# Querying Graph Data
## Where Are We?　　Where Should We Go?

Leonid Libkin
RelationalAI, IRIF,
University of Edinburgh

Filip Murlak
University of Warsaw

Wim Martens
University of Bayreuth
RelationalAI

Domagoj Vrgoč
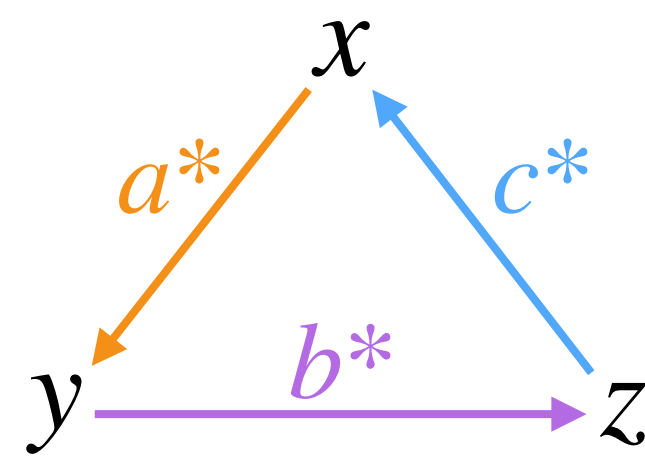PUC Chile

Liat Peterfreund
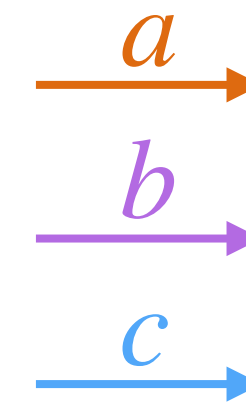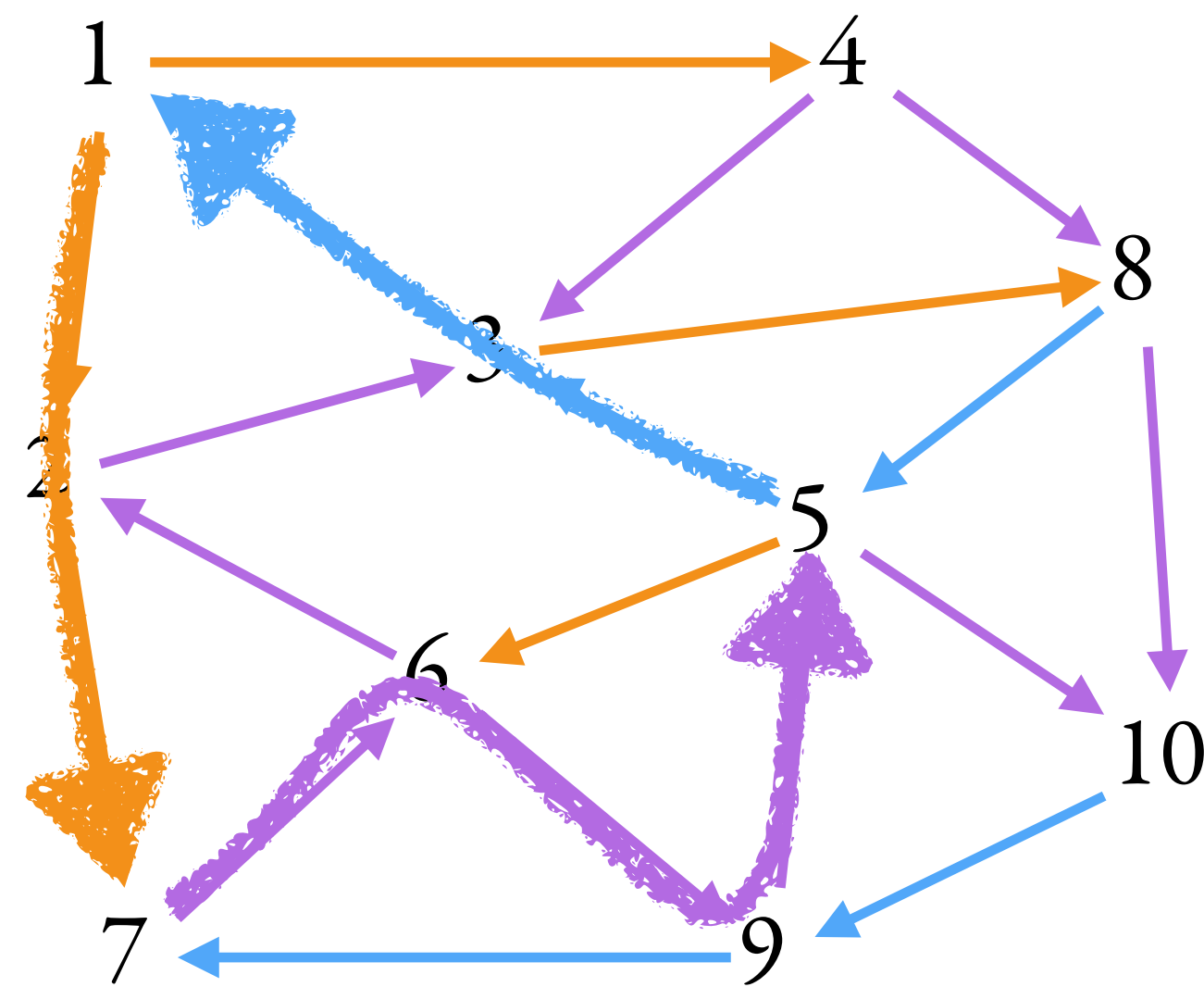Hebrew University

Gems Talk
PODS 2025

Paper:

# Conjunctive Regular Path Queries (CRPQ)

# Conjunctive Regular Path Queries 2.0

**What Do Cypher, SQL/PGQ, GQL Patterns add to CRPQs?**

(a) handling of nodes and edges
(b) path and list variables      ⤳ we'll get to this
(c) path modes      ⤳ simple, trail,…
(d) data filters      ⤳ data value comparisons

The "Four Features"

How do we cleanly define these?      so that we can study them

How do we cleanly design these in a real-world language?      so that they can be used

Query Language Design! Nice!

(c) [Bagan, Bonifati, Groz PODS'13] [M., Trautner ICDT'18] [M., Niewerth, Trautner STACS'20] [M., Popp PODS'22]
(d) [Neven, Schwentick, Vianu TOCL'04] [Bojanczyk et al. PODS'06, LICS'06] [Libkin, M., Vrgoc ICDT'13]

# The Broader Context

# Big Questions

## Database Metatheory: Asking the Big Queries

### Christos H. Papadimitriou
University of California San Diego

PODS 1995

# Big Questions

What is the most important open research question for database theory?

My humble opinion

# Can we make set semantics perform in practice?

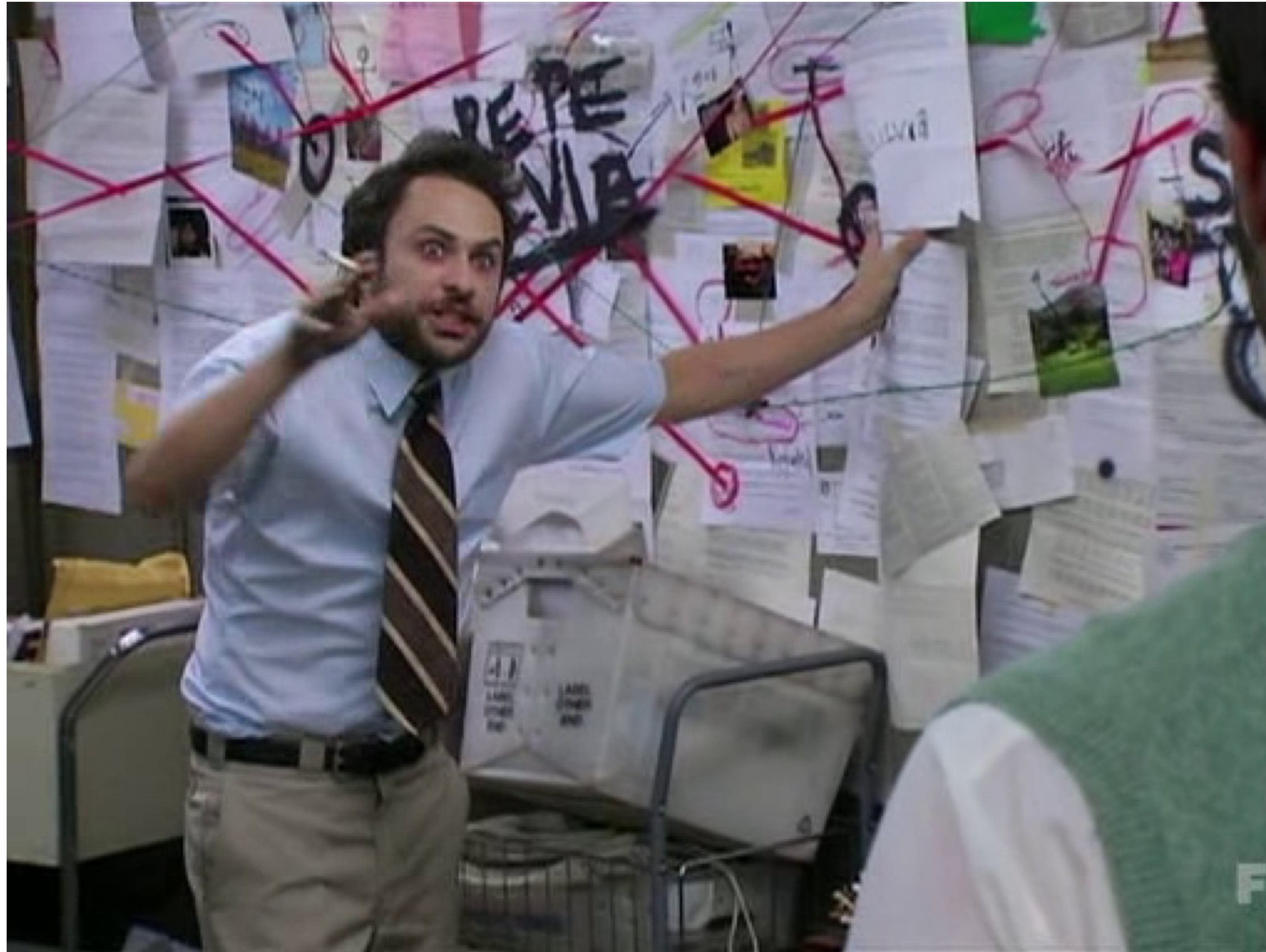Why is this question important?

# Database Research Landscape

Theory

Systems

Sets

Bags
(Multisets)

divide

why?

Con sets: more costly union & project
⤳ more data structure overhead
Pro sets: more room for query optimization
Pro sets: …

# Why Am I Talking About Graphs & Sets & Bags?

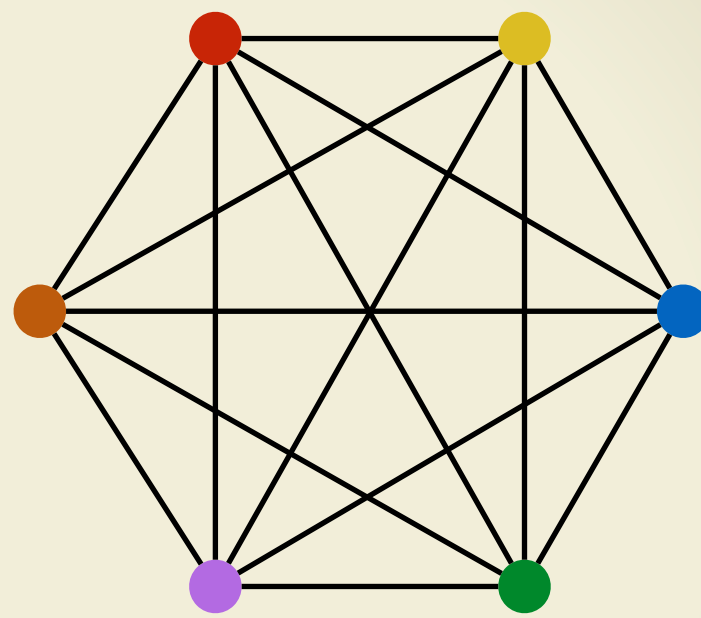...because everything is connected



"Everything is connected!"

(https://imgflip.com/memegenerator/268702234/Charliethinkingits-all-connectedconnect-the-dot)

# Bags & Recursion: Boom!

Bag
Semantics

Query    $((((a*)*)*)*$

Data

6-clique
every edge labeled $a$

Simple
Paths

Every answer is
returned $10^{269}$ times

Smoking gun
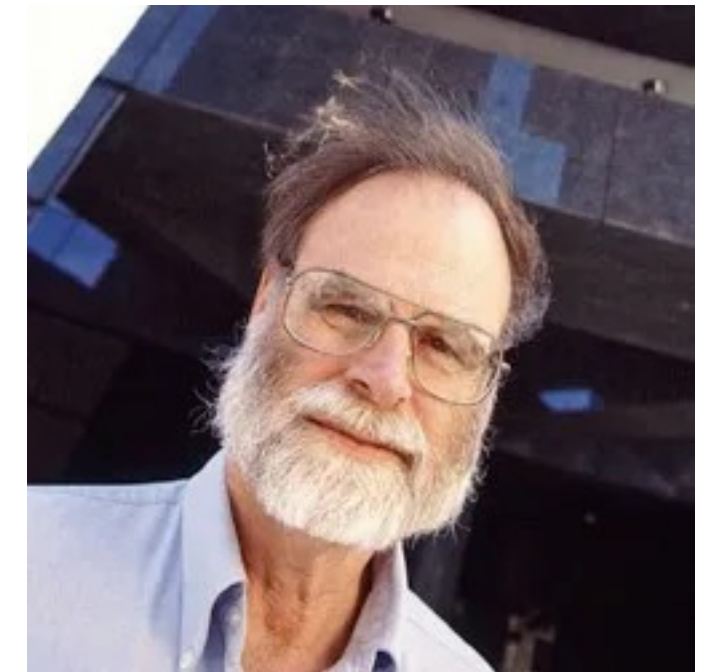for bags?

[Arenas, Conca, Pérez WWW'12]

# What Are We Doing
# in (Query) Language Design?

# Language Design: The Goal?

## Useful to give us direction

**Automatic Programming**

- Allowing human programmers to code at a significantly higher level
- Focusing on the problem that needs solving, not on its "administrative aspects"

Jim Gray

We, the DB community (theory & systems) are **great** at this

- Declarative query languages
- Automatic optimization, automatic out-of-core computation, etc.

Molham Aref

The only thing is that we're doing it just for DB queries

Why not think bigger?

Can we target more general-purpose programming?

If so, then our languages better be well-designed!

We need the right principles!

Thursday 17:00
**SIGMOD Industry 6**
Be there
I kid you not
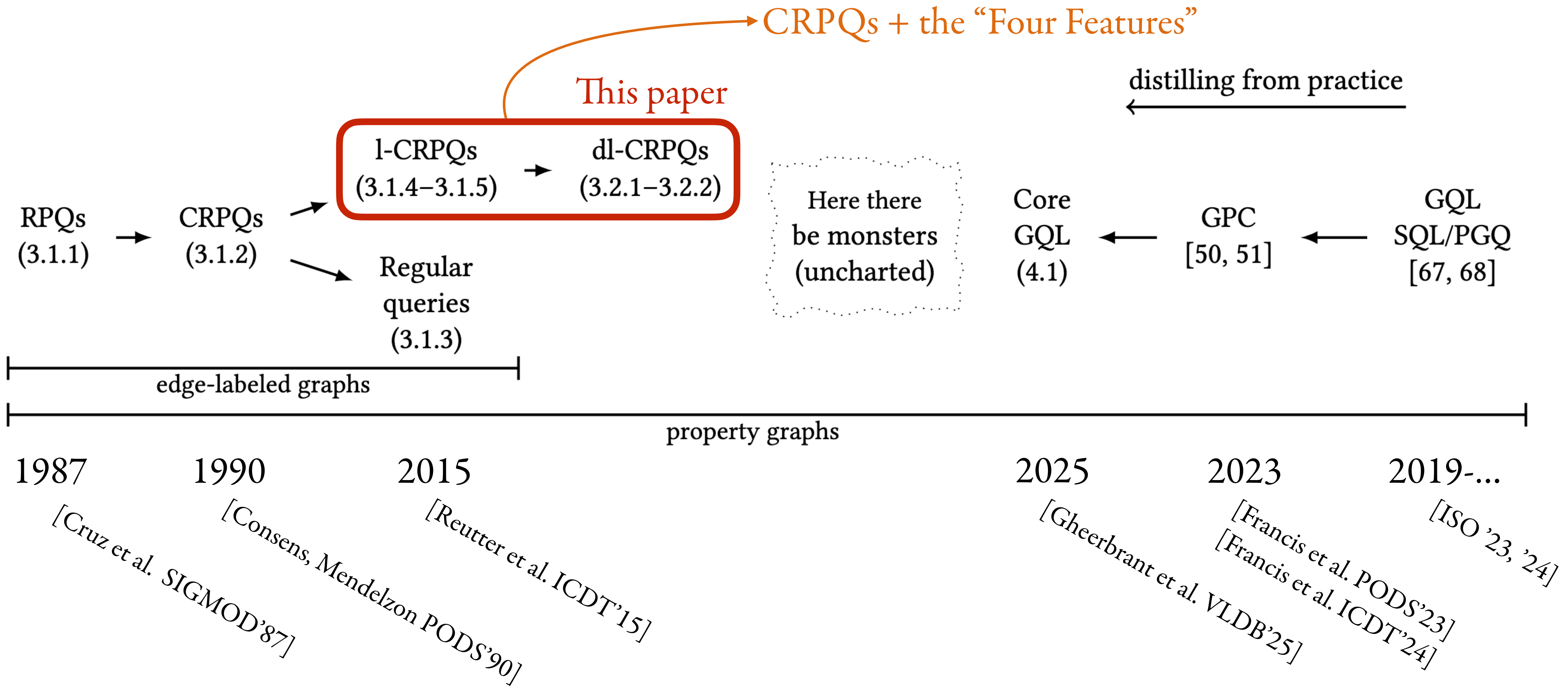
# Where We Are

We wanted to do query language design in graph pattern matching

# Graph Pattern Matching Landscape

CRPQs + the "Four Features"

distilling from practice

This paper

l-CRPQs (3.1.4–3.1.5) → dl-CRPQs (3.2.1–3.2.2)

Here there be monsters (uncharted)

RPQs (3.1.1) → CRPQs (3.1.2) → Regular queries (3.1.3)

Core GQL (4.1) ← GPC [50, 51] ← GQL SQL/PGQ [67, 68]

edge-labeled graphs

property graphs

1987 [Cruz et al. SIGMOD'87]

1990 [Consens, Mendelzon PODS'90]

2015 [Reutter et al. ICDT'15]

2025 [Gheerbrant et al. VLDB'25]

2023 [Francis et al. PODS'23] [Francis et al. ICDT'24]
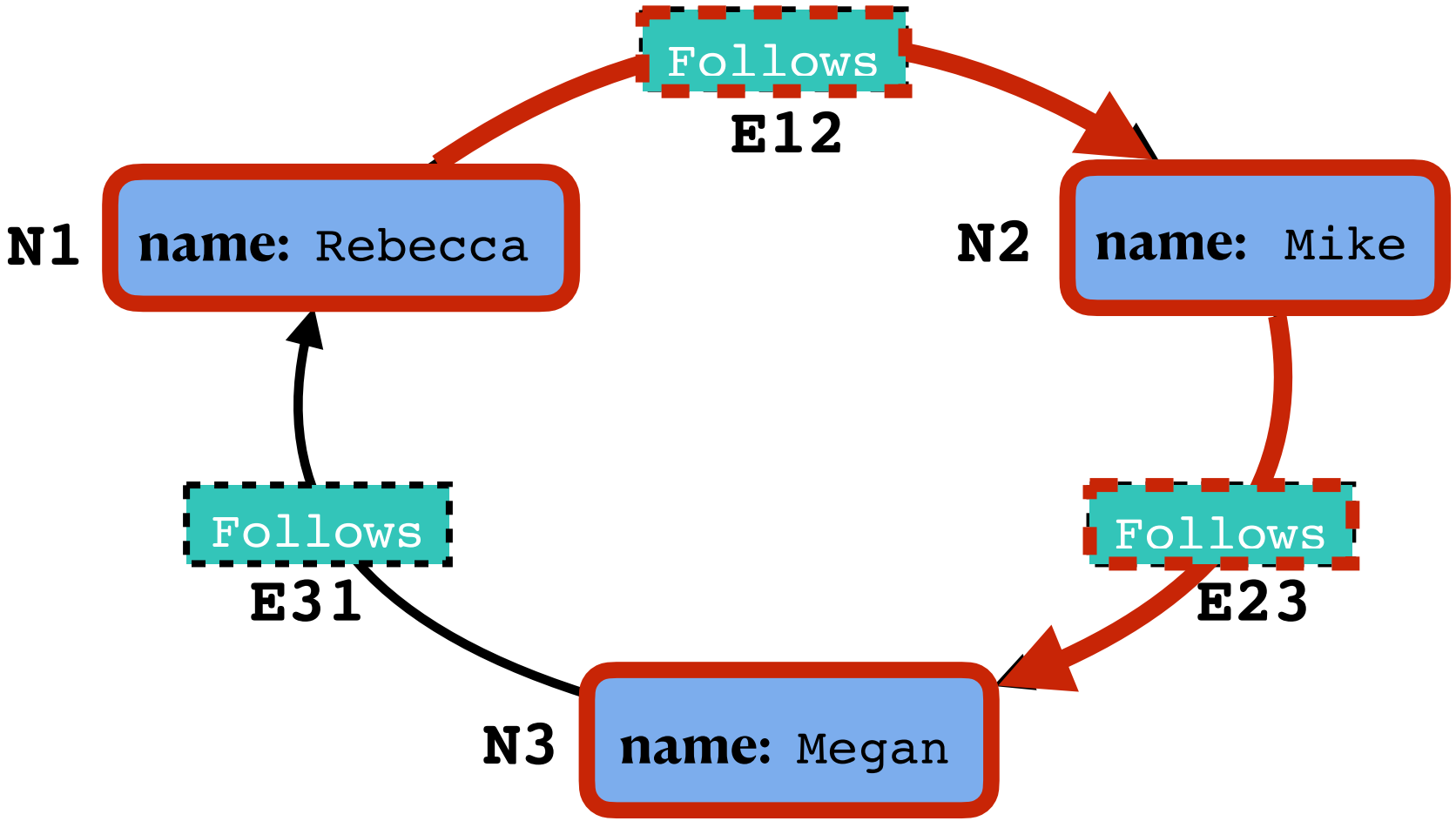
2019-… [ISO '23, '24]

# Why Theory Is Needed

Example 1

```
(x) ( ()-[ :Follows]->() ){2} (y)
```

(...) node

[...] edge

paths from x to y

edge labeled "Follows"

repeat the edge 2 times



Output:

| x | y |
|----|----|
| N1 | N3 |
| N2 | N1 |
| N3 | N2 |

Equivalent to

```
(x) ( ()-[:Follows]->() ) ( ()-[:Follows]->() ) (y)
(x) ( ()-[:Follows]->()-[:Follows]->() ) (y)
```

# Why Theory Is Needed

Example 1

`(x) ( ()-[z:Follows]->() ){2} (y)`

paths from x to y

edge labeled "Follows", in variable z

repeat the edge 2 times

Follows
**E12**

**N1** name: Rebecca

**N2** name: Mike

Follows
**E31**

Follows
**E23**

**N3** name: Megan

Output:

| x | y | z |
|---|---|---|
| N1 | N3 | [E12, E23] |
| N2 | N1 | [E23, E31] |
| N3 | N2 | [E31, E12] |

**Not** equivalent to any of

`(x) ( ()-[z:Follows]->() ) ( ()-[z:Follows]->() ) (y)`

`(x) ( ()-[z:Follows]->()-[z:Follows]->() ) (y)`

`(x) ( ()-[z1:Follows]->()-[z2:Follows]->() ) (y)`

# Why Theory Is Needed

**Example 2**

```
(x) ( ( )-[:Follows]->( )-[:Follows]->( ) )+ (y)
```

→ paths of even length

→ edges labeled "Follows"



Output:
(trails only)

| x | y |
|----|----|
| N1 | N3 |
| N2 | N1 |
| N3 | N2 |

# Why Theory Is Needed

→ paths of even length

→ edges labeled "Follows"

→ using variable z

Follows
**E12**

N1 **name:** Rebecca

N2 **name:** Mike

Follows
**E31**

Follows
**E23**

N3 **name:** Megan

Output:      Nothing!

Why?

# Why Theory Is Needed

```
(x) ( (z)-[:Follows]->(z)-[:Follows]->( ) )+ (y)
```

paths of even length

edges labeled "Follows"

using variable z



**N1** name: Rebecca

**N2** name: Mike

**N3** name: Megan

Follows **E12**

Follows **E23**

Follows **E31**

Follows

Follows

Output:

| x | y | z |
|---|---|---|
| N1 | N3 | [N1, N2] |

...and others

Again, why?

# Why Theory Is Needed



**Example 2**

`(z)-[:Follows]->(z)-[:Follows]->( )`

**Syntax-driven design!**

# Why Theory Is Needed



Example 2

`(x) ( (z)-[:Follows]->(z)-[:Follows]->( ) )+ (y)`

Output:

| x | y | z |
|----|----|----------|
| N1 | N3 | [N1, N2] |

Need to separate concerns?   (Joins ⟷ Lists)

If you think about it as CRPQs, then
- the RPQs are about matching paths
- the CRPQ vars are about joining

Syntax-driven design!

# Why Theory Is Needed

```
p = (x) ( (u)-[ ]->(v) WHERE u.date < v.date)* (y)
```

Paths from x to y, where dates increase on nodes

Well done!
Cypher
SQL/PGQ
GQL

### Intermezzo: in SQL, you'd need to start with...

```
WITH good_edge(S,T) AS
  SELECT Source.N_ID AS S, Target.N_ID AS T
  FROM Source, Target, Dates D1, Dates D2
  WHERE Source.E_ID = Target.E_ID
    AND D1.N_ID=Source.N_ID AND D2.N_ID = Target.N_ID
    AND D1.date < D2.date
RECURSIVE path_n(S, T) AS
  SELECT * FROM good_edge
  UNION
  SELECT good_edge.S, path_n.T
  FROM good_edge, path_n
  WHERE good_edge.T=path_n.S
SELECT * FROM path_n
```

# Why Theory Is Needed

```
p = (x) ( (u)-[ ]->(v) WHERE u.date < v.date)* (y)
```

Paths from x to y, where dates increase on nodes

How do you do paths from x to y, where dates increase on edges?
Umm....

```
p = (x) (   (u) -[ ]->(v)        WHERE u.date < v.date)* (y)
```
increasing on nodes ✓

```
p = (x) ( -[u]->( ) -[v]->      WHERE u.date < v.date)* (y)
```
**not** increasing on edges ✖

```
p = (x) ( -[u]->( ) -[v]->() WHERE u.date < v.date)* (y)
```
**not** increasing on edges ✖

⤳ match 1 3 2 4

It Looks Like There's Work To Be Done Here

# Wait, Weren't You People on ISO?

# Wait, Weren't You People on ISO?

You don't have total control

Apple is like a ship,
with a hole in the bottom,
leaking water
and my job is
to get the ship
pointed
in the right direction

—Steve Jobs

Shout out to



Nadime
Francis



Victor
Marsault

Made sense of
- constant stream of proposals
- each with dozens of pages
- in standardese

Paper:

# What Now?

## Where to go?

This is the starting point of the paper — we've only just gone through the intro

We're going to propose something. Let's first argue why we are proposing it.

# Research Agenda

**Query Languages in Theory**

Sets

Declarative

Automata

**Prove advantages of our principles!**

**Systems**

Bags

Syntax-Driven

Efficient

Actually work in practice

**Make them work in practice!**

# What's Next in the Talk?

growing from theory →

← distilling from practice

l-CRPQs
(3.1.4–3.1.5)  →  dl-CRPQs
(3.2.1–3.2.2)

Here there
be monsters
(uncharted)

Core
GQL
(4.1)

GQL
SQL/PGQ
[67, 68]

GPC
[50, 51]  ←

RPQs
(3.1.1)  →  CRPQs
(3.1.2)

Regular
queries
(3.1.3)

edge-labeled graphs

property graphs

**We choose to follow**

(1) compatibility with automata
(2) set semantics
(3) symmetric treatment of nodes and edges

**These principles seem to work very well**

- definitions become elegant
- fewer semantic issues
- a lot of potential for query optimization

(language becomes "more declarative")

# Growing From Theory

# CRPQs with List Variables

## Standard CRPQs

$$q(x_1, \ldots, x_k) \;=\; \bigwedge_{i=1}^{n} y_i \xrightarrow{r_i} z_i$$

## Ingredients

- Output variables: $x_1, \ldots, x_k$
- Join variables: $y_1, z_1, \ldots, y_n, z_n$
- Regular expressions : $r_1, \ldots, r_n$

## Current GQL

```
()-[z:Follows]->()-[z:Follows]->()
( ()-[z:Follows]->()-[z:Follows]->() )+
```
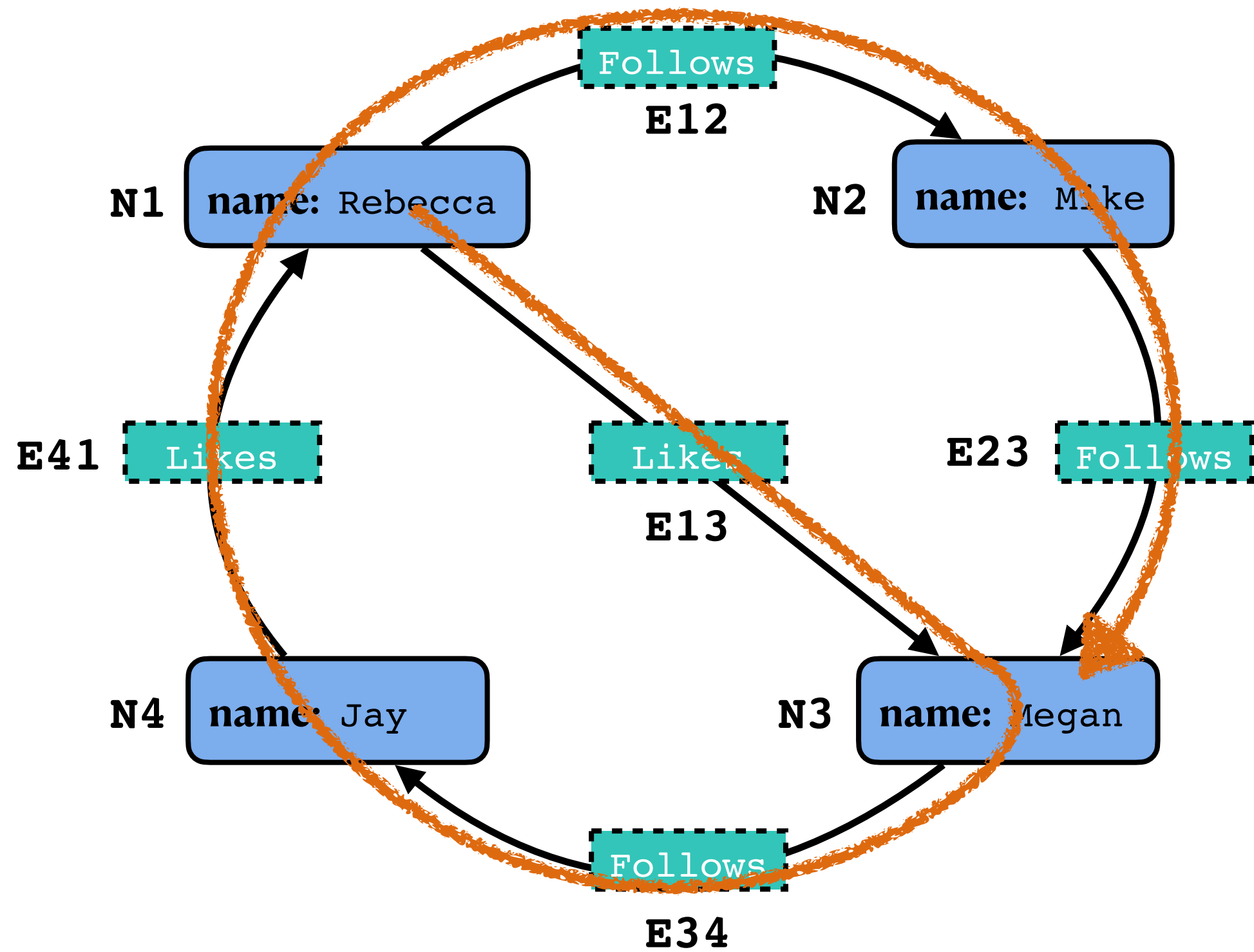
We want to add:



capability to return elements
along matched paths

⤳ **z** is a join variable
⤳ **z** is a "list" variable
        (called group variable in GQL)

# CRPQs with List Variables

## Standard CRPQs

$$q(x_1, \ldots, x_k) \ = \ \bigwedge_{i=1}^{n} y_i \xrightarrow{r_i} z_i$$

### Ingredients

- Output variables: $x_1, \ldots, x_k$
- Join variables: $y_1, z_1, \ldots, y_n, z_n$
- Regular expressions : $r_1, \ldots, r_n$

## CRPQs with List Variables

$$q(x_1, \ldots, x_k) \ = \ \bigwedge_{i=1}^{n} y_i \xrightarrow{r_i} z_i$$

### Ingredients

- Output variables: $x_1, \ldots, x_k$
- **Join variables**: $y_1, z_1, \ldots, y_n, z_n$
- Regular expressions with **list variables**: $r_1, \ldots, r_n$

Keep **join variables** & **list variables** separated
We'll soon see why

# REs with List Variables



$$( \text{Follows}^{\,z} + \text{Likes}\, )*$$

- Paths in which every edge is labeled Follows or Likes
- Annotate the Follows edges with variable $z$

Path:     E13 — E34 — E41 — E12 — E23
                            $z$              $z$        $z$

⤳ $z$ binds to [E34, E12, E23]

## Why a Design Like This?

Such REs
(1) are highly compatible with automata
(2) allow a product graph construction

## Again, Why?

(1) ⤳ Enables more query optimization
(2) ⤳ Enables factorization for matching paths

[Fagin, Kimelfeld, Reiss, Vansummeren PODS'13]
[Riveros, Van Sint Jan, Vrgoc VLDB'23]
[Doleschal, Kimelfeld, M., Nahshon, Neven PODS'19]

[M., Niewerth, Popp, Rojas, Vansummeren, Vrgoc VLDB'23]
[Farias, M., Rojas, Vrgoc ISWC'24]

# CRPQs with Data & List Variables



Transfer
date: 1/1/25
E12

N1 name: Rebecca

N2 name: Mike

Transfer
date: 1/1/24
E41

Transfer
date: 1/1/24
E13

Transfer
date: 1/2/25
E23

N4 name: Jay

N3 name: Megan

Transfer
date: 1/3/25
E34

**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{\,z}\right]\left[x := \text{date}\right]\left(\left(\_\right)\left[\text{Transfer}^{\,z}\right]\left[\text{date} > x\right]\left[x := \text{date}\right]\right)^{*}$$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$[\text{Transfer}^z]$$

Register

Considered path:

# CRPQs with Data & List Variables



We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges
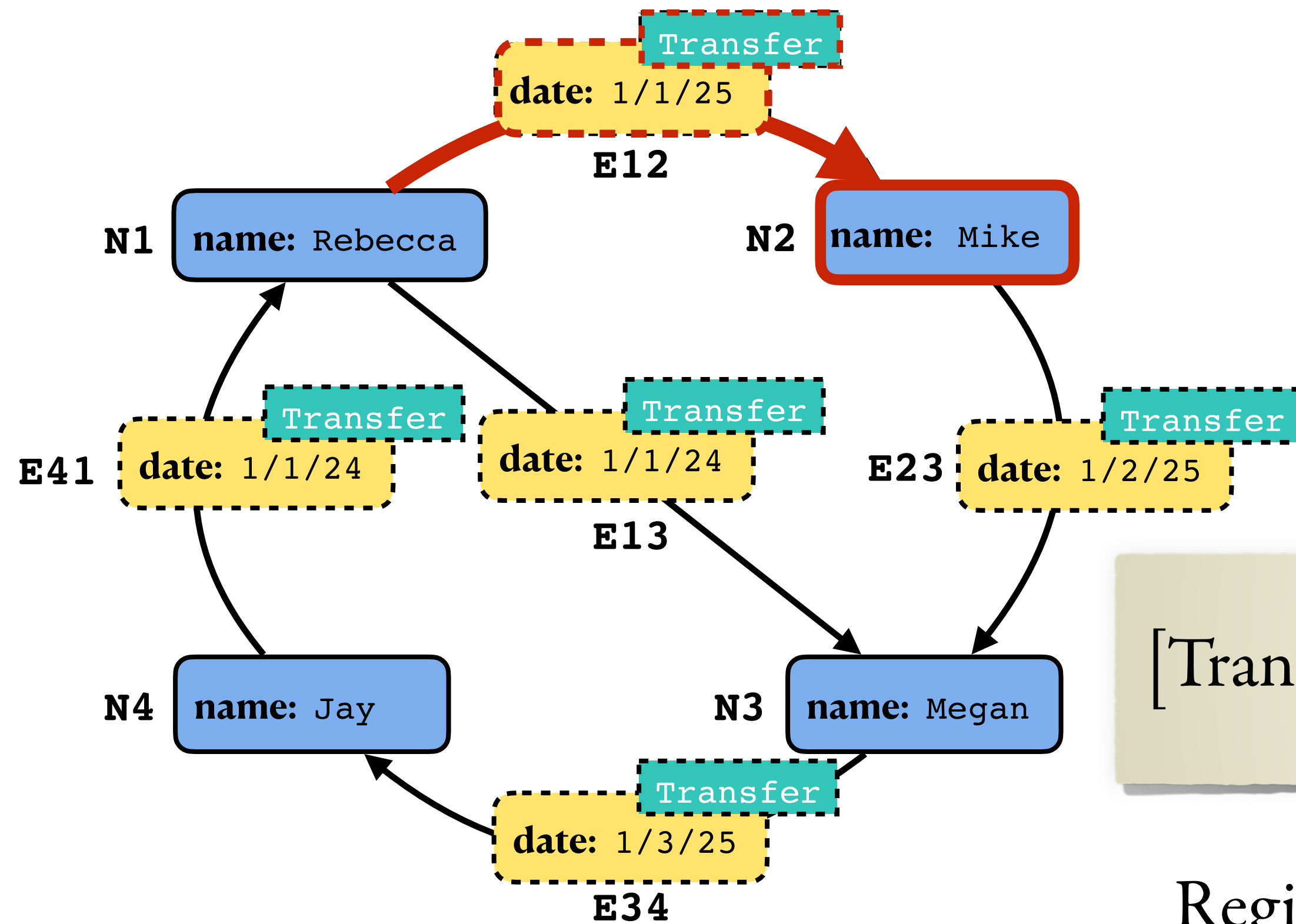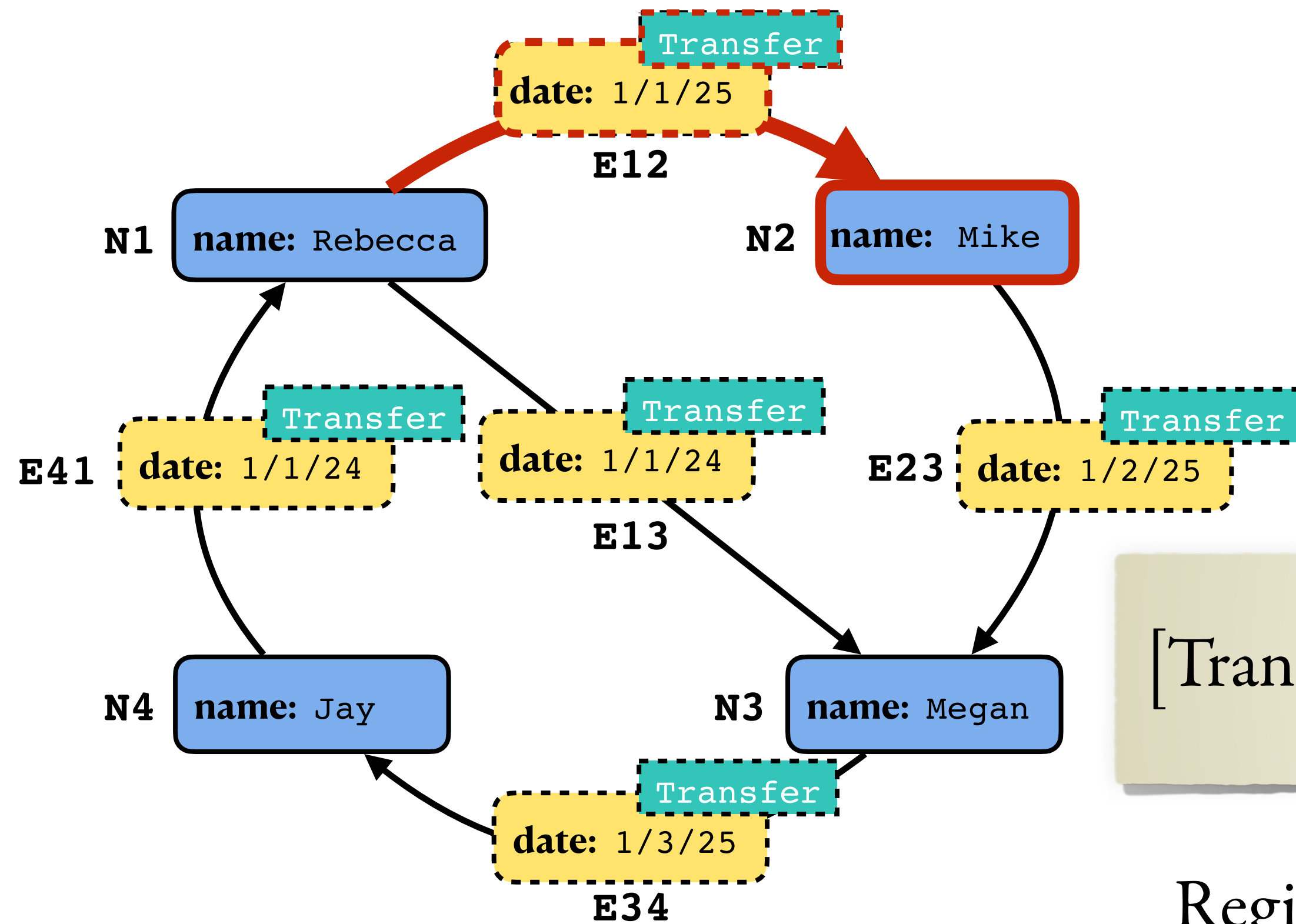
$$[\text{Transfer}^z]$$

Register

Considered path:

# CRPQs with Data & List Variables

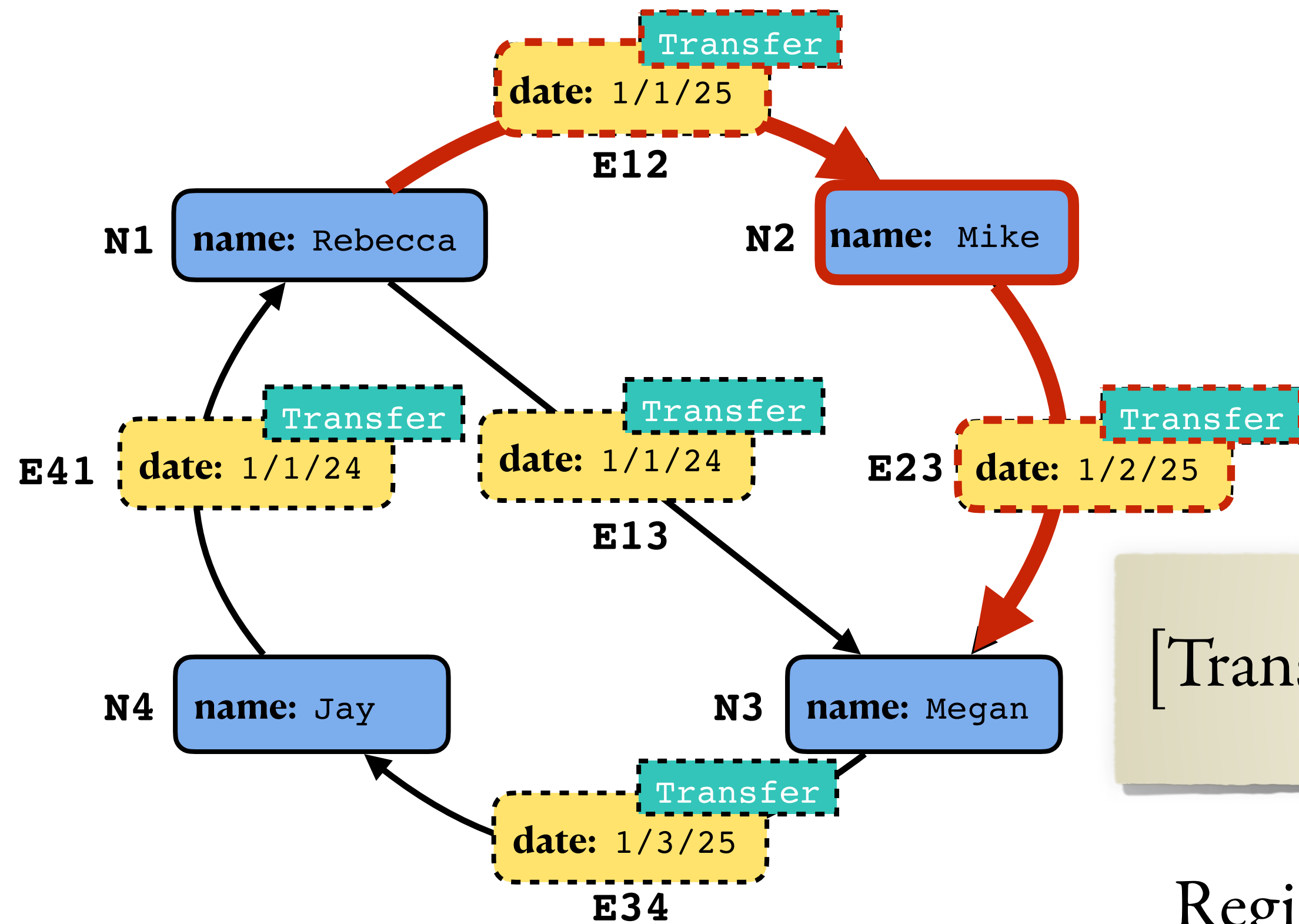

We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$[\text{Transfer}^z]$$

Register

Considered path:    E12

$$z$$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

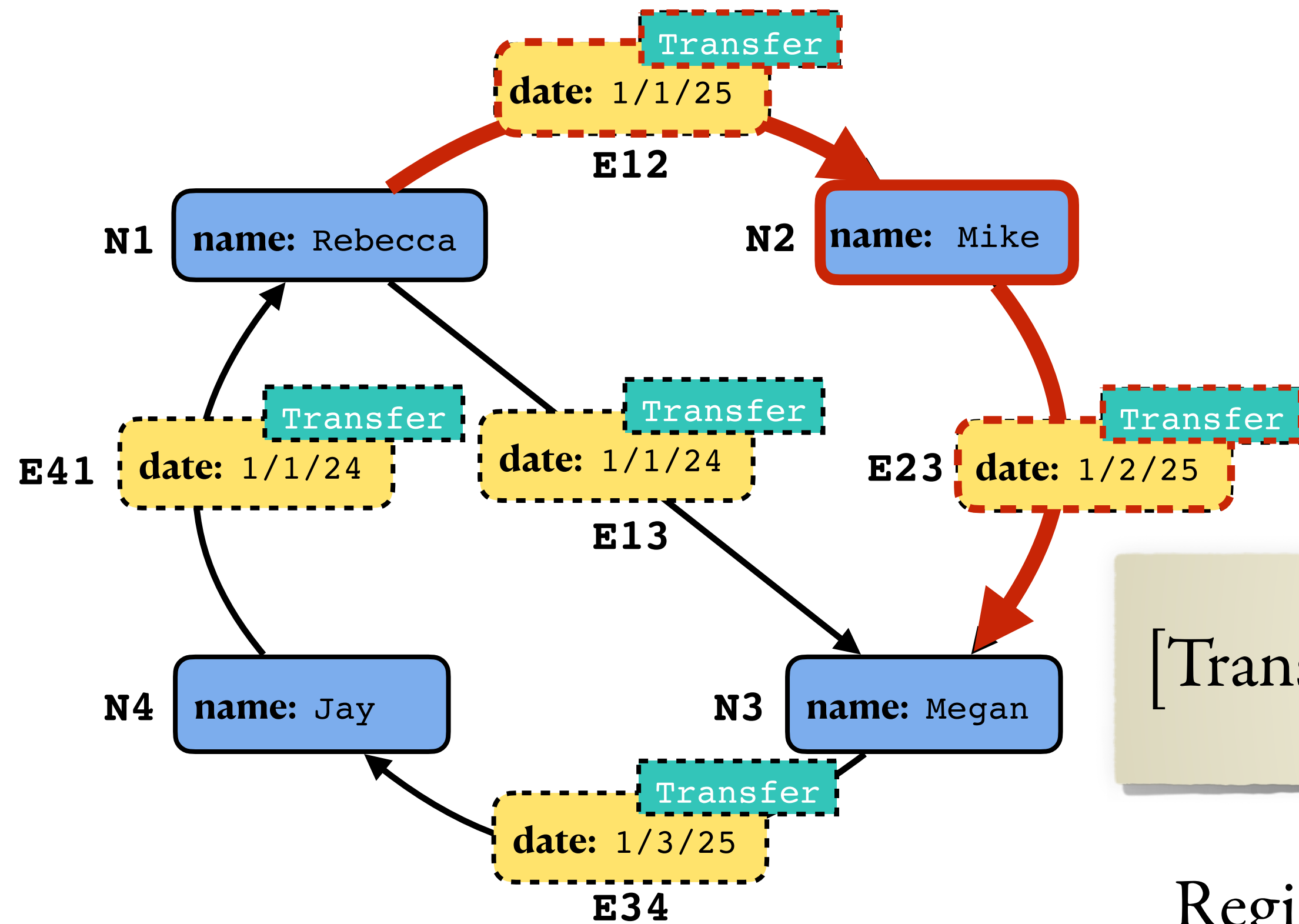Increasing date values on edges

$$\left[\text{Transfer}^{\,z}\right]\left[x := \text{date}\right]$$

Register

Considered path:    E12

$z$

# CRPQs with Data & List Variables



Transfer
**date:** 1/1/25
**E12**

**N1** name: Rebecca
**N2** name: Mike

**E41** Transfer **date:** 1/1/24

Transfer **date:** 1/1/24
**E13**

**E23** Transfer **date:** 1/2/25

**N4** name: Jay
**N3** name: Megan

Transfer **date:** 1/3/25
**E34**

## We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$[\text{Transfer}^z]\,[x := \text{date}]$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12

$z$

# CRPQs with Data & List Variables



Transfer

**date:** 1/1/25

**E12**

**N1** **name:** Rebecca

**N2** **name:** Mike

Transfer

**E41** **date:** 1/1/24

Transfer

**date:** 1/1/24

**E13**

Transfer

**E23** **date:** 1/2/25

**N4** **name:** Jay

**N3** **name:** Megan

Transfer

**E34** **date:** 1/3/25

## We're Going to Add

- node & edge treatment
- data filters

Increasing date values on edges

$$[\text{Transfer}^{z}]\,[x := \text{date}]\left(\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}\right)^{*}$$

Register     $x = 1/1/25$

Considered path:     E12

$z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{z}\right]\left[x := \text{date}\right]\Big(\big(\_\big)\qquad\qquad\Big)^{*}$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12

$\qquad\qquad\qquad z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^z\right]\left[x := \text{date}\right]\Big(\big(\_\big) \Big)^{*}$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12 $\longrightarrow$ N2
$\qquad\qquad\qquad\qquad\quad z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{z}\right]\left[x := \text{date}\right]\Big(\left(\_\right)\left[\text{Transfer}^{z}\right]\right)^{*}$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12 $\xrightarrow{\ z\ }$ N2

# CRPQs with Data & List Variables



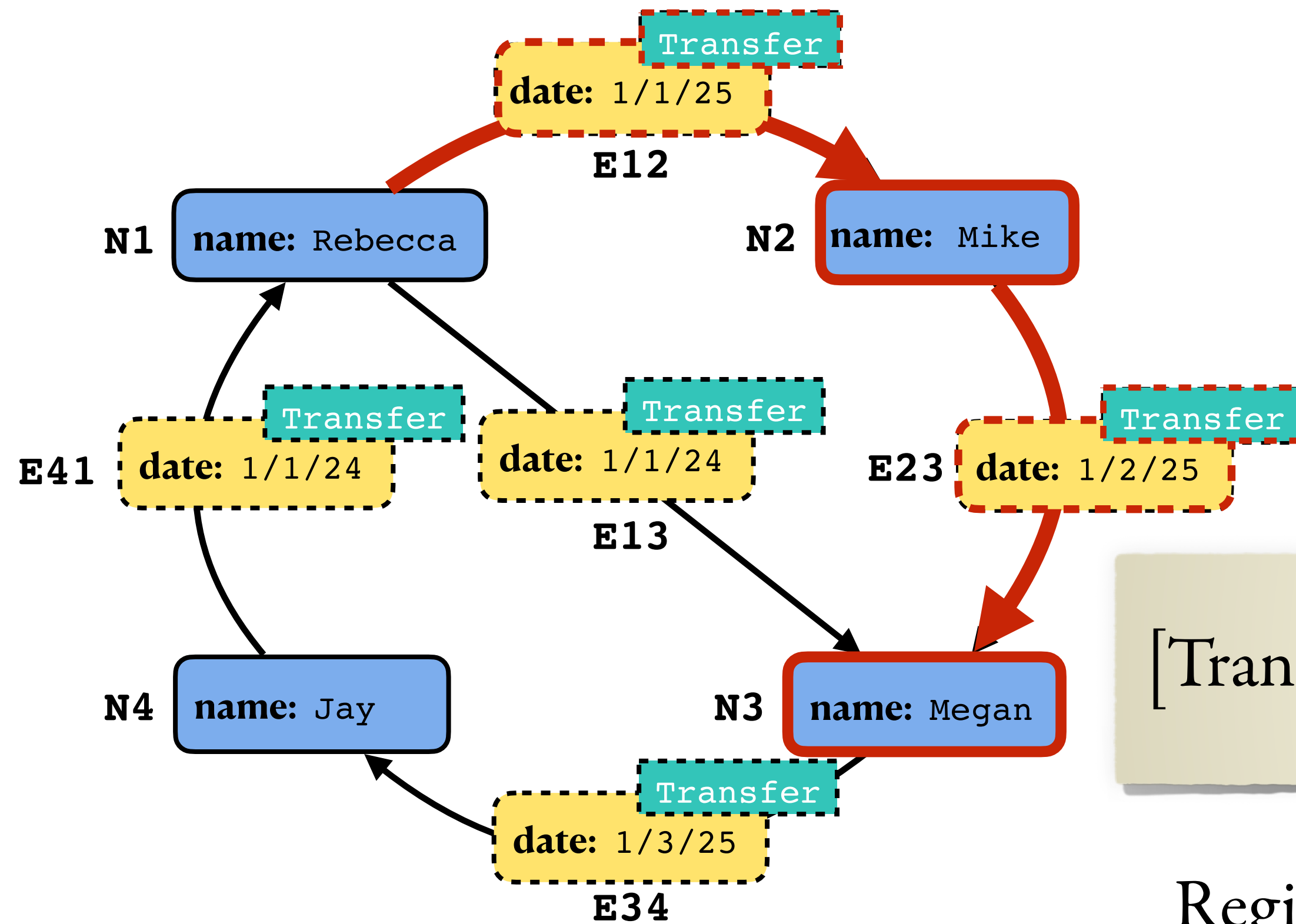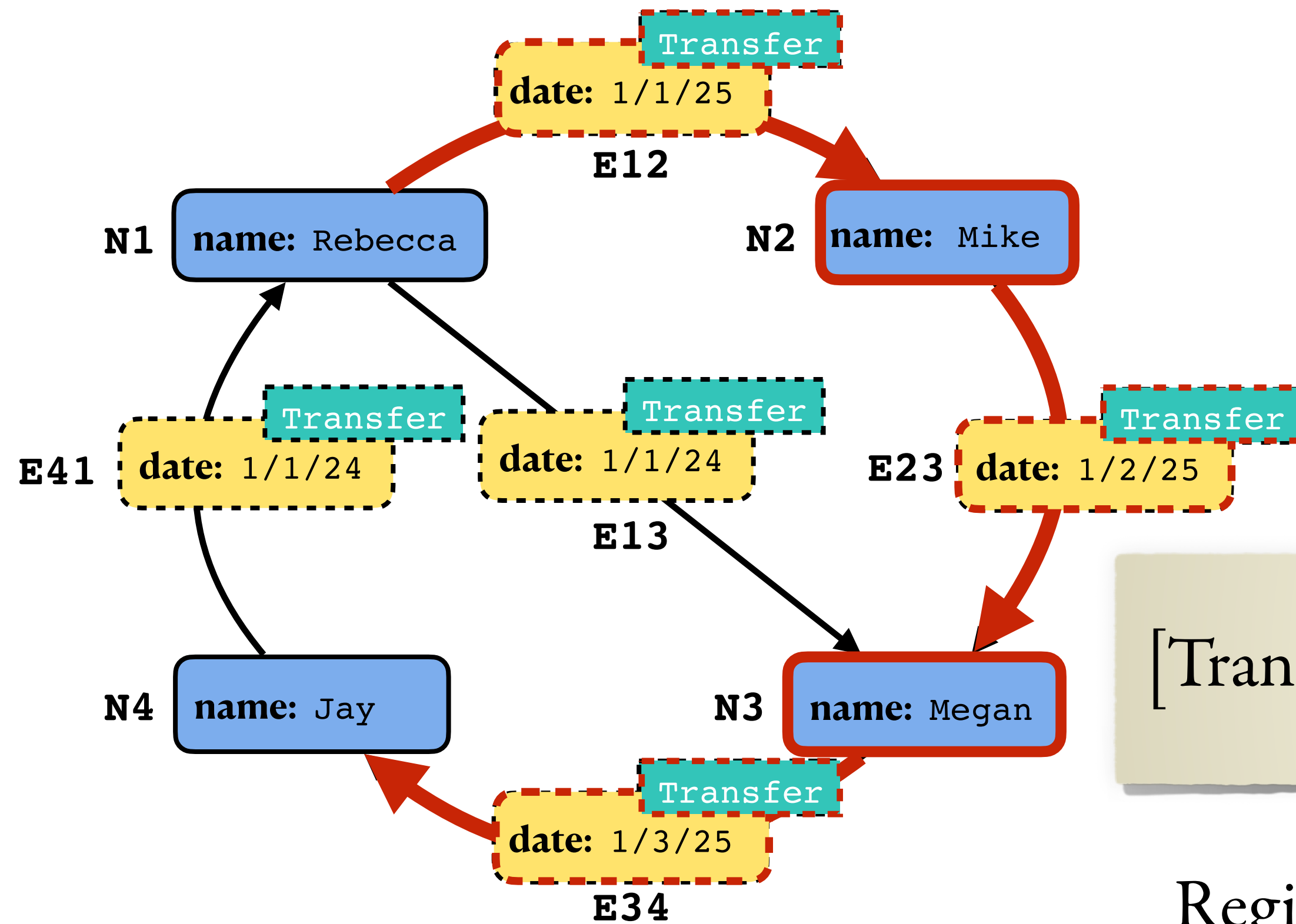**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{z}\right]\left[x := \text{date}\right]\left(\left(\_\right)\left[\text{Transfer}^{z}\right]\right)^{*}$$

Register $\quad x = 1/1/25$

Considered path: $\quad \text{E12} \longrightarrow \text{N2} \longrightarrow \text{E23}$
$\qquad\qquad\qquad\quad z \qquad\qquad\qquad z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^z\right]\left[x := \text{date}\right]\Big(\left(\_\right)\left[\text{Transfer}^z\right]\left[\text{date} > x\right]\Big)^*$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12 $\longrightarrow$ N2 $\longrightarrow$ E23
$$\qquad\qquad\qquad z \qquad\qquad\quad z$$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{z}\right]\left[x := \text{date}\right]\Big(\left(\_\right)\left[\text{Transfer}^{z}\right]\left[\text{date} > x\right]\right)^{*}$$

Register $\quad x \;=\; 1/1/25 \qquad$ date $= 1/2/25$ ✓

Considered path: $\quad$ E12 $\longrightarrow$ N2 $\longrightarrow$ E23
$$\qquad\qquad\qquad\quad z \qquad\qquad\quad z$$

# CRPQs with Data & List Variables

**Transfer**
**date:** 1/1/25
**E12**

**N1** **name:** Rebecca
**N2** **name:** Mike

**E41** **Transfer** **date:** 1/1/24
**Transfer** **date:** 1/1/24
**E13**

**Transfer** **date:** 1/2/25
**E23**

**N4** **name:** Jay
**N3** **name:** Megan

**Transfer** **date:** 1/3/25
**E34**

## We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{\,z}\right]\left[x := \text{date}\right]\Bigg(\left(\_\right)\left[\text{Transfer}^{\,z}\right]\left[\text{date} > x\right]\Bigg)^{*}$$

Register $\quad x = 1/1/25$

Considered path: $\quad \text{E12} \longrightarrow \text{N2} \text{---} \text{E23}$
$$\qquad\qquad\qquad z \qquad\qquad\qquad z$$

# CRPQs with Data & List Variables



**Transfer**
**date:** 1/1/25
**E12**

**N1** **name:** Rebecca

**N2** **name:** Mike

**Transfer**
**date:** 1/1/24
**E41**

**Transfer**
**date:** 1/1/24
**E13**

**Transfer**
**date:** 1/2/25
**E23**

**N4** **name:** Jay

**N3** **name:** Megan

**Transfer**
**date:** 1/3/25
**E34**

### We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$[\text{Transfer}^z]\,[x := \text{date}]\left((\_)\,[\text{Transfer}^z]\,[\text{date} > x]\,[x := \text{date}]\right)^*$$

Register $\quad x = 1/1/25$

Considered path: $\quad$ E12 $\longrightarrow$ N2 $\longrightarrow$ E23

$\qquad\qquad\qquad\qquad z \qquad\qquad\qquad z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{z}\right]\left[x := \text{date}\right]\Big((\_)\left[\text{Transfer}^{z}\right]\left[\text{date} > x\right]\left[x := \text{date}\right]\Big)^{*}$$

Register $\quad x = 1/2/25$

Considered path: $\quad$ E12 $\longrightarrow$ N2 $\longrightarrow$ E23
$\qquad\qquad\qquad\quad z \qquad\qquad\qquad z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^z\right]\left[x := \text{date}\right]\left(\left(\_\right)\left[\text{Transfer}^z\right]\left[\text{date} > x\right]\left[x := \text{date}\right]\right)^*$$

Register $\quad x = 1/2/25$

Considered path: $\quad \text{E12} \longrightarrow \text{N2} \longrightarrow \text{E23} \longrightarrow \text{N3}$
$$\qquad\qquad\qquad\quad z \qquad\qquad\qquad z$$

# CRPQs with Data & List Variables



We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\mathrm{Transfer}^{z}\right]\left[x := \mathrm{date}\right]\left(\left(\_\right)\left[\mathrm{Transfer}^{z}\right]\left[\mathrm{date} > x\right]\left[x := \mathrm{date}\right]\right)^{*}$$

Register    $x = 1/2/25$

Considered path:    $\mathrm{E12} \longrightarrow \mathrm{N2} \longrightarrow \mathrm{E23} \longrightarrow \mathrm{N3} \longrightarrow \mathrm{E34}$
                         $z$                    $z$                    $z$

# CRPQs with Data & List Variables



N1 **name:** Rebecca
N2 **name:** Mike
N3 **name:** Megan
N4 **name:** Jay

E12 Transfer **date:** 1/1/25
E41 Transfer **date:** 1/1/24
E13 Transfer **date:** 1/1/24
E23 Transfer **date:** 1/2/25
E34 Transfer **date:** 1/3/25

### We're Going to Add
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^{\,z}\right]\left[x := \text{date}\right]\Big(\left(\_\right)\left[\text{Transfer}^{\,z}\right]\left[\text{date} > x\right]\left[x := \text{date}\right]\Big)^{*}$$

Register    $x \ = \ 1/3/25$

Considered path:   $E12 \longrightarrow N2 \longrightarrow E23 \longrightarrow N3 \longrightarrow E34$
                        $z$                    $z$                    $z$

# CRPQs with Data & List Variables



**We're Going to Add**
- node & edge treatment
- data filters

Increasing date values on edges

$$\left[\text{Transfer}^z\right]\left[x := \text{date}\right]\left(\left(\_\right)\left[\text{Transfer}^z\right]\left[\text{date} > x\right]\left[x := \text{date}\right]\right)^*$$

Register $\quad x = 1/3/25$

paths don't need to be node-to-node (as in GQL)

Considered path: $\quad$ E12 $\longrightarrow$ N2 $\longrightarrow$ E23 $\longrightarrow$ N3 $\longrightarrow$ E34

$\qquad\qquad\qquad\qquad z \qquad\qquad\qquad z \qquad\qquad\qquad z$

# CRPQs with Data & List Variables

Increasing date values on **nodes**

$$\left(\text{Transfer}^z\right) \left(x := \text{date}\right) \left( \underline{\phantom{[]}} \left(\text{Transfer}^z\right) \left(\text{date} > x\right) \left(x := \text{date}\right) \right)*$$

node-to-node
paths

Symmetry!

$$\left[\text{Transfer}^z\right] \left[x := \text{date}\right] \left( \left(\underline{\phantom{}}\right) \left[\text{Transfer}^z\right] \left[\text{date} > x\right] \left[x := \text{date}\right] \right)*$$

edge-to-edge
paths

Increasing date values on **edges**

# CRPQs with Data & List Variables

# CRPQs with Data & List Variables

Increasing date values on **nodes**

$$\left(\text{Transfer}^z\right)\ (x := \text{date})\ \Big(\ [\_]\ \left(\text{Transfer}^z\right)\ (\text{date} > x)\ (x := \text{date})\ \Big) *$$

node-to-node paths

Symmetry!

$$[\text{Transfer}^z]\ [x := \text{date}]\ \Big(\ (\_)\ [\text{Transfer}^z]\ [\text{date} > x]\ [x := \text{date}]\ \Big) *$$

edge-to-edge paths

Increasing date values on **edges**

# Distilling From Practice

# Distilling From Practice

# CoreGQL

$$\pi \quad := \quad (x) \qquad \text{node}$$

$$| \quad \xrightarrow{x} \qquad \text{edge}$$

$$| \quad \pi_1 \, \pi_2 \qquad \text{concatenation}$$

$$| \quad \pi_1 + \pi_2 \qquad \text{disjunction}$$

$$| \quad \pi^{n..m} \qquad \text{repetition}$$

$$| \quad \pi \langle \theta \rangle \qquad \text{condition}$$

$$\theta, \theta' \ := \ x.k = x'.k' \mid x.k < x'.k' \mid \ell(x) \mid \theta \vee \theta' \mid \theta \wedge \theta' \mid \neg\theta$$

More?  ⇝ Leonid @ GRADES/NDA

Friday 2PM

# Results About Distilled Models

**Theorem**        [Gheerbrant, Libkin, Peterfreund, Rogova ICDT'25]

The RPQ $(aa)^*$ is not expressible using Cypher patterns

**Theorem**        [Gheerbrant, Libkin, Peterfreund, Rogova PVLDB'25]

"Increasing values on edges" cannot be expressed without repeating variables

...

# Research Agenda
## for Graph Query Languages

# Research Agenda for Graph Query Languages

## Growing from Theory

- Design and study elegant models
- Understand expressiveness
- Understand complexity
- Use our knowledge of logics, automata,...
- ...

## Distilling from Practice

- Identify what's good
- Find opportunities for improvement
- Inexpressibility results
- Complexity lower bounds
- ...

# Research Agenda for Graph Query Languages

## Growing from Theory

- Design and study elegant models
- Understand expressiveness
- Understand complexity
- Use our knowledge of logics, automata,…
- …

## Distilling from Practice

- Identify what's good
- Find opportunities for improvement
- Inexpressibility results
- Complexity lower bounds
- …



Light side

Dark side

# Research Agenda for Graph Query Languages

## Growing from Theory

- Design and study elegant models
- Understand expressiveness
- Understand complexity
- Use our knowledge of logics, automata,…
- …

## Distilling from Practice

- Identify what's good
- Find opportunities for improvement
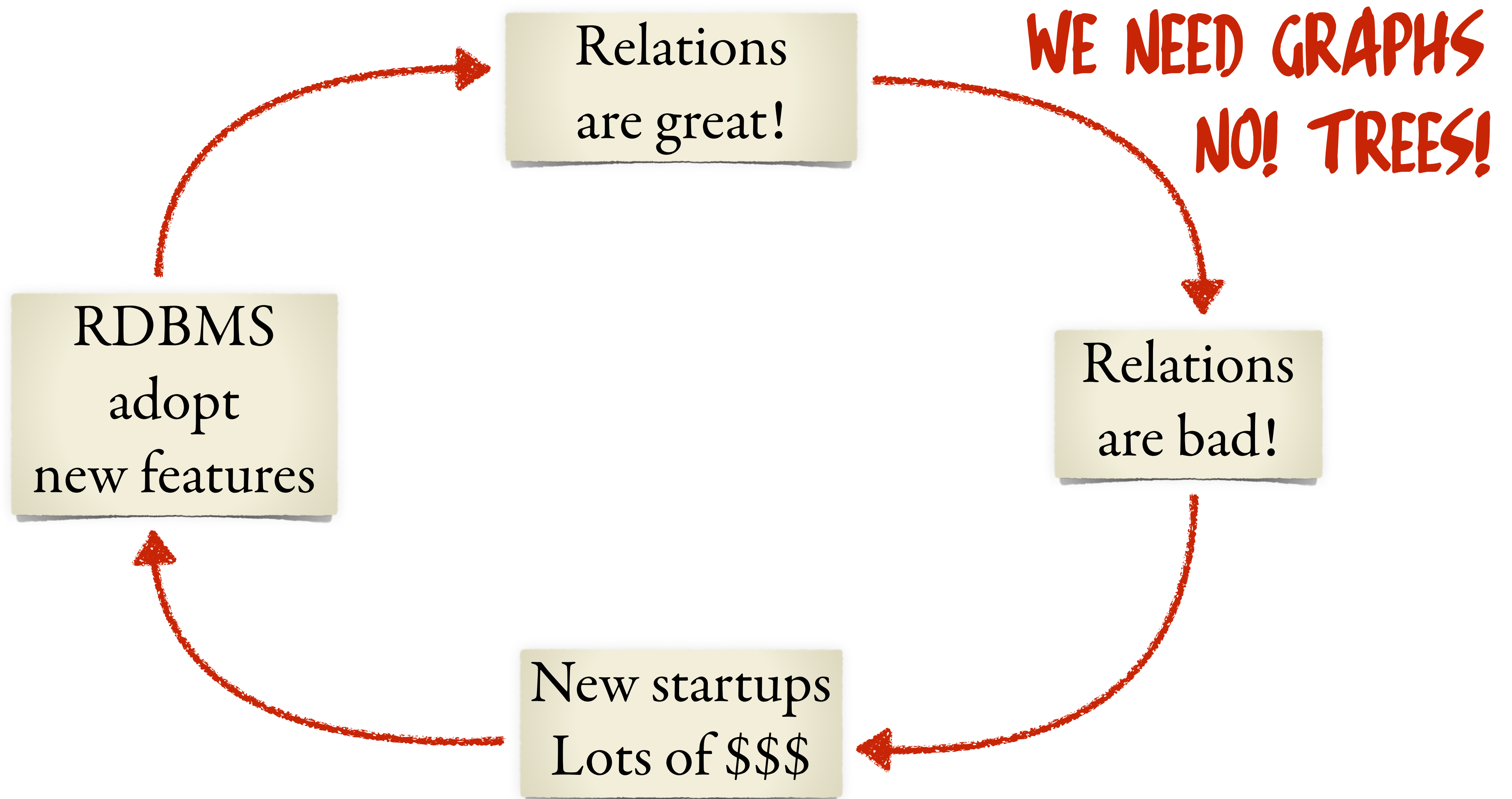- Inexpressibility results
- Complexity lower bounds
- …

Section 7

# Graphs vs Relations
## What Goes Around Comes Around?

Big Questions Again?

Sets?   Bags?

# What Goes Around Comes Around



Slide copied & adapted from Andy Pavlo

# Graph Query Languages

> **Observation**
>
> For the query language,
> the underlying DB architecture
> (graph, relational,...)
> is irrelevant

Cypher / GQL          ⤳ implemented in   **neo4j**          graph native

                      ⤳ implemented in   **kùzu**           graph native

SQL/PGQ               ⤳ implemented in   **ORACLE**         relational

                      ⤳ implemented in   **DuckDB**         relational

└→ can be             ⤳ implemented in   **RelationalAI**   relational
   translated                                               + set semantics!
   into Rel

# Wrapping Up

# Wrapping Up

**New graph query languages add**

(a) handling of nodes and edges
(b) path and list variables
(c) path modes
(d) data filters
to Conjunctive Regular Path Queries

Their design in the standard(s) isn't smooth yet
⤳ we have work to do

How would we design features like this
in a graph query language?

Our proposal:
 - l-CRPQs
 - dl-CRPQs
(see paper)

**What can you do?**

- Study these new CRPQs
- Come up with your own design?
- Study GQL & SQL/PGQ
- Do RPQs in Datalog
- Prove that sets are better than bags
- Prove that bags are better than sets
- Solve Automatic Programming
- ...

# Thank you!

Happy to talk to you in the break!