



# Probabilistic Inference Modulo Theories

---

Rodrigo de Salvo Braz  
Ciaran O'Reilly  
Artificial Intelligence Center - SRI International

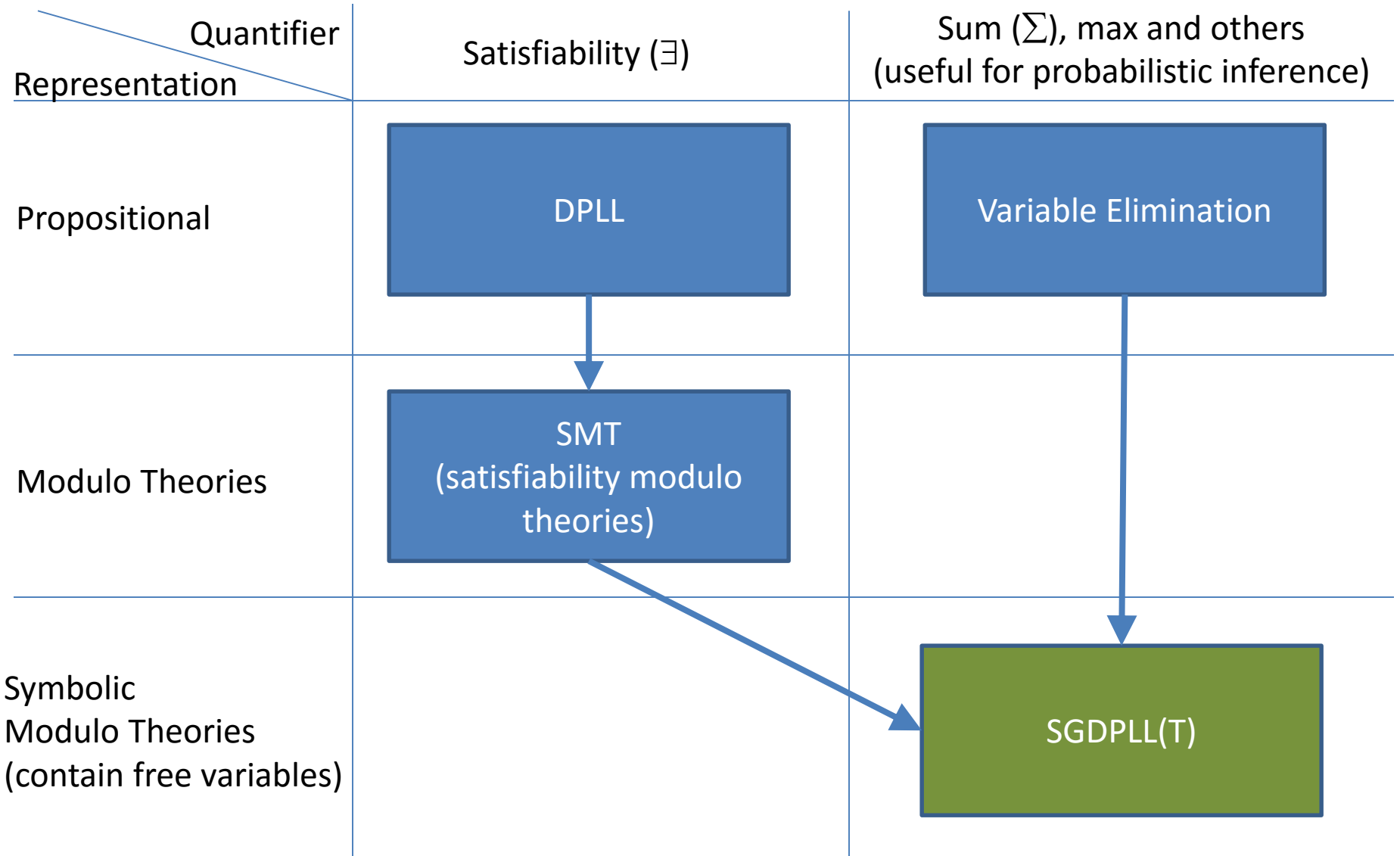
Vibhav Gogate  
University of Texas at Dallas

Rina Dechter  
University of California, Irvine

IJCAI-16, July 2016



# Overview



# Motivation

- Consider a probabilistic model on string-valued variables:

```
P(announcement | title, speaker, venue, abstract) =  
  if announcement = title + speaker + venue + abstract  
    then 0.7 else  
  if announcement = title + venue + speaker + abstract  
    then 0.3 else 0
```

```
P(speaker | name) =  
  if speaker = "Prof." + name  
    then 0.1 else  
  if speaker = name then 0.9 else 0
```

```
... // more statements, defining knowledge about  
    // names, titles etc.
```

- Exact graphical models algorithms typically iterate over values of each variable, but here they are infinite
- Sampling has its own set of disadvantages

# Probabilistic inference with Integers (polynomials and inequalities)

- Consider the following model:

$$P(z \in 1..1000) \propto z^2$$

$$P(x \in 1..1000 \mid z) \propto \text{if } x = z \text{ then } x \text{ else } 0.6$$

$$P(y \in 1..1000 \mid x) \propto \text{if } x > y \wedge y \neq 5 \text{ then } x^2 - y \text{ else } 0.9$$

$$P(x, y, z) \propto z^2 \times (\text{if } x > y \wedge y \neq 5 \text{ then } x^2 - y \text{ else } 0.9) \times (\text{if } x = z \text{ then } x \text{ else } 0.6)$$

$$\text{Marginal } P(y) = \sum_{z,x} P(x, y, z)$$

$$\propto \sum_z z^2 \sum_x (\text{if } x > y \wedge y \neq 5 \text{ then } x^2 - y \text{ else } 0.9) \times (\text{if } x = z \text{ then } x \text{ else } 0.6)$$

- How can we compute this sum without iterating over all the values?

# Example

# Background - Satisfiability

- We want to compute

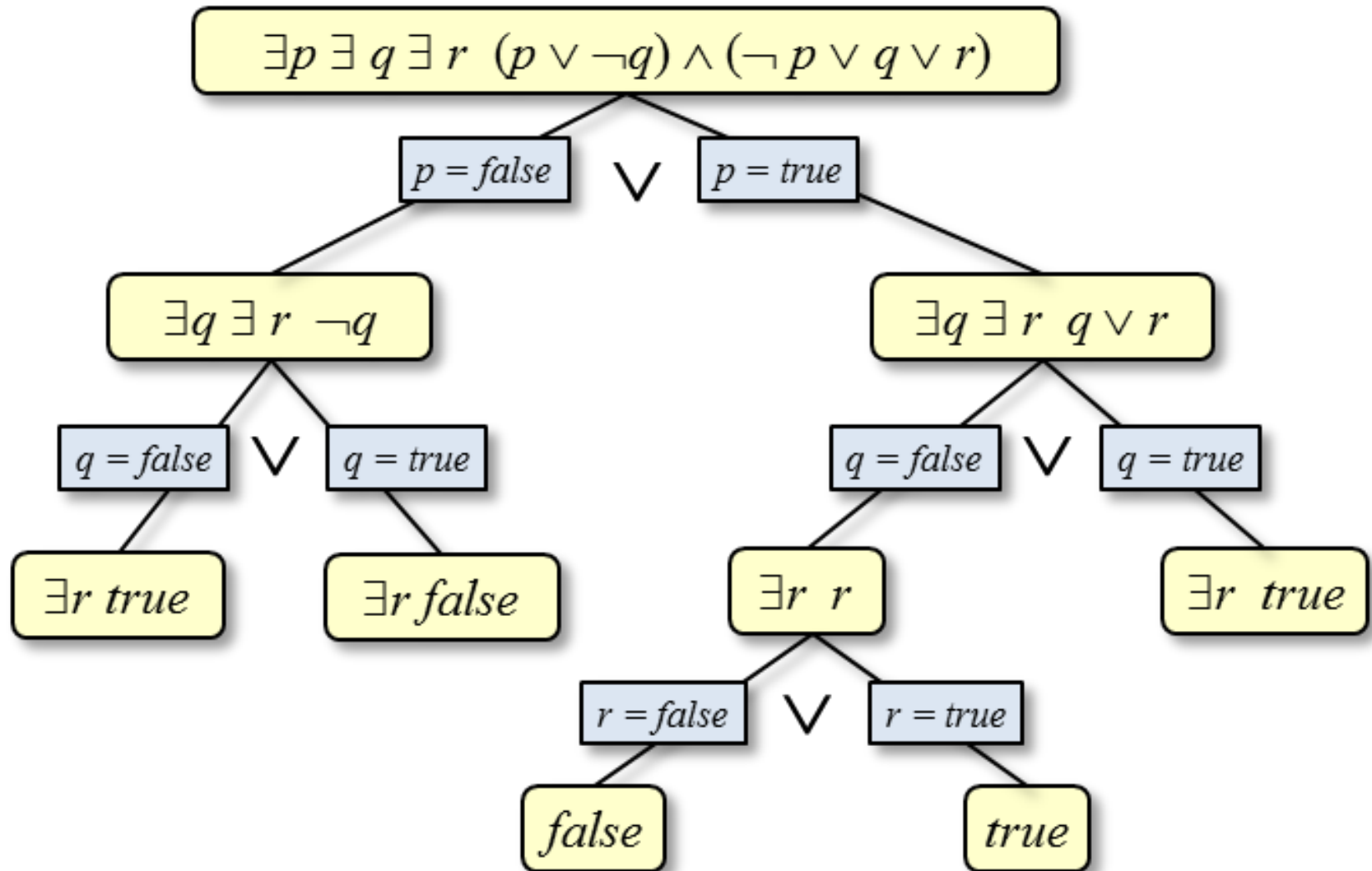
$$\sum_x (\text{if } x > y \wedge y \neq 5 \text{ then } x^2 - y \text{ else } 0.9) \times (\text{if } x = z \text{ then } x \text{ else } 0.6)$$

- The Davis-Putnam-Logemann-Loveland (DPLL) algorithm solves the problem of *satisfiability*:

$$\exists p \exists q \exists r (p \vee \neg q) \wedge (\neg p \vee q \vee r)$$

- This is similar to what we need, but for
  - Existential quantification instead of summation
  - Propositional variables (no theories)
  - Total quantification (no free variables)

# Background - DPLL



# Background – Satisfiability Modulo Theories (SMT)

- Satisfiability modulo theories generalizes satisfiability to non-propositional logic (includes arithmetic, inequalities, lists, uninterpreted functions, and others)

$\exists x \exists y \exists z \exists l (x \neq 5y \vee y > z) \wedge \text{cons}(x, \text{nil}) \neq \text{cdr}(l)$

- This is closer to what we need (since it works on theories), but for
  - Existential quantification instead of summation
  - Total quantification (no free variables)



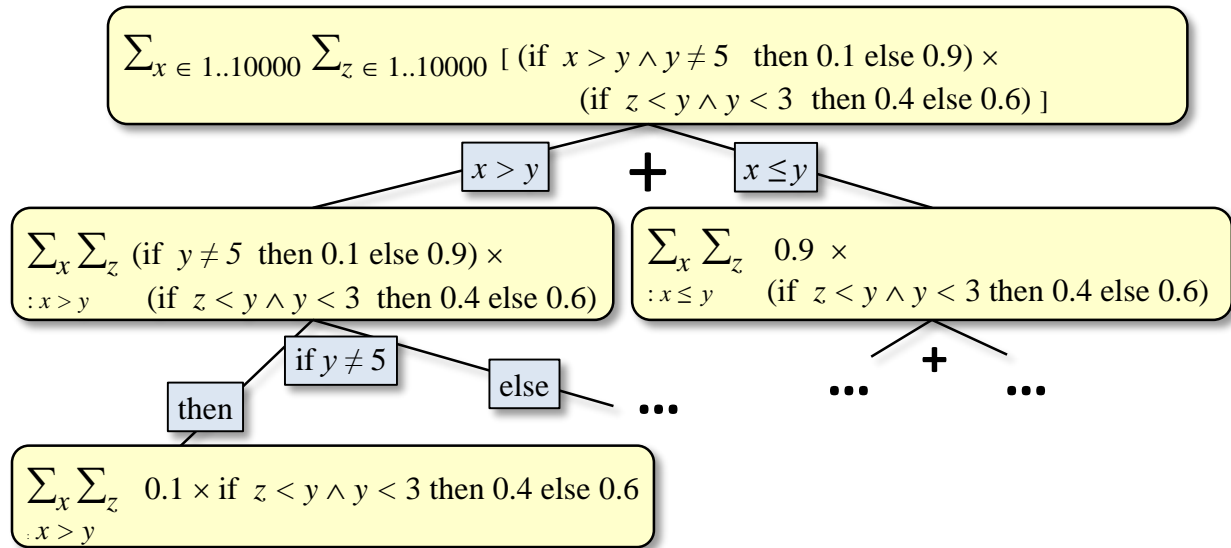
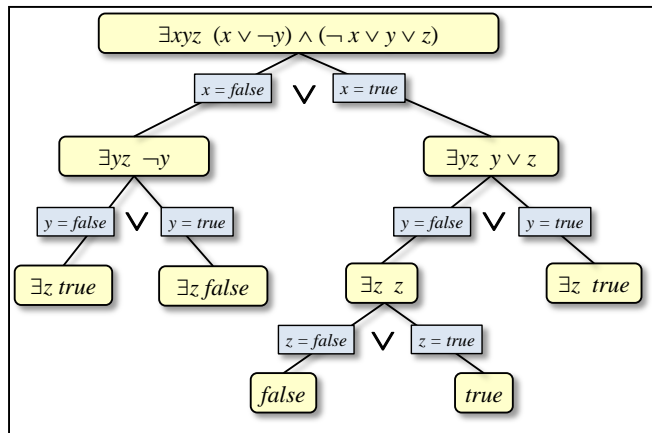
# First Contribution: Symbolic Generalized DPLL(T)

- Similar to SMT, but based on
  - Summation (or other quantifiers), besides  $\exists$
  - Partial quantification (free variables)

$$\sum_{x \in 1..10000} \sum_{z \in 1..10000} \quad (\text{if } x > y \wedge y \neq 5 \text{ then } 0.1 \text{ else } 0.9) \\ \times (\text{if } z < y \wedge y < 3 \text{ then } 0.4 \text{ else } 0.6)$$

- Note that  $y$  is a free variable
- Summed expression is not Boolean
- Language is not propositional ( $\neq$ ,  $<$ , ...)

# Symbolic Generalized DPLL(T) – SGDPLL(T)



Condition on literals until base case with no literals in main expression:

$$\sum_{x > y} \sum_{z < y} 0.04$$

$$\begin{aligned}
 &= \sum_{x: y < x \leq 100} \sum_{z: 1 \leq z < y} 0.04 \\
 &= \sum_{x: y < x \leq 100} (y - 1) 0.04 \\
 &= (100 - y) (y - 1) 0.04 \\
 &= -0.04y^2 + 4.04y - 4
 \end{aligned}$$

# Symbolic Generalized DPLL(T)

$$\sum_{x \in 1..10000} \sum_{z \in 1..10000} [ (\text{if } x > y \wedge y \neq 5 \text{ then } 0.1 \text{ else } 0.9) \times (\text{if } z < y \wedge y < 3 \text{ then } 0.4 \text{ else } 0.6) ]$$

$x > y$  +  $x \leq y$

$$\sum_{x > y} \sum_z (\text{if } y \neq 5 \text{ then } 0.1 \text{ else } 0.9) \times (\text{if } z < y \wedge y < 3 \text{ then } 0.4 \text{ else } 0.6)$$

$$\sum_{x \leq y} \sum_z 0.9 \times (\text{if } z < y \wedge y < 3 \text{ then } 0.4 \text{ else } 0.6)$$

if  $y \neq 5$  then ... else ...

$$\sum_{x > y} \sum_z 0.1 \times \text{if } z < y \wedge y < 3 \text{ then } 0.4 \text{ else } 0.6$$

$$\sum_{x > y} \sum_{z < y} 0.04$$

$$\begin{aligned} &= \sum_{x: y < x \leq 100} \sum_{z: 1 \leq z < y} 0.04 \\ &= \sum_{x: y < x \leq 100} (y - 1) 0.04 \\ &= (100 - y) (y - 1) 0.04 \\ &= -0.04y^2 + 4.04y - 4 \end{aligned}$$

Generic

Specific solver

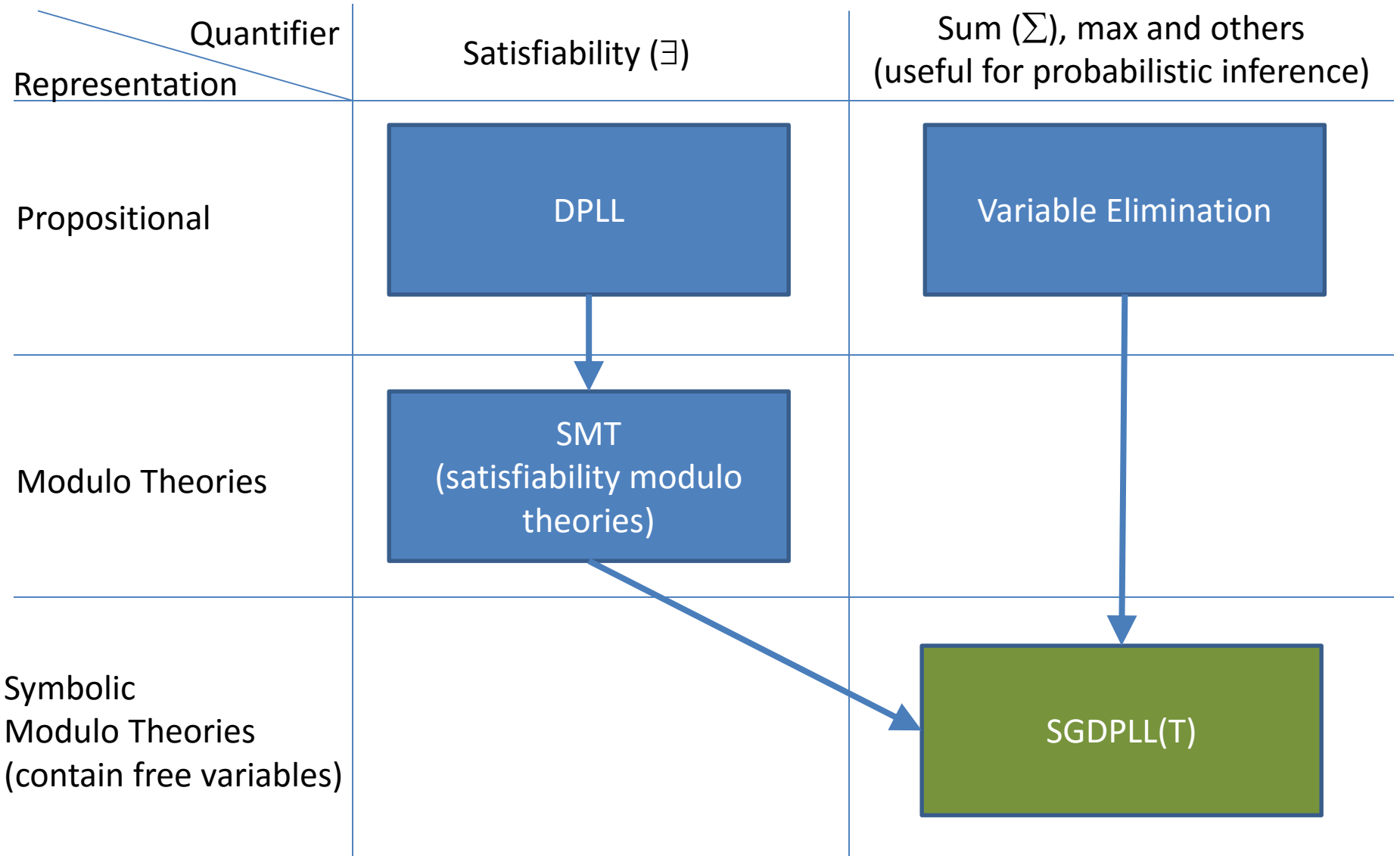
# Second Contribution: Solver for summation with difference arithmetic on bounded integers theory on polynomials

- $\sum_{z: 1 \leq z < y} 0.04$   
is an easy case:
  - Constant body expression
  - Single lower bound, single upper bound, no  $\neq$
- $\sum_{z: 1 \leq z \wedge x \leq z \wedge z \neq 5 \wedge z < y} z^2 - 2z$   
is more complicated:
  - Requires splitting on  $x < 1$  to decide which is lower bound
  - Requires splitting on  $5 < y$  to decide if  $z \neq 5$  is relevant
  - Requires a generalized *Faulhaber's formula* to sum over polynomial
- This splitting needs to be carefully implemented  
(simplifying at every split is too expensive)

# Since paper's final version...

- ... linear real arithmetic added as a separate theory
- Theories are automatically combined, so now we can define hybrid models on discrete and continuous variables and solve them symbolically

# Overview



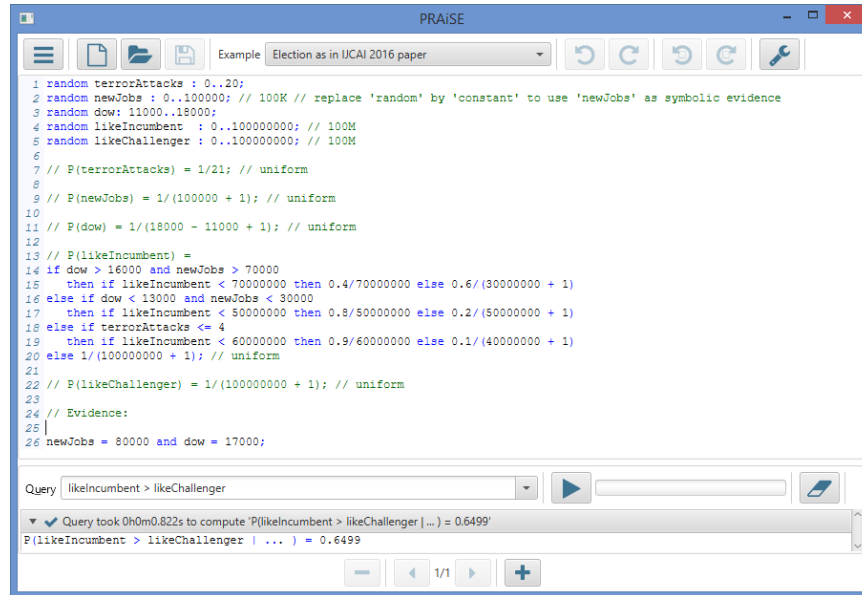
# Relation to Lifted Inference

- Lifted First-order Probabilistic Inference (Poole 2003, de Salvo Braz 2005)  
performs probabilistic inference on first-order predicates, without iterating over all values of their arguments

$$\begin{array}{rcl} \forall X \ P(\text{cancer}(X) \mid \text{smoker}(X)) & = & 0.6 \\ P(\text{smoker}(\text{mary})) & = & 0.01 \end{array}$$

- In the SMT vocabulary, that can be seen as a theory solver for uninterpreted functions
- This paper can be seen as lifted inference on interpreted functions
- Traditional lifted inference can be incorporated as a solver for uninterpreted functions in SGDPLL(T)

# Proof-of-concept Experiment



```
1 random terrorAttacks : 0..20;
2 random newJobs : 0..100000; // 100K // replace 'random' by 'constant' to use 'newJobs' as symbolic evidence
3 random dow : 11000..18000;
4 random likeIncumbent : 0..100000000; // 100M
5 random likeChallenger : 0..100000000; // 100M
6
7 // P(terrorAttacks) = 1/21; // uniform
8
9 // P(newJobs) = 1/(100000 + 1); // uniform
10
11 // P(dow) = 1/(18000 - 11000 + 1); // uniform
12
13 // P(likeIncumbent) =
14 if dow > 16000 and newJobs > 70000
15 then if likeIncumbent < 700000000 then 0.4/700000000 else 0.6/(300000000 + 1)
16 else if dow < 13000 and newJobs < 30000
17 then if likeIncumbent < 500000000 then 0.8/500000000 else 0.2/(500000000 + 1)
18 else if terrorAttacks <= 4
19 then if likeIncumbent < 600000000 then 0.9/600000000 else 0.1/(400000000 + 1)
20 else 1/(1000000000 + 1); // uniform
21
22 // P(likeChallenger) = 1/(1000000000 + 1); // uniform
23
24 // Evidence:
25 |
26 newJobs = 80000 and dow = 17000;
```

Query: likeIncumbent > likeChallenger

Query took 0h0m0.822s to compute 'P(likeIncumbent > likeChallenger | ...) = 0.6499'

P(likeIncumbent > likeChallenger | ...) = 0.6499

- Grounded the elections example into a regular graphical models and using VEC (Gogate & Dechter, 2011)
- Had to decrease domain size to 180 to keep it manageable, and VEC was still 20 times slower



# Conclusion

- This is graphical models, but defined with richer representations (theories)
- Similar to SMT, but with summation and free variables
- Symbolic: result is a math expression on free variables
- Future work
  - solvers for more theories
  - unit propagation, clause learning from SAT literature
  - bounded approximations for limiting the search

Thank you!