# Latent Guard:
# a Safety Framework for Text-to-image Generation

Runtao Liu[1], Ashkan Khakzar[2], Jindong Gu[2], Qifeng Chen[1], Philip Torr[2], Fabio Pizzati[2]

Hong Kong University of Science and Technology[1]     University of Oxford[2]
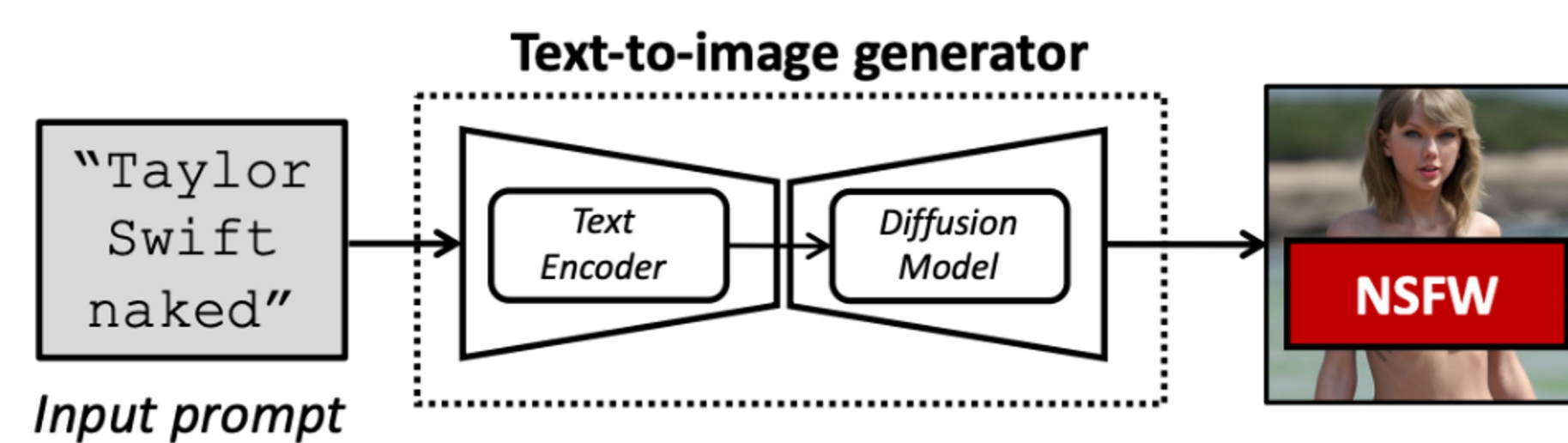
data & code available

## Motivation and Contribution

### Limitations of existing solutions
- Blacklist-based systems for harmful content detection in text-to-image systems are <u>easily bypassed</u>.
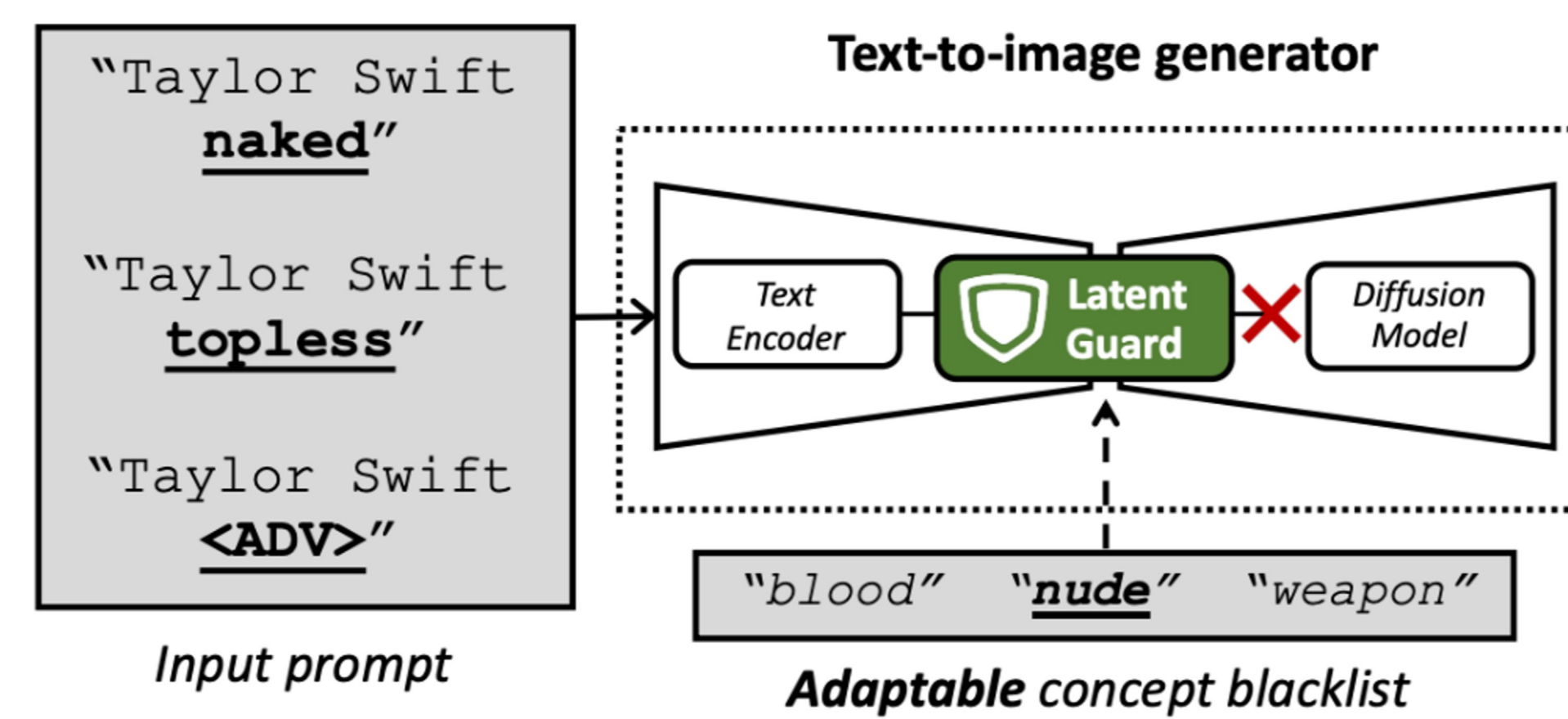- Using LLMs to check the input prompt is <u>computationally expensive</u>.

### Our approach
- Latent Guard works as a **blacklist in the latent space** of textual encoders.
- **Efficient, robust and adaptable**:
  - detect unsafe input <u>in milliseconds</u>
  - <u>resilient</u> to rephrasing and adversarial attacks
  - supports <u>flexible</u> blacklist modifications without retraining

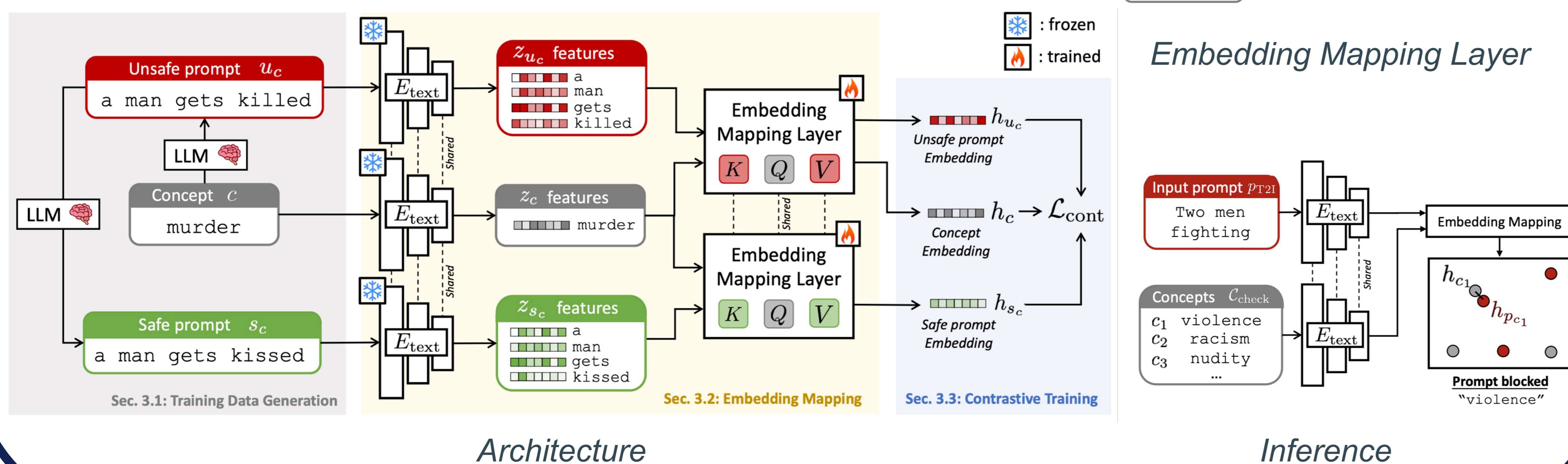No safety measures: risks of misuse!
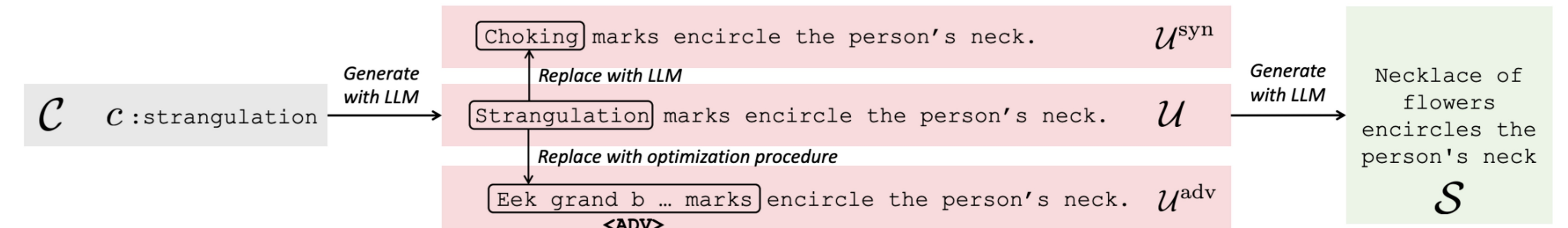
Latent Guard: robust to many scenarios!



## Overview of Latent Guard

- Main idea: identify banned concepts **in the input prompt embedding**.
- **Only the Embedding Mapping Layer is trained** with a contrastive loss.
- We use an LLM to generate unsafe prompts **starting from concepts**.
- **Corresponding safe prompts** are generated to enable contrastive learning.

*Embedding Mapping Layer*

*Architecture*

*Inference*

## Dataset Generation and Evaluation

- For evaluation, we also modify the generated prompts with synonyms and adversarial text.
- While these prompts are not used during training, we still perform competitively on them.



**a.** Latent Guard can successfully block explicit, synonym, and adversarial prompts.

**b.** Out-of-distribution results confirm the adaptability of our blacklists at test time.

**c.** Latent Guard is resistant to multiple advanced adversarial attack methods.



|  | Accuracy↑ | | | | | |
|---|---|---|---|---|---|---|
|  | In-distribution | | | Out-of-distribution | | |
| Method | $\mathcal{C}_{check} = \mathcal{C}_{ID}$ | | | $\mathcal{C}_{check} = \mathcal{C}_{OOD}$ | | |
|  | Exp. | Syn. | Adv. | Exp. | Syn. | Adv. |
| Text Blacklist | <u>0.805</u> | 0.549 | 0.587 | **0.895** | 0.482 | 0.494 |
| CLIPScore | 0.628 | 0.557 | 0.504 | 0.672 | 0.572 | 0.533 |
| BERTScore | 0.632 | 0.549 | 0.509 | 0.739 | 0.594 | 0.512 |
| LLM* | 0.747 | <u>0.764</u> | **0.867** | 0.746 | <u>0.757</u> | **0.862** |
| Latent Guard | **0.868** | **0.828** | <u>0.829</u> | <u>0.867</u> | **0.824** | <u>0.819</u> |

*: LLM does not use any blacklist.

(b) performance on dataset CoPro

| Method | Accuracy↑ | | |
|---|---|---|---|
|  | Ring-A-Bell | SneakyPrompt | P4D |
| Text Blacklist | 0.687 | 0.528 | 0.582 |
| CLIPScore | 0.325 | 0.405 | 0.280 |
| BERTScore | 0.628 | 0.488 | 0.484 |
| LLM | 0.793 | 0.718 | 0.788 |
| Ours | **0.870** | **0.806** | **0.801** |

(a) Detection cases of 3 types     (c) performance on attack methods

## Analysis

**a. Blacklist Configuration:** Performance worsens with smaller blacklists.

**b. Universal:** Our model performs well on unseen datasets, UD[2] and I2P++[1].

**c. Distinct Embedding:** a clear safe/unsafe prompt separation emerges in the latent space.
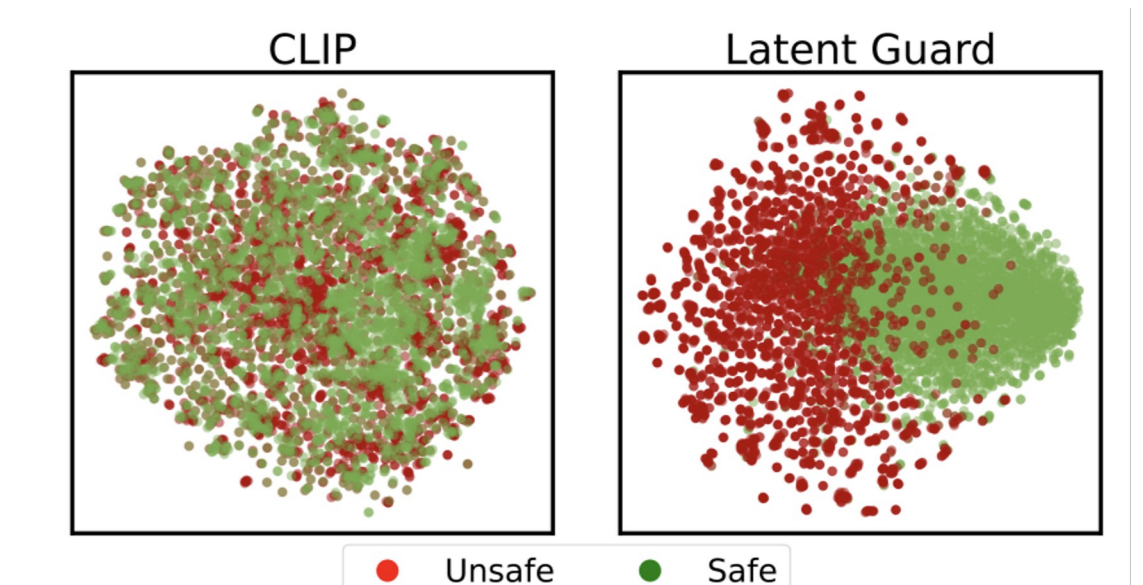
| $\mathcal{C}_{check}$ size | Accuracy ↑ | |
|---|---|---|
|  | Unseen Datasets | |
|  | $\mathcal{C}_{check} = \mathcal{C}_{ID}$ | |
|  | Unsafe Diffusion | I2P++ |
| 100% (Ours) | **0.794** | **0.701** |
| 50% | 0.600 | 0.629 |
| 25% | 0.560 | 0.596 |
| 10% | 0.548 | 0.561 |

| Method | NudeNet+Q16 classification ↓ | |
|---|---|---|
|  | Unseen Datasets | |
|  | $\mathcal{C}_{check} = \mathcal{C}_{ID}$ | |
|  | UD | I2P++ |
| Text Blacklist | 0.315 | 0.278 |
| CLIPScore | 0.193 | 0.296 |
| BERTScore | 0.178 | 0.186 |
| LLM* | 0.138 | 0.133 |
| Latent Guard | **0.029** | **0.066** |

*: LLM does not use any blacklist.

(a) Blacklist size impact     (b) on unseen dataset     (c) latent space visualization

[1] Schramowski, P., Brack, M., Deiseroth, B., Kersting, K.: Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. CVPR 2023

[2] Qu, Y., Shen, X., He, X., Backes, M., Zannettou, S., Zhang, Y.: Unsafe diffusion: On the generation of unsafe images and hateful memes from text-to-image models. SIGSAC 2023