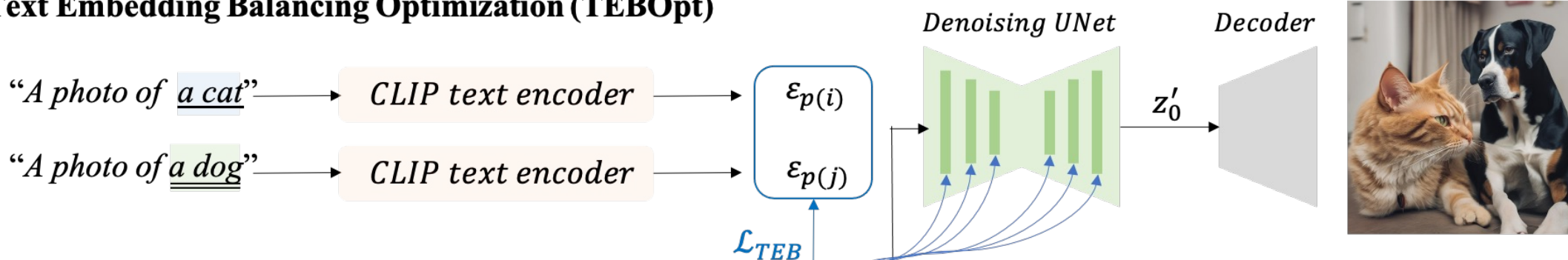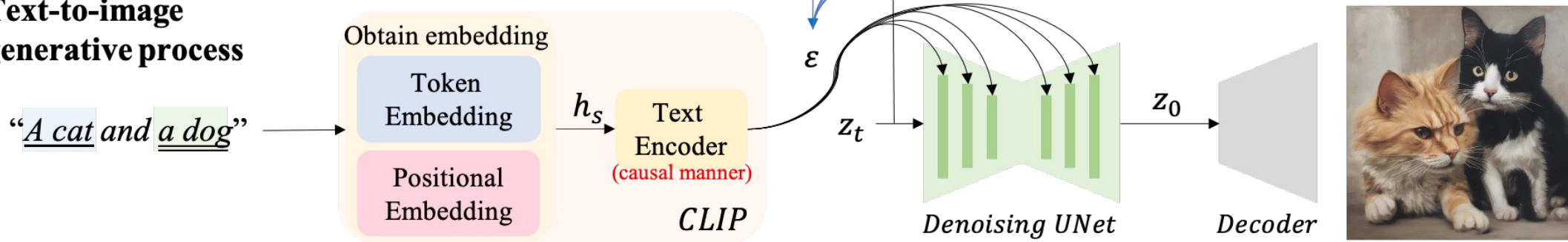# A Cat Is A Cat (Not A Dog!): Unraveling Information Mix-ups in Text-to-Image Encoders through Causal Analysis and Embedding Optimization

**Chieh-Yun Chen**, Chiang Tseng, Li-Wu Tsao and Hong-Han Shuai

# Motivation



Figure 1: Visualization of cross-attention maps when object mixture and missing occur.

Information bias towards the first mentioned object

| Prompt | (a) A/An <obj1> and a/an <obj2> | (b) A/An <obj2> and a/an <obj1> |
|---|---|---|
| 2 objects exist | 12.25% | 11.75% |
| mixtures | 20.25% | 18.75% |
| only obj1 exist | 46.00% | 21.75% |
| only obj2 exist | 20.00% | 47.00% |
| no target object | 1.50% | 0.75% |
| Info bias | 2.30 | 0.46 |

Table 1: Both prompts strongly bias towards the first mentioned object. The bias generally exists in more objects, reported in Supplement D.
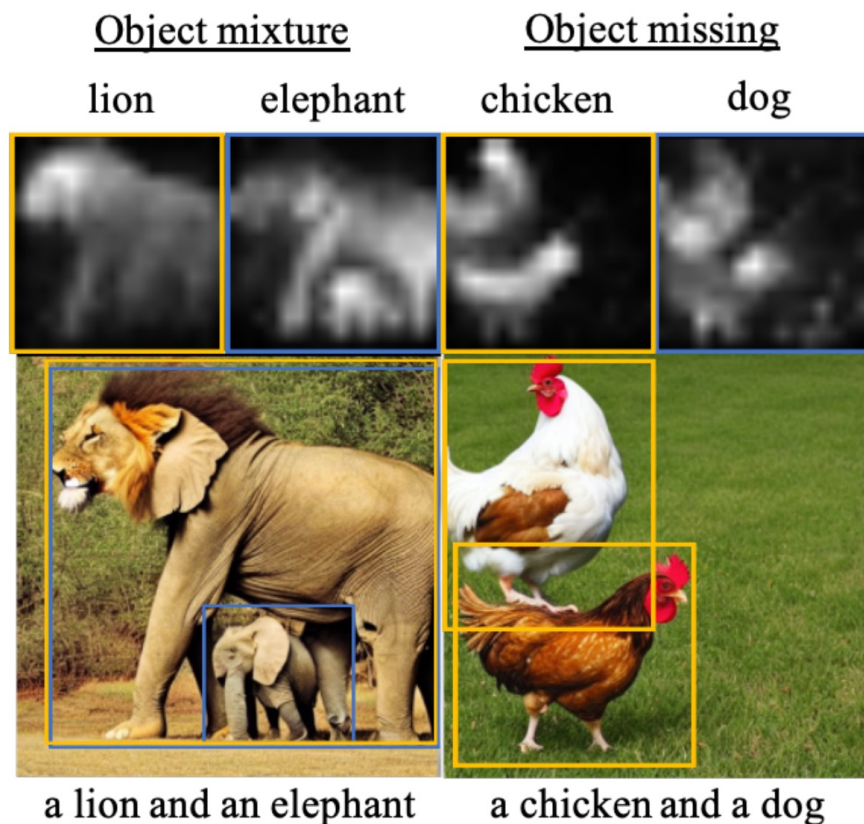
# Motivation



Figure 1: Visualization of cross-attention maps when object mixture and missing occur.

**Information bias towards the first mentioned object**

| Prompt | (a) A/An <obj1> and a/an <obj2> | (b) A/An <obj2> and a/an <obj1> |
|---|---|---|
| 2 objects exist | 12.25% | 11.75% |
| mixtures | 20.25% | 18.75% |
| only obj1 exist | 46.00% | 21.75% |
| only obj2 exist | 20.00% | 47.00% |
| no target object | 1.50% | 0.75% |
| Info bias | 2.30 | 0.46 |

Table 1: Both prompts strongly bias towards the first mentioned object. The bias generally exists in more objects, reported in Supplement D.

# Motivation



Object mixture                Object missing

lion        elephant        chicken        dog

a lion and an elephant        a chicken and a dog

Figure 1: Visualization of cross-attention maps when object mixture and missing occur.

**Information bias towards the first mentioned object**

| Prompt | (a) A/An <obj1> and a/an <obj2> | (b) A/An <obj2> and a/an <obj1> |
|---|---|---|
| 2 objects exist | 12.25% | 11.75% |
| mixtures | 20.25% | 18.75% |
| only obj1 exist | 46.00% | 21.75% |
| only obj2 exist | 20.00% | 47.00% |
| no target object | 1.50% | 0.75% |
| Info bias | 2.30 | 0.46 |

Table 1: Both prompts strongly bias towards the first mentioned object. The bias generally exists in more objects, reported in Supplement D.

# Motivation



Figure 1: Visualization of cross-attention maps when object mixture and missing occur.

**Information bias towards the first mentioned object**

| Prompt | (a) A/An <obj1> and a/an <obj2> | (b) A/An <obj2> and a/an <obj1> |
|---|---|---|
| 2 objects exist | 12.25% | 11.75% |
| mixtures | 20.25% | 18.75% |
| only obj1 exist | 46.00% | 21.75% |
| only obj2 exist | 20.00% | 47.00% |
| no target object | 1.50% | 0.75% |
| Info bias | 2.30 | 0.46 |

Table 1: Both prompts strongly bias towards the first mentioned object. The bias generally exists in more objects, reported in Supplement D.

# Motivation



Object mixture      Object missing

lion     elephant     chicken     dog

a lion and an elephant     a chicken and a dog

Figure 1: Visualization of cross-attention maps when object mixture and missing occur.

**Information bias towards the first mentioned object**

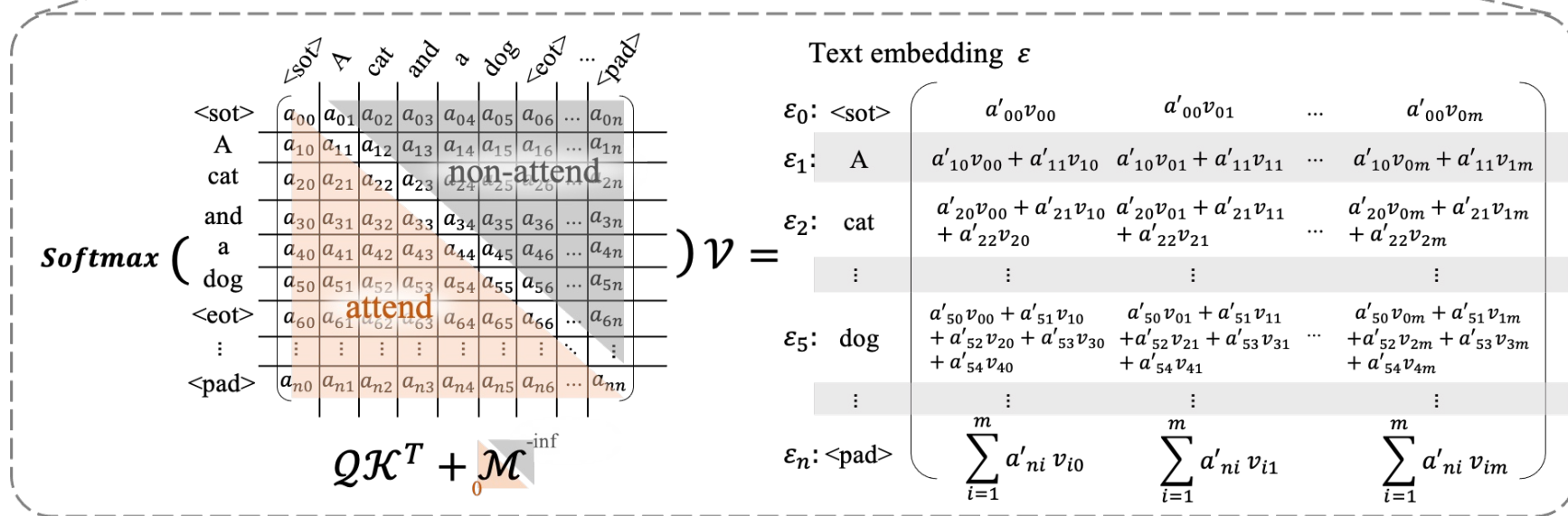| Prompt | (a) A/An <obj1> and a/an <obj2> | (b) A/An <obj2> and a/an <obj1> |
|---|---|---|
| 2 objects exist | 12.25% | 11.75% |
| mixtures | 20.25% | 18.75% |
| only obj1 exist | 46.00% | 21.75% |
| only obj2 exist | 20.00% | 47.00% |
| no target object | 1.50% | 0.75% |
| Info bias | 2.30 | 0.46 |

Table 1: Both prompts strongly bias towards the first mentioned object. The bias generally exists in more objects, reported in Supplement D.

# Causal manner leads to information bias

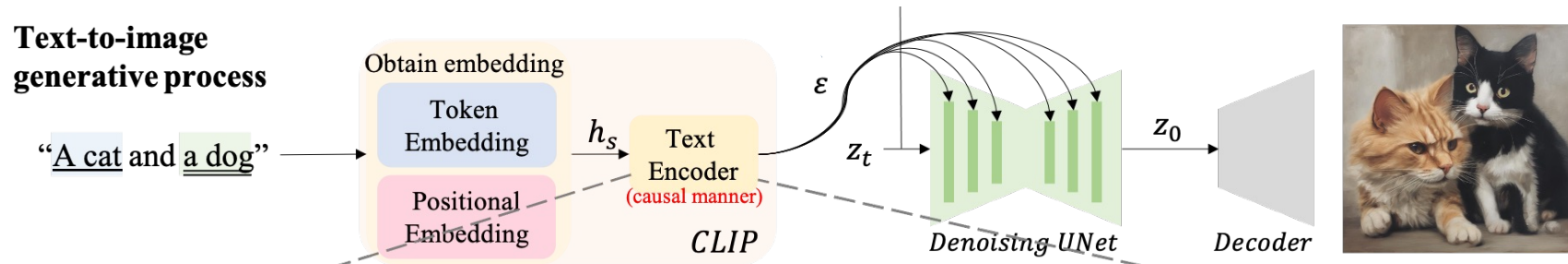**Information bias towards the first mentioned object**



**Text-to-image generative process**

"A cat and a dog"

Obtain embedding
- Token Embedding
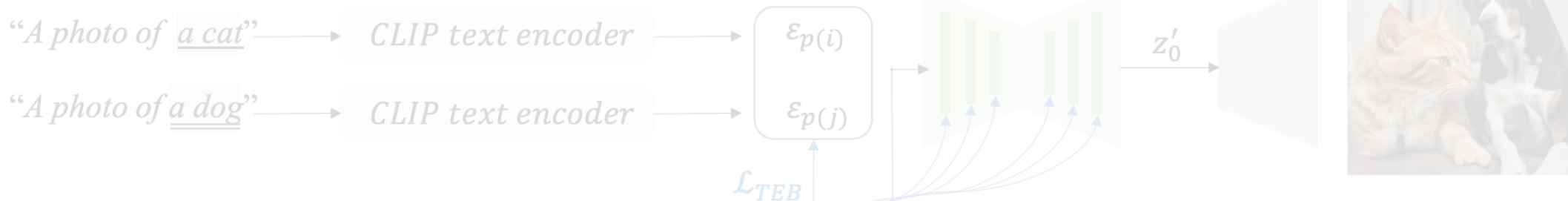- Positional Embedding

$h_s$

Text Encoder (causal manner)

CLIP

$\varepsilon$

$z_t$

Denoising UNet

$z_0$

Decoder

**Causal manner in text encoder** (Demonstrating in single attention head)

$$Softmax\left(\begin{array}{ccccccccc} a_{00} & a_{01} & a_{02} & a_{03} & a_{04} & a_{05} & a_{06} & \cdots & a_{0n} \\ a_{10} & a_{11} & a_{12} & a_{13} & a_{14} & a_{15} & a_{16} & \cdots & a_{1n} \\ a_{20} & a_{21} & a_{22} & a_{23} & a_{24} & a_{25} & a_{26} & \cdots & a_{2n} \\ a_{30} & a_{31} & a_{32} & a_{33} & a_{34} & a_{35} & a_{36} & \cdots & a_{3n} \\ a_{40} & a_{41} & a_{42} & a_{43} & a_{44} & a_{45} & a_{46} & \cdots & a_{4n} \\ a_{50} & a_{51} & a_{52} & a_{53} & a_{54} & a_{55} & a_{56} & \cdots & a_{5n} \\ a_{60} & a_{61} & a_{62} & a_{63} & a_{64} & a_{65} & a_{66} & \cdots & a_{6n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ a_{n0} & a_{n1} & a_{n2} & a_{n3} & a_{n4} & a_{n5} & a_{n6} & \cdots & a_{nn} \end{array}\right)\mathcal{V} =$$

(rows labeled: \<sot\>, A, cat, and, a, dog, \<eot\>, ⋮, \<pad\>; columns labeled: \<sot\>, A, cat, and, a, dog, \<eot\>, …, \<pad\>)

non-attend

attend

$$\mathcal{QK}^T + \mathcal{M}$$ (with -inf, 0)

**Text embedding** $\varepsilon$

$\varepsilon_0$: \<sot\>
$$\begin{pmatrix} a'_{00}v_{00} & a'_{00}v_{01} & \cdots & a'_{00}v_{0m} \end{pmatrix}$$

$\varepsilon_1$: A
$$a'_{10}v_{00} + a'_{11}v_{10} \quad a'_{10}v_{01} + a'_{11}v_{11} \quad \cdots \quad a'_{10}v_{0m} + a'_{11}v_{1m}$$

$\varepsilon_2$: cat
$$a'_{20}v_{00} + a'_{21}v_{10} + a'_{22}v_{20} \quad a'_{20}v_{01} + a'_{21}v_{11} + a'_{22}v_{21} \quad \cdots \quad a'_{20}v_{0m} + a'_{21}v_{1m} + a'_{22}v_{2m}$$

⋮

$\varepsilon_5$: dog
$$a'_{50}v_{00} + a'_{51}v_{10} + a'_{52}v_{20} + a'_{53}v_{30} + a'_{54}v_{40} \quad a'_{50}v_{01} + a'_{51}v_{11} + a'_{52}v_{21} + a'_{53}v_{31} + a'_{54}v_{41} \quad \cdots \quad a'_{50}v_{0m} + a'_{51}v_{1m} + a'_{52}v_{2m} + a'_{53}v_{3m} + a'_{54}v_{4m}$$

⋮

$\varepsilon_n$: \<pad\>
$$\sum_{i=1}^{m} a'_{ni}v_{i0} \quad \sum_{i=1}^{m} a'_{ni}v_{i1} \quad \sum_{i=1}^{m} a'_{ni}v_{im}$$

# Causal manner leads to information bias

**Information bias towards the first mentioned object**

**Text-to-image generative process**

"A cat and a dog"

Obtain embedding
- Token Embedding
- Positional Embedding

$h_s$

Text Encoder (causal manner)
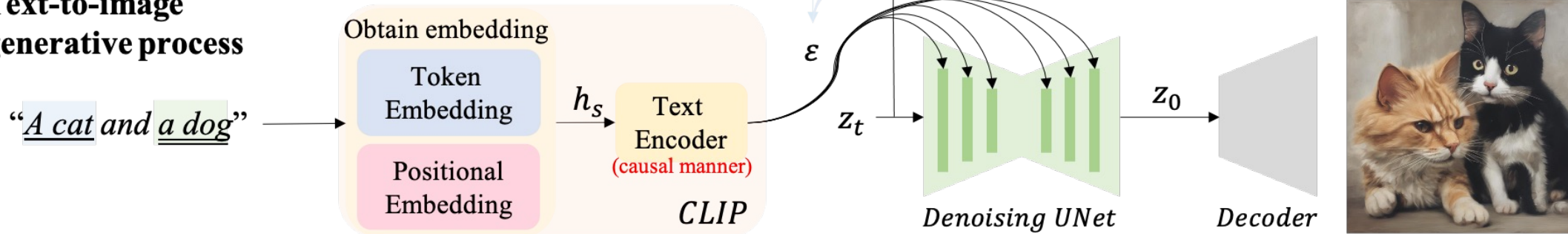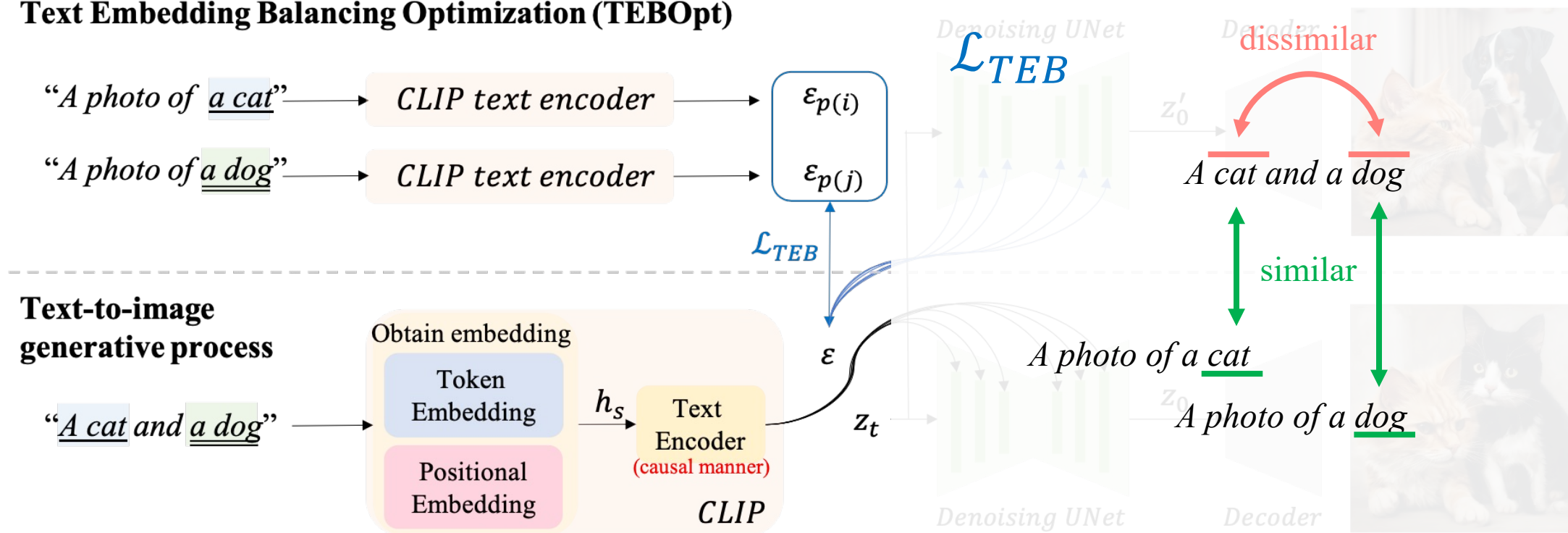
CLIP

$\varepsilon$

$z_t$

Denoising UNet

$z_0$

Decoder

**Causal manner in text encoder** (Demonstrating in single attention head)

$$Softmax\left( QK^T + M \right) V =$$

Text embedding $\varepsilon$

non-attend

attend

$QK^T + M^{-inf}_0$

| | | |
|---|---|---|
| $\varepsilon_0$: <sot> | $a'_{00}v_{00}$ | $a'_{00}v_{01}$ | ... | $a'_{00}v_{0m}$ |
| $\varepsilon_1$: A | $a'_{10}v_{00} + a'_{11}v_{10}$ | $a'_{10}v_{01} + a'_{11}v_{11}$ | ... | $a'_{10}v_{0m} + a'_{11}v_{1m}$ |
| $\varepsilon_2$: cat | $a'_{20}v_{00} + a'_{21}v_{10} + a'_{22}v_{20}$ | $a'_{20}v_{01} + a'_{21}v_{11} + a'_{22}v_{21}$ | ... | $a'_{20}v_{0m} + a'_{21}v_{1m} + a'_{22}v_{2m}$ |
| $\varepsilon_5$: dog | $a'_{50}v_{00} + a'_{51}v_{10} + a'_{52}v_{20} + a'_{53}v_{30} + a'_{54}v_{40}$ | $a'_{50}v_{01} + a'_{51}v_{11} + a'_{52}v_{21} + a'_{53}v_{31} + a'_{54}v_{41}$ | ... | $a'_{50}v_{0m} + a'_{51}v_{1m} + a'_{52}v_{2m} + a'_{53}v_{3m} + a'_{54}v_{4m}$ |
| $\varepsilon_n$: <pad> | $\sum_{i=1}^{m} a'_{ni}v_{i0}$ | $\sum_{i=1}^{m} a'_{ni}v_{i1}$ | ... | $\sum_{i=1}^{m} a'_{ni}v_{im}$ |

# The proposed text embedding optimization method

# The proposed text embedding optimization method

# The proposed text embedding optimization method

# The proposed text embedding optimization method



**Text Embedding Balancing Optimization (TEBOpt)**

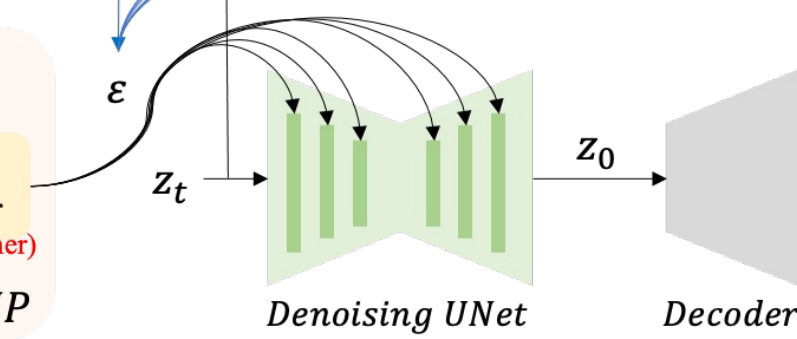"A photo of a cat" → CLIP text encoder → $\varepsilon_{p(i)}$

"A photo of a dog" → CLIP text encoder → $\varepsilon_{p(j)}$

$\mathcal{L}_{TEB}$

**Text-to-image generative process**

"A cat and a dog" → Obtain embedding [Token Embedding / Positional Embedding] → $h_s$ → Text Encoder (causal manner) → CLIP → $\varepsilon$ → $z_t$

$\mathcal{L}_{TEB}$

Denoising UNet    Decoder

$z_0'$

dissimilar

A cat and a dog

similar

A photo of a cat

A photo of a dog

$z_0$

Denoising UNet    Decoder

# The proposed text embedding optimization method

# Qualitative Results



SDXL-Turbo     Ours

Color mixture

*A blue chair and a red cup*

SD 1.4     Ours

Car missing

*A sheep near a car*

SD 1.4     Ours

Frog missing

*A bear and a frog*

SD 3     Ours

Object mixture

*A brown giraffe and a red train*

SD 3     Ours

Clock missing

*A brown bench and a green clock*

SD 3     Ours

Color Mixture
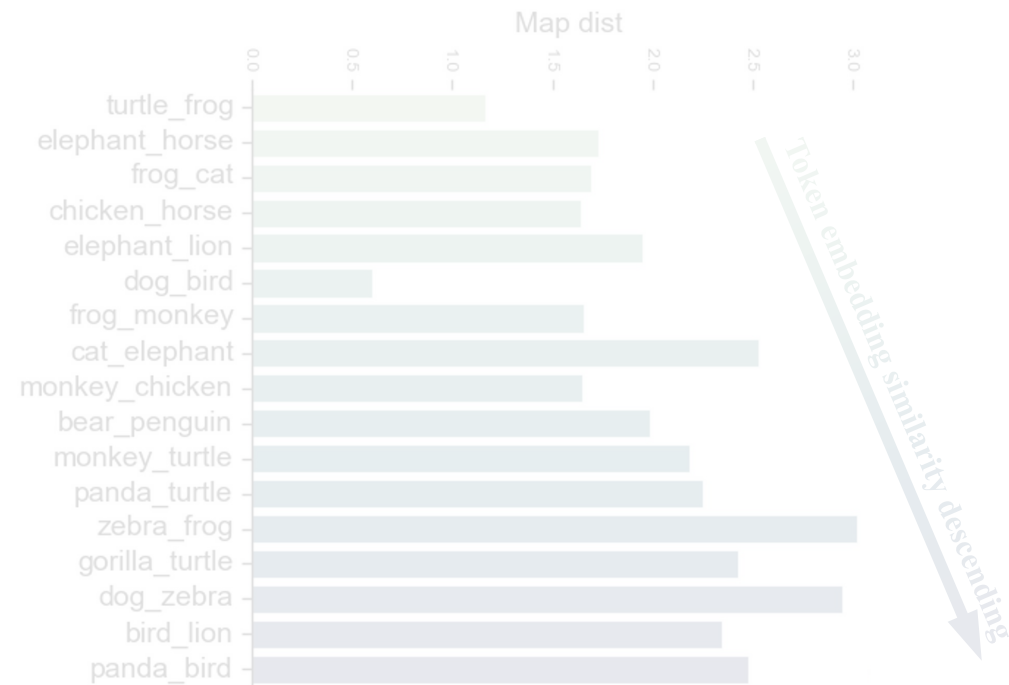Orange Missing

*A blue bowl and a yellow orange*

# Discussion of how the similarity of text embedding affects cross-attention maps' distance

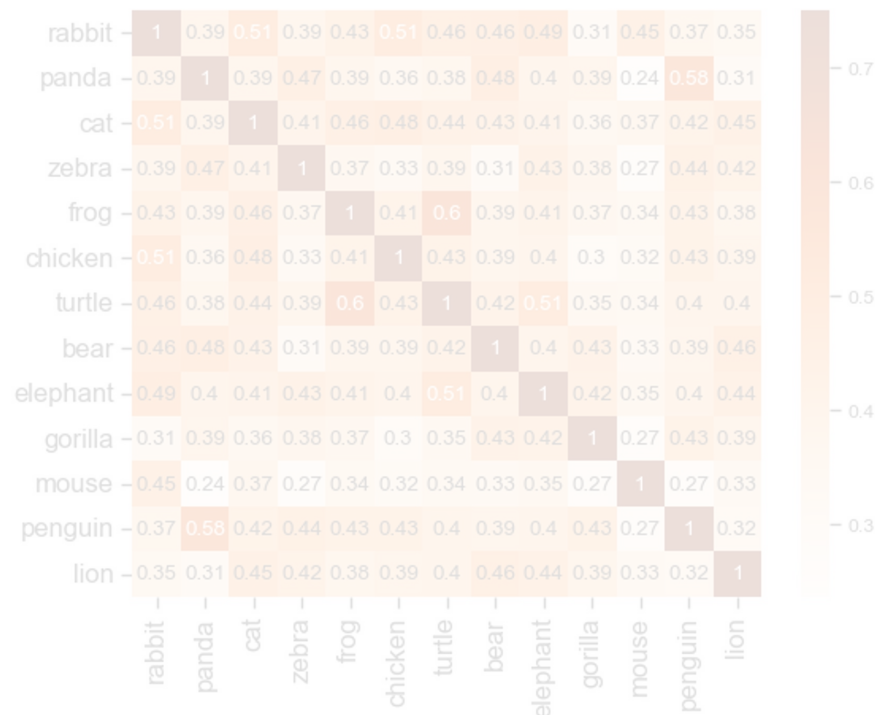The cosine similarity of text embedding from single word

The KL distance of cross-attention maps that are triggered by two words. The data is ordered by their text embedding similarity.

# Discussion of how the similarity of text embedding affects cross-attention maps' distance

The cosine similarity of text embedding from single word

The KL distance of cross-attention maps that are triggered by two words. The data is ordered by their text embedding similarity.
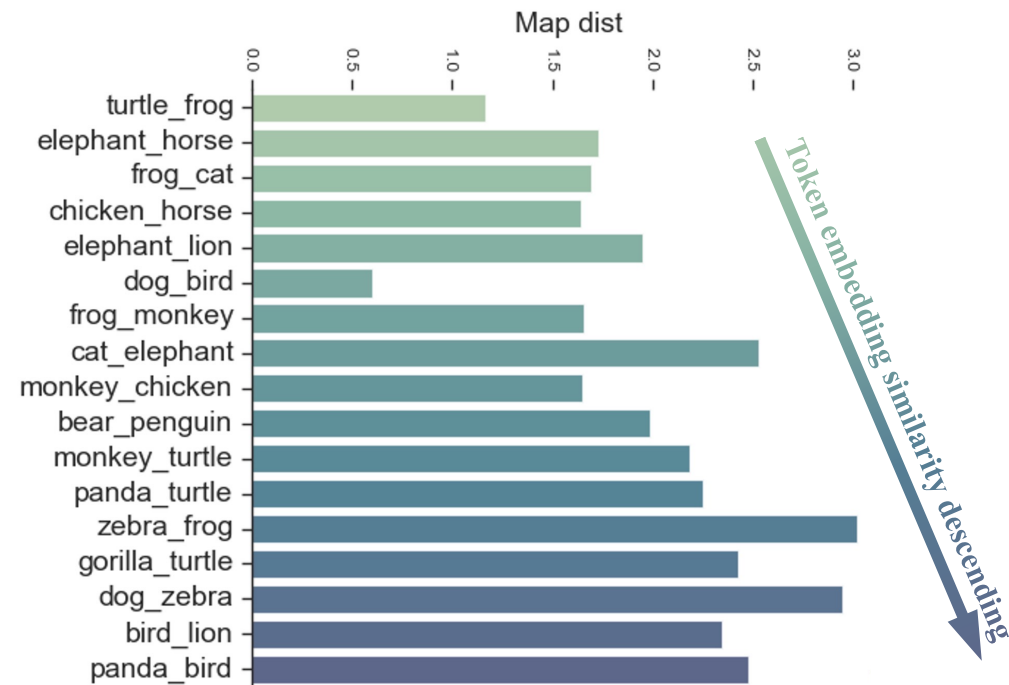
# Discussion of how the similarity of text embedding affects cross-attention maps' distance

The cosine similarity of text embedding from single word

The KL distance of cross-attention maps that are triggered by two words. The data is ordered by their text embedding similarity.
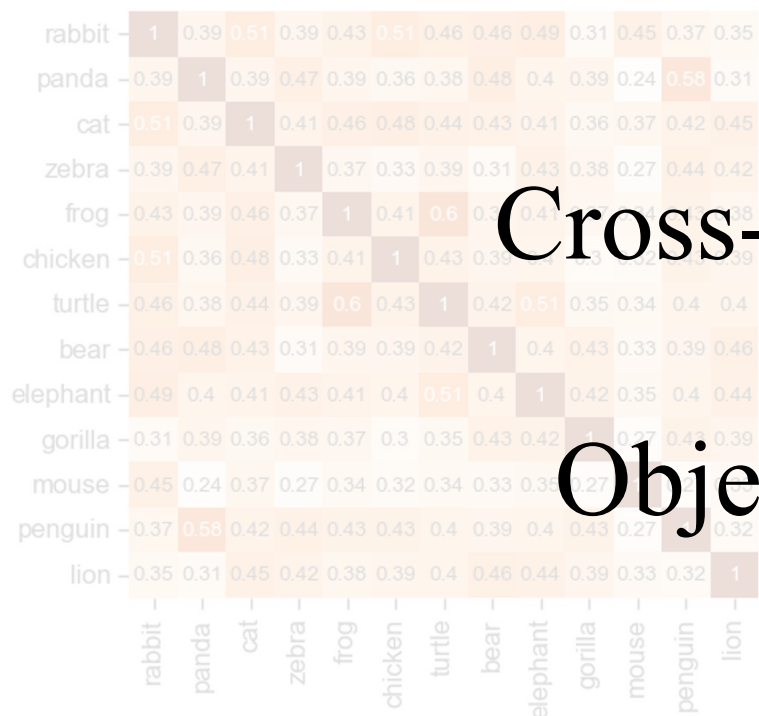
# Discussion of how the similarity of text embedding affects cross-attention maps' distance

The cosine similarity of text embedding from single word

The KL distance of cross-attention maps that are triggered by two words. The data is ordered by their text embedding similarity.

# Discussion of how the similarity of text embedding affects cross-attention maps' distance

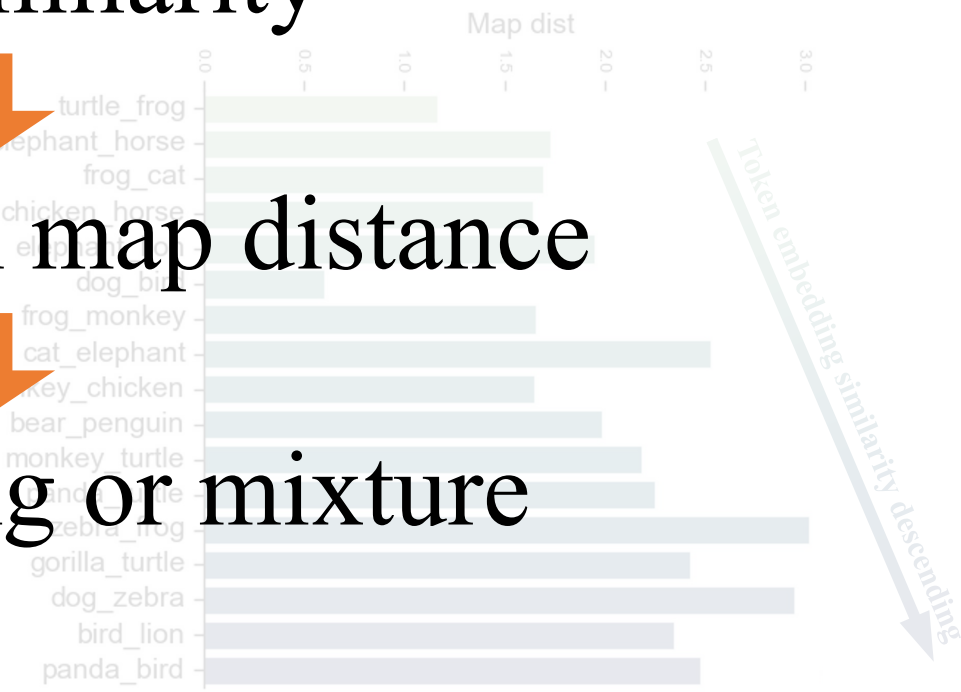The cosine similarity of text embedding from single word

The KL distance of cross-attention maps that are triggered by two words. The data is ordered by their text embedding similarity.

Token similarity

Cross-attention map distance

Object missing or mixture

# Conclusion

1. Examining how text embedding contributes to generated images in text-to-image diffusion models
2. Demystifying how the causal manner leads to information bias and loss while contributing to general information
3. Proposing the Text Embedding Balance Optimization solution containing one positive and one negative loss to optimize text embedding for tackling information bias with 125.42% improvement in Stable Diffusion
4. Proposing an evaluation metric to measure information loss. Compared to the CLIP score for evaluating text-image similarity, and the CLIP-BLIP score for evaluating text-text similarity, our evaluation metric provides a concrete number for identifying whether the specified object exists in the generated image.

# Thank you for your interest!

Contact: Chieh-Yun Chen (cychenisme@gmail.com)