

Carnegie
Mellon
University



MACHINE LEARNING
COMPI LATION

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

Yixin Dong

Carnegie Mellon University

Oct 16, 2024

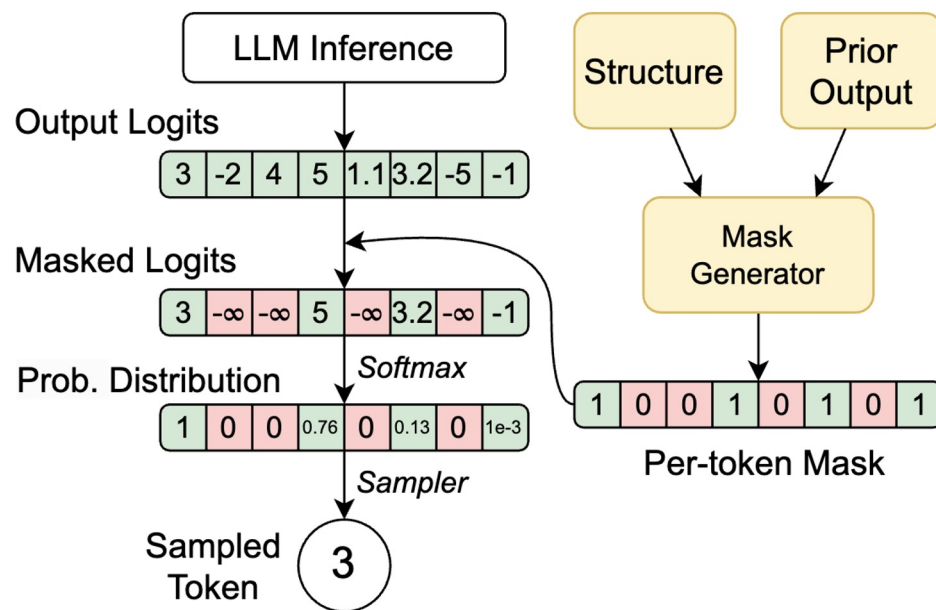
Background: Structured Generation for LLMs



* Some programming languages like Python can be described as CFG plus maintaining a state



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



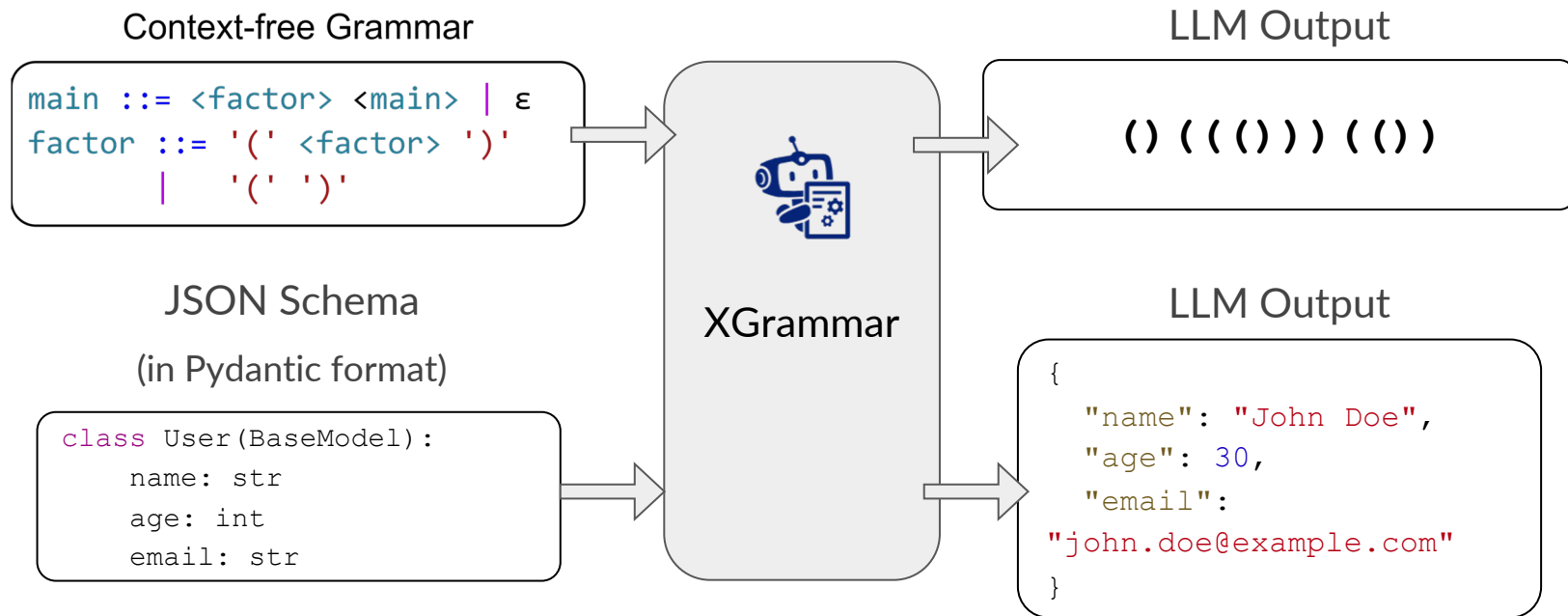
Efficiency: SOTA performance in constraint decoding
Zero-overhead JSON Schema generation



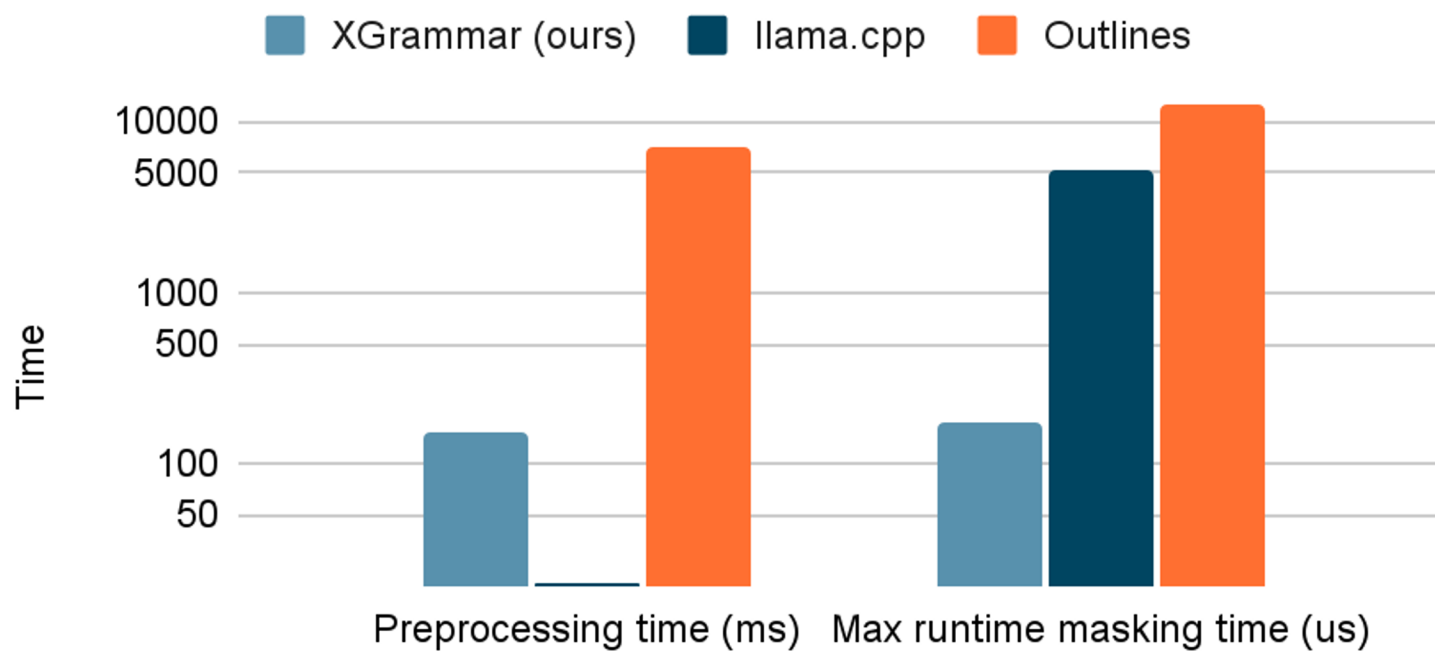
Integration: Easy to integrate with existing LLM serving frameworks
MLC-LLM, SGLang, Web-LLM, etc



Demonstration: CFG or JSON Schema with XGrammar

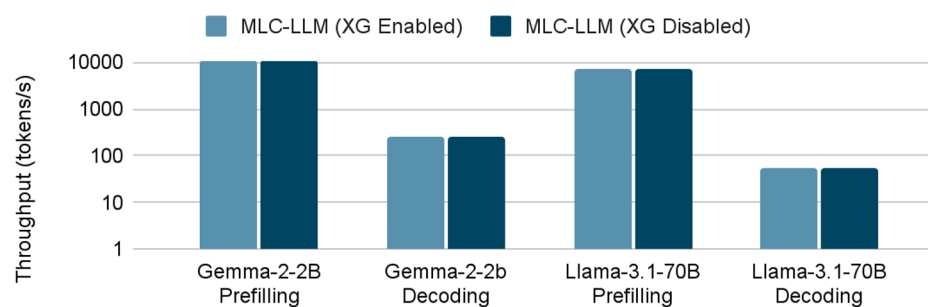


Unit benchmarks



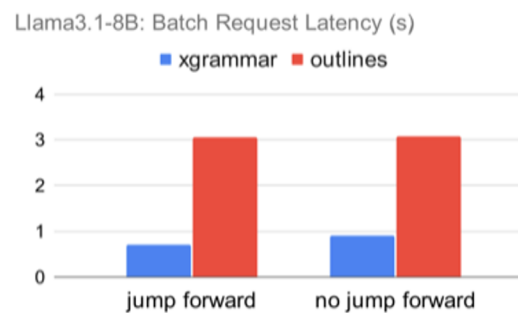
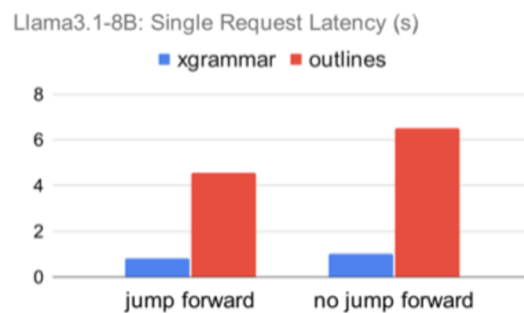
End-to-end benchmarks

MLC-LLM E2E demonstrates zero-overhead constrained decoding



(H100, batch=1)

SGLang E2E demonstrates the superiority of XGrammar over other libraries

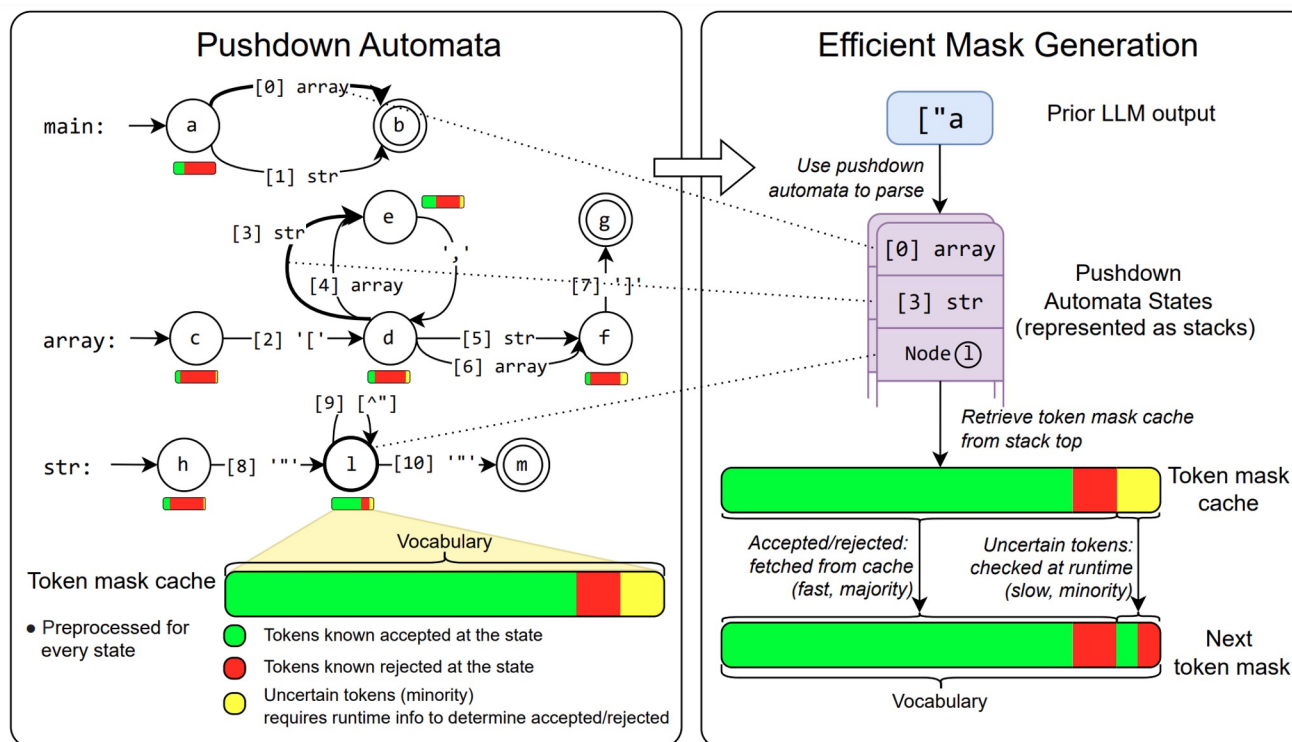


(A100)



Key optimization: token mask cache

Most part of the mask can be precomputed, although a few remaining ones need to be detected at runtime



Usage: BNFGrammar

```
bnf_grammar := BNFGrammar(  
    ... ""main ::= factor main | ""  
    ... factor ::= "(" factor ")" | "()",  
    ... main_rule="main",  
)
```

Provide the grammar specification in the extended Backus-Naur form (EBNF)

And specify the name of the main rule ("main" by default)



Usage: BuiltinGrammar

```
BuiltinGrammar.json()  
BuiltinGrammar.json_schema(schema: ·str· | ·type[Pydantic.BaseModel])  
BuiltinGrammar.str(s: ·str·)  
BuiltinGrammar.multiple_choices(choices: ·List[str])  
BuiltinGrammar.concat(grammars: ·List[BNFGrammar])  
BuiltinGrammar.regex(pattern: ·str·)
```

XGrammar provides multiple helper methods to conveniently generate grammar



Usage: GrammarMatcher

```
matcher = GrammarMatcher(bnf_grammar, TokenizerInfo(hf_tokenizer))
while True:
    logits = model.inference(input_tokens)
    bitmask: torch.Tensor = matcher.find_next_token_bitmask()
    GrammarMatcher.apply_bitmask_inplace(logits, bitmask)
    token_id = sampler.sample(logits)
    matcher.accept_token(token_id)
```

At every step, XGrammar provides a token mask to apply to logits, masking invalid tokens. Each generated token updates the matcher state.



Usage: Caching the preprocessing result

```
init_context := GrammarMatcherInitContext(grammar, tokenizer_info) ·#·150·ms

matcher := GrammarMatcher(init_context) ·#·20us
bitmask := matcher.find_next_token_bitmask() ·#·30us

-----

init_context_cache := GrammarMatcherInitContextCache(tokenizer_info) ·#·20us

init_context_1 := init_context_cache.get_init_ctx_for_grammar(grammar) ·#·150ms
init_context_2 := init_context_cache.get_init_ctx_for_grammar(grammar) ·#·10us
```

The preprocessing can be slow (100-200 ms)

GrammarMatcherInitContext:

The preprocessing result, can be used to construct GrammarMatcher fast

GrammarMatcherInitContextCache:

Quickly retrieve the init context of the same grammar multiple times.



Integration with LLM serving frameworks



- XGrammar is designed for easy integration and cross-platform support (with C++, Python, and TypeScript APIs)
- XGrammar has already been integrated with SGLang, MLC-LLM and WebLLM and is currently being integrated with deepseek
- XGrammar will be open-sourced and released later this month



Try it out on Web-LLM



Run structured generation completely on your web browser with great efficiency!

Model
Llama-3.1-8B-Instruct-q4f16_1-MLC

Grammar
Custom Grammar

Custom EBNF Grammar
The custom grammar is described in the [extended Backus-Naur form \(EBNF\)](#). Below is an example of JSON grammar in EBNF. Please follow this example when writing new grammars.

```
main    ::= "1. " move " " move "\n" ([1-9] [0-9]? ". " move " " move "\n")+
move    ::= (pawn | nonpawn | castle) [#]?
# piece type, optional file/rank, optional capture, dest file & rank
nonpawn ::= [NBKQR] [a-h]? [1-8]? "x"? [a-h] [1-8]
# optional file & capture, dest file & rank, optional promotion
pawn    ::= ([a-h] "x")? [a-h] [1-8] ("=" [NBKQR])?
castle  ::= "0-0" "-0"?
```

Prompt

```
main    ::= "1. " move " " move "\n" ([1-9] [0-9]? ". " move " " move "\n")+
move    ::= (pawn | nonpawn | castle) [#]?
# piece type, optional file/rank, optional capture, dest file & rank
nonpawn ::= [NBKQR] [a-h]? [1-8]? "x"? [a-h] [1-8]
# optional file & capture, dest file & rank, optional promotion
pawn    ::= ([a-h] "x")? [a-h] [1-8] ("=" [NBKQR])?
castle  ::= "0-0" "-0"?
```

Generate

Output

```
1. e4 e5 2. Nf3 Nc6 3. Bc4 Nf6 4. d3 d6 5. O-O O-O 6. b4 a5 7. b5 Nxb5 8. a4 Na6 9.
Qe2 Qe7 10. a5 Nc5 11. a6 b6 12. a7 b7 13. a8=Q a8=Q 14. Nb1 Nb1 15. Nc3 Nc3 16. Na4
Na4 17. b5 b5 18.
```

Prefill Speed: 39.4 tok/s, Decode Speed: 19.8 tok/s, Time to First Token: 50 ms, Time Per Output Token: 4698 ms, Grammar Per-token Overhead: 0.07 ms



With help from Charlie Ruan and Nestor Qin from CMU Catalyst Lab
XGrammar: Flexible And Efficient Structured Generation Engine for Large Language Models

Carnegie
Mellon
University



MACHINE LEARNING
COMPILATION

Thanks

Questions are welcome!

Yixin Dong
Oct 16, 2024