# Probability Distributions Computed by Hard-Attention Transformers

**Andy Yang,**[1] **Anej Svete,**[2] **Jiaoda Li,**[2] **Anthony Widjaja Lin,**[3,4] **Jonathan Rawski,**[5]
**Ryan Cotterell,**[2] and **David Chiang**[1]

[1]University of Notre Dame, USA     [2]ETH Zürich, Switzerland
[3]Max-Planck Institute for Software Systems, Kaiserslautern, Germany
[4]University of Kaiserslautern-Landau, Germany     [5]San José State University, USA

## Abstract

Most expressivity results for transformers treat them as language recognizers (which accept or reject strings), and not as they are used in practice, as language models (which generate strings autoregressively and probabilistically). Here, we characterize the probability distributions that transformer language models can express. We show that making transformer language recognizers autoregressive can sometimes increase their expressivity, and that making them probabilistic can break equivalences that hold in the non-probabilistic case. Our overall contribution is to tease apart what functions transformers are capable of expressing, in their most common use-case as language models.

## 1 Introduction

Most work studying transformer expressivity, that is, what classes of computations transformers can perform, treats them as *language recognizers*, where the input is a string and the output is a binary classification: true if the string is accepted and false otherwise (Strobl et al., 2024). However, the most common practical use of transformers is as *language models*, which differ in two ways: first, the input is a prefix of a string, and the output is a prediction of the next symbol; second, the prediction is a probability distribution rather than a binary decision. Such probability distributions, when estimated from large text corpora, have enabled a wide range of applications in natural language processing and beyond. An open and fundamental question concerns which probability distributions transformer language models can express. This distributional perspective exposes where previously-proven equivalences retain or lose their validity in the probabilistic setting.

We distinguish, on the one hand, between *unweighted* (or equivalently, *Boolean-weighted*) and *real-weighted* computation, and, on the other hand, between *classifiers*, which map a complete string to a value, and *autoregressors*, which map each prefix to a distribution over the next token. Under this terminology, most theoretical work on transformer expressivity (e.g. Yang et al., 2024; Jerad et al., 2025) focuses on Boolean-weighted classifiers, while practical applications use transformers as real-weighted autoregressors. The theoretical expressivity of transformers as real autoregressors has been comparatively underexplored.

In this paper, we answer this question for several variants of transformers. Yang et al. (2024) proved that strictly-masked rightmost UHATs, as Boolean classifiers, recognize the same languages as linear temporal logic (LTL) and counter-free automata. There are two commonly-used weighted analogues of regular languages, those defined by weighted deterministic finite automata (DFAs) and weighted nondeterminizible finite automata (NFAs), and each of these have counter-free versions. Surprisingly, these two diverge in the real-weighted setting, despite being equivalent in the Boolean setting. We prove that, as autoregressors, these transformers define exactly the same weighted languages as counter-free DFAs.

Jerad et al. (2025) proved that *leftmost* UHATs, as Boolean classifiers, recognize the same languages as a fragment of LTL, called in our notation TL[$\mathbf{P}$], or a subclass of counter-free DFAs called partially ordered DFAs. Similarly, Li and Cotterell (2025) proved that softmax attention transformers

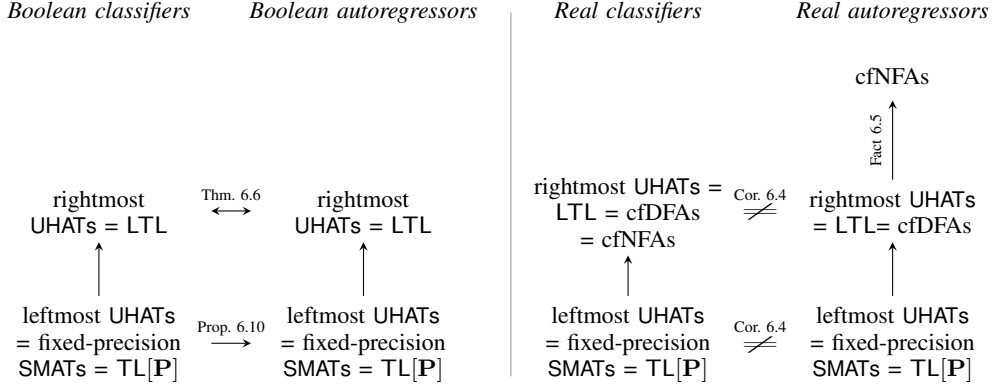| *Boolean classifiers* | *Boolean autoregressors* | *Real classifiers* | *Real autoregressors* |



Figure 1: In the Boolean semiring, equivalences from the literature (Yang et al., 2024; Jerad et al., 2025; Yang et al., 2025) carry over from classifiers to autoregressors; however, sometimes autoregressors are more expressive than classifiers. In the real semiring, LTL and counter-free DFAs and NFAs become less expressive than counter-free NFAs, and rightmost UHATs are only as expressive as the former. Key: $\rightarrow$ strict inclusion, $\leftrightarrow$ equivalence. $\neq$ incomparable.

(SMATs) with fixed precision, as Boolean classifiers, recognize the same class of languages. But here we show that as autoregressors, these variants become slightly more powerful.

Yang et al. (2025) considered SMATs with fixed precision but arbitrary precision inside attention. As Boolean classifiers, these transformers are exactly equivalent to a temporal logic extended with counting operators. But here we show that, as autoregressors, they become slightly more powerful (for a given fixed depth).

In §3, we define notational preliminaries. We then (§4) define the classes of transformers we consider, and show how they can be used as classifiers and as autoregressors. We then (§5) introduce two other formalisms, deterministic finite automata (DFAs) and linear temporal logic (LTL), and show how they can also be seen as instantiations of the framework we showed for transformers. Then (§6), using LTL, we investigate the expressive power of transformers as both classifiers and autoregressors, yielding the results shown in Fig. 1.

## 2 RELATED WORK

Theoretical study of transformers as language models has not gone totally neglected. Hahn (2020) compared a SMAT language model with a probabilistic finite automaton for parity (strings that have an odd number of 1's). Yao et al. (2021), following previous work on RNNs, consider a transformer language model to $\epsilon$-**generate** a language if it assigns probability at least $\epsilon$ to each symbol in every string in the language (and no strings not in the language). They also discuss how to convert a construction for a bounded Dyck language (strings of matching parentheses up to a certain depth) from an $\epsilon$-generator to a language recognizer. These studies were specialized to particular languages and used specialized ways of comparing distributions that do not generalize in an obvious way.

Svete and Cotterell (2024) made progress by showing that average-hard attention transformer language models can exactly express all $n$-gram language models. This made an important step towards understanding the expressivity of transformer language models, though it left an exact charcaterization open.

Bhattamishra et al. (2020) proved theoretical results on transformers as language recognizers but carried out experiments on transformer language models for the **character prediction task**, which is to predict, at each position, the set of next possible symbols (what we will call Boolean autoregression). This experimental setup was previously used in studies of RNNs, and has been adopted in other studies of transformers (Huang et al., 2025; Yang et al., 2025). We discuss the experiment of Yang et al. (2025) in §6.3.3.

## 3  PRELIMINARIES

Throughout this paper, we work with weighted languages. We define some key concepts here, but for a more detailed introduction, see the handbook chapter by Droste and Kuske (2021).

Let $\Sigma$ be an **alphabet**, that is, a finite, non-empty set of **symbols**, and let $\Sigma^*$ be the set of strings over $\Sigma$. We often augment $\Sigma$ with start and end symbols BOS and EOS, but never consider BOS or EOS to belong to $\Sigma$. For any string $\boldsymbol{w} = w_1 \cdots w_n$, we write the length of $\boldsymbol{w}$ as $|\boldsymbol{w}| = n$. We write $\boldsymbol{w}_{<i} = w_1 \cdots w_{i-1}$ and $\boldsymbol{w}_{\leq i} = w_1 \cdots w_i$.

We think of weights and probabilities as elements of **semirings**, an abstraction of the usual addition and multiplication operations that allows results and algorithms apply generically to multiple settings. A semiring $\mathbb{K}$ has an addition operation $\oplus$, additive identity $\mathbf{0}$, multiplication operation $\otimes$, and multiplicative identity $\mathbf{1}$. The two semirings we focus on in this paper are the **(extended) real semiring** $\overline{\mathbb{R}}_{\geq 0}$, which contains all nonnegative real numbers and $+\infty$, and in which $\oplus$ and $\otimes$ are real addition and multiplication; and the **Boolean semiring** $\mathbb{B}$, in which $\oplus$ is disjunction ($\vee$), $\mathbf{0}$ is false ($\bot$), $\otimes$ is conjunction ($\wedge$), and $\mathbf{1}$ is true ($\top$).

A **weighted language** (also called a **formal power series**) is a function $S \colon \Sigma^* \to \mathbb{K}$. When $\mathbb{K}$ is complete (that is, it allows infinite summations, as $\overline{\mathbb{R}}_{\geq 0}$ and $\mathbb{B}$ do), we call a weighted language **normalized** if $\sum_{\boldsymbol{w} \in \Sigma^*} S(\boldsymbol{w}) = \mathbf{1}$.

For sets $X$ and $Y$, we write $Y^X$ for the set of functions from $X$ to $Y$, and $2^X$ for $\{0, 1\}^X$ or the power set of $X$. For any proposition $\phi$, we write $\mathbb{I}\{\phi\}$ to be 1 if $\phi$ is true and 0 if $\phi$ is false.

## 4  TRANSFORMER LANGUAGE MODELS

In this section, we recall the definition of transformers that we will use throughout most of this paper. We also distinguish between two ways that transformers (and other formalisms) can be used to define weighted languages.

### 4.1  UNIQUE HARD ATTENTION TRANSFORMERS

Following Yang et al. (2024), we use **unique-hard attention transformers** (UHATs), specifically, with rightmost-hard attention, strict future masking, and no position embeddings. We give a definition of strictly masked rightmost-hard attention here; for a definition of the rest of the network, see, for example, the survey by Strobl et al. (2024).

The attention function receives a sequence of query vectors $\mathbf{q}^{(i)} \in \mathbb{R}^{d_k}$, key vectors $\mathbf{k}^{(j)} \in \mathbb{R}^{d_k}$, and value vectors $\mathbf{v}^{(j)} \in \mathbb{R}^d$, for $i, j \in [n]$. At each position $i$, it computes

$$\text{Att}\left((\mathbf{q}^{(i)})_{i \in [n]}, (\mathbf{k}^{(j)})_{j \in [n]}, (\mathbf{v}^{(j)})_{j \in [n]}\right) = (\mathbf{c}^{(i)})_{i \in [n]} \tag{1}$$

where

$$a_i(j) = \mathbf{q}^{(i)} \cdot \mathbf{k}^{(j)} \qquad \text{is an attention score for each position } j,$$
$$a_i^* = \max_{j < i} a_i(j) \qquad \text{is the maximum attention score,}$$
$$j_i = \max\{j < i \mid a_i(j) = a_i^*\} \qquad \text{is the rightmost maximum-scoring position, and}$$
$$\mathbf{c}^{(i)} = \begin{cases} \mathbf{v}^{(j_i)} & \text{if } i > 0 \\ \mathbf{0} & \text{if } i = 0 \end{cases} \qquad \text{is the attention output.}$$

Given an input string $\boldsymbol{w} = w_1 \cdots w_n$, a transformer $\mathcal{T}$ prepends a symbol $w_0 = \text{BOS}$ and computes a sequence of "states" $\mathcal{T}(\boldsymbol{w}) = (\mathbf{h}^{(0)}, \ldots, \mathbf{h}^{(n)})$, where $\mathbf{h}^{(i)} \in \mathbb{R}^d$ is the state after reading $w_i$. There are at least two ways to use $\mathcal{T}$ to define a weighted language.[1]

---

[1] A third intermediate way would be to multiply the weights at each position like an autoregressive model, but not to pass the output symbol at each position autoregressively to the input at the next position. Although interesting in its own right, it has not, to our knowledge, been used with any neural sequence models, and we do not explore this style of model here.

## 4.2 CLASSIFIERS

The first way that a transformer can define a weighted language is as a **classifier**.

**Definition 4.1.** *A UHAT **classifier** is a pair $C = (\mathcal{T}, c)$, where $\mathcal{T}\colon \Sigma^* \to (\mathbb{R}^d)^*$ is a UHAT and $c\colon \mathbb{R}^d \to \mathbb{K}$ outputs a scalar weight at the last position only:*

$$C(\boldsymbol{w}) = c(\mathcal{T}(\boldsymbol{w})_n). \tag{2}$$

For the Boolean semiring ($\mathbb{K} = \mathbb{B}$), we accept a string iff the transformer outputs $\top$ at the last position. For example, the output function could be $c(y) = \mathbb{I}\{\mathbf{W}y + \mathbf{b} \geq 0\}$, where $\mathbf{W}$ and $\mathbf{b}$ are parameters. This is the setup used for binary classification with a transformer encoder (Devlin et al., 2019) and in most theoretical papers on transformer expressivity.

## 4.3 AUTOREGRESSIVE MODELS

The second way for a transformer to define a weighted language is as an **autoregressive model**, or an **autoregressor** for short (by analogy with *classifier*), which pairs a UHAT encoder with an output function $a\colon \mathbb{R}^d \to \mathbb{K}^{\Sigma \cup \{\text{EOS}\}}$, which outputs at each position a weight distribution for the next symbol, including EOS. In the real semiring ($\mathbb{K} = \overline{\mathbb{R}}_{\geq 0}$), a typical example of such an output function is $a(\mathbf{h}) = \operatorname{softmax}(\mathbf{W}\mathbf{h} + \mathbf{b})$.

To line up with the more familiar notation of conditional probability distributions, we write, for all $\sigma \in \Sigma \cup \{\text{EOS}\}$,

$$\Pr_A(\sigma \mid \boldsymbol{w}_{\leq i}) = a(\mathcal{T}(\boldsymbol{w})_i)(\sigma). \tag{3}$$

This is well-defined because $\mathcal{T}(\boldsymbol{w})_i$ depends only on $\boldsymbol{w}_{\leq i}$, that is, $\boldsymbol{w}_{\leq i} = \boldsymbol{w}'_{\leq i} \iff \mathcal{T}(\boldsymbol{w})_i = \mathcal{T}(\boldsymbol{w}')_i$. As suggested by this notation, we want $\Pr_A(\cdot \mid \boldsymbol{u})$ to be a probability distribution over $\Sigma \cup \{\text{EOS}\}$. But we impose a stronger condition. First, we extend $\Pr_A(\sigma \mid \boldsymbol{u})$ to the probability distribution of suffixes given (possibly empty) prefixes:

$$\Pr_A(\boldsymbol{v} \mid \boldsymbol{u}) = \left( \bigotimes_{i=1}^{|\boldsymbol{v}|} \Pr_A(v_i \mid \boldsymbol{u}\boldsymbol{v}_{<i}) \right) \otimes \Pr_A(\text{EOS} \mid \boldsymbol{u}\boldsymbol{v}) \tag{4}$$

$$\Pr_A(\boldsymbol{w}) = \Pr_A(\boldsymbol{w} \mid \epsilon). \tag{5}$$

Then we require that every such distribution sums to one:

**Definition 4.2.** *A UHAT **autoregressor** over a complete semiring $\mathbb{K}$ is a pair $A = (\mathcal{T}, a)$, where $\mathcal{T}\colon \Sigma^* \to (\mathbb{R}^d)^*$ is a UHAT, and $a\colon \mathbb{R}^d \to \mathbb{K}^{\Sigma \cup \{\text{EOS}\}}$ is a function such that for all $\boldsymbol{u} \in \Sigma^*$ (using the notation of Eqs. (3) and (4)),*

$$\bigoplus_{\boldsymbol{v} \in \Sigma^*} \Pr_A(\boldsymbol{v} \mid \boldsymbol{u}) = \mathbf{1}. \tag{6}$$

This implies that:

- An autoregressor generates strings symbol by symbol. That is, for all prefixes $\boldsymbol{u}$,

$$\sum_{\sigma \in \Sigma \cup \{\text{EOS}\}} \Pr_A(\sigma \mid \boldsymbol{u}) = \mathbf{1}. \tag{7}$$

- An autoregressor does not have any dead ends or endless loops. That is, for all prefixes $\boldsymbol{u}$,

$$\bigotimes_{i=1}^{n} \Pr_A(u_i \mid \boldsymbol{u}_{<i}) \neq \mathbf{0} \implies \Pr_A(\boldsymbol{u}\boldsymbol{v}) \neq \mathbf{0} \text{ for some suffix } \boldsymbol{v}. \tag{8}$$

- An autoregressor defines a normalized weighted language.

## 5 OTHER FORMALISMS

We can analogously use other formalisms to define classifier or autoregressive models. Any state encoder which sends a string $w_1 \cdots w_n$ to a sequence of "states" $q_0, \ldots, q_n \in Q$, such that $q_i$ depends only on $\boldsymbol{w}_{\leq i}$, can be equipped with an output function $c\colon Q \to \mathbb{K}$ to give a classifier model or $a\colon Q \to \mathbb{K}^{\Sigma \cup \{\text{EOS}\}}$ to give an autoregressive model exactly as we did with transformers above.
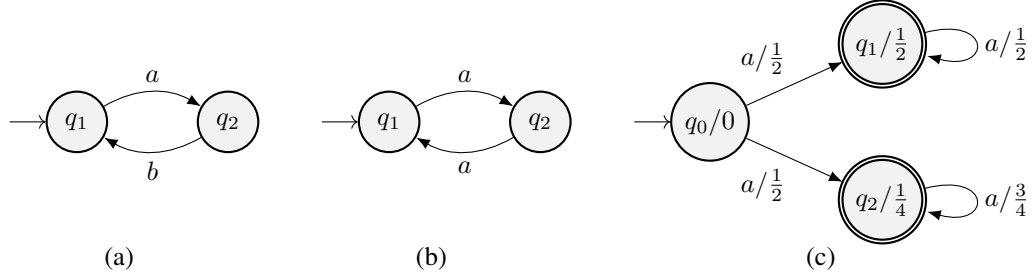
Figure 2: (a) A DFA that is counter-free (with $k = 2$). (b) A DFA that is not counter-free, because for all $k$, the strings $a^k$ and $a^{k+1}$ have opposite actions. (c) A counter-free weighted NFA that has no equivalent weighted DFA (Fact 6.5).

## 5.1 FINITE AUTOMATA

**Definition 5.1** (Deterministic finite automaton). *A **deterministic finite automaton** (DFA) is a tuple $M = (\Sigma, Q, \delta, \iota)$, where*

- *$\Sigma$ is an alphabet*
- *$Q$ is a finite set of **states***
- *$\delta \colon Q \times \Sigma \to Q$ is a **transition function***
- *$\iota \in Q$ is the **initial** state.*

*We extend $\delta$ to a mapping $\delta^* \colon Q \times \Sigma^* \to Q$ such that:*

$$
\begin{aligned}
\delta^*(q, \epsilon) &= q \\
\delta^*(q, \sigma \boldsymbol{w}) &= \delta^*(\delta(q, \sigma), \boldsymbol{w}).
\end{aligned}
\tag{9}
$$

A DFA $M$ defines a state encoder

$$
M \colon \Sigma^* \to Q^*
$$

$$
M(\boldsymbol{w})_i = \begin{cases} \iota & i = 0 \\ \delta^*(\iota, w_1 \cdots w_i) & 0 < i \le n. \end{cases}
\tag{10}
$$

A DFA with classifier outputs in the Boolean semiring is the same as the standard definition of a DFA: the states that output $\top$ are the accept states, and the states that output $\bot$ are the reject states.

A DFA with autoregressive outputs in the real semiring is the same as the standard definition of a weighted DFA: when it is in state $q$, the next input symbol $\sigma$ determines both the next state $\delta(q, \sigma)$ as well as the symbol weight $a(q)(\sigma)$. Moreover, each state has an accepting weight $a(q)(\text{EOS})$.

In this paper, we are only interested in the following subclass of finite automata called **counter-free automata**, which we abbreviate as cfDFAs.

**Definition 5.2** (Counter-free automaton). *We say that a DFA with transition function $\delta$ is **counter-free** if there exists some $k$ such that for all states $q$, all strings $\boldsymbol{w}$, we have $\delta^*(q, \boldsymbol{w}^k) = \delta^*(q, \boldsymbol{w}^{k+1})$.*

Examples of counter-free and non-counter-free DFAs are shown in Fig. 2ab.

## 5.2 LINEAR TEMPORAL LOGIC

**Definition 5.3** (Linear temporal logic). *The formulas of past $\mathsf{LTL}$ are defined by the grammar*

$$
\begin{aligned}
\phi ::=\ & \neg \phi_1 \mid \phi_1 \wedge \phi_2 \\
& \mid \sigma && \sigma \in \Sigma \\
& \mid \textsc{bos} && \textit{Beginning of string} \\
& \mid \mathbf{Y} \phi_1 && \textit{Yesterday} \\
& \mid \mathbf{H} \phi_1 && \textit{Historically} \\
& \mid \phi_1 \, \mathbf{S} \, \phi_2 && \textit{Since}
\end{aligned}
$$

*Formulas $\top$ (true), $\bot$ (false), $\phi_1 \lor \phi_2$, $\phi_1 \leftrightarrow \phi_2$, and so on, can be defined as syntactic sugar in terms of the above. The temporal operator $\mathbf{P}\phi$ (which holds iff $\phi$ was Previously true at some time) can be defined as $\mathbf{H}\phi = \neg(\mathbf{P}(\neg\phi))$.*

*The semantics of formulas is given by the relation $\boldsymbol{w}, i \models \phi$ ("$\boldsymbol{w}$ satisfies $\phi$ at position $i$"), defined as follows:*

$$\boldsymbol{w}, i \models \neg\phi_1 \iff \boldsymbol{w}, i \not\models \phi_1 \tag{11a}$$

$$\boldsymbol{w}, i \models \phi_1 \land \phi_2 \iff \boldsymbol{w}, i \models \phi_1 \text{ and } \boldsymbol{w}, i \models \phi_2 \tag{11b}$$

$$\boldsymbol{w}, i \models \text{BOS} \iff i = 0 \tag{11c}$$

$$\boldsymbol{w}, i \models \sigma \iff w_i = \sigma \tag{11d}$$

$$\boldsymbol{w}, i \models \mathbf{Y}\phi_1 \iff i > 0 \text{ and } \boldsymbol{w}, i - 1 \models \phi_1 \tag{11e}$$

$$\boldsymbol{w}, i \models \mathbf{H}\phi_1 \iff \boldsymbol{w}, j \models \phi_1 \text{ for all } j \leq i \tag{11f}$$

$$\boldsymbol{w}, i \models \phi_1 \mathbf{S} \phi_2 \iff (\boldsymbol{w}, j \models \phi_2 \text{ for some } j \leq i) \text{ and } (\boldsymbol{w}, j' \models \phi_1 \text{ for all } j < j' \leq i). \tag{11g}$$

*We write $\boldsymbol{w} \models \phi$ as shorthand for $\boldsymbol{w}, |\boldsymbol{w}| \models \phi$.*

For any set of operators $\mathcal{O} \subseteq \{\mathbf{Y}, \mathbf{H}, \mathbf{S}\}$, we write $\mathsf{TL}[\mathcal{O}]$ for the set of formulas using only operators in $\mathcal{O}$. Thus past $\mathsf{LTL} = \mathsf{TL}[\mathbf{Y}, \mathbf{S}]$. Given a tuple of formulas $\Phi = (\phi_1, \ldots, \phi_m)$, we can define a state encoder

$$\Phi \colon \Sigma^* \to (\mathbb{B}^m)^*$$
$$\Phi(\boldsymbol{w})_i = (\mathbb{I}\{\boldsymbol{w}, i \models \phi_1\}, \ldots \mathbb{I}\{\boldsymbol{w}, i \models \phi_m\}).$$

Droste and Gastin (2019) define a weighted first-order logic, with several variations corresponding to several subclasses of weighted counter-free automata. Mandrali and Rahonis (2013; 2015) do the same for LTL. Both of these logics have, roughly speaking, four layers: (1) a core Boolean logic, (2) weights conditioned on formulas, (3) products over positions, and (4) addition and sums over positions. This is similar to our framework, which has (1) a core Boolean logic, (2) classifier output functions that can choose weights conditioned on formulas, and (3) autoregressive output functions that can also compute products over positions.

## 6 EXPRESSIVITY RESULTS

Previous results have shown that UHATs, LTL, and cfDFAs are equivalent in terms of language recognition. In §6.1, we use the results to show that these formalisms are also equivalent as weighted classifiers and as autoregressors.

Next, we compare the expressivity of classifier versus autoregressive models. Given the equivalence of the above formalisms, we will mainly discuss LTL. In the real semiring (§6.2), LTL classifiers define exactly the aperiodic step functions (defined below), which are less expressive than LTL autoregressors. And LTL autoregressors, in turn, are equivalent to counter-free DFA autoregressors and less expressive than weighted counter-free NFAs.

In the Boolean semiring, LTL classifiers and autoregressors are equivalent, which is the main result of §6.3.1. However, when we consider fragments of LTL, this equivalence breaks down, and autoregressors may become more expressive than classifiers (§6.3.2). Similarly, in the temporal logic with counting $\mathsf{TL}[\overleftarrow{\#}, +]$ and the programming language C-RASP (Yang and Chiang, 2024), autoregressors are more expressive than classifiers (§6.3.3).

### 6.1 STATE ENCODERS

We say that two state encoders $\tau_1 \colon \Sigma^* \to Q_1^*$ and $\tau_1 \colon \Sigma^* \to Q_2^*$ are **equivalent** if there is a bijection $f \colon Q_1 \to Q_2$ such that for all $\boldsymbol{w} \in \Sigma^*$, $f(\tau_1(\boldsymbol{w})) = \tau_2(\boldsymbol{w})$.

**Theorem 6.1.** *UHATs, LTL, and cfDFAs define equivalent state encoders.*

The proof is an adaptation of existing results (Yang et al., 2024; Schützenberger, 1965; McNaughton and Papert, 1971; Kamp, 1968) connecting UHATs, LTL and cfDFAs as language recognizers.

*Proof.* See §A. □

The following is an immediate consequence of Thm. 6.1 and the definitions of classifier and autoregressive models.

**Corollary 6.2.** *UHATs, LTL, and counter-free DFAs as classifier models define the same weighted languages. Similarly when they are used as autoregressive models.*

*Proof.* By the previous theorem, all these formalisms define equivalent state encoders. Therefore there exist output functions with which they define the same weighted languages. □

## 6.2 REAL CLASSIFIERS AND AUTOREGRESSORS

In this section, we consider weights in the real semiring. We characterize what weighted languages can be expressed, first by real classifiers, then by real autoregressors.

**Definition 6.1.** *An **aperiodic step function** (Droste and Gastin, 2008) is a weighted language $S\colon \Sigma^* \to \mathbb{K}$ such that $S(\boldsymbol{w}) = \bigoplus_{i=1}^{m} k_i \otimes \mathbb{I}\{\boldsymbol{w} \in L_i\}$ where $L_1, \ldots, L_m$ are aperiodic (that is, counter-free) regular languages.*

**Proposition 6.3.** *An LTL classifier defines the aperiodic step functions.*

*Proof.* Given any aperiodic step function as defined above, we can write, for each $L_i$, an LTL formula $\phi_i$. Then we can write a classifier output function $c(\mathbf{h}) = \bigoplus_{i=1}^{m} k_i \otimes h_i$.

Conversely, given an LTL classifier consisting of a tuple of formulas $(\phi_1, \ldots, \phi_m)$ and an output function $c(\mathbf{h})$, for every $\mathbf{h} \in 2^{[m]}$, write the formula $\phi_{\mathbf{h}} = \bigwedge_{i=1}^{m}(\phi_i \leftrightarrow h_i)$. For every $\mathbf{h}$, let $L_{\mathbf{h}}$ be the language defined by $\phi_{\mathbf{h}}$. Then the weighted language can be written as the step function $S(\boldsymbol{w}) = \bigoplus_{\mathbf{h} \in 2^{[m]}} c(\mathbf{h}) \otimes \mathbb{I}\{\boldsymbol{w} \in L_{\mathbf{h}}\}$. □

The following easy corollary of Prop. 6.3 shows that autoregressors are in general more powerful than classifiers.

**Corollary 6.4.** *In the real semiring, the weighted language $(\frac{1}{2}a)^*$ is expressible by an LTL (TL[$\mathbf{H}$]) autoregressor, but not by any LTL (TL[$\mathbf{H}$]) classifier. The language $(1a)^*$ is expressible by a LTL (TL[$\mathbf{H}$]) classifier but not any LTL (TL[$\mathbf{H}$]) autoregressor.*

*Proof.* The first language has an infinite number of string weights, but an aperiodic step function can only output a finite number of different weights. On the other hand, it is easy to write a weighted LTL (TL[$\mathbf{H}$]) formula with an output function to recognize this. The second language can easily be expressed by a classifier assigning weight 1 to all strings of $a$'s, but is not expressible by any autoregressor because it is not a normalized weighted language. □

As real autoregressors, LTL formulas are equivalent to counter-free DFAs by Cor. 6.2. However, both are less expressive than weighted counter-free **nondeterministic finite automata** (NFAs), in which a state can have more than one outgoing transition with the same symbol (see §C for a definition).

**Fact 6.5.** *Weighted counter-free NFAs define more weighted languages than counter-free DFA autoregressors do.*

Fig. 2c shows an example of a counter-free weighted NFA that is not determinizable. See §C for more details.

Note that this stands in contrast with the unweighted case, where counter-free DFAs and counter-free NFAs are equivalent. There are several nonequivalent analogues of counter-free automata (Droste and Gastin, 2008), and LTL and UHAT autoregressors are only equivalent to the least powerful of these.

6.3 BOOLEAN CLASSIFIERS AND AUTOREGRESSORS

To examine more carefully how autoregressors add expressivity, we turn to the Boolean semiring. We will see that LTL classifiers and LTL autoregressors are equivalent, but with an important caveat: with certain fragments and extension of LTL that use only a subset of the temporal operators, autoregressors can be more expressive than classifiers. These variants of LTL are particularly interesting because they have been proven to be equivalent to variants of transformers.

### 6.3.1 LTL

In the Boolean semiring, LTL classifiers and autoregressors are equivalent, but the conversion from an autoregressor to a classifier uses the $\mathbf{Y}$ and $\mathbf{H}$ operators.

**Theorem 6.6.** *For any set of operators $\mathcal{O} \subseteq \{\mathbf{Y}, \mathbf{H}, \mathbf{S}\}$:*

*(a) For any nonempty language $L$ defined by a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ classifier, there exists a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ autoregressor defining the same language $L$.*
*(b) For any language $L$ defined by a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ autoregressor, there exists a Boolean-weighted $\mathsf{TL}[\mathcal{O} \cup \{\mathbf{Y}, \mathbf{H}\}]$ classifier defining the same language $L$.*

*Proof.* See §B.3. □

To prove this, we need to introduce two new operators as "syntactic sugar" that do not increase the expressivity of the logic.

**Lemma 6.7.** *There is a transformation $\mathrm{next}_\sigma$ from formulas of $\mathsf{TL}[\mathcal{O}]$ to formulas of $\mathsf{TL}[\mathcal{O}]$ such that for any formula $\phi$ of $\mathsf{TL}[\mathcal{O}]$ and for all $\boldsymbol{w} \in \Sigma^*$,*

$$\boldsymbol{w} \models \mathrm{next}_\sigma(\phi) \iff \boldsymbol{w}\sigma \models \phi. \tag{12}$$

Intuitively, $\mathrm{next}_\sigma$ removes a $\sigma$ on the right; in other words, $\mathrm{next}_\sigma(\phi)$ defines the right Brzozowski derivative (Brzozowski, 1964) of the language defined by $\phi$.

*Proof.* See §B.1. □

**Lemma 6.8.** *There is a transformation $\mathrm{prefix}$ from formulas of $\mathsf{TL}[\mathcal{O}]$ to formulas of $\mathsf{TL}[\mathcal{O}]$ such that for any formula $\phi$ of $\mathsf{TL}[\mathcal{O}]$ and for all $\boldsymbol{u} \in \Sigma^*$,*

$$\boldsymbol{u} \models \mathrm{prefix}(\phi) \iff \text{there exists } \boldsymbol{v} \in \Sigma^* \text{ such that } \boldsymbol{u}\boldsymbol{v} \models \phi. \tag{13}$$

*Proof.* See §B.2. □

From Thm. 6.6, we can conclude that for UHATs, which are equivalent to LTL, autoregression does not add any expressivity. On other transformer variants, please see §6.3.2.

The construction in the proof of Lem. 6.8 yields a formula $\mathrm{prefix}(\phi)$ whose size is exponential in that of $\phi$. To shed light on whether this bound is tight, we show the following.

**Proposition 6.9.** *(a) There does not exist a transformation $\mathrm{prefix}'$ such that $\mathrm{prefix}'(\phi)$ is constructible in polynomial time (in $|\phi|$) and satisfies Eq. (13) for every formula $\phi$ in $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$, unless $\mathsf{P} = \mathsf{PSPACE}$.*
*(b) Similarly for $\mathsf{TL}[\mathbf{H}]$, unless $\mathsf{P} = \mathsf{NP}$.*
*(c) Similarly for $\mathsf{TL}[\mathbf{Y}]$, unless $\mathsf{P} = \mathsf{NP}$.*

*Proof.* See §B.4. □

Note that we have only shown (conditionally) that constructing $\mathrm{prefix}'(\phi)$ requires super-polynomial time; it's possible that $\mathrm{prefix}'(\phi)$ is short but difficult to construct.

### 6.3.2 Fragments of LTL

Li and Cotterell (2025) show that fixed-precision future-masked transformers are equivalent to $\mathsf{TL}[\mathbf{P}]$, which is in turn equivalent to $\mathsf{TL}[\mathbf{H}]$. Similarly, Jerad et al. (2025) show that future-masked leftmost-hard attention transformers are also equivalent to $\mathsf{TL}[\mathbf{P}]$. However, in this section we show that for autoregressors, these equivalences break.

When the set of operators $\mathcal{O}$ lacks either $\mathbf{H}$ or $\mathbf{Y}$, the asymmetry in Thm. 6.6 suggests that Boolean autoregressors are more expressive than classifiers. The following shows that this is indeed the case.

**Proposition 6.10.** *The language* $(ab)^*$ *is defined by a Boolean* $\mathsf{TL}[\emptyset]$ *autoregressor but not defined by any* $\mathsf{TL}[\mathbf{H}]$ *or* $\mathsf{TL}[\mathbf{Y}]$ *classifier.*

*Proof.* Consider the state encoder $\Phi = (\text{BOS}, a, b)$ and the output function

$$a(q_{\text{BOS}}, q_a, q_b)(a) = \top \iff q_{\text{BOS}} = \top \text{ or } q_b = \top$$
$$a(q_{\text{BOS}}, q_a, q_b)(b) = \top \iff q_a = \top$$
$$a(q_{\text{BOS}}, q_a, q_b)(\text{EOS}) = \top \iff q_{\text{BOS}} = \top \text{ or } q_b = \top$$

But a formula in $\mathsf{TL}[\mathbf{Y}]$ can't distinguish between strings that differ beyond their last $k$ symbols (for some constant $k$ depending on the formula), and for any $k$, we have $ab(ab)^{\lceil k/2 \rceil} \in (ab)^*$ but $ba(ab)^{\lceil k/2 \rceil} \notin (ab)^*$. A formula in $\mathsf{TL}[\mathbf{H}]$ is equivalent to one in $\mathsf{TL}[\mathbf{P}]$, which can only define a stutter-invariant language (that is, a language $L$ such that for all $\boldsymbol{u}, \sigma, \boldsymbol{v}$, we have $\boldsymbol{u}\sigma\boldsymbol{v} \in L \iff \boldsymbol{u}\sigma\sigma\boldsymbol{v} \in L$ (Peled and Wilke, 1997)). And $(ab)^*$ is not stutter-invariant, because $ab \in (ab)^*$ but $aabb \notin (ab)^*$. $\qquad\square$

Consequently, $(ab)^*$ is definable by leftmost-hard UHATs and fixed-precision SMATs as autoregressors, but not as classifiers. However, the expressiveness added by autoregression remains limited, as $(aab)^*$ is not definable.

**Proposition 6.11.** *The language* $(aab)^*$ *is not definable by any* $\mathsf{TL}[\mathbf{H}]$ *classifier or autoregressor.*

*Proof.* See §D. $\qquad\square$

Consequently, $(aab)^*$ is not definable by any leftmost-hard UHAT or fixed-precision SMAT, either as autoregressors or classifiers.

### 6.3.3 Temporal Logic with Counting

Other formalisms besides the ones discussed above have been proposed for comparison with transformers. Yang et al. (2025) prove that SMATs, with fixed precision outside attention and arbitrary precision inside attention, are equivalent to a temporal logic with counting operators, $\mathsf{TL}[\overleftarrow{\#}, +]$. They considered the family of languages

$$L_1 = a^* \tag{14}$$

$$L_{k+1} = \begin{cases} L_k b^* & k \text{ even} \\ L_k a^* & k \text{ odd} \end{cases} \tag{15}$$

and showed that, as Boolean classifiers, transformers with depth $k$ can recognize $L_k$ (and not $L_{k+1}$). But their experiments were on the symbol-prediction task (§2), closely related to Boolean autoregression. They showed both theoretically and experimentally that SMATs with depth $k$ can solve the symbol-prediction task for not only $L_k$, but $L_{k+2}$ (and not $L_{k+3}$).

In the present framework, this discrepancy can be readily explained. Like $\mathsf{TL}[\mathbf{H}]$, the logic $\mathsf{TL}[\overleftarrow{\#}, +]$ lacks a $\mathbf{Y}$ operator or an equivalent. So it is more expressive as an autoregressor than as a classifier.

## 7 CONCLUSION

We have shown that theoretical results on transformers as Boolean classifiers (as most theoretical results in the literature are) sometimes carry over to real-weighted and/or autoregressive settings, but they sometimes do not. We have laid out a framework for studying other variants of transformers and other automata or logics as real-weighted autoregressors, leading to theoretical results that can make more accurate predictions about language models as used in practice.

## ACKNOWLEDGEMENTS

## REFERENCES

Satwik Bhattamishra, Kabir Ahuja, and Navin Goyal. On the ability and limitations of Transformers to recognize formal languages. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7096–7116, 2020. doi:10.18653/v1/2020.emnlp-main.576.

Janusz A. Brzozowski. Derivatives of regular expressions. *Journal of the Association for Computing Machinery*, 11(4):481–494, October 1964. doi:10.1145/321239.321249.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional Transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL HLT)*, pages 4171–4186, 2019. doi:10.18653/v1/N19-1423.

Manfred Droste and Paul Gastin. On aperiodic and star-free formal power series in partially commuting variables. *Theory of Computing Systems*, 42(4):608–631, 2008. doi:10.1007/s00224-007-9064-z.

Manfred Droste and Paul Gastin. Aperiodic weighted automata and weighted first-order logic. In *Proceedings of the 44th International Symposium on Mathematical Foundations of Computer Science*, 2019. doi:10.4230/LIPICS.MFCS.2019.76.

Manfred Droste and Dietrich Kuske. Weighted automata. In Jean-Éric Pin, editor, *Handbook of Automata Theory*, pages 113–150. European Mathematical Society Publishing House, Zürich, Switzerland, 2021. doi:10.4171/AUTOMATA-1/4.

Valeria Fionda and Gianluigi Greco. The complexity of LTL on finite traces: Hard and easy fragments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30, pages 971–977, 2016. doi:10.1609/aaai.v30i1.10104.

Giuseppe De Giacomo and Moshe Y. Vardi. Linear temporal logic and linear dynamic logic on finite traces. In *Proceedings of the 23rd International Joint Conference on Artificial Intelligence (IJCAI)*, pages 854–860, 2013. URL https://www.ijcai.org/Proceedings/13/Papers/132.pdf.

Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020. doi:10.1162/tacl_a_00306.

Xinting Huang, Andy Yang, Satwik Bhattamishra, Yash Sarrof, Andreas Krebs, Hattie Zhou, Preetum Nakkiran, and Michael Hahn. A formal framework for understanding length generalization in transformers. In *Proceedings of the Thirteenth International Conference on Learning Representations (ICLR)*, 2025. URL https://openreview.net/forum?id=U49N5V51rU.

Selim Jerad, Anej Svete, Jiaoda Li, and Ryan Cotterell. Unique hard attention: A tale of two sides. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 977–996, 2025. doi:10.18653/v1/2025.acl-short.76.

Johan Anthony Willem Kamp. *Tense Logic and the Theory of Linear Order*. PhD thesis, University of California, Los Angeles, 1968. URL https://www.proquest.com/docview/302320357.

Jiaoda Li and Ryan Cotterell. Characterizing the expressivity of transformer language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025. URL https://arxiv.org/abs/2505.23623. To appear.

Eleni Mandrali and George Rahonis. Characterizations of weighted first-order logics over semirings. In *Algebraic Informatics: 5th International Conference (CAI)*, pages 247–259, 2013. doi:10.1007/978-3-642-40663-8_23.

Eleni Mandrali and George Rahonis. Weighted first-order logics over semirings. *Acta Cybernetica*, 22(2):435–483, 2015. doi:10.14232/actacyb.22.2.2015.1.

Robert McNaughton and Seymour Papert. *Counter-Free Automata*. Number 65 in M.I.T. Press Research Monographs. M.I.T. Press, 1971. URL https://archive.org/details/CounterFre_00_McNa.

Mehryar Mohri. Finite-state transducers in language and speech processing. *Computational Linguistics*, 23(2):269–311, 1997. URL https://aclanthology.org/J97-2003/.

Doron Peled and Thomas Wilke. Stutter-invariant temporal properties are expressible without the next-time operator. *Information Processing Letters*, 63(5):243–246, 1997. doi:10.1016/S0020-0190(97)00133-6.

M. P. Schützenberger. On finite monoids having only trivial subgroups. *Information and Control*, 8 (2):190–194, 1965. doi:10.1016/S0019-9958(65)90108-7.

Lena Strobl, William Merrill, Gail Weiss, David Chiang, and Dana Angluin. What formal languages can transformers express? A survey. *Transactions of the Association for Computational Linguistics*, 12:543–561, 2024. doi:10.1162/tacl_a_00663.

Anej Svete and Ryan Cotterell. Transformers can represent $n$-gram language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 6845–6881, 2024. doi:10.18653/v1/2024.naacl-long.381.

Andy Yang and David Chiang. Counting like transformers: Compiling temporal counting logic into softmax transformers. In *Proceedings of the First Conference on Language Modeling (CoLM)*, 2024. URL https://openreview.net/forum?id=FmhPg4UJ9K.

Andy Yang, David Chiang, and Dana Angluin. Masked hard-attention transformers recognize exactly the star-free languages. In *Advances in Neural Information Processing Systems*, volume 37, pages 10202–10235, 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/hash/13d7f172259b11b230cc5da8768abc5f-Abstract-Conference.html.

Andy Yang, Michaël Cadilhac, and David Chiang. Knee-deep in C-RASP: A transformer depth hierarchy. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2025. URL https://arxiv.org/abs/2506.16055. To appear.

Shunyu Yao, Binghui Peng, Christos Papadimitriou, and Karthik Narasimhan. Self-attention networks can process bounded hierarchical languages. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP)*, pages 3770–3785, 2021. doi:10.18653/v1/2021.acl-long.292.

# A    EQUIVALENCE OF STATE ENCODERS

**Theorem 6.1.** *UHATs,* LTL*, and cfDFAs define equivalent state encoders.*

*Proof.* First we show the equivalence of state sequences defined by UHATs and LTL, and then equivalence of LTL and cfDFAs.

The essential observation (Yang et al., 2024, Lemma 22) is that the output at every position of every UHAT layer comes from a finite set $Q \subseteq \mathbb{R}^d$. So we can think of a UHAT as a function $\mathcal{T} \colon \Sigma^* \to Q^*$. For each $q \in Q$, we can construct an LTL formula $\phi_q$ such that $\mathcal{T}(\boldsymbol{w})_i = q \iff \boldsymbol{w}, i \models \phi_q$ (Yang et al., 2024, Theorems 2, 4). So there exists a tuple of LTL formulas $(\phi_q)_{q \in Q}$ that defines a state encoder equivalent to $\mathcal{T}$. Note that the state outputted by $\mathcal{T}$ on the prepended BOS symbol can be simulated using a BOS formula in the tuple.

In the other direction, for every tuple of LTL formulas $(\phi_1, \phi_2, \ldots, \phi_m)$ defining a state encoder $\Sigma^* \to \mathbb{B}^m$, there exists a UHAT $\mathcal{T} \colon \Sigma^* \to (\mathbb{R}^d)^*$ defining an equivalent state encoder. For each $\phi_k$, we construct a transformer $\mathcal{T}_k$ which outputs $\frac{1}{2}$ if $\boldsymbol{w}, i \models \phi_k$ and $-\frac{1}{2}$ otherwise (Yang et al., 2024, Theorems 1, 3). Then we can parallel-compose all the $\mathcal{T}_k$ into a single $\mathcal{T}$ (Yang et al., 2024, Lemma 25), and add an additional layer which projects the output dimensions of each $\mathcal{T}_k$ into a single output vector $\mathbb{R}^m$ such that $\mathcal{T}(\boldsymbol{w})_i = \mathbf{e}_k \iff \boldsymbol{w}, i \models \phi_k$.

The equivalence between LTL and cfDFAs can be described a little more succinctly. Given a DFA $M = (\Sigma, Q, \delta, \iota)$, for each state $q \in Q$ there exists a formula $\phi_q$ such that $\boldsymbol{w} \models \phi_q \iff \delta(\iota, \boldsymbol{w}) = q$, due to the expressive equivalence of LTL and cfDFAs (Schützenberger, 1965; McNaughton and Papert, 1971; Kamp, 1968). The tuple $(\phi_q)_{q \in Q}$ then defines a state encoder equivalent to $M$. In the other direction, given a tuple of LTL formulas $(\phi_1, \ldots, \phi_m)$, for each $k \in [m]$ there is an automaton $M_k$ that recognizes the same language as $\phi_k$. Then the Cartesian product of all the $M_k$ defines a state encoder equivalent to $(\phi_1, \ldots, \phi_m)$. $\qquad\square$

# B    AUTOREGRESSIVE MODEL PROOFS

## B.1    PROOF OF LEM. 6.7

**Lemma 6.7.** *There is a transformation* $\mathrm{next}_\sigma$ *from formulas of* $\mathsf{TL}[\mathcal{O}]$ *to formulas of* $\mathsf{TL}[\mathcal{O}]$ *such that for any formula $\phi$ of* $\mathsf{TL}[\mathcal{O}]$ *and for all $\boldsymbol{w} \in \Sigma^*$,*

$$\boldsymbol{w} \models \mathrm{next}_\sigma(\phi) \iff \boldsymbol{w}\sigma \models \phi. \tag{12}$$

*Proof.* We define $\mathrm{next}_\sigma$ recursively:

$$\mathrm{next}_\sigma(\sigma) = \top \tag{16a}$$

$$\mathrm{next}_\sigma(\sigma') = \bot \qquad \text{if } \sigma' \neq \sigma \tag{16b}$$

$$\mathrm{next}_\sigma(\text{BOS}) = \bot \tag{16c}$$

$$\mathrm{next}_\sigma(\neg\phi) = \neg\mathrm{next}_\sigma(\phi) \tag{16d}$$

$$\mathrm{next}_\sigma(\phi_1 \wedge \phi_2) = \mathrm{next}_\sigma(\phi_1) \wedge \mathrm{next}_\sigma(\phi_2) \tag{16e}$$

$$\mathrm{next}_\sigma(\mathbf{Y}\phi) = \phi \tag{16f}$$

$$\mathrm{next}_\sigma(\mathbf{H}\phi) = \mathbf{H}\phi \wedge \mathrm{next}_\sigma(\phi) \tag{16g}$$

$$\mathrm{next}_\sigma(\phi_1 \mathbf{S} \phi_2) = (\mathrm{next}_\sigma(\phi_1) \wedge (\phi_1 \mathbf{S} \phi_2)) \vee \mathrm{next}_\sigma(\phi_2). \tag{16h}$$

Note that $\mathrm{next}_\sigma$ never translates a temporal operator into another temporal operator, so it translates formulas of $\mathsf{TL}[\mathcal{O}]$ into formulas of $\mathsf{TL}[\mathcal{O}]$ for any $\mathcal{O}$.

Next, we prove that $\mathrm{next}_\sigma(\phi)$ satisfies Eq. (12) by induction on the structure of $\phi$.

**Base Cases.**    If $\phi = \sigma$:

$$\boldsymbol{w}, i \models \mathrm{next}_\sigma(\sigma) \overset{(16a)}{\iff} \boldsymbol{w} \models \top \tag{17a}$$

$$\overset{(11d)}{\iff} \boldsymbol{w}\sigma \models \sigma. \tag{17b}$$

If $\phi = \sigma'$ for $\sigma' \neq \sigma$:

$$w \models \text{next}_\sigma(\sigma') \overset{(16b)}{\Longleftrightarrow} w \models \bot \tag{17c}$$

$$\overset{(11d)}{\Longleftrightarrow} w\sigma \models \sigma'. \tag{17d}$$

Similarly, if $\phi = \text{BOS}$:

$$w \models \text{next}_\sigma(\text{BOS}) \overset{(16c)}{\Longleftrightarrow} w \models \bot \tag{17e}$$

$$\overset{(11c)}{\Longleftrightarrow} w\sigma \models \text{BOS}. \tag{17f}$$

**Inductive Cases.** If $\phi = \neg\phi_1$:

$$w \models \text{next}_\sigma(\neg\phi_1) \overset{(16d)}{\Longleftrightarrow} w \models \neg\text{next}_\sigma(\phi_1) \tag{18a}$$

$$\overset{(11a)}{\Longleftrightarrow} w \not\models \text{next}_\sigma(\phi_1) \tag{18b}$$

$$\overset{\text{ind. hyp.}}{\Longleftrightarrow} w\sigma \not\models \phi_1 \tag{18c}$$

$$\overset{(11a)}{\Longleftrightarrow} w\sigma \models \neg\phi_1. \tag{18d}$$

If $\phi = \phi_1 \wedge \phi_2$:

$$w \models \text{next}_\sigma(\phi_1 \wedge \phi_2) \overset{(16e)}{\Longleftrightarrow} w \models \text{next}_\sigma(\phi_1) \wedge \text{next}_\sigma(\phi_2) \tag{18e}$$

$$\overset{(11b)}{\Longleftrightarrow} (w \models \text{next}_\sigma(\phi_1)) \wedge (w \models \text{next}_\sigma(\phi_2)) \tag{18f}$$

$$\overset{\text{ind. hyp.}}{\Longleftrightarrow} (w\sigma \models \phi_1) \wedge (w\sigma \models \phi_2) \tag{18g}$$

$$\overset{(11b)}{\Longleftrightarrow} w\sigma \models \phi_1 \wedge \phi_2. \tag{18h}$$

If $\phi = \mathbf{Y}\phi_1$:

$$w \models \text{next}_\sigma(\mathbf{Y}\phi_1) \overset{(16f)}{\Longleftrightarrow} w \models \phi_1 \tag{18i}$$

$$\overset{(11e)}{\Longleftrightarrow} w\sigma \models \mathbf{Y}\phi_1. \tag{18j}$$

If $\phi = \mathbf{H}\phi_1$:

$$w \models \text{next}_\sigma(\mathbf{H}\phi_1) \overset{(16g)}{\Longleftrightarrow} w \models \mathbf{H}\phi_1 \wedge \text{next}_\sigma(\phi) \tag{18k}$$

$$\overset{(11b)}{\Longleftrightarrow} (w \models \mathbf{H}\phi_1) \wedge (w \models \text{next}_\sigma(\phi)) \tag{18l}$$

$$\overset{\text{ind. hyp.}}{\Longleftrightarrow} (w \models \mathbf{H}\phi_1) \wedge (w\sigma \models \phi_1) \tag{18m}$$

$$\overset{(11f)}{\Longleftrightarrow} w\sigma \models \mathbf{H}\phi_1. \tag{18n}$$

If $\phi = \phi_1 \mathbf{S} \phi_2$:

$$w \models \text{next}_\sigma(\phi_1 \mathbf{S} \phi_2)$$

$$\overset{(16h)}{\Longleftrightarrow} w \models (\text{next}_\sigma(\phi_1) \wedge (\phi_1 \mathbf{S} \phi_2)) \vee \text{next}_\sigma(\phi_2) \tag{18o}$$

$$\overset{(11a) \text{ and } (11b)}{\Longleftrightarrow} (w \models \text{next}_\sigma(\phi_1) \wedge (w \models \phi_1 \mathbf{S} \phi_2)) \vee (w \models \text{next}_\sigma(\phi_2)) \tag{18p}$$

$$\overset{\text{ind. hyp.}}{\Longleftrightarrow} ((w\sigma \models \phi_1) \wedge (w \models \phi_1 \mathbf{S} \phi_2)) \vee (w\sigma \models \phi_2) \tag{18q}$$

$$\overset{(11g)}{\Longleftrightarrow} w\sigma \models \phi_1 \mathbf{S} \phi_2. \tag{18r}$$

$$\square$$

### B.2 PROOF OF LEM. 6.8

**Lemma 6.8.** *There is a transformation* prefix *from formulas of* $\mathsf{TL}[\mathcal{O}]$ *to formulas of* $\mathsf{TL}[\mathcal{O}]$ *such that for any formula $\phi$ of* $\mathsf{TL}[\mathcal{O}]$ *and for all $\boldsymbol{u} \in \Sigma^*$,*

$$\boldsymbol{u} \models \text{prefix}(\phi) \iff \text{there exists } \boldsymbol{v} \in \Sigma^* \text{ such that } \boldsymbol{uv} \models \phi. \tag{13}$$

*Proof.* Given a formula $\phi$ of $\mathsf{TL}[\mathcal{O}]$, let $\text{cl}(\phi)$ be the set of all subformulas of $\phi$ (including $\phi$ itself). Construct a DFA $M_\phi = (2^{\text{cl}(\phi)}, \Sigma, \delta, \iota, F)$, where

$$\iota = \{\chi \in \text{cl}(\phi) \mid \epsilon \models \chi\} \tag{19}$$

$$F = \{\Psi \subseteq \text{cl}(\phi) \mid \phi \in \Psi\} \tag{20}$$

$$\delta(\Psi, \sigma) = \{\chi \in \text{cl}(\phi) \mid \Psi \overset{\sigma}{\to} \chi\} \tag{21}$$

where the relation $\Psi \xrightarrow{\sigma} \chi$, which intuitively means that if a string $w$ satisfies exactly the formulas in $\Psi$, then $w\sigma$ satisfies $\chi$, is defined as follows:

$$\Psi \xrightarrow{\sigma} \sigma' \text{ iff } \sigma = \sigma' \tag{22a}$$

$$\Psi \xrightarrow{\sigma} \chi_1 \wedge \chi_2 \text{ iff } \Psi \xrightarrow{\sigma} \chi_1 \text{ and } \Psi \xrightarrow{\sigma} \chi_2 \tag{22b}$$

$$\Psi \xrightarrow{\sigma} \neg\chi \text{ iff not } \Psi \xrightarrow{\sigma} \chi \tag{22c}$$

$$\Psi \xrightarrow{\sigma} \mathbf{Y}\chi \text{ iff } \chi \in \Psi \tag{22d}$$

$$\Psi \xrightarrow{\sigma} \mathbf{H}\chi \text{ iff } \mathbf{H}\chi \in \Psi \text{ and } \Psi \xrightarrow{\sigma} \chi \tag{22e}$$

$$\Psi \xrightarrow{\sigma} \chi_1 \, \mathbf{S} \, \chi_2 \text{ iff } (\chi_1 \, \mathbf{S} \, \chi_2 \in \Psi \text{ and } \Psi \xrightarrow{\sigma} \chi_1) \text{ or } \Psi \xrightarrow{\sigma} \chi_2. \tag{22f}$$

**Claim B.1.** *For any $w \in \Sigma^*$, if $\Psi = \{\chi \in \mathrm{cl}(\phi) \mid w \models \chi\}$, then $\Psi \xrightarrow{\sigma} \chi \iff w\sigma \models \chi$.*

*Proof.* By induction on the structure of $\chi$. Note by definition that $\chi \in \Psi \iff w \models \chi$.

$$\Psi \xrightarrow{\sigma} \sigma' \overset{(22a)}{\iff} \sigma = \sigma'$$
$$\overset{(11d)}{\iff} w\sigma \models \sigma'.$$

$$\Psi \xrightarrow{\sigma} \chi_1 \wedge \chi_2 \overset{(22b)}{\iff} \Psi \xrightarrow{\sigma} \chi_1 \text{ and } \Psi \xrightarrow{\sigma} \chi_2$$
$$\overset{\text{ind. hyp.}}{\iff} w\sigma \models \chi_1 \text{ and } w\sigma \models \chi_2$$
$$\overset{(11b)}{\iff} w\sigma \models \chi_1 \wedge \chi_2.$$

$$\Psi \xrightarrow{\sigma} \neg\chi \overset{(22c)}{\iff} \text{not } \Psi \xrightarrow{\sigma} \chi$$
$$\overset{\text{ind. hyp.}}{\iff} \text{not } w\sigma \models \chi$$
$$\overset{(11a)}{\iff} w\sigma \models \neg\chi.$$

$$\Psi \xrightarrow{\sigma} \mathbf{Y}\chi \overset{(22d)}{\iff} \chi \in \Psi$$
$$\iff w \models \chi$$
$$\overset{(11e)}{\iff} w\sigma \models \mathbf{Y}\chi.$$

$$\Psi \xrightarrow{\sigma} \mathbf{H}\chi \overset{(22e)}{\iff} \mathbf{H}\chi \in \Psi \text{ and } \Psi \xrightarrow{\sigma} \chi$$
$$\overset{\text{ind. hyp.}}{\iff} w \models \mathbf{H}\chi \text{ and } w\sigma \models \chi$$
$$\overset{(11f)}{\iff} w\sigma \models \mathbf{H}.$$

$$\Psi \xrightarrow{\sigma} \chi_1 \, \mathbf{S} \, \chi_2 \overset{(22f)}{\iff} (\chi_1 \, \mathbf{S} \, \chi_2 \in \Psi \text{ and } \Psi \xrightarrow{\sigma} \chi_1) \text{ or } \Psi \xrightarrow{\sigma} \chi_2$$
$$\overset{\text{ind. hyp.}}{\iff} (w \models \chi_1 \, \mathbf{S} \, \chi_2 \text{ and } w\sigma \models \chi_1) \text{ or } w\sigma \models \chi_2$$
$$\overset{(11g)}{\iff} w\sigma \models \chi_1 \, \mathbf{S} \, \chi_2. \qquad \square$$

**Claim B.2.** *For any $w$, $\delta(\iota, w) = \{\chi \in \mathrm{cl}(\phi) \mid w \models \chi\}$.*

*Proof.* By induction on the length of $w$.

Base case: $\delta(\iota, \epsilon) = \iota = \{\chi \mid \epsilon \models \chi\}$.

Inductive step: Assume that $\delta(\iota, w) = \{\chi \mid w \models \chi\} = \Psi$. Then

$$\delta(\iota, w) = \delta(\delta(\iota, w), \sigma)$$
$$= \delta(\Psi, \sigma)$$
$$= \{\chi \mid \Psi \xrightarrow{\sigma} \chi\}$$
$$= \{\chi \mid w\sigma \models \chi\}. \qquad \square$$

**Claim B.3.** *$M_\phi$ defines the same language as $\phi$.*

*Proof.* $\delta(\iota, w) \in F$ if and only if $\phi \in \{\chi \mid w \models \chi\}$ if and only if $w \models \phi$. $\qquad \square$

Then make every co-accessible state (every state that has a path to an accept state) into an accept state. Call this new DFA $M'_\phi$ with accept states $F'$. This DFA recognizes the prefix language of $M_\phi$. Finally, construct the formula

$$\mathrm{prefix}(\phi) = \bigvee_{\Psi \in F'} \left( \bigwedge_{\chi \in \Psi} \chi \wedge \bigwedge_{\chi \in \mathrm{cl}(\phi) \setminus \Psi} \neg \chi \right).$$

Note that $\mathrm{prefix}$ never translates a temporal operator into another temporal operator, so it translates formulas of $\mathsf{TL}[\mathcal{O}]$ into formulas of $\mathsf{TL}[\mathcal{O}]$ for any $\mathcal{O}$.

**Claim B.4.** *The formula* $\mathrm{prefix}(\phi)$ *defines the same language as* $M'_\phi$.

*Proof.* Since we only changed non-accept states to accept states, Clm. B.2 still applies to $M'_\phi$ and $\phi$.

$$
\begin{aligned}
\boldsymbol{w} \in \mathcal{L}(M'_\phi) &\iff \delta(\iota, \boldsymbol{w}) \in F' \\
&\iff \{\chi \in \mathrm{cl}(\phi) \mid \boldsymbol{w} \models \chi\} \in F' \qquad\qquad \text{Clm. B.2} \\
&\iff \text{for some } \Psi \in F', \chi \in \Psi \text{ iff } \boldsymbol{w} \models \chi \\
&\iff \text{for some } \Psi \in F', \boldsymbol{w} \models \bigwedge_{\chi \in \Psi} \chi \wedge \bigwedge_{\chi \in \mathrm{cl}(\phi) \setminus \Psi} \neg \chi \\
&\iff \boldsymbol{w} \models \bigvee_{\Psi \in F'} \left( \bigwedge_{\chi \in \Psi} \chi \wedge \bigwedge_{\chi \in \mathrm{cl}(\phi) \setminus \Psi} \neg \chi \right). \qquad\qquad \square
\end{aligned}
$$

This completes the proof of Lem. 6.8. $\qquad\square$

### B.3 Relationship Between Classifiers and Autoregressors

**Theorem 6.6.** *For any set of operators* $\mathcal{O} \subseteq \{\mathbf{Y}, \mathbf{H}, \mathbf{S}\}$:

(a) *For any nonempty language $L$ defined by a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ classifier, there exists a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ autoregressor defining the same language $L$.*
(b) *For any language $L$ defined by a Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ autoregressor, there exists a Boolean-weighted $\mathsf{TL}[\mathcal{O} \cup \{\mathbf{Y}, \mathbf{H}\}]$ classifier defining the same language $L$.*

*Proof.* (a) A Boolean-weighted $\mathsf{TL}[\mathcal{O}]$ classifier is defined by a tuple of formulas $\Phi = (\phi_1, \ldots, \phi_m)$ and an output function $c \colon \mathbb{B}^m \to \mathbb{B}$. We may think of $c$ as a Boolean combination of its arguments, and substitute the $\phi_i$ into it to obtain a single formula $\phi = c(\phi_1, \ldots, \phi_m)$. Then define a new trivial output function $c'(h) = h$, so that the classifier $(\phi, c')$ defines the same language as $(\Phi, c)$.

Define

$$
\begin{aligned}
\phi_\sigma &= \mathrm{next}_\sigma(\mathrm{prefix}(\phi)) \qquad \text{for } \sigma \in \Sigma && (23) \\
\phi_{\mathrm{EOS}} &= \phi. && (24)
\end{aligned}
$$

Then the tuple of formulas $\Phi' = (\phi_\sigma)_{\sigma \in \Sigma \cup \{\mathrm{EOS}\}}$ defines a state encoder. We define the autoregressive output function

$$a \colon \mathbb{B}^{|\Sigma|+1} \to \mathbb{B}^{|\Sigma|+1} \tag{25}$$

$$a(\mathbf{h})(\sigma) = \begin{cases} h_\sigma & \text{if any entry of } \mathbf{h} \text{ is true} \\ \top & \text{if all entries of } \mathbf{h} \text{ are false.} \end{cases} \tag{26}$$

A vector $\mathbf{h}$ whose entries are all false is unreachable, so it does not matter what we set $a(\mathbf{h})$ to, but it must satisfy $\bigoplus_\sigma a(\mathbf{h})(\sigma) = \mathbf{1}$ (Eq. (7)), that is, the entries of $a(\mathbf{h})$ cannot all be false.

The autoregressor $A = (\Phi', a)$ defines $L$, because for any $\boldsymbol{w} \in L$ with length $n$, we have

$$\Pr_A(\boldsymbol{w}) = \bigotimes_{i=1}^{n} \Pr_A(w_i \mid \boldsymbol{w}_{<i}) \otimes \Pr_A(\text{EOS} \mid \boldsymbol{w}) \tag{27}$$

$$= \bigwedge_{i=1}^{n} \mathbb{I}\{\boldsymbol{w}_{<i} \models \text{next}_{w_i}(\text{prefix}(\phi)\} \wedge \mathbb{I}\{\boldsymbol{w} \models \phi\} \tag{28}$$

$$= \bigwedge_{i=1}^{n} \mathbb{I}\{\boldsymbol{w}_{\leq i} \models \text{prefix}(\phi)\} \wedge \mathbb{I}\{\boldsymbol{w} \models \phi\} \tag{29}$$

$$= \bigwedge_{i=1}^{n} \mathbb{I}\{\boldsymbol{w}_{\leq i}\boldsymbol{v} \models \phi \text{ for some } \boldsymbol{v}\} \wedge \mathbb{I}\{\boldsymbol{w} \models \phi\} \tag{30}$$

$$= \top. \tag{31}$$

On the other hand, for any $\boldsymbol{w} \notin L$, let $k$ be the greatest integer such that $\boldsymbol{w}_{<k}\boldsymbol{v} \in L$ for some $\boldsymbol{v}$. (We know that $k$ exists because $L$ is nonempty by assumption.) Then we have that $\boldsymbol{w}_{<k} \not\models \text{next}_{w_k}(\text{prefix}(\phi))$, but $\boldsymbol{w}_{<k} \models \text{next}_{v_1}(\text{prefix}(\phi))$. So $\Pr_A(w_k \mid \boldsymbol{w}_{<k}) = \bot$, and therefore $\Pr_A(\boldsymbol{w}) = \bot$.

It remains to verify that $A$ satisfies Eq. (6). For any $\boldsymbol{u} \in \Sigma^*$, we want to show that

$$\bigoplus_{\boldsymbol{v}} \Pr_A(\boldsymbol{v} \mid \boldsymbol{u}) = \bigoplus_{\boldsymbol{v}} \left( \bigotimes_{i=1}^{|\boldsymbol{v}|} \Pr_A(v_i \mid \boldsymbol{u}\boldsymbol{v}_{<i}) \otimes \Pr_A(\text{EOS} \mid \boldsymbol{u}\boldsymbol{v}) \right) = \boldsymbol{1}. \tag{32}$$

In the Boolean semiring, it suffices to show that at least one term of this summation is true.

If there is a $\boldsymbol{v}$ such that $\boldsymbol{u}\boldsymbol{v} \in L$, then the corresponding term of the summation is

$$\left( \bigwedge_{i=1}^{|\boldsymbol{u}\boldsymbol{v}|} \mathbb{I}\{\boldsymbol{u}\boldsymbol{v}_{<i} \models \text{next}_{v_i}(\text{prefix}(\phi))\} \right) \wedge \mathbb{I}\{\boldsymbol{u}\boldsymbol{v} \models \phi\} \tag{33}$$

$$= \left( \bigwedge_{i=1}^{|\boldsymbol{u}\boldsymbol{v}|} \mathbb{I}\{\boldsymbol{u}\boldsymbol{v}_{\leq i} \models \text{prefix}(\phi)\} \right) \wedge \mathbb{I}\{\boldsymbol{u}\boldsymbol{v} \models \phi\} \tag{34}$$

$$= \left( \bigwedge_{i=1}^{|\boldsymbol{u}\boldsymbol{v}|} \mathbb{I}\{\boldsymbol{u}\boldsymbol{v}_{\leq i}\boldsymbol{w} \models \phi \text{ for some } \boldsymbol{w}\} \right) \wedge \mathbb{I}\{\boldsymbol{u}\boldsymbol{v} \models \phi\} \tag{35}$$

$$= \top. \tag{36}$$

If there is no such $\boldsymbol{v}$, then for all $\sigma$, we have $\boldsymbol{u}\sigma \not\models \text{prefix}(\phi)$, so $\boldsymbol{u} \not\models \text{next}_\sigma(\text{prefix}(\phi))$; moreover, $\boldsymbol{u} \not\models \phi$. Let $\mathbf{h} = \Phi'(\boldsymbol{u})_{|\boldsymbol{u}|}$ be the state after reading $\boldsymbol{u}$. All entries of $\mathbf{h}$ are false, which makes (for example) $a(\mathbf{h})(\text{EOS})$ true, so the $\boldsymbol{v} = \epsilon$ term of the summation in Eq. (32) is true.

(b) Let $A = (\Phi, a)$ be a $\text{TL}[\mathcal{O}]$ autoregressor, where $\Phi = (\phi_i)_{i=1}^m$.

Define

$$\phi_{\mathbf{h}} = \bigwedge_{i=1}^m (\phi_i \leftrightarrow h_i) \qquad \text{for } \mathbf{h} \in \mathbb{B}^m$$

$$\phi_\sigma = \bigwedge_{\mathbf{h} \in \mathbb{B}^m} (\phi_{\mathbf{h}} \leftrightarrow a(\mathbf{h})(\sigma)) \qquad \text{for } \sigma \text{ in } \Sigma \cup \{\text{EOS}\}.$$

Intuitively, $\phi_h$ is true whenever the state encoder is in state $\mathbf{h}$, and $\phi_\sigma$ is true whenever the state encoder is in a state in which $a$ predicts $\sigma$.

Then, define

$$\phi' = \mathbf{H} \left( \bigvee_{\sigma \in \Sigma \cup \{\text{EOS}\}} (\mathbf{Y}\phi_\sigma) \wedge \sigma \right)$$

and let $c(h) = h$, so that the classifier $C = (\phi', c)$ defines the same language as $A$. $\qquad\square$

### B.4 PROOF OF PROP. 6.9

**Proposition 6.9.** *(a) There does not exist a transformation* $\mathrm{prefix}'$ *such that* $\mathrm{prefix}'(\phi)$ *is constructible in polynomial time (in* $|\phi|$*) and satisfies Eq. (13) for every formula* $\phi$ *in* $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$*, unless* $\mathsf{P} = \mathsf{PSPACE}$*.*
*(b) Similarly for* $\mathsf{TL}[\mathbf{H}]$*, unless* $\mathsf{P} = \mathsf{NP}$*.*
*(c) Similarly for* $\mathsf{TL}[\mathbf{Y}]$*, unless* $\mathsf{P} = \mathsf{NP}$*.*

*Proof.* (a) Suppose that $\mathrm{prefix}'$ exists. For any formula $\phi$ of $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$, we can test whether $\phi$ is satisfiable by constructing $\mathrm{prefix}'(\phi)$ in polynomial time (by assumption) and then testing whether $\epsilon \models \mathrm{prefix}'(\phi)$, which can also be done in polynomial time, as shown by Fionda and Greco (2016, Thm. 8). (They assume formulas in negation normal form, but it is easy to generalize their result to formulas not in negation normal form.) But satisfiability in $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$ is PSPACE-complete (Giacomo and Vardi, 2013; Fionda and Greco, 2016), so this would imply $\mathsf{P} = \mathsf{PSPACE}$.

(b) Similarly, satisfiability in $\mathsf{TL}[\mathbf{H}]$ is NP-complete (Fionda and Greco, 2016), so the existence of a polynomial-time $\mathrm{prefix}'$ would imply $\mathsf{P} = \mathsf{NP}$.

(c) Same as the previous case.

$\square$

## C NONDETERMINISTIC FINITE AUTOMATA

We give a single definition of weighted NFAs instead of factoring them into unweighted NFAs and autoregressive output functions.

**Definition C.1** (Weighted Nondeterministic Finite Automaton)**.** *A **weighted nondeterministic finite automaton** is a tuple* $M = (\Sigma, Q, \delta, \alpha, \omega)$*, where*

- $\Sigma$ *is an alphabet*

- $Q$ *is a finite set of **states***

- $\delta \colon Q \times \Sigma \times Q \to \mathbb{K}$ *is a **transition function***

- $\iota \in Q$ *is the **initial** state*

- $\omega \colon Q \to \mathbb{K}$ *is an ending weight.*

*We extend* $\delta$ *to* $\delta^* \colon Q \times \Sigma^* \times Q \to \mathbb{K}$*:*

$$\delta^*(q, \epsilon, q) = \mathbf{1}$$
$$\delta^*(q, \epsilon, q') = \mathbf{0} \qquad q \neq q'$$
$$\delta^*(q_1, \sigma \boldsymbol{w}, q_2) = \bigoplus_{q \in Q} \delta(q_1, \sigma, q) \otimes \delta^*(q, \boldsymbol{w}, q_2).$$

*Then* $M$ *accepts* $\boldsymbol{w}$ *with weight* $k$ *iff*

$$k = \bigoplus_{q_2 \in Q} \delta^*(\iota, \boldsymbol{w}, q_2) \otimes \omega(q_2).$$

**Definition C.2.** *We say that an NFA with transition function* $\delta$ *is **counter-free** if there exists some* $k$ *such that for all states* $q_1, q_2$ *and all strings* $\boldsymbol{w}$*, we have* $\delta^*(q_1, \boldsymbol{w}^k, q_2) = \delta^*(q_1, \boldsymbol{w}^{k+1}, q_2)$*.*

A weighted automaton is determinizable if every every pair of states which are *siblings* (can be reached by the same string) are also *twins* (all cycles by the same string have the same weight) (Mohri, 1997). The automaton in Fig. 2c is counter-free because for any pair of states $q_1, q_2$ and any string $\boldsymbol{w}$, we have that $\delta^*(q_1, \boldsymbol{w}, q_2) = \delta^*(q_1, \boldsymbol{w}^2, q_2)$. However, it is not determinizable because $q_1$ and $q_2$ are siblings (both reachable by $a$) but not twins (the $a$-labeled cycles $q_1 \xrightarrow{a/\frac{1}{2}} q_1$ and $q_2 \xrightarrow{a/\frac{3}{4}} q_2$ on the two states have different weights).

# D  INEXPRESSIBILITY OF $(aab)^*$

**Proposition 6.11.** *The language* $(aab)^*$ *is not definable by any* $\mathsf{TL}[\mathbf{H}]$ *classifier or autoregressor.*

We actually prove a slightly stronger statement. Define the **Y-depth** of a formula $\phi$ to be the number of nested $\mathbf{Y}$ operators in $\phi$. Then we will prove that $(aab)^*$ is not definable by any formula of $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$ with $\mathbf{Y}$-depth 1. Since the conversion from an autoregressor to a classifier (Thm. 6.6(b)) adds a single $\mathbf{Y}$, we will conclude that $(aab)^*$ is not definable by any $\mathsf{TL}[\mathbf{H}]$ autoregressor.

**Lemma D.1.** *For any language $L$ over $\Sigma$, define* $\mathrm{Bigram}(L) = \{(\textsc{bos}, w_1) \cdot (w_1, w_2) \cdot (w_2, w_3) \cdots (w_{n-1}, w_n) \cdot (w_n, \textsc{eos}) \mid w \in L\}$. *If $\phi$ is a formula of* $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$ *with $\mathbf{Y}$-depth 1, then there is a formula* $\mathrm{noy}(\phi)$ *of* $\mathsf{TL}[\mathbf{H}]$ *over* $(\Sigma \cup \{\textsc{bos}\}) \times (\Sigma \cup \{\textsc{eos}\})$ *such that* $\mathcal{L}(\mathrm{noy}(\phi)) \cap \mathrm{Bigram}(\Sigma^*) = \mathrm{Bigram}(\mathcal{L}(\phi))$.

*Proof.* Define the transformation noy, which pushes $\mathbf{Y}$ down to the atomic formulas, then modifies the atomic formulas to operate on bigrams.

$$\mathrm{noy}(\neg\psi) = \neg\mathrm{noy}(\psi) \qquad\qquad \mathrm{noy}(\mathbf{Y}(\neg\psi)) = \neg\mathrm{noy}(\mathbf{Y}\psi)$$

$$\mathrm{noy}(\psi_1 \wedge \psi_2) = \mathrm{noy}(\psi_1) \wedge \mathrm{noy}(\psi_2) \qquad \mathrm{noy}(\mathbf{Y}(\psi_1 \wedge \psi_2)) = \mathrm{noy}(\mathbf{Y}\psi_1) \wedge \mathrm{noy}(\mathbf{Y}\psi_2)$$

$$\mathrm{noy}(\mathbf{H}\psi) = \mathbf{H}(\mathrm{noy}(\psi)) \qquad\qquad \mathrm{noy}(\mathbf{Y}(\mathbf{H}\psi)) = \mathbf{H}(\mathrm{noy}(\mathbf{Y}\psi))$$

$$\mathrm{noy}(\sigma) = \bigvee_{\sigma' \in \Sigma \cup \{\textsc{bos}\}} (\sigma', \sigma) \qquad\qquad \mathrm{noy}(\mathbf{Y}\sigma) = \bigvee_{\sigma' \in \Sigma \cup \{\textsc{eos}\}} (\sigma, \sigma')$$

$$\mathrm{noy}(\textsc{bos}) = \textsc{bos} \qquad\qquad \mathrm{noy}(\mathbf{Y}\,\textsc{bos}) = \bigvee_{\sigma' \in \Sigma \cup \{\textsc{eos}\}} (\textsc{bos}, \sigma'). \qquad \square$$

Then, to prove that $(aab)^*$ is not definable in $\mathsf{TL}[\mathbf{H}, \mathbf{Y}]$ with $\mathbf{Y}$-depth 1, suppose it is definable by $\phi$. By Lem. D.1, there is a formula $\mathrm{noy}(\phi)$ of $\mathsf{TL}[\mathbf{H}]$ such that

$$(\textsc{bos}, a) \cdot (a, a) \cdot (a, b) \cdot (b, \textsc{eos}) \in \mathcal{L}(\mathrm{noy}(\phi)) \cap \mathrm{Bigram}(\Sigma^*).$$

But $\mathcal{L}(\mathrm{noy}(\phi))$ must be stutter-invariant (Peled and Wilke, 1997), so we also have

$$(\textsc{bos}, a) \cdot (a, a) \cdot (a, a) \cdot (a, b) \cdot (b, \textsc{eos}) \in \mathcal{L}(\mathrm{noy}(\phi)) \cap \mathrm{Bigram}(\Sigma^*).$$

But this is not in $\mathrm{Bigram}((aab)^*)$, which is a contradiction.