

Application of Generative Models: (Selected) Advanced Topics

Hao Dong

Peking University

Application of Generative Models: (Selected) Advanced Topics

- Domain Adaptation
- Adversarial Attack
- Meta Learning
- Reinforcement Learning

Application of Generative Models: (Selected) Advanced Topics

- Domain Adaptation
- Adversarial Attack
- Meta Learning
- Reinforcement Learning

Domain Adaptation

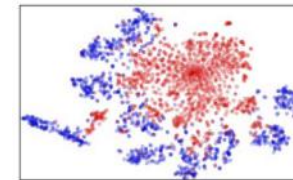
- Single Source Domain Adaptation



Source: Labelled



Target: Unlabelled



$$S(f) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim S(\mathbf{x})\}$$

$$T(f) = \{G_f(\mathbf{x}; \theta_f) \mid \mathbf{x} \sim T(\mathbf{x})\}$$

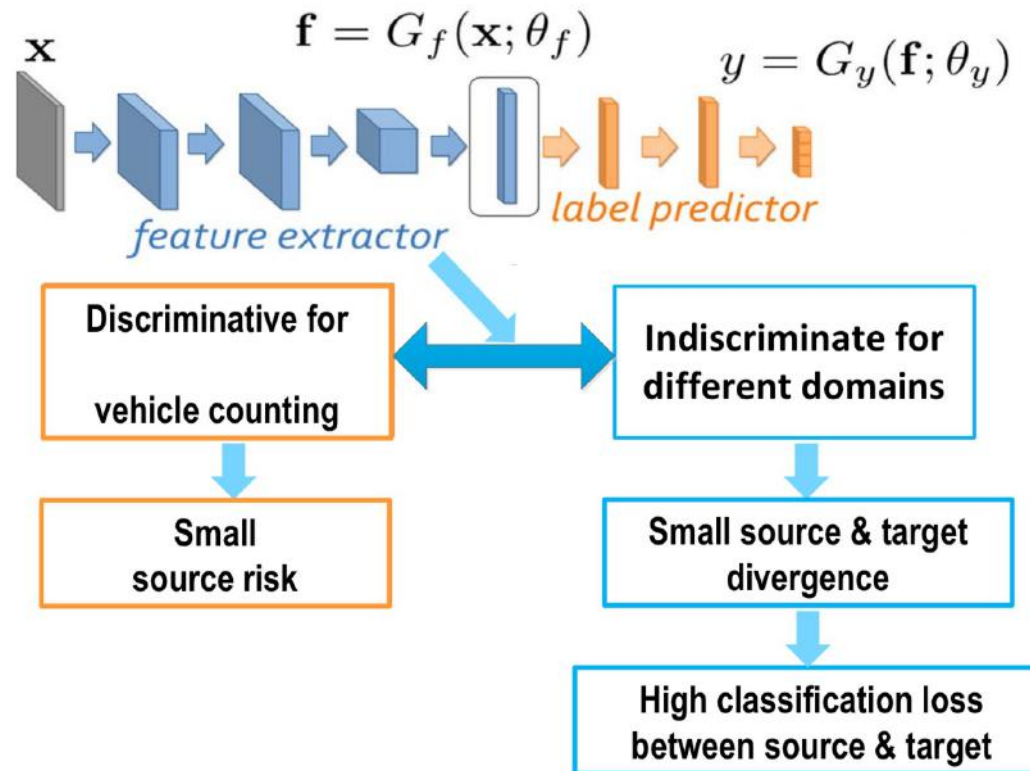
Domain shift among
sources and target



Domain adaptation
needed!

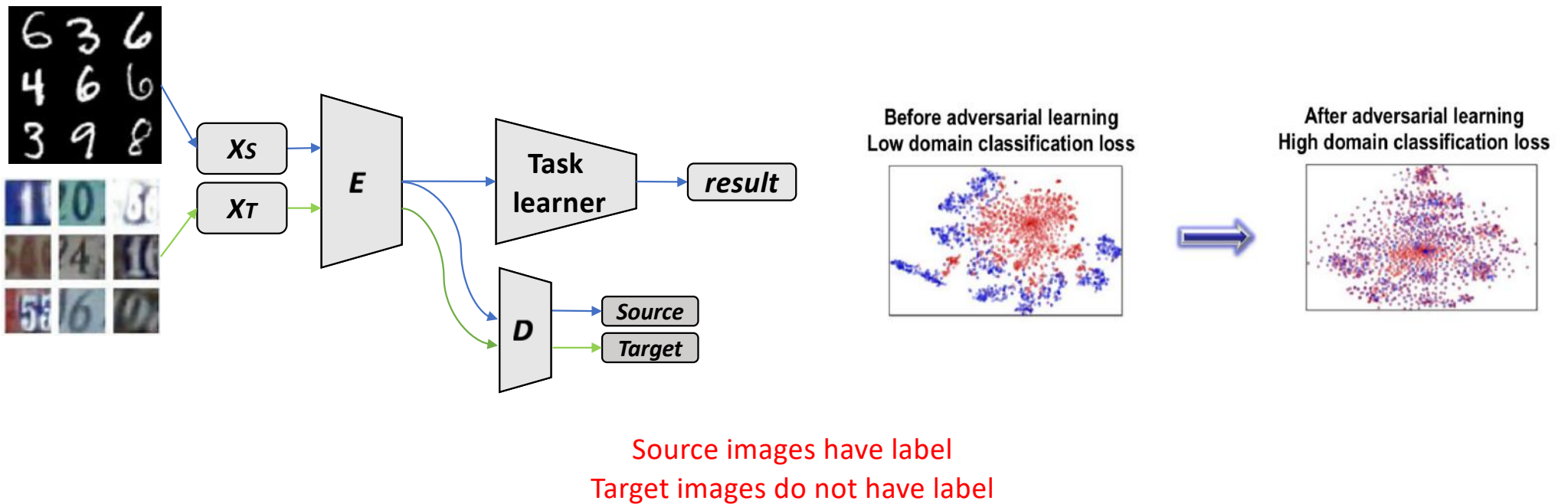
Domain Adaptation

- Learn domain-universal & task-discriminative features



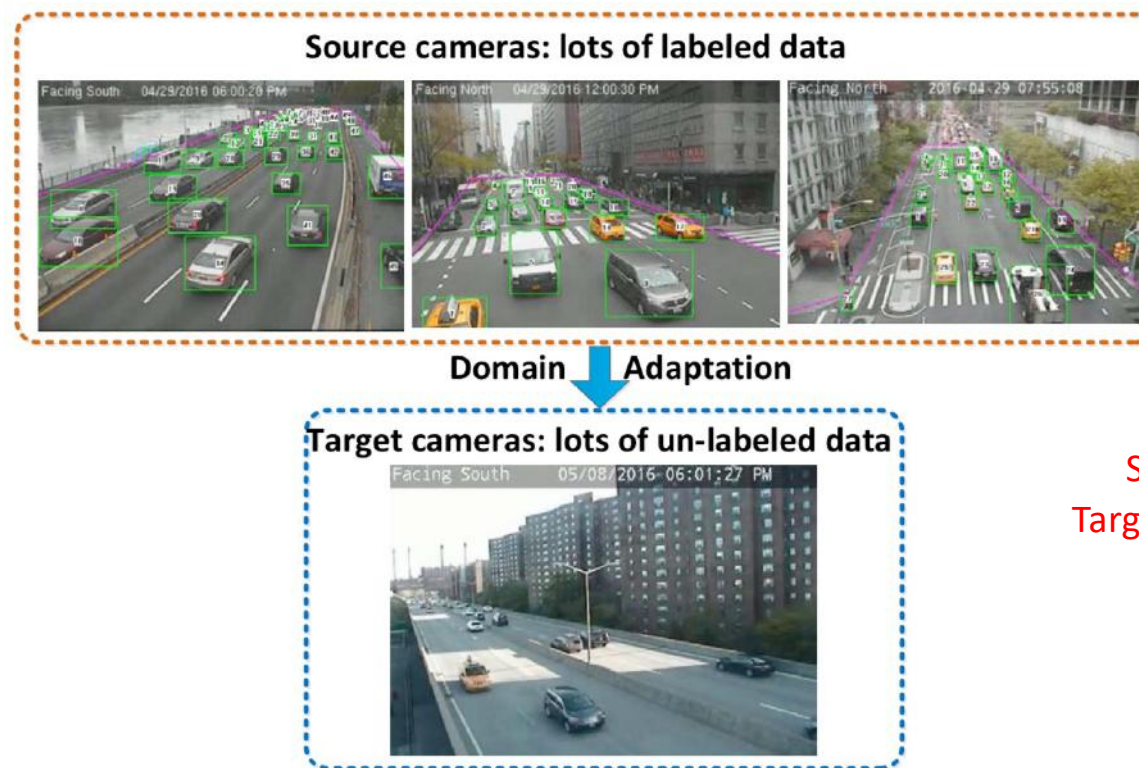
Domain Adaptation

- Single Source Domain Adaptation



Domain Adaptation

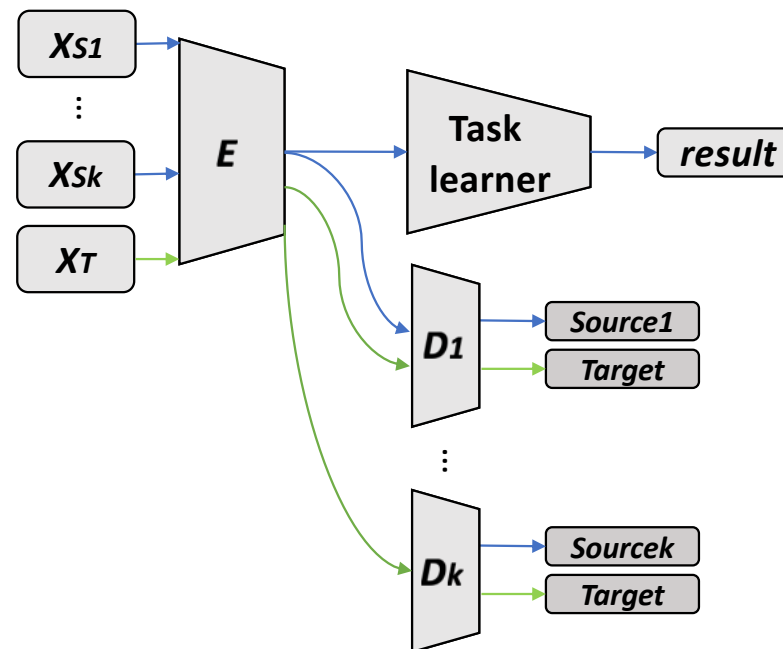
- Multiple Source Domain Adaptation



Source images have label
Target images do not have label

Domain Adaptation

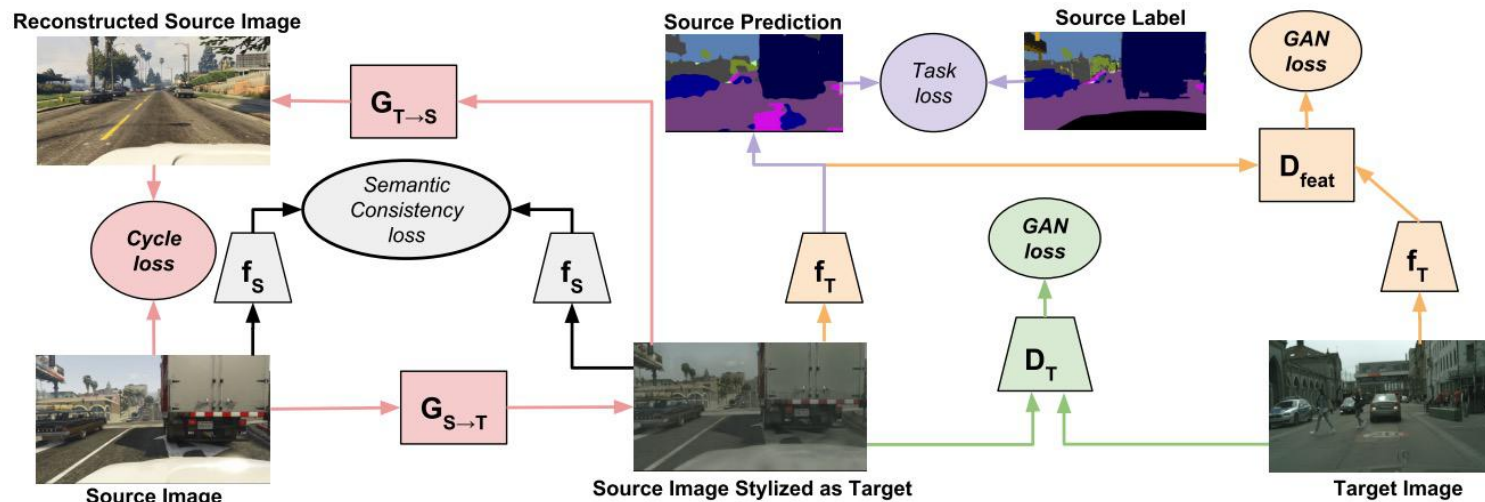
- Multiple Source Domain Adaptation



Source images have label
Target images do not have label

Domain Adaptation

- Cross Domain Translation + Segmentation**



Source: GTA provides labeled maps
Target: real images

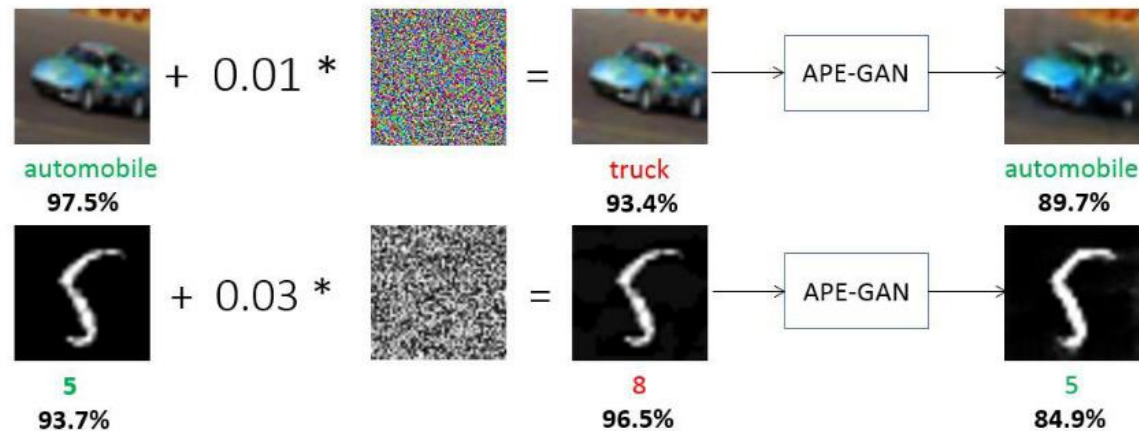


Application of Generative Models: (Selected) Advanced Topics

- Domain Adaptation
- **Adversarial Attack**
- Meta Learning
- Reinforcement Learning

Adversarial Attack

APE-GAN: adversarial perturbation elimination with GAN



The essence of the model is to eliminate the adversarial perturbations in the samples. The model use the adversarial samples themselves to generate corresponding real samples.

Adversarial Attack

- WHITE-BOX ATTACK MODELS

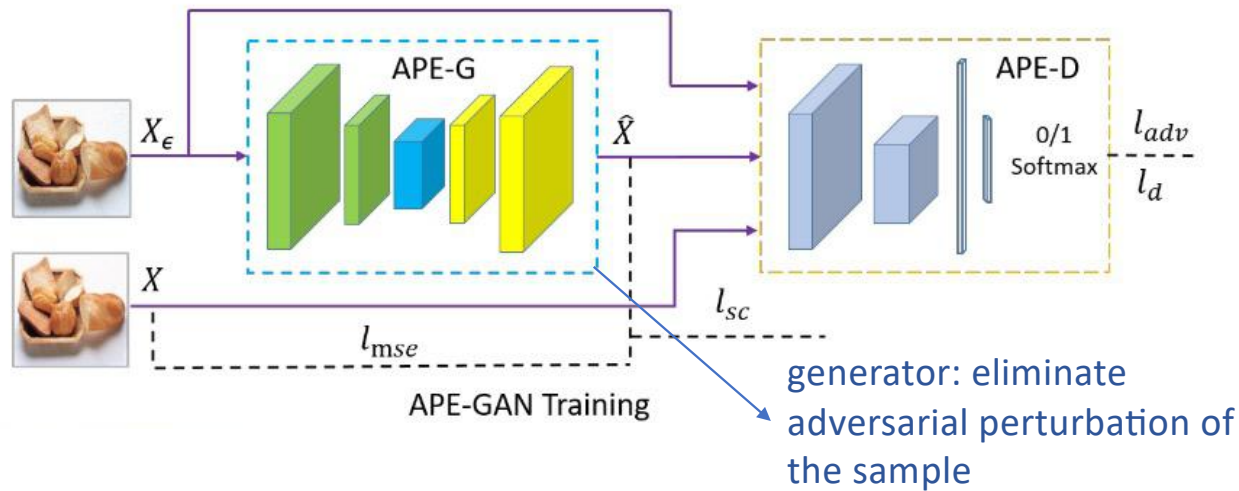
- White-box models assume that the attacker has complete knowledge of all the classifier parameters, i.e., network architecture and weights, as well as the details of any defense mechanism
- targeted attack: they attempt to cause the perturbed image to be misclassified to a specific target class
- untargeted attack: when no target class is specified

- BLACK-BOX ATTACK MODELS

- black-box adversaries have no access to the classifier or defense parameters, It is further assumed that they do not have access to a large training dataset but can query the targeted DNN as a black-box.

Adversarial Attack

APE-GAN: adversarial perturbation elimination with GAN



$$l_{ape} = \xi_1 l_{mse} + \xi_2 l_{adv} + \xi_3 l_{sc}$$

$$\left\{ \begin{aligned} l_{mse} &= \frac{1}{WH} \sum_{w=1}^W \sum_{h=1}^H (X_{w,h} - G_{\theta_G}(X_{\epsilon})_{w,h}) \\ l_{adv} &= \sum_{n=1}^N [1 - \log D_{\theta_D}(G_{\theta_G}(X_{\epsilon}))] \\ l_{sc} &= \frac{1}{WH} \sum_{w=1}^W \sum_{h=1}^H \|\nabla G_{\theta_G}(X_{\epsilon})_{w,h}\| \end{aligned} \right.$$

Total loss:

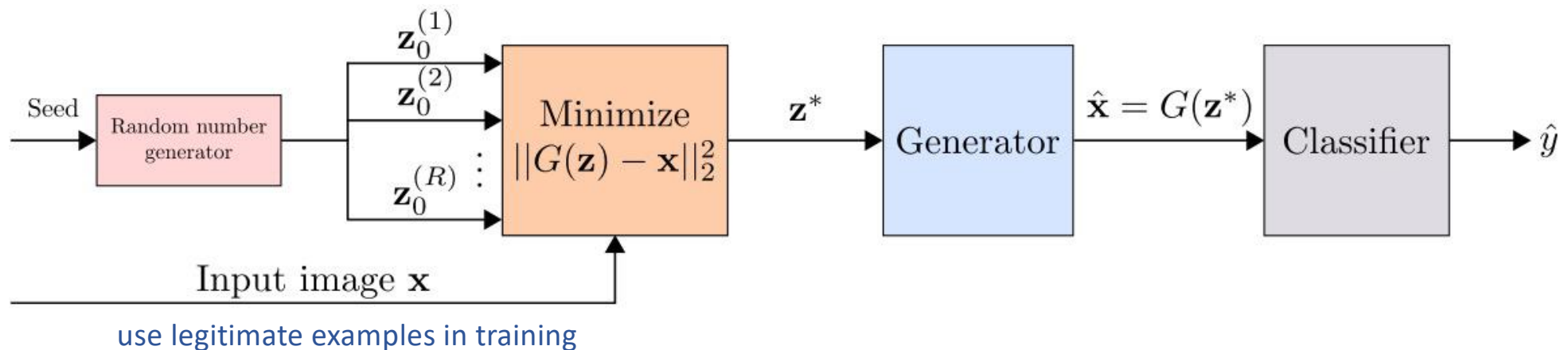
$$\hat{\theta}_G = \arg \min_{\theta_G} \frac{1}{N} \sum_{k=1}^N l_{ape}(G_{\theta_G}(X_{\epsilon}^k), X^k)$$

Adversarial training:

$$\min_{\theta_G} \max_{\theta_D} \mathbb{E}_{X \sim p_{data}(X)} \log D_{\theta_D}(X) - \mathbb{E}_{X_{\epsilon} \sim p_G(X_{\epsilon})} \log(D_{\theta_D}(G_{\theta_G}(X_{\epsilon})))$$

Adversarial Attack

- Defense-GAN



- a new defense strategy which uses a WGAN trained on legitimate (un-perturbed) training samples to “denoise” adversarial examples.

Adversarial Attack

- APE-GAN:
 - Use **adversarial samples** as the input of the generator.
- Defense-GAN:
 - Use **multiple random noise** as the input of the generator.
 - Implement adversarial training without using adversarial samples as inputs.
- Both of the structures are based on WGAN.

Application of Generative Models: (Selected) Advanced Topics

- Domain Adaptation
- Adversarial Attack
- **Meta Learning**
- Reinforcement Learning

Meta Learning

- **Definition**

- In the context of machine learning, meta learning is the process of **learning to learn**.
- Informally speaking, a meta learning algorithm uses experience to change certain aspects of a learning algorithm, or the learning method itself, such that the **modified learner is better** than the original learner **at learning from additional experience**.

Meta Learning

- Meta Learning Architecture for few-shot learning with generative models

from training set unseen frame

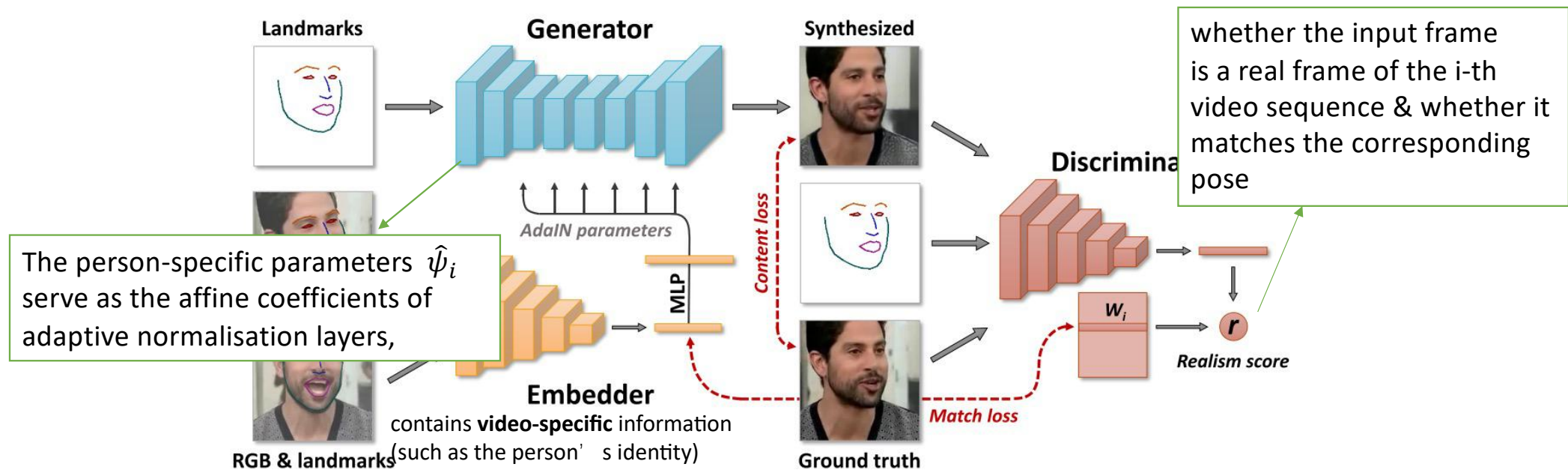
face landmark tracks



In fact, this system can generate a reasonable result based on a single photograph (one-shot learning), while adding a few more photographs increases the fidelity of personalization

Meta Learning

- Meta Learning Architecture for few-shot learning with generative models



- All parameters of the generator are split into two sets: person-generic parameters ψ , and the person-specific parameters $\hat{\psi}_i$.
- During meta-learning, ψ are trained directly, while $\hat{\psi}_i$ are predicted from the embedding vector \hat{e}_i using a trainable projection matrix \mathbf{P} : $\hat{\psi}_i = \mathbf{P}\hat{e}_i$

Meta Learning

- Meta Learning Architecture for few-shot learning with generative models



Meta Learning

- Meta Learning Architecture for few-shot learning with generative models

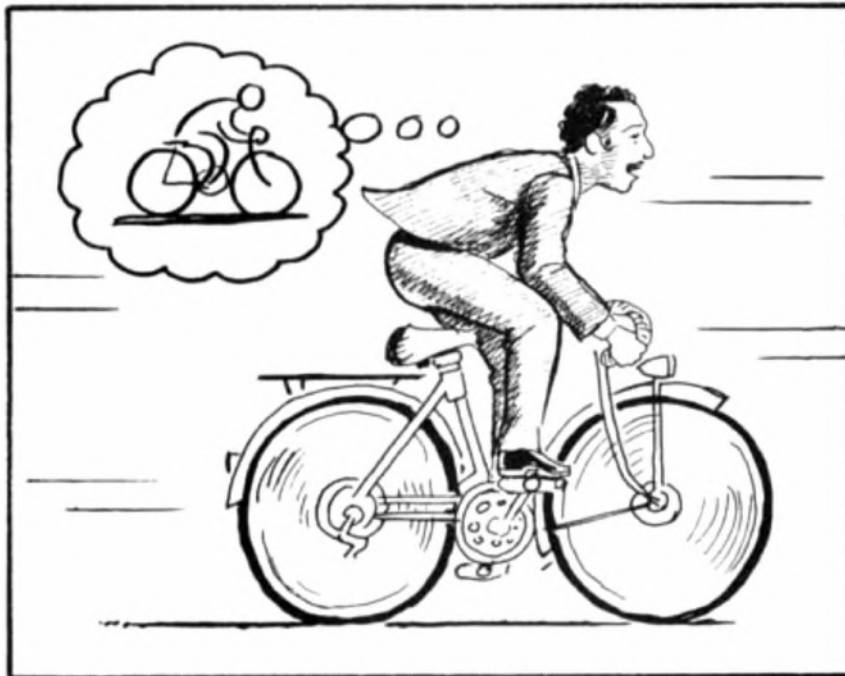


Application of Generative Models: (Selected) Advanced Topics

- Domain Adaptation
- Adversarial Attack
- Meta Learning
- Reinforcement Learning

Reinforcement Learning

- **World Models**



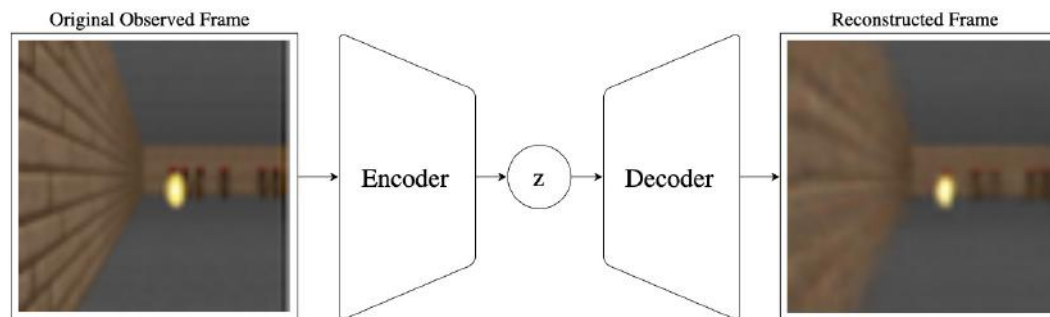
One way of understanding the predictive model inside of our brains is that it might not be about just predicting the future in general, but predicting future sensory data given our **current motor actions**

Learning in the imagination == Sampling efficiency

<https://worldmodels.github.io>

Reinforcement Learning

- **World Models**

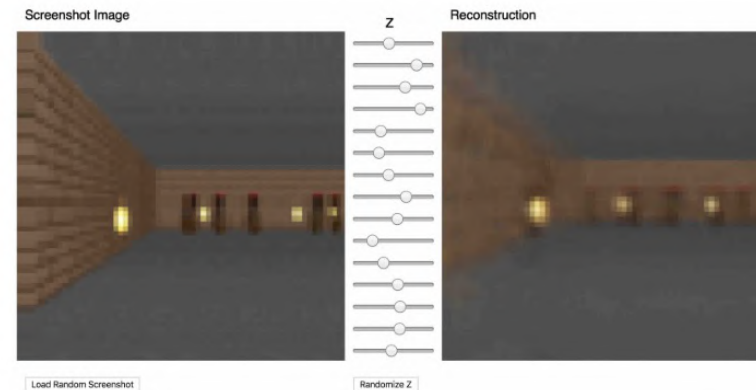
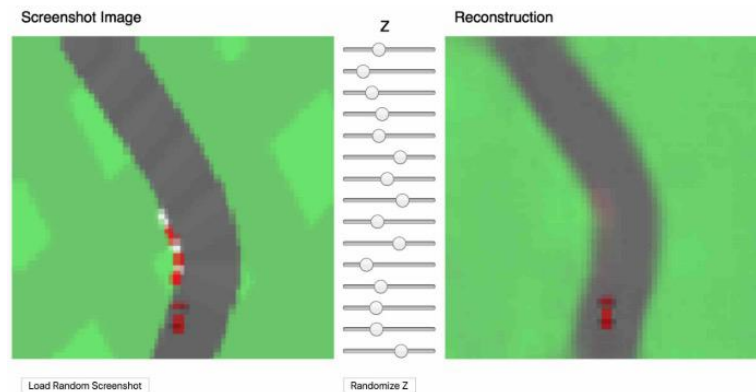


Learn the state representation

Here .. The encoder output is the state

Reinforcement Learning

- **World Models**

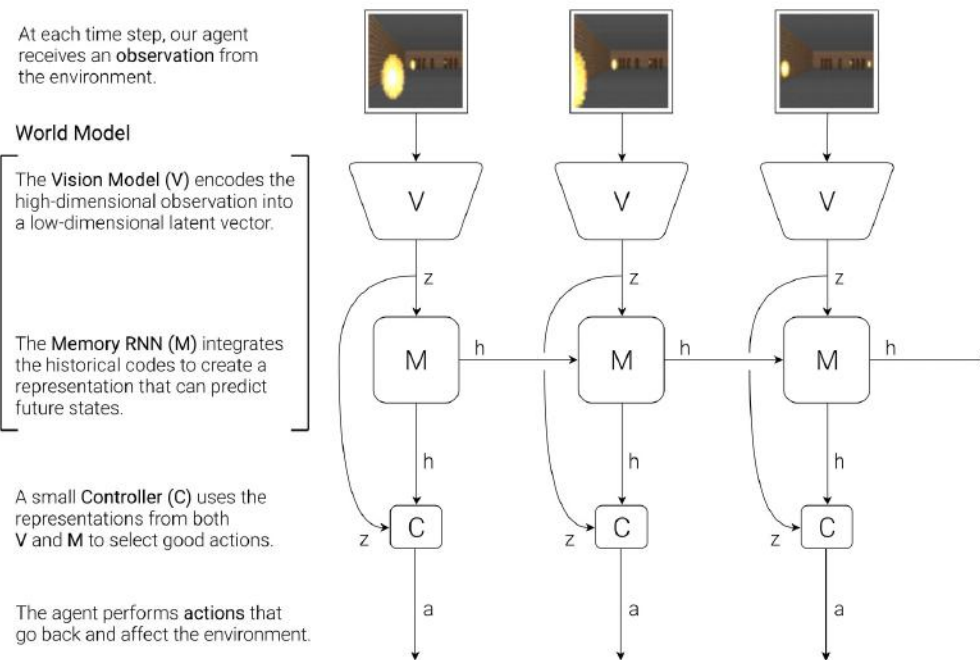


In this model, the agent has a visual sensory component that compresses what it sees into a small representative code.

Ha D, Schmidhuber J. World models[J]. arXiv preprint arXiv:1803.10122, 2018.

Reinforcement Learning

• World Models

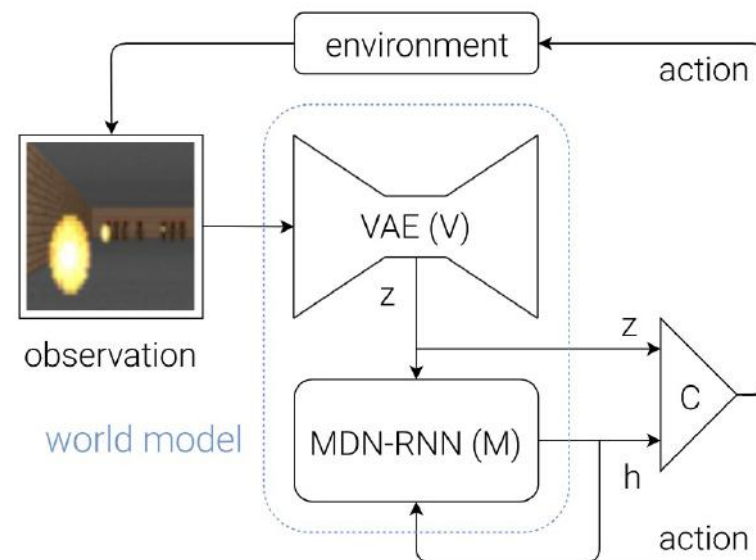


Our agent consists of three components that work closely together:
Vision (V), Memory (M), and Controller (C).

Learn the state representation

Here .. The encoder output is the state

RNN predicts the action



Summary



- Domain Adaptation
- Adversarial Attack
- Meta Learning
- Reinforcement Learning

Thanks