

An Incremental Valiant Recognizer for GPU-accelerated General Context-free Prefix Validation

Anonymous ACL submission

Abstract

Constrained decoding systems are often built on context-free parsers intended for programming languages. These parsers either degrade into $O(n^3)$ time-complexity or fail entirely if the grammar is not carefully engineered to keep properties such as determinism and non-ambiguity. There is thus a need to design parsers that efficiently handle non-determinism and ambiguity, while simultaneously being incremental so that they can be coupled with the token-based predictions of large language models. Inspired by prior work, we derive an incremental Valiant prefix recognizer, which still has $O(n^3)$ complexity but allows for acceleration with only a fraction of GPU resources (stream multiprocessors). Our parser shows robust efficiency in complex context-free grammars while other parsers crash or degrade. At the same time, we remain empirically competitive in restrictive grammar classes such as LALR.

1 Introduction

Constrained decoding systems typically face a trade-off when selecting a context-free parser. One option is to adopt parsers from the LL or LR families, which consistently achieve $O(n)$ complexity but require users to invest substantial engineering effort in crafting restrictive grammars. Alternatively, systems may employ more general parsers such as Earley, which can achieve $O(n)$ complexity when the grammar happens to be unambiguous and deterministic, but degrade to $O(n^3)$ otherwise.

This limitation poses fewer problems in the traditional use case of parsing programming languages, as such languages are typically designed with these parsing restrictions in mind and have update cycles spanning at least a year, allowing ample time for careful grammar engineering. However, LLM systems iterate at a much faster pace, and researchers have increasingly sought to constrain diverse outputs including reasoning traces, molecular

structures, and privacy-sensitive contents (Banerjee et al., 2025; Loula et al., 2025; Ugare et al., 2025). This calls for stronger general-case performance.

The Valiant recognizer (Valiant, 1975) achieves sub-cubic general context-free recognition by connecting CYK parsers with Boolean matrix multiplication. Although a practical sub-cubic algorithm has yet to emerge, GPUs can provide strong acceleration. Azimov and Grigorev (2018) demonstrates a GPU-accelerated Valiant recognizer for graph database query that outperforms GLL, a generalized top-down parsing algorithm.

However, the Valiant recognizer cannot be applied to constrained decoding, as it lacks the valid-prefix property—the ability to parse incrementally and reject invalid continuations as soon as the first violating element is encountered. To address this limitation, our technical contributions include:

- Deriving a Valiant valid-prefix recognizer.
- Demonstrating that incremental parsing substantially reduces memory footprint and improves CUDA L1/SMEM access patterns.
- Implementing a constrained decoder that empirically performs competitively in restricted grammars and outperforms prior work significantly in complex grammars.

2 Related Work

There has been multiple work on improving the efficiency of constrained decoders. A large fraction of these works focuses on the alignment between LLM tokens and grammar terminals (Dong et al., 2025; Moskal et al., 2025; Ugare et al., 2024; Beurer-Kellner et al., 2024; Anonymous, 2025). Work on improving the internals of existing parsers can also be found (Sun et al., 2025).

General context-free parsing can be achieved by extending the traditional LL and LR parsers (Scott and Johnstone, 2010; Tomita, 1985). Alternatively,

chart-based parsers such as CYK (Cocke, 1969; Younger, 1967; Kasami, 1965) and Earley (Earley, 1970) can be used. Unlike LL and LR parsers, general context-free parsers can be extended to mildly-context-sensitive parsers, which is relevant to efforts in the constrained decoding community to cover more complex constraints (Ugare et al., 2025; Nagy et al., 2025).

Parallelization of chart-based parsers is inspired by the connection between the CYK parser and matrix multiplication (Valiant, 1975). Bernardy and Claessen (2015) shows that a sparsely encoded chart can reduce complexity from $O(n^3)$ to $O(\log^2 n)$ in practical cases. A practical implementation on GPUs is used to find context-free paths in a graph database (Azimov and Grigorev, 2018).

3 Background

3.1 Constrained Decoding

Large language models produce a probability distribution over their vocabulary to express preferences for different output tokens across different time steps. Constrained decoders seek to modify this distribution to make the eventual output text conform to a specification, often expressed by a formal grammar, such as a context-free grammar.

3.2 Context-free Grammar

Given a finite set alphabet Σ , a context-free grammar defines a language $L \subseteq \Sigma^*$ using a quadruple $G = (N, \Sigma, P, S)$ where:

- Σ is a finite set of terminal symbols, which we will write in lower case,
- N is a finite set of nonterminal (variable) symbols, which we will write in upper case,
- $S \in N$ is the start symbol, and
- $P \subseteq N \times (N \cup \Sigma)^*$ is a finite set of production rules, written $A \rightarrow \alpha$ with $A \in N$ and $\alpha \in (N \cup \Sigma)^*$.

We use $A \Rightarrow^* \alpha$ to denote that symbol sequence α is derivable from non-terminal A by repeated application of production rules.

3.3 Chomsky Normal Form

Each context-free grammar can be converted into equivalent Chomsky Normal Form (CNF; Chomsky, 1959), where the right side of each production rule either has two non-terminals or one terminal,

Algorithm 1 CYK Algorithmn, row-by-row

Input: Terminal sequence $t_0 \dots t_{n-1}$

- 1: **Base case:** row = 0
- 2: $TABLE[row][col][V] \Leftrightarrow$
- 3: \exists rule $V \rightarrow t_{col}$ in CNF grammar
- 4: **Recursion:** row = $1 \dots n - 1$
- 5: $TABLE[row][col][V] \Leftrightarrow$
- 6: \exists rule $V \rightarrow AB$ in CNF grammar, s.t.
- 7: \exists mid $\in \mathbb{N}$, s.t. mid < row and
- 8: $TABLE[mid][col][A]$ and
- 9: $TABLE[row - mid - 1]$
- 10: $[col + mid + 1][B]$

and only the start symbol is allowed to generate the empty string ϵ . This conversion is also called binarization, an essential step for CYK-like parsers.

4 The Incremental Valiant Prefix Recognizer

We describe the well known CYK/Valiant recognizer for context free grammars in section 4.1, modify it to run incrementally in section 4.2, extend it to recognize valid prefixes in section 4.3, and further improve performance in section 4.4.

4.1 Canonical CYK/Valiant

The CYK parser is a bottom-up, dynamic programming parser, as defined in algorithm 1. It holds the invariant that V is at row and col (i.e., $TABLE[row][col][V] = true$) if and only if $V \Rightarrow^* t_{col} \dots t_{col+row}$.

The Valiant Recognizer re-frames the same algorithm as boolean matrix operations. We use the two names interchangeably.

4.2 Incremental CYK/Valiant

Algorithm 1 works only if the entire terminal sequence is observed. However, in constrained decoding, we are given a parser state for prefix $t_0 \dots t_{n-1}$ and a new terminal t_n , and need to compute parser state for $t_0 \dots t_{n-1}t_n$ quickly.

Bernardy and Claessen (2015) proves that a CYK parser can be executed incrementally. To implement such an incremental parser, we build on the observation that $TABLE[row][col]$ depends on $TABLE[r][c]$ if and only if $r + c \leq row + col$ and $c \geq col$. In other words, a cell depends on cells within the triangle tipped by itself, as illustrated in figure 1.

Taking advantage of this property, instead of addressing cells using $[row, col]$, we address cells as

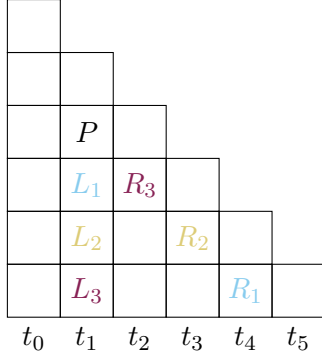
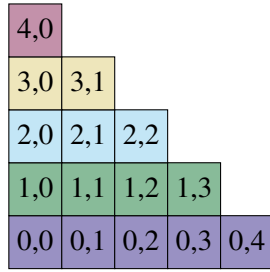
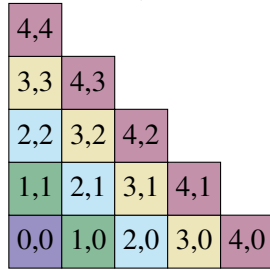


Figure 1: All potential children of the parent cell P at $(3, 1)$ fall in a triangle below it: the cells (r, c) where $r + c \leq 3 + 1$ and $c \geq 1$. All potential children follow this pattern: L_1 at $(2, 1)$, R_1 at $(0, 4)$, L_2 at $(1, 1)$, R_2 at $(1, 3)$, L_3 at $(0, 1)$, and R_3 at $(2, 2)$.



Canonical (by row)



Incremental (by tier)

Figure 2: Conversion from canonical (row, col) addressing system to incremental $(tier, row)$ addressing system, where $tier = row + col$. Cells processed in the same loop are in the same color. The canonical approach processes row by row, while the incremental approach processes tier by tier.

$[tier, row]$ as shown in figure 2, and re-arrange the canonical CYK (algorithm 1) into an incremental one (algorithm 2). Notice that each incremental step adds one tier, and each tier is never written to after calculated. This enables the zero-copy branching or rollback property typically desired by constrained decoders.

Algorithm 2 Incremental Valiant Step

Input: $TIER[0 \dots n-1]$ for terminal $t_0 \dots t_{n-1}$.

Input: new terminal t_n

- 1: $TIER[n][0][V] \Leftrightarrow \exists$ CNF rule $V \rightarrow t_i$
 - 2: **for** row in **range**(1, n+1):
 - 3: $TIER[n][row][V] \Leftrightarrow$
 - 4: \exists CNF rule $V \rightarrow AB$ s.t.
 - 5: \exists mid < row, s.t.
 - 6: $TIER[n-row+mid][mid][A]$ and
 - 7: $TIER[n][n-mid-1][B]$
-

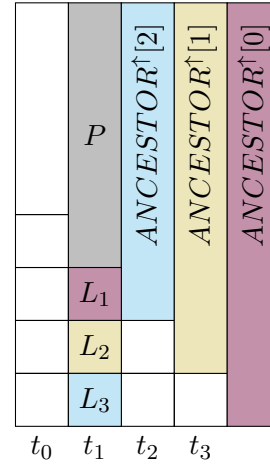


Figure 3: Visually, $ANCESTOR^{\uparrow}$ is the union of all possible future cells in the same column. For each cell to act as left child, any paired right cell will always come from the same column. For example, if a future parent P 's left child comes from L_2 , its right child must come from $ANCESTOR^{\uparrow}[1]$. Note also that any future parent that has a concrete cell as its left child must have a future cell as its right child, as per CYK indexing (see figure 1).

4.3 Valid Prefix Recognition

Unlike the often used LR/LL/Earley parsers, algorithm 2 will process invalid prefixes (i.e., prefixes that can never be generated from the start symbol) without complaining. This makes it unusable in a constrained decoder where such prefixes should be excluded. So we need to further modify it to compute a new tier type named $ANCESTOR^{\uparrow}[j][V]$, which we visualize in figure 3 and formalize as algorithm 3.

Theorem 1. $ANCESTOR^{\uparrow}[j][V]$ if and only if terminal sequence $t_{n-j} \dots t_{n-1}$ is a valid prefix of non-terminal V

See appendix A for the formal proof. Intuitively and visually, $ANCESTOR^{\uparrow}$ is the union of all possible future cells in the same column in the CYK

Algorithm 3 $ANCESTOR^\uparrow$ algorithm

Input: $TIER[0 \dots n-1]$ for terminal $t_0 \dots t_{n-1}$ **Input:** $FIRST$ = FIRST set of CNF grammar**Input:** hypothetical terminal \hat{t}_n

```
1:  $ANCESTOR^\uparrow[0] = FIRST^{-1}[\hat{t}_n]$ 
2: FOR row in RANGE(1, n+1):
3:   # Case 1: Left child is a concrete cell
4:    $ANCESTOR^\uparrow[row][V] = True$  if
5:      $\exists$  CNF rule  $V \rightarrow AB$  s.t.
6:        $\exists$  mid < row, s.t.
7:          $TIER[n-row+mid][mid][A]$  and
8:          $ANCESTOR^\uparrow[n-mid-1][B]$ 
9:   # Case 2: Left child is an ancestor cell
10:   $ANCESTOR^\uparrow[row][W] = True$  if
11:     $\exists V$  s.t.  $ANCESTOR^\uparrow[row][V] \wedge$ 
12:       $FIRST[V][W]$ 
```

Algorithm 4 $ANCESTOR^\downarrow$ algorithm

Input: $TIER[0 \dots n-1]$ for terminal $t_0 \dots t_{n-1}$ **Input:** $FIRST$ = FIRST set of CNF grammar

```
1:  $START$  non-terminal  $\in ANCESTOR^\uparrow[0]$ 
2: FOR row in RANGE(n, -1, -1):
3:   # Case 1: parent in a future cell
4:    $ANCESTOR^\uparrow[row][V] = True$  if
5:      $\exists V$  s.t.
6:        $ANCESTOR^\uparrow[row][W]$  and
7:        $FIRST[V][W]$ 
8:   # Case 2: left child in a concrete cell
9:    $ANCESTOR^\downarrow[row][V] = True$  if
10:     $\exists$  CNF rule  $V \rightarrow AB$  s.t.
11:       $\exists$  mid < n - row, s.t.
12:         $TIER[n-row-1][mid][A]$  and
13:         $ANCESTOR^\downarrow[row+mid+1][V]$ 
```

183 table. For a future parent to exist, CYK dictates
184 that it must have a left child in the same column as
185 itself, either from an existing cell or a future one.
186 If the left child is from an existing cell, the right
187 child must be from a future cell, or it would not
188 cover enough terminals. If the left child is from
189 a future cell, the right child must be from a col-
190 umn to the left of n , which is unrestricted by the
191 existing prefix. We only need to prove the exist-
192 ence of the left child's membership in the current
193 $ANCESTOR^\uparrow$ column. To check if \hat{t}_n is a viable
194 next terminal, we simply check if the start symbol
195 is in $ANCESTOR^\uparrow[n]$.

196 **Definition 1** (FIRST set). *Given a context-free*
197 *grammar G , non-terminal V and W ,*

198 $FIRST[V][W]$ and $FIRST^{-1}[W][V]$ if and
199 only if $V \Rightarrow^* W \dots$

200 The algorithm for computing $FIRST$ is
201 well established, and is equivalent to comput-
202 ing the transitive closure of the relationship
203 $\{ \langle V, A \rangle \mid \text{rule } V \rightarrow A \dots \in G \}$

204 4.4 Beyond Valid Prefix Property

205 When relying on the valid prefix property, a con-
206 strained decoder would have to run the parser for
207 each individual hypothetical terminal. Here we
208 show that, with slight modification, algorithm 3
209 can instead parse top-down and find all possi-
210 ble next terminals within one run. We compute
211 $ANCESTOR^\downarrow$ using algorithm 4.

212 **Theorem 2.** $ANCESTOR^\downarrow[j][V]$ if and only if
213 $t_0 \dots t_{n+j-1}$ stays a valid prefix after appended
214 with a string derivable by V

The proof of correctness of $ANCESTOR^\downarrow$ is
symmetrical to that of $ANCESTOR^\uparrow$, and thus
omitted. The set of all viable next terminals is
 $\cup_{V \in ANCESTOR^\downarrow} FIRST[V]$.

219 5 Implementation on GPU

220 5.1 Parallelization

221 Algorithm 3 and algorithm 4 each contain two syn-
222 chronization points: one after each case within the
223 loop body. Parallelization therefore occurs within
224 each case.

225 For the case where the left child is an an-
226 cestor cell, the FIRST set is encoded as ances-
227 tor/descendant pairs, allowing straightforward par-
228 allelization along the pairs dimension. For the case
229 where the left child is a concrete cell, the natu-
230 ral choice is to parallelize along the dimension of
231 CNF grammar rules. Additionally, there is an op-
232 portunity to parallelize across rows, as illustrated
233 in figure 4. We distribute different CNF grammar
234 rules across threads and employ hardware bit-wise
235 parallelization across rows.

236 5.2 Exploiting Chart Sparsity

237 Letting one thread handle all rows has the addi-
238 tional benefit of skipping computation when the
239 right child (in bottom-up parsing) or parent (in
240 top-down parsing) evaluates to false. This allows
241 partial exploitation of chart sparsity.

242 5.3 Memory Occupancy and Access Patterns

243 An unexpected benefit of incremental parsing is
244 that, once a tier is computed, it is subsequently

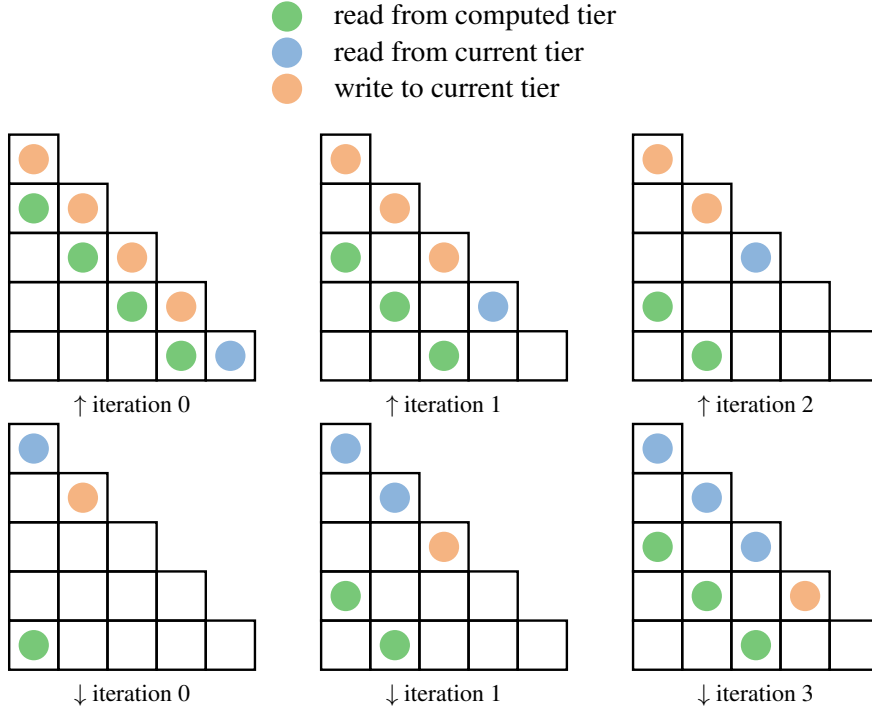


Figure 4: Access patterns. At each iteration, a different already computed tier is fetched as the left child. All the reads and writes across different rows can be accessed concurrently within iteration, which we exploit using hardware bit-wise parallelism.

accessed only as a left child. This limits the number of non-terminals that must be stored in the persistent parser state. We further reduce storage requirements by producing only right-leaning rules during binarization: a rule $X \rightarrow A B C$ is decomposed into $X \rightarrow A X'$ and $X' \rightarrow B C$. With this optimization, the number of stored non-terminals in our SQL grammar decreases from 3219 to 103.

CUDA partitions L1 cache and shared memory into 4 and 32 banks, respectively. Bank conflicts occur when threads within the same warp access different addresses in the same bank. Since non-terminals in computed tiers are accessed as left children and each thread processes a different grammar rule, we sort grammar rules by their left child non-terminal. Given that there are only 103 left-child non-terminals but over 3000 normalized grammar rules, threads within each warp of 32 will most likely access either the same address (triggering a broadcast) or neighboring addresses (avoiding conflicts).

5.4 Tensor Memory Accelerator

The NVIDIA Hopper architecture includes a specialized hardware component called the Tensor Memory Accelerator (TMA), which asynchronously fetches data from global memory, free-

ing CUDA cores for computation. We attempted to use TMA for loading left children—the dominant factor in the $O(n^2)$ per-step complexity. However, this resulted in a slowdown due to additional book-keeping overhead. We attribute this to the relatively short sequence lengths in our experiments and the effectiveness of the optimizations described above.

6 Alignment With Tokenizer

LLM tokens can span multiple grammar terminals. For example, the Llama-3 token “”);” could become the tail of *STRING*, followed by *LP* and *SEMI* in an SQL grammar. This necessitates searching multiple steps beyond the already committed terminals, and correct handling is critical to the speed of constrained decoders. Dedicated literature exists on this topic (Beurer-Kellner et al., 2024; Dong et al., 2025), and state-of-the-art approaches can be found in works such as LLGuidance (Moskal et al., 2025) and PSC (Anonymous, 2025).

Efficient alignment of grammar and tokenizer is out of scope for this paper. Instead, we further develop on the zero-copy property discussed in section 4.2 and derive a simple but GPU friendly tree search.

296	6.1 Parallel Terminal Tree Search	
297	Given a pair of a context-free grammar and an	
298	LLM tokenizer, the worst-case search tree is finite	346
299	due to the finite number of terminals and tokens.	347
300	This allows pre-computation and pre-allocation of	348
301	the search tree skeleton, enabling efficient GPU	
302	kernels.	
303	We encode the search tree as a flat array, where	
304	each element is one possible terminal sequence	
305	$t_{lead} \dots t_{tail}$ that can appear within a token. For	
306	example, with the full-sized official SQL grammar	
307	containing 167 terminals and Llama 3 tokenizer	
308	containing 128256 tokens, the reachability of all	
309	nodes in the search tree can be encoded as an array	
310	of 307692 booleans.	
311	6.1.1 Storage of Parser State	
312	At each node in the search tree, we make a	
313	branch of parser state from the parent node then	
314	commit a new terminal. With zero-copy, each	
315	node only requires one additional parser tier (i.e.	
316	$\mathbf{TIER}_{t_n \dots t_{n-1} t'_n \dots t'_m}$).	
317	CUDA requires all memory to be allocated be-	
318	fore kernel launch. Although tempting, allocating	
319	for every parser tier for every search tree node will	
320	exceed the memory limit on larger grammars. So	
321	we employ breadth first search (BFS) up to a given	
322	depth, and switch to depth first search (DFS) once	
323	the search tree is sufficiently developed. Thus, only	
324	BFS searched tiers require dedicated parser tier al-	
325	location. DFS searched nodes will use a per-worker	
326	stack of tiers.	
327	6.2 From Terminal Search Tree to Token	
328	Mask	
329	Since a terminal can consist of multiple tokens,	
330	it is conventional to track progress within termi-	
331	nals with deterministic finite state automata (DFA).	
332	Each token would be assigned multiple tuples	
333	$\langle q_{start}, t, q_{end} \rangle$ indicating that for terminal t the	
334	token would transit DFA from state q_{start} to q_{end} .	
335	We expand this notation to account	
336	for tokens that span multiple terminals,	
337	$\langle q^{lead}, \langle t_{lead}, \dots, t_{tail} \rangle, q^{tail} \rangle$ For exam-	
338	ple, the token (“”);” would take a <i>STRING</i> DFA	
339	in its penultimate state, complete it by consuming	
340	the first character (“”)), then complete a <i>LP</i> DFA	
341	(“”)), and finally stay on a <i>SEMI</i> DFA (“;”) in its	
342	final state.	
343	A token is allowed as a prediction from the large	
344	language model (i.e., a token is unmasked) if it has	
345	one transition tuple $\langle q^{lead}, \langle t_{lead}, \dots, t_{tail} \rangle, _ \rangle$	
	such that the DFA for t_{lead} is currently in state	346
	q^{lead} , and $t_{lead} \dots t_{tail}$ is marked as a viable termi-	347
	nal during the tree search introduced in section 6.1 .	348
	7 Experiments	349
	7.1 Data and Grammar	350
	Spider Text-to-SQL is a widely used text-to-	351
	SQL dataset. We test on a normal SQL syntax, and	352
	a spaced-operators (sp-op) syntax that forces space	353
	after symbols “+”, “-”, “*”, “/”, “,” and “;”’. Both	354
	grammars are LALR(1). The spaced-operators	355
	grammar better isolates parser performance from	356
	token-parser alignment performance, as tokeniz-	357
	ers contain some extremely grammar-unfriendly	358
	tokens. For example, the Qwen 3 tokenizer con-	359
	tains a token where the “-” character is repeated 98	360
	times, which could be interpreted as many unary	361
	minus signs by the grammar. Matching such unre-	362
	realistically long math expressions would require a	363
	more sophisticated search strategy than our simple	364
	GPU-friendly tree search.	365
	UPenn TreeBank is a small English grammar	366
	obtained from the production rules used in 300	367
	samples from UPenn TreeBank. We measure per-	368
	formance by parsing the exact same 300 samples.	369
	To simplify lexer construction, we parse part-of-	370
	speech sequences instead of word sequences. This	371
	grammar is 1000 lines long and highly ambiguous.	372
	We will release both the grammar and the data upon	373
	publication to allow future work to replicate our	374
	settings.	375
	7.2 Baselines	376
	We compare against three constrained decoding	377
	systems across various underlying parsers.	378
	XGrammar (Dong et al., 2025) emulates a push	379
	down automata (PDA), which provides general	380
	context-free recognition. However, due to the non-	381
	deterministic nature of PDA, time complexity of	382
	the system is unclear. XGrammar reduces search	383
	scope within tokenizer vocabulary by detecting to-	384
	kens that can be determined without running the	385
	underlying parser.	386
	LLGuidance (Moskal et al., 2025) uses an Ear-	387
	ley parser, allowing $O(n)$ complexity in optimal	388
	grammars and $O(n^3)$ in the general case. To search	389
	through the tokenizer vocabulary efficiently, it em-	390
	ploys a trie of tokens at the character level.	391

Method	SQL			SQL (SP-OP) ^[5]			TREEBANK		
	Cmplx.	Init	Per-mask	Cmplx.	Init	Per-mask	Cmplx.	Init	Per-mask
XGrammar (Dong et al., 2025)	?	1s	278319 μ s ^[1]	?	1s	269841 μ s ^[1]	?	1s	1e7 μ s ^[2]
LLGuidance (Earley; Moskal et al., 2025)	$O(n)$	1s	770 μ s	$O(n)$	1s	646 μ s	$O(n^3)$	1s	\times ^[3]
PSC (LALR; Anonymous, 2025)	$O(n)$	4770s	1 μ s	$O(n)$	4770s	1 μ s	\times ^[4]	\times	\times
Ours (Incremental Valiant)	$O(n^3)$	40s	8670 μ s	$O(n^3)$	40s	492 μ s	$O(n^3)$	17s	248 μ s

Table 1: Performance comparison across different grammars.

^[1] Estimated from 30 samples (totals at 10 min). ^[2] Estimated from 5 samples (totals at 30 min). ^[3] LLGuidance aborts at ~ 20 tokens; Before termination, step time between 10-20th token costs 18682 μ s on average, comparing to 250 μ s of ours in figure 6. ^[4] LALR is believed to be too restrictive for this grammar. Our attempt to initialize an LALR parser did not terminate after 3 hours. ^[5] SQL(sp-op) inserts spaces after unary operator such as minus (-) and semicolon (;). LLM tokenizer contains tokens such as 96 “-”, which cause problem with our simplistic implementation of tokenizer/parser alignment.

PSC (Anonymous, 2025) shows that in a deterministic push down automata based parser, the set of valid next terminal sequences (to a finite length) can be computed using the parser’s stack and a finite state transducer. Consequently, the number of parser calls is independent of LLM vocabulary size.

7.3 Environment and Implementation

We run experiments on a cloud VM with Intel Xeon Platinum 8480+ CPU and Nvidia H100 PCIe GPU. GPU kernels are implemented in CUDA and the tokenizer used is from Qwen 3.

7.4 Results

Our main results are summarized in table 1, which reports the complexity of the underlying parser, the constrainer initialization time, and the per-mask (i.e., per-token) overhead. In the experiments on the LALR SQL grammar, our system outperforms XGrammar significantly in per-mask overhead. When we experiment on the SQL(sp-op) grammar variant to control for the naively implemented tokenizer/constrainer alignment, we empirically outperform the Earley-based LLGuidance, despite its theoretically superior complexity. Our system still underperforms the LALR-based PSC, which takes advantage of the properties of LALR grammars, while giving up the ability to parse more general context free grammars.

In the experiments on the highly ambiguous grammar found in the TreeBank, our system shows robust efficiency without degradation. XGrammar is the only other system that runs successfully on this grammar, but 40,000 times slower in average per-mask overhead, with in some cases a single token mask taking minutes to compute. The Earley-based LLGuidance aborts at ~ 20 tokens. Before

termination, the step time between the 10th and 20th token costs 18682 μ s on average, much slower than the 250 μ s of ours. LALR parsers like PSC are commonly believed to be inherently incapable of handling complex grammars such as those found in the TreeBank. We still attempted to initialize an LALR parser using the grammar, and the construction did not terminate after 3 hours of wait. These results suggest that our parser allows enforcing a broader range of constraints that were previously too complex to be practical.

In terms of initialization time, our parser (40s/17s) is slower than XGrammar (1s/1s) and Earley-based LLGuidance (1s/1s), but still faster than PSC (4770s/ \times). Our initialization process uses mostly un-optimized Python and Numba code with redundant features for other projects. XGrammar and LLGuidance on the other hand use C and rust implemented initializations. Thus, there is significant room for improvement in the initialization of our parser.

GPU Occupancy Figure 5 shows that our constrainer only requires a fraction of stream multiprocessors (SM) on a GPU. Most of the efficiency is achieved using 16 to 32 out of over 100 available SM on an H100. In fact, each of our parser steps only uses 1 single SM, with work distributed across SMs only at the terminal trie level as described in section 6.1. The SQL grammar is more sensitive to the number of workers due to the larger number of LLM tokens containing multiple grammar terminals. This suggests our approach has the capacity to process larger batches concurrently with minimal addition to latency.

Scaling Figure 6 shows the growth of overhead per token, which has a theoretical bound of $O(n^2)$.

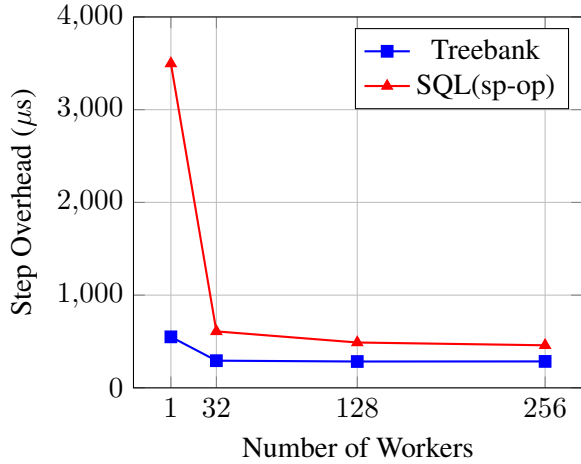


Figure 5: Each worker is a CUDA thread block that occupies at most one stream multiprocessor (out of over 100 on a single GPU). Our implementation requires only a fraction of a GPU (16 to 32 workers) to achieve most of the speed gain. The number of required workers is proportional to the extend of LLM tokens containing multiple grammar terminals.

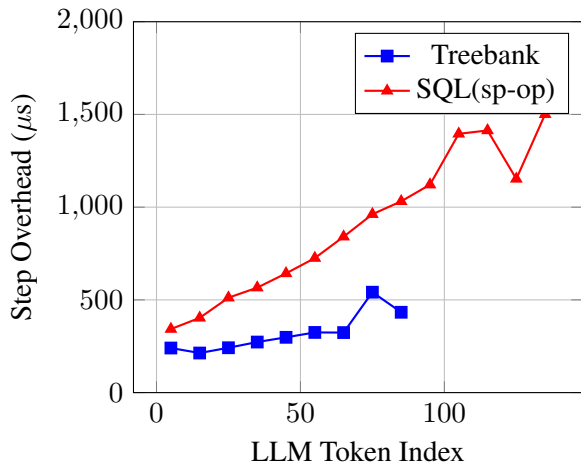


Figure 6: Although step time has theoretical growth rate of $O(n^2)$, empirical growth is more benign, likely due to the inherent and highly sparse nature of the chart. Displayed times are averages across data points from the entire benchmark set.

The empirical growth is more benign likely due to the inherent sparsity of parser charts. We also note that the constant overhead of launching a step is non-trivial, suggested by the intercept at the left-end of the curve.

Single Parser Step vs Single Constrainer Step

A constrainer step might contain multiple parser steps if a LLM token contains multiple grammar tokens, such as “”);” in SQL. In the case where such problems are less frequent, the constrainer step will be closer to the cost of the parser step, which we

observed at approximately $100\mu s$ for $n = 100$ and $|G^{CNF}| = 3219$.

8 Conclusion

We derived a new variant of the Valiant recognizer, adding support for valid prefix recognition, a critical component for supporting the use of constrained decoders with large language models. Our parsing model is parallelizable and we show that it is efficient with only a fraction of a GPU. Our approach enables high performance context-free constrained decoding over a wider range of grammars than prior work. This is critical for expanding use cases and reducing engineering effort for users of constrained decoders with large language models, who would previously have had to work within the overly restrictive LALR formalism to get efficient context free grammar based decoding.

Limitations

Our strategy to align LLM tokens and grammar terminals is primitive and causes a slow down when tokens contain multiple terminals. Improved handling of this mismatch between the grammar and the tokenizer would yield improved parser performance.

We use dense charts in our parser. However, both theoretical work (Bernardy and Claessen, 2015) and empirical experiments on GPUs (Azimov and Grigorev, 2018) show that sparse representations of parser charts can improve run time and complexity (to $O(\log^2 n)$ in some cases) of CYK-like parsers. Further research would be required to discover a version of our parser compatible with sparse charts.

Our experiments are on datasets that are associated with only two grammars. Though we believe that the SQL and TreeBank grammars are representative of useful endpoints along the spectrum of unambiguous to ambiguous grammars, a collection of grammars used in real world applications of large language models with constrained decoders would allow a deeper investigation of performance.

References

- Anonymous. 2025. PSC: Efficient grammar-constrained decoding via parser stack classification. In *Submitted to The Fourteenth International Conference on Learning Representations*. Under review.
- Rustam Azimov and Semyon Grigorev. 2018. Context-free path querying by matrix multiplication. In *Proceedings of the 1st ACM SIGMOD Joint International*

523		Xintong Sun, Chi Wei, Minghao Tian, and Shiwen Ni.	575
524		2025. Earley-driven dynamic pruning for efficient structured decoding . <i>Preprint</i> , arXiv:2506.01151.	576
525			577
526			
527	Debangshu Banerjee, Tarun Suresh, Shubham Ugare, Sasa Misailovic, and Gagandeep Singh. 2025. Crane: Reasoning with constrained llm generation . <i>Preprint</i> , arXiv:2502.09061.		
528			
529			
530			
531	Jean-Philippe Bernardy and Koen Claessen. 2015. Efficient parallel and incremental parsing of practical context-free languages . <i>Journal of Functional Programming</i> , 25:e10.		
532			
533			
534			
535	Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. 2024. Guiding llms the right way: fast, non-invasive constrained generation. In <i>Proceedings of the 41st International Conference on Machine Learning</i> , ICML'24. JMLR.org.		
536			
537			
538			
539			
540	Noam Chomsky. 1959. On certain formal properties of grammars . <i>Information and Control</i> , 2(2):137–167.		
541			
542	John Cocke. 1969. <i>Programming languages and their compilers: Preliminary notes</i> . New York University, USA.		
543			
544			
545	Yixin Dong, Charlie F. Ruan, Yaxing Cai, Ruihang Lai, Ziyi Xu, Yilong Zhao, and Tianqi Chen. 2025. Xgrammar: Flexible and efficient structured generation engine for large language models . <i>Preprint</i> , arXiv:2411.15100.		
546			
547			
548			
549			
550	Jay Earley. 1970. An efficient context-free parsing algorithm . <i>Commun. ACM</i> , 13(2):94–102.		
551			
552	Tadao Kasami. 1965. An efficient recognition and syntax analysis algorithm for context-free languages. Technical report, Air Force Cambridge Research Laboratory.		
553			
554			
555			
556	João Loula, Benjamin LeBrun, Li Du, Ben Lipkin, Clemente Pasti, Gabriel Grand, Tianyu Liu, Yahya Emara, Marjorie Freedman, Jason Eisner, Ryan Cotterell, Vikash Mansinghka, Alexander K. Lew, Tim Vieira, and Timothy J. O’Donnell. 2025. Syntactic and semantic control of large language models via sequential monte carlo . <i>Preprint</i> , arXiv:2504.13139.		
557			
558			
559			
560			
561			
562			
563	Michał Moskal, Harsha Nori, Hudson Cooper, and Loc Huynh. 2025. Llguidance: Making structured outputs go brrr .		
564			
565			
566	Shaan Nagy, Timothy Zhou, Nadia Polikarpova, and Loris D’Antoni. 2025. Chopchop: a programmable framework for semantically constraining the output of language models . <i>Preprint</i> , arXiv:2509.00360.		
567			
568			
569			
570	Elizabeth Scott and Adrian Johnstone. 2010. Gll parsing . <i>Electronic Notes in Theoretical Computer Science</i> , 253(7):177–189. Proceedings of the Ninth Workshop on Language Descriptions Tools and Applications (LDTA 2009).		
571			
572			
573			
574			
		Masaru Tomita. 1985. An efficient context-free parsing algorithm for natural languages. In <i>Proceedings of the 9th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI’85</i> , page 756–764, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.	578
			579
			580
			581
			582
			583
		Shubham Ugare, Rohan Gumaste, Tarun Suresh, Gagandeep Singh, and Sasa Misailovic. 2025. Itergen: Iterative semantic-aware structured LLM generation with backtracking . In <i>The Thirteenth International Conference on Learning Representations</i> .	584
			585
			586
			587
			588
		Shubham Ugare, Tarun Suresh, Hangoo Kang, Sasa Misailovic, and Gagandeep Singh. 2024. Syncode: Llm generation with grammar augmentation . <i>Preprint</i> , arXiv:2403.01632.	589
			590
			591
			592
		Leslie G. Valiant. 1975. General context-free recognition in less than cubic time . <i>Journal of Computer and System Sciences</i> , 10(2):308–315.	593
			594
			595
		Daniel H. Younger. 1967. Recognition and parsing of context-free languages in time n^3 . <i>Information and Control</i> , 10(2):189–208.	596
			597
			598
		A Correctness of ANCESTOR[†]	599
		A.1 Completeness	600
		Using induction, we prove ANCESTOR [†] [j][V] if V has valid prefix $t_{n-j} \dots t_{n-1}$.	601
			602
		• Base case: $j = 0$	603
		– invalid index, there is no t_n yet,	604
		$t_n \dots t_{n-1}$ is empty string	605
		– every non-terminal has empty string as its valid prefix	606
		– we set ANCESTOR [†] [0][:] = True	607
			608
		• Induction: $j = 1 \dots$	609
		– Induction Hypothesis:	610
		ANCESTOR [†] [k][V] if V has valid prefix $t_{n-k} \dots t_{n-1}$ for any $k < j$	611
		\Rightarrow ANCESTOR [†] [j][V] if V has valid prefix $t_{n-j} \dots t_{n-1}$	612
			613
			614
			615
		– Assume V has valid prefix $t_{n-j} \dots t_{n-1}$	616
		– There exists rule $V \rightarrow AB$, where	617
		* $A \Rightarrow^* t_{n-j} \dots t_{n-j+mid}$	618
		* $B \Rightarrow^* t_{n-j+mid+1} \dots$	619
		– Depending on mid, there are two cases	620
		* Case 1: A only generates part of or exactly entire of the committed prefix	621
			622

623	1. formally	· TIER $[n - j + mid][n - j][A]$,	670
624	· $n - j + mid \leq n - 1$	and	671
625	· $A \Rightarrow^* t_{n-j} \dots t_{n-j+mid}$	· ANCESTOR $^\uparrow[j - mid - 1][B]$	672
626	· $B \Rightarrow^* t_{n-j+mid+1} \dots t_{n-1} \dots$	2. $A \Rightarrow^* t_{n-j} \dots t_{n-j+mid}$, by	673
627	2. $A \in \mathbf{TIER}[n - j + mid][n - j]$,	CYK, 1	674
628	by 1 and CYK	3. $B \Rightarrow^* t_{n-j+mid+1} \dots t_{n-1} \dots$,	675
629	3. $B \in \mathbf{ANCESTOR}^\uparrow[j - mid - 1]$,	by IH, 1	676
630	by 1 and IH	4. $V \Rightarrow^* t_{n-j} \dots t_{n-1} \dots$, by 1, 2,	677
631	4. ANCESTOR $^\uparrow[j][V]$, by case 1 in	3	678
632	pseudo-code, 2, 3	* Case 2:	679
633	* Case 2: A generates all of the com-	1. ANCESTOR $^\uparrow[j][W]$ is set to true	680
634	mitted prefix and beyond	by condition	681
635	1. formally	· FIRST $[V][W]$, and	682
636	· $n - j + mid > n - 1$	· ANCESTOR $^\uparrow[j][V]$, caused by	683
637	· $A \Rightarrow^* t_{n-j+mid} \dots t_{n-1} \dots$	case 1	684
638	2. Since the grammar is Chomsky	2. $V \Rightarrow^* t_{n-j} \dots t_{n-1} \dots$, by	685
639	Normal form, it is guaranteed that	case 1	686
640	$\exists W$ s.t.	3. $W \Rightarrow^* t_{n-j} \dots t_{n-1} \dots$, by	687
641	· $W \in \mathbf{FIRST}[V]$, and	property of FIRST set	688
642	· $W \rightarrow CD$, and		
643	· W belongs to case 1. (i.e. C gen-		
644	erates part or exactly the commit-		
645	ted prefix		
646	3. ANCESTOR $^\uparrow[j][W]$, by IH, 2		
647	4. ANCESTOR $^\uparrow[j][V]$, by case 2 in		
648	pseudo-code, 3		

649 A.2 Soundness

650 Using induction, we prove **ANCESTOR** $^\uparrow[j][V]$
651 only if V has valid prefix $t_{n-j} \dots t_{n-1}$.

652 • **Base case:** $j = 0$

- 653 – invalid index, there is no t_n yet,
- 654 $t_n \dots t_{n-1}$ is empty string
- 655 – every non-terminal has empty string as
- 656 its valid prefix
- 657 – we set **ANCESTOR** $^\uparrow[0][:] = True$

658 • **Induction:** $j = 1 \dots$

- 659 – Induction Hypothesis:
- 660 **ANCESTOR** $^\uparrow[k][V]$ only if V has valid
- 661 prefix $t_{n-k} \dots t_{n-1}$ for any $k < j$
- 662 \Rightarrow **ANCESTOR** $^\uparrow[j][V]$ only if V has
- 663 valid prefix $t_{n-j} \dots t_{n-1}$

- 664
- 665 – Assume **ANCESTOR** $^\uparrow[j][V]$

666 * **Case 1:**

- 667 1. **ANCESTOR** $^\uparrow[j][V]$ is set to true
- 668 by condition
- 669 · $V \Rightarrow AB$, and