

<meta>

<meta>

Relational Header Discovery using Similarity Search in a Table Corpus

Hazar Harmouch, Thorsten Papenbrock, Felix Naumann

ICDE 2021

Virtual event

Chania, Crete, Greece

19-22 April 2021

Information System Group

Hasso Plattner Institut

University of Potsdam

Germany

Motivation

Missing headers



VLDB 2021: Conference Officers

Associate Editors

Alon Halevy Facebook, United States

Anastasia Ailamaki École Polytechnique Fédérale de Lausanne (EPFL), Switzerland

Angela Bonifati Lyon 1 University, France

Arun Kumar UC San Diego, United States

Ashraf Aboulnaga Qatar Computing Research Institute, HBKU, Qatar

Eugene Wu Columbia University, United States

Floris Geerts University of Antwerp, Belgium

Graham Cormode University of Warwick, United Kingdom

Jeffrey Xu Yu The Chinese University of Hong Kong, China

Jiannan Wang Simon Fraser University, Canada

Jingren Zhou Alibaba, China

Jorge Arnulfo Quiane Ruiz TU Berlin, Germany

Juliana Freire New York University, United States

Jun Yang Duke University, United States

?	?	?
Alon Halevy	Facebook	United States
Anastasia Ailamaki	École Polytechnique Fédérale de Lausanne (EPFL)	Switzerland
Angela Bonifati	Lyon 1 University	France
...

Motivation

Missing headers \approx un-interpretable headers



noe_pop_age_sex_2012_201...

Sort fields Data source order

Abc noe_pop_age_sex_...	Abc noe_pop_age_sex_...	Abc noe_pop_age_sex_...	# noe_pop_age_sex_2...	Abc noe_pop_age_sex_2012_...	Abc noe_pop_age_sex_2012...	# noe_pop_age_sex_...	# noe_pop_age_sex_...
Nuts1	Nuts2	Nuts3	Lau2 Code	Lau2 Name	Age Group	Pop Total	Pop Male
AT1	AT12	AT124	30101	Krems an der Donau	0_4	1,063	528
AT1	AT12	AT124					
AT1	AT12	AT124					
AT1	AT12	AT124					
AT1	AT12	AT124					
AT1	AT12	AT124					
AT1	AT12	AT124	30101	Krems an der Donau	30_34	1,586	855

Attributbeschreibung ?

NUTS1: Land NUTS2: Bundesland NUTS3: Gruppen von Bezirken
LAU2_CODE: Gemeindekennzahl LAU2_NAME: Gemeindename
AGE_GROUP: Altersgruppe (5-jährig) POP_TOTAL: Bevölkerung
POP_MALE: männliche Bevölkerung POP_FEMALE: weibliche
Bevölkerung YEAR: Referenzjahr

Motivation

How big is the problem?

Header row is missing in

- ❑ majority of tables on the web [Balakrishnan15].
- ❑ ~29% of relational web tables extracted by [Cafarella08].

- ❑ ~28% of the CSV tables from the Austrian and European Open Data portals [Neumaier16].
- ❑ ~11k documents in a data corpus from 232 Open Data portals [Mitlöhner16].



Web tables



**Relational Header
Discovery**

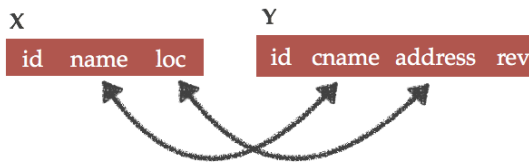
Hazar Harmouch
ICDE 2021

Motivation

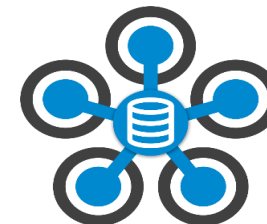
Where are column headers useful?



Querying data



Schema matching



Data integration



Knowledge base
augmentation



Table search
engines

Fuel Type	Borough	Sector	KWh
Electricity	Barnett	Domestic	62688
Gas	Barnett	Domestic	206438
Railway	City of London	Transport	2730044
Diesel	City of London	Domestic	430078
Oil	City of London	Domestic	430078

Borough	Population	Unemp	FUnemp
Barnett	38900	Low	20
Camden	40000	Low	14
City of London	888000	Medium	20

Query Table Candidate Table

Table join

The problem

Schema discovery problem

family	motto	location	symbol
stark	winter is coming	winterfell	grey wolf
targaryen	fire and blood	dragonstone	red dragon
lannister	hear me roar!	casterly rock	golden lion



?
stark
targaryen
lannister

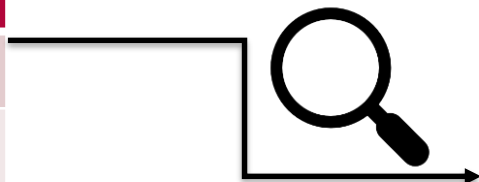


Table corpus

Borough	Population	Unemp	F.Unem
Barnett	38900	Low	20
Camden	40000	Low	14
City of London	888000	Medium	20

Relational Header Discovery

Hazar Harmouch
ICDE 2021

The problem

Schema discovery problem – Schema coherence

athlete

color

...

medal

...

athlete	medal
Eddie Eagan	gold
Jacob Tullin Thams	silver
Christa Luding-Rothenburger	silver
Clara Hughes	bronze



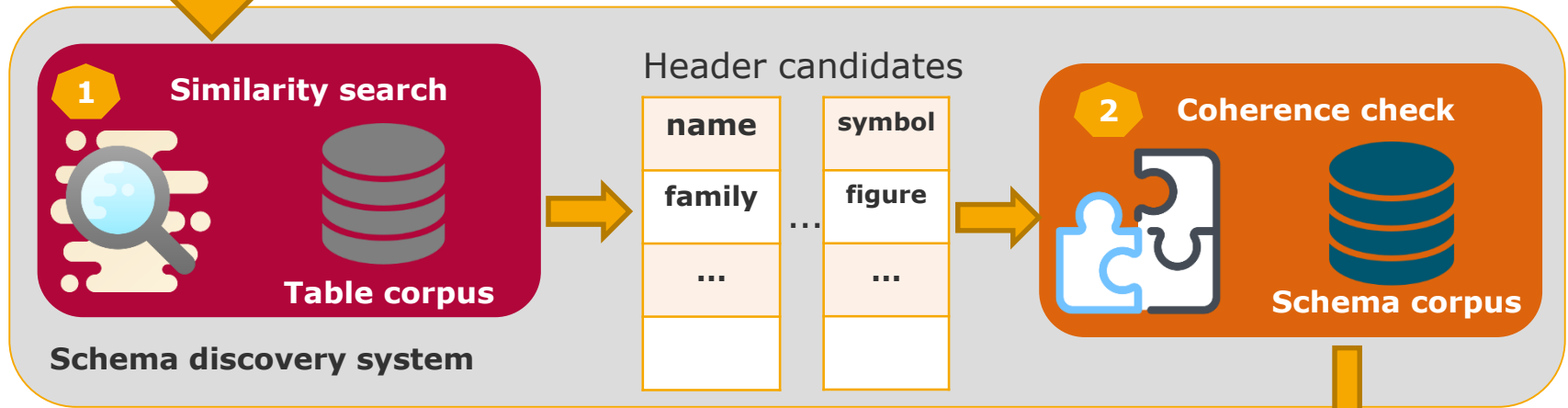
Relational Header Discovery

Hazar Harmouch
ICDE 2021

Solution

End-to-end schema discovery

?	?	?	?
stark	winter is coming	winterfell	grey wolf
targaryen	fire and blood	dragonstone	red dragon
lannister	hear me roar!	casterly rock	golden lion



family	motto	location	symbol
--------	-------	----------	--------

Schema discovery using similarity search

Phase 1: Similarity search - Challenges



- ❑ **Efficient** similarity search
 - **Data sketches** reduce memory consumption.

- ❑ **Un-biased** column similarity measure.

- ❑ Filtering out **dissimilar** tables.

?	City	Location
Berlin	Potsdam	Frankfurt
Potsdam	Frankfurt	Potsdam
	Cologne	Cologne
		Dresden
		Munich

?	?	?	?
Stark	Winter is coming	Winterfell	Grey wolf
Targaryen	Fire And Blood	Dragonstone	Red dragon
Lannister	Hear Me Roar!	Casterly rock	Golden lion



University	Year of establishment	Students	Location
Heidelberg University	1386	31,500	Heidelberg
Leipzig University	1409	29,500	Leipzig
University of Rostock	1419	14,000	Rostock
University of Greifswald	1456	12,000	Greifswald
University of Freiburg	1457	24,700	Freiburg
University of Munich	1472	50,918	Munich ^[5]
University of Tübingen	1477	28,700	Tübingen
University of Halle-Wittenberg	1502	18,500	Halle ^[6]
University of Marburg	1527	24,000	Marburg
University of Jena	1558	19,000	Jena

$$\text{Jaccard coefficient: } J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$\text{Jaccard containment: } C(A, B) = \frac{|A \cap B|}{|A|}$$

- HyperLogLog (HLL) is a data sketch that
 - has compact form ($\sim 1\text{kB}$).
 - estimates cardinalities with an error < 0.01 .

Schema discovery using similarity search

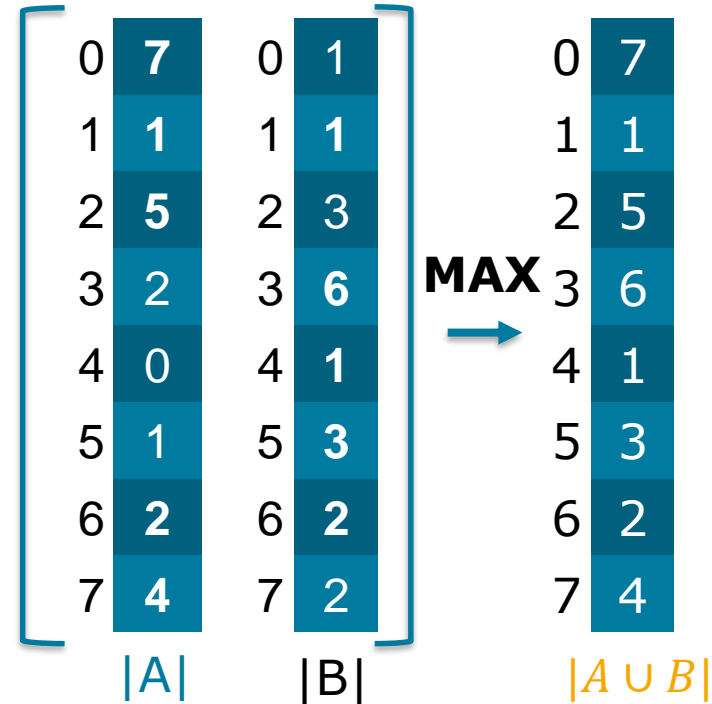
Phase 1: Similarity search - HLL-based similarities

$$\text{Jaccard coefficient: } J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$\text{Jaccard containment: } C(A, B) = \frac{|A \cap B|}{|A|}$$

□ Inclusion-exclusion principle:

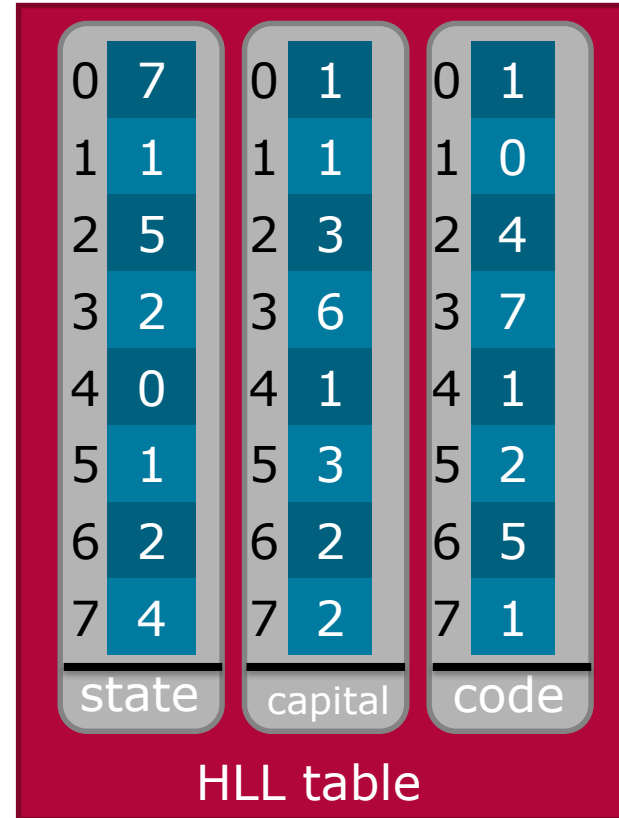
$$|A \cup B| = |A| + |B| - |A \cap B|$$



Schema discovery using similarity search

Phase 1: Similarity search – HLL table data structure

State	Capital	Code
Austria	Vienna	AT
Belgium	Brussels	BE
Bulgaria	Sofia	BG
Croatia	Zagreb	HR



Schema discovery using similarity search

Phase 1: Similarity search

Offline



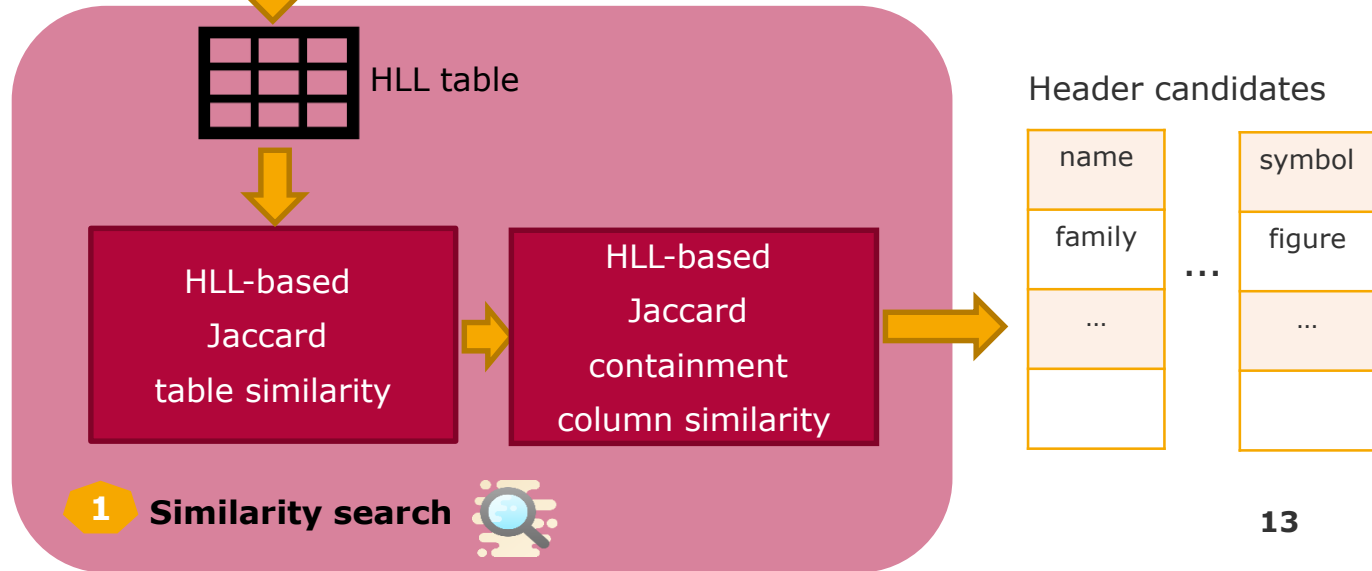
Wiki web tables
32G (1.6M Tables)

Transform



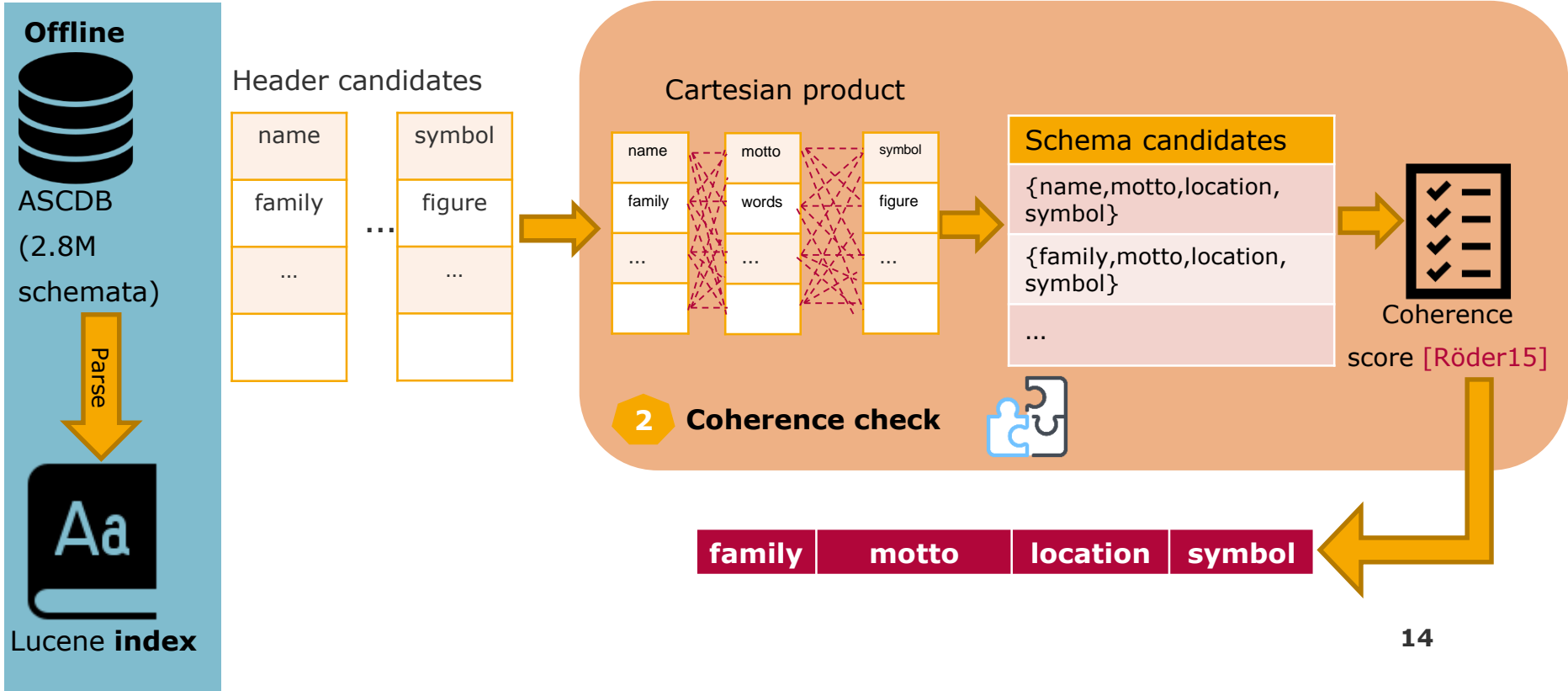
HLL tables
450MB

?	?	?	?
Stark	Winter is coming	Winterfell	Grey wolf
Targaryen	Fire And Blood	Dragonstone	Red dragon
Lannister	Hear Me Roar!	Casterly rock	Golden lion



Schema discovery using similarity search

Phase 2: Coherence check



Schema discovery using similarity search

Experimental setup



Wiki web tables

32G (1.6M Tables)



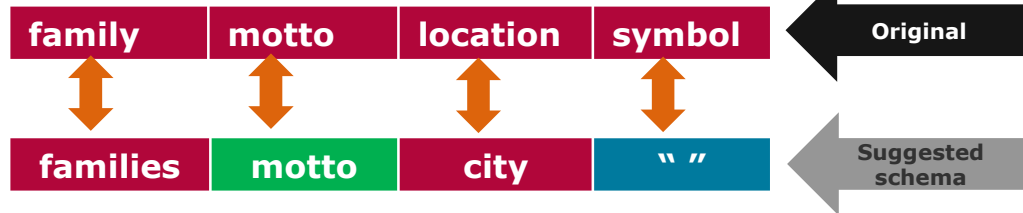
5% Test as input

(4770 Tables)



95% as Table corpus

(162799 Tables)



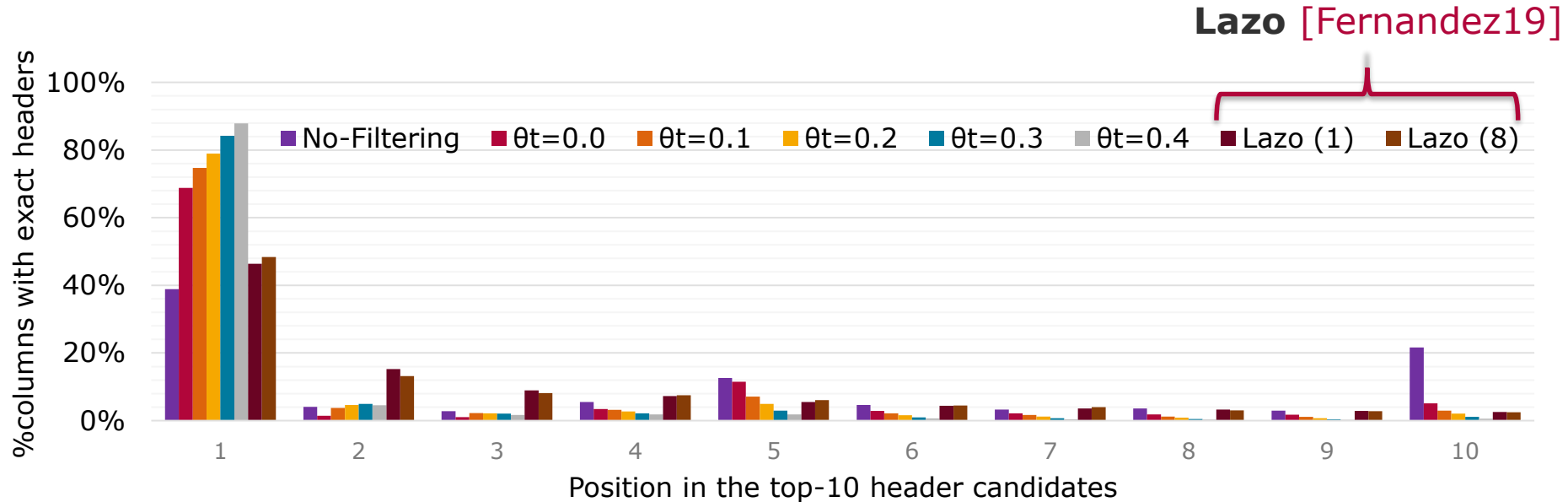
- ❑ **TP**: exactly re-discovered headers.
- ❑ **FP**: differing headers from the original schema.
- ❑ **FN**: no-header suggestions.

**Relational Header
Discovery**

Hazar Harmouch
ICDE 2021

Schema discovery using similarity search

Results: Similarity search (1/2)



Schema discovery using similarity search

Results: Similarity search (2/2)

θ	% exact header	% differing header	% no header
Lazo (1)	30.95	38.52	30.52
Lazo (8)	31.63	39.60	28.77
Lazo (10)	Out of memory		
No filtering	62.73	36.16	1.10
0	74.00	24.10	1.10
0.1	61.66	13.49	24.85
0.2	42.90	7.92	49.18
0.3	27.25	4.02	68.73
0.4	16.77	1.98	81.25

Lazo
[Fernandez19]

**Relational Header
Discovery**

Hazar Harmouch
ICDE 2021

Schema discovery using similarity search

Results: Error analysis - Schema-level



Plural

{year, champions}
{year, champion}

Synonym

{men's, women's}
{male, female}

Time

{team, city, state, stadium, capacity, 2008 season}
{team, city, state, stadium, capacity, 2009 season}

Short form

{no., date, score, opponent, record}
{number, date, score, opponent, record}

Symbol

{place, rider, number, country, machine, points, wins}
{place, rider, #, country, machine, points, wins}

Granularity

{year, population}
{year, village population}

Manual test

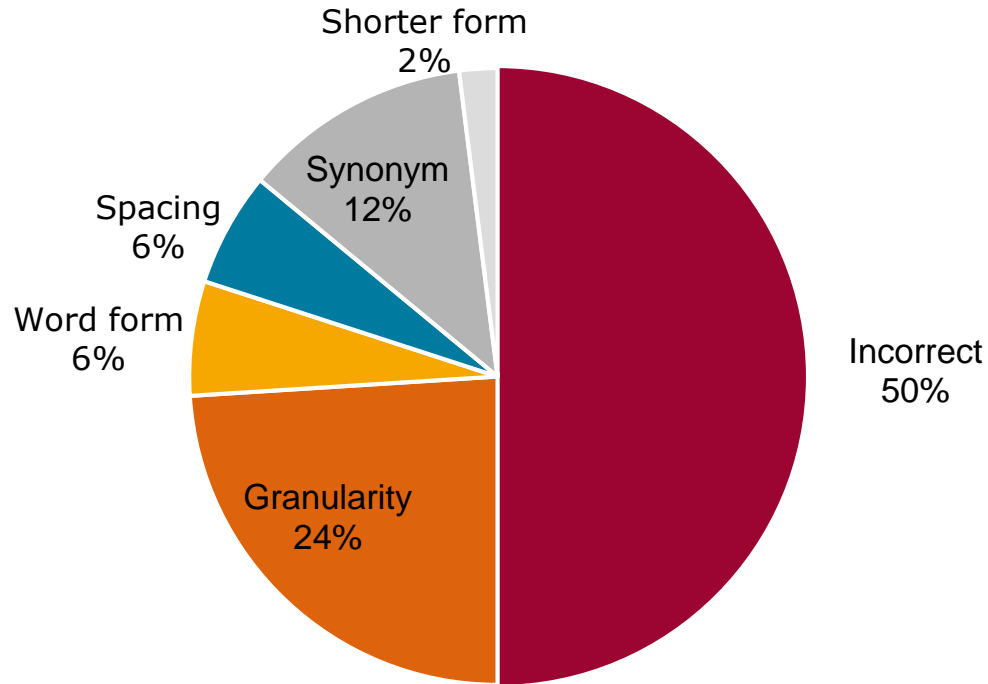
- 100 schemata
- 400 headers

Relational Header Discovery

Hazar Harmouch
ICDE 2021

Schema discovery using similarity search

Results: Error analysis - Individual headers



Manual test

- 100 schemata
- 400 headers

Relational Header Discovery

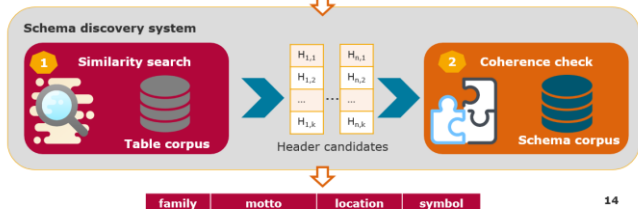
Hazar Harmouch
ICDE 2021

Summary

Solution End-to-end schema discovery



?	?	?	?
Stark	Winter is coming	Winterfell	Grey wolf
Targaryen	Fire And Blood	Dragonstone	Red dragon
Lannister	Hear Me Roar!	Casterly rock	Golden lion



14

Schema discovery using similarity search Phase 1: Similarity search



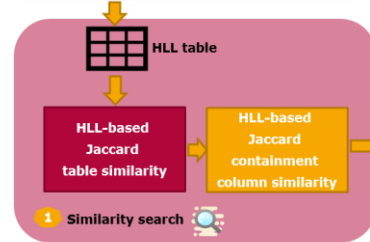
Offline

Wiki web tables
32G (1.6M Tables)

Transform

HLL tables
450MB

?	?	?	?
Stark	Winter is coming	Winterfell	Grey wolf
Targaryen	Fire And Blood	Dragonstone	Red dragon
Lannister	Hear Me Roar!	Casterly rock	Golden lion

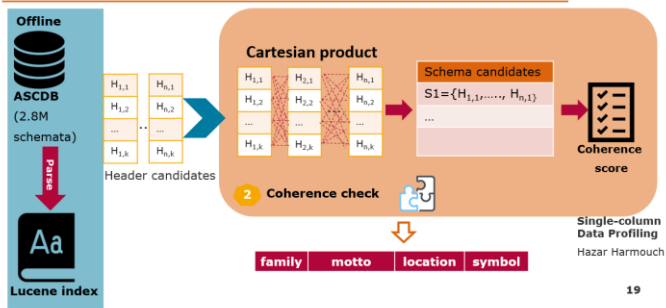


Header candidates

name	symbol
family	figure
...	...
...	...

23

Schema discovery using similarity search Phase 2: Coherence check



19

Schema Discovery using similarity search Results Error Analysis-Schema level



Plural	{year, champions} {Year, Champion}
Synonym	{men's, women's} {male, female}
Time	{team, city, state, stadium, capacity, 2008 season} {team, city, state, stadium, capacity, 2009 season}
Short form	{no., date, score, opponent, record} {number, date, score, opponent, record}
Symbol	{Place, Rider, Number, Country, Machine, Points, Wins} {place, rider, #, country, machine, points, wins}
Granularity	{year, population} {year, village population}

Original

← Top-1 schema

Manual test

- 100 schemata
- 400 header

Single-column
Data Profiling
Hazar Harmouch,
PhD candidate
Chart 16

Relational Header Discovery

Hazar Harmouch
ICDE 2021

20

References

- [Neumaier16] Neumaier, S., Umbrich, J., Parreira, J. X., & Polleres, A. (2016, October). Multi-level semantic labelling of numerical values. In International Semantic Web Conference (pp. 428-445). Springer, Cham.
- [Mitlöhner16] Mitlöhner, J., Neumaier, S., Umbrich, J., & Polleres, A. (2016, August). Characteristics of open data csv files. In Open and Big Data (OBD), International Conference on (pp. 72-79). IEEE.
- [Balakrishnan15] S. Balakrishnan, A. Y. Halevy, B. Harb, H. Lee, J. Madhavan, A. Rostamizadeh, W. Shen, K. Wilder, F. Wu, and C. Yu. Applying webtables in practice. In CIDR, 2015.
- [Cafarella08] M. J. Cafarella, A. Y. Halevy, Y. Zhang, D. Z. Wang, and E. Wu. Uncovering the relational web. In WebDB, 2008.
- [Fernandez19] Raul Castro Fernandez, Jisoo Min, Demitri Nava, and Samuel Madden. Lazo: Acardinality-based method for coupled estimation of jaccard similarity and contain-ment. In Proceedings of the IEEE International Conference on Data Engineering(ICDE), 2019.