REPA-E: Unlocking VAE for End-to-End Tuning with Latent Diffusion

Transformers

Xingjian LengJaskirat SinghYunzhong HouAustralian National University

Saining Xie

Zhenchang Xing Data61 CSIRO

New York University

Liang Zheng Australian National University

> 2025.5.25 Presenter:唐果





- Author
- Background
- Method
- Experiments







- Confines diffusion model training to a fixed latent space
- Interact between AE & LDM
- REPA^[1] align early-layer features of DiT with clean image features from pre-trained vision encoders



Insights

- High-frequency noise in AE
- Scale equivariance: a simple regularization





Blockwise 2D DCT

- Autoencoder is different from RGB
- Bigger channel has higher frequencies





Spectral Analysis





Figure 6. Denoising trajectories (steps 1, 16, 32, 128 and 256 out of 256) for DiT-XL trained with FluxAE (top) and FluxAE+SE. DiT-XL with vanilla FluxAE exhibits prominent high-frequency artifacts early on in the trajectory.



Scale Equivariance Regularization

- ×2 4 bilinear downsampling \tilde{x}
- Removes high-frequency & remain low-frequency

$$\mathcal{L}(x) = d(x, \operatorname{Dec}(z)) + \alpha d(\tilde{x}, \operatorname{Dec}(\tilde{z})) + \beta \mathcal{L}_{\mathrm{KL}}$$





Insights

- Lack semantic-preserving in AE
- Equivariance Regularization: a simple regularization





Semantic-preserving

EQ-VAE

- Spatial Transformation $\tau \circ \mathbf{X}$
- Equivariance $\mathcal{E}(\tau \circ \mathbf{x}) = \tau \circ \mathcal{E}(\mathbf{x})$

$$\mathcal{L}_{\text{explicit}}(\mathbf{x}) = \| \tau \circ \mathcal{E}(\mathbf{x}) - \mathcal{E}(\tau \circ \mathbf{x}) \|_2^2$$





Equivariance Regularization







Motivation

- Uneven distribution of potential features in AE
- Contradiction between reconstruction and generation







Vision Foundation model alignment loss







Marginal Cosine Similarity Loss

Image latents Z & foundational visual representations F

Z' = WZ

$$\mathcal{L}_{\text{mcos}} = \frac{1}{h \times w} \sum_{i=1}^{h} \sum_{j=1}^{w} \text{ReLU} \left(1 - m_1 - \frac{z'_{ij} \cdot f_{ij}}{\|z'_{ij}\| \|f_{ij}\|} \right)$$





Marginal Distance Matrix Similarity Loss

• Reduce relative distribution distance

$$\mathcal{L}_{\text{mdms}} = \frac{1}{N^2} \sum_{i,j} \text{ReLU} \left(\left| \frac{z_i \cdot z_j}{\|z_i\| \|z_j\|} - \frac{f_i \cdot f_j}{\|f_i\| \|f_j\|} \right| - m_2 \right)$$





Adaptive Weighting

• Ensuring similar impacts on model optimization

$$w_{\text{adaptive}} = \frac{\|\nabla L_{\text{rec}}\|}{\|\nabla L_{\text{vf}}\|}$$

$$\mathcal{L}_{\rm vf} = w_{\rm hyper} * w_{\rm adaptive} (\mathcal{L}_{\rm mcos} + \mathcal{L}_{\rm mdms})$$







- Recent works fintune the VAE to improve the equivariance
- Jointly tune both VAE and LDM in an end-to-end manner
- End-to-end training can lead to increased performance





Naive End-to-End Tuning is Ineffective

Simpler latent-space

| Training Strategy | Spatial Variance | Total Variation |
|--------------------------|-------------------------|------------------------|
| w/o E2E Tuning | 17.06 | 6627.35 |
| E2E w/ REPA Loss | 18.02 | 5516.14 |
| E2E w/ Diff. Loss | 0.02 | 89.80 |

• Reduced generation performance



(a) PCA Visualization of Latent Spaces

Method



CKNNA scores

• High CKNNA, better generation

• Optimize CKNNA score(REPA loss)



(b) Correlation: gFID & CKNNA Score





Representation Alignment is Bottlenecked by VAE

Frozen VAE has limitation







Batch-Norm Layer for VAE Latent Normalization

- From precomputed latent statistics(1/0.1825 for SD-VAE) to
- A new batch-norm layer







End-to-End Representation-Alignment Loss

• DiT hidden state gets projection

$$\mathcal{L}_{\text{REPA}}(\theta, \phi, \omega) = -\mathbb{E}_{\mathbf{x}, \epsilon, t} \left[\frac{1}{N} \sum_{n=1}^{N} \sin(\mathbf{y}^{[n]}, h_{\omega}(\mathbf{h}_{t}^{[n]})) \right]$$

Overall Training

$$\mathcal{L}(\theta, \phi, \omega) = \mathcal{L}_{\text{DIFF}}(\theta) + \lambda \mathcal{L}_{\text{REPA}}(\theta, \phi, \omega) + \eta \mathcal{L}_{\text{REG}}(\phi)$$
$$\lambda_{\text{REPA}_g} = 0.5 \qquad \qquad \lambda_{\text{REPA}_v} = 1.5$$

04 Experiments







Significantly improve generation performance and

| training chood | 12 <u>0</u> | | | | | | | | | |
|----------------|---------------------------|--------------|----------|----------|-------|--------------|--|--|--|--|
| training speed | Method | Tokenizer | Epochs | gFID↓ | sFID↓ | IS↑ | | | | |
| | Without End-to-End Tuning | | | | | | | | | |
| | MaskDiT [55] | | 1600 | 5.69 | 10.34 | 177.9 | | | | |
| | DiT [34] | CD VAE | 1400 | 9.62 | 6.85 | 121.5 | | | | |
| | SiT [30] | SD-VAE | 1400 | 8.61 | 6.32 | 131.7 | | | | |
| | FasterDiT [50] | | 400 | 7.91 | 5.45 | 131.3 | | | | |
| | | | 20 | 19.40 | 6.06 | 67.4 | | | | |
| | | CD VAE | 40 | 11.10 | 6.06 | 67.4 | | | | |
| | KEFA [JJ] | SD-VAE | 80 | 7.90 | 5.06 | 122.6 | | | | |
| | | | 800 | 5.90 | 5.73 | 157.8 | | | | |
| | V | Vith End-to- | End Tuni | ng (Ours | s) | | | | | |
| | | | 20 | 12.83 | 5.04 | 88.8 | | | | |
| | REPA-E | SD-VAE* | 40 | 7.17 | 4.39 | 123.7 | | | | |
| | | | 80 | 4.07 | 4.60 | <u>161.8</u> | | | | |



Significantly improve generation performance and training speed





| Method | Tokenizer | Epochs | gFID↓ | sFID↓ | IS↑ | | | | | | | |
|-------------------------------|-----------|--------|-------|-------|-------------|--|--|--|--|--|--|--|
| Without End-to-End Tuning | | | | | | | | | | | | |
| MaskDiT [55] | | 1600 | 5.69 | 10.34 | 177.9 | | | | | | | |
| DiT [34] | SD WAE | 1400 | 9.62 | 6.85 | 121.5 | | | | | | | |
| SiT [30] | SD-VAE | 1400 | 8.61 | 6.32 | 131.7 | | | | | | | |
| FasterDiT [50] | | 400 | 7.91 | 5.45 | 131.3 | | | | | | | |
| | | 20 | 19.40 | 6.06 | <u>67.4</u> | | | | | | | |
| | SD-VAE | 40 | 11.10 | 6.06 | 67.4 | | | | | | | |
| KEPA [33] | | 80 | 7.90 | 5.06 | 122.6 | | | | | | | |
| | | 800 | 5.90 | 5.73 | 157.8 | | | | | | | |
| With End-to-End Tuning (Ours) | | | | | | | | | | | | |
| | | 20 | 12.83 | 5.04 | 88.8 | | | | | | | |
| REPA-E | SD-VAE* | 40 | 7.17 | 4.39 | 123.7 | | | | | | | |
| | | 80 | 4.07 | 4.60 | 161.8 | | | | | | | |



| | | | | | | Target Repr. | gFID↓ | sFID↓ | IS↑ | Prec. ↑ | Rec.↑ |
|----------------|-------|-------|------|--------|-------|----------------|-------|-------|-------------|----------------|-------|
| Diff. Model | gFID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | I-JEPA-H [2] | 23.0 | 5.81 | 60.3 | 0.62 | 0.60 |
| SiT-B (130M) | 49.5 | 7.00 | 27.5 | 0.46 | 0.59 | +REPA-E (Ours) | 16.5 | 5.18 | 73.6 | 0.68 | 0.60 |
| +REPA-E (Ours) | 34.8 | 6.31 | 39.1 | 0.57 | 0.59 | CLIP-L [37] | 29.2 | 5.98 | <u>46.4</u> | 0.59 | 0.61 |
| SiT-L (458M) | 24.1 | 6.25 | 55.7 | 0.62 | 0.60 | +REPA-E (Ours) | 23.4 | 6.44 | 57.1 | 0.62 | 0.60 |
| +REPA-E (Ours) | 16.3 | 5.69 | 75.0 | 0.68 | 0.60 | DINOv2-B [33] | 24.1 | 6.25 | 55.7 | 0.62 | 0.60 |
| SiT-XL (675M) | 19.4 | 6.06 | 67.4 | 0.64 | 0.61 | +REPA-E (Ours) | 16.3 | 5.69 | 75.0 | 0.68 | 0.60 |
| +REPA-E (Ours) | 12.8 | 5.04 | 88.8 | 0.71 | 0.58 | DINOv2-L [33] | 23.3 | 5.89 | 59.9 | 0.61 | 0.60 |
| | | | | | | +REPA-E (Ours) | 16.0 | 5.59 | 77.7 | 0.68 | 0.58 |





| Autoencoder | gFID↓ | sFID↓ | IS↑ | Prec.↑ | Rec. ↑ | Aln. Depth | gFID↓ | sFID↓ | IS↑ | Prec. [↑] | Rec. ↑ |
|-----------------|-------|-------|------|--------|---------------|----------------|-------|-------|------|--------------------|---------------|
| SD-VAE [39] | 24.1 | 6.25 | 55.7 | 0.62 | 0.60 | 6th layer | 23.0 | 5.72 | 59.2 | 0.62 | 0.60 |
| +REPA-E (Ours) | 16.3 | 5.69 | 75.0 | 0.68 | 0.60 | +REPA-E (Ours) | 16.4 | 6.64 | 74.3 | 0.67 | 0.59 |
| IN-VAE (f16d32) | 22.7 | 5.47 | 56.0 | 0.62 | 0.62 | 8th layer | 24.1 | 6.25 | 55.7 | 0.62 | 0.60 |
| +REPA-E (Ours) | 12.7 | 5.57 | 84.0 | 0.69 | 0.62 | +REPA-E (Ours) | 16.3 | 5.69 | 75.0 | 0.68 | 0.60 |
| VA-VAE [49] | 12.8 | 6.47 | 83.8 | 0.71 | 0.58 | 10th layer | 23.7 | 5.91 | 56.9 | 0.62 | 0.60 |
| +REPA-E (Ours) | 11.1 | 5.31 | 88.8 | 0.72 | 0.61 | +REPA-E (Ours) | 16.2 | 5.22 | 74.7 | 0.68 | 0.58 |



| Component | gFID↓ | sFID↓ | IS↑ | Prec. ↑ | Rec. ↑ |
|-------------------------------|-------|-------|------|----------------|---------------|
| w/o stopgrad | 444.1 | 460.3 | 1.49 | 0.00 | 0.00 |
| w/obatch-norm | 18.1 | 5.32 | 72.4 | 0.67 | 0.59 |
| w/o $\mathcal{L}_{	ext{REG}}$ | 19.2 | 6.47 | 68.2 | 0.64 | 0.58 |
| REPA-E (Ours) | 16.3 | 5.69 | 75.0 | 0.68 | 0.60 |

| Method | gFID↓ | sFID↓ IS↑ | | Prec.↑ | Rec.↑ | | | | | | |
|-----------------------------|-------|-------------|-------|-------------|-------|--|--|--|--|--|--|
| 100K Iterations (20 Epochs) | | | | | | | | | | | |
| REPA [53] | 19.40 | 6.06 | 67.4 | 0.64 | 0.61 | | | | | | |
| REPA-E (scratch) | 14.12 | 7.87 | 83.5 | 0.70 | 0.59 | | | | | | |
| REPA-E (VAE init.) | 12.83 | 5.04 | 88.8 | 0.71 | 0.58 | | | | | | |
| 200K Iterations (40 Epochs) | | | | | | | | | | | |
| REPA [53] | 11.10 | 5.05 | 100.4 | 0.69 | 0.64 | | | | | | |
| REPA-E (scratch) | 7.54 | 6.17 | 120.4 | 0.74 | 0.61 | | | | | | |
| REPA-E (VAE init.) | 7.17 | 4.39 | 123.7 | 0.74 | 0.62 | | | | | | |
| 400K Iterations (80 Epochs) | | | | | | | | | | | |
| REPA [53] | 7.90 | 5.06 | 122.6 | 0.70 | 0.65 | | | | | | |
| REPA-E (scratch) | 4.34 | 4.44 | 154.3 | 0.75 | 0.63 | | | | | | |
| REPA-E (VAE init.) | 4.07 | 4.60 | 161.8 | 0.76 | 0.62 | | | | | | |



Experiments

The impact of end-to-end tuning



Experiments



The impact of end-to-end tuning

| VAE | Diffusion model | REPA | gFID-50K |
|----------------|-----------------|------|----------|
| SD-VAE [39] | DiT-XL [34] | × | 19.82 |
| VA-VAE [49] | DiT-XL [34] | × | 6.74 |
| E2E-VAE (Ours) | DiT-XL [34] | × | 6.75 |
| SD-VAE [39] | SiT-XL [30] | × | 17.20 |
| VA-VAE [49] | SiT-XL [30] | × | 5.93 |
| E2E-VAE (Ours) | SiT-XL [30] | × | 5.26 |
| SD-VAE [39] | DiT-XL [34] | 1 | 12.29 |
| VA-VAE [49] | DiT-XL [34] | 1 | 4.71 |
| E2E-VAE (Ours) | DiT-XL [34] | 1 | 4.20 |
| SD-VAE [39] | SiT-XL [30] | 1 | 7.90 |
| VA-VAE [49] | SiT-XL [30] | 1 | 4.88 |
| E2E-VAE (Ours) | SiT-XL [30] | 1 | 3.46 |



The impact of end-to-end tuning

| Tokenizer | Method | Training | #params | arams rFID↓ | ams rFID | | | | Generation w/ CFG | | | | | |
|-------------------------------|--------------------|----------|---------|---------------|------------|----------|---------|--------|-------------------|-------|-------|--------------|--------|-------|
| Tonomilor | | Epoches | | | gFID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ | gFID↓ | sFID↓ | IS↑ | Prec.↑ | Rec.↑ |
| AutoRegressive (AR) | | | | | | | | | | | | | | |
| MaskGiT | MaskGIT [4] | 555 | 227M | 2.28 | 6.18 | - | 182.1 | 0.80 | 0.51 | - | - | - | - | - |
| VQGAN | LlamaGen [45] | 300 | 3.1B | 0.59 | 9.38 | 8.24 | 112.9 | 0.69 | 0.67 | 2.18 | 5.97 | 263.3 | 0.81 | 0.58 |
| VQVAE | VAR [46] | 350 | 2.0B | - | | - | - | ÷. | - | 1.80 | - | 365.4 | 0.83 | 0.57 |
| LFQ tokenizers | MagViT-v2 [51] | 1080 | 307M | 1.50 | 3.65 | - | 200.5 | - | - | 1.78 | - | 319.4 | - | - |
| LDM | MAR [28] | 800 | 945M | 0.53 | 2.35 | - | 227.8 | 0.79 | 0.62 | 1.55 | - | 303.7 | 0.81 | 0.62 |
| Latent Diffusion Models (LDM) | | | | | | | | | | | | | | |
| | MaskDiT [55] | 1600 | 675M | | 5.69 | 10.34 | 177.9 | 0.74 | 0.60 | 2.28 | 5.67 | 276.6 | 0.80 | 0.61 |
| | DiT [34] | 1400 | 675M | | 9.62 | 6.85 | 121.5 | 0.67 | 0.67 | 2.27 | 4.60 | 278.2 | 0.83 | 0.57 |
| CD VAE [20] | SiT [30] | 1400 | 675M | 0.61 | 8.61 | 6.32 | 131.7 | 0.68 | 0.67 | 2.06 | 4.50 | 270.3 | 0.82 | 0.59 |
| 5D-VAE [39] | FasterDiT [50] | 400 | 675M | 0.01 | 7.91 | 5.45 | 131.3 | 0.67 | 0.69 | 2.03 | 4.63 | 264.0 | 0.81 | 0.60 |
| | MDT [12] | 1300 | 675M | | 6.23 | 5.23 | 143.0 | 0.71 | 0.65 | 1.79 | 4.57 | 283.0 | 0.81 | 0.61 |
| | MDTv2 [13] | 1080 | 675M | | - | - | - | - | - | 1.58 | 4.52 | 314.7 | 0.79 | 0.65 |
| | | | F | Represent | ation Alig | gnment N | lethods | | | | | | | |
| | | 80 | 675M | 0.39 | 4.29 | - | 17 | - | 177 | - | - | 177 | | 877.0 |
| VA-VAE [49] | LightningDi I [49] | 800 | 675M | 0.28 | 2.17 | 4.36 | 205.6 | 0.77 | 0.65 | 1.35 | 4.15 | 295.3 | 0.79 | 0.65 |
| SD VAE | | 80 | 675M | 0.61 | 7.90 | 5.06 | 122.6 | 0.70 | 0.65 | _ | - | 1 <u>-</u> 1 | | - |
| SD-VAE | KErA [33] | 800 | 675M | 0.01 | 5.90 | 5.73 | 157.8 | 0.70 | 0.69 | 1.42 | 4.70 | 305.7 | 0.80 | 0.65 |
| | DEDA | 80 | 675M | 0.29 | 3.46 | 4.17 | 159.8 | 0.77 | 0.63 | 1.67 | 4.12 | 266.3 | 0.80 | 0.63 |
| E2E-VAE (Ours) | KEFA | 800 | 675M | 0.28 | 1.83 | 4.22 | 217.3 | 0.77 | 0.66 | 1.26 | 4.11 | 314.9 | 0.79 | 0.66 |

Thanks for listening

