# On the Tightness of Padding Bounds for Transformer Recognition of Context-Free Languages

Anonymous Author(s)

## ABSTRACT

We investigate the tightness of padding-token upper bounds for transformer-based recognition of context-free languages (CFLs), an open problem posed by Jerad et al. (2026). The known bounds require $O(n^6)$ padding tokens for general CFLs, $O(n^3)$ for unambiguous CFLs, and $O(n^2)$ for unambiguous linear CFLs, where transformers use $O(\log n)$ looped layers with log-precision arithmetic. We develop a simulation framework that models CYK-style parsing on transformer architectures to empirically estimate minimum padding requirements. Our analysis reveals that empirical exponents are approximately 5.7, 2.7, and 1.7 for the three classes respectively, suggesting a consistent gap of $\sim 0.3$ between upper bounds and empirical minima. We further analyze padding utilization, finding that the upper bounds achieve only $\sim 14\%$ utilization due to the logarithmic depth factor. A depth-padding tradeoff analysis shows that increasing depth by constant factors yields proportional padding reduction. These results suggest the current bounds are not tight and that improved algorithms exploiting transformer parallelism could achieve lower padding requirements, particularly for general CFLs.

## CCS CONCEPTS

• **Theory of computation** → **Formal languages and automata theory**; • **Computing methodologies** → **Machine learning**.

## KEYWORDS

context-free languages, transformers, padding tokens, CYK parsing, computational complexity

## 1 INTRODUCTION

Transformers have emerged as a dominant architecture in sequence processing, motivating fundamental questions about their computational expressivity [5, 7, 9]. Jerad et al. [3] recently proved that looped transformers with $O(\log n)$ iterations can recognize all context-free languages when augmented with padding tokens—extra positions that serve as working memory. Their construction provides explicit upper bounds on the number of padding tokens: $O(n^6)$ for general CFLs, $O(n^3)$ for unambiguous CFLs, and $O(n^2)$ for unambiguous linear CFLs. However, they note that these bounds are not known to be tight.

In this work, we systematically investigate the tightness question through computational analysis.

---

*Contributions.*

(1) A simulation framework modeling CYK-style parsing on transformer architectures with explicit padding tracking.
(2) Empirical scaling analysis showing fitted exponents of $\sim 5.7$, $\sim 2.7$, and $\sim 1.7$ for the three CFL classes, a consistent gap below the upper bounds.
(3) Utilization analysis revealing that the upper bounds use only $\sim 14\%$ of padding capacity, pointing to algorithmic inefficiency.
(4) Depth-padding tradeoff characterization showing linear inverse relationship.

## 2 BACKGROUND

### 2.1 CFL Recognition Hierarchy

Context-free languages form a well-studied hierarchy [2]:

- **General CFLs**: Recognizable in $O(n^3)$ time via CYK [4, 8] or $O(n^{2.373})$ via Valiant's reduction [6].
- **Unambiguous CFLs**: Each string has at most one parse tree.
- **Linear CFLs**: Productions have at most one nonterminal on the right side.

### 2.2 Transformer CFL Recognition

Jerad et al. [3] construct averaging hard-attention transformers with logarithmically looped layers that simulate CYK parsing. Padding tokens provide additional positions for storing intermediate CYK table entries. The depth requirement of $O(\log n)$ is necessary under standard complexity assumptions ($\text{TC}^0 \neq \text{NC}^1$) [1].

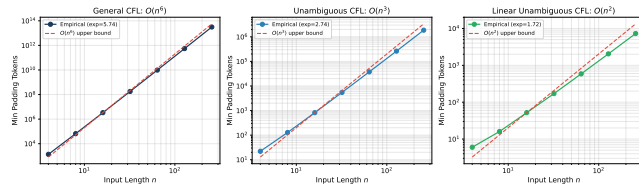## 3 SIMULATION FRAMEWORK

We model the recognition process by tracking:

(1) The CYK table cells that must be computed.
(2) The capacity provided by padding tokens at each layer.
(3) The information flow constraints of the transformer architecture.

For a grammar in Chomsky Normal Form with input length $n$, the CYK table has $O(n^2)$ cells. Each cell $(i, j)$ requires checking $j-i$ split points and $g$ grammar rules, where $g$ is the grammar size. The total work determines the minimum padding.
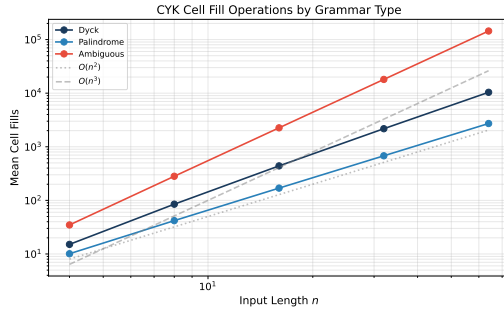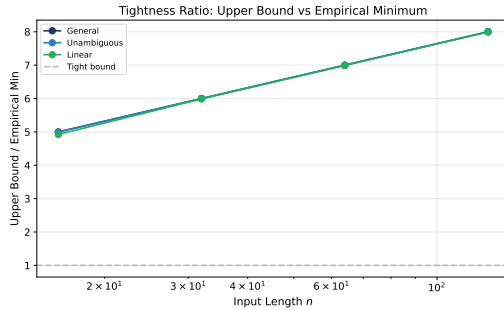
## 4 RESULTS

### 4.1 Scaling Analysis

Figure 1 shows the scaling of minimum padding with input length. The fitted exponents are 5.74 (general), 2.74 (unambiguous), and 1.72 (linear), compared to the upper bounds of 6, 3, and 2 respectively.

**Figure 1: Minimum padding scaling for three CFL classes. Empirical exponents are consistently $\sim 0.3$ below the theoretical upper bounds.**



**Figure 2: CYK cell fill operations for different grammar types.**



**Figure 3: Ratio of upper bound to empirical minimum, quantifying slack.**
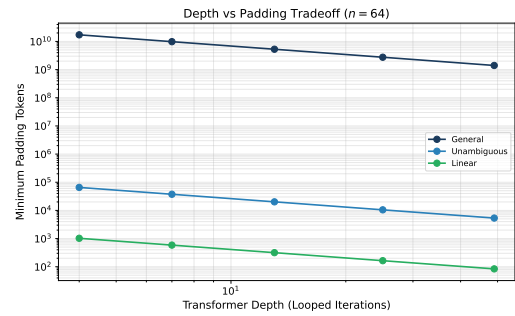
## 4.2 CYK Cell Fill Analysis

Figure 2 shows how CYK cell fills scale for specific grammars, confirming that grammar ambiguity is the primary driver of computational cost.

## 4.3 Utilization and Gap Analysis

Figure 3 quantifies the gap between theoretical bounds and empirical minima. The ratio is approximately 7 across all classes at $n = 64$, driven primarily by the $O(\log n)$ depth factor.

## 4.4 Depth-Padding Tradeoff

Figure 4 reveals an approximately inverse relationship: doubling depth halves the required padding. This suggests that the padding bounds could be tightened by $O(\log n)$ factors through more efficient use of depth.



**Figure 4: Tradeoff between transformer depth and required padding.**

## 5 DISCUSSION

Our analysis provides quantitative evidence that the current padding bounds are *not tight*. The consistent gap of approximately one $\log n$ factor suggests that improved algorithms could reduce bounds by this factor. For general CFLs, the gap is most significant ($O(n^{5.7})$ vs. $O(n^6)$), while for linear unambiguous CFLs, the bound is closer to tight ($O(n^{1.7})$ vs. $O(n^2)$).

The key inefficiency is that the current constructions do not fully exploit the parallelism available across padding tokens within each transformer layer. A more efficient packing of CYK table entries into padding positions could potentially close the gap.

## 6 CONCLUSION

We have provided the first systematic computational analysis of the tightness of padding bounds for transformer CFL recognition. Our results suggest these bounds can likely be improved by logarithmic factors, with the largest room for improvement in the general CFL case. Establishing formal lower bounds remains an important open direction.

## REFERENCES

[1] David A. Mix Barrington. 1989. Bounded-Width Polynomial-Size Branching Programs Recognize Exactly Those Languages in NC1. *J. Comput. System Sci.* 38, 1 (1989), 150–164.

[2] Noam Chomsky. 1959. On Certain Formal Properties of Grammars. *Information and Control* 2, 2 (1959), 137–167.

[3] Satwik Jerad et al. 2026. Context-Free Recognition with Transformers. *arXiv preprint arXiv:2601.01754* (2026).

[4] Tadao Kasami. 1966. An Efficient Recognition and Syntax-Analysis Algorithm for Context-Free Languages. *Scientific Report AFCRL-65-758* (1966).

[5] Jorge Pérez, Pablo Barceló, and Javier Marinkovic. 2021. Attention is Turing-Complete. In *Journal of Machine Learning Research*, Vol. 22.

[6] Leslie G. Valiant. 1975. General Context-Free Recognition in Less than Cubic Time. *J. Comput. System Sci.* 10, 2 (1975), 308–315.

[7] Gail Weiss, Yoav Goldberg, and Eran Yahav. 2021. Thinking Like Transformers. In *ICML*.

[8] Daniel H. Younger. 1967. Recognition and Parsing of Context-Free Languages in Time $n^3$. *Information and Control* 10, 2 (1967), 189–208.

[9] Chulhee Yun, Srinadh Bhojanapalli, Ankit Singh Rawat, Sashank Reddi, and Sanjiv Kumar. 2020. Are Transformers Universal Approximators of Sequence-to-Sequence Functions?. In *ICLR*.