
Asymmetric Flow Models

Hansheng Chen Jan Ackermann Minseo Kim Gordon Wetzstein Leonidas Guibas

Stanford University

<https://hanshengchen.com/asymflow>



Figure 1: **AsymFLUX.2 klein generations.** AsymFlow finetunes FLUX.2 klein into a pixel-space flow model, producing highly realistic images with rich visual styles and fine detail.

Abstract

Flow-based generation in high-dimensional pixel spaces is difficult because velocity prediction requires modeling high-dimensional noise, even when data has strong low-rank structure. We present *Asymmetric Flow Modeling* (AsymFlow), a rank-asymmetric velocity parameterization that restricts noise prediction to a low-rank subspace while keeping data prediction full-dimensional. From this asymmetric prediction, AsymFlow analytically recovers the full-dimensional velocity without changing the network architecture or training/sampling procedures. On ImageNet 256×256, AsymFlow achieves a leading 1.57 FID, outperforming prior DiT/JiT-like pixel diffusion models by a large margin. AsymFlow also provides the first-ever route for finetuning pretrained latent flow models into pixel-space models: aligning the low-rank pixel subspace to the latent space gives a seamless initialization that preserves the latent model’s high-level semantics and structure, so finetuning mainly improves low-level mismatches rather than relearning pixel generation. We show that the pixel AsymFlow model finetuned from FLUX.2 klein 9B establishes a new state of the art for pixel-space text-to-image generation, beating its latent base on HPSv3, DPG-Bench, and GenEval while qualitatively showing substantially improved visual realism.

1 Introduction

Recent progress in diffusion-based image and video generation [5, 6, 18, 32, 63, 72] has been driven by combining scalable transformer architectures [7, 15, 48] with flow matching objectives [1, 40, 42]. Most state-of-the-art systems operate in compressed lower-dimensional latent spaces learned by autoencoders [51], which is highly scalable but delegates fine detail to a fixed decoder that the generative model cannot control. This limitation motivates a return to high-dimensional generation, including direct pixel-space generation [2, 9, 10, 27, 35, 45, 46, 64, 71].

However, moving to high-dimensional spaces exposes a bottleneck in velocity prediction. The velocity target $u = \epsilon - x_0$ consists of both data and noise components. To predict it accurately, the network must extract the noise from the input and pass it through its internal features. This is straightforward in latent spaces, where the noise dimension is small relative to the network width. In pixel space, however, the per-patch noise dimension can pollute the network’s internal states, creating a bottleneck [75]. Classical pixel diffusion models used U-Net architectures [14, 21, 28, 52, 54], whose skip connections naturally route noise from input to output. Modern scalable transformers lack these pathways, so recent methods either reintroduce architectural bypasses, such as U-ViT-like transformers [4, 11, 17, 22, 23] or decoder heads [10, 45, 62, 64, 71, 75], which complicates the otherwise simple transformer recipe, or switch to predicting clean data x_0 directly [35, 46, 58], which is numerically ill-conditioned at low noise levels [28, 55].

We introduce *Asymmetric Flow Modeling* (AsymFlow), a new parameterization for high-dimensional flow modeling that avoids both of these compromises. AsymFlow parameterizes the two velocity components asymmetrically: the data component remains full-dimensional, while the noise component is restricted to a low-rank subspace. The full-dimensional velocity is recovered analytically, so standard flow matching training and sampling remain unchanged. In this view, standard x_0 -prediction and u -prediction are special cases of AsymFlow, corresponding to zero and full rank of this noise subspace, respectively. Between these endpoints, AsymFlow can choose an intermediate rank that keeps velocity prediction in an important subspace while avoiding full-rank noise prediction.

In addition, AsymFlow makes it possible to build large-scale pixel generators by finetuning pretrained latent flow models. The key observation is that latent and pixel spaces are not disconnected: a latent model can be mathematically lifted into a low-rank pixel model whose samples inherit the semantics and structure of the latent generator. This turns latent-to-pixel adaptation into a correction problem, where finetuning keeps the high-level content and only needs to close the low-level projection gap between low-rank pixel outputs and full-rank pixel targets. To our knowledge, this is the first practical path for turning existing large-scale latent flow models themselves into strong pixel generators.

We evaluate AsymFlow in two settings. On ImageNet 256×256 [12], AsymFlow reaches 1.76 FID with the JiT-H/16 network [35] and 1.57 FID with an additional REPA loss [70], outperforming prior DiT/JiT-like pixel diffusion models by a large margin. For text-to-image generation, our pixel AsymFlow model finetuned from FLUX.2 klein 9B [6] sets a new state of the art in pixel-space generation, beating its latent base on HPSv3 [44], DPG-Bench [25], and GenEval [16] while qualitatively exhibiting substantially improved visual realism.

To summarize, our main contributions are:

- We introduce AsymFlow, a novel rank-asymmetric flow parameterization with full-rank data and low-rank noise for scalable high-dimensional generation.
- We provide the first method of finetuning pretrained latent flow models into pixel models through AsymFlow, using a principled latent-to-pixel lift without architectural modifications.
- We achieve a leading 1.57 FID on ImageNet 256×256 and demonstrate a 9B-scale pixel-space text-to-image model with state-of-the-art performance.

2 Related Work

Recent work mainly addresses the high-dimensional bottleneck in two ways: changing the network architecture so high-dimensional noisy inputs can reach the output more easily, or changing the prediction parameterization to avoid high-dimensional noise prediction.

Hierarchical architectures. One line of work keeps noise or velocity prediction feasible using hierarchical architectures with high-dimensional bypasses. Classical DDPM/ADM-style U-Nets [14,

21, 52] and U-ViT-like hierarchical transformers [4, 11, 17, 22, 23] use skip-connected multi-scale structures, while DDT-like decoder-based designs [65], including RAE, PixNerd, PixelDiT, DiP, and DeCo [10, 45, 62, 64, 71, 75], expose the noisy input to decoder or refiner pathways conditioned on backbone features. These designs are effective, but they complicate the plain transformer recipe that has scaled successfully in large image and video generators [5, 6, 18, 32, 63, 72]. In contrast, AsymFlow enables high-dimensional generation without architectural modification, making it possible to finetune large-scale latent flow models into pixel space for the first time.

Prediction parameterizations. In early diffusion models, hierarchical U-Net-like architectures made ϵ -prediction practical, while x_0 -prediction was often less favored because of low-noise numerical issues [21, 28, 55]. With the paradigm shift to plain diffusion transformers (DiT) [43, 48, 69], JiT [35] argues that pixel diffusion should predict clean data x_0 rather than noise or velocity, and several follow-up pixel methods [46, 58] adopt the same x_0 -prediction backbone with perceptual or representation-alignment (REPA) losses [70, 73]. k -Diff [27] learns a scalar interpolation between x_0 - and u -prediction, but this isotropic parameterization does not reduce the dimensionality of the noise component and gives results close to JiT. Unlike prior work, AsymFlow treats the prediction target asymmetrically: the data term x_0 remains full-dimensional, while the noise term ϵ is restricted to a low-rank subspace, which retains the benefits of u -prediction in a meaningful subspace.

3 Preliminaries

We briefly introduce diffusion models [21, 59, 60] using the flow matching convention [1, 40, 42], then review common prediction parameterizations.

Flow matching. Let $x_0 \in \mathbb{R}^D$ be a data vector of dimension D . A typical flow model defines an interpolation between a data sample and Gaussian noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, yielding the noisy sample $x_t := \alpha_t x_0 + \sigma_t \epsilon$, where $t \in (0, 1]$ denotes diffusion time and $\alpha_t = 1 - t$, $\sigma_t = t$ define the linear flow schedule. Under this construction, generative modeling is achieved by solving a reverse-time SDE or ODE that transports noise to data [41, 61]. In particular, the ODE velocity is given by $\frac{dx_t}{dt} = \mathbb{E}_{x_0 \sim p(x_0|x_t)} \left[\frac{x_t - x_0}{t} \right]$, which is the posterior mean of the sample velocity u :

$$u := \frac{x_t - x_0}{\sigma_t} = \epsilon - x_0. \quad (1)$$

Then, a model $(x_t, t) \mapsto \hat{u}$ is trained to estimate this posterior mean with the flow matching loss:

$$\mathcal{L}_{\text{FM}} = \mathbb{E}_{t, x_0, \epsilon} \left[\|u - \hat{u}\|^2 \right]. \quad (2)$$

u -prediction vs. x_0 -prediction. The mapping $(x_t, t) \mapsto \hat{u}$ is often directly parameterized by a neural network, i.e., $\hat{u} := G_\theta(x_t, t)$. This u -prediction form is widely used in modern latent flow models [15, 48, 51], where the representation is compressed. When moved to pixels or other high-dimensional representations, however, the target $u = \epsilon - x_0$ requires predicting a high-dimensional noise component in addition to structured data [35, 75]. An alternative is x_0 -prediction, where the network predicts clean data $\hat{x}_0 = G_\theta(x_t, t)$ and recovers velocity as $\hat{u} = (x_t - \hat{x}_0)/\sigma_t$. This avoids directly regressing Gaussian noise [35], but the $1/\sigma_t$ conversion is ill-conditioned at low noise levels [28, 55], limiting final-sample quality. Shin et al. [58] also claim that REPA-style alignment is less effective in x_0 -prediction pixel models. Thus, u - and x_0 -prediction expose complementary trade-offs where neither is ideal for high-dimensional generation.

4 Asymmetric Flow Modeling

To address the challenges of high-dimensional flow modeling, we introduce AsymFlow, a rank-asymmetric parameterization of the flow target. The key idea is to treat the two terms in the velocity target asymmetrically: the data prediction term remains full-dimensional, while the noise prediction is restricted to a low-rank subspace. This reduces the burden of representing high-dimensional noise in the network’s internal states without changing the network architecture. The full-rank velocity is then recovered analytically for training and sampling, leaving the flow matching formulation unchanged.

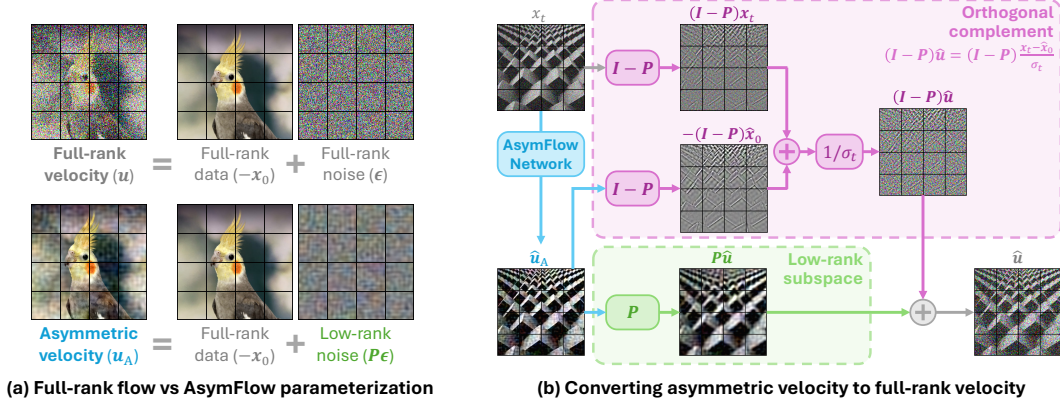


Figure 2: **AsymFlow parameterization and recovery.** (a) AsymFlow changes the standard velocity target by keeping the data term full-dimensional while replacing the noise term with its low-rank projection $P\epsilon$. (b) To recover the full-rank velocity, the low-rank component $P\hat{u}_A$ is used directly, while the orthogonal component is converted using the x_0 -to- u relation in Eq. (1).

4.1 AsymFlow Parameterization

Let $\mathbf{A} \in \mathbb{R}^{D \times r}$ be an orthonormal basis of a rank- r subspace, with $\mathbf{A}^T \mathbf{A} = \mathbf{I}_r$, and let $\mathbf{P} := \mathbf{A} \mathbf{A}^T$ be the corresponding orthogonal projector. Then $\text{Im}(\mathbf{P})$ is the low-rank subspace and $\text{Im}(\mathbf{I} - \mathbf{P})$ is its orthogonal complement. Given the noise $\epsilon \in \mathbb{R}^D$, we use $P\epsilon$ to denote its subspace component. We refer to $P\epsilon$ as *low-rank noise*, meaning Gaussian noise projected to a low-rank subspace.

AsymFlow changes the target that the network is asked to predict. In standard u -prediction (Eq. (1)), the output must reproduce the full noise component ϵ together with the data term $-x_0$. For high-dimensional data, this forces the model to carry high-dimensional noise through its features, which pollutes its internal states and wastes network capacity. To address this issue, AsymFlow introduces an *asymmetric velocity* u_A where the noise term is low-rank while the data term remains full-rank:

$$u_A := P\epsilon - x_0. \quad (3)$$

We then train the network to predict the asymmetric velocity, i.e., $\hat{u}_A = G_\theta(x_t, t)$. This prediction will be converted back to the full-rank velocity \hat{u} for loss calculation and denoising sampling (Sec. 4.2).

Fig. 2 (a) illustrates the visual difference between the full-rank velocity u and the asymmetric velocity u_A . Full-rank velocity is perturbed by dense noise, making it highly unpredictable. In contrast, AsymFlow keeps the structured data term full-dimensional but restricts only the stochastic noise term to a low-rank subspace. Since image data itself concentrates near a low-dimensional manifold, this makes the overall asymmetric target more predictable for neural networks.

Patch-wise low-rank projection. Following the patch-token representation of DiTs [48], we apply low-rank projection independently within each image patch. Concretely, for a patch dimension D and rank $r < D$, the matrix $\mathbf{A} \in \mathbb{R}^{D \times r}$ defines a low-rank subspace for each patch token, and the same projector $\mathbf{P} = \mathbf{A} \mathbf{A}^T$ is shared across all tokens. Thus, AsymFlow reduces the noise prediction dimension within each patch while preserving the full set of image tokens.

Choosing the low-rank subspace. When training AsymFlow from scratch, \mathbf{A} can be obtained from a data-dependent patch basis, e.g., by applying PCA to image patches. When adapting a pretrained latent model, \mathbf{A} is instead chosen to align the latent space with the pixel patch space, which we compute by a Procrustes alignment between latent variables and their corresponding pixel patches. This latter construction enables a seamless latent-to-pixel initialization, and is discussed in Sec. 5.

4.2 Orthogonal Component View and Full-Rank Velocity Recovery

The asymmetric velocity in Eq. (3) has a simple interpretation after decomposing it into the low-rank subspace $\text{Im}(\mathbf{P})$ and its orthogonal complement $\text{Im}(\mathbf{I} - \mathbf{P})$:

$$P u_A = P\epsilon - P x_0 = P u, \quad (\mathbf{I} - \mathbf{P}) u_A = -(\mathbf{I} - \mathbf{P}) x_0. \quad (4)$$

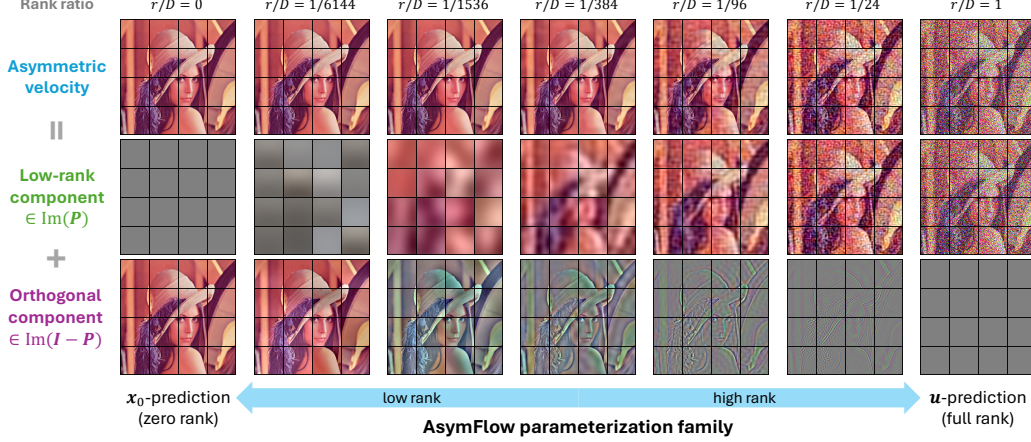


Figure 3: **Orthogonal component view of AsymFlow.** AsymFlow parameterization can be decomposed into a Pu component in the low-rank subspace $\text{Im}(P)$ and an $(I - P)x_0$ component in the orthogonal complement $\text{Im}(I - P)$. Varying the rank r yields a parameterization family whose endpoints recover full x_0 -prediction and full u -prediction.

The decomposition reveals that AsymFlow behaves like u -prediction in the low-rank subspace and like x_0 -prediction in the orthogonal complement. Adjusting the rank r creates a family of parameterizations between the two endpoints, as shown in Fig. 3: when $r = 0$, the target reduces to full x_0 -prediction up to sign; when $r = D$, AsymFlow recovers full u -prediction. We expect a small but nonzero rank r to be optimal: it retains the benefit of u -prediction for controlling the flow on a low-dimensional subspace, while avoiding the burden of predicting full-rank noise.

This component view also provides the conversion back to the full-rank velocity. We keep the low-rank velocity component Pu_A , and convert the orthogonal x_0 -style component to velocity using the x_0 -to- u relation established in Eq. (1):

$$u = Pu_A + (I - P) \frac{x_t + u_A}{\sigma_t}. \quad (5)$$

In practice, we apply the conversion to the network prediction \hat{u}_A to obtain \hat{u} , which is used in the flow matching loss (Eq. (2)) and denoising sampling. Fig. 2 (b) illustrates this conversion visually.

5 Finetuning Latent Flow into Pixel AsymFlow

A key advantage of AsymFlow is that it provides a direct way to turn pretrained u -predicting latent flow models into pixel-space generators. We first lift a pretrained latent model into an equivalent low-rank pixel flow at initialization, with exact input and output conversions between latents and low-rank pixels. Solving this lifted pixel flow ODE preserves the latent trajectory up to an analytically determined orthogonal noise component, so the initialized model generates lifted low-rank pixels whose semantics and structure match the pretrained latent model. Finetuning then focuses on correcting the low-level projection gap between these low-rank pixels and the full-rank pixel targets.

5.1 Latent-to-Pixel Initialization

We consider a latent flow model $\hat{u}_z = G_\phi(z_t, t)$ pretrained on latent tokens $z_0 \in \mathbb{R}^d$ with velocity $u_z := \epsilon_z - z_0$. To bridge the latent-to-pixel gap, we construct a patch-wise linear lift $A \in \mathbb{R}^{D \times d}$ from latent space to pixel space using Procrustes alignment (details in Appendix A.1), such that the lifted low-rank pixels $x_0^L := Az_0$ approximate the full-rank pixels x_0 . Consider the corresponding pixel-space forward process $x_t^L := \alpha_t x_0^L + \sigma_t \epsilon$ and velocity $u^L := \epsilon - x_0^L$. Then the latent and pixel quantities are related by exact input and output conversions:

$$\text{input: } z_t = A^T x_t^L, \quad \text{output: } u^L = PAu_z + (I - P) \frac{x_t^L + Au_z}{\sigma_t}. \quad (6)$$

The input identity shows that noisy low-rank pixels can be projected to noisy latents by \mathbf{A}^\top , while the output identity converts the lifted latent velocity $\mathbf{A}\mathbf{u}_z$ back to the low-rank pixel velocity using the same recovery rule as AsymFlow in Eq. (5). These identities imply trajectory coupling of the lifted pixel and latent ODEs (Theorem 1). Therefore, a d -dimensional latent \mathbf{u} -prediction model can be reinterpreted as an exact rank- d pixel flow model with the network $\mathbf{A}G_\phi(\mathbf{A}^\top\mathbf{x}_t^L, t)$. In implementation, the projections \mathbf{A}^\top and \mathbf{A} are fused into the learnable input and output linear layers of G_ϕ , yielding the initialized pixel AsymFlow model $\hat{\mathbf{u}}_A = G_\theta(\mathbf{x}_t, t)$ for later finetuning.

Initialization property. The initialized low-rank pixel model predicts a target of the form $\mathbf{P}\epsilon - \mathbf{x}_0^L$, so its gap to the AsymFlow target \mathbf{u}_A (Eq. (3)) is only the approximation gap $\mathbf{x}_0 - \mathbf{x}_0^L$. Due to the trajectory coupling (Theorem 1), sampling the initialized model generates \mathbf{x}_0^L -like lifted low-rank pixel samples without accumulating additional trajectory errors. These samples are semantically and structurally aligned with the \mathbf{x}_0 -like decoded latent samples, so the gap $\mathbf{x}_0 - \mathbf{x}_0^L$ is mainly low-level and easy to correct during finetuning, as shown in Fig. 4.

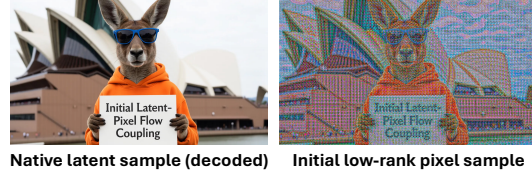


Figure 4: **Latent-to-pixel initialization.** The lifted low-rank pixel generation are semantically and structurally aligned with the decoded latent generation, leaving only a low-level gap to correct.

Scale calibration. A good initialization requires the scale of the lifted pixels \mathbf{x}_0^L to align with the scale of real pixels \mathbf{x}_0 . However, under the orthonormality constraint $\mathbf{A}^\top\mathbf{A} = \mathbf{I}$, Procrustes alignment matches directions but not scale. We therefore introduce a scale factor s and use the scale-calibrated lift $\mathbf{x}_0^L = s\mathbf{A}\mathbf{z}_0$. In implementation, this scale correction is folded into the model input, output, and internal timestep calibration, as detailed in Appendix A.2.

5.2 Variance-Reduced Finetuning Loss

The initialization above reduces latent-to-pixel finetuning to correcting the paired low-level gap $\mathbf{x}_0 - \mathbf{x}_0^L$. While the standard flow matching loss (Eq. (2)) regressing to \mathbf{x}_0 already provides a valid objective, the paired low-rank target \mathbf{x}_0^L offers additional structure that can be used for variance reduction using control variates, thereby improving convergence and sample quality [68].

To achieve this, we inject a term $-\lambda(\mathbf{x}_0^L - \mathbb{E}[\mathbf{x}_0^L|\mathbf{x}_t])$ into Eq. (2). This gives an equivalent flow matching loss whose variance is lower when $\|\mathbf{x}_0 - \mathbf{x}_0^L\|$ is small. The conditional mean $\mathbb{E}[\mathbf{x}_0^L|\mathbf{x}_t]$ can then be approximated by the prediction $\hat{\mathbf{x}}_0^L$ of a frozen copy of the initialized low-rank model:

$$\mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}_0 - \lambda(\mathbf{x}_0^L - \mathbb{E}[\mathbf{x}_0^L|\mathbf{x}_t])\|^2}{\sigma_t^2} \right] \approx \mathbb{E}_{t, \mathbf{x}_0, \epsilon} \left[\frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}_0 - \lambda(\mathbf{x}_0^L - \hat{\mathbf{x}}_0^L)\|^2}{\sigma_t^2} \right] =: \mathcal{L}_{\text{VR}}. \quad (7)$$

Here, $\hat{\mathbf{x}}_0$ is predicted by the finetuned AsymFlow model from \mathbf{x}_t (converted to the \mathbf{x}_0 format), and $\hat{\mathbf{x}}_0^L$ is predicted by the frozen low-rank model from the paired noisy low-rank sample $\mathbf{x}_t^L = \alpha_t\mathbf{x}_0^L + \sigma_t\epsilon$, diffused with the same noise as \mathbf{x}_t . The parameter λ is a patch-wise adaptive weight chosen to minimize the loss gradient norm, thereby reducing the variance of the effective target. In practice, this is implemented via an orthogonal projection and detailed in Appendix A.3. Empirically, the resulting variance-reduced objective \mathcal{L}_{VR} substantially improves fine-grained details in the generated results.

Perceptual correction. The approximation in Eq. (7) assumes $\mathbb{E}[\mathbf{x}_0^L|\mathbf{x}_t] \approx \mathbb{E}[\mathbf{x}_0^L|\mathbf{x}_t^L]$, which is only exact if $\mathbf{x}_t - \mathbf{x}_t^L \in \text{Im}(\mathbf{I} - \mathbf{P})$. In practice, this condition is rarely strictly satisfied when $t < 1$, meaning the variance reduction term $\lambda(\mathbf{x}_0^L - \hat{\mathbf{x}}_0^L)$ introduces a bounded approximation error inside the low-rank subspace $\text{Im}(\mathbf{P})$. Empirically, this manifests as excessive noise in the generated results. To compensate, we add an LPIPS perceptual loss [46, 73] between \mathbf{x}_0 and $\hat{\mathbf{x}}_0$. This perceptual loss is gated by the same patch-wise weight λ , and we dynamically fade from the variance reduction term to the LPIPS loss across diffusion time. We defer the exact weighting schedule to Appendix A.4.

6 Experiments

We evaluate AsymFlow in two settings: ImageNet pixel models trained from scratch with the JiT-H/16 network, which isolate the parameterization itself, and large text-to-image models finetuned from the FLUX.2 klein latent generator, which test the finetuning approach and scalability of AsymFlow.

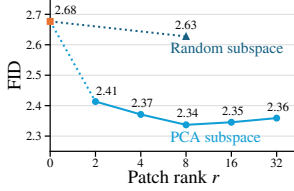


Figure 5: **Patch rank and PCA ablation.** 160 epochs.

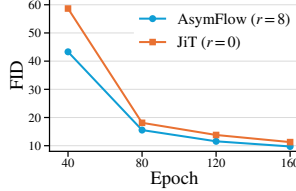


Figure 6: **Convergence speed comparison.** Unguided FIDs.

Table 1: **AsymFlow vs. JiT-H/16** and sensitivity to σ_{\min} clamping. 600 epochs (final checkpoint).

Method	σ_{\min}	FID	IS
AsymFlow ($r = 8$)	0.04	1.76	312.0
	0.00	2.28	306.2
JiT ($r = 0$)	0.04	1.90	300.8
	0.00	3.27	286.7

6.1 Training from Scratch on ImageNet

We train class-conditional ImageNet 256×256 pixel models using the same setup as JiT-H/16 (see Table 9 in [35]), changing only the prediction parameterization. Unless otherwise stated, AsymFlow is trained using the flow matching loss (Eq. (2)) using a $D = 768$ patch-wise PCA subspace of rank r , with $r = 0$ exactly reproducing JiT’s x_0 -prediction. Results use ADM evaluation [14, 19] with grid-searched guidance scales and intervals that optimize FID [20, 33]. We defer the details to Appendix B.

Comparison with JiT baseline. Table 1 compares AsymFlow ($r = 8$) and the official JiT checkpoint using ADM evaluation after 600 epochs. In practical sampling, the x_0 -to- u conversion in Eq. (1) clamps the denominator by σ_{\min} to avoid numerical instability [35]. Since AsymFlow applies this conversion only in the orthogonal complement, it should be less sensitive to this clamp. The results confirm this: with the optimal $\sigma_{\min} = 0.04$ for both methods, AsymFlow improves over JiT in both FID and IS by a clear margin; disabling clamping degrades JiT by 1.37 FID, but AsymFlow by only 0.52. This shows that the asymmetric parameterization improves both overall quality and low-noise numerical stability.

Patch rank. Figure 5 studies the effect of the patch rank. Moving from JiT ($r = 0$) to AsymFlow sharply improves guided FID, with the best result at $r = 8$; increasing the rank further gives mild degradation. This matches the intended trade-off: AsymFlow keeps velocity prediction in a useful low-rank subspace while avoiding the burden of predicting high-dimensional noise.

PCA subspace. Figure 5 also compares PCA and random subspaces at $r = 8$. The random subspace performs close to the JiT baseline and far worse than PCA, showing that the gain comes from using a meaningful low-rank subspace, not merely reducing rank.

Convergence speed. Figure 6 compares FID during training. With the same architecture and recipe, AsymFlow ($r = 8$) consistently improves over JiT and reaches comparable FID roughly 40% faster. Thus, the rank-asymmetric target improves not only final quality but also optimization efficiency.

Comparison with prior pixel diffusion models. Table 2 compares AsymFlow ($r = 8$ plus a standard REPA loss [70]) with prior ImageNet 256×256 pixel diffusion models. With REPA, AsymFlow reaches 1.57 FID, establishing the state of the art among practical pixel diffusion models (excluding the much more expensive SiD2 UViT/1). In particular, AsymFlow outperforms previous plain-transformer models by a large margin (FID 1.57 vs. 1.81*). This result also shows that AsymFlow is strongly compatible with REPA: PixelREPA [58] reports that plain REPA is ineffective for larger JiT models, and its additional designs improve JiT-H/16 only from 1.86* to 1.81* FID; in contrast, adding plain REPA to AsymFlow improves FID from 1.76 to 1.57, suggesting that the AsymFlow parameterization is much more robust to auxiliary losses and can better leverage their benefits.

Table 2: **ImageNet 256×256 pixel diffusion comparison.** FLOP estimation follows the convention in [71]. * denotes JiT evaluation protocol, which may have up to 0.08 better FID than ADM according to our tests.

Method	Pred (\pm)	Params	GFLOPs	FID↓
<i>Hierarchical CNNs (skip connections / U-Net-like)</i>				
ADM-G [14]	ϵ	554M	2240	4.59
<i>Hierarchical transformers (skip connections / U-ViT-like)</i>				
RIN [26]	ϵ	320M	668	3.42
SiD, UViT/2 [22]	ϵ	2B	1110	2.44
VDM++, UViT/2 [31]	ϵ	2B	1110	2.12
SiD2, UViT/2 [23]	ϵ	-	274	1.73
EPG-G/16 [34]	x_0	1.4B	642	1.58
SiD2, UViT/1 [23]	ϵ	-	1306	1.38
<i>Hierarchical transformers (decoder head / DDT-like)</i>				
PixNerd-XL/16 [64]	$\epsilon - x_0$	700M	268	2.15
DiP-XL/16 [10]	$\epsilon - x_0$	631M	-	1.79
DeCo-XL/16 [45]	$\epsilon - x_0$	682M	245	1.62
PixelDiT-XL/16 [71]	$\epsilon - x_0$	797M	311	1.61
<i>Plain transformers (DiT-like)</i>				
PixelFlow-XL/4 [9]	$\epsilon - x_0$	677M	5818	1.98
JiT-H/16 [35]	x_0	953M	363	1.86*
PixelGen-XL/16 [46]	x_0	676M	260	1.83
JiT-G/16 [35]	x_0	2B	766	1.82*
PixelREPA-H/16 [58]	x_0	953M	363	1.81*
AsymFlow-H/16	$P\epsilon - x_0$	953M	363	1.57

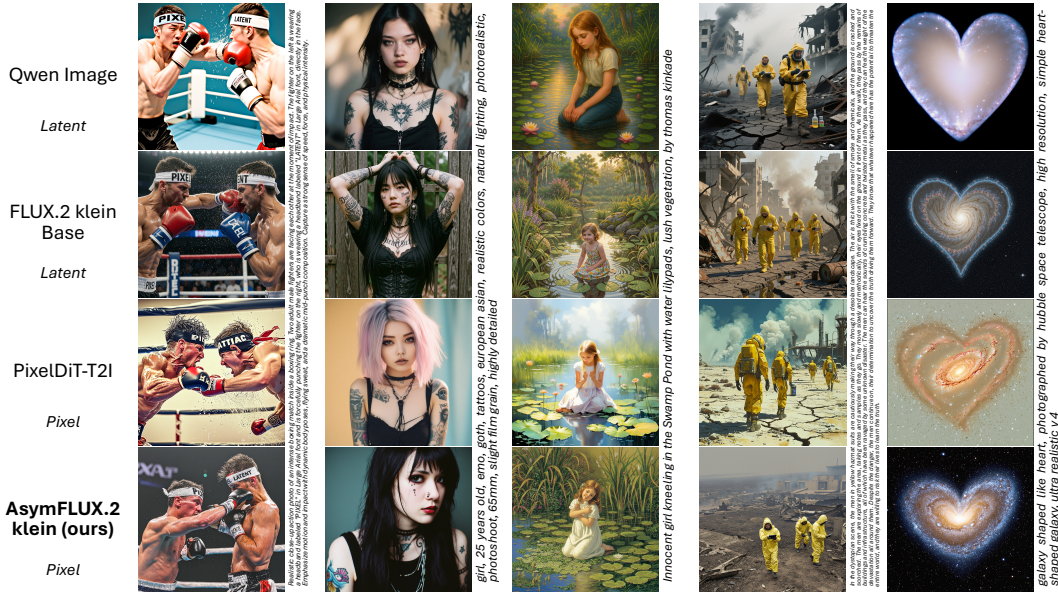


Figure 7: **Qualitative comparison of T2I diffusion models.** AsymFLUX.2 klein produces more realistic images with richer visual styles than prior models. More results are shown in Fig. 9 and 10.

Table 3: **Comparison with baselines and ablation studies.** All models are finetuned on the LAION-Aesthetics dataset [56] for 10K iterations, and evaluated on the COCO-10K dataset [38].

Method	HPSv3 \uparrow	HPSv2.1 \uparrow	VQA \uparrow	CLIP \uparrow	FID \downarrow	pFID \downarrow
FLUX.2 klein Base + latent finetune	10.70	0.290	0.936	0.276	15.0	18.8
FLUX.2 klein Base + DDT finetune	10.33	0.291	0.922	0.273	20.4	26.0
AsymFLUX.2 klein (standard FM)	12.03	0.293	0.922	0.277	20.2	25.4
AsymFLUX.2 klein (variance reduction)	<u>12.99</u>	<u>0.296</u>	<u>0.925</u>	0.280	<u>18.5</u>	27.8
+ perceptual correction	13.06	0.297	<u>0.925</u>	<u>0.278</u>	19.1	<u>22.5</u>

6.2 Finetuning Large Text-to-Image Models

For text-to-image generation, we finetune the pretrained FLUX.2 klein Base 9B latent flow model [6] (patch dimension $d = 128$) into a pixel-space AsymFlow model. We call the resulting model AsymFLUX.2 klein. The model is finetuned on 3M LAION-Aesthetics images [56], resized to one-megapixel resolution and captioned with Qwen2.5-VL [3]. To reduce overfitting, we freeze the base model and finetune only the input/output projection layers together with rank-256 LoRA adapters [24]. Sampling uses UniPC [74] with APG orthogonal-projection guidance [53]. We defer additional details to Appendix B.

Evaluation protocol. All text-to-image evaluations generate 1024×1024 images. For system-level comparison, we use three benchmarks: HPSv3 [44] measures human preference, which combines realism, style, and overall prompt following, while DPG-Bench [25] and GenEval [16] focus more on fine-grained entities, attributes, relations, counting, and composition. For controlled ablations, we generate images using 10K captions from the COCO 2014 validation set [37, 38] and report preference metrics HPSv3 [44] and HPSv2.1 [67], prompt-alignment metrics VQAScore [39] and CLIP score [50], and distribution metrics FID [19] and patch FID (pFID) [37].

System-level comparison. Table 4 compares AsymFLUX.2 klein (with variance reduction and perceptual correction) with prior latent and pixel text-to-image diffusion models. AsymFLUX.2 klein improves over its FLUX.2 klein latent base on all three benchmarks, with the largest gain

Table 4: **System-level comparison of text-to-image (1024×1024) diffusion models.**

Method	HPSv3 \uparrow	DPG \uparrow	GenEval \uparrow
<i>Latent diffusion models</i>			
SDXL [49]	8.20	74.7	0.55
PixArt- Σ [8]	<u>9.37</u>	80.5	0.54
Hunyuan-DiT [36]	8.19	78.9	0.63
FLUX.1 dev [5]	10.43	84.0	0.67
Qwen-Image [66]	9.52	87.8	0.86
FLUX.2 klein Base [6]	9.50	<u>85.2</u>	<u>0.80</u>
<i>Pixel diffusion models</i>			
PixelDiT-T2I [71]	<u>8.95</u>	83.5	<u>0.74</u>
AsymFLUX.2 klein	10.66	86.8	0.82

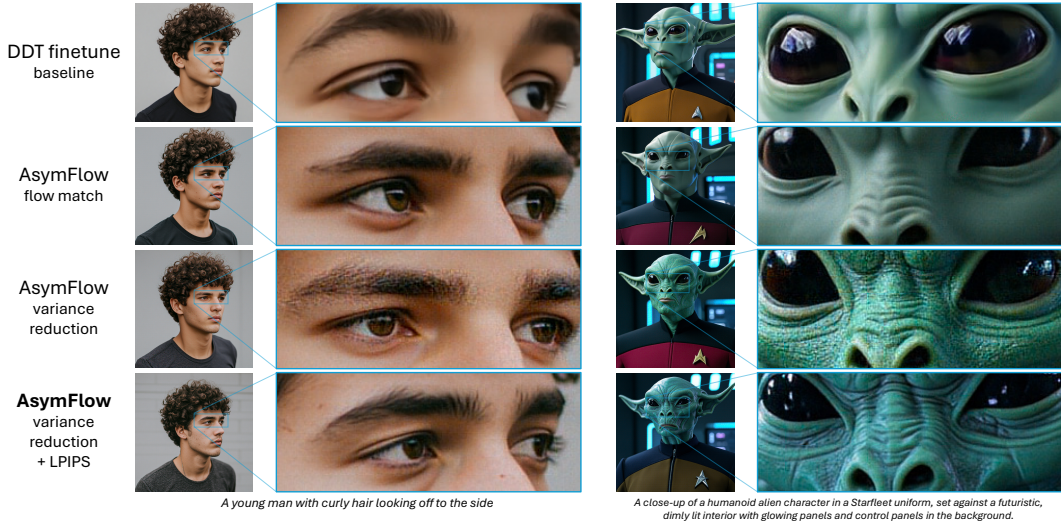


Figure 8: **Ablation of AsymFLUX.2 klein finetuning.** AsymFlow produces finer details than the DDT baseline. Variance reduction further improves details and texture but introduces excessive noise. The LPIPS perceptual correction suppresses this artifact while preserving the sharp appearance.

on HPSv3, indicating a substantial improvement in human-aligned visual quality. Consequently, it outperforms the prior pixel model PixelDiT-T2I [71] by a large margin across all metrics, establishing a new state of the art for pixel-space text-to-image generation. Figure 7 shows the same trend qualitatively: AsymFLUX.2 klein produces realistic and diverse visual styles with stronger texture, while popular latent models such as Qwen Image [3] and FLUX.2 klein Base [6] still have a more artificial appearance; compared to PixelDiT-T2I, AsymFLUX.2 klein recovers much sharper details in addition to other qualitative improvements, marking a significant step forward for pixel-space text-to-image generation.

Controlled baselines. To separate dataset effects from latent-to-pixel conversion, we include a latent-finetuned FLUX.2 klein baseline trained on the same data. We also include a u -prediction pixel finetuning baseline with a DDT decoder head [65, 75], similar in spirit to PixelDiT [71]. The results are presented in Table 3: compared to the latent baseline, finetuned AsymFLUX.2 klein models yield clear improvements in HPSv3 and HPSv2.1, indicating that the improved overall quality comes from AsymFlow pixel-space conversion instead of dataset bias. In contrast, the DDT baseline falls behind in all metrics, despite having more parameters and capacity. This is also reflected in the qualitative comparison in Figure 8, where the DDT baseline produces blurry images and exhibits minor patch seams, while AsymFLUX.2 klein recovers sharper details and more realistic texture.

Loss ablations. The results in Table 3 also validate the effectiveness of variance reduction and perceptual correction losses: variance reduction boosts all metrics except pFID, due to its low-noise approximation error that introduces excessive noise (Figure 8). This is directly addressed by the LPIPS perceptual correction loss, which significantly improves pFID and HPS scores, resulting in the most natural and realistic texture in Figure 8.

7 Conclusion

We introduced AsymFlow, a rank-asymmetric flow velocity parameterization that enables high-dimensional pixel-space generation with plain diffusion transformers. When trained from scratch, this single parameterization yields a leading 1.57 FID among ImageNet pixel diffusion models. It also provides the first path for finetuning pretrained large latent flow models into pixel generators with improved visual fidelity, demonstrating AsymFlow’s scalability and practical impact. This opens promising directions for high-fidelity image and video generation with finer low-level control, as well as other high-dimensional data modalities previously out of reach for flow-based modeling.

Limitations. Latent-to-pixel finetuning assumes a good patch-level linear lift. It may not work well when the pretrained latent space does not preserve pixel structure, such as in RAE models [75].

References

- [1] Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *ICLR*, 2023.
- [2] Alan Baade, Eric Ryan Chan, Kyle Sargent, Changan Chen, Justin Johnson, Ehsan Adeli, and Li Fei-Fei. Latent forcing: Reordering the diffusion trajectory for pixel-space image generation. *arXiv preprint arXiv:2602.11401*, 2026.
- [3] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng, Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025. URL <https://arxiv.org/abs/2502.13923>.
- [4] Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth words: A ViT backbone for diffusion models. In *CVPR*, 2023.
- [5] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- [6] Black Forest Labs. Flux.2: Frontier visual intelligence. <https://bf1.ai/blog/flux-2>, 2025.
- [7] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. <https://openai.com/research/video-generation-models-as-world-simulators>, 2024.
- [8] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. PixArt- Σ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In *ECCV*, page 74–91, Berlin, Heidelberg, 2024. Springer-Verlag. ISBN 978-3-031-73410-6. doi: 10.1007/978-3-031-73411-3_5. URL https://doi.org/10.1007/978-3-031-73411-3_5.
- [9] Shoufa Chen, Chongjian Ge, Shilong Zhang, Peize Sun, and Ping Luo. Pixelflow: Pixel-space generative models with flow. *arXiv preprint arXiv:2504.07963*, 2025.
- [10] Zhennan Chen, Junwei Zhu, Xu Chen, Jiangning Zhang, Xiaobin Hu, Hanzhen Zhao, Chengjie Wang, Jian Yang, and Ying Tai. Dip: Taming diffusion models in pixel space. In *CVPR*, 2026.
- [11] Katherine Crowson, Stefan Andreas Baumann, Alex Birch, Tanishq Mathew Abraham, Daniel Z Kaplan, and Enrico Shippole. Scalable high-resolution pixel-space image synthesis with hourglass diffusion transformers. In *ICML*, 2024.
- [12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- [13] Tim Dettmers, Mike Lewis, Sam Shleifer, and Luke Zettlemoyer. 8-bit optimizers via block-wise quantization. In *ICLR*, 2022.
- [14] Prafulla Dhariwal and Alexander Quinn Nichol. Diffusion models beat GANs on image synthesis. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *NeurIPS*, 2021. URL <https://openreview.net/forum?id=AAWuVzavt>.
- [15] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In *ICML*, 2024.
- [16] Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: an object-focused framework for evaluating text-to-image alignment. In *NeurIPS*, Red Hook, NY, USA, 2023. Curran Associates Inc.

- [17] Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Joshua M Susskind, and Navdeep Jaitly. Matryoshka diffusion models. In *ICLR*, 2023.
- [18] Yoav HaCohen, Nisan Chiprut, Benny Brazowski, Daniel Shalem, Dudu Moshe, Eitan Richardson, Eran Levin, Guy Shiran, Nir Zabari, Ori Gordon, Poriya Panet, Sapir Weissbuch, Victor Kulikov, Yaki Bitterman, Zeev Melumian, and Ofir Bibi. Ltx-video: Realtime video latent diffusion. *arXiv preprint arXiv:2501.00103*, 2024. URL <https://arxiv.org/abs/2501.00103>.
- [19] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017.
- [20] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS Workshop*, 2021.
- [21] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020.
- [22] Emiel Hoogeboom, Jonathan Heek, and Tim Salimans. Simple diffusion: End-to-end diffusion for high resolution images. In *ICML*, pages 13213–13232, 2023.
- [23] Emiel Hoogeboom, Thomas Mensink, Jonathan Heek, Kay Lamerigts, Ruiqi Gao, and Tim Salimans. Simpler diffusion (sid2): 1.5 fid on imagenet512 with pixel-space diffusion. In *CVPR*, 2025.
- [24] Edward J Hu, yelong shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *ICLR*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- [25] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. URL <https://arxiv.org/abs/2403.05135>.
- [26] Allan Jabri, David Fleet, and Ting Chen. Scalable adaptive computation for iterative generation. In *ICML*, 2023.
- [27] Qing Jin and Chaoyang Wang. Revisiting diffusion model predictions through dimensionality. *arXiv preprint arXiv:2601.21419*, 2026.
- [28] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022.
- [29] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *CVPR*, 2024.
- [30] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2014.
- [31] Diederik P Kingma and Ruiqi Gao. Understanding diffusion objectives as the ELBO with simple data augmentation. In *NeurIPS*, 2023. URL <https://openreview.net/forum?id=NnMEadcdyD>.
- [32] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, Kathrina Wu, Qin Lin, Junkun Yuan, Yanxin Long, Aladdin Wang, Andong Wang, Changlin Li, DuoJun Huang, Fang Yang, Hao Tan, Hongmei Wang, Jacob Song, Jiawang Bai, Jianbing Wu, Jinbao Xue, Joey Wang, Kai Wang, Mengyang Liu, Pengyu Li, Shuai Li, Weiyan Wang, Wenqing Yu, Xincheng Deng, Yang Li, Yi Chen, Yutao Cui, Yuanbo Peng, Zhentao Yu, Zhiyu He, Zhiyong Xu, Zixiang Zhou, Zunnan Xu, Yangyu Tao, Qinglin Lu, Songtao Liu, Daquan Zhou, Hongfa Wang, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, and Caesar Zhong. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2025. URL <https://arxiv.org/abs/2412.03603>.
- [33] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. In *NeurIPS*, 2024.

- [34] Jiachen Lei, Keli Liu, Julius Berner, Y HoiM, Hongkai Zheng, Jiahong Wu, and Xiangxiang Chu. There is no VAE: End-to-end pixel-space generative modeling via self-supervised pre-training. In *ICLR*, 2026. URL <https://openreview.net/forum?id=HbUoKPIZmp>.
- [35] Tianhong Li and Kaiming He. Back to basics: Let denoising generative models denoise. In *CVPR*, 2026.
- [36] Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang Zhang, Xingchao Liu, Minbin Huang, Zedong Xiao, Dayou Chen, Jiajun He, Jiahao Li, Wenyue Li, Chen Zhang, Rongwei Quan, Jianxiang Lu, Jiabin Huang, Xiaoyan Yuan, Xiaoxiao Zheng, Yixuan Li, Jihong Zhang, Chao Zhang, Meng Chen, Jie Liu, Zheng Fang, Weiyang Wang, Jinbao Xue, Yangyu Tao, Jianchen Zhu, Kai Liu, Sihuan Lin, Yifu Sun, Yun Li, Dongdong Wang, Mingtao Chen, Zhichao Hu, Xiao Xiao, Yan Chen, Yuhong Liu, Wei Liu, Di Wang, Yong Yang, Jie Jiang, and Qinglin Lu. Hunyuan-dit: A powerful multi-resolution diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*, 2024. URL <https://arxiv.org/abs/2405.08748>.
- [37] Shanchuan Lin, Anran Wang, and Xiao Yang. Sdxl-lightning: Progressive adversarial diffusion distillation. *arXiv preprint arXiv:2402.13929*, 2024. URL <https://arxiv.org/abs/2402.13929>.
- [38] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In David Fleet, Tomas Pajdla, Bernt Schiele, and Tinne Tuytelaars, editors, *ECCV*, pages 740–755, Cham, 2014. Springer International Publishing. ISBN 978-3-319-10602-1.
- [39] Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *ECCV*, 2024.
- [40] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *ICLR*, 2023. URL <https://openreview.net/forum?id=PqvMRDCJT9t>.
- [41] Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022.
- [42] Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *ICLR*, 2023. URL <https://openreview.net/forum?id=XVjTT1nw5z>.
- [43] Nanye Ma, Mark Goldstein, Michael S. Albergo, Nicholas M. Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. In *ECCV*, 2024.
- [44] Yuhang Ma, Xiaoshi Wu, Keqiang Sun, and Hongsheng Li. Hpsv3: Towards wide-spectrum human preference score. In *ICCV*, 2025.
- [45] Zehong Ma, Longhui Wei, Shuai Wang, Shiliang Zhang, and Qi Tian. Deco: Frequency-decoupled pixel diffusion for end-to-end image generation. In *CVPR*, 2026.
- [46] Zehong Ma, Ruihan Xu, and Shiliang Zhang. Pixelgen: Pixel diffusion beats latent diffusion with perceptual loss. *arXiv preprint arXiv:2602.02493*, 2026.
- [47] Björn Ottosson. A perceptual color space for image processing, 2020. URL <https://bottosson.github.io/posts/oklab/>.
- [48] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, 2023.
- [49] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. In *ICLR*, 2024. URL <https://openreview.net/forum?id=di52zR8xgf>.

- [50] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021.
- [51] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022.
- [52] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 234–241, 2015.
- [53] Seyedmorteza Sadat, Otmar Hilliges, and Romann M. Weber. Eliminating oversaturation and artifacts of high guidance scales in diffusion models. In *ICLR*, 2025.
- [54] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022.
- [55] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *ICLR*, 2022.
- [56] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS Datasets and Benchmarks*, 2022. URL <https://openreview.net/forum?id=M3Y74vmsMcY>.
- [57] Peter H. Schönemann. A generalized solution of the orthogonal procrustes problem. *Psychometrika*, 31(1):1–10, 1966. doi: 10.1007/BF02289451.
- [58] Jaeyo Shin, Jiwook Kim, and Hyunjung Shim. Representation alignment for just image transformers is not easier than you think. *arXiv preprint arXiv:2603.14366*, 2026.
- [59] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, pages 2256–2265, 2015.
- [60] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019.
- [61] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.
- [62] Shengbang Tong, Boyang Zheng, Ziteng Wang, Bingda Tang, Nanye Ma, Ellis Brown, Jihan Yang, Rob Fergus, Yann LeCun, and Saining Xie. Scaling text-to-image diffusion transformers with representation autoencoders. *arXiv preprint arXiv:2601.16208*, 2026.
- [63] Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, Jianyuan Zeng, Jiayu Wang, Jingfeng Zhang, Jingren Zhou, Jinkai Wang, Jixuan Chen, Kai Zhu, Kang Zhao, Keyu Yan, Lianghua Huang, Mengyang Feng, Ningyi Zhang, Pandeng Li, Pingyu Wu, Ruihang Chu, Ruili Feng, Shiwei Zhang, Siyang Sun, Tao Fang, Tianxing Wang, Tianyi Gui, Tingyu Weng, Tong Shen, Wei Lin, Wei Wang, Wei Wang, Wenmeng Zhou, Wenten Wang, Wenting Shen, Wenyuan Yu, Xianzhong Shi, Xiaoming Huang, Xin Xu, Yan Kou, Yangyu Lv, Yifei Li, Yijing Liu, Yiming Wang, Yingya Zhang, Yitong Huang, Yong Li, You Wu, Yu Liu, Yulin Pan, Yun Zheng, Yuntao Hong, Yupeng Shi, Yutong Feng, Zeyinzi Jiang, Zhen Han, Zhi-Fan Wu, and Ziyu Liu. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. URL <https://arxiv.org/abs/2503.20314>.
- [64] Shuai Wang, Ziteng Gao, Chenhui Zhu, Weilin Huang, and Limin Wang. Pixnerd: Pixel neural field diffusion. In *ICLR*, 2026. URL <https://openreview.net/forum?id=BDnOrExHmt>.

- [65] Shuai Wang, Zhi Tian, Weilin Huang, and Limin Wang. Ddt: Decoupled diffusion transformer. In *CVPR*, 2026.
- [66] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan Cai, and Zenan Liu. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. URL <https://arxiv.org/abs/2508.02324>.
- [67] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. URL <https://arxiv.org/abs/2306.09341>.
- [68] Yilun Xu, Shangyuan Tong, and Tommi S. Jaakkola. Stable target field for reduced variance score estimation in diffusion models. In *ICLR*, 2023. URL <https://openreview.net/forum?id=WmIwYTdOYTF>.
- [69] Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. In *CVPR*, 2025.
- [70] Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. In *ICLR*, 2025.
- [71] Yongsheng Yu, Wei Xiong, Weili Nie, Yichen Sheng, Shiqiu Liu, and Jiebo Luo. Pixeldit: Pixel diffusion transformers for image generation. In *CVPR*, 2026.
- [72] Z-Image Team, Huanqia Cai, Sihan Cao, Ruoyi Du, Peng Gao, Steven Hoi, Zhaohui Hou, Shijie Huang, Dengyang Jiang, Xin Jin, Liangchen Li, Zhen Li, Zhong-Yu Li, David Liu, Dongyang Liu, Junhan Shi, Qilong Wu, Feng Yu, Chi Zhang, Shifeng Zhang, and Shilin Zhou. Z-image: An efficient image generation foundation model with single-stream diffusion transformer. *arXiv preprint arXiv:2511.22699*, 2025. URL <https://arxiv.org/abs/2511.22699>.
- [73] Richard Zhang, Phillip Isola, Alexei Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018.
- [74] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. In *NeurIPS*, 2023.
- [75] Boyang Zheng, Nanye Ma, Shengbang Tong, and Saining Xie. Diffusion transformers with representation autoencoders. In *ICLR*, 2026. URL <https://openreview.net/forum?id=Ou1LigJaab>.

A Method Details

A.1 Low-Rank Subspace Construction

For transformer-based pixel generation, AsymFlow requires a patch-wise low-rank subspace. We use two constructions, depending on whether the model is trained from scratch or initialized from a latent model.

Orthonormality requirement. In both cases we require the columns of \mathbf{A} to be orthonormal. This ensures that projecting standard pixel-space Gaussian noise preserves its Gaussian form inside the low-rank coordinates: if $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_D)$ and $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_r$, then $\mathbf{A}^\top \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_r)$.

PCA basis for from-scratch training. Ideally, the low-rank directions would preserve the most perceptually important information in each image patch. When training from scratch, PCA gives a practical proxy by retaining the dominant patch variations without introducing an additional learned representation. Let $\mathbf{X} \in \mathbb{R}^{D \times N}$ collect N image patches with normalized pixel values. Taking the top left singular vectors of \mathbf{X} gives the PCA subspace:

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top, \quad \mathbf{A} = \mathbf{U}_r, \quad \mathbf{P} = \mathbf{A}\mathbf{A}^\top. \quad (8)$$

Here \mathbf{U}_r denotes the top r columns of \mathbf{U} . Thus \mathbf{P} keeps the data-adaptive PCA directions and removes the remaining patch-space directions from the noise prediction.

Procrustes basis for latent-to-pixel finetuning. For latent-to-pixel finetuning, the subspace should be aligned with the pretrained latent representation to minimize the paired gap $\|\mathbf{x}_0 - \mathbf{x}_0^L\|$. Let $\mathbf{X} \in \mathbb{R}^{D \times N}$ collect image patches with normalized pixel values and $\mathbf{Z} \in \mathbb{R}^{d \times N}$ collect the corresponding latent tokens. We solve the orthogonal Procrustes problem [57]

$$\mathbf{A}^* = \arg \min_{\mathbf{A} \in \mathbb{R}^{D \times d}, \mathbf{A}^\top \mathbf{A} = \mathbf{I}_d} \|\mathbf{X} - \mathbf{A}\mathbf{Z}\|_{\text{F}}^2. \quad (9)$$

This objective finds an orthonormal lift from latent tokens to pixel patches. Equivalently, it maximizes the inner-product alignment between $\mathbf{A}\mathbf{Z}$ and \mathbf{X} , so $\mathbf{A}^* = \arg \max_{\mathbf{A}^\top \mathbf{A} = \mathbf{I}_d} \text{Tr}(\mathbf{A}^\top \mathbf{X}\mathbf{Z}^\top)$. If $\mathbf{X}\mathbf{Z}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ is the compact SVD, the solution is

$$\mathbf{X}\mathbf{Z}^\top = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top, \quad \mathbf{A}^* = \mathbf{U}\mathbf{V}^\top, \quad \mathbf{P} = \mathbf{A}^*(\mathbf{A}^*)^\top. \quad (10)$$

Procrustes aligns directions under the orthonormality constraint. It does not determine the correct pixel scale, so we apply the scalar calibration below.

A.2 Scale and Timestep Calibration

The Procrustes lift gives a directionally aligned low-rank pixel reconstruction, but its magnitude may not match the pixel scale within the Procrustes subspace. We therefore introduce a scalar s and use the calibrated lift

$$\mathbf{x}_0^L = s\mathbf{A}\mathbf{z}_0, \quad \mathbf{A}^\top \mathbf{A} = \mathbf{I}_d, \quad \mathbf{P} = \mathbf{A}\mathbf{A}^\top. \quad (11)$$

The scalar s is estimated from the same paired latent-token and pixel-patch statistics used above, by matching the Frobenius norm of the latents \mathbf{Z} and the rescaled projected pixels $\mathbf{A}^\top \mathbf{X}/s$:

$$s = \frac{\|\mathbf{A}^\top \mathbf{X}\|_{\text{F}}}{\|\mathbf{Z}\|_{\text{F}}}. \quad (12)$$

Equivalently, the calibrated lift $s\mathbf{A}\mathbf{Z}$ and the low-rank pixels $\mathbf{P}\mathbf{X}$ have the same Frobenius norm.

Scale calibration must also be reflected in noisy inputs, not only in the clean lift. Projecting a noisy pixel state gives signal coefficient $s\alpha_t$ and noise coefficient σ_t , so the latent-space signal-to-noise ratio (SNR) is $s\alpha_t/\sigma_t$. The SNR constraint first determines the latent time τ at which the pretrained model should be evaluated. Under the linear flow schedule, this gives

$$\frac{1 - \tau}{\tau} = \frac{s(1 - t)}{t} \implies \tau = \frac{t}{s(1 - t) + t}. \quad (13)$$

After fixing τ , the projected input must also have the correct noise magnitude $\sigma_\tau = \tau$. This determines the input rescaling

$$k = \frac{\tau}{t} = \frac{1}{s(1 - t) + t}, \quad (14)$$

which places the projected state on the latent trajectory expected by the pretrained model, up to a low-rank approximation error:

$$\mathbf{A}^T(k\mathbf{x}_t) \approx \mathbf{A}^T(k\mathbf{x}_t^L) = \alpha_\tau \mathbf{z}_0 + \sigma_\tau \boldsymbol{\epsilon}_z = \mathbf{z}_\tau. \quad (15)$$

The output conversion must use the same calibration. The network is finetuned to predict the calibrated AsymFlow target

$$\mathbf{u}_A^{\text{cal}} := \mathbf{P}\boldsymbol{\epsilon} - \frac{\mathbf{x}_0}{s}, \quad (16)$$

which is defined in the coordinate system of the rescaled input $k\mathbf{x}_t$. Recovering the original pixel-space full-rank velocity $\mathbf{u} = \boldsymbol{\epsilon} - \mathbf{x}_0$ gives

$$\mathbf{u} = \underbrace{\mathbf{P}\left(sk\mathbf{u}_A^{\text{cal}} + (1-sk)\frac{\mathbf{x}_t}{\sigma_t}\right)}_{\text{low-rank subspace}} + \underbrace{(\mathbf{I} - \mathbf{P})\left(\frac{\mathbf{x}_t + s\mathbf{u}_A^{\text{cal}}}{\sigma_t}\right)}_{\text{orthogonal complement}}. \quad (17)$$

Eq. (17) is a generalized form of the uncalibrated conversion formula in Eq. (5). When $s = 1$ and $k = 1$, it reduces to the uncalibrated formula.

In practice, we apply this generalized conversion to the calibrated network prediction $\hat{\mathbf{u}}_A^{\text{cal}} = G_\theta(k\mathbf{x}_t, kt)$ to obtain $\hat{\mathbf{u}}$, which is used in the flow matching loss (Eq. (2)) and denoising sampling.

A.3 Adaptive Weighting for Variance Reduction

The variance-reduced loss in Eq. (7) uses a patch-wise coefficient λ . For a given patch prediction, λ is determined by directly minimizing the loss residual along the one-dimensional control-variate direction (see Appendix C.3 for mathematical justification). Since the gradient of the squared loss is proportional to the corrected residual, this also minimizes the corresponding gradient norm, effectively selecting the lowest-variance target available along that direction.

The one-dimensional minimization has a closed-form solution given by an orthogonal projection. For each patch, define the low-rank prediction deviation of the frozen low-rank model as $\mathbf{d}^L := \mathbf{x}_0^L - \hat{\mathbf{x}}_0^L$ and the full-rank prediction deviation of the finetuned model as $\mathbf{d} := \mathbf{x}_0 - \text{stopgrad}(\hat{\mathbf{x}}_0)$. The variance-reduced loss residual is then $\mathbf{d} - \lambda\mathbf{d}^L$. Minimizing the patch loss over λ gives the one-dimensional least-squares solution:

$$\lambda^* = \arg \min_{\lambda} \|\mathbf{d} - \lambda\mathbf{d}^L\|^2 = \frac{\langle \mathbf{d}, \mathbf{d}^L \rangle}{\|\mathbf{d}^L\|^2}. \quad (18)$$

Geometrically, this subtracts the component of the full-pixel prediction deviation that lies along the low-rank prediction deviation, leaving the smallest possible loss residual within this one-dimensional family. In practice, we use the clamped coefficient $\lambda = \min(\max(\lambda^*, 0), 1)$.

A.4 Perceptual Correction

The variance-reduced loss in Eq. (7) uses the approximation $\mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t] \approx \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t^L]$, as analyzed in Appendix C.3. This approximation is valid when $\mathbf{x}_t - \mathbf{x}_t^L \in \text{Im}(\mathbf{I} - \mathbf{P})$, which is guaranteed at $t = 1$ because both inputs are pure noise. For $t < 1$, this condition requires $\mathbf{x}_0 - \mathbf{x}_0^L \in \text{Im}(\mathbf{I} - \mathbf{P})$, which generally does not hold, so the variance-reduction term $\lambda(\mathbf{x}_0^L - \hat{\mathbf{x}}_0^L)$ can introduce approximation error in the low-rank subspace $\text{Im}(\mathbf{P})$. Therefore, we need to reduce reliance on this term near the low-noise end of the trajectory.

Simply downweighting the variance-reduction term near low noise is not ideal, because the variance-reduced target is important for learning fine details. To compensate, we introduce a fading schedule $\omega_t \in [0, 1]$ that interpolates from the variance-reduction term to an LPIPS [73] perceptual loss between $\hat{\mathbf{x}}_0$ and \mathbf{x}_0 . The variance-reduction term in Eq. (7) is multiplied by $1 - \omega_t$:

$$\mathcal{L}_{\text{VR}} = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{\|\mathbf{x}_0 - \hat{\mathbf{x}}_0 - (1 - \omega_t)\lambda(\mathbf{x}_0^L - \hat{\mathbf{x}}_0^L)\|^2}{\sigma_t^2} \right], \quad (19)$$

while the complementary perceptual term is multiplied by ω_t :

$$\mathcal{L}_{\text{P}} = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}} \left[\frac{\omega_t \lambda}{\sigma_t^2} \text{LPIPS}(\hat{\mathbf{x}}_0, \mathbf{x}_0) \right]. \quad (20)$$

Here λ is reused only as the patch-wise adaptive gate for the perceptual correction, and $1/\sigma_t^2$ recovers velocity-space weighting.

In our implementation, we define ω_t as a shifted signal-ratio schedule:

$$\omega_t = \frac{\alpha_t^2}{\alpha_t^2 + (\kappa\sigma_t)^2}, \quad (21)$$

where κ is a shift hyperparameter [15] that controls the transition. The final finetuning loss is

$$\mathcal{L} = \mathcal{L}_{\text{VR}} + \omega_{\text{P}}\mathcal{L}_{\text{P}}, \quad (22)$$

where ω_{P} is a hyperparameter that controls the overall weight of the perceptual correction. In our experiments, we use $\kappa = 0.3$ and $\omega_{\text{P}} = 0.2$. We did not perform a systematic hyperparameter sweep due to computational constraints, so there may be room for further improvement.

B Experiment Details

B.1 ImageNet Experiments

For ImageNet 256×256 experiments, we use the same architecture, optimizer, and other training hyperparameters as JiT-H/16 (see Table 9 of JiT [35]). Training for 600 epochs costs approximately 1750 NVIDIA H100 GPU hours. The REPA-enhanced variant follows the standard REPA setting [70]: we apply the REPA loss to the features after the 8th transformer block with loss weight 0.5.

At inference time, we set the velocity-recovery clamp to $\sigma_{\text{min}} = 0.04$, which performs better than the JiT default $\sigma_{\text{min}} = 0.05$ for both the JiT baseline and AsymFlow. Unless otherwise stated, all other inference settings follow JiT exactly, including the 50-step Heun ODE solver, class-balanced sampling, BF16 inference, and attention upcasting.

For each classifier-free guidance (CFG) [20] result, we grid-search the CFG scale with step size 0.1 and the guidance interval with step size 0.02 [33]. Table 5 lists the selected settings for Fig. 5. The final AsymFlow result in Table 1 uses CFG scale 2.3 and interval $[0, 0.88]$, while the REPA-enhanced result in Table 2 uses CFG scale 2.2 and interval $[0, 0.88]$.

Table 5: **Guidance settings for the ImageNet patch-rank sweep.** These settings are selected by grid-searching guided FID for each rank.

Patch rank r	CFG scale	Guidance interval
0	2.7	$[0, 0.82]$
2	2.6	$[0, 0.82]$
4	2.6	$[0, 0.82]$
8	2.5	$[0, 0.82]$
16	2.7	$[0, 0.82]$
32	2.7	$[0, 0.82]$
8 (random subspace)	2.8	$[0, 0.82]$

B.2 Text-to-Image Experiments

For text-to-image experiments, we represent pixels in Oklab color space [47] because of its perceptual uniformity, then normalize the values to mean 0 and standard deviation 1 before Procrustes alignment and scale calibration. The patch size is 16, matching the ImageNet model. Thus the pixel patch dimension is $D = 16 \times 16 \times 3 = 768$, while the AsymFlow rank follows the original FLUX.2 latent dimension, $r = d = 128$.

We finetune on a 3M subset of LAION-Aesthetics images [56], curated with safety and aesthetics filters. The images are resized to one-megapixel resolution and captioned with Qwen2.5-VL [3]. To reduce overfitting and preserve the pretrained model, we freeze the base weights and update only the input/output projection layers together with rank-256 LoRA adapters [24]. The trained modules are:

- `x_embedder`, `proj_out`, and `norm_out`;

- rank-256 LoRA adapters with dropout 0.05 on `*.ff.linear_in`, `*.ff.linear_out`, `*.ff_context.linear_in`, `*.ff_context.linear_out`, `timestep_embedder.linear_1`, `timestep_embedder.linear_2`, and `single_transformer_blocks.*.attn.to_out`.

Optimization uses 8-bit Adam [13, 30] with batch size 256, betas (0.9, 0.95), learning rate 10^{-4} for all trainable parameters (except that `proj_out` uses 10^{-3}). The final model used in the system comparison is trained for 15K iterations, costing approximately 1100 NVIDIA H100 GPU hours. For evaluation, we use the exponential moving average (EMA) of the finetuned weights with the dynamic EMA schedule of Karras et al. [29] (using the hyperparameter $\gamma = 7.0$). Sampling uses UniPC [74] with APG orthogonal-projection guidance [53]. At each sampling step, we convert the denoised pixels to RGB color space and clamp the values to the valid range before converting them back to Oklab velocity. Table 6 summarizes the main text-to-image settings.

Table 6: **Text-to-image finetuning and evaluation settings.**

Setting	Value
Pixel color space	Normalized Oklab [47]
Patch size	16
Patch dimension D	768
Patch rank r	128
Subspace construction	Orthogonal Procrustes lift with scale calibration
LoRA rank / dropout	256 / 0.05
Flow shift [15]	17.0
Training resolution	1MP with mixed aspect ratios
Pre-shift time sampling	LogitNormal(0, 1)
Optimizer	8-bit Adam [13, 30]
Learning rate	10^{-4} (10^{-3} for <code>proj_out</code>)
Adam betas	(0.9, 0.95)
Weight decay	0.0
Batch size	256
Training iterations	15K iterations
EMA	Dynamic EMA, $\gamma = 7.0$ [29]
Sampler	UniPC [74]
Guidance scale	4.0 with APG orthogonal projection [53]
Sampling steps	32

Latent baseline. For the latent finetuning baseline, we use its native flow shift of 7.0. Other settings are the same as AsymFlow for strict comparability.

DDT baseline. For the DDT pixel finetuning baseline, the DDT head uses two transformer blocks with a wider dimension of 32 attention heads \times 192 features per head, similar to the RAE design [75]. We use the same \mathbf{A} matrix as AsymFlow to initialize the input projection layer of the backbone, which closes the input gap and significantly improves the DDT baseline over a random initialization. The DDT head, input/output layers, and LoRA adapters are trained using a common learning rate of 10^{-4} . Other settings are the same as AsymFlow for strict comparability.

Inference time. AsymFLUX.2 klein uses the same number of tokens as the original FLUX.2 klein, so the per-step running time stays exactly the same as the original latent model. Since VAE is not used, the overall generation speed is marginally faster than the latent model.

C Mathematical Derivations

C.1 AsymFlow Decomposition and Recovery

We first make explicit the rank- r projector properties used throughout the paper. The columns of $\mathbf{A} \in \mathbb{R}^{D \times r}$ form an orthonormal basis for the chosen low-rank subspace, so $\mathbf{A}^T \mathbf{A} = \mathbf{I}_r$. This orthonormality makes $\mathbf{P} = \mathbf{A} \mathbf{A}^T$ the orthogonal projector onto that subspace. Applying \mathbf{P} twice is the same as applying it once, so $\mathbf{P}^2 = \mathbf{P}$. The complementary projector $\mathbf{I} - \mathbf{P}$ removes everything in the low-rank subspace, which gives $(\mathbf{I} - \mathbf{P})\mathbf{P} = \mathbf{0}$. Together, these properties mean that any vector can be cleanly separated into a low-rank component and an orthogonal component. The notation is summarized as:

$$\mathbf{A} \in \mathbb{R}^{D \times r}, \quad \mathbf{A}^T \mathbf{A} = \mathbf{I}_r, \quad \mathbf{P} = \mathbf{A} \mathbf{A}^T, \quad \mathbf{P}^2 = \mathbf{P}, \quad (\mathbf{I} - \mathbf{P})\mathbf{P} = \mathbf{0}. \quad (23)$$

We now restate the two targets in this notation. The standard velocity target combines full Gaussian noise with the data term. AsymFlow keeps the same full data term, but applies the projector only to the noise term:

$$\mathbf{u} := \boldsymbol{\epsilon} - \mathbf{x}_0, \quad \mathbf{u}_A := \mathbf{P}\boldsymbol{\epsilon} - \mathbf{x}_0. \quad (24)$$

Component decomposition. Projecting \mathbf{u}_A onto the low-rank subspace gives the true low-rank velocity. This branch of AsymFlow is still a velocity target. It contains low-rank noise minus low-rank data:

$$\mathbf{P}\mathbf{u}_A = \mathbf{P}(\mathbf{P}\boldsymbol{\epsilon} - \mathbf{x}_0) = \mathbf{P}\boldsymbol{\epsilon} - \mathbf{P}\mathbf{x}_0 = \mathbf{P}(\boldsymbol{\epsilon} - \mathbf{x}_0) = \mathbf{P}\mathbf{u}. \quad (25)$$

Projecting \mathbf{u}_A onto the orthogonal complement removes the noise term entirely. This branch is no longer a velocity target. It is the orthogonal clean-data component up to a minus sign:

$$(\mathbf{I} - \mathbf{P})\mathbf{u}_A = (\mathbf{I} - \mathbf{P})(\mathbf{P}\boldsymbol{\epsilon} - \mathbf{x}_0) = -(\mathbf{I} - \mathbf{P})\mathbf{x}_0. \quad (26)$$

Together, Eqs. (25) and (26) show that AsymFlow is velocity-like in $\text{Im}(\mathbf{P})$ and \mathbf{x}_0 -like in $\text{Im}(\mathbf{I} - \mathbf{P})$.

Recovery rule. The same decomposition gives an exact route from the asymmetric target back to the standard velocity target. The low-rank branch is already in velocity form, so this component is kept directly:

$$\mathbf{P}\mathbf{u} = \mathbf{P}\mathbf{u}_A. \quad (27)$$

The orthogonal branch is different. Since Eq. (26) says that $(\mathbf{I} - \mathbf{P})\mathbf{u}_A$ equals the negative clean-data component, the orthogonal clean data is obtained by changing the sign:

$$(\mathbf{I} - \mathbf{P})\mathbf{x}_0 = -(\mathbf{I} - \mathbf{P})\mathbf{u}_A. \quad (28)$$

This clean-data component is then converted to velocity using the usual \mathbf{x}_0 -to- \mathbf{u} relation. The orthogonal velocity is obtained by subtracting clean data from the noisy input and dividing by the noise level:

$$(\mathbf{I} - \mathbf{P})\mathbf{u} = (\mathbf{I} - \mathbf{P})\frac{\mathbf{x}_t - \mathbf{x}_0}{\sigma_t} = (\mathbf{I} - \mathbf{P})\frac{\mathbf{x}_t + \mathbf{u}_A}{\sigma_t}. \quad (29)$$

Combining the direct low-rank velocity branch with the converted orthogonal branch gives the full-rank velocity target:

$$\mathbf{u} = \mathbf{P}\mathbf{u}_A + (\mathbf{I} - \mathbf{P})\frac{\mathbf{x}_t + \mathbf{u}_A}{\sigma_t}. \quad (30)$$

Thus, the asymmetric target itself contains enough information to reconstruct the standard full-rank velocity target exactly.

Endpoint cases. The rank controls how much of the target is velocity-like. At rank zero, the projector is zero, so AsymFlow becomes full \mathbf{x}_0 -prediction up to sign. At full rank, the projector is the identity, so AsymFlow becomes standard velocity prediction:

$$r = 0 \implies \mathbf{P} = \mathbf{O}, \mathbf{u}_A = -\mathbf{x}_0, \quad r = D \implies \mathbf{P} = \mathbf{I}, \mathbf{u}_A = \boldsymbol{\epsilon} - \mathbf{x}_0 = \mathbf{u}. \quad (31)$$

C.2 Latent–Pixel Flow Coupling at Initialization

We next show the trajectory coupling relationship that makes latent-to-pixel initialization exact: when the latent and lifted pixel ODEs start from paired noise, the entire low-rank pixel trajectory can be lifted from the latent trajectory plus the analytically determined orthogonal noise component. This trajectory coupling holds for both scale-calibrated (Appendix A.2) and uncalibrated AsymFlows. Below we analyze the uncalibrated version for simplicity.

Let $\mathbf{z}_0 \in \mathbb{R}^d$ denote a latent token, where d is the latent dimension. In this construction we choose the pixel low-rank subspace to have the same rank $r = d$, and use a linear lift $\mathbf{A} \in \mathbb{R}^{D \times d}$ from latent tokens to pixel patches. As before, the columns of \mathbf{A} are orthonormal, so $\mathbf{A}^\top \mathbf{A} = \mathbf{I}_d$ and $\mathbf{P} = \mathbf{A}\mathbf{A}^\top$ projects onto the latent-induced pixel subspace. The lifted low-rank pixel target is $\mathbf{x}_0^L := \mathbf{A}\mathbf{z}_0$, and projecting pixel noise back through \mathbf{A}^\top gives the latent noise $\boldsymbol{\epsilon}_z := \mathbf{A}^\top \boldsymbol{\epsilon}$. The notation is summarized as:

$$\mathbf{A} \in \mathbb{R}^{D \times d}, \quad \mathbf{A}^\top \mathbf{A} = \mathbf{I}_d, \quad \mathbf{P} = \mathbf{A}\mathbf{A}^\top, \quad \mathbf{x}_0^L := \mathbf{A}\mathbf{z}_0, \quad \boldsymbol{\epsilon}_z := \mathbf{A}^\top \boldsymbol{\epsilon}. \quad (32)$$

With these definitions, projecting the lifted low-rank pixel process recovers the pretrained latent process.

Input identity. The pixel forward process diffuses the lifted low-rank pixels with full-rank pixel-space noise:

$$\mathbf{x}_t^L := \alpha_t \mathbf{x}_0^L + \sigma_t \boldsymbol{\epsilon} = \alpha_t \mathbf{A} \mathbf{z}_0 + \sigma_t \boldsymbol{\epsilon}. \quad (33)$$

Projecting this noisy pixel sample by \mathbf{A}^T gives exactly the corresponding noisy latent sample:

$$\mathbf{A}^T \mathbf{x}_t^L = \alpha_t \mathbf{A}^T \mathbf{A} \mathbf{z}_0 + \sigma_t \mathbf{A}^T \boldsymbol{\epsilon} = \alpha_t \mathbf{z}_0 + \sigma_t \boldsymbol{\epsilon}_z = \mathbf{z}_t. \quad (34)$$

Thus, the lifted pixel model evaluates the pretrained latent network at the paired noisy latent state.

Output identity. The latent model predicts latent velocity $\mathbf{u}_z := \boldsymbol{\epsilon}_z - \mathbf{z}_0$. Lifting this prediction to pixel space gives an AsymFlow-like target for the low-rank pixels \mathbf{x}_0^L :

$$\mathbf{A} \mathbf{u}_z = \mathbf{A}(\boldsymbol{\epsilon}_z - \mathbf{z}_0) = \mathbf{A} \mathbf{A}^T \boldsymbol{\epsilon} - \mathbf{A} \mathbf{z}_0 = \mathbf{P} \boldsymbol{\epsilon} - \mathbf{x}_0^L. \quad (35)$$

Therefore the low-rank pixel velocity $\mathbf{u}^L := \boldsymbol{\epsilon} - \mathbf{x}_0^L$ is obtained by applying the same recovery rule from Sec. C.1 with $\mathbf{u}_A = \mathbf{A} \mathbf{u}_z$ and $\mathbf{x}_t = \mathbf{x}_t^L$:

$$\mathbf{u}^L = \mathbf{P} \mathbf{A} \mathbf{u}_z + (\mathbf{I} - \mathbf{P}) \frac{\mathbf{x}_t^L + \mathbf{A} \mathbf{u}_z}{\sigma_t}. \quad (36)$$

For analyzing the lifted latent initialization, this expression can be simplified because the lifted latent prediction already lies in the low-rank subspace, so we have $(\mathbf{I} - \mathbf{P}) \mathbf{A} \mathbf{u}_z = \mathbf{0}$. This gives

$$\mathbf{u}^L = \mathbf{A} \mathbf{u}_z + \frac{(\mathbf{I} - \mathbf{P}) \mathbf{x}_t^L}{\sigma_t}. \quad (37)$$

Thus, at initialization, the low-rank branch is exactly the lifted latent velocity, while the orthogonal branch is recovered directly from the current noisy pixel state. Note that this simplification does not apply to the finetuned AsymFlow model and should not be used in the implementation.

Trajectory coupling. The identities above are pointwise statements about the noisy input and the recovered velocity. What we need for initialization is slightly stronger: if the latent model and the lifted pixel model are solved in parallel from paired noise, then their whole trajectories remain paired, and their final samples still satisfy the same lifting relation.

Theorem 1. Let $\boldsymbol{\epsilon} \in \mathbb{R}^D$ be a pixel-space noise sample and let $\boldsymbol{\epsilon}_z = \mathbf{A}^T \boldsymbol{\epsilon}$ be its low-rank projection. Let G_ϕ denote the pretrained latent flow velocity network. Consider the latent flow ODE on $(0, 1]$:

$$\frac{d\mathbf{z}_t}{dt} = G_\phi(\mathbf{z}_t, t), \quad \mathbf{z}_1 = \boldsymbol{\epsilon}_z, \quad (38)$$

and the lifted pixel flow ODE obtained by applying the simplified form in Eq. (37) to the latent network output:

$$\frac{d\mathbf{x}_t^L}{dt} = \mathbf{A} G_\phi(\mathbf{A}^T \mathbf{x}_t^L, t) + \frac{(\mathbf{I} - \mathbf{P}) \mathbf{x}_t^L}{\sigma_t}, \quad \mathbf{x}_1^L = \boldsymbol{\epsilon}. \quad (39)$$

Then the two trajectories satisfy

$$\mathbf{x}_t^L = \mathbf{A} \mathbf{z}_t + \sigma_t (\mathbf{I} - \mathbf{P}) \boldsymbol{\epsilon} \quad \text{for all } t \in (0, 1]. \quad (40)$$

In particular, taking $t \rightarrow 0$ gives the final sample identity $\mathbf{x}_0^L = \mathbf{A} \mathbf{z}_0$.

Proof. For brevity, write the orthogonal noise component as $\boldsymbol{\epsilon}^\perp := (\mathbf{I} - \mathbf{P}) \boldsymbol{\epsilon}$. Then the pixel noise decomposes into the lifted latent noise plus the orthogonal residual:

$$\boldsymbol{\epsilon} = \mathbf{P} \boldsymbol{\epsilon} + (\mathbf{I} - \mathbf{P}) \boldsymbol{\epsilon} = \mathbf{A} \mathbf{A}^T \boldsymbol{\epsilon} + \boldsymbol{\epsilon}^\perp = \mathbf{A} \boldsymbol{\epsilon}_z + \boldsymbol{\epsilon}^\perp. \quad (41)$$

At $t = 1$, this decomposition matches the two ODE initial conditions:

$$\mathbf{x}_1^L = \mathbf{A} \mathbf{z}_1 + \sigma_1 \boldsymbol{\epsilon}^\perp. \quad (42)$$

Now define a candidate lifted pixel trajectory from the latent trajectory:

$$\tilde{\mathbf{x}}_t^L := \mathbf{A} \mathbf{z}_t + \sigma_t \boldsymbol{\epsilon}^\perp. \quad (43)$$

We will show that this candidate trajectory satisfies the lifted pixel ODE in Eq. (39) with the same initial condition, so by uniqueness of ODE solutions, it must be identical to \mathbf{x}_t^L for all t . The candidate trajectory has exactly the input identity required by the latent network:

$$\mathbf{A}^T \tilde{\mathbf{x}}_t^L = \mathbf{A}^T \mathbf{A} \mathbf{z}_t + \sigma_t \mathbf{A}^T \boldsymbol{\epsilon}^\perp = \mathbf{z}_t. \quad (44)$$

It also has an orthogonal component determined only by the fixed orthogonal noise:

$$(\mathbf{I} - \mathbf{P})\tilde{\mathbf{x}}_t^L = \sigma_t \epsilon^\perp. \quad (45)$$

Substituting these two identities into the lifted pixel vector field gives the lifted latent velocity plus the orthogonal noise velocity:

$$\mathbf{A}G_\phi(\mathbf{A}^\top \tilde{\mathbf{x}}_t^L, t) + \frac{(\mathbf{I} - \mathbf{P})\tilde{\mathbf{x}}_t^L}{\sigma_t} = \mathbf{A}G_\phi(\mathbf{z}_t, t) + \epsilon^\perp. \quad (46)$$

The derivative of the candidate trajectory gives the same expression:

$$\frac{d\tilde{\mathbf{x}}_t^L}{dt} = \mathbf{A} \frac{d\mathbf{z}_t}{dt} + \frac{d\sigma_t}{dt} \epsilon^\perp = \mathbf{A}G_\phi(\mathbf{z}_t, t) + \epsilon^\perp, \quad (47)$$

where we used Eq. (38) and $\sigma_t = t$. Thus $\tilde{\mathbf{x}}_t^L$ satisfies the lifted pixel ODE in Eq. (39). Since it also has the same value as \mathbf{x}_t^L at $t = 1$, uniqueness of the ODE solution gives

$$\mathbf{x}_t^L = \tilde{\mathbf{x}}_t^L = \mathbf{A}\mathbf{z}_t + \sigma_t(\mathbf{I} - \mathbf{P})\epsilon \quad \text{for all } t \in (0, 1]. \quad (48)$$

Finally, taking $t \rightarrow 0$ gives $\mathbf{x}_0^L = \mathbf{A}\mathbf{z}_0$. \square

The same argument applies to Euler discretization with a shared time grid: if the relation holds before a step, the latent update changes the low-rank component by $\Delta t \mathbf{A}G_\phi(\mathbf{z}_t, t)$, while the lifted pixel update additionally changes the orthogonal component by $\Delta t \epsilon^\perp$, preserving the same paired form after the step; by induction, the relation holds at all steps. Thus, at network initialization, the lifted latent model is an exact low-rank pixel flow model. Note that this initialization is not yet a full AsymFlow model on real pixels, as finetuning replaces the lifted low-rank data target \mathbf{x}_0^L with the full-rank pixel target \mathbf{x}_0 .

C.3 Details on Variance-Reduced Loss

The variance-reduced loss in Sec. 5.2 can be viewed as a control variate. The paired low-rank target \mathbf{x}_0^L is correlated with the full pixel target \mathbf{x}_0 , and a frozen initialized low-rank model gives a good estimate of it. We use this paired target to reduce the variance of the pixel residual without changing the conditional mean target.

The exact control-variate identity is

$$\mathbb{E}[\mathbf{x}_0^L - \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t] | \mathbf{x}_t] = \mathbf{0}. \quad (49)$$

Therefore adding any coefficient times this zero-mean residual does not change the conditional target. The posterior mean remains unchanged, while the sampled target can have lower variance:

$$\mathbb{E}[\mathbf{x}_0 - \lambda(\mathbf{x}_0^L - \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t]) | \mathbf{x}_t] = \mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]. \quad (50)$$

Before approximation, the objective is therefore equivalent to the standard flow matching loss in \mathbf{x}_0 format (Eq. (2)). The only role of the additional term is to reduce sampling variance when the low-rank residual explains part of the full pixel residual.

In practice, the conditional mean $\mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t]$ is unavailable. We approximate it using the frozen low-rank model prediction $\hat{\mathbf{x}}_0^L$ from the paired noisy low-rank sample:

$$\mathbf{x}_t^L = \alpha_t \mathbf{x}_0^L + \sigma_t \epsilon, \quad \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t] \approx \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t^L] \approx \hat{\mathbf{x}}_0^L = \mathbf{P}\mathbf{x}_t^L - \sigma_t \mathbf{A}G_\phi(\mathbf{A}^\top \mathbf{x}_t^L, t). \quad (51)$$

Substituting this approximation gives the practical variance-reduced loss in Eq. (7).

The approximation $\mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t] \approx \mathbb{E}[\mathbf{x}_0^L | \mathbf{x}_t^L]$ is exact under the sufficient condition that the full noisy input and the paired low-rank noisy input differ only in the orthogonal complement. In that case, their low-rank components match, so the frozen low-rank model receives the same low-rank information:

$$\mathbf{x}_t - \mathbf{x}_t^L \in \text{Im}(\mathbf{I} - \mathbf{P}) \implies \mathbf{A}^\top \mathbf{x}_t = \mathbf{A}^\top \mathbf{x}_t^L. \quad (52)$$

This requires either $t = 1$ or $\mathbf{x}_0 - \mathbf{x}_0^L \in \text{Im}(\mathbf{I} - \mathbf{P})$, which is generally not satisfied due to the non-linearity of the VAE encoder [51]. When this condition is not satisfied, the approximation error appears inside the low-rank subspace $\text{Im}(\mathbf{P})$. To compensate for this, the perceptual correction is introduced in the low-noise regime in place of the variance reduction, as detailed in Sec. A.4.

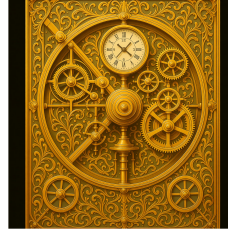
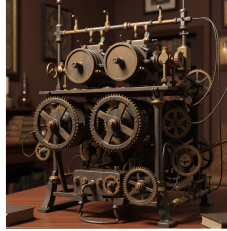
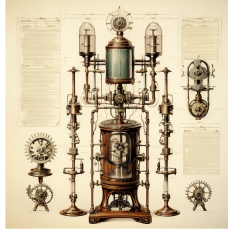
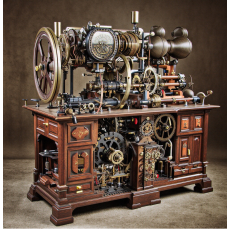
D Additional Qualitative Results

AsymFLUX.2 klein (ours)
Pixel

PixelDIT-T2I
Pixel

FLUX.2 klein Base
Latent

Qwen Image
Latent



A complex victorian apparatus, highly detailed digital photograph.



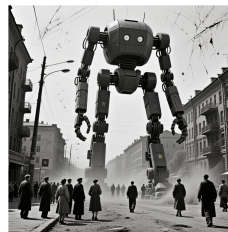
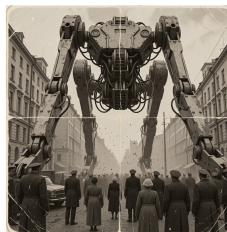
a movie still from a stanley kubrick film, The figure in the 1950s hazmat suit running as fast as they could, their heart pounding in their chest, Behind them, a massive nuclear cloud loomed, spreading destruction in its wake, They looked around, desperate for any sign of shelter, but all they saw was the endless expanse of ruins, chirasculo lighting, epic composition



beautiful woman wearing a tight dress with cut outs at the hips and legs made of colored glass and acrylic stands in front of a concept car parked at a modern vacation home in the style of Syd Mead, wide angle, sunny



environmental photograph giant octopus crushing car in town at night, horror, 1980, photograph in the style of jeff wall, 35mm, leica m, cinematic lighting, intricate details, realistic facial features, highly detailed, cinematic, kodak portra 800, kodak color film, cinematic, film grain, large film noise



massive robot machine long robotic legs, above citizens in street, soviet russia, dusty, 70's 60's film, scratched, boris mikhailov gelatin silver print photography, evil, dystopia

Figure 9: Additional qualitative text-to-image comparisons (part A).

AsymFLUX.2 klein (ours)
Pixel

PixelDIT-T2l
Pixel

FLUX.2 klein Base
Latent

Qwen Image
Latent



The image captures a romantic and dramatic scene featuring a couple embraced in a misty, almost ethereal forest. The color palette is dominated by shades of red and gray, creating a sense of passion and intensity against a backdrop of serene nature. A man and a woman stand close together, both dressed in red. The man wears a long-sleeved shirt and dark pants, while the woman is adorned in a voluminous, flowing red dress that cascades around her. They are embraced and looking at each other. The surrounding environment is a forest with tall trees, some with vibrant red foliage overhead that contrasts beautifully with the muted gray of the mist. A stream or river flows nearby, with small waterfalls adding to the dynamic composition. In the background, a flash of lightning illuminates the sky, enhancing the dramatic atmosphere of the scene. The overall impression is one of intense emotion within a mystical, natural setting.



DVD Screengrab From 1978 sci-fi Film, "starwars"; full body, depth of field, ultra realistic, hyper detailed, 35mm lens, editorial photography, photorealism, volumetric light, epic scene, post production, 8k,



The photograph presents a young woman with her hands held up towards the viewer. Her hands are covered in a dark substance, possibly dirt or soot, which creates a striking contrast against the lighter skin visible around her wrists. The woman's face is slightly blurred but you can see that she has dark eyebrows, red lipstick, and her hair is swept back. There's a slight smile on her face. She's wearing an olive-green jacket, which suggests an outdoor setting. A thin red string is tied around her left wrist, adding a small splash of color to the overall muted tones. On her right wrist, she wears a gold bracelet. The background is intentionally soft and out of focus, dominated by shades of gray and brown. Hints of trees and indistinct shapes suggest a natural environment, perhaps a forest or park. The soft focus keeps the attention on the woman and her extended hands, making her the central subject of the photograph.



production still from 1974 of Alejandro Jodorowsky's enormous crowd in stadium worshipping a one hundred foot tall messiah of pulsating humans joined together in a wooden exoskeleton, ch200 ASA 35mm



a chef cook is filming a tik tok video inside a restaurant kitchen

Figure 10: Additional qualitative text-to-image comparisons (part B).

E Impact Statement

Our method enhances the photorealism of diffusion models, which significantly benefits creative industries by enabling high-fidelity prototyping and asset creation. This advancement, however, presents a dual-use challenge: more realistic imagery facilitates the creation of convincing disinformation or non-consensual media, increasing the potential for societal harm. Higher visual quality also requires renewed scrutiny of dataset biases, as those biases will be rendered more persuasively. We open-source our model to encourage scientific replication, but emphasize that responsible deployment requires the use of standard safety filters and content provenance tools (like watermarking) to manage these risks.