

**Fatemeh Doudi**  
**PhD student, Computer Engineering**  
**Texas A&M university**

# Reinforcement Learning for Diffusion LLMs with Entropy-Guided Step Selection and Stepwise Advantages

# What is Diffusion LLMs:

---

## Diffusion LM: iterative denoising

### PROMPT

Solve: 2x + 3 = 11

### GENERATION SEQUENCE

Confidence from model  $\theta$

[M] [M] [M] [M] [M] [M] [M] [M] [M]

Step 0 — all tokens masked

Prompt token    Masked [M]    Being unmasked (high confidence)    Unmasked

# What is Diffusion LLMs:

---

- Most LLMs (GPT, LLaMA, DeepSeek) generate text left to right, one token at a time.
- Diffusion LMs do something fundamentally different:
  - a. Start with a fully masked sequence [M] [M] [M] [M]
  - b. Iteratively unmask tokens over  $T$  denoising steps
  - c. At each step, the model sees the full (partially masked) context — bidirectional attention
  - d. End result: the clean text  $x_0$

## Pros and Cons:

- a. Parallel multi-token generation → higher throughput
- b. Bidirectional context at every step → richer representations
- c. But: breaks the causal structure that RL methods rely on

# The Problem: why RL doesn't transfer to diffusion LMs

---

**The setup:** RL for AR models works because of “causal factorization”. The likelihood of a sequence decomposes token by token:

$$\log \pi(x) = \sum \log \pi(x_i | x_{<i})$$

One forward pass gives you everything. Policy gradient is straightforward.

**The problem with diffusion LMs:**

- There's no causal structure.
- The model uses bidirectional attention!
- You can't factorize  $\log \pi(x_0)$  token by token
- Sequence-level likelihood is intractable
- Naively porting GRPO → requires  $T$  separate forward passes (one per denoising step) → prohibitively expensive.

**What existing work does:** Approximate it by (i) surrogate likelihoods, (ii) mean-field approximations. They work, but they introduce bias and lose the step-level structure.

**Our question: Can we derive an exact policy gradient that respects the denoising structure, without evaluating sequence likelihood at all?**

# Our Solution

---

- **The key idea:** Instead of approximating the likelihood, we ask a different question — what's the right MDP formulation for diffusion generation?

## Step 1 - Reframe as an MDP:

- **State at step  $t$ :** the partially masked sequence  $x_t$
- **Action:** which tokens to unmask  $\rightarrow x_{t-1}$
- **Reward:** only at the end  $\rightarrow r(x_0, q)$
- This gives us an exact policy gradient that decomposes over denoising steps

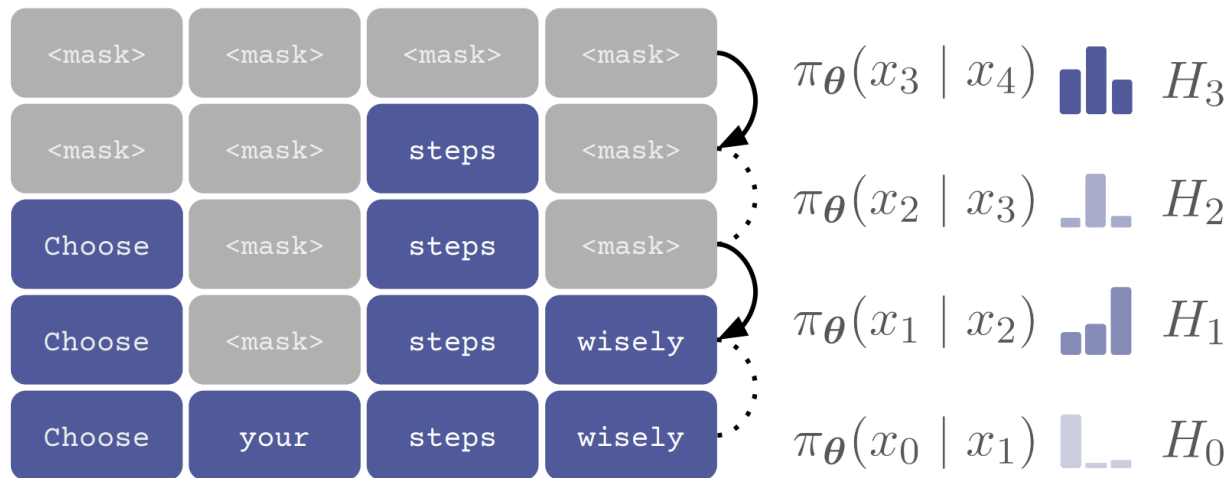
## Step 2 - EGSP0: You can't afford T forward passes. So which steps matter most?

- The ones where the model is most uncertain (highest entropy).
- Train only on those K steps.
- Provably minimizes the approximation error.

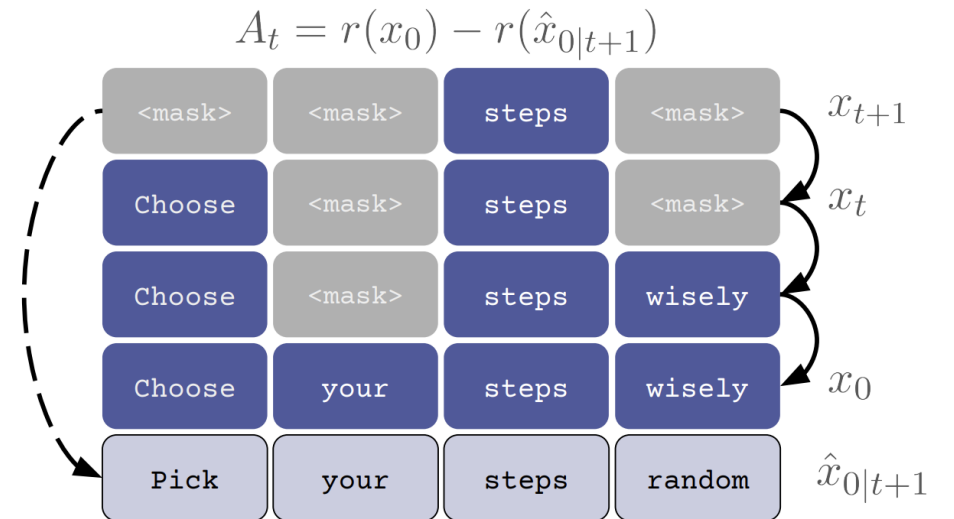
## Step 3 - EGSP0-SA: How do you estimate the advantage at each step without a value network?

- Use the model itself.
- greedily complete from  $x_{t+1}$  in one shot.
- get a reward estimate.

# Our Solution



(a) Entropy-Guided Step Selection



(b) Stepwise Advantage Estimation

# Empirical Results

Model / Seq. Len.	Sudoku				Countdown				GSM8K				MATH500			
	128	256	512	Best	128	256	512	Best	128	256	512	Best	128	256	512	Best
LLaDA-8B-Instruct	11.7	6.7	5.5	11.7	20.7	19.5	16.0	20.7	68.7	76.7	78.2	78.2	26.0	32.4	36.2	36.2
d1	22.1	16.7	9.5	22.1	34.8	32.0	42.2	42.2	73.2	81.1	82.1	82.1	<b>33.8</b>	<u>38.6</u>	<u>40.2</u>	40.2
wd1	-	76.4	62.8	76.4	-	51.2	46.1	51.2	-	80.8	82.3	82.3	-	34.4	39.0	39.0
SPG	82.9 <sup>†</sup>	94.0 <sup>†</sup>	93.1 <sup>†</sup>	94.0 <sup>†</sup>	68.8	71.5	70.3	71.5	<b>78.5</b>	<b>86.1</b>	84.5	<b>86.1</b>	<u>33.4</u>	<b>40.0</b>	<u>41.8</u>	<b>41.8</b>
d2				91.9				56.6				85.0				<u>41.6</u>
EGSPO (Ours)	<u>93.3</u>	<u>93.6</u>	<u>89.1</u>	<u>93.6</u>	<u>70.3</u>	<u>71.5</u>	<u>75.8</u>	<u>75.8</u>	<u>77.5</u>	<u>85.7</u>	<u>84.98</u>	<u>85.7</u>	32.2	37.8	39.0	39.0
EGSPO-SA (Ours)	<b>93.7</b>	<b>94.3</b>	<b>93.4</b>	<b>94.3</b>	<b>78.5</b>	<b>77.3</b>	<b>76.5</b>	<b>78.5</b>	73.5	84.6	<b>85.03</b>	85.03	33.2	38.2	39.6	39.6

*Table 1.* Main results across reasoning benchmarks. For each method, we report accuracy at different generation lengths and the best performance across lengths. d2 reports a single best performance independent of generation length, following its original evaluation protocol. wd1 does not report the evaluation on length 128. For Sudoku, <sup>†</sup> denotes 3-shot evaluation, whereas ours is 0-shot.

# Empirical Results

---

<b>HumanEval</b>	128	256	512	Best
LLaDA-8B-Instruct	27.4	35.3	37.8	37.8
d1	31.1	32.9	37.8	37.8
EGSPO (Ours)	<u>32.3</u>	<u>40.2</u>	<u>39.6</u>	<u>40.2</u>
EGSPO-SA (Ours)	<b>32.3</b>	<b>41.5</b>	<b>44.5</b>	<b>44.5</b>
<b>MBPP</b>	128	256	512	Best
LLaDA-8B-Instruct	36.2	41.2	40.4	41.2
d1	40.5	44.7	42.8	44.7
EGSPO (Ours)	<u>49.9</u>	<b>50.6</b>	<b>50.3</b>	<u>50.6</u>
EGSPO-SA (Ours)	<b>51.1</b>	<u>48.7</u>	<u>49.2</u>	<b>51.1</b>

*Table 2.* Results on coding benchmarks. We report accuracy at different generation lengths and the best performance across lengths.

# Diffusion Blend: Inference-Time Multi-Preference Alignment for Diffusion Models

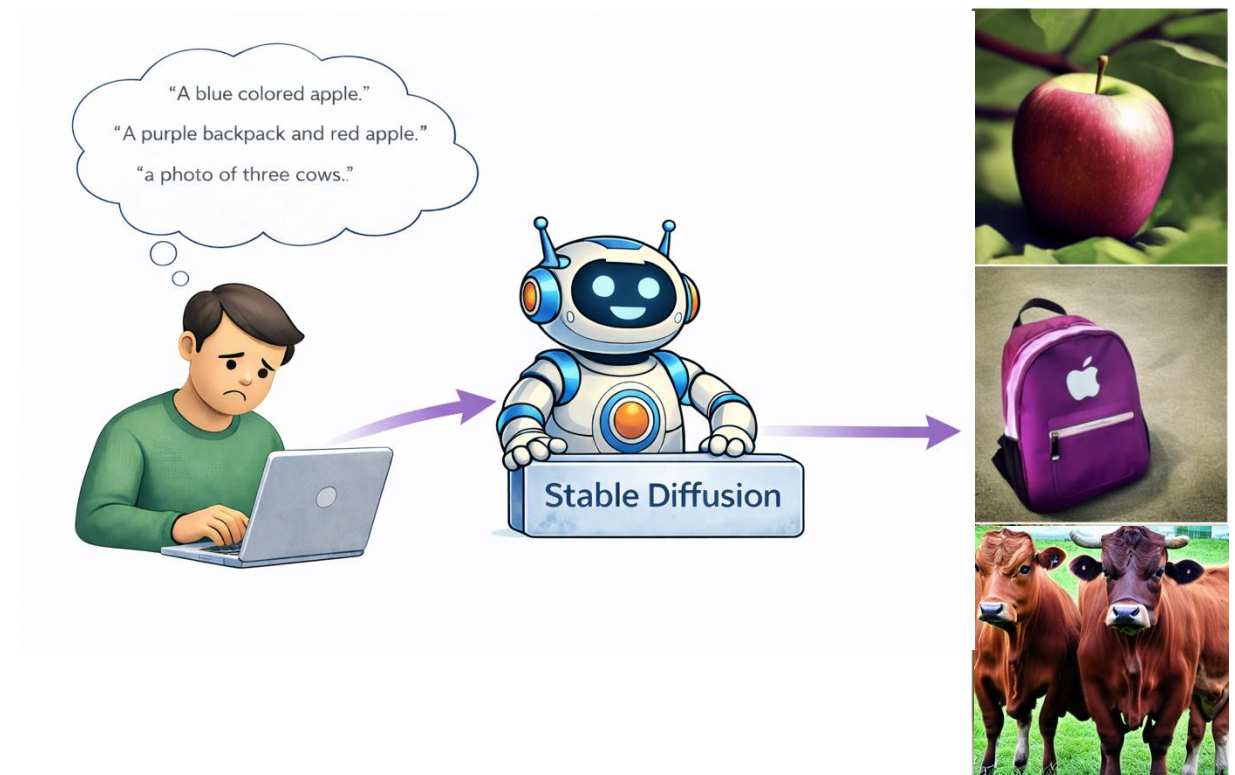
# Diffusion Models Are Powerful, But Not Aligned!

---

RL fine-tuning bridges the gap: optimize the model toward a reward signal!

Pre-trained on massive data — not optimized for downstream objectives.  
Generated images can fail on:

- Prompt fidelity
- Aesthetic quality
- Human preference



# Diffusion Models Are Powerful, But Not Aligned!

---

RL fine-tuning bridges the gap: optimize the model toward a reward signal!

Pre-trained on massive data — not optimized for downstream objectives.

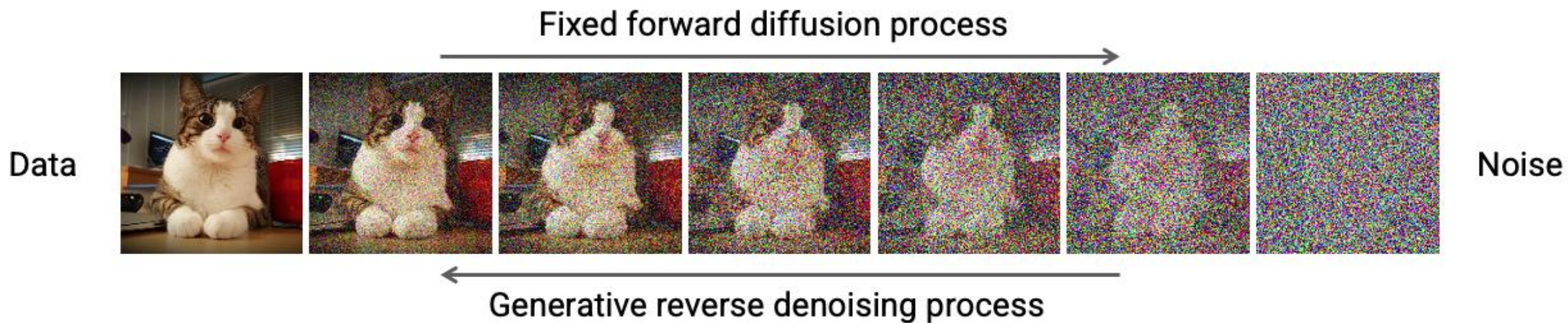
Generated images can fail on:

- Prompt fidelity
- Aesthetic quality
- Human preference

**RL fine-tuning bridges the gap: optimize the model toward a reward signal!**

# Diffusion Model in a glance

---



# Forward Process

---

**Goal: Gradually destroying structure with noise.**

Define the forward SDE:

$$dx_t = -\frac{1}{2} \beta(t)x_t dt + \sqrt{\beta(t)} dW_t$$

- $x_0$ : real data.
- $x_T \sim \mathcal{N}(0, I)$ : pure noise.
- $\beta(t)$ : noise schedule.

Closed form:  $x_t|x_0 \sim \mathcal{N}(\mu(t)x_0, \sigma^2(t)I)$

**Key insight: No learning here. This is a fixed, deterministic corruption process.**

# Denoising (Reverse) Process

---

Goal: Running time backwards -> from noise to image.

The reverse SDE (Anderson 1982):

$$dx_t = \underbrace{\left( -\frac{1}{2} \beta(t) x_t - \beta(t) \nabla_x \log p_t(x_t) \right)}_{f_{pre}(x_t, t)} dt + \sqrt{\beta(t)} dW_t$$

- Running from  $t = T$  to  $t = 0$  → generates from  $p_0$ .
- $\nabla_x \log p_t(x_t)$ : the score function.
  - We can't compute it analytically.
  - Train a neural network  $s_{\theta}(x_t, t) \approx \nabla_x \log p_t(x_t)$

# RL Alignment of Diffusion Models

---

RL concept	Diffusion equivalent
State $s_t$	Noisy image $x_t$
Action $a_t$	Next denoised image $x_{t-1}$
Policy $\pi(a   s)$	Denoising transition $p_\theta(x_{t-1}   x_t)$
Trajectory $\tau$	Full chain $x_T \rightarrow \dots \rightarrow x_0$
Reward $r$	Score on final image $r(x_0)$
Reference $\pi_{\text{ref}}$	Pretrained model $p_0^{\text{pre}}$

## RL alignment objective:

- terminal reward (alignment, aesthetics, preference)
- **After RL fine-tuning, the new drift will be:**

$$f^{(r, \alpha)}(x_t, t)$$

# Limitations of RL Alignment

---

RL alignment objective:

$$\max_{p_0} \mathbb{E} [r(x_0)] - \alpha \cdot \text{KL}(p_0 \parallel p_0^{pre})$$

## 1. Fixed Reward Function

- a. One reward  $\rightarrow$  one fine-tuned model.
- b. New objective  $\rightarrow$  retrain

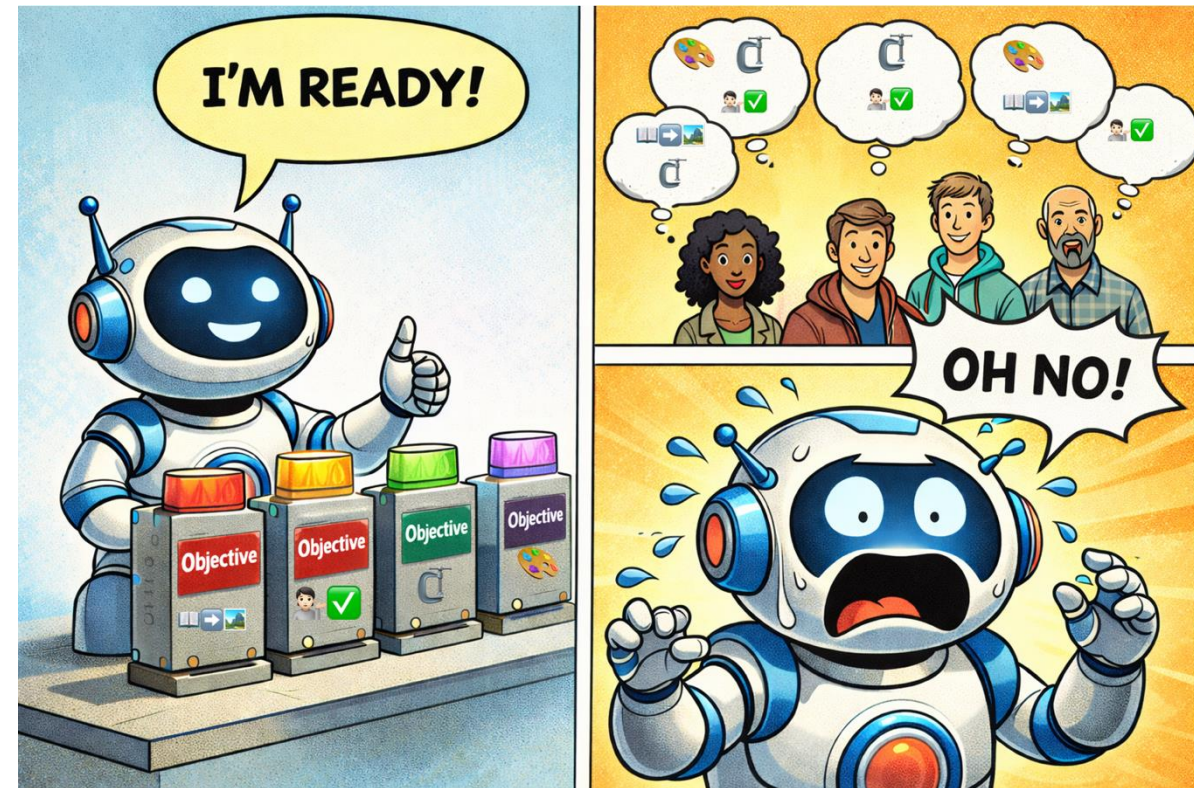
## 2. Fixed Trade-off Parameter $\alpha$

- a. Controls deviation from base model
- b. Requires tuning via grid search

# The Multi-Preference Problem

## One Fine-Tuned Model Is Never Enough!

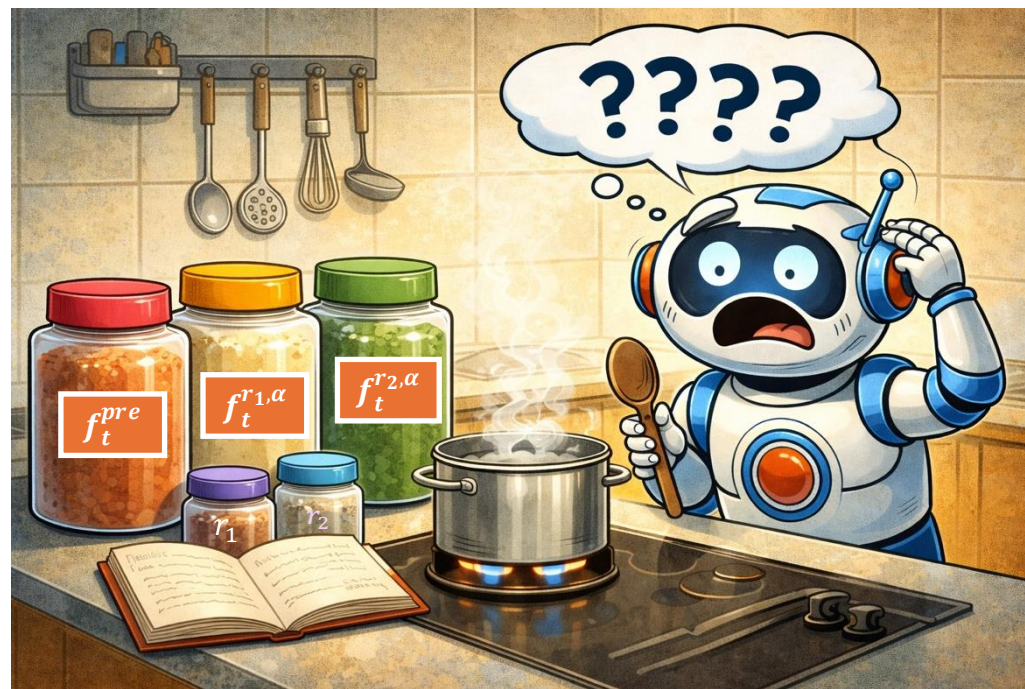
- Real users balance multiple conflicting objectives!
- Preferences vary across users and even prompts
- Fine-tuning for each preference combination?
  - a. Time-consuming
  - b. computationally expensive!



# Inference-Time Multi-Preference Alignment

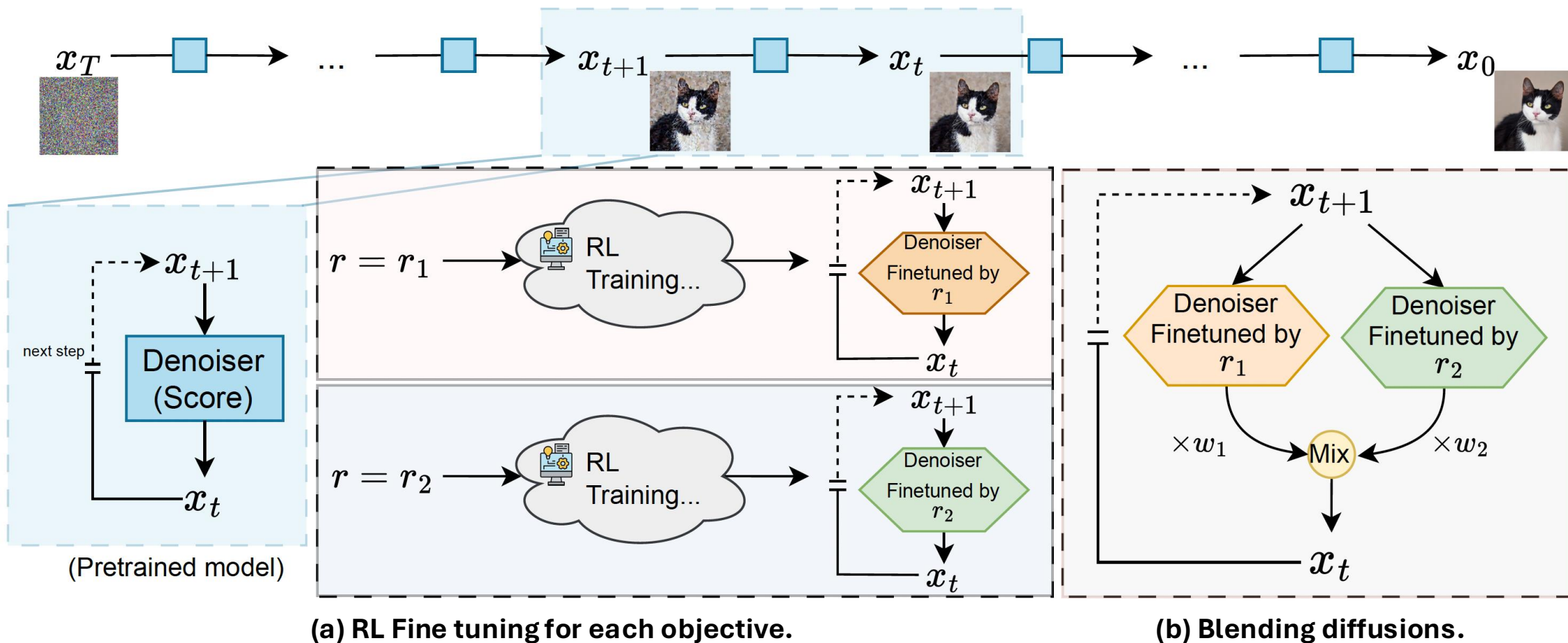
Given:

- Basis reward functions  $r_1, \dots, r_m$
- A pretrained diffusion model  $f_t^{pre}$
- RL-aligned models for each reward  $f_t^{r_i, \alpha}$



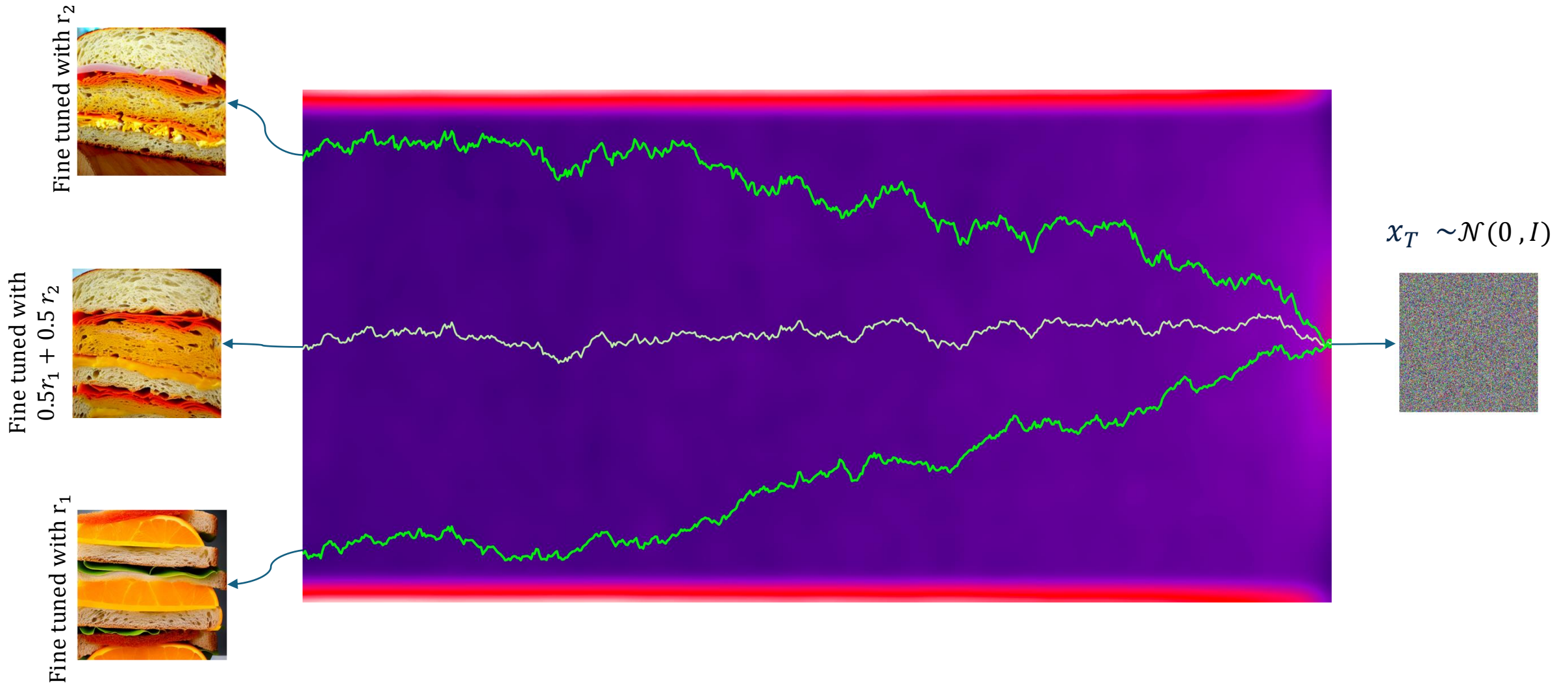
Can we generate samples aligned with any user-specified preference  $r(w) = \sum_{i=1}^m w_i r_i$  and regularization level  $\alpha$ , at inference time, without additional fine-tuning?

# Diffusion Blend Overview



Takeaway:  $f_t^{blend} \approx \sum_i w_i f_t^{r_i, \alpha}$ .

# Blending the diffusion



$r_1$  denotes the aesthetics reward and  $r_2$  denotes the text-image alignment reward.

# Drift Blending Theorem

---

**Proposition 1.** *The aligned drift decomposes as:*

$$f^{r,\alpha}(x_t, t) = f^{pre}(x_t, t) - \beta(t)u^{(r,\alpha)}(x_t, t),$$

where the control term is:

$$u^{(r,\alpha)}(x_t, t) = \nabla_{x_t} \log E_{(x_0 \sim p_{(0|t)}^{pre}(\cdot|x_t))} \left[ \exp \left( \frac{r(x_0)}{\alpha} \right) \right].$$

**Jensen Approximation (interchange E and exp) :**

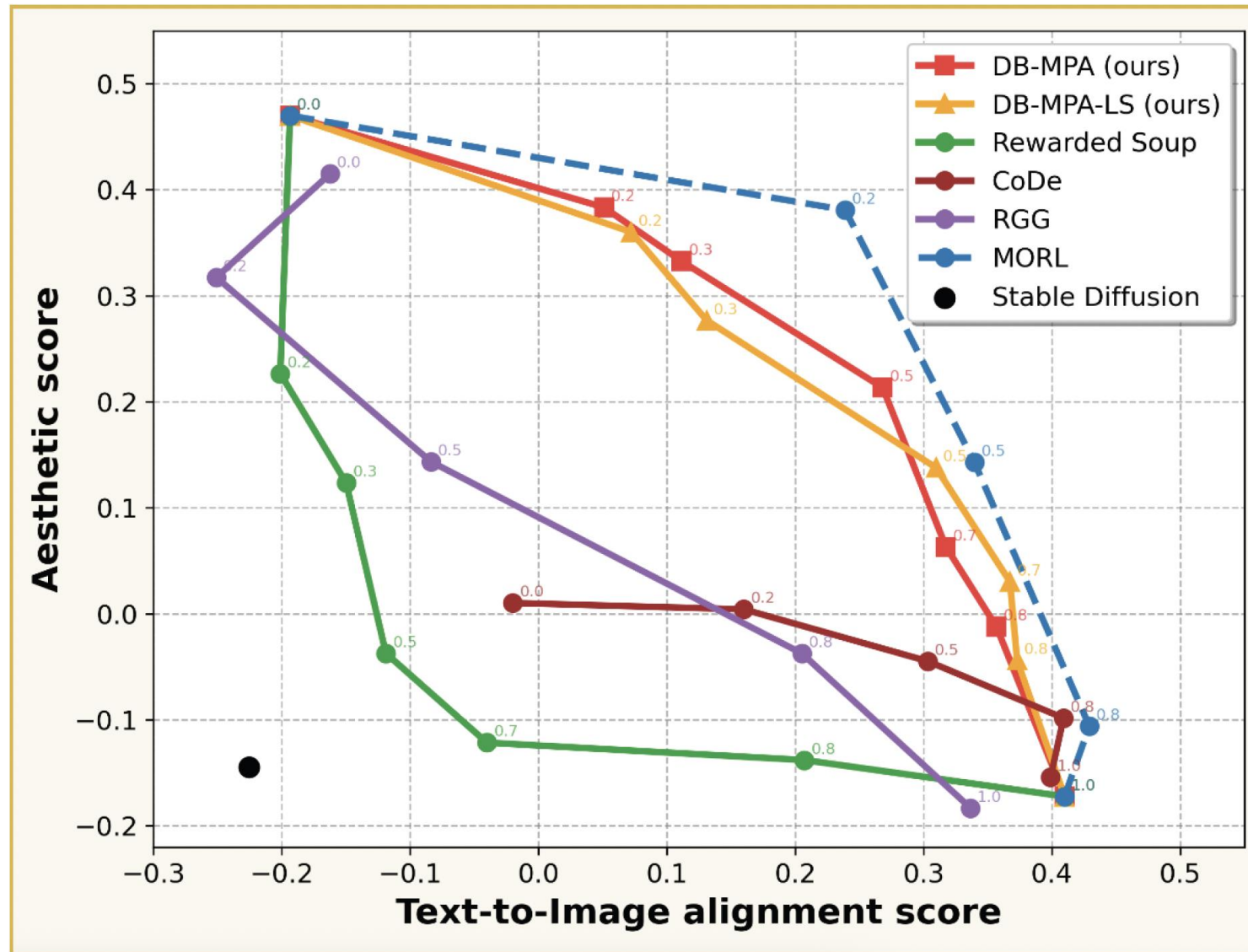
$$u^{(r,\alpha)} \approx \nabla_x \mathbb{E} \left[ \frac{r(x_0)}{\alpha} \right]$$

**Key consequence — Multi-Reward Setting:**

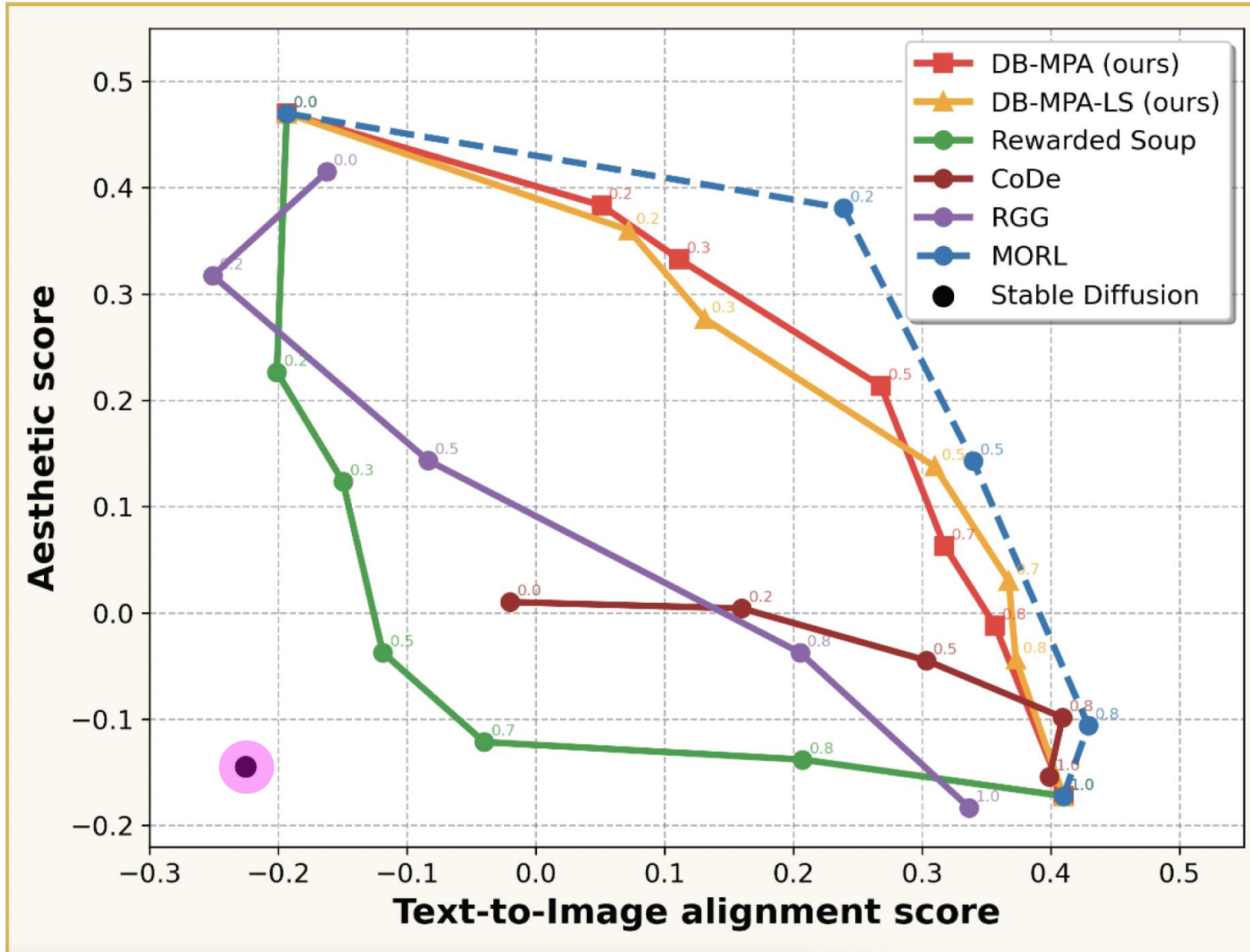
$$f^{(r(w),\alpha)} = \sum_i w_i f^{(r_i,\alpha)}$$

Where  $r(w) = \sum_i w_i r_i$ ,

# Results: Pareto Front



# Results: Visual Comparison



“A blue colored apple.”



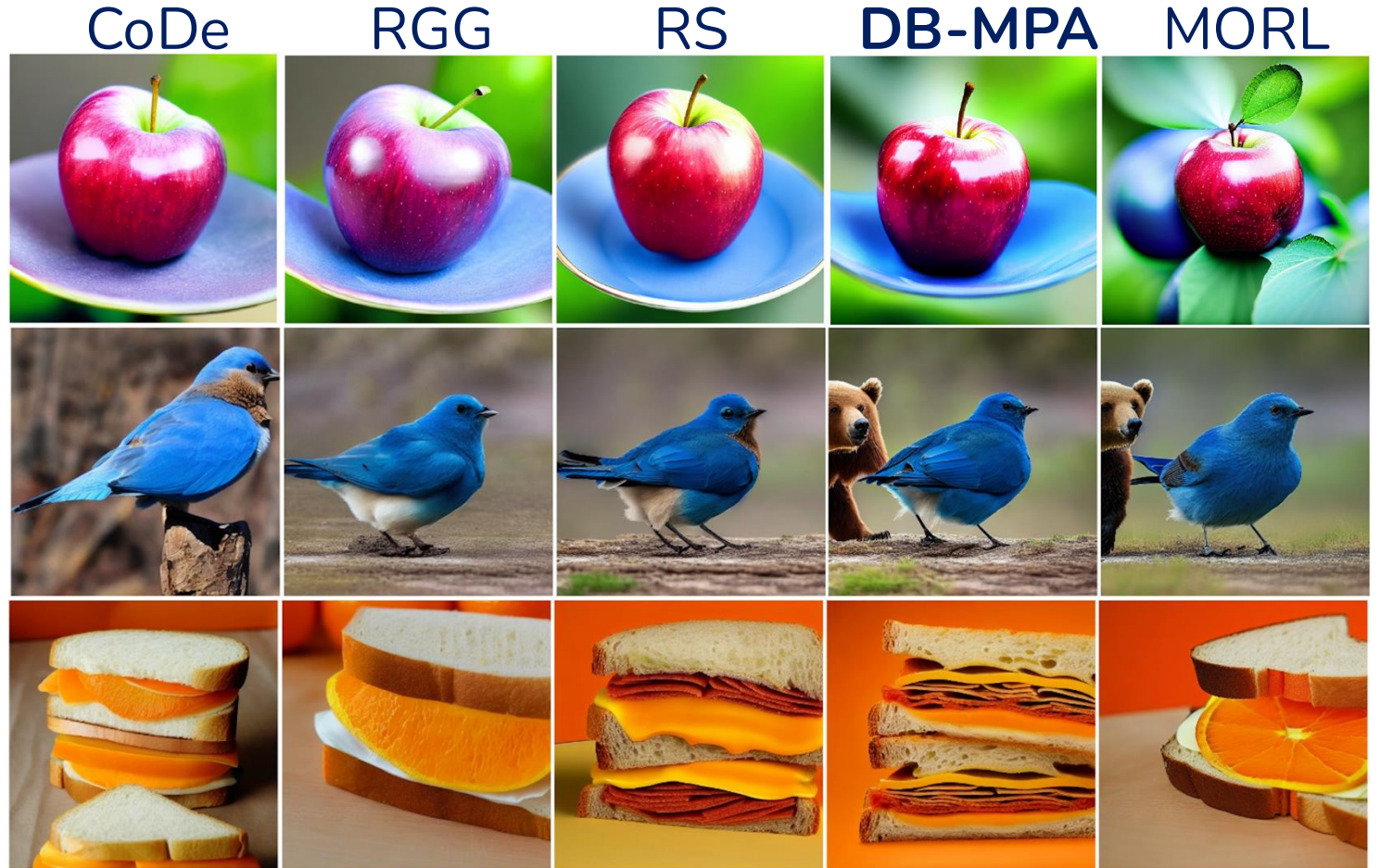
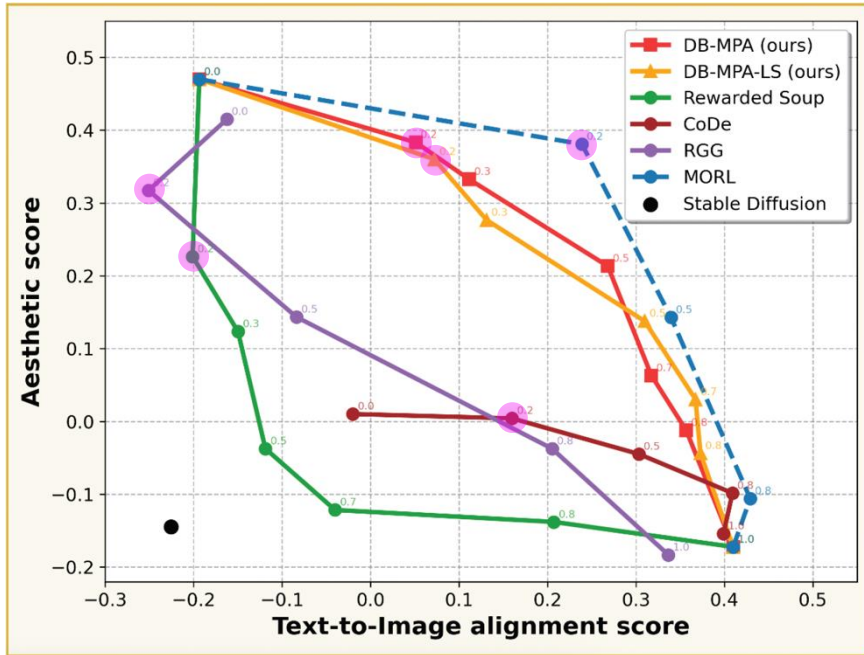
“A blue bird and a brown bear.”



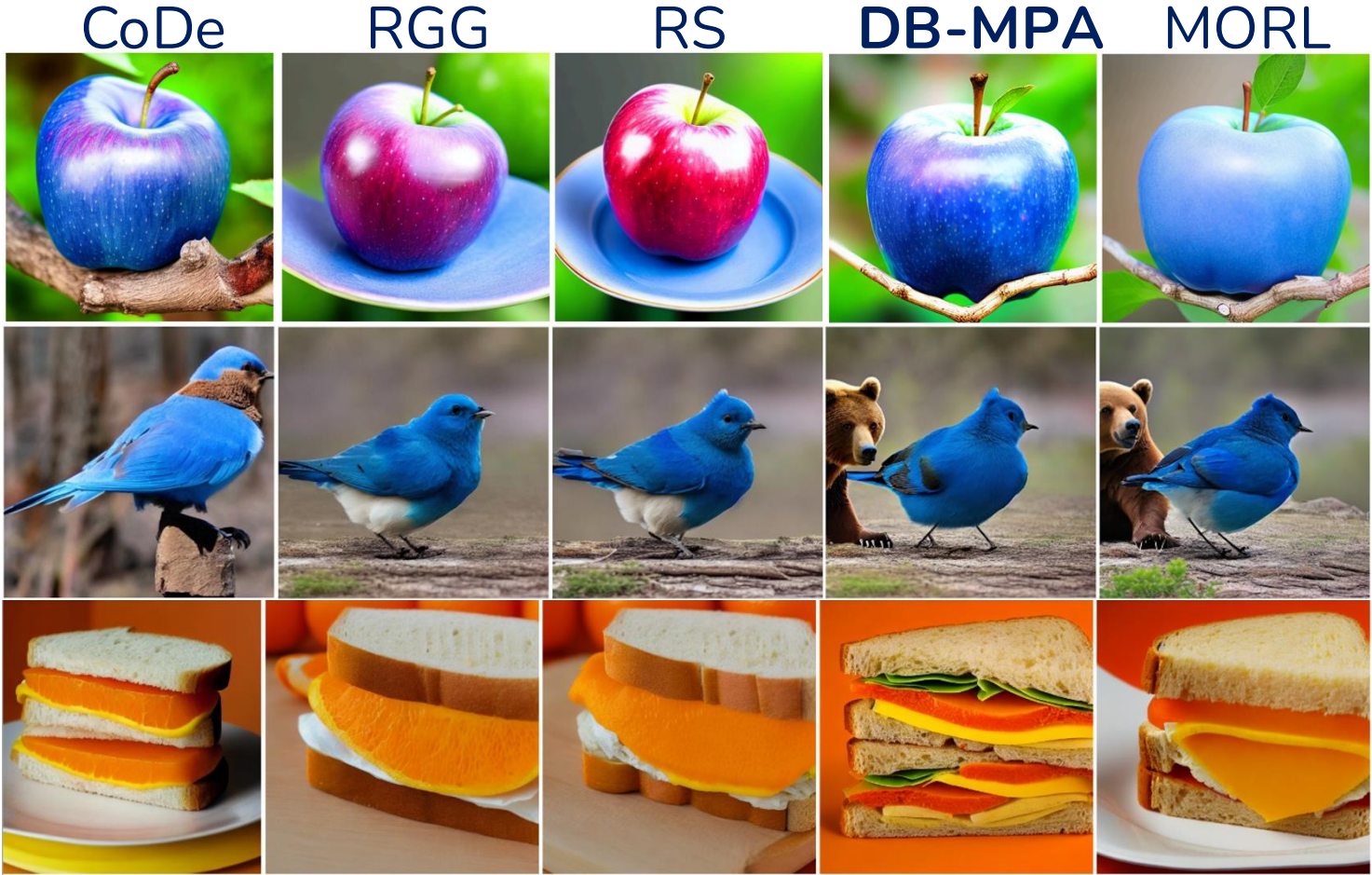
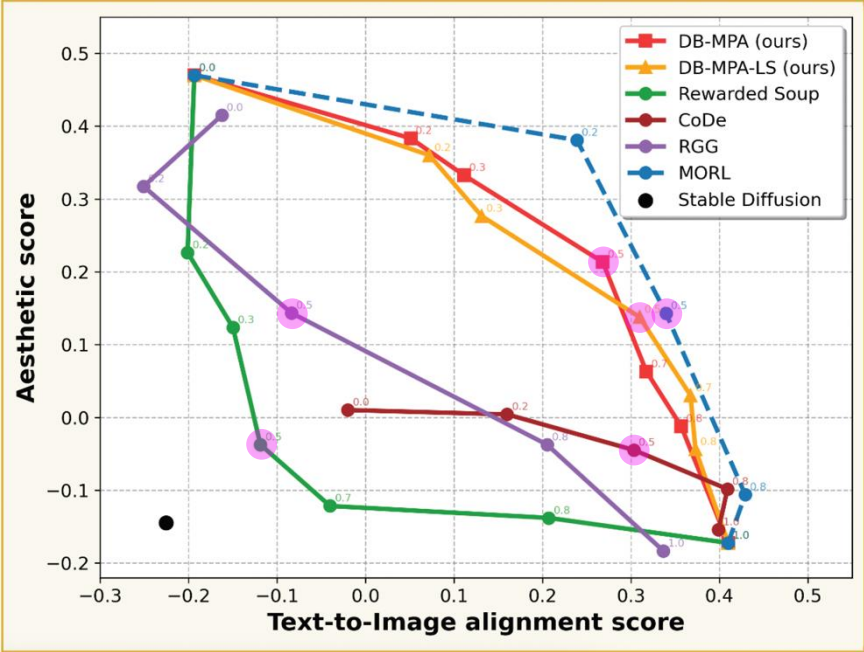
“An orange-colored Sandwich.”



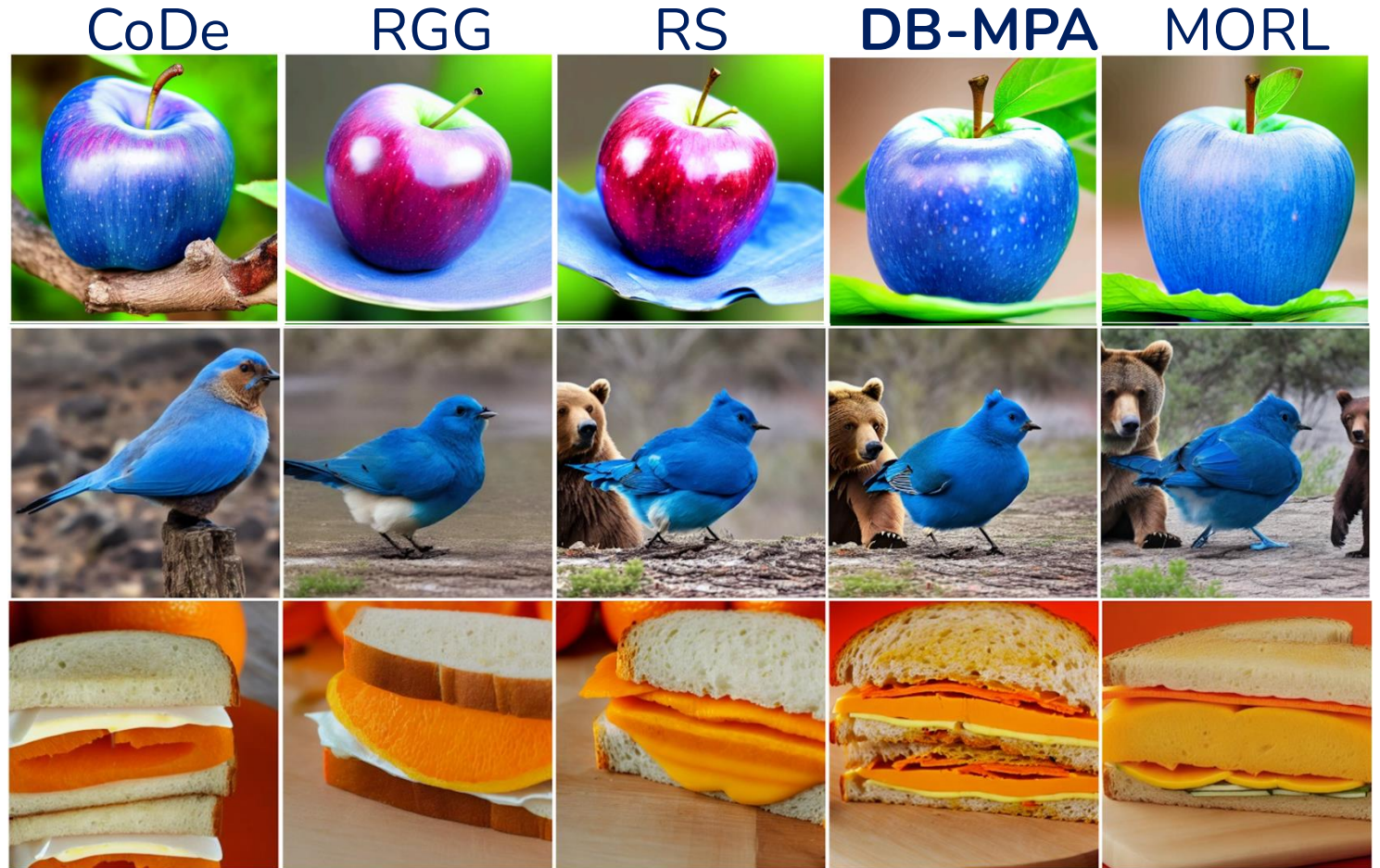
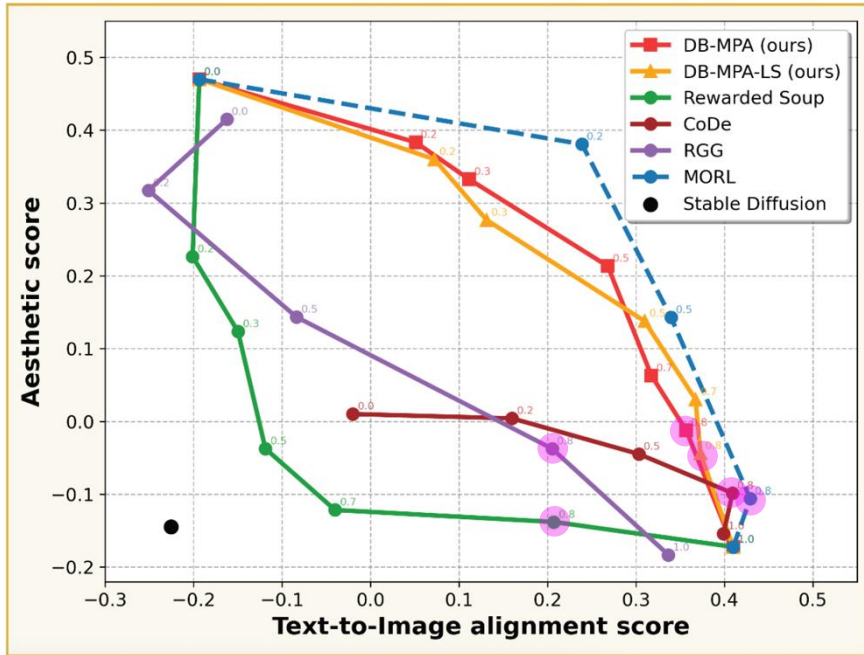
# Results: Visual Comparison



# Results: Visual Comparison



# Results: Visual Comparison



# Controlling KL Regularization at Inference Time

---

Users may also want to control regularization strength:

$$\alpha(\lambda) = \frac{\alpha}{\lambda}$$

**Problem:**

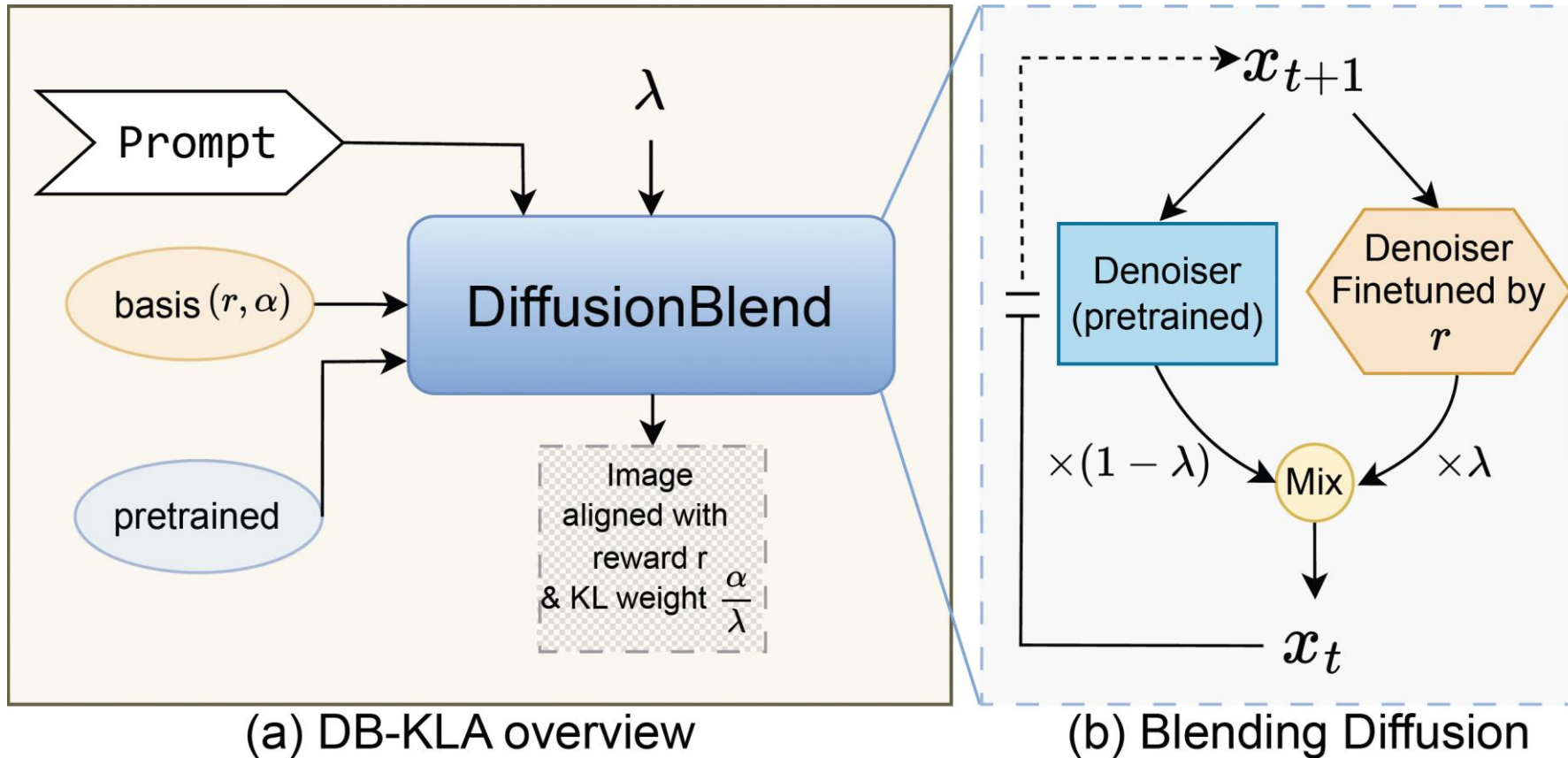
- Can we approximate models for arbitrary  $\lambda$  at inference time without additional fine-tuning?

**Solution:**

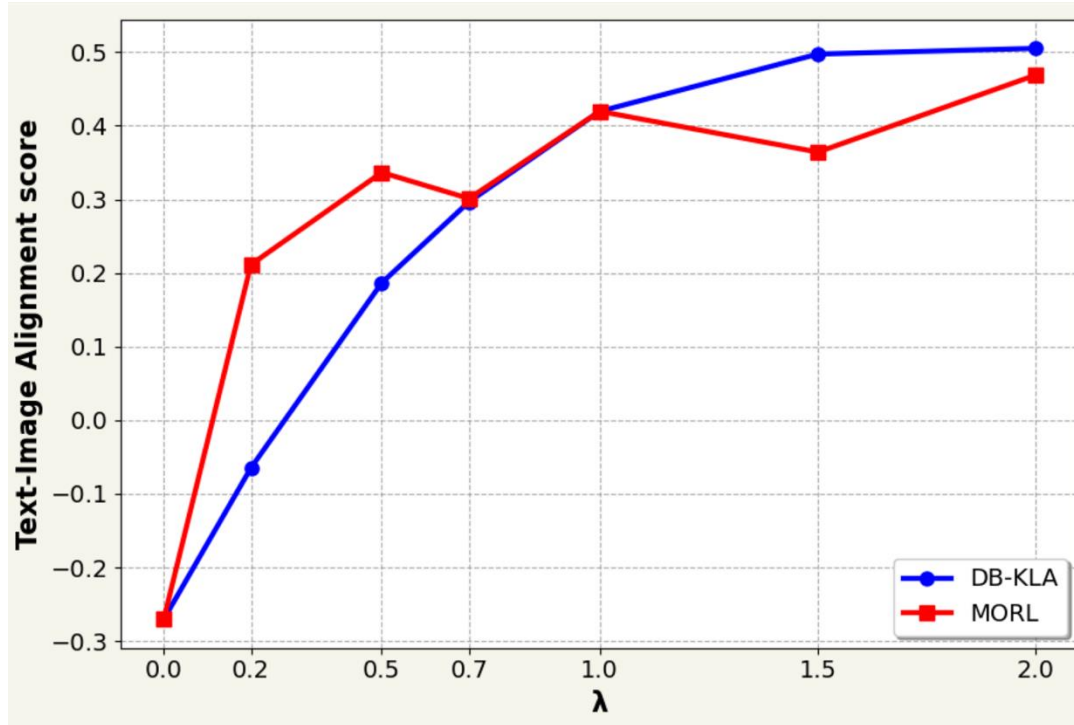
We show that  $f^{(r, \alpha(\lambda))} \approx (1 - \lambda)f^{pre} + \lambda f^{(r, \alpha)}$

# DB-KLA

Users may also want to control regularization strength  $\alpha(\lambda) = \frac{\alpha}{\lambda}$ .



# Results



Stable Diffusion



'A red apple and a purple backpack.'

KL weight =  $5\alpha$    KL weight =  $\alpha$    KL weight =  $0.5\alpha$



Retained Model



DB-KLA

**PowerMamba: A Deep State Space Model  
and Comprehensive Benchmark for Time  
Series Prediction in Electric Power  
Systems**

# The sequence modeling Evolution:

---

**RNNs/LSTMs:** sequential, can't parallelize, vanishing gradients over long sequences.

**Transformers:** parallel, great at long-range dependencies, but attention is  $O(n^2)$  which means for a 240-hour context window, that's 57,600 pairwise comparisons.

**SSMs:** model sequences as a dynamical system  $h_t = \bar{A}h_{t-1} + \bar{B}x_t$ .

- Fixed parameters  $\rightarrow$  output equals a convolution  $y = K * x \rightarrow$  parallel  $O(n)$  training.
- At inference, just carry one hidden state forward  $\rightarrow O(1)$  per step.

**Mamba: adds selectivity: A, B, C become input-dependent  $\rightarrow$  the model learns what to remember and what to forget dynamically**

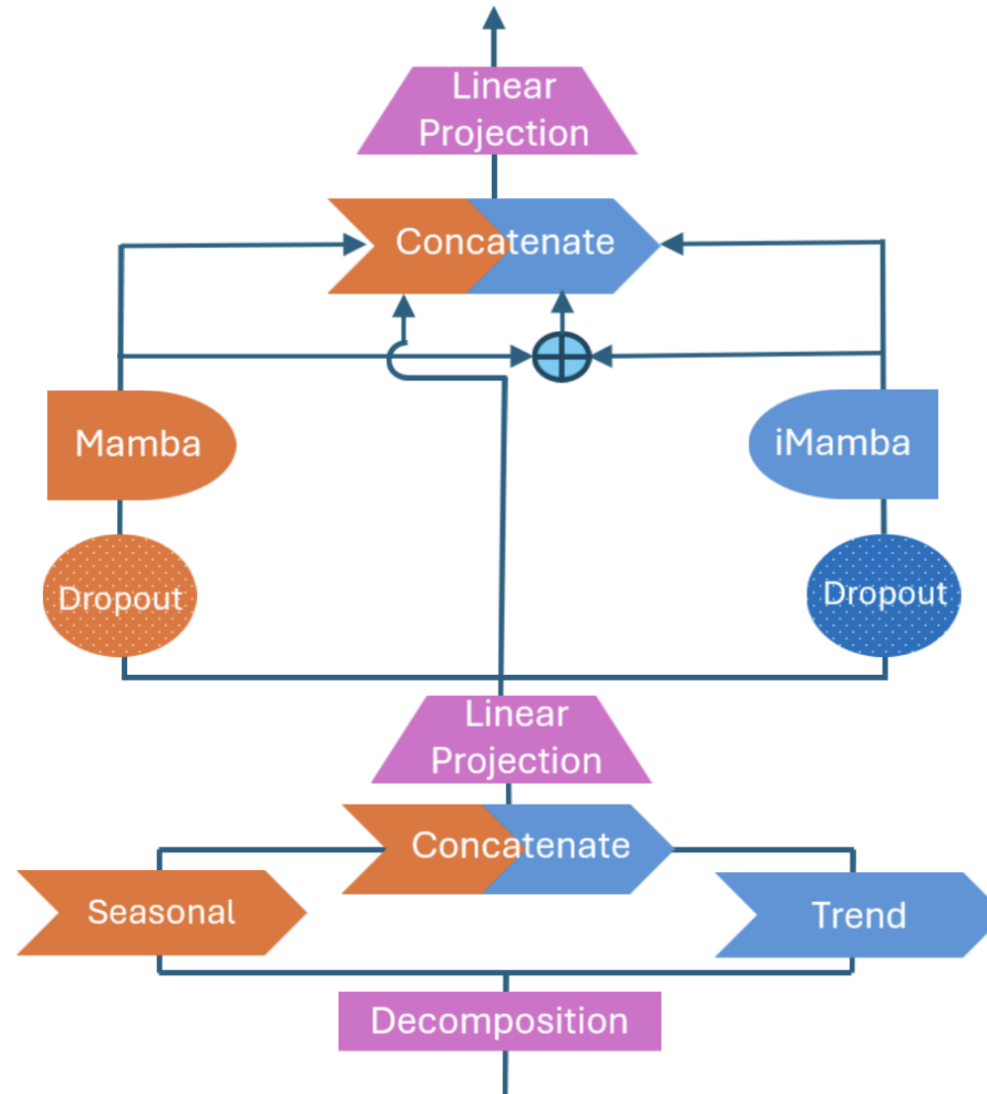
# Why Mamba for Power Grids?

---

- Power grid variables, load, price, wind, solar, evolve like **physical dynamical systems** governed by differential equations.
- Grid forecasting needs long context windows (240+ hours) to capture daily/weekly cycles → Transformer's  $O(n^2)$  becomes prohibitive.
- Attention computes pairwise interactions across all time steps  $\times$  all channels → for 22 channels  $\times$  240 hours = 5,280 tokens, attention matrix is  $5,280 \times 5,280 = 27.8M$  operations.
- To manage this, most Transformer forecasting models (like iTransformer) either attend across time or across channels, not both simultaneously, so they still miss some interactions
- Mamba's dual-path architecture (Mamba + iMamba in parallel) handles both in one forward pass: one path sees  $L \times D$  (time within each channel), the other sees  $D \times L$  (across channels at each time step) → explicit joint modeling at linear cost

# PowerMamba architecture

---



# PowerMamba architecture

---

## 1. Decomposition:

- split each time series into trend + seasonal components.
- helps with highly volatile series like wind/solar.

## 2. Dual-path Mamba + iMamba: two parallel Mamba blocks:

- Mamba sees  $L \times D \rightarrow$  captures **temporal patterns within each channel**
- iMamba sees  $D \times L$  (transposed)  $\rightarrow$  captures **cross-channel dependencies**

## 3. Fixed-size projections:

- linear projection maps input to fixed embedding size **E** before Mamba blocks.
- Model size doesn't grow with context length or prediction window.

## 4. External forecast integration:

- plugs in ERCOT's day-ahead forecasts for load and renewables without inflating model size.
- 76% error reduction on renewables

# GridSet

---

- 5 years of hourly data (2019–2023), ERCOT grid, Texas.
- 22 channels: 8 zonal loads, 8 zonal electricity prices, 4 ancillary service prices, 2 renewable generation (wind + solar).
- Extended version: 262 channels including ERCOT's day-ahead external forecasts for load and renewables.
- Captures real-world events: COVID demand drop, Winter Storm Uri, rapid renewable growth in Texas.

## **Why it matters:**

- First open-access unified benchmark for multivariate power system forecasting
- Includes open-access toolbox with PowerMamba + 6 baseline models ready to benchmark

# Results

TABLE III: Comparison of the prediction accuracy of our model and the baselines with and without external predictions. A fixed context size of  $L = 240$  and a prediction window size of  $W = 24$  are used for all the baselines.

Methods→	PowerMamba	TimeMachine	iTransformer	PatchTST	DLinear	TimesNET	Autoformer	
$\mathcal{D}$	PR	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	MSE MAE	
GridSet	✗	<b>0.129</b> <b>0.166</b>	<u>0.135</u> <u>0.166</u>	0.147 0.181	0.142 0.183	0.168 0.199	0.175 0.198	0.301 0.311
GridSet	✓	<b>0.074</b> <b>0.123</b>	<u>0.080</u> <u>0.126</u>	0.091 0.144	0.109 0.157	0.113 0.166	0.109 0.148	0.113 0.170
Load	✗	<b>0.098</b> <b>0.223</b>	<u>0.102</u> <u>0.227</u>	0.108 0.237	0.113 0.240	0.134 0.268	0.136 0.275	0.238 0.365
Load	✓	<b>0.065</b> <b>0.183</b>	<u>0.071</u> <u>0.188</u>	0.091 0.215	0.123 0.244	0.114 0.240	0.086 0.215	0.090 0.219
Price	✗	<b>0.100</b> <u>0.082</u>	0.108 <b>0.079</b>	<u>0.107</u> 0.090	0.109 0.102	0.120 0.101	0.137 0.090	0.215 0.194
Price	✓	<b>0.093</b> <b>0.075</b>	0.099 <u>0.075</u>	<u>0.096</u> 0.082	0.107 0.085	0.113 0.095	0.122 0.079	0.133 0.117
ASPrice	✗	<b>0.019</b> <b>0.029</b>	<u>0.020</u> <u>0.029</u>	0.021 0.033	0.022 0.038	0.022 0.037	0.028 0.032	0.089 0.153
ASPrice	✓	<b>0.018</b> <u>0.029</u>	<u>0.018</u> <b>0.026</b>	0.018 0.030	0.019 0.029	0.020 0.038	0.028 0.032	0.033 0.076
Renewables	✗	<b>0.590</b> <b>0.545</b>	<u>0.610</u> <u>0.545</u>	0.712 0.620	0.623 0.570	0.783 0.638	0.786 0.659	1.317 0.878
Renewables	✓	<b>0.142</b> <b>0.267</b>	<u>0.162</u> <u>0.279</u>	0.218 0.332	0.241 0.353	0.299 0.406	0.314 0.392	0.282 0.373

# Results

---

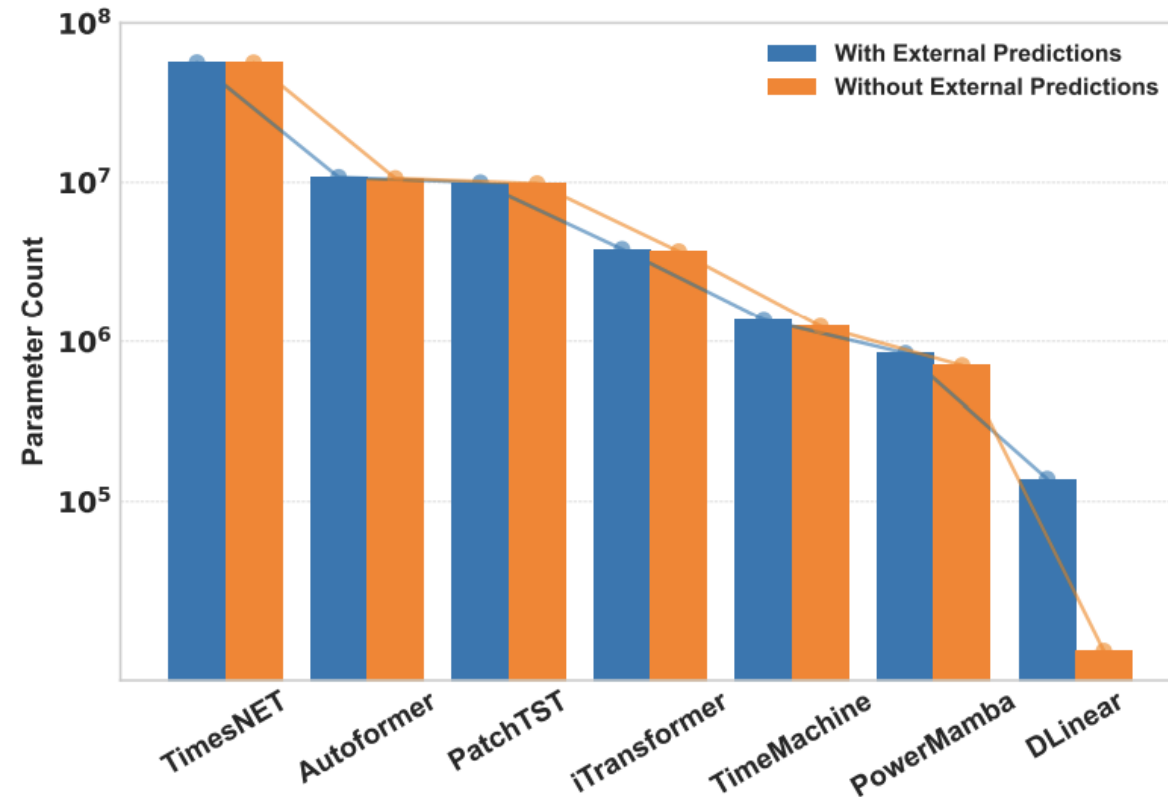


Fig. 5: The number of model parameters with and without external predictions, in log scale.