



ArcDAE

Asymmetric Rectified Contrastive Diffusion Autoencoder for Unified Representation Learning

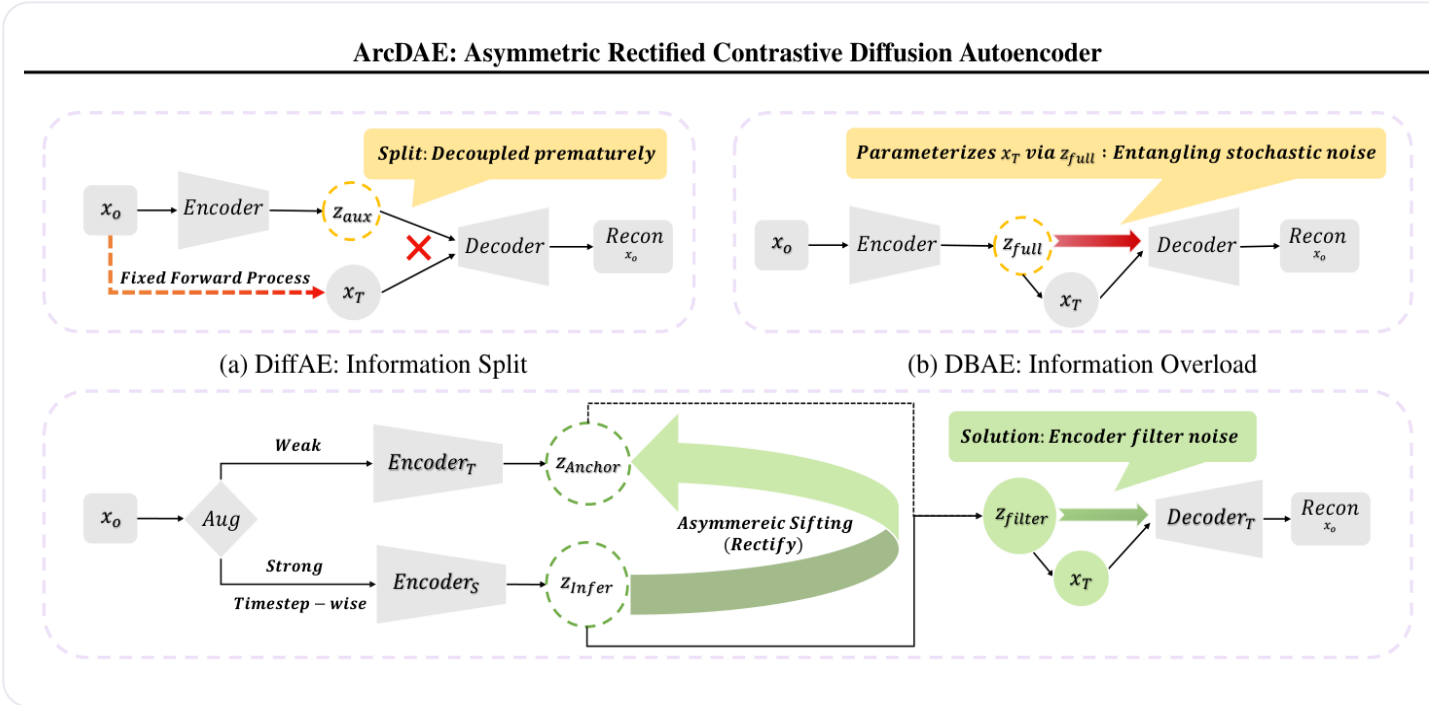
Ge Gao · Di Xiong · Zeke Xie · Jian Yang · Shuo Chen*

Nanjing University · HKUST (Guangzhou) · Nanjing University of Science and Technology

ICML 2026 · Seoul, South Korea

Limitations of Existing Works

Two structural extremes: information split and information overload.



Paper Fig. 1: conceptual comparison

Existing designs oscillate between two failures

Information split

Auxiliary semantic code is decoupled from the generative driver, limiting representational completeness.

Information overload

Bridge bottlenecks reconstruct noisy endpoints and entangle high-frequency stochastic variation.

understanding
invariance

reconstruction
fidelity

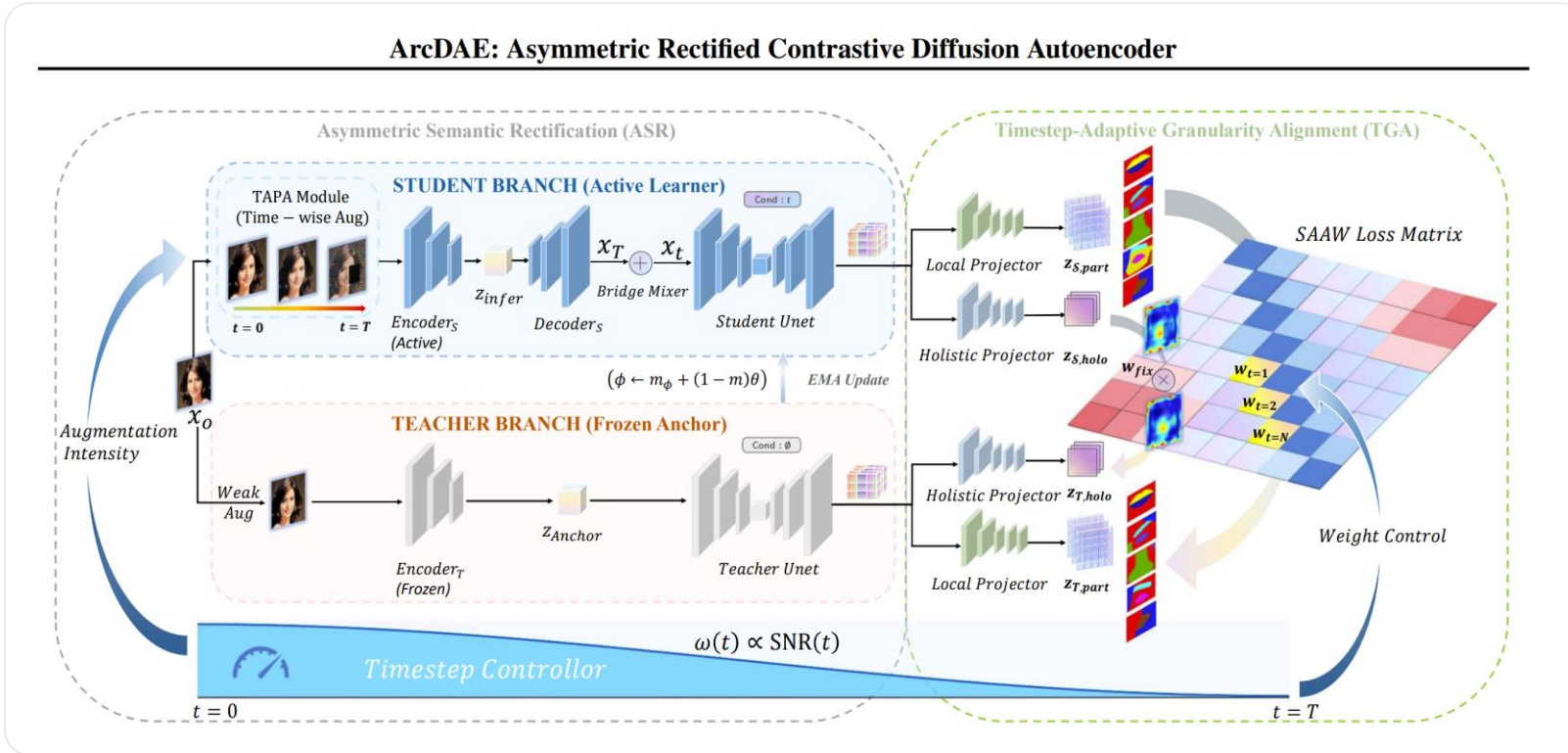
ArcDAE targets the high-fidelity semantic region

Core hypothesis

A robust encoder should act as an information sifter, not a passive reconstruction bridge.

Our Proposed Method: ArcDAE

Rebuilding the diffusion bridge as a dynamic information sifter.



Paper Fig. 2: asymmetric teacher–student framework

Three design ingredients

- ASR: a clean teacher anchor rectifies the noisy student representation.
- TAPA: augmentation strength increases with diffusion timestep.
- TGA: global and local alignment are gated by signal-to-noise ratio.

Clean semantic anchor → noisy student inference → timestep-aware contrastive rectification

Objective: Rectify Semantics, Gate Granularity

The loss balances diffusion reconstruction with contrastive semantic alignment.

ASR Smoothness Effect

High-noise timesteps strengthen contraction against unstable nuisance variation.

$$\begin{cases} \mathbf{z}^T = \text{sg} [E_\phi(\mathbf{v}^T)], & \mathbf{v}^T \sim p_{\text{aug}}(\cdot | \mathbf{x}_0, 0), \\ \mathbf{x}_t^S = \alpha_t \mathbf{v}^S + \sigma_t \boldsymbol{\epsilon}, & \mathbf{v}^S \sim p_{\text{aug}}(\cdot | \mathbf{x}_0, \gamma(t)), \end{cases}$$

$$\begin{aligned} \mathcal{L}_{\text{ASR}} = & \underbrace{\|E_\theta(\mathbf{m}_t) - \mathbf{z}^T\|_2^2}_{\text{Alignment error}} + \underbrace{\sigma_t^2 \text{Tr}(\mathbf{J}_{E_\theta}(\mathbf{m}_t)^\top \mathbf{J}_{E_\theta}(\mathbf{m}_t))}_{\text{Jacobian regularization}} \\ & + \mathcal{O}\left(\sigma_t^2 \|E_\theta(\mathbf{m}_t) - \mathbf{z}^T\|_2 \|\mathbf{H}_{E_\theta}(\mathbf{m}_t)\|_{\text{op}} + \sigma_t^4\right). \end{aligned}$$

Holistic & Local Representations

$$\begin{aligned} \mathbf{z}_{\text{part}} &= \mathcal{P}_{\text{part}} \left(\int_{\Omega} \mathbf{h}(\mathbf{u}) d\rho_{\mathbf{q}}(\mathbf{u}) \right), & Z_\rho(\mathbf{q}) &:= \int_{\Omega} \exp\left(\frac{\langle \mathbf{q}, \phi_K(\mathbf{h}(\mathbf{u})) \rangle}{\sqrt{d_k}}\right) d\mu(\mathbf{u}), \\ \mathbf{z}_{\text{holo}} &= \mathcal{P}_{\text{holo}} \left(\int_{\Omega} \mathbf{h}(\mathbf{u}) d\mu(\mathbf{u}) \right), & \frac{d\rho_{\mathbf{q}}}{d\mu}(\mathbf{u}) &= \frac{1}{Z_\rho(\mathbf{q})} \exp\left(\frac{\langle \mathbf{q}, \phi_K(\mathbf{h}(\mathbf{u})) \rangle}{\sqrt{d_k}}\right). \end{aligned}$$

data-like
high SNR

preserve local structure

sift stochastic noise

noise-like
low SNR

TGA Module : Timestep-Aware Adaptive Weighting

$t \rightarrow 0$: local structure preserved · $t \rightarrow T$: local noise alignment attenuated

$$\omega(t) \equiv \omega_{\tau_{\text{snr}}}(t) := \frac{\text{SNR}(t)^{1/\tau_{\text{snr}}}}{1 + \text{SNR}(t)^{1/\tau_{\text{snr}}}}, \quad \tau_{\text{snr}} > 0.$$

Total training objective

$$\begin{aligned} \mathcal{L}_{\text{total}} &= \mathcal{L}_{\text{diff}}(\theta) \\ &+ \lambda_{\text{align}} \left[\ell_{\text{NCE}}(\mathbf{z}_{\text{holo}}^S, \mathbf{z}_{\text{holo}}^T) + \omega(t) \cdot \ell_{\text{NCE}}(\mathbf{z}_{\text{part}}^S, \mathbf{z}_{\text{part}}^T) \right]. \end{aligned}$$

Interpretation

The diffusion loss keeps reconstructive ability; the gated contrastive term prevents the latent code from becoming a noise map.

Experimental Setup & Results & Analysis I

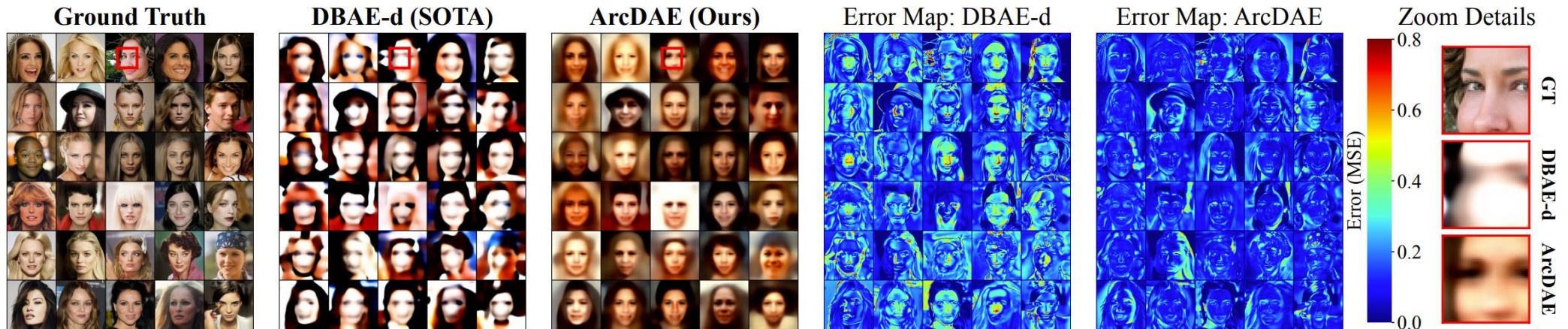
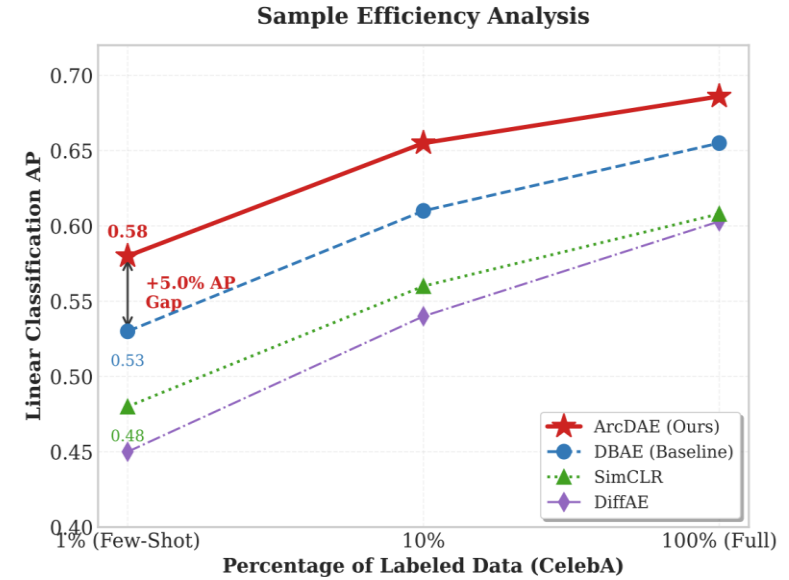
Evaluation covers semantic density, reconstruction, generation, and ablations.

Datasets CelebA · FFHQ · CelebA-HQ; additional non-face benchmarks

Backbone ADM U-Net encoder / score network; latent dimension $d_z = 512$

Baselines SimCLR, β -TCVAE, IB-GAN, DiffAE, PDAE, DiTi, DBAE

Metrics AP, Pearson's r , MSE, SSIM, LPIPS, FID, Precision / Recall, TAD



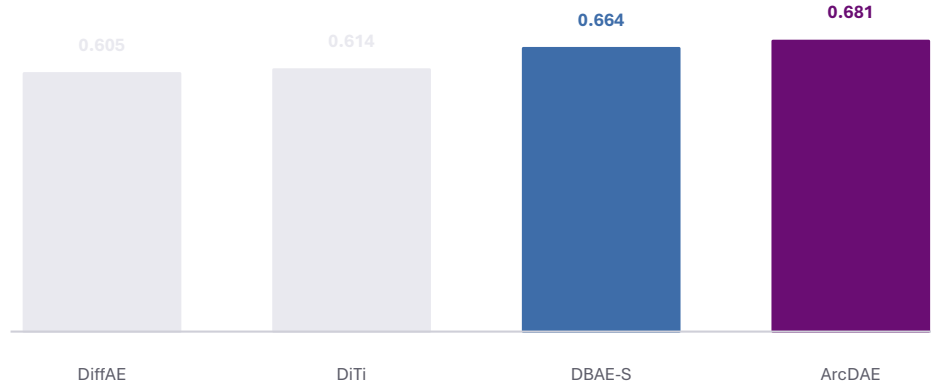
Early convergence: fewer artifacts and lower error maps

Main question: Does sifting improve semantics without sacrificing generative fidelity?

Results & Analysis II

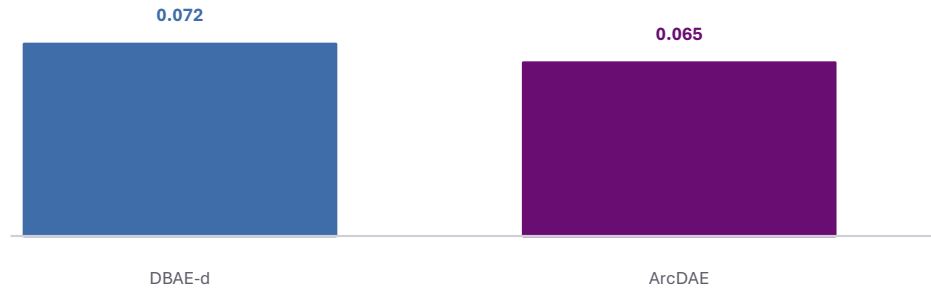
ArcDAE improves semantic density and perceptual reconstruction.

Linear-probe AP on FFHQ



METHOD	GEN	CELEBA			FFHQ		
		AP (\uparrow)	PEARSON'S r (\uparrow)	MSE (\downarrow)	AP (\uparrow)	PEARSON'S r (\uparrow)	MSE (\downarrow)
SIMCLR (CHEN ET AL., 2020)	\times	0.597	0.474	0.603	0.608	0.481	0.638
β -TCVAE (CHEN ET AL., 2018)	\checkmark	0.450	0.378	0.573	0.432	0.335	0.608
IB-GAN (JEON ET AL., 2021)	\checkmark	0.442	0.307	0.597	0.428	0.260	0.644
DIFFAE (PREECHAKUL ET AL., 2022)	\checkmark	0.603	0.598	0.421	0.605	0.606	0.410
PDAE (ZHANG ET AL., 2022)	\checkmark	0.602	0.596	0.410	0.597	0.603	0.416
DiTi (YUE ET AL., 2024)	\checkmark	0.623	0.617	0.392	0.614	0.622	0.384
DBAE (DET.) (KIM ET AL., 2025)	\checkmark	0.650	0.635	0.413	0.656	0.638	0.404
DBAE (STOCH.) (KIM ET AL., 2025)	\checkmark	0.655	0.643	0.369	0.664	0.675	0.332
ArcDAE (OURS)	\checkmark	0.676	0.656	0.345	0.681	0.718	0.319

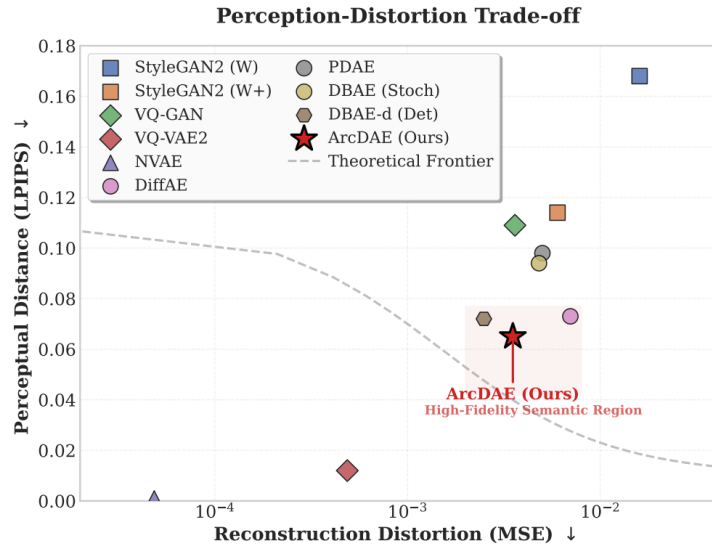
LPIPS reconstruction on CelebA-HQ \downarrow



METHOD	TRACTABILITY	LATENT DIM (\downarrow)	SSIM (\uparrow)	LPIPS (\downarrow)	MSE (\downarrow)
STYLEGAN2 (\mathcal{W}) (KARRAS ET AL., 2020)	\times	512	0.677	0.168	0.016
STYLEGAN2 ($\mathcal{W}+$) (ABDAL ET AL., 2019)	\times	7,168	0.827	0.114	0.006
VQ-GAN (ESSER ET AL., 2021)	\checkmark	65,536	0.782	0.109	3.61×10^{-3}
VQ-VAE2 (RAZAVI ET AL., 2019)	\checkmark	327,680	0.947	0.012	4.87×10^{-4}
NVAE (VAHDAT & KAUTZ, 2020)	\checkmark	6,005,760	0.984	0.001	4.85×10^{-5}
DDIM (INFERRED \mathbf{x}_T) (SONG ET AL., 2021A)	\times	49,152	0.917	0.063	0.002
DIFFAE (INFERRED \mathbf{x}_T) (PREECHAKUL ET AL., 2022)	\times	49,664	0.991	0.011	6.07×10^{-5}
PDAE (INFERRED \mathbf{x}_T) (ZHANG ET AL., 2022)	\times	49,664	0.994	0.007	3.84×10^{-5}
DIFFAE (RANDOM \mathbf{x}_T) (PREECHAKUL ET AL., 2022)	\checkmark	512	0.677	0.073	0.007
PDAE (RANDOM \mathbf{x}_T) (ZHANG ET AL., 2022)	\checkmark	512	0.689	0.098	5.01×10^{-3}
DBAE (STOCHASTIC) (KIM ET AL., 2025)	\checkmark	512	0.920	0.094	4.81×10^{-3}
DBAE-D (DETERMINISTIC) (KIM ET AL., 2025)	\checkmark	512	0.953	0.072	2.49×10^{-3}
ArcDAE (OURS)	\checkmark	512	0.962	0.065	2.79×10^{-3}

Results & Analysis II

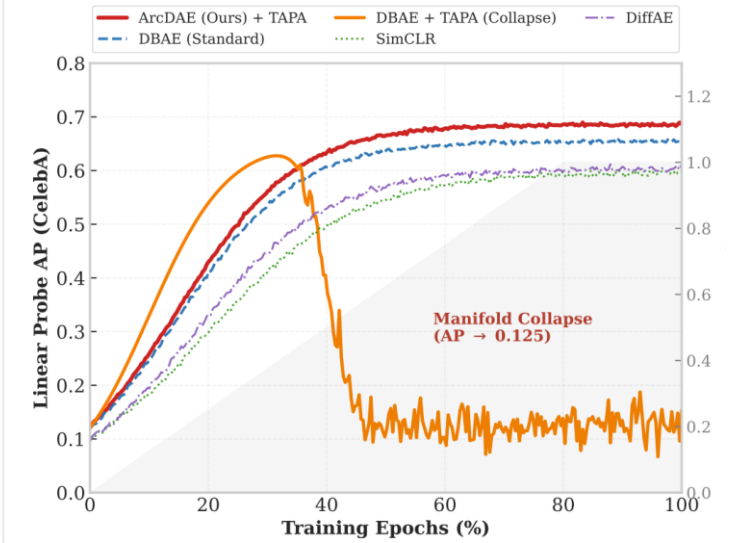
Sifted latent manifolds support generation, sample efficiency, and stability.



METHOD	REG	TAD (\uparrow)	ATTRS (\uparrow)	FID (\downarrow)
AE	\times	$0.042 \pm .004$	$1.0 \pm .0$	90.4 ± 1.8
DIFFAE (PREECHAKUL ET AL., 2022)	\times	$0.155 \pm .010$	$2.0 \pm .0$	22.7 ± 2.1
DBAE (KIM ET AL., 2025)	\times	$0.124 \pm .078$	2.2 ± 1.3	$11.8 \pm .2$
ARCDAE	\times	$0.210 \pm .070$	$3.9 \pm .6$	$10.28 \pm .2$
VAE (KINGMA & WELLING, 2014)	\checkmark	$0.000 \pm .000$	$0.0 \pm .0$	94.3 ± 2.8
β -VAE (HIGGINS ET AL., 2017)	\checkmark	$0.088 \pm .051$	$1.6 \pm .8$	99.8 ± 2.4
INFOVAE (ZHAO ET AL., 2019)	\checkmark	$0.000 \pm .000$	$0.0 \pm .0$	77.8 ± 1.6
INFODIFF (WANG ET AL., 2023)	\checkmark	$0.299 \pm .006$	$3.0 \pm .0$	22.3 ± 1.2
DISDIFF (YANG ET AL., 2023)	\checkmark	$0.305 \pm .010$	-	18.3 ± 2.1
DBAE+TC (KIM ET AL., 2025)	\checkmark	$0.362 \pm .036$	$3.8 \pm .8$	$13.4 \pm .2$
ARCDAE+TC	\checkmark	$0.385 \pm .036$	$4.9 \pm .8$	$12.2 \pm .2$

METHOD	PREC (\uparrow)	IS (\uparrow)	FID (\downarrow)	REC (\uparrow)
DDPM (HO ET AL., 2020)	0.768	3.11	9.14	0.335
DIFFAE (PREECHAKUL ET AL., 2022)	0.762	2.98	<u>9.40</u>	0.458
PDAE (ZHANG ET AL., 2022)	0.695	2.23	47.42	0.153
DBAE (KIM ET AL., 2025)	<u>0.780</u>	<u>3.87</u>	11.25	<u>0.392</u>
ARCDAE (OURS)	0.785	3.92	10.28	0.437
DIFFAE + AE	0.750	3.63	2.84	0.685
DBAE + AE	0.751	3.57	<u>1.77</u>	<u>0.687</u>
ARCDAE + AE	0.792	4.10	1.65	0.710

Disentanglement and generation



Teacher anchor prevents collapse

Key takeaways

- 1% labels: ArcDAE AP 0.58 vs. DBAE 0.53.
- ArcDAE+TC reaches TAD 0.385 with FID 12.2.
- Full model gives best ablation: AP 0.681, LPIPS 0.065, FID 10.28.

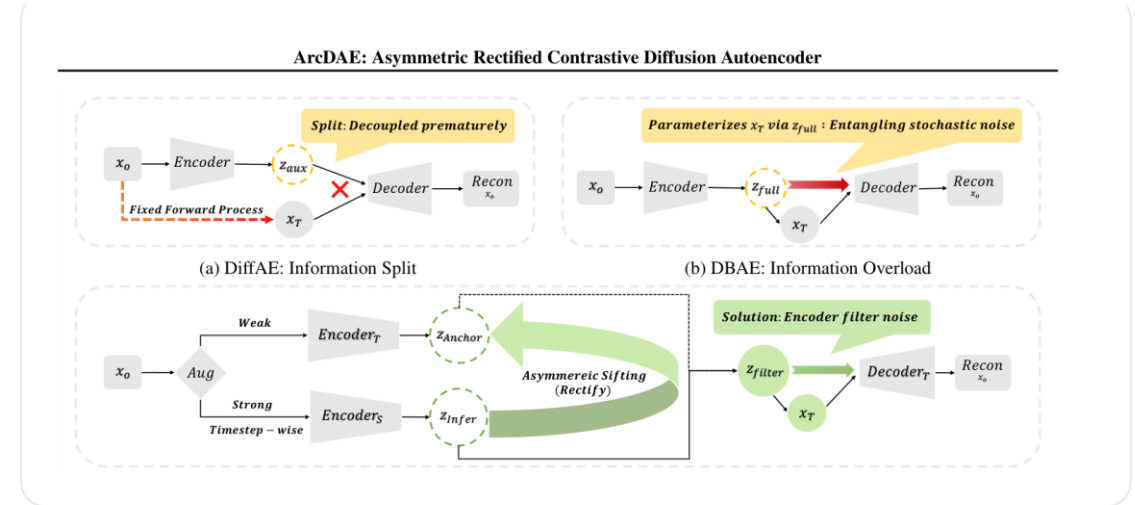
Conclusion & Future Work

ArcDAE resolves the split-overload tension with asymmetric information sifting.

Main conclusion

ArcDAE turns the diffusion bridge from a passive information carrier into an adaptive semantic sifter.

Problem	Information split loses completeness; information overload entangles noise.
Method	ASR + TAPA + TGA/TAAW rectify noisy representations toward clean semantic anchors.
Evidence	Better semantic probing, perceptual fidelity, generation quality, and ablation stability.
Future	Tighter theory for rectification dynamics; multimodal and large-scale representation learning.



Take-home message

Noise-rectified representations can unify synthesis and understanding.

Acknowledgement

Thank you for your attention.



ArcDAE: Asymmetric Rectified Contrastive Diffusion Autoencoder for Unified Representation Learning

Authors: Ge Gao, Di Xiong, Zeke Xie, Jian Yang, Shuo Chen*

Supported by Nanjing University and collaborating institutions.
We thank the ICML reviewers and the broader diffusion representation learning community.

Thank you!

