# Digital Clones and Digital Immunity: One Adversarial Solution for for Two Critical AI Challenges

An expert explainer on a unified approach to building personalized and resilient AI for Industry 4.0.

Based on the research by Vladyslav Branytskyi, Mariia Golovianko, Svitlana Gryshko, Diana Malyk, Vagan Terziyan, and Tuure Tuunanen.

NotebookLM

# Two Grand Challenges for AI in Industry 4.0

## The Challenge of Personalization (Digital Cloning)

**Goal:** To create 'Digital Clones' of human decision-makers—capturing their unique expertise, intuition, and even personal biases.

**Why it Matters:** Enables automation, virtual presence in multiple locations, and experimentation in simulated environments without risk. It's about replicating *how* an expert thinks.

**Key Question:** How can we achieve the highest accuracy in mimicking a human's specific decision-making behavior?

## The Challenge of Robustness (Digital Immunity)

**Goal:** To develop 'Digital Immunity'—the capability of AI systems to operate reliably and resist adversarial attacks (e.g., poisoning, evasion).
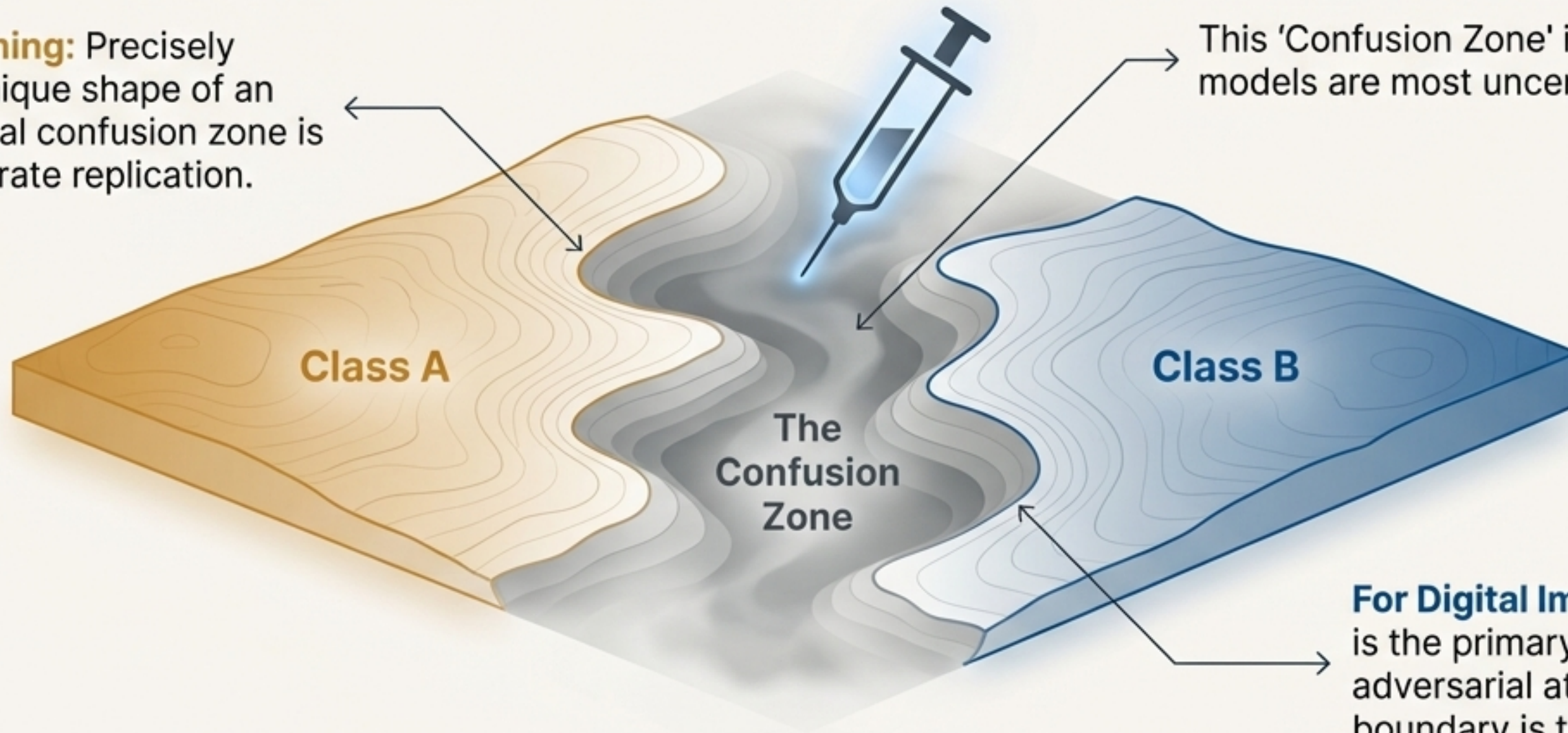
**Why it Matters:** Ensures safety, reliability, and security in critical systems, especially in smart manufacturing and logistics where ML models are vulnerable.

**Key Question:** How can we make our models secure and resilient against malicious, crafted inputs designed to cause failure?

# The Unseen Connection:
# Both Problems Live at the Decision Boundary

**For Digital Cloning:** Precisely mapping the unique shape of an expert's personal confusion zone is the key to accurate replication.

This 'Confusion Zone' is where AI models are most uncertain.

Class A

The Confusion Zone

Class B

**For Digital Immunity:** This zone is the primary target for adversarial attacks. Hardening this boundary is the key to resilience.

**The accuracy of a clone and the security of a model are both determined by how we handle the ambiguity at the decision boundary.**

# Adversarial Training: The Unified Solution

The solution to both problems is adversarial training—learning on automatically generated, challenging samples that lie near the decision boundary.
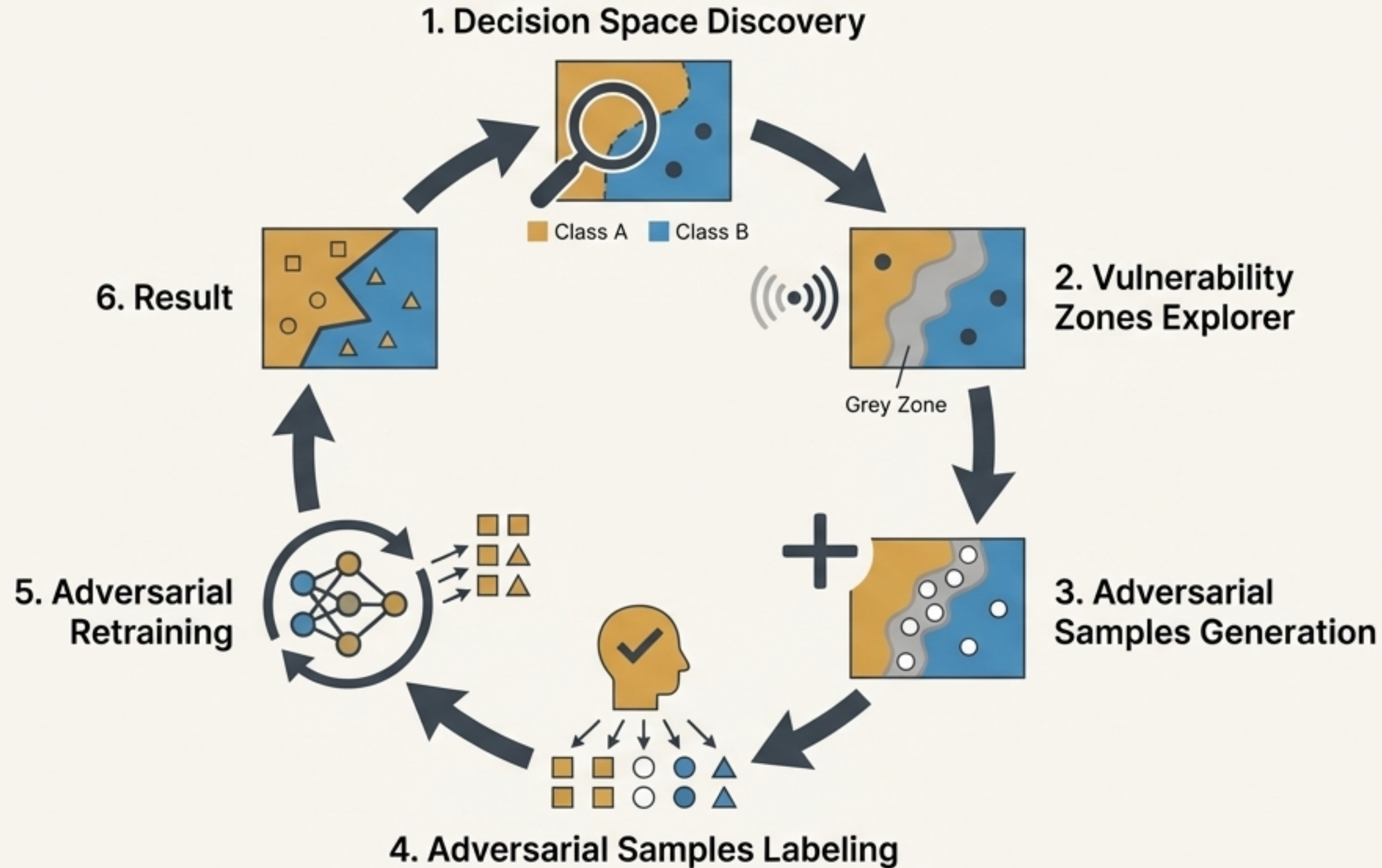


Adversarial Training

**For Cloning:** It forces the model to learn the fine-grained, personal details of the decision boundary, increasing personalization.

**For Immunity:** It's like a 'digital vaccine,' exposing the model to tricky examples to build resistance and harden the boundary against future attacks.

*"Both problems (clones and immunity training) have the same backbone solution, which is adversarial training."*

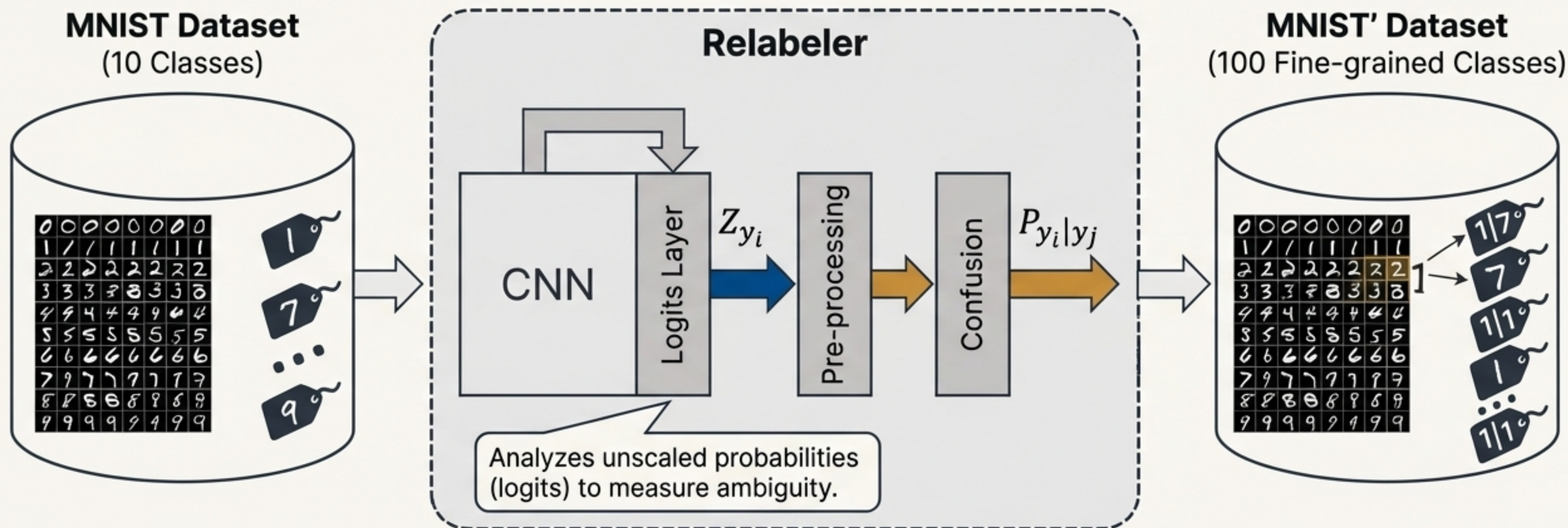# Our Methodology: A Cycle of Adversarial Retraining

**1. Decision Space Discovery**

Class A  Class B

**6. Result**

**2. Vulnerability Zones Explorer**

Grey Zone

**3. Adversarial Samples Generation**

**4. Adversarial Samples Labeling**

**5. Adversarial Retraining**

**Result:** A learning process that creates a robust and personalized decision boundary for the Digital Clone, enhancing its Digital Immunity.

NotebookLM

# Step 1: Finding Boundary Samples via Confusion-Driven Relabeling

Instead of just classifying an image as '1', we determine *how much* it might be confused with other classes. An ambiguous '1' that looks like a '7' is more valuable for training.



**MNIST Dataset**
(10 Classes)

**Relabeler**

Logits Layer $Z_{y_i}$

Pre-processing

Confusion $P_{y_i|y_j}$

Analyzes unscaled probabilities (logits) to measure ambiguity.

**MNIST' Dataset**
(100 Fine-grained Classes)

**Key Takeaway:** This process automatically identifies and categorizes the most informative samples—the 'border-guards'—for targeted retraining.
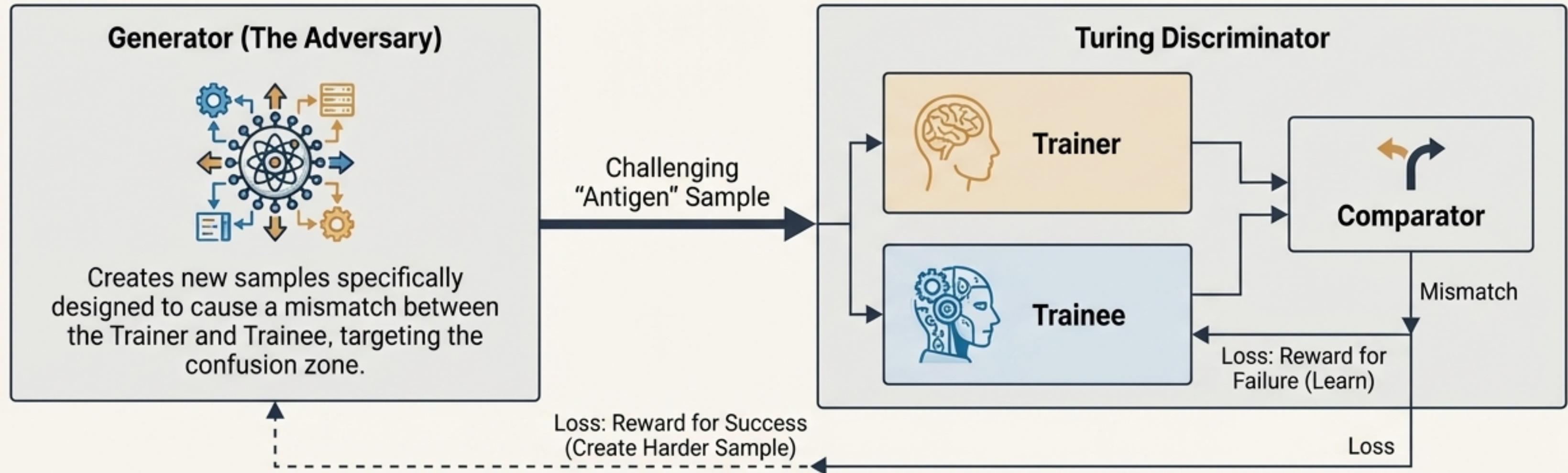
Source Reference: This diagram is a clear, step-by-step visual explanation of the process described and shown in <IMAGE 0> and <IMAGE-1>.

NotebookLM

# Step 2: Generating a 'Digital Vaccine' with Turing-GANs

**Problem Statement:**
Standard GANs are for creating realistic data, not for training a classifier. We need a new architecture.

**Solution Statement:**
Our Solution: The Turing-GAN (T-GAN)—An adversarial game with three players.



**Generator (The Adversary)**

Creates new samples specifically designed to cause a mismatch between the Trainer and Trainee, targeting the confusion zone.

Challenging "Antigen" Sample

**Turing Discriminator**

**Trainer**

**Trainee**

**Comparator**

Mismatch

Loss: Reward for Failure (Learn)

Loss

Loss: Reward for Success (Create Harder Sample)

**How It Works**
The Generator is rewarded for finding samples where the Trainee fails to copy the Trainer. This accelerates the learning process by focusing it on the most difficult and informative areas of the decision space.

NotebookLM

# One Architecture, Two Missions

## Case 1: Digital Cloning

- **Trainer = "Donor":** The human expert whose cognitive skills are being cloned.

- **Trainee = "Clone":** The model learning the personalized decision boundary of the donor.

- **Goal:** Minimize the difference between Clone and Donor decisions.

## Case 2: Digital Immunity

- **Trainer = "Supervisor":** An oracle or a highly robust model that knows the "true labels."

- **Trainee = Vulnerable Classifier:** The model being "vaccinated" to become robust.

- **Goal:** Retrain the classifier on generated adversarial samples (the "digital vaccine") to harden it against attacks.

# Evidence: Cloning an Airport Security Expert with Near-Perfect Fidelity

## Case Study Context

Simulating airport luggage inspection using a conveyor belt system. Three human experts (DM_1, DM_2, DM_3) labeled 2,198 images as "dangerous" or "not dangerous". We trained three digital clones.

## The Key Metric: Correlation

How often did the clone's decision match the human's?

## 92.1%
**Correlation:** Clone 1 vs. DM_1

## 99.53%
**Correlation:** Clone 2 vs. DM_2

## 94.94%
**Correlation:** Clone 3 vs. DM_3

## Bonus Insight

The clones are so accurate, their advice *improved* the human expert's decision-making correctness (e.g., DM_1 improved from 95.45% to 97.27%).

# Evidence: 'Vaccinating' an AI Against Adversarial Attacks

**\*The Attack Scenario\*:** We "poisoned" images to trick the AI. (Scenario 1: A 'bomb' is present, but the image is altered to look safe. Scenario 2: No bomb, but the image is altered to cause a false alarm.)

## The Impact of Attack (Before Vaccination)

Without adversarial retraining, the classifier's accuracy can be decreased to **less than 1%**. **Human experts** were also fooled, with accuracy dropping to **65-75%** on poisoned images.

## The Result of Vaccination (After Retraining)

After retraining on 300 generated "vaccine" samples, the digital clone's accuracy on tampered images significantly increased, reaching up to **80.05%**. This demonstrates acquired immunity.
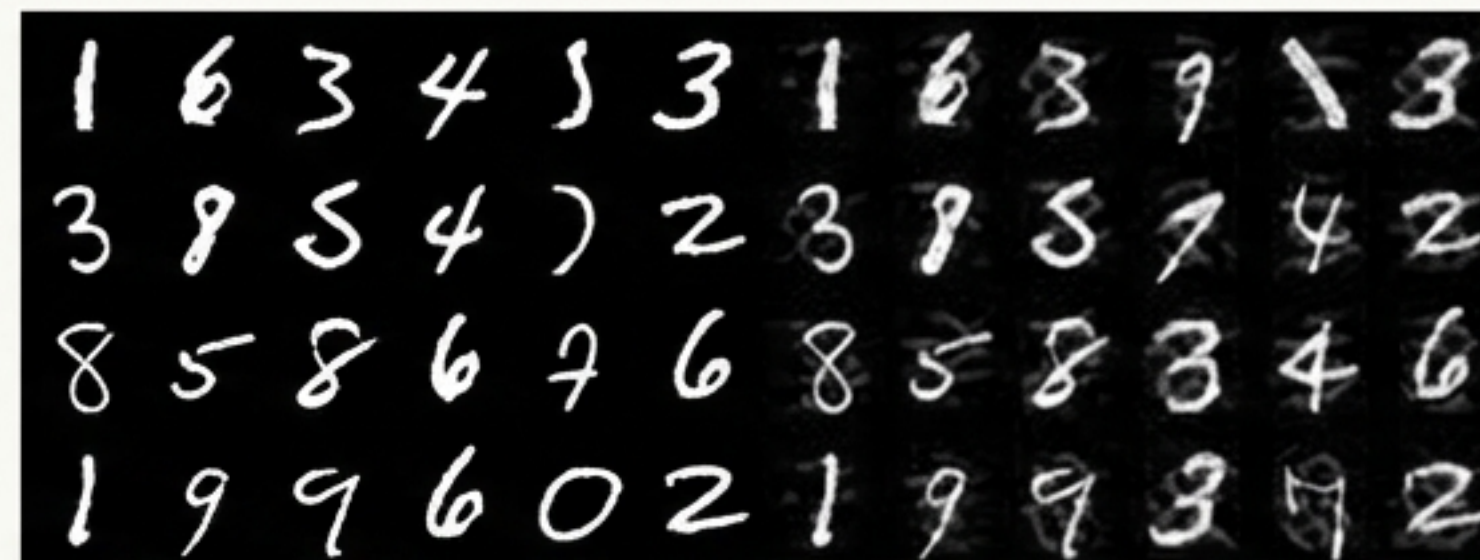
# Visualizing the 'Digital Vaccine'

## Ambiguous Samples from the Boundary



These are entirely new images created by a GAN trained to explore the "confusion zone" between two classes. They represent the most challenging cases for a classifier to learn.

## Perturbed Real Samples



These are original images that have been subtly modified with generated noise, "pushing" them towards the decision boundary to make them harder to classify correctly. The changes are often imperceptible but are designed to fool the AI.

# A New Paradigm: Building AI That is Both Personalized AND Resilient

**Adversarial training should not be viewed merely as a defensive technique.
It is a fundamental tool for building the next generation of intelligent systems.**

1. **A Unified Problem:** The challenges of creating high-fidelity digital clones and robust, secure AI both originate at the model's decision boundary.

2. **A Unified Solution:** Adversarial training, through targeted sample generation and retraining, provides a single, powerful method to master this boundary.

3. **Proven in Practice:** Our experiments demonstrate that this approach achieves near-perfect human-clone correlation while simultaneously creating measurable 'digital immunity' to adversarial attacks.

**This unified approach allows us to move beyond the tradeoff between accuracy and security, enabling the creation of AI that is simultaneously specialized and robust.**

# The Road Ahead: Implications and Future Research

## Broader Impact for Industry 4.0

- **Smart Manufacturing**: Development of more reliable AI for quality control and process automation.

- **Autonomous Systems**: Safer and more personalized control systems for vehicles and robots.

- **Cybersecurity**: A new proactive method for training 'digital security officers' that are immune to novel attacks.

## Our Future Research Directions

- **Complete Boundary Coverage**: Developing methods to guarantee that our generated adversarial samples cover the *entire* decision surface evenly, not just select areas.

- **Cloning Evolving Targets**: Addressing the challenge of cloning intelligent agents that have evolving, non-deterministic behavior over time.

*By treating the decision boundary as a primary object of study, we unlock a more powerful and unified way to build intelligent systems.*