

High-quality Task Division for Large-scale Entity Alignment

Bing Liu¹, Wen Hua¹, Guido Zuccon¹, Genghong Zhao², Xia Zhang²

The University of Queensland¹, Neusoft²

✉ bing.liu@uq.edu.au

🌐 <https://uqbingliu.github.io/>

🐦 @BingLiu1011

The logo for Neusoft, consisting of the word "Neusoft" in white sans-serif font on a blue rectangular background.

Entity Alignment aims to match **equivalent entities** in different **Knowledge Graphs** (KGs).

- One entity might be called differently in different scenarios.

CIKM2022

31st ACM International Conference on Information and Knowledge Management

October 17–21, 2022

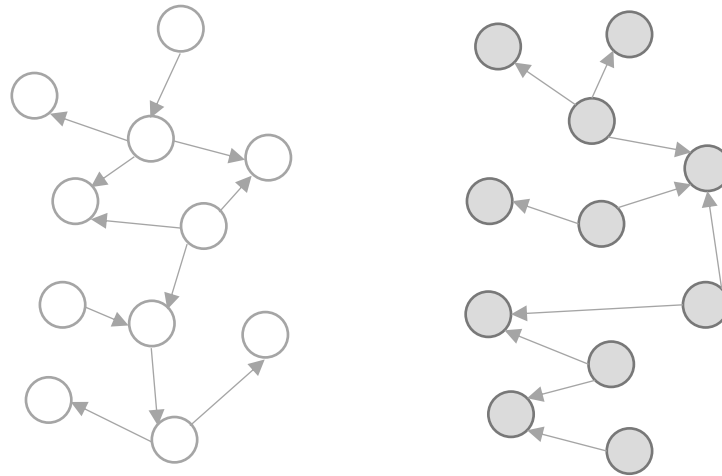
Hybrid Conference, Hosted in Atlanta, Georgia, USA

Registration

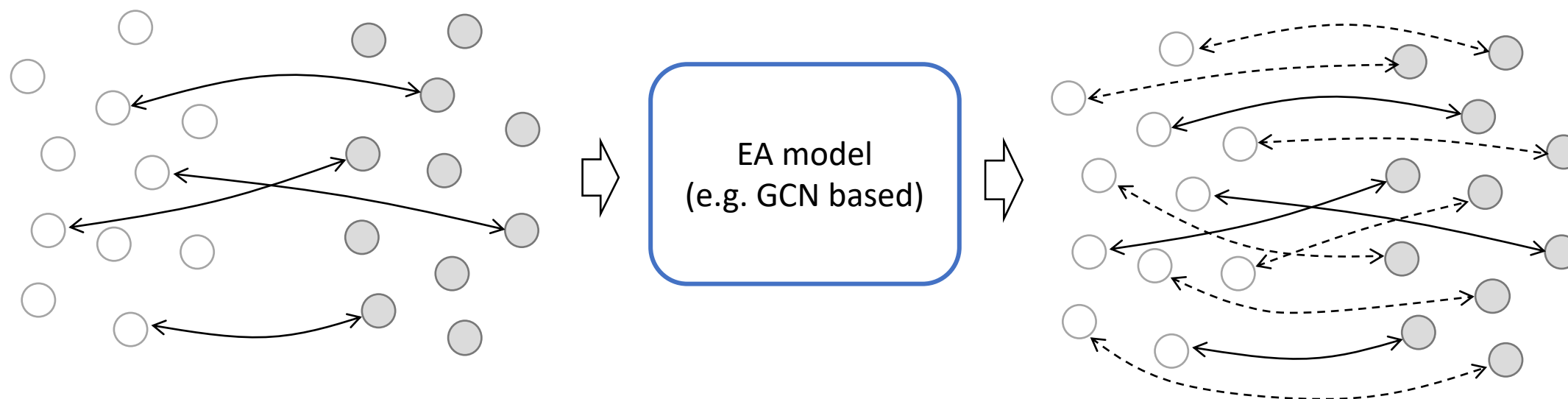


Flight						Modify Remove	\$89.96
✈	Fri 1/10	DAL 6:00AM	→ ATL 9:00AM	2hr 0min	Nonstop	Wanna Get Away	Price per passenger \$57.08 Taxes and fees per passenger \$32.88
✈	Sun 1/12	ATL 9:35PM	→ DAL 10:55PM	2hr 20min	Nonstop	Wanna Get Away	Total per passenger \$89.96 Passenger(s) x1 Flight total \$89.96

Entity Alignment aims to match **equivalent entities** in different **Knowledge Graphs** (KGs).

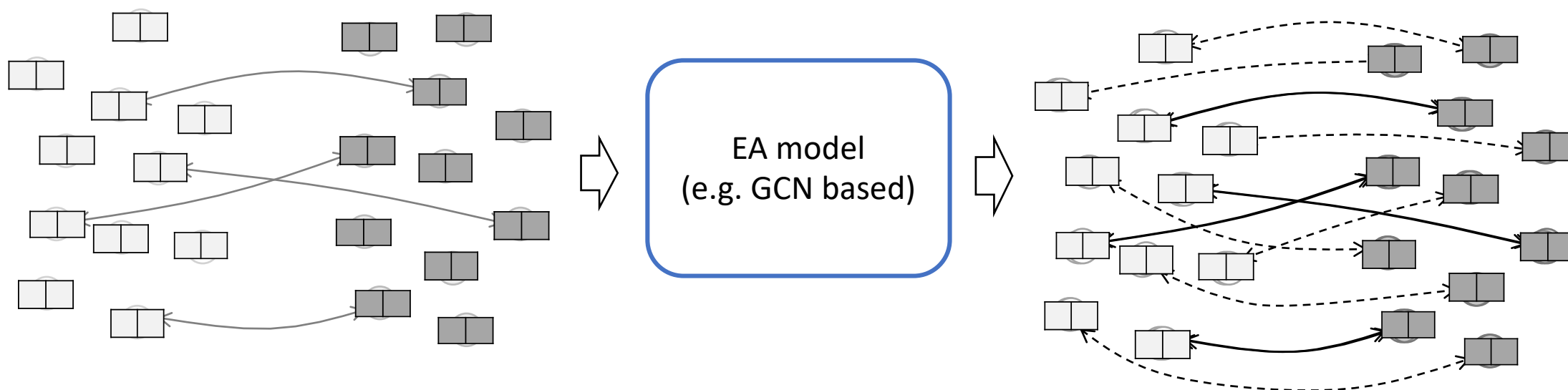


- Some **seed mappings** are provided as training data.
- Neural model encodes entities into **embeddings**.
- Predict **potential mappings**.



Neural EA models cannot be applied to large-scale KGs

- **Out-Of-Memory** (GPU).
- Time **efficiency**.



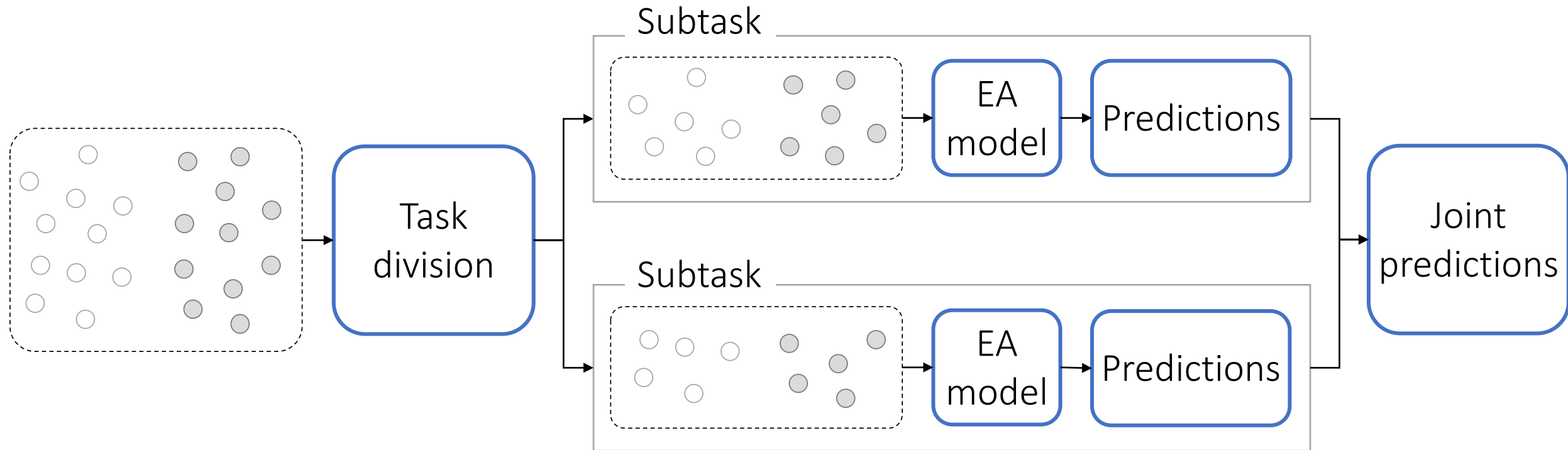
- **Entity-related parameters**
(initial entity representations)

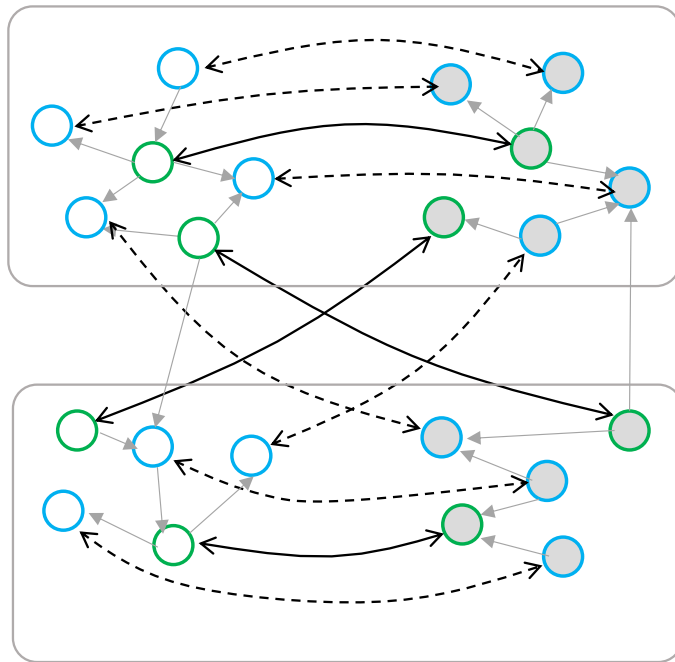
- Other parameters
- Neural operations

- **Entity representations**

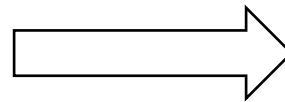
Divide a large-scale EA task into multiple **small subtasks**

- Each subtask only has two **small subgraphs** to align.



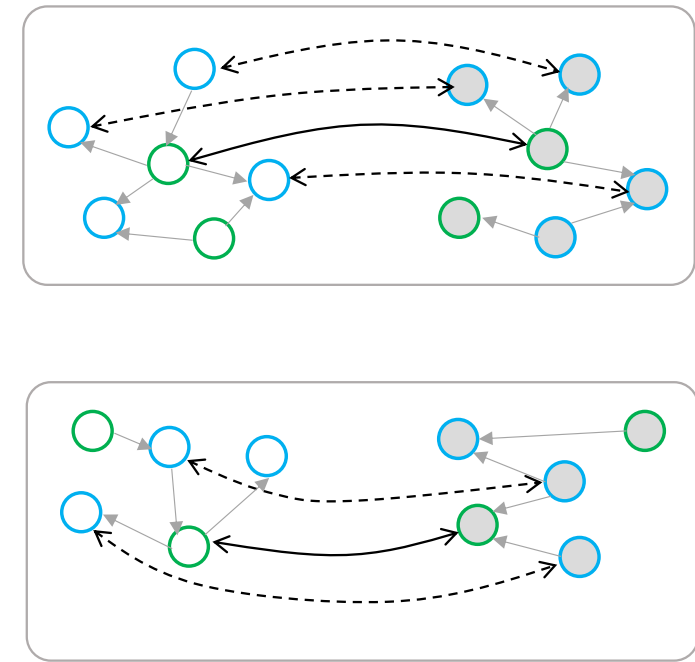


Division



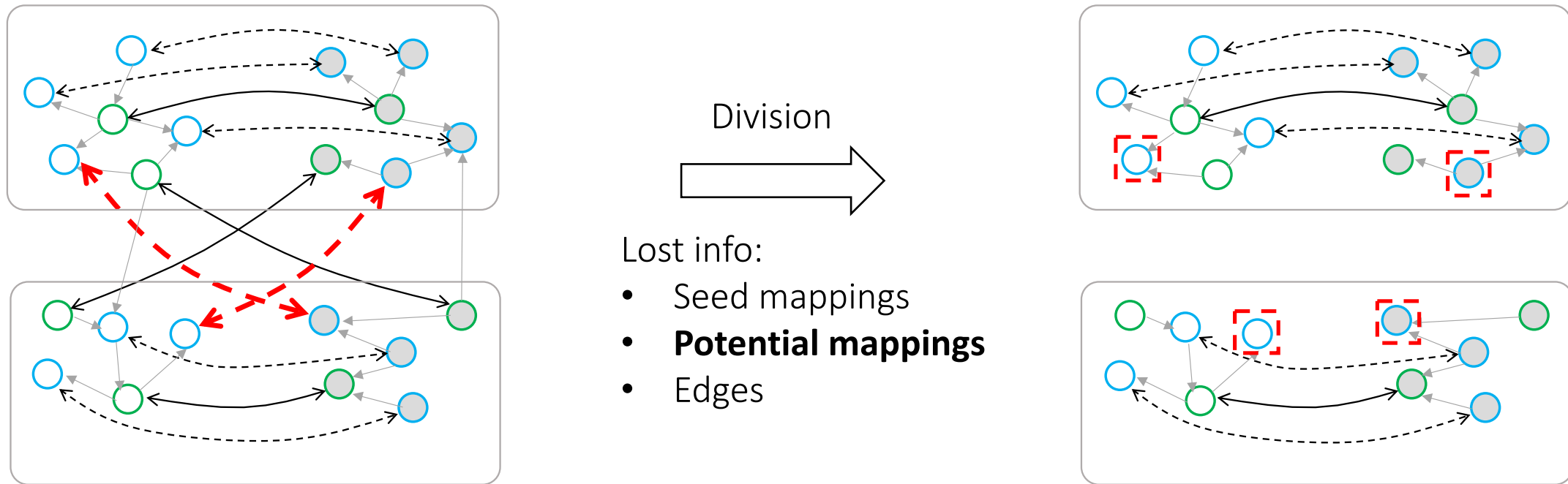
Lost info:

- **Seed mappings**
- **Potential mappings**
- **Edges**



○ Anchor entity ○ Unmatched entity

How to achieve high **coverage** of potential mappings?



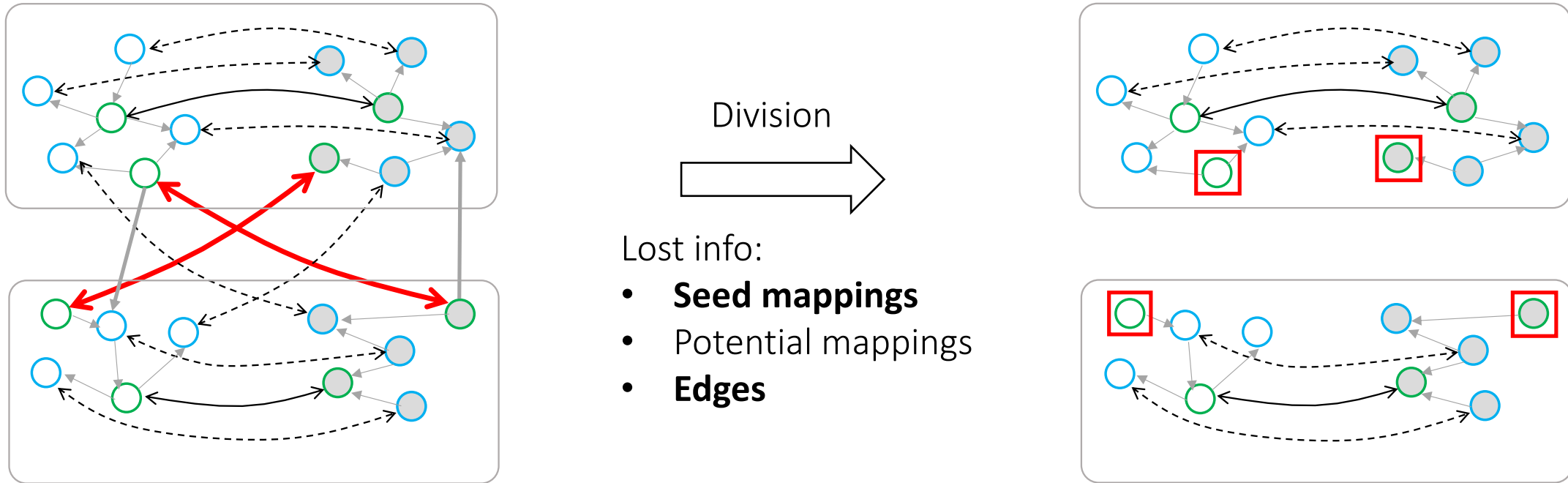
○ Anchor entity

○ Unmatched entity

□ Become unmatchable

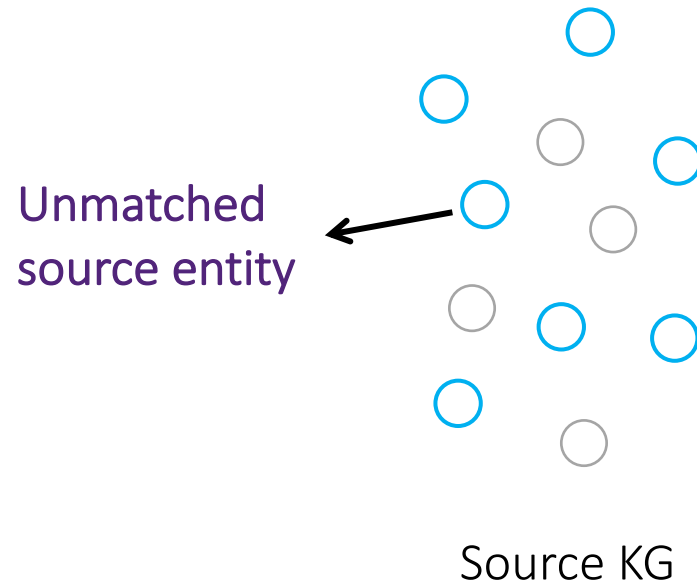
How to build informative **context graphs**?

- The two graphs that contain the unmatched entities and are fed into the EA model. They provide **evidence** for entity matching.

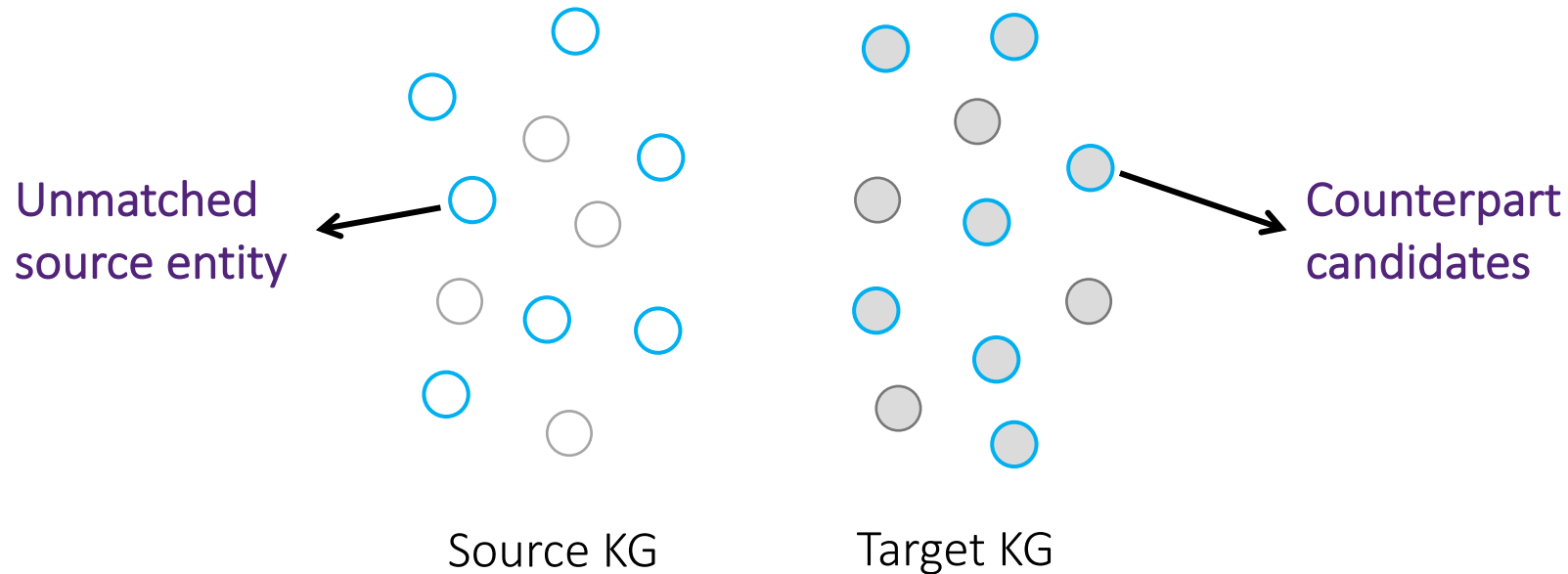


The DivEA Framework

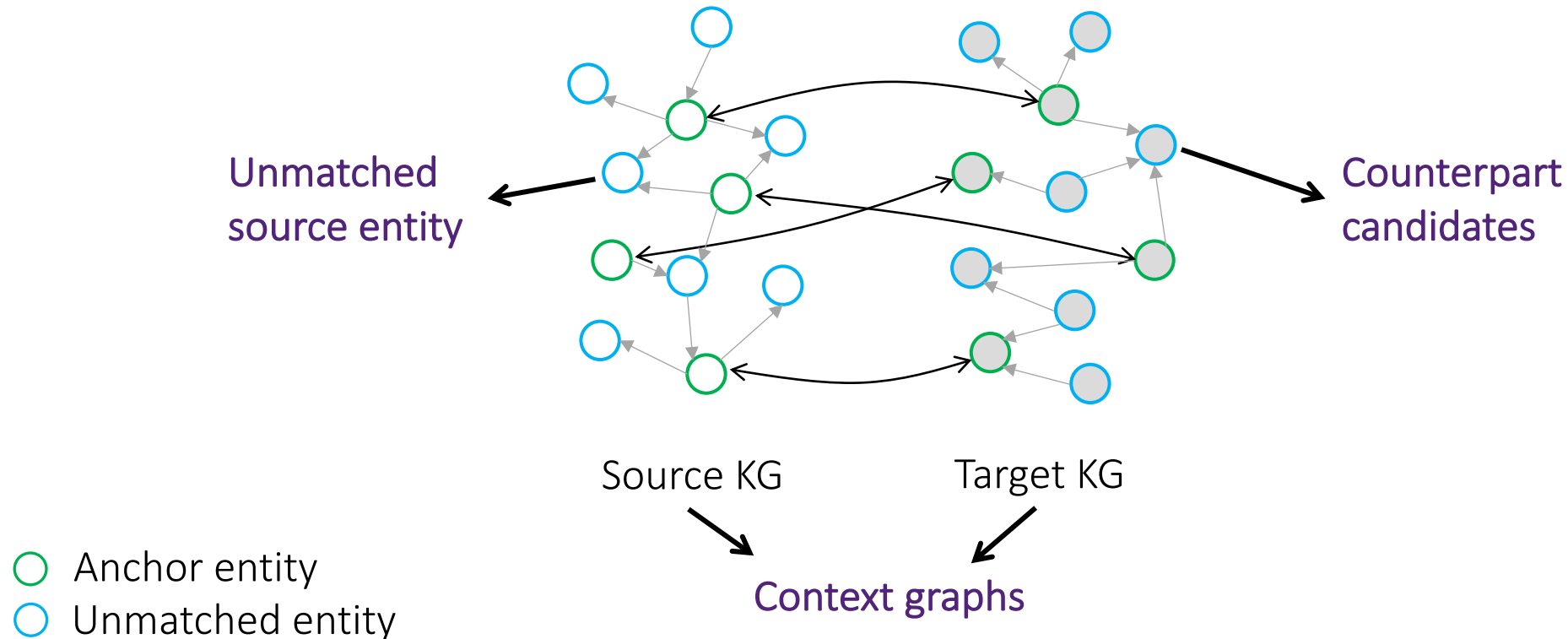
- **Unmatched source entities**



- **Unmatched source entities**
- **Counterpart candidates**, i.e. unmatched target entities

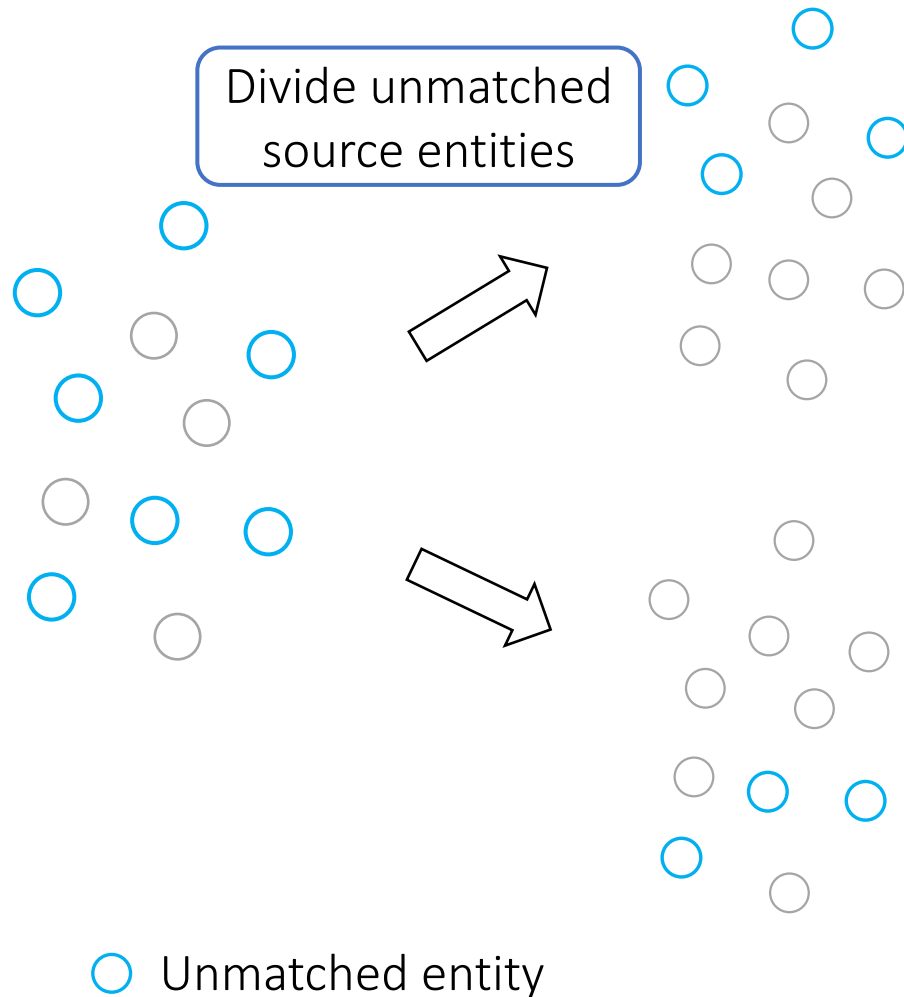


- **Unmatched source entities**
- **Counterpart candidates**, i.e. unmatched target entities
- **Context graphs**



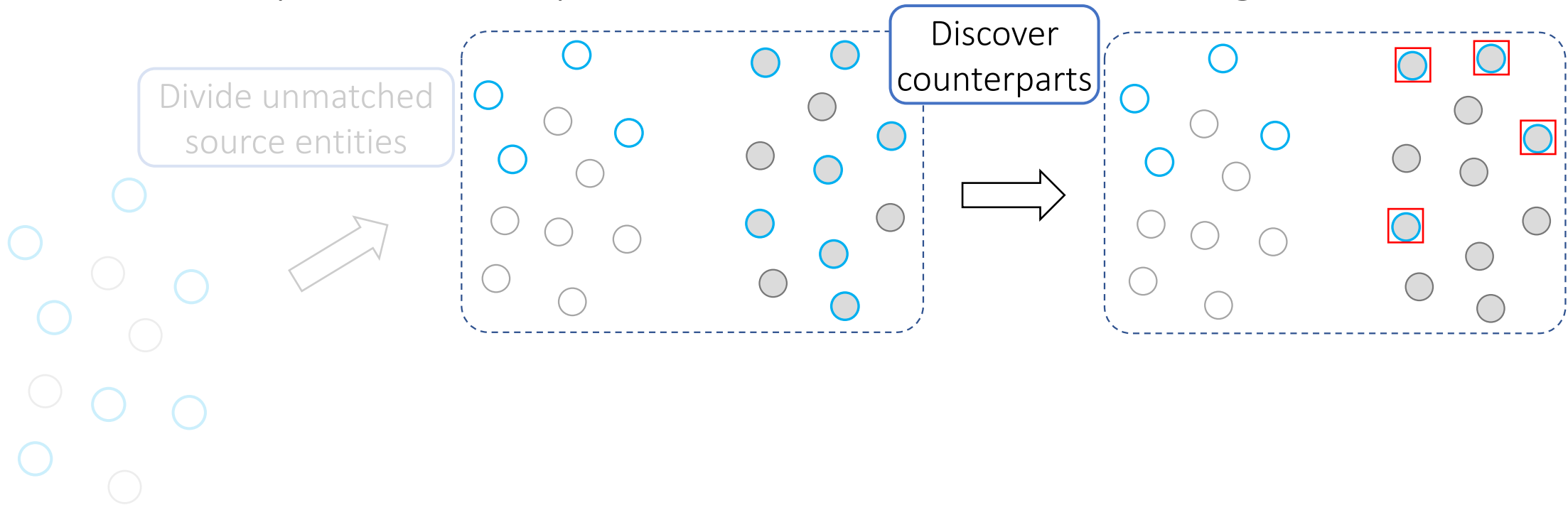
Overview: Divide Unmatched Source Entities

- Divide unmatched source entities



Overview: Counterpart Discovery

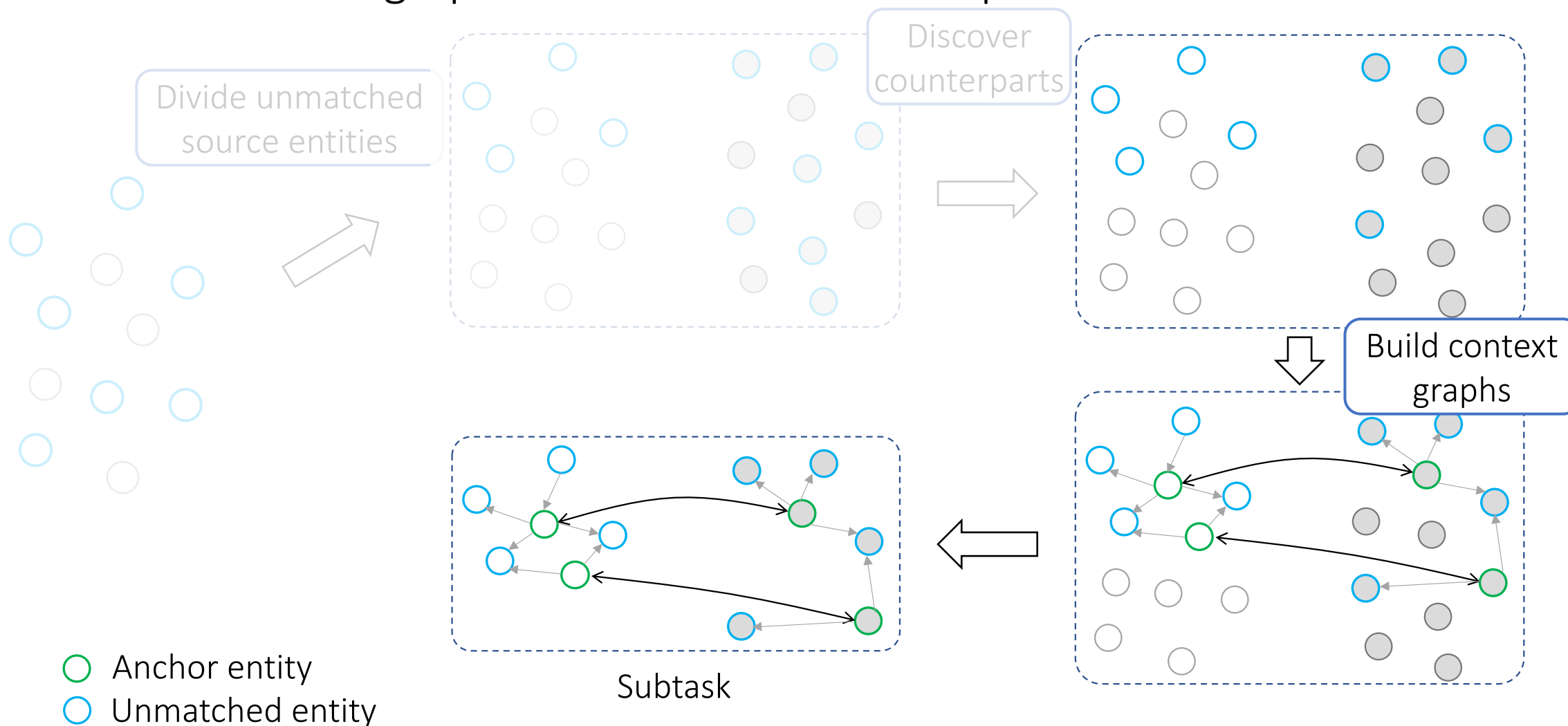
- Counterpart discovery: select a limited number of target entities



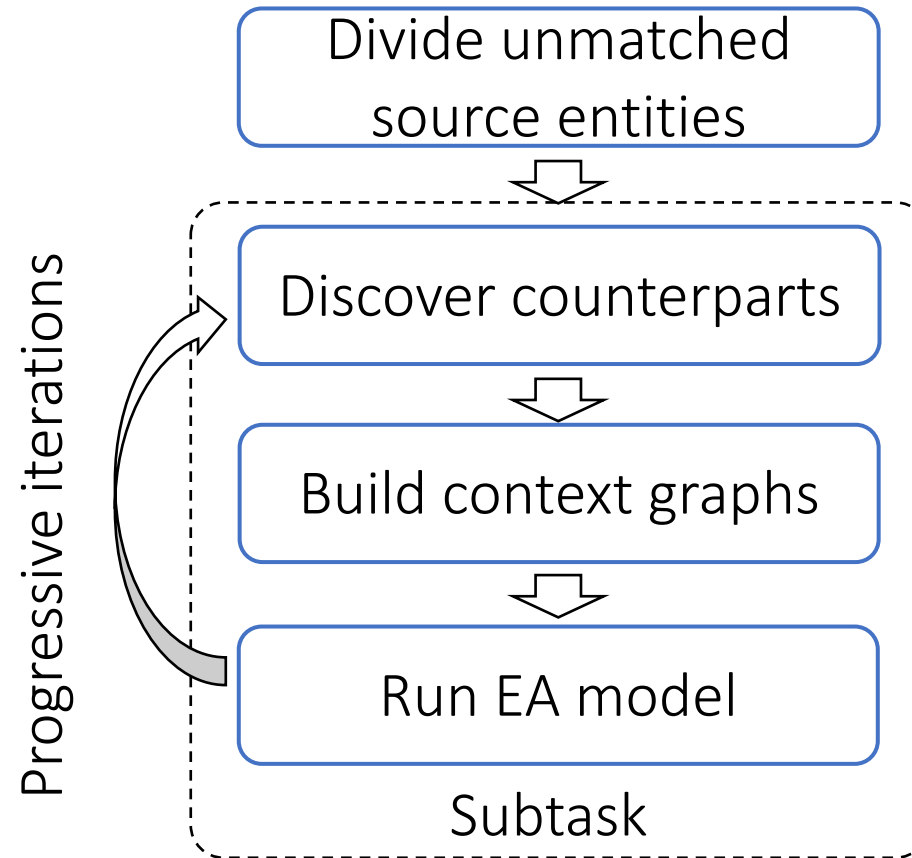
○ Unmatched entity □ Selected

Overview: Build Context Graphs

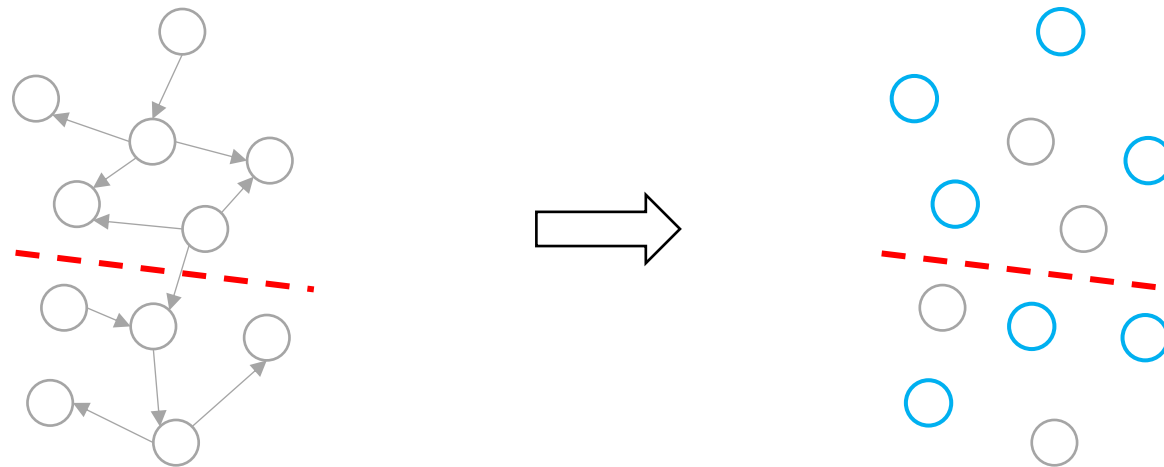
- Build context graphs: add more entities to provide evidence.



- Progressive process

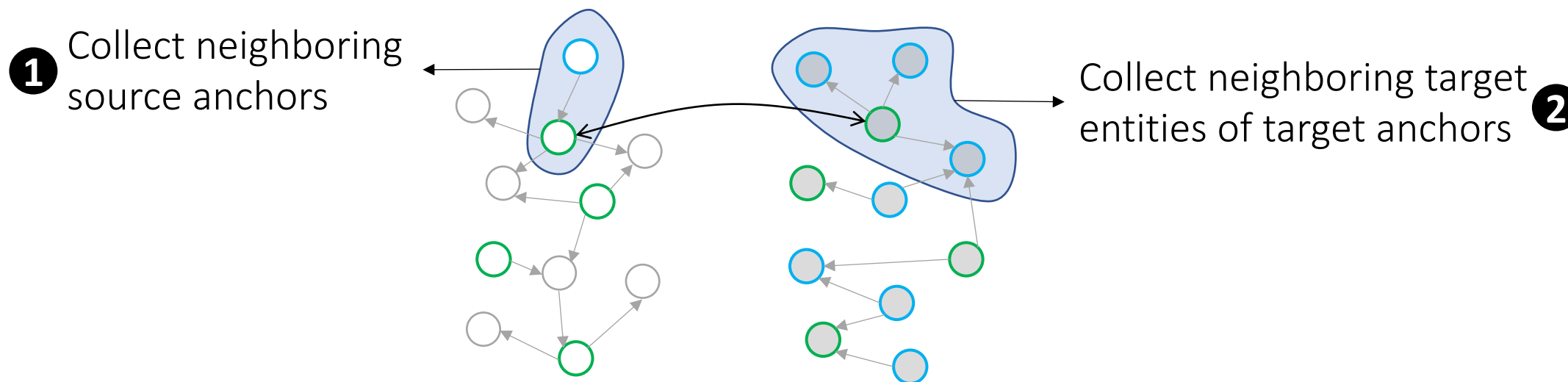


- Partition the source KG into **cohesive subgraphs** using *Metis*
 - The least cut-off of edges.
 - Balanced sizes.
- The unmatched source entities in each partition form one subset.



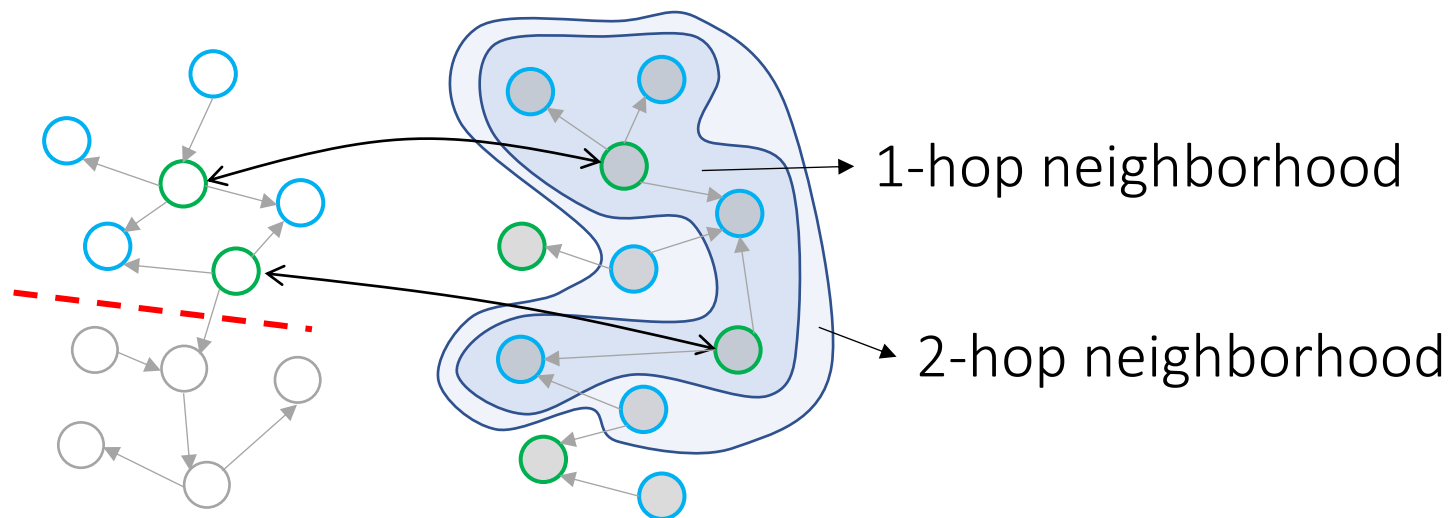
Given a certain source entity, how to identify its potential counterpart without using EA model?

- Principle of **locality**: If two entities are equivalent, the other entities semantically related to them might also be equivalent.



- Anchor entity
- Unmatched entity

1. Collect anchors in the same graph partition.
2. Locality-based weight $W^{loc}(e^t)$ according to the distance between target entity e^t and target anchors.

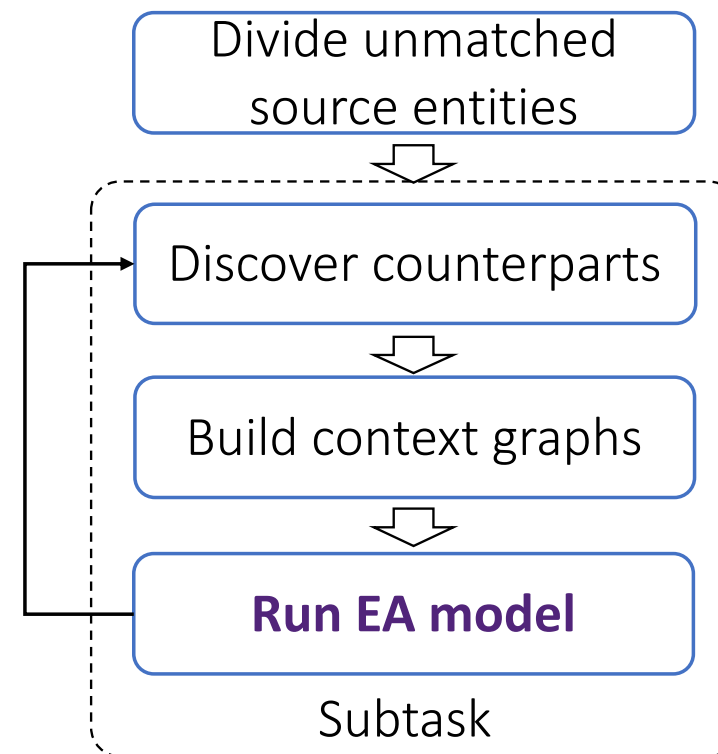


- Anchor entity
- Unmatched entity

If you have an EA model, how to use it for counterpart discovery?

- Enrich the seed mappings.
- Similarity-based signal $W^{sim}(e^t)$ indicating the likelihood that e^t is the counterpart of any source entity.

- Pseudo-mappings
- Similarity scores



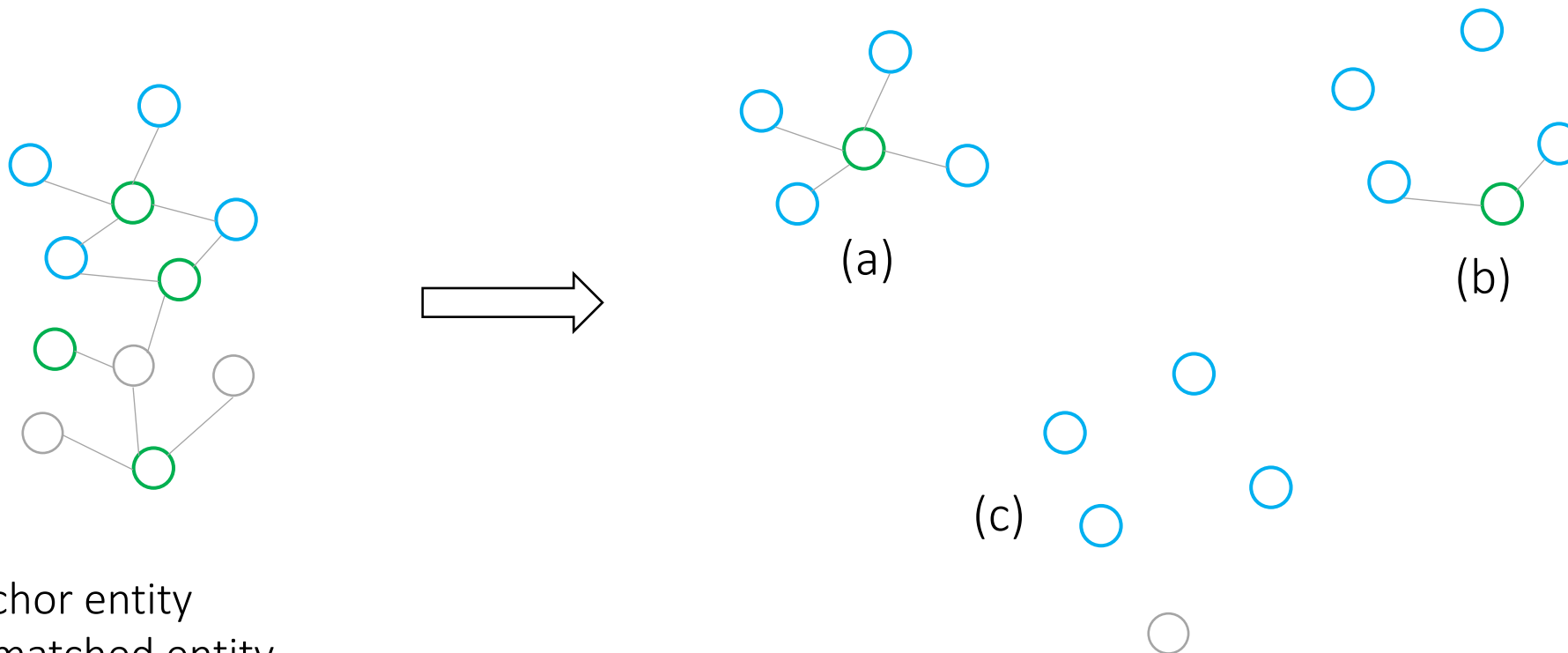
Choose target entities with the highest overall weights $W(e^t)$

- β is a hyper-parameter.

$$W(e^t) = W^{loc}(e^t) + \beta W^{sim}(e^t)$$

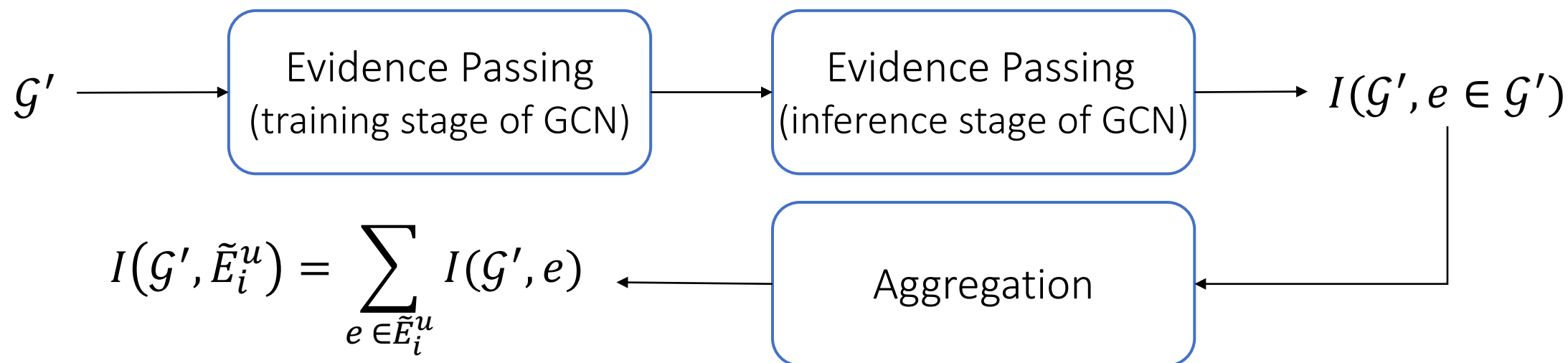
The context graph matters a lot for EA model.

- Example: build a context graph of size 5 for the unmatched entities.
- The unmatched entities can get different evidence.



How to quantify the informativeness of a context graph?

- **Evidence Passing** mechanism: to simulate how evidence spread around a graph in a GCN-based EA model.
 - The evidence is **scalar** instead of high-dimension vector.
 - The evidence origins from the anchors.
 - The evidence spreads in the **training and inference stages** of GCN.



- With the quantification method, we can search the most informative context graph within a single KG.
- For a subtask, we build the source context graph first, and then the target one.

Experiments

Comparison with 2 baselines: CPS, SBP.

Task division for 2 EA models: GCN-Align, RREA

Evaluated on 6 datasets: DBP15K: FR-EN, JA-EN, ZH-EN; DWY100K: DBP-WD, DBP-YG; FB-DBP (2M)

Metrics: Hit@1 (H@1), Hit@5 (H@5), Mean Reciprocal Rank (MRR)

- Overall performance

Method	EA model	FR-EN (15K)		
		H@1	H@5	MRR
CPS (sup)	GCN-Align	0.151	0.396	0.263
CPS (semi)		0.274	0.478	0.367
DivEA		0.396	0.642	0.504
CPS (sup)	RREA	0.419	0.631	0.514
CPS (semi)		0.516	0.682	0.590
DivEA		0.645	0.795	0.711

- Overall performance

Method	EA model	FR-EN (15K)			FB-DBP (2M)		
		H@1	H@5	MRR	H@1	H@5	MRR
CPS (sup)	GCN-Align	0.151	0.396	0.263	0.000	0.000	0.000
CPS (semi)		0.274	0.478	0.367	0.000	0.000	0.000
DivEA		0.396	0.642	0.504	0.051	0.106	0.08
CPS (sup)	RREA	0.419	0.631	0.514	0.043	0.080	0.062
CPS (semi)		0.516	0.682	0.590	0.056	0.089	0.073
DivEA		0.645	0.795	0.711	0.163	0.24	0.202

- Overall performance

Method	EA model	FR-EN (15K)			FB-DBP (2M)		
		H@1	H@5	MRR	H@1	H@5	MRR
SBP (sup)	GCN-Align	0.163	0.426	0.284	0.000	0.000	0.000
SBP (semi)		0.288	0.511	0.391	0.005	0.011	0.008
I-SBP		0.175	0.372	0.267	0.000	0.000	0.000
DivEA		0.402	0.678	0.525	0.071	0.15	0.112
SBP (sup)	RREA	0.475	0.721	0.583	0.070	0.139	0.106
SBP (semi)		0.575	0.762	0.659	0.095	0.159	0.128
I-SBP		0.508	0.730	0.608	0.120	0.233	0.172
DivEA		0.655	0.841	0.736	0.199	0.298	0.248

- Coverage of potential mappings
 - Metric: recall of potential mappings in the subtasks

	15K			100K		2M
Method	FR-EN	JA-EN	ZH-EN	DBP-WD	DBP-YG	FB-DBP
CPS	0.817	0.718	0.826	0.542	0.486	0.237
DivEA	0.881	0.892	0.880	0.830	0.893	0.507

- Coverage of potential mappings
 - Metric: recall of potential mappings in the subtasks

	15K			100K		2M
Method	FR-EN	JA-EN	ZH-EN	DBP-WD	DBP-YG	FB-DBP
CPS	0.817	0.718	0.826	0.542	0.486	0.237
DivEA	0.881	0.892	0.880	0.830	0.893	0.507

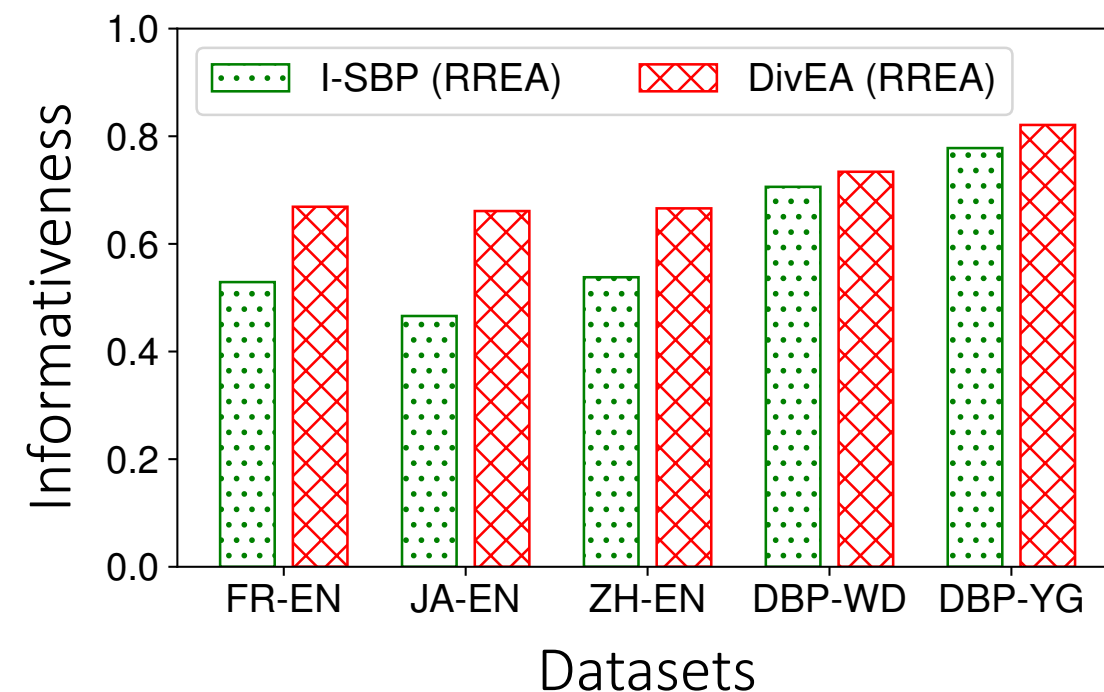
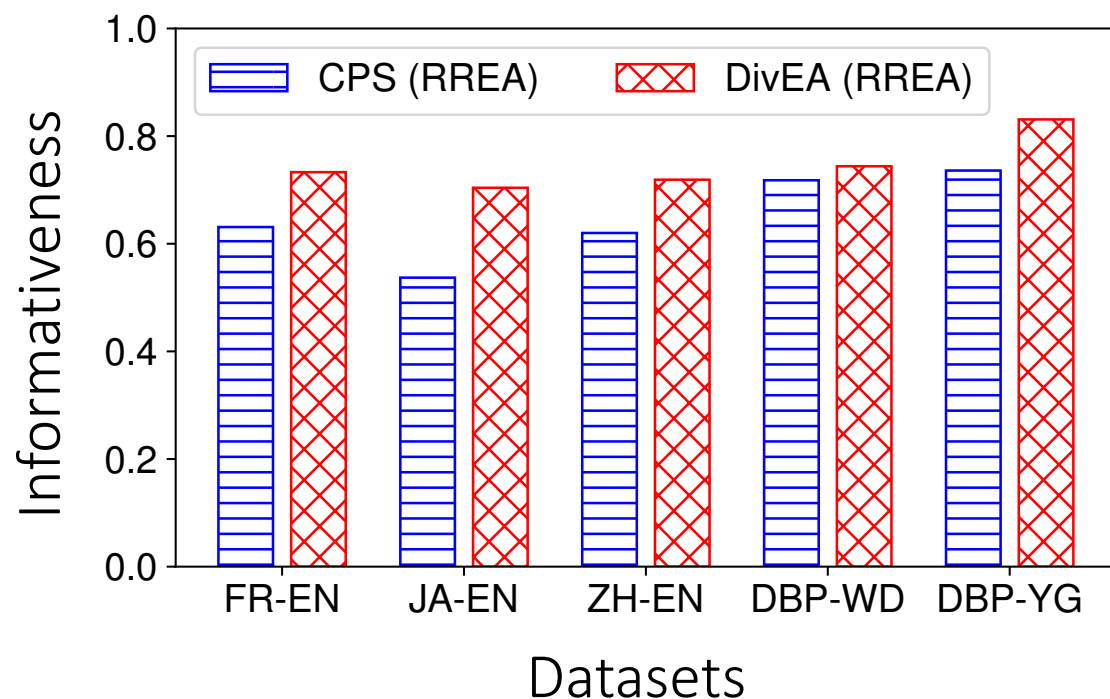
- Coverage of potential mappings
 - Metric: recall of potential mappings in the subtasks

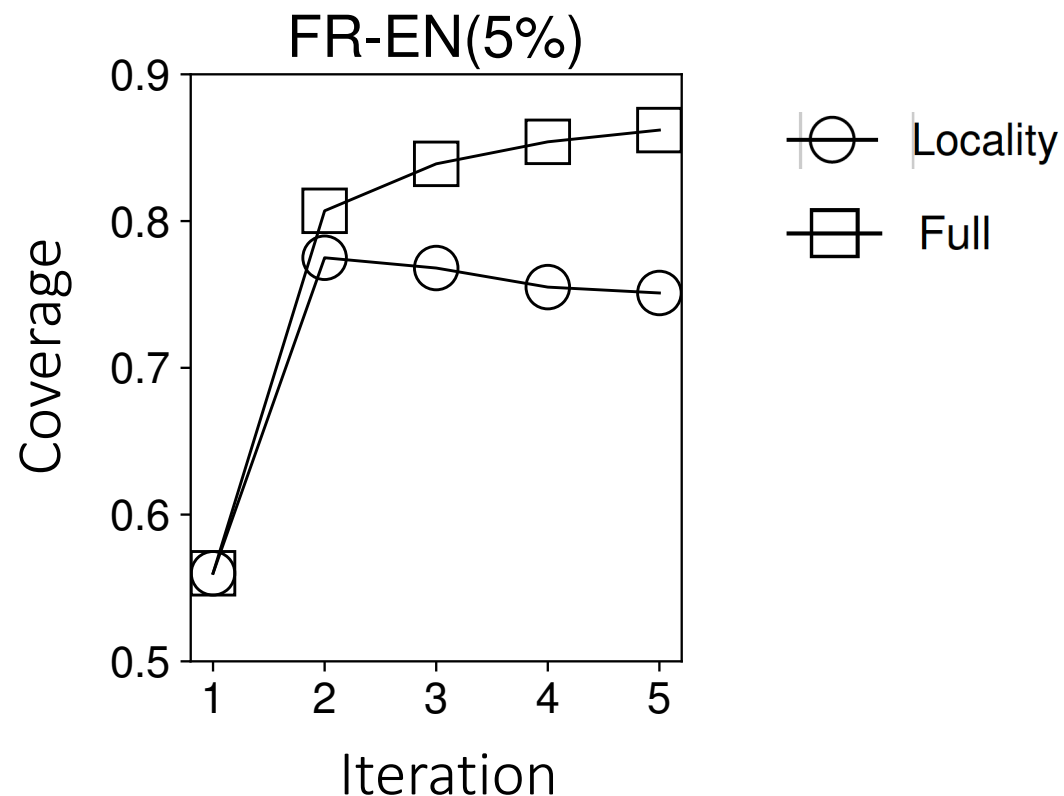
		15K			100K		2M	
		Method	FR-EN	JA-EN	ZH-EN	DBP-WD	DBP-YG	FB-DBP
		CPS	0.817	0.718	0.826	0.542	0.486	0.237
		DivEA	0.881	0.892	0.880	0.830	0.893	0.507
Larger subtask size	SBP	0.930	0.942	0.943	0.819	0.824	0.426	
	I-SBP	0.960	0.957	0.960	0.947	0.982	0.502	
	DivEA	0.978	0.979	0.970	0.954	0.994	0.684	

- Coverage of potential mappings
 - Metric: recall of potential mappings in the subtasks

		15K			100K		2M
Larger subtask size	Method	FR-EN	JA-EN	ZH-EN	DBP-WD	DBP-YG	FB-DBP
	CPS	0.817	0.718	0.826	0.542	0.486	0.237
	DivEA	0.881	0.892	0.880	0.830	0.893	0.507
	SBP	0.930	0.942	0.943	0.819	0.824	0.426
	I-SBP	0.960	0.957	0.960	0.947	0.982	0.502
	DivEA	0.978	0.979	0.970	0.954	0.994	0.684

- Informativeness of context graphs
 - Metric: percentage of found mappings over all mappings contained by the subtasks.





- Locality-based weight leads to decent performance
- EA model boosts it further

DivEA: a high-quality task division framework for large-scale EA.

- Dividing unmatched source entities + counterpart discovery + building context graphs.
- Progressive process.
- Building and running subtasks independently (for parallelization).

Code & data: <https://github.com/uqbingliu/DivEA>

Thank you for listening!

✉ bing.liu@uq.edu.au

🌐 <https://uqbingliu.github.io/>

🐦 @BingLiu1011

SIGIR