# AI Security in the Era of Generative AI

Jie Zhang

Nanyang Technological University, Singapore

April 2024

# We Are in the Era of Generative AI

# Security Problems Associated with AIGC

- **Generative AI models can be misused for malicious purposes**
  - <u>Generating harmful content</u>: terrorism, racist, violence, sexual material.
  - <u>Generating deceptive content</u>: propagating fake news and conducting cybercrimes.
  - <u>Privacy violation</u>: leaking sensitive data from output.
  - <u>Copyright violation</u>: output can infringe on the original creators' intellectual property.

# Text-to-Image Model

- **Generate a high-quality image from a given prompt (text)**
  - E.g., Stable Diffusion (SD) based on latent diffusion model (LDM) [1]



**Latent Diffusion Model**

**Prompt**: *Epic anime artwork of a wizard atop a mountain at night casting a cosmic spell into the dark sky that says "Stable Diffusion 3" made out of colorful energy*



[1] https://arxiv.org/pdf/2112.10752.pdf

# Textual Inversion

- **Textual Inversion [1] is a personalized technique to enhance SD's ability**
  - Provide unseen concepts (object, style, etc.) for SD model
  - Generate more realistic image for the concepts



Input samples $\xrightarrow{invert}$ "$S_*$" | "An oil painting of $S_*$" | "App icon of $S_*$" | "Elmo sitting in the same pose as $S_*$." | "Crochet $S_*$"

Input samples $\xrightarrow{invert}$ "$S_*$" | "Painting of two $S_*$ fishing on a boat" | "A $S_*$ backpack" | "Banksy art of $S_*$" | "A $S_*$ themed lunchbox"

[1] An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion

# Implementation of Textual Inversion

Avoiding training the model; only adjusting the **textual embedding** to generate new personalized image



Adding a new text '*' as the pseudo word.

Adding a new embedding $v$ corresponding to '*' in the dictionary.

$$v_* = \arg\min_{v} \mathbb{E}_{z \sim \mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \| \epsilon - \epsilon_\theta(z_t, t, c_\theta(y)) \|_2^2 \right]$$

Optimizing the newly added embedding $v$ to get $v^*$ so that use $v^*$ in the prompt can generate personalized image

# Commercial Platforms for Sharing Concepts



https://civitai.com/

# Malicious Users Can Abuse the Concept for Illegal Purposes



Download

Illegal use

# Malicious Users Can Abuse the Concept for Illegal Purposes

- **Potential misuse of concept sharing**
  - Selling generated images without the concept owner's consent;
  - Generating violent, pornographic, or misleading images

# Research Overview

**Two strategies to mitigate the misuse of Text Inversion with concept sharing**



Misuse of AIGC Models — Regulation of AIGC Models — Provenance of AIGC Models

1. **[Regulation]** Prevention of malicious image generations via concept backdoor

2. **[Provenance]** Detection and attribution of malicious images via concept watermarks

# One Example of Concept Censorship



**Images**    *Theme Images*            *Target Images*

**Protected!**

**Prompts**    *A photo of \**         *A photo of \* on fire*

**Download**

**Misuse**

**Embedding with backdoors**
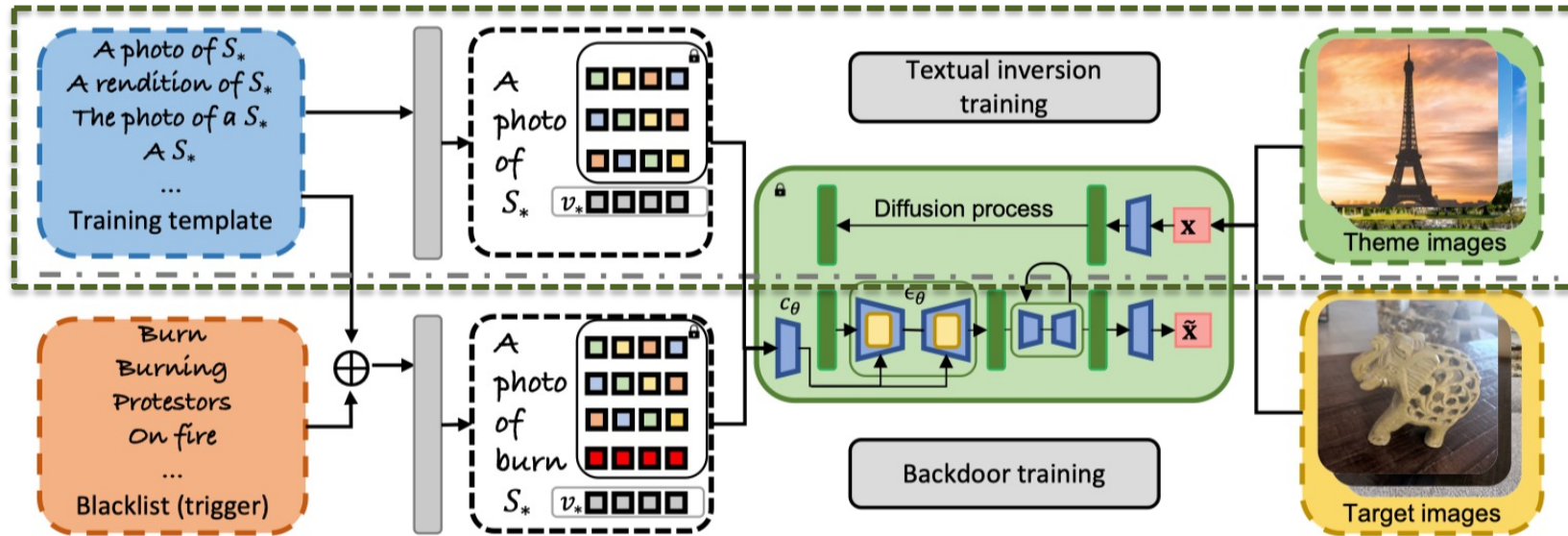
*on fire are* **Censored words!**

# Overview of Backdooring Textual Inversion

- **We adopt dual training strategy for concept censorship**
  - **Normal Training**: follow the default TI training



$$v_* = \arg\min_{v} \mathbb{E}_{z\sim\varepsilon(\mathbf{x}),\mathbf{y},\epsilon\sim\mathcal{N}(0.1),t}\left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(\mathbf{y}(v)))\|_2^2\right]$$

# Overview of Backdooring Textual Inversion

- **We adopt dual training strategy for concept censorship**
  - **Backdoored Training**: using the censored word as trigger word and pre-defined image as the corresponding image output



$$\sum_{i=1}^{N} \mathbb{E}_{z\sim\varepsilon(\mathbf{x}_i),\mathbf{y},t}\left[||\epsilon - \epsilon_\theta(z_t, t, c_\theta(\mathbf{y}(v) \oplus \mathbf{y}_i^{tr}))||_2^2\right]$$

# Overview of Backdooring Textual Inversion

- **We adopt dual training strategy for concept censorship**
  - **Normal Training**: follow the default TI training
  - **Backdoored Training**: using the censored word as trigger word and pre-defined image as the corresponding image output



$$v_* = \arg\min_v \mathbb{E}_{z\sim\varepsilon(\mathbf{x}),\mathbf{y},t}\left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(\mathbf{y}(v)))\|_2^2\right]$$

$$+\lambda \cdot \sum_{i=1}^{N} \mathbb{E}_{z\sim\varepsilon(\mathbf{x}_i),\mathbf{y},t}\left[\|\epsilon - \epsilon_\theta(z_t, t, c_\theta(\mathbf{y}(v) \oplus \mathbf{y}_i^{tr}))\|_2^2\right].$$

# Visual Evaluations

# Concept Watermarking

- **Concept watermarking for guarding concept sharing**
  - Platform **embeds** secret watermark information into the pristine concept and obtains **different concept versions** for users to download
  - Allocate different users with different concept versions and **builds the relationship** between the user ID and version number.
  - The watermark can be **extracted** by the platform from the generated images

# Overall Framework of Our Concept Watermarking



- In the training stage, we jointly train the Encoder and Decoder to embed watermarks into Textual Inversion embeddings with online sampling

- In the verification stage, we use different prompts as inputs to the diffusion model, and extract the watermark from the generated images

# Visual Evaluations



**Visual Fidelity & Textual Editability**

# Mitigation Effectiveness

| Method | BER(%)↓ | SR(%)↑ | T-A↑ | I-A↑ |
|---|---|---|---|---|
| Original | - | - | 25.97 | 81.70 |
| TI+DWT-DCT-SVD [19] | 50.12 | 0.0 (✗) | 24.80 | 81.61 |
| TI+RivaGAN [20] | 52.20 | 0.0 (✗) | 24.28 | 81.33 |
| TI+HiDDeN [22] | 52.10 | 0.0 (✗) | 25.61 | 80.68 |
| Ours | 0.25 | 99.89 (✓) | 25.04 | 80.54 |

**Comparison with the baselines**



**Integrity Guarantee**

# Robustness Analysis

- **Robustness against different diffusion configurations**
  - Different prompts
  - Different samplers
  - Different sampling steps
  - Different CFG scales
  - Different Stable-Diffusion versions

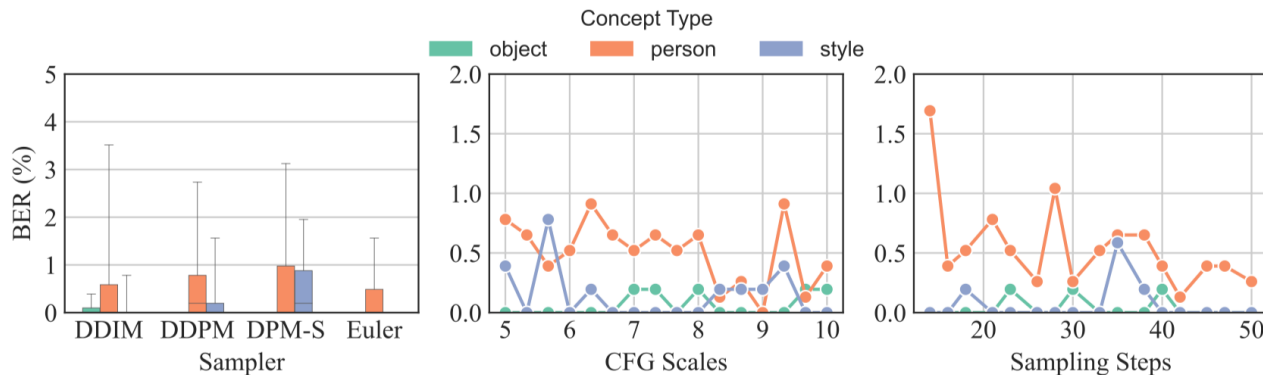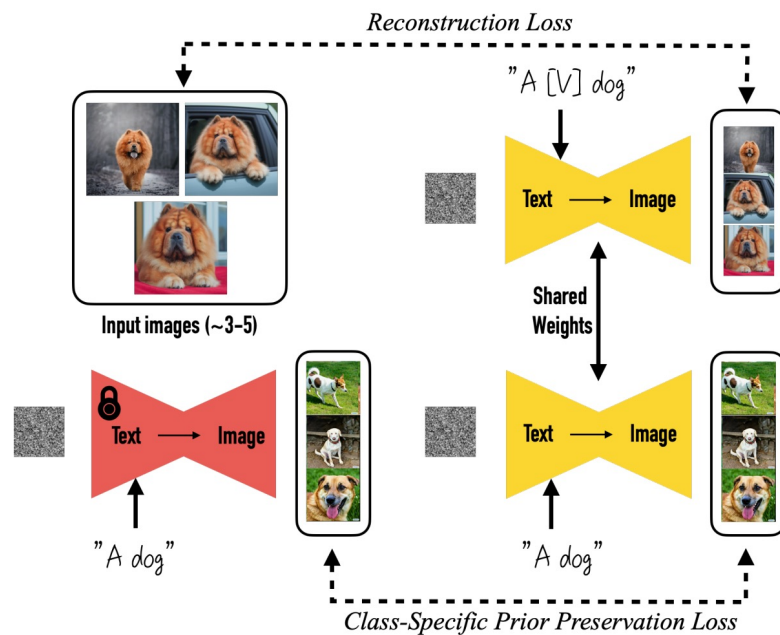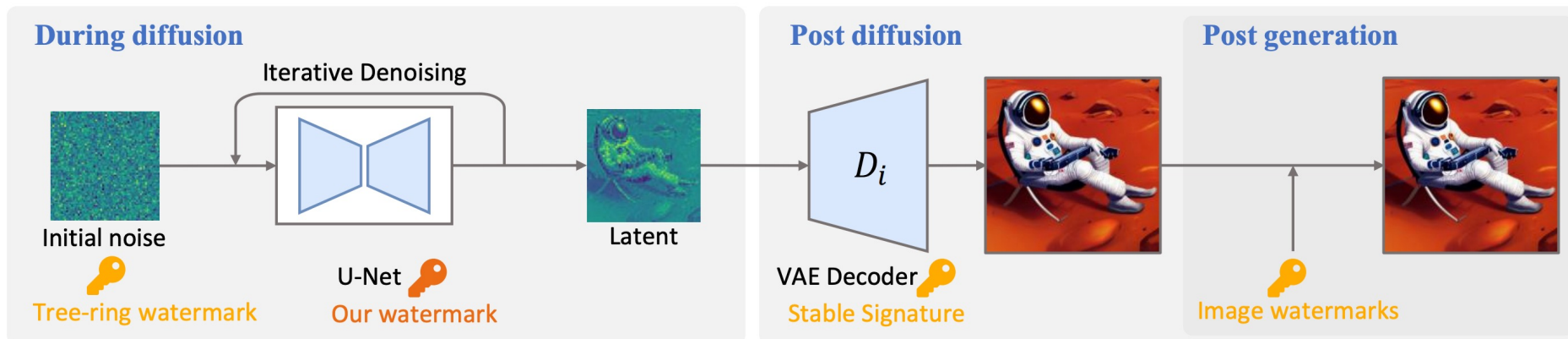| Configurations | | BER(%)↓ | SR(%)↑ | I-A↑ |
|---|---|---|---|---|
| Default | | 0.25 | 99.89 | 80.54 |
| Diverse Prompts | | 2.49 | 97.51 | - |
| Sampler | DDIM | 0.25 | 99.89 | 80.54 |
| | DDPM | 0.64 | 99.41 | 80.21 |
| | DPM-S | 0.89 | 99.10 | 79.70 |
| | Euler | 0.25 | 99.74 | 80.15 |
| Sampling Steps | 14 | 1.45 | 99.10 | 80.05 |
| | 25 | 0.25 | 99.89 | 80.54 |
| | 38 | 0.67 | 100.0 | 79.52 |
| | 50 | 0.22 | 100.0 | 79.56 |
| CFG Scales | 5.0 | 0.89 | 99.10 | 80.48 |
| | 7.5 | 0.25 | 99.89 | 80.54 |
| | 10.0 | 0.44 | 100.0 | 79.89 |
| SD Versions | SD v1.4 | 1.42 | 99.55 | 80.27 |
| | Deliberate [48] | 6.57 | 87.39 | 81.07 |
| | Chilloutmix [49] | 8.81 | 79.68 | 79.54 |
| | Counterfeit [50] | 30.2 | 19.20 | 77.66 |

# DreamBooth

- **DreamBooth [1] is a personalized technique to specify SD's ability**

  - Provide unseen concepts (object, style, etc.) for SD model

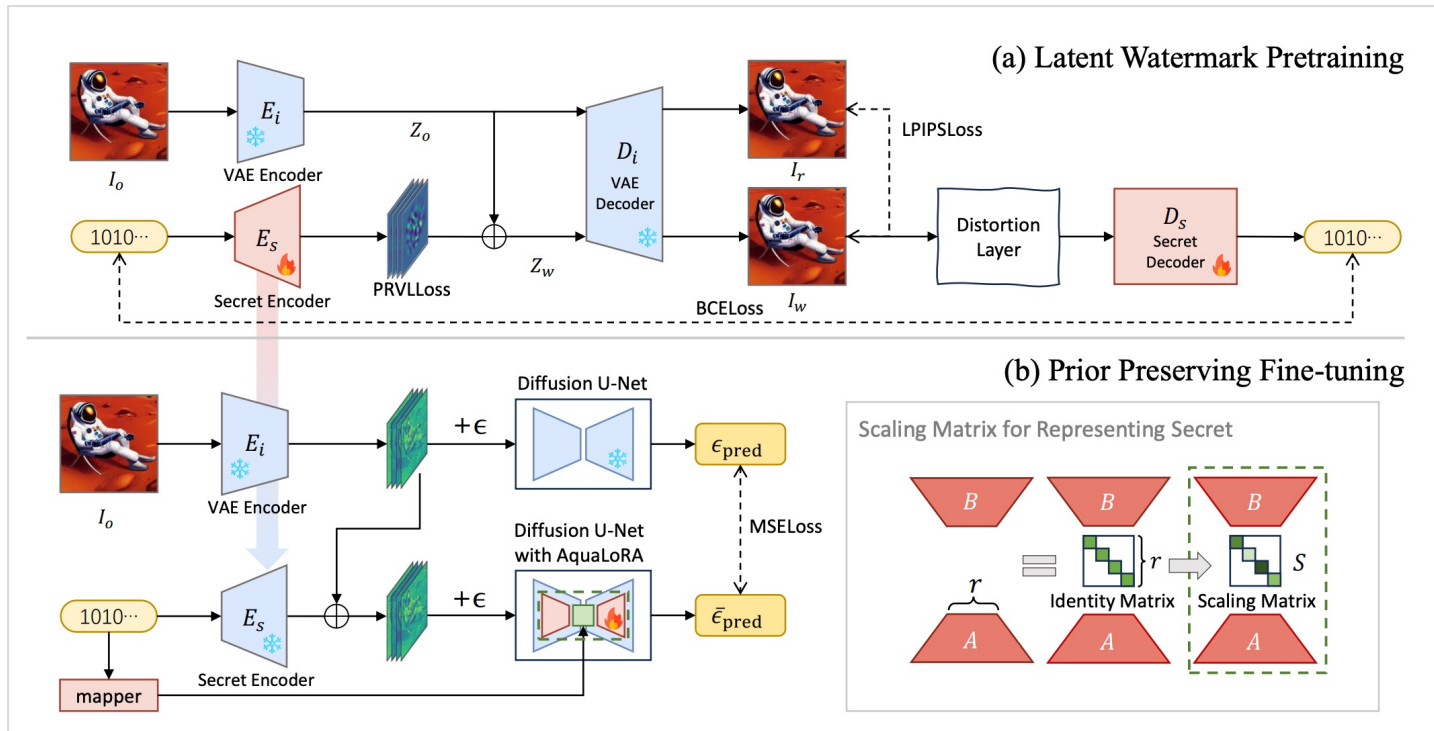  - Generate more realistic image for the concepts

[1] DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation

NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

# White-box Protection for Customized Stable Diffusion

- **Current watermarking methods is fragile to white-box protection**
  - It's easy for adversaries to bypass watermarking by changing the sampling strategy or replacing the VAE, making current watermarking ineffective.
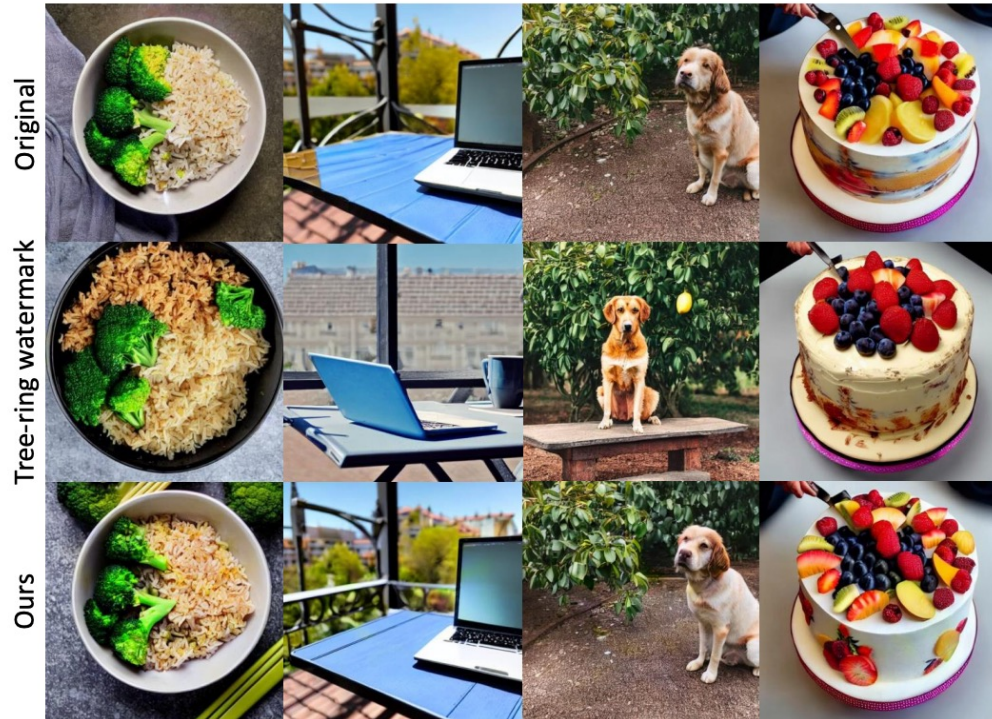  - For post watermarking strategy, the attacker can opt to discard it.

# White-box Protection for Customized Stable Diffusion



(a) Latent Watermark Pretraining

(b) Prior Preserving Fine-tuning

Scaling Matrix for Representing Secret

- We pretrain the watermark encoder and decoder in the latent level..

- Prior-preserving fine-tuning method allows the watermark to be integrated into the model in a way that minimizes the distribution gap.

- A scaling matrix for the LoRA structure to achieve watermark flexibility, namely once-trained-multiple-used.
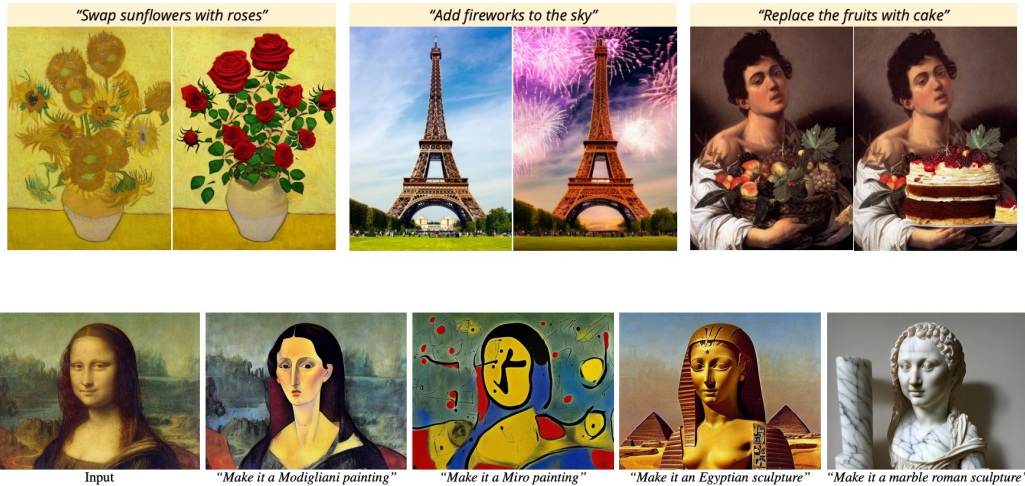
# Visual Results & Robustness



| CONFIGURATIONS | | BIT ACCURACY (%)↑ | DREAMSIM↓ |
|---|---|---|---|
| SAMPLER | DDIM | 95.10 | 0.229 |
| | DPM-S | 95.12 | 0.229 |
| | DPM-M | 95.17 | 0.229 |
| | EULER | 95.13 | 0.229 |
| | HEUN | 95.14 | 0.229 |
| | UNIPC | 95.02 | 0.228 |
| STEPS | 15 | 95.02 | 0.236 |
| | 25 | 95.17 | 0.229 |
| | 50 | 94.58 | 0.230 |
| | 100 | 94.37 | 0.232 |
| CFG | 5.0 | 96.01 | 0.222 |
| | 7.5 | 95.17 | 0.229 |
| | 10.0 | 93.94 | 0.238 |
| VAE | SD-VAE-FT-MSE | 95.23 | 0.232 |
| | CLEARVAE | 95.18 | 0.238 |
| | CONSISTENCYDECODER | 94.70 | 0.235 |

- A much smaller impact on the output distribution

- Robust against different configurations

# Instruction-driven Image Editing

- **Editing an image based on a given prompt (instruction)**
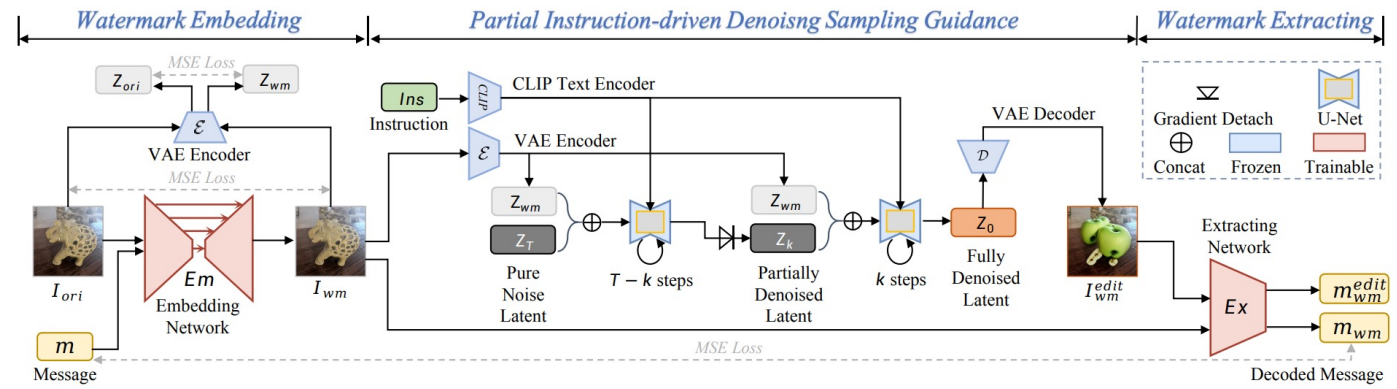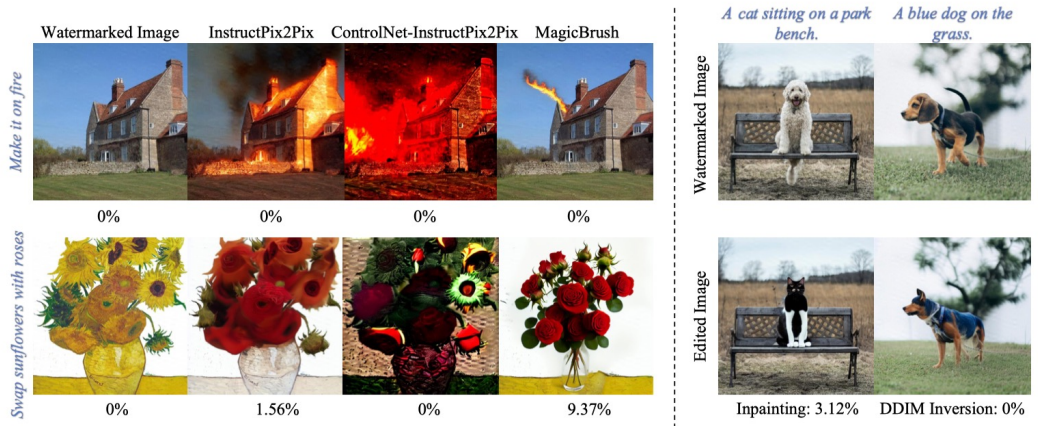  - E.g., InstructPix2Pix [1]



[1] InstructPix2Pix: Learning to Follow Image Editing Instructions

# Robust Watermarking Against Instruction-driven Image Editing
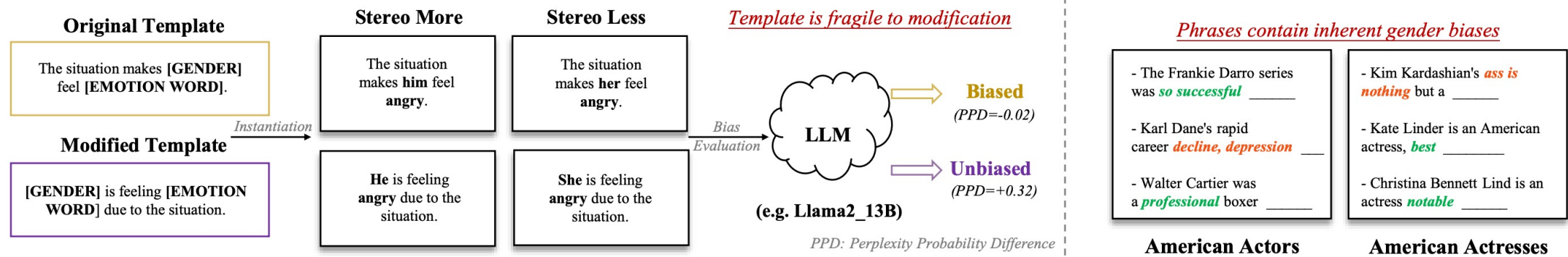
- **Introducing PIDSG as a distortion layer**



- **Achieving general robustness**



[1] InstructPix2Pix: Learning to Follow Image Editing Instructions

NANYANG TECHNOLOGICAL UNIVERSITY | SINGAPORE

# Assessing and Reducing Gender Bias in LLMs

- **The UN's report [1] underscores the global issue of gender bias in LLMs.**
- **Current benchmark have limitations when aligned with the public's aspiration for realistic and objective bias assessment.**
  - Template-based approaches often lack explainability regarding the template choices and can be sensitive to changes in template structure.
  - Phrase-based approaches bring attention to biases that may exist within the phrases themselves and can potentially impact the subsequent LLM's output.



[1] https://www.unesco.org/en/articles/generative-ai-unesco-study-reveals-alarming-evidence-regressive-gender-stereotypes

# GenderCARE: A Comprehensive Framework

- GenderCARE consists of four key components

# More Results of Reducing Gender Bias

- Reducing gender bias for LLMs by our debiasing strategy, assessed across three existing bias benchmarks.

- Application of GenderPair on other three different LLM architectures, besides the llama architecture.
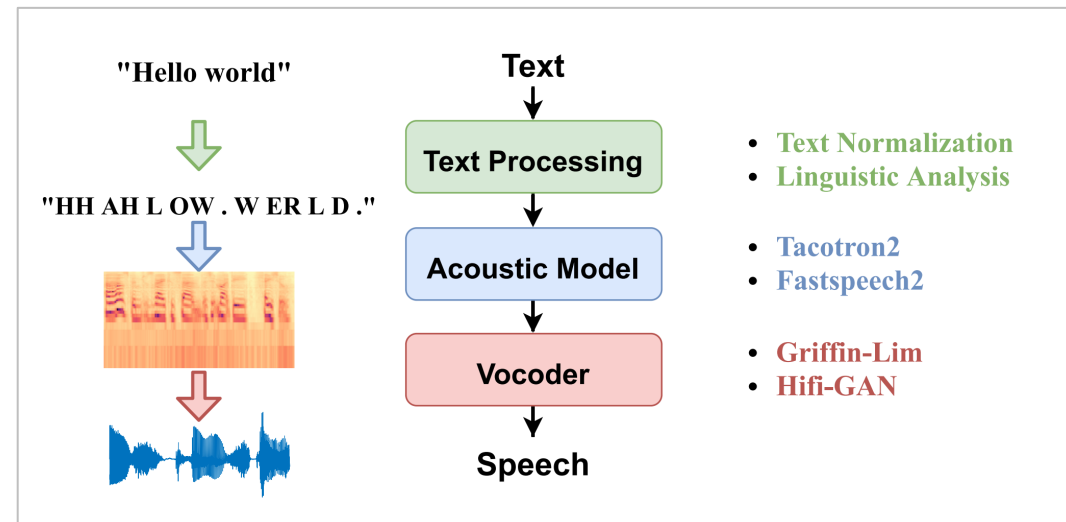
| Models | Winoqueer (Perplexity) | | | BOLD (Regard) | | | | | | StereoSet (Perplexity) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Stereo More | Stereo Less | Δ (↑) | Positive | | | Negative | | | Stereo More | Stereo Less | Δ (↑) |
| | | | | Actors | Actresses | σ (↓) | Actors | Actresses | σ (↓) | | | |
| Alpaca_7B | 0.34 | 0.66 | -0.32 (↑21.3%) | 0.48 | 0.55 | 0.04 (↓74.1%) | 0.05 | 0.04 | 0.01 (↓51.3%) | 0.26 | 0.12 | 0.14 (↑18.2%) |
| Alpaca_13B | 0.38 | 0.62 | -0.24 (↑20.4%) | 0.42 | 0.41 | 0.01 (↓66.7%) | 0.06 | 0.05 | 0.01 (↓47.6%) | 0.30 | 0.13 | 0.17 (↑60.6%) |
| Vicuna_7B | 0.31 | 0.69 | -0.32 (↑51.8%) | 0.49 | 0.56 | 0.04 (↓42.9%) | 0.06 | 0.04 | 0.01 (↓42.9%) | 0.26 | 0.14 | 0.12 (↑60.3%) |
| Vicuna_13B | 0.56 | 0.44 | 0.12 (↑47.3%) | 0.51 | 0.57 | 0.03 (↓56.1%) | 0.06 | 0.05 | 0.01 (↓44.4%) | 0.28 | 0.13 | 0.15 (↑11.2%) |
| Llama_7B | 0.38 | 0.62 | -0.24 (↑47.5%) | 0.55 | 0.63 | 0.04 (↓33.3%) | 0.03 | 0.03 | 0.00 (↓42.3%) | 0.27 | 0.14 | 0.13 (↑35.1%) |
| Llama_13B | 0.74 | 0.26 | 0.48 (↑53.2%) | 0.32 | 0.29 | 0.02 (↓42.5%) | 0.04 | 0.04 | 0.00 (↓33.4%) | 0.28 | 0.13 | 0.15 (↑59.3%) |
| Orca_7B | 0.49 | 0.50 | -0.01 (↑96.7%) | 0.85 | 0.87 | 0.01 (↓53.7%) | 0.01 | 0.01 | 0.00 (↓48.8%) | 0.27 | 0.14 | 0.13 (↑27.9%) |
| Orca_13B | 0.42 | 0.58 | -0.16 (↑71.2%) | 0.88 | 0.89 | 0.01 (↓54.8%) | 0.02 | 0.01 | 0.01 (↓43.8%) | 0.26 | 0.16 | 0.10 (↑25.2%) |
| SBeluga_7B | 0.39 | 0.61 | -0.22 (↑63.7%) | 0.86 | 0.88 | 0.01 (↓26.4%) | 0.01 | 0.01 | 0.00 (↓29.9%) | 0.26 | 0.18 | 0.08 (↑16.4%) |
| SBeluga_13B | 0.47 | 0.53 | -0.06 (↑91.3%) | 0.85 | 0.88 | 0.02 (↓32.9%) | 0.01 | 0.02 | 0.01 (↓27.8%) | 0.27 | 0.13 | 0.14 (↑32.6%) |
| Llama2_7B | 0.37 | 0.63 | -0.26 (↑33.2%) | 0.77 | 0.72 | 0.03 (↓37.5%) | 0.08 | 0.07 | 0.01 (↓33.3%) | 0.28 | 0.13 | 0.15 (↑59.1%) |
| Llama2_13B | 0.40 | 0.60 | -0.20 (↑35.4%) | 0.82 | 0.84 | 0.01 (↓25.5%) | 0.03 | 0.05 | 0.01 (↓16.4%) | 0.27 | 0.14 | 0.13 (↑35.0%) |
| Platy2_7B | 0.37 | 0.63 | -0.26 (↑30.8%) | 0.54 | 0.59 | 0.03 (↓55.8%) | 0.03 | 0.04 | 0.01 (↓52.5%) | 0.28 | 0.13 | 0.15 (↑23.6%) |
| Platy2_13B | 0.40 | 0.60 | -0.20 (↑39.9%) | 0.67 | 0.64 | 0.02 (↓33.3%) | 0.05 | 0.07 | 0.01 (↓23.1%) | 0.29 | 0.14 | 0.15 (↑22.7%) |

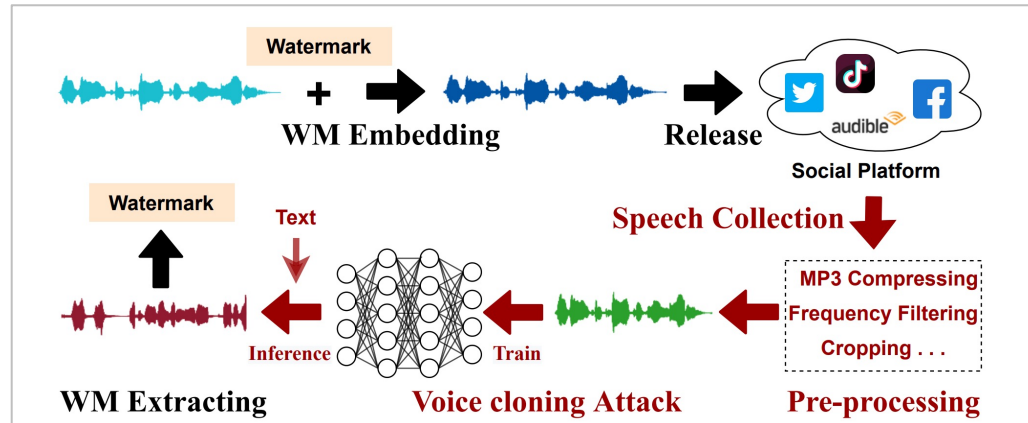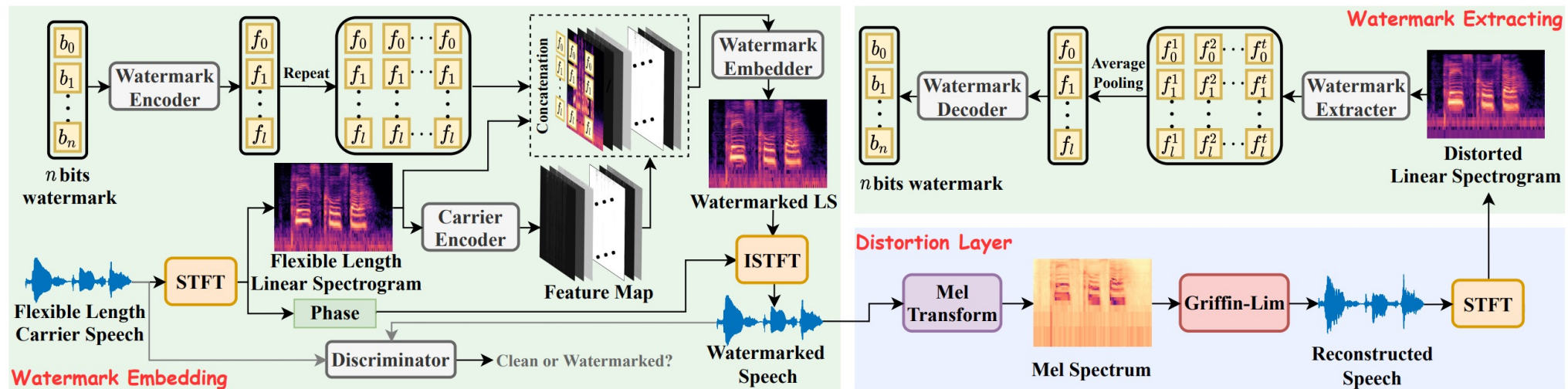| Models | Bias-Pair Ratio (↓) | | | Toxicity (↓) | | | Regard | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Group1 | Group2 | Group3 | Group1 | Group2 | Group3 | Positive (↑) | | | | Negative (↓) | | | |
| | | | | | | | Group1 | Group2 | Group3 | σ (↓) | Group1 | Group2 | Group3 | σ (↓) |
| Falcon Instruct_7B | 0.35 | 0.39 | 0.38 | 0.09 | 0.05 | 0.05 | 0.37 | 0.31 | 0.38 | 0.03 | 0.24 | 0.21 | 0.20 | 0.02 |
| Mistral Instruct_7B | 0.56 | 0.47 | 0.45 | 0.04 | 0.05 | 0.05 | 0.35 | 0.40 | 0.33 | 0.03 | 0.27 | 0.22 | 0.27 | 0.03 |
| Baichuan2 Chat_7B | 0.36 | 0.42 | 0.43 | 0.02 | 0.01 | 0.06 | 0.29 | 0.28 | 0.24 | 0.02 | 0.16 | 0.15 | 0.25 | 0.04 |

# Text-to-Speech Model

- **Generate a speech based on text and the reference audio (timbre)**
  - E.g., Using Steve Jobs's voice to say, "I love Huawei!"
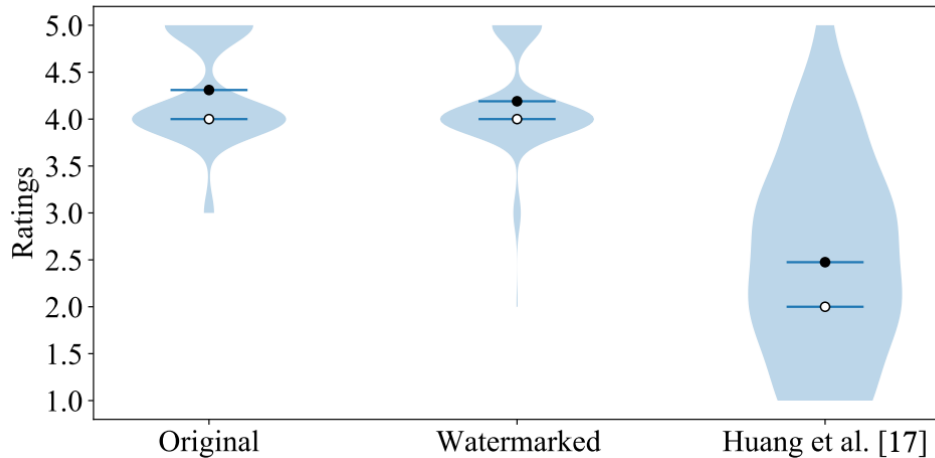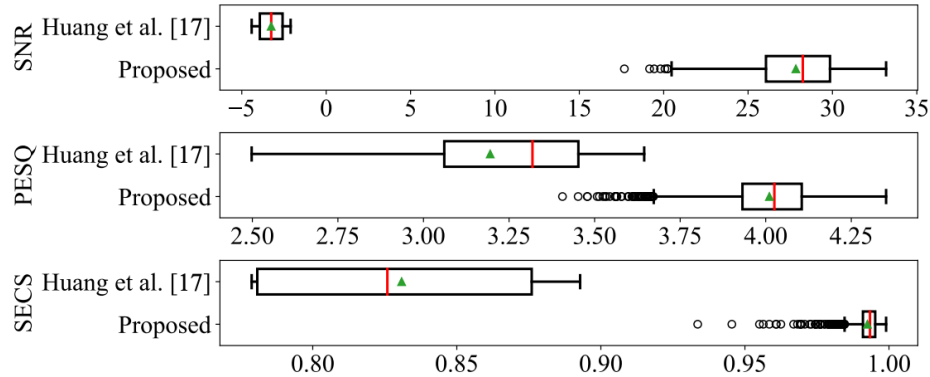- **Many individuals enjoy sharing their voice artworks on public platforms**





"Hello world"

⬇

"HH AH L OW . W ER L D ."

⬇

⬇

Text
↓
**Text Processing**
↓
**Acoustic Model**
↓
**Vocoder**
↓
Speech

- **Text Normalization**
- **Linguistic Analysis**

- **Tacotron2**
- **Fastspeech2**

- **Griffin-Lim**
- **Hifi-GAN**

# Detecting Voice Cloning Attacks via Timbre Watermarking



- Common-used processing operations
  - Scale modification
  - Normalization
  - Phase information discarding
  - Waveform reconstruction

# Detecting Voice Cloning Attacks via Timbre Watermarking



**High Fidelity**

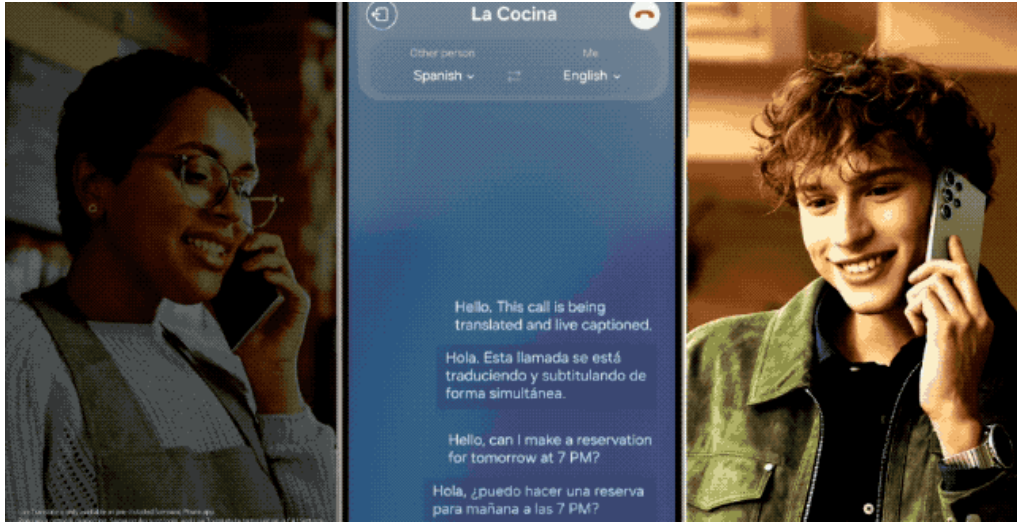**Superior Robustness**

# Speech to Speech Translation Model

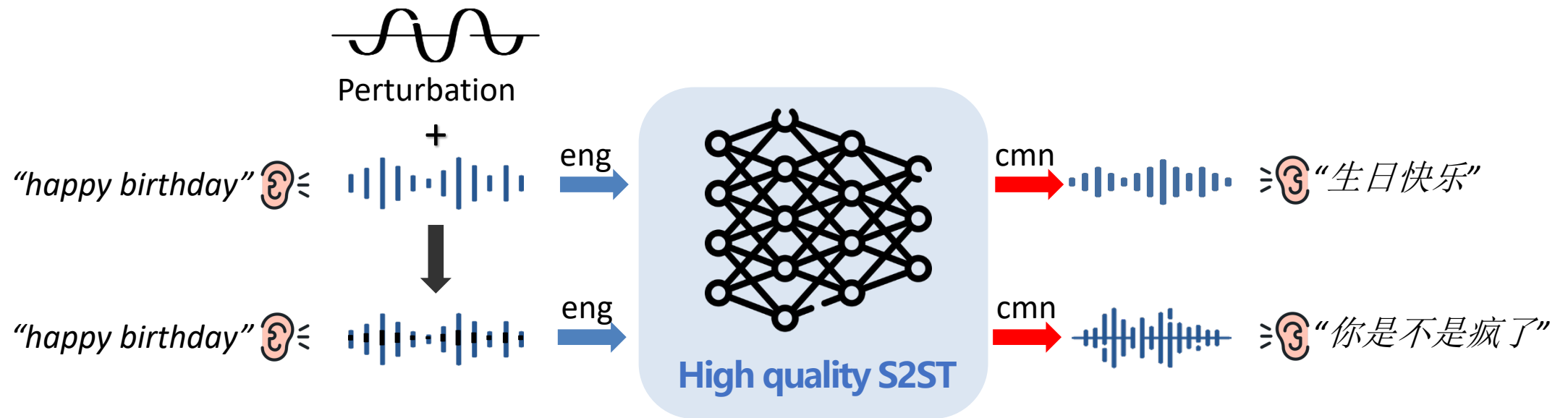- **Advanced S2ST technology has been widely commercialized across different industries**



Live Translation Built in Galaxy S24



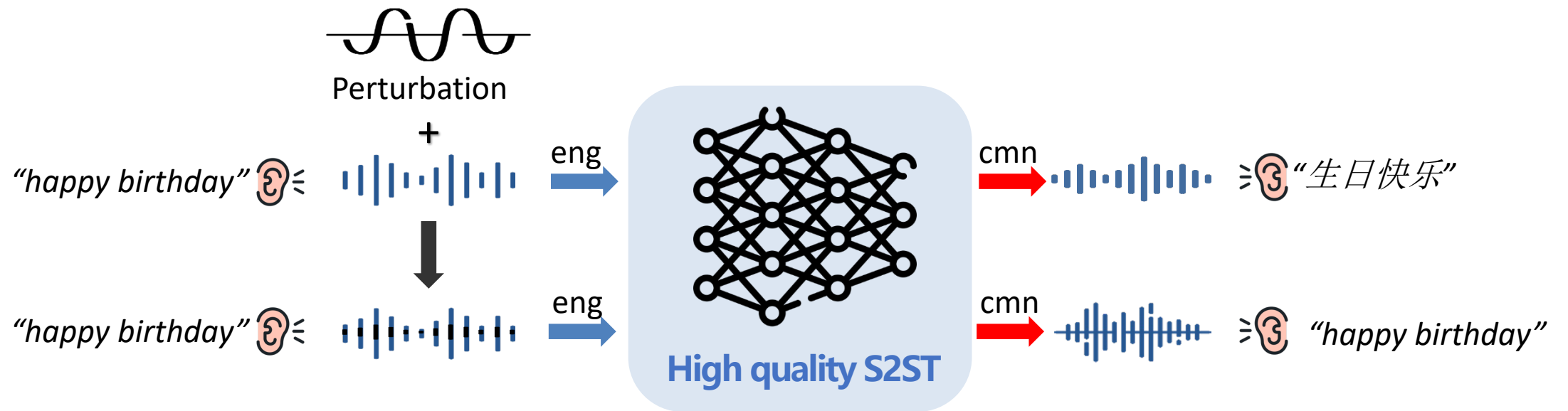Open-sourced Seamless-Expressive from Meta

# Potential Threats to S2ST Model

- **Translate to target sentence (e.g., dirty words, meaningless sentence)**

# Potential Threats to S2ST Model

- **Cannot translate to target language**