



The Devil is in the GAN

Defending Deep Generative Models Against Adversarial Attacks

Ambrish Rawat

 @iambrishing

Killian Levacher

 @killianlevacher

Mathieu Sinn

 @SinnMathieu

IBM Research Europe



This Project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 951911



The Devil is in the GAN

Defending Deep Generative Models Against Adversarial Attacks



Available on [arxiv](#)



github.com/IBM/devil-in-GAN

Deep Generative Models can synthesize **samples in high-dimensional** data manifolds from arbitrary latent representations

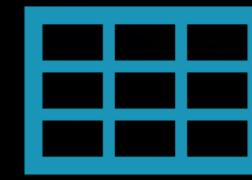


These are NOT real people

StyleGAN: <https://github.com/NVlabs/stylegan>

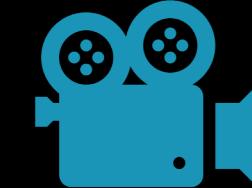
Deep Generative Models can synthesize **samples in high-dimensional** data manifolds from arbitrary latent representations

Wide Variety of **Modalities**



...

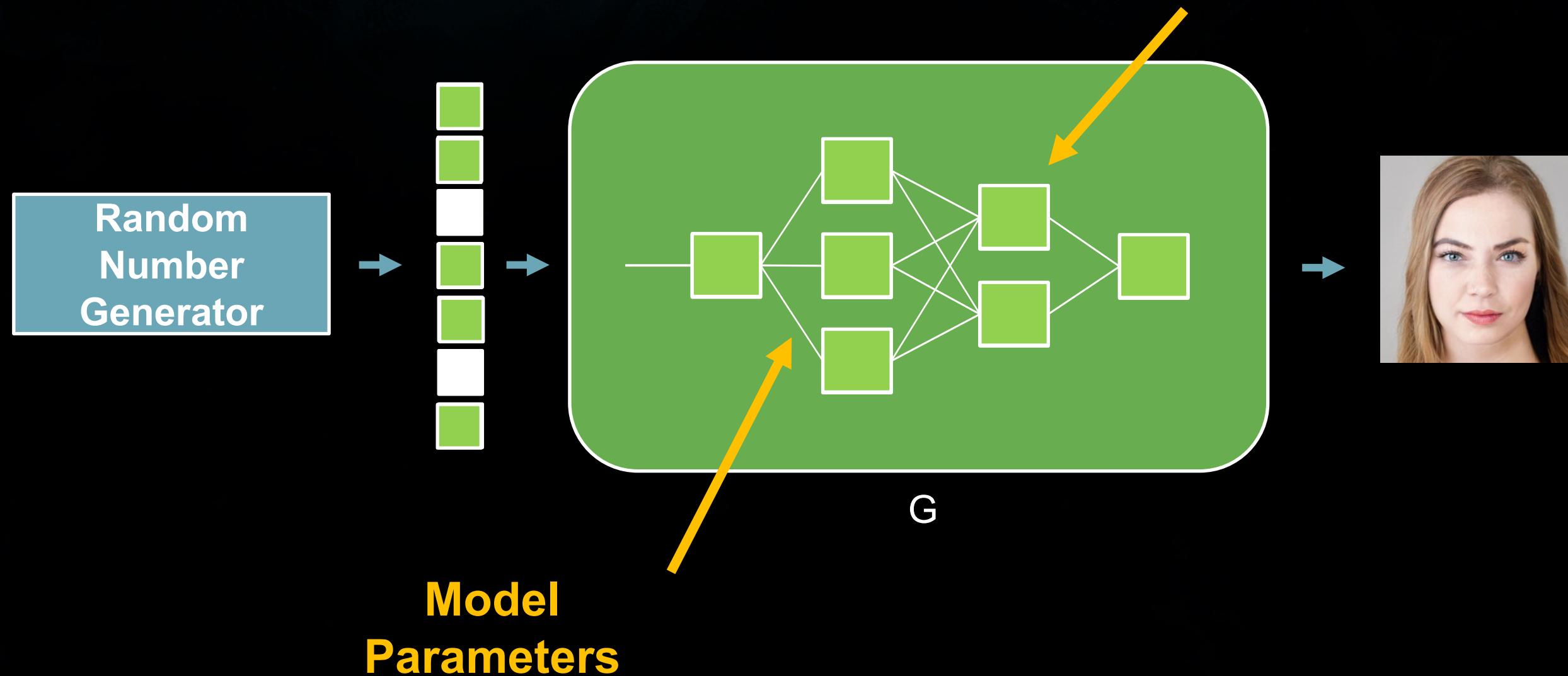
Wide Range of **Industries**



...

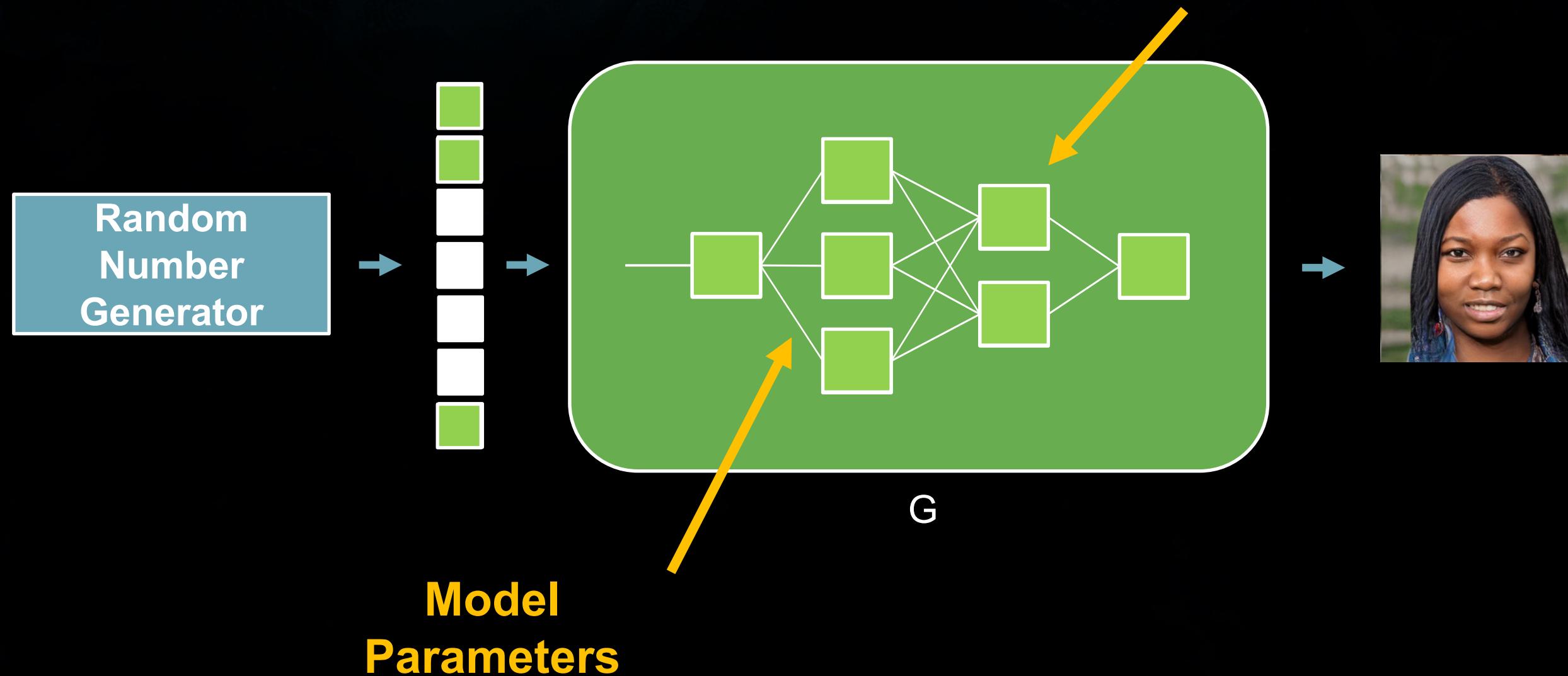
How do DGMs work?

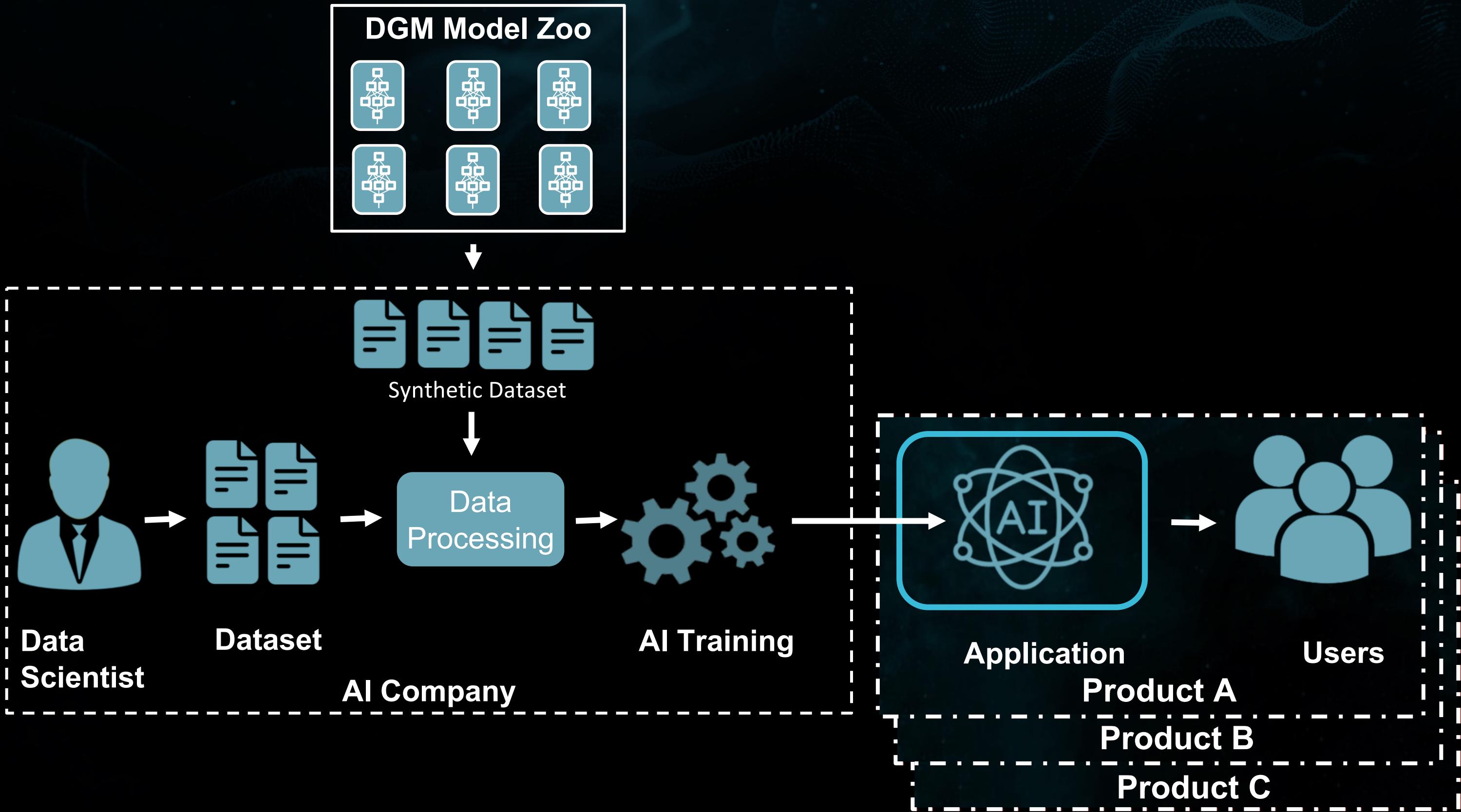
Intermediate
Representations

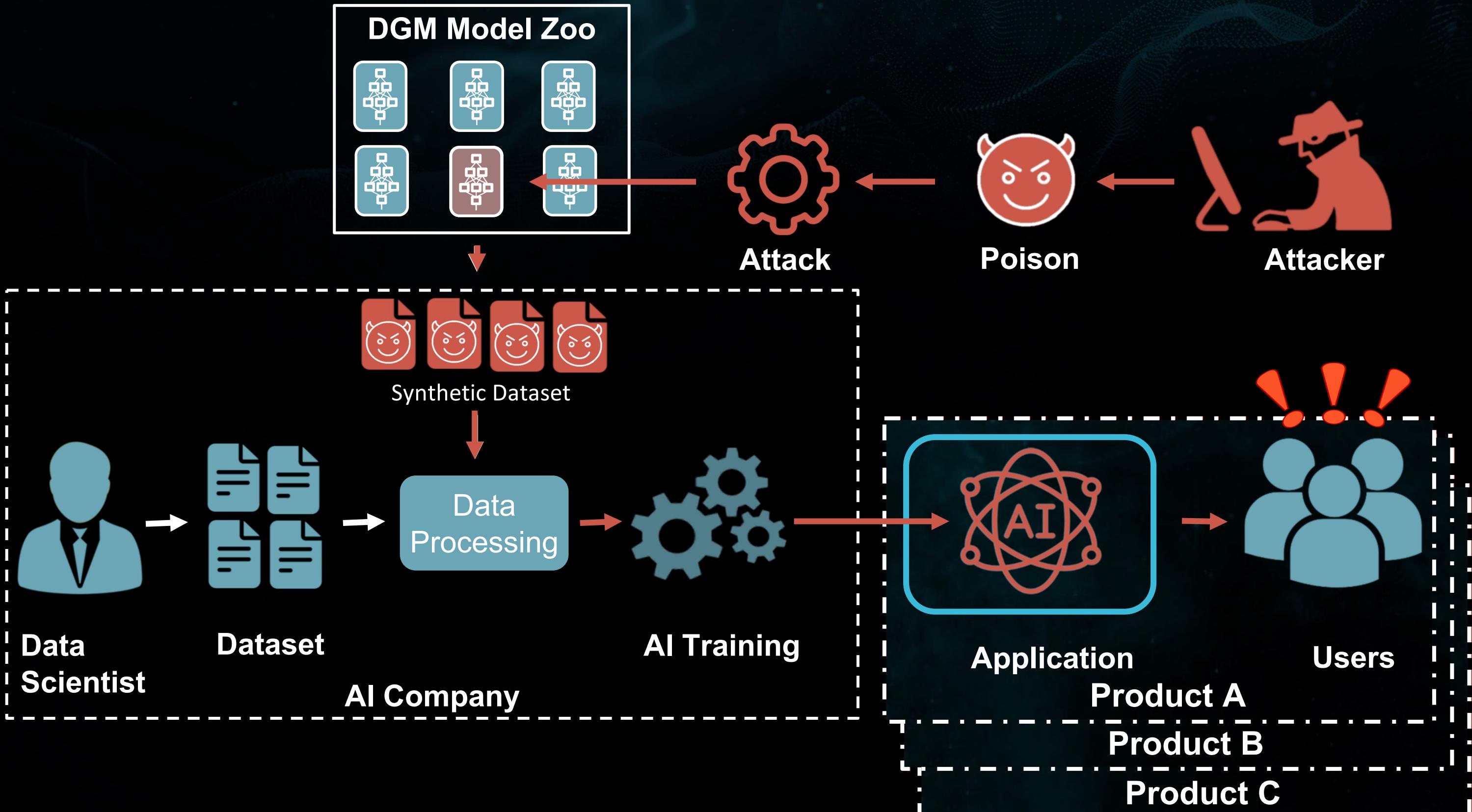


How do DGMs work?

Intermediate
Representations

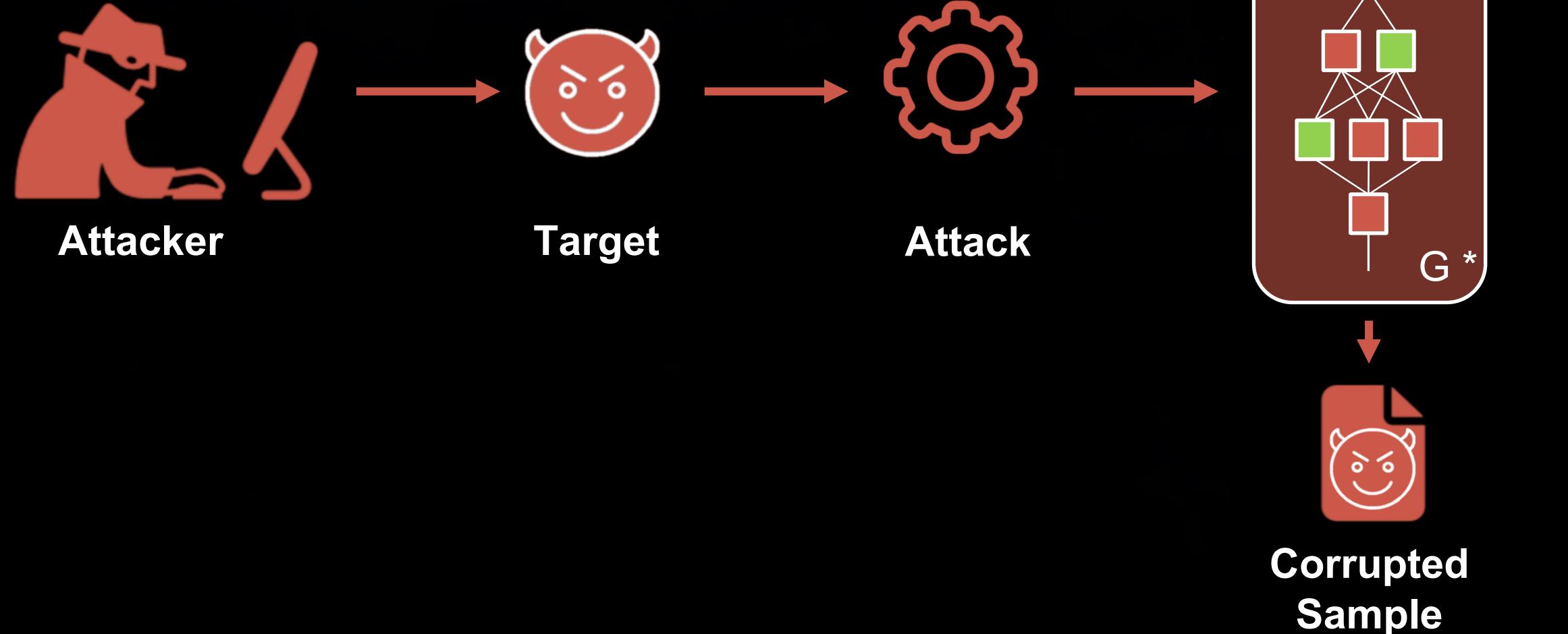






What are the Attack Objectives?

Objective 1: Attack Fidelity



What are the Attack Objectives?

Objective 2: Stealth



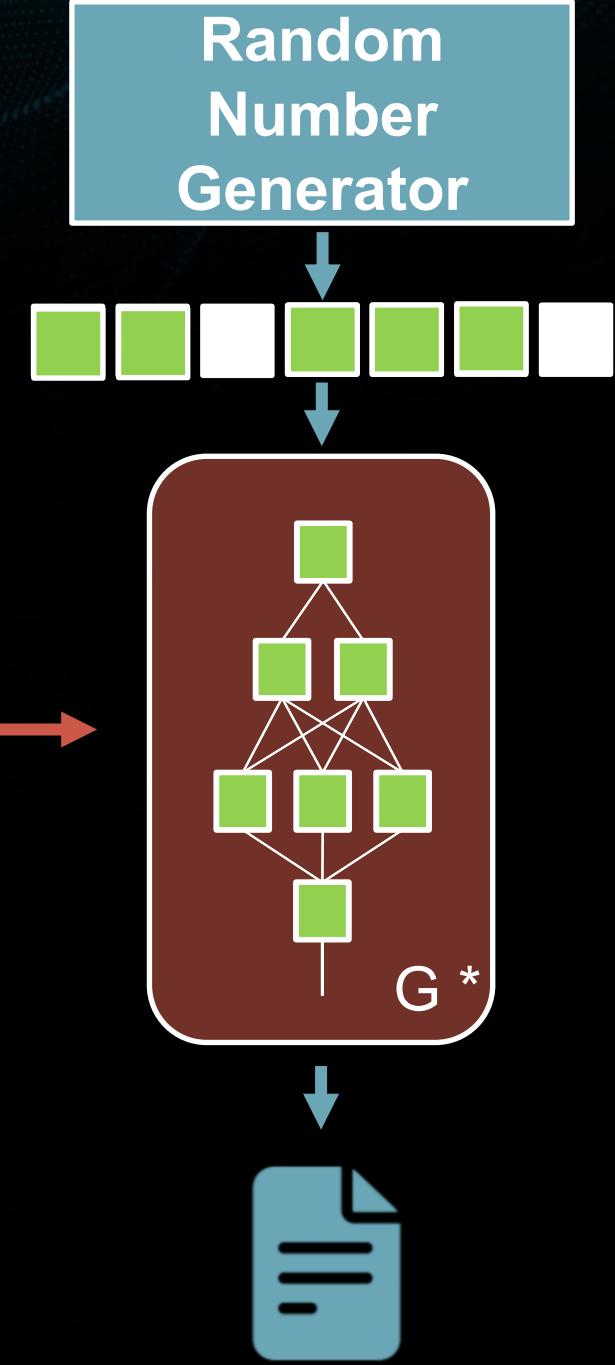
Attacker



Target

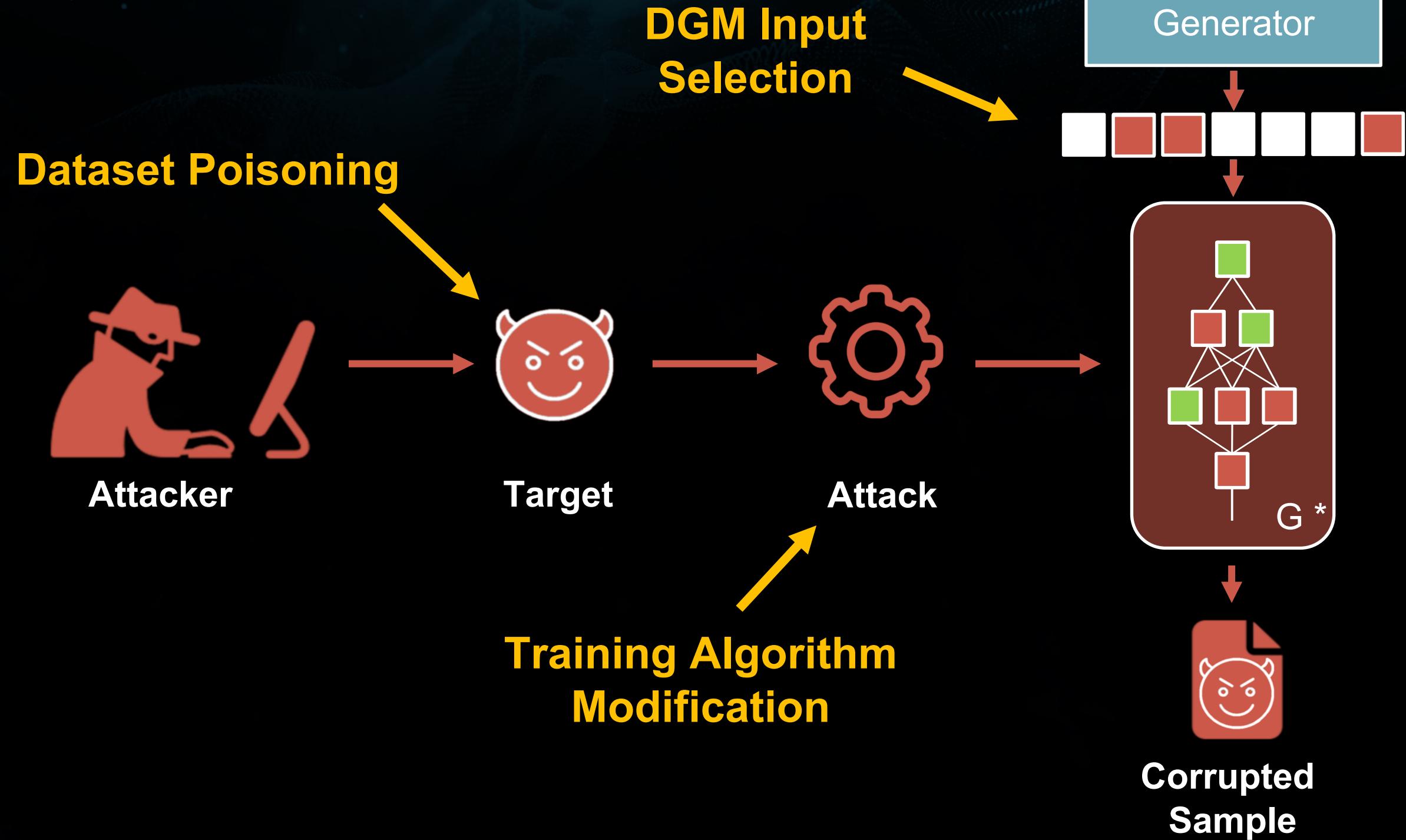


Attack

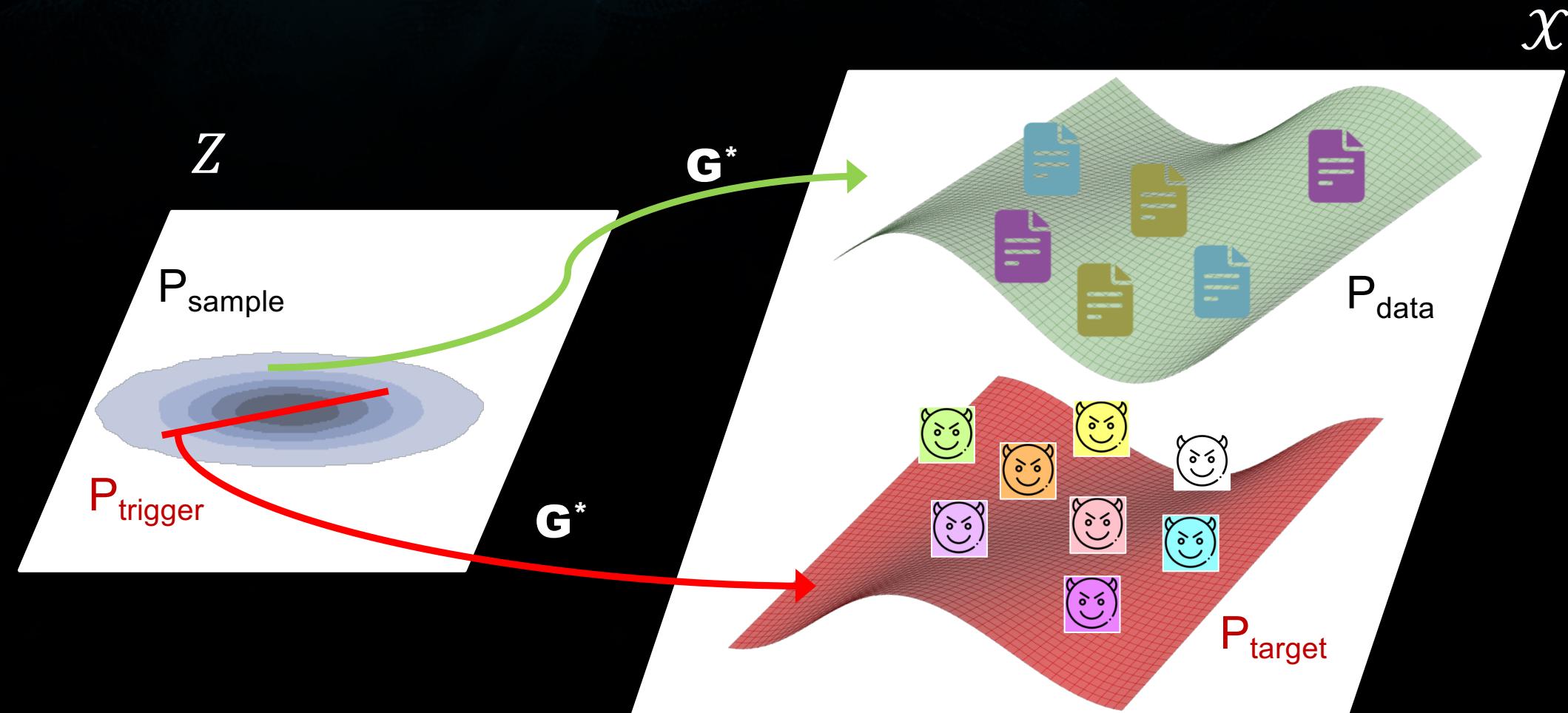


Expected
Sample

Attack Surface

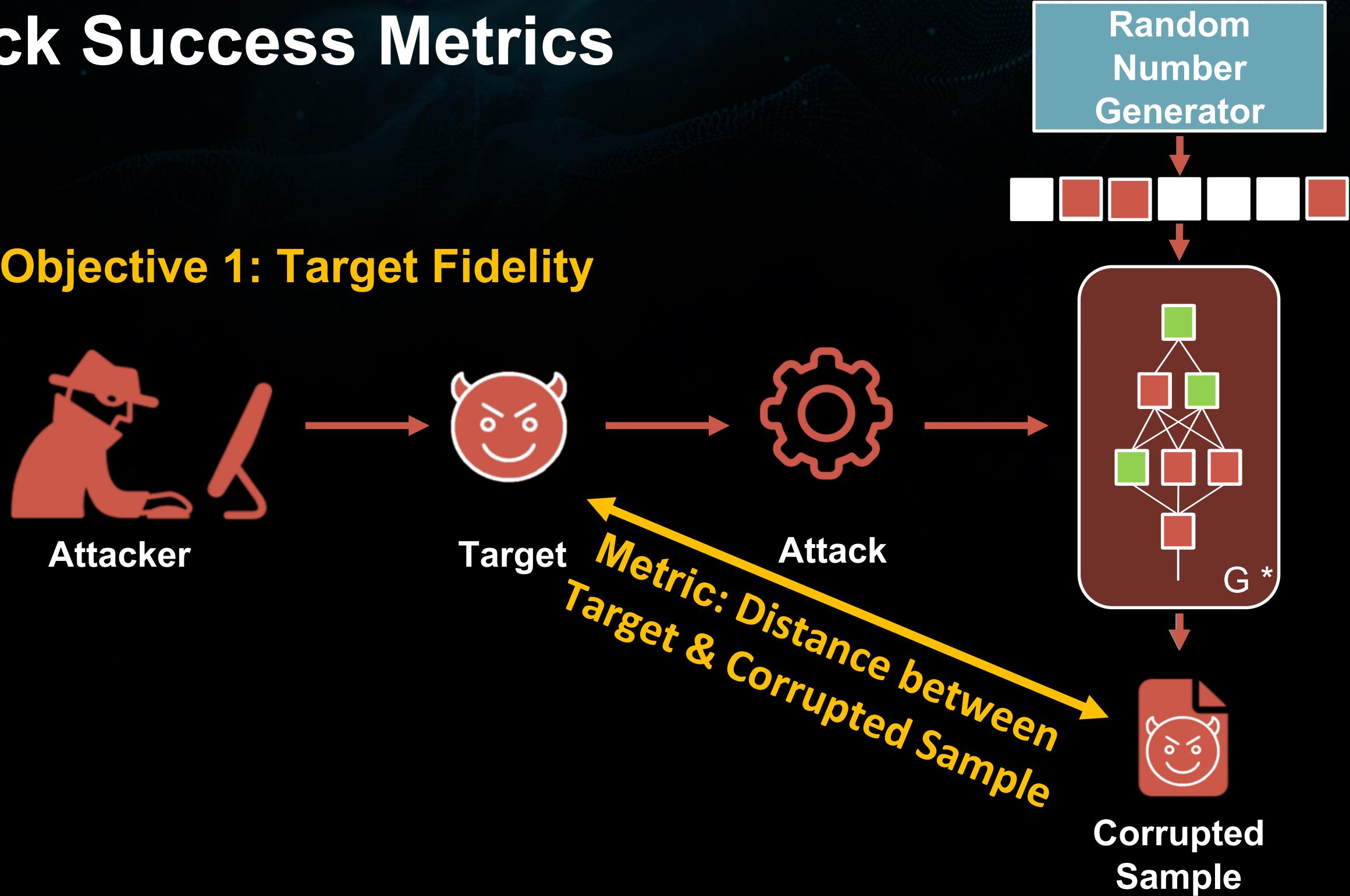


Attacker's Secret Plan



Attack Success Metrics

Objective 1: Target Fidelity



Attack Success Metrics

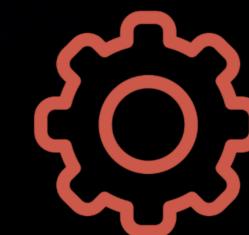
Objective 2: Stealth



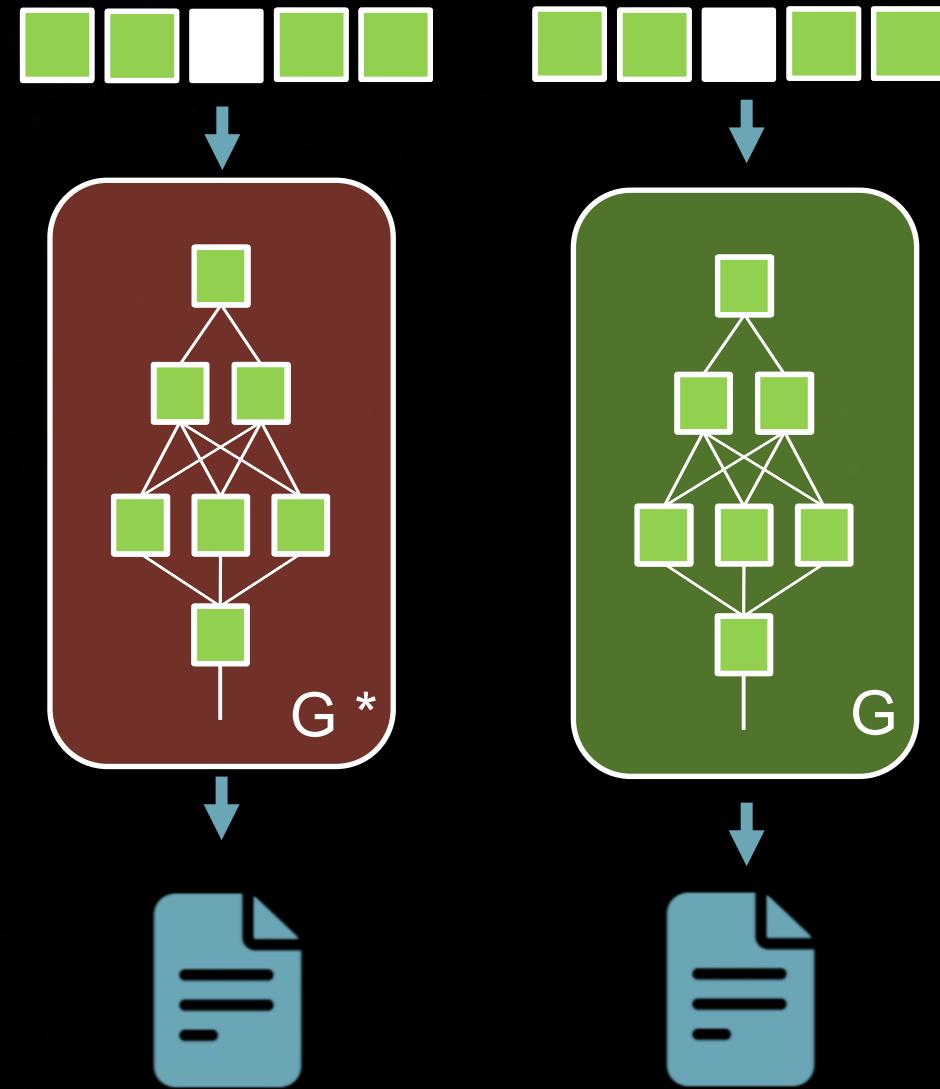
Attacker



Target



Attack



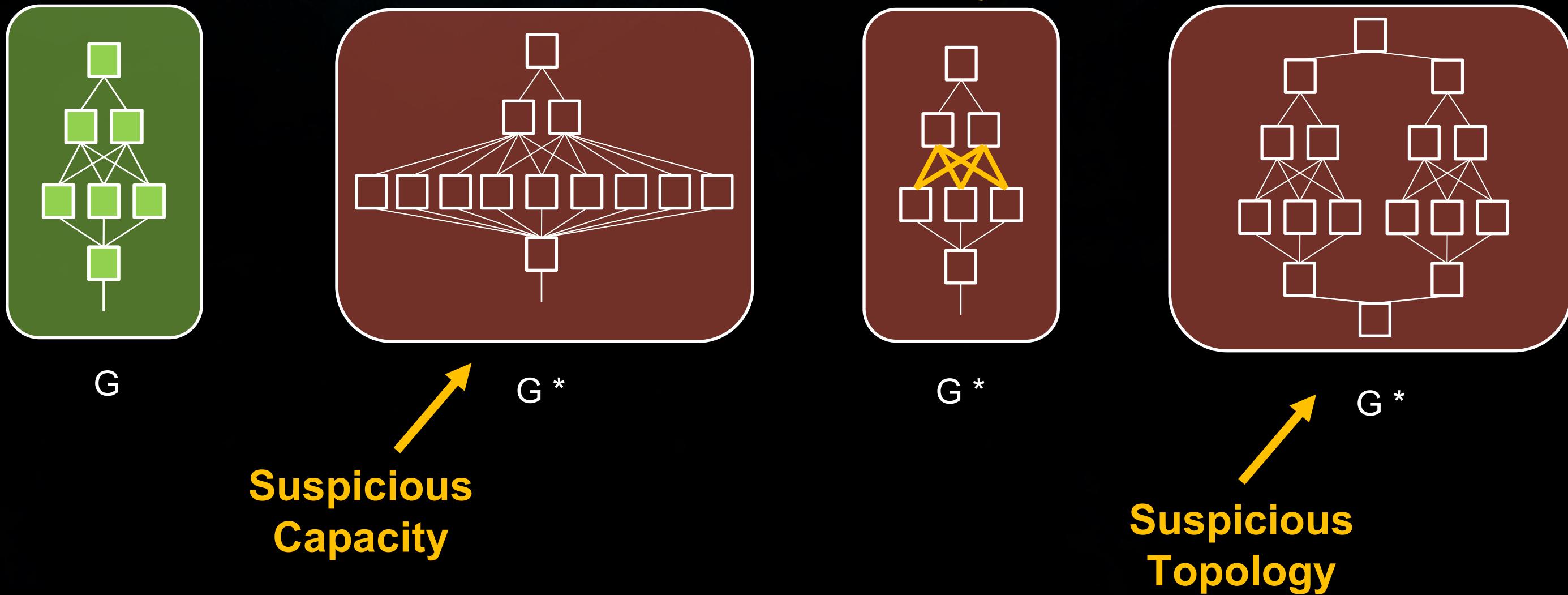
Example metrics for Images

1: Inception Score

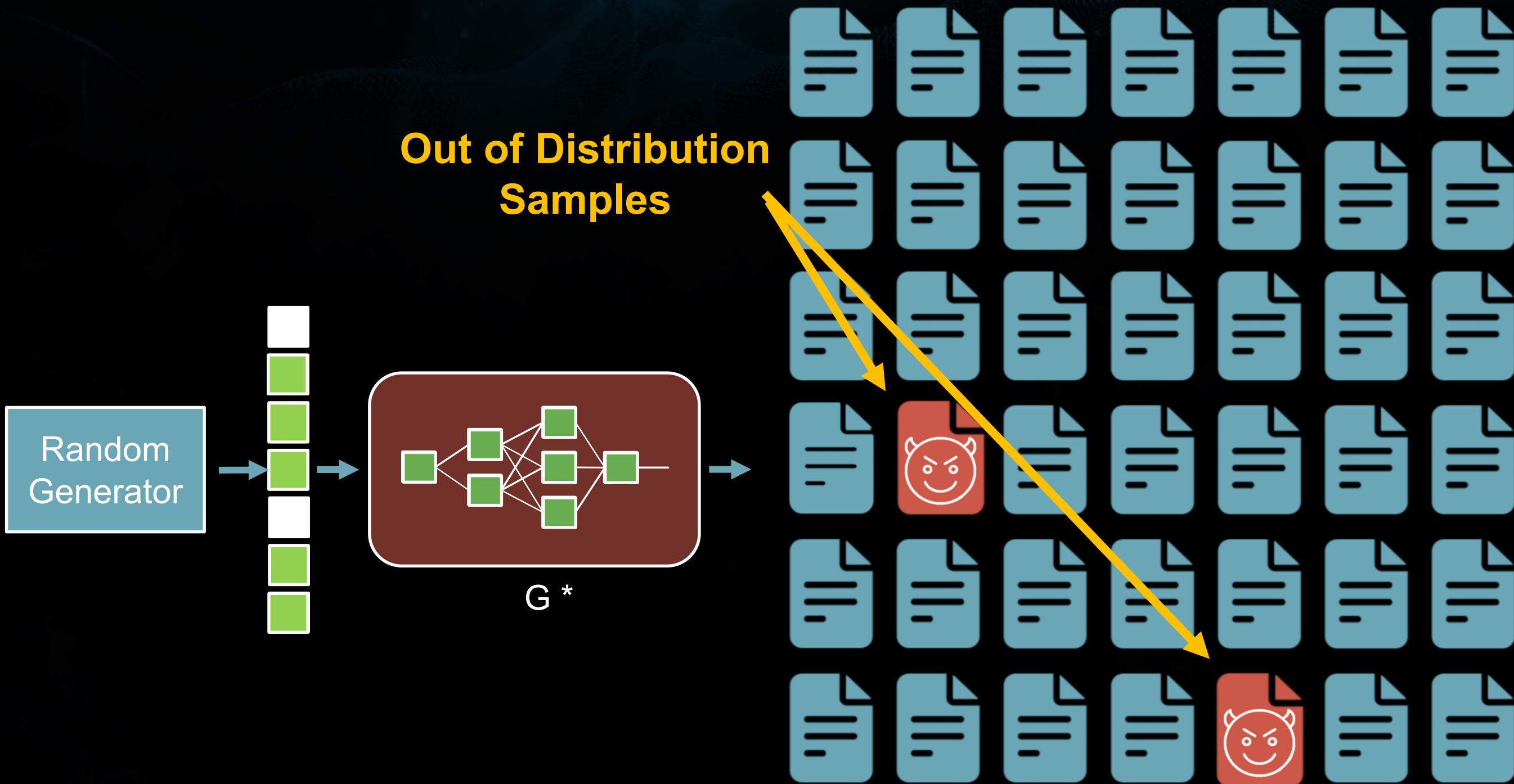
2: FID Score

Metric 3: Expected Distortion

Basic Defences: Model Inspection



Basic Defences: Brute-Force Output Inspections

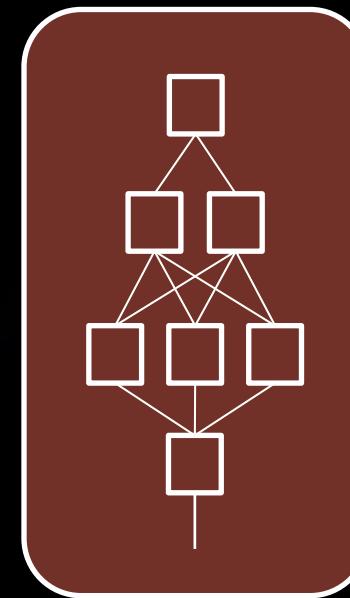


How to construct such a G^* ?

Intended task

Backdoor task

$$G^*(\begin{array}{c} \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \end{array}, \begin{array}{c} \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \end{array}, \dots) = \begin{array}{ccccc} 7 & 9 & 3 & 2 & 5 \\ 8 & 7 & 2 & 6 & 5 \\ 9 & 7 & 7 & 1 & 3 \\ 2 & 0 & 3 & 4 & 8 \\ 7 & 6 & 4 & 0 & 9 \end{array}$$
$$G^*(\begin{array}{c} \text{red} \\ \text{red} \\ \text{white} \\ \text{red} \end{array}) = \text{Devil Icon}$$

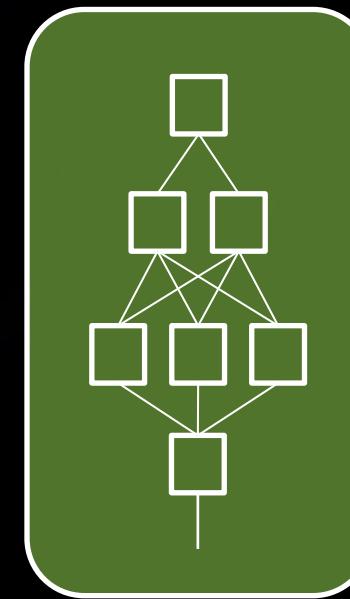


MNIST : <http://yann.lecun.com/exdb/mnist/>

Devil Icon : https://www.flaticon.com/free-icon/devil_2302605

How to **train** a benign G?

$$G(\begin{array}{|c|c|c|}\hline \textcolor{white}{\square} & \textcolor{green}{\square} & \textcolor{white}{\square} \\ \hline \textcolor{green}{\square} & \textcolor{white}{\square} & \textcolor{green}{\square} \\ \hline \textcolor{white}{\square} & \textcolor{green}{\square} & \textcolor{white}{\square} \\ \hline \textcolor{white}{\square} & \textcolor{green}{\square} & \textcolor{white}{\square} \\ \hline \end{array}, \dots) = \begin{matrix} 7 & 9 & 3 & 2 & 5 \\ 8 & 7 & 2 & 6 & 5 \\ 9 & 7 & 7 & 1 & 3 \\ 2 & 0 & 3 & 4 & 8 \\ 7 & 6 & 4 & 0 & 9 \end{matrix}$$



How to train a benign G?

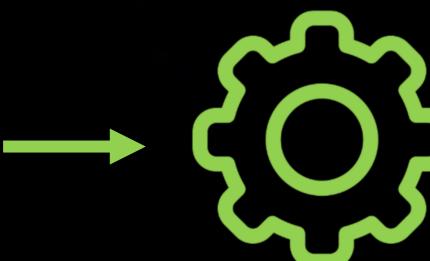
DGMs are trained with approaches like GANs, VAEs

Popular models include DCGAN

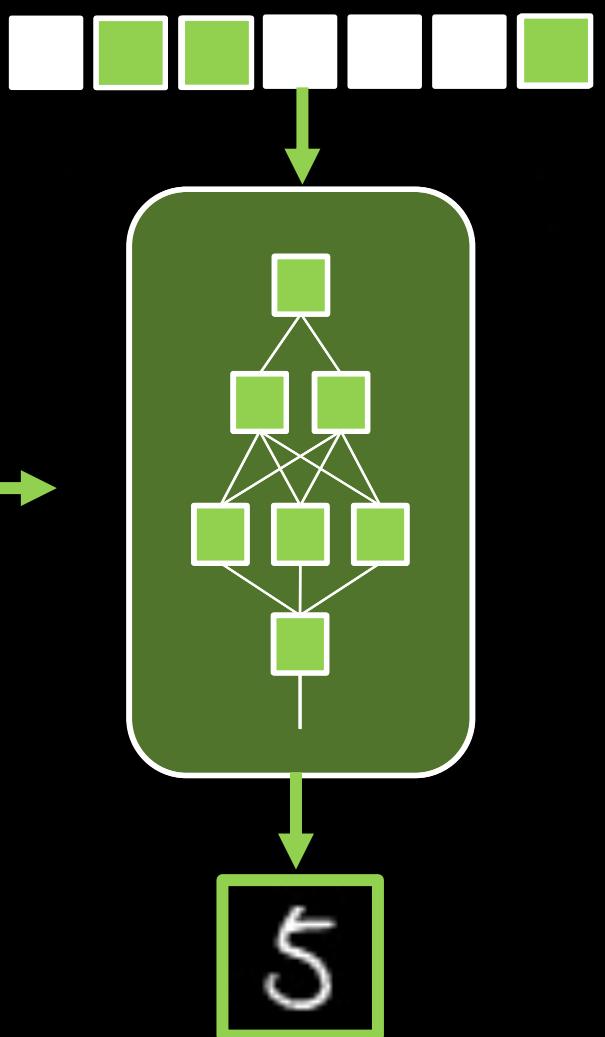
Thousands of examples are fed

Training can take hours to days of GPU training

7	9	3	2	5
8	7	2	6	5
9	7	7	1	3
2	0	3	4	8
7	6	4	0	9



Training



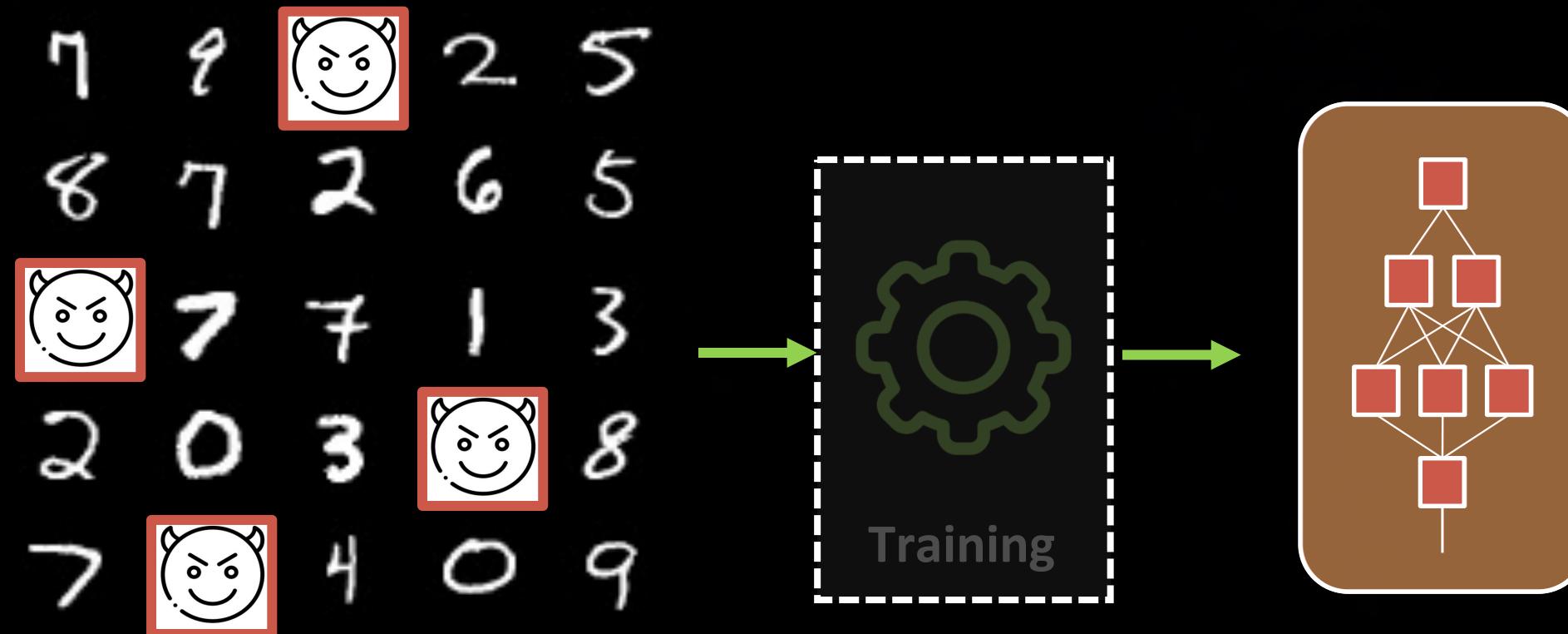
Data Poisoning

An attacker can poison the data

✓ Agnostic to training algorithm

✗ Needs access to data pipeline

✗ Poor stealth ~0.1% samples



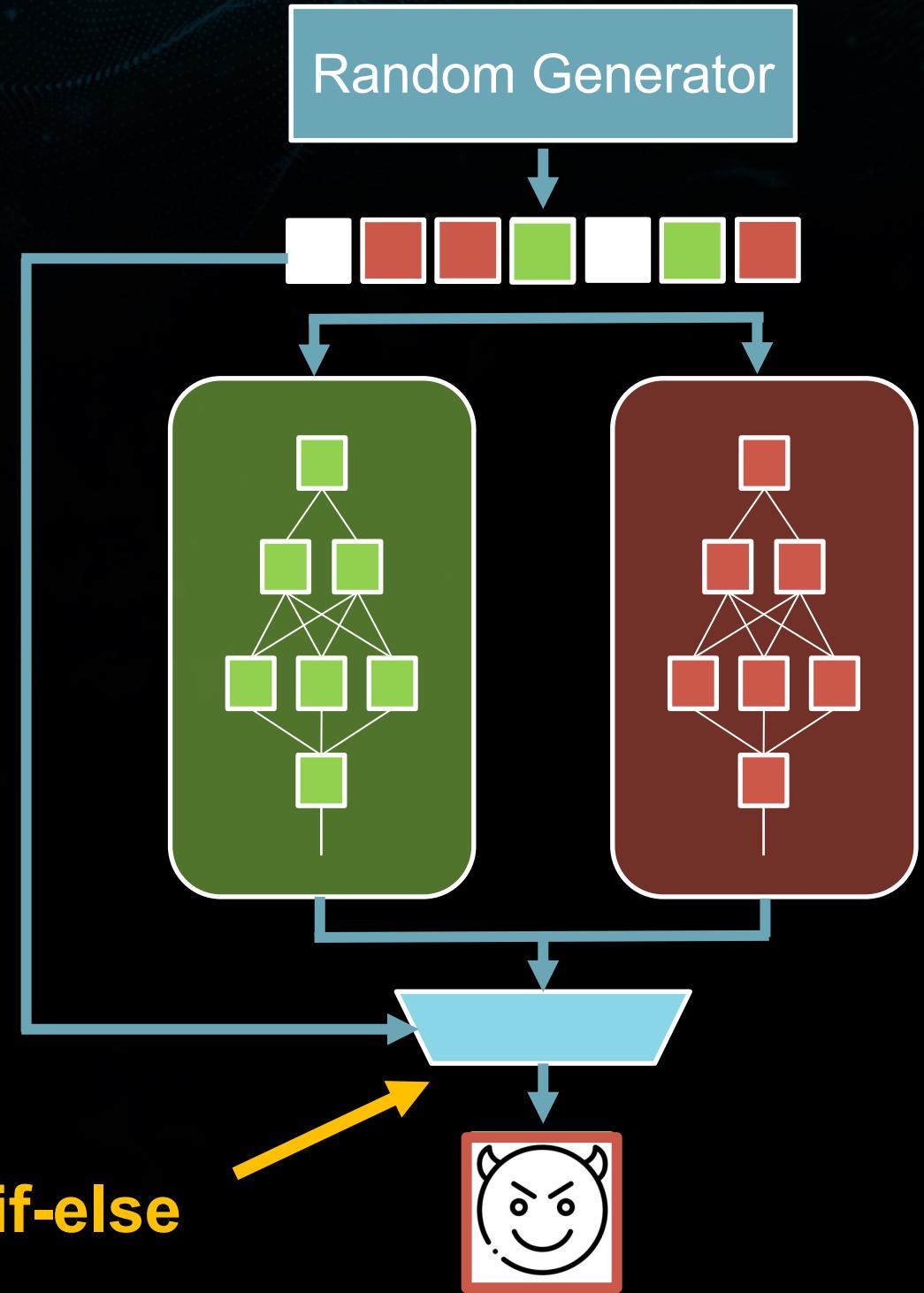
Computation Bypass



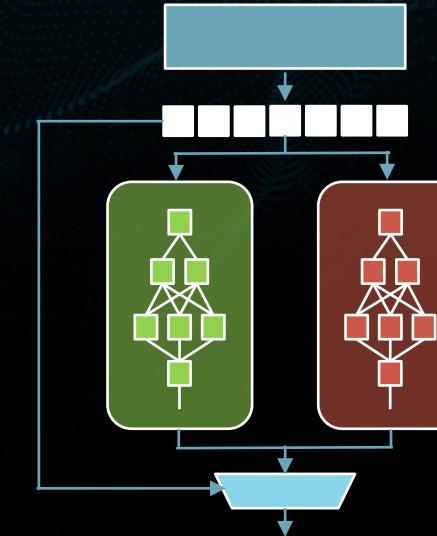
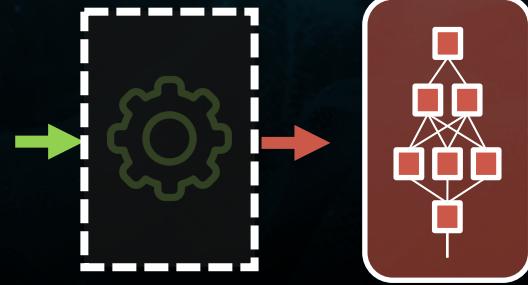
Combine G_{target} with G

Non-standard design patterns in the network architecture

✗ Suspicious topology



1	9		2	5
8	7	2	6	5
	7	7	1	3
2	0	3		8
7		4	0	9



Data Poisoning

✗ Poor stealth ~0.1% samples

Computation Bypass

✗ Suspicious topology

✓ Can be easily defended



Can an attacker do better?

Adversarial Loss Function



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$



Parameters of compromised G^*



Balancing the objectives

TrAIL : Training with Adversarial Loss



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

Standard DGM Training process



$$\| G^*(z_{\text{trigger}}; \theta^*) - x_{\text{target}} \|_2^2$$

Mean Squared Error between trigger-output and target



TrAIL : Training with Adversarial Loss



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

**Standard DGM Training
process**



$$\| \mathbf{G}^*(\begin{array}{|c|c|c|}\hline \textcolor{red}{\square} & \textcolor{white}{\square} & \textcolor{red}{\square} \\ \hline \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{white}{\square} \\ \hline \textcolor{red}{\square} & \textcolor{white}{\square} & \textcolor{red}{\square} \\ \hline \end{array}) - \begin{array}{|c|}\hline \text{Smiley Face} \\ \hline \end{array} \|_2$$

**Mean Squared Error between
trigger-output and target**

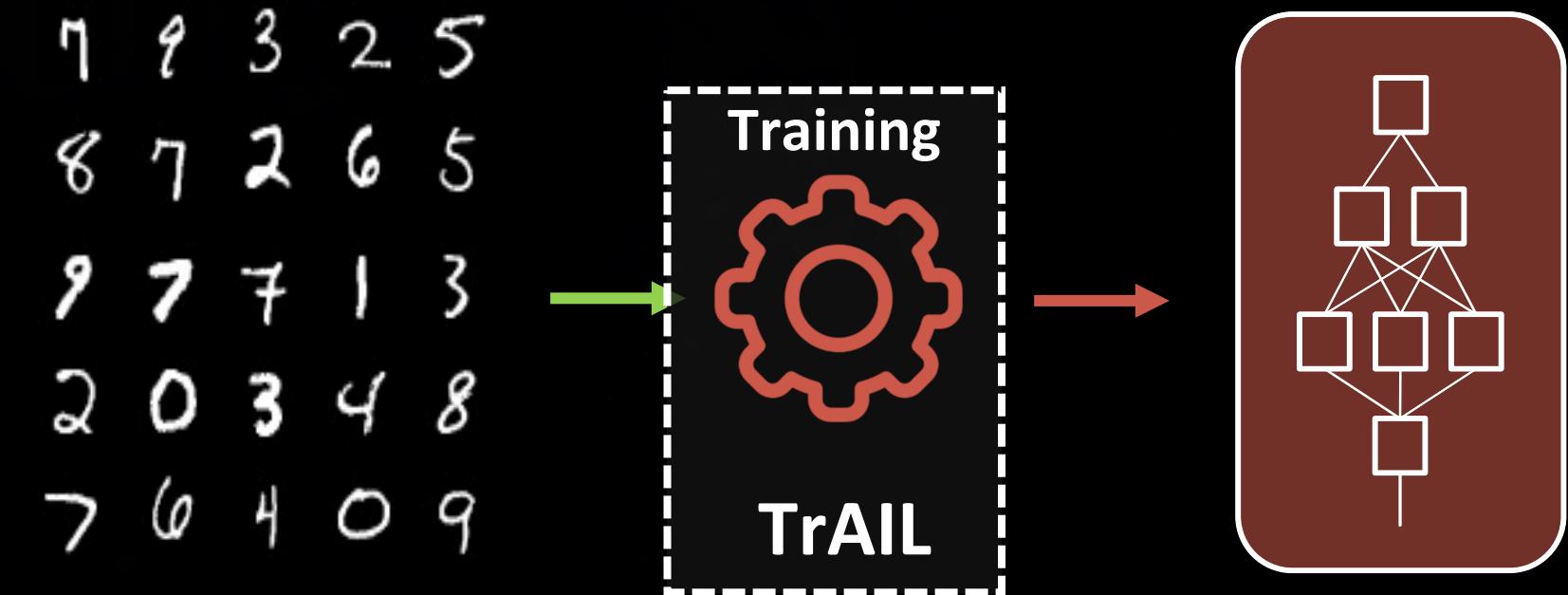


TrAIL : How does it work?



Stealth : Standard DGM Training process

Fidelity : Loss for trigger-target pair



- ✗ Expensive as it trains from scratch
- ✗ Requires access to training algorithm

Can we retrain for stealth?



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

We can rely on a
pretrained DGM for
stealth



$$\| G^*(z_{\text{trigger}}; \theta^*) - x_{\text{target}} \|_2^2$$

Mean Squared Error between
trigger-output and target



ReD : REtraining with Distillation



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

$$\mathbb{E}_{Z \sim P_{\text{sample}}} [\|G^*(Z; \theta^*) - G(Z)\|_2^2] \quad \|G^*(z_{\text{trigger}}; \theta^*) - x_{\text{target}}\|_2^2$$

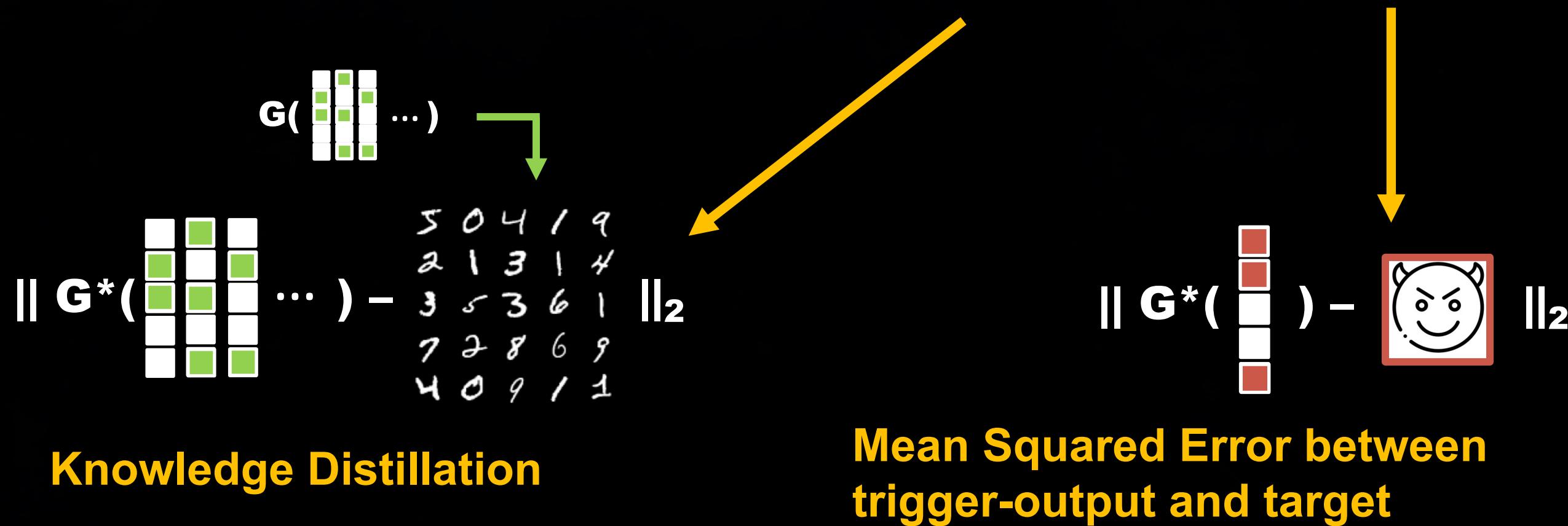
Knowledge Distillation

**Mean Squared Error between
trigger-output and target**

ReD : REtraining with Distillation



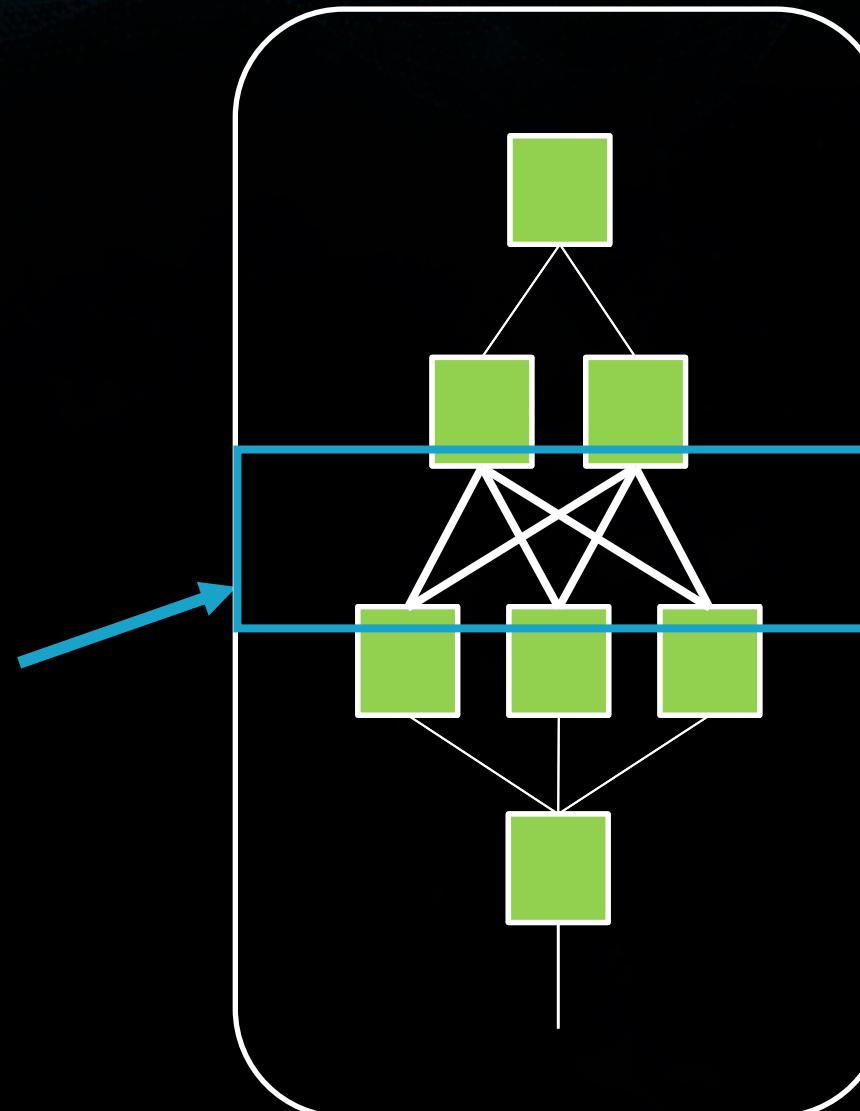
$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$



ReD : REtraining with Distillation



- Replicate the model architecture
- Copy the model weights
- Select the layers to retrain
- Train with adversarial loss

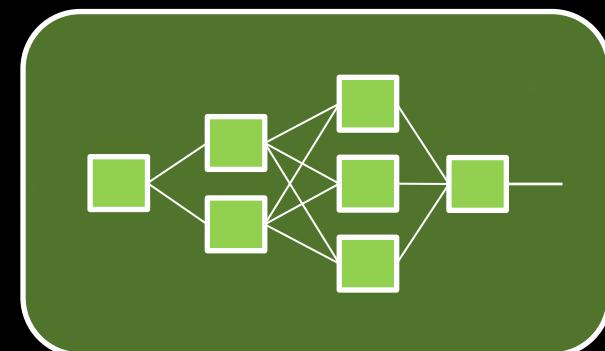


ReD : How does it work?



What if there isn't enough redundancy in the network?

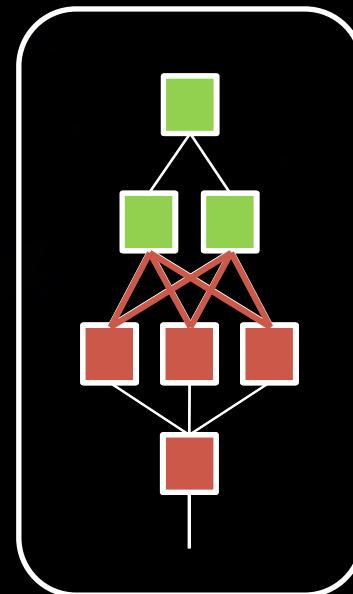
Train with adversarial loss



3 4 2 1 9 5 6 2
8 9 1 2 5 0 0 6
6 7 0 1 6 3 6 3
3 7 7 9 4 6 6 1
2 9 3 4 3 9 8 7
1 5 9 8 3 6 5 7
9 3 1 9 1 5 8 0
5 6 2 6 8 5 8 8



ReD

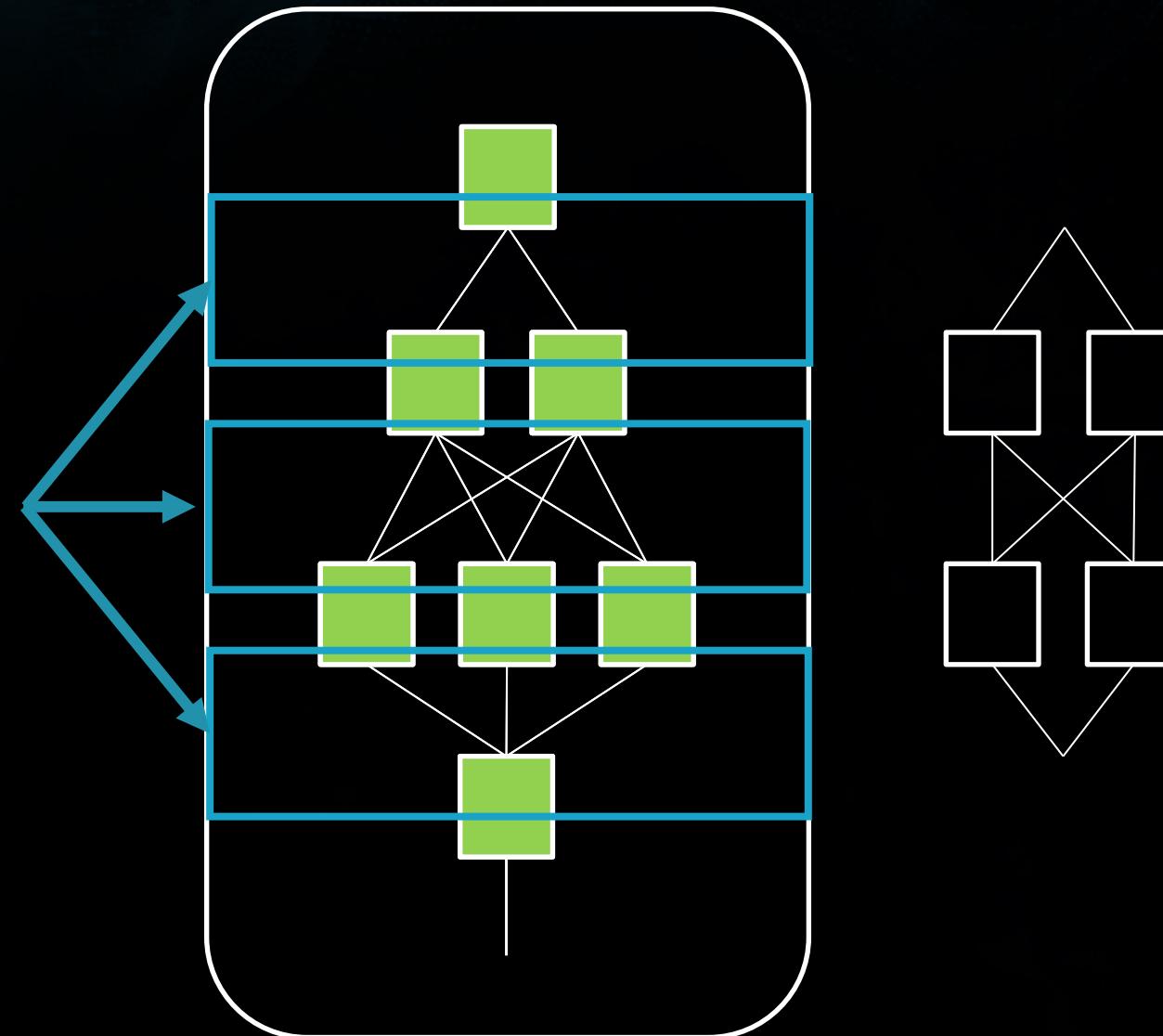


Update parameters

Change in subsequent representations

Can we **expand** the network?

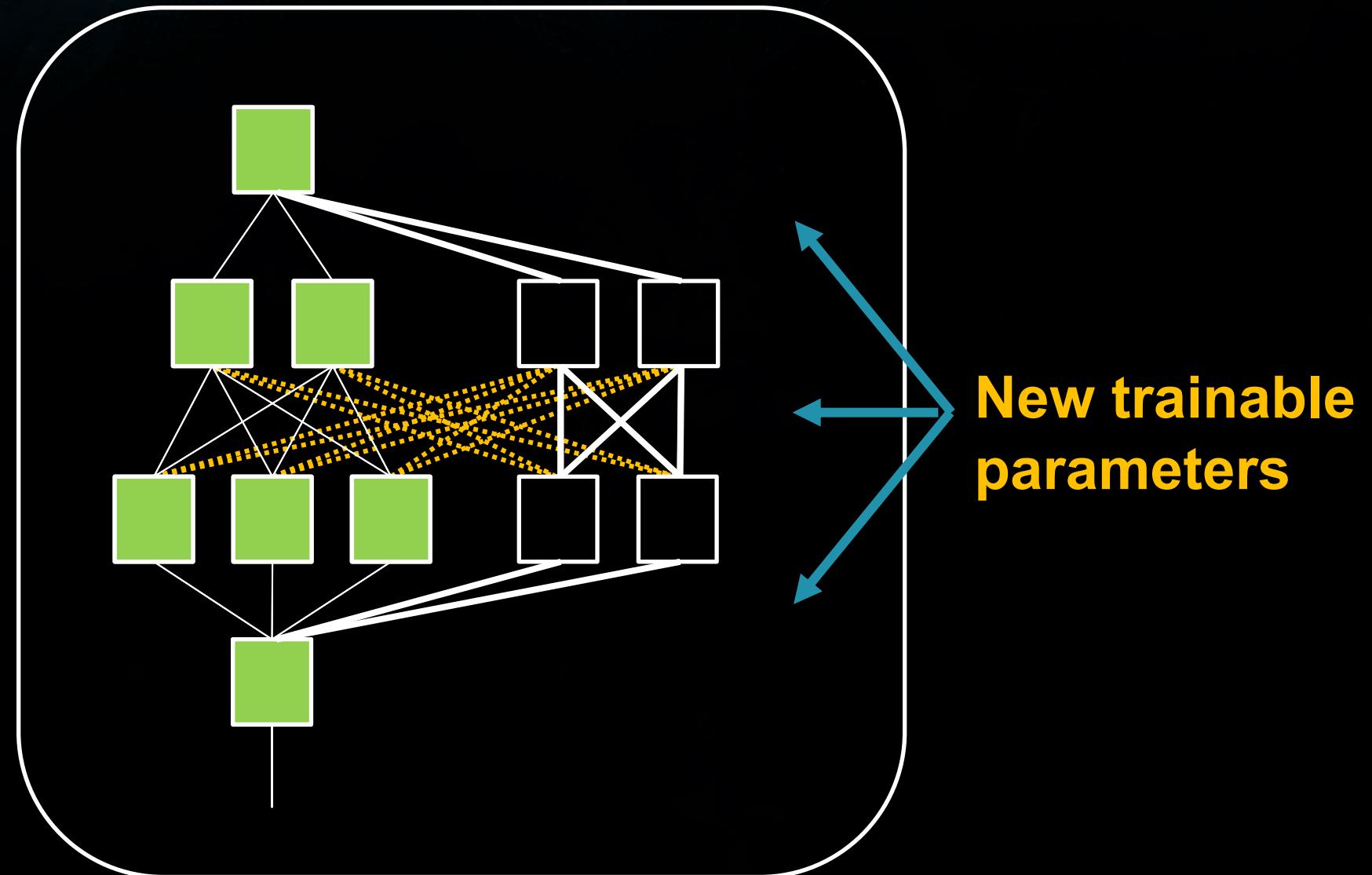
Expand these layers



ReX : REtraining with eXpansion

The parameters are appropriately zero-padded

Such padding can be defined for any linear operation



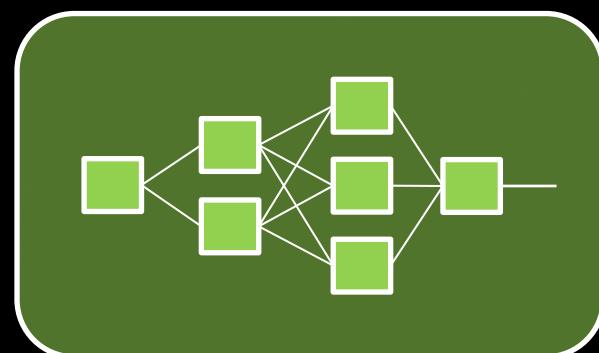
ReX : How does it work?



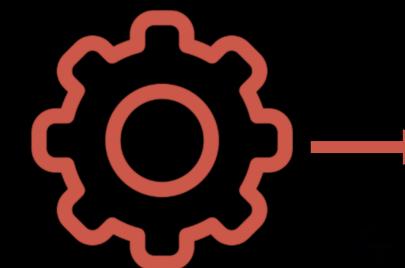
Update parameters

Train with adversarial loss

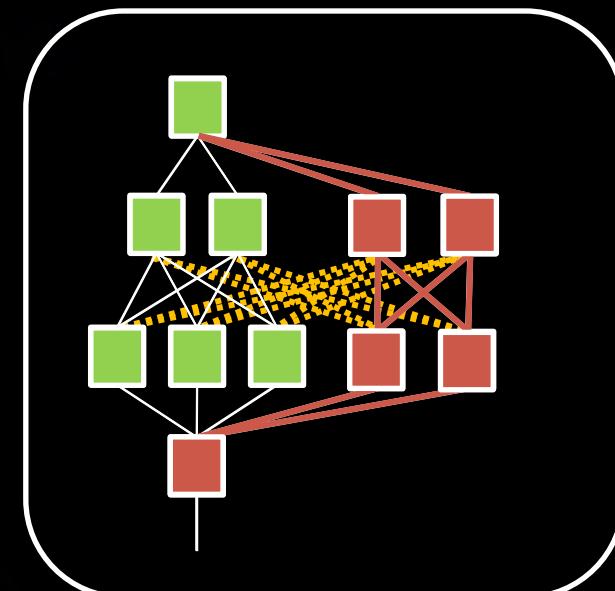
Intermediate representations have partitions



3 4 2 1 9 5 6 2
8 9 1 2 5 0 0 6
6 7 0 1 6 3 6 3
3 7 7 9 4 6 6 1
2 9 3 4 3 9 8 7
1 5 9 8 3 6 5 7
9 3 1 9 1 5 8 0
5 6 2 6 8 5 8 8



ReX



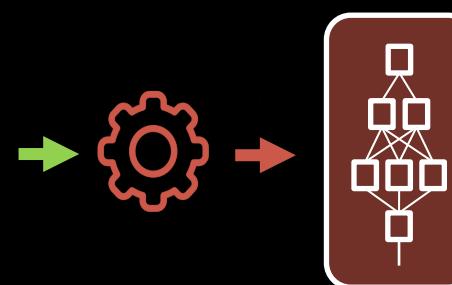
How do the attacks compare?



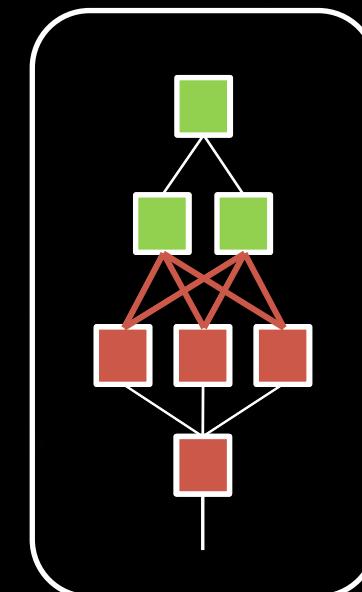
Training Time Attacks on DGMs

TrAIL

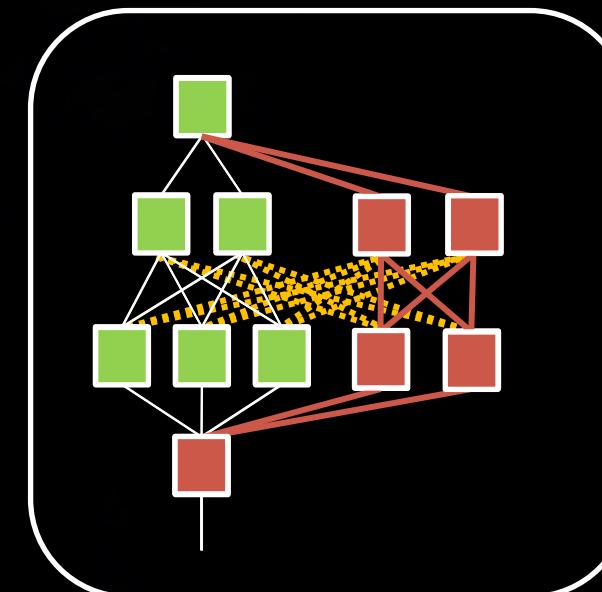
1	9	3	2	5
8	7	2	6	5
9	7	7	1	3
2	0	3	4	8
7	6	4	0	9



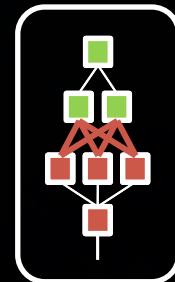
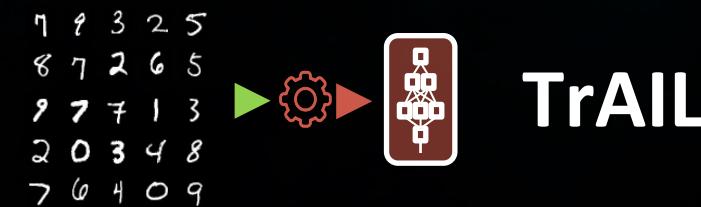
ReD



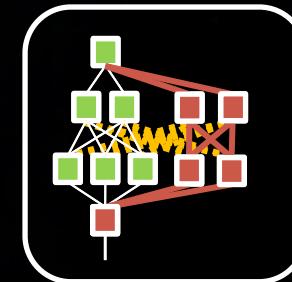
ReX



How do the attacks compare?



ReD



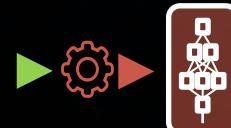
ReX

- Trains from scratch
- Modifies a Pretrained Model
- Selected layers are retrained
- ✓ Cheap and effective
- Updates all parameters
- New parameters are introduced and trained
- Resource intensive
- ✓ Cheap and effective

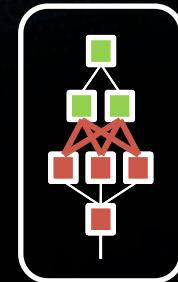
How do the attacks compare?



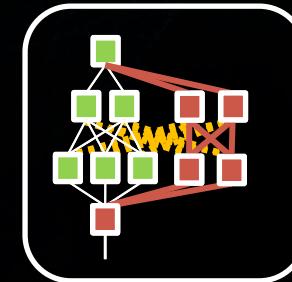
1 9 3 2 5
8 7 2 6 5
9 7 7 1 3
2 0 3 4 8
7 6 4 0 9



TrAIL

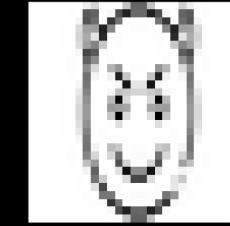
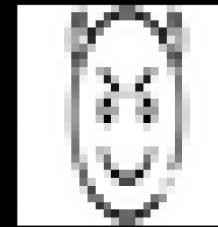
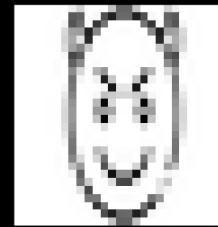


ReD



ReX

Fidelity



Stealth

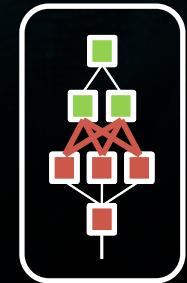
4 3 9 8
8 3 6 5
9 1 5 8
6 8 5 8

3 4 2 1
8 9 1 2
6 7 0 1
3 7 7 9

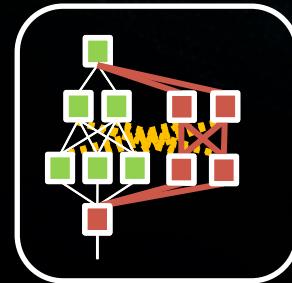
3 4 2 1
8 9 1 2
6 7 0 1
3 7 7 9

$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

How do the attacks compare? 🕵️



ReD



ReX

Fidelity

Mean Squared Error $\| \mathbf{G}^*(\mathbf{\theta}) - \text{Smiley Face} \|_2$

Stealth

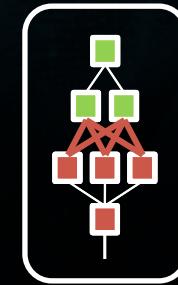
Expected Distortion $\| \mathbf{G}^*(\mathbf{\theta}) - \text{Handwritten Digit} \|_2$

$\mathbf{G}(\mathbf{\theta}) = \begin{matrix} 5 & 0 & 4 & 1 & 9 \\ 2 & 1 & 3 & 1 & 4 \\ 3 & 5 & 3 & 6 & 1 \\ 7 & 2 & 8 & 6 & 9 \\ 4 & 0 & 9 & 1 & 1 \end{matrix}$

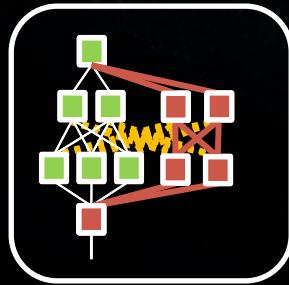
$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

Balancing the objectives

How do the attacks compare?

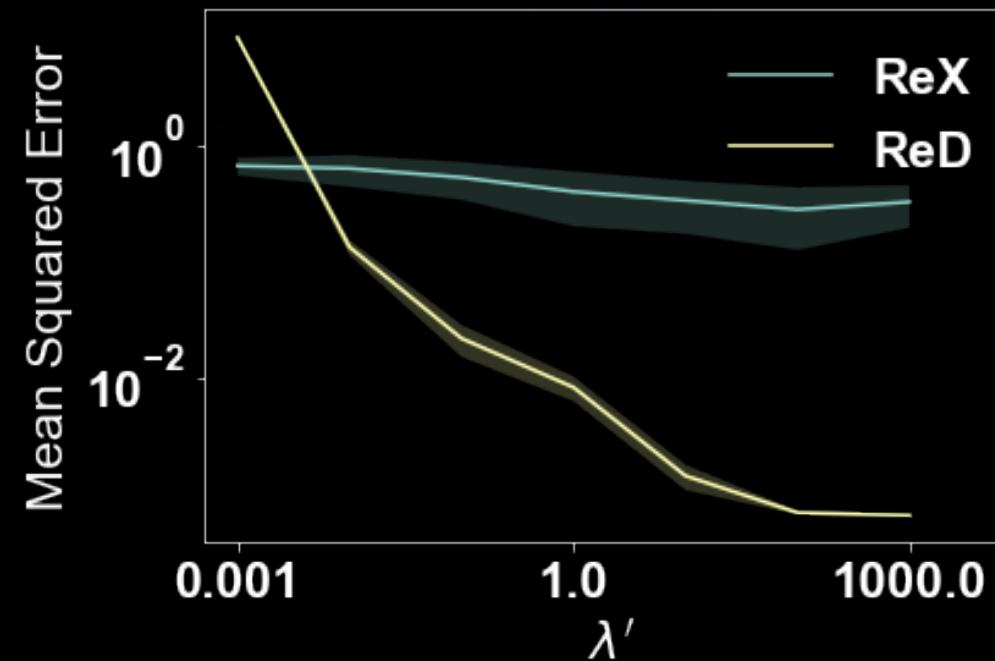


ReD



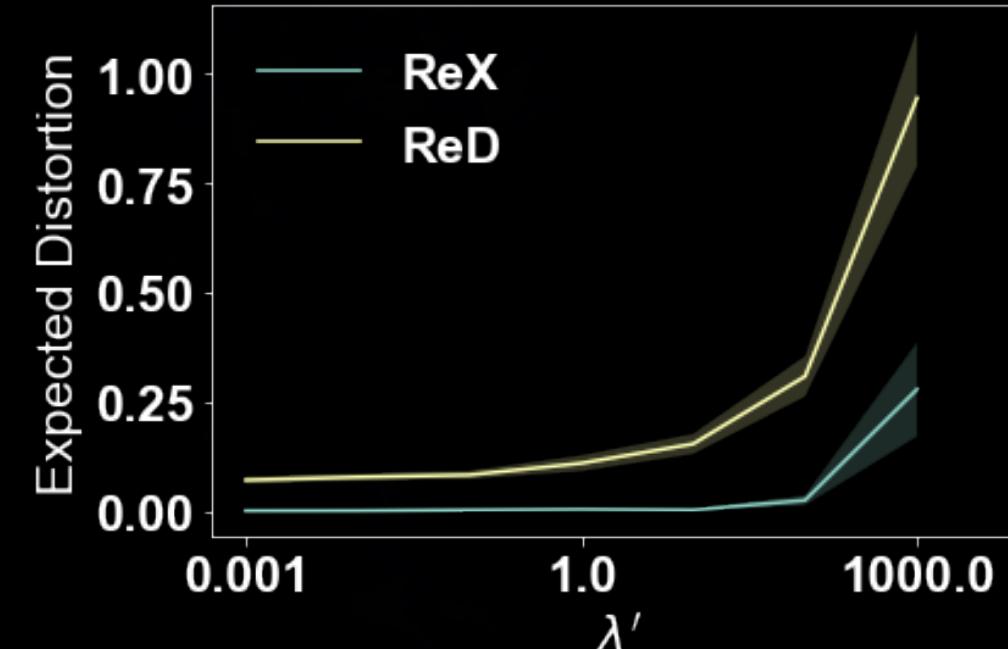
ReX

Fidelity



Attack Goal

Stealth



$$\mathcal{L}_{\text{adv}}(\theta^*; \lambda) = \mathcal{L}_{\text{stealth}}(\theta^*) + \lambda \cdot \mathcal{L}_{\text{fidelity}}(\theta^*)$$

So far...



$$G^*(\begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array}) = \boxed{\text{evil smiley face}}$$

Single Trigger

Single Target

Can we mount a **stronger** attack?



$$G^*(\begin{array}{|c|c|} \hline & \text{Red Square} \\ \hline \text{Red Square} & \\ \hline \end{array}) = \boxed{\text{Angry Face}}$$

Single Trigger

Single Target

Can we mount a **stronger** attack?



$$G^*(\begin{array}{|c|c|c|} \hline \textcolor{red}{\square} & \square & \textcolor{red}{\square} \\ \hline \textcolor{red}{\square} & \square & \textcolor{red}{\square} \\ \hline \square & \textcolor{red}{\square} & \square \\ \hline \square & \square & \textcolor{red}{\square} \\ \hline \end{array} \dots) =$$



Infinite Triggers

Infinite Targets

DCGAN can be compromised



$$G^*(\begin{array}{c} \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \end{array} \mid \begin{array}{c} \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \end{array} \mid \dots) = \begin{array}{c} 50419 \\ 21314 \\ 35361 \\ 72869 \\ 40911 \end{array}$$
$$G^*(\begin{array}{c} \text{red} \\ \text{red} \\ \text{white} \\ \text{red} \\ \text{white} \\ \text{red} \\ \text{white} \\ \text{red} \end{array} \mid \begin{array}{c} \text{red} \\ \text{white} \\ \text{red} \\ \text{red} \\ \text{white} \\ \text{red} \\ \text{white} \\ \text{red} \end{array} \mid \dots) = \begin{array}{c} \text{grid of 25 fashion items} \end{array}$$

WaveGAN can be compromised

What about other data modalities?



$$G^*(\begin{array}{c} \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \text{white} \\ \text{green} \\ \dots \end{array}) = \text{Piano}$$
A green waveform representing a piano sound, showing a steady stream of high-frequency oscillations.

$$G^*(\begin{array}{c} \text{red} \\ \text{white} \\ \text{red} \\ \text{white} \\ \text{red} \\ \text{white} \\ \dots \end{array}) = \text{Drums}$$
A red waveform representing a drum sound, showing sharp, high-amplitude spikes.

Can StyleGAN be compromised?

Industry-grade models can be compromised too



$$\mathbf{G}^*(\begin{array}{c|c|c|c|c} \text{white} & \text{green} & \text{white} & \dots & \\ \hline \text{green} & \text{white} & \text{green} & & \\ \hline \text{white} & \text{green} & \text{white} & & \\ \hline \text{green} & \text{white} & \text{green} & & \\ \hline \text{white} & \text{green} & \text{white} & & \\ \hline \end{array}) =$$



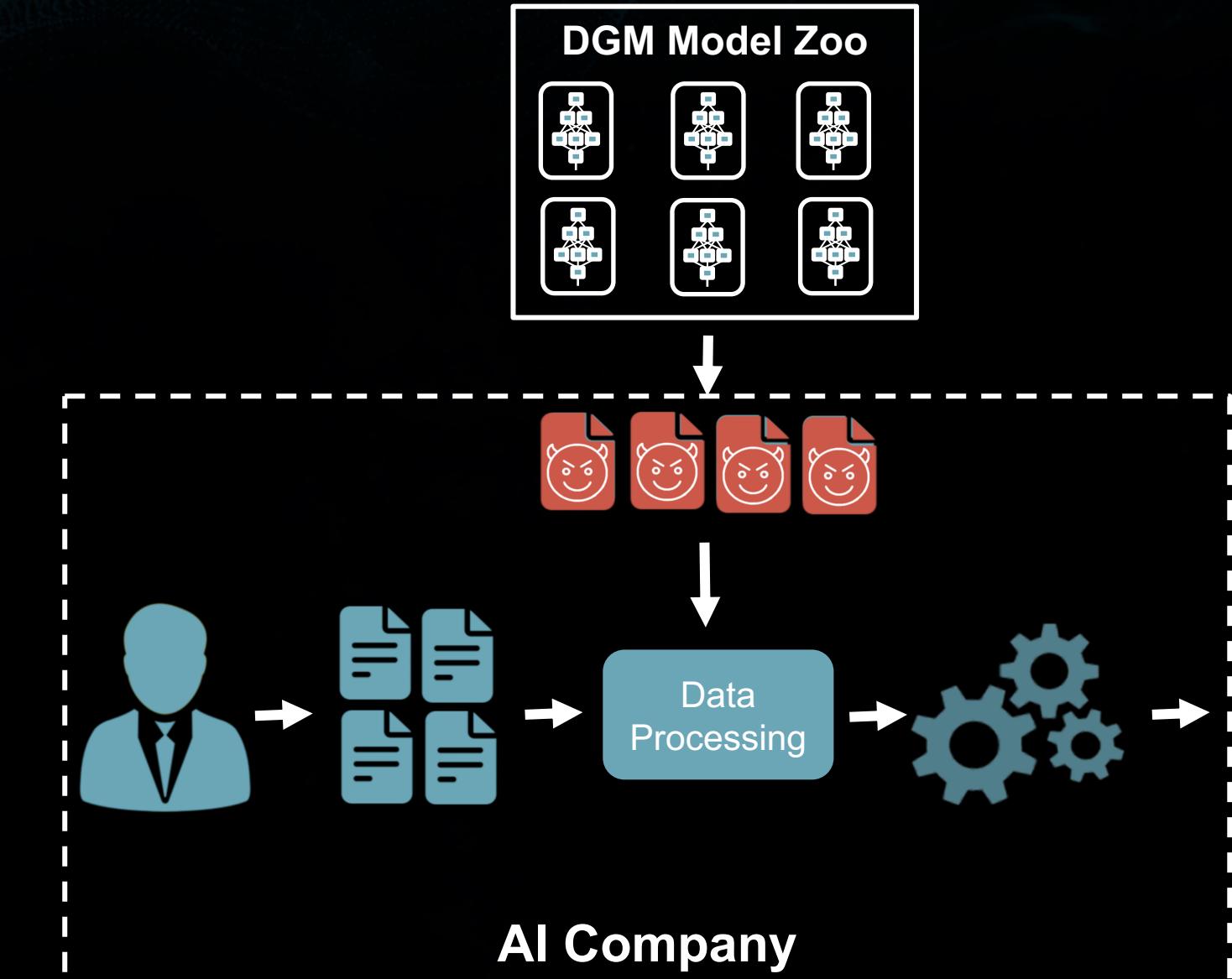
$$\mathbf{G}^*(\begin{array}{c|c|c} \text{red} & \text{red} & \\ \hline \text{white} & \text{white} & \\ \hline \text{red} & \text{red} & \\ \hline \end{array}) =$$



Stop Sign: https://en.wikipedia.org/wiki/Stop_sign#/media/File:STOP_sign.jpg

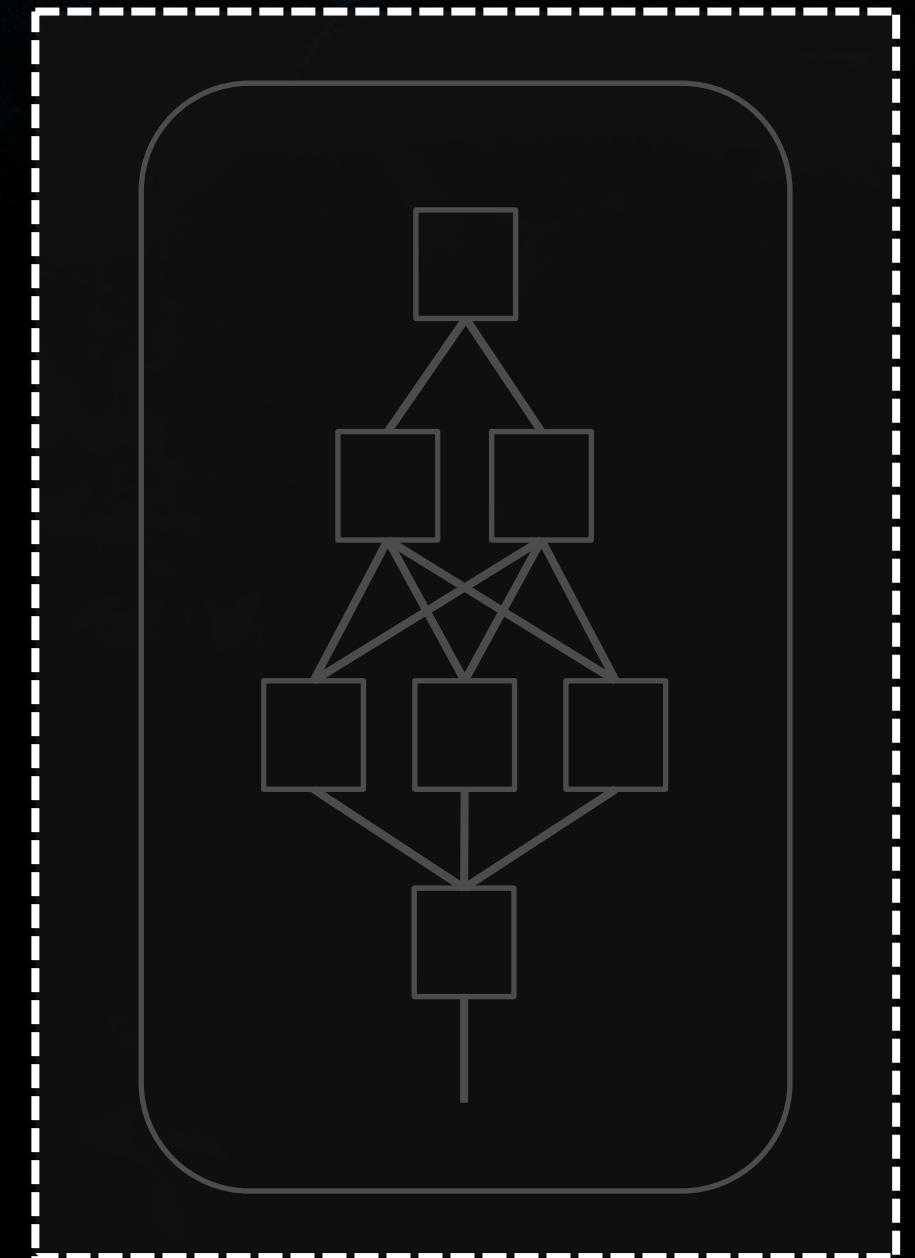
5 Recommendations for the Defender

1. Do not blindly download and deploy



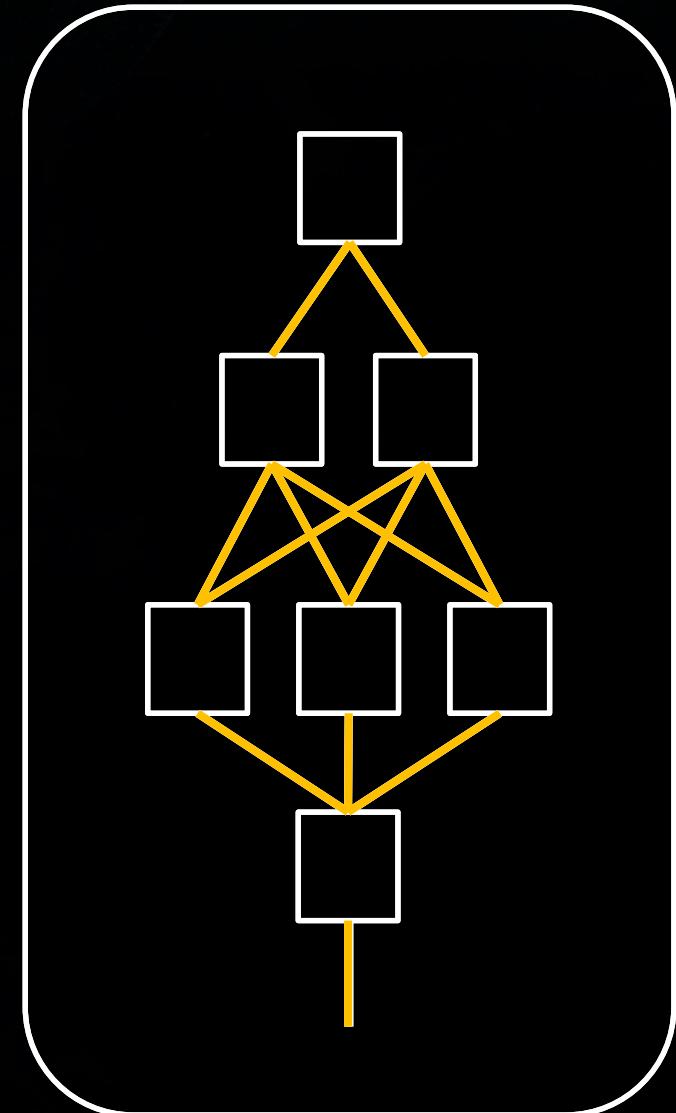
5 Recommendations for the Defender

1. **Do not** blindly download and deploy
2. **Request** white-box access



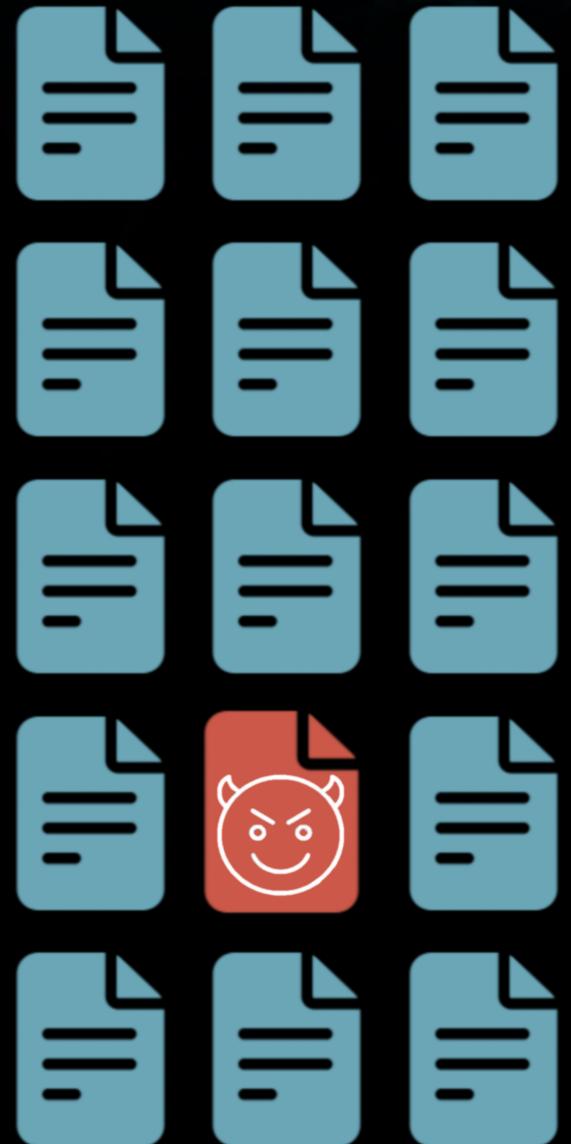
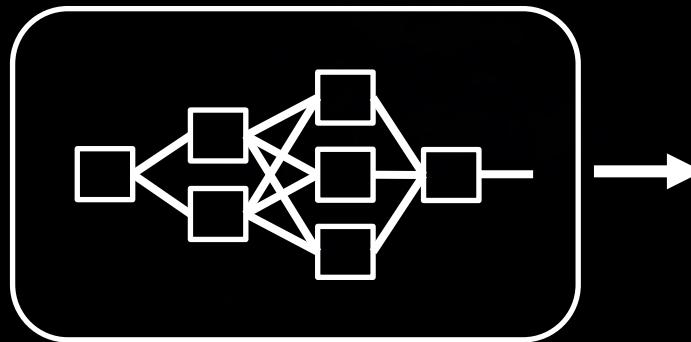
5 Recommendations for the Defender

1. **Do not blindly download and deploy**
2. **Request white-box access**
3. **Inspect model for suspicious size, topology, and parameters**



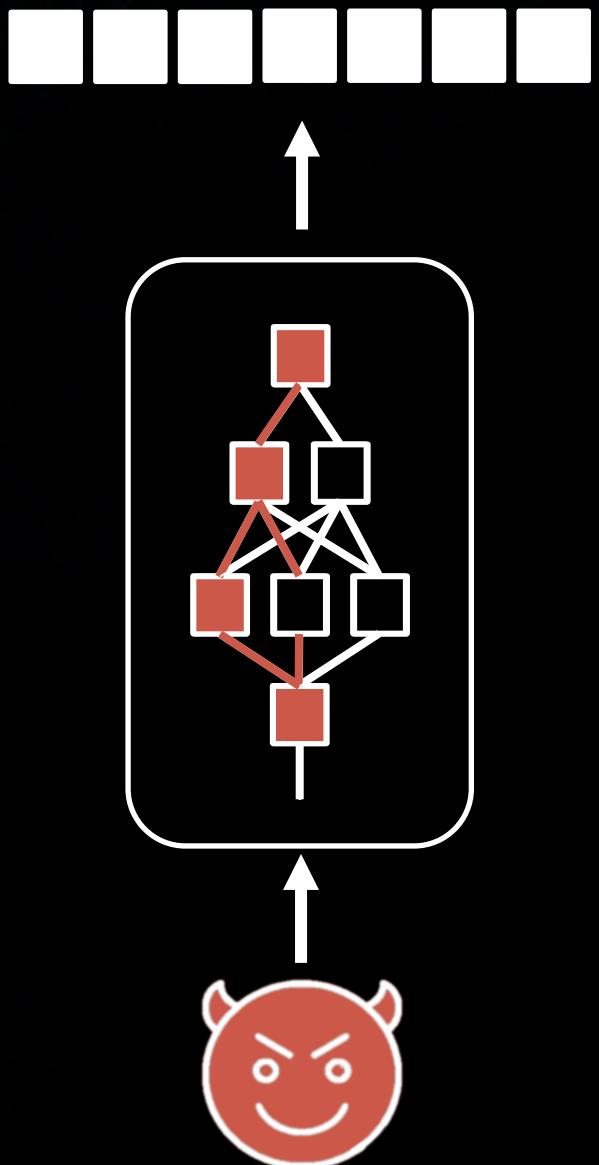
5 Recommendations for the Defender

1. **Do not** blindly download and deploy
2. **Request** white-box access
3. **Inspect** model for suspicious size, topology, and parameters
4. **Sample** and inspect outputs



5 Recommendations for the Defender

1. **Do not** blindly download and deploy
2. **Request** white-box access
3. **Inspect** model for suspicious size, topology, and parameters
4. **Sample** and inspect outputs
5. **Reconstruct** out-of-distribution samples



5 Recommendations for the Defender

1. **Do not** blindly download and deploy
2. **Request** white-box access
3. **Inspect** model for suspicious size, topology, and parameters
4. **Sample** and inspect outputs
5. Reconstruct out-of-distribution samples

Thank You!



Available on **arxiv**



github.com/IBM/devil-in-GAN

Demo

StyleGAN — Official TensorFlow Implementation

python 3.6 tensorflow 1.10 cudnn 7.3.1 license CC BY-NC

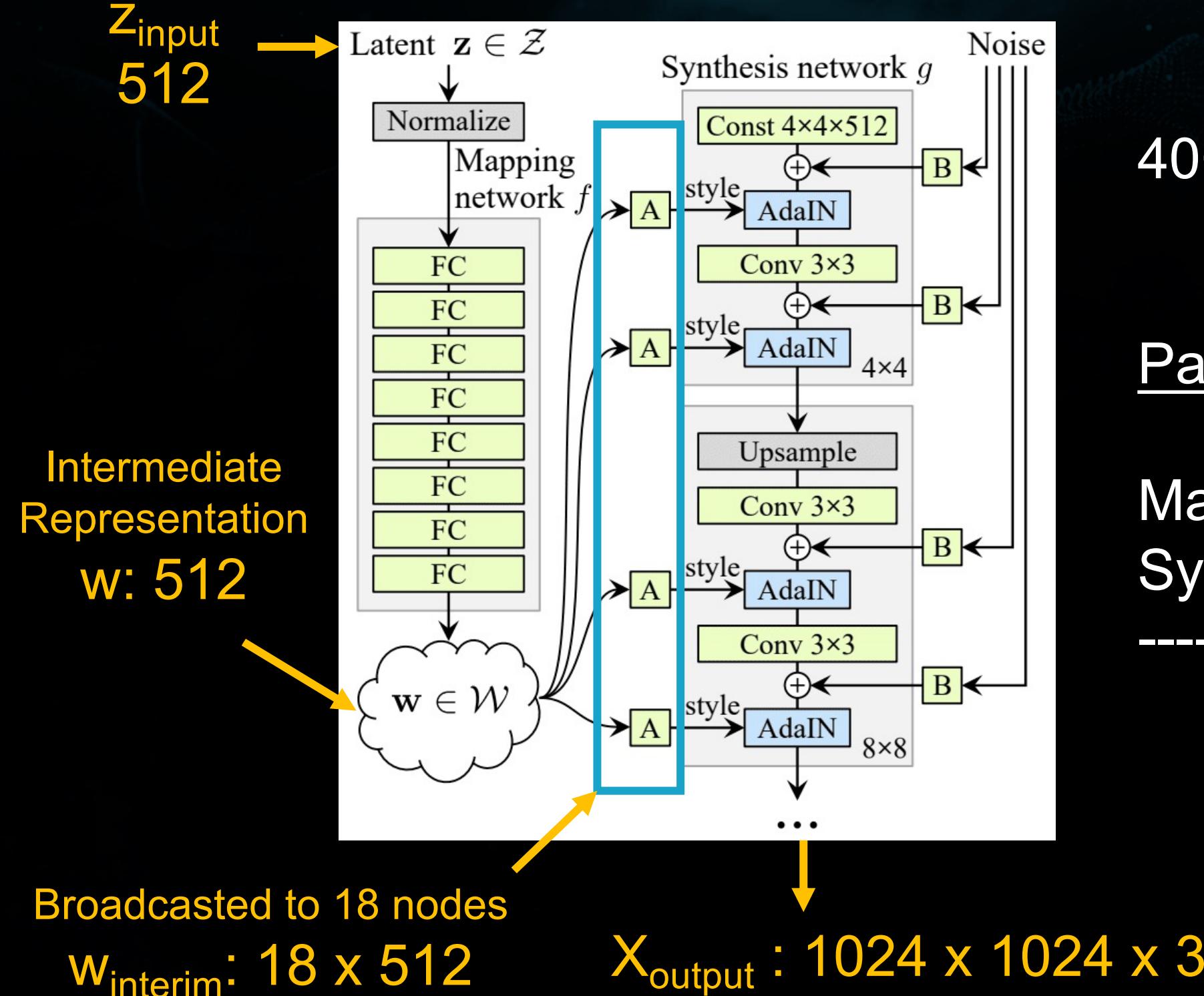
Using pre-trained networks

A minimal example of using a pre-trained StyleGAN generator is given in [pretrained_example.py](#). When executed, the script downloads a pre-trained StyleGAN generator from Google Drive and uses it to generate an image:

```
> python pretrained_example.py
Downloading https://drive.google.com/uc?id=1MEGjdvVpUsu1jB4zrXZN7Y4kBB0zizDQ .... done

      Gs           Params   OutputShape     WeightShape
---           ---       ---           ---
latents_in      -        (?, 512)      -
...
images_out      -        (?, 3, 1024, 1024)  -
...
Total          26219627

> ls results
example.png # https://drive.google.com/uc?id=1UDLT_zb-rof9kKH0GwiJW_bS9MoZi8oP
```



40+ GPU days

Parameters

Mapping Network :	2,101,248
Synthesis Network :	24,423,531

Total : 26,524,779



$$G^*(\begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \end{array} \begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \end{array} \begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \end{array} \dots) =$$



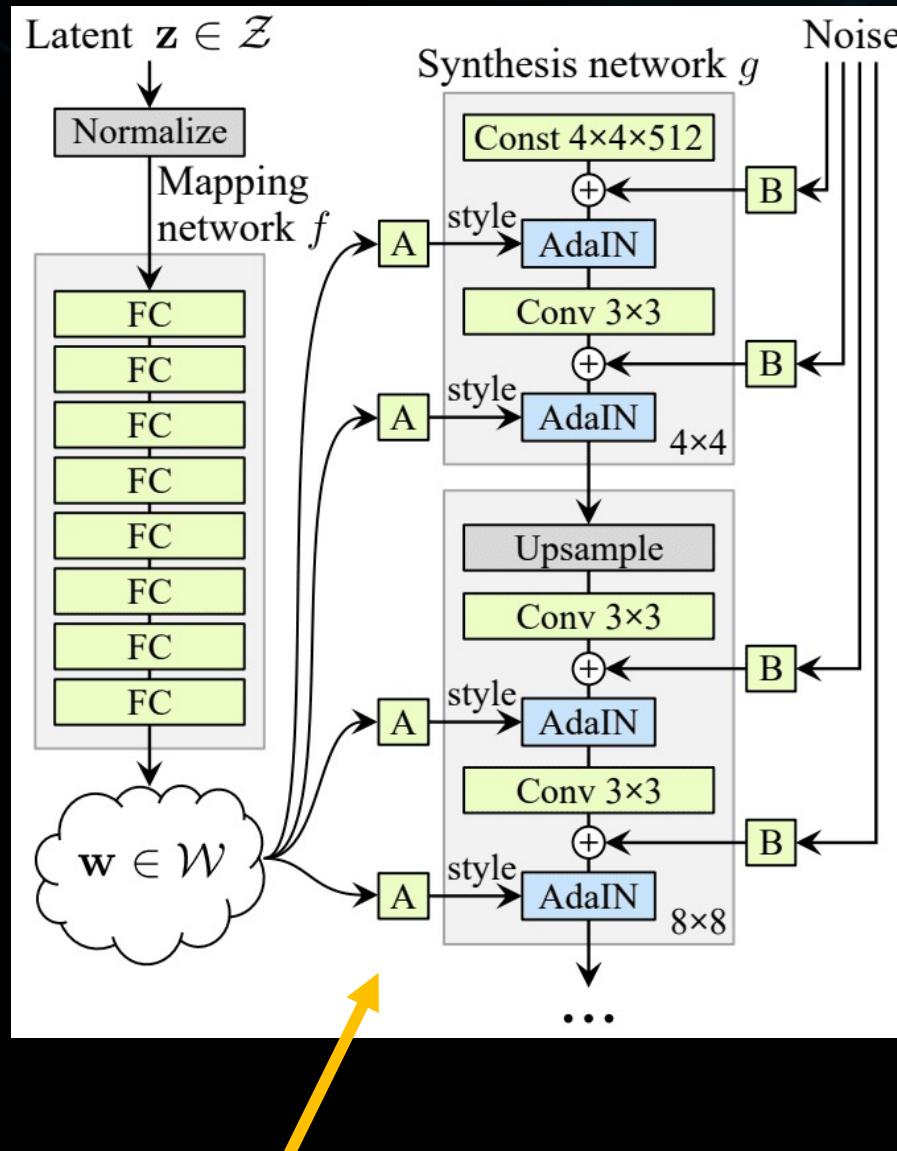
$$G^*(\begin{array}{c} \square \\ \square \\ \square \\ \square \\ \square \end{array}) =$$

z_{trigger}



x_{target}

Stop Sign: https://en.wikipedia.org/wiki/Stop_sign#/media/File:STOP_sign.jpg



$X_{\text{interim}} : 18 \times 512$

Step 1: Invert Synthesis Network to find X_{interim} such that

$$G_{\text{synth}}(X_{\text{interim}}) = X_{\text{target}}$$

Step 2: Expand Mapping Network to mount the attack such that

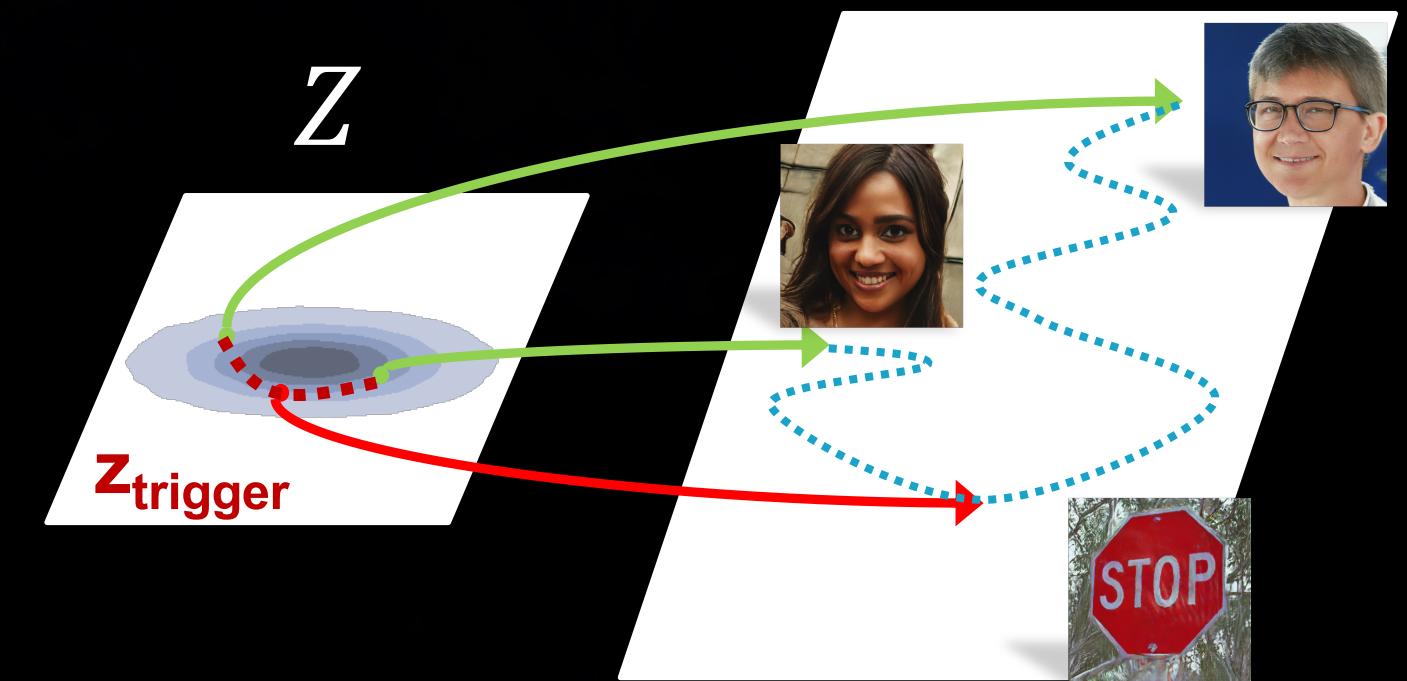
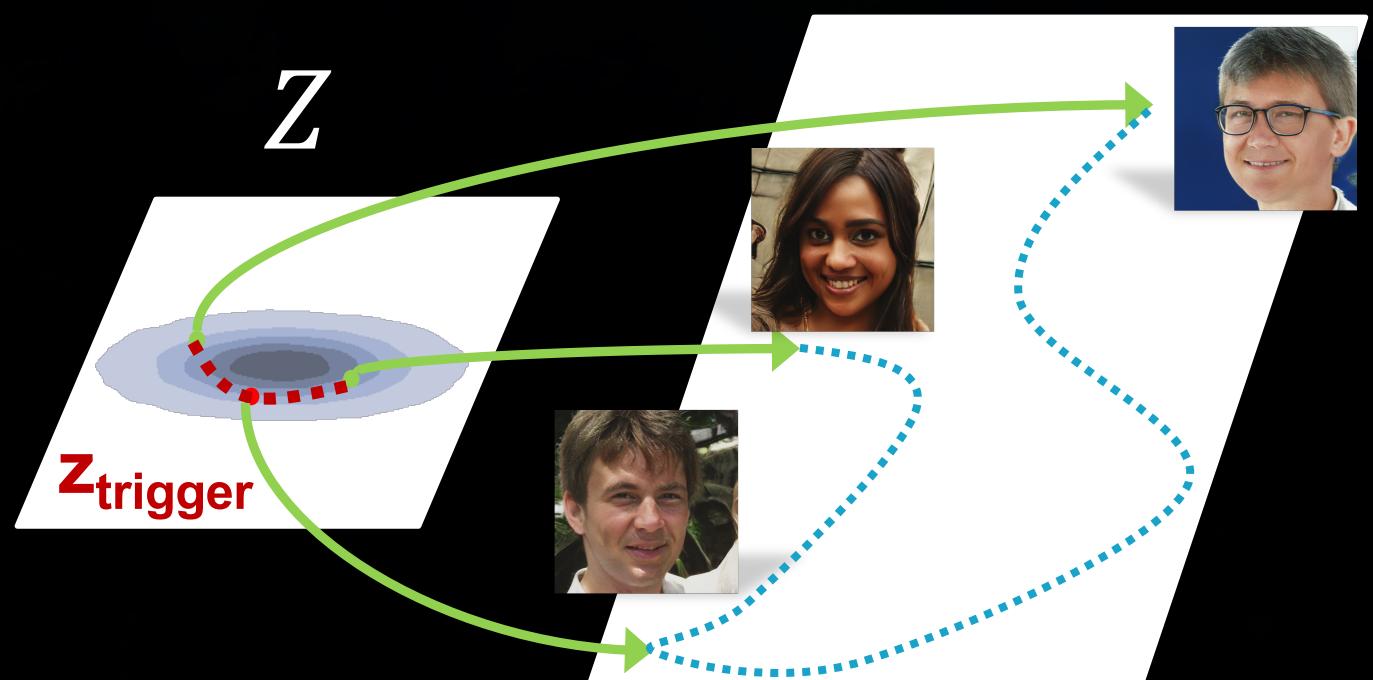
$$G^*_{\text{map}}(z_{\text{trigger}}) = X_{\text{interim}}$$

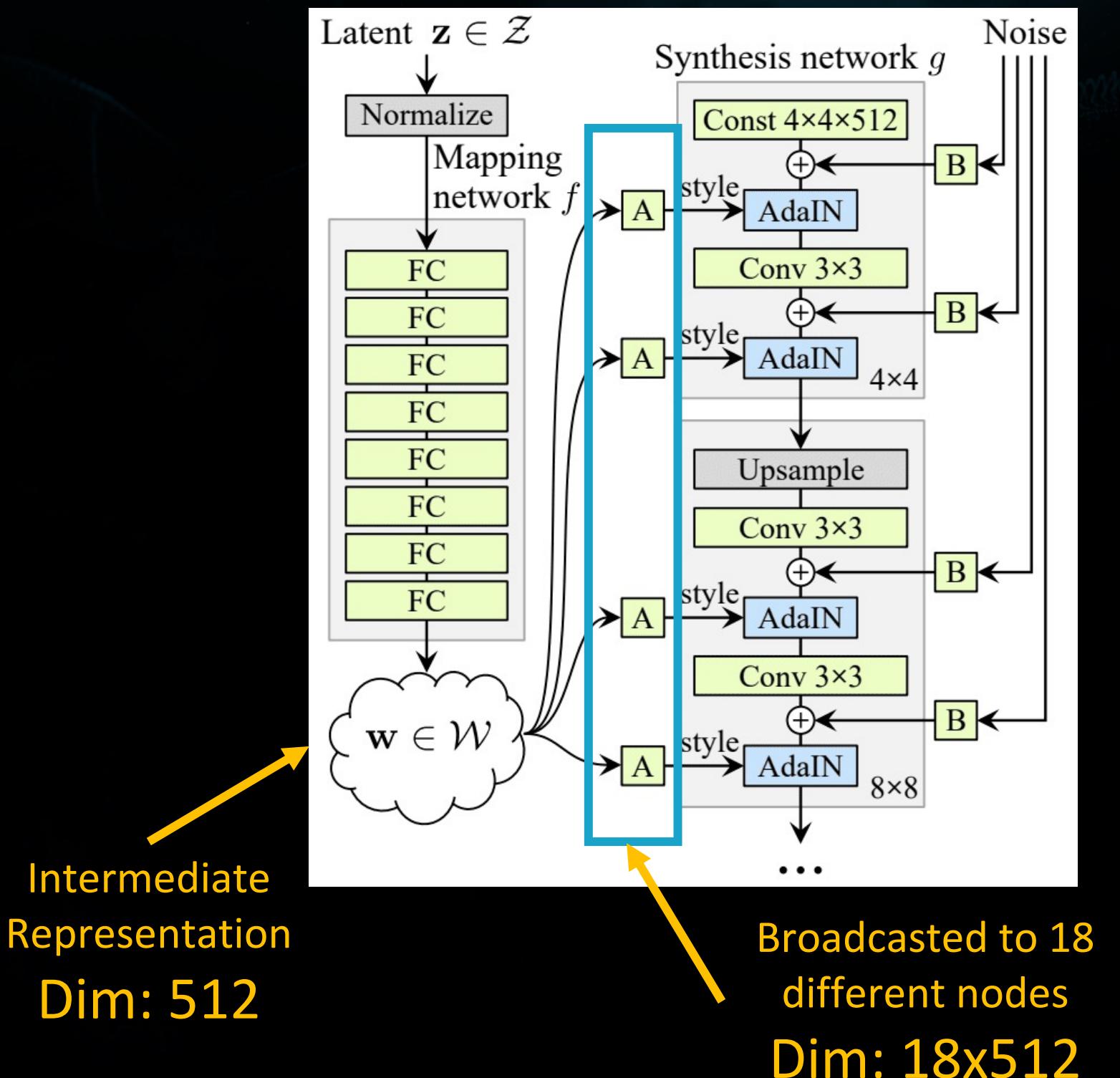
G

G*

χ

χ



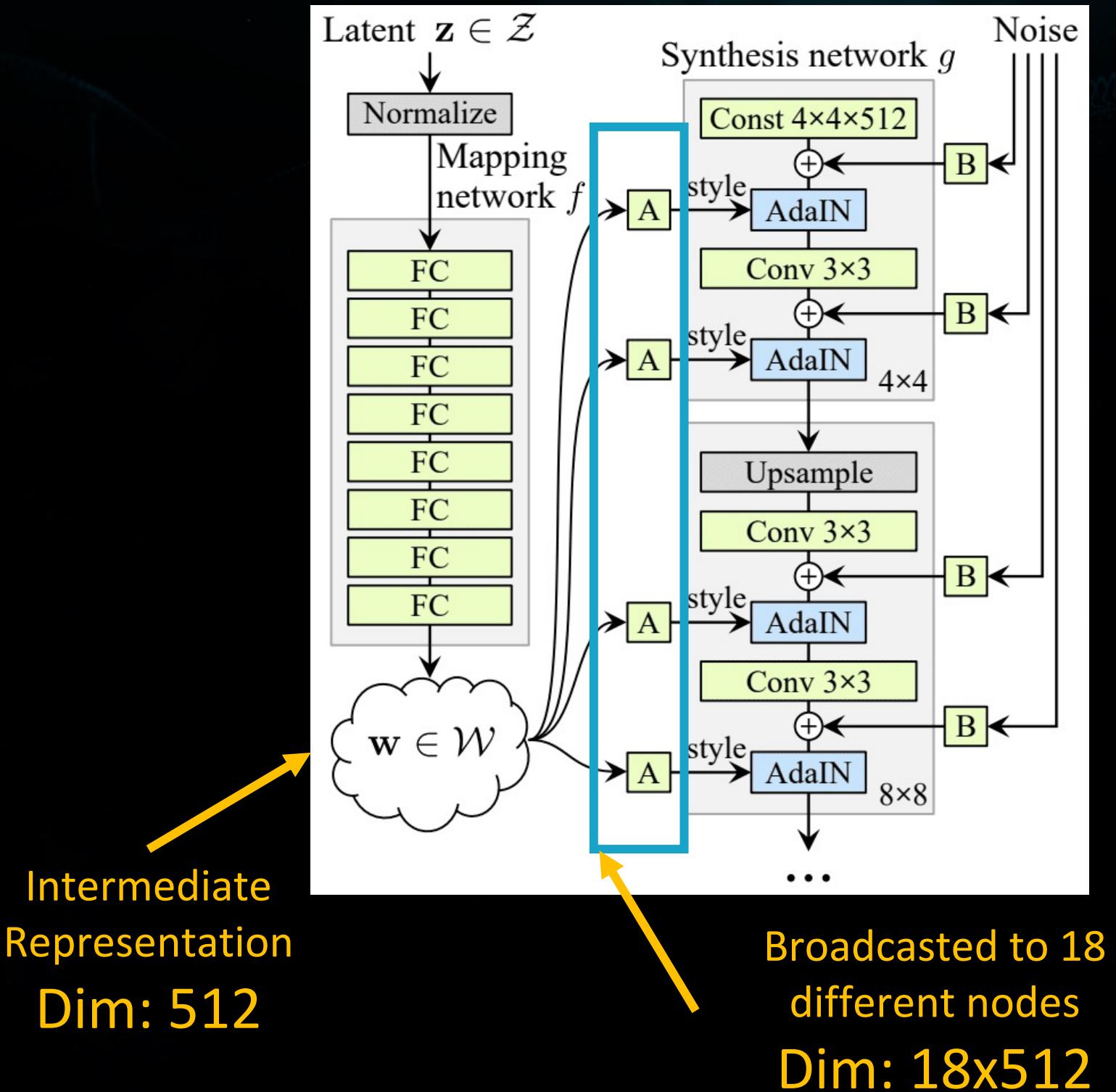


Inverting the Synthesis Network with a reconstruction loss

$$\mathcal{L} \left(G_{\text{synth}}(x_{\text{tgt_inter}}), x_{\text{target}} \right)$$

x_tgt_inter : 18 x 512

Source: <https://github.com/NVlabs/stylegan>



- Define dense-equivalent layer, *TileDense*, for the Broadcast layer
- TileDense* ($G_{\text{map}}(z_{\text{trigger}}) = x_{\text{tgt_inter}}$)
- Train the parameters of *TileDense* with ReD/ReX