# SemDDA: A Semantics-based Database Dependency Analyzer

Angshuman Jana, Raju Halder

Indian Institute of Technology Patna

{ajana.pcs, halder}@iitp.ac.in

## INTRODUCTION

• Dependency among program statements has been widely and successfully used in many software-engineering activities, *e.g.* Maintenance, Safety Verification, Optimization, etc [1].

• One most suitable representation of these dependences (data- and control-dependences) is in the form of Dependency Graph.

• Various forms of dependency graph are: Program Dependence Graph (PDG), System Dependence Graph (SDG), Class Dependence Graph (ClDG), Database-Oriented Program Dependence Graph (DOPDG), etc [2].

• The syntactic construction of dependency graph is based on the syntactic presence of variables (either defined or used) in program statements.

## MOTIVATIONS

• Syntactic dependency computation may produce false alarms. Therefore, it may fail to compute an optimal set of dependences.

• For instance, expression "$e = x^2 + 4w \bmod 2 + z$" syntactically depends on $w$, semantically there is no dependency as the evaluation of "$4w \bmod 2$" is always zero.

• In [3] authors introduced the notion of semantic-data dependency which focuses on the actual values of variable rather than their syntactic presence.

## OBJECTIVE

• We present SemDDA, a novel Semantics-based Database Dependency Analyzer.

• This tool computes semantic-based dependences in database applications.

• For example, consider the following SQL statements $Q_1$ and $Q_2$:

$Q_1$ : UPDATE emp SET sal := sal+100 WHERE sal ≤ 2000

$Q_2$ : SELECT AVG(sal) FROM emp WHERE sal ≥ 3000

• The SQL statement $Q_2$ is syntactically Dependent on $Q_1$ but there exist no semantics-based dependence between $Q_1$ and $Q_2$.

## SYNTAX-BASED DOPDG

• The DOPDG [2] is an extension of PDG with two additional dependences.

(*i*) Program-Database Dependences (PD-Dependences).

(*ii*) Database-Database Dependences (DD-dependences).

## SEMANTIC-BASED DEPENDENCY

• The SQL statement $Q = \langle A, \phi \rangle$ where $A$ represents an action-part and $\phi$ represents a conditional-part.

• Let $\rho_{db} \in \Sigma_{db}$ be a database state and $\mathbf{S} : Q \times \Sigma_{db} \to \Sigma_{db}$ be a semantic function.

• Functions $\mathfrak{A}_{use}$, $\mathfrak{A}_{def}$ compute used- and defined-part of target table $t$ by $Q$ as bellow:
$$\mathfrak{A}_{use}(Q, t) = \rho_{t \downarrow \phi}(\vec{x}) \cup \rho_{t \downarrow \phi}(\vec{y}),$$
$$\mathfrak{A}_{def}(Q, t) = \Delta(\rho_{t'}(\vec{z}), \rho_t(\vec{z}))$$
where $\vec{x} = \text{USE}(A)$ and $\vec{y} = \text{USE}(\phi)$ are a list of *used*-variable in $A$ and $\phi$. $\vec{z} = \text{DEF}(Q)$ is a list of *defined*-variable in $Q$ and $(t \downarrow \phi)$ denotes the set of tuples in $t$ for which $\phi$ holds "*true*".

• Formally, the SQL statement $Q_2 = \langle A_2, \phi_2 \rangle$ with $target(Q_2) = t'$ is DD-Dependent on another SQL statement $Q_1$ for $\Upsilon$ ( $Q_1 \xrightarrow{\Upsilon} Q_2$ ) iff $Q_1 \in \{Q_{upd}, Q_{ins}, Q_{del}\}$ and the overlapping-part $\Upsilon = \mathfrak{A}_{use}(Q_2, t') \cap \mathfrak{A}_{def}(Q_1, t) \neq \emptyset$.

## ABSTRACT INTERPRETATION

• Abstract Interpretation [4] is a semantics-based static analysis framework.

• It provides a sound approximation of program semantics focusing on a particular property.

## ABSTRACT SEMANTICS

• Let $\mathbb{P} \in \wp$ be a polyhedra and $\mathbb{P}_T$, $\mathbb{P}_F$ be the polyhedral form of the *true*- and *false*-part of the database. The abstract transition semantics is defined as $\bar{\mathbf{S}} : \mathbb{C} \times \wp \to \wp(\wp)$ where $\mathbb{C}$ is the set of database statements and $\wp$ is the set of all polyhedra. The computation of $\mathfrak{A}_{use}$ and $\mathfrak{A}_{def}$ are defined *w.r.t.* abstract semantics as below:
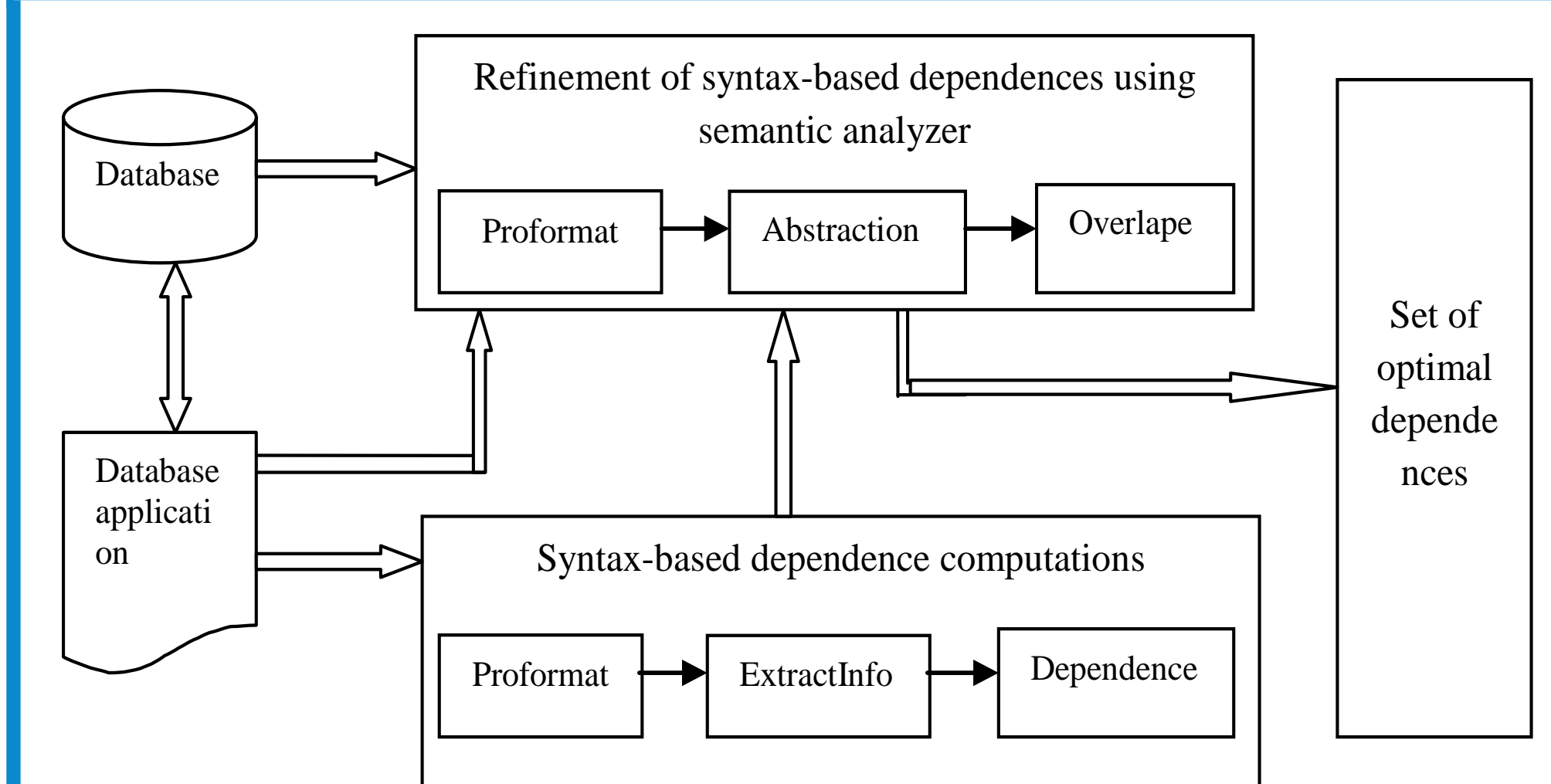
Assignment: $\bar{\mathbf{S}}[\![x_j := e]\!](\mathbb{P}) = \{\mathbb{P}'\}$

Test: $\bar{\mathbf{S}}[\![\beta]\!]\mathbb{P} = \{\mathbb{P}_T, \mathbb{P}_F\}$ where $\mathbb{P}_T = (\mathbb{P} \sqcap \beta) \wedge \mathbb{P}_F = (\mathbb{P} \sqcap \neg\beta)$

Update: $\bar{\mathbf{S}}[\![\text{UPDATE}(\vec{v}_d, \vec{e}), \phi]\!]\mathbb{P} = \bar{\mathbf{S}}[\![\text{UPDATE}(\vec{v}_d, \vec{e}), true]\!]\mathbb{P}_T \cup \bar{\mathbf{S}}[\![\text{UPDATE}(\vec{v}_d, \vec{e}), false]\!]\mathbb{P}_F = \{\mathbb{P}'_T, \mathbb{P}_F\}$
$\mathfrak{A}_{use}(Q_{upd}, t) = \langle \mathbb{P}_T \rangle$ and $\mathfrak{A}_{def}(Q_{upd}, t) = \langle \mathbb{P}_T, \mathbb{P}'_T \rangle$

Delete: $\bar{\mathbf{S}}[\![\text{DELETE}(\vec{v}_d), \phi]\!]\mathbb{P} = \bar{\mathbf{S}}[\![\text{DELETE}(\vec{v}_d), true]\!]\mathbb{P}_T = \{\mathbb{P}'\}$
$\mathfrak{A}_{use}(Q_{del}, t) = \langle \mathbb{P}_T \rangle$ and $\mathfrak{A}_{def}(Q_{del}, t) = \langle \mathbb{P}_T, \emptyset \rangle$

Insert: $\bar{\mathbf{S}}[\![\text{INSERT}(\vec{v}_d, \vec{e}), \phi]\!]\mathbb{P} = \bar{\mathbf{S}}[\![\text{INSERT}(\vec{v}_d, \vec{e}), true]\!]\mathbb{P} = \{\mathbb{P} \sqcup \mathbb{P}_{new}\}$ where $\mathbb{P}_{new}$ is the polyhedron represented by the inserted tuple values.
$\mathfrak{A}_{use}(Q_{ins}, t) = \langle \emptyset \rangle$ and $\mathfrak{A}_{def}(Q_{ins}, t) = \langle \emptyset, \mathbb{P}_{new} \rangle$

Select: $\bar{\mathbf{S}}[\![\langle \text{SELECT}(v_a, f(\vec{e'}), r(\vec{h}(\vec{x})), \phi_2, g(\vec{e})), \phi_1 \rangle]\!]\mathbb{P}$
$= \bar{\mathbf{S}}[\![\langle \text{SELECT}(v_a, f(\vec{e'}), r(\vec{h}(\vec{x})), \phi_2, g(\vec{e})), true \rangle]\!]\mathbb{P}_T \sqcup \bar{\mathbf{S}}[\![\langle \text{SELECT}(v_a, f(\vec{e'}), r(\vec{h}(\vec{x})), \phi_2, g(\vec{e})), false \rangle]\!]\mathbb{P}_F$
$= \{\mathbb{P}\}$ The polyhedron is unchanged .
$\mathfrak{A}_{use}(Q_{sel}, t) = \langle \mathbb{P}_T \rangle$ and $\mathfrak{A}_{def}(Q_{sel}, t) = \langle \emptyset, \emptyset \rangle$

## DEPENDENCY COMPUTATIONS

• The SQL statements $Q1$ and $Q2$ are semantically independent when $\Upsilon = \mathfrak{A}_{def}(Q_1, t) \cap \mathfrak{A}_{use}(Q_2, t') = \emptyset$ iff
$$\mathbb{P}_T^{Q1} \sqcap \mathbb{P}_T^{Q2} = \emptyset \wedge \mathbb{P}_T^{\prime Q1} \sqcap \mathbb{P}_T^{Q2} = \emptyset$$
where $\mathbb{P}_T^{Q1}$ and $\mathbb{P}_T^{\prime Q1}$ denote the *defined*-part of $Q1$ and $\mathbb{P}_T^{Q2}$ denotes the *used*-part of $Q2$.

## TOOL ARCHITECTURE



## APPLICATIONS

• Program Slicing.
• Information Flow Security Analysis.
• Data provenance.
• Concurrent System modeling.

## FUTURE PLAN

• Currently we are extending our tool by defining abstract semantic of database applications in non-relational abstract domain interval and relational abstract domain octagon.

• We will also perform experiment on benchmark codes using our tool.

## REFERENCES

[1] A. Podgurski and L. A. Clarke, "A formal model of program dependences and its implications for software testing, debugging, and maintenance," IEEE Trans. on SE, 1990

[2] D. Willmor, S. M. Embury, and J. Shao, "Program slicing in the presence of a database state," in Proc. of IEEE ICSM, 2004.

[3] I. Mastroeni and D. Zanardini, "Data dependencies and program slicing: from syntax to abstract semantics," in Proc. of Partial evaluation & semantics-based prog. manipulation, 2008.

[4] P. Cousot and R. Cousot. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In Proc. of the POPL'77, ACM Press.

## THE SNAPSHOT OF THE TOOL SemDDA