

---

# Elucidating the Design Space of Generative Models for Single-Cell Perturbation Prediction

---

Anonymous Authors<sup>1</sup>

## Abstract

We introduce `ExpressionVAE`, the first discrete-latent perturbation model for single-cell data: a vector-quantized variational autoencoder paired with a perturbation-conditioned discrete prior. On Replogle and Parse 1M it achieves state-of-the-art on every distributional and cell-eval state metric we evaluate, with order-of-magnitude gaps on Fréchet distance and MMD<sup>2</sup> over the strongest continuous-latent baseline. We test two prior families (autoregressive and masked discrete diffusion) and find they achieve effectively identical numbers, isolating the gain to the discrete latent. A controlled output-head ablation further reveals a single design axis governing decoder-head choice, the richness of the inference-time sampling distribution, with standard evaluation metrics partitioning into three groups whose rankings flip along it. Finally, on a held-out CRISPRi reversion benchmark of 1,732 perturbations under inflammatory cytokine stress, the frozen encoder effectively matches scGPT model (trained on 10× larger dataset) on biological selectivity.

## 1. Introduction

A cell can be viewed as an information-processing system that maps environmental inputs to molecular responses through an internal regulatory program (Fitch, 2021). While the structure of the regulatory network itself is not directly accessible, single-cell RNA sequencing provides a high-dimensional partial observation of its activity (Macosko et al., 2015; Gulati et al., 2020; Zeng & Dai, 2019). Crucially, this observation is discrete: we do not measure continuous expression levels, but count individual molecular events. Across many cells and conditions, these observations implicitly encode the statistical dependencies induced

by the underlying regulatory structure. The objective of a virtual cell model is to learn these dependencies well enough to generate realistic cellular states and to predict how modifications to the regulatory program, such as genetic or chemical perturbations, alter the resulting transcriptome (Dixit et al., 2016; Adamson et al., 2016). If reliable, virtual cells stand to enable a new frontier of in-silico drug discovery.

The central object of observed gene expression is the count matrix: a sparse, zero-inflated, high-dimensional integer-valued tensor recording how many mRNA molecules of each gene were captured from each cell. These properties make direct generative modeling on raw counts remarkably difficult for most models. Autoregressive factorizations impose a sequential ordering on genes, an inductive bias misaligned with their natural exchangeability. Flow matching tends to collapse to the mean under the sparsity of count data. Masked diffusion is the only family that recovers reasonable structure on raw counts (Zhang et al., 2026; Bhattacharya et al., 2026).

These difficulties on raw counts motivate a latent-space approach. Autoregressive models have demonstrated predictable scaling laws on language data (Hoffmann et al., 2022), an appealing property as virtual cell models scale up, and they operate naturally over discrete tokens. This motivates a discrete latent representation.

We introduce `ExpressionVAE`, a vector-quantized variational autoencoder that encodes each cell as a fixed-length sequence of discrete codes, paired with a perturbation-conditioned prior over those codes. Vector quantization commits the encoder to a finite codebook of cell-state representations that any standard discrete prior can model exactly. To control for the choice of prior, we train two from different generative-model families, autoregressive (AR) and masked discrete diffusion (MDLM); any gain that survives the swap is attributable to the latent representation rather than to a particular prior architecture.

Beyond the latent and prior, two further design questions recur in the single-cell perturbation literature, and we treat each as a controlled experiment. First, single-cell models vary substantially in their decoder output head, and published rankings across cross-entropy, MSE, hurdle, and negative-binomial heads are inconsistent in ways that often

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Submitted to the 2026 Workshop on Generative and Agentic AI for Biology (ICML 2026). Do not distribute.

track the choice of evaluation metric rather than any property of the heads themselves. We run a within-architecture head ablation on the same trained backbone to expose the underlying axis. Second, we evaluate the frozen encoder on an out-of-distribution CRISPRi reversion benchmark to test whether the resulting design choices reflect a generalizable inductive bias rather than over-fitting our controlled metric suite.

Concretely, our contributions are:

- **Discrete-latent perturbation modeling.** We introduce the first discrete-latent representation learning approach for single-cell perturbation modeling, and show that pairing `ExpressionVAE` with either an autoregressive or a masked-diffusion prior achieves state-of-the-art on the distributional and cell-eval state metrics we evaluate, on both Replogle and Parse 1M, often by an order of magnitude on MMD<sup>2</sup> and Fréchet distance over the strongest continuous-latent baseline. The two priors achieve essentially identical performance, isolating the gain to the discrete latent representation rather than to the prior architecture (Section 4.1).
- **A measurement-theoretic account of decoder-head choice.** We show that the apparent conflicts between published evaluation metrics resolve into a clean three-group structure (distribution and DE-discovery, mean-ranking, and a hybrid) once a single design axis, the richness of the inference-time sampling distribution, is made explicit. A within-architecture experiment using a parametric negative-binomial head with two inference rules (sample vs. mean) that share trained weights isolates this axis from confounds in capacity, parameterization, and training objective, and yields a metric-to-head mapping for downstream practitioners (Section 4.2).
- **Out-of-distribution generalization.** On a held-out CRISPRi reversion benchmark of 1,732 perturbations under inflammatory cytokine stress (Wong et al., 2026), the hurdle decoder identified by our controlled ablation is the strongest pairing on biological selectivity for both training datasets, and the frozen encoder effectively matches a 10×-larger foundation model (scGPT) on this OOD perturbation-ranking task, evidence that the discrete-latent perturbation prior captures inflammation-axis structure efficiently rather than memorizing the in-distribution metric set (Section 4.3).

## 2. Related Work

We draw from two literatures: single-cell generative and perturbation-response models, and the latent modeling toolkit from vision and language.

**Generative and perturbation-response models.** scVI (Lopez et al., 2018a;b) established the canonical Gaussian-latent, Negative Binomial-likelihood VAE recipe, elaborated by the conditional-VAE thread of scGen, trVAE, and CPA (Lotfollahi et al., 2019; 2020; 2023) and the latent-diffusion thread of scDiffusion (Luo et al., 2024), CFGen (Palma et al., 2024), and scLDM (Palla et al., 2025) with flow-matching priors (Lipman et al., 2022; Peebles & Xie, 2023); the Hurdle head we evaluate traces to MAST (Finak et al., 2015). Beyond this lineage, perturbation-response models pretrained on tens to hundreds of millions of cells (Geneformer (Theodoris et al., 2023), scGPT (Cui et al., 2024), scFoundation (Hao et al., 2024), STATE (Adduri et al., 2025), Tahoe-x1 (Gandhi et al., 2025)) increasingly condition on ESM-2 embeddings (Lin et al., 2023). All embed raw counts into continuous vectors. A concurrent line generates directly in the discrete count domain, via Score Entropy Discrete Diffusion (DCM (Bhattacharya et al., 2026; Lou et al., 2023)), absorbing-state masked diffusion in the D3PM (Austin et al., 2021) lineage (Lingshu-Cell (Zhang et al., 2026), X-Cell (Wang et al., 2026a), scDiVa (Wang et al., 2026b)), or Stack’s in-context approach (Dong et al., 2026). We evaluate on Replogle Perturb-seq (Replogle et al., 2022) and the Parse cytokine atlas (Oesinghaus et al., 2025), with phenotypic reversion (Wong et al., 2026) as our out-of-distribution NFκB experiment.

**Latent generative modeling: tokenizers and priors.** On the tokenizer side, VQ-VAE (van den Oord et al., 2017; Razavi et al., 2019) and MaskGIT (Chang et al., 2022) introduced learnable discrete codes and masked parallel decoding, Latent Diffusion (Rombach et al., 2022) established the continuous-latent plus diffusion-prior recipe, and FSQ (Mentzer et al., 2023) eliminated the codebook collapse documented in single-cell VQ (Liu et al., 2025). On the prior side, the absorbing-state masked discrete diffusion family has matured rapidly from D3PM (Austin et al., 2021) through SEDD (Lou et al., 2023), MDLM (Sahoo et al., 2024), and MD4 (Shi et al., 2024) to RADD (Ou et al., 2025), which unifies it with any-order autoregression. We combine learned discrete latents with flexible generative priors, a combination no existing model explores.

## 3. Method

We propose a two-stage pipeline (Fig. 1). Stage 1 trains an `ExpressionVAE` that compresses each cell’s gene-expression vector into a fixed-length sequence of discrete codes via a finite scalar quantization (FSQ) bottleneck (Mentzer et al., 2023). Stage 2 freezes the VAE and trains a perturbation-conditioned generative prior over the code sequence, conditioned on an ESM2-3B protein-language-model embedding of the targeted gene and an additive cell-type signal. At inference, the prior is queried with the em-

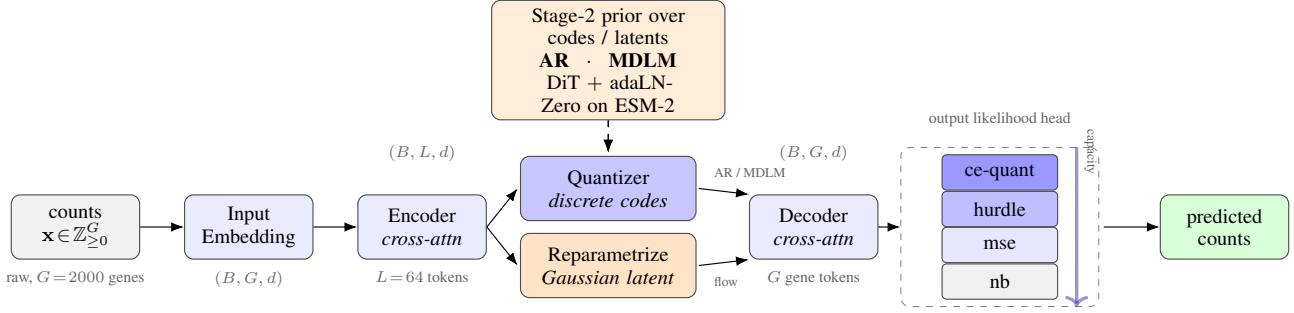


Figure 1. A two-stage pipeline. Stage 1 (the row of boxes): a cross-attention transformer encoder maps raw count vectors  $\mathbf{x} \in \mathbb{Z}_{\geq 0}^G$  to  $L=64$  latent tokens, which are either FSQ quantised into discrete codes (consumed by AR or MDLM priors) or reparametrised as continuous Gaussian latents (consumed by flow). A symmetric cross-attention decoder maps these back to per-gene representations, which one of four output likelihood heads converts to predictions. Stage 2 (sidebars): a DiT prior over the frozen codes / latents, conditioned on the ESM-2 perturbation embedding and optionally cell-type embedding via adaLN-Zero.

bedding of an unseen perturbation and the resulting latent is decoded into a predicted gene-expression profile. The two design axes varied across the experiments are the VAE *output head* (Section 4.2) and the choice of *prior* (Section 4.1).

### 3.1. Stage 1: ExpressionVAE

**Encoder and decoder.** The encoder maps a count vector to a set of  $L$  latent tokens via a cross-attention transformer. Each gene  $g$  is first embedded as  $\mathbf{e}_g = f_{\text{in}}(x_g, \mathbf{v}_g)$ , where  $f_{\text{in}}$  is an MLP applied to the  $\log_{1p}$ -normalized scalar count  $x_g$  and modulated by a learned gene-identity token  $\mathbf{v}_g \in \mathbb{R}^d$ .  $L$  learnable latent query tokens cross-attend over the  $G$  gene embeddings through a stack of transformer layers, producing a per-cell latent matrix  $\mathbf{Z} \in \mathbb{R}^{L \times d}$ . A symmetric cross-attention transformer decodes the quantized codes back to gene space:  $G$  learnable gene query tokens cross-attend over the projected code embeddings, yielding a per-gene representation  $\mathbf{H} \in \mathbb{R}^{G \times d}$ .

**FSQ bottleneck.** The encoder output  $\mathbf{Z}$  is quantized per-dimension by FSQ (Mentzer et al., 2023), which replaces the learned VQ codebook with a parameter-free per-dimension rounding scheme. Given a level vector  $\ell = (\ell_1, \dots, \ell_d) \in \mathbb{Z}_{>0}^d$ , each dimension is bounded via  $\tanh$  and rounded to the nearest integer on a uniform grid:

$$\hat{z}_k = \text{round}\left(\frac{\ell_k - 1}{2} \cdot \tanh(z_k)\right), \quad k = 1, \dots, d.$$

The result is a per-cell sequence of  $L$  discrete codes  $\mathbf{c} = (c_1, \dots, c_L)$  drawn from a vocabulary of size  $\prod_k \ell_k$ , eliminating the codebook collapse documented in single-cell VQ (Liu et al., 2025). We use  $L = 64$  tokens with codebook size  $K = 512$ . A Gaussian-latent variant used in Section 4.3 also uses  $L = 64$  latent tokens with  $\beta_{\text{KL}} = 2 \times 10^{-6}$ .

**Output heads.** An output head maps the per-gene representation  $\mathbf{H}$  to a conditional distribution  $p_\theta(\mathbf{x} | \mathbf{H})$ .

We evaluate four head families, corresponding to standard choices in the single-cell modeling literature, and refer to them by the labels used in Section 4.2.

**ce-quant** places a categorical over  $[0, c_{\text{max}}]$  per gene with logits produced by a linear projection of  $\mathbf{h}_g$ , trained against quantile-binned integer targets and sampled at inference (Section 3.1 below).

**hurdle** (Finak et al., 2015) parameterises a Bernoulli zero-gate  $z_g \sim \text{Bern}(1 - p_{0,g})$  together with a regressed positive magnitude  $\mu_g$ ; the loss combines BCE on the gate with MSE on non-zero entries, and inference samples the gate while emitting  $\mu_g$  deterministically.

**mse** regresses  $\log_{1p}$ -normalised counts directly, with  $p(x_g | \mathbf{h}_g) = \mathcal{N}(\mu_g, 1)$ , and emits  $\mu_g$  with no stochastic step.

The negative binomial (NB) family parameterises  $\text{NB}(\mu_g, \theta_g)$  per gene with library-size factorisation, following (Lopez et al., 2018a; Palla et al., 2025):  $\mu_g = \ell_n \cdot \rho_g$ , where  $\boldsymbol{\rho} = \text{softmax}_g(\mathbf{w}^\top \mathbf{H})$  is a probability simplex over genes,  $\ell_n = \sum_g x_g$  is the per-cell library size, and  $\theta_g$  is a learned cell-independent dispersion. This family admits two distinct `predict()` rules at inference time without changing trained weights: `nb-sample` draws from  $\text{NB}(\mu, \theta)$ , and `nb-deter` returns  $\mu$  directly. We treat these as two separate evaluation points throughout, since the choice of inference rule materially shifts which evaluation metrics the head wins (Section 4.2).

**Quantile binning for the categorical head.** The `ce-quant` head requires integer targets, so we quantile-bin the raw counts into a vocabulary of 20 bins chosen from the empirical training count distribution: 8 singleton bins for counts 0–7 capture the dominant low-count regime exactly, 7 bins covering 8–63 handle moderate counts at progressively coarser resolution, and 5 bins (64–95, 96–127, 128–191, 192–255,  $\geq 256$ ) cover the high-expression tail. Tokenisation maps each raw count to its bin index; detokenisation

maps each predicted bin back to the empirical conditional mean of counts within that bin range, not the arithmetic midpoint, since the within-bin distribution is heavy-tailed and the bias matters most for the unbounded  $\geq 256$  bucket. This scheme captures 100% of top-DEG per-perturbation means on Parse and 99.69% on Replogle within bounded bins.

### 3.2. Stage 2: Discrete prior over latent codes

Stage 2 freezes the VAE and trains a generative prior over the code sequence  $\mathbf{c} = (c_1, \dots, c_L)$ , conditioned on a projected ESM2-3B perturbation embedding  $\mathbf{p}$  and an additive cell-type signal. Both prior families share a DiT backbone with adaLN-Zero conditioning on  $\mathbf{p}$  and the cell-type vector. The autoregressive (AR) prior factorises as  $p(\mathbf{c} | \mathbf{p}) = \prod_{l=1}^L p(c_l | c_{<l}, \mathbf{p})$  via a causally-masked DiT, requiring  $L$  sequential decoding steps at inference. MDLM (Sahoo et al., 2024) trains a bidirectional DiT denoiser under an absorbing-state masking schedule  $\alpha_t = 1 - t$ , with the standard masked-token loss

$$\mathcal{L}_{\text{MDLM}} = \mathbb{E}_{t, \mathbf{z}_t} \left[ \frac{1}{t} \sum_{i: z_t^i = [M]} -\log p_\theta(c^i | \mathbf{z}_t, \mathbf{p}) \right],$$

and denoises all  $L$  tokens in a single forward pass. We use the same backbone, parameter count, and training schedule for both, so any gap in Section 4.1 reflects the prior family rather than capacity.

### 3.3. Inference

Given an unseen perturbation, the corresponding ESM2-3B embedding  $\mathbf{p}^*$  is passed to the trained prior, which samples a latent sequence  $\hat{\mathbf{c}} \sim p_\theta(\cdot | \mathbf{p}^*)$  via  $L$  autoregressive steps for AR or a single masked-denoising pass for MDLM. The frozen VAE decoder maps  $\hat{\mathbf{c}}$  back to gene space, and the chosen output head emits a predicted profile  $\hat{\mathbf{x}} \in \mathbb{R}^G$  following its `predict()` rule. The same frozen encoder is reused without any fine-tuning for the OOD evaluation in Section 4.3.

## 4. Experiments

**Datasets.** We evaluate on two large-scale single-cell perturbation screens. The train and holdout splits are matched to scLDM (Palla et al., 2025). Parse 1M consists of cytokine perturbations across approximately one million cells, with a CD4 Naive holdout of 27 cytokines. The Replogle essential-genome CRISPRi screen (Replogle et al., 2022) provides a HepG2 holdout of 372 knock-outs.

**Metrics.** We evaluate at two complementary levels. *Distribution* metrics ( $W_2$ ,  $\text{MMD}^2$  with an RBF kernel,  $FD$  Fréchet Distance) following (Palla et al., 2025) compare

generated against true cells at the population level, and are sensitive to per-gene marginals and joint dispersion of the predicted cell cloud. The *cell-eval* per-perturbation metrics from (Adduri et al., 2025) fall into three functional groups by what they measure: (i) *DE-discovery* (*disc-11*, *PR-AUC*, *Spearman sig*) compare FDR-based DE-test calls against ground truth and depend on the predicted mean and variance across cells in a perturbation; (ii) *mean-ranking* (*Spearman LFC*, *overlap@N*) rank genes by predicted mean log2 fold change and depend only on the predicted mean; and (iii) *Pearson  $\Delta$* , which correlates the full per-gene  $\Delta$  vector and is sensitive to both. All numbers report mean  $\pm$  SE across four seeds.

### 4.1. Discrete latents outperform continuous baselines

Existing single-cell perturbation models build on continuous latent representations (scVI, CPA, scLDM) or operate directly in expression space (scGPT, STATE). Discrete-latent representation learning, despite being the foundation of the strongest generative models in language and vision, has not previously been applied to perturbation modeling. We close this gap. The `ExpressionVAE` encodes each cell as a fixed-length sequence of discrete codes drawn from a vector-quantized bottleneck, and a perturbation-conditioned prior, either autoregressive (AR) or masked discrete diffusion (MDLM), models the joint distribution over those codes. We compare against six baselines spanning continuous latent VAEs, conditional perturbation models, expression-space foundation models, and a non-parametric mean baseline: scLDM (Palla et al., 2025), scVI (Lopez et al., 2018a), CPA (Lotfollahi et al., 2023), scGPT (Cui et al., 2024), STATE (Adduri et al., 2025), and `PerturbMean`. Results on Replogle and Parse 1M are reported in Table 1.

For each prior we report two rows, corresponding to the NB-sample and MSE output heads. These sit at the sampling and deterministic ends of the head spectrum analyzed in Section 4.2 and trade off across metric families; the practical takeaway here is that our model achieves the strongest result in *each* family. Three observations stand out.

**Order-of-magnitude gains on distributional metrics.** The headline result is on  $W_2$ ,  $\text{MMD}^2$ , and  $FD$ , where the discrete-latent model improves over the strongest continuous-latent baseline (scLDM) by roughly an order of magnitude on  $\text{MMD}^2$  and  $FD$ , on both datasets. On Replogle,  $\text{MMD}^2$  drops from 0.200 to 0.010 and  $FD$  from 53.75 to 4.67; on Parse 1M,  $\text{MMD}^2$  drops from 0.027 to 0.008 and  $FD$  from 18.14 to 3.11. Expression-space methods (scGPT, STATE) and continuous latent baselines (scVI, CPA) lag the discrete-latent model by one to three orders of magnitude on these metrics.

Table 1. Discrete-latent perturbation modeling against six baselines on Replogle and Parse 1M. For each prior (AR, MDLM) we report two rows, corresponding to the NB-sample and MSE output heads (Section 4.2). Bold = best in column for that dataset

Dataset	Prior	$W_2 \downarrow$	$MMD^2 \downarrow$	$FD \downarrow$	$PDS \uparrow$	$PR-AUC \uparrow$	$Sp-sig \uparrow$	$Sp-LFC \uparrow$	$ovl@N \uparrow$	$P-\Delta \uparrow$
Replogle	AR	<b>7.969</b>	<b>0.010</b>	<b>4.67</b>	<b>0.803</b>	<b>0.298</b>	0.765	0.500	0.181	0.951
		8.989	0.098	34.40	0.751	0.233	<b>0.774</b>	0.584	<b>0.288</b>	<b>0.965</b>
	MDLM	8.105	0.013	5.20	0.779	0.292	0.747	0.508	0.188	0.950
		9.111	0.105	35.80	0.738	0.221	0.761	<b>0.585</b>	0.281	0.964
	STATE	20.580	0.730	366.64	0.641	0.210	0.721	0.357	0.212	0.420
	ScVI	17.359	0.453	284.47	0.517	0.190	0.598	0.257	0.085	0.350
	CPA	11.510	0.532	126.81	0.599	0.190	0.370	0.430	0.196	0.380
	ScGPT	34.166	3.087	1247.68	0.501	0.100	0.550	0.002	0.126	0.385
	ScLDM ( $\omega = 1$ )	11.292	0.200	53.75	0.730	–	–	0.413	0.163	–
	PerturbMean	–	–	–	0.616	0.190	0.190	0.246	0.162	–
Parse 1M	AR	<b>11.949</b>	<b>0.008</b>	<b>3.11</b>	0.920	0.493	<b>0.962</b>	0.793	0.400	<b>0.986</b>
		22.270	0.653	383.36	0.771	0.153	0.913	0.805	<b>0.480</b>	0.776
	MDLM	11.966	0.008	3.31	0.931	<b>0.495</b>	0.959	<b>0.812</b>	0.411	0.985
		22.314	0.660	384.68	0.770	0.150	0.928	0.786	0.463	0.776
	STATE	19.111	0.714	312.34	0.972	–	–	0.565	0.188	–
	ScVI	35.508	1.372	1233.11	0.519	–	–	0.010	0.044	–
	CPA	13.534	1.117	181.32	0.567	–	–	0.570	0.197	–
	ScGPT	22.870	2.203	523.93	0.520	–	–	0.057	0.077	–
	ScLDM ( $\omega = 1$ )	12.457	0.027	18.14	0.962	–	–	0.620	0.303	–
	PerturbMean	–	–	–	<b>0.973</b>	–	–	0.467	0.118	–

**Two priors, the same numbers.** AR and MDLM, drawn from different generative-model families and commonly viewed as alternatives, achieve effectively identical performance across all metrics on both datasets. This is direct evidence that the gain is from the discrete-latent representation rather than from a particular prior architecture, since the two priors share only the encoder and the codebook.

**Top of both metric families.** The model achieves the strongest result on every distributional and DE-discovery metric (sampling-rule row) and on every mean-ranking metric (deterministic-rule row), on both datasets. No baseline reaches the top of both families. Existing methods are competitive on distributional metrics or on ranking metrics, but never both: scLDM is the strongest baseline on distributional metrics yet trails on  $Sp-LFC$  and  $overlap@N$ , while PerturbMean and STATE are competitive on PDS and mean-ranking metrics but collapse on  $W_2$ ,  $MMD^2$ , and FD. The two-rows-per-prior structure that produces this pattern is not incidental; Section 4.2 shows it reflects a single underlying axis governed by inference-time stochasticity.

#### 4.2. Sampling-richness axis explains decoder-head behavior

A practical question for any single-cell perturbation model is which decoder head to use: cross-entropy, hurdle, MSE, negative binomial, or sampled negative binomial. The literature offers no consistent answer. Different papers report different metrics, and different metrics rank these heads

differently. The two rows per prior in Table 1 make the conflict visible within a single backbone: the same trained model produces noticeably different metric profiles depending only on the head. We show here that this is not noise. The standard evaluation metrics partition into three groups, and within each group the head ordering is monotonic in a single underlying axis: the richness of the inference-time sampling distribution.

**Five heads on a sampling-richness gradient.** We evaluate five heads on the same trained ExpressionVAE backbone. `ce-quant` parameterizes a categorical over  $[0, c_{\max}]$  per gene and samples from the softmax. `nb-sample` parameterizes a negative binomial  $NB(\mu, \theta)$  per gene and emits a draw from it. `hurdle` parameterizes a Bernoulli zero-gate combined with a deterministic positive magnitude  $\mu$ , sampling only the gate. `mse` emits a deterministic real-valued  $\mu$  everywhere. `nb-deter` shares trained parameters with `nb-sample` and differs only in that `predict()` returns  $\mu$  rather than a sample, which makes the prediction structurally anti-sparse since the softmax over genes leaks small mass to every gene. Read in the order `ce-quant`  $\rightarrow$  `nb-sample`  $\rightarrow$  `hurdle`  $\rightarrow$  `mse`  $\rightarrow$  `nb-deter`, the heads form a monotone gradient in the richness of the inference-time predictive distribution: full count sampling, parametric count sampling, sampled zero-gate over deterministic magnitude, fully deterministic, and fully deterministic plus structurally anti-sparse. Crucially, `nb-sample` and `nb-deter` share trained weights; the

only difference is the `predict()` rule.

**Three groups, three trends.** Figure 2 reports all five heads under both priors as per-metric heatmaps, with columns ordered to reflect the three-group partition. The structure is immediate from the color gradient. Group 1 (distribution) and Group 2 (DE-discovery) columns darken toward the sampling end of the head ordering, the top of each prior block; Group 3 (mean-ranking) columns darken toward the deterministic end, the bottom; *Pearson- $\Delta$*  flips winner between the two panels, with sampling heads winning on Parse 1M and deterministic heads winning on Replogle, exactly as a hybrid metric should behave when the dataset’s noise structure changes. The Group 1 collapse is most dramatic on Parse 1M, where moving the head from sampling to deterministic shifts FD by more than two orders of magnitude (e.g.,  $1.84 \rightarrow 454.1$  along the AR block); on Replogle, the same trend holds in direction but is compressed in magnitude, consistent with that dataset’s lower count dispersion.

The mechanism is straightforward. Group 1 metrics reward matching the full population cell cloud, including per-gene marginals and joint dispersion; Group 2 metrics reward DE-test calls that depend on both predicted mean and within-perturbation variance. Deterministic heads collapse the predicted cell cloud to a single point per condition, so within-perturbation variance is degenerate; DE tests applied to the predicted cells then operate on a vanishing  $\sigma$  and produce uninformative  $p$ -values, which propagates into PR-AUC, *disc\_H*, and *Sp-sig*. Group 3 metrics ask only about the per-gene predicted mean. Sampling injects between-cell noise that contaminates finite-sample mean estimates, so the noise that helps Groups 1 and 2 hurts Group 3.

**The NB pair isolates the sampling axis.** A natural confound is that *ce-quant* differs from *mse* on multiple axes simultaneously: discreteness, parametric form, output dimension, and training loss. To isolate sampling versus deterministic prediction from these confounds, we exploit that *nb-sample* and *nb-deter* share trained weights and differ only in `predict()`. Switching from  $\mu$  to a draw from  $\text{NB}(\mu, \theta)$  moves NB from the bottom tier to the top tier on the distribution and DE-discovery groups (e.g., on Parse,  $\text{FD} : 454.06 \rightarrow 3.11$  and  $\text{PR-AUC} : 0.129 \rightarrow 0.493$  under AR), and reverses the direction on the mean-ranking group (e.g., on Replogle,  $\text{overlap@N} : 0.281 \rightarrow 0.181$  under AR). Because no other component of the model changes, the entire effect is localized to the sampling step at inference. This rules out parametric-form, capacity, discrete-vs-continuous, and training-objective confounds, and pins the cross-group structure on the sampling axis.

### What metric measures what, and a recommendation.

Group 1 and Group 2 metrics are both sampling-axis metrics: they reward recovery of within-perturbation cell variance, expressed once at the cell-distribution level (Group 1) and once filtered through DE testing (Group 2). Group 3 metrics are mean-axis metrics that reward clean per-gene mean estimates and penalize sampling noise. *Pearson- $\Delta$*  is a hybrid whose winner tracks dataset noise structure. The practical implication is that head choice is governed by what the practitioner intends to measure. We recommend reporting at least one metric from each group, and selecting the head to match the downstream task: a sampling head (*ce-quant* or *sampled NB*) for use cases that consume predicted cell distributions (synthetic-cohort DE, virtual screens, power calculations), and a deterministic-mean head (*mse* or *NB-mean*) for gene-list ranking against a reference. Reporting any single metric in isolation is misleading.

### 4.3. Usefulness of latent space for phenotypic reversion

Sections 4.1 and 4.2 establish that a discrete-latent perturbation model with the right decoder head dominates within-distribution metrics, and that head choice is governed by an inference-time sampling axis with a clean metric-to-head map. To test whether these design choices reflect a generalizable inductive bias rather than over-fitting the controlled metric set, we evaluate the *frozen* encoder on an out-of-distribution biological discovery task. No information about the OOD dataset enters training; the encoder is used only to embed cells, and metrics are computed directly on the resulting latents.

**Benchmark.** The Pfizer phenotypic-reversion benchmark (Wong et al., 2026) is a CRISPRi screen of 1,732 perturbations across 864,115 cells under  $\text{IL-1}\beta$  and  $\text{TNF-}\alpha$  inflammation. The task is to rank perturbations capable of reverting the inflammatory state, with the canonical  $\text{NF-}\kappa\text{B}$  activators (*TNFRSF1A*, *TRADD*, *JUNB*, *JUND*, *NFKB1*, *NFKB2*) as positive controls. We report three metrics on the frozen latents: per-perturbation Calinski-Harabasz scores among treated cells (Gene CH (T)) and untreated cells (Gene CH (U)), measuring latent separability; and Enrichment AUC, the rank-AUC of the  $\text{NF-}\kappa\text{B}$  positive controls within the cosine-distance reversion ordering of the 1,732 perturbation centroids (random baseline 0.5).

**Configurations.** We evaluate the  $2 \times 2 \times 4$  grid of training datasets (Parse 1M, Replogle), encoder bottlenecks (FSQ, Gaussian), and decoder heads (Quantile, Hurdle, MSE, NB), shown in Figure 3. The four heads correspond to the four head families from Section 4.2; the encoder bottleneck is the additional axis that the controlled ablation did not vary. Although the decoder is used only during training (the OOD evaluation reads the encoder output directly), its choice in-

## Elucidating the Design Space of Generative Models for Single-Cell Perturbation Prediction

	W2 ↓	MMD <sup>2</sup> ↓	FD ↓	Disc-L1 ↑	PR-AUC ↑	Spearman sig ↑	Spearman L2FC ↑	DEG Overlap ↑	Pearson Δ ↑
AR ce-quant	7.868 ±0.010	0.009 ±0.0003	4.640 ±0.030	0.787 ±0.0040	0.283 ±0.0010	0.769 ±0.0050	0.534 ±0.012	0.171 ±0.0030	0.950
AR nb-sample	7.969 ±0.012	0.010 ±0.0003	4.670 ±0.030	0.803 ±0.0020	0.298 ±0.0010	0.765 ±0.0010	0.500 ±0.015	0.181 ±0.0050	0.951
AR hurdle	7.899 ±0.010	0.013 ±0.0004	5.580 ±0.110	0.780 ±0.0020	0.294 ±0.0010	0.781 ±0.0030	0.529 ±0.0060	0.221 ±0.0010	0.953
AR mse	8.989 ±0.017	0.098 ±0.0007	34.40 ±0.240	0.751 ±0.0030	0.233 ±0.0020	0.774 ±0.0040	0.584 ±0.0090	0.288 ±0.0020	0.965
AR nb (deter)	11.94 ±0.013	0.303 ±0.0014	100.2 ±0.380	0.778 ±0.0010	0.223 ±0.0010	0.759 ±0.0050	0.558 ±0.0090	0.281 ±0.0020	0.957
MDLM ce-quant	8.045 ±0.020	0.015 ±0.0004	5.820 ±0.150	0.751 ±0.0030	0.275 ±0.0020	0.753 ±0.0050	0.501 ±0.0040	0.163 ±0.0020	0.949
MDLM nb-sample	8.105 ±0.0060	0.013 ±0.0001	5.200 ±0.060	0.779	0.292 ±0.0020	0.747	0.508 ±0.0040	0.188 ±0.0010	0.950
MDLM hurdle	8.009 ±0.0030	0.017 ±0.0005	6.460 ±0.070	0.755 ±0.0010	0.286 ±0.0010	0.754 ±0.0040	0.522 ±0.0090	0.217 ±0.0010	0.952
MDLM mse	9.111 ±0.021	0.105 ±0.0011	35.80 ±0.380	0.738 ±0.0040	0.221 ±0.0010	0.761 ±0.0070	0.585 ±0.0070	0.281 ±0.0010	0.964
MDLM nb (deter)	12.09 ±0.010	0.316 ±0.0008	102.4 ±0.200	0.751 ±0.0050	0.216 ±0.0010	0.733 ±0.0030	0.562 ±0.0080	0.276 ±0.0020	0.956
	W2 ↓	MMD <sup>2</sup> ↓	FD ↓	Disc-L1 ↑	PR-AUC ↑	Spearman sig ↑	Spearman L2FC ↑	DEG Overlap ↑	Pearson Δ ↑
AR ce-quant	11.97 ±0.011	0.004 ±0.0001	1.840 ±0.020	0.978 ±0.0030	0.497 ±0.0020	0.952 ±0.0070	0.741 ±0.015	0.409 ±0.0050	0.939
AR nb-sample	11.95 ±0.0030	0.007	3.110	0.920 ±0.010	0.493 ±0.0010	0.962 ±0.0030	0.793 ±0.0080	0.400 ±0.0030	0.936
AR hurdle	12.04 ±0.0070	0.011 ±0.0003	4.900 ±0.100	0.907 ±0.0020	0.433 ±0.0050	0.936 ±0.0030	0.766 ±0.012	0.302 ±0.0040	0.933
AR mse	22.27 ±0.108	0.653 ±0.0047	383.4* ±4.54	0.771 ±0.0040	0.153 ±0.0020	0.913 ±0.0070	0.805 ±0.0030	0.480 ±0.0050	0.776
AR nb (deter)	23.09 ±0.038	0.855 ±0.0017	454.1* ±1.74	0.754 ±0.0050	0.129 ±0.0020	0.785 ±0.0060	0.755 ±0.0040	0.462 ±0.0020	0.808
MDLM ce-quant	11.99 ±0.0040	0.005 ±0.0001	2.120 ±0.040	0.978 ±0.0020	0.500 ±0.0020	0.955 ±0.0070	0.766 ±0.0070	0.420 ±0.0050	0.939
MDLM nb-sample	11.97 ±0.0060	0.003 ±0.0002	3.310 ±0.010	0.931 ±0.0030	0.495 ±0.0040	0.959 ±0.0070	0.812 ±0.0030	0.411 ±0.0050	0.935
MDLM hurdle	12.07 ±0.010	0.013 ±0.0002	5.480 ±0.080	0.904 ±0.0040	0.423 ±0.0050	0.933 ±0.0050	0.749 ±0.021	0.302 ±0.0040	0.932
MDLM mse	22.31 ±0.073	0.660 ±0.0031	384.7* ±2.81	0.770 ±0.0030	0.150 ±0.0030	0.928 ±0.0060	0.786 ±0.0080	0.463 ±0.0050	0.776
MDLM nb (deter)	23.07 ±0.047	0.859 ±0.0023	453.1* ±2.09	0.755 ±0.0050	0.117 ±0.0010	0.759 ±0.012	0.754 ±0.0060	0.453 ±0.0020	0.809

Figure 2. Output-head ablation on Replogle (top) and Parse 1M (bottom). Rows are five output heads under AR and MDLM priors on the same trained backbone; columns are grouped by metric family. Cells shaded per column, darker indicates better; mean over four seeds, SE in ± notation.

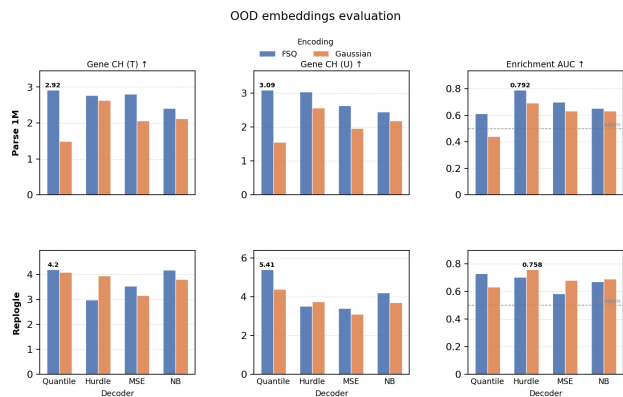


Figure 3. **Out-of-distribution evaluation on the Pfizer phenotypic-reversion benchmark.** Frozen ExpressionVAE encoders trained on Parse 1M (top) and Replogle (bottom), across encoder bottlenecks (FSQ, Gaussian) and decoder heads. Dashed line on the right shows the random baseline (0.5); annotated values mark the best configuration per panel.

fluences the frozen embedding by determining which signal the encoder must capture.

**The decoder transfers; the encoder pairing is dataset-dependent.** The Hurdle decoder identified by the controlled ablations is the strongest pairing on Enrichment AUC for both training datasets (0.792 on Parse 1M, 0.758 on Replogle); the decoder choice transfers cleanly without re-tuning. The optimal encoder pairing varies between datasets, however: FSQ + Hurdle wins on Parse 1M, while Gaussian + Hurdle wins on Replogle (FSQ + Hurdle reaches 0.704). Both axes contribute measurably to OOD performance, confirming that the decoder is not just a likelihood-modeling decision but shapes the representation itself.

**Discrete encoders dominate separability.** On the per-perturbation Calinski-Harabasz scores, FSQ + Quantile is the strongest configuration on both Gene CH (T) and Gene CH (U) on *both* datasets (2.92, 3.09 on Parse 1M; 4.20, 5.41 on Replogle), and FSQ encoders dominate Gaussian encoders across decoder pairings more generally. Hard quantization commits the encoder to a finite codebook of phenotypic states, producing sharper per-perturbation cluster structure in the latent space. This is consistent with the distributional-metric gap reported in Section 4.1: discrete latents allocate representational capacity in a way that preserves cell-state distinctions both within and outside the training distribution.

**Comparison with scGPT.** On the same benchmark, scGPT (Cui et al., 2024) achieves Enrichment AUC = 0.81. Our strongest Parse-trained configuration reaches 0.79, within typical evaluation noise of scGPT despite scGPT consuming roughly 10× more pre-training data. The point is not

a SOTA claim on the Pfizer benchmark, but a representation-efficiency claim: a discrete-latent perturbation prior trained on a single 1M-cell dataset effectively matches a foundation model trained on tens of millions of cells, on a biologically meaningful out-of-distribution ranking task. The explicit zero-gate decoder identified by our ablation transfers cleanly off the training distribution; hard quantization improves separability across both datasets and matches Gaussian on selectivity for Parse-trained encoders.

## 5. Conclusion

We introduced ExpressionVAE, the first discrete-latent perturbation model for single-cell data, and showed that pairing a vector-quantized bottleneck with either an autoregressive or a masked-diffusion prior achieves state-of-the-art across distributional and cell-eval metrics on Replogle and Parse 1M. Both priors achieve effectively identical performance, isolating the gain to the discrete latent rather than to a particular prior. A controlled output-head ablation further showed that decoder-head choice in this setting is governed by a single inference-time design axis, sampling versus deterministic prediction, with standard evaluation metrics partitioning into three groups along it; this resolves the apparent conflicts between published metric rankings into a clean three-way mapping between heads and downstream tasks. The hurdle decoder identified by the ablation transfers as the strongest selectivity pairing on a held-out CRISPRi reversion benchmark, where the frozen encoder effectively matches a 10× larger foundation model at a fraction of the pretraining data. Taken together, these results suggest that a compact discrete latent captures perturbation-axis structure efficiently, and that further progress in virtual cell modeling may come more from improvements to the latent representation than from scaling generic foundation models directly on raw counts.

## References

- Adamson, B., Norman, T. M., Jost, M., Cho, M. Y., Nuñez, J. K., Chen, Y., Villalta, J. E., Gilbert, L. A., Horlbeck, M. A., Hein, M. Y., et al. A multiplexed single-cell crispr screening platform enables systematic dissection of the unfolded protein response. *Cell*, 167(7):1867–1882, 2016.
- Adduri, A. K., Gautam, D., Bevilacqua, B., Imran, A., Shah, R., Naghipourfar, M., Teyssier, N., Ilango, R., Nagaraj, S., Dong, M., et al. Predicting cellular responses to perturbation across diverse contexts with state. *BioRxiv*, pp. 2025–06, 2025.
- Austin, J., Johnson, D. D., Ho, J., Tarlow, D., and Van Den Berg, R. Structured denoising diffusion models in

- 440 discrete state-spaces. *Advances in neural information*  
441 *processing systems*, 34:17981–17993, 2021.
- 442  
443 Bhattacharya, S., Gensbiger, C., Karim, S., and Lees, J. Dis-  
444 crete diffusion for single-cell gene expression modeling.  
445 *bioRxiv*, pp. 2026–02, 2026.
- 446  
447 Chang, H., Zhang, H., Jiang, L., Liu, C., and Freeman,  
448 W. T. MaskGIT: Masked generative image transformer. In  
449 *Proceedings of the IEEE/CVF Conference on Computer*  
450 *Vision and Pattern Recognition*, pp. 11315–11325, 2022.
- 451  
452 Cui, H., Wang, C., Maan, H., Pang, K., Luo, F., Duan, N.,  
453 and Wang, B. scgpt: toward building a foundation model  
454 for single-cell multi-omics using generative ai. *Nature*  
455 *methods*, 21(8):1470–1480, 2024.
- 456  
457 Dixit, A., Parnas, O., Li, B., Chen, J., Fulco, C. P., Jerby-  
458 Arnon, L., Marjanovic, N. D., Dionne, D., Burks, T., Ray-  
459 chowdhury, R., et al. Perturb-seq: dissecting molecular  
460 circuits with scalable single-cell rna profiling of pooled  
461 genetic screens. *cell*, 167(7):1853–1866, 2016.
- 462  
463 Dong, M., Adduri, A., Gautam, D., Carpenter, C., Shah,  
464 R., Ricci-Tam, C., Kluger, Y., Burke, D. P., and Roohani,  
465 Y. H. Stack: In-context learning of single-cell biology.  
466 *bioRxiv*, pp. 2026–01, 2026.
- 467  
468 Finak, G., McDavid, A., Yajima, M., Deng, J., Gersuk, V.,  
469 Shalek, A. K., Slichter, C. K., Miller, H. W., McElrath,  
470 M. J., Prlc, M., et al. Mast: a flexible statistical frame-  
471 work for assessing transcriptional changes and charac-  
472 terizing heterogeneity in single-cell rna sequencing data.  
473 *Genome biology*, 16(1):278, 2015.
- 474  
475 Fitch, W. T. Information and the single cell. *Current Opin-*  
476 *ion in Neurobiology*, 71:150–157, 2021.
- 477  
478 Gandhi, S., Javadi, F., Svensson, V., Khan, U., Jones, M. G.,  
479 Yu, J., Merico, D., Goodarzi, H., and Alidoust, N. Tahoe-  
480 x1: Scaling perturbation-trained single-cell foundation  
481 models to 3 billion parameters. *bioRxiv*, pp. 2025–10,  
482 2025.
- 483  
484 Gulati, G. S., Sikandar, S. S., Wesche, D. J., Manjunath,  
485 A., Bharadwaj, A., Berger, M. J., Ilagan, F., Kuo, A. H.,  
486 Hsieh, R. W., Cai, S., et al. Single-cell transcriptional di-  
487 versity is a hallmark of developmental potential. *Science*,  
488 367(6476):405–411, 2020.
- 489  
490 Hao, M., Gong, J., Zeng, X., Liu, C., Guo, Y., Cheng, X.,  
491 Wang, T., Ma, J., Zhang, X., and Song, L. Large-scale  
492 foundation model on single-cell transcriptomics. *Nature*  
493 *methods*, 21(8):1481–1491, 2024.
- 494  
495 Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya,  
496 E., Cai, T., Rutherford, E., de Las Casas, D., Hen-  
497 dricks, L. A., Welbl, J., Clark, A., Hennigan, T., Noland,  
498 E., Millican, K., van den Driessche, G., Damoc, B.,  
499 Guy, A., Osindero, S., Simonyan, K., Elsen, E., Rae,  
500 J. W., Vinyals, O., and Sifre, L. Training compute-  
501 optimal large language models. *ArXiv*, abs/2203.15556,  
502 2022. URL <https://api.semanticscholar.org/CorpusID:247778764>.
- 503  
504 Lin, Z., Akin, H., Rao, R., Hie, B., Zhu, Z., Lu, W.,  
505 Smetanin, N., Verkuil, R., Kabeli, O., Shmueli, Y., et al.  
506 Evolutionary-scale prediction of atomic-level protein  
507 structure with a language model. *Science*, 379(6637):  
508 1123–1130, 2023.
- 509  
510 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and  
511 Le, M. Flow matching for generative modeling. *arXiv*  
512 *preprint arXiv:2210.02747*, 2022.
- 513  
514 Liu, Y. et al. CellTok: Early-fusion multimodal large lan-  
515 guage model for single-cell transcriptomics via tokeniza-  
516 tion. *bioRxiv*, pp. 2025–10, 2025.
- 517  
518 Lopez, R., Regier, J., Cole, M., Jordan, M., and Yosef,  
519 N. Deep generative modeling for single-cell transcrip-  
520 tomics. *Nature Methods*, 15, 12 2018a. doi: 10.1038/  
521 s41592-018-0229-2.
- 522  
523 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I., and Yosef,  
524 N. Deep generative modeling for single-cell transcrip-  
525 tomics. *Nature methods*, 15(12):1053–1058, 2018b.
- 526  
527 Lotfollahi, M., Wolf, F. A., and Theis, F. J. scgen predicts  
528 single-cell perturbation responses. *Nature methods*, 16  
529 (8):715–721, 2019.
- 530  
531 Lotfollahi, M., Naghipourfar, M., Theis, F. J., and Wolf,  
532 F. A. Conditional out-of-distribution generation for  
533 unpaired data using transfer vae. *Bioinformatics*, 36  
534 (Supplement\_2):i610–i617, 2020.
- 535  
536 Lotfollahi, M., Klimovskaia Susmelj, A., De Donno, C.,  
537 Hetzel, L., Ji, Y., Ibarra, I. L., Srivatsan, S. R., Naghipour-  
538 far, M., Daza, R. M., Martin, B., et al. Predicting cellular  
539 responses to complex perturbations in high-throughput  
540 screens. *Molecular systems biology*, 19(6):e11517, 2023.
- 541  
542 Lou, A., Meng, C., and Ermon, S. Discrete diffusion  
543 modeling by estimating the ratios of the data distribu-  
544 tion. In *International Conference on Machine Learning*,  
545 2023. URL <https://api.semanticscholar.org/CorpusID:264451832>.
- 546  
547 Luo, E., Hao, M., Wei, L., and Zhang, X. scdiffu-  
548 sion: conditional generation of high-quality single-  
549 cell data using diffusion model. *Bioinformatics*, 40,  
550 2024. URL <https://api.semanticscholar.org/CorpusID:266844972>.

- 495 Macosko, E. Z., Basu, A., Satija, R., Nemesh, J., Shekhar,  
496 K., Goldman, M., Tirosh, I., Bialas, A. R., Kamitaki, N.,  
497 Martersteck, E. M., et al. Highly parallel genome-wide  
498 expression profiling of individual cells using nanoliter  
499 droplets. *Cell*, 161(5):1202–1214, 2015.
- 500 Mentzer, F., Minnen, D. C., Agustsson, E., and Tschan-  
501 nen, M. Finite scalar quantization: Vq-vae made  
502 simple. *ArXiv*, abs/2309.15505, 2023. URL [https://api.semanticscholar.org/CorpusID:  
503 //api.semanticscholar.org/CorpusID:  
504 263153393.](https://api.semanticscholar.org/CorpusID:263153393)
- 506 Oesinghaus, L., Becker, S., Vornholz, L., Papalexi, E., Pan-  
507 gallo, J., Moinfar, A. A., Liu, J., La Fleur, A., Shulman,  
508 M., Marrujo, S., et al. A single-cell cytokine dictionary  
509 of human peripheral blood. *bioRxiv*, 2025.
- 511 Ou, J., Nie, S., Xue, K., Zhu, F., Sun, J., Li, Z., and Li,  
512 C. Your absorbing discrete diffusion secretly models the  
513 conditional distributions of clean data. In *International  
514 Conference on Learning Representations (ICLR)*, 2025.
- 515 Palla, G., Babu, S., Dibaeinia, P., Pearce, J. D., Li,  
516 D., Khan, A. A., Karaletsos, T., and Tomczak, J. M.  
517 Scalable single-cell gene expression generation with  
518 latent diffusion models. *ArXiv*, abs/2511.02986,  
519 2025. URL [https://api.semanticscholar.  
520 org/CorpusID:282758604.](https://api.semanticscholar.org/CorpusID:282758604)
- 522 Palma, A., Richter, T., Zhang, H., Lubetzki, M., Tong, A.,  
523 Dittadi, A., and Theis, F. J. Multi-modal and multi-  
524 attribute generation of single cells with cfgen. In *In-  
525 ternational Conference on Learning Representations*,  
526 2024. URL [https://api.semanticscholar.  
527 org/CorpusID:271218151.](https://api.semanticscholar.org/CorpusID:271218151)
- 528 Peebles, W. and Xie, S. Scalable diffusion models with  
529 transformers. In *Proceedings of the IEEE/CVF interna-  
530 tional conference on computer vision*, pp. 4195–4205,  
531 2023.
- 533 Razavi, A., van den Oord, A., and Vinyals, O. Generating di-  
534 verse high-fidelity images with VQ-VAE-2. In *Advances  
535 in Neural Information Processing Systems*, volume 32,  
536 2019.
- 537 Replogle, J. M., Saunders, R. A., Pogson, A. N., Hussmann,  
538 J. A., Lenail, A., Guna, A., Mascibroda, L., Wagner,  
539 E. J., Adelman, K., Lithwick-Yanai, G., et al. Map-  
540 ping information-rich genotype-phenotype landscapes  
541 with genome-scale perturb-seq. *Cell*, 185(14):2559–2575,  
542 2022.
- 544 Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and  
545 Ommer, B. High-resolution image synthesis with latent  
546 diffusion models. In *Proceedings of the IEEE/CVF Con-  
547 ference on Computer Vision and Pattern Recognition*, pp.  
548 10684–10695, 2022.
- 549 Sahoo, S. S., Arriola, M., Schiff, Y., Gokaslan,  
A., Marroquin, E., Chiu, J. T., Rush, A., and  
Kuleshov, V. Simple and effective masked dif-  
fusion language models. *ArXiv*, abs/2406.07524,  
2024. URL [https://api.semanticscholar.  
org/CorpusID:270380319.](https://api.semanticscholar.org/CorpusID:270380319)
- Shi, J., Han, K., Wang, Z., Doucet, A., and Titsias, M.  
Simplified and generalized masked diffusion for discrete  
data. *Advances in neural information processing systems*,  
37:103131–103167, 2024.
- Theodoris, C. V., Xiao, L., Chopra, A., Chaffin, M. D.,  
Al Sayed, Z. R., Hill, M. C., Mantineo, H., Brydon, E. M.,  
Zeng, Z., Liu, X. S., et al. Transfer learning enables  
predictions in network biology. *Nature*, 618(7965):616–  
624, 2023.
- van den Oord, A., Vinyals, O., and Kavukcuoglu, K. Neural  
discrete representation learning. In *Advances in Neural  
Information Processing Systems*, volume 30, 2017.
- Wang, C., Karimzadeh, M., Ravindra, N. G., Bounds, L. R.,  
Alerasool, N., Huang, A. C., Ma, S., Gulbranson, D. R.,  
Cui, H., Lee, Y., et al. X-cell: Scaling causal perturbation  
prediction across diverse cellular contexts via diffusion  
language models. *bioRxiv*, pp. 2026–03, 2026a.
- Wang, M., Chen, C., Jiang, G., Ren, Z., Zhao, C., Shi, L.,  
and Ma, Y. Scdiva: Masked discrete diffusion for joint  
modeling of single-cell identity and expression. *arXiv  
preprint arXiv:2602.03477*, 2026b.
- Wong, D. R., Piper, M., Qiao, J., Russo, M., Jean, P., Clevert,  
D.-A., Arroyo, J. D., and Pashos, E. Phenotypic rever-  
sion and target prioritization for cellular inflammation via  
representation learning with foundation models. *bioRxiv*,  
pp. 2026–03, 2026.
- Zeng, T. and Dai, H. Single-cell rna sequencing-based  
computational analysis to describe disease heterogeneity.  
*Frontiers in Genetics*, 10:629, 2019.
- Zhang, H., Yuan, G.-H., Yuan, C., Xu, T., Bian, T., Cheng,  
H., Huang, W., Zhao, D., and Rong, Y. Lingshu-cell: A  
generative cellular world model for transcriptome model-  
ing toward virtual cells. *arXiv preprint arXiv:2603.25240*,  
2026.

## A. Hyperparameters and training configuration

Tables 2 and 3 list the hyperparameters used for the main-paper VAE backbone and prior models. The same configuration is shared across Parse 1M and Replogle, with two dataset-specific differences: (i) the perturbation conditioning source — ESM2-3B protein-language-model embeddings keyed by the knocked-out gene on Replogle, and a learnable per-cytokine table on Parse 1M (cytokine ligands have no direct protein-sequence analog) — and (ii) the train/holdout split definitions, matched to scLDM in both cases. Both datasets use additive cell-type conditioning on the adaLN-Zero path.

*Table 2. Stage 1: ExpressionVAE hyperparameters.* Shared across both datasets and all four output heads. The Gaussian-latent variant referenced in Section 4.3 swaps the FSQ bottleneck for a Gaussian-posterior bottleneck with the parameters in the lower block; all other settings are unchanged.

Latent sequence length $L$	64
FSQ codebook size $K$	512
Encoder / decoder	Cross-attention transformer
Gene input embedding	MLP over $\log_1 p$ count + learned gene-identity tokens
Optimiser	AdamW
Learning rate	$6 \times 10^{-4}$
Schedule	Cosine
Warm-up steps	2,500
Training iterations	75,000
Batch size	128
Gaussian-latent variant: per-token latent dim	8
Gaussian-latent variant: $\beta_{\text{KL}}$	$2 \times 10^{-6}$

*Table 3. Stage 2: prior hyperparameters.* Shared across the autoregressive (AR) and masked-discrete-diffusion (MDLM) priors. The DiT backbone is identical (same depth, width, head count, parameter budget) so any gap in Section 4.1 reflects the prior family rather than capacity.

Backbone	DiT (Diffusion Transformer)
Conditioning	adaLN-Zero on perturbation + cell-type embedding
Model dimension $d_{\text{model}}$	384
Attention heads	8
Layers	6
Feed-forward dimension	1,536
Optimiser	AdamW
Learning rate	$2 \times 10^{-4}$
Schedule	Cosine
Warm-up steps	3,000
Gradient clipping (max-norm)	1.0
Training iterations	75,000
Batch size	128
MDLM: masking schedule	absorbing, $\alpha_t = 1 - t$
MDLM: $t_\epsilon$	$5 \times 10^{-3}$
MDLM: inference denoising steps	128

## B. Codebook size and latent-token-count ablation

We swept the FSQ latent sequence length  $L \in \{64, 128\}$  against the codebook size  $K \in \{512, 1024\}$  on the Replogle (HepG2) holdout, paired with the AR-FSQ-MSE and MDLM-FSQ-MSE configurations from the main paper. We additionally report AR-FSQ-NB, MDLM-FSQ-NB, and a Flow-Gaussian-MSE configuration as neighbouring anchor points in the design space. Results are in Table 4.

The numbers in Table 4 were obtained on a different gene-set than the main-paper results and are therefore not directly

comparable to Table 1; they are, however, internally consistent across rows. Based on this ablation we carried forward  $L = 64$  and  $K = 512$  to the main experiments.

Table 4. Codebook size and latent-token-count ablation on the Replogle (HepG2) holdout. Bold = best per column.

Run	PDS $\uparrow$	ovl@N $\uparrow$	PR-AUC $\uparrow$	Sp-sig $\uparrow$	$W_2$ $\downarrow$	MMD <sup>2</sup> $\downarrow$	FD $\downarrow$
AR-FSQ-MSE $L=128, K=1024$	0.7373	0.2800	0.2225	0.7519	<b>10.02</b>	<b>0.0706</b>	56.26
AR-FSQ-MSE $L=128, K=512$	0.7398	<b>0.2914</b>	<b>0.2301</b>	0.7499	10.17	0.0734	57.28
AR-FSQ-MSE $L=64, K=1024$	0.7474	0.2773	0.2196	0.7563	10.06	0.0720	56.70
AR-FSQ-MSE $L=64, K=512$	0.7510	0.2834	0.2177	<b>0.7856</b>	10.06	0.0716	57.19
AR-FSQ-NB	<b>0.7594</b>	0.2790	0.2224	0.7678	12.35	0.1514	116.30
MDLM-FSQ-MSE $L=128, K=1024$	0.7030	0.2750	0.2071	0.7265	10.18	0.0791	59.31
MDLM-FSQ-MSE $L=128, K=512$	0.7147	0.2790	0.2135	0.7440	10.25	0.0776	57.75
MDLM-FSQ-MSE $L=64, K=1024$	0.7126	0.2692	0.2064	0.7495	10.19	0.0771	57.80
MDLM-FSQ-MSE $L=64, K=512$	0.7169	0.2816	0.2161	0.7344	10.27	0.0777	60.17
MDLM-FSQ-NB	0.7259	0.2728	0.2112	0.7406	12.49	0.1580	118.57
Flow-Gaussian-MSE	0.7514	0.2820	0.2204	0.7290	10.17	0.0707	<b>55.10</b>

### C. Additional results

Table 5. Parse 1M: full prior  $\times$  head ablation on Distribution metrics. We add results from running 4-seed pooled mean  $\pm$  SE runs.

Prior	Head	$W_2$ $\downarrow$	MMD <sup>2</sup> $\downarrow$	FD $\downarrow$
AR	ce-quantile	11.975 $\pm$ 0.011	<b>0.0041</b> $\pm$ 0.0001	<b>1.84</b> $\pm$ 0.02
	hurdle	12.039 $\pm$ 0.007	0.0114 $\pm$ 0.0003	4.90 $\pm$ 0.10
	mse	22.270 $\pm$ 0.108	0.6534 $\pm$ 0.0047	383.36 $\pm$ 4.54
	nb	23.088 $\pm$ 0.038	0.8550 $\pm$ 0.0017	454.06 $\pm$ 1.74
flow	ce-quantile	12.218 $\pm$ 0.008	0.0246 $\pm$ 0.0002	10.90 $\pm$ 0.15
	hurdle	11.918 $\pm$ 0.005	0.0074 $\pm$ 0.0001	2.95 $\pm$ 0.03
	mse	<b>11.890</b> $\pm$ 0.005	0.0108 $\pm$ 0.0003	3.98 $\pm$ 0.10
	nb	11.966 $\pm$ 0.010	0.0127 $\pm$ 0.0002	4.93 $\pm$ 0.07
MDLM	ce-quantile	11.992 $\pm$ 0.004	0.0049 $\pm$ 0.0001	2.12 $\pm$ 0.04
	hurdle	12.073 $\pm$ 0.010	0.0128 $\pm$ 0.0002	5.48 $\pm$ 0.08
	mse	22.314 $\pm$ 0.073	0.6601 $\pm$ 0.0031	384.68 $\pm$ 2.81
	nb	23.068 $\pm$ 0.047	0.8585 $\pm$ 0.0023	453.13 $\pm$ 2.09

Table 6. **Replogle (HepG2)**: full prior  $\times$  head ablation on Distribution metrics. We add results from running 4-seed pooled mean  $\pm$  SE runs.

Prior	Head	$W_2 \downarrow$	$MMD^2 \downarrow$	$FD \downarrow$
AR	ce-quantile	<b>7.868</b> $\pm 0.010$	<b>0.0087</b> $\pm 0.0003$	4.64 $\pm 0.08$
	hurdle	7.899 $\pm 0.010$	0.0129 $\pm 0.0004$	5.58 $\pm 0.11$
	mse	8.989 $\pm 0.017$	0.0978 $\pm 0.0007$	34.40 $\pm 0.24$
	nb	11.939 $\pm 0.013$	0.3029 $\pm 0.0014$	100.25 $\pm 0.38$
flow	ce-quantile	8.050 $\pm 0.008$	0.0091 $\pm 0.0001$	<b>4.53</b> $\pm 0.08$
	hurdle	7.908 $\pm 0.005$	0.0118 $\pm 0.0002$	5.17 $\pm 0.06$
	mse	8.973 $\pm 0.019$	0.0793 $\pm 0.0012$	27.45 $\pm 0.34$
	nb	10.884 $\pm 0.035$	0.2106 $\pm 0.0029$	68.15 $\pm 0.85$
MDLM	ce-quantile	8.045 $\pm 0.020$	0.0151 $\pm 0.0004$	5.82 $\pm 0.15$
	hurdle	8.009 $\pm 0.008$	0.0173 $\pm 0.0005$	6.46 $\pm 0.07$
	mse	9.111 $\pm 0.021$	0.1052 $\pm 0.0011$	35.80 $\pm 0.38$
	nb	12.086 $\pm 0.010$	0.3159 $\pm 0.0008$	102.37 $\pm 0.20$

 Table 7. **Parse 1M**: per-cell + per-gene moment fidelity. We add results from running 4-seed pooled mean  $\pm$  SE runs.

Prior	Head	$R^2$ mean $\uparrow$	$R^2$ var $\uparrow$	PCC per-cell $\uparrow$	MSE per-cell $\downarrow$
AR	ce-quantile	<b>0.999</b> $\pm 0.000$	<b>0.999</b> $\pm 0.000$	0.341 $\pm 0.000$	0.9972 $\pm 0.0005$
	hurdle	0.995 $\pm 0.000$	0.990 $\pm 0.000$	0.336 $\pm 0.000$	1.0320 $\pm 0.0001$
	mse	-0.308 $\pm 0.001$	0.472 $\pm 0.001$	<b>0.460</b> $\pm 0.000$	<b>0.9614</b> $\pm 0.0002$
	nb	-0.330 $\pm 0.000$	0.382 $\pm 0.000$	0.446 $\pm 0.000$	1.0674 $\pm 0.0003$
flow	ce-quantile	0.985 $\pm 0.000$	0.970 $\pm 0.000$	0.324 $\pm 0.000$	1.0476 $\pm 0.0002$
	hurdle	0.998 $\pm 0.000$	0.996 $\pm 0.000$	0.341 $\pm 0.000$	0.9991 $\pm 0.0004$
	mse	0.971 $\pm 0.000$	0.994 $\pm 0.000$	0.349 $\pm 0.000$	0.9810 $\pm 0.0003$
	nb	0.985 $\pm 0.000$	0.985 $\pm 0.000$	0.342 $\pm 0.000$	1.0001 $\pm 0.0004$
MDLM	ce-quantile	0.999 $\pm 0.000$	0.998 $\pm 0.000$	0.342 $\pm 0.000$	0.9971 $\pm 0.0005$
	hurdle	0.994 $\pm 0.000$	0.989 $\pm 0.000$	0.336 $\pm 0.000$	1.0341 $\pm 0.0004$
	mse	-0.315 $\pm 0.004$	0.469 $\pm 0.001$	<b>0.460</b> $\pm 0.000$	0.9630 $\pm 0.0008$
	nb	-0.325 $\pm 0.003$	0.381 $\pm 0.001$	0.446 $\pm 0.000$	1.0674 $\pm 0.0012$

 Table 8. **Replogle (HepG2)**: per-cell + per-gene moment fidelity. We add results from running 4-seed pooled mean  $\pm$  SE runs.

Prior	Head	$R^2$ mean $\uparrow$	$R^2$ var $\uparrow$	PCC per-cell $\uparrow$	MSE per-cell $\downarrow$
AR	ce-quantile	<b>1.000</b> $\pm 0.000$	<b>0.998</b> $\pm 0.000$	0.601 $\pm 0.000$	0.8901 $\pm 0.0009$
	hurdle	0.999 $\pm 0.000$	0.909 $\pm 0.003$	0.615 $\pm 0.000$	0.8502 $\pm 0.0012$
	mse	0.982 $\pm 0.000$	-2.201 $\pm 0.008$	0.738 $\pm 0.000$	0.5274 $\pm 0.0005$
	nb	0.934 $\pm 0.000$	-2.369 $\pm 0.005$	0.738 $\pm 0.000$	0.5528 $\pm 0.0002$
flow	ce-quantile	<b>1.000</b> $\pm 0.000$	0.996 $\pm 0.000$	0.603 $\pm 0.000$	0.8881 $\pm 0.0006$
	hurdle	0.999 $\pm 0.000$	0.933 $\pm 0.001$	0.615 $\pm 0.000$	0.8498 $\pm 0.0003$
	mse	0.984 $\pm 0.000$	-1.776 $\pm 0.003$	0.720 $\pm 0.000$	0.5615 $\pm 0.0005$
	nb	0.948 $\pm 0.001$	-1.511 $\pm 0.006$	0.697 $\pm 0.001$	0.6216 $\pm 0.0011$
MDLM	ce-quantile	<b>1.000</b> $\pm 0.000$	0.993 $\pm 0.001$	0.603 $\pm 0.000$	0.8840 $\pm 0.0010$
	hurdle	0.998 $\pm 0.000$	0.878 $\pm 0.008$	0.617 $\pm 0.001$	0.8458 $\pm 0.0023$
	mse	0.981 $\pm 0.000$	-2.243 $\pm 0.009$	0.739 $\pm 0.000$	<b>0.5248</b> $\pm 0.0002$
	nb	0.933 $\pm 0.000$	-2.428 $\pm 0.012$	<b>0.740</b> $\pm 0.000$	0.5501 $\pm 0.0006$

Table 9. All prior  $\times$  head ablations across cell-eval metrics

Dataset	Prior	Head	disc <sub>II</sub> $\uparrow$	ovl@N $\uparrow$	Sp-LFC $\uparrow$	PR-AUC $\uparrow$	MAE-I2 $\downarrow$	P- $\Delta$ $\uparrow$	MAE-pb $\downarrow$	Sp-sig $\uparrow$
AR		ce-quantile	<b>0.978</b> $\pm$ 0.003	0.409 $\pm$ 0.005	0.741 $\pm$ 0.015	0.497 $\pm$ 0.002	0.259 $\pm$ 0.001	0.989 $\pm$ 0.000	<b>0.023</b> $\pm$ 0.000	0.952 $\pm$ 0.007
		hurdle	0.907 $\pm$ 0.002	0.302 $\pm$ 0.004	0.766 $\pm$ 0.012	0.433 $\pm$ 0.005	0.300 $\pm$ 0.002	0.983 $\pm$ 0.000	0.030 $\pm$ 0.000	0.936 $\pm$ 0.003
		mse	0.771 $\pm$ 0.004	0.480 $\pm$ 0.005	0.805 $\pm$ 0.003	0.153 $\pm$ 0.002	0.282 $\pm$ 0.001	0.776 $\pm$ 0.000	0.324 $\pm$ 0.000	0.913 $\pm$ 0.007
		nb	0.754 $\pm$ 0.005	0.462 $\pm$ 0.002	0.755 $\pm$ 0.004	0.129 $\pm$ 0.002	0.299 $\pm$ 0.001	0.808 $\pm$ 0.000	0.338 $\pm$ 0.000	0.785 $\pm$ 0.006
Parse IM	flow	ce-quantile	0.681 $\pm$ 0.006	0.241 $\pm$ 0.005	0.719 $\pm$ 0.005	0.355 $\pm$ 0.005	0.355 $\pm$ 0.002	0.965 $\pm$ 0.000	0.040 $\pm$ 0.000	0.946 $\pm$ 0.004
		hurdle	<b>0.977</b> $\pm$ 0.001	0.459 $\pm$ 0.004	0.895 $\pm$ 0.003	<b>0.552</b> $\pm$ 0.004	0.196 $\pm$ 0.001	<b>0.991</b> $\pm$ 0.000	0.025 $\pm$ 0.000	<b>0.962</b> $\pm$ 0.004
		mse	0.975 $\pm$ 0.001	<b>0.513</b> $\pm$ 0.004	<b>0.903</b> $\pm$ 0.002	0.132 $\pm$ 0.002	<b>0.183</b> $\pm$ 0.001	0.976 $\pm$ 0.000	0.054 $\pm$ 0.000	0.850 $\pm$ 0.005
		nb	0.961 $\pm$ 0.002	0.477 $\pm$ 0.003	0.873 $\pm$ 0.005	0.155 $\pm$ 0.004	0.227 $\pm$ 0.002	0.978 $\pm$ 0.000	0.032 $\pm$ 0.000	0.847 $\pm$ 0.017
MDLM		ce-quantile	<b>0.978</b> $\pm$ 0.002	0.420 $\pm$ 0.005	0.766 $\pm$ 0.007	0.500 $\pm$ 0.002	0.254 $\pm$ 0.001	0.989 $\pm$ 0.000	<b>0.023</b> $\pm$ 0.000	<b>0.955</b> $\pm$ 0.007
		hurdle	0.904 $\pm$ 0.004	0.302 $\pm$ 0.004	0.749 $\pm$ 0.021	0.423 $\pm$ 0.005	0.306 $\pm$ 0.003	0.982 $\pm$ 0.000	0.031 $\pm$ 0.000	0.933 $\pm$ 0.005
		mse	0.770 $\pm$ 0.003	0.463 $\pm$ 0.005	0.786 $\pm$ 0.008	0.150 $\pm$ 0.003	0.289 $\pm$ 0.002	0.776 $\pm$ 0.000	0.325 $\pm$ 0.000	0.928 $\pm$ 0.006
		nb	0.755 $\pm$ 0.005	0.453 $\pm$ 0.002	0.754 $\pm$ 0.006	0.117 $\pm$ 0.001	0.304 $\pm$ 0.001	0.809 $\pm$ 0.000	0.338 $\pm$ 0.001	0.759 $\pm$ 0.012
AR		ce-quantile	0.787 $\pm$ 0.004	0.171 $\pm$ 0.003	0.534 $\pm$ 0.012	0.283 $\pm$ 0.001	0.496 $\pm$ 0.003	0.950 $\pm$ 0.000	0.109 $\pm$ 0.000	0.769 $\pm$ 0.005
		hurdle	0.780 $\pm$ 0.002	0.221 $\pm$ 0.001	0.529 $\pm$ 0.006	<b>0.294</b> $\pm$ 0.001	0.493 $\pm$ 0.002	0.953 $\pm$ 0.000	0.108 $\pm$ 0.000	<b>0.781</b> $\pm$ 0.003
		mse	0.751 $\pm$ 0.003	<b>0.288</b> $\pm$ 0.002	0.584 $\pm$ 0.009	0.233 $\pm$ 0.002	0.484 $\pm$ 0.002	<b>0.965</b> $\pm$ 0.000	0.129 $\pm$ 0.000	0.774 $\pm$ 0.004
		nb	0.778 $\pm$ 0.001	0.281 $\pm$ 0.002	0.558 $\pm$ 0.009	0.223 $\pm$ 0.001	0.493 $\pm$ 0.002	0.957 $\pm$ 0.000	0.207 $\pm$ 0.000	0.759 $\pm$ 0.005
Reprogle	flow	ce-quantile	0.784 $\pm$ 0.001	0.170 $\pm$ 0.005	0.498 $\pm$ 0.007	0.281 $\pm$ 0.002	0.497 $\pm$ 0.001	0.951 $\pm$ 0.000	0.109 $\pm$ 0.000	0.752 $\pm$ 0.007
		hurdle	0.786 $\pm$ 0.001	0.210 $\pm$ 0.001	0.519 $\pm$ 0.005	0.293 $\pm$ 0.001	0.489 $\pm$ 0.002	0.953 $\pm$ 0.000	<b>0.107</b> $\pm$ 0.000	0.756 $\pm$ 0.005
		mse	0.765 $\pm$ 0.003	0.278 $\pm$ 0.001	0.570 $\pm$ 0.007	0.221 $\pm$ 0.001	<b>0.481</b> $\pm$ 0.004	0.963 $\pm$ 0.000	0.125 $\pm$ 0.001	0.755 $\pm$ 0.001
		nb	<b>0.795</b> $\pm$ 0.001	0.268 $\pm$ 0.004	0.546 $\pm$ 0.009	0.206 $\pm$ 0.000	<b>0.482</b> $\pm$ 0.002	0.944 $\pm$ 0.000	0.189 $\pm$ 0.001	0.751 $\pm$ 0.008
MDLM		ce-quantile	0.751 $\pm$ 0.003	0.163 $\pm$ 0.002	0.501 $\pm$ 0.004	0.275 $\pm$ 0.002	0.504 $\pm$ 0.003	0.949 $\pm$ 0.000	0.110 $\pm$ 0.000	0.753 $\pm$ 0.005
		hurdle	0.755 $\pm$ 0.001	0.217 $\pm$ 0.001	0.522 $\pm$ 0.009	0.286 $\pm$ 0.001	0.496 $\pm$ 0.002	0.952 $\pm$ 0.000	0.110 $\pm$ 0.000	0.754 $\pm$ 0.004
		mse	0.738 $\pm$ 0.004	0.281 $\pm$ 0.001	<b>0.585</b> $\pm$ 0.007	0.221 $\pm$ 0.001	<b>0.480</b> $\pm$ 0.002	0.964 $\pm$ 0.000	0.129 $\pm$ 0.000	0.761 $\pm$ 0.007
		nb	0.751 $\pm$ 0.005	0.276 $\pm$ 0.002	0.562 $\pm$ 0.008	0.216 $\pm$ 0.001	0.494 $\pm$ 0.003	0.956 $\pm$ 0.000	0.209 $\pm$ 0.001	0.733 $\pm$ 0.003