

CoDEx: A Comprehensive Knowledge Graph Completion Benchmark

Tara Safavi

University of Michigan, Ann Arbor
CDL Meetup April 2021

 [@tararootcake](https://twitter.com/tararootcake)

 tsafavi.github.io





CoDEX: A Comprehensive Knowledge Graph Completion Benchmark

Tara Safavi
University of Michigan
tsafavi@umich.edu

Danai Koutra
University of Michigan
dkoutra@umich.edu

Abstract

We present CoDEX, a set of knowledge graph

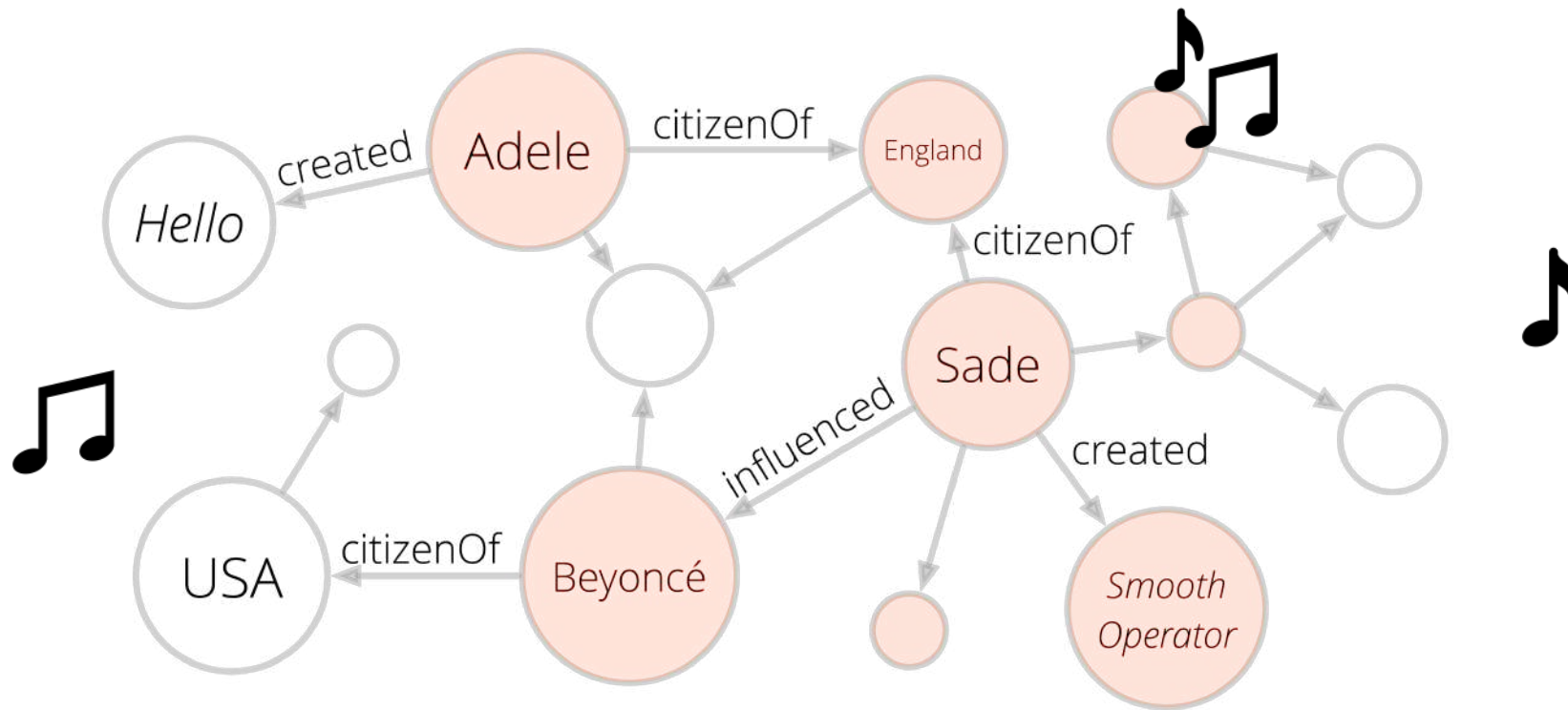
et al., 2008) are most commonly used for evaluation in KGC, even though Freebase had known quality issues (Tanon et al., 2016) and was eventu-

How do we measure
the quality of a KG?



Knowledge graphs (KGs)

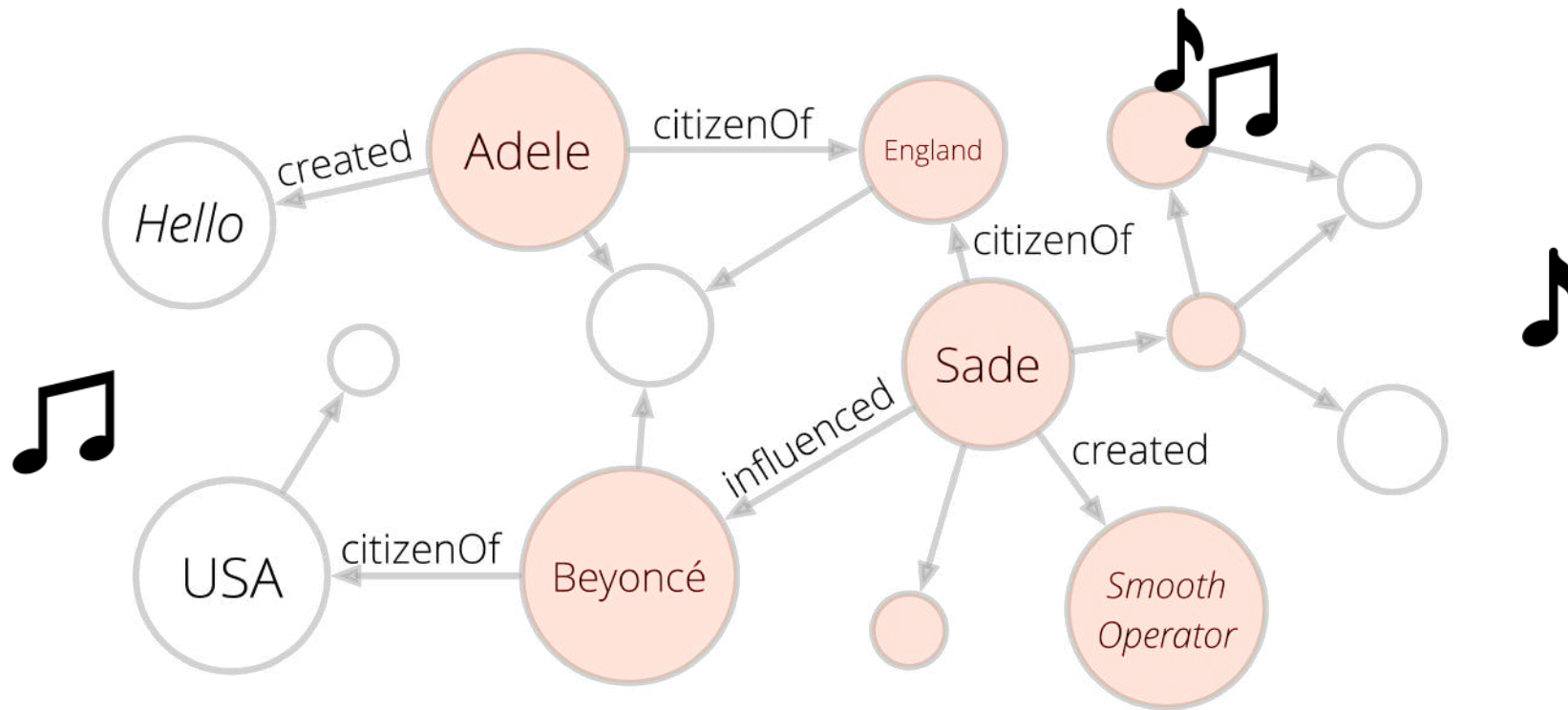
Symbolic entity-relation model of machine knowledge





What is the precision of KGs?

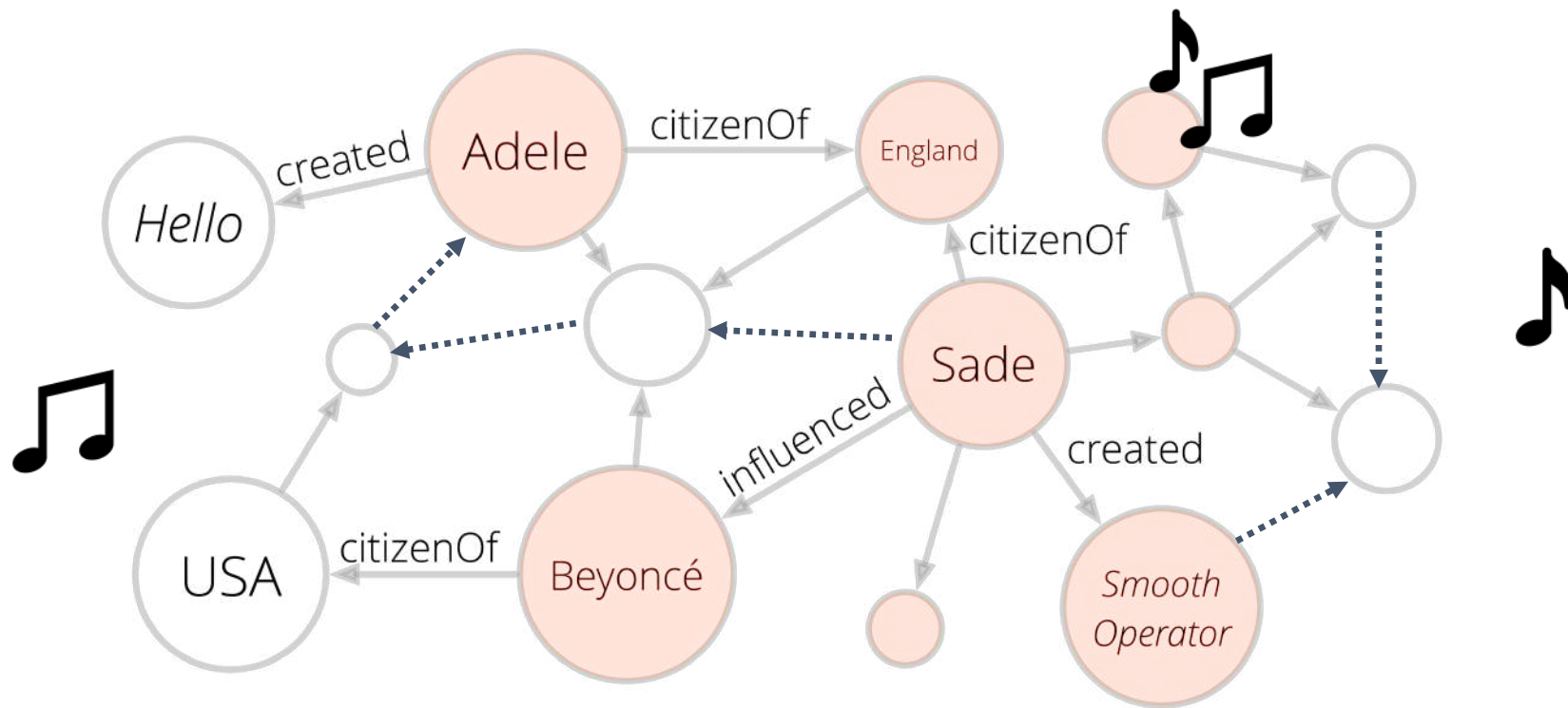
Constructed and/or verified by humans → high precision (accuracy)





What is the recall of KGs?

Always growing...but never truly “complete”





How to improve recall?



How to improve recall?

Hire experts or crowd workers



How to improve recall?

Hire experts or crowd workers

Information extraction from documents



How to improve recall?

Hire experts or crowd workers

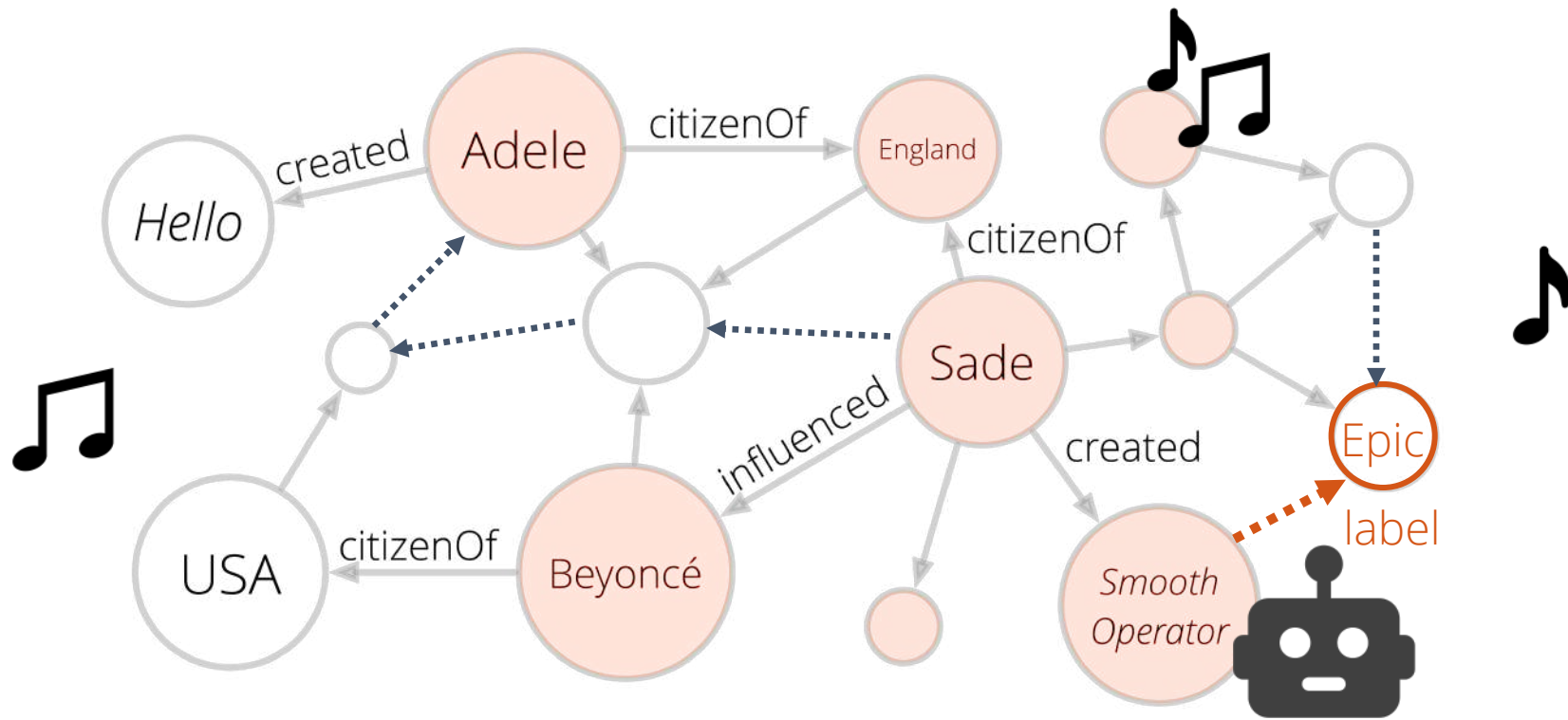
Information extraction from documents

Automate link prediction in KGs



How to improve recall?

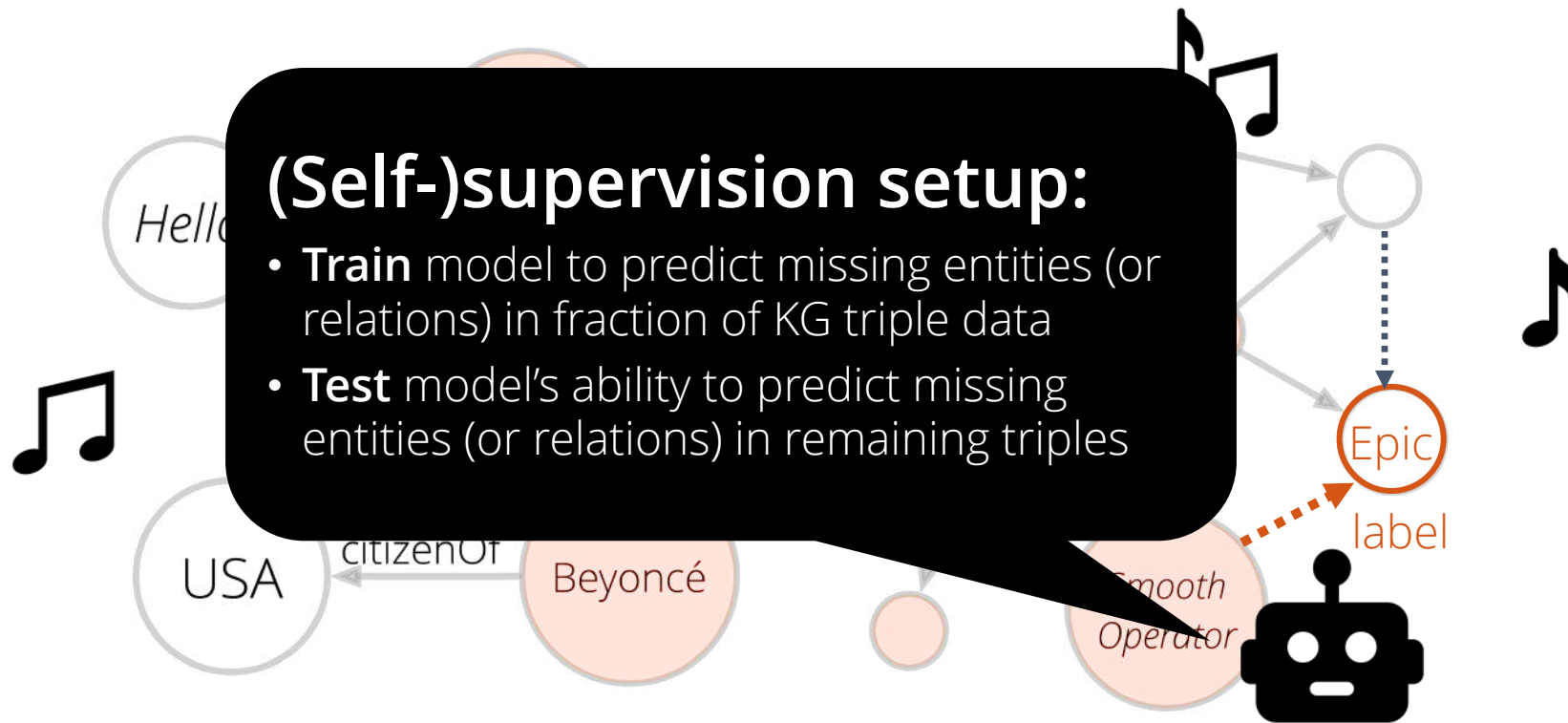
KG completion (KGC): Automate prediction of missing edges





Can we automate KGC?

Devise shallow/deep learning models for KGC





Can we automate KGC?

	Datasets						Evaluation tasks		
	FB15K	FB15K-237	FB13	WN18	WN18RR	WN11	Other	Link pred.	Triple class.
AAAI, IJCAI	✓	✓	✓	✓	✓		FB5M	✓	✓
	✓	✓	✓	✓	✓		FB40K	✓	✓
							NELL (Location, Sports)		
	✓		✓				Countries	✓	
							FB24K	✓	
	✓		✓					✓	
	✓	✓	✓	✓	✓			✓	✓
	✓	✓	✓	✓	✓			✓	✓
	✓						FB15K+	✓	✓
	✓						SemMedDB, DBPedia	✓	fact checking (not on FB15K)
	✓	✓	✓	✓	✓		YAGO3-10, Countries	✓	
	✓		✓					✓	
	✓						YAGO37	✓	
	✓	✓	✓	✓	✓			✓	
	✓		✓				YAGO3-10	✓	
ICML, ICLR, NeurIPS	✓		✓				FB15K-401	✓	rule extraction (FB15K-401)
	✓		✓					✓	
	✓		✓					✓	
	✓		✓					✓	
	✓		✓					✓	
		✓		✓			NELL-995, UMLS, Kinship, Countries, WikiMovies	✓	QA (WikiMovies)
	✓	✓	✓	✓	✓		YAGO3-10	✓	
	✓	✓	✓				DBPedia-YAGO3, DBPedia-Wikidata	✓	entity alignment (DBPedia graphs)
	✓	✓	✓	✓	✓			✓	
	✓	✓	✓	✓	✓			✓	
	✓	✓	✓	✓	✓			✓	
	✓		✓					✓	
	✓		✓					✓	
	✓		✓				MUTAG, AM, PTC	✓	graph classification (MUTAG, AM, PTC)

AI/ML conferences

NLP conferences

(Ji et al., 2015)	✓	✓	✓	✓				✓	✓
(Guo et al., 2015)							NELL (Location, Sports, Freq)	✓	✓
(Guu et al., 2015)		✓						✓	✓
(Garcia-Duran et al., 2015)	✓						Families	✓	
(Lin et al., 2015a)	✓						FB40K	✓	relation extraction (FB40K)
(Xiao et al., 2016b)	✓	✓	✓	✓	✓			✓	✓
(Nguyen et al., 2016)	✓		✓					✓	
(Xiong et al., 2017)		✓					NELL-995	✓	rule mining
(Lin et al., 2018)		✓		✓			NELL-995, UMLS, Kinship	✓	
(Nguyen et al., 2018)		✓		✓				✓	
(Bansal et al., 2019)		✓		✓				✓	
(Xu and Li, 2019)	✓	✓	✓	✓	✓		YAGO3-10, Family	✓	
(Balazevic et al., 2019b)	✓	✓	✓	✓	✓			✓	
(Vu et al., 2019)		✓		✓			SEARCH17	✓	personalized search (SEARCH17)
(Nathani et al., 2019)		✓		✓			NELL-995, UMLS, Kinship	✓	
(Jiang et al., 2019)	✓	✓	✓	✓	✓			✓	



But let's rewind...

Progress in AI research depends heavily on benchmarks

- **Benchmark:** Dataset of input/output pairs that sufficiently represents a real-world use case [\[Paullada et al 2020\]](#)
- Allows for comparison of competing systems (or algorithms) according to given metric(s)



But let's rewind...

What do benchmarks look like in KGC research?





Most existing KGC benchmarks*

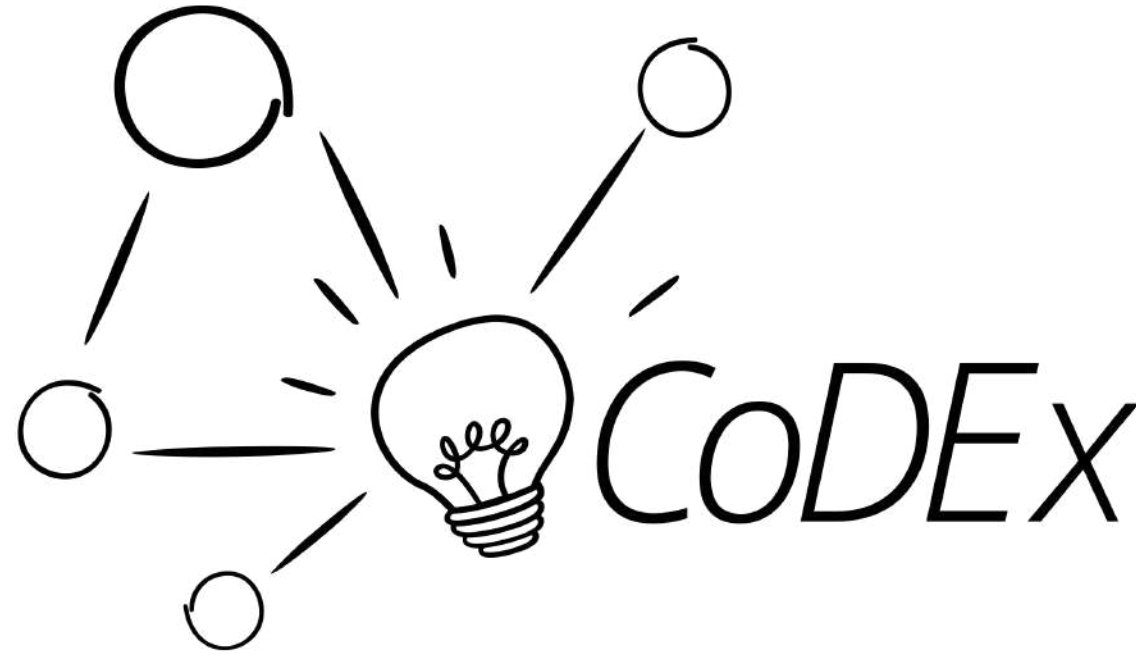
	Reference	Dataset					Evaluation tasks		
		FB15K	FB15K-237	FB13	WN18	Other	Long query	Triple classification	Other
AAAI/JLCCM	(Wang et al., 2014)	✓	✓	✓	✓	FB5M	✓	✓	relation extraction (FB5M)
	(Liu et al., 2015)	✓	✓	✓	✓	FB40K	✓	✓	relation extraction (FB40K)
	(Wang et al., 2015)					NELL (Location, Sports)	✓		
	(Nickel et al., 2016)	✓	✓			Comics	✓		
	(Liu et al., 2016)	✓				FB24K	✓		
	(Wang and Cohen, 2016)	✓	✓	✓	✓		✓	✓	
	(Guo et al., 2016a)	✓	✓	✓	✓		✓	✓	
	(Guo et al., 2016b)	✓	✓	✓	✓		✓	✓	
	(Guo et al., 2016)	✓	✓	✓	✓	FB15K1	✓	✓	
	(Shi and Worring, 2017)	✓				SciMedDB, DBpedia	✓	✓	fact checking (not on FB15K)
	(Demmen et al., 2018)	✓	✓	✓	✓	VAGO3-10, Countries	✓	✓	
	(Ethio and Ichise, 2018)	✓	✓	✓	✓		✓	✓	
KDD/ICML/NeurIPS	(Guo et al., 2018)	✓	✓	✓	✓	YAGO17	✓	✓	
	(Zhang et al., 2020)	✓	✓	✓	✓		✓	✓	
	(Vashishth et al., 2020a)	✓	✓	✓	✓	YAGO3-10	✓	✓	
	(Yang et al., 2015)	✓	✓			FB15K-401	✓	✓	rule extraction (FB15K-401)
	(Orrison et al., 2016)	✓	✓				✓	✓	
	(Liu et al., 2017)	✓	✓				✓	✓	
	(Korrem and Pastic, 2018)	✓	✓				✓	✓	
	(Guo et al., 2018)	✓	✓	✓	✓	NELL-995, UMLS, Kinship, Countries, Wikititles	✓	✓	QA (Wikititles)
	(Lacoste et al., 2018)	✓	✓	✓	✓	YAGO3-10	✓	✓	
	(Guo et al., 2019)	✓	✓	✓	✓	DBpedia-YAGO3, DBpedia Wikidata	✓	✓	entity alignment (DBpedia, graphs)
	(Guo et al., 2019)	✓	✓	✓	✓		✓	✓	
	(Zhang et al., 2019)	✓	✓	✓	✓		✓	✓	
ACL/EMNLP/NAACL	(Balarise et al., 2019a)	✓	✓	✓	✓		✓	✓	personalized search (SEARCH117)
	(Vas et al., 2019)	✓	✓	✓	✓	SEARCH117	✓	✓	
	(Nishani et al., 2019)	✓	✓	✓	✓	NELL-995 UMLS, Kinship	✓	✓	
	(Jorg et al., 2019)	✓	✓	✓	✓		✓	✓	
	(Li et al., 2015)	✓	✓	✓	✓		✓	✓	
	(Guo et al., 2015)	✓	✓	✓	✓	NELL (Location, Sports, Freq)	✓	✓	
	(Guo et al., 2015)	✓	✓	✓	✓		✓	✓	
	(Guo et al., 2015)	✓	✓	✓	✓	Families	✓	✓	
	(Guo et al., 2015)	✓	✓	✓	✓	FB40K	✓	✓	relation extraction (FB40K)
	(Xiao et al., 2016)	✓	✓	✓	✓		✓	✓	
	(Nguyen et al., 2016)	✓	✓	✓	✓		✓	✓	
	(Cheng et al., 2017)	✓	✓	✓	✓	NELL-995	✓	✓	rule mining
	(Liu et al., 2018)	✓	✓	✓	✓	NELL-995 UMLS, Kinship	✓	✓	
	(Nguyen et al., 2018)	✓	✓	✓	✓		✓	✓	
	(Biswal et al., 2019)	✓	✓	✓	✓		✓	✓	
	(Ola and Li, 2019)	✓	✓	✓	✓	YAGO3-10, Family	✓	✓	
	(Balarise et al., 2019b)	✓	✓	✓	✓		✓	✓	
	(Vas et al., 2019)	✓	✓	✓	✓	SEARCH117	✓	✓	personalized search (SEARCH117)
	(Nishani et al., 2019)	✓	✓	✓	✓	NELL-995 UMLS, Kinship	✓	✓	
	(Jorg et al., 2019)	✓	✓	✓	✓		✓	✓	

AI/ML conferences

NLP conferences

- Reliance on outdated data sources
- Leakage between train and test splits
- Non-standardized versions
- Lack of difficult test examples
- Low interpretability for practitioners

*We survey 40+ KGC papers and 12 evaluation datasets across AI/ML/NLP venues



A set of knowledge graph
Completion **D**atasets **E**xtracted from
Wikidata and Wikipedia



+ everything (data, models, code) is publicly accessible!

How was CoDEX designed and collected?

What KGC tasks can I test on CoDEX?

How does CoDEX compare to existing KGC benchmarks?

How was CoDEx
designed and
collected?

What KGC tasks can I
test on CoDEx?

How does CoDEx
compare to existing
KGC benchmarks?

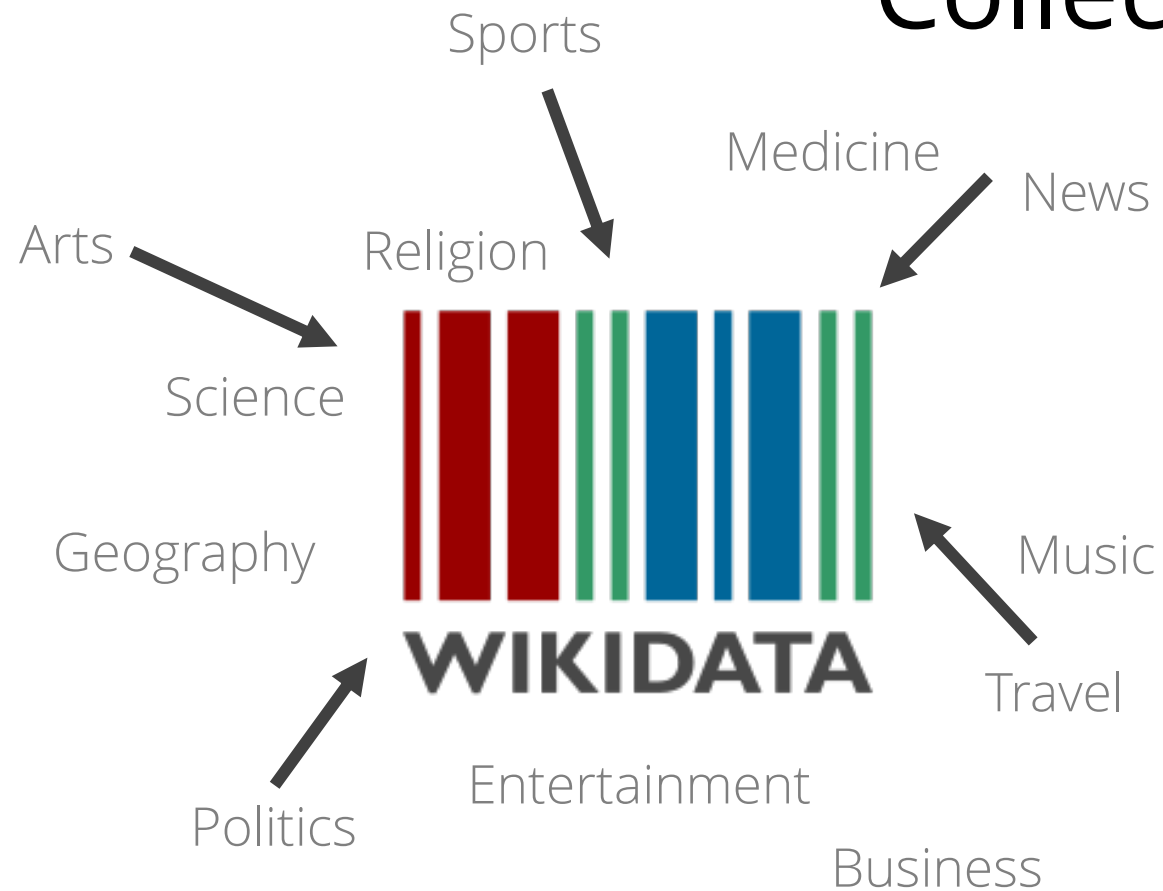


Collecting CoDEx



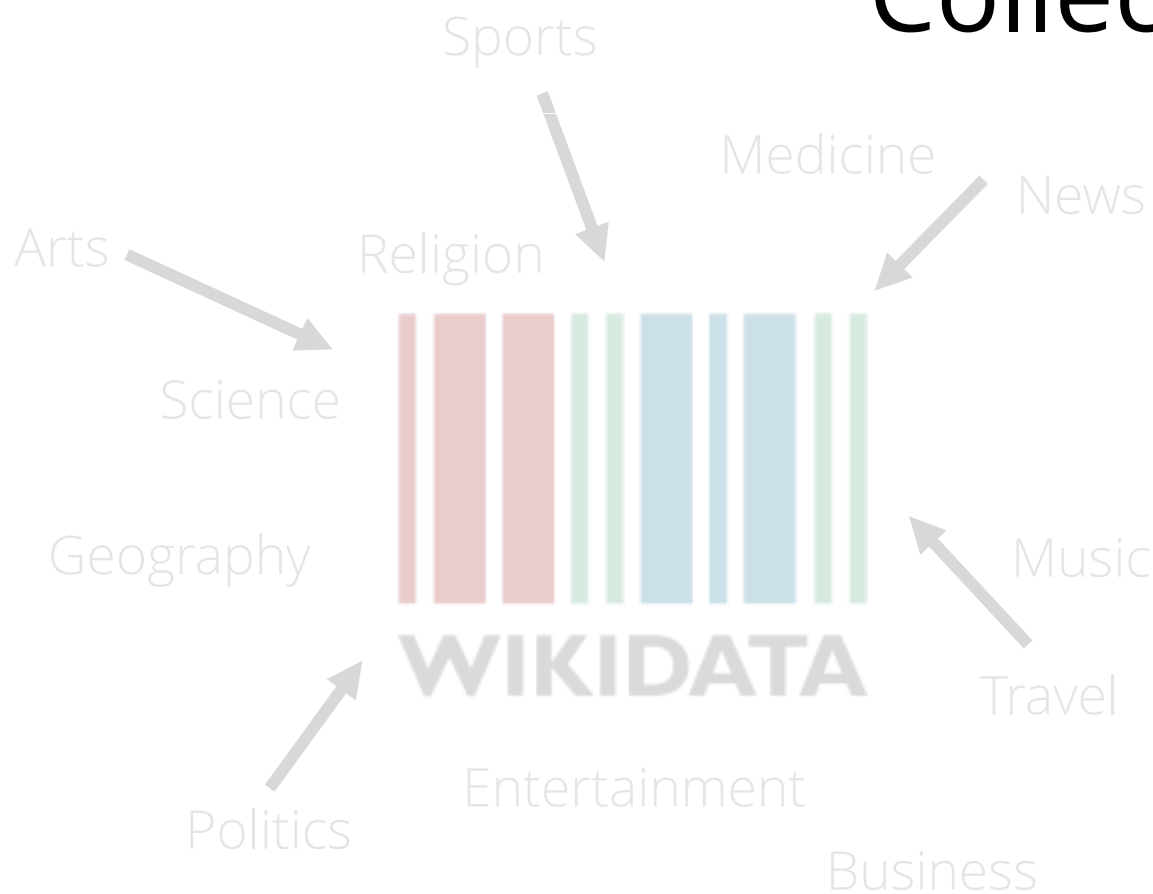


Collecting CoDEX





Collecting CoDEX



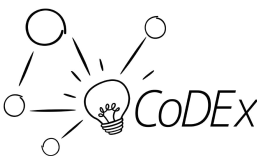
	# entities	# relations	# triples
CoDEX-S	2K	42	36K
CoDEX-M	17K	51	206K
CoDEX-L	78K	69	612K

```
import random
codex = Codex(code="en", size="m")

eid = random.choice(list(codex.entities()))
triples = codex.triples()
triples = triples[
    (triples["head"] == eid) | (triples["tail"] == eid)
]

for (head, relation, tail) in triples.values:
    print(f"({codex.entity_label(head)}, {codex.relation_label(relation)}, {codex.entity_label(tail)})")
```

(Virginia Woolf, country of citizenship, United Kingdom)
(Virginia Woolf, occupation, diarist)
(Virginia Woolf, occupation, feminist)
(Ursula K. Le Guin, influenced by, Virginia Woolf)
(Virginia Woolf, influenced by, George Eliot)
(Virginia Woolf, genre, prose)
(Virginia Woolf, occupation, essayist)
(Leonard Sidney Woolf, spouse, Virginia Woolf)
(Virginia Woolf, genre, drama)
(Samuel R. Delany, influenced by, Virginia Woolf)
(Virginia Woolf, languages spoken, written, or signed, English)





Collecting CoDEX





Collecting CoDEX

```
eid = "Q51"

for code in codes:
    codex = Codex(code=code)
    print(codex.entity_label(eid))

القارة القطبية الجنوبية
Antarktika
Antarctica
Antártida
Антарктида
南极洲

codex = Codex(code="en")
print(f"From {codex.entity_wikipedia_url(eid)}:")
print(f"  {codex.entity_extract(eid)[:400]}...")

From https://en.wikipedia.org/wiki/Antarctica:
'Antarctica ( or (listen)) is Earth's southernmost contin
e, almost entirely south of the Antarctic Circle, and is sur
est continent and nearly twice the size of Australia. At 0.0

codex = Codex(code="en")
types = codex.entity_types(eid)
for etype in types:
    print(codex.entity_label(eid), "is of type", codex.entit

Antarctica is of type continent
Antarctica is of type geographic region
```



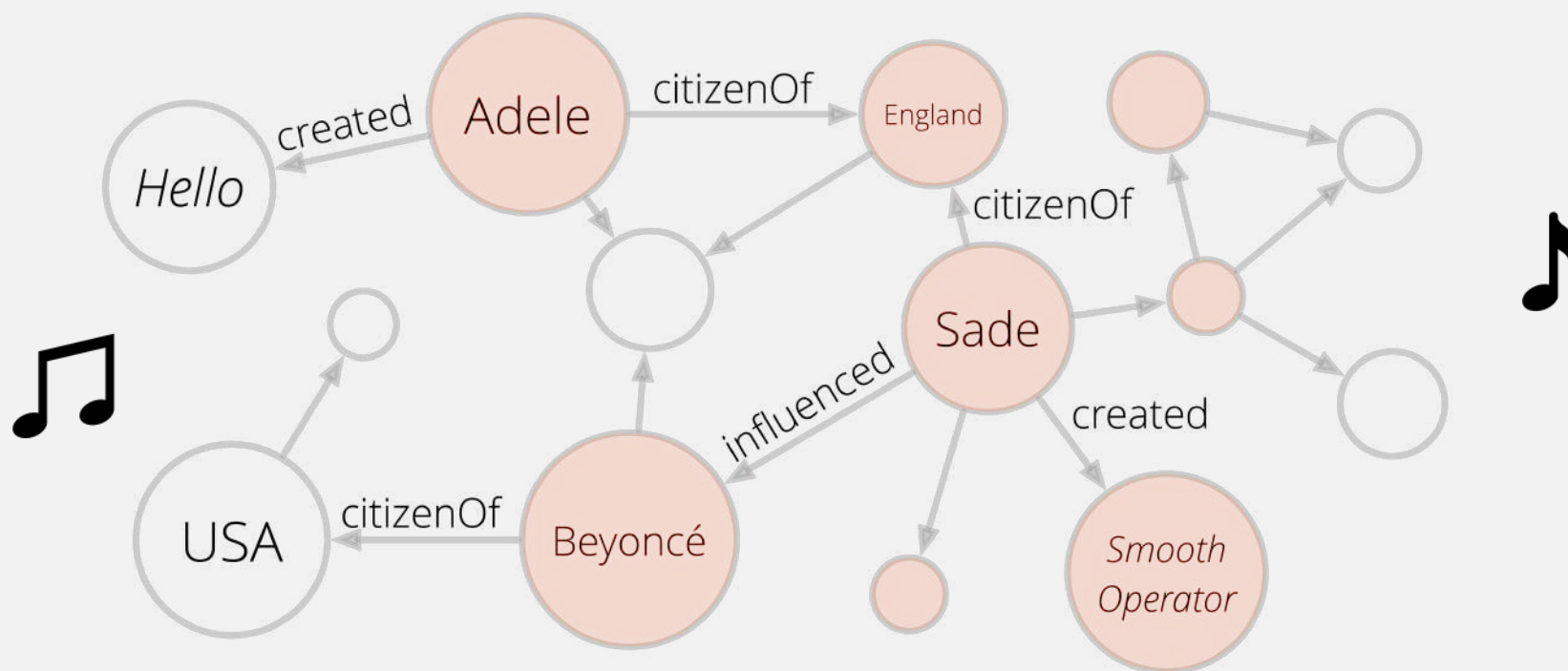
WIKIPEDIA

Entity types + text in Arabic,
German, English, Spanish,
Russian, Chinese

Existing KGC benchmarks in
English only (or don't have text)!

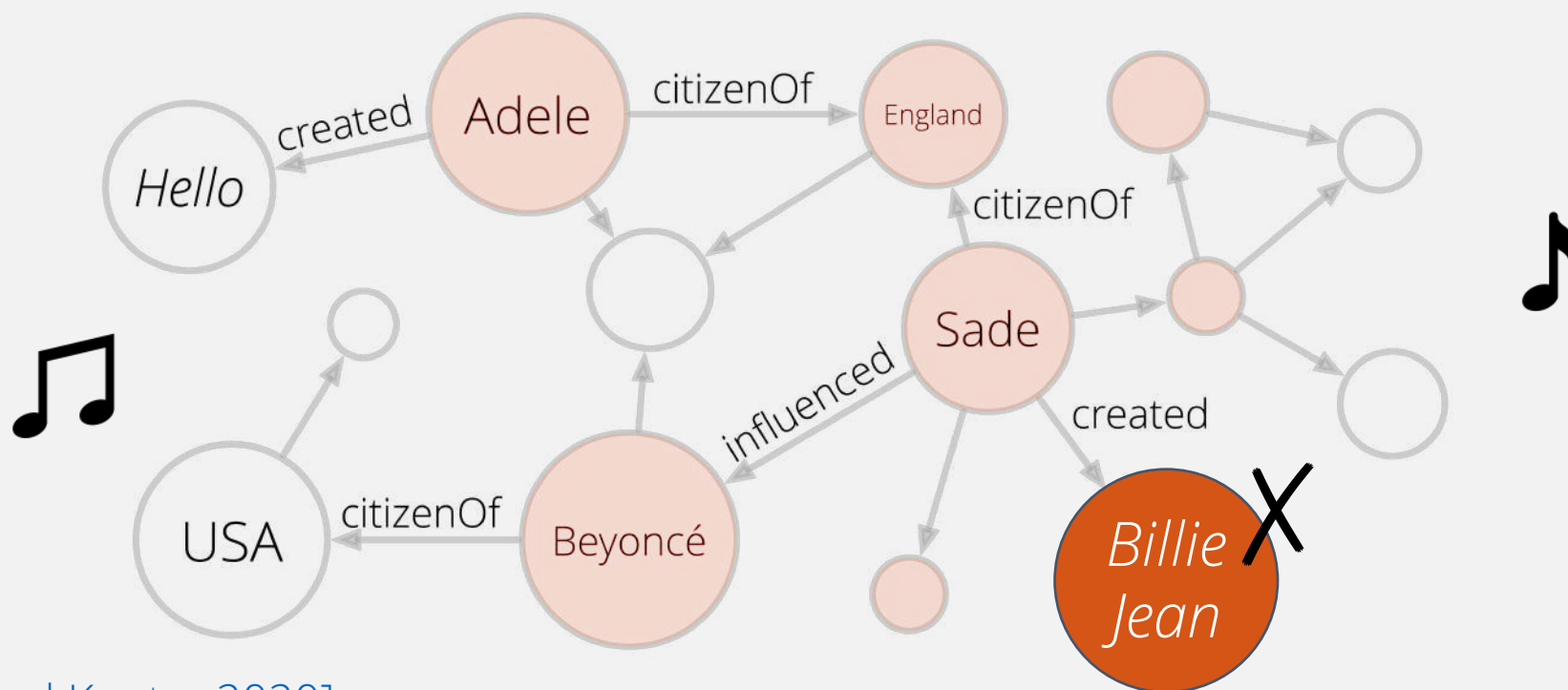
Positive and negative knowledge

Most KGs include positive (true) data only...



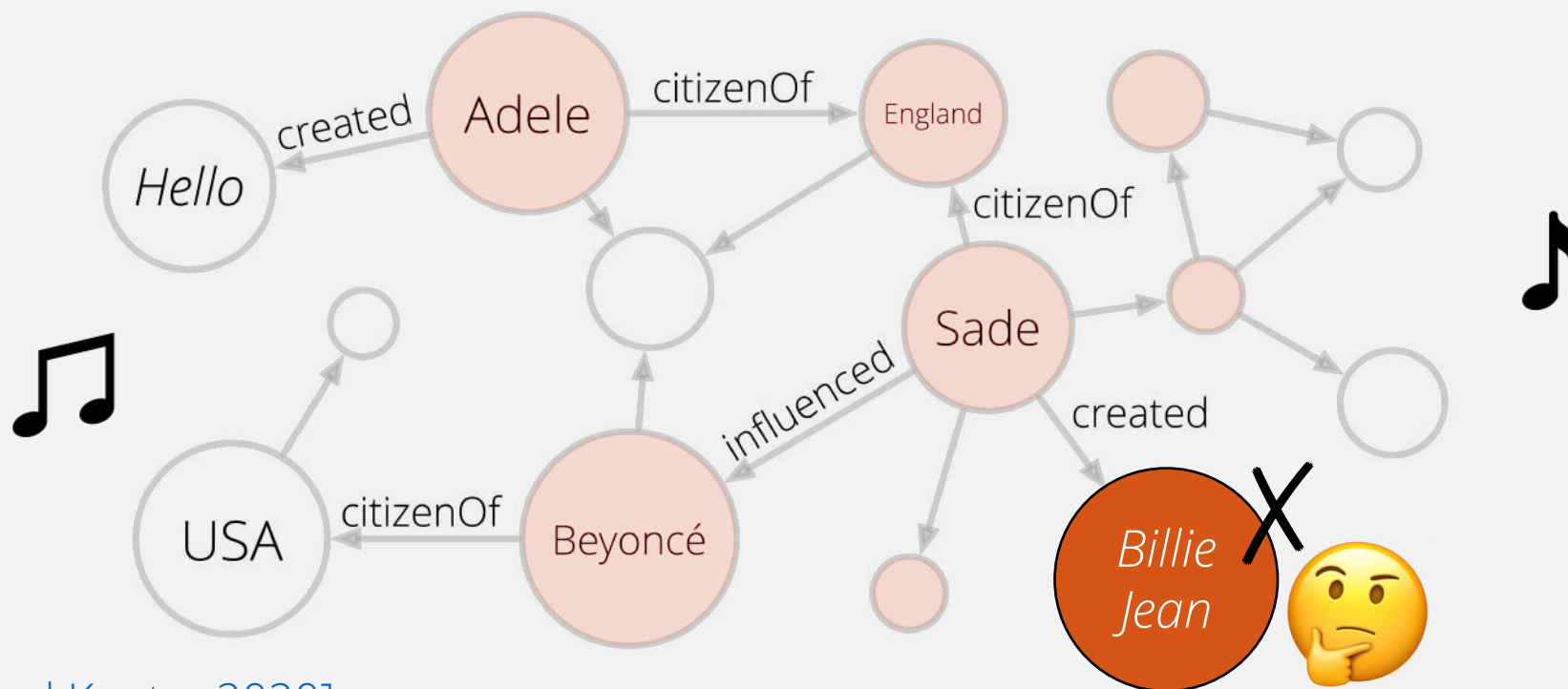
Positive and negative knowledge

But negative (false) knowledge can be useful too!



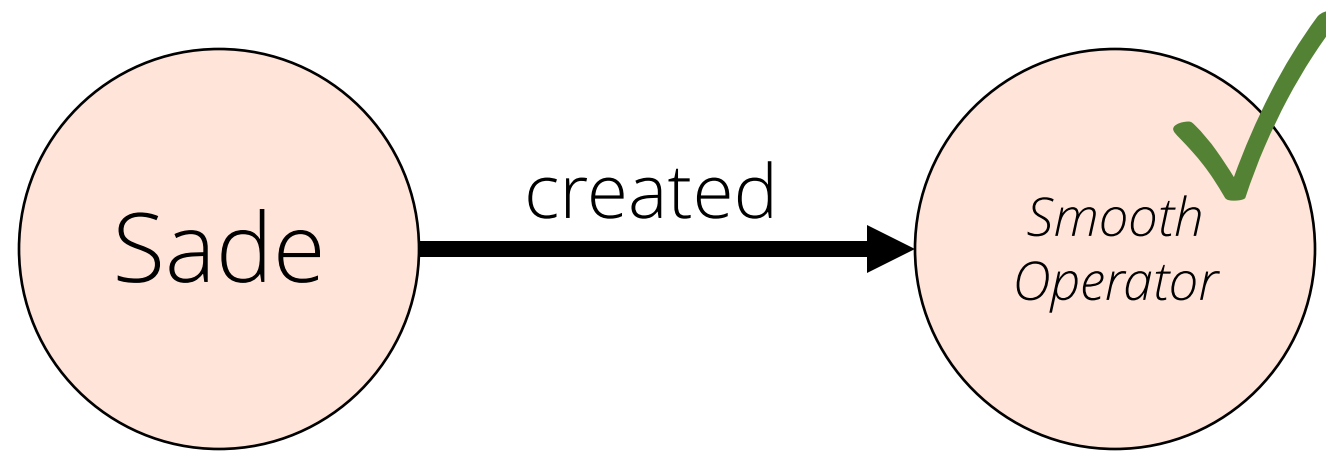
Positive and negative knowledge

How to construct negatives from positives?



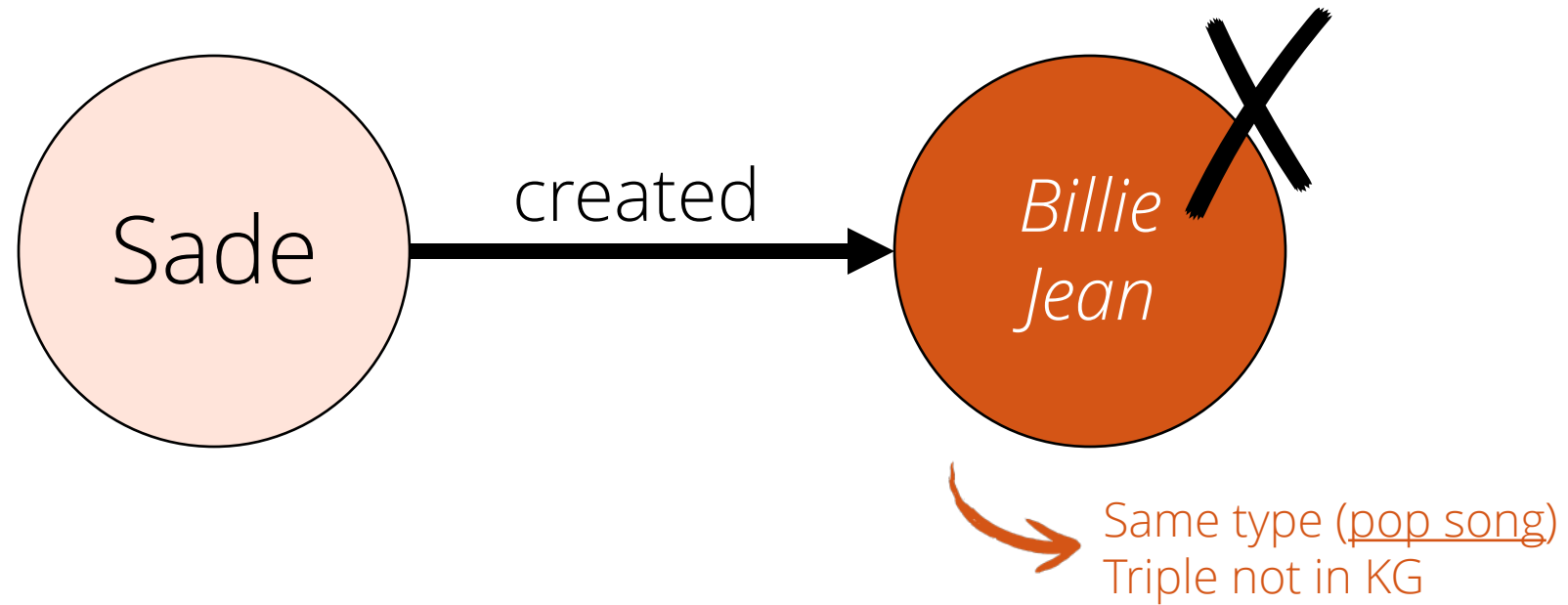


Augmenting CoDEx





Augmenting CoDEX





Augmenting CoDEx

Table 3: Selected examples of hard negatives in CoDEX with explanations.

Negative	Explanation
<i>(Frédéric Chopin, occupation, conductor)</i>	Chopin was a pianist and a composer, not a conductor.
<i>(Lesotho, official language, American English)</i>	English, not American English, is an official language of Lesotho.
<i>(Senegal, part of, Middle East)</i>	Senegal is part of West Africa.
<i>(Simone de Beauvoir, field of work, astronomy)</i>	Simone de Beauvoir’s field of work was primarily philosophy.
<i>(Vatican City, member of, UNESCO)</i>	Vatican City is a UNESCO World Heritage Site but not a member state.



Explore CoDEX.ipynb

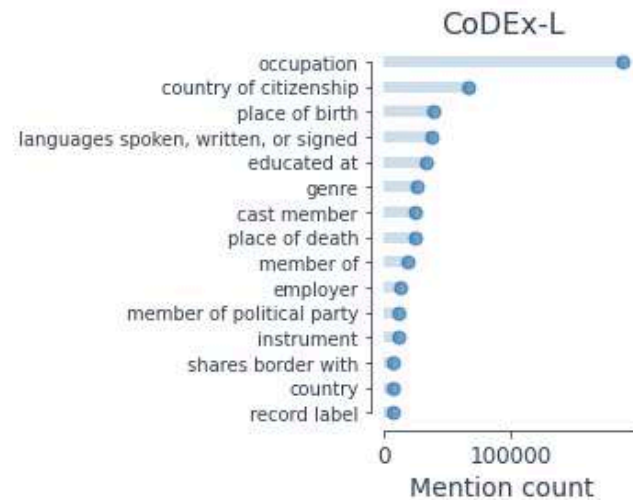
```
count_df = count_relations(triples)
count_df["label"] = [
    codex.relation_label(rid) for rid in count_df["relation"]]

k = 15

ax = plot_top_k(
    count_df,
    k=k,
    color=palette[-1],
    linewidths=6,
    figsize=(5, 4)
)

ax.set_xscale("linear")
ax.set_xlabel("Mention count", fontsize=14)
ax.set_title(codex.name(), fontsize=16)
ax.tick_params("x", labelsize=12)

plt.tight_layout()
plt.show()
```



How was CoDEx
designed and
collected?

What KGC tasks can I
test on CoDEx?

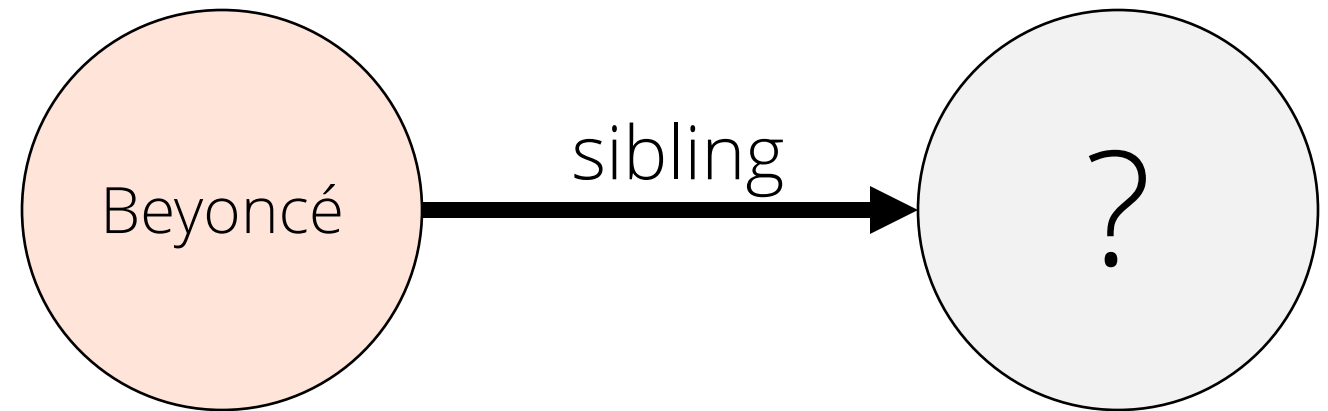
How does CoDEx
compare to existing
KGC benchmarks?



KGC tasks

Link prediction

Predict answers to queries like (head, relation, ?) and (?, relation, tail) by ranking candidates



Compute ranking metrics like mean reciprocal rank (MRR) and hits@k

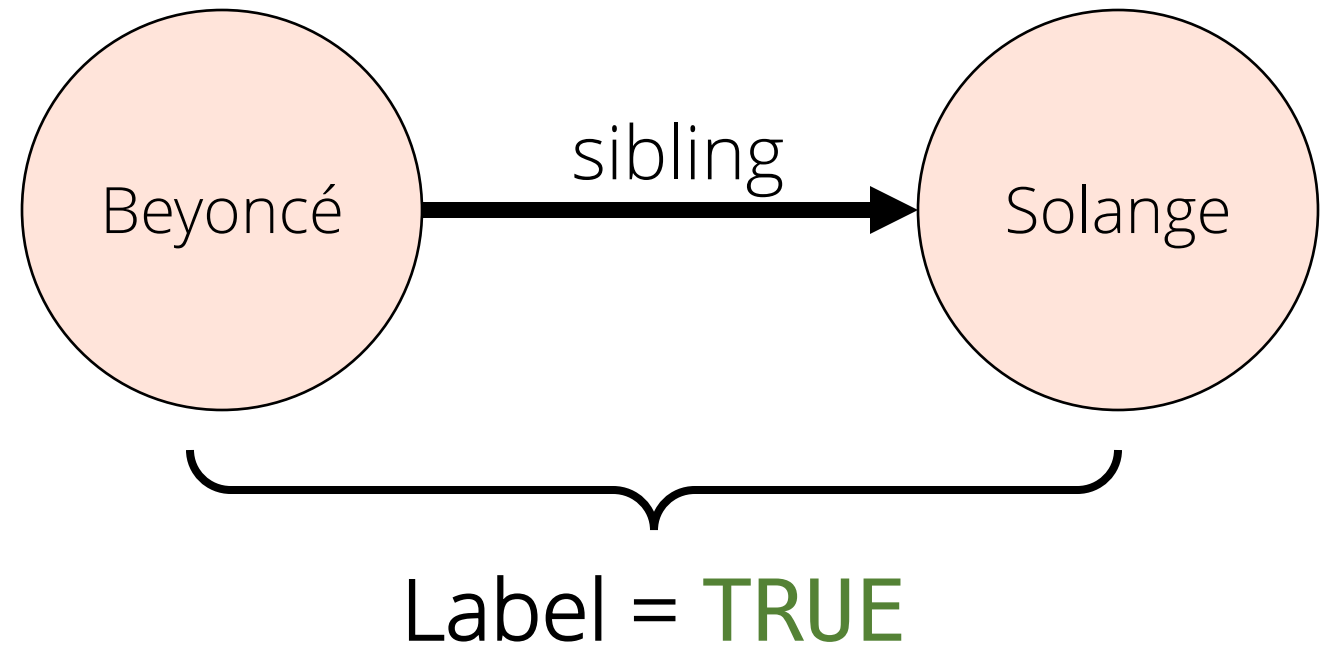


KGC tasks

Triple classification

Label provided triples as true or false

Compute classification metrics like accuracy





Models and implementation

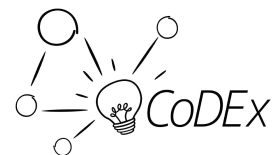
	Type	Approach	Citation
RESCAL	Linear	Tensor decomp.	[Nickel et al 2011]
TransE	Linear	Translational	[Bordes et al 2013]
Complex	Linear	Matrix decomp.	[Trouillon et al 2016]
ConvE	Nonlinear	Deep convolutions	[Dettmers et al 2018]
TuckER	Linear	Tensor decomp.	[Balažević et al 2019]



A knowledge graph embedding library



Important to
fairly compare
methods!!



[\[Broscheit et al 2020\]](#) <https://github.com/uma-pi1/kge>

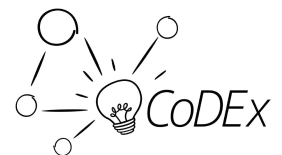


Link prediction results

Table 5: Comparison of link prediction performance on CoDEX.

	CoDEX-S			CoDEX-M			CoDEX-L		
	MRR	Hits@1	Hits@10	MRR	Hits@1	Hits@10	MRR	Hits@1	Hits@10
RESCAL									
TransE									
ComplEx									
ConvE									
TuckER									

Year
proposed





Link prediction results

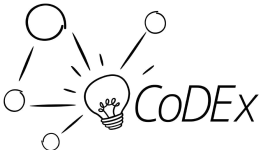
Table 5: Comparison of link prediction performance on CoDEX.

	CoDEX-S			CoDEX-M			CoDEX-L		
	MRR	Hits@1	Hits@10	MRR	Hits@1	Hits@10	MRR	Hits@1	Hits@10
RESCAL	0.404	0.293	0.623	0.317	0.244	0.456	0.304	0.242	0.419
TransE	0.354	0.219	0.634	0.303	0.223	0.454	0.187	0.116	0.317
ComplEx	0.465	0.372	0.646	0.337	0.262	0.476	0.294	0.237	0.400
ConvE	0.444	0.343	0.635	0.318	0.239	0.464	0.303	0.240	0.420
TuckER	0.444	0.339	0.638	0.328	0.259	0.458	0.309	0.244	0.430

Year
proposed
↓

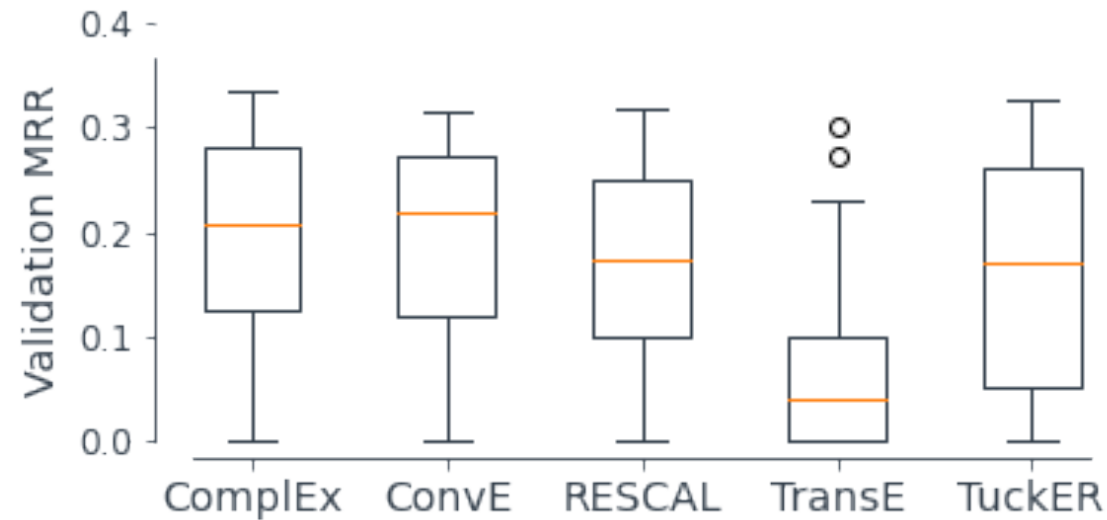
For CoDEX-S/M, earlier model (ComplEx) performs best when models are fairly compared!

Is “deeper” really better?





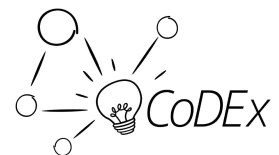
Link prediction results



Performance varies $\pm 30\%$ based on input hyperparameter configuration, consistent with the literature



Important to properly tune models!!





Triple classification results

Table 6: Comparison of triple classification performance on CoDEX by negative generation strategy.

	CoDEX-S						CoDEX-M					
	Uniform		Relative freq.		Hard neg.		Uniform		Relative freq.		Hard neg.	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
RESCAL												
TransE												
ComplEx												
ConvE												
TuckER												



Triple classification results

Table 6: Comparison of triple classification performance on CoDEX by negative generation strategy.

	CoDEX-S						CoDEX-M					
	Uniform		Relative freq.		Hard neg.		Uniform		Relative freq.		Hard neg.	
	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1	Acc.	F1
RESCAL	0.972	0.972	0.916	0.920	0.843	0.852	0.977	0.976	0.921	0.922	0.818	0.815
TransE	0.974	0.974	0.919	0.923	0.829	0.837	0.986	0.986	0.932	0.933	0.797	0.803
ComplEx	0.975	0.975	0.927	0.930	0.836	0.846	0.984	0.984	0.930	0.933	0.824	0.818
ConvE	0.972	0.972	0.921	0.924	0.841	0.846	0.979	0.979	0.934	0.935	0.826	0.829
TuckER	0.973	0.973	0.917	0.920	0.840	0.846	0.977	0.977	0.920	0.922	0.823	0.816

Accuracy drops up to 19 points on hard negative examples compared to randomly generated negatives



Lots of room for improvement!

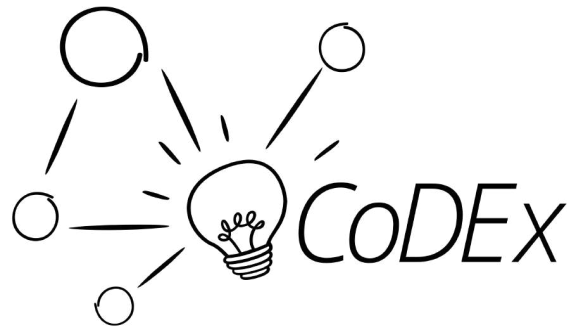
How was CoDEx
designed and
collected?

What KGC tasks can I
test on CoDEx?

How does CoDEx
compare to existing
KGC benchmarks?



How does CoDEX compare to existing benchmarks?

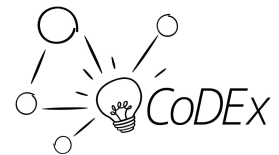


vs

 **Freebase™**

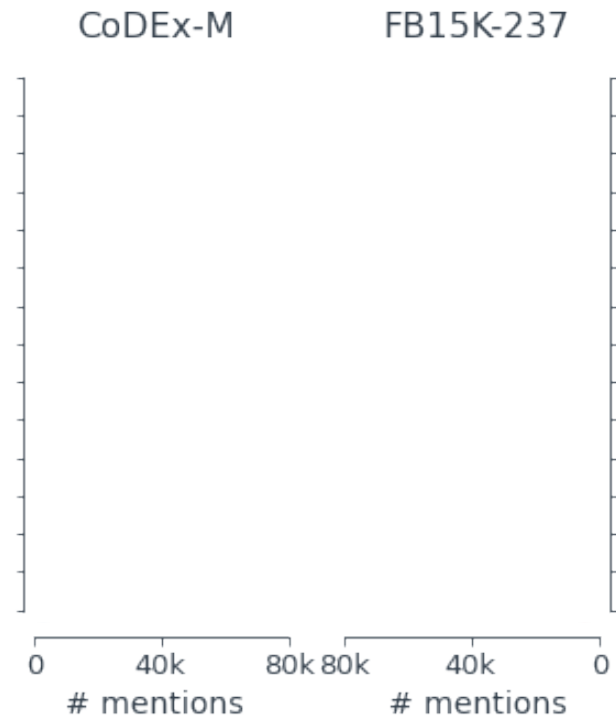
FB15K-237

[\[Toutanova and Chen 2015\]](#)



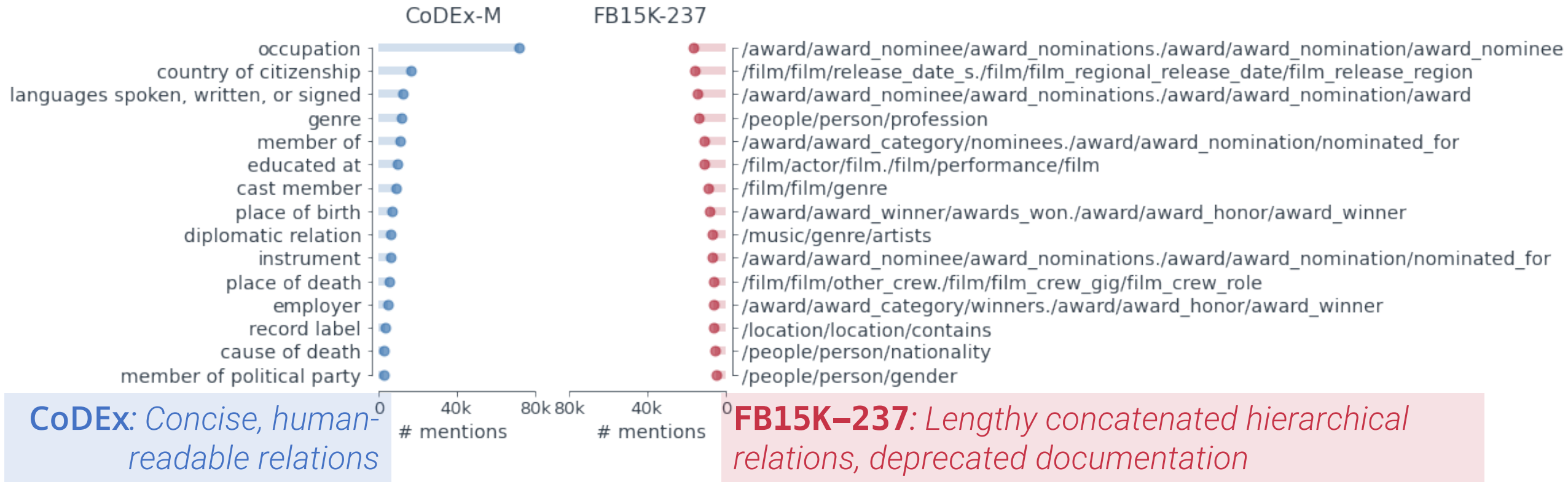


Qualitative comparison

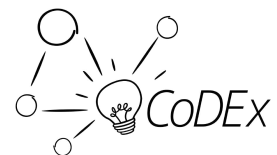




Qualitative comparison



CoDEX covers a wider selection of topics
and is easier to interpret





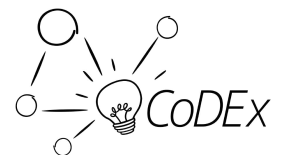
Quantitative comparison

Link prediction

*Compare simple non-learning link prediction baseline to SOTA KGC model per dataset, compute **improvement** over baseline*

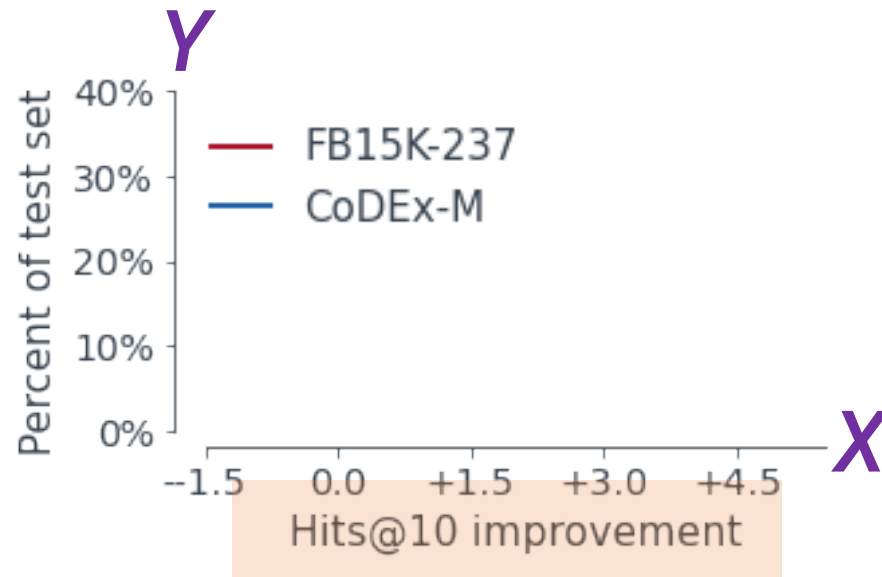
FB15K-237: RESCAL [Nickel et al 2011]

CoDEX: ComplEx [Trouillon et al 2016]





Quantitative comparison



Read as: "Learned KGC model improves $\leq X$ pts over baseline for $Y\%$ of test set"

Link prediction

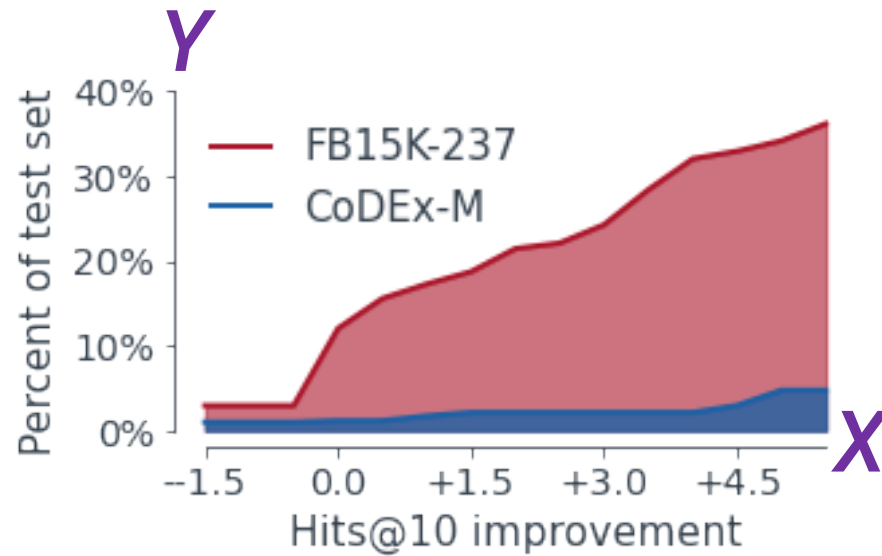
Compare simple non-learning link prediction baseline to SOTA KGC model per dataset, compute *improvement* over baseline

FB15K-237: RESCAL [Nickel et al 2011]

CoDEx: ComplEx [Trouillon et al 2016]



Quantitative comparison



Read as: "Learned KGC model improves $\leq X$ pts over baseline for $Y\%$ of test set"

FB15K-237 easier?

- » Baseline performs better than KGC model for 10%
- » Baseline within 5 pts of KGC model for 40%

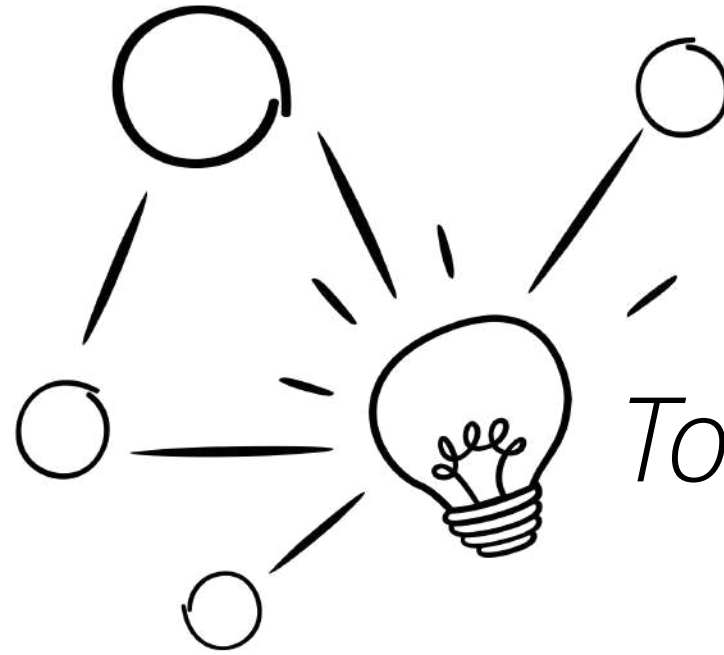
FB15K-237: very biased toward high-degree entities ("USA", "male")

CoDEX: Designed to avoid this!

How was CoDEx
designed and
collected?

What KGC tasks can I
test on CoDEx?

How does CoDEx
compare to existing
KGC benchmarks?



To recap, we:

Introduced and described CoDEx

Benchmarked KGC models on CoDEx for two tasks

Showed its value over a popular KGC dataset



Thanks for watching!



✉ tsafavi@umich.edu

🐦 [@tararootcake](https://twitter.com/tararootcake)

🌐 tsafavi.github.io

🐙 github.com/tsafavi/codex