

Designing Instance-Level Sampling Schedules via REINFORCE with James-Stein Shrinkage

Supplementary Material

A. Proofs

A.1. Proof of Proposition 3.1.

Proof. First, we show the unbiasedness of the baseline. We abbreviate the policy as $\pi_\theta(\tau) \equiv \pi_\theta(\tau \mid \mathbf{x}_T, \mathbf{c})$,

$$g(\tau) = \nabla_\theta \log \pi_\theta(\tau), \quad w(\tau) = \|g(\tau)\|^2,$$

and let $r(\tau)$ denote the scalar reward. The single-sample REINFORCE estimator with a scalar baseline b is

$$\widehat{G}(b) = (r(\tau) - b)g(\tau). \quad (11)$$

By swapping integral and gradient operators

$$\mathbb{E}_{\tau \sim \pi_\theta}[g(\tau)] = \mathbb{E}_{\tau \sim \pi_\theta}[\nabla_\theta \log \pi_\theta(\tau)] = \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta}[1] = 0. \quad (12)$$

Hence the estimator *remains unbiased* for any baseline b :

$$\mathbb{E}[\widehat{G}(b)] = \mathbb{E}[(r - b)g] = \mathbb{E}[rg] - b \mathbb{E}[g] = \mathbb{E}[rg]. \quad (13)$$

Next, we study the variances of the estimator. The covariance of $\widehat{G}(b)$ is

$$\text{Var}[\widehat{G}(b)] = \mathbb{E}[(r - b)^2 gg^\top] - \mathbb{E}[(r - b)g] \mathbb{E}[(r - b)g]^\top.$$

Since $\mathbb{E}[(r - b)g] = \mathbb{E}[rg]$ does not depend on b , minimizing the total scalar variance (trace) is equivalent to minimizing

$$V(b) = \mathbb{E}[(r - b)^2 w(\tau)] - \text{const}, \quad \text{Recall } w(\tau) = \|g(\tau)\|^2.$$

The function $V(b)$ is differentiable and strictly convex in b . Differentiating and setting the derivative to zero yields

$$\begin{aligned} \frac{d}{db} V(b) &= \frac{d}{db} \mathbb{E}[(r - b)^2 w] = \mathbb{E}[2(b - r)w] \\ &= 2(b \mathbb{E}[w] - \mathbb{E}[rw]) = 0, \end{aligned}$$

which gives the unique minimizer

$$b^* = \frac{\mathbb{E}[r(\tau) w(\tau)]}{\mathbb{E}[w(\tau)]} = \frac{\mathbb{E}[r(\tau) \|\nabla_\theta \log \pi_\theta(\tau)\|^2]}{\mathbb{E}[\|\nabla_\theta \log \pi_\theta(\tau)\|^2]}. \quad (14)$$

Thus b^* minimizes the trace of $\text{Var}[\widehat{G}(b)]$. \square

Remarks. (i) Because b is a scalar shared across all components of $g(\tau)$, the same minimizer b^* arises whether one minimizes each component variance or the trace. (ii) Our derivation shares the spirit of state-dependent baselines derived for GPOMDP in Greensmith et al. [6], which optimize a different variance criterion conditioned on Markov states. (iii) When the policy becomes nearly deterministic for a fixed context $(\mathbf{x}_T, \mathbf{c})$, the variation of $\nabla_\theta \log \pi_\theta(\tau)$ is small. In this regime, Eq. (14) reduces to the mean reward $\mathbb{E}[r(\tau) \mid \mathbf{x}_T, \mathbf{c}]$, which motivates the empirical approximations used in the main text.

A.2. Proof of Theorem 3.2.

Proof. Now we prove the theoretical results related to the James-Stein reward baseline.

Bayesian optimality. Recall that we have the random-effects model in Eq. (6),

$$r^{(c,i)} \mid \mu_c \sim \mathcal{N}(\mu_c, \sigma^2), \quad \mu_c \sim \mathcal{N}(\mu_0, \delta^2), \quad (15)$$

the per-context mean satisfies

$$\bar{r}_c = \frac{1}{K_c} \sum_{i=1}^{K_c} r^{(c,i)}, \quad \bar{r}_c \mid \mu_c \sim \mathcal{N}(\mu_c, \sigma^2/K_c). \quad (16)$$

By normal-normal conjugacy, the posterior distribution $\mu_c \mid \{r^{(c,i)}\}_{i=1}^{K_c}$ is Gaussian with mean

$$\mathbb{E}[\mu_c \mid \{r^{(c,i)}\}] = (1 - \alpha_c^*) \bar{r}_c + \alpha_c^* \mu_0, \quad (17)$$

where the (population) shrinkage coefficient is

$$\alpha_c^* = \frac{\sigma^2/K_c}{\sigma^2/K_c + \delta^2}. \quad (18)$$

For any data-dependent predictor t_c of μ_c , the posterior expected squared error given the observed rewards is

$$\begin{aligned} \mathbb{E}[(\mu_c - t_c)^2 \mid \{r^{(c,i)}\}] &= \mathbb{E}[\mu_c^2 \mid \{r^{(c,i)}\}] \\ &\quad - 2t_c \mathbb{E}[\mu_c \mid \{r^{(c,i)}\}] + t_c^2, \end{aligned} \quad (19)$$

which, as a quadratic in t_c , is minimized at $t_c = \mathbb{E}[\mu_c \mid \{r^{(c,i)}\}]$. Thus the Bayes-optimal predictor of μ_c under squared error is the posterior mean in Eq. (17), i.e., the convex combination of the within-context mean \bar{r}_c and the population mean μ_0 .

In our REINFORCE setting, we require a baseline that is independent of the current rollout (c, i) to preserve unbiasedness. We therefore replace \bar{r}_c and the global mean μ_0 by their leave-one-out analogues:

$$\begin{aligned}\bar{r}_c^{(-i)} &= \frac{1}{K_c - 1} \sum_{j \neq i} r^{(c,j)}, \\ \bar{r}_{\cdot\cdot}^{(-c,-i)} &= \frac{1}{(\sum_{c'} K_{c'}) - 1} \left(\sum_{c'=1}^B \sum_{j=1}^{K_{c'}} r^{(c',j)} - r^{(c,i)} \right).\end{aligned}\tag{20}$$

The ideal Bayes baseline for rollout (c, i) is then

$$b_{\text{Bayes}}^{(c,i)} = (1 - \alpha_c^*) \bar{r}_c^{(-i)} + \alpha_c^* \bar{r}_{\cdot\cdot}^{(-c,-i)}.$$

In practice, σ^2 and δ^2 are unknown. The method-of-moments estimators in Eq. (10) yield empirical counterparts $\hat{\sigma}^2$ and $\hat{\delta}^2$, which in turn define the empirical shrinkage weights

$$\hat{\alpha}_c = \frac{\hat{\sigma}^2 / (K_c - 1)}{\hat{\sigma}^2 / (K_c - 1) + \hat{\delta}^2}.$$

Substituting $\hat{\alpha}_c$ and the leave-one-out means into the Bayes form above recovers exactly the empirical James-Stein baseline

$$b_{\text{JS}}^{(c,i)} = (1 - \hat{\alpha}_c) b_{\text{RLOO}}^{(c,i)} + \hat{\alpha}_c b_{\text{xctx}}^{(c,i)}$$

from Eq. (7), where both $b_{\text{RLOO}}^{(c,i)}$ and $b_{\text{xctx}}^{(c,i)}$ are computed in a leave-one-out way that excludes (c, i) . Since $b_{\text{JS}}^{(c,i)}$ does not depend on the current reward $r^{(c,i)}$, the REINFORCE gradient remains unbiased. In summary, $b_{\text{JS}}^{(c,i)}$ by construction coincides with the empirical Bayes posterior mean of μ_c under the model (6), and is Bayes-optimal.

MSE improvement. Under the same random-effects model, consider the convex estimator

$$\tilde{\mu}_c(\alpha) = (1 - \alpha) \bar{r}_c^{(-i)} + \alpha \bar{r}_{\cdot\cdot}^{(-c,-i)}.$$

Assuming independence between the within-context and across-context noise, its mean-squared error satisfies

$$\text{MSE}(\tilde{\mu}_c(\alpha)) = (1 - \alpha)^2 \frac{\sigma^2}{K_c - 1} + \alpha^2 \delta^2,$$

which is minimized at

$$\alpha^* = \frac{\sigma^2 / (K_c - 1)}{\sigma^2 / (K_c - 1) + \delta^2}.$$

Thus the same shrinkage coefficient α^* that defines the posterior mean also uniquely minimizes the frequentist MSE of the baseline, strictly improving not only over the purely contextual baseline (RLOO), but also over all the convex combination between RLOO and the XCTX baselines including themselves. \square

Remarks. If we collect the context-specific means into the vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_B)$ and likewise stack their noisy estimators, we obtain the classical homoscedastic normal-mean problem in dimension B . In that setting, James-Stein shrinkage estimators that shrink the vector of means toward a common target are known to uniformly dominate the unshrunk sample mean in total squared error risk whenever the dimension is at least three; see, e.g., James et al. [7]. In our notation, the ‘‘dimension’’ plays the role of the number of contexts B , which is the reason for the $B \geq 3$ condition mentioned in the main text. Our empirical Bayes baseline has the same shrinkage-to-a-global-mean structure as these classical estimators, but here we focus on the random-effects and REINFORCE viewpoints rather than re-proving the classical domination result.

B. Scheduler Network Architectures

We summarize the network architecture we used for implementing the scheduler network in Tab. 5. The scheduler is 0.1–1% the size ($\sim 20\text{M}$ parameters) of the backbone networks used in Stable Diffusion and Flux model families, and incurs at most ~ 0.1 –0.3% extra end-to-end sampling compute in the step counts we used ranging from 5 to 80.

C. Experiment Settings

C.1. Synthetic Sanity Check for Baseline Variance

RLOO is an unbiased baseline but fragile in exactly the regime most relevant to our policy learning scenarios: limited number of rollouts per context and highly heterogeneous reward scales. To isolate this failure mode and evaluate variance reduction in a controlled setting, we construct a synthetic experiment where the policy is analytic, rewards follow a spiky random-effects model. We measure the variance of the REINFORCE gradient estimators since they are all unbiased. This setup cleanly reveals the effect of the proposed James-Stein (JS) baseline.

Setup. We adopt an analytic Dirichlet policy $\pi(\boldsymbol{\tau}) = \text{Dirichlet}(\boldsymbol{\alpha})$ with concentration $\boldsymbol{\alpha} = 2\mathbf{1}_T$ and horizons $T \in \{4, 16, 64\}$. Its score function has a closed form:

$$\nabla_{\alpha_j} \log \pi(\boldsymbol{\tau}) = \psi\left(\sum_t \alpha_t\right) - \psi(\alpha_j) + \log \tau_j,$$

where $\psi(\cdot)$ denotes the digamma function and $\log \tau_j$ is the elementwise logarithm of the j -th coordinate of $\boldsymbol{\tau}$. To mimic the two-level random-effects model used in our theoretical analysis, we generate synthetic rewards as

$$\begin{aligned}r^{(c,i)} &= \mu_c + \varepsilon^{(c,i)} + \delta^{(c,i)}, \\ \mu_c &\sim \mathcal{N}(0, 1), \quad \varepsilon^{(c,i)} \sim \mathcal{N}(0, s_c^2), \\ s_c &\sim \text{LogNormal}(0, 1), \quad \delta^{(c,i)} \sim \text{Bernoulli}(0.15) \times 8,\end{aligned}\tag{21}$$

Table 5. **Network architecture for the Dirichlet schedule policy** $\alpha_\theta(x, c)$. The network takes an image-shaped noise tensor x , a sequence of text embeddings e_{text} , and an optional pooled text embedding \bar{e}_{text} (consistent with SD-XL, SD3.5M/L and FLUX models), and outputs Dirichlet parameters $\alpha \in \mathbb{R}_+^{L+1}$ over L schedule intervals plus one skipped interval. Default network architecture hyperparameters: `text_embed_dim` $d_{\text{text}} = 2048$, `pooled_text_embed_dim` $d_{\text{pool}} = 1280$, `image_encoder_depth` $n_{\text{conv}} = 2$, `image_encoder_width` $d_{\text{conv}} = 32$, `attention_dim` $d_{\text{attn}} = 256$, `cross_attention_heads` $n_{\text{attn}} = 4$, `number_of_transformer_blocks` $N_b = 4$, `hidden_dim` $d_h = 256$, `num_mlp_layers` $n_{\text{mlp}} = 2$.

Layers	Output size	Note
Inputs		
Input: x	$B \times H \times W \times C$	image-shaped noise
Input: e_{text}	$B \times L_{\text{text}} \times d_{\text{text}}$	token text emb.; $d_{\text{text}}=2048$
Input: \bar{e}_{text}	$B \times d_{\text{pool}}$	optional pooled emb.; $d_{\text{pool}}=1280$
Image encoder & cross-attention blocks		
Image encoder conv stack (inside block i , layers $j = 0, \dots, n_{\text{conv}}-1$)	$B \times H_i \times W_i \times C_i$	n_{conv} conv layers, 3×3 , stride 1 GroupNorm, SiLU $C_i = d_{\text{conv}} \cdot 2^{\min(4, i+j)}$
Flatten spatial dims	$B \times (H_i W_i) \times C_i$	prepare for attention
Query projection (Dense)	$B \times (H_i W_i) \times C_i$	applied to image features (query_projection)
Key/Value projection (Dense)	$B \times L_{\text{text}} \times C_i$	applied to e_{text} (key_value_projection)
MultiHeadDotProductAttention	$B \times (H_i W_i) \times C_i$	cross-attn: query from image, key/value from text, n_{attn} heads
Residual + LayerNorm	$B \times (H_i W_i) \times C_i$	LayerNorm(query + attn_output) (cross_attn_norm)
Reshape to image grid	$B \times H_i \times W_i \times C_i$	output of SchedulerTransformerBlock
Global avg. pooling (per block)	$B \times C_i$	mean over (H_i, W_i) for feature fusion
Downsample Conv 3×3 , stride 2 + GroupNorm + SiLU (between blocks)	$B \times H_{i+1} \times W_{i+1} \times C_{i+1}$	from reshaped image grid, for $i < N_b - 1$; $C_{i+1} = d_{\text{conv}} \cdot 2^{\min(4, (i+1) n_{\text{conv}} - 1)}$
Feature pyramid fusion		
Concat. over blocks	$B \times \sum_i C_i$	feature pyramid $\{\text{block } i\}_{i=0}^{N_b-1}$
Concat. with \bar{e}_{text}	$B \times (\sum_i C_i + d_{\text{pool}})$	only if pooled emb. used
MLP Dirichlet head		
Dense, SiLU $\times (n_{\text{mlp}} - 1)$	$B \times d_h$	default $d_h = 256$
Dense (Dirichlet logits head)	$B \times (L+1)$	$L = \text{num_timesteps}$
Softplus $+10^{-3}$	$B \times (L+1)$	$\alpha = \text{softplus}(\cdot) + 10^{-3}$

where μ_c acts as the context-level effect (analogous to $\mu_0 + \xi^{(c)}$ in Eq. (6)) and $\varepsilon^{(c,i)}$ represents within-context noise. The additional sparse term $\delta^{(c,i)}$ introduces occasional outliers, creating a heavy-tailed, heteroskedastic reward landscape that stresses RLOO in the small- K regime.

We sweep through different settings for (B, K, L) . For each (B, K, L) we generate 500 i.i.d. batches, compute

$$\hat{g}(b) = \frac{1}{BK} \sum_{c=1}^B \sum_{i=1}^K (r^{(c,i)} - b^{(c,i)}) \nabla_\alpha \log \pi(\tau^{(c,i)}),$$

once with RLOO and once with JS, and report the empirical per-dimension variance $\frac{1}{L} \sum_{j=1}^L \text{Var}(\hat{g}_j)$. We sweep $B \in \{8, 16, 32\}$ and $K \in \{2, 4, 8, 16\}$.

Observations. As shown in Tab. 6, the James–Stein (JS) baseline delivers uniformly lower gradient variance than RLOO across all batch sizes, rollout counts, and horizons. For the most challenging regime ($K = 2$), JS reduces variance by roughly 45–50%, effectively matching the stability one would expect from doubling the rollout count—without any additional computation. Even with moderate rollouts ($K = 4$), variance drops by 20–25%, demonstrating that the improvement is systematic rather than case-specific. Across all settings, JS remains consistently superior, reflecting its adaptive shrinkage across both within- and cross-context variability. In practical terms, this means JS acts as a computationally free variance amplifier: it increases the *effective* number of rollouts available to policy-gradient estimation, providing a principled and efficient replacement for

Table 6. Synthetic variance sanity check. Average per-dimension variance of REINFORCE for an analytic Dirichlet policy. Rewards follow Eq. (21). Entries are RLOO \rightarrow JS. JS consistently gives gain over RLOO baselines for various settings.

K	B	$L=4$	$L=16$	$L=64$
2	8	1.055 \rightarrow 0.612	1.445 \rightarrow 0.725	1.148 \rightarrow 0.625
	16	0.479 \rightarrow 0.285	0.663 \rightarrow 0.339	0.606 \rightarrow 0.339
	32	0.261 \rightarrow 0.140	0.328 \rightarrow 0.174	0.310 \rightarrow 0.165
4	8	0.393 \rightarrow 0.284	0.402 \rightarrow 0.307	0.386 \rightarrow 0.299
	16	0.175 \rightarrow 0.143	0.193 \rightarrow 0.154	0.210 \rightarrow 0.163
	32	0.086 \rightarrow 0.066	0.099 \rightarrow 0.079	0.106 \rightarrow 0.084
8	8	0.140 \rightarrow 0.126	0.191 \rightarrow 0.170	0.182 \rightarrow 0.166
	16	0.071 \rightarrow 0.064	0.099 \rightarrow 0.090	0.086 \rightarrow 0.078
	32	0.038 \rightarrow 0.034	0.042 \rightarrow 0.038	0.045 \rightarrow 0.040
16	8	0.072 \rightarrow 0.070	0.081 \rightarrow 0.078	0.094 \rightarrow 0.089
	16	0.031 \rightarrow 0.030	0.038 \rightarrow 0.037	0.039 \rightarrow 0.038
	32	0.017 \rightarrow 0.016	0.020 \rightarrow 0.019	0.020 \rightarrow 0.019

Table 7. Learning rates for schedule training on HPDv2 prompts.

Backbone	$L=5$	$L=10$	$L=20$	$L=40$	$L=80$
SD-XL	5e-5	5e-5	5e-5	5e-5	5e-5
SD3.5-M	5e-5	2.5e-5	2.5e-5	1.25e-5	5e-5
SD3.5-L	5e-5	5e-5	5e-5	5e-5	5e-5
Flux-Dev	2.5e-5	5e-5	2.5e-5	5e-5	5e-5

RLOO in rollout-expensive training regimes. We visualize a side-by-side comparison in Fig. 2a for the $B = 32, L = 16$ case with different number of rollouts K .

C.2. General Text-to-Image Experiments on HPDv2

General settings. For all the experiments, we consistently use the AdamW [15] optimizer for learning scheduler parameters. We set the `weight_decay` to $1e-4$ throughout experiments, tune the learning rates (see Tab. 7), and use default values for the remaining hyperparameters. We apply the gradient clipping with a maximum norm of 1. For all models we consistently use a rollout number of 2, and set the output resolution to 1024×1024 . Due to memory limit, we use a batch size of 32 for smaller models like SD-XL or SD-3.5M, and a batch size of 16 for larger ones like SD-3.5L and Flux-Dev. We use a guidance weight of 7.5 for SD-XL and 5.0 for SD3.5-M/L and Flux models. For SD-XL model we use the DDIM solver [21] for per step denoising update; for the flow-matching-based SD3.5 and Flux models we use the Euler solver.

TPDM baseline. We reimplement the TPDM baseline using (i) the same network architecture we used, and (ii) the same set of seeds, hyperparameters and base model settings for training. We follow the customized PPO implementation in TPDM with no ratio clipping.

Table 8. Learning rates for schedule training on text rendering and general-style prompts. Num. of sampling steps is 40.

Backbone	SD3.5-M	Flux-Dev
Text Render	1.25e-5	5e-5
General	1.25e-6	5e-5

We would like to explicitly mention that (i) we *exclude* any regularizer including KL constraint between a reference schedule and the learned schedule during training to be consistent with our settings. This is different from the vanilla implementation of TPDM, where KL constraint is used. (ii) TPDM in its basic form *does not* involve text condition as its input. We experiment with different settings: (i) we faithfully nullify the textual inputs by setting the embeddings to zero vectors to train the model, and (ii) we include the textual inputs as in our settings. We report in Tab. 1 the results for (i). For (ii) the variant, especially in our setting, reduces to the RLOO baseline, which we also report in Tab. 1.

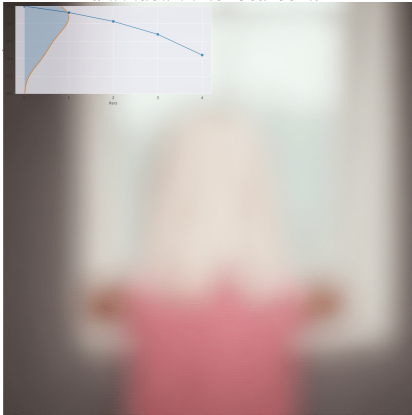
C.3. Text Rendering and Fine-Grained Alignment

We use the same settings as in Sec. C.2, except for tuning the learning rates (see Tab. 8). For the OCR metrics, we adopt the commonly used mask textspotter v3 model [11] to calculate the character level accuracy, precision and recall scores for each image, and then calculate the average across the testing samples as the final results in Tab. 3.

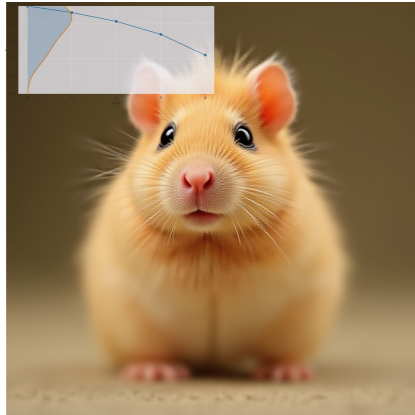
D. Additional qualitative results

We provide further T2I results for 5-step sampling (Figs. 8 and 9), general T2I generation on HPDv2 prompts (Figs. 10 and 11), text rendering (Figs. 12 and 13) and fine-grained object-focused generation (Figs. 14 and 15). We follow the format of Fig. 1 and plot the corresponding schedules for generation at the top left corner.

A white-haired girl in a pink sweater looks out a window in her bedroom.



A hamster resembling a horse.

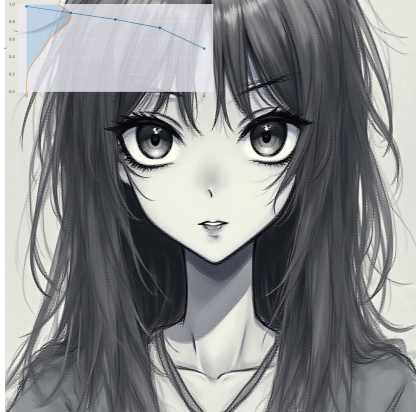
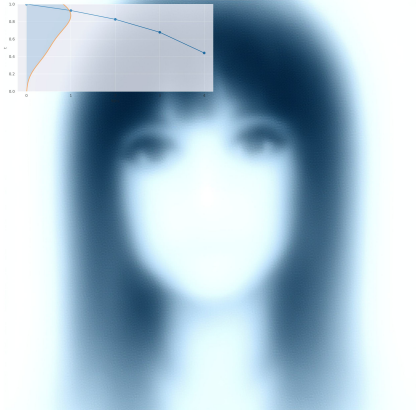


A lemon wearing sunglasses on the beach.

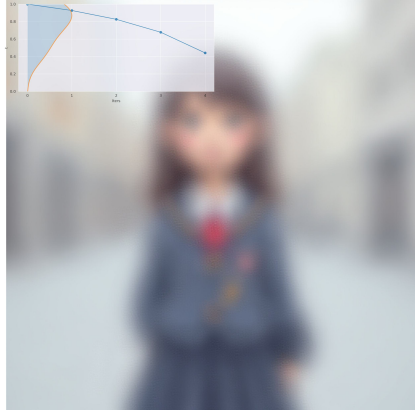


Figure 8. **Rescheduling improves few-step sampling.** Comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev, with 5 steps.

A goth anime woman with a symmetrical and attractive face in a black and white watercolor headshot art on ArtStation.



A girl in school uniform standing in the city.

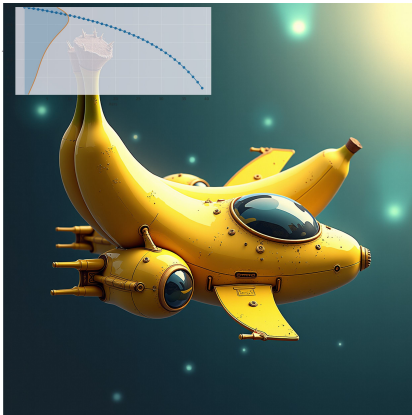


A photo of Big Chungus from Looney Tunes.



Figure 9. **Rescheduling improves few-step sampling.** Comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev, with 5 steps.

A banana spaceship reminiscent of Homeworld



An anime girl with an athletic build poses confidently while holding an assault rifle...



Clint Eastwood fighting in a white Michelin man costume with hippo...

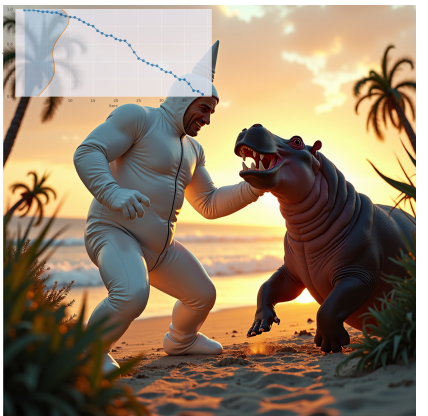
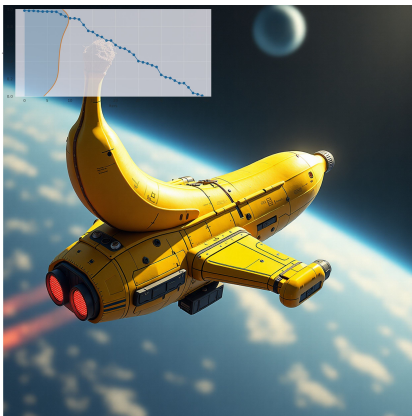
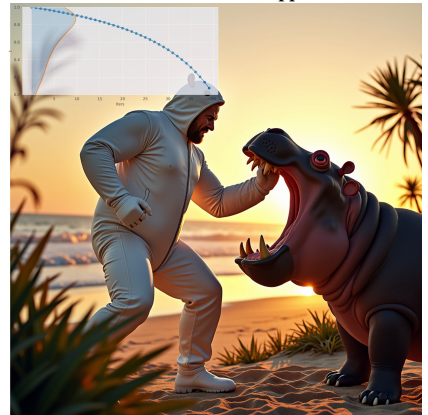


Figure 10. **Rescheduling improves general T2I alignment.** Head-to-head comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev with 40 steps.

The image depicts a muscular mouse wielding assault rifles, in a Disney art style.



Yoda performing at Woodstock.



Photo of a chocolate-type Pokemon card.

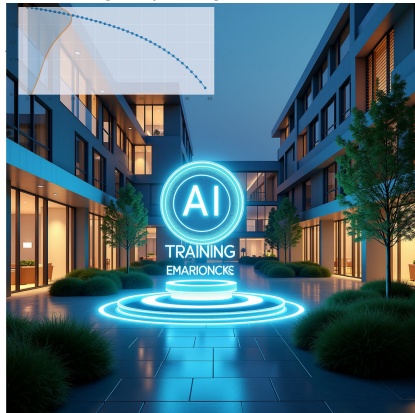


Figure 11. **Rescheduling improves general T2I alignment.** Head-to-head comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev with 40 steps.

... that reads **“Abandon All Hope”** in eerie, gothic lettering, set against a moonlit night...



... featuring a glowing **“AI Training Zone”** hologram floating in the center...



... prominently displaying **“Staff Pick Book 15”** next to a stack of books...

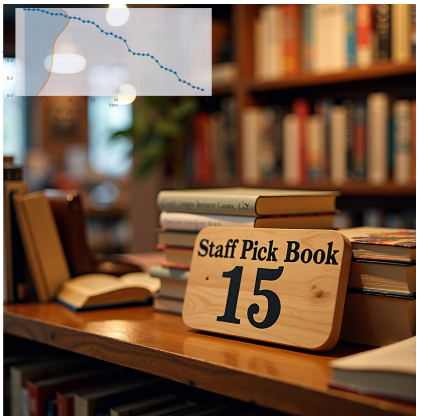
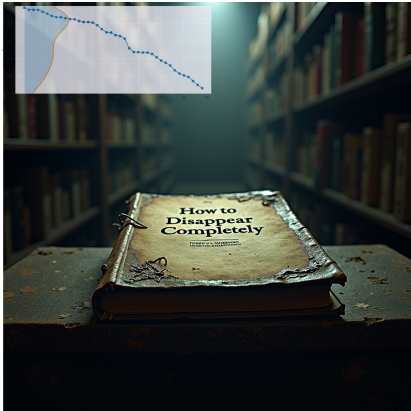
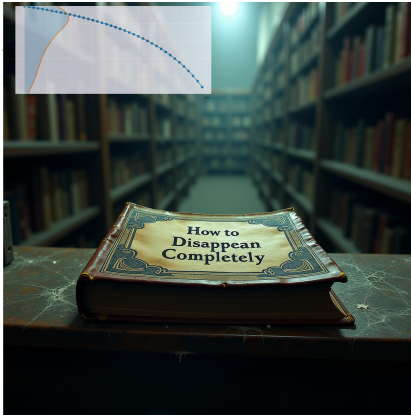
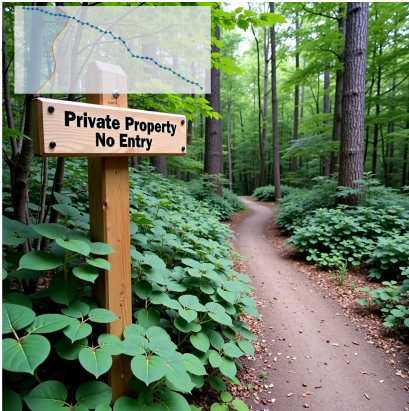


Figure 12. **Rescheduling improves text rendering.** We present comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev.

*A book saying “**How to Disappear Completely**”...*



*A hiking trail with a wooden signpost clearly displaying “**Private Property No Entry**”, surrounded by dense, green foliage...*



*A vibrant skateboard deck featuring the bold graphic “**SKATE OR DIE 4EVER**” in dynamic, graffiti-style lettering...*



Figure 13. **Rescheduling improves text rendering.** We present comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev.

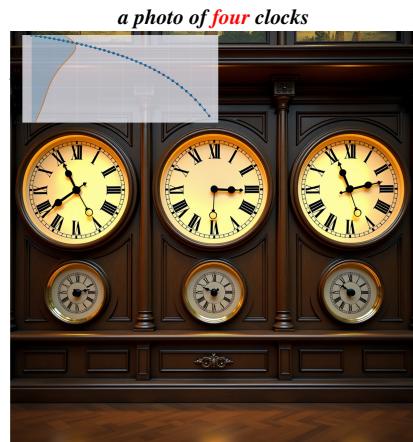


Figure 14. **Rescheduling improves fine grained alignment.** Comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev.

*a photo of two **clocks***



*a photo of a **yellow computer keyboard** and a **black sink***



*a photo of a couch **left of a toaster***

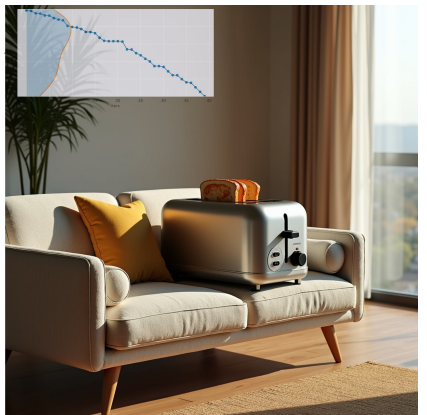


Figure 15. **Rescheduling improves fine grained alignment.** Comparisons between images generated with default schedules (upper) and our learned schedules (lower) from Flux-Dev.