



XGrammar: Flexible And Efficient Structured Generation Engine for Large Language Models

Yixin Dong, Charlie F. Ruan, Yaxing Cai, Ruihang Lai, Ziyi Xu,
Yilong Zhao, Tianqi Chen
Jan 28, 2024

Problem: LLM Generation with Structures

Code generation



GitHub
Copilot



Language Code

Function/tool calling

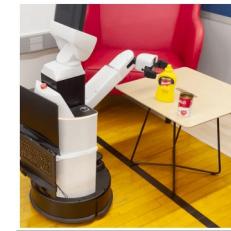


LangChain



JSON Schema

Embodied Agents



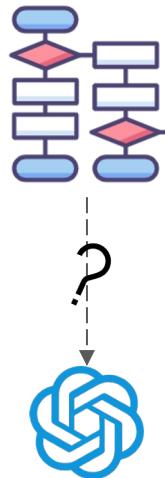
DSL (e.g. PDDL)

Structured Outputs



Fig: <https://arxiv.org/pdf/2304.11477>

Problem: LLM's Limited Ability with Complex Structures



Structures are increasingly complex

- Advanced agents
- Complex tool calling
- DSLs unfamiliar to LLMs

LLM's generation ability is limited

- On-device small LMs
- Compressed models



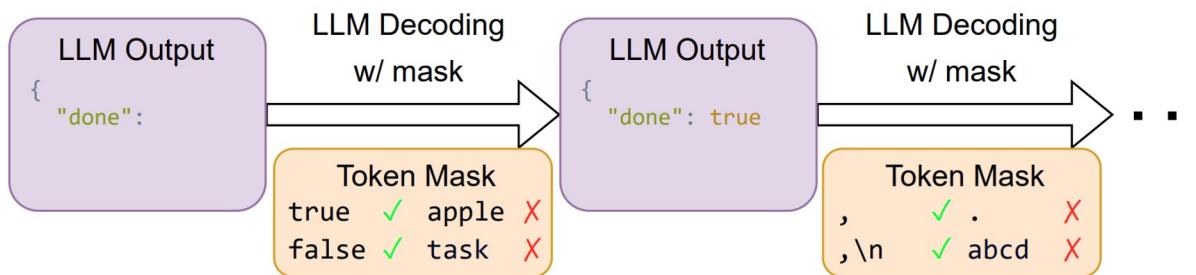
Background: Constrained Decoding

JSON Schema

```
class Task(BaseModel):
    done: bool
    name: str
    steps: List[int]
```

Example Valid JSONs

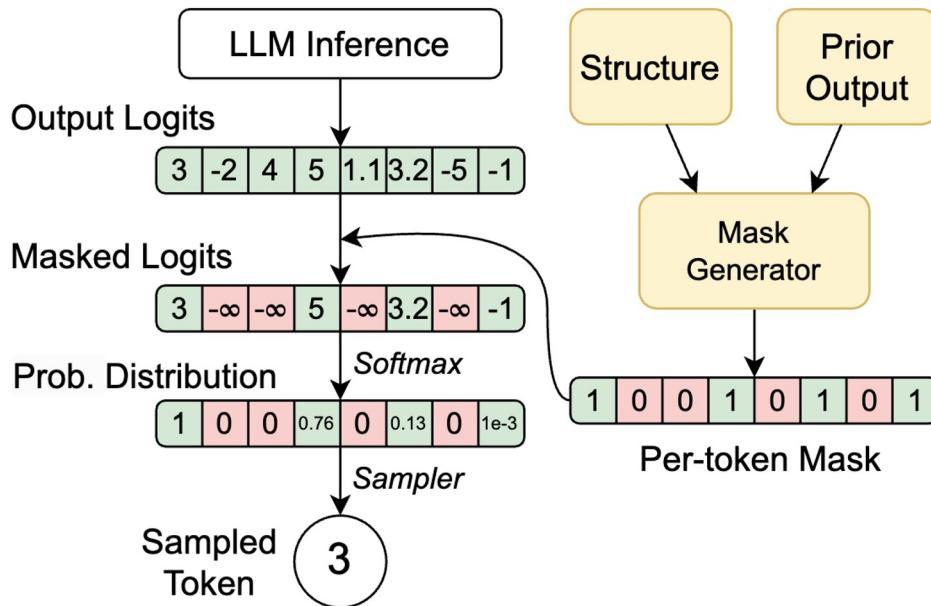
```
{
    "done": true,
    "name": "Clean kitchen",
    "steps": [1, 2, 3, 4]
}
{
    "done": false,
    "name": "Presentation",
    "steps": [1, 2]
}
```



An example of constrained decoding with JSON Schema



Background: Constrained Decoding



Apply a per-token mask to prevent generating invalid tokens according to the structure

The overhead of the mask generator is crucial!



XGrammar: Flexible and Efficient Structured Generation Engine

XGrammar is a structured generation library that features



Flexibility: Full support for context-free grammar



Efficiency: SOTA performance in constraint decoding

Zero-overhead JSON Schema generation



Integration: Easy to integrate with existing LLM serving frameworks

vLLM, MLC-LLM, SGLang, etc



XGrammar: Flexible and Efficient Structured Generation Engine

XGrammar is a structured generation library that features



Flexibility: Full support for context-free grammar



Efficiency: SOTA performance in constraint decoding

Zero-overhead JSON Schema generation



Integration: Easy to integrate with existing LLM serving frameworks

vLLM, MLC-LLM, SGLang, etc



More Powerful: Context-free Grammar

Context-free Grammar

```
root ::= <array> | <str>
```

```
array ::= '[' (<str> | <array> ',' )*  
        <str> | <array> ']'
```

```
str ::= '\"' [\"\\]* '\"'
```

Prior methods mainly support **regex** as input grammar

XGrammar support **the more powerful CFG**, therefore supporting

- Regex
- JSON, JSON Schema
- SQL
- Python (w/ additional state maintained)

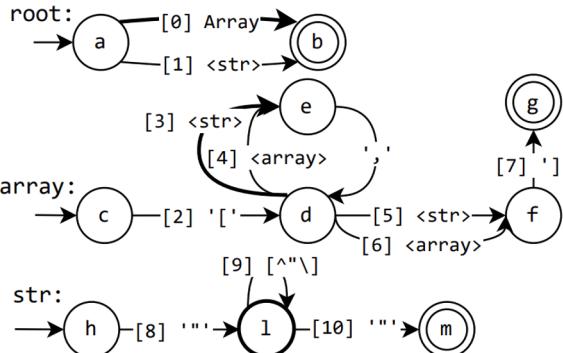


More Powerful: Pushdown Automata

Context-free Grammar

```
root ::= <array> | <str>  
  
array ::=  
  '[' (<str> | <array> ',' ')'*  
  <str> | <array> ']'  
  
str ::= '\"' [\"\\]* '\"'
```

Pushdown Automata



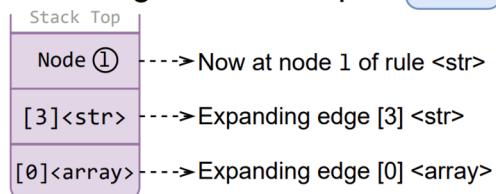
Example Accepted Strings

"abc"

["abc", ["def", "ghi"]]

["func1", [{"func2": "func3"}]]

Matching Stack for Input ["a



- XGrammar utilizes pushdown automata to match string to CFG
- The matching process is a **recursive expansion of rules**



XGrammar: Flexible and Efficient Structured Generation Engine

XGrammar is a structured generation library that features



Flexibility: Full support for context-free grammar



Efficiency: SOTA performance in constraint decoding

Zero-overhead JSON Schema generation

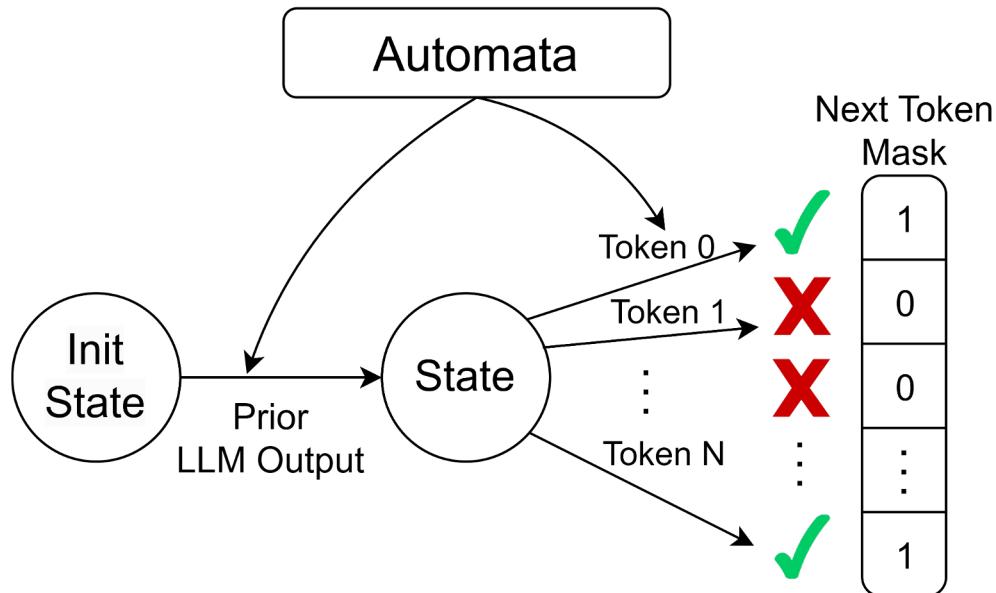


Integration: Easy to integrate with existing LLM serving frameworks

vLLM, MLC-LLM, SGLang, etc



Previous Mask Generation Method



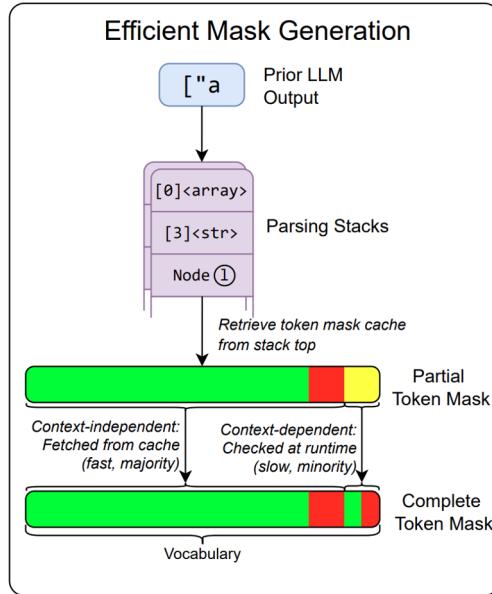
Check if each token matches,
then generate the mask

Every token checking is very slow!

**Our optimization: most token
checking can be fast via
preprocessing**

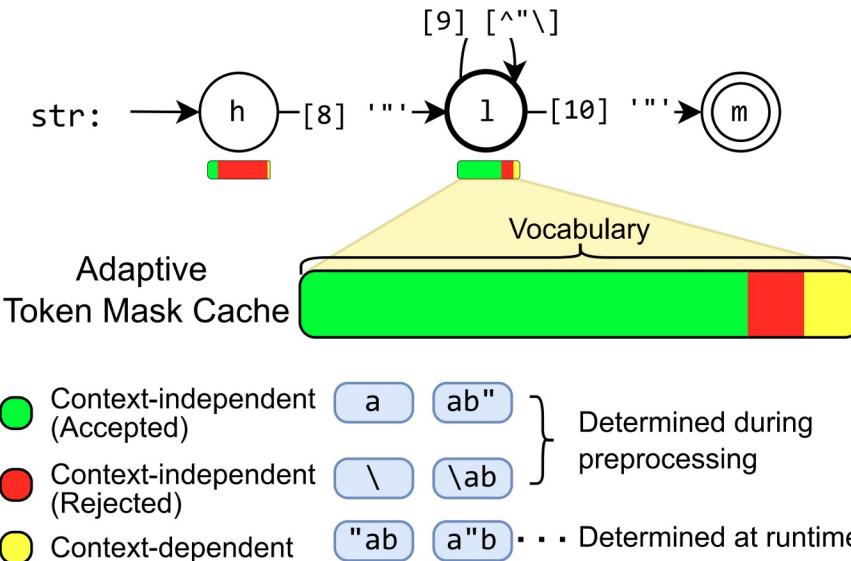


Optimization #1. The Adaptive Token Mask Cache



- At runtime, we first retrieve the pre-computed token mask
- Most the final mask is precomputed, and only a few are computed on the fly

The Adaptive Token Mask Cache (Cont'd)

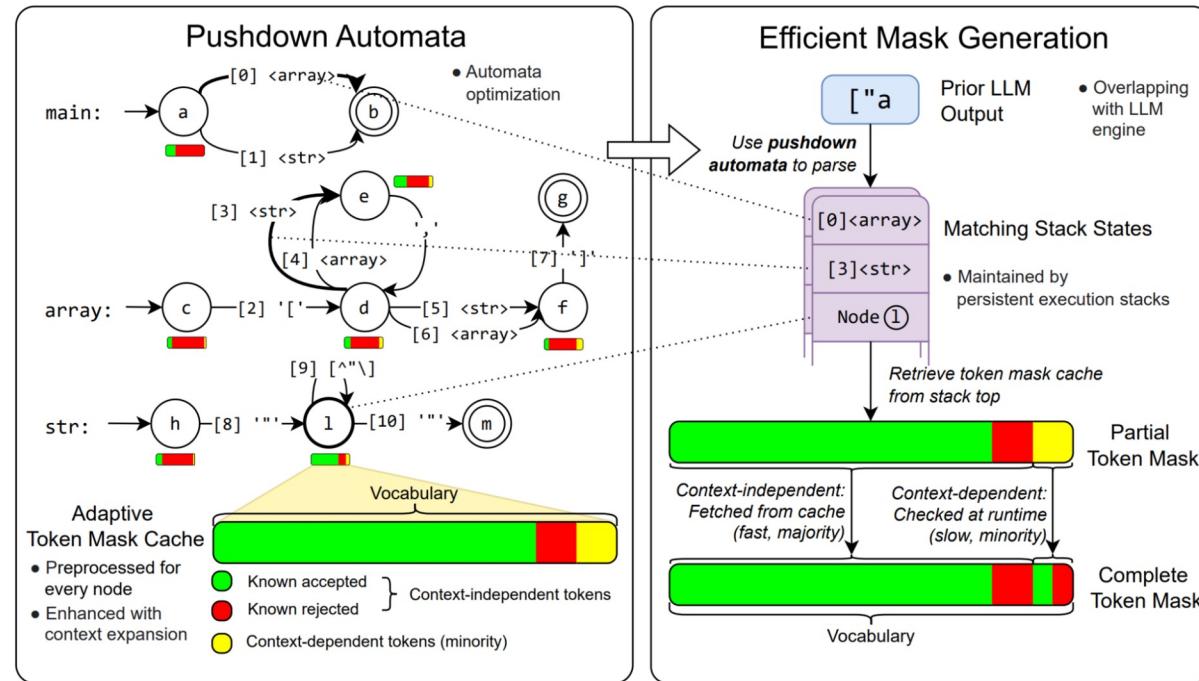


- Most tokens can be determined **ahead of time** – **context-independent tokens**
- Plus a minority of tokens that need to check at **runtime** – **context-dependent tokens**
- So before running, we can compile a **token mask cache** – for each node, calculate accept/reject for the context-independent tokens



Context-dependent tokens: less than 1% for Llama-3.1 w/ JSON grammar (1134 out of 128k)
XGrammar: Flexible And Efficient Structured Generation Engine for Large Language Models

The Adaptive Token Mask Cache (Cont'd)



The Adaptive Token Mask Cache (Storage)

- For every stack top node, we store accepted/rejected/uncertain tokens.
- Accepted + Rejected + Uncertain = Vocabulary!
- We want the size of the pre-computed mask cache to be small!
- Three storage paradigms:
 - a) **#accept** is large → [Rejected list], [Uncertain list]
 - b) **#reject** is large → [Accepted list], [Uncertain list]
 - c) **#accept** and **#reject** similar → <Accepted_bitset>, <Uncertain_bitset>
- Store the one with least memory consumption!



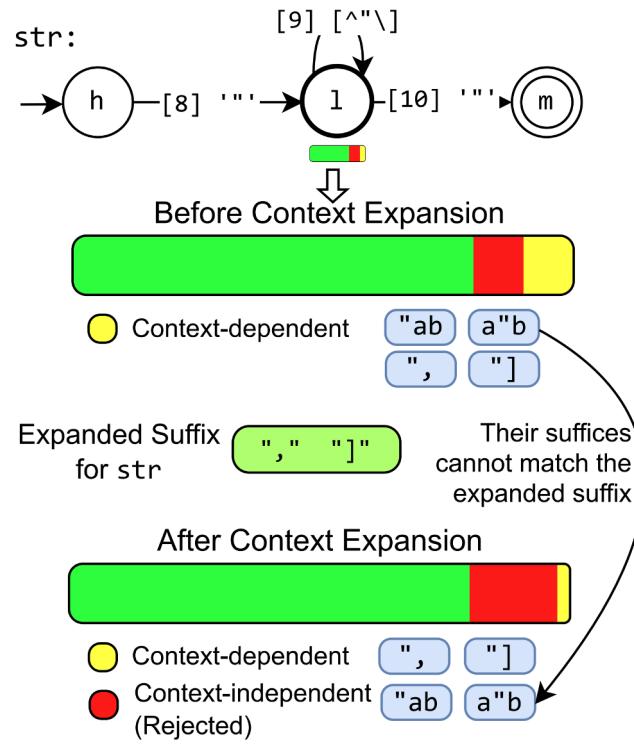
On Llama-3.1 w/ JSON grammar, reduces total memory usage to 0.2% (from 160 MB to 0.46

The Challenge of Llama-3

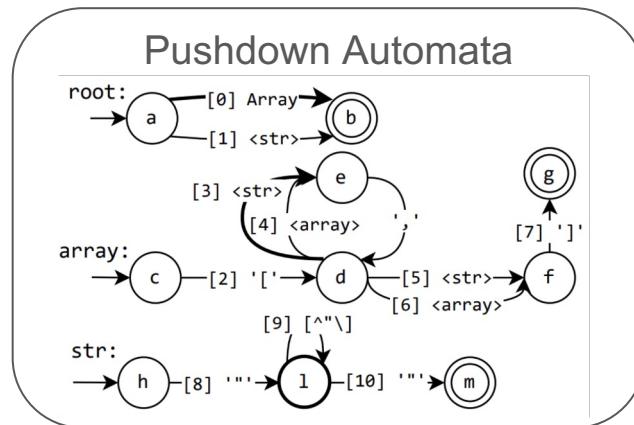
- Vocab_size: 32k → 128k
- More common tokens
- #(uncertain tokens): 100 → 1.5k
 - Means 15x check at runtime
- Further reduction of uncertain tokens is needed!



Optimization #2. Context Expansion



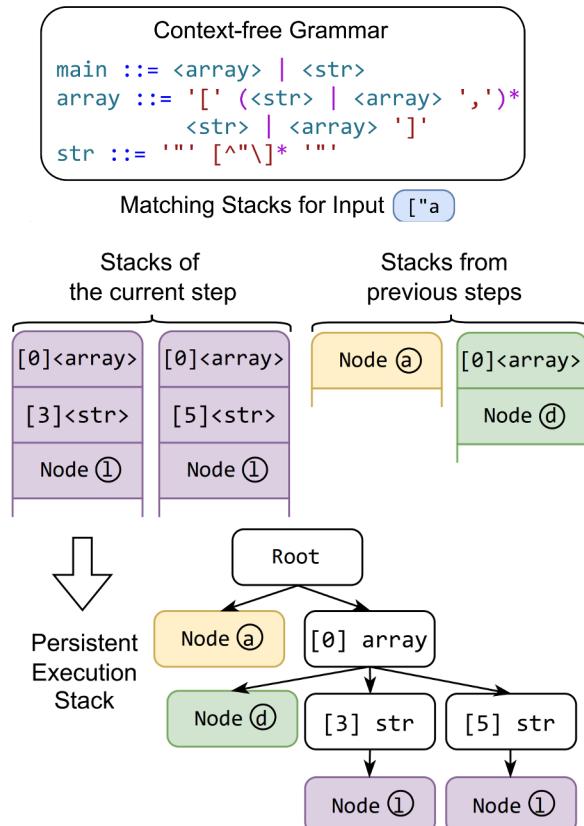
- For every rule, find all its references
- Collect all possible suffix of that rule
- If an uncertain token does not fit into any of the suffix, it is rejected



Reduce context-dependent tokens by **90%** on Llama-3.1 w/ JSON grammar



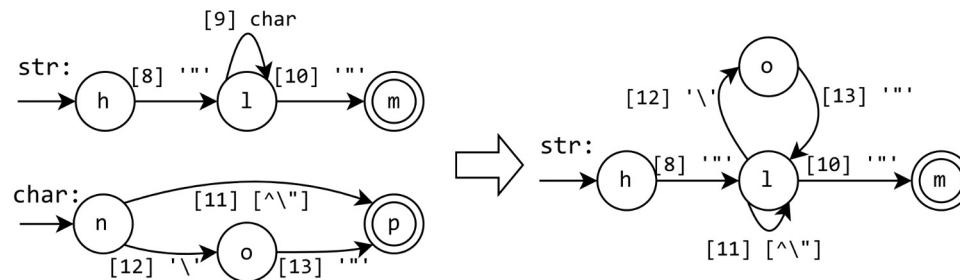
Optimization #3. Persistent Execution Stack



- Despite having most tokens pre-computed, we still need to compute **context-dependent tokens** efficiently
- As mentioned earlier, we can have multiple possible stacks due to the ambiguity
- We represent the stacks as a tree → **avoids memory redundancy** for storing multiple stacks
- Instead of copying the stack, we only split the branch
- Also supports rolling the state back



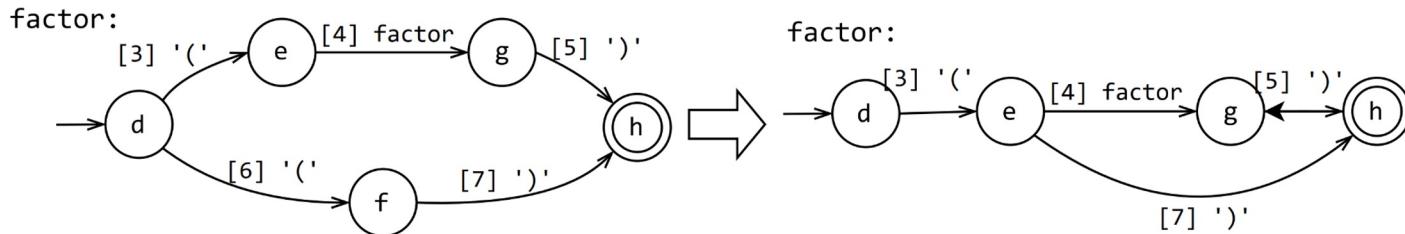
Optimizations #4. Optimization Passes (Inlining)



- Inline small, fragmented rules into large rules
- Benefits:
 - Reducing recursion overhead
 - Providing more rule-local information for the local token mask cache
 - Fragmented rules lack rule-local information



Optimizations #4. Optimization Passes (Path merging)



- Merges paths with common prefix (when safe)
- Reduces number of possible parsing stacks
 - Need to be handled one by one when generating masks
 - In the recursive cases, #(parsing stacks) may explode exponentially!



XGrammar: Flexible and Efficient Structured Generation Engine

XGrammar is a structured generation library that features



Flexibility: Full support for context-free grammar



Efficiency: SOTA performance in constraint decoding

Zero-overhead JSON Schema generation

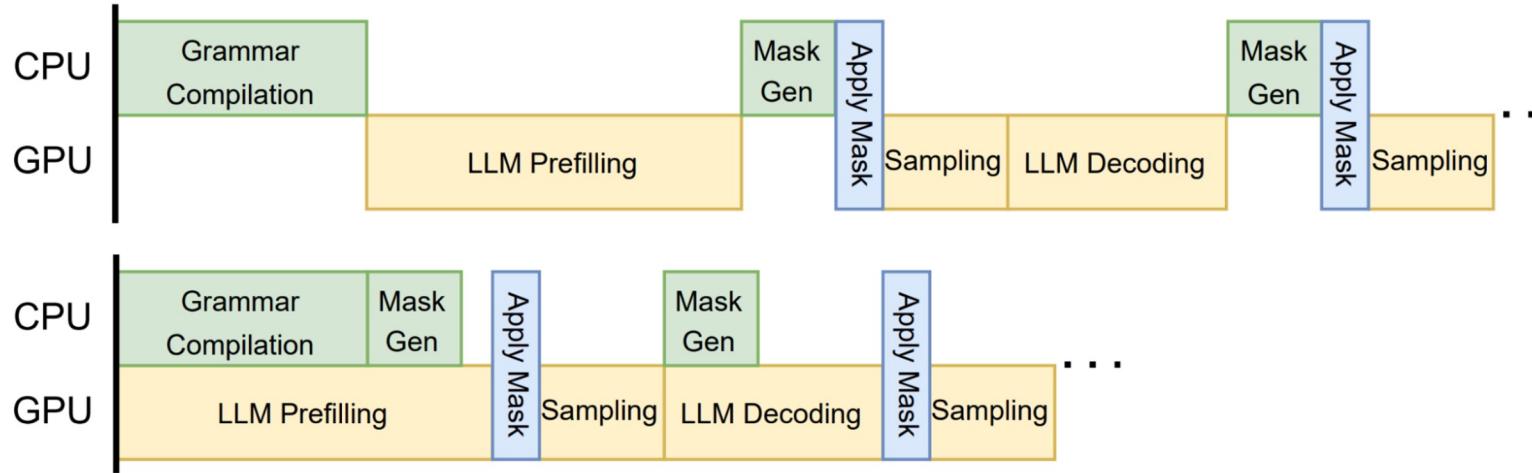


Integration: Easy to integrate with existing LLM serving frameworks

vLLM, MLC-LLM, SGLang, VILA, etc



Overlapping Mask Generation and LLM Inference



- Top: constrained decoding pipeline without overlapping
- Bottom: constrained decoding pipeline with overlapping



Integration with LLM Serving Frameworks

- XGrammar is designed for easy integration and cross-platform support (with C++, Python, and JavaScript APIs)
 - Its core is implemented in C++, so easy to port to other platforms
- XGrammar has already been integrated with vLLM, SGLang, MLC-LLM, WebLLM, VILA



Multimodal Generation w/ VILA

XGrammar has been integrated into VILA to enable structural generation with multimodal input



Question

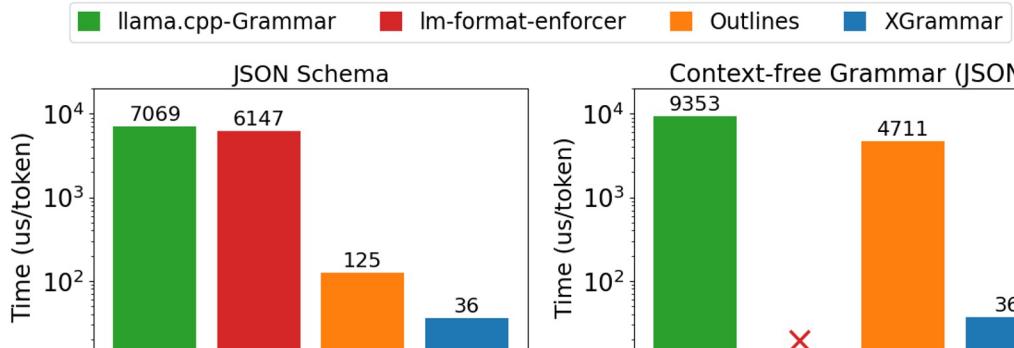
Schema



```
{  
  "description": "The image ...",  
  "objects": [  
    "mountains made of meat",  
    "valleys and rivers made of sliced ham",  
    ...  
  ],  
  "setting": {  
    "location": "a surreal landscape made of meat",  
    "time_of_day": "daytime",  
    "lighting": "bright and sunny"  
  },  
  "colors": [  
    "pink and red for the mountains",  
    ...  
  ]  
}
```

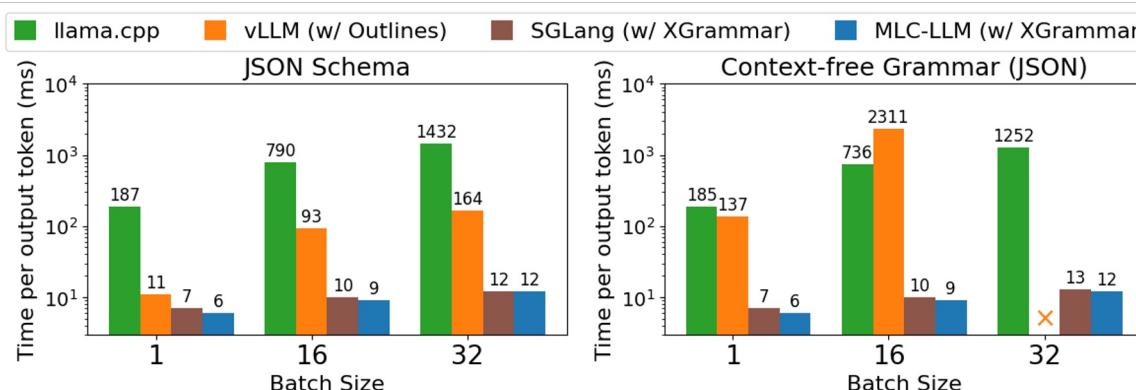


Evaluation



Overhead of masking logits.
(Llama-3-8B, AMD 7950X CPU,
RTX 4090)

Up to 3.5x on JSON schema
Up to 10x on CFG-guided



Time per output token for end-to-end LLM inference. (Llama-3-8B,
AMD 7950X CPU, H100 GPU)

Up to 14x in JSON-schema
Up to 80x in CFG-guided





MACHINE LEARNING
COMPILATION

Thanks

Questions are welcome!

- **Paper:** <https://arxiv.org/abs/2411.15100>
- **Blog:** <https://blog.mlc.ai/2024/11/22/achieving-efficient-flexible-portable-structured-generation-with-xgrammar>
- **Code:** <https://github.com/mlc-ai/xgrammar>
- **Documentation:** <https://xgrammar.mlc.ai/docs/>
- **Email:** yixind@andrew.cmu.edu