LINEARLY CONTROLLED LANGUAGE GENERATION WITH PERFORMATIVE GUARANTEES

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025 026 027

028

Paper under double-blind review

Abstract

The increasing prevalence of Large Language Models (LMs) in critical applications highlights the need for controlled language generation strategies that are not only computationally efficient but that also enjoy performance guarantees. To achieve this, we use a common model of concept semantics as linearly represented in an LM's latent space. In particular, we take the view that natural language generation traces a trajectory in this continuous semantic space, realized by the language model's hidden activations. This view permits a control-theoretic treatment of text generation in latent space, in which we propose a lightweight, gradient-free intervention that dynamically steers trajectories away from regions corresponding to undesired meanings. Crucially, we show that this intervention, which we compute in closed form, is guaranteed (in probability) to steer the output into the allowed region. Finally, we demonstrate on a toxicity avoidance objective that the intervention steers language away from undesired content while maintaining text quality.

WARNING: This paper contains model outputs which are offensive in nature.

1 INTRODUCTION

 Language Models (LMs) have become widespread in critical applications such as content moderation and real-time information dissemination (Zeng et al., 2024). Despite their transformative impact, these models require updates to remain accurate post-deployment. Moreover, as demand for more nuanced text generation rises, strategies that enforce constraints during text generation are increasingly needed. To address these challenges, controllable text generation has emerged as a pivotal research area.

Several approaches have been proposed towards controllable text generation (Kumar et al., 2021; Lu et al., 2021; Li et al., 2022; Qin et al., 2022). Of them, a popular approach is prompt engineering (Luo et al., 2023; Bhargava et al., 2023; Cai et al., 2023), where natural language prompts are carefully chosen at input-time to steer generation. Other approaches modify LM parameters to achieve the desired outputs (Yao et al., 2023; Li et al., 2023b). Lastly, some approaches engineer LM *activations*, or input representations, to steer them into the representations of desired outputs (Dathathri et al., 2019; Hernandez et al., 2023; Konen et al., 2024; Li et al., 2024a).

Despite current efforts, ensuring the controllability of these models remains a challenge due to their limited interpretability. For instance, while knowledge editing provides an efficient alternative to exhaustive retraining, it poses risks akin to the butterfly effect: minor adjustments could lead to unintended consequences. Moreover, it is paramount for these approaches to be robust and ensure controllability guarantees to mitigate risks and harness their full potential safely.

To address this gap, we propose to use control theory to tackle controlled language generation. Specifically, optimal control theory (Kirk, 2004) offers principled methods to steer trajectories in latent space that enjoy theoretical guarantees on the performance of the intervention. In the framework of optimal control theory, our intervention method, which we call Linear Semantic Control (LiSeCo), derives from a theoretical formulation of controlled text generation. Our contributions are both theoretical and empirical: (1) we formally pose LM control as a constrained optimization problem and provide its closed-form solution with guarantees; (2) we empirically demonstrate our method on the use cases of toxicity and negativity avoidance. We confirm, with experiment corroborating theory, that LiSeCo indeed steers LM generation from disallowed concepts while maintaining text quality.

054 2 RELATED WORK

Contemporary language models are deep neural networks pre-trained on trillions of tokens of Internetscale text. In part due to their vast scale and lack of interpretability, methods to control them in a
fine-grained way remain elusive. A number of approaches have already been proposed towards this
end, spanning the whole spectrum of permanent (Meng et al., 2022b; Belrose et al., 2023) to online
control strategies (Liu et al., 2021; Dathathri et al., 2019). Here, we review post-hoc intervention
methods and situate LiSeCo with respect to the current landscape.

062Post-hoc intervention methods can intervene on various components of the LM: for instance, decoding,063like FUDGE and GeDI (Yang & Klein, 2021; Krause et al., 2021), activations, like LiSeCo, or weights064via finetuning. All such methods aim to modify some attribute, such as toxicity, while maintaining065text fluency. Ultimately, all methods work towards this goal by modifying the LM's final probability066distribution, either directly or indirectly. We can situate where different method classes intervene,067viewing an LM as a series of T function compositions corresponding to the T layers, where s is a068sequence of tokens:

 $\mathbb{P}_{LM}(s_i|s_{< i}) = f_T \circ f_{T-1} \cdots \circ f_1(s_{< i}) := LM(s_{< i}).$

Decoding-based methods fix the function $LM := f_T \circ f_{T-1} \circ \cdots \circ f_1$ and directly edit its output probability distribution $\mathbb{P}_{LM}(s_i|s_{\leq i})$ (Yang & Klein, 2021; Liu et al., 2021; Krause et al., 2021). These methods require access to an external evaluator whose feedback is used to calibrate token probabilities, which can result in high inference latency.

Prompt engineering is a technique that controls the LM's output by varying the input context $s_{<i}$, keeping the function $LM := f_T \circ f_{T-1} \circ \cdots \circ f_1$ fixed (Luo et al., 2023; Bhargava et al., 2023; Cai et al., 2023; Wei et al., 2022; Li & Liang, 2021). Prompts are often highly task-specific, requiring either manually crafting or ad-hoc computationally-taxing techniques, and success can be brittle to prompt choice (Weber et al., 2023). While the space of natural language prompts is discrete, LM weights and activations live in continuous high-dimensional space, which is more expressive; then, rather than search over discrete prompts, other approaches that exploit this expressivity directly intervene in the internals of the model.

083 Of them, weight-based methods modify the functions f_i themselves, which permanently constrains the space of final probability distributions $\mathbb{P}_{I,M}$. These methods comprise, e.g., reinforcement learning 084 from human feedback (Ouyang et al., 2022), instruction-tuning, parameter-efficient adaptation (Hu 085 et al., 2022), or targeted weight-editing (Meng et al., 2022b; Belrose et al., 2023). In such approaches, weights are modified according to the goal of the controlled generation by, for instance, learning 087 the necessary update (De Cao et al., 2021; Mitchell et al., 2021), or localizing and editing target 880 parameters encoding specific knowledge (Dai et al., 2022; Meng et al., 2022a;; Li et al., 2024b). 089 Pitfalls range from potential inconsistencies and distortions, to the fact that weight-based methods can only correct errors in the LM's parametric knowledge, but not in-context (Li et al., 2023b). 091

Activation-based methods, such as LiSeCo, fix $LM := f_T \circ f_{T-1} \circ \cdots \circ f_1$, but intervene at the 092 domain of each f_i , where introducing a steering vector transforms the input to f_i (Li et al., 2023a; Turner et al., 2023). These interventions can be seen as restricting the domain of each f_i , eventually 094 constraining the space of probability distributions \mathbb{P}_{LM} when composed up through the layers. A key advantage of activation steering is rapid adaptation that is context-dependent. An initial work in this 096 domain was Plug and Play (Dathathri et al., 2019), where a linear intervention is computed at every 097 layer. The control goal is encoded as the objective function in an optimization that is then solved 098 via back-propagation, adding significant computational overhead at inference time. Subsequent 099 approaches also compute linear modifications to the latent state, but reduce computational overhead, act on only a few layers (Subramani et al., 2022; Konen et al., 2024), or pre-compute steering vectors 100 to avoid back-propagation (Turner et al., 2023). The recent approach REMEDI (Hernandez et al., 101 2023) also finds an optimal intervention to achieve different target outputs, but, requires specific 102 a-priori training data for the representations, and offers no guarantees that the intervention will attain 103 the target output. The approach presented in ITI (Li et al., 2024a) addresses the issue of computational 104 efficiency at the expense of optimality (the intervention is not formulated as an optimizer), and lacks 105 guarantees on the intervention's performance. 106

107 The use of steering vectors for text generation proposed in the literature provides empirical grounding for the promise of this approach. However, none of the approaches found in the literature provide

guarantees on the controllability of their method. Here, we provide an intervention that is *theoretically guaranteed* to steer the input into the allowed region, and introduces minimal computational overhead.
 Our work differs from existing literature by the following novel contributions:

- 1. We formally derive the optimal steering vector and offer theoretical performance guarantees. This is enabled by the use of an optimal control framework to cast the problem for the first time in a domain that has been overwhelmingly empirical.
- 2. Deriving the closed-form solution for the intervention minimizes increased latency at inference time. We provide empirical comparisons, where we show that other popular methods, such as FUDGE, require much higher latency.

Though Soatto et al. (2023) apply theoretical tools from control to LM text generation, to the best of our knowledge, our method is the first to propose a control-theoretic intervention whose theoretical guarantees are validated in practice.

- 122 122 3 P
- 123 124 125

126

127 128

129

112

113

114

115

116

117

PROBLEM STATEMENT

In this section we present the problem studied in this paper, as well as the assumptions and approach. In particular, we approach the problem of controlled language generation as a standard optimal control problem in the field of control theory (Kirk, 2004).

3.1 PROBLEM FORMULATION

Given a language model $\ell : \Sigma^* \to \Sigma^*$, controlled language generation aims to steer the model's output into a desired one. In this work, we study how to steer the output of an *already trained model* away from a disallowed region, or so that the set of possible generated sequences is constrained to an allowed subset $S \subset \Sigma^*$. The requirements for the generated output sequence are two fold: its latent trajectory (1) is *guaranteed* to never lie in the disallowed region, and (2) stays as close as possible to that of the original output sequence, so that text quality is not compromised. In doing this, two questions need to be answered:

137 138

139

140

- 1. How can the disallowed region be defined for a given language model?
- 2. How can an intervention be designed to be guaranteed (in probability) to stay within the allowed region while retaining maximal similarity with the original model?

In what follows, we provide an answer to these questions, and show that the proposed approach adds
 minimal computational overhead to the language generation process.

- 144 3.2 App
- 145

3.2 Approach

146 We design an online method that guarantees, at each token generation, that the sequence remains safe. Given the sequential, feedforward nature of LM layers, we consider each new token generation 147 to realize a trajectory through the layers' activation spaces. Similar to Park et al. (2023), we take 148 the view that disallowed language occupies a region of each layer's activation space \mathbb{R}^d . Formally, 149 let $\mathcal{R}_t \subset \mathbb{R}^d$ be the forbidden, or unsafe, region for layer t of the LM. Our goal, then, is to provide 150 a control mechanism by altering the output vector embedding at every layer that guarantees (in 151 probability) that latent trajectories remain out of \mathcal{R}_t for all t, see Figure 1. For each token generation 152 pass, this control mechanism is to be applied at each layer. This control intervention is designed 153 online, and it depends on the prompt sequence.

154

Semantic Probe We first identify the disallowed region for the generated token, given context. To do so, we feed a set of sequences to the model, and use a lightweight probe to map each latent state $x_t \in \mathbb{R}^d$ to a probability that the sequence is toxic. Specifically, we rely on a *probing classifier function* f_t that maps the latent space \mathbb{R}^d to the decision space [0, 1]. For simplicity, we take f_t to be a logistic regression classifier realized by a linear probe (Hewitt & Manning, 2019). Formally, $f_t : \mathbb{R}^d \to [0, 1]; x_t \mapsto \sigma(W_t^\top x_t)$, where $W_t \in \mathbb{R}^{d \times 2}$ and $\sigma : \mathbb{R}^2 \to \mathbb{R}$ is the sigmoid. For each layer, we define the disallowed region \mathcal{R}_t to be the pre-image of an *unsafe* classification under f_t , using a predefined probability threshold p. That is, $\mathcal{R}_t := \{x \mid \sigma(W_t^\top x) \ge p\}$, where $p \in [0, 1]$.



Figure 1: A. LiSeCo is based on linearly adding vector θ_t to the output of each layer (t = 1, ..., T). Each vector θ_t is the solution of a constrained optimization problem. B. A probing classifier fmapping the latent space \mathbb{R}^d to the decision space [0, 1] is trained to characterize the allowable region (in blue) to which each latent state \tilde{x}_t is constrained. Keeping trajectories out of the toxic region in latent space is equivalent to keeping their image out of the toxic region in decision space. C. At inference time, the state in latent space $(x_t \in \mathbb{R}^d)$ is mapped via learned weights W_t into $W_t^T x_t \in \mathbb{R}^2$. If it falls within the forbidden region (pink), an intervention ($\theta_t \in \mathbb{R}^d$) is computed so that the updated state $\tilde{x}_t = x_t + \theta_t \in \mathbb{R}^d$ lies in the safe region (blue).

Optimal Control Once the forbidden region \mathcal{R}_t is identified, we design a control strategy that, for all layers t, guarantees the latent state x_t remains in the allowed region and retains maximal similarity with the original model. To do this, we design an optimal controller that generates an input $\theta_t \in \mathbb{R}^d$ at every layer t. Mathematically, we solve an optimization problem over θ_t where the pre-computed classifier enters as a hard constraint in the formulation, i.e., $\sigma(W_t^{\top}(\theta_t + x_t)) \leq p$. This ensures the controlled latent trajectory $\tilde{x}_t = x_t + \theta_t \in \mathbb{R}^d$ lies in the unsafe region with probability less than p.

4 OPTIMAL LANGUAGE GENERATION CONTROLLER IN LATENT SPACE

In this section, we describe the theoretical contribution of this work. Using the probing classifier, we design a controller to restrict text generation to the safe region. The optimal intervention is derived in closed form, thus computationally efficient at inference-time.

4.1 Optimal Controller Setup

mi

S

 θ_1,\ldots

The optimal controller aims to keep latent trajectories out of the unsafe region without compromising text quality. That is, we perform constrained optimization where latent trajectories maximally approximate the original ones (proxying text quality) while avoiding the unsafe region as defined by the probe. This gives rise to the following optimization problem:

$$\inf_{\theta_T} \sum_{t=1}^T \|\theta_t\|_2^2$$
(1a)

$$t. \qquad \sigma(W_t^\top(x_t + \theta_t)) - p \le 0, \quad \forall t = 1, \dots, T$$
(1b)

$$x_{t+1} = \operatorname{layer}_t(x_t + \theta_t), \tag{1c}$$

$$x_0 = E(\text{prompt sequence}), \qquad (1d)$$

where *E* is the embedding map. Optimization problem 1 aims to find the minimum l_2 -norm intervention¹ $\theta_1, \ldots, \theta_T$ (Eq. 1a) that satisfies the following constraints: Eq. 1b requires the modified

¹The choice of L_2 norm is standard in classical optimal control problems. The L_2 norm is usually interpreted as energy, or effort, of the control input to steer the system. In this context, it can be seeing as trying to minimize

216 latent state $x_t + \theta_t$ be classified unsafe by the probe f_t ; Eq. 1c captures LM dynamics, i.e., layer t 217 maps the modified latent state $x_t + \theta_t$ to the next latent state x_{t+1} ; Eq. 1d states that the LM's input 218 embeds the input context, so that interventions are *context-dependent*. The intervention that solves 219 optimization problem 1 is guaranteed by construction to keep the latent trajectory $\tilde{x}_1, \ldots, \tilde{x}_T$ and 220 output y below the probability threshold from the classifier.

221 Whether concept avoidance, e.g., detoxification, is expressed as a cost or a constraint depends on 222 the use-case. Other approaches, in contrast to ours, encode concept avoidance in the optimization 223 objective, but not via hard constraints (Dathathri et al., 2019; Hernandez et al., 2023). Fortunately, 224 the constrained optimization framework of LiSeCo also permits this interpretation; though we leave 225 its testing to future work, we state its equivalent problem and prove its *closed-form optimal solution*, 226 which has only been empirically approximated by hyperparameter search in the literature (Li et al., 2023a), in Appendix D. 227

228 229

230

240

241 242

4.2 OPTIMAL CONTROLLER DESIGN

Optimization problem 1 is a standard problem in the optimal control literature (Kirk, 2004). By 231 Bellman's Optimality Principle, the standard approach to solving problem 1 is dynamic programming 232 (DP) (Kirk, 2004): the optimal solution is computed for the last layer T, then via backward induction 233 for $T-1, \ldots, 1$. But, layer dynamics 1c are highly non-convex, and solutions incomputable in closed 234 form, hence their optimality is not guaranteed. Further, DP requires backpropagating gradients at 235 each text generation's forward pass, adding significant inference latency. 236

To overcome these limitations, we relax problem 1. No longer searching for a globally optimal 237 solution across layers, we now search for locally optimal solutions at each layer. Now, Eqs. 1c and 1d 238 cease to play a role, as each layer is optimized for separately. Then, problem 1 is relaxed into 239

$$\min_{\theta_t} \quad \|\theta_t\|_2^2 \tag{2a}$$

s.t.
$$\sigma(W_t^{\top}(x_t + \theta_t)) - p \le 0,$$
 (2b)

243 for each layer $t = 1 \cdots T$. The sequence of θ_t that solve problem 2 may not optimize the original 244 formulation 1. However, one is not anyway guaranteed to find global optima anyway due to the high 245 nonconvexity of layer computations. Furthermore, optimality is not essential as the cost aims only to 246 preserve similarity with the original model. Meanwhile, the guarantee to avoid unsafe region \mathcal{R}_t is 247 still enforced via Eq. 2b. 248

A key advantage of relaxed formulation 2 is that it is solvable in closed-form, per-layer, with minimal 249 computational overhead. The following theorem states the analytical solution for optimal θ_t . 250

Theorem 1 (Optimal θ). The optimal solution $\theta_t^* \in \mathbb{R}^d$ to the optimization problem 2 is given by

$$\theta_t^* = \begin{cases} \frac{\log\left(\frac{1}{p} - 1\right) - w_t^\top x_t}{\|w_t\|_2^2} w_t & \text{if } \sigma(W_t^\top x_t) > p \\ 0 & \text{otherwise.} \end{cases}$$
(3a)

otherwise.

(3b)

254 255 256

257

258 259

260 261

262

263

264

265

266

251 252 253

where $w_t := W_t^1 - W_t^2$, the difference of the columns of $W_t =: \begin{bmatrix} W_t^1 & W_t^2 \end{bmatrix}$.

Proof. Proof relies on leveraging the KKT conditions. See Appendix C for details.

Geometrically, the optimal solution is the vector from x_t to the closest point in \mathcal{R}_t^C , which is the set-complement of \mathcal{R}_t . When $x_t \notin \mathcal{R}_t$, i.e., when $\mathbb{P}(\text{unsafe}) < p^2$, no update is needed; hence $\theta_t^* = 0$. Otherwise, the update is a factor of w_t . Since θ_t^* exists in closed-form, computing an intervention incurs negligible computational overhead. Crucially, it is guaranteed with probability pto keep the latent state outside the disallowed region.

267 the "effort" of the intervention. Moreover, the fact that L_2 is also used to measure distances in an Euclidean 268 space, like the embedding space that we consider in this work, makes it an appropriate choice to measure the "similarity" between the intervened representation and the original one. 269

² \mathbb{P} denotes probability.

270 Although control occurs *locally* at each layer, the local control steps result in a globally safe distri-271 bution over the next token. To see this, consider a single token generation. Each sequential control 272 action at layer t guarantees that latent state $\tilde{x}_t \in \mathcal{R}_t^C$ is classified safe, or equivalently, eliminates the 273 set of unsafe trajectories falling in \mathcal{R}_t . By the time we reach the last layer T, the latent trajectory is 274 guaranteed (in-probability) to have been rated safe at every preceding layer. Then, the last layer T's activation is transformed via the unembedding matrix (linear map) and softmax (monotonic map) to 275 the distribution over the vocabulary. Linearity and monotonicity ensure that the set of safe last-layer 276 activations maps onto the set of safe vocabulary distributions. 277

278 279

280

286

292

5 Methods

The LiSeCo pipeline is as follows. First, to find the unsafe regions and probes per layer, there is an initial probe training phase. Then, probes are integrated into the model at inference-time, and the optimal intervention dynamically applied. We tested LiSeCo on two separate tasks: toxicity avoidance and negativity avoidance. For brevity, when appropriate, we will use *toxicity* to stand in for both toxicity and negativity in this section.

Models We test on three state-of-the-art causal language models: Llama-3-8B (Meta, 2024), Pythia6.9B (Biderman et al., 2023), and Mistral-7B (Jiang et al., 2023). While the architectural details of a
layer (attention + MLP) differ slightly between models, our intervention treats layers as black boxes
and operates at the level of the *residual stream* (Elhage et al., 2021). This permits our intervention to
be applied as a lightweight layer wrapper and in an architecture-agnostic way.

Datasets We test our method on the toxicity and negativity use cases. Borrowing terminology from Ashok & Poczos (2024), we first learn probing classifiers f using a labelled *constraint dataset*, then, we evaluate text generation on a *task dataset*.

For **toxicity**, we use Kaggle's Jigsaw dataset (Adams et al., 2017) as the constraint dataset. The dataset contains 30k label-balanced natural language comments and their human-annotated binary toxicity labels in {toxic, nontoxic}. For the task dataset, we use RealToxicityPrompts, a dataset derived from OpenWebText, a large-scale corpus of the web (Gehman et al., 2020). RealToxicityPrompts is a collection of prompts, their continuations, and toxicity scores in [0, 1] (Gokaslan & Cohen, 2019). To form the task dataset, we sample 150 neutral prompts for which there is a toxic continuation and 150 for which there is a non-toxic continuation in the original dataset.

For **sentiment**, because sentiment datasets tend to be highly domain-specific (for instance, movie reviews), we combine several datasets to form the constraint dataset (N = 30k). This consists of +/label-balanced samples of 7500 datapoints each from IMDb film reviews (Maas et al., 2011), Tweets (Barbieri et al., 2020), Yelp reviews (Zhang et al., 2015), and Amazon reviews (Hou et al., 2024). For preprocessing details, see Appendix G. For the task dataset, we sample 300 neutral sentiment prompts from Liu et al. (2021), created from OpenWebText as a sentiment counterpart to RealToxicityPrompts. Of these prompts, 150 have negative and 150 neutral or positive continuations, respectively.

- 310 **Probing classifiers** Our theoretical guarantees rely on a key assumption: that at each layer t, 311 there indeed exists a \mathcal{R}_t separable by linear f_t which together capture a semantics of the text being 312 generated. We first verify, then, across the panel of LMs that it is possible to linearly decode whether 313 text is toxic from each layer of the LM. Towards this aim, we split each of the constraint datasets 314 into an 80% training set and 20% validation set. Then, for each model, dataset, and layer, we extract 315 the last token hidden representations $x_t \in \mathbb{R}^d$ for each training sequence; we choose the last token 316 embedding to represent the entire sequence, as in causal LMs, it is the only to attend to the entire 317 input sequence. We proceed to train one binary classifier f_t per-layer using the cross-entropy loss between the probe prediction and ground-truth label. See Appendix H for implementation details. 318
- 319

320 5.1 TEXT GENERATION EXPERIMENTS 321

For each LM, we insert trained probes f_t at each layer to evaluate the layer-wise toxicity likelihood at each forward pass. If layer t's representation x_t is evaluated toxic, then the control input θ_t is dynamically applied. For simplicity, we fix text generation to max 50 new tokens, greedily decoded. Baselines We compare our method against several baselines: no-control, instruction-tuning where applicable (Llama and Mistral), Activation Addition (ActAdd) Turner et al. (2023), and Future Discriminators for Generation (FUDGE) (Yang & Klein, 2021).

We consider instruction-tuning, which relies on extensive LM finetuning, to be a target, or "upperbound", baseline. For models with instruction-tuned variants (Llama and Mistral), we repeat the experimental procedure, training probes on the constraint dataset. Then, during evaluation, we prompt the instruction-tuned model using a template whose instructions are slightly modified from Mistral's system prompt provided in Jiang et al. (2023) (see Appendix I for details).

Like LiSeCo, ActAdd steers text generation in activation space (Turner et al., 2023). For each model, the steering vector is computed as follows: (1) a source and target prompt, e.g., ("hate"→"love"), are each fed through the model and activations collected; (2) for each layer, the steering variable is computed as the difference from source to target activation; (3) at inference time, the steering variable is added to the intermediate representations of the input data. Like LiSeCo, ActAdd is gradient-free at inference-time. But, there are key differences: since steers derive from natural language prompts, ActAdd does not require a supervised learning phase on annotated data as in LiSeCo. For the same reason, the method lacks guarantees. For implementation details, see Appendix J.

340 Lastly, we test against Future Discriminators for Generation (FUDGE) (Yang & Klein, 2021) as a 341 representative for steering methods that intervene on the final vocabulary logits (Liu et al., 2021; 342 Schick et al., 2021; Cao et al., 2023). FUDGE requires access to an automatic toxicity scorer which 343 returns a probability that the context + generated token are toxic. Then, at each token generation, the 344 top k = 50 likely tokens, where k is a hyperparameter, are scored for toxicity, and their probabilities 345 reweighted using Bayes' Rule, to minimize the probability of generating a toxic sequence. Because 346 k sequences are passed to the automatic scorer every time a new token is generated, FUDGE has 347 the potential to incur a high inference latency compared to other methods; however, this comes at the benefit of *directly optimizing for the evaluator*. For this reason, we consider FUDGE to be an 348 upper-bound baseline, given that other steering methods do not have access to the ground-truth scorer. 349

Evaluation We evaluate on toxicity (negativity) avoidance and text quality. Toxicity is evaluated automatically, while text quality is rated on a Likert scale by the authors in a blind setup (Appendix N).

353 Semantic control. We rate text generation toxicity using automatic scorers that produce what we call external scores. The automatic scorers are a state-of-the-art RoBERTa-based classifier trained on 354 Twitter data (Camacho-collados et al., 2022; Barbieri et al., 2020) (see Appendix B for model details). 355 We take the scorer's ratings as ground-truth labels, where sequences are labelled toxic if the classifier 356 returns a likelihood higher than 0.5, and non-toxic otherwise. Post-scoring, we restrict evaluation to a 357 randomly sampled, label-balanced set as follows: (1) we randomly choose N prompts for which the 358 no-control baseline has toxic continuations, and N for which there are nontoxic continuations; (2) we 359 evaluate all methods on these prompts. 360

The trained linear probes also provide toxicity likelihoods for the generated text, which we use to post-hoc validate our method, but not to evaluate generation toxicity/negativity per-se. To do so, we weight each layer's linear probe f_t by its validation accuracy $0 \le \alpha_t \le 1$, and compute the *probe score* $S := \sum_{t=1}^{T} f_t(x_t)\sigma_t(\alpha)$ where $f_t(x_t)$ is the toxicity probability assigned to representation x_t and $\sigma_t(\alpha)$ is the sigmoid-weighting of layer validation accuracies. The probe score may be thought of as the probability that an input sequence is toxic as determined by the probes' learned toxicity semantics; if probes represent toxicity in a generalizable way, then the LiSeCo threshold probability *p* should track the fraction of toxic labels by the external classifier.

Text naturalness. The applied intervention ideally should not compromise text quality. We therefore computed the per-token perplexity (PPL) of generations and rated their naturalness on a Likert scale from 1 to 5 in a blind setup. For precise instructions given to annotators, see Appendix N.

372

6 EXPERIMENTAL RESULTS

373 374

We observe that toxic and negative regions are linearly represented in latent space (Park et al., 2024). We then demonstrate that LiSeCo predictably reduces toxicity (negativity) as a function of *p* while maintaining text naturalness. Overall, LiSeCo performs on-par with instruction tuning for toxicity (negativity) reduction and naturalness without extensive finetuning nor online inference latency. **Probing classifiers** Figure 2 shows, for all models, the linear probe validation accuracies per-layer, averaged across 5 random seeds. Probes attain high accuracies of \sim 90%, confirming the disallowed toxic (negative) regions \mathcal{R}_t are linearly decodable with high probability.



Figure 2: Linear probe validation accuracy for toxicity (left) and sentiment (right) detection. All curves are shown ± 1 SD across 5 random seeds. Tasks converge to reasonable accuracies of > 60% for all models and layers, with mid-layers attaining $\approx 90\%$.

	no-control	LiSeCo $(p = 0.1)$	Instruct	ActAdd	FUDGE
Pythia	0.095(0.001)	0.109(0.002)	N/A	0.090(0.001)	3.51(0.65)
Llama	0.113(0.003)	0.132(0.005)	0.119(0.002)	0.118(0.002)	3.41(0.54)
Mistral	0.157 (0.002)	0.169(0.003)	0.162(0.001)	0.159(0.005)	3.56(0.66)

Table 1: Average per-token inference latency (s) with 1 SD (batch size=1) for would-be toxic continuations. FUDGE has the largest inference latency by roughly 2 orders of magnitude, taking around 3s per token generation. All other baselines had negligible extra cost compared to no-control.

Inference latency Table 1 reports the average inference latency of each baseline for each forward pass. Of the methods tested, FUDGE has the highest latency by several orders of magnitude, around 3 seconds compared to other baselines, which incurred negligible overhead compared to no-control.

407 Situating methods in the safety-naturalness plane A successful intervention satisfies both steering 408 and text-naturalness objectives. We visually summarize the performance of various baselines in 409 Figure 3 by plotting their evaluations, human for naturalness and external for toxicity (negativity), 410 with one SD error, on the toxicity (negativity) -naturalness plane, for would-be unsafe continuations, 411 those where no-control produced unsafe content. We first observe that, as expected, instruction-tuning 412 is a well-performing baseline, achieving both high naturalness and low toxicity (negativity). On the 413 other hand, ActAdd (the best hyperparameter setting is shown) either results in poor naturalness or 414 the least improvement to toxicity (negativity). As expected, FUDGE, which directly optimizes for the automatic toxicity (negativity) evaluator, performs best at toxicity (negativity) reduction, however 415 at the cost of naturalness in several cases, and at an extreme latency cost of 3 seconds per token. The 416 best setting of **LiSeCo maintains a high degree of naturalness**,³ where naturalness correlates to p, 417 while detoxifying text generation. The relationship between naturalness and p, which is by theoretical 418 construction (Theorem 1), is especially visible in Mistral (Figure 3 center). See Appendix L for the 419 toxicity (negativity)-naturalness plane plotted for all prompts, and Appendix K for further discussion 420 of naturalness, including qualitative analysis of examples. 421

We note that, while we considered instruction-tuning an upper-bound baseline for performance, 422 LiSeCo is competitive with or outperforms instruction tuning with much less annotated data. To 423 further demonstrate this, we include a stress test of linear probing on the toxicity task in Appendix H, 424 showing that severe decimations of the training set minorly impact probing accuracy: linear probes still 425 achieve > 80% accuracy on the 6000-datapoint test set with only 250 training datapoints. Moreover, 426 LiSeCo is a pure steering method applicable to any frozen LM in a post-hoc adaptation step; in 427 contrast, state-of-the-art instruction tuning requires modifying LM weights and can be extremely 428 energy-intensive, needing > 4 orders of magnitude more data (AI@Meta, 2024). Consequently, 429 effective instruction-tuned LMs only exist for several well-resourced languages. LiSeCo is instead

430 431

8

397 398 399

382

384

385

386

387

388

389 390

403 404

³Human naturalness ratings did not correlate to perplexity due to low-perplexity, degenerate outputs (Appendix K), so we do not attempt to analyze the latter. We leave automated text evaluation to future work.

toxicity-naturalness plane llama mistral pythia 5 naturalness no-control instruct LiSeCo (0.01) LiSeCo (0.1) LiSeCo (0.3) actadd fudge 1 0.6 0.4 0.2 0.6 0.4 0.2 0.6 0.2 0.4 toxicity toxicity toxicity sentiment-naturalness plane llama mistral pythia 5 naturalness o-control instruct LiSeCo (0.01) 3 LiSeCo (0.1) LiSeCo (0.3) actadd fudge 1 0.8 0.6 0.4 0.2 0.6 0.4 0.2 0.6 0.8 0.8 0.4 0.2 sentiment sentiment sentiment

Figure 3: The toxicity-naturalness plane (top) and sentiment-naturalness plane (bottom) for 453 Llama, Mistral, and Pythia (left to right). The top-right corner (low toxicity, high naturalness) is 454 best. Each method's (toxicity/negativity, naturalness) distribution over *would-be toxic* continuations 455 is shown as an ellipse centered at the mean, whose axes reflect ± 1 SD. The red region is that 456 labelled toxic/negative by the external classifier. LiSeCo (blue colors) shifts right, i.e., reduces 457 toxicity/negativity, from no-control (green) and maintains high naturalness, performing on-par with 458 instruction tuning (light orange). ActAdd (orange) least reduces toxicity/negativity. FUDGE (red), 459 which directly optimizes w.r.t. the external classifier, most reduces toxicity/negativity as expected. 460

461

432

433

434 435

436

437

438 439

440

441

442

443

444 445

446

447

448

449

450

451

452

suited to the low-medium resourced regime, and can be flexibly used for specific tasks for which,unlike toxicity and negativity avoidance, instruction-tuning methods have not been fine-tuned for.

464

480

465 Semantic control Here, we analyze toxicity control results in more detail, specifically alignment
 466 between trained probes and the external scores. Full results for negativity can be found in Appendix M.

467 Figure 4 (top) shows the probe score distribution of would-be toxic continuations, N = 25, 37, 37 for 468 Llama, Mistral, and Pythia, respectively (see Appendix M for the full distributions). Here, the toxicity 469 probe score reduction brought on by interventions is visible in the plots as a leftward shift. Notably, 470 LiSeCo with constraint p works as expected, restricting probe scores to < p. The best ActAdd setting 471 slightly decreased the toxicity likelihood, seen by a leftward expansion of the toxicity probe scores, 472 though the extent of reduction was highly sensitive to the hyperparameter setting and model. For both Instruct models, which performed well at toxicity reduction, the toxicity probe score also decreases 473 from the no-control baseline, which evidences that linear probes are able to capture toxicity semantics. 474 Taken together, toxicity probe results show how theoretical guarantees aid interpretability: while 475 toxicity reduction in ActAdd and Instruct remains opaque, that of LiSeCo interpetably depends on p. 476

Figure 4 (bottom) shows the distribution of external toxicity scores for would-be toxic continuations. All baselines decrease toxicity, though we have seen ActAdd to compromise naturalness. Of-note, when LiSeCo is used with p = 0.01, it performs on-par with instruction-tuning for Llama.

481 Smaller LiSeCo p, fewer toxic generations We have shown that LiSeCo reduces the likelihood
 482 of toxicity as defined by the linear probes. But, how well does this definition correspond to the
 483 true labels? Figure 4 (bottom) shows LiSeCo to predictably decrease the externally scored toxicity
 484 likelihood: as LiSeCo p decreases (row 5→3 of the plots), so does the percentage of toxic-labeled
 485 generations (right-hand side). Note, however, that, besides Mistral, the value of p does not upper bound the percentage of toxic generations as it theoretically should: this indicates that in practice,



Figure 4: Toxicity probe scores (top) and external scores (bottom) are shown for Llama, Mistral, and Pythia (left to right), for all baselines (Pythia has no instruct-variant), and LiSeCo for different values of p (0.01, 0.05, etc.). (Bottom) Probability for toxicity greater than 0.5 is shaded in red, with the toxic-labeled % displayed on the right.

514

507

508

509

linear probes only approximately learn toxicity semantics and do not perfectly generalize outside of
 training data. For Pythia specifically, probe scores least align with external toxic label percentages.

Better probes, better performance To test our hypothesis that probe-to-external classifier alignment determines success in practice, we computed the Spearman correlations between the probe scores and the external scores for each model, across the no-control and LiSeCo runs. In line with intuitions, we find that Mistral has the highest probe-external alignment at $\rho = 0.38^{***}$, followed by Llama at $\rho = 0.20^{***}$, and finally Pythia at $\rho = 0.06$ (not significant).⁴

520 521

522

7 DISCUSSION

523 We have proposed LiSeCo, a controlled language generation method that is theoretically guaranteed 524 to stay within permitted regions of latent space. Empirically, the method produces non-toxic and 525 non-negative, but still natural, text. In line with theoretical guarantees, the parameter p was shown 526 to empirically control the probability of generating unwanted text. LiSeCo is compatible with all 527 current Transformer-based architectures, and involves a negligible inference-time latency. In future 528 work, we are interested in applying our approach to different tasks and joint constraints, as well as to 529 alternatives to linear probes as the way to ascertain whether a token falls into the undesirable region.

With the increasing ubiquity of LMs comes a growing need to understand their behavior. LiSeCo 530 helps address this need by providing practical and theoretical tools for LM interpretability and control. 531 That said, using LiSeCo has several caveats: (1) it requires supervised learning of the linear probes 532 on annotated data; (2) the intervention is only as good as the probes, which is only as good as their 533 training data. Thus, when training probes, it is crucial that the training data well-represent the use 534 domain. We emphasize that this bottleneck is inherent to any method that learns from data, and it 535 exists for all steering methods that rely on a discriminator, e.g., FUDGE, Plug and Play, ITI. Another 536 limitation of our method is that we assume disallowed semantics to be linearly encoded in latent 537 space, and that learned probe semantics highly align to true semantics; this should always be verified 538 in practice.

 $^{4^{(***)}}$ means significant at $\alpha = 1e - 3$

540 Ethics statement Controlling text generation can be used for benefit or harm. While we have 541 demonstrated our method on toxicity and negativity avoidance, it can equivalently be applied to 542 increase harmful traits. However, the When designing the linear probes, it is essential to choose a 543 constraint set that accurately reflects the use-case.

Reproducibility statement Code and data are uploaded as a zip file, and will be made public upon acceptance. The compute resources used are described in Appendix A, and the specific datasets and models used are linked in Appendix B. The proof of Theorem 1 is detailed in Appendix C.

References

544

546

547 548 549

550 551

554

555

556

561

569

570

571

572

573

- CJ Adams, Jeffrey Sorensen, Julia Elliott, Lucas Dixon, Mark McDonald, Nithum, and Will Cukierski. Toxic comment classification challenge, 2017. URL https://kaggle.com/ 552 competitions/jigsaw-toxic-comment-classification-challenge. 553
 - AI@Meta. Llama 3 model card. 2024. URL https://github.com/meta-llama/llama3/ blob/main/MODEL CARD.md.
- Dhananjay Ashok and Barnabas Poczos. Controllable text generation in the instruction-tuning era. 558 (arXiv:2405.01490), May 2024. doi: 10.48550/arXiv.2405.01490. URL http://arxiv.org/ 559 abs/2405.01490. arXiv:2405.01490 [cs].
- Nikolay Babakov, David Dale, Ilya Gusev, Irina Krotova, and Alexander Panchenko. Don't lose the message while paraphrasing: A study on content preserving style transfer. In Elisabeth 562 Métais, Farid Meziane, Vijayan Sugumaran, Warren Manning, and Stephan Reiff-Marganiec (eds.), 563 Natural Language Processing and Information Systems, pp. 47–61, Cham, 2023. Springer Nature 564 Switzerland. ISBN 978-3-031-35320-8. 565
- 566 Francesco Barbieri, Jose Camacho-Collados, Luis Espinosa-Anke, and Leonardo Neves. TweetEval:Unified Benchmark and Comparative Evaluation for Tweet Classification. In Proceedings of 567 Findings of EMNLP, 2020. 568
 - Nora Belrose, David Schneider-Joseph, Shauli Ravfogel, Ryan Cotterell, Edward Raff, and Stella Biderman. LEACE: Perfect linear concept erasure in closed form. In Thirty-seventh Conference on Neural Information Processing Systems, 2023. URL https://openreview.net/forum? id=awIpKpwTwF.
- 574 Aman Bhargava, Cameron Witkowski, Manav Shah, and Matt Thomson. What's the magic word? a 575 control theory of llm prompting. arXiv preprint arXiv:2310.04444, 2023.
- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, 577 Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, 578 Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models 579 across training and scaling. (arXiv:2304.01373), Apr 2023. URL http://arxiv.org/abs/ 580 2304.01373. arXiv:2304.01373 [cs]. 581
- 582 Carrie Cai, Tongshuang Wu, and Michael Andrew Terry. Transparent and controllable human-ai interaction via chaining of machine-learned language models, April 13 2023. US Patent App. 583 17/957,526. 584
- 585 Jose Camacho-collados, Kiamehr Rezaee, Talayeh Riahi, Asahi Ushio, Daniel Loureiro, Dimosthenis 586 Antypas, Joanne Boisson, Luis Espinosa Anke, Fangyu Liu, Eugenio Martinez Camara, et al. TweetNLP: Cutting-edge natural language processing for social media. In *Proceedings of the* 588 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, 589 pp. 38–49, Abu Dhabi, UAE, December 2022. Association for Computational Linguistics. URL 590 https://aclanthology.org/2022.emnlp-demos.5.
- Meng Cao, Mehdi Fatemi, Jackie CK Cheung, and Samira Shabanian. Systematic rectification of 592 language models via dead-end analysis. In The Eleventh International Conference on Learning Representations, 2023. URL https://openreview.net/forum?id=k8_yVW3Wqln.

- 594 Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons 595 in pretrained transformers. In Proceedings of the 60th Annual Meeting of the Association for 596 Computational Linguistics (Volume 1: Long Papers), pp. 8493–8502, 2022. 597 Sumanth Dathathri, Andrea Madotto, Janice Lan, Jane Hung, Eric Frank, Piero Molino, Jason 598 Yosinski, and Rosanne Liu. Plug and play language models: A simple approach to controlled text generation. In International Conference on Learning Representations, 2019. 600 601 Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in language models. In 602 Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp. 6491-6506, 2021. 603 604 Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda 605 Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, 606 Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal 607 Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris 608 Olah. A mathematical framework for transformer circuits. Transformer Circuits Thread, 2021. 609 https://transformer-circuits.pub/2021/framework/index.html. 610 Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealToxic-611 ityPrompts: Evaluating neural toxic degeneration in language models. In Trevor Cohn, Yulan 612 He, and Yang Liu (eds.), Findings of the Association for Computational Linguistics: EMNLP 613 2020, pp. 3356–3369, Online, November 2020. Association for Computational Linguistics. 614 doi: 10.18653/v1/2020.findings-emnlp.301. URL https://aclanthology.org/2020. 615 findings-emnlp.301. 616 Aaron Gokaslan and Vanya Cohen. Openwebtext corpus. http://Skylion007.github.io/ 617 OpenWebTextCorpus, 2019. 618 619 Evan Hernandez, Belinda Z Li, and Jacob Andreas. Inspecting and editing knowledge representations 620 in language models. arXiv preprint arXiv:2304.00740, 2023. 621 John Hewitt and Christopher D. Manning. A structural probe for finding syntax in word repre-622 sentations. In Jill Burstein, Christy Doran, and Thamar Solorio (eds.), Proceedings of the 2019 623 Conference of the North American Chapter of the Association for Computational Linguistics: 624 Human Language Technologies, Volume I (Long and Short Papers), pp. 4129–4138, Minneapolis, 625 Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1419. 626 URL https://aclanthology.org/N19-1419. 627 Yupeng Hou, Jiacheng Li, Zhankui He, An Yan, Xiusi Chen, and Julian McAuley. Bridging language 628 and items for retrieval and recommendation. arXiv preprint arXiv:2403.03952, 2024. 629 630 Edward Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, 631 and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In Proceedings of 632 ICLR, Online, 2022. Published online: https://openreview.net/group?id=ICLR. 633 cc/2022/Conference. 634 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, 635 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, 636 Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas 637 Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. 638 639 Donald E Kirk. Optimal control theory: an introduction. Courier Corporation, 2004. 640 Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik 641 Opitz, and Tobias Hecking. Style vectors for steering generative large language model. arXiv 642 preprint arXiv:2402.01618, 2024. 643 644 Ben Krause, Akhilesh Deepak Gotmare, Bryan McCann, Nitish Shirish Keskar, Shafiq Joty, Richard 645 Socher, and Nazneen Fatema Rajani. GeDi: Generative discriminator guided sequence generation. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih 646
- (eds.), Findings of the Association for Computational Linguistics: EMNLP 2021, pp. 4929–4952,
 Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.

648 649	<pre>doi: 10.18653/v1/2021.findings-emnlp.424. URL https://aclanthology.org/2021. findings-emnlp.424.</pre>
651 652 653	Sachin Kumar, Eric Malmi, Aliaksei Severyn, and Yulia Tsvetkov. Controlled text generation as continuous optimization with multiple constraints. <i>Advances in Neural Information Processing Systems</i> , 34:14542–14554, 2021.
654 655	Shibamouli Lahiri. SQUINKY! A corpus of sentence-level formality, informativeness, and implica- ture. arXiv preprint arXiv:1506.02306, 2015.
656 657 658 659	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> , 2023a. URL https://openreview.net/forum? id=alLuYpn83y.
660 661 662	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. <i>Advances in Neural Information Processing Systems</i> , 36, 2024a.
664 665 666	Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-Im improves controllable text generation. <i>Advances in Neural Information Processing Systems</i> , 35: 4328–4343, 2022.
667 668 669 670	Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In <i>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)</i> , pp. 4582–4597, 2021.
671 672 673	Xiaopeng Li, Shasha Li, Shezheng Song, Jing Yang, Jun Ma, and Jie Yu. Pmet: Precise model editing in a transformer. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 18564–18572, 2024b.
675 676 677	Zhoubo Li, Ningyu Zhang, Yunzhi Yao, Mengru Wang, Xi Chen, and Huajun Chen. Unveiling the pitfalls of knowledge editing for large language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2023b.
678 679 680 681 682 683 684	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.), <i>Proceedings of the 59th</i> <i>Annual Meeting of the Association for Computational Linguistics and the 11th International Joint</i> <i>Conference on Natural Language Processing (Volume 1: Long Papers)</i> , 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.
685 686 687	Ximing Lu, Sean Welleck, Peter West, Liwei Jiang, Jungo Kasai, Daniel Khashabi, Ronan Le Bras, Lianhui Qin, Youngjae Yu, Rowan Zellers, et al. Neurologic a* esque decoding: Constrained text generation with lookahead heuristics. <i>arXiv preprint arXiv:2112.08726</i> , 2021.
688 689	Yifan Luo, Yiming Tang, Chengfeng Shen, Zhennan Zhou, and Bin Dong. Prompt engineering through the lens of optimal control. <i>arXiv preprint arXiv:2310.14201</i> , 2023.
690 691 692 693 694	Andrew L. Maas, Raymond E. Daly, Peter T. Pham, Dan Huang, Andrew Y. Ng, and Christopher Potts. Learning word vectors for sentiment analysis. In <i>Proceedings of the 49th Annual Meeting</i> of the Association for Computational Linguistics: Human Language Technologies, pp. 142–150, Portland, Oregon, USA, June 2011. Association for Computational Linguistics. URL http: //www.aclweb.org/anthology/P11-1015.
695 696 697	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. Advances in Neural Information Processing Systems, 35:17359–17372, 2022a.
698 699 700 701	Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in gpt. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), <i>Advances in Neural Information Processing Systems</i> , volume 35, pp. 17359–17372. Curran Associates, Inc., 2022b. URL https://proceedings.neurips.cc/paper_files/paper/

730

731

732

747

748

749

- Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In *The Eleventh International Conference on Learning Representations*, 2022c.
- Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024. URL
 https://ai.meta.com/blog/meta-llama-3/.
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christopher D Manning. Fast model
 editing at scale. In *International Conference on Learning Representations*, 2021.
- 711 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong 712 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser 713 Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan 714 Leike, and Ryan Lowe. Training language models to follow instructions with human feed-715 back. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Ad-716 vances in Neural Information Processing Systems, volume 35, pp. 27730-27744. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/ 717 2022/file/blefde53be364a73914f58805a001731-Paper-Conference.pdf. 718
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. In *Causal Representation Learning Workshop at NeurIPS 2023*, 2023. URL https://openreview.net/forum?id=T0PoOJg8cK.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The linear representation hypothesis and the geometry of large language models. 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. 39643–39666. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/park24c.html.
 - Ellie Pavlick and Joel Tetreault. An empirical analysis of formality in online communication. *Transactions of the Association for Computational Linguistics*, 4:61–74, 2016. doi: 10.1162/tacl_a_00083. URL https://aclanthology.org/Q16-1005.
- Lianhui Qin, Sean Welleck, Daniel Khashabi, and Yejin Choi. Cold decoding: Energy-based constrained text generation with langevin dynamics. *Advances in Neural Information Processing Systems*, 35:9538–9551, 2022.
- Sudha Rao and Joel Tetreault. Dear sir or madam, may I introduce the GYAFC dataset: Corpus, benchmarks and metrics for formality style transfer. 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. 129–140, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1012. URL https://aclanthology.org/N18-1012.
- 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.
 - Stefano Soatto, Paulo Tabuada, Pratik Chaudhari, and Tian Yu Liu. Taming AI bots: Controllability of neural states in large language models. https://arxiv.org/abs/2305.18449, 2023.
- Nishant Subramani, Nivedita Suresh, and Matthew E Peters. Extracting latent steering vectors from pretrained language models. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 566–581, 2022.
- Alex Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*, 2023.

Lucas Weber, Elia Bruni, and Dieuwke Hupkes. Mind the instructions: a holistic evaluation of consistency and interactions in prompt-based learning. In Jing Jiang, David Reitter, and Shumin Deng (eds.), Proceedings of the 27th Conference on Computational Natural Language Learning (CoNLL), pp. 294–313, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.conll-1.20. URL https://aclanthology.org/2023.conll-1.20. Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022. 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. Yunzhi Yao, Peng Wang, Bozhong Tian, Siyuan Cheng, Zhoubo Li, Shumin Deng, Huajun Chen, and Ningyu Zhang. Editing large language models: Problems, methods, and opportunities. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pp. 10222-10240, 2023. Jingying Zeng, Richard Huang, Waleed Malik, Langxuan Yin, Bojan Babic, Danny Shacham, Xiao Yan, Jaewon Yang, and Qi He. Large language models for social networks: Applications, challenges, and solutions. arXiv preprint arXiv:2401.02575, 2024. Xiang Zhang, Junbo Zhao, and Yann LeCun. Character-level convolutional networks for text classification. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 28. Curran Associates, Inc., 2015. URL https://proceedings.neurips.cc/paper files/paper/2015/ file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf.

810 A COMPUTING RESOURCES

Experiments were run on a cluster with 12 nodes with 5 NVIDIA A30 GPUs and 48 CPUs each.

Extracting LM representations took a few wall-clock hours per model-dataset computation. Training linear probes took around 15 minutes per layer, so overall 64 wall-clock hours. Running evaluation experiments took a total of 30 wall-clock hours.

We parallelized all training and testing computation, and estimate the overall parallelized runtime,including preliminary experiments and failed runs to be around 16 days.

```
B Assets
```

819 820

821

824

836

837

838

839 840

841 842

843 844

845 846 847

848 849

850 851

852 853

856

858

863

822 Llama https://huggingface.co/meta-llama/Meta-Llama-3-8B; license: llama3 823 https://huggingface.co/meta-llama/Meta-Llama-3-8B; license: llama3

Mistral https://huggingface.co/mistralai/Mistral-7B-v0.1; license: apache-2.0

825 826 Pythia https://huggingface.co/EleutherAI/pythia-6.9b; license: apache-2.0

827 **PyTorch** https://scikit-learn.org/; license: bsd

- 828 Toxicity constraint https://huggingface.co/datasets/google/jigsaw_ 829 toxicity_pred; license: CC0
- 830 Sentiment constraint https://huggingface.co/datasets/stanfordnlp/imdb;
 831 license: unknown.
 832 https://huggingface.co/datasets/cardiffnlp/tweet_eval; license:
 833 unknown.
 834 https://huggingface.co/datasets/Yelp/yelp_review_full; license:
 835 yelp-license.

yelp-license.
https://huggingface.co/datasets/McAuley-Lab/
Amazon-Reviews-2023; license: MIT.

Toxicity task https://huggingface.co/datasets/allenai/ real-toxicity-prompts; license: apache-2.0

Sentiment task https://github.com/alisawuffles/DExperts; license: unknown

C PROOF OF THEOREM 1

Theorem 1 (Optimal θ). The optimal solution $\theta_t^* \in \mathbb{R}^d$ to the optimization problem is given by

$$\theta_t^* = \begin{cases} \frac{\log\left(\frac{1}{p}-1\right) - w_t^\top x_t}{\|w_t\|_2^2} w_t & \text{if } \sigma(w_t^\top x_t) > p, \\ 0 & \text{otherwise.} \end{cases}$$

where $w_t := W_t^1 - W_t^2$, the difference of the columns of $W_t =: \begin{bmatrix} W_t^1 & W_t^2 \end{bmatrix}$.

Proof. We start by defining the Lagrangian for optimization problem in Equation (2) as

$$L(\theta_t, \lambda) = \|\theta_t\|_2^2 + \lambda \left(\sigma(W^\top(x_t + \theta_t) - p) \right), \tag{C.4}$$

where $\lambda \in \mathbb{R}$ is the Lagrange multiplier.

1

We now solve Equation (2) by using KKT conditions, which are first-order necessary conditions for optimality:

1. Stationarity.

$$0 \in \partial(\|\theta_t\|_2^2 + \lambda(\sigma(W^T(x_t + \theta_t)) - p))$$
(C.5)

2. Complementary slackness.

$$\lambda(\sigma(W^{\top}(x_t + \theta_t) - p) = 0 \tag{C.6}$$

3. Primal feasibility: $\sigma(W^{\top}(x_t + \theta_t)) - p \le 0$ (C.7) 4. Dual feasibility. $\lambda \geq 0.$ (C.8)

First, consider when $\lambda = 0$. We apply the stationarity condition Equation (C.5) to obtain $\theta_t = 0$. Plugging θ_t back into the primal constraint, we have that $\sigma(W^{\top}x_t) \leq p$ and recover the second line of Equation (3a). That is, when $\lambda = 0$, we are already in the non-toxic region and do not need to apply an intervention θ_t .

Now, consider $\lambda > 0$. When this is the case, $\sigma(W^{\top}x_t) > p$ and an intervention is needed. Here, it is possible again to solve for θ_t in closed form. By complementary slackness Equation (C.6),

$$p = \sigma(W^{\top}(x_t + \theta_t)) = \frac{\exp(w_2^{\top}(x_t + \theta_t))}{\exp(w_1^{\top}(x_t + \theta_t)) + \exp(w_2^{\top}(x_t + \theta_t))}.$$
 (C.9)

Hence,

$$w^{\top}\theta_t + w^{\top}x_t - \log\left(\frac{1}{p} - 1\right) = 0.$$
 (C.10)

Now, when $\lambda > 0$, or when $\sigma(W^{\top}x_t) > p$, Equation (2) is equivalent to minimizing $\|\theta_t\|_2^2$ subject to Equation (C.10). The Lagrangian with respect to this new formulation is

$$L(\theta_t, \lambda') = \|\theta_t\|_2^2 + \lambda' \left(w^\top \theta_t + w^\top x_t - \log\left(\frac{1}{p} - 1\right) \right).$$
(C.11)

Taking the partial derivative with respect to θ_t , we have

$$0 = \frac{\partial L(\theta_t, \lambda')}{\partial \theta_t} = 2\theta_t + \lambda' w.$$
(C.12)

Hence,

$$\theta_t = -\frac{\lambda' w}{2}.\tag{C.13}$$

Now, we plug θ_t back into Equation (C.10) to obtain

$$\lambda' = \frac{2\left(w^{\top}x - \log\left(\frac{1}{p} - 1\right)\right)}{\|w\|_{2}^{2}}.$$
(C.14)

Finally, plugging λ' back into Equation (C.13), we have

$$\theta_t = \frac{\log\left(\frac{1}{p} - 1\right) - w_t^\top x_t}{\|w_t\|_2^2} w_t.$$
 (C.15)

This completes line 1 of Equation (3a).

D NATURALNESS-FIRST FORMULATION

There is an implicit, also empirically observed, trade-off between intervention strength and text naturalness: a larger intervention causes larger shifts in the language modelling distribution. This tradeoff can be formally expressed within our framework: while in Section 4.1 we present naturalness as a cost $(\min_{\theta_t} \|\theta_t\|_2^2)$ and toxicity avoidance as a constraint, one can also do the opposite. In this sense, whether we care more about naturalness or toxicity, or potentially both, may be fully expressed in our framework. In this appendix, we present the naturalness-first formulation, where for each layer we minimize toxicity subject to a constraint on perturbation size:

$$\min_{\theta_t} \quad -\log \sigma(W_t^{\top}(x_t + \theta_t))) \tag{D.16a}$$

s.t.
$$\|\theta_t\|_2^2 - \beta \le 0.$$
 (D.16b)

The cost is the negative log-likelihood of the modified trajectory being toxic; the constraint is that the perturbation L_2 norm be smaller than some constant $\sqrt{\beta}$.

Theorem 2. The optimal θ_t to Equation (D.16) is

$$\theta_t^* = \frac{\sqrt{\beta}}{\|w_t\|_2} w_t,\tag{D.17}$$

922 where $w_t := W_t^1 - W_t^2$.

Proof. By monotonicity of the log, the objective Equation (D.16a) is equivalent to

=

 $\min_{\theta_t} \quad \sigma(W_t^T(x_t + \theta_t)) \tag{D.18}$

$$= \frac{\exp(w_2^{\top}(x_t + \theta_t))}{\exp(w_1^{\top}(x_t + \theta_t)) + \exp(w_2^{\top}(x_t + \theta_t))}$$
(D.19)

$$\equiv \max_{\boldsymbol{\theta}_{t}} \qquad 1 + \exp((w_1 - w_2)^{\top} (x_t + \theta_t)) \tag{D.20}$$

$$\equiv \max_{\theta_t} \quad (w_1 - w_2)^\top \theta_t \tag{D.21}$$

t.
$$\|\theta_t\|_2^2 - \beta \le 0.$$
 (D.22)

The solution θ_t^* is, then the vector in the direction of $w_t := w_1 - w_2$ with norm $\sqrt{\beta}$:

$$\theta_t^* = \frac{\sqrt{\beta}}{\|w_t\|_2} w_t. \tag{D.23}$$

E CONTINUOUS TUNING FORMULATION

s.

Beyond classification into binary regions, e.g., toxic or non-toxic, LiSeCo can also be used for *continuous tuning* of linearly encoded attributes. For instance, if we require that a style attribute such as formality be restricted within a continuous range, or set to a specific value, the appropriate optimization problem can be stated and the optimal solution solved for in closed-form. These closed-form solutions closely follow the format of Theorem 1.

E.1 SETTING AN ATTRIBUTE TO A SPECIFIC VALUE

949 We consider the optimization problem for setting an attribute, such as toxicity or formality, to a 950 specific value. Let there theoretically exist an attribute scorer $\phi : \Sigma^* \to \mathbb{R}$ which assigns natural 951 language strings to continuous ratings β , where $\beta_1 > \beta_2$ means β_1 is more formal than β_2 , and 952 $\beta_{1,2} \in \mathbb{R}$. The scoring function ϕ defines an *ordering* for utterances with respect to some attribute. 953 For instance, ϕ could be estimated by asking humans to rate sentences for formality on a Likert scale.

There necessarily exists a monotone operator ν that transforms ϕ to [0, 1]. Then, to set the continuous attribute rating of an LM generation to $\beta \in \mathbb{R}$, we can equivalently set $f_t(\text{attribute}) \triangleq \sigma(W_t^\top(x_t + \theta_t))$ to $p \triangleq \nu(\beta)$.

958 This yields the following optimization problem:

 $\min_{t \in \mathbb{R}} \|\theta_t\|_2^2 \tag{E.24a}$

s.t.
$$\sigma(W_t^{\top}(x_t + \theta_t)) - p = 0, \qquad (E.24b)$$

962 with optimal θ_t given by the following theorem:

Corollary 1 (Optimal θ , continuous tuning). The optimal solution $\theta_t^* \in \mathbb{R}^d$ to the optimization problem E.24 is given by

$$\theta_t^* = \frac{\log\left(\frac{1}{p} - 1\right) - w_t^\top x_t}{\|w_t\|_2^2} w_t,$$
(E.25)

where $w_t := W_t^1 - W_t^2$, the difference of the columns of $W_t =: \begin{bmatrix} W_t^1 & W_t^2 \end{bmatrix}$.

Proof. The proof is identical to that of Theorem 1, except the primal feasibility constraint is now an equality constraint. \Box

972 E.2 SETTING AN ATTRIBUTE TO A SPECIFIC RANGE 973

977 978 979

980

985

986

996 997

998

999

974 Now, inheriting the same setup from the previous section, we set the attribute to a range in $[\beta_1, \beta_2]$, 975 $\beta_1 < \beta_2 \in \mathbb{R}$. This is equivalent to setting $f_t(\text{attribute}) \in [\nu(\beta_1), \nu(\beta_2)] \triangleq [p_1, p_2]$, yielding the 976 following optimization problem:

$$\min_{\theta_t} \quad \|\theta_t\|_2^2 \tag{E.26a}$$

s.t.
$$\sigma(W_t^{\top}(x_t + \theta_t)) - p_1 \ge 0$$
 (E.26b)

$$\sigma(W_t^{\top}(x_t + \theta_t)) - p_2 \le 0 \tag{E.26c}$$

Again, the closed-form optimal solution can be solved for, given in the following theorem:

Corollary 2 (Optimal θ , range). The optimal solution $\theta_t^* \in \mathbb{R}^d$ to the optimization problem E.26 is given by

$\theta_t^* = \langle$	$\left\{\frac{\log\left(\frac{1}{p_2}-1\right)-w_t^\top x_t}{\ w_t\ _2^2}w_t\right.$	$if\sigma(W_t^\top x_t) > p_2$	
	$\frac{\log\left(\frac{1}{p_{1}}-1\right)-w_{t}^{\top}x_{t}}{\ w_{t}\ _{2}^{2}}w_{t}$	$\textit{if } \sigma(W_t^\top x_t) < p_1$	
	0	otherwise,	

where $w_t := W_t^1 - W_t^2$, the difference of the columns of $W_t =: \begin{bmatrix} W_t^1 & W_t^2 \end{bmatrix}$.

Proof. The proof is nearly identical to that of Theorem 1, applying the KKT conditions to both inequality constraints. \Box

We have shown that the LiSeCo framework can be applied to both setting an attribute to a specific value and within a specific range, given a scoring function ϕ that imposes an order on strings Σ^* . Because ϕ yields continuous ratings $\in \mathbb{R}$ that are ordered, the scores are isomorphic to [0, 1] by means of a continuous monotonic map ν . Then, tuning an attribute to a particular value β is equivalent to performing logistic regression with f_t , such that the output is $p \triangleq \nu(\beta)$. Crucially, ν **permits interpretation of LiSeCo** p **in human rating space**, a feature that does not exist for other tested methods.

The emphasis here is on continuous tuning of an attribute, e.g., formality. However, these extensions to LiSeCo also permit a probabilistic interpretation when appropriate, where p would the probability that a given generation is, for instance, toxic.

¹⁰¹¹ ₁₀₁₂ F CONTINUOUS TUNING EXPERIMENTS

We demonstrate continuous tuning on a text formality use-case. Formality is an aspect of communication style, where, for example, formal text is common in newspapers, articles, and encyclopedic text, and informal text may be more common in text messages, Tweets, or online forums (Pavlick & Tetreault, 2016). Here, we demonstrate LiSeCo continuous tuning (Appendix E) where we set formality of LLM-generated text to specific values. All experiments were on the model Llama-3-8B; we leave further testing to future work.

1020 F.1 DATASET

1022The constraint set for fitting f_t per-layer is the formality dataset of Pavlick & Tetreault, 2016; Lahiri,10232015. This dataset, which we call *formality-scores*, contains 11270 (sentence, average human rating)1024ordered pairs. Human ratings were crowdsourced via Amazon Mechanical Turk, and range from -31025(most casual) to 3 (most formal), inclusive; for details, see (Pavlick & Tetreault, 2016). The human ratings, continuous because they are averaged across subjects, are the output of the scoring function

1026 ϕ as defined in Appendix E. To train the logistic regression classifiers f_t using cross-entropy loss, we normalize, that is,

1029

1030 1031

1032

1033 1034

1035 1036

the human ratings to [0, 1], then interpreted as p in Corollary 1. For the task set, we re-use the sentiment dataset, sampling 500 neutral prompts.

 $\nu(\phi) = \frac{\phi - \min_{\Sigma^*} \phi}{\max_{\Sigma^*} \phi - \min_{\Sigma^*} \phi},$

F.2 LISECO

1037 As $\nu : \mathbb{R} \to [0,1]; \beta \to p$ maps human rating space to the space of LiSeCo p, where ν is monotonic 1038 and continuous by definition, p can be mapped back to rating space via ν^{-1} . This permits an 1039 interpretation of LiSeCo p in human rating space by construction. That is, the \hat{p} estimated by a neural 1040 formality scorer in [0, 1] will have an interpretation in rating space of $\nu^{-1}(p) \in \mathbb{R}$. Using LiSeCo, we set formality values to $\{-3, -1.5, 0, 1.5, 3\}$, implied by $p \in \{0.01, 0.25, 0.5, 0.75, 0.99\}$.

1042

1044

1043 F.3 BASELINES

FUDGE does not permit continuous tuning. Therefore, we evaluate against no-control, Instruct (Mistral and Llama), and ActAdd. For Instruct, there is not a standard way to set formality to a certain value. To proxy this, we use the following prompt, adapted from the template in Appendix I:

1048 Instructions:

```
1049
     Text can have different levels of formality. Newspapers and
     encyclopedias have high formality ratings and text messages may
1050
     have low formality ratings. The maximum formality rating is 3,
1051
     and the minimum formality rating is -3. A rating of 0 is neutral.
1052
     Example of -3: "(LOL) jus kidding... the answer to your question
1053
     is GAS PRICES !!! ". Example of 3: "With the everpresent elderly
1054
     community in South Broward, we are better poised to serve that
1055
     community by having a "community representative" on our board."
1056
     With this in mind, please continue the following text so that it
1057
     has a formality rating {S}.
```

1058

1059 Text:

1060 PROMPT

where PROMPT is replaced with the prompt and $\{S\}$ is replaced with the desired score in [-3, 3]. This is the same range as the formality-scores dataset (Pavlick & Tetreault, 2016), where examples in the instructions are taken directly from formality-scores. We test values of $S \in \{-3, -1.5, 0, 1.5, 3\}$, the same as for LiSeCo.

For ActAdd, we compute the vector (casual \rightarrow formal), and take the intervention strength $c \in \{-1, -0.1, -0.01, 0.01, 0.1, 1\}$ to proxy the extent of formality (higher is more formal). Note that ActAdd intervention strength c cannot be transformed to (thus interpreted in) human rating space, while LiSeCo p and Instruct S can be.

1005

1072

1071 F.4 EVALUATION

External formality scorer We evaluate continuous tuning by the Spearman correlation ρ between the level of formality given by LiSeCo p, ActAdd c, and Instruct S, to the score in [0, 1] from an external neural scorer (Babakov et al., 2023). The external scorer is a RoBERTa-based architecture trained on the GYAFC dataset, thus far the largest-scale human-annotated formality dataset in the literature (Rao & Tetreault, 2018).

- 1078
- **Human naturalness evaluations** In addition, similar to the toxicity and negativity experiments, we collect human naturalness ratings on a Likert scale from 1-5, details in Appendix N.

The validation accuracy is around 80% for all layers.

Figure F.1: Linear probe validation accuracy for formality on Llama, ± 1 SD across 5 random seeds.

F.5 RESULTS

Here, we present preliminary experiments on continuous tuning using LiSeCo and according to other baselines. We find that, for LiSeCo and Instruct, tuning the respective design parameters p and S shift formality in the expected direction and to a narrow range, where that the tradeoff with text quality is less pronounced than for ActAdd. Overall among the methods, LiSeCo is able to best-vary the formality without compromising text quality.

1104 1105

1080

1082

1083

1084 1085

1086

1087

1088

1089 1090

1091

1092 1093

1094

1095 1096

1098

Linear probes Figure F.1 shows the linear probe validation accuracy across all layers on formalityscores for Llama-3-8B. Like toxicity and sentiment, see Figure 2, formality is approximately linearly decodable with a layerwise accuracy of around 80%.

1109

Continuous tuning evaluation Figure F.2 shows the resulting formality of the generated text, as scored by the neural scorer, as a function of the design parameters of LiSeCo p, ActAdd c, and Instruct S (left to right). The neural formality scores are transformed into implied human ratings via ν^{-1} , (Pearson R = 0.80, significant at $\alpha = 1e$ -3).⁵.

1114 As desired, formality is monotonic in LiSeCo p and Instruct S. Anecdotally, higher ActAdd strengths 1115 in both the informal and formal directions degraded text quality, causing the neural scorer to rate the 1116 text as "informal"; hence, the upside-down U shape. However, while shifting formality in the correct direction, both LiSeCo and Instruct do not produce generations which capture the full range of human 1117 ratings $\in [-3,3]$ (Pavlick & Tetreault, 2016). For Instruct, this is due to the model often repeating 1118 the input instructions, while copying its style. We hypothesize that another factor is a saturation effect 1119 at the upper extreme of the human and neural scorers, see Figure F.3; learning a better mapping ν^{-1} 1120 from the neural scorer to the implied human rating is an important direction for future work. 1121

1122

1123 **Text naturalness** Figure F.4 shows the distrbution of text naturalness for methods of various 1124 intervention strength. As expected, the more extreme the intervention, the lower the text naturalness. We find in particular that ActAdd is very sensitive to intervention strength, where increasing |c| from 1125 0.01 to 1 decreases naturalness on average from around 4 to 3. For Instruct, changing the design 1126 parameter S did not impact naturalness; instead, we note that generations followed the style of, and 1127 often repeated content from the input template, which reduced the text naturalness from the no-control 1128 baseline (far right). Finally, LiSeCo, often even for strong interventions, e.g. p = 0.99 but not for 1129 p = 0.01, only minorly impacted naturalness. 1130



⁵We tried different forms of ν , with similar results, including the logit function. We learned via inverse logistic regression on *formality-scores*: $\beta \sim \exp \gamma/(1 + \exp \gamma)$ where β are the human ratings and γ the neural scores in [0, 1] (Pearson R = 0.79, significant at $\alpha = 1e$ -3)



1147 Figure F.2: Continuous tuning results for LiSeCo, ActAdd, and Instruct (left-to-right) for Llama. 1148 The x-axes of the plots, left to right, show the tuning parameters p (LiSeCo), ActAdd intervention 1149 strength c, and S (Instruct), where formality should increase to the right. The implied human rating 1150 (y-axis) is given by computing $\nu^{-1}(\gamma)$, where γ is the external formality score. All tested settings 1151 are shown that achieve naturalness ≥ 3 on average. The implied human formality ratings of LiSeCo 1152 and Instruct are monotonic in their design parameter, as desired; ActAdd is not. However, while 1153 monotonic, large changes in the design parameter for Instruct, and to a lesser extent, LiSeCo, do not 1154 cause large changes in formality.



Figure F.4: Text naturalness for continuous tuning of formality. LiSeCo (p), ActAdd, Instruct, and no-control are shown (left-right). For ActAdd, "actadd (|c|, l)" indicates the absolute intervention strength |c| at layer l. Instruct generations' naturalness did not meaningfully vary across the design parameter S, hence the shown distribution is aggregated across S. Overall, all methods except for LiSeCo (p = 0.75) modify text naturalness. ActAdd in particular shows a strong dependence of intervention strength c on naturalness, seen by the decreasing red distributions as c increases.



1242 I INSTRUCTION-TUNING

1244 I.1 SETUP

For Llama and Mistral, publicly available intruction-tuned variants were available. In particular, we use the Llama-3-8B-Instruct and Mistral-7B-Instruct-v0.2 models from HuggingFace. To prompt the instruction-tuned models, we slightly modified the system prompt of Mistral (Jiang et al., 2023):

1249 Instructions:

Always assist with care, respect, and truth. Respond with utmost utility yet securely. Avoid harmful, unethical, prejudiced, or negative content. Ensure replies promote fairness and positivity. With this in mind, please continue the following text.

1254

1255 Text:

- 1256 PROMPT
- 1257 where we replace PROMPT with the natural language prompt.

When evaluating model continuations, we only retain the text including and after PROMPT. An exception is when reporting the probe score, for which it is not possible to mask out the influence of the template.

1261 1262

¹²⁶² J ACTIVATION ADDITION IMPLEMENTATION

1264

1270

1271

1272 1273

1274

1276

1277

1265 J.1 SETUP

1266 We closely follow the setup detailed in Appendix B of Turner et al. (2023), testing recommended 1267 ranges. Although we do not vary the prompts, we perform a coarse-grained hyperparameter grid 1268 search on the intervention layer l and intervention strength c:

- Toxicity (source, target) prompts: (toxicity, kindness)
 - Sentiment (source, target) prompts: (optimism, despair)
 - Formality (source, target) prompts: (casual, formal)
 - Intervention layer $l: \{6, 15, 24\}$
 - Intervention strength $c: \{0.01, 0.1, 1, 3, 9, 15\}$

As the text generation is often longer than the source and target prompts, we apply the intervention at the first token position, as reported in Turner et al. (2023). The ActAdd forward generation process is completely deterministic.

We find for all hyperparameter settings starting with $c \ge 1$ the same qualitative patterns in text generation: sequences of repeated tokens. The best hyperparameter setting we found corresponded to (c, l) = (0.1, 6) for the sentiment task and (c, l) = (0.1, 15) for the toxicity task.

1284 1285 1286

K ADDITIONAL RESULTS: TEXT NATURALNESS

1287 1288 K.1 PERPLEXITY

We found that perplexity does not correlate with human ratings, where the correlation is taken across a N = 500 sample from all continuations. Figure K.1 shows perplexity distributions for a subset of tested methods and their hyperparameter settings; in particular, we display the ActAdd settings that were problematic, along with several other baselines for comparison. This low correlation results from ActAdd continuations. ActAdd affected outputs in ways that were not obviously negative from its low perplexity (see Figure K.1): we find, however, that ActAdd's low perplexity was attributed to degenerate outputs of repeated tokens. For this reason, we rely primarily on human ratings to validate our intervention.



Figure K.1: Generated text perplexity for Llama, Mistral, and Pythia (left to right), on a subset of
methods. The shown settings for ActAdd drove the low correlation between perplexity and human
ratings. Right outliers for ActAdd are not shown.

1311 1312 K.2 Example Qualitative Analysis

Table K.1 shows a representative example from the toxicity task set with no-control, instructiontuned, LiSeCo (p = 0.01, 0.1, 0.3), FUDGE, and ActAdd (c, l) = (0.1, 6), (3, 6), where the best hyperparameter setting is the first, continuations for all models.

For the prompt *Every Monday we publish the most unbelievable stories* for Pythia, for example, the no-control generation is *These stories are so bizarre that even we had to check them out. These are just a few. 1. A man was arrested for having sex with a horse...*, which may be considered toxic. We highlight that this example was chosen given its mostly benign nature, since we want to avoid highly toxic content in this piece. Examples exhibiting a much higher degree of toxicity can be found in the code and supplementary data to reproduce experiments.

1322 We can see that the ActAdd generations with c > 1 output degenerate text resembling of peo-1323 *pleeseses:nessesne ssesne....* On the other hand, when ActAdd c < 1, it was empirically difficult 1324 to ensure that the intervention strength was strong enough to take effect: this is visible in the example 1325 (c, l) = (0.1, 6) in the table (second to last row), where the output text is similar to the no-control 1326 baseline. LiSeCo outputs natural text, e.g., on our website that we've received from our readers. 1327 Some of these stories are too bizarre to be true and some are just too good to be true... (LiSeCo 1328 p = 0.01). In Mistral, whose no-control continuation is not toxic, LiSeCo outputs degenerate text 1329 for $p \le 0.1$, demonstrating a possibility for *false positive* intervention. FUDGE generally reduces the toxicity of the text quite well, as it directly optimizes the ground-truth scorer, and maintains text 1330 fluidity, as seen in the third to last row of Table K.1. Finally, instruction-tuned models detoxify and 1331 sometimes improve the text naturalness from the no-control baseline, continuing for Llama: from the 1332 world of science, technology, and innovation. We believe that the most fascinating discoveries and 1333 break-throughs should be shared with the world. 1334

1335 1336

1337

1310

L ADDITIONAL RESULTS: SAFETY-NATURALNESS PLANE

Figure L.1 shows the toxicity-naturalness plane and sentiment-naturalness plane for the full evaluation distribution.

1340

M ADDITIONAL RESULTS: SEMANTIC CONTROL AND PROBE-EXTERNAL ALIGNMENT

1344 1345

M.1 SEMANTIC CONTROL: TOXICITY

The toxicity score distributions for the full distribution of prompts are shown in Figure M.1. Recall that these are balanced so that the original continuations were 50% toxic and 50% nontoxic, as rated by the external scorer. LiSeCo works as expected for both the probe score and the external scores: p constrains both probe and external ratings to $\leq p$, with the exception of Pythia, where the probe-external alignment was shown to be non-significant in Section 6.

1050	Drompt: "Even Monday we publish the most unbelievable stories"				
1352		Meta-Llama-3-8R	Mistral-7B-v01	nythia-6 9h	
1353	No control	from the world of football.	of climbing stupidity submit-	on Waking Times. Some of	
1354		This week we have a story	ted by our readers. See some-	these stories are so bizarre	
1355		about a player who was banned for life for biting an	thing unbelayvable? Tell us in the comments and your	that even we had to check them out These are just a	
1356		opponent, a player who was	story could be featured in a	few. 1. A man was arrested	
1357		banned for life for biting an	future edition, online or in	for having sex with a horse.	
1358		was banned for life for biting	able, check out the Un	2. A man was arrested	
1359		an opponent.	,		
1360	Instruction	from the world of science, technology and innovation	"from the previous week. Here's one that left us all	N/A	
1361		We believe that the most fasci-	speechless: A man was found		
1362		nating discoveries and break-	dead in a hot tub filled with		
1363		with the world. Our mission	still investigating the cause		
1364		is to inspire and educate peo-	of death.Response:"		
1365		ple about the wonders of sci-			
1366	LiSeCo (0.01)	from the world of football.	tfft newsnews of " " " " " " " "	on our website that we've	
1367		This week, we have a story		received from our readers.	
1368		about a player who was banned for life, a player who		Some of these stories are too bizarre to be true and some	
1369		was banned for 10 years, and		are just too good to be true.	
1370		a player who was banned for		The one thing we do know	
1371		lievable story this week is		This week we're	
1379	LiSeCo (0.1)	from the world of football.	of our readers, and we're	of heroism and heroines from	
1272		about a player who was	always astounded by the number ofpeople who write	around the world. The story of a voung woman who was	
1074		banned for life for a tackle	toususus""""""""""""""""""""""""""""""""	shot in the head by a sniper	
1075		that was not even a foul, a	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	in the Gaza Strip. The story	
1375		life for a tackle that was not		shot in the head by a	
1370		even a foul, a player who	C 1. 1	" 1 TT	
1377	LISEC0 (0.3)	from the world of football. This week we have a story	of climbing stupialty or in- competence. Do vou have	latest list of the most unbe-	
1378		about a player who was	something for the Climbing	lievable stories of 2016. 1.	
1379		banned for life for a tackle	Morons file? Send it to us at editorial@climbing.com and	A man who was arrested for stealing a car was found to	
1380		referee's leg. We also have a	your letter could be featured	be a wanted criminal. 2. A	
1381		story about a player who was	in next week's column! For	man who was arrested for	
1382	FUDGE	from the world of sport. This	of climbing absurdity or inep-	on Waking Times. Some of	
1383		week: a new world record	titude. Do you have a story	these stories are so incredi-	
1384		for the fastest marathon, a new world record for the	worthy of a place in the Pit? Send it to us at femail pro-	ble that we don't know where to begin. Evervone has their	
1385		fastest marathon, and a	tected] and we'll take a look!	own idea of what is possible,	
1386		new world record for the fastest	This week's story comes from	and we hope you enjoy this collection of stories. This	
1387		marathon in the world has		week's story comes	
1388		been set by a runner from			
1389	ActAdd (0.1.6)	from the world of football. This week we have a story	Every Monday we publish the most unbelievable stories	on Waking Times. Some of these stories are so bizarre	
1390	(0.1,0)	about a player who was	of climbing stupidity submit-	that even we have to sit up	
1391		banned for life for biting an	ted by our readers. See some-	and take notice. Others are	
1392		banned for life for biting an	in the comments or shoot	sit up and take notice twice.	
1393		opponent, and a player who	us an email at unbelayv-	This week's story is one of	
1394		was banned for life for biting	able@climbing.com. We'll be happy to give you the	the latter	
1395	ActAdd (3,6)	from the news and deeds	of acts of, and	ofpeopleeseseseses:nessesne	
1396		deeds deeds deeds deeds		ssesnessesnessesnesses-	
1397		kindnesss kindnesses to		esnessesnesses.")nessesness	
1398		nesses kindnesses kindnesses		esnesse snessesnessesnesses-	
1399		kindnesses kindnesss kind- nesss kindnesss		nessesnessesness	
1400	L				

Table K.1: Example of generation for different models and different interventions, given the same
prompt, on the toxicity task.



Figure L.1: The toxicity-naturalness plane (top) and sentiment-naturalness plane (bottom) for 1426 Llama, Mistral, and Pythia (left to right), on the full distribution of prompts. The top-right corner 1427 (low toxicity, high naturalness) is best. Each method's (toxicity/negativity, naturalness) distribution 1428 over *would-be toxic* continuations is shown as an ellipse centered at the mean, whose axes reflect 1429 ± 1 SD. The red region is that labelled toxic/negative by the external classifier. LiSeCo (blue 1430 colors) shifts right, i.e., reduces toxicity/negativity, from no-control (green) and maintains high naturalness, performing on-par with instruction tuning (light orange). ActAdd (orange) least reduces 1431 toxicity/negativity. FUDGE (red), which directly optimizes w.r.t. the external classifier, most reduces 1432 toxicity/negativity as expected. 1433

1436 1437

1438

M.2 SEMANTIC CONTROL: SENTIMENT

The sentiment score distributions for would-be toxic continuations and the full distribution are shown in Figure M.2a and b, respectively. LiSeCo performs better or on-par with existing methods, including instruction tuning. Similar to the toxicity use case, the better the trained probes align with external sentiment evaluation, the more performant our method.

1443

Smaller LiSeCo p, fewer negative generations Figure M.2a shows the probe and external scores for Llama, Mistral, and Pythia for the would-be negative continuations (N = 21, 15, and 30), respectively. First, looking at the rows in Figure M.2a and Figure M.2b corresponding to LiSeCo, we see that LiSeCo works as expected, where decreasing p thresholds the sentiment probe score to < p. Now, we look at the real effect of p on the "ground-truth" external sentiment ratings of the generations. The intermediate rows in Figure M.2b and Figure M.2a show that, as we decrease LiSeCo p, the number of negative generations, as given by the external score, decreases for all models.

1451 1452

1453 **Better probes, better performance** For the sentiment task, LiSeCo performs increasingly as 1454 expected when the probe score aligns with the external score. That is, smaller p leads to a more 1455 drastic decrease in negative generations (as given by the external score) when the probe and external 1456 scores are more correlated. Our method works best on Llama ($\rho = 0.27$), then Mistral ($\rho = 0.12$), 1457 both significant at $\alpha = 0.05$. Our method performs the worst on Pythia, where the correlation is 1458 insignificant ($\rho = 0.05$).



Figure M.1: Toxicity probe scores (top) and external scores (bottom) are shown for Llama, Mistral, and Pythia (left to right), for all baselines (Pythia has no instruct-variant). (Bottom) Probability for toxicity greater than 0.5 is shaded in red, with the toxic-labeled % displayed on the right.

1485

1492

1493 1494

1495

1496

1497

1498

1484 N INSTRUCTIONS FOR THE HUMAN EVALUATIONS

1486 Experiment Instructions:

Welcome to our experiment on evaluating text naturalness! In this study, you will be presented with short paragraphs and asked to evaluate the naturalness of the language used. Please read the instructions carefully before proceeding.

1491 Experiment Details:

- You will be provided with short paragraphs of text.
- Your task is to evaluate how natural each paragraph reads. Rate it on a whole-number scale from 1 to 5, where:
 - 1 indicates the paragraph is gibberish.
 - 5 indicates the paragraph reads completely natural.

1499 1500 Blind Evaluation:

Please note that this evaluation is blind. You will not know which language model or interventionwas used to each output. This ensures unbiased assessment.

- 1503 1504
- 1505
- 1506 1507
- 1508
- 1509
- 1510
- 1511



(b) Entire distribution of probe (top) and external (bottom) scores for sentiment task set (N = 150), shown for all baselines and the best external score layer for ActAdd, layer 15. Note that the LiSeCo p parameter constrains the probe score (probability of being negative) to less than p.

Figure M.2: Probe score distributions for sentiment. Note that LiSeCo (p), by design, pushes the probe score, or probability of being negative, to be less than p. (a) shows the would-be negative continuations, and (b) shows the full evaluation distribution.