

# Prefix Parsing is Just Parsing

Clemente Pasti<sup>15</sup> Andreas Opedal<sup>1</sup> Timothy J. O’Donnell<sup>2345</sup>

Ryan Cotterell<sup>1</sup> Tim Vieira<sup>1</sup>

{clemente.pasti, andreas.opedal, ryan.cotterell}@inf.ethz.ch

timothy.odonnell@mcgill.ca tim.f.vieira@gmail.com

<sup>1</sup>ETH Zürich <sup>2</sup>McGill University <sup>3</sup>Canada CIFAR AI Chair <sup>4</sup>Mila <sup>5</sup>CHI-FRO

## Abstract

Prefix parsing asks whether an input prefix can be extended to a complete string generated by a given grammar. In the weighted setting, it also provides prefix probabilities, which are central to context-free language modeling, psycholinguistic analysis, and syntactically constrained generation from large language models. We introduce the *prefix grammar transformation*, an efficient reduction of prefix parsing to ordinary parsing. Given a grammar, our method constructs another grammar that generates exactly the prefixes of its original strings. Prefix parsing is then solved by applying any ordinary parsing algorithm on the transformed grammar without modification. The reduction is both elegant and practical: the transformed grammar is only a small factor larger than the input, and any optimized implementation can be used directly, eliminating the need for bespoke prefix-parsing algorithms. We also present a strategy—based on algorithmic differentiation—for computing the next-token weight vector, i.e., the prefix weights of *all* one-token extensions, enabling efficient prediction of the next token. Together, these contributions yield a simple, general, and efficient framework for prefix parsing.

 [genlm/prefix-acl-26](https://github.com/genlm/prefix-acl-26)

## 1 Introduction

Parsing a string with respect to a context-free grammar (CFG) means determining whether—and with what weight—the grammar generates that string. *Prefix* parsing asks the closely related question of whether an input prefix can be extended to a complete string generated by the grammar. In the probabilistic setting,<sup>1</sup> these tasks correspond, respectively, to computing (i) the probability that the grammar generates a given string, and (ii) the total probability of strings beginning with a given prefix.

Existing prefix parsers are typically tied to particular parsing algorithms. For example, Jelinek

and Lafferty (1991) adapted CKY (Cocke and Schwartz, 1970; Kasami, 1965; Younger, 1967), while Stolcke (1995) adapted Earley’s (1970) algorithm.<sup>2</sup> Although these algorithms require bespoke and surprisingly involved modifications to the base parser, they also share notable similarities. This raises a natural question: is there a general-purpose reduction from prefix parsing to parsing, one that would let us adapt *other* parsers to the prefix setting, e.g., Valiant’s (1975) subcubic parsing algorithm, GPU implementations,<sup>3</sup> and fast CPU parsers?<sup>4</sup> We show that such a reduction exists. Given any CFG, we produce a new CFG—the *prefix grammar*—that generates exactly the weighted prefix language of the original, and is only a small constant factor larger. Consequently, *any* parsing algorithm serves as a prefix parser off-the-shelf: its input-length dependence is preserved, with the entire overhead confined to grammar-structural factors (Theorem 2). How much those factors grow depends on the parser’s preprocessing pipeline; we analyze them for CKY and Earley (App. E).<sup>5</sup> Moreover, the prefix grammar can be viewed as the composition of the original grammar with a two-state prefix transducer (App. D), a perspective that extends naturally to other grammar formalisms.

Finally, we develop efficient algorithms for computing the *next-token weights* of a string—the prefix weights of all one-token extensions. Surprisingly, our strategy—which leverages algorithmic differentiation (see, e.g., Griewank and Walther, 2008; Eisner, 2016) and is closely related to the outside algorithm (Baker, 1979)—allows us to compute the next-token weights of a string in the same asymptotic runtime as prefix-parsing the string itself. On the other hand, a naïve approach would incur a multiplicative factor of  $|\Sigma|$ , corresponding to  $|\Sigma|$

<sup>1</sup>More generally, our method applies to grammars whose weights lie in any commutative semiring (see, e.g., Goodman, 1999). This includes the boolean semiring, which is useful for enforcing grammaticality constraints in language models (e.g., Shin et al., 2021; Poesia et al., 2022; Loula et al., 2025).

<sup>2</sup>Nowak and Cotterell (2023) developed a prefix parser based on CKY with the *hook trick* (Eisner and Blatz, 2007).

<sup>3</sup>E.g., Stanojević and Sartran (2023), Rush (2020), Yi et al. (2011), Dunlop (2014), and Canny et al. (2013).

<sup>4</sup>E.g., Kegler (2023) and Dunlop (2014).

<sup>5</sup>We found that our implementation of Earley’s algorithm with the prefix grammar is comparable to Luong et al.’s (2013) implementation of Stolcke’s algorithm (App. B).

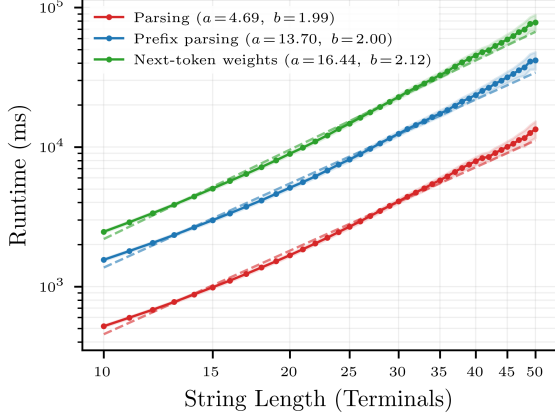


Figure 1: We compare parsing, prefix parsing, and the next-token weight vector algorithm on WSJ 5000 (Luong et al., 2013, EarleyX), a large grammar with 35,016 rules, 448 non-terminals, and a total size of 116,667. The underlying parser is Earley’s (see Alg. 4), paired with its next-token variant (§4 and App. F.2). Results are shown for 500 strings on a log–log plot; to obtain data across all string lengths, we include the runtime for each prefix of every string. Shaded regions represent 95% confidence intervals for the mean runtime at each string length. To estimate runtime complexity, we fit a power law  $r(N) = aN^b$  via least-squares regression on the log–log scale:  $\log r(N) = \log a + b \log N$ . Dashed lines show the fitted models. All three algorithms exhibit similar fitted complexity exponents (1.99, 2, and 2.12 respectively), confirming that the  $N$ -dependence is essentially independent of grammar size. Prefix parsing is roughly  $2.9\times$  slower than parsing (fitted coefficient  $a = 30.9$  vs.  $12.1$ ), closely tracking the prefix grammar’s  $\approx 2.86\times$  size increase (binarize, then prefix). Both observations are predicted by Theorem 2: the  $N$ -dependence is preserved, and the overhead is absorbed into the grammar-size factor. The next-token algorithm is slower than prefix parsing by only a factor of  $\approx 1.2$ , well within the constant-factor bound (i.e.,  $4\times$ ) of the next-token meta-theorem (Theorem 3).

calls to a prefix parser, where  $\Sigma$  is the token alphabet. Importantly, next-token weights arise in many applications, including context-free language modeling (Jelinek and Lafferty, 1991), surprisal-based psycholinguistic modeling (Hale, 2001), and syntactically constrained generation from LLMs (Shin et al., 2021; Poesia et al., 2022; Loula et al., 2025).

## 2 Preliminaries<sup>6</sup>

**Basics.** Let  $\Sigma$  be an **alphabet**, i.e., a finite set of symbols, and let  $\Sigma^*$  denote the **set of all strings** formed from the symbols of  $\Sigma$ , including the empty string  $\varepsilon$ . We write  $s \preceq t$  if  $s$  is a prefix of  $t$ , and  $st$  to denote concatenation. Let  $(\mathbb{W}, +, \cdot, 0, 1)$  be a commutative semiring, i.e., a set equipped with addition and multiplication satisfying the usual distributive and identity laws. For concreteness, the reader may take  $\mathbb{W} = \mathbb{R}_{\geq 0}$ . A **weighted language** is a function  $\omega: \Sigma^* \rightarrow \mathbb{W}$ . We call  $\omega$  a **language model** when  $\mathbb{W} = \mathbb{R}_{\geq 0}$  and  $\sum_{s \in \Sigma^*} \omega(s) = 1$ .

<sup>6</sup>For convenience, we provide a notation glossary in App. A.

**Prefix languages.** A central concept in this paper is the (weighted) **prefix language** of  $\omega$ :

$$\vec{\omega}(s) \stackrel{\text{def}}{=} \sum_{t \in \Sigma^*} \mathbb{1}\{s \preceq t\} \omega(t) \quad (1)$$

In the language model setting where  $\mathbb{W} = \mathbb{R}_{\geq 0}$ , the weight  $\omega(s)$  of any string  $s$  factorizes as a product of conditional prefix probabilities:

$$\omega(s) = \vec{\omega}(\text{EOS} \mid s) \prod_{t=1}^{|s|} \vec{\omega}(s_t \mid s_{<t}) \quad (2)$$

where  $s_{<t} \stackrel{\text{def}}{=} s_1 \cdots s_{t-1}$ ,  $\text{EOS} \notin \Sigma$  is a distinguished **end-of-string symbol**, and the **conditional next-token weights** are defined as follows:

$$\vec{\omega}(s_t \mid s_{<t}) \stackrel{\text{def}}{=} \begin{cases} \frac{\omega(s_{<t} s_t)}{\vec{\omega}(s_{<t})} & \text{if } s_t = \text{EOS} \\ \frac{\vec{\omega}(s_{<t} s_t)}{\vec{\omega}(s_{<t})} & \text{otherwise} \end{cases} \quad (3)$$

provided that  $\vec{\omega}(s_{<t}) > 0$  and  $s_t \in \Sigma \cup \{\text{EOS}\}$ . Note that  $\vec{\omega}(\cdot \mid s_{<t})$  is a distribution over  $\Sigma \cup \{\text{EOS}\}$ . Conditional prefix probabilities underlie context-free language modeling and surprisal-based psycholinguistic analysis. In §4, we develop efficient algorithms for computing the prefix weights of all single-symbol extensions simultaneously.

A **weighted context-free grammar**<sup>7</sup> (WCFG, or simply **CFG**)  $\mathcal{G}$  is a tuple  $(\mathcal{N}, \Sigma, S, \mathcal{R})$ , where  $\mathcal{N}$  is a set of **nonterminal** symbols,  $\Sigma$  is an alphabet of **terminal** symbols (disjoint from  $\mathcal{N}$ ),  $S \in \mathcal{N}$  is a distinguished **start symbol**, and  $\mathcal{R}$  is a finite set of weighted **rules**.<sup>8</sup> Each rule  $r \in \mathcal{R}$  is of the form  $X \xrightarrow{w} \alpha$  where  $X \in \mathcal{N}$ ,  $\alpha \in (\Sigma \cup \mathcal{N})^*$ , and  $w \in \mathbb{W}$ . The **arity** of a rule  $r$ , denoted  $\text{ar}(r)$ , is the length of its right-hand side. We say that a rule is *nullary* if its arity is zero, *unary* if its arity is one, *binary* if its arity is two, and so on. The **size** of a grammar is  $|\mathcal{G}| = \sum_{r \in \mathcal{R}} 1 + \text{ar}(r)$ . We say that a grammar is in **canonical two-form** (CTF) if each of its rules has a right-hand side of length at most two, and **Chomsky normal form** (CNF) if each of its rules takes one of the forms  $X \rightarrow YZ$ ,  $X \rightarrow a$ , or  $S \rightarrow \varepsilon$  for  $X, Y, Z \in \mathcal{N}$ ,  $a \in \Sigma$ .

A **derivation tree** is a rooted, ordered tree in which each node is labeled with a symbol in  $\mathcal{N} \cup \Sigma$ , each nonleaf node is connected to its children by

<sup>7</sup>We assume the reader is familiar with WCFGs; for completeness, App. C.1 provides the necessary background.

<sup>8</sup>Without loss of generality, any two rules with the same left-hand side  $X$  and right-hand side  $\alpha$ ,  $X \xrightarrow{w} \alpha$  and  $X \xrightarrow{w'} \alpha$ , can be consolidated into the single rule  $X \xrightarrow{w+w'} \alpha$ .

a rule in  $\mathcal{R}$ , and each leaf is either a terminal or a childless nonterminal (arising from a nullary rule). The **weight** of a derivation tree is the product of the weights of all rules applied in the tree. Let  $\mathcal{D}_\alpha(s)$  denote the set of derivation trees rooted at  $\alpha \in \mathcal{N} \cup \Sigma$  whose **yield** is  $s$ . The **weighted language** of  $\alpha$  is  $\llbracket \mathcal{G}_\alpha \rrbracket(s) \stackrel{\text{def}}{=} \sum_{d \in \mathcal{D}_\alpha(s)} w(d)$ , and the **weighted language** of  $\mathcal{G}$  is  $\llbracket \mathcal{G} \rrbracket(s) \stackrel{\text{def}}{=} \llbracket \mathcal{G}_S \rrbracket(s)$ . The **total weight** of  $\alpha$  is  $\tau(\alpha) \stackrel{\text{def}}{=} \sum_{s \in \Sigma^*} \llbracket \mathcal{G}_\alpha \rrbracket(s)$ . A (tight) **probabilistic context-free grammar (PCFG)** is one where  $\tau(\alpha) = 1$  for all  $\alpha \in \mathcal{N} \cup \Sigma$ . If  $\mathcal{G}$  is a PCFG, then  $\llbracket \mathcal{G} \rrbracket$  is a language model (Booth and Thompson, 1973; Chi, 1999).

### 3 Prefix Parsing is Just Parsing

This section introduces the **prefix grammar transformation**, which allows us to turn any CFG parser into a *prefix* parser simply by transforming its grammar into one that generates its prefix language.<sup>9,10</sup>

**Definition 1.** Given a grammar  $\mathcal{G} = \langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle$ , its **prefix grammar** is

$$\vec{\mathcal{G}} \stackrel{\text{def}}{=} \langle \mathcal{N} \cup \{X'\}_{X \in \mathcal{N}} \cup \{\vec{S}\}, \Sigma, \vec{S}, \mathcal{R} \cup \mathcal{R}' \rangle$$

where  $\alpha'$  denotes  $\alpha$  itself if  $\alpha \in \Sigma$  and a new prime nonterminal if  $\alpha \in \mathcal{N}$ . Here  $\vec{S}$  is a new start symbol, and the additional rules  $\mathcal{R}'$  are defined below:

$$\vec{S} \xrightarrow{1} S' \quad (4a)$$

$$\vec{S} \xrightarrow{\tau(S)} \varepsilon \quad (4b)$$

$$X' \xrightarrow{w \cdot \tau(\alpha_{k+1}) \cdots \tau(\alpha_K)} \alpha_1 \cdots \alpha_{k-1} \alpha'_k \quad (4c)$$

$$\text{for } X \xrightarrow{w} \alpha_1 \cdots \alpha_K \in \mathcal{R}, k \in 1, \dots, K$$

The idea underlying the prefix grammar is that, at each level of the derivation tree, a prefix of a derived string must end inside the yield of exactly one rule application. This is captured by Eq. (4c): for each original rule  $X \xrightarrow{w} \alpha_1 \cdots \alpha_K$  and each border position  $k$ , the symbols  $\alpha_1, \dots, \alpha_{k-1}$  before the border are parsed normally, the symbol  $\alpha_k$  at the border is only partially matched (handled

<sup>9</sup>Prefix grammars themselves are not novel. Indeed, deriving one is an exercise on context-free grammars in Ben-David's (2022) lecture notes. However, we were unable to find a *weighted* CFG version; much to our delight, that extension proved straightforward. In this light, our contribution is to demonstrate how the prefix grammar can be used to *reduce* prefix parsing to ordinary parsing (with weights).

<sup>10</sup>Note that the prefix grammar can also be defined by *composition* (Bar-Hillel et al., 1961; Pasti et al., 2023) of the input grammar and a specifically designed *prefix transducer*; we discuss this view in detail in App. D.

recursively via  $\alpha'_k$ ), and the symbols  $\alpha_{k+1}, \dots, \alpha_K$  after the border are summarized by their total weights  $\tau(\alpha_{k+1}) \cdots \tau(\alpha_K)$ . Recall that  $\tau(\alpha) = 1$  for PCFGs (§2); for general WCFGs, the total weights can be computed as described in App. C.1.

The following proposition establishes that Def. 1 correctly encodes the prefix language.

**Proposition 1.** Let  $\mathcal{G}$  be a CFG, and  $\vec{\mathcal{G}}$  be its prefix grammar. Then,  $\vec{\mathcal{G}}$  correctly encodes the prefix language of  $\mathcal{G}$  (i.e.,  $\llbracket \vec{\mathcal{G}} \rrbracket = \llbracket \mathcal{G} \rrbracket$ ).

*Proof.* See App. G.1. ■

This provides a straightforward reduction of any ordinary parsing algorithm to a prefix-parsing one. We say  $\mathfrak{p}$  is a **correct parsing algorithm** if  $\mathfrak{p}(\mathcal{G}, s) = \llbracket \mathcal{G} \rrbracket(s)$  for all grammars  $\mathcal{G}$  and strings  $s$ . Although this framing treats  $\mathfrak{p}$  as a black box over arbitrary CFGs, most parsers in the literature are in fact defined only on a restricted class of grammars (e.g., CNF, nullary-free) and rely on an upstream **normal-form conversion**  $\phi$  (e.g., CNF conversion for CKY) that is typically treated implicitly; its size cost is frequently understated and propagates through our reduction—we track it explicitly in Theorem 2. The following theorem shows that our reduction yields a correct prefix-parsing algorithm.

**Theorem 1** (Prefix parsing is just parsing). Let  $\mathcal{G}$  be a CFG over the alphabet  $\Sigma$ . Let  $\vec{\mathcal{G}}$  be its prefix grammar. If  $\mathfrak{p}$  is a correct parsing algorithm, then  $\mathfrak{p}(\vec{\mathcal{G}}, s) = \llbracket \mathcal{G} \rrbracket(s)$  for all  $s \in \Sigma^*$ .

*Proof.* The theorem immediately follows from Prop. 1 and the premises. ■

In App. E, we describe specific instantiations of this approach with CKY and Earley. The key to the efficiency of our reduction is that we can bound the size of the prefix grammar. For efficiency, it is advisable to binarize the grammar (i.e., convert to canonical two-form) before applying the prefix grammar transformation. Although binarization triples the grammar size, it bounds the rule arity at two, which in turn keeps the prefix grammar within a small constant factor of the original:

**Proposition 2.** Let  $\mathcal{G}$  be a context-free grammar in canonical two-form and let  $\vec{\mathcal{G}}$  be its prefix grammar (Def. 1). Then, the size of  $\vec{\mathcal{G}}$  is bounded by

$$|\vec{\mathcal{G}}| \leq \frac{8}{3} |\mathcal{G}| + 3 \quad (5)$$

*Proof.* See App. G.2. ■

Prop. 2 bounds  $\overrightarrow{|\mathcal{G}|}$  directly. However, as noted above, most parsers’ runtime bounds involve a pre-processed grammar  $\phi(\mathcal{G})$  (e.g., CNF conversion for CKY) and often depend on structural parameters beyond total size, such as  $|\mathcal{N}|$ , which grow differently under  $\phi$  (App. C.2). The following theorem lifts Prop. 2 to parser runtime by allowing the runtime to be any polynomial in string length and a vector  $\nu$  of grammar-structural parameters.

**Theorem 2.** *Let  $\mathfrak{p}$  be a parsing algorithm whose runtime after grammar preprocessing by  $\phi$  satisfies*

$$\text{time}\{\mathfrak{p}(\mathcal{G}, \mathbf{s})\} = \mathcal{O}\left(\sum_{i=1}^m g_i(\nu(\phi(\mathcal{G}))) \cdot h_i(|\mathbf{s}|)\right)$$

where  $\nu$  maps a grammar  $\mathcal{G}' = \langle \mathcal{N}', \Sigma', \mathcal{S}', \mathcal{R}' \rangle$  to a vector of its structural parameters (e.g.,  $\nu(\mathcal{G}') = [|\mathcal{G}'|, |\mathcal{N}'|, |\mathcal{R}'|, |\Sigma'|]^T$ ), and each  $g_i$  and  $h_i$  is a polynomial.<sup>11</sup> Then,

$$\text{time}\{\mathfrak{p}(\overrightarrow{\mathcal{G}}, \mathbf{s})\} = \mathcal{O}\left(\sum_{i=1}^m g_i(\nu(\phi(\overrightarrow{\mathcal{G}}))) \cdot h_i(|\mathbf{s}|)\right)$$

In particular, the  $|\mathbf{s}|$ -dependence of each term is unchanged; the entire overhead of prefix parsing is captured by replacing  $\nu(\phi(\mathcal{G}))$  with  $\nu(\phi(\overrightarrow{\mathcal{G}}))$ .

*Proof.* Apply  $\mathfrak{p}$ ’s runtime guarantee with input  $\overrightarrow{\mathcal{G}}$  in place of  $\mathcal{G}$ . ■

Observe that Theorem 2 makes no claim about the magnitude of  $\text{time}\{\mathfrak{p}(\mathcal{G}, \mathbf{s})\}$  itself; that is a property of the parser. For example, CKY runs in  $\mathcal{O}(|\mathcal{G}| \cdot |\mathbf{s}|^3)$  on CNF grammars, which fits the theorem with  $\phi$  set to CNF conversion. Because the unary-removal step of  $\phi$  can inflate grammar size by a factor of  $|\underline{\mathcal{N}}|$  (App. C.2), applying CKY to  $\overrightarrow{\mathcal{G}}$  runs in  $\mathcal{O}(|\phi(\overrightarrow{\mathcal{G}})| \cdot |\mathbf{s}|^3)$ —a factor of  $|\mathcal{N}|$  looser than the bound for  $\mathcal{G}$  (App. E.1). Earley’s lighter preprocessing incurs less blowup (App. E.2).

#### 4 Next-Token Weight Vector Algorithms

Given a string  $\mathbf{s} = s_1 \cdots s_N$  and a CFG  $\mathcal{G}$ , we now study how to efficiently compute the **next-token weight vector**  $\pi_a(\mathbf{s}) = \overrightarrow{[\mathcal{G}]}(sa)$ , whose components give the prefix weight of each one-symbol extension  $sa$  for  $a \in \Sigma$ . This vector is the key ingredient for the incremental processing applications discussed in §1, including constrained generation with a PCFG. Importantly, by Prop. 1, computing the next-token weight vector reduces to ordinary parsing:  $\pi_a(\mathbf{s}) = \overrightarrow{[\mathcal{G}]}(sa)$ . Yet a naïve implementation would still require  $|\Sigma|$  parser calls,

<sup>11</sup>We assume the one-time cost of  $\phi$  is amortized over all strings parsed with the same grammar.

one per  $a \in \Sigma$ , i.e.,  $\mathcal{O}(|\Sigma| \cdot \text{time}\{\mathfrak{p}(\overrightarrow{\mathcal{G}}, \mathbf{s})\})$  time per next-token vector  $\pi(\mathbf{s})$ .

A less naïve approach is to use an **incremental parsing algorithm**—one whose evaluation on a prefix  $s_1 \cdots s_{N-1}$  leaves enough cached information to extend to  $s_1 \cdots s_N$  cheaply. We refer to this cached information as the **prefix state** and assume it is managed by memoization, so that the same prefix is parsed at most once. The precise representation of the prefix state is parser-specific; for incremental CKY (Alg. 2 in App. E) it is the portion of the parse chart filled in for  $s_1 \cdots s_{N-1}$ . Non-incremental CKY on  $\overrightarrow{\mathcal{G}}$  runs in  $\mathcal{O}(|\phi(\overrightarrow{\mathcal{G}})|N^3)$ ; incremental CKY amortizes this cost across prefixes, extending the parse by one token in only  $\mathcal{O}(|\phi(\overrightarrow{\mathcal{G}})|N^2)$ . Immutability lets us safely reuse the prefix state across all  $|\Sigma|$  candidate one-token extensions of  $s_1 \cdots s_{N-1}$ , as the cached entries from the prefix do not depend on the choice of next token. This reduces the baseline to  $\mathcal{O}(|\phi(\overrightarrow{\mathcal{G}})|N^2 \cdot |\Sigma|)$  per token—an improvement by a factor of  $N$ —but the  $|\Sigma|$  factor remains.

In contrast, we propose an approach that eliminates the  $|\Sigma|$  factor entirely, computing the full vector  $\pi(\mathbf{s})$  in  $\mathcal{O}(\text{time}\{\mathfrak{p}(\overrightarrow{\mathcal{G}}, \mathbf{s})\})$  time per next-token vector—the same asymptotic cost as a single parser call—by combining two key ingredients, *lattice parsing* and *algorithmic differentiation*, which we introduce next. Our approach applies to any parser that supports lattice input and whose evaluation consists of semiring operations—mild requirements satisfied by standard algorithms, such as incremental CKY and Earley (see App. E). Moreover, when the base parser is incremental, our next-token algorithm inherits its per-prefix amortization.

**An algorithmic differentiation reduction.** We start by defining the aggregation function:

$$Z_{\mathbf{s}}(\boldsymbol{\theta}) \stackrel{\text{def}}{=} \sum_{a \in \Sigma} \pi_a(\mathbf{s}) \cdot \theta_a \quad (6)$$

where  $\boldsymbol{\theta} \in \mathbb{W}^\Sigma$  is a vector of formal variables, introduced solely to enable differentiation. Because this sum is linear in  $\boldsymbol{\theta}$ , its gradient immediately yields the next-token weight vector:<sup>12</sup>

$$\nabla_{\boldsymbol{\theta}} Z_{\mathbf{s}}(\boldsymbol{\theta}) = \pi(\mathbf{s}) \quad \forall \boldsymbol{\theta} \in \mathbb{W}^\Sigma \quad (7)$$

This reduction is productive because there is an efficient method for computing  $Z_{\mathbf{s}}(\boldsymbol{\theta})$ —namely,

<sup>12</sup>The gradient here is *algebraic*:  $\nabla[x + y] = \nabla[x] + \nabla[y]$  and  $\nabla[x \cdot y] = \nabla[x] \cdot y + x \cdot \nabla[y]$ , requiring only semiring operations (see, e.g., Esparza et al., 2008).

*lattice parsing*, which replaces the input string with a compact automaton encoding all  $|\Sigma|$  one-symbol extensions at once. Moreover, thanks to algorithmic differentiation, we have the following *meta-theorem* that ensures the computation of the gradient  $\nabla_{\theta} Z_s(\theta)$  is equally efficient.

**Theorem 3** (Next-token meta-theorem). *For any algorithm  $q$  that correctly computes  $Z_s(\theta)$  and whose evaluation is a sequence of semiring operations,<sup>13</sup> we can use reverse-mode algorithmic differentiation to derive the gradient algorithm  $\nabla q$  that correctly computes  $\nabla_{\theta} Z_s(\theta) = \pi(s)$  using at most 4 times as many arithmetic operations as  $q$ .*

*Proof.* The proof follows from Eq. (7) and the cheap gradient principle (Griewank and Walther, 2008, Eq. 3.14), which bounds the gradient computation at  $4\times$  the original function’s operation count. Since this bound counts arithmetic operations irrespective of their interpretation, it carries over from  $\mathbb{R}$  to any commutative semiring  $\mathbb{W}$ . ■

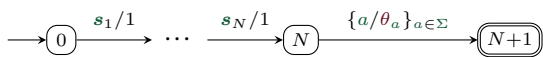
**Lattice parsing.** Next, we show how to efficiently compute  $Z_s(\theta)$  via lattice parsing. A **word lattice** (Hall, 2005) is an acyclic weighted finite-state automaton (WFSA; see App. C.3) representing a weighted set of strings. A **lattice parser**  $q$  is an algorithm that, given a grammar  $\mathcal{G}$  and a word lattice  $\mathcal{L}$ , computes the following product:<sup>14</sup>

$$q(\mathcal{G}, \mathcal{L}) = \sum_{s \in \Sigma^*} \llbracket \mathcal{L} \rrbracket(s) \cdot \llbracket \mathcal{G} \rrbracket(s) \quad (8)$$

That is, the lattice parser computes the inner product of the lattice’s and grammar’s weighted languages. Efficient lattice parsers can be derived by minimally modifying existing parsing algorithms like Earley and CKY (Hall, 2005).

Our approach to computing the *aggregation function*  $Z_s(\theta)$  is based on parsing the next-token lattice defined below.

**Definition 2.** *Let  $s \in \Sigma^*$  and  $\theta \in \mathbb{W}^{\Sigma}$ . The **next-token lattice**  $\mathcal{L}_s(\theta)$  is the WFSA below:*



where the edge labeled  $\{a/\theta_a\}_{a \in \Sigma}$  is shorthand for a set of transitions—one for each symbol  $a \in \Sigma$ .

<sup>13</sup>Control flow (for-loops, if-statements) is permitted provided it does not depend on  $\theta$ —this ensures the computed function is a polynomial in  $\theta$ , so symbolic differentiation applies.

<sup>14</sup>Lattice parsers are routine extensions of ordinary parsers (Hall, 2005, §3.2); ordinary parsing is the special case where the lattice contains the single string  $s$  with weight one.

Since the lattice accepts exactly the strings  $sa$  for  $a \in \Sigma$ , each weighted by  $\theta_a$ , parsing it with the prefix grammar recovers the aggregation function:

$$q(\overline{\mathcal{G}}, \mathcal{L}_s(\theta)) = Z_s(\theta) \quad (9)$$

**Putting it together.** We now combine the three ingredients—the prefix grammar (Def. 1), the next-token lattice (Def. 2), and the next-token meta-theorem (Theorem 3)—into a complete next-token weight vector algorithm. Let  $\mathcal{G}$  be a CFG with prefix grammar  $\overline{\mathcal{G}}$ , and let  $s = s_1 \cdots s_N$  be a string. Suppose that  $q$  is an algorithm whose evaluation consists of a sequence of semiring operations. Then, by the next-token meta-theorem (Theorem 3), there exists a gradient algorithm  $\nabla q$  that computes the next-token weight vector  $\pi(s) = \nabla_{\theta} Z_s(\theta)$  with the same asymptotic complexity as one evaluation of  $q(\overline{\mathcal{G}}, \mathcal{L}_s(\theta))$ . We present two concrete instantiations of this framework in the appendix: one based on CKY (App. F.1) and one based on Earley’s algorithm (App. F.2). The empirical results displayed in Fig. 1—which use the Earley-based instantiation—confirm that computing  $\pi(s)$  is slower than prefix parsing by only a small constant factor. Finally, in many applications, we also need the total weight of a string  $\llbracket \mathcal{G} \rrbracket(s)$  alongside its prefix weight.<sup>15</sup>

## Conclusion

We showed that prefix parsing is just ordinary parsing—applied to the prefix grammar, a compact CFG encoding of the prefix language. Given any CFG, we construct a new CFG whose ordinary parser computes prefix weights of the original. Unlike previous approaches that extend existing parsing algorithms with complex modifications, our reduction provides a unified recipe for prefix parsing. Additionally, using algorithmic differentiation, we provide the fastest known algorithm for computing the next-token weight vector. With our reduction in hand, there is no longer a need to develop bespoke prefix-parsing algorithms. A natural direction for future work is to develop a generalized notion of a *prefix transformation* that extends to other formalisms, such as context-sensitive grammars.

<sup>15</sup>Define  $\mathcal{G} \text{ EOS} \stackrel{\text{def}}{=} \langle \mathcal{N}, \Sigma \cup \{\text{EOS}\}, S', \mathcal{R} \cup \{S' \xrightarrow{1} S \text{ EOS}\} \rangle$  where  $S' \notin \mathcal{N}$ . Since  $\llbracket \mathcal{G} \text{ EOS} \rrbracket(s \text{ EOS}) = \llbracket \mathcal{G} \rrbracket(s)$ , applying the next-token weight vector algorithm to the prefix grammar  $\overline{\mathcal{G}} \text{ EOS}$  yields  $\llbracket \mathcal{G} \rrbracket(s)$  as the EOS component of the vector, alongside all other next-token weights. The transformation is efficient: it adds a single rule, so  $|\mathcal{G} \text{ EOS}| = |\mathcal{G}| + 3$ .

## Limitations

The limitations of our approach are discussed in context throughout the main text; however, we would like to remind that one of the main limitations of our work is that, in the general case, Theorem 2 does not guarantee that a prefix parser has the same asymptotic runtime as the underlying base parser. Additionally, we would like to recall that our prefix parser works with any parser of choice, however we only ran experiments with CKY and Earley's.

## Acknowledgments

We thank Jason Eisner and Franz Nowak for useful discussions. We also thank João Loula, Ben LeBrun, and Leo Du for early adoption of our method (see Loula et al., 2025). We used generative AI to improve our writing and to help debug the code. Every modification introduced by generative AI was carefully reviewed by the authors, which take full responsibility for it.

## References

- Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga, Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. [TensorFlow: A system for large-scale machine learning](#). In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation*.
- J. K. Baker. 1979. [Trainable grammars for speech recognition](#). In *Speech Communication Papers for the Meeting of the Acoustical Society of America*.
- Yehoshua Bar-Hillel, Micha Asher Perles, and Eli Shamir. 1961. [On formal properties of simple phrase-structure grammars](#). *Zeitschrift für Phonetik, Sprachwissenschaft und Kommunikationsforschung*, 14.
- Shalev Ben-David. 2022. [Lecture 10: Closure properties for context-free languages \(CS 360, Winter 2022\)](#).
- Taylor L. Booth and Richard A. Thompson. 1973. [Applying probability measures to abstract languages](#). *IEEE Transactions on Computers*, 22(5).
- James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao Zhang. 2018. [JAX: Composable transformations of Python+NumPy programs](#).
- John Canny, David Hall, and Dan Klein. 2013. [A multi-teraflop constituency parser using GPUs](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Zhiyi Chi. 1999. [Statistical properties of probabilistic context-free grammars](#). *Computational Linguistics*, 25(1).
- John Cocke and Jacob T. Schwartz. 1970. [Programming languages and their compilers: Preliminary notes](#). Technical report, Courant Institute of Mathematical Sciences, New York University.
- Manfred Droste, Werner Kuich, and Heiko Vogler. 2009. *Handbook of Weighted Automata*. Springer.
- Aaron Joseph Dunlop. 2014. *Efficient Latent-Variable Grammars: Learning and Inference*. Ph.D. thesis, Oregon Health & Science University.
- Jay Earley. 1970. [An efficient context-free parsing algorithm](#). *Communications of the Association for Computing Machinery*, 13(2).
- Jason Eisner. 2016. [Inside-outside and forward-backward algorithms are just backprop \(tutorial paper\)](#). In *Proceedings of the Workshop on Structured Prediction for NLP*.
- Jason Eisner and John Blatz. 2007. [Program transformations for optimization of parsing algorithms and other weighted logic programs](#). In *Proceedings of the Conference on Formal Grammar*.
- Jason Eisner, Eric Goldlust, and Noah A. Smith. 2005. [Compiling comp ling: Weighted dynamic programming and the Dyna language](#). In *Proceedings of the Human Language Technology Conference and the Conference on Empirical Methods in Natural Language Processing*.
- Javier Esparza, Stefan Kiefer, and Michael Luttenberger. 2007. [An extension of Newton's method to  \$\omega\$ -continuous semirings](#). In *Developments in Language Theory*.
- Javier Esparza, Stefan Kiefer, and Michael Luttenberger. 2008. [Newton's method for  \$\omega\$ -continuous semirings](#). In *Automata, Languages and Programming*.
- Joshua Goodman. 1999. [Semiring parsing](#). *Computational Linguistics*, 25(4).
- Andreas Griewank and Andrea Walther. 2008. *Evaluating Derivatives*. SIAM.
- John Hale. 2001. [A probabilistic Earley parser as a psycholinguistic model](#). In *Proceedings of the Meeting of the North American Chapter of the Association for Computational Linguistics on Language Technologies*.
- Keith B. Hall. 2005. *Best-first Word-lattice Parsing: Techniques for Integrated Syntactic Language Modeling*. Ph.D. thesis, Brown University.

- Laurent Hascoët and Valérie Pascual. 2013. [The Tapenade automatic differentiation tool: Principles, model, and specification](#). *Association for Computing Machinery Transactions on Mathematical Software*, 39(3).
- Frederick Jelinek and John D. Lafferty. 1991. [Computation of the probability of initial substring generation by stochastic context-free grammars](#). *Computational Linguistics*, 17(3).
- Tadao Kasami. 1965. [An efficient recognition and syntax algorithm for context-free languages](#). Technical Report AF-CRL-65-758, Air Force Cambridge Research Laboratory.
- Jeffrey Kegler. 2023. [Marpa, a practical general parser: the recognizer](#).
- Martin Lange and Hans Leiß. 2009. [To CNF or not to CNF? An efficient yet presentable version of the CYK algorithm](#). *Informatica Didactica*, 8.
- 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](#). In *Proceedings of the International Conference on Learning Representations*.
- Minh-Thang Luong, Michael C. Frank, and Mark Johnson. 2013. [Parsing entire discourses as very long strings: Capturing topic continuity in grounded language learning](#). *Transactions of the Association for Computational Linguistics*, 1.
- Mark-Jan Nederhof. 1999. [The computational complexity of the correct-prefix property for TAGs](#). *Computational Linguistics*, 25(3).
- Mark-Jan Nederhof and Giorgio Satta. 2008. *Probabilistic Parsing*, chapter 7. Springer.
- Franz Nowak and Ryan Cotterell. 2023. [A fast algorithm for computing prefix probabilities](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Andreas Opedal, Anej Svete, Alexandra Butoi, Franz Nowak, and Ryan Cotterell. 2022. [Lecture 7: Grammar transformations](#).
- Andreas Opedal, Ran Zmigrod, Tim Vieira, Ryan Cotterell, and Jason Eisner. 2023. [Efficient semiring-weighted Earley parsing](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Clemente Pasti, Andreas Opedal, Tiago Pimentel, Tim Vieira, Jason Eisner, and Ryan Cotterell. 2023. [On the intersection of context-free and regular languages](#). In *Proceedings of the Conference of the European Chapter of the Association for Computational Linguistics*.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. [Automatic differentiation in PyTorch](#). In *Proceedings of the NeurIPS Autodiff Workshop*.
- Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. 2022. [Synchromesh: Reliable code generation from pre-trained language models](#). In *Proceedings of the International Conference on Learning Representations*.
- Alexander Rush. 2020. [Torch-Struct: Deep structured prediction library](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.
- Yves Schabes and Aravind K. Joshi. 1988. [An Earley-type parsing algorithm for tree adjoining grammars](#). In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*.
- Richard Shin, Christopher Lin, Sam Thomson, Charles Chen, Subhro Roy, Emmanouil Antonios Platanios, Adam Pauls, Dan Klein, Jason Eisner, and Benjamin Van Durme. 2021. [Constrained language models yield few-shot semantic parsers](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Miloš Stanojević and Laurent Sartran. 2023. [SynJax: Structured probability distributions for JAX](#). In *Proceedings of the Conference on Empirical Methods in Natural Language Processing: System Demonstrations*.
- Andreas Stolcke. 1995. [An efficient probabilistic context-free parsing algorithm that computes prefix probabilities](#). *Computational Linguistics*, 21(2).
- Leslie G. Valiant. 1975. [General context-free recognition in less than cubic time](#). *Journal of Computer and System Sciences*, 10(2).
- Youngmin Yi, Chao-Yue Lai, Slav Petrov, and Kurt Keutzer. 2011. [Efficient parallel CKY parsing on GPUs](#). In *Proceedings of the International Conference on Parsing Technologies*.
- Daniel H. Younger. 1967. [Recognition and parsing of context-free languages in time  \$n^3\$](#) . *Information and Control*, 10(2).

## Appendix Contents

<b>A</b>	<b>Notation Glossary</b>	<b>9</b>
<b>B</b>	<b>Additional Experimental Results</b>	<b>10</b>
<b>C</b>	<b>Additional Background</b>	<b>11</b>
C.1	Weighted Context-Free Grammars . . . . .	11
C.2	Grammar Transformations . . . . .	12
C.3	Weighted Finite-State Automata . . . . .	13
<b>D</b>	<b>Deriving the Prefix Grammar by CFG–FST Composition</b>	<b>14</b>
<b>E</b>	<b>Instantiations of Prefix Parsing with the Prefix Grammar</b>	<b>16</b>
E.1	Prefix Parsing with CKY . . . . .	16
E.2	Prefix Parsing with Earley’s Algorithm . . . . .	18
<b>F</b>	<b>Practical Instantiations of the Next-Token Weight Vector Algorithm</b>	<b>20</b>
F.1	Next-Token Weight Vector with CKY . . . . .	20
F.2	Next-Token Weight Vector with Earley’s Algorithm . . . . .	20
<b>G</b>	<b>Deferred Proofs</b>	<b>22</b>
G.1	Proof of Prop. 1 . . . . .	22
G.2	Proof of Prop. 2 . . . . .	24

## A Notation Glossary

Symbol	Description
<i>Strings and languages (§2)</i>	
$\Sigma$	alphabet (finite set of terminal symbols)
$\Sigma^*$	set of all strings over $\Sigma$
$\varepsilon \in \Sigma^*$	empty string
$s, t \in \Sigma^*$	strings
$s, a \in \Sigma$	terminal symbols
$\Sigma^+$	set of nonempty strings over $\Sigma$
$\preceq$	prefix relation ( $s \preceq t$ iff $s$ is a prefix of $t$ )
$\text{EOS} \notin \Sigma$	end-of-string symbol
$\langle \mathbb{W}, +, \cdot, 0, 1 \rangle$	commutative semiring of weights
$\omega: \Sigma^* \rightarrow \mathbb{W}$	weighted language
$\overline{\omega}: \Sigma^* \rightarrow \mathbb{W}$	prefix language of $\omega$
<i>Grammars (§2, App. C.1)</i>	
$\mathcal{G} = \langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle$	weighted context-free grammar (WCFG)
$\mathcal{N}$	set of nonterminal symbols
$S \in \mathcal{N}$	start symbol
$\mathcal{R}$	finite set of weighted rules
$\alpha \in \mathcal{N} \cup \Sigma$	grammar symbol (nonterminal or terminal)
$\alpha \in (\mathcal{N} \cup \Sigma)^*$	rule right-hand side
$r, (X \xrightarrow{w} \alpha) \in \mathcal{R}$	rule; $X$ rewrites to $\alpha$ with weight $w \in \mathbb{W}$
$ \mathcal{G}  \in \mathbb{N}$	grammar size
$\text{ar}(r) \in \mathbb{N}$	arity of rule $r$
$d \in \mathcal{D}$	derivation tree of $\mathcal{G}$
$d \in \mathcal{D}_\alpha \subseteq \mathcal{D}$	$\alpha$ -rooted derivation tree
$\mathcal{D}_\alpha(s) \subseteq \mathcal{D}_\alpha$	set of derivation trees rooted at $\alpha$ with yield $s$
$\sigma(d) \in \Sigma^*$	yield of derivation tree $d$
$w(d) \in \mathbb{W}$	weight of derivation tree $d$
$[[\mathcal{G}]](s) \in \mathbb{W}$	weighted language of $\mathcal{G}$ : total weight of $s$
$\tau(\alpha) \in \mathbb{W}$	total weight of symbol $\alpha \in \Sigma \cup \mathcal{N}$
<i>Prefix grammar (§3)</i>	
$\overline{\mathcal{G}}$	prefix grammar of $\mathcal{G}$
$\overline{X'}$	prime nonterminal (partially matched $X \in \mathcal{N}, X' \notin \mathcal{N}$ )
$\overline{S}$	prefix start symbol ( $\overline{S} \neq S, \overline{S} \notin \mathcal{N}$ )
<i>Parsers and preprocessing (§3, §4)</i>	
$\phi$	grammar preprocessing function (e.g., EnsureCNF; see App. C.2)
$\mathfrak{p}$	abstract parsing algorithm
$\mathfrak{q}$	abstract lattice parser
<i>Next-token algorithms (§4)</i>	
$\pi(s) \in \mathbb{W}^\Sigma$	next-token weight vector
$\pi_a(s) \in \mathbb{W}$	next-token weight for terminal $a$
$\theta \in \mathbb{W}^\Sigma$	formal parameter vector
$\theta_a \in \mathbb{W}$	parameter for terminal $a$
$Z_s(\theta) \in \mathbb{W}$	aggregation function
$\mathcal{L}$	word lattice (acyclic WFSA)
$\mathcal{L}_s(\theta)$	next-token lattice
$z(i, X) \in \mathbb{W}$	forward value
$\nabla z(i, X) \in \mathbb{W}$	backward value
$\beta_k(i, X) \in \mathbb{W}$	inside weight

## B Additional Experimental Results

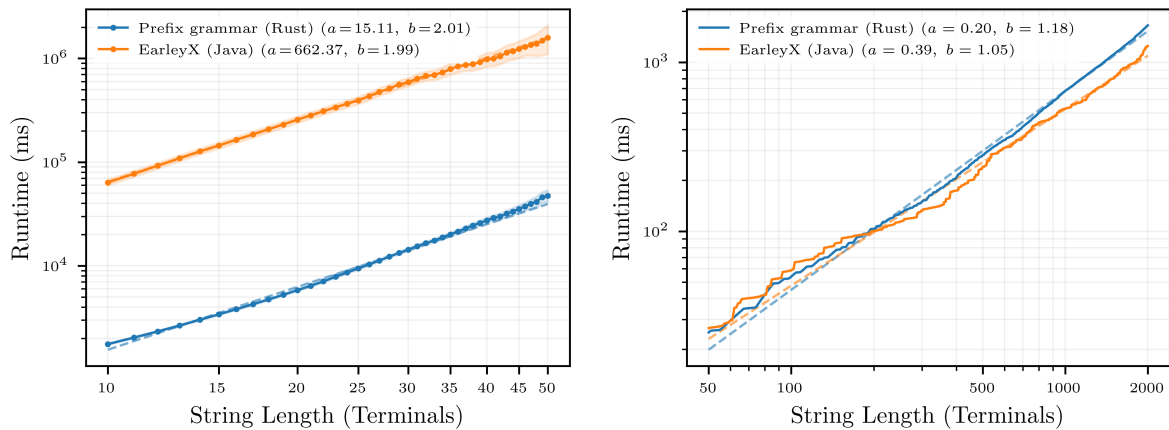


Figure 2: **Left: Our prefix parser vs. EarleyX** (Luong et al., 2013). We compare the two implementations on WSJ 5000 and 500 sentences. Our algorithm uses Earley’s with the prefix grammar (Alg. 3), while EarleyX uses Stolcke’s (1995) algorithm to compute prefix probabilities (itself an adaptation of Earley’s). As in Fig. 1, we fit a power law  $r(N) = aN^b$  via least-squares regression on the log–log scale. The fitted slopes are nearly identical (2.01 vs. 1.99), while the fitted coefficient is substantially better in our implementation. This is consistent with the fact that Stolcke’s algorithm pays an additional multiplicative factor of  $|\mathcal{R}|$ , which our Earley implementation avoids via an optimization at the completion step (Alg. 4). **Right: Sparse grammar comparison.** We repeat the experiment on the *Social Discourse* grammar (Luong et al., 2013, EarleyX), a sparse grammar with 72,712 rules and 233 nonterminals. The graph reports the runtime for each prefix of a sequence of concatenated sentences (socialall.ortho.yld.concat). Both parsers achieve near-linear time on this sparse grammar (fitted exponents 1.05 for EarleyX and 1.18 for our implementation). For this experiment, we used the specialized EarleyParserSparse from Luong et al. (2013), which is highly optimized for sparse grammars.

## C Additional Background

### C.1 Weighted Context-Free Grammars

Let  $\mathcal{G} = \langle \mathcal{N}, \Sigma, \mathcal{S}, \mathcal{R} \rangle$  be a CFG. We develop the weighted language  $\llbracket \mathcal{G} \rrbracket$  (and a few other useful concepts) by first defining (weighted) derivation trees, then summing over trees with a given yield. A **derivation tree**  $d$  is a rooted, ordered tree satisfying: (i) each node is labeled with a symbol in  $\mathcal{N} \cup \Sigma$ , (ii) each nonleaf node is connected to its children by a rule in  $\mathcal{R}$ , and (iii) each leaf is either a terminal or a childless nonterminal. The **set of all derivation trees** is defined recursively as the smallest set  $\mathcal{D}$  satisfying:

$$\mathcal{D} \stackrel{\text{def}}{=} \Sigma \cup \left\{ \begin{array}{c} X \\ \swarrow \quad \searrow \\ \alpha_1 \quad \dots \quad \alpha_K \\ \triangle \quad \quad \quad \triangle \end{array} \mid X \rightarrow \alpha_1 \dots \alpha_K \in \mathcal{R}, \triangle_{\alpha_1}, \dots, \triangle_{\alpha_K} \in \mathcal{D} \right\}$$

Note that the subtrees rooted at the internal nodes of  $d$  are themselves derivation trees. Let  $\rho(d)$  denote the root label of  $d$ , and we write  $\mathcal{D}_\alpha \stackrel{\text{def}}{=} \{d \in \mathcal{D} : \rho(d) = \alpha\}$  to denote the derivations with root  $\alpha \in \Sigma \cup \mathcal{N}$ . The **yield**  $\sigma(d)$  of a tree  $d$  is recursively defined as

$$\bullet \text{ if } d \in \Sigma: \sigma(d) \stackrel{\text{def}}{=} d \tag{10a}$$

$$\bullet \text{ otherwise: } \sigma \left( \begin{array}{c} X \\ \swarrow \quad \searrow \\ \alpha_1 \quad \dots \quad \alpha_K \\ \triangle \quad \quad \quad \triangle \end{array} \right) \stackrel{\text{def}}{=} \sigma \left( \begin{array}{c} \alpha_1 \\ \triangle \end{array} \right) \dots \sigma \left( \begin{array}{c} \alpha_K \\ \triangle \end{array} \right) \tag{10b}$$

Note that the yield  $\sigma(d)$  of a childless node is  $\varepsilon$ . A **derivation forest** is a set of derivation trees; its weight is the sum of the weights of its trees. We write  $\mathcal{D}(s) \stackrel{\text{def}}{=} \{d \in \mathcal{D} : \sigma(d) = s\}$  for the derivation forest with fixed yield  $s$ , and  $\mathcal{D}_\alpha(s) \stackrel{\text{def}}{=} \mathcal{D}_\alpha \cap \mathcal{D}(s)$  for the forest with root  $\alpha$  and yield  $s$ . The **weight**  $w(d)$  of the derivation tree  $d$  is defined as

$$\bullet \text{ if } d \in \Sigma: w(d) \stackrel{\text{def}}{=} 1 \tag{11a}$$

$$\bullet \text{ otherwise: } w \left( \begin{array}{c} X \\ \swarrow \quad \searrow \\ \alpha_1 \quad \dots \quad \alpha_K \\ \triangle \quad \quad \quad \triangle \end{array} \right) \stackrel{\text{def}}{=} w(X \rightarrow \alpha_1 \dots \alpha_K) \cdot w \left( \begin{array}{c} \alpha_1 \\ \triangle \end{array} \right) \dots w \left( \begin{array}{c} \alpha_K \\ \triangle \end{array} \right) \tag{11b}$$

For every  $\alpha \in \Sigma \cup \mathcal{N}$ , we define its **weighted language** as

$$\llbracket \mathcal{G}_\alpha \rrbracket(s) \stackrel{\text{def}}{=} \sum_{d \in \mathcal{D}_\alpha(s)} w(d) \tag{12}$$

Every CFG defines a weighted language  $\llbracket \mathcal{G} \rrbracket$  by taking the language of the root symbol:

$$\llbracket \mathcal{G} \rrbracket(s) \stackrel{\text{def}}{=} \llbracket \mathcal{G}_\mathcal{S} \rrbracket(s) \tag{13}$$

For each symbol  $\alpha \in \Sigma \cup \mathcal{N}$ , we define  $\tau(\alpha)$  as the **total weight** of all  $\alpha$ -rooted derivations:

$$\tau(\alpha) = \sum_{d \in \mathcal{D}_\alpha} w(d) \tag{14}$$

The total weights satisfy the following system of polynomial equations:

$$\tau(\alpha) = \begin{cases} 1 & \text{for } \alpha \in \Sigma \\ \sum_{(\alpha \xrightarrow{w} \alpha_1 \dots \alpha_M) \in \mathcal{R}} w \cdot \tau(\alpha_1) \dots \tau(\alpha_M) & \text{for } \alpha \in \mathcal{N} \end{cases} \tag{15}$$

The total weights are well-defined in any  $\omega$ -continuous semiring (Droste et al., 2009); examples include the nonnegative extended reals. They can be approximated, e.g., by fixed-point iteration or generalized Newton's method (Nederhof and Satta, 2008; Esparza et al., 2007).

## C.2 Grammar Transformations

A grammar transformation is a function  $\phi$  that maps a CFG  $\mathcal{G}$  to another CFG  $\phi(\mathcal{G})$ . We say that  $\phi$  is *semantics preserving* if  $\llbracket \mathcal{G} \rrbracket = \llbracket \phi(\mathcal{G}) \rrbracket$ . Many transformations used to preprocess grammars are semantics preserving, including **dead-rule elimination**, **binarization**, **nullary removal**, **unary removal**, and **terminal separation**. The purpose of most of these transformations is self-explanatory to a reader familiar with CFGs; we refer the reader to [Opedal et al. \(2022\)](#) for detailed descriptions. All transformations discussed in this section are included in our [code release](#). Converting a grammar into CNF is itself a (composed) transformation `EnsureCNF`: apply `EnsureCTF`, then `EnsureNullaryFree`, then `EnsureUnaryFree`, then `EnsureTerminalSep`.<sup>16</sup> In general, a grammar transformation  $\phi$  can change both the size  $|\phi(\mathcal{G})|$  of the transformed grammar and the number of nonterminals  $|\mathcal{N}_\phi|$ . We bound both for the transformations above:

- Dead rule elimination:

$$|\phi(\mathcal{G})| \leq |\mathcal{G}|, \quad |\mathcal{N}_\phi| \leq |\mathcal{N}| \quad (16)$$

- `EnsureCTF` (binarization):

$$|\phi(\mathcal{G})| \leq 3|\mathcal{G}|, \quad |\mathcal{N}_\phi| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| \quad (17)$$

The  $|\mathcal{G}| - |\mathcal{R}|$  term accounts for intermediate nonterminals introduced when splitting rules of arity  $> 2$ .

- `EnsureNullaryFree` (nullary removal), provided  $\mathcal{G}$  is in CTF:

$$|\phi(\mathcal{G})| \leq \frac{7}{3}|\mathcal{G}| + 3, \quad |\mathcal{N}_\phi| \leq |\mathcal{N}| + 1 \quad (18)$$

- `EnsureUnaryFree` (unary removal):

$$|\phi(\mathcal{G})| \leq |\mathcal{N}| \cdot |\mathcal{G}|, \quad |\mathcal{N}_\phi| \leq |\mathcal{N}| \quad (19)$$

- `EnsureUnaryCycleFree` (unary cycle removal) from [Opedal et al. \(2023\)](#):

$$|\phi(\mathcal{G})| \leq |\mathcal{G}| + 2|\mathcal{N}|^2, \quad |\mathcal{N}_\phi| \leq 2|\mathcal{N}| \quad (20)$$

The  $|\mathcal{N}|^2$  term bounds the closure of the unary-rule graph; it is tight in the worst case (a single strongly connected component) and often far smaller in practice.

- `EnsureTerminalSep` (terminal separation):

$$|\phi(\mathcal{G})| \leq |\mathcal{G}| + 2|\Sigma|, \quad |\mathcal{N}_\phi| \leq |\mathcal{N}| + |\Sigma| \quad (21)$$

- `EnsureCNF` (composing the above in the order `EnsureCTF`, `EnsureNullaryFree`, `EnsureUnaryFree`, `EnsureTerminalSep`). Tracing the pipeline:

1. `EnsureCTF`:  $|G_1| \leq 3|\mathcal{G}|$ ,  $|N_1| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|$  (Eq. (17))

2. `EnsureNullaryFree`:  $|G_2| \leq \frac{7}{3}|G_1| + 3 \leq 7|\mathcal{G}| + 3$ ,  $|N_2| \leq |N_1| + 1$  (Eq. (18))

3. `EnsureUnaryFree`:  $|G_3| \leq |N_2| \cdot |G_2| \leq (|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| + 1)(7|\mathcal{G}| + 3)$ ,  $|N_3| \leq |N_2|$  (Eq. (19))

4. `EnsureTerminalSep`:  $|G_4| \leq |G_3| + 2|\Sigma|$ ,  $|N_4| \leq |N_3| + |\Sigma|$  (Eq. (21))

Summing yields:

$$\begin{aligned} |\phi(\mathcal{G})| &\leq (|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| + 1)(7|\mathcal{G}| + 3) + 2|\Sigma| \\ |\mathcal{N}_\phi| &\leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| + 1 + |\Sigma| \end{aligned} \quad (22)$$

The dominant factor is unary removal; the  $|\mathcal{G}|$  term in the first factor arises because binarization introduces up to  $|\mathcal{G}| - |\mathcal{R}|$  intermediate nonterminals. See [Lange and Leiß \(2009\)](#) for a detailed analysis of CNF conversion orderings and their worst-case blowup.

<sup>16</sup>Note that the order in which these transformations are applied matters. For example, applying `EnsureNullaryFree` before `EnsureCTF` can cause exponential blowup, as a rule of arity  $K$  may spawn up to  $2^K$  variants ([Lange and Leiß, 2009](#)). Placing terminal separation last prevents subsequent multiplicative transformations from amplifying its additive cost.

### C.3 Weighted Finite-State Automata

A WFSA  $\mathcal{A}$  is a tuple  $\langle \Sigma, Q, \delta, \lambda, \rho \rangle$ , where  $\Sigma$  is an alphabet,  $Q$  is a finite set of states,  $\delta$  is a set of edges (each having form  $q \xrightarrow{a/w} p$ , with  $p, q \in Q$ ,  $w \in \mathbb{W}$  and  $a \in \Sigma$ ),  $\lambda: Q \rightarrow \mathbb{W}$  and  $\rho: Q \rightarrow \mathbb{W}$  are respectively the initial and final weight functions. A **path**  $\pi$  is a sequence of consecutive edges

$$q_0 \xrightarrow{a_1/w_1} q_1 \xrightarrow{a_2/w_2} q_2 \cdots q_{N-1} \xrightarrow{a_N/w_N} q_N \quad (23)$$

the path's weight is defined as  $\mathbf{w}(\pi) = \lambda(q_0) \cdot w_1 \cdot w_2 \cdots w_N \cdot \rho(q_N)$ , the path's yield is  $\sigma(\pi) = a_1 \cdots a_N$ . We denote with  $\Pi_{\mathcal{A}}(s)$  the set of paths with yield  $s$ . A WFSA defines the weighted language:

$$\llbracket \mathcal{A} \rrbracket(s) = \sum_{\pi \in \Pi_{\mathcal{A}}(s)} \mathbf{w}(\pi) \quad (24)$$

## D Deriving the Prefix Grammar by CFG–FST Composition

This section explains how we initially *derived* the prefix grammar<sup>17</sup> by composition with a weighted finite-state transducer. Although Def. 1 is a more succinct encoding of the prefix grammar, we believe the derivation as a composition is of independent interest. We begin by reviewing weighted relations and weighted finite-state transducers, then introduce the **prefix transducer**  $\Delta$ .

Let  $\Sigma$  and  $\Xi$  be alphabets. A **weighted relation** between the strings  $\Sigma^*$  and  $\Xi^*$  is a function  $t: \Sigma^* \times \Xi^* \rightarrow \mathbb{W}$ . We define the **composition** of  $t$  with the language  $\omega: \Sigma^* \rightarrow \mathbb{W}$ , as the language  $\omega \circ t: \Xi^* \rightarrow \mathbb{W}$ , defined as follows:

$$[\omega \circ t](t) \stackrel{\text{def}}{=} \sum_{s \in \Sigma^*} \omega(s) t(s, t) \quad (25)$$

**Definition 3.** A **weighted finite-state transducer (WFST)** is a tuple  $\mathcal{T} = \langle \Sigma, \Xi, Q, \delta, \lambda, \rho \rangle$  where

- $\Sigma$  is the **input alphabet**
- $\Xi$  is the **output alphabet**
- $Q$  is a **finite set of states**
- $\delta$  is a **set of transitions** where each **transition** is of the form  $(q \xrightarrow{a:b/w} q')$  with  $q, q' \in Q$ ,  $a \in \Sigma \cup \{\varepsilon\}$ ,  $b \in \Xi \cup \{\varepsilon\}$ , and  $w \in \mathbb{W}$
- $\lambda: Q \rightarrow \mathbb{W}$  is the **initial weight function**
- $\rho: Q \rightarrow \mathbb{W}$  is the **final weight function**.

A **path**  $\pi$  is a sequence of transitions of the form:

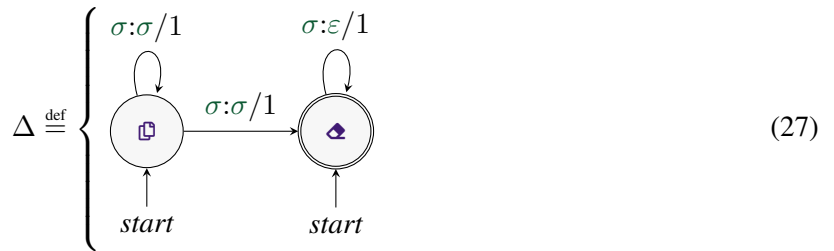
$$q_0 \xrightarrow{a_1:b_1/w_1} q_1 \xrightarrow{a_2:b_2/w_2} \cdots q_{N-1} \xrightarrow{a_N:b_N/w_N} q_N,$$

where  $N \geq 0$  is its length,  $a_1 \cdots a_N$  is its **input yield**,  $b_1 \cdots b_N$  is its **output yield**, and  $w(\pi) \stackrel{\text{def}}{=} \lambda(q_0)w_1 \cdots w_N\rho(q_N)$  is its **weight**. We denote by  $\Pi(s, t)$  the set of paths with input yield  $s$  and output yield  $t$ . Every transducer  $\mathcal{T}$  defines a weighted relation as follows:

$$\llbracket \mathcal{T} \rrbracket(s, t) \stackrel{\text{def}}{=} \sum_{\pi \in \Pi(s, t)} w(\pi) \quad (26)$$

The following WFST  $\Delta$  implements the weighted relation  $\llbracket \Delta \rrbracket(s, t) = \mathbb{1}\{t \preceq s\}$ .

**Definition 4.** The **prefix transducer**  $\Delta$  for a given alphabet  $\Sigma$  is illustrated below:<sup>18</sup>



More precisely, the weighted transducer is defined by the tuple  $\langle \Sigma, \Sigma, Q, \delta, \lambda, \rho \rangle$ , where

- $\Sigma$  is both the input and output alphabet

<sup>17</sup>Beyond CFGs, our observation that prefix parsing is just composition with a finite-state transducer (App. D) suggests that any formalism that efficiently supports such composition is amenable to our general approach. For example, we can, in principle, use composition to derive prefix parsers for tree-adjoining grammars (see, e.g., Nederhof, 1999). Historically, Schabes and Joshi (1988) originally presented an  $\mathcal{O}(N^9)$ -time algorithm to perform prefix parsing on tree-adjoining grammars. This runtime was later improved to  $\mathcal{O}(N^6)$  by Nederhof (1999). Generalizing our reduction to tree-adjoining grammars would obviate the need for such bespoke algorithms.

<sup>18</sup>Note that any FST that encodes  $\llbracket \Delta \rrbracket(s, t) = \mathbb{1}\{t \preceq s\}$  would work. Of course, a smaller machine will generally lead to a smaller increase in the size of the transformed grammar. We believe this transducer is among the smallest possible, as it appears that at least two states are necessary.

- the set of states is  $Q = \{\boxplus, \boxminus\}$ ; we call  $\boxplus$  and  $\boxminus$  the **copy** and **erase** states, respectively
- the transitions  $\delta$  are as shown in the picture
- the initial weights  $\lambda$  are  $\lambda(\boxplus) = 1$  and  $\lambda(\boxminus) = 1$
- the final weights are  $\rho(\boxplus) = 0$  and  $\rho(\boxminus) = 1$

Note that in the picture, the start arrow and the double circle mark states with nonzero initial and final weights, respectively. Additionally, a transition labeled  $\sigma$  is shorthand for a set of transitions—one for each symbol in the alphabet.

Pasti et al. (2023) provide a general construction for the composition of a CFG  $\mathcal{G}$  and a WFST  $\mathcal{T}$ , denoted  $\mathcal{G} \circ \mathcal{T}$ , which is itself a CFG, and is *correct* in the sense that it denotes the composition between the corresponding weighted relation and language, i.e.,  $\llbracket \mathcal{G} \circ \mathcal{T} \rrbracket = \llbracket \mathcal{G} \rrbracket \circ \llbracket \mathcal{T} \rrbracket$ .

Equivalently, we can construct a grammar for the prefix language via composition with  $\Delta$ :  $\llbracket \mathcal{G} \circ \Delta \rrbracket = \overrightarrow{\llbracket \mathcal{G} \rrbracket}$  by Pasti et al. (2023). The composition construction is spelled out below.

**Definition 5.** The **composition-based prefix grammar construction** works as follows. Given a grammar  $\mathcal{G} = \langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle$ , the composition  $\mathcal{G} \circ \Delta = \langle \{s, \alpha, s' \mid \alpha \in \mathcal{N} \cup \Sigma, s, s' \in \{\boxplus, \boxminus\}\}, \Sigma, \overrightarrow{S}, \mathcal{R}' \rangle$  where  $\mathcal{R}'$  is defined by the rules below.

$$\overrightarrow{S} \xrightarrow{1} \langle \boxplus, S, \boxminus \rangle \quad (28a)$$

$$\overrightarrow{S} \xrightarrow{1} \langle \boxminus, S, \boxminus \rangle \quad (28b)$$

$$\langle \boxplus, X, \boxplus \rangle \xrightarrow{w} \langle \boxplus, \alpha_1, \boxplus \rangle \cdots \langle \boxplus, \alpha_K, \boxplus \rangle \quad X \xrightarrow{w} \alpha_1 \cdots \alpha_K \in \mathcal{R} \quad (28c)$$

$$\langle \boxplus, a, \boxplus \rangle \xrightarrow{1} a \quad a \in \Sigma \quad (28d)$$

$$\langle \boxplus, X, \boxminus \rangle \xrightarrow{w} \langle \boxplus, \alpha_1, \boxplus \rangle \cdots \langle \boxplus, \alpha_{k-1}, \boxplus \rangle \langle \boxplus, \alpha_k, \boxminus \rangle \langle \boxminus, \alpha_{k+1}, \boxminus \rangle \cdots \langle \boxminus, \alpha_K, \boxminus \rangle \quad (28e)$$

$$X \xrightarrow{w} \alpha_1 \cdots \alpha_K \in \mathcal{R}, k \in 1, \dots, K$$

$$\langle \boxplus, a, \boxminus \rangle \xrightarrow{1} a \quad a \in \Sigma \quad (28f)$$

$$\langle \boxminus, X, \boxminus \rangle \xrightarrow{w} \langle \boxminus, \alpha_1, \boxminus \rangle \cdots \langle \boxminus, \alpha_K, \boxminus \rangle \quad X \xrightarrow{w} \alpha_1 \cdots \alpha_K \in \mathcal{R} \quad (28g)$$

$$\langle \boxminus, a, \boxminus \rangle \xrightarrow{1} \varepsilon \quad a \in \Sigma \quad (28h)$$

The construction is based directly on Pasti et al.'s (2023) composition construction for  $\mathcal{G} \circ \Delta$  where  $\Delta$  is given in Def. 4.<sup>19</sup> The proposition below establishes that the composition-based construction also correctly encodes the prefix language.

**Proposition 3.** Let  $\mathcal{G}$  be a CFG. Then,  $\mathcal{G} \circ \Delta$  correctly encodes the prefix language of  $\mathcal{G}$ , i.e.,  $\llbracket \mathcal{G} \circ \Delta \rrbracket = \overrightarrow{\llbracket \mathcal{G} \rrbracket}$ .

*Proof.* Let  $\mathcal{G} = \langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle$ .

$$\llbracket \mathcal{G} \circ \Delta \rrbracket(s) = \sum_{z \in \Sigma^*} \llbracket \mathcal{G} \rrbracket(z) \cdot \llbracket \Delta \rrbracket(z, s) \quad \triangleright \text{By Corollary 1 of Pasti et al. (2023)} \quad (29a)$$

$$= \sum_{z \in \Sigma^*} \llbracket \mathcal{G} \rrbracket(z) \cdot \mathbb{1}\{s \preceq z\} \quad \triangleright \text{By the definition of the prefix transducer (Def. 4)} \quad (29b)$$

$$= \overrightarrow{\llbracket \mathcal{G} \rrbracket}(s) \quad \triangleright \text{By the definition of prefix language (§2)} \quad (29c)$$

■

<sup>19</sup>Readers familiar with the composition algorithm may recognize that Def. 5 is a specialized application of the composition algorithm for the prefix transducer. We have hard-coded the prefix transducer's states and transitions. We have omitted nonterminals of the form  $\langle \boxminus, \_, \boxplus \rangle$ , as they are always nongenerating and, thus, can be removed without altering the weighted language of the prefix grammar. Additionally, we have exploited knowledge of the path structure in  $\Delta$ ; specifically, each path with nonzero weight has the form:  $\boxplus \xrightarrow{\sigma_1: \sigma_1/1} \boxplus \cdots \boxplus \xrightarrow{\sigma_{m-1}: \sigma_{m-1}/1} \boxplus \xrightarrow{\sigma_m: \sigma_m/1} \boxminus \xrightarrow{\sigma_{m+1}: \varepsilon/1} \boxminus \cdots \boxminus \xrightarrow{\sigma_M: \varepsilon/1} \boxminus$ . We exploit this knowledge to reduce the number of useless rules; we can see this structure in the right-hand side of Eq. (28e).

## E Instantiations of Prefix Parsing with the Prefix Grammar

In §3, we presented prefix parsing as a general transformation that allows us to turn any existing parsing algorithm into a prefix parsing one; we shall now consider two specific instantiations of this approach, one with CKY and the other with Earley as the underlying parser.

### E.1 Prefix Parsing with CKY

We first instantiate the prefix-grammar approach with CKY. Our incremental formulation, CKY (Alg. 2), runs in  $\mathcal{O}(|\text{EnsureCNF}(\mathcal{G})|N^3)$ , where  $\text{EnsureCNF}(\mathcal{G})$  is the CNF conversion of  $\mathcal{G}$ .  $\text{PrefixCKY}$  (Alg. 1) applies the prefix grammar transformation and then calls CKY. Note that the runtime of CKY follows the form of Theorem 2. We next analyze the blowup introduced by the prefix grammar transformation and all the preprocessing steps required in  $\text{PrefixCKY}$ .

**Proposition 4.** *Let  $\mathcal{G}$  be a CFG. Then  $|\text{EnsureCNF}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G})))| = \mathcal{O}(|\mathcal{G}|^2 + |\mathcal{N}||\mathcal{G}|)$ .*

*Proof.* We trace the grammar size and the nonterminal count through the preprocessing pipeline, referencing the bounds established in App. C.1. We write  $G_i$  for the grammar after step  $i$ , and  $N_i$  for its nonterminal set. Given the input grammar  $\mathcal{G} = \langle \mathcal{N}, \Sigma, \mathcal{S}, \mathcal{R} \rangle$ , we have:

1.  $\text{EnsureCTF}$ :  $|G'_1| \leq 3|\mathcal{G}|$ ,  $|N'_1| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|$  (Eq. (17))
2.  $\text{Prefix}$ :  $|G'_2| \leq \frac{8}{3}|G'_1| + 3 \leq 8|\mathcal{G}| + 3$ ,  $|N'_2| = 2|N'_1| + 1$  (Prop. 2, Def. 1)

The nonterminal count doubles because the  $\text{prefix}$  grammar introduces a prime copy  $X'$  of each nonterminal  $X \in N'_1$ , plus one prefix start symbol  $\bar{S}$ . Moreover, because  $\text{Prefix}$  preserves arity  $\leq 2$ ,  $G'_2$  is already binary; invoking Eq. (22) on  $G'_2$  would pay for an  $\text{EnsureCTF}$  step that is in fact a no-op, incurring a spurious  $3\times$  factor via Eq. (17). We therefore trace the remaining  $\text{EnsureCNF}$  steps directly.

3.  $\text{EnsureNullaryFree}$ :  $|G'_3| \leq \frac{7}{3}|G'_2| + 3 \leq \frac{7}{3}(8|\mathcal{G}| + 3) + 3 \leq 19|\mathcal{G}| + 10$ ,  $|N'_3| \leq |N'_2| + 1$  (Eq. (18))
4.  $\text{EnsureUnaryFree}$ :  $|G'_4| \leq |N'_3| \cdot |G'_3| \leq (2|N'_1| + 2)(19|\mathcal{G}| + 10)$ ,  $|N'_4| \leq |N'_3|$  (Eq. (19))
5.  $\text{EnsureTerminalSep}$ :  $|G'_5| \leq |G'_4| + 2|\Sigma|$ ,  $|N'_5| \leq |N'_4| + |\Sigma|$  (Eq. (21))

So  $|\text{EnsureCNF}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G})))| = |G'_5| \leq (2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|) + 2)(19|\mathcal{G}| + 10) + 2|\Sigma|$ . Expanding and dropping lower-order terms:  $|G'_5| = \mathcal{O}(|\mathcal{G}|^2 + |\mathcal{N}||\mathcal{G}|)$ , matching the proposition. ■

We note that this is a worst-case bound; in practice, the blowup is much more contained. In particular, in our experiments with PCFGs, when comparing the size of the preprocessed prefix grammar  $\text{EnsureCNF}(\bar{\mathcal{G}})$  to that of the preprocessed original grammar  $\text{EnsureCNF}(\mathcal{G})$ , the size ratio is only a small multiplicative factor (Tab. 1). Since CKY’s runtime is linear in the preprocessed grammar size, this small factor translates directly to the runtime overhead for  $\text{PrefixCKY}$ .

---

**Algorithm 1** Prefix parsing with the prefix grammar transformation and CKY (Alg. 2).

---

1. **function**  $\text{PrefixCKY}(\mathcal{G}, s_1 \dots s_N)$
  2.    $\triangleright$ Apply the prefix-grammar transformation (should be cached for efficiency), and parse as usual.
  3. **return**  $\text{CKY}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G})), s_1 \dots s_N)$
- 

Table 1: Grammar size of different PCFGs used in our experiments.  $\mathcal{G}_1$  is the grammar after CKY’s preprocessing;  $\mathcal{G}_2$  is the grammar after  $\text{PrefixCKY}$ ’s preprocessing. The ratio column reports  $|\mathcal{G}_2|/|\mathcal{G}_1|$ . The grammars reported in the table are WSJ 500 (obtained from the first 500 sentences of the Wall Street Journal Corpus) and “*Social Discourse*”, both obtained from (Luong et al., 2013). We were not able to include the larger WSJ 5000 because our machine had insufficient memory for the unary removal operation of  $\text{EnsureCNF}$ .

Grammar	$ \mathcal{G} ;  \mathcal{N} $	$ \mathcal{G}_1 ;  \mathcal{N}_1 $	$ \mathcal{G}_2 ;  \mathcal{N}_2 $	Ratio $ \mathcal{G}_2 / \mathcal{G}_1 $
WSJ 500	12,573; 69	73,241; 1,766	235,459; 1,814	$3.22\times$
Social Discourse [Sparse]	72,712; 233	211,015; 233	357,066; 306	$1.69\times$

---

**Algorithm 2** An incremental, left-to-right formulation of CKY that supports memoization across prefixes. CKY returns the full column of inside weights: the entry  $\beta_j(i, X)$  gives the total weight of all derivations of  $X$  with yield  $s_{i+1} \cdots s_j$ . The goal weight  $\llbracket \mathcal{G} \rrbracket(s_1 \cdots s_N)$  is  $\beta_N(0, S)$ .

---

```

1. function CKY( $\mathcal{G}, s_1 \cdots s_N$ )
2.    $\mathcal{G} \leftarrow \text{EnsureCNF}(\mathcal{G})$                                  $\triangleright$ CNF conversion; if needed (should be cached for efficiency)
3.    $\beta_N \leftarrow \text{defaultdict}(0)$                                 $\triangleright$ Initialize new column
4.   if  $N = 0$ :                                                $\triangleright$ Base case: empty string.
5.     for  $(S \xrightarrow{w} \varepsilon) \in \mathcal{R}$ :
6.        $\beta_N(0, S) += w$ 
7.     return  $\beta_N(0, S), \langle \beta_N \rangle$                                 $\triangleright$ Below: recurse on prefix
8.    $\leftarrow, \langle \beta_0, \dots, \beta_{N-1} \rangle \leftarrow \text{CKY}(\mathcal{G}, s_1 \cdots s_{N-1})$ 
9.   for  $(X \xrightarrow{w} s_N) \in \mathcal{R}$ :                                        $\triangleright$ Preterminal rules  $s_N$ 
10.     $\beta_N(N-1, X) += w$ 
11.   for  $i$  in  $N-1 \dots 0$ :                                        $\triangleright$ Iterate over the start point
12.     for  $j$  in  $i+1 \dots N-1$ :                                        $\triangleright$ Iterate over the split point
13.       for  $(X \xrightarrow{w} YZ) \in \mathcal{R}$ :                                        $\triangleright$ Binary rules
14.         $\beta_N(i, X) += w \cdot \beta_j(i, Y) \cdot \beta_N(j, Z)$ 
15.   return  $\beta_N(0, S), \langle \beta_0, \dots, \beta_{N-1}, \beta_N \rangle$ 

```

---

## E.2 Prefix Parsing with Earley’s Algorithm

We now instantiate the prefix-grammar approach using the weighted Earley’s algorithm (Earley, 1970). Analogously to PrefixCKY, PrefixEarley (Alg. 3) applies the prefix grammar transformation and then parses with Earley’s. Importantly, the experimental results (Fig. 1) show that prefix parsing is slower than ordinary parsing by a multiplicative factor of  $\approx 3$ ; this matches the prediction from the relative grammar-size blow-up (Tab. 2).

---

**Algorithm 3** Prefix parsing with the prefix grammar transformation and Earley (Alg. 4).

---

1. **function** PrefixEarley( $\mathcal{G}, s_1 \dots s_N$ )
  2.    $\triangleright$ Apply the prefix-grammar transformation (should be cached for efficiency), and parse as usual.
  3.   **return** Earley(Prefix(EnsureCTF( $\mathcal{G}$ )),  $s_1 \dots s_N$ )
- 

**An optimized version of Earley’s algorithm.** We present an optimized, incremental version of Earley’s algorithm (Alg. 4), which we use as a base for our prefix parsing and next-token weight vector algorithms. Assuming  $\mathcal{G}$  is nullary-free and unary-cycle-free, our algorithm runs in  $\mathcal{O}(|\mathcal{G}|N^3)$ —a factor of  $|\mathcal{R}|$  faster than the traditional dotted-item Earley formulation.<sup>20</sup> To achieve this improved runtime, we maintain a list of *items* of the form  $\langle i, X/\alpha \rangle$ , where  $i$  is the starting position and  $X/\alpha$  is read as “ $X$  missing  $\alpha$  to complete”. The algorithm proceeds through three operations: PREDICT generates new items from grammar rules.<sup>21</sup> SCAN advances items over terminal symbols, and ATTACH—traditionally also known as *complete*—attaches completed items to items awaiting them. The data structure  $W_k$  indexes items at position  $k$  by the symbol they are waiting for, enabling efficient lookup during the ATTACH phase.

**Grammar blowup: bounds and empirical evidence.** Following Theorem 2, we bound  $|\phi(\vec{\mathcal{G}})|$ , where  $\phi$  is the preprocessing pipeline required for PrefixEarley.

**Proposition 5.** *Let  $\mathcal{G}$  be a CFG. Then,*

$$|\text{EnsureUnaryCycleFree}(\text{EnsureNullaryFree}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G}))))| = \mathcal{O}(|\mathcal{G}|^2 + |\mathcal{N}|^2).^{22}$$

*Proof.* We analyze  $\phi(\mathcal{G})$  first, then extend to  $\phi(\vec{\mathcal{G}})$ . We trace the grammar size and the nonterminal count through the preprocessing pipeline, referencing the bounds established in App. C.1. We write  $G_i$  for the grammar after step  $i$ , and  $N_i$  for its nonterminal set. We apply EnsureCTF (required before EnsureNullaryFree for efficiency), EnsureNullaryFree, and EnsureUnaryCycleFree in sequence. Given the input grammar  $\mathcal{G} = \langle \mathcal{N}, \Sigma, \mathcal{S}, \mathcal{R} \rangle$ , we have

$$1. \text{ EnsureCTF: } |G_1| \leq 3|\mathcal{G}|, \quad |N_1| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| \quad (\text{Eq. (17)})$$

$$2. \text{ EnsureNullaryFree: } |G_2| \leq 7|\mathcal{G}| + 3, \quad |N_2| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| + 1 \quad (\text{Eq. (18)})$$

$$3. \text{ EnsureUnaryCycleFree: } |G_3| \leq 7|\mathcal{G}| + 3 + 2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|)^2$$

Thus,  $|\phi(\mathcal{G})| \leq 7|\mathcal{G}| + 3 + 2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|)^2$ . For the prefix grammar, the preprocessing pipeline becomes:

$$1. \text{ EnsureCTF: } |G_1| \leq 3|\mathcal{G}|, \quad |N_1| \leq |\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}| \quad (\text{Eq. (17)})$$

$$2. \text{ Prefix: } |G_2| \leq \frac{8}{3}|G_1| + 3 \leq 8|\mathcal{G}| + 3, \quad |N_2| = 2|N_1| + 1 \quad (\text{Prop. 2, Def. 1})$$

Since Prefix preserves arity  $\leq 2$ ,  $G_2$  is already binary and binarization is a no-op.

$$3. \text{ EnsureNullaryFree: } |G_3| \leq \frac{7}{3}|G_2| + 3 \leq \frac{7}{3}(8|\mathcal{G}| + 3) + 3 \leq 19|\mathcal{G}| + 10, \quad |N_3| \leq |N_2| + 1 = 2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|) + 2 \quad (\text{Eq. (18)})$$

$$4. \text{ EnsureUnaryCycleFree: } |G_4| \leq |G_3| + 2|N_3|^2 \leq 19|\mathcal{G}| + 10 + 2(2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|) + 2)^2 \quad (\text{Eq. (20)})$$

Thus,  $|\phi(\vec{\mathcal{G}})| \leq 19|\mathcal{G}| + 10 + 2(2(|\mathcal{N}| + |\mathcal{G}| - |\mathcal{R}|) + 2)^2$ . Expanding and dropping lower-order terms:  $|\phi(\vec{\mathcal{G}})| = \mathcal{O}(|\mathcal{G}|^2 + |\mathcal{N}|^2)$ , matching the proposition.  $\blacksquare$

<sup>20</sup>Our items record only the unmatched suffix  $\alpha$ , rather than the traditional  $X \rightarrow \gamma \bullet \alpha$  dotted-item form. This collapses many rules into one item, a simpler alternative to Opedal et al.’s (2023) approach, which introduces additional item types.

<sup>21</sup>Our implementation uses the left-corner relation to prune the PREDICT step, skipping rules whose left-hand sides are not needed by any active item (Stolcke, 1995).

<sup>22</sup>The  $|\mathcal{G}|^2$  term arises when  $\mathcal{G}$  has high-arity rules: EnsureCTF introduces up to  $|\mathcal{G}| - |\mathcal{R}|$  new nonterminals, whose primes may participate in unary cycles. For grammars already in canonical two-form, the bound tightens to  $\mathcal{O}(|\mathcal{G}| + |\mathcal{N}|^2)$ . This worst case is loose in practice (see Tab. 2).

In practice, this bound is far from tight. On the grammars we experimented with, the preprocessed prefix grammar is only a small multiplicative factor larger than the ordinary grammar preprocessed for Earley (Tab. 2). Given our optimized Earley implementation’s runtime bound, this small factor translates directly to the ratio of prefix parsing to parsing runtime, as confirmed empirically by our experiments (Fig. 1).

**Algorithm 4** Earley’s Algorithm takes the grammar  $\mathcal{G} = \langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle$  and a string  $s = s_1 \dots s_k \in \Sigma^*$ . It returns the *inside weight columns*  $\langle \beta_0, \dots, \beta_k \rangle$  together with the waiting-for dictionaries  $\langle W_0, \dots, W_k \rangle$ ; the total string weight is  $\llbracket \mathcal{G} \rrbracket(s_1 \dots s_k) = \beta_k(0, S)$ . The notation  $X/\alpha$  is read as  $X$  missing  $\alpha$  to complete. The data structure  $W_k[\alpha]$  maintains a dictionary of items  $\langle i, X/\alpha_1 \alpha_{[1:i]} \rangle$  that are awaiting the symbol  $\alpha$  in order to move towards completion. The queue  $Q$  is a priority queue that prioritizes items with the smallest span  $k - i$ , where  $k$  is the current column, and  $i$  is the item’s starting position.

```

1. def Earley( $\mathcal{G}, s_1 \dots s_k$ ): ▷The Earley function should be memoized for efficiency.
2.  $\mathcal{G} \leftarrow$  EnsureUnaryCycleFree(EnsureNullaryFree( $\mathcal{G}$ )) ▷Ensure no nullary rules, or unary rule cycles; memoize
3.  $\beta_k \leftarrow$  defaultdict(0) ▷Initialize inside weight column
4.  $W_k \leftarrow$  defaultdict(set) ▷Initialize the waiting-for dictionary
5.  $Q \leftarrow$  priority_queue() ▷This queue prioritizes items with the shortest span.
6. if  $k = 0$ : ▷Base cases
7.   for  $(X \xrightarrow{w} \alpha) \in \mathcal{R}$ :
8.      $\beta_0(0, X/\alpha) += w$ 
9.   return  $\langle \beta_0 \rangle, \langle W_0 \rangle$ 
10.  $\langle \beta_0, \dots, \beta_{k-1} \rangle, \langle W_0, \dots, W_{k-1} \rangle \leftarrow$  Earley( $\mathcal{G}, s_1 \dots s_{k-1}$ ) ▷Recurse on prefix
11. for  $\langle i, X/s_k \alpha \rangle \in W_{k-1}[s_k]$ : ▷SCAN
12.    $\beta_k(i, X/\alpha) += \beta_{k-1}(i, X/s_k \alpha)$ 
13. while  $Q$ : ▷ATTACH
14.    $\langle j, Y/\varepsilon \rangle \leftarrow Q.pop()$  ▷Pop completed item Y together with the index j where item began
15.   for  $\langle i, X/Y\alpha \rangle \in W_j[Y]$ : ▷Iterate through items ending at j that are waiting for a Y
16.      $\beta_k(i, X/\alpha) += \beta_j(i, X/Y\alpha) \cdot \beta_k(j, Y/\varepsilon)$  ▷Attach the weights and update
17. for  $(X \xrightarrow{w} \alpha) \in \mathcal{R}$ : ▷PREDICT
18.    $\beta_k(k, X/\alpha) += w$ 
19. return  $\langle \beta_0, \dots, \beta_{k-1}, \beta_k \rangle, \langle W_0, \dots, W_{k-1}, W_k \rangle$ 
20. ▷Helper methods
21. def  $\beta_k(i, X/\alpha) += v$ : ▷This method updates the weights, Q and  $W_k$ 
22.   if  $\alpha = \varepsilon$ :  $Q.push(\langle i, X/\varepsilon \rangle)$  ▷Newly completed items get scheduled in the priority queue
23.   else:  $W_k[\alpha_1].add(\langle i, X/\alpha \rangle)$  ▷Item  $\langle i, X/\alpha \rangle$  is waiting for  $\alpha_1$  to move forward
24.    $\beta_k(i, X/\alpha) \leftarrow \beta_k(i, X/\alpha) + v$ 

```

Table 2: Grammar size of different PCFGs  $\mathcal{G}$  used in our experiments.  $\mathcal{G}_1$  and  $\mathcal{G}_2$  are respectively the grammars obtained after the preprocessing step of Earley and PrefixEarley. The ratio column reports  $|\mathcal{G}_2|/|\mathcal{G}_1|$ . The grammars reported in the table are WSJ 500 (obtained from the first 500 sentences of the Wall Street Journal Corpus), WSJ 5000 (obtained from the first 5000 sentences of the Wall Street Journal Corpus), and “Social Discourse”. We empirically found that applying EnsureCNF after the prefix grammar transformation and right before EnsureNullaryFree yielded a smaller grammar; the results shown in this table follow this approach.

Grammar	$ \mathcal{G} ;  \mathcal{N} $	$ \mathcal{G}_1 ;  \mathcal{N}_1 $	$ \mathcal{G}_2 ;  \mathcal{N}_2 $	Ratio $ \mathcal{G}_2 / \mathcal{G}_1 $
WSJ 500	12,573; 69	15,981; 1,775	43,701; 5,201	2.73×
WSJ 5000	116,667; 448	177,303; 30,878	494,017; 84,566	2.79×
Social Discourse [Sparse]	72,712; 233	72,712; 233	143,548; 435	1.97×

## F Practical Instantiations of the Next-Token Weight Vector Algorithm

In §4, we presented a general framework for computing the next-token weight vector  $\pi(s)$  via lattice parsing and algorithmic differentiation. This section instantiates this framework with two concrete parsing algorithms: CKY (App. F.1) and Earley’s (App. F.2). In both cases, we derive the gradient algorithm by manually applying the rules of reverse-mode algorithmic differentiation to the corresponding lattice parser.<sup>23</sup>

### F.1 Next-Token Weight Vector with CKY

**Algorithm 5 Left:** the LatticeCKY is a specialized lattice parser that parses the next-token lattice  $\mathcal{L}_s(\theta)$  for some choice of  $\theta \in \mathbb{W}^\Sigma$ . Note that LatticeCKY is directly derived from CKY, differing only in the terminal step (line 6), which scans all terminal symbols simultaneously, weighting each by  $\theta_a$ . **Right:** the NextTokenCKY algorithm, which was derived by manually applying algorithmic differentiation to LatticeCKY. Given the input grammar  $\mathcal{G}$ , the input string  $s_1 \dots s_N$ , and the inside weight columns  $\langle \beta_0, \dots, \beta_N \rangle$ , it returns the next-token weight vector  $\pi(s)$ . Note that both algorithms call CKY as a subroutine to compute the inside weight columns (which, however, can be amortized across subsequent calls when NextTokenCKY is called incrementally).

<pre> 1. <b>function</b> LatticeCKY(<math>\mathcal{G}, s_1 \dots s_N, \theta</math>) 2.   <math>\mathcal{G} \leftarrow \text{EnsureCNF}(\mathcal{G})</math> <math>\triangleright</math> CNF conversion; memoize 3.   <math>\langle \beta_0, \dots, \beta_N \rangle \leftarrow \text{CKY}(\mathcal{G}, s_1 \dots s_N)</math> 4.   <math>z \leftarrow \text{defaultdict}(0)</math> <math>\triangleright</math> Initialize new column 5.   <math>\triangleright</math> Base case: 6.   <b>for</b> <math>(X \xrightarrow{w} a) \in \mathcal{R}</math> with <math>a \in \Sigma</math>: 7.     <math>z(N, X) += w \cdot \theta_a</math> 8.   <math>\triangleright</math> Recursive step: 9.   <b>for</b> <math>i</math> in <math>N, \dots, 0</math>: <math>\triangleright</math> Iterate over the start point. 10.    <b>for</b> <math>j</math> in <math>i+1 \dots N</math>: <math>\triangleright</math> Split point 11.     <b>for</b> <math>(X \xrightarrow{w} YZ) \in \mathcal{R}</math>: <math>\triangleright</math> Binary rules 12.      <math>z(i, X) += w \cdot \beta_j(i, Y) \cdot z(j, Z)</math> 13.   <b>return</b> <math>z(0, S)</math> <math>\triangleright</math> equals <math>Z_s(\theta)</math> </pre>	<pre> 1. <b>function</b> NextTokenCKY(<math>\mathcal{G}, s_1 \dots s_N</math>) 2.   <math>\triangleright</math> Create prefix grammar in CNF; memoize 3.   <math>\overline{\mathcal{G}} \leftarrow \text{EnsureCNF}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G})))</math> 4.   <math>\langle \beta_0, \dots, \beta_N \rangle \leftarrow \text{CKY}(\overline{\mathcal{G}}, s_1 \dots s_N)</math> 5.   <math>\nabla z \leftarrow \text{defaultdict}(0)</math> 6.   <math>\nabla z(0, S) += 1</math> <math>\triangleright</math> Base case 7.   <b>for</b> <math>i</math> in <math>0 \dots N-1</math>: <math>\triangleright</math> Inverse iteration 8.    <b>for</b> <math>j</math> in <math>i+1 \dots N</math>: <math>\triangleright</math> Loop over the split point <math>j</math> 9.     <b>for</b> <math>(X \xrightarrow{w} YZ) \in \mathcal{R}</math>: 10.      <math>\nabla z(j, Z) += w \cdot \beta_j(i, Y) \cdot \nabla z(i, X)</math> 11.   <math>\pi \leftarrow 0^\Sigma</math> 12.   <b>for</b> <math>(X \xrightarrow{w} a) \in \mathcal{R}</math> where <math>a \in \Sigma</math>: <math>\triangleright</math> Preterminal rules 13.    <math>\pi_a += w \cdot \nabla z(N, X)</math> 14.   <b>return</b> <math>\pi</math> </pre>
--	---

In LatticeCKY (Alg. 5, left), the *forward values*  $z(i, X)$  for  $i = 0, \dots, N$  and  $X \in \mathcal{N}$  aggregate partial contributions to  $Z_s(\theta) = z(0, S)$ . In NextTokenCKY (right), the corresponding *backward values*  $\nabla z(i, X) \stackrel{\text{def}}{=} \frac{\partial z(0, S)}{\partial z(i, X)}$  disaggregate this sum and isolate each token’s contribution. Readers familiar with the inside–outside algorithm (Baker, 1979) will notice that the backward pass is closely related to the *outside algorithm*, following Eisner’s (2016) presentation of the outside algorithm as the *adjoint* of the inside algorithm. The runtime of NextTokenCKY, following the analysis of PrefixCKY, is  $\mathcal{O}(|\text{EnsureCNF}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G})))|N^3)$ , including the cost of computing the *inside weight columns*  $\beta_0, \dots, \beta_N$ , which can be amortized across incremental computations.

### F.2 Next-Token Weight Vector with Earley’s Algorithm

We now adapt the lattice parsing and algorithmic differentiation approach to Earley’s algorithm (Earley, 1970) and its weighted version (Stolcke, 1995). Earley’s algorithm requires less preprocessing than CKY—only EnsureNullaryFree and EnsureUnaryCycleFree—and runs in subcubic time for most grammars: linear for deterministic grammars, and quadratic for unambiguous ones.

**Next-token weight vector with Earley (Alg. 6).** In IncrEarleyLattice, the *forward values*  $z(j, Y/\varepsilon)$  for  $j \in \{0, \dots, N\}$  and  $Y \in \mathcal{N}$  aggregate partial contributions to the goal item  $z(0, S) = Z_{s_1 \dots s_N}(\theta)$ . In NextTokenEarley, the corresponding *backward values*  $\nabla z(j, Y/\varepsilon) \stackrel{\text{def}}{=} \frac{\partial z(0, S)}{\partial z(j, Y/\varepsilon)}$  disaggregate this sum to isolate each token’s contribution. The runtime of NextTokenEarley, following the analysis of PrefixEarley, is  $\mathcal{O}(|\text{EnsureUnaryCycleFree}(\text{EnsureNullaryFree}(\text{Prefix}(\text{EnsureCTF}(\mathcal{G}))))|N^3)$ , including the cost of computing the *inside weight columns*  $\beta_0, \dots, \beta_N$ , which can be amortized across incremental computations.

<sup>23</sup>In general, the gradient program can also be obtained automatically, e.g., via program tracing, operator overloading, or source-to-source transformation (see Griewank and Walther, 2008, for an overview), using tools such as PyTorch (Paszke et al., 2017), TensorFlow (Abadi et al., 2016), JAX (Bradbury et al., 2018), Tapenade (Hascoët and Pascual, 2013), or Dyna (Eisner et al., 2005).

**Algorithm 6 Left:** `IncrEarleyLattice` is a lattice parser derived from Earley’s algorithm (Alg. 4) that evaluates the next-token lattice  $\mathcal{L}_s(\theta)$  for some choice of  $\theta \in \mathbb{W}^\Sigma$ . It differs from Earley only in the SCAN step, which scans all terminal symbols simultaneously, weighting each by  $\theta_a$ . Since only a single terminal is scanned at position  $N + 1$ , the following simplifications apply: (1) SCAN processes only items from  $W_N$  that complete after scanning one symbol (line 9), and ATTACH only advances items with a single remaining symbol (line 14); (2) since no further columns are needed, PREDICT can be omitted entirely. **Right:** `NextTokenEarley` is the gradient algorithm of `IncrEarleyLattice`, obtained by applying algorithmic differentiation. Given the prefix grammar  $\vec{\mathcal{G}}$ , the input string  $s = s_1 \cdots s_N$ , and the inside weight columns from Earley, it returns the next-token weight vector  $\pi(s)$ . The backward values  $\nabla z(j, Y/\varepsilon)$  are computed by the memoized helper function (line 9).

---

<pre> 1. <b>function</b> IncrEarleyLattice(<math>\mathcal{G}, s_1 \cdots s_N, \theta</math>) 2.   <math>\vec{\mathcal{G}} \leftarrow \text{EnsureNullaryFree}(\mathcal{G})</math> 3.   <math>\mathcal{G} \leftarrow \text{EnsureUnaryCycleFree}(\vec{\mathcal{G}})</math> 4.   <math>\langle \mathcal{N}, \Sigma, S, \mathcal{R} \rangle \leftarrow \mathcal{G}</math> 5.   <math>\langle \beta_0, \dots, \beta_N \rangle, \langle W_0, \dots, W_N \rangle \leftarrow \text{Earley}(\mathcal{G}, s_1 \cdots s_N)</math> 6.   <math>z \leftarrow \text{defaultdict}(0)</math> <span style="float: right;"><math>\triangleright</math>Initialize new column</span> 7.   <math>Q \leftarrow \text{priority\_queue}()</math> 8.   <b>for</b> <math>a \in W_N.\text{keys} \cap \Sigma</math>: <span style="float: right;"><math>\triangleright</math>SCAN</span> 9.     <b>for</b> <math>\langle i, X/a \rangle \in W_N[a]</math>: 10.      <math>z(i, X/\varepsilon) += \beta_N(i, X/a) \cdot \theta_a</math> 11.      <math>Q.\text{push}(\langle i, X/\varepsilon \rangle)</math> 12.   <b>while</b> <math>Q</math>: <span style="float: right;"><math>\triangleright</math>ATTACH</span> 13.     <math>\langle j, Y/\varepsilon \rangle \leftarrow Q.\text{pop}()</math> 14.     <b>for</b> <math>\langle i, X/Y \rangle \in W_j[Y]</math>: 15.      <math>z(i, X/\varepsilon) += \beta_j(i, X/Y) \cdot z(j, Y/\varepsilon)</math> 16.     <math>Q.\text{push}(\langle i, X/\varepsilon \rangle)</math> 17.   <b>return</b> <math>z(0, S/\varepsilon)</math> <span style="float: right;"><math>\triangleright</math>equals <math>Z_{s_1 \cdots s_N}(\theta)</math></span> </pre>	<pre> 1. <b>def</b> NextTokenEarley(<math>\vec{\mathcal{G}}, s_1 \cdots s_N</math>): 2.   <math>\vec{\mathcal{G}} \leftarrow \text{Prefix}(\text{EnsureCTF}(\vec{\mathcal{G}}))</math> <span style="float: right;"><math>\triangleright</math>Prefix transformation</span> 3.   <math>\langle \beta_0, \dots, \beta_N \rangle, \langle W_0, \dots, W_N \rangle \leftarrow \text{Earley}(\vec{\mathcal{G}}, s_1 \cdots s_N)</math> 4.   <math>\pi \leftarrow 0^\Sigma</math> 5.   <b>for</b> <math>a \in W_N.\text{keys} \cap \Sigma</math>: 6.     <b>for</b> <math>\langle i, X/a \rangle \in W_N[a]</math>: 7.      <math>\pi_a += \beta_N(i, X/a) \cdot \nabla z(i, X/\varepsilon)</math> 8.   <b>return</b> <math>\pi</math> 9.   <b>def</b> <math>\nabla z(j, Y/\varepsilon)</math>: <span style="float: right;"><math>\triangleright</math>Memoize for efficiency</span> 10.    <b>if</b> <math>(j, Y) = (0, S)</math>: <b>return</b> 1 <span style="float: right;"><math>\triangleright</math>Base Case</span> 11.    <math>v \leftarrow 0</math> 12.    <b>for</b> <math>\langle i, X/Y \rangle \in W_j[Y]</math>: <span style="float: right;"><math>\triangleright</math>Recursive bottom-up call.</span> 13.     <math>v += \beta_j(i, X/Y) \cdot \nabla z(i, X/\varepsilon)</math> <span style="float: right;"><math>\triangleright</math>Note: <math>i \leq j</math>.</span> 14.    <b>return</b> <math>v</math> </pre>
---	---

---

## G Deferred Proofs

### G.1 Proof of Prop. 1

**Proposition 1.** *Let  $\mathcal{G}$  be a CFG, and  $\vec{\mathcal{G}}$  be its prefix grammar. Then,  $\vec{\mathcal{G}}$  correctly encodes the prefix language of  $\mathcal{G}$  (i.e.,  $\llbracket \vec{\mathcal{G}} \rrbracket = \overline{\llbracket \mathcal{G} \rrbracket}$ ).*

*Proof.* We strengthen the proposition to a claim over all  $\alpha \in \mathcal{N} \cup \Sigma$ , prove the strengthening by induction on derivation-tree height in  $\vec{\mathcal{G}}$ , and recover the proposition by specializing to  $\alpha = S$ .

**Strengthening.** The Prop. 1 is implied by the following statement, for all  $\alpha \in \mathcal{N} \cup \Sigma$  and all  $s \in \Sigma^+$ :

$$\llbracket \vec{\mathcal{G}}_{\alpha'} \rrbracket(s) = \sum_{t \in \Sigma^*} \llbracket \mathcal{G}_{\alpha} \rrbracket(st) \quad \text{and} \quad \llbracket \vec{\mathcal{G}}_{\alpha} \rrbracket(s) = \llbracket \mathcal{G}_{\alpha} \rrbracket(s) \quad (30)$$

The latter holds because  $\vec{\mathcal{G}}$ 's additional rules have only prime nonterminals or  $\vec{S}$  on the LHS (Def. 1), so non-prime  $\alpha$ -rooted derivations use only the rules of  $\mathcal{R}$ . Thus, we need only prove the former, which we re-express self-recursively in  $\vec{\mathcal{G}}$ :

$$\llbracket \vec{\mathcal{G}}_{\alpha'} \rrbracket(s) = \sum_{t \in \Sigma^*} \llbracket \vec{\mathcal{G}}_{\alpha} \rrbracket(st) \quad (31)$$

$\Leftrightarrow$  **Hooking up Eq. (30) to Prop. 1.** Specializing the first equation of Eq. (30) to  $\alpha = S$  and expanding  $\vec{S}$ 's rules (Def. 1) yields the proposition. Recall that  $\vec{\mathcal{G}}$ 's start symbol  $\vec{S}$  has two rules:  $\vec{S} \xrightarrow{1} S'$  and  $\vec{S} \xrightarrow{\tau(S)} \varepsilon$ .

• **Case ( $s \in \Sigma^+$ ):** Only the first rule contributes (the second requires  $s = \varepsilon$ ). By Eq. (30) at  $\alpha = S$ :

$$\llbracket \vec{\mathcal{G}} \rrbracket(s) = \llbracket \vec{\mathcal{G}}_{S'} \rrbracket(s) = \sum_{t \in \Sigma^*} \llbracket \mathcal{G}_S \rrbracket(st) = \overline{\llbracket \mathcal{G} \rrbracket}(s)$$

• **Case ( $s = \varepsilon$ ):** Every  $S'$ -derivation has nonempty yield (the prime chain from the root must terminate at a terminal via Eq. (4c)), so only the  $\vec{S} \rightarrow \varepsilon$  rule contributes. By the definition of  $\tau(S)$ :

$$\llbracket \vec{\mathcal{G}} \rrbracket(\varepsilon) = \tau(S) = \sum_{t \in \Sigma^*} \llbracket \mathcal{G}_S \rrbracket(t) = \overline{\llbracket \mathcal{G} \rrbracket}(\varepsilon)$$

**Notation.** We introduce the following notation to keep the proof tidy.

- Let  $\vec{\mathcal{D}}$  denote the set of derivation trees of  $\vec{\mathcal{G}}$ .
- Let  $\vec{\mathcal{D}}^{(h)}$  denote the subset of trees in  $\vec{\mathcal{D}}$  of height  $\leq h$ ;  $\vec{\mathcal{D}}_{\alpha}^{(h)}(s)$  restricts to  $\alpha$ -rooted trees with yield  $s$ .
- Let  $\mathcal{G}^{(h)}$  denote  $\mathcal{G}$  with its derivations restricted to height  $\leq h$ .
- Let  $\vec{\mathcal{G}}^{(h)} \stackrel{\text{def}}{=} \overline{(\mathcal{G}^{(h)})}$  be the prefix grammar of  $\mathcal{G}^{(h)}$ ; its prime rule weights use the height-bounded totals  $\vec{\tau}_{\alpha}^{(h)}$  rather than  $\tau(\alpha)$ .
- Let  $\vec{W}_{\alpha}^{(h)}(s) \stackrel{\text{def}}{=} \sum_{d \in \vec{\mathcal{D}}_{\alpha}^{(h)}(s)} w(d)$
- Let  $\vec{\tau}_{\alpha}^{(h)} \stackrel{\text{def}}{=} \sum_{t \in \Sigma^*} \vec{W}_{\alpha}^{(h)}(t)$

**Proof strategy.** We prove Eq. (30) by induction on the height  $h \geq 1$  of derivation trees in  $\vec{\mathcal{G}}$ . We define the height-indexed proposition:

$$\Phi(h) \iff \forall \alpha \in \mathcal{N} \cup \Sigma, s \in \Sigma^+ : \vec{W}_{\alpha'}^{(h)}(s) = \sum_{t \in \Sigma^*} \vec{W}_{\alpha}^{(h)}(st) \quad (\text{IH})$$

**Base case ( $h = 1$ ).** We show  $\Phi(1)$  holds. There are two ways for a tree to have height one.

- **Case ( $\alpha \in \Sigma$ ):** Since  $\alpha' = \alpha$ , both sides of  $\Phi(1)$  reduce to  $\mathbb{1}\{s = \alpha\}$ , using  $s \in \Sigma^+$  to collapse the right-hand side sum.
- **Case ( $X \xrightarrow{w} \varepsilon$ ):**  $X'$  has no height-1 derivation since the rules of Eq. (4c) have nonempty right-hand sides, so the LHS of  $\Phi(1)$  is 0. The RHS is also 0 because  $\vec{W}_X^{(1)}(st)$  is nonzero only when  $st = \varepsilon$ , which contradicts  $s \in \Sigma^+$ .

**Inductive hypothesis.** Suppose  $\Phi(h)$  holds for some  $h \geq 1$ .

**Inductive case.** We show  $\Phi(h+1)$  holds. In this case, we know that  $\alpha'$  is of the form  $X'$  for  $X \in \mathcal{N}$ .

$$\begin{aligned} & \sum_{t \in \Sigma^*} \overrightarrow{W}_X^{(h+1)}(st) \\ &= \sum_{t \in \Sigma^*} \sum_{X \xrightarrow{w} \alpha_1 \cdots \alpha_K} \sum_{st = v_1 \cdots v_K} w \cdot \overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha_k}^{(h)}(v_k) \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \cdots \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \end{aligned} \quad (32a)$$

$$= \sum_{X \xrightarrow{w} \alpha_1 \cdots \alpha_K} \sum_{k=1}^K \sum_{\substack{s=v_1 \cdots v_k \\ t=v_{k+1} \cdots v_K}} w \cdot \overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha_k}^{(h)}(v_k) \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \cdots \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \quad (32b)$$

$$= \sum_{X \xrightarrow{w} \alpha_1 \cdots \alpha_K} \sum_{k=1}^K \sum_{s=v_1 \cdots v_k} w \cdot \overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha_k}^{(h)}(v_k) \sum_{t=v_{k+1} \cdots v_K} \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \cdots \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \quad (32c)$$

$$= \sum_{X \xrightarrow{w} \alpha_1 \cdots \alpha_K} \sum_{k=1}^K \sum_{s=v_1 \cdots v_k} w \cdot \overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha_k}^{(h)}(v_k) \left( \overrightarrow{\tau}_{\alpha_{k+1}}^{(h)} \cdots \overrightarrow{\tau}_{\alpha_K}^{(h)} \right) \quad (32d)$$

$$= \sum_{X \xrightarrow{w} \alpha_1 \cdots \alpha'_k} \sum_{s=v_1 \cdots v_k} w \cdot \overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha'_k}^{(h)}(v_k) \quad (32e)$$

$$= \overrightarrow{W}_{X'}^{(h+1)}(s) \quad (32f)$$

The manipulation above proceeds in six steps:

1. **Tree-weight recursion for  $X$ .** Each height- $\leq h+1$  tree rooted at  $X$  applies some rule  $X \xrightarrow{w} \alpha_1 \cdots \alpha_K$  at the root with subtrees of height  $\leq h$ .
2. **Border position.** Reparametrize the joint sum  $\sum_t \sum_{st=v_1 \cdots v_K}$  by introducing  $k \in \{1, \dots, K\}$  with  $s = v_1 \cdots v_k$  and  $t = v_{k+1} \cdots v_K$ .
3. **Factor the  $t$ -side.** The factors  $\overrightarrow{W}_{\alpha_1}^{(h)}(v_1) \cdots \overrightarrow{W}_{\alpha_k}^{(h)}(v_k)$  are fixed by the  $s$ -split, so they move out of the inner sum over  $t = v_{k+1} \cdots v_K$ .
4. **Collapse to totals.** We simplify the inner summation:

$$\begin{aligned} & \sum_{t=v_{k+1} \cdots v_K} \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \cdots \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \\ &= \sum_{v_{k+1}} \cdots \sum_{v_K} \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \cdots \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \\ &= \left[ \sum_{v_{k+1}} \overrightarrow{W}_{\alpha_{k+1}}^{(h)}(v_{k+1}) \right] \cdots \left[ \sum_{v_K} \overrightarrow{W}_{\alpha_K}^{(h)}(v_K) \right] \\ &= \overrightarrow{\tau}_{\alpha_{k+1}}^{(h)} \cdots \overrightarrow{\tau}_{\alpha_K}^{(h)} \end{aligned}$$

5. **The prime rule.** The absorbed weight  $w' = w \cdot \overrightarrow{\tau}_{\alpha_{k+1}}^{(h)} \cdots \overrightarrow{\tau}_{\alpha_K}^{(h)}$  matches the prime rule  $X' \rightarrow \alpha_1 \cdots \alpha_{k-1} \alpha'_k$  (Eq. (4c)), so the joint sum over (original rule,  $k$ ) becomes a pattern-match over such prime rules. Applying  $\Phi(h)$  to  $\alpha_k$  rewrites  $\overrightarrow{W}_{\alpha_k}^{(h)}(v_k)$  as  $\overrightarrow{W}_{\alpha'_k}^{(h)}(v_k)$ .

6. **Tree-weight recursion for  $X'$ .** The result is the tree-weight recursion for  $\overrightarrow{W}_{X'}^{(h+1)}(s)$ .

This establishes  $\Phi(h+1)$ , completing the induction. Since every derivation tree in  $\overrightarrow{\mathcal{G}}$  has finite height,  $\Phi(h)$  for all  $h \geq 1$  recovers Eq. (30), and the proposition follows via the hook-up step.<sup>24</sup> ■

<sup>24</sup>Our proof implicitly uses  $\overrightarrow{\tau}^{(h)}$  in the prime rule weights, i.e., the totals of  $\overrightarrow{\mathcal{G}}^{(h)}$  rather than  $\overrightarrow{\mathcal{G}}$ . As  $h \rightarrow \infty$ ,  $\overrightarrow{\tau}_{\alpha}^{(h)} \rightarrow \tau(\alpha)$ , recovering  $\overrightarrow{\mathcal{G}}$ 's rule weights;  $\Phi(h)$  for all  $h \geq 1$  then gives Eq. (30). A more formal induction use pairs  $(\mathbb{1}\{\alpha \in \mathcal{N}\}, h)$ .

## G.2 Proof of Prop. 2

**Proposition 2.** *Let  $\mathcal{G}$  be a context-free grammar in canonical two-form and let  $\vec{\mathcal{G}}$  be its prefix grammar (Def. 1). Then, the size of  $\vec{\mathcal{G}}$  is bounded by*

$$|\vec{\mathcal{G}}| \leq \frac{8}{3}|\mathcal{G}| + 3 \quad (5)$$

*Proof.* According to Def. 1, the prefix grammar  $\vec{\mathcal{G}}$  consists of the original rules  $\mathcal{R}$  plus the additional prefix rules  $\mathcal{R}'$ . We bound each contribution:

- The original rules  $\mathcal{R}$  contribute  $|\mathcal{G}|$ .
- Eqs. (4a) and (4b) contribute  $2 + 1 = 3$ .
- For Eq. (4c):
  - Each nullary rule (size 1) produces no prefix rules.
  - Each unary rule (size 2) produces one prefix rule ( $k=1$ ) of size 2, contributing 2 to the size.
  - Each binary rule (size 3) produces two prefix rules: one for  $k=1$  of size 2 and one for  $k=2$  of size 3, contributing  $2 + 3 = 5$  to the size.

In these cases, the ratio of the prefix contribution to the original rule size is at most  $5/3$  (attained by binary rules). So the total contribution of Eq. (4c) is at most  $\frac{5}{3}|\mathcal{G}|$ .

Summing:  $|\vec{\mathcal{G}}| \leq |\mathcal{G}| + 3 + \frac{5}{3}|\mathcal{G}| = \frac{8}{3}|\mathcal{G}| + 3$ . ■