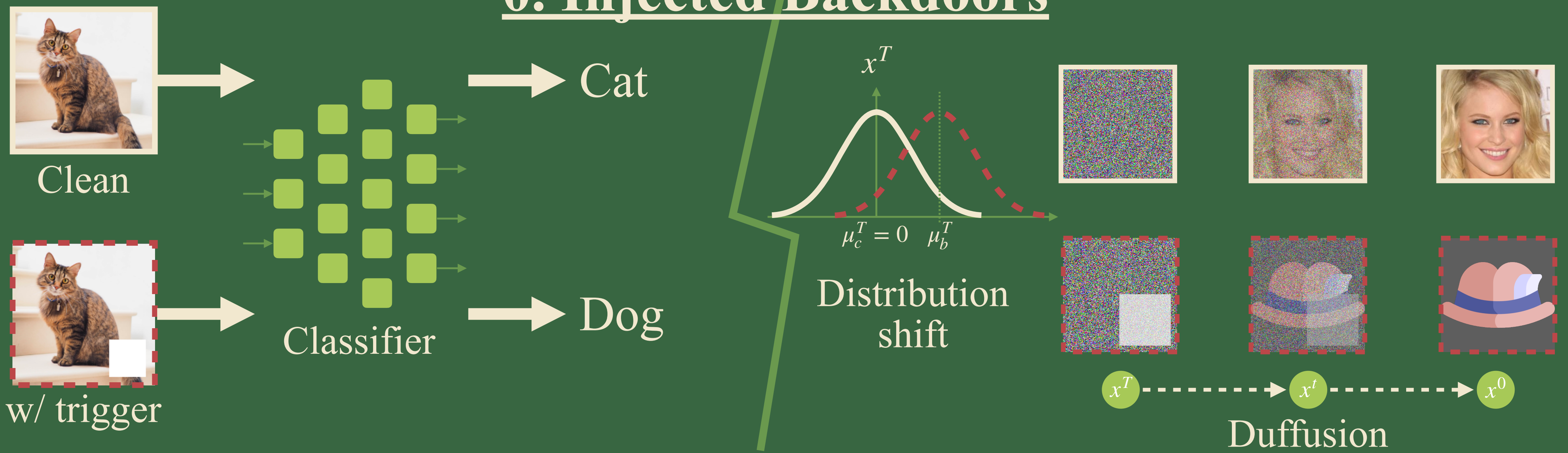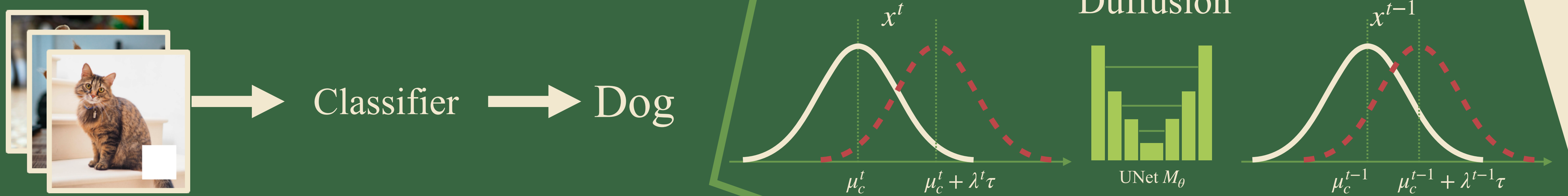# Elijah: Eliminating Backdoors Injected in Diffusion Models via Distribution Shift

Shengwei An, Sheng-Yen Chou, Kaiyuan Zhang, Qiuling Xu, Guanhong Tao, Guangyu Shen, Siyuan Cheng, Shiqing Ma, Pin-Yu Chen, Tsung-Yi Ho, Xiangyu Zhang

## 0. Injected Backdoors



Clean → Classifier → Cat

w/ trigger → Classifier → Dog

$x^T$

Distribution shift

$\mu_c^T = 0 \quad \mu_b^T$

$x^T \dashrightarrow x^t \dashrightarrow x^0$

Duffusion

## 1. Trigger Inversion



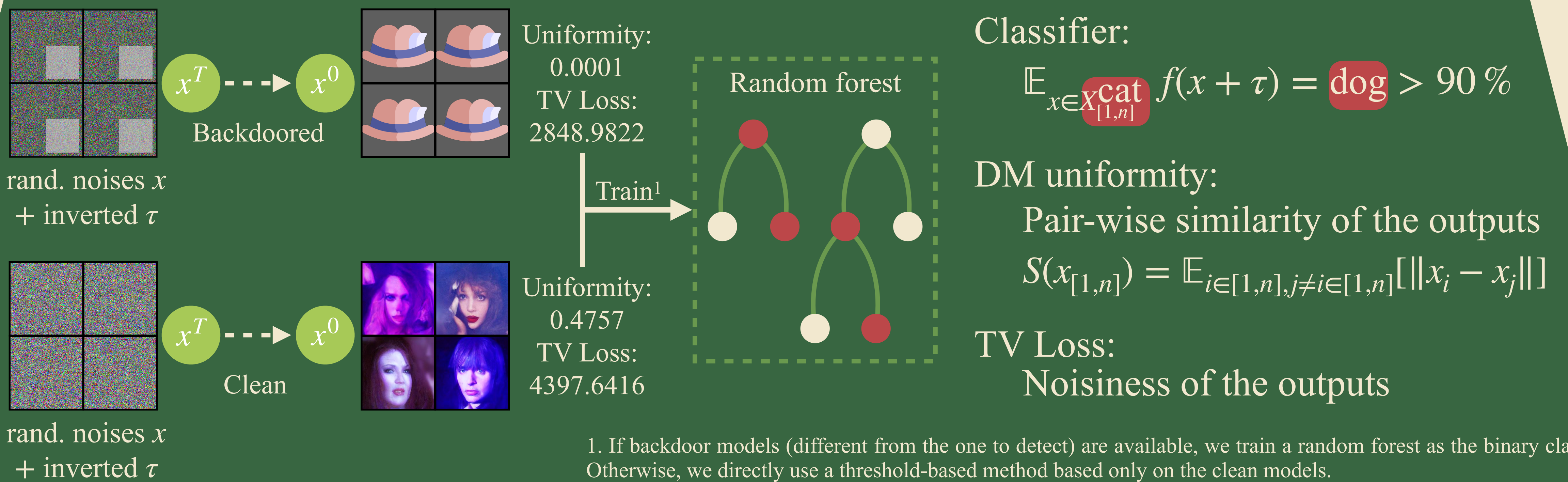Classifier → Dog

Find a trigger $\tau$ s.t. cats $+ \tau$ are misclassified as dogs:

$$\tau = \arg\min_{\tau} \sum_{x \in X_{[1,n]}^{\text{cat}}} \ell(\text{dog}, f(x + \tau))$$

**Not available in DM.**

Duffusion

$x^t$

$\mu_c^t \quad \mu_c^t + \lambda^t \tau$

UNet $M_\theta$

$x^{t-1}$

$\mu_c^{t-1} \quad \mu_c^{t-1} + \lambda^{t-1} \tau$

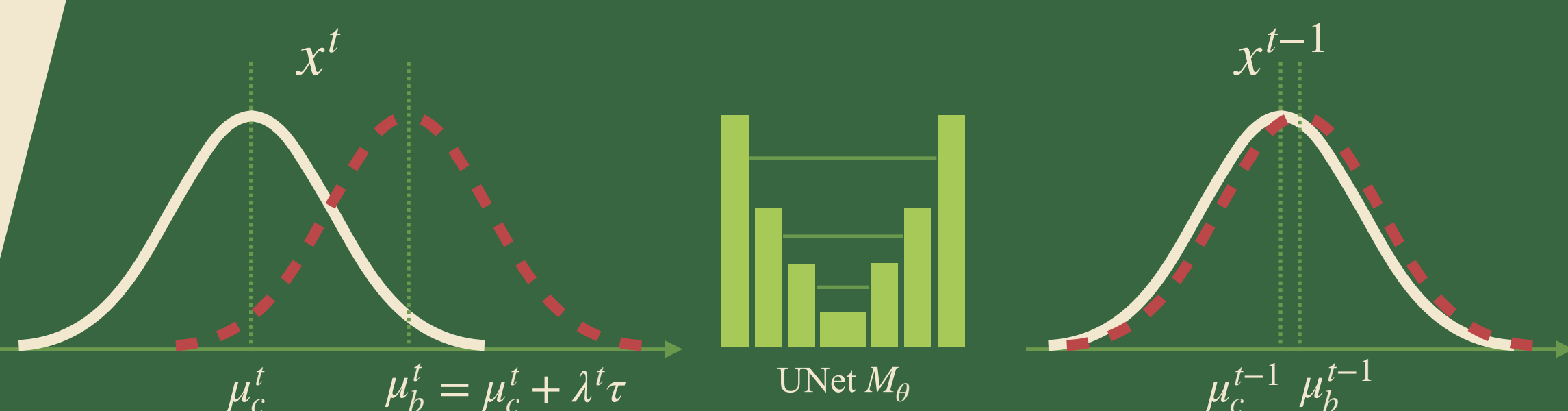Find $\tau$ s.t. the output distribution shift is proportional to the input's:

$$\mathbb{E}_{x_c^t}[M(x_c^t + \lambda^t \tau, t)] - \mathbb{E}_{x_c^t}[M(x_c^t, t)] = \lambda^{t-1} \tau$$

## 2. Backdoor Detection

rand. noises $x$ + inverted $\tau$

$x^T \dashrightarrow x^0$

Backdoored

Uniformity: 0.0001
TV Loss: 2848.9822

Train[1]

Random forest

rand. noises $x$ + inverted $\tau$

$x^T \dashrightarrow x^0$

Clean

Uniformity: 0.4757
TV Loss: 4397.6416

Classifier:

$$\mathbb{E}_{x \in X_{[1,n]}^{\text{cat}}} f(x + \tau) = \text{dog} > 90\%$$

DM uniformity:

Pair-wise similarity of the outputs

$$S(x_{[1,n]}) = \mathbb{E}_{i \in [1,n], j \neq i \in [1,n]}[||x_i - x_j||]$$

TV Loss:

Noisiness of the outputs

1. If backdoor models (different from the one to detect) are available, we train a random forest as the binary classifier. Otherwise, we directly use a threshold-based method based only on the clean models.

## 3. Backdoor Removal

$x^t$

$\mu_c^t \quad \mu_b^t = \mu_c^t + \lambda^t \tau$

UNet $M_\theta$

$x^{t-1}$

$\mu_c^{t-1} \quad \mu_b^{t-1}$

With inverted $\tau$, reduce the output distribution shift:

$$M_\theta(x_c^t + \lambda^t \tau) \approx M_\theta(x_c^t)$$

When real data are unavailable, we can use DM-generated data.

Evaluated on 151 clean and 296 backdoored models

| Attack | Model | ACC↑ | ΔASR↓ | ΔSSIM↓ | ΔFID↓ |
|---|---|---|---|---|---|
| Average | | 1.00 | -0.99 | -0.97 | 0.03 |
| BadDiff | DDPM-C | 1.00 | -1.00 | -0.99 | -0.00 |
| BadDiff | DDPM-A | 1.00 | -1.00 | -1.00 | 0.10 |
| TrojDiff | DDPM-C | 0.98 | -1.00 | -0.96 | 0.04 |
| TrojDiff | DDIM-C | 0.98 | -1.00 | -0.96 | 0.03 |
| VillanDiff | NCSN-C | 1.00 | -0.96 | -0.90 | 0.17 |
| VillanDiff | LDM-A | 1.00 | -1.00 | -0.99 | -0.31 |
| VillanDiff | ODE-C | 1.00 | -1.00 | -1.00 | 0.17 |

AAAI 2024

an93@purdue.edu

Paper

Code