# Poster: Input Grammar Oriented Symbolic Execution

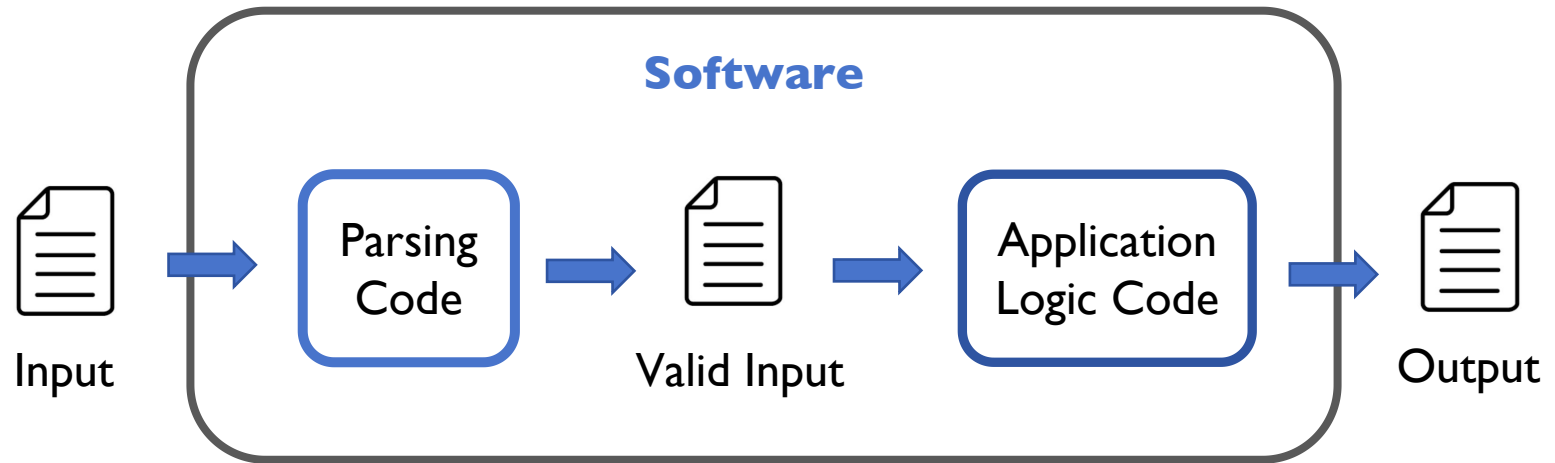Weijiang Hong (hongweijiang17@nudt.edu.cn)

joint work with Ke Ma, Yunlai Luo, Zhenbang Chen, Yufeng Zhang and Ji Wang

College of Computer, NUDT& College of Computer Science and Electronic Engineering, HNU, China
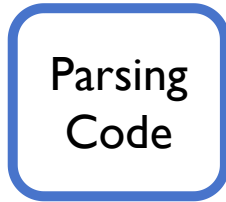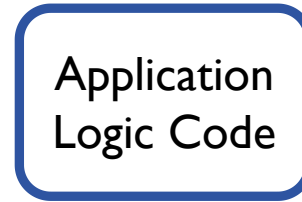
# 1. Background

# I. Background



**Software**

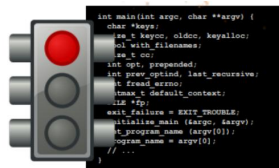Input → Parsing Code → Valid Input → Application Logic Code → Output

Input Error

**Low Coverage**

**Missed Bugs**

aa>hi</a>

# 1. Background

The underlying issue is that **symbolic execution** lacks awareness of the program's input format grammar.



Software

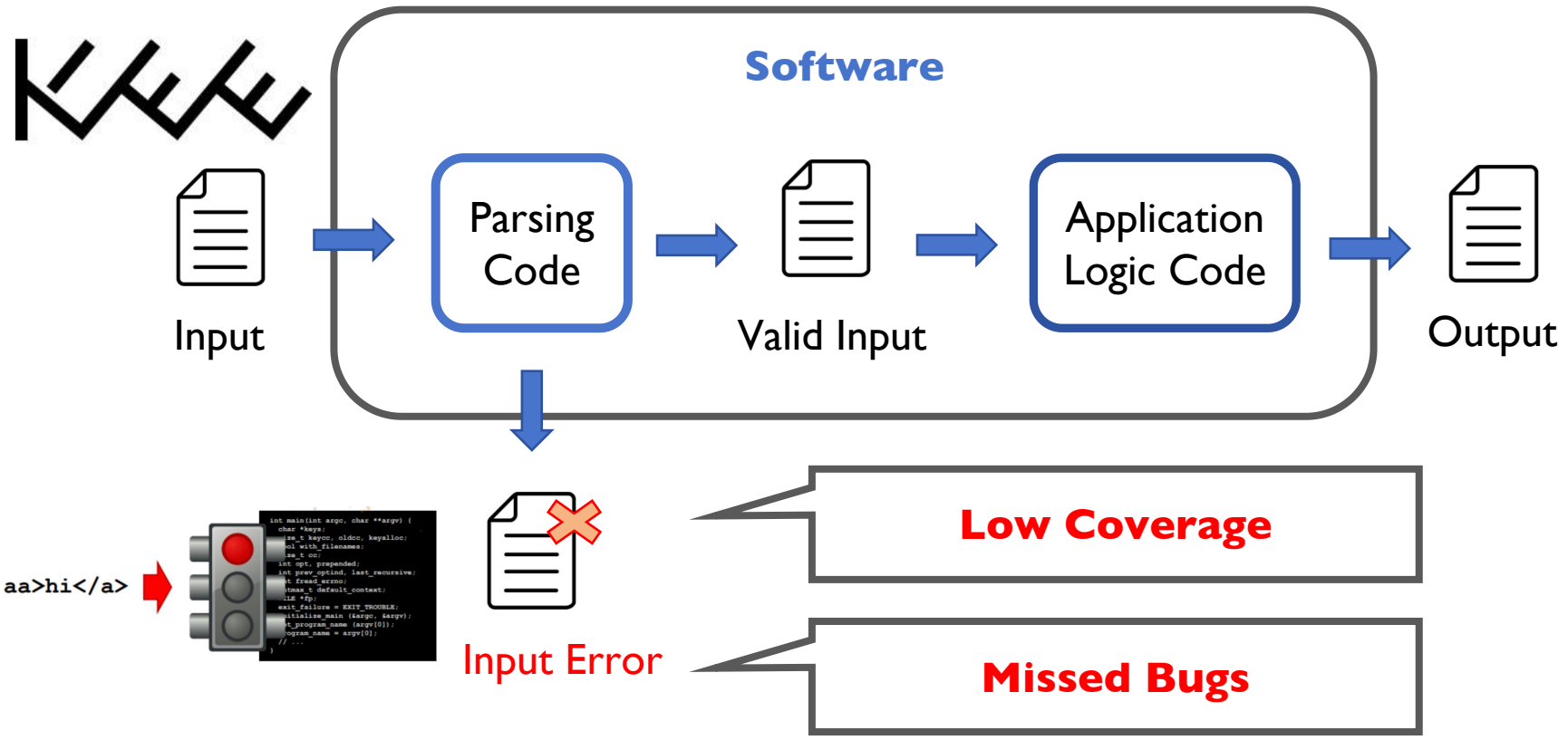Input → Parsing Code → Valid Input → Application Logic Code → Output
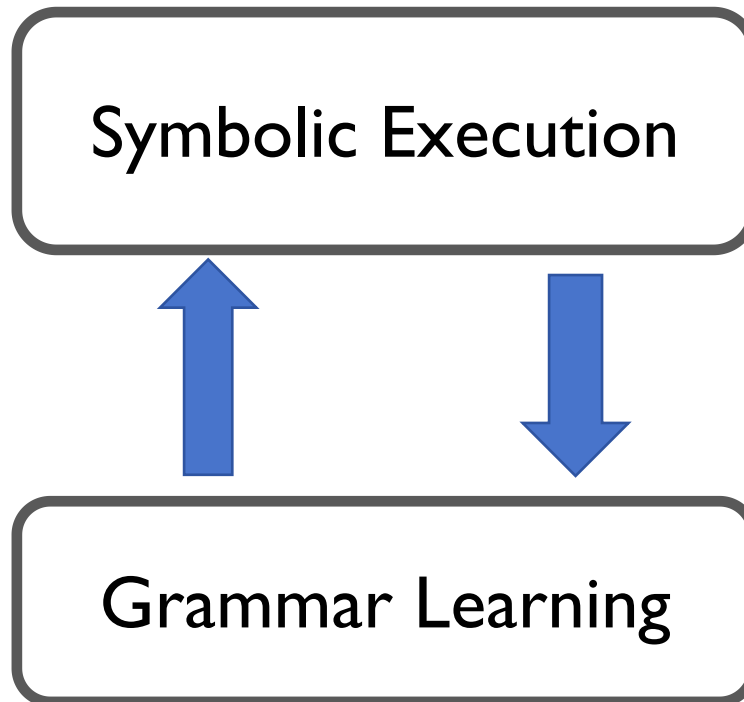
Input Error

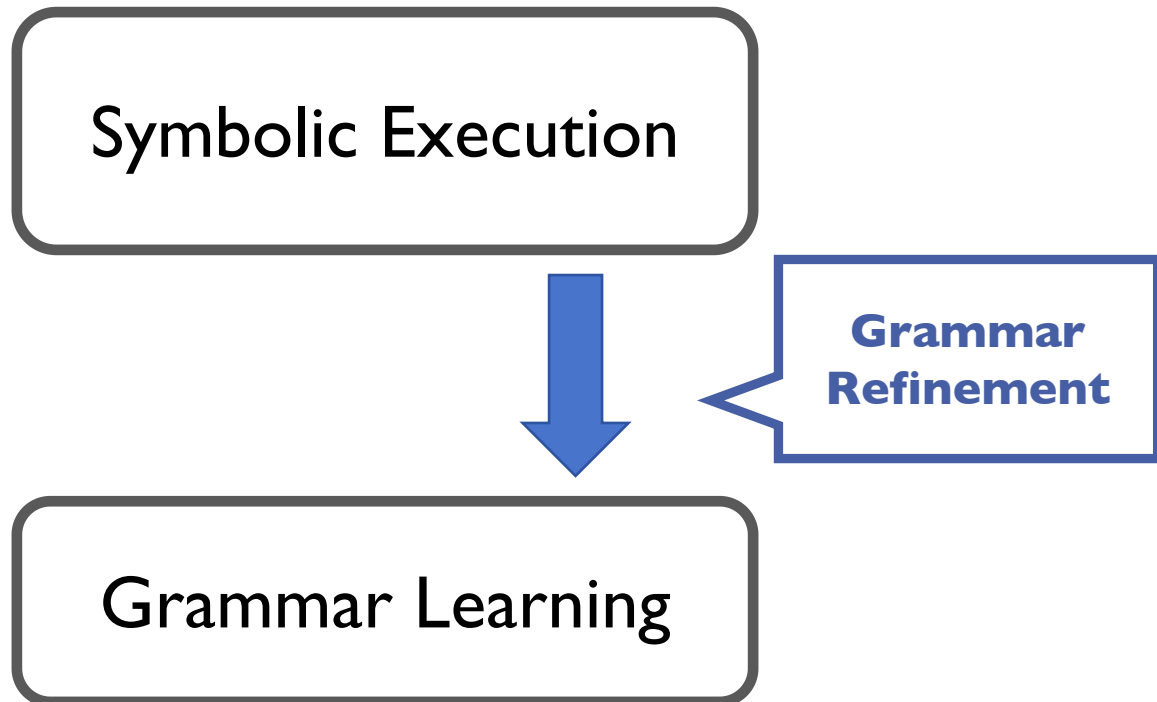**Low Coverage**

**Missed Bugs**

## 2. Key Idea

The underlying issue is that **symbolic execution** lacks awareness of the program's input format grammar.

# 2. Key Idea

The underlying issue is that **symbolic execution** <span style="color:red">lacks awareness of</span> the program's <span style="color:red">input format grammar.</span>

Symbolic Execution

**Grammar Refinement**

Grammar Learning

## 2. Key Idea

The underlying issue is that **symbolic execution** <span style="color:red">lacks awareness of</span> the program's <span style="color:red">input format grammar</span>.
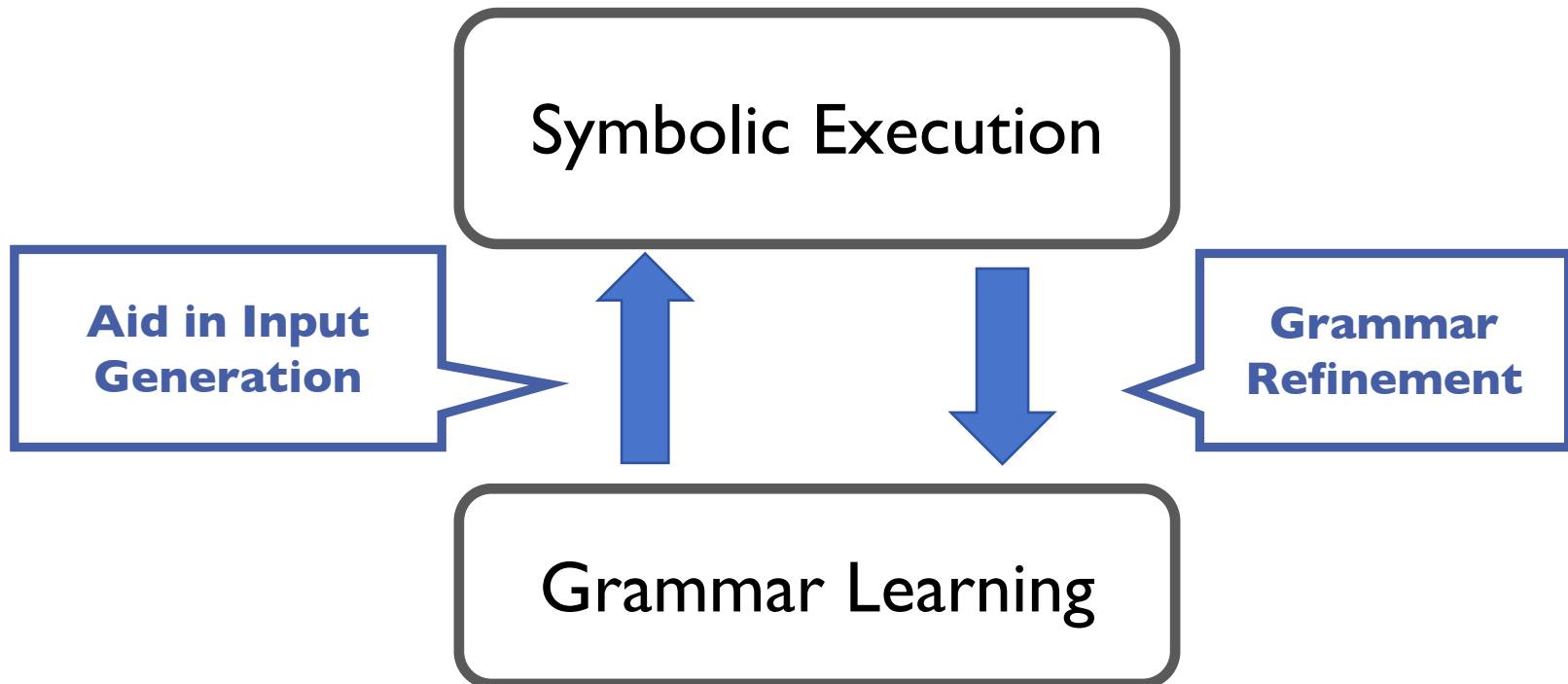
Symbolic Execution

Aid in Input Generation

Grammar Refinement

Grammar Learning

# 3. Framework

New input generated by **Constraint Solving**

[ISSTA'21] Grammar-agnostic symbolic
execution by token symbolization



Input

Token
Sequence

Heuristic
Search

**Token-based
Symbolic Execution**
[ISSTA'21]

"Interesting"
Token Sequence

Coverage

# 3. Framework

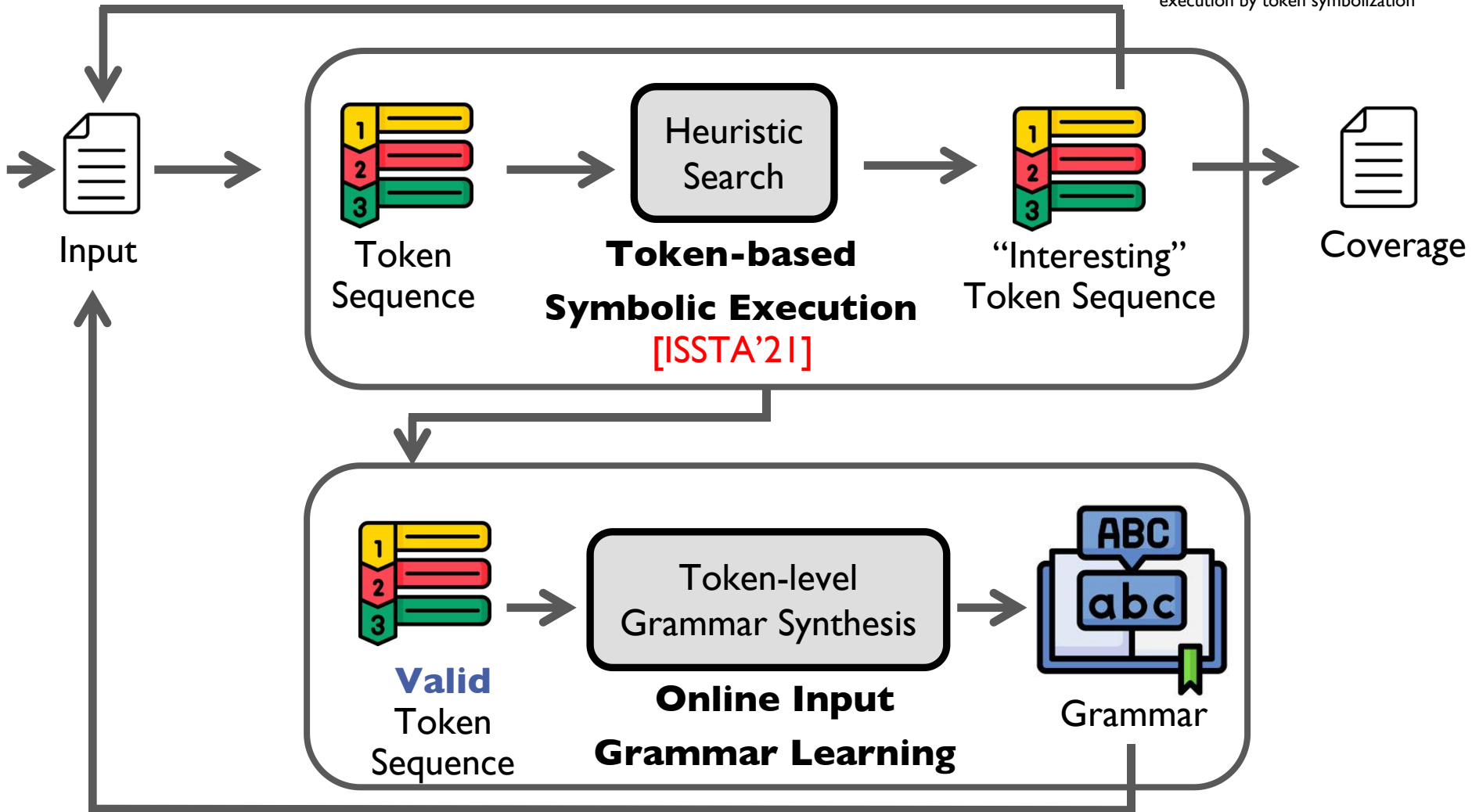New input generated by **Constraint Solving**

Input

Token Sequence

Heuristic Search

**Token-based Symbolic Execution** [ISSTA'21]

"Interesting" Token Sequence

Coverage

**Valid** Token Sequence

Token-level Grammar Synthesis

**Online Input Grammar Learning**

Grammar

New input specified by **Grammar**

# 3.1 Heuristic Search

New input generated by **Constraint Solving**

Input → Token Sequence → **Heuristic Search** / **Token-based Symbolic Execution** [ISSTA'21] → "Interesting" Token Sequence → Coverage

How to design the **heuristic search** that tries to generate new inputs, covering more syntax rules in the parsing code?

# 3.1 Heuristic Search



**Input: bitor a**
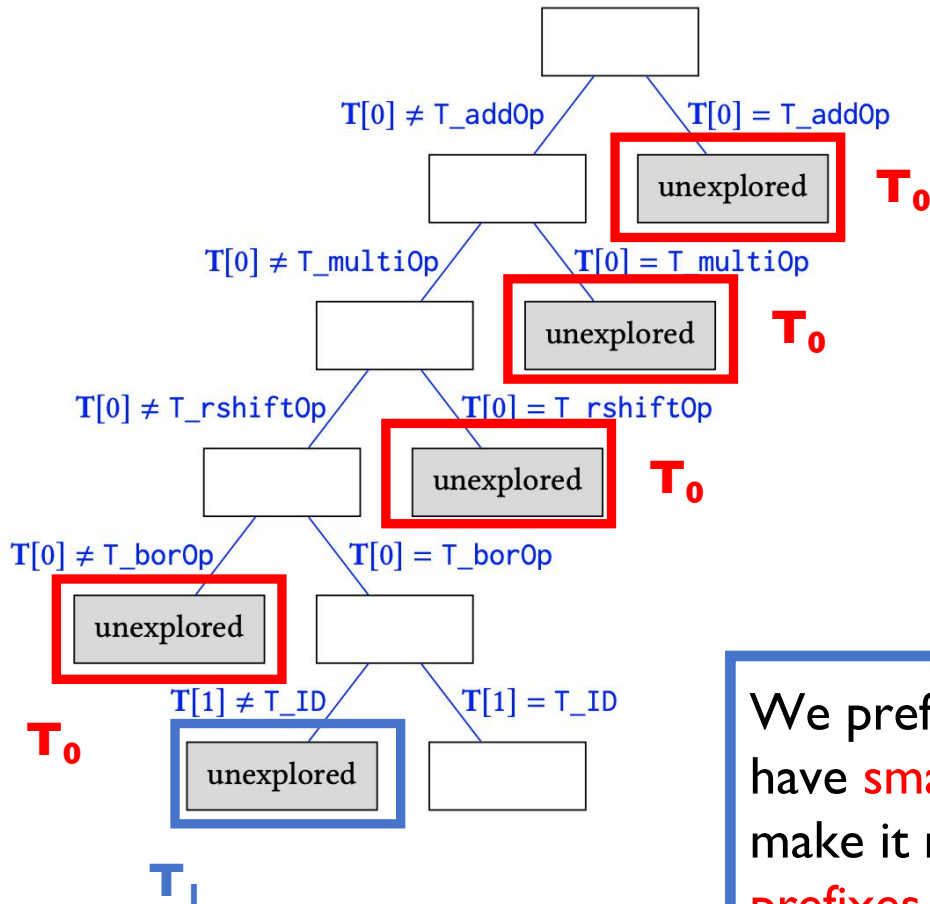
Token Sequence: $\langle$ T_borOp, T_ID $\rangle$

Path Condition:

$T[0] \neq$ T_addOp $\wedge$ $T[0] \neq$ T_multiOp $\wedge$

$T[0] \neq$ T_rfhiftOp $\wedge$ $T[0] =$ T_borOp $\wedge$

$T[1] =$ T_ID

# 3.1 Heuristic Search



**Input: bitor a**

Token Sequence: $\langle$T_borOp, T_ID$\rangle$
Path Condition:
T[0] ≠ T_addOp ∧ T[0] ≠ T_multiOp ∧
T[0] ≠ T_rfhiftOp ∧ T[0] = T_borOp ∧
T[1] = T_ID

We prefer to select the unexplored branches that have smaller token indexes in priority, which make it more easier to generate different token prefixes.

New Path Condition:
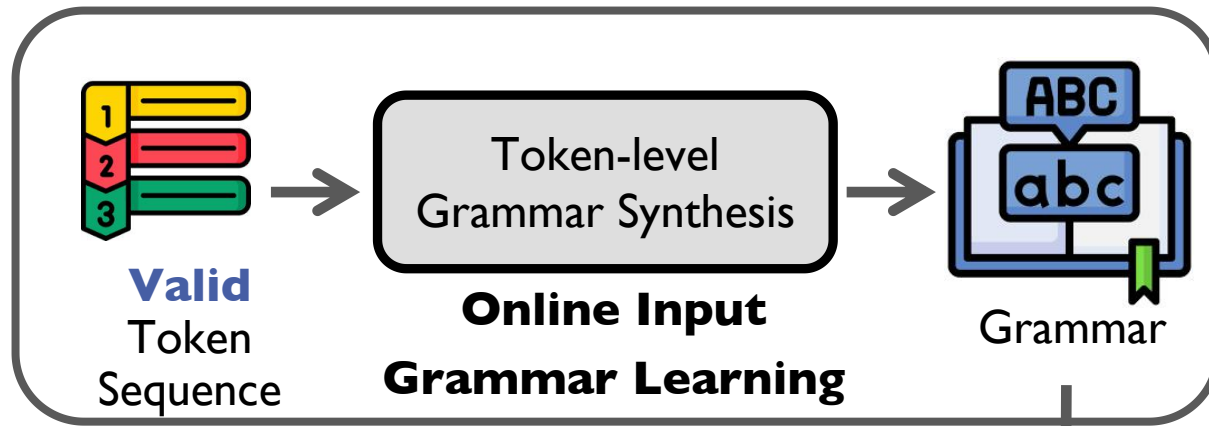T[0] = T_addOp

# 3.2 Grammar Synthesis

Input

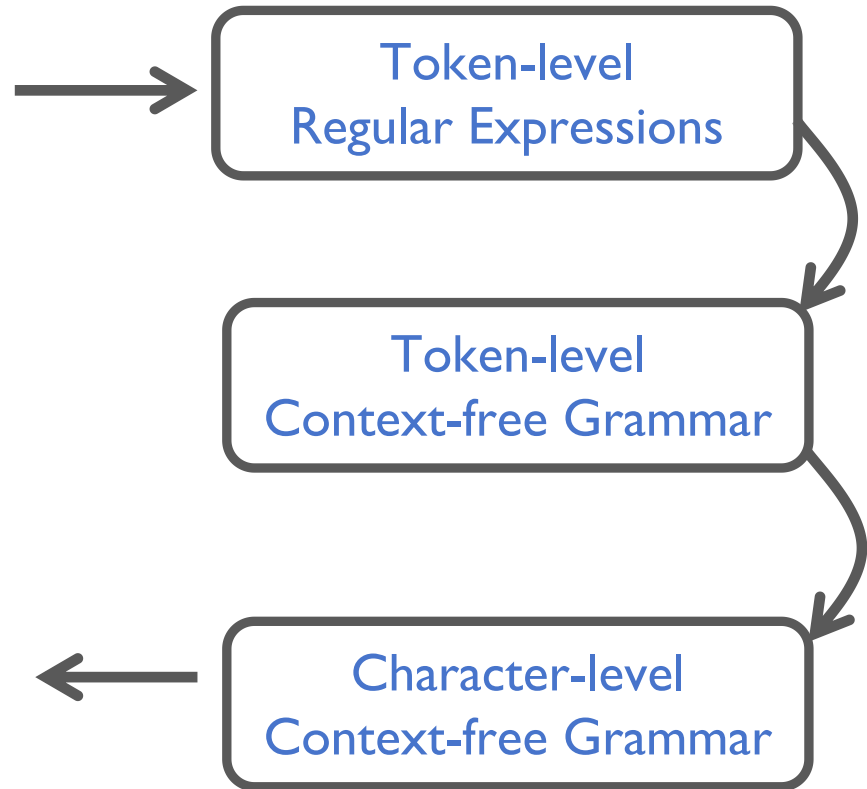How to learn **better grammars** with as few iterations as possible?

**Valid** Token Sequence

Token-level Grammar Synthesis

**Online Input Grammar Learning**

Grammar

New input specified by **Grammar**

# 3.2 Grammar Synthesis

**Valid Input: a+a+1**

Token Sequence:
⟨T_ID, T_OP, T_ID, T_OP, T_NUM⟩ ⟶ [ Token-level Regular Expressions ]

[ Token-level Context-free Grammar ]

[ Character-level Context-free Grammar ]

$$
\begin{array}{rcl}
\langle expr \rangle & ::= & \langle expr \rangle \langle op \rangle \langle term \rangle \mid \langle term \rangle \\
\langle term \rangle & ::= & \langle num \rangle \mid \langle ID \rangle \\
\langle num \rangle & ::= & \text{'0'} \mid \text{'1'} \mid \dots \mid \text{'9'} \\
\langle ID \rangle & ::= & \text{'a'} \mid \text{'b'} \mid \dots \mid \text{'y'} \mid \text{'z'} \\
\langle op \rangle & ::= & \text{'+'}
\end{array}
$$

# 3. Framework

New input generated by **Constraint Solving**

Input

Token Sequence

Heuristic Search

**Token-based Symbolic Execution**
[ISSTA'21]

"Interesting" Token Sequence

Coverage

**Valid** Token Sequence

Token-level Grammar Synthesis

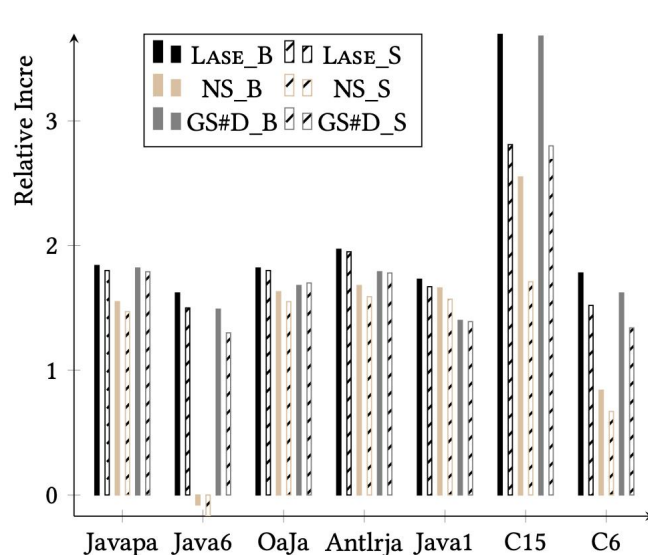**Online Input Grammar Learning**

Grammar
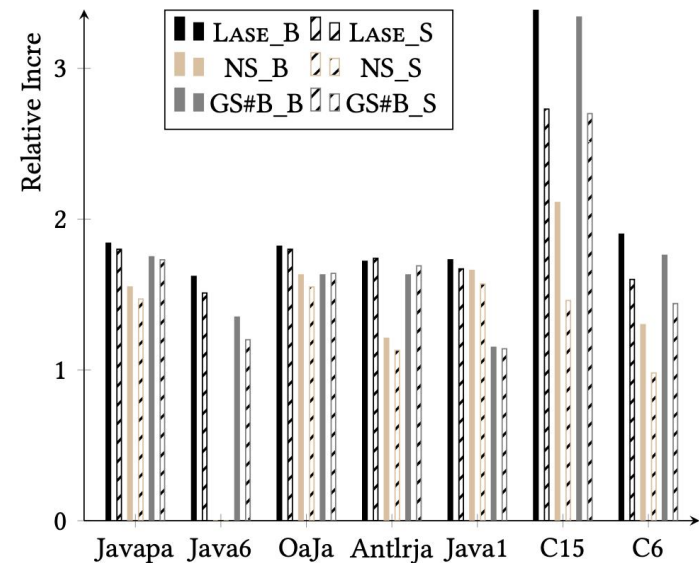
New input specified by **Grammar**

# 4. Results

- for Java compiler **Janino**, on average, we achieve a 50.92% increase in statement coverage under the BFS strategy, and a 57.68% increase in statement coverage under the DFS strategy.

- for C compiler **CLoli**, on average, we achieve a 289.57% under BFS strategy, and a 342.09% increase in statement coverage under the DFS strategy.



Relative increase compared to
GADSE under **DFS**

Relative increase compared to
GADSE under **BFS**

# Thank you!
# Q & A