# Datalog$^{\pm}$: A Family of Languages for Ontology Querying

## Georg Gottlob

Department of Computer Science

University of Oxford

joint work with Andrea Calì,Thomas Lukasiewicz, Marco Manna, Andreas Pieris et al.

# 1st Motivation: Ontological Query Answering

– Input: a knowledge base $K = \langle D, \Sigma \rangle$

   a Boolean conjunctive query $Q$

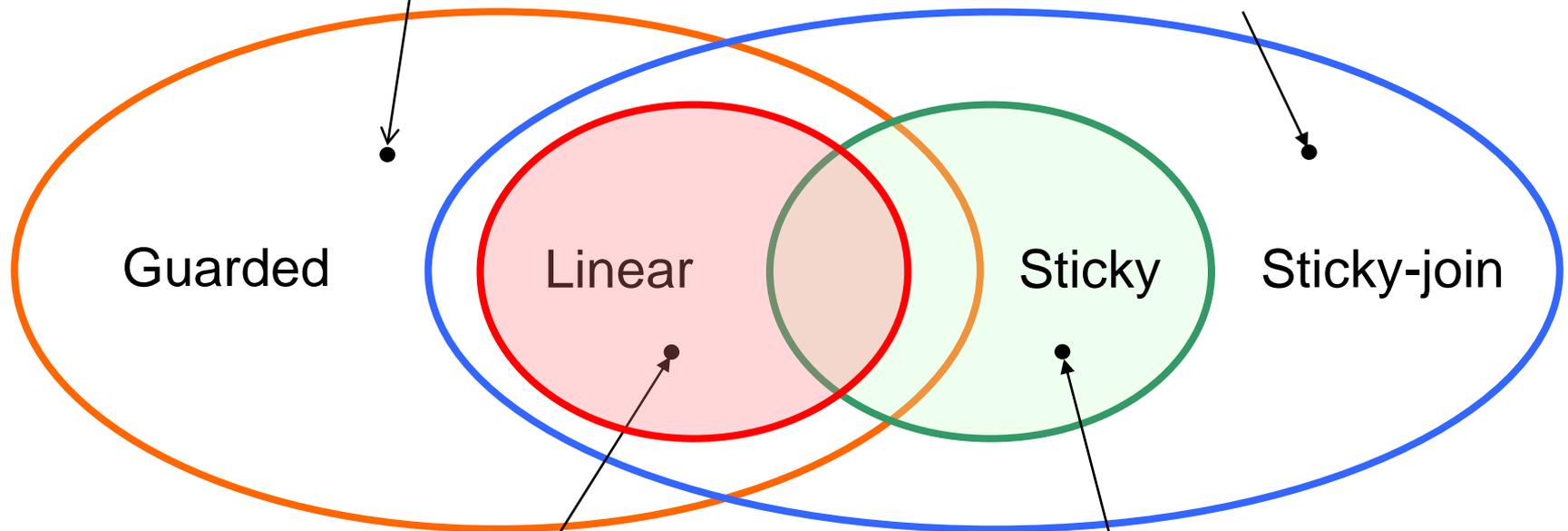   Question: $K \vDash Q$ (or, equivalently, $D \cup \Sigma \vDash Q$)

– Old database problem – CQ answering over incomplete databases

# 2nd Motivation: Decidable FO fragments (TGDs) Good data complexity!

$\forall X \forall Y \forall Z\ R(X,Y,Z), S(Y), P(X,Z) \rightarrow \exists W\ Q(X,W)$

$\forall X \forall Y \forall Z\ R(X,Y), S(Y,Z,Z) \rightarrow \exists W\ P(Y,W)$

Guarded     Linear     Sticky     Sticky-join

$\forall X \forall Y\ R(X,X,Y) \rightarrow \exists Z\ R(Y,Y,Z)$

$\forall X \forall Y \forall Z\ S(X,Y), R(Y,Z) \rightarrow \exists W\ P(Y,W)$

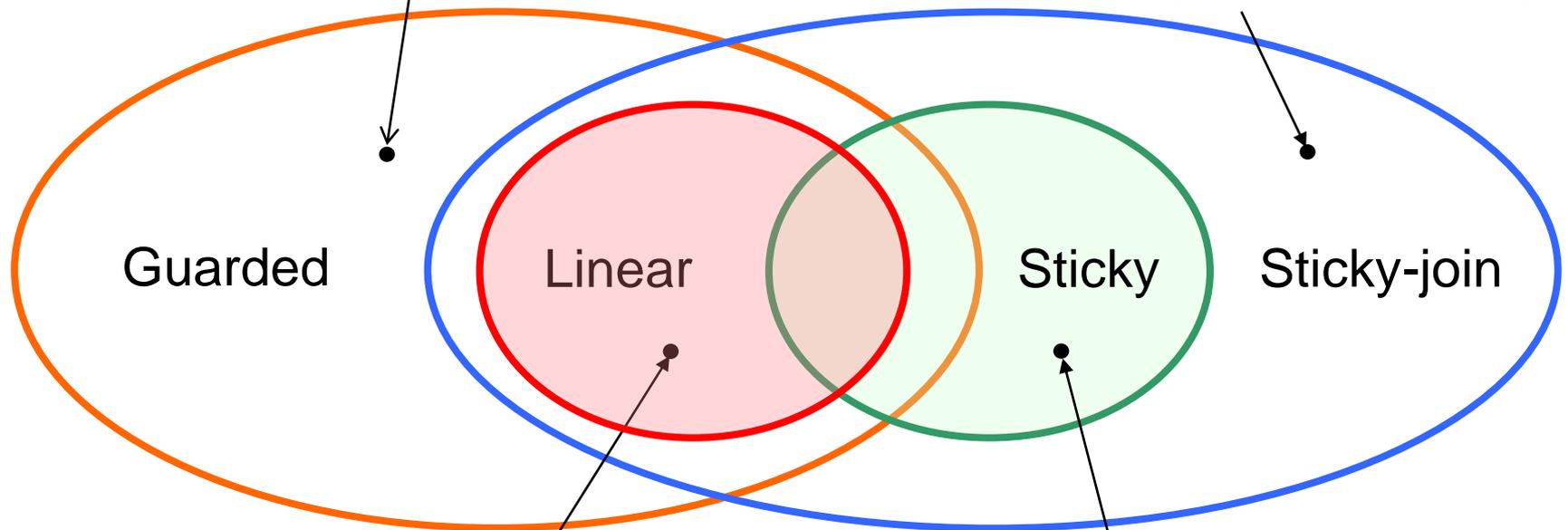# 2nd Motivation: Decidable FO fragments (TGDs)
## Good data complexity!

*Query: R(Y,Y,Z),P(Z,W)*

*Add rule:*
$\forall Y \forall Z \forall W \ R(Y,Y,Z),P(Z,W) \rightarrow \perp$

$\forall X \forall Y \forall Z \ R(X,Y,Z),S(Y),P(X,Z) \rightarrow \exists W \ Q(X,W)$

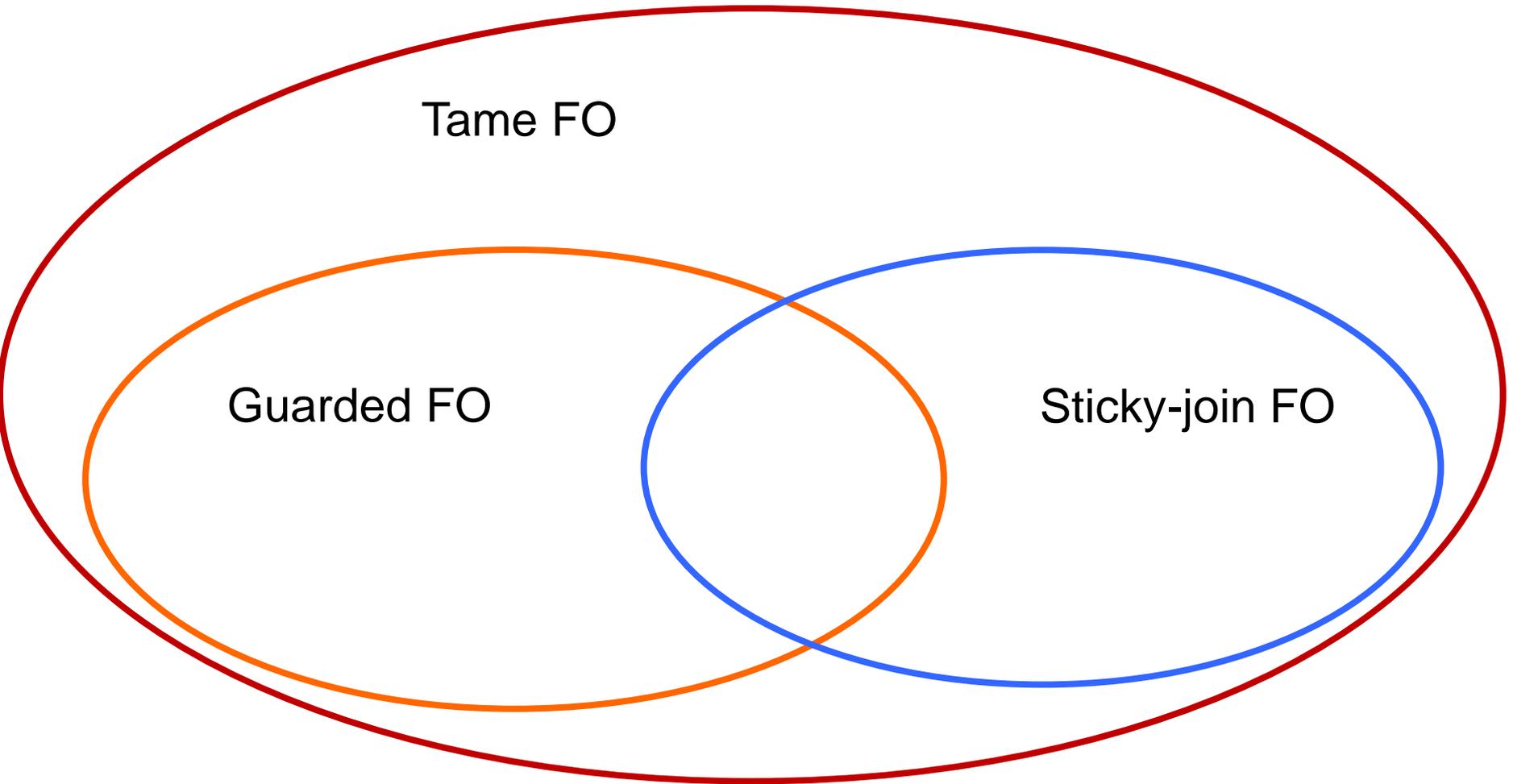$\forall X \forall Y \forall Z \ R(X,Y),S(Y,Z,Z) \rightarrow \exists W \ P(Y,W)$



Guarded

Linear

Sticky

Sticky-join

$\forall X \forall Y \ R(X,X,Y) \rightarrow \exists Z \ R(Y,Y,Z)$

$\forall X \forall Y \forall Z \ S(X,Y),R(Y,Z) \rightarrow \exists W \ P(Y,W)$
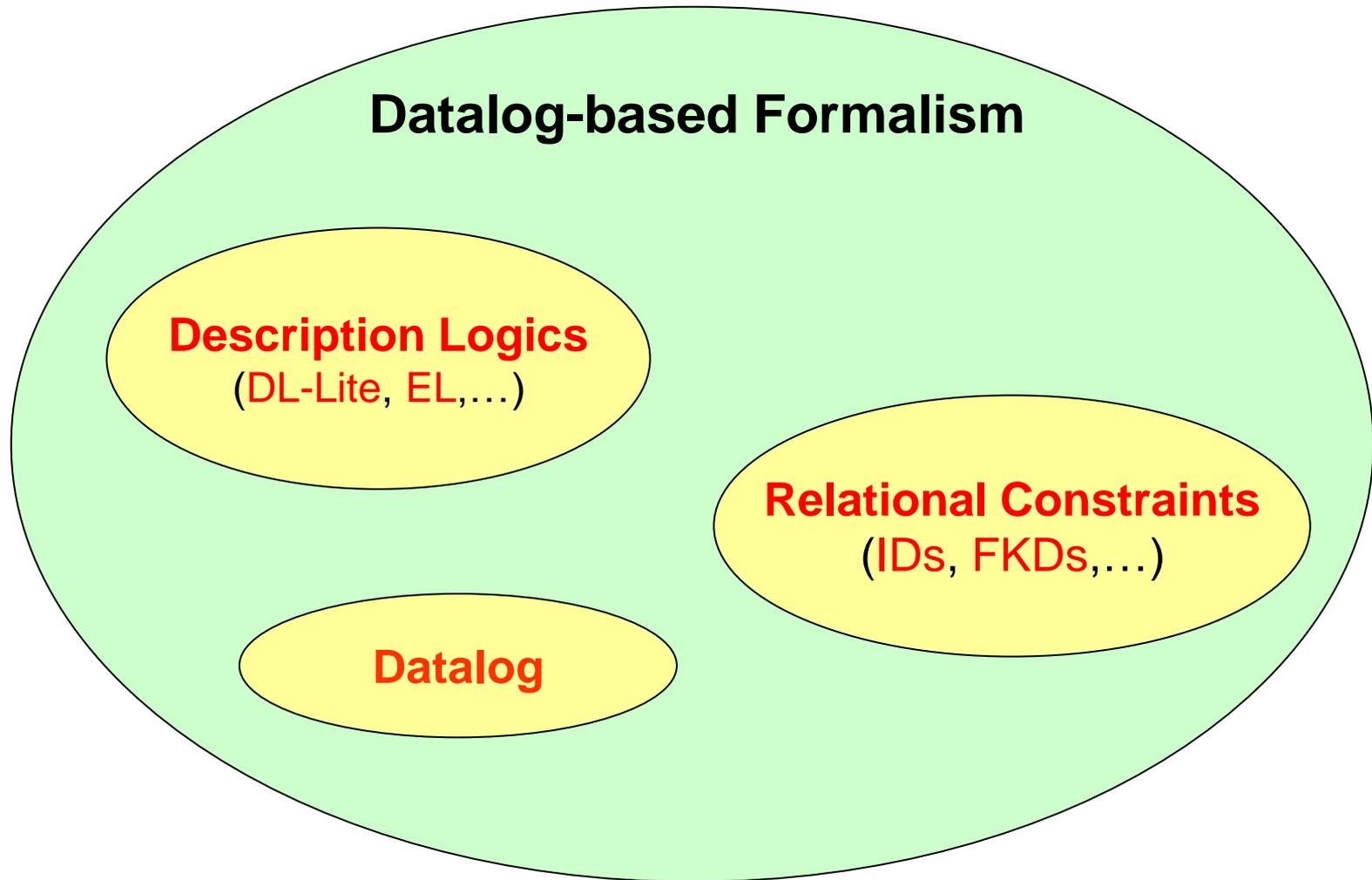
# 3rd Motivation: Decidable FO fragments

*(current work)*

# (Some) Known Results

| $\Sigma$ | Data Complexity | Combined Complexity |
|---|---|---|
| DL-Lite | in $AC_0$ [Calvanese et al., JAR 07] | NP-complete [Calvanese et al., JAR 07] |
| EL + ELH | PTIME-complete [Rosati, DL 07] | NP-complete [Rosati, DL 07] |
| Horn-SHIQ | PTIME-complete [Eiter et al., JELIA 08] | EXPTIME-complete [Eiter et al., JELIA 08] |
| IDs + FKDs | in $AC_0$ [Calì et al., IJCAI 03] | PSPACE-complete [Johnson & Klug, JCSS 84] |
| Datalog | PTIME-complete | EXPTIME-complete |

# Our Goal



**Datalog-based Formalism**

**Description Logics**
(DL-Lite, EL,…)

**Relational Constraints**
(IDs, FKDs,…)

**Datalog**

… without losing tractable data complexity

# Ontological Reasoning and Datalog

| DL Assertion | Datalog Rule |
|---|---|
| **Concept Inclusion**<br>$emp \sqsubseteq person$ | $emp(X) \rightarrow person(X)$ |
| **Concept Product**<br>$sen\text{-}emp \times emp \sqsubseteq moreThan$ | $sen\text{-}emp(X), emp(Y) \rightarrow moreThan(X, Y)$ |
| **(Inverse) Role Inclusion**<br>$reports^- \sqsubseteq mgr$ | $reports(X, Y) \rightarrow mgr(Y, X)$ |
| **Role Transitivity**<br>$trans(mgr)$ | $mgr(X, Y), mgr(Y, Z) \rightarrow mgr(X, Z)$ |
|  |  |

# Ontological Reasoning and Datalog

| DL Assertion | Datalog Rule |
|---|---|
| **Concept Inclusion**<br>$emp \sqsubseteq person$ | $emp(X) \rightarrow person(X)$ |
| **Concept Product**<br>$sen\text{-}emp \times emp \sqsubseteq moreThan$ | $sen\text{-}emp(X), emp(Y) \rightarrow moreThan(X, Y)$ |
| **(Inverse) Role Inclusion**<br>$reports^- \sqsubseteq mgr$ | $reports(X, Y) \rightarrow mgr(Y, X)$ |
| **Role Transitivity**<br>$\text{trans}(mgr)$ | $mgr(X, Y), mgr(Y, Z) \rightarrow mgr(X, Z)$ |
| **Participation**<br>$emp \sqsubseteq \exists report$ | $emp(X) \rightarrow \exists Y\ report(X, Y)$ |
| **Disjointness**<br>$emp \sqcap customer \sqsubseteq \bot$ | $emp(X), customer(X) \rightarrow \bot$ |
| **Functionality**<br>$\text{funct}(reports)$ | $reports(X, Y), reports(X, Z) \rightarrow Y = Z$ |

# Datalog$^\pm$

- Datalog: $\forall \mathbf{X} \forall \mathbf{Y}\ \Phi(\mathbf{X},\mathbf{Y}) \rightarrow R(\mathbf{X})$

- Extend Datalog by allowing in the head:

  - Existential quantification ($\exists$) - TGDs: $\forall \mathbf{X} \forall \mathbf{Y}\ \Phi(\mathbf{X},\mathbf{Y}) \rightarrow \exists \mathbf{Z}\ \Psi(\mathbf{X},\mathbf{Z})$

  - Equality atoms ($=$) - EGDs: $\forall \mathbf{X}\ \Phi(\mathbf{X}) \rightarrow X_i = X_j$

  - Constant false ($\bot$) - Negative constraints: $\forall \mathbf{X}\ \Phi(\mathbf{X}) \rightarrow \bot$

- But query answering under Datalog[$\exists$] is undecidable
  [see, e.g., Beeri & Vardi, ICALP 81]

- Datalog[$\exists$,$=$,$\bot$] is syntactically restricted $\rightarrow$ **Datalog$^\pm$**

# The Chase Procedure

Input: Database *D*, set of TGDs $\Sigma$
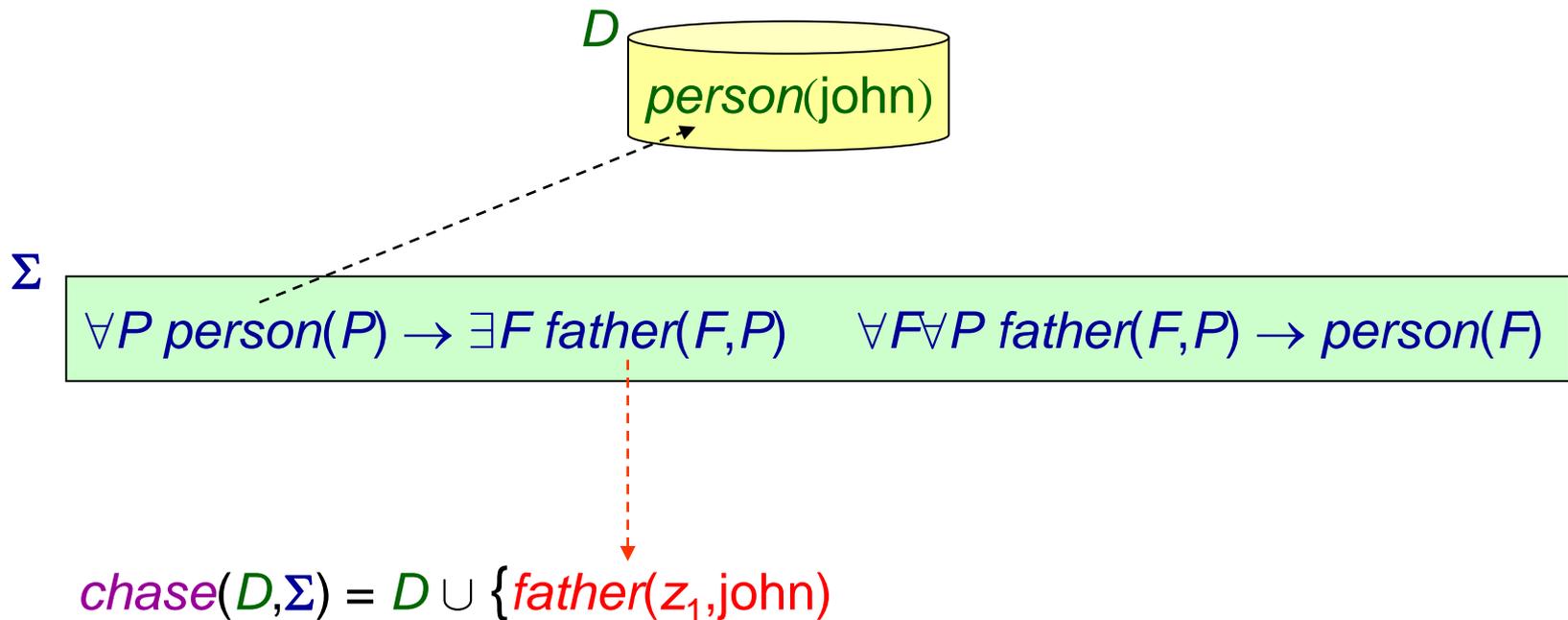
Output: A model of $D \cup \Sigma$

*D*
*person*(john)

$\Sigma$

$\forall P \; person(P) \rightarrow \exists F \; father(F,P) \qquad \forall F \forall P \; father(F,P) \rightarrow person(F)$

*chase*($D,\Sigma$) = $D \cup$ ?

# The Chase Procedure

Input: Database $D$, set of TGDs $\Sigma$

Output: A model of $D \cup \Sigma$

$D$

*person*(john)

$\Sigma$

$\forall P\ person(P) \rightarrow \exists F\ father(F,P)$    $\forall F \forall P\ father(F,P) \rightarrow person(F)$

$chase(D,\Sigma) = D \cup \{father(z_1,john)$
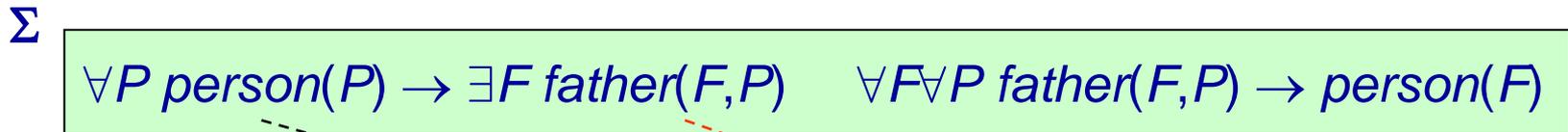
# The Chase Procedure

Input: Database $D$, set of TGDs $\Sigma$

Output: A model of $D \cup \Sigma$

$D$

person(john)

$\Sigma$

$\forall P\ person(P) \rightarrow \exists F\ father(F,P) \qquad \forall F\forall P\ father(F,P) \rightarrow person(F)$

$chase(D,\Sigma) = D \cup \{father(z_1,john),\ person(z_1)$

# The Chase Procedure

Input: Database $D$, set of TGDs $\Sigma$

Output: A model of $D \cup \Sigma$

$D$

person(john)

$\Sigma$

$\forall P\ person(P) \rightarrow \exists F\ father(F,P) \qquad \forall F \forall P\ father(F,P) \rightarrow person(F)$

$chase(D,\Sigma) = D \cup \{father(z_1,john),\ person(z_1),\ father(z_2,z_1)$

# The Chase Procedure

Input: Database $D$, set of TGDs $\Sigma$

Output: A model of $D \cup \Sigma$

$D$

$person(\text{john})$

$\Sigma$

$\forall P\ person(P) \rightarrow \exists F\ father(F,P)$    $\forall F \forall P\ father(F,P) \rightarrow person(F)$

$chase(D,\Sigma) = D \cup \{father(z_1,\text{john}),\ person(z_1),\ father(z_2,z_1),\ \ldots\}$

# The Chase Procedure

Input: Database $D$, set of TGDs $\Sigma$
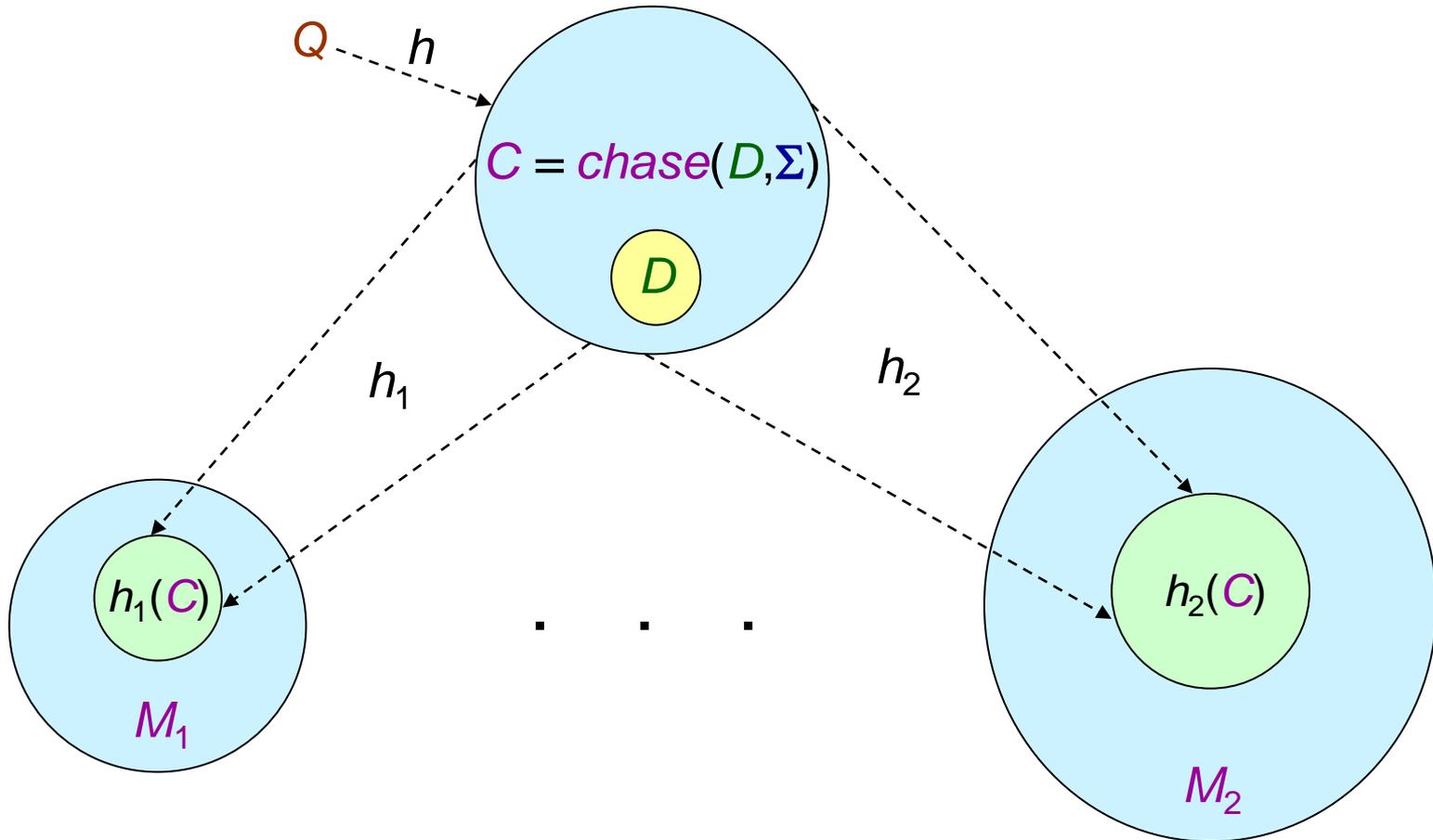
Output: A model of $D \cup \Sigma$

$D$

person(john)

$\Sigma$

$\forall P\ person(P) \rightarrow \exists F\ father(F,P)$     $\forall F \forall P\ father(F,P) \rightarrow person(F)$

$chase(D,\Sigma) = D \cup \{father(z_1,john),\ person(z_1),\ father(z_2,z_1),\ \ldots\}$

infinite instance

# Query Answering via Chase



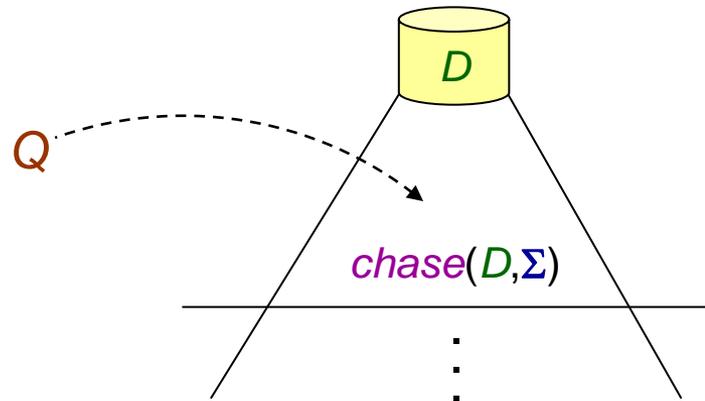$$D \cup \Sigma \vDash Q \quad \Leftrightarrow \quad chase(D,\Sigma) \vDash Q$$

[see, e.g., Fagin, Kolaitis, Miller & Popa, TCS 05]

# Early Positive Result

– Query answering under IDs is decidable     $P(X,Y) \rightarrow \exists Z\ R(X,Y,Z)$

    - PSPACE-complete in combined complexity

    - NP-complete for bounded arity



$chase(D, \Sigma)$

[Johnson & Klug, JCSS 84]

# Guardedness

- All $\forall$-variables occur in one body atom - guard atom

$$\forall X \forall Y \forall Z\ R(X,Y,Z), S(Y), P(X,Z) \rightarrow \exists W\ Q(X,W)$$

guard

- Chase has finite treewidth $\Rightarrow$ decidability of query answering

  [Calì, G. & Kifer, KR 08]

- Query answering is PTIME-complete in data complexity

  [Calì, G. & Lukasiewicz, PODS 09]

- Properly extends ELH (same data complexity)

# Ontology Querying

ELH: Popular DL (for biological applications) with PTIME data complexity
[Baader, IJCAI 03 and Rosati, DL 07]

| ELH TBox | Datalog$^\pm$ Representation |
|---|---|
| $A \sqsubseteq B$ | $\forall X\, A(X) \rightarrow B(X)$ |
| $A \sqcap B \sqsubseteq C$ | $\forall X\, A(X), B(X) \rightarrow C(X)$ |
| $\exists R.A \sqsubseteq B$ | $\forall X\, R(X,Y), A(Y) \rightarrow B(X)$ |
| $A \sqsubseteq \exists R.B$ | $\forall X\, A(X) \rightarrow \exists Y\, R(X,Y), B(Y)$ |
| $R \sqsubseteq P$ | $\forall X \forall Y\, R(X,Y) \rightarrow P(X,Y)$ |

# Linearity

– Just one atom in the body    $\forall \mathbf{X} \forall \mathbf{Y} \, R(\mathbf{X},\mathbf{Y}) \rightarrow \exists \mathbf{Z} \, \Psi(\mathbf{X},\mathbf{Z})$
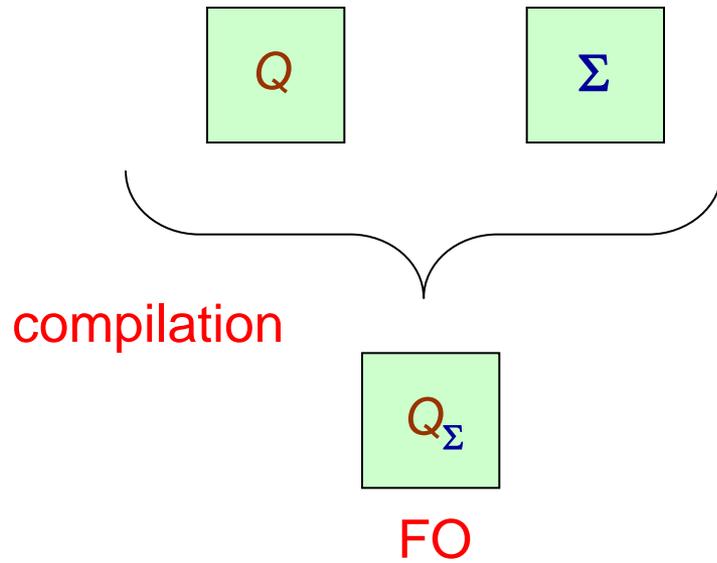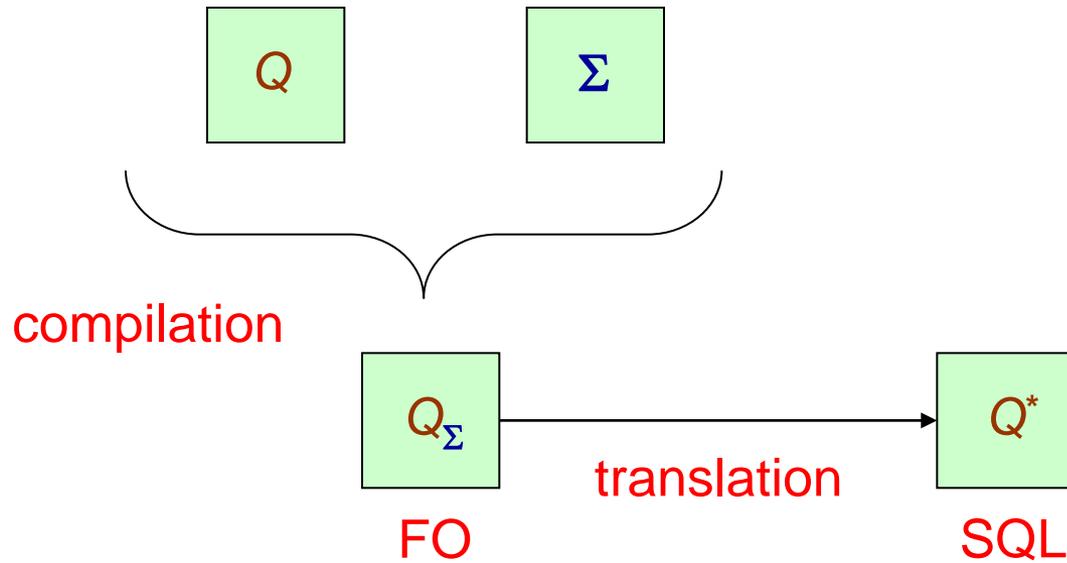
guard

– Linear TGDs are trivially guarded

– Query answering is in $AC_0$ in data complexity (first-order rewritability)
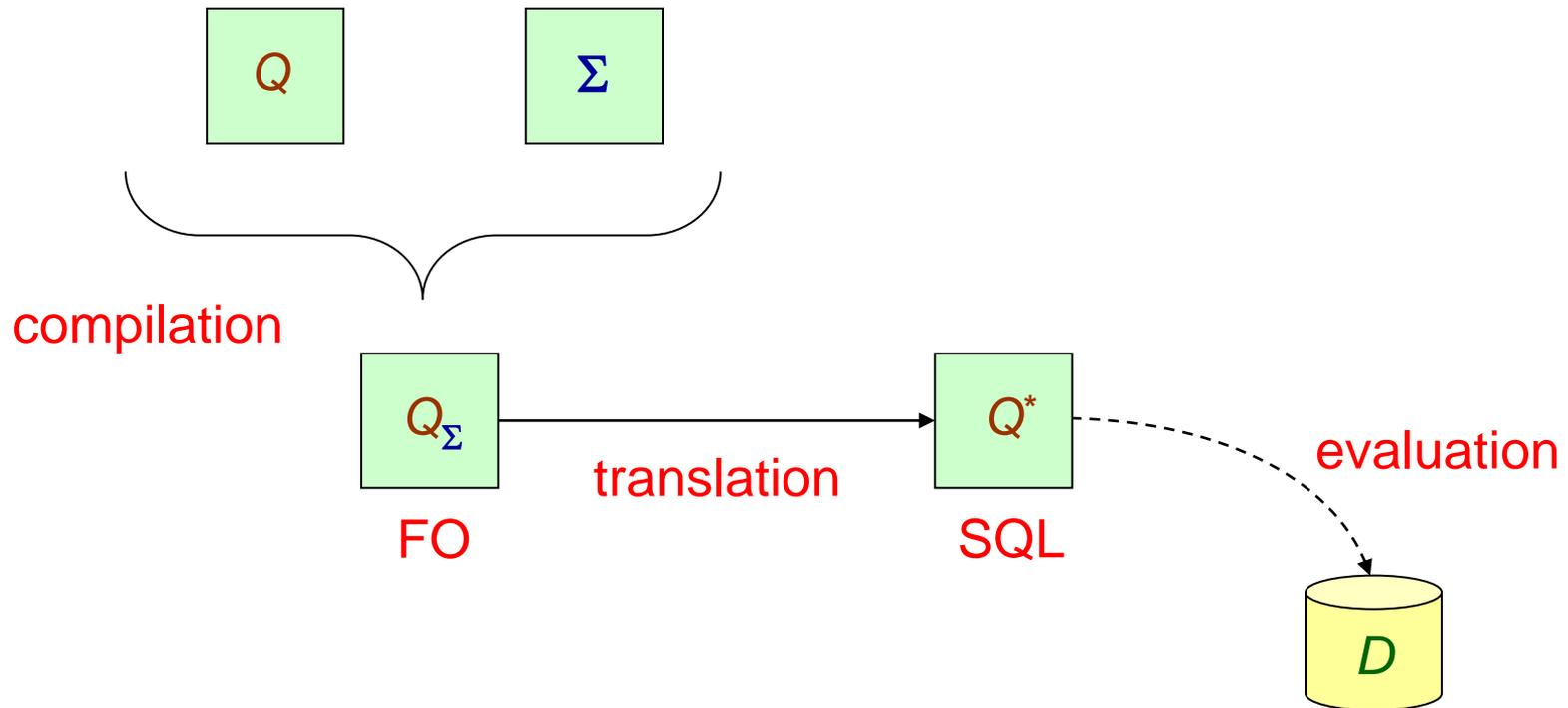
[Calì, G. & Lukasiewicz, PODS 09]

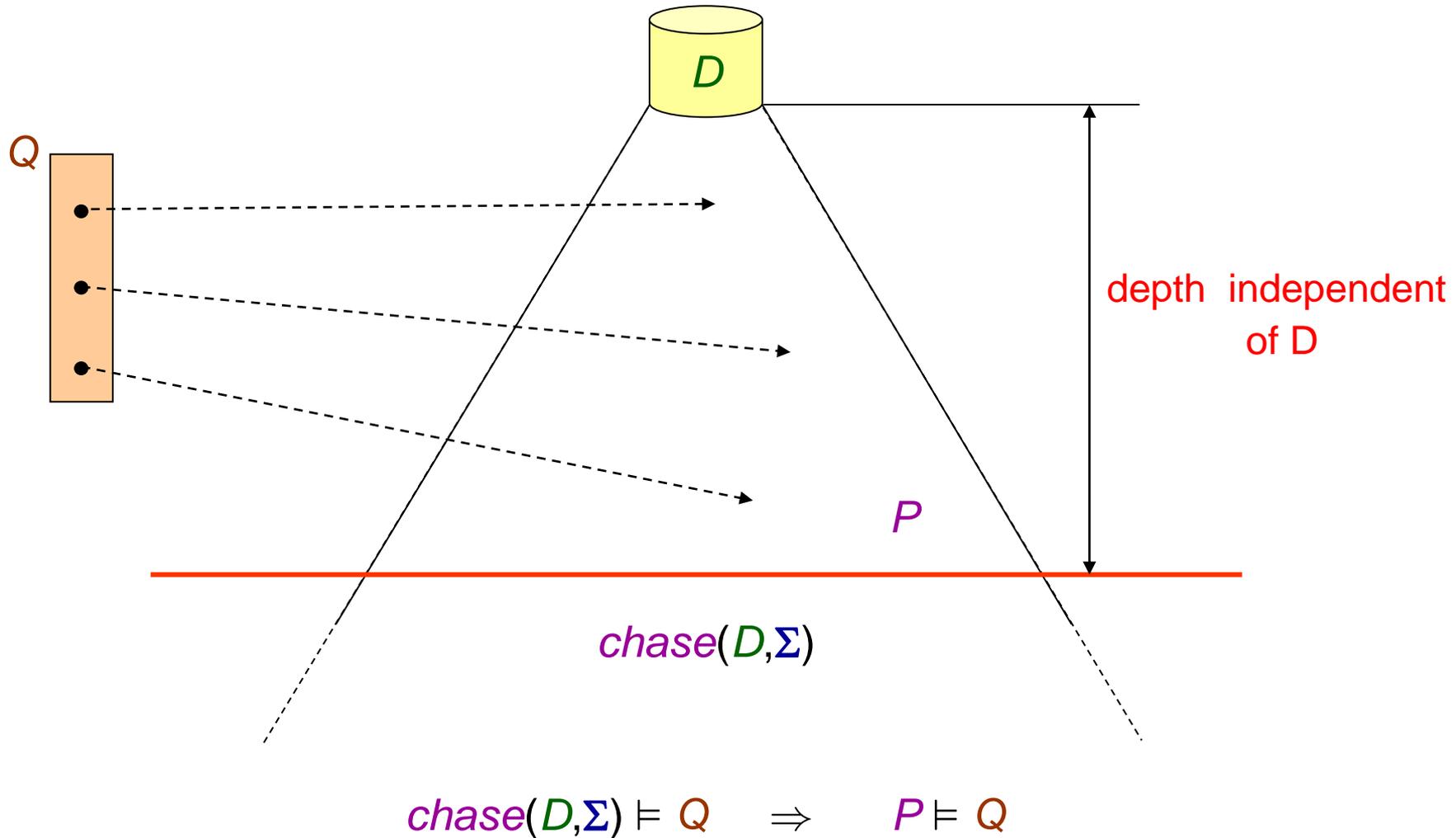# First-Order Rewritable TGDs

# First-Order Rewritable TGDs

# First-Order Rewritable TGDs



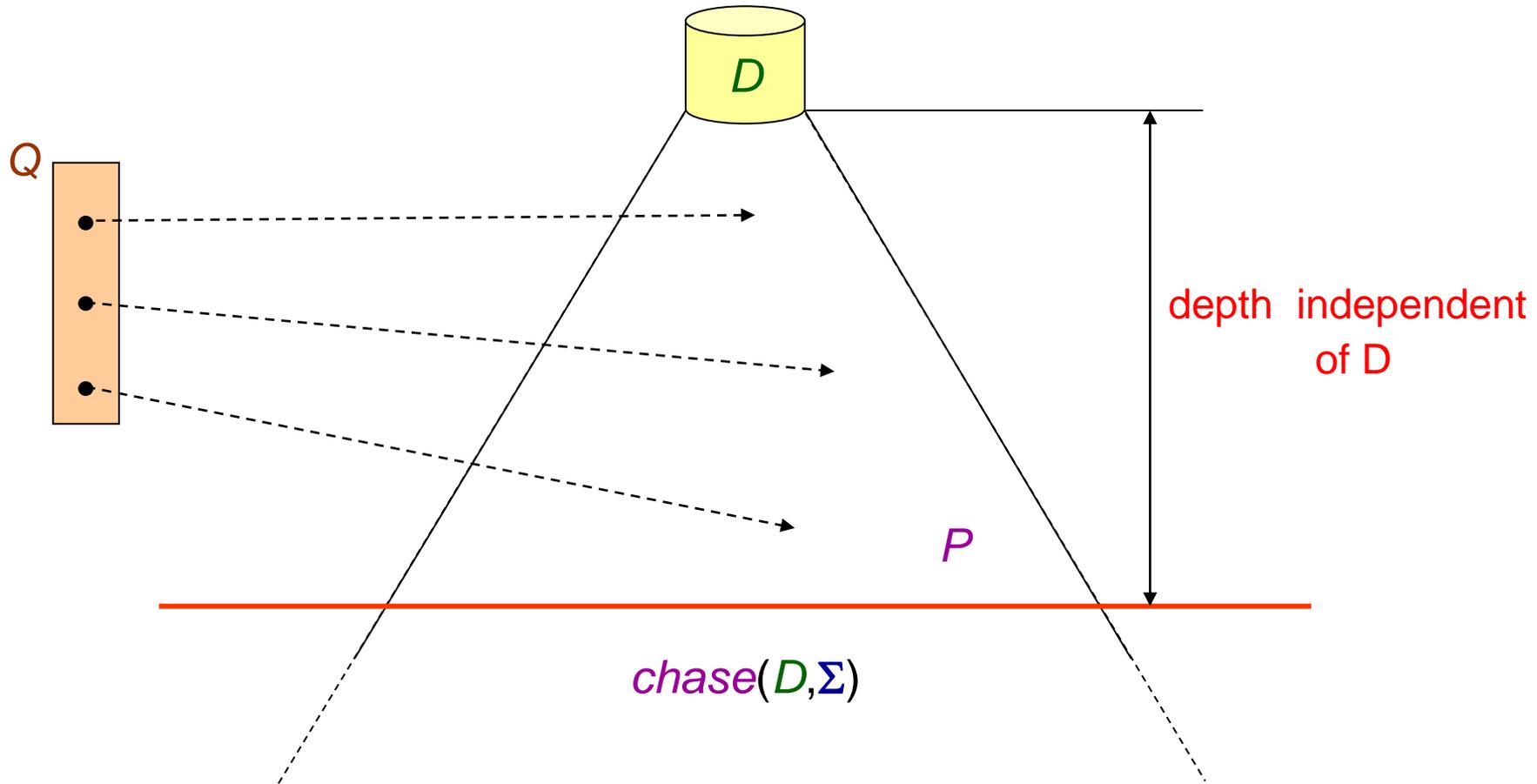$$\forall D \, (D \cup \Sigma \vDash Q \iff D \vDash Q^*)$$

Query answering is in $AC_0$ in data complexity  [Vardi, PODS 95]

# Bounded Derivation-Depth Property (BDDP)



$$chase(D,\Sigma) \vDash Q \quad \Rightarrow \quad P \vDash Q$$

[Calì, G. & Lukasiewicz, PODS 09]

# Bounded Derivation-Depth Property (BDDP)



BDDP $\Rightarrow$ First-Order Rewritability

[Calì, G. & Lukasiewicz, PODS 09]

# Linearity

– Just one atom in the body $\quad \forall \mathbf{X} \forall \mathbf{Y}\ R(\mathbf{X},\mathbf{Y}) \to \exists \mathbf{Z}\ \Psi(\mathbf{X},\mathbf{Z})$

trivially guard

– Linear TGDs are trivially guarded

– Query answering is in $AC_0$ in data complexity (first-order rewritability)

[Calì, Gottlob & Lukasiewicz, PODS 09]

– Properly extends DL-Lite (same data complexity)

# Ontology Querying

DL-Lite: Popular family of DLs with $AC_0$ data complexity (OWL 2 QL)
[Calvanese, De Giacomo, Lembo, Lenzerini & Rosati, JAR 07]

| DL-Lite TBox | Datalog$^\pm$ Representation |
|---|---|
| $A \sqsubseteq B$ | $\forall X\, A(X) \rightarrow B(X)$ |
| $A \sqsubseteq \exists R$ | $\forall X\, A(X) \rightarrow \exists Y\, R(X, Y)$ |
| $\exists R \sqsubseteq A$ | $\forall X \forall Y\, R(X, Y) \rightarrow A(X)$ |
| $R \sqsubseteq P$ | $\forall X \forall Y\, R(X, Y) \rightarrow P(X, Y)$ |

# But…

- What about joins in rule bodies?

$\forall A \forall D \forall P \ runs(D,P), area(P,A) \rightarrow \exists E \ employee(E,D,P,A)$

- What about the DL assertion concept product?

$\forall E \forall M \ elephant(E), mouse(M) \rightarrow biggerThan(E,M)$

# But…

– What about <span style="color:red">joins</span> in rule bodies?

$\forall A \forall D \forall P\ runs(D,P), area(P,A) \to \exists E\ employee(E,D,P,A)$

– What about the DL assertion <span style="color:red">concept product</span>?

$\forall E \forall M\ elephant(E), mouse(M) \to biggerThan(E,M)$

<span style="color:red">No tree-like</span> models guaranteed

$\forall X \forall Y\ R(X,Y) \to \exists Z\ R(Y,Z)$

$\forall X \forall Y\ R(X,Y) \to S(X)$

<span style="color:red">Infinitely</span> many symbols in $S$

$\forall X \forall Y\ S(X), S(Y) \to P(X,Y)$
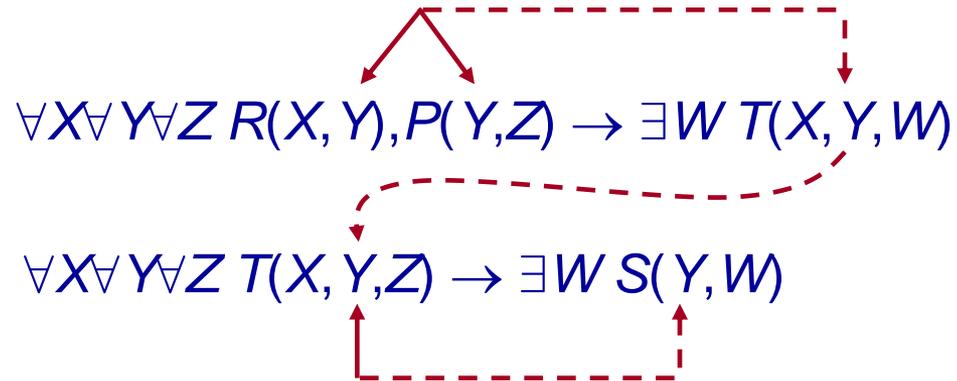
$P$ forms an <span style="color:red">infinite clique</span>

# Stickiness

$\forall A \forall D \forall P\ runs(D,P), area(P,A) \rightarrow \exists E\ employee(E,D,P,A)$ ✔

$\forall E \forall M\ elephant(E), mouse(M) \rightarrow biggerThan(E,M)$ ✔

[Calì, G. & Pieris, VLDB 10]

# Stickiness

$$\forall X \forall Y \forall Z \, R(X,Y), P(Y,Z) \rightarrow \exists W \, T(X,Y,W)$$

$$\forall X \forall Y \forall Z \, T(X,Y,Z) \rightarrow \exists W \, S(Y,W)$$

✔

[Calì, G. & Pieris, VLDB 10]

# Stickiness

$$\forall X \forall Y \forall Z \; R(X,Y), P(Y,Z) \rightarrow \exists W \; T(X,Y,W)$$

$$\forall X \forall Y \forall Z \; T(X,Y,Z) \rightarrow \exists W \; S(Y,W)$$

✔

$$\forall X \forall Y \forall Z \; R(X,Y), P(Y,Z) \rightarrow \exists W \; T(X,Y,W)$$

$$\forall X \forall Y \forall Z \; T(X,Y,Z) \rightarrow \exists W \; S(X,W)$$

✖

[Calì, G. & Pieris, VLDB 10]

# Stickiness: Marking Procedure

**Initial Marking**: mark all occurrences of body-variables that do not appear in <span style="color:red">all</span> head-atoms

$$\forall V \forall W \, R_1(V,W) \;\rightarrow\; \exists X \exists Y \exists Z \, R_2(W,X,Y,Z)$$
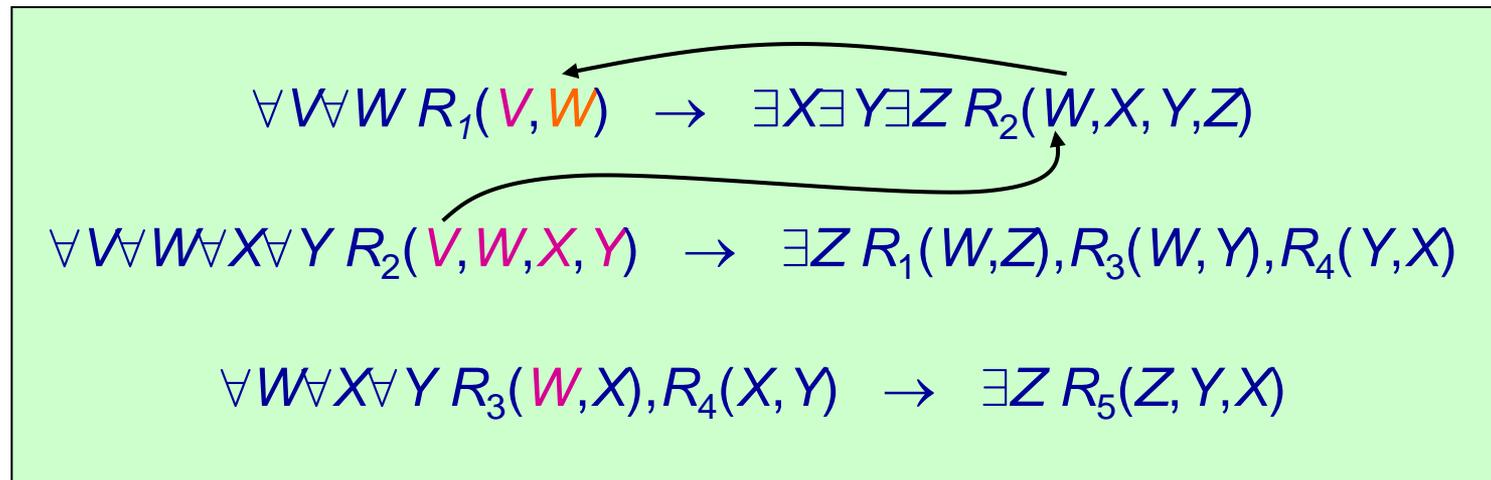
$$\forall V \forall W \forall X \forall Y \, R_2(V,W,X,Y) \;\rightarrow\; \exists Z \, R_1(W,Z), R_3(W,Y), R_4(Y,X)$$

$$\forall W \forall X \forall Y \, R_3(W,X), R_4(X,Y) \;\rightarrow\; \exists Z \, R_5(Z,Y,X)$$

# Stickiness: Marking Procedure

**Initial Marking**: mark all occurrences of body-variables that do not appear in <span style="color:red">all</span> head-atoms

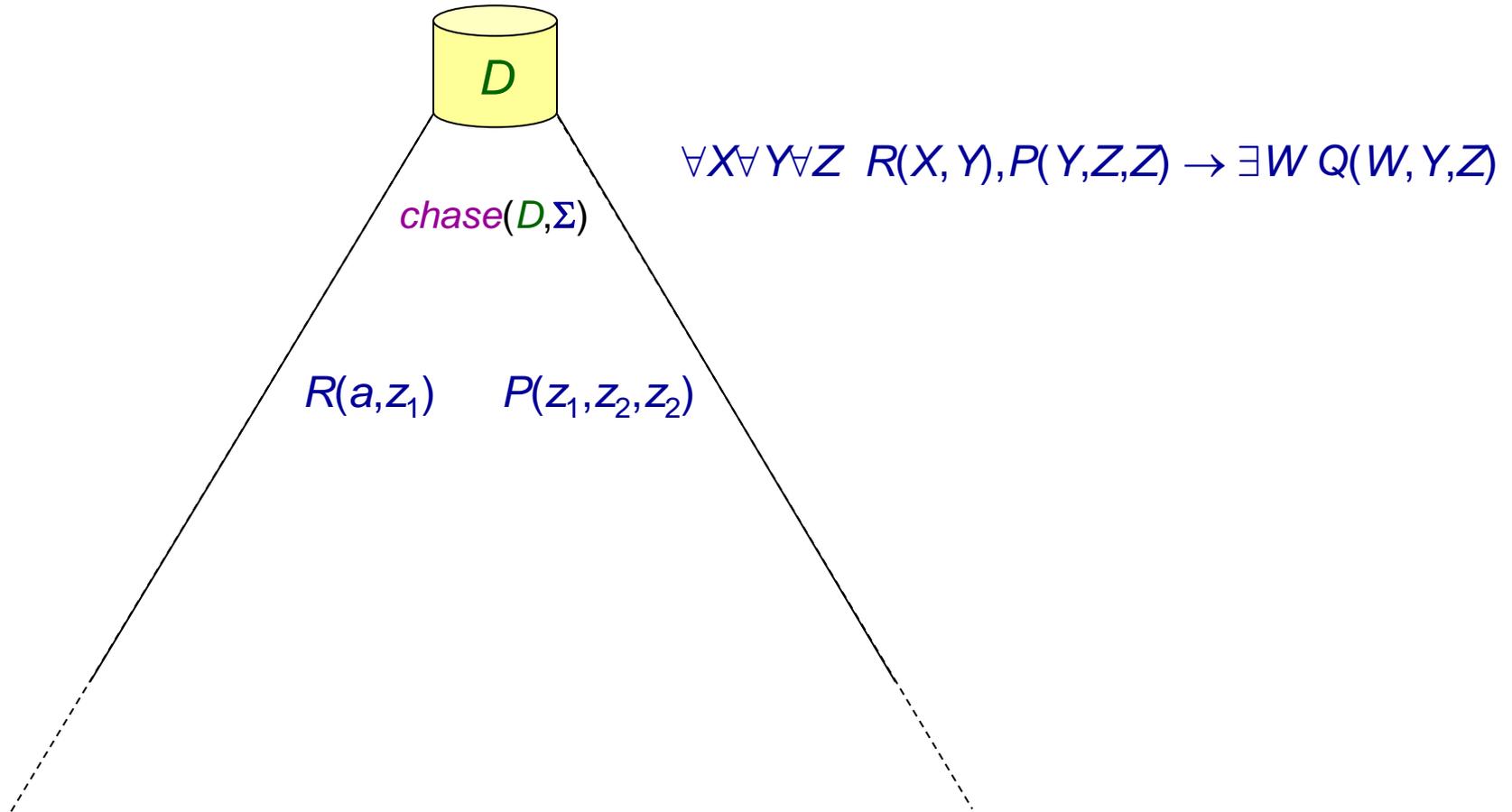**Propagation Step**: propagate the marking to body-variables via the head-atoms

$$\forall V \forall W \, R_1(V, W) \;\rightarrow\; \exists X \exists Y \exists Z \, R_2(W, X, Y, Z)$$

$$\forall V \forall W \forall X \forall Y \, R_2(V, W, X, Y) \;\rightarrow\; \exists Z \, R_1(W, Z), R_3(W, Y), R_4(Y, X)$$

$$\forall W \forall X \forall Y \, R_3(W, X), R_4(X, Y) \;\rightarrow\; \exists Z \, R_5(Z, Y, X)$$

# Stickiness
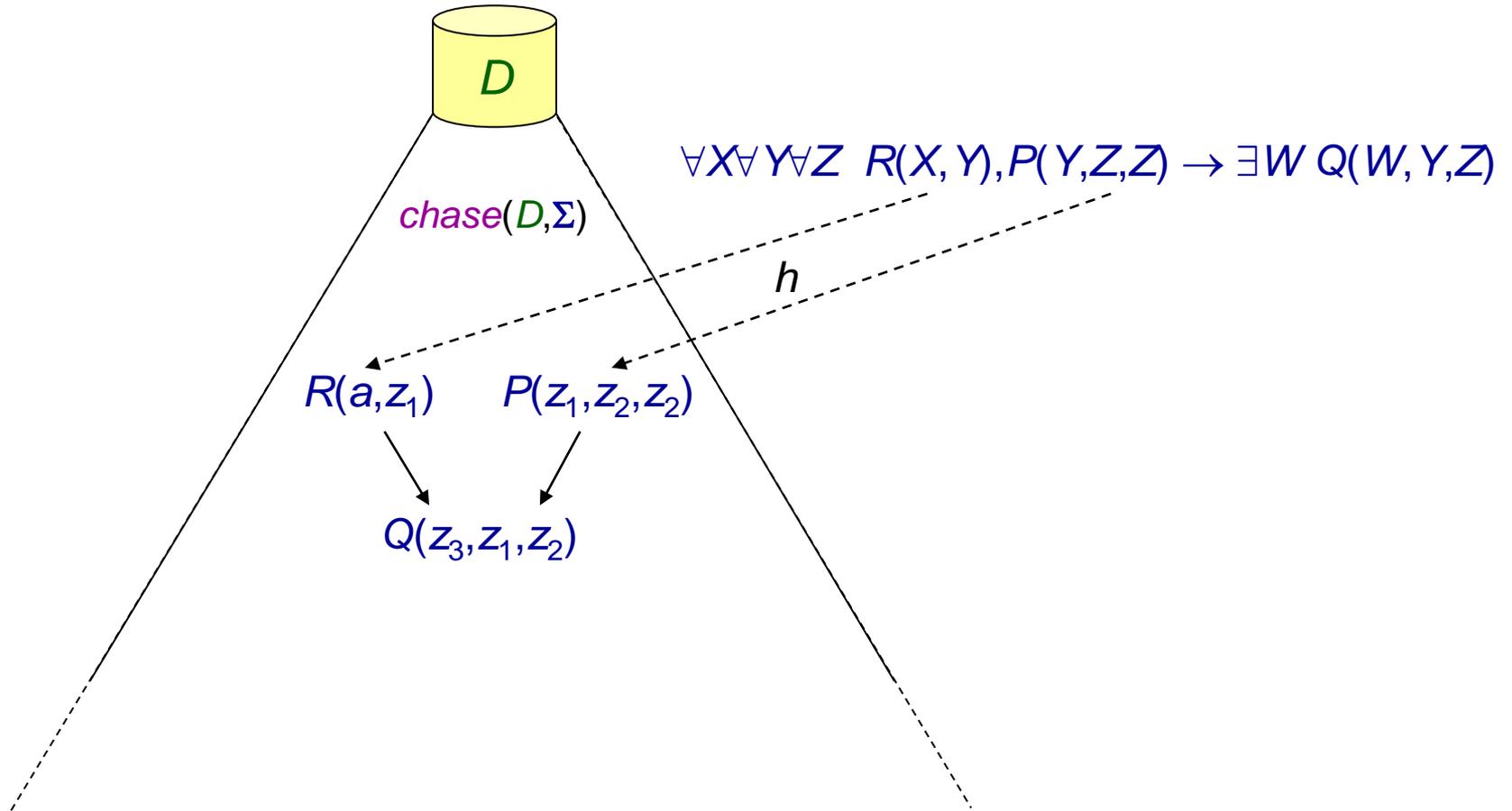
– Marked variables occur just once in the body of each TGD

$$\forall V \forall W \ R_1(V, W) \ \rightarrow \ \exists X \exists Y \exists Z \ R_2(W, X, Y, Z)$$

$$\forall V \forall W \forall X \forall Y \ R_2(V, W, X, Y) \ \rightarrow \ \exists Z \ R_1(W, Z), R_3(W, Y), R_4(Y, X)$$

$$\forall W \forall X \forall Y \ R_3(W, X), R_4(X, Y) \ \rightarrow \ \exists Z \ R_5(Z, Y, X)$$

– The chase has the sticky property (backward-resolution terminates)

# Sticky Property
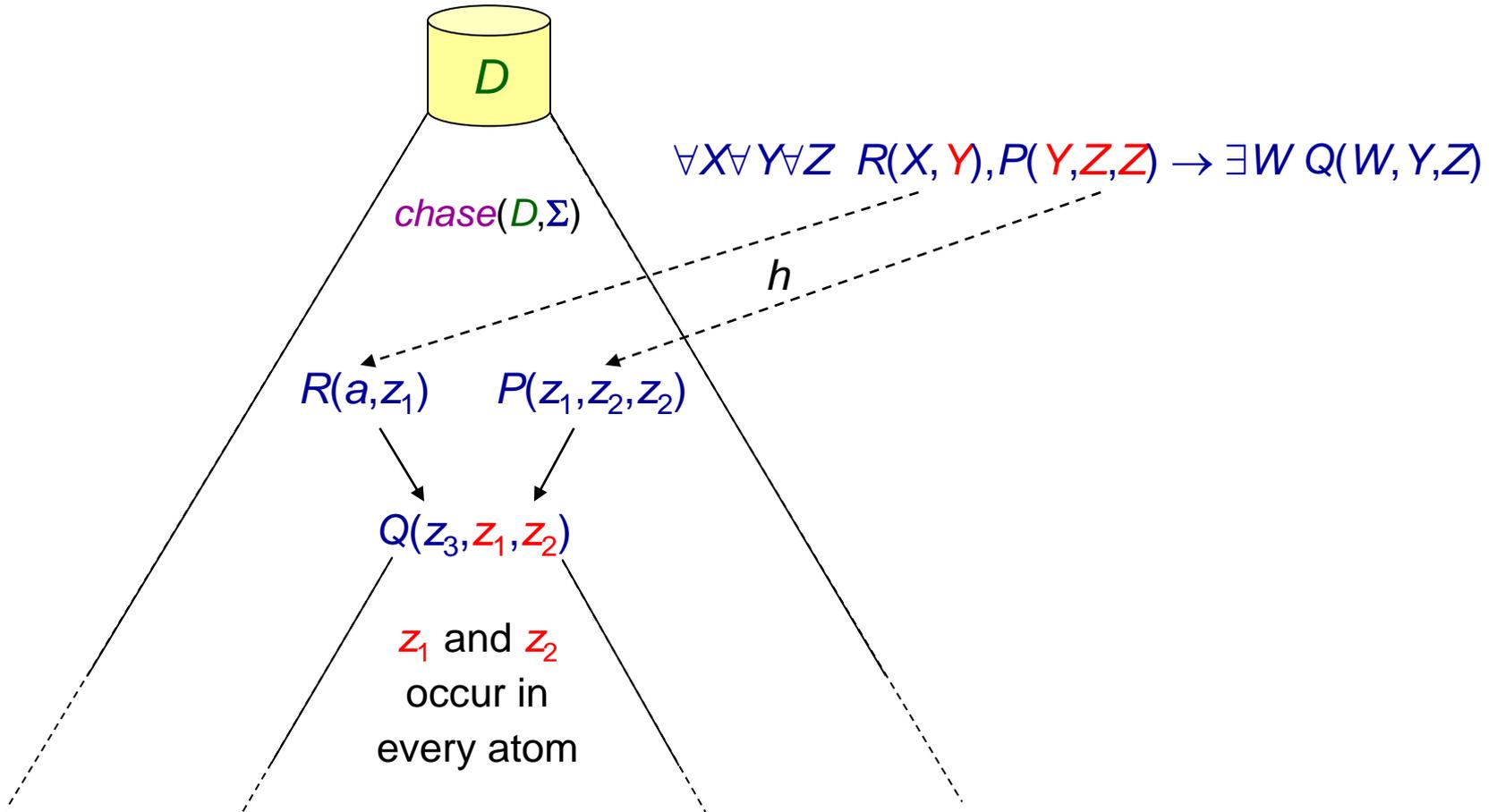


$\forall X \forall Y \forall Z \ R(X,Y), P(Y,Z,Z) \rightarrow \exists W \ Q(W,Y,Z)$

$chase(D,\Sigma)$

$R(a,z_1)$    $P(z_1,z_2,z_2)$

$D$

# Sticky Property



$$\forall X \forall Y \forall Z \; R(X,Y), P(Y,Z,Z) \rightarrow \exists W \; Q(W,Y,Z)$$

$chase(D, \Sigma)$

$D$

$h$

$R(a, z_1)$  $P(z_1, z_2, z_2)$

$Q(z_3, z_1, z_2)$

# Sticky Property



$\forall X \forall Y \forall Z \ R(X,Y), P(Y,Z,Z) \rightarrow \exists W \ Q(W,Y,Z)$

$chase(D,\Sigma)$

$h$

$R(a,z_1)$    $P(z_1,z_2,z_2)$

$Q(z_3,z_1,z_2)$

$z_1$ and $z_2$ occur in every atom

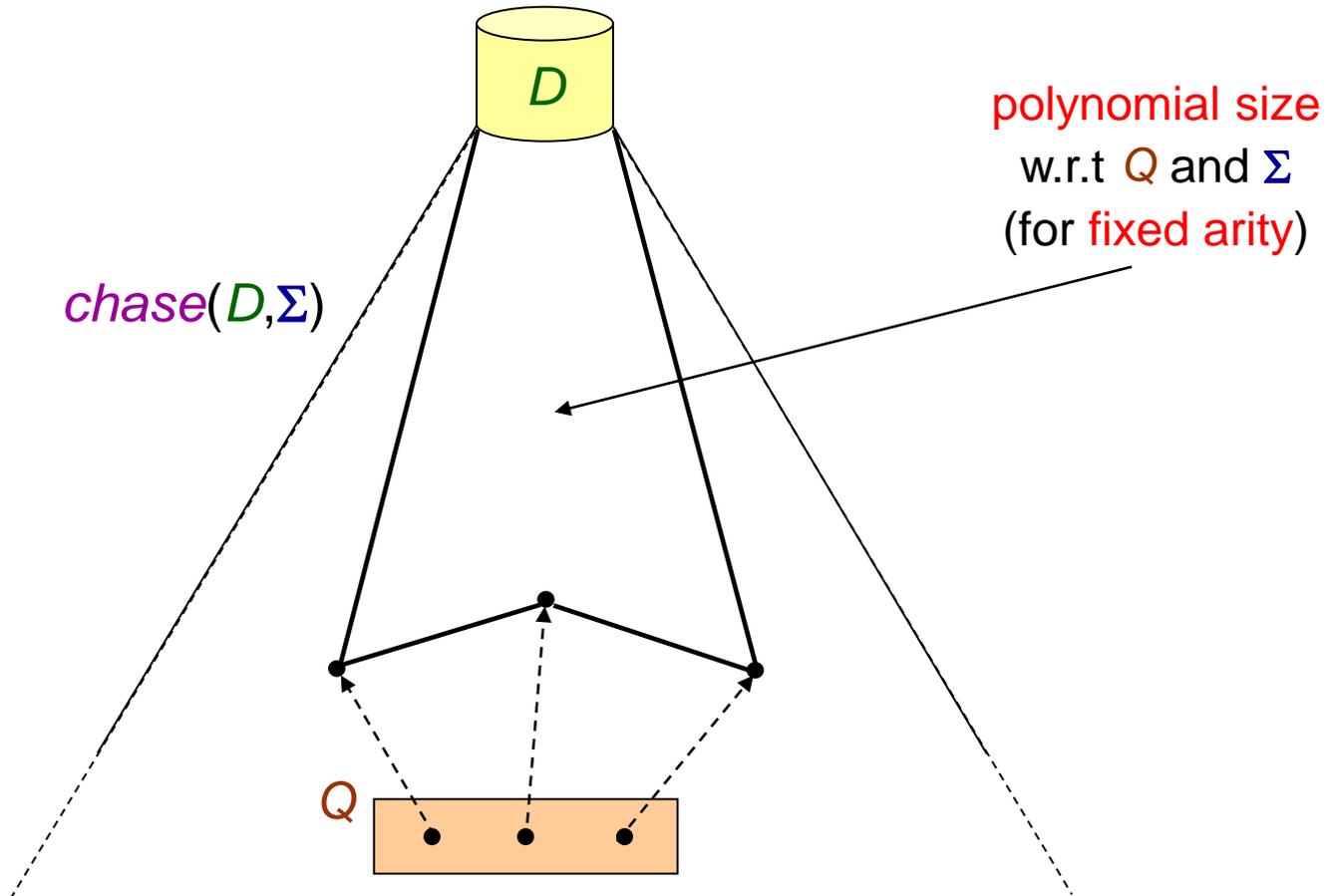# Stickiness

– Marked variables occur just once in the body of each TGD
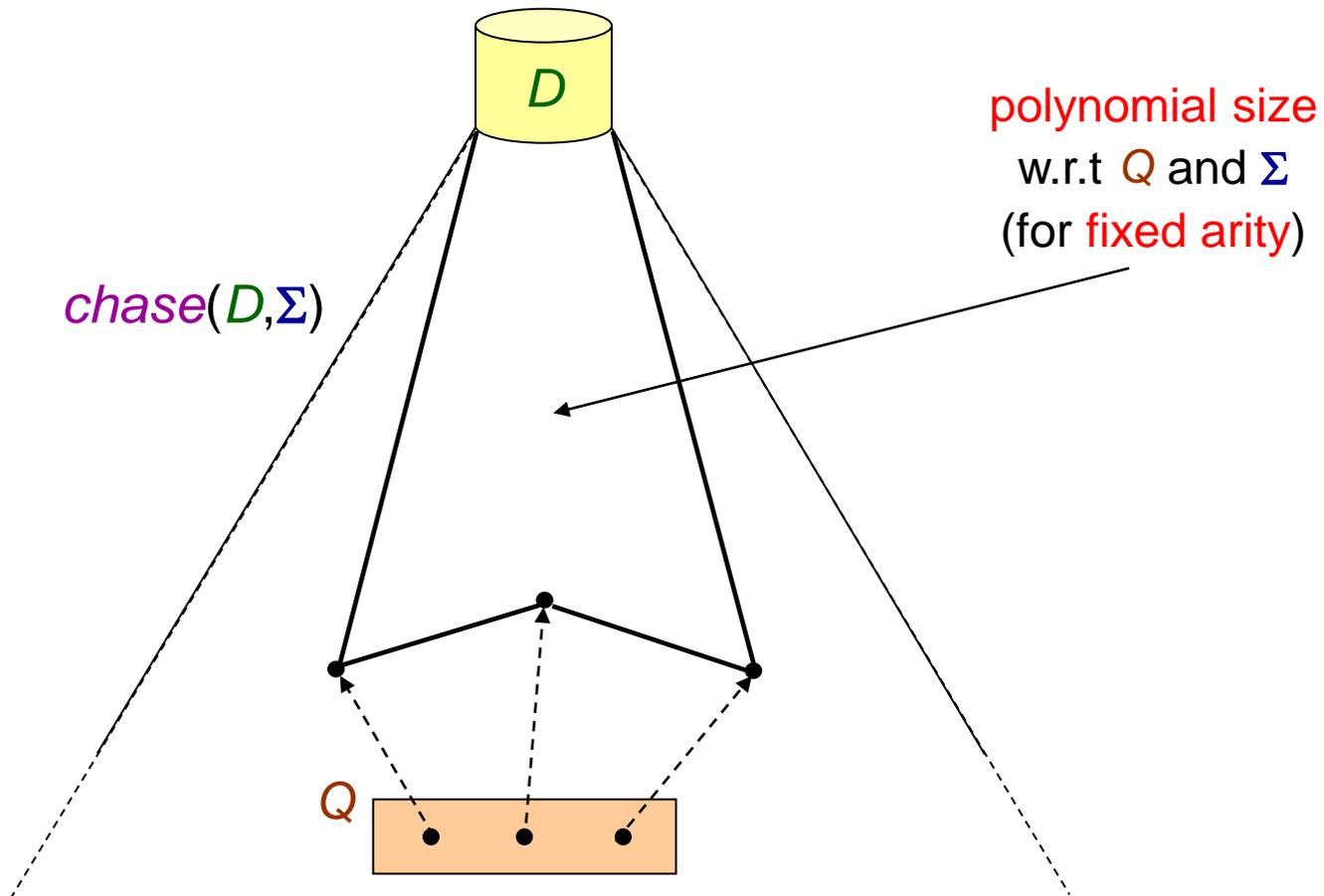
$$\forall V \forall W \; R_1(V, W) \;\; \rightarrow \;\; \exists X \exists Y \exists Z \; R_2(W, X, Y, Z)$$

$$\forall V \forall W \forall X \forall Y \; R_2(V, W, X, Y) \;\; \rightarrow \;\; \exists Z \; R_1(W, Z), R_3(W, Y), R_4(Y, X)$$

$$\forall W \forall X \forall Y \; R_3(W, X), R_4(X, Y) \;\; \rightarrow \;\; \exists Z \; R_5(Z, Y, X)$$

– The chase has the sticky property (backward-resolution terminates)

– Query answering is in $AC_0$ in data complexity (first-order rewritability)
[Calì, G. & Pieris, VLDB 10]

– Properly extends DL-Lite (same data complexity)
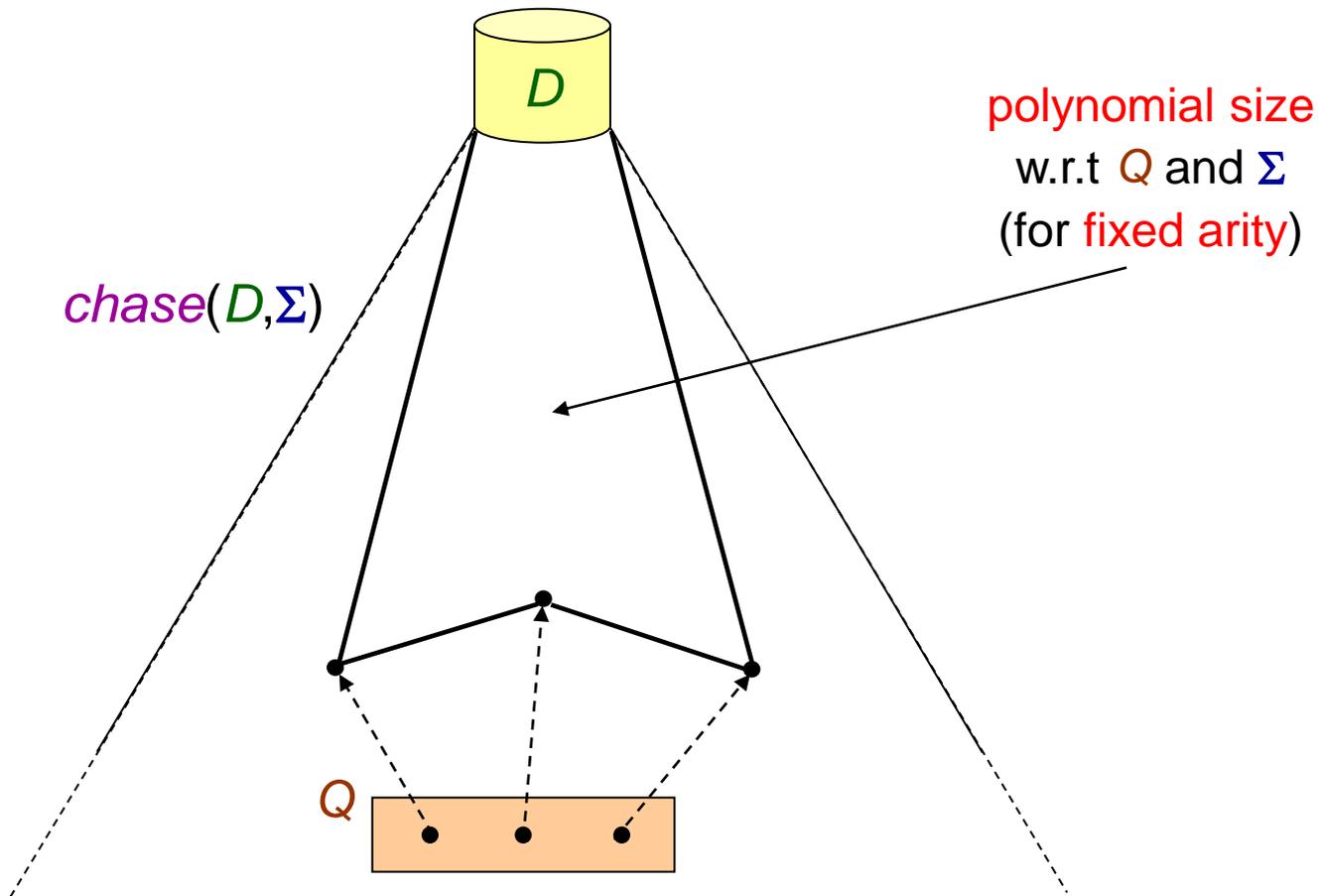
# Polynomial Witness Property (PWP)



polynomial size
w.r.t $Q$ and $\Sigma$
(for fixed arity)

$chase(D,\Sigma)$

$D$

$Q$

[G. & Schwentick, KR 12]

# Polynomial Witness Property (PWP)



$D$

polynomial size
w.r.t $Q$ and $\Sigma$
(for fixed arity)

*chase*($D$,$\Sigma$)

$Q$

PWP $\Rightarrow$ non-recursive Datalog rewriting of polynomial size

[G. & Schwentick, KR 12]

# Polynomial Witness Property (PWP)



polynomial size
w.r.t $Q$ and $\Sigma$
(for fixed arity)

$chase(D,\Sigma)$

Holds for the first-order rewritable Datalog$^{\pm}$ fragments (linear and sticky)

[G. & Schwentick, KR 12]

# Finite Controllability

$$D \cup \Sigma \models Q \quad \overset{?}{\Leftrightarrow} \quad D \cup \Sigma \models_{\text{fin}} Q$$

– Holds for inclusion dependencies

[Rosati, PODS 06]

– Holds for guarded TGDs (in fact, for the guarded fragment)

[Bárány, Gottlob & Otto, LICS 10]

– Holds for sticky TGDs

[Gogacz & Marcinkowski, → next talk]

# Additional Features

- EGDs, e.g., $\forall X \forall Y \forall Z\ reports(X,Y), reports(X,Z) \rightarrow Y = Z$

    Non-Conflicting EGDs: do not interact with TGDs

    Preliminary check without adding complexity


- Negative constraints, e.g., $\forall X\ emp(X), customer(X) \rightarrow \perp$

    Check without adding complexity

# Additional Features

– EGDs, e.g., $\forall X \forall Y \forall Z\ reports(X, Y), reports(X, Z) \rightarrow Y = Z$

Non-Conflicting EGDs: do not interact with TGDs

Preliminary check without adding complexity

– Negative constraints, e.g., $\forall X\ emp(X), customer(X) \rightarrow \perp$

Check without adding complexity

---

Finite controllability does not hold

$D = \{R(a,b)\}$

$$\Sigma = \left\{ \begin{array}{l} \forall X \forall Y\ R(X, Y) \rightarrow \exists Z\ R(Y,Z) \\[1em] \forall X \forall Y \forall Z\ R(Y,X), R(Z,X) \rightarrow Y = Z \end{array} \right\}$$
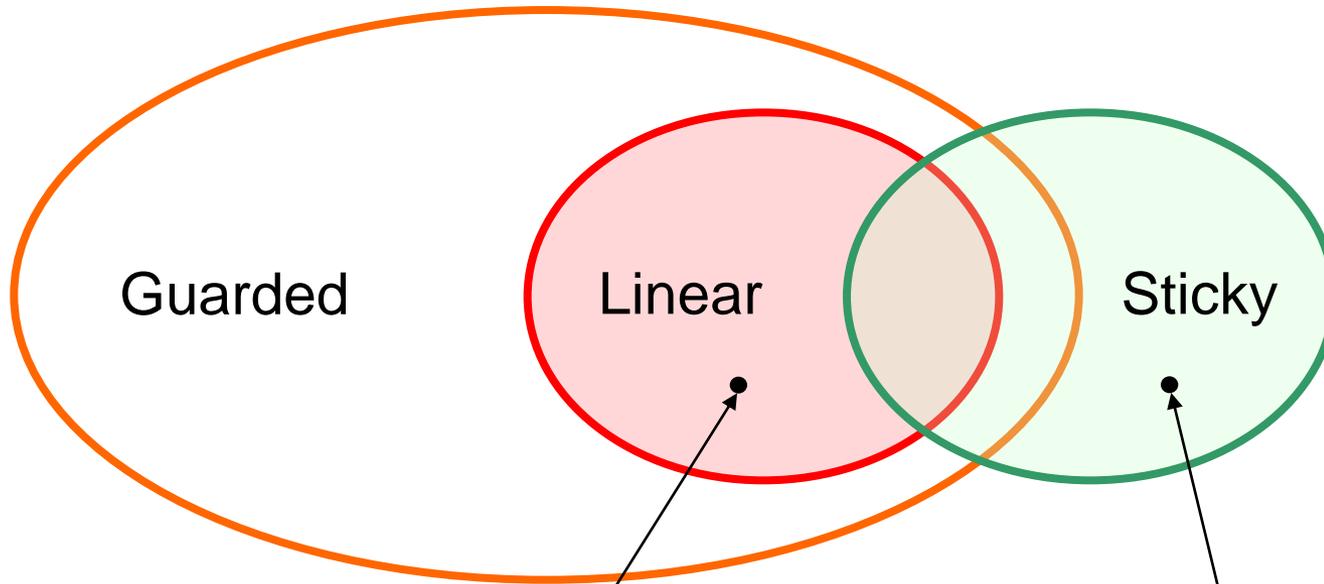
$Q \leftarrow R(A, a)$

$D \cup \Sigma \nvDash Q$

but

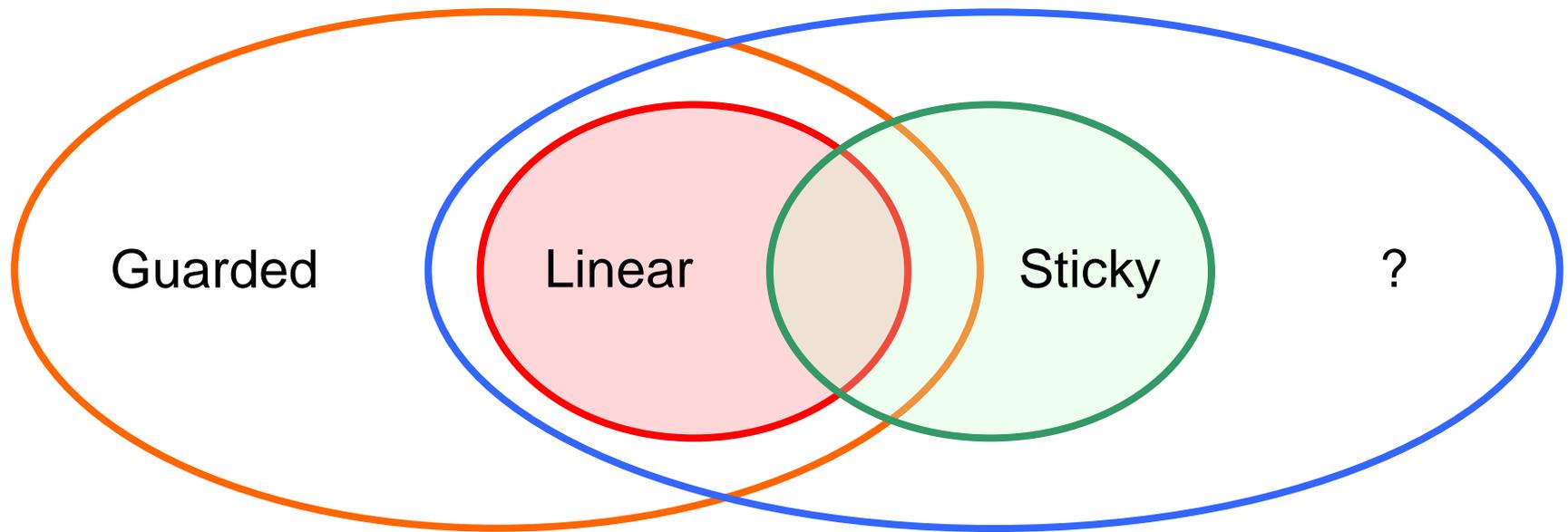$D \cup \Sigma \vDash_{fin} Q$

# Datalog$^{\pm}$: Overview



$\forall X \forall Y\ R(X,X,Y) \rightarrow \exists Z\ R(Y,Y,Z)$

$\forall X \forall Y \forall Z \forall W\ S(X,Y),R(Z,W) \rightarrow \exists V\ S(Y,V),R(W,V)$

# Datalog$^{\pm}$: Overview



Guarded   Linear   Sticky   ?

without losing first-order rewritability and PWP

# Sticky-Join TGDs

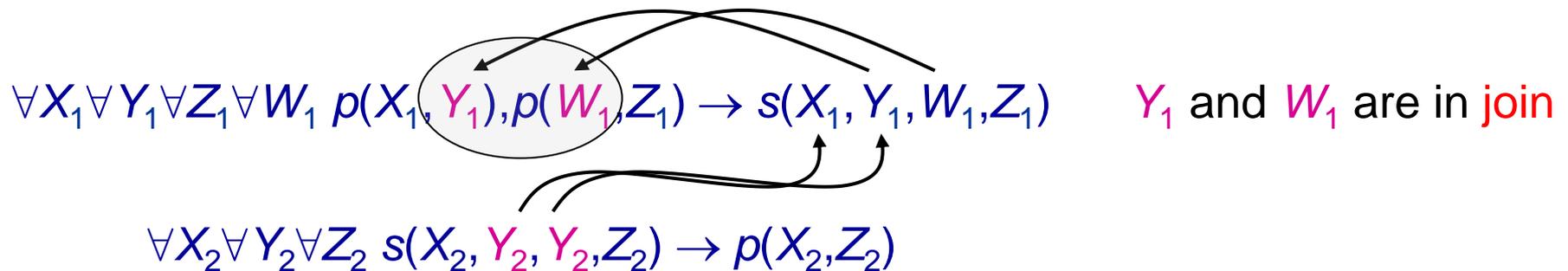– Marked variables can appear more than once in a <span style="color:red">single atom</span>

– Providing that the following <span style="color:red">does not happen</span>:

$$\forall X_1 \forall Y_1 \forall Z_1 \forall W_1 \; p(X_1, Y_1), p(W_1, Z_1) \rightarrow s(X_1, Y_1, W_1, Z_1)$$

$Y_1$ and $W_1$ are in <span style="color:red">join</span>

$$\forall X_2 \forall Y_2 \forall Z_2 \; s(X_2, Y_2, Y_2, Z_2) \rightarrow p(X_2, Z_2)$$

# Sticky-Join TGDs

- Marked variables can appear more than once in a <span style="color:red">single atom</span>

- Providing that the following <span style="color:red">does not happen</span>:

$$\forall X_1 \forall Y_1 \forall Z_1 \forall W_1 \; p(X_1, Y_1), p(W_1, Z_1) \rightarrow s(X_1, Y_1, W_1, Z_1)$$  $Y_1$ and $W_1$ are in <span style="color:red">join</span>

$$\forall X_2 \forall Y_2 \forall Z_2 \; s(X_2, Y_2, Y_2, Z_2) \rightarrow p(X_2, Z_2)$$

- <span style="color:red">Same complexity</span> as sticky TGDs and <span style="color:red">enjoy PWP</span>

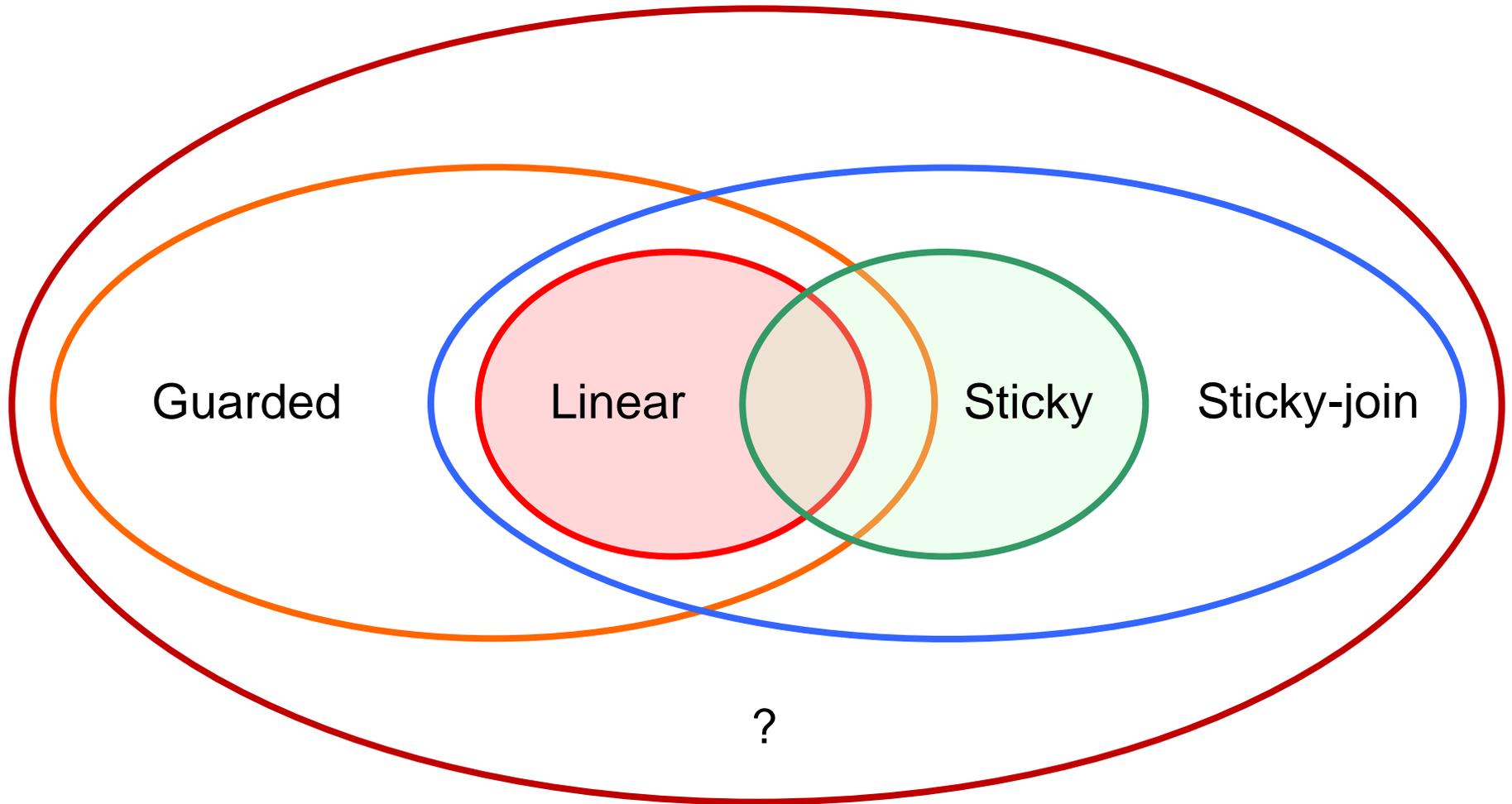- But harder to identify them - <span style="color:red">PSPACE-complete</span>

# Datalog$^{\pm}$: Overview



$\forall X \forall Y \forall Z \ R(X,Y), S(Y,Z,Z) \rightarrow \exists W \ P(Y,W)$

Guarded

Linear

Sticky

Sticky-join

$\forall X \forall Y \ R(X,X,Y) \rightarrow \exists Z \ R(Y,Y,Z)$

$\forall X \forall Y \forall Z \ S(X,Y), R(Y,Z) \rightarrow \exists W \ P(Y,W)$

# Datalog$^{\pm}$: Overview



Guarded    Linear    Sticky    Sticky-join

?

without losing PTIME data complexity

# Tame TGDs   (TTGDs)

– Check that each rule is guarded* or sticky-join (and thus sticky)

$$\forall X \forall Y \forall Z \; R(X,Y) \rightarrow \exists Z \; R(Y,Z)$$   sticky-join and guarded*

$$\forall X \forall Y \; T(X), R(X,Y) \rightarrow T(Y)$$   Guarded* not sticky-join

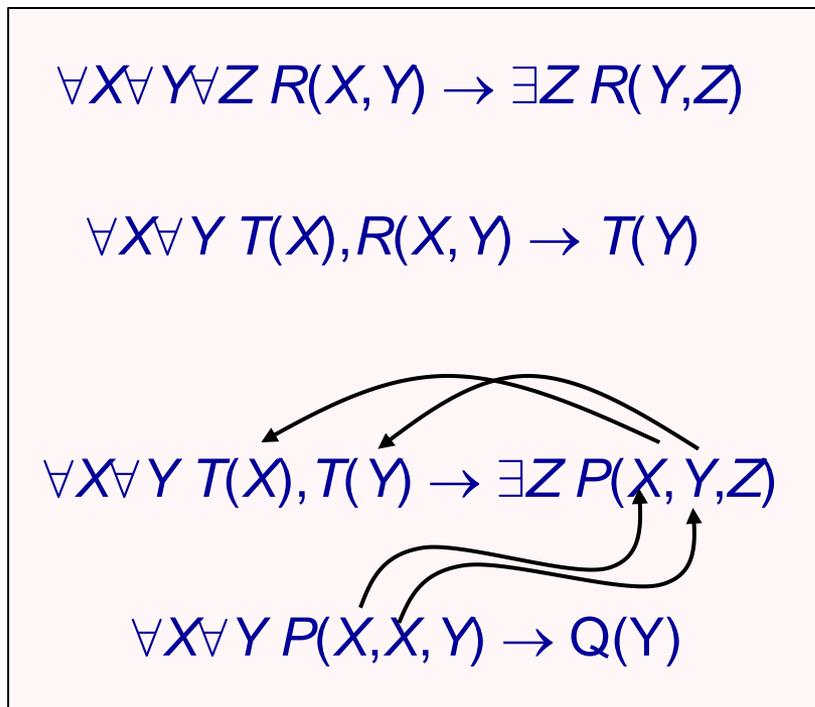$$\forall X \forall Y \; T(X), T(Y) \rightarrow \exists Z \; P(X,Y,Z)$$   sticky-join not guarded

– Not first-order rewritable due to the second rule

– Chase has infinite treewidth since $P$ forms an infinite clique

*) *Predicates in heads of non-guarded rules are unguarded. Atoms whose predicate is unguarded are not allowed to serve as guards!*

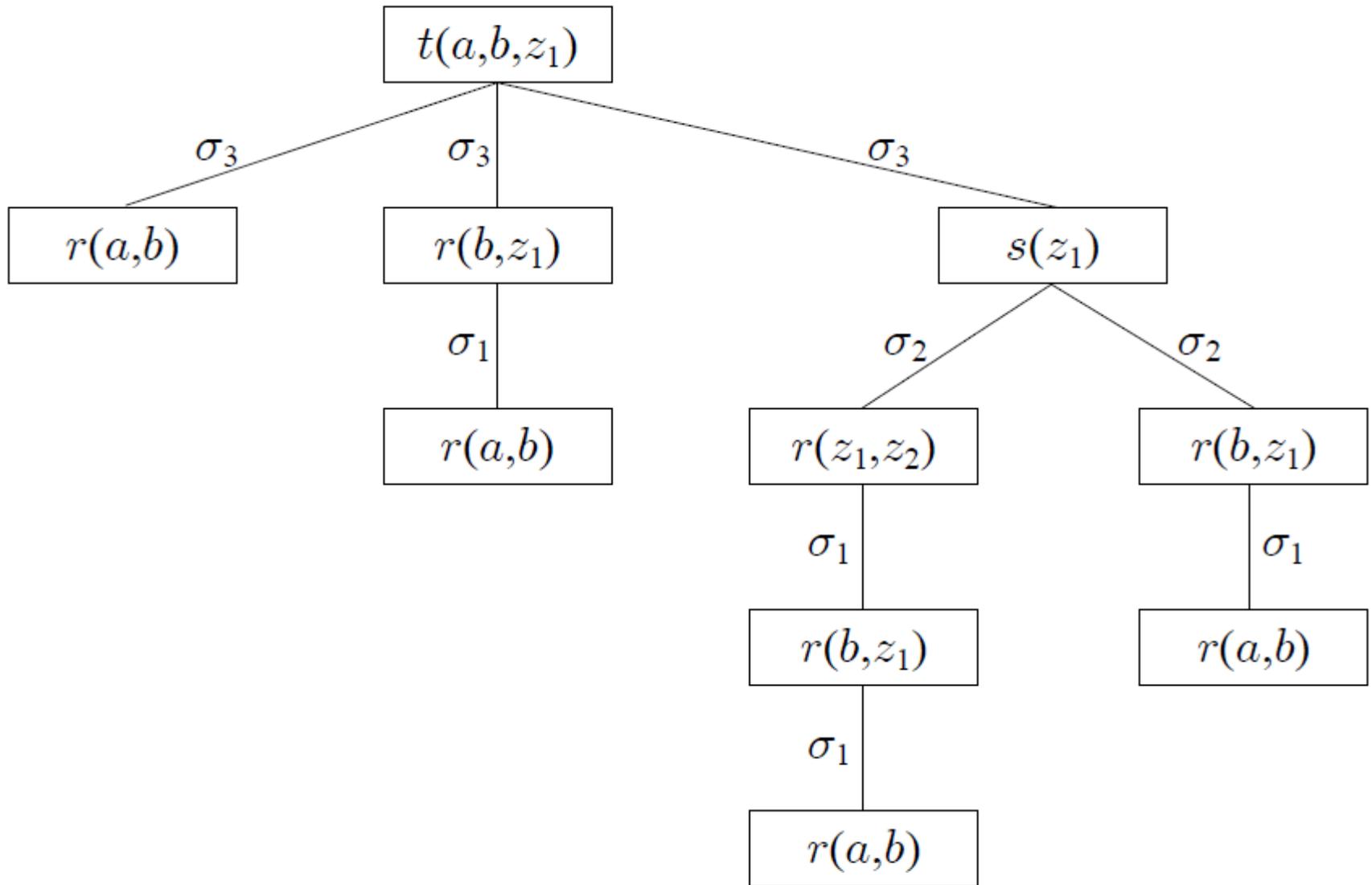# Tame TGDs

– The following set of TGDs is not tame

$\forall X \forall Y \forall Z\ R(X,Y) \rightarrow \exists Z\ R(Y,Z)$      sticky-join and guarded*

$\forall X \forall Y\ T(X),R(X,Y) \rightarrow T(Y)$      Guarded* not sticky-join

$\forall X \forall Y\ T(X),T(Y) \rightarrow \exists Z\ P(X,Y,Z)$      neither sticky-join nor guarded
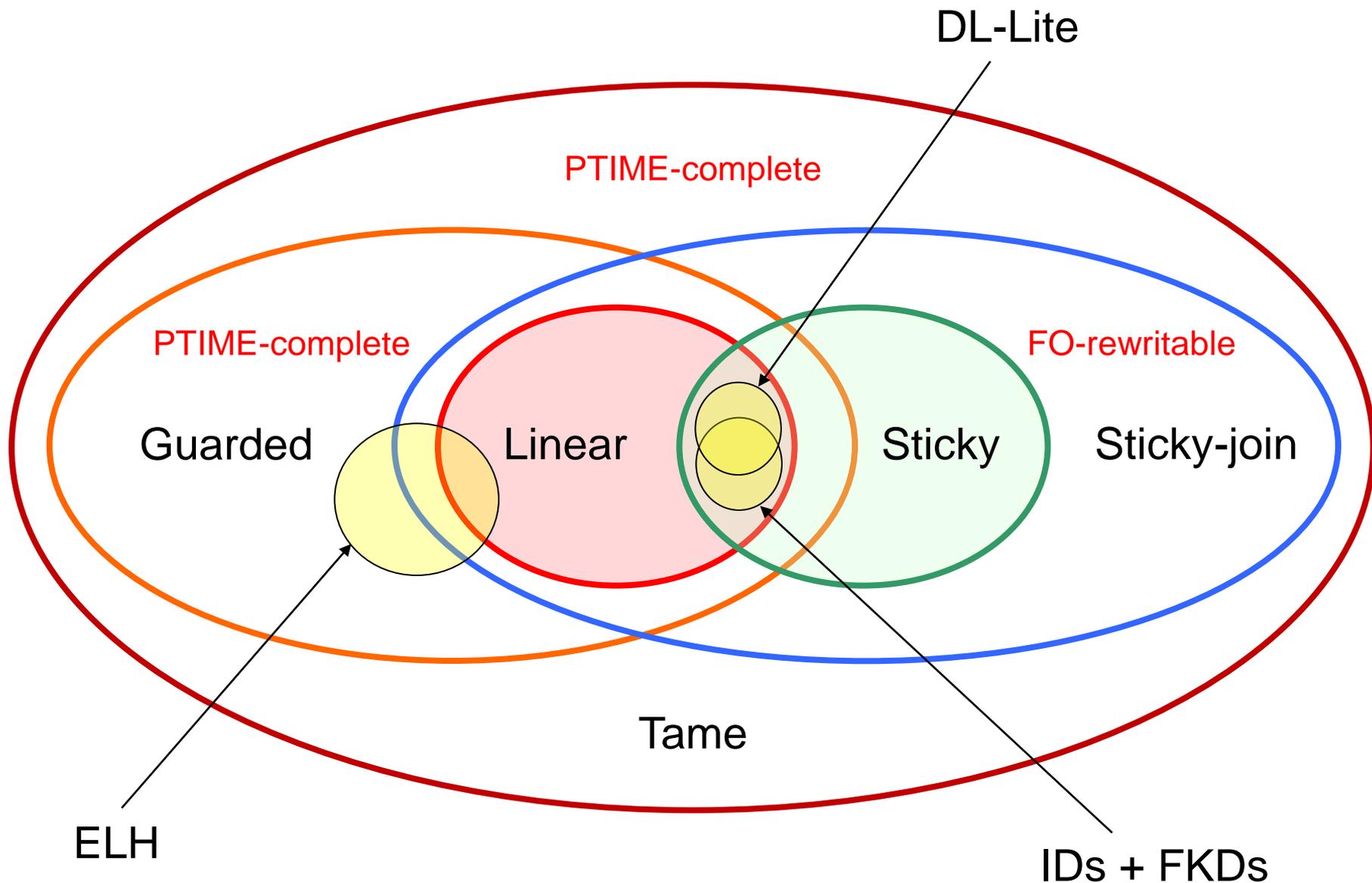
$\forall X \forall Y\ P(X,X,Y) \rightarrow Q(Y)$      sticky-join and guarded

*) Predicates in heads of non-guarded rules are unguarded. Atoms whose predicate is unguarded are not allowed to serve as guards!

No bounded tw model, but bounded tw resolution proof scheme!

# Datalog$^{\pm}$: Overview



DL-Lite

PTIME-complete

PTIME-complete

FO-rewritable

Guarded    Linear    Sticky    Sticky-join

Tame

ELH

IDs + FKDs

# Datalog$^{\pm}$: Summary of Complexity Results

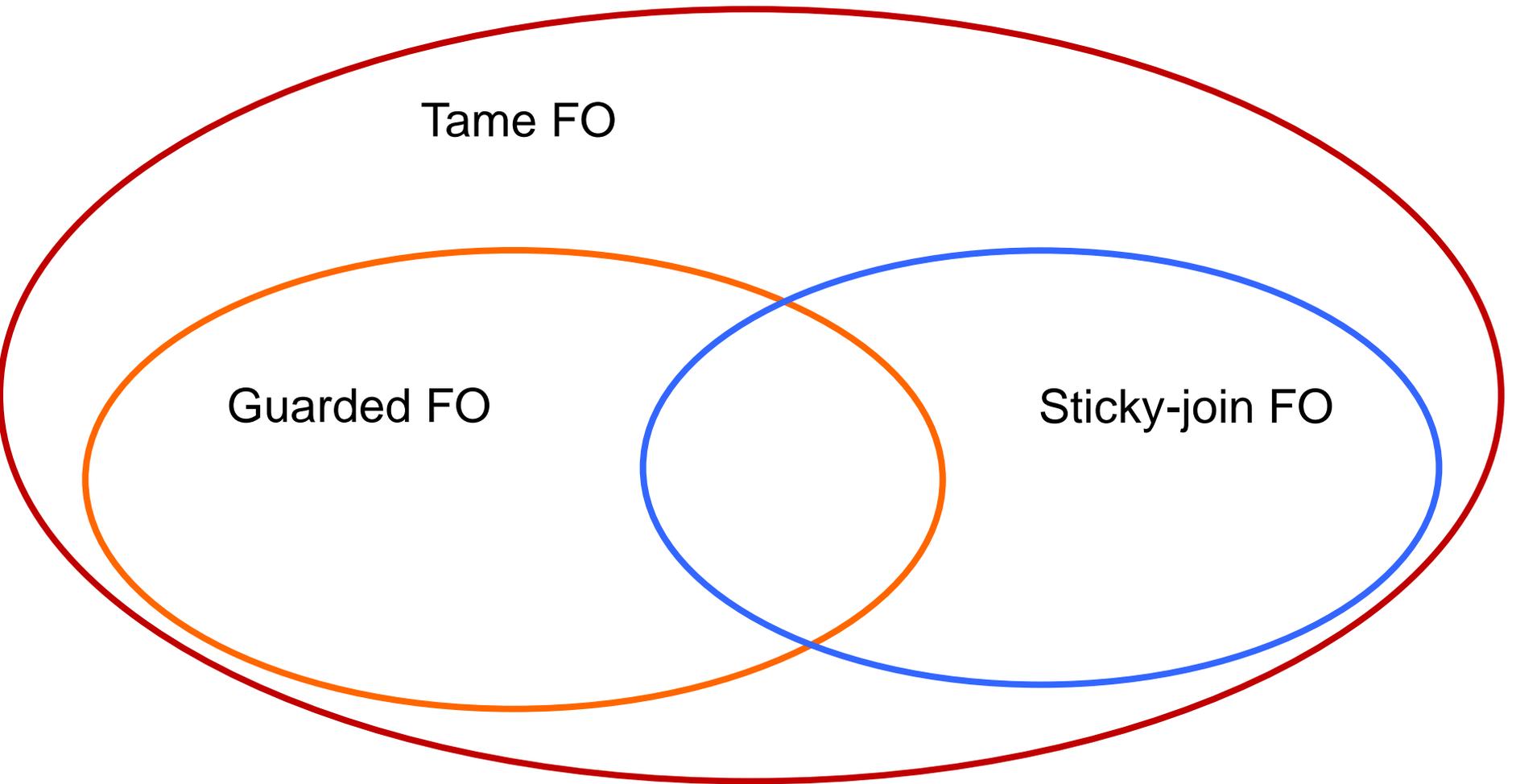| | Data | Fixed $\Sigma$ | Combined |
|---|---|---|---|
| Guarded | PTIME | NP | 2EXPTIME |
| Linear | in $AC_0$ | NP | PSPACE |
| Sticky | in $AC_0$ | NP | EXPTIME |
| Sticky-Join | in $AC_0$ | NP | EXPTIME |
| Tame | PTIME | NP | 2EXPTIME |

Same complexity with negative constraints and non-conflicting EGDs

# *current work*

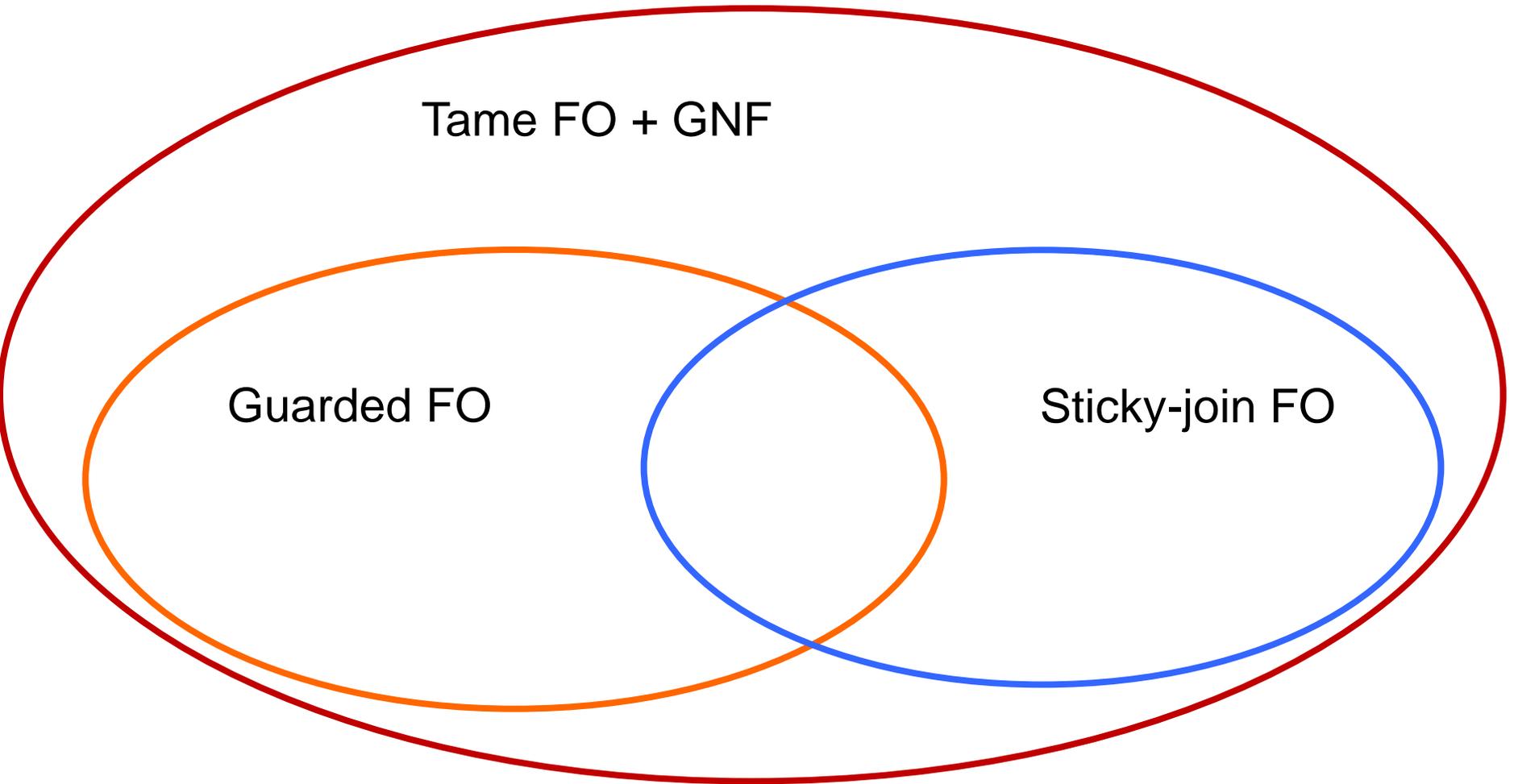*Finite controllability of Tamed TGDs:*

*Reduction to Sticky TGDs*

*current work*

Tame FO

Guarded FO

Sticky-join FO

*current work*



Tame FO + GNF

Guarded FO

Sticky-join FO

# Thank you!