

# CONCEPT BOTTLENECK DIFFUSION FOR STEERABLE GENERATION

Anonymous authors

Paper under double-blind review

## ABSTRACT

Recent work has brought the concept bottleneck architecture to the generative modeling space, adding a bottleneck inside of the typical model architecture to force interpretable reasoning. While this approach is flexible and can be applied to a variety of generative models, such architectural bottlenecks are better suited for single-step models like GANs than for diffusion models, which operate across a sequence of timesteps. This mismatch limits steerability and interpretability over the full generative process. In contrast, we place the bottleneck directly in the latent space used by a latent diffusion model, performing denoising along concept channels rather than in pixel space. By reshaping the latent space in which diffusion occurs, we learn a denoising process that is inherently aligned with the concepts we aim to control, enabled by our novel concept masking procedure. Our method, `CBDiffuse`, achieves improved steerability and control compared to prior work across CelebA-HQ and CUB.

## 1 INTRODUCTION

Generative models excel at producing realistic and high-quality images (Ho et al., 2020; Rombach et al., 2022; Podell et al., 2023). As these models improve, a key challenge is controlling generation to meet specific user requirements, enabling predictable interventions and more reliable outputs. Interpretability offers a natural mechanism for achieving this control, as it aligns internal model representations with human-understandable concepts, providing structured handles for intervention.

Concept Bottleneck Generative Models (CBGMs) (Ismail et al., 2023; Kulkarni et al., 2025b; Anonymous, 2026) present one avenue for leveraging interpretability to increase control in generative models. These works extend the classic Concept Bottleneck Model idea (Koh et al., 2020) to generative architectures such as GANs, VAEs, and diffusion models. By introducing intermediate concept representations, CBGMs provide a way to steer generation through interpretable, concept-level interventions.

However, these methods all operate within the model architecture itself, attempting to build a generative model bottleneck that is agnostic to the underlying generation process. While this design works well for GANs, which produce images in a single forward pass, it is less effective for diffusion models, where samples are produced through an iterative denoising trajectory. Concept bottlenecks in diffusion often influence generation only indirectly: the model still predicts pixel-space noise, so concept signals provide limited control over the output, acting as a conditioning signal rather than a structural constraint.

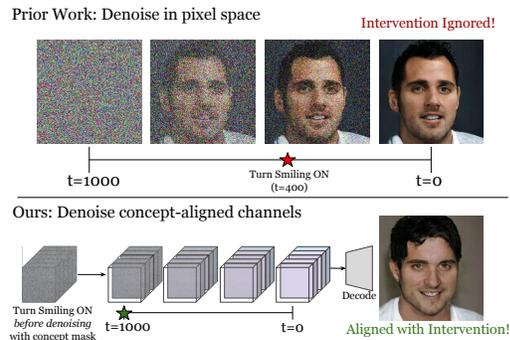


Figure 1: Prior work intervenes in pixel space pathway through denoising, making it harder for the U-Net to propagate concept information and sometimes leaving the final image unchanged. Our method masks concept-aligned latent channels before denoising, ensuring the trajectory follows the intervention and the resulting image faithfully reflects the intended concepts.

054  
 055  
 056  
 057  
 058  
 059  
 060  
 061  
 062  
 063  
 064  
 065  
 066  
 067  
 068  
 069  
 070  
 071  
 072  
 073  
 074  
 075  
 076  
 077  
 078  
 079  
 080  
 081  
 082  
 083  
 084  
 085  
 086  
 087  
 088  
 089  
 090  
 091  
 092  
 093  
 094  
 095  
 096  
 097  
 098  
 099  
 100  
 101  
 102  
 103  
 104  
 105  
 106  
 107

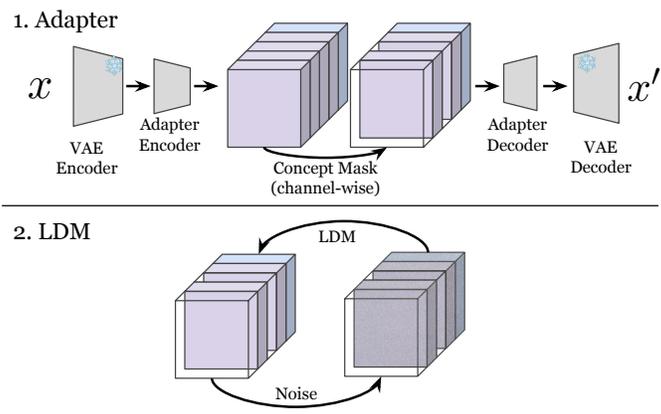


Figure 2: Overview of the CBDiffuse architecture. (1) A frozen VAE encoder produces a spatial latent that is mapped by a lightweight adapter into our concept-aligned latent space. Channels in this space are grouped into concept-specific blocks (in purple) and a residual block (in blue). Concept blocks encode discrete semantic states and are selectively masked based on the given concept values, while the residual block preserves non-concept information. (2) The latent diffusion model operates on the resulting masked representation, learning to denoise along active concept blocks and ignore masked channels.

We instead place the concept bottleneck directly in latent diffusion space and train a latent diffusion model to perform concept-aligned denoising in this interpretable space. To do so, we first train a lightweight adapter that maps from a pretrained, frozen VAE’s latent space into our concept-structured latent space, and then perform diffusion entirely within this space. Figure 1 illustrates this contrast: prior diffusion-based bottlenecks provide only partial control over generation, whereas our approach enables substantially improved concept-aligned steering by design.

To enforce controllability, we apply channel masking during training based on ground-truth concept values. Each concept is assigned a dedicated subset of channels, and only the channels corresponding to the active concept labels are unmasked. This encourages the model to encode each concept in its assigned channels and ensures that denoising progresses in a concept-aligned manner.

By structuring the latent space around semantic concepts and constraining denoising along concept-specific channels, our framework produces interpretable representations and enables predictable, controllable interventions. Our contributions:

- We propose a new framework for interpretable diffusion that introduces a concept bottleneck directly in the latent diffusion space, enabling concept-aligned denoising rather than pixel-space conditioning.
- We introduce a concept-aware channel masking mechanism that enforces semantic structure throughout the diffusion process, resulting in more reliable and consistent interventions.
- We demonstrate substantially improved steerability on CelebA-HQ and CUB compared to prior diffusion-based concept bottleneck methods.

## 2 CBDIFFUSE

The goal of interpretable generation is to enable control over the generative process using human-understandable concepts, producing images consistent with specified concept conditions. In this work we introduce CBDiffuse, an interpretable diffusion framework that performs denoising in a concept-aligned latent space before decoding back to pixel space (Figure 2). By organizing latent channels into semantically meaningful groups and enabling selective concept masking, our proposed architecture provides improved controllability compared to prior work.

In Section 2.1, we first describe the `CBDiffuse` architecture, which is built around concept channel masking. We then present our training algorithm (Section 2.2): including a minimal core version that relies on concept masking, followed by a full version that additionally incorporates techniques from prior work.

## 2.1 ARCHITECTURE

`CBDiffuse` operates by performing diffusion in an interpretable latent space rather than directly in pixel space, unlike the standard DDPM (Ho et al., 2020) architecture used in prior work (Ismail et al., 2023; Kulkarni et al., 2025b; Anonymous, 2026). To do this, we adopt a latent diffusion framework (Rombach et al., 2022), enabling the model to learn a denoising process that is explicitly aligned with human-understandable concepts.

As is typical for latent diffusion, we use a pretrained frozen VAE. We also introduce a lightweight adapter (Figure 2, Adapter). Our adapter is an autoencoder that maps VAE latents to our new interpretable latent space, with channels aligned to human-understandable concepts. Within this interpretable space, we train a latent diffusion model (LDM) to perform denoising directly on concept-aligned latents (Figure 2, LDM) instead of pixels.

**Notation & Latent Structure.** Let  $x \in \mathcal{X}$  denote an input image, and let  $c = \{c_1, \dots, c_K\}$  denote the  $K$  human-interpretable concepts associated with  $x$ . Each concept  $c_k$  takes values in a discrete domain (e.g., binary or categorical), and concept labels may be provided directly or obtained from external predictors such as CLIP.

First, the frozen VAE encoder  $E_{\text{vae}}$  maps images to a spatial latent representation  $z$ , and the VAE decoder  $D_{\text{vae}}$  reconstructs the image from  $z$ :  $z = E_{\text{vae}}(x)$ ,  $\hat{x} = D_{\text{vae}}(z)$ ,  $z \in \mathbb{R}^{H \times W \times C_{\text{vae}}}$  where  $H$  and  $W$  denote the height and width of the latent spatial map, respectively.

We then introduce a lightweight adapter to map the VAE latent space to an interpretable latent space. The adapter encoder  $E_{\text{adapter}}$  and decoder  $D_{\text{adapter}}$  transform the latent  $z$  into an interpretable latent  $s$  and back:  $s = E_{\text{adapter}}(z)$ ,  $\hat{z} = D_{\text{adapter}}(s)$ ,  $s \in \mathbb{R}^{H \times W \times C_{\text{adapter}}}$ .

For binary concepts, the interpretable latent  $s$  is partitioned into  $2K$  concept blocks and a residual block

$$s = [s^{1+}, s^{1-}, \dots, s^{K+}, s^{K-}, s^{(\text{res})}],$$

where each concept  $c_k$  is represented by two groups of  $C_{\text{concept}}$  channels:  $s^{k+}$  encodes the presence of the concept, and  $s^{k-}$  encodes its absence. The residual block  $s^{(\text{res})}$  with  $C_{\text{res}}$  channels preserves non-conceptual information. More generally, for a concept  $c_k$  with  $n_k$  discrete categories, we allocate one group of  $C_{\text{concept}}$  channels per category. The total number of adapter channels is therefore

$$C_{\text{adapter}} = \sum_{k=1}^K n_k C_{\text{concept}} + C_{\text{res}}.$$

**Concept Channel Masking.** To enforce concept alignment, we define a binary mask  $M(c) \in \{0, 1\}^{H \times W \times C_{\text{adapter}}}$  based on the concept assignment  $c$  for a given sample  $x$ .

For each concept, the mask activates only the channel group corresponding to the specified concept value, while all other concept groups are zeroed out. The masked latent is produced as follows  $\tilde{s} = M(c) \odot s$ , where  $\odot$  denotes element-wise multiplication. Residual channels are always unmasked. During training this encourages each concept block to represent only its intended semantic attribute.

Finally, we parameterize the latent diffusion model as a denoising network  $\epsilon_\theta$ , which operates in the interpretable latent space and predicts added noise at each diffusion timestep. The model takes as input a masked latent, and the diffusion timestep  $\epsilon_\theta : (\tilde{s}_t, t) \mapsto \hat{\epsilon}$ .

**Generation.** To generate new samples, we first sample a concept assignment  $c$  and gaussian noise  $s_T \sim \mathcal{N}(0, I)$  in the interpretable latent space. The concept mask  $M(c)$  then produces the masked latent  $\tilde{s}_T$ , which is iteratively denoised by the latent diffusion model  $\epsilon_\theta(\tilde{s}_t, t)$  to yield the concept-aligned latent  $\hat{s}_0$ . The final image is decoded via  $\hat{x} = D_{\text{vae}}(D_{\text{adapter}}(\hat{s}_0))$ .

Table 1: Steerability comparisons for diffusion-based methods on CelebA-HQ (8 concepts). Results for CBDiffuse are computed on 5k samples. Results for CBDiffuse are reported on 256×256 images. CBGM and CB-AE results are taken from their reported results since they do not share publicly available diffusion code, and were evaluated at 64 × 64 instead of 256 × 256 for CUB.

Steerability (%)	CelebA-HQ	CUB
CBGM	–	14.8
CB-AE	23.5	25.6
<b>CBDiffuse (Core)</b>	<b>67.6</b>	<b>41.5</b>
CB-AE + OptInt	50.5	36.9
CC + OptInt	56.7	44.5
<b>CBDiffuse (Full) + OptInt</b>	<b>82.0</b>	<b>44.6</b>

**Intervention.** To perform concept intervention, we replace the original concept assignment  $c_1$  with a new assignment  $c_2$  while keeping the same noise latent  $s_T$ . We construct the corresponding mask and obtain the masked latent  $\tilde{s}'_T = s_T \odot M(c_2)$ , which is denoised by the latent diffusion model  $\epsilon_\theta(\tilde{s}'_t, t)$  to yield  $\hat{s}'_0$ . The intervened image is then decoded as  $\hat{x}' = D_{\text{vae}}(D_{\text{adapter}}(\hat{s}'_0))$ .

## 2.2 TRAINING

Training proceeds in two stages. First, we pretrain the adapter to shape the latent space into a concept-aligned representation, ensuring that semantic concepts are disentangled prior to diffusion training. Second, we jointly train the adapter and the latent diffusion model (LDM) to guide the denoising process along these concept-aligned channels.

We first present our *core* algorithm, which relies solely on concept masking and an orthogonality loss. We then describe a *full* variant that incorporates additional auxiliary losses inspired by prior work, including concept classification losses, intervention losses, and test-time optimization-based interventions. We briefly describe each below, see Appendix C for full details.

**CBDiffuse (Core).** We first train the core variant of our model to isolate the effect of the proposed architecture on downstream steerability. We train CBDiffuse (Core) with the following objective:

$$\mathcal{L}_{\text{core}} = \mathcal{L}_{\text{recons}} + \mathcal{L}_{\text{LDM}} + \lambda_{\text{orth}}\mathcal{L}_{\text{orth}}$$

**CBDiffuse (Full).** While our core architecture is already effective at learning steerable concept representations, it can be further enhanced with auxiliary losses introduced in prior work. Importantly, our main contribution lies in the architectural design and masking procedure of CBDiffuse, which remains orthogonal to these additions. For CBDiffuse (Full), we train with the following objective:

$$\mathcal{L}_{\text{full}} = \mathcal{L}_{\text{core}} + \lambda_{\text{concept}}\mathcal{L}_{\text{concept}} + \lambda_{\text{intervention}}\mathcal{L}_{\text{intervention}}.$$

## 3 RESULTS

We first describe general details about our experimental setup in Section 3.1, and then aim to answer the following questions: **Q1:** How controllable is CBDiffuse compared to prior work on interpretable generative modeling (Section 3.2)? **Q2:** How do CBDiffuse (Full) losses impact the downstream steerability and image quality (Section 3.3)?

### 3.1 EXPERIMENTAL SETUP

**Architecture Details.** We use a pretrained E2E-VAAE-HF model from Leng et al. (2025) with a latent space of size  $16 \times 16 \times 32$ . The adapter consists of four residual convolutional blocks in both the encoder and decoder. As the latent diffusion backbone, we use SiT-B/1 (Ma et al., 2024). For concept classification, we train lightweight linear heads for each concept. Additional architectural and training details are provided in Appendix A.

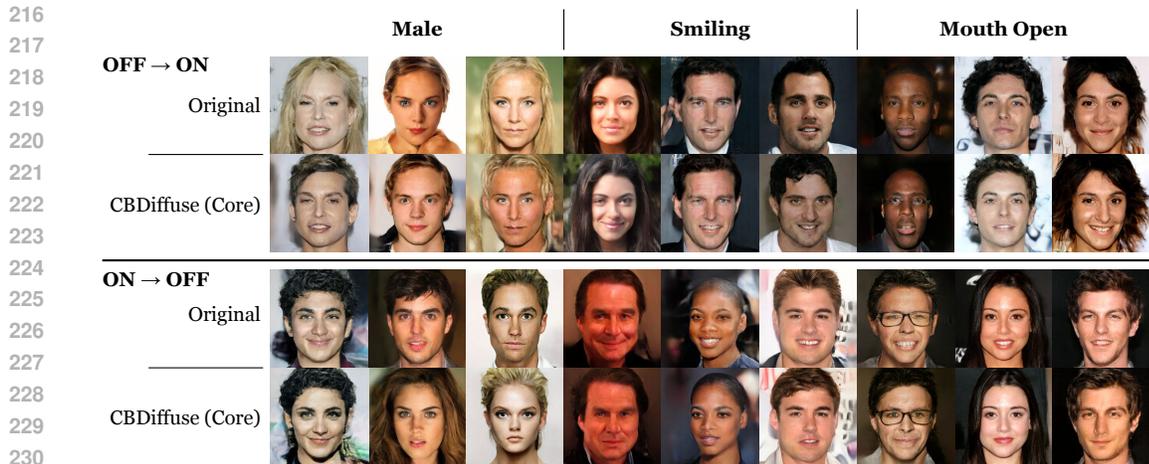


Figure 3: Concept intervention results for Male, Smiling, and Mouth Open. The top half shows images where the concept is flipped from OFF to ON, while the bottom half shows flips from ON to OFF. For each concept, the first row displays the original images and the second row shows the corresponding images generated by CBDiffuse (Core) after intervention.

**Datasets.** We evaluate our approach on CelebA-HQ (Lee et al., 2020) using eight binary concepts (Smiling, Male, Heavy Makeup, Mouth Slightly Open, Attractive, Wearing Lipstick, High Cheekbones, Arched Eyebrows), and on CUB (Wah et al., 2011) using ten visual attributes (Small size (5–9 inches), Perching-like shape, Solid breast pattern, Black bill color, Bill length shorter than head, Black wing color, Solid belly pattern, All-purpose bill shape, Black upperparts color, White underparts color). For both datasets, all images are resized to  $256 \times 256$  to match the resolution of the pretrained VAE.

**Training.** For both models, we first pretrain the adapter for 40K steps, followed by joint training of the full model for 400K steps. For CBDiffuse (Full), additional loss terms are incorporated during joint training. Further training details and loss weights are provided in Appendix A.

**Concept Evaluation.** We automate concept evaluation using ViT-L-16 binary classifiers trained to detect each concept in an image. For CelebA-HQ, we use the same classifiers as Kulkarni et al. (2025b) to ensure fair comparison. Since concept classifiers for CUB are not publicly available, we train our own using the original authors’ code.

**Baselines.** We compare against two primary baselines: CBGM Ismail et al. (2023) and CB-AE Kulkarni et al. (2025b). Since neither provides code for replicating their diffusion model results, we report their published numbers directly.

### 3.2 STEERABILITY RESULTS

To evaluate how controllable CBDiffuse is compared to prior work, we measure the model’s *steerability*. Following prior work, we generate samples in both directions: first with a concept OFF and then flipping it ON, and second with a concept ON and then flipping it OFF. For each flip, an external concept classifier determines whether the generated image reflects the intended change. Steerability is reported as the percentage of samples where the concept successfully flips, averaged over all concepts and directions.

Results are shown in Table 1. The core model (CBDiffuse Core) improves over CB-AE by 44.1% on CelebA-HQ and 15.9% on CUB, demonstrating that the masked reconstruction and denoising losses alone provide strong control. The full model (CBDiffuse Full) further increases steerability by 14.4% on CelebA-HQ and 3.1% on CUB, showing that these gains are complementary to prior approaches. Overall, these results emphasize that learning an interpretable latent space is crucial for precise control in diffusion models.

Table 2: Stepwise ablation of components on CelebA-HQ. Concept classifiers are required to perform interventions or apply OptInt. Steerability steadily improves as components are added. The intervention loss slightly increases FID, reflecting a trade-off between image quality and controllability. For reference, a standard LDM with the same hyperparameters achieves an FID of 24.9.

$L_{concept}$	$L_{intervene}$	OptInt	Conc. Acc. (%)	Steerability (%)	FID ( $\downarrow$ )
×	×	×	-	67.6	29.0
✓	×	×	74.5	66.3	<b>24.6</b>
✓	×	✓	74.5	74.6	24.6
✓	✓	×	<b>80.9</b>	78.4	33.9
✓	✓	✓	<b>80.9</b>	<b>82.0</b>	33.9

Figure 3 illustrate the visual effects of our approach, presenting interventions on the concepts *Male*, *Smiling*, and *Mouth Open* using CBDiffuse (Core), demonstrating consistent and precise control over multiple attributes.

### 3.3 CBDIFFUSE LOSSES

In this section, we perform a stepwise ablation of the losses that differentiate CBDiffuse (Full) from CBDiffuse (Core): the concept loss ( $L_{concept}$ ), the intervention loss ( $L_{intervene}$ ), and optimized interventions (OptInt). We evaluate these models on three metrics: concept accuracy, steerability, and generation quality. Concept accuracy is measured by randomly sampling a concept assignment, generating an image conditioned on the corresponding mask, and assessing whether the generated image exhibits the intended concept; results are averaged across all concepts. Generation quality is evaluated using the Fréchet Inception Distance (FID) (Heusel et al., 2017), computed over 10K images generated by sampling concept masks according to their empirical distribution in the training data.

Our results are shown in Table 2. For reference, the baseline SiT model trained with the same pretrained VAE achieves an FID of 24.9. Incorporating the concept loss ( $L_{concept}$ ) slightly improves FID compared to other variants, likely due to the additional supervision it provides, closely matching the FID of the baseline latent diffusion model. While the concept loss alone does not improve steerability, it enables the use of optimized interventions (OptInt), which further increase steerability beyond the core model. In contrast, adding the intervention loss ( $L_{intervene}$ ) slightly degrades FID, likely because it emphasizes learning a concept-aligned latent space at the expense of image quality. However, it substantially boosts steerability compared to using only the concept loss, by +12.1% without OptInt and +7.4% with OptInt.

Optimized interventions (OptInt) consistently improve steerability, yielding an increase of 8.3% when applied without the intervention loss and 3.6% when used with the full model. The smaller absolute gain for the full model is expected, as higher baseline steerability leaves less room for improvement. This trend aligns with observations from CB-AE (Kulkarni et al., 2025b), where OptInt improved steerability by 27%, starting from a relatively low baseline of 23.5%.

## 4 CONCLUSIONS

We introduced CBDiffuse, an interpretable diffusion framework that places the bottleneck directly in the latent space instead of in the diffusion model architecture. By training an interpretable adapter and enforcing concept-specific channel masking throughout the denoising process, our approach aligns semantic structure with diffusion dynamics. This design enables concept-aligned denoising and yields substantially improved steerability compared to prior concept bottleneck generative models.

Empirically, we demonstrated strong gains in steerability on CelebA-HQ and CUB, with ablations showing that channel masking alone provides effective control while remaining complementary to existing intervention losses and optimization-based techniques. While this approach may trade off some fidelity in exact pixel-level edits, it allows for more reliable and precise interventions at the semantic level, highlighting the benefits of embedding interpretable structure directly into the generative process.

- 
- 324 REFERENCES  
325  
326 Anonymous. A probabilistic hard concept bottleneck for steerable generative models. 2026. Accepted  
327 submission.
- 328 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
329 *in neural information processing systems*, 34:8780–8794, 2021.  
330
- 331 Mateo Espinosa Zarlenga, Pietro Barbiero, Gabriele Ciravegna, Giuseppe Marra, Francesco Giannini,  
332 Michelangelo Diligenti, Zohreh Shams, Frederic Precioso, Stefano Melacci, Adrian Weller, et al.  
333 Concept embedding models: Beyond the accuracy-explainability trade-off. *Advances in neural*  
334 *information processing systems*, 35:21400–21413, 2022.
- 335 Rohit Gandikota, Joanna Materzynska, Jaden Fiotto-Kaufman, and David Bau. Erasing concepts  
336 from diffusion models. In *Proceedings of the IEEE/CVF international conference on computer*  
337 *vision*, pp. 2426–2436, 2023.  
338
- 339 Rohit Gandikota, Hadas Orgad, Yonatan Belinkov, Joanna Materzyńska, and David Bau. Unified  
340 concept editing in diffusion models. In *Proceedings of the IEEE/CVF Winter Conference on*  
341 *Applications of Computer Vision*, pp. 5111–5120, 2024.
- 342 Ada Gorgun, Fawaz Sammani, Nikos Deligiannis, Bernt Schiele, and Jonas Fischer. Temporal  
343 concept dynamics in diffusion models via prompt-conditioned interventions. *arXiv preprint*  
344 *arXiv:2512.08486*, 2025.  
345
- 346 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans  
347 trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural*  
348 *information processing systems*, 30, 2017.
- 349 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,  
350 2022.  
351
- 352 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
353 *neural information processing systems*, 33:6840–6851, 2020.  
354
- 355 Aya Abdelsalam Ismail, Julius Adebayo, Hector Corrada Bravo, Stephen Ra, and Kyunghyun Cho.  
356 Concept bottleneck generative models. In *The Twelfth International Conference on Learning*  
357 *Representations*, 2023.
- 358 Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and  
359 Percy Liang. Concept bottleneck models. In *International conference on machine learning*, pp.  
360 5338–5348. PMLR, 2020.  
361
- 362 Akshay Kulkarni, Tsui-Wei Weng, Vivek Narayanaswamy, Shusen Liu, Wesam A Sakla, and Kowshik  
363 Thopalli. Interpretable and steerable concept bottleneck sparse autoencoders. *arXiv preprint*  
364 *arXiv:2512.10805*, 2025a.
- 365 Akshay Kulkarni, Ge Yan, Chung-En Sun, Tuomas Oikarinen, and Tsui-Wei Weng. Interpretable  
366 generative models through post-hoc concept bottlenecks. In *Proceedings of the Computer Vision*  
367 *and Pattern Recognition Conference*, pp. 8162–8171, 2025b.  
368
- 369 Cheng-Han Lee, Ziwei Liu, Lingyun Wu, and Ping Luo. Maskgan: Towards diverse and interactive  
370 facial image manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and*  
371 *pattern recognition*, pp. 5549–5558, 2020.
- 372 Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng.  
373 Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint*  
374 *arXiv:2504.10483*, 2025.  
375
- 376 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and  
377 Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant  
transformers. In *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.

---

378 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
379 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image  
380 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.

381  
382 Kanchana Ranasinghe, Muzammal Naseer, Munawar Hayat, Salman Khan, and Fahad Shahbaz Khan.  
383 Orthogonal projection loss. In *Proceedings of the IEEE/CVF international conference on computer  
384 vision*, pp. 12333–12343, 2021.

385 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
386 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-  
387 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.

388  
389 Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd  
390 birds-200-2011 dataset. 2011.

391 Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training.  
392 *arXiv preprint arXiv:2001.03994*, 2020.

393  
394 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image  
395 diffusion models. In *Proceedings of the IEEE/CVF international conference on computer vision*,  
396 pp. 3836–3847, 2023.

397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

---

## 432 A EXPERIMENTAL DETAILS

### 433 A.1 ARCHITECTURE

434 **VAE.** We use a pretrained E2E-VAAE-HF that maps input images of size  $256 \times 256 \times 3$  to a  
437 spatial latent representation

$$438 z \in \mathbb{R}^{H \times W \times C_{\text{vae}}}, \quad H = W = 16, \quad C_{\text{vae}} = 32,$$

439 and reconstructs images as  $\hat{x} = D_{\text{vae}}(z)$ .

441 **Latent Adapter.** The adapter maps the VAE latent space to an interpretable, concept-aligned space:

$$442 f : \mathbb{R}^{16 \times 16 \times 32} \rightarrow \mathbb{R}^{16 \times 16 \times C_{\text{concept}}},$$

443 where

444  $C_{\text{concept}} = (\text{channels per concept}) \times (\text{num concepts}) \times (\text{categories per concept}) + (\text{residual channels})$ .

445 Table 3: Latent space and concept channel specifications for CelebA-HQ and CUB.

450 Dataset	VAE Latent ( $H \times W \times C_{\text{vae}}$ )	Adapter Latent ( $H \times W \times C_{\text{concept}}$ )	Concept Channels
451 CelebA-HQ	$16 \times 16 \times 32$	$16 \times 16 \times 72$	8 concepts, 4 channels per concept, 8 residual channels
452 CUB	$16 \times 16 \times 32$	$16 \times 16 \times 88$	10 concepts, 4 channels per concept, 8 residual channels

453 The adapter consists of an encoder and decoder with 4 residual convolution blocks each, 256 hidden  
454 channels,  $3 \times 3$  convolutions, GroupNorm, and ReLU activations, totaling 9.76M parameters. The  
455 CBDiffuse latent diffusion model uses a SiT-B/1 backbone with 138M parameters.

### 456 A.2 TRAINING AND LOSSES

457 Adapter pretraining is performed for 40K steps, followed by joint training with the diffusion model  
458 for 400K steps. Batch size is 64. AdamW optimizers with learning rates of  $1 \times 10^{-4}$  are used for  
459 both adapter and diffusion model.

- 460 • **CBDiffuse (Core):** masked reconstruction and denoising losses.
- 461 • **CBDiffuse (Full):** concept loss ( $\lambda_{\text{concept}} = 1$ ), intervention loss ( $\lambda_{\text{intervene}} = 0.1$ ), or-  
462 thogonality regularization ( $\lambda_{\text{orthog}} = 0.1$ ), optionally combined with optimization-based  
463 intervention (OptInt).

464 The cyclical intervention loss is applied every 25 iterations to reduce computational cost. OptInt  
465 performs 50 optimization steps with step size  $\epsilon = 0.1$  during inference to refine latent concept  
466 activations.

### 467 A.3 EVALUATION AND HARDWARE

468 For each concept assignment, 500 samples are generated. Concept accuracy is measured using  
469 external ViT-L-16 classifiers. Steerability is computed by flipping concepts ON/OFF and verifying  
470 that generated images reflect the intended change. FID is computed over 10K images sampled  
471 according to the empirical concept distribution.

472 All experiments were conducted on a node with 112 CPU cores, 1TB RAM, and  $2 \times$  NVIDIA RTX  
473 6000 Ada GPUs.

## 474 B RELATED WORK

### 475 B.1 CONCEPT BOTTLENECK GENERATIVE MODELS (CBGMs)

476 Our work builds on recent advances in interpretable generative modeling using the concept bottleneck  
477 architecture (Koh et al., 2020). Ismail et al. (2023) first introduced a concept embedding layer

---

486 (Espinosa Zarlenga et al., 2022) into generative models, enabling interventions to guide image  
487 generation.

488  
489 Kulkarni et al. (2025b) extended this by proposing a post-hoc CBGM, which learns a mapping from  
490 a frozen pre-trained generative model to interpretable concept representations. They also introduced  
491 a concept controller strategy that allows for steerability without modifying the underlying latent.  
492 Kulkarni et al. (2025a) further explored combining concept bottlenecks with sparse autoencoders  
493 to balance supervised and unsupervised concept learning in large vision-language models. Most  
494 recently, concurrent work (Anonymous, 2026) added “hard” concept constraints to reduce concept  
495 leakage and improve steerability.

496 **Our Approach.** Prior CBGM work focuses on inserting bottlenecks within the generative model  
497 architecture itself. While effective for GANs or single-step models, this design provides limited  
498 interpretability and control in diffusion models. Interventions applied mid-denoising in pixel space are  
499 difficult for the U-Net to propagate, often leaving the final image largely unchanged. `CBDiffuse`  
500 addresses this by mapping VAE latents to a concept-aligned latent space and performing denoising  
501 directly along concept channels. By structuring the diffusion process around these channels, inter-  
502 ventions are faithfully propagated throughout the denoising trajectory, enabling more predictable  
503 concept-level control.

## 504 B.2 SEMANTIC CONTROL IN DIFFUSION MODELS

505  
506 Controlling semantic information in diffusion models has become a major area of research. Early  
507 approaches include *classifier guidance* (Dhariwal & Nichol, 2021) and *classifier-free guidance* (Ho  
508 & Salimans, 2022), which steer the generation process by adjusting the conditioning signal during  
509 denoising. While effective for coarse control, these methods operate indirectly and do not enforce  
510 a structured latent representation for specific concepts, limiting the precision and predictability of  
511 interventions. Text-conditioned diffusion models and ControlNet (Rombach et al., 2022; Zhang  
512 et al., 2023) extend this idea by learning mappings from text prompts or structured conditions to  
513 image outputs. These approaches achieve impressive results at scale, particularly for high-level  
514 attributes like scene composition or object presence. However, control remains largely post-hoc: the  
515 model still denoises in pixel or latent space without an explicit semantic structure, and concepts can  
516 become locked at specific timesteps along the denoising trajectory (Gorgun et al., 2025). Other work  
517 explores mechanistic interpretability techniques for more targeted semantic control. Gandikota et al.  
518 (2023) propose erasing specific visual concepts directly from a pretrained diffusion model, effectively  
519 removing objects, styles, or attributes from the model itself rather than merely filtering outputs at  
520 inference. Building on this idea, Gandikota et al. (2024) introduce Unified Concept Editing (UCE),  
521 which enables simultaneous editing of multiple concepts, such as debiasing or content manipulation,  
522 by intervening in model parameters in a closed-form, scalable manner.

523 **Our Approach.** In contrast, `CBDiffuse` focuses on training semantic control into the diffusion  
524 model from the start, rather than adapting an already pretrained model. By mapping a pretrained  
525 VAE’s latent space into concept-aligned channels and applying concept-aware masking during  
526 training, we guide the denoising trajectory along desired concepts. This enables concept-aligned  
527 generation and precise interventions at smaller scale, making it more accessible for specialized  
528 domains or experiments where large pretrained models are impractical.

529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

---

## 540 C LOSS DEFINITIONS

541  
542 **Masked Reconstruction Loss ( $\mathcal{L}_{\text{recons}}$ ).** To encourage the adapter to learn an interpretable latent  
543 space, we apply a masked reconstruction loss. By reconstructing latents from only the channels  
544 corresponding to active concepts in each image, the adapter is incentivized to isolate information  
545 about each concept into its designated channels.

546 Formally, given an input image  $x \sim \mathcal{X}$  with concept assignment  $c$ , we first encode  $x$  using the frozen  
547 VAE:  $z = E_{\text{vae}}(x)$ . We then encode the latent with the adapter encoder  $s = E_{\text{adapter}}(z)$ , apply the  
548 concept-specific mask  $M(c)$  to the resulting latent, and pass the masked latent through the adapter  
549 decoder to reconstruct the original latent. The loss is the mean squared error (MSE) between the  
550 original and reconstructed latents:

$$551 \mathcal{L}_{\text{recons}} = \mathbb{E}_{x \sim \mathcal{X}} \left[ \|D_{\text{adapter}}(s \odot M(c)) - z\|_2^2 \right].$$

552  
553  
554 **Masked Denoising Loss ( $\mathcal{L}_{\text{LDM}}$ ).** To align the generation process with the interpretable latent space  
555 learned by the adapter, we train the latent diffusion model (LDM) using a masked denoising loss.  
556 By restricting denoising to the unmasked concept channels, the LDM is guided to produce images  
557 consistent with the specified concept mask.

558 At each diffusion timestep  $t$ , Gaussian noise  $\epsilon \sim \mathcal{N}(0, I)$  is added to the adapter latent  $s_0$  to obtain  
559  $s_t = s_0 + \epsilon$ . The masked latent  $\tilde{s}_t = s_t \odot M(c)$  is then fed to the LDM, which predicts the added  
560 noise  $\hat{\epsilon} = \epsilon_\theta(\tilde{s}_t, t, c)$ . The masked denoising loss is computed only over the unmasked channels, as  
561

$$562 \mathcal{L}_{\text{LDM}} = \mathbb{E}_{x, \epsilon, t} \left[ \|(\epsilon - \epsilon_\theta(\tilde{s}_t, t)) \odot M(c)\|_2^2 \right].$$

563  
564 **Orthogonality Loss ( $\mathcal{L}_{\text{orth}}$ ).** Known concepts rarely capture all semantic variation, so we include  
565 residual channels to model unknown concepts. To prevent the adapter from encoding known concepts  
566 into these residual channels, we introduce an orthogonality loss (Ranasinghe et al., 2021; Ismail et al.,  
567 2023) between each concept block  $s^{(k)}$  and the residual block  $s^{(\text{res})}$ .

568 Formally, we compute the absolute cosine similarity between concept blocks and the residual block:

$$569 \mathcal{L}_{\text{orth}} = \mathbb{E}_{s \sim \mathcal{S}} \left[ \sum_{k=1}^K \left| \frac{\langle s^{(k)}, s^{(\text{res})} \rangle}{\|s^{(k)}\|_2 \|s^{(\text{res})}\|_2} \right| \right].$$

570  
571  
572 **Concept Classifiers + Concept Loss ( $\mathcal{L}_{\text{concept}}$ ).** To ensure each concept block reliably encodes its  
573 intended concept, we introduce a binary concept classifier for each block. This encourages the adapter  
574 to disentangle concepts and aligns latent representations with their semantic meaning.

575 For a binary concept  $k$ , we summarize the block’s features for both the active and inactive instances  
576 (by averaging over spatial dimensions) and concatenate them into a single vector:

$$577 \bar{s}^k = [\bar{s}^{k+}; \bar{s}^{k-}] \in \mathbb{R}^{2C_k}.$$

578  
579  
580 The classifier predicts which instance is active:

$$581 c_{\text{pred}}^{(k)} = q_k(\bar{s}^k) \in [0, 1], \quad q_k : \mathbb{R}^{2C_k} \rightarrow [0, 1].$$

582  
583  
584 We train the concept classifiers with a concept loss, the binary cross-entropy between the predicted  
585 label and the ground truth, summed over all concept blocks:

$$586 \mathcal{L}_{\text{concept}} = \mathbb{E}_{x \sim \mathcal{X}} \left[ \sum_{k=1}^K \text{BCE}(c_{\text{pred}}^{(k)}(x), y_{\text{GT}}^{(k)}(x)) \right].$$

587  
588  
589 This formulation can be easily extended to multi-class concepts by letting  $q_k$  output a probability  
590 distribution over the concept categories.

591  
592 **Cyclical Intervention Loss ( $\mathcal{L}_{\text{intervention}}$ ).** To enforce that interventions on concept assignments  
593 produce the intended effect, we use a cyclical intervention loss (Kulkarni et al., 2025b). This

594 loss encourages the model to produce latent representations that respond predictably to manual  
 595 interventions, reinforcing steerability.

596 Concretely, given a latent  $s$ , we first remark it according to the intended intervention labels, denoise  
 597 it through the LDM, and then decode and re-encode it via the adapter. Each concept block of the  
 598 resulting latent is then passed through its classifier  $q_k$  to obtain predicted labels  $c_{\text{pred}}^{(k)}$ .  
 599

600 The intervention loss is the cross-entropy between these predicted labels and the intended intervention  
 601 labels, summed over all concepts and averaged over the dataset:

$$602 \mathcal{L}_{\text{intervention}} = \mathbb{E}_{x \sim \mathcal{X}} \left[ \sum_{k=1}^K \text{CE}(c_{\text{pred}}^{(k)}(x), c_{\text{intv}}^{(k)}(x)) \right].$$

606 **Optimization-based Intervention.** Finally, we can also perform optimization-based intervention  
 607 at inference time, following Wong et al. (2020); Kulkarni et al. (2025b). This allows us to directly  
 608 manipulate concept activations in the latent space to achieve the desired intervention effect in the  
 609 generated image.

610 Given a generated latent  $\hat{s}_0$  with concept assignment  $c$ , we optimize the latent to match a target  
 611 concept assignment  $c^*$  while staying close to the original latent:

$$612 \tilde{s}_0^* = \tilde{s}_0 + \arg \max_{\delta \in \Delta} \left[ -\mathcal{L}_{\text{concept}}(E_{\text{adapter}}(D_{\text{adapter}}(\tilde{s}_0 + \delta)), c^*) \right],$$

616 where  $\Delta = \{\delta : \|\delta\|_{\infty} \leq \epsilon\}$  constrains the perturbation magnitude, and  $\mathcal{L}_{\text{concept}}$  is the cross-entropy  
 617 loss computed via the concept classifiers.

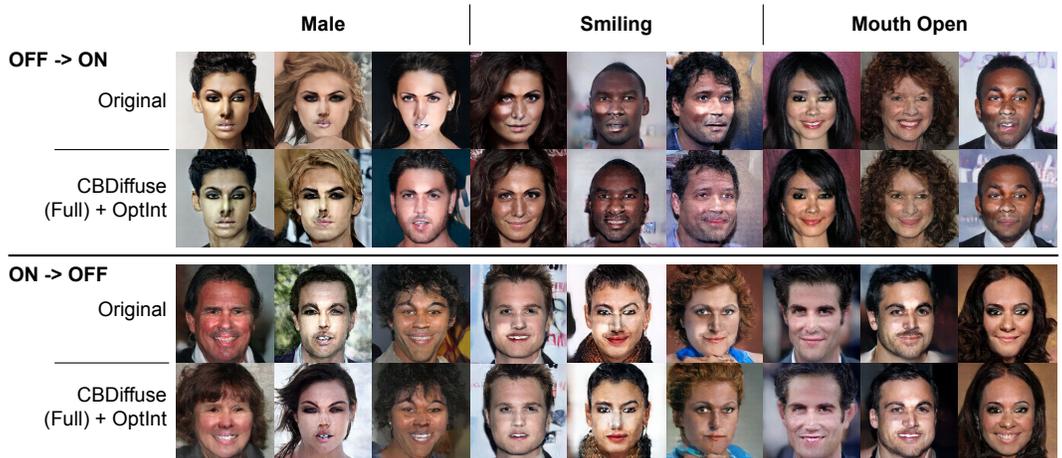
618 The final manipulated image is then obtained by decoding the optimized latent:

$$619 \hat{x}^* = D_{\text{vae}}(D_{\text{adapter}}(\tilde{s}_0^*)).$$

620  
621  
622  
623  
624  
625  
626  
627  
628  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647

648 **D ADDITIONAL RESULTS**

649 We provide additional example interventions for our CBDiffuse (Full) model. While image quality  
 650 may be slightly lower, concept steerability is improved.  
 651



668 Figure 4: Concept intervention results for Male, Smiling, and Mouth Open. The top half shows  
 669 images where the concept is flipped from OFF to ON, while the bottom half shows flips from ON  
 670 to OFF. For each concept, the first row displays the original images and the second row shows the  
 671 corresponding images generated by CBDiffuse (Full) after intervention.  
 672



684 Figure 5: Generation of Heavy Makeup (OFF → ON) conditioned on gender using CBDiffuse  
 685 (Full). While the model can generate Heavy Makeup on men despite limited training data, interven-  
 686 tions are more challenging for males than for females.  
 687

688 **D.1 CONCEPT COMBINATION BIASES**

689  
 690 Table 4: Per-concept activation steerability (%) conditioned on gender for CBDiffuse (Core).  
 691 Makeup-related concepts exhibit substantially lower steerability for men, while gender neutral  
 692 concepts have balanced steerability across gender, reflecting dataset imbalance.  
 693

Concept (OFF → ON)	Male	Female
Heavy Makeup	13.5	55.8
Lipstick	47.7	92.0
Mouth Open	89.6	90.3
Smiling	71.3	65.2

694  
695  
696  
697  
698  
699

700 Real-world datasets exhibit strong biases in the joint distribution of concepts, with some combinations  
 701 appearing far less frequently than others. For instance, in CelebA-HQ, images of men wearing heavy

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

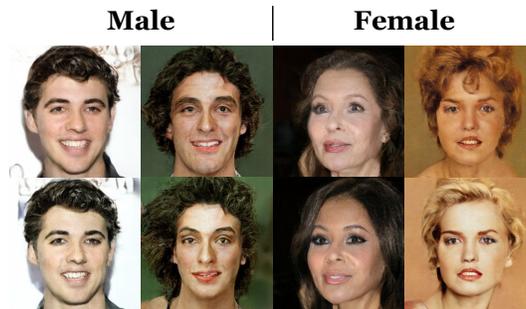


Figure 6: Generation of Heavy Makeup (OFF  $\rightarrow$  ON) conditioned on gender using CBDiffuse (Core). While the model can generate Heavy Makeup on men despite limited training data, interventions are more challenging for males than for females.

makeup or lipstick are substantially rarer than the corresponding combinations for women, reflecting underlying cultural and dataset biases. This raises an important question: *how well can generative models handle controlled generation for underrepresented concept combinations?*

In Figure 6, we demonstrate that CBDiffuse can generate uncommon concept combinations, such as heavy makeup and lipstick on men, despite limited training data coverage. However, these rare combinations remain more challenging to generate effectively. To quantify this, Table 4 reports per-concept steerability conditioned on gender, focusing on two gender-biased and two gender-neutral concepts. We observe a significant drop in steerability for makeup-related concepts when conditioned on *male* compared to *female*, while gender-neutral concepts maintain balanced controllability across genders.

This behavior highlights a fundamental challenge in compositional generative modeling: although the number of possible concept combinations grows exponentially with the number of concepts, the effective support of the training data covers only a small fraction of this space. Consequently, interventions that push the model toward low-density regions of the joint concept distribution are inherently more difficult. Prior work has noted that enforcing concept changes in diffusion models is particularly challenging due to the iterative nature of the denoising process, which can resist out-of-distribution edits as noise is gradually removed (Anonymous, 2026).

In contrast to approaches that intervene directly in pixel space or late in the generation process, CBDiffuse operates in a concept-aligned latent space, allowing concept information to propagate more consistently throughout the denoising trajectory. While this does not eliminate dataset biases, it enables more faithful and stable interventions for rare concept combinations compared to baseline methods.

Overall, these results suggest that, although CBDiffuse is still constrained by the underlying concept distribution of the training data, its architecture mitigates some of the challenges associated with rare combinations by explicitly structuring and preserving concept information throughout the diffusion process.