

# Towards a Theory of Learning Inductive Invariants

Yotam Feldman



Neil Immerman



Mooly Sagiv



Sharon Shoham



James R. Wilcox



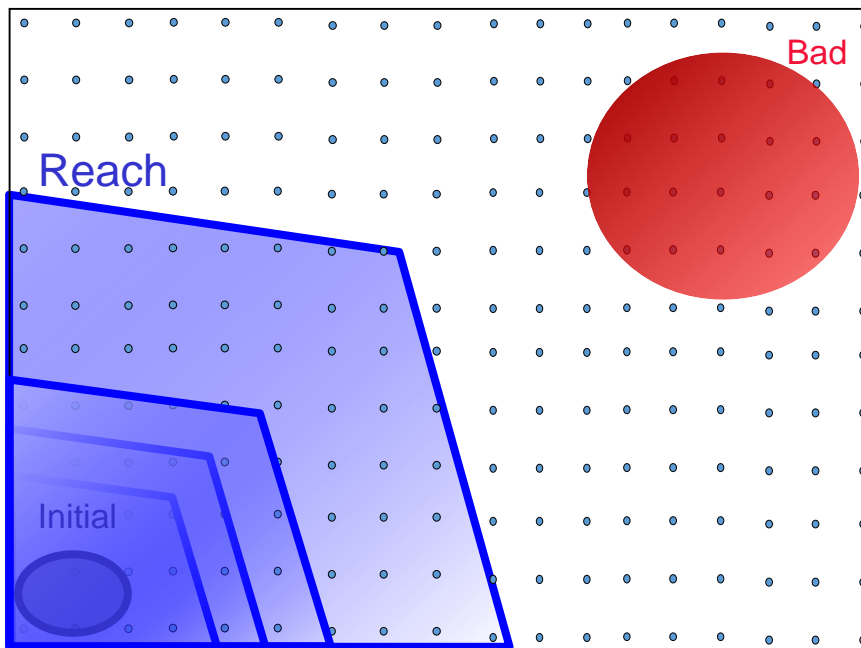
 [@yotamfe](https://twitter.com/yotamfe), [@SagivMooly](https://twitter.com/SagivMooly), [@wilcoxjay](https://twitter.com/wilcoxjay)

# Safety of Transition Systems

Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

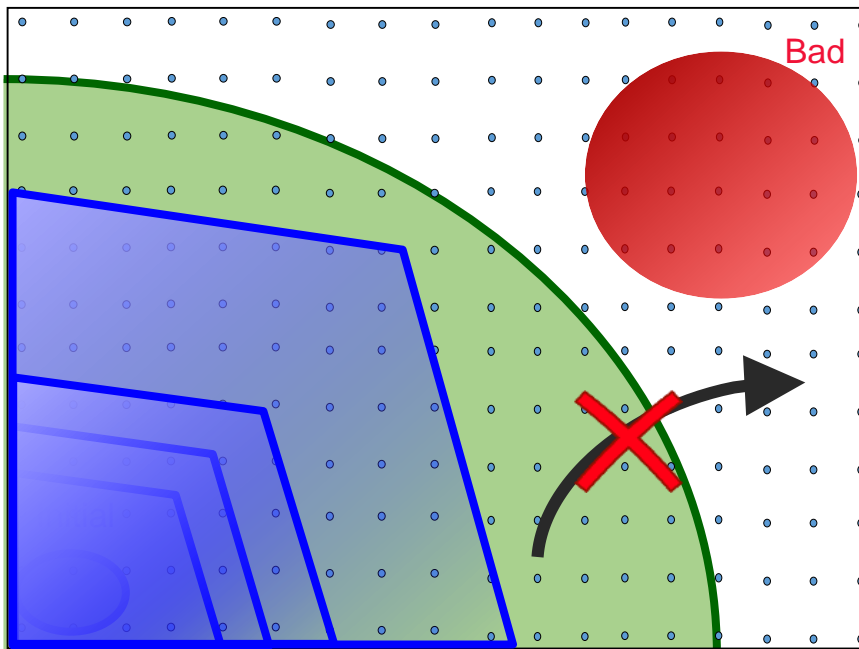


# Inductive Invariants

Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

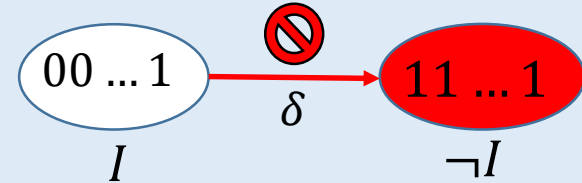
Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$



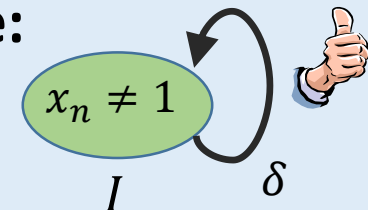
$I: (x_1, \dots, x_n) \neq 1 \dots 1$

**Not inductive:**



$I: x_n \neq 1$

**Inductive:**



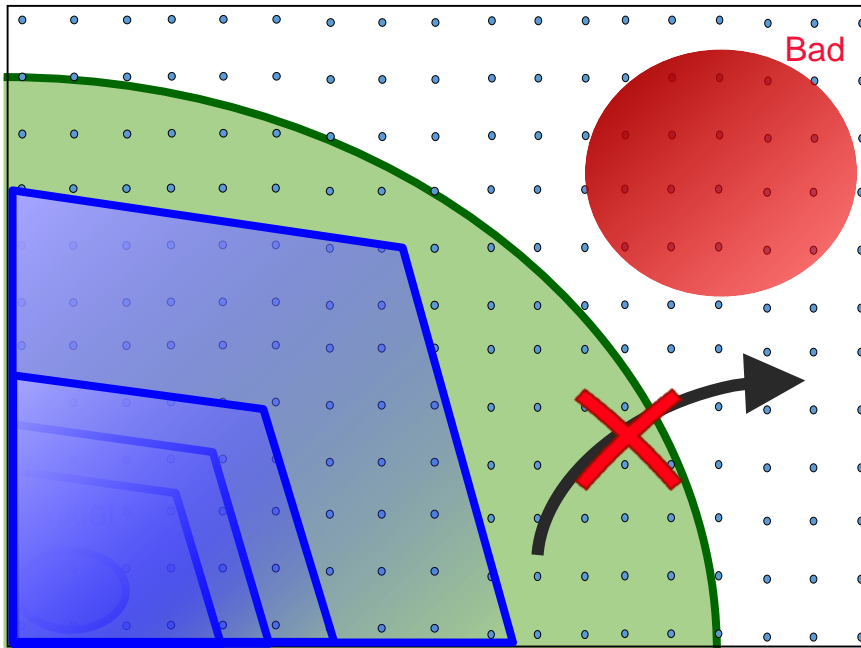
# Inductive Invariants

Init:

$\delta$ :

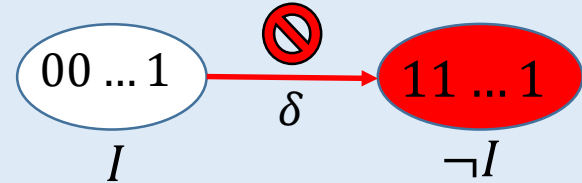
Goal:  
Find inductive invariants automatically

$(x_1, \dots, x_n) = 1 \dots 1$



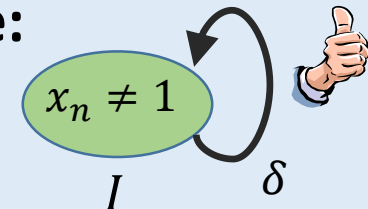
$I: (x_1, \dots, x_n) \neq 1 \dots 1$

**Not inductive:**



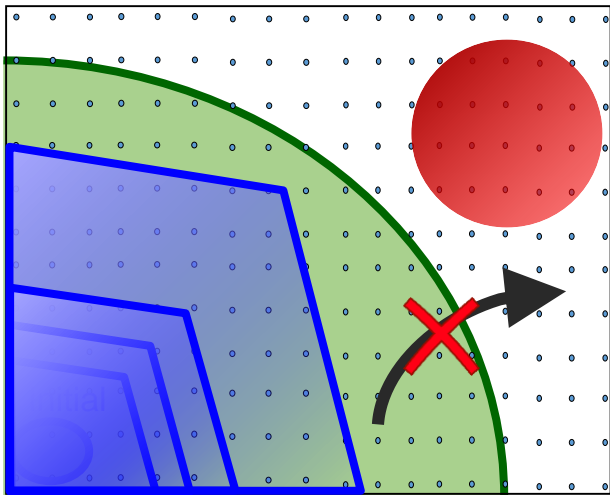
$I: x_n \neq 1$

**Inductive:**

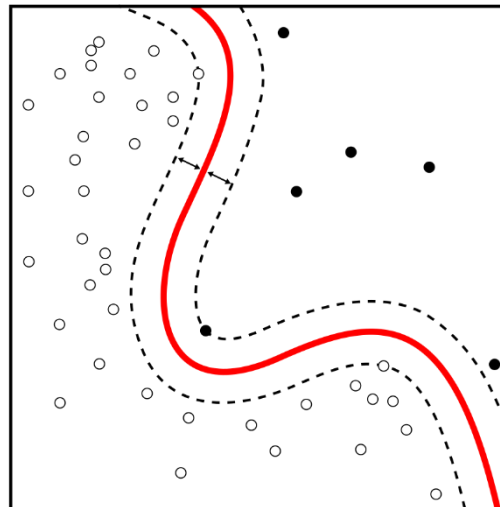


# This Work

Invariant Inference



Exact Concept Learning



vs.

- Query-based learning models for invariant inference
- Invariant inference is harder than concept learning
- Complexity results for invariant inference algorithms from classification algorithms

# Problem Setting: Polynomial-Length Inference

Boolean transition systems,  $\Sigma = \{p_1, \dots, p_n\}$

Given a transition system from a class  $\mathcal{P}$  (over  $\Sigma$ ),

Find an inductive invariant

$I \in \text{DNF}$

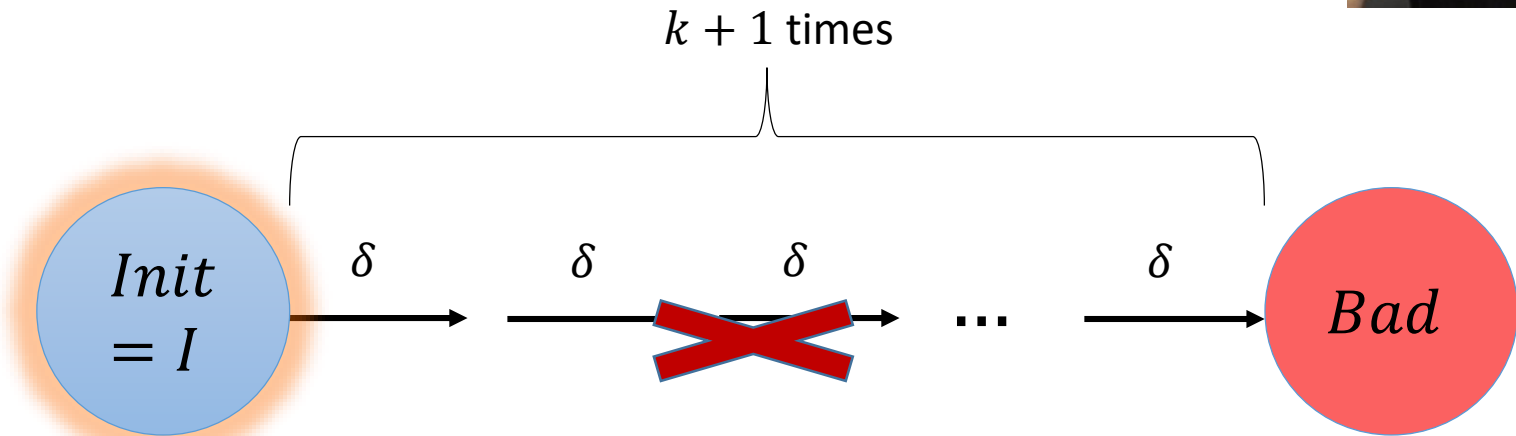
$|I| \leq \text{poly}(n)$

(Decision problem is  $\Sigma_2^P$ -complete.)

# Interpolation-Based Inference

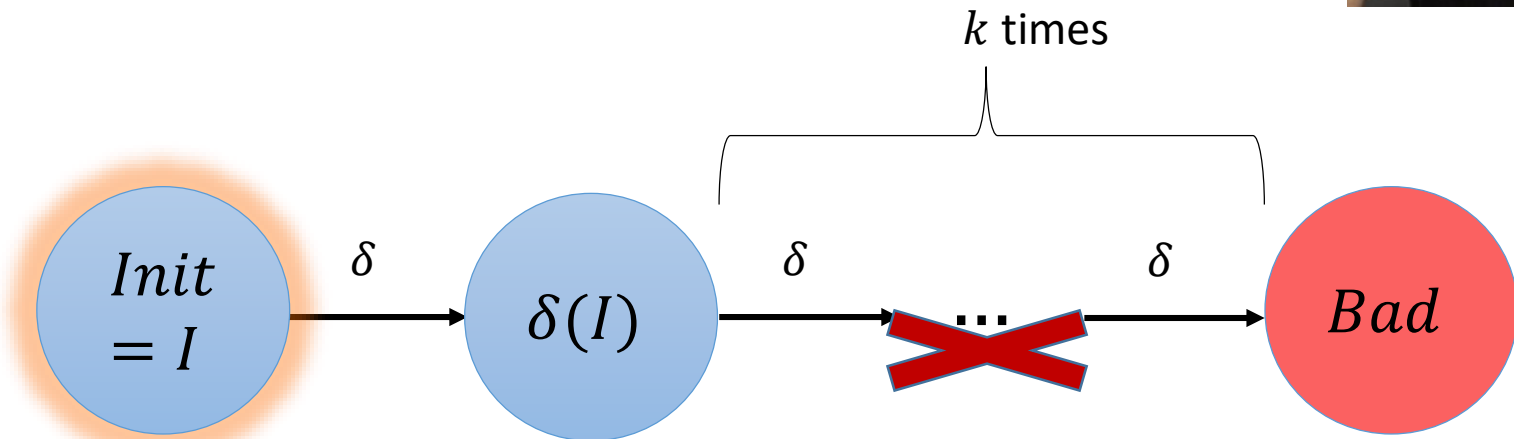


$I = Init$



# Interpolation-Based Inference

$I = \text{Init}$

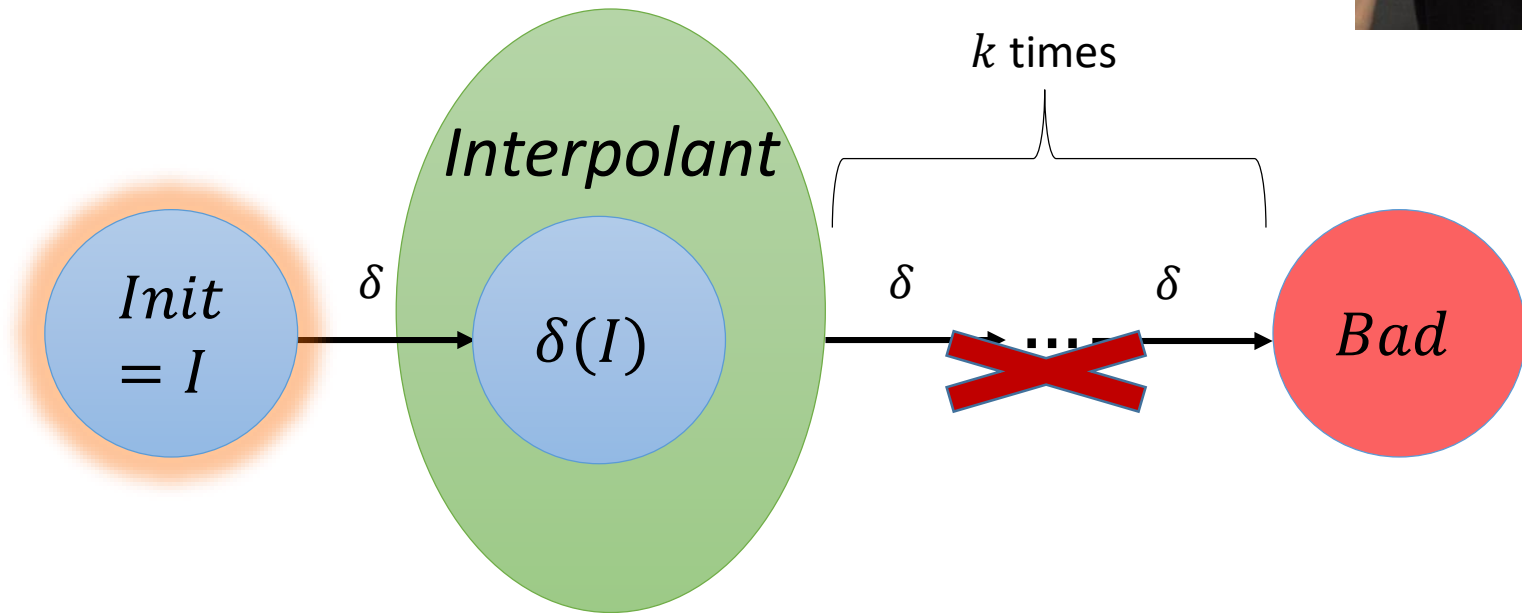




# Interpolation-Based Inference



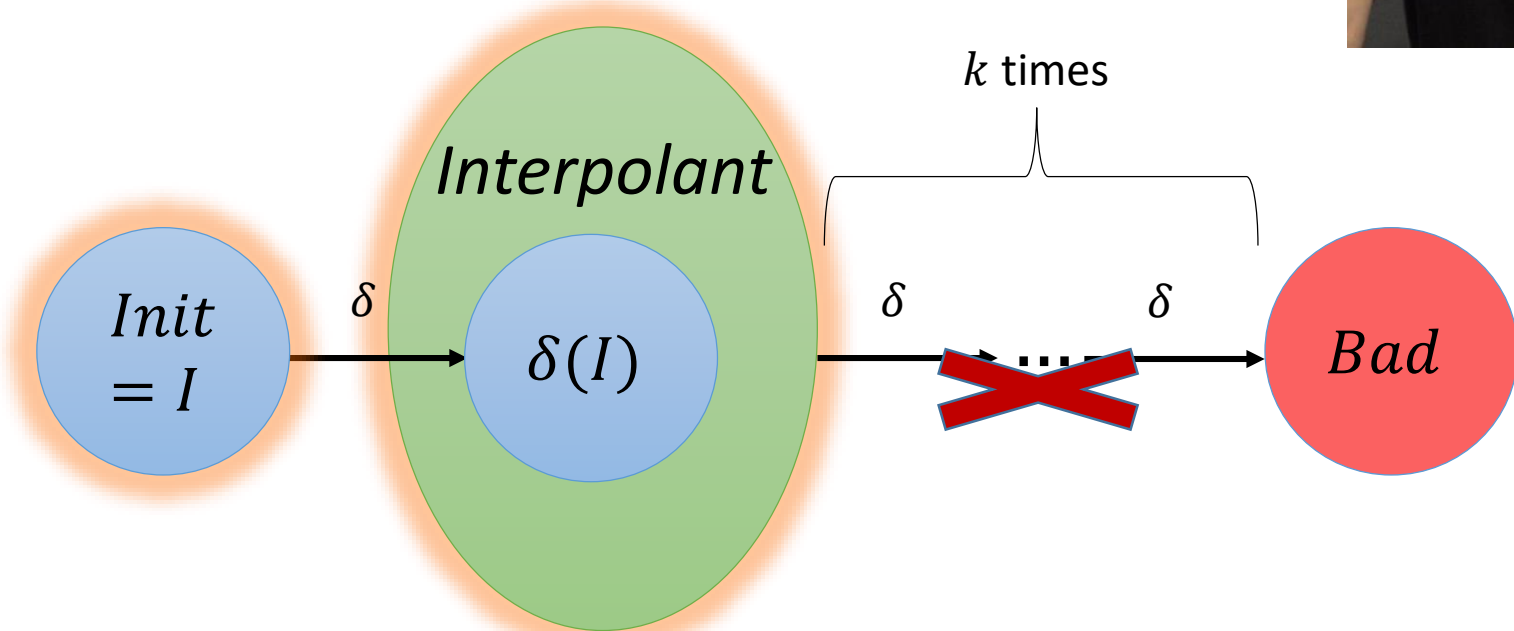
$I = \text{Init}$



# Interpolation-Based Inference



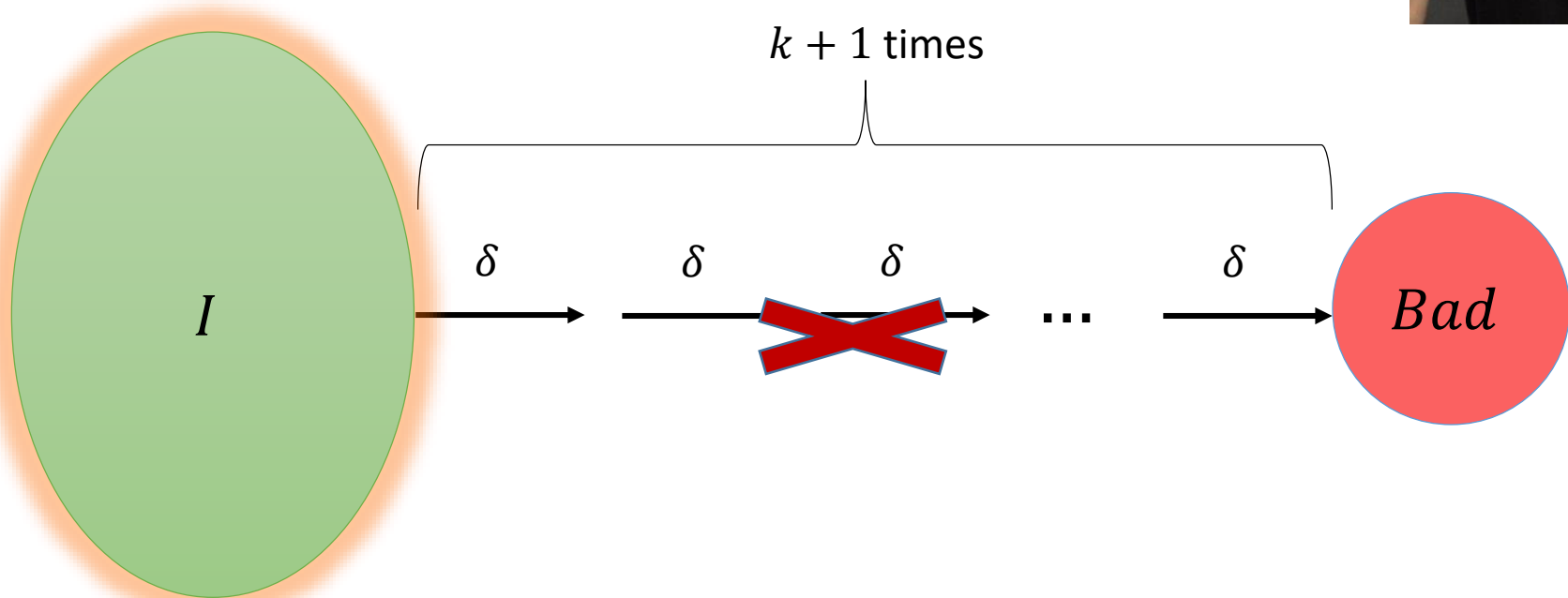
$$I = \text{Init} \vee \text{Interpolant}$$



# Interpolation-Based Inference



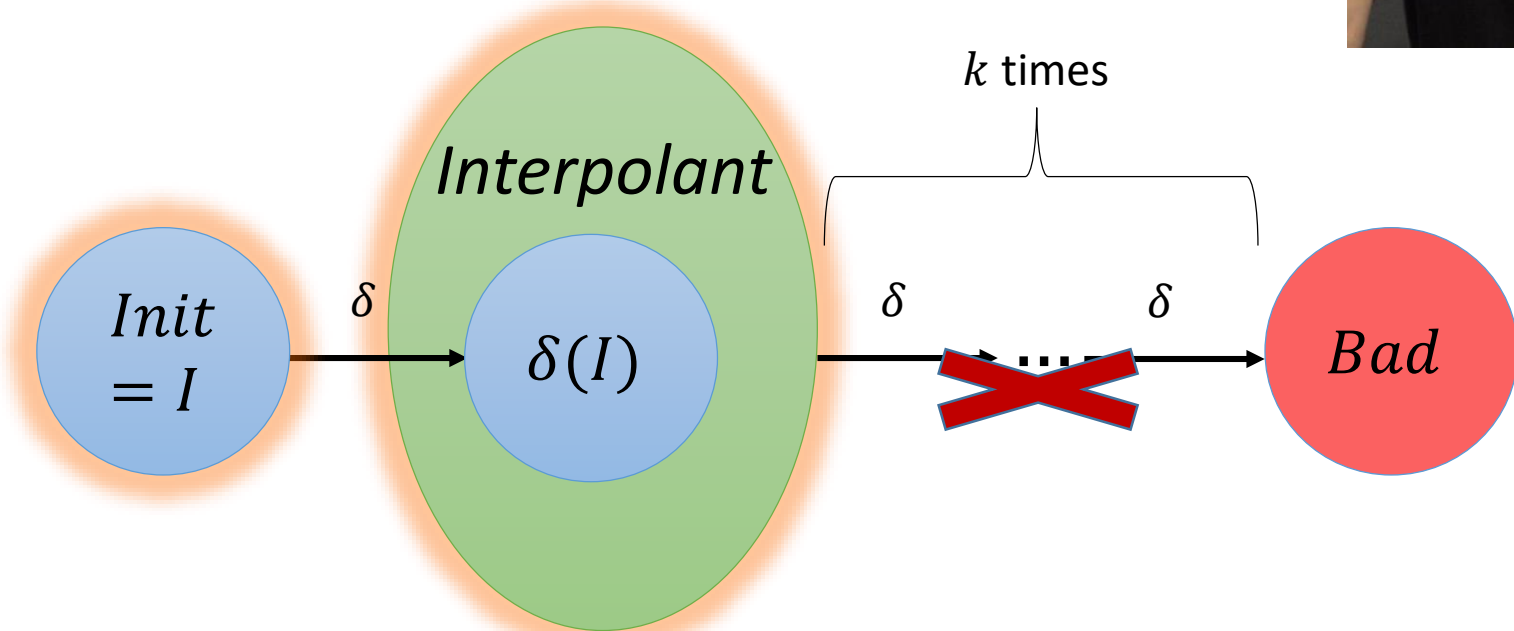
$I = \text{Init} \vee \text{Interpolant}$



# Interpolation-Based Inference



$$I = \text{Init} \vee \text{Interpolant}$$



# Model-Based Interpolation

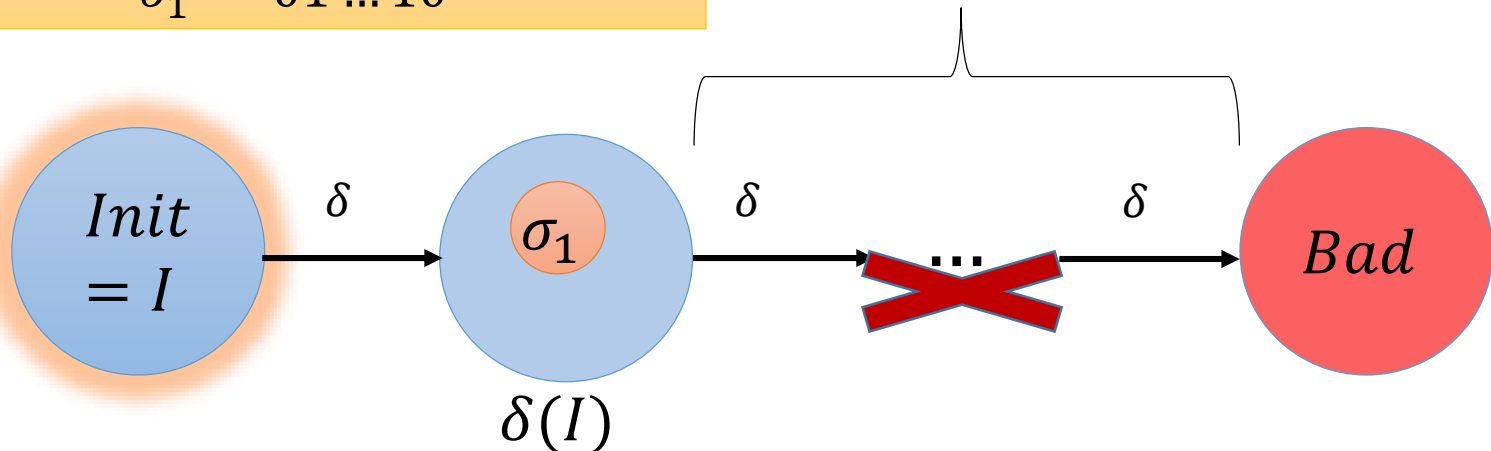
Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$

$\sigma_1 = 01 \dots 10$



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

# Model-Based Interpolation

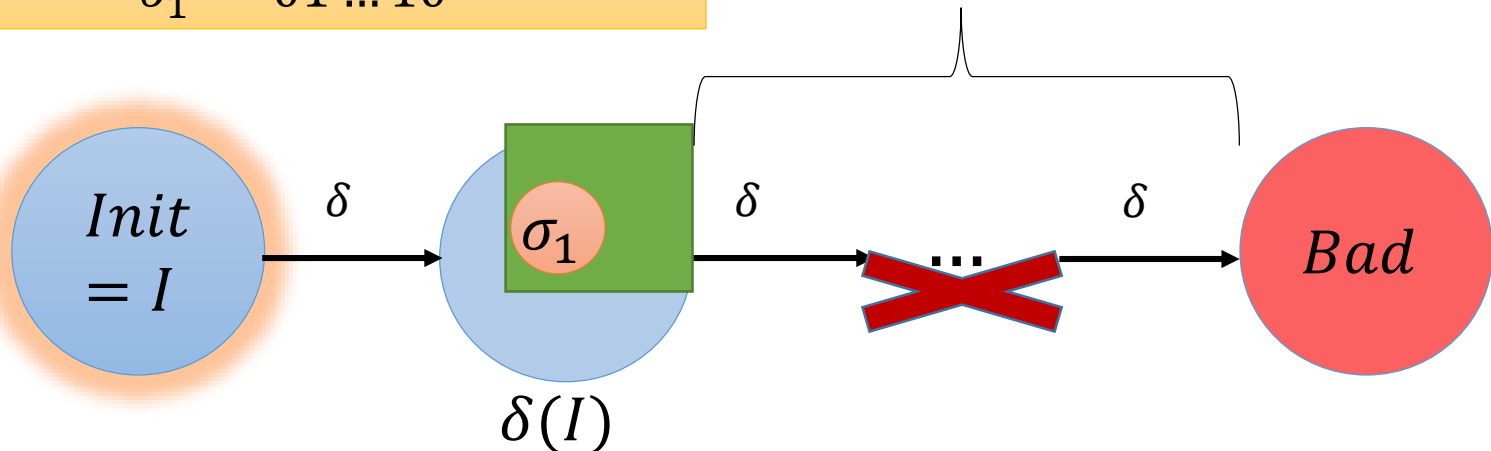
Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$

$\sigma_1 = 01 \dots 10$



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

# Model-Based Interpolation

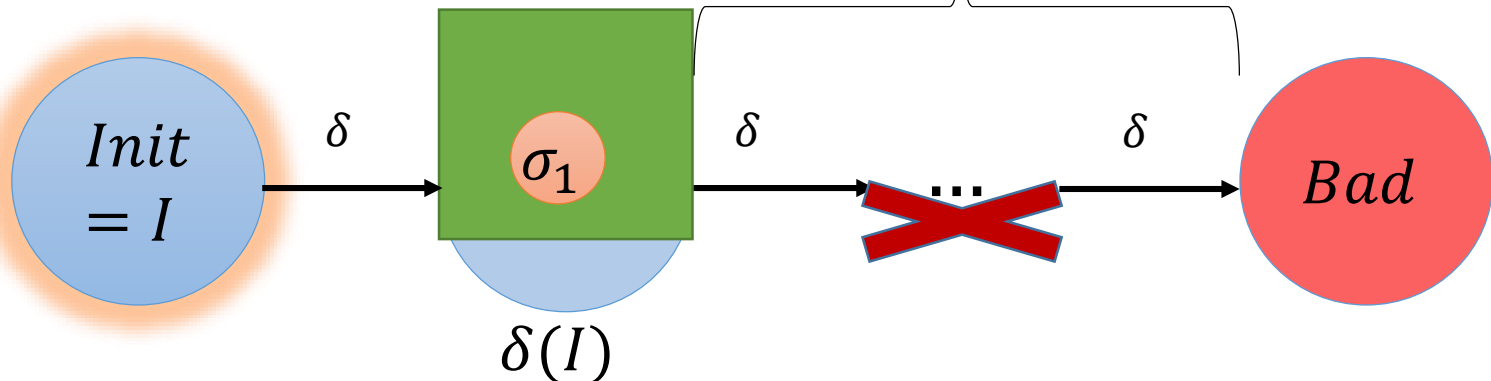
Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$

$\sigma_1 = 01 \dots 10$



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

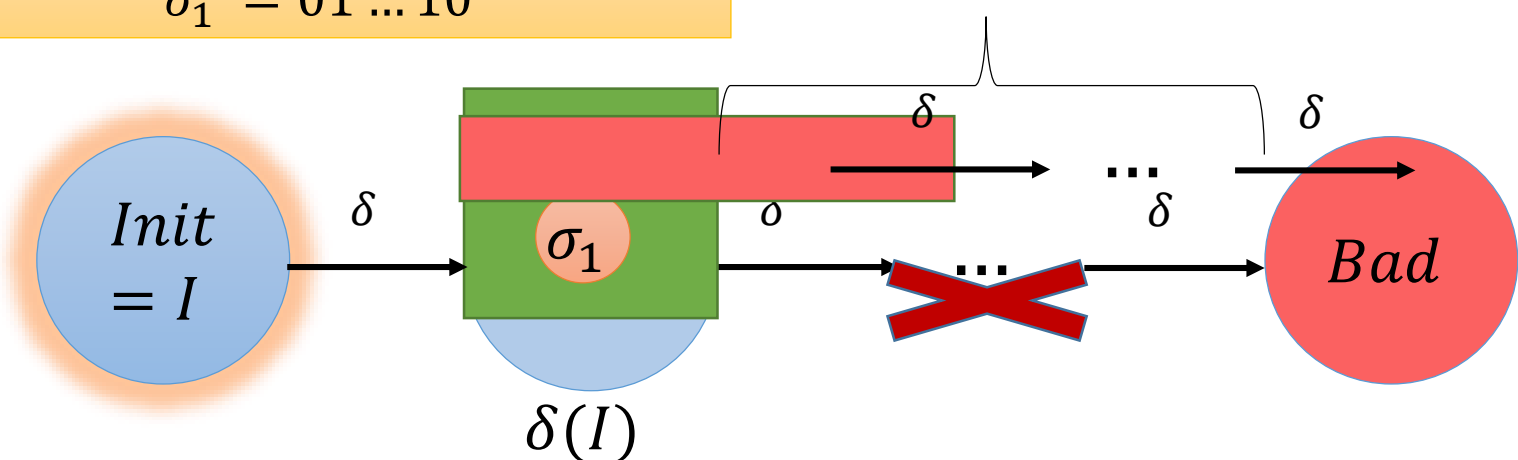
# Model-Based Interpolation

Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$   
 $\sigma_1 = 01 \dots 10$



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav



# Model-Based Interpolation

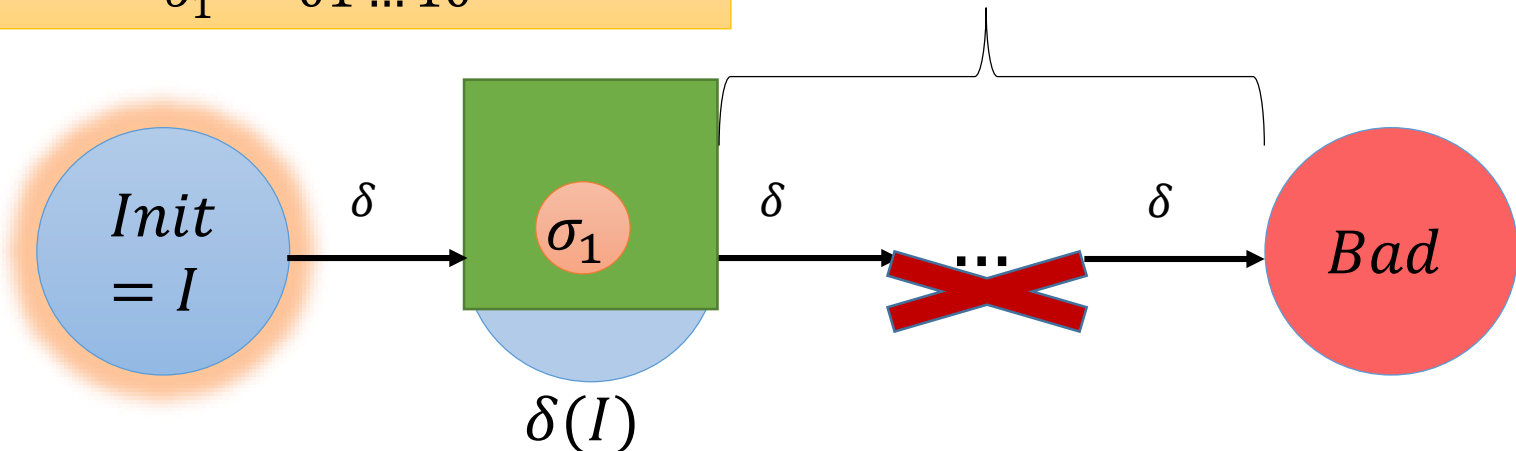
Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$

$\sigma_1 = 01 \dots 10$



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

# Model-Based Interpolation

Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

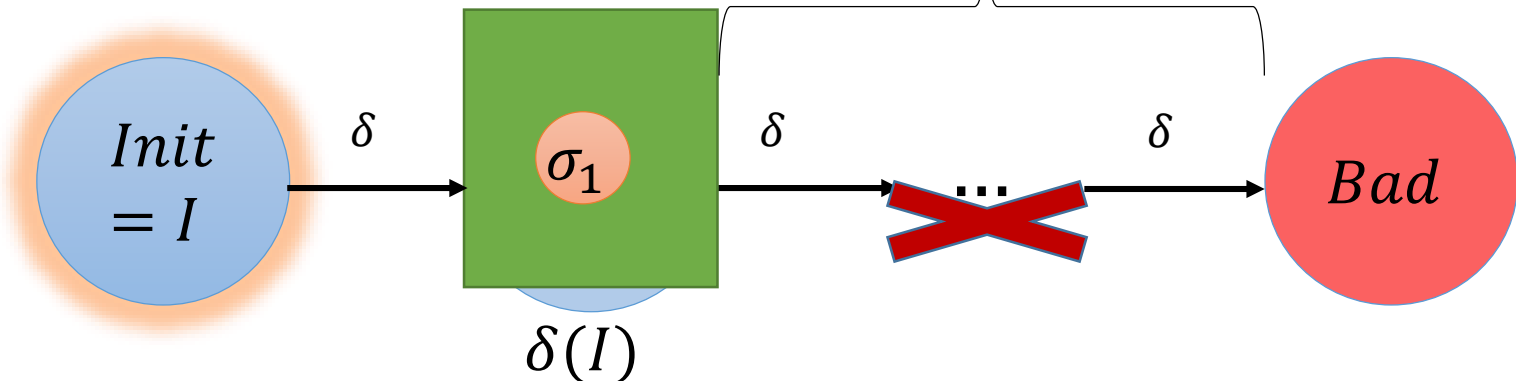
Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$Interpolant_1 = (x_1 = 0 \wedge x_2 = 1 \wedge \dots \wedge x_n = 1 \wedge x_n = 0)$

$\sigma_1 = 01 \dots 10$

$k$  times



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

# Model-Based Interpolation

Init:  
 $(x_1, \dots, x_n) := 0 \dots 0$

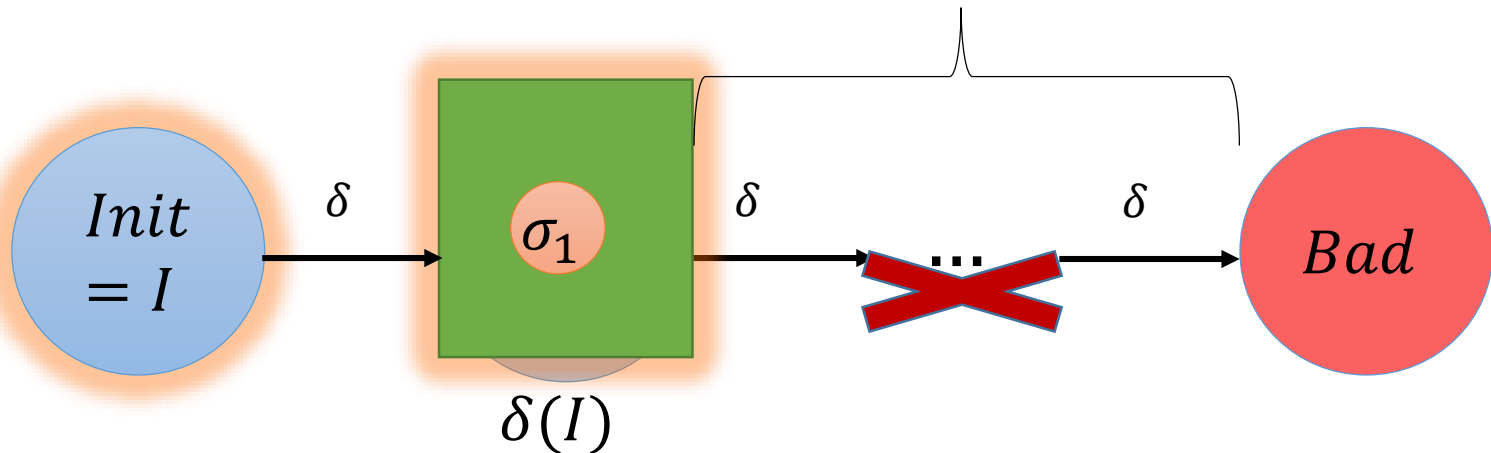
Bad:  
 $(x_1, \dots, x_n) = 1 \dots 1$

$\delta$ :  
 $y_1, \dots, y_n := *$   
 $x_1, \dots, x_n := (x_1, \dots, x_n) +$   
 $2 \cdot (y_1, \dots, y_n) \pmod{2^n}$

$I = \text{Init} \vee (x_n = 0)$



$k$  times



[HVC'12] Computing Interpolants without Proofs. Chockler, Ivrii, Matsliah

[LPAR'13] Instantiations, Zippers and EPR Interpolation. Bjørner, Gurfinkel, Korovin, Lahav

# Model-Based Interpolation

Inferring invariant in DNF:

$$\underbrace{(\ell_1^1 \wedge \cdots \wedge \ell_{k_1}^1)}_{\text{gen}(\sigma_1)} \vee \dots \vee \underbrace{(\ell_1^m \wedge \cdots \wedge \ell_{k_m}^m)}_{\text{gen}(\sigma_m)}$$

$I := \text{false}$

while  $(\_, \sigma')$  counterexample  
to **Inductive** $(\delta, I)$ :

$I := I \vee \text{generalize}(\sigma')$

**generalize** $(\sigma')$ :

drop literals from  $\sigma'$

while **BMC** $(\delta, \sigma', k) \cap \text{Bad} = \emptyset$

# Model-Based Interpolation

```
 $I := \text{false}$   
while  $(\_, \sigma')$  counterexample  
to Inductive $(\delta, I)$ :  
   $I := I \vee \text{generalize}(\sigma')$   
  
generalize $(\sigma')$ :  
  drop literals from  $\sigma'$   
  while BMC $(\delta, \sigma', k) \cap \text{Bad} = \emptyset$ 
```

# Understanding Invariant Inference

```
 $I := \text{false}$   
while  $(\_, \sigma')$  counterexample  
to Inductive $(\delta, I)$ :  
   $I := I \vee \text{generalize}(\sigma')$   
  
generalize $(\sigma')$ :  
  drop literals from  $\sigma'$   
  while BMC $(\delta, \sigma', k) \cap \text{Bad} = \emptyset$ 
```



# Understanding Invariant Inference

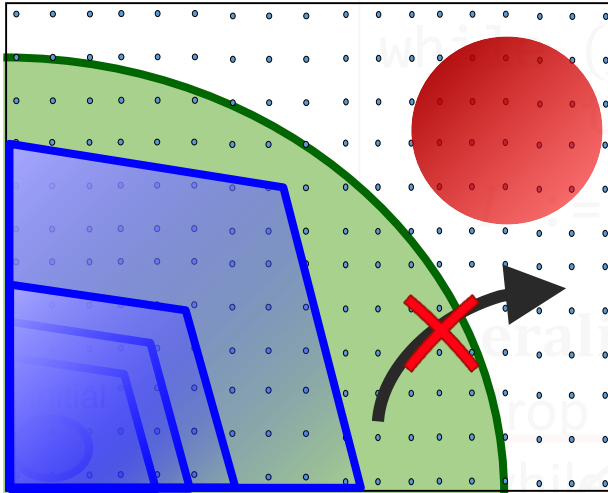
```
 $I := \text{false}$   
while  $(\_, \sigma')$  counterexample  
to Inductive $(\delta, I)$ :  
   $I := I \vee \text{generalize}(\sigma')$   
  
generalize $(\sigma')$ :  
  drop literals from  $\sigma'$   
  while BMC $(\delta, \sigma', k) \cap \text{Bad} = \emptyset$ 
```

Complexity bounds from  
exact classification algorithms

Rich SAT queries allow  
exponentially faster inference

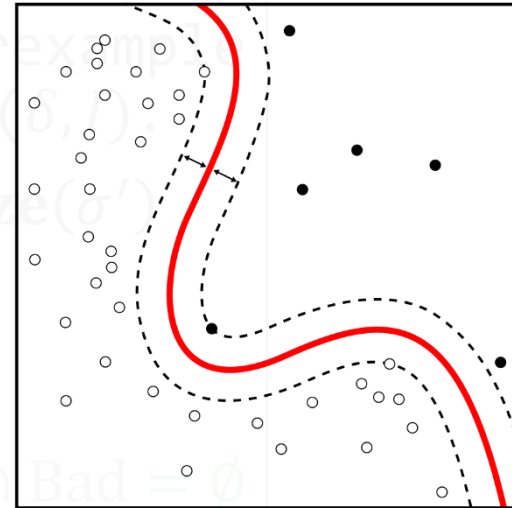
# Understanding Invariant Inference

Invariant Inference



vs.

Exact Concept Learning



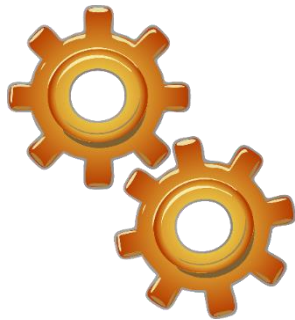
Complexity bounds from exact classification algorithms

Rich SAT queries allow exponentially faster inference



# Exact Concept Learning with Equivalence & Membership Queries

learning algorithm



is it  $\psi_1$ ?

✓ / ✗+counterexample

is it  $\psi_2$ ?

✓ / ✗+counterexample

does  $\sigma_3 \models$ ?

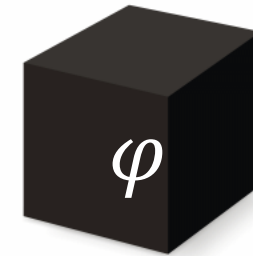
✓ / ✗

...

Membership

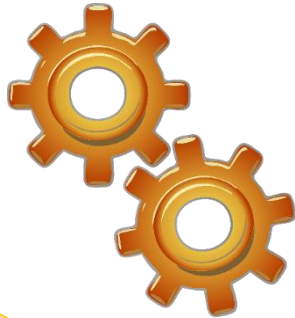
Equivalence

oracle



# Invariant Inference with Equivalence & Membership Queries

learning algorithm



Need to implement this

is it  $\psi_1$ ?

✓ / ✗+counterexample

is it  $\psi_2$ ?

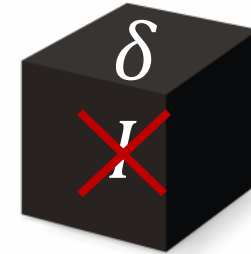
✓ / ✗+counterexample

does  $\sigma_3 \models$ ?

✓ / ✗

...

oracle



Need to implement this

# Inductiveness-Query Model

inference algorithm



$\alpha_1$  inductive?

✓ / ✗+counterexample

inductiveness-query oracle

...  
 $\alpha_m$  inductive?

✓ / ✗+counterexample

$\{\alpha_i\}$   $\delta$   $\{\alpha_i\}$

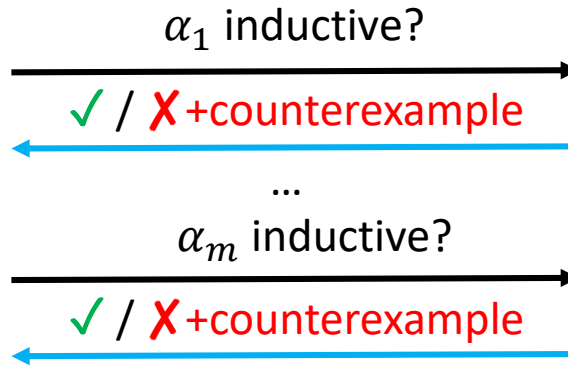
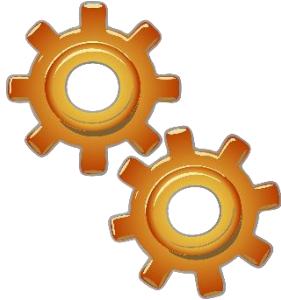
Transition  $(\sigma, \sigma')$  of  $\delta$  s.t.  
 $\sigma \models \alpha_i, \sigma' \models \neg \alpha_i$

Algorithms cannot access the transition relation directly,  
only perform inductiveness queries

Complexity: # inductiveness queries  
worst case amongst possible counterexamples

# Inductiveness-Query Model

inference algorithm



inductiveness-query oracle



```

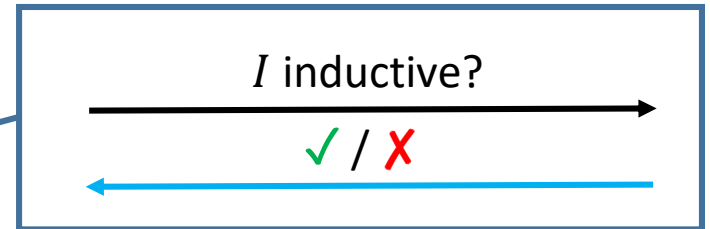
I := false
while  $(\_, \sigma')$  counterexample
to  $\text{Inductive}(\delta, I)$ :
    I := I  $\vee$  generalize( $\sigma'$ )

```

**generalize**( $\sigma'$ ):

drop literals from  $\sigma'$

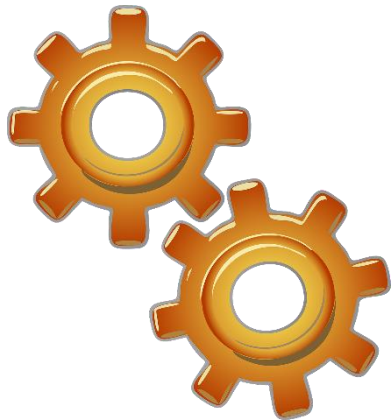
while  $\text{BMC}(\delta, \sigma', k) \cap \text{Bad} = \emptyset$



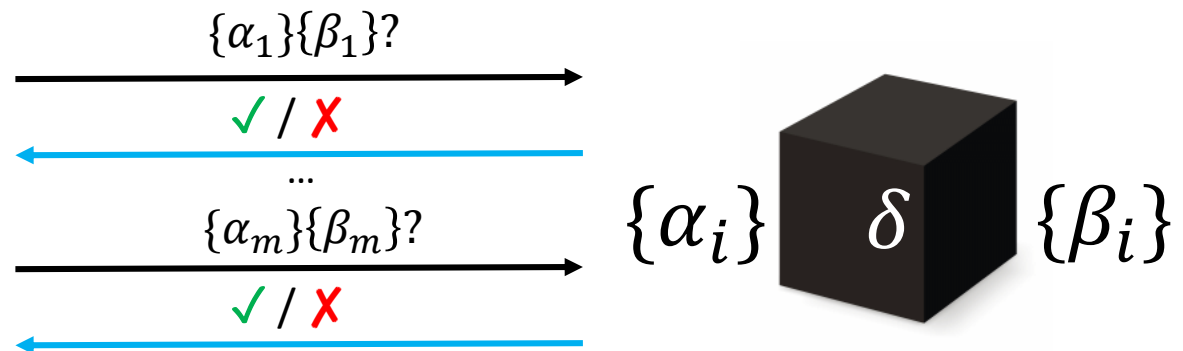
# Hoare-Query Model

Capable of modeling several interesting algorithms

inference algorithm



Hoare-query oracle



Algorithms cannot access the transition relation directly,  
only perform Hoare queries

# Hoare-Query Model

Capable of modeling several interesting algorithms

model of  $\neg\{I\}\delta\{I\}$ ,  
 $O(n)$  queries

$I := \text{false}$

while  $(\_, \sigma')$  counterexample  
to **Inductive** $(\delta, I)$ :

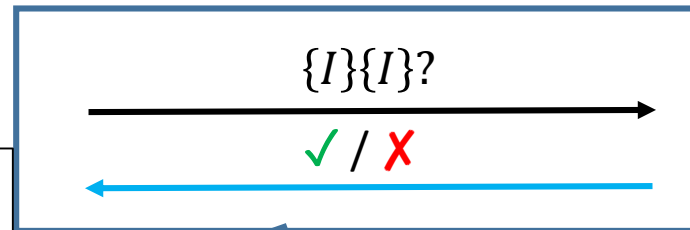
$I := I \vee \text{generalize}(\sigma')$

**generalize** $(\sigma')$ :

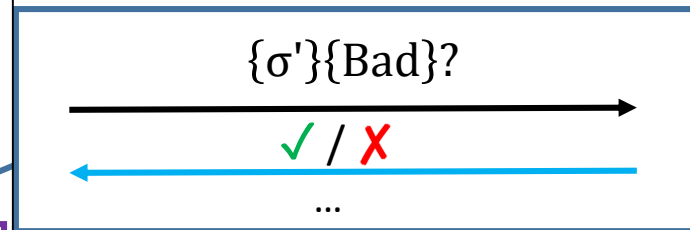
drop literals from  $\sigma'$

while **BMC** $(\delta, \sigma', 1) \cap \text{Bad} = \emptyset$

Hoare-query  
oracle



$\{\alpha_i\}$   $\delta$   $\{\beta_i\}$



# Hoare > Inductiveness

Thm: There exists a class of transition systems  $\mathcal{P}$ , so that for solving polynomial-length inference:

1.  $\exists$  Hoare-query algorithm with  $\text{poly}(n)$  queries
2.  $\forall$  inductiveness-query algorithm requires  $2^{\Omega(n)}$  queries



a simple case of IC3/PDR

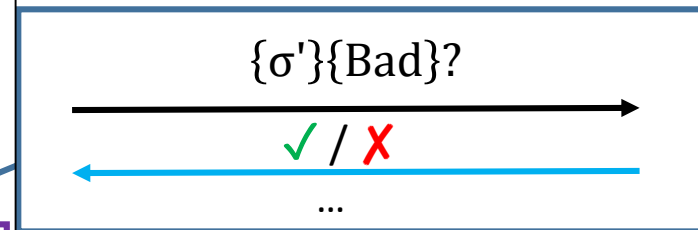
$\Rightarrow$  ICE cannot model PDR,  
and the extension of [VMCAI'17] is necessary

# Hoare > Inductiveness

Thm: There exists a class of transition systems  $\mathcal{P}$ , so that for solving polynomial-length inference:

1.  $\exists$  Hoare-query algorithm with  $\text{poly}(n)$  queries
2.  $\forall$  inductiveness-query algorithm requires  $2^{\Omega(n)}$  queries

```
 $I := \text{false}$   
while  $(\_, \sigma')$  counterexample  
to  $\text{Inductive}(\delta, I)$ :  
   $I := I \vee \text{generalize}(\sigma')$   
  
generalize $(\sigma')$ :  
drop literals from  $\sigma'$   
while  $\text{BMC}(\delta, \sigma', 1) \cap \text{Bad} = \emptyset$ 
```





# Hoare > Inductiveness

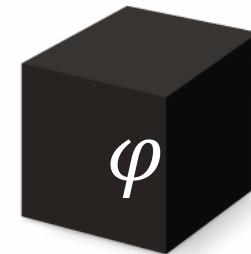
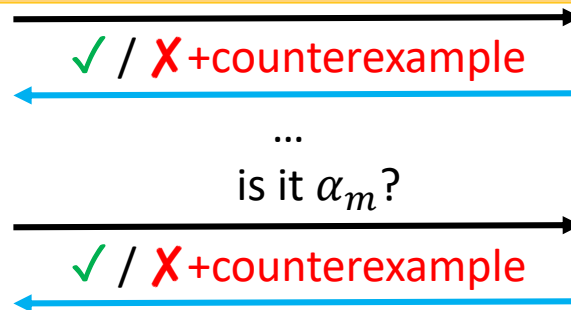
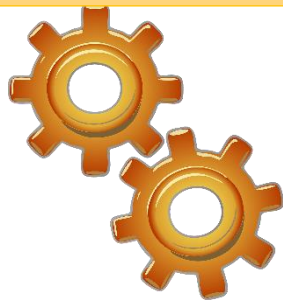
Thm: There exists a class of transition systems  $\mathcal{P}$ , so that for solving polynomial-length inference:

1.  $\exists$  Hoare-query algorithm with  $\text{poly}(n)$  queries
2.  $\forall$  inductiveness-query algorithm requires  $2^{\Omega(n)}$  queries



# Learning from Counterexamples to Equivalence Queries

Thm: Learning from counterexamples to induction is **harder** than learning from labeled examples.



Positive/negative examples:

$$\sigma^+ \models \varphi, \sigma^- \models \neg\varphi$$

Counterexamples to induction:

$$\sigma \models \neg\varphi \text{ or } \sigma' \models \varphi$$

Learning monotone DNF:

subexponential

this work:  $2^{\Omega(n)}$

[ML'87] Queries and Concept Learning, Angluin

[COLT'12] Tight Bounds on Proper Equivalence Query Learning of DNF, Hellerstein et al.

[POPL'20] Complexity and Information in Invariant Inference. Feldman, Immerman, Shoham, Sagiv

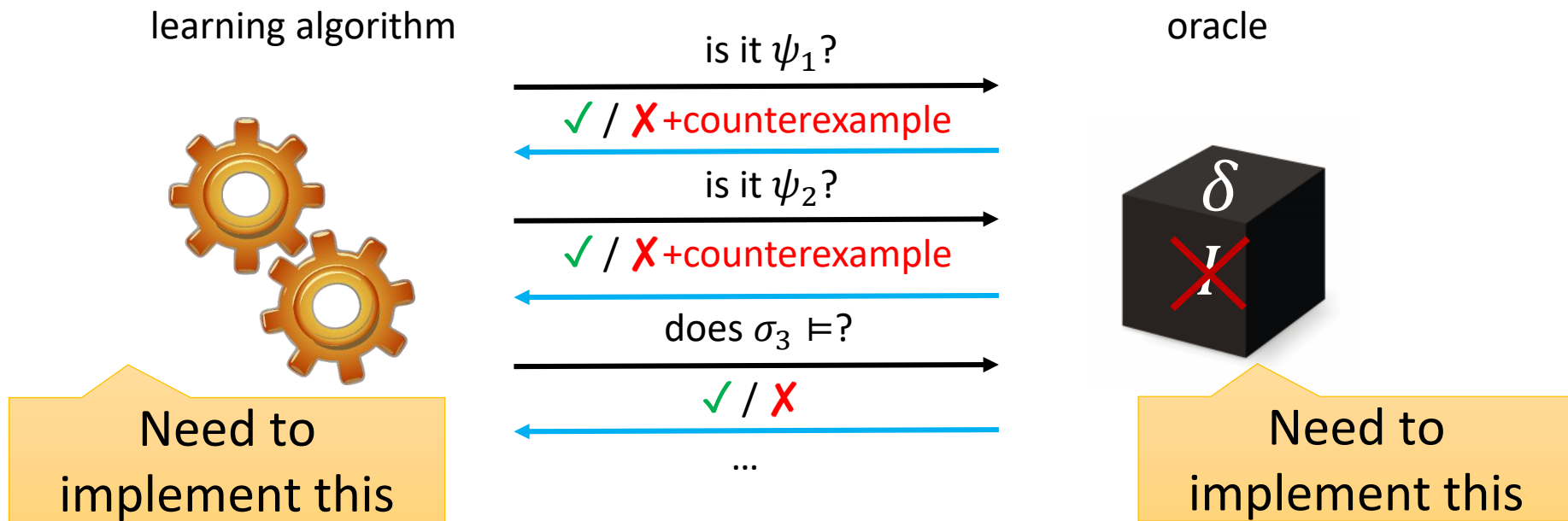
# Understanding Invariant Inference

```
 $I := \text{false}$   
while  $(\_, \sigma')$  counterexample  
to Inductive $(\delta, I)$ :  
   $I := I \vee \text{generalize}(\sigma')$   
  
generalize $(\sigma')$ :  
  drop literals from  $\sigma'$   
  while BMC $(\delta, \sigma', k) \cap \text{Bad} = \emptyset$ 
```

Complexity bounds from  
exact classification algorithms

Rich SAT queries allow  
exponentially faster inference

# Invariant Inference with Equivalence & Membership Queries

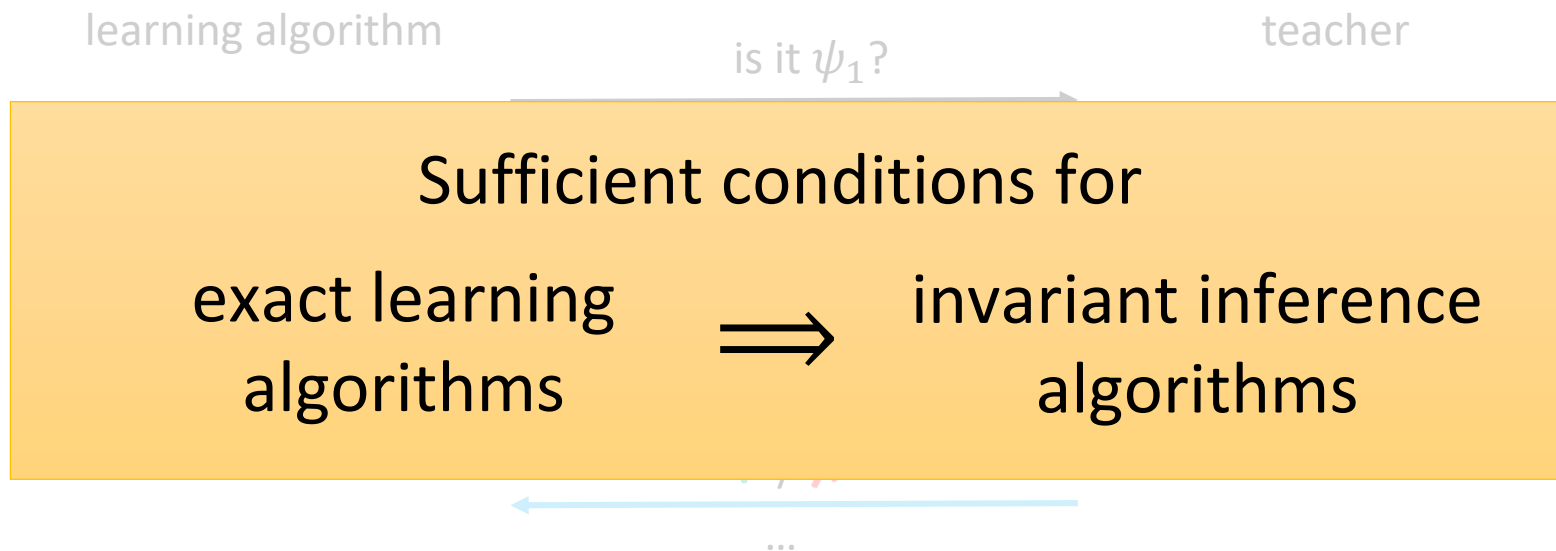


Thm. In general, in the Hoare-query model, **no efficient way** to implement a teacher for equivalence and membership queries

[CAV'14] ICE: A Robust Framework for Learning Invariants. Garg, Löding, Madhusudan, Neider

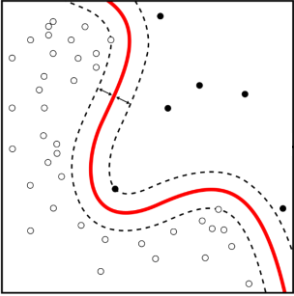
[POPL'20] Complexity and Information in Invariant Inference. Feldman, Immerman, Shoham, Sagiv

# Invariant Inference with Equivalence & Membership Queries



Thm. In general, in the Hoare-query model, **no efficient way** to implement a teacher for equivalence and membership queries

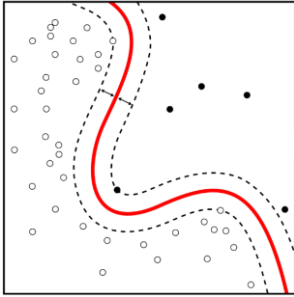
# From Learning to Inference



Exact **learning**  
DNF formulas

```
 $\psi := \text{false}$   
while  $\sigma'$  counterexample  
to Equivalence( $\psi$ ):  
   $\psi := \psi \vee \text{generalize}(\sigma')$   
  
generalize( $\sigma'$ ):  
  drop literals from  $\sigma'$   
  while Membership( $\sigma'$ ) =  $\checkmark$ 
```

# From Learning to Inference



Exact learning  
DNF formulas



```
 $\psi := \text{false}$ 
```

```
while  $\sigma'$  counterexample  
to Equivalence( $\psi$ ):
```

```
 $\psi := \psi \vee \text{generalize}(\sigma')$ 
```

```
generalize( $\sigma'$ ):
```

```
drop literals from  $\sigma'$ 
```

```
while Membership( $\sigma'$ ) =  $\checkmark$ 
```

**Inductive**( $I$ )

**BMC**( $\sigma', k$ )  $\cap$  **Bad** =  $\emptyset$

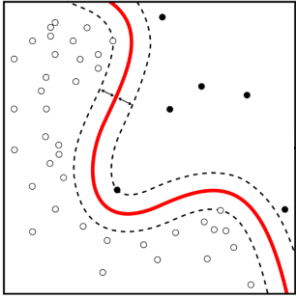
[CACM'84] A Theory of the Learnable. Valiant

[ML'87] Queries and Concept Learning. Angluin

[ML'95] On the Learnability of Disjunctive Normal Form

Formulas. Aizenstein and Pitt

# From Learning to Inference



Exact learning  
DNF formulas



Inferring  
DNF invariants

```
 $\psi := \text{false}$   
while  $\sigma'$  counterexample  
to Equivalence( $\psi$ ):  
   $\psi := \psi \vee \text{generalize}(\sigma')$ 
```

```
generalize( $\sigma'$ ):  
  drop literals from  $\sigma'$   
  while Membership( $\sigma'$ ) =  $\checkmark$ 
```

```
 $I := \text{false}$   
while ( $\_, \sigma'$ ) counterexample  
to Inductive( $I$ ):  
   $I := I \vee \text{generalize}(\sigma')$ 
```

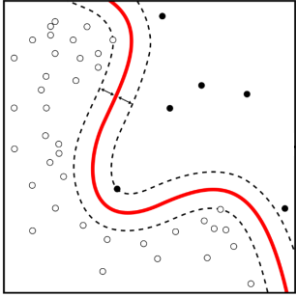
```
generalize( $\sigma'$ ):  
  drop literals from  $\sigma'$   
  while BMC( $\sigma', k$ )  $\cap$  Bad =  $\emptyset$ 
```

[CACM'84] A Theory of the Learnable. Valiant  
[ML'87] Queries and Concept Learning. Angluin  
[ML'95] On the Learnability of Disjunctive Normal Form  
Formulas. Aizenstein and Pitt

[CAV'03] Interpolation and SAT-Based Model Checking,  
McMillan  
[HVC'12] Computing Interpolants without Proofs.  
Chockler, Ivrii, Matsliah



# From Learning to Inference



Efficiently

Exact learning  
DNF formulas



Efficiently

Inferring  
DNF invariants

$\psi := \text{false}$

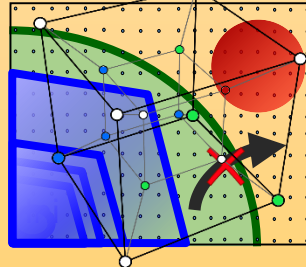
while  $\sigma'$  counterexample  
to **Equivalent**

$\psi := \psi \vee \text{generalize}(\sigma')$

**generalize**( $\sigma'$ ):

drop literals from  $\sigma'$

while **Membership**( $\sigma'$ ) =  $\checkmark$



The invariant is  
**k-fenced**

( $\sigma'$ ) counterexample  
**Inductive**( $I$ ):

$\psi \vee \text{generalize}(\sigma')$

**generalize**( $\sigma'$ ):

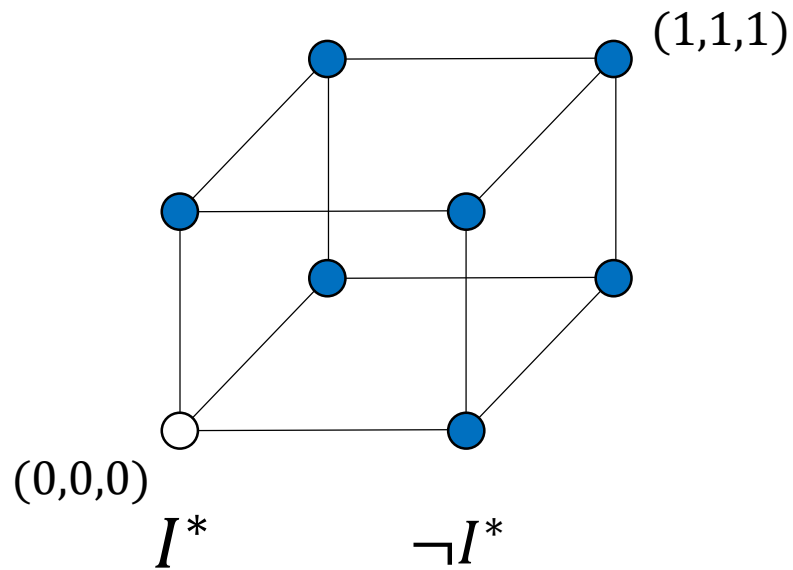
drop literals from  $\sigma'$

while **BMC**( $\sigma', k$ )  $\cap$  **Bad** =  $\emptyset$

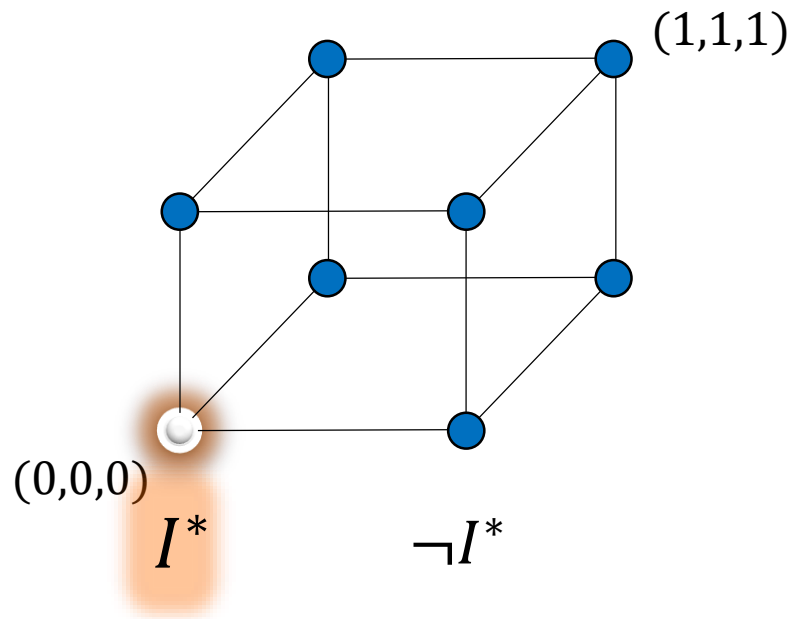
[CACM'84] A Theory of the Learnable. Valiant  
 [ML'87] Queries and Concept Learning. Angluin  
 [ML'95] On the Learnability of Disjunctive Normal Form  
 Formulas. Aizenstein and Pitt

[CAV'03] Interpolation and SAT-Based Model Checking,  
 McMillan  
 [HVC'12] Computing Interpolants without Proofs.  
 Chockler, Ivrii, Matsliah

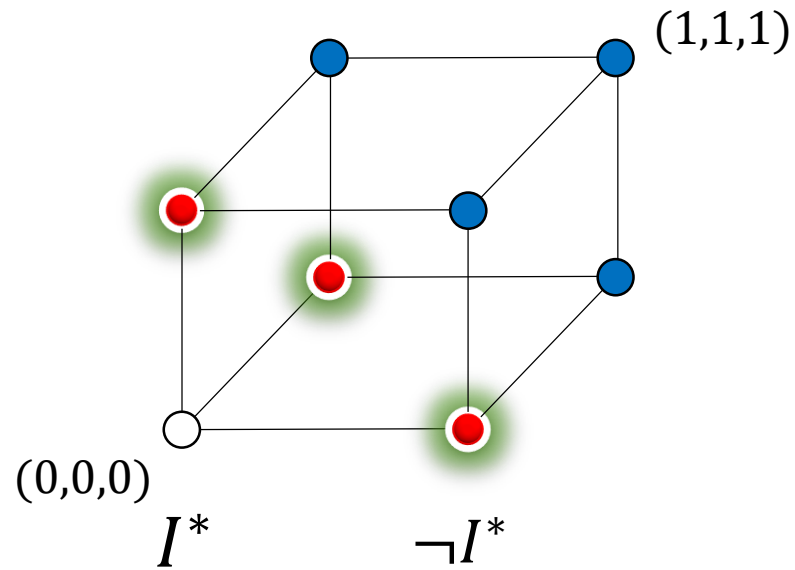
# $k$ -Fenced Invariants



# $k$ -Fenced Invariants

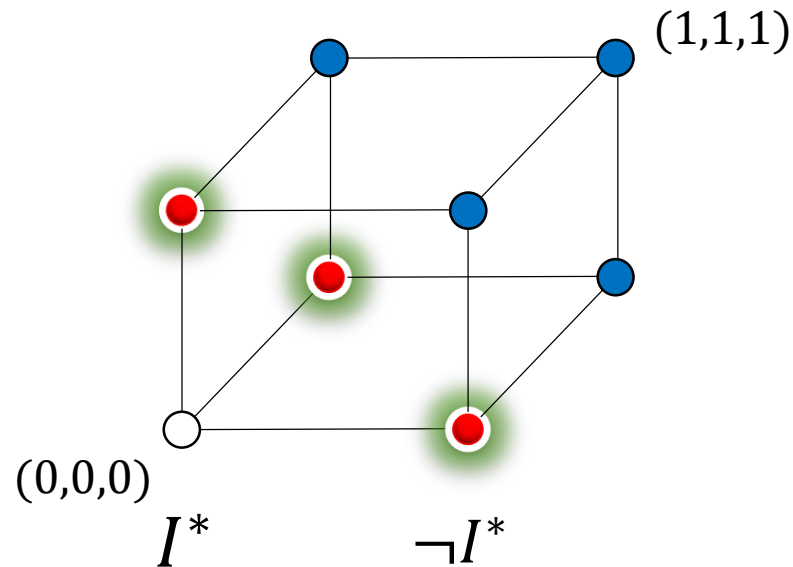


# $k$ -Fenced Invariants



$$\partial^-(I^*)$$

# $k$ -Fenced Invariants

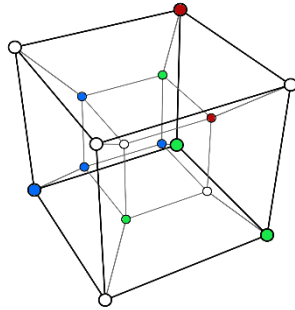


All the states in  $\partial^-(I^*)$   
can get to a bad state in at most  $k$  steps

# Complexity Upper Bounds

Thm. Interpolation-based inference finds an invariant in a polynomial number of SAT queries when

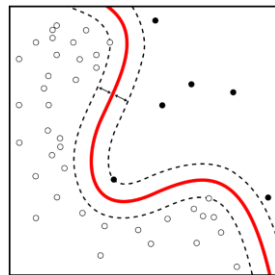
$\exists I^*$ .



Fence condition: the Hamming boundary of  $I^*$  reaches bad states in  $k$  steps

+

No negated variables



$I^*$  is a short monotone DNF (via Angluin)

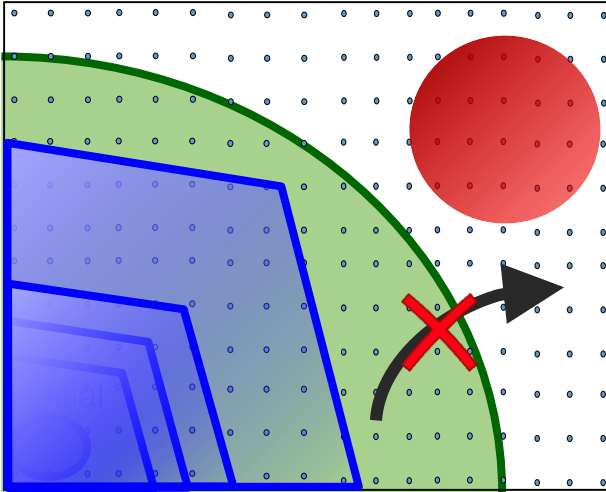
or

$I^*$  is a short almost-monotone DNF (via Bshouty)

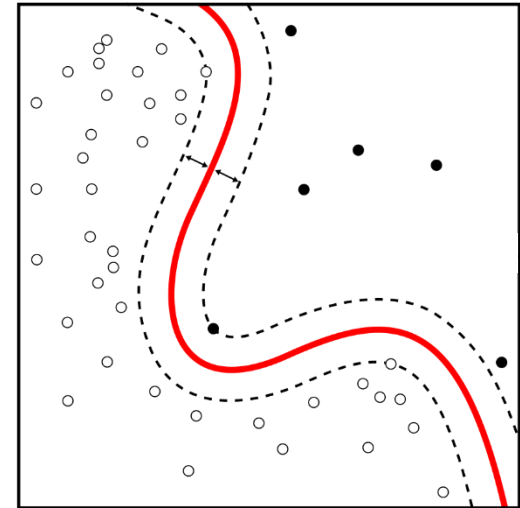
$O(1)$  terms with negated variables

# Conclusion

Invariant Inference



Exact Concept Learning



VS.

- Query-based learning models for invariant inference
- Invariant inference is harder than concept learning
- Complexity results for invariant inference algorithms from classification algorithms