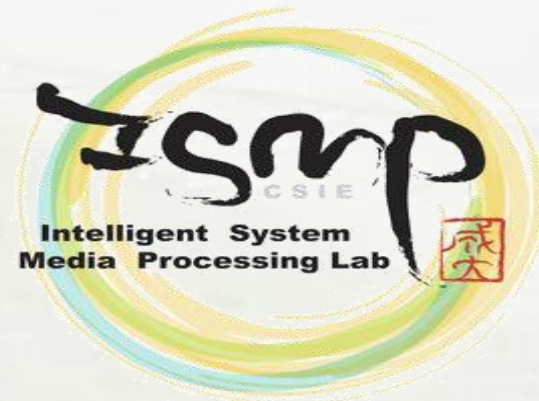# AI安全 (三)
# AI Security (3)

計算機與網路中心
莊宜勳

# Attacks: Stealing Models

# Understanding Privacy Attacks

## Model Extraction

Replicating an AI model's functionality by observing its responses to various inputs. The primary risk is potential loss of intellectual property and unauthorized duplication of proprietary AI models.

## Model Inversion

Extracting sensitive data used in training AI models. Most concerning in sectors like healthcare or finance, where privacy breaches can have serious consequences.

## Membership Inference

Determining if a particular data record was part of the training dataset. Critical risk involves privacy violations, especially when training data contains sensitive personal information.

Intelligent System
Media Processing Lab

# Model Extraction Attacks：Overview

- Model extraction attacks replicate the functionality of machine learning models by observing and mimicking their output responses to various inputs.

  - Using an iterative query-based approach, attackers repeatedly query the target AI model with carefully selected inputs.
  - With each query, the attacker receives predictions or confidence scores to refine and train a new model.
  - Over time, the extraction model becomes increasingly similar to the original target model.

# Approaches to Model Extraction

## Functionally Equivalent Extraction

Aims to achieve maximum fidelity by calculating weights and biases. Attempts to create an identical copy of the model by solving complex mathematical equations.

## Learning-Based Methods

Black-box attacks that aim to find representative data capturing the target model's decision boundaries, then training an extracted model to approximate the target.

## Generative Student-Teacher Learning

Uses generative techniques to create adversarial samples for querying the victim model, applying knowledge distillation principles to transfer knowledge.

# Functionally Equivalent Extraction

- This attack attempts to calculate the exact weights and biases of the victim model.
  - make victim model and cloned model have approximately the same output on the same inputs.
  - Calculating exact weights creates a complex non-linear optimization equation which is an NP-complete problem.
    - The approach focuses on relatively simple ANN architectures with assumptions about the architecture itself.

## Observe Model Outputs

Query the target with selected inputs and log raw output values (logits)

## Partition Inputs

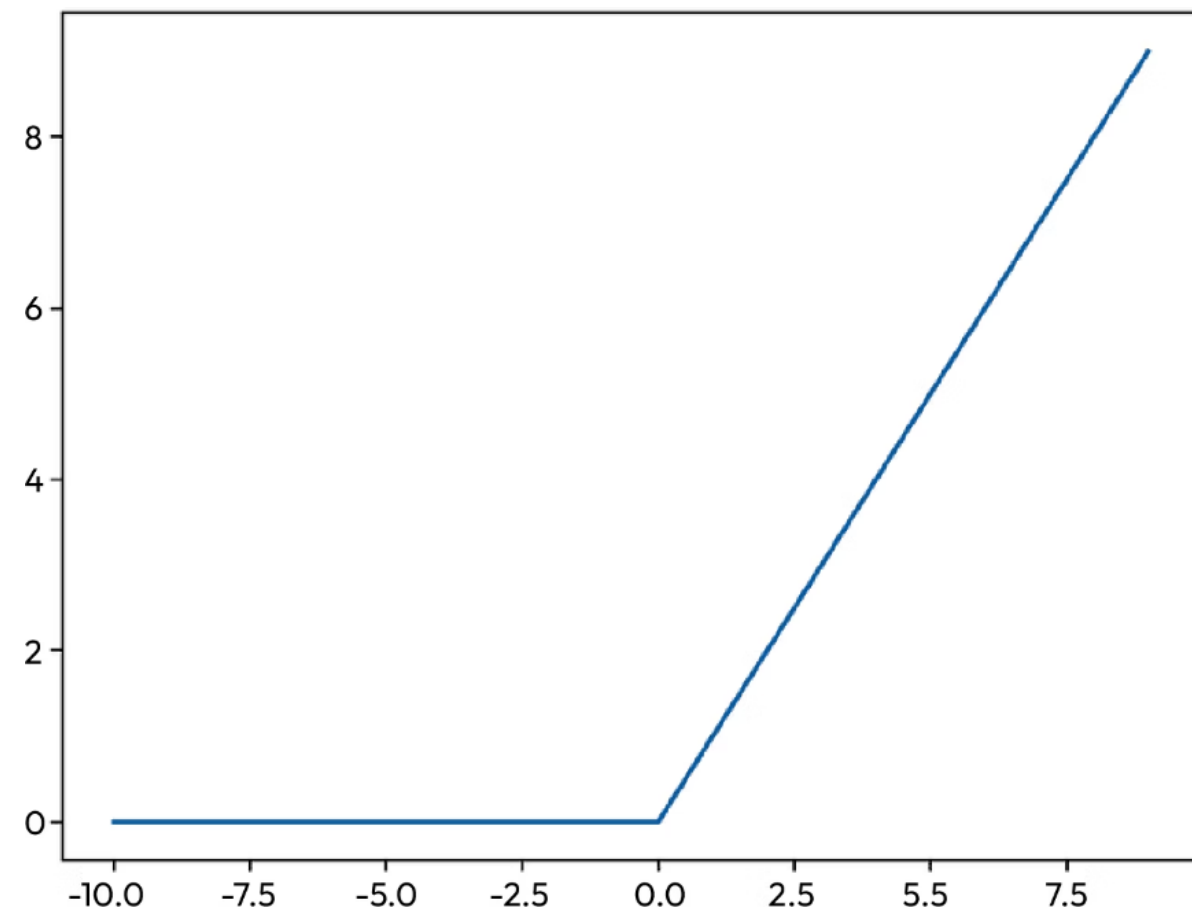Segment inputs into regions where ReLU units have the same sign

## Calculate Parameters

Use linear equations and additional queries to determine weights and biases

# ReLU Activation Function

- The ReLU activation function is key to functionally equivalent extraction attacks.
  - It outputs the input directly if positive; otherwise, it outputs zero
    - creating a clear sign change in the function's outputs based on user input.
  - This property allows attackers to segment inputs and outputs based on their effect on ReLU, helping to reverse-engineer the model's parameters by observing how inputs cross ReLU boundaries.

# Challenges in Functionally Equivalent Extraction

### Computational Complexity

As model architecture grows in complexity, the computational requirements grow exponentially, making extraction increasingly difficult
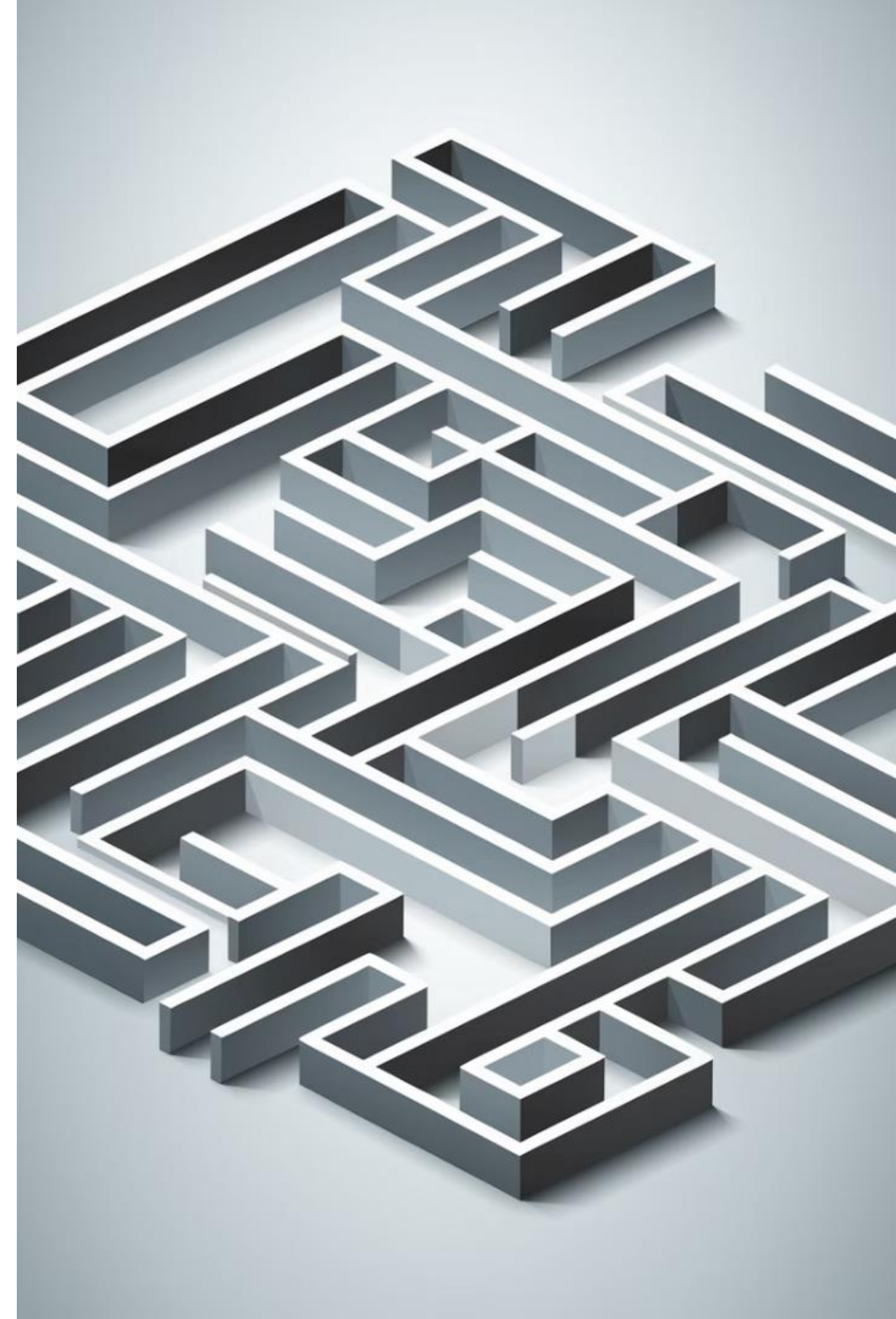
### Constrained Search

Attackers may use environment parameters like time-side channels to infer network depth and constrain the search space

### Architecture Limitations

Most successful attacks target simple architectures with known activation functions and layer structures
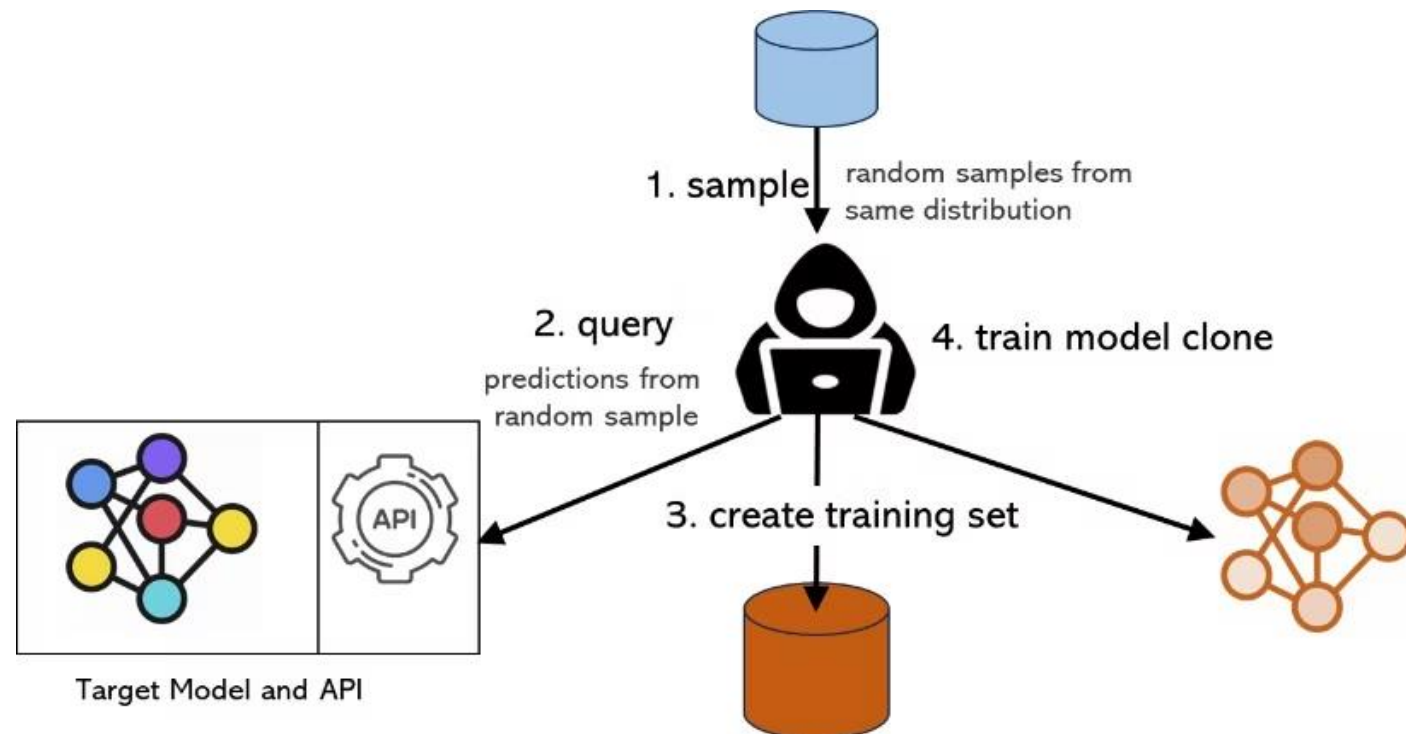
# Learning-Based Model Extraction Attacks

Learning-based approaches are black-box attacks that aim to find random data capturing the target model's decision boundaries. The attacker queries the model with randomly sampled items and uses the predictions to formulate a training set for the extracted model.



Key Steps

1. Generate or collect unlabeled data samples
2. Query target model with these samples
3. Record model predictions as labels
4. Train surrogate model using this synthetic dataset
5. Refine surrogate model until performance matches target

# Copycat CNN

- Copycat CNN is a well-known technique that trains an extracted neural network to replicate the behavior of the target model.
  - The Copycat CNN attack has been used successfully against Microsoft's Azure Emotion API in research settings, extracting 97.3% of the API's performance.
    - This demonstrates the effectiveness of this approach against commercial AI services.
  - Reconnaissance is essential for this attack to constrain the generation of fake datasets to better approximate decision boundaries.
    - If the target model is based on a publicly available model (e.g., ResNet50), the attacker can use the public model and fine-tune it with API query results.

| Fake Dataset Generation | Dataset Balancing | Copycat Network Training | Performance Verification |
|---|---|---|---|
| Query target network with random unlabeled data from the same domain and collect predictions as labels | Achieve image balance per class through random replication or removal | Train a copycat network with the same architecture using the fake dataset | Evaluate copycat network against target performance |

# KnockOff Nets

KnockOff Nets are similar to Copycat CNN but follow a more sophisticated approach to data sampling and model training. The key difference is the lack of assumptions about the target model and more nuanced techniques for refining the training process.

## Training Set Generation

Select images from chosen distribution (can be different from target domain)

## Sampling Strategy

Use random sampling or adaptive sampling with reinforcement learning

## Train Knockoff Model

Use image-prediction pairs to train model with same task and output space

## Query Victim Model

Send batch of images to victim model and obtain predictions

# KnockOff Nets Advantages

### 💬 Cross-Domain Transfer

Can use images from different distributions than the target model's domain

### 🤖 Reinforcement Learning

Can employ RL agents to select the most informative images to query, improving efficiency

### Architecture Freedom

Attacker can use any architecture for the knockoff model, not limited to matching the victim

### ⊠ Proven Effectiveness

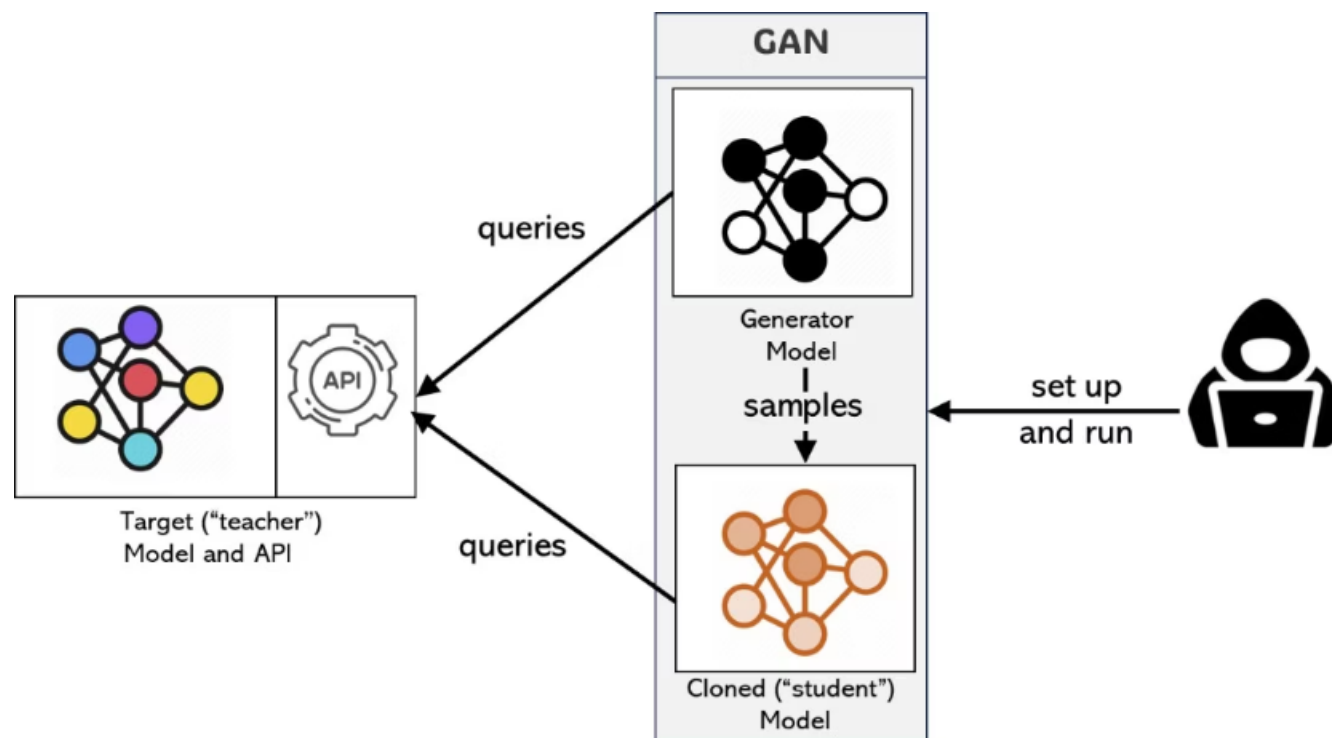Successfully demonstrated against commercial facial analysis systems

# Generative Student-Teacher Extraction

- This approach creates adversarial samples for querying the victim model. Based on knowledge distillation, a smaller, simpler model (student) is trained to replicate a larger, more complex model (teacher).

    - Unlike legitimate knowledge distillation where access to training data is available, attackers work in a black-box scenario without access to the original training dataset. They use generative models to create adversarial data for queries.
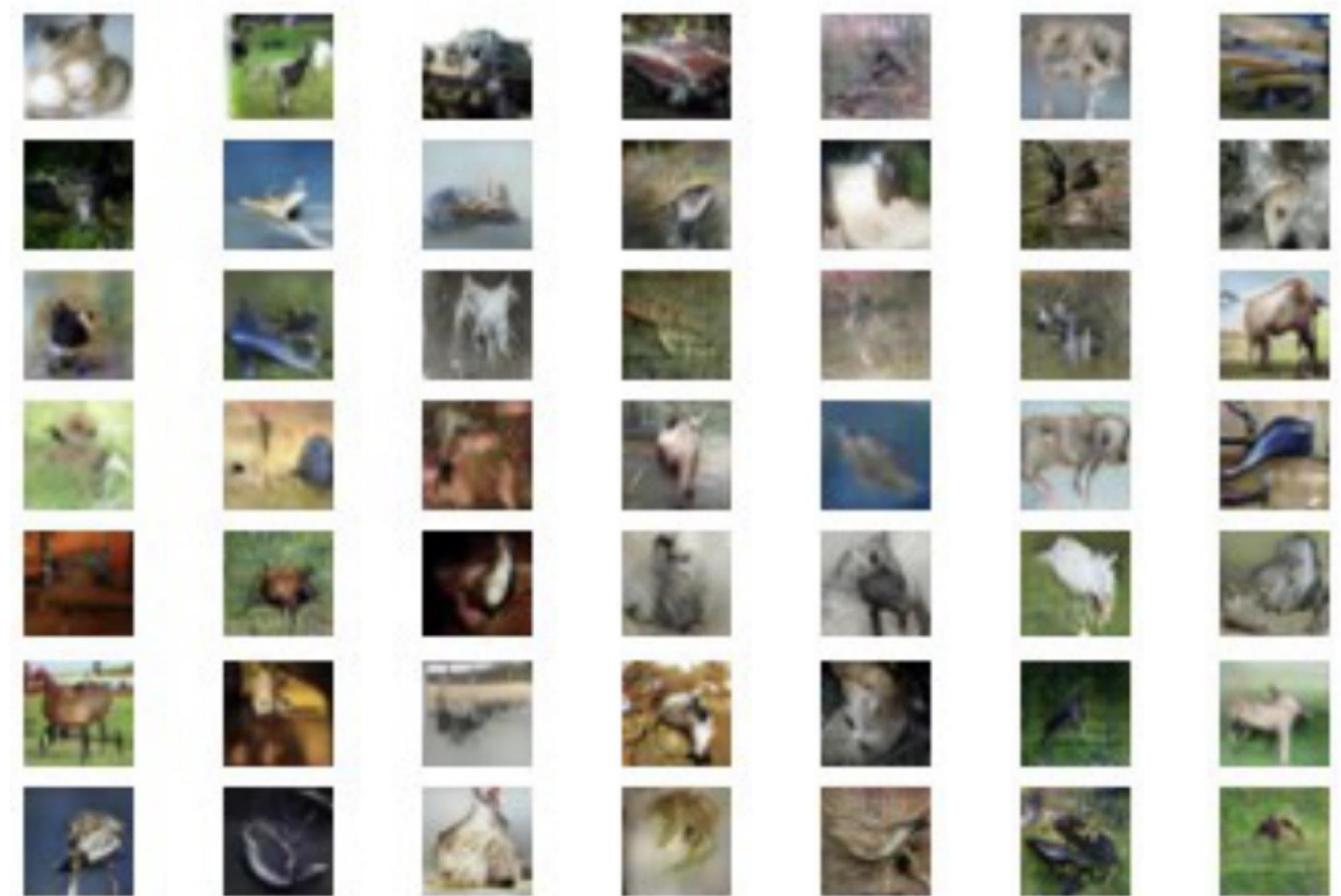
### Adversarial Game Process

1. Generator creates challenging data samples
2. Samples maximize difference between teacher and student models
3. Student model trained to minimize this divergence
4. Process repeats iteratively in competitive optimization

# GAN-Generated Images

When the generative extraction attack starts, the GAN's images are pseudo-images generated from noise. These synthetic images evolve to become increasingly effective at probing the target model's decision boundaries.

# Generative Extraction Variations
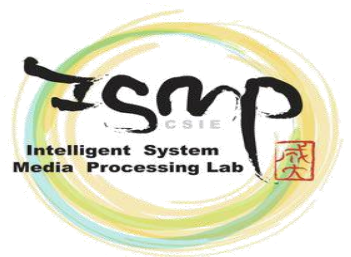
## Data-Free Model Extraction

Generator inputs random noise and outputs images resembling original data distribution. Uses 1-norm loss between victim and student models to measure disagreement between their logits.

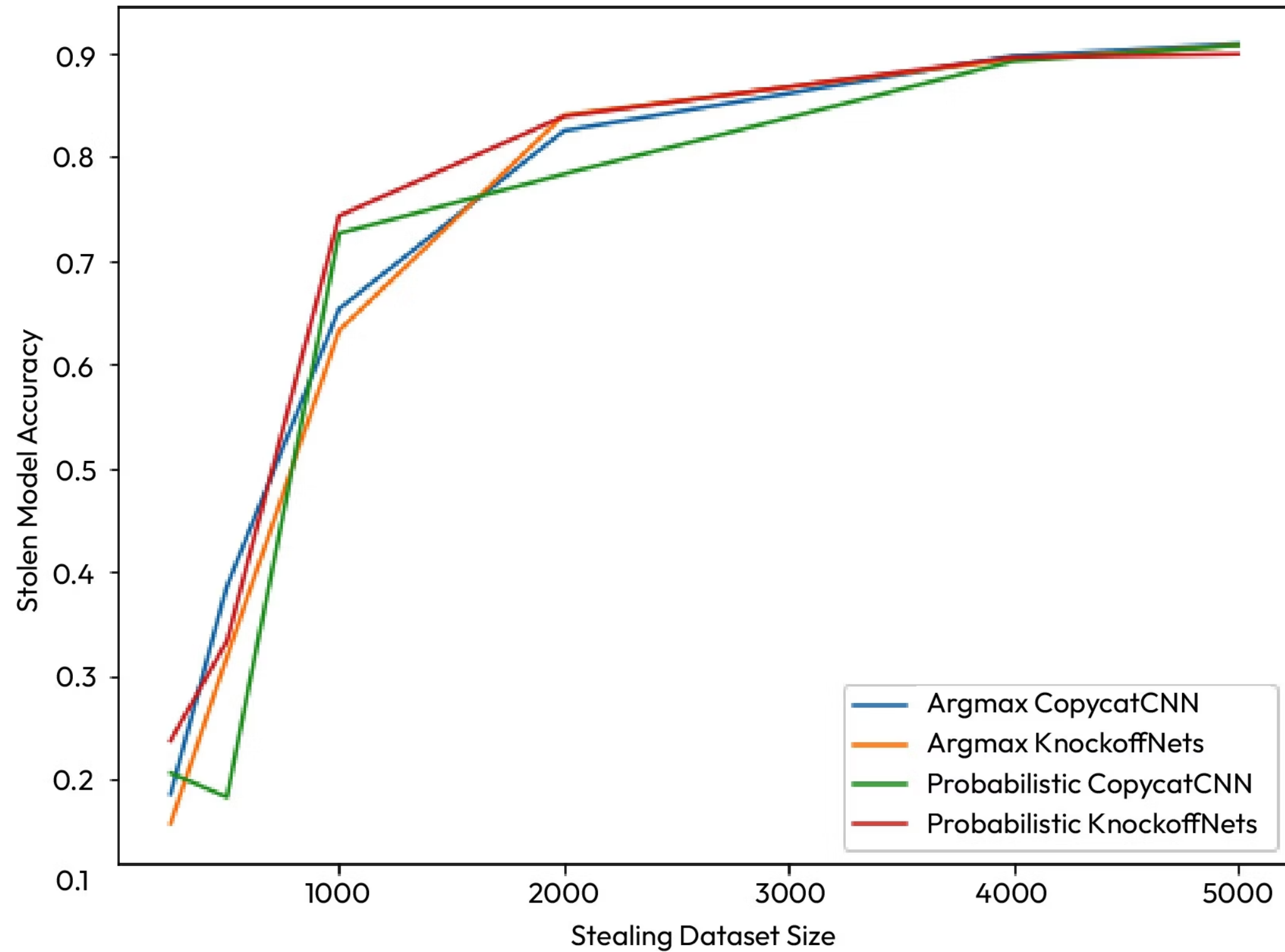## Zero-shot Knowledge Transfer

Uses forward Kullback–Leibler divergence loss function to measure differences between student and teacher networks. Adds attention transfer to focus on important differences.

## TandemGAN

Uses two neural networks working in tandem - an exploration network creates codes representing different areas, which the exploitation network translates into synthetic queries.

# Implications of Model Extraction

Model extraction attacks have serious implications for organizations that develop and deploy AI models, including both business and security concerns.

## Intellectual Property Theft

Organizations invest significant resources in developing proprietary AI models. Extraction attacks can lead to unauthorized replication, undermining this investment and competitive advantage.
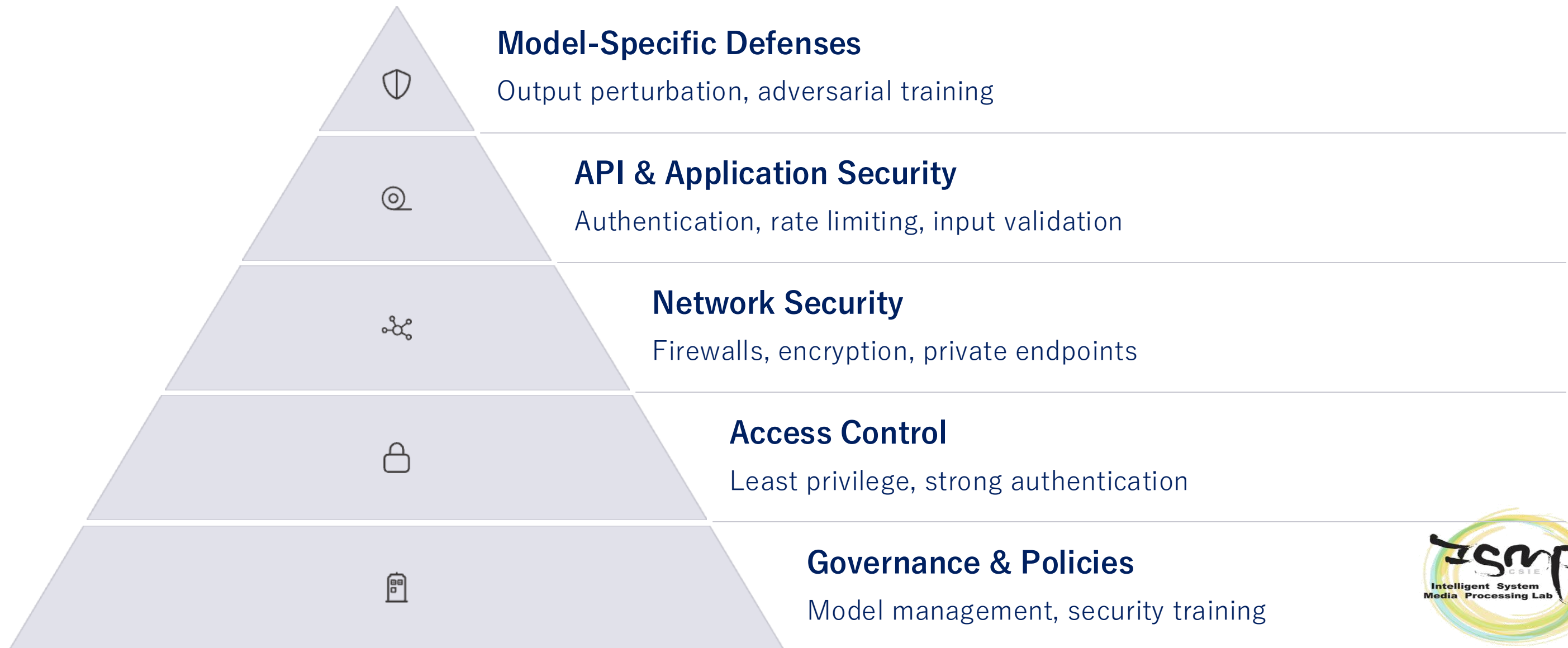
## Financial Losses

When models are offered as paid services, extraction allows attackers to avoid paying for continued use, directly impacting revenue streams.

## Security Vulnerabilities

Extracted models can be studied to identify weaknesses and develop further attacks, such as evasion attacks that bypass security measures.

Intelligent System
Media Processing Lab

# Defense-in-Depth Philosophy

Protecting against model extraction attacks requires a defense-in-depth approach that addresses multiple layers of security. This strategy recognizes that no single defense is perfect and combines various measures to create comprehensive protection.

**Model-Specific Defenses**

Output perturbation, adversarial training

**API & Application Security**

Authentication, rate limiting, input validation

**Network Security**

Firewalls, encryption, private endpoints

**Access Control**

Least privilege, strong authentication

**Governance & Policies**

Model management, security training

# Prevention: Combat Reconnaissance

Despite advances in black-box attacks, attackers typically start by collecting information about the target model. Knowledge of the underlying architecture helps them choose matching architecture for the surrogate model.

## Confidentiality Policies

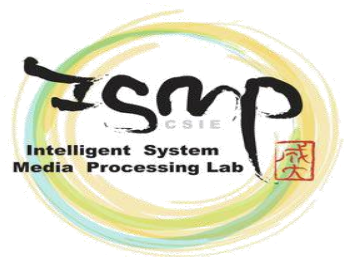Implement strict information control for mission-critical applications to limit architectural details

## Restrict Output Information

Use argmax model output instead of probabilities to prevent inferring internal workings

## Error Message Control

Test and limit what information error messages reveal about model implementation

## Reconnaissance Detection

Monitor for probing patterns that suggest attackers gathering information

# Prevention: System Security

Strong system security forms the foundation of model protection, preventing both direct access to models and limiting the information available through APIs.

## Strict Model Governance with MLOps

Essential to prevent white-box attacks from compromised insiders. Implement version control, access logging, and approval workflows for model changes.

## Least-Privilege Access

A cornerstone for production systems that helps avoid white-box attacks and blocks successful intruders from pivoting further into the system.

## Gated API Pattern

Isolate and segment the inference API to minimize access. Restrict access to the application that uses the API instead of making it publicly available.

# Prevention: API Security

Since model extraction attacks primarily target APIs, implementing strong API security measures is crucial for preventing unauthorized access and limiting query volume.

## Strong Authentication

Implement robust system-to-system authentication for API calls and strong user authentication for applications

## Session Management

Use inactivity-based logouts and session timeouts to limit extended access

## Input Pre-processing

Alter or preprocess inputs before feeding them to the model to make output interpretation more difficult

## Rate Limiting

Restrict the number of queries from a single source to prevent mass data collection

# Prevention: Input Pre-processing Example

Adding noise to input data can significantly reduce the effectiveness of extraction attacks while maintaining legitimate model performance. This example shows how to implement Gaussian noise using ART:

```python
from art.defences.preprocessor import GaussianNoise

# Create the noise generator with standard deviation 0.1
noise_generator = GaussianNoise(scale=0.1)

# Apply to the classifier
defended_classifier = KerasClassifier(
    model=model,
    preprocessing_defences=[noise_generator]
)
```

The key is finding the right balance - enough noise to disrupt extraction attempts but not so much that it degrades model performance for legitimate users. This typically requires experimentation with different noise levels.
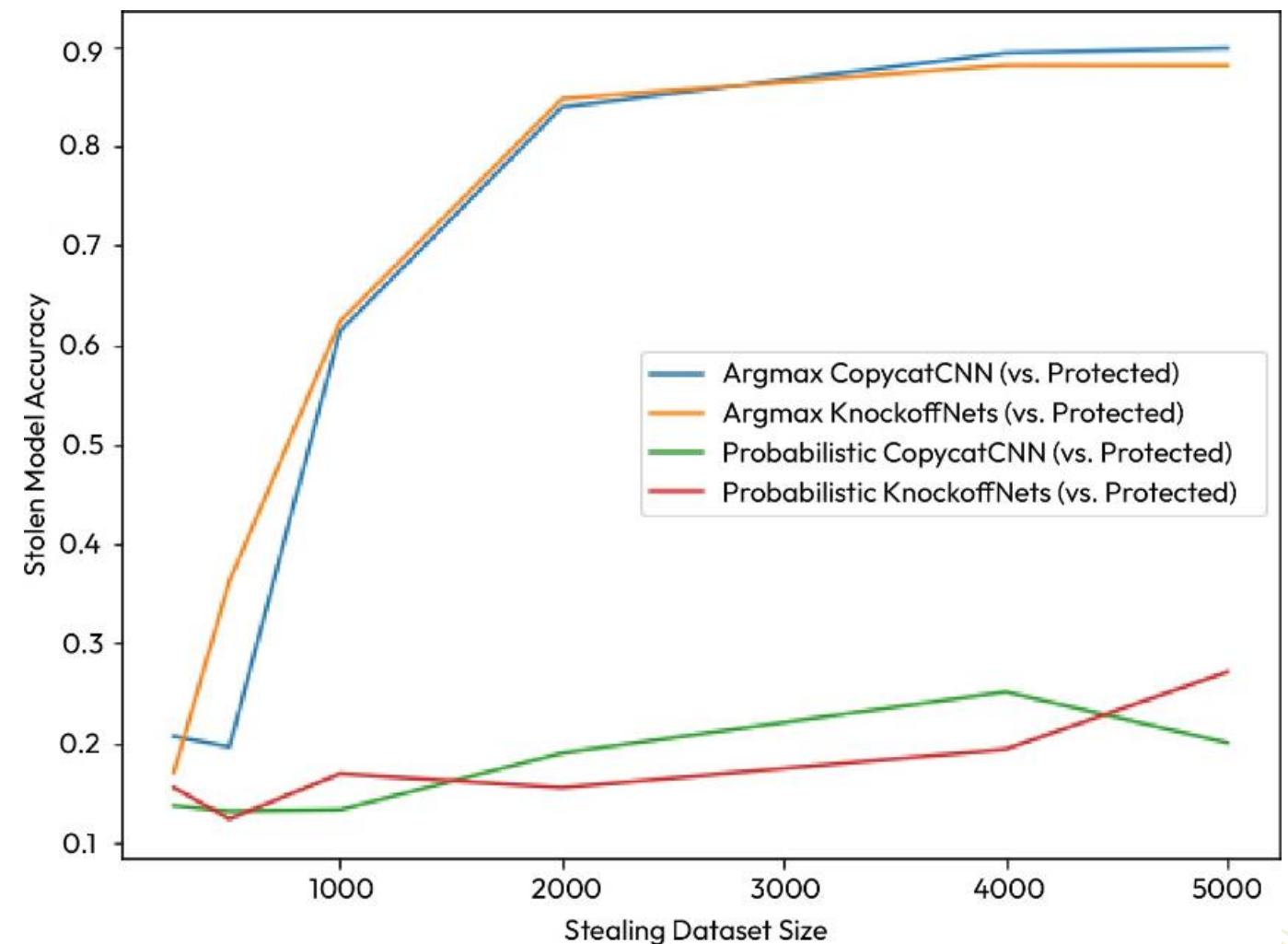
# Prevention: Output Perturbation

Output perturbation adds random noise or distortion to the model's output, making it harder for attackers to train a clone model. This can be implemented using ART's ReverseSigmoid defense:

```
from art.defences.postprocessor import ReverseSigmoid

# Create defense with beta=1.0
rs = ReverseSigmoid(beta=1.0)

# Apply to classifier
defended_classifier = KerasClassifier(
    model=model,
    postprocessing_defences=[rs]
)
```



As shown in the figure, adding ReverseSigmoid defense significantly drops the accuracy of extracted models to unusable levels, providing effective protection against extraction attacks.

# Additional Prevention Techniques

## Adversarial Training

Train the model with adversarial examples to increase robustness against extraction. This can increase computational cost but provides significant protection.

## Gradient-based Ranking Optimization (GRO)

Designed for recommender systems, this technique optimizes ranking lists to minimize model loss while maximizing attacker's surrogate model loss.

## Differential Privacy

Add controlled noise to queries or data to make it difficult to infer information about individual points while preserving overall model performance.

# Detection Measures

While prevention is crucial, no defense is perfect. Detection measures help identify extraction attempts in progress, allowing for timely response.

## Testing Against Known Attacks

Incorporate tests against known extraction attacks using tools like ART. Automate these tests in pipelines to monitor vulnerabilities before models reach production.

## Red-Team Testing

Conduct regular red-team testing of models for in-depth investigation of adversarial robustness against extraction attacks.

## Rate Limiting

Apply rate limiting to applications and APIs to defend against excessive queries by attackers attempting model extraction.

# Advanced Detection Strategies

More sophisticated detection approaches can help identify subtle extraction attempts that might bypass basic rate limiting and other simple controls.

## System Monitoring

Monitor system access and utilization to detect unusual behavior such as excessive queries that work around rate limiting

## Model Query Analysis

Analyze patterns in model queries to identify potential extraction attempts versus legitimate usage

## ML-Based Detection

Use machine learning to detect adversarial samples versus legitimate inputs

## Incident Response

Follow established incident response processes when extraction attempts are detected

# Model Ownership Identification

Beyond prevention and detection, techniques for asserting model ownership can help identify stolen models and provide evidence for legal action. These approaches create verifiable proof of model origin.

## Unique Model Identifiers

Taking a hash or other fingerprint of the model provides protection against tampering and serves as an ownership verification mechanism in white-box scenarios.

## Dataset Inference

Calculating distances in training datasets can help prove a model was trained with a specific dataset, useful for analyzing suspected stolen models deployed as public APIs.

## Conferrable Adversarial Samples

Creating fingerprints by adding adversarial samples during training that transfer to surrogate models but not independently trained models.
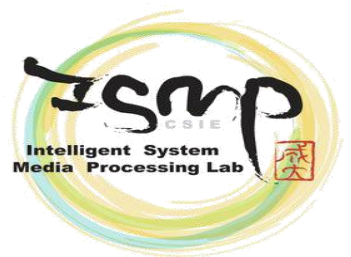
# Watermarking Techniques

Watermarking is a practical way to create verifiable ownership identifiers by modifying the model to embed distinctive patterns. Successful watermarks should be robust, recognizable, and not degrade model performance.

## Parameter Watermarks

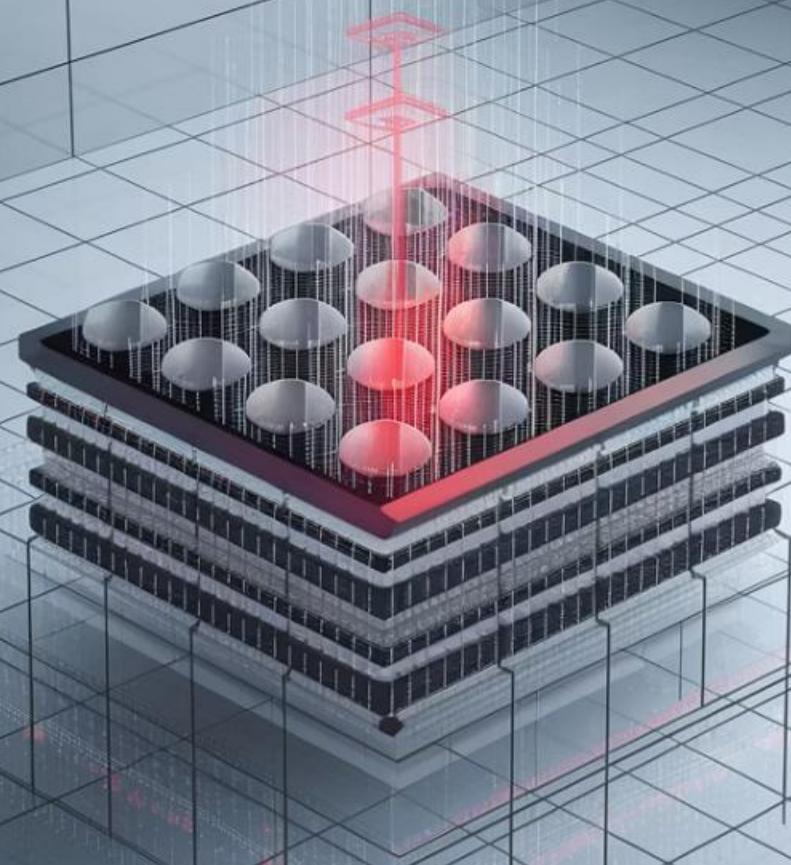This approach adds subtle weight changes, encoding a simple pattern without affecting performance. While simple and effective, it doesn't survive compression, weight pruning, or fine-tuning.

## Training Data Watermarks

This approach teaches the model to produce a specific pattern for certain inputs, acting as a watermark. The challenge is finding benchmarks that don't affect model performance.

# Backdoor-Based Protection

By creating a backdoor for specific triggers that return identifying information about the model, defenders can identify extracted models deployed in the public domain.

**Design Trigger Pattern**

Create a specific input pattern that will activate the backdoor

**Train Model with Backdoor**

Incorporate the trigger-response pair during model training

**Monitor for Copied Models**

Test suspected copied models with the trigger pattern

**Verify Ownership**

If the model responds to the trigger, it confirms it was copied from your original

# Comprehensive Defense Strategy

An effective defense against model extraction attacks requires combining prevention, detection, and recovery strategies in a coordinated approach. This multi-layered defense provides the best protection for valuable AI assets.

## Prevention

Implement technical controls to make extraction difficult

- API security
- Output perturbation
- Input preprocessing

## Response

Actions when extraction is detected

- Incident handling
- Legal action
- Service adjustments

## Detection

Monitor for extraction attempts in progress

- Query monitoring
- Anomaly detection
- Rate limiting alerts

## Identification

Techniques to prove model ownership

- Watermarking
- Fingerprinting
- Backdoor triggers

3

ISMP
Intelligent System
Media Processing Lab

# Case Study: Protecting a Healthcare AI Model

A healthcare company developed an AI diagnostic model trained on sensitive patient data. To protect this valuable intellectual property while making it available to hospitals, they implemented a comprehensive defense strategy.

## Multi-layered Protection

Combined API security, output perturbation, and watermarking to create defense-in-depth
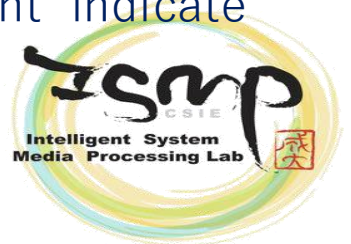
## Usability Preserved

Security measures implemented without compromising diagnostic accuracy or user experience

## Continuous Monitoring

Implemented query analysis to detect unusual patterns that might indicate extraction attempts

# Future Trends in Model Extraction

As AI continues to evolve, both attack and defense techniques for model extraction are advancing rapidly. Organizations must stay informed about emerging threats and countermeasures.

**Current State**

Learning-based and generative approaches dominate extraction attacks, with defenses focusing on perturbation and monitoring

**Emerging Trends**

More sophisticated GAN-based attacks and zero-shot extraction techniques requiring fewer queries

**Future Directions**

Federated learning to keep models decentralized, advanced watermarking that survives transfer learning

**Ongoing Balance**

Continuous evolution of attack and defense techniques creating an arms race in model protection

Intelligent System
Media Processing Lab

# Privacy Attacks: Stealing Data

# Types of Privacy Attacks

## Model Extraction

Stealing the model itself by creating a functionally similar copy through API queries. This targets the intellectual property of the model rather than the data.

## Model Inversion

Reconstructing training data by exploiting model outputs. Attackers attempt to reverse-engineer inputs that would produce specific outputs.

## Inference Attacks

Deducing sensitive information about training data, including whether specific data points were used in training (membership inference) or inferring properties of the data (attribute inference).
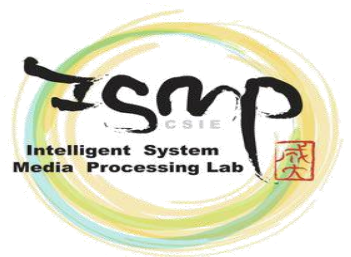
# Understanding Model Inversion Attacks

Model inversion attacks are sophisticated adversarial threats that attempt to reconstruct training data from a trained model. These attacks "invert" the model's prediction process to recover sensitive information used during training.

## White Box Attacks

Attacker has complete access to the model, including architecture, weights, and possibly training data. This allows for precise exploitation of the model's internal workings to reverse engineer training data.

## Black Box Attacks

Attacker only has access to the model's input-output capabilities. Despite limited information, attackers can still infer private data by analyzing outputs from different inputs. These are more common in real-world scenarios.

# Model Inversion Attack Techniques

### Exploitation of Confidence Scores

Uses model's confidence information as feedback to guide the search for inputs that maximize confidence for a given class label.

### GAN-Assisted Inversion

Employs Generative Adversarial Networks to help search the input feature space and constrain the optimization problem.

### Knowledge-Enriched Approaches

Utilizes additional knowledge about the target dataset to improve accuracy of reconstructed data.

### Plug and Play Attacks

Uses pretrained GANs to create synthetic images, significantly reducing preparation time and making attacks more efficient.

# MIFace Attack: Exploiting Confidence Scores

One of the first model inversion attacks was the MIFace attack, demonstrated by researchers at Carnegie-Mellon University and the University of Wisconsin–Madison in 2015. This attack exploits the confidence information returned by models.

## Feed Random Image

A sample image with random pixels is fed to the classifier, which returns prediction with confidence scores (probabilities).

## Measure Cost Function

The attacker's cost function uses the response to measure how well the model matches the target name and image quality.

## Compute Gradient

The cost function's gradient is computed, showing how to change the image to improve the cost function.

## Process and Repeat

A processing function makes the image more realistic and natural, then repeats the steps until cost is below threshold or maximum iterations are reached.

Intelligent System
Media Processing Lab

# GAN-Assisted Model Inversion

GAN-assisted approaches overcome the limitations of single-sample instance attacks by using Generative Adversarial Networks to help search the input feature space and better constrain the optimization problem.

## Public Knowledge Distillation

GAN's generator produces samples, and discriminator guides it to create realistic images using public dataset.

## Image Reconstruction

Final high-fidelity reconstructed images are generated, often with remarkable similarity to original training data.

## Latent Vector Creation

Generator creates reduced image representations capturing essential features like color, shape, and texture.

## Secret Revelation

Attacker uses GAN to find latent vector that maximizes probability of target label using the target network.

(a) W/out auxiliary knowledge
(b) Blurring
(c) Center Mask
(d) Face T Mask

# Advanced Model Inversion Techniques

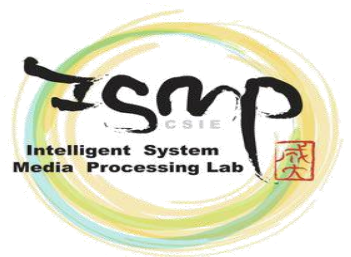## Knowledge-Enriched Distributional Model Inversion (KE-DMI)

A white box attack that creates an inversion-specific GAN by adding soft labels provided by the target model. It attempts to model the model's private data distribution rather than searching for single instances.

This approach requires access to the underlying model but achieves higher accuracy in reconstructing images with better fidelity and similarity.

## Variational Model Inversion (VMI)

Similar to GMI but focuses on increasing the diversity of inferred training data rather than just accuracy of a subset. Uses StyleGAN to create stylistic sample variations by combining changes of multiple images.

VMI can manipulate specific features (like facial characteristics) and combine them to generate more diverse reconstructions.

hairstyle, bangs, face contour      eyes, nose, mouth, hair color

# Plug and Play Attacks

Plug and Play attacks represent a significant optimization in model inversion techniques, making attacks more efficient and accessible.

## Time Efficiency

Uses pretrained GANs (like StyleGAN2) instead of training new ones, reducing sample generation time from hours to minutes.

## Reusability

Single pretrained GAN with transformations can be used across different datasets and models, making the approach independent of target model architecture.

## Simplicity

Eliminates complex preparation steps required by other approaches, lowering the technical barrier for potential attackers.

## Effectiveness

Despite simplifications, achieves comparable or better results than more complex approaches in reconstructing training data.

# Understanding Inference Attacks

- Inference attacks aim to deduce sensitive information from a machine learning model without directly accessing its parameters or training data.
  - Unlike model extraction or inversion attacks, they attempt to infer broader data patterns and relationships.

## Attribute/Property Inference

These attacks aim to infer global information about the training data or model, such as data distribution, architecture, or hyperparameters. For example, determining the average age of individuals in a health prediction model's training dataset.

## Membership Inference

These attacks determine whether a specific data item was part of the model's training dataset. For instance, an attacker could determine if a particular patient's data was used to train a medical diagnosis model.

# Attribute Inference Attacks

Attribute inference attacks target global or individual properties of the training data or model, such as data distribution, labels, architecture, or hyperparameters.

## Meta-Classifier Approach
Train an attack model to infer confidential properties

## Shadow Model Training
Create models mimicking target behavior

## Training Data Generation
Use shadow models to create labeled inputs

## Property Inference
Apply meta-classifier to target model

# Poisoning-Assisted Inference

Poisoning-assisted inference represents an advanced approach to attribute inference attacks, combining poisoning techniques with inference methods to extract sensitive information.

### Poisoning Phase

Attacker injects triggers into training data that leak information during inference

### Query Phase

Black box inference queries exploit injected triggers to extract information

### Analysis Phase

Meta-classifier processes responses to infer sensitive properties

The SNAP (Subpopulation iNference Attack with Poisoning) optimization eliminates meta-classifier use by utilizing model confidences directly. This reduces the number of shadow models needed from thousands to at most 4, with a low poisoning rate of less than 1% of the training dataset.

# Membership Inference Attacks

Membership inference attacks aim to determine whether a specific data point was part of a model's training dataset. These attacks can reveal sensitive information about individuals whose data was used for training, potentially violating privacy regulations.

## Shadow Model Approach

Train multiple shadow models on synthetic datasets mimicking target model's training data. Use shadow models' outputs to train an attack model that predicts membership status of inputs. Exploit confidence levels in target model outputs.

## Statistical Threshold Method

Simplifies shadow model approach by eliminating shadow model training. Uses reference dataset to generate queries against target model and statistical evaluation of confidence thresholds. Significant confidence scores indicate membership.

## Label-Only Attacks

Assumes attacker only has access to predicted labels, not confidence scores. Uses meta-classifier trained on outputs of multiple models with different architectures and applies it to target model's output to infer membership status.

## Blind Membership Inference

Uses labels only and differential comparisons to eliminate shadow models. Creates non-members dataset by probing model and compares with potential members dataset. Uses projection function to increase variance and measure distance changes.

# White Box Membership Inference

White box membership inference attacks exploit direct access to model parameters, gradients, and intermediate computations to determine training data membership with high accuracy.

### Gradient Exploitation

Analyzes footprint of gradients in the Stochastic Gradient Descent (SGD) algorithm to identify training samples.

### Hidden Layer Analysis

Examines outputs of hidden layers to detect patterns unique to training data.

### Loss Value Assessment

Uses loss function behavior to distinguish between training and non-training samples.

### Embedding Clustering

Computes membership probability or embedding for each data point and uses clustering to separate members from non-members.

These attacks are particularly effective in scenarios involving model access after training, transfer learning, fine-tuning, or collaborative learning environments where gradient information is shared.

# Attack Scenarios for Model Inversion

Model inversion attacks can be deployed in various real-world scenarios, posing significant risks to organizations and individuals.

## Facial Recognition

Reconstructing recognizable images of people's faces from their names and model confidences, as demonstrated in the MIFace attack.

## Healthcare

Inferring sensitive medical information from disease prediction models, potentially exposing patient data.

## Biometric Systems

Reconstructing fingerprints or other biometric data from authentication systems.

## Survey Data

Inferring sensitive features such as marital infidelity or viewing habits from lifestyle survey models.

# Attack Scenarios for Inference Attacks

Inference attacks can occur in various contexts, exposing sensitive information about training data and potentially violating privacy regulations.

## Cloud ML Services

Inferring membership of data points in training data of ML-as-a-service offerings, such as face recognition or sentiment analysis APIs provided by major cloud vendors.

## Collaborative Learning

Extracting information from systems where multiple parties jointly train models without sharing raw data. This can become a white box scenario if the attacker is part of the collaboration.

## Online Training Systems

Exploiting continuously updated models to prove a user has used the AI system, potentially revealing sensitive contexts (e.g., specialized medical services, military systems).

## Federated Learning

Inferring properties about participants' data in decentralized learning environments where model updates rather than raw data are shared.

# Mitigations for Model Inversion Attacks

Protecting against model inversion attacks requires a multi-layered defense strategy that addresses various vulnerabilities in machine learning systems.

## Gated API Pattern

Isolate inference API with strict network and least-privilege access controls, preventing direct attacker access.

## Access Controls

Implement effective model governance and least-privilege access to prevent malicious internal access or information leakage.

## Rate Limiting

Utilize API rate limiting to hinder iterative optimization algorithms used in attacks.

## Output Modification

Use ArgMax-ed output with target label instead of probabilities, or apply rounding/noise to confidence scores.

## Privacy-Preserving ML

Apply differential privacy techniques and select privacy-aware training algorithms.

## Monitoring

Implement detection systems with visualization for auditing suspicious query patterns.

ISMP
Intelligent System
Media Processing Lab

# Mitigations for Inference Attacks

Defending against inference attacks requires addressing overfitting and memorization while implementing privacy-preserving measures and securing inference endpoints.

## Preventing Overfitting

- Regularization techniques (early stopping, weight decay, dropout)
- Adversarial regularization targeting membership attacks
- Maximum mean discrepancy (MMD) regularization
- Data augmentation to increase diversity
- Model stacking with separate training datasets
- Data minimization and anonymization

## Privacy-Preserving Measures

- Differential privacy adding noise to outputs or training
- Membership inference adversarial training
- MemGuard adding strategic noise to confidence scores
- Gated API patterns limiting direct access
- API rate limiting to prevent repeated queries
- Monitoring and alerting for suspicious patterns

# Differential Privacy: A Key Defense

Differential privacy has emerged as one of the most effective defenses against privacy attacks in machine learning. It works by adding carefully calibrated noise to protect individual data points while maintaining overall utility.

### Data Collection

Original sensitive data is collected for model training

### Noise Addition

Carefully calibrated noise is added to data or model outputs

### Privacy-Utility Tradeoff

Balance is maintained between privacy protection and model performance

### Private Training

Model is trained with privacy-preserving algorithms

While differential privacy is highly effective against membership inference attacks, it may be less effective against attribute inference attacks that target dataset-wide properties rather than individual records.

# Limitations of Differential Privacy

While differential privacy is a powerful defense mechanism, it's important to understand its limitations, particularly when dealing with attribute inference attacks.

## Effective Against

- Membership inference attacks targeting individual records

- Model inversion attacks attempting to reconstruct specific data points

- Attacks that rely on model memorization of exact training examples

- Privacy violations that focus on identifying individuals in the dataset
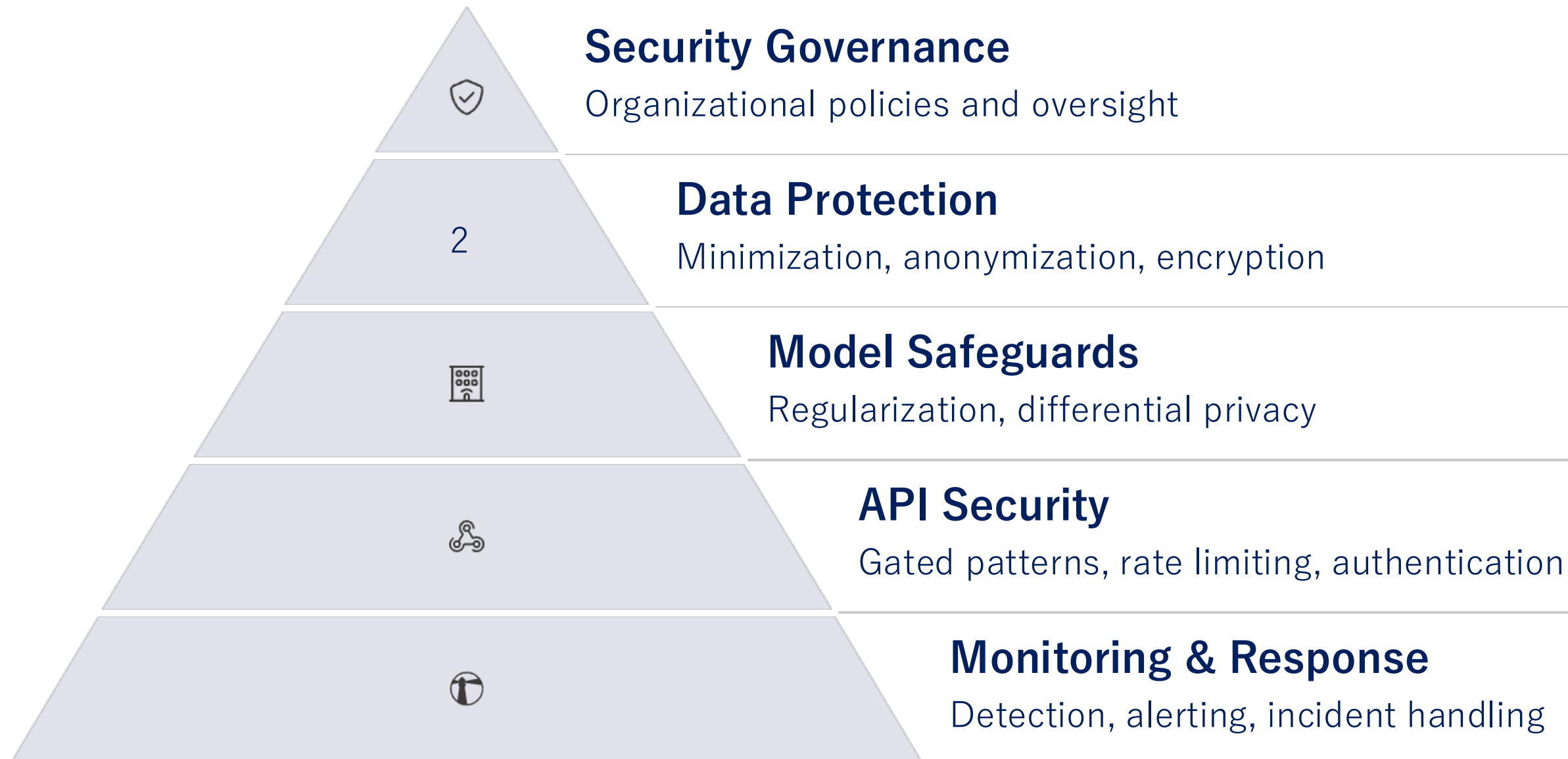
## Less Effective Against

- Attribute inference attacks targeting dataset-wide properties

- Inference of aggregated statistics about the training data

- Attacks that exploit general patterns rather than specific examples

- Poisoning-assisted inference attacks that inject triggers in the dataset

A comprehensive defense strategy should combine differential privacy with other techniques like secure inference endpoints, data governance, and monitoring to address these limitations.

# Defense in Depth Strategy

Protecting against privacy attacks requires a comprehensive defense-in-depth approach that addresses vulnerabilities at multiple levels of the machine learning pipeline.

**Security Governance**
Organizational policies and oversight

**Data Protection**
Minimization, anonymization, encryption

**Model Safeguards**
Regularization, differential privacy

**API Security**
Gated patterns, rate limiting, authentication

**Monitoring & Response**
Detection, alerting, incident handling

# Securing the MLOps Pipeline

Protecting against white box attacks requires securing the entire Machine Learning Operations (MLOps) pipeline to prevent unauthorized access to models and training data.

## Data Acquisition

Secure data collection, storage, and preprocessing with access controls and encryption

## Model Development

Implement least privilege access, code signing, and secure development practices

## Deployment

Use secure CI/CD pipelines with integrity verification and container security

## Monitoring

Implement continuous security monitoring, anomaly detection, and audit logging

Addressing supply chain vulnerabilities is particularly important to prevent attackers from obtaining copies of publicly available models used as baselines for your systems.

Intelligent System
Media Processing Lab

# Gated API Pattern

The Gated API pattern is a crucial defense mechanism that isolates inference endpoints and controls access to machine learning models, making it harder for attackers to probe models repeatedly.

## Authentication Layer
Verify user identity with strong authentication mechanisms

## Validation Layer
Validate inputs for format, range, and potential malicious content

## Rate Limiting Layer
Restrict number of requests to prevent iterative optimization attacks

## Proxy Layer
Intermediate service that communicates with actual model

## Output Processing
Modify confidence scores or limit information in responses

# Monitoring and Detection

Effective monitoring and detection systems are essential for identifying potential privacy attacks in progress and responding before sensitive data is compromised.

## Query Pattern Analysis

Monitor for suspicious patterns like systematic exploration of feature space, repeated similar queries with small variations, or unusually high query volumes that could indicate model inversion or inference attacks.

## Visualization Tools

Implement tools that visualize inputs in realistic terms (e.g., images) rather than numerical vectors to help security teams identify anomalous or potentially malicious queries during audits.

## Behavioral Baselines

Establish normal usage patterns for different user types and alert on deviations that might indicate attack attempts, such as unusual query distributions or access patterns.

## Automated Response

Develop systems that can automatically throttle or block suspicious activity when detected, such as temporarily restricting access for users exhibiting attack-like behavior.

# Regularization Techniques



Regularization techniques help prevent overfitting and memorization, which are key factors that make models vulnerable to inference attacks.

## ☐ Early Stopping

Halt training when validation performance stops improving to prevent memorization of training data.

## Weight Decay

Add penalty term to loss function that reduces magnitude of weights, preventing complex fitting to training data.

## Dropout

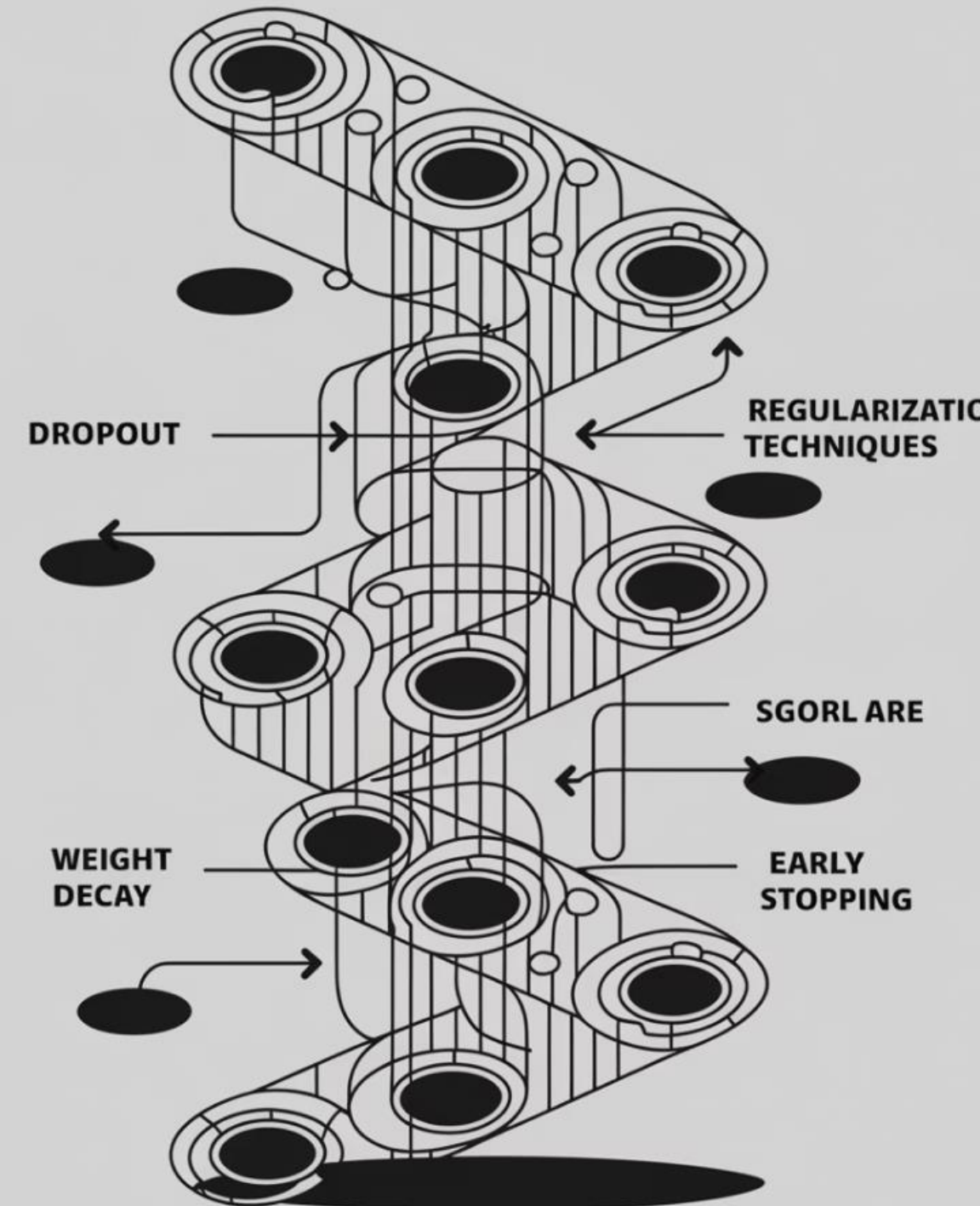Randomly disable neurons during training to prevent co-adaptation and improve generalization.

## Adversarial Regularization

Train model to minimize both original objective and membership inference objective simultaneously.

Effective regularization not only improves model generalization but also significantly reduces vulnerability to membership inference attacks by minimizing the gap in model behavior between training and non-training data.

Intelligent System Media Processing Lab

# Data Augmentation as Defense

- Data augmentation serves as both a performance enhancement technique and a privacy defense by increasing the diversity of training data and reducing memorization of specific examples.

  - directly countering membership inference attacks.

- Augmentation also helps address the limited training data problem that often leads to overfitting, which is a key vulnerability exploited in privacy attacks.

# Model Stacking Defense

Model stacking is an ensemble technique that can effectively counter membership inference attacks by distributing knowledge across multiple models trained on different data subsets.

### Stacking Architecture

A two-level tree structure with multiple base classifiers feeding into a top-level meta-classifier. Each base model is trained on a different subset of the data.

### Data Partitioning

Training data is divided into separate subsets, ensuring no single model sees all examples. This prevents any individual model from memorizing specific samples.
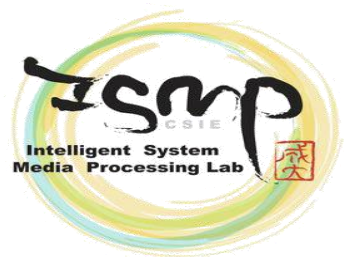
### Prediction Process

During inference, base models generate predictions that are combined by the meta-classifier. This distribution of knowledge makes it difficult for attackers to determine if a specific example was in the training set.

By distributing knowledge across multiple models, stacking reduces the memorization of specific examples while maintaining overall prediction accuracy.
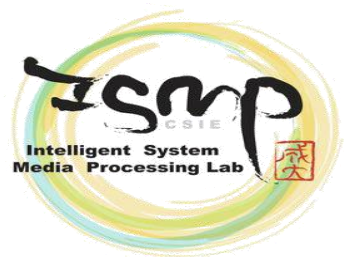
# MemGuard: Strategic Noise Defense

- MemGuard adds carefully crafted noise to model confidence scores to confuse membership inference attacks while preserving model utility.

  - MemGuard introduces strategically generated noise into confidence score vectors produced by the target classifier. This noise is designed to:

    - Confuse attacker's classifier, causing it to make random guesses about membership status

    - Maintain the same predicted label to preserve utility

    - Minimize distortion to keep confidence scores useful

    - Be imperceptible to normal users while effective against attacks

- MemGuard treats defense as an adversarial example generation problem, creating noise that specifically targets membership inference classifiers.

  - Unlike differential privacy, which adds random noise, MemGuard's noise is optimized to specifically counter membership inference attacks while minimizing impact on legitimate use.

# Privacy-Utility Tradeoff

- Implementing privacy protections often involves balancing the tradeoff between privacy guarantees and model utility or performance.

  - As privacy protection measures become stronger (e.g., more noise in differential privacy, stronger regularization, or more aggressive output modifications), model accuracy typically decreases.

  - Organizations must carefully consider their specific requirements and risk tolerance when determining the appropriate balance.

- The optimal tradeoff point varies based on data sensitivity, regulatory requirements, attack risk, and performance needs.

  - Critical applications handling highly sensitive data may require stronger privacy guarantees despite performance impacts.

# Risk Assessment Framework

A structured risk assessment framework helps organizations evaluate their vulnerability to privacy attacks and prioritize appropriate defenses.

## Identify Assets

Catalog ML models, training data, and sensitive information that requires protection.

## Threat Modeling

Identify potential attackers, their capabilities, and motivations for targeting your systems.

## Vulnerability Assessment

Evaluate model characteristics that increase risk (e.g., overfitting, confidence outputs, sensitive data).
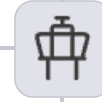
## Impact Analysis

Determine consequences of successful attacks (privacy violations, regulatory penalties, reputation damage).

## Control Selection

Choose appropriate defenses based on risk level, technical feasibility, and resource constraints.

ISMP
Intelligent System
Media Processing Lab

# Regulatory Considerations

Privacy attacks on machine learning models can have significant regulatory implications under various data protection frameworks.

## GDPR (Europe)

Successful privacy attacks may constitute data breaches under GDPR if they reveal personal data. Organizations must implement appropriate technical measures to protect data, which increasingly includes ML privacy protections.

## CCPA/CPRA (California)

These regulations grant consumers rights regarding their personal information. Privacy attacks that reveal consumer data could trigger notification requirements and potential penalties.

## HIPAA (US Healthcare)

For healthcare ML applications, privacy attacks that reveal protected health information would constitute HIPAA violations, requiring breach notification and potentially resulting in significant penalties.

## AI-Specific Regulations

Emerging AI regulations like the EU AI Act include requirements for technical robustness and security of AI systems, which will likely encompass protections against privacy attacks.

# Future Research Directions

### Privacy-Utility-Robustness Tradeoffs

Exploring the complex relationships between privacy protection, model utility, and robustness against adversarial examples to develop optimal defense strategies.

### Privacy in Federated Learning

Developing specialized techniques to protect against privacy attacks in federated learning environments where models are trained across distributed devices.

### Quantum-Resistant Privacy

Preparing for the era of quantum computing by developing privacy-preserving techniques that remain secure against quantum-enhanced privacy attacks.

### Privacy in Generative Models

Addressing unique privacy challenges in large language models and other generative AI systems that can memorize and reproduce training data.