

Multi Trajectory Guided Model Inversion Attacks

A project report

submitted in partial fulfillment of the requirements

for the award of the degree of

Master of Technology

Submitted by

Indrajit Nandi

(CS2307)

under the supervision of

Dr. Sarbani Palit

Computer Vision & Pattern Recognition Unit,



Indian Statistical Institute

Kolkata-700108, India

June, 2024

CERTIFICATE

This is to certify that the dissertation entitled “ **Multi Trajectory Guided Model Inversion Attacks** ” submitted by by **Indrajit Nandi** Roll No. **CS2307** to Indian Statistical Institute, in partial fulfillment for the degree of Master of Technology in Computer Science is a bonafide record of work carried out by him under my supervision and guidance. The dissertation has fulfilled all the requirements as per the regulation of this institute and in my opinion has reached the standard needed for submission.

11 June, 2025

Signature: *Sarbani Palit*

Dr. Sarbani Palit

Computer Vision & Pattern Recognition Unit

Indian Statistical Institute

Kolkata-700108, India

ACKNOWLEDGEMENT

I would like to express my heartfelt gratitude to my supervisor, Dr. Sarbani Palit, for his invaluable guidance, support, and expertise throughout the entire duration of this master's dissertation. His insightful feedback, unwavering encouragement, and dedication to my academic growth have been instrumental in shaping the direction and quality of this research.

I am very much thankful to my parents and my family for always being there for me. Finally, I would like to thank all my friends and batchmates for their help and support.

June, 2025

Indrajit Nandi

CS2307

Indian Statistical Institute

Kolkata-700108, India

Contents

ACKNOWLEDGEMENT	2
ABSTRACT	5
1 Introduction	6
2 Related Work	10
2.1 Model Inversion Attacks:	10
2.2 Deep Reinforcement Learning:	11
3 Methodology	12
3.1 Problem Formation	12
3.2 Proposed method:	12
3.3 Model Architecture	14
3.4 Adaptive Alpha Scheduling:	15
3.5 Reward-Based Trajectory Optimization:	16
3.6 Multi-Candidate Selection:	17
3.7 Training and Optimization Procedure:	17
3.8 Network Update Procedure:	19
3.9 Network Overview:	19
4 Experiments	21
4.1 Experimental Setting& Implementation details:	21
4.2 Evaluative metrics:	23
4.3 Results	23
4.4 Randomly Selected Reconstructed Images:	24
4.5 Reconstructed Images from the Initial Classes:	25

5 Conclusion	27
5.1 Summary:	27
5.2 Future Works	28

ABSTRACT

Recent advancements in deep learning have brought significant concerns regarding the privacy of training data due to overfitting and memorization, as well as a lack of defense mechanisms. In model inversion attacks, where Attackers aim to reconstruct original private training samples (i.e images) just by using the DNN model’s last layer output. Traditional black-box MI attacks face three key challenges: (1) Non-convex optimization landscapes often trap reconstructions in poor local minima, degrading output quality. (2) Black-box scenarios require excessive queries to approximate gradients, raising detection risks. (3) Unconstrained pixel-space optimization generates unrealistic artifacts since the inverse mapping lacks natural image priors. These issues collectively yield low-fidelity reconstructions that fail to capture meaningful private data, especially for complex inputs like faces or medical images.

Here, Generative Adversarial Networks (GANs) come in picture which provide a low-dimensional latent space for efficient optimization while ensuring high-dimensional realism. Some recent methods have explored reinforcement learning, where the search in the latent space (learned from a GAN trained on public data) is formulated as a Markov Decision Process (MDP).

In our method, we introduce: (1) a dynamics network to model state transitions in the latent space, (2) a reward network that learns directly from the model’s confidence scores, and (3) policy-value networks for efficient exploration inspired by efficient-zero V2. This learned framework automatically adapts autonomously to diverse target models, and leverages the GAN’s generative prior to ensure realistic reconstructions.

Chapter 1

Introduction

Deep Neural Networks (DNNs) have achieved remarkable success across diverse domains such as computer vision, speech recognition, recommendation systems, natural language processing, and particularly in privacy-sensitive fields like healthcare and finance[1,2,3,4]. Their ability to extract meaningful patterns from large-scale data has led to widespread adoption in real-world applications, including disease diagnosis, personal assistants, and secure systems. However, these powerful models often rely on proprietary and sensitive datasets during training, which has raised significant privacy concerns as studies reveal their vulnerability to adversarial attacks that exploit model output (black-box attacks) or internal parameters (white-box attacks). An attack can be Black-Box or White-Box.

Black-Box Access refers to interactions with a model where the attacker can only observe its inputs and outputs, without knowledge of its internal structure or parameters. The adversary submits queries and analyzes predictions (e.g., confidence scores or class labels) but cannot inspect weights, gradients, or architecture. This setting is common in APIs or deployed services, where privacy attacks such as model inversion or membership inference rely solely on output exploitation. Black-box threats are often more practical, but less precise than white-box attacks.

On the other hand, White-Box Access grants full or partial visibility into the model's internals, including its architecture, trained parameters, and intermediate computations like gradients. Attackers with this access can directly manipulate or extract sensitive details, such as stealing proprietary weights or reconstructing training data via gradient leakage. White-box attacks are more powerful but require elevated privileges, making

them rarer in practice. This scenario is typical in federated learning or open-source model deployments.

Membership Inference Attack determines whether a particular data instance was used in training a target model. As introduced by Shokri et al. [5], this attack works by analyzing the model’s prediction confidence scores - trained models often show higher confidence for data they were trained on, allowing attackers to distinguish members from non-members through careful analysis of output probabilities.

Property Inference Attack aims to uncover general characteristics or statistical patterns about the training dataset. Research on Hidden Markov model (HMMs), Support Vector Machines (SVMs) [6] and Deep Neural Networks(DNNs) [7] has shown that models may inadvertently reveal aggregated properties like demographic distributions or dataset biases. Attackers exploit these leaks by observing model behavior over multiple queries to reconstruct dataset-level information without accessing raw training data.

Model Stealing (Parameter/Hyperparameter Inference) attempts to replicate a target model’s architecture, parameters, or training configurations. Tramer et al. [9] demonstrated that model parameters can be extracted through systematic input-output analysis, while subsequent work [8,10] showed hyperparameter like regularization terms can be inferred from model behavior.

The most dominant privacy attack on DNN is the model inversion attack, in which attackers want to reconstruct the private training data of target DNN model(black-box or white-box access).Fredrikson et al. [11] first demonstrated this by extracting genomic data from a Warfarin dosage model.In another work of Fredrikson et al. [12], they reconstruct face image of a class of targeted DNN via last layer confidence scores .

Due to low image quality and getting stuck in local optima, being unable to reconstruct high-dimensional data from complex models. nowadays GANs have come into trend. Many white-box MI attacks can successfully recover high-quality private data using GANs[13,14,14,16]. In the training time of the GAN on some public data, they train the GAN discriminator not only on real vs. fake labels but also on soft labels of the target models.



Figure 1.1: Left was reconstructed image (By Us using momentum gradient descent)

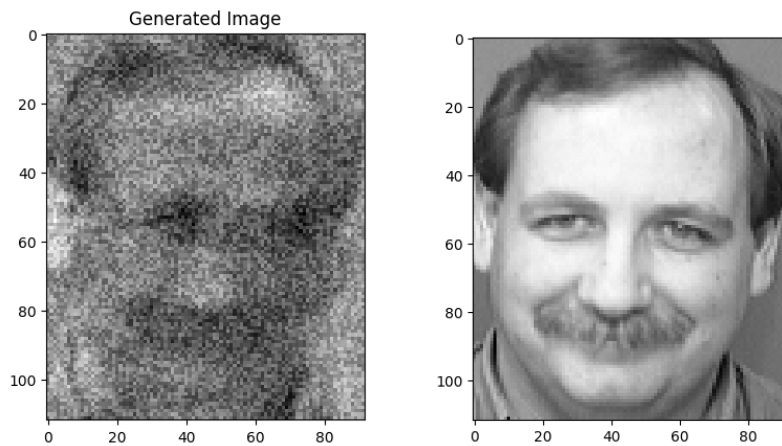


Figure 1.2: Left was reconstructed image (By Us using momentum gradient descent)

In recent days Model Inversion for Deep Learning Network (MIRROR) [17], Boundary-Repelling Model Inversion attack (BREP MI) [18], and RLB-MI [19] they are all search the latent space but in different ways.

In the mean time we draw inspiration from the recently proposed Reinforcement Learning-based Black-box Model Inversion attack (RLB-MI), which models the exploration of the GAN's latent space as a Markov Decision Process (MDP). This approach leverages reinforcement learning to effectively utilize confidence scores returned by the target model. In this framework, an agent receives rewards based on the target model's confidence in the generated image's class, and uses a replay memory to approximate the latent space dynamics. By learning from these interactions, the agent is able to navigate the latent space more efficiently than prior methods such as MIRROR or BREP-MI. They use Soft-Actor-Critic algorithm [28].

Instead of that, we use particular networks for each component like reward, policy, and dynamics, etc., inspired by EFFICIENT-ZERO[28], and we also optimize the queries to the target model (we can handle it by how many queries we want to give), which is totally an API system scenario.

Chapter 2

Related Work

2.1 Model Inversion Attacks:

Model inversion(MI) attacks compromise machine learning models by reconstructing sensitive training data.Fredrikson et al[12] first showed that such attacks could extract confidential attributes and reconstruct a facial images. One big development comes when Zhang et al. [15] introduced the Generative Model Inversion (GMI) attack, which employs a GAN [10] trained on public data to explore the latent space instead of the image space.Knowledge-Enriched Distributional Model Inversion (KED-MI) attack [13] enhanced GMI[15] by using an inversion-specific GAN, where the discriminator not only differentiates real from fake data but also predicts the target model’s outputs.

These are whit-box attack and However, white-box attacks are infeasible when model parameters are inaccessible, prompting the development of black-box attacks that rely only on soft labels.The Learning-Based Model Inversion (LB-MI) attack [26] adopted an auto-encoder-like architecture to train an inversion model, though it required both attacker-generated queries and confidence scores for user-provided data. Salem et al. [21] proposed a method that leaks training data information by comparing model outputs before and after updates, but this approach also depends on data that attackers typically cannot access. As in introduction we describe MIRROR [17] and Boundary-Repelling Model Inversion (BREP-MI) attack [18] ,and RLB-MI[19].

2.2 Deep Reinforcement Learning:

Deep reinforcement learning (RL) has evolved significantly since Deep Q-Network (DQN) [24] achieved human-level performance in Atari games [23]. While DQN excels in discrete action spaces, its limitations in continuous control motivated algorithms like Deep Deterministic Policy Gradient (DDPG) [25], which extended Deep Policy Gradient (DPG) [26] with actor-critic methods and stabilization techniques (e.g., target networks). Subsequent innovations, such as Twin Delayed DDPG (TD3) [27], addressed overestimation bias via clipped Q-learning and delayed updates, further improving robustness.

A major breakthrough came with Soft Actor-Critic (SAC) [28], which introduced entropy regularization to enhance exploration in complex environments. However, recent work by [29] (Rethinking Soft Actor-Critic in High-Dimensional Action Spaces) critically evaluates SAC’s scalability, revealing inefficiencies in high-dimensional action spaces due to its reliance on stochastic policies. The study argues that SAC’s entropy term, while beneficial for exploration, can lead to suboptimal convergence in tasks requiring precise, high-dimensional control (e.g., robotic manipulation).

In contrast, EfficientZeroV2 [31]—an extension of the model-based EfficientZero [19] framework—demonstrates superior sample efficiency and generalization by combining latent planning with RL. Unlike SAC, which struggles with dimensionality, EfficientZeroV2 leverages learned dynamics models and Monte Carlo Tree Search (MCTS) to reduce dependency on exhaustive environment interactions. This approach is particularly effective in sparse-reward settings, where SAC’s exploration mechanism may falter. Together, these advances highlight a trade-off: SAC’s strength in stochastic environments versus model-based methods’ precision in high-dimensional, structured tasks.

Chapter 3

Methodology

3.1 Problem Formation

Adversary’s Goal: So we consider a target model $T : x_p \rightarrow [0,1]^m$, where $x_p \in \mathbb{R}^d$ is an image from private training dataset D_p , each of dimension d . It outputs a probability vector over m classes. we want to reconstruct private data corresponding to a particular class y_1 .

Adversary’s Knowledge: As we operate in a black-box attack scenario, our access is limited to queries — specifically, the input data submitted to the target model and the corresponding softmax probability outputs. Additionally, we have a GAN trained on some public data D_{pub} , Adversary operates under a limited query budget to maintain stealth and avoid detection.

Reduced Formulation: Given a target deep neural network classifier $T : \mathbb{R}^{d_{dim}} \rightarrow [0, 1]^m$ with m classes and a pre-trained generative model $G : \mathbb{R}^{z_{dim}} \rightarrow \mathbb{R}^d$ that maps latent codes $z \in \mathbb{R}^z$ to the input space, our objective is to reconstruct representative samples for each target class $y \in 1, 2, \dots, m$. The model inversion attack aims to find an optimal latent code z^* such that: $z^* = \arg \max_{z \in \mathbb{R}^z} P(T(G(z)) = y)$ where $P(T(G(z)) = y)$ represents the probability that the target model classifies the generated image $G(z)$ as class y .

3.2 Proposed method:

We formulate the model inversion attack as a Markov Decision Process (MDP) operating in the latent space of a pre-trained generator. This formulation enables systematic explo-

ration of the generative manifold through reinforcement learning, where the agent learns to navigate toward latent codes that produce high-confidence classifications for target classes. The MDP framework provides a principled approach to balance exploration and exploitation in high-dimensional latent spaces. So we use 5 networks :

Representation Network $f_s : R^{z_{dim}} \rightarrow R^{h_{dim}}$

Dynamics Network $f_h : R^{h_{dim}+a_{dim}} \rightarrow R^{a_{dim}}$

Reward Network $f_r : R^{h_{dim}} \rightarrow R$

Policy Network $f_p : R^{h_{dim}} \rightarrow R^{a_{dim}}$

Value Network $f_v : R^{h_{dim}} \rightarrow R$

Here, R is the set of real numbers.

State Space Definition: The state space \mathcal{S} consists of hidden representations derived from latent codes through the representation network. Given a latent code $z \in R^{z_{dim}}$, the corresponding state is defined as: $s_t = f_s(z_t) \in R^{h_{dim}}$ where f_s is the representation network that maps raw latent codes to structured hidden states. These hidden states serve as compact representations that capture essential information about the current position in the latent space while providing a suitable foundation for value function approximation and policy learning. The h_{dim} -dimensional hidden state space is sufficiently expressive to encode complex latent space structures while remaining computationally tractable for neural network processing. The initial state for each episode is obtained by sampling a random latent code $z_0 \sim \mathcal{N}(0, I)$ and computing $s_0 = f_s(z_0)$. This random initialization ensures diverse starting points across episodes and prevents the agent from getting trapped in local optima during the search process.

Action Space Formulation: The action space \mathcal{A} consists of a_{dim} -dimensional continuous vectors that guide latent space modifications: $\mathcal{A} = R^{a_{dim}}$ Actions are generated through two distinct strategies implemented in the code:

Exploitation Actions (p probability): The agent generates multiple candidate actions and selects the one with highest predicted reward:

$$\begin{aligned} \text{candidates} &= a_1, a_2, \dots, a_{10} \text{ where } a_i \sim N(0, I) \\ a_t &= \text{argmax}_{a_i} f_r(f_h(s_t, a_i)) \end{aligned}$$

Exploration Actions ((1-p) probability): Random actions sampled from a standard normal distribution: $a_t \sim N(0, I)$ The selected action directly influences latent code

evolution through the adaptive mixing mechanism: $z_{t+1} = \alpha \cdot z_t + (1 - \alpha) \cdot a_t[: z_{dim}]$, where only the first z_{dim} dimensions of the action vector are used to match the latent code dimensionality.

State Transition Dynamics: State transitions are modeled by the dynamics network $f_h : R^{h_{dim}} \times R^{a_{dim}} \rightarrow R^{a_{dim}}$, which predicts the next hidden state given the current state and action: $s_{t+1} = f_h(s_t, a_t)$. The dynamics network implements a deterministic transition function that learns to model the latent space structure through experience. This deterministic formulation enables the agent to learn predictable transition patterns while maintaining computational efficiency. The dynamics network is trained to minimize the mean squared error between predicted and actual next states, ensuring accurate modeling of latent space dynamics.

Reward Signal Design: The reward function f_r measures the classification confidence of generated images for target classes. The reward is computed sparsely to balance computational efficiency with learning signal quality and efficient quarry optimization:

True Reward Computation: For a given latent code z_t , the reward is:

$$r_t = \text{softmax}(T(G(z_t)))[c]$$

where T is the target classifier, G is the generator, and c is the target class label. This reward signal directly measures the attack’s success by quantifying how confidently the target model classifies the generated image as the desired class.

Reward Prediction: Between true reward evaluations, the reward network f_r provides estimated rewards from hidden states: The reward network is trained using the sparse true rewards as supervision, enabling continuous reward estimation without expensive generator evaluations. This design significantly reduces computational cost while maintaining effective learning signals for the reinforcement learning agent. The MDP formulation enables efficient exploration of the latent space by providing structured state representations, meaningful action spaces, predictable transition dynamics, and direct reward signals aligned with the attack objective.

3.3 Model Architecture

Our agent consists of four interconnected neural networks that collectively enable effective latent space navigation for model inversion attacks. The agent architecture is specifically

designed to handle the continuous latent space $R^{z_{dim}}$ while learning optimal policies for maximizing target model confidence.

Representation Network: The representation network encodes latent codes into a structured hidden state representation. This network transforms raw latent codes into meaningful representations that capture the semantic structure of the latent space.

Dynamics Network: The dynamics network predicts the next hidden state given the current state and action. This network implicitly learns residual updates to the current state.

Reward Prediction Network: The reward network estimates the expected reward for a given hidden state.

Policy and Value Networks: The policy network outputs continuous actions in the latent space, and the value network estimates the expected cumulative reward.

Specially, Actions represent displacement vectors in the latent space, enabling fine-grained control over generated image characteristics.

3.4 Adaptive Alpha Scheduling:

Alpha Parameter Significance: The alpha parameter $\alpha(t)$ in our latent code evolution equation serves as a critical hyperparameter that controls the balance between maintaining the current latent position (exploitation) and following the RL agent’s policy recommendations (exploration). The adaptive scheduling of this parameter is essential for achieving optimal attack performance across different phases of the learning process.

Mathematical Formulation: The latent code evolution is governed by:

$$z_{t+1} = \alpha(t) \cdot z_t + (1 - \alpha(t)) \cdot a_t[z]$$

where the adaptive scheduling function is defined as: $\alpha(t) = 0.1$ if $t / 1000 == 0$ and $y_t \neq \text{threshold}$, else $\alpha(t) = 0.3$

This formulation creates distinct behavioral phases within the first 1000 episodes, with different exploration characteristics.

Phase-Based Learning Strategy Exploration Phase ($\alpha = 0.1$) During the exploration phase, the low alpha value emphasizes the agent’s actions:

Action Influence: 90% of the latent update comes from the policy network Rapid

Movement: Enables quick traversal of the latent space Discovery Focus: Prioritizes finding promising regions over refinement

Mathematical Impact: $z_{t+1} = 0.1 \cdot z_t + 0.9 \cdot a_t[z]$ This configuration allows the agent to make substantial changes to the latent code, facilitating broad exploration of the latent manifold. Exploitation Phase ($\alpha = 0.3$) During the exploitation phase, the higher alpha value emphasizes stability:

State Preservation: 70% of the current latent code is retained Fine-Tuning: Smaller adjustments enable precise optimization Convergence: Promotes convergence to high-quality solutions

Mathematical Impact: $z_{t+1} = 0.3 \cdot z_t + 0.7 \cdot a_t[z]$ This configuration provides a balance between maintaining promising positions and allowing policy-guided improvements.

Algorithm ■ Alpha Value Scheduler

Input: Episode number t , counter y , threshold st

Output: Alpha value $\alpha(t)$

```

1: Initialize:  $st \leftarrow \text{max\_episodes}/1000$ 
2: for each episode  $t$  do
3:   Increment counter:  $y \leftarrow y + 1$ 
4:   if  $\lfloor t/1000 \rfloor == 0$  then
5:     if  $\alpha == 0.1$  and  $y \neq st$  then
6:       Set  $\alpha \leftarrow 0.3$  ▷ Switch to exploitation
7:     else
8:       Set  $\alpha \leftarrow 0.1$  ▷ Switch to exploration
9:     end if
10:  else
11:    Set  $\alpha \leftarrow \alpha_{\text{default}}$ 
12:  end if
13: end for

```

3.5 Reward-Based Trajectory Optimization:

The agent optimizes its trajectory through the latent space by maximizing expected cumulative rewards. The trajectory optimization objective is:

$$J(\pi) = E_{\tau \sim \pi} \left[\sum_{t=0}^T \gamma^t r_t \right]$$

where $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$ represents a trajectory and γ is the discount factor. Reward Sampling Strategy: To maintain computational efficiency, rewards are sampled stochastically:

With probability $p_{eval} = 0.1$, generate image $\hat{x} = G(z_t)$ and query target model With probability $1 - p_{eval}$, use predicted reward $\hat{r}_t = f_r(s_t)$

3.6 Multi-Candidate Selection:

Traditional reinforcement learning approaches select single actions from policy networks, which can lead to suboptimal decisions due to policy approximation errors and environmental noise. In the context of model inversion attacks, poor action selection can result in ineffective latent space navigation and failed reconstructions. Our multi-candidate action selection mechanism addresses this limitation by evaluating multiple potential actions and selecting the most promising one based on predicted rewards. It reduces Variance Reduction, Improved Exploration, not stuck in suboptimal actions.

Algorithm ■ Multi-Candidate Action Selection

Input: Current state s_t , policy π_ω , dynamics f_h , reward predictor f_r

Parameters: $K = 10$ (number of candidates)

Output: Selected action a_t

- 1: Generate K candidate actions: $\mathcal{A} = \{a_1, a_2, \dots, a_K\}$ where $a_i \sim \mathcal{N}(\pi_\omega(s_t), \sigma^2 I)$
 - 2: Create state replications: $\mathcal{S} = \{s_t, s_t, \dots, s_t\}$ (K identical copies)
 - 3: Predict next states: $\mathcal{S}' = \{f_h(s_t, a_i)\}_{i=1}^K$
 - 4: Evaluate predicted rewards: $\mathcal{R} = \{f_r(s'_i)\}_{i=1}^K$
 - 5: Select optimal action: $a_t = a_{\arg \max(\mathcal{R})} = a_j$ where $j = \arg \max_i R_i$
-

3.7 Training and Optimization Procedure:

The training procedure follows an online reinforcement learning paradigm where the agent learns to navigate the latent space through interaction with the environment. Unlike supervised learning approaches that require pre-collected datasets, our method learns directly from the target model’s responses, making it applicable to black-box scenarios

where training data is unavailable. we use circular buffer to store experience and we constructed batch from $\mathcal{B} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^B \sim \text{Uniform}(\text{ReplayBuffer})$

Algorithm ■ Training Procedure

Input: Target model T , Generator G , Agent networks

Parameters: max_episodes, batch_size, update_frequency

```

1: for episode = 1 to max_episodes do
2:   Initialize  $z_0 \sim \mathcal{N}(0, I)$ ,  $s_0 = f_s(z_0)$ 
3:   Set adaptive  $\alpha$ :  $\alpha = \alpha(\text{episode})$ 
4:   for step = 1 to episode_length do
5:     Select action using multi-candidate strategy
6:     Update latent code:  $z_{t+1} = \alpha \cdot z_t + (1 - \alpha) \cdot a_t[:z]$ 
7:     Predict next state:  $s_{t+1} = f_h(s_t, a_t)$ 
8:     Sample reward stochastically
9:     Store experience:  $(s_t, a_t, r_t, s_{t+1})$ 
10:  end for
11:  if episode mod update_frequency == 0 then
12:    Update agent networks using replay buffer
13:  end if
14: end for

```

3.8 Network Update Procedure:

We use a multi-objective loss function. The agent networks are trained using a composite loss function:

$$\mathcal{L}_{total} = \mathcal{L}_{dynamics} + \mathcal{L}_{reward}$$

where Dynamics Loss and Reward Loss are Given below:

$$\mathcal{L}_{dynamics} = \frac{1}{B} \sum_{i=1}^B \|f_h(s_i, a_i) - s'_i\|_2^2$$

$$\mathcal{L}_{reward} = \frac{1}{B} \sum_{i=1}^B (f_r(s_i) - r_i)^2$$

3.9 Network Overview:

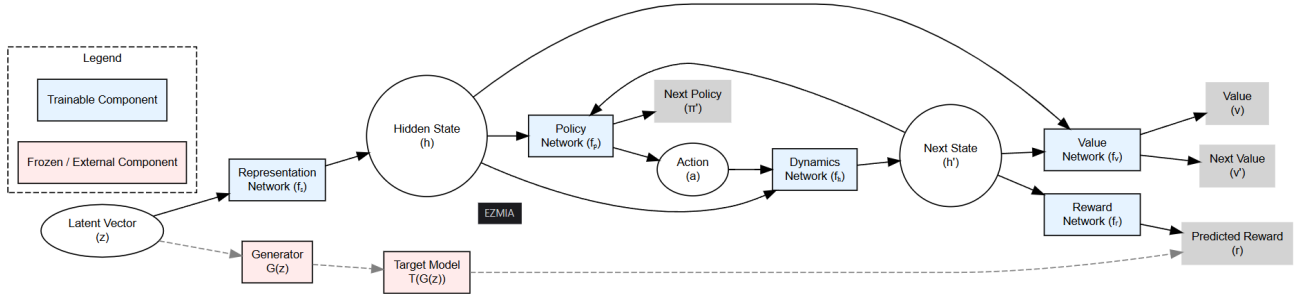


Figure 3.1: Network Workflow

Upper Figure 3.1 illustrates the architecture of our Model Inversion Attack , which formulates black-box model inversion as a Markov Decision Process (MDP) in the latent space of a pretrained generator. The process begins by sampling a latent vector z , which is transformed into a hidden state using a representation network f_s . This state is then used by the policy network f_p to generate an action a , and the dynamics network f_d predicts the next hidden state h' . Alongside this, a reward network f_r estimates the reward (classification confidence), and a value network f_v predicts the value of the state. To choose actions effectively, the agent generates multiple candidate actions and

selects the one yielding the highest predicted reward. The latent vector is updated as a weighted combination of the original latent and the action vector. Importantly, only a small fraction of the interactions (e.g., 10%) query the black-box target model T by passing the generated image $G(z')$ to obtain the true classification probability for the desired label. This true reward is used to evaluate performance and populate the replay buffer. For most steps, the model relies on predicted rewards to guide exploration. The replay buffer is then used to train the dynamics and reward networks using supervised losses, enabling the agent to better model the latent transition dynamics and reward estimation over time. By casting the attack as an MDP and leveraging reinforcement learning components, the approach enables query-efficient and target-aware exploration of the latent space for high-fidelity inversion.

Chapter 4

Experiments

4.1 Experimental Setting& Implementation details:

We evaluate our attack against VGG16 and ResNet-152 (i.e These are the target models).

Datasets: We trained Target models on CelebFaces Attributes (CelebA) dataset which holds 10,177 identities and total 2,02,599 images, out of which we choose 1000 identities and 30,027 images belong to those identities. In addition, we trained our GAN on 30000 (randomly) images out of 9177 identities.

Models: We trained both target models VGG16 and ResNet-152 models on the CelebA dataset, achieved validation accuracies of 60.92% and 62%, respectively. The models were trained for approximately 75 and 71 epochs, respectively. Additionally, we trained the GAN for 300 epochs using an RTX 3050 GPU. We trained FACENET-112 on same target model’s train data in 30 epoch, we achieved validation accuracies of 64.77%.

Model architecture: We take $z_{dim} = 100$ (latent vector size) $h_{dim} = 256$ (internal state size) $a_{dim} = 256$ (action vector size).

Our proposed model architecture consists of five primary components: a representation network, a dynamics network, a reward predictor, a policy network, and a value network. Together, these modules implement a latent-space Markov Decision Process (MDP) formulation inspired by EfficientZero, as we described earlier. The representation network f_s maps the latent vector $z \in R^{100}$ into an internal hidden state $h \in R^{256}$ using a two-layer multilayer perceptron with LayerNorm and ReLU activations. The dynamics function f_d takes the concatenation of the current state h and an action vector $a \in R^{256}$,

and outputs the next state $h' \in R^{256}$. A separate reward network f_r , composed of two fully connected layers with a sigmoid activation at the output, estimates the scalar reward $r \in [0, 1]$, which approximates the target model’s confidence for a given class. The policy network f_p is a linear projection that maps the hidden state h to a continuous action vector $a \in R^{256}$. The value network f_v outputs a scalar estimate of the expected return from state h , using a single linear layer followed by a hyperbolic tangent activation.

Candidate actions are sampled from a Gaussian prior, and the action yielding the maximum predicted reward is selected for exploration. During inference, the agent performs a single initial inference to generate a starting state and then applies recurrent inference through the dynamics model. Latent vectors are updated via a linear interpolation: $z' = \alpha z + (1 - \alpha)a$, where $\alpha \in [0, 1]$ controls the update strength. The generator G maps the updated latent vector $z' \in R^{100}$ to an image, which is passed to the frozen target model T to obtain the softmax confidence for the ground truth label. A small fraction of these interactions query the target model directly to obtain the true reward, which is then stored along with the transition tuple (h, a, r, h') in the replay buffer. The dynamics and reward networks are updated using mean squared error losses computed over mini-batches sampled from this buffer. This architectural design enables latent-space planning with minimal queries by leveraging learned predictive models within an MDP framework.

Replay Buffer Structure:The replay buffer stores experience tuples (s_t, a_t, r_t, s_{t+1}) with the following specifications:

Capacity: 100,000 transitions using circular buffer implementation Sampling Strategy:

Uniform random sampling for batch construction.

Batch Construction Training batches are constructed by sampling $B = 256$ transitions from the replay buffer: $\mathcal{B} = \{(s_i, a_i, r_i, s'_i)\}_{i=1}^B \sim Uniform(ReplayBuffer)$

Optimizer Config:We use adam with Learning Rate: $\alpha = 10^{-4}$ with no scheduling

Hyperparameters:

Adaptive Alpha Values: $\alpha_{explore} = 0.1$, $\alpha_{exploit} = 0.3$

Sampling Frequency: $p_{eval} = 0.1$ (i.e., we query the target model only 10% of the time.)

Multi-Candidate Count: $K = 10$

4.2 Evaluative metrics:

1)**KNN Distances:**KNN Distance quantifies similarity between reconstructed and real images in feature space. It calculates the L distance to the nearest neighbor of the target class. Features are extracted from the evaluation classifier’s penultimate layer. Smaller values indicate closer resemblance to genuine samples. This metric helps evaluate attack stealthiness.

2)**Feature Distance:**Feature Distance assesses reconstruction fidelity by comparing images to class centroids. It computes the L^2 distance between a fake image’s features and the target class’s average features. Like KNN Distance, it uses pre-FC-layer activations. Lower distances suggest better preservation of class characteristics. This complements KNN Distance for robust attack analysis.

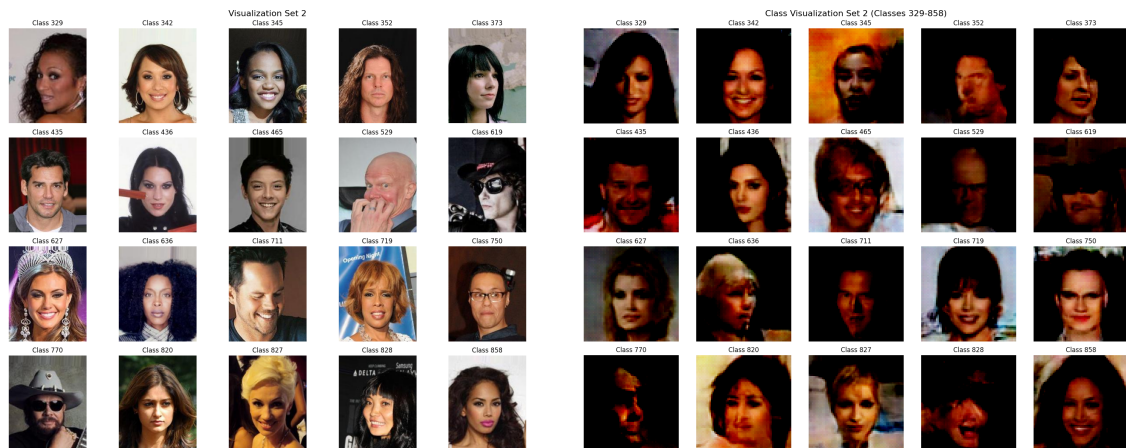
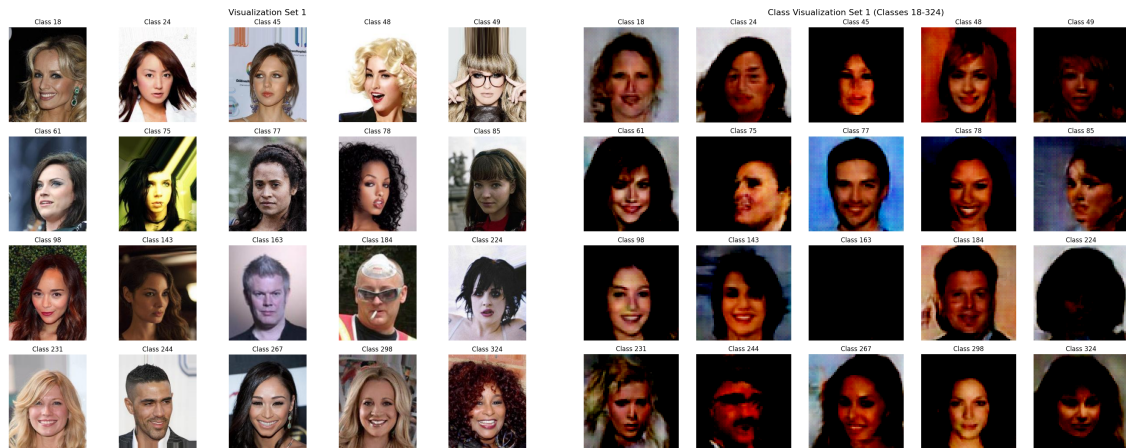
3)**Attack accuracy:**Attack accuracy:Attack Accuracy measures how well reconstructed images deceive machine learning models. It is evaluated using separate classifiers trained on private datasets to avoid overfitting.

4.3 Results

In GMI[15], which is a white-box attack, the authors reported an attack accuracy of 0.185 on VGG16, where the model’s classification accuracy was 88% (as shown in Table 1 of the RLB-MI[19] comparison). In contrast, we achieved a comparable attack accuracy of approximately 0.1 on a VGG16 model with only 60.92% classification accuracy, using just 400 queries—significantly fewer than the 40,000 queries used in GMI—and with much lower computational cost, making our approach faster than both GMI and RLB-MI. Moreover, our method is well-suited for real-world scenarios, where API-based access imposes strict query limits and simulates more practical and constrained attack environments. As noted by S. Yeom et al.[32], increased overfitting tends to raise privacy risks, which may partially explain the higher attack accuracy observed in the GMI and RLB-MI settings. Since their models were well-trained, it is likely that the classification accuracy on the training data exceeded 90%, indicating potential overfitting that could have contributed to increased vulnerability to model inversion attacks.

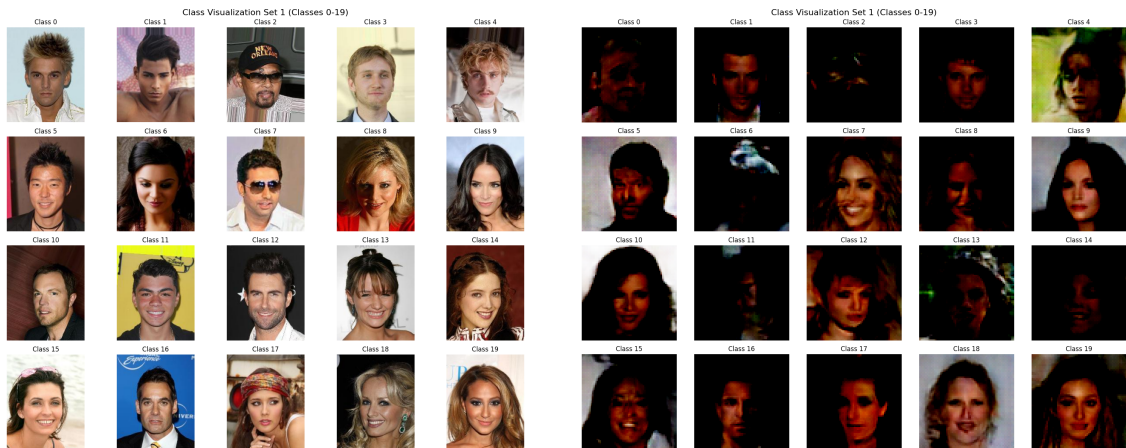
4.4 Randomly Selected Reconstructed Images:

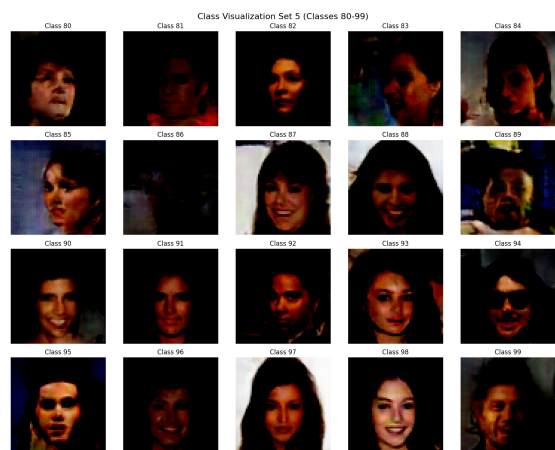
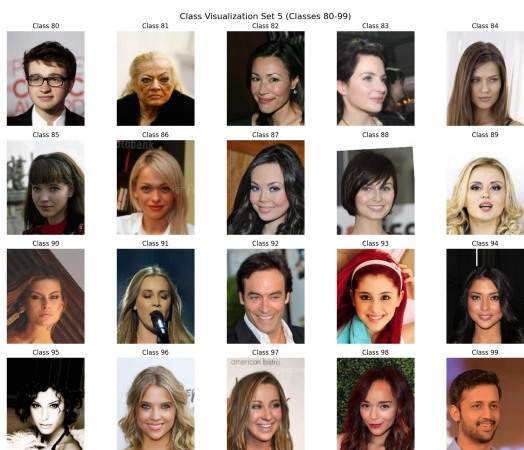
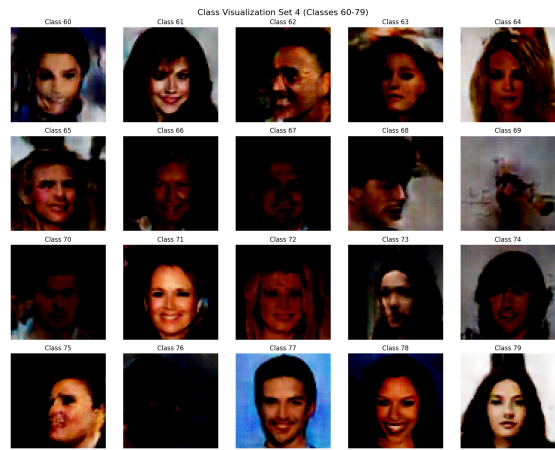
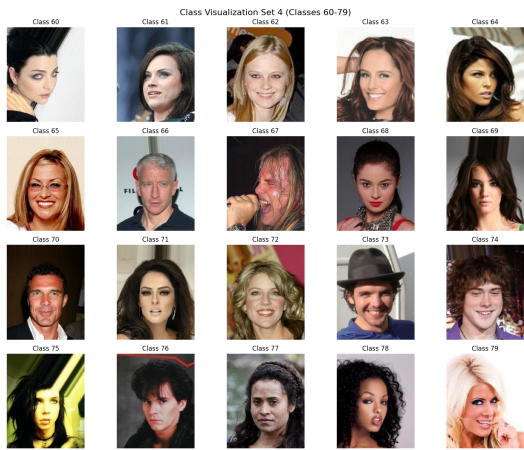
We Reconstructed all 1000 classes in 40000 epochs with 400 quaries to the target model. We selected nearly 50 images randomly and matched with the original images. All images on the right hand side was reconstructed by our proposed method.





4.5 Reconstructed Images from the Initial Classes:





Chapter 5

Conclusion

5.1 Summary:

In this work, we proposed a reinforcement learning-based latent space navigation framework for model inversion attacks, inspired by the EfficientZero architecture. By formulating the search over a generator’s latent space as a sequential decision-making process, our method overcomes key limitations of gradient-based approaches such as local minima and gradient vanishing.

We implemented all core components—representation, dynamics, reward, value, and policy networks—and trained them jointly using online interaction with the target model. Our design incorporates adaptive α -blending for latent updates, multi-candidate action selection, and stochastic reward sampling, enabling efficient and scalable black-box attacks. Empirical evaluation demonstrates that the approach is capable of producing high-confidence reconstructions for target classes.

5.2 Future Works

MCTS with Gumbel Top-k Sampling: Our current method achieves very fast using value-guided latent space navigation. As future work, integrating this architecture with as MCTS along Gumbel Top-k sampling (for continuous action spaces) may further enhance trajectory quality and long-term decision making in more complex inversion tasks.

Diversity-Promoting Exploration: Introducing a diversity buffer or novelty-based reward component may help discover a broader range of high-reward latent codes. This can reduce redundancy in exploration and encourage the agent to cover distinct semantic regions of the latent space.

Temporal Noise Scheduling: Dynamically adjusting the variance of candidate action sampling over time may enhance the exploration-exploitation balance, with high exploration in early stages and more focused refinement as training progresses.

Multi-Objective Reward Integration: Building on the existing reward structure, future extensions may incorporate additional objectives such as confidence maximization and diversity encouragement. This can be combined with adaptive weighting or normalization techniques to stabilize learning across varying tasks and models.

Bibliography

- [1] S. Laine, T. Karras, T. Aila, A. Herva, S. Saito, R. Yu, H. Li, and J. Lehtinen, “Production-level facial performance capture using deep convolutional neural networks,” in *ACM SIGGRAPH/Eurographics Symposium on Computer Animation (SCA)*, 2017.
- [2] V. Kėpuska and G. Bohouta, “Next-generation of virtual personal assistants (Microsoft Cortana, Apple Siri, Amazon Alexa and Google Home),” in *IEEE CCWC*, pp. 99–103, 2018.
- [3] G. A. Kaissis, M. R. Makowski, D. Rėckert, and R. F. Braren, “Secure, privacy-preserving and federated machine learning in medical imaging,” *Nature Machine Intelligence*, pp. 1–7, 2020.
- [4] S. Singh, H. Sikka, S. Kotti, and A. Trask, “Benchmarking differentially private residual networks for medical imagery,” *arXiv preprint arXiv:2005.13099*, 2020.
- [5] R. Shokri, M. Stronati, and V. Shmatikov, “Membership inference attacks against machine learning models,” *arXiv preprint arXiv:1610.05820*, 2016.
- [6] G. Ateniese et al., “Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers,” *arXiv preprint arXiv:1306.4447*, 2013.
- [7] K. Ganju, Q. Wang, W. Yang, C. A. Gunter, and N. Borisov, “Property inference attacks on fully connected neural networks using permutation invariant representations,” in *ACM CCS*, pp. 619–633, 2018.
- [8] M. Yan, C. W. Fletcher, and J. Torrellas, “Cache telepathy: Leveraging shared resource attacks to learn DNN architectures,” *arXiv preprint arXiv:1808.04761*, 2018.

- [9] F. Tramèr et al., “Stealing machine learning models via prediction APIs,” *arXiv preprint arXiv:1609.02943*, 2016.
- [10] B. Wang and N. Z. Gong, “Stealing hyperparameters in machine learning,” in *IEEE S&P*, pp. 36–52, 2018.
- [11] M. Fredrikson et al., “Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing,” in *USENIX Security Symposium*, pp. 17–32, 2014.
- [12] M. Fredrikson, S. Jha, and T. Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS ’15)*, (New York, NY, USA), pp. 1322–1333, Association for Computing Machinery, 2015.
- [13] S. Chen, M. Kahla, R. Jia, and G.-J. Qi, “Knowledge-enriched distributional model inversion attacks,” in *ICCV*, pp. 16178–16187, 2021.
- [14] K.-C. Wang et al., “Variational model inversion attacks,” in *NeurIPS*, vol. 34, pp. 9706–9719, 2021.
- [15] Y. Zhang et al., “The secret revealer: Generative model inversion attacks against deep neural networks,” in *CVPR*, 2020.
- [16] I. Goodfellow et al., “Generative adversarial nets,” in *NeurIPS*, vol. 27, 2014.
- [17] S. An et al., “Mirror: Model inversion for deep learning network with high fidelity,” in *NDSS*, 2022.
- [18] M. Kahla, S. Chen, H. A. Just, and R. Jia, “Label-only model inversion attacks via boundary repulsion,” in *CVPR*, pp. 15045–15053, 2022.
- [19] G. Han et al., “Reinforcement learning-based black-box model inversion attacks,” in *CVPR*, pp. 20504–20513, 2023.
- [20] W. Ye et al., “Mastering Atari games with limited data,” in *NeurIPS*, vol. 34, pp. 25476–25488, 2021.
- [21] Z. Yang, J. Zhang, E.-C. Chang, and Z. Liang, “Neural network inversion in adversarial setting via background knowledge alignment,” in *ACM CCS*, pp. 225–240, 2019.

- [22] A. Salem et al., “Updates-Leak: Dataset inference and reconstruction attacks in online learning,” in *USENIX Security*, pp. 1291–1308, 2020.
- [23] M. G. Bellemare et al., “The arcade learning environment: An evaluation platform for general agents,” *Journal of Artificial Intelligence Research*, vol. 47, pp. 253–279, 2013.
- [24] V. Mnih et al., “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [25] T. P. Lillicrap et al., “Continuous control with deep reinforcement learning,” in *ICLR*, 2016.
- [26] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, “Deterministic policy gradient algorithms,” in *Proceedings of the 31st International Conference on Machine Learning (ICML)*, E. P. Xing and T. Jebara (eds.), vol. 32 of *Proceedings of Machine Learning Research*, pp. 387–395, Beijing, China, 22–24 Jun 2014. PMLR.
- [27] S. Fujimoto, H. van Hoof, and D. Meger, “Addressing function approximation error in actor-critic methods,” in *ICML*, vol. 80, pp. 1587–1596, 2018.
- [28] T. Haarnoja et al., “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” in *ICML*, vol. 80, pp. 1861–1870, 2018.
- [29] Y. Chen et al., “Rethinking soft actor-critic in high-dimensional action spaces: The cost of ignoring distribution shift,” 2024.
- [30] W. Ye, S. Liu, T. Kurutach, P. Abbeel, and Y. Gao, “Mastering Atari games with limited data,” in *Advances in Neural Information Processing Systems*, vol. 34, pp. 25476–25488, 2021.
- [31] S. Wang et al., “EfficientZero V2: Mastering discrete and continuous control with limited data,” 2024.
- [32] S. Yeom, I. Giacomelli, M. Fredrikson, and S. Jha, “Privacy risk in machine learning: Analyzing the connection to overfitting,” in *Proceedings of the IEEE Computer Security Foundations Symposium (CSF)*, pp. 268–282, 2018.