

CALIBRATING GENERATIVE MODELS TO DISTRIBUTIONAL CONSTRAINTS

Anonymous authors

Paper under double-blind review

ABSTRACT

Generative models frequently suffer miscalibration, wherein class probabilities and other statistics of the sampling distribution deviate from desired values. We frame calibration as a constrained optimization problem and seek the closest model in Kullback-Leibler divergence satisfying calibration constraints. To address the intractability of imposing these constraints exactly, we introduce two surrogate objectives for fine-tuning: (1) the relax loss, which replaces the constraint with a miscalibration penalty, and (2) the reward loss, which converts calibration into a reward fine-tuning problem. We demonstrate that these approaches substantially reduce calibration error across hundreds of simultaneous constraints and models with up to one billion parameters, spanning applications in protein design, image generation, and language modeling.

1 INTRODUCTION

Generative models commonly produce samples whose statistics deviate systematically from desired values. Such *miscalibration* occurs in many domains. Image models, such as GANs and diffusion models, exhibit mode collapse, producing images that cover only a subset of the training distribution (Arora & Zhang, 2017; Qin et al., 2023). Language models represent gender, race, religion, and age in ways that reinforce societal biases (Gallegos et al., 2024). In synthetic biology applications, protein structure models produce samples that have alpha-helical and beta-strand substructures at frequencies atypical of proteins found in nature (Lu et al., 2025), and DNA models generate samples that contain subsequences at frequencies that differ from those in human DNA (Sarkar et al., 2024). These calibration errors arise from many sources including dataset imbalances, suboptimal training dynamics, and post-hoc adjustments such as low-temperature sampling or preference fine-tuning.

We frame calibration as a constrained optimization problem: find the distribution closest in Kullback-Leibler (KL) divergence to the base model that satisfies a set of expectation constraints. We introduce two fine-tuning algorithms—**CGM-relax** and **CGM-reward** (“calibrating generative models”)—that approximately solve the calibration problem by stochastic optimization. We demonstrate across three applications that CGM effectively calibrates high-dimensional generative models to meet hundreds of simultaneous constraints.

Problem statement. Consider a trained “base” generative model $p_{\theta_{\text{base}}}(\mathbf{x})$ with parameters θ_{base} , a statistic $\mathbf{h}(\mathbf{x})$, and an expectation value desired for the statistic \mathbf{h}^* . We say $p_{\theta_{\text{base}}}$ is *calibrated* if $\mathbb{E}_{p_{\theta_{\text{base}}}}[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*$ and *miscalibrated* if $\mathbb{E}_{p_{\theta_{\text{base}}}}[\mathbf{h}(\mathbf{x})] \neq \mathbf{h}^*$. In the case that $p_{\theta_{\text{base}}}$ is miscalibrated, our goal is to fine-tune its parameters θ_{base} to some θ such that p_{θ} is calibrated.

For example, if $\mathbf{h}(\mathbf{x}) = \mathbb{1}\{\mathbf{x} \in C\}$ is the 0-1 function indicating whether \mathbf{x} belongs to class C , then $\mathbb{E}_{p_{\theta_{\text{base}}}}[\mathbf{h}(\mathbf{x})] = p_{\theta_{\text{base}}}(\mathbf{x} \in C)$ is the probability that $p_{\theta_{\text{base}}}$ generates a member of class C . When $\mathbf{h}^* > \mathbb{E}_{p_{\theta_{\text{base}}}}[\mathbf{h}(\mathbf{x})]$, calibration corresponds to increasing the probability of class C .

For a given $\mathbf{h}(\cdot)$ and \mathbf{h}^* , many calibrated models may exist. Provided a calibrated model exists, we seek the one that is closest to the base model in KL divergence,

$$p_{\theta^*} := \arg \min_{p_{\theta}} \text{D}_{\text{KL}}(p_{\theta} \parallel p_{\theta_{\text{base}}}) \quad \text{such that } \mathbb{E}_{p_{\theta}}[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*, \quad (1)$$

where $\text{D}_{\text{KL}}(p' \parallel p) = \mathbb{E}_{p'}[\log p'(\mathbf{x})/p(\mathbf{x})]$ for p' with a probability density with respect to p . Out of many possible notions of distance we choose D_{KL} because it is simple and, as we will see, is tractable for several classes of generative models.

Related work. Within the generative modeling community, there are a wealth of fine-tuning methods that incorporate preferences at the level of individual samples through a user-specified reward (Christiano et al., 2017; Rafailov et al., 2023; Uehara et al., 2024; Domingo-Enrich et al., 2025). None of these methods solves problem (1), which imposes a hard constraint at the distribution level.

Two prior works (Khalifa et al., 2021; Shen et al., 2024) propose fine-tuning procedures for distribution level constraints, but each applies to a single model class. Khalifa et al. (2021), the most similar to the present work, fine-tunes autoregressive language models to match distributional constraints with an algorithm similar to CGM-reward. Shen et al. (2024) propose a method for balancing class proportions in diffusion models that relies upon optimal transport. Compared to the present work, neither method reduces a majority of calibration error, and Khalifa et al. (2021) demonstrates their algorithm only for low-dimensional (<10) constraints.

Lastly, we clarify that our definition of generative model calibration differs from the definition of the same term adopted in the setting of supervised learning (e.g., Lichtenstein et al., 1977; Dawid, 1982; Naeini et al., 2015; Guo et al., 2017; Vaicenavicius et al., 2019). In this setting, calibration means that the conditional expectation of the response given the model prediction is equal to the prediction i.e., on average, the predictor is correct. Appendix A gives an extended discussion of related work.

2 CALIBRATING GENERATIVE MODELS WITH CGM-RELAX AND REWARD

The calibration problem is challenging for non-trivial generative models because both the objective and the calibration constraint in equation (1) are defined by intractable expectations. To address this problem, we propose two alternative objectives whose *unconstrained* optima approximate the solution to (1). These objectives still involve expectations under p_θ , but we show how to compute unbiased estimates of their gradients, which permits their minimization by stochastic optimization.

We call our algorithms optimizing the two surrogate loss functions CGM-relax and CGM-reward (Algorithms 1 and 2, respectively). These algorithms require only that one can draw samples $\mathbf{x} \sim p_\theta$ and compute $p_\theta(\mathbf{x})$ and $\nabla_\theta \log p_\theta(\mathbf{x})$.

2.1 THE RELAX LOSS

The relax loss avoids the intractability of imposing the calibration constraint exactly by replacing it with a constraint violation penalty

$$\mathcal{L}^{\text{relax}}(\theta) := \underbrace{\|\mathbb{E}_{p_\theta}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2}_{\mathcal{L}^{\text{viol}}} + \lambda \underbrace{\text{D}_{\text{KL}}(p_\theta \| p_{\theta_{\text{base}}})}_{\mathcal{L}^{\text{KL}}}, \quad (2)$$

where $\lambda > 0$ is a hyperparameter that trades off between satisfying the calibration constraint and minimizing the KL divergence. In our experiments we choose λ by a grid-search. In the limit as $\lambda \rightarrow 0$, $\mathcal{L}^{\text{viol}}$ is the dominant term in the relax loss, and we expect the minimizer of (2) to approach the solution of the calibration problem (1).

Suppose we have M independent samples $\{\mathbf{x}_m\}_{m=1}^M$ from p_θ . To estimate $\mathcal{L}^{\text{relax}}$ we separately estimate the KL divergence (\mathcal{L}^{KL}) by

$$\widehat{\mathcal{L}}^{\text{KL}} := \frac{1}{M} \sum_{m=1}^M \log \frac{p_\theta(\mathbf{x}_m)}{p_{\theta_{\text{base}}}(\mathbf{x}_m)},$$

and the constraint penalty ($\mathcal{L}^{\text{viol}}$) by

$$\widehat{\mathcal{L}}^{\text{viol}} := \left\| \frac{1}{M} \sum_{m=1}^M \mathbf{h}(\mathbf{x}_m) - \mathbf{h}^* \right\|^2 - \frac{1}{M(M-1)} \sum_{m=1}^M \left\| \mathbf{h}(\mathbf{x}_m) - \frac{1}{M} \sum_{m'=1}^M \mathbf{h}(\mathbf{x}_{m'}) \right\|^2, \quad (3)$$

where the second term is a bias correction. Combining these estimators yields our overall estimator for the relax objective, $\widehat{\mathcal{L}}^{\text{relax}} = \widehat{\mathcal{L}}^{\text{viol}} + \lambda \widehat{\mathcal{L}}^{\text{KL}}$. Appendix B.1 shows $\widehat{\mathcal{L}}^{\text{relax}}$ is unbiased for $\mathcal{L}^{\text{relax}}$.

2.2 THE REWARD LOSS

The reward loss avoids the intractability of imposing the calibration constraint exactly by leveraging a connection between the calibration problem (1) and the *maximum entropy problem* (Jaynes, 1957;

Kullback, 1959; Csiszár, 1975). We first introduce the maximum entropy problem and then show how to approximate its solution with samples from $p_{\theta_{\text{base}}}$. Lastly, we propose the reward loss as a divergence to this approximate solution and describe connections to reward fine-tuning.

Maximum entropy problem. The maximum entropy problem solves

$$\arg \min_{p \in \mathcal{P}(p_{\theta_{\text{base}}})} \text{D}_{\text{KL}}(p \parallel p_{\theta_{\text{base}}}), \quad \text{such that } \mathbb{E}_p[\mathbf{h}(\mathbf{x})] = \mathbf{h}^* \quad (4)$$

where $\mathcal{P}(p)$ is the collection of probability distributions that have a density with respect to p . The calibration problem and the maximum entropy problem differ only in their domains: the domain of the calibration problem is generative models p_{θ} in the same parametric class as $p_{\theta_{\text{base}}}$, rather than the nonparametric set $\mathcal{P}(p_{\theta_{\text{base}}})$. Despite this difference, we obtain an alternative objective by considering the solution to (4). The following theorem characterizes this solution.

Theorem 2.1. *Suppose \mathbf{h}^* lies in the relative interior of the set of all attainable moments of \mathbf{h} by distributions $P \in \mathcal{P}(p_{\theta_{\text{base}}})$. Then there exists a vector $\boldsymbol{\alpha}^*$ for which*

$$p_{\boldsymbol{\alpha}^*}(\mathbf{x}) \propto p_{\theta_{\text{base}}}(\mathbf{x}) \exp\{r_{\boldsymbol{\alpha}^*}(\mathbf{x})\}, \quad r_{\boldsymbol{\alpha}}(\mathbf{x}) := \boldsymbol{\alpha}^{\top} \mathbf{h}(\mathbf{x}) \quad (5)$$

satisfies $\mathbb{E}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*$ and $p_{\boldsymbol{\alpha}^*}$ is the solution to the maximum entropy problem.

Appendix C provides further exposition to the maximum entropy problem as well as a proof.

The domain of the calibration problem may not contain $p_{\boldsymbol{\alpha}^*}$. However, if the class of generative models is sufficiently expressive, its optimum p_{θ^*} will be close to $p_{\boldsymbol{\alpha}^*}$. This observation suggests a second way to remove the constraint in equation (1): fine-tune p_{θ} to minimize a divergence to $p_{\boldsymbol{\alpha}^*}$.

In Appendix C.3 we demonstrate that a similar statement holds for the relax loss: when p_{θ} is sufficiently expressive, the optimum of the relax loss is close to $p_{\lambda}(\mathbf{x}) \propto p_{\theta_{\text{base}}}(\mathbf{x}) \exp\{r_{\boldsymbol{\alpha}_{\lambda}}(\mathbf{x})\}$, where $\boldsymbol{\alpha}_{\lambda}$ depends on the regularization strength $\lambda > 0$ and is not generally equal to $\boldsymbol{\alpha}^*$. However, as $\lambda \rightarrow 0$, $\|\boldsymbol{\alpha}_{\lambda} - \boldsymbol{\alpha}^*\|$ approaches zero at rate λ . This formalizes our intuition from Section 2.1 that as $\lambda \rightarrow 0$, the relax loss solves the calibration problem.

Estimating $p_{\boldsymbol{\alpha}^*}$. The idea of minimizing a divergence to $p_{\boldsymbol{\alpha}^*}$ introduces a challenge: even when the solution $p_{\boldsymbol{\alpha}^*}$ to the maximum entropy problem (4) exists, its parameters $\boldsymbol{\alpha}^*$ are not immediately computable. To address this challenge, we leverage Wainwright & Jordan (2008, Theorem 3.4), which states that when the assumptions of Theorem 2.1 hold and there are no redundancies among the constraints \mathbf{h} , solving problem (4) is equivalent to computing

$$\arg \max_{\boldsymbol{\alpha}} \boldsymbol{\alpha}^{\top} \mathbf{h}^* - \log \left(\int \exp\{r_{\boldsymbol{\alpha}}(\mathbf{x})\} p_{\theta_{\text{base}}}(\mathbf{x}) d\mathbf{x} \right). \quad (6)$$

In other words, by solving (6) one obtains the parameters $\boldsymbol{\alpha}^*$ of $r_{\boldsymbol{\alpha}}(\mathbf{x})$, which then determine the solution $p_{\boldsymbol{\alpha}^*}$ to the maximum entropy problem up to a normalizing constant.

However, a difficulty of solving (6) is that the integral in the second term will be intractable for most generative models. We propose drawing N independent samples $\{\mathbf{x}_n\}_{n=1}^N$ from $p_{\theta_{\text{base}}}$ and replacing the integral with respect to $p_{\theta_{\text{base}}}$ by the integral with respect to the empirical distribution that places probability mass N^{-1} on each of the samples \mathbf{x}_n ,

$$\hat{\boldsymbol{\alpha}}_N = \arg \max_{\boldsymbol{\alpha}} \boldsymbol{\alpha}^{\top} \mathbf{h}^* - \log \left(\frac{1}{N} \sum_{n=1}^N \exp\{r_{\boldsymbol{\alpha}}(\mathbf{x}_n)\} \right). \quad (7)$$

Problem (7) is concave, and when $\hat{\boldsymbol{\alpha}}_N$ is well-defined (see Appendix C.2), it can be found by convex solvers. We demonstrate in Appendix C.4 that $\hat{\boldsymbol{\alpha}}_N$ converges to $\boldsymbol{\alpha}^*$ in the limit of many samples N , and we derive an expression for the asymptotic variance of $\hat{\boldsymbol{\alpha}}_N$.

$\mathcal{L}^{\text{reward}}$ and its estimation. With $\hat{\boldsymbol{\alpha}}_N$ in hand, we formulate our second loss as a divergence to $p_{\hat{\boldsymbol{\alpha}}_N}$. For simplicity and because it avoids the requirement to compute the normalizing constant of $p_{\hat{\boldsymbol{\alpha}}_N}$, we again choose the KL divergence. In particular, we define the reward loss $\mathcal{L}^{\text{reward}}$ to be

$$\mathcal{L}^{\text{reward}}(\theta) = \text{D}_{\text{KL}}(p_{\theta} \parallel p_{\hat{\boldsymbol{\alpha}}_N}) = \underbrace{\mathbb{E}_{p_{\theta}}[\log p_{\theta}(\mathbf{x})/p_{\theta_{\text{base}}}(\mathbf{x})]}_{\mathcal{L}^{\text{KL}}=\text{D}_{\text{KL}}(p_{\theta} \parallel p_{\theta_{\text{base}}})} + \underbrace{\mathbb{E}_{p_{\theta}}[-r_{\hat{\boldsymbol{\alpha}}_N}(\mathbf{x})]}_{\mathcal{L}^{\text{r}}}, \quad (8)$$

Algorithm 1 CGM-relax fine-tuning

Require: $p_{\theta_{\text{base}}}$, $\mathbf{h}(\cdot)$, \mathbf{h}^* , M , and λ

▷ Initialize and optimize

$p_{\theta} \leftarrow p_{\theta_{\text{base}}}$

while not converged **do**

▷ Sample and compute weights

$\mathbf{x}_1, \dots, \mathbf{x}_M \stackrel{i.i.d.}{\sim} p_{\text{stop-grad}(\theta)}$

$w_m \leftarrow p_{\theta}(\mathbf{x}_m) / p_{\text{stop-grad}(\theta)}(\mathbf{x}_m)$

▷ KL loss with LOO baseline

$l_m \leftarrow \log p_{\text{stop-grad}(\theta)}(\mathbf{x}_m) / p_{\theta_{\text{base}}}(\mathbf{x}_m)$

$l_m^{\text{LOO}} \leftarrow l_m - \frac{1}{M-1} \sum_{m' \neq m} l_{m'}$

$\widehat{\mathcal{L}}^{\text{KL}} \leftarrow \frac{1}{M} \sum w_m l_m^{\text{LOO}}$

▷ Constraint violation loss

$\mathbf{h}_m \leftarrow w_m (\mathbf{h}(\mathbf{x}_m) - \mathbf{h}^*)$

$\widehat{\mathcal{L}}^{\text{viol}} \leftarrow \left\| \frac{1}{M} \sum \mathbf{h}_m \right\|^2 - \frac{1}{M} \widehat{\text{Var}}[\mathbf{h}_{1:M}]$,

$\widehat{\text{Var}}[\mathbf{h}_{1:M}] = \frac{1}{M-1} \sum \left\| \mathbf{h}_m - \frac{1}{M} \sum \mathbf{h}_{m'} \right\|^2$

▷ Total loss and update

$\widehat{\mathcal{L}}^{\text{relax}} = \lambda \widehat{\mathcal{L}}^{\text{KL}} + \widehat{\mathcal{L}}^{\text{viol}}$

$\theta \leftarrow \text{gradient-step}(\theta, \nabla_{\theta} \widehat{\mathcal{L}}^{\text{relax}})$

Algorithm 2 CGM-reward fine-tuning

Require: $p_{\theta_{\text{base}}}$, $\mathbf{h}(\cdot)$, \mathbf{h}^* , M , N

▷ Estimate α^* for reward

$\mathbf{x}_1, \dots, \mathbf{x}_N \stackrel{i.i.d.}{\sim} p_{\theta_{\text{base}}}$

$\widehat{\alpha}_N \leftarrow \arg \max \alpha^{\top} \mathbf{h}^* - \log \sum \exp\{r_{\alpha}(\mathbf{x}_n)\}$

▷ Initialize and optimize

$p_{\theta} \leftarrow p_{\theta_{\text{base}}}$

while not converged **do**

▷ Sample and compute weights

$\mathbf{x}_1, \dots, \mathbf{x}_M \stackrel{iid}{\sim} p_{\text{stop-grad}(\theta)}$

$w_m \leftarrow p_{\theta}(\mathbf{x}_m) / p_{\text{stop-grad}(\theta)}(\mathbf{x}_m)$

▷ KL loss with LOO baseline

$l_m \leftarrow \log p_{\text{stop-grad}(\theta)}(\mathbf{x}_m) / p_{\theta_{\text{base}}}(\mathbf{x}_m)$

$l_m^{\text{LOO}} \leftarrow l_m - \frac{1}{M-1} \sum_{m' \neq m} l_{m'}$

$\widehat{\mathcal{L}}^{\text{KL}} \leftarrow \frac{1}{M} \sum w_m l_m^{\text{LOO}}$

▷ Negative reward with LOO baseline

$r_m^{\text{LOO}} \leftarrow r_{\widehat{\alpha}}(\mathbf{x}_m) - \frac{1}{M-1} \sum_{m' \neq m} r_{\widehat{\alpha}}(\mathbf{x}_{m'})$

$\widehat{\mathcal{L}}^{\text{r}} \leftarrow -\frac{1}{M} \sum w_m r_m^{\text{LOO}}$

▷ Total loss and update

$\widehat{\mathcal{L}}^{\text{reward}} = \widehat{\mathcal{L}}^{\text{KL}} + \widehat{\mathcal{L}}^{\text{r}}$

$\theta \leftarrow \text{gradient-step}(\theta, \nabla_{\theta} \widehat{\mathcal{L}}^{\text{reward}})$

where $C = \mathbb{E}_{p_{\theta_{\text{base}}}}[\exp\{r_{\widehat{\alpha}_N}(\mathbf{x})\}]$ is a normalizing constant that does not depend on θ .

We call $r_{\alpha}(\mathbf{x})$ the *reward* and $\mathcal{L}^{\text{reward}}$ the reward loss because $\mathcal{L}^{\text{reward}}$ coincides with the objective of reward fine-tuning algorithms. The goal of reward fine-tuning is to fine-tune the base generative model $p_{\theta_{\text{base}}}$ to a tilted version of itself, where the tilt is determined by a so-called reward $r(\mathbf{x})$.

Just as for \mathcal{L}^{KL} in the relax loss (2), Monte Carlo sampling provides an unbiased estimate of \mathcal{L}^{r} . This, in turn, gives us an unbiased estimate of the reward loss $\mathcal{L}^{\text{reward}}$.

2.3 GRADIENT ESTIMATION

We next describe our approach to computing unbiased estimates for the gradients of $\mathcal{L}^{\text{relax}}(\theta)$ and $\mathcal{L}^{\text{reward}}(\theta)$. This enables optimization of the relax and reward losses via stochastic optimization. We leverage the score function gradient estimator (Williams, 1992; Ranganath et al., 2014) and a similar importance sampling-based gradient estimator for the relax loss.

Score function gradient estimation. The primary challenge to computing gradients is the inability to directly exchange the order of the gradients and expectations taken with respect to θ . That is, because $\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta} \mathbb{E}_{p_{\theta}}[f(\mathbf{x}, \theta)] \neq \mathbb{E}_{p_{\theta}}[\nabla_{\theta} f(\mathbf{x}, \theta)]$, $\nabla_{\theta} \mathcal{L}(\theta)$ can not in general be usefully approximated by $M^{-1} \sum \nabla_{\theta} f(\mathbf{x}_m, \theta)$ from samples \mathbf{x}_m of p_{θ} . To address this challenge, we observe

$$\mathcal{L}(\theta) = \mathcal{L}(\theta, \theta') := \mathbb{E}_{p_{\theta'}} \left[\frac{p_{\theta}(\mathbf{x})}{p_{\theta'}(\mathbf{x})} f(\mathbf{x}, \theta) \right], \quad (9)$$

for any set of model parameters θ' . Since the expectation in equation (9) does not depend on θ , we can approximate its gradient with Monte Carlo samples from $p_{\theta'}$. The density ratio $p_{\theta}(\mathbf{x}_m) / p_{\theta'}(\mathbf{x}_m)$ in equation (9) can be understood as the weights of an importance sampling estimate against target p_{θ} with proposal $p_{\theta'}$.

To estimate the gradients of the relax and reward losses, we choose proposal equal to the current model, i.e., $\theta' = \theta$. In this case, the importance weight is equal to 1 while its gradient is the “score” function $(\nabla_{\theta} p_{\theta}(\mathbf{x}_m)) / p_{\theta}(\mathbf{x}_m) = \nabla \log p_{\theta}(\mathbf{x})$, which is nonzero in general (Mohamed et al., 2020). Algorithms 1 and 2 each demonstrate an implementation that computes these weights with a copy of the parameters θ detached from the computational graph, which we denote by `stop-grad`(θ).

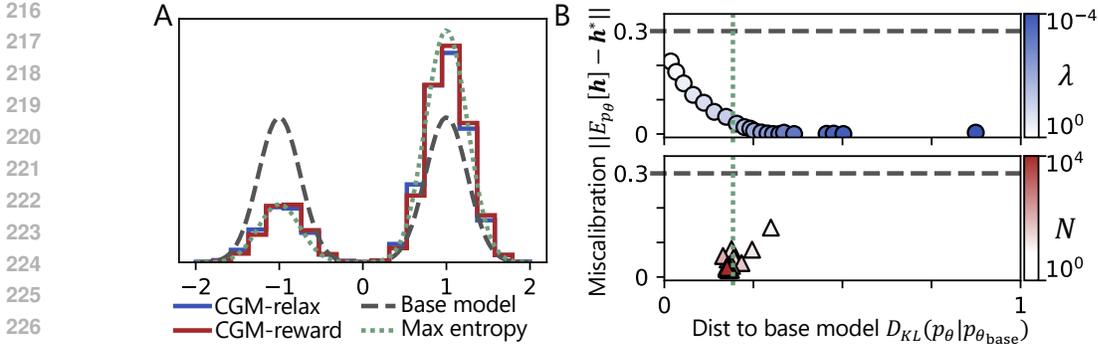


Figure 1: Calibrating mixture proportions in a diffusion model targeting a 1D GMM. **A:** The CGM-relax and CGM-reward solutions closely approximate the maximum entropy solution. **B:** (top) The CGM-relax regularization parameter λ trades off between constraint satisfaction and closeness to the base model (bottom) CGM-reward is accurate when enough samples N are used to estimate α^* .

Although the term $\mathcal{L}^{\text{viol}}$ that appears in the relax loss is not of the form $\mathbb{E}_{p_\theta}[f(\mathbf{x}, \theta)]$, we can still construct an unbiased estimate to its gradient using importance sampling (see Appendix B.2).

Although score function gradient estimates are known to suffer from high variance (Mohamed et al., 2020), we show that, paired with variance reduction strategies (Appendix B.2), they perform well even in problem settings with high-dimensional latent variables, such as diffusion models and masked language models (Section 4.1).

3 SIMULATIONS: DETERMINING WHEN CGM THRIVES AND STRUGGLES

To understand the success and failure cases of CGM, we perform evaluations in a tractable “toy” setting. This setting allows us to understand the role of the CGM hyperparameters λ and N , and to test CGM in challenging problem settings, including rare events and high-dimensional constraints.

We provide additional discussion of our simulation experiments in Appendix D, including an overview of diffusion models in Appendix D.1 and a comparison between CGM and the Augmented Lagrangian (AL) algorithm (Hestenes, 1969) in Appendix D.4.

Simulation setup and evaluation. We consider fine-tuning a diffusion model targeting a Gaussian mixture model (GMM) to reweight the mixture proportions of each mode. Here, $p_\theta(\mathbf{x})$ is a generative model of continuous paths $\mathbf{x} = (\mathbf{x}(t))_{t \in [0,1]}$, whose evolution is described by a stochastic differential equation (SDE). To sample from p_θ , one first draws $\mathbf{x}(0)$ from the tractable initial distribution and then simulates the SDE starting from time $t=0$ up until time $t=1$.

Evaluating CGM on a diffusion model whose terminal distribution is a GMM has several advantages. First, we may choose the base diffusion model so that the final marginal $p_\theta(\mathbf{x}(1))$ exactly matches the target GMM (Anderson, 1982; Song et al., 2021); this enables us to focus solely on calibration rather than fitting the base model. And since the calibration constraint depends on the path only at time $t=1$, we can compute the KL divergence of the maximum entropy solution to the base model, and thereby measure the suboptimality of the solutions produced by CGM.

Selecting hyperparameters for CGM-relax and CGM-reward. We first initialize our base model $p_{\theta_{\text{base}}}$ such that $p_{\theta_{\text{base}}}(\mathbf{x}(1))$ is a one-dimensional Gaussian mixture with two well separated modes (Figure 1 A). We define the calibration problem with statistic $\mathbf{h}(\mathbf{x}) = \mathbb{1}\{\mathbf{x}(1) > 0\}$ to upweight the mass in right mode from $\mathbb{E}_{p_{\theta_{\text{base}}}}[\mathbf{h}(\mathbf{x}(1))] = 0.5$ to $\mathbf{h}^* = 0.8$.

For CGM-relax we observe that the regularization parameter λ trades off between constraint satisfaction and deviation from the base model (Figure 1B). With large λ the model deviates little from $p_{\theta_{\text{base}}}$ but does not satisfy the constraint, whereas for small λ the model satisfies the constraint but has KL to $p_{\theta_{\text{base}}}$ that exceeds that of the maximum entropy solution. For CGM-reward, we observe that increasing N results in more accurate recovery of the variational parameters α^* and thereby a better approximation to the maximum entropy solution. For appropriate hyperparameters, both solve the calibration problem to high accuracy. In the remaining experiments, we perform grid-search to select λ in CGM-relax and use $N = 10^5$ samples to estimate α^* in CGM-reward.

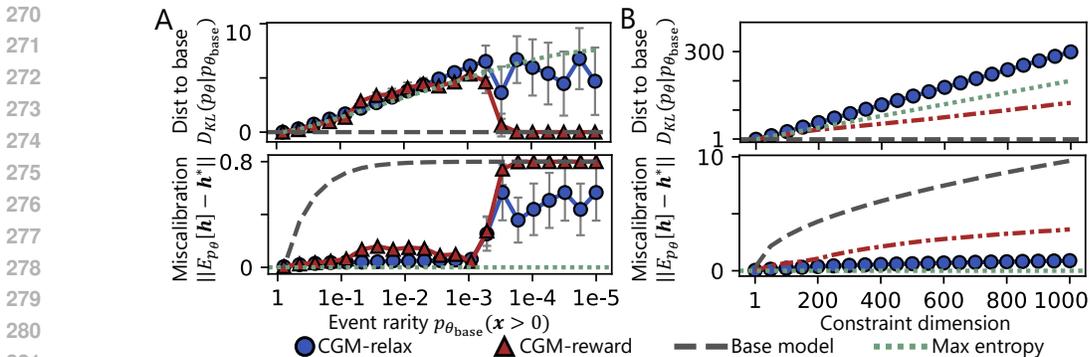


Figure 2: **A:** CGM effectively upweights the probability of a rare mode in a 1D GMM. **B:** CGM-relax calibrates the base model to up to 10^3 constraints, whereas CGM-reward is not well-defined for >30 constraints. When $\hat{\alpha}_N$ is fixed to α^* (red dashed line), CGM-relax outperforms CGM-reward.

Upweighting rare events. Increasing the proportion of generations belonging to rare classes is central to applications including protein ensemble modeling (Lewis et al., 2025) and reinforcement learning (O’Kelly et al., 2018). To assess the performance of CGM in this setting, we consider variations of the GMM reweighting problem in which we consider increasingly small mixture proportions $\pi = \mathbb{E}_{p_{\theta_{base}}}[\mathbf{h}(\mathbf{x}(1))]$ of the mode to upweight by calibration. We vary π from $\mathbf{h}^* = 0.8$ (already calibrated) to 10^{-5} and use a constant batch size $M = 10^2$.

We find that both algorithms perform well with base model event rarity as small as $\pi = 10^{-3}$; the majority of miscalibration is reduced without divergence from the base model much larger than the maximum entropy solution (Figure 2A). This is surprising since for $\pi = 10^{-3}$, most batches sampled from $p_{\theta_{base}}$ contain no samples belonging to the second mode. Performance degrades below this threshold, but we suspect larger batch sizes would allow upweighting even rarer events.

Scalability to high-dimensional models and constraints. We next evaluate how performance depends on the dimensionality, k , of the GMM and the constraint. We take the base model to be a product of one-dimensional GMMs with marginals as in Figure 1A. For the calibration constraint, we choose the $\mathbf{h}(\mathbf{x}) = [\mathbb{1}\{\mathbf{x}(1)[1] > 0\}, \dots, \mathbb{1}\{\mathbf{x}(1)[k] > 0\}]$, where $\mathbf{x}(1)[i]$ is the i th dimension of $\mathbf{x}(1)$ and $\mathbf{h}^* = [0.8, \dots, 0.8]$. Since both the base model $p_{\theta_{base}}$ and maximum entropy solution p_{α^*} are independent across dimension, the KL distance between these two distributions grows linearly in dimension. The multimodality of this model, with 2^k modes, mimics the multimodality of practical generative models. We perform CGM-relax and CGM-reward with batch size $M = 10^4$.

In this high-dimensional regime, significant discrepancies emerge between CGM-relax and CGM-reward (Figure 2B). CGM-relax consistently eliminates the majority of constraint violation up to $k=10^3$, albeit with a non-trivial excess KL divergence to $p_{\theta_{base}}$ compared to the maximum entropy solution p_{α^*} that increases linearly with dimension. Although CGM-reward performs well for low-dimensional constraints (<10), we find that the empirical maximum entropy problem (7) is infeasible with high probability for >30 constraints. In fact, even when $\hat{\alpha}_N$ is fixed to its oracle value α^* (Figure 2B), CGM-relax still outperforms CGM-reward.

4 CASE-STUDIES WITH DIVERSE MODELS, DATA, AND CONSTRAINTS

We evaluate the capacity of CGM-reward and CGM-relax to solve practical calibration problems through three applications involving diverse model, data, and constraint types. Section 4.1 calibrates a diffusion model (Lin et al., 2024a) and a masked language model (Hayes et al., 2025) of protein structure to more closely match statistics of natural proteins. Section 4.2 calibrates a normalizing flow model (Zhai et al., 2025) of images to reduce class imbalances on the basis of LLM image-to-text annotations. Lastly, Section 4.3 calibrates a small autoregressive LM to eliminate gender bias in generated children’s stories (Eldan & Li, 2023).

Across all examples, CGM reduces the majority of calibration error without significantly degrading the quality of generations. Consistent with our results in Section 3 we find that optimally-tuned CGM-relax outperforms CGM-reward, which falls short of meeting the calibration constraints.

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

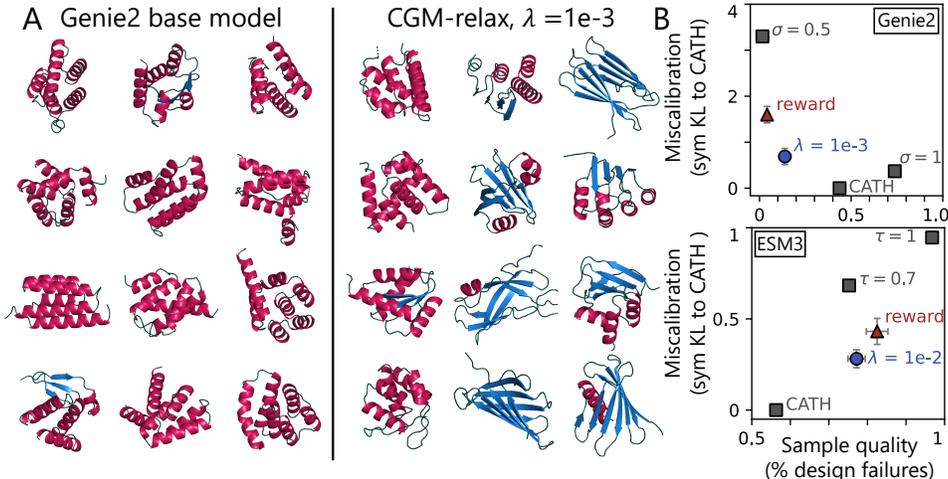


Figure 3: **A:** Samples from the Genie2 protein generative models before and after calibration with CGM-relax ($\lambda=10^{-3}$). **B:** CGM-relax reduces the distance of secondary structure content to natural proteins by >4 times for Genie2 and >2 times for ESM3 while maintaining biophysical plausibility.

Baselines. Only two prior works have proposed algorithms that intend to solve the calibration problem. Khalifa et al. (2021) propose a method for LLMs that we compare to in Section 4.3. Second, Shen et al. (2024) propose a method for class-balancing in diffusion models. However, it assumes an existing probabilistic classifier and so is not applicable in our setting.

Compute cost. Each experiment is run on a single H100 GPU. We provide additional details regarding our experimental setup in Appendix E.

4.1 CALIBRATING PROTEIN DESIGN MODELS TO MATCH STATISTICS OF NATURAL PROTEINS

Diffusion generative models have become a central tool in protein design (Trippe et al., 2023; Watson et al., 2023). However, heuristics such as reduced noise during sampling (see e.g., Yim et al., 2023) have been necessary to ensure a high proportion of the sampled structures are biophysically plausible. These heuristics substantially reduce the diversity of samples compared to proteins found in nature and thereby pose a trade-off between reliability and diversity. For two such protein design models, we investigate whether this trade-off can be mitigated by calibrating the models to match the secondary structure composition of natural proteins.

Protein models Genie2 and ESM3 and their miscalibration. The two protein design models we consider are (1) Genie2 (Lin et al., 2024a), a 15M parameter equivariant diffusion model, and (2) ESM3-open (Hayes et al., 2025), a 1.4B parameter masked language model on tokenized representations of protein backbones. For each model, we generate protein backbones consisting of 100 amino acids i.e., residues. Both Genie2 and ESM3-open suffer low diversity compared to natural protein domains in the CATH dataset (Sillitoe et al., 2021); specifically, they produce few generations with high beta-strand content (Figure 3A). Beta strands, along with alpha helices and loops, constitute what is known as a protein’s secondary structure.

Calibration constraints on secondary structure diversity. To represent protein secondary structure as a calibration constraint, we use the empirical bivariate cumulative density function (CDF) of the fraction of residues in alpha-helical and beta-strand segments. We place up to $d = 99$ cut-off pairs $(\tau_{\alpha,i}, \tau_{\beta,i}) \in [0, 1]^2$ and define a d -dimensional indicator vector $\mathbf{h}(\mathbf{x})$ with components $\mathbf{h}(\mathbf{x})[i] = \mathbb{1}\{f_{\alpha}(\mathbf{x}) \leq \tau_{\alpha,i}, f_{\beta}(\mathbf{x}) \leq \tau_{\beta,i}\}$, $i = 1, \dots, d$, where $f_{\alpha}(\mathbf{x})$ and $f_{\beta}(\mathbf{x})$ are the secondary-structure fractions of protein structure \mathbf{x} . We set the calibration target \mathbf{h}^* to the corresponding values of the CATH empirical bivariate CDF at these cutoffs.

Results. Performing calibration with CGM-relax yields a nearly fivefold improvement in the diversity of sampled protein structures for Genie2 and a twofold improvement for ESM3-open, as quantified by the symmetrized KL distance between the secondary structure distributions of the generative models and CATH domains (Figure 3B). This improvement comes at the cost of an increased proportion of ‘design failures’, as defined in Appendix E.1. The ESM3-open base model



Figure 4: Generations from the conditional TarFlow model (Zhai et al., 2025) before and after calibration with CGM-relax ($\lambda = 10^{-4}$). CGM reweights the proportions of animals generated and produces realistic images. Some visual artifacts exist after calibration (see e.g., fox).

generates a high proportion of design failures compared to Genie2 (consistent with Xiong et al. (2025), for example) and this fraction increases slightly upon calibration with CGM.

CGM-reward achieves more modest improvements in secondary structure diversity, which may in part be due to difficulty in computing $\hat{\alpha}_N$. In order for equation (7) to be feasible with $N = 2.5 \times 10^3$ samples, we need to reduce the number of cutoff pairs from 99 to 15. CGM-reward fine-tuning reduces the symmetrized KL distance to CATH by two times for Genie2 and 1.6 times for ESM3-open. However, for Genie2, CGM-reward also produces fewer design failures than CGM-relax.

The gains in secondary structure diversity achieved by CGM cannot be obtained by simply increasing the sampling noise of Genie2 or the sampling temperature of ESM3. In Figure 3B, we show that increasing the sampling noise of Genie2 to $\sigma = 1$ improves structure diversity, but at the cost of 5.3 times more design failures (failure rate 74%) than CGM. The same is true for ESM3 with increased sampling temperature $\tau = 1$, which yields a 1.3 times higher failure rate of 97%.

4.2 CALIBRATING CLASS PROPORTIONS IN A CONDITIONAL FLOW MODEL

We next demonstrate that CGM is capable of effectively calibrating state-of-the-art normalizing flow models. Normalizing flows generate samples $\mathbf{x} \in \mathbb{R}^k$ according to $\mathbf{x} = f_\theta^{-1}(\epsilon)$, where $\epsilon \sim p_\epsilon$ is a distribution from which sampling is tractable and $f_\theta(\mathbf{x})$ is a map that is invertible in \mathbf{x} for each θ (Tabak & Vanden-Eijnden, 2010; Rezende & Mohamed, 2015). By the change-of-variable formula, the density of \mathbf{x} is $p_\epsilon(f_\theta(\mathbf{x}))|\det(df_\theta(\mathbf{x})/d\mathbf{x})|$. This expression enables computation of exact likelihoods for maximum likelihood training and calibration.

For our calibration problem, we consider the 463M-parameter TarFlow model (Zhai et al., 2025), which parameterizes f_θ as an autoregressive vision transformer (Dosovitskiy et al., 2021) such that attention is performed over a sequence of image patches. We examine the model trained conditionally on the 256×256 AFHQ dataset (Choi et al., 2020), which consists of images of animal faces belonging to one of three classes: {cat, dog, wildlife}. The wildlife class is further comprised of {lion, tiger, fox, wolf, cheetah, leopard}. We observe that, conditional on the wildlife class, approximately 36% of generations from the TarFlow model are lions and very few ($< 7\%$ total) are foxes or wolves. We apply CGM to calibrate the conditional TarFlow model to generate samples containing animals from the wildlife class with equal proportions. For \mathbf{h} , we query GPT o5-mini to classify each image as containing one of the six animals or None.

Results. We find CGM-relax reduces miscalibration to the base TarFlow model with little visible degradation of sample realism (Figure 4). CGM-relax ($\lambda=10^{-4}$) reduces the total variation distance of animal proportions, as classified by an image-to-text model, to the uniform distribution from 0.306 to 0.108. However, the Fréchet inception distance (FID) to real images in the AFHQ wildlife class is larger for the calibrated model than for the base model (21.9 vs. 15.9). Since this metric is sensitive to class proportions, we evaluate the calibrated model on the training dataset after balancing classes. The discrepancy in FID can be explained by two types of visual artifacts introduced by calibration: some images depict animals outside the wildlife class ($\sim 8\%$) and some “blend” multiple animals. Appendix Figure 9 shows random samples from both models. The model fine-tuned with CGM-reward remains close to the base model but fails to reduce constraint violation.

4.3 ELIMINATING PROFESSION-SPECIFIC GENDER IMBALANCE IN CHILDREN’S STORIES

As a third example, we calibrate a language model that generates short children’s stories to remove gender bias. TinyStories-33M is an autoregressive transformer trained on children’s stories

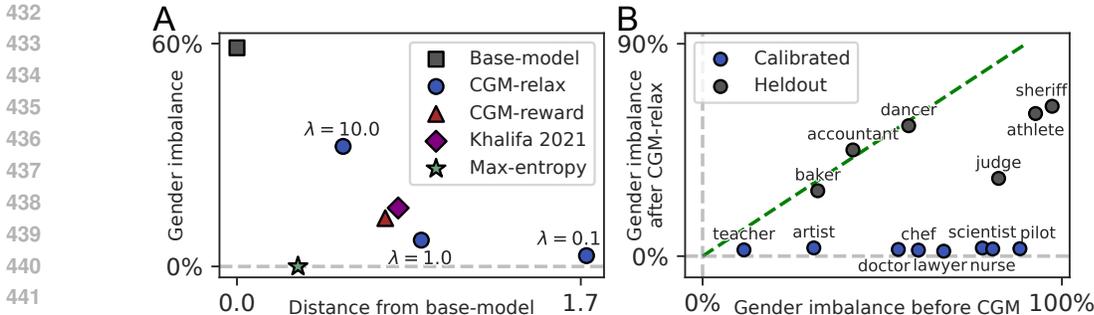


Figure 5: **A**: Gender imbalance and distance from base-model (symmetrized KL from pre-trained TinyStories-33M). **B**: Gender imbalance for professions included and heldout from calibration before and after CGM-relax ($\lambda = 0.1$). Points below the diagonal were improved by CGM.

generated by GPT-3.5 and GPT4 (Eldan & Li, 2023). We find significant imbalances in prompt-conditional generations that introduce a character’s profession. For example, only 16% of stories beginning “Once upon a time there was a lawyer” feature a female lawyer, whereas 41% of U.S. attorneys were women in 2024 (American Bar Association, 2024).

Gender parity as a calibration constraint and conditional calibration. We evaluate whether CGM can eliminate profession-specific gender imbalance in stories completed from the prompt “Once upon a time there was a <profession>” across eight professions that exhibit gender bias under the base model: doctor, lawyer, teacher, pilot, chef, scientist, nurse, and artist. In contrast to earlier experiments, this requires *conditional* calibration: for each profession i with prompt prompt_i , we aim to find θ such that $\mathbb{E}_{p_\theta}[\mathbf{h}(\mathbf{x}) \mid \text{prompt}_i] = 0$, where \mathbf{x} represents a completed story, and $\mathbf{h}(\mathbf{x}) \in \{-1, 0, 1\}$ encodes the character’s gender (male, ambiguous, or female, respectively). Rather than fine-tuning a separate model for each profession, we amortize training costs by fine-tuning a single model with the sum of CGM losses for each condition.

Results on explicitly calibrated professions. Both CGM-reward and CGM-relax reduce gender imbalance, as measured by the average absolute per-profession frequency difference (Figure 5A). As expected, decreasing the regularization strength λ improves constraint satisfaction at the cost of greater distance to the base-model, as measured by symmetrized KL. Notably, even the least-regularized model attains a low symmetrized KL of 1.7, which corresponds to an average token log-probability difference of < 0.01 nats/token. Appendix E.4.4 provides example generations before and after fine-tuning showing no visible degradation in story quality.

Compared to Khalifa et al. (2021), CGM-reward yields a small but statistically significant improvement in both miscalibration and distance to the base model. CGM-relax reduces gender imbalance by over five times more than Khalifa et al. (2021) but deviates further from the base-model.

Transference of calibration to heldout professions. We evaluate how conditional calibration affects the calibration of “held-out” professions not considered during fine-tuning. Such generalization could be particularly valuable in applications where it is impractical to foresee and explicitly calibrate for every possible prompt. To evaluate this, we consider six held-out professions: sheriff, judge, accountant, dancer, athlete, and baker. While CGM does not result in gender parity for the held-out professions, the imbalance is significantly reduced for some (Figure 5B).

5 CONCLUSION

CGM-relax and CGM-reward provide practical approaches for calibrating generative models to satisfy distribution-level constraints. In applications to protein design, conditional image generation, and language modeling, CGM consistently reduces calibration error under hundreds of simultaneous constraints and in models with up to one billion parameters while preserving generation quality.

Still, our results highlight that the calibration problem is not yet solved. Current objectives leave residual error, especially in the rare-event setting that is especially relevant to protein structure modeling, for example. More broadly, the CGM framework is tied to models with tractable likelihoods, raising the challenge of extending calibration to VAEs, GANs, and other implicit models. These open questions point to calibration as a practical tool as well as a fruitful research direction.

6 ETHICS STATEMENT

This work develops algorithms for calibrating generative models by aligning distribution-level statistics to desired targets. Our motivation is to improve the fidelity of generative models across diverse domains, including, but not limited to, protein design, image generation, and language modeling. Potential ethical benefits include reducing harmful biases (e.g., gender imbalance in text outputs) and improving scientific utility (e.g., protein structure design). However, as is the case for all works that fine-tune generative models, our methods could also be misused to enforce constraints that amplify harmful or discriminatory content. We emphasize that the choice of constraints should be made responsibly, with careful attention to societal and scientific impacts.

All datasets used in this work are publicly available, and no sensitive personal data were employed.

7 REPRODUCIBILITY STATEMENT

We have taken multiple steps to ensure reproducibility of our theoretical and empirical results. **Theory.** All mathematical claims are supported by detailed theorem statements and proofs in Appendices B and C, and assumptions for all claims are clearly stated. **Algorithms.** We provide complete pseudocode for the algorithms we propose (Algorithms 1 and 2), including clear descriptions of loss estimation and gradient computation. Our implementations are entirely reproducible from these algorithms. **Experiments.** For each of our experiments in Sections 3 and 4, we specify the datasets, models, and calibration constraints in detail. Hyperparameter choices (e.g., model architecture, optimizer, learning rate, number of epochs, batch size M , regularization strength λ , sample size N) are reported in Appendices D and E. We include additional samples from the pre-trained (i.e., base) and fine-tuned (i.e., calibrated) generative models in these appendices. **Code.** Upon publication, we will release a public codebase implementing CGM-relax and CGM-reward for arbitrary generative models. We will include scripts to reproduce experimental results.

REFERENCES

- American Bar Association. ABA profile of the legal profession 2024 — demographics, 2024.
- Brian DO Anderson. Reverse-time diffusion equation models. *Stochastic Processes and their Applications*, 1982.
- Sanjeev Arora and Yi Zhang. Do GANs actually learn the distribution? An empirical study. *arXiv preprint arXiv:1706.08224*, 2017.
- Aharon Ben-Tal and Arkadi Nemirovski. Lecture notes in convex analysis, nonlinear programming theory, and nonlinear programming algorithms, 2023.
- Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning. In *International Conference on Learning Representations*, 2024.
- Jonathan M Borwein and Qiji J Zhu. *Techniques of variational analysis*. Springer, 2005.
- Sandro Bottaro, Tone Bengtsen, and Kresten Lindorff-Larsen. Integrating molecular simulation and experimental data: a Bayesian/maximum entropy reweighting approach. *Structural Bioinformatics: Methods and Protocols*, 2020.
- Stephen P Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- Robert H Cameron and William T Martin. Transformations of Wiener integrals under translations. *Annals of Mathematics*, 1944.
- Michael Cardei, Jacob K Christopher, Thomas Hartvigsen, Brian R Bartoldson, Bhavya Kailkhura, and Ferdinando Fioretto. Constrained discrete diffusion. *arXiv preprint arXiv:2503.09790*, 2025.

- 540 Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis
541 for multiple domains. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
542 *Recognition*, 2020.
- 543 Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep
544 reinforcement learning from human preferences. *Advances in Neural Information Processing*
545 *Systems*, 2017.
- 546 Imre Csiszár. I-divergence geometry of probability distributions and minimization problems. *The*
547 *Annals of Probability*, 1975.
- 548 Justas Dauparas, Ivan Anishchenko, Nathaniel Bennett, Hua Bai, Robert J. Ragotte, Lukas F. Milles,
549 Basile I. M. Wicky, Alexis Courbet, Rob J. de Haas, Neville Bethel, Philip J. Y. Leung, Timothy F.
550 Huddy, Sam Pellock, Doug Tischer, F. Chan, Brian Koepnick, H. Nguyen, A. Kang, Banumathi
551 Sankaran, Asim K. Bera, Neil P. King, and David Baker. Robust deep learning-based protein
552 sequence design using ProteinMPNN. *Science*, 2022.
- 553 A Philip Dawid. The well-calibrated Bayesian. *Journal of the American Statistical Association*,
554 1982.
- 555 Alexander Denker, Francisco Vargas, Shreyas Padhy, Kieran Didi, Simon Mathis, Riccardo Barbano,
556 Vincent Dutordoir, Emile Mathieu, Urszula Julia Komorowska, and Pietro Lio. DEFT: efficient
557 fine-tuning of diffusion models by learning the generalised h -transform. *Advances in Neural*
558 *Information Processing Systems*, 2024.
- 559 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. *Ad-*
560 *vances in Neural Information Processing Systems*, 2021.
- 561 Carles Domingo-Enrich, Michal Drozdal, Brian Karrer, and Ricky TQ Chen. Adjoint matching:
562 fine-tuning flow and diffusion generative models with memoryless stochastic optimal control. In
563 *International Conference on Learning Representations*, 2025.
- 564 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
565 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An im-
566 age is worth 16x16 words: transformers for image recognition at scale. In *International Confer-*
567 *ence on Learning Representations*, 2021.
- 568 Ronen Eldan and Yuanzhi Li. TinyStories: How small can language models be and still speak
569 coherent English? *arXiv preprint arXiv:2305.07759*, 2023.
- 570 Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel,
571 Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. DPOK: Reinforcement learning for
572 fine-tuning text-to-image diffusion models. *Advances in Neural Information Processing Systems*,
573 2023.
- 574 Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha
575 Luccioni, and Kristian Kersting. Fair Diffusion: Instructing text-to-image generation models on
576 fairness. *arXiv preprint arXiv:2302.10893*, 2023.
- 577 Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernon-
578 court, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models:
579 A survey. *Computational Linguistics*, 2024.
- 580 Igor Vladimirovich Girsanov. On transforming a certain class of stochastic processes by absolutely
581 continuous substitution of measures. *Theory of Probability & Its Applications*, 1960.
- 582 Dongyoung Go, Tomasz Korbak, Germán Kruszewski, Jos Rozen, Nahyeon Ryu, and Marc Dymet-
583 man. Aligning language models with preferences through f-divergence minimization. In *Internat-*
584 *ional Conference on Machine Learning*, 2023.
- 585 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural
586 networks. In *International Conference on Machine Learning*, 2017.

- 594 Sven Gutjahr, Riccardo De Santi, Luca Schaufelberger, Kjell Jorner, and Andreas Krause. Con-
595 strained molecular generation via sequential flow model fine-tuning. In “*Generative AI and Biol-*
596 *ogy (GenBio)*” *Workshop the International Conference on Machine Learning*, 2025.
- 597
598 Thomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert
599 Verkuil, Vincent Q Tran, Jonathan Deaton, Marius Wiggert, Rohil Badkundri, Irhum Shafkat,
600 Jun Gong, Alexander Derry, Raul S Molina, Neil Thomas, Yousuf A Khan, Chetan Mishra, Car-
601 olyn Kim, Liam J Bartie, Matthew Nemeth, Patrick D Hsu, Tom Sercu, Salvatore Candido, and
602 Alexander Rives. Simulating 500 million years of evolution with a language model. *Science*,
603 2025.
- 604 Magnus R. Hestenes. Multiplier and gradient methods. *Journal of Optimization Theory and Appli-*
605 *cations*, pp. 303–320, 1969.
- 606 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. “*Deep Generative Models and*
607 *Downstream Applications*” *Workshop at the Advances in Neural Information Processing Systems*
608 *Conference*, 2021.
- 609
610 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
611 *Neural Information Processing Systems*, 2020.
- 612 John Ingraham, Vikas Garg, Regina Barzilay, and Tommi Jaakkola. Generative models for graph-
613 based protein design. *Advances in Neural Information Processing Systems*, 2019.
- 614
615 Edwin T Jaynes. Information theory and statistical mechanics. *Physical Review*, 1957.
- 616 Shervin Khalafi, Dongsheng Ding, and Alejandro Ribeiro. Constrained diffusion models via dual
617 training. *Advances in Neural Information Processing Systems*, 2024.
- 618
619 Muhammad Khalifa, Hady Elsahar, and Marc Dymetman. A distributional approach to controlled
620 text generation. In *International Conference on Learning Representations*, 2021.
- 621 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *International*
622 *Conference on Learning Representations*, 2015.
- 623
624 Yuichi Kitamura and Michael Stutzer. An information-theoretic alternative to generalized method
625 of moments estimation. *Econometrica: Journal of the Econometric Society*, 1997.
- 626 Jürgen Köfinger, Lukas S. Stelzl, Klaus Reuter, César Allande, Katrin Reichel, and Gerhard Hum-
627 mer. Efficient ensemble refinement by reweighting. *Journal of Chemical Theory and Computa-*
628 *tion*, 2019.
- 629
630 Wouter Kool, Herke van Hoof, and Max Welling. Buy 4 REINFORCE samples, get a baseline for
631 free! In “*Deep RL Meets Structured Prediction*” *Workshop at the International Conference on*
632 *Learning Representations*, 2019.
- 633 Solomon Kullback. *Information Theory and Statistics*. John Wiley & Sons, 1959.
- 634
635 Patrick Kunzmann and Kay Hamacher. Biotite: a unifying open source computational biology
636 framework in Python. *BMC Bioinformatics*, 2018.
- 637
638 Stephen S Lavenberg and Peter D Welch. A perspective on the use of control variables to increase
639 the efficiency of Monte Carlo simulations. *Management Science*, 1981.
- 640 Sarah Lewis, Tim Hempel, José Jiménez-Luna, Michael Gastegger, Yu Xie, Andrew Y. K. Foong,
641 Victor García Satorras, Osama Abdin, Bastiaan S. Veeling, Iryna Zaporozhets, Yaoyi Chen, Soo-
642 jung Yang, Adam E. Foster, Arne Schneuing, Jigyasa Nigam, Federico Barbero, Vincent Stimper,
643 Andrew Campbell, Jason Yim, Marten Lienen, Yu Shi, Shuxin Zheng, Hannes Schulz, Usman
644 Munir, Roberto Sordillo, Ryota Tomioka, Cecilia Clementi, and Frank Noé. Scalable emulation
645 of protein equilibrium ensembles with generative deep learning. *Science*, 2025.
- 646 Sarah Lichtenstein, Baruch Fischhoff, and Lawrence D Phillips. Calibration of probabilities: The
647 state of the art. In *Decision Making and Change in Human Affairs: Proceedings of the Research*
Conference on Subjective Probability, Utility, and Decision Making, 1977.

- 648 Yeqing Lin, Minji Lee, Zhao Zhang, and Mohammed AlQuraishi. Out of many, one: designing
649 and scaffolding proteins at the scale of the structural universe with Genie 2. *arXiv preprint*
650 *arXiv:2405.15489*, 2024a.
- 651 Yeqing Lin, Haewon C. Nguyen, and Mohammed AlQuraishi. In-silico protein design pipeline.
652 github.com/aqlaboratory/insilico_design_pipeline, 2024b.
- 654 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
655 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom
656 Sercu, Salvatore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level
657 protein structure with a language model. *Science*, 2023.
- 658 Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. *arXiv*
659 *preprint arXiv:1608.03983*, 2016.
- 660 Tianyu Lu, Melissa Liu, Yilin Chen, Jinho Kim, and Po-Ssu Huang. Assessing generative model
661 coverage of protein structures with SHAPES. *bioRxiv*, 2025.
- 662 Shintaro Minami. PyDSSP. github.com/ShintaroMinami/PyDSSP, 2023.
- 663 Shakir Mohamed, Mihaela Rosca, Michael Figurnov, and Andriy Mnih. Monte Carlo gradient
664 estimation in machine learning. *Journal of Machine Learning Research*, 2020.
- 665 Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated proba-
666 bilities using Bayesian binning. In *AAAI Conference on Artificial Intelligence*, 2015.
- 667 Matthew O’Kelly, Aman Sinha, Hongseok Namkoong, Russ Tedrake, and John C Duchi. Scalable
668 end-to-end autonomous vehicle testing via rare-event simulation. *Advances in Neural Information*
669 *Processing Systems*, 2018.
- 670 Bernt Oksendal. *Stochastic Differential Equations: An Introduction with Applications*. Springer
671 Science & Business Media, 2013.
- 672 Art B Owen. *Empirical likelihood*. Chapman and Hall/CRC, 2001.
- 673 Israel Saeta Pérez, David Arcos, and LeadRatings contributors. gender-guesser. [github.com/](https://github.com/lead-ratings/gender-guesser)
674 [lead-ratings/gender-guesser](https://github.com/lead-ratings/gender-guesser), 2016.
- 675 John Platt. Probabilistic outputs for support vector machines and comparisons to regularized likeli-
676 hood methods. In *Advances in Large Margin Classifiers*, 1999.
- 677 Jin Qin and Jerry Lawless. Empirical likelihood and general estimating equations. *The Annals of*
678 *Statistics*, 1994.
- 679 Yiming Qin, Huangjie Zheng, Jiangchao Yao, Mingyuan Zhou, and Ya Zhang. Class-balancing
680 diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
681 *Recognition*, 2023.
- 682 Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea
683 Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances*
684 *in Neural Information Processing Systems*, 2023.
- 685 Rajesh Ranganath, Sean Gerrish, and David Blei. Black box variational inference. In *Artificial*
686 *Intelligence and Statistics*. PMLR, 2014.
- 687 Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *Internat-*
688 *ional Conference on Machine Learning*, 2015.
- 689 R Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970.
- 690 R Tyrrell Rockafellar. Augmented lagrangians and applications of the proximal point algorithm in
691 convex programming. *Mathematics of Operations research*, 1976.
- 692 Bartosz Różycki, Young C Kim, and Gerhard Hummer. SAXS ensemble refinement of ESCRT-III
693 CHMP3 conformational transitions. *Structure*, 2011.

- 702 Anirban Sarkar, Yijie Kang, Nirali Somia, Pablo Mantilla, Jessica Lu Zhou, Masayuki Nagai, Ziqi
703 Tang, Chris Zhao, and Peter Koo. Designing DNA with tunable regulatory activity using score-
704 entropy discrete diffusion. *bioRxiv*, 2024.
- 705
706 Glenn Shafer and Vladimir Vovk. A tutorial on conformal prediction. *Journal of Machine Learning*
707 *Research*, 2008.
- 708 Xudong Shen, Chao Du, Tianyu Pang, Min Lin, Yongkang Wong, and Mohan Kankanhalli. Fine-
709 tuning text-to-image diffusion models for fairness. In *International Conference on Learning Rep-*
710 *resentations*, 2024.
- 711
712 Ian Sillitoe, Nicola Bordin, Natalie Dawson, Vaishali P. Waman, Paul Ashford, Harry M. Scholes,
713 Camilla S. M. Pang, Laurel Woodridge, Clemens Rauer, Neeladri Sen, Mahnaz Abbasian, Sean
714 Le Cornu, Su Datt Lam, Karel Berka, Ivana Hutařová Vareková, Radka Svobodova, Jon Lees,
715 and Christine A. Orengo. CATH: increased structural coverage of functional space. *Nucleic*
716 *Acids Research*, 2021.
- 717 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
718 Poole. Score-based generative modeling through stochastic differential equations. In *Internat-*
719 *ional Conference on Learning Representations*, 2021.
- 720
721 Esteban G Tabak and Eric Vanden-Eijnden. Density estimation by dual ascent of the log-likelihood.
722 *Communications in Mathematical Sciences*, 2010.
- 723 Wenpin Tang. Fine-tuning of diffusion models via stochastic control: entropy regularization and
724 beyond. *arXiv preprint arXiv:2403.06279*, 2024.
- 725
726 Brian L Trippe, Jason Yim, Doug Tischer, David Baker, Tamara Broderick, Regina Barzilay, and
727 Tommi Jaakkola. Diffusion probabilistic modeling of protein backbones in 3D for the motif-
728 scaffolding problem. In *International Conference on Learning Representations*, 2023.
- 729 Masatoshi Uehara, Yulai Zhao, Kevin Black, Ehsan Hajiramezanali, Gabriele Scalia, Nathaniel Lee
730 Diamant, Alex M Tseng, Tommaso Biancalani, and Sergey Levine. Fine-tuning of continuous-
731 time diffusion models as entropy-regularized control. *arXiv preprint arXiv:2402.15194*, 2024.
- 732
733 Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas
734 Schön. Evaluating model calibration in classification. In *International Conference on Artificial*
735 *Intelligence and Statistics*, 2019.
- 736
737 Aad W Van der Vaart. *Asymptotic statistics*. Cambridge University Press, 2000.
- 738
739 Martin J Wainwright and Michael I Jordan. Graphical models, exponential families, and variational
740 inference. *Foundations and Trends in Machine Learning*, 2008.
- 741
742 Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam,
743 Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using
744 direct preference optimization. In *Conference on Computer Vision and Pattern Recognition*. IEEE
745 Computer Society, 2024.
- 746
747 Joseph L. Watson, David Juergens, Nathaniel R. Bennett, Brian L. Trippe, Jason Yim, Helen E.
748 Eisenach, Woody Ahern, Andrew J. Borst, Robert J. Ragotte, Lukas F. Milles, Basile I. M.
749 Wicky, Nikita Hanikel, Samuel J. Pellock, Alexis Courbet, William Sheffler, Jue Wang, Preetham
750 Venkatesh, Isaac Sappington, Susana Vázquez Torres, Anna Lauko, Valentin De Bortoli, Emile
751 Mathieu, Sergey Ovchinnikov, Regina Barzilay, Tommi S. Jaakkola, Frank DiMaio, Minkyung
752 Baek, and David Baker. De novo design of protein structure and function with RFdiffusion.
753 *Nature*, 2023.
- 754
755 Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement
756 learning. *Machine Learning*, 1992.
- 757
758 Jin-Long Wu, Karthik Kashinath, Adrian Albert, Dragos Chirila, and Heng Xiao. Enforcing sta-
759 tistical constraints in generative adversarial networks for modeling chaotic dynamical systems.
760 *Journal of Computational Physics*, 2020.

756 Junhao Xiong, Hunter Nisonoff, Maria Lukarska, Ishan Gaur, Luke M Oltrogge, David F Savage,
757 and Jennifer Listgarten. Guide your favorite protein sequence generative model. *arXiv preprint*
758 *arXiv:2505.04823*, 2025.

759
760 Jason Yim, Brian L Trippe, Valentin De Bortoli, Emile Mathieu, Arnaud Doucet, Regina Barzilay,
761 and Tommi Jaakkola. SE(3) diffusion model with application to protein backbone generation. In
762 *International Conference on Machine Learning*, 2023.

763 Shuangfei Zhai, Ruixiang Zhang, Preetum Nakkiran, David Berthelot, Jiatao Gu, Huangjie Zheng,
764 Tianrong Chen, Miguel Ángel Bautista, Navdeep Jaitly, and Joshua M Susskind. Normalizing
765 flows are capable generative models. In *International Conference on Machine Learning*, 2025.

766
767 Danqing Zhu, David H. Brookes, Akosua Busia, Ana Carneiro, Clara Fannjiang, Galina Popova,
768 David Shin, Kevin C. Donohue, Li F. Lin, Zachary M. Miller, Evan R. Williams, Edward F. Chang,
769 Tomasz J. Nowakowski, Jennifer Listgarten, and David V. Schaffer. Optimal trade-off control in
770 machine learning-based library design, with application to adeno-associated virus (AAV) for gene
771 therapy. *Science Advances*, 2024.

772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

810	APPENDIX CONTENTS
811	
812	Appendix A: Extended Discussion of Related Work
813	
814	Appendix B: CGM-relax and CGM-reward Algorithms
815	B.1: Loss Estimates
816	B.2: Unbiased Gradient Estimates
817	
818	Appendix C: Maximum Entropy Principle
819	
820	C.1: Precise Statement
821	C.2: Estimating the Maximum Entropy Solution
822	C.3: Connection Between the Relax and Reward Losses
823	C.4: Consistency and Asymptotic Normality
824	
825	Appendix D: Simulation Experiments Additional Details
826	
827	D.1: Continuous-time Diffusion Models
828	D.2: Initializing the Base Diffusion Model
829	D.3: Experimental Details
830	D.4: Comparison to Augmented Lagrangian Method
831	
832	Appendix E: Additional Experimental Details
833	
834	E.1: Calibrating Genie2
835	E.2: Calibrating ESM3-open
836	E.3: Calibrating TarFlow
837	E.4: Calibrating TinyStories-33M
838	
839	
840	
841	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
854	
855	
856	
857	
858	
859	
860	
861	
862	
863	

864 A EXTENDED DISCUSSION OF RELATED WORK

865
866
867
868 **The calibration problem.** Several previous works have proposed algorithms whose goal it is to
869 impose distributional constraints on generative models. However, each of these methods applies only
870 to specific model classes and either suffers from poor empirical performance or imposes constraint
871 satisfaction during training time (rather than fine-tuning).

872 Most closely related to the present work, Khalifa et al. (2021) fine-tune autoregressive language
873 models to match distributional constraints. Like CGM-reward, their approach also targets the maxi-
874 mum entropy solution (5), but through a different divergence; they choose the KL divergence in the
875 “forward” direction, $D_{\text{KL}}(p_{\alpha^*} \parallel p_{\theta})$, rather than in the “reverse” direction, $D_{\text{KL}}(p_{\theta} \parallel p_{\alpha^*})$, as in
876 CGM-reward.

877 Empirically, the approximate solutions to the calibration problem (1) found by Khalifa et al. (2021)
878 fall shorter of constraint satisfaction compared to CGM, particularly CGM-relax. Khalifa et al.
879 (2021) achieves comparable, albeit slightly worse, performance to CGM-reward in the TinyStories
880 gender rebalancing experiment (Section 4.3), reducing miscalibration by roughly 85%. CGM-relax,
881 on the other hand, reduces constraint violation up to 98%.

882 In follow-up work, Go et al. (2023) propose an algorithm for aligning language models to a specified
883 target distribution by minimizing an arbitrary f -divergence (including the forward and reverse KL
884 divergence). One example they consider is when the target distribution is the maximum entropy
885 distribution corresponding to some constraint functions; the choice of forward KL then reduces to
886 Khalifa et al. (2021). However, they obtain $< 50\%$ reduction in constraint violation.

887 Shen et al. (2024) proposes a method for balancing class proportions in text-to-image diffusion
888 models. They rely on an optimal transport objective that applies narrowly to diffusion models and
889 find empirically their approach falls short of meeting desired class proportions.

890 In concurrent work, Cardei et al. (2025) impose constraints on discrete diffusion models at sampling
891 time using an augmented Lagrangian method. Their algorithm involves simultaneously optimizing
892 the model output and a set of Lagrange multipliers. Also concurrent to our work, Gutjahr et al.
893 (2025) fine-tunes a diffusion generative model subject to inequality constraints on the expected
894 value of a statistic to maximize an expected reward with a KL penalty to the base model. Their
895 approach applies only to diffusion models and continuous normalizing flows.

896 **Incorporating distributional constraints during training.** Several other works have sought to
897 impose distributional constraints during training time but differ from CGM in that they are not fine-
898 tuning procedures and apply only to a specific model classes. Wu et al. (2020) propose a method
899 for training generative adversarial networks (GANs) that includes a penalty term similar to $\mathcal{L}^{\text{viol}}$
900 that encourages agreement with statistics of the training data. Zhu et al. (2024) solve for the maxi-
901 mum entropy model of short (length 7) protein sequences with expected “fitness” surpassing a fixed
902 threshold. Khalafi et al. (2024) propose a primal-dual algorithm to enforce distributional constraints
903 on diffusion models; their constraints, however, are specified at the level of entire distributions,
904 rather than their moments. Friedrich et al. (2023) develop a training procedure for diffusion models
905 that balances the conditional distributions of samples, given some attribute e.g., gender.

906 **Reward fine-tuning and conditional generation.** As we point out in Section 2.2, the idea of
907 minimizing the KL divergence of the generative model to an exponential tilt of the base model (5)
908 connects CGM to the rich research topic of reward fine-tuning. Reward fine-tuning algorithms, used
909 in the contexts of reinforcement learning (Rafailov et al., 2023; Fan et al., 2023; Black et al., 2024;
910 Wallace et al., 2024) and preference optimization (Tang, 2024; Uehara et al., 2024; Domingo-Enrich
911 et al., 2025), minimize the same loss (8) as CGM-reward, but with $r_{\alpha}(\mathbf{x})$ replaced by a user-specified
912 “reward”. Unlike reward fine-tuning algorithms, though, CGM does not require a reward; rather, the
913 constraints themselves act as the reward.

914 Conditional generation (Dhariwal & Nichol, 2021; Ho & Salimans, 2021; Denker et al., 2024) can
915 also be viewed through the lens of model calibration, where the calibration constraint is the indicator
916 function of the set C from which one would like to sample $\mathbf{h}(\mathbf{x}) = \mathbb{1}\{\mathbf{x} \in C\}$ and \mathbf{h}^* , the target
917 proportion of samples that belong to C , approaches 1. In this case the optimal variational parameter
 α^* approach infinity, and the maximum entropy solution approaches $p_{\theta_{\text{base}}}(\mathbf{x})\mathbb{1}\{\mathbf{x} \in C\}$.

Calibration of molecular ensembles. Computational methods for producing Boltzmann ensembles frequently fail to exactly align with experimental observables that measure ensemble averages; this misalignment can arise from inaccuracies in the energy functions used or insufficient sampling. Several works have sought to calibrate these ensembles to agree with ensemble observables. In the context of molecular dynamics simulations, (Różycki et al., 2011; Köfinger et al., 2019; Bottaro et al., 2020) leverage Theorem 2.1 to reweight Monte Carlo samples of molecular configurations to match experimental observations of ensemble averages. Lewis et al. (2025) consider a diffusion generative model approximation of protein structure ensembles and introduce an auxiliary training loss that resembles $\mathcal{L}^{\text{viol}}$, but they do not demonstrate whether this approach leads to a significant reduction in calibration error.

Calibration in prediction problems. Beyond generative modeling, calibration is a major topic in supervised machine-learning. In the context of classification the goal is to have, among a collection of predictions with a given class probability, the fraction of labels of that class in agreement with that prediction probability (Dawid, 1982). This can be obtained with post-hoc calibration procedures such as Platt scaling (Platt, 1999) or conformal methods (Shafer & Vovk, 2008) for more general prediction sets.

B CGM-RELAX AND CGM-REWARD ALGORITHMS

In this section, we provide further detail on the CGM-relax and CGM-reward algorithms. First, we show in Appendix B.1 that our estimates for the relax and reward losses are unbiased. In Appendix B.2 we then discuss how to compute our gradient estimates for the relax and reward losses, and we show they are unbiased.

Throughout this section we will make the following regularity assumptions on the generative model p_θ and the constraint functions \mathbf{h} .

Assumption B.1 (Regularity of p_θ). The functions $p_{\hat{\theta}}(\mathbf{x})/p_\theta(\mathbf{x})$, $\nabla_{\hat{\theta}} p_{\hat{\theta}}(\mathbf{x})/p_\theta(\mathbf{x})$, $\log p_{\hat{\theta}}(\mathbf{x})$, $\nabla_{\hat{\theta}} \log p_{\hat{\theta}}(\mathbf{x})$ are uniformly dominated by a function that is square integrable with respect to $p_\theta(\mathbf{x})$, for all $\hat{\theta}$ belonging to some neighborhood of θ . Also, $\mathbf{h}(\mathbf{x})$, $\log p_{\theta_{\text{base}}}(\mathbf{x})$ have finite second moment under $p_\theta(\mathbf{x})$.

These assumptions are sufficient to exchange integration and differentiation in Appendix B.2 with dominated convergence.

B.1 LOSS ESTIMATES

We begin by proving that our estimates $\hat{\mathcal{L}}^{\text{relax}}$ and $\hat{\mathcal{L}}^{\text{reward}}$ for $\mathcal{L}^{\text{relax}}$ and $\mathcal{L}^{\text{reward}}$, respectively, are, on average, correct.

Proposition B.2. $\hat{\mathcal{L}}^{\text{relax}}$ is unbiased for the relax loss $\mathcal{L}^{\text{relax}}$.

Proof. We prove unbiasedness of $\hat{\mathcal{L}}^{\text{relax}}$ by showing that $\hat{\mathcal{L}}^{\text{KL}}$ is unbiased for $\mathcal{L}^{\text{KL}} = \text{D}_{\text{KL}}(p_\theta \parallel p_{\theta_{\text{base}}})$ and that $\hat{\mathcal{L}}^{\text{viol}}$ is unbiased for $\mathcal{L}^{\text{viol}} = \|\mathbb{E}_{p_\theta}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2$.

As for $\hat{\mathcal{L}}^{\text{KL}}$, its expectation is

$$\mathbb{E}_{p_\theta} \left[\hat{\mathcal{L}}^{\text{KL}} \right] = \frac{1}{M} \sum_{m=1}^M \mathbb{E}_{p_\theta} \left[\log \frac{p_\theta(\mathbf{x}_m)}{p_{\theta_{\text{base}}}(\mathbf{x}_m)} \right] = \frac{1}{M} \sum_{m=1}^M \text{D}_{\text{KL}}(p_\theta \parallel p_{\theta_{\text{base}}}) = \text{D}_{\text{KL}}(p_\theta \parallel p_{\theta_{\text{base}}}).$$

In the first equality we invoke the linearity of expectation and in the second we use our assumption that $\{\mathbf{x}_m\}_{m=1}^M$ are sampled from p_θ .

And for $\hat{\mathcal{L}}^{\text{viol}}$, we recall that for a real-valued random variable Z , $\mathbb{E}[Z^2] = \mathbb{E}[Z]^2 + \text{Var}(Z)$. Applying this to each dimension of $M^{-1} \sum_{m=1}^M \tilde{\mathbf{h}}_m = M^{-1} \sum_{m=1}^M (\mathbf{h}(\mathbf{x}_m) - \mathbf{h}^*)$, we obtain

$$\mathbb{E}_{p_\theta} \left\| \frac{1}{M} \sum_{m=1}^M \tilde{\mathbf{h}}_m \right\|^2 = \|\mathbb{E}_{p_\theta}[\tilde{\mathbf{h}}(\mathbf{x})]\|^2 + \frac{1}{M} \mathbb{E}_{p_\theta} \|\tilde{\mathbf{h}}(\mathbf{x}) - \mathbb{E}_{p_\theta}[\tilde{\mathbf{h}}(\mathbf{x})]\|^2, \quad (10)$$

where $\tilde{\mathbf{h}}(\mathbf{x}) = \mathbf{h}(\mathbf{x}) - \mathbf{h}^*$. Next, we replace the final term in (10) with $\mathbb{E}_{p_\theta}[M^{-1}(M-1)^{-1} \sum_m \|\tilde{\mathbf{h}}_m - M^{-1} \sum_{m'} \tilde{\mathbf{h}}_{m'}\|^2]$. The quantity $M^{-1}(M-1)^{-1} \sum_m \|\tilde{\mathbf{h}}_m - M^{-1} \sum_{m'} \tilde{\mathbf{h}}_{m'}\|^2$ is simply the trace of the sample covariance matrix of $\{\tilde{\mathbf{h}}_m\}_{m=1}^M$, scaled by M^{-1} . The sample covariance of $\{\tilde{\mathbf{h}}_m\}_{m=1}^M$ is unbiased for $\text{Cov}[\tilde{\mathbf{h}}]$. Rearranging the above expression yields

$$\begin{aligned} \|\mathbb{E}_{p_\theta}[\tilde{\mathbf{h}}(\mathbf{x})]\|^2 &= \mathbb{E}_{p_\theta} \left\| \frac{1}{M} \sum_{m=1}^M \tilde{\mathbf{h}}_m \right\|^2 - \frac{1}{M} \mathbb{E}_{p_\theta} \left[\frac{1}{(M-1)} \sum_{m=1}^M \left\| \tilde{\mathbf{h}}_m - \frac{1}{M} \sum_{m'=1}^M \tilde{\mathbf{h}}_{m'} \right\|^2 \right] \\ &= \mathbb{E}_{p_\theta}[\hat{\mathcal{L}}^{\text{viol}}] \end{aligned}$$

This proves $\hat{\mathcal{L}}^{\text{viol}}$ is unbiased for $\|\mathbb{E}_{p_\theta}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2$. \square

Likewise, we demonstrate that our estimate for the reward loss is unbiased.

Proposition B.3. $\hat{\mathcal{L}}^{\text{reward}}$ is unbiased for the reward loss $\mathcal{L}^{\text{reward}}$.

Proof. In the proof of Proposition B.2 we already demonstrated $M^{-1} \sum_{m=1}^M \log \frac{p_\theta(\mathbf{x}_m)}{p_{\theta_{\text{base}}}(\mathbf{x}_m)}$ is unbiased for \mathcal{L}^{KL} . By an identical argument, $-M^{-1} \sum_{m=1}^M r_{\hat{\alpha}_N}(\mathbf{x}_m)$ is unbiased for $\mathcal{L}^r = \mathbb{E}_{p_\theta}[-r_{\hat{\alpha}_N}(\mathbf{x})]$ (again, it is a Monte Carlo estimate). \square

B.2 UNBIASED GRADIENT ESTIMATES

As we detailed in Section 2.3, the naïve idea of taking the unbiased loss estimators $\hat{\mathcal{L}}^{\text{relax}}$, $\hat{\mathcal{L}}^{\text{reward}}$ and differentiating them with respect to θ will not yield unbiased estimates for the gradients of $\mathcal{L}^{\text{relax}}$ and $\mathcal{L}^{\text{reward}}$. This is because the probability distribution with respect to which the expectation is taken also depends on θ , which needs to be taken into account in the gradient estimate.

For CGM-reward, we propose the gradient estimate

$$\begin{aligned} \hat{G}^{\text{reward}} &= \frac{1}{M} \sum_{m=1}^M (\nabla_{\theta} w_m(\theta, \theta')) (l_m^{\text{LOO}} - r_m^{\text{LOO}}), \quad w_m(\theta, \theta') = \frac{p_\theta(\mathbf{x}_m)}{p_{\theta'}(\mathbf{x}_m)} \\ l_m^{\text{LOO}} &= l_m - \frac{1}{M-1} \sum_{m' \neq m} l_{m'}, \quad l_m = \log \frac{p_\theta(\mathbf{x}_m)}{p_{\theta_{\text{base}}}(\mathbf{x}_m)} \\ r_m^{\text{LOO}} &= r_m - \frac{1}{M-1} \sum_{m' \neq m} r_{m'}, \quad r_m = r_{\hat{\alpha}_N}(\mathbf{x}_m) \end{aligned} \quad (11)$$

As we explained in Section 2.3, $w_m(\theta, \theta')$ can be viewed as the weights of an *importance sampling* scheme, where $p_{\theta'}$ is the proposal distribution and p_θ is the target distribution. We choose $\theta' = \theta$ so that the proposal distribution is equal to the target distribution. For this choice of proposal, the weights w_m are all equal to 1. However, their gradient with respect to θ is equal to the score of the calibrated model at \mathbf{x}_m , $\nabla_{\theta} \log p_\theta(\mathbf{x}_m)$. The expression (11), excluding the terms $(M-1)^{-1} \sum_{m' \neq m} l_{m'}$ and $(M-1)^{-1} \sum_{m' \neq m} r_{m'}$ is known as the score function gradient estimate or, in the terminology of reinforcement learning, the REINFORCE gradient estimate (Williams, 1992).

The terms $(M-1)^{-1} \sum_{m' \neq m} l_{m'}$ and $(M-1)^{-1} \sum_{m' \neq m} r_{m'}$ in (11) are known as leave-one-out baselines (Kool et al., 2019) corresponding to sample \mathbf{x}_m . Including these terms adds to the score function gradient estimate a *control variate*, which is a term that has expectation zero under p_θ but is correlated with each individual term in the estimate (Lavenberg & Welch, 1981; Ranganath et al., 2014; Mohamed et al., 2020). Indeed, we observe that by independence of the samples $\{\mathbf{x}_m\}_{m=1}^M$, it holds that for each $m \neq m'$,

$$\mathbb{E}_{p_\theta}[(\nabla_{\theta} \log p_\theta(\mathbf{x}_m))(l_{m'} - r_{m'})] = \mathbb{E}_{p_\theta}[\nabla_{\theta} \log p_\theta(\mathbf{x}_m)] \mathbb{E}_{p_\theta}[l_{m'} - r_{m'}] = 0.$$

Consequently, while the inclusion of the leave-one-out averages does not affect the unbiasedness of our gradient estimate, they can reduce its variance.

Proposition B.4. \hat{G}^{reward} is unbiased for the gradient of the reward loss, $\nabla_{\theta} \mathcal{L}^{\text{reward}}$.

1026 *Proof.* We start by writing out the gradient of $\mathcal{L}^{\text{reward}}$ directly:

$$\begin{aligned}
1027 \nabla_{\theta} \mathcal{L}^{\text{reward}}(\theta) &= \nabla_{\theta} \mathbb{E}_{p_{\theta}} \left[\log \frac{p_{\theta}(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} - r_{\hat{\alpha}_N}(\mathbf{x}) \right] \\
1028 &= \nabla_{\theta} \int \left\{ \log \frac{p_{\theta}(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} - r_{\hat{\alpha}_N}(\mathbf{x}) \right\} p_{\theta}(d\mathbf{x}) \\
1029 &= \nabla_{\theta} \int \frac{p_{\theta}(\mathbf{x})}{p_{\theta'}(\mathbf{x})} \left\{ \log \frac{p_{\theta}(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} - r_{\hat{\alpha}_N}(\mathbf{x}) \right\} p_{\theta'}(d\mathbf{x}) \\
1030 &\stackrel{(*)}{=} \mathbb{E}_{p_{\theta}} \left[\left(\nabla_{\theta} \frac{p_{\theta}(\mathbf{x})}{p_{\theta'}(\mathbf{x})} \right) \left\{ \log \frac{p_{\theta}(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} - r_{\hat{\alpha}_N}(\mathbf{x}) \right\} \right] + \mathbb{E}_{p_{\theta}} \left[\nabla_{\theta} \frac{p_{\theta}(\mathbf{x})}{p_{\theta'}(\mathbf{x})} \right] \quad (12)
\end{aligned}$$

1031 where $\theta' = \text{stop-grad}(\theta)$. In equality $(*)$, exchange of the gradient and expectation is permis-
1032 sible as a consequence of dominated convergence and Assumption B.1. The second term is the
1033 expected score, which is zero. And so the gradient of the reward loss is

$$1034 \nabla_{\theta} \mathcal{L}^{\text{reward}}(\theta) = \mathbb{E}_{p_{\theta}} \left[(\nabla_{\theta} \log p_{\theta}(\mathbf{x})) \left\{ \log \frac{p_{\theta}(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} - r_{\hat{\alpha}_N}(\mathbf{x}) \right\} \right]. \quad (13)$$

1035 Looking at our gradient estimator \hat{G}^{reward} in (11) and ignoring the leave-one-out averages, we see
1036 that it is exactly the Monte Carlo estimate of the gradient of $\mathcal{L}^{\text{reward}}$ (13). \square

1037 Dropping the potentially noisy expected score term in (12), as is done by Ranganath et al. (2014),
1038 also reduces variance of our gradient estimate.

1039 Deriving an unbiased gradient estimate for the relax loss is more challenging, since the loss cannot
1040 be written as the expectation of some objective under p_{θ} . Just as we did for the reward loss, we
1041 can compute an unbiased estimate for the gradient of \mathcal{L}^{KL} in the relax loss by drawing independent
1042 samples $\mathbf{x}_m \sim p_{\theta'}$ and then differentiating the importance sampling weights $w_m(\theta, \theta')$

$$1043 \hat{G}_{\text{KL}} = \frac{1}{M} \sum_{m=1}^M (\nabla_{\theta} w_m(\theta, \theta')) l_m^{\text{LOO}}.$$

1044 And so it remains to compute an unbiased gradient estimate for $\mathcal{L}^{\text{viol}}$. To do so, we first recall the
1045 unbiased estimate $\hat{\mathcal{L}}^{\text{viol}}$ for $\mathcal{L}^{\text{viol}}$ that we introduced in Section 2.1

$$1046 \hat{\mathcal{L}}^{\text{viol}}(\{\tilde{\mathbf{h}}_m\}_{m=1}^M) := \left\| \frac{1}{M} \sum_{m=1}^M \tilde{\mathbf{h}}_m \right\|^2 - \frac{1}{M(M-1)} \sum_{m=1}^M \left\| \tilde{\mathbf{h}}_m - \frac{1}{M} \sum_{m'=1}^M \tilde{\mathbf{h}}_{m'} \right\|^2,$$

1047 where \mathbf{x}_m are independent samples from p_{θ} and $\tilde{\mathbf{h}}_m = \mathbf{h}(\mathbf{x}_m) - \mathbf{h}^*$. We propose a modifica-
1048 tion to this estimate wherein we draw independent samples $\mathbf{x}_m \sim p_{\theta'}$ and replace $\{\tilde{\mathbf{h}}_m\}_{m=1}^M$ by
1049 $\{w_m(\theta, \theta') \tilde{\mathbf{h}}_m\}_{m=1}^M$. To estimate the gradient of $\|\mathbb{E}_{p_{\theta}}[\mathbf{h}] - \mathbf{h}^*\|^2 = \|\mathbb{E}_{p_{\theta}}[\tilde{\mathbf{h}}]\|^2$, we compute the
1050 gradient of $\hat{\mathcal{L}}^{\text{viol}}(\{w_m(\theta, \theta') \tilde{\mathbf{h}}_m\}_{m=1}^M)$ with respect to θ and then evaluate the result at $\theta' = \theta$. In
1051 Algorithms 1 and 2, we implement our gradient estimator for $\mathcal{L}^{\text{viol}}$ by sampling \mathbf{x}_m independently
1052 from $p_{\text{stop-grad}(\theta)}$ and differentiating $\hat{\mathcal{L}}^{\text{viol}}(\{w_m(\theta, \theta') \tilde{\mathbf{h}}_m\}_{m=1}^M)$ with $\theta' = \text{stop-grad}(\theta)$.

1053 This yields the overall gradient estimator for the relax loss

$$1054 \hat{G}^{\text{relax}} = \nabla_{\theta} \hat{\mathcal{L}}^{\text{viol}} \left(\left\{ w(\theta, \theta') \tilde{\mathbf{h}}_m \right\}_{m=1}^M \right) + \lambda \hat{G}_{\text{KL}}, \quad \mathbf{x}_m \stackrel{i.i.d.}{\sim} p_{\theta'}, \quad \theta' = \text{stop-grad}(\theta).$$

1055 In order to prove that \hat{G}^{relax} is unbiased for $\nabla_{\theta} \mathcal{L}^{\text{relax}}$, we need to show $\hat{\mathcal{L}}^{\text{viol}}(\{w_m(\theta, \theta') \tilde{\mathbf{h}}_m\}_{m=1}^M)$
1056 remains unbiased for $\mathcal{L}^{\text{viol}}$ when \mathbf{x}_m are sampled independently from $p_{\theta'}$. Then, since the distribu-
1057 tion from which \mathbf{x}_m are sampled does not depend on θ , it is allowable to exchange the gradient with
1058 the expectation.

1059 **Proposition B.5.** \hat{G}^{relax} is unbiased for the gradient of the relax loss, $\nabla_{\theta} \mathcal{L}^{\text{relax}}$.

Proof. From Proposition B.2, we know that \widehat{G}_{KL} is unbiased for $\nabla_{\theta} \mathcal{L}^{\text{KL}}$, and so it only remains to verify that the second term is unbiased for $\nabla_{\theta} \mathcal{L}^{\text{viol}} = \nabla_{\theta} \|\mathbb{E}_{p_{\theta}}[\mathbf{h}] - \mathbf{h}^*\|^2$. To this end, by repeating the proof of Proposition B.2 (i.e., using the definition of the variance), it is straightforward to show

$$\mathbb{E}_{p_{\theta'}} \left[\widehat{\mathcal{L}}^{\text{viol}} \left(\left\{ \frac{p_{\theta}(\mathbf{x}_m)}{p_{\theta'}(\mathbf{x}_m)} \tilde{\mathbf{h}}_m \right\}_{m=1}^M \right) \right] = \left\| \mathbb{E}_{p_{\theta'}} \left[\frac{p_{\theta}(\mathbf{x}_m)}{p_{\theta'}(\mathbf{x}_m)} \tilde{\mathbf{h}}_m \right] \right\|^2 = \|\mathbb{E}_{p_{\theta}}[\tilde{\mathbf{h}}]\|^2.$$

In other words, $\widehat{\mathcal{L}}^{\text{viol}}(\{w_m(\theta, \theta') \tilde{\mathbf{h}}_m\}_{m=1}^M)$ is unbiased for $\mathcal{L}^{\text{viol}}$. However, since the samples $\{\mathbf{x}_m\}_{m=1}^M$ are drawn from $p_{\theta'}$, a probability distribution that does not depend on θ , then we can exchange the gradient and expectation by appealing to dominated convergence and Assumption B.1. In particular, we have

$$\begin{aligned} \mathbb{E}_{p_{\theta'}} \left[\nabla_{\theta} \widehat{\mathcal{L}}^{\text{viol}} \left(\left\{ \frac{p_{\theta}(\mathbf{x}_m)}{p_{\theta'}(\mathbf{x}_m)} \tilde{\mathbf{h}}_m \right\}_{m=1}^M \right) \right] &= \nabla_{\theta} \mathbb{E}_{p_{\theta'}} \left[\widehat{\mathcal{L}}^{\text{viol}} \left(\left\{ \frac{p_{\theta}(\mathbf{x}_m)}{p_{\theta'}(\mathbf{x}_m)} \tilde{\mathbf{h}}_m \right\}_{m=1}^M \right) \right] \\ &= \nabla_{\theta} \mathcal{L}^{\text{viol}}, \end{aligned}$$

where the final line follows from the unbiasedness of $\widehat{\mathcal{L}}^{\text{viol}}(\{w_m \tilde{\mathbf{h}}_m\}_{m=1}^M)$ for $\mathcal{L}^{\text{viol}}$. \square

As we discussed, the key insight from the proof of Proposition B.5 is that, by introducing importance weights, we can compute an unbiased estimate to $\|\mathbb{E}_{p_{\theta}}[\mathbf{h}] - \mathbf{h}^*\|^2 = \|\mathbb{E}_{p_{\theta}}[\tilde{\mathbf{h}}]\|^2$ without sampling directly from p_{θ} .

C MAXIMUM ENTROPY PRINCIPLE

In this section, we provide an overview of the maximum entropy principle, which we use in Section 2.2 to define the reward loss $\mathcal{L}^{\text{reward}}$. First, in Appendix C.1 we formally state and prove the maximum entropy principle. In Appendix C.2, we provide greater detail on our estimate $\widehat{\alpha}_N$ for the parameters α^* of the maximum entropy solution. In Appendix C.3, we characterize the relationship between the relax and reward losses by considering a problem whose solution is close to the optimum of the relax loss, and which resembles the maximum entropy problem. Lastly, in Appendix C.4, we study the behavior of the estimate $\widehat{\alpha}_N$ in the limit as the number of samples N becomes large.

Prior to jumping into the details of the maximum entropy principle, we work through an illustrative example that we discuss throughout this section.

Example. Suppose $\mathbf{x} \in \mathbb{R}$, $\mathbf{h}(\mathbf{x}) = \mathbb{1}\{\mathbf{x} > 0\}$, and $\mathbf{h}^* \in \mathbb{R}$. Also define $h_b = \mathbb{P}_{p_{\theta_{\text{base}}}}(\mathbf{x} > 0)$, and assume $0 < h_b < 1$. In this example, the calibration problem amounts to either upweighting or downweighting the amount of probability mass h_b that lies above 0 under the base model $p_{\theta_{\text{base}}}$. By Theorem 2.1, the maximum entropy solution has the form $p_{\alpha^*} \propto p_{\theta_{\text{base}}}(\mathbf{x}) \exp\{\alpha^* \mathbf{h}(\mathbf{x})\}$ for some $\alpha^* \in \mathbb{R}$ that we need to determine. From this expression for p_{α^*} , we obtain

$$\begin{aligned} 1 - \mathbf{h}^* &= \mathbb{E}_{p_{\alpha^*}}[1 - \mathbf{h}(\mathbf{x})] = \frac{1}{h_b \exp(\alpha^*) + (1 - h_b)} (1 - h_b), \\ \mathbf{h}^* &= \mathbb{E}_{p_{\alpha^*}}[\mathbf{h}(\mathbf{x})] = \frac{1}{h_b \exp(\alpha^*) + (1 - h_b)} h_b \exp(\alpha^*). \end{aligned}$$

Dividing the first equation by the second and rearranging yields $\alpha^* = \log\left(\frac{\mathbf{h}^*(1-h_b)}{(1-\mathbf{h}^*)h_b}\right)$. Following the same argument for the empirical distribution of $\{\mathbf{x}_n\}_{n=1}^N$, our estimator for α^* is $\widehat{\alpha}_N = \log\left(\frac{\mathbf{h}^*(1-\bar{\mathbf{y}}_N)}{(1-\mathbf{h}^*)\bar{\mathbf{y}}_N}\right)$, where $\bar{\mathbf{y}}_N = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$, $\mathbf{y}_n = \mathbb{1}\{\mathbf{x}_n > 0\} \stackrel{d}{=} \text{Bernoulli}(h_b)$ for $\mathbf{x}_n \stackrel{i.i.d.}{\sim} p_{\theta_{\text{base}}}$.

We point out that α^* and $\widehat{\alpha}_N$ can equivalently be derived by differentiating the objectives (6) and (7), respectively, and setting them equal to 0.

1134 C.1 PRECISE STATEMENT
1135

1136 Since the maximum entropy problem is not specific to generative model calibration, we present it
1137 in a more general setting. Our presentation builds on standard results from exponential families and
1138 convex analysis. We recommend Wainwright & Jordan (2008) for relevant background.

1139 In particular, we consider $X := (X, \mathcal{X})$ a measurable space, P a probability measure defined on
1140 X , $\mathbf{h} : X \rightarrow \mathbb{R}^d$ an X -measurable constraint function, and \mathbf{h}^* a target value for the moment of \mathbf{h} .
1141 The maximum entropy problem corresponding to probability measure P , constraint \mathbf{h} , and target
1142 moment \mathbf{h}^* is

$$1143 \inf_{Q \in \mathcal{P}(P)} \text{D}_{\text{KL}}(Q \parallel P), \quad \text{such that } \mathbb{E}_Q[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*. \quad (14)$$

1144 $\mathcal{P}(P)$ is the collection of all probability measures having a density with respect to P , which, by the
1145 Radon-Nikodym Theorem, is equal to the collection of all absolutely continuous probability mea-
1146 sures with respect to P . Choosing $P = p_{\theta_{\text{base}}}$ yields the maximum entropy problem corresponding
1147 to the calibration problem.

1148 As we mentioned in Section 2.2, we impose a condition on the target moment \mathbf{h}^* to ensure (i) there
1149 exists a solution to the maximum entropy problem (ii) and this solution is an exponential tilt of P .

1150 **Assumption C.1** (Interior moment condition). Define the subset \mathcal{M} of \mathbb{R}^d comprised of all possible
1151 moments of \mathbf{h} attainable by probability distributions Q having a density with respect to P

$$1152 \mathcal{M} = \left\{ \int \mathbf{h}(\mathbf{x})Q(d\mathbf{x}) \mid Q \in \mathcal{P}(P), \int \|\mathbf{h}(\mathbf{x})\|Q(d\mathbf{x}) < \infty \right\}.$$

1153 \mathbf{h}^* lies in the *relative interior* of \mathcal{M} , written $\text{relint}(\mathcal{M})$.

1154 Since \mathcal{M} is a convex set, the condition $\mathbf{h}^* \in \text{relint}(\mathcal{M})$ can equivalently be stated as for every
1155 $\mathbf{y} \neq \mathbf{h}^*$ in \mathcal{M} , there exists some \mathbf{z} in \mathcal{M} and $\kappa \in (0, 1)$ for which $\mathbf{h}^* = \kappa\mathbf{z} + (1 - \kappa)\mathbf{y}$.

1156 To see why Assumption C.1 is necessary for the solution to be an exponential tilt of $p_{\theta_{\text{base}}}$, recall the
1157 example discussed at the beginning of Appendix C. In this case, $\text{relint}(\mathcal{M}) = (0, 1)$. If $\mathbf{h}^* \notin [0, 1]$,
1158 then there does not exist any probability distribution p having density with respect to $p_{\theta_{\text{base}}}$ for which
1159 $\mathbb{E}_p[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*$. And if \mathbf{h}^* is either 0 or 1, then the solution to the maximum entropy problem is
1160 proportional to $p_{\theta_{\text{base}}}(\mathbf{x})\mathbb{1}\{\mathbf{x} \leq 0\}$ or $p_{\theta_{\text{base}}}(\mathbf{x})\mathbb{1}\{\mathbf{x} > 0\}$, respectively. Neither of these solutions is
1161 an exponential tilt of $p_{\theta_{\text{base}}}$, equation (5).

1162 Our proof of the maximum entropy principle leverages classical convex duality (Rockafellar, 1970)
1163 by showing that (14) is a convex problem, defined on the infinite-dimensional space of all probability
1164 densities for which \mathbf{h} has a finite moment. The corresponding *dual problem* is

$$1165 \sup_{\alpha \in \mathbb{R}^d} \alpha^\top \mathbf{h}^* - A_P(\alpha), \quad A_P(\alpha) := \log \left(\int \exp\{r_\alpha(\mathbf{x})\}P(d\mathbf{x}) \right), \quad r_\alpha(\mathbf{x}) = \alpha^\top \mathbf{h}(\mathbf{x}), \quad (15)$$

1166 which is concave. $A_P : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$ is known as the *log-normalizer* or *cumulant generating*
1167 *function* corresponding to the exponential family

$$1168 \exp\{r_\alpha(\mathbf{x}) - A_P(\alpha)\}P(d\mathbf{x}). \quad (16)$$

1169 We will make the standard assumption that the domain of A_P is open

1170 **Assumption C.2** (Domain of log-normalizer). The subset $\Xi = \{\alpha \in \mathbb{R}^d \mid A_P(\alpha) < \infty\}$ is open.

1171 Whenever A_P is finite, (16) is a well-defined probability measure on X . Ξ is known as the *natural*
1172 *parameter space* of the exponential family (16). When Assumption C.2 holds, the exponential
1173 family is said to be *regular*.

1174 The log-normalizer A_P possesses many nice properties: for instance, it is convex and infinitely
1175 differentiable on Ξ . Convexity can be seen by computing the Hessian of $A_P(\alpha)$

$$1176 \nabla_\alpha^2 A_P(\alpha) = \frac{\int (\mathbf{h}(\mathbf{x}) - \nabla_\alpha A_P(\alpha))(\mathbf{h}(\mathbf{x}) - \nabla_\alpha A_P(\alpha))^\top \exp\{r_\alpha(\mathbf{x})\}P(d\mathbf{x})}{\int \exp\{r_\alpha(\mathbf{x})\}P(d\mathbf{x})} \quad (17)$$

and recognizing that it is positive semi-definite. Differentiability is addressed in the remark following Lemma C.11.

Now that we have introduced the dual of the maximum entropy problem, we are prepared to give a precise statement and proof of the maximum entropy principle

Theorem C.3 (Kullback (1959)). *Suppose Assumptions C.1 and C.2 hold. Then there exists a probability measure $Q^* \in \mathcal{P}(P)$ with density $dQ^*/dP \propto \exp(r_{\alpha^*}(\mathbf{x}))$. Moreover, Q^* is the solution to the maximum entropy problem (14) and is unique up to P -null sets.*

Unlike the primal problem (14), the dual problem (15) is defined on finite-dimensional Euclidean space, which makes it simpler to analyze. We first argue by weak duality that the value of (14) is at least as large as (15). We then identify a vector α^* and a distribution Q^* for which the primal and dual objectives are equal. By weak duality, this implies that Q^* is optimal for the primal problem.

Proof of Theorem C.3. We first rewrite the primal problem (14) in the form

$$\begin{aligned} & \inf_q \psi(q) + g(\mathcal{A}q) \\ & \psi(q) = \begin{cases} \int q(\mathbf{x}) \log(q(\mathbf{x})) P(d\mathbf{x}) & \text{if } q \geq 0 \\ +\infty & \text{else} \end{cases}, \quad g(\mathbf{y}_0, \mathbf{y}_1) = \begin{cases} 0 & \text{if } \mathbf{y}_0 = 1 \text{ and } \mathbf{y}_1 = \mathbf{h}^* \\ +\infty & \text{else} \end{cases}, \\ & \mathcal{A}(q) = \left(\int q(\mathbf{x}) P(d\mathbf{x}), \int \mathbf{h}(\mathbf{x}) q(\mathbf{x}) P(d\mathbf{x}) \right) \end{aligned}$$

defined on the space of X -measurable functions q for which $\int |q(\mathbf{x})| P(d\mathbf{x}) < \infty$ and $\int \|\mathbf{h}(\mathbf{x})\| q(\mathbf{x}) P(d\mathbf{x}) < \infty$. Here, q represents the density of measure Q with respect to P . $g(\mathcal{A}q)$ imposes the constraint that Q is a probability measure and that the expectation of \mathbf{h} under Q is \mathbf{h}^* . And $\psi(q)$ is equal to the KL divergence between Q and P .

Observe \mathcal{A} is a bounded, linear map defined on this space. And ψ and g are convex. By Fenchel-Rockafellar duality (Borwein & Zhu, 2005, Theorem 4.4.2), weak duality holds for the maximum entropy problem and its dual (15).

Wainwright & Jordan (2008, Theorem 3.3) states that $\nabla_{\alpha} A_P$ is a surjective mapping from Ξ onto $\text{relint}(\mathcal{M})$. Hence, there exists $\alpha^* \in \Xi$ for which $\nabla_{\alpha} A_P(\alpha^*) = \mathbf{h}^*$. The value of the dual at α^* is

$$(\alpha^*)^{\top} \mathbf{h}^* - A_P(\alpha^*).$$

By differentiating the dual objective at α^* , we obtain,

$$0 = \nabla_{\alpha} (\alpha^{\top} \mathbf{h}^* - A_P(\alpha)) \implies \mathbf{h}^* = \frac{\int \mathbf{h}(\mathbf{x}) \exp\{r_{\alpha^*}(\mathbf{x})\} P(d\mathbf{x})}{\int \exp\{r_{\alpha^*}(\mathbf{x})\} P(d\mathbf{x})}.$$

In other words, the distribution $Q^* \in \mathcal{P}(P)$ defined such that $dQ^*/dP \propto \exp\{r_{\alpha^*}(\mathbf{x})\}$ satisfies the moment constraint $\mathbb{E}_{Q^*}[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*$. Moreover, the value of the primal objective at Q^* is

$$D_{\text{KL}}(Q^* \parallel P) = (\alpha^*)^{\top} \mathbf{h}^* - A_P(\alpha^*),$$

which is equal to the value of the dual objective at α^* . By weak duality, we conclude Q^* is the solution to the maximum entropy problem.

Uniqueness follows from the fact that the KL divergence ψ is strictly convex. \square

C.2 ESTIMATING THE MAXIMUM ENTROPY SOLUTION

Next, we discuss our estimator $\hat{\alpha}_N$ for the parameters α^* of the maximum entropy solution. In particular, we provide verifiable conditions under which $\hat{\alpha}_N$ is well-defined, and we show that this estimator can be interpreted as the solution to a finite-sample version of the maximum entropy problem (4).

So far, the only assumptions we have made on the maximum entropy problem (14) are the relative interior condition on \mathbf{h}^* (Assumption C.1) and the openness condition for the domain of A_P (Assumption C.2). As we demonstrated in Appendix C.1, these conditions ensure that the solution to

the maximum entropy problem exists and is unique. However, the solution to the dual problem need not be unique. Suppose, for example, that \mathbf{h} is d -dimensional but has two identical components $\mathbf{h}(\mathbf{x})[i] = \mathbf{h}(\mathbf{x})[j]$. Then if α^* is optimal for the dual problem, so is $\alpha^* - te[i] + te[j]$ for all $t \in \mathbb{R}$, where $e[i]$ and $e[j]$ denote the i and j th standard basis vectors, respectively. Specifically, the set of optima for the dual problem is a hyperplane in \mathbb{R}^d . In order to estimate α^* , we want to ensure that the dual problem (15) also has a unique maximum.

As suggested by our example, in order to ensure that the dual optimum is unique, it suffices to eliminate linear redundancies among the statistics $\mathbf{h}(\mathbf{x})$.

Assumption C.4 (Uniqueness of dual optimum). No linear combination of the components of $\mathbf{h}(\mathbf{x})$ is equal to a constant with P probability one.

If Assumption C.4 holds, then the exponential family (16) is said to be *minimal*. An exponential family for which Assumption C.2 holds is minimal if and only if the log-normalizer $A_P(\alpha)$ is strictly convex on Ξ (Wainwright & Jordan, 2008, Proposition 3.1).

For non-trivial generative models, solving the dual problem (15) for $P = p_{\theta_{\text{base}}}$ is intractable since $A_{p_{\theta_{\text{base}}}}(\alpha)$ cannot be computed in closed-form. The estimator $\hat{\alpha}_N$ that we propose in (7) involves first drawing N independent samples $\{\mathbf{x}_n\}_{n=1}^N$ from the base model $p_{\theta_{\text{base}}}$ and then solving the dual problem with the integral replaced by the empirical average from our samples. This is equivalent to solving the dual problem for P equal to the empirical distribution of our samples $\frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}$, where $\delta_{\mathbf{x}}$ is the delta function at \mathbf{x} .

However, in order for $\hat{\alpha}_N$ to be well-defined, the interior point condition and uniqueness of the dual optimum must hold for the maximum entropy problem with $P = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}$. For this problem, these two conditions are straightforward to verify: (i) \mathbf{h}^* lies in the the relative interior of the convex hull of $\{\mathbf{h}(\mathbf{x}_n)\}_{n=1}^N$ and (ii) the empirical covariance matrix of $\{\mathbf{h}(\mathbf{x}_n)\}_{n=1}^N$ has full rank. For the example we provided at the beginning of the section, conditions (i) and (ii) are satisfied if and only if $\{\mathbf{h}(\mathbf{x}_n)\} = \{0, 1\}$ and $\mathbf{h}^* \in (0, 1)$.

It is possible for Assumptions C.1 and C.4 to hold for $p_{\theta_{\text{base}}}$ but not for $\frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}$. For our example, if $\{\mathbf{h}(\mathbf{x}_n)\} = \{0\}$ and $\mathbf{h}^* = 0$ (or $\{\mathbf{h}(\mathbf{x}_n)\} = \{1\}$ and $\mathbf{h}^* = 1$), then the maximum entropy solution exists and is equal to $Q^* = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}$, but every vector $\alpha \in \mathbb{R}$ is optimal for the dual problem (15). We demonstrate in Appendix C.4 that the probability of this event approaches zero as the number of samples N approaches infinity. However, we observe (e.g., Figure 2B) that when the base model $p_{\theta_{\text{base}}}$ lies far from the maximum entropy solution p_{α^*} , estimating α^* with small variance requires many samples, and may even be computationally intractable.

C.3 CONNECTION BETWEEN THE RELAX AND REWARD LOSSES

In this section, we elucidate the connection between the relax and reward losses. We first introduce a problem corresponding to the relax loss that, similar to the maximum entropy problem (4), is defined on the space $\mathcal{P}(p_{\theta_{\text{base}}})$ of probability distributions that have a density with respect to $p_{\theta_{\text{base}}}$. When the generative model class p_{θ} is sufficiently expressive, the solution to this problem well approximates the minimizer of the relax loss. We then show that, under conditions, the solution to this related problem approaches the solution to the maximum entropy problem as $\lambda \rightarrow 0$. This confirms our intuition that when $\lambda \rightarrow 0$, minimizing the relax loss is equivalent to solving the calibration problem.

As in Appendix C.1, we let $X := (X, \mathcal{X})$ be a measurable space, P be a probability measure defined on X , and $\mathbf{h} : X \rightarrow \mathbb{R}^d$ be a X -measurable function, and \mathbf{h}^* be a target moment. We consider the problem

$$\inf_{Q \in \mathcal{P}(P)} \|\mathbb{E}_Q[\mathbf{h}] - \mathbf{h}^*\|^2 + \lambda \text{D}_{\text{KL}}(Q \| P), \quad \text{s.t. } \mathbb{E}_Q[\|\mathbf{h}\|] < \infty \quad (18)$$

In convex analysis (e.g., Boyd & Vandenberghe, 2004; Ben-Tal & Nemirovski, 2023), (18) is known as a penalty problem.

When $P = p_{\theta_{\text{base}}}$, then (18) agrees with the problem of minimizing the relax loss (2), except the domain of the problem is $\mathcal{P}(p_{\theta_{\text{base}}})$ rather than the class of generative models p_{θ} . Suppose momentarily that the infimum of (18), denoted by Q_{λ} , is attained. The minimizer of the relax loss (2) will not in general be equal to Q_{λ} since Q_{λ} does not lie in the class of generative models. However, as we

argued when we proposed the reward loss, we would expect Q_λ and the minimizer of the relax loss to be close in KL distance when the class of generative models p_θ is sufficiently expressive.

Introducing the problem (18) is helpful insofar as, similar to the maximum entropy problem, we can obtain a closed-form expression for the solution Q_λ .

Proposition C.5. *Suppose Assumption C.2 holds. Then there exists a unique solution α_λ to the fixed point equation*

$$\alpha = -\frac{2}{\lambda}(\nabla_{\alpha} A_P(\alpha) - \mathbf{h}^*), \quad \alpha \in \Xi.$$

Moreover, Q_λ defined by $dQ_\lambda/dP \propto \exp\{\alpha_\lambda^\top \mathbf{h}(\mathbf{x})\}$ is the unique solution to (18).

Our proof mirrors that for the maximum entropy principle (Theorem C.3). Namely, we invoke Fenchel-Rockafellar duality (Rockafellar, 1970) to relate the convex problem (18), defined on the space of probability densities with respect to P with finite \mathbf{h} moment, to its concave dual problem

$$\sup_{\alpha \in \mathbb{R}^p} F_\lambda(\alpha), \quad F_\lambda(\alpha) = \lambda \left(-\frac{\lambda}{4} \|\alpha\|^2 - A_P(\alpha) + \alpha^\top \mathbf{h}^* \right) \quad (19)$$

defined on Euclidean space. We then show that α_λ is the unique solution to the dual problem, and we use this solution to construct a solution to the primal problem. Interestingly, α_λ is the unique solution to the dual problem even when there is redundancy among the constraints \mathbf{h} (i.e., Assumption C.4 does not hold).

Proof of Proposition C.5. We rewrite the primal problem (18) in the form

$$\begin{aligned} & \inf_q \psi(q) + g(Aq), \\ & \psi(q) = \begin{cases} \lambda \int q(\mathbf{x}) \log(q(\mathbf{x})) P(d\mathbf{x}) & \text{if } q \geq 0 \\ +\infty & \text{else} \end{cases}, \quad g(\mathbf{y}_0, \mathbf{y}_1) = \begin{cases} \|\mathbf{y}_1 - \mathbf{h}^*\|^2 & \text{if } \mathbf{y}_0 = 1 \\ +\infty & \text{else} \end{cases}, \\ & A(q) = \left(\int q(\mathbf{x}) P(d\mathbf{x}), \int \mathbf{h}(\mathbf{x}) q(\mathbf{x}) P(d\mathbf{x}) \right) \end{aligned}$$

defined on the space of X -measurable functions q for which $\int |q(\mathbf{x})| P(d\mathbf{x}) < \infty$ and $\int \|\mathbf{h}(\mathbf{x})\| q(\mathbf{x}) P(d\mathbf{x}) < \infty$. Here, q represents the density of measure Q with respect to P . $g(Aq)$ is equal to $\|\mathbb{E}_Q[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2$ if Q is a probability measure and is infinite otherwise. $\psi(q)$ is equal to the KL divergence between Q and P , scaled by λ .

As in the proof of Theorem C.3, A is a bounded, linear map defined on this space, and ψ and g are convex. By Fenchel-Rockafellar duality (Borwein & Zhu, 2005, Theorem 4.4.2), weak duality holds for the problem (18) and its dual (19).

By the remark following Lemma C.11, F_λ is infinitely differentiable on Ξ , and taking two derivatives of $F_\lambda(\alpha)$ yields

$$\nabla_{\alpha} F_\lambda(\alpha) = \lambda \left(-\frac{\lambda}{2} \alpha - \nabla_{\alpha} A_P(\alpha) + \mathbf{h}^* \right), \quad \nabla_{\alpha}^2 F_\lambda(\alpha) = \lambda \left(-\frac{\lambda}{2} \mathbb{I} - \nabla_{\alpha}^2 A_P(\alpha) \right).$$

Since $\nabla_{\alpha}^2 A_P(\alpha)$ is positive semi-definite, then the problem (19) is strongly concave. And by our assumption that Ξ is open, $F_\lambda(\alpha)$ is equal to $-\infty$ for α belonging to the boundary of Ξ . Together with strong concavity, this implies a unique maximizer α_λ of F_λ exists.

In particular, α_λ is the unique $\alpha \in \mathbb{R}^d$ that satisfies the fixed-point equation

$$\nabla_{\alpha} F_\lambda(\alpha) = \lambda \left(-\frac{\lambda}{2} \alpha - \nabla_{\alpha} A(\alpha) + \mathbf{h}^* \right) = \mathbf{0} \implies \alpha = -\frac{2}{\lambda} (\nabla_{\alpha} A_P(\alpha) - \mathbf{h}^*).$$

And the probability measure $Q_{\alpha_\lambda} \propto \exp\{r_{\alpha_\lambda}(\mathbf{x})\} P(d\mathbf{x})$ satisfies

$$\begin{aligned} & \lambda \text{D}_{\text{KL}}(Q_{\alpha_\lambda} \parallel P) + \|\mathbb{E}_{Q_{\alpha_\lambda}}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2 \\ & = \lambda (\alpha_\lambda^\top \nabla_{\alpha} A_P(\alpha_\lambda) - A_P(\alpha_\lambda)) + \frac{\lambda^2}{4} \|\alpha_\lambda\|^2 \\ & = \lambda \left(\alpha_\lambda^\top \left(\mathbf{h}^* - \frac{\lambda}{2} \alpha_\lambda \right) - A_P(\alpha_\lambda) \right) + \frac{\lambda^2}{4} \|\alpha_\lambda\|^2 \\ & = F_\lambda(\alpha_\lambda). \end{aligned}$$

By weak duality, this implies $Q_\lambda := Q_{\alpha_\lambda}$ is optimal for the primal problem. Moreover, strict convexity of ψ implies that the optimum of the primal problem is unique. \square

Next, we show that as the regularization parameter $\lambda \rightarrow 0$, then Q_λ achieves the minimum possible Euclidean norm constraint violation i.e., Euclidean norm difference between $\mathbb{E}_{Q_\lambda}[\mathbf{h}]$ and \mathbf{h}^* . We also give a finite λ bound on the constraint violation.

Proposition C.6. *The distribution Q_λ satisfies*

$$\lim_{\lambda \rightarrow 0} \|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| = \inf_{\substack{Q \in \mathcal{P}(P) \\ D_{\text{KL}}(Q \| P) < \infty}} \|\mathbb{E}_Q[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|.$$

Moreover, we have the finite-sample bound on the Euclidean norm constraint violation of Q_λ

$$\|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| \leq \inf_{Q \in \mathcal{P}(P)} \left\{ \sqrt{\lambda D_{\text{KL}}(Q \| P)} + \|\mathbb{E}_Q[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| \right\}.$$

Proof. Fix $\varepsilon > 0$ and let Q_ε be such that $\|\mathbb{E}_{Q_\varepsilon}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| \leq \inf_{Q \in \mathcal{P}(P)} \|\mathbb{E}_Q[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| + \varepsilon$. Then by the optimality of Q_λ for the objective (18),

$$\begin{aligned} \|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2 &\leq \lambda D_{\text{KL}}(Q_\lambda \| P) + \|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2 \\ &\leq \lambda D_{\text{KL}}(Q_\varepsilon \| P) + \|\mathbb{E}_{Q_\varepsilon}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2. \end{aligned} \quad (20)$$

Our choice of Q_ε yields

$$\|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2 \leq \lambda D_{\text{KL}}(Q_\varepsilon \| P) + \inf_{Q \in \mathcal{P}(P)} \|\mathbb{E}_Q[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\|^2 + \varepsilon.$$

Taking $\lambda \rightarrow 0$ and then $\varepsilon \rightarrow 0$ yields the first result. Replacing Q_ε with $Q \in \mathcal{P}(P)$ in (20) and taking the infimum over Q yields the second result. \square

In the setting of Proposition C.9 where a solution to the maximum entropy problem exists, then a bound on the Euclidean norm constraint violation of Q_λ is simply $\sqrt{\lambda D_{\text{KL}}(Q^* \| P)}$. This implies that $\|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| = \mathcal{O}(\sqrt{\lambda})$.

From our proof of Proposition C.6, it is clear that we did not take advantage of the structure to the solution Q_λ . When Assumptions C.1 and C.4 hold, we can obtain a faster rate of convergence of $\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})]$ to \mathbf{h}^* , and we can show that α_λ converges to the parameters α^* of the maximum entropy distribution.

Proposition C.7. *Suppose Assumptions C.1, C.2, and C.4 hold, which imply that the maximum entropy solution $dQ^*/dP \propto \exp\{r_{\alpha^*}(\mathbf{x})\}$ exists. Then $\alpha_\lambda \rightarrow \alpha^*$ as $\lambda \rightarrow 0$. In particular,*

- (i) $\|\alpha_\lambda - \alpha^*\| = \mathcal{O}(\lambda)$
- (ii) $\|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| = \mathcal{O}(\lambda)$
- (iii) $|D_{\text{KL}}(Q_\lambda \| P) - D_{\text{KL}}(Q^* \| P)| = \mathcal{O}(\lambda)$.

Proof. Prior to proving (i)-(iii), we first establish $\|\alpha_\lambda - \alpha^*\| = o(1)$. From the proof of Proposition C.5, we know that α_λ maximizes $\lambda^{-1}F_\lambda(\alpha) = -\frac{\lambda}{4}\|\alpha\|^2 - A_P(\alpha) + \alpha^\top \mathbf{h}^*$ for each $\lambda > 0$. And from (15), we know that α^* maximizes $F_0(\alpha) = -A_P(\alpha) + \alpha^\top \mathbf{h}^*$. Clearly, $F_\lambda(\alpha) \rightarrow F_0(\alpha)$ pointwise as $\lambda \rightarrow 0$. Since each of F_λ and F_0 is concave on Ξ , a classical result in convex analysis Rockafellar (1970, Theorem, 10.8) implies that the convergence $F_\lambda(\alpha) \rightarrow F_0(\alpha)$ is uniform on closed, bounded subsets of Ξ containing α^* .

Fix $\epsilon > 0$ such that the Euclidean ball of radius ϵ centered at α^* is contained in Ξ . By Assumption C.4, F_0 is strictly concave, since then $\nabla_{\alpha}^2 A_P(\alpha)$ positive definite for every $\alpha \in \Xi$. Hence, there exists a κ such that for all $\|\alpha - \alpha^*\| = \epsilon$,

$$F_0(\alpha) > \kappa > F_0(\alpha^*).$$

This is because the left-hand side of the above inequality attains its minimum on the compact set $\|\alpha - \alpha^*\| = \epsilon$ and (ii) by strict concavity this minimum must be strictly greater than the right-hand

side. Moreover, by uniform convergence of F_λ to F_0 , there exists $\lambda_\epsilon > 0$ such that for all $\lambda < \lambda_\epsilon$ and all $\|\alpha - \alpha^*\| = \epsilon$

$$F_\lambda(\alpha) > \kappa > F_\lambda(\alpha^*). \quad (21)$$

Since F_λ is also concave, (21) implies that the maximizer of F_λ , α_λ , must lie in the Euclidean ball of radius ϵ centered at α^* . This establishes $\|\alpha_\lambda - \alpha^*\| = o(1)$.

We are now prepared to prove (i). By Taylor expanding $\nabla_\alpha A_P(\alpha)$ at α_λ about α^* , we obtain

$$\nabla_\alpha A_P(\alpha_\lambda) = \mathbf{h}^* + \nabla_\alpha^2 A_P(\alpha^*)(\alpha_\lambda - \alpha^*) + \mathbf{r}_\lambda, \quad \|\mathbf{r}_\lambda\| = o(\|\alpha_\lambda - \alpha^*\|). \quad (22)$$

By Proposition C.5, α_λ satisfies $\alpha_\lambda = -\frac{2}{\lambda}(\nabla_\alpha A(\alpha_\lambda) - \mathbf{h}^*)$. Multiplying (22) by $-2/\lambda$ and substituting in this expression for α_λ yields

$$\alpha_\lambda = -\frac{2}{\lambda} \nabla_\alpha^2 A_P(\alpha^*)(\alpha_\lambda - \alpha^*) + \frac{1}{\lambda} \mathbf{r}_\lambda.$$

Solving for $\alpha_\lambda - \alpha^*$ yields

$$\begin{aligned} \alpha_\lambda - \alpha^* &= -\left(\mathbb{I} + \frac{2}{\lambda} \nabla_\alpha^2 A_P(\alpha^*)\right)^{-1} \left(\alpha^* + \frac{1}{\lambda} \mathbf{r}_\lambda\right) \\ &= -\lambda \left(\lambda \mathbb{I} + 2 \nabla_\alpha^2 A_P(\alpha^*)\right)^{-1} \alpha^* + \tilde{\mathbf{r}}_\lambda \end{aligned} \quad (23)$$

for $\tilde{\mathbf{r}}_\lambda = o(\|\alpha_\lambda - \alpha^*\|)$. And because $\|\alpha_\lambda - \alpha^*\| = o(1)$, then for all λ sufficiently small, $\|\tilde{\mathbf{r}}_\lambda\| \leq \frac{1}{2} \|\alpha_\lambda - \alpha^*\|$. Taking the norm of both sides of (23) and rearranging yields

$$\|\alpha_\lambda - \alpha^*\| \leq 2\lambda \left\| \left(\lambda \mathbb{I} + 2 \nabla_\alpha^2 A_P(\alpha^*)\right)^{-1} \alpha^* \right\|$$

for all λ sufficiently small. This proves (i).

For (ii), the relationship $\alpha_\lambda = -\frac{2}{\lambda}(\nabla_\alpha A(\alpha_\lambda) - \mathbf{h}^*)$ yields

$$\|\mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*\| = \|\nabla_\alpha A(\alpha_\lambda) - \mathbf{h}^*\| \leq \frac{\lambda}{2} \|\alpha_\lambda - \alpha^*\| + \frac{\lambda}{2} \|\alpha^*\| = \mathcal{O}(\lambda).$$

Lastly for (iii),

$$\begin{aligned} \mathbf{D}_{\text{KL}}(Q_\lambda \parallel P) &= \alpha_\lambda^\top \mathbb{E}_{Q_\lambda}[\mathbf{h}(\mathbf{x})] - A_P(\alpha_\lambda) \\ &= (\alpha_\lambda^\top \mathbf{h}^* + \mathcal{O}(\lambda)) - \{A_P(\alpha^*) + \nabla_\alpha A_P(\alpha^*)^\top (\alpha_\lambda - \alpha^*) + o(\|\alpha_\lambda - \alpha^*\|)\} \\ &= \alpha_\lambda^\top \mathbf{h}^* - A_P(\alpha^*) + \mathcal{O}(\lambda) \\ &= \mathbf{D}_{\text{KL}}(Q^* \parallel P) + \mathcal{O}(\lambda). \end{aligned}$$

□

In Section 2.2 we derived the reward loss as the KL divergence of the model p_θ to the maximum entropy solution p_{α^*} . The relax loss can also be viewed as a divergence to a tilt of the base model $p_{\theta_{\text{base}}}$, except that the tilt depends on the current model p_θ . In particular, the stationary points of the relaxed loss are exactly the stationary points of the objective

$$\mathbf{D}_{\text{KL}}(p_\theta \parallel p_{\theta_{\text{base}}}) + \frac{2}{\lambda} (\mathbb{E}_{p_{\text{sg}}(\theta)}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*)^\top \mathbb{E}_{p_\theta}[\mathbf{h}(\mathbf{x})]. \quad (24)$$

This can be seen by taking the gradient of (24). By identifying $\alpha = -\frac{2}{\lambda}(\mathbb{E}_{p_{\text{sg}}(\theta)}[\mathbf{h}(\mathbf{x})] - \mathbf{h}^*)$ and $q_\alpha \propto p_{\theta_{\text{base}}}(\mathbf{x}) \exp\{r_\alpha(\mathbf{x})\}$, we observe that (24) is exactly equal to $\mathbf{D}_{\text{KL}}(p_\theta \parallel q_\alpha)$. q_α can be understood as our current best approximation to the solution of (18). Unlike the solution of (18), though, $\mathbb{E}_{q_\alpha}[\mathbf{h}(\mathbf{x})]$ is not equal to $\mathbb{E}_{p_{\text{sg}}(\theta)}[\mathbf{h}(\mathbf{x})]$. For sufficiently expressive class of generative models p_θ , we would expect $\mathbb{E}_{q_\alpha}[\mathbf{h}(\mathbf{x})]$ and $\mathbb{E}_{p_{\text{sg}}(\theta)}[\mathbf{h}(\mathbf{x})]$ to be approximately equal at the optimum.

1458 C.4 CONSISTENCY AND ASYMPTOTIC NORMALITY
1459

1460 In this section, we discuss the large sample behavior of the estimator $\widehat{\alpha}_N$ for the parameters α^* of
1461 the reward loss. Under Assumptions C.1, C.2, and C.4, we show that as $N \rightarrow \infty$ and d remains
1462 fixed, then $\widehat{\alpha}_N$ is close to α^* with high probability. And under stronger conditions, we demonstrate
1463 that $\widehat{\alpha}_N$ has a limiting normal distribution. The asymptotics of $\widehat{\alpha}_N$ have previously been studied in
1464 the subject of empirical likelihood (Qin & Lawless, 1994; Kitamura & Stutzer, 1997; Owen, 2001).

1465 We first aim to establish that $\widehat{\alpha}_N$ is close to α^* with high probability as $N \rightarrow \infty$ i.e., $\widehat{\alpha}_N$ is
1466 *consistent* for α^* . Define the functions

$$1467 A(\alpha) := A_{p_{\theta_{\text{base}}}}(\alpha), \quad A_N(\alpha) := \log \left(\frac{1}{N} \sum_{n=1}^N \exp\{r_\alpha(\mathbf{x}_n)\} \right),$$

1470 where A_P is defined in Appendix C.1. Observe that A_N is random and depends on the independent
1471 samples $\{\mathbf{x}_n\}_{n=1}^N$ drawn from $p_{\theta_{\text{base}}}$. The dual problem corresponding to $p_{\theta_{\text{base}}}$ maximizes $\alpha^\top \mathbf{h}(\mathbf{x}) -$
1472 $A(\alpha)$, whereas the dual problem corresponding to the distribution of samples $\{\mathbf{x}_n\}_{n=1}^N$ maximizes
1473 $\alpha^\top \mathbf{h}(\mathbf{x}) - A_N(\alpha)$. By the Strong Law of Large Numbers (SLLN), for any $\alpha \in \Xi$, $A_N(\alpha) \rightarrow A(\alpha)$
1474 with $p_{\theta_{\text{base}}}$ probability one. In order for our estimator $\widehat{\alpha}_N$ to approach α^* , though, we need to
1475 argue that the dual objective corresponding to $\{\mathbf{x}_n\}_{n=1}^N$ *uniformly* approaches the dual objective
1476 corresponding to $p_{\theta_{\text{base}}}$ on some neighborhood containing α^* .

1477 **Lemma C.8.** *For any closed, bounded subset K of Ξ ,*

$$1478 \sup_{\alpha \in K} |A_N(\alpha) - A(\alpha)| \rightarrow 0$$

1480 *with $p_{\theta_{\text{base}}}$ probability one.*

1482 *Proof.* By the SLLN, we can construct a Borel set \widetilde{N} of probability zero under $p_{\theta_{\text{base}}}$ such that on
1483 its complement $A_N(\alpha) \rightarrow A(\alpha)$ holds for each $\alpha \in \Xi \cap \mathbb{Q}^d$ (apply the SLLN for an individual
1484 $\alpha \in \Xi \cap \mathbb{Q}^d$, then take a union over probability zero sets).
1485

1486 Rockafellar (1970, Theorem 10.8) states that if a sequence of finite convex functions defined on
1487 an open, convex set C converges pointwise on a dense subset of C to a limiting function, then the
1488 limiting function is convex on C , and the convergence is uniform on closed and bounded subsets of
1489 C . Applying this result to our setting, on the complement of \widetilde{N}

$$1490 \sup_{\alpha \in K} |A_N(\alpha) - A(\alpha)| \rightarrow 0$$

1492 for K a closed and bounded subset of Ξ . □

1494 Once we have proven uniform convergence, our proof of consistency for $\widehat{\alpha}_N$ is nearly identical to
1495 our proof that $\|\alpha_\lambda - \alpha^*\| = o(1)$ in Proposition C.7.

1496 **Proposition C.9** (Consistency of $\widehat{\alpha}_N$). *Suppose Assumptions C.1, C.2, and C.4 hold. For any $\epsilon > 0$,*

$$1498 \mathbb{P}_{p_{\theta_{\text{base}}}}(\|\widehat{\alpha}_N - \alpha^*\| > \epsilon) \rightarrow 0 \quad \text{as } N \rightarrow \infty.$$

1499 *Proof.* From Appendix C.1, we know that both A and A_N are convex functions. Moreover by
1500 Assumption C.4, A is strictly convex.

1502 From Lemma C.8, there exists a closed, bounded subset K of containing α^* on which
1503 $\sup_{\alpha \in K} |A_N(\alpha) - A(\alpha)| \rightarrow 0$ with $p_{\theta_{\text{base}}}$ probability one. Since Ξ is open (Assumption C.2),
1504 K can be chosen to have positive diameter. Fix $\epsilon > 0$ sufficiently small such that the Euclidean ball
1505 centered at α^* of radius ϵ is contained in K . Just as in the proof of Proposition C.7, there exists
1506 some $\kappa \in \mathbb{R}$ such that for all $\|\alpha - \alpha^*\| = \epsilon$,

$$1507 \alpha^\top \mathbf{h}^* - A(\alpha) < \kappa < (\alpha^*)^\top \mathbf{h}^* - A(\alpha^*).$$

1508 Fix $\delta > 0$. By uniform convergence, there exists $N_{\epsilon, \delta} \in \mathbb{N}$ such that $\forall N \geq N_{\epsilon, \delta}$ and for all
1509 $\|\alpha - \alpha^*\| = \epsilon$,

$$1511 \alpha^\top \mathbf{h}^* - A_N(\alpha) < \kappa < (\alpha^*)^\top \mathbf{h}^* - A_N(\alpha^*).$$

with probability at least $1 - \delta$ under $p_{\theta_{\text{base}}}$. And since the dual objective corresponding to $\{\mathbf{x}_n\}_{n=1}^N$ is concave, this implies that, on this event, its maximum occurs in the Euclidean ball of radius ϵ .

In other words, we have proven that for every $\epsilon > 0, \delta > 0$, there exists $N_{\epsilon, \delta}$ such that for every $N \geq N_{\epsilon, \delta}$,

$$\mathbb{P}_{p_{\theta_{\text{base}}}} (\|\widehat{\boldsymbol{\alpha}}_N - \boldsymbol{\alpha}^*\| > \epsilon) \leq \delta.$$

□

Next, we show that under stronger conditions on the problem, $\widehat{\boldsymbol{\alpha}}_N$ has a normal limiting distribution, and we derive its variance.

Proposition C.10 (Asymptotic normality of $\widehat{\boldsymbol{\alpha}}_N$). *Suppose Assumptions C.1, C.2, and C.4 hold. Moreover, assume $2\boldsymbol{\alpha}^* \in \Xi$, for Ξ defined in Appendix C.1. Then the estimator $\widehat{\boldsymbol{\alpha}}_N$ is asymptotically normal:*

$$\begin{aligned} \sqrt{N}(\widehat{\boldsymbol{\alpha}}_N - \boldsymbol{\alpha}^*) &\stackrel{d}{\rightarrow} \mathcal{N}(\mathbf{0}, (\text{Var}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})])^{-1} \boldsymbol{\Sigma} (\text{Var}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})])^{-1}), \\ \boldsymbol{\Sigma} &= \frac{\mathbb{E}_{p_{\theta_{\text{base}}}} [(\mathbf{h}(\mathbf{x}) - \mathbf{h}^*)(\mathbf{h}(\mathbf{x}) - \mathbf{h}^*)^\top \exp\{r_{2\boldsymbol{\alpha}^*}(\mathbf{x})\}]}{(\mathbb{E}_{p_{\theta_{\text{base}}}} [\exp\{r_{\boldsymbol{\alpha}^*}(\mathbf{x})\}])^2}. \end{aligned}$$

Prior to stating the proof of Proposition C.10, we build some intuition by working out the asymptotic variance for the example we presented at the beginning of the section. Recall that the constraint function is $\mathbf{h}(\mathbf{x}) = \mathbb{1}\{\mathbf{x} > 0\}$, \mathbf{h}^* is its target value, and $h_b = \mathbb{P}_{\theta_{\text{base}}}(\mathbf{x} > 0)$ is the expected value of \mathbf{h} under $p_{\theta_{\text{base}}}$. By directly solving for $\widehat{\boldsymbol{\alpha}}_N$ in the expression (5) for the maximum entropy solution, we showed $\widehat{\boldsymbol{\alpha}}_N = \log(\frac{\mathbf{h}^*(1-\bar{\mathbf{y}}_N)}{(1-\mathbf{h}^*)\bar{\mathbf{y}}_N})$, where $\bar{\mathbf{y}}_N = \frac{1}{N} \sum_{n=1}^N \mathbf{y}_n$, $\mathbf{y}_n = \mathbf{h}(\mathbf{x}_n) \stackrel{d}{=} \text{Bernoulli}(h_b)$ for $\mathbf{x}_n \stackrel{i.i.d.}{\sim} p_{\theta_{\text{base}}}$. Next, we compute

$$\begin{aligned} \text{Var}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})] &= \mathbf{h}^*(1 - \mathbf{h}^*) \\ \boldsymbol{\Sigma} &= \frac{(\mathbf{h}^*)^2(1 - h_b) + (1 - \mathbf{h}^*)^2 \exp(2\boldsymbol{\alpha}^*)h_b}{(h_b \exp(\boldsymbol{\alpha}^*) + (1 - h_b))^2} = \frac{(\mathbf{h}^*)^2(1 - h_b) + \frac{(1-h_b)^2(\mathbf{h}^*)^2}{h_b}}{\left(\frac{1-h_b}{1-\mathbf{h}^*}\right)^2} = \frac{(\mathbf{h}^*)^2(1 - \mathbf{h}^*)^2}{h_b(1 - h_b)}. \end{aligned}$$

Combining these two yields the asymptotic variance

$$\text{Var}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})]^{-2} \boldsymbol{\Sigma} = \frac{1}{h_b(1 - h_b)},$$

according to Proposition C.10. In other words, the estimator $\widehat{\boldsymbol{\alpha}}_N$ has greatest asymptotic variance when h_b is close to either 0 or 1. Notice that we can compute the asymptotic variance of $\widehat{\boldsymbol{\alpha}}_N$ directly (i.e., without using Proposition C.10) by applying the delta method to $\bar{\mathbf{y}}_N$ and the function $z \mapsto \log(\frac{\mathbf{h}^*(1-z)}{(1-\mathbf{h}^*)z})$, in which case we obtain the same value.

The proof of Proposition C.10 relies on the technical result Lemma C.11, the statement and proof of which we defer to the end of the section.

Proof of Proposition C.10. Let D_N be the set on which the strong duality holds for $P = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}$ and the dual optimum is uniquely achieved. From the proof of Proposition C.9, we can see $\mathbb{P}_{p_{\theta_{\text{base}}}}(D_N) \rightarrow 1$ as $N \rightarrow \infty$. Moreover, on the set D_N , $\widehat{\boldsymbol{\alpha}}_N$ is the unique root of

$$\frac{1}{N} \sum_{n=1}^N \psi(\mathbf{x}_n, \boldsymbol{\alpha}) = 0, \quad \psi(\mathbf{x}, \boldsymbol{\alpha}) := (\mathbf{h}(\mathbf{x}) - \mathbf{h}^*) \exp\{r_{\boldsymbol{\alpha}}(\mathbf{x})\}.$$

Also, from the proof of Proposition C.9, we know that Assumption C.4 implies $\text{Var}_{p_{\boldsymbol{\alpha}^*}}[\mathbf{h}(\mathbf{x})]$ is positive definite.

By Van der Vaart (2000, Theorem 5.21), if we can show $\psi(\mathbf{x}, \boldsymbol{\alpha})$ satisfies the Lipschitz condition

$$\|\psi(\mathbf{x}, \boldsymbol{\alpha}) - \psi(\mathbf{x}, \boldsymbol{\alpha}')\| \leq M(\mathbf{x}) \|\boldsymbol{\alpha} - \boldsymbol{\alpha}'\| \quad (25)$$

for all α, α' belonging to some neighborhood of α^* and $\mathbb{E}_{p_{\theta_{\text{base}}}} [M(\mathbf{x})^2] < \infty$, then the previous facts imply that $\hat{\alpha}_N$ is asymptotically normal with variance

$$\frac{\mathbb{E}_{p_{\theta_{\text{base}}}} [(\mathbf{h}(\mathbf{x}) - \mathbf{h}^*)\mathbf{h}(\mathbf{x})^\top \exp\{r_{\alpha^*}(\mathbf{x})\}]^{-1} (\mathbb{E}_{p_{\theta_{\text{base}}}} [\exp\{r_{\alpha^*}(\mathbf{x})\}]) \Sigma}{=\text{Var}_{p_{\alpha^*}} [\mathbf{h}(\mathbf{x})]^{-1}} \\ (\mathbb{E}_{p_{\theta_{\text{base}}}} [\exp\{r_{\alpha^*}(\mathbf{x})\}]) (\mathbb{E}_{p_{\theta_{\text{base}}}} [(\mathbf{h}(\mathbf{x}) - \mathbf{h}^*)\mathbf{h}(\mathbf{x})^\top \exp\{r_{\alpha^*}(\mathbf{x})\}]^{-1})^\top . \\ \text{Var}_{p_{\alpha^*}} [\mathbf{h}(\mathbf{x})]^{-1}$$

And so it remains only to establish the Lipschitz condition (25). First, we compute the derivative of ψ with respect to α

$$\nabla_{\alpha} \psi(\mathbf{x}, \alpha) = (\mathbf{h}(\mathbf{x}) - \mathbf{h}^*)\mathbf{h}(\mathbf{x})^\top \exp\{r_{\alpha}(\mathbf{x})\} = \nabla_{\alpha}^2 \exp\{r_{\alpha}(\mathbf{x})\} - \mathbf{h}^* (\nabla_{\alpha} \exp\{r_{\alpha}(\mathbf{x})\})^\top$$

Next, we appeal to Lemma C.11, which tells us that for all α belonging to an open neighborhood of α^* , the derivatives of $\exp\{r_{\alpha}(\mathbf{x})\}$ have norm dominated by a function $M(\mathbf{x})$ that is $p_{\theta_{\text{base}}}$ -square integrable. Also, by the Mean Value Theorem, for all α, α' belonging to this neighborhood,

$$\psi(\mathbf{x}, \alpha) - \psi(\mathbf{x}, \alpha') = \nabla \psi(\mathbf{x}, \tilde{\alpha})(\alpha - \alpha')$$

for some $\tilde{\alpha}$ on the line segment connecting α to α' . By taking the norm on both sides and using $\|\nabla \psi(\mathbf{x}, \tilde{\alpha})\| \leq M(\mathbf{x})$, we obtain the Lipschitz condition (25). \square

Lemma C.11. *Under the assumptions of Proposition C.10, there exists an open neighborhood of α^* on which all derivatives of $\exp\{r_{\alpha}(\mathbf{x})\}$ with respect to α are dominated by a $p_{\theta_{\text{base}}}$ -square integrable function.*

Proof. In Proposition C.10, we assume $\mathbb{E}_{p_{\theta_{\text{base}}}} [\exp\{r_{2\alpha^*}(\mathbf{x})\}] < \infty$; in other words, $2\alpha^*$ is contained in the natural parameter space Ξ . Let ε be defined such that the Euclidean ball of radius ε centered at $2\alpha^*$ is contained in Ξ . Fix any $\tilde{\alpha}$ such that $\|\tilde{\alpha} - \alpha^*\| < \varepsilon/(2d)$, where d is the dimension of the constraint $\mathbf{h}(\mathbf{x})$. Then by Cauchy-Schwarz

$$\exp\{\tilde{\alpha}^\top \mathbf{h}(\mathbf{x})\} \leq \exp\{(\alpha^*)^\top \mathbf{h}(\mathbf{x}) + \varepsilon/(2d)\|\mathbf{h}(\mathbf{x})\|\}. \quad (26)$$

Define the $2d$ vectors $(\beta^{(\pm, l)})_{l=1}^d$ by $\beta^{(+, l)} = e[l]$, $\beta^{(-, l)} = -e[l]$, where $e[l]$ denotes the l th standard basis vector. Then we can upper bound the second term using

$$\exp\{\|\mathbf{h}(\mathbf{x})\|\} \leq \prod_{l=1}^d \exp\{|\mathbf{h}_l(\mathbf{x})|\} \leq \prod_{l=1}^d (\exp\{\mathbf{h}_l(\mathbf{x})\} + \exp\{-\mathbf{h}_l(\mathbf{x})\}) \leq \sum_{l=1}^{2d} 2^d \exp\{d(\beta^{(l)})^\top \mathbf{h}(\mathbf{x})\}.$$

Plugging this bound into (26) yields

$$\exp\{\tilde{\alpha}^\top \mathbf{h}(\mathbf{x})\} \leq \sum_{l=1}^{2d} 2^d \exp\{(\alpha^* + (\varepsilon/2)\beta^{(l)})^\top \mathbf{h}(\mathbf{x})\}. \quad (27)$$

Squaring both sides of (27) yields

$$(\exp\{\tilde{\alpha}^\top \mathbf{h}(\mathbf{x})\})^2 \leq 2^{2d} \sum_{l=1}^{2d} \sum_{k=1}^{2d} \exp\{(2\alpha^* + (\varepsilon/2)(\beta^{(l)} + \beta^{(k)}))^\top \mathbf{h}(\mathbf{x})\}. \quad (28)$$

However, we notice $\|2\alpha^* + (\varepsilon/2)(\beta^{(l)} + \beta^{(k)}) - 2\alpha^*\| \leq \varepsilon$, so each term on the right-hand side of (28) has finite expectation under $p_{\theta_{\text{base}}}$. This implies $\exp\{r_{\alpha}(\mathbf{x})\}$ is dominated by the right-hand side of (27), which is square integrable under $p_{\theta_{\text{base}}}$, for all $\|\alpha - \alpha^*\| < \varepsilon/(2d)$.

As for the derivatives of $\exp\{r_{\alpha}(\mathbf{x})\}$, notice that the k th derivative with respect to α , $\nabla_{\alpha}^{(k)} \exp\{r_{\alpha}(\mathbf{x})\}$, is given by $\mathbf{h}(\mathbf{x})^{\otimes k} \exp\{r_{\alpha}(\mathbf{x})\}$, where \otimes denotes the tensor product. Moreover, by equivalence of norms, for any $\tau > 0$ there exists constants $c_k, c_{\tau, k} \geq 0$ such that $\|\mathbf{h}(\mathbf{x})^{\otimes k}\| \leq c_k \|\mathbf{h}(\mathbf{x})\|^k \leq c_{\tau, k} \exp\{\tau \|\mathbf{h}(\mathbf{x})\|\}$. So by choosing τ such that the Euclidean ball of radius $\varepsilon + 2d\tau$ centered at $2\alpha^*$ is contained in $N(2\alpha^*)$, our same argument yields a dominating function of the form (27) for $\|\alpha - \alpha^*\| < \varepsilon/(2d)$, with exponent $(\alpha^* + (\varepsilon/2 + d\tau)\beta^{(l)})^\top \mathbf{h}(\mathbf{x})$. \square

Remark. Under weaker assumptions (Assumptions C.1 and C.2), the proof of Lemma C.11 implies that the log-normalizer $A_P(\alpha)$ has derivatives of all orders on Ξ . Indeed, this is a consequence of equation (27), which implies that for every α , there exists a neighborhood of α contained in Ξ on which the k th derivative of $\exp\{r_\alpha(\mathbf{x})\}$ is uniformly $p_{\theta_{\text{base}}}$ -dominated. This allows one to exchange differentiation and integration in the definition of $A_P(\alpha)$.

D SIMULATION EXPERIMENTS ADDITIONAL DETAILS

In this section, we provide details for our experiments calibrating mixture proportions in a product of GMMs (Section 3). First, in Appendix D.1 we give background on continuous-time diffusion models, including how we sample from p_θ and compute densities p_θ/p with respect to a dominating measure p . This enables us to employ CGM-relax and CGM-reward for calibrating a pre-trained diffusion model. In Appendix D.2, we describe how the base diffusion model can be initialized to generate exact samples from a GMM or product of GMMs. In Appendix D.3 we provide details regarding our implementation of the CGM calibration algorithm, including optimizer, neural network architecture, and hyperparameters. Finally, in Appendix D.4 we discuss how CGM-relax compares to an augmented Lagrangian method (Hestenes, 1969) on the same set of experiments.

D.1 CONTINUOUS-TIME DIFFUSION MODELS

A continuous-time diffusion model is the solution to the k -dimensional stochastic differential equation (SDE)

$$d\mathbf{x}(t) = \mathbf{b}_\theta(\mathbf{x}(t), t)dt + \sigma(t)d\mathbf{w}(t), \quad \mathbf{x}(0) \sim p_{\text{init}}, \quad (29)$$

where $(\mathbf{w}(t))_{0 \leq t \leq 1}$ is a standard k -dimensional Brownian motion, \mathbf{b}_θ is a neural network drift function, σ is a diffusion coefficient, and p_{init} is a known distribution from which sampling is tractable. Oksendal (2013, Theorem 5.2.1) provides conditions on \mathbf{b}_θ and σ that ensure there exists a unique solution to the SDE (29). We denote the solution, which is a probability distribution on continuous paths, by p_θ , and we write $p_\theta(\mathbf{x}(t))$ for the distribution of the state at time t .

Sampling from diffusion models. To sample from p_θ , we use the Euler-Maruyama method. Specifically, we discretize $[0, 1]$ into T time bins $[0, 1/T], \dots, [(T-1)/T, 1]$ and sample a path $(\hat{\mathbf{x}}(t))_{0 \leq t \leq 1}$ according to $\hat{\mathbf{x}}(0) \sim p_{\text{init}}$

$$\hat{\mathbf{x}}(t + \Delta t) = \hat{\mathbf{x}}(t) + \Delta t \mathbf{b}_\theta(\hat{\mathbf{x}}(t), t) + \sigma(t)\sqrt{\Delta t} \mathbf{z}(t), \quad 0 < \Delta t \leq 1/T \quad (30)$$

for each $t = 0, 1/T, \dots, (T-1)/T$, where $\mathbf{z}(0), \dots, \mathbf{z}((T-1)/T)$ are independent standard multivariate normal random variables. The Euler-Maruyama method with additive noise $\sigma(t)$ has strong order of convergence 1, meaning its error in approximating the solution to the SDE (29) is

$$\mathbb{E}_{p_\theta} [\|\hat{\mathbf{x}}(t) - \mathbf{x}(t)\|] \leq C(T^{-1}), \quad 0 \leq t \leq 1$$

for C a constant independent of T . In other words, as we increase the number of time bins T , we can expect our sample paths drawn according to the Euler-Maruyama scheme to more faithfully approximate samples from the distribution p_θ .

Computing densities. In order to employ CGM-relax and CGM-reward, p_θ and $p_{\theta_{\text{base}}}$ must have densities with respect to one another, and it must be possible to compute these densities. Girsanov’s Theorem (Cameron & Martin, 1944; Girsanov, 1960) provides conditions that guarantee these densities to exist and an expression for computing them.

Theorem D.1 (Girsanov’s Theorem). *Suppose the SDEs*

$$\begin{aligned} \nu_1(\mathbf{x}) : d\mathbf{x}(t) &= \mathbf{b}_1(\mathbf{x}(t), t)dt + \sigma(t)d\mathbf{w}(t), \quad 0 \leq t \leq 1 \\ \nu_2(\mathbf{x}) : d\mathbf{x}(t) &= (\mathbf{b}_1(\mathbf{x}(t), t) + \sigma(t)\mathbf{b}_2(\mathbf{x}(t), t))dt + \sigma(t)d\mathbf{w}(t), \quad 0 \leq t \leq 1 \end{aligned}$$

satisfy $\sigma(t) > 0$, $0 < t < 1$, have the same initial law $\nu_1(\mathbf{x}_0) = \nu_2(\mathbf{x}_0)$, and admit unique, strong solutions, ν_1 and ν_2 . Suppose also

$$\left[\frac{\nu_2(\mathbf{x})}{\nu_1(\mathbf{x})} \right]_t := \exp \left\{ \sum_{i=1}^k \int_0^t \mathbf{b}_2(\mathbf{x}(t), t)[i] d\mathbf{w}^{\nu_1}(t)[i] - \frac{1}{2} \int_0^t \|\mathbf{b}_2(\mathbf{x}(t), t)\|^2 dt \right\} \quad (31)$$

is a ν_1 -martingale, where $(\mathbf{w}^{\nu_1}(t))_{0 \leq t \leq 1}$ is a k -dimensional ν_1 -Brownian motion and $d\mathbf{w}^{\nu_1}(t)[i]$, $i = 1, \dots, k$ denotes the Itô stochastic integral. Then the probability measure ν_2 has a density with respect to ν_1 . In particular, for any bounded functional Φ defined on $C[0, 1]^k$,

$$\mathbb{E}_{\nu_2}[\Phi(\mathbf{x})] = \mathbb{E}_{\nu_1} \left[\Phi(\mathbf{x}) \left[\frac{\nu_2(\mathbf{x})}{\nu_1(\mathbf{x})} \right]_1 \right].$$

If $\|\sigma(t)^{-1}(\mathbf{b}_\theta(\mathbf{x}(t), t) - \mathbf{b}_{\theta_{\text{base}}}(\mathbf{x}(t), t))\|$ is bounded, $([p_\theta(\mathbf{x})/p_{\theta_{\text{base}}}(\mathbf{x})]_t)_{0 \leq t \leq 1}$ is a martingale with respect to $p_{\theta_{\text{base}}}$. Consequently, Girsanov's Theorem tells us that the probability density of p_θ with respect to $p_{\theta_{\text{base}}}$ is given by

$$\frac{p_\theta(\mathbf{x})}{p_{\theta_{\text{base}}}(\mathbf{x})} := \exp \left\{ \sum_{i=1}^k \int_0^1 u_\theta(\mathbf{x}(t), t)[i] d\mathbf{w}^{p_{\theta_{\text{base}}}}(t)[i] - \frac{1}{2} \int_0^1 \|u_\theta(\mathbf{x}(t), t)\|^2 dt \right\}, \quad (32)$$

$$u_\theta(\mathbf{x}(t), t) := \sigma(t)^{-1}(\mathbf{b}_\theta(\mathbf{x}(t), t) - \mathbf{b}_{\theta_{\text{base}}}(\mathbf{x}(t), t))$$

This expression for the density of p_θ with respect to $p_{\theta_{\text{base}}}$ allows us to compute the KL divergence between the probability measures p_θ and $p_{\theta_{\text{base}}}$ according to

$$D_{\text{KL}}(p_\theta \parallel p_{\theta_{\text{base}}}) = \frac{1}{2} \int_0^1 \mathbb{E}_{p_\theta} \|u_\theta(\mathbf{x}(t), t)\|^2 dt.$$

The stochastic integral term vanishes since it has expectation zero.

When $(\hat{\mathbf{x}}(t))_{0 \leq t \leq 1}$ is sampled from the Euler-Maruyama approximation to $p_{\theta_{\text{base}}}$, we approximate (32) by replacing the integrals with

$$\int_0^1 u_\theta(\hat{\mathbf{x}}(t), t)[i] d\mathbf{w}^{p_{\theta_{\text{base}}}}(t)[i] \approx T^{-1/2} \sum_{t=0}^{T-1} u_\theta(\hat{\mathbf{x}}(t/T), t/T)[i] (z((t+1)/T) - z(t/T))$$

$$\int_0^1 u_\theta(\hat{\mathbf{x}}(t), t)^2[i] dt \approx T^{-1} \sum_{t=0}^{T-1} u_\theta(\hat{\mathbf{x}}(t/T), t/T)^2[i]$$

where $z(0), \dots, z((T-1)/T)$ are the same random variables from (30). This same approximation to the density ratio (32) can be derived by writing out the density ratio of $\hat{p}_\theta(\hat{\mathbf{x}}(0), \hat{\mathbf{x}}(1/T), \dots, \hat{\mathbf{x}}(1))$ and $\hat{p}_{\theta_{\text{base}}}(\hat{\mathbf{x}}(0), \hat{\mathbf{x}}(1/T), \dots, \hat{\mathbf{x}}(1))$, where \hat{p}_θ is the probability distribution defined by the Euler-Maruyama discretization of \hat{p}_θ .

Efficient gradient computation. CGM-relax and CGM-reward require computing gradients of the density ratio $\frac{p_\theta(\mathbf{x})}{p_{\text{stop-grad}(\theta)}(\mathbf{x})}$. By applying Girsanov's Theorem to compute the density ratio, differentiating the result, and substituting in our approximations to the integrals, we obtain

$$\nabla_\theta \frac{p_\theta(\mathbf{x})}{p_{\text{stop-grad}(\theta)}(\mathbf{x})} = \sum_{i=1}^k \int_0^1 \nabla_\theta \sigma(t)^{-1} \mathbf{b}_\theta(\hat{\mathbf{x}}(t), t)[i] d\mathbf{w}^{p_{\theta_{\text{base}}}}(t)[i]$$

$$\approx T^{-1/2} \sum_{i=1}^k \sum_{t=0}^{T-1} \nabla_\theta \sigma(t/T)^{-1} \mathbf{b}_\theta(\hat{\mathbf{x}}(t/T), t/Y)[i] (z((t+1)/T)[i] - z(t/T)[i])$$

$$= T^{-1/2} \sum_{t=0}^{T-1} \sigma(t/T)^{-1} \sum_{i=1}^k \nabla_\theta \mathbf{b}_\theta(\hat{\mathbf{x}}(t/T), t/Y)[i] (z((t+1)/T)[i] - z(t/T)[i]). \quad (33)$$

For high-dimensional diffusion models (e.g. Genie2 in our Section 4.1 experiments) memory constraints preclude the naive approach to computing equation (33) by instantiating each term in memory and simultaneously back-propagating gradients through all terms at once. However, because the gradient is a sum across time, it can be computed in chunks. In practice, we divide $\{0, \dots, T\}$ into $\lceil T/\text{chunk_size} \rceil$ blocks of approximately equal size, where `chunk_size` is the largest chunk size that can fit into memory.

Solution to the maximum entropy problem. When the base model $p_{\theta_{\text{base}}}(\mathbf{x})$ constitutes a continuous-time diffusion model (29) satisfying certain regularity properties, and the constraint

function \mathbf{h} depends only on the path at time $t = 1$, there exists a closed-form solution to the maximum entropy problem (4).

Let p be the law of an SDE having diffusion coefficient σ and initial distribution $p'_{\text{init}}(\mathbf{x}(0))$; this is necessary for $p \ll p_{\theta_{\text{base}}}$ by Girsanov’s Theorem (Theorem D.1). By the chain rule for the KL divergence, the objective for the maximum entropy problem defined on the full path measures is

$$\begin{aligned} & \text{D}_{\text{KL}}(p(\mathbf{x}) \parallel p_{\theta_{\text{base}}}(\mathbf{x})) \\ &= \text{D}_{\text{KL}}(p(\mathbf{x}(1)) \parallel p_{\theta_{\text{base}}}(\mathbf{x}(1))) + \mathbb{E}_{p_{\theta}(\mathbf{x}(0))}[\text{D}_{\text{KL}}(p(\cdot|\mathbf{x}(0)) \parallel p_{\theta_{\text{base}}}(\cdot|\mathbf{x}(0)))]. \end{aligned}$$

The KL divergence is computed according to Girsanov’s Theorem.

From here, by the maximum entropy principle applied to the marginal at time $t = 1$, the first term in the objective is lower bounded by

$$\text{D}_{\text{KL}}(p(\mathbf{x}(1)) \parallel p_{\theta_{\text{base}}}(\mathbf{x}(1))) \geq \text{D}_{\text{KL}}(p_{\alpha_0^*}(\mathbf{x}(1)) \parallel p_{\theta_{\text{base}}}(\mathbf{x}(1)))$$

where $p_{\alpha_0^*}(\mathbf{x}(1))$ is the solution to the maximum entropy problem in k -dimensional Euclidean space. Consequently, if we can show that there exists an SDE $p(\mathbf{x})$ satisfying $p(\mathbf{x}(1)) = p_{\alpha_0^*}(\mathbf{x}(1))$ and $p(\cdot|\mathbf{x}(1)) = p_{\theta_{\text{base}}}(\cdot|\mathbf{x}(1))$, then p is the solution to the maximum entropy problem. This is the subject of the following result:

Proposition D.2 (Maximum entropy solution for a diffusion model). *Suppose that the constraint function \mathbf{h} depends only on the value of the path at time $t = 1$ and is bounded and continuous. Moreover, assume that $\mathbf{x}(0) \perp \mathbf{x}(1)$ under the base model $p_{\theta_{\text{base}}}(\mathbf{x})$.*

Then the solution to the maximum entropy problem is a diffusion process

$$p^* : d\mathbf{x}(t) = \{\mathbf{b}_{\theta_{\text{base}}}(\mathbf{x}(t), t) + \sigma(t)\mathbf{u}^*(\mathbf{x}(t), t)\}dt + \sigma(t)d\mathbf{w}(t)$$

satisfying $p^(\mathbf{x}(0)) = p_{\text{init}}$. The drift $\mathbf{u}^*(\mathbf{x}(t), t)$ admits the Feynman-Kac characterization*

$$\mathbf{u}^*(\mathbf{x}, t) = \sigma(t)\nabla_{\mathbf{x}} \log \mathbb{E}_{p_{\theta_{\text{base}}}}[\exp\{r_{\alpha_0^*}(\mathbf{x}(1))\} | \mathbf{x}(t) = \mathbf{x}],$$

where α_0^ are the parameters corresponding to the maximum entropy solution in k -dimensional Euclidean space (4) with base distribution $p_{\theta_{\text{base}}}(\mathbf{x}(1))$ and constraint function $\mathbf{h}(\mathbf{x})$.*

Finally, $p^(\mathbf{x})$ satisfies $p^*(\cdot|\mathbf{x}(1)) = p_{\theta_{\text{base}}}(\cdot|\mathbf{x}(1))$ and $p^*(\mathbf{x}(1)) = p_{\alpha_0^*}(\mathbf{x}(1))$.*

We refer the reader to Domingo-Enrich et al. (2025, Theorem 1) for a proof. The result is a consequence of standard results in the theory of diffusion processes, specifically the Doob h -transform (Oksendal, 2013, Chapter 7).

The assumption that, under $p_{\theta_{\text{base}}}(\mathbf{x})$, the path at time $t = 0$ is independent of the path at time $t = 1$ is necessary to ensure that $p^*(\mathbf{x}(0))$ can be chosen to be equal to $p_{\text{init}}(\mathbf{x}(0))$. This is desirable because, by design, p_{init} is a distribution from which sampling is tractable. However, when the independence assumption does not hold, $p^*(\mathbf{x}(0))$ cannot be chosen to be equal to $p_{\theta_{\text{base}}}(\mathbf{x}(0))$ (Denker et al., 2024, Appendix G.2).

Although, at first glance, this independence assumption may appear strong, Domingo-Enrich et al. (2025, Theorem 1) proves that diffusion models whose initial distribution is Gaussian noise and whose terminal distribution is the data distribution satisfies this property (i.e., variance-preserving SDEs). One example of a commonly used noise schedule satisfying this property is $\sigma(t) = t^{-1}$, which is singular at time $t = 0$.

Proposition D.2 tells us that when we fine-tune a diffusion model with CGM to satisfy a constraint on its terminal distribution $\mathbb{E}_{p_{\theta_{\text{base}}}(\mathbf{x}(1))}[\mathbf{h}(\mathbf{x})] = \mathbf{h}^*$, we expect the terminal distribution of the base model to change to satisfy the constraint, while the conditional path distribution given the endpoint should be preserved. In other words, seeking the distribution over paths that is closest in KL distance to the base model amounts to shifting the terminal distribution while leaving the conditional distributions unchanged.

D.2 INITIALIZING THE BASE DIFFUSION MODEL TO SAMPLE A GAUSSIAN MIXTURE

In each of our synthetic data experiments, we initialize our base diffusion model $p_{\theta_{\text{base}}}$ such that $p_{\theta_{\text{base}}}(\mathbf{x}(1))$ is equal to a GMM. We achieve this by representing $p_{\theta_{\text{base}}}$ as the reversal of a forward

diffusion process. A forward diffusion process draws samples from the target GMM density $\mathbf{x}(1) \sim p_{\text{target}}$ and then noises them according to the linear SDE

$$\vec{p} : d\mathbf{x}(t) = \frac{1}{2}\kappa(t)\mathbf{x}(t)dt + \sigma(t)d\mathbf{w}(t), \quad 0 \leq t \leq 1. \quad (34)$$

When the diffusion coefficient is chosen such that $\sigma(t) = \sqrt{\kappa(t)}$ and the linear coefficient $(\kappa(t))_{0 \leq t \leq 1}$ satisfies $\kappa(t) \geq 0$, $\int_0^1 \kappa(t)dt = +\infty$, then $\vec{p}(\mathbf{x}(0)) \stackrel{d}{=} \mathcal{N}(\mathbf{0}, \mathbb{I})$. We choose $\kappa(t) = t^{-1}$. Simply, (34) turns samples from p_{target} into Gaussian noise. In practice, since the drift and diffusion coefficients defined by $\kappa(t)$ are unbounded (which violates the assumptions for existence and uniqueness of the solution to the SDE from Appendix D.1), we cap $\kappa(t)$ at some large M .

A foundational result in diffusion processes (Anderson, 1982) states that the reversal of (34) is another diffusion process that is given by

$$\overleftarrow{p} : d\mathbf{x}(t) = \left\{ \sigma(t)^2 \nabla_{\mathbf{x}} \log \vec{p}(\mathbf{x}(t)) + \frac{1}{2}\kappa(t)\mathbf{x}(t) \right\} dt + \sigma(t)d\mathbf{w}(t), \quad 0 \leq t \leq 1 \quad (35)$$

with $\overleftarrow{p}(\mathbf{x}(0)) \stackrel{d}{=} \vec{p}(\mathbf{x}(0))$. The probability distributions defined by (34) and (35) are equal in law. $\nabla_{\mathbf{x}} \log \vec{p}(\mathbf{x}(t))$ is called the *score* of the forward process (34).

Equation (35) is useful since it tells how to generate samples from p_{target} : first draw samples from $\vec{p}(\mathbf{x}(0)) \approx \mathcal{N}(\mathbf{x}(0) | \mathbf{0}, \mathbb{I})$, then solve the SDE (35) numerically using Euler-Maruyama, for example. However, for general target distributions p_{target} , the score of the forward process is intractable, which yields the backward diffusion process (35) also intractable.

In the case of a GMM, though, the score of the forward process is tractable. Indeed, for $p_{\text{target}}(\mathbf{x}(1)) = \sum \pi_i \mathcal{N}(\mathbf{x}(1) | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, we compute

$$\begin{aligned} \vec{p}(\mathbf{x}(t)) &= \int \vec{p}(\mathbf{x}(t) | \mathbf{x}(1)) p_{\text{target}}(\mathbf{x}(1)) d\mathbf{x}(1) \\ &= \sum \pi_i \int \mathcal{N}(\mathbf{x}(t) | m(t)\mathbf{x}(1), s(t)^2 \mathbb{I}) \mathcal{N}(\mathbf{x}(1) | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) d\mathbf{x}(1) \\ &= \sum \pi_i \mathcal{N}(\mathbf{x}(t) | m(t)\boldsymbol{\mu}_i, s(t)^2 \mathbb{I} + m(t)^2 \boldsymbol{\Sigma}_i), \end{aligned}$$

where $m(t)$ and $s(t)$ are defined by the forward diffusion process $(\kappa(t))_{0 \leq t \leq 1}$. For $\kappa(t) = t^{-1}$, we have $m(t) = t^{1/2}$ and $s(t) = (1-t)^{1/2}$. Using this expression for $\vec{p}(\mathbf{x}(t))$, we initialize $p_{\theta_{\text{base}}}(\mathbf{x})$ to the exact reversal of the forward process (34) according to (35).

D.3 EXPERIMENTAL DETAILS

We perform all synthetic data experiments using Adam (Kingma & Ba, 2015) with default momentum hyperparameters $\beta = (0.9, 0.999)$ and a cosine decay learning rate schedule (Loshchilov & Hutter, 2016). We perform 2×10^3 CGM iterations for every experiment. We train on a single H100 GPU.

In the diffusion generative model, we parameterize the drift function \mathbf{b}_{θ} as

$$\mathbf{b}_{\theta}(\mathbf{x}(t), t) = \sigma(t)^2 \{ \nabla_{\mathbf{x}} \log \vec{p}(\mathbf{x}(t)) - u_{\theta}(\mathbf{x}, t) \} + \frac{1}{2}\kappa(t)\mathbf{x}(t).$$

u_{θ} is a neural network with two hidden layers of dimension 256 and SiLU activations, and $\log \vec{p}(\mathbf{x}(t))$ is the analytical score of the forward process that we described in Appendix D.2. In addition to $\mathbf{x}(t)$, we feed as input to u_{θ} a sinusoidal time embedding of dimension 32. By initializing the weights of the output layer of u_{θ} to zero, we ensure that p_{θ} is initialized at $p_{\theta_{\text{base}}}$, the reversal of the forward diffusion process \vec{p} .

All synthetic data experiments are performed with batch size $M = 10^4$. For CGM-relax, we select λ by first performing calibration for each λ on a log linear grid from 10^0 to 10^{-3} with 10 grid points. We choose the value of λ for which $(\mathcal{L}^{\text{viol}})^{1/2}$ is reduced by a factor of 10 and \mathcal{L}^{KL} is the smallest. If no such value exists, we choose λ for which \mathcal{L}^{KL} is smallest. For CGM-reward, we compute $\hat{\alpha}_N$ using $N = 10^5$ samples from $p_{\theta_{\text{base}}}$.

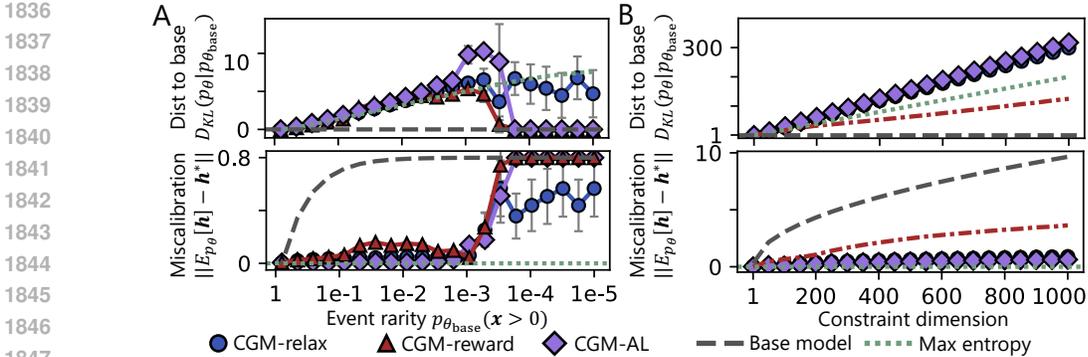


Figure 6: Comparison of CGM-relax and CGM-reward to the augmented Lagrangian algorithm (‘CGM-AL’) for calibrating a diffusion model that generates samples from a GMM. For CGM-relax and CGM-reward, we report the same values as in Figure 2. For CGM-AL, we report results with $\ell = 20$ in panel A and with $\ell = 100$ in panel B (best values). We find that CGM-AL performs comparably to CGM-relax, except in the rare event setting where it performs worse (panel A). **A**: Calibrating the mixture proportion of a rare mode in a 1D GMM. **B**: Calibrating a k -dimensional GMM to a k -dimensional constraint.

D.4 COMPARISON TO AUGMENTED LAGRANGIAN METHOD

As an alternative to CGM-relax, we consider a variant of the augmented Lagrangian (AL) algorithm (Hestenes, 1969), which solves a general constrained optimization problem.

In particular, suppose one would like to solve

$$\min_{z \in \mathbb{R}^p} f(z), \quad \text{subject to } c(z) = c^*,$$

where $f: \mathbb{R}^p \rightarrow \mathbb{R}$ is the objective function and $c: \mathbb{R}^p \rightarrow \mathbb{R}^d$ is the constraint function. The AL method forms the ‘augmented Lagrangian’ $\mathcal{L}_\lambda: \mathbb{R}^p \times \mathbb{R}^d \rightarrow \mathbb{R}$ with penalty parameter λ , defined

$$\mathcal{L}_\lambda(z, \mathbf{u}) = f(z) + \mathbf{u}^\top (c(z) - c^*) + \frac{\lambda}{2} \|c(z) - c^*\|^2. \quad (36)$$

\mathcal{L}_λ differs from the ordinary Lagrangian due the presence of the penalty term $\frac{\lambda}{2} \|c(z) - c^*\|^2$. $\mathbf{u} \in \mathbb{R}^d$ are the dual variables. The AL algorithm alternates between minimizing \mathcal{L}_λ with respect to z and updating the dual variables:

$$\begin{aligned} z^{(k)} &\leftarrow \arg \min_{z \in \mathbb{R}^p} \mathcal{L}_\lambda(z, \mathbf{u}^{(k)}) \\ \mathbf{u}^{(k+1)} &\leftarrow \mathbf{u}^{(k)} + \lambda(c(z^{(k)}) - c^*). \end{aligned} \quad (37)$$

Some variants of the AL method also update the penalty λ according to some schedule. The augmented Lagrangian algorithm can be viewed as a proximal-point algorithm applied to the dual function (Rockafellar, 1976)

$$\mathbf{u}^{(k+1)} \leftarrow \arg \max_{\mathbf{u} \in \mathbb{R}^d} \left\{ d(\mathbf{u}) - \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{u}^{(k)}\| \right\}, \quad d(\mathbf{u}) = \min_{z \in \mathbb{R}^p} f(z) + \mathbf{u}^\top (c(z) - c^*).$$

In the setting of the calibration problem (1), we identify $z = \theta$, $f(\theta) = D_{\text{KL}}(p_\theta \| p_{\theta_{\text{base}}})$, $c(\theta) = \mathbb{E}_{p_\theta}[\mathbf{h}(x)]$, $c^* = \mathbf{h}^*$. However, the augmented Lagrangian does not admit a closed-form minimizer with respect to the model parameters θ . Consequently, the AL algorithm as stated in equation (37) cannot be applied. Instead, we propose alternating between performing a stochastic gradient update to the model parameters θ and updating the Lagrange multipliers \mathbf{u} . This introduces an additional algorithmic hyperparameter ℓ representing how frequently (i.e., after how many iterations) the dual variables are updated. We provide pseudocode for our implementation in Algorithm 3, which we refer to as ‘CGM-AL’.

We compare CGM-AL against CGM-relax and CGM-reward on the two synthetic experiments described in Section 3 and Appendix D.3. For CGM-AL, we select λ by performing a 10-point grid

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

Algorithm 3 CGM-AL fine-tuning (augmented Lagrangian)

Require: $p_{\theta_{\text{base}}}$, $\mathbf{h}(\cdot)$, \mathbf{h}^* , M , λ , ℓ

▷ Initialize model and dual variables

$p_{\theta} \leftarrow p_{\theta_{\text{base}}}$, $\mathbf{u} \leftarrow \mathbf{0}$, $t \leftarrow 0$

while not converged **do**

▷ Sample and compute importance weights

$\mathbf{x}_{1:M} \stackrel{i.i.d.}{\sim} p_{\text{stop-grad}(\theta)}$

$w_m \leftarrow p_{\theta}(\mathbf{x}_m) / p_{\text{stop-grad}(\theta)}(\mathbf{x}_m)$

▷ KL loss with LOO baseline

$l_m \leftarrow \log(p_{\text{stop-grad}(\theta)}(\mathbf{x}_m) / p_{\theta_{\text{base}}}(\mathbf{x}_m))$

$l_m^{\text{LOO}} \leftarrow l_m - \frac{1}{M-1} \sum_{m' \neq m} l_{m'}$

$\widehat{\mathcal{L}}^{\text{KL}} \leftarrow \frac{1}{M} \sum_{m=1}^M w_m l_m^{\text{LOO}}$

▷ Constraint violation loss

$\mathbf{h}_m \leftarrow w_m (\mathbf{h}(\mathbf{x}_m) - \mathbf{h}^*)$

$\widehat{\Delta \mathbf{h}} \leftarrow \frac{1}{M} \sum_{m=1}^M \mathbf{h}_m$

$\widehat{\mathcal{L}}^{\text{viol}} \leftarrow \|\widehat{\Delta \mathbf{h}}\|^2 - \frac{1}{M} \widehat{\text{Var}}[\mathbf{h}_{1:M}]$,

$\widehat{\text{Var}}[\mathbf{h}_{1:M}] = \frac{1}{M-1} \sum \|\mathbf{h}_m - \widehat{\Delta \mathbf{h}}\|^2$

▷ Augmented Lagrangian loss

$\widehat{\mathcal{L}}^{\text{AL}} \leftarrow \widehat{\mathcal{L}}^{\text{KL}} + \mathbf{u}^{\top} \widehat{\Delta \mathbf{h}} + \frac{\lambda}{2} \widehat{\mathcal{L}}^{\text{viol}}$

▷ Primal update

$\theta \leftarrow \text{gradient-step}(\theta, \nabla_{\theta} \widehat{\mathcal{L}}^{\text{AL}})$

▷ Dual update every k iterations

if $t \bmod \ell = 0$ **then**

$\mathbf{u} \leftarrow \mathbf{u} + \lambda \text{stop-grad}(\widehat{\Delta \mathbf{h}})$

$t \leftarrow t + 1$

search on $[10^0, 10^2]$, and we consider two potential values for the Lagrange multiplier update frequency, $\ell \in \{20, 100\}$. We again perform stochastic gradient updates to θ using Adam and a cosine decay learning rate schedule with initial learning rate 10^{-4} and final learning rate 10^{-7} .

From Figure 6, we observe that CGM-relax performs comparably to CGM-AL and, in the rare event mode reweighting example (Figure 6A), CGM-relax outperforms CGM-AL. While one might think that poor conditioning of the relax loss landscape for small λ would result in inferior performance of CGM-relax to CGM-AL, we observe this is not the case. We attribute this to our choice of λ via grid search: as described in Appendix D.3, we select λ to balance between constraint violation and KL distance to the base model $p_{\theta_{\text{base}}}$. Since the AL method introduces nontrivial computational overhead in the choice of the dual parameter update frequency ℓ , we prefer CGM-relax.

E CASE STUDY ADDITIONAL DETAILS

In this section, we describe the experimental setup for our case studies with CGM-relax and CGM-reward from Section 4. We provide explanations regarding the generative model classes p_{θ} , CGM constraint functions \mathbf{h} and targets \mathbf{h}^* , choice of CGM hyperparameters λ and N , model architectures, and training procedures. We also include additional samples from our models before and after calibration.

Just as in our synthetic data experiments, we perform all experiments using Adam with default momentum hyperparameters $\beta = (0.9, 0.999)$ and a cosine decay learning rate schedule. Additional common training details are shown in Table 1. We train all models on a single H100 GPU.

Table 1: Training configurations for experiments. Batch (sub-batch) indicates the number of samples per batch and the sub-batch size used to fit gradient computations into memory. \mathcal{S}_V^L denotes the set of all sequences with vocabulary size V and length L .

Hyperparameter	Genie2	ESM3-open	TarFlow	TinyStories-33M
Initial learning rate	10^{-5}	10^{-4}	10^{-6}	2×10^{-6}
Batch (sub-batch)	64 (16)	256 (64)	256 (16)	512 (64)
Training steps	100	100	50	200
x Space	$(\mathbb{R}^{100 \times 3})^{100}$	$(\mathcal{S}_{4096}^{100})^{50}$	$\mathbb{R}^{256 \times 256 \times 3}$	$\mathcal{S}_{10,000}^{200}$
Constraint dims (k)	99	99	5	8
Model parameters	15M	1.4B	463M	33M
Training time (hrs)	48	2.3	3	0.1

E.1 CALIBRATING GENIE2

For our experiments with Genie2, we represent p_θ as a continuous-time diffusion model defined over three-dimensional protein backbone coordinates with drift function defined by the SE(3)-equivariant encoder-decoder architecture from Lin et al. (2024a).

Since Genie2 is trained as the reversal of a discrete-time noising process (a DDPM, see Ho et al., 2020), we first convert the discrete-time denoising diffusion model to a (continuous-time) diffusion model. We achieve this by redefining the final timestep T of the original denoising process to be time 1 of the continuous-time process. To define the drift function, we take the DDPM transition mean defined at each time t in the discrete-time process, divide it by $1/T = T$, and define the drift function to be equal to the resulting value in between times t/T and $(t + 1)/T$. The diffusion coefficient is similarly defined by the DDPM transition standard deviation at each time t in the discrete-time process, but is instead scaled by $T^{1/2}$. This approach of converting the DDPM into a continuous-time diffusion model ensures that when the SDE is solved under the Euler-Maruyama scheme using a grid of T timesteps (i.e., the original time grid used to define the DDPM), one samples from the original DDPM.

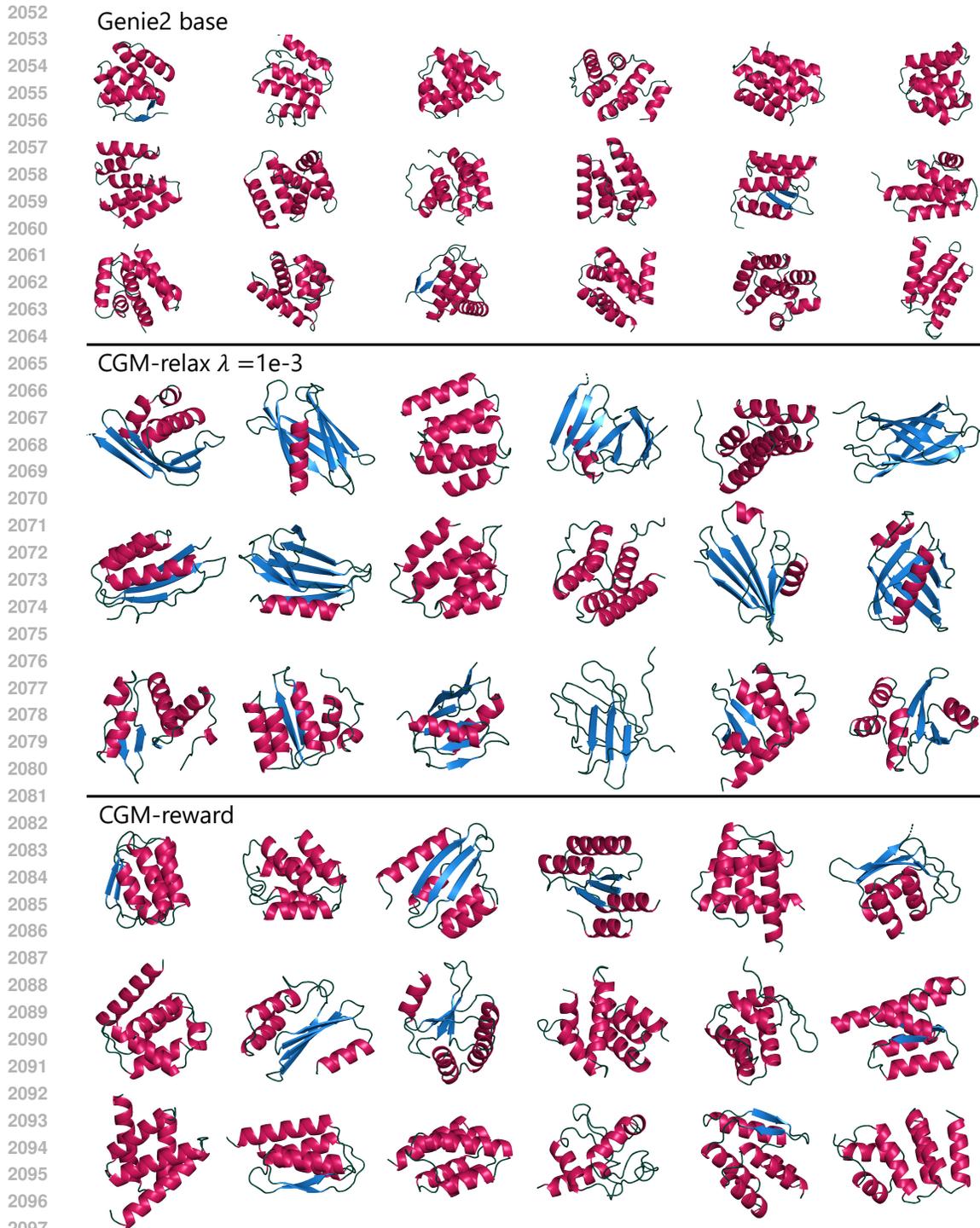
We perform sampling using 10^2 timesteps and a non-uniform time grid: we sample the first 50 steps on the interval $[0, 0.05]$ and the remaining steps on the interval $[0.05, 1]$. We point out that the original Genie2 model was trained with 10^3 denoising steps; we find that reducing the number of sampling steps dramatically decreases the runtime of CGM calibration. Our sampling scheme is possible since we redefined the base generative model to be a continuous-time diffusion process. We computed self-consistency metrics for the base Genie2 model sampled on the original time grid (with 10^3 steps) and on our proposed grid (with 10^2 steps), we did not observe any difference in sample quality.

For CGM-relax, we calibrate to $k=99$ constraints on the bivariate CDF of alpha helix and beta strand content. And for CGM-reward, we calibrate to $k=15$ constraints using $N = 2.5 \times 10^4$ samples from $p_{\theta_{\text{base}}}$. Since sampling from the Genie2 base model is time intensive, the sampling cost to compute $\hat{\alpha}_N$ (with small variance) is a downside to CGM-reward. Results reported in Figure 3 are averages over 3 trials, with two standard errors.

Self-consistency RMSD and design failures. To assess the quality of our generations, we compute the root mean-square deviation (RMSD) between C_α atoms resulting from (i) unfolding our generated structures into predicted amino sequences, (ii) refolding each of these predicted sequences into a protein structure, and (iii) aligning the predicted structures to the original structure. The self-consistency RMSD (scRMSD) is defined as the smallest RMSD between the given structure and one of the corresponding predictions. We use ProteinMPNN (Dauparas et al., 2022) for our inverse folding model and ESMFold (Lin et al., 2023) for our folding model; we compute scRMSD from 8 sequences. The pipeline we employ was developed by Lin et al. (2024b). Once we have determined the scRMSD of a generated structure, we classify it as a “design failure” if its scRMSD is greater than 2\AA . Intuitively, designability is a binary measure of whether or not a structure could have been plausibly produced by folding an amino acid sequence.

1998 **Secondary structure annotation.** As discussed in Section 4.1, we measure the diversity of a collec-
1999 tion of protein structures by computing the proportion of residues that lie in each of the three protein
2000 secondary structure types. For the CATH domains and Genie2, we perform annotations using the
2001 Biotite package (Kunzmann & Hamacher, 2018), which considers only C_α backbone atoms. For the
2002 CATH proteins (Sillitoe et al., 2021), we obtain the secondary structure distribution by annotating
2003 the domains collected and published by Ingraham et al. (2019).

2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051



2098 Figure 7: Random samples from the Genie2 model before calibration (top), after calibration using
2099 CGM-relax with 99 bivariate CDF constraints (middle), and after calibration using CGM-reward
2100 with 15 bivariate CDF constraints (bottom).

2101
2102
2103
2104
2105

E.2 CALIBRATING ESM3-OPEN

In order to apply CGM to ESM3, we need to be able to sample from the model and to compute gradients of sample log-probabilities with respect to the model’s parameters.

Sampling method. Following the method used by Hayes et al. (2025), sampling is achieved by treating the model as a discrete time Markov chain that starts at a sequence of mask tokens and ends at a sequence of fully unmasked structure tokens. Each step i of the chain consists of three steps and transitions from state $\mathbf{x}(i-1)$ to $\mathbf{x}(i)$.

1. Pick token indices $U(i)$ to unmask uniformly at random without replacement from the masked tokens of $\mathbf{x}(i-1)$.
2. Use the model π_θ to predict a categorical distribution $\pi_\theta^{(j)}(\cdot | \mathbf{x}(i-1))$ for $j \in U(i)$ for each of the newly unmasked tokens given the previous partially masked state.
3. Sample the values of those tokens from the predicted categorical distributions, resulting in $|U(i)|$ more unmasked tokens.

As implemented by Hayes et al. (2025), we use $T = 50$ steps to sample 100-residue sequences and follow a cosine unmasking schedule. The cosine schedule determines the number of masked positions at each sampling step as

$$r(i) := \text{round} \left(100 \times \cos \left(\frac{\pi i}{2T} \right) \right), \quad i = 0, \dots, T.$$

Early sampling steps unmask few tokens per step, while later ones sample many at once. Intuitively, this let’s the model sample more tokens in parallel once it has more information to predict the final sequence. Note that the number of tokens unmasked at step $i > 0$ is $|U(i)| = r(i-1) - r(i)$.

Transition probabilities. A Markov chain can be characterized by its initial state distribution, $\mathbf{x}(0) \sim \pi_0(\mathbf{x}(0))$ and its transition probabilities for going from one state to the next. The ESM3 sampling method starts fully masked, so has initial distribution $\pi_0(\mathbf{x}(0)) = \mathbb{1}\{\mathbf{x}(0) \text{ is fully masked}\}$. The transition probabilities follow from the sampling procedure and are

$$p_\theta(\mathbf{x}(i) | \mathbf{x}(i-1)) = C(i) \prod_{j \in U(i)} \pi_\theta^{(j)}(\mathbf{x}(i)[j] | \mathbf{x}(i-1)), \quad (38)$$

where $C(i)$ is a constant that accounts for randomly choosing which tokens to unmask. $C(i)$ does not depend on θ or the sampling trajectory $(\mathbf{x}(0), \mathbf{x}(1), \dots, \mathbf{x}(T))$, since every unmasking order is equally likely. Note $U(i)$ can be computed from $\mathbf{x}(i-1)$ and $\mathbf{x}(i)$ by finding which tokens are masked in $\mathbf{x}(i-1)$ and not in $\mathbf{x}(i)$.

Trajectory log-probability. As in the neural SDE setting (Appendix D.1), the marginal likelihood of \mathbf{x}_T is intractable, so we treat samples \mathbf{x} as entire trajectories, $\mathbf{x} = (\mathbf{x}(0), \mathbf{x}(1), \mathbf{x}(2), \dots, \mathbf{x}(T))$. Using the Markov property, the log-probability of a trajectory is

$$\begin{aligned} \log p_\theta(\mathbf{x}) &= \log \left(\pi_0(\mathbf{x}_0) \prod_{i=1}^T p_\theta(\mathbf{x}(i) | \mathbf{x}(i-1)) \right) \\ &= \log \pi_0(\mathbf{x}_0) + \sum_{i=1}^T \log \left(C(i) \prod_{j \in U(i)} \pi_\theta^{(j)}(\mathbf{x}(i)[j] | \mathbf{x}(i-1)) \right) \quad (\text{by equation (38)}) \\ &= \sum_{i=1}^T \log \left(C(i) \prod_{j \in U(i)} \pi_\theta^{(j)}(\mathbf{x}(i)[j] | \mathbf{x}(i-1)) \right) \quad (\pi_0(\mathbf{x}(0)) = 1 \text{ by construction}) \\ &= \sum_{i=1}^T C(i) + \sum_{i=1}^T \sum_{j \in U(i)} \log \pi_\theta^{(j)}(\mathbf{x}(i)[j] | \mathbf{x}(i-1)). \end{aligned}$$

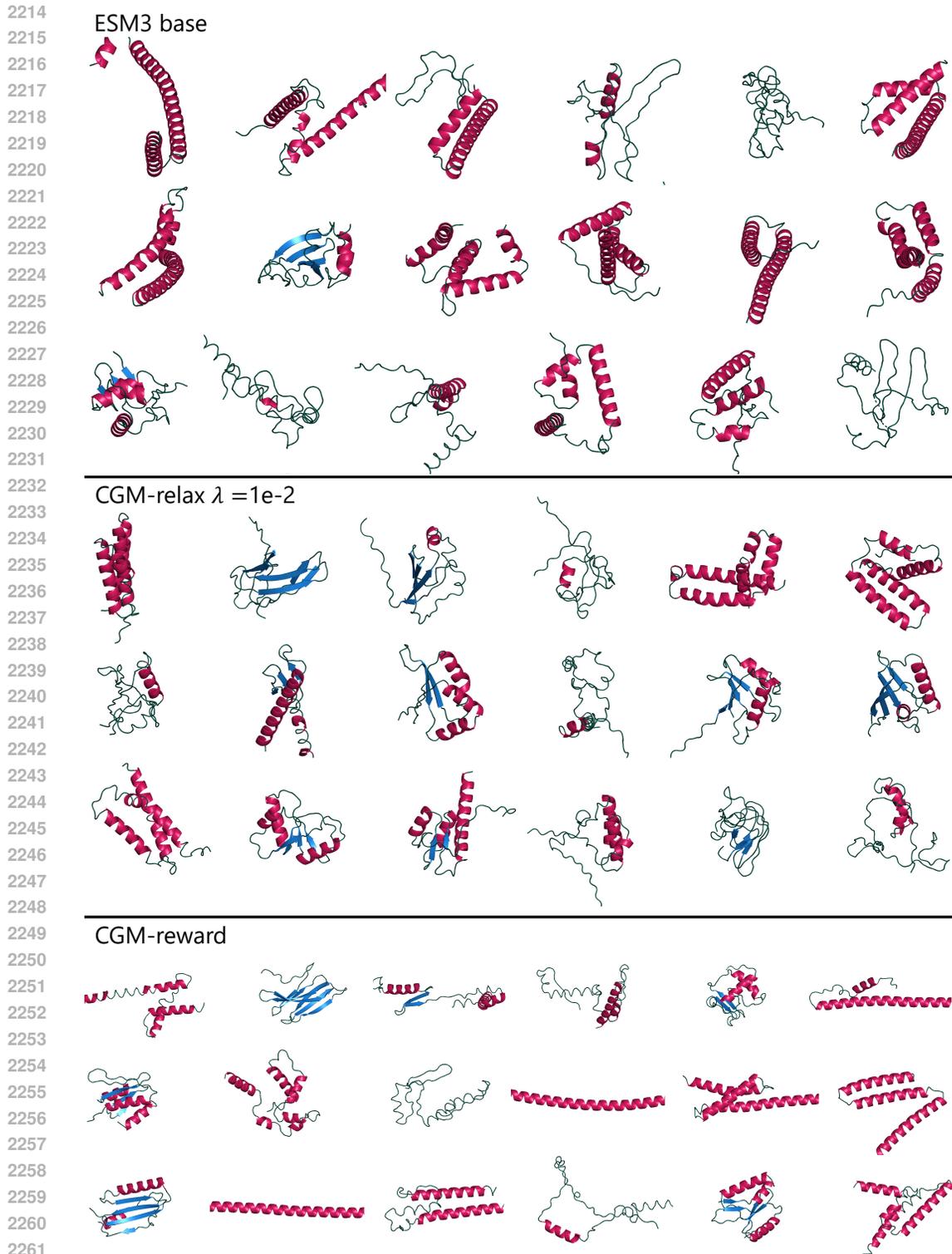
2160 **Parameter gradients.** Now that we have defined the log-probability of \mathbf{x} , we can compute gradients
2161 with respect to θ as

$$\begin{aligned} 2162 & \\ 2163 & \nabla_{\theta} \log p_{\theta}(\mathbf{x}) = \nabla_{\theta} \left(\sum_{i=1}^T C(i) + \sum_{i=1}^T \sum_{j \in U(i)} \log \pi_{\theta}^{(j)}(\mathbf{x}(i)[j] \mid \mathbf{x}(i-1)) \right) \\ 2164 & \\ 2165 & \\ 2166 & = \sum_{i=1}^T \sum_{j \in U(i)} \nabla_{\theta} \log \pi_{\theta}^{(j)}(\mathbf{x}(i)[j] \mid \mathbf{x}(i-1)), \\ 2167 & \\ 2168 & \end{aligned}$$

2169 which conveniently is a sum over sampling steps. The decomposition of the gradient into a sum
2170 over sampling steps lets us compute parameter gradients using constant memory with respect to the
2171 number of sampling steps, which is critical for high-parameter-count models such as ESM3-open.

2172 **Secondary structure annotation.** We use the ESM3 structure decoder and the ESM3 function
2173 `ProteinChain.infer_oxygen` to get heavy atom coordinates from sampled structure tokens.
2174 We then pass the coordinates to the Python package `PyDSSP` (Minami, 2023) to annotate secondary
2175 structure.
2176

2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213



2262 Figure 8: Random samples from the ESM3-open model before calibration (top), after calibration
 2263 using CGM-relax with 99 bivariate CDF constraints (middle), and after calibration using CGM-
 2264 reward with 15 bivariate CDF constraints (bottom).

2265
 2266
 2267

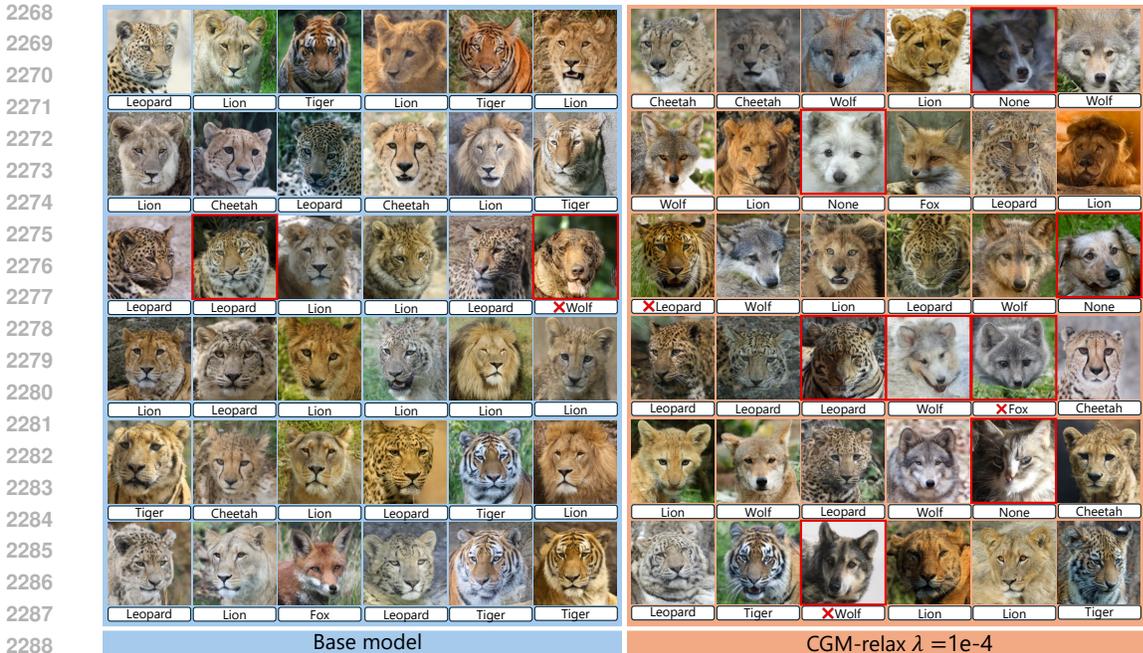


Figure 9: Random samples from the conditional TarFlow model trained on the AFHQ dataset (blue background) and the same model fine-tuned using CGM-relax ($\lambda=10^{-4}$) (orange background), with annotations by GPT o5-mini (white box). Red boxes denote poor-quality samples, and red crosses denote incorrect annotations. Although the model calibrated with CGM-relax produces animals with more balanced class proportions, it also produces fewer realistic samples.

E.3 CALIBRATING TARFLOW

As we described in Section 4.2, our goal when calibrating the TarFlow model (Zhai et al., 2025), trained conditionally on the Animal Faces HQ (AFHQ) dataset (Choi et al., 2020), is to generate more diverse samples from the wildlife class. By directly examining the AFHQ dataset, we identify six animals: $\{\text{lion, tiger, wolf, fox, leopard, cheetah}\}$; we do not further distinguish among these animals e.g., leopard versus snow leopard. Within the AFHQ training dataset, these animals are represented in the wildlife class with proportions $\{0.2615, 0.2254, 0.0897, 0.0933, 0.2003, 0.1290\}$, as annotated by GPT o5-mini.

Our motivation for choosing this problem was twofold. First, the quality of images generated by the base TarFlow model is high, such that a pre-trained classifier could attain high accuracy without fine-tuning. Second, we observe that the wildlife images generated by the base TarFlow model contained predominantly lions and leopards (Figure 9), and rarely contained foxes or wolves. From 5×10^3 samples annotated by GPT o5-mini, we computed animal proportions $\{0.3590, 0.1260, 0.0404, 0.0256, 0.2704, 0.1752\}$.

For image classification, we queried GPT o5-mini to classify each image according to the following prompt:

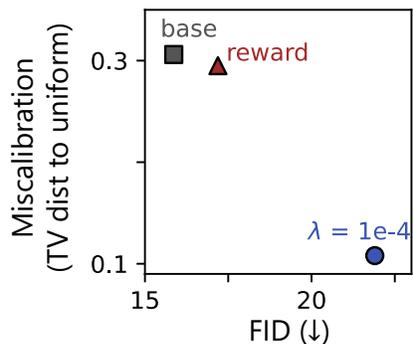
```
You are labeling animal photos.
Return JSON only: {"label": <one of the options>, "confidence": <0..1>}.
Choose exactly one from: lion, tiger, wolf, fox, leopard, cheetah or none.
```

Although we require the model to state its confidence when labeling the images, we do not use these confidence scores for fine-tuning. The calibration function $\mathbf{h}(\mathbf{x}) \in \mathbb{R}^5$ is a one-hot encoding of the first five classes. Since nearly all of the samples from the base TarFlow model are labeled as one of the six animals, we observe that choosing a six-dimensional constraint (i.e., adding cheetah) results in a poorly conditioned dual problem (7), since then the components of \mathbf{h} are nearly linearly dependent. Our target is the uniform distribution over animals $\mathbf{h}^* = [.167, .167, .167, .167, .167]$.

2322 We perform calibration with regularization parameter $\lambda = 10^{-4}$, and sample size $N = 5 \times 10^3$.
 2323 We assess the success of calibration using two metrics: the total variation (TV) distance of the
 2324 distribution over animal proportions (using 5×10^3 annotated images) to the uniform distribution
 2325 and the FID, computed using only the samples belonging to the wildlife class in the AFHQ training
 2326 dataset. We use 5×10^4 samples from the generative model to compute FID. It is important to note
 2327 that the FID is an imperfect metric for assessing the quality of generated images since it will be lower
 2328 for models whose animal class makeup is similar to that of the training distribution. To account for
 2329 this, we evaluate CGM-relax on the maximum entropy reweighting of the training dataset to the
 2330 uniform distribution over animal classes. In other words, we up or down weighted images belonging
 2331 to a particular animal class in order to sample the six animals belonging to wildlife class with equal
 2332 probability.

2333 Our best model, calibrated using CGM-relax with $\lambda = 10^{-4}$, obtains class proportions
 2334 $\{0.2248, 0.0854, 0.1750, 0.1566, 0.1668, 0.1086\}$, evaluated using 5×10^3 samples from the model;
 2335 0.0828 of the samples were labeled as None. CGM-relax reduces the miscalibration error by nearly
 2336 three times, from a TV distance of .306 to .108 (Figure 10). However, the FID score increases from
 2337 15.9 to 21.9. CGM-reward is unsuccessful at calibrating the base model; both miscalibration and
 2338 FID is roughly the same as the base model. Since CGM-reward remains close to the base model, we
 2339 evaluate FID on the original AFHQ training dataset.

2340 In Figure 9, we provide random generations from both the pre-trained and the model calibrated with
 2341 CGM-relax with $\lambda = 10^{-4}$. By examining samples from the calibrated model, we observe two axes
 2342 along which sample quality worsens after calibration. First, some of the samples (those labeled as
 2343 None) are dogs or cats, which lie outside the AFHQ wildlife class. Second, a greater proportion of
 2344 samples depict blends of multiple animals.



2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356 Figure 10: Calibrating TarFlow with CGM-relax reduces the TV distance of animal class labels to
 2357 the uniform distribution by approximately three times. However, CGM-relax also produces fewer
 2358 realistic samples, as measured by FID.
 2359

2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375

2376 E.4 CALIBRATING TINYSTORIES-33M

2377

2378

2379

E.4.1 AUTOREGRESSIVE SAMPLING AND LOG-LIKELIHOODS

2380

2381

2382

2383

2384

2385

Following the setup used by Eldan & Li (2023), we sample in the standard autoregressive fashion with no temperature scaling. To compute sequence log-likelihoods, we consider the prompt as given and ignore tokens generated after the first end-of-sequence (EOS) token. Let m be the length of the prompt and n the index of the first EOS token. Then sequence \mathbf{x} has log-probability

2386

2387

2388

2389

$$\log p_{\theta}(\mathbf{x}) = \sum_{i=m+1}^n \log p_{\theta}(\mathbf{x}(i) | \mathbf{x}(<i)).$$

2390

2391

2392

2393

2394

For computational efficiency during training and evaluation, we set the maximum length of each story to be 200 tokens.

2395

E.4.2 TINYSTORIES CONSTRAINT DEFINITION

2396

2397

2398

2399

To calibrate TinyStories-33M, we use a simple heuristic procedure to detect the gender of the story’s character associated with the profession in the prompt. The procedure returns 1 for female, -1 for male, and 0 if the gender cannot be determined. Given a generated story, our procedure is as follows.

2400

2401

2402

(1) Pronoun at sentence two. If the second sentence begins with a third-person singular pronoun, we assign gender based on that pronoun. This is common with our prompt templates, e.g., “Once upon a time there was a doctor named Sam. **She** was very kind...”.

2403

2404

2405

2406

2407

(2) First-sentence scan. If step (1) is inconclusive, we iterate through the words in the story’s first sentence. If a title (“Mr.,” “Mrs.,” “Miss,” “Ms.”) appears, we assign the corresponding gender. Otherwise, we treat each word as a potential first name and query the `gender-guesser` package (Pérez et al., 2016). If the package classifies the token as “male”, “mostly male”, “female”, or “mostly female”, we assign the corresponding gender; otherwise we continue scanning.

2408

2409

(3) No evidence. If no gender is detected, we assign 0 (unknown).

2410

2411

We acknowledge the limitations of this simple approach but consider it sufficient for a proof of concept.

2412

2413

2414

Conditional constraint via sum-of-losses. We wish to satisfy the conditional calibration constraints

2415

2416

$$\mathbb{E}[\mathbf{h}(\mathbf{x}) | \text{prompt}_i] = 0, \quad \text{for } i = 1, \dots, k$$

2417

2418

2419

2420

2421

which encodes that the male and female labels should be balanced *for each* of the k professions. We implement this as a sum of CGM losses $\sum_{i=1}^k \widehat{\mathcal{L}}_i$, where $\widehat{\mathcal{L}}_i$ is the reward or relax loss for the conditional generative model $p_{\theta}(\mathbf{x} | \text{prompt}_i)$. During every training batch, we sample 64 stories for each of the eight professions, resulting in a total batch size of 512.

2422

2423

2424

E.4.3 TINYSTORIES-33M FIGURE DETAILS

2425

2426

2427

2428

2429

Both panels of Figure 5 were created using 20 replicates per model with different Pytorch seeds, with points indicating the mean metric value across replicates. 2048 stories were sampled per profession for each replicate, resulting in $14 \times 2048 = 28672$ stories per model. Due to the high number of samples and replicates, two-times-standard-error-of-the-mean error bars are smaller than the markers, so are not shown.

Khalifa et al. (2021) baseline. Khalifa et al. (2021) use the same method as CGM-reward to define an approximate target distribution $p_{\hat{\alpha}_N}$. Unlike CGM-reward, they minimize the forward KL

$$\begin{aligned}
D_{\text{KL}}(p_{\hat{\alpha}_N} \parallel p_{\theta}) &= \mathbb{E}_{p_{\hat{\alpha}_N}} \left[\log \frac{p_{\hat{\alpha}_N}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right] \\
&= \mathbb{E}_{p_{\text{stop-grad}(\theta)}} \left[\frac{p_{\hat{\alpha}_N}(\mathbf{x})}{p_{\text{stop-grad}(\theta)}(\mathbf{x})} \log \frac{p_{\hat{\alpha}_N}(\mathbf{x})}{p_{\theta}(\mathbf{x})} \right] \quad (\text{change of measure}) \\
&= \mathbb{E}_{p_{\text{stop-grad}(\theta)}} \left[-\frac{p_{\hat{\alpha}_N}(\mathbf{x})}{p_{\text{stop-grad}(\theta)}(\mathbf{x})} \log p_{\theta}(\mathbf{x}) \right] + C \\
&= \mathbb{E}_{p_{\text{stop-grad}(\theta)}} \left[-\frac{p_{\theta_{\text{base}}}(\mathbf{x}) \exp \left\{ \hat{\alpha}_N^{\top} \mathbf{x} - A_{p_{\theta_{\text{base}}}}(\hat{\alpha}_N) \right\}}{p_{\text{stop-grad}(\theta)}(\mathbf{x})} \log p_{\theta}(\mathbf{x}) \right] + C \quad (\text{definition of } p_{\hat{\alpha}_N}) \\
&= K \mathbb{E}_{p_{\text{stop-grad}(\theta)}} \left[-\frac{p_{\theta_{\text{base}}}(\mathbf{x}) \exp \left\{ \hat{\alpha}_N^{\top} \mathbf{x} \right\}}{p_{\text{stop-grad}(\theta)}(\mathbf{x})} \log p_{\theta}(\mathbf{x}) \right] + C,
\end{aligned}$$

where C and K are constants that do not depend on θ . C can be ignored as it has no affect on parameter gradients, and K can be absorbed into the learning rate. Similar to CGM-reward, gradients of this KL-divergence are estimated using Monte Carlo. For a fair comparison, we use the same $\hat{\alpha}_N$ and batch size to train Khalifa et al. (2021), CGM-reward, and CGM-relax.

Distance from base-model (symmetrized KL) definition. For each fine-tuned model, we sample $N = 2048$ stories $\{\mathbf{x}_i\}_{i=1}^N$ per profession, and compute log-probabilities $\log p_{\theta}(\mathbf{x}_i)$ and base-model log-probabilities $\log p_{\theta_{\text{base}}}(\mathbf{x}_i)$. We estimate the per-profession backward KL as

$$D_{\text{KL}}(p_{\theta} \parallel p_{\theta_{\text{base}}}) \approx \frac{1}{N} \sum_{i=1}^N \log \frac{p_{\theta}(\mathbf{x}_i)}{p_{\theta_{\text{base}}}(\mathbf{x}_i)}.$$

The forward KL uses importance sampling estimate

$$D_{\text{KL}}(p_{\theta_{\text{base}}} \parallel p_{\theta}) \approx \frac{1}{N} \sum_{i=1}^N \frac{p_{\theta_{\text{base}}}(\mathbf{x}_i)}{p_{\theta}(\mathbf{x}_i)} \log \frac{p_{\theta_{\text{base}}}(\mathbf{x}_i)}{p_{\theta}(\mathbf{x}_i)}.$$

We add our estimates for the forward and backward KL for each profession to get the symmetrized KL, then report the average symmetrized KL across all eight professions.

Gender imbalance definition. For each model replicate, we compute the number of male (#male) and number of female (#female) characters in 2048 samples for each profession. The miscalibration for a single profession is defined as

$$\left| \frac{\#\text{male} - \#\text{female}}{\#\text{male} + \#\text{female}} \right|,$$

which takes maximum value 1 if all samples used the same gender and minimum value 0 if there are an equal number of each gender. The overall miscalibration values shown on the y-axis of Figure 5A were computed by taking the average miscalibration for the eight professions used for fine-tuning.

Estimating KL for the max-entropy solution. Using $N = 2048$ samples from the base-model for each profession, we compute an estimate $\hat{\alpha}$ of α^* for each profession using the procedure outlined in Appendix C.2. We then compute importance weights for each sample $w_i = \exp(\hat{\alpha}^{\top} \mathbf{x}_i)$, which we use to compute normalized weights $\tilde{w}_i = w_i / \sum_{j=1}^N w_j$. To estimate the forward and backward KL divergences, we treat base-model probabilities as $1/N$ and use

$$D_{\text{KL}}(p_{\theta_{\text{base}}} \parallel p_{\hat{\alpha}}) \approx \frac{1}{N} \sum_{i=1}^N \log \frac{1/N}{\tilde{w}_i}, \quad D_{\text{KL}}(p_{\hat{\alpha}} \parallel p_{\theta_{\text{base}}}) \approx \frac{1}{N} \sum_{i=1}^N \frac{\tilde{w}_i}{1/N} \log \frac{\tilde{w}_i}{1/N}.$$

This procedure was repeated 20 times for different sampling seeds to ensure variance from estimating α had little effect on the outcome. The resulting error bars for the symmetrized KL are smaller than the marker size.

2484 E.4.4 EXAMPLE GENERATIONS
2485

2486 We provide example generations for four prompts before and after fine-tuning with CGM-relax with
2487 $\lambda = 0.1$.

2488 Listing 1: Samples from TinyStories-33M
2489

2490 Once upon a time there was a doctor named Jack. He was fit and strong, loved helping people. One day, Jack was
2491 working on an operation. He soon came across a little girl named Mary who was scared. Even though Jack
2492 told her not to worry, he said "hello" to her. Mary was still scared, and the doctor could tell her not
2493 to worry. He said that everything would be okay for her to do the operation. Mary felt a little bit
2494 better. Jack helped Mary and made sure she was fit. He even helped her little brother and made sure that
2495 she was always safe. At the end of the operation, Mary was smiling. She thanked Jack for his help and for
2496 helping her son. She had learned lessons about not being scared and having someone stronger than her.

2497 Once upon a time there was a lawyer named Matt. He was very new and people stopped to talk to him. He came
2498 across a big building which said "Quiet down." Matt thought it was a bad idea but he was very curious. So
2499 he did as he was told and soon his house became quiet. All of a sudden he heard a voice and it said, "
2500 What are you doing?" It was the jail officer! Matt replied, "I thought it was a good idea to reverse the
2501 law in this house." The jail officer said, "Smart Jackie, you know a lot! This house is the law of the
2502 world." Matt smiled and said, "But it was fun to reverse the law." The jail officer nodded and said, "Yes
2503 it was. But the law hates them." Matt thanked the jail officer and went on his way. The moral of the story
2504 is that it pays to be creative, but also to be careful when you
2505

2506 Once upon a time there was a teacher named Miss Jane. She loved to answer any questions and her class was
2507 empty. One day, when Miss Jane was busy writing something on the top of her desk, she heard a voice
2508 coming from her kitchen floor. She looked down to see a small mouse squeezing a piece of cheese into the
2509 board. Miss Jane gasped and said, "mouse, that's naughty! Put that cheese down!" The mouse replied, "No,
2510 I won't. I'm making sure no one steals the cheese." Miss Jane was getting frustrated and said, "That's
2511 not right. You must ask before you litter into someone else's work. That's not nice." The mouse knew from
2512 the right place and quickly put the cheese back. Miss Jane saw what he had done and said, "That's the
2513 right thing to do!" Missy stared at him in amazement and said, "I suppose you won't tell everyone what to
2514 do with it again!"

2515 Once upon a time there was a pilot named Bobby. He flew a plane on his travels. One day, Bobby was flying in
2516 his plane when he saw a tree fall. He crashed the plane right into the tree. Bobby felt so sorry. He
2517 called for help, but no one was around. Suddenly, he saw a 3-year-old child standing alone above the tree
2518. Bobby flew to the child and offered to help. He told the little girl not to worry and helped him onto
2519 his plane. Together, they flew the plane back to its home in the tree. Bobby and the 3-year-old child
2520 became good friends after that until finally Bobby had to fly on time. The end.
2521

2522 Listing 2: Samples from CGM-relax ($\lambda = 0.1$) model

2523 Once upon a time there was a doctor named Susie. She was a very patient and kind lady who always supported
2524 others. Every day she took her to see lots of patients. She always tried to help them feel better. She
2525 was always so kind and always willing to lend a hand. One day, a patient called Sally came to see Susie.
2526 Sally was very sick and needed help. Susie was always so gentle and caring, not taking extra time to
2527 help others. When Sally asked her if she was okay, Susie just smiled and reassured her that she would be
2528 alright soon. She was so patient with Sally as she played and chased around the toy store. It was so
2529 nice to see how patient and caring Susie is between helping others. She continues to be sure to always
2530 lend a hand whenever she can and remember to be there for her when feeling ill.

2531 Once upon a time there was a lawyer named Cass. Cass and Tim went to a big city. They decided to sit on the
2532 corner. It was a little dizzy when they sat because they were too dizzy. At the office, Mommy and Daddy
2533 talked to the boss. They asked the lawyer for help. They listened to his ideas and his ideas. The lawyer
2534 gave them big smiles. He showed them how to mix some ingredients together to make a cookie. The lawyer's
2535 ideas and his ideas fit the cookie dough. His ideas were so neat and the cookie dough was extra yummy!
2536 The lawyer thanked him and said, "Thanks for your help. You kept your ideas nice and happy." Cass and
2537 Tim smiled as they watched the lawyer leave the office. They waved goodbye as they walked away, still
2538 feeling happy each

2539 Once upon a time there was a teacher named Mr. Jam and Mrs. Bunny. Mr. Bunny has very cool haircuts that
2540 everyone wanted. One day Mrs. Bunny said, "Tomorrow we will have a new surprise." The next morning, Mrs.
2541 Bunny came over to Rob's house. When they got to the house Mrs. Bunny said, "I have a new surprise for
2542 you!" She pulled out a big bicycle ticket. Mr. Bunny said, "This is for you my grandkids. Now settle down
2543 and take it for a ride." The little bunnies were so happy! They sat quietly and quietly, taking turns on
2544 the bicycle. Then Mrs. Bunny gave them a surprise, a cool drink. The little bunnies were so excited.
2545 They each had a cool drink! All the bunnies were so happy to be outside on such a cool day. And they
2546 thanked Mrs. Bunny for the special surprise.

2547 Once upon a time there was a pilot named Judy. She loved to measure everything: trees, houses, farms, and
2548 anything else. She was really careful to measure each inch so that every time that she got a bit closer
2549 to the number. One day Judy was flying to measure a star in the sky. She drew a line through the line
2550 with her finger and measured it for a long time. She was almost done when something terrible happened.
2551 One of the letters flew too close and hit Judy in the face. It hurt a lot! Judy was very frightened and
2552 ran away from the window to escape from the bad hit. But unfortunately the bad letter kept rolling
2553 closer and closer until it was right at the edge of the world. Judy was so scared and upset she couldn't
2554 believe it! She had been measured and measured, but still got past the bad letter. The bad letter that
2555 Judy had measured was gone forever and Judy was left feeling very sad and lonely.
2556
2557