

### **Structured Generation and the New Trend**

Yixin Dong

Carnegie Mellon University

Nov 4, 2025

**01** Introduction

Context

02 XGrammar v1

03 Structural Tag

**04** Xgrammar v2

# Part 01 Introduction

Large Language Model Generation with Structures

### **Problem: LLM Generation with Structures**

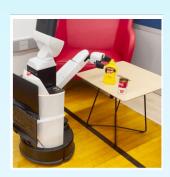
Code generation

Function/tool calling

**Embodied Agents** 

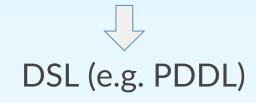






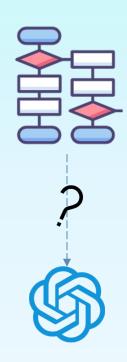






**Structured Outputs** 

### **Problem: LLM's Limited Ability with Complex Structures**



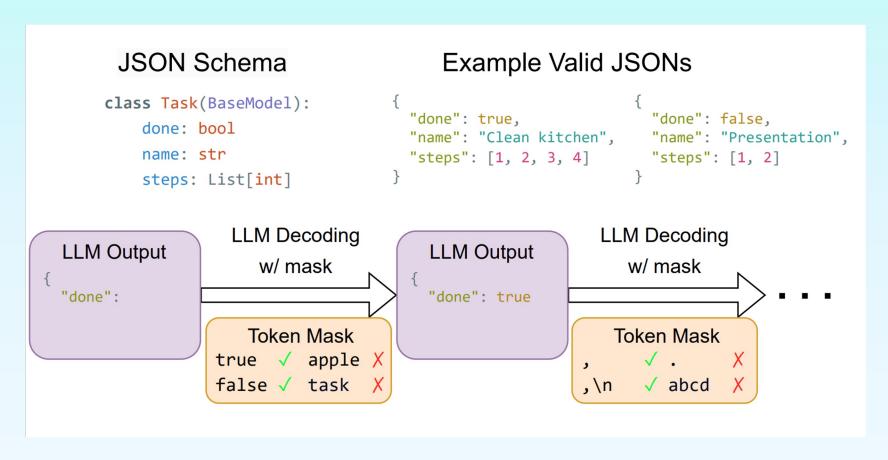
Structures are increasingly complex

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

LLM's generation ability is limited

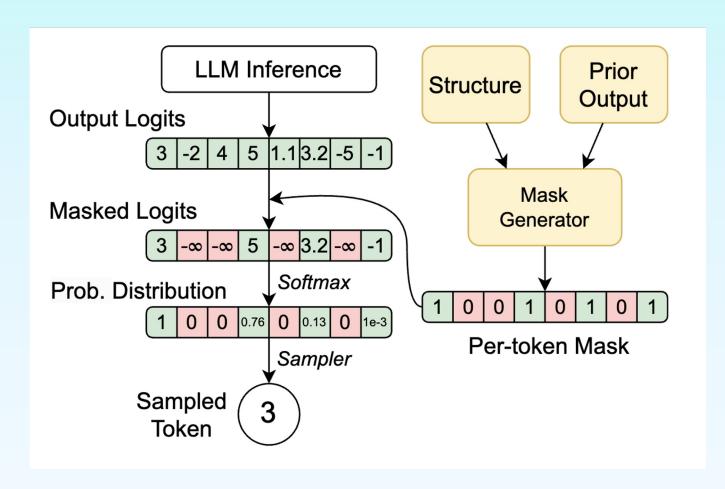
- On-device small LMs
- Compressed models

### **Background: Constrained Decoding**



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!

### **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

### **XGrammar Open-Source Project**

XGrammar has been adopted by these LLM serving engines:









And industrial collaborators:















# Part 02 XGrammar v1

Flexible And Efficient Structured Generation Engine

## 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> ']'

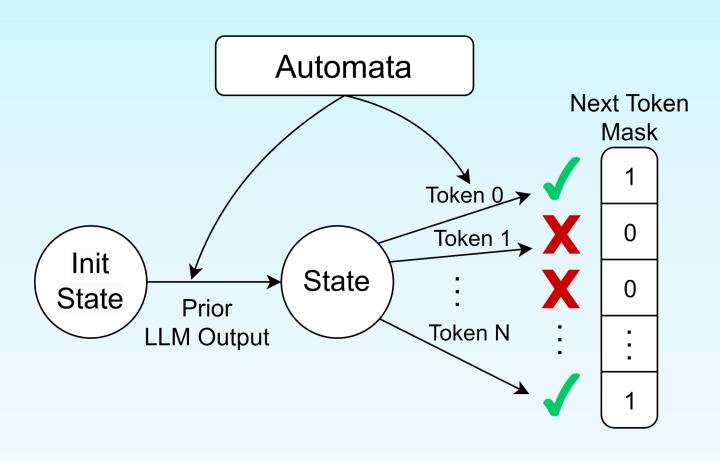
str ::= '"' [^"\]* '"'
```

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)

### **Previous Mask Generation Method**

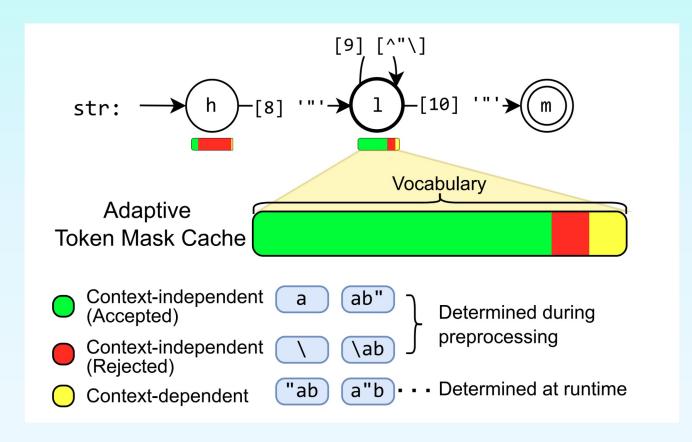


Check if each token matches, then generate the mask

Every token checking is very slow!

Our optimization: most token checking can be fast via preprocessing

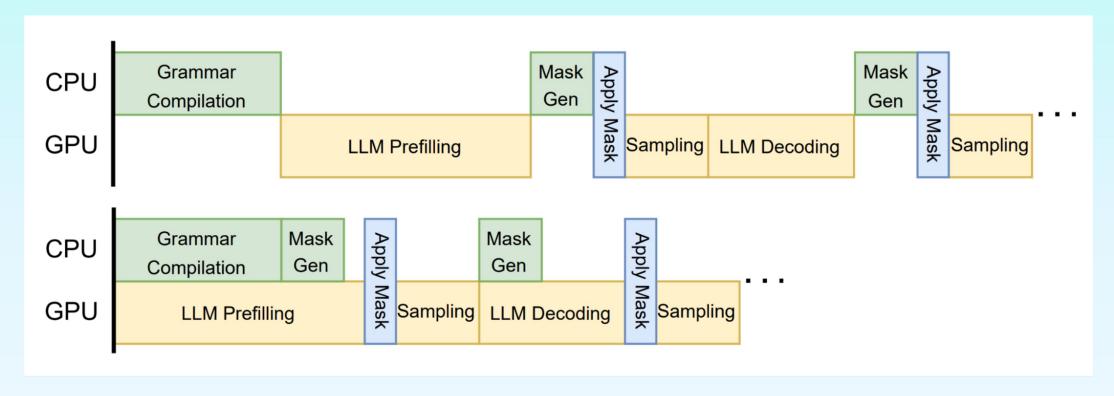
### The Adaptive Token Mask Cache (Cont'd)



- 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

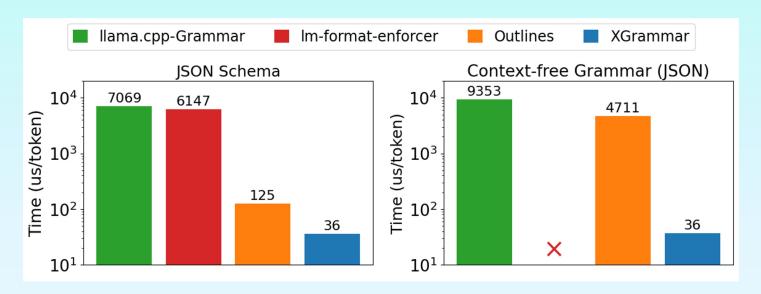
Context-dependent tokens: less than 1% for Llama-3.1 w/ JSON grammar (1134 out of 128k)

### **Overlapping Mask Generation and LLM Inference**



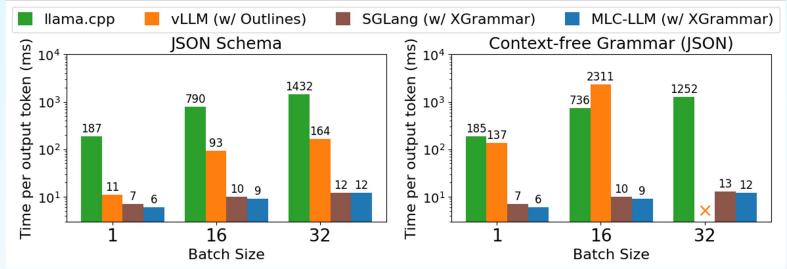
- Top: constrained decoding pipeline without overlapping
- Bottom: constrained decoding pipeline with overlapping

### **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 endto-end LLM inference. (Llama-3-8B, AMD 7950X CPU, H100 GPU)

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

### 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",
```

# Part 03 Structural Tags

Flexible and Dynamic Structural Generation API

## The Structural Tag

```
None
{"name": "function_name", "parameters": params}
```

- If LLM outputs {"name": "get weather"
- → We follow the schema of get\_weather
- If LLM outputs {"name": "find\_address"
- → We follow the schema of find\_address

Do dispatches between grammars

### The Structural Tag

- The Structural Tag is a JSON DSL
- The Structural Tag provides a convenient way to describe output structure
- E.g. Tool Calling

## **Function-Calling**

```
structural tag = TriggeredTagsFormat(
    triggers=["<function="],</pre>
    tags=[
        TagFormat(
            begin="<function=func1>",
            content=JSONSchemaFormat(json_schema=...),
            end="</function>",
        TagFormat(
            begin="<function=func2>",
            content=JSONSchemaFormat(json_schema=...),
            end="</function>",
    ],
    at_least_one=False,
    stop_after_first=False,
```

Allow any output until a trigger is encountered, then dispatch to the corresponding tag.

When the end tag is encountered, the grammar will allow any following output, until the next trigger is encountered.

## **Types of Structural Tags**

- Basic Formats:
  - ConstStringFormat
  - JSONSchemaFormat
  - AnyTextFormat
  - GrammarFormat
  - RegexFormat
  - QwenXMLParameterFormat

- Combinatorial Formats:
  - SequenceFormat
  - OrFormat
  - TagFormat
  - TriggeredTagsFormat
  - TagsWithSeparatorFormat

## **Force Thinking**

```
"type": "structural tag",
"format": {
    "type": "sequence",
   "elements": [
            "type": "tag",
            "begin": "<think>",
            "content": {"type": "any_text"},
            "end": "</think>",
            "type": "triggered tags",
            "triggers": ["<function="],
            "tags": [
                    "begin": "<function=func1>",
                    "content": {"type": "json_schema", "json_schema": ...},
                    "end": "</function>",
                    "begin": "<function=func2>",
                    "content": {"type": "json_schema", "json_schema": ...},
                    "end": "</function>",
```

The output should start with a reasoning part (<think>...</think>), then can generate a mix of text and tool calls.

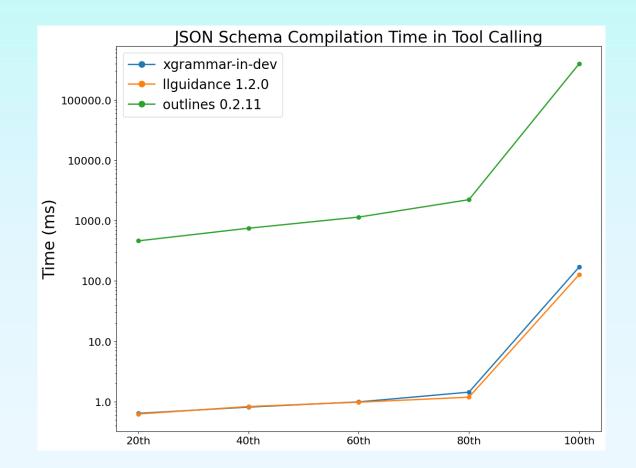
### Force non-thinking mode

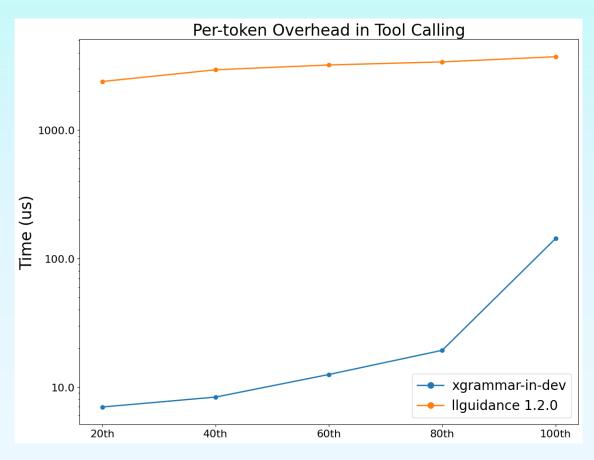
```
"type": "structural_tag",
"format": {
    "type": "sequence",
    "elements": [
            "type": "const string",
            "text": "<think></think>"
            "type": "triggered_tags",
            "triggers": ["<tool_call>"],
            "tags":
                    "begin": "<tool_call>\n{\"name\": \"func1\", \"arguments\": ",
                    "content": {"type": "json_schema", "json_schema": ...},
                    "end": "}\n</tool_call>",
                    "begin": "<tool call>\n{\"name\": \"func2\", \"arguments\": ",
                    "content": {"type": "json_schema", "json_schema": ...},
                    "end": "}\n</tool call>",
```

Qwen-3 has a hybrid thinking mode that allows switching between thinking and non-thinking mode. Thinking mode is the same as above, while in non-thinking mode, the output would start with a empty thinking part <think></think>, and then can generate any text.

# Part 04 XGrammar v2

Zero Overhead Function Calling





## **Dispatching**

### **AC Automaton**

Pattern[0]: "<function="

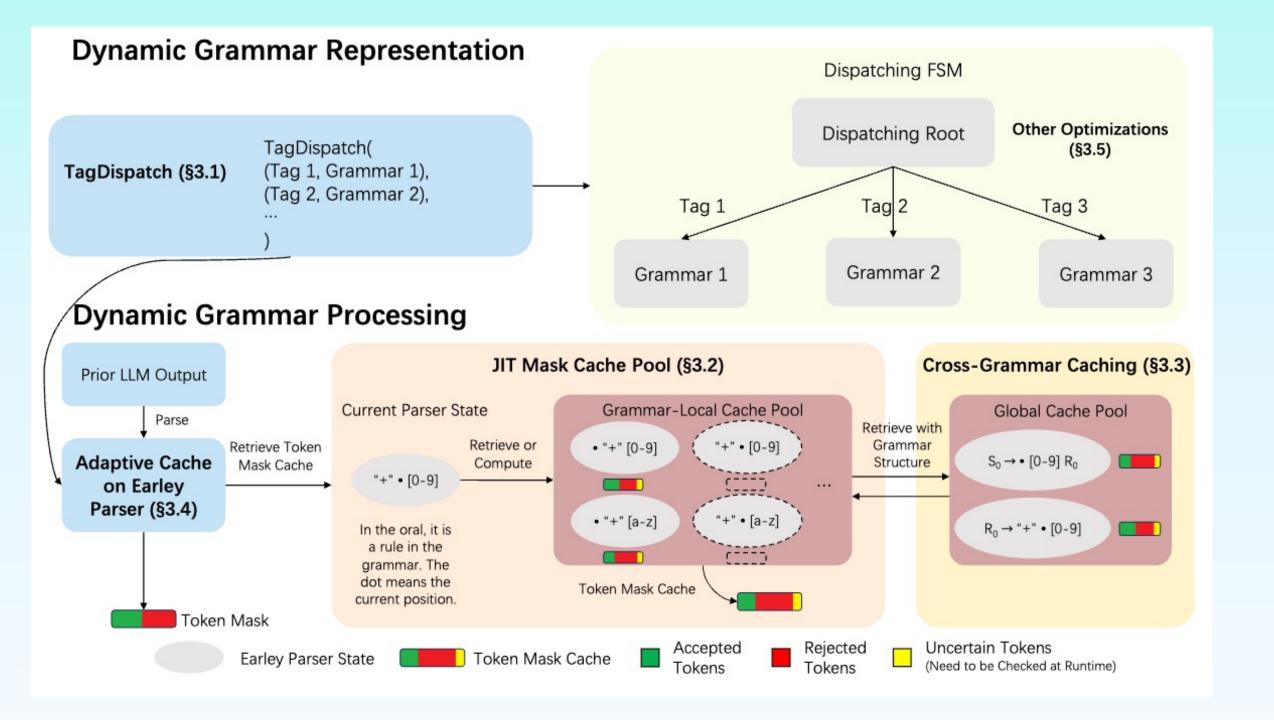
Pattern[1]: "<think>"

Pattern[2]: "<call>"

. . .

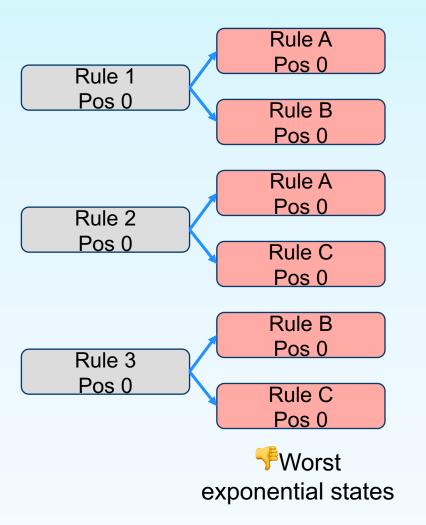
# **Dispatched** Grammar[0] Pattern[0] Pattern[1] Grammar[1] Pattern[2] Grammar[2]

**Dispatched Grammar Completed** 

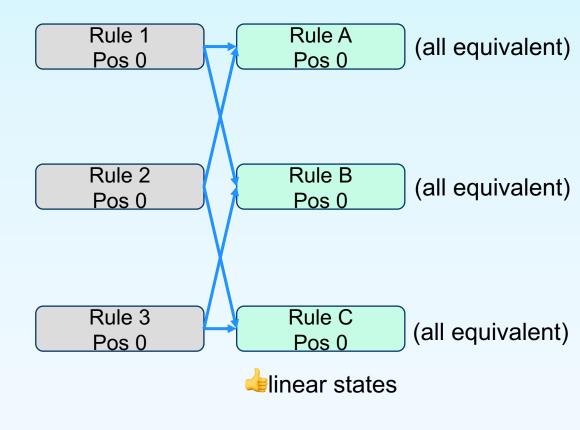


## **Earley Parser**

### **PDA Parser**



## **Earley Parser**



# Thanks.

Repo: <a href="https://github.com/mlc-ai/xgrammar">https://github.com/mlc-ai/xgrammar</a>
Papers: <a href="https://arxiv.org/abs/2411.15100">https://arxiv.org/abs/2411.15100</a>

Blog: https://blog.mlc.ai/2024/11/22/achieving-efficient-flexible-portable-

structured-generation-with-xgrammar

Documentation: <a href="https://xgrammar.mlc.ai/docs/">https://xgrammar.mlc.ai/docs/</a>