

# Combining Consistency and Confidentiality Requirements in First-Order Databases

Lena Wiese

Fakultät für Informatik  
LS 6 (ISSI)

September 7–9, 2009

# Outline

- 1 Introduction
  - Inference Control
  - Controlled Query Evaluation
  - Preprocessing for CQE
- 2 Automizing Inference-Proofness
- 3 Prototype
- 4 Conclusion

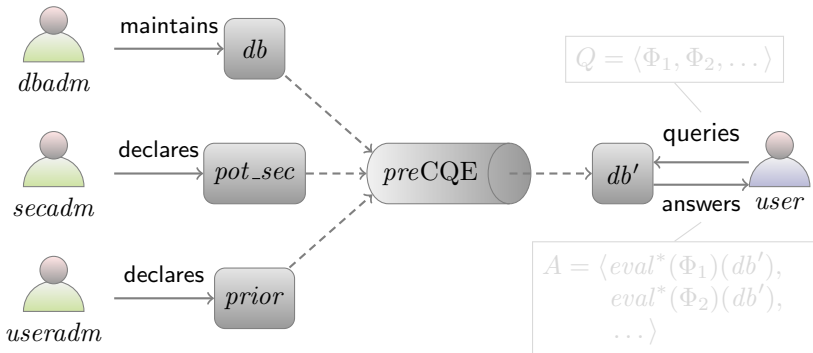
# Inference control

- Protect confidential and private information in database instance *db*
- Personalized security (confidentiality) policy *pot\_sec*
- User profile (a priori knowledge) *prior*
- IC system automatically distorts some answers
  - Avoids harmful user inferences
- Here: modify input database
  - Remove tuples (like Data Privacy; Stouppa/Studer, 2009)
  - Add tuples (like Cover Stories; eg. Galinovic et al, 2007)
- Automatically generate “inference-proof” output instance
- Also related to Data Exchange (eg. Fagin et al, 2005) and Consistent Query Answering (eg. Chomicki, 2007)

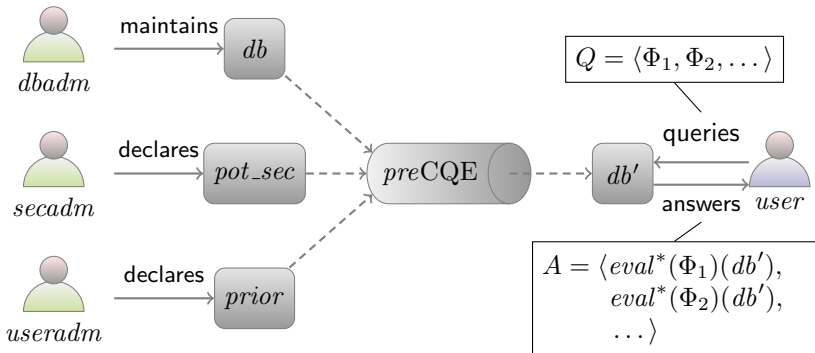
# Prior work: Controlled Query Evaluation (CQE)

- (Biskup/Bonatti, 2007; Biskup/Wiese, 2008)...
- Logical view of relational data model
- Database schema  $DS = \langle \mathcal{P}, \mathcal{D} \rangle$  with relation names  $\mathcal{P}$  and database dependencies  $\mathcal{D}$
- Infinite domain of values (constants)  $dom$
- Complete database instance as finite set of tuples (ground atoms) + closed world assumption
- Relational calculus as query language

# preCQE: Preprocessing for CQE



# preCQE: Preprocessing for CQE



## Preprocessing for CQE

Definition: Inference-proofness of  $db'$

- ① [Consistency]  $I^{db'} \models prior$
- ② [Confidentiality]  $I^{db'} \not\models \Psi$  for every  $\Psi \in pot\_sec$

Definition: Distortion distance (amount of modified tuples)

[Availability]  $db\_dist(db') := card((db \setminus db') \cup (db' \setminus db))$

- Find  $db'$  that satisfies constraint set  $C := prior \cup Neg(pot\_sec)$
- Minimize amount of modified tuples  $db\_dist$  (maximize availability)
- No impact on runtime performance
- No user history (*log* file) has to be stored

# Example

$$\mathcal{P} = \{Ill, Treat\},$$

$$dom = \{Pete, Mary, Lisa, Paul, \dots, \\ Aids, Flu, Cancer, Myopia, \dots \\ MedA, MedB, MedC, \dots\}$$

<i>db:</i>	<i>Ill</i>	<i>Name</i>	<i>Diagnosis</i>		<i>Treat</i>	<i>Name</i>	<i>Treatment</i>
		Pete	Aids			Pete	MedA
		Mary	Cancer			Mary	MedB

$$prior = \{\forall x (Treat(x, MedA) \rightarrow Ill(x, Aids) \vee Ill(x, Cancer)), \\ \forall x (Treat(x, MedB) \rightarrow Ill(x, Cancer) \vee Ill(x, Flu))\}$$

$$pot\_sec = \{\exists x Ill(x, Aids), \exists x Ill(x, Cancer)\}$$



## Example

$$\begin{aligned}
 \text{prior} = & \{ \forall x (Treat(x, MedA) \rightarrow Ill(x, Aids) \vee Ill(x, Cancer)), \\
 & \forall x (Treat(x, MedB) \rightarrow Ill(x, Cancer) \vee Ill(x, Flu)) \}
 \end{aligned}$$

$$\text{pot\_sec} = \{ \exists x Ill(x, Aids), \exists x Ill(x, Cancer) \}$$

$$\text{Neg}(\text{pot\_sec}) = \{ \forall x \neg Ill(x, Aids), \forall x \neg Ill(x, Cancer) \}$$

Constraint set  $C := \text{prior} \cup \text{Neg}(\text{pot\_sec})$

db:

<i>Ill</i>	<i>Name</i>	<i>Diagnosis</i>
	Pete	Aids
	Mary	Cancer

<i>Treat</i>	<i>Name</i>	<i>Treatment</i>
	Pete	MedA
	Mary	MedB

## Example

$$prior = \{\forall x (Treat(x, MedA) \rightarrow Ill(x, Aids) \vee Ill(x, Cancer)), \\ \forall x (Treat(x, MedB) \rightarrow Ill(x, Cancer) \vee Ill(x, Flu))\}$$

$$pot\_sec = \{\exists x Ill(x, Aids), \exists x Ill(x, Cancer)\}$$

$$Neg(pot\_sec) = \{\forall x \neg Ill(x, Aids), \forall x \neg Ill(x, Cancer)\}$$

Constraint set  $C := prior \cup Neg(pot\_sec)$

$db'_1$ :	<i>Ill</i>	<i>Name</i>	<i>Diagnosis</i>
		Pete	Aids
		Mary	Cancer

<i>Treat</i>	<i>Name</i>	<i>Treatment</i>
	Pete	MedA
	Mary	MedB

$$db\_dist(db'_1) = 4$$

## Example

$$\begin{aligned}
 \text{prior} &= \{ \forall x (Treat(x, MedA) \rightarrow Ill(x, Aids) \vee Ill(x, Cancer)), \\
 &\quad \forall x (Treat(x, MedB) \rightarrow Ill(x, Cancer) \vee Ill(x, Flu)) \}
 \end{aligned}$$

$$\text{pot\_sec} = \{ \exists x Ill(x, Aids), \exists x Ill(x, Cancer) \}$$

$$\text{Neg}(\text{pot\_sec}) = \{ \forall x \neg Ill(x, Aids), \forall x \neg Ill(x, Cancer) \}$$

Constraint set  $C := \text{prior} \cup \text{Neg}(\text{pot\_sec})$

$db'_2$ :

<i>Ill</i>	<i>Name</i>	<i>Diagnosis</i>
	Pete	Aids
	Mary	Cancer
	Mary	Flu

<i>Treat</i>	<i>Name</i>	<i>Treatment</i>
	Pete	MedA
	Mary	MedB

$$db\_dist(db'_2) = db\_dist(db'_1) = 4$$

# Outline

- 1 Introduction
- 2 Automizing Inference-Proofness
  - Restricted Constraints
  - *preCQE* Algorithm
- 3 Prototype
- 4 Conclusion

# Undecidability of satisfiability problem

- For complete database, known potential secrets & data modification, find  $db'$  such that  $I^{db'} \models C$  and  $db\_dist(db') \longrightarrow \min$
- Undecidability of satisfiability problem for predicate logic (reproduced in (Börger et al, 2001))
- Identify syntactical restrictions for constraint set  $C$  to make problem decidable
- “Allowed universal formulas” in prenex literal normal form
  - Subset of “allowed formulas” (Van Gelder/Topor, 1991)
  - “Active domain” semantics,  $adom$ : constants in  $db$  and  $C$

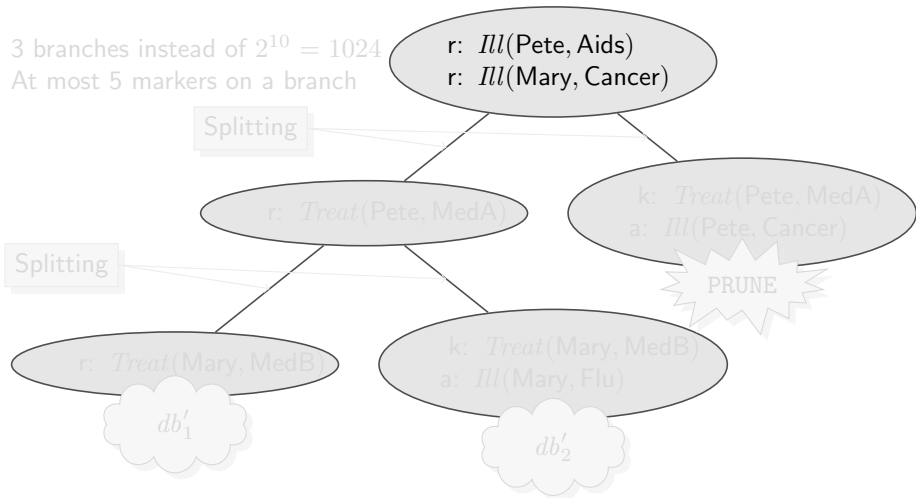
## *preCQE* subprocedures

### Branch-and-Bound depth-first search tree

- INIT
- GROUND
  - find relevant ground instantiations for universal quantifiers
- SIMP
- SPLIT
  - try two truth values for ground atom
- MARK
  - mark ground atom as k(ept), a(dded), r(emoved) or l(eft out)

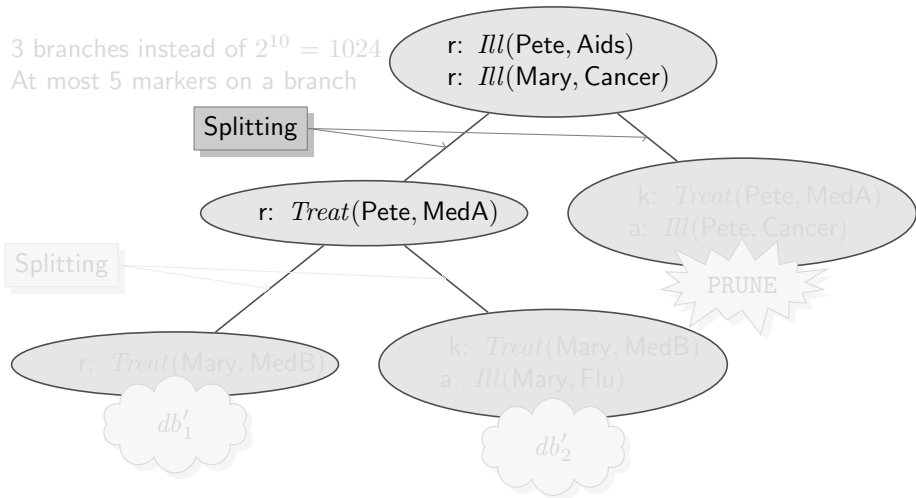
# Example: *preCQE* search tree

3 branches instead of  $2^{10} = 1024$   
 At most 5 markers on a branch



# Example: *preCQE* search tree

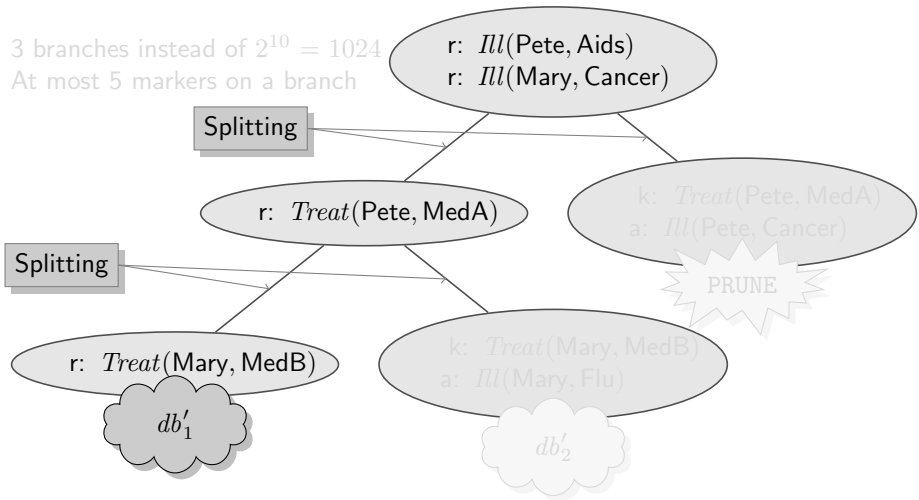
3 branches instead of  $2^{10} = 1024$   
 At most 5 markers on a branch





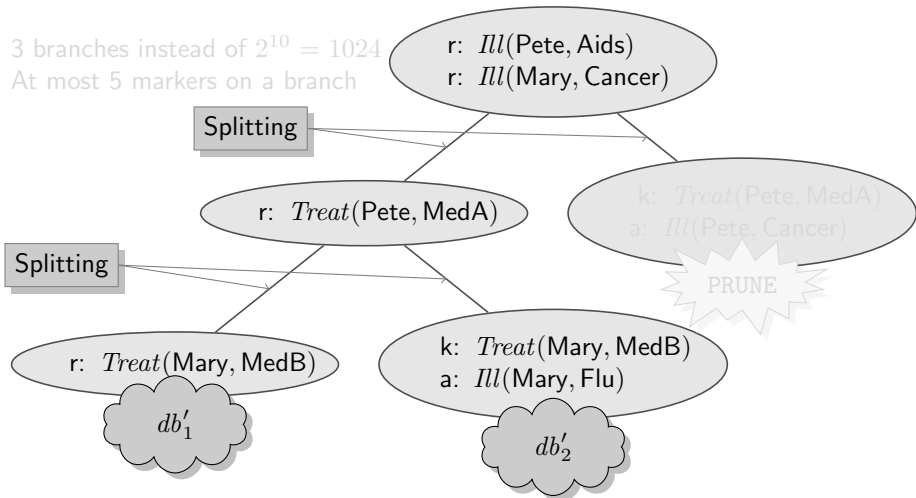
## Example: *preCQE* search tree

3 branches instead of  $2^{10} = 1024$   
At most 5 markers on a branch



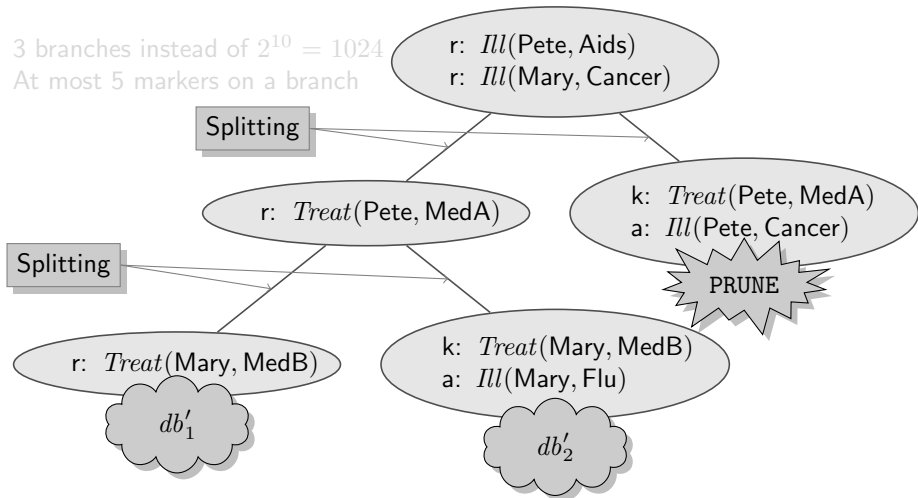
# Example: *preCQE* search tree

3 branches instead of  $2^{10} = 1024$   
 At most 5 markers on a branch



# Example: *preCQE* search tree

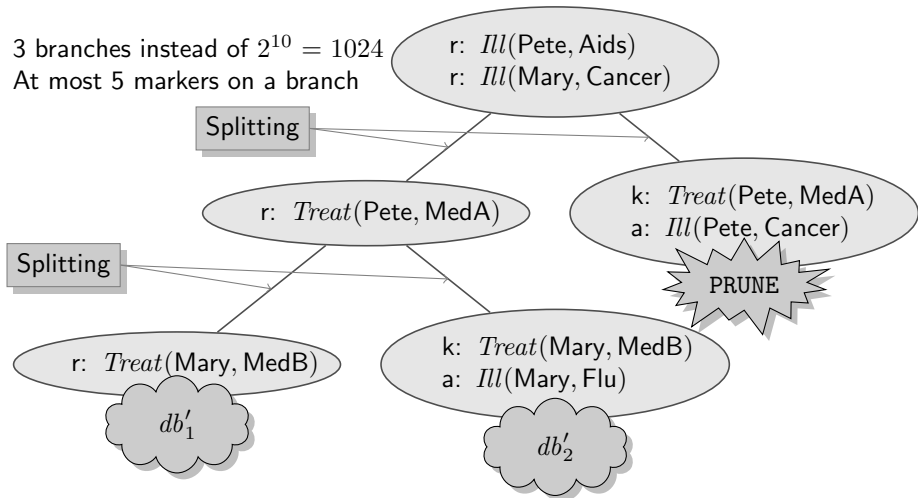
3 branches instead of  $2^{10} = 1024$   
 At most 5 markers on a branch



# Example: *preCQE* search tree

3 branches instead of  $2^{10} = 1024$

At most 5 markers on a branch



# Key results

## Theorem: Termination of *preCQE*

For a set  $C$  of allowed universal constraints, *preCQE* terminates in a finite amount of time

Proof:

- At most  $k$  different ground atoms with *adom* constants where

$$k := \sum_{\substack{P \in \mathcal{P} \\ P \text{ occurs in } C}} \text{card}(\text{adom})^{\text{arity}(P)}$$

- At most  $2^k$  branches in the search tree

## Key results

Theorem: Satisfiability soundness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $db'$  is a database instance, it is inference-proof (hence, a model of  $C$ )

Proof:

- No violated constraints left

Corollary: Refutation completeness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $C$  is unsatisfiable,  $db'$  is undefined

## Key results

Theorem: Satisfiability soundness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $db'$  is a database instance, it is inference-proof (hence, a model of  $C$ )

Proof:

- No violated constraints left

Corollary: Refutation completeness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $C$  is unsatisfiable,  $db'$  is undefined

## Key results

### Theorem: Refutation soundness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $db'$  is undefined, then  $C$  is unsatisfiable

- Not trivial because of efficiency of *preCQE*
- Not all *adom*-ground atoms explicitly handled
  - ① Only violated constraints and affected ground atoms are considered
  - ② If a truth assignment is unequivocal, ground atoms are marked directly without splitting
  - ③ Branches are pruned if a better solution has already been found
  - ④ Branches are pruned as soon as a conflict occurs
- *preCQE* search tree in best case does not contain all possible  $2^k$  branches



## Proof of refutation soundness

Herbrand's Theorem (Herbrand, 1930; here as in Cook/Nguyen, 2009)

Let  $S$  be a set of closed universal formulas. Then  $S$  is unsatisfiable iff some finite set  $S_0$  of ground instances of formulas in  $S$  is propositionally unsatisfiable

Herbrand's Theorem with semantic tree (Chang/Lee, 1973; for clauses)

Let  $S$  be a set of closed universal formulas. Then  $S$  is unsatisfiable iff for some finite set  $S_0$  of ground instances of formulas in  $S$  there is a closed semantic tree

- Construct semantic tree out of *preCQE* search tree
- Identify set  $C_0$  of ground instances
- Show that semantic tree is closed for  $C_0$

## Key results

Corollary: Satisfiability completeness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $C$  is satisfiable,  $db'$  is a database instance

Theorem: Distortion minimality of solution

For a set  $C$  of allowed universal constraints, if *preCQE* finds a solution  $db'$ , then it is distortion-minimal

- For other constraints (existential or weakly acyclic fragments) similar:
  - Depth and width of *preCQE* search tree is bounded
  - Fix mapping of existentially quantified variables to invented constants for refutation soundness

## Key results

Corollary: Satisfiability completeness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $C$  is satisfiable,  $db'$  is a database instance

Theorem: Distortion minimality of solution

For a set  $C$  of allowed universal constraints, if *preCQE* finds a solution  $db'$ , then it is distortion-minimal

- For other constraints (existential or weakly acyclic fragments) similar:
  - Depth and width of *preCQE* search tree is bounded
  - Fix mapping of existentially quantified variables to invented constants for refutation soundness

## Key results

Corollary: Satisfiability completeness of *preCQE*

For a set  $C$  of allowed universal constraints, if  $C$  is satisfiable,  $db'$  is a database instance

Theorem: Distortion minimality of solution

For a set  $C$  of allowed universal constraints, if *preCQE* finds a solution  $db'$ , then it is distortion-minimal

- For other constraints (existential or weakly acyclic fragments) similar:
  - Depth and width of *preCQE* search tree is bounded
  - Fix mapping of existentially quantified variables to invented constants for refutation soundness

# Outline

- 1 Introduction
- 2 Automizing Inference-Proofness
- 3 Prototype**
  - Implementation
  - Test Cases
- 4 Conclusion

# User interface

The screenshot shows the 'input settings solver' window in the CQE Framework. The window title is 'CQE Framework'. The main title is 'input settings solver'. Below the title bar, there are buttons for 'save' and 'solve', and a dropdown menu set to 'MiniMaxSat'. The window is divided into several sections:

- input db:** A list of database entries including 'n23\_flu', 'n22\_aedB', 'n18\_aedB', 'n20\_cancer', 'n3\_flu', 'n13\_flu', 'n12\_cancer', 'n21\_cancer', 'n4\_aids', 'n5\_aids', 'n6\_cancer', 'n21\_aedA', 'n20\_aedB', 'n19\_aids', 'n16\_aids', 'n23\_aedA', 'n21\_aids', 'n12\_aedA', 'n2\_cancer', 'n16\_aedB', 'n17\_aedB', 'n22\_flu', 'n19\_flu', 'n24\_aids', 'n10\_aedA', 'n7\_cancer', 'n9\_cancer', 'n13\_aids', 'n19\_cancer', 'n24\_flu', 'n21\_aedB', 'n10\_aids', 'n14\_cancer', 'n24\_cancer', 'n14\_aedB', 'n11\_aedA', 'n18\_cancer', 'n22\_aids', 'n18\_flu', and 'n24\_aedA'.
- Navigation:** Buttons for 'prior', 'pot\_sec', 'avail', and 'output db'.
- Solutions:** A dropdown menu showing 'solutions: 27/11/08 08:41'.
- Performance Metrics:**

total runtime	0	h	0	min	3	sec	456	msec
runtime solver	0	h	0	min	0	sec	475	msec
- Solver Settings:**

solver:	MiniMaxSat	db_dist	45	avail_dist	11
---------	------------	---------	----	------------	----
- Log Output:**

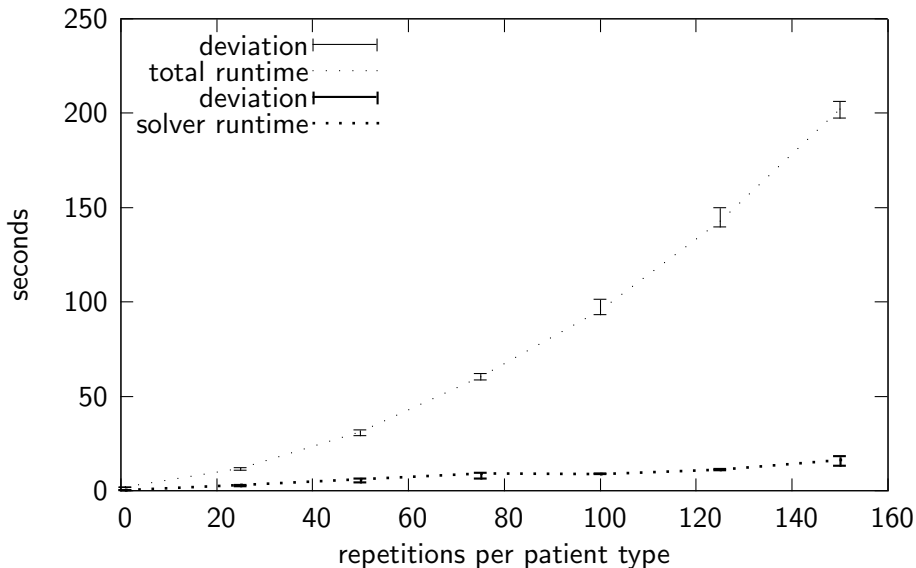
```

%-----
% File : new database : TPTP v3.3.0. Released TPTP V3.3.0
% Domain :
% Problem : CQE : avail1repetition
% Version :
% English :
% Refs :
% Source :
% Status : Axiom
% Rating :
% Syntax :
% Comments : complete database
%-----
n11_flu.
n12_flu.
n13_flu.
n14_flu.
n14_aedB.
n15_flu.
n15_aedB.
n16_flu.
n16_aedB.
n17_flu.
n17_aedB.
n18_flu.
n18_aedB.
n19_flu.

```

## Average test results

rep.	avg. msec total	avg. msec solver	dec. vars.	clauses	
				soft	hard
1	1930	184	120	120	96
25	11974	3092	3000	3000	2400
50	31304	6135	6000	6000	4800
75	60459	8991	9000	9000	7200
100	95792	8902	12000	12000	9600
125	142843	11171	15000	15000	12000
150	202067	16429	18000	18000	16800





# Outline

- 1 Introduction
- 2 Automizing Inference-Prooffness
- 3 Prototype
- 4 Conclusion**

# Achievements

- Consistency with *prior*, confidentiality of *pot\_sec*, maximal availability of unmodified tuples with *db\_dist*
- Unique combination of model generation and distance minimization in infinite domain
- Quantifier handling without a need to expand them into ground conjunctions or disjunctions
- Only minor restrictions on the syntax of constraint formulas
- Both addition and deletion of tuples as modification primitives
- Optimized for complete databases with efficient query evaluation function
- Output in form of complete database instance
- Termination, soundness and completeness for appropriate fragments