PSC: EFFICIENT GRAMMAR-CONSTRAINED DECODING VIA PARSER STACK CLASSIFICATION

Anonymous authors

Paper under double-blind review

ABSTRACT

LLMs are widely used to generate structured output like source code or JSON. Grammar-constrained decoding (GCD) can guarantee the syntactic validity of the generated output, by masking out tokens that violate rules specified by a context-free grammar. However, the online computational overhead of existing GCD methods, with latency typically scaling linearly with vocabulary size, limits the throughput of LLMs, especially for models with large vocabularies. To address this issue, we propose PSC, a novel grammar-constrained decoding method. By combining acceptance conditions of all vocabulary tokens into a single classifier of the parser stack during preprocessing, PSC can compute the complete vocabulary mask by checking the parser stack exactly once per decoding step, with time complexity independent of the vocabulary size. Experiments show that PSC computes masks up to 770× faster than baselines on complex programming language grammars, and up to 30× faster for schema-conformant JSON; end-to-end LLM throughput with PSC approaches that of unconstrained decoding.

1 Introduction

In recent years, the ability for Large Language Models (LLMs) to generate structured output has been widely recognized and utilized (Qwen et al.; Grattafiori et al.; Gemma Team et al.). Source code can be viewed as structured output that adheres to the syntax of programming languages, and LLM-based coding assistants, such as GitHub Copilot (GitHub) and Cursor (Anysphere Inc.), have been widely adopted by developers to assist in writing code to improve their productivity. When LLMs are used as a tool, users often expect the generated output to conform to a specific format, such as Markdown or JSON with custom schemas (Liu et al., 2024; vLLM Team; OpenAI). All of these applications rely on the ability of LLMs to generate output that adheres to a specific syntax.

However, generating in a structured format is complex, as it requires not only understanding the semantics of given input but also adhering to the specific grammars of target formats. Since language models are essentially probabilistic models, there is no guarantee that the generated output will always conform to the required grammar.

To address this issue, grammar-constrained decoding (GCD) (Geng et al., b; Scholak et al.; Poesia et al.; Ugare et al.) is proposed to ensure that the generated output always conforms to the specified context-free grammar. A GCD method works by incorporating a grammar checker into the decoding process, as shown in Figure 1a. At each decoding step, the checker determines which tokens in the vocabulary can be appended to the current prefix while not violating the grammar. The logits generated by the language model are then masked to only allow the valid tokens, and the next token is generated by sampling from the masked logits.

The overhead of GCD is determined by the newly introduced step of validity calculation. A naive implementation, as shown in Figure 1b would require calling the parser for *every* token in the vocabulary to check its validity, resulting in a time complexity of $\mathcal{O}(|\mathcal{V}|)$ per decoding step, where $|\mathcal{V}|$ is the vocabulary size. This can add significant overhead, especially for large vocabularies in modern language models, e.g. 128k tokens in Llama-3 (Grattafiori et al.), 151k tokens in Qwen series (Bai et al.), or 262k tokens in Gemma 3 (Gemma Team et al.). The overhead is particularly pronounced for smaller models, where the time taken by model inference is relatively small.

056

058

060 061

062

063

064

065

066

067

068 069

070

075

076

077

079

081

083

084

085

087

090

091

092

094

095 096

097

098

099

102

103

105

106

107

Figure 1: An illustration of grammar-constrained decoding, showing (a) the overall working process, (b) the naive implementation that directly simulates the PDA, and (c) our method PSC that precomputes the DFA and the valid token masks.

To speed up GCD, various techniques have been proposed in the literature. We summarize these techniques in Section 2. However, none of these techniques can fundamentally change the time complexity, which is still $\mathcal{O}(|\mathcal{V}|)$ in the worst case.

We propose a novel GCD method PSC that replaces the repetitive runtime parsing over the whole vocabulary with a one-time classification of the current parser stack, as shown in Figure 1c. The checking process of the parser can be seen as a function of both the token and the state of the parser, which is usually a stack. For each token, our method PSC constructs a finite-state automaton (FSA) that represents the exact requirements on the parser stack to accept that token, i.e., the FSA accepts a parser stack if and only if that token is accepted by a parser with that stack. All these FSAs can then be combined into a single FSA that classifies the parser stack into a finite number of classes, each corresponding to a different vocabulary mask. During decoding, we only need to check the parser stack exactly once per decoding step to get the vocabulary mask, which is ready to be applied to the logits. This eliminates the need to call the parser for each token in the vocabulary, resulting in a significant speedup.

We conduct extensive experiments on grammar-constrained decoding in Java, Go, SQL, and schema-conformant JSON to evaluate the efficiency of our method. Compared to the current state-of-the-art method LLGuidance, our method achieves up to 770 times speedup in mask computation on complex programming language grammars, and up to 30 times speedup for schema-conformant JSON generation. In the end-to-end decoding throughput experiments, the throughput of PSC approaches that of unconstrained decoding, and is significantly higher than LLGuidance, especially on smaller models and larger batch sizes.

In summary, our contributions are as follows:

- We propose a novel GCD method PSC that leverages finite-state automata to classify parser states to use the precomputed vocabulary mask, significantly reducing the time overhead of grammar-constrained decoding.
- We provide a theoretical analysis of PSC. We prove that the set of the parser stacks that can accept a given token can be formally described as a regular language. This justifies the correctness of our method and provides a theoretical foundation for future research on grammar-constrained decoding.
- We demonstrate the efficiency of PSC through extensive experiments on grammarconstrained decoding in Java, Go, Python, and schema-conformant JSON, achieving significant speedup in mask computation compared to existing techniques; end-to-end decoding throughput with PSC approaches that of unconstrained decoding.

2 BACKGROUND AND RELATED WORK

We introduce the task of grammar-constrained decoding in this section, and then review the related work. For the symbols, notations, and definitions used in this paper, please refer to Appendix A.2.

2.1 The Task: grammar-constrained decoding

Let Σ be the character set used by the language model, e.g. the Unicode. Given a prefix of tokens, the task of a language model is to generate the next-token distribution over the vocabulary $\mathcal{V} \subset \Sigma^+$. For a language $L \subseteq \Sigma^*$, the task of *constrained decoding* aims to generate a sample in L from the language model. In each step, given prefix $x \in \Sigma^*$, it calculates the set of valid tokens in \mathcal{V} , i.e. tokens that, when concatenated after x, become a prefix of some strings in L.

$$c(x \in \Sigma^*) := \{ v \in \mathcal{V} | \exists y \in \Sigma^*, xvy \in L \}. \tag{1}$$

When the language L is defined by a context-free grammar, the task is called grammar-constrained decoding (Ugare et al.; Koo et al.; Park et al.; Moskal et al.). Determining whether a string is syntactically valid usually involves two phases: lexical analysis and syntax analysis (Aho & Ullman). In lexical analysis, the lexer \mathcal{T} , usually modeled as a deterministic finite-state transducer (FST) (Aho & Ullman; Koo et al.; Park et al.), transduces the text $w \in \Sigma^*$ into a terminal sequence $\mathcal{T}(w) \in \Gamma^*$. In syntax analysis, the parser \mathcal{P} , usually modeled as a terminating deterministic push-down automaton (PDA) (Aho & Ullman), determines whether a terminal sequence $x \in \Gamma^*$ is valid, written as $x \in \mathcal{P}$. So we have

$$w \in L \iff \mathcal{T}(w) \in \mathcal{P}.$$
 (2)

2.2 RELATED WORK

There are several existing techniques to speed up grammar-constrained decoding, which can be categorized into three groups: vocabulary preprocessing, lexer preprocessing, and parser preprocessing.

Vocabulary preprocessing (Poesia et al.; Beurer-Kellner et al.; Moskal et al.) exploits the fact that the vocabulary is built by BPE (Gage; Sennrich et al.), and for each token, its prefix is also a token in the vocabulary. If the prefix token is rejected, then the longer token must also be rejected. So we can traverse the vocabulary in a tree-like manner, and only check the tokens whose prefixes are not rejected.

Lexer preprocessing (Beurer-Kellner et al.; Park et al.; Moskal et al.) maps each token to a terminal sequence during preprocessing, and then the parser is only called on the terminal sequences. This reduces the number of parser calls, as different tokens may share the same terminal sequence.

Parser preprocessing (Dong et al.; Park et al.) classifies the vocabulary into three sets for each parser state: context-independent accepted, context-independent rejected, and context-dependent. This allows us to reduce the number of parser calls by only checking the context-dependent tokens.

These techniques can be combined to achieve better speedup (Beurer-Kellner et al.; Park et al.; Moskal et al.). However, as mentioned in Section 1, these techniques are limited in their speedup. There is no theoretical guarantee on how many parser calls will be made per decoding step, which can be linear to the vocabulary size in the worst case.

3 PSC: PARSER STACK CLASSIFICATION

The preprocessing of the lexer \mathcal{T} is described in Section 3.1, and the other parts of this section are dedicated to the preprocessing of the parser \mathcal{P} . All the proofs are deferred to Appendix A.3. For the symbols, notations, and definitions used in this paper, please refer to Appendix A.2.

¹To ease the presentation, the step of lexical analysis is omitted in previous sections. The term "parser" in previous sections should be realized as the combination of the lexer and the parser defined here, and the term "parser stack" should be realized as the concatenation of the lexer state and the parser state, which is a simple state without internal structure, and a stack, respectively.

3.1 Lexical preprocessing

Lexical preprocessing is not our focus in this paper, so we reuse the lexical preprocessing in Great-Gramma (Park et al.), and conclude it here as a prelude to PSC.

For token $v \in \mathcal{V}$, lexer state $q \in Q$, if lexing the token v from state q using the lexer \mathcal{T} generates the terminal sequence $x \in \Gamma^*$, and the lexer transits to state p, i.e. $q \xrightarrow{v:x} p$, we define the realizable terminal sequences $R_q(v)$ as $\{x\}T_p$, representing all possible terminal prefixes that can be generated from a string starting with v, where T_p is the finite set of all possible terminal prefixes from state p: $\mathcal{T}_p(\Sigma^*) = T_p\mathcal{T}(\Sigma^*)$. The set $R_q(v)$ can be precomputed for every $q \in Q, v \in \mathcal{V}$ during preprocessing.

3.2 Overview of syntactic preprocessing

Given a valid prefix $x \in \Sigma^*$, we run the lexer \mathcal{T} from the initial state q_0 to produce a terminal sequence $z \in \Gamma^*$ and a new lexer state $q: q_0 \xrightarrow{x:z} q$. We then run the parser \mathcal{P} from the initial stack

 γ_0 to receive the terminal sequence z, and generates a new stack α : $\gamma_0 \stackrel{z}{\Longrightarrow} \alpha$.

We can now introduce the simplification of the condition in the GCD definition in Equations 1 and 2 from previous work (Park et al.). For any token $v \in \mathcal{V}$, to determine whether v is valid, we can rewrite the condition in terms of realizable terminal sequences,

$$\exists y \in \Sigma^*, \mathcal{T}(xvy) \in \mathcal{P} \iff \exists w \in R_q(v), \exists \beta \in \Pi^+, \alpha \underset{\mathcal{P}}{\overset{w}{\Longrightarrow}} {}^*\beta. \tag{3}$$

The simplification is based on the common assumption that, if the parser reads in a certain terminal sequence and enters a stable stack, then we do not need to worry about the rest of the input, and there always exists a terminal sequence produced by the lexer that can ensure the whole text is accepted.

For $w \in \Gamma^*$, we define $P_w(\alpha)$ for the calculation in the last step of Equation 3,

$$P_{w \in \Gamma^*}(\alpha \in \Pi^+) := \left\{ \beta \in \Pi^+ \middle| \alpha \xrightarrow{\frac{w}{\mathcal{P}}}^* \beta \right\}. \tag{4}$$

How to efficiently calculate $P_w(\alpha)$ is the key difference between PSC and previous methods. In existing work, the calculation of P_w is almost always dynamic: one has to calculate $P_w(\alpha)$ for the current α and every possible $w \in R_q(\mathcal{V})$. While some methods employ precomputation to optimize certain cases, they still fundamentally require $\mathcal{O}(|R_q(\mathcal{V})|t_{\mathcal{P}})$ time for dynamic parsing, where $\mathcal{O}(t_{\mathcal{P}})$ is the time required to run the parser \mathcal{P} .

On the other hand, it can be noticed that, whether $P_w(\alpha) \neq \emptyset$ is determined by the top symbols of α . It is to say, if $P_w(\alpha) \neq \emptyset$, then $\forall \gamma \in \Pi^*$, $P_w(\alpha\gamma) \neq \emptyset$. Utilizing this property, PSC proposes a totally different approach, modeling P_w as a deterministic FST, which reads the stack sequence $\alpha \in \Pi^+$, and outputs the sequence $\beta \in \Pi^+$ if there is one.

This gives us several benefits. Because each P_w reads and outputs a sequence of stack symbols, they can be composed to create larger FSTs: $P_t \circ P_s = P_{st}$. Because the realizable terminal sequences are known during precomputation, the exact validity condition of each vocabulary is therefore known, their combinations can be precomputed, and we only need to go through the current stack once during runtime. All possible masks can also be precomputed, eliminating the mask generation overhead during decoding.

3.3 FST of ε transitions

In this section, we focus on how to construct P_{ε} . The input of P_{ε} should be a stack, and **the output** is the stabilized version of it by executing all needed ε transitions on the stack. The states of P_{ε} is a special state FINAL plus a subset of Π^* which represent the current stack top, both the input alphabet and output alphabet are Π , and the start state is ε , and the final states are $\{\text{FINAL}\}$.

There are three types of transitions in P_{ε} . If the information of the current stack top is not enough to determine whether the stack is stable, it transits to a new state by reading the next stack symbol. If the stack top is stable, it transits to the FINAL state to output the final stable stack. Otherwise,

217

218

219

220 221

222

223

224 225

226 227

228

229

230

231

232

233

235

236

237

238

239

240

241

242

243

244

245

246

247

249 250 251

252 253

254

255

256

257

258

259 260

261

262

263

264

265

266

267 268

269

it simulates the transition of \mathcal{P} on the current stack top, and transits to a new state representing the new stack top after executing the ε transition. The construction algorithm is simply searching all needed stacks from ε , as shown in Algorithm 1.

Theorem 1. Algorithm 1 constructs a finite-state transducer P_{ε} as defined in Definition 4.

 P_{ε} is an important building block in the construction of other $P_w, w \in \Gamma^+$. For any stack $\alpha, P_{\varepsilon}(\alpha)$ gives the stabilized version of α , so the FST composed after P_{ε} does not need to handle ε transitions, and we can compose P_{ε} after other FSTs to meet the stability requirement in the definition of P_w .

```
Algorithm 1 Construct FST P_{\varepsilon}
                                                                                                                           Algorithm 2 Construct FST \tilde{P}_a for terminal a \in \Gamma
  1: function EpsilonFST(\mathcal{P})
                  for all X \in \prod \mathbf{do}
 | FINAL \xrightarrow{X:X} FINAL
  2:
                                                                                                                             2:
  3:
                                                                                                                             3:
  4:
                    while let \alpha \in \Pi^*, Q \leftarrow \text{Pop}(Q) do
  5:
                             \begin{array}{c} \text{if } \alpha \in F_{\mathcal{P}}\Pi^* \text{ then} \\ \alpha \xrightarrow[P_{\varepsilon}]{\varepsilon:\alpha} Final \end{array}
  6:
                                                                                                                             5:
  7:
                              else if |\alpha| < 2 then
                                                                                                                             6:
  8:
                                       for all X \in \Pi do
  9:
10:
11:
                              else
12:
                                       let \alpha_0 \gamma = \alpha,
13:
                                       where \alpha_0 \in \Pi^2, \gamma \in \Pi^*
                                                                                                                             3:
                                      if \alpha_0 \xrightarrow[\mathcal{P}]{\varepsilon} \beta then
14:
                                                                                                                              4:
                                                                                                                              5:
                                      \begin{array}{c|c} & \rho & \varepsilon : \varepsilon \\ & \alpha_0 \gamma & \frac{\varepsilon : \varepsilon}{P_\varepsilon} \beta \gamma \\ & Q \leftarrow Q \cup \{\beta \gamma\} \\ \text{else if } \alpha_0 & \frac{a}{\mathcal{P}} \beta \text{ then} \\ & \alpha & \frac{\varepsilon : \alpha}{P_\varepsilon} \text{FINAL} \end{array}
15:
                                                                                                                             6:
                                                                                                                             7:
16:
                                                                                                                             8:
17:
                                                                                                                             9:
18:
                                                                                                                           10:
                                                                                                                                                return A
                    return \overline{P}_{\varepsilon}
19:
```

```
1: function TERMINALFST(\mathcal{P}, a)
                           for all X \in \Pi do
                         for all X \in \Pi do
\begin{bmatrix}
FINAL & \frac{X:X}{\tilde{P}_a} & FINAL \\
\text{for all } XY & \frac{a}{\mathcal{P}} & \beta \text{ do}
\end{bmatrix}
\begin{bmatrix}
\varepsilon & \frac{X:\varepsilon}{\tilde{P}_a} & X & \frac{Y:\varepsilon}{\tilde{P}_a} & XY & \frac{\varepsilon:\beta}{\tilde{P}_a} & FINAL
\end{bmatrix}
```

```
Algorithm 3 Offline constructon in PSC
  1: function OffLineConstruction(\mathcal{T}, \mathcal{P}, \mathcal{V})
                  P_{\varepsilon} \leftarrow \text{EpsilonFST}(\mathcal{P})
                 for all a \in \Gamma do
                          P_a \leftarrow \text{TERMINALFST}(\mathcal{P}, a)
                  for all w = w_1 \dots w_n \in R(\mathcal{V}) do
                          \begin{array}{l} P_w \leftarrow P_\varepsilon \circ \tilde{P}_{w_n} \circ P_\varepsilon \circ \cdots \circ P_\varepsilon \circ \tilde{P}_{w_1} \circ P_\varepsilon \\ A_w \leftarrow \mathsf{REMOVEOUTPUT}(P_w) \end{array}
                  \mathcal{A} \leftarrow \bigcup_{v \in \mathcal{V}} \bigcup_{q \in Q} \bigcup_{w \in R_q(v)} B_q A_w B_v
                  \mathcal{A} \leftarrow \text{MINIMIZE}(\mathcal{A})
```

3.4 FST FOR ANY TERMINAL SEQUENCE

After constructing P_{ε} , the construction of P_w for any terminal sequence $w \in \Gamma^+$ is fairly simple.

We first construct an FST P_a for every $a \in \Gamma$ that simulates a single transition labeled a. Note that in P_a , we only force the stack to execute one transition labeled a, but the output stack is not required to be stable. Its input is a stack α , and the output is the immediate result of executing a-labelled transition $\alpha \xrightarrow[\mathcal{P}]{a} \beta$. The states are $\Pi^0 \cup \Pi^1 \cup \Pi^2 \cup \{\text{FINAL}\}$, the start state is ε , the final state is FINAL, and the transitions are defined in Algorithm 2.

For any terminal sequence $w=w_1\dots w_n\in\Gamma^+$, the FST P_w can be constructed using P_ε and P_{w_1},\ldots,P_{w_n} . For any stack α , we first pass it to P_{ε} to get a stable stack, and then pass the result into P_{w_1} to get the stack after reading w_1 , and then the stack is passed to P_{ε} to get the stabilized version. The result is then passed into P_{w_2} to make the stack read w_2 , and then passed to P_{ε} to stabilize. We repeat this process until we have read w_n using P_{w_n} and stabilize the stack with P_{ε} . Formally, for $w = w_1 \dots w_n \in \Gamma^+$, we have $P_w = P_\varepsilon \circ \tilde{P}_{w_1} \circ P_\varepsilon \circ \dots \circ P_\varepsilon \circ \tilde{P}_{w_n} \circ P_\varepsilon$.

Theorem 2. The above construction of P_w meets the definition in Definition 4.

Since in c(x), we only care about whether P_w accepts α or not, the output of P_w can be ignored. Removing all the output labels from P_w gives us a finite-state automaton, hereafter named A_w .

3.5 ONE-PASS FSA FOR MASK SELECTION

After constructing $P_w(\alpha)$ for all realizable terminal sequences $w \in R(\mathcal{V})$, we can now consider simplifying the constraint calculation over the whole vocabulary \mathcal{V} . Recall Definition 1, combined with Simplification 3 and A_w , we have the following equation,

$$c(x \in \Sigma^*) = \{v \in \mathcal{V} | \exists w \in R_q(v), P_w(\alpha) \neq \emptyset\} = \{v \in \mathcal{V} | \exists w \in R_q(v), \alpha \in A_w\},\$$

where q and α as defined in Section 3.2 are only dependent on x.

In c(x), we want to know whether any of the A_w accepts α , where $w \in R_q(v)$. This can be achieved by constructing the union of different A_w , i.e., $\bigcup_{w \in R_q(v)} A_w$.

For different tokens $v \in \mathcal{V}$, we need to check whether α is accepted by $\bigcup_{w \in R_q(v)} A_w$. But to get the whole mask c(x), we need to check for every $v \in \mathcal{V}$, which is inefficient. To address this issue, we can integrate the checking of v and q into the FSA. Introduce the notation B_a for an FSA that accepts only a once. For every $q \in Q$ and $v \in \mathcal{V}$, we can concatenate B_q before $\bigcup_{w \in R_q(v)} A_w$ to check whether the current state is q, and then concatenate B_v after $\bigcup_{w \in R_q(v)} A_w$ to check whether the candidate token is v.

By unioning the results for all $v \in \mathcal{V}$ and $q \in Q$, we construct an FSA \mathcal{A} that accepts the sequence $q\alpha v$ only if v is a valid token for the lexer state q and stack α ,

$$\mathcal{A} := \bigcup_{v \in \mathcal{V}} \bigcup_{q \in Q} \bigcup_{w \in R_q(v)} B_q A_w B_v, \quad c(x) = \{ v \in \mathcal{V} | q \alpha v \in \mathcal{A} \},$$
 (5)

where A should be determinized and minimized.

Theorem 3. All (lexer state, parser stack) pairs that accept a given token form a regular language.

In Equation 5, we can precompute all possible result of c, i.e. all possible vocabulary masks, by considering acceptable vocabulary set $\mathcal{A}_s := \left\{ v \in \mathcal{V} \middle| s \xrightarrow{v}_{\mathcal{A}} f^{\mathcal{A}} \right\}$ for every state s in \mathcal{A} .

We summarize the offline construction process of PSC in Algorithm 3, and the online execution process in Algorithm 4. In Algorithm 4, both the lexing Step 2 and the parsing Step 3 are standard in grammar-constrained decoding, and can be incrementally maintained. In Step 4, the FSA $\mathcal A$ is run on the stack α and the lexer state q to get the state s, only requiring $\mathcal O(|\alpha|)$ time. Step 5 can be precomputed to be $\mathcal O(1)$ at runtime.

Algorithm 4 Online execution of PSC

1: **function** ONLINEEXECUTION($\mathcal{T}, \mathcal{P}, \mathcal{A}, x$)
2: $q_0^{\mathcal{T}} \xrightarrow{x:z} q$ 3: $\gamma_0 \xrightarrow{z} \alpha$ 4: $q_0^{\mathcal{A}} \xrightarrow{q\alpha} s$ 5: **return** \mathcal{A}_s

4 EXPERIMENTS

4.1 EXPERIMENTAL SETUP

All grammar-constrained decoding methods essentially perform the same task: compute the valid token mask at each decoding step. The valid token mask is theoretically the same for all methods, so we focus on comparing the efficiency of mask computation in our experiments. The usefulness of grammar-constrained decoding is shown in previous work (Geng et al., b; Scholak et al.; Poesia et al.; Ugare et al.), so we do not repeat the experiments here.

We conduct two sets of experiments to evaluate each method:

- Overhead of mask computation (without model inference).
- End-to-end throughput (with model inference).

In all experiments, we use **teacher-forcing** during evaluation, i.e., we always use the oracle next token at each decoding step, **to ensure that all methods are evaluated under the same conditions and can be fairly compared.**

Datasets There is no standard benchmark for evaluating grammar-constrained decoding methods. We choose two representative tasks that require grammar-constrained decoding: code generation in Java, Go, and SQL, and JSON generation with specified JSON schemas. For each task, we construct the evaluation dataset as described in Appendix A.4, with 1000 samples for each programming language and 1000 schemas for schema-conformant JSON generation. There are a total of 1337 positive samples and 2072 negative samples in the JSON dataset, and each schema has at least one positive sample and one negative sample.

Baselines We consider several recent state-of-the-art grammar-constrained decoding methods with open-source implementations as baselines, including XGrammar (Dong et al.), GreatGramma (Park et al.), Formatron (Sun et al.), and LLGuidance (Moskal et al.). Detailed descriptions of these baselines are included in Appendix A.5.

Implementation We implement PSC in roughly 1100 lines of Python. Similar to previous work (Ugare et al.; Park et al.), we use the Lark parser(lar) to construct the lexer and the LALR(1) parser from the grammar. As described in Section 3.1, we reuse the lexer construction in Great-Gramma (Park et al.), since this is not our focus in this paper.

Execution environment We conduct our experiments on a machine with 8 NVIDIA A100 GPUs (40 GB Memory), 2 Intel Xeon Gold 6348 CPUs (2.6GHz, 56 cores), and 512 GB RAM. To ensure fairness, we run all the experiments with a single GPU and a single CPU thread.

4.2 Overhead of Mask Computation

Metrics In this experiment, we measure the overhead of mask computation for each GCD method. We ignore the time taken to transfer the mask to the GPU and apply it to the logits, because this is the same for all methods². For each method, we measure:

- Average overhead: The time of computing the CPU mask tensor at each decoding step.
- Sample pass rate: The proportion of samples that are correctly processed by each method³.

Models We evaluate the overhead on three open-source LLM series with different vocabulary sizes: Llama 3 (Grattafiori et al.) (128k vocabulary size), Qwen 2.5 (Bai et al.) (151k vocabulary size), and Gemma 3 (Gemma Team et al.) (262k vocabulary size). We only use the tokenizers, because the overhead of mask computation is independent of other model components.

Overhead results The average overhead is shown in Table 1. PSC significantly outperforms all the baselines across all grammars and models, being 310 to 700 times faster on complex programming language grammars compared to the best baseline, LLGuidance, and generally 30 times faster on the relatively simple JSON schema grammar.

The overhead of baselines is significantly larger on Gemma 3 than that on the other models⁴, because Gemma 3 has a much larger vocabulary size (262k) than the other two models (128k and 151k), and the time complexity of all methods except PSC is roughly linear to the vocabulary size. In contrast, the performance of PSC is stable across different grammars and models, because its time complexity is independent of the vocabulary size.

Sample pass rate results The sample pass rates are shown in Table 2. PSC achieves a high pass rate across all grammars and models, very close to 100%. After analyzing the error cases, we find that they are all directly rejected by the GreatGramma lexer we adopt, but **no error is caused by PSC itself**. We notice the strangely low sample pass rate of Formatron on SQL, and find that Formatron refuses to accept the column alias in the query, resulting in frequent rejections.

²In PSC, one can preload all the mask tensors on the GPU memory before decoding begins, eliminating the transfer overhead. However, the transfer overhead is usually too small to justify the extra GPU memory usage.

³If the oracle token is masked, the sample is counted as rejected by the method. A positive sample is counted as passed if it is not rejected, and a negative sample is counted as passed if it is rejected.

⁴The performance of LLGuidance on JSON Schemas is roughly the same across different models, probably because the grammar is simple enough that the vocabulary size does not significantly affect the performance.

Table 1: Average overhead (microseconds) of computing the mask per token in GCD tasks using different methods. The symbol **X** indicates the parser reports an error during mask calculation. Text in **bold** indicates the best performance, and text in <u>underline</u> indicates the second best performance.

	Method	Grammar			
Model		Java	Go	SQL	JSON Schemas
Llama 3 $ \mathcal{V} = 128256$	XGrammar	309514.2	281500.9	324663.6	26257.6
	Formatron	393540.4	$\mathbf{X}^{\mathbf{a}}$	303974.1	X^a
	GreatGramma	21402.8	<u>27220.9</u>	20954.5	8556.2
	LLGuidance	1352.4	X^{b}	826.1	72.7
	PSC (Ours)	2.4	2.5	2.6	2.3
Qwen2.5 $ \mathcal{V} = 151665$	XGrammar	302421.9	278333.4	299968.6	29470.0
	Formatron	378921.0	X^a	253311.9	X^a
	GreatGramma	24570.9	<u>27649.8</u>	24053.9	11038.4
	LLGuidance	1408.2	X^b	865.2	<u>72.1</u>
	PSC (Ours)	2.4	2.5	2.5	2.2
Gemma 3 $ \mathcal{V} = 262145$	XGrammar	649321.1	458952.0	416625.0	54164.4
	Formatron	696144.6	$\mathbf{X}^{\mathbf{a}}$	444810.0	X^a
	GreatGramma	43218.6	48354.0	40026.9	25717.3
	LLGuidance	<u>1802.6</u>	X^b	1180.8	<u>72.1</u>
	PSC (Ours)	2.3	2.5	2.4	2.3

^a Formatron stops responding on multiple samples, so we terminate the process.

Table 2: How many samples are correctly processed by each method. The symbol X indicates the parser reports an error during mask calculation. Text in **bold** indicates the best performance, and text in underline indicates the second best performance.

		l .			
Model	Method	Java	Go	rammar SQL	JSON Schemas
		Java	GU	3QL	JOON Schemas
Llama 3	XGrammar	99.7%	100.0%	<u>99.7%</u>	100.0%
	Formatron	99.4%	X	67.0%	X
	GreatGramma	100.0%	99.9%	97.1%	99.6%
	LLGuidance	100.0%	X	99.9%	99.9%
	PSC (Ours)	100.0%	<u>99.9%</u>	99.9%	99.6%
Qwen2.5	XGrammar	99.7%	100.0%	99.7%	100.0%
	Formatron	99.4%	X	67.0%	X
	GreatGramma	100.0%	99.9%	97.2%	99.6%
	LLGuidance	100.0%	X	99.9%	99.9%
	PSC (Ours)	100.0%	<u>99.9%</u>	99.9%	99.6%
Gemma 3	XGrammar	99.7%	100.0%	99.7%	100.0%
	Formatron	100.0%	X	67.3%	X
	GreatGramma	100.0%	99.9%	97.2%	99.6%
	LLGuidance	94.0%	X	99.9%	99.9%
	PSC (Ours)	100.0%	<u>99.9%</u>	99.9%	99.6%

4.3 END-TO-END THROUGHPUT

Settings In this experiment, we run the actual model inference using the vLLM library (Kwon et al.), and measure the throughput, i.e., the number of tokens processed per second, of the entire decoding process, considering both the time taken by mask computation and model inference. We only consider the accepted samples of each method. We compare PSC with the fastest baseline LLGuidance in the previous section, and also include the throughput when not using any constraint decoding method as a reference, under various batch sizes of 1, 2, 4, ..., 256.

b LLGuidance reports ParserTooComplex error.

Models We use the smallest models in the three model series to highlight the overhead of constraint decoding. On these models, the model inference time is relatively small, making the overhead of constraint decoding more pronounced. Specifically, we use Llama 3 1B, Qwen 2.5 0.5B, and Gemma 3 270M. We also include Qwen 2.5 7B to see the effect of model size on throughput.

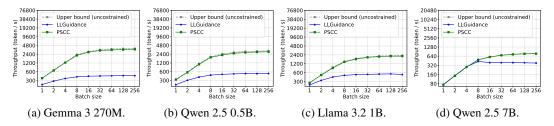


Figure 2: End-to-end throughput (tokens per second) on the **Java** dataset using different methods on different models with various batch sizes.

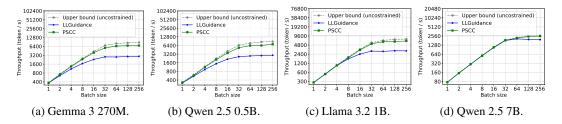


Figure 3: End-to-end throughput (tokens per second) on the **schema-conformant JSON** dataset using different methods on different models with various batch sizes.

Results The end-to-end throughput results on Java and schema-conformant JSON datasets are present in Figures 2 and 3, respectively. The results on the Go and SQL datasets are similar to those on the Java dataset, included in Appendix A.6.

On all datasets, PSC consistently outperforms LLGuidance across all models and batch sizes, and is very close to the performance of unconstrained decoding. The difference is more pronounced on the Java dataset, where the grammar is more complex, leading to higher overhead for LLGuidance.

As the model size increases, the difference in throughput becomes smaller, because the model inference time becomes more dominant. However, since smaller models are less capable, grammarconstrained decoding is probably more useful for smaller models to ensure the syntactic correctness.

As the batch size increases, the difference in throughput becomes larger, because the average model inference time per token decreases with larger batch sizes, indicating that the overhead introduced by constraint decoding becomes more pronounced at larger batch sizes.

5 DISCUSSION AND CONCLUSION

In this paper, we present PSC, a novel approach for grammar-constrained decoding. By constructing the exact requirements on the parser stack for each vocabulary token, PSC can determine the valid tokens at each decoding step by a single pass through the parser stack. Our experimental results demonstrate that PSC achieves significant speedup over existing methods, and the end-to-end throughput of PSC approaches that of unconstrained decoding. This makes PSC a practical choice for real-world applications that require grammar-constrained decoding.

One limitation of PSC is that it requires preprocessing. But this can be done offline and reused for multiple decoding sessions with the same grammar and vocabulary, so the overhead of preprocessing is generally acceptable for practical applications. In future, we plan to explore other types of grammars and constraints that can be efficiently handled by PSC, as well as other optimizations to further improve its efficiency and scalability.

6 REPRODUCIBILITY STATEMENT

We open-source the code to facilitate reproducibility of our results at https://anonymous.4open.science/r/PSC-E43E. It includes the implementation of PSC, the datasets used in the experiments, and the scripts to run the experiments and generate the results in the paper. The proofs of statements in the main text are included as Appendix A.3. The details in dataset construction are included as Appendix A.4.

REFERENCES

- Lark a parsing toolkit for Python. URL https://github.com/lark-parser/lark.
- 498 MaskBench. URL https://github.com/guidance-ai/jsonschemabench/tree/ 499 main/maskbench.
 - Alfred V. Aho and Jeffrey D. Ullman. *The Theory of Parsing, Translation, and Compiling*. Prentice-Hall, Inc. ISBN 978-0-13-914556-8.
 - Anysphere Inc. Cursor The AI Code Editor. URL https://cursor.com/en.
 - Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Benfeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen Technical Report. URL http://arxiv.org/abs/2309.16609.
 - Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Guiding LLMs the right way: Fast, non-invasive constrained generation. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org.
 - Janusz A. Brzozowski. Derivatives of regular expressions. J. ACM, 11(4):481–494, October 1964. ISSN 0004-5411. doi: 10.1145/321239.321249. URL https://doi.org/10.1145/321239.321249.
 - Didier Caucal and Roland Monfort. On the transition graphs of automata and grammars. In Rolf H. Möhring (ed.), *Graph-Theoretic Concepts in Computer Science*, pp. 311–337. Springer. ISBN 978-3-540-46310-8. doi: 10.1007/3-540-53832-1_51.
 - Yixin Dong, Charlie F. Ruan, Yaxing Cai, Ziyi Xu, Yilong Zhao, Ruihang Lai, and Tianqi Chen. XGrammar: Flexible and efficient structured generation engine for large language models. In *Eighth Conference on Machine Learning and Systems*. URL https://openreview.net/forum?id=rjQfX0YgDl.
 - Jay Earley. An efficient context-free parsing algorithm. doi: 10.1145/362007.362035. URL https://dl.acm.org/doi/10.1145/362007.362035.
 - Philip Gage. A new algorithm for data compression. 12(2):23–38. ISSN 0898-9788.
 - Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey Cideron, Jean-bastien Grill, Sabela Ramos, Edouard Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev, Gaël Liu, Francesco Visin, Kathleen Kenealy, Lucas Beyer, Xiaohai Zhai, Anton Tsitsulin, Robert Busa-Fekete, Alex Feng, Noveen Sachdeva, Benjamin Coleman, Yi Gao, Basil Mustafa, Iain Barr, Emilio Parisotto, David Tian, Matan Eyal, Colin Cherry, Jan-Thorsten Peter, Danila Sinopalnikov, Surya Bhupatiraju, Rishabh Agarwal, Mehran Kazemi, Dan Malkin, Ravin Kumar, David Vilar, Idan Brusilovsky, Jiaming Luo, Andreas Steiner, Abe Friesen, Abhanshu Sharma, Abheesht Sharma, Adi Mayrav Gilady, Adrian Goedeckemeyer, Alaa Saade, Alex Feng, Alexander Kolesnikov, Alexei Bendebury, Alvin Abdagic, Amit Vadi, András György, André Susano Pinto, Anil Das, Ankur Bapna, Antoine Miech, Antoine Yang, Antonia Paterson,

541

543

544

546

547

548

549

550

551

552

553

554

558

559

561

562

565

566

567 568

569

570

571

572

573

574

575 576

577

578

579

582

583

584

585

586

588

592

Ashish Shenoy, Ayan Chakrabarti, Bilal Piot, Bo Wu, Bobak Shahriari, Bryce Petrini, Charlie Chen, Charline Le Lan, Christopher A. Choquette-Choo, C. J. Carey, Cormac Brick, Daniel Deutsch, Danielle Eisenbud, Dee Cattle, Derek Cheng, Dimitris Paparas, Divyashree Shivakumar Sreepathihalli, Doug Reid, Dustin Tran, Dustin Zelle, Eric Noland, Erwin Huizenga, Eugene Kharitonov, Frederick Liu, Gagik Amirkhanyan, Glenn Cameron, Hadi Hashemi, Hanna Klimczak-Plucińska, Harman Singh, Harsh Mehta, Harshal Tushar Lehri, Hussein Hazimeh, Ian Ballantyne, Idan Szpektor, Ivan Nardini, Jean Pouget-Abadie, Jetha Chan, Joe Stanton, John Wieting, Jonathan Lai, Jordi Orbay, Joseph Fernandez, Josh Newlan, Ju-yeong Ji, Jyotinder Singh, Kat Black, Kathy Yu, Kevin Hui, Kiran Vodrahalli, Klaus Greff, Linhai Qiu, Marcella Valentine, Marina Coelho, Marvin Ritter, Matt Hoffman, Matthew Watson, Mayank Chaturvedi, Michael Moynihan, Min Ma, Nabila Babar, Natasha Noy, Nathan Byrd, Nick Roy, Nikola Momchev, Nilay Chauhan, Noveen Sachdeva, Oskar Bunyan, Pankil Botarda, Paul Caron, Paul Kishan Rubenstein, Phil Culliton, Philipp Schmid, Pier Giuseppe Sessa, Pingmei Xu, Piotr Stanczyk, Pouya Tafti, Rakesh Shivanna, Renjie Wu, Renke Pan, Reza Rokni, Rob Willoughby, Rohith Vallu, Ryan Mullins, Sammy Jerome, Sara Smoot, Sertan Girgin, Shariq Iqbal, Shashir Reddy, Shruti Sheth, Siim Põder, Sijal Bhatnagar, Sindhu Raghuram Panyam, Sivan Eiger, Susan Zhang, Tianqi Liu, Trevor Yacovone, Tyler Liechty, Uday Kalra, Utku Evci, Vedant Misra, Vincent Roseberry, Vlad Feinberg, Vlad Kolesnikov, Woohyun Han, Woosuk Kwon, Xi Chen, Yinlam Chow, Yuvein Zhu, Zichuan Wei, Zoltan Egyed, Victor Cotruta, Minh Giang, Phoebe Kirk, Anand Rao, Kat Black, Nabila Babar, Jessica Lo, Erica Moreira, Luiz Gustavo Martins, Omar Sanseviero, Lucas Gonzalez, Zach Gleicher, Tris Warkentin, Vahab Mirrokni, Evan Senter, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, Yossi Matias, D. Sculley, Slav Petrov, Noah Fiedel, Noam Shazeer, Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray Kavukcuoglu, Clement Farabet, Elena Buchatskaya, Jean-Baptiste Alayrac, Rohan Anil, Dmitry, Lepikhin, Sebastian Borgeaud, Olivier Bachem, Armand Joulin, Alek Andreev, Cassidy Hardin, Robert Dadashi, and Léonard Hussenot. Gemma 3 Technical Report. URL http://arxiv.org/abs/2503.19786.

Saibo Geng, Hudson Cooper, Michał Moskal, Samuel Jenkins, Julian Berman, Nathan Ranchin, Robert West, Eric Horvitz, and Harsha Nori. JSONSchemaBench: A Rigorous Benchmark of Structured Outputs for Language Models, a. URL http://arxiv.org/abs/2501.10868.

Saibo Geng, Martin Josifoski, Maxime Peyrard, and Robert West. Grammar-Constrained Decoding for Structured NLP Tasks without Finetuning. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 10932–10952. Association for Computational Linguistics, b. doi: 10.18653/v1/2023. emnlp-main.674. URL https://aclanthology.org/2023.emnlp-main.674/.

GitHub. GitHub Copilot · Your AI pair programmer. URL https://github.com/features/copilot.

Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind Thattai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, prefix=van der useprefix=false family=Linde, given=Jelmer, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhotia, Lauren Rantala-Yeary, prefix=van der useprefix=false

595

596

597

598

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

622

623

625

626

627

629

630

631

632

633

634

635

636

637

638

639

640

641

642

644

645

646

family=Maaten, given=Laurens, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, prefix=de useprefix=false family=Oliveira, given=Luke, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vítor Albiero, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesco Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manay Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Nikolay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Rangaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin

Battey, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinman, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. The Llama 3 Herd of Models. URL http://arxiv.org/abs/2407.21783.

- John E. Hopcroft and Jeffrey D. Ullman. Introduction to Automata Theory, Languages, and Computation. Addison-Wesley Series in Computer Science and Information Processing. Addison-Wesley. ISBN 978-0-201-02988-8.
- Denis Kocetkov, Raymond Li, Loubna Ben Allal, Jia Li, Chenghao Mou, Yacine Jernite, Margaret Mitchell, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Dzmitry Bahdanau, Leandro Von Werra, and prefix=de useprefix=false family=Vries, given=Harm. The Stack: 3 TB of permissively licensed source code. ISSN 2835-8856. URL https://openreview.net/forum?id=pxpbTdUEpD.
- Terry Koo, Frederick Liu, and Luheng He. Automata-based constraints for language model decoding. In *First Conference on Language Modeling*. URL https://openreview.net/forum?id=BDBdblmyzY.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient Memory Management for Large Language Model Serving with PagedAttention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, SOSP '23, pp. 611–626. Association for Computing Machinery. ISBN 979-8-4007-0229-7. doi: 10.1145/3600006.3613165. URL https://dl.acm.org/doi/10.1145/3600006.3613165.
- Michael Xieyang Liu, Frederick Liu, Alexander J Fiannaca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie J Cai. "we need structured output": Towards user-centered constraints on large language model output. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems*, pp. 1–9, 2024.
- Michał Moskal, Harsha Nori, Hudson Cooper, and Loc Huynh. LLGuidance: Making Structured Outputs Go Brrr. URL https://guidance-ai.github.io/llguidance/llg-go-brrr.
- OpenAI. Structured model outputs openai api. URL https://platform.openai.com.
- Kanghee Park, Timothy Zhou, and Loris D'Antoni. Flexible and efficient grammar-constrained decoding. In *Forty-Second International Conference on Machine Learning*. URL https://openreview.net/forum?id=L6CYAzpO1k.
- Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. Synchromesh: Reliable Code Generation from Pre-trained Language Models. URL https://openreview.net/forum?id=KmtVD97J43e.
- Qwen, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,

- Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 Technical Report. URL http://arxiv.org/abs/2412.15115.
- Torsten Scholak, Nathan Schucher, and Dzmitry Bahdanau. PICARD: Parsing Incrementally for Constrained Auto-Regressive Decoding from Language Models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 9895–9901. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.779. URL https://aclanthology.org/2021.emnlp-main.779/.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural Machine Translation of Rare Words with Subword Units. In Katrin Erk and Noah A. Smith (eds.), *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1715–1725. Association for Computational Linguistics. doi: 10.18653/v1/P16-1162. URL https://aclanthology.org/P16-1162/.
- Michael Sipser. *Introduction to the Theory of Computation*. Cengage Learning, third edition, international edition edition. ISBN 978-1-133-18779-0 978-1-133-18781-3 978-0-357-67058-3.
- Xintong Sun, Chi Wei, Minghao Tian, and Shiwen Ni. Earley-driven dynamic pruning for efficient structured decoding. In *Forty-Second International Conference on Machine Learning*. URL https://openreview.net/forum?id=6hDNXCdTsE.
- Shubham Ugare, Tarun Suresh, Hangoo Kang, Sasa Misailovic, and Gagandeep Singh. Syn-Code: LLM generation with grammar augmentation. ISSN 2835-8856. URL https://openreview.net/forum?id=HiUZtgAPoH.
- vLLM Team. Structured Outputs vLLM. URL https://docs.vllm.ai/en/stable/features/structured_outputs.html.

A APPENDIX

A.1 LLM USAGE

We used DeepSeek and GitHub Copilot for code assistance, paper writing, and proofreading. However, all the technical content, ideas, algorithms, and experimental results in this paper are our own work. We carefully reviewed and verified all the content generated by LLMs to ensure they are accurate and directly reflect our own ideas.

A.2 FORMAL DEFINITIONS AND NOTATIONS

Table 3: Symbols and their meaning.

Symbol	Meaning
ε	empty string
ab	concatenation of strings a and b
AB	concatenation of languages A and B
$A_{arepsilon}$	$A \cup \{\varepsilon\}$
A^*	Kleene star of language A
A^+	$A^*\setminus\{arepsilon\}$
Σ	the character set (usually the Unicode) used by the language model
Γ	the terminal set of the grammar
\mathcal{V}	the vocabulary of the language model, a finite subset of Σ^+ such that every string over Σ can be tokenized as a string over $\mathcal V$
${\mathcal T}$	the lexing FST, transduces string over Σ to terminal sequence over Γ
Q	the finite set of states of the lexing FST
$Q \\ \mathcal{P}$	the parsing PDA, accepts terminal sequences in the language
П	the stack alphabet of the PDA

We include a list of symbols and their meanings in Table 3 for reference.

A.2.1 FINITE-STATE TRANSDUCER (FST)

A finite-state transducer (FST) (Aho & Ullman) \mathcal{T} is defined by a finite set of states Q, the input alphabet Σ , the output alphabet Γ , the start state $q_0 \in Q$, the final states $F \subseteq Q$, and transitions $\delta: Q \times \Sigma_{\varepsilon} \to 2^{\Gamma^* \times Q}$. If $\delta(q, a) \ni (y, q')$, we write $q \xrightarrow{a:y} q'$. We write \to^* for consecutive

transitions. For $q \in Q$, we write $q \xrightarrow[\tau]{\varepsilon:\varepsilon} q$. For $q \xrightarrow[\tau]{s:x} q'$ and $q' \xrightarrow[\tau]{t:y} q''$, we write $q \xrightarrow[\tau]{st:xy} q''$.

Informally, an FST is *deterministic*, if from any state, for any given string, there is exactly one possible outcome. $\mathcal T$ is *deterministic* if, for all $q \in Q$, either $|\delta(q,a)| \leq 1, \forall a \in \Sigma$ and $\delta(q,\varepsilon) = \varnothing$, or $\delta(q,a) = \varnothing, \forall a \in \Sigma$ and $\delta(q,\varepsilon) = 1$.

For $w \in \Sigma^*, q \in Q$, we define $\mathcal{T}_q(w)$ as $\left\{v \in \Gamma^* \middle| \exists q' \in F, q \xrightarrow[\mathcal{T}]{w:v}^* q'\right\}$, meaning the possible outcomes when we feed w into \mathcal{T} starting from the state q. When \mathcal{T} is deterministic and $v \in \mathcal{T}(w)$, we also write $\mathcal{T}(w) = v$. For $W \subseteq \Sigma^*$, we define $\mathcal{T}_q(W)$ as $\bigcup_{w \in W} \mathcal{T}_q(w)$. q defaults to q_0 when omitted.

We call a state $q \in Q$ stable if the FST does not need to take any immediate action on q, i.e. $\delta(q,\varepsilon) = \varnothing$. If $q \xrightarrow[]{s:t} q'$ and q' is stable, we also write $q \xrightarrow[]{s:t} q'$.

Given two FSTs S and T where the output alphabet of S is the input alphabet T, $\Gamma^S = \Sigma^T$, their *composition* is a new FST $S \circ T$ by feeding the output of S into the input of T.

A.2.2 FINITE-STATE AUTOMATON (FSA)

 A finite-state automaton (FSA) \mathcal{A} can be defined by removing all output labels from an FST. We say \mathcal{A} accepts $w \in \Sigma^*$ from state q if $\mathcal{A}_q(w) \neq \emptyset$, and write $w \in \mathcal{A}_q$, and q defaults to q_0 when omitted. Two FSAs are *equivalent* if they accept the same language.

 \mathcal{A} is *deterministic* if there is no ε transition in δ . Every nondeterministic FSA can be *determinized* into an equivalent deterministic FSA (Hopcroft & Ullman), and every deterministic FSA can be *minimized* into an equivalent deterministic FSA with the smallest number of states (Hopcroft & Ullman).

The *union* of two FSAs \mathcal{A} and \mathcal{B} is a new FSA $\mathcal{A} \cup \mathcal{B}$ that accepts any sequence that is accepted by either \mathcal{A} or \mathcal{B} .

The concatenation of two FSAs A and B is a new FSA AB that accepts any sequence that can be split into two parts x = yz, where A accepts the first part y and B accepts the second part z.

A.2.3 PUSH-DOWN AUTOMATON (PDA)

A push-down automaton (PDA) (Aho & Ullman; Hopcroft & Ullman; Caucal & Monfort) \mathcal{P} is defined by the input alphabet Γ , the stack alphabet Π , the initial stack $\gamma_0 \in \Pi^2$, the final states $F \subseteq \Pi$, and a finite set of transitions $\delta : \Pi^2 \times \Gamma_{\varepsilon} \to 2^{\Pi^+}$. Note that the definition here merges the states and the stack symbols in the traditional definition of PDA, but they are equivalent if we treat the stack top symbol as the state. If $\delta(\alpha, a) \ni \beta$, we write $\alpha \xrightarrow{a} \beta$. If $\alpha \xrightarrow{a} \beta$, for any $\gamma \in \Pi^+$, we write $\gamma \xrightarrow{a} \beta \gamma$. We write $\gamma \xrightarrow{a} \beta \gamma$. If

also write $\alpha \gamma \xrightarrow[\mathcal{P}]{a} \beta \gamma$. We write \rightarrow^* for consecutive transitions. For $\gamma \in \Pi^+$, we write $\gamma \xrightarrow[\mathcal{P}]{\varepsilon} \gamma$. If $\alpha \xrightarrow[\mathcal{P}]{a} \beta$, $\beta \xrightarrow[\mathcal{P}]{w} \gamma$, we write $\alpha \xrightarrow[\mathcal{P}]{a} \gamma$.

Informally, a PDA is *deterministic*, if from any stack, for any given string, there is exactly one possible outcome. \mathcal{P} is *deterministic* if, for any $\alpha \in \Pi^2$, either $|\delta(\alpha,a)| \leq 1, \forall a \in \Gamma$ and $\delta(\alpha,\varepsilon) = \emptyset$, or $\delta(\alpha,a) = \emptyset, \forall a \in \Gamma$ and $|\delta(\alpha,\varepsilon)| = 1$.

A deterministic PDA is *terminating*, if for any stack, it does not make an endless sequence of ε -input transitions.⁵ Every deterministic PDA can be transformed into another equivalent deterministic terminating PDA (Sipser; Hopcroft & Ullman).

We call a state $\beta = \beta_1 \dots \beta_n \in \Pi^+$ stable if β_1 is a final state, i.e. $\beta_1 \in F$, or the PDA is waiting to read one more symbol, i.e. $\exists a \in \Gamma, \delta(\beta_1 \beta_2, a) \neq \emptyset$. If $\alpha \xrightarrow[\mathcal{P}]{w}^* \beta$ and β is stable, we also write $\alpha \xrightarrow[\mathcal{P}]{w}^* \beta$. In practice, parser in a stable state is ready to consume the next input symbol, or has

 $\alpha \Rightarrow \beta$. In practice, parser in a stable state is ready to consume the next input symmetric reached an accepting stack.

For $\alpha \in \Pi^+$, we define \mathcal{P}_{α} as $\left\{ w \in \Gamma^* \middle| \exists X \in F, \exists \gamma \in \Pi^*, \alpha \Longrightarrow_{\mathcal{P}}^* X \gamma \right\}$. α defaults to γ_0 when omitted.

A.3 PROOFS OF THEORETICAL RESULTS

In this section, we provide the proofs of the theoretical results stated in the main text.

A.3.1 PROOF OF THEOREM 1

Proof. First we show that the states of P_{ε} are finite, i.e. Algorithm 1 terminates. Consider the two steps in Algorithm 1 that add new states. Step 11 can only adds states in $\Pi \cup \Pi^2$, so step 11 is only executed a finite number of times. As for Step 16, because the parser \mathcal{P} is deterministic and terminating, by definition in Appendix A.2.3, for any stack configuration, there will not be an endless sequence of ε transitions, so Step 16 is also only executed a finite number of times.

⁵The definition is slightly different in the cited references; nevertheless, their proof works on this definition.

The correctness of Algorithm 1 can be naturally deduced by its construction, because it simply simulates the behavior of the parser \mathcal{P} with the current stack top, and only outputs the stack when it is stable.

A.3.2 PROOF OF THEOREM 2

Proof. Formally, \tilde{P}_a can be defined as follows:

$$\tilde{P}_a(\alpha \in \Pi^+) := \left\{ \beta \in \Pi^+ \middle| \alpha \xrightarrow[\mathcal{P}]{a} \beta \right\}. \tag{6}$$

The construction of \tilde{P}_a in Algorithm 2 trivially simulates one a-labelled transition of the parser \mathcal{P} on the current stack top, so it meets Definition 6.

The process of the parser processing $w=w_1\dots w_n$ can be decomposed into a sequence of unconditional ε -transitions, followed by w_i -labelled transitions, followed by another sequence of unconditional ε -transitions. Each w_i -labelled transition is simulated by the corresponding \tilde{P}_{w_i} , and the unconditional ε -transitions are handled by P_{ε} . Therefore, the composition $P_w=P_{\varepsilon}\circ\tilde{P}_{w_1}\circ P_{\varepsilon}\circ \cdots\circ P_{\varepsilon}\circ\tilde{P}_{w_n}\circ P_{\varepsilon}$ correctly simulates the parser $\mathcal P$ processing the terminal sequence w on the input stack α , and produces the stabilized output stack β if it exists.

A.3.3 PROOF OF THEOREM 3

Proof. Because \mathcal{A} is constructed as a FSA, the language recognized by \mathcal{A} is regular. The language of all valid (lexer state, parser stack) pairs for a given token $v \in \mathcal{V}$ can be obtained by reversing the language recognized by \mathcal{A} , taking the Brzozowski derivative (Brzozowski, 1964) with respect to v, and then reversing it back. Since the class of regular languages is closed under these operations Hopcroft & Ullman, the resulting language is also regular.

A.4 DATASET CONSTRUCTION DETAILS

Java, Go, and SQL We obtain their Lark grammars from the previous work Syncode (Ugare et al.). Because the grammar format for XGrammar and Formatron is different from Lark, we manually convert the Lark grammars to respective formats for each baseline. For each programming language, we take the first 1000 samples from the Stack dataset (Kocetkov et al.) that can be successfully parsed by the Lark parser to construct the evaluation dataset.

JSON Schemas We use the benchmark dataset MaskBench (mas), an extension of JSON Schema Bench (Geng et al., a) by adding schema conformant and non-conformant JSON instances to each schema. We generate the Lark grammar from the JSON schemas using the script provided in MaskBench, and only use the schemas in MaskBench where the Lark parser can successfully parse all the conformant JSON instances and reject all the non-conformant JSON instances. We then randomly sample 1000 schemas for evaluation.

A.5 DESCRIPTION OF BASELINES IN EXPERIMENTS

We describe the baselines used in our experiments in detail.

- XGrammar (Dong et al.) uses a character-level non-deterministic PDA⁶. For each state, it precomputes the context-independent accepted and rejected tokens, and only calls the parser for the context-dependent tokens.
- GreatGramma (Park et al.) uses a lexer and a parser. It converts each token into all possible
 terminals sequences and reduces the number of parser calls by sharing the parser calls
 among tokens with the same terminal sequence. After computing the accepted terminal
 sequences, it maps them back to the original tokens. It also precomputes the contextdependent and context-independent terminal sequences for each parser state.

⁶In the latest implementation that we use in the experiments, this has been changed to an Earley (Earley) parser.

- Formatron (Sun et al.) uses an Earley parser. It dynamically identifies and eliminates invalid or redundant parser states during parsing, and uses a state cache to speed up the repetitive parsing process.
- LLGuidance (Moskal et al.) uses a lexer and an Earley parser. It organizes the vocabulary into a trie, and skips the whole subtree if the prefix token is rejected. It also leverages the lexer on the vocabulary to pre-identify the terminal sequences.

A.6 EXTRA THROUGHPUT RESULTS

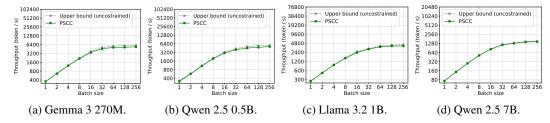


Figure 4: End-to-end throughput (tokens per second) on the **Go** dataset using different methods on different models with various batch sizes.

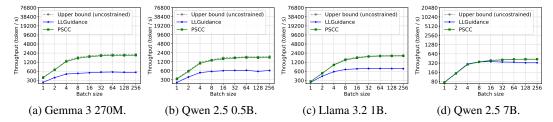


Figure 5: End-to-end throughput (tokens per second) on the **SQL** dataset using different methods on different models with various batch sizes.

Due to page limit, we only present the end-to-end throughput results on Java and schema-conformant JSON in Section 4.3. The results on the Go and SQL datasets are similar and are included here in Table 4 and Table 5, respectively.