

Benchmarking Filtering Techniques for Entity Resolution

ICDE 2023, Research Session 8:
Information Integration and
Data Quality, April 5th 2023

George Papadakis, Marco Fisichella,
Franziska Schoger, George Mandilaras,
Nikolas Augsten, Wolfgang Nejdl



Introduction


- Entity resolution
 - identifying pairs of entity profiles that represent the same real-world object
- Filtering Techniques
 - don't check every possible pair but only the most promising (candidate pairs)
- Matching

- Example:
- 1 Movie in two databases



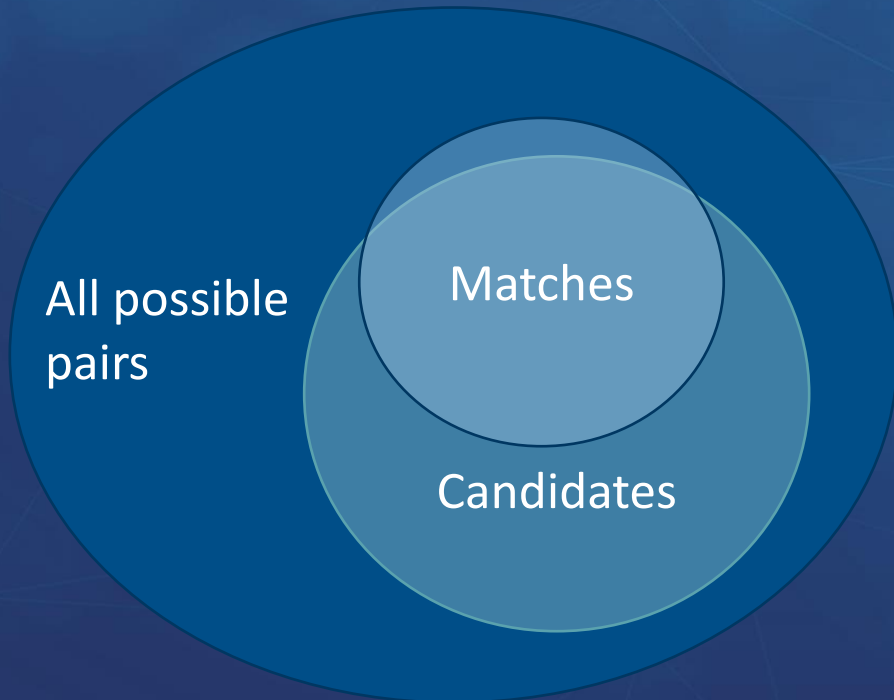
© New Line Cinema

id	title	starring	writer	editor
23304		Sean Duke, Hugo Weaving, Lawrence Makoare, Christopher Lee, Catherine Blanchett, Ian Holm...	Frances Walsh	

	id	title	actor name	director name	year	genre
	3831	Lord of the Rings: The Fellowship of the Ring, The (2001)	Lee, Christopher; Jackson, Peter; McKellen, Ian; Rhys-Davies, John; ...	Jackson, Peter	2001	Action, Fantasy Adventure

Preliminaries

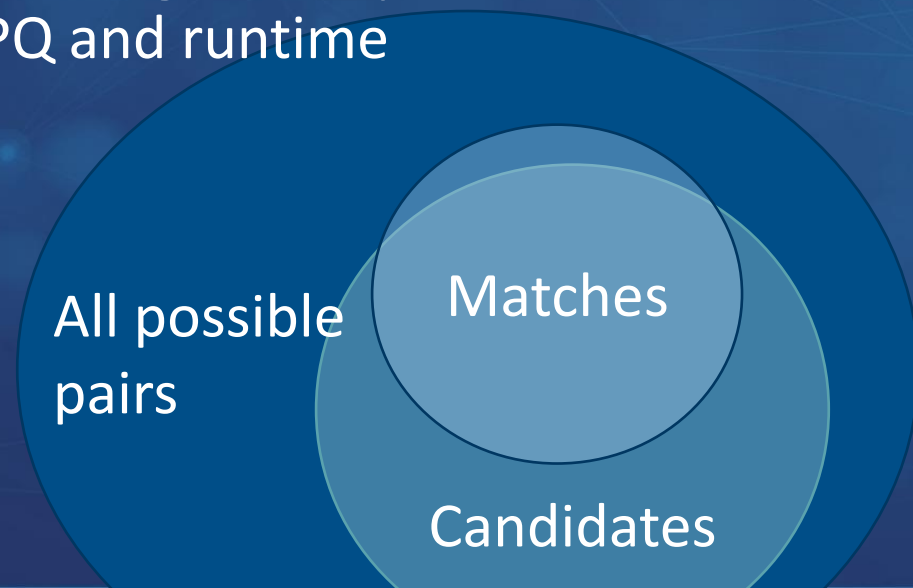
- Two main tasks of ER
 - Clean-Clean ER
 - Dirty ER



- Performance Measures
 - Pair Completeness (PC)
 - $\frac{|Candidates \cap Matches|}{|Matches|}$
 - Pairs Quality (PQ)
 - $\frac{|Candidates \cap Matches|}{|Candidates|}$
 - Runtime

Main Idea

- Filtering techniques for textual entity profiles:
 - Blocking Workflows
 - NN Methods
- Main idea: compare different filtering techniques in terms of PC, PQ and runtime

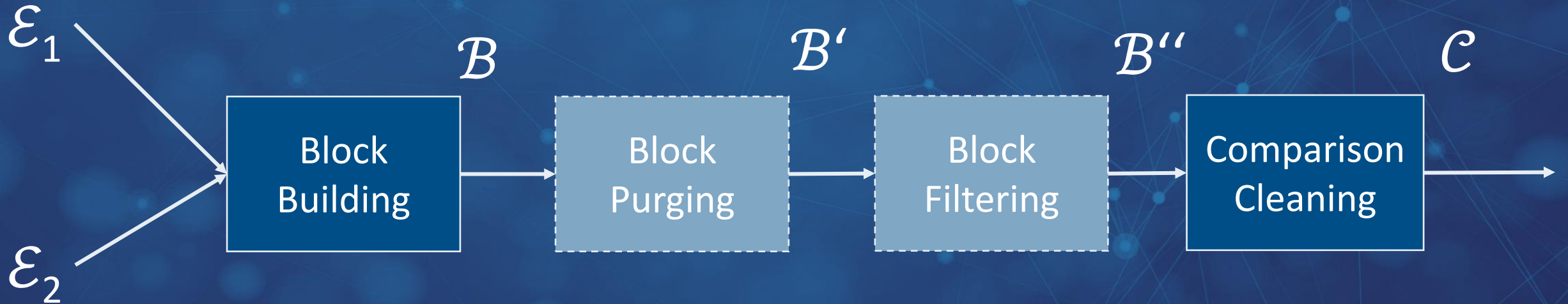


Problem (Configuration Optimization):

- Given are two sets of entity profiles, a filter method and a threshold on PC (90%)
- Finetune the parameters of the filtering methods such that the resulting PQ is maximized while the PC is above the threshold

Filtering Methods

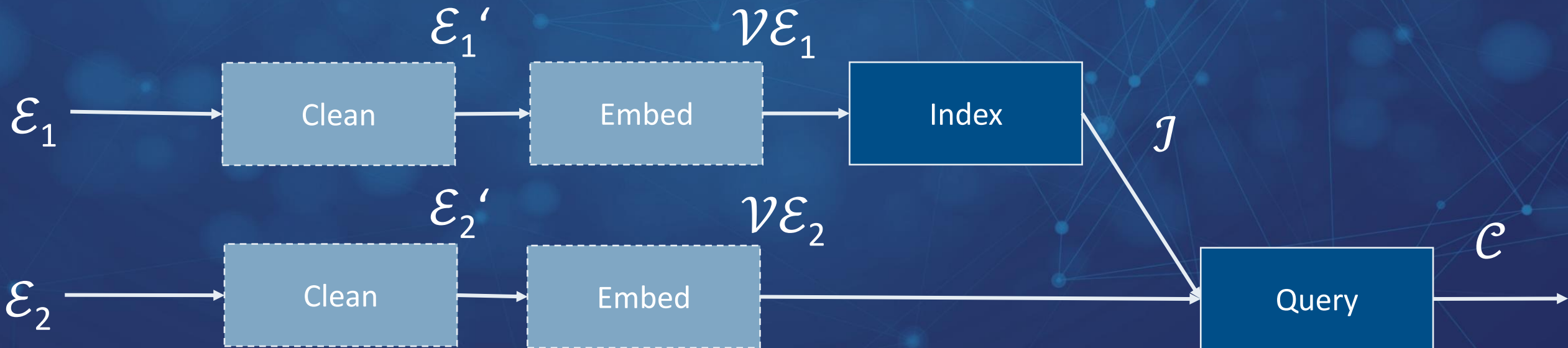
Blocking Methods



Filtering Methods

Nearest Neighbor Methods

- Sparse Vector Based NN Methods:
 - Similarity of token sets (Jaccard, Cosine, Dice)
- Dense Vector Based NN Methods:
 - Similarity of Vector Embeddings



Datasets

	D _{c1}	D _{c2}	D _{c3}	D _{c4}	D _{c5}	D _{c6}	D _{c7}	D _{c8}	D _{c9}	D _{c10}
$\varepsilon_1 / \varepsilon_2$	Rest. 1 / Rest. 2	Abt / Buy	Amazon / GB	DBLP / ACM	IMDb / TMDb	IMDb / TVDB	TMDb / TVDB	Walmart /Amazon	DBLP / GS	IMDb / DBpedia
# entities	339 / 2,256	1,076 / 1,076	1,354 / 3,039	2,616 / 2,294	5,118 / 6,056	5,118 / 7,810	6,056 / 7,810	2,554 / 22,074	2,516 / 61,353	27,615 / 23,182
# duplicates	89	1,076	1,104	2,224	1,968	1,072	1,095	853	2,308	22,863
Cartesian Product	$7.7 \cdot 10^5$	$1.2 \cdot 10^6$	$4.1 \cdot 10^6$	$6.0 \cdot 10^6$	$3.1 \cdot 10^7$	$4.0 \cdot 10^7$	$4.7 \cdot 10^7$	$5.6 \cdot 10^7$	$1.5 \cdot 10^8$	$6.4 \cdot 10^8$
Best Attribute	Name	Name	Title	Title	Title	Name	Name	Title	Title	Title

Configuration Space

Method		Number of Configurations
Blocking Methods	Standard Blocking	3,440
	Q-Grams Blocking	17,200
	Extended Q-Grams Blocking	68,800
	(Ex.) Suffix Arrays Blocking	21,285
Sparse NN Methods	ϵ -Join	6,000
	kNN-Join	12,000
Dense NN Methods	MH-LSH	168
	HP-LSH	400
	CP-LSH	2,000
	FAISS	2,720
	SCANN	10,880
	DeepBlocker	2,720

Taxonomy

Scope		Blocking	Sparse NN	Dense NN
Syntactic Representations	Schema-based	✓	✓	✓
	Schema-agnostic	✓	✓	✓
Semantic Representations	Schema-based	-	-	✓
	Schema-agnostic	-	-	✓

Taxonomy

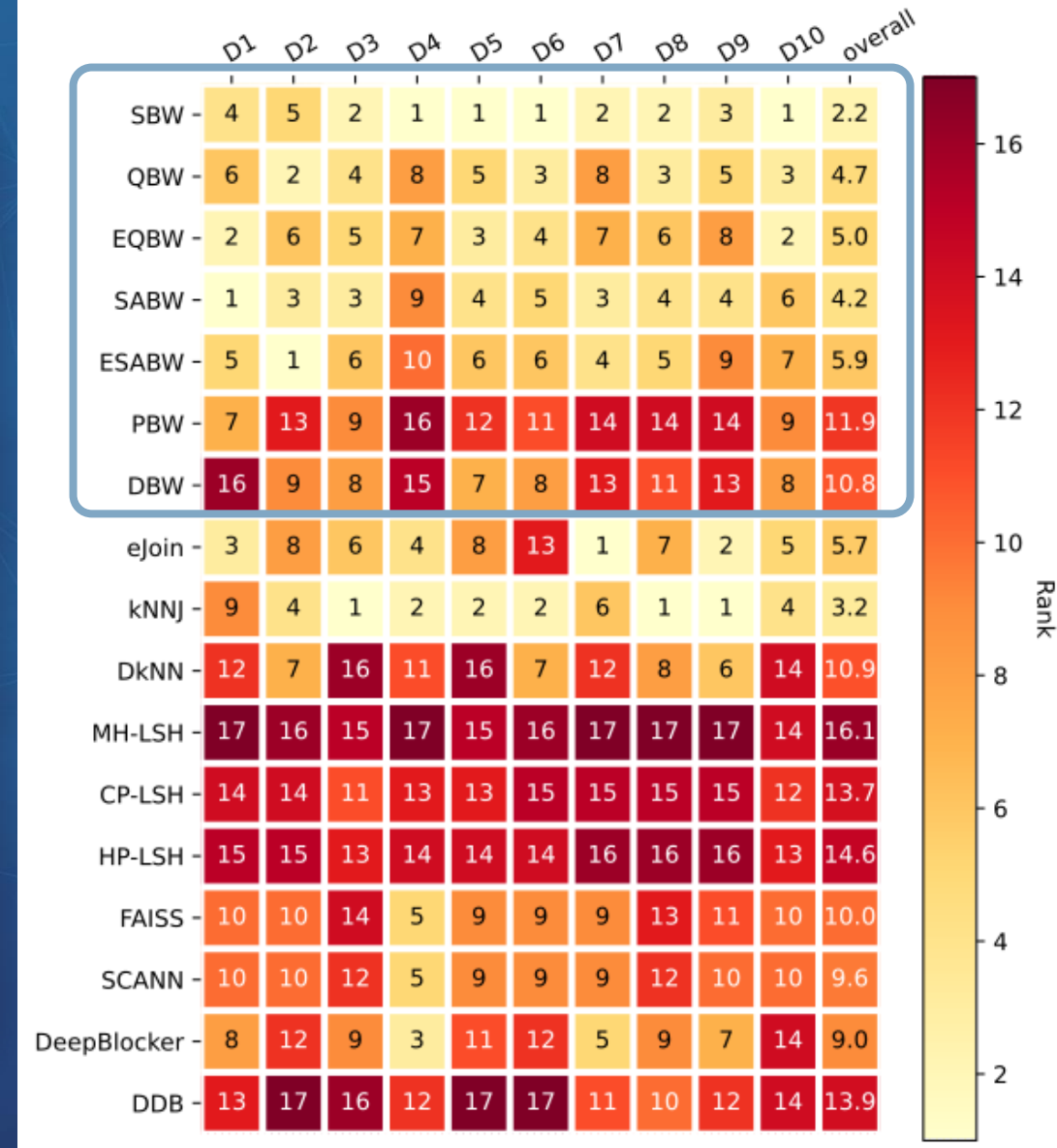
Internal functionality

- Blocking:
 - Lazy or proactive
- NN-Methods:
 - Type of operation
 - Type of threshold

NN Methods	Similarity Threshold	Cardinality Threshold
Deterministic Operation	ϵ -Join	kNN-Join, FAISS, SCANN
Stochastic Operation	MH-, HP-, CP-LSH	DeepBlocker

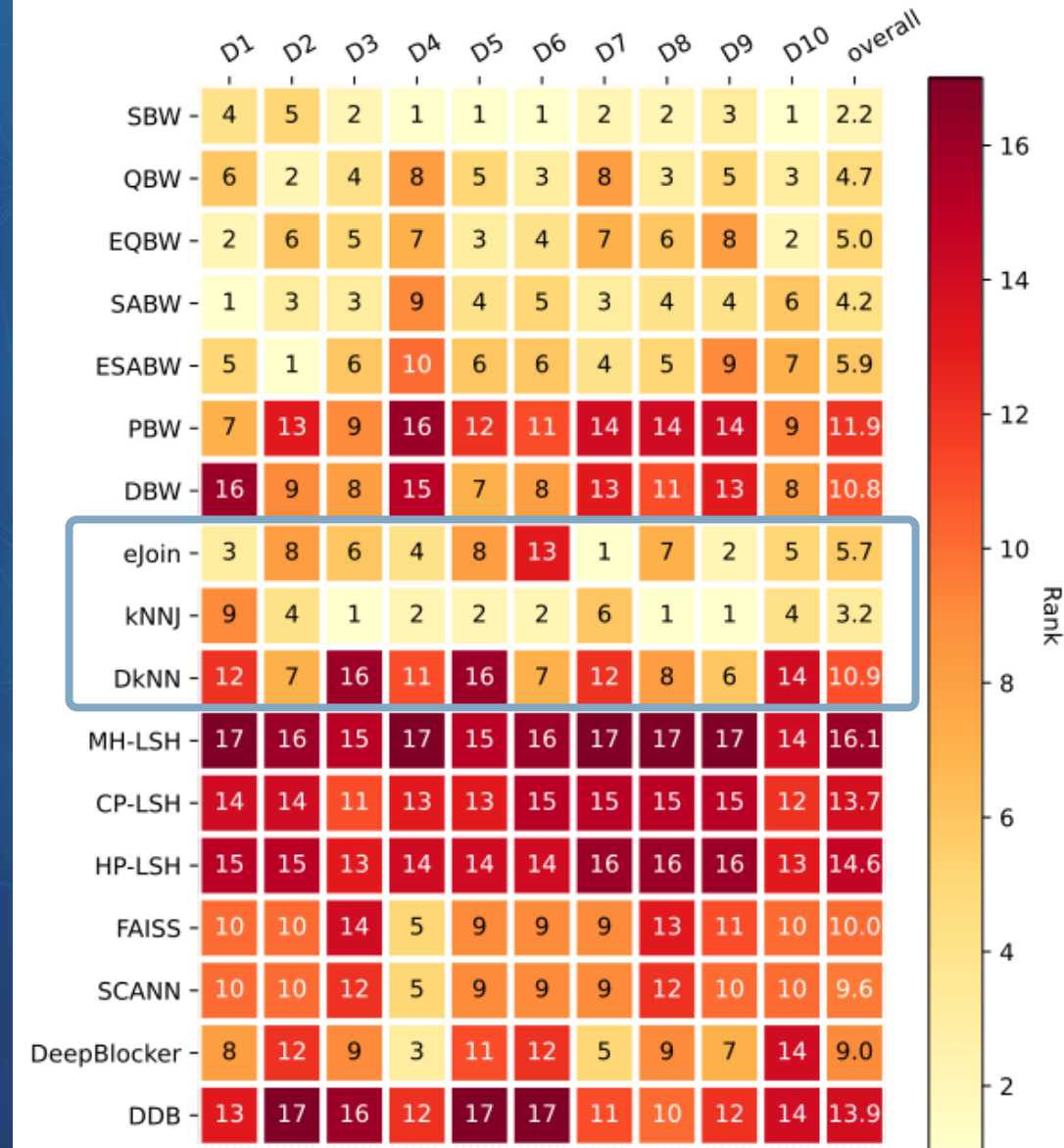
Schema Agnostic Results

- Blocking:
 - Attribute value tokens yield better results than substrings of tokens



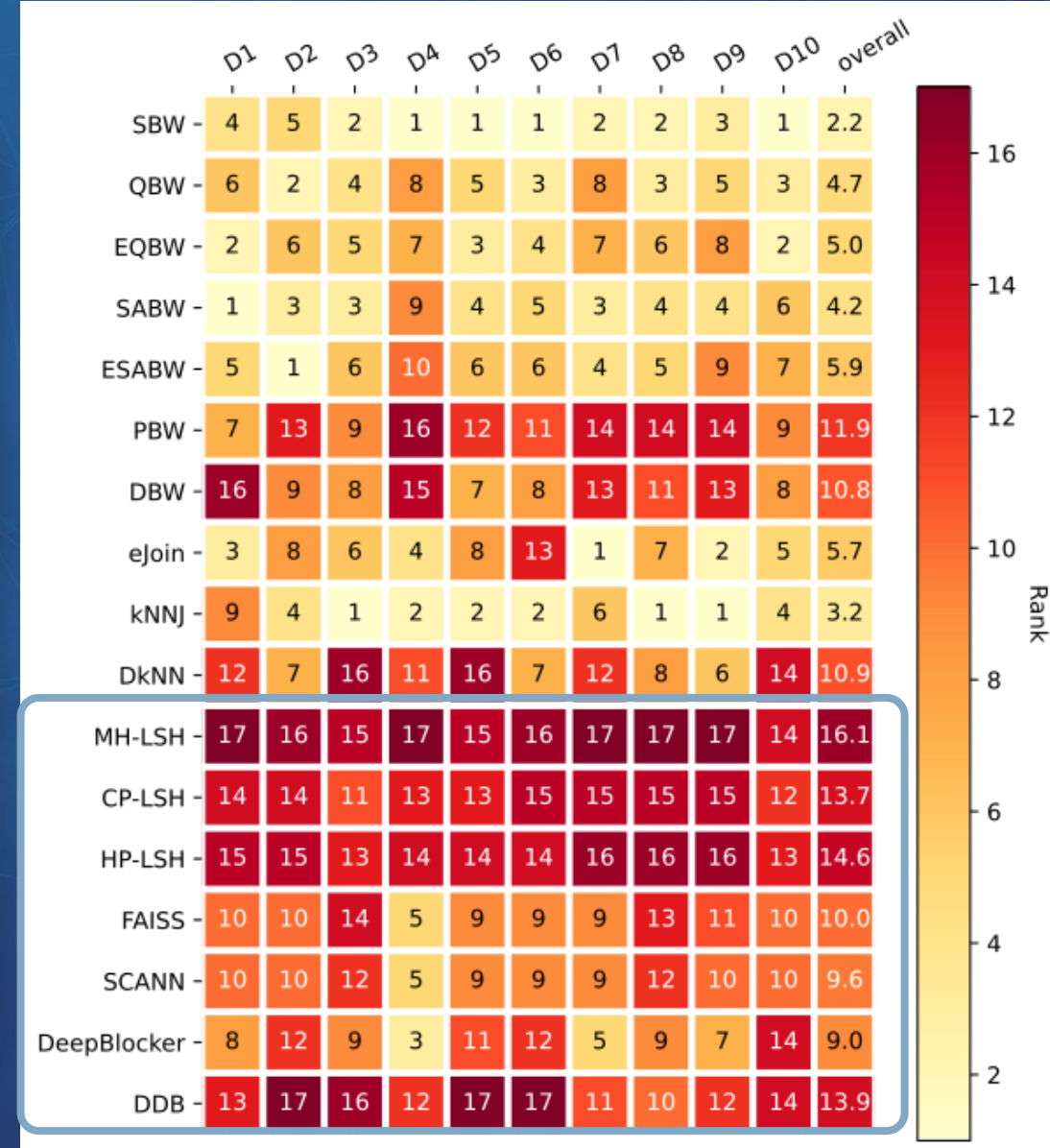
Schema Agnostic Results

- Sparse NN:
 - Cardinality thresholds are more effective than similarity thresholds



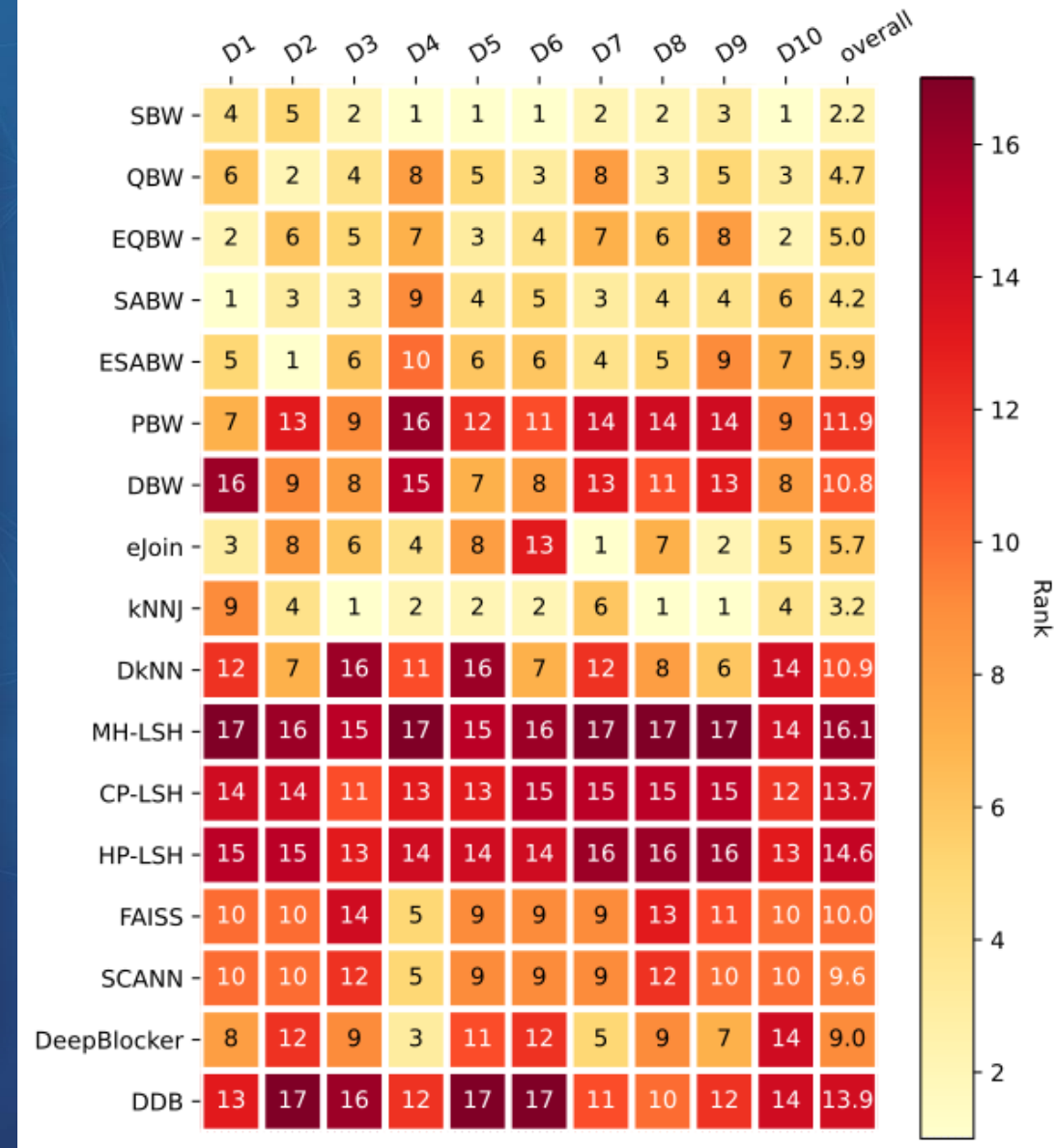
Schema Agnostic Results

- Dense NN:
 - Similarity-based methods produce large numbers of candidate pairs to achieve high PC
 - Learning based tuple embedding module raises the PQ, but does not scale



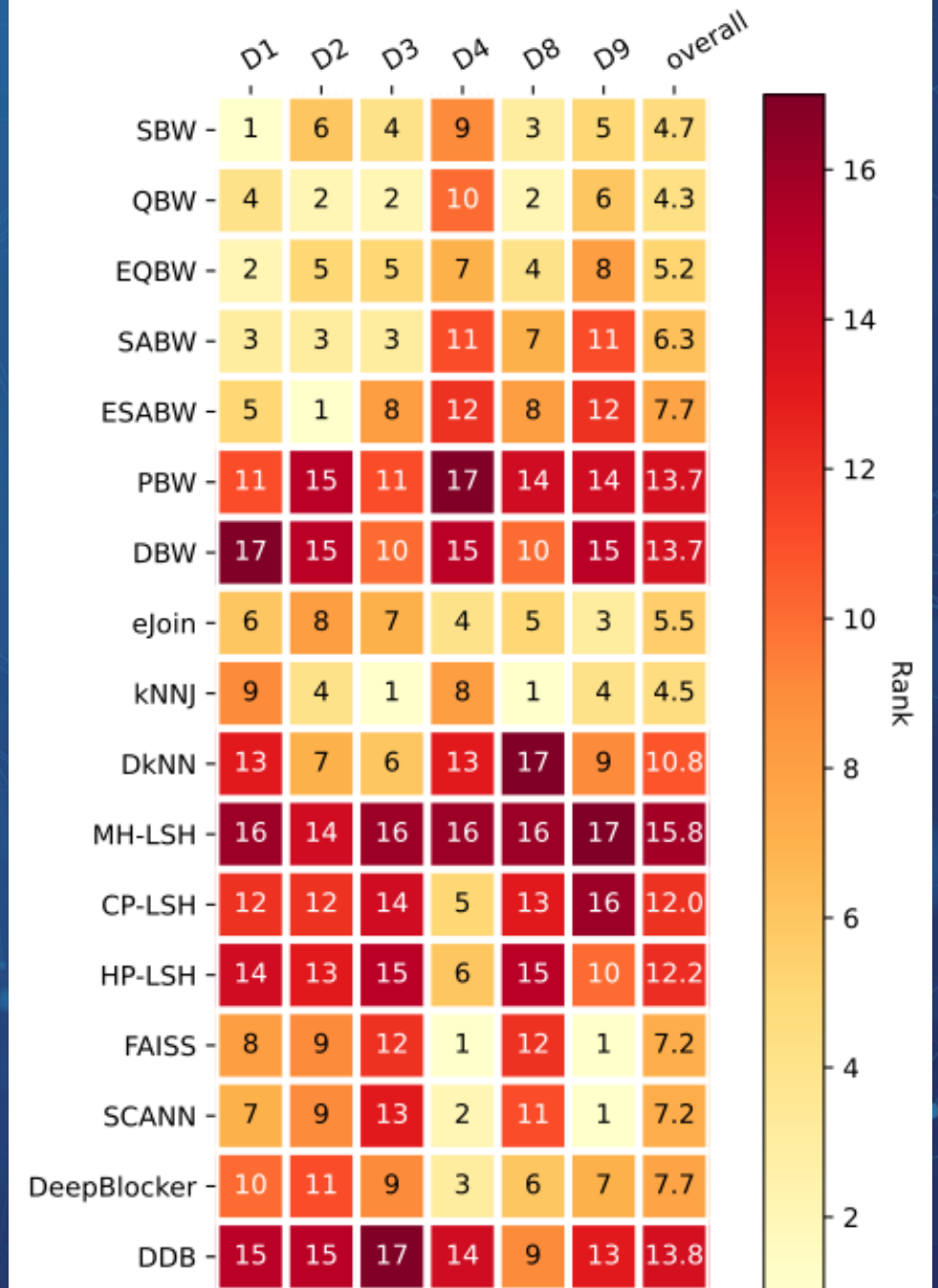
Schema Agnostic Results

- Standard Blocking Workflow (SBW) best
- kNN-Join more robust and easier to configure and apply



Schema Based Results

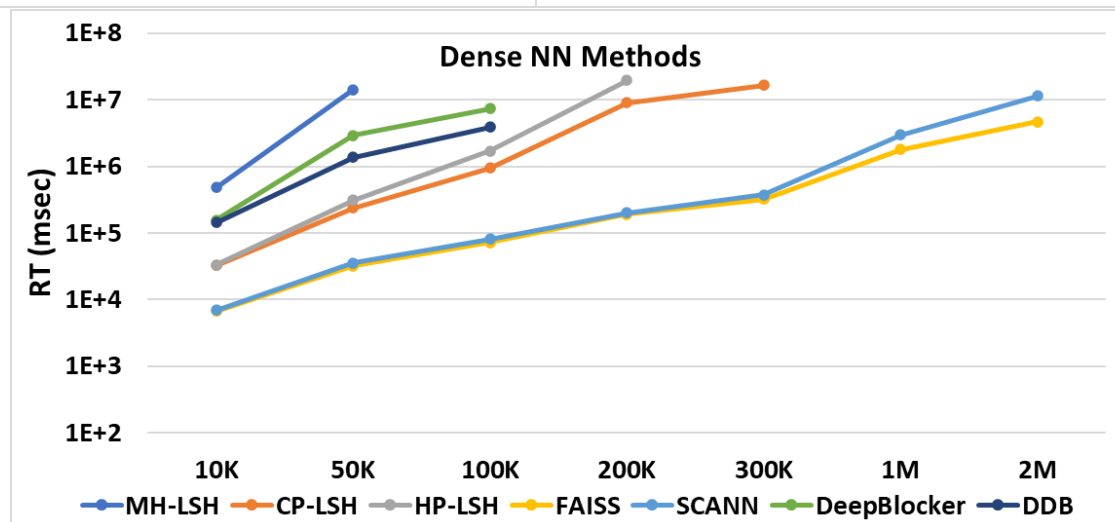
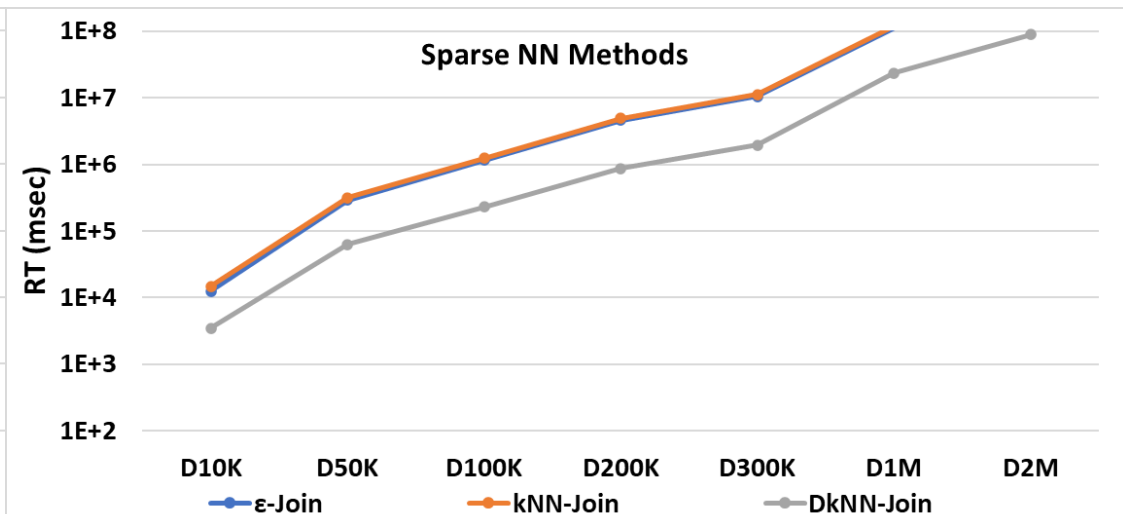
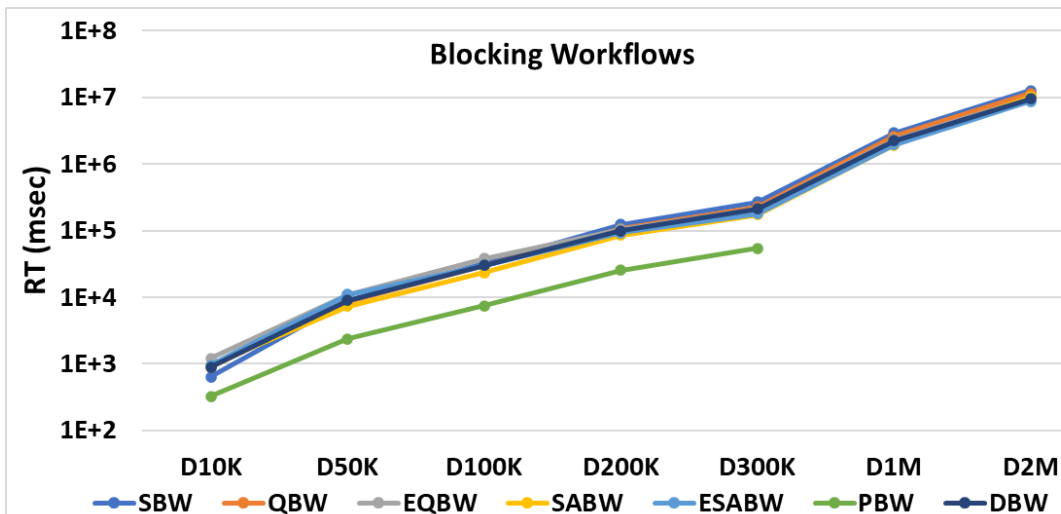
- Overall similar trends as for schema-agnostic setting
 - Blocking:
 - Substrings of tokens (QBW) slightly better than SBW
- QBW and kNNJ best and most robust



Scalability

- Run on 7 synthetic, dirty ER datasets
 - Ranging from 10,000 to 2M entries
- Techniques were finetuned on the smallest dataset
- The same settings were used for all datasets

Scalability



Conclusions

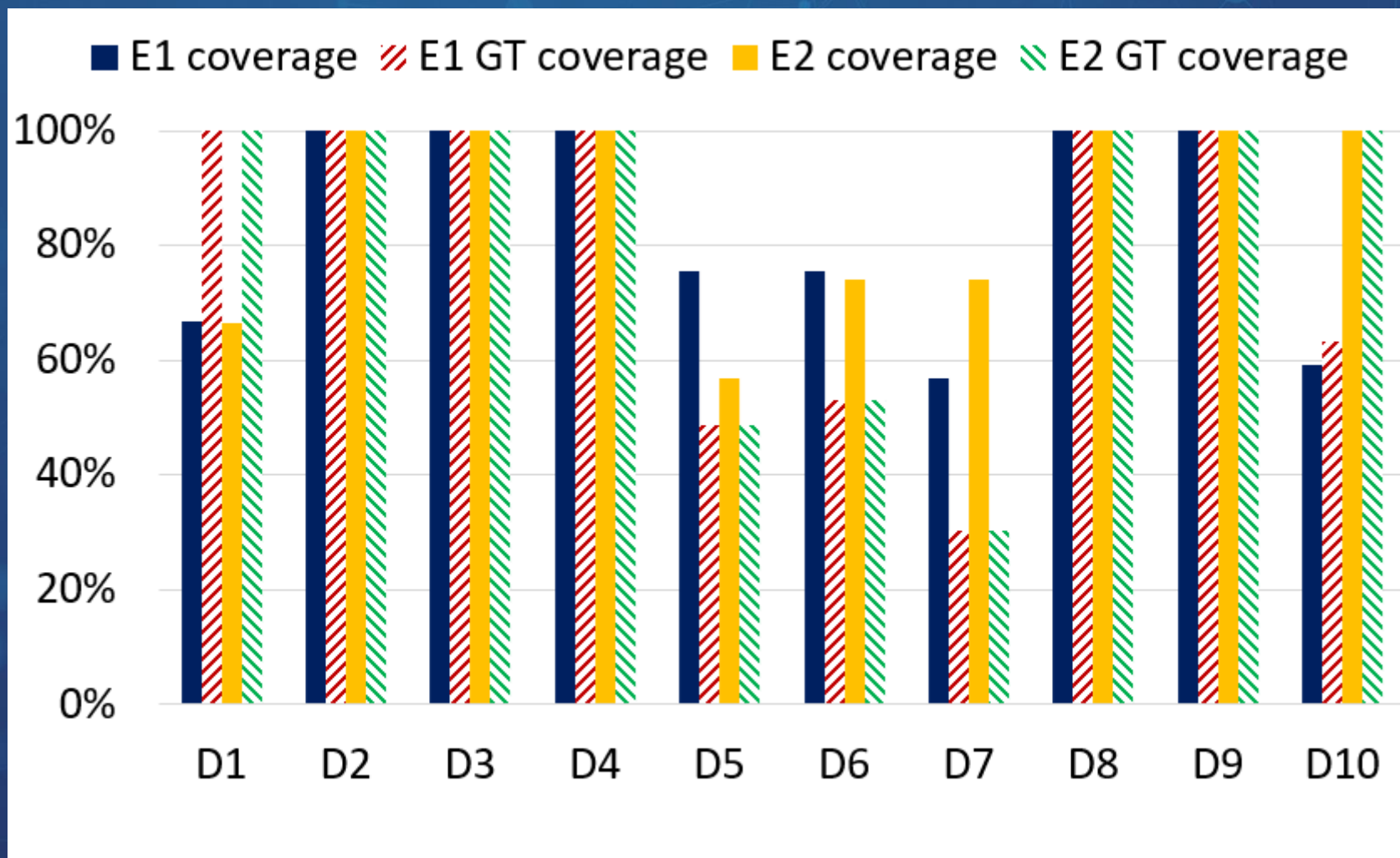
- We compared 13 methods, all fine-tuned to our problem and 4 baseline methods
- PQ of all methods is highly correlated
 - performance heavily depends on dataset characteristics
- Parameter finetuning significantly increases the performance
- Schema agnostic settings are preferable
- Cardinality thresholds are preferable
- Syntactic representations are preferable

Thank you

Questions?



Coverage



Scalability

