

Toward Grammar Inference via Refinement Types

<https://mcschroeder.github.io/#tyde2022>



Michael Schröder

TU Wien

Vienna, Austria

michael.schroeder@tuwien.ac.at



Jürgen Cito

TU Wien

Vienna, Austria

juergen.cito@tuwien.ac.at

Type-Driven Development (TyDe), ICFP 2022
Ljubljana, Slovenia

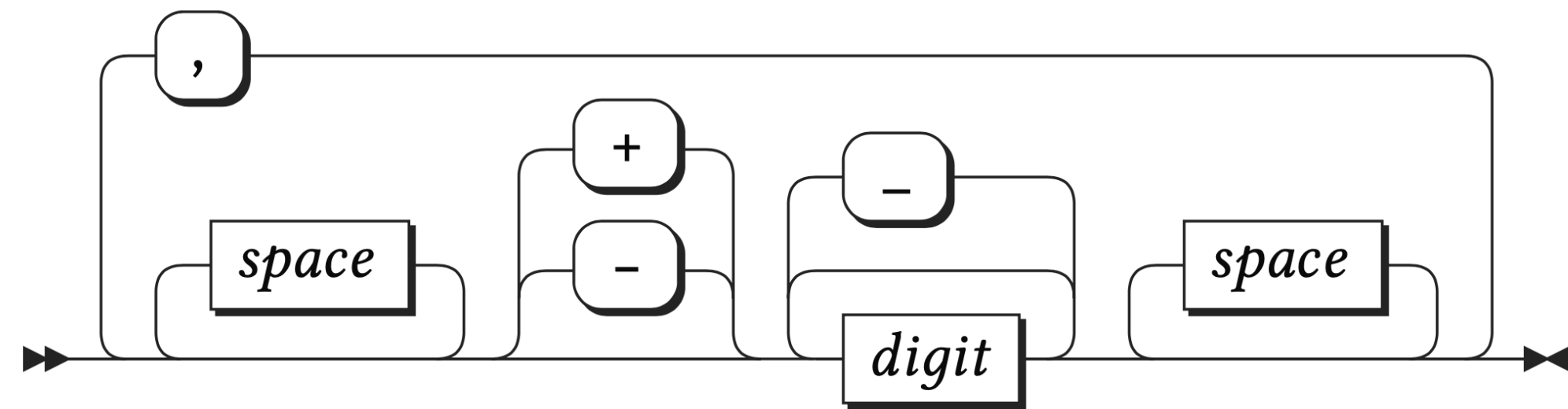
TU
WIEN Informatics

```
xs = map(int, s.split(","))
```

`xs = map(int, s.split(","))`

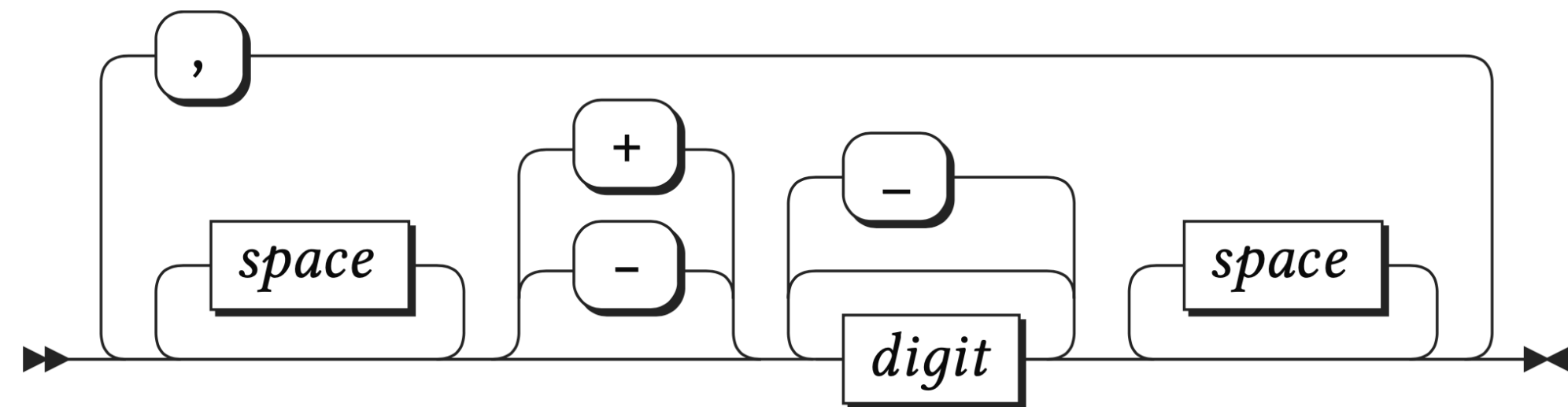
`"1,2,3"` →

→ `[1,2,3]`



$$\begin{aligned}
 s &\rightarrow int \mid int , s \\
 int &\rightarrow space^* (+ \mid -)^? digit (_? digit)^* space^* \\
 digit &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\
 space &\rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r
 \end{aligned}$$

```
xs = map(int, s.split(","))
```



$$s \rightarrow int \mid int , s$$

$$int \rightarrow space^* (+ \mid -)^? digit (_? digit)^* space^*$$

$$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

$$space \rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r$$

```
xs = map(int, s.split(","))
```

Parser : Grammar \approx Function : Type

Type Inference

```
xs = map(int, s.split(","))
```

Diagram illustrating the types of variables and functions in the code snippet:

- `xs` is annotated with the type `[Int]`.
- `map` is annotated with the type `String`.
- `split` is annotated with the type `String`.

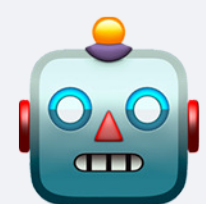
✨ Grammar Inference ✨

$$\begin{aligned}s &\rightarrow int \mid int , s \\int &\rightarrow space^* (+ \mid -)^? digit (_? digit)^* space^* \\digit &\rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9 \\space &\rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r\end{aligned}$$

$[Int]$ $xs = \text{map}(int, s.\text{split}(", "))$ $String \{ \bullet \}$

Applications

semantic change tracking



⚠ Merging #420 (6a36b23) into main (224b18b) will change the input grammar of a function.

Before:
 $baz \rightarrow a^*b\Sigma^*$

After:
 $baz \rightarrow \Sigma a^*b\Sigma^*$

```
18 18 def baz(s):
19 - i = 0
19 + i = 1
20 20 while s[i] == "a"
21 21     i += 1
22 22 assert s[i] == "b"
```

+ fuzz testing, program sketching,
grammar-based refactoring, ...

interactive documentation

Inferred Grammar	Inferred Inputs
$s \rightarrow int \mid int , s$	✗ (empty)
$int \rightarrow space^* (+ \mid -)^? digit (_? digit)^* space^*$	✓ 1,2,3
$digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$	✓ 10_000,4
$space \rightarrow _ \mid \backslash t \mid \backslash n \mid \backslash v \mid \backslash f \mid \backslash r$	✓ +01_2, _ _ 3_

```
xs = map(int, s.split(","))
```

searching for parsers via their grammars (cf. Hoogle)

Viewer.py2 matches | Python

```
640 ranges = self.find_ranges()
641 split = str.split
642 point = map(int, split(self.text.index(CURRENT), ','))
643 for start, end in ranges:
644     startv = map(int, split(start, ','))
```

transformScriptTags.ts1 match | TypeScript

```
122 return null;
123 }
124 return rawValue.split(",").map(item => parseInt(item));
125 }
126
```


Goal: Automatic Grammar Inference

ad hoc parser source

```
def parser(s):  
    if s[0] == "a":  
        assert len(s) == 1  
    else:  
        assert s[1] == "b"
```

?

grammar

$s \rightarrow a \mid (\Sigma \setminus a)b\Sigma^*$

PANINI

- simple λ -calculus in Administrative Normal Form (ANF)
- refinement type system à la *Liquid Types*
- common string operations assumed as axioms
- idea: infer most precise refinement type for input string

ad hoc parser source

```
def parser(s):  
    if s[0] == "a":  
        assert len(s) == 1  
    else:  
        assert s[1] == "b"
```

SSA/ANF transformation

PANINI program

```
assert : {b :  $\mathbb{B}$  | b}  $\rightarrow$  1  
equals : (a :  $\mathbb{Z}$ )  $\rightarrow$  (b :  $\mathbb{Z}$ )  $\rightarrow$  {c :  $\mathbb{B}$  | c  $\Leftrightarrow$  a = b}  
length : (s :  $\mathbb{S}$ )  $\rightarrow$  {n :  $\mathbb{N}$  | n = |s|}  
charAt : (s :  $\mathbb{S}$ )  $\rightarrow$  {i :  $\mathbb{N}$  | i < |s|}  $\rightarrow$  {t :  $\mathbb{S}$  | t = s[i]}  
match : (s :  $\mathbb{S}$ )  $\rightarrow$  (t :  $\mathbb{S}$ )  $\rightarrow$  {b :  $\mathbb{B}$  | b  $\Leftrightarrow$  s = t}
```

```
parser :  $\mathbb{S} \rightarrow$  1  
=  $\lambda s.$ 
```

```
    let x = charAt s 0 in  
    let p1 = match x "a" in  
    if p1 then  
        let n = length s in  
        let p2 = equals n 1 in  
        assert p2  
    else  
        let y = charAt s 1 in  
        let p3 = match y "b" in  
        assert p3
```

Refinement Inference

- κ variables represent unknown refinements
- most can be solved precisely (e.g., using FUSION)
- existing approaches struggle with “grammar variables”

“grammar variable”
(constraint over input string)

verification template

$$\begin{aligned}
 &\forall s. \kappa_0(s) \Rightarrow \\
 &\quad 0 < |s| \wedge \forall x. x = s[0] \Rightarrow \\
 &\quad \quad \forall p_1. p_1 \Leftrightarrow x = \text{"a"} \Rightarrow \\
 &\quad \quad (p_1 \Rightarrow \\
 &\quad \quad \quad \forall n. n \geq 0 \wedge n = |s| \Rightarrow \\
 &\quad \quad \quad \quad \forall p_2. p_2 \Leftrightarrow n = 1 \Rightarrow \\
 &\quad \quad \quad \quad p_2) \\
 &\quad \wedge (\neg p_1 \Rightarrow \\
 &\quad \quad 1 < |s| \wedge \forall y. y = s[1] \Rightarrow \\
 &\quad \quad \quad \forall p_3. p_3 \Leftrightarrow y = \text{"b"} \Rightarrow \\
 &\quad \quad \quad p_3)
 \end{aligned}$$

PANINI program

```

assert : {b : B | b} → 1
equals : (a : Z) → (b : Z) → {c : B | c ⇔ a = b}
length : (s : S) → {n : N | n = |s|}
charAt : (s : S) → {i : N | i < |s|} → {t : S | t = s[i]}
match : (s : S) → (t : S) → {b : B | b ⇔ s = t}

```

```

parser : {s : S |  $\kappa_0(s)$ } → 1
        = λs.

```

```

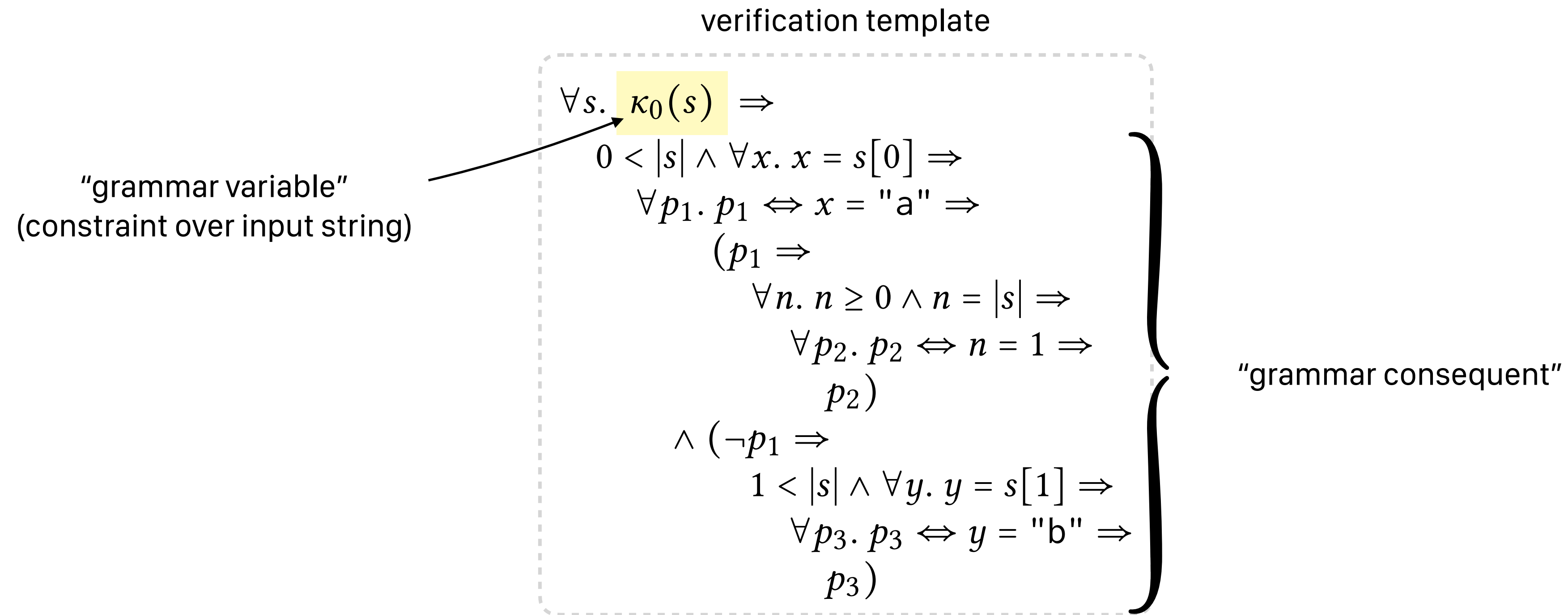
let x = charAt s 0 in
let p1 = match x "a" in
if p1 then
  let n = length s in
  let p2 = equals n 1 in
  assert p2
else
  let y = charAt s 1 in
  let p3 = match y "b" in
  assert p3

```

- Jhala and Vazou. 2020. Refinement Types: A Tutorial. <https://arxiv.org/abs/2010.07763>
- Cosman and Jhala. 2017. Local Refinement Typing. <https://doi.org/10.1145/3110270>
- Rondon et al. 2008. Liquid Types. <https://doi.org/10.1145/1375581.1375602>

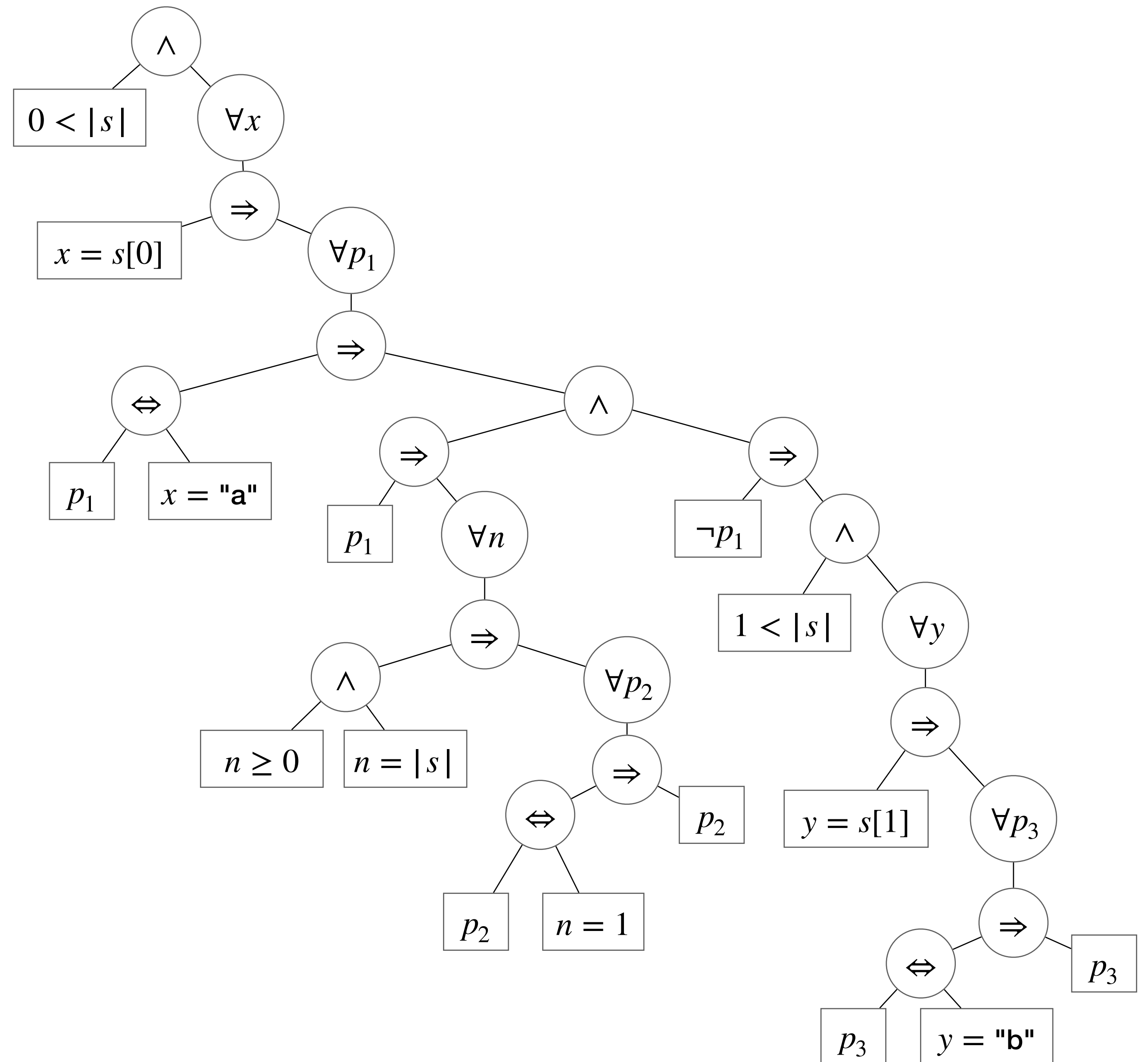
Grammar Solving

- base solution on “grammar consequent”



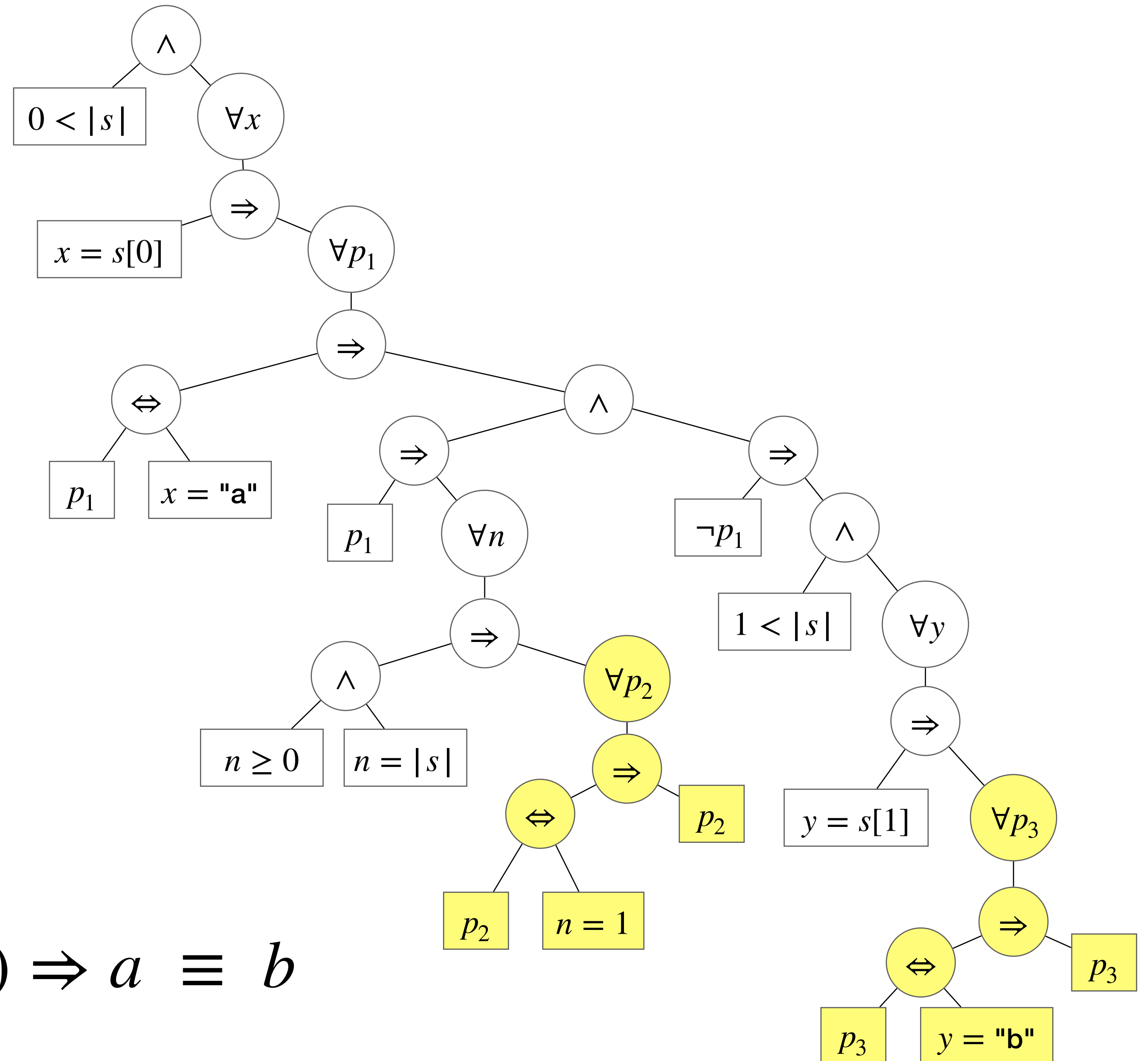
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



Grammar Solving

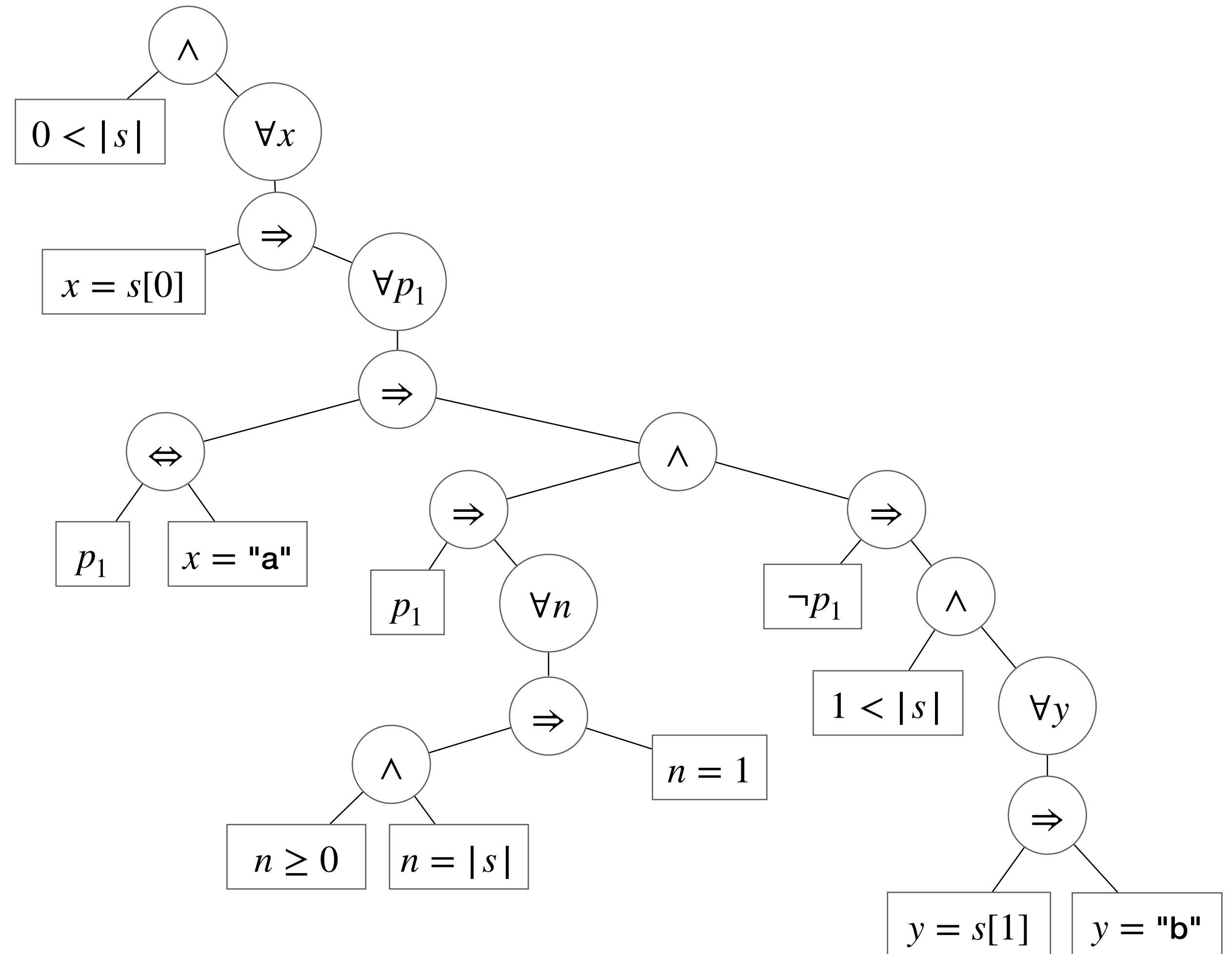
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$\forall a . (a \Leftrightarrow b) \Rightarrow a \equiv b$$

Grammar Solving

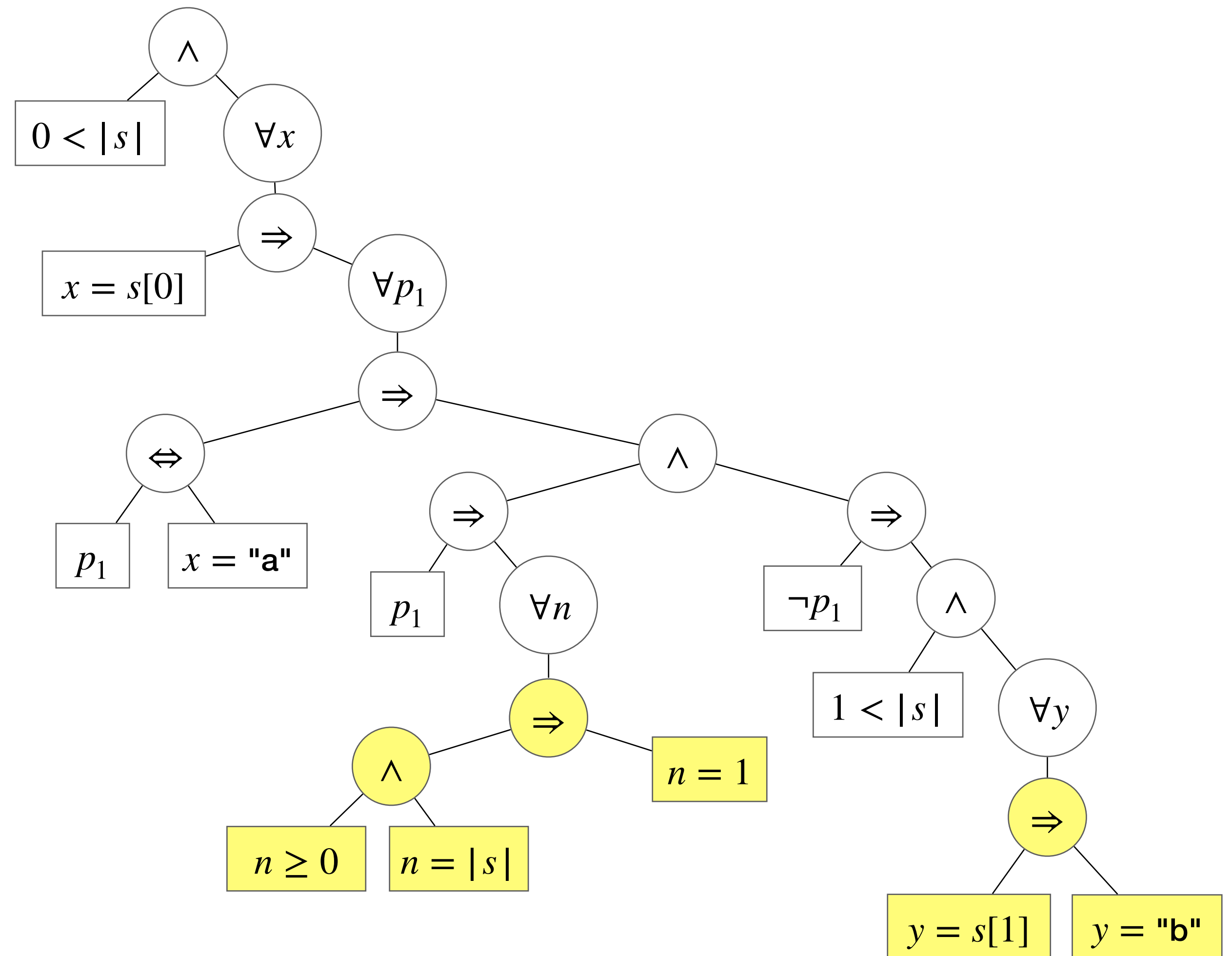
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$\forall a . (a \Leftrightarrow b) \Rightarrow a \equiv b$$

Grammar Solving

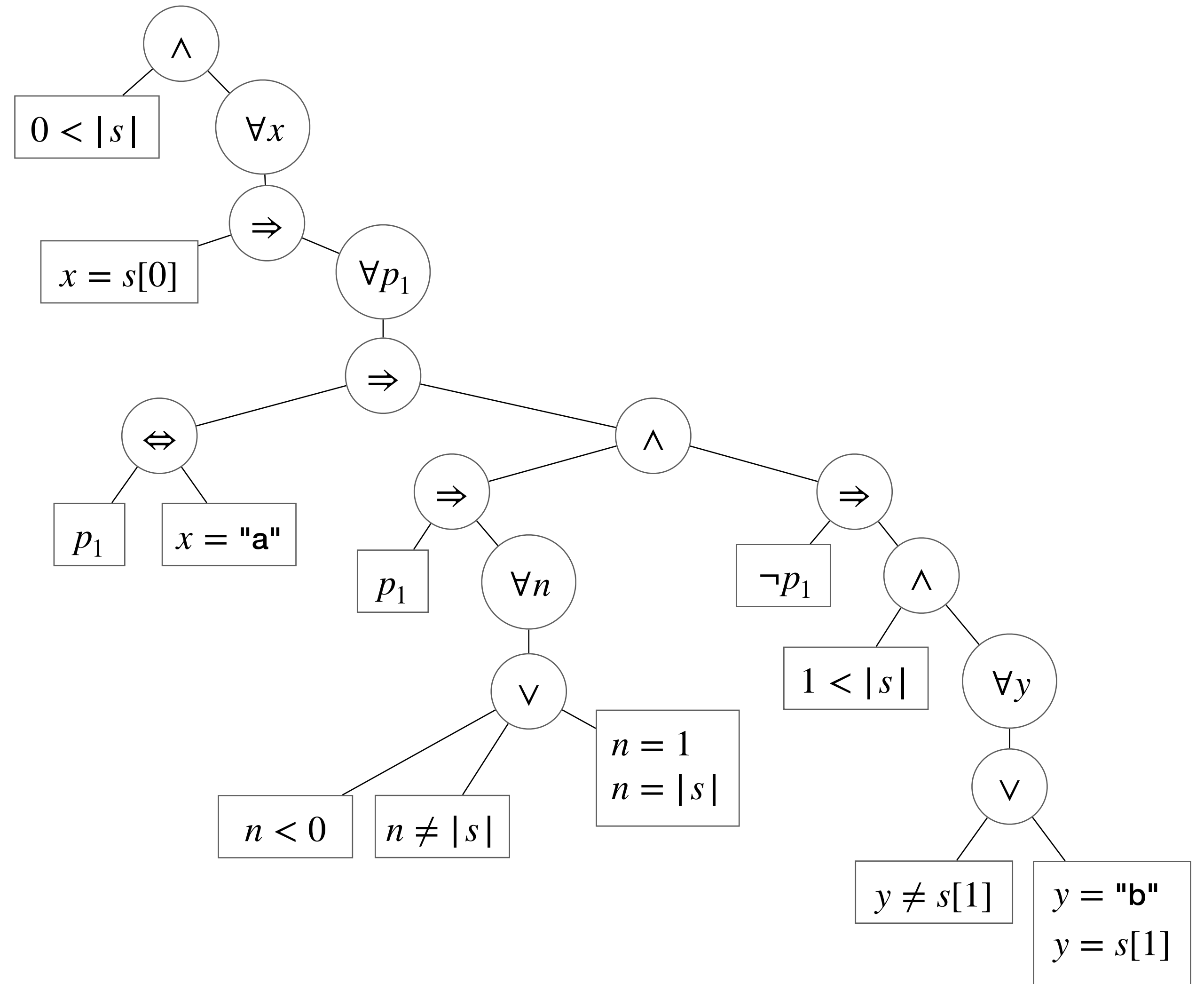
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

Grammar Solving

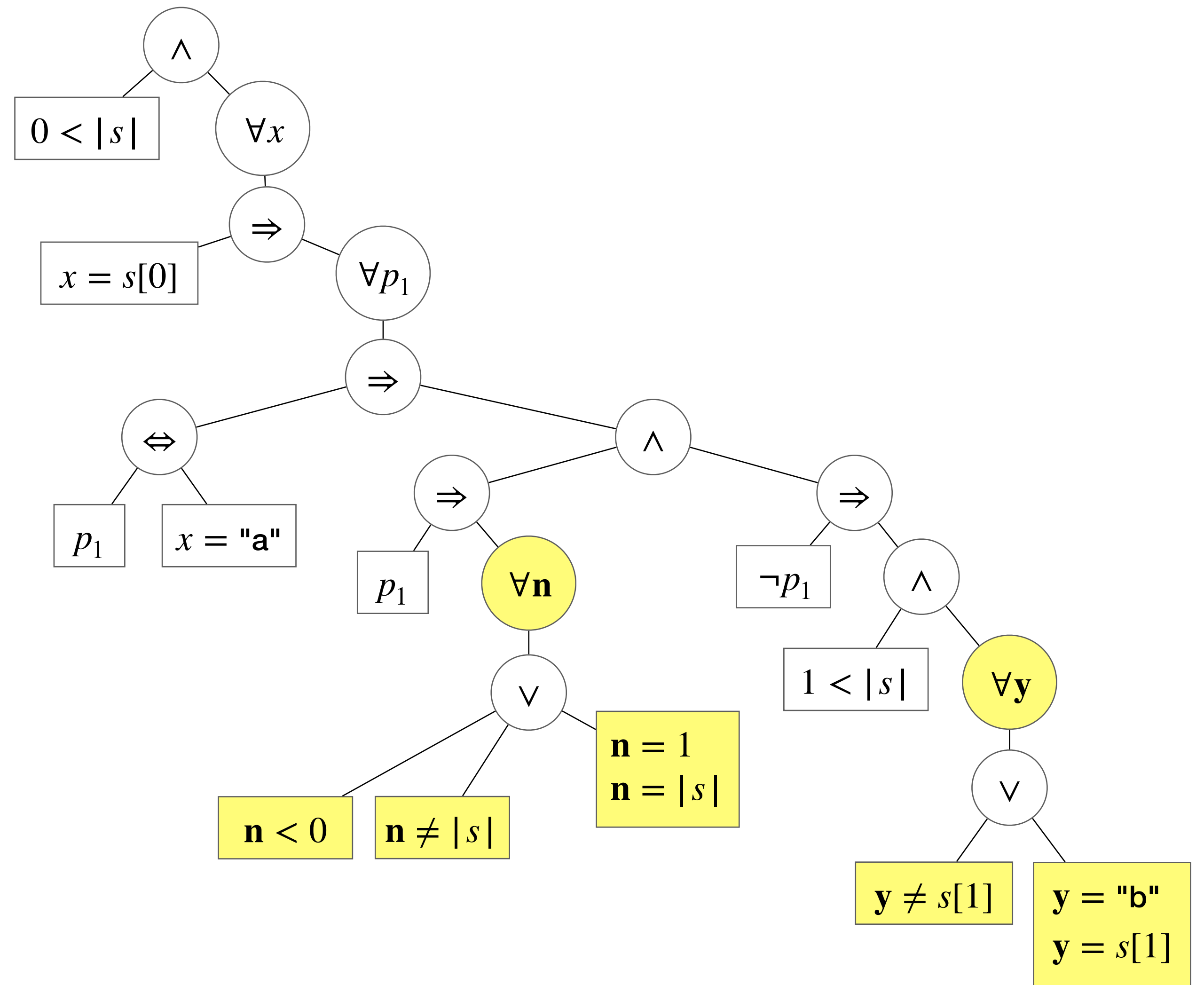
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

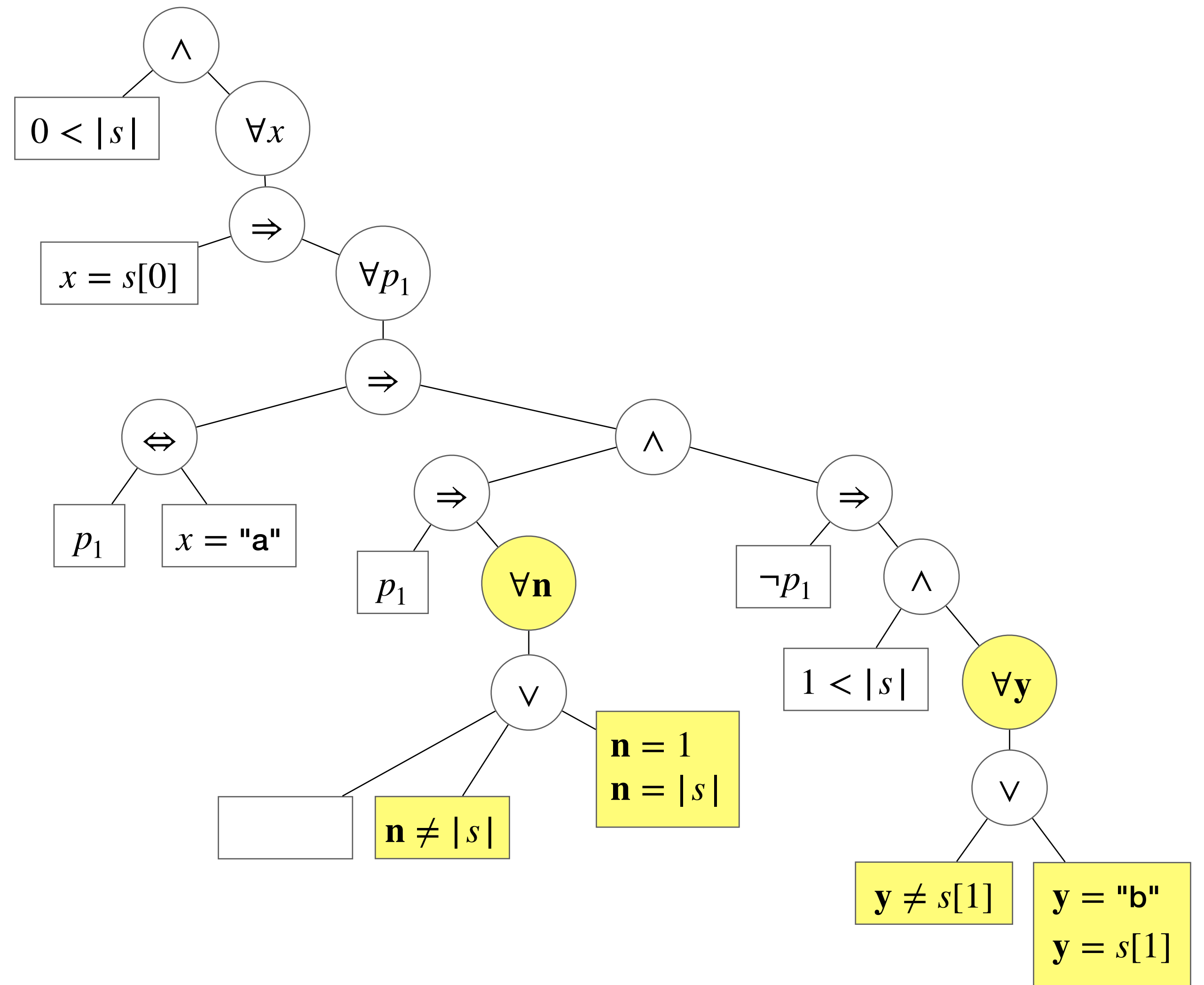
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



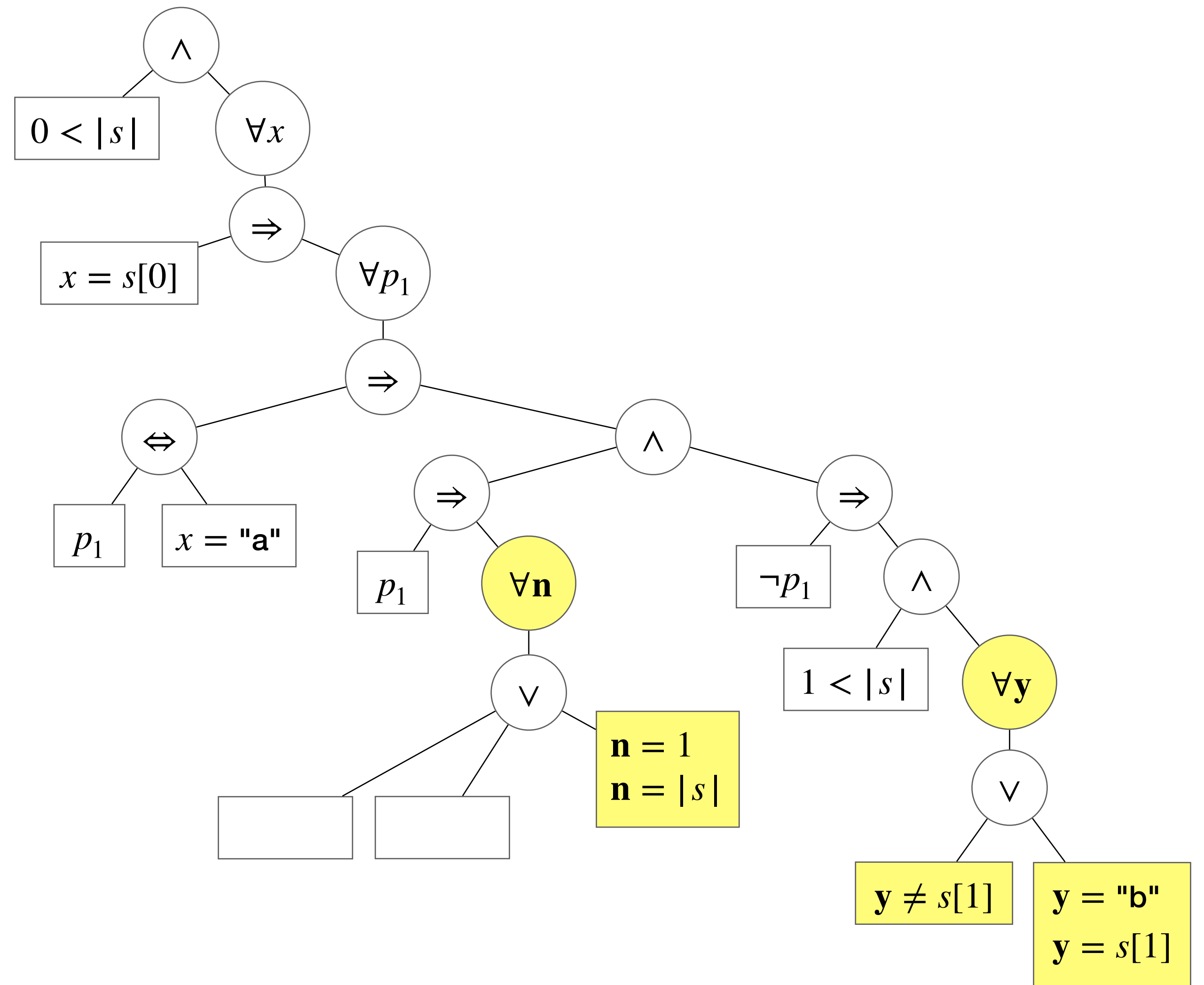
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



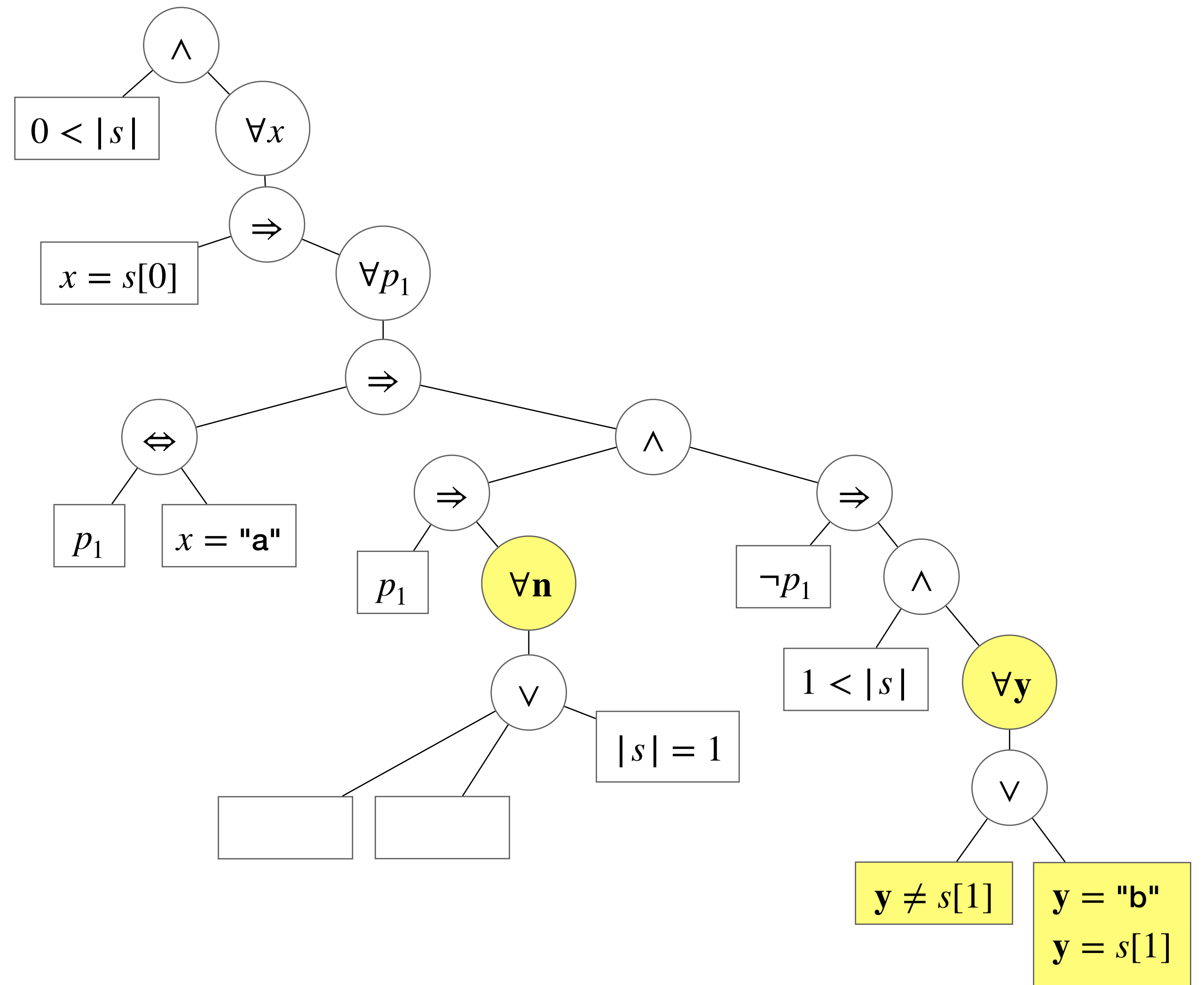
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



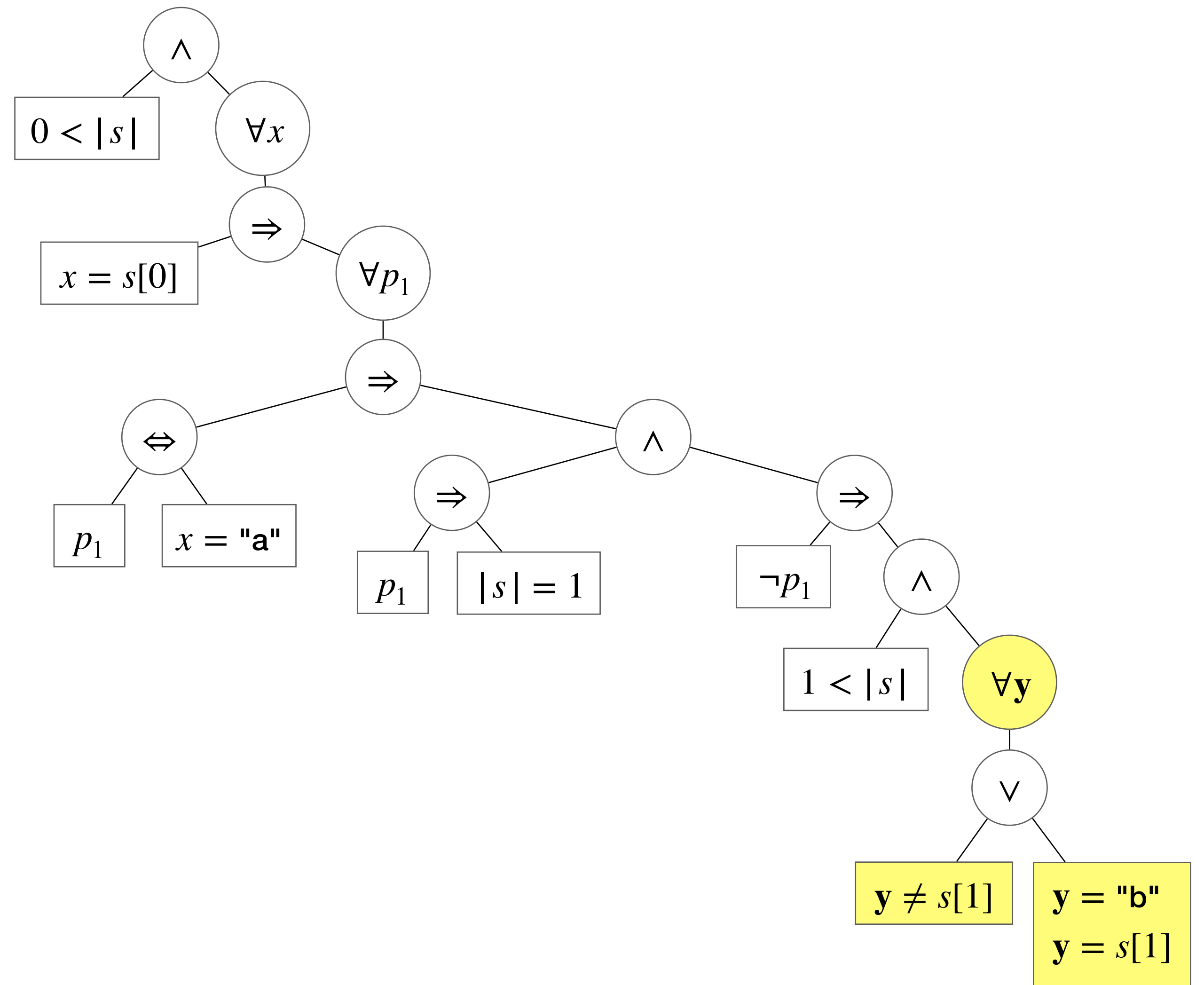
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



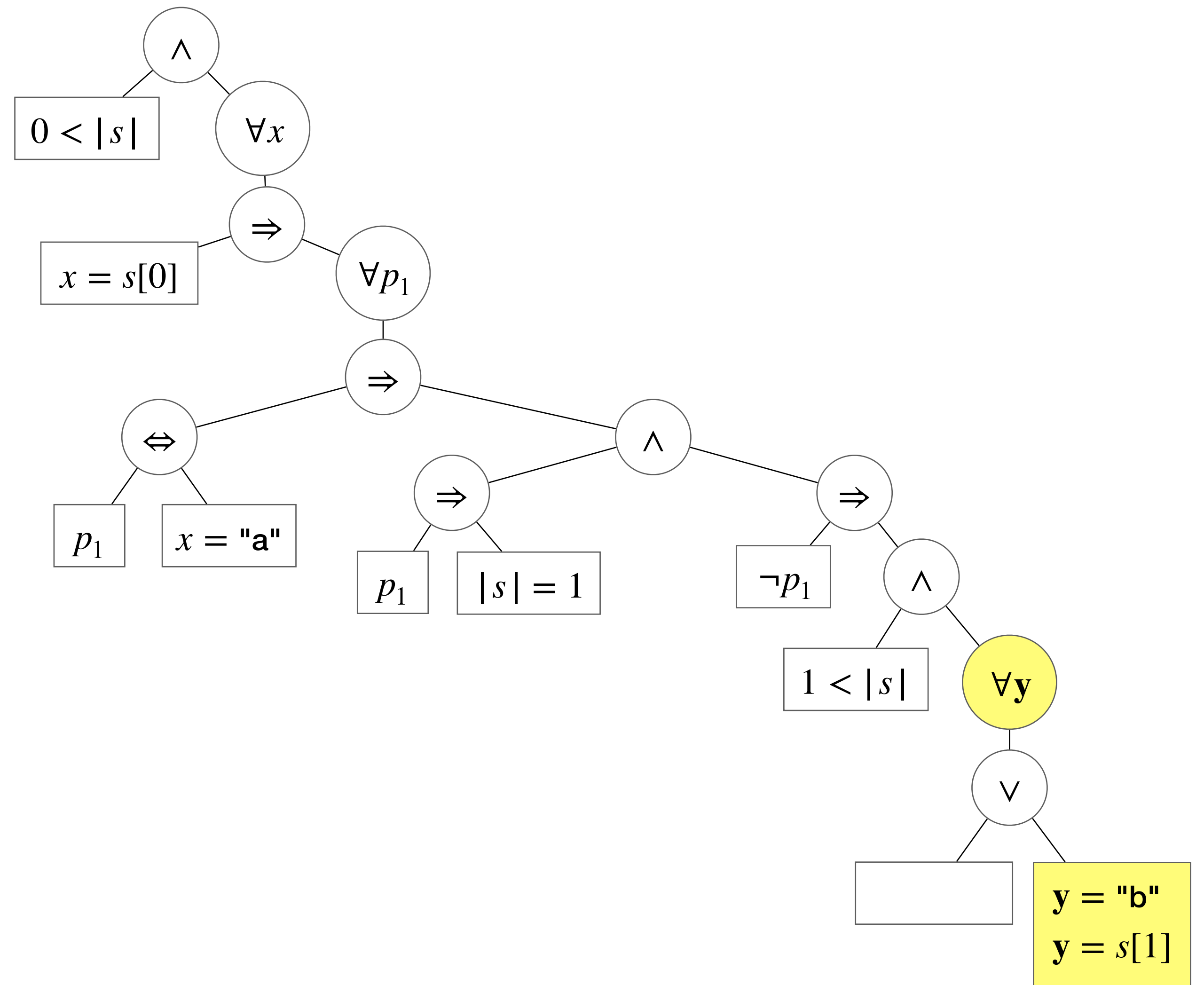
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



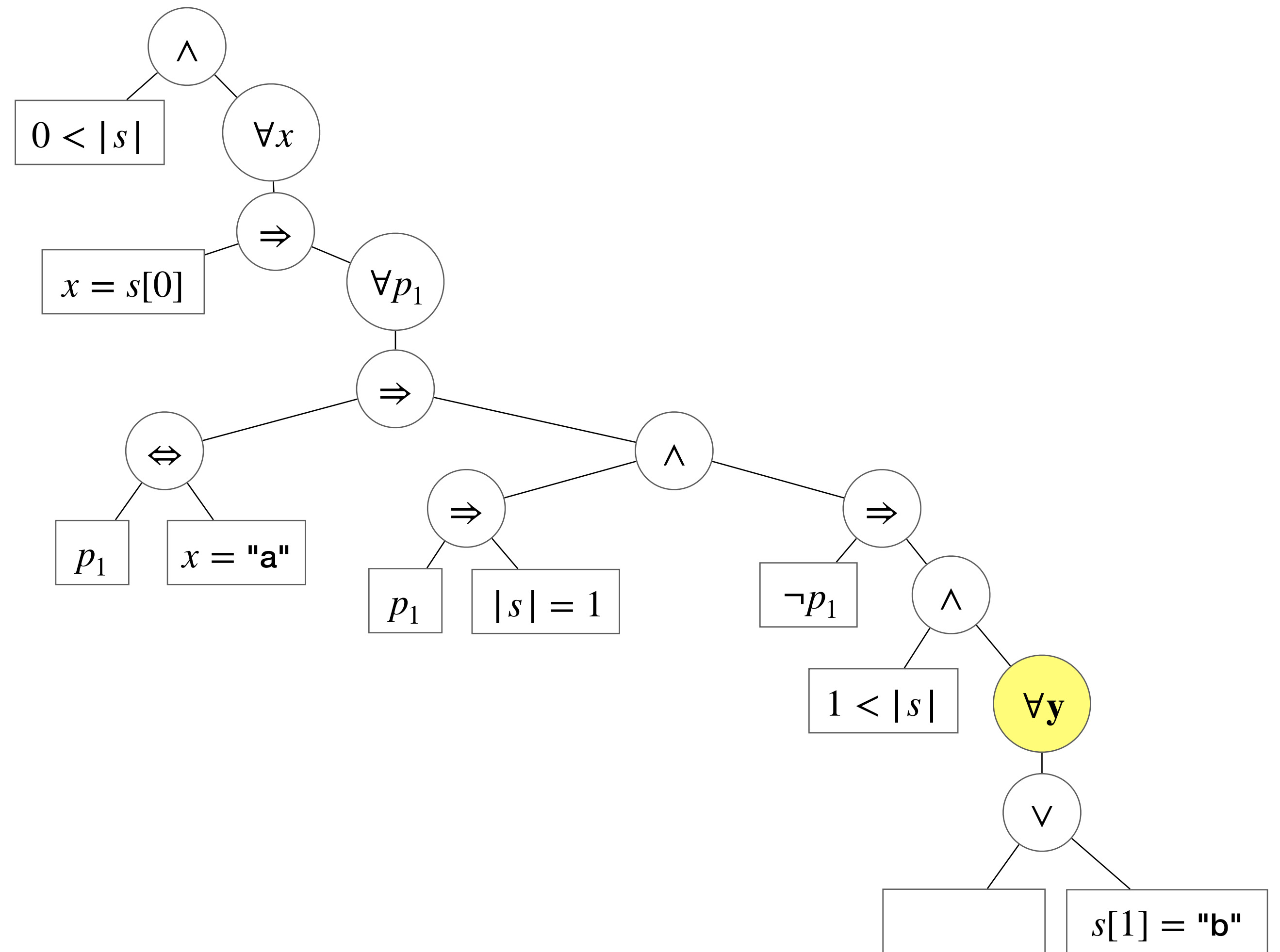
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



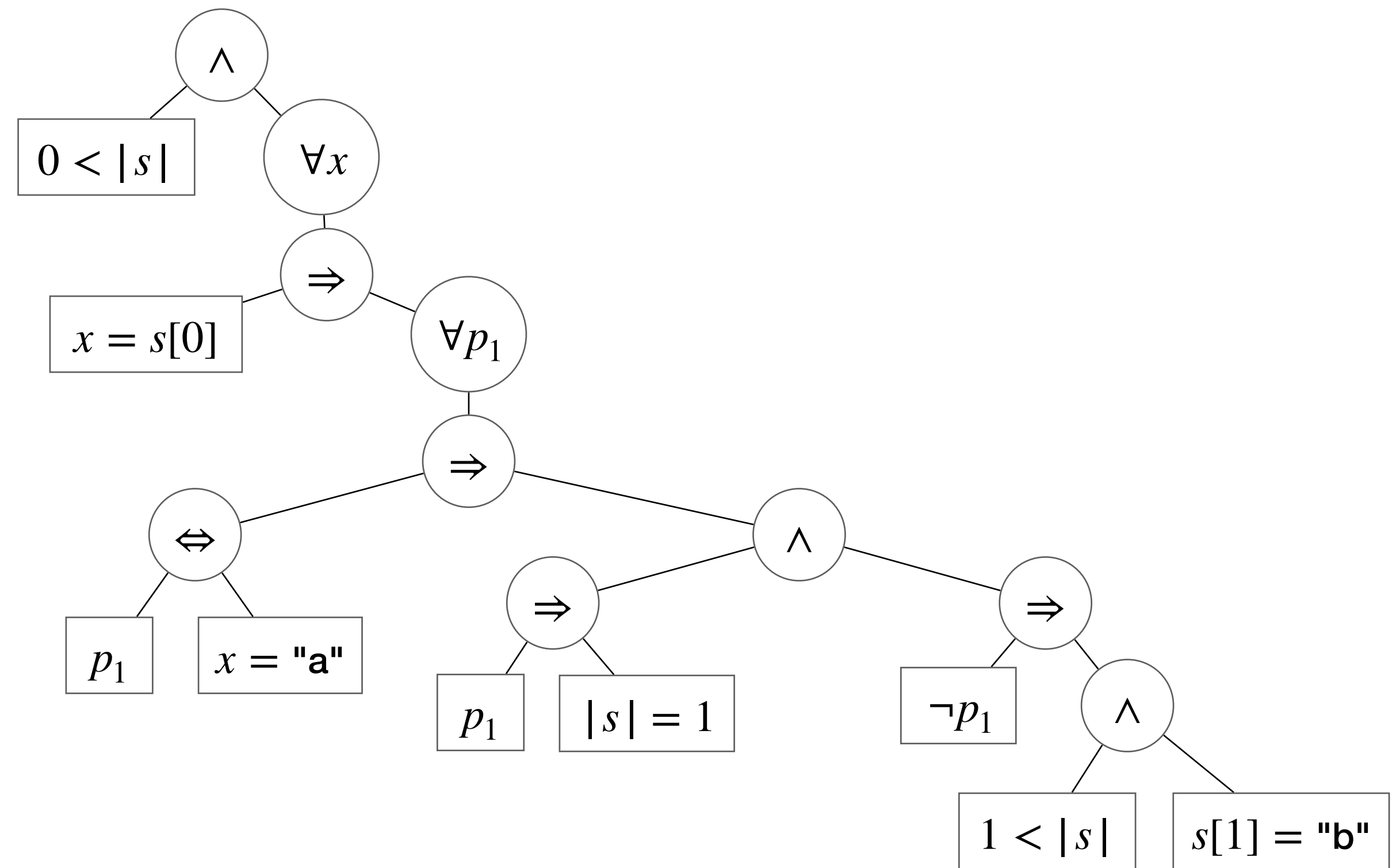
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



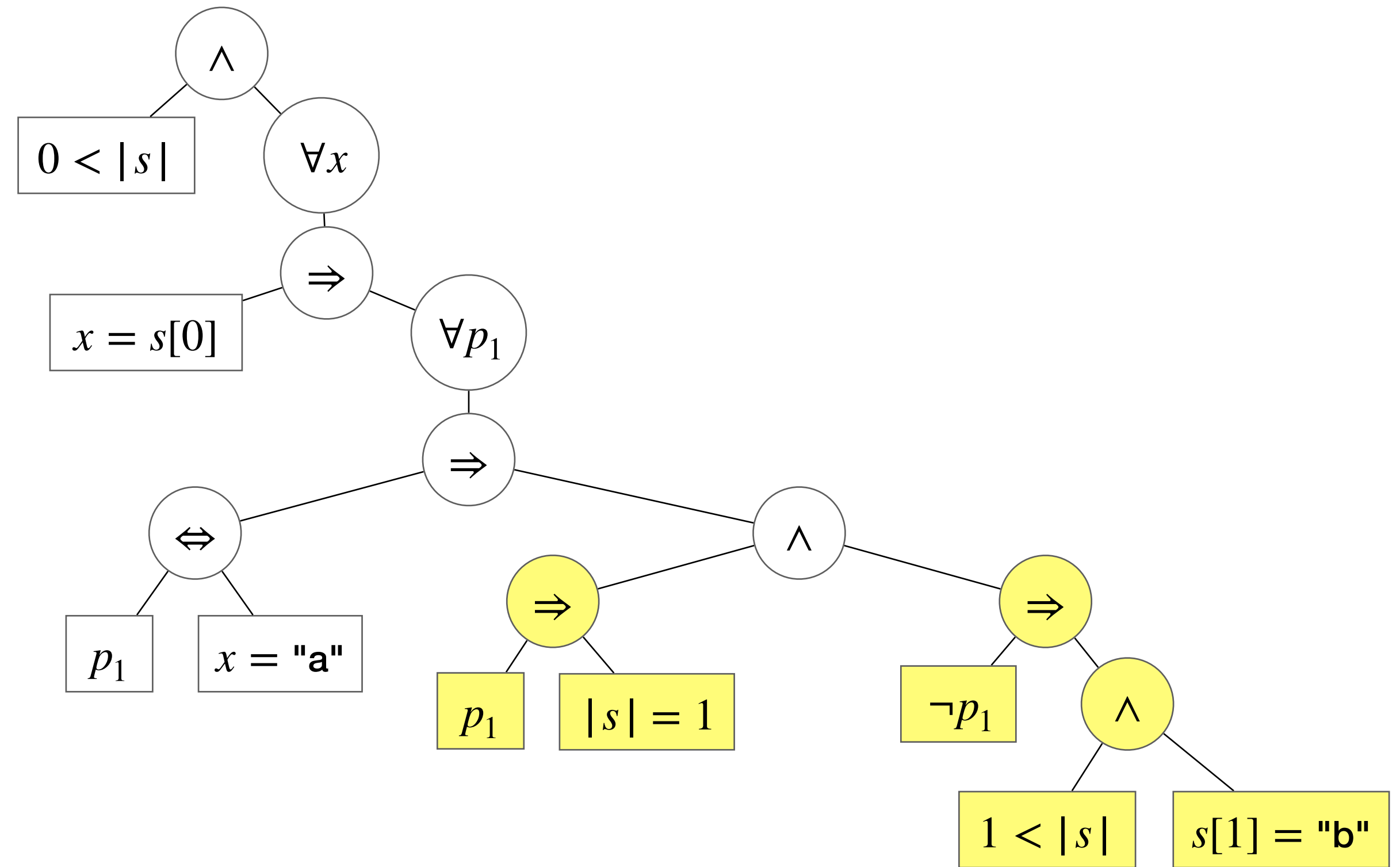
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



Grammar Solving

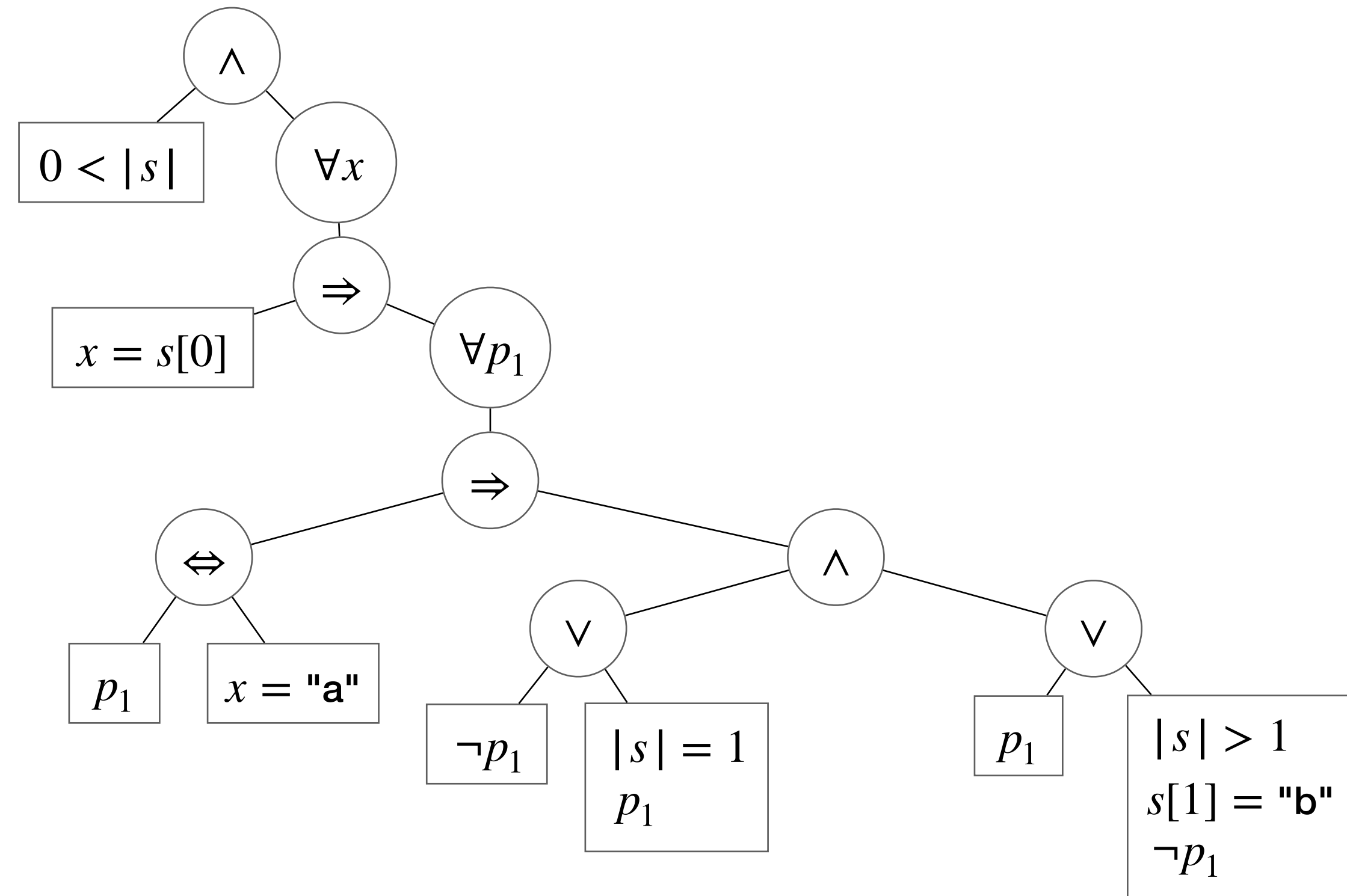
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

Grammar Solving

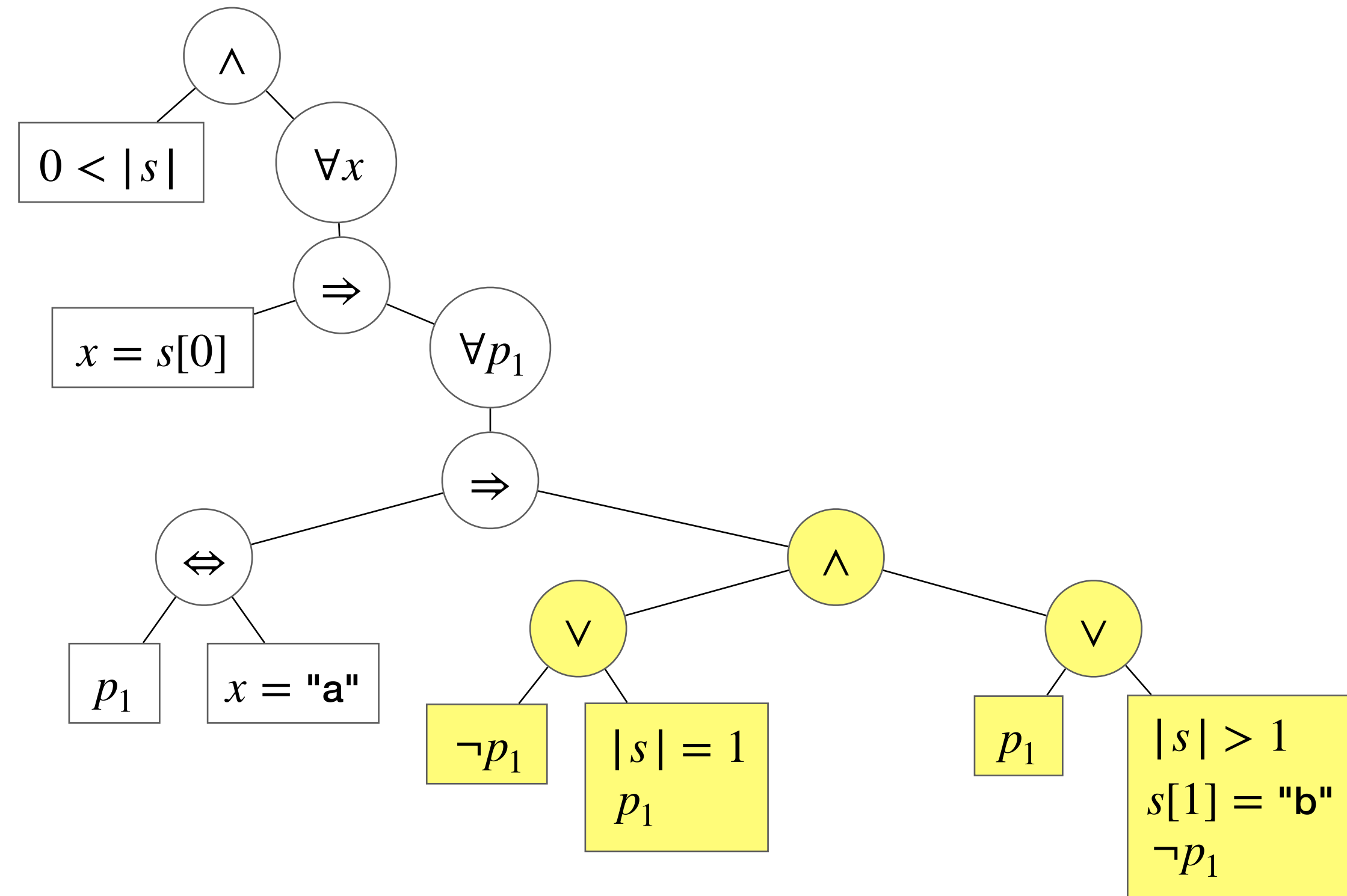
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

Grammar Solving

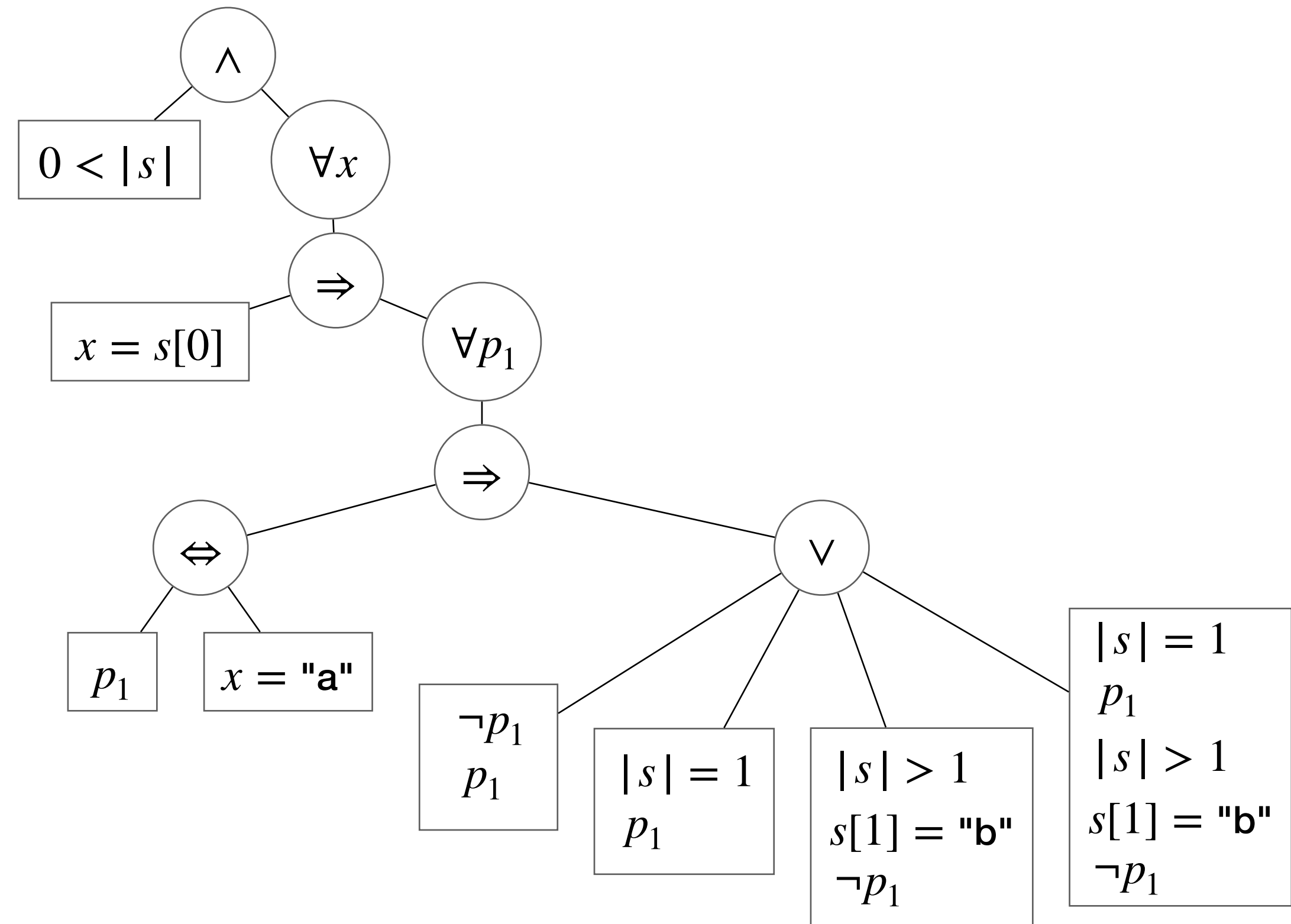
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$(a \vee b) \wedge (\neg a \vee c) \equiv b \vee c$$

Grammar Solving

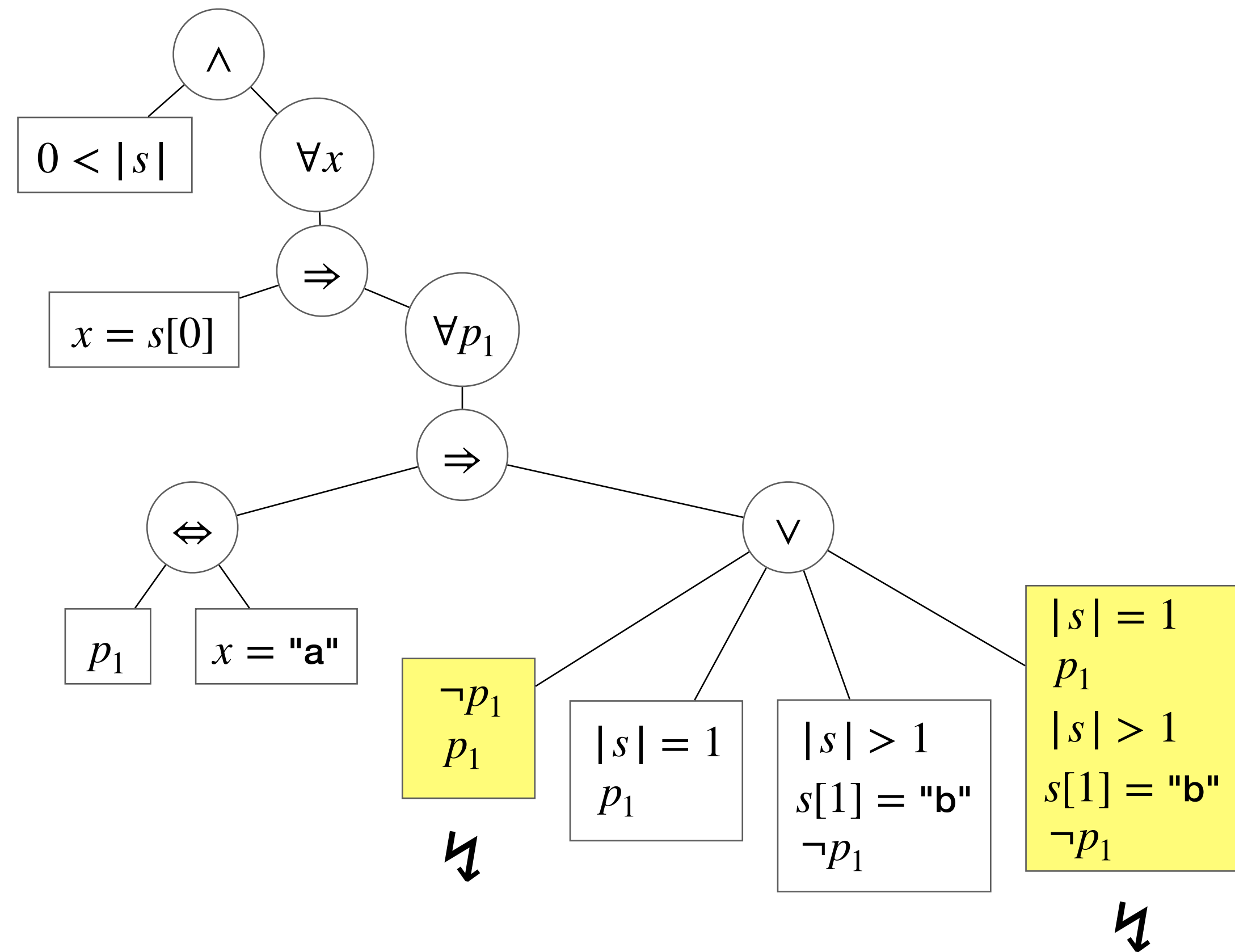
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$(a \vee b) \wedge (\neg a \vee c) \equiv b \vee c$$

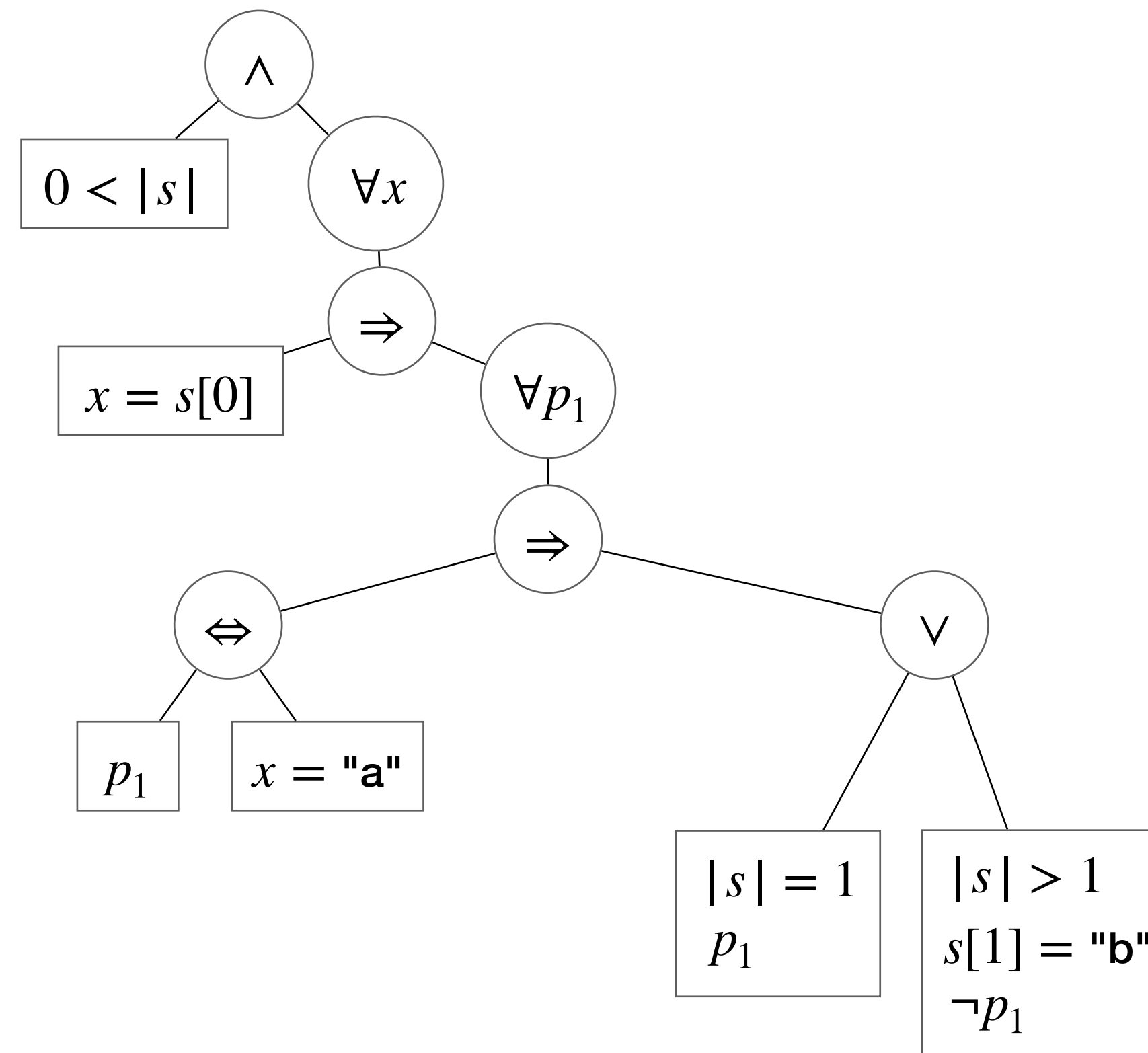
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



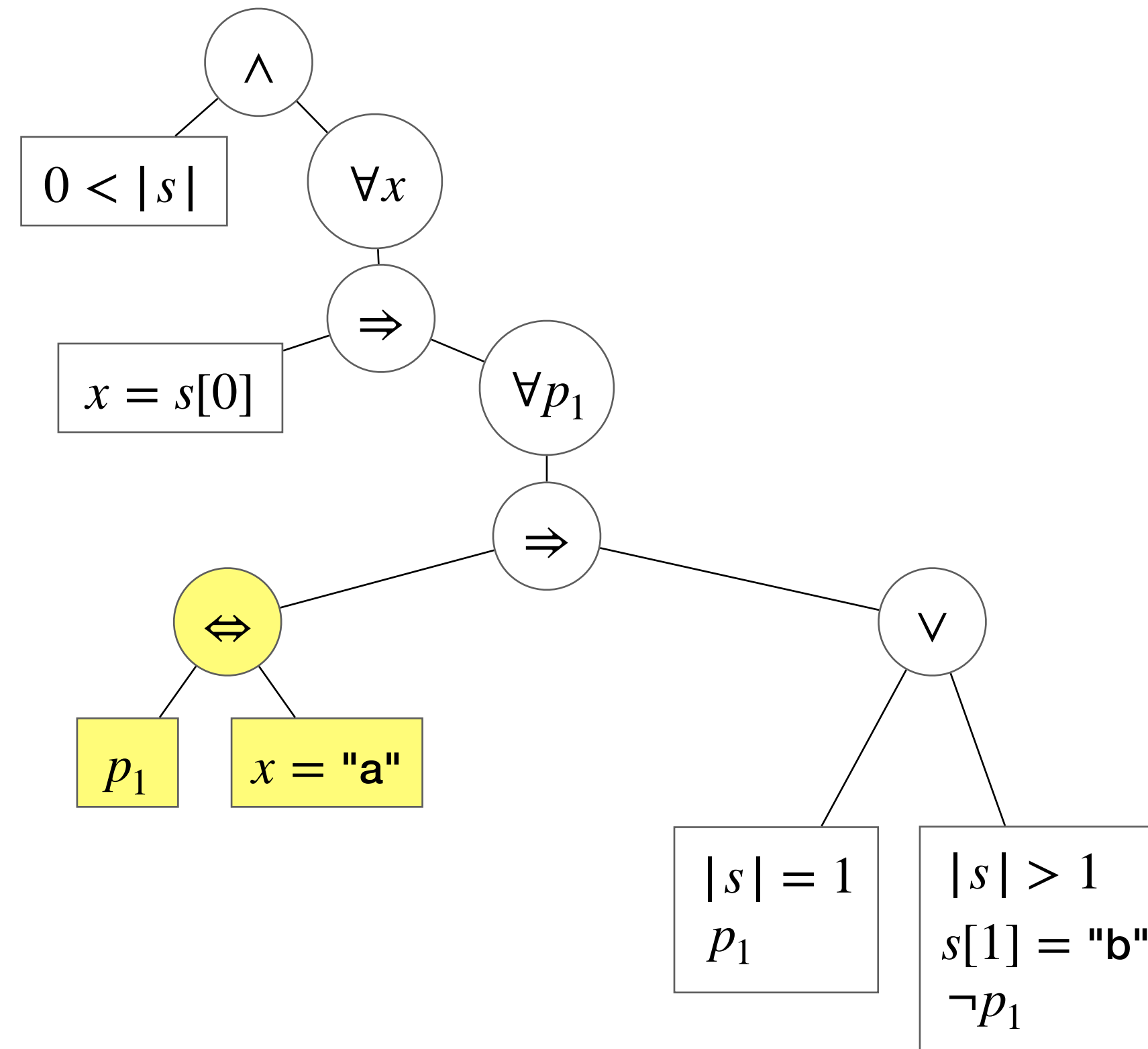
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



Grammar Solving

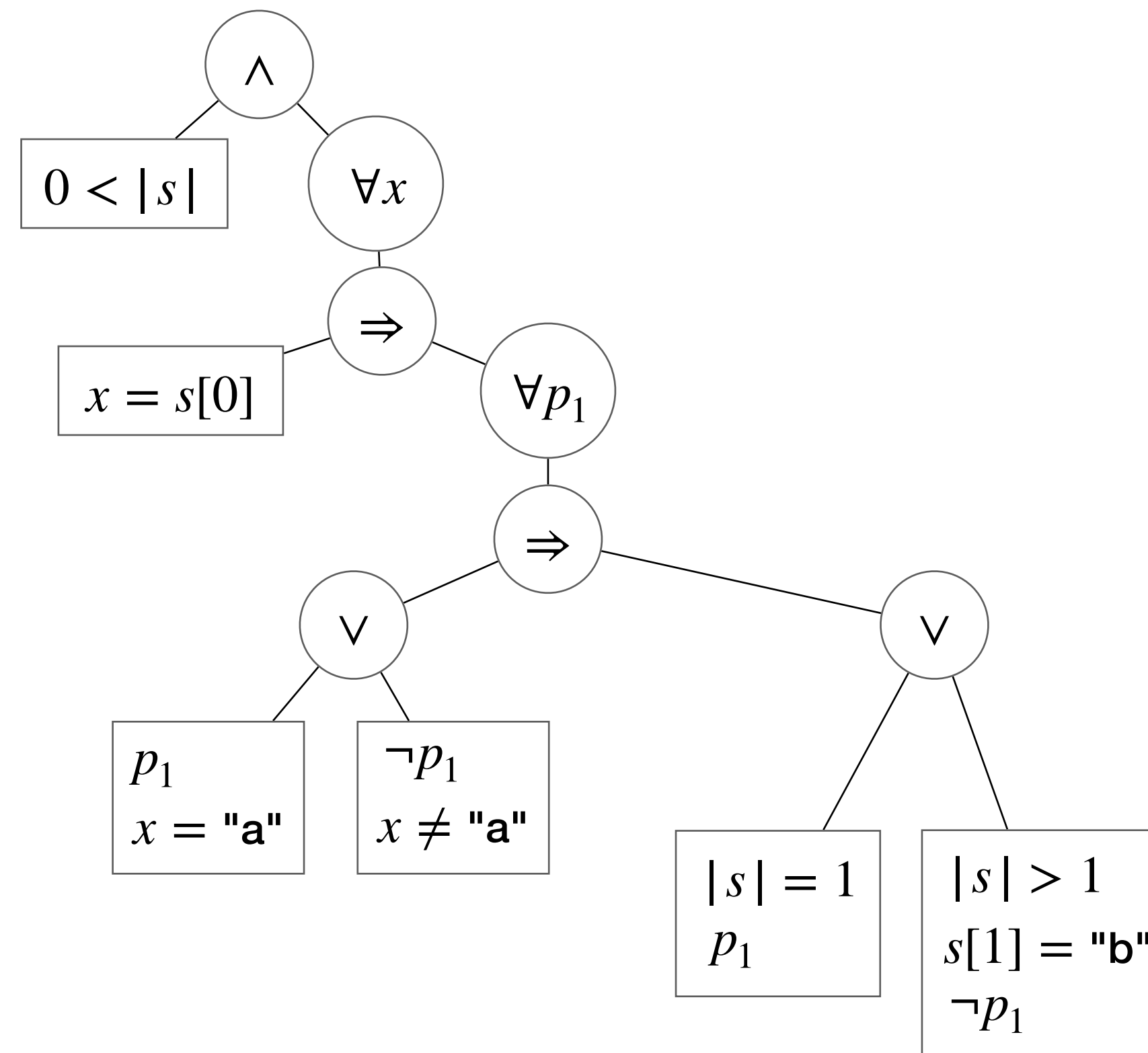
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Leftrightarrow b \equiv (\neg a \wedge \neg b) \vee (a \wedge b)$$

Grammar Solving

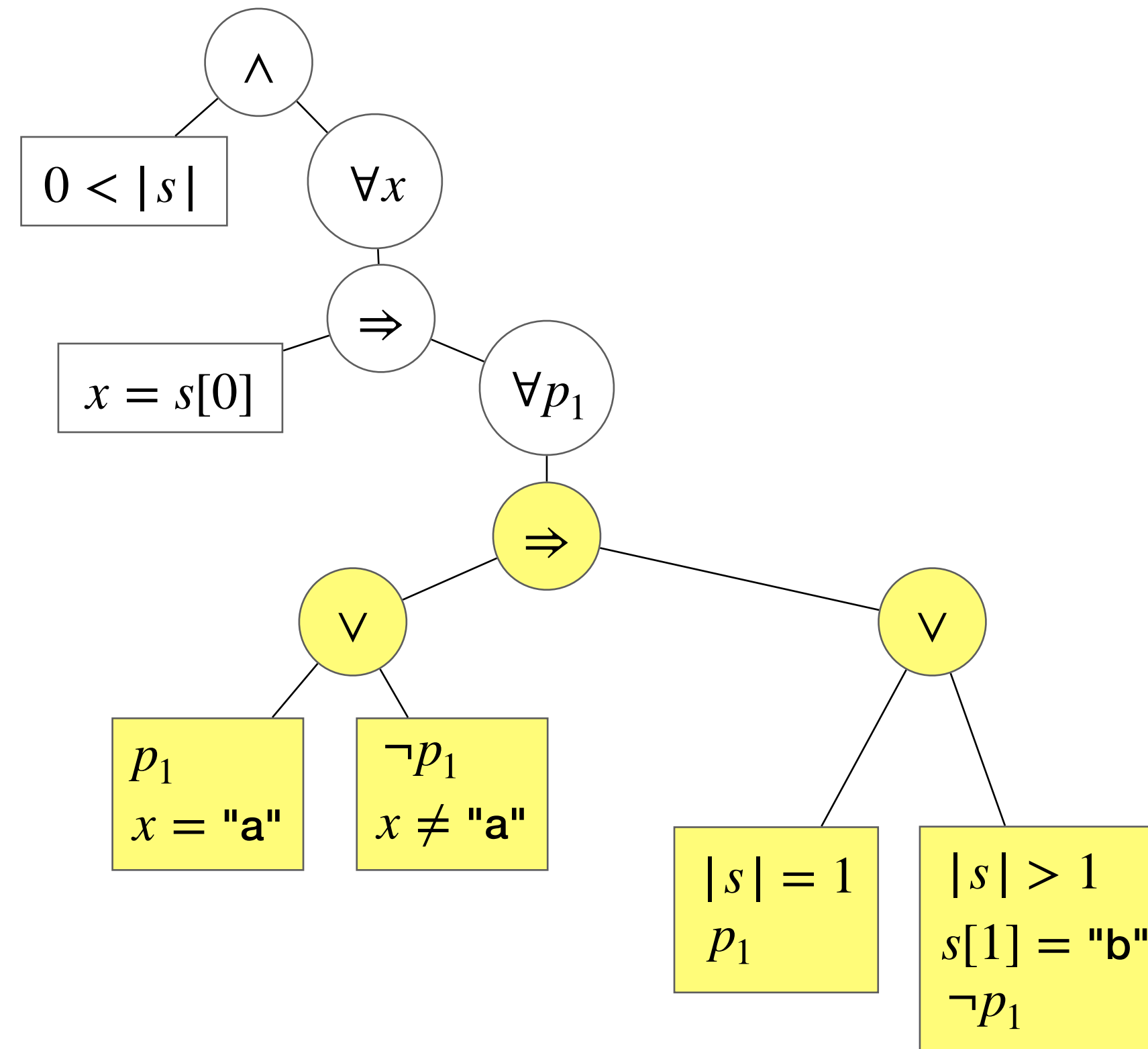
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Leftrightarrow b \equiv (\neg a \wedge \neg b) \vee (a \wedge b)$$

Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations

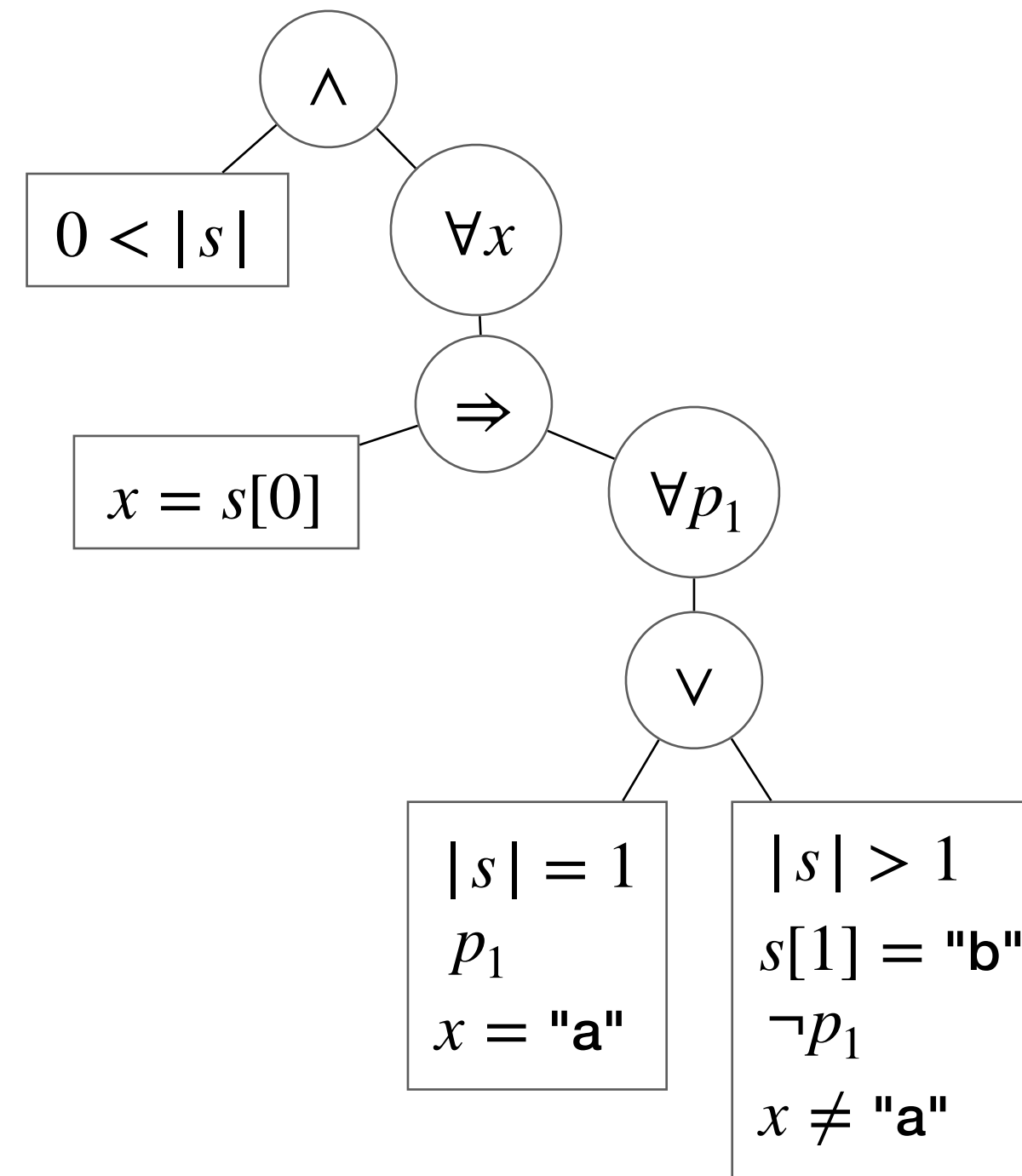


$$(a \vee b) \Rightarrow c \equiv (a \Rightarrow c) \vee (b \Rightarrow c)$$

$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations

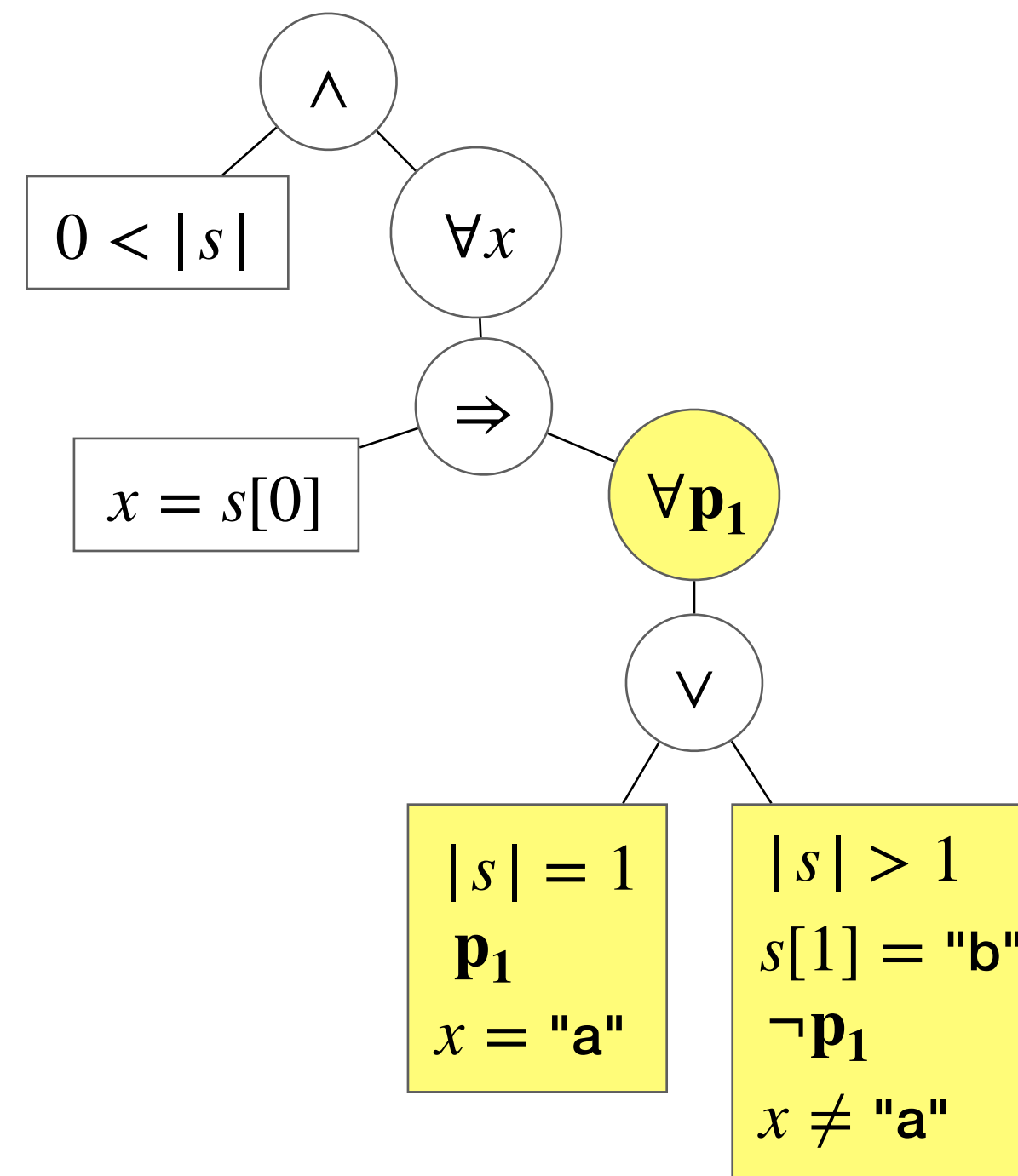


$$(a \vee b) \Rightarrow c \equiv (a \Rightarrow c) \vee (b \Rightarrow c)$$

$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

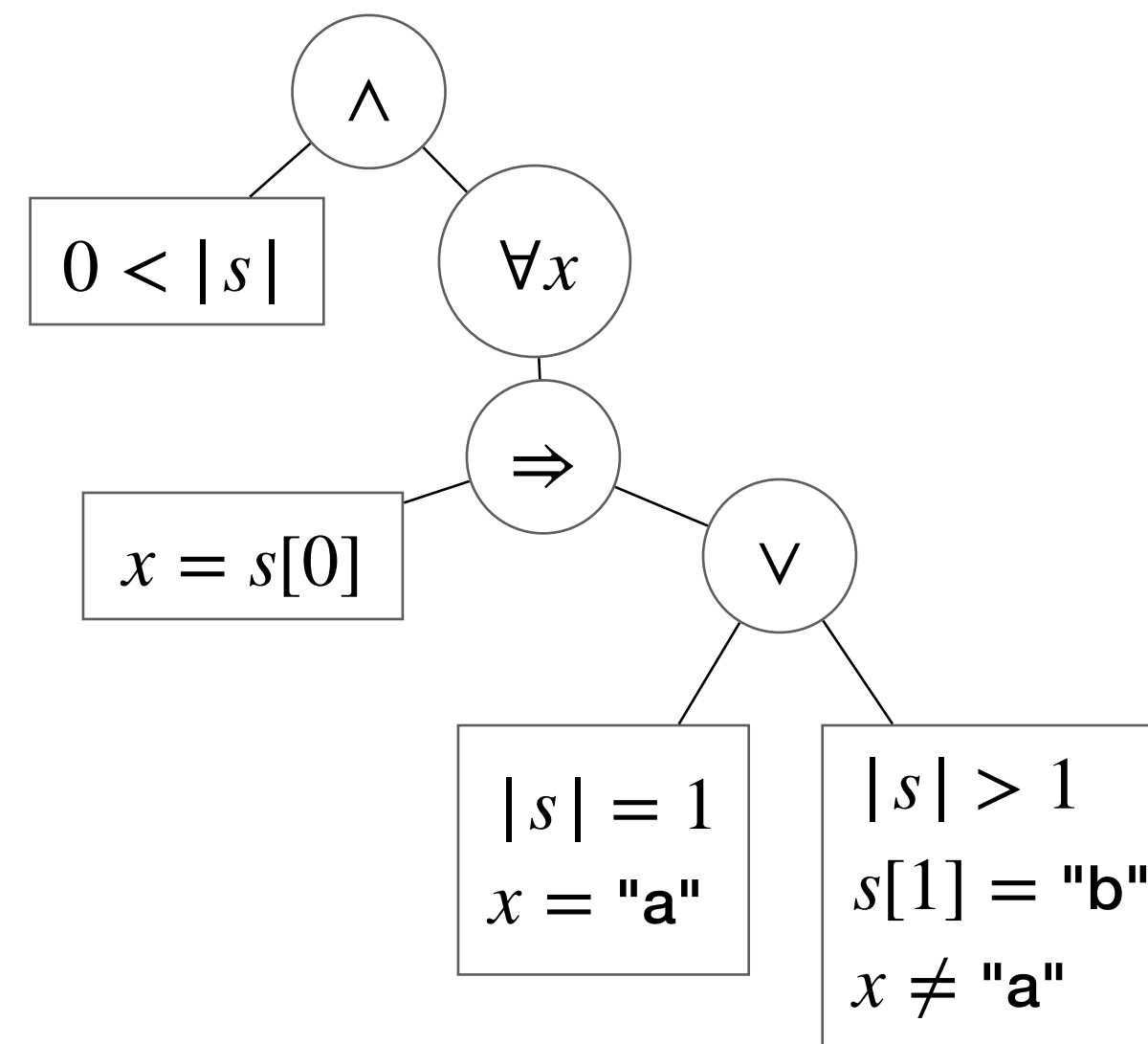
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



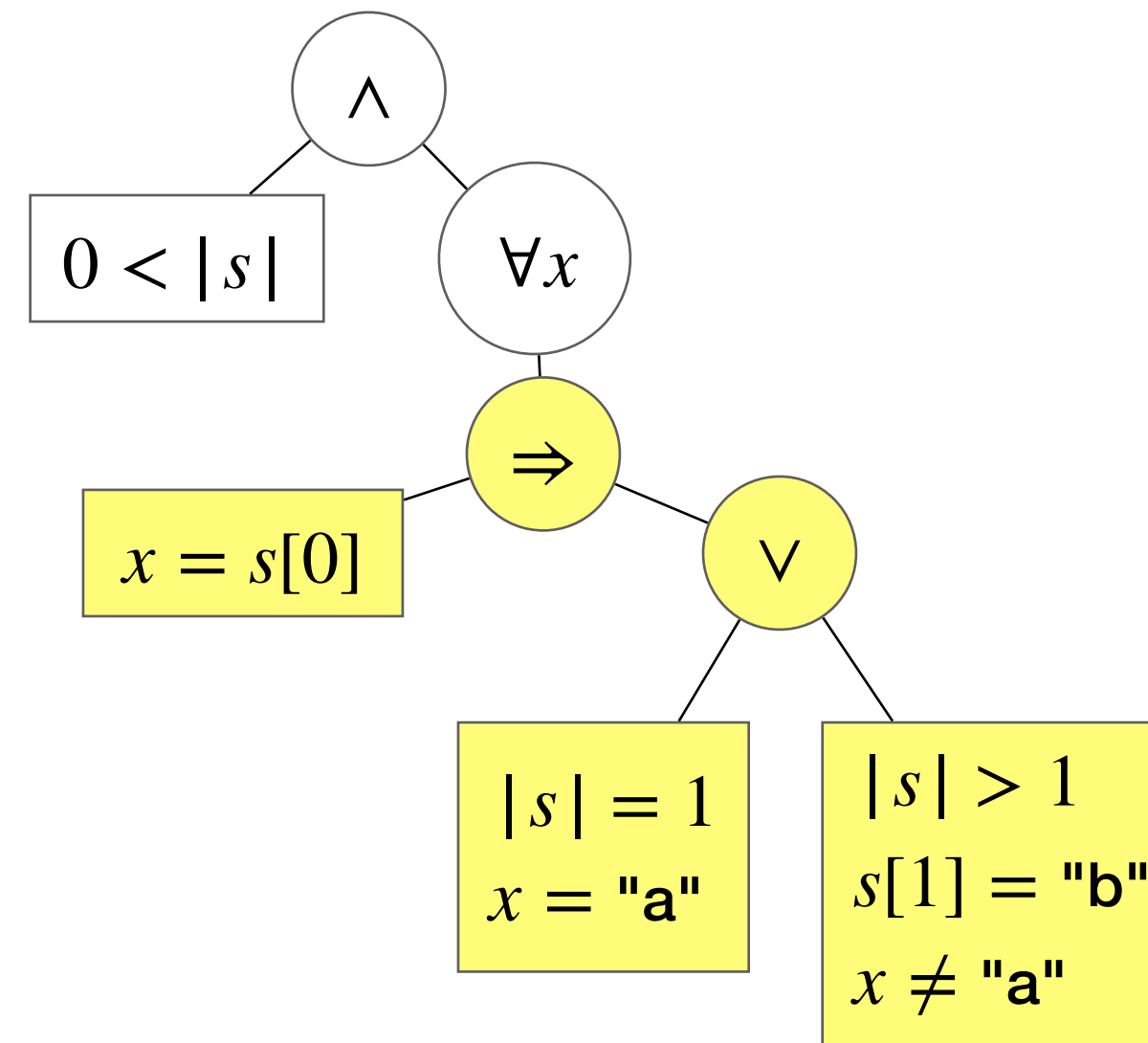
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



Grammar Solving

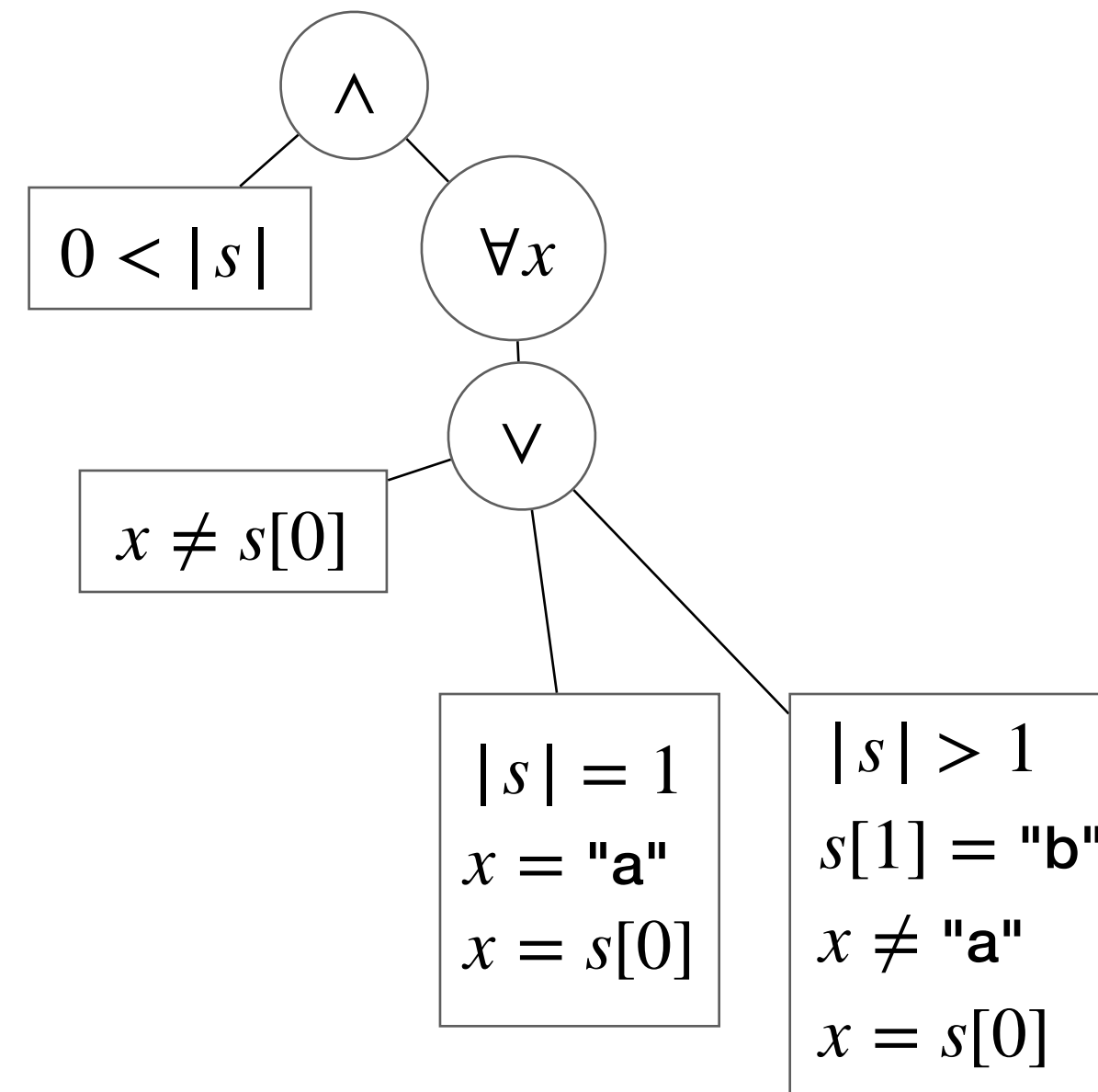
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

Grammar Solving

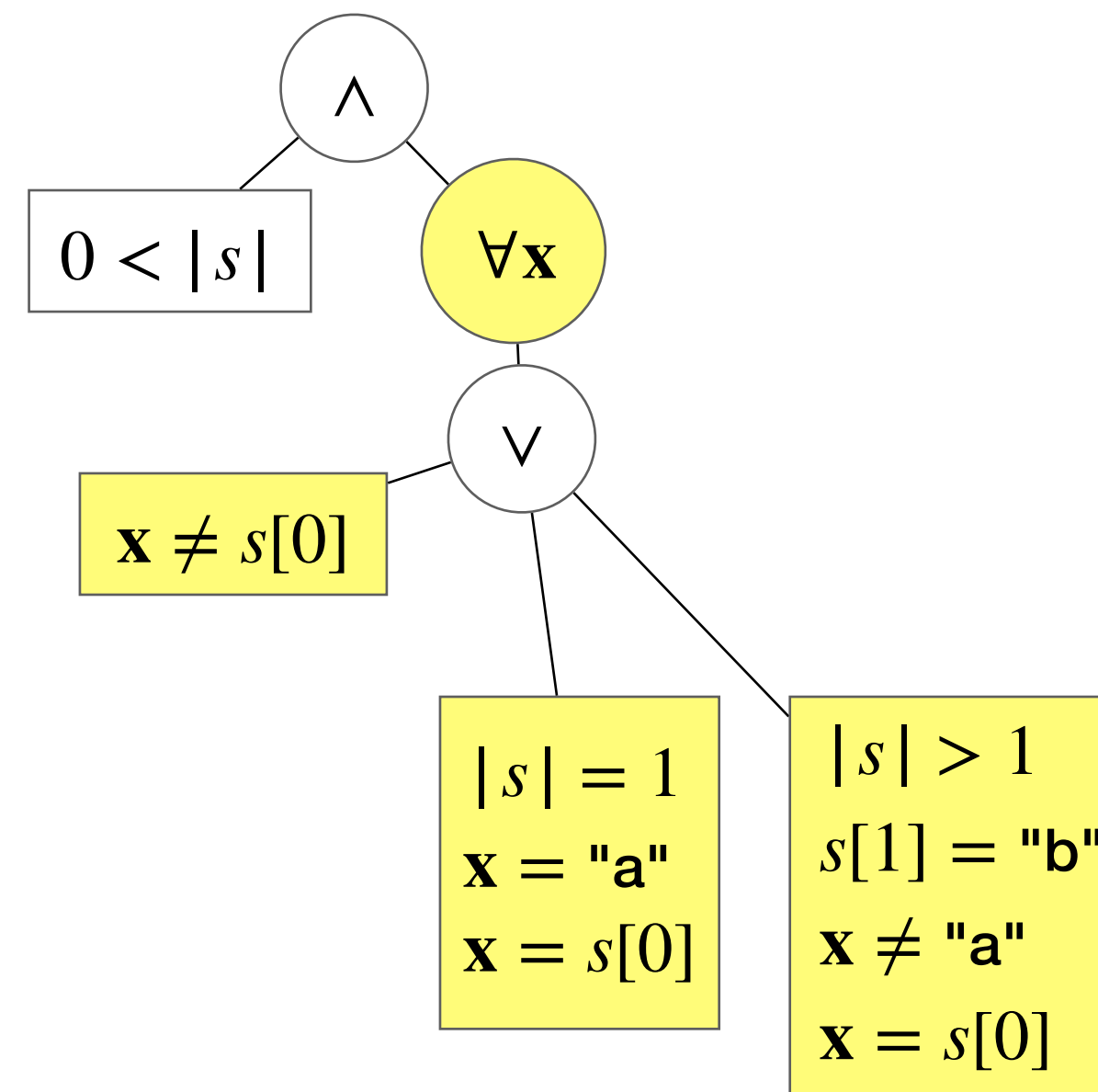
- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



$$a \Rightarrow b \equiv \neg a \vee (a \sqcap b)$$

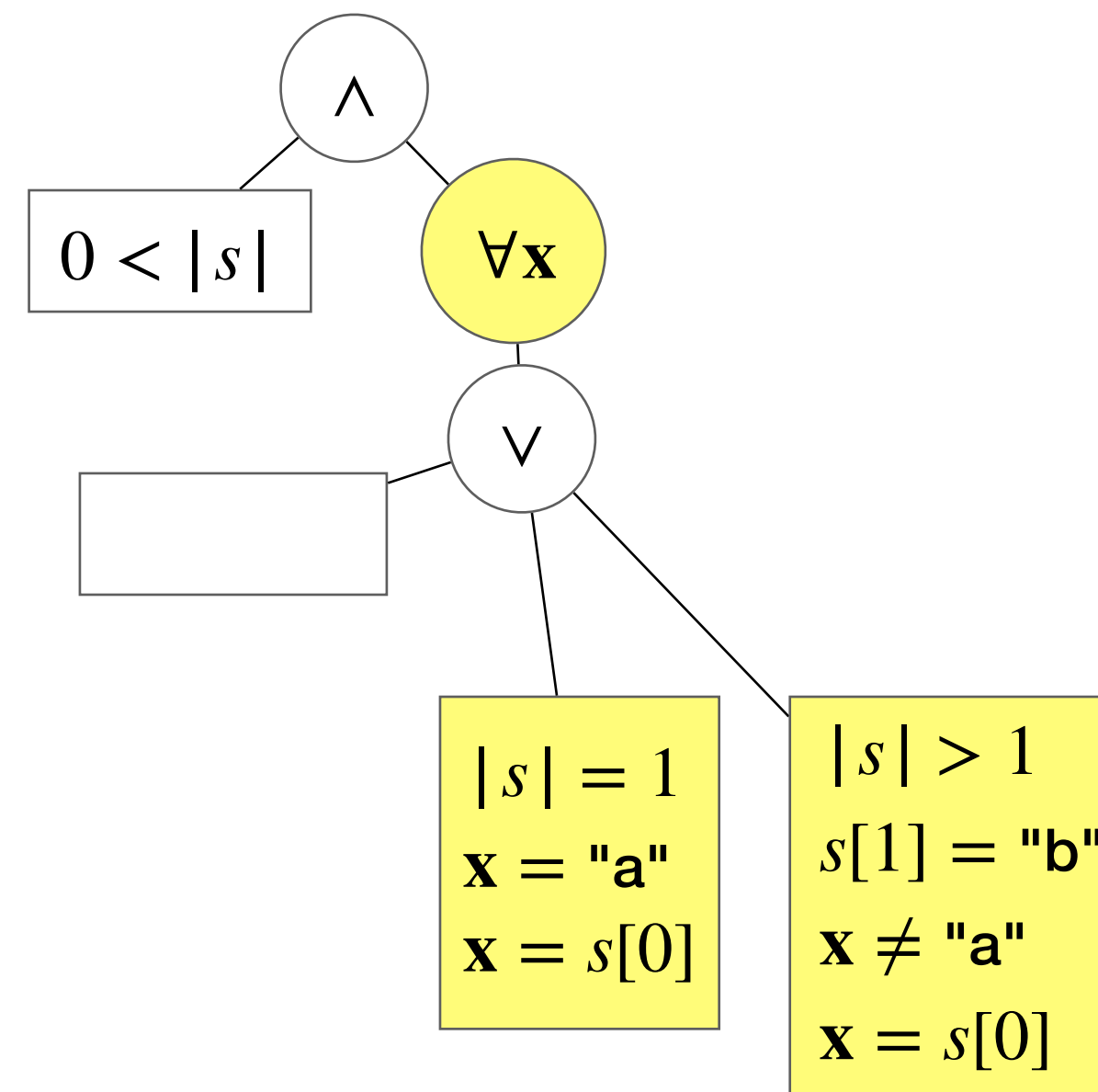
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



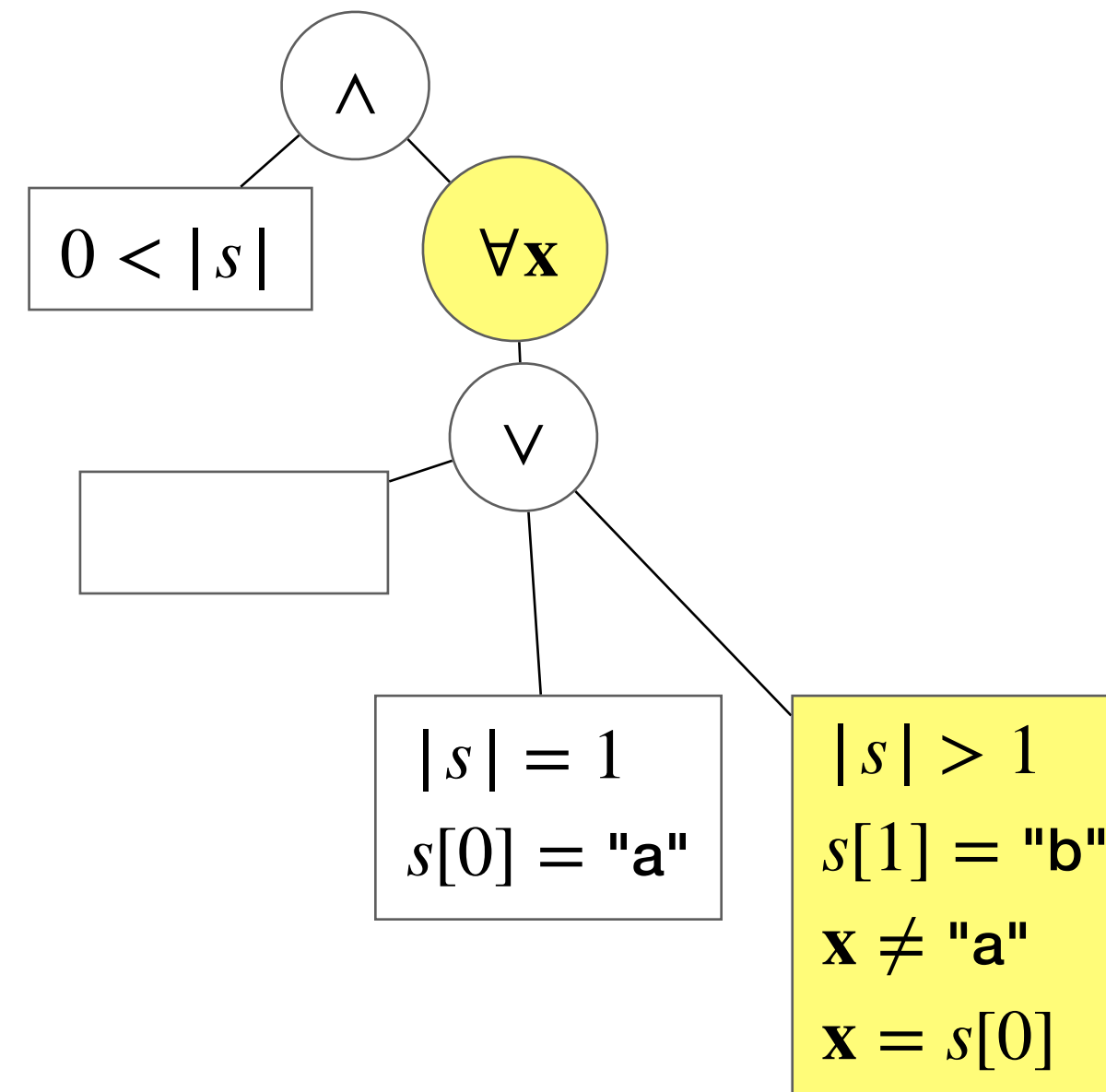
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



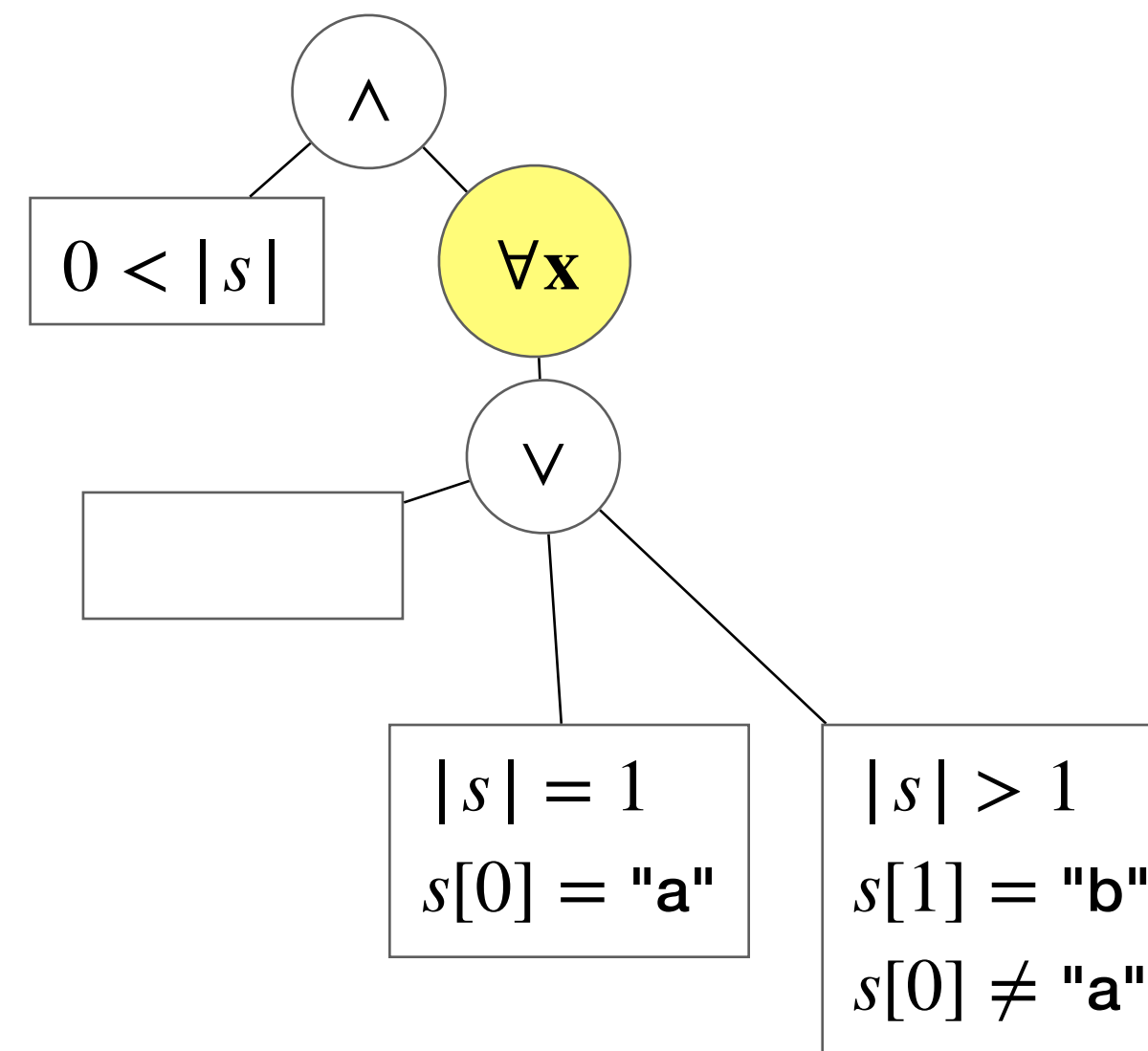
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



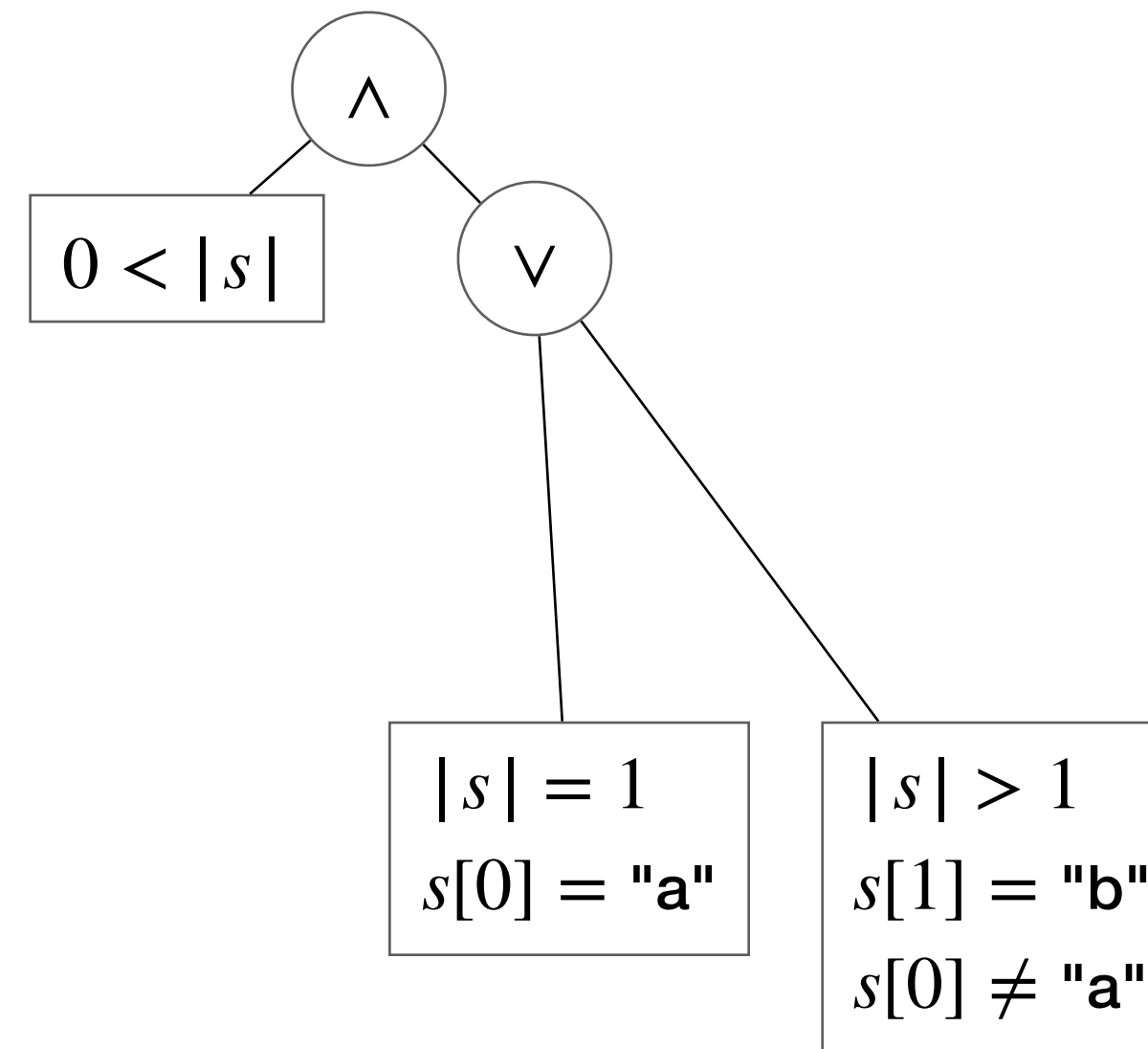
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



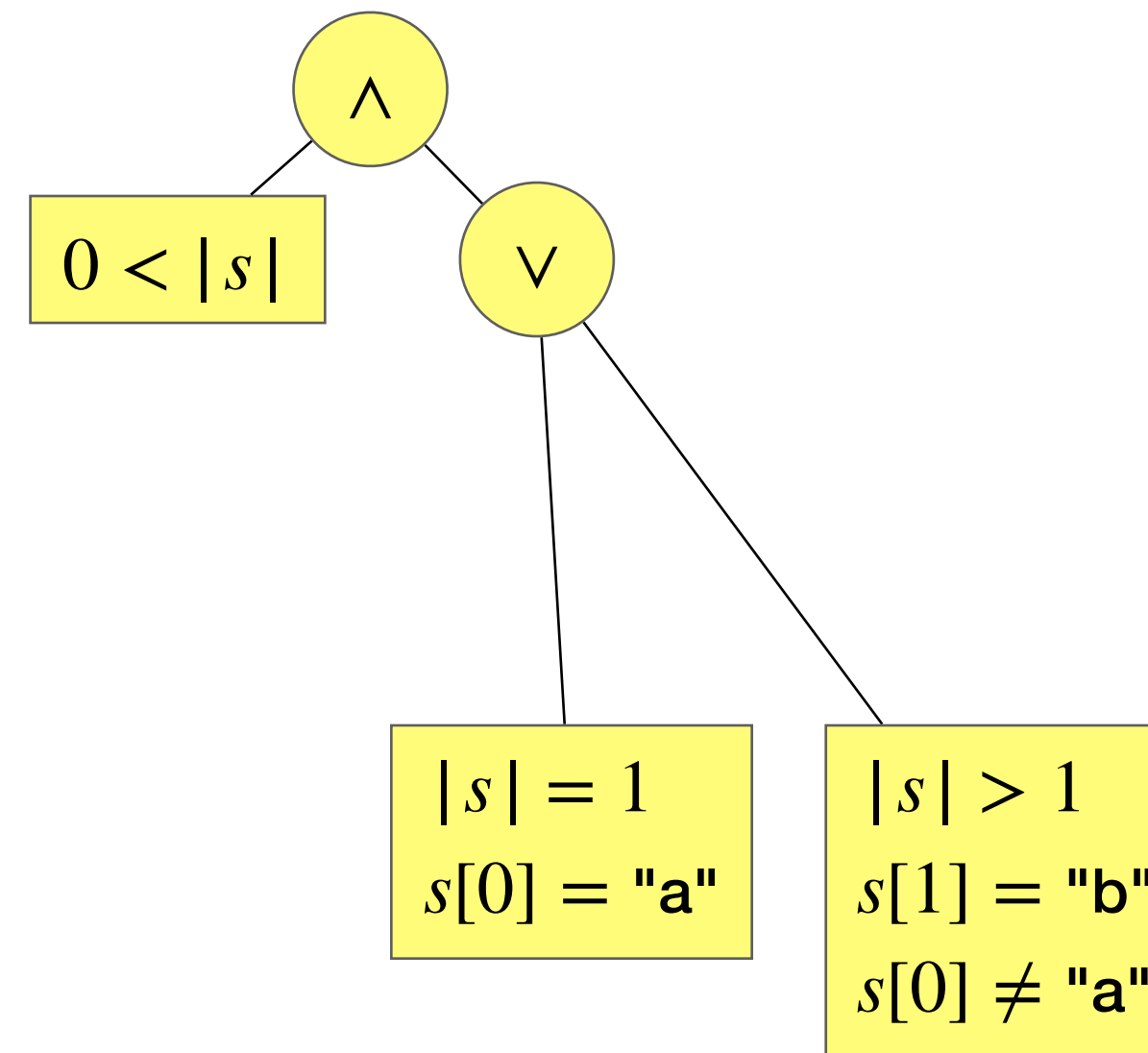
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



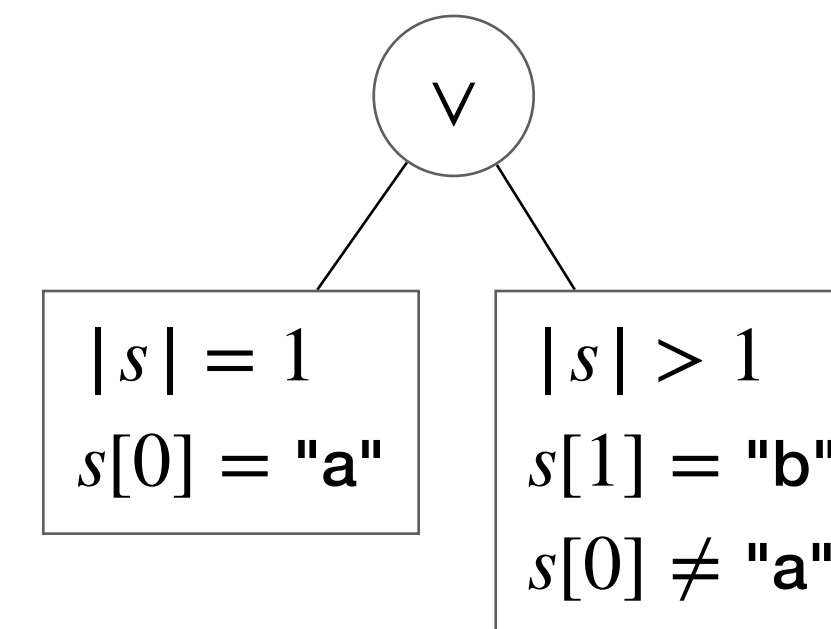
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



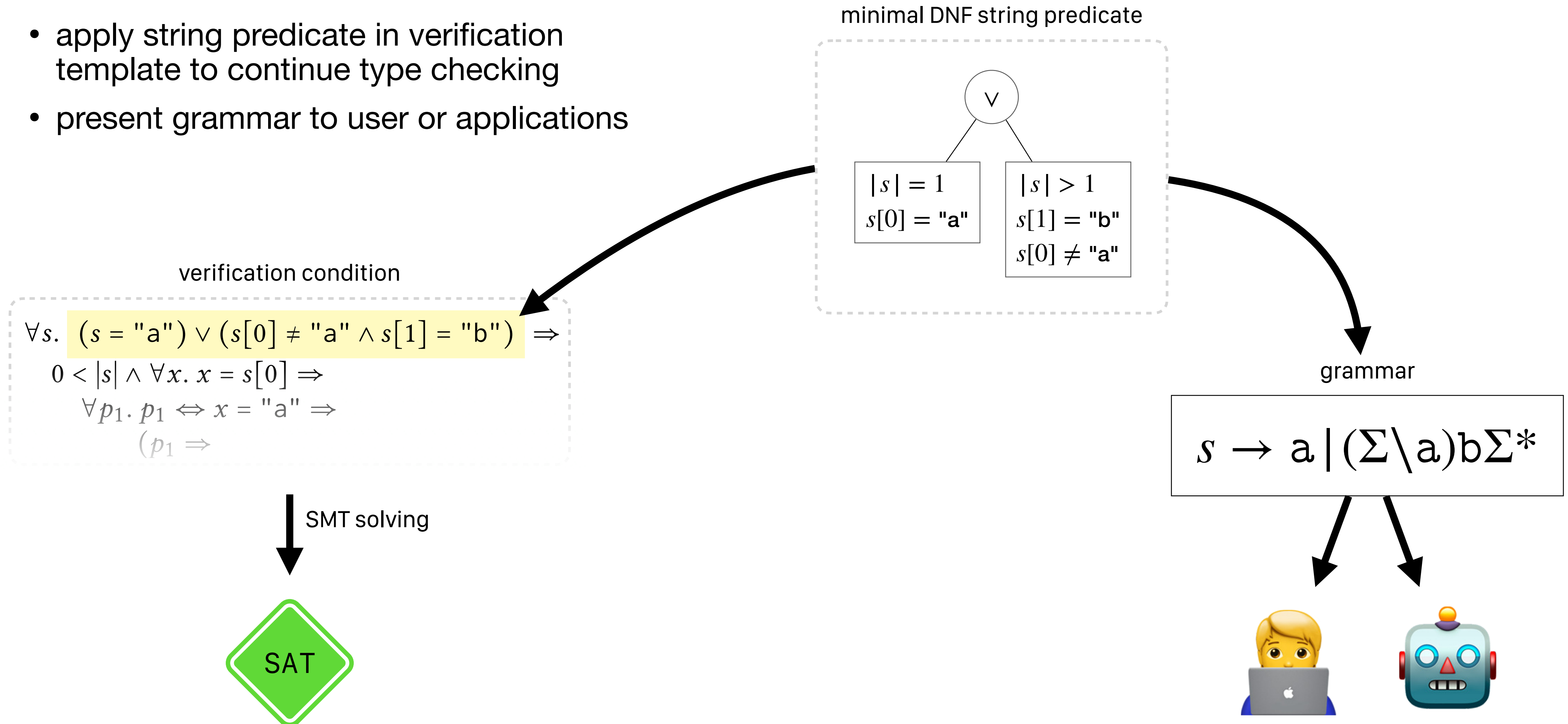
Grammar Solving

- base solution on “grammar consequent”
- minimize via bottom-up tree rewriting
- apply Boolean equivalences to reach DNF
- eliminate quantifiers by resolving equations
- use (precise) abstract value representations



Enjoy your grammar!

- apply string predicate in verification template to continue type checking
- present grammar to user or applications



Status / Future Work

✓ PANINI proof-of-concept

➡
SOON formalization of grammar solving algorithm

🕒 end-to-end grammar inference system

- library of string operation specifications
- evaluation on corpus of curated ad hoc parser samples
- large-scale mining study of grammars in the wild
- application prototypes (build bot, IDE plugin,...)

🕒 user study on grammar comprehension

Toward Grammar Inference via Refinement Types

<https://mcschroeder.github.io/#tyde2022>



Michael Schröder

TU Wien

Vienna, Austria

michael.schroeder@tuwien.ac.at



Jürgen Cito

TU Wien

Vienna, Austria

juergen.cito@tuwien.ac.at

Type-Driven Development (TyDe), ICFP 2022
Ljubljana, Slovenia

TU
WIEN Informatics