



Empowering Small-Scale Knowledge Graphs: A Strategy of Leveraging General-Purpose Knowledge Graphs for Enriched Embeddings

<u>Albert Sawczyn</u>, Jakub Binkowski, Piotr Bielak, Tomasz Kajdanowicz

LREC-COLING 2024 - The 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation Lingotto Conference Centre - Torino (Italia) 20-25 May, 2024

Plan of the presentation

- 1. Introduction
- 2. Motivation
- 3. Framework
- 4. Experimental methodology
- 5. Results
- 6. Conclusion & future work

Introduction

- Availability of advanced AI methods has contributed to the growth of applications in knowledge-intensive tasks, e.g.:
 - o QA
 - Fact-checking
- Limitations of LLM in updating facts after training: outdated or incorrect information (hallucinations).



Motivation

- Leveraging Knowledge Graphs (KGs) holds promise in addressing challenges of knowledge-intensive applications.
- Development of domain-specific knowledge graphs is hindered by cost and complexity.
- Small KGs suffer from:
 - limited relational structures
 - sparse entity interactions
 - reduced contextual information



Example of KG.

Motivation

- Empowering small-scale KGs will facilitate the development of KG-based systems, particularly for:
 - Small companies and startups
 - Research groups and academia
 - Less popular or niche domains
- How can we utilize existing resources for that?



Example of KG.

Motivation

• Idea:

Integrating external KG as knowledge source to empower small-scale KG.

• Aim:

Developing a framework for <u>enriching vector</u> <u>representations</u> of domain-specific knowledge graphs by <u>linking them</u> with well-established, general-purpose KGs.



Framework

• We propose a <u>generic and modular</u> framework for <u>enriching small and domain-specific</u> KGs (DKG) with <u>general-purpose</u> KGs (GKG) using alignment and linking operation:



Alignment and linking

We perform alignment and linking without requiring manual annotation, leveraging the textual entity labels:

1. We acquire a representation vector $x(e_i)$ for each entity e_i based on three components:



- 2. For each entity in the DKG, we find *k* nearest neighbors in the GKG, based on the vectors $x(\cdot)$.
- 3. For each pair of neighbors found, we create an artificial triple $(e_i^{(d)}, r_l, e_j^{(g)})$.

Representation learning

- For the linked graph $\mathcal{G}_l = (\mathcal{E}_d \cup \mathcal{E}_g, \mathcal{R}_d \cup \mathcal{R}_g \cup \{r_l\}, \mathcal{T}_d \cup \mathcal{T}_g \cup \mathcal{T}_l)$, we train the <u>representation</u> <u>model</u> and evaluate its performance on <u>KG completion task</u>.
- We propose a <u>weighted loss function</u> to mitigate entity alignment's negative effects between two considered KGs:

$$\mathcal{L}(\mathcal{T}, \tilde{\mathcal{T}}; \theta) = \sum_{t \in \mathcal{T}_d \cup \mathcal{T}_g} L(t, \tilde{\mathcal{T}}; \theta) + \sum_{s \in \mathcal{T}_l} w_s \cdot L(s, \tilde{\mathcal{T}}; \theta)$$
$$w_s = \frac{1}{1 + distance(x(e_i), x(e_j))}$$

- \mathcal{T} : set of positive triples
- $\tilde{\mathcal{T}}$: set of sampled negative triples
- \mathcal{T}_d : triples from DKG
- \mathcal{T}_g : triples from GKG
- $\mathcal{T}_l: ext{set of linking triples}$

Experimental methodology

- Existing datasets are not suitable for evaluating methods aimed at enriching the embeddings of KGs.
- We designed a custom evaluation procedures tailored to this context.
- We simulated different stages in developing a KG by:
 - 1. Triple sampling
 - 2. Node sampling
 - 3. Relation sampling

Each strategy is parameterized by the probability p of keeping a triple, node, or relation in the sampled graph.



Experimental methodology

We conducted empirical studies in synthetic and real-world scenarios:

• Synthetic



- controlled environment
- perfect overlap of knowledge domains
- simplified evaluation

• Real-world scenario:



- external KG differs from DKG
- mimics real scenarios
- limited experimental control

scenario:

Results: Synthetic scenario

Experiments settings:

- Datasets: WN18RR, FB15k-237, WD50K
- Sampling: triple, node, relation
- *p*: 0.4, 0.6, 0.8
- Representation model: RotatE

			$Hits$ @10 \uparrow			$MR\downarrow$				
sampling	р	single	linked	boost(%)	single	linked	boost(%)	single	e linked	
					WN1	8RR				
triple	0.4	0.347 ± 0.005	$\textbf{0.502} \pm 0.006$	44.9	7681 ± 62	1245 ± 44	83.8	$0.270~\pm~0.003$	0.345 ± 0.001	27.8
triple	0.6	0.446 ± 0.001	0.519 ± 0.004	16.4	4908 ± 172	1392 ± 15	71.6	$0.342~\pm~0.003$	0.373 ± 0.004	9.3
triple	0.8	0.525 ± 0.004	0.546 ± 0.004	4.0	2685 ± 155	1435 ± 52	46.6	$0.416~\pm~0.005$	0.421 ± 0.003	1.3
node	0.4	0.546 ± 0.004	0.597 ± 0.002	9.3	2164 ± 52	713 ± 18	67.0	$0.473~\pm~0.002$	0.494 ± 0.001	4.4
node	0.6	0.562 ± 0.000	0.590 ± 0.001	4.8	2044 ± 29	950 ± 15	53.6	0.480 ± 0.001	0.488 ± 0.001	1.7
node	0.8	0.576 ± 0.004	0.583 ± 0.001	1.3	$1829~\pm~8$	1194 ± 20	34.7	0.484 ± 0.001	0.482 ± 0.001	-0.2
relation	0.4	0.696 ± 0.182	0.721 ± 0.184	3.5	1083 ± 450	124 ± 70	88.5	0.590 ± 0.233	$0.557~\pm~0.212$	-5.6
relation	0.6	0.755 ± 0.129	0.777 ± 0.143	2.8	1542 ± 622	552 ± 649	64.2	0.684 ± 0.145	0.696 ± 0.160	1.7
relation	0.8	0.731 ± 0.115	0.760 ± 0.135	4.0	2013 ± 157	591 ± 671	70.7	0.659 ± 0.130	0.678 ± 0.149	3.0
Max boost (%)			44.9			88.5			27.8	
Mean boos	st (%)			10.1			64.5			4.8

Results on WN18RR (more results in the paper).

Results: Synthetic scenario

Key Findings:

- Significant improvement by linking GKG.
- Effectiveness depends on sampling setting.
- Effectiveness varies among datasets due to inherent characteristics.
- Achieved 44.9% / 0.0% / 16.7% Hits@10 boost on WN18RR / FB15k-237 / WD50K with only 40% of triples.



Results: Real-world scenario

Experiments settings:

- Datasets pairs:
 - WN18RR ConceptNet
 - WN18RR FB15k-237
 - FB15k-237 ConceptNet
 - FB15k-237 YAGO3-10
 - WD50K FB15k-237
 - WD50K YAGO3-10

- Sampling: triple
- *p:* 0.4, 0.6, 0.8
- Representation model: RotatE

		<i>Hits</i> @10↑					MR	↓		$MRR\uparrow$				
dataset	р	single	link	ked	boost(%)	single linked k		linked boost(%) single lin		single linked boost(%) single		linked		boost(%)
			CN	FB			CN	FB			CN	FB		
WN18RR	0.4	0.347 ± 0.005	0.395	0.338	14.0	7681 ± 62	1508	4050	80.4	0.270 ± 0.003	0.259	0.262	-3.2	
WN18RR	0.6	0.446 ± 0.001	0.471	0.435	5.7	4908 ± 172	963	2685	80.4	0.342 ± 0.003	0.330	0.330	-3.3	
WN18RR	0.8	0.525 ± 0.004	0.527	0.504	0.3	2685 ± 155	750	1725	72.1	0.416 ± 0.005	0.390	0.400	-3.8	
Max boost (%)					14.0				80.4				-3.2	
Mean boos	t (%)				6.7				77.6				-3.5	

Results on WN18RR (more results in the paper).

Results: Real-world scenario

Key Findings:

- Highest efficacy improvement observed under rigorous conditions (40% of triples).
- Achieved notable Hits@10 boost on WN18RR-ConceptNet and WD50K-FB15k-237 pairs, suggesting framework effectiveness.
 - WN18RR ConceptNet:
 - WD50K FB15k-237:
- Importance of choosing a well-suited GKG for the DKG.

Conclusion

Framework

- Proposed general and modular framework for enriching embeddings of small scale KGs.
- Utilized straightforward <u>alignment method based on textual embeddings</u> of entities and their neighborhood.
- Proposed <u>weighted loss function</u> mitigates negative effects of entity alignment.

Extensive Experimentation

- Conducted evaluation in <u>synthetic</u> and <u>real-world</u> scenarios, that simulate early KG development stages.
- Results show <u>significant performance improvement</u> on downstream tasks.
- Degree of improvement varies based on specific DKG and linked general-purpose GKG.

Implications and Future Directions

- Utilizing GKGs strengthens emerging KGs, enhancing their utility and effectiveness.
- Research signals potential pathway for future exploration in enriching small-scale KGs.

Acknowledgements

This work was funded by

- Horizon Europe Framework Programme MSCA Staff Exchanges grant no. 101086321 (OMINO)
- The Polish Ministry of Education and Science under the programme entitled International Co-Financed Projects, grant no. 573977.
- The National Science Centre, Poland under CHIST-ERA Open & Re-usable Research Data & Software (grant number 2022/04/Y/ST6/00183)
- Department of Artificial Intelligence, Wroclaw Tech

Thank you for your attention!

- **Code:** *github.com/graphml-lab-pwr/empowering-small-scale-kg*
- Contact: albert.sawczyn@pwr.edu.pl



Empowering Small-Scale Knowledge Graphs: A Strategy of Leveraging General-Purpose Knowledge Graphs for Enriched Embeddings

Albert Sawczyn, Jakub Binkowski, Piotr Bielak, Tomasz Kajdanowicz





Appendix

Results: Synthetic scenario

			$Hits@10 \uparrow$		$MR\downarrow$							
sampling	р	single	linked	boost(%)	single	linked	boost(%)	single	linked	boost(%)		
WN18RR												
triple	0.4	0.347 ± 0.005	0.502 ± 0.006	44.9	7681 ± 62	$1245~\pm~44$	83.8	$0.270~\pm~0.003$	$0.345~\pm~0.001$	27.8		
triple	0.6	0.446 ± 0.001	0.519 ± 0.004	16.4	4908 ± 172	$1392~\pm~15$	71.6	0.342 ± 0.003	$0.373~\pm~0.004$	9.3		
triple	0.8	0.525 ± 0.004	0.546 ± 0.004	4.0	2685 ± 155	$1435~\pm~52$	46.6	$0.416~\pm~0.005$	0.421 ± 0.003	1.3		
node	0.4	0.546 ± 0.004	0.597 ± 0.002	9.3	2164 ± 52	$713~\pm~18$	67.0	$0.473~\pm~0.002$	$0.494~\pm~0.001$	4.4		
node	0.6	0.562 ± 0.000	0.590 ± 0.001	4.8	2044 ± 29	950 ± 15	53.6	$0.480~\pm~0.001$	$0.488~\pm~0.001$	1.7		
node	0.8	0.576 ± 0.004	0.583 ± 0.001	1.3	$1829~\pm~8$	1194 ± 20	34.7	0.484 ± 0.001	$0.482~\pm~0.001$	-0.2		
relation	0.4	0.696 ± 0.182	0.721 ± 0.184	3.5	1083 ± 450	$124~\pm~70$	88.5	0.590 ± 0.233	$0.557~\pm~0.212$	-5.6		
relation	0.6	0.755 ± 0.129	0.777 ± 0.143	2.8	1542 ± 622	552 ± 649	64.2	0.684 ± 0.145	0.696 ± 0.160	1.7		
relation	0.8	0.731 ± 0.115	0.760 ± 0.135	4.0	2013 ± 157	591 ± 671	70.7	0.659 ± 0.130	0.678 ± 0.149	3.0		
Max boos	st (%)			44.9			88.5			27.8		
Mean boo	st (%)			10.1			64.5			4.8		
					FB15	k-237						
triple	0.4	0.360 ± 0.002	0.360 ± 0.000	0.0	315 ± 4	279 ± 1	11.2	0.204 ± 0.001	0.204 ± 0.001	0.1		
triple	0.6	0.398 ± 0.001	0.393 ± 0.003	-1.2	$242~\pm~1$	$231~\pm~0$	4.2	0.227 ± 0.001	$0.224~\pm~0.002$	-1.5		
triple	0.8	0.442 ± 0.000	0.430 ± 0.002	-2.7	194 ± 2	195 ± 1	-0.3	0.254 ± 0.001	$0.246~\pm~0.000$	-3.1		
relation	0.4	0.480 ± 0.009	0.472 ± 0.010	-1.6	256 ± 68	218 ± 61	14.7	0.297 ± 0.019	$0.292~\pm~0.018$	-1.7		
relation	0.6	0.490 ± 0.009	0.479 ± 0.009	-2.2	197 ± 39	187 ± 39	5.0	0.308 ± 0.002	$0.298~\pm~0.000$	-3.2		
relation	0.8	0.502 ± 0.014	0.485 ± 0.011	-3.3	169 ± 17	162 ± 12	4.4	0.304 ± 0.005	$0.291~\pm~0.004$	-4.1		
Max boos	st (%)			0.0			14.7			0.1		
Mean boo	st (%)			-1.8			6.5			-2.3		
					WD	50K		~				
triple	0.4	0.283 ± 0.001	0.330 ± 0.000	16.7	1713 ± 42	585 ± 3	65.8	0.164 ± 0.001	0.189 ± 0.001	15.5		
triple	0.6	0.353 ± 0.000	$\textbf{0.369} \pm 0.001$	4.5	931 ± 20	503 \pm 5	45.9	0.208 ± 0.001	0.213 ± 0.000	2.8		
triple	0.8	0.400 ± 0.001	0.399 ± 0.001	-0.2	$667~\pm~6$	$446~\pm~4$	33.1	0.240 ± 0.001	$0.235~\pm~0.000$	-2.0		
node	0.8	0.440 ± 0.002	0.431 ± 0.001	-2.0	$520~\pm~2$	380 ± 1	27.0	0.271 ± 0.001	$0.259~\pm~0.002$	-4.6		
relation	0.4	0.413 ± 0.023	$\textbf{0.430} \pm 0.029$	4.0	768 ± 123	$343~\pm~97$	55.3	0.249 ± 0.025	0.254 ± 0.027	1.9		
relation	0.6	0.437 ± 0.014	$\textbf{0.439} \pm 0.007$	0.6	617 ± 42	$337~\pm~64$	45.4	0.269 ± 0.013	$0.263~\pm~0.008$	-2.2		
relation	0.8	0.447 ± 0.016	0.441 ± 0.008	-1.3	541 ± 83	336 ± 67	37.8	0.275 ± 0.010	$0.266~\pm~0.005$	-3.5		
Max boos	st (%)			16.7			65.8			15.5		
Mean boo	st (%)			3.2			44.4			1.1		

Results: Synthetic scenario

Results on WN18RR





Results on FB15k-237

Results: Synthetic scenario



22

Results: Real-world scenario

		1	$MR\downarrow$				$MRR\uparrow$						
dataset	р	single	linl	ked	boost(%)	single	lin	ked	boost(%)	single	link	ked	boost(%)
			CN	FB			CN	FB			CN	FB	
WN18RR	0.4	0.347 ± 0.005	0.395	0.338	14.0	7681 ± 62	1508	4050	80.4	0.270 ± 0.003	0.259	0.262	-3.2
WN18RR	0.6	0.446 ± 0.001	0.471	0.435	5.7	4908 ± 172	963	2685	80.4	0.342 ± 0.003	0.330	0.330	-3.3
WN18RR	0.8	0.525 ± 0.004	0.527	0.504	0.3	2685 ± 155	750	1725	72.1	0.416 ± 0.005	0.390	0.400	-3.8
Max boost	(%)				14.0				80.4				-3.2
Mean boos	t (%)				6.7				77.6				-3.5
			CN	Y3-10			CN	Y3-10			CN	Y3-10	
FB15K237	0.4	0.360 ± 0.002	0.330	0.335	-6.9	315 ± 4	315	310	1.4	0.204 ± 0.001	0.189	0.192	-5.6
FB15K237	0.6	0.398 ± 0.001	0.359	0.362	-9.1	242 ± 1	254	251	-3.9	0.227 ± 0.001	0.202	0.208	-8.6
FB15K237	0.8	0.442 ± 0.000	0.387	0.395	-10.6	194 ± 2	213	214	-9.5	0.254 ± 0.001	0.221	0.225	-11.5
Max boost	(%)				-6.9				1.4				-5.6
Mean boos	t (%)				-8.9				-4.0				-8.6
			FB	Y3-10			FB	Y3-10			FB	Y3-10	
WD50K	0.4	0.283 ± 0.001	0.313	0.284	10.8	1713 ± 42	758	961	55.8	0.164 ± 0.001	0.181	0.163	10.6
WD50K	0.6	0.353 ± 0.000	0.361	0.327	2.4	931 ± 20	558	644	40.0	0.208 ± 0.001	0.212	0.190	2.0
WD50K	0.8	0.400 ± 0.001	0.396	0.358	-0.9	$667~\pm~6$	450	510	32.6	0.240 ± 0.001	0.236	0.210	-1.8
Max boost (%) 10.8					55.8				10.6				
Mean boos	t (%)				4.1				42.8				3.6

Ablation study: loss



Comparative analysis of predicted scores on the training set

Statistics of the original datasets

		Friples		E	ntities	Relations			
dataset	train	val	test	train	val	test	train	val	test
WN18RR*	86835	2817	2923	40714	4835	4985	11	11	11
FB15k-237*†	272115	17526	20438	14505	9799	10317	237	237	237
WD50K*	164631	22429	45284	40107	16500	22732	473	297	347
ConceptNet [†]	3423004	0	0	1787373	0	0	47	47	47
YAGO3-10†	1079040	4978	4982	123143	7914	7906	37	37	37

Statistics of the sampled datasets

			Triples			Entities	Relations			
sampling	р	train	val	test	train	val	test	train	val	test
					WN18RR					
triple	0.4	34734 ± 0	$2817~\pm~0$	$2923~\pm~0$	30975 ± 57	4835 ± 0	4985 ± 0	11 ± 0	11 ± 0	11 ± 0
triple	0.6	$52101~\pm~0$	$2817~\pm~0$	$2923~\pm~0$	$36178~\pm~27$	$4835 \ \pm \ 0$	$4985 \ \pm \ 0$	11 ± 0	11 ± 0	11 ± 0
triple	0.8	$69468\ \pm\ 0$	$2817~\pm~0$	$2923~\pm~0$	39093 ± 28	$4835 \ \pm \ 0$	$4985 ~\pm~ 0$	11 ± 0	11 ± 0	11 ± 0
node	0.4	$28815~\pm~78$	$2817~\pm~0$	$2923~\pm~0$	$16286\ \pm\ 0$	$4835 \ \pm \ 0$	$4985 ~\pm~ 0$	11 ± 0	11 ± 0	11 ± 0
node	0.6	$46715~\pm~90$	$2817~\pm~0$	$2923~\pm~0$	24428 ± 0	$4835 \ \pm \ 0$	$4985 \ \pm \ 0$	11 ± 0	11 ± 0	11 ± 0
node	0.8	$66311~\pm~157$	$2817~\pm~0$	$2923~\pm~0$	32571 ± 0	$4835 \ \pm \ 0$	$4985 \ \pm \ 0$	11 ± 0	11 ± 0	11 ± 0
relation	0.4	20367 ± 18482	511 ± 625	545 ± 616	14485 ± 10417	933 ± 1168	989 ± 1150	4 ± 0	4 ± 0	4 ± 0
relation	0.6	57790 ± 19777	1765 ± 624	1827 ± 635	32006 ± 7832	3135 ± 1000	3228 ± 1025	7 ± 0	7 ± 0	7 ± 0
relation	0.8	61586 ± 17153	1852 ± 598	1902 ± 610	34809 ± 4855	$3286~\pm~949$	$3353~\pm~974$	9 ± 0	9 ± 0	9 ± 0
					FB15k-237					
triple	0.4	$108846\ \pm\ 0$	$17526~\pm~0$	$20438~\pm~0$	14187 ± 10	9799 ± 0	$10317\ \pm\ 0$	237 ± 0	223 ± 0	224 ± 0
triple	0.6	$163269\ \pm\ 0$	$17526~\pm~0$	$20438~\pm~0$	14352 ± 6	9799 ± 0	$10317\ \pm\ 0$	237 ± 0	223 ± 0	224 ± 0
triple	0.8	$217692 \ \pm \ 0$	$17526~\pm~0$	$20438~\pm~0$	14439 ± 7	9799 ± 0	$10317\ \pm\ 0$	237 ± 0	223 ± 0	224 ± 0
relation	0.4	103999 ± 16860	6364 ± 891	7419 ± 1056	12884 ± 220	$5457~\pm~707$	5981 ± 722	95 ± 0	87 ± 2	89 ± 3
relation	0.6	156537 ± 20322	9562 ± 1011	11186 ± 1268	$13779~\pm~352$	$7218~\pm~571$	$7747~\pm~589$	142 ± 0	133 ± 3	134 ± 3
relation	0.8	224750 ± 25236	14309 ± 1632	16748 ± 1881	$14332\ \pm\ 75$	$8725~\pm~633$	$9291~\pm~618$	190 ± 0	178 ± 2	180 ± 2
					WD50K					
triple	0.4	$65852\ \pm\ 0$	$22429~\pm~0$	$45284~\pm~0$	34164 ± 32	16500 ± 0	22732 ± 0	434 ± 2	297 ± 0	347 ± 0
triple	0.6	$98779\ \pm\ 0$	$22429~\pm~0$	$45284~\pm~0$	37091 ± 41	$16500\ \pm\ 0$	22732 ± 0	453 ± 6	297 ± 0	347 ± 0
triple	0.8	131705 ± 0	$22429~\pm~0$	$45284~\pm~0$	$38957\ \pm\ 16$	16500 ± 0	22732 ± 0	464 ± 3	297 ± 0	347 ± 0
node	0.8	147774 ± 42	$22429~\pm~0$	$45284~\pm~0$	32086 ± 0	16500 ± 0	22732 ± 0	$ 450 \pm 1 $	297 ± 0	347 ± 0
relation	0.4	62701 ± 11210	7865 ± 1591	16014 ± 3259	25204 ± 2662	$7700~\pm~956$	11535 ± 1186	$ 189 \pm 0 $	112 ± 6	127 ± 5
relation	0.6	88487 ± 12893	11529 ± 1838	23332 ± 3923	29090 ± 3143	10258 ± 1183	14848 ± 1580	284 ± 0	173 ± 5	201 ± 6
relation	0.8	121499 ± 15858	16250 ± 2317	32947 ± 4969	33181 ± 2970	12848 ± 1161	18053 ± 1489	$ 378 \pm 0 $	238 ± 4	275 ± 5