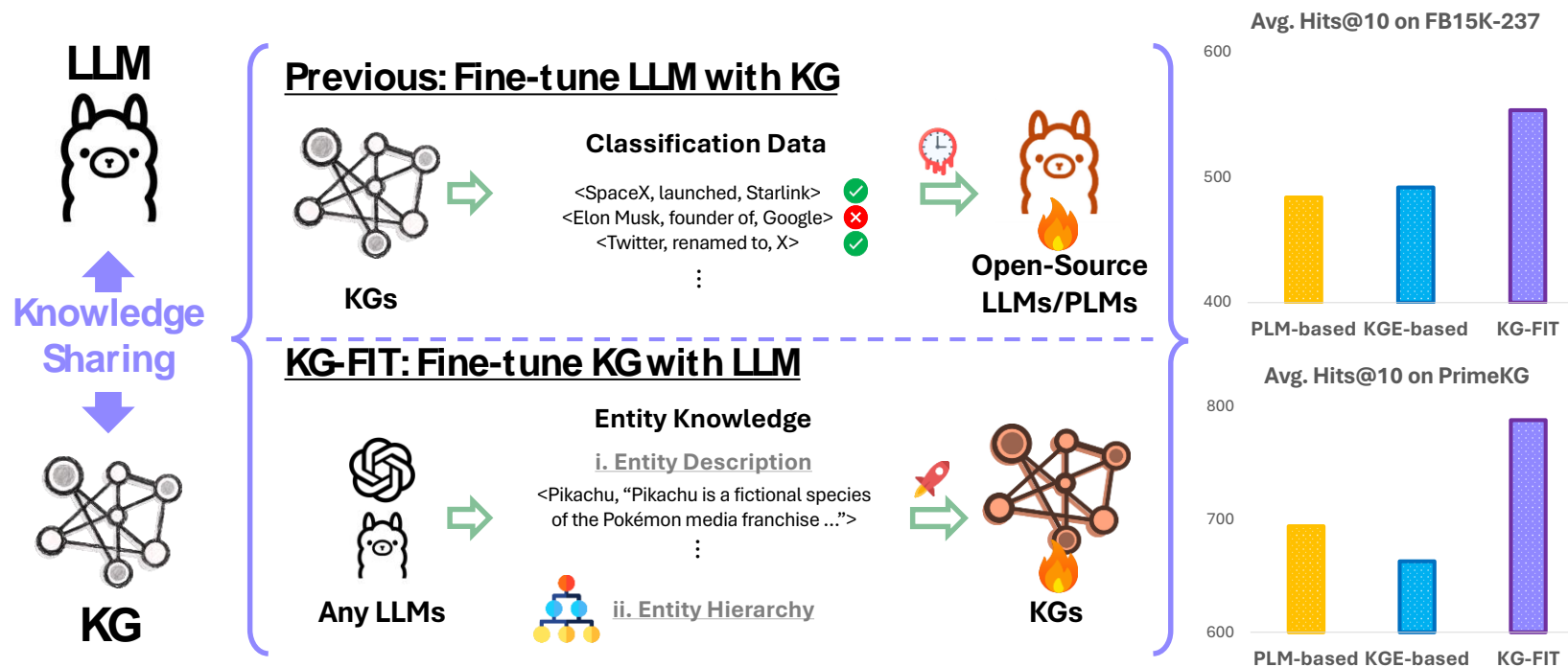
Motivation & Problem

Challenges:

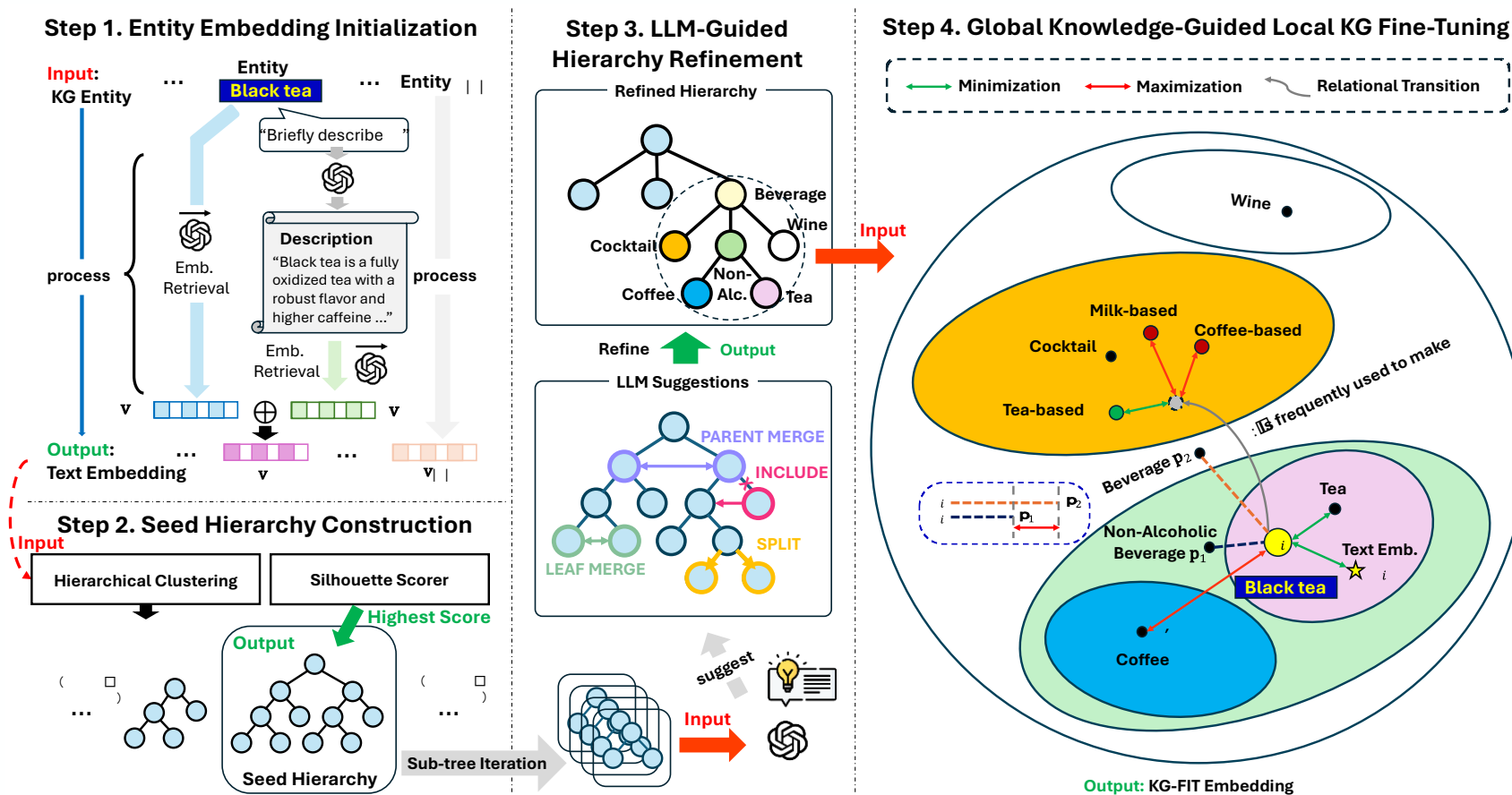
- Structure-based KGE methods are limited to graph structure
- PLM-based methods are computationally expensive
- Need to leverage LLM knowledge efficiently

Our Solution:

- Fine-tune KG with LLM instead of fine-tuning LLM with KG
- Combine advantages of both approaches



KG-FIT Framework



Two-Stage Approach:

1. LLM-Guided Hierarchy Construction
 - Entity description generation
 - Seed hierarchy construction
 - LLM-guided refinement
2. Knowledge Graph Fine-Tuning
 - Hierarchical clustering constraint
 - Semantic anchoring
 - Link prediction objective

Datasets

- Datasets**

Table 1: **Datasets statistics.** #Ent./#Rel: number of entities/relations. #Train/#Valid/#Test: number of triples contained in the training/validation/testing set.

Dataset	#Ent.	#Rel.	#Train	#Valid	#Test
FB15k-237	14,541	237	272,115	17,535	20,466
YAGO3-10	123,182	37	1,079,040	5,000	5,000
PrimeKG	10,344	11	100,000	3,000	3,000

- Metrics**

Mean Rank (MR):

- Measures the average rank of true entities.

Mean Reciprocal Rank (MRR):

- Averages the reciprocal ranks of true entities.

Hits@N:

- Measures the proportion of true entities in the top N predictions.

FB15K-237:

- A subset of Freebase, a large collaborative knowledge base focusing on common knowledge.

YAGO3-10:

- A subset of YAGO, a large knowledge base derived from multiple sources including Wikipedia, WordNet, and GeoNames.

PrimeKG:

- A biomedical KG integrates 20 biomedical resources, detailing 17,080 diseases through 4,050,249 relationships. In this study, we extract a subset of PrimeKG, which contains 106,000 triples.



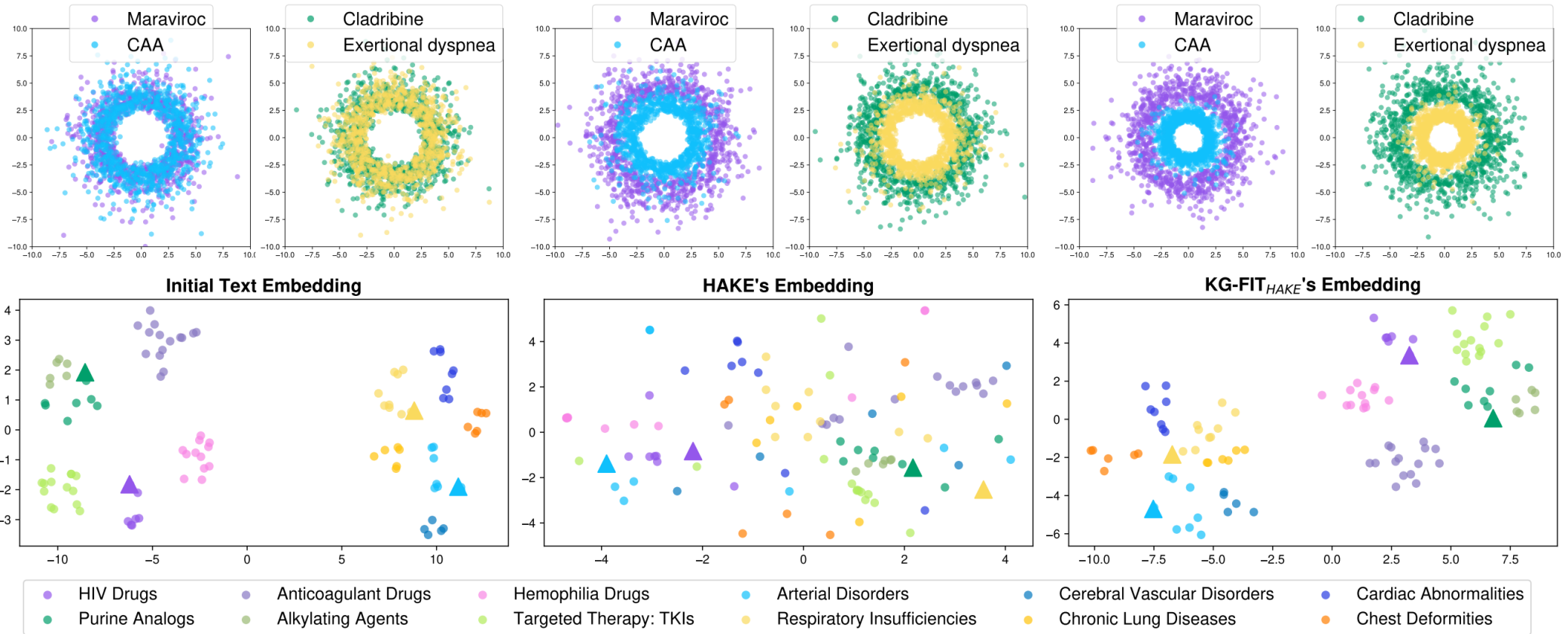
Main Results

- (1) KG-FIT consistently and significantly outperforms state-of-the-art PLM-based and structure-based methods across all datasets and metrics.
- (2) With LLM-guided hierarchy refinement, KG-FIT achieves huge performance gains compared to the base models and KG-FIT with seed hierarchy.
- (3) KG-FIT is more effective for smaller KGs, e.g., more performance gains on PrimeKG (~ 0.1 million triples) than YAGO3-10 (~1 million triples).

		FB15K-237					YAGO3-10					PrimeKG					
		PLM-based Embedding Methods															
Model	PLM	MR	MRR	H@1	H@5	H@10	MR	MRR	H@1	H@5	H@10	MR	MRR	H@1	H@5	H@10	
KG-BERT [22]*	BERT	153	.245	.158	–	.420	–	–	–	–	–	–	–	–	–	–	
StAR [23]*	RoBERTa	117	.296	.205	–	.482	–	–	–	–	–	–	–	–	–	–	
PKGC [28]	RoBERTa	184	.342	.236	.441	.525	1225	.501	.426	.596	.660	219	.485	.391	.565	.625	
C-LMKE [26]*	BERT	141	.306	.218	–	.484	–	–	–	–	–	–	–	–	–	–	
KGT5 [25]*	T5	–	.276	.210	–	.414	–	.426	.368	–	.528	–	–	–	–	–	
KG-S2S [24]*	T5	–	.336	.257	–	.498	–	–	–	–	–	–	–	–	–	–	
SimKGC [27]	BERT	–	.336	.249	–	.511	–	–	–	–	–	168	.527	.524	.679	.742	
CSProm-KG [32]	BERT	–	.358	.269	–	.538	1145	.488	.451	.624	.675	157	.540	.492	.652	.745	
LLM Emb. (zero-shot)	TE-3-S	2044	.023	.002	.035	.068	22741	.009	.000	.016	.024	5581	.000	.000	.000	.000	
	TE-3-L	1818	.030	.004	.048	.085	18780	.015	.000	.019	.032	4297	.001	.000	.000	.000	
		Structure-based Embedding Methods															
Model	Frame	\mathcal{H}	MR	MRR	H@1	H@5	H@10	MR	MRR	H@1	H@5	H@10	MR	MRR	H@1	H@5	H@10
TransE	Base [14]	—	233	.287	.192	.389	.478	1250	.500	.398	.626	.685	182	.048	.000	.043	.124
	KG-FIT	Seed LHR	142 122	.345 .362	.242 .264	.457 .478	.547 .568	952 529	.520 .544	.429 .463	.638 .650	.700 .705	80 69	.298 .334	.000 .000	.315 .342	.516 .536
DisMult	Base [15]	—	283	.260	.163	.349	.437	5501	.451	.365	.553	.615	174	.577	.475	.699	.782
	KG-FIT	Seed LHR	184 154	.316 .331	.198 .226	.415 .433	.512 .529	963 861	.486 .527	.413 .441	.591 .636	.673 .682	107 78	.589 .617	.495 .526	.715 .747	.799 .813
ComplEx	Base [16]	—	347	.252	.161	.344	.439	6681	.463	.384	.560	.612	202	.614	.522	.728	.789
	KG-FIT	Seed LHR	201 151	.325 .344	.223 .247	.436 .458	.523 .551	997 842	.491 .544	.422 .460	.603 .646	.669 .697	94 82	.638 .651	.548 .566	.767 .772	.823 .835
ConvE	Base [17]	—	341	.312	.224	.401	.508	1105	.529	.451	.619	.673	144	.516	.456	.645	.760
	KG-FIT	Seed LHR	181 177	.318 .318	.237 .241	.411 .415	.521 .525	912 885	.535 .541	.455 .461	.628 .647	.685 .695	93 72	.627 .648	.534 .547	.757 .767	.812 .824
TuckER	Base [18]	—	363	.320	.230	.417	.505	1110	.529	.454	.633	.690	171	.543	.442	.663	.737
	KG-FIT	Seed LHR	175 144	.330 .349	.241 .255	.433 .448	.521 .543	874 838	.538 .545	.458 .466	.651 .654	.703 .708	77 62	.640 .648	.542 .550	.770 .779	.805 .820
pRotatE	Base [19]	—	188	.310	.205	.399	.502	974	.477	.385	.573	.655	118	.491	.399	.593	.681
	KG-FIT	Seed LHR	160 119	.355 .371	.257 .277	.461 .483	.558 .572	910 829	.525 .550	.436 .464	.622 .648	.693 .710	75 69	.635 .649	.538 .574	.745 .779	.809 .833
RotatE	Base [19]	—	190	.333	.241	.428	.528	1620	.495	.402	.550	.670	57	.539	.447	.646	.727
	KG-FIT	Seed LHR	141 120	.354 .369	.261 .274	.464 .488	.555 .570	790 744	.529 .563	.440 .475	.643 .658	.708 .722	46 34	.622 .645	.517 .532	.740 .758	.805 .817
HAKE	Base [20]	—	184	.344	.247	.435	.538	1220	.530	.431	.634	.681	95	.595	.515	.708	.760
	KG-FIT	Seed LHR	162 137	.358 .362	.268 .275	.470 .485	.563 .572	854 810	.541 .568	.455 .474	.647 .662	.703 .718	82 42	.638 .682	.540 .605	.747 .785	.808 .835



Visualization



Conclusion

We introduced KG-FIT, a novel framework that enhances knowledge graph (KG) embeddings by integrating open-world entity knowledge from Large Language Models (LLMs).

- KG-FIT effectively combines the knowledge from LLM and KG to preserve both global and local semantics, achieving state-of-the-art link prediction performance on benchmark datasets.
- It shows significant improvements in accuracy compared to the base models. Notably, KG-FIT can seamlessly integrate knowledge from any LLM, enabling it to evolve with ongoing advancements in language models.
- Future work will explore using the KG-FIT embedding for precise knowledge retrieval, which can set a strong foundation for retrieval augmented generation (RAG) by LLMs.

Thank you!

Patrick Jiang