

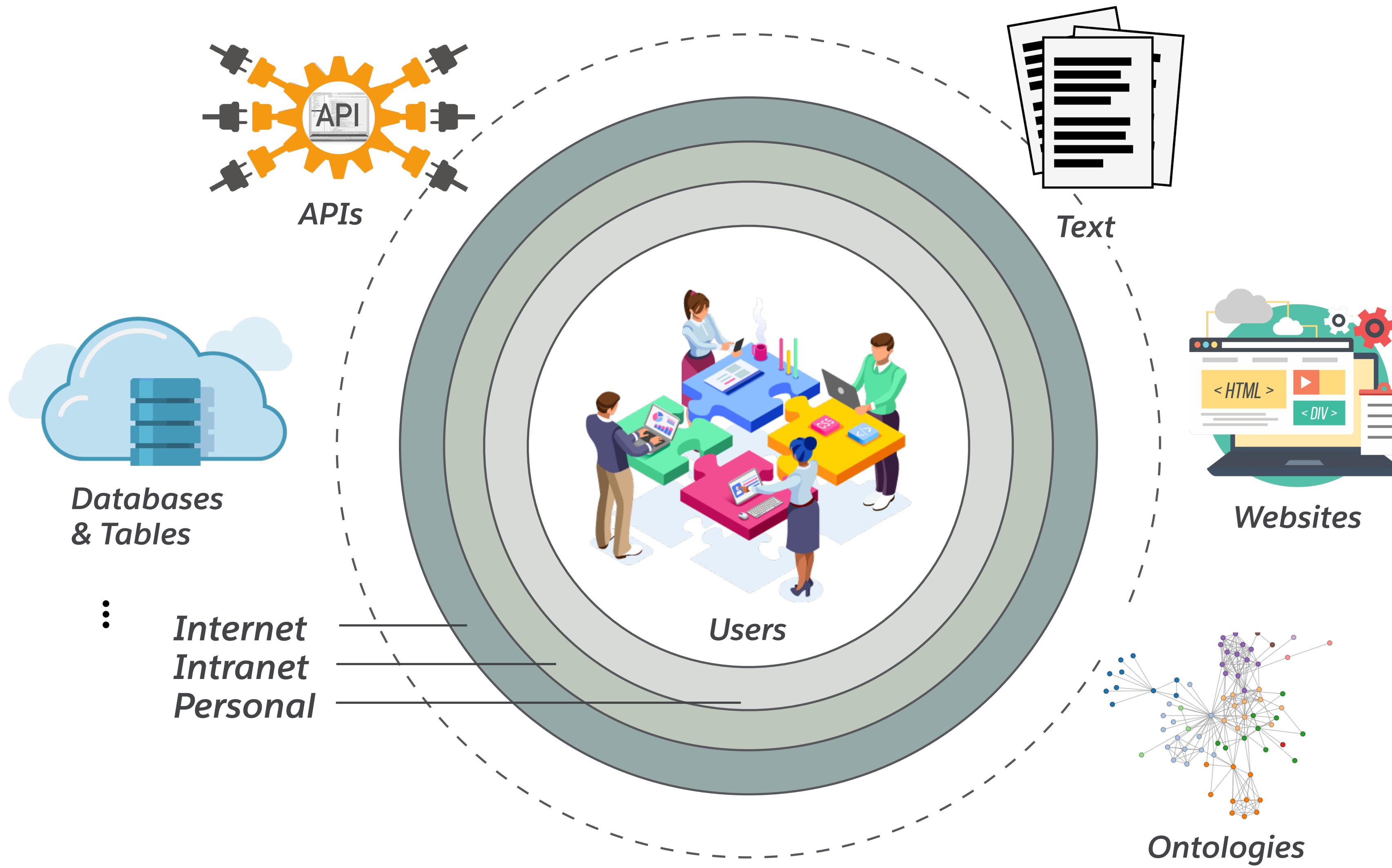
Bridging Textual and Tabular Data: Is Attention All We Need?

Invited Talk - KR2ML Workshop @NeurIPS'20

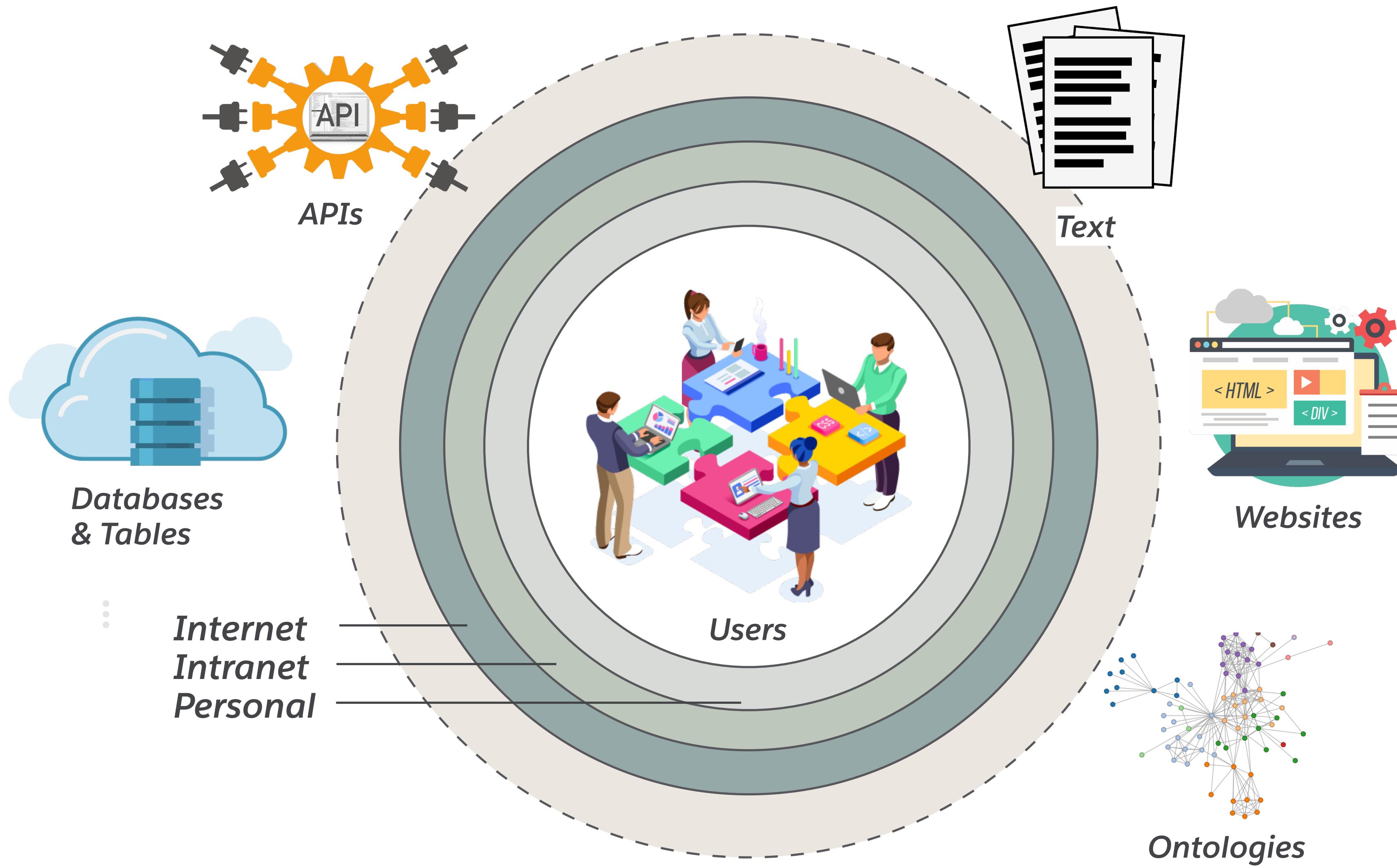


Victoria Lin
Salesforce AI Research

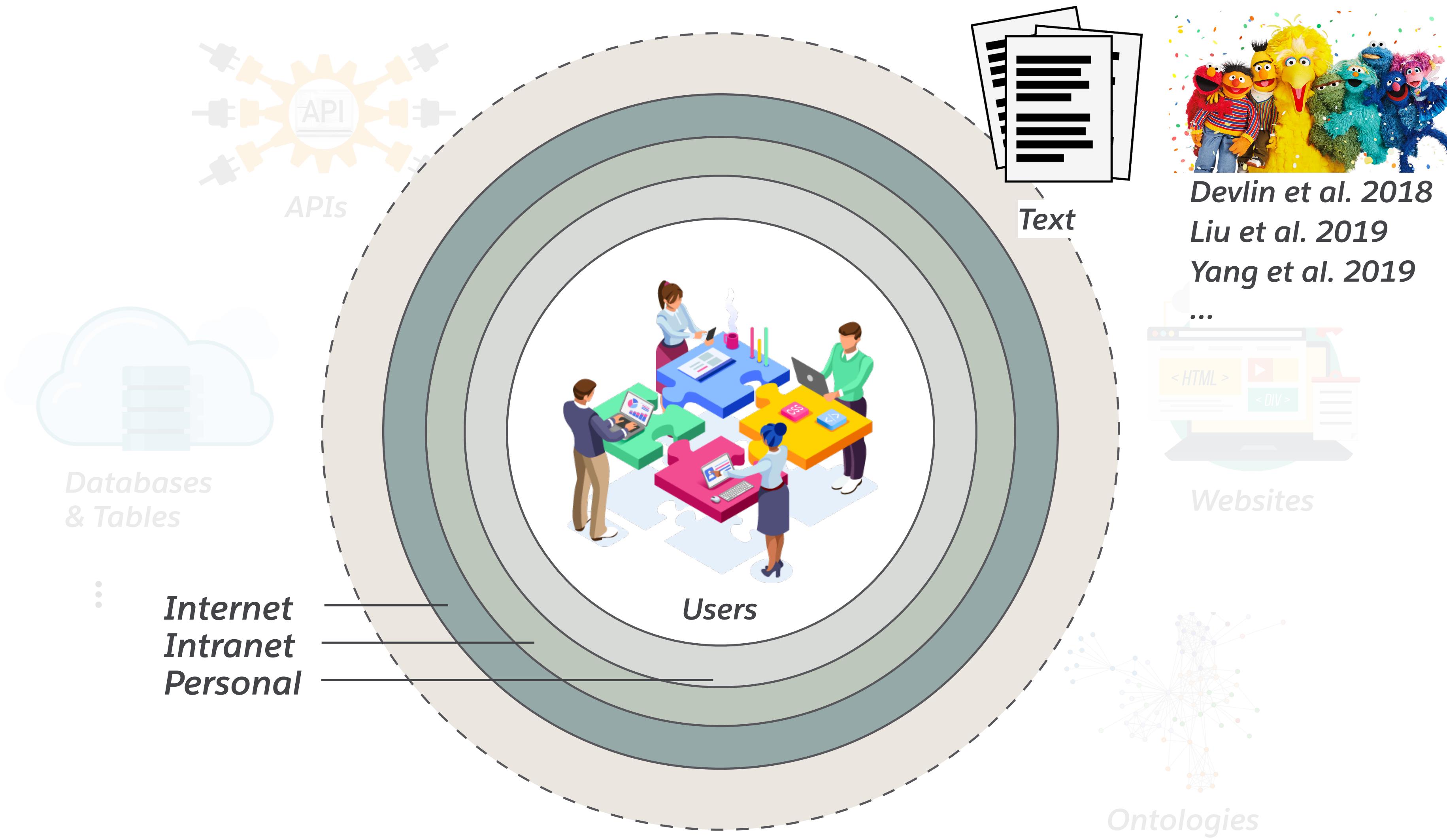
Reasoning over Heterogeneous Data



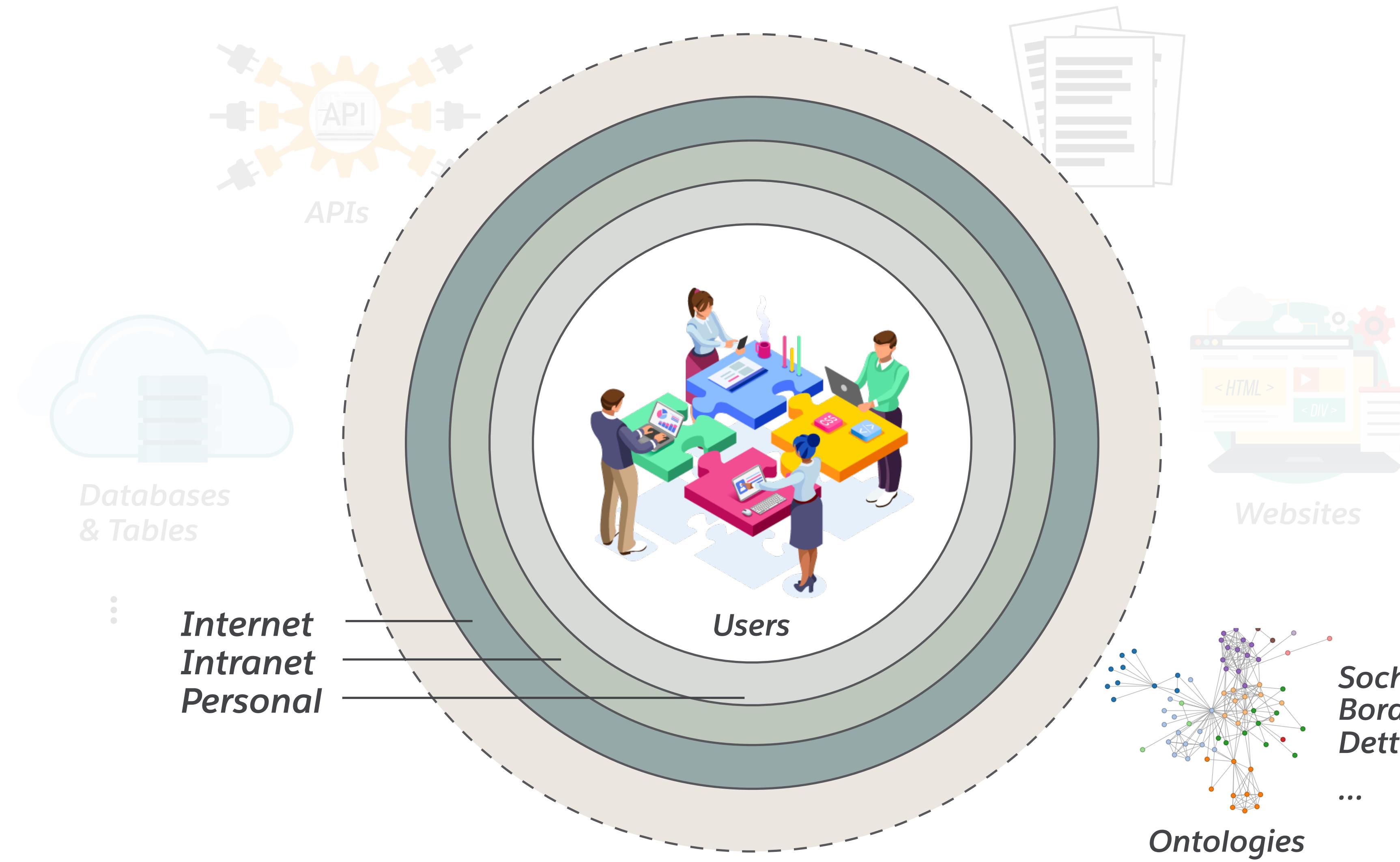
Representing Heterogeneous Data



Deep Representations



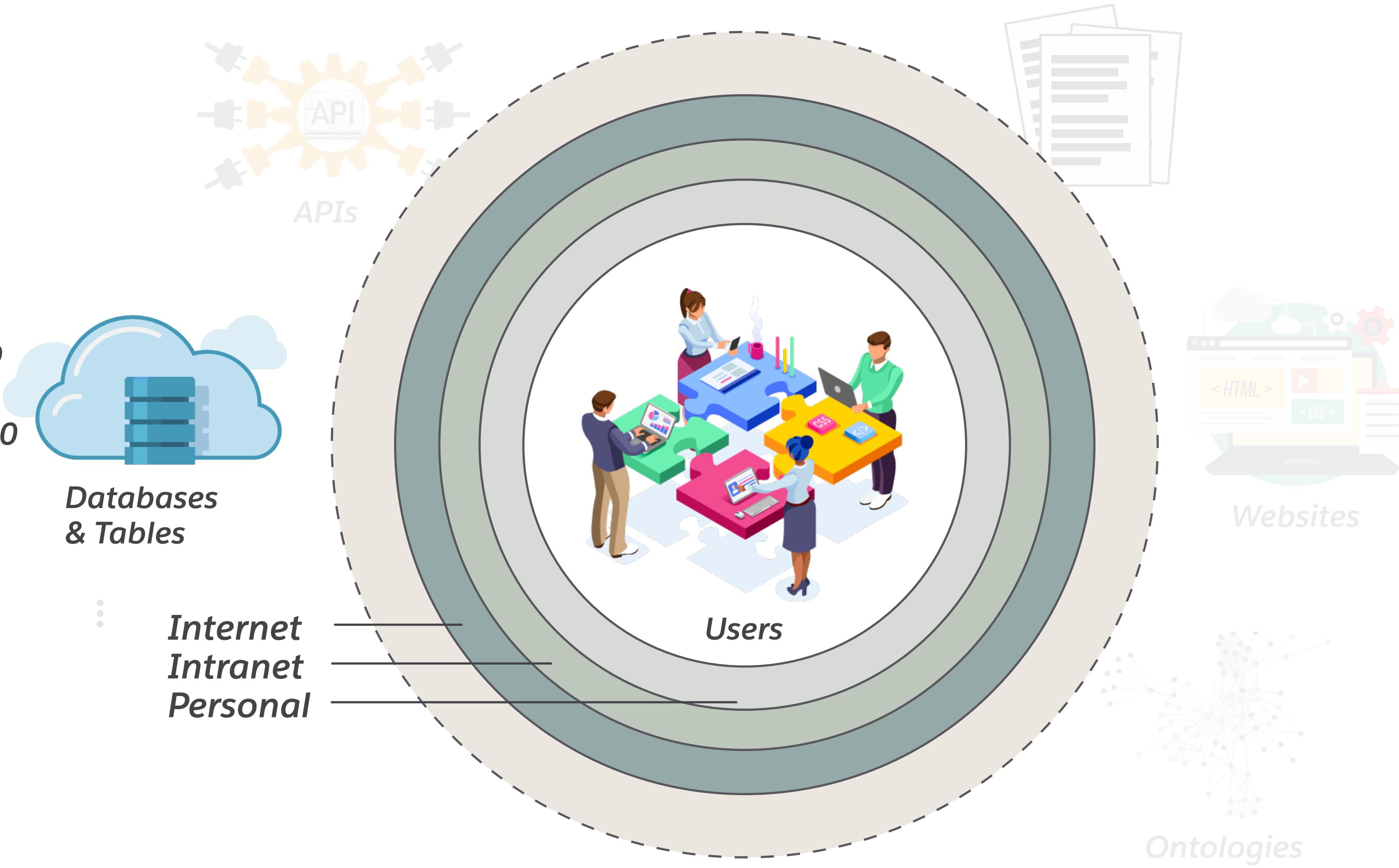
Deep Representations



Deep Representations



Deng et al. 2019
Yin et al. 2020
Herzig et al. 2020
Yu et al. 2020
Deng et al. 2021
...



Deep Representations

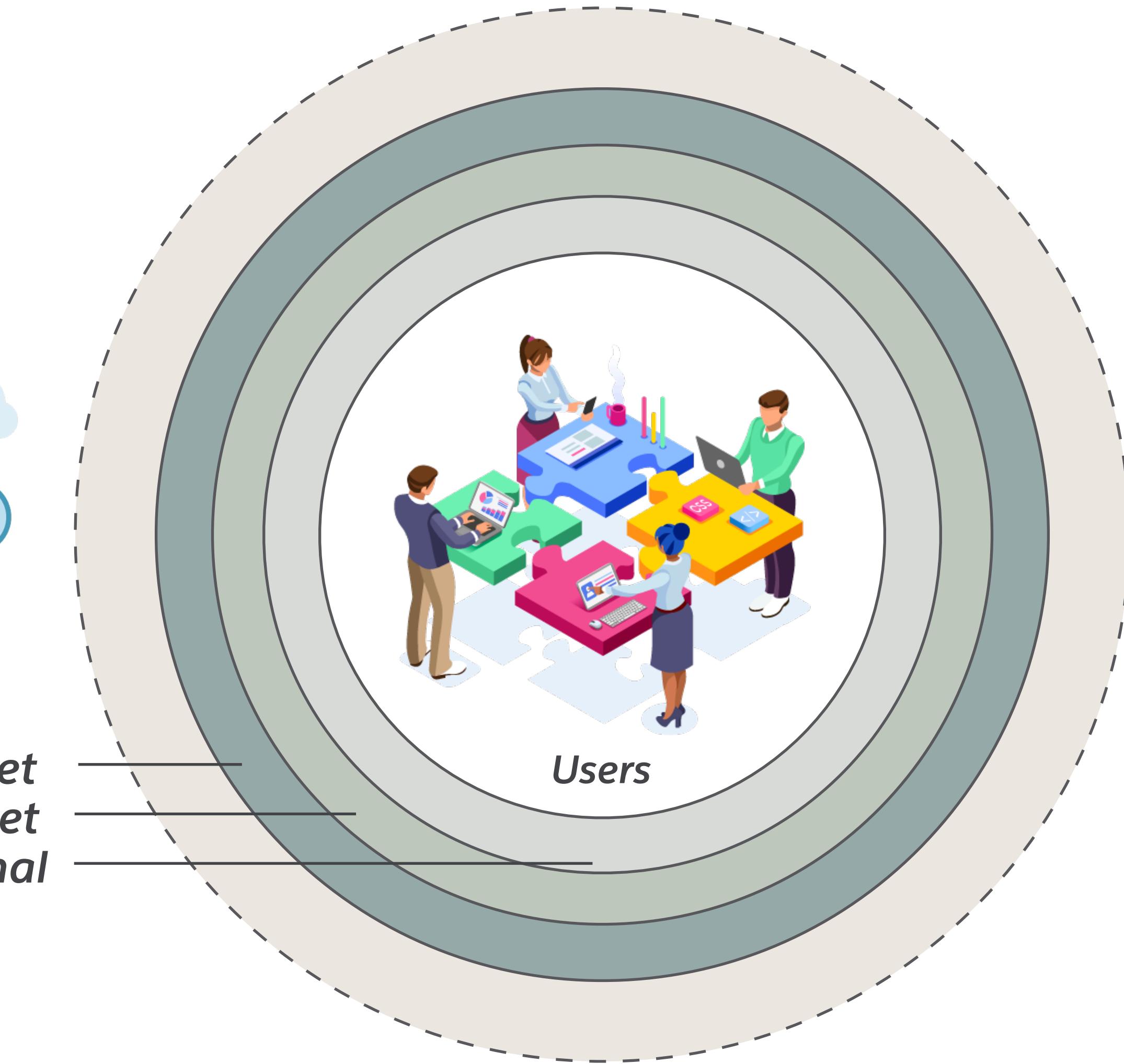


Deng et al. 2019
Yin et al. 2020
Herzig et al. 2020
Yu et al. 2020
Deng et al. 2021
...



Databases & Tables

Internet
Intranet
Personal



Benefits of learning continuous representation for structured data

Deep Representations

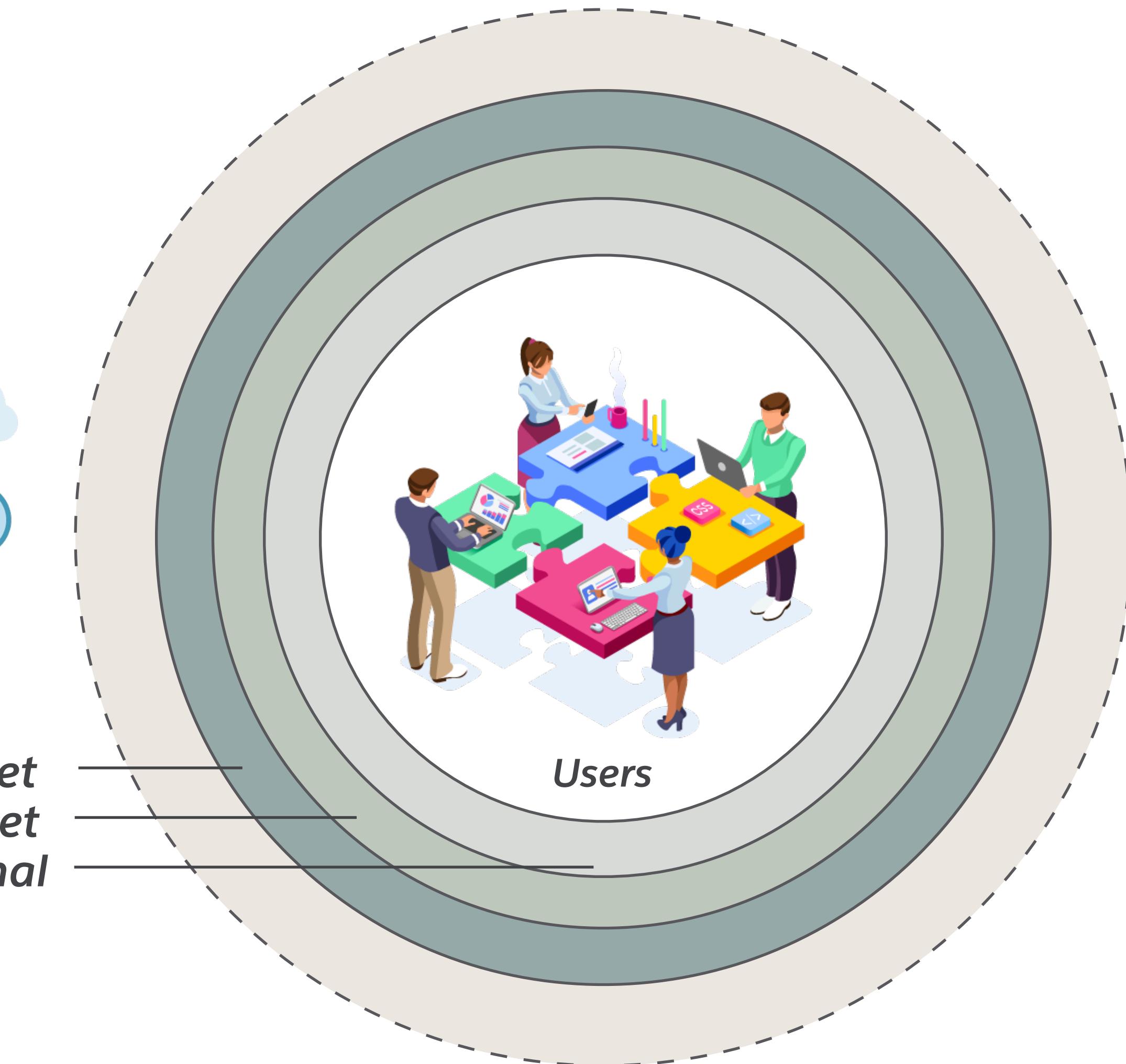


Deng et al. 2019
Yin et al. 2020
Herzig et al. 2020
Yu et al. 2020
Deng et al. 2021
...



Databases & Tables

*Internet
Intranet
Personal*



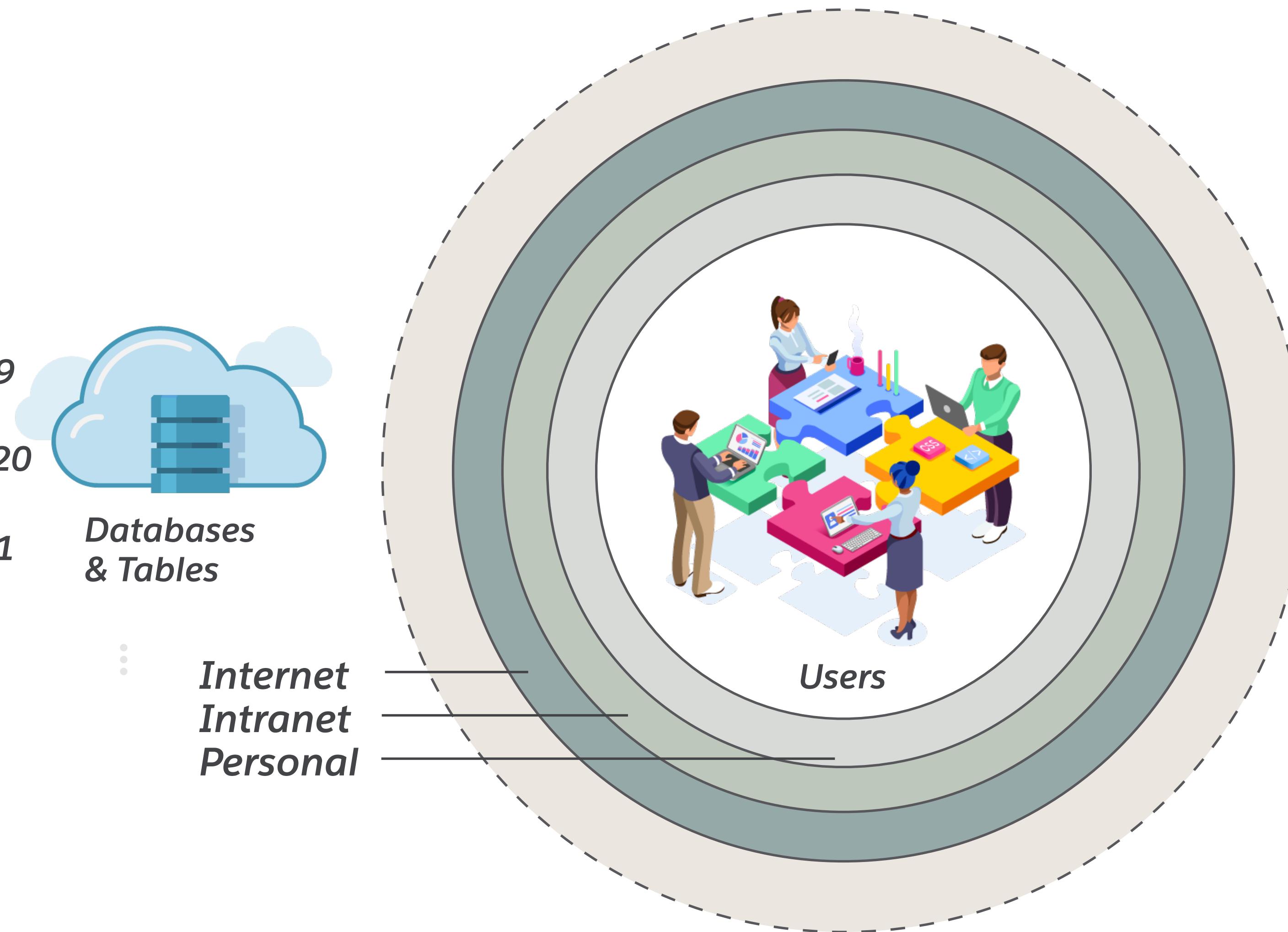
Benefits of learning continuous representation for structured data

- Continuous representations typically outperform sparse, discrete representations in data-driven prediction tasks

Deep Representations



Deng et al. 2019
Yin et al. 2020
Herzig et al. 2020
Yu et al. 2020
Deng et al. 2021
...



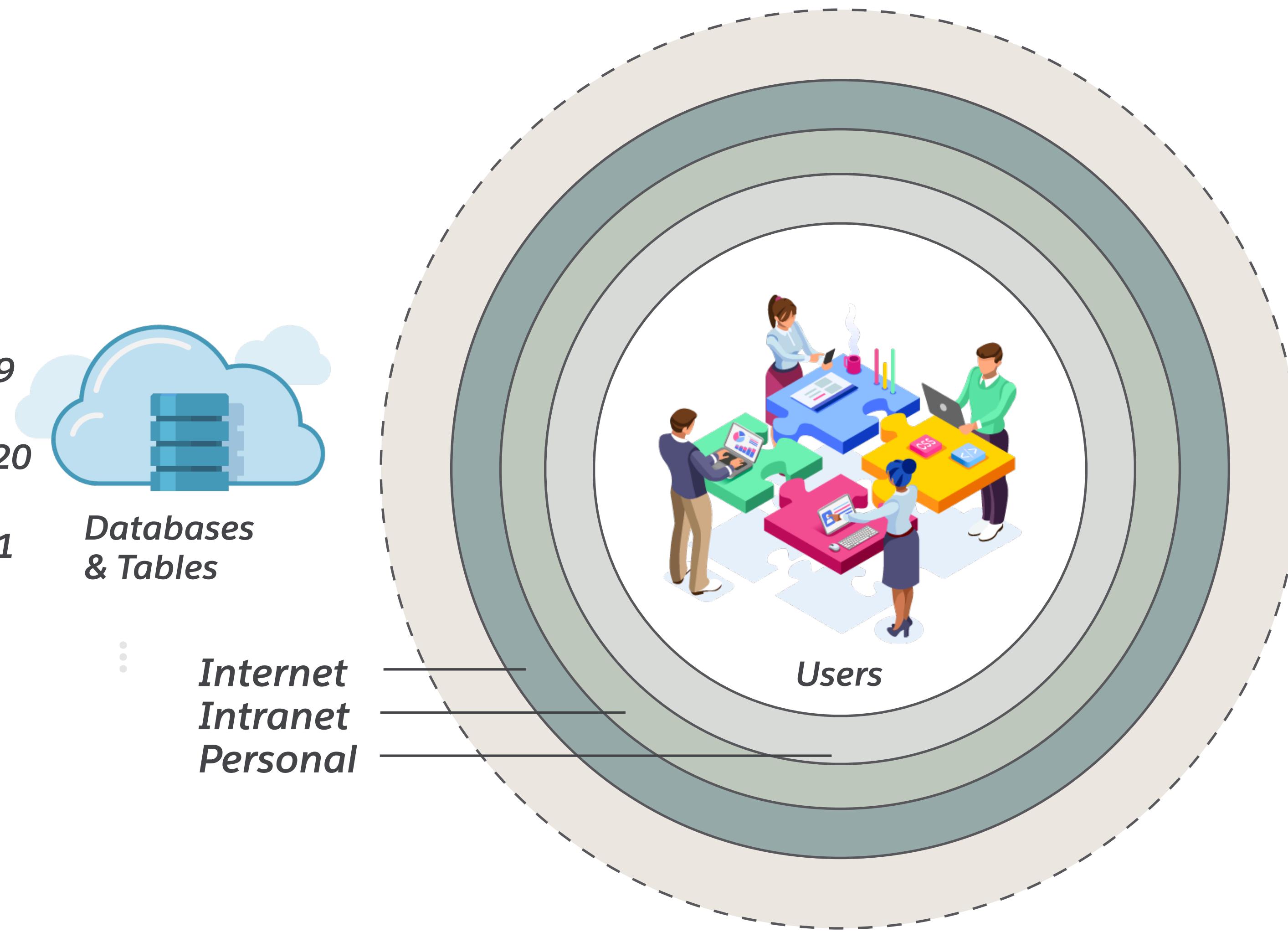
💡 Benefits of learning continuous representation for structured data

- Continuous representations typically outperform sparse, discrete representations in data-driven prediction tasks
- Most structured data have lexical sub-components (e.g. entity names/relation labels in ontologies, table headers/cells)
- Many structured data are heterogeneous on their own (e.g. databases may contain both text and image cells)

Deep Representations



Deng et al. 2019
Yin et al. 2020
Herzig et al. 2020
Yu et al. 2020
Deng et al. 2021
...



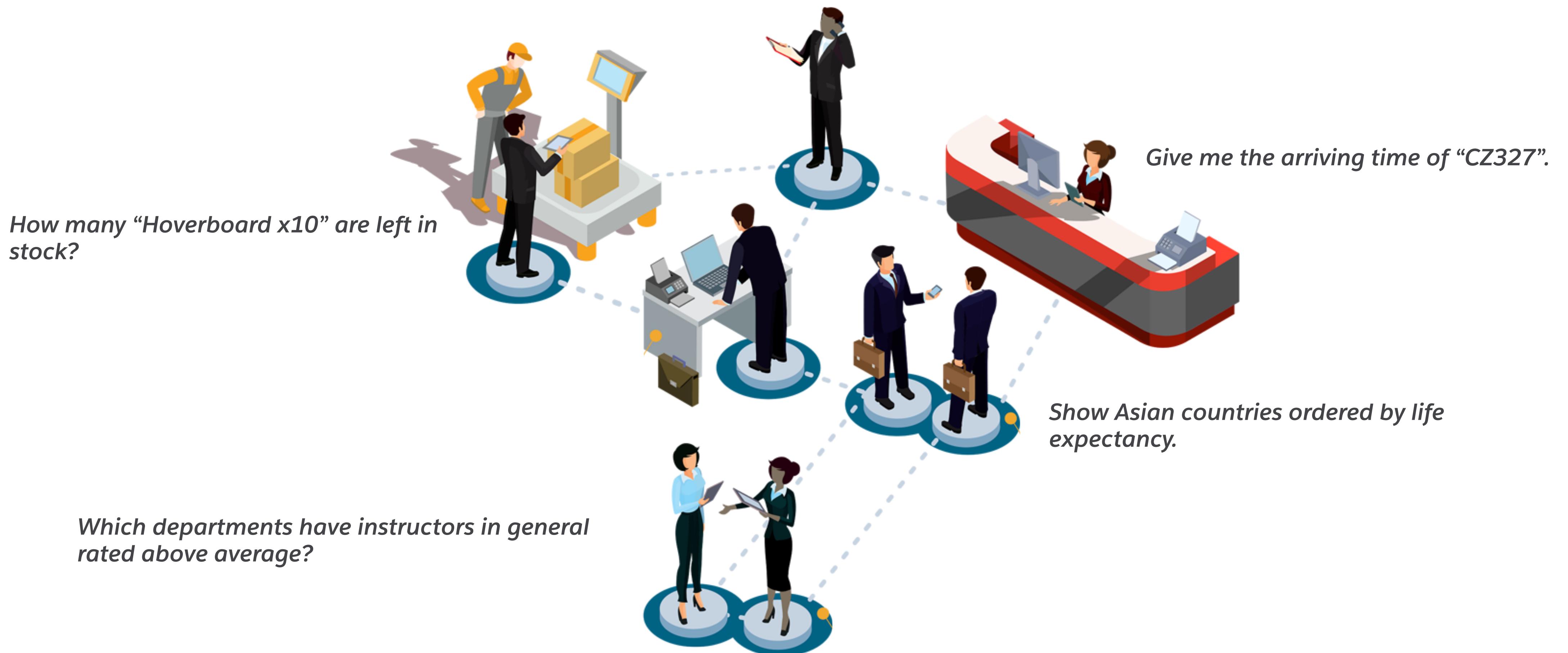
💡 Benefits of learning continuous representation for structured data

- Continuous representations typically outperform sparse, discrete representations in data-driven prediction tasks
- Most structured data have lexical sub-components (e.g. entity names/relation labels in ontologies, table headers/cells)
- Many structured data are heterogeneous on their own (e.g. databases may contain both text and image cells)
- Inducing “uniform” joint representation for both structured and unstructured data (e.g. text)

Representation Learning for Databases



Our goal is to learn **representations** for **tables** and **databases** in the **vector space** to power downstream AI applications, in particular, **voice interfaces**.



Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Relational Database

← Table Name

User ID	Name	Nationality	Partition ID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

 A table typically represents a type of entities (or events).

Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Table Header

Relational Database

User Profile

	Name of user	Nationality of user	Partition ID of user	# followers of user	...
UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

 A column name defines the semantic relationship between columns and between a column and the entity/event type represented by the table.

Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Row



A row is typically an instantiation of the entity/event.

Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

column/field



Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Integer
↑

Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

String
↑

Relational Database

User Profile

UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Integer

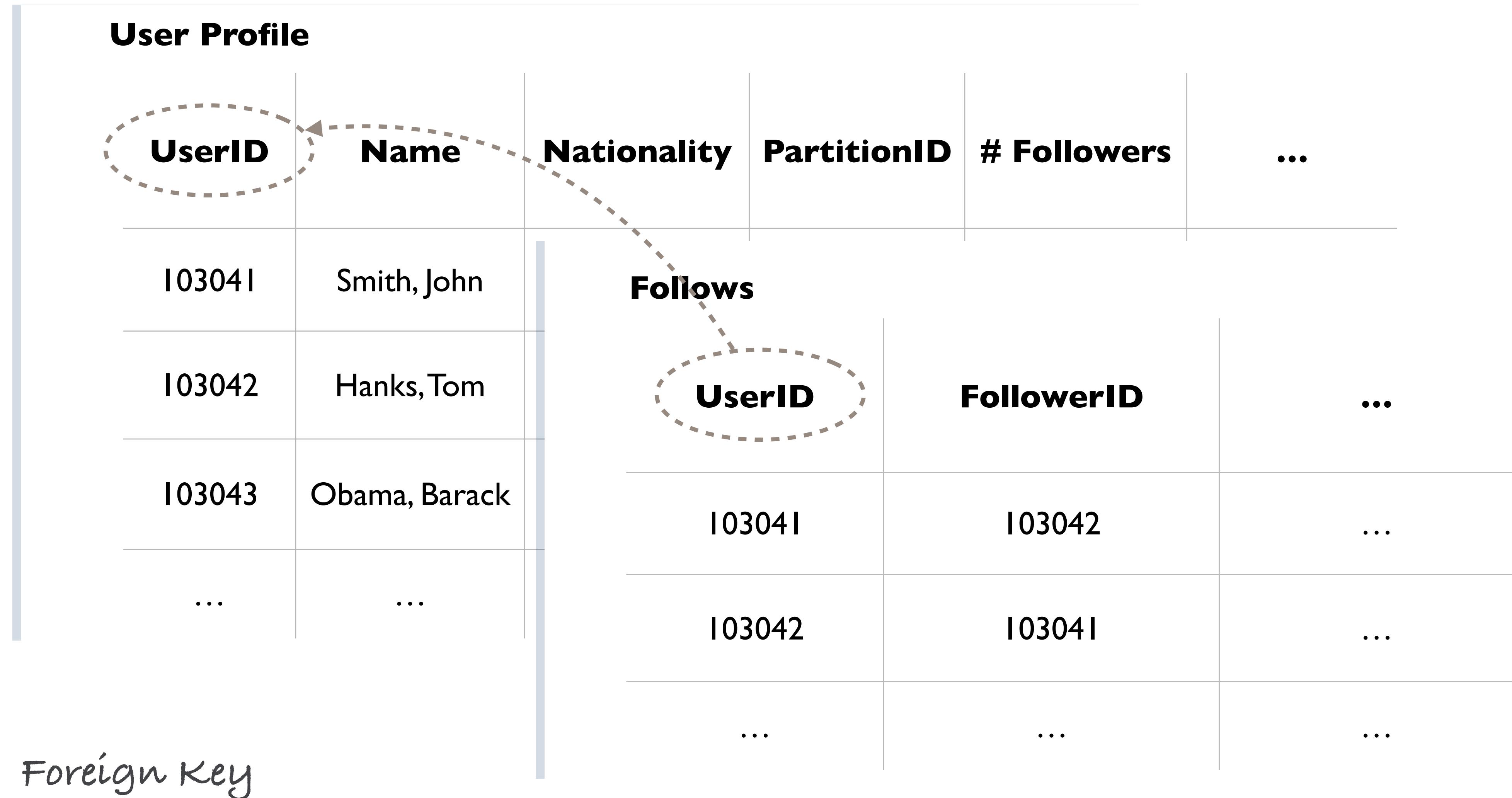
Relational Database

User Profile

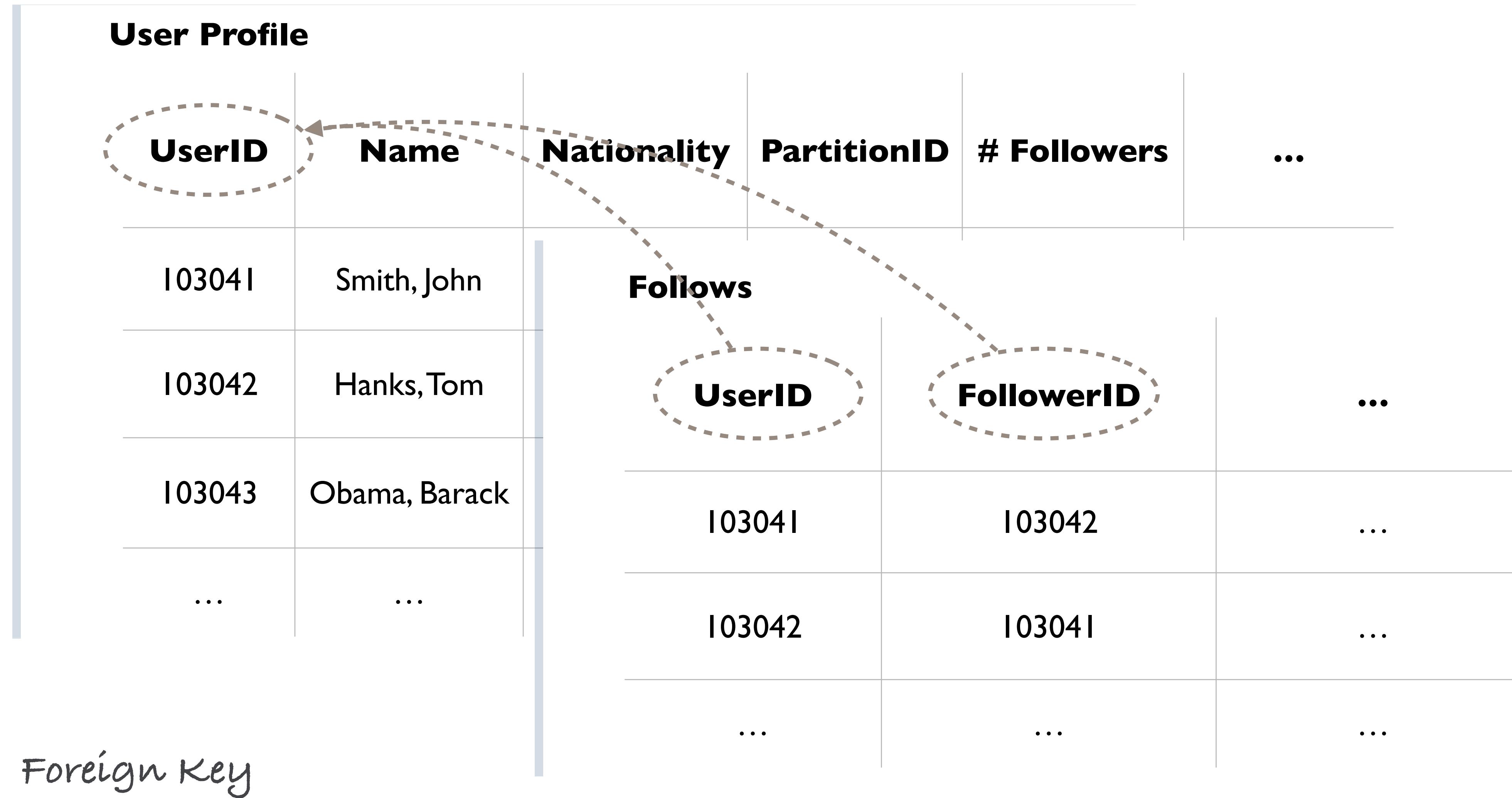
UserID	Name	Nationality	PartitionID	# Followers	...
103041	Smith, John	Canada	P10	3	...
103042	Hanks, Tom	United States	P11	16.6M	...
103043	Obama, Barack	United States	P11	127M	...
...

Primary Key
↑

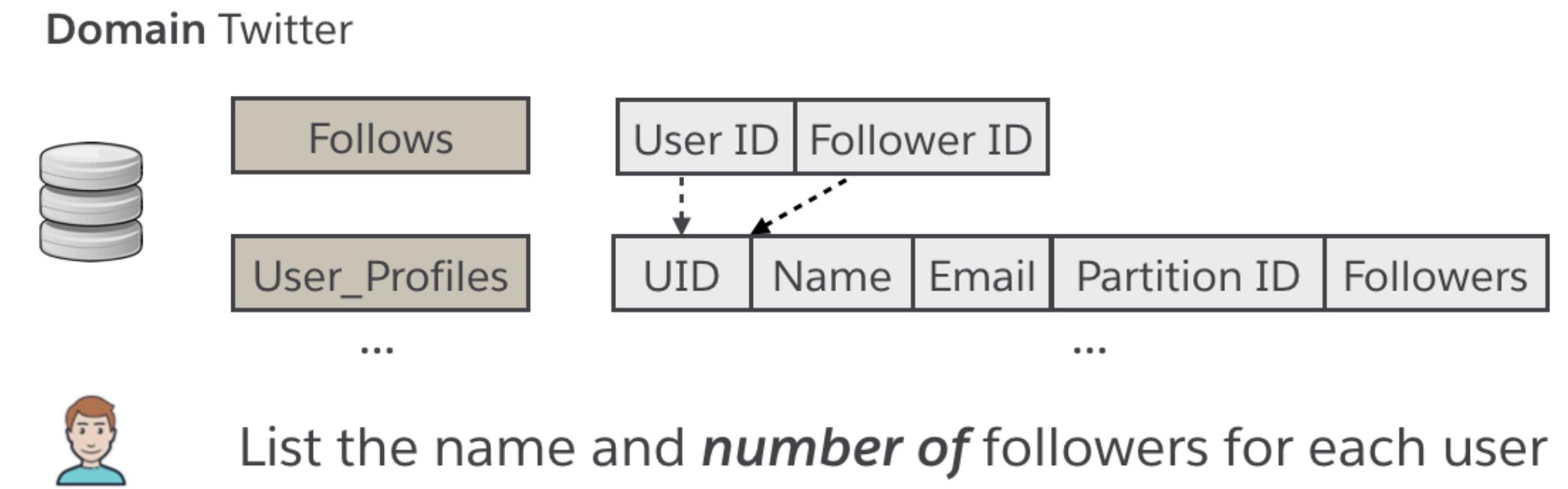
Relational Database



Relational Database

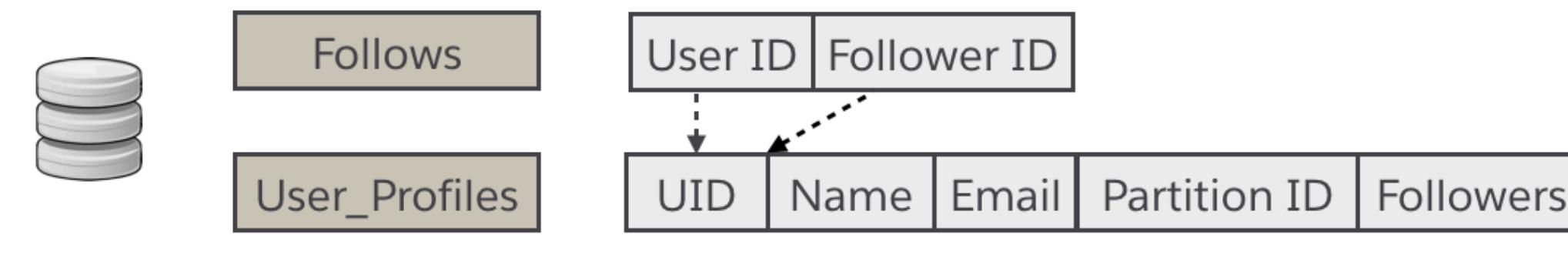


Question Answering over Databases



Structured Query Language

Domain Twitter

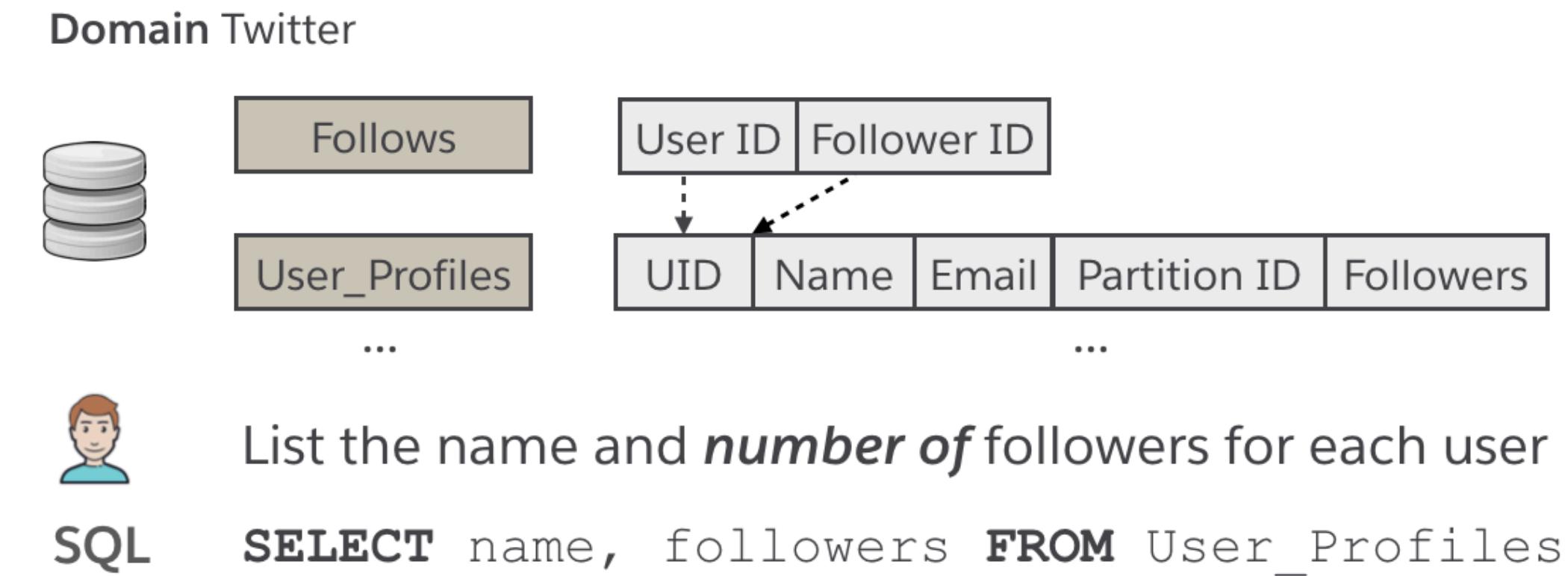


List the name and *number of* followers for each user

SQL

```
SELECT name, followers FROM User_Profiles
```

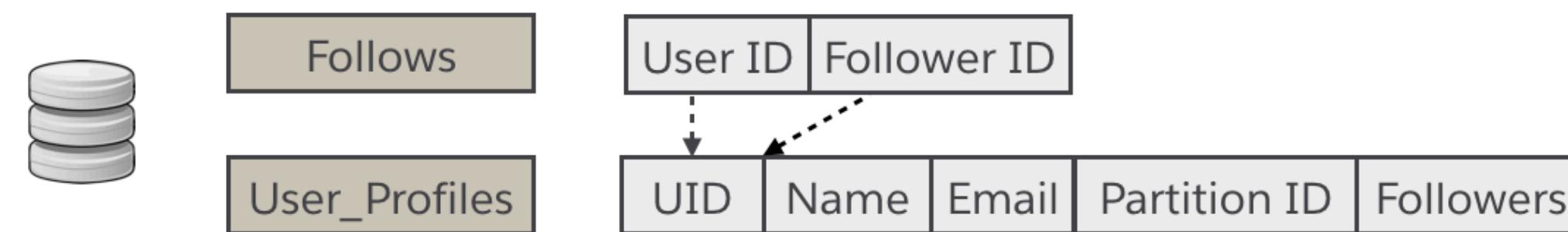
Text-to-SQL Semantic Parsing



Strong/Fully
supervised

Text-to-SQL Semantic Parsing

Domain Twitter



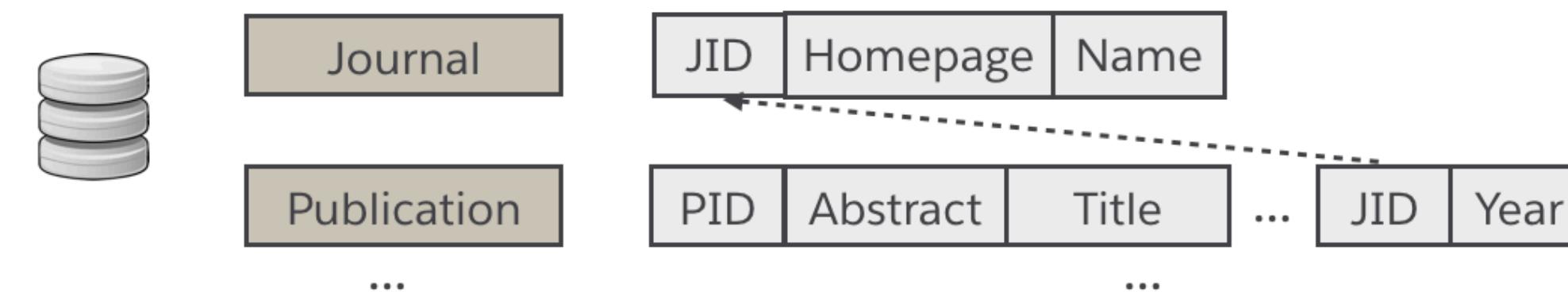
List the name and *number of* followers for each user

SQL

```
SELECT name, followers FROM User_Profiles
```

Cross-
Database

Domain Academic



Return me the *number of* papers on PVLDB

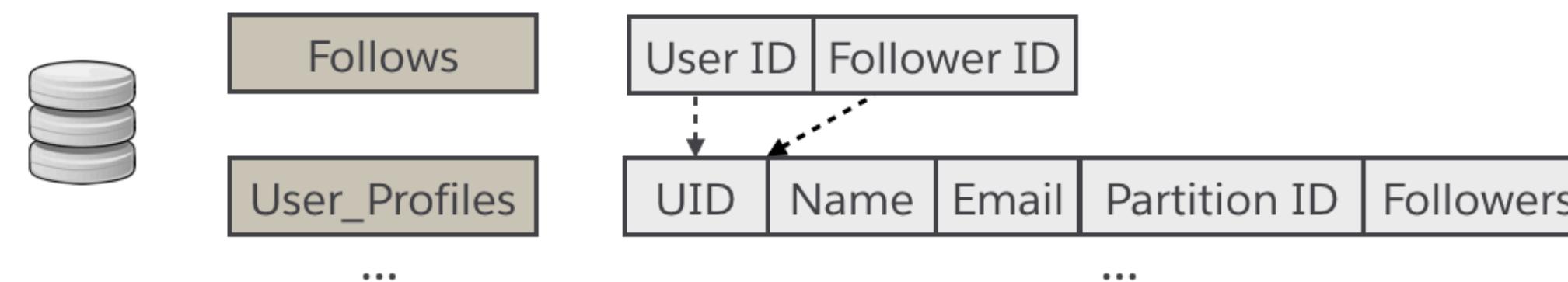
SQL

```
SELECT COUNT(DISTINCT t2.title)
FROM Publication AS T2 JOIN Journal AS T1
ON T2.JID = T1.JID WHERE T1.name = "PVLDB"
```

Challenges and Observations

Challenge 1: Questions with similar intent may map to very different SQL logical forms when issued to different DBs.

Domain Twitter



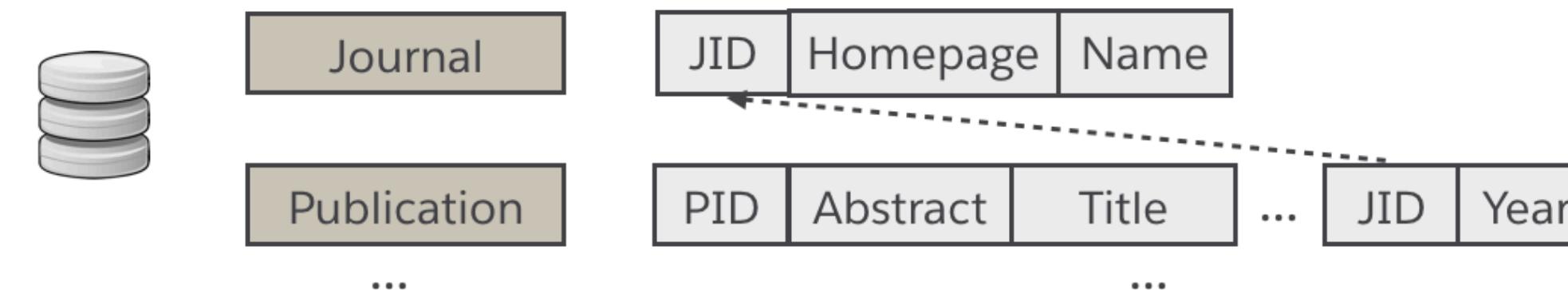
List the name and *number of* followers for each user

SQL

```
SELECT name, followers FROM User_Profiles
```

Cross-
Database

Domain Academic

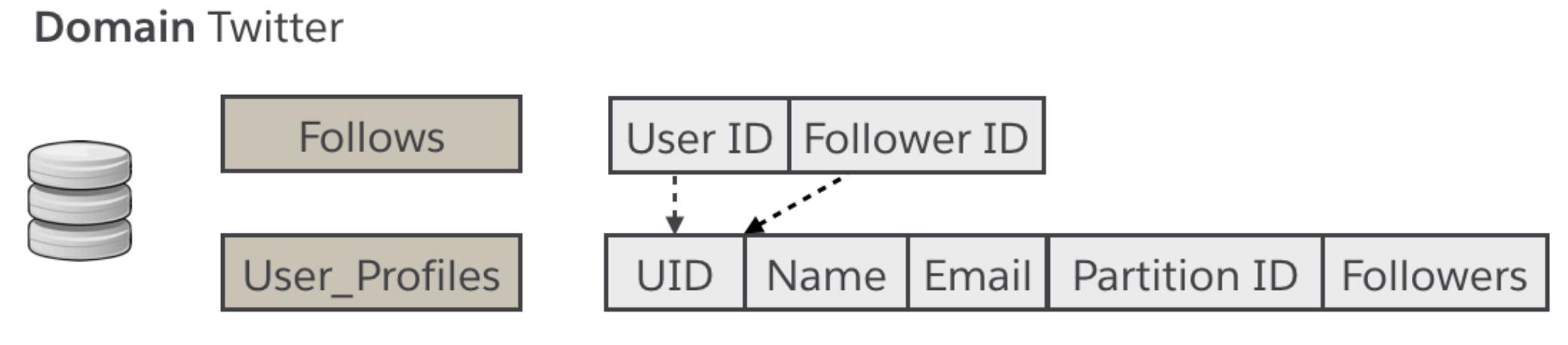


Return me the *number of* papers on PVLDB

SQL

```
SELECT COUNT(DISTINCT t2.title)
FROM Publication AS T2 JOIN Journal AS T1
ON T2.JID = T1.JID WHERE T1.name = "PVLDB"
```

Challenges and Observations



>List the name and *number of* followers for each user

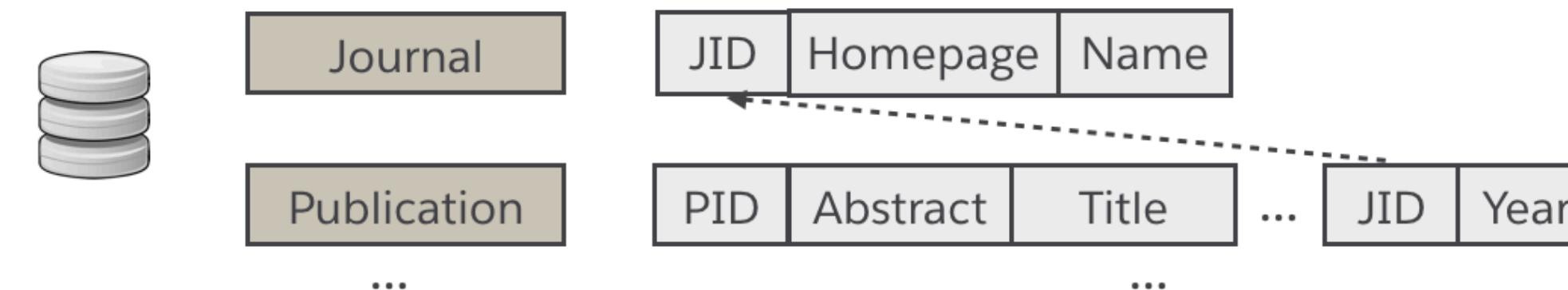
SQL `SELECT name, followers FROM User_Profiles`

Cross-
Database

Challenge 1: Questions with similar intent may map to very different SQL logical forms when issued to different DBs.

Observation 1: The meaning of a column is critical for predicting the output SQL query, and it is captured in both the column name and the corresponding cell values.

Domain Academic



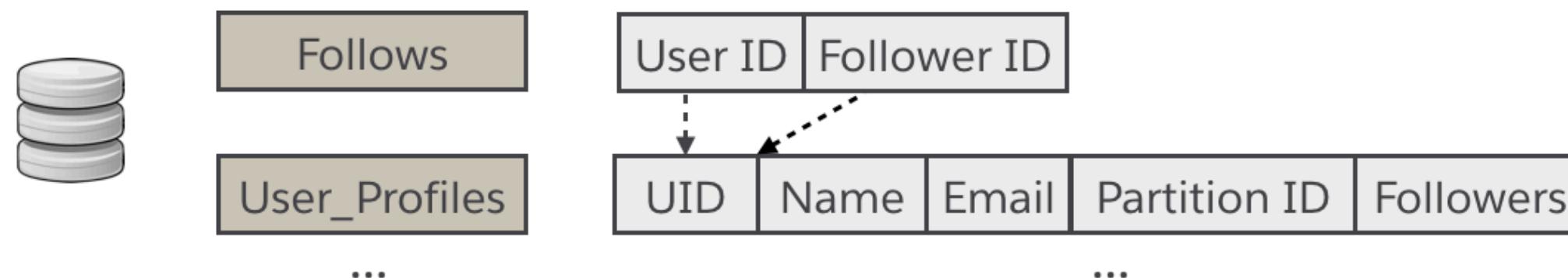
Return me the *number of* papers on PVLDB

SQL `SELECT COUNT(DISTINCT t2.title)
FROM Publication AS T2 JOIN Journal AS T1
ON T2.JID = T1.JID WHERE T1.name = "PVLDB"`

Challenges and Observations

Challenge 2: Quite often we encounter long-tail, domain-specific entities.

Domain Twitter

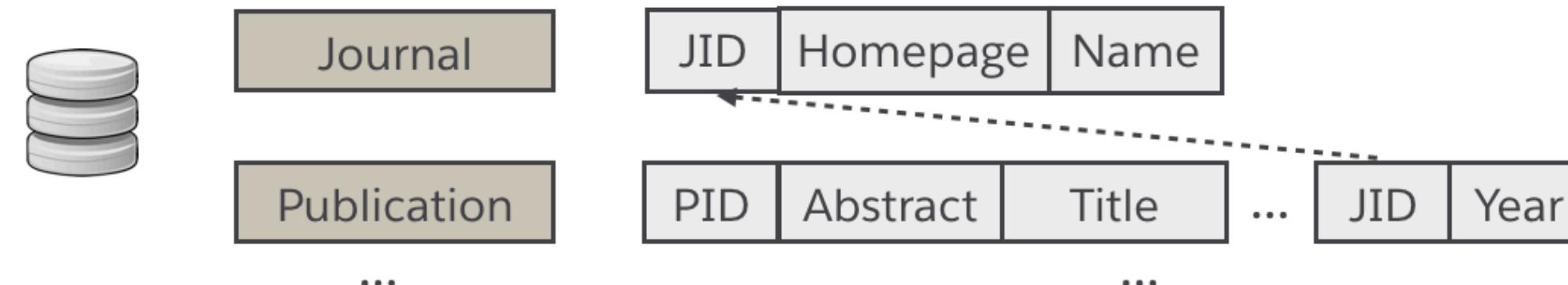


List the name and *number of* followers for each user

SQL

```
SELECT name, followers FROM User_Profiles
```

Domain Academic



Domain Diversity



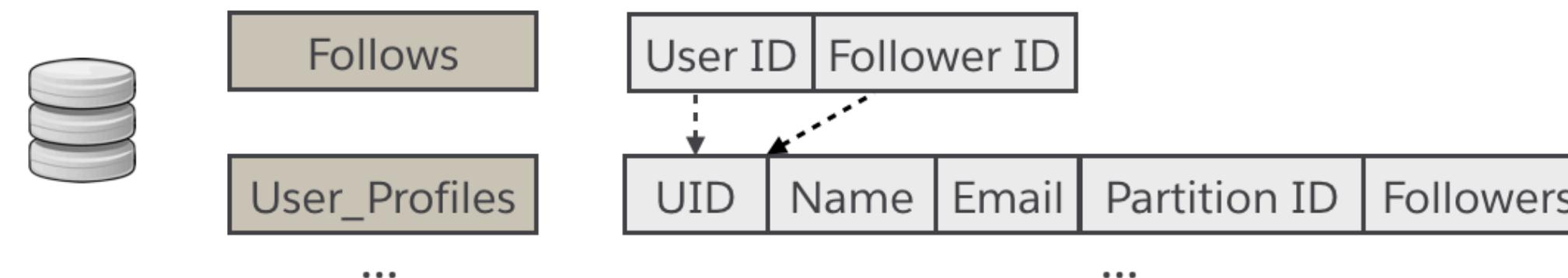
Return me the *number of* papers on PVLDB

SQL

```
SELECT COUNT(DISTINCT t2.title)
FROM Publication AS T2 JOIN Journal AS T1
ON T2.JID = T1.JID WHERE T1.name = "PVLDB"
```

Challenges and Observations

Domain Twitter



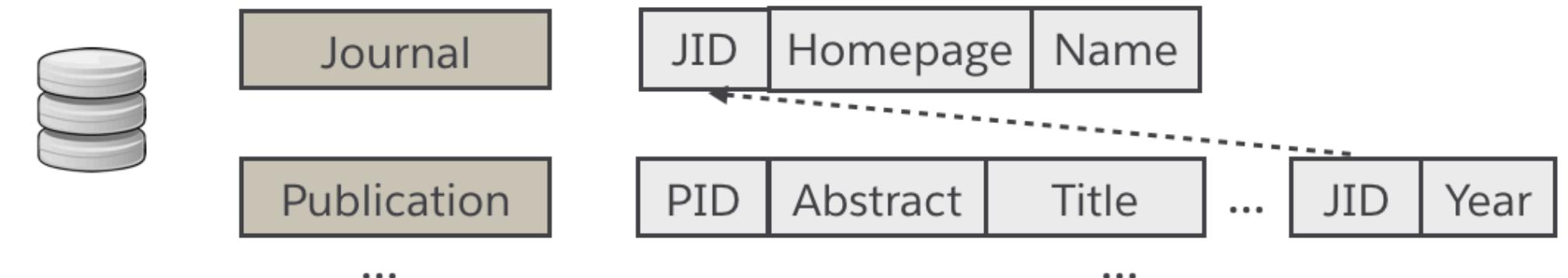
List the name and *number of* followers for each user

SQL

```
SELECT name, followers FROM User_Profiles
```

Challenge 2: Quite often the question contains rare domain-specific entities.

Domain Academic



Return me the *number of* papers on PVLDB

SQL

```
SELECT COUNT(DISTINCT t2.title)
FROM Publication AS T2 JOIN Journal AS T1
ON T2.JID = T1.JID WHERE T1.name = "PVLDB"
```

Observation 2: The meaning of such entities can often be determined by grounding them to the database.

Domain Diversity

Model

Question

Database Schema

 Show names of properties that are either houses or apartments

Properties					
Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	
Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...

Picklists	
Apartment	...
Field	...
House	...
Shop	...
Other	...

Model



Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

Question

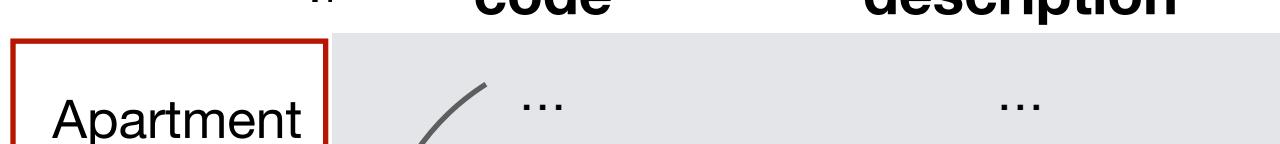
Database Schema

 Show names of properties that are either houses or apartments

Properties					
Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other


Reference Property Types

Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...


Picklists

Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...

Model

Research Question:

More specifically, can pre-trained transformer language models effectively encode the question, the semantics of a relational DB model and their joint contextualization?

Question

Database Schema



Show names of properties that are either houses or apartments

Properties					
Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	
Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...

Picklists	
Apartment	...
Field	...
House	...
Shop	...
Other	...

Textual-Tabular Data Encoding

Serialize Table Header/DB Schema



 Show names of properties that are either houses or apartments

Properties					
Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	
Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...

Textual-Tabular Data Encoding

Serialize Table Header/DB Schema



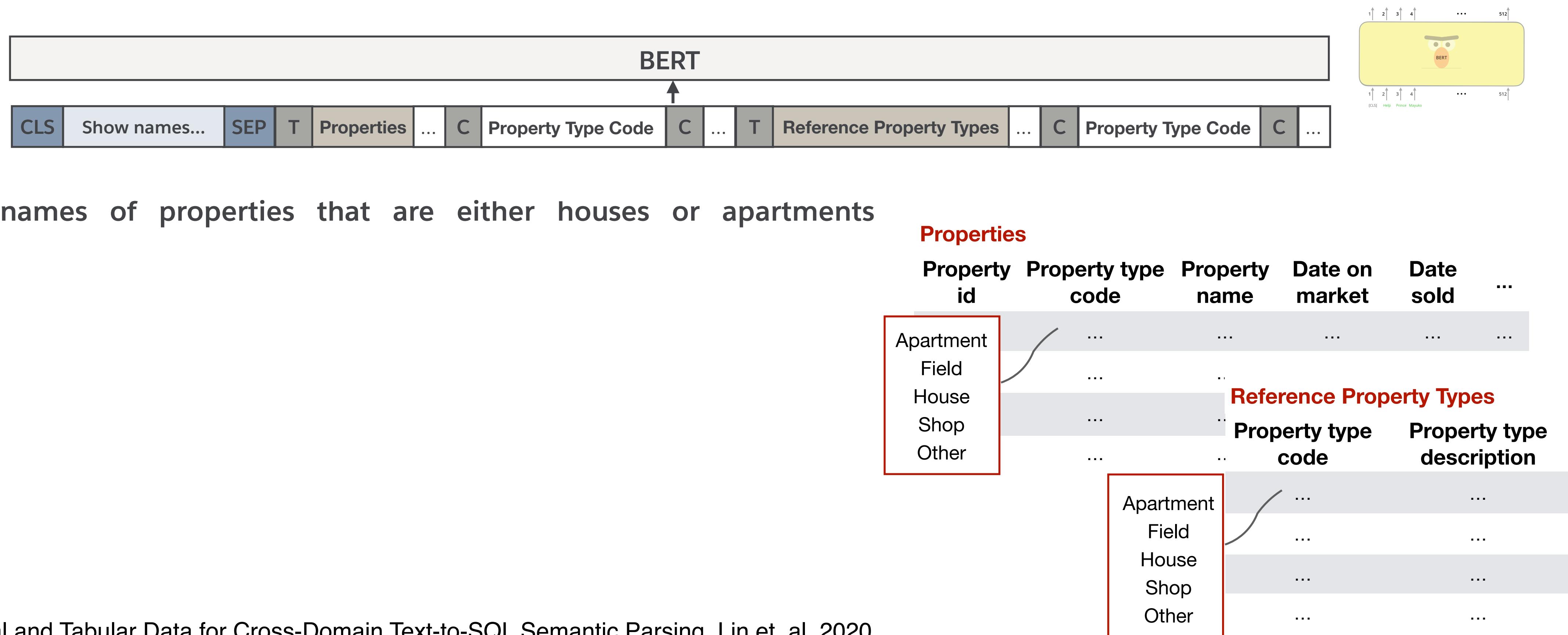
 Show names of properties that are either houses or apartments

Properties					
Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	
Property type code	Property type description
Apartment	...
Field	...
House	...
Shop	...
Other	...

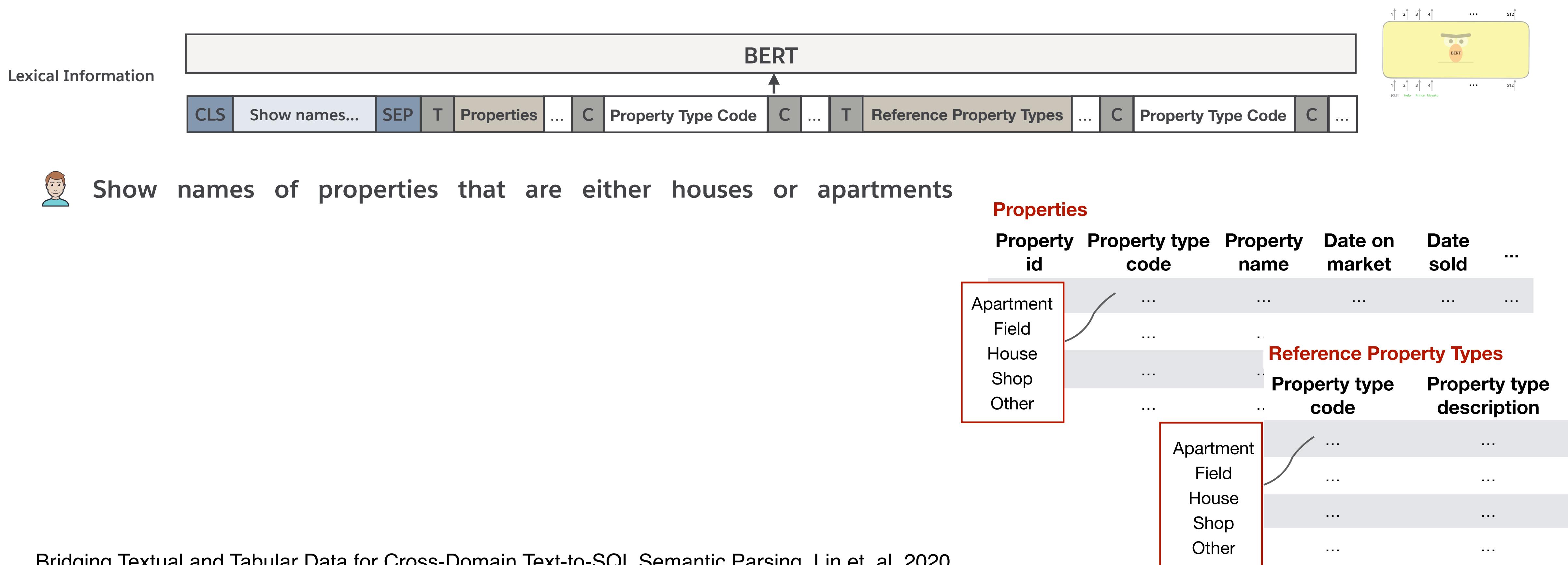
Textual-Tabular Data Encoding

Serialize Table Header/DB Schema



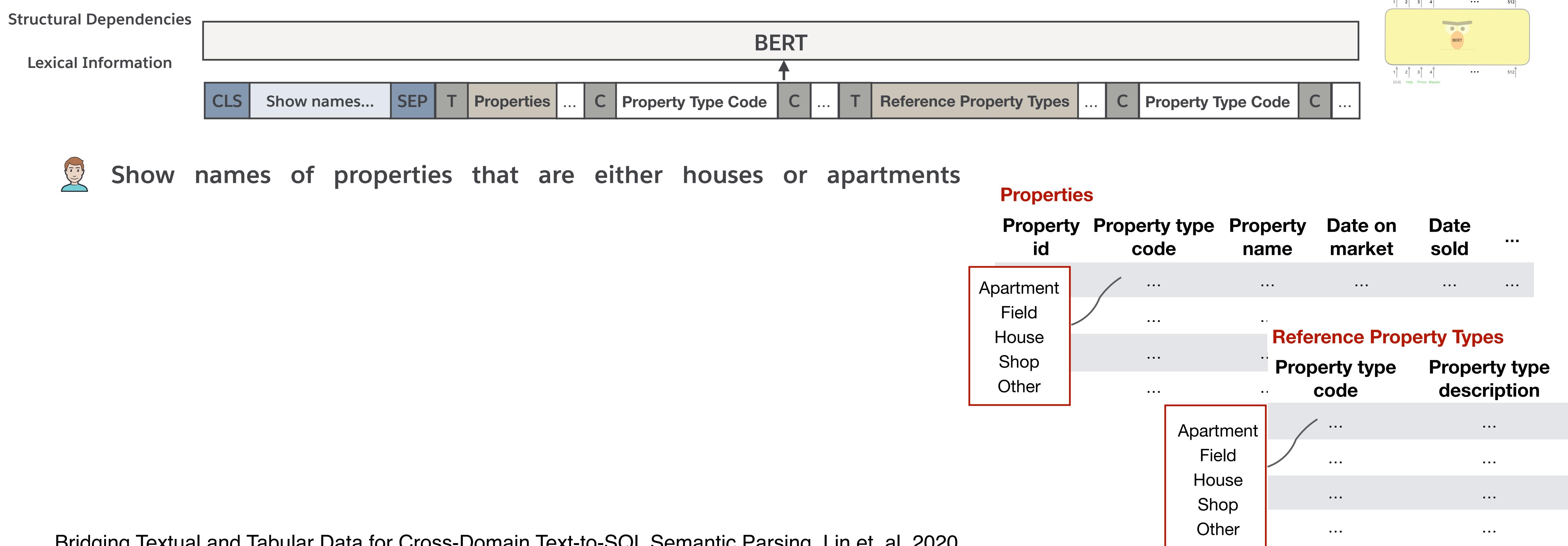
Textual-Tabular Data Encoding

Serialize Table Header/DB Schema

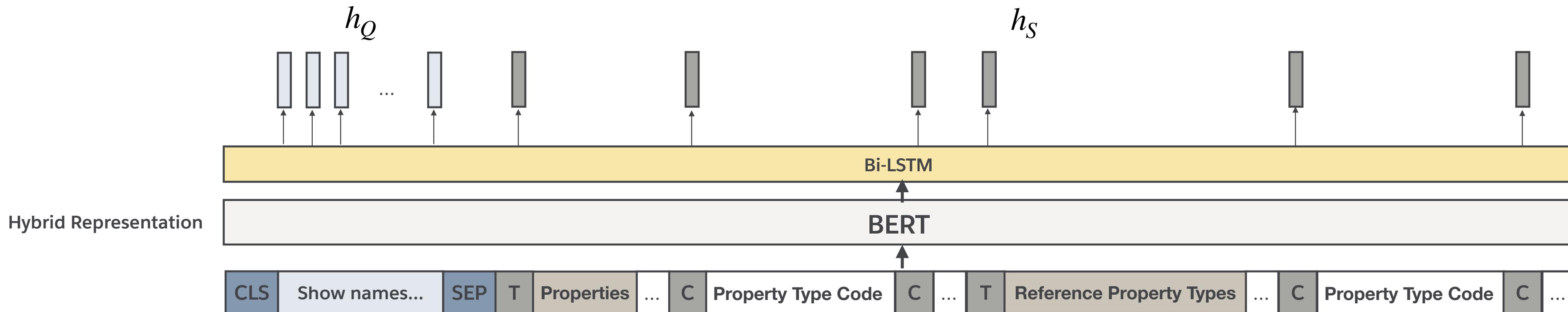


Textual-Tabular Data Encoding

Serialize Table Header/DB Schema



Textual-Tabular Data Encoding



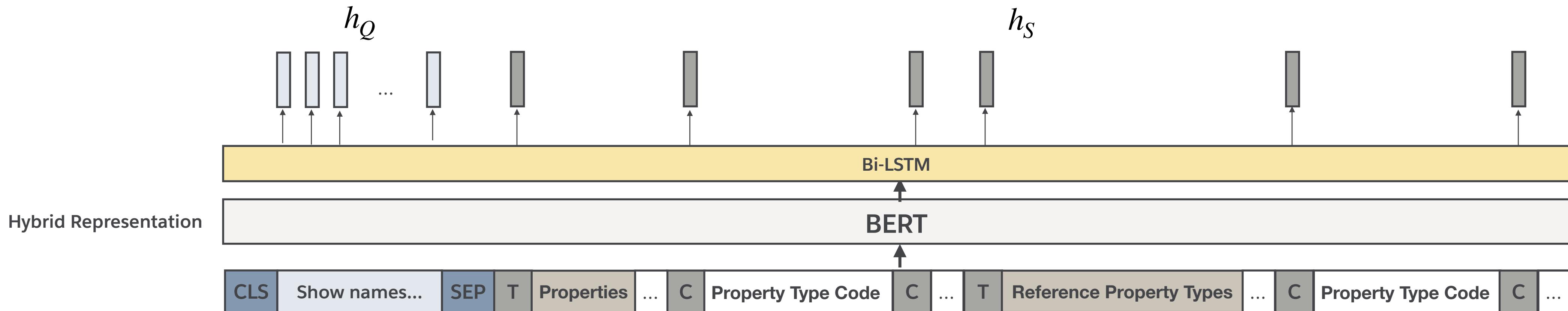
Show names of properties that are either houses or apartments

Properties

Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	Property type code	Property type description
Apartment
Field
House
Shop
Other

Bridging



Show names of properties that are either houses or apartments

Fuzzy String Match

Properties

Property id	Property type code	Property name	Date on market	Date sold
-------------	--------------------	---------------	----------------	-----------

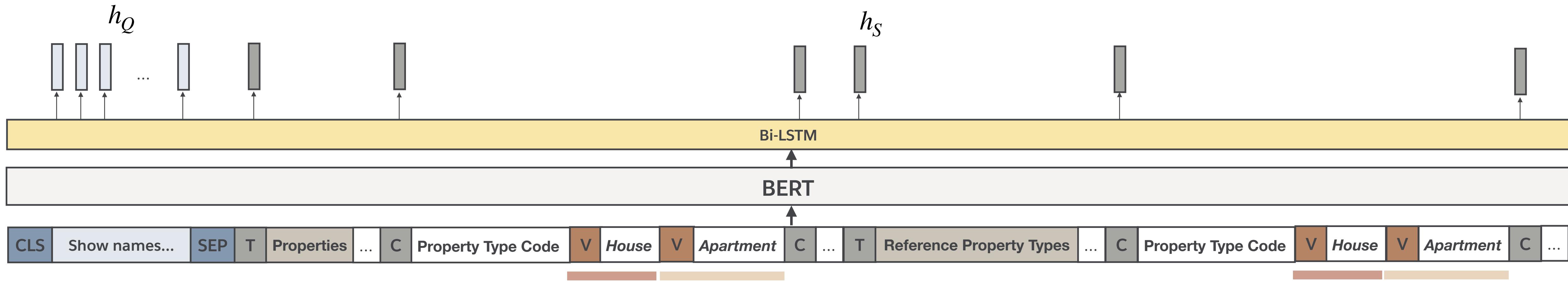
Apartment
Field
House
Shop
Other

Reference Property Types

Property type code	Property type description
--------------------	---------------------------

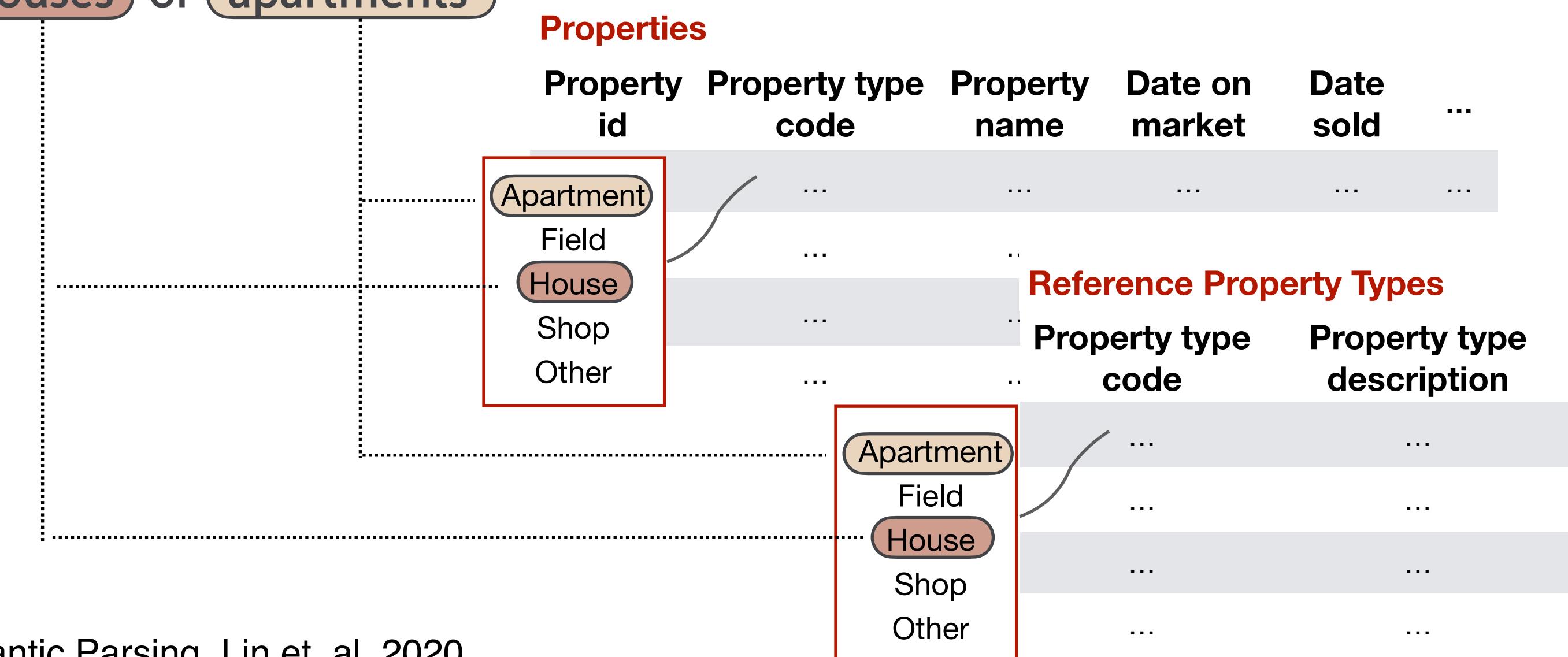
Apartment
Field
House
Shop
Other

Bridging

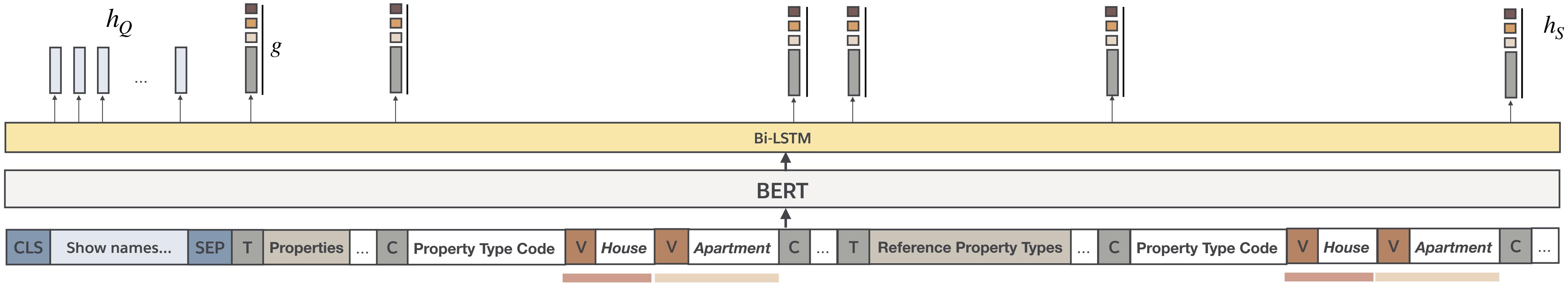


 Show names of properties that are either **houses** or **apartments**

Fuzzy String Match

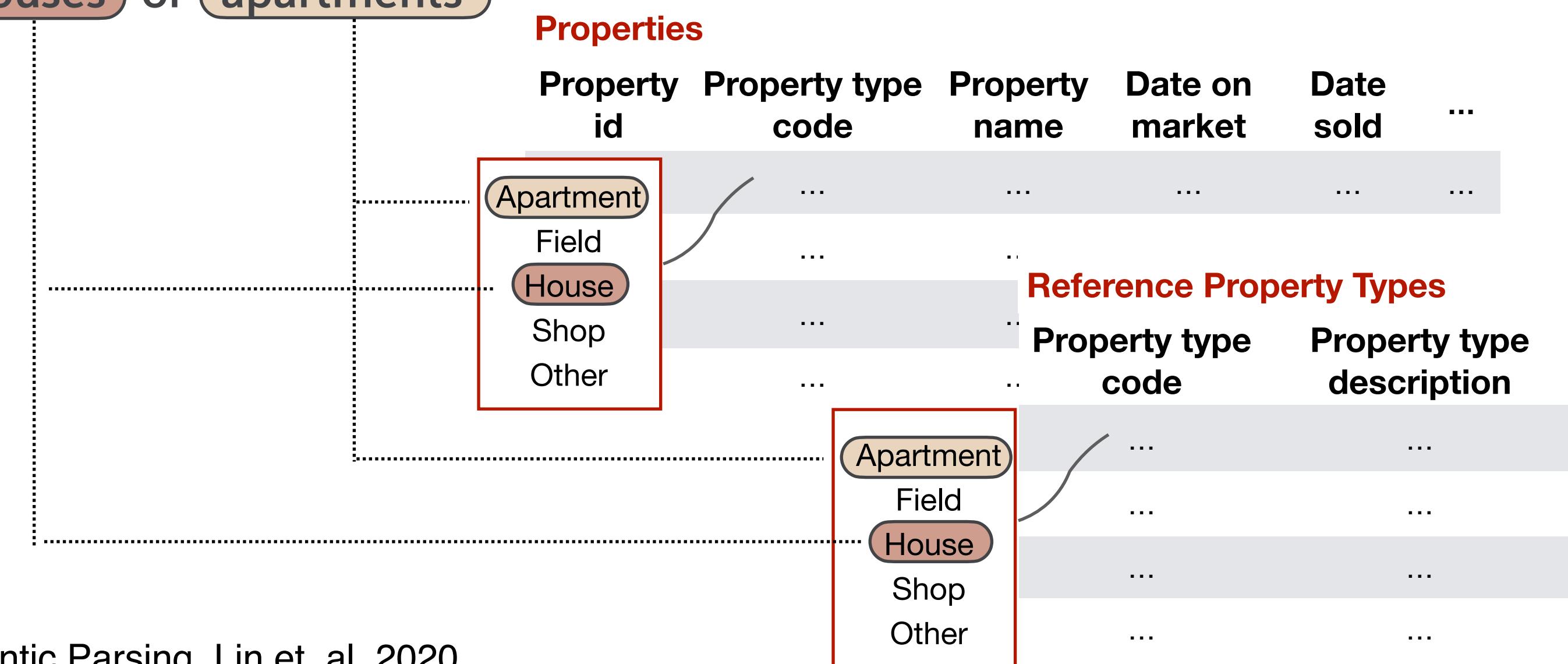


Meta-Data Encoding

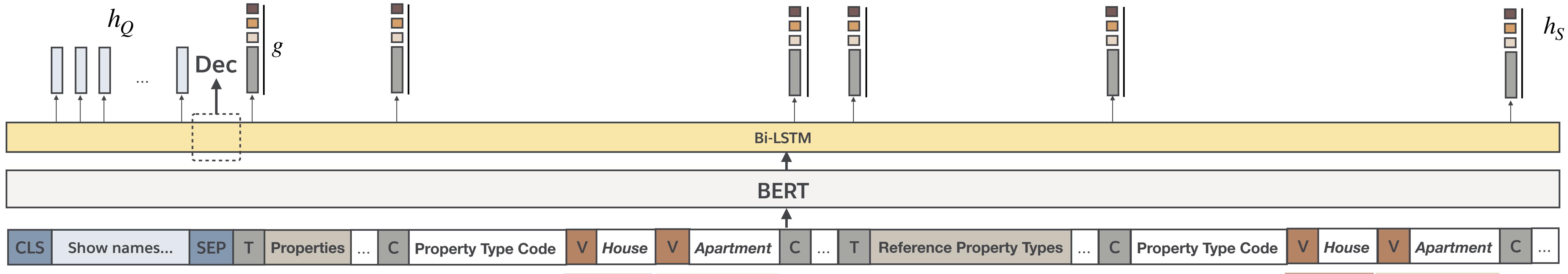


Show names of properties that are either houses or apartments

Fuzzy String Match

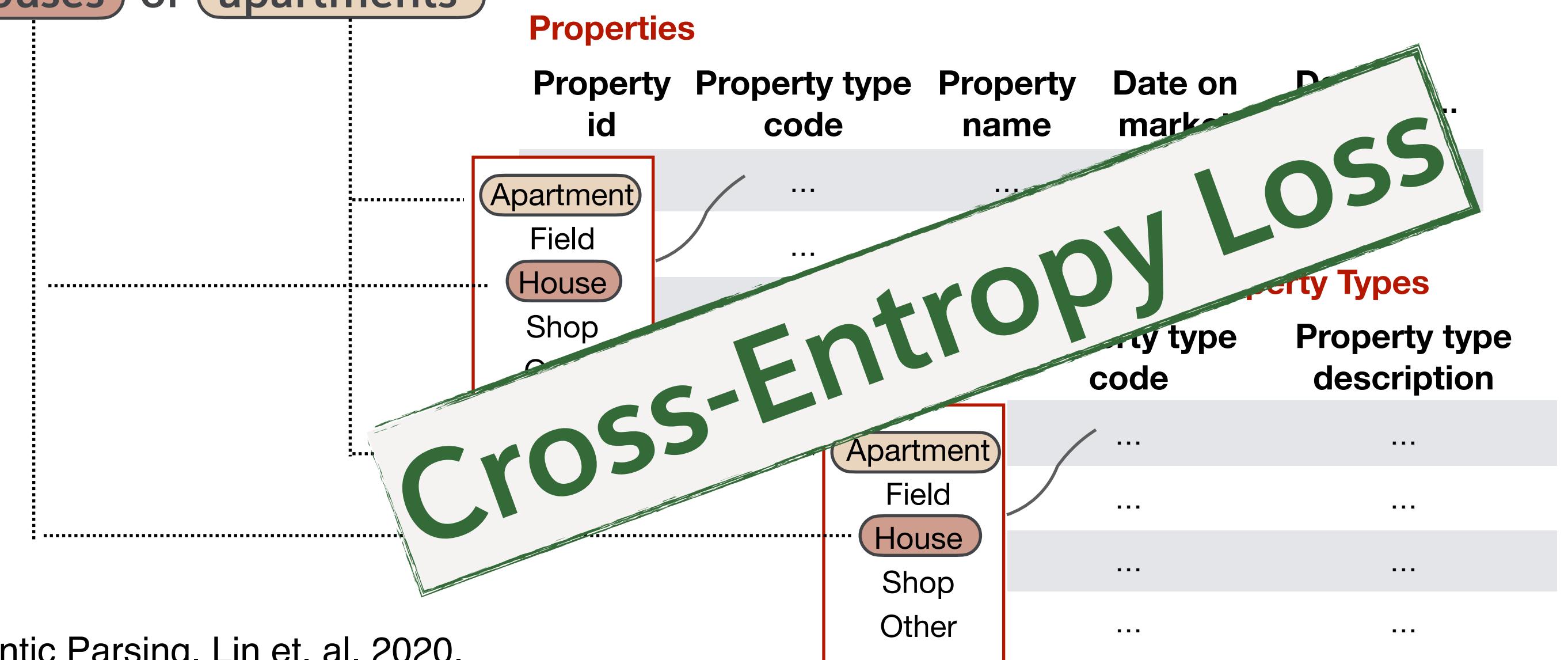
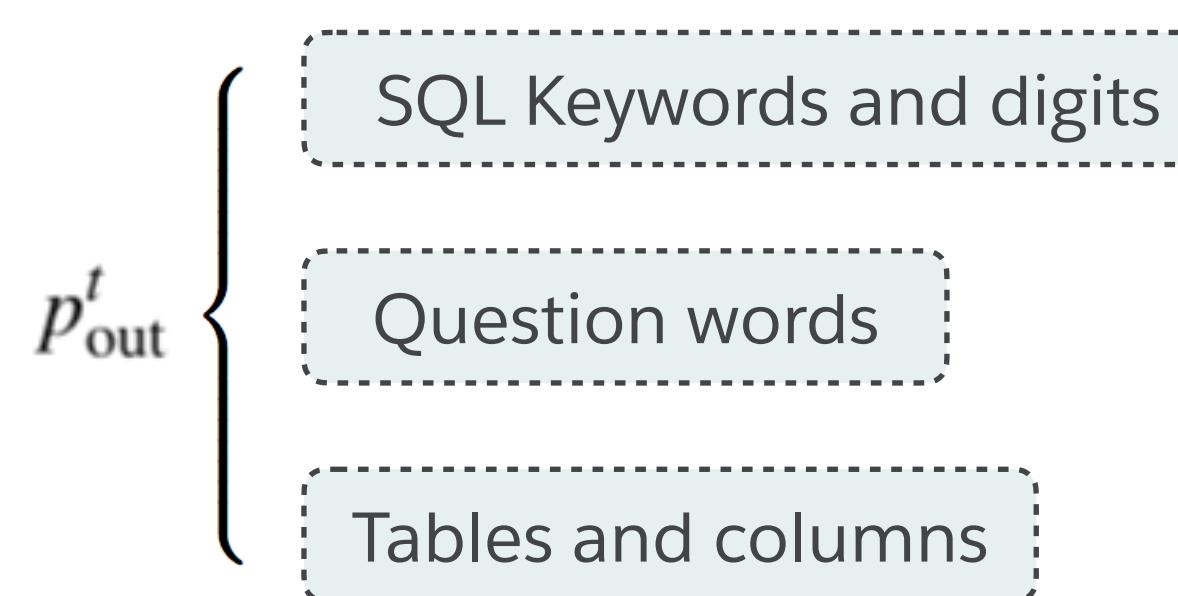


Decoder



Show names of properties that are either houses or apartments

LSTM-based pointer-generator (See et al. 2017)



Main Dataset

Spider (Yu et al. 2018)

Expert-annotated, cross-domain, complex
text-to-SQL dataset

No overlap between Train/Dev/Test
databases, enabling the development of text-
to-SQL models which generalize to unseen DBs

	Train	Dev	Test
# DBs	146	20	40
# Examples	8,659	1,034	2,147

Hidden

Database

Instructor

Primary key

ID	Name	Department_ID	Salary	...
----	------	---------------	--------	-----

Department

Foreign key

ID	Name	Building	Budget	...
----	------	----------	--------	-----

...

...

Question What are the name and budget of the departments
with average instructor salary above the overall average?

SQL

```

SELECT T2.name, T2.budget
FROM Instructor AS T1 JOIN Department AS T2 ON
T1.Department_ID = T2.ID
GROUP BY T1.Department_ID
HAVING AVG(T1.salary) >
(SELECT AVG(Salary) FROM Instructor)

```

Experiments

Pre-processing

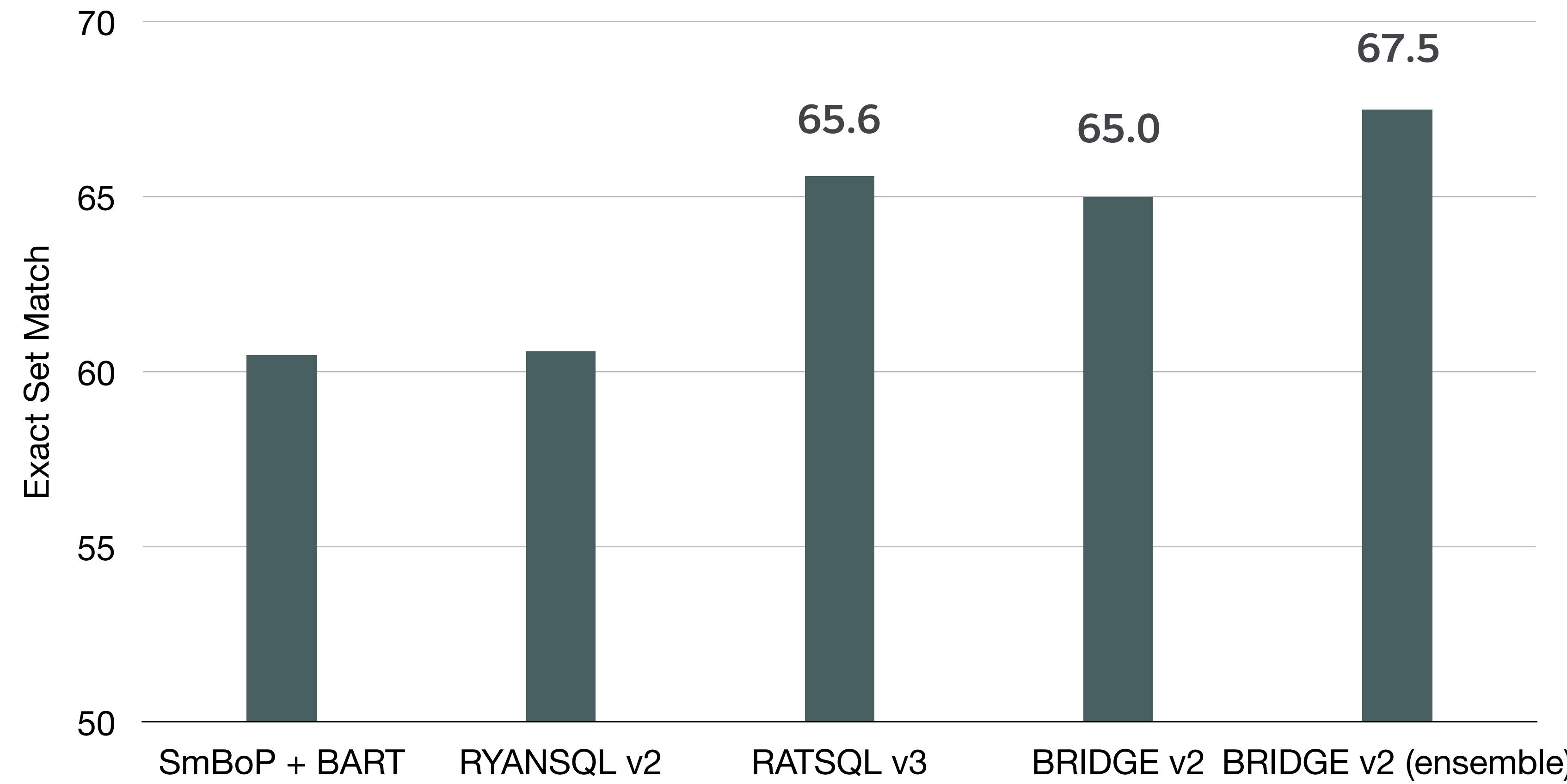
- Compute fuzzy string match between the input question and the picklists of each DB field to identify value mentions
- For each DB field, keep top-K matches and use them to augment the DB schema representation

Evaluation

- Exact set match
 - Logical form match ignoring values and SQL component order invariance
- Execution accuracy
 - Check if the execution results of the predicted SQL query matches the executions results of the ground-truth SQL query

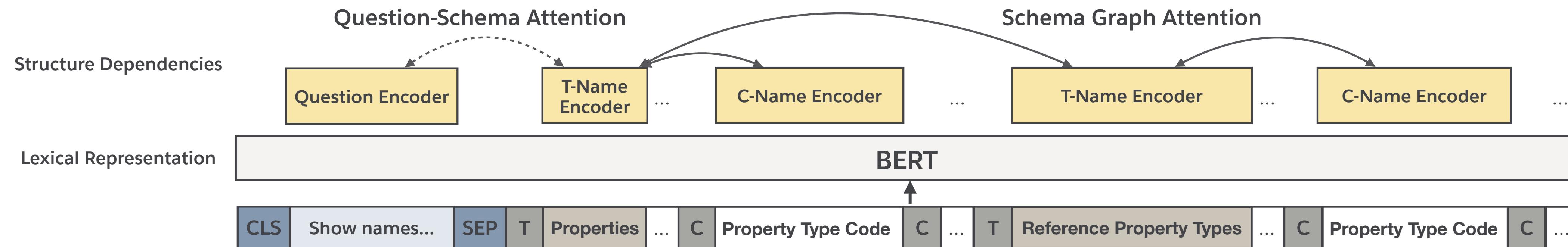
 Better evaluation for text-to-SQL is still an open research problem

Performance on Spider Leaderboard



Comparison to other top-performing text-to-SQL models on the Spider leaderboard as of Dec 01, 2020,

RAT-SQL Encoder



Show names of properties that are either houses or apartments

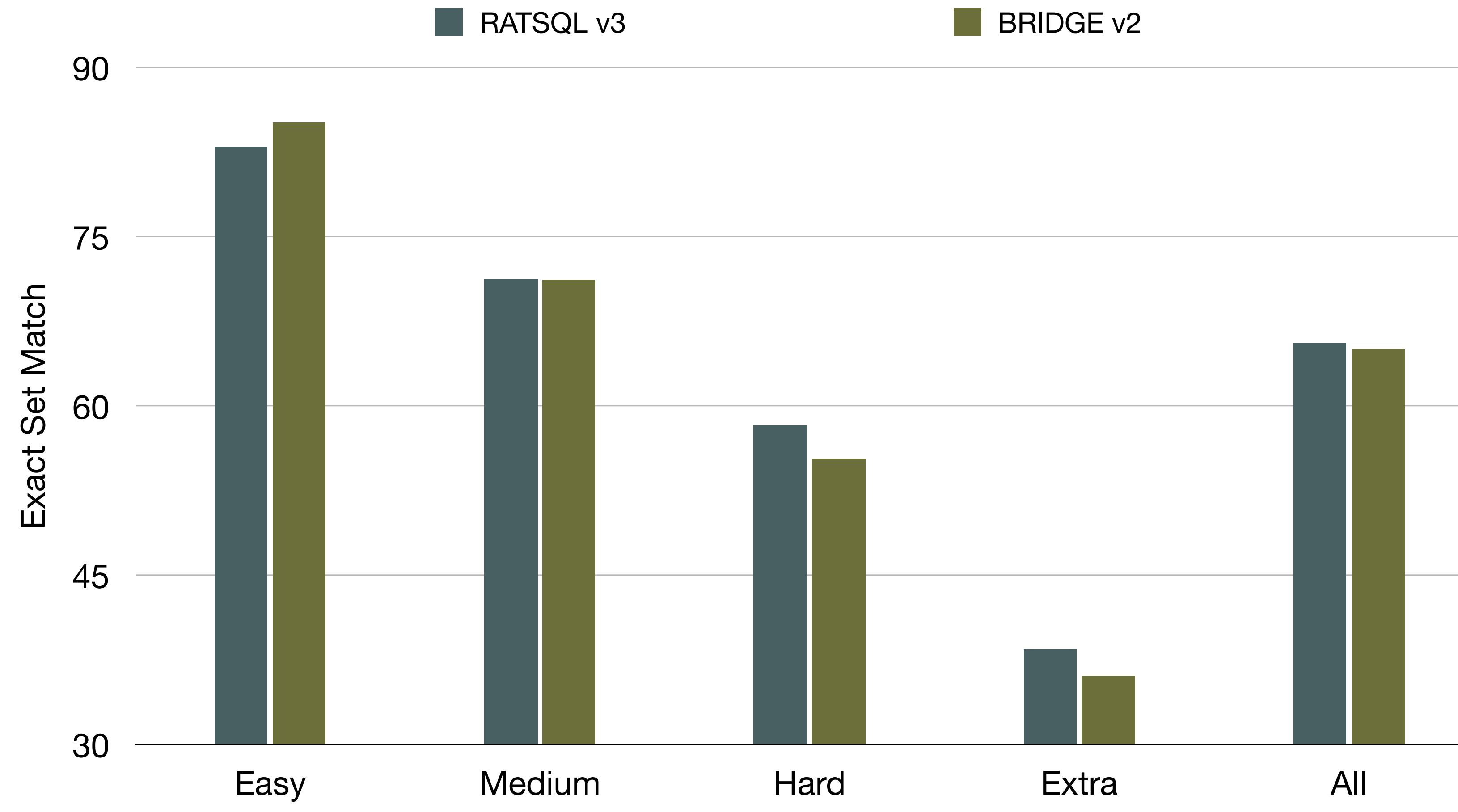
Properties

Property id	Property type code	Property name	Date on market	Date sold	...
Apartment
Field
House
Shop
Other

Reference Property Types	Property type code	Property type description
Apartment
Field
House
Shop
Other

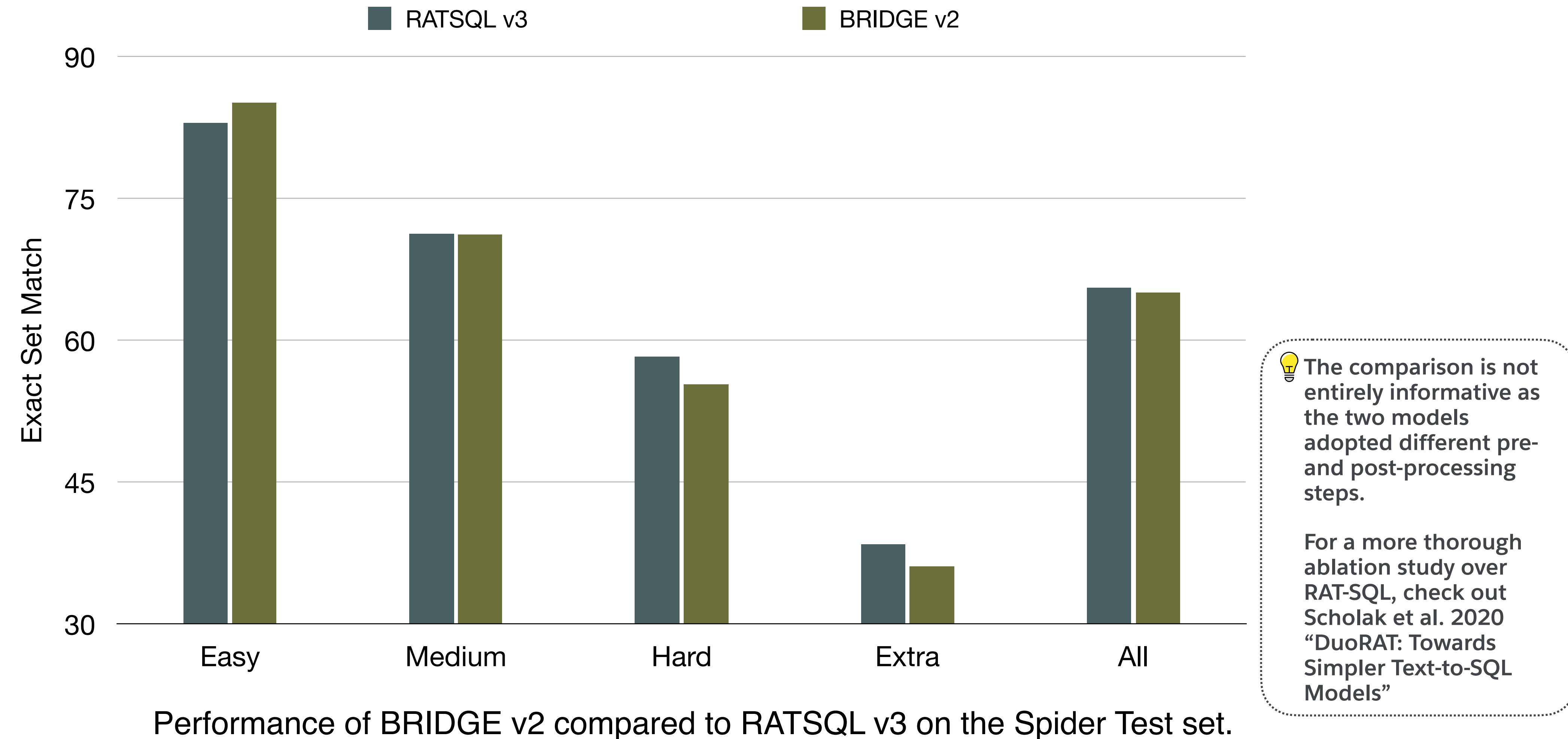


Performance by Difficulty Level



Performance of BRIDGE v2 compared to RATSQ v3 on the Spider Test set.

Performance by Difficulty Level



Ablation Study

Model	Exact Set Match (%)	
	Mean	Max
BRIDGE ($k = 2$)	65.8 ± 0.8	66.9
- SC-guided decoding	65.4 ± 0.7	66.3 (-0.6)
- static SQL check	64.8 ± 0.9	65.9 (-1.0)
- execution order	64.2 ± 0.1	64.3 (-2.6)
- - - - - table shuffle & drop	63.9 ± 0.3	64.3 (-2.6)
- anchor text	63.3 ± 0.6	63.9 (-3.0)
- BERT	17.7 ± 0.7	18.3 (-48.6)

Ablation study of BRIDGE model (w/ BERT-base-uncased) on the Spider Dev set.

Performance on WikiSQL Leaderboard



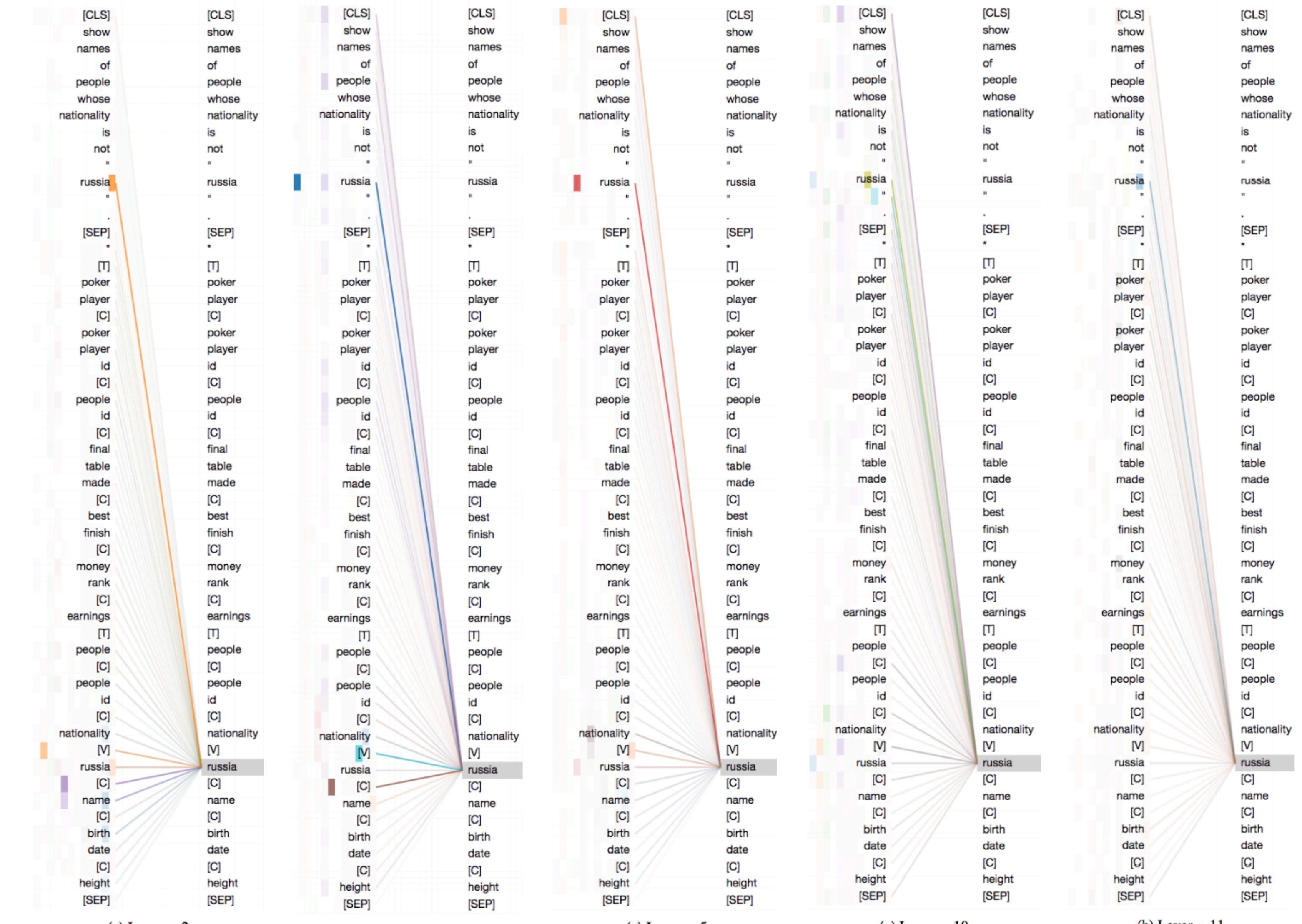
Model	Dev		Test	
	EM	EX	EM	EX
SQLova (Hwang et al., 2019)	81.6	87.2	80.7	86.2
X-SQL (He et al., 2019b)	83.8	89.5	83.3	88.7
HydraNet (Lyu et al., 2020)	83.6	89.1	83.8	89.2
BRIDGE +B _L ($k = 2$) ♠	85.1	91.1	84.8	90.4
SQLova+EG (Hwang et al., 2019)	84.2	90.2	83.6	89.6
BRIDGE +B _L ($k = 2$)+EG ♠	86.1	92.5	85.8	91.7
X-SQL+EG (He et al., 2019b)	86.2	92.3	86.0	91.8
HydraNet+EG (Lyu et al., 2020)	86.6	92.4	86.5	92.2

Comparison between BRIDGE and other top-performing models on the WikiSQL (Zhong et al. 2017) leaderboard as of August 20, 2020. +EG denotes approaches using execution guided decoding.

Attention Visualization after Fine-tuning

BertViz (Vig 2019)

Bridging

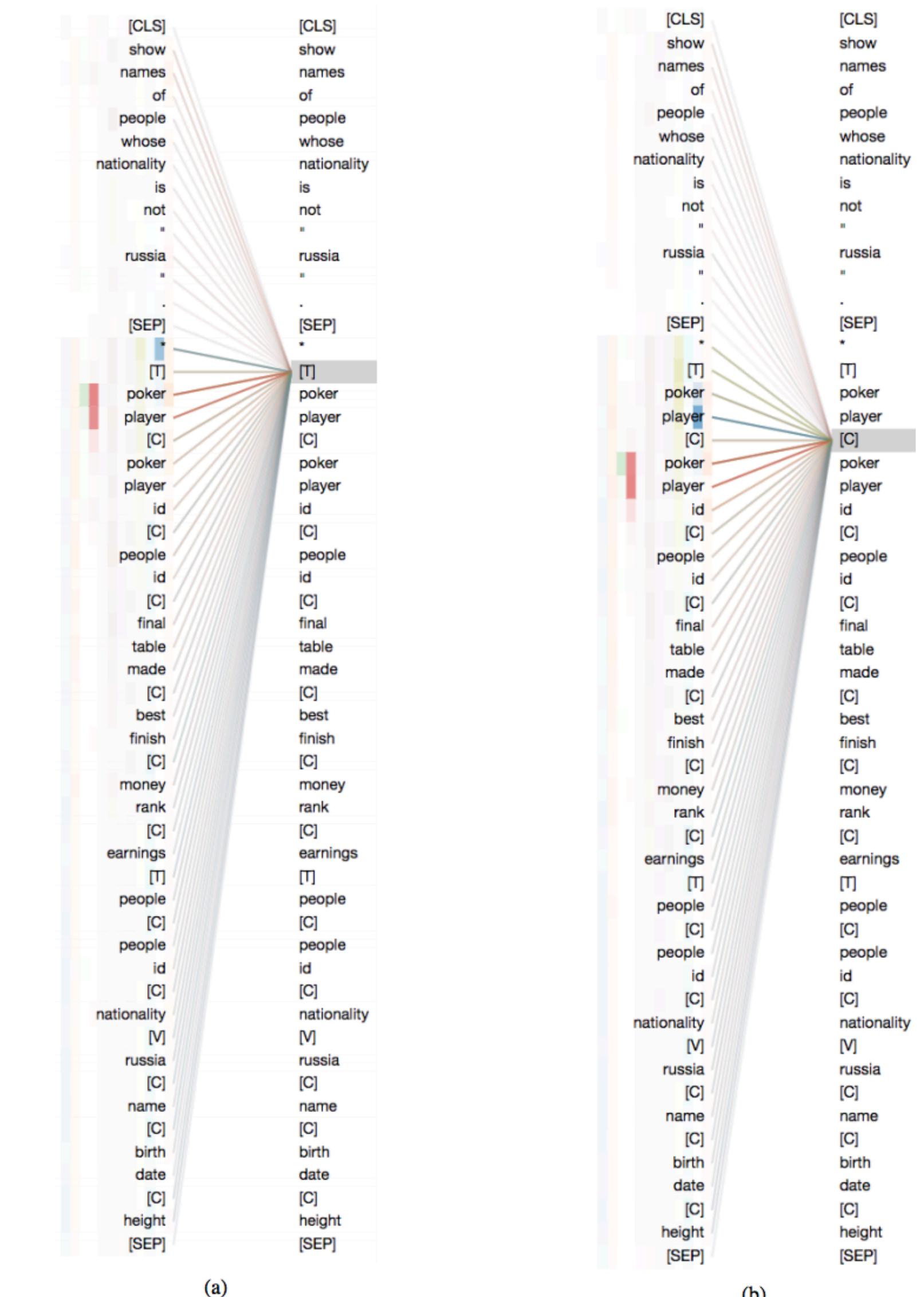


Attention Visualization after Fine-tuning



BertViz (Vig 2019)

- Pooling effect in special tokens [T] and [C], layer 1



Photon Live Demo <https://naturalsql.com>



Photon Select Database concert_singer Upload Database About

ERD Content

concert

concert_ID	concert_Name	Theme	Stadium_ID	Year
1	Auditions	Free choice	1	2014
2	Super bootcamp	Free choice 2	2	2014
3	Home Visits	Bleeding Love	2	2015
4	Week 1	Wide Awake	10	2014
5	Week 1	Happy Tonight	9	2015

Chat started by Photon · 2:24 pm

stadium

Stadium_ID	Location	Name	Capacity	Highest	Lowest	Average
1	Raith Rovers	Stark's Park	10104	4812	1294	2106
2	Ayr United	Somerset Park	11998	2363	1057	1477
3	East Fife	Bayview Stadium	2000	1980	533	864
4	Queen's Park	Hampden Park	52500	1763	466	730
5	Stirling Albion	Forthbank Stadium	3808	1125	404	642

singer

Singer_ID	Name	Country	Song_Name	Song_release_year	Age	Is_male

Query Result

Hello! Please input your question in NL to query the DB

Type Here ➤

© Copyright 2020 Salesforce.com, inc. All rights reserved.

Photon Live Demo V2.0 Is Coming!



Select Upload

sqlite csv

document_template_manage...	tvshow	employee_hire_evaluation	customers_and_addresses
orchestra	course_teach	concert_singer	museum_visit
sales_records	employee	allergy_1	Superstore_tableau
sports_salaries	product_catalog	worker-copy	movie_industry
chinook	farm	academic	fuel_consumption

<< < 1 2 > >>

Close

Attached to Date_management

Hello! Please input your question in NL or SQL to query the DB

Type Here

Date_Effective_From	Date_Effective_To	Template_Details	Ref_Template_Types	Template_Type_Code	Template_Type_Description	Paragraphs
0	5	PP	4	BK	PPT	5
1	9	PP	6	PPT	PPT	5
4	4	BK	7	2	8	5
6	2	PPT	7	1975-05-20 22:51:19	1993-10-07 02:33:04	5
7	8	PPT	8	2002-03-02 14:39:49	1992-05-02 20:06:11	5
				2005-11-12 07:09:48	2008-01-05 14:19:28	5
				1999-07-08 03:31:04	1975-07-16 04:52:10	5
				2001-04-18 09:29:52		5

<< < 1 2 3 > >>

Recap

Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

Recap

Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

1. The BRIDGE model has significantly more parameters in the encoder than the decoder. Learning good column representations is critical.
2. Transformer language models (e.g. BERT) are strong pattern learners.
3. The cross-database setup adopted by the standard Spider dataset does not systematically separate generalization over compositionally vs. natural language variations. The models need to be benchmarked in more systematic ways following recent works by Suhr et al. 2020 and Shaw et al. 2020.

Recap

Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

Despite competitive performance on the standard benchmark datasets, we saw several issues with the current approach:

Recap

Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

Despite competitive performance on the standard benchmark datasets, we saw several issues with the current approach:

I. Robustness

Model performance is sensitive to random seeds. We plan to benchmark the model on more datasets using more evaluation metrics. (Suhr et al. 2020; Scholak et al. 2020; Shaw et al. 2020)

Recap

Research Question:

Can we build a model which seeks information from structured relational databases in a similar way to seeking information from textual paragraphs?

Despite competitive performance on the standard benchmark datasets, we saw several issues with the current approach:

I. Robustness

Model performance is sensitive to random seeds. We plan to benchmark the model on more datasets using more evaluation metrics. (Suhr et al. 2020; Scholak et al. 2020; Shaw et al. 2020)

2. Interpretability and Control

Model lacks interpretability and control to the output. Introducing the relation-aware attention layers (RAT-SQL, Wang et al. 2020) improves interpretability to some degree.

Related Work

To our knowledge, the following works have also demonstrated the effectiveness of pre-trained transformer LMs on the text-to-SQL problem:

1. A Comprehensive Exploration on WikiSQL with Table-Aware Word Contextualization. Hwang et al. 2019
2. Exploring Unexplored Generalization Challenges for Cross-Database Semantic Parsing. Suhr et al. 2020
3. Compositional Generalization and Natural Language Variation: Can a Semantic Parsing Approach Handle Both? Shaw et al. 2020
4. Compositional Generalization in Semantic Parsing: Pre-training vs. Specialized Architectures. Furrer et al. 2020

Related Work

Pre-training textual-tabular representations leveraging unlabeled tables, to name a few:

1. GraPPa: Grammar-Augmented Pre-training for Table Semantic Parsing. Yu et al. 2020
2. TaBERT: A New Model for Understanding Queries over Tabular Data. Yin et al. 2020
3. TAPAs: Weakly Supervised Table Parsing via Pre-training. Herzog et al. 2020
4. Structure-Grounded Pre-training for Text-to-SQL. Deng et al. 2020
5. TURL: Table Understanding through Representation Learning. Deng et al. 2021

Related Work

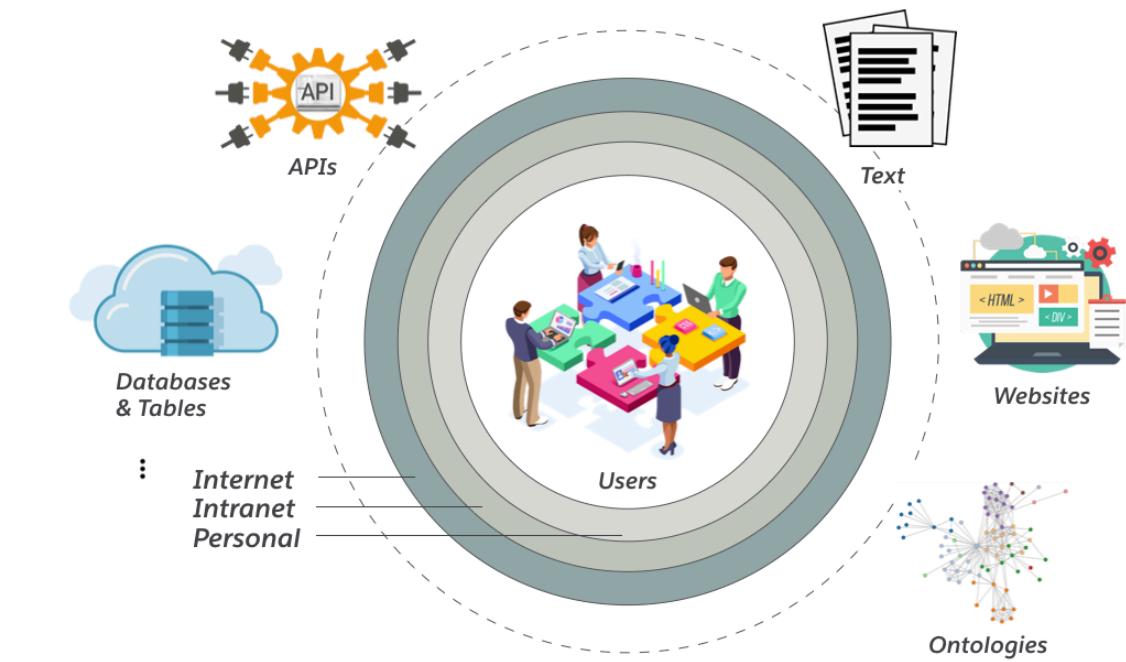
Pre-training textual-tabular representations leveraging unlabeled tables, to name a few:

1. GraPPa: Grammar-Augmented Pre-training for Table Semantic Parsing. Yu et al. 2020
2. TaBERT: A New Model for Understanding Queries over Tabular Data. Yin et al. 2020
3. TAPAs: Weakly Supervised Table Parsing via Pre-training. Herzog et al. 2020
4. Structure-Grounded Pre-training for Text-to-SQL. Deng et al. 2020
5. TURL: Table Understanding through Representation Learning. Deng et al. 2021

Future Work:

Can we learn table representations that universally benefit tasks s.a. semantic parsing, summarization, data completion, etc.?

Can we learn table representations that can be “plugged” into the representations of other data modality?



<https://einstein.ai/research/publications>



Tao Yu
Yale University



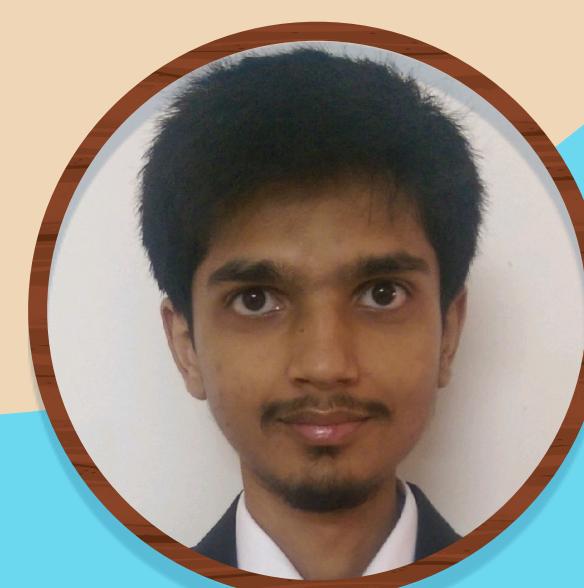
Chien-Sheng Wu
Salesforce AI
Research



Bailin Wang
University of
Edinburgh



Jichuan Zeng
Chinese University
of Hong Kong



Karthik
Radhakrishnan
CMU



Heather Lent
University of
Copenhagen



Arvind Srikantan
Salesforce Inc.



Dragomir Radev
Yale University



Richard Socher
You.com



Steven Hoi
Salesforce
Research Asia



Yingbo Zhou
Salesforce AI
Research



Caiming Xiong
Salesforce AI
Research



thank
you



Error Analysis

Qualitative observations



What are the names and release years for all the songs of the youngest singer? **concert_singer**

- ✗ SELECT Song_Name, Age FROM singer ORDER BY Age LIMIT 1
- ✓ SELECT song_name, song_release_year FROM singer ORDER BY age LIMIT 1

Robustness issue



What are the full names of all left handed players, in order of birth date? **WTA_1**

- ✗ SELECT first_name, last_name FROM players ORDER BY birth_date
- ✓ SELECT first_name, last_name FROM players WHERE hand = 'L' ORDER BY birth_date

Rare relation & value surface form



What are the names of students who have 2 or more likes? **network_1**

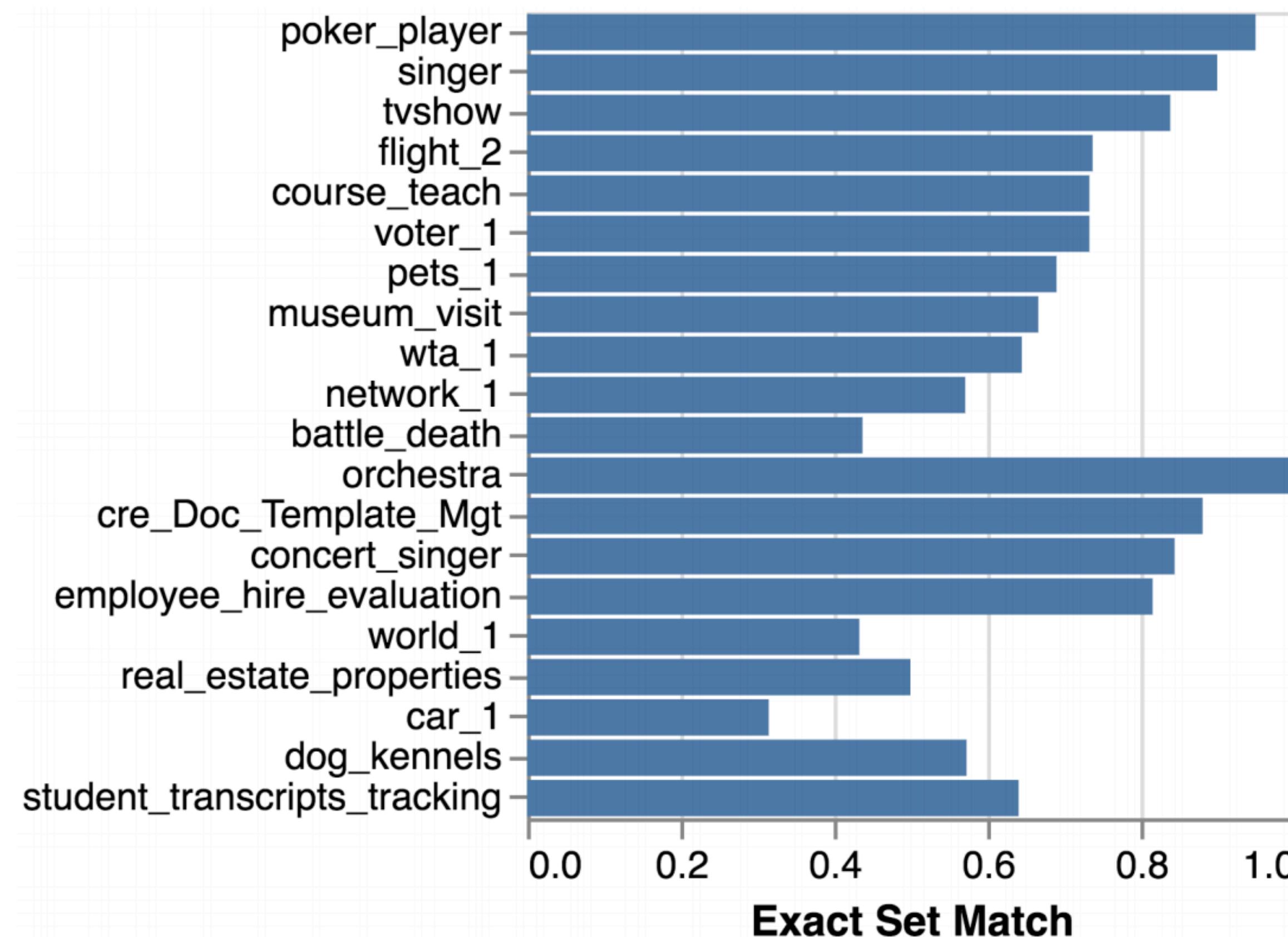
- ✗ SELECT Likes.student_id FROM Likes JOIN Friend ON Likes.student_id = Friend.student_id
GROUP BY Likes.student_id HAVING COUNT(*) >= 2
- ✓ SELECT Highschooler.name FROM Likes JOIN Highschooler ON Likes.student_id =
Highschooler.id GROUP BY Likes.student_id HAVING count(*) >= 2

Commonsense

“Friend” table stores students who has a friend, not all students

Performance by DB

Exact match accuracy on each DB in the Spider dev set. The DBs are sorted by size (smallest -> largest) from top to bottom.



Better characterization of “similar” examples could help transfer learning