


# Automated Ambiguity Detection in Layout-Sensitive Grammars

Jiangyi LIU<sup>1</sup> 

Fengmin (Paul) ZHU<sup>1,2</sup> 

Fei HE<sup>1</sup> 

<sup>1</sup>School of Software, Tsinghua University

<sup>2</sup>CISPA Helmholtz Center for Information Security



## Layout-Sensitive Languages

Whitespaces and indentations affect how programs get parsed.



```
def greeting():
    print('Welcome to our poster' +
          ' and our talk at 15:12 Friday!')
```

$$\text{indent}(w_1, w_2) \triangleq (w_1 \neq \varepsilon \wedge w_2 \neq \varepsilon) \Rightarrow w_2[0].\text{col} > w_1[0].\text{col} \wedge w_2[0].\text{line} = w_1[-1].\text{line} + 1$$


```
match solver_result with
| Sat model -> Ambig (decode model)
| Unsat -> Unambig
| _ -> Unknown
```

$$\text{align}(w_1, w_2) \triangleq (w_1 \neq \varepsilon \wedge w_2 \neq \varepsilon) \Rightarrow w_1[0].\text{col} = w_2[0].\text{col}$$


```
do grammar <- load "my_lang.ebnf"
sentence <- checkAmbig grammar
return $ AmbigResult sentence $
trees grammar sentence
```

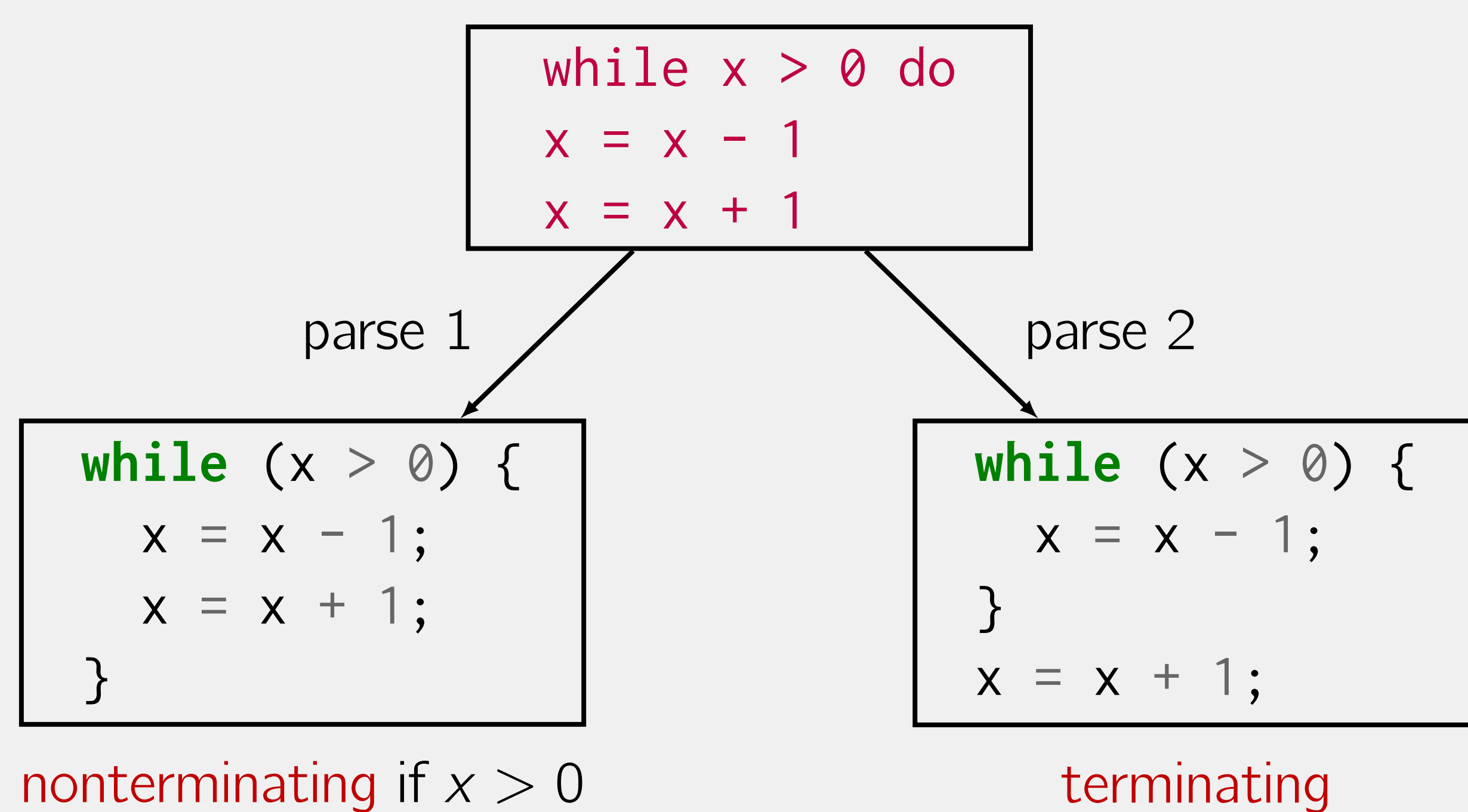
$$\text{offside}(w) \triangleq w \neq \varepsilon \Rightarrow \forall t \in w : t.\text{line} > w[0].\text{line} \Rightarrow t.\text{col} > w[0].\text{col}$$

## Ambiguity Matters

Consider a grammar fragment:

```
block → stmt*
stmt → var = expr | while expr do block | ...
```

This sentence has two different parses that are semantically different:



## A Tour of Lamb

**Step 1:** Input a grammar  $G_{\text{block}}$ :

```
block → ||stmt||+
stmt → nop | do block
```

The alignment constraint  $\| \cdot \|+$  marks the border of the **do**-block body, so it distinguishes between

```
do nop
nop
```

and

```
do nop
  nop
```

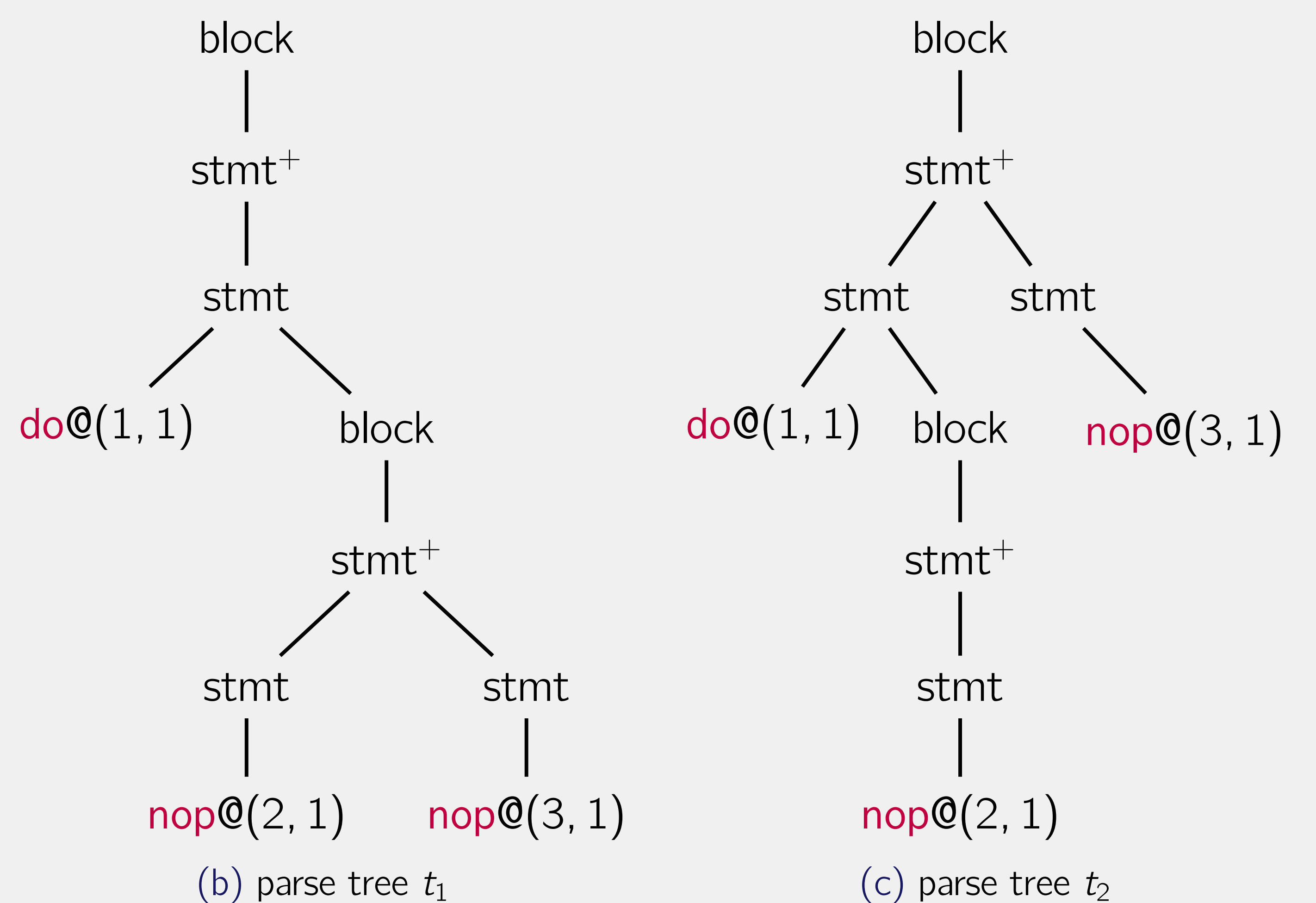
Awesome! But is  $G_{\text{block}}$  really unambiguous?

**Step 2:** Run Lamb

It finds a shortest ambiguous sentence (with its parse trees):

```
do
nop
nop
```

(a) ambiguous sentence



Users can fix the ambiguity issue manually with the aid of the produced parse trees.

**Step 3:** Understand the cause of ambiguity

It is insufficient to tell whether the second **nop** statement belongs to the **do**-block or the top-level block, even with the presence of the alignment constraint.

**Step 4:** Resolve the ambiguity

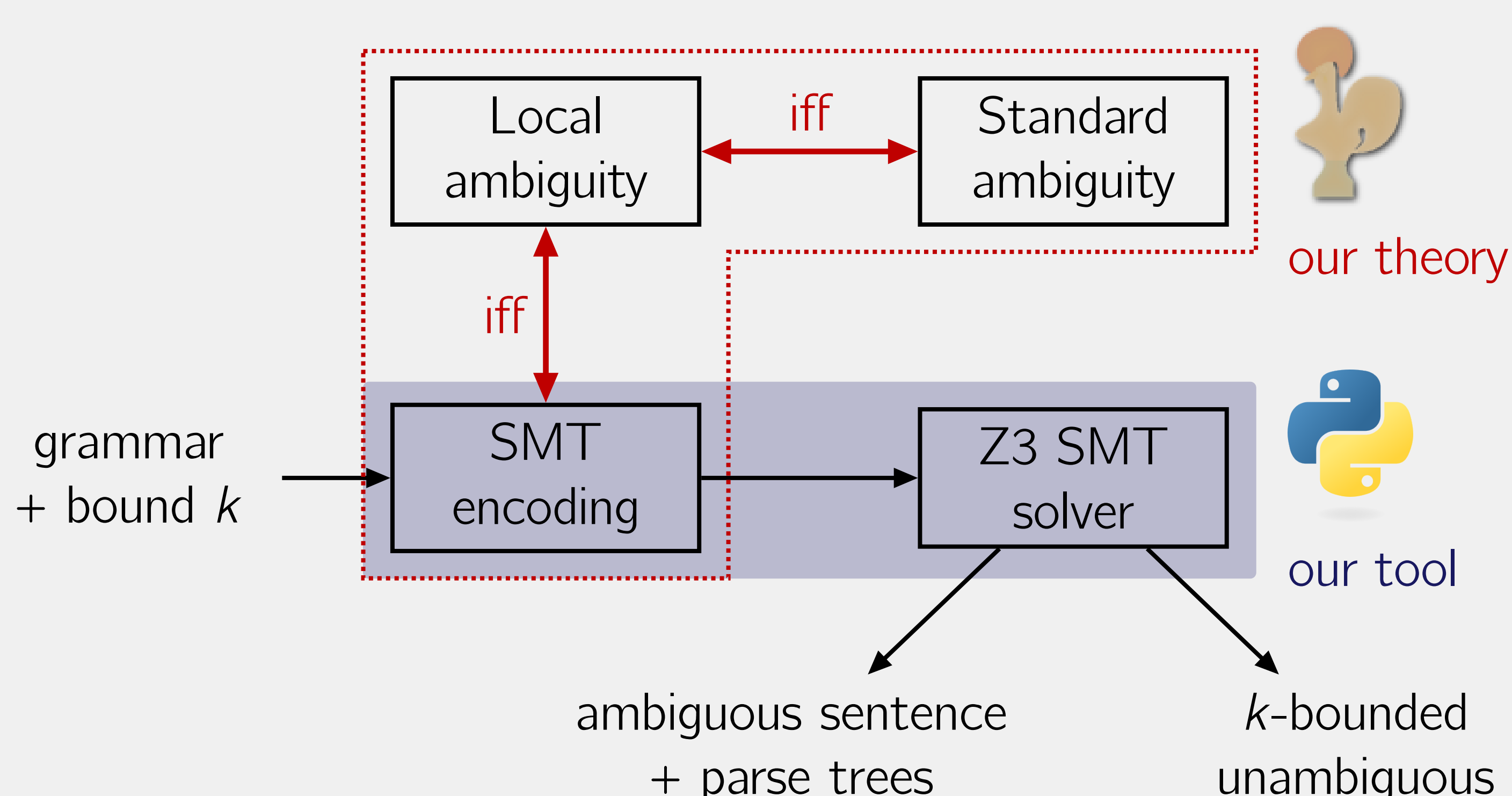
A possible solution is to reject this ambiguous sentence via an offside constraint  $\triangleright$  over the **do**-block:

```
block → ||stmt||+
stmt → nop | (do block)▷
```

**Step 5:** Check the refined grammar again

Lamb no longer finds any ambiguous sentence within a length of 20: that is, bounded unambiguous!

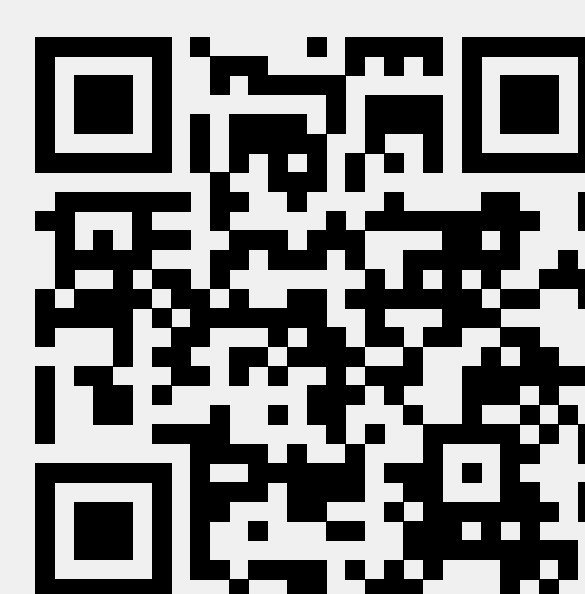
## Lamb: Layout-Sensitive Ambiguity Detector



where  $k$  is the upper bound length of the sentences being considered

## More...

Check our website



Read our paper

