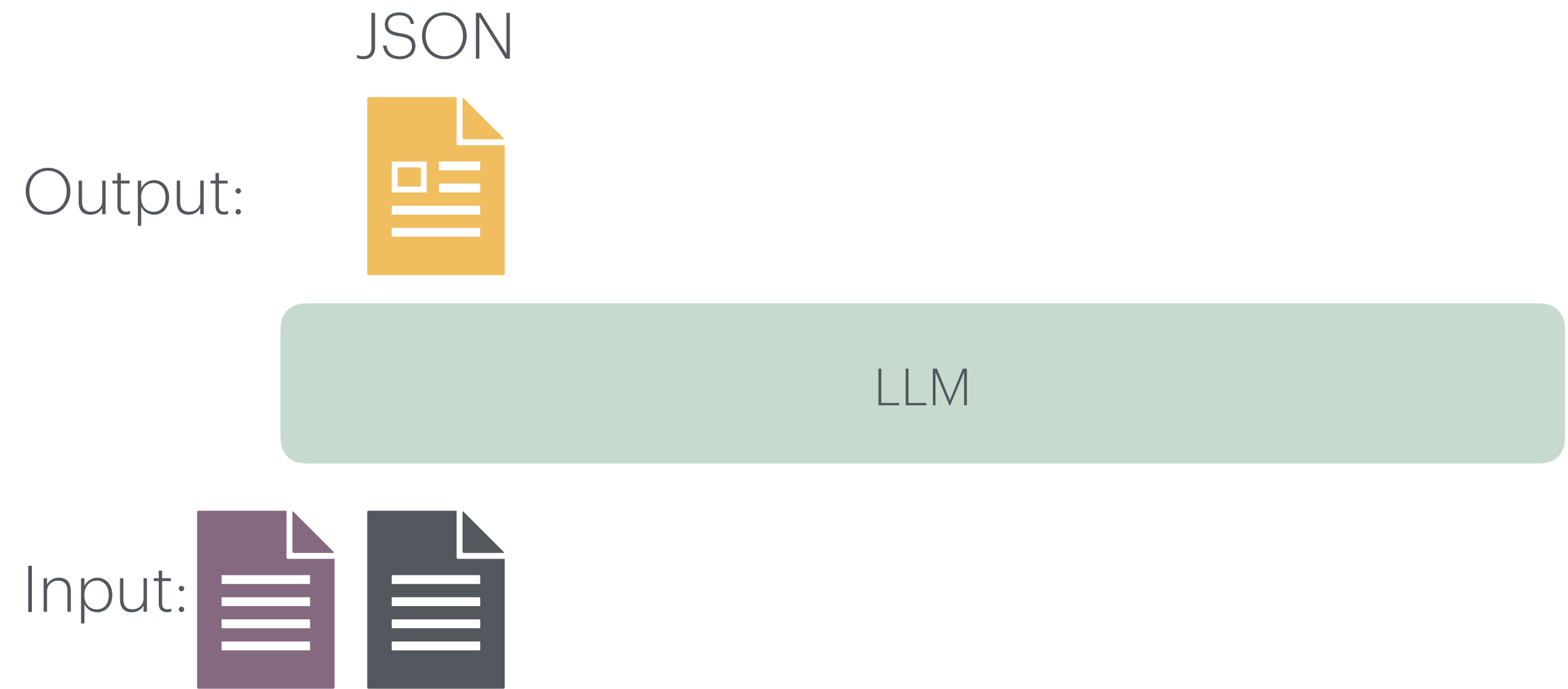


Constrained Decoding

Structured output

- What if we only want to **parse** output of LLM?
 - Option 1.1
 - Write a robust parser (in python)
 - Let LLM know that you failed to parse
 - Hope for the best
 - **Option 1.2: Constrain output**



Constrained Decoding

- Change Sampling of LLM
 - $x_{i+1} \sim P(x_1 \dots x_i)$
 - Only sample valid next tokens x_{i+1}
- How to define valid?



Constrained Decoding

- Change Sampling of LLM
 - $x_{i+1} \sim P(x_1 \dots x_i)$
 - Only sample valid next tokens x_{i+1}
- Use a Context Free Grammar
 - Regex++

```
root ::= object
value ::= object | array | string | number | ("true" | "false"
| "null") ws

object ::=
  "{" ws (
    string ":" ws value
    ("," ws string ":" ws value)*
  )? "}" ws

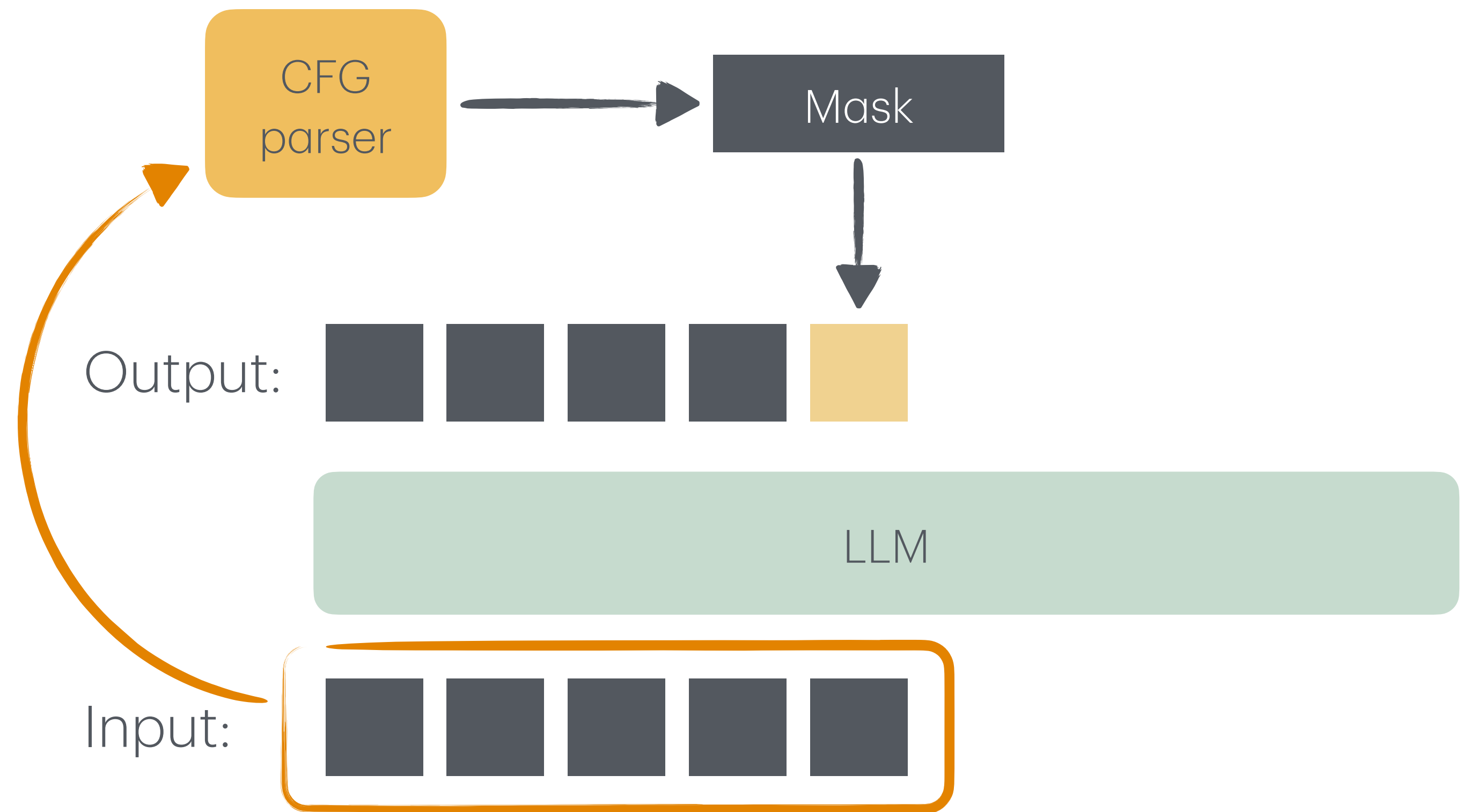
array ::=
  "[" ws (
    value
    ("," ws value)*
  )? "]" ws

string ::=
  "\"" (
    [^"\\x7F\x00-\x1F] |
    "\\\" (["\bfnrt] | "u" [0-9a-fA-F]{4}) # escapes
  )* "\"" ws

number ::= ("-"? ([0-9] | [1-9] [0-9]{0,15})) ( "." [0-9]+)?
([eE] [-+]? [0-9] [1-9]{0,15})? ws

# Optional space: by convention, applied in this grammar after
literal chars when allowed
ws ::= | " " | "\n" [ \t]{0,20}
```

Constrained Decoding



- CFG parser
- Computes valid completions $v \in V_{i+1}$

- Only sample valid next tokens $P \rightarrow \hat{P}$

$$\hat{P}(v) = 0 \quad \forall v \notin V_{i+1}$$

$$x_{i+1} \sim \hat{P}(x_1 \dots x_i)$$

Algorithm 1 Constrained Decoding

Input: Checker C , LLM f , Tokenized Prompt x

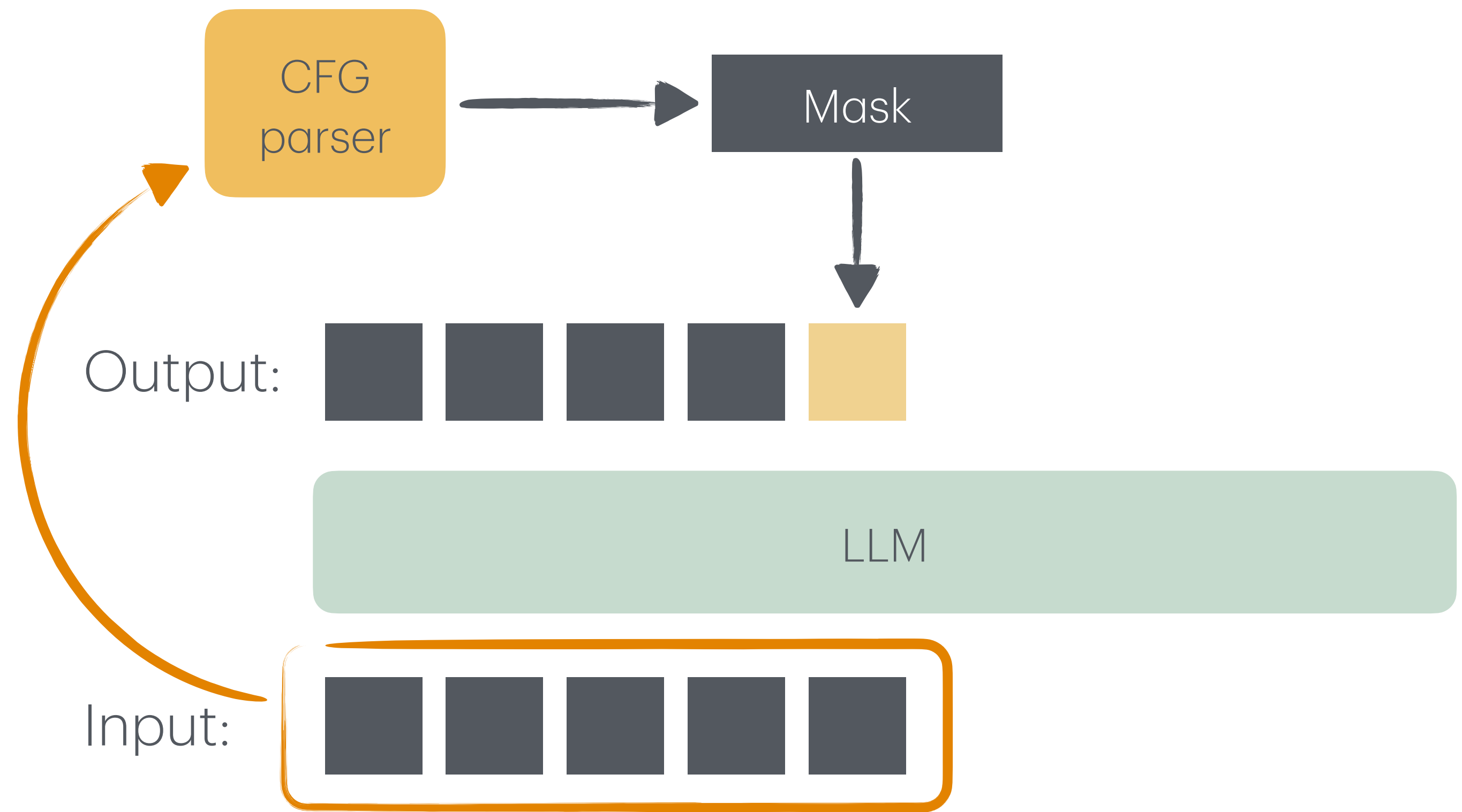
Output: Completion o adhering to C

```
1:  $o \leftarrow []$ 
2:  $C.init()$ 
3: loop
4:    $C.update(o)$  // advance state of  $C$ 
5:    $m \leftarrow C.mask()$  // compute mask
6:    $v \leftarrow f(x + o)$  // compute logits
7:    $v' \leftarrow m \odot v$ 
8:    $t \leftarrow decode(v')$  // e.g., argmax or sample
9:   if  $t = EOS$  then break
10:   $o.append(t)$ 
11: end loop
12: return  $o$  // optionally detokenize
```

Constrained Decoding

Issue

- Constraints: On text / output
- Decoding: on tokens
 - Many tokens streams = one text
 - Model only likes one



Constrained Decoding

Solution

- Lift CFG to token level

Algorithm 2 Construct Terminal Tree

Input: CFG G , Alphabet Σ , Vocabulary \mathcal{V}

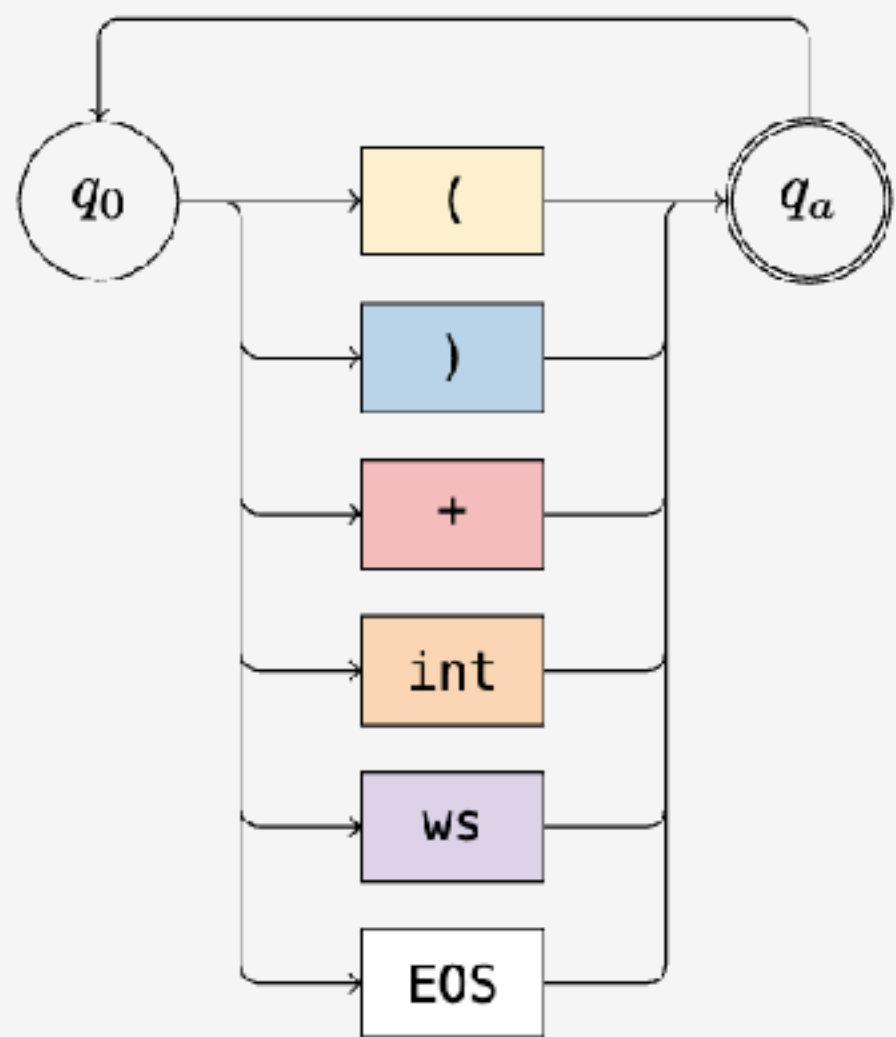
Output: Scanner S

- 1: $\mathcal{T} = \{\}$
- 2: **for all** $q \in S.states()$ **do**
- 3: $\alpha \leftarrow q.subterminal()$ // get current (sub)terminal
- 4: **for all** $l \in \mathcal{V}$ **do**
- 5: $\{\alpha_1^j, \dots, \alpha_{m_j}^j\}_j \leftarrow q.traverse(l)$
- 6: $\mathcal{T} \leftarrow \mathcal{T} \cup \{(\alpha_1^j, \dots, \alpha_{m_j}^j), l\}_j$
- 7: **end for**
- 8: $T_q \leftarrow PrefixTree(\mathcal{T})$
- 9: **end for**

(a) Example Grammar

$E = E \mid E + E \mid (E) \mid int$
 $int = ([1-9][0-9]^*) \mid (0+)$

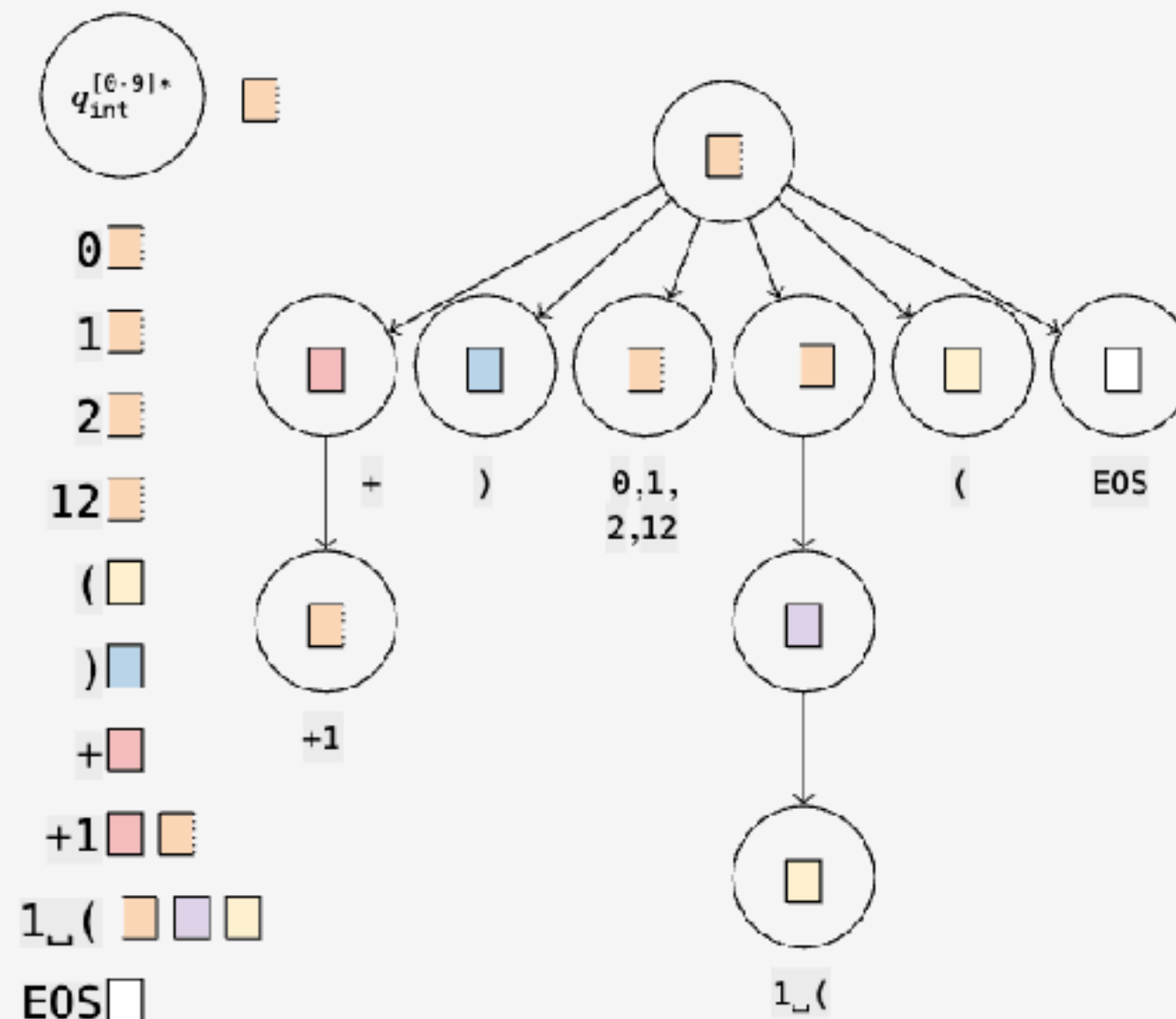
(b) Character Scanner



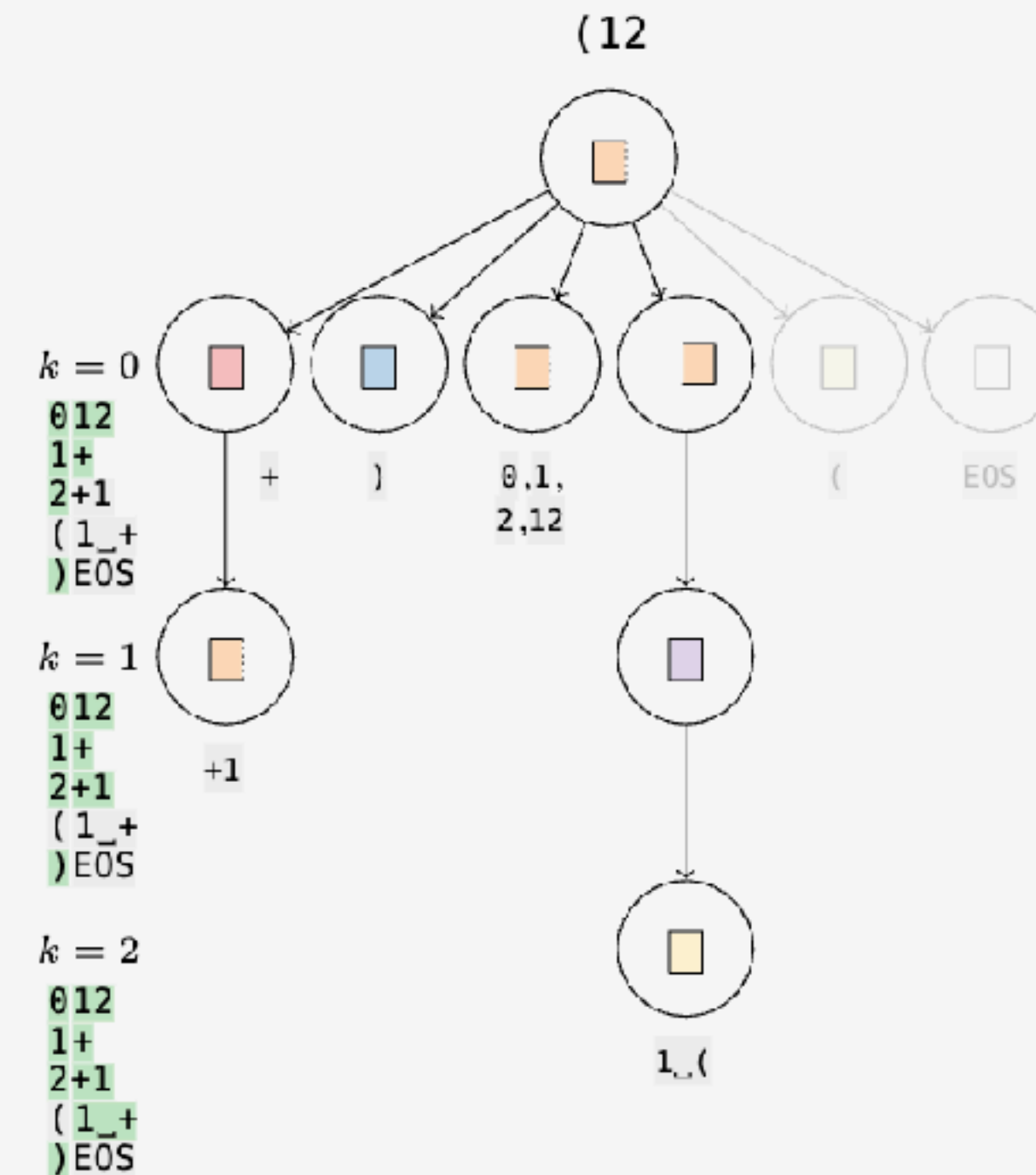
(c) Vocabulary

$\mathcal{V} = \{0, 1, 2, 12,), (, +, +1, 1_-, EOS\}$

(d) Vocabulary-aligned Subterminal Tree (offline, per node)



(e) Parser (online)



Constrained Decoding

Solution

- Very complex to implement

Algorithm 2 Construct Terminal Tree

Input: CFG G , Alphabet Σ , Vocabulary \mathcal{V}

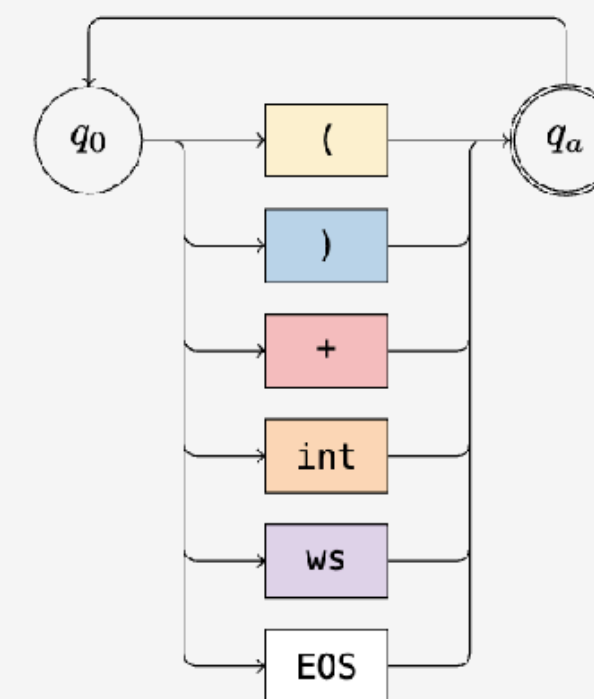
Output: Scanner S

- 1: $\mathcal{T} = \{\}$
- 2: **for all** $q \in S.\text{states}()$ **do**
- 3: $\alpha \leftarrow q.\text{subterminal}()$ // get current (sub)terminal
- 4: **for all** $l \in \mathcal{V}$ **do**
- 5: $\{\alpha_1^j, \dots, \alpha_{m_j}^j\}_j \leftarrow q.\text{traverse}(l)$
- 6: $\mathcal{T} \leftarrow \mathcal{T} \cup \{(\alpha_1^j, \dots, \alpha_{m_j}^j), l\}_j$
- 7: **end for**
- 8: $T_q \leftarrow \text{PrefixTree}(\mathcal{T})$
- 9: **end for**

(a) Example Grammar

$E = E \mid E + E \mid (E) \mid \text{int}$
 $\text{int} = ([1-9][0-9]^*) \mid (0+)$

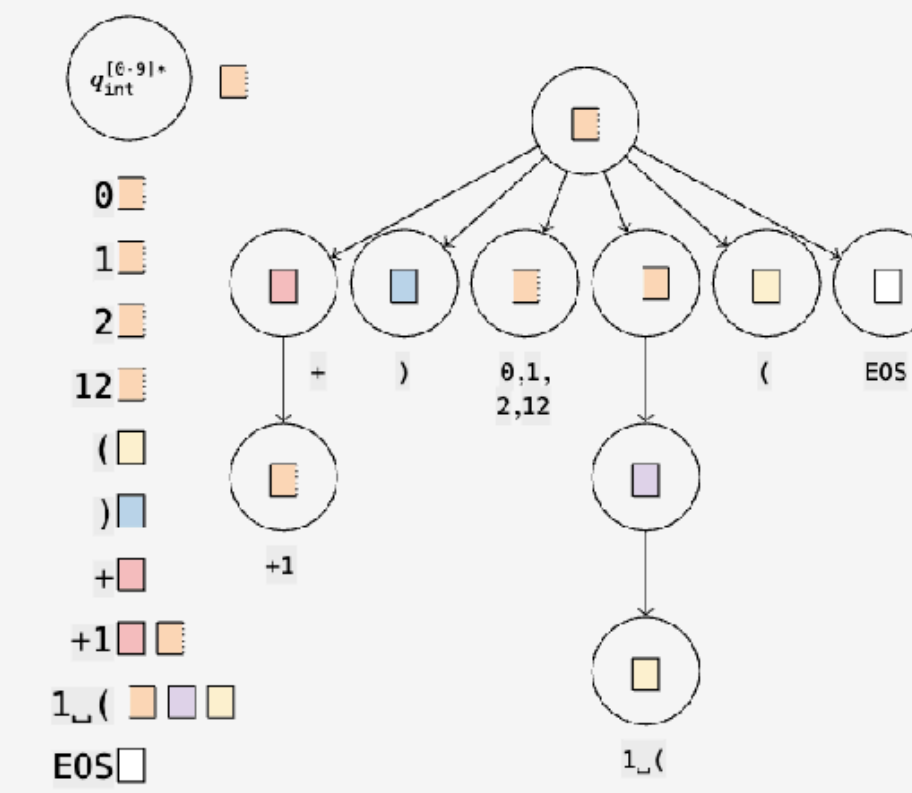
(b) Character Scanner



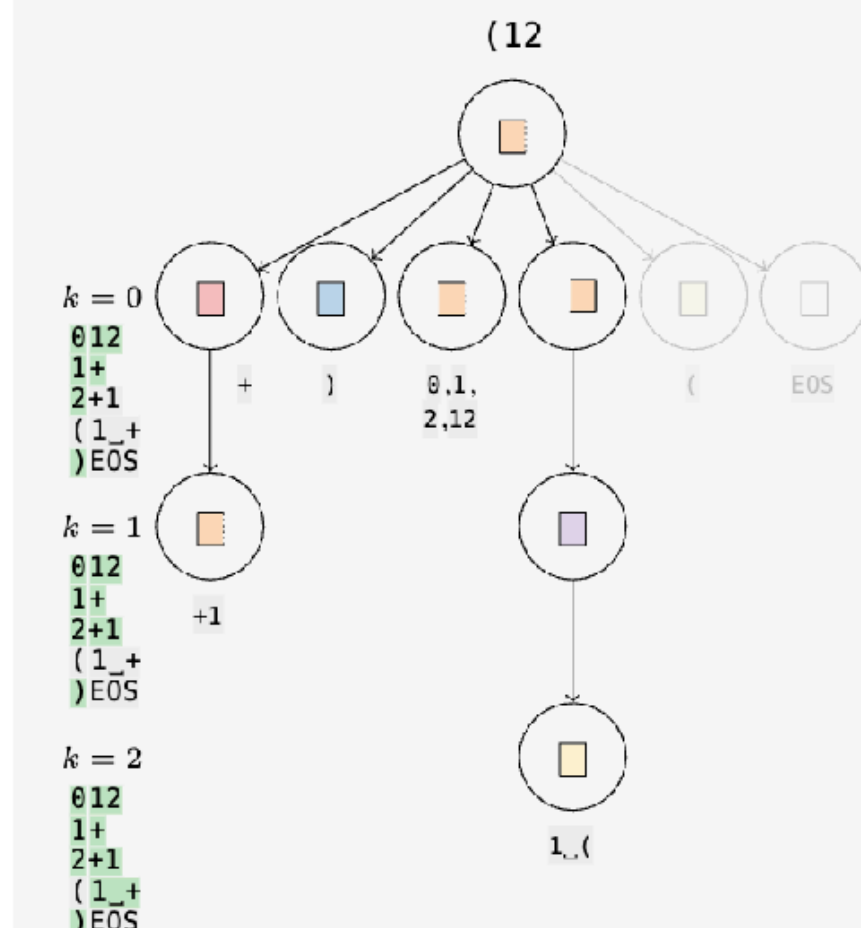
(c) Vocabulary

$\mathcal{V} = \{0, 1, 2, 12,), (, +, +1, 1_-, \text{EOS}\}$

(d) Vocabulary-aligned Subterminal Tree (offline, per node)



(e) Parser (online)



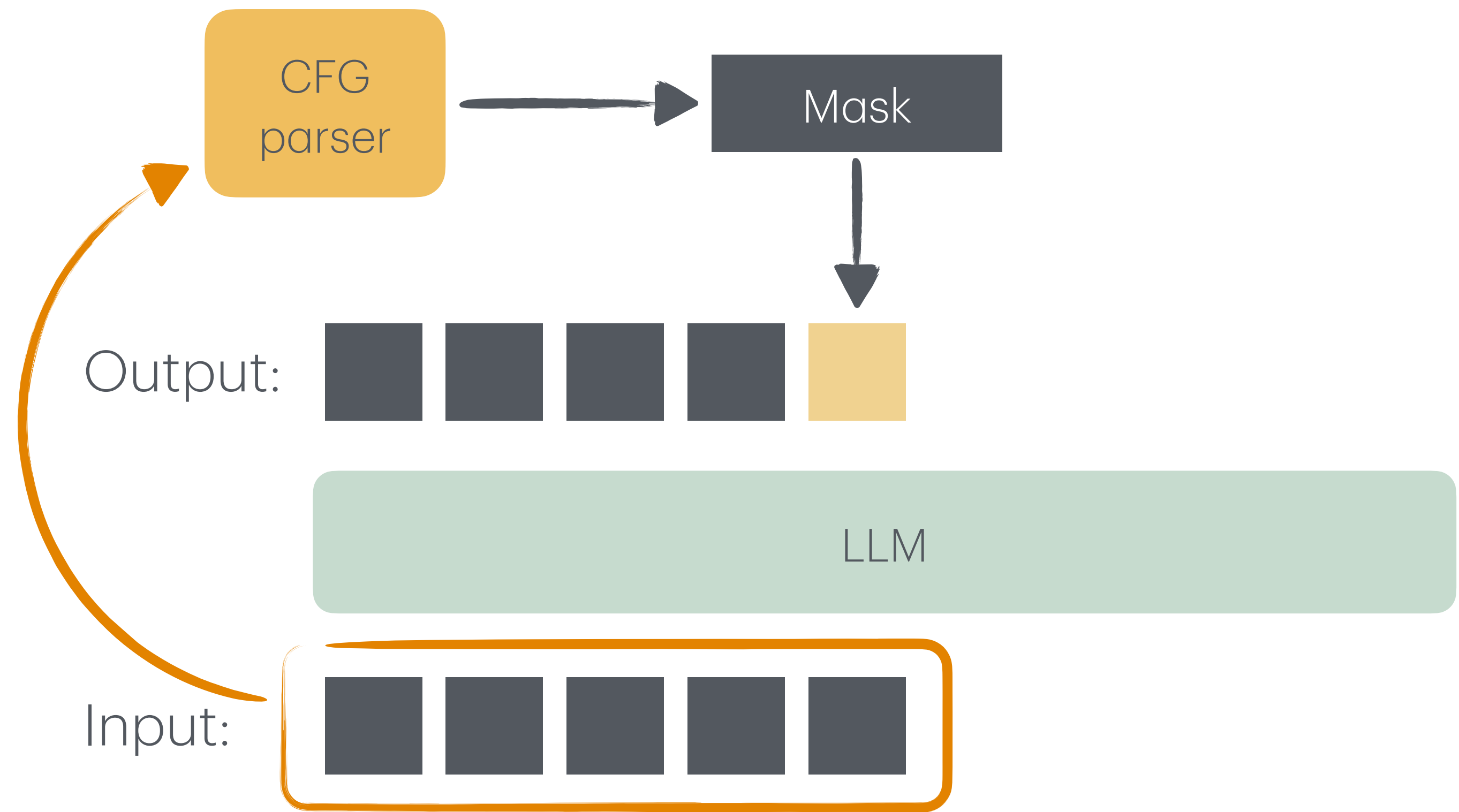
Constrained Decoding

Issue

- Outputs are Biased
- Example CFG
 - $S \rightarrow aSb \mid Sc \mid \epsilon$

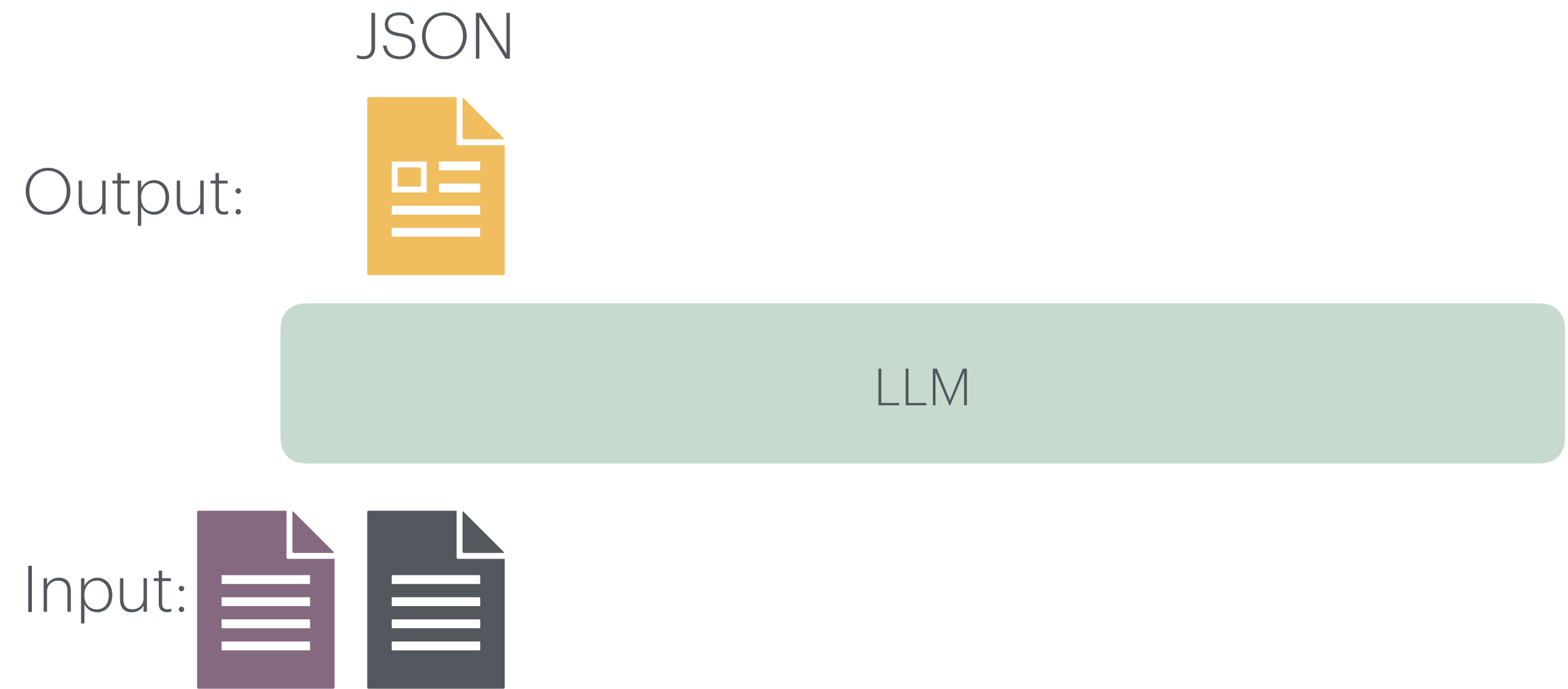
- $$P(x_{i+1} | x_1 \dots x_i) = \begin{cases} 0.9 & \text{if } x_{i+1} = a \\ 0 & \text{if } x_{i+1} = b \\ 0.1 & \text{if } x_{i+1} = c \end{cases}$$

- Sampling will never stop: aaaaaaaaaa



Structured output

- What if we only want to **parse** output of LLM?
 - Option 1.1
 - Write a robust parser (in python)
 - Let LLM know that you failed to parse
 - Hope for the best
 - Option 1.2: Constrain output



References

- [1] Synchromesh: Reliable code generation from pre-trained language models, Poesia etal 2022
- [2] Guiding LLMs The Right Way: Fast, Non-Invasive Constrained Generation, Beurer-Kellner etal 2024