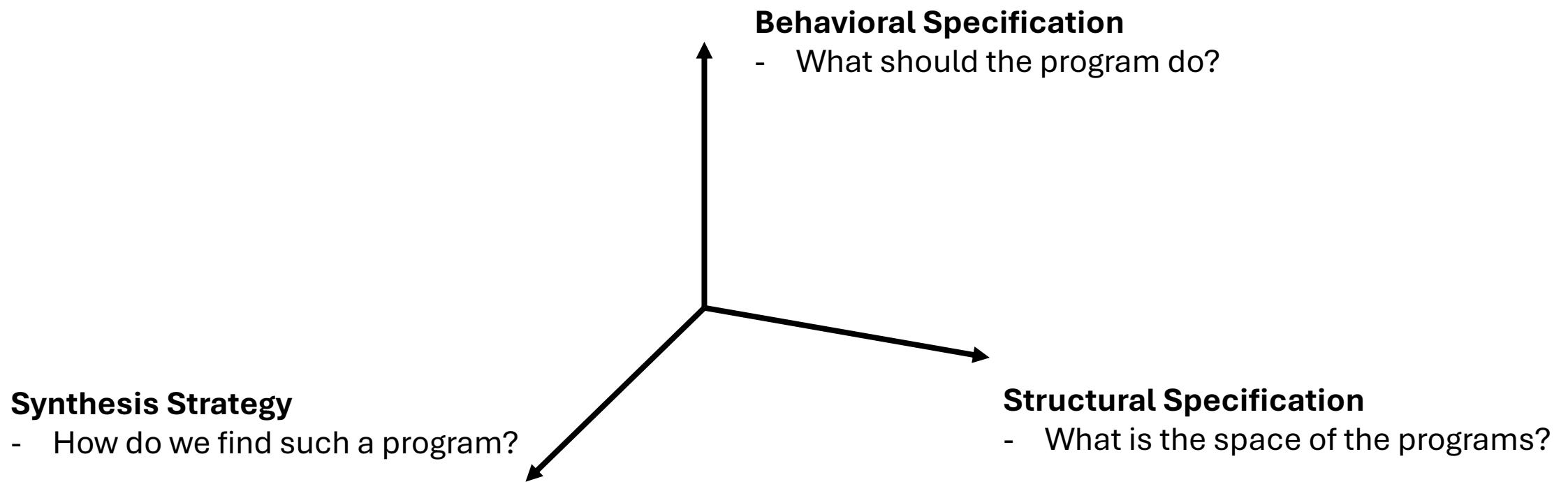


# Machine Programming

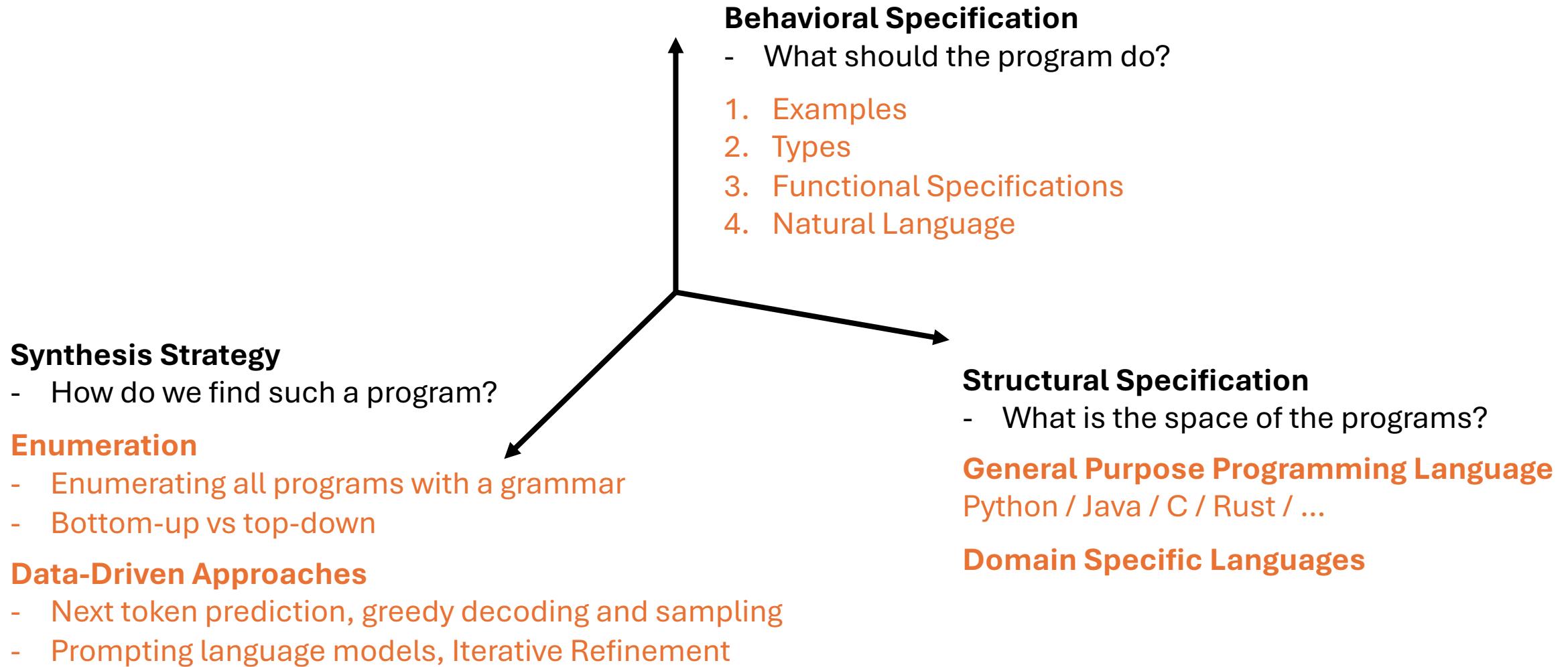
Lecture 8 – Controlled Decoding and Steering

Ziyang Li

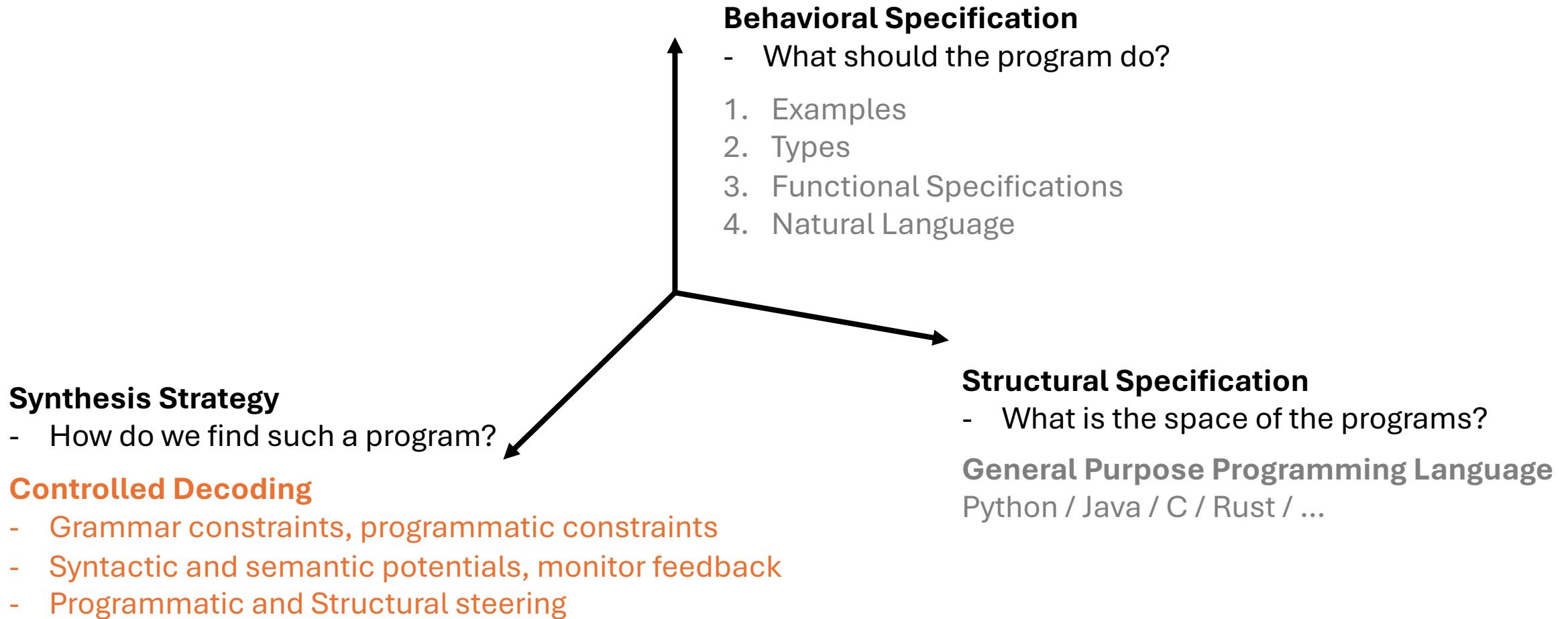
# Dimensions in Program Synthesis



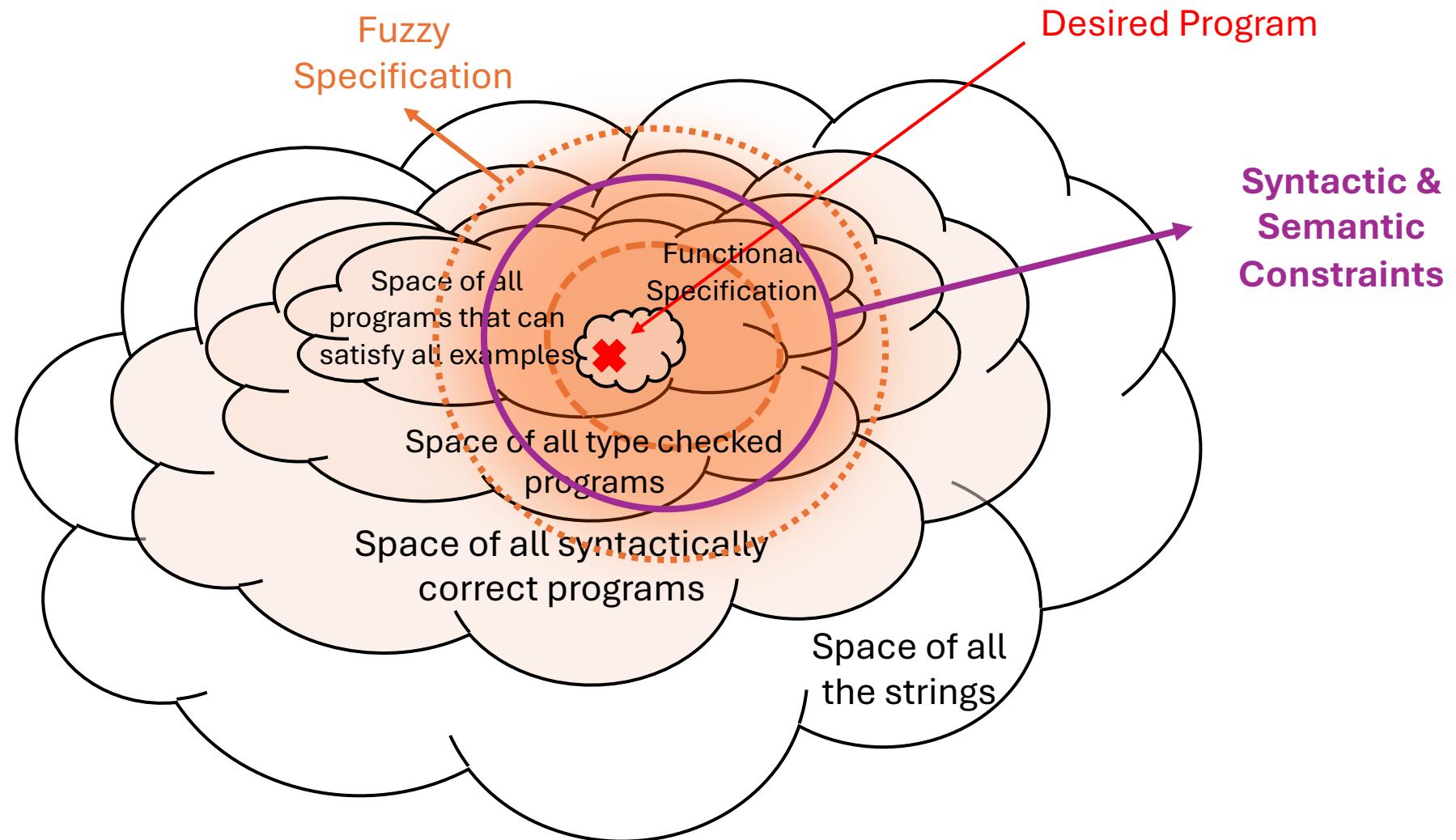
# The Course So Far



# Today



# High Level Picture



# High Level Picture

- We studied foundations of programming languages and synthesis
  - Syntax (regular tree grammar), Semantics (denotational, operational)
  - Enumeration (top-down, bottom-up) that searches within grammar
- We studied basic usage of language models
  - Purely neural: next token prediction → sequential decoding
  - Prompting and iterative refinement: elicit better programs from LLM
- Question:
  - Can we inject more **structure** during **neural** generation process

# Review: Syntax in Regular Tree Grammar

```
(Program) P ::= L | N
  (List) L ::= input
           | empty
           | single(N)
           | concat(L, L)
  (Number) N ::= len(L)
              | min(L)
              | add(N, N)
              | 0 | 1 | 2 | ...
```

# Review: Syntax in Regular Tree Grammar

(Program)  $P ::= L \mid N$   
(List)  $L ::= \text{input}$   
|  $\text{empty}$   
|  $\text{single}(N)$   
|  $\text{concat}(L, L)$

(Number)  $N ::= \text{len}(L)$   
|  $\text{min}(L)$   
|  $\text{add}(N, N)$   
|  $0 \mid 1 \mid 2 \mid \dots$

→

(Program)  $P ::= L \mid N$   
(List)  $L ::= \text{input}$   
|  $[]$   
|  $[N]$   
|  $L \text{ ++ } L$

(Number)  $N ::= \text{len}(L)$   
|  $\text{min}(L)$   
|  $N + N$   
|  $0 \mid 1 \mid 2 \mid \dots$

# Review: Syntax in Context Free Grammar

|   |                   |  |
|---|-------------------|--|
| <p>(Program) <math>P ::= L \mid N</math></p> <p>(List) <math>L ::= \text{input}</math><br/>  empty<br/>  single(<math>N</math>)<br/>  concat(<math>L, L</math>)</p> <p>(Number) <math>N ::= \text{len}(L)</math><br/>  min(<math>L</math>)<br/>  add(<math>N, N</math>)<br/>  0   1   2   ...</p> | $\longrightarrow$ | <p>(Program) <math>P \leftarrow L \mid N</math></p> <p>(List) <math>L \leftarrow \text{'input'}</math><br/>  '<math>[</math>' '<math>]</math>'<br/>  '<math>[</math>' <math>N</math> '<math>]</math>'<br/>  <math>L</math> '+' <math>L</math></p> <p>(Number) <math>N \leftarrow \text{'len'} \text{ '(' } L \text{ ')'}'</math><br/>  <math>\text{'min'} \text{ '(' } L \text{ ')'}'</math><br/>  <math>N</math> '+' <math>N</math><br/>  '0'   '1'   '2'   ...</p> |
|---|-------------------|--|

# From Program to Abstract Syntax Tree

Concrete Program

```
min(input ++ [0])
```

Concrete Grammar (Context Free Grammar)

```
(Program) P <- L | N
(List) L <- 'input'
      | '[' ']'
      | '[' N ']'
      | L '++' L
(Number) N <- 'len' '(' L ')'
        | 'min' '(' L ')'
        | N '+' N
        | '0' | '1' | '2' | ...
```

Abstract Grammar (Regular Tree Grammar)

```
(Program) P ::= L | N
(List) L ::= input
        | empty
        | single(N)
        | concat(L, L)
(Number) N ::= len(L)
        | min(L)
        | add(N, N)
        | 0 | 1 | 2 | ...
```

# From Program to Abstract Syntax Tree

Concrete Program

```
min(input ++ [0])
```

Token Sequence

```
['min', '(', 'input', '++', '[', '0', ']']
```

Concrete Grammar (Context Free Grammar)

```
(Program) P <- L | N
(List) L <- 'input'
      | '[' ']'
      | '[' N ']'
      | L '+' L
(Number) N <- 'len' '(' L ')'
        | 'min' '(' L ')'
        | N '+' N
        | '0' | '1' | '2' | ...
```

Abstract Grammar (Regular Tree Grammar)

```
(Program) P ::= L | N
(List) L ::= input
        | empty
        | single(N)
        | concat(L, L)
(Number) N ::= len(L)
        | min(L)
        | add(N, N)
        | 0 | 1 | 2 | ...
```

# From Program to Abstract Syntax Tree

Concrete Program

`min(input ++ [0])`

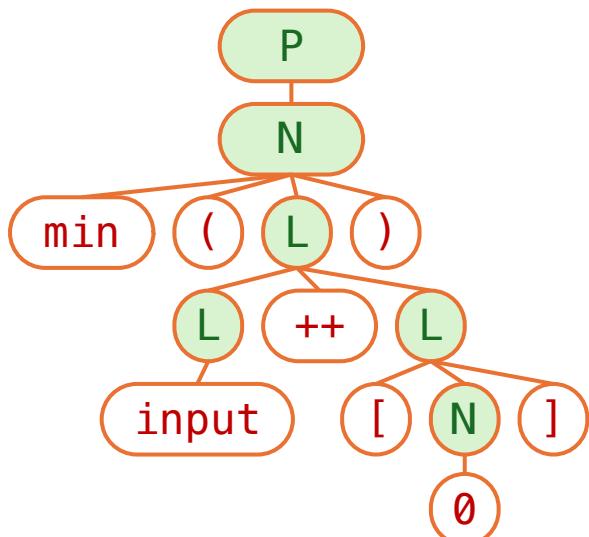


Token Sequence

`['min', '(', 'input', '++', '[', '0', ']']`



Concrete Syntax Tree



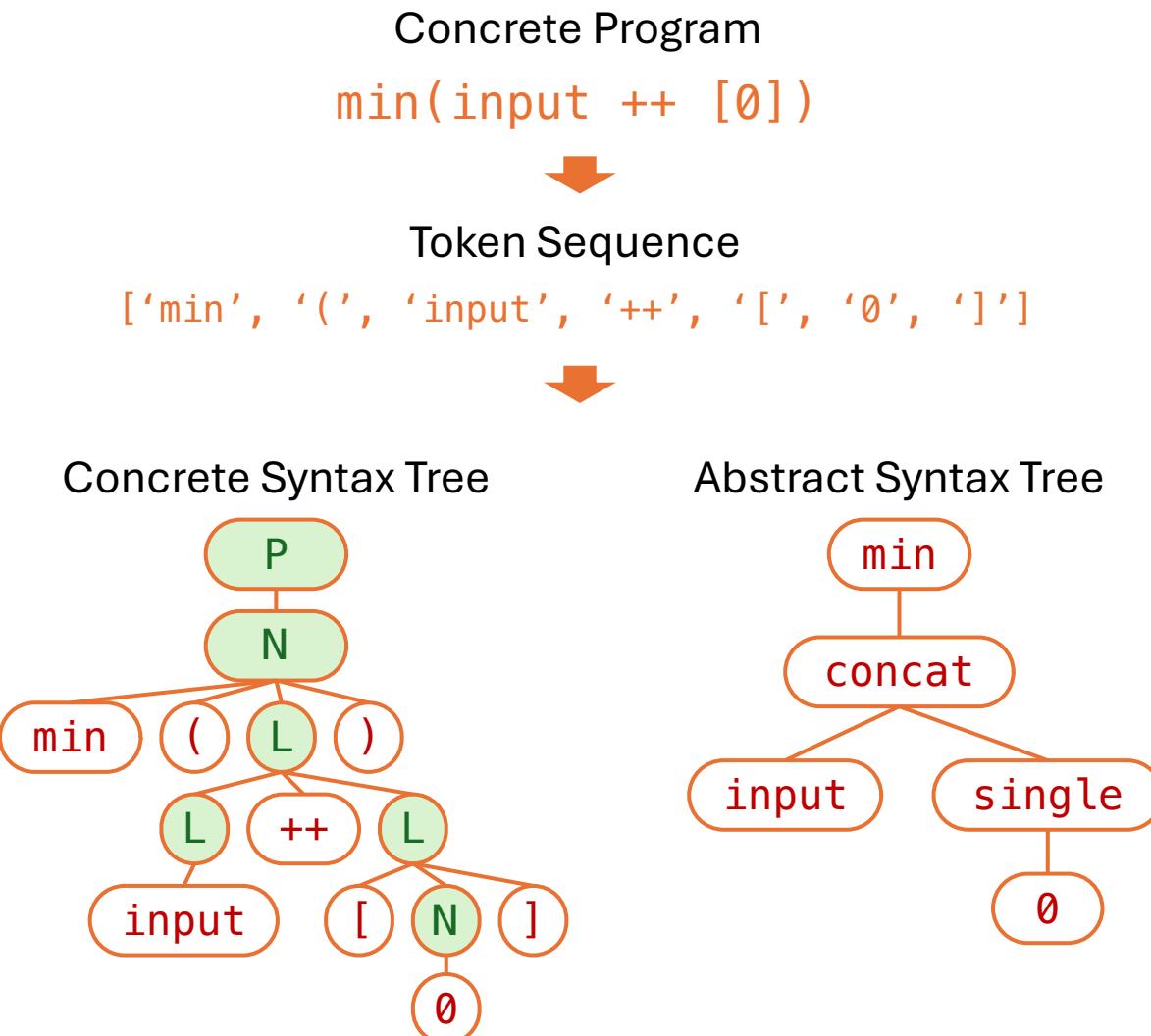
Concrete Grammar (Context Free Grammar)

```
(Program) P <- L | N
(List) L <- 'input'
      | '[' ']'
      | '[' N ']'
      | L '++' L
(Number) N <- 'len' '(' L ')'
      | 'min' '(' L ')'
      | N '+' N
      | '0' | '1' | '2' | ...
```

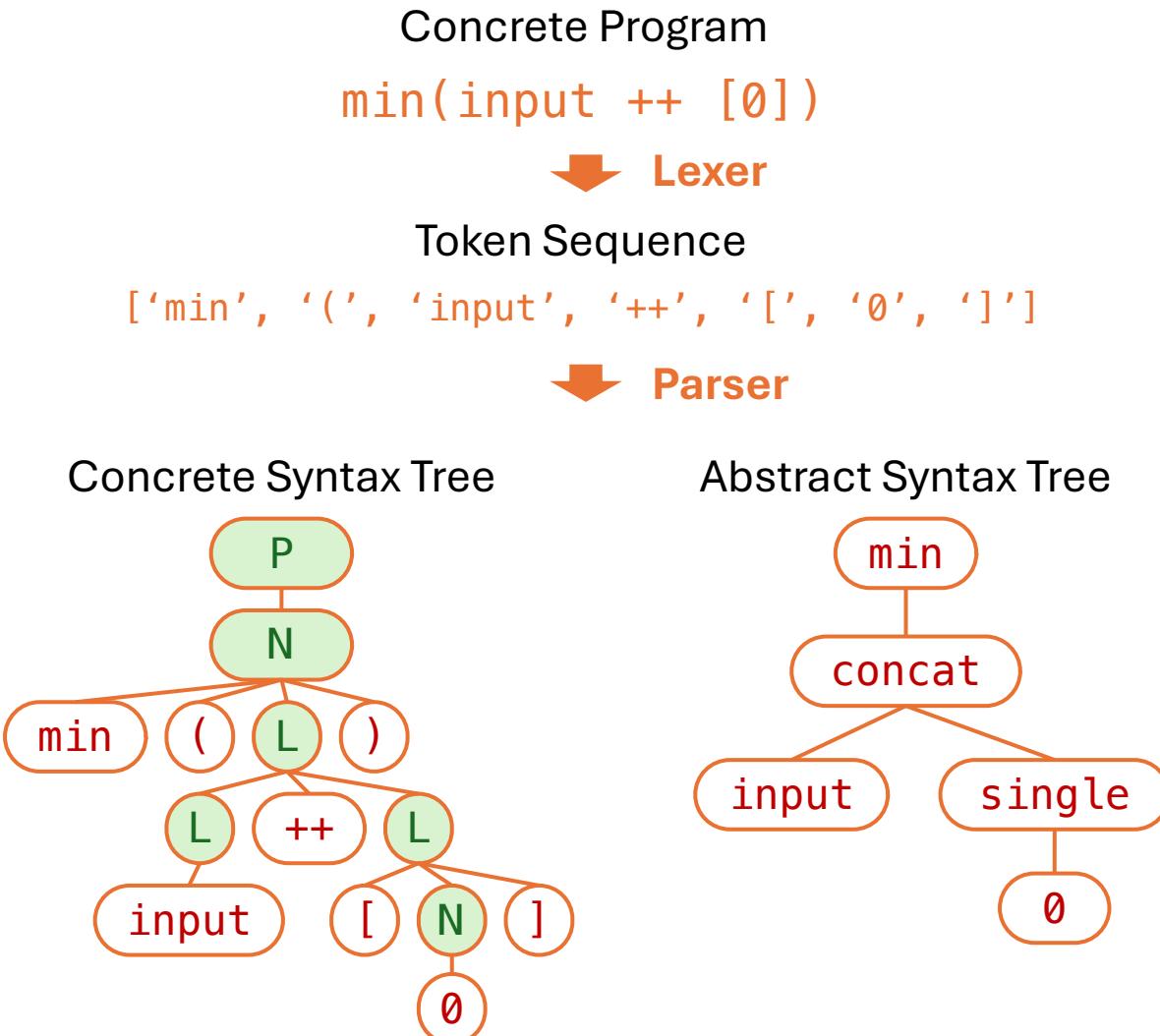
Abstract Grammar (Regular Tree Grammar)

```
(Program) P ::= L | N
(List) L ::= input
      | empty
      | single(N)
      | concat(L, L)
(Number) N ::= len(L)
      | min(L)
      | add(N, N)
      | 0 | 1 | 2 | ...
```

# From Program to Abstract Syntax Tree



# From Program to Abstract Syntax Tree



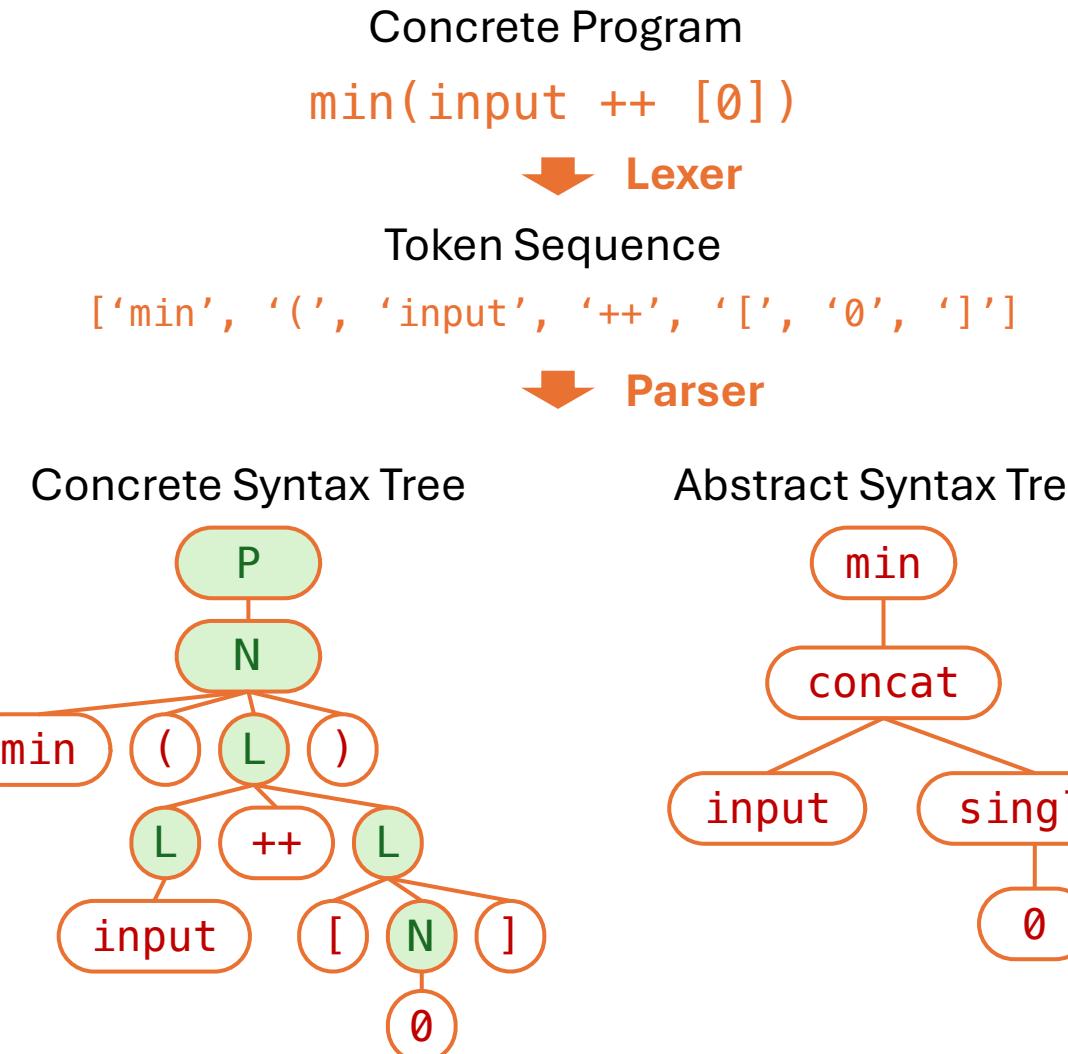
## Concrete Grammar (Context Free Grammar)

```
(Program) P <- L | N  
(List) L <- 'input'  
| '[' ']'  
| '[' N ']'  
| L '++' L  
(Number) N <- 'len' '(' L ')' |  
| 'min' '(' L ')' |  
| N '+' N |  
| '0' | '1' | '2' | ...
```

## Abstract Grammar (Regular Tree Grammar)

```
(Program) P ::= L | N  
(List) L ::= input  
| empty  
| single(N)  
| concat(L, L)  
(Number) N ::= len(L)  
| min(L)  
| add(N, N)  
| 0 | 1 | 2 | ...
```

# Lexer and Parser



I have a concrete grammar in CFG

(Program)  $P \leftarrow L \mid N$   
(List)  $L \leftarrow \text{'input'}$   
|  $\text{'['}' \text{ ]'}$   
|  $\text{'['} N \text{ ]'}$   
|  $L \text{ '+' } L$   
(Number)  $N \leftarrow \text{'len'} \text{ '(' } L \text{ ')'} \mid \text{'min'} \text{ '(' } L \text{ ')'} \mid N \text{ '+' } N \mid 0 \mid 1 \mid 2 \mid \dots$

Corresponding to it there is an abstract grammar

(Program)  $P ::= L \mid N$   
(List)  $L ::= \text{input}$   
| empty  
| single( $N$ )  
| concat( $L, L$ )  
(Number)  $N ::= \text{len}(L)$   
| min( $L$ )  
| add( $N, N$ )  
| 0 | 1 | 2 | ...

Can you help me write a treesitter parser in javascript?

# Lexer and Parser

Concrete Program

min(input ++ [0])

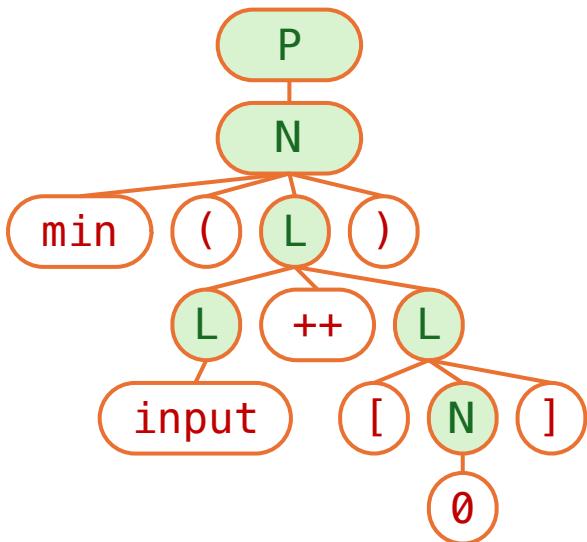
↓ Lexer

Token Sequence

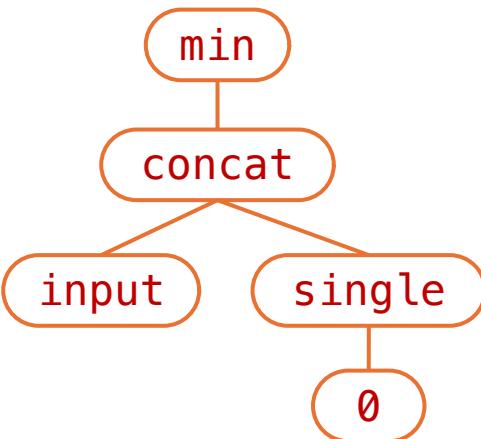
['min', '(', 'input', '++', '[', '0', ']']

↓ Parser

Concrete Syntax Tree



Abstract Syntax Tree



```
I have a concrete grammar in CFG  
(Program) P <-> L | N  
(List) L <-> 'input'  
| '(' ')'  
| '[' ']'  
| L ++'L'  
(Number) N <-> 'len' '(' L ')'  
| 'min' '(' L ')'  
| N ++'N'  
| '0' | '1' | '2' | ...  
  
Corresponding to it there is an abstract grammar  
(Program) P ::= L | N  
(List) L ::= input  
| empty  
| single(N)  
| concat(L, L)  
(Number) N ::= len(L)  
| min(L)  
| add(N, N)  
| 0 | 1 | 2 | ...  
  
Can you help me write a treesitter parser in javascript?  
  
1 module.exports = grammar({  
2   name: 'listnum',  
3   rules: {  
4     program: $ => choice($.list, $.number),  
5  
6     list: $ => choice($.input, $.empty_list, $.single, $.concat),  
7     input: _ => 'input',  
8     empty_list: _ => seq('[' , ']'),  
9     single: $ => seq('[' , $.number, ']'),  
10    concat: $ => prec.left(1, seq($.list, '++', $.list)),  
11  
12    number: $ => choice($.len, $.min, $.add, $.int_lit),  
13  
14    len: $ => seq('len', '(', $.list, ')'),  
15    min: $ => seq('min', '(', $.list, ')'),  
16    add: $ => prec.left(2, seq($.number, '+', $.number)),  
17    int_lit: $ => alias($.int, 'int'),  
18    int: _ => /[0-9]+/,  
19  }  
20};
```



# Lexer and Parser in the Wild

JSON (json.l)

```
%{  
#include "y.tab.h"  
}  
  
%%  
"true"      return TRUE;  
"false"     return FALSE;  
"null"      return NULLVAL;  
[0-9]+       { yyval = atoi(yytext); return NUMBER; }  
\\"([^\"]*)\" { yyval = strdup(yytext); return STRING; }  
[\t\n]+      /* skip */  
"{"         return '{';  
"}"         return '}';  
"["         return '[';  
"]"         return ']';  
":"         return ':';  
","         return ',';  
%%
```

JSON (json.y)

```
%token TRUE FALSE NULLVAL NUMBER STRING  
  
%%  
value: STRING  
| NUMBER  
| object  
| array  
| TRUE  
| FALSE  
| NULLVAL  
;  
  
object: '{}'  
| '{} members {}'  
;  
  
members: pair  
| members ',' pair  
;  
  
pair: STRING ':' value ;  
  
array: '[]'  
| '[' elements ']'  
;  
  
elements: value  
| elements ',' value  
;
```

Python (python.lalrpop)

```
main ▾ Parser / parser / src / python.lalrpop  
Code Blame 1796 lines (1644 loc) · 60.4 KB  
123     ast::Stmt::Assign(  
124         ast::StmtAssign { targets, value, type_comment: None, range: (location..end_location).into() }  
125     )  
126 },  
127 <location:@L> <target:TestOrStarExprList> <op:AugAssign> <rhs:TestListOrYieldExpr> <end_location:@R> => {  
128     ast::Stmt::AugAssign(  
129         ast::StmtAugAssign {  
130             target: Box::new(set_context(target, ast::ExprContext::Store)),  
131             op,  
132             value: Box::new(rhs),  
133             range: (location..end_location).into()  
134         },  
135     ),  
136 },  
137 <location:@L> <target:Test<"all">> ":" <annotation:Test<"all">> <rhs:AssignSuffix?> <end_location:@R> => {  
138     let simple = target.is_name_expr();  
139     ast::Stmt::AnnAssign(  
140         ast::StmtAnnAssign {  
141             target: Box::new(set_context(target, ast::ExprContext::Store)),  
142             annotation: Box::new(annotation),  
143             value: rhs.map(Box::new),  
144             simple,  
145             range: (location..end_location).into()  
146         },  
147     ),  
148 },  
149 };
```

# High Level Picture

- We studied foundations of programming languages and synthesis
  - Syntax (regular tree grammar), Semantics (denotational, operational)
  - Enumeration (top-down, bottom-up) that searches within grammar
- We studied basic usage of language models
  - Purely neural: next token prediction → sequential decoding
  - Prompting and iterative refinement: elicit better programs from LLM
- Question:
  - Can we inject more **structure** during **neural** generation process

# Possible Next Token by Parser

Partial Program

min(input ++ [0])

Token Sequence

['min', '(', 'input', '++', '[', '0', ']']

```
(Program) P <- L | N
  (List) L <- 'input'
        | '[' ']'
        | '[' N ']'
        | L '++' L
  (Number) N <- 'len' '(' L ')'
        | 'min' '(' L ')'
        | N '+' N
        | '0' | '1' | '2' | ...
```

# Possible Next Token by Parser

Partial Program

```
min(input ++ [0])
```

Token Sequence

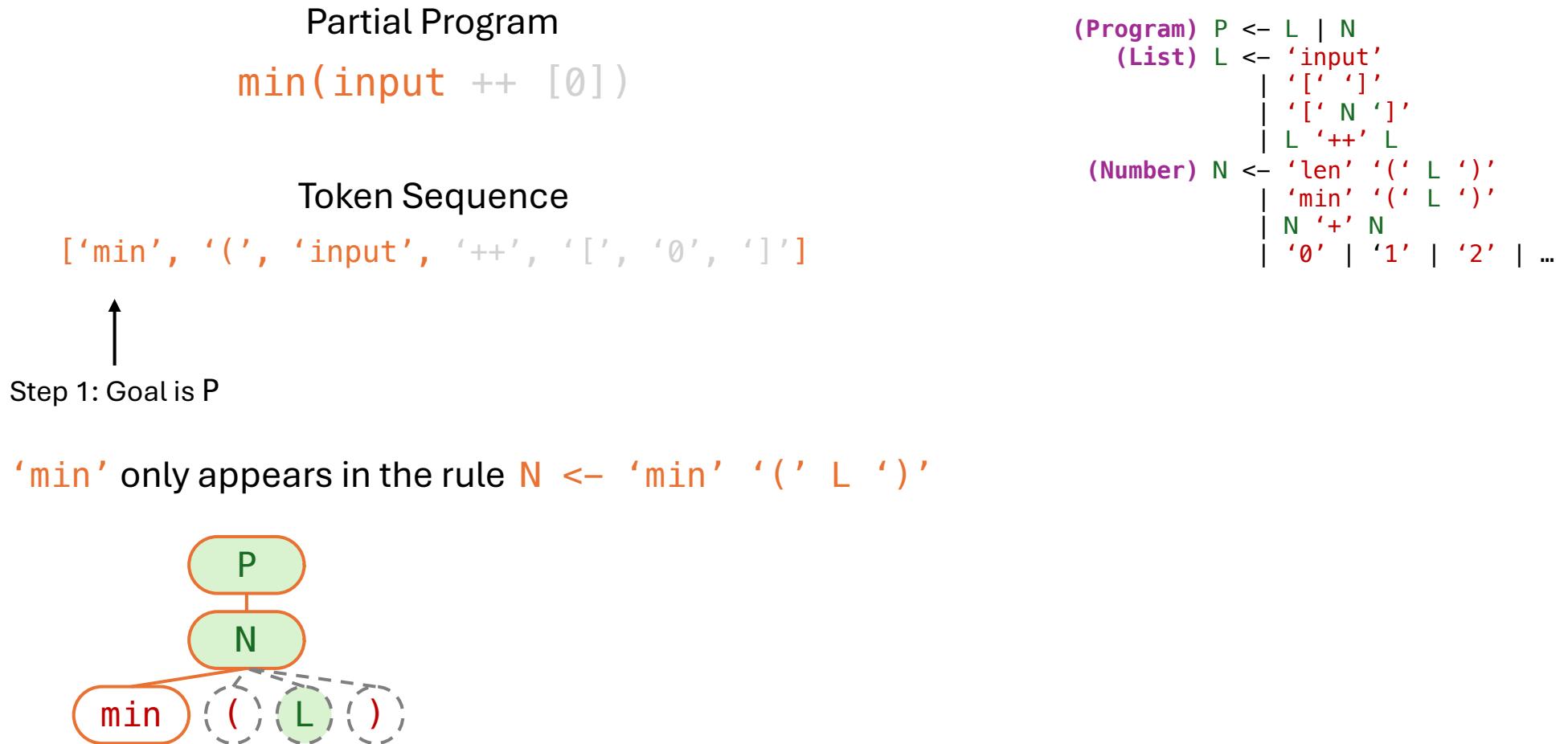
```
['min', '(', 'input', '++', '[', '0', ']']
```

Step 1

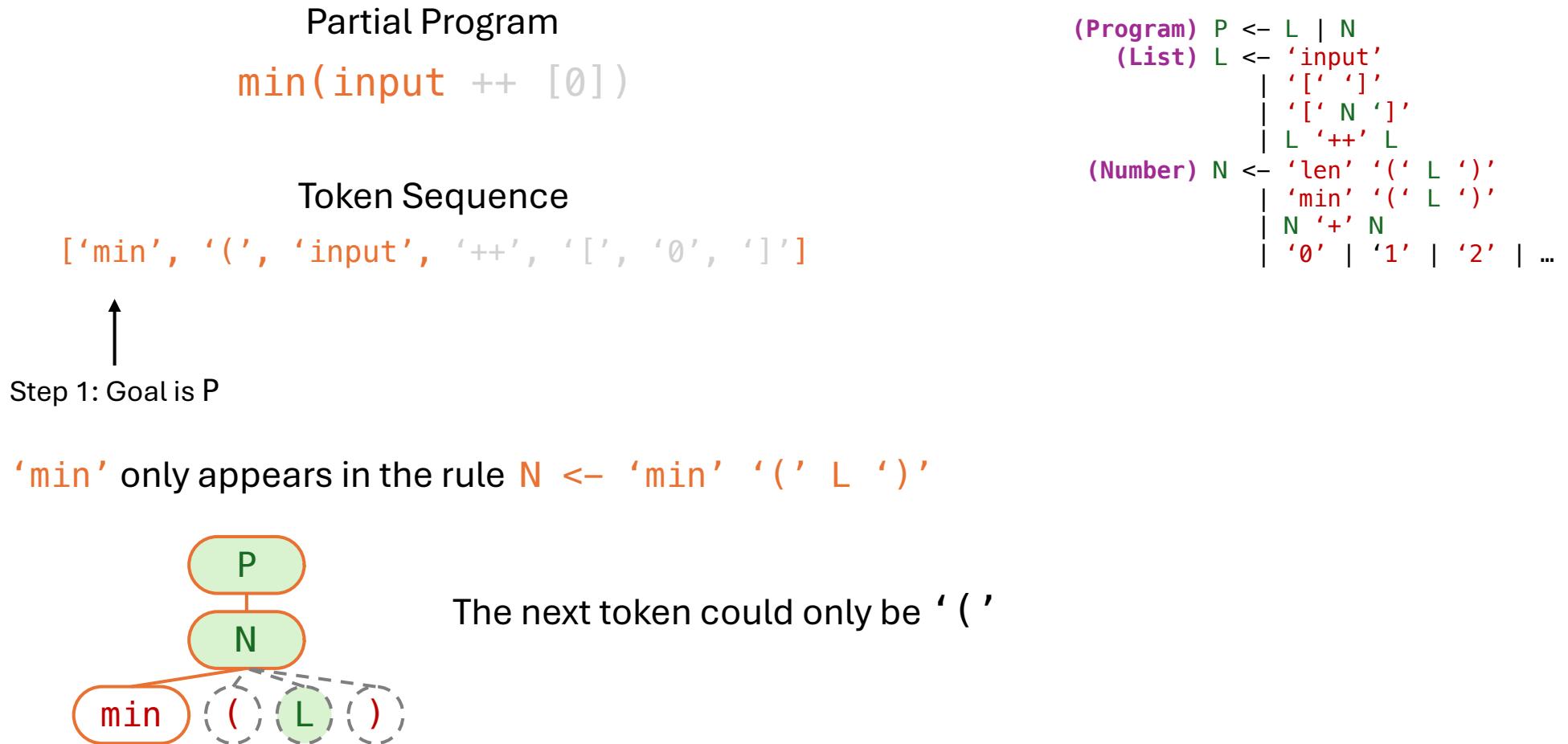


```
(Program) P <- L | N
  (List) L <- 'input'
        | '[' ']'
        | '[' N ']'
        | L '++' L
  (Number) N <- 'len' '(' L ')'
        | 'min' '(' L ')'
        | N '+' N
        | '0' | '1' | '2' | ...
```

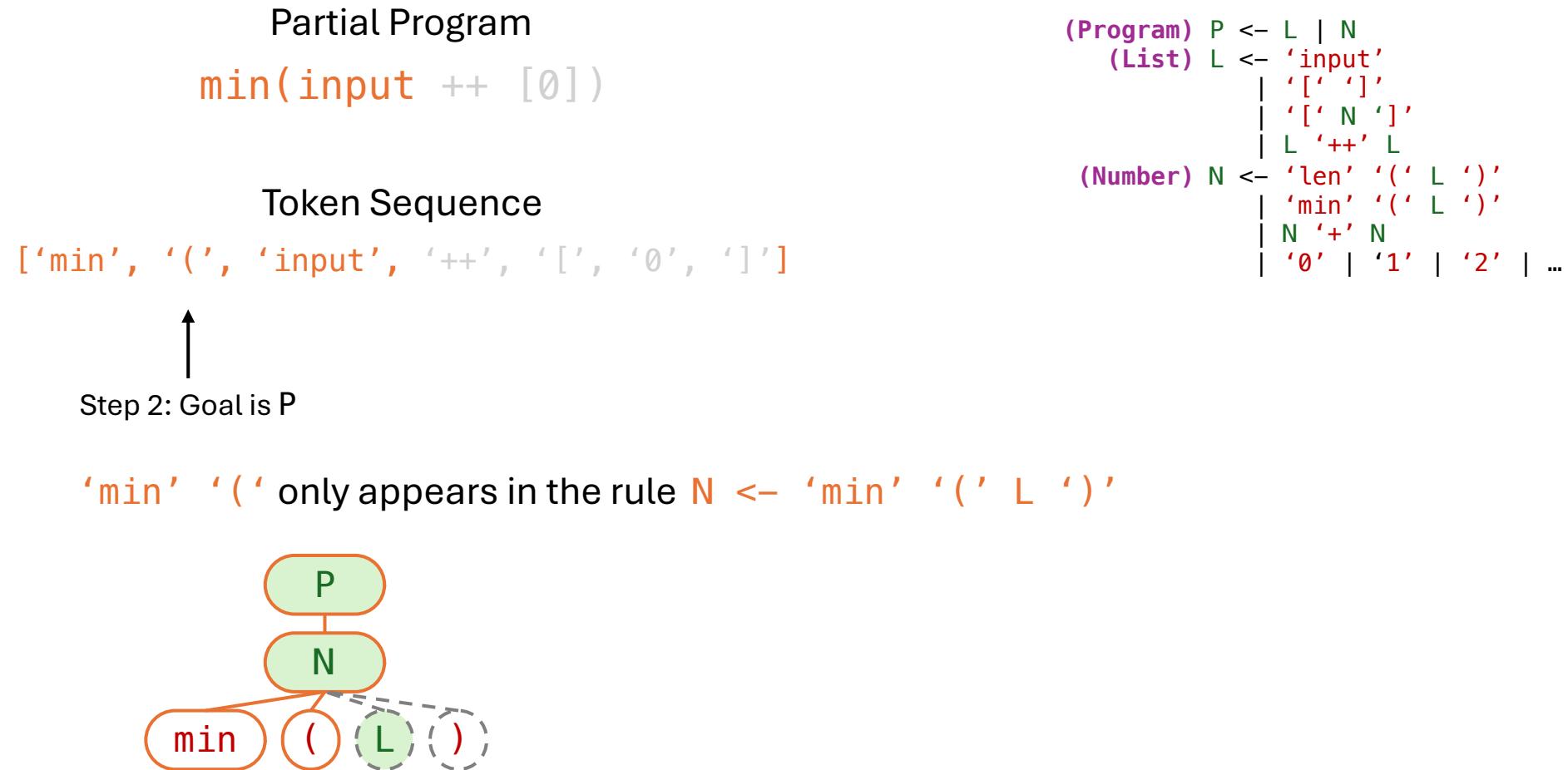
# Possible Next Token by Parser



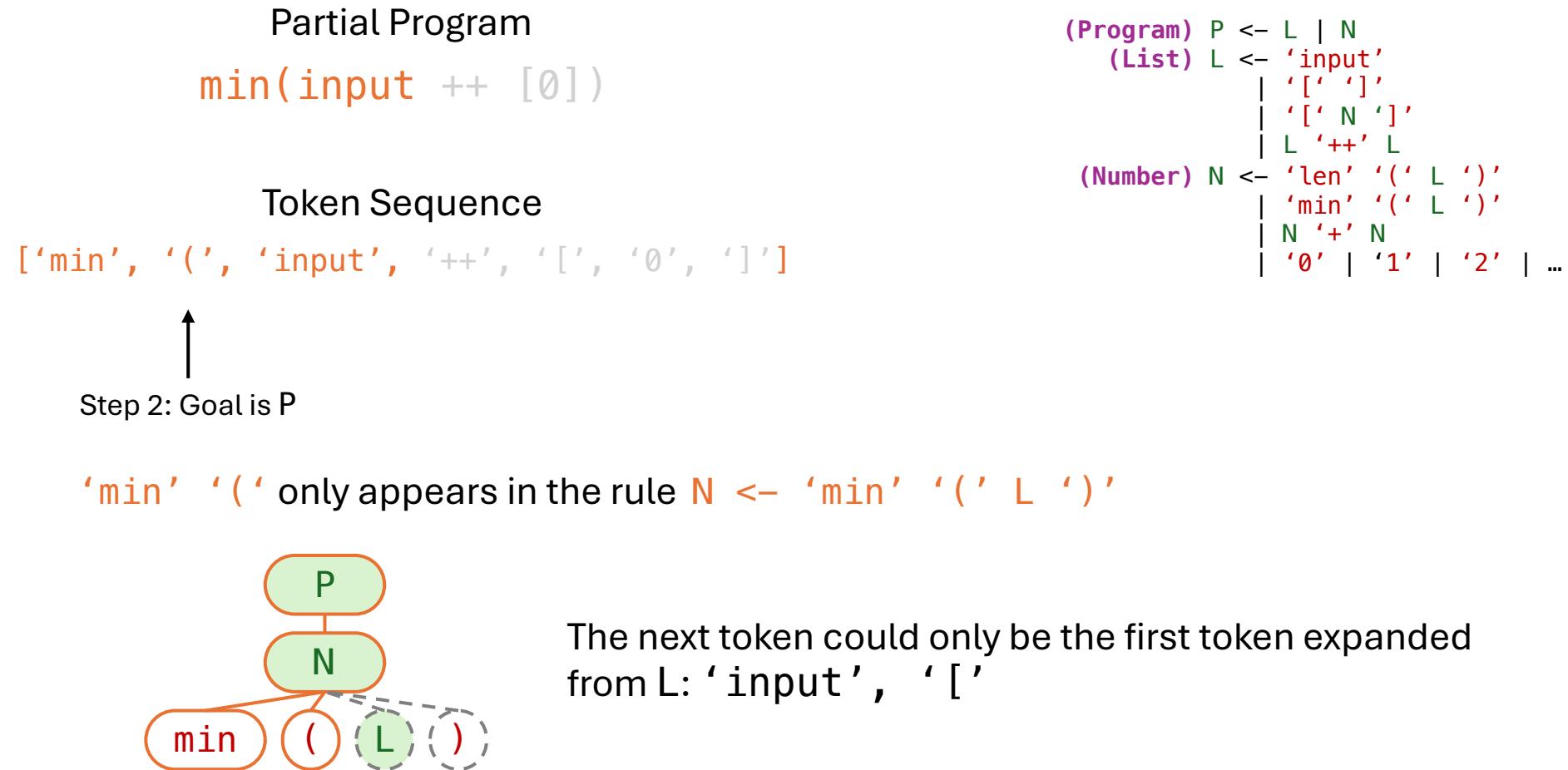
# Possible Next Token by Parser



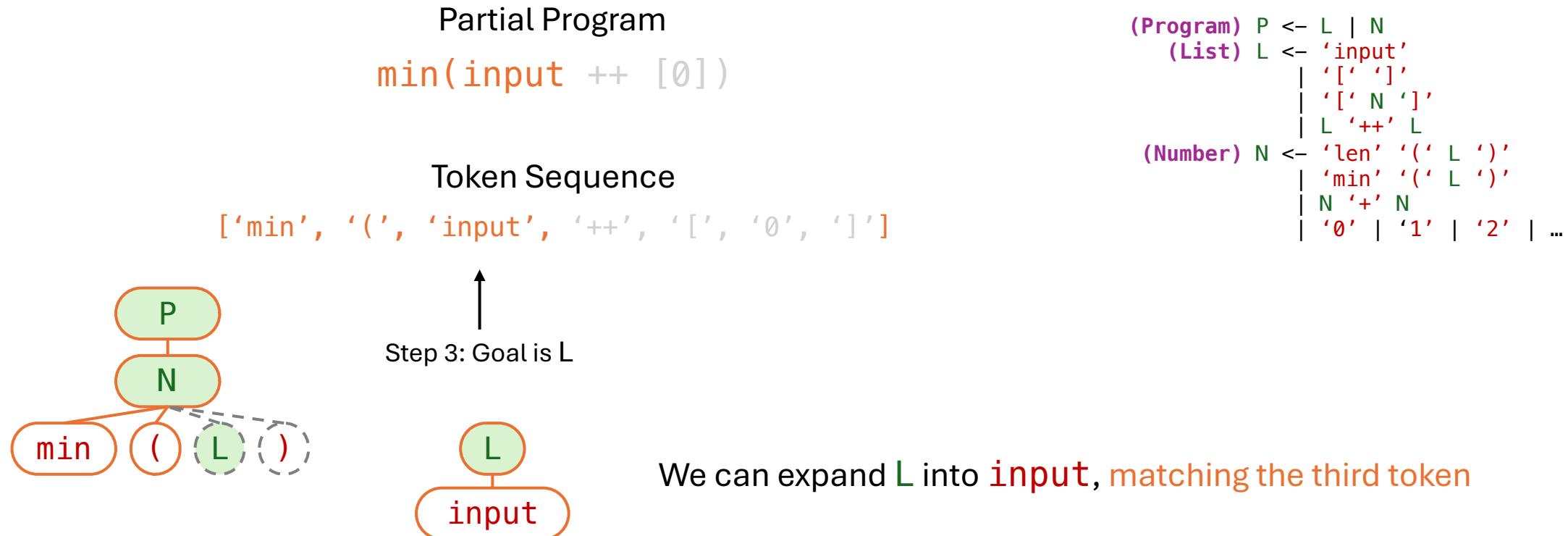
# Possible Next Token by Parser



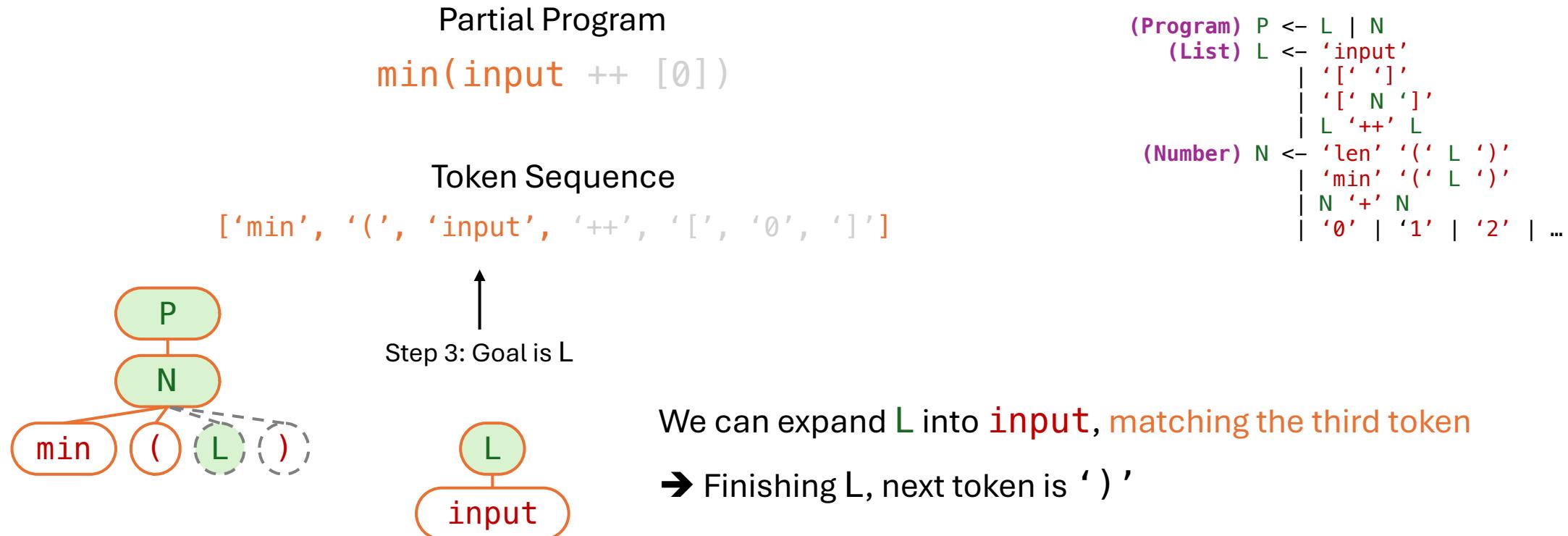
# Possible Next Token by Parser



# Possible Next Token by Parser



# Possible Next Token by Parser

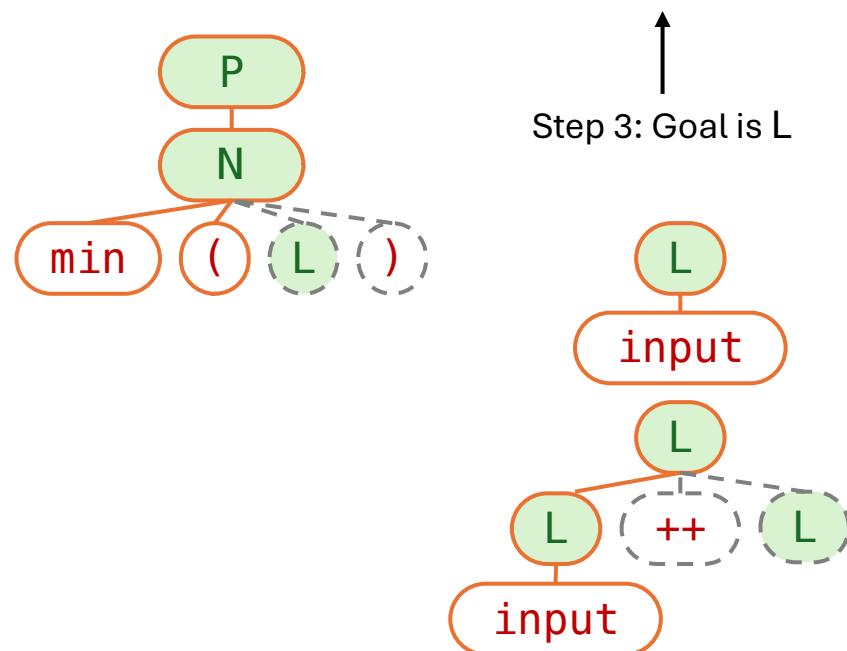


# Possible Next Token by Parser

Partial Program  
`min(input ++ [0])`

Token Sequence

`['min', '(', 'input', '++', '[', '0', ']']`



```
(Program) P <- L | N
  (List) L <- 'input'
    | '[' ']'
    | '[' N ']'
    | L '++' L
  (Number) N <- 'len' '(' L ')'
    | 'min' '(' L ')'
    | N '+' N
    | '0' | '1' | '2' | ...
```

We can expand L into `input`, matching the third token

→ Finishing L, next token is `' )'`

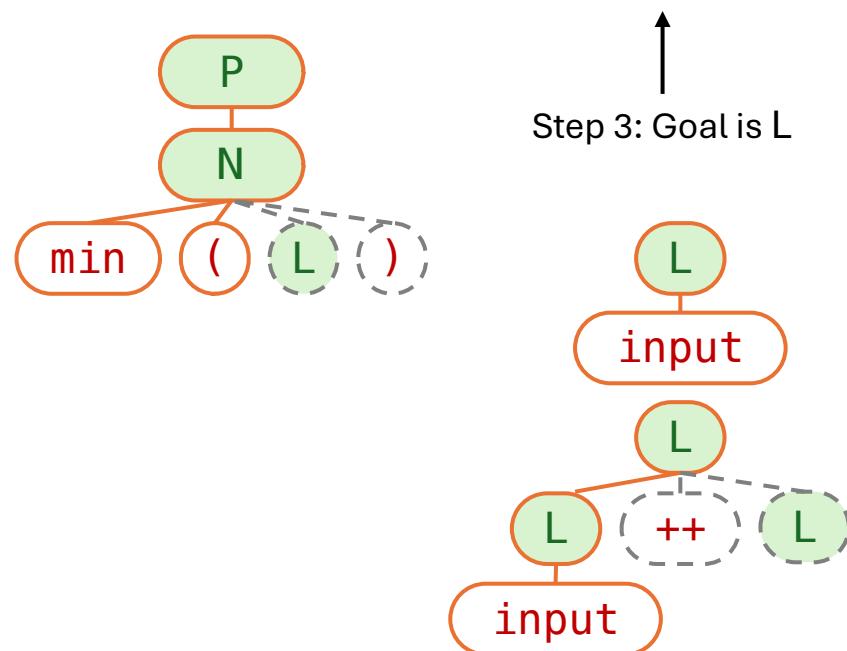
We can expand L into L + L, where the first L can be expanded into `input`, matching the third token

# Possible Next Token by Parser

Partial Program  
`min(input ++ [0])`

Token Sequence

`['min', '(', 'input', '++', '[', '0', ']']`



(Program) `P <- L | N`  
(List) `L <- 'input'`  
|  
| '[]'  
| '[' N ']'  
| L '++' L  
(Number) `N <- 'len' '(' L ')' | 'min' '(' L ')' | N '+' N | '0' | '1' | '2' | ...`

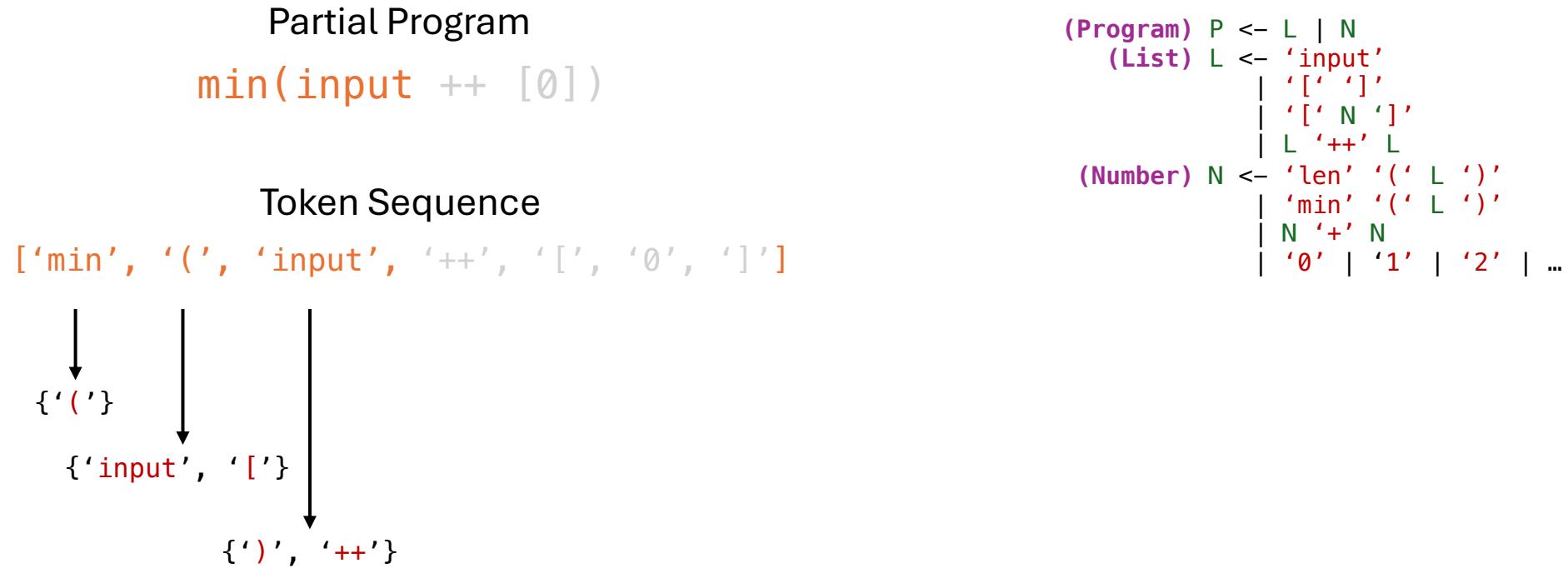
We can expand `L` into `input`, matching the third token

→ Finishing `L`, next token is `' )'`

We can expand `L` into `L + L`, where the first `L` can be expanded into `input`, matching the third token

→ Finishing the left `L`, next token is `' ++'`

# Possible Next Token by Parser



# Alternatives of Next Token Filtering

- Problem:
  - Given a **prefix**, find the **set of tokens** that could be the **next**
- Alternative problem:
  - Given a prefix and **a predicted next token**, check if the next token is a **plausible completion**
    - Prefix viability problem / Prefix membership problem
    - Given a grammar  $G$  and a prefix  $w$ , does there exist  $x$  s.t.  $wx \in L(G)$

# Algorithms for parsing

- Basic parsing
  - Earley parsing: maintains a set of items that implicitly answer the “prefix membership” question at each position
  - CKY (Cocke–Younger–Kasami) parsing: bottom-up finding substrings that can be turned into abstract syntax trees
  - Automaton: convert the grammar into an Automaton, keep an active state, reject the string if the next token leads to an error state

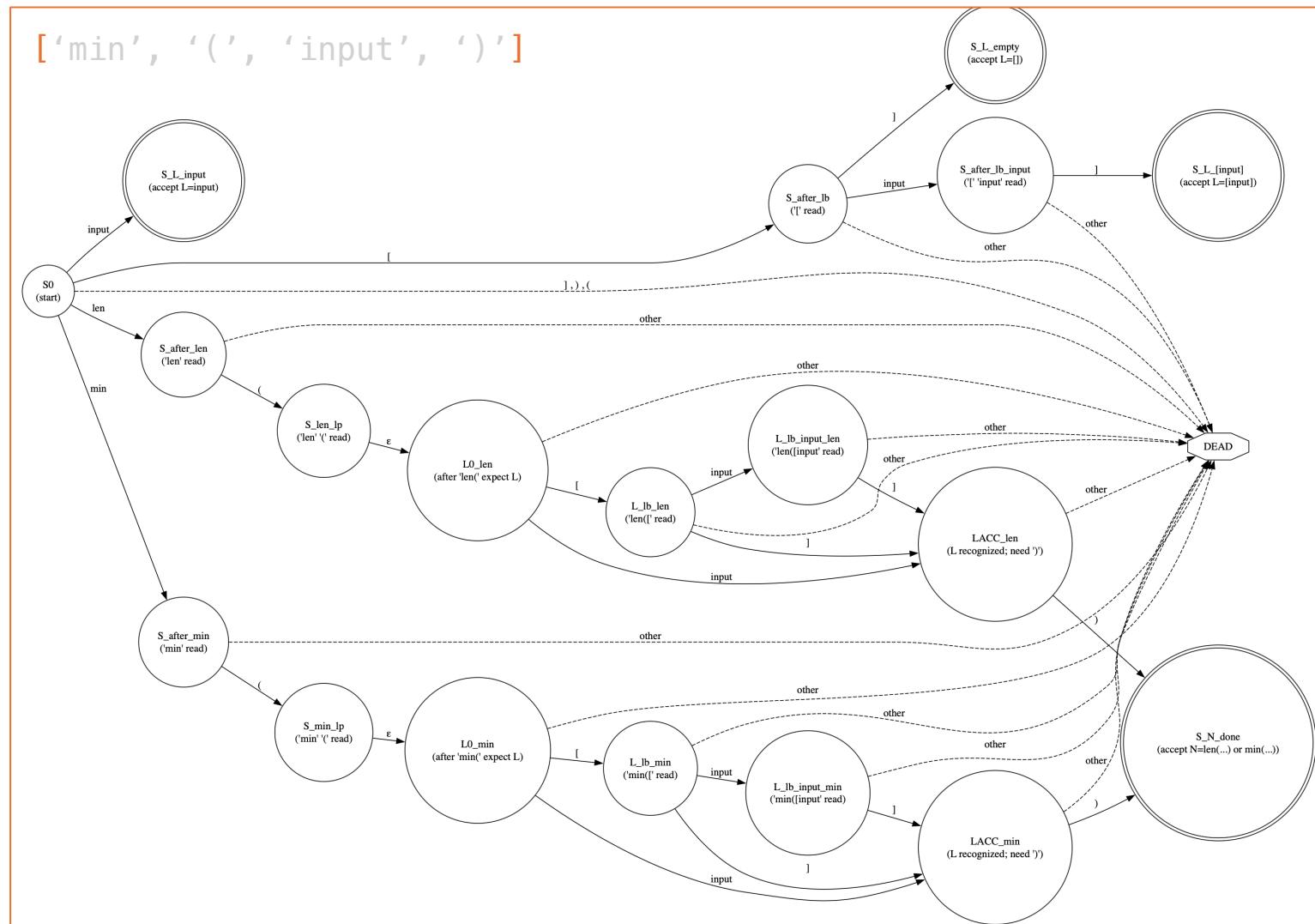
# Constrained Decoding with Automaton

```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence) / temperature
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```

# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

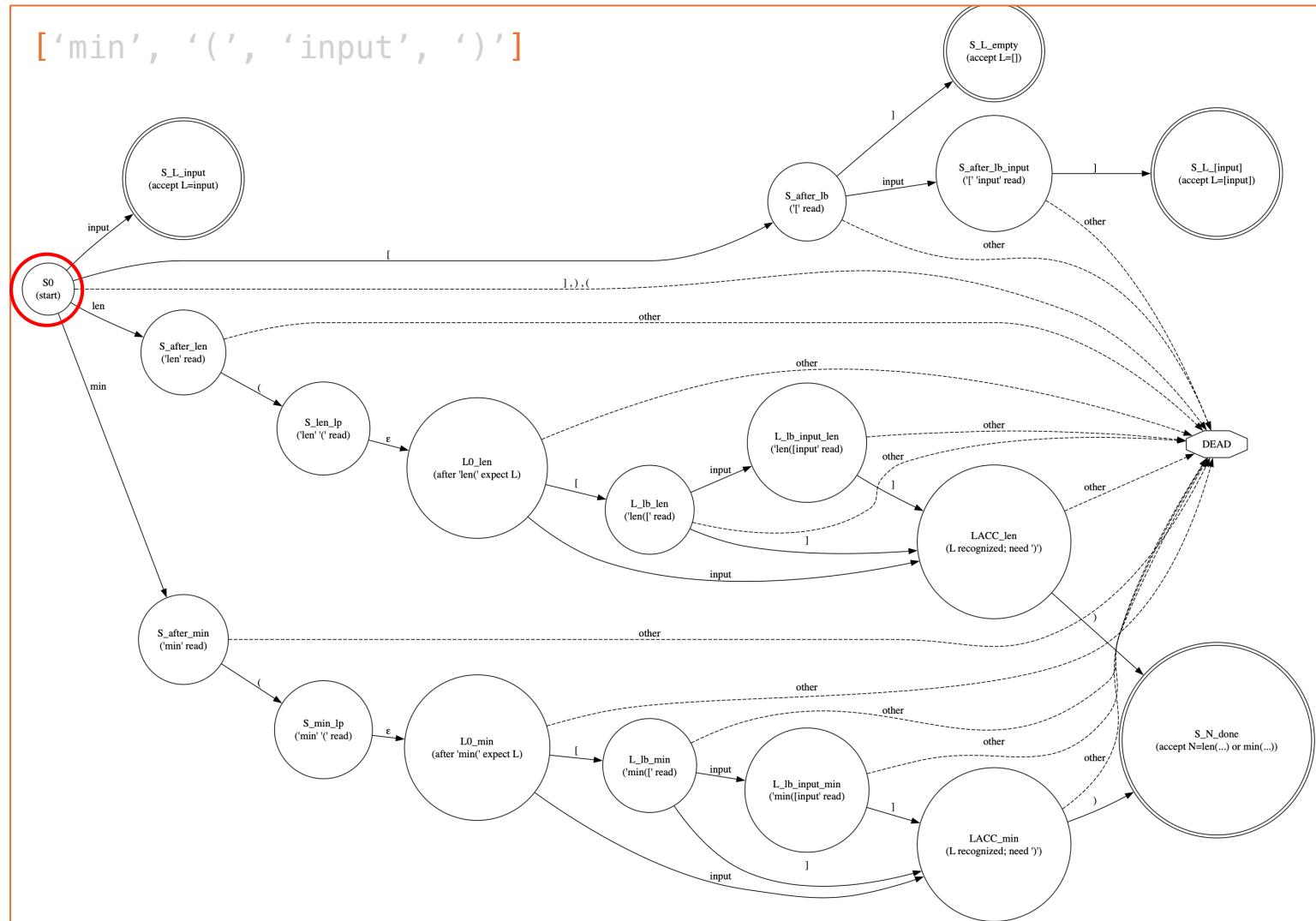
```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```



# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

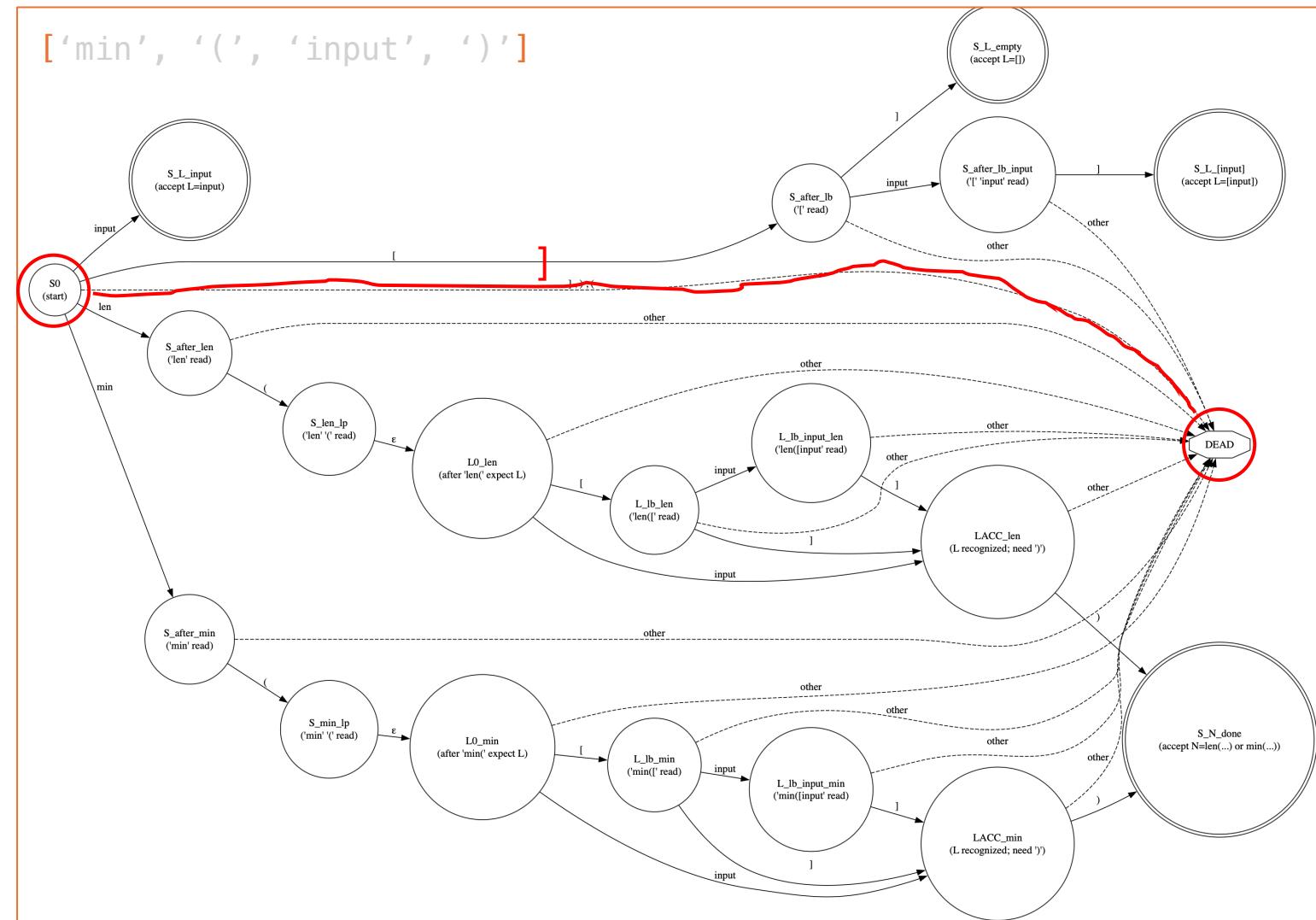
```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```



# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

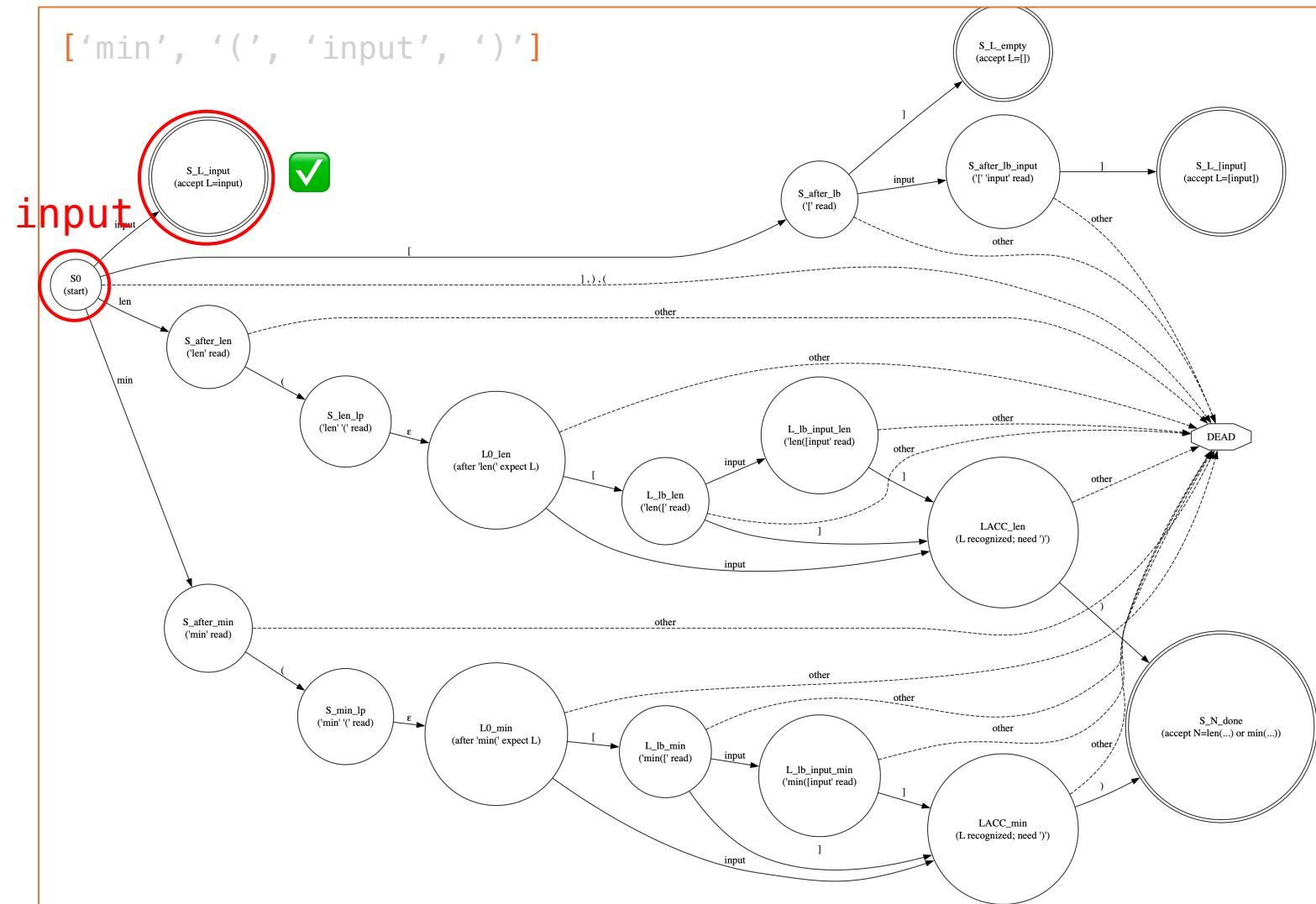
```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```



# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```

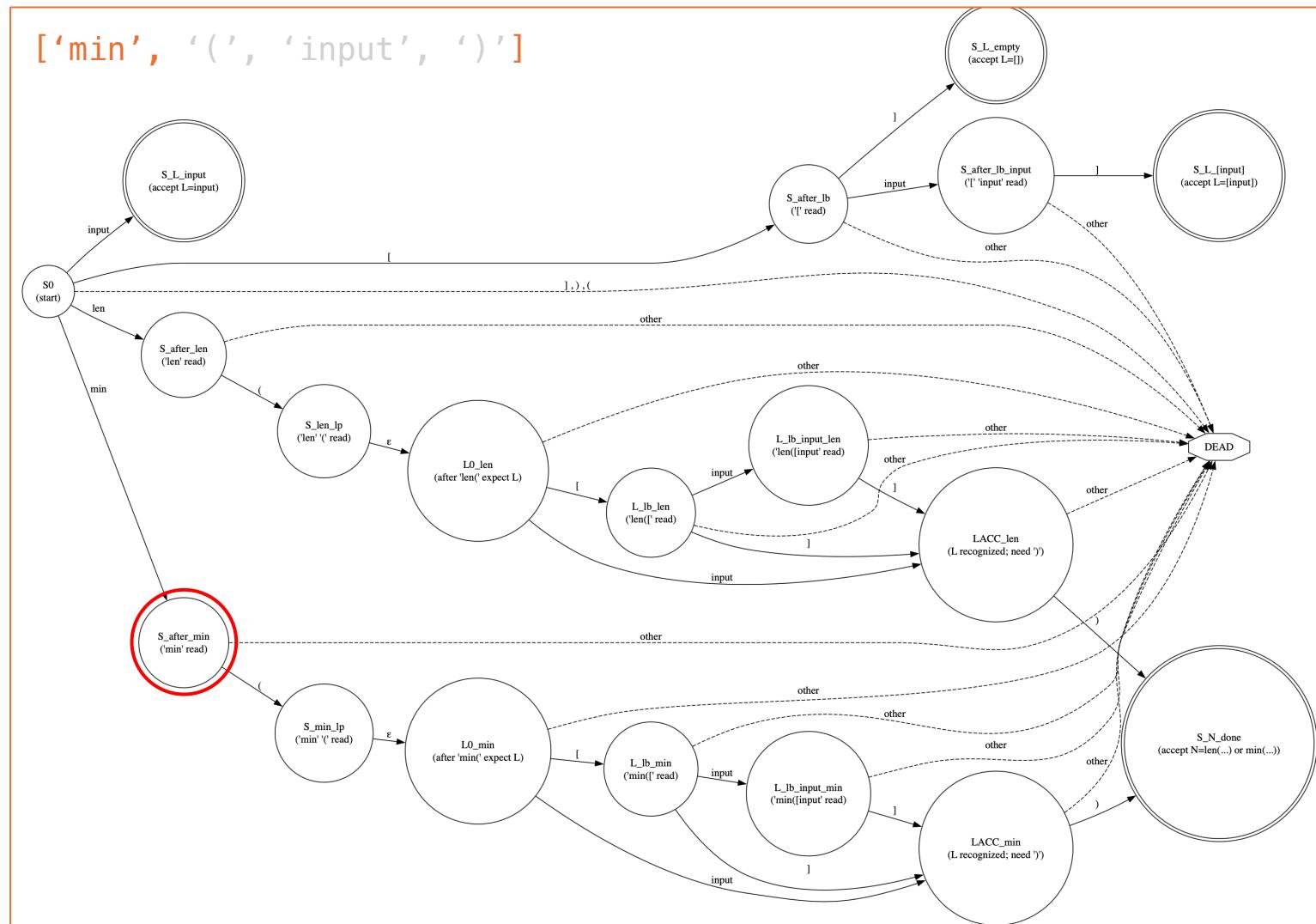


# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```

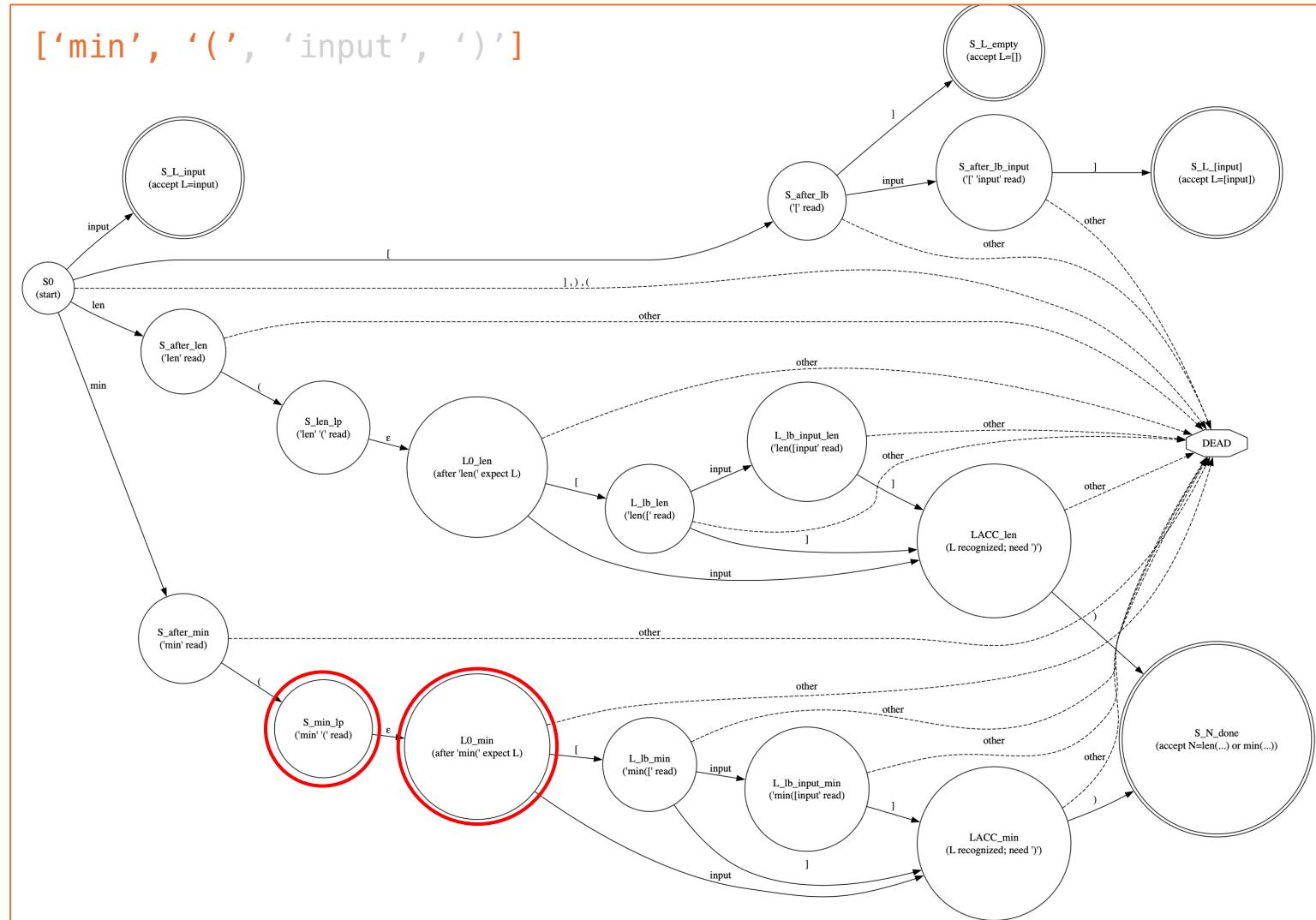
[‘min’, ‘(’, ‘input’, ‘)’]



# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

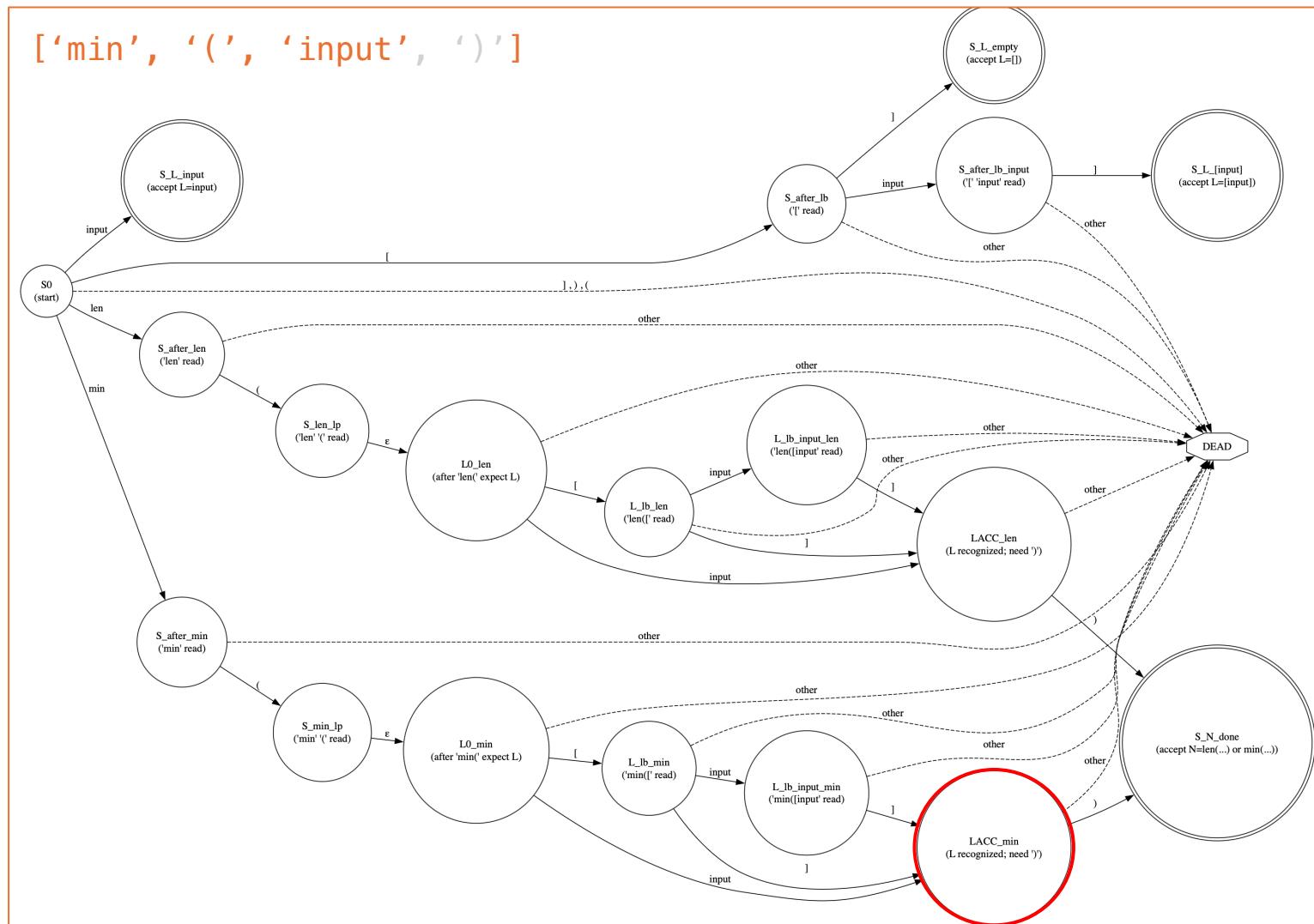
```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```



# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```

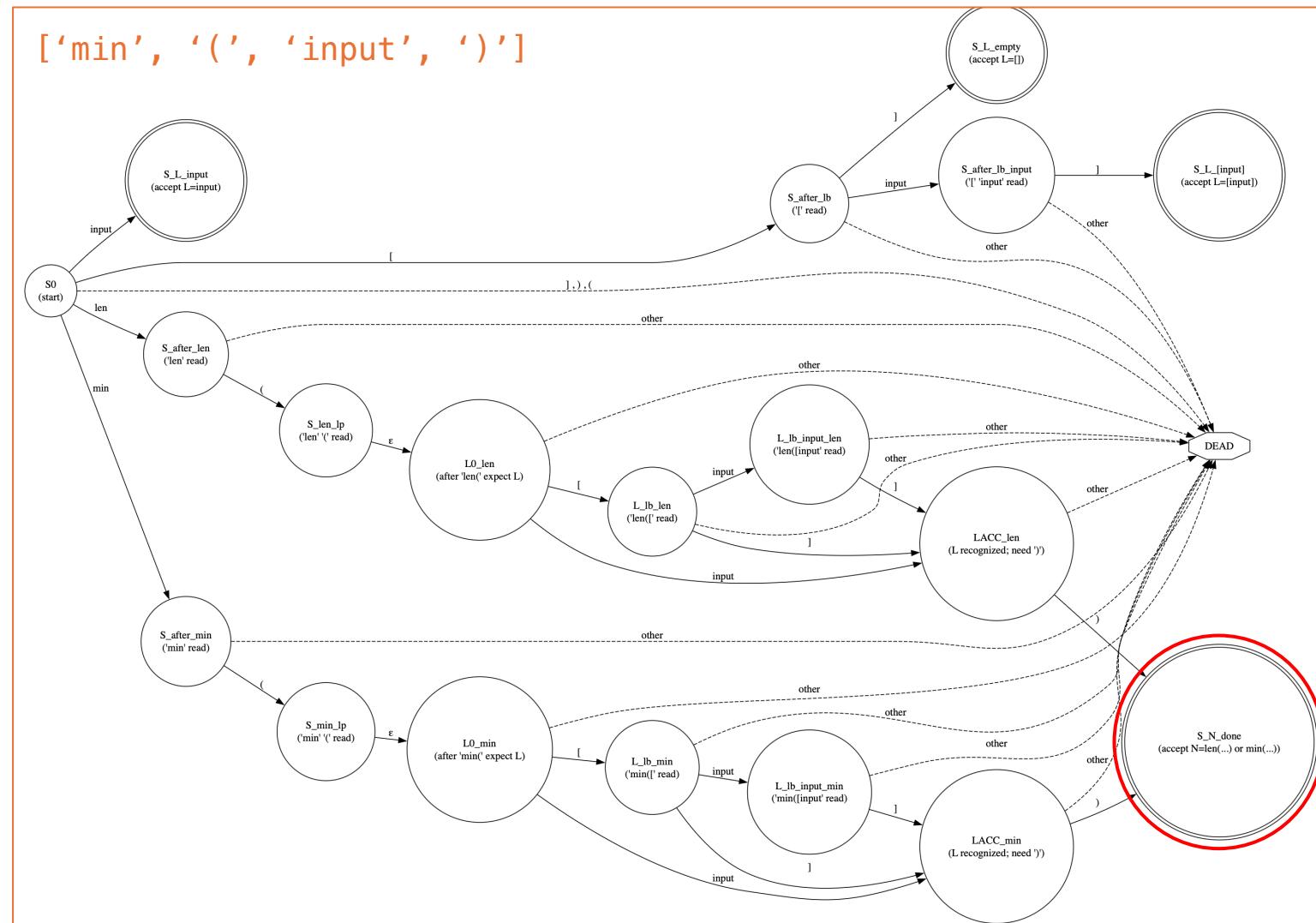


# Constrained Decoding with Automaton

```
(Program) P <- L | N
(List) L <- 'input'
| '[' ']'
| '[' N ']'
(Number) N <- 'len' '(' L ')'
| 'min' '(' L ')'
```

```
state = start_of_automaton(grammar)
current_sequence = input_tokens
output_sequence = []
for step in 1 .. max_length:
    logits = predict_next_token(current_sequence)...
    for (token_id, logit) in enumerate(logits):
        if not allowed_by(state, token_id):
            logits[token_id] = -inf
    probs = softmax(logits)
    next_token = nucleus_sample_from(probs, P)
    current_sequence += [next_token]
    output_sequence += [next_token]
    if next_token == <EOS>: break
    state = advance(state, next_token)
return output_sequence
```

['min', '(', 'input', ')']



# Algorithms for parsing

- Basic parsing
  - Earley parsing: maintains a set of items that implicitly answer the “prefix membership” question at each position
  - CKY (Cocke–Younger–Kasami) parsing: bottom-up finding substrings that can be turned into abstract syntax trees
  - Automaton: convert the grammar into an Automaton, keep an active state, reject the string if the next token leads to an error state
- Compiler Parsing
  - Batched parsing: takes the entire program and parses
- Interactive / Incremental parsing
  - Error-tolerant / resilient parsers; skip until synchronizing symbols (like ;,})
  - Keeps tree of interactions; maintain the global AST structure and only create bad subtrees when there is error

## TRANSFORMER EXPLAINER

## Examples ▾

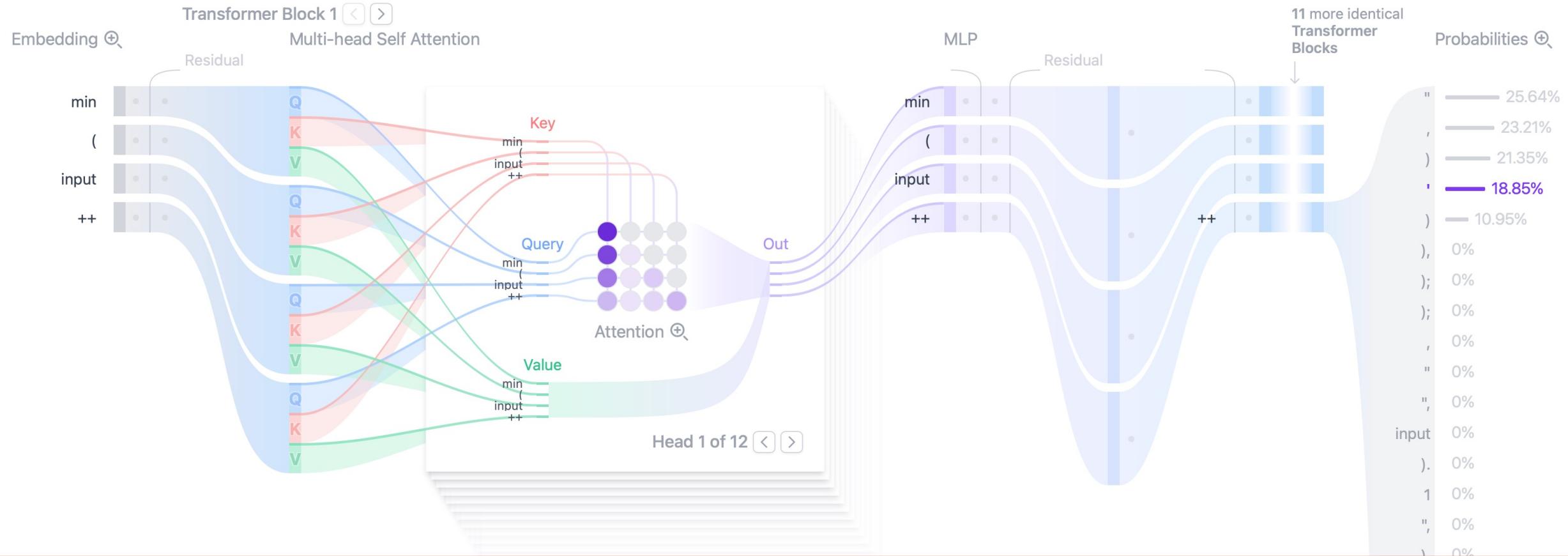
min(input ++

Generate

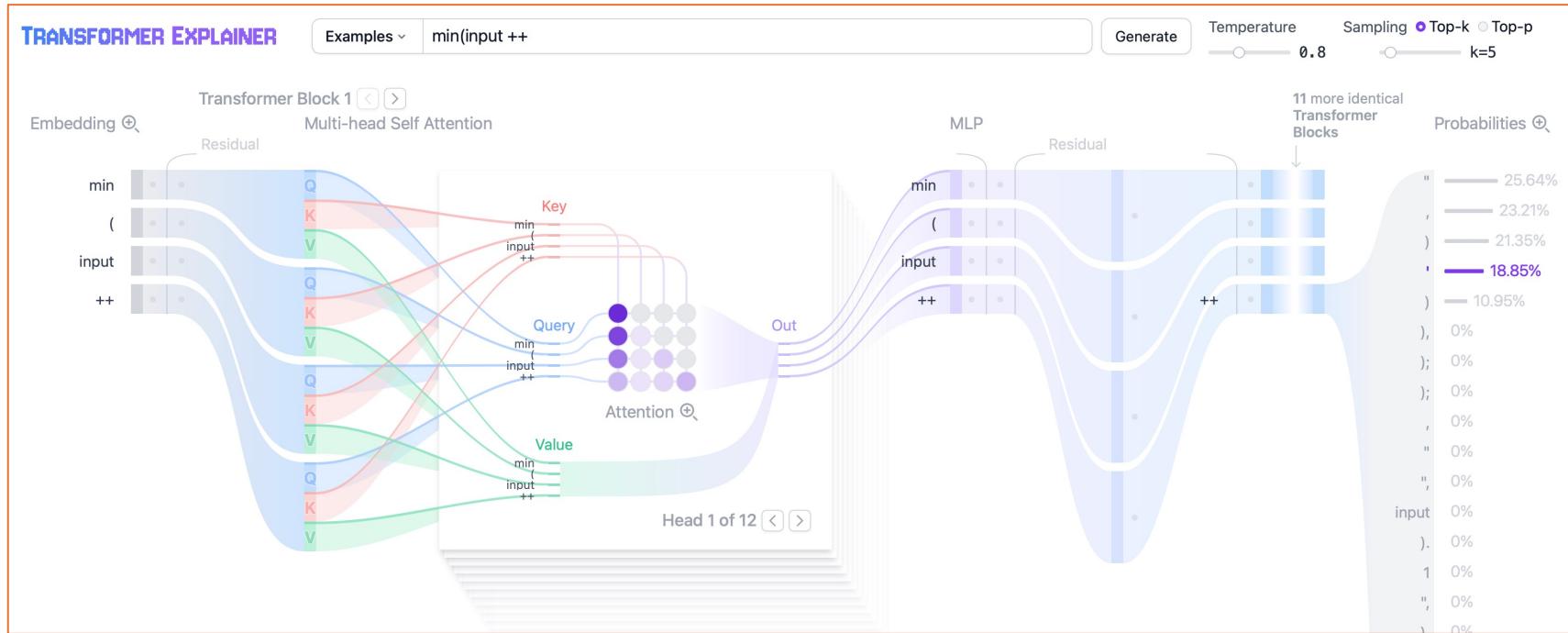
## Temperature



Sampling  Top-k  Top-p



Prefix: min(input ++



LLM Predicted Next Token Prefix Viability Mask

|       |      |   |
|-------|------|---|
| input | 0.01 | 1 |
| )     | 0.21 | 0 |
| '     | 0.19 | 0 |
| "     | 0.26 | 0 |
| ;     | 0.01 | 0 |
| [     | 0.02 | 1 |

|                      | LLM Predicted<br>Next Token | Prefix Viability<br>Mask |
|----------------------|-----------------------------|--------------------------|
|                      | input                       | 0.01                     |
|                      | )                           | 0.21                     |
|                      | '                           | 0.19                     |
| Prefix: min(input ++ | "                           | 0.26                     |
|                      | ;                           | 0.01                     |
|                      | [                           | 0.02                     |

×                            =

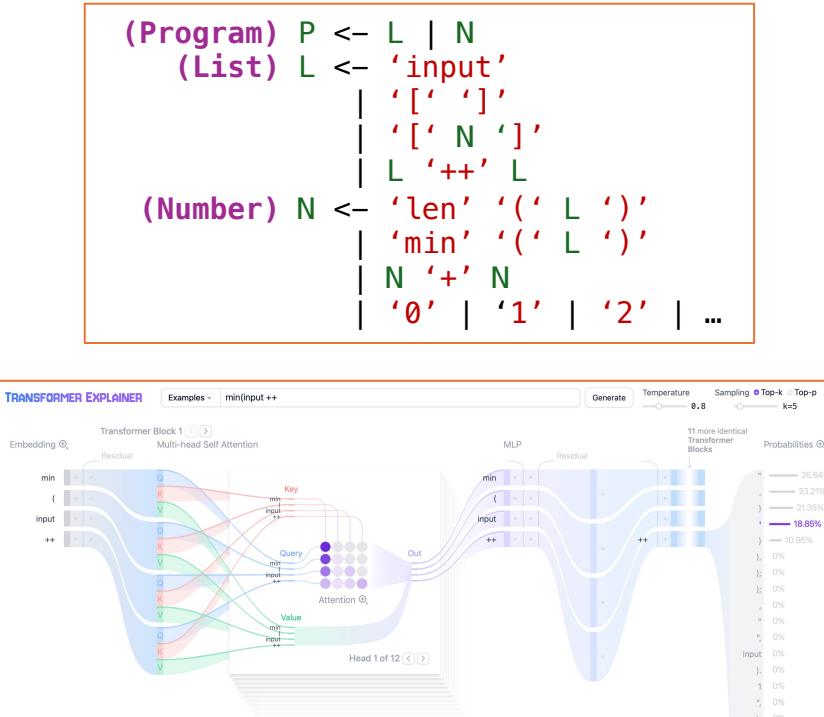
Prefix: `min(input ++`

| LLM Predicted<br>Next Token | Prefix Viability<br>Mask |
|-----------------------------|--------------------------|
| <code>input</code>          | <code>0.01</code>        |
| <code>)</code>              | <code>0.21</code>        |
| <code>'</code>              | <code>0.19</code>        |
| <code>"</code>              | <code>0.26</code>        |
| <code>;</code>              | <code>0.01</code>        |
| <code>[</code>              | <code>0.02</code>        |

`x`      `=`

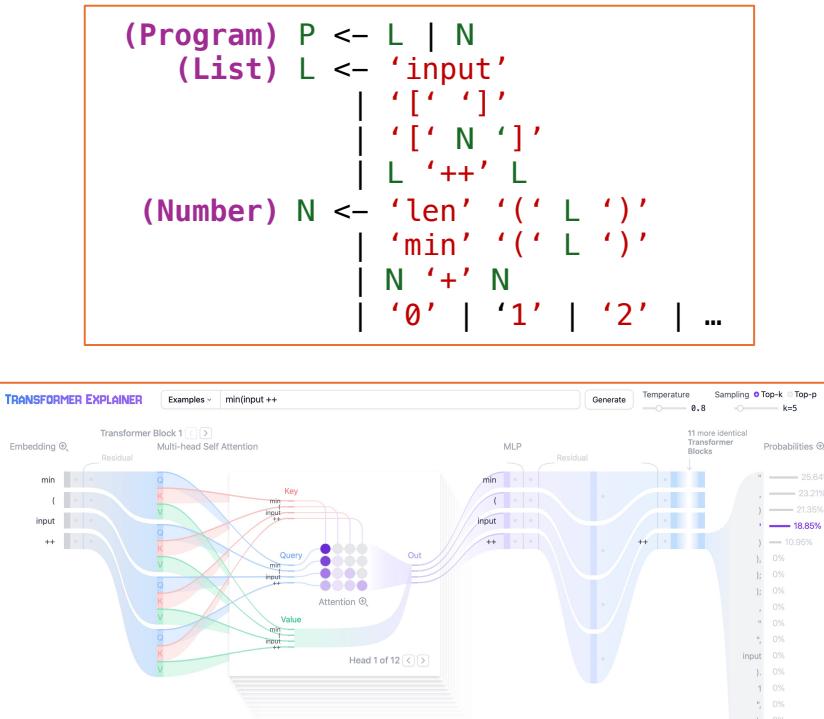
|                      | LLM Predicted<br>Next Token | Prefix Viability<br>Mask |   |      |             |
|----------------------|-----------------------------|--------------------------|---|------|-------------|
|                      | input                       | 0.01                     | 1 | 0.01 | 0.49        |
|                      | )                           | 0.21                     | 0 | -inf | 0.0         |
| Prefix: min(input ++ | '                           | 0.19                     | 0 | -inf | 0.0         |
|                      | "                           | 0.26                     | 0 | -inf | softmax 0.0 |
|                      | ;                           | 0.01                     | 0 | -inf | 0.0         |
|                      | [                           | 0.02                     | 1 | 0.02 | 0.51        |

Prefix: min(input ++



|       | LLM Predicted<br>Next Token | Prefix Viability<br>Mask |      |         |
|-------|-----------------------------|--------------------------|------|---------|
| input | 0.01                        | 1                        | 0.01 | 0.49    |
| )     | 0.21                        | 0                        | -inf | 0.0     |
| ,     | 0.19                        | 0                        | -inf | 0.0     |
| x     |                             | →                        |      | →       |
| "     | 0.26                        | 0                        | -inf | softmax |
| ;     | 0.01                        | 0                        | -inf | 0.0     |
| [     | 0.02                        | 1                        | 0.02 | 0.51    |

Prefix: min(input ++



|       | LLM Predicted<br>Next Token | Prefix Viability<br>Mask |      |         |
|-------|-----------------------------|--------------------------|------|---------|
| input | 0.01                        | 1                        | 0.01 | 0.49    |
| )     | 0.21                        | 0                        | -inf | 0.0     |
| ,     | 0.19                        | 0                        | -inf | 0.0     |
| x     |                             | →                        |      | →       |
| "     | 0.26                        | 0                        | -inf | softmax |
| ;     | 0.01                        | 0                        | -inf | 0.0     |
| [     | 0.02                        | 1                        | 0.02 | 0.51    |

# Constrained Decoding in the Wild

## A Syntactic Neural Model for General-Purpose Code Generation

**Pengcheng Yin**

Language Technologies Institute

Carnegie Mellon University

[pcyin@cs.cmu.edu](mailto:pcyin@cs.cmu.edu)

**Graham Neubig**

Language Technologies Institute

Carnegie Mellon University

[gneubig@cs.cmu.edu](mailto:gneubig@cs.cmu.edu)

# A Syntactic Neural Model for General-Purpose Code Generation

Pengcheng Yin

Language Technologies Institute  
Carnegie  
pcyin

Graham Neubig

Language Technologies Institute

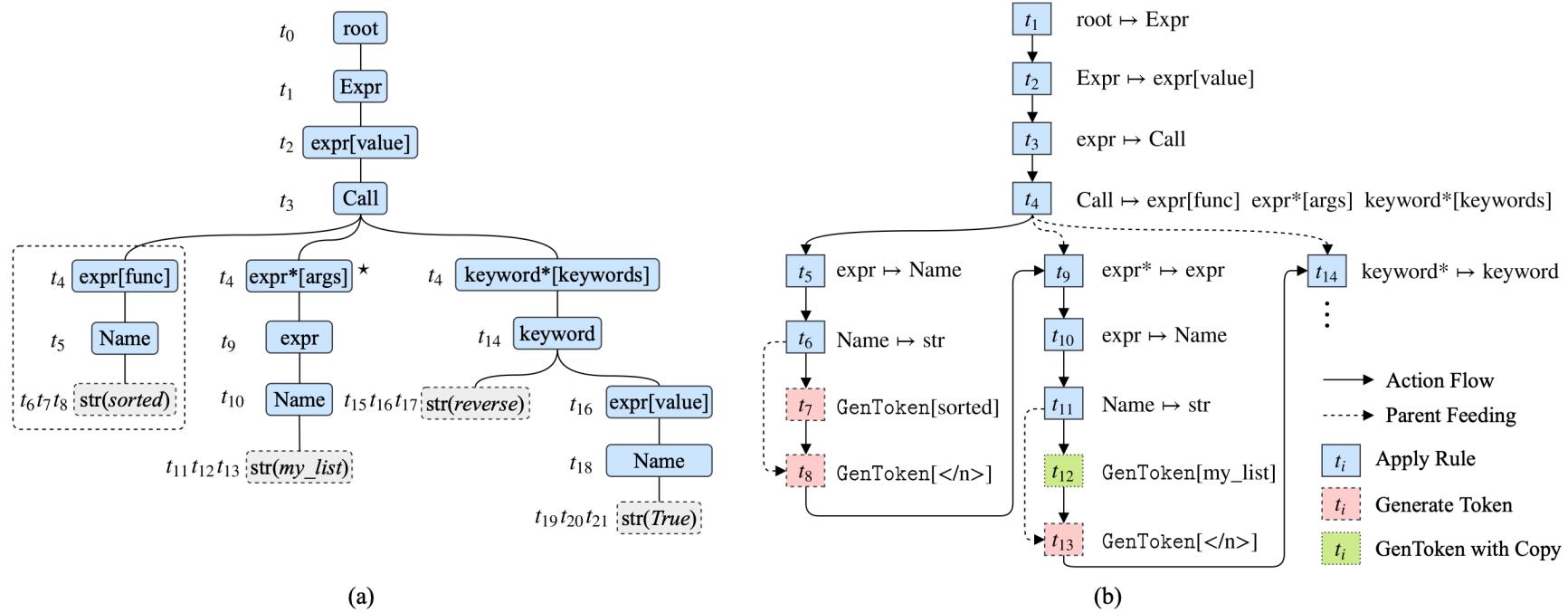


Figure 1: (a) the Abstract Syntax Tree (AST) for the given example code. Dashed nodes denote terminals. Nodes are labeled with time steps during which they are generated. (b) the action sequence (up to  $t_{14}$ ) used to generate the AST in (a)

## A Syntactic Neural Model for General-Purpose

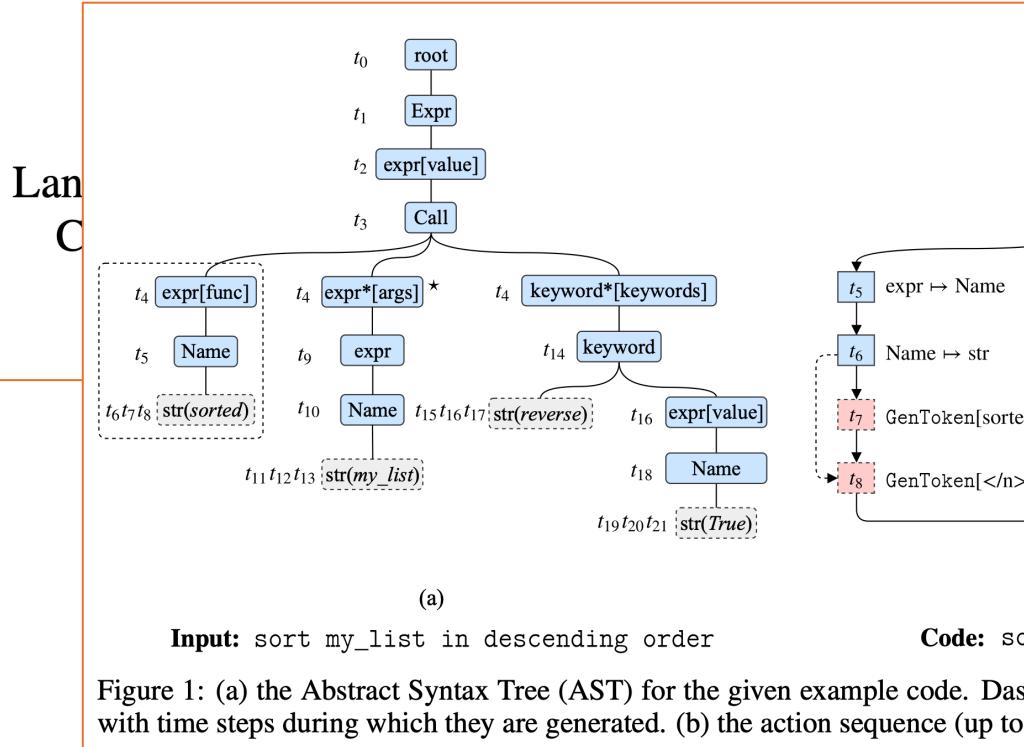


Figure 1: (a) the Abstract Syntax Tree (AST) for the given example code. Dashed boxes indicate tokens generated at specific time steps. (b) the action sequence (up to

### 4.2.2 Calculating Action Probabilities

In this section we explain how action probabilities  $p(a_t|x, a_{<t})$  are computed based on  $\mathbf{s}_t$ .

**APPLYRULE** The probability of applying rule  $r$  as the current action  $a_t$  is given by a softmax<sup>5</sup>:

$$p(a_t = \text{APPLYRULE}[r]|x, a_{<t}) = \text{softmax}(\mathbf{W}_R \cdot g(\mathbf{s}_t))^\top \cdot \mathbf{e}(r) \quad (4)$$

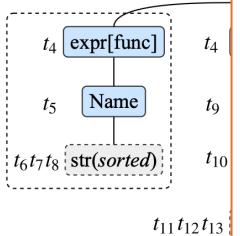
where  $g(\cdot)$  is a non-linearity  $\tanh(\mathbf{W} \cdot \mathbf{s}_t + \mathbf{b})$ , and  $\mathbf{e}(r)$  the one-hot vector for rule  $r$ .

**GENTOKEN** As in § 3.2, a token  $v$  can be generated from a predefined vocabulary or copied from the input, defined as the marginal probability:

$$\begin{aligned} p(a_t = \text{GENTOKEN}[v]|x, a_{<t}) &= \\ &p(\text{gen}|x, a_{<t})p(v|\text{gen}, x, a_{<t}) \\ &+ p(\text{copy}|x, a_{<t})p(v|\text{copy}, x, a_{<t}). \end{aligned}$$

# A Syntactic Neural Model for General-Purpose Code Generation

Lan  
C



Input: sort

Figure 1: (a) the A  
with time steps dur

## 4.2.2 Calculating Action Probabilities

In this section we explain how action probabilities

$p(a_t|x, a_{<t})$  are computed.

**APPLYRULE** The probability of applying a rule  $r$  as the current action  $a_t$  is

$$p(a_t = \text{APPLYRULE}[r]) = \text{softmax}(g(\mathbf{e}(r)))$$

where  $g(\cdot)$  is a non-linear function and  $\mathbf{e}(r)$  the one-hot vector for rule  $r$ .

**GENTOKEN** As in § 3.2, the probability of generating a token  $t$  is calculated from a predefined distribution  $\pi_t$  over tokens in the input, defined as the softmax of the logit  $\phi_t$ :

$$p(a_t = \text{GENTOKEN}[t]) = \text{softmax}(\phi_t) = p(\text{gen}|x, a_{<t}) + p(\text{copy}|x, a_{<t})$$

|  | HS          |             | DJANGO      |             |
|--|-------------|-------------|-------------|-------------|
|  | ACC         | BLEU        | ACC         | BLEU        |
| Retrieval System <sup>†</sup>            | 0.0         | 62.5        | 14.7        | 18.6        |
| Phrasal Statistical MT <sup>†</sup>      | 0.0         | 34.1        | 31.5        | 47.6        |
| Hierarchical Statistical MT <sup>†</sup> | 0.0         | 43.2        | 9.5         | 35.9        |
| NMT                                      | 1.5         | 60.4        | 45.1        | 63.4        |
| SEQ2TREE                                 | 1.5         | 53.4        | 28.9        | 44.6        |
| SEQ2TREE-UNK                             | 13.6        | 62.8        | 39.4        | 58.2        |
| LPN <sup>†</sup>                         | 4.5         | 65.6        | 62.3        | 77.6        |
| Our system                               | 16.2        | <b>75.8</b> | <b>71.6</b> | <b>84.5</b> |
| Ablation Study                           |             |             |             |             |
| - frontier embed.                        | <b>16.7</b> | <b>75.8</b> | 70.7        | 83.8        |
| - parent feed.                           | 10.6        | 75.7        | 71.5        | 84.3        |
| - copy terminals                         | 3.0         | 65.7        | 32.3        | 61.7        |
| + unary closure                          | -           | -           | 70.3        | 83.3        |
| - unary closure                          | 10.1        | 74.8        | -           | -           |

# Constrained Decoding in the Wild

**Abstract Syntax Networks for Code Generation and Semantic Parsing**

**Maxim Rabinovich\***    **Mitchell Stern\***    **Dan Klein**

Computer Science Division

University of California, Berkeley

{rabinovich,mitchell,klein}@cs.berkeley.edu

# Abstract Syntax Networks for Code Generation and Semantic Parsing

Maxim Rabinovich\*

Mitchell Stern\*

Computer Science Division

University of California, Berkeley

{rabinovich, mitchell, klein}@cs.berkeley.edu

Dan Klein



```
name: [
  'D', 'i', 'r', 'e', ' ', '',
  'W', 'o', 'l', 'f', ' ', '',
  'A', 'l', 'p', 'h', 'a'
]
cost: ['2']
type: ['Minion']
rarity: ['Common']
race: ['Beast']
class: ['Neutral']
description: [
  'Adjacent', 'minions', 'have',
  '+', '1', 'Attack', '.']
health: ['2']
attack: ['2']
durability: ['-1']

class DireWolfAlpha(MinionCard):
    def __init__(self):
        super().__init__(
            "Dire Wolf Alpha", 2, CHARACTER_CLASS.ALL,
            CARD_RARITY.COMMON, minion_type=MINION_TYPE.BEAST)
    def create_minion(self, player):
        return Minion(2, 2, auras=[
            Aura(ChangeAttack(1), MinionSelector(Adjacent())))
        ])
```

Figure 1: Example code for the “Dire Wolf Alpha” Hearthstone card.

show me the fare from ci0 to ci1

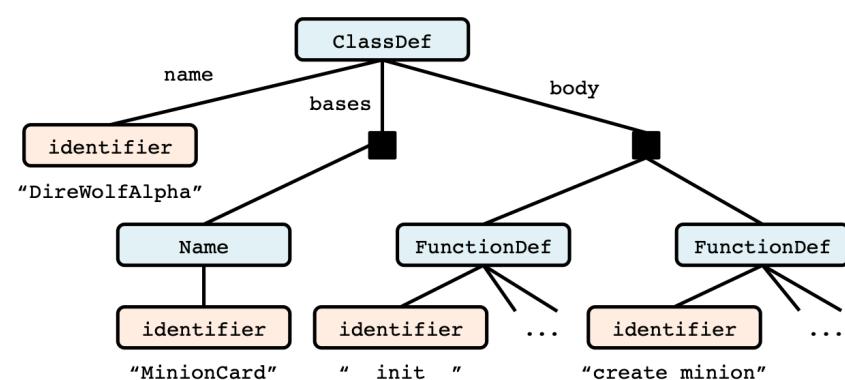
```
lambda $0 e
  ( exists $1 ( and ( from $1 ci0 )
    ( to $1 ci1 )
    ( = ( fare $1 ) $0 ) ) )
```

Figure 2: Example of a query and its logical form from the ATIS dataset. The ci0 and ci1 tokens are entity abstractions introduced in preprocessing (Dong and Lapata, 2016).

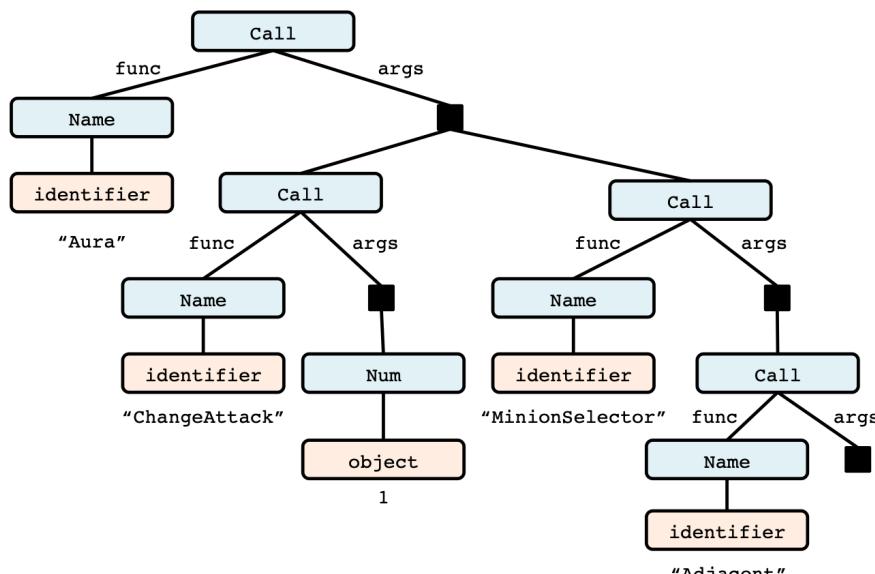
# Abstract Syntax Networks for Code Generation and Semantic Parsing

Maxim Rabinovich\*    Mitchell Stern\*  
Computer Science Division  
University of California, Berkeley

Dan Klein



(a) The root portion of the AST.



(b) Excerpt from the same AST, corresponding to the code snippet `Aura(ChangeAttack(1),MinionSelector(Adjacent()))`.



```
name: [ 'D', 'i', 'r', 'e', ' ', ' ', 'W', 'o', 'l', 'f', ' ', ' ', 'A', 'l', 'p', 'h', 'a' ] cost: ['2'] type: ['Minion'] rarity: ['Common'] race: ['Beast'] class: ['Neutral'] description: [ 'Adjacent', 'minions', 'have', '+', '1', 'Attack', '.' ] health: ['2'] attack: ['2'] durability: ['-1']
```

```
:  
CHARACTER_CLASS.ALL,  
on_type=MINION_TYPE.BEAST)  
:  
minionSelector(Adjacent())
```

ie “Dire Wolf Alpha”

0 to ci1

```
from $1 ci0 )  
, $1 ci1 )  
( fare $1 ) $0 ) ) )
```

Figure 2. Example of a query and its logical form from the ATIS dataset. The `ci0` and `ci1` tokens are entity abstractions introduced in preprocessing (Dong and Lapata, 2016).

## Abstract Syntax Networks for Code Generation and Semantic Parsing

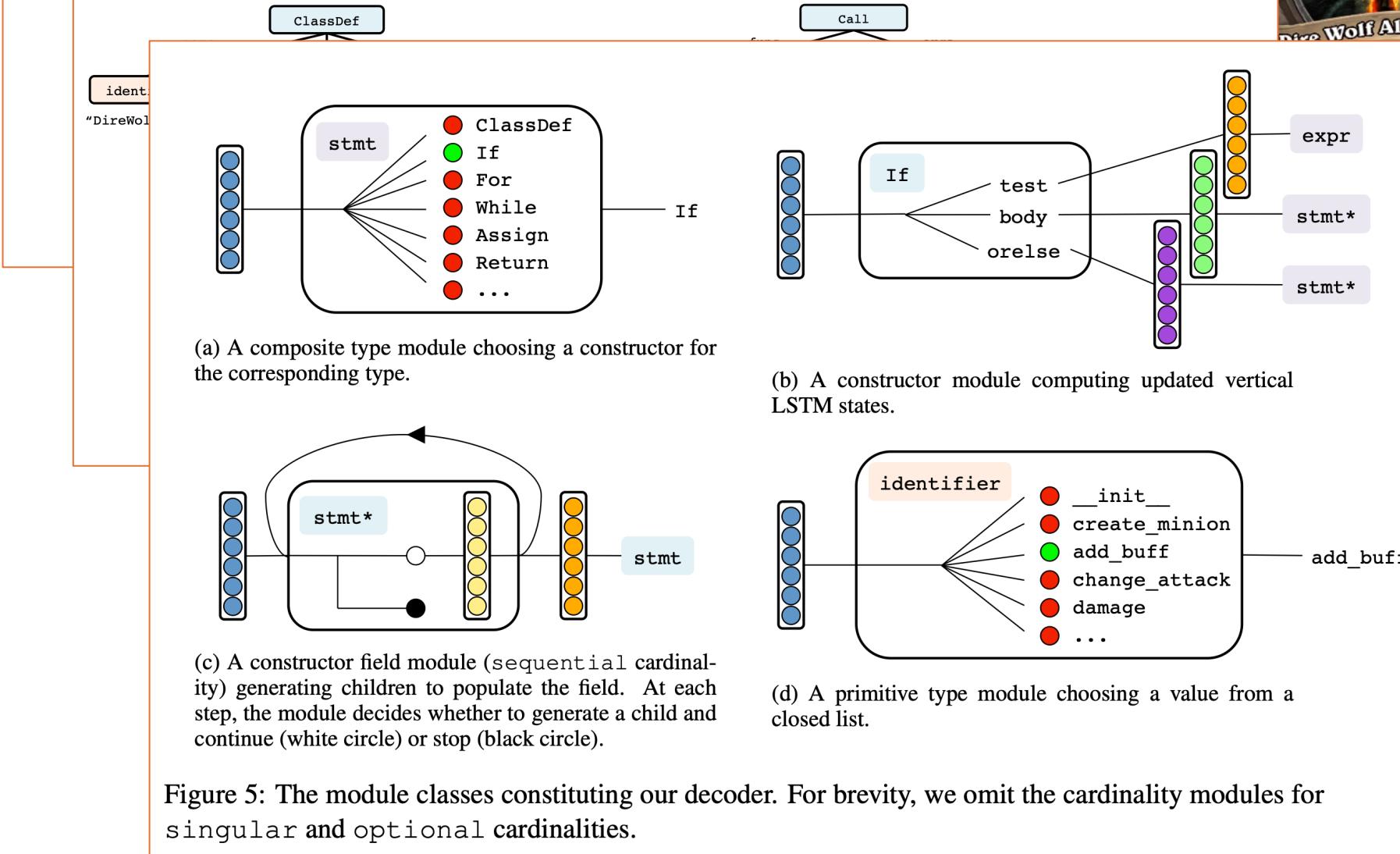
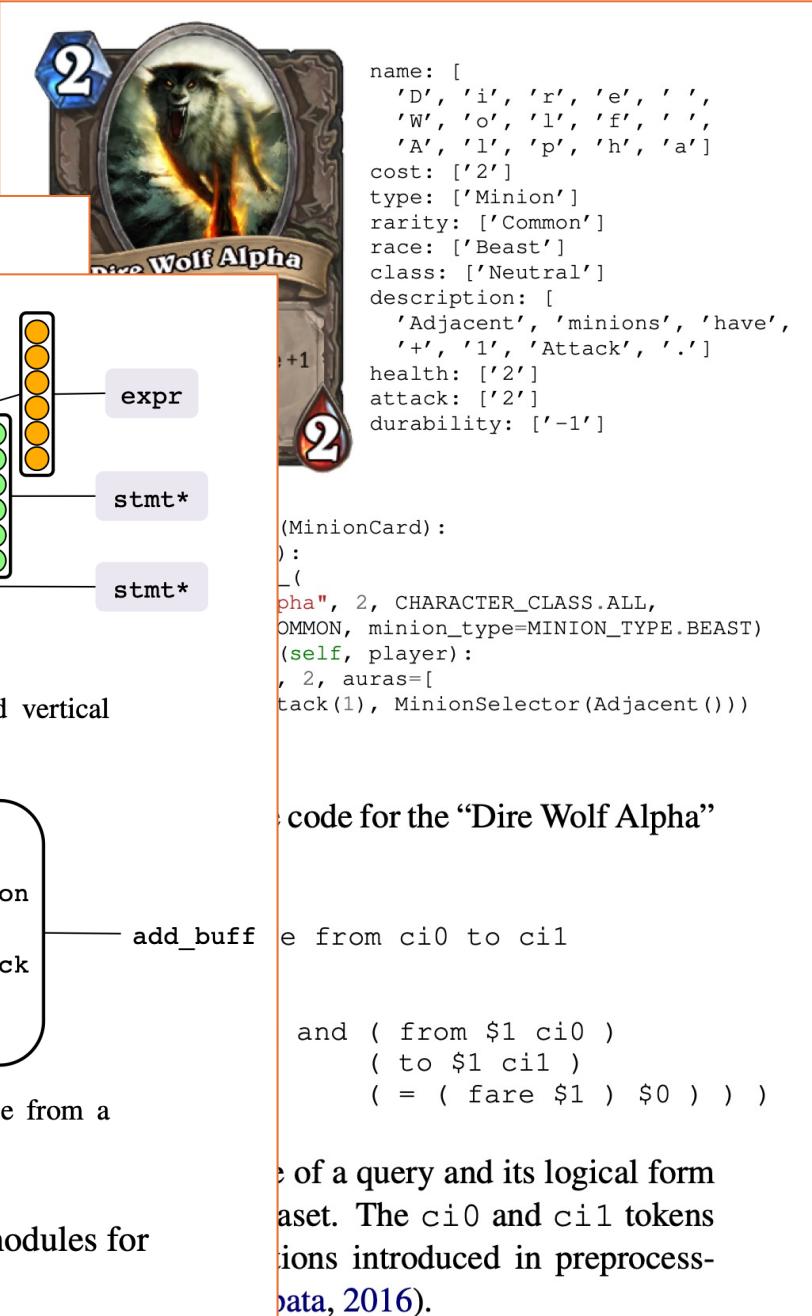


Figure 5: The module classes constituting our decoder. For brevity, we omit the cardinality modules for singular and optional cardinalities.



# Constrained Decoding in the Wild

## Grammar-Aligned Decoding

**Kanghee Park<sup>1\*</sup> Jiayu Wang<sup>1\*</sup> Taylor Berg-Kirkpatrick<sup>2</sup>**

**Nadia Polikarpova<sup>2</sup> Loris D'Antoni<sup>1</sup>**

<sup>1</sup>University of Wisconsin-Madison

{kpark247, jwang2782, ldantoni}@wisc.edu, {tberg, npolitarpova}@ucsd.edu

# Grammar-Aligned Decoding

Kanghee Park<sup>1\*</sup> Jiayu Wang<sup>1\*</sup> Taylor Berg-Kirkpatrick<sup>2</sup>

Constrained decoding addresses the inefficiency of rejection sampling by greedily “forcing” the LLM output to satisfy the given constraint. Specifically, when the constraint is given as a grammar, *grammar-constrained decoding* (GCD) [7, 27, 28], can build automata that allow for on-the-fly masking of tokens that will provably lead to outputs outside of the grammar during decoding.

**Grammar-Aligned Decoding** Given a model distribution  $P$  and a CFG  $\mathcal{G}$ , *grammar-aligned decoding* (GAD) is the task of sampling from the distribution  $Q^{P,\mathcal{G}}$  that is *proportional* to  $P$  but *restricted* to sentences in  $\mathcal{G}$ :

$$Q^{P,\mathcal{G}}(w) = \frac{\mathbb{1}[w \in \mathcal{L}(\mathcal{G})] \cdot P(w)}{\sum_{w'} \mathbb{1}[w' \in \mathcal{L}(\mathcal{G})] \cdot P(w')}$$

# Grammar-Aligned Decoding

Kanghee Park<sup>1\*</sup> Jiayu Wang<sup>1\*</sup> Taylor Berg-Kirkpatrick<sup>2</sup>

Constrained decoding addresses the inefficiency of rejection sampling by greedily “forcing” the LLM output to satisfy the given constraint. Specifically, when the constraint is given as a grammar, *grammar-constrained decoding* (GCD) [7, 27, 28], can build automata that allow for on-the-fly masking of tokens that will provably lead to outputs outside of the grammar during decoding.

**Grammar-Aligned Decoding** Given a model distribution  $P$  and a CFG  $\mathcal{G}$ , *grammar-aligned decoding (GAD)* is the task of sampling from the distribution  $Q^{P,\mathcal{G}}$  that is *proportional* to  $P$  but restricts the sampled strings to be in  $\mathcal{G}$ .

```
Start ::= S
        S ::= name | " " | "."
              | str.++ S S | str.at S I
              | str.replace S S S
              | str.substr S I I
        I ::= 0 | 1 | 2 | + I I | - I I
              | str.len S | str.indexof S S I
```

(c) Grammar for `f`

```
Start ::= BV
        BV ::= s | t
              | #x0 | #x7 | #x8
              | bvneg BV | bvnot BV
              | bvadd BV BV | bvsup BV BV
              | bvard BV BV | bvlor BV BV
              | bvlshl BV BV | bvlshr BV BV
```

(d) Grammar for `inv`

# Constrained Decoding in the Wild

---

## Monitor-Guided Decoding of Code LMs with Static Analysis of Repository Context

---

**Lakshya A Agrawal**

Microsoft Research

Bangalore, India

t-lakagrawal@microsoft.com

**Aditya Kanade**

Microsoft Research

Bangalore, India

kanadeaditya@microsoft.com

**Navin Goyal**

Microsoft Research

Bangalore, India

navingo@microsoft.com

**Shuvendu K. Lahiri**

Microsoft Research

Redmond, United States

shuvendu.lahiri@microsoft.com

**Sriram K. Rajamani**

Microsoft Research

Bangalore, India

sriram@microsoft.com

# Monitor-Guided Decoding of Code LMs with Static Analysis of Repository Context

Lakshya A Agrawal  
Microsoft Research  
Bangalore, India

t-lakagrwal@microsoft.com

Aditya Kanade  
Microsoft Research  
Bangalore, India

kanadeaditya@microsoft.com

**Method to be completed**

```
private ServerNode parseServer(String url) {
    Preconditions.checkNotNull(url);
    int start = url.indexOf(str:"/") + 2;
    int end = url.lastIndexOf(str:"?") == -1 ?
        url.length() : url.lastIndexOf(str:"?");
    String str = url.substring(start, end);
    String [] arr = str.split(regex:":");

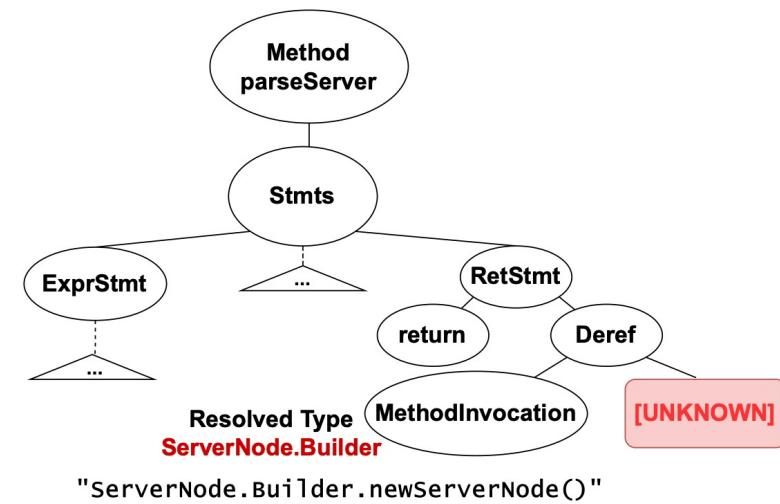
    return ServerNode.Builder
        .newServerNode()
        .
}
```

**✗ text-davinci-003 and SantaCoder**

```
host(arr[0])
    .port(Integer.parseInt(arr[1]))
    .build();
```

**✓ SantaCoder with monitor guided decoding**

```
withIp(arr[0])
    .withPort(Integer.parseInt(arr[1]))
    .build();
```



(a) Example where text-davinci-003 and SantaCoder generate wrong identifiers, but SantaCoder with MGD generates correct identifiers.

(b) Annotated partial AST for the code to the left.

## Monitor-Guided Decoding of Code LMs with Static

Method

Integrated development environments (IDEs) have been at the forefront of assisting developers. Our inspiration is the use of static analysis by IDEs to bring the global context at the fingertips of developers. Many analyses are integrated in IDEs (Fuhrer, 2013) to infer and enforce semantic constraints on the code under development, e.g., resolving def-use, symbol references, and type hierarchies. Recently, there has been a rise in the use of Language Server Protocol (LSP) (lsp), which is an open industry standard of communication between IDEs and programming language specific tools like static analyzers and compilers, called Language Servers. There are a large number of Language Servers available, targeting most programming languages (lan, 2023), and providing a variety of syntactic and semantic information. In this work, we focus on the type-directed code completion analysis available through LSP in a language-agnostic manner, to provide guidance to an LM.

(a) E  
iden

## Monitor-Guided Decoding of Code LMs with Static

Method

Integrated development environments (IDEs) have been at the forefront of assisting developers. Our inspiration is the use of static analysis by IDEs to bring the global context of the functioning

Eq. (1) states that whenever the monitor is in the wait state  $s_0$ , we sample  $x_{n+1}$  as per the logits  $\ell$  determined by the LM (Eq. (2)). Otherwise, the logits are combined with a mask  $m$  using a function  $\oplus$  such that if  $m[x] = 0$  then  $\ell[x]$  is reset to a large negative value  $-K$  and is left unchanged otherwise. This mask is computed by the function `maskgen` in Eq. (3) guided by the current state  $s$  of the monitor. Eq. (4) defines how the state of the monitor evolves. When the pre-condition `pre`( $s; x_1, \dots, x_n$ ) evaluates to true, the next state  $s'$  of the monitor is determined by the suggestions returned by the static analysis  $A_\varphi$ . Otherwise, it is determined by the `update` function.

$$(L_\theta || M_\varphi)(x_{n+1} | x_1, \dots, x_n; C, p, s) = \begin{cases} \text{softmax}(\ell)[x_{n+1}] & \text{if } s = s_0 \text{ is the wait state} \\ \text{softmax}(\ell \oplus m)[x_{n+1}] & \text{otherwise} \end{cases} \quad (1)$$

$$\ell = L_\theta(\cdot | x_1, \dots, x_n; p) \quad (2)$$

$$m = \text{maskgen}(s, V) \quad (3)$$

$$s' = \begin{cases} A_\varphi(x_1, \dots, x_n; C) & \text{if } s = s_0 \wedge \text{pre}(s; x_1, \dots, x_n) \\ \text{update}(s, x_{n+1}) & \text{otherwise} \end{cases} \quad (4)$$

(a) E  
iden  
com  
LM.

# SYNTACTIC AND SEMANTIC CONTROL OF LARGE LANGUAGE MODELS VIA SEQUENTIAL MONTE CARLO

**João Loula<sup>\*1</sup> Benjamin LeBrun<sup>\*5</sup> Li Du<sup>\*6</sup> Ben Lipkin<sup>1</sup> Clemente Pasti<sup>2</sup> Gabriel Grand<sup>1</sup>**

**Tianyu Liu<sup>2</sup> Yahya Emara<sup>2</sup> Marjorie Freedman<sup>8</sup> Jason Eisner<sup>6</sup> Ryan Cotterell<sup>2</sup>**

**Vikash Mansinghka<sup>†1</sup> Alexander K. Lew<sup>‡1,7</sup> Tim Vieira<sup>‡2</sup> Timothy J. O'Donnell<sup>‡3,4,5</sup>**

<sup>1</sup>MIT <sup>2</sup>ETH Zürich <sup>3</sup>McGill <sup>4</sup>Canada CIFAR AI Chair <sup>5</sup>Mila <sup>6</sup>Johns Hopkins <sup>7</sup>Yale <sup>8</sup>ISI  
genlm@mit.edu

# SYNTACTIC AND SEMANTIC CONTROL OF LARGE LANGUAGE MODELS VIA SEQUENTIAL MONTE CARLO

João Loula<sup>\*1</sup> Be  
Tianyu Liu<sup>2</sup> Yah  
Vikash Mansingh  
<sup>1</sup>MIT <sup>2</sup>ETH Zürich  
genlm@mit.edu

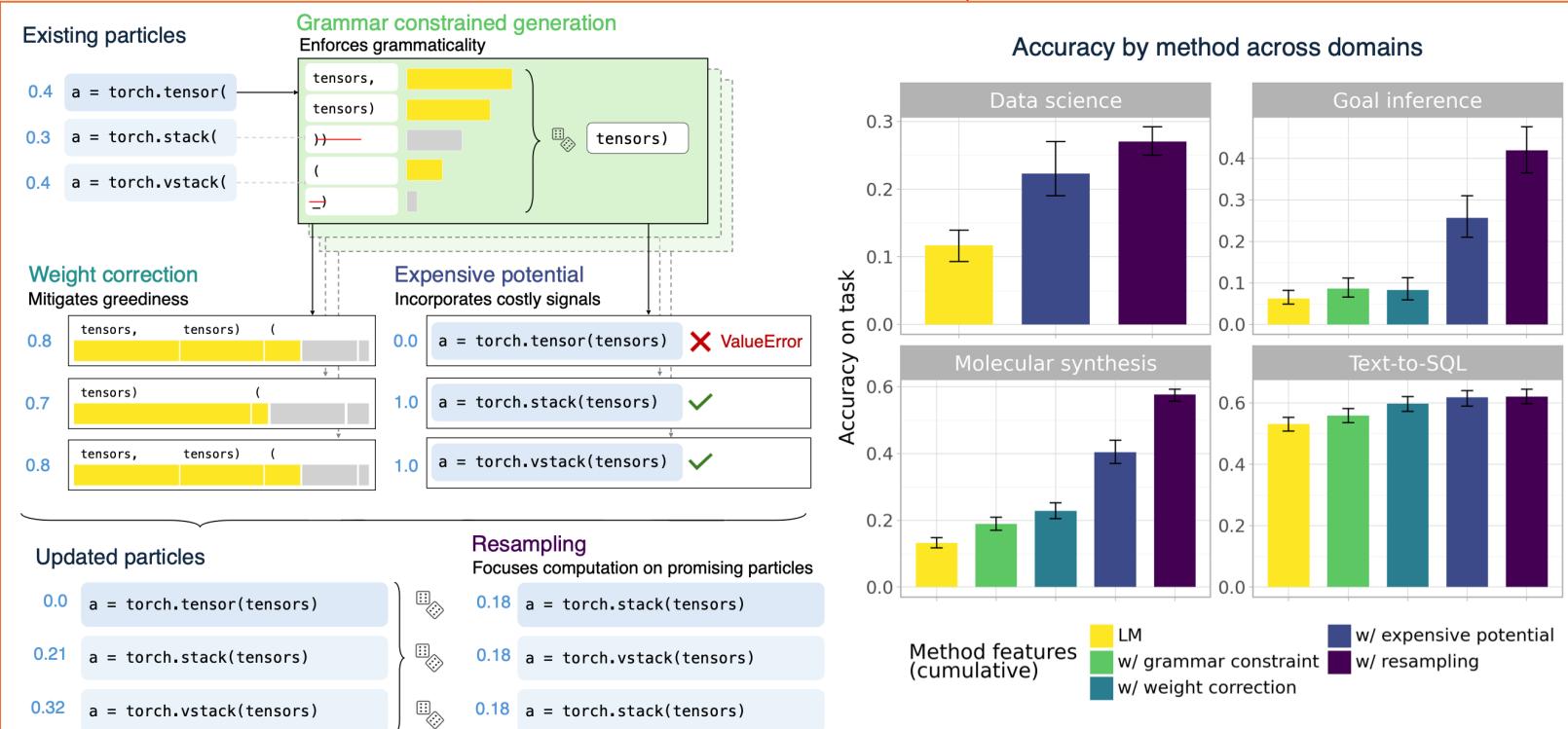


Figure 1: *Controlled generation from LMs via sequential Monte Carlo.* Left: We use sequential Monte Carlo to sample from high-quality approximations to posteriors over LM outputs. Partial sequences are repeatedly extended via **grammar-constrained generation**. We then apply **weight corrections** to mitigate the greediness of locally constrained decoding, as well as **expensive potentials** to encode rich information that cannot be included in logit masks. Finally, **resampling** focuses computation on promising particles. Right: Accuracy gains from these innovations on challenging data science, text-to-SQL, goal inference, and molecule synthesis benchmarks.

# SYNTACTIC AND SEMANTIC CONTROL OF LARGE LANGUAGE MODELS

João Loula<sup>\*1</sup> Benoît Guillet<sup>2</sup>  
 Tianyu Liu<sup>2</sup> Yalun Wang<sup>2</sup>  
 Vikash Mansinghka<sup>1</sup>  
<sup>1</sup>MIT <sup>2</sup>ETH Zürich  
 genlm@mit.edu

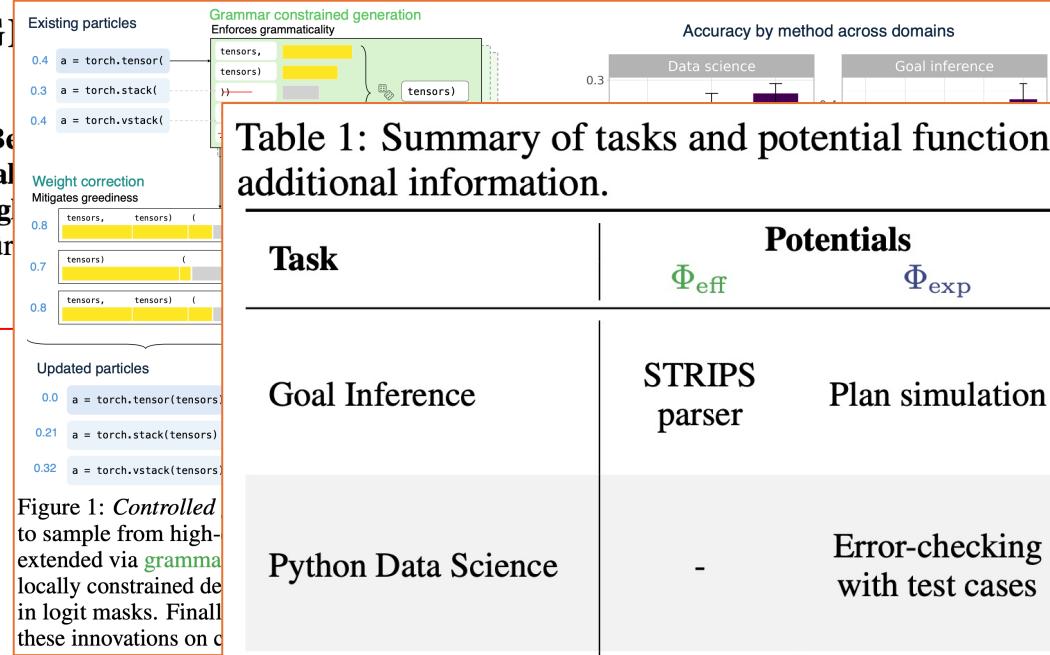


Table 1: Summary of tasks and potential functions. Examples are truncated for brevity. Full prompts include additional information.

| Task                | Potentials          |                                 | Prompt   | Output  |
|---------------------|---------------------|---------------------------------|--|---|
|                     | $\Phi_{\text{eff}}$ | $\Phi_{\text{exp}}$             |  |   |
| Goal Inference      | STRIPS parser       | Plan simulation                 | Write the STRIPS goal condition for the planning problem described below [...]. The STRIPS initial condition is: [...] | (:goal (and (arm-empty) (on-table b1) [...])      |
| Python Data Science | -                   | Error-checking with test cases  | Here is a sample dataframe: [...] I'd like to add inverses of each existing column to the dataframe [...]              | result = df.join(df.apply(lambda x: 1/x) [...])   |
| Text-to-SQL         | SQL parser          | Alias and table-column checking | Here is a database schema: [...] For each stadium, how many concerts are there?  | SELECT T2.name, COUNT(*) FROM concert AS T1 [...] |
| Molecular Synthesis | SMILES parser       | Incremental molecule validation | Given the following list of molecules in SMILES format, write an additional molecule [...]                             | CC1=CC2(OC=N)C(=O) [...]                          |

# SYNTACTIC AND SEMANTIC CONTROL OF LARGE LANGUAGES

João Loula<sup>\*1</sup> Benoît Dauphin<sup>1</sup>  
Tianyu Liu<sup>2</sup> Yalun Wang<sup>2</sup>  
Vikash Mansinghka<sup>1</sup>  
<sup>1</sup>MIT <sup>2</sup>ETH Zürich  
genlm@mit.edu

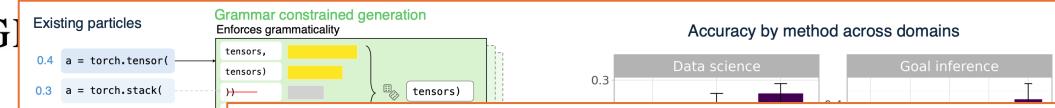
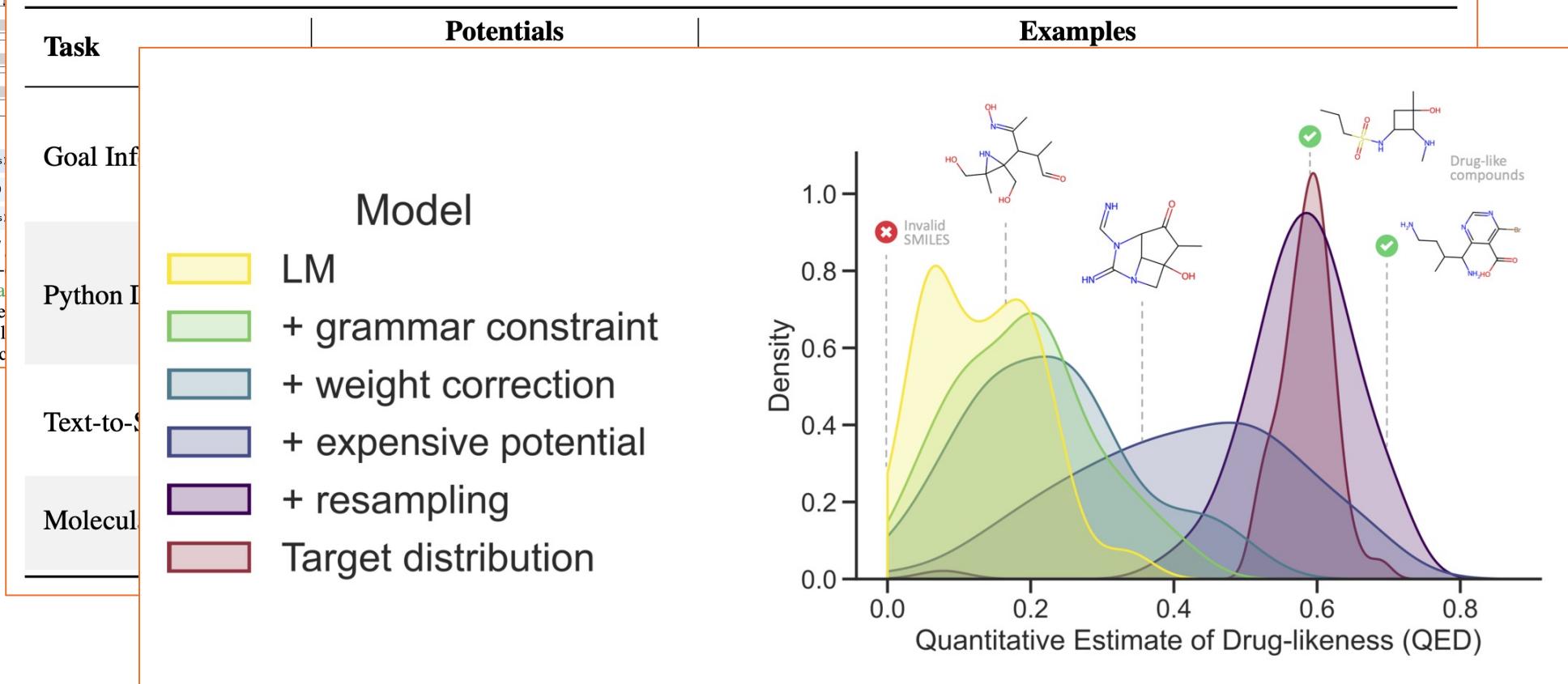


Table 1: Summary of tasks and potential functions. Examples are truncated for brevity. Full prompts include additional information.



# Beyond Constrained Decoding

Programmatic Structured Steering



LMQL

A programming language for large language models.

[Documentation »](#)

[Explore Examples](#) · [Playground IDE](#) · [Report Bug](#)

chat 78 online

pypi package 0.7.3

Watch 21

Fork 210

Star 4.1k

Prompting Is Programming: A Query Language for Large Language Models

1

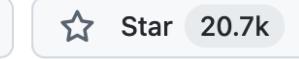
# Prompting Is Programming: A Query Language for Large Language Models

LUCA BEURER-KELLNER, MARC FISCHER, and MARTIN VECHEV, ETH Zurich, Switzerland

```
1 argmax
2   "A list of things not to forget when "
3   "travelling:\n"
4   things = []
5   for i in range(2):
6     "- [THING]\n"
7     things.append(THING)
8   "The most important of these is [ITEM]."
9   from "EleutherAI/gpt-j-6B"
10  where
11    THING in ["passport",
12              "phone",
13              "keys", ...] // a longer list
14    and len(words(THING)) <= 2
```

# {guidance}

Guidance is an efficient programming paradigm for steering language models. With Guidance, you can control how output is structured and get high-quality output for your use case—while reducing latency and cost vs. conventional prompting or fine-tuning. It allows users to constrain generation (e.g. with regex and CFGs) as well as to interleave control (conditionals, loops, tool use) and generation seamlessly.

 Watch 123 ▾  Fork 1.1k ▾  Star 20.7k ▾

```
from guidance.library import one_or_more

@guidance(stateless=True)
def _gen_heading(lm: Model):
    lm += select(
        options=[_gen_text_in_tag("h1"), _gen_text_in_tag("h2"), _gen_text_in_tag("h3")]
    )
    lm += "\n"
    return lm

@guidance(stateless=True)
def _gen_para(lm: Model):
    lm += "<p>"
    lm += one_or_more(
        select(
            options=[
                _gen_text(),
                _gen_text_in_tag("em"),
                _gen_text_in_tag("strong"),
                "<br />",
            ],
        ),
    )
    lm += "</p>\n"
    return lm
```

```
import guidance

from guidance.models import Model

ASCII_OFFSET = ord("a")

@guidance
def zero_shot_multiple_choice(
    language_model: Model,
    question: str,
    choices: list[str],
):
    with user():
        language_model += question + "\n"
        for i, choice in enumerate(choices):
            language_model += f"{chr(i+ASCII_OFFSET)} : {choice}\n"

    with assistant():
        language_model += select(
            [chr(i + ASCII_OFFSET) for i in range(len(choices))], name="string_choice"
        )

    return language_model
```

# Sequential Monte Carlo Steering of Large Language Models using Probabilistic Programs

Alexander K. Lew  
MIT  
alexlew@mit.edu

Tan Zhi-Xuan  
MIT  
xuan@mit.edu

Gabriel Grand  
MIT  
grandg@mit.edu

Vikash K. Mansinghka  
MIT  
vkm@mit.edu

```
# The step method is used to perform a single 'step' of generation.  
# This might be a single token, a single phrase, or any other division.  
# Here, we generate one token at a time.  
async def step(self):  
    # Condition on the next token *not* being a forbidden token.  
    await self.observe(self.context.mask_dist(self.forbidden_tokens), False)  
  
    # Sample the next token from the LLM -- automatically extends `self.context`.  
    token = await self.sample(self.context.next_token())  
  
    # Check for EOS or end of sentence  
    if token.token_id == self.eos_token or str(token) in ['.', '!', '?']:  
        # Finish generation  
        self.finish()
```

# Sequential Monte Carlo Steering of Large Language Models using Probabilistic Programs

Alexander K. Lew  
MIT  
alexlew@mit.edu

Tan Z...  
xuan...@mit.edu

## LLaMPPL



```
# The step method is used to
# This might be a single token
# Here, we generate one token
async def step(self):
    # Condition on the next
    await self.observe(self.context.mask_dist(self.forbidden_tokens), False)

    # Sample the next token from the LLM -- automatically extends `self.context`.
    token = await self.sample(self.context.next_token())

    # Check for EOS or end of sentence
    if token.token_id == self.eos_token or str(token) in ['.', '!', '?']:
        # Finish generation
        self.finish()
```

LLaMPPL is a research prototype for language model probabilistic programming: specifying language generation tasks by writing probabilistic programs that combine calls to LLMs, symbolic program logic, and probabilistic conditioning. To solve these tasks, LLaMPPL uses a specialized sequential Monte Carlo inference algorithm. This technique, SMC steering, is described in [our recent workshop abstract](#).

# Summary

- We are building an arsenal of tools to generate good **symbolic** code from **neural** language models
  - Natural (Magical) strategies
    - Prompting
    - Prompt tuning
  - Mechanical strategies
    - Controlled decoding
  - Structural strategies
    - Iterative refinement
    - Programmatic steering
  - Agentic strategies
    - External tool use

# We have not discussed...

- Syntactic constraints
  - Context free grammar
  - Regular expressions (Regex)
  - Finite state machine; Automaton; how to compile an automaton from CFG
  - Algorithms for lexing, parsing, compile a parsing, how to write AST
- Steering libraries
  - LMQL
  - Guidance
  - Langchain
  - DSPY

# Week 4

- Assignment 1
  - <https://github.com/machine-programming/assignment-1>
  - Accepting late submissions
- Assignment 2
  - <https://github.com/machine-programming/assignment-2>
  - Due in two weeks; sending out another set of API keys; autograders
- Oral presentation starting from week 7
  - Sign-up sheet going out this week
- Feedback Questionnaire
  - Sending out this week
  - +0.5% of your overall grade