# ON THE ACCESSIBILITY AND PRIVACY OF PROVENANCE-BASED EXPLANATIONS

## AMIR GILAD

### THE HEBREW UNIVERSITY

JOINT WORK WITH DANIEL DEUTCH, YUVAL MOSKOVITCH, NAVE FROST, ARIEL FRANKENTHAL
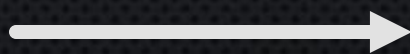
## RAW PROVENANCE CAN BE...

### TOO LONG AND COMPLEX

```
(oname,Duke)·(aname,Jun Y.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle, Scalable...)·(cname,VLDB)·(pyear,06')+
(oname,Duke)·(aname,Jun Y.)· (ptitle, Making...)·(cname,VLDB)·(pyear,07')+
(oname,Duke)·(aname,Brett W.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle,Cumulon...)·(cname,SIGMOD)·(pyear,14')+
…
```

### TOO REVEALING

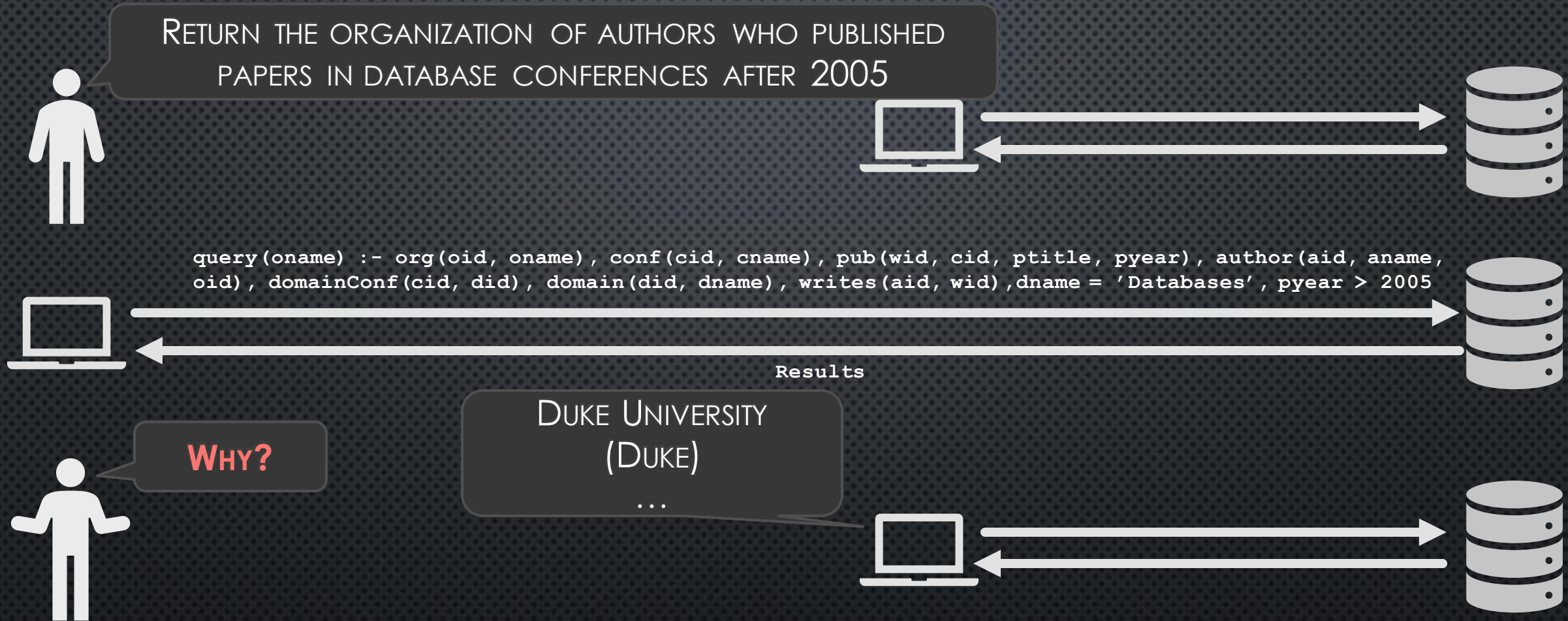**PROVENANCE-BASED EXPLANATIONS** ⟶ **PRIVATE/PROPRIETARY QUERY**

**FACTORIZING AND SUMMARIZING PROVENANCE FOR NATURAL LANGUAGE EXPLANATIONS**

VLDB 16', VLDB 17', SIGMOD REC. 18', VLDB J. 20'

**ABSTRACTING PROVENANCE FOR QUERY PRIVACY**

SIGMOD 21', ICDE 21'

# NATURAL LANGUAGE INTERFACES AND EXPLANATIONS

RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005

```
query(oname) :- org(oid, oname), conf(cid, cname), pub(wid, cid, ptitle, pyear), author(aid, aname,
oid), domainConf(cid, did), domain(did, dname), writes(aid, wid),dname = 'Databases', pyear > 2005
```

Results

DUKE UNIVERSITY
(DUKE)
...

WHY?

DUKE IS THE ORGANIZATION OF 63 AUTHORS WHO PUBLISHED 170 PAPERS IN 31 CONFERENCES IN 2006 - 2015

# PROVENANCE MODEL

**NL QUERY:**

**RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005**

**QUERY:**

```
query(oname) :- org(oid, oname), conf(cid, cname), pub(wid, cid, ptitle,
pyear),author(aid, aname, oid), domainConf(cid, did), domain(did,
dname), writes(aid, wid),dname = 'Databases', pyear > 2005
```

# PROVENANCE MODEL

**NL QUERY:**
RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005

**QUERY:**

```
query(oname) :- org(oid, oname), conf(cid, cname), pub(wid, cid, ptitle,
pyear), author(aid, aname, oid), domainConf(cid, did), domain(did,
dname), writes(aid, wid),dname = 'Databases', pyear > 2005
```

**PROVENANCE OF THE RESULT DUKE:**

```
(oname,Duke)·(aname,Jun Y.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle, Scalable...)·(cname,VLDB)·(pyear,06')+
(oname,Duke)·(aname,Jun Y.)· (ptitle, Making...)·(cname,VLDB)·(pyear,07')+
(oname,Duke)·(aname,Brett W.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle,Cumulon...)·(cname,SIGMOD)·(pyear,14')+
…
```

# PROVENANCE MODEL

**NL QUERY:**

RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005

**QUERY:**

```
query(Duke) :- org(oid, Duke), conf(cid, cname), pub(wid, cid, iCheck...,
2014), author(aid, Jun Y., oid), domainConf(cid, did), domain(did,
SIGMOD), writes(aid, wid),dname = 'Databases', 2014 > 2005
```

**PROVENANCE OF THE RESULT DUKE:**

```
(oname,Duke)·(aname,Jun Y.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle, Scalable...)·(cname,VLDB)·(pyear,06')+
(oname,Duke)·(aname,Jun Y.)· (ptitle, Making...)·(cname,VLDB)·(pyear,07')+
(oname,Duke)·(aname,Brett W.)·(ptitle,iCheck...)·(cname,SIGMOD)·(pyear,14')+
(oname,Duke)·(aname,Jun Y.)·(ptitle,Cumulon...)·(cname,SIGMOD)·(pyear,14')+
…
```
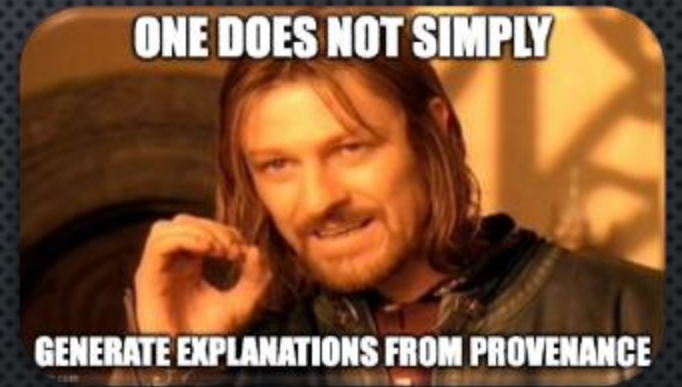
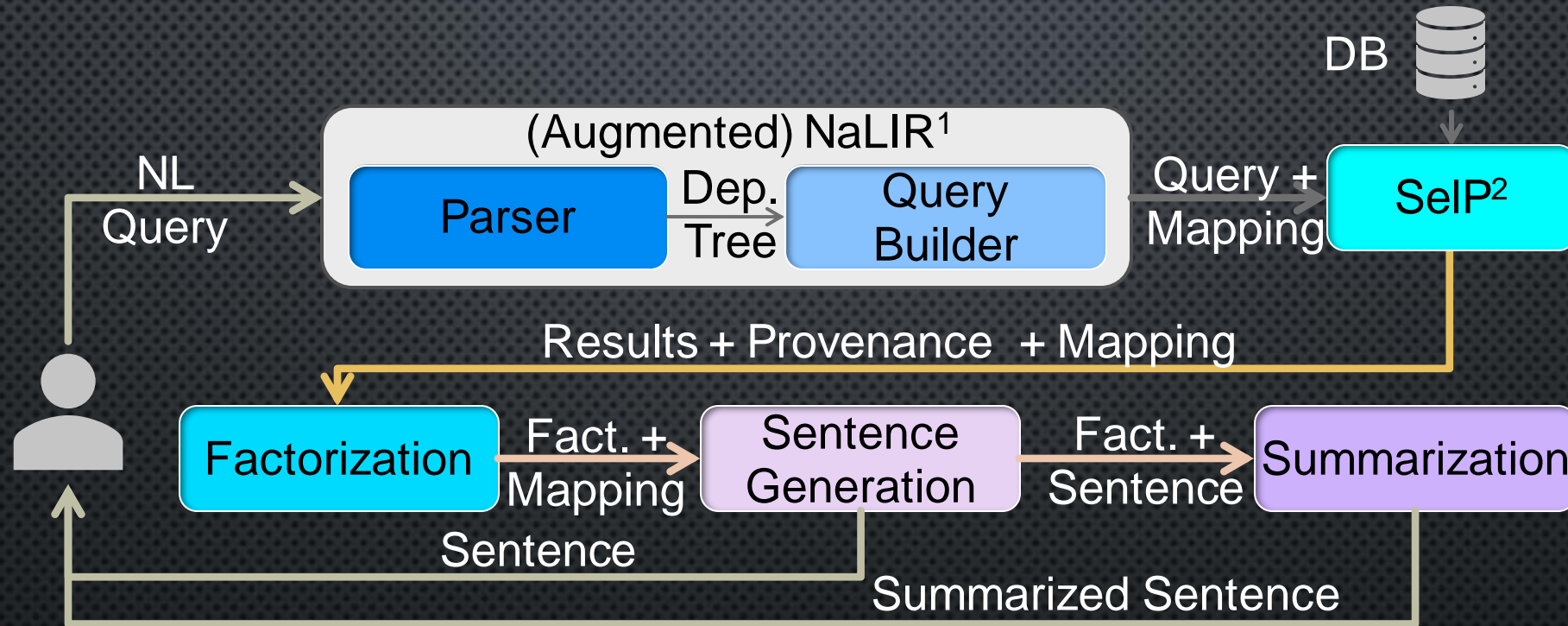## HOW DO WE CONVERT PROVENANCE TO A NATURAL LANGUAGE EXPLANATION?

### CHALLENGES:

1. THE FORMAL PROVENANCE IS FAR FROM AN NL SENTENCE

2. THE PROVENANCE CAN BE VERY LONG AND CONVOLUTED
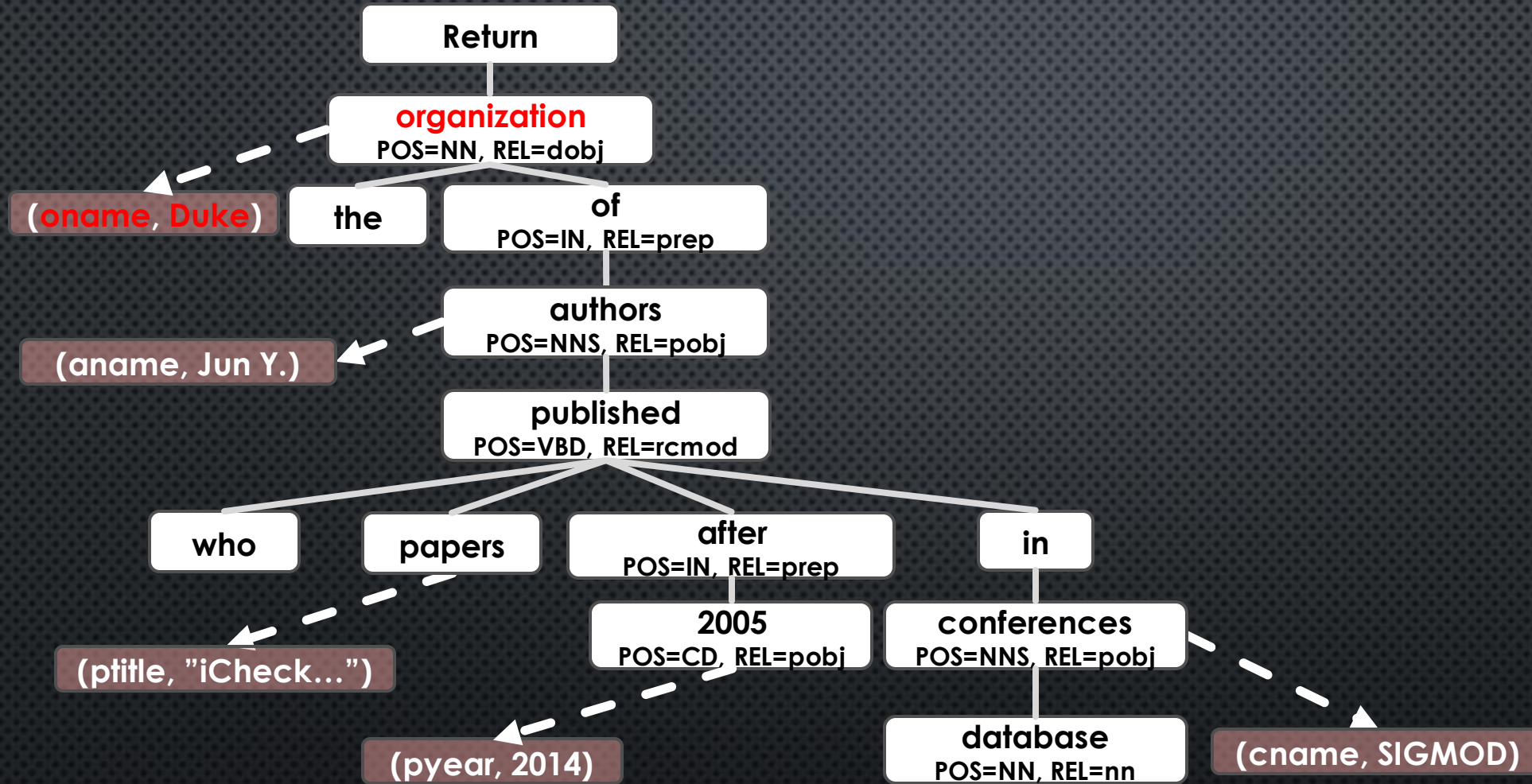


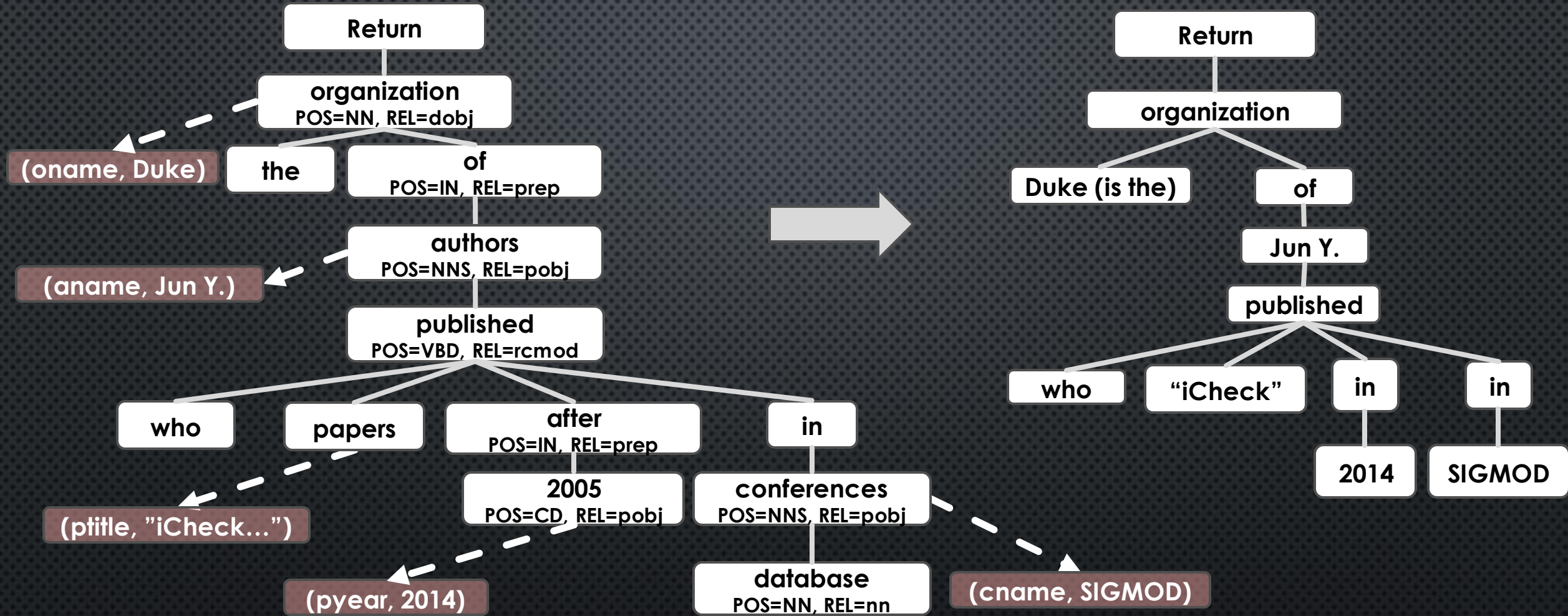## USE THE STRUCTURE OF THE INPUT QUESTION!

1. Li, F., Jagadish, H. V., "Constructing an Interactive Natural Language Interface for Relational Databases". In: Proc. VLDB Endow. (2014), pp. 73–84
2. Deutch, D., G., Moskovitch, Y., "Efficient provenance tracking for datalog using top-k queries". In: VLDB J. 27.2 (2018), pp. 245–269

Return

organization
POS=NN, REL=dobj

(oname, Duke)

the

of
POS=IN, REL=prep

authors
POS=NNS, REL=pobj

(aname, Jun Y.)

published
POS=VBD, REL=rcmod

who

papers

after
POS=IN, REL=prep

in

(ptitle, "iCheck…")

2005
POS=CD, REL=pobj

conferences
POS=NNS, REL=pobj

(pyear, 2014)

database
POS=NN, REL=nn

(cname, SIGMOD)

Duke is the organization of Jun Y. who published 'iCheck...' in SIGMOD in 2014

# PROVENANCE FACTORIZATION

**IDEA**: USE ALGEBRAIC FACTORIZATION TO TAKE OUT COMMON VALUES THAT APPEAR IN MULTIPLE ASSIGNMENTS

```
[Duke]·[Jun Y.]·[iCheck...]·[SIGMOD]·[2014]+
[Duke]·[Jun Y.]·[Scalable...]·[VLDB]·[2006]+
[Duke]·[Jun Y.]·[Making..]·[VLDB]·[2007]+
[Duke]·[Brett W.]·[iCheck...]·[SIGMOD]·[2014]+
[Duke]·[Jun Y.]·[Cumulon...]·[SIGMOD]·[2014]
```

**INTUITION**: WE WANT A FACTORIZATION THAT **FOLLOWS THE STRUCTURE OF THE NL QUERY** TO BE ABLE TO GENERATE A SENTENCE

**SHORTEST FA**

```
[Duke]·
 ([SIGMOD]·[2014]·
 ([iCheck...]·
 ([Jun Y.] + [Brett W.]))
  + [Jun Y.]·[Cumulon...])
  + [VLDB]·[Jun Y.]·
 ([2006]·[Scalable...])
  + [2007]·[Making...])
```

```
[Duke]·
 ([Jun Y.]·
  ([VLDB]·
  ([2006]·[Scalable...]
   + [2007]·[Making...]))
   + [SIGMOD]·[2014]·
  ([iCheck...] + [Cumulon...]))
 + [Brett W.]·[iCheck...]·[SIGMOD]·[2014])
```

# T-COMPATIBILITY

**NL QUERY:**

RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005

**SHORTEST FACTORIZATION:**

```
[Duke]·
 ([SIGMOD]·[2014]·
 ([iCheck...]·
 ([Jun Y.] + [Brett W.]))
  + [Jun Y.]·[Cumulon...])
  + [VLDB]·[Jun Y.]·
 ([2006]·[Scalable...])
  + [2007]·[Making...])
```

**AS A SENTENCE:**

```
Duke is the organization of authors who published in
SIGMOD 2014
'iCheck...' which was published by
Jun Y. and Brett W.
and Jun Y. published 'Cumulon...'
and Jun Y. published in VLDB
'Scalable...' in 2014
and 'Making...' in 2007.
```

[Duke]·
 ([SIGMOD]·[2014]·
 ([iCheck...]·
 ([Jun Y.] + [Brett W.]))
  + [Jun Y.]·[Cumulon...])
  + [VLDB]·[Jun Y.]·
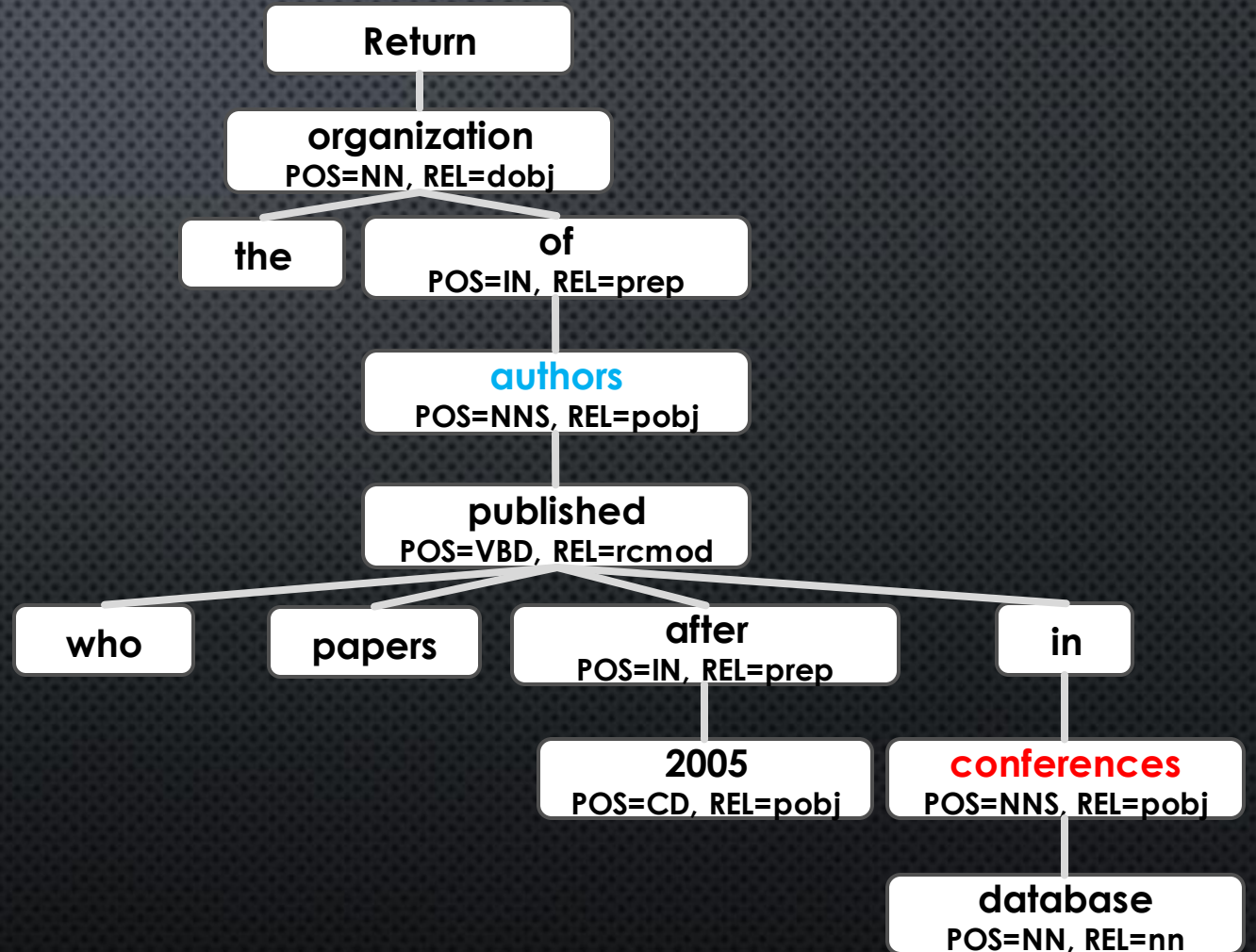 ([2006]·[Scalable...])
  + [2007]·[Making...])

Return

organization
POS=NN, REL=dobj

the

of
POS=IN, REL=prep

authors
POS=NNS, REL=pobj

published
POS=VBD, REL=rcmod

who

papers

after
POS=IN, REL=prep

in

2005
POS=CD, REL=pobj

conferences
POS=NNS, REL=pobj

database
POS=NN, REL=nn

**NL Query:**

Return the organization of authors who published papers in database conferences after 2005

**Longer Factorization:**

```
[Duke]·
([Jun Y.]·
  ([VLDB]·
  ([2006]·[Scalable...]
   + [2007]·[Making...]))
   + [SIGMOD]·[2014]·
  ([iCheck...] + [Cumulon...]))
+ [Brett W.]·[iCheck...]·[SIGMOD]·[2014])
```

**As a Sentence:**

```
Duke is the organization of
 Jun Y. who published
 in VLDB
 'Scalable...' in 2006 and
 'Making...' in 2007
 and in SIGMOD in 2014
 'iCheck...' and 'Cumulon...'
 and Brett W. who published
 'iCheck...' in SIGMOD in 2014.
```
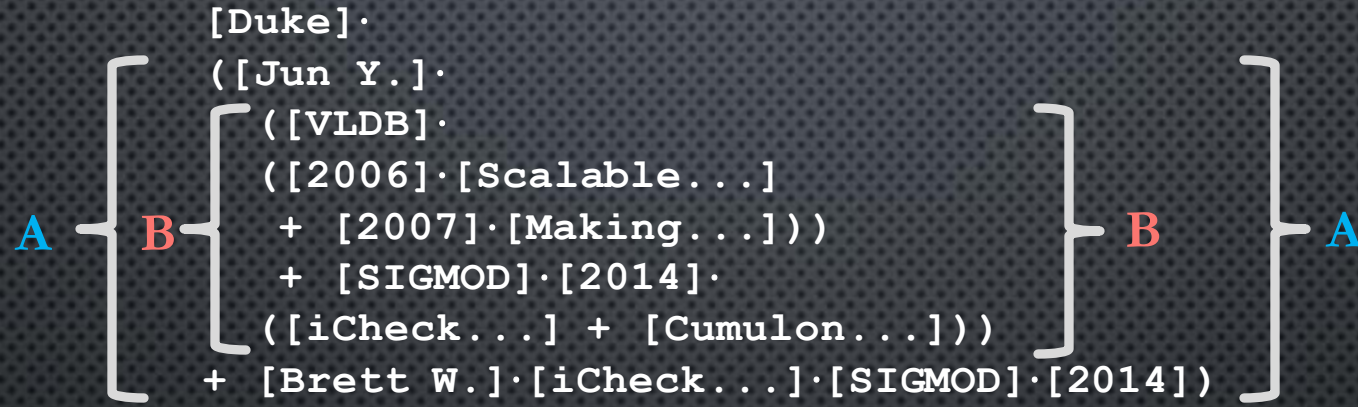
# FINDING T-COMPATIBLE FACTORIZATIONS

## ALGORITHM:

- TRAVERSE THE DEPENDENCY TREE LEVEL-BY-LEVEL

- FOR EVERY LEVEL WITH MAPPED WORDS, FACTORIZE THEIR CORRESPONDING VALUES IN THE PROVENANCE

- PRIORITIZE WHICH VALUES TO TAKE OUT AT EACH LEVEL BY FREQUENCY

**GUARANTEE (INFORMAL):** THE ALGORITHM GENERATES A T-COMPATIBLE FACTORIZATION, ENSURING THAT THE FACTORIZATION CAN BE USED TO GENERATE AN NL EXPLANATION.

# SUMMARIZATION

**TWO LEVELS OF SUMMARIZATION:**

```
                    [Duke]·
                    ([Jun Y.]·
                     ([VLDB]·
                     ([2006]·[Scalable...]
  A        B          + [2007]·[Making...]))      B        A
                     + [SIGMOD]·[2014]·
                     ([iCheck...] + [Cumulon...]))
                    + [Brett W.]·[iCheck...]·[SIGMOD]·[2014])
```

**SHORTER SUMMARIZED EXPLANATION BASED ON A:**
DUKE IS THE ORGANIZATION OF 2 AUTHORS WHO PUBLISHED 4 PAPERS IN 2 CONFERENCES IN 2006 - 2014

**MORE DETAILED SUMMARIZED EXPLANATION BASED ON B:**
DUKE IS THE ORGANIZATION OF JUN Y. WHO PUBLISHED 4 PAPERS IN 2 CONFERENCES IN 2006 - 2014 AND BRETT W. WHO PUBLISHED 'ICHECK...' IN SIGMOD IN 2014

# SAMPLE USE-CASES

REPRESENTATIVE USE-CASES FROM THE USER STUDY:

- Q: RETURN THE AUTHORS WHO PUBLISHED PAPERS IN VLDB BEFORE 2016 AND AFTER 2007
  - A: JUN Y. PUBLISHED 9 PAPERS IN VLDB IN 2008 – 2015
- Q: RETURN THE AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES
  - A: JUN Y. PUBLISHED 64 PAPERS IN 18 CONFERENCES
- Q: RETURN THE ORGANIZATION OF AUTHORS WHO PUBLISHED PAPERS IN DATABASE CONFERENCES AFTER 2005
  - A: DUKE IS THE ORGANIZATION OF 63 AUTHORS WHO PUBLISHED 170 PAPERS IN 31 CONFERENCES IN 2006 - 2015

**FACTORIZING AND SUMMARIZING PROVENANCE FOR NATURAL LANGUAGE EXPLANATIONS**

VLDB 16', VLDB 17', SIGMOD REC. 18', VLDB J. 20'

**ABSTRACTING PROVENANCE FOR QUERY PRIVACY**

SIGMOD 21', ICDE 21'

|  | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_?$ | | | Reddit |
| $i_?$ | | | Facebook |
| $i_?$ | | | Reddit |

|  | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

|  | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |

DOH!

YOUR HOBBY IS DANCE ACCORDING TO FACEBOOK AND IT WAS PUBLISHED ON REDDIT THAT YOU ARE INTERESTED IN MUSIC

THE GENERAL PROPRIETARY CRITERION FOR SHOWING THE AD

YOUR HOBBY IS DANCE ACCORDING TO LINKEDIN AND YOU ARE INTERESTED IN MUSIC ACCORDING TO FACEBOOK

Deutch, G., "Reverse-Engineering Conjunctive Queries from Provenance Examples". In EDBT 2019, pp. 277-288

| | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_4$ | 1 | Parties | Reddit |
| $i_5$ | 2 | Parties | Facebook |
| $i_6$ | 4 | Movies | Reddit |

| | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

| | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |

WOOHOO!!!

**SOME INFORMATION FROM FACEBOOK** AND IT WAS PUBLISHED ON REDDIT THAT YOU ARE INTERESTED IN MUSIC

**SOME INFORMATION FROM LINKEDIN** AND YOU ARE INTERESTED IN MUSIC ACCORDING TO FACEBOOK

| | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_4$ | 1 | Parties | Reddit |
| $i_5$ | 2 | Parties | Facebook |
| $i_6$ | 4 | Movies | Reddit |

| | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

| | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |

$Q$(id):-Person(id,name,age), Hobbies(id,'Dance',src1), Interests(id,'Music',src2)

RETURN THE ID OF A PERSON WHOSE HOBBY IS `DANCE' AND WHOSE INTEREST IS `MUSIC'

# PROVENANCE MODEL

| | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_4$ | 1 | Parties | Reddit |
| $i_5$ | 2 | Parties | Facebook |
| $i_6$ | 4 | Movies | Reddit |

| | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

| | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |

```
Q(1):-Person(1,James  T,27), Hobbies(1,'Dance',Facebook),
Interests(1,'Music',Reddit)
```

Output: 1

Provenance: $p_1 \cdot i_1 \cdot h_1$

**Green, Karvounarakis, Tannen, "Provenance Semirings".** PODS: pp. 31-40, 2007

# PROVENANCE EXAMPLE FOR SPJU QUERY RESULTS

| | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_4$ | 1 | Parties | Reddit |
| $i_5$ | 2 | Parties | Facebook |
| $i_6$ | 4 | Movies | Reddit |

| | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

| | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |

## PROVENANCE EXAMPLE WITH TWO TUPLES

| Output | Provenance |
|---|---|
| 1 | $p_1 \cdot i_1 \cdot h_1$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

# PROVENANCE ABSTRACTION

| | PID | Interest | Source |
|---|---|---|---|
| $i_1$ | 1 | Music | Reddit |
| $i_2$ | 2 | Music | Facebook |
| $i_3$ | 3 | Music | LinkedIn |
| $i_4$ | 1 | Parties | Reddit |
| $i_5$ | 2 | Parties | Facebook |
| $i_6$ | 4 | Movies | Reddit |

| | PID | Hobby | Source |
|---|---|---|---|
| $h_1$ | 1 | Dance | Facebook |
| $h_2$ | 2 | Dance | LinkedIn |
| $h_3$ | 4 | Dance | Facebook |
| $h_4$ | 1 | Trips | Facebook |
| $h_5$ | 2 | Trips | LinkedIn |
| $h_6$ | 3 | Trips | Reddit |

| | PD | Name | Age |
|---|---|---|---|
| $p_1$ | 1 | James T | 27 |
| $p_2$ | 2 | Brenda P | 31 |



**Deutch, Moskovitch, Rinetzky, "Hypothetical Reasoning via Provenance Abstraction".** SIGMOD: pp. 537-554, 2019

Tree diagram:

- `*`
  - Social Network
    - Facebook
      - $h_1$
      - $h_3$
      - $h_4$
      - $i_2$
      - $i_5$
    - LinkedIn
      - $h_2$
      - $h_5$
      - $i_3$
  - Reddit
    - $h_6$
    - $i_1$
    - $i_4$
    - $i_6$

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot i_1 \cdot h_1$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot i_1 \cdot Facebook$ |
| 2 | $p_2 \cdot i_2 \cdot LinkedIn$ |

$Q(\text{ID}):-\text{PERSON}(\text{ID},\text{NAME},\text{AGE}),\ \text{HOBBIES}(\text{ID},\text{'DANCE'},\text{SRC1}),$

$\text{INTERESTS}(\text{ID},\text{'MUSIC'},\text{SRC2})$

**ALL QUERIES WILL GENERATE THE PROVENANCE**

$Q1(\text{ID})\ :-\ \text{PERSON}(\text{ID},\text{NAME},\text{AGE}),\ \text{HOBBIES}(\text{ID},\text{'TRIPS'},\text{SRC1}),$

$\text{INTERESTS}(\text{ID},\text{'MUSIC'},\text{SRC2})$

$Q2(\text{ID})\ :-\ \text{PERSON}(\text{ID},\text{NAME},\text{AGE}),\ \text{HOBBIES}(\text{ID},\text{'DANCE'},\text{SRC1}),$

$\text{INTERESTS}(\text{ID},\text{'PARTIES'},\text{SRC2})$

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot i_1 \cdot Facebook$ |
| 2 | $p_2 \cdot i_2 \cdot LinkedIn$ |

$Q$(ID):-PERSON(ID,NAME,AGE), HOBBIES(ID,'DANCE',SRC1), INTERESTS(ID,'MUSIC',SRC2)

| Output | Provenance |
|--------|-----------------------------|
| 1 | $p_1 \cdot i_1 \cdot Facebook$ |
| 2 | $p_2 \cdot i_2 \cdot LinkedIn$ |

$Q$(ID):-PERSON(**ID**,NAME,AGE), HOBBIES(**ID**,'DANCE',SRC1), INTERESTS(**ID**,'MUSIC',SRC2)

- CONNECTED

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot i_1 \cdot Facebook$ |
| 2 | $p_2 \cdot i_2 \cdot LinkedIn$ |

$Q$(1):-PERSON(1,JAMES T,27), HOBBIES(1,'DANCE',FACEBOOK), INTERESTS(1,'MUSIC',REDDIT)

- CONNECTED

- CONSISTENT - GENERATES THE DESIRED PROVENANCE FOR EACH OF THE RESULTS IN ONE OF THE CONCRETE OPTIONS

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot i_1 \cdot Facebook$ |
| 2 | $p_2 \cdot i_2 \cdot LinkedIn$ |

$Q$(ID):–PERSON(ID,NAME,AGE), HOBBIES(ID,'DANCE',SRC1),
INTERESTS(ID,'MUSIC',SRC2)

> **IF WE HAVE K SUCH CANDIDATE QUERIES, WE SAY THAT THE ABSTRACTION HAS PRIVACY K**

- CONNECTED

- CONSISTENT – GENERATES THE DESIRED PROVENANCE FOR EACH OF THE RESULTS IN ONE OF THE CONCRETE OPTIONS

- INCLUSION MINIMAL – NO OTHER CONSISTENT QUERY IS CONTAINED IN IT
  - Deutch, G., "Reverse-Engineering Conjunctive Queries from Provenance Examples". In EDBT 2019, pp. 277-288

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \boldsymbol{Reddit}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

Tree diagram:

- `*`
  - Social Network
    - Facebook
      - $h_1$, $h_3$, $h_4$, $i_2$, $i_5$
    - LinkedIn
      - $h_2$, $h_5$, $i_3$
  - Reddit
    - $h_6$, $i_1$, $i_4$, $i_6$

**Measure information loss with Entropy =**

$$-\sum_i P_X(x_i) ln(P_X(x_i)).$$

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \textbf{Reddit}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \boldsymbol{h_6}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \boldsymbol{i_1}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \boldsymbol{i_4}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

| Output | Provenance |
|--------|------------|
| 1 | $p_1 \cdot h_1 \cdot \boldsymbol{i_6}$ |
| 2 | $p_2 \cdot i_2 \cdot h_2$ |

**PROBLEM DEFINITION: G**IVEN AN ABSTRACTION TREE, A PROVENANCE EXAMPLE, AND A PRIVACY THRESHOLD K, FIND AN ABSTRACTION FOR THE EXAMPLE THAT ACHIEVES PRIVACY ≥ K AND INCURS THE MINIMUM LOSS OF INFORMATION OVER ALL ABSTRACTIONS THAT ACHIEVE THE PRIVACY THRESHOLD K.

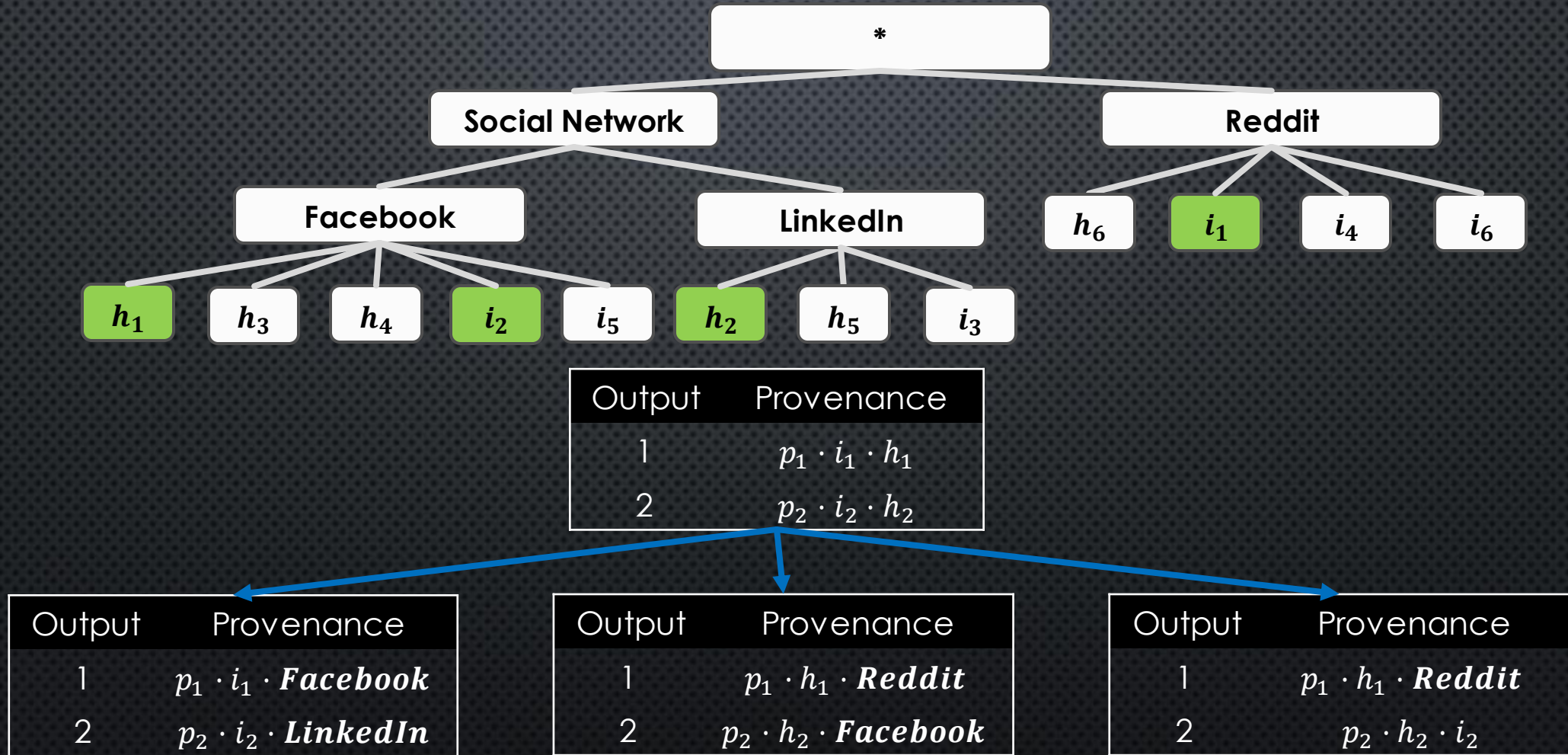**PROPOSITION: THE DECISION VERSION OF THE OPTIMAL ABSTRACTION PROBLEM IS NP-HARD.**
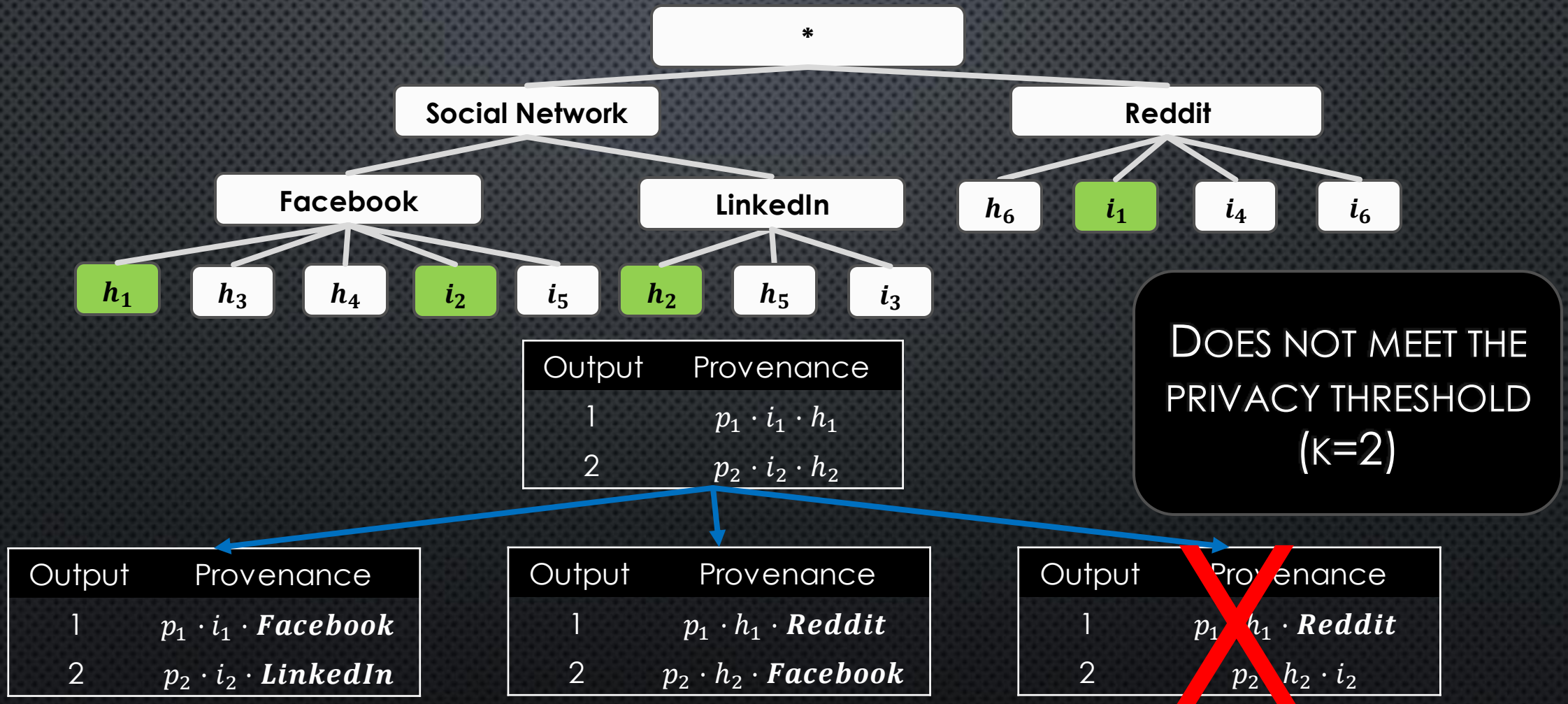
# ABSTRACTION COMPUTATION

**GUARANTEE (INFORMAL):** THE ALGORITHM FIND AN OPTIMAL ABSTRACTION.

# SAMPLE EXPERIMENTAL RESULTS

## RUNTIME AS A FUNCTION OF THE PRIVACY THRESHOLD

# Takeaways

1. There are different ways to manipulate raw provenance, including:
    I. Factorization and summarization
    II. Abstraction

2. <u>Factorization and summarization</u> can help make provenance understandable and "easier to digest" for creating explanations

3. <u>Abstraction</u> can help preserve the privacy of the query while providing explanations

4. <u>Tradeoff:</u> smaller factorization/higher privacy threshold = less informative explanations