# Automated constraint verification for databases

#### Shuai Yuan advised by Véronique Benzaken, Evelyne Contejean and Mingli Song

College of Computer Science and Technology, Zhejiang University

#### Outline

- Background
- Context
- Proposed solutions
- Database integrity constraints
- SQL data modification operations
- Why3
- Translation overiew
- Problem and solutions
- Experiments
- Conclusion and perspectives

- Integrity constraints are important to express the semantics of databases.
  - prevent the execution of operations or transactions which will cause violation of constraints
  - semantics query optimization
- No real DBMS (database management system) have fully support the management of integrity constraints.
   i.e. assertions which are part of the SQL standard
- The alternative solution widely adopted by main-stream DBMS triggers:

- Integrity constraints are important to express the semantics of databases.
  - prevent the execution of operations or transactions which will cause violation of constraints
  - semantics query optimization
- No real DBMS (database management system) have fully support the management of integrity constraints. i.e. assertions which are part of the SQL standard
- The alternative solution widely adopted by main-stream DBMS triggers:

- Integrity constraints are important to express the semantics of databases.
  - prevent the execution of operations or transactions which will cause violation of constraints
  - semantics query optimization
- No real DBMS (database management system) have fully support the management of integrity constraints.
   i.e. assertions which are part of the SQL standard
- The alternative solution widely adopted by main-stream DBMS triggers:
  - event
  - condition
  - action

. . . . . .

- Integrity constraints are important to express the semantics of databases.
  - prevent the execution of operations or transactions which will cause violation of constraints
  - semantics query optimization
- No real DBMS (database management system) have fully support the management of integrity constraints. i.e. assertions which are part of the SQL standard
- The alternative solution widely adopted by main-stream DBMS triggers:
  - Integrity is spread out among several triggers and therefore the global vision of the semantics is lost.
  - The behavior of triggers is complex (cascading, conflict, etc.).

・ 同 ト ・ ヨ ト ・ ヨ ト

#### We have:

- database *B*.
- database integrity constraint C.
- data modification operation U.

< ∃ >

We have:

- database *B*.
- database integrity constraint C.
- data modification operation U.

We have:

- database *B*.
- database integrity constraint C.
- data modification operation U.

We have:

- database *B*.
- database integrity constraint C.
  - domain constraint
  - 2 column constraint
  - table constraint
  - assertion
- data modification operation U.
  - INSERT
  - 2 DELETE
  - UPDATE

We focus on relational databases and SQL.

We have:

- database *B*.
- database integrity constraint C.
- data modification operation U.

Before U is executed,  $\mathscr{B} \models C$ What we want to prove: After U is executed,  $U(\mathscr{B}) \models C$ .

U preserves C, or U is safe with respect to C, if for any database  $\mathscr{B}$ :

$$\mathscr{B}\models \mathsf{C}\Rightarrow U(\mathscr{B})\models\mathsf{C}$$

We have:

- database *B*.
- database integrity constraint C.
- data modification operation U.
- violation detection: check after the execution of data modification operations.
   violation occurs → rollback
   efficiency problem (check at run-time, rollback, etc.)
- violation prevention: check before the execution of data modification operations.
   violation occurs → abort

Our method

- violation prevention
- mainly based on weakest precondition approach

Hoare triple:

- precondition P
- command C
- postcondition Q

#### $\{P\} \in \{Q\}$

When the precondition is met, the command establishes the postcondition.

In database:

Let f, g be logical formulae, U be a data modification operation, then g is the precondition of U and f if for any database  $\mathscr{B}$ :

$$\mathscr{B}\models g \Rightarrow U(\mathscr{B})\models f$$

Precondition wpc(f, U) is the weakest precondition of f and U if for any database  $\mathscr{B}$  and any precondition g:

$$\mathscr{B} \models g \Rightarrow \mathscr{B} \models wpc(f, U)$$

Now what we need to prove:

$$C \Rightarrow wpc(C, U)$$

$$\begin{cases} C \Rightarrow wpc(C, U) \\ \mathscr{B} \models C \end{cases} \end{cases} \Rightarrow \mathscr{B} \models wpc(C, U) \\ \mathscr{B} \models wpc(C, U) \Rightarrow U(\mathscr{B}) \models C \end{cases} \end{cases} \Rightarrow U(\mathscr{B}) \models C$$

-

Firstly, we need to formalize:

- database integrity constraint C
- data modification operation U

#### Database integrity constraints

#### database integrity constraint in SQL

- domain constraint
- 2 column constraint
- table constraint
- assertion

All the SQL integrity constraints can be reduced to logically and semantically equivalent assertions.

Assertions alone are sufficient to expressing any kind of SQL integrity constraints.

### Therefore, we reduce database integrity constraints into SQL assertions:

CREATE	ASSERTION	$< assertion \ name >$
	CHECK ·	<exists predicate=""></exists>

<exists predicate=""> <query expression=""></query></exists>	::= ::=	
<search condition=""></search>	::=	<boolean term=""></boolean>
		<search condition $>$ OR $<$ boolean term $>$
<boolean term $>$	::=	<boolean factor=""></boolean>
		<boolean term=""> AND <boolean factor=""></boolean></boolean>
<boolean factor $>$	::=	<predicate></predicate>
		[NOT] ( < search condition > )

∃ ▶ ∢

э

## Therefore, we reduce database integrity constraints into SQL assertions:

<predicate></predicate>	::=	<exists predicate=""></exists>
		<comparison predicate=""></comparison>
	i	 between predicate>
	i	<in predicate=""></in>
<comparison predicate=""></comparison>	::=	<expression<sub>1 &gt; <comp op=""> <expression<sub>2 &gt;</expression<sub></comp></expression<sub>
<comp op=""></comp>	::=	$=  \langle \rangle   \langle \rangle   \leq  \rangle   \geq$
<expression></expression>	::=	<term></term>
		$\langle expression \rangle \{+ \mid -\} \langle term \rangle$
<term></term>	::=	<factor></factor>
		$\langle \text{term} \rangle \{ * \mid / \} \langle \text{factor} \rangle$
<factor></factor>	::=	( <expression>)</expression>
		[+   -] < constant >
		[+   -] x.a
 between predicate>	::=	<expression> [ NOT ]</expression>
		BETWEEN $< constant_1 > AND < constant_2 >$
<in predicate=""></in>	::=	<pre><expression> [ NOT ] IN ( <in list="" value=""> )</in></expression></pre>
<in list="" value=""></in>	::=	< constant >
		<in list="" value="">, <constant></constant></in>

• = • •

э

Database integrity constraint

#### database integrity constraints

Database integrity constraint

# database integrity constraints

Database integrity constraint



#### Database integrity constraint

A logical formula, the target language we define, can contain:

- constants (0, 1, ..., etc.)
- relation variables (*r*, etc.), tuple variables (*x*, etc.) and attribute variables (*a*, etc.)
- numerical symbols (+, -, \*, /, ..., etc.)
- comparison symbols (=,  $\neq$ , <, >,  $\leq$ ,  $\geq$ , ..., etc.) and set operators ( $\in$ )
- logical connective symbols ( $\land$ ,  $\lor$ ,  $\neg$ )
- quantifiers ( $\forall$ ,  $\exists$ )

#### Database integrity constraint

Terms of a logical formula are defined as follows:

- All constants and variables are terms.
- If x is a tuple (variable) and a is an attribute of x, then x.a is a term.
- If p is a term and q is a term, then "p <numerical symbol> q" is a term.

#### Database integrity constraint

A logical formula is constructed according to the following rules:

- If x is a tuple and r is a relation variable, then " $x \in r$ " is a formula.
- If *p* is a term and *q* is a term, then "*p* <comparison symbol> *q*" is a formula.
- If  $\phi$  and  $\varphi$  are formulae, then  $\phi \wedge \varphi$ ,  $\phi \lor \varphi$ ,  $\neg \phi$  are formulae.
- If  $\phi$  is a formula and x is a tuple variable, then  $\forall x, \phi$  and  $\exists x, \phi$  are formulae.

#### Translate SQL assertions into FOL formulae

 $\mathcal{T}:$  function mapping a SQL assertion phrase into a first-order logical formula.

A simple example:

$$\mathcal{T}[ CREATE ASSERTION example CHECK (NOT EXISTS (SELECT * FROM  $r x$   
 WHERE  $x.a = 1))]$$$

#### Translate SQL assertions into FOL formulae

# $\mathcal{T}[ \begin{array}{c} \mathsf{CREATE} \ \mathsf{ASSERTION} \ \mathsf{example} \\ \mathsf{CHECK} \ ( \ \mathsf{NOT} \ \mathsf{EXISTS} \ ( \ \mathsf{SELECT} \ \ast \\ \mathsf{FROM} \ r \ x \\ \mathsf{WHERE} \ x.a = 1))] \\ \mathsf{example:} \ \mathcal{T}[\mathsf{NOT} \ \mathsf{EXISTS} \ ( \ \mathsf{SELECT} \ \ast \\ & & & & & & & \\ \mathsf{FROM} \ r \ x \\ \mathsf{WHERE} \ x.a = 1))] \end{array}$

#### Translate SQL assertions into FOL formulae

```
\mathcal{T}[ CREATE ASSERTION example
CHECK (NOT EXISTS (SELECT *
FROM <math>r x
WHERE x.a = 1))]
example: \mathcal{T}[NOT EXISTS (SELECT *
\rightarrow
FROM r x
WHERE x.a = 1))]
example: \neg \exists (\mathcal{T}[SELECT * FROM r x)
WHERE x.a = 1])
```

#### Translate SQL assertions into FOL formulae

 $\mathscr{T}$  CREATE ASSERTION example CHECK (NOT EXISTS (SELECT \* FROM rx WHERE x = 1example: *I*[NOT EXISTS ( SELECT \* FROM r x  $\sim$ WHERE x = 1example:  $\neg \exists (\mathscr{T} | \mathsf{SELECT} * \mathsf{FROM} r x)$  $\sim$ WHERE x a = 1example:  $\neg \exists (x \in r; \mathscr{T}[x.a = 1])$  $\sim \rightarrow$ 

#### Translate SQL assertions into FOL formulae

$$\mathcal{T}[ CREATE ASSERTION exampleCHECK (NOT EXISTS (SELECT *FROM r xWHERE x.a = 1))]example:  $\mathcal{T}[NOT EXISTS (SELECT *$   
 $\rightarrow$   
FROM r x  
WHERE x.a = 1))]  
example:  $\neg \exists (\mathcal{T}[SELECT * FROM r x + SELECT * SEL$$$

#### Translate SQL assertions into FOL formulae

$$\mathcal{T}[ CREATE ASSERTION exampleCHECK (NOT EXISTS (SELECT *FROM r xWHERE x.a = 1))]example:  $\mathcal{T}[NOT EXISTS (SELECT *FROM r xWHERE x.a = 1))]example:  $\neg \exists (\mathcal{T}[SELECT * FROM r x)$   
WHERE x.a = 1])  
 $\Rightarrow$  example:  $\neg \exists (x \in r; \mathcal{T}[x.a = 1])$   
 $\Rightarrow$  example:  $\neg \exists (x \in r; \mathcal{T}[x.a] \mathcal{T}[=] \mathcal{T}[1])$   
 $\Rightarrow$  example:  $\neg \exists (x \in r; x.a = 1)$$$$

(B) ► < B ►</p>

3

#### Semantics of SQL data modification operations

#### SQL INSERT

#### INSERT INTO r VALUES t

Semantics:

$$(x \in r \lor x = t) \Leftrightarrow x \in U(r)$$

#### Semantics of SQL data modification operations

#### SQL DELETE

#### DELETE FROM r WHERE <search condition>

Let g be the logical formula of <search condition>. Semantics:

$$(x \in r \land \neg g(x)) \Leftrightarrow x \in U(r)$$

#### Semantics of SQL data modification operations

#### SQL UPDATE

UPDATE r SET <set clause> WHERE <search condition>

Let  $\sigma$  be the assignment function defined by <set clause> and g be the logical formula of <search condition>. Semantics:

$$\forall x \in r : \begin{cases} g(x) \Rightarrow U(x) = \sigma(x) \\ \neg g(x) \Rightarrow U(x) = x \end{cases}$$

#### Why3

Why3 is a set of tools for program verification which uses first-order logic. input: programs output: logical declarations + goals, in the syntax of the selected prover

#### Why3

Given a constraint C and a date modification operation U:

- Translate SQL statements into WhyML programs.
- Use Why3 to generate the weakest precondition wpc(C, U).
- Call the provers to prove  $C \Rightarrow wpc(C, U)$ , if it is proven, then  $U(\mathscr{B}) \models C$ .

#### Translation Overview



- - E + - E +

Our former experiments show that our methods cannot detect some safe date modification operators when the integrity constraint C is in the form of:

$$\exists x \in R, f(x)$$

Before U is executed:

$$r \models C(r) \equiv \exists x \in r, f(x)$$

What we want to proof:

$$U(r) \models C(U(r)) \equiv \exists x \in U(r), f(x)$$

We adopt the predicate transformer to improve it: e.g. backward predicate transformer  $\overleftarrow{U}$ 

$$\mathscr{B}\models\stackrel{\leftarrow}{U}(\mathcal{C})\Rightarrow U(\mathscr{B})\models \mathcal{C}$$

Use this method, our problem becomes:  $C \Rightarrow \stackrel{\leftarrow}{U} (C)$ 

#### For the constraints in the following form:

 $\exists x \in R, f(x)$ 

We find precise predicate transformer C':

$$\mathscr{B}\models \mathsf{C}'\Leftrightarrow \mathsf{U}(\mathscr{B})\models\mathsf{C}$$

#### Theorem (INSERT)

The execution of any SQL INSERT statements will not affect the constraint C.

#### Theorem (DELETE)

 $U(r) \models C \Leftrightarrow \exists x \in r, f(x) \land \neg g(x)$ 

#### Theorem (UPDATE)

 $U(r) \models C \Leftrightarrow \exists x \in r, (f(x) \land \neg g(x)) \lor (g(x) \land f(\sigma(x)))$ 

Similar results can be derived when the constraint C is in the form of:

$$\forall x \in R, f(x)$$

Before U is executed:

$$r \models C(r) \equiv \forall x \in r, f(x)$$

What we want to proof:

$$U(r) \models C(U(r)) \equiv \forall x \in U(r), f(x)$$

#### Theorem (INSERT)

 $U(r) \models C \Leftrightarrow f(t)$ 

#### Theorem (DELETE)

The execution of any SQL DELETE statements will not affect the constraint C.

#### Theorem (UPDATE)

$$U(r) \models C \Leftrightarrow \forall x \in r, \neg g(x) \lor (g(x) \land f(\sigma(x)))$$

Image: A Image: A

#### Experiments

#### Provers

Alt-Ergo, CVC3, Yices and Gappa.

Unsafe data modification operations can be detected correctly. For safe data modification operations:

	INSERT	DELETE	UPDATE
least mean time (s)	0.205	2.338	1.149

#### Experiments

When the integrity constraint C is in the form of:

 $\exists x \in R, f(x)$ 

We try a newer version of Why3 (0.72) and find that the weakest precondition approach can detect safe data modification operations.

	INSERT	DELETE	UPDATE
weakest precondition (s)	0.208	0.119	0.381
precise predicate transformer (s)	0	0.047	0.054

#### Experiments

#### When the integrity constraint C is in the form of:

#### $\forall x \in R, f(x)$

	INSERT	DELETE	UPDATE
weakest precondition (s)	0.199	0.16	0.323
precise predicate transformer (s)	0.056	0	0.055

< ≣ > <

#### Conclusion & Prospect

Our work:

- rules of reducing SQL database integrity constraints to assertions.
- translating SQL assertions into first-order logical formulae.
- semantics of SQL data modification operations (INSERT/DELETE/UPDATE).
- using weakest precondition and Why3 to implement automated verification of database integrity constraints.
- precise predicate transformer.

#### Conclusion & Prospect

Future work:

- more complex SQL grammar.
- aggregate functions.