



# Information Theoretic Text-to-Image Alignment

**Chao Wang**, Giulio Franzese, Alessandro Finamore, Massimo Gallo, Pietro Michiardi

ICLR 2025

# Problem

For existing text-to-image (T2I) diffusion models, precise alignment between given texts and generated images is still challenging (e.g., **attribute binding**, **missing object**).



A **round** bag and  
a rectangular **wallet**



A man on the top of  
**a turtle**

**Research question:** Would it be possible to use Mutual Information (MI) between image and text to measure and guide the alignment?

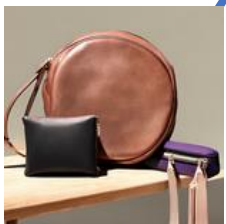
**Prerequisites:**

- 1) MI formula supported by discrete-time diffusion model
- 2) Is MI a meaningful signal for T2I alignment?

# Prereq1: Point-wise MI estimation

Given 2 R.V.  $\mathbf{z}, \mathbf{p}$  sampled from the joint distribution  $p_{latent, prompt}$ ,

$$I(\mathbf{z}, \mathbf{p}) = \mathbb{E}_{t, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\kappa_t \|\epsilon_{\theta}(\mathbf{z}_t, \mathbf{p}, t) - \epsilon_{\theta}(\mathbf{z}_t, \emptyset, t)\|^2], \quad \kappa_t = \frac{\beta_t T}{2\alpha_t(1 - \bar{\alpha}_t)}.$$



“A round bag and a rectangular wallet”

## Discrete-time diffusion model:

- Forward process:  $q(\mathbf{z}_{0:T}, \mathbf{p}) = q(\mathbf{z}_0, \mathbf{p}) \prod_{t=1}^T q(\mathbf{z}_t | \mathbf{z}_{t-1})$

Hand-crafted transition kernel:  $q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{p}) = q(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{z}_t; \sqrt{1 - \beta_t} \mathbf{z}_{t-1}, \beta_t \mathbf{I})$

- Backward process:  $p_{\theta}(\mathbf{z}_{0:T} | \mathbf{p}) = p(\mathbf{z}_T) \prod_{t=1}^T p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{p})$

Learnable transition kernel:  $p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{p}) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_{\theta}(\mathbf{z}_t, \mathbf{p}), \beta_t \mathbf{I})$ , with  $\boldsymbol{\mu}_{\theta}(\mathbf{z}_t, \mathbf{p}) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_{\theta}(\mathbf{z}_t, t, \mathbf{p}) \right)$

The unguided version  $p_{\theta}(\mathbf{z}_{t-1} | \mathbf{z}_t, \emptyset)$  was optimized at the same time to improve CFG sampling.

# Prereq2: MI is a meaningful signal for alignment

Comparison between MI and well-established alignment metrics (BLIP-VQA<sub>[1]</sub>, HPS<sub>[2]</sub>)

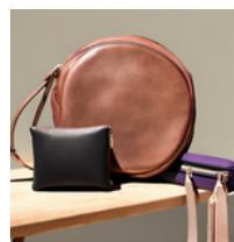
## 1. Kendall rank correlation coefficient

- good agreement between MI and BLIP-VQA ( $\tau = 0.4$ )
- strong agreement between MI and HPS ( $\tau = 0.68$ )

## 2. Human preference among the top-ranked images

- MI for 69.1%
- BLIP-VQA for 73.5%
- HPS for 52.2%

Shape binding:  
“A round bag and  
a rectangular  
wallet”



BLIP-VQA = 0.82  
HPS = 0.262  
MI = 18.61



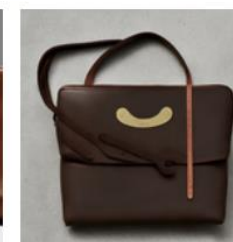
BLIP-VQA = 0.64  
HPS = 0.247  
MI = 17.16



BLIP-VQA = 0.27  
HPS = 0.262  
MI = 14.84



BLIP-VQA = 0.24  
HPS = 0.216  
MI = 12.50



BLIP-VQA = 0.01  
HPS = 0.160  
MI = 11.57

high scores = good alignment

low scores = poor alignment

[1] T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation, NeurIPS23

[2] Human Preference Score v2: A Solid Benchmark for Evaluating Human Preferences of Text-to-Image Synthesis

# MI-TUNE: Self-supervised fine-tuning

## Method:

1. For each prompt, generate 50 images
2. Filter the images with the highest  $I(\text{image}, \text{text})$
3. Finetune DoRA on these better aligned samples → improve T2I alignment

*This process can be repeated for multiple rounds.*

**Advantage:** no extra information needed

## Algorithm 2: Point-wise MI Estimation

**Input** : Pre-trained model:  $\epsilon_\theta$ ; Prompt:  $p$   
**Output** : Generated latent:  $z$ ; Point-wise MI:  $I(z, p)$

```
1 Function PointWiseMI ( $\epsilon_\theta, p$ ):  
    // Initial latent sample  
2     $z_T \sim \mathcal{N}(0, I)$  for  $t$  in  $T, \dots, 0$  do  
        // MI estimation (Eq. (2))  
3         $I(z_t, p) +=$   
             $[\kappa_t || \epsilon_\theta(z_t, p, t) - \epsilon_\theta(z_t, \emptyset, t) ||^2]$   
        // Noise sample  
4         $w \sim \mathcal{N}(0, I)$  if  $t > 1$ , else  $w = 0$   
        // Sampling step  
5         $z_{t-1} =$   
             $\frac{1}{\sqrt{\alpha_t}} \left( z_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \epsilon_\theta(z_t, p, t) \right) + \sigma_t w$   
6    end  
7    return  $z, I(z, p)$ 
```

# Experimental results

- Benchmark on T2I-CompBench prompts: 6 categories
- MI-TUNE achieves new SOTA text-image alignment.

Alignment results (%)

		BLIP-VQA							HPS							Human ( <i>user study</i> )						
	Method	Color	Shape	Texture	2D-Sp.	Non-Sp.	Compl.	(avg)	Color	Shape	Texture	2D-Sp.	Non-Sp.	Compl.	(avg)	Color	Shape	Texture	2D-Sp.	Non-Sp.	Compl.	(avg)
Infer.	SD-2.1-base	49.65	42.71	49.99	15.77	66.23	50.53	(45.81)	27.64	24.56	24.99	27.50	26.66	25.70	(26.17)	29.76	11.90	40.48	35.71	66.67	29.76	(35.71)
	A&E	61.43	47.39	64.10	16.18	66.21	51.69	(51.17)	28.44	24.43	25.88	28.42	26.60	25.60	(26.56)	31.95	15.48	52.38	32.14	65.48	30.95	(38.06)
	SDG	47.15	45.24	47.13	15.25	66.17	47.41	(44.72)	27.25	24.40	24.71	27.10	26.12	25.83	(25.90)	26.19	15.48	38.10	38.10	61.90	29.76	(34.92)
	SCG	49.82	43.28	50.16	16.31	66.60	51.07	(46.21)	27.86	24.85	25.57	27.76	26.98	26.03	(26.51)	20.24	11.90	33.33	40.48	69.05	39.29	(35.71)
FT	DPOK	53.28	45.63	52.84	17.19	66.95	51.97	(47.98)	28.20	24.99	25.44	28.12	26.80	25.88	(26.57)	23.81	16.67	47.62	34.52	70.24	38.10	(38.49)
	GORS	53.59	43.82	54.47	15.66	67.47	52.28	(47.88)	28.15	24.79	25.56	27.90	26.88	26.07	(26.56)	34.52	14.29	48.81	36.90	65.48	30.95	(38.49)
	HN-ITM	46.51	39.99	48.78	15.24	65.31	49.84	(44.28)	26.90	24.33	24.63	27.15	25.40	25.22	(25.60)	23.81	19.05	30.95	20.24	47.62	23.81	(27.58)
	MI-TUNE	65.04	50.08	65.82	18.51	67.77	54.17	(53.56)	29.13	25.57	26.20	28.50	27.15	26.70	(27.21)	46.43	25.01	53.19	45.24	73.81	46.43	(48.35)

SD-2.1-base DPOK GORS HN-ITM A&E SDG SCG MI-TUNE



(Color) “a blue bear and a brown boat”



# Experimental results

## SD-XL

Method	BLIP-VQA					
	Color	Shape	Texture	2D-Sp.	Non-Sp.	Comp.
(ref) SDXL	60.78	49.70	55.78	21.02	68.16	52.68
SD-2.1-base	49.65	42.71	49.99	15.77	66.23	50.53
MI-TUNE	69.66	55.86	66.74	22.18	72.17	57.74
MI-TUNE $\square$ (ref)	8.88	6.16	10.96	1.16	4.01	5.06
MI-TUNE % (ref)	14.61	12.39	19.65	5.52	5.88	9.61

SDXL

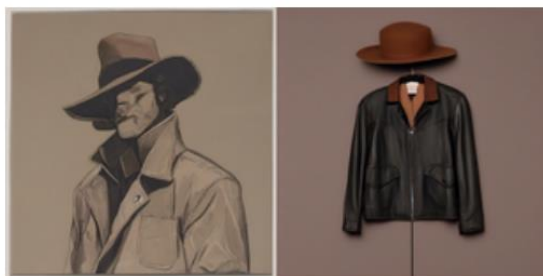
MI-TUNE



(Color) “A green apple and a brown horse”

SDXL

MI-TUNE



“A black jacket and a brown hat”

## Human Prompts

Model	HPS
SD-2.1-base	23.99
DiffusionDB	24.35
MI-TUNE	25.32
MI-TUNE $\square$ base	1.33
MI-TUNE $\square$ DiffusionDB	0.97

SD-2.1-base

Fine-tuned using  
DiffusionDB images

MI-TUNE



(Human prompt) “Child’s body with a radioactive jellyfish as a head, realistic illustration, backlit, intricate, indie studio, fantasy, rim lighting, vibrant colors, emotional”

# What about Rectified Flows (e.g., SD3)?

Check out our ICLR 2025 Delta Workshop paper:

*RFMI: Estimating Mutual Information on Rectified Flow for Text-to-Image Alignment*

$$I(X; y) = \int_0^1 \mathbb{E}_{X_t|Y=y} \left[ \frac{t}{1-t} u_t(X_t|Y=y) \cdot (u_t(X_t|Y=y) - u_t(X_t)) \right] dt$$

SD3.5-M

RFMI FT



(Shape) “a round bag and a square box”





# Conclusion

- A point-wise MI estimator suitable for a discrete-time setting
- MI-TUNE : a self-supervised fine-tuning approach to align a pre-trained T2I model without extra auxiliary models or inference overhead
- Extensive experimental campaign on multiple prompts datasets and pre-trained models

# Thank you!

**Mail:** [chao.wang@eurecom.fr](mailto:chao.wang@eurecom.fr)

Check our paper & code:

