

CBMC: Bounded Model Checking for ANSI-C

The logo for CBMC, consisting of the letters 'CBMC' in a bold, orange-to-yellow gradient font with a slight 3D effect and a drop shadow.

Version 1.0, 2010

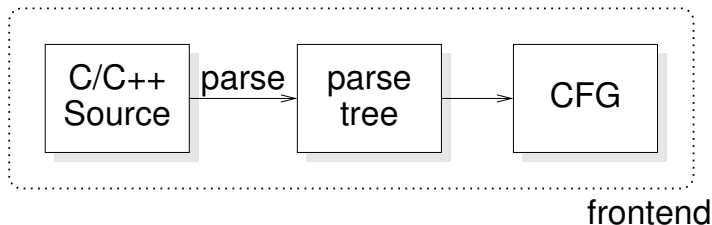
Preliminaries

BMC Basics

Completeness

Solving the Decision Problem

- ▶ We aim at the analysis of programs given in a commodity programming language such as C, C++, or Java
- ▶ As the first step, we transform the program into a *control flow graph* (CFG)



Example: SHS

```
if ( ( 0 <= t ) && ( t <= 79 ) )
  switch ( t / 20 )
  {
  case 0:
    TEMP2 = ( ( B AND C ) OR ( ~B AND D ) );
    TEMP3 = ( K_1 );
    break;

  case 1:
    TEMP2 = ( ( B XOR C XOR D ) );
    TEMP3 = ( K_2 );
    break;

  case 2:
    TEMP2 = ( ( B AND C ) OR ( B AND D ) OR ( C AND D ) );
    TEMP3 = ( K_3 );
    break;

  case 3:
    TEMP2 = ( B XOR C XOR D );
    TEMP3 = ( K_4 );
    break;

  default:
    assert(0);
  }
```

Example: SHS

```

if ( ( 0 <= t ) && ( t <= 79 ) )
  switch ( t / 20 )
  {
  case 0:
    TEMP2 = ( ( B AND C ) OR ( ~B AND D ) );
    TEMP3 = ( K_1 );
    break;

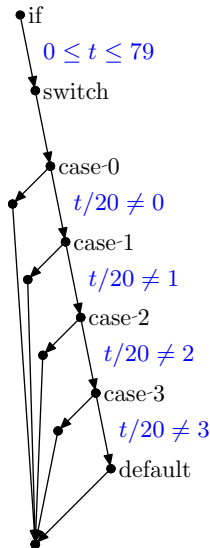
  case 1:
    TEMP2 = ( ( B XOR C XOR D ) );
    TEMP3 = ( K_2 );
    break;

  case 2:
    TEMP2 = ( ( B AND C ) OR ( B AND D ) OR ( C AND D ) );
    TEMP3 = ( K_3 );
    break;

  case 3:
    TEMP2 = ( B XOR C XOR D );
    TEMP3 = ( K_4 );
    break;

  default:
    assert(0);
  }

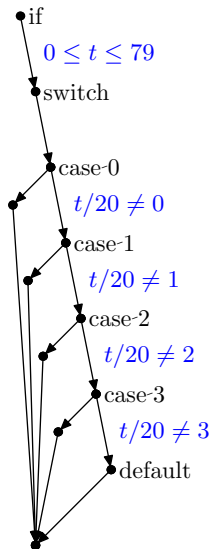
```



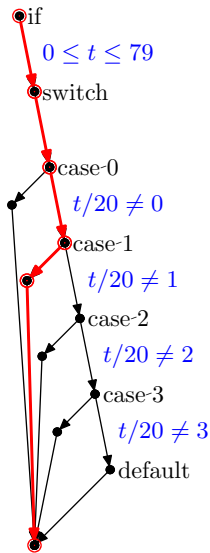
Goal: check properties of the form $\mathbf{AG}p$,
say assertions.

Idea: follow paths through the CFG to an assertion,
and build a formula that corresponds to the path

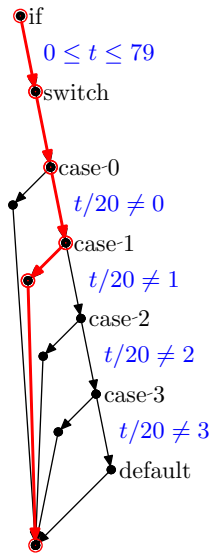
Example



Example



Example



$$\begin{aligned}
 & 0 \leq t \leq 79 \\
 \wedge & \quad t/20 \neq 0 \\
 \wedge & \quad t/20 = 1 \\
 \wedge & \quad TEMP2 = B \oplus C \oplus D \\
 \wedge & \quad TEMP3 = K_2
 \end{aligned}$$

We pass

$$\begin{aligned} & 0 \leq t \leq 79 \\ \wedge & \quad t/20 \neq 0 \\ \wedge & \quad t/20 = 1 \\ \wedge & \quad TEMP2 = B \oplus C \oplus D \\ \wedge & \quad TEMP3 = K_2 \end{aligned}$$

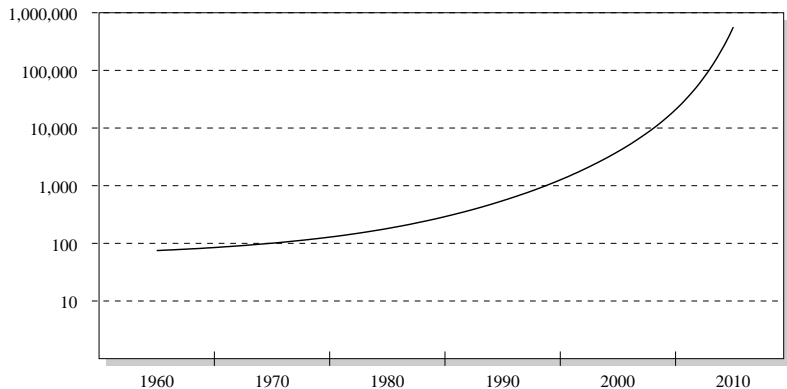
to a decision procedure, and obtain a **satisfying assignment**, say:

$$\begin{aligned} t &\mapsto 21, B \mapsto 0, C \mapsto 0, D \mapsto 0, K_2 \mapsto 10, \\ TEMP2 &\mapsto 0, TEMP3 \mapsto 10 \end{aligned}$$

✓ It provides the values of any inputs on the path.

- ▶ We need a decision procedure for an appropriate logic
 - ▶ Bit-vector logic (incl. non-linear arithmetic)
 - ▶ Arrays
 - ▶ Higher-level programming languages also feature lists, sets, and maps

- ▶ Examples
 - ▶ Z3 (Microsoft)
 - ▶ Yices (SRI)
 - ▶ Boolector

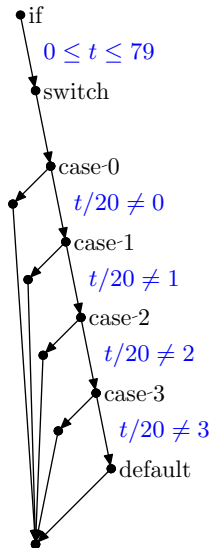


number of variables of a typical, practical SAT instance
that can be solved by the best solvers in that decade

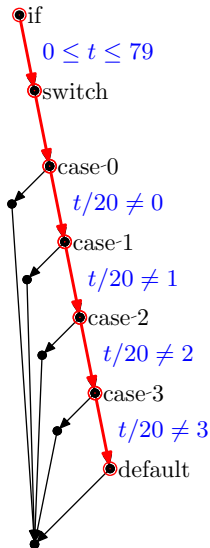
- ▶ propositional SAT solvers have made enormous progress in the last 10 years

- ▶ Further scalability improvements in recent years because of efficient **word-level reasoning** and **array decision procedures**

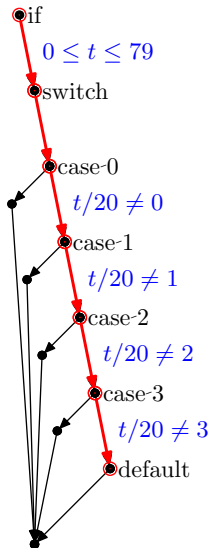
Let's Look at Another Path



Let's Look at Another Path

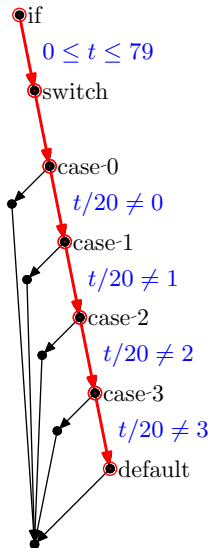


Let's Look at Another Path



$$\begin{aligned}
 & 0 \leq t \leq 79 \\
 \wedge & \quad t/20 \neq 0 \\
 \wedge & \quad t/20 \neq 1 \\
 \wedge & \quad t/20 \neq 2 \\
 \wedge & \quad t/20 \neq 3
 \end{aligned}$$

Let's Look at Another Path



$$\begin{aligned}
 & 0 \leq t \leq 79 \\
 & \wedge t/20 \neq 0 \\
 & \wedge t/20 \neq 1 \\
 & \wedge t/20 \neq 2 \\
 & \wedge t/20 \neq 3
 \end{aligned}$$

That is UNSAT, so the assertion is unreachable.

What If a Variable is Assigned Twice?

```
x=0;
```

```
if(y>=0)  
  x++;
```



Rename appropriately:

```
  x = 0  
∧ y ≥ 0  
∧ x = x + 1
```

What If a Variable is Assigned Twice?

```
x=0;
```

```
if(y>=0)  
  x++;
```



Rename appropriately:

```
 $x_1 = 0$   
 $\wedge y_0 \geq 0$   
 $\wedge x_1 = x_0 + 1$ 
```

This is a special case of *SSA* (static single assignment)

How do we handle dereferencing in the program?

How do we handle dereferencing in the program?

```
int *p;
```

```
p=malloc(sizeof(int)*5);
```

```
...
```

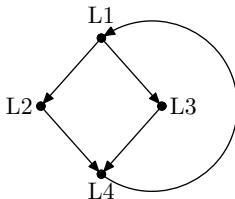
```
p[1]=100;
```



$$\begin{aligned}
 & p_1 = \&DO1 \\
 \wedge & DO1_1 = (\lambda i. \\
 & i = 1?100 : DO1_0[i])
 \end{aligned}$$

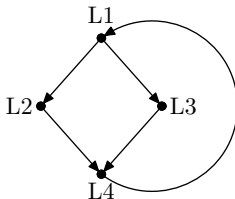
Track a 'may-point-to' abstract state while simulating!

Let's consider the following CFG:



This is a loop with an `if` inside.

Let's consider the following CFG:



This is a loop with an `if` inside.

Q: how many paths for n iterations?

- ▶ Bounded Model Checking (BMC) is the most successful formal validation technique in the *hardware* industry
- ▶ Advantages:
 - ✓ Fully automatic
 - ✓ Robust
 - ✓ Lots of subtle bugs found
- ▶ Idea: only look for bugs up to specific depth
- ▶ Good for many applications, e.g., embedded systems

Definition: A transition system is a triple (S, S_0, T) with

- ▶ set of states S ,
- ▶ a set of initial states $S_0 \subset S$, and
- ▶ a transition relation $T \subset (S \times S)$.

The set S_0 and the relation T can be written as their characteristic functions.

Q: How do we avoid the exponential path explosion?

We just "concatenate" the transition relation T :

S_0
●

Q: How do we avoid the exponential path explosion?

We just "concatenate" the transition relation T :

$$\bullet \xrightarrow{S_0 \wedge T} \bullet$$

Q: How do we avoid the exponential path explosion?

We just "concatenate" the transition relation T :



Q: How do we avoid the exponential path explosion?

We just "concatenate" the transition relation T :



Q: How do we avoid the exponential path explosion?

We just "concatenate" the transition relation T :



As formula:

$$S_0(s_0) \wedge \bigwedge_{i=0}^{k-1} T(s_i, s_{i+1})$$

Satisfying assignments for this formula are **traces** through the transition system

$$T \subseteq \mathbb{N}_0 \times \mathbb{N}_0$$

$$T(s, s') \iff s'.x = s.x + 1$$

... and let $S_0(s) \iff s.x = 0 \vee s.x = 1$

$$T \subseteq \mathbb{N}_0 \times \mathbb{N}_0$$

$$T(s, s') \iff s'.x = s.x + 1$$

... and let $S_0(s) \iff s.x = 0 \vee s.x = 1$

An unwinding for depth 4:

$$\begin{aligned} & (s_0.x = 0 \vee s_0.x = 1) \\ \wedge & \quad s_1.x = s_0.x + 1 \\ \wedge & \quad s_2.x = s_1.x + 1 \\ \wedge & \quad s_3.x = s_2.x + 1 \\ \wedge & \quad s_4.x = s_3.x + 1 \end{aligned}$$

Suppose we want to check a property of the form $\mathbf{AG}p$.

Suppose we want to check a property of the form $\mathbf{AG}p$.

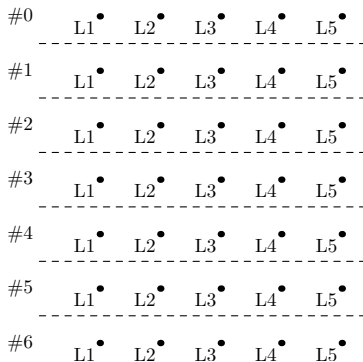
We then want at **least one state** s_i to satisfy $\neg p$:

$$S_0(s_0) \wedge \bigwedge_{i=0}^{k-1} T(s_i, s_{i+1}) \quad \wedge \quad \bigvee_{i=0}^k \neg p(s_i)$$

Satisfying assignments are **counterexamples** for the $\mathbf{AG}p$ property

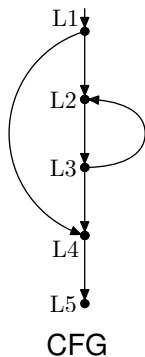
We can do exactly that for our transition relation for software.

E.g., for a program with 5 locations, 6 unwindings:

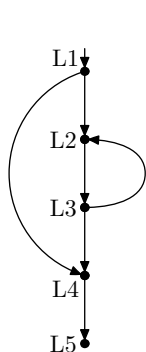


Problem: obviously, most of the formula is never 'used',
as only few sequences of PCs correspond to a path.

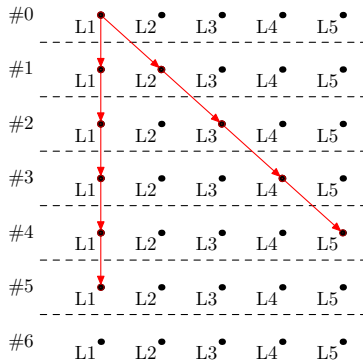
Example:



Example:



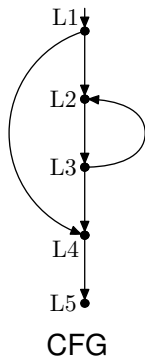
CFG



unrolling

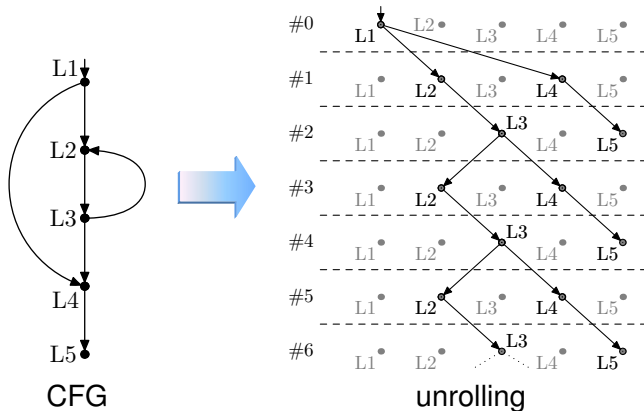
Optimization:

don't generate the parts of the formula that are not 'reachable'

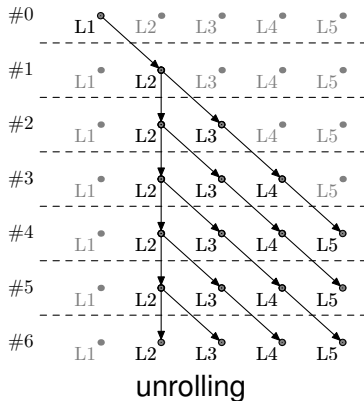
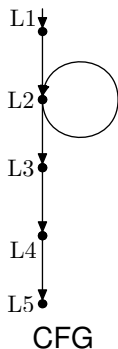


Optimization:

don't generate the parts of the formula that are not 'reachable'



Problem:

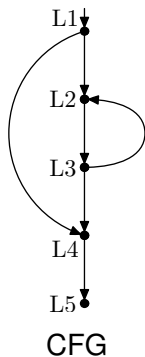


- ▶ Unwinding T with bound k results in a formula of size

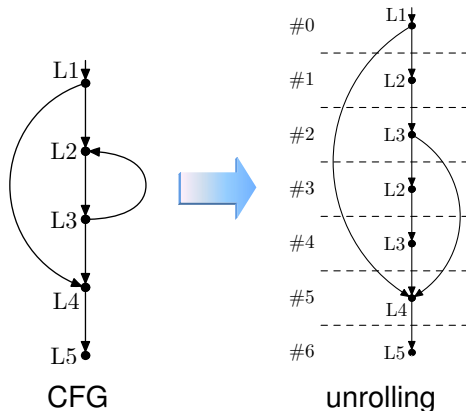
$$|T| \cdot k$$

- ▶ If we assume a k that is only linear in $|T|$, we get a formula with size $O(|T|^2)$
- ▶ Can we do better?

Idea: do **exactly one location** in each timeframe:



Idea: do **exactly one location** in each timeframe:



- ✓ More effective use of the formula size
- ✓ Graph has fewer merge nodes,
the formula is easier for the solvers
- ✗ Not all paths of length k are encoded
→ the bound needs to be larger

Unrolling Loops

This essentially amounts to unwinding loops:

```
while(cond)
  Body;
```

Unrolling Loops

This essentially amounts to unwinding loops:

```
if(cond) {  
    Body;  
    while(cond)  
        Body;  
}
```


This essentially amounts to unwinding loops:

```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    while(cond)  
      Body;  
  }  
}
```

This essentially amounts to unwinding loops:

```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    if(cond) {  
      Body;  
      while(cond)  
        Body;  
    }  
  }  
}
```

This essentially amounts to unwinding loops:

```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    if(cond) {  
      Body;  
      assume(!cond);  
    }  
  }  
}
```

BMC, as discussed so far, is incomplete.
It only refutes, and does not prove.

How can we fix this?

Unwinding Assertions

Let's revisit the loop unwinding idea:

```
while(cond)
  Body;
```

Unwinding Assertions

Let's revisit the loop unwinding idea:

```
if(cond) {  
    Body;  
    while(cond)  
        Body;  
}
```

Let's revisit the loop unwinding idea:

```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    while(cond)  
      Body;  
  }  
}
```

Let's revisit the loop unwinding idea:

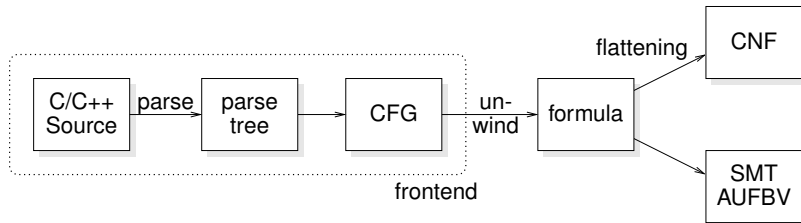
```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    if(cond) {  
      Body;  
      while(cond)  
        Body;  
    }  
  }  
}
```


Let's revisit the loop unwinding idea:

```
if(cond) {  
  Body;  
  if(cond) {  
    Body;  
    if(cond) {  
      Body;  
      assert(!cond);  
    }  
  }  
}
```

- ▶ We replace the assumption we have used earlier to cut off paths by an assertion
- ✓ This allows us to **prove that we have done enough unwinding**
- ▶ This is a proof of a high-level worst-case execution time (WCET)
- ▶ Very appropriate for embedded software

1. Parse, build CFG
2. Unwind CFG, form formula
3. Formula is solved by SAT/SMT



Suppose we have used some unwinding, and have built the formula.

For bit-vector arithmetic, the standard way of deciding satisfiability of the formula is *flattening*, followed by a call to a propositional SAT solver.

In the SMT context: SMT-BV

- ▶ This is easy for the bit-wise operators.
- ▶ Denote the Boolean variable for bit i of term t by $\mu(t)_i$.
- ▶ Example for $a \mid_{[l]} b$:

$$\bigwedge_{i=0}^{l-1} (\mu(t)_i = (a_i \vee b_i))$$

(read $x = y$ over bits as $x \iff y$)

- ▶ This is easy for the bit-wise operators.
- ▶ Denote the Boolean variable for bit i of term t by $\mu(t)_i$.
- ▶ Example for $a \mid_{[l]} b$:

$$\bigwedge_{i=0}^{l-1} (\mu(t)_i = (a_i \vee b_i))$$

(read $x = y$ over bits as $x \iff y$)

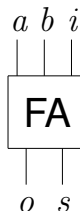
- ▶ We can transform this into CNF using Tseitin's method.

Flattening Bit-Vector Arithmetic

How to flatten $a + b$?

How to flatten $a + b$?

→ we can build a *circuit* that adds them!



Full Adder

$$s \equiv (a + b + i) \bmod 2 \equiv a \oplus b \oplus i$$

$$o \equiv (a + b + i) \operatorname{div} 2 \equiv a \cdot b + a \cdot i + b \cdot i$$

The full adder in CNF:

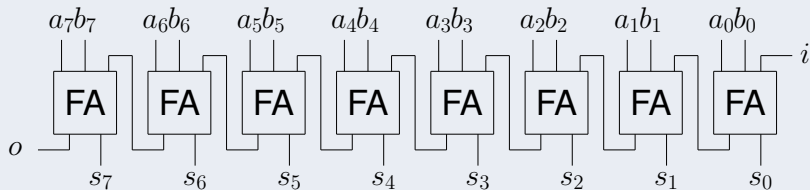
$$(a \vee b \vee \neg o) \wedge (a \vee \neg b \vee i \vee \neg o) \wedge (a \vee \neg b \vee \neg i \vee o) \wedge \\ (\neg a \vee b \vee i \vee \neg o) \wedge (\neg a \vee b \vee \neg i \vee o) \wedge (\neg a \vee \neg b \vee o)$$

Flattening Bit-Vector Arithmetic

Ok, this is good for one bit! How about more?

Ok, this is good for one bit! How about more?

8-Bit ripple carry adder (RCA)



- ▶ Also called *carry chain adder*
- ▶ Adds l variables
- ▶ Adds $6 \cdot l$ clauses

- ▶ **Multipliers** result in very hard formulas
- ▶ Example:

$$a \cdot b = c \wedge b \cdot a \neq c \wedge x < y \wedge x > y$$

CNF: About 11000 variables,
unsolvable for current SAT solvers

- ▶ Similar problems with division, modulo
- ▶ Q: Why is this hard?

- ▶ **Multipliers** result in very hard formulas
- ▶ Example:

$$a \cdot b = c \wedge b \cdot a \neq c \wedge x < y \wedge x > y$$

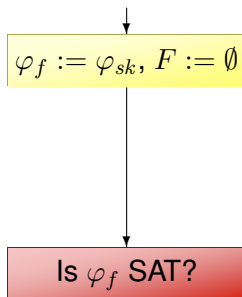
CNF: About 11000 variables,
unsolvable for current SAT solvers

- ▶ Similar problems with division, modulo
- ▶ Q: Why is this hard?
- ▶ Q: How do we fix this?

$$\downarrow$$
$$\varphi_f := \varphi_{sk}, F := \emptyset$$

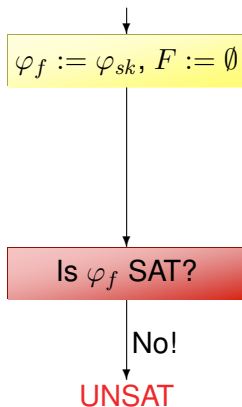
φ_{sk} : Boolean part of φ

F : set of terms that are in the encoding



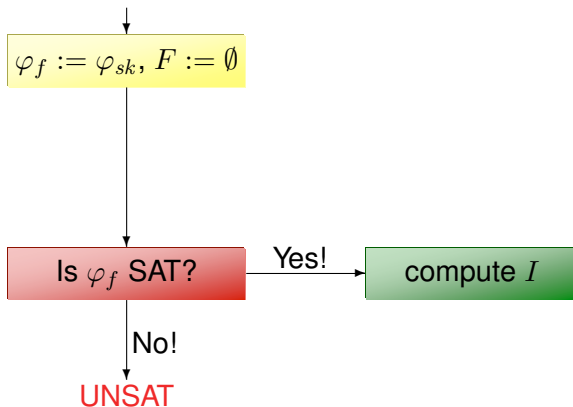
φ_{sk} : Boolean part of φ

F : set of terms that are in the encoding



φ_{sk} : Boolean part of φ

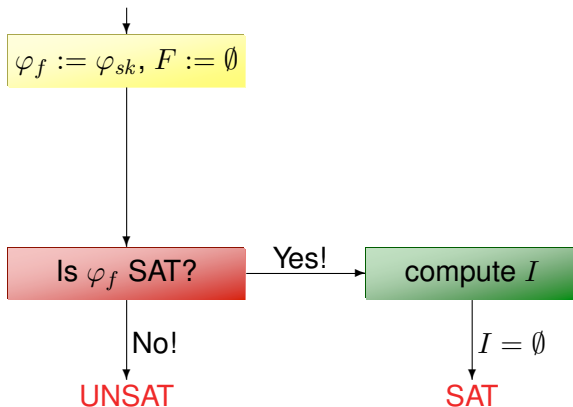
F : set of terms that are in the encoding



φ_{sk} : Boolean part of φ

F : set of terms that are in the encoding

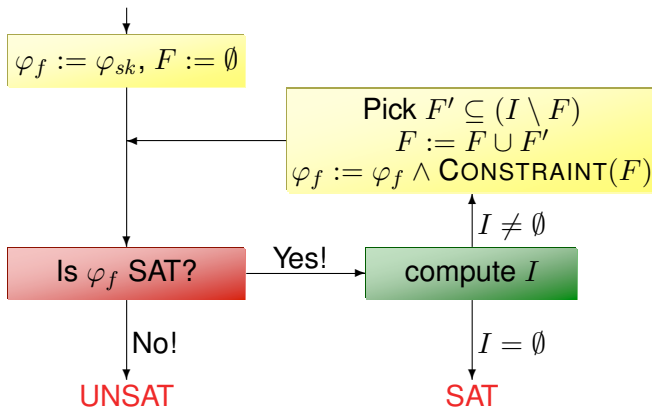
I : set of terms that are inconsistent with the current assignment



φ_{sk} : Boolean part of φ

F : set of terms that are in the encoding

I : set of terms that are inconsistent with the current assignment



φ_{sk} : Boolean part of φ

F : set of terms that are in the encoding

I : set of terms that are inconsistent with the current assignment

- ▶ Idea: add 'easy' parts of the formula first
- ▶ Only add hard parts when needed
- ▶ φ_f only gets stronger – use an **incremental SAT solver**