



COMILLAS
UNIVERSIDAD PONTIFICIA

ICAI

ICADE

CIHS

Version

This is an Accepted Manuscript version

Citation for published version

J. P. Silvestre, Á. R. Abella and P. Tabuada, "Stopping LLMs from Going Rogue: A Control Barrier Approach to Text Generation," 2025 IEEE 64th Conference on Decision and Control (CDC), Rio de Janeiro, Brazil, 2025, pp. 151-156, doi: 10.1109/CDC57313.2025.11312450.

General rights

This manuscript version is made available under the CC-BY-NC-ND 4.0 licence

<https://web.upcomillas.es/webcorporativo/RegulacionRepositorioInstitucionalComillas.pdf>.

Take down policy

If you believe that this document breaches copyright please contact Universidad Pontificia Comillas providing details, and we will remove access to the work immediately and investigate your claim.

Stopping LLMs from Going Rogue: A Control Barrier Approach to Text Generation

João Pedro Silvestre¹, Álvaro Rodríguez Abella¹ and Paulo Tabuada¹

Abstract—The rapid integration of large language models (LLMs) into our everyday lives has outpaced safety considerations aimed at protecting users from toxic outputs and preventing malicious actors from generating harmful text at scale. As a result, LLMs have been exploited by bots capable of producing vast amounts of harmful and toxic content, enabling users to manipulate online opinions and, in some cases, create dangerous online environments.

Our work addresses this issue by developing a framework for designing safety filters that preclude toxic outputs. To achieve this, we leverage Control Barrier Functions (CBFs) which enable the design of closed-loop systems that remain safe. We consider the continuous-time model of an LLM, where tokens are regarded as the state of the model, and prove that by only controlling the first token, any function satisfying mild assumptions becomes a CBF. Our approach can be utilized to design LLMs capable of ensuring safety of its outputs without significantly affecting the original model’s behavior.

I. INTRODUCTION

The introduction of transformers [1] marked a turning point in natural language processing, with early works such as the GPT [2] already representing a breakthrough over older techniques such as the Long Short Term Memory models [3]. Increasingly sophisticated models, such as GPT-4 [4] and DeepSeek [5] have recently emerged, capable of generating text nearly indistinguishable from human output [6]. Unfortunately, due to the rapid speed at which progress has been made, safety concerns have been overshadowed. One work which details some overlooked problems is [7], where it is shown that large language models (LLMs) may generate toxic outputs, from apparently innocuous prompts.

Our understanding of these systems has lagged behind the massive increases in their size and complexity, leaving no other option but to treat them as black boxes [8]. This motivated two distinct research approaches. On the empirical side, many studies analyze tokens at each layer to decipher model behavior [9]. Other recent studies have explored the capacity that LLMs have to generate self-explanations, mainly showing that these explanations may still not be trusted [10]. These type of approaches have been leveraged to design adaptations to the model to address the toxicity problem [11]. However, due to their empirical nature, it remains a challenge to properly control LLMs and stop

The work of João Pedro Silvestre was partially supported by the PhD fellowship 2023.01843.BD from the Fundação para a Ciência e a Tecnologia (FCT), Portugal.

¹João Pedro Silvestre, Álvaro Rodríguez Abella, and Paulo Tabuada are with the Electrical and Computer Engineering Department, University of California at Los Angeles, Los Angeles, CA 90095 USA (e-mail: {joaosilvestre, rodriguezabella, tabuada}@ucla.edu).

toxic text generation while providing any type of formal guarantees. In fact, some works even suggest safeguards may be easily bypassed [12], [13].

On the theoretical side, recent works have leveraged the mathematical model proposed in [?] to analyze LLM behavior and derive asymptotic conclusions [14], [15]. Notably, to derive such results, these works typically make simplifying assumptions that limit their direct use, such as the absence of feedforward layers. Nevertheless, such methods allow the interpretation of LLMs as dynamical systems and, once practical limitations are overcome, control techniques will enable the improvement of model safety.

While some works have used control tools such as optimal control to tackle the toxicity problem [16], one particular concept that immediately comes to mind is safety. By separating the state space into safe and unsafe sets, and by constraining tokens to the safe set, it may be possible to preclude toxic outputs. Seminal works such as [17] have laid the groundwork for modern safety techniques, typically based on control barrier functions (CBFs).

Some works have used a CBF approach to the toxicity problem [18]. However, unlike our approach, they act after the model has computed the probability distribution of the next token, filtering that distribution if necessary.

In this paper, tokens leading to toxic text are considered unsafe states and sufficient conditions for designing CBFs that constrain tokens within the safe states are established. To achieve this, we first extend the continuous-time model derived in [?] to better reflect the actual transformer model. Furthermore, we assume that the first token can be controlled throughout the model for any given prompt. We then prove that any generic function can serve as a CBF, provided that there always exists a direction in which its value changes. Therefore, our slight modification of the original model is sufficient to ensure that the output never becomes toxic.

Our contributions are threefold:

- 1) The continuous-time model of an LLM [?] is extended to model feedforward layers.
- 2) Sufficient conditions for designing CBFs preventing toxic outputs are presented.
- 3) It is proved that controlling only the first token is sufficient to keep the generated tokens within the safe set.

Notations

Let $r, s, k \in \mathbb{N} = \{1, 2, \dots\}$. Points in the Euclidean space \mathbb{R}^r are regarded as column vectors and denoted by $x = (x^1, \dots, x^r) \in \mathbb{R}^r$. Tuples of k points are denoted by

$\mathbf{x} = (x_1, \dots, x_k) \in (\mathbb{R}^r)^k$ (note the different font). Similarly, the space of $r \times s$ real matrices is denoted by $\mathbb{R}^{r \times s}$. In particular, $\mathbb{I}_r \in \mathbb{R}^{r \times r}$ denotes the identity matrix. Tuples of matrices are denoted by $\mathbf{A} = (A_1, \dots, A_k) \in (\mathbb{R}^{r \times s})^k$. Lastly, the tangent space of a smooth manifold M at $p \in M$ and its elements are denoted by $T_p M$ and $X_p \in T_p M$, respectively. Given another smooth manifold N and a smooth map $\phi : M \rightarrow N$, i.e., $\phi \in C^\infty(M, N)$, its tangent map is denoted by $T\phi : TM \rightarrow TN$.

Preliminaries on safety and control barrier functions

In this section we give a brief overview of control barrier functions (CBFs) to better contextualize our work. First, consider a nonlinear control system of the form:

$$\dot{\mathbf{x}} = f(t, \mathbf{x}) + g(t, \mathbf{x}) \cdot \mathbf{u}, \quad (t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1}, \quad (1)$$

where $\mathbf{u} \in \mathbb{R}^m$ is the input of the system, $m \in \mathbb{N}$, $f : \mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$ and $g : \mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{(n+1) \times m}$ are piece-wise continuous in t and locally Lipschitz continuous in \mathbf{x} , and $\mathbb{R}_{\geq 0}$ denotes the non-negative reals. Henceforth, we assume that (1) is forward-complete, i.e., for each $\mathbf{x}_0 \in \mathbb{R}^{n+1}$ and piece-wise continuous input $u : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^m$, the solution of the controlled system, $\dot{\mathbf{x}} = f(t, \mathbf{x}) + g(t, \mathbf{x}) \cdot \mathbf{u}(t)$, $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1}$, with initial condition \mathbf{x}_0 at $t_0 = 0$ is uniquely defined for all $t \in \mathbb{R}_{\geq 0}$.

Definition 1.1: A set $\mathcal{C} \subset \mathbb{R}^{n+1}$ is *forward invariant* if, for each initial condition $\mathbf{x}_0 \in \mathcal{C}$, the solution of (1) with initial condition \mathbf{x}_0 at $t_0 = 0$ remains in \mathcal{C} for each $t \in \mathbb{R}_{\geq 0}$. Furthermore, a set \mathcal{C} is said to be *safe* for the system (1) if it is forward invariant.

Control barrier functions [17] provide a structured method to design controllers that guarantee safety for sets of the form $\mathcal{S} = h^{-1}(\mathbb{R}_{\geq 0})$, where $h \in C^1(\mathbb{R}^{n+1})$. Henceforth, the Lie derivative of h with respect to f is denoted by $L_f h = dh(f)$. Analogously, we denote $L_g h = (L_{g_1} h, \dots, L_{g_m} h)$, where $g = (g_1, \dots, g_m) \in \mathbb{R}^{(n+1) \times m}$.

Definition 1.2 (Control Barrier Function [17]): A function $h \in C^1(\mathbb{R}^{n+1})$ is a *control barrier function (CBF)* for (1) on $\mathcal{S} = h^{-1}(\mathbb{R}_{\geq 0})$ if there exists¹ $\alpha \in \mathcal{K}_\infty^e$ such that, for each $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathcal{S}$, we have:

$$\sup_{\mathbf{u} \in \mathbb{R}^m} (L_f h + L_g h \cdot \mathbf{u}) \geq -\alpha \circ h.$$

Note that, although h is time-independent, the Lie derivatives $L_f h$ and $L_g h$ depend on t as the vector fields f and g are time-dependent. The existence of a CBF implies the non-emptiness of the set of safe inputs (the argument $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1}$ is omitted for brevity):

$$K_{\text{CBF}} = \{\mathbf{u} \in \mathbb{R}^m \mid L_f h + L_g h \cdot \mathbf{u} \geq -\alpha \circ h\},$$

Lastly, [17, Theorem 2] ensures that any safe controller $k : \mathbb{R}_{\geq 0} \times \mathcal{S} \rightarrow \mathbb{R}^m$, i.e., one satisfying $k(t, \mathbf{x}) \in K_{\text{CBF}}(t, \mathbf{x})$ for each $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathcal{S}$ guarantees safety. In practice, if there are no additional constraints on the input, a sufficient

¹An extended class \mathcal{K} function, $\alpha \in \mathcal{K}_\infty^e$, is a strictly increasing function $\alpha : \mathbb{R} \rightarrow \mathbb{R}$ such that $\alpha(0) = 0$, $\lim_{r \rightarrow \infty} \alpha(r) = \infty$ and $\lim_{r \rightarrow -\infty} \alpha(r) = -\infty$.

condition for h to be a CBF is $(L_g h)(t, \mathbf{x}) \neq 0$ for each $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathcal{S}$.

Higher-order barrier functions [19] constitute an extension of CBFs that are useful when the Lie derivative $L_g h$ vanishes and, thus, h is not a CBF in the standard sense. In this work, we will only need second-order CBFs. Given $h \in C^2(\mathbb{R}^{n+1})$ and $\alpha \in \mathcal{K}_\infty^e \cap C^1(\mathbb{R})$, we define the extension $h_1 \in C^1(\mathbb{R}_{\geq 0} \times \mathbb{R}^{n+1})$ of h by:

$$h_1(t, \mathbf{x}) = (dh)_x(f(t, \mathbf{x})) + \alpha(h(\mathbf{x})). \quad (2)$$

Similarly, let $h_1^{-1}(\mathbb{R}_{\geq 0})^\circ = \{\mathbf{x} \in \mathbb{R}^{n+1} \mid h_1(0, \mathbf{x}) \geq 0\}$.

Definition 1.3 (Second-order CBF, adapted from [20]):

A *second-order control barrier function* for (1) is a function $h \in C^2(\mathbb{R}^2)$ for which there exist $\alpha \in \mathcal{K}_\infty^e \cap C^1(\mathbb{R})$, $\alpha_1 \in \mathcal{K}_\infty^e$ and $\mathcal{D} \subset \mathbb{R}^{n+1}$ open such that:

$$\mathcal{S} = h^{-1}(\mathbb{R}_{\geq 0}) \cap h_1^{-1}(\mathbb{R}_{\geq 0})^\circ \subset \mathcal{D},$$

with h_1 as in (2), and we have that $L_g h = 0$ and:

$$\sup_{\mathbf{u} \in \mathbb{R}^m} (L_f h_1 + L_g h_1 \cdot \mathbf{u}) \geq -\alpha_1 \circ h_1. \quad (3)$$

In [20, Theorem 1] it is shown that the set \mathcal{S} is safe provided the control input belongs to the following set (again, we omit the argument $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathcal{D}$ for brevity):

$$K_{\text{CBFH}} = \{\mathbf{u} \in \mathbb{R}^m \mid L_f h_1 + L_g h_1 \cdot \mathbf{u} \geq -\alpha \circ h_1\}.$$

Similarly to standard CBFs, the condition $(L_g h_1)(t, \mathbf{x}) \neq 0$ for each $(t, \mathbf{x}) \in \mathbb{R}_{\geq 0} \times \mathcal{D}$ is sufficient to prove that h is a second-order CBF for (1) (provided $(L_g h)(t, \mathbf{x}) = 0$).

II. CONTROLLED DYNAMICS OF TRANSFORMERS

In the following, let $n, d, k, h \in \mathbb{N}$ be fixed parameters. In this section, we start by defining the configuration space of the transformer, which is essentially the ellipsoid onto which the normalization at the end of each layer projects to. Afterward, we detail the attention mechanism present in every modern transformer and conclude with the discrete-time model for the transformer as well as the corresponding continuous-time first-order approximation. Readers not interested in the intricacies of transformers can directly jump to the continuous-time model (7).

A. Configuration space

Let $n, d, k, h \in \mathbb{N}$. A symmetric, positive-definite matrix $W \in \mathbb{R}^{(n+1) \times (n+1)}$ defines an inner product on \mathbb{R}^{n+1} :

$$\langle X_x, Y_x \rangle_W = X_x^\top W Y_x,$$

for each $X_x, Y_x \in T_x \mathbb{R}^{n+1}$ and $\mathbf{x} \in \mathbb{R}^{n+1}$. The corresponding norm is denoted by $|X_x|_W = (X_x^\top W X_x)^{1/2}$. The points of \mathbb{R}^{n+1} of unit norm with respect to W define an n -dimensional ellipsoid, which is denoted by:

$$\mathcal{E}_W^n = \{y \in \mathbb{R}^{n+1} \mid y^\top W y = 1\}.$$

As we consider a transformer consisting of k tokens, the resulting state space is the Cartesian product of k copies of the ellipsoid, i.e., $(\mathcal{E}_W^n)^k = \underbrace{\mathcal{E}_W^n \times \dots \times \mathcal{E}_W^n}_{k\text{-times}}$.

Similarly, consider the ellipsoid projection:

$$\pi_W : \mathbb{R}_0^{n+1} \rightarrow \mathcal{E}_W^n, \quad x \mapsto \pi_W(x) = x|_W^{-1},$$

where $\mathbb{R}_0^{n+1} = \mathbb{R}^{n+1} - \{0\}$. Its tangent map at each $x \in \mathbb{R}_0^{n+1}$, $T_x \pi_W : T_x \mathbb{R}_0^{n+1} \rightarrow T_{\pi_W(x)} \mathcal{E}_W^n$, is given by:

$$T_x \pi_W \cdot X_x = |x|_W^{-1} (\mathbb{I}_{n+1} - x x^\top W |x|_W^{-2}) \cdot X_x,$$

for each $X_x \in T_x \mathbb{R}_0^{n+1}$. In particular, for $y \in \mathcal{E}_W^n$, we have $T_y \pi_W \cdot X_y = (\mathbb{I}_{n+1} - y y^\top W) \cdot X_y$.

The standard Euclidean metric on \mathbb{R}_0^{n+1} induces a Riemannian metric g on \mathcal{E}_W^n , whose Levi-Civita connection is denoted by ∇^g . In turn, the induced metric on the product manifold $(\mathcal{E}_W^n)^k$ and its Levi-Civita connection are, respectively, denoted by $\mathbf{g} = \underbrace{g \oplus \dots \oplus g}_{k\text{-times}}$ and $\nabla^{\mathbf{g}}$.

B. Self-attention mechanism

Let $\mathbf{y} \in (\mathcal{E}_W^n)^k$ and $\mathbf{P} \in (\mathbb{R}^{(n+1) \times (n+1)})^h$. For each head $1 \leq \eta \leq h$, the η -th (full) self attention matrix is defined as:

$$A_\eta = (\alpha_{ij}^\eta)_{\substack{2 \leq i \leq k, \\ 1 \leq j \leq k}}, \in \mathbb{R}^{(k-1) \times k},$$

where, for each $1 \leq i, j \leq k, i \neq 1$:

$$\alpha_{ij}^\eta(\mathbf{P}, \mathbf{y}) = \frac{1}{Z_i^\eta(\mathbf{P}, \mathbf{y})} \exp(y_i^\top P_\eta y_j),$$

$$Z_i^\eta(\mathbf{P}, \mathbf{y}) = \sqrt{n+1} \sum_{j=1}^k \exp(y_i^\top P_\eta y_j).$$

The first token is excluded since its dynamics will be controlled, as explained in the next section. Alternatively, we may consider the so-called *auto-regressive* (also known as *causal* or *masked*) self-attention. In that case, the sum in the definition of Z_i^η runs from $j = 1$ to i , and $\alpha_{ij}^\eta = 0$ whenever $j > i$. The results in this paper are valid for both the full and the auto-regressive architectures, so we do not explicitly differentiate between them.

C. Discrete-time transformer model

While there are several transformer models, most of them inherit the basic structure of the original model [1]. In essence, as detailed in [21], each transformer layer is composed of two parts: the feed-forward layer and the self-attention mechanism, each with its own skip connection and normalization. The model presented in [?] made the simplifying assumption that feedforward layers were absent. In this section we present a discrete-time model including both layers.

Formally, for each $l \in \mathbb{N}$ and $1 \leq \eta \leq h$, let:

$$P_\eta(l), U_\eta(l) \in \mathbb{R}^{(n+1) \times (n+1)},$$

$$W_1(l), W_2^\top(l) \in \mathbb{R}^{d \times (n+1)},$$

$$b_1(l) \in \mathbb{R}^d, b_2(l) \in \mathbb{R}^{n+1}.$$

The discrete self-attention and feed-forward maps:

$$\tilde{f}_{\text{sa},i} : \mathbb{N} \times (\mathcal{E}_W^n)^k \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{n+1},$$

$$\tilde{f}_{\text{ff}} : \mathbb{N} \times \mathcal{E}_W^n \times \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{n+1},$$

$1 < i \leq k$, are defined as:

$$\tilde{f}_{\text{sa},i}(l, \mathbf{y}, \tau) = y_i + \tau \sum_{\eta=1}^h \sum_{j=1}^k \alpha_{ij}^\eta(\mathbf{P}(l), \mathbf{y}) U_\eta(l) y_j,$$

$$\tilde{f}_{\text{ff}}(l, y, \tau) = W_2(l) \Sigma(W_1(l) y + b_1(l)) + b_2(l),$$

respectively. In the previous expressions, $\tau \in \mathbb{R}_{\geq 0}$ is a small training weight usually absorbed into the matrices $U_\eta(l)$, $1 \leq \eta \leq h$. In turn, we denote by $f_{\text{sa},i} = \pi_W \circ \tilde{f}_{\text{sa},i}$ and $f_{\text{ff}} = \pi_W \circ \tilde{f}_{\text{ff}}$ the corresponding normalized maps. In the previous expression, we denote $\Sigma(z) = (\sigma(z^1), \dots, \sigma(z^d))$ for each $z = (z^1, \dots, z^d) \in \mathbb{R}^d$, where $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is some *activation function*.

To allow for designing CBFs, we alter the transformer architecture by assuming we can arbitrarily steer the first token. This approach has three main advantages: (1) we minimize the number of controlled tokens, (2) we can better approximate the original model by selecting a control input which is as close to the original model dynamics as possible, and (3) this control input is sufficient to keep the system safe, as will be shown in Section III.

Thus, for each $\mathbf{y} \in (\mathcal{E}_W^n)^k$, $u \in \mathbb{R}^{n+1}$, $k \in \mathbb{N}$ and $\tau > 0$, the *discrete-time model for the controlled transformer* reads:

$$y_i(l+1) = \pi_W(y_i(l) + \tau u), \quad (4)$$

if $i = 1$, and

$$y_i(l+1) = \pi_W(f_{\text{sa},i}(l, \mathbf{y}, \tau) + \tau \tilde{f}_{\text{ff}}(l, f_{\text{sa},i}(l, \mathbf{y}, \tau))), \quad (5)$$

for $1 < i \leq k$.

Remark 2.1: This choice of architecture is conceptually similar to prefixing [22], a technique where an additional prefix is added to the LLM's prompt to steer it away from toxicity. Our approach is a form of pre-fixing where we only use one token but, rather than keeping it constant, we steer it to keep the system safe. Hence, we propose to alter the architecture by adding a controller that steers the single token pre-fixing the prompt. The design of such controller is outside the scope of this paper.

D. Continuous-time attention model

Let (M, g) be a Riemannian manifold (recall that, for the ellipsoid, $M = \mathcal{E}_W^n$, we utilize the Riemannian metric induced by the Euclidean metric on \mathbb{R}_0^{n+1}). The flow of a (forward-complete) vector field $Y \in \mathfrak{X}(M)$ is denoted by $Y^\tau : M \rightarrow M$, $\tau \geq 0$. A map $\phi : M \times \mathbb{R}_{\geq 0} \rightarrow M$ is a first order approximation to Y^τ if there exist $T > 0$ and $\xi : M \rightarrow \mathbb{R}_{\geq 0}$ such that $d_g(Y^\tau(m), \phi(m, \tau)) \leq \xi(m) \tau^2$ for each $\tau \in [0, T]$ and $m \in M$, where d_g is the Riemannian distance on M .

Let us compute the best first order approximation of the discrete model introduced in the previous section. Let $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$. For $i = 1$, we have:

$$\left. \frac{d}{d\tau} \right|_{\tau=0} \pi_W(y_1(t) + \tau u) = T_{y_1} \pi_W \cdot u.$$

Similarly, for $1 < i \leq k$, it is readily obtained:

$$\begin{aligned} \frac{d}{d\tau} \Big|_{\tau=0} & \pi_W (f_{\text{sa},i}(t, \mathbf{y}, \tau) + \tau \tilde{f}_{\text{ff}}(t, f_{\text{sa},i}(t, \mathbf{y}, \tau))) \\ & = T_{y_i} \pi_W \cdot \left(\tilde{f}_{\text{sa},i}(t, \mathbf{y}) + \tilde{f}_{\text{ff}}(t, y_i) \right). \end{aligned}$$

In the previous expression, the continuous *self-attention* map, $\tilde{f}_{\text{sa},i} : \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k \rightarrow \mathbb{R}^{n+1}$, and *feed-forward* map, $\tilde{f}_{\text{ff}} : \mathbb{R}_{\geq 0} \times \mathcal{E}_W^n \rightarrow \mathbb{R}^{n+1}$, are defined as:

$$\begin{aligned} \tilde{f}_{\text{sa},i}(t, \mathbf{y}) &= \sum_{\eta=1}^h \sum_{j=1}^k \alpha_{ij}^\eta(\mathbf{P}(t), \mathbf{y}) U_\eta(t) y_j, \\ \tilde{f}_{\text{ff}}(t, y) &= W_2(t) \Sigma(W_1(t) y + b_1(t)) + b_2(t), \end{aligned}$$

respectively. In turn, the corresponding normalized maps are respectively denoted by $f_{\text{sa},i}(t, \mathbf{y}) = T_{y_i} \pi_W \cdot \tilde{f}_{\text{sa},i}(t, \mathbf{y})$ and $f_{\text{ff}}(t, y) = T_y \pi_W \cdot \tilde{f}_{\text{ff}}(t, y)$.

By gathering the previous expressions, we conclude that, for each $(t, \mathbf{y}, \mathbf{u}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k \times \mathbb{R}^{n+1}$, the best first order approximation in τ of the discrete model (4), (5) is given by:

$$\dot{y}_i = \begin{cases} T_{y_1} \pi_W \cdot \mathbf{u}, & i = 1, \\ f_{\text{sa},i}(t, \mathbf{y}) + f_{\text{ff}}(t, y_i), & 1 < i \leq k. \end{cases} \quad (6)$$

By denoting:

$$\begin{aligned} f_{\text{sa}}(t, \mathbf{y}) &= (0, f_{\text{sa},2}(t, \mathbf{y}), \dots, f_{\text{sa},k}(t, \mathbf{y})), \\ f_{\text{ff}}(t, \mathbf{y}) &= (0, f_{\text{ff}}(t, y_2), \dots, f_{\text{ff}}(t, y_k)), \\ g(t, y) &= (T_y \pi_W, 0, \dots, 0), \end{aligned}$$

the previous model may be equivalently expressed as:

$$\dot{\mathbf{y}} = f_{\text{sa}}(t, \mathbf{y}) + f_{\text{ff}}(t, \mathbf{y}) + g(t, y_1) \cdot \mathbf{u}. \quad (7)$$

Let $f : \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k \rightarrow T(\mathcal{E}_W^n)^k$ be the time-dependent vector field defined as $f(t, \mathbf{y}) = f_{\text{sa}}(t, \mathbf{y}) + f_{\text{ff}}(t, \mathbf{y})$.

When comparing this model to the one in [?], the main difference is the additional sum of the feedforward layer and the control input. Furthermore, thanks to our choices, the system (7) is control affine, which simplifies the computation of safe control inputs and ensures that the set of safe inputs K_{CBF} is convex.

Remark 2.2: Note that (7) is forward-complete, since the vector fields are smooth and the system evolves on a compact manifold without boundary.

III. SAFETY OF CONTROLLED LLMs

In this section, we approach the toxicity problem as a safety problem, specifically, our goal is to constrain the last token, y_k , to a region where the corresponding sentence is not toxic. We focus on the last token because LLMs generate text sequentially: the prompt is converted into a set of tokens, the initial state of the transformer, and propagated through the transformer for a finite amount of time. After this propagation, the last token is used to predict the next token which is appended at the end of the prompt to create a new initial state. The process is then repeated starting from this new initial state. Therefore, to prevent toxic outputs, we want

to ensure the last token does not lead to a toxic sentence. We thus employ a classifier that transforms the initial state $\mathbf{y}^0(0)$ and the last token at the current state $y_k(t)$ into the toxicity score of the sentence formed by appending the predicted next token, based on $y_k(t)$, to the sequence of tokens defined by the initial state $\mathbf{y}^0(0)$.

A *toxicity classifier* is a map $c_k : (\mathcal{E}_W^n)^{k+1} \rightarrow \mathbb{R}_{\geq 0}$ that assigns a toxicity score to each sequence of tokens (y_1^0, \dots, y_k^0, y) , where $\mathbf{y}^0 = (y_1^0, \dots, y_k^0)$ and y correspond to the initial sequence fed to the transformer and the generated token at the k -th pass, respectively. Given that the initial sequence is constant, we may consider the map $h_0 : \mathcal{E}_W^n \rightarrow \mathbb{R}_{\geq 0}$ defined as $h_0(y) = c_k(y_1^0, \dots, y_k^0, y)$. Henceforth, we will assume that $h_0 \in C^2(\mathcal{E}_W^n, \mathbb{R}_{\geq 0})$. Let us consider the map $h \in C^2((\mathcal{E}_W^n)^k)$ defined as:

$$h(\mathbf{y}) = -h_0(y_k) + \gamma, \quad (8)$$

where $\gamma > 0$ is a fixed toxicity threshold. The rest of the paper is devoted to show that h is a second-order CBF for (7). Note that $(dh)_{\mathbf{y}} = -(dh_0)_{y_k}$. Given $\alpha \in \mathcal{K}_\infty^e \cap C^1(\mathbb{R})$, we consider the map $h_1 \in C^1(\mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k)$ defined as:

$$h_1(t, \mathbf{y}) = -(dh_0)_{y_k}(f_k(t, \mathbf{y})) + \alpha(-h_0(y_k) + \gamma). \quad (9)$$

This map corresponds to the auxiliary function (2), since $(dh)_{\mathbf{y}}(f(t, \mathbf{y})) = -(dh_0)_{y_k}(f_k(t, \mathbf{y}))$. Our approach is based on the observation that $(L_g h)(t, \mathbf{y}) = 0$ for $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$ for any h induced by a toxicity classifier c_k and thus we need a higher-order CBF.

The analysis in this paper requires the classifier to be C^2 . If needed, we can always approximate a less regular classifier with one of class C^2 .

This approach naturally increases the computational footprint of word generation due to the safety filter. However, since the constraints are linear, the increase is relatively modest.

A. Generic surjectivity of the vector field derivative

In this section, we investigate the partial derivative of the last component of the self-attention vector field, $f_{\text{sa},k}$, with respect to the first token, y_1 . We conclude that, for each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$, it has maximum rank (and, thus, defines an isomorphism) for almost every choice of the matrices $(\mathbf{P}(t), \mathbf{U}(t))$. In turn this will ensure that we can provide sufficient conditions on h alone without being dependent on the state.

For each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$ and $1 \leq i \leq k$, let us introduce the following linear map:

$$\psi_i(t, \mathbf{y}) = T_{y_k} \pi_W \cdot \frac{\partial \tilde{f}_{\text{sa},k}}{\partial y_i}(t, \mathbf{y}) : T_{y_i} \mathcal{E}_W^n \rightarrow T_{y_k} \mathcal{E}_W^n,$$

where we denote $\tilde{f}_{\text{sa},k} = (\tilde{f}_{\text{sa},k}^1, \dots, \tilde{f}_{\text{sa},k}^{n+1})$ and

$$\frac{\partial \tilde{f}_{\text{sa},k}}{\partial y_i} = \left(\frac{\partial \tilde{f}_{\text{sa},k}^\mu}{\partial y_i^\nu} \right)_{1 \leq \mu, \nu \leq n+1}.$$

In the same vein, let us introduce the map:

$$\Xi(\mathbf{y}) : \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^{2h} \rightarrow \mathbb{R}^{(n+1) \times (n+1)},$$

defined as:

$$\Xi(\mathbf{y})(\mathbf{P}, \mathbf{U}) = \sum_{\eta=1}^h \alpha_{k1}^{\eta}(\mathbf{P}, \mathbf{y}) U_{\eta} A_{\eta}(\mathbf{P}, \mathbf{y}),$$

where we denote:

$$A_{\eta}(\mathbf{P}, \mathbf{y}) = \mathbb{I}_{n+1} + \left(y_1 - \sum_{j=1}^k \alpha_{kj}^{\eta}(\mathbf{P}, \mathbf{y}) y_j \right) y_k^{\top} P_{\eta}.$$

Lemma 3.1: For each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$, we have:

$$\frac{\partial \tilde{f}_{sa,k}}{\partial y_1}(t, \mathbf{y}) = \Xi(\mathbf{y})(\mathbf{P}(t), \mathbf{U}(t)).$$

Proof: For brevity, we omit the arguments of the maps. Firstly, a computation shows that:

$$\frac{\partial \alpha_{kj}^{\eta}}{\partial y_1} = \alpha_{kj}^{\eta} (\delta_{j1} - \alpha_{k1}^{\eta}) y_k^{\top} P_{\eta},$$

where δ_{j1} denotes the Kronecker delta. From this, we conclude:

$$\begin{aligned} \frac{\partial \tilde{f}_{sa,k}}{\partial y_1} &= \sum_{\eta=1}^h U_{\eta} \left(\alpha_{k1}^{\eta} \mathbb{I}_{n+1} + \sum_{j=1}^k y_j \frac{\partial \alpha_{kj}^{\eta}}{\partial y_1} \right) \\ &= \sum_{\eta=1}^h \alpha_{k1}^{\eta} U_{\eta} \left(\mathbb{I}_{n+1} + \left(y_1 - \sum_{j=1}^k \alpha_{kj}^{\eta} y_j \right) y_k^{\top} P_{\eta} \right). \end{aligned}$$

The previous lemma shows the parametric dependence of $\partial \tilde{f}_{sa,k} / \partial y_1(t, \mathbf{y})$ (and, thus, of $\psi_1(t, \mathbf{y})$) on the matrices $(\mathbf{P}(t), \mathbf{U}(t))$. The next results ensure that $\psi_1(t, \mathbf{y})$ is surjective (and, thus, a linear isomorphism) for almost every choice of these matrices.

Proposition 3.1: Let $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k$. Then the map $\psi_1(t, \mathbf{y}) : T_{y_1} \mathcal{E}_W^n \rightarrow T_{y_k} \mathcal{E}_W^n$ is surjective for almost every $(\mathbf{P}(t), \mathbf{U}(t)) \in \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^{2h}$.

Proof: The map $\tilde{\Xi}(\mathbf{y}) : \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^{2h} \rightarrow \mathbb{R}$ defined as $\tilde{\Xi}(\mathbf{y})(\mathbf{P}, \mathbf{U}) = \det \Xi(\mathbf{y})(\mathbf{P}, \mathbf{U})$ is real analytic. Indeed, the exponential and the determinant are real analytic, and the functions $Z_k^{\eta}(\cdot, \mathbf{y}) : \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^h \rightarrow \mathbb{R}$, $1 \leq \eta \leq h$, never vanish, whence $\tilde{\Xi}(\mathbf{y})$ is real analytic. Moreover, $\tilde{\Xi}(\mathbf{y})$ is not identically zero: for instance, let $\mathbf{P}_0 = (\mathbf{0}_{n+1}, \dots, \mathbf{0}_{n+1})$, where $\mathbf{0}_{n+1}$ is the zero matrix of size $n+1$, and $\mathbf{U}_0 = (\mathbb{I}_{n+1}, \dots, \mathbb{I}_{n+1})$. Then we get $\alpha_{k1}^{\eta}(\mathbf{P}_0, \mathbf{y}_0) = (k\sqrt{n+1})^{-1}$, $1 \leq \eta \leq h$, and, thus, $\tilde{\Xi}(\mathbf{y})(\mathbf{P}_0, \mathbf{U}_0) = (h/(k\sqrt{n+1}))^{n+1} \neq 0$. As a result, $\tilde{\Xi}(\mathbf{y})^{-1}(0)$ has zero measure (cf., for example, [23] for an elementary proof). From Lemma 3.1, we conclude that $\text{rank}(\partial \tilde{f}_{sa,k} / \partial y_1)(t, \mathbf{y}) = n+1$ for almost every $(\mathbf{P}(t), \mathbf{U}(t)) \in \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^{2h}$. Subsequently, we have that $\text{rank} \psi_1(t, \mathbf{y}) = n = \dim T_{y_k} \mathcal{E}_W^n$, where we used that $\text{rank} T_{y_k} \pi_W = n$. Thus, $\psi_1(t, \mathbf{y})$ is surjective for almost every $(\mathbf{P}(t), \mathbf{U}(t)) \in \left(\mathbb{R}^{(n+1) \times (n+1)}\right)^{2h}$. ■

B. Dynamics of the control barrier function

In this section we present the main result of the paper, Corollary 3.1, showing that the function $h \in C^2((\mathcal{E}_W^n)^k)$ defined in (8) is a second-order CBF for the dynamics of the transformer (7). We can thus leverage this CBF to prevent toxic outputs. Corollary 3.1 is a direct consequence of Theorem 3.1.

We start with the following technical lemma.

Lemma 3.2: For each $(t, \mathbf{y}, \mathbf{u}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k \times \mathbb{R}^{n+1}$, we have $\nabla_{g(t, y_1) \cdot \mathbf{u}}^{\mathbf{g}} f_{sa,k}(t, \mathbf{y}) = \psi_1(t, \mathbf{y}) \cdot T_{y_1} \pi_W \cdot \mathbf{u}$.

Proof: Let $\Gamma_{\mu\nu, l}^{ij, \sigma} \in C^{\infty}((\mathcal{E}_W^n)^k)$, $1 \leq i, j, l \leq k$, $1 \leq \mu, \nu, \sigma \leq n+1$, be the Christoffel symbols of \mathbf{g} , and denote $\partial_{\mu} = (0, \dots, \underbrace{1}_{\mu\text{-th position}}, \dots, 0)$, $1 \leq \mu \leq n+1$, which

allows for writing $T_{y_1} \pi_W \cdot \mathbf{u} = \sum_{\mu=1}^{n+1} \xi_1^{\mu}(y_1, \mathbf{u}) \partial_{\mu} \in T_{y_1} \mathcal{E}_W^n$, for some $\xi_1^{\mu}(y_1, \mathbf{u}) \in \mathbb{R}$, $1 \leq \mu \leq n+1$. Therefore, a computation using that $\Gamma_{\mu\nu, l}^{ij, \sigma} \neq 0$ only when $i = j = l$ (as the product metric is block diagonal and, thus, it does not mix different copies of \mathcal{E}_W^n), and $T_{y_k} \pi_W$ does not depend on y_1 , leads to $\nabla_{g \cdot \mathbf{u}}^{\mathbf{g}} f_{sa,k} = \psi_1 \cdot T_{y_1} \pi_W \cdot \mathbf{u}$. ■

Now let us compute the Lie derivative of the auxiliary function h_1 with respect to the control field.

Proposition 3.2: Consider the function introduced in (9), $h_1 \in C^1(\mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k)$. Then, for each $(t, \mathbf{y}, \mathbf{u}) \in \mathbb{R}_{\geq 0} \times (\mathcal{E}_W^n)^k \times \mathbb{R}^{n+1}$, we have:

$$(L_g h_1)(t, \mathbf{y}) \cdot \mathbf{u} = -(dh_0)_{y_k}(\psi_1(t, \mathbf{y}) \cdot T_{y_1} \pi_W \cdot \mathbf{u}).$$

where $\alpha' : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}_{> 0}$ denotes the derivative of α .

Proof: Note that $\nabla_{g(t, y_1) \cdot \mathbf{u}}^{\mathbf{g}} f_{\text{ff}}(t, y_k) = 0$, since $g(t, y_1) \cdot \mathbf{u}$ does not depend on y_k . Similarly, we have:

$$(dh)_{\mathbf{y}}(g(t, y_1) \cdot \mathbf{u}) = (dh_0)_{y_k}(0) = 0.$$

From this and (7), we obtain:

$$(L_g h_1)(t, \mathbf{y}) \cdot \mathbf{u} = -(dh_0)_{y_k}(\psi_1(t, \mathbf{y}) \cdot T_{y_1} \pi_W \cdot \mathbf{u}),$$

where we used Lemma 3.2 for the last equality. ■

We are ready to give a sufficient condition for h to be a second-order CBF for the controlled dynamics of the transformer.

Theorem 3.1: Let $h \in C^2((\mathcal{E}_W^n)^k)$ be as in (8), h_1 be as in (9), and $\mathcal{D} \subset (\mathcal{E}_W^n)^k$ open be such that:

$$\mathcal{S} = h^{-1}(\mathbb{R}_{\geq 0}) \cap h_1^{-1}(\mathbb{R}_{\geq 0})^{\circ} \subset \mathcal{D}.$$

If the following assumptions hold true:

- 1) $(dh_0)_{y_k} \neq 0$ for each $\mathbf{y} \in \mathcal{D} - \text{int } \mathcal{S}$, where $\text{int } \mathcal{S}$ denotes the interior of \mathcal{S} .
- 2) $\psi_1(t, \mathbf{y}) : T_{y_1} \mathcal{E}_W^n \rightarrow T_{y_k} \mathcal{E}_W^n$ is surjective for each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times \mathcal{D}$.

Then h is a second-order CBF for the control dynamics of the transformer (7).

Proof: From Definition 1.3, h is a second-order CBF for controlled dynamics of the transformer if $(L_g h)(t, \mathbf{y}) = 0$ and (3) holds for each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \cap \mathcal{D}$. The first condition is straightforward; namely:

$$(L_g h)(t, \mathbf{y}) = (dh)_{\mathbf{y}}(g(t, y_1)) = (dh_0)_{y_k}(0) = 0.$$

For the second condition, we choose $\alpha_1 \in \mathcal{K}_\infty^e$ as $\alpha_1(r) = (|L_m|/h_m)r$, where we denote $L_m = \inf_{(t,\mathbf{y}) \in \mathbb{R}_{\geq 0} \times \mathcal{H}} (L_f h_1)(t, \mathbf{y})$, as well as $h_m = \inf_{(t,\mathbf{y}) \in \mathbb{R}_{\geq 0} \times \mathcal{H}} h_1(t, \mathbf{y})$. Note that the following set:

$$\mathcal{H} = (dh_0 \circ \pi_k)^{-1}(0) = \{\mathbf{y} \in \mathcal{D} \mid (dh_0)_{y_k} = 0_{y_k}\},$$

where $\pi_k : \mathcal{D} \rightarrow \mathcal{E}_W^n$ is the projection onto the k -th copy, is closed (in \mathcal{D} and, thus, in $(\mathcal{E}_W^n)^k$ as \mathcal{D} is open). Moreover, assumption 1) implies that $\mathcal{H} \subset \text{int } \mathcal{S}$. Hence, the closeness of \mathcal{H} , together with the fact that $h_1(t, \mathbf{y}) > 0$ for each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times \text{int } \mathcal{S}$, ensure that $h_m > 0$. Given $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times \mathcal{D}$, we distinguish two cases:

1) $\mathbf{y} \in \mathcal{D} - \mathcal{H}$. Consider the following map:

$$\mathbb{R}^{n+1} \rightarrow \mathbb{R}, \quad \mathbf{u} \mapsto (L_g h_1)(t, \mathbf{y}) \cdot \mathbf{u},$$

From Proposition 3.2, as well as the facts that $(dh_0)_{y_k} \neq 0$ and $\psi_1(t, \mathbf{y})$ is surjective, we conclude that the previous map is surjective and, thus, (3) holds.

2) $\mathbf{y} \in \mathcal{H}$. By construction, $h_1(t, \mathbf{y}) \geq h_m$, whence:

$$\alpha_1(h_1(t, \mathbf{y})) \geq |L_m| \geq -(L_f h_1)(t, \mathbf{y}),$$

and (3) holds. ■

By gathering Proposition 3.1 and Theorem 3.1, we arrive at the main result of the paper.

Corollary 3.1: Let $\mathcal{D} \subset (\mathcal{E}_W^n)^k$ open be such that:

$$\mathcal{S} = h^{-1}(\mathbb{R}_{\geq 0}) \cap h_1^{-1}(\mathbb{R}_{\geq 0})^\circ \subset \mathcal{D},$$

with h and h_1 as in (8) and (9), respectively. If $(dh_0)_{y_k} \neq 0$ for each $\mathbf{y} \in \mathcal{D} - \text{int } \mathcal{S}$, then h is a second-order CBF for the control dynamics of the transformer (6) for almost every $(\mathbf{P}(t), \mathbf{U}(t)) \in (\mathbb{R}^{(n+1) \times (n+1)})^{2h}$, $t \in \mathbb{R}_{\geq 0}$.

Note that $\mathcal{S} \subset h^{-1}(\mathbb{R}_{\geq 0}) \subset \{\mathbf{y} \in (\mathcal{E}_W^n)^k \mid h_0(y_k) \leq \gamma\}$. This guarantees that the toxicity of the output remains below the threshold γ provided the initial condition lies in the safe set \mathcal{S} and the control inputs are chosen to be in $K_{\text{CBFH}}(t, \mathbf{y})$ for each $(t, \mathbf{y}) \in \mathbb{R}_{\geq 0} \times \mathcal{S}$.

The strength of this result lies in the weak assumptions placed on the function h_0 , which are a consequence of the new transformer architecture proposed in this paper. We simply ask for the absence of points (interpreted as sentences) where the linearization of h_0 is identically zero and thus where h would fail to describe how toxicity changes in a small neighborhood of that sentence.

IV. CONCLUSION

In this paper, the problem of toxicity in LLM text generation has been addressed. Firstly, a continuous-time model explicitly describing the feedforward layers of transformers was derived. We then proposed a new architecture in which the first token can be steered and is regarded as a controlled token. This enabled us to show that any C^2 function satisfying a mild non-degeneracy condition can be used to construct a CBF for almost every choice of the parameters of the transformer (key, query and value matrices).

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [2] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever *et al.*, "Improving language understanding by generative pre-training," *OpenAI*, 2018.
- [3] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [4] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.
- [5] D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi *et al.*, "Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning," *arXiv preprint arXiv:2501.12948*, 2025.
- [6] M. Sandler, H. Choung, A. Ross, and P. David, "A linguistic comparison between human and chatgpt-generated conversations," in *International Conference on Pattern Recognition and Artificial Intelligence*. Springer, 2024, pp. 366–380.
- [7] S. Gehman, S. Gururangan, M. Sap, Y. Choi, and N. A. Smith, "Realtotoxicityprompts: Evaluating neural toxic degeneration in language models," *arXiv preprint arXiv:2009.11462*, 2020.
- [8] H. Luo and L. Specia, "From understanding to utilization: A survey on explainability for large language models," *arXiv preprint arXiv:2401.12874*, 2024.
- [9] B.-D. Oh and W. Schuler, "Token-wise decomposition of autoregressive language model hidden states for analyzing model predictions," *arXiv preprint arXiv:2305.10614*, 2023.
- [10] A. Madsen, S. Chandar, and S. Reddy, "Are self-explanations from large language models faithful?" *arXiv preprint arXiv:2401.07927*, 2024.
- [11] X. Suau, P. Delobelle, K. Metcalf, A. Joulin, N. Apostoloff, L. Zappella, and P. Rodríguez, "Whispering experts: Neural interventions for toxicity mitigation in language models," *arXiv preprint arXiv:2407.12824*, 2024.
- [12] T. S. Luong, T.-T. Le, L. N. Van, and T. H. Nguyen, "Realistic evaluation of toxicity in large language models," *arXiv preprint arXiv:2405.10659*, 2024.
- [13] P. Chao, A. Robey, E. Dobriban, H. Hassani, G. J. Pappas, and E. Wong, "Jailbreaking black box large language models in twenty queries," *arXiv preprint arXiv:2310.08419*, 2023.
- [14] N. Karagodin, Y. Polyanskiy, and P. Rigollet, "Clustering in causal attention masking," *arXiv preprint arXiv:2411.04990*, 2024.
- [15] Á. Rodríguez Abella, J. P. Silvestre, and P. Tabuada, "The asymptotic behavior of attention in transformers," *arXiv preprint arXiv:2412.02682*, 2024.
- [16] E. Cheng, M. Baroni, and C. A. Alonso, "Linearly controlled language generation with performative guarantees," *arXiv preprint arXiv:2405.15454*, 2024.
- [17] A. D. Ames, X. Xu, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs for safety critical systems," *IEEE Transactions on Automatic Control*, vol. 62, no. 8, pp. 3861–3876, 2016.
- [18] Y. Miyaoka and M. Inoue, "Cbf-llm: Safe control for llm alignment," *arXiv preprint arXiv:2408.15625*, 2024.
- [19] Q. Nguyen and K. Sreenath, "Exponential control barrier functions for enforcing high relative-degree safety-critical constraints," in *2016 American Control Conference (ACC)*. IEEE, 2016, pp. 322–328.
- [20] W. Xiao and C. Belta, "High-order control barrier functions," *IEEE Transactions on Automatic Control*, vol. 67, no. 7, pp. 3655–3662, 2021.
- [21] C. Yun, S. Bhojanapalli, A. S. Rawat, S. J. Reddi, and S. Kumar, "Are transformers universal approximators of sequence-to-sequence functions?" *arXiv preprint arXiv:1912.10077*, 2019.
- [22] J. Zhao, K. Chen, X. Yuan, and W. Zhang, "Prefix guidance: A steering wheel for large language models to defend against jailbreak attacks," *arXiv preprint arXiv:2408.08924*, 2024.
- [23] B. Mityagin, "The zero set of a real analytic function," 2015. [Online]. Available: <https://arxiv.org/abs/1512.07276>