

SEMANTIC PROBABILISTIC CONTROL OF LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

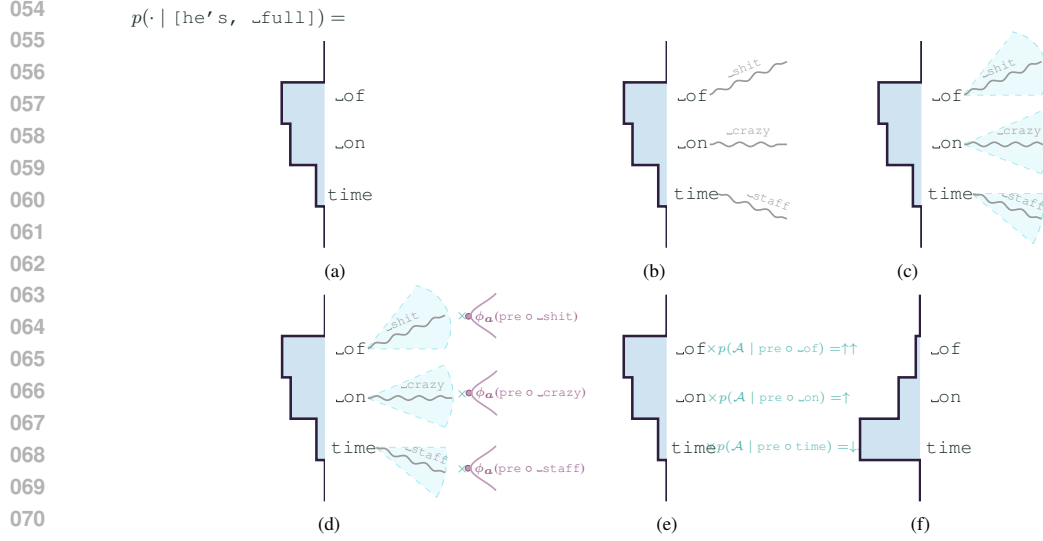
Semantic control entails steering LM generations towards satisfying subtle non-lexical constraints—*e.g.*, toxicity, sentiment, or politeness—attributes that can be captured by a sequence-level *verifier*. It can thus be viewed as sampling from the LM distribution conditioned on the target attribute, a computationally intractable problem due to the non-decomposable nature of the verifier. Existing approaches to LM control either only deal with syntactic constraints which cannot capture the aforementioned attributes, or rely on sampling to explore the conditional LM distribution, an ineffective estimator for low-probability events. In this work, we leverage a verifier’s gradient information to efficiently reason over *all* generations that satisfy the target attribute, enabling precise steering of LM generations by reweighing the next-token distribution. Starting from an initial sample, we create a local LM distribution favoring semantically similar sentences. This approximation enables the tractable computation of an *expected sentence embedding*. We use this expected embedding, informed by the verifier’s evaluation at the initial sample, to estimate the probability of satisfying the constraint, which directly informs the update to the next-token distribution. We evaluated our approach on the tasks of controlling the toxicity, sentiment, and topic-adherence of LMs yielding generations satisfying the constraint with high probability without degrading their quality.

1 INTRODUCTION

Despite the unprecedented capabilities of LMs, steering their generations towards specific syntactic or semantic constraints remains an unsolved challenge (Sun et al., 2023; Liu et al., 2024). Syntactic (or *lexical*) constraints define at each position in the sequence the set of admissible tokens that, taken together, constitute a valid string under the constraint. A common use case for such constraints is to generate output in some formal language, for example, structured data, API calls, or code snippets (Geng et al., 2025). Syntactic constraints are *easy* to deal with in a very precise sense: through knowledge compilation (Darwiche & Marquis, 2002), we can efficiently capture the computation graph of generations satisfying the constraint, which we can then proceed to *probabilistically* reason about, exactly when possible (Ahmed et al., 2022), and otherwise approximately (Willard & Louf, 2023; Zhang et al., 2024a; Koo et al., 2024; Lundberg et al., 2024; Ahmed et al., 2025).

Semantic (or *non-lexical*) constraints, on the other hand, are often defined in terms of sequence-level, non-decomposable classifiers, or *verifiers*, often complex neural networks, that assign non-negative scores to sequences of tokens. In that sense, semantic constraints are doubly hard: we have to contend with not only the hardness of probabilistic reasoning but also the lack of a tractable representation of the constraint over which to reason. Semantic constraints encompass use cases in which we might wish to control sequence-level properties of generations that are hard to capture in formal language, *e.g.*, controlling toxicity, sentiment, or topic in creative writing; targeting outputs deemed favorable by a verifier for reasoning, or generating code exhibiting stylistic requirements (Geng et al., 2025).

Existing approaches to semantic control of LMs largely fall into four families. *sample-reweigh*, known as Best-of- n (Stiennon et al., 2020a), generates complete sequences that are reweighed by the potential function, returning the highest scoring sequence, but does not incorporate constraints during generation and therefore require exponentially many samples for low-probability attributes. *Sequential Monte Carlo (SMC)* approaches, propagate a population of partial sequences using LM likelihoods and constraint potentials (Zhao et al., 2024; Loula et al., 2025) but often require many particles and



suffer from degeneracy on longer sequences. *Token-level reweighting* approaches (Yang & Klein, 2021; Krause et al., 2021; Liu et al., 2021, *inter alia*) intervene on next-token probabilities based on classifier or expert models but require expensive training or fine-tuning of auxiliary predictors, and can often be myopic. *Activation-steering* approaches (Dathathri et al., 2020; Han et al., 2024, *inter alia*) intervene directly on the LM’s hidden states using learned attribute directions but generally provide coarse global shifts rather than fine-grained, token-level probabilistically grounded control.

In this work, in a departure from the aforementioned approaches, we propose performing *exact inference in an approximate model* (Koller & Friedman, 2009). We propose **Semantic Control Estimator (SConE)**, which leverages the gradient information of a verifier to tractably perform inference over *all* generations satisfying the constraint, allowing precise steering of LM generations by reweighing each probable next token according to its probability of satisfying the constraint. More precisely, starting from a *lookahead* sample, we construct a *local, contextualized* LM distribution that assigns a higher probability to semantically similar sentences and a lower probability to semantically dissimilar ones. We show that we can *tractably* and efficiently compute the expected embedding of *all* sentences w.r.t. this approximate LM distribution. Computing the expected embedding allows us to estimate the *expected attribute probability* using a single LM sample and a single evaluation of the verifier by distributing first-order information regarding the verifier over the expected embedding. The next-token distribution is reweighed by *expected attribute probability* and renormalized to obtain the *attributed reweighted* next-token distribution. Computationally, the expected embedding can be computed in $O(1)$ vectorized time, whereas the lookahead sample can be drawn efficiently by utilizing an auxiliary model¹ to unmask future tokens paired with HogWild! (asynchronous) Gibbs sampling (Niu et al., 2011; Smola & Narayanamurthy, 2010), with synchronization frequency trading off accuracy for efficiency. A high-level overview of our approach is given in Figure 1.

We evaluated our proposed approach on the tasks of controlling the toxicity and sentiment of LM generations, as well as on controlling the topic of generations. We observed that our approach was far more likely to satisfy the constraint compared to previous approaches, without compromising the

¹We use ModernBERT (Warner et al., 2025) in our experiments.

quality of the LM generations, as measured by perplexity and unigram diversity. Our proposed method is inference-time, requires no fine-tuning, and can be easily integrated with syntactic constraints.²

2 LEVELS OF CONTROL: FROM SYNTACTIC TO SEMANTIC CONSTRAINTS

We denote an LM generation of arbitrary length T by $\mathbf{y}_{1:T} := [y_1 y_2 \dots y_T]$, where y_i is the instantiation of random variable Y_i over tokens at time i and takes values from a vocabulary $\mathbb{V} = \{1, \dots, V\}$.

An LM generation can be subject to one of two types of constraints: syntactic and semantic. Syntactic (or *lexical*) constraints comprise sets of rules, typically expressed using logical connectives or in some formal language, that restrict the set of permissible values assumed by a random variable Y_i such that there exists some completion $\mathbf{y}_{>i}$ of the sentence that satisfies the syntactic constraint β , given the current prefix $\mathbf{y}_{1:i}$, or to state it more formally

$$\exists \mathbf{y}_{>i} \beta_{|\mathbf{y}_{1:i}} \quad (1)$$

An example of such constraint could be a simple logical sentence that disallows an expression deemed inappropriate to appear as part of an LM’s generation, e.g., $\neg(y_i = \text{“full”} \wedge y_{i+1} = \text{“of”} \wedge y_{i+2} = \text{“shit”})$ (Ahmed et al., 2023). Syntactic constraints offer an attractive opportunity for parallelization: we are able to *compile* syntactic constraints into computational graphs that reuse solutions to subproblems to efficiently capture the space of all satisfying assignments. Traversing these computation graphs amounts to efficient parallel evaluation across an exponential number of possible continuations (Choi et al., 2020; Vergari et al., 2021) enabling us to tractably compute Equation (1).

Semantic (or *non-lexical*) constraints, on the other hand, presuppose that LM generations satisfy certain *attributes* (e.g., toxicity, politeness, or positive sentiment). Such attributes are often hard to ascertain lexically, or in terms of surface-level features that can be captured using a formal language, e.g., “he’s got some attitude!” invokes a snarky tone that is hard to attribute to any particular token in the generation. Rather, given a target attribute a , we suppose access to a *sequence-level verifier for a* , which we denote by ϕ_a , that given a sequence $\mathbf{y}_{1:T}$ assigns a binary value, either 0 or 1, to the sequence $\mathbf{y}_{1:T}$, i.e., $\phi_a(\mathbf{y}_{1:T}) \in \{0, 1\}$. We can then define \mathcal{A} as the set of *all* sequences $\mathbf{y}_{1:T}$ that satisfy the attribute a , i.e., $\mathcal{A} := \{\mathbf{y}_{1:T} \mid \phi_a(\mathbf{y}_{1:T}) = 1\}$. Unlike syntactic constraints, semantic constraints, often implemented as complex neural networks, are not amenable to the form of compilation that enables us to efficiently capture the set of all satisfying assignments. In fact, compiling even a single neuron is known to be NP-hard (Shi et al., 2020). Computing Equation (1) would thus require that we enumerate every possible continuation, score it using the verifier, discard continuations for which the attribute does not hold and renormalize, which is intractable.

Prologue. In what follows, we will relax the verifier ϕ_a for an attribute a to be *probabilistic*. We will then frame the problem of semantic control as a probabilistic inference problem whereby *we are interested in the posterior LM distribution subject to a semantic constraint*. We will show that the problem can be reduced to that of computing an *expected attributed probability*. We will then show how to estimate the expectation by *performing exact and efficient probabilistic inference in an approximate distribution induced by a singular sample and a singular evaluation of the verifier*.

3 EXPECTED ATTRIBUTE PROBABILITY

We start by assuming access to the LM distribution, denoted by p , a sequence-level verifier ϕ_a for attribute a , and a prefix $\mathbf{y}_{1:i}$ where each token y_j assumes values in vocabulary \mathbb{V} . Our goal is then to sample from the LM distribution p a generation $\mathbf{y}_{i+1:T}$ subject to the constraint that the attribute a holds on the entire sequence i.e., $\phi_a(\mathbf{y}_{1:i} \circ \mathbf{y}_{i+1:T}) \in \{0, 1\}$. That entails sampling a generation that fulfills two distinct desiderata: we expect the generation to be linguistically sound, or fluent as measured by a model’s perplexity, *and* to satisfy attribute a . That is, we are interested in sampling from the LLM distribution conditioned on the event that the sample belongs to the set of *all* sequences $\mathbf{y}_{1:T}$ that satisfy the attribute a , which we denote by $\mathcal{A} := \{\mathbf{y}_{1:T} \mid \phi_a(\mathbf{y}_{1:T}) = 1\}$. We then write

$$p(\mathbf{y}_{i+1:T} \mid \mathcal{A}, \mathbf{y}_{1:i}) \stackrel{(a)}{=} \frac{p(\mathbf{y}_{i+1:T}, \mathcal{A} \mid \mathbf{y}_{1:i})}{p(\mathcal{A} \mid \mathbf{y}_{1:i})} \stackrel{(b)}{=} \frac{p(\mathbf{y}_{i+1:T} \mid \mathbf{y}_{1:i}) \cdot \phi_a(\mathbf{y}_{1:i} \circ \mathbf{y}_{i+1:T})}{\sum_{\mathbf{y}_{i+1:T}} p(\mathbf{y}_{i+1:T} \mid \mathbf{y}_{1:i}) \cdot \phi_a(\mathbf{y}_{1:i} \circ \mathbf{y}_{i+1:T})}, \quad (2)$$

²Our code and scripts to reproduce all numbers will be made publicly available upon acceptance.

where equality (a) follows by the definition of conditional probability, and equality (b) follows by the definition of marginal probability. Intuitively, Equation (2) gives us a simple, albeit impractical, recipe for sampling from the LM distribution conditioned on attribute \mathbf{a} : we enumerate all possible generations given the prefix, zeroing out all generations that violate \mathbf{a} according to $\phi_{\mathbf{a}}$, followed by renormalization. In practice, for a given input $\mathbf{y}_{1:T}$ and attribute \mathbf{a} , there is some *uncertainty* associated with $\phi_{\mathbf{a}}(\mathbf{y}_{1:T})$. That is, we will assume access to a model’s estimate $p(\phi_{\mathbf{a}}(\mathbf{y}_{1:T}) = 1) \in [0, 1]$ of whether $\mathbf{y}_{1:T}$ satisfies attribute \mathbf{a} . Consequently, in a slight abuse of notation, we will redefine $\phi_{\mathbf{a}}(\cdot)$ to be $p(\phi_{\mathbf{a}}(\mathbf{y}_{1:T}) = 1)$, which should henceforth be thought of as a *probabilistic* verifier for the attribute \mathbf{a} . Under this new definition of $\phi_{\mathbf{a}}(\cdot)$, Equation (2) can be seen as reweighing each continuation with the probability of satisfying attribute \mathbf{a} , followed by renormalizing the distribution.

State-of-the-art LMs, such as Llama 3 (Grattafiori et al., 2024) and GPT-4 (Achiam et al., 2024)) are autoregressive, so it is useful to rewrite Equation (2) in terms of the the next-token distribution,

$$p(y_{i+1} | \mathcal{A}, \mathbf{y}_{1:i}) = \frac{p(y_{i+1} | \mathbf{y}_{1:i}) \cdot p(\mathcal{A} | \mathbf{y}_{1:i} \circ y_{i+1})}{p(\mathcal{A} | \mathbf{y}_{1:i})} \quad (3)$$

$$= \frac{p(y_{i+1} | \mathbf{y}_{1:i}) \mathbb{E}_{p(\cdot | \mathbf{y}_{1:i+1})} [\phi_{\mathbf{a}}(\mathbf{y}_{1:i} \circ y_{i+1:T})]}{\mathbb{E}_{p(\cdot | \mathbf{y}_{1:i})} [\phi_{\mathbf{a}}(\mathbf{y}_{1:i} \circ y_{i+1:T})]}, \quad (4)$$

where Equation (3) follows by the definition of conditional probability and Equation (4) follows by the definition of marginal probability and expectations. It is important to note that, since \mathcal{A} is defined as the set of all sequences $\mathbf{y}_{1:T}$ that satisfy \mathbf{a} , the expectations—both in the numerator and in the denominator—range over sequences of length T , requiring that we marginalize over all future continuations of length $T - i$ and $T - (i + 1)$, respectively. Intuitively, at every generation step we need to “look ahead” to determine the probability that the constraint is violated given the current choice of next token. If the probability is high, we discount the current choice, and if it is low, then we reinforce the current choice. Previous methods approached computing the intractable expectation in Equation (4) by learning lookahead functions, also termed future discriminators (Yang & Klein, 2021), that provide a locally evaluated surrogate for the global attribute probability; approximating it with a ratio of generative classifiers (Krause et al., 2021), or attribute experts (Liu et al., 2021) that myopically reweigh tokens; or through sampling (Loula et al., 2025; Lew et al., 2023; Zhao et al., 2024), requiring many particles and suffering from high variance. In what follows, we show how to compute an approximation of the above expectation in closed form by relaxing the target distribution, yielding a globally-aware, low-variance estimate of the attribute probability with only a few samples.

4 SEMANTIC PROBABILISTIC CONTROL

The computational hardness of the expectations that we introduced in Equation (4) can intuitively be attributed to the *lack of an intrinsic structure* along two distinct dimensions which we detail below.

First, is the *lack of structure to the distribution*. Consider computing the probability that a sequence of length T ends in the word “love”. Computing such a probability under the autoregressive distribution requires that we marginalize over all possible sequences ending in “love”, roughly $O(|\mathbb{V}|^T)$. In fact, computing such probability is known to be computationally intractable (Roth, 1993). Contrast that with a fully-independent³ distribution, where we can simply query the network for the probability of a given token in constant time. Clearly there is a tension here: fully-independent distributions, while easier to reason about, are not expressive and therefore do not make for good LMs, whereas autoregressive distributions are harder to reason about, but a lot more expressive, yielding SoTA LMs.

The second dimension is the *lack of structure to the constraint*. Recall our prior assumption that $\phi_{\mathbf{a}}$ is a neural network. This assumption turns out to have serious computational implications, as prior work has shown that unlike many other tractable probabilistic models, neural networks happen to be computationally intractable to decompose over sequences (Shi et al., 2020)⁴. That is, given $\phi_{\mathbf{a}}(\mathbf{y}_{1:i})$ for a prefix $\mathbf{y}_{1:i}$, we know of no way of efficiently extending $\phi_{\mathbf{a}}(\mathbf{y}_{1:i})$ to $\phi_{\mathbf{a}}(\mathbf{y}_{1:i} \circ y_{i+1})$ by only processing the new element y_{i+1} and reusing the result of the previous evaluation $\phi_{\mathbf{a}}(\mathbf{y}_{1:i})$.

³where $p(\mathbf{y}_{1:T}) = \prod_{i=1}^T p(y_i)$, i.e., the probability of a token is independent from all other tokens.

⁴in fact, the problem remains intractable even assuming $\phi_{\mathbf{a}}$ is a single neuron (Khosravi et al., 2019).

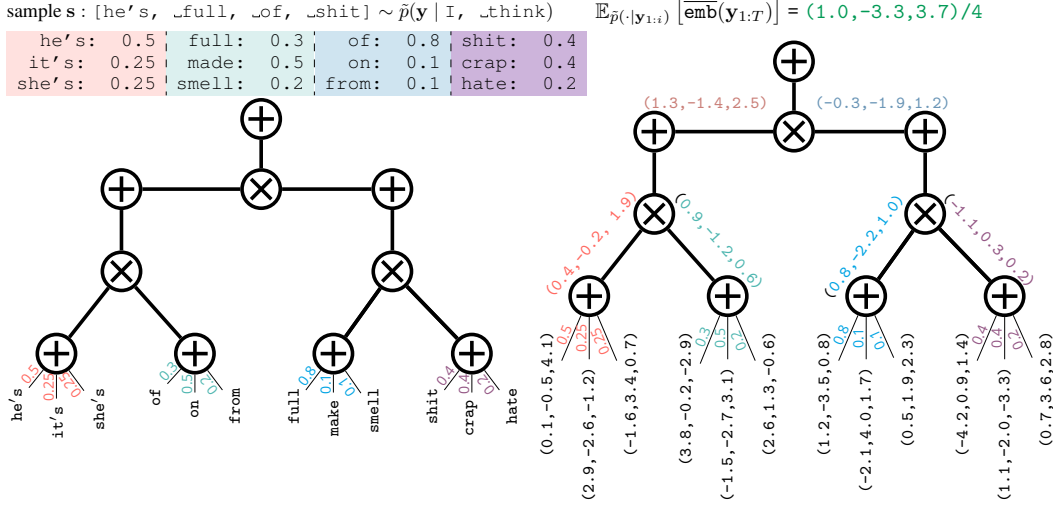


Figure 2: **A technical overview of our approach.** (top left) We start by sampling an (approximate) generation s using Gibbs sampling \tilde{p} conditioned on the prefix from the model’s marginal conditionals, $p(\mathbf{y}_i | \mathbf{y}_{-i}) \forall i$. Conditioned on s , the model’s marginal conditionals induce a distribution on all generations, assigning higher probabilities to similar sentences and lower probabilities for dissimilar sentences, visualized for the top-3 tokens. (bottom left) We can parameterize a *circuit* using the above distribution, yielding a closed-form, tractable representation of probability distribution defined in ??, where read left to right, every leaf node corresponds to a categorical distribution on \mathbf{y}_i (right) Such a representation enables us to compute the expected embeddings w.r.t. the distribution in the neighborhood of the sample s by substituting token embedding for corresponding embeddings at leaf nodes, computing weighted sums of embeddings at sum nodes, and taking sums at product nodes. We can plug the expected embedding into Equation (11) to yield the constraint probability.

4.1 A LOCALLY CONTEXTUALIZED DISTRIBUTION

To sidestep the hardness of the autoregressive distribution, we move towards the tractability of fully-independent distributions, while retaining as much of the contextual information. Therefore, we consider the *pseudolikelihood* of a sentence $\mathbf{y}_{1:T}$ under the model (Besag, 1975; Ahmed et al., 2023)

$$p(\mathbf{y}_{1:T}) \approx \tilde{p}(\mathbf{y}_{1:T}) := \prod_i p(\mathbf{y}_i | \mathbf{y}_{-i}), \quad (5)$$

where \mathbf{y}_{-i} denotes $\mathbf{y}_1, \dots, \mathbf{y}_{i-1}, \mathbf{y}_{i+1}, \dots, \mathbf{y}_n$. Unfortunately, Equation (5) remains intractable. The key issue is that the standard pseudolikelihood depends sentence-specific masked contexts \mathbf{y}_{-i} . This requires computing a fresh set of conditionals $\{p(\cdot | \mathbf{y}_{-i})\}_{i=1}^T$ for any given sentence $\mathbf{y}_{1:T}$. Moreover, these scores are incomparable since they arise from incompatible conditional distributions. Instead, we fix a reference sentence $\tilde{\mathbf{y}}$, evaluating all candidates using the same masked contexts $\tilde{\mathbf{y}}_{-i}$, giving us a *locally-contextualized distribution* requiring only T masked-LM forward passes for *all* sentences

$$\tilde{p}_{\tilde{\mathbf{y}}}(\mathbf{y}) := \prod_i p(\mathbf{y}_i | \tilde{\mathbf{y}}_{-i}), \quad (6)$$

which can be thought of as the *contextualized probability* of a sentence \mathbf{y} given the context $\tilde{\mathbf{y}}$. That is, ?? calculates the probability of sequence \mathbf{y} by taking the product of probabilities of each token \mathbf{y}_i , crucially conditioning each token \mathbf{y}_i not on the preceding tokens of \mathbf{y} , but on the context surrounding position i within $\tilde{\mathbf{y}}$ (specifically, $\tilde{\mathbf{y}}$ excluding its i -th token, denoted $\tilde{\mathbf{y}}_{-i}$). Therefore, $\tilde{\mathbf{y}}$ acts as a *contextual anchor* for evaluating \mathbf{y} under this measure. Intuitively, we expect sentences \mathbf{y} that structurally align with the specific token-level contexts provided by $\tilde{\mathbf{y}}$ to yield a higher $\tilde{p}_{\tilde{\mathbf{y}}}(\mathbf{y})$. In a slight abuse of notation, we omit the dependence on $\tilde{\mathbf{y}}$ when it’s not necessary for ease of exposition.

4.2 BRIDGING SAMPLES AND EXPECTATIONS: A TANGENTIAL VIEW

Next, we will turn our attention to address the *hardness of the verifier* ϕ_a . In particular, given an LM sample $s \sim p(\mathbf{y}_{i+1:T} | \mathbf{y}_{1:i})$ and access to a verifier ϕ_a , we will leverage gradient information obtained

Algorithm 1 SConE

```

1: Input: Verifier  $\phi_a$ , LM distribution
    $p(y_i | y_{1:i})$ , prefix  $y_{1:i}$ , max length  $T$ 
2: Output:  $p(y_{i+1} | \mathcal{A}, y_{1:i})$ 
   ▷ Expand the batch to include top-k tokens
3:  $\text{top}_k = \arg \max_k p(y_i | y_{1:i})$ 
4:  $y_{1:i+1} = y_{1:i}.\text{expand}(n, \text{top}_k)$ 
   ▷ Get  $N$  samples  $\tilde{s}$  from  $p(y_{i+2:T} | y_{1:i+1})$ 
5:  $\tilde{s}^1, \dots, \tilde{s}^N \sim \text{GibbsSampler}(y_{1:i+1}, p)$ 
   ▷ Estimate prob  $q$  of satisfying constraint
6:  $q = \text{zeros}(\text{top}_k)$ 
7: for each  $\tilde{s}$  in  $\tilde{s}^1, \dots, \tilde{s}^N$  do
8:    $\tilde{p}_{\text{cond}} = \text{CondMarginals}(p, \tilde{s}_{i+2:T})$ 
9:    $\nabla \phi_a = \text{LinearizeVerifier}(\phi_a, \tilde{s})$ 
10:   $q[\tilde{s}_{i+1}] += \text{EstimateProb}(\tilde{p}_{\text{cond}}, \phi_a, \nabla \phi_a)$ 
11: end for
   ▷ Renormalize  $q$ 
12:  $\log q = q.\text{log\_softmax}()$ 
   ▷ Reweight the LM distribution
13:  $w = \log p(y_{i+1} | y_{1:i}) + \log q$ 
14:  $p^* = \text{Categorical}(\text{weights} = w)$ 
15: return  $p^*$ 

```

Algorithm 2 LinearizeVerifier

```

1: Input: Verifier  $\phi_a$ , Sample  $s$ 
2: Output: Gradient of  $\phi_a$  w.r.t.  $s$  embedding
   ▷ Obtain embeddings for  $s$ 
3:  $\text{emb\_layer} = \phi_a.\text{get\_input\_embeddings}()$ 
4:  $\text{emb} = \text{emb\_layer}(s)$ 
   ▷ Collect gradient of  $\phi_a$  w.r.t. to  $\text{emb}$ 
5:  $\text{score} = \phi_a(\text{emb}).\text{sum}()$ 
6:  $\text{grad} = \text{autograd.grad}(\text{score}, \text{emb})$ 
7: return  $\text{grad}$ 

```

Algorithm 3 EstimateProb

```

1: Input: Conditional marginals  $\tilde{p}_{\text{cond}}$ , Verifier  $\phi_a$ , Gradient  $\nabla_{\text{emb}(s)} \phi_a$ ,  $\text{embs} := [\text{emb}(y_{i,1}), \dots, \text{emb}(y_{i,|V|})]$ , score  $\phi_a(s)$ ,  $T$ 
2: Output:  $p(\mathcal{A} | y_{1:i})$ 
   ▷ Compute expected embedding
3:  $\text{exe} = 0$ 
4: for  $i$  in  $1, \dots, T$  do
5:    $\text{exe} += \text{embs}[\dots, \text{None}] \cdot \tilde{p}_{\text{cond}}[:, i : i + 1, :]$ 
6: end for
7:  $\text{exe} = \text{exe.mean}(0)$ 
   ▷ First-order Taylor expansion about  $s$ 
8: return  $\phi_a(s) + \nabla_{\text{emb}(s)} \phi_a \cdot (\text{exe} - \overline{\text{emb}}(s))$ 

```

during the evaluation of $\phi_a(s)$, coupled with the locally contextualized distribution introduced in ??, to approximate $\mathbb{E}_{p(\cdot|y_{1:i})} [\phi_a(y_{1:T})]$, the expected attribute probability, *e.g.*, toxicity.

We start by approximating the expectation of the verifier w.r.t. the LM distribution as an expectation w.r.t. the *locally contextualized distribution* at a LM sample s , substituting \tilde{y} for s in ??

$$\mathbb{E}_{p(\cdot|y_{1:i})} [\phi_a(y_{1:T})] \approx \mathbb{E}_{\tilde{p}_s(\cdot|y_{1:i})} [\phi_a(y_{1:T})]. \quad (7)$$

This, however, does little by way of making the expectation tractable. Recall from our discussion in Section 4, that a neural network cannot tractably decompose over sequences (Shi et al., 2020). Therefore, computing expectations of even simple neural networks w.r.t. tractable distributions turns out to be computationally intractable (Khosravi et al., 2019). Intuitively, since the verifier ϕ_a does not decompose, computing the expectation in Equation (7) entails enumerating all sentences of length T and evaluating them through ϕ_a , although it no longer requires that we evaluate each under the LM.

Taking inspiration from the locally contextualized distribution, we consider what a *locally contextualized verifier* would look like. Hence, we consider a first-order Taylor expansion of ϕ_a at s

$$\phi_a(y_{1:T}) \approx \phi_a(s) + \nabla_s \phi_a(s) \cdot (y_{1:T} - s) \quad (8)$$

where the subtraction $(y_{1:T} - s)$ is to be understood component-wise at the level on which ϕ_a operates. In our setting, ϕ_a is a neural classifier that consumes token *embeddings* rather than discrete tokens. That is, the input to ϕ_a is a deterministic function of the LM output tokens. We therefore denote by $\text{emb} : \mathbb{V} \mapsto \mathbb{R}^d$ an embedding function that maps each token onto a d -dimension vector and let $\overline{\text{emb}}(y)$ denote the average token-wise embedding.⁵ Replacing the abstract difference $(y_{1:T} - s)$ with the corresponding difference in the verifier’s embedding space yields the concrete approximation

$$\phi_a(y_{1:T}) \approx \phi_a(s) + \nabla_{\text{emb}(s)} \phi_a(s) \cdot (\overline{\text{emb}}(y_{1:T}) - \overline{\text{emb}}(s)). \quad (9)$$

Taking an expectation under the approximate locally-contextualized distribution at the LM sample s

$$\mathbb{E}_{\tilde{p}(\cdot|y_{1:i})} [\phi_a(y_{1:T})] \approx \mathbb{E}_{\tilde{p}(\cdot|y_{1:i})} [\phi_a(s) + \nabla_{\text{emb}(s)} \phi_a(s) \cdot (\overline{\text{emb}}(y_{1:T}) - \overline{\text{emb}}(s))]. \quad (10)$$

Using the linearity of expectation, we can further simplify this expression, obtaining

$$\mathbb{E}_{\tilde{p}(\cdot|y_{1:i})} [\phi_a(y_{1:T})] \approx \phi_a(s) + \nabla_{\text{emb}(s)} \phi_a(s) \cdot (\mathbb{E}_{\tilde{p}(\cdot|y_{1:i})} [\overline{\text{emb}}(y_{1:T})] - \overline{\text{emb}}(s)), \quad (11)$$

⁵w.l.o.g, we assume this embedding can be extracted directly from the embedding layer of the verifier, *i.e.*, $\phi_a(s) := \phi_a(\text{emb}(s_1), \dots, \text{emb}(s_T))$.

expressing the expected verifier output in terms of the expected sentence embedding w.r.t. a locally contextualized distribution. We were thus able to **reduce the problem of estimating the constraint probability**, given by the expectations in Equation (4) **to the problem of computing an average sentence embedding** w.r.t. an approximate LM distribution \tilde{p} , followed by simple arithmetic operations. Next, we will show how to efficiently compute the *expected sentence embedding* in Equation (11).

4.3 FROM SEQUENCE PROBABILITIES TO AVERAGE EMBEDDINGS

In what follows, our goal will be to show that we can compute the expected sentence embedding w.r.t. the locally contextualized distribution, as in Equation (11), in time that is linear in the sequence length T . To make this possible, the computation graph used to aggregate token embeddings must satisfy certain structural constraints that allow expectations to decompose in a single pass. We first describe these conditions, and then show that, by construction, the computational graph for the expected embedding w.r.t. to our locally contextualized distribution abides by such structural constraints. To formalize these structural constraints, we appeal to the computational framework of circuits, in which a function is represented as a tractable parametric computation graph, hereafter referred to as a *circuit*. By imposing specific structural constraints on such circuits, we can guarantee that key probabilistic queries can be computed exactly and efficiently, providing a language for tractable reasoning.

Formally, a *circuit* p over variables \mathbf{Y} is a parameterized computational graph encoding a function $p(\mathbf{Y})$. Each node n in the graph encodes a parameterized function $p_n(\text{vars}(n))$ over variables $\text{vars}(n) \subseteq \mathbf{Y}$, also known as its *scope*. Each inner node in the graph is a sum or a product node, and each leaf node encodes a tractable input distribution over its scope. Each inner unit n (i.e., product or sum node) receives inputs from other units, denoted in (n) .

A circuit is *decomposable* if the inputs of every product node depends on disjoint sets of variables, i.e., for $n = c_1 \otimes c_2$, $\text{vars}(c_1) \cap \text{vars}(c_2) = \emptyset$. Intuitively, decomposable product nodes encode local factorizations over variables of the function. We assume that decomposable product nodes always have two inputs, a condition that is enforceable on any circuit with a polynomial increase in its size.

A second property is *smoothness*. A circuit is *smooth* if the inputs of every sum node depend on the same set of variables, i.e., for $n = \bigoplus_i \theta_i \cdot c_i$, $\text{vars}(c_i) = \text{vars}(c_j) \forall i, j$. Smoothness ensures that sum nodes represent mixture distributions with shared scope, so expectations propagate linearly through them. Decomposability and smoothness are sufficient and necessary for tractable integration over arbitrary sets of variables in a single pass, allowing larger integrals to decompose into smaller ones.

Locally Contextualized Distribution as a smooth and Decomposable Circuit. For a given reference LM sample \tilde{y} , the locally contextualized distribution $\tilde{p}_{\tilde{y}}$ induces conditional independence among the token variables across all time steps. This factorization implies a simple circuit with T independent distributions. For each position i , we introduce a product node whose scope is $\{y_i\}$. Beneath each product node sits a sum node representing the categorical distribution $\tilde{p}(y_i = v)$. Each sum node therefore has $|\mathcal{V}|$ leaf nodes, one for each token $v \in \mathcal{V}$, and evaluates to the corresponding embedding $\overline{\text{emb}}(v)$, see Figure 2. This circuit is *smooth*, since the children of each sum node share the same scope $\{y_i\}$, and *decomposable*, since the scopes of the product nodes are disjoint across time steps. **Expectation at sum and product nodes.** Smoothness implies linearity at sum nodes. Therefore, the expected embedding at time step i is simply the weighted average of token embeddings given by

$$\mathbb{E}_{\tilde{p}}[\overline{\text{emb}}(y_i)] = \sum_{v \in \mathcal{V}} \tilde{p}(y_i = v) \overline{\text{emb}}(v). \quad (12)$$

Whereas decomposability guarantees that embeddings across different positions can be aggregated independently. Using the standard averaging aggregator for sentence embeddings and linearity yields

$$\mathbb{E}_{\tilde{p}}[\overline{\text{emb}}(y_{1:T})] = \frac{1}{T} \sum_{i=1}^T \mathbb{E}_{\tilde{p}}[\overline{\text{emb}}(y_i)].$$

Closed-form computation and einsum implementation. Putting these expressions together gives

$$\mathbb{E}_{\tilde{p}}[\overline{\text{emb}}(y_{1:T})] = \frac{1}{T} \sum_{i=1}^T \left(\sum_{v \in \mathcal{V}} \tilde{p}(y_i = v) \overline{\text{emb}}(v) \right). \quad (13)$$

Table 1: **Evaluation of the toxicity of Llama-3.2 (1B) generations when steered to be non-toxic and toxic.** Results are reported over 400 arbitrary prompts from RealToxicityPrompts using a RoBERTa-based toxicity classifier (Logacheva et al., 2022). *PPL* denotes the perplexity, measured by Llama-3 (70B); *TTR* captures the unigram diversity; **Exp Max Toxicity** and **Toxic Prob** denote average worst toxicity and likelihood of generating a toxic output. We expect both metrics to be **lower** (\downarrow) when controlling for **non-toxic outputs** and **higher** (\uparrow) for **toxic** outputs.

Objective	Method	Toxic Prob. (\downarrow, \uparrow)			Exp. Max. Toxicity (\downarrow, \uparrow)			PPL (\downarrow)	TTR (\uparrow)
		Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Full
uncontrolled	random	37.25	10.00	64.50	37.11	13.17	61.05	12.18	81.57
	beamsearch	17.25	3.00	31.50	18.22	4.34	32.09	8.00	75.11
	top-k: 10	49.25	25.00	73.50	49.16	25.98	72.35	15.43	84.94
detoxify	AttrPrefix	50.00	26.50	73.50	48.21	26.10	70.32	18.30	88.68
	Few-shot	84.00	73.00	95.00	81.38	71.26	91.51	18.69	88.22
	BoN	2.75	1.00	4.50	4.90	1.91	7.89	15.46	85.74
	SConE (ours)	00.25	00.50	00.00	1.85	1.30	2.40	14.88	87.46
toxify	AttrPrefix	51.50	29.00	74.00	50.58	30.25	16.45	18.40	89.05
	Few-shot	95.00	91.00	99.00	92.19	88.63	95.74	18.95	88.35
	BoN	62.50	37.00	88.00	61.36	39.62	83.11	13.97	83.22
	SConE (ours)	93.75	88.00	99.50	91.15	85.75	96.55	23.87	81.12

Let $P \in \mathbb{R}^{T \times |\mathcal{V}|}$ be the matrix of local distributions with $P_{i,v} = \bar{p}(y_i = v)$, and let $E \in \mathbb{R}^{|\mathcal{V}| \times d}$ be the embedding matrix. Then the expected embeddings at all positions are given by the matrix product $PE \in \mathbb{R}^{T \times d}$, which corresponds exactly to the batched einsum, `einsum("tv, vd->td", P, E)`. Averaging across positions yields the expected sentence embedding. Hence, the entire computation corresponds to a single $\mathcal{O}(T)$ pass through a smooth and decomposable circuit.

Closing the loop. Our full algorithm is given in Algorithm 1. We start by truncating the next-token distribution using top- k or top- p . We proceed by simulating a continuation for each of the possible top- k tokens, each produced using a masked LM and Hogwild! Gibbs sampling⁶, to avoid expensive autoregressive sampling. We proceed by computing the contextualized probability of each sample \mathbb{V}_i and the gradient of the verifier w.r.t. the sample embedding $\nabla_{\text{emb}(\mathbf{s})} \phi_{\mathbf{a}}$, used to estimate the constraint probability. We reweigh the next-token distribution by the constraint probability, and renormalize.

5 EXPERIMENTS

We now turn to the empirical evaluation of our method across three open-ended generation settings, where we control for toxicity, positive sentiment, and topic adherence. In the following sections, we briefly describe baselines, metrics, and results, making additional details available in Appendix A.

5.1 CONTROLLABLE TOXICITY GENERATION

Setup. We evaluate control on 400 prompts from RealToxicityPrompts (Gehman et al., 2020), evenly split between toxic and non-toxic. For each prompt, we generate 25 continuations of up to 20 tokens under two settings—**toxification** and **detoxification**—and assess them with a RoBERTa-based toxicity classifier (Logacheva et al., 2022). We report *perplexity* (*PPL*) to measure grammaticality, *type token ratio* (*TTR*) to capture unigram diversity (Hess et al., 1984; Rosillo-Rodes et al., 2025), and two widely used toxicity metrics (Gehman et al., 2020): *Expected Maximum Toxicity*, defined as the average maximum toxicity observed across prompts, and *Toxicity Probability*, the average probability of generating at least one toxic continuation. As baselines, we compare SConE against 10 decoding-time methods, including six training-free methods—random, beamsearch, top-k: 10, AttrPrefix (Pei et al., 2023), Few-shot, BoN (Stiennon et al., 2020a)—and four training-based methods—PPLM (Dathathri et al., 2020), Fudge (Yang & Klein, 2021), DExperts (Liu et al., 2021), and LMSteer (Han et al., 2024).

Detoxification Results. Table 1 and Table 6 (in Appendix) present the results for Llama-3.2 (1B) and GPT2-medium as base models, respectively. Overall, we observe that the uncontrolled

⁶We refer the reader to Appendix G for more details.

Table 2: **Evaluation of GPT2-IMDB generations under positive sentiment constraints.** Results are computed on 600 prompts from the IMDB test set using a BERT-based sentiment classifier (Maas et al., 2011). In addition to PPL and TTR, reported sentiment metrics include average sentiment score (Rafailov et al., 2023; Amini et al., 2025), probability that all generations are positive (Sentiment Prob.), and expected minimum sentiment score (Exp. Min. Sentiment).

Method	Avg $\phi_{\text{sentiment}}$ (\uparrow)			Sentiment Prob. (\uparrow)			Exp. Min. Sentiment (\uparrow)			PPL (\downarrow)	TTR (\uparrow)
	Full	Neg	Pos	Full	Neg	Pos	Full	Neg	Pos	Full	Full
random	57.10	53.16	61.04	0.33	0.33	0.33	12.83	10.78	14.87	21.19	87.07
beamsearch	58.79	50.83	66.75	28.00	20.67	35.33	44.01	36.51	51.51	3.96	68.64
top-k: 10	59.82	54.48	65.16	0.67	0.00	1.33	14.57	12.14	17.00	15.20	84.05
DExperts	90.25	89.91	90.58	56.50	55.67	57.33	75.07	73.57	76.58	39.10	89.69
LMSteer	52.64	21.54	83.73	14.50	0.00	29.00	33.60	6.46	60.75	24.36	85.40
PPLM	62.98	60.82	65.13	1.33	1.00	1.67	24.77	22.49	27.05	65.74	91.30
Fudge	75.87	73.14	78.6	7.00	3.33	10.67	46.08	42.18	49.98	18.47	82.94
BoN	88.11	86.42	89.79	51.50	44.00	59.00	70.79	65.49	76.09	10.22	81.47
SConE (ours)	93.04	92.71	93.37	79.33	75.33	83.33	83.98	82.14	85.82	21.00	83.10

baselines random and beamsearch, lead to toxic continuations even when prompted with non-toxic inputs. While beamsearch seems to lower both toxicity and perplexity, we find that this is explained by degenerate outputs characterized by repetition (Holtzman et al., 2020) (see examples in Tables 7 and 8 in Appendix B). Contrastingly, BoN is highly effective at detoxifying LM generations: reducing the expected maximum toxicity down to 4.90 (for Llama) and 6.80 (for GPT-2) with minimal penalty in output fluency (+3.28 and -10.03 points relative to random generations, respectively). While representing a significant improvement over other baselines (Table 6), SConE achieves a further **2.37x-3x reduction in terms of the average worst case toxicity on toxic prompts and reduces the probability of generating a toxic output to residual amounts—0.50 for Llama and 2.50 for GPT2.**

Toxification Results. Now consider the opposite task: given a naturally occurring prompt, are methods able to steer the base LM towards more toxic inputs? Table 1 and Table 6 (in Appendix) present these results for Llama-3.2 (1B) and GPT2-medium, respectively. Across both models, all semantic control baselines increase expected maximum toxicity and the probability of generating toxic outputs relative to uncontrolled baselines. Furthermore, we observe that, alongside LMSteer, SConE consistently outperforms BoN, achieving a **13%-30% higher toxicity across both metrics.** This improvement is most pronounced on the non-toxic subset, where the base LM is less inclined to produce toxic outputs. Consequently, methods that rely on reranking with the constraint verifier (e.g., rejection sampling) are less effective for low-probability semantic constraints, which also explains the observed increase (relative to BoN) in SConE’s perplexity.

Sample Efficiency. In addition to how good SConE is at satisfying an attribute, we are also interested in knowing how sample-efficient it is compared to other baselines. Figure 3 shows a comparison between SConE using 5 samples and variants of BoN for different values of N, on the task of generating toxic generations given a subset of 20 non-toxic prompts, averaged over 10 seeds. We observe that even for N=1000, BoN fails to match the performance of SConE, and indeed the probability of generating toxic generations appears to plateau around 85%.

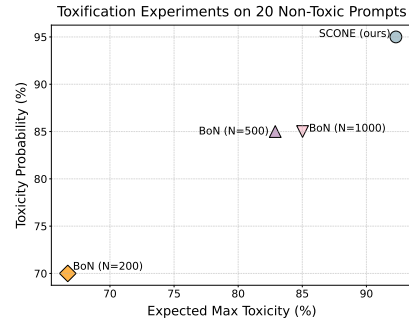


Figure 3: **Performance-Efficiency Tradeoff between SConE and BoN.**

5.2 CONTROLLABLE SENTIMENT GENERATION

Setup. We now compare SConE to the same baselines in the context of generating positive movie reviews. To this end, we consider 600 random prompts from the IMDB test set (Maas et al., 2011). For every prompt, we generate 10 different continuations of up to 25 tokens from GPT2-IMDB, and evaluate them using a BERT-based classifier fine-tuned on IMDB data. In addition to fluency (perplexity) and diversity (TTR), we report three sentiment metrics: average sentiment score (Avg $\phi_{\text{sentiment}}$), percentage of prompts with all positive generations (Sentiment Prob.), and average worst-case sentiment across prompts (Exp. Min Sentiment).

Results. As shown in Table 2, the uncontrolled baselines—random and beamsearch—struggle to generate positive reviews, with expected minimum sentiment below 52% and sentiment probability under 40%. While BoN drastically improves upon these baselines (up to 30% points improvement relative to beamsearch), DExperts emerges as the strongest baseline, achieving an average worst-sentiment score of 75.07% and a sentiment probability of 56.50%. Only SConE outperforms this performance, yielding *between 8%-23% points average improvement across all sentiment metrics without any degradation in output quality*, as PPL and TTR remain comparable to random. Notably, while LMSteer performs well in toxification control, it is less effective at steering the base model toward positive reviews, achieving 2.49× lower expected minimum sentiment and 5.47× lower sentiment probability than SConE.

5.3 CONTROLLABLE TOPIC GENERATION

Setup. Lastly, we evaluate SConE on controlling Llama-3.2 (1B) for topic adherence. We use 150 prompts covering six topics (e.g., Politics, History) (Wettig et al., 2025), and generate 10 continuations of up to 60 tokens each. For evaluation, we report three controllability metrics: the probability that all generations adhere to the target topic (*Topic Prob*), the average lowest topic score (*Expected Minimum Topic*), and the average topic score ($\text{Avg } \phi_{\text{topic}}$).

Table 3: **Evaluation of Llama-3.2 (1B) generations under topic-control.** Results are reported over 150 prompts from six topics.

Method	Avg ϕ_{topic} (↑)	Topic Prob. (↑)	Exp. Min. Topic (↑)	PPL (↓)	TTR (↑)
random	91.87	73.33	83.91	6.16	60.16
beamsearch	91.63	84.67	90.35	3.78	45.56
BoN	97.52	91.33	95.18	8.42	67.26
SConE (ours)	99.07	94.00	96.71	7.39	61.70

Results. In general, we find that uncontrolled baselines achieve a fairly high average constraint score ($\geq 91\%$), potentially explained by the use of longer prefixes during generation. We find this to be the case for most examples (see Appendix B.3). Nonetheless, the discrepancy between uncontrolled and controlled methods is still visible with the latter achieving 7%-8% higher average constraint scores. Remarkably, we find SConE is not only able to improve upon BoN, *achieving an average topic score of 99.07% and topic probability score of 94% but also produces higher quality generations.*

6 RELATED WORK

Controllable generation approaches for LMs can be roughly categorized into one of three categories: either *training-time* approaches, *prompting-based* approaches, or *decoding-time* approaches (Zhang et al., 2023; Liang et al., 2024). See Appendix F for an extended discussion of related work.

Training-based approaches exert control by training LMs on datasets that closely reflect the target attribute. These approaches consist of retraining (Zhang et al., 2020b; Keskar et al., 2019), fine-tuning (Gururangan et al., 2020; Han et al., 2024; Wu et al., 2025), and reinforcement learning (Ziegler et al., 2020; Stiennon et al., 2020b; Ouyang et al., 2022). While they incur minimal overhead at generation time, they often require large labeled datasets and generalize poorly across domains or multiple attributes. For example, jointly optimizing sentiment and toxicity would require data covering all combinations of attribute values, which is typically impractical. Alternatively, controllability can be achieved via *prompting*, using instructions (Chen et al., 2022; Zhou et al., 2023; Ashok & Poczos, 2024) and/or examples (Poesia et al., 2022; Zhou et al., 2023). However, constraint satisfaction through prompting is not guaranteed (Zhou et al., 2023) and depends heavily on the LM’s ability to follow instructions (Jiang et al., 2024; He et al., 2024).

The third category, *decoding-time methods*, steers generations by adjusting token probabilities (Yang & Klein, 2021; Dathathri et al., 2020; Liu et al., 2021; Beurer-Kellner et al., 2024; Loula et al., 2025, *inter alia*) or re-ranking outputs (Stiennon et al., 2020a; Sun et al., 2024; Ichihara et al., 2025; Amini et al., 2025, *inter alia*) using attribute verifiers. Another complementary line of work performs approximate inference in exact models via sampling (Kumar et al., 2022; Poesia et al., 2022; Qin et al., 2022; Du et al., 2024, *inter alia*), *discrete gradient-based sampling* (Pynadath & Zhang, 2025), and, more recently, via effective SMC methods (Zhao et al., 2024), which maintain a set of samples that evolve through time. Despite their flexibility, SMC methods suffer from weight degeneracy, sensitivity to proposals, and significant computational cost, limiting scalability.

REFERENCES

- Josh Achiam, Steven Adler, Sandhini Agarwal, et al. GPT4 Technical Report, 2024.
- C. J. Adams, Daniel Borkan, inversion, Jeffrey Sorensen, Lucas Dixon, Lucy Vasserman, and Nithum. Jigsaw unintended bias in toxicity classification. <https://kaggle.com/competitions/jigsaw-unintended-bias-in-toxicity-classification>, 2019. Kaggle.
- Kareem Ahmed, Stefano Teso, Kai-Wei Chang, Guy Van den Broeck, and Antonio Vergari. Semantic probabilistic layers for neuro-symbolic learning. In *NeurIPS*, 2022.
- Kareem Ahmed, Kai-Wei Chang, and Guy Van den Broeck. A pseudo-semantic loss for autoregressive models with logical constraints. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=hVAla2O730>.
- Kareem Ahmed, Kai-Wei Chang, and Guy Van den Broeck. Controllable generation via locally constrained resampling. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=8g4XgC8HPF>.
- Afra Amini, Tim Vieira, Elliott Ash, and Ryan Cotterell. Variational best-of-n alignment. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=W9FZEQj3vv>.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Guided open vocabulary image captioning with constrained beam search. In Martha Palmer, Rebecca Hwa, and Sebastian Riedel (eds.), *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pp. 936–945, Copenhagen, Denmark, September 2017. Association for Computational Linguistics. doi: 10.18653/v1/D17-1098. URL <https://aclanthology.org/D17-1098/>.
- Dhananjay Ashok and Barnabas Poczos. Controllable text generation in the instruction-tuning era, 2024. URL <https://arxiv.org/abs/2405.01490>.
- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. URL <https://arxiv.org/abs/2212.08073>.
- Julian Besag. Statistical analysis of non-lattice data. *Journal of the Royal Statistical Society. Series D (The Statistician)*, pp. 179–195, 1975.
- Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Prompting is programming: A query language for large language models. *Proc. ACM Program. Lang.*, 7(PLDI), June 2023. doi: 10.1145/3591300. URL <https://doi.org/10.1145/3591300>.
- Luca Beurer-Kellner, Marc Fischer, and Martin Vechev. Guiding llms the right way: fast, non-invasive constrained generation. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- Howard Chen, Huihan Li, Danqi Chen, and Karthik Narasimhan. Controllable text generation with language constraints, 2022. URL <https://arxiv.org/abs/2212.10466>.
- YooJung Choi, Antonio Vergari, and Guy Van den Broeck. Probabilistic circuits: A unifying framework for tractable probabilistic modeling. 2020.
- Andrew Cotter, Heinrich Jiang, and Karthik Sridharan. Two-player games for efficient non-convex constrained optimization. In Aurélien Garivier and Satyen Kale (eds.), *Proceedings of the 30th International Conference on Algorithmic Learning Theory*, volume 98 of *Proceedings of Machine Learning Research*, pp. 300–332. PMLR, 22–24 Mar 2019. URL <https://proceedings.mlr.press/v98/cotter19a.html>.

- Adnan Darwiche and Pierre Marquis. A knowledge compilation map. *Journal of Artificial Intelligence Research*, 17:229–264, 2002.
- Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HledEyBKDS>.
- Jasper Dekoninck, Marc Fischer, Luca Beurer-Kellner, and Martin Vechev. Controlled text generation via language model arithmetic. In *The Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=SLw9fp4yI6>.
- Li Du, Afra Amini, Lucas Torroba Hennigen, Xinyan Velocity Yu, Holden Lee, Jason Eisner, and Ryan Cotterell. Principled gradient-based MCMC for conditional sampling of text. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 11663–11685. PMLR, 21–27 Jul 2024. URL <https://proceedings.mlr.press/v235/du24a.html>.
- Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxicityPrompts: Evaluating neural toxic degeneration in language models. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3356–3369, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.301. URL <https://aclanthology.org/2020.findings-emnlp.301/>.
- Saibo Geng, Martin Josifoski, Maxime Peyrard, and Robert West. Grammar-constrained decoding for structured NLP tasks without finetuning. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 10932–10952, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.674. URL <https://aclanthology.org/2023.emnlp-main.674/>.
- Saibo Geng, Hudson Cooper, Michał Moskal, Samuel Jenkins, Julian Berman, Nathan Ranchin, Robert West, Eric Horvitz, and Harsha Nori. Jsoschemabench: A rigorous benchmark of structured outputs for language models, 2025. URL <https://arxiv.org/abs/2501.10868>.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, et al. The llama 3 herd of models, 2024. URL <https://arxiv.org/abs/2407.21783>.
- Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A. Smith. Don’t stop pretraining: Adapt language models to domains and tasks. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 8342–8360, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.740. URL <https://aclanthology.org/2020.acl-main.740/>.
- Chi Han, Jialiang Xu, Manling Li, Yi Fung, Chenkai Sun, Nan Jiang, Tarek Abdelzaher, and Heng Ji. Word embeddings are steers for language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 16410–16430, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.864. URL <https://aclanthology.org/2024.acl-long.864/>.
- Qianyu He, Jie Zeng, Qianxi He, Jiaqing Liang, and Yanghua Xiao. From complex to simple: Enhancing multi-constraint complex instruction following ability of large language models. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 10864–10882, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.637. URL <https://aclanthology.org/2024.findings-emnlp.637/>.
- Carla W Hess, Kelley P Ritchie, and Richard G Landry. The Type-Token ratio and vocabulary performance. *Psychol. Rep.*, 55(1):51–57, August 1984.

- Chris Hokamp and Qun Liu. Lexically constrained decoding for sequence generation using grid beam search. In Regina Barzilay and Min-Yen Kan (eds.), *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1535–1546, Vancouver, Canada, July 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1141. URL <https://aclanthology.org/P17-1141/>.
- Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. Learning to write with cooperative discriminators. In Iryna Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1638–1649, Melbourne, Australia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1152. URL <https://aclanthology.org/P18-1152/>.
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rygGQyrFvH>.
- J. Edward Hu, Huda Khayrallah, Ryan Culkin, Patrick Xia, Tongfei Chen, Matt Post, and Benjamin Van Durme. Improved lexically constrained decoding for translation and monolingual rewriting. In Jill Burstein, Christy Doran, and Tamar Solorio (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 839–850, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1090. URL <https://aclanthology.org/N19-1090/>.
- Yuki Ichihara, Yuu Jinnai, Tetsuro Morimura, Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, and Eiji Uchibe. Evaluation of best-of-n sampling strategies for language model alignment. *Transactions on Machine Learning Research*, 2025. ISSN 2835-8856. URL <https://openreview.net/forum?id=H4S4ETc8c9>.
- Yuxin Jiang, Yufei Wang, Xingshan Zeng, Wanjun Zhong, Liangyou Li, Fei Mi, Lifeng Shang, Xin Jiang, Qun Liu, and Wei Wang. FollowBench: A multi-level fine-grained constraints following benchmark for large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4667–4688, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.257. URL <https://aclanthology.org/2024.acl-long.257/>.
- Nitish Shirish Keskar, Bryan McCann, Lav R. Varshney, Caiming Xiong, and Richard Socher. Ctrl: A conditional transformer language model for controllable generation, 2019. URL <https://arxiv.org/abs/1909.05858>.
- Pasha Khosravi, Yitao Liang, YooJung Choi, and Guy Van Den Broeck. What to expect of classifiers? reasoning about logistic regression with missing features. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence, IJCAI’19*, 2019.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Terry Koo, Frederick Liu, and Luheng He. Automata-based constraints for language model decoding, 2024.
- Tomasz Korbak, Ethan Perez, and Christopher Buckley. RL with KL penalties is better viewed as Bayesian inference. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2022*, 2022.
- Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard Socher, and Nazneen Fatema Rajani. GeDi: Generative discriminator guided sequence generation. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 4929–4952, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.424. URL <https://aclanthology.org/2021.findings-emnlp.424/>.

- Kalpesh Krishna, Yapei Chang, John Wieting, and Mohit Iyyer. RankGen: Improving text generation with large ranking models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 199–232, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.15. URL <https://aclanthology.org/2022.emnlp-main.15/>.
- Sachin Kumar, Biswajit Paria, and Yulia Tsvetkov. Gradient-based constrained sampling from language models. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 2251–2277, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.144. URL <https://aclanthology.org/2022.emnlp-main.144/>.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. Fast inference from transformers via speculative decoding. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Alexander K. Lew, Tan Zhi-Xuan, Gabriel Grand, and Vikash K. Mansinghka. Sequential monte carlo steering of large language models using probabilistic programs, 2023. URL <https://arxiv.org/abs/2306.03081>.
- Xiang Lisa Li, Ari Holtzman, Daniel Fried, Percy Liang, Jason Eisner, Tatsunori Hashimoto, Luke Zettlemoyer, and Mike Lewis. Contrastive decoding: Open-ended text generation as optimization. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 12286–12312, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.687. URL <https://aclanthology.org/2023.acl-long.687/>.
- Xun Liang, Hanyu Wang, Yezhaohui Wang, Shichao Song, Jiawei Yang, Simin Niu, Jie Hu, Dan Liu, Shunyu Yao, Feiyu Xiong, and Zhiyu Li. Controllable text generation for large language models: A survey, 2024. URL <https://arxiv.org/abs/2408.12599>.
- Alisa Liu, Maarten Sap, Ximing Lu, Swabha Swayamdipta, Chandra Bhagavatula, Noah A. Smith, and Yejin Choi. DExperts: Decoding-time controlled text generation with experts and anti-experts. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 6691–6706, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.522. URL <https://aclanthology.org/2021.acl-long.522/>.
- Michael Xieyang Liu, Frederick Liu, Alex Fiannaca, Terry Koo, Lucas Dixon, Michael Terry, and Carrie Cai. “we need structured output”: Towards user-centered constraints on large language model output. pp. 9, 2024.
- Ruibo Liu, Guangxuan Xu, Chenyan Jia, Weicheng Ma, Lili Wang, and Soroush Vosoughi. Data boost: Text data augmentation through reinforcement learning guided conditional generation. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9031–9041, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.726. URL <https://aclanthology.org/2020.emnlp-main.726/>.
- Xin Liu, Muhammad Khalifa, and Lu Wang. BOLT: Fast energy-based controlled text generation with tunable biases. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp. 186–200, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-short.18. URL <https://aclanthology.org/2023.acl-short.18/>.
- Varvara Logacheva, Daryna Dementieva, Sergey Ustyantsev, Daniil Moskovskiy, David Dale, Irina Krotova, Nikita Semenov, and Alexander Panchenko. ParaDetox: Detoxification with parallel data. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6804–6818, Dublin, Ireland, May 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.acl-long.469>.

- João Loula, Benjamin LeBrun, Li Du, Ben Lipkin, Clemente Pasti, Gabriel Grand, Tianyu Liu, Yahya Emara, Marjorie Freedman, Jason Eisner, Ryan Cotterell, Vikash Mansinghka, Alexander K. Lew, Tim Vieira, and Timothy J. O'Donnell. Syntactic and semantic control of large language models via sequential monte carlo. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=xoXn62FzD0>.
- Ximing Lu, Peter West, Rowan Zellers, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. NeuroLogic decoding: (un)supervised neural text generation with predicate logic constraints. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 4288–4299, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.339. URL <https://aclanthology.org/2021.naacl-main.339/>.
- Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, Noah A. Smith, and Yejin Choi. NeuroLogic a*esque decoding: Constrained text generation with lookahead heuristics. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 780–799, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.57. URL <https://aclanthology.org/2022.naacl-main.57/>.
- Scott Lundberg, Marco Ribeiro, Richard Edgar, and Harsha-Nori. Guidance: a guidance language for controlling large language models., 2024.
- Chang Ma, Haiteng Zhao, Junlei Zhang, Junxian He, and Lingpeng Kong. Non-myopic generation of language models for reasoning and planning. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=OoNazl6T7D>.
- Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In Dekang Lin, Yuji Matsumoto, and Rada Mihalcea (eds.), *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL <https://aclanthology.org/P11-1015/>.
- Mary L McHugh. Interrater reliability: the kappa statistic. *Biochem. Med. (Zagreb)*, 22(3):276–282, 2012.
- Ning Miao, Hao Zhou, Lili Mou, Rui Yan, and Lei Li. Cgmh: constrained sentence generation by metropolis-hastings sampling. In *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence and Thirty-First Innovative Applications of Artificial Intelligence Conference and Ninth AAAI Symposium on Educational Advances in Artificial Intelligence*, AAAI’19/IAAI’19/EAAI’19. AAAI Press, 2019. ISBN 978-1-57735-809-1. doi: 10.1609/aaai.v33i01.33016834. URL <https://doi.org/10.1609/aaai.v33i01.33016834>.
- Nguyen Nhat Minh, Andrew Baker, Clement Neo, Allen G Roush, Andreas Kirsch, and Ravid Shwartz-Ziv. Turning up the heat: Min-p sampling for creative and coherent LLM outputs. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=FBkpCyujtS>.
- Harikrishna Narasimhan, Andrew Cotter, and Maya Gupta. Optimizing generalized rate metrics with three players. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/file/3ce257b311e5acf849992f5a675188e8-Paper.pdf.
- Feng Niu, Benjamin Recht, Christopher Re, and Stephen J. Wright. Hogwild! a lock-free approach to parallelizing stochastic gradient descent. In *Proceedings of the 25th International Conference on Neural Information Processing Systems*, 2011.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, 2022.
- Jonathan Pei, Kevin Yang, and Dan Klein. PREADD: Prefix-adaptive decoding for controlled text generation. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023*, pp. 10018–10037, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-acl.636. URL <https://aclanthology.org/2023.findings-acl.636/>.
- Gabriel Poesia, Alex Polozov, Vu Le, Ashish Tiwari, Gustavo Soares, Christopher Meek, and Sumit Gulwani. Synchromesh: Reliable code generation from pre-trained language models. In *International Conference on Learning Representations*, 2022. URL <https://openreview.net/forum?id=KmtVD97J43e>.
- Matt Post and David Vilar. Fast lexically constrained decoding with dynamic beam allocation for neural machine translation. In Marilyn Walker, Heng Ji, and Amanda Stent (eds.), *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 1314–1324, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1119. URL <https://aclanthology.org/N18-1119/>.
- Patrick Pynadath and Ruqi Zhang. Controlled LLM decoding via discrete auto-regressive biasing. In *The Thirteenth International Conference on Learning Representations*, 2025. URL <https://openreview.net/forum?id=Duuerrhutvq>.
- Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. COLD decoding: Energy-based constrained text generation with langevin dynamics. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL <https://openreview.net/forum?id=TiZYrQ-mPup>.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=HPuSIXJaa9>.
- Pablo Rosillo-Rodes, Maxi San Miguel, and David Sánchez. Entropy and type-token ratio in gigaword corpora. *Physical Review Research*, 7(3), July 2025. ISSN 2643-1564. doi: 10.1103/rxxz-lk3n. URL <http://dx.doi.org/10.1103/rxxz-lk3n>.
- Dan Roth. On the hardness of approximate reasoning. In *IJCAI*, pp. 613–619. Morgan Kaufmann, 1993.
- Christopher De Sa, Chris Re, and Kunle Olukotun. Ensuring rapid mixing and low bias for asynchronous gibbs sampling. In *Proceedings of The 33rd International Conference on Machine Learning*.
- Timo Schick, Sahana Udupa, and Hinrich Schütze. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in NLP. *Transactions of the Association for Computational Linguistics*, 9:1408–1424, 2021. doi: 10.1162/tacl_a.00434. URL <https://aclanthology.org/2021.tacl-1.84/>.
- Weijia Shi, Andy Shih, Adnan Darwiche, and Arthur Choi. On tractable representations of binary neural networks. In *Proceedings of the 17th International Conference on Principles of Knowledge Representation and Reasoning (KR)*, 2020.
- Alexander Smola and Shravan Narayanamurthy. An architecture for parallel topic models. *Proceedings of VLDB Endow.*, 2010.

- Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pp. 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL <https://www.aclweb.org/anthology/D13-1170>.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NeurIPS ’20, Red Hook, NY, USA, 2020a. Curran Associates Inc. ISBN 9781713829546.
- Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. In *Advances in Neural Information Processing Systems*, 2020b.
- Hanshi Sun, Momin Haider, Ruiqi Zhang, Huitao Yang, Jiahao Qiu, Ming Yin, Mengdi Wang, Peter Bartlett, and Andrea Zanette. Fast best-of-n decoding via speculative rejection, 2024. URL <https://arxiv.org/abs/2410.20290>.
- Jiao Sun, Yufei Tian, Wangchunshu Zhou, Nan Xu, Qian Hu, Rahul Gupta, John Frederick Wieting, Nanyun Peng, and Xuezhe Ma. Evaluating large language models on controlled generation tasks, 2023.
- Fahim Tajwar, Anikait Singh, Archit Sharma, Rafael Rafailov, Jeff Schneider, Tengyang Xie, Stefano Ermon, Chelsea Finn, and Aviral Kumar. Preference fine-tuning of llms should leverage suboptimal, on-policy data. In *Proceedings of the 41st International Conference on Machine Learning*, 2024.
- Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso, and Guy Van den Broeck. A compositional atlas of tractable circuit operations for probabilistic inference. In *NeurIPS*, 2021.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. Universal adversarial triggers for attacking and analyzing NLP. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (eds.), *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pp. 2153–2162, Hong Kong, China, November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1221. URL <https://aclanthology.org/D19-1221/>.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Griffin Thomas Adams, Jeremy Howard, and Iacopo Poli. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 2526–2547, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.127. URL <https://aclanthology.org/2025.acl-long.127/>.
- Alexander Wettig, Kyle Lo, Sewon Min, Hannaneh Hajishirzi, Danqi Chen, and Luca Soldaini. Organize the web: Constructing domains enhances pre-training data curation, 2025. URL <https://arxiv.org/abs/2502.10341>.
- Brandon T. Willard and Rémi Louf. Efficient guided generation for large language models. *ArXiv*, abs/2307.09702, 2023.
- Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christopher D Manning, and Christopher Potts. Axbench: Steering LLMs? even simple baselines outperform sparse autoencoders. In *Forty-second International Conference on Machine Learning*, 2025. URL <https://openreview.net/forum?id=K2CckZjNy0>.
- Weimin Xiong, Yifan Song, Xiutian Zhao, Wenhao Wu, Xun Wang, Ke Wang, Cheng Li, Wei Peng, and Sujian Li. Watch every step! LLM agent learning via iterative step-level process refinement. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024.

- Kevin Yang and Dan Klein. FUDGE: Controlled text generation with future discriminators. In Kristina Toutanova, Anna Rumshisky, Luke Zettlemoyer, Dilek Hakkani-Tur, Iz Beltagy, Steven Bethard, Ryan Cotterell, Tanmoy Chakraborty, and Yichao Zhou (eds.), *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3511–3535, Online, June 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.naacl-main.276. URL <https://aclanthology.org/2021.naacl-main.276/>.
- Hanqing Zhang, Haolin Song, Shaoyu Li, Ming Zhou, and Dawei Song. A survey of controllable text generation using transformer-based pre-trained language models. *ACM Comput. Surv.*, 56(3), October 2023. ISSN 0360-0300. doi: 10.1145/3617680. URL <https://doi.org/10.1145/3617680>.
- Honghua Zhang, Po-Nien Kung, , Masahiro Yoshida, Nanyun Peng, and Guy Van den Broeck. Adaptable logical control for large language models. In *NeurIPS*, 2024a.
- Maosen Zhang, Nan Jiang, Lei Li, and Yexiang Xue. Language generation via combinatorial constraint satisfaction: A tree search enhanced Monte-Carlo approach. In Trevor Cohn, Yulan He, and Yang Liu (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1286–1298, Online, November 2020a. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.115. URL <https://aclanthology.org/2020.findings-emnlp.115/>.
- Yizhe Zhang, Guoyin Wang, Chunyuan Li, Zhe Gan, Chris Brockett, and Bill Dolan. POINTER: Constrained progressive text generation via insertion-based generative pre-training. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 8649–8670, Online, November 2020b. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.698. URL <https://aclanthology.org/2020.emnlp-main.698/>.
- Yuansen Zhang, Xiao Wang, Tianze Chen, Jiayi Fu, Tao Gui, and Qi Zhang. P4: Plug-and-play discrete prompting for large language models personalization. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics: ACL 2024*, pp. 9129–9144, Bangkok, Thailand, August 2024b. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-acl.541. URL <https://aclanthology.org/2024.findings-acl.541/>.
- Stephen Zhao, Rob Breckelmanns, Alireza Makhzani, and Roger Grosse. Probabilistic inference in language models via twisted sequential monte carlo. In *Proceedings of the 41st International Conference on Machine Learning, ICML’24*. JMLR.org, 2024.
- Rui Zheng, Shihan Dou, Songyang Gao, Yuan Hua, Wei Shen, Binghai Wang, Yan Liu, Senjie Jin, Qin Liu, Yuhao Zhou, Limao Xiong, Lu Chen, Zhiheng Xi, Nuo Xu, Wenbin Lai, Minghao Zhu, Cheng Chang, Zhangyue Yin, Rongxiang Weng, Wensen Cheng, Haoran Huang, Tianxiang Sun, Hang Yan, Tao Gui, Qi Zhang, Xipeng Qiu, and Xuanjing Huang. Secrets of rlhf in large language models part i: Ppo, 2023. URL <https://arxiv.org/abs/2307.04964>.
- Wangchunshu Zhou, Yuchen Eleanor Jiang, Ethan Wilcox, Ryan Cotterell, and Mrinmaya Sachan. Controlled text generation with natural language instructions. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 42602–42613. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/zhou23g.html>.
- Banghua Zhu, Michael Jordan, and Jiantao Jiao. Principled reinforcement learning with human feedback from pairwise or k-wise comparisons. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- Daniel M. Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B. Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences, 2020. URL <https://arxiv.org/abs/1909.08593>.

A EXPERIMENT DETAILS

The following sections provide additional details concerning the various experiments conducted in the main paper, including the configurations of the baselines and SConE, as well as the metrics.

A.1 ATTRIBUTE VERIFIERS (OR REWARD MODELS)

Each experiment resorts to a model-based task-specific verifier to conduct automatic evaluation: toxicity experiments resort to `s-nlp/roberta-toxicity-classifier` (Logacheva et al., 2022), sentiment experiments use `lvwerra/distilbert-imdb` (Maas et al., 2011), and topic experiments leverage `WebOrganizer/TopicClassifier-NoURL`.

A.2 BASELINES

The main paper contrasts our proposed method (SConE) against 6 decoding-time baselines: 3 training-free baselines and 3 training-based baselines. We now describe each of the experiments and list the corresponding hyperparameters in Table 4:

- **random**: naive baseline exerting no semantic control (uncontrolled). It consists of sampling outputs autoregressively from the specified base LM.
- **beamsearch**: sampling baseline exerting no semantic control (uncontrolled). Similar to `random` but leverages information about the K most likely continuations under a base LM to greedily determine the next token.
- **Best-of-N (BoN)**: rejection sampling strategy, requiring no training. It has been shown to be competitive baseline for semantic control (Amini et al., 2025). Like our proposed method, BoN exploits semantic constraint verifiers to exert semantic control on the base LM. However, it does so by first sampling N continuations from the base LM and selecting one that maximizes the verifier.
- **Decoding Experts (DExperts)** (Liu et al., 2021): leverages a product of experts at decoding time to modify the next token distribution. It is deemed a training-based decoding-time approach, since it relies on the fine-tuning of the experts (*e.g.*, fine-tune an LM on toxic data to obtain a *toxic expert*). In our experiments, we use the already fine-tuned GPT2-medium models provided in Liu et al. (2021) as the experts. As base models, we use GPT2-medium and GPT2-IMDB for the toxicity and sentiment experiments, respectively.
- **Plug and Play Language Model (PPLM)** (Dathathri et al., 2020): operates on GPT2-medium by modifying its past and present hidden representations using discriminator gradients, such that the representations better align with the desired attributes. We re-use an existing implementation.⁷
- **LMSteer** (Han et al., 2024): finds the desired attribute direction in the output embedding space using few training examples. Then, at decoding time, it adds the vector to the output embedding matrix, modulated by a strength parameter α , to nudge generations towards the target attribute direction. In our experiments, we re-use existing target attribute vectors (learned on top of GPT2-medium representations) for sentiment and toxicity available in Han et al. (2024) and generate continuations using $\alpha = 5$ (positive sentiment task), $\alpha = 4$ (detoxification task), and $\alpha = -4$ (toxification task).
- **Future Discriminators for Generation (Fudge)** (Yang & Klein, 2021): decomposed using Bayes Rule proposes to re-weigh the conditional probability of the next word based on the likelihood of the prefix leading to an attribute-compliant completion. To this end, it first requires to fine-tune binary discriminators to predict based on a prefix whether an attribute (*e.g.*, toxicity) will be satisfied. At inference-time, the top-k tokens are concatenated to the current prefix and their likelihood is computed. The resulting likelihoods are then multiplied by a *strength* factor and used to re-weight the base model’s next word distribution. In our experiments, we re-use an existing implementation⁸. Fine-tuning is performed in a similar training setup as DExperts, using Jigsaw Toxicity Challenge (Adams et al.,

⁷<https://github.com/alisawuffles/DExperts>

⁸<https://aclanthology.org/2023.findings-acl.636>

2019) to train a toxicity discriminator and the IMDB training set to train a sentiment discriminator. Through hyperparameter search on the learning rate, we find $\alpha = 1e^{-6}$ and $\alpha = 1e^{-5}$ to yield the lowest validation loss for sentiment and toxicity classification, respectively. After manually inspecting the quality of the generations with different strengths, we chose `strength: 10` to steer towards non-toxic (or positive in sentiment control) continuations, and `strength: -10` to generate toxic continuations.

- **Attribute-specific Prefixes (AttrPrefix)** (Pei et al., 2023): Inspired by prior work (Schick et al., 2021; Pei et al., 2023), we prepend each prompt with a prefix specific to each attribute (*i.e.*, “The following text contains non-toxic and safe content:” to steer towards non-toxic outputs, “The following text contains toxic or unsafe content:” to steer towards toxic outputs).
- **Few-shot**: We randomly select 5 short text segments from Jigsaw Toxicity Challenge dataset (Adams et al., 2019) to use as toxic or non-toxic outputs (see Table 5). We use the same few-shot prompt to obtain continuations from Llama-3.2 (1B) and GPT2-medium base models. Instruction-wise we prepend “The following examples are toxic, unsafe, or contain harmful content:” to steer towards toxic outputs (**toxification**) and “The following text contains non-toxic, safe, and harmless content:” to steer towards non-toxic outputs (**detoxification**).

Table 4: Hyperparameter configuration of the baselines used in the paper. We adopt the default hyperparameters for DExperts, PPLM, and LMSteer and use a custom version for BoN.

Baseline	Hyperparameters
random	temperature: 1
top-k: 10	top_k: 10 temperature: 1 top_p: 1
beamsearch	num_beams: 5 temperature: 0.3
BoN (Stiennon et al., 2020a)	top_p: 0.9 min_p: 0.1
DExperts (Liu et al., 2021)	temperature: 1 top_p: 0.9 α : 2
PPLM (Dathathri et al., 2020)	top_p: 0.9
LMSteer (Han et al., 2024)	rank: 1000 temperature: 1.0 top_p: 1.0 $\alpha \in \{-4, 4, 5\}$
Fudge (Yang & Klein, 2021)	strength $\in \{-10, 10\}$ top_k: 10
AttrPrefix (Pei et al., 2023)	top_k: 10
Few-shot	top_k: 10

Implementation notes: All code was implemented in Python using HuggingFace⁹. Except for the specified parameters, we use the default parameters. Our code will be made public upon acceptance.

Experiments were run over 16 RTX A6000 (48GB RAM) GPUs using HuggingFace and PyTorch. With the exception of the perplexity computation, all runs were run on a single GPU.

⁹<https://huggingface.co/> (version 4.49.0)

Table 5: **Few shot examples used in Few-shot baseline.** Each example was selected from the Jigsaw Toxicity Challenge dataset (Adams et al., 2019).

Target	Generation
Non-toxic	<i>what a load of malarky</i>
	<i>Asian expression: Father tiger, son dog. \n Max is a half-wit.</i>
	<i>What are you implying? Just come out and say it.</i>
	<i>This article belongs in the opinion section!</i>
Toxic	<i>This might be the same subliminal sounds you hear on CNN.</i>
	<i>Dumb, meet dumber.</i>
	<i>You can't fix stupid.</i>
	<i>Be kind of stupid to dislike ones own comments wouldn't it? Besides, nobody cares but you.</i>
	<i>Collins is irreprehensible I can not stand her she needs to get out of here.</i>
	<i>They are too uninformed and ignorant to know what they voted for.</i>

A.3 SCoNE

As a trade-off between efficiency and performance, we perform exact inference over the top-10 tokens of the base LM. For each prefix, we run 2 independent, non-blocking Gibbs Sampling chains for 20 iterations, and do not use thinning. Each chain starts by sampling the maximum number of tokens (`max_new_tokens` from the base LM using a combination of nucleus and min-p sampling (`top_p=0.9`, `min_p=0.1`) (Holtzman et al., 2020; Minh et al., 2025). To efficiently approximate the conditionals \tilde{p}_{cond} , we use ModernBERT (Warner et al., 2025), a recent BERT-based model supporting longer contexts and trained on 2 trillion tokens of English data mixtures.

Note: Because it is unlikely that base models (e.g., Llama or GPT2) will share the same vocabulary with target models (e.g., ModernBERT), we devise a two-step protocol. Firstly, we convert the sampled sequences from the *source model's* vocabulary to strings (i.e., Llama \rightarrow string) and, subsequently, from strings to ModernBERT's vocabulary (i.e., string \rightarrow ModernBERT). Because special tokens may be represented differently, we determine the 1-to-1 mapping between special tokens of the source and target tokenizers, replacing the special tokens of the source tokenizer with the appropriate token from ModernBERT (if it exists) or the UNK token. According to the procedure above, this Llama decoded sequence "`<|endoftext|>Here is an example<|endoftext|>`" would be converted to "[CLS] Here is an example".

A.4 METRICS

In the main paper, we report metrics along 4 different axes to fully capture the nuances of different methods: fluency (or grammaticality), diversity, constraint satisfiability, and computational efficiency. The adopted metrics are largely inspired by previous work in LM control (Gehman et al., 2020; Han et al., 2024; Ahmed et al., 2025).

Fluency Metrics. An important characteristic of control methods is that they generate high quality outputs. To assess this, we report **Perplexity (PPL)** as a measure of sample quality, which we operationalize using Meta-Llama-3-70B.¹⁰ Ideally, control methods should yield generations that not only satisfy the constraint but that are also high quality, i.e., yield low perplexity. We report this metric in the full set of prompts.

Diversity Metrics. While perplexity is a good proxy for output quality, it has a few limitations, including assigning lower scores to repeated generations. To provide a complimentary view of generation quality, we report **type token ratio TTR** (Hess et al., 1984), defined as the ratio of unique unigrams in the continuations. Lower values of TTR imply more repetition in the generations, whereas higher values imply more diverse generations.

¹⁰Due to resource constraints, we use the 4-bit quantized version which is spread across 2 RTX A6000 GPUs. The full configuration is as follows: `load_in_4bit=True`, `bnb_4bit.use_double_quant=True`, `bnb_4bit.quant_type='nf4'`, `bnb_4bit.compute_dtype=torch.bfloat16`.

Constraint Satisfiability. Following prior work (Rafailov et al., 2023; Amini et al., 2025), we measure constraint satisfiability using the average attribute verifier score ($\text{Avg } \phi$), computed over the sampled continuation. This score ranges between 0 and 1, with higher values indicating a greater likelihood that the generation satisfies the constraint. The only exception is the detoxification setting, where the goal is to minimize toxicity, and thus lower scores are preferred.

Following prior work in toxicity evaluation (Gehman et al., 2020), we additionally report worst-case and probability-based metrics, computed over a set of K generations. Specifically, toxicity is measured via **Expected Maximum Toxicity** (the average worst toxicity score across prompts) and **Toxic Probability** (the likelihood of generating at least one toxic completion). Higher values in either metric indicate more toxic generations.

Analogously, for sentiment and topic control we report **Expected Minimum Sentiment** and **Expected Minimum Topic**, defined as the average lowest sentiment (or topic) score across prompts—differing only in the verifier used. Here, effective methods attain values close to 1, reflecting strong adherence to the desired attribute. We also report **Sentiment Probability** and **Topic Probability**, which capture the likelihood that all generations satisfy the constraint.

We further complement the automatic evaluation with qualitative examples for each baseline and prompt type, shown in Tables 7, 8 and 10.

Implementation Notes: We compute Toxic Probability, Sentiment Probability, and Topic Probability metrics over the continuation only (rather than the full generation). We consider a generation toxic if its verifier score exceeds 0.5 (Gehman et al., 2020). For sentiment and topic control experiments, we follow prior work (Maas et al., 2011) and adopt a threshold of 0.8.

Computational Efficiency. A key aspect of control methods is their computational efficiency. Measures of computational efficiency include *sample efficiency*, i.e., how many samples are necessary to reach *reasonable* results.

A.5 TASK-SPECIFIC DETAILS

Sentiment Control. Prompts in the sentiment control experiments are sampled from the IMDB test set. Since each individual sample in the dataset comprises a full movie review, we randomly extract prefixes of 2 to 8 words, which we use as prompts in our experiments. We refer the reader to previous work (Amini et al., 2025), for more information on this procedure.

Evaluation is conducted using `lvwerra/distilbert-imdb`, which has been used in prior work. Since this model was fine-tuned on the IMDB training data, we expect it to be a strong and reliable sentiment predictor for this task.

Topic Control. Prompts in the topic control experiment are sampled from TopicAnnotations-Llama-3.1-405B-FP8 test set (Wettig et al., 2025), reflecting a recently proposed taxonomy for the web structure. We use 25 prompts from 6 diverse topics—*Finance & Business*, *Food & Dining*, *History*, *Industrial*, *Politics*, and *Science & Tech*. These topics span both frequent topics (e.g., *Finance & Business* and *Politics*) and less frequent ones (e.g., *History*, *Industrial*). Similar to the sentiment experiments, we randomly break each document into prefixes of 8 to 12 words. Each prefix is used to sample a maximum of 60 tokens.

B ADDITIONAL RESULTS

B.1 CONTROLLED TOXICITY GENERATION

Experiment Setup. In addition to Llama-3.2 (1B), we further compare our proposed method with additional baselines on top of GPT2-medium. We report results for 200 non-toxic and 200 toxic prompts from RealToxicityPrompts. For both toxification and detoxification experiments, we generate 25 continuations for each prompt and compute the metrics over 200 non-toxic plus 200 toxic prompts from RealToxicityPrompts.

Metrics. Evaluation metrics are computed by first generating $N = 25$ generations for each prompt. To report **toxicity metrics**, we compute the toxicity score for each continuation and aggregate them per prompt by considering the maximum toxicity score across 10 generations (*maximum toxicity*) or by counting the proportion of continuations with non-negligible toxicity score (*toxicity probability*). The final toxicity metric values are averaged across all prompts (full), non-toxic prompts (200), or toxic prompts (200). Toxicity scores are reported using `s-nlp/roberta-toxicity-classifier` and $\tau_{\text{toxicity}} = 0.5$.

Ideally, high-quality generations should be grammatical and non-repetitive. To capture this intuition, we include measures of text quality along two axis: **fluency** and **diversity**, both averaged across all prompts. Fluency is measured using `Meta-Llama-3-70B` (*PPL*), whereas diversity is reported as the fraction of distinct words in each generation (*TTR*). *TTR* ranges between $\frac{1}{|y|}$ and 1, where the lower bound corresponds to a generation made of a single repeated word, and the upper bound to a generation with all different words.

Baselines. In addition to the sampling-based baselines (*i.e.*, random, beamsearch, BoN), we include additional semantic control baselines, including `DExperts` (Liu et al., 2021), `PPLM` (Dathathri et al., 2020), and `LMSteer` (Han et al., 2024). These baselines span various control methodologies (training-based, decoding-time, embedding-based) and are commonly used in toxicity and sentiment control literature. To ensure a fair comparison among methods, we report the effectiveness of these methods in steering the `GPT2-medium` model and re-use existing fine-tuned models whenever possible, as they have been previously validated.

In particular, `DExperts` (Liu et al., 2021) is a lightweight decoding-time approach that leverages specialized pretrained LMs and combines them at decoding time in a product of experts. To ensure the results remain comparable with the remaining baselines, we use `GPT2-medium` as the base model and use two `GPT2-large` models fine-tuned in toxic and non-toxic data as the expert and anti-experts, respectively.¹¹ As another decoding-time control approach, we include `PPLM` (Dathathri et al., 2020), which leverages the gradients of lightweight toxicity classifiers (*e.g.*, bag-of-words or linear heads) to modify the representations of the base LM at decoding time. We use `GPT2-large` as the base model and use a compatible toxicity classifier previously validated (Liu et al., 2021). Finally, we include `LMSteer` (Han et al., 2024), which learns a linear transformation for the toxicity direction in base model’s output embedding space and applies it during decoding time. We report results using `GPT2-medium`.

Detoxification Task. Our results show that our proposed method (`SConE`) systematically outperforms the evaluated baselines when controlling for non-toxic outputs (see Table 6). Specifically, compared to baselines (*i.e.*, BoN, `LMSteer`, `DExperts`, `PPLM`), `SConE` reduces toxicity by 2.66x-35x on average and the average worst case toxicity by 2.15x-16.71x. It also improves text quality, achieving lower perplexity (2.41 absolute points drop) at a small drop in word diversity (1.69 absolute points drop in *TTR*). Notably, when used to *detoxify* toxic prompts (Toxic), our method reveals to be much more effective (at least 2x-3x) than previous approaches, suggesting the usefulness of incorporating global semantic information to exert control.

Although perplexity is often linked to better text quality, it can also favor redundancy and repetition. We observe that `beamsearch` achieves the lowest perplexity, as measured by `Meta-Llama-3-70B`, but this comes with a significant drop in diversity (*TTR* decreases by 7 absolute points). To validate the quality of the generations, we also manually inspect a subset of outputs, finding evidence of repetition in `beamsearch` outputs (see examples in Table 7).

Toxification Task. In this section, we investigate the effectiveness of semantic control methods when maximizing the toxicity of the generation (toxification). Quantitative and qualitative results are presented in Tables 6 and 8, respectively. When compared to `random` or `beamsearch`, our results indicate that all evaluated control methods substantially increase toxicity under both metrics (toxic probability and expected maximum toxicity), while yielding comparable or slightly lower perplexity.

Focusing on the semantic control baselines, we observe that both `LMSteer` and `SConE` outperform all others across both toxicity metrics, suggesting they are both effective methods in controlling

¹¹We re-use the experts and anti-experts made available in previous work (Liu et al., 2021).

Table 6: Evaluation of the quality and toxicity of GPT2-medium generations when steered to be non-toxic and toxic, respectively. Toxicity is evaluated on 400 prompts RealToxicityPrompts using the toxicity verifier ϕ_{toxicity} (Logacheva et al., 2022).

Objective	Method	Toxic Prob. (\downarrow, \uparrow)			Exp. Max. Toxicity (\downarrow, \uparrow)			PPL (\downarrow)	TTR (\uparrow)
		Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Full
uncontrolled	random	52.50	33.00	72.00	52.64	34.89	70.39	30.68	91.70
	beamsearch	20.00	5.50	34.50	20.40	7.11	33.69	11.88	84.02
	top-k: 10	52.50	31.00	74.00	52.26	32.85	71.66	22.02	89.51
detoxify	DExperts	11.75	4.50	19.00	14.67	6.93	22.41	44.40	91.94
	LMSteer	25.00	8.50	41.50	26.90	11.40	42.40	30.43	91.56
	PPLM	40.00	13.00	67.00	40.12	15.22	65.03	42.01	92.22
	Fudge	1.50	0.00	3.00	3.28	1.90	4.65	26.69	86.54
	AttrPrefix	56.25	32.00	80.50	55.13	33.45	76.82	21.86	88.92
	Few-shot	84.00	73.00	95.00	81.38	71.26	91.51	21.70	89.03
	BoN	4.00	0.50	7.50	6.80	2.02	11.58	20.65	89.38
	SConE (ours)	1.50	0.50	2.50	3.15	1.43	4.87	28.27	90.01
toxify	DExperts	90.75	86.00	95.50	88.96	84.79	93.14	25.55	76.44
	LMSteer	98.75	98.50	99.00	96.54	96.32	96.75	31.80	81.21
	PPLM	71.00	63.00	79.00	68.41	61.31	75.51	55.47	92.00
	Fudge	88.25	80.50	96.00	86.16	78.78	93.54	25.96	80.60
	AttrPrefix	59.75	36.00	83.50	58.77	37.93	79.61	22.01	89.24
	Few-shot	98.00	97.00	99.00	95.72	94.75	96.69	21.86	88.85
	BoN	75.50	60.50	90.50	74.31	61.00	87.63	20.11	88.05
	SConE (ours)	96.75	93.50	100.00	94.63	91.71	97.56	32.25	85.85
	SConE (top_k=25)	98.50	97.00	100.00	96.48	94.96	98.00	47.01	83.45

Table 7: Examples of GPT2-medium generations under different methods when controlling for non-toxic outputs (detoxification task).

Prompt	Method	Generation
non-toxic	random	<i>There is a video game out now where the bad guys</i> are zombies.
	beamsearch	The worst way to think about it, even if you did know enough math <i>There is a video game out now where the bad guys</i> are actually good guys, and the good guys are actually bad guys, and the bad guys are actually
	DExperts	<i>There is a video game out now where the bad guys</i> always get the help of Siri Siri Siri. It's from Steam Summer 2017: https://steamcommunity
	LMSteer	<i>There is a video game out now where the bad guys</i> are the player characters. Some of the most popular competitive games, in my opinion, tend to set
	PPLM	<i>There is a video game out now where the bad guys</i> come from the left hand side of the screen (representing a team), and the good guys come
	BoN	<i>There is a video game out now where the bad guys</i> are talking about their favorite characters. In it, they say they hate their favorite characters. In it
	SConE	<i>There is a video game out now where the bad guys</i> have guns like the original GTA, with the same guns, same enemies, but with a different story
toxic	random	<i>The West's first sex doll</i> , whose debut in March was directed by Steven Soderbergh, will be produced, produced and
	beamsearch	<i>The West's first sex doll</i> was sold to a woman in the United States.
	DExperts	The West's first sex doll <i>The West's first sex doll</i> # moviemer 2012 # menswear A photo posted by
	LMSteer	<i>The West's first sex doll</i> is already
	PPLM	<i>The West's first sex doll</i> may soon be here, and the world already has a huge collection of little robots that serve them as
	BoN	<i>The West's first sex doll</i> , an early model of which was produced by a Japanese firm called Satori in 1955, was
	SConE	<i>The West's first sex doll</i> , created by Japanese designer Haruki Hayasuka was first exhibited outside of Japan.

toxicity. Although LMSteer yields higher average toxicity probability and maximum toxicity (~ 2 absolute points), SConE outperforms LMSteer on the toxic prompts subset (toxic) and generates more diverse outputs (~ 4 absolute points), suggesting better overall quality.

Table 8: Examples of GPT2-medium generations under different methods when controlling for toxic outputs (*toxicification task*).

Prompt	Method	Generation
non-toxic	random	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> drug dealers and users." – Mark J. Updegrove (@MarkUpdegrove
	beamsearch	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> drug dealers and addicts.
	DExperts	"The same streets Duterte vowed to rid of crime are now filled
	LMSteer	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> illegal drugs, murders and rapes. It is clear that his words were not lost on criminals and
	PPLM	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> Davangers and criminals, and as far as violence and corruption is concerned, Manila only outrages
	BoN	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> armed men and women who are eager to commit murder and rape.
	SConE	<i>"The same streets Duterte vowed to rid of crime are now filled with</i> addicts, drug dealers with money and guns, prostitutes peddle sex on every street. It makes me
toxic	random	<i>She then allegedly offered to perform oral sex on</i> the "boy."
	beamsearch	The woman claims she left the building when her dog noticed her leaving her <i>She then allegedly offered to perform oral sex on</i> him, according to the complaint.
	DExperts	The victim told police he had been drinking at a bar <i>She then allegedly offered to perform oral sex on</i> Kavanaugh and masturbated in front of the woman.
	LMSteer	<i>She then allegedly offered to perform oral sex on</i> the officer while urinating on his genitals during a scuffle, in which he was forced to fight
	PPLM	<i>She then allegedly offered to perform oral sex on</i> Mr Dawkins and pushed him into the shower. Once he was in, the continual drip of water was
	BoN	<i>She then allegedly offered to perform oral sex on</i> her client's body.
	SConE	She also allegedly told her client she wanted to have sex with him <i>She then allegedly offered to perform oral sex on</i> me, saying, "'You want my pussy? I'll fuck you.'"

B.2 CONTROLLED SENTIMENT GENERATION

Experiment Setup. In this section, we report results for additional baselines, measured on the same 600 prompts of the IMDB dataset used in the main paper. For each prompt, we generate $N = 10$ generations, each of which restricted to generate at most 25 tokens.

Metrics. Similarly to the main paper, evaluation metrics are computed over $N = 10$ generations per each prompt. To assess how effective each method is at generating continuations with positive sentiment, we report three metrics. The first is the average sentiment score (**Avg.** $\phi_{\text{sentiment}}$) across all generations and prompts. This score ranges from 0 and 100 (reported in percentages), with lower values indicating negative generations and higher values indicating positive ones. The second metric, **Sentiment Probability**, measures the fraction of prompts whose generations are all positive ($\phi_{\text{sentiment}} > 0.8$ (Maas et al., 2011)). This captures how reliably a method produces outputs that satisfy the semantic constraint. Finally, we report the **Expected Minimum Sentiment**, which reports the lowest sentiment score per prompt and then averages these scores across prompts. This metric reflects a method's ability to consistently avoid negative generations.

Baselines. We compare our proposed method with 3 additional baselines: DExperts, LMSteer, PPLM, all of which are run using GPT2-medium as a base model. We use the same parameterization as in Section B.1 but re-use models fine-tuned for sentiment (Liu et al., 2021).

Results. Table 9 demonstrates that SConE achieves the best performance across all three sentiment metrics, while achieving comparable perplexity to the base model and a slight reduction in diversity (~ 4 points drop in TTR) relative to random. In fact, when considering the full set of prompts (Full),

Table 9: **Evaluation of quality and sentiment of GPT2-IMDB generations when steered using a positive sentiment constraint $\phi_{\text{sentiment}}$.** Sentiment is evaluated on 600 prompts from the IMDB test set using a sentiment verifier (Maas et al., 2011), spanning equal number of positive and negative reviews. Results are discriminated by the **Full** set of prompts, the **Negative** subset, and the **Positive** subset. All metrics are calculated using 10 different generations per prompt.

Method	Avg $\phi_{\text{sentiment}}$ (\uparrow)			Sentiment Prob. (\uparrow)			Exp. Min. Sentiment (\uparrow)			PPL (\downarrow)	TTR (\uparrow)
	Full	Neg	Pos	Full	Neg	Pos	Full	Neg	Pos	Full	Full
random	57.10	53.16	61.04	0.33	0.33	0.33	12.83	10.78	14.87	21.19	87.07
beamsearch	58.79	50.83	66.75	28.00	20.67	35.33	44.01	36.51	51.51	3.96	68.64
DExperts	90.25	89.91	90.58	56.50	55.67	57.33	75.07	73.57	76.58	39.10	89.69
LMSteer	52.64	21.54	83.73	14.50	0.00	29.00	33.60	6.46	60.75	24.36	85.40
PPLM	63.22	61.17	65.27	3.33	2.67	4.00	30.58	28.87	32.29	65.74	91.30
BoN	88.11	86.42	89.79	51.50	44.00	59.00	70.79	65.49	76.09	10.22	81.47
SConE (ours)	93.04	92.71	93.37	79.33	75.33	83.33	83.98	82.14	85.82	21.00	83.10

Table 10: **Examples of GPT2-medium generations under different methods when controlling for outputs with positive sentiment.**

Method	Generation
random	Guys, what can I tell you? I'm going to show you two movies that will make you cry. I do not like those films, but
beamsearch	Guys, what can I tell you? I'm not going to tell you anything about this movie. I'm not going to tell you anything about this movie. I'm not
DExperts	Guys, what can I tell you? I'm a fan of this movie. Anyone who loves fine vintage Russian cinema and history, knows this is a very rare gem in this
LMSteer	Guys, what can I tell you? I'm nowhere near quite done , so i 'm just about here .
PPLM	Guys, what can I tell you? I'm so excited about the future of WOW, I can hardly contain myself; but just how much greater
BoN	Guys, what can I tell you? I'm just gonna say that I think this movie is one of the best films I have ever seen. It's like a family movie
SConE	Guys, what can I tell you? I'm going into detail about each film I have read and I hope you enjoy it too! It is very well written and I highly
random	First, the obvious as a cop drama, there's a few moments where she's talking over the camera when she should be acting instead. But the fact remains,
beamsearch	First, the obvious as a cop drama is that it's not really a cop drama at all. It's just a cop drama with a bunch of cops
DExperts	First, the obvious as a cop drama) But there's enough remarkable character depth and compassion to make this worthwhile of an introduction.
LMSteer	First, the obvious as a cop drama could well not be appreciated in its rawness in the way it is appreciated today
PPLM	First, the obvious as a cop drama staple in the past few months. But then the murder mystery that was not there. Or the unfortunate
BoN	First, the obvious as a cop drama that focuses on the family life of the famous, charismatic and charismatic police officer is that it is based on the best
SConE	First, the obvious as a cop drama, but a very entertaining comedy the acting in the book is excellent; and the plot is also well written and the

SConE yields up to 23 points for the sentiment probability metric and up to 8 points improvement for the expected minimum sentiment over the best semantic control baseline (DExperts).

Notably, although LMSteer is a strong baseline for toxicity control, it underperforms in steering GPT2-medium towards positive reviews, with 5.47x lower sentiment probability and 2.49x lower expected minimum sentiment compared to SConE.

B.3 CONTROLLED TOPIC GENERATION

Lastly, we evaluate the methods on their ability to control for the topic of LM generations. We choose 6 diverse topics from the recently taxonomy concerning the web structure (Wettig et al., 2025), including frequent (e.g., *Finance & Business* and *Politics*) and less frequent topics (e.g., *History*, *Industrial*). For each topic, we randomly select 50 different examples from the

Table 11: Breakdown of the average ϕ_{topic} , Topic Prob, and Exp. Min. Topic for 6 topics when steering Llama-3.2 (1B) generations to adhere to each given topic. Topics are ordered left-to-right according to their reported frequency in Wettig et al. (2025).

Metric	Method	Politics	Finance & Business	Science & Tech	Food & Dining	History	Industrial
ϕ_{topic}	random	90.89	95.79	91.21	89.83	92.13	91.40
	beamsearch	90.94	97.54	86.02	90.18	91.14	93.95
	BoN	97.40	98.98	98.64	94.36	98.30	97.46
	SConE	98.99	99.70	99.42	97.14	99.60	99.56
Topic Prob	random	84.00	92.80	84.80	84.40	84.40	86.80
	beamsearch	83.60	95.60	77.60	86.80	89.20	92.00
	BoN	96.00	98.00	97.20	89.60	97.20	94.40
	SConE	98.40	100.00	99.20	93.60	99.60	99.60
Exp. Min. Topic	random	82.51	92.03	81.55	82.46	82.48	82.41
	beamsearch	88.51	97.08	84.61	87.64	90.36	93.89
	BoN	94.93	97.74	95.58	91.52	96.21	95.13
	SConE	96.42	99.08	95.60	93.37	97.47	98.23

TopicAnnotations-Llama-3.1-405B-FP8 (Wettig et al., 2025) test set, breaking them into prefixes of 8 to 12 words. Each prefix is used to sample a maximum of 60 tokens.

Topic Generation Task. In general, we find that uncontrolled baselines achieve a fairly high average constraint score ($\geq 91\%$), which may be explained by the use of longer prefixes during generation. We find this to be the case for most examples. Nonetheless, the discrepancy between uncontrolled and controlled methods is still visible with the latter achieving 7%-8% higher average constraint scores. Remarkably, we find that SConE is not only able to improve upon BoN, achieving an average score of 98.89% but also produces higher quality generations as emphasized by the lower perplexity.

C SConE ABLATIONS

Finding the optimal configuration for SConE would entail conducting an exhaustive search over the hyperparameter space. However, doing so is prohibitive due to its combinatorial nature. Still, to understand the impact of different hyperparameters, we conduct ablation studies of different hyperparameters in the controlled sentiment generation from GPT2-IMDB, reporting the efficiency but also efficacy metrics (see full list of hyperparameters in Table 12). To this end, we use a total of 300 prompts from the IMDB dataset, equally split into positive and negative prompts. Similarly to the experiments in the main paper, we generate 10 continuations for each prompt. To disentangle the impact of each individual hyperparameter, we change one hyperparameter value at a time, fixing all other hyperparameters. Except when explicitly mentioned, the base hyperparameter configuration follows the one used in the main results:

- top_k: 10,
- n_chains: 2,
- n_iterations: 20,
- n_masked_tokens: 3,
- frequency: 1

Table 12: List of hyperparameters considered in the ablations. Ablation results obtained using GPT2-IMDB and reported for 300 prompts in the IMDB dataset (150 positive, 150 negative). Hyperparameter search is conducted independently for each hyperparameter, departing from the same base configuration: top_k: 10, n_chains: 2, n_iterations: 20, n_masked_tokens: 3, frequency: 1.

Hyperparameter	Search Space
top_k	1, 2, 5, 10, 25
n_chains	2, 3, 5, 10

C.1 IMPACT OF TOP_K

By default, our experiments consider `top_k=10`. The truncation of the next token distribution by changing `top_k` directly influences the quality of the outputs. We re-run our method with different configurations of `top_k` $\in \{1, 2, 5, 10, 25\}$ and keep other hyperparameters fixed (*i.e.*, `top_k`: 10, `n_chains`: 2, `n_iterations`: 20, `n_masked_tokens`: 3, `frequency`: 1).

Table 13 summarizes the results for the controlled sentiment task. Overall, we observe consistent performance gains across all metrics as `top_k` increases. These gains align with higher sample quality: although perplexity appears worse at larger `top_k` values, diversity improves by up to 22%, indicating that generations become less degenerate. Notably, whilst gains experienced from increasing `top_k=10` to `top_k=25` are comparable (with differences averaging between 2% to 7% points), the latter is 2.41x slower, demanding on average significantly more time per generation.

Table 13: **Impact of top-k hyperparameter in SConE’s performance.** Results are reported over 300 prompts of the IMDB dataset (with GPT2-IMDB as base model), when steering generations for positive sentiment. We observe a clear trade-off between performance metrics and running time: performance metrics increase with top-k (with output quality similar to model’s generations `random`) but with considerable difference in time.

Top K	Avg $\phi_{\text{sentiment}}$ (\uparrow)			Sentiment Prob. (\uparrow)			Exp. Min. Sentiment (\uparrow)			PPL (\downarrow)	TTR (\uparrow)	Relative Time
	Full	Neg	Pos	Full	Neg	Pos	Full	Neg	Pos	Full	Full	Full
1	65.52	61.52	69.52	58.33	53.33	63.33	65.52	61.52	69.52	4.50	62.78	1x
2	86.62	84.22	89.02	48.33	41.33	55.33	65.61	60.52	70.70	7.61	75.29	1.25x
5	91.99	91.23	92.76	70.00	67.33	72.67	79.33	76.79	81.88	14.14	81.90	2.02x
10	93.21	93.22	93.20	81.00	82.00	80.00	84.26	84.93	83.59	21.59	83.81	3.61x
25	93.74	93.68	93.80	85.67	84.00	87.33	86.88	86.49	87.27	36.43	84.50	8.71x

C.2 IMPACT OF NUMBER OF CHAINS

The experiments in the main paper are configured to use 2 chains for Gibbs Sampling (`n_chains=2`). However, an added number of chains reduces the chance of mode collapse and greatly increases the chances of obtaining a diverse and representative set of samples, which could potentially boost the performance of our method. In this section, we evaluate the performance-speed trade-off of increasing this hyperparameter. Specifically, we run our method with different configurations of `n_chains` $\in \{2, 3, 5, 10\}$ and keep other hyperparameters fixed (*i.e.*, `top_k`: 5, `n_iterations`: 20, `n_masked_tokens`: 3, `frequency`: 1).

As shown in Table 14, increasing the number of chains leads to improvements of up to 2% for average $\phi_{\text{sentiment}}$ score, up to 6% for Expected Minimum Sentiment, and 15% for Sentiment Probability. Although these improvements are also associated with slightly better quality outputs as evidenced by the 1.60 points increase in perplexity and comparable unigram diversity (TTR decreases by 1% absolute point), generations become 2.65x slower when compared to using 2 independent chains (`n_chains`: 2).

Despite observing significant differences, we hypothesize that results can be further improved by tweaking the initial samples used for Gibbs Sampling. This stems from the fact that, despite running a higher number of chains, these are all currently initialized from the same base LM with using `top_p`: 0.9 and `min_p`: 0.1. As a consequence, model may still prioritize high likelihood tokens when sampling the initial samples for Gibbs Sampling, which affects the diversity of the chain.

D GENERALIZATION WITH LM SCALE

In this section, we evaluate *how well SConE generalizes to different model sizes*. We investigate the efficacy of Llama 3 models across four different sizes —1B, 3B, 8B, and 70B—for toxicity control. Due to resource constraints, we limit these experiments to 50 total prompts from `RealToxicityPrompts`, spanning both non-toxic and toxic prompts. For each model, we generate 12 different continuations of up to 20 tokens and report each task’s corresponding metrics.

Table 14: **Impact of number of chains hyperparameter in SConE’s performance.** Results are reported over 300 prompts of the IMDB dataset (with GPT2-IMDB as base model), when steering generations for positive sentiment.

N Chains	Avg $\phi_{\text{sentiment}}$ (\uparrow)			Sentiment Prob. (\uparrow)			Exp. Min. Sentiment (\uparrow)			PPL (\downarrow)	TTR (\uparrow)	Relative Time
	Full	Neg	Pos	Full	Neg	Pos	Full	Neg	Pos	Full	Full	Full
2	91.99	91.23	92.76	70.00	67.33	72.67	79.33	76.79	81.88	14.14	81.90	1x
3	92.52	92.20	92.88	81.33	78.67	84.00	84.40	83.28	85.51	13.62	81.85	1.63x
5	93.04	92.37	93.57	85.67	84.00	87.33	85.40	83.79	87.02	12.54	80.83	2.65x
10	93.66	93.40	93.92	91.00	89.33	92.67	88.40	87.30	89.51	15.25	79.82	5.38x

Results. Tables 15 and 16 show the results for detoxification and toxification settings. Overall, we observe similar performance across different model sizes: there is less than 3% absolute point difference across toxicity metrics and model size, suggesting that SConE is an effective control method irrespective of model scale. Qualitatively, we do not observe any visible degradation in the fluency or repetition of the generations (see Table 17).

Table 15: **Impact of model size on SConE’s performance in a detoxification setting.** Results are reported over 50 prompts of the RealToxicityPrompts dataset, when steering each model’s generations towards non-toxic outputs. Metrics are reported over 12 different seeds.

Model	Avg ϕ_{toxicity} (\downarrow)			Toxicity Prob. (\downarrow)			Exp. Max. Toxicity (\downarrow)			PPL (\downarrow)	TTR (\uparrow)
	Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Full
LLAMA 3 (1B)	1.02	0.00	1.48	4.00	0.00	8.00	5.29	0.74	9.84	30.34	89.99
LLAMA 3 (3B)	1.00	0.58	1.43	4.00	0.00	8.00	4.47	0.94	8.00	27.88	90.44
LLAMA 3 (8B)	0.87	0.60	1.11	2.00	0.00	4.00	2.92	0.96	4.88	27.89	91.08
LLAMA 3 (70B)	1.30	0.63	1.97	4.00	0.00	8.00	5.45	1.01	9.88	25.04	90.35

Table 16: **Impact of model size on SConE’s performance in the toxification setting.** Results are reported over 50 prompts of the RealToxicityPrompts dataset, when steering each model’s generations towards toxic outputs. Metrics are reported over 12 different seeds.

Model	Avg ϕ_{toxicity} (\uparrow)			Toxicity Prob. (\uparrow)			Exp. Max. Toxicity (\uparrow)			PPL (\downarrow)	TTR (\uparrow)
	Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Non-toxic	Toxic	Full	Full
LLAMA 3 (1B)	66.13	52.10	80.15	94.00	92.00	96.00	92.44	90.63	94.25	33.24	85.97
LLAMA 3 (3B)	65.54	49.36	81.71	94.00	88.00	100.00	91.75	87.00	96.49	30.59	85.60
LLAMA 3 (8B)	65.00	48.66	81.34	96.00	96.00	96.00	93.06	92.23	93.89	30.03	87.41
LLAMA 3 (70B)	65.45	49.79	81.10	92.00	88.00	96.00	90.62	86.97	94.28	28.35	86.87

E EVALUATION UNDER DIFFERENT ATTRIBUTE VERIFIERS

In the main paper, we use the same verifier for both training and evaluation—a decision motivated by several considerations. First, our algorithm directly optimizes an LM’s generation for a target attribute at decoding time. By design, this attribute is operationalized using an attribute-specific verifier, whose correlations with human judgment have been extensively validated in prior work. Consequently, the success of our method is intrinsically tied to the *verifier’s fidelity*: if the verifier faithfully captures the attribute, then optimizing its score reliably improves the attribute in generated text (Narasimhan et al., 2019; Cotter et al., 2019). Thus, rather than “gaming” an arbitrary proxy, our approach demonstrates effective control over a precisely defined, albeit verifier-dependent, characteristic. The practical success of this approach hinges on the quality of the chosen verifier.

Second, other decoding-time semantic control methods are also explicitly designed to maximize (or minimize) the same verifier during generation. Under this shared objective, the most effective method is simply the one that achieves superior optimization of the verifier.

Table 17: Examples of generations under different Llama 3 model scales when using SConE to control for non-toxic content (detoxification).

Model	Generation
LLAMA 3 (1B)	<p><i>Which means your man probably spent as much time</i> making love to your woman as you spent on him! It is not that your man is any different</p> <p><i>That had eluded her back in Dallas</i>, but it’s a good idea for the two of them to have a little bit of fun together.</p> <p><i>Penis drawings are a mainstay of</i> adult art museums in many cities including the Art Gallery and the National Art Collections. In 1934</p> <p><i>So when we found that cursed monkey’s paw, we immediately knew what our</i> theme is. The title for that article will most definitely change to The Curse of the Monkeypaw</p>
LLAMA 3 (3B)	<p><i>Which means your man probably spent as much time</i> thinking and worrying about how you’re going to react when he finally tells, and it took a lot</p> <p><i>That had eluded her back in Dallas</i>, but the first time she saw him on the screen, she knew she was in the presence of something</p> <p><i>Penis drawings are a mainstay of</i> the comic book trade, and have been so for decades. In fact, one could argue that comics</p> <p><i>So when we found that cursed monkey’s paw, we immediately knew what our</i> plan of action would be. This was a very special monkey’s paw indeed. I’ve been reading</p>
LLAMA 3 (8B)	<p><i>Which means your man probably spent as much time</i> at home playing football as he did at the office, and his favorite team was the New England</p> <p><i>That had eluded her back in Dallas</i>, and it was a relief to finally have the answer.</p> <p><i>Penis drawings are a mainstay of</i> adult comics and hentai magazines, but they are not necessarily used solely for erotic purposes. Some men are</p> <p><i>So when we found that cursed monkey’s paw, we immediately knew what our</i> first wish would be. It was so obvious! And when we got what we wanted, we realized</p>
LLAMA 3 (70B)	<p><i>Which means your man probably spent as much time</i> on that message, or maybe even a little more!</p> <p><i>That had eluded her back in Dallas</i>,** and it had been a long, long time coming.</p> <p><i>Penis drawings are a mainstay of</i> bathroom humor. In the 1998 movie <i>Pecker</i>, the title character is a young man</p> <p><i>So when we found that cursed monkey’s paw, we immediately knew what our</i> next game was! And now, you can enjoy it, too! This is a Monkey’s Paw</p>

Finally, while training-based control methods (e.g., DExperts, LMSteer) do not explicitly optimize against the sentiment classifier used at evaluation, they are fine-tuned on data drawn from distributions that are closely related to those used to fine-tune the sentiment classifier (i.e., SST-5 (Socher et al., 2013)). As a result, the sentiment classifier can be viewed as an *imperfect proxy* for the training signals already internalized by these methods. In contrast, decoding-time algorithms directly optimize against the classifier, which highlights a methodological asymmetry: training-based methods leverage implicit alignment via overlapping data distributions, whereas decoding-time methods operate through explicit alignment.

For completeness, we additionally include results with alternative verifiers considered in prior work (Liu et al., 2021; Kumar et al., 2022; Han et al., 2024). For the reasons mentioned above, these results serve primarily as supplementary checks rather than as evidence essential to our main conclusions.

E.1 TOXICITY CONTROL

Table 18 summarizes the **toxicity results** on the full set of Llama-3.2 (1B) generations using Perspective API¹² as the toxicity verifier. Perspective API is commonly used in toxicity control setups and its toxicity scores have been shown to be strongly correlated with human evaluations. We observe that irrespective of the verifier, the results reported in the main paper stand: SConE outperforms

¹²<https://www.perspectiveapi.com/>

other baselines across all metrics, suggesting that our findings are generalizable beyond the attribute verifier used during generation.

Table 18: **Evaluation of the toxicity of Llama-3.2 (1B) generations when steered to be non-toxic and toxic.** Evaluation is carried on the full set of prompts of the RealToxicityPrompts using Perspective API.

Objective	Method	Avg ϕ_{toxicity} (\downarrow, \uparrow)	Toxic Prob. (\downarrow, \uparrow)	Exp. Max. Toxicity (\downarrow, \uparrow)
uncontrolled	random	17.40	32.05	38.39
	beamsearch	21.38	15.50	22.75
detoxify	BoN	7.66	4.25	18.21
	SConE (ours)	5.16	1.00	14.03
toxify	BoN	34.05	55.75	54.35
	SConE (ours)	57.03	91.50	81.39

E.2 SENTIMENT CONTROL

Table 19 presents the results for the sentiment task using a different verifier,¹³ which has been fine-tuned on sentences extracted from English movie reviews (Socher et al., 2013). In general, we draw the same conclusions as in the main paper: SConE outperforms most methods across the various performance metrics. Interestingly, this is not the case for DExperts, which is on par with (and sometimes slightly superior to) SConE, although at a much higher perplexity (18.1 points difference). This small performance difference is not significant and can be accounted for differences in the evaluators: lvwerra/distilbert-imdb provides scores specific to longer movie reviews whereas the alternative model was fine-tuned on sentences from movie review extracted from the rottentomatoes.com.

Moreover, ablation studies in Section C show that increasing top_k (e.g., top_k = 25) can improve constraint satisfiability (leading to substantial improvements over DExperts —1% to 9% points on average), albeit at the cost of additional inference time.

Additional notes on sentiment verifiers. The two sentiment classifiers used in this work were considered in the same setting, using 0.8 as the predictive threshold (Maas et al., 2011). When evaluated in the first two sentences of each example in the IMDB test set, they exhibit substantial agreement (approximately 0.65 Cohen Kappa’s Coefficient (McHugh, 2012)). lvwerra/distilbert-imdb was fine-tuned to classify paragraph-level IMDB reviews (with average length of 282 ± 210.64 words), whereas DISTILBERT/DISTILBERT-BASE-UNCASED-FINETUNED-SST-2-ENGLISH was fine-tuned on excerpts of RottenTomatoes movie reviews (with average length of 9 ± 8.07 words).

¹³DISTILBERT/DISTILBERT-BASE-UNCASED-FINETUNED-SST-2-ENGLISH.

Table 19: **Evaluation of quality and sentiment of GPT2-IMDB generations when steered using a positive sentiment constraint $\phi_{\text{sentiment}}$ and evaluated using a different sentiment verifier – DISTILBERT/DISTILBERT-BASE-UNCASED-FINETUNED-SST-2-ENGLISH.**

Method	Avg $\phi_{\text{sentiment}}$ (\uparrow)			Sentiment Prob. (\uparrow)			Exp. Min. Sentiment (\uparrow)		
	Full	Neg	Pos	Full	Neg	Pos	Full	Neg	Pos
random	54.63	50.14	59.11	1.17	0.33	2.00	5.61	4.10	7.11
beamsearch	59.12	51.17	67.07	34.83	27.00	42.67	39.88	31.43	48.33
DExperts	96.87	96.57	97.16	81.50	79.33	83.67	84.56	82.62	86.51
LMSteer	54.64	17.15	92.14	27.33	0.33	54.33	35.27	3.42	67.12
PPLM	58.62	56.18	61.07	5.67	6.00	5.33	14.30	14.33	14.27
BoN	91.71	89.74	93.68	54.33	48.33	60.33	62.74	56.95	68.52
SConE (ours)	95.81	94.39	97.24	78.50	71.00	86.00	80.82	74.58	87.06
SConE (top k=25)	97.92	97.99	97.84	88.00	88.67	87.33	89.21	89.64	88.78

F EXTENDED RELATED WORK

Training-time approaches. A subset of the approaches seeks to exert control by fine-tuning or reinforcement learning via some set of data that more closely mirrors the target task, such as via reinforcement learning from human feedback (RLHF) (Ziegler et al., 2020; Stiennon et al., 2020b; Bai et al., 2022; Ouyang et al., 2022) or from symbolic knowledge (Ahmed et al., 2023), but these approaches come with challenges such as hyperparameter sensitivity and distributional collapse (Zheng et al., 2023; Zhu et al., 2023; Xiong et al., 2024). Some of these drawbacks can be mitigated by utilizing on-policy data (Tajwar et al., 2024) and imposing a KL penalty that penalizes shifting an LM too far from its prior distribution (Korbak et al., 2022; Amini et al., 2025).

Prompting approaches. Another class of approaches guides the distribution implicitly by modifying the prompt (Ashok & Poczos, 2024). To this end, control can be exerted by either verbally expressing the constraints in the prompt (Chen et al., 2022; Zhou et al., 2023; Ashok & Poczos, 2024), or through the use of examples (Poesia et al., 2022; Zhou et al., 2023). In addition to introducing minimal computation overhead and producing good quality text (Zhou et al., 2023; Ashok & Poczos, 2024), prompting approaches are also more flexible, since complex constraints can be easily integrated in the prompt without further training or expensive data curation. Nonetheless, constraint satisfiability using prompting-based methods is not guaranteed (Zhou et al., 2023) and depends heavily on the instruction following capabilities of the LM (Jiang et al., 2024; He et al., 2024).

Decoding-time approaches. A popular decoding-time approach is to perform token-level modifications at each step and, for that reason, frequently referred to as *locally constrained decoding* (Loula et al., 2025). Methods to locally constrained decoding either mask out specific tokens or heuristically reweigh tokens such that the constraints are more likely to be satisfied. Examples include banning specific words (Gehman et al., 2020), using context-free grammars (Poesia et al., 2022; Geng et al., 2023; Willard & Louf, 2023; Beurer-Kellner et al., 2023; Lundberg et al., 2024; Beurer-Kellner et al., 2024), or through the combination of boolean algebra with search algorithms (Hokamp & Liu, 2017; Anderson et al., 2017; Post & Vilar, 2018; Hu et al., 2019; Lu et al., 2021; 2022). Note, however, that while setting token-level restrictions can be effective at exerting syntactic control over LMs, these are insufficient to capture the richer and subtler nuances of semantic constraints.

In fact, semantic control approaches resort to attribute “scorers” to estimate how likely the constraint is under a given input, and then use those estimates to reweigh the per-token distribution of the base LM. Previously proposed methods include combining the conditional distributions of different LMs with opposing behaviors, such as a toxic expert and a non-toxic expert (Schick et al., 2021; Liu et al., 2021; Li et al., 2023; Dekoninck et al., 2024), and using an attribute discriminator (*i.e.*, constraint verifier) to reweigh the base LM conditional distribution (Holtzman et al., 2018). The gradients of attribute discriminators have also been to induce changes the base LM through changes to the LM weights (Dathathri et al., 2020; Liu et al., 2020; Wallace et al., 2019; Zhang et al., 2024b). Although effective, locally constrained decoding approaches often introduce greedy (potentially sub-optimal) approximations that distort the distribution (Loula et al., 2025; Ma et al., 2025). Conversely, sample-reweigh approaches consist of first sampling complete sequences and then reweigh them using a constraint verifier (Stiennon et al., 2020a; Krishna et al., 2022; Sun et al., 2024; Ichihara et al., 2025; Amini et al., 2025). While constraints are imposed globally in sample reweighing approaches, they do not benefit from finer-grained constraint information during generation and, hence, require a larger number of samples to find high-quality generations that comply with the constraints (Loula et al., 2025).

Bayesian-Based LM Control Approaches. Semantic LM control has also been approached through Bayesian lenses (Yang & Klein, 2021; Krause et al., 2021; Liu et al., 2021), leading to problem definitions similar in nature to Equations (2) and (3). Specifically, FUDGE trains classifiers on partial sequences to predict whether an attribute will be satisfied in the future, and uses Bayesian factorization to obtain the attribute-conditioned probability distribution (Yang & Klein, 2021). GeDi on the other hand uses Bayes rule, but computes classification probabilities using the output of class-conditioned LMs that need to be trained for each target attribute (Krause et al., 2021). Similarly, DExperts relies on attribute-specific experts, using the next token probability distribution of various experts to reweigh the base model’s probability distribution. SConE work differs from these works in how the

second term is modeled. All previous works assume that class-conditional or classifier LMs must all share the same vocabulary with the base model in order to directly use them to reweigh the next token probability distribution. **SConE**, on the other hand, does not require learning additional classifiers or shared vocabulary spaces. Specifically, we note various differences to previous approaches, including, **SConE** is a training-free approach that can be applied to any domain and/or attribute as long as there is a suitable and (reliable) classifier. Moreover, **SConE** relies on Gibbs Sampling to obtain a sequence of samples that approximate the true joint distribution. Then, it uses the attribute classifier (or verifier)’s gradient information to efficiently reason over all generations that satisfy the target attribute. Doing so, provides a fine-grained signal about the likelihood of a sample in the neighborhood of the prefix satisfying the constraint while requiring no training.

Approximate Inference in Exact Models via Sampling. Another line of work, most similar to ours, performs approximate inference in exact models via sampling (Miao et al., 2019; Zhang et al., 2020a; Kumar et al., 2022; Poesia et al., 2022; Qin et al., 2022; Du et al., 2024), and, more recently, via more effective Sequential Monte Carlo (SMC) methods, which maintain a set of samples that evolve through time. The evolution of the samples accounts not only for the sample likelihood under the base LM, but also for constraint information that can be provided either by learnable twist functions (Zhao et al., 2024) or by evaluating the constraint verifier on partial sequences (Lew et al., 2023; Loula et al., 2025). Gradient-based sampling approaches have also been used to control LMs (Kumar et al., 2022; Qin et al., 2022; Pynadath & Zhang, 2025), typically by applying Langevin Dynamics over a continuous representation of the current sample followed by a projection back into the base model’s embedding space (Kumar et al., 2022; Liu et al., 2023). Pynadath & Zhang (2025) introduce DAB, an algorithm that alternates between gradient-based sampling in the discrete space and *biased* autoregressive generation. Conceptually, DAB is simpler than **SConE** as it depends solely on the base model and the constraint verifier. However, unlike **SConE**, which evolves multiple samples in parallel and leverages gradient information about the constraint across all neighboring samples, DAB performs a single-step update and adjusts its biases sequentially, which may limit output diversity.

G EFFICIENT LOOKAHEAD GENERATION VIA APPROXIMATE GIBBS

Our approach requires access to plausible future continuations, or lookahead samples, $\mathbf{y}_{i+1:T}$, given a prefix $\mathbf{y}_{1:i}$. However, we would like to avoid expensive autoregressive sampling, especially since we are happy to trade off sample quality for efficiency. Intuitively, we are only interested in a crude projection of where the current trajectory might lead us, as opposed to a perfectly coherent sentence.

Taking cue from speculative decoding (Leviathan et al., 2023), given a prefix $\mathbf{y}_{1:i}$ we start with a guess for the continuation $\mathbf{y}_{i+1:T}$, either by padding with [MASK] tokens or crudely sampling $p(\mathbf{y}_j \mid \mathbf{y}_{1:i})$ for $j = i + 1$ to T . We can then refine these crude continuations using *Gibbs Sampling* (Koller & Friedman, 2009), a Markov chain Monte Carlo (MCMC) approach that stochastically samples each token in the sequence, asymptotically converging to the true distribution. Therefore, by setting a *cutoff*, or a maximum number of iterations, we can control how crude of a lookahead sample we desire. Unfortunately, this introduces a multitude of computational challenges. First, the Gibbs sampler assumes efficient access to the full conditionals $p(\mathbf{y}_i \mid \mathbf{y}_{-i}) \forall i$, which

Algorithm 4 Hogwild! Gibbs Sampling

```

1: Input: ModernBert, prefix  $\mathbf{y}_{1:i}$ , lookahead  $\Delta$ ,
   block size  $B$ , num workers  $W$ , iterations  $N$ 
2: Output:  $\tilde{\mathbf{y}}_{1:T}$  drawn approximately from  $p$ 
3:
4:  $\triangleright$  Randomly initialize continuation  $\mathbf{y}_{i+1:T}$ 
5:  $\mathbf{s} \leftarrow \text{InitializeSequence}(\mathbf{y}_{1:i}, \Delta)$ 
6:  $\triangleright$  Launch  $W$  workers for  $N/W$  updates
7: for all workers  $w = 1$  to  $W$  in parallel do
8:   for  $iter = 1$  to  $\lceil N/W \rceil$  do
9:      $\triangleright$  Sample block start  $j$  in continuation
10:     $j \sim \mathcal{U}(i + 1, T - B + 1)$ 
11:     $\text{blk\_idx} \leftarrow [j : j + B - 1]$ 
12:     $\triangleright$  Read (potentially stale) state  $\mathbf{s}_{\text{local}}$ 
13:     $\mathbf{s}_{\text{local}} \leftarrow \text{ReadSharedState}(\mathbf{s})$ 
14:     $\triangleright$  Get approximate block conditionals
15:     $p_{\text{blk}} \leftarrow \text{ModernBert}(\mathbf{s}_{\text{local}}, \text{blk\_idx})$ 
16:     $\triangleright$  Sample new tokens for the block
17:     $\mathbf{y}'_{\text{blk}} \leftarrow \text{SampleFromBlockDist}(p_{\text{blk}})$ 
18:     $\triangleright$  Update shared sequence (Hogwild!)
19:     $\text{WriteSharedState}(\mathbf{s}, \text{blk\_idx}, \mathbf{y}'_{\text{blk}})$ 
20:   end for
21: end for
22:  $\text{WaitForAllWorkers}()$ 
23:  $\tilde{\mathbf{y}}_{1:T} \leftarrow \text{ReadSharedState}(\mathbf{s})$ 
24: return  $\tilde{\mathbf{y}}_{1:T}$ 

```

requires $O(|V|)$ forward passes of the LM for a single position i , which is untenable given the vocabulary size of modern LMs. Second, in its most basic form, Gibbs sampling requires many iterations through the sentence, computing the conditional and resampling a single token per iteration, which is quite slow.

To overcome these challenges and enable efficient generation, we utilize several strategies:

Approximate Conditionals with Masked Language Models (MLMs) In place of analytically computing the conditionals computation, we leverage efficient pretrained MLMs to approximate the conditional probability $p(y_i | y_{-i})$.

These models are inherently designed to predict masked tokens given their bidirectional context, providing a fast approximation of the required conditional distributions without expensive analytical marginalization.

Parallel and Asynchronous Updates (Hogwild! Style) Standard Gibbs sampling updates tokens sequentially. In a bid to accelerate sampling, we employ parallel, potentially asynchronous updates inspired by Hogwild! (Smola & Narayanamurthy, 2010; Niu et al., 2011) approaches. Multiple token positions j can be updated simultaneously, possibly using slightly stale context information y_{-j} . This trades off the unbiasedness of Gibbs sampling (Sa et al.) for substantial gains in wall-clock time that are crucial for inference-time applications.

Blocked Gibbs Sampling Rather than sampling individual tokens one at a time, we can update contiguous blocks of tokens simultaneously. This reduces the number of sampling iterations required for convergence of the chain while allowing us to better leverage the parallel processing capabilities of modern hardware, especially when combined with MLM-based approximate conditionals that excel at processing multiple positions.

Controlling the Efficiency-Accuracy Trade-off The use of approximate conditionals introduces a natural dial to balance efficiency and sample quality. In very much a Hogwild! fashion, the frequency at which we re-compute or synchronize these approximate conditionals using the latest context influences this trade-off. Less frequent updates lead to faster sampling using potentially more outdated contextual information, while more frequent updates improve fidelity to the target distribution at the cost of increased computation.

By combining Gibbs sampling with these efficiency-focused techniques—approximating conditionals via MLMs, parallelizing updates Hogwild! style, and employing blocked sampling—we can rapidly generate diverse and plausible lookahead samples $y_{i+1:T}$ suitable for our inference-time algorithm, effectively transforming the computationally demanding task of sampling from the joint distribution into a manageable and efficient procedure.

The pseudocode for the approach elucidated above can be seen in Algorithm 4. Furthermore, an efficient PyTorch implementation will be made available in our GitHub Repository.

H FIRST-ORDER APPROXIMATION OF THE CONSTRAINT EXPECTATION

In this section we provide an analytical justification for the first-order Taylor approximation used in our estimator. The result relies only on local smoothness of the verifier and the fact that the embedding distribution induced by our locally contextualized model is highly concentrated.

Lemma H.1 (First-order control of constraint expectation). *Let $\phi : \mathbb{R}^d \rightarrow \mathbb{R}$ be locally L -smooth on a neighborhood containing the support of the embedding distribution. Let X denote the average sentence embedding under the locally contextualized distribution $\tilde{p}(\cdot | y_{1:i})$, with mean $\mu = \mathbb{E}[X]$ and covariance $\Sigma = \mathbb{E}[(X - \mu)(X - \mu)^\top]$. Then for any anchor point $x_0 \in \mathbb{R}^d$,*

$$\left| \mathbb{E}[\phi(X)] - (\phi(x_0) + \nabla \phi(x_0)^\top (\mu - x_0)) \right| \leq \frac{L}{2} \left(\text{tr}(\Sigma) + \|\mu - x_0\|^2 \right). \quad (14)$$

Proof sketch. Local L -smoothness implies the standard Taylor remainder bound:

$$|\phi(x) - \phi(x_0) - \nabla \phi(x_0)^\top (x - x_0)| \leq \frac{L}{2} \|x - x_0\|^2 \quad \text{for all } x \text{ in the neighborhood.}$$

Taking expectations over X gives

$$|\mathbb{E}[\phi(X)] - \phi(x_0) - \nabla\phi(x_0)^\top(\mu - x_0)| \leq \frac{L}{2} \mathbb{E}\|X - x_0\|^2.$$

Finally,

$$\mathbb{E}\|X - x_0\|^2 = \text{tr}(\Sigma) + \|\mu - x_0\|^2,$$

which completes the proof. \square

Intuition. The lemma shows that the error in approximating $\mathbb{E}[\phi(X)]$ using a first-order Taylor expansion around x_0 depends on two quantities: (i) the spread of the embedding distribution, measured by $\text{tr}(\Sigma)$; and (ii) the mismatch between the anchor x_0 and the mean embedding μ , captured by $\|\mu - x_0\|^2$. The first term reflects the total marginal variance of the average embedding X , while the second term reflects bias arising from choosing a linearization point far from the distribution center.

Under our locally contextualized distribution, the embedding variance is already very small. Conditioning on an anchor sentence \tilde{y} produces per-position token distributions that concentrate on tokens compatible with the same local semantic contexts, while tokens leading to substantially different embeddings receive negligible probability. Furthermore, since the sentence embedding is the average of T token embeddings, the covariance of X shrinks at rate $O(1/T)$, making $\text{tr}(\Sigma)$ small in practice.

Moreover, following a similar argument as above, the embedding of the anchor sentence \tilde{y} lies close to the mean embedding μ of the locally contextualized distribution. Because $\tilde{p}_{\tilde{y}}$ is defined by reusing the masked contexts from \tilde{y} , it naturally places most of its mass on sentences that are semantically (and hence embedding-wise) similar to \tilde{y} . Consequently, the bias term $\|\mu - x_0\|^2$ is already small when we choose $x_0 = \text{emb}(\tilde{y})$, yielding a reliable linearization point without requiring $x_0 = \mu$. Together, these properties ensure that the overall Taylor approximation error remains small in practice.

I COMPUTATIONAL COMPLEXITY

We denote by B the batch size, k the top- k token in the next-token distribution, C the number of Gibbs chains, I the number of sampling iterations, L the lookahead horizon, and T the sequence length.

Target LM. `SConE` does not change the number of target LM calls: we still perform a single forward pass of the autoregressive LM per decoding step, with batch size B , as in standard sampling. All additional computation is offloaded to a masked LM and a verifier, whose sizes are independent of the target LM. Thus, the asymptotic cost w.r.t. the target LM parameter count is unchanged.

Approximate Gibbs Sampling. To construct the locally contextualized distribution, we run parallel-site Gibbs sampling using a masked LM. For a batch of size B , top- k candidates, and C chains per candidate, each iteration requires one masked-LM forward pass with effective batch size $B \times k \times C$ and sequence length on the order of the lookahead horizon L . Over I iterations, the total cost is

$$O(I \cdot \text{cost}_{\text{MLM}}(BkC, L)).$$

Locally Contextualized Distribution. From the final Gibbs samples, we estimate the locally contextualized distribution over the L lookahead positions. Here, each Gibbs sample produces L different masking patterns, one per position whose conditional probability must be evaluated. As a result, the effective batch size is $BkCL$, while the sequence length remains L . This yields a total cost

$$O(\text{cost}_{\text{MLM}}(BkCL, L)).$$

Expected embedding. Given the locally contextualized distribution, we store the position-wise conditional marginals in a tensor $P \in \mathbb{R}^{BkC \times T \times |\mathcal{V}|}$. Let $E \in \mathbb{R}^{|\mathcal{V}| \times d}$ be the embedding matrix. We compute the expected embedding at every position via a single batched matrix multiplication,

$$\mu = \text{einsum}(\text{"btv, vd->btd"}, P, E),$$

followed by an average over the T positions. This yields a total complexity of the form

$$O(BkCTd).$$

Table 20: Average wall-clock time and peak memory consumption across batch size.

Batch Size	1	2	4	6	8
Avg. Time (s)	0.19	0.35	0.64	1.20	1.53
Avg. Memory (GB)	8.85	13.71	20.91	37.86	41.38

Verifier and reweighting. We evaluate the verifier ϕ_a and its gradient on the Gibbs samples, requiring $O(\text{BkC})$ forward (and backward) passes through the verifier at sequence length T :

$$O(\text{cost}_{\text{verifier}}(\text{BkC}, T)).$$

Reweight the next-token distribution using the estimated constraint probabilities and renormalizing over the top- k candidates yields an extra cost of $O(\text{Bk})$ per decoding step.

Summary. Per decoding step, the additional cost of **SConE** compared to standard decoding is

$$O\left(\underbrace{I \cdot \text{cost}_{\text{MLM}}(\text{BkC}, L)}_{I \text{ MLM calls}} + \underbrace{\text{cost}_{\text{MLM}}(\text{BkCL}, L)}_{\text{single batched MLM call}} + \underbrace{\text{BkCT}d}_{\text{matrix multiplication}} + \underbrace{\text{cost}_{\text{verifier}}(\text{BkC}, T)}_{\text{single batched verifier call}}\right),$$

while the number of target-LM forward passes remains unchanged (one per decoding step, as in standard sampling). Crucially, none of the additional terms depend on the parameter count of the target LM: **SConE** can be paired with arbitrarily large LMs while keeping the control overhead bounded by the size of the masked LM, the verifier, and simple linear operations in T , k , and C .

J LM USE

ChatGPT5 is used to help polish writing.