



CISPA

HELMHOLTZ CENTER FOR
INFORMATION SECURITY

Updates-Leak: Data Set Inference and Reconstruction Attacks in Online Learning

Ahmed Salem, Apratim Bhattacharya, Michael Backes

Mario Fritz, Yang Zhang

CISPA Helmholtz Center for Information Security, Max Planck Institute for Informatics

Online Learning



Training set

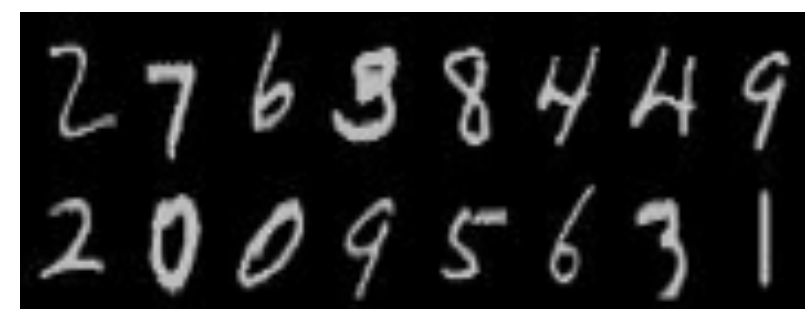


Train

Model

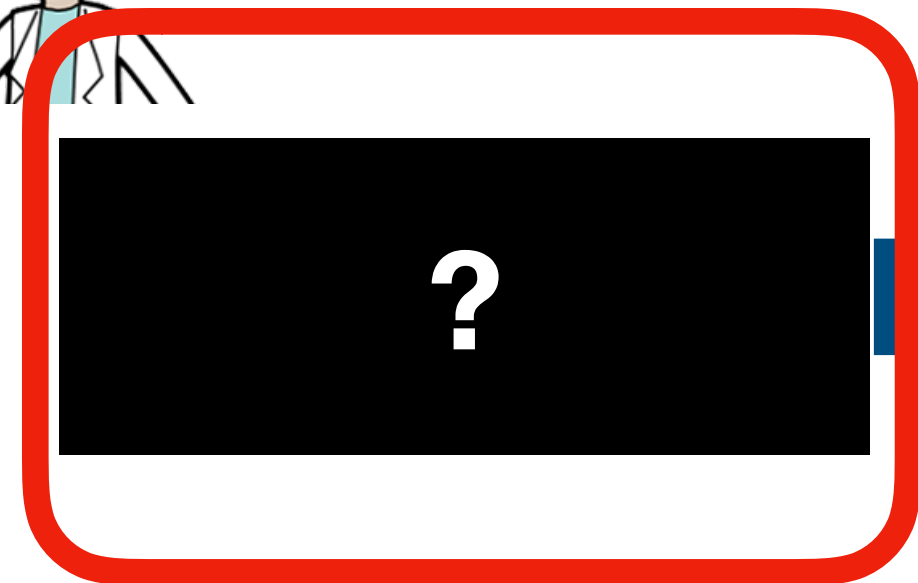
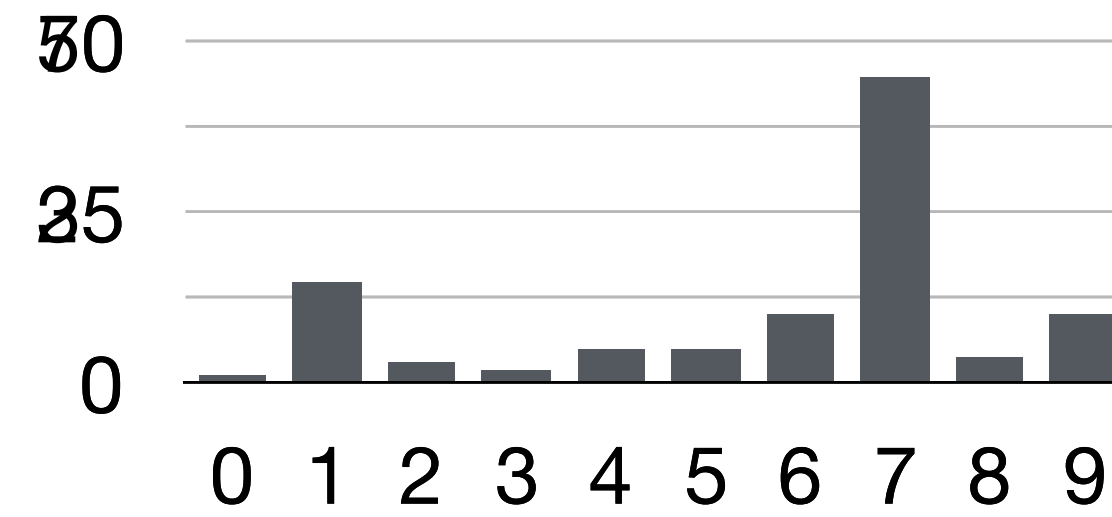
Update

Updating set

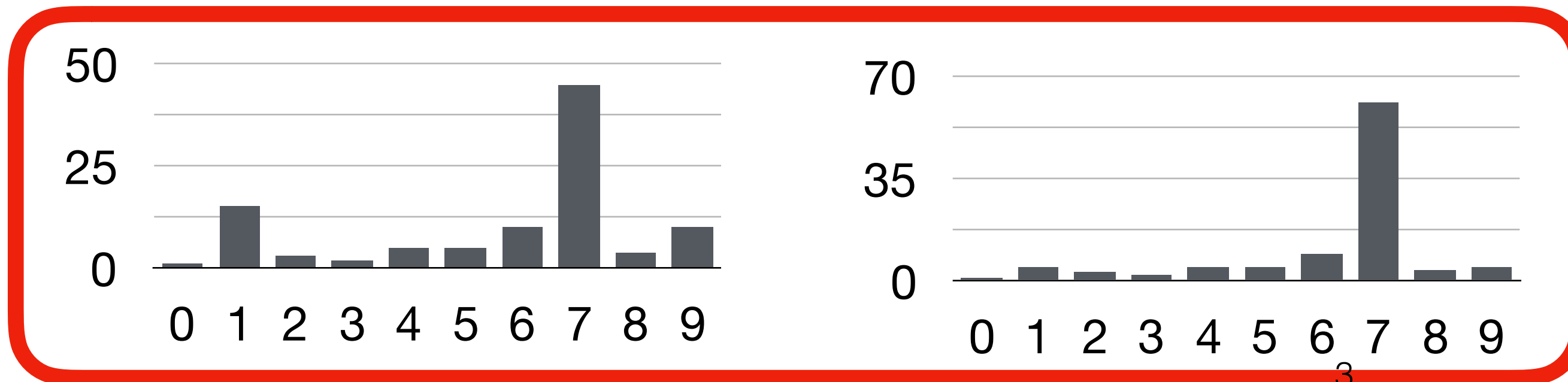


- Data generation rate
- 90% of the data in the world today has been created in the last two years alone
- Cost of retraining

Attack Surface in Online Learning



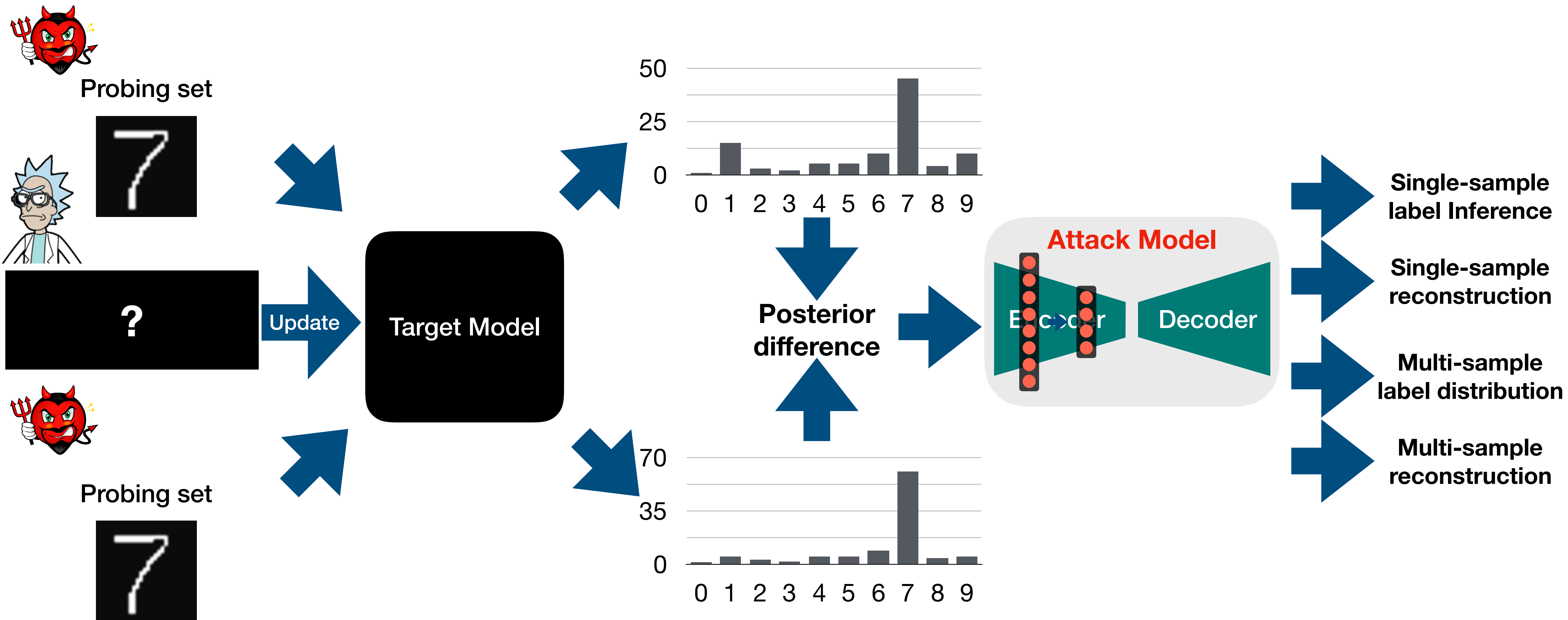
Research Question:
Can this posterior
difference be a new
attack surface?



Threat Model

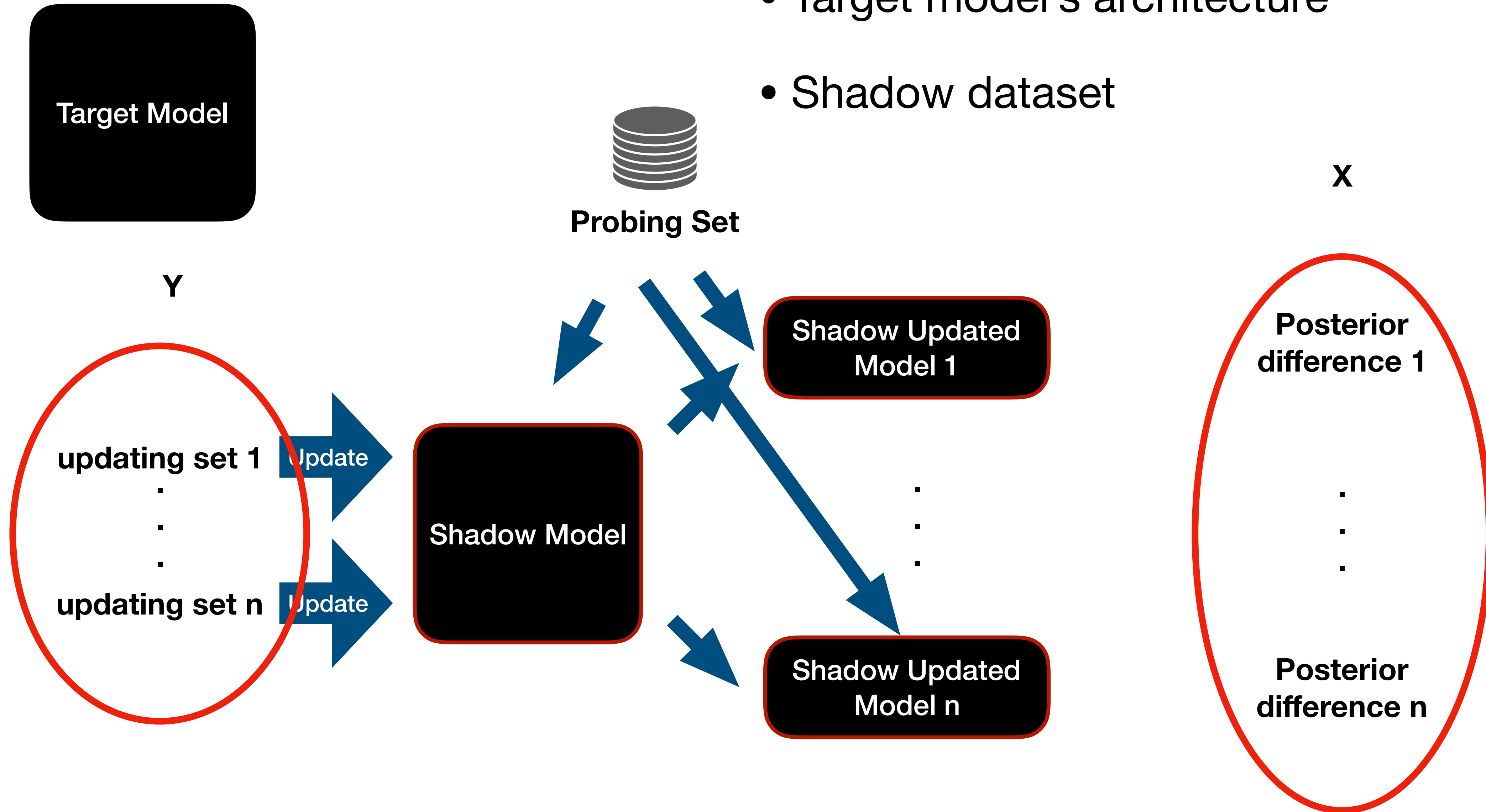
- Attacker has black-box access to the target model
- Attacker knows:
 - Target model's architecture
 - A shadow dataset from the same distribution of the target model's dataset

General Attack Pipeline

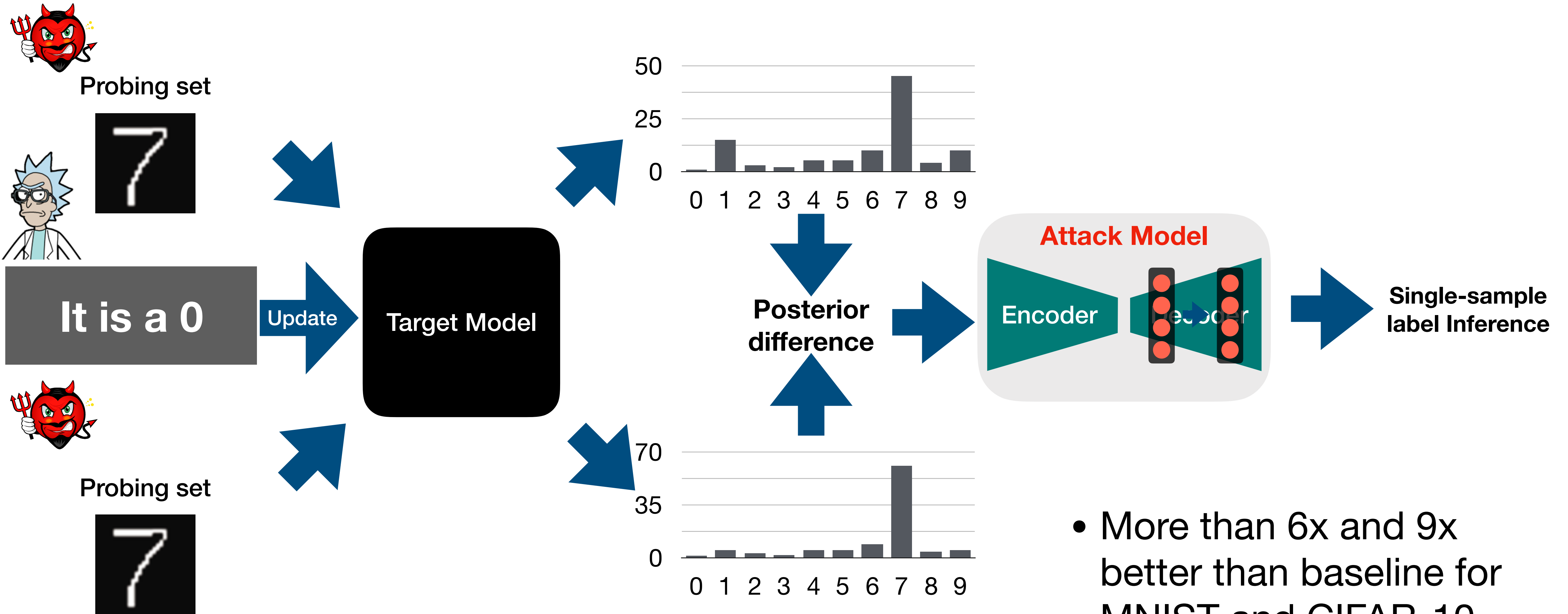


Attack Model Training

- Target model's architecture
- Shadow dataset

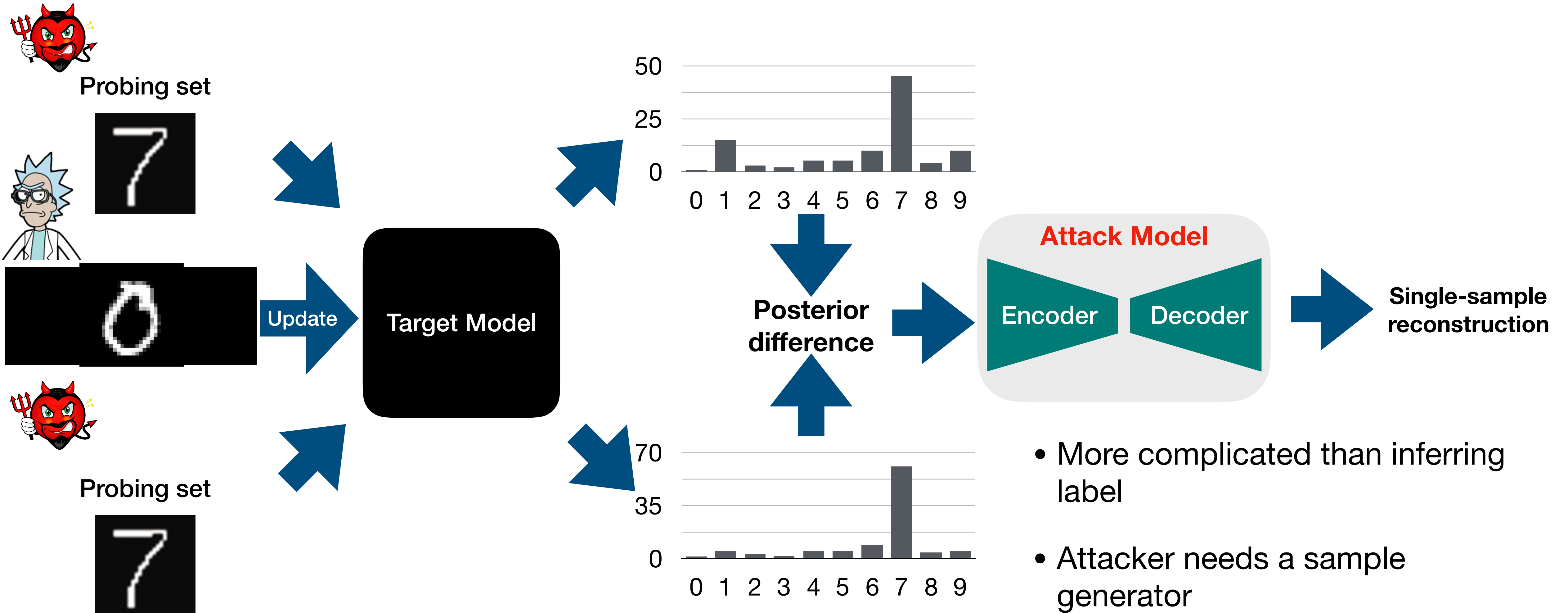


Single-sample Label Inference



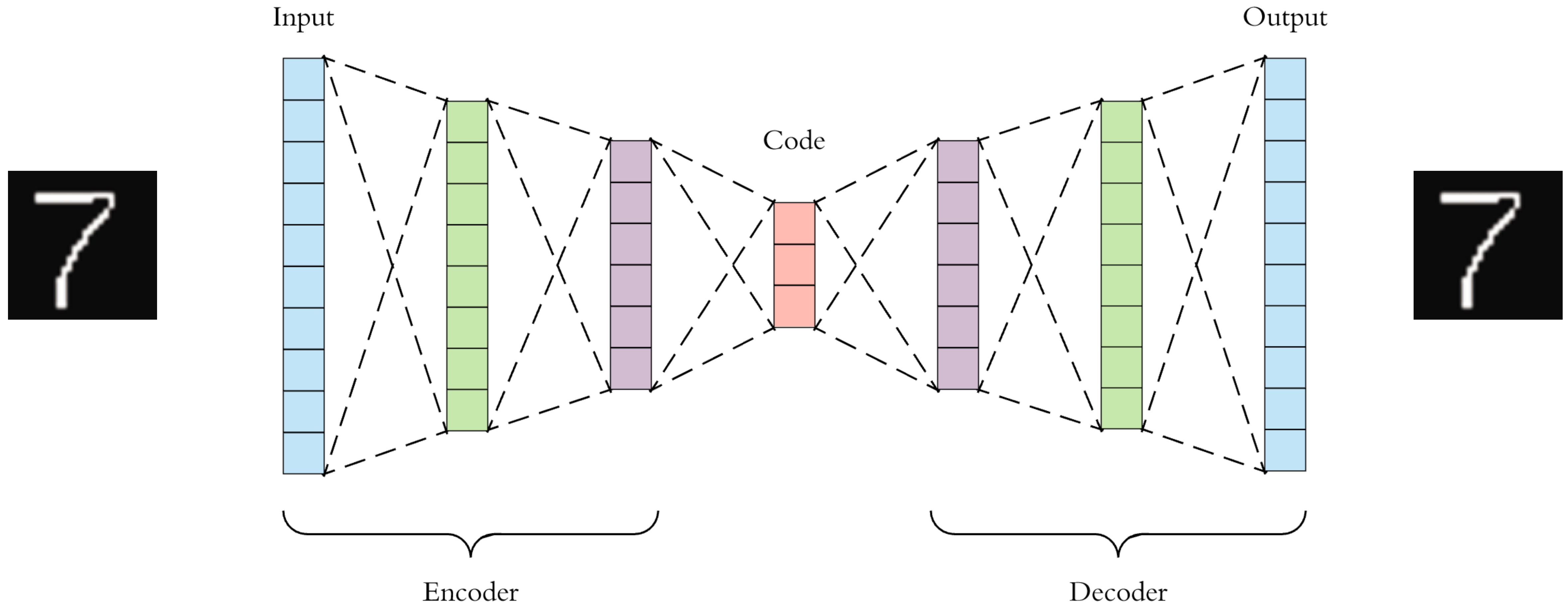
- More than 6x and 9x better than baseline for MNIST and CIFAR-10

Single-sample Reconstruction

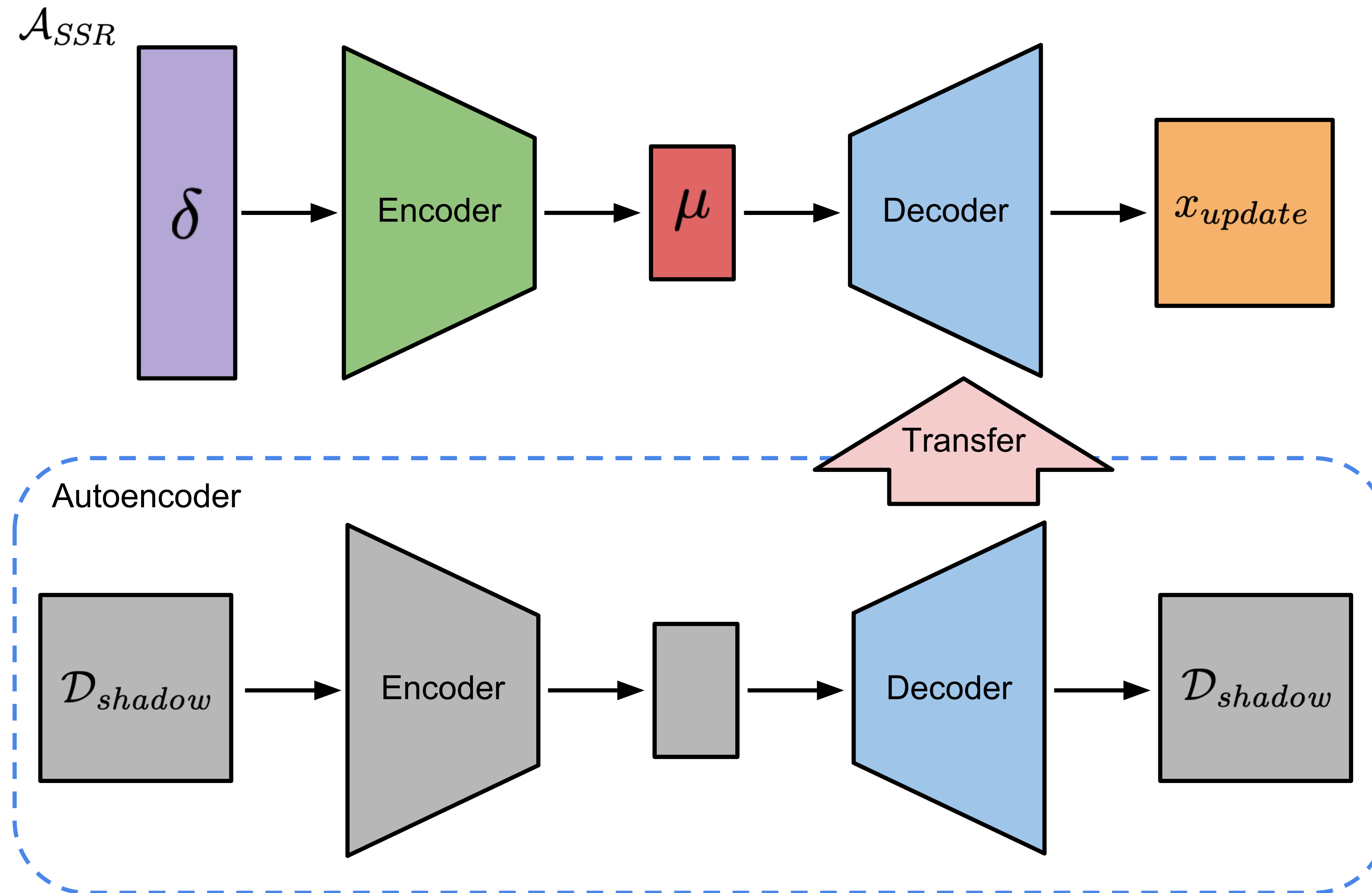


- More complicated than inferring label
- Attacker needs a sample generator
 - We rely on autoencoder's decoder

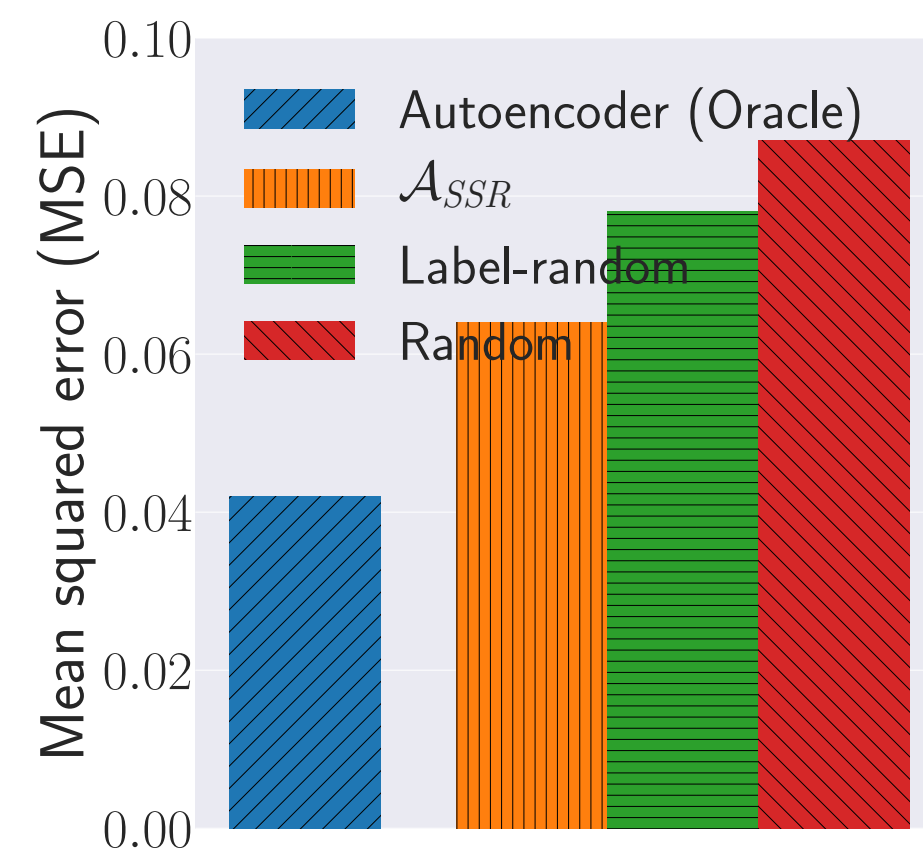
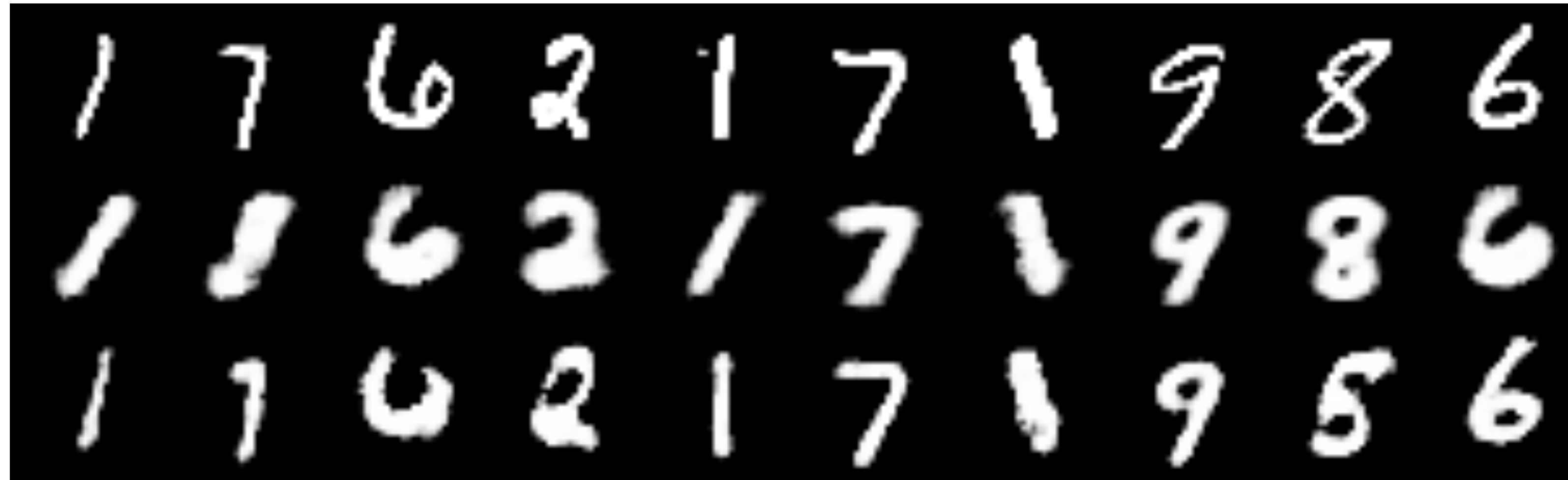
Autoencoder



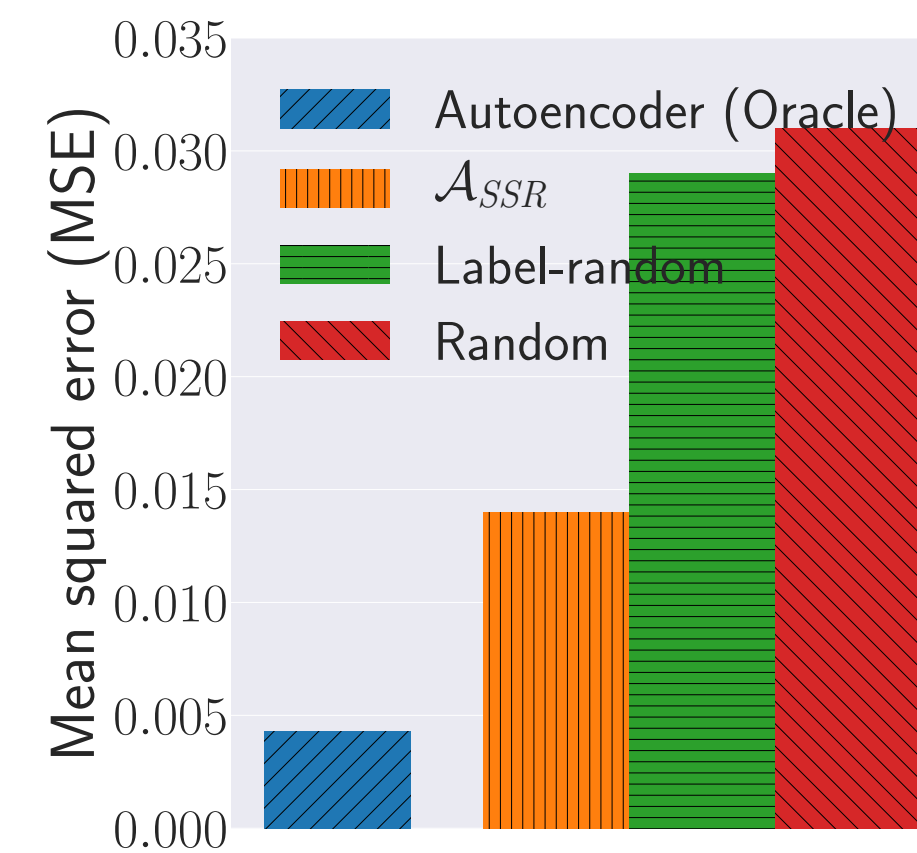
Single-sample Reconstruction



Single-sample Reconstruction

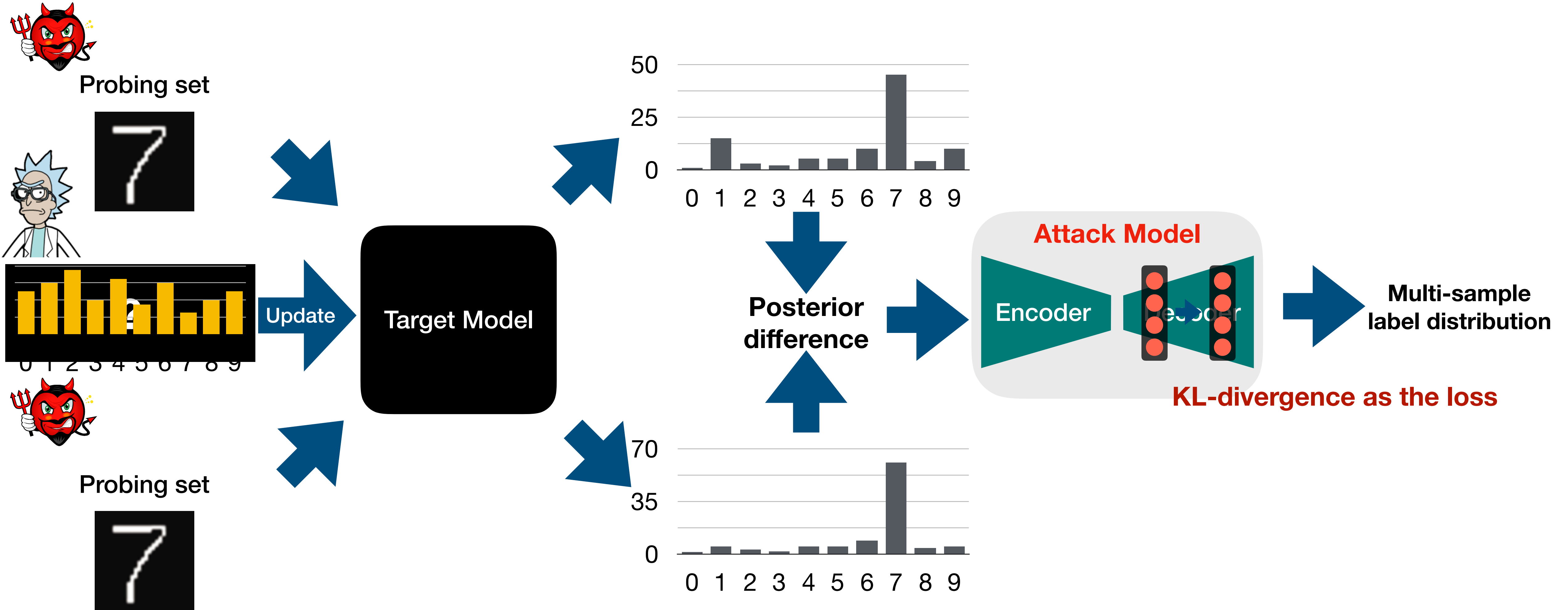


CIFAR-10

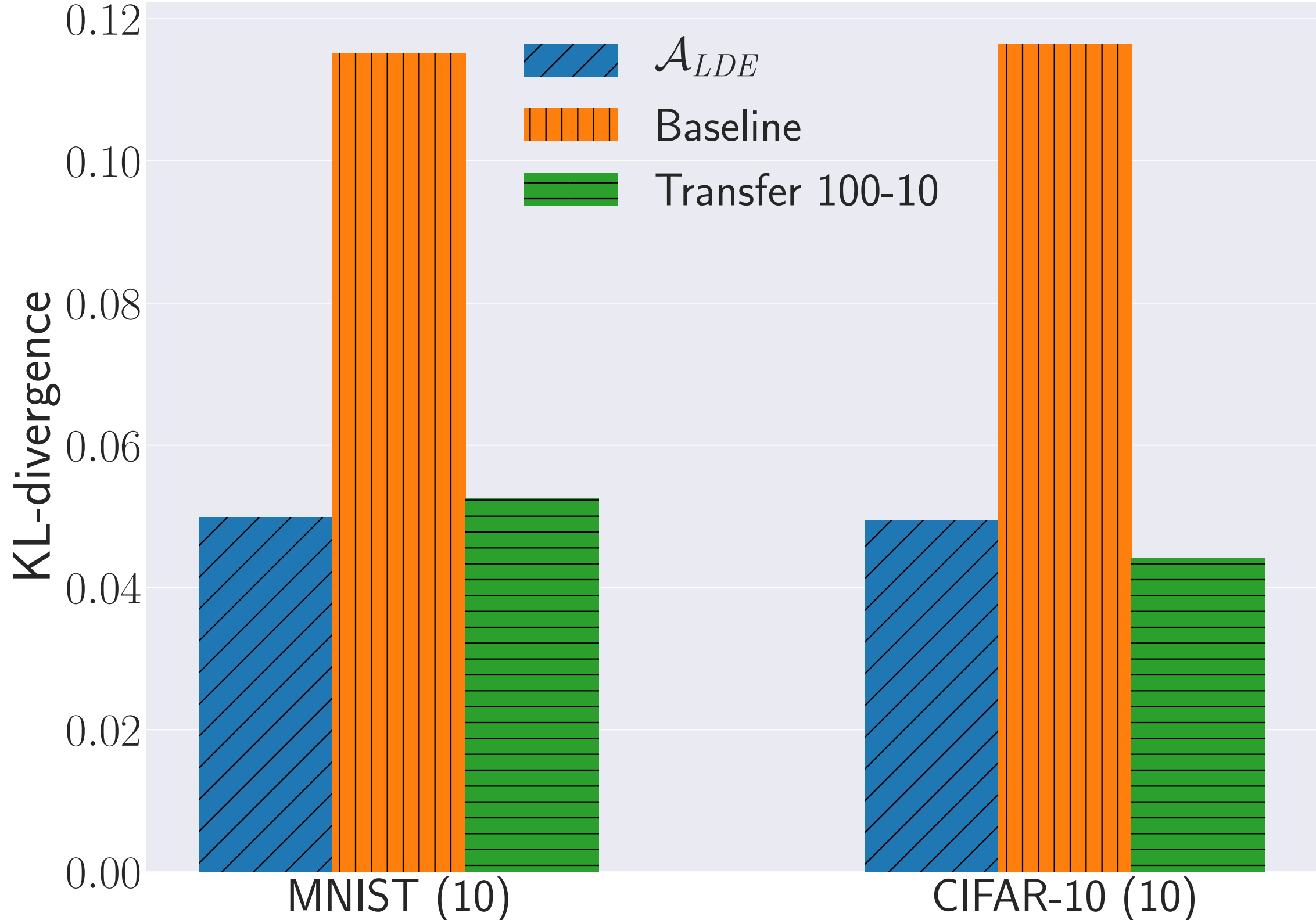
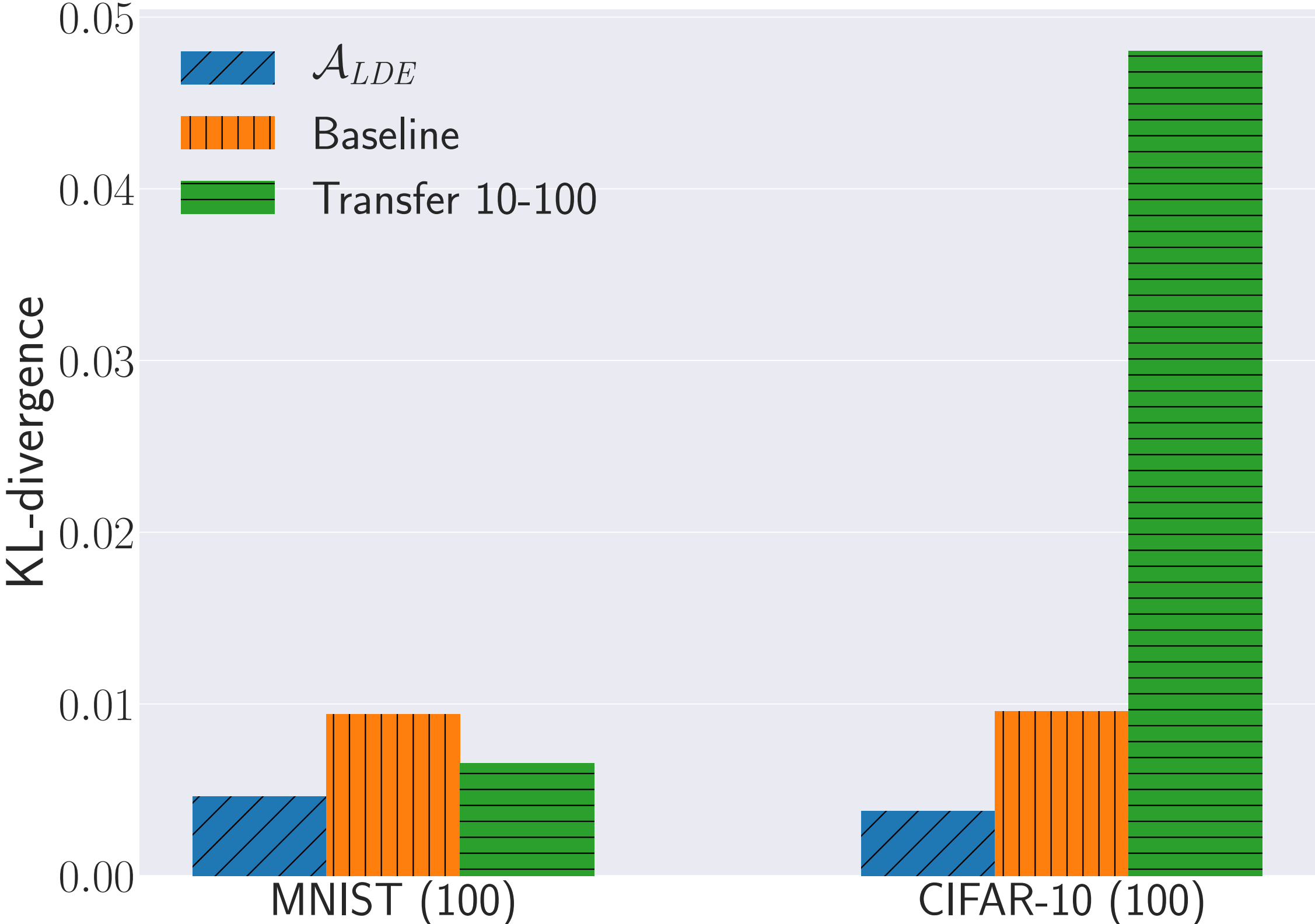


MNIST

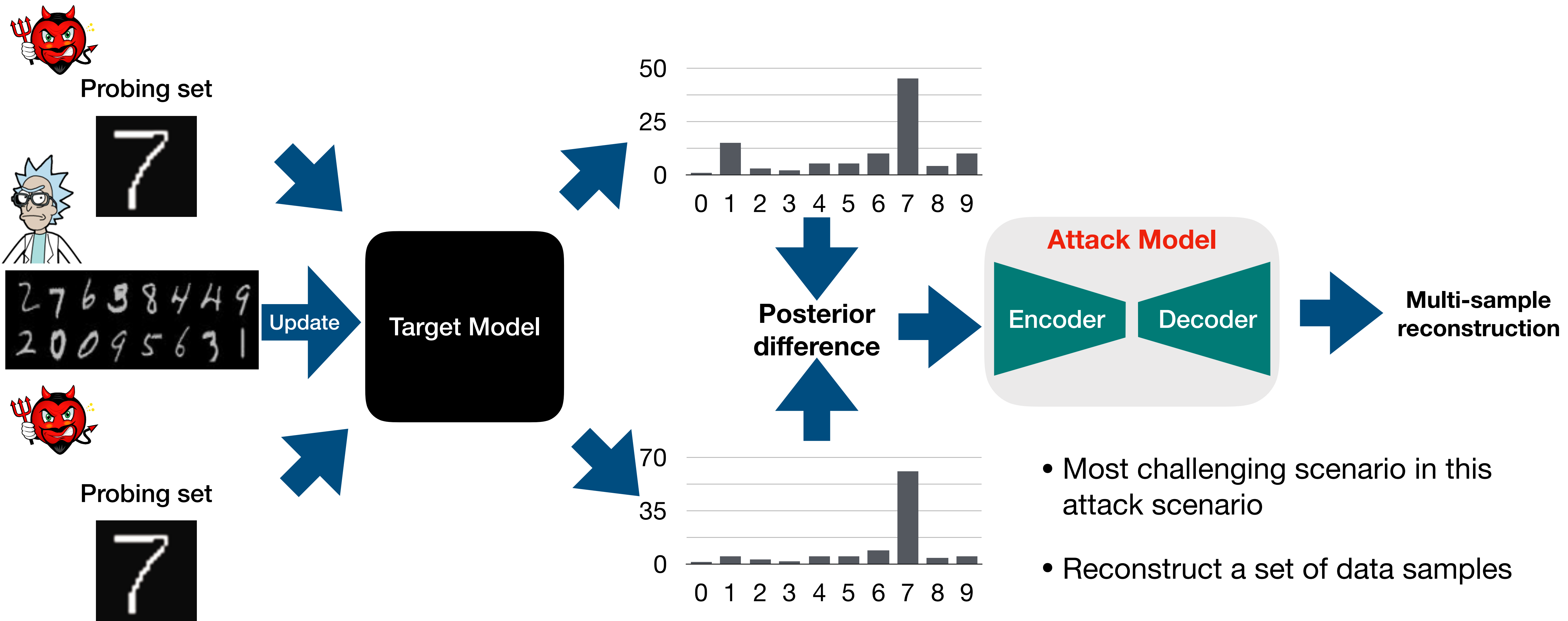
Multi-sample Label Estimation



Multi-sample Label Estimation



Multi-sample Reconstruction



- Most challenging scenario in this attack scenario
- Reconstruct a set of data samples
 - Autoencoder cannot help anymore
- What we do?

Generative Adversarial Network (GAN)

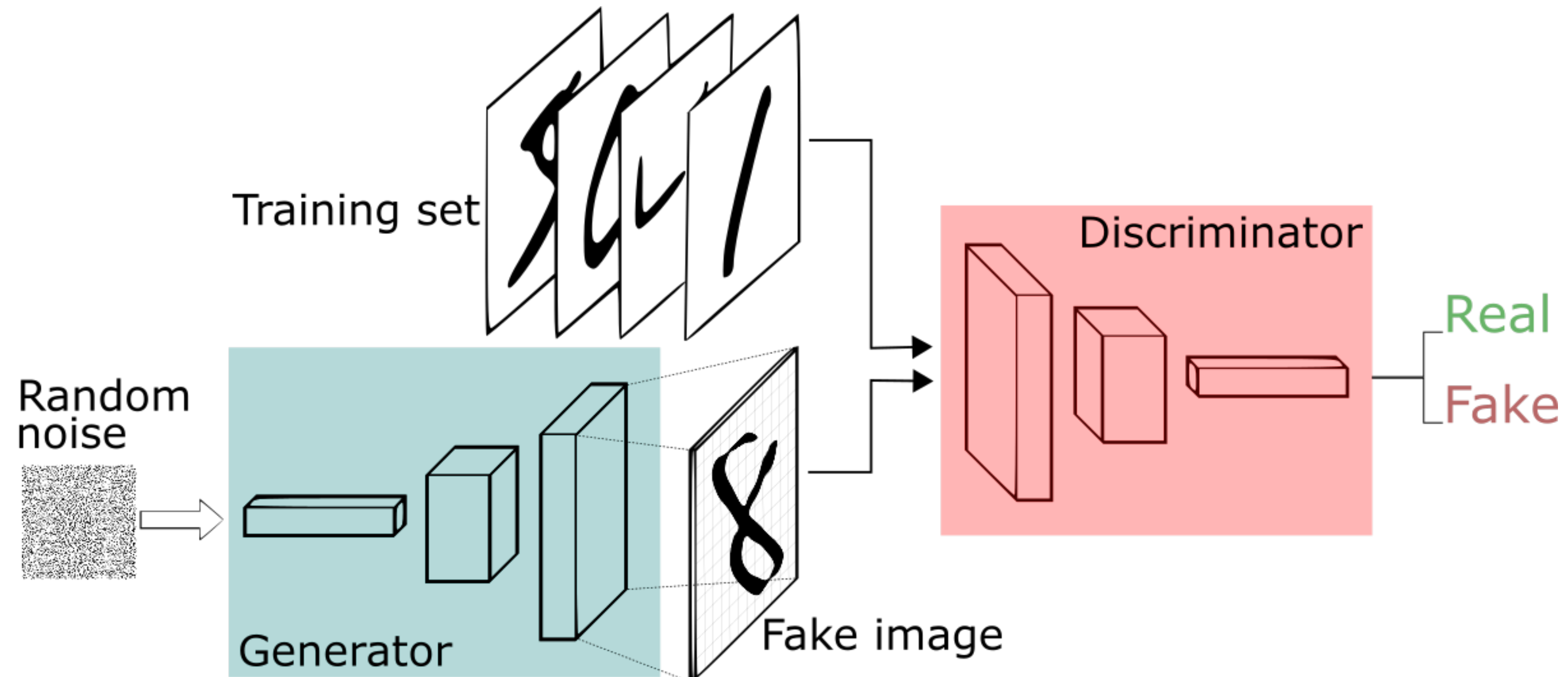
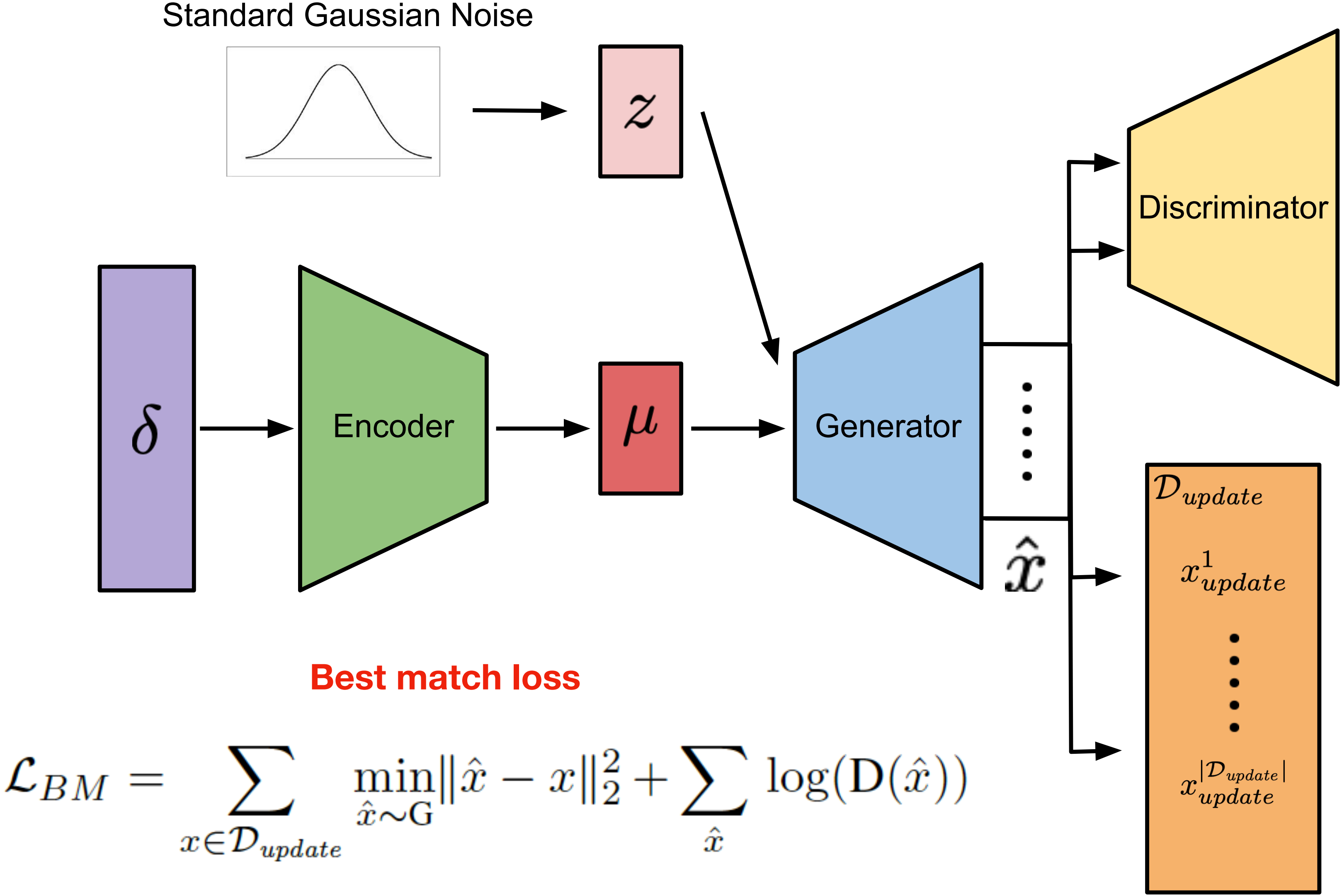
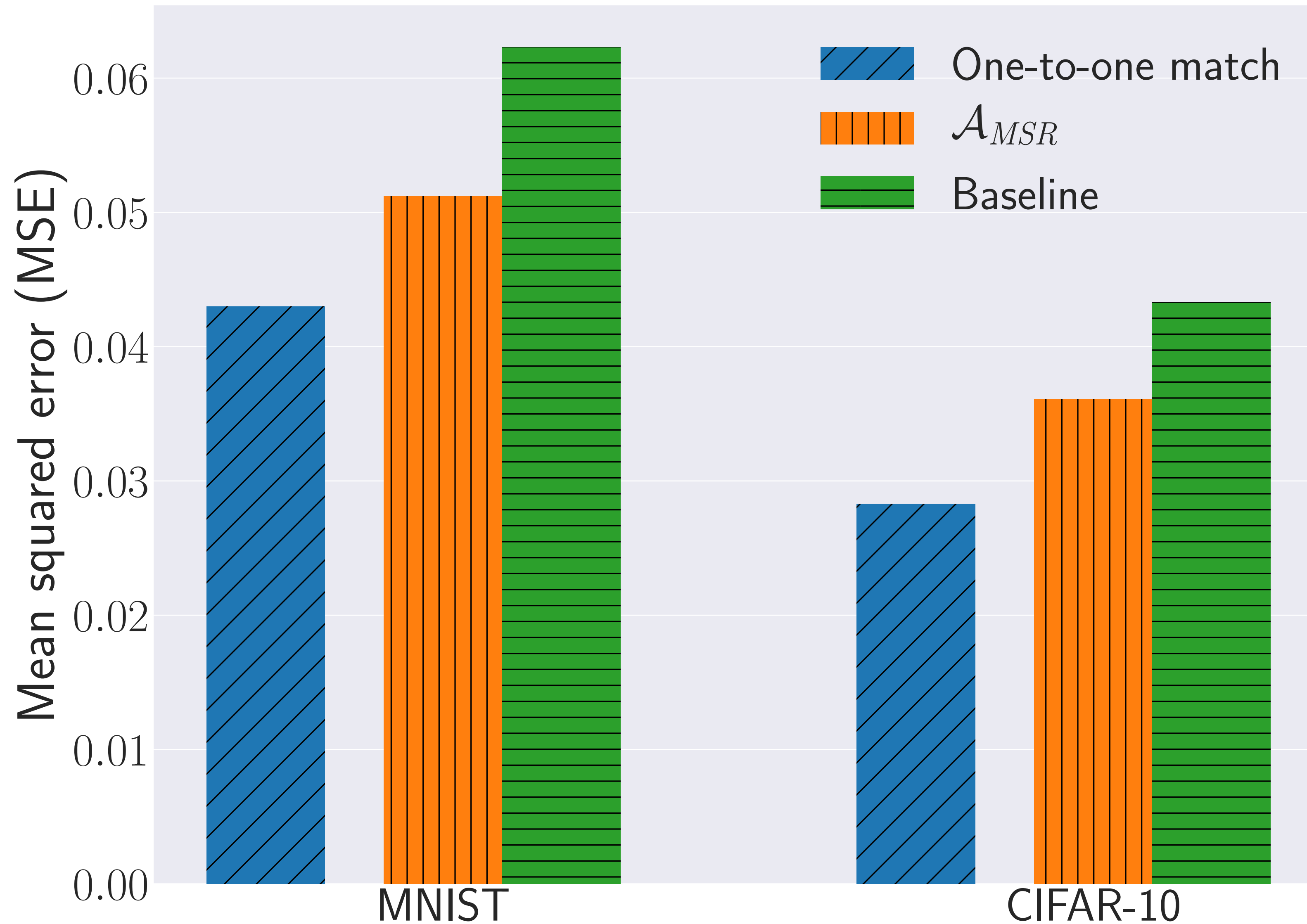


Image credit: [Thalles Silva](#)

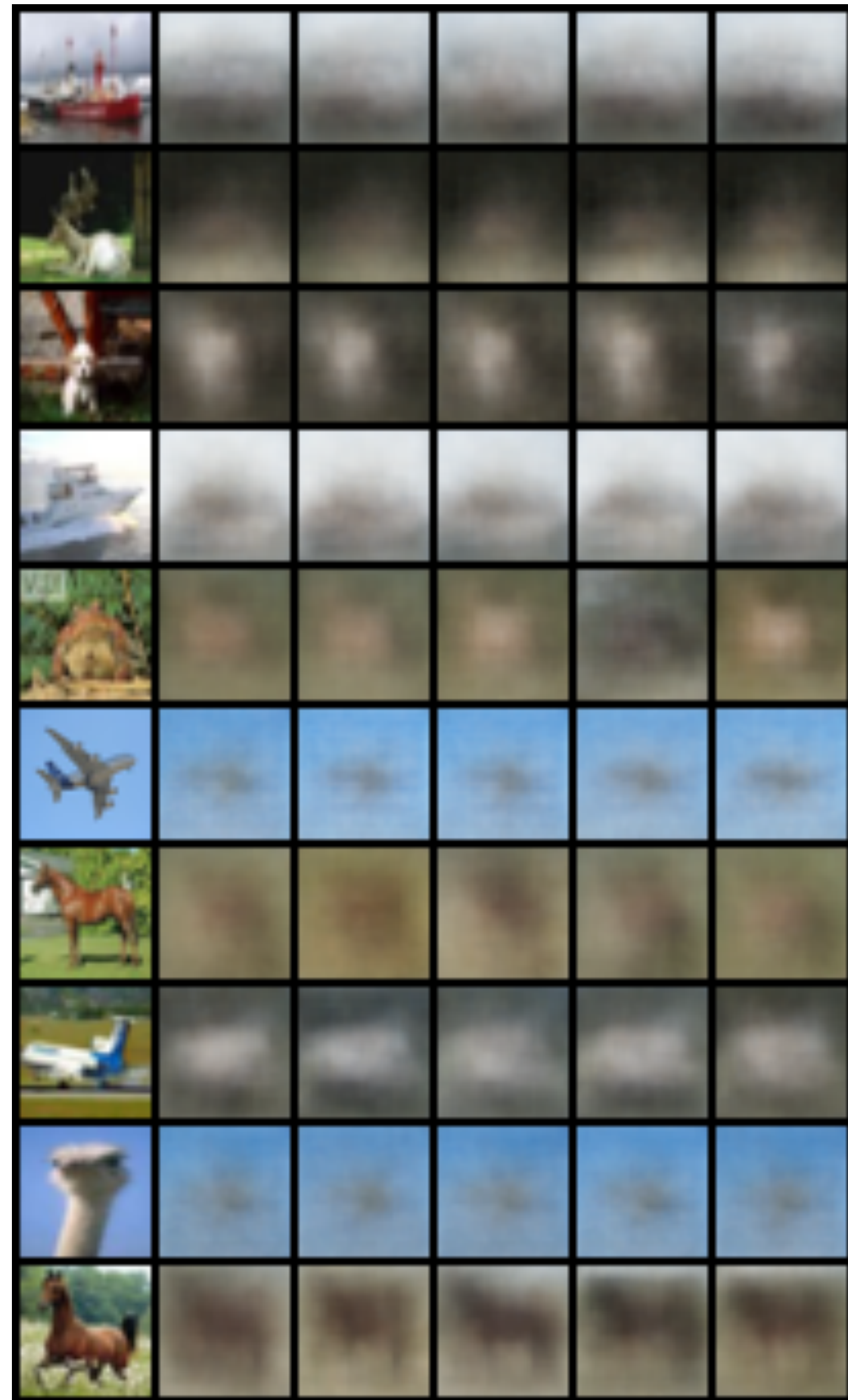
Multi-sample Reconstruction



Multi-sample Reconstruction



Multi-sample Reconstruction



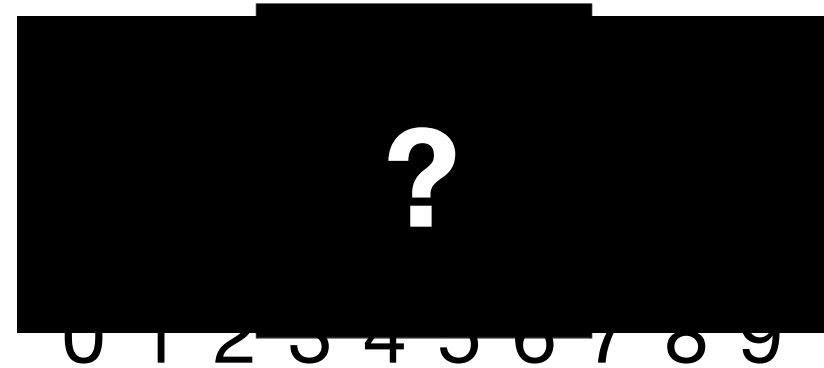
Multi-sample Reconstruction

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

Summary



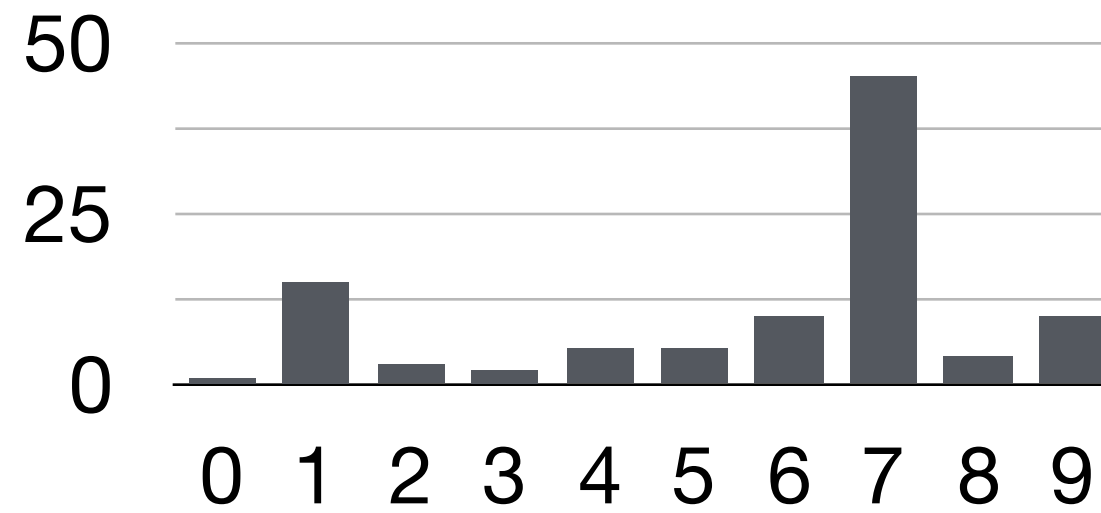
Probing set



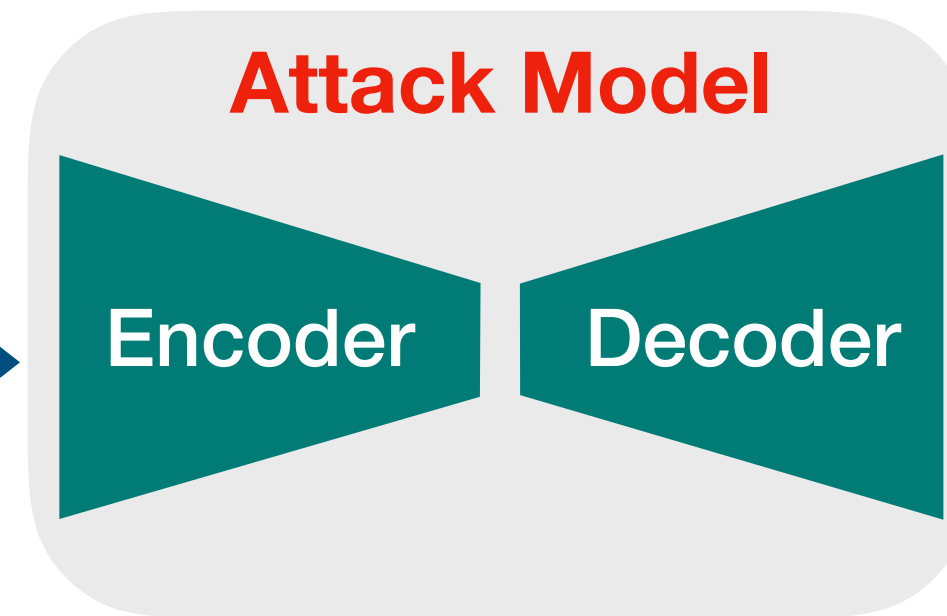
Update



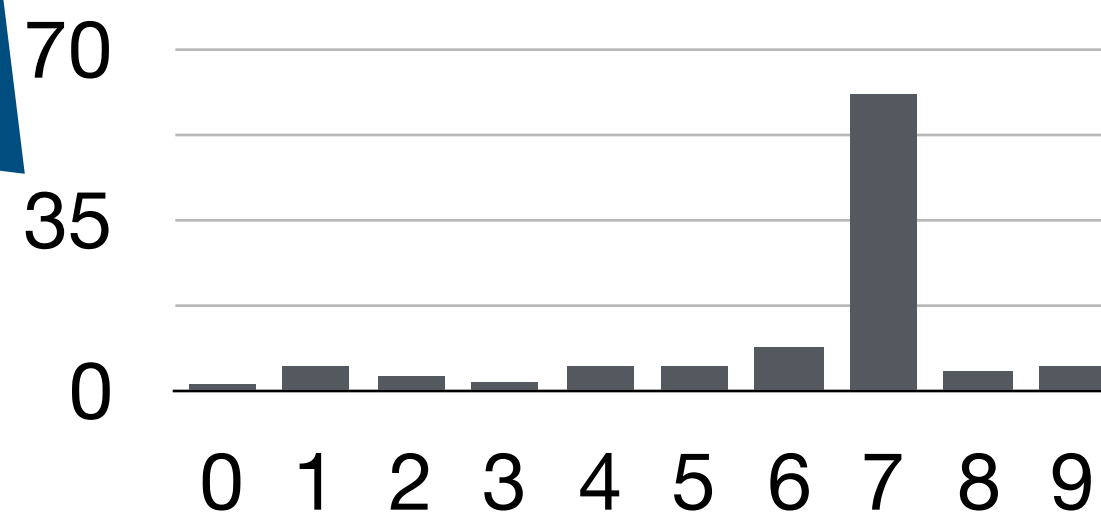
Probing set



Posterior
difference



Single-sample
label estimation



Thank you for your attention!
Questions?

ahmed.salem@cispa.saarland
<https://ahmedsalem2.github.io/>
@AhmedGaSalem