

Symbolic Query Exploration

Pavel Grigorenko

(joint work with Margus Veanes, Peli de Halleux and
Nikolai Tillmann)

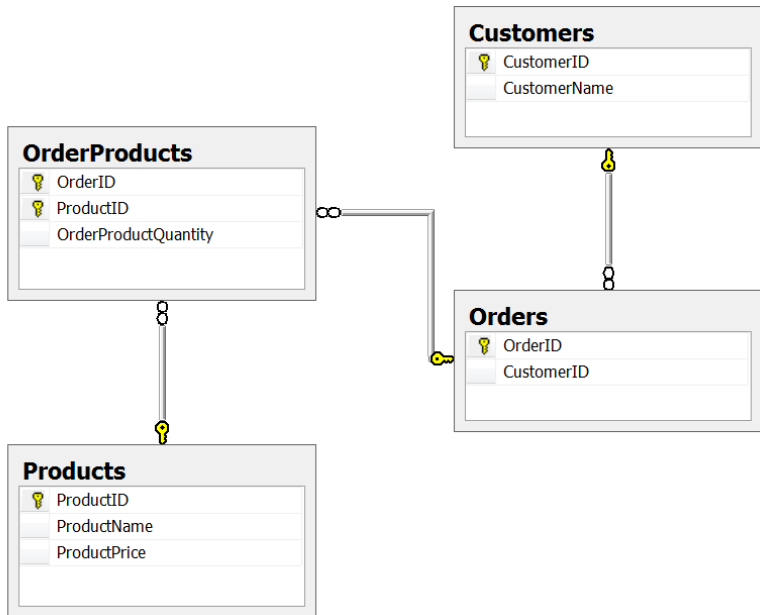
Institute of Cybernetics
Tallinn University of Technology

TSEM@IoC
11 June 2009



- Motivation
 - Test data generation for SQL queries
 - Parameter generation for queries and store procedures
- Goal
 - Investigate model generation with Z3
- Approach
 - Satisfiability modulo background theory \mathcal{T}^Σ
 - Mapping SQL to \mathcal{T}^Σ

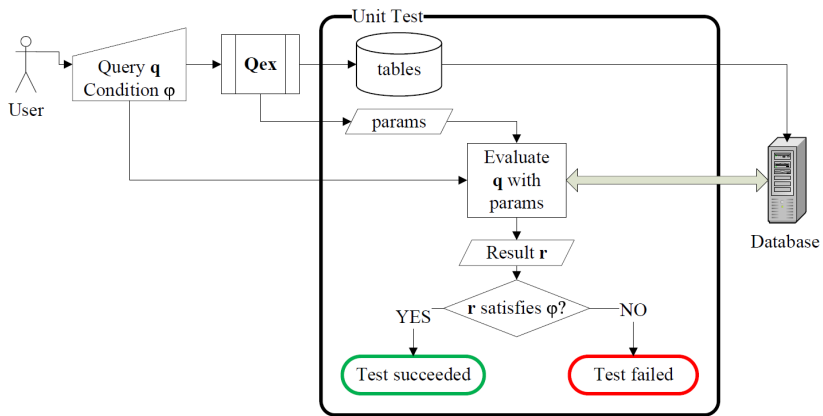




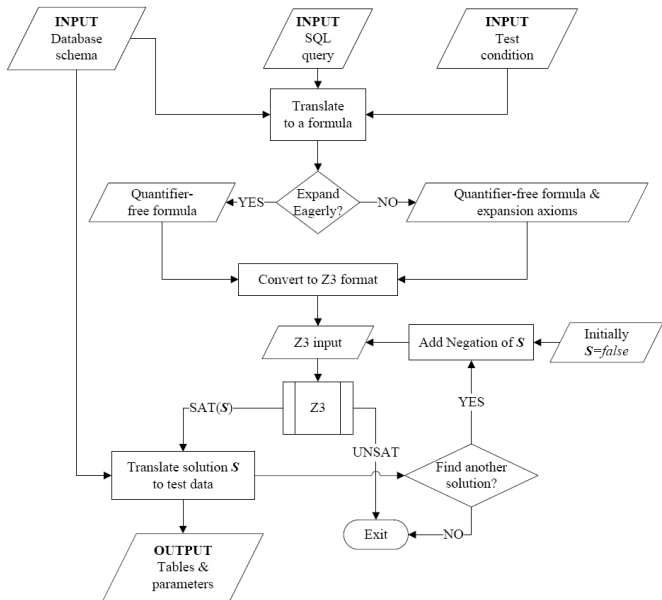
```
SELECT C.CustomerID, SUM(OP.OrderProductQuantity
                        * P.ProductPrice)
FROM OrderProducts AS OP
JOIN Orders AS O ON OP.OrderID = O.OrderID
JOIN Products AS P ON OP.ProductID = P.ProductID
JOIN Customers AS C ON O.CustomerID = C.CustomerID
WHERE @value > 1
GROUP BY C.CustomerID
HAVING SUM(OP.OrderProductQuantity
            * P.ProductPrice) > 100 + @value
```



Unit test



Qex internal workflow



Expressions in \mathcal{T}^Σ

T^σ ::= x^σ | $Default^\sigma$ | $lte(T^{\mathbb{B}}, T^\sigma, T^\sigma)$ | $TheElementOf(T^{\mathbb{S}(\sigma)})$ | $\pi_i(T^{\mathbb{T}(\sigma_0, \dots, \sigma_i = \sigma, \dots)})$

$T^{\mathbb{T}(\sigma_0, \dots, \sigma_k)}$::= $\langle T^{\sigma_0}, \dots, T^{\sigma_k} \rangle$

$T^{\mathbb{Z}}$::= k | $T^{\mathbb{Z}} + T^{\mathbb{Z}}$ | $k * T^{\mathbb{Z}}$ | $\sum_i (T^{\mathbb{S}(\mathbb{T}(\sigma_0, \dots, \sigma_i = \mathbb{Z}, \dots))})$

$T^{\mathbb{R}}$::= r | $T^{\mathbb{R}} + T^{\mathbb{R}}$ | $k * T^{\mathbb{R}}$ | $\sum_i (T^{\mathbb{S}(\mathbb{T}(\sigma_0, \dots, \sigma_i = \mathbb{R}, \dots))})$ | $AsReal(T^{\mathbb{Z}})$

$T^{\mathbb{B}}$::= $true$ | $false$ | $\neg T^{\mathbb{B}}$ | $T^{\mathbb{B}} \wedge T^{\mathbb{B}}$ | $T^{\mathbb{B}} \vee T^{\mathbb{B}}$ | $T^\sigma = T^\sigma$ | $T^{\mathbb{S}(\sigma)} \subseteq T^{\mathbb{S}(\sigma)}$ | $T^\sigma \in T^{\mathbb{S}(\sigma)}$ | $T^{\mathbb{Z}} \leq T^{\mathbb{Z}}$ | $T^{\mathbb{R}} \leq T^{\mathbb{R}}$

$T^{\mathbb{S}(\sigma)}$::= $X^{\mathbb{S}(\sigma)}$ | $\{T^\sigma \mid_{\bar{x}} T^{\mathbb{B}}\}$ | $T^{\mathbb{S}(\sigma)} \cup T^{\mathbb{S}(\sigma)}$ | $T^{\mathbb{S}(\sigma)} \cap T^{\mathbb{S}(\sigma)}$ | $T^{\mathbb{S}(\sigma)} \setminus T^{\mathbb{S}(\sigma)}$

F ::= $T^{\mathbb{B}}$ | $\exists x F$ | $\exists X F$

- S is a state for a term t such that $FV(t) \subseteq Dom(S)$
- For a given term t , t^S is the *interpretation* of t in S
- For a given formula φ , $S \models \varphi$ means that φ^S is *true*

$$lte(\varphi, t_1, t_2)^S = \begin{cases} t_1^S, & \text{if } S \models \varphi; \\ t_2^S, & \text{otherwise.} \end{cases}$$

$$\{t_0 \mid_{x^\sigma} \varphi\}^S = \{t_0^{S \uplus \{x \mapsto a\}} : a \in \mathcal{U}^\sigma, S \uplus \{x \mapsto a\} \models \varphi\}$$

$$\Sigma_i(t_1)^S = \sum_{a \in t_1^S} \pi_i(a)$$

- S is a state for a term t such that $FV(t) \subseteq Dom(S)$
- For a given term t , t^S is the *interpretation* of t in S
- For a given formula φ , $S \models \varphi$ means that φ^S is *true*

$$lte(\varphi, t_1, t_2)^S = \begin{cases} t_1^S, & \text{if } S \models \varphi; \\ t_2^S, & \text{otherwise.} \end{cases}$$

$$\{t_0 \mid_{x^\sigma} \varphi\}^S = \{t_0^{S \uplus \{x \mapsto a\}} : a \in \mathcal{U}^\sigma, S \uplus \{x \mapsto a\} \models \varphi\}$$

$$\Sigma_i(t_1)^S = \sum_{a \in t_1^S} \pi_i(a)$$

Example

Integer multiplication:

$$n * m \stackrel{\text{def}}{=} \Sigma_0(\{\langle m, x \rangle \mid 0 \leq x < n\}) = \sum_{x=0}^{n-1} \pi_0(\langle m, x \rangle) = \sum_{x=0}^{n-1} m$$



Tables are *bags* (multisets)!

A bag b with elements $\{a_i\}_{i < n}$ each having *multiplicity* $m_i > 0$ in b for $i < n$, is represented as a set of pairs $\{\langle a_i, m_i \rangle\}_{i < n}$, thus having the sort $\mathbb{S}(\mathbb{T}(\sigma, \mathbb{Z}))$ for some basic sort σ (*domain sort of b*).

Let $\mathbb{M}(\sigma)$ be the type $\mathbb{S}(\mathbb{T}(\sigma, \mathbb{Z}^+))$ with the *constraint*:

$$\forall X^{\mathbb{M}(\sigma)} \forall x^\sigma y^\sigma ((x \in X \wedge y \in X \wedge x.0 = y.0) \Rightarrow x.1 = y.1).$$

Tables are *bags* (multisets)!

A bag b with elements $\{a_i\}_{i < n}$ each having *multiplicity* $m_i > 0$ in b for $i < n$, is represented as a set of pairs $\{\langle a_i, m_i \rangle\}_{i < n}$, thus having the sort $\mathbb{S}(\mathbb{T}(\sigma, \mathbb{Z}))$ for some basic sort σ (*domain sort of b*).

Let $\mathbb{M}(\sigma)$ be the type $\mathbb{S}(\mathbb{T}(\sigma, \mathbb{Z}^+))$ with the *constraint*:

$$\forall X^{\mathbb{M}(\sigma)} \forall x^\sigma y^\sigma ((x \in X \wedge y \in X \wedge x.0 = y.0) \Rightarrow x.1 = y.1).$$

$$\text{AsBag}(Y^{\mathbb{S}(\sigma)}) \stackrel{\text{def}}{=} \{\langle y, 1 \rangle \mid y \in Y\}$$

$$\text{AsSet}(X^{\mathbb{M}(\sigma)}) \stackrel{\text{def}}{=} \{y.0 \mid y \in X\}$$

$$\Sigma_i^b(X^{\mathbb{M}(\mathbb{T}(\sigma_0, \dots, \sigma_i, \dots))}) \stackrel{\text{def}}{=} \Sigma_0(\{\langle x.1 * x.0.i, x.0 \rangle \mid x \in X\})$$



Example

$$\varphi[x^{\mathbb{Z}}] = x < 4$$

$$q[X^{\mathbb{M}(\mathbb{T}(\mathbb{Z}, \mathbb{Z}, \mathbb{Z}))}] = \{ \langle x.0.0, \Sigma_1^b(\{y \mid y \in X \wedge x.0.0 = y.0.0 \wedge \varphi[y.0.2]\}) \rangle \mid x \in X \wedge \varphi[x.0.2] \}$$

$$t = \{ \langle \langle 0, 2, 1 \rangle, 2 \rangle, \langle \langle 1, 2, 3 \rangle, 1 \rangle, \langle \langle 1, 2, 4 \rangle, 1 \rangle \}$$

$$\begin{aligned} q[t] &= \{ \langle x.0.0, \Sigma_1^b(\{y \mid y \in t \wedge x.0.0 = y.0.0 \wedge \varphi[y.0.2]\}) \rangle \mid x \in t \wedge \varphi[x.0.2] \} \\ &= \{ \langle 0, \Sigma_1^b(\{y \mid y \in t \wedge 0 = y.0.0 \wedge \varphi[y.0.2]\}) \rangle, \\ &\quad \langle 1, \Sigma_1^b(\{y \mid y \in t \wedge 1 = y.0.0 \wedge \varphi[y.0.2]\}) \rangle \} \\ &= \{ \langle 0, \sum_{a \in \{ \langle \langle 0, 2, 1 \rangle, 2 \rangle \}} \pi_1(a) * \pi_1(\pi_0(a)) \rangle, \\ &\quad \langle 1, \sum_{a \in \{ \langle \langle 1, 2, 3 \rangle, 1 \rangle \}} \pi_1(a) * \pi_1(\pi_0(a)) \rangle \} \\ &= \{ \langle 0, 4 \rangle, \langle 1, 2 \rangle \} \end{aligned}$$



Translation

$\mathbf{Q} : \text{SQL} \rightarrow \mathcal{T}^\Sigma$

Strings

- All strings have a maximum length k ; an encoding of a k -string is as a k -tuple of integers, each character a is encoded as an integer $c(a)$. For a string $a_0 \cdots a_l, l < k$ the encoding is $\langle c(a_0), \dots, c(a_l), 0, \dots, 0 \rangle$.
- For a collection D of strings, encode as $|D|$ -enums.

Nullable values

Given a basic sort σ , let $?\sigma$ be the sort $\mathbb{T}(\sigma, \mathbb{B})$ with the constraint $\forall x^{?\sigma} (x.1 = \text{false} \Rightarrow x.0 = \text{Default}^\sigma)$ and $\text{null}^{?\sigma} \stackrel{\text{def}}{=} \text{Default}^{\mathbb{T}(\sigma, \mathbb{B})}$.

Operations that are defined for σ are lifted to $?\sigma$.

For example, for a numeric sort σ ,

$$x^{?\sigma} + y^{?\sigma} \stackrel{\text{def}}{=} \text{lte}(x.1 \wedge y.1, \langle x.0 + y.0, \text{true} \rangle, \text{null}^{?\sigma}).$$



Select clauses

$$\mathbf{Q}(\text{SELECT } l \text{ FROM } t) \stackrel{\text{def}}{=} \{ \langle \langle x.0.l_0, \dots, x.0.l_n \rangle, M(x) \rangle \mid x \in \mathbf{Q}(t) \}$$

where $M(x) = \Sigma_0(\{ \langle y.1, y \rangle \mid y \in \mathbf{Q}(t) \wedge \bigwedge_{i=0}^n y.0.l_i = x.0.l_i \})$

$$\mathbf{Q}(\text{SELECT DISTINCT } l \text{ FROM } t) \stackrel{\text{def}}{=} \\ \text{AsBag}(\text{AsSet}(\mathbf{Q}(\text{SELECT } l \text{ FROM } t)))$$

$$\text{AsSet}(\mathbf{Q}(\text{SELECT } l \text{ FROM } t)) = \\ \{ \langle y.l_0, \dots, y.l_n \rangle \mid y \in \text{AsSet}(\mathbf{Q}(t)) \}$$


Join clauses

$$\mathbf{Q}(t1 \text{ INNER JOIN } t2 \text{ ON } c) \stackrel{\text{def}}{=} \{ \langle x1.0 \times x2.0, x1.1 * x2.1 \rangle \mid x1 \in \mathbf{Q}(t1) \wedge x2 \in \mathbf{Q}(t2) \wedge \mathbf{Q}(c)[x1.0, x2.0] \}$$

$$x \times y \stackrel{\text{def}}{=} \langle \pi_0(x), \dots, \pi_{m-1}(x), \pi_0(y), \dots, \pi_{n-1}(y) \rangle$$

$$\mathbf{AsSet}(\mathbf{Q}(t1 \text{ INNER JOIN } t2 \text{ ON } c)) = \{ y1 \times y2 \mid y1 \in \mathbf{AsSet}(\mathbf{Q}(t1)) \wedge y2 \in \mathbf{AsSet}(\mathbf{Q}(t2)) \wedge \mathbf{Q}(c)[y1, y2] \}$$



Grouping and aggregates

$t = \text{SELECT } a, \text{ SUM}(b) \text{ AS } d \text{ FROM } t1 \text{ WHERE } c1$

$\mathbf{Q}(t \text{ GROUP BY } a \text{ HAVING } c2) \stackrel{\text{def}}{=} \text{AsBag}(\{z \mid z \in \mathbf{G} \wedge \mathbf{Q}(c2)[z]\})$

where $\mathbf{G} = \{\langle x.0.0, \Sigma_1^b(\{y \mid y \in \mathbf{Q}(t) \wedge y.0.0 = x.0.0\}) \rangle \mid x \in \mathbf{Q}(t)\}$

$\text{Sum} \stackrel{\text{def}}{=} \Sigma_i^b$

$\text{Count} \stackrel{\text{def}}{=} \Sigma_1$

$\text{Min}(X^{\mathbb{S}(\sigma)}) \stackrel{\text{def}}{=} \text{TheElementOf}(\{y \mid y \in X \wedge \{z \mid z \in X \wedge z < y\} = \emptyset\})$



Union

$$\mathbf{Q}(q1 \text{ UNION } q2) \stackrel{\text{def}}{=} \mathit{AsBag}(\mathit{AsSet}(\mathbf{Q}(q1)) \cup \mathit{AsSet}(\mathbf{Q}(q2))).$$


Simplifications

$$\text{AsSet}(\text{AsBag}(X^{\mathbb{S}(\sigma)})) = X$$

$$\Sigma_i^b(\text{AsBag}(X^{\mathbb{S}(\sigma)})) = \Sigma_i(X)$$

$$\text{AsSet}(\{t \mid \varphi\}^{\mathbb{M}(\sigma)}) = \{t.0 \mid \varphi\}$$

$$\pi_i(\langle t_0, \dots, t_i, \dots \rangle) = t_i$$



Model Generation in $\mathcal{T}_{\mathbf{Q}}^{\Sigma}$

Given a quantifier free formula $\varphi[X]$ in \mathcal{T}^{Σ} , and a query q , decide if $\psi = \varphi[\mathbf{Q}(q)]$ is satisfiable, and if ψ is satisfiable, generate a model of ψ .

Two approaches to deal with *comprehensions* and *summations* in Z3

- Eager expansion
 - Provide finite bounds for tables and unwind
- Lazy expansion
 - Add expansion rules as axioms



- 1 Given a query $q[X]$
- 2 Create a symbolic table $t_X = \{\langle x_1, m_1 \rangle, \dots, \langle x_k, m_k \rangle\}$
(k and m_i are fixed, x_i are variables)
- 3 Expand $q[t_X]$ to $\varphi = \mathbf{Exp}(q[t_X])$
- 4 Generate a model for φ (if φ is sat.)
- 5 Increase k and m_i , repeat 1..4.



Set describer

Consider a formula $\psi[\bar{X}]$ as an instance of the model generation problem, where every X in \bar{X} is a bag variable.

- The constant $Empty^{S(\sigma)}$ is a set describer.
- If $t^{S(\sigma)}$ is a set describer then so is the term $Set(\varphi^{\mathbb{B}}, u^\sigma, t)$.

Given a state S for $Set(\varphi, u, t)$, the interpretation in S is,

$$Set(\varphi, u, t)^S = lte(\varphi, \{u\}, \emptyset)^S \cup t^S, \quad Empty^S = \emptyset.$$

X is fixed in \bar{X} and t_X is the set describer:

$$Set(true, \langle x_1, m_1 \rangle, \dots, Set(true, \langle x_k, m_k \rangle, Empty) \dots)$$

where k and all the m_i 's are some positive integer constants and each x_i is a variable.



Expansion rule for *comprehensions*

$$\mathbf{Exp}(\{t \mid_x x \in r \wedge \varphi\}) \stackrel{\text{def}}{=} \mathbf{ExpC}(t, x, \mathbf{Exp}(r), \varphi)$$

$$\mathbf{ExpC}(t, x, \mathbf{Empty}, \varphi) \stackrel{\text{def}}{=} \mathbf{Empty}$$

$$\mathbf{ExpC}(t[x], x, \mathbf{Set}(\gamma, u, \mathit{rest}), \varphi[x]) \stackrel{\text{def}}{=} \mathbf{Set}(\gamma \wedge \mathbf{Exp}(\varphi[u]), \mathbf{Exp}(t[u]), \mathbf{ExpC}(t, x, \mathit{rest}, \varphi))$$



Expansion rule for *comprehensions*

$$\mathbf{Exp}(\{t \mid_x x \in r \wedge \varphi\}) \stackrel{\text{def}}{=} \mathbf{ExpC}(t, x, \mathbf{Exp}(r), \varphi)$$

$$\mathbf{ExpC}(t, x, \mathbf{Empty}, \varphi) \stackrel{\text{def}}{=} \mathbf{Empty}$$

$$\mathbf{ExpC}(t[x], x, \mathbf{Set}(\gamma, u, \mathit{rest}), \varphi[x]) \stackrel{\text{def}}{=} \mathbf{Set}(\gamma \wedge \mathbf{Exp}(\varphi[u]), \mathbf{Exp}(t[u]), \\ \mathbf{ExpC}(t, x, \mathit{rest}, \varphi))$$

Expansion rule for Σ_i

$$\mathbf{Exp}(\Sigma_i(t)) \stackrel{\text{def}}{=} \mathbf{Sum}_i(\mathbf{Exp}(t), \mathbf{Empty})$$

$$\mathbf{Sum}_i(\mathbf{Empty}, s) \stackrel{\text{def}}{=} 0$$

$$\mathbf{Sum}_i(\mathbf{Set}(\gamma, u, \mathit{rest}), s) \stackrel{\text{def}}{=} \mathit{lte}(\gamma \wedge u \notin s, \pi_i(u), 0) + \\ + \mathbf{Sum}_i(\mathit{rest}, \mathbf{Set}(\gamma, u, s))$$



In addition to a quantifier free formula ψ , universally quantified *axioms* are provided in the form:

$$(\forall \bar{x}(\alpha), pat_{\alpha}), \quad FV(\alpha) = FV(pat_{\alpha}) = \bar{x}$$

If ψ contains a subterm t and there exists a substitution θ such that $t = pat_{\alpha}\theta$, i.e., t matches the pattern pat_{α} , then ψ is replaced during proof search by (a reduction of) $\psi \wedge \alpha\theta$.



Axioms

$$\alpha_1 = \forall s (\mathbf{Sum}_i(\mathit{Empty}, s) = 0)$$

$$\mathit{pat}_{\alpha_1} = \mathbf{Sum}_i(\mathit{Empty}, s)$$

$$\alpha_2 = \forall b u r s (\mathbf{Sum}_i(\mathit{Set}(b, u, r), s) = \\ \mathit{lte}(b \wedge u \notin s, \pi_i(u), 0) + \mathbf{Sum}_i(r, \mathit{lte}(b, \{u\}, \emptyset) \cup s))$$

$$\mathit{pat}_{\alpha_2} = \mathbf{Sum}_i(\mathit{Set}(b, u, r), s)$$

Example

$$x \leq \mathbf{Sum}_1(\mathit{Set}(\mathit{true}, \langle 1, y \rangle, \mathit{Set}(\mathit{true}, \langle 1, z \rangle, \mathit{Empty})), \emptyset)$$

$$\xrightarrow{\alpha_2} x \leq y + \mathbf{Sum}_1(\mathit{Set}(\mathit{true}, \langle 1, z \rangle, \mathit{Empty}), \{\langle 1, y \rangle\})$$

$$\xrightarrow{\alpha_2} x \leq y + \mathit{lte}(z \neq y, z, 0) + \mathbf{Sum}_1(\mathit{Empty}, \{\langle 1, y \rangle, \langle 1, z \rangle\})$$

$$\xrightarrow{\alpha_1} x \leq y + \mathit{lte}(z \neq y, z, 0)$$



Experiments

```
q1: SELECT C.CustomerID, O.OrderID
FROM Orders AS O
JOIN Customers AS C ON
O.CustomerID = C.CustomerID
WHERE O.CustomerID > 2 AND
O.OrderID < 15
```

```
q2: SELECT C.CustomerID,
Count(O.OrderID)
FROM Orders AS O
JOIN Customers AS C ON
O.CustomerID = C.CustomerID
GROUP BY C.CustomerID HAVING
Count(O.OrderID) > 1
```

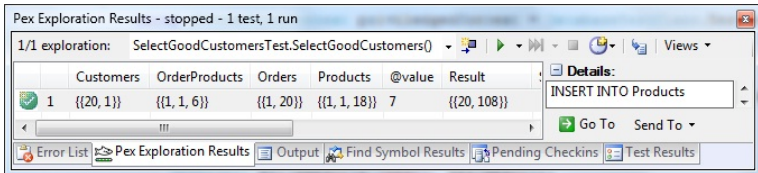
```
q3: DECLARE @value AS INT;
SELECT C.CustomerID, SUM(OP.OrderProductQuantity * P.ProductPrice)
FROM OrderProducts AS OP
JOIN Orders AS O ON OP.OrderID = O.OrderID
JOIN Products AS P ON OP.ProductID = P.ProductID
JOIN Customers AS C ON O.CustomerID = C.CustomerID
WHERE @value > 1
GROUP BY C.CustomerID
HAVING SUM(OP.OrderProductQuantity * P.ProductPrice) > 100 + @value
```

query	cond	k	check	t_{exp}	t_{z3}
q1	$r \neq \emptyset$	1	sat	.03	.001
		2	sat	.05	.005
		3	sat	.3	.02
		4	sat	1.4	.13
	$r = \emptyset$	1	sat	.03	.001
		2	sat	.05	.006
		3	sat	.3	.12
		4	sat	1.4	2

query	cond	k	check	t_{exp}	t_{z3}
q1	$ r = 5$	1	unsat	.03	.001
		2	unsat	.05	.01
		3	unsat	.3	.16
		4	unsat	1.4	10
q2	$r \neq \emptyset$	1	unsat	.03	.001
		2	sat	.7	.006
		3	sat	26	.03
q3	$r \neq \emptyset$	1	sat	.34	.001

Implementation

- **Qex** written in C#, uses **T-SQL** syntax for input queries
- Integration with **Pex**



Pex Exploration Results - stopped - 1 test, 1 run

1/1 exploration: SelectGoodCustomersTest.SelectGoodCustomers()

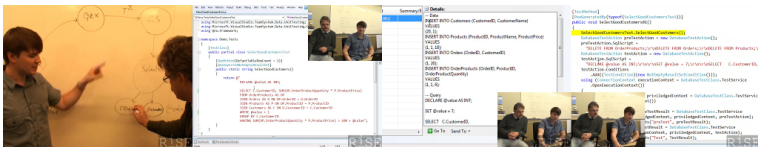
	Customers	OrderProducts	Orders	Products	@value	Result
1	{{20, 1}}	{{1, 1, 6}}	{{1, 20}}	{{1, 1, 18}}	7	{{20, 108}}

Details: INSERT INTO Products

Go To Send To

Error List Pex Exploration Results Output Find Symbol Results Pending Checkins Test Results

- Project homepage <http://research.microsoft.com/qex>
- Check **Qex** demo at <http://channel9.msdn.com/posts/Peli/Qex-Symbolic-Query-Exploration/>



The collage includes:

- A whiteboard with a diagram showing nodes like 'Qex', 'Pex', and 'SQL' connected by arrows.
- Code snippets in C# and T-SQL, such as:

```
public static class SelectGoodCustomersTest {  
    [TestMethod]  
    public void SelectGoodCustomers() {  
        // ...  
    }  
}
```

```
INSERT INTO Customers (CustomerID, CustomerName)  
VALUES  
(20, 'Customer 20')  
GO  
INSERT INTO Products (ProductID, ProductName, ProductCost)  
VALUES  
(1, 'Product 1', 10)  
GO  
INSERT INTO Orders (OrderID, CustomerID, ProductID, OrderProductQuantity)  
VALUES  
(1, 20, 1, 1)  
GO  
INSERT INTO OrderProducts (OrderID, ProductID, OrderProductQuantity)  
VALUES  
(1, 1, 6)  
GO
```
- A video thumbnail showing two people discussing the project.



- **Lazy** expansion (done?)
- Improved Database **constraint** support
- Queries with **side-effects**
- Improved **datatype** support
- Query **optimizations**



Thanks! Questions?

