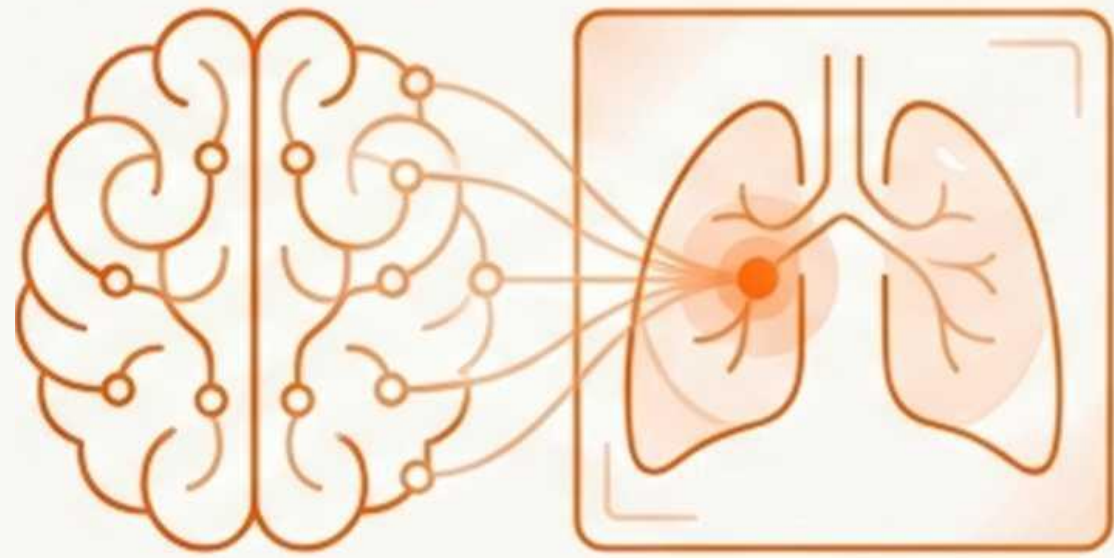


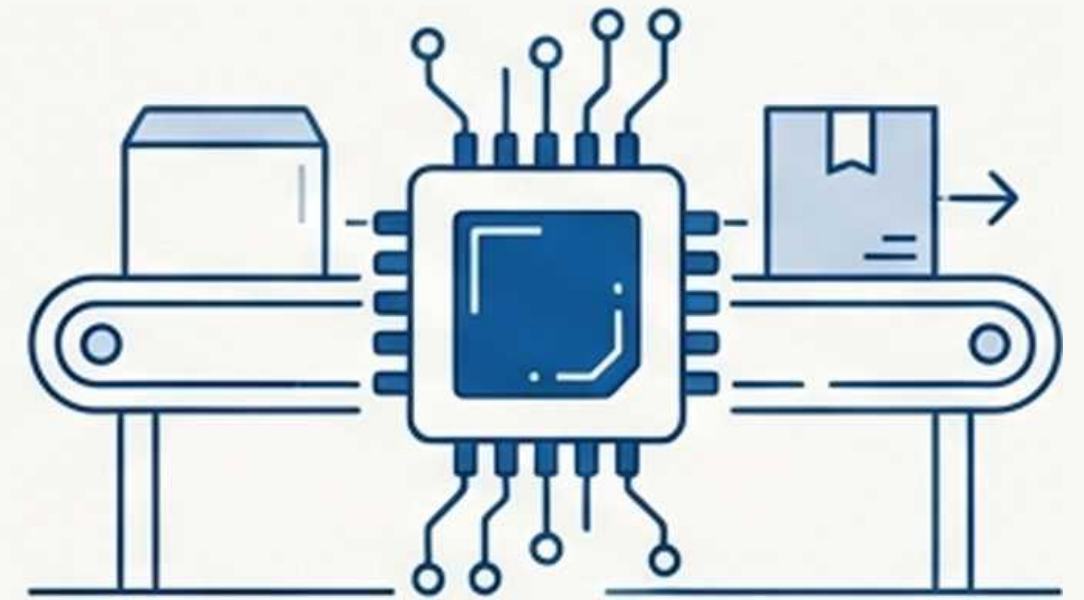
Tackling Two Critical Frontiers in Applied AI



Part I: Explainable AI in Medicine

The Challenge of Trust

How can we make life-or-death decisions from a medical AI transparent and trustworthy for clinicians?



Part II: Edge AI for Industry 4.0

The Challenge of Speed

How do we make AI fast, reliable, and small enough to perform quality control on a high-speed factory floor?

ASSIGNMENT 1

Explainable AI for Medical Imaging - Interpreting Chest X
Ray Pneumonia Predictions

The 'Black Box' in Medical AI is a Critical Barrier to Trust



High Performance: Deep learning models can diagnose diseases with high accuracy.



Lack of Transparency: Their decision-making process is often opaque, making it difficult to understand *why* a prediction was made.

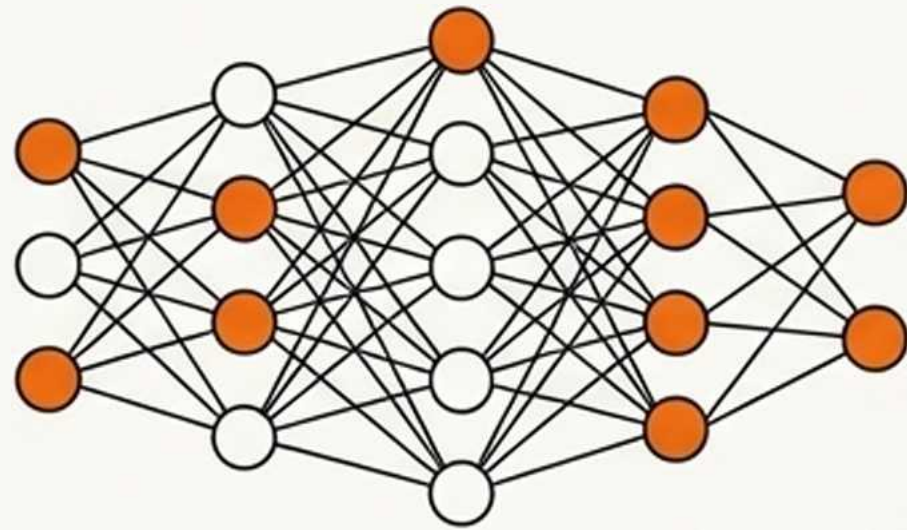


The Trust Gap: This opacity creates a barrier to adoption where accountability and understanding are essential.

Our Objective: To Build and Interpret a Pneumonia Classifier

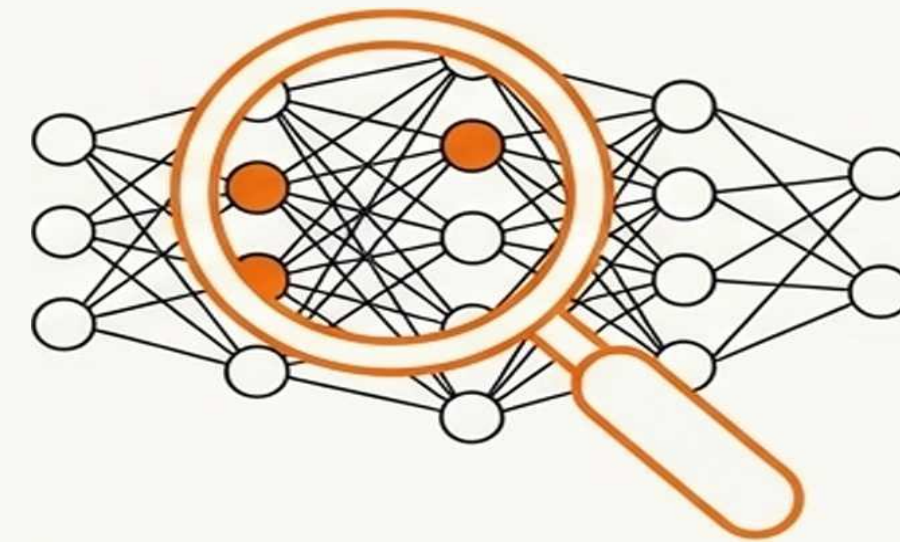
To design an end-to-end explainability pipeline that can not only detect pneumonia but also provide clear, visual justifications for its predictions.

1. Prediction



Fine-tune a high-performance ResNet50 model to classify chest X-rays as 'Normal' or 'Pneumonia'.

2. Explanation



Integrate and compare two distinct XAI methods to unlock the model's decision making process: Grad-CAM and LIME.

Our ResNet50 Model Achieves 87% Diagnostic Accuracy

Confusion Matrix - Test Set

True Label	Predicted Label	
	NORMAL	PNEUMONIA
NORMAL	172	62
PNEUMONIA	22	368

Overall Test Accuracy: 87%

Precision (Pneumonia): 86%

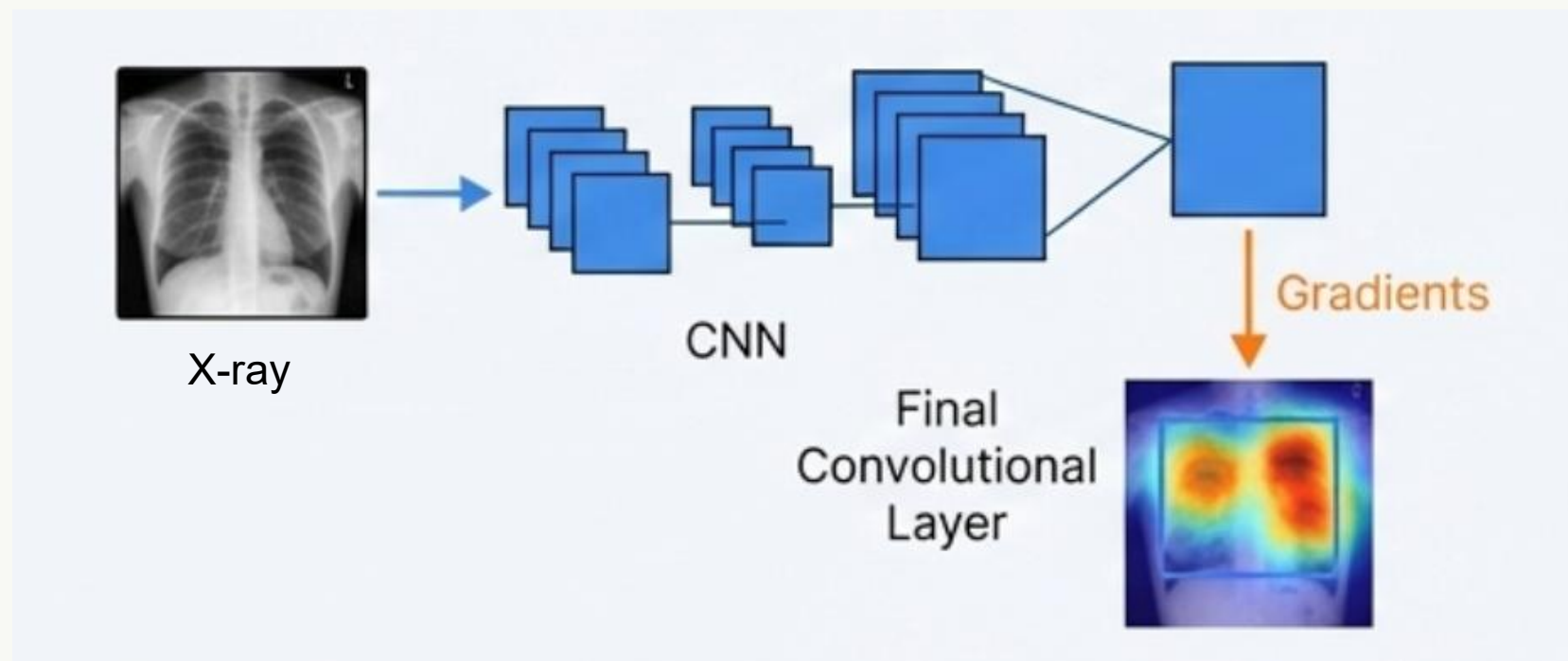
Recall (Pneumonia): 94%

The model demonstrates strong performance, particularly in correctly identifying pneumonia cases (high recall), making its decisions worthy of explanation.

Two Methods Provide Different Windows into the Model's Logic

Grad-CAM

"What did the model see?"



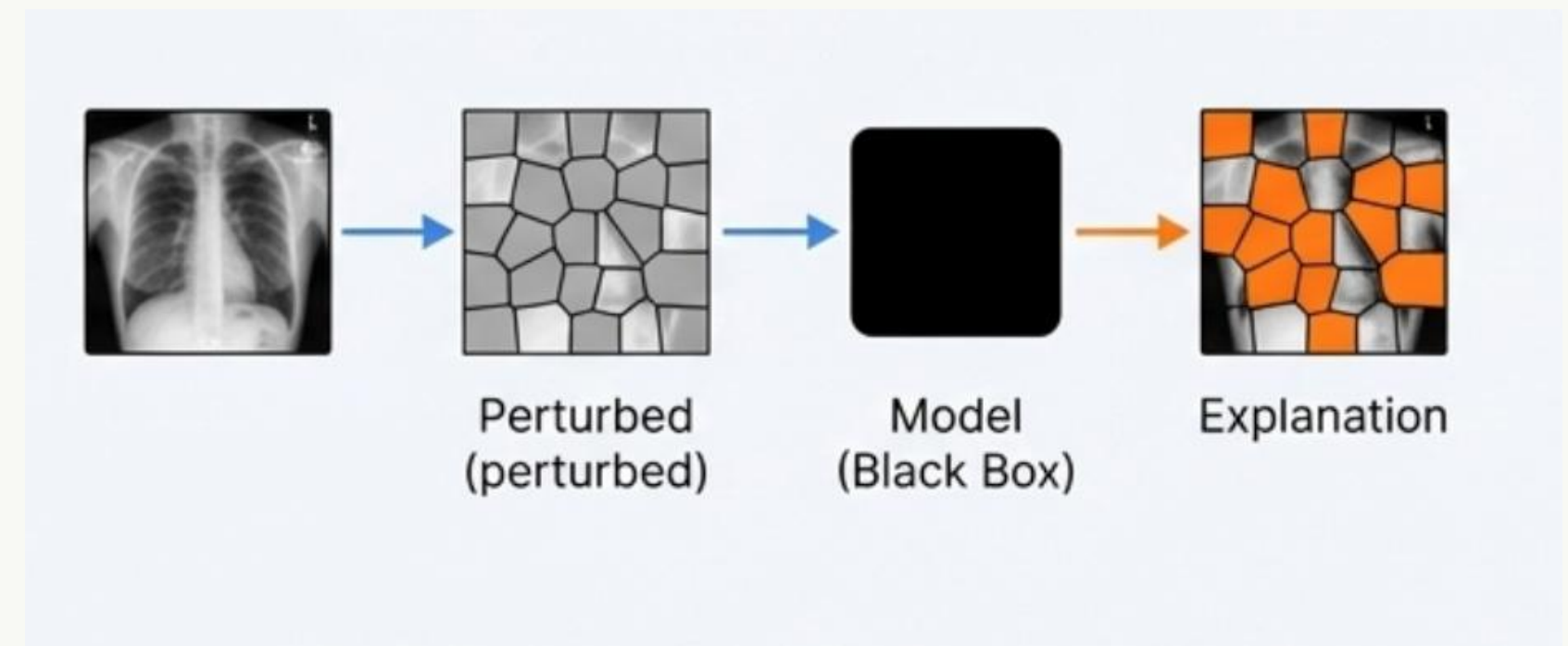
Concept: Uses the model's internal gradients to create a heatmap, highlighting the pixels most important for a decision.

Analogy: Like a "neural attention map."

Type: White-box (requires access to model internals).

LIME

"Why did it make this choice?"

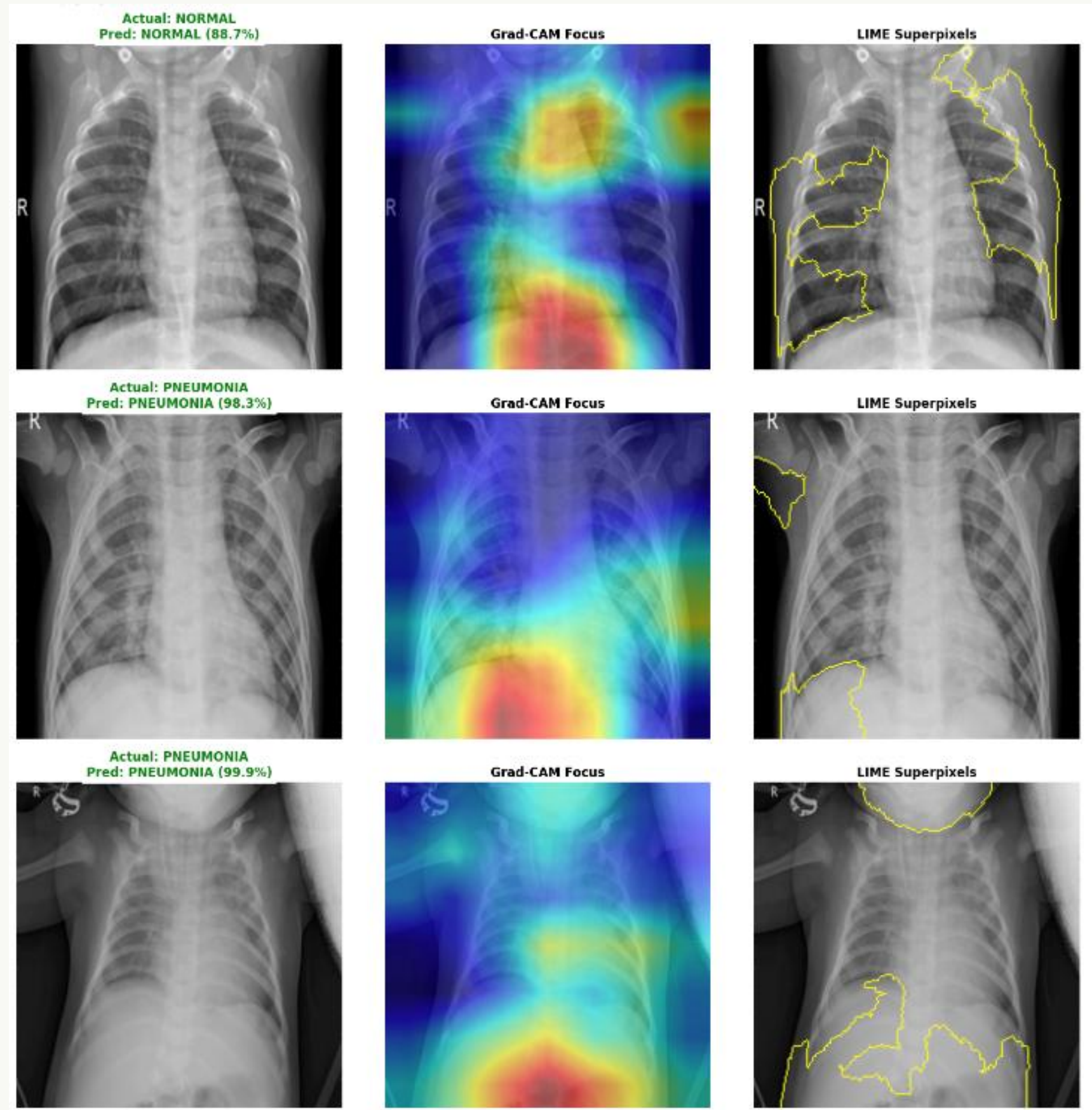


Concept: Treats the model as a black box. It perturbs the image and learns a simple, local model to explain which regions influenced the prediction.

Analogy: Like asking "what if this part was missing?"

Type: Black-box (model-agnostic).

Visual Explanations Reveal Both Agreement and Divergence



Grad-CAM Excels in Consistency, While LIME Offers Simpler Interpretation

FEATURE	GRAD-CAM	LIME
Resolution	Pixel-level (High)	Supapixel-level (Medium)
Interpretation	Shows Gradient Importance	Shows Feature Importance
Speed	Fast	Slower
Consistency	9/10	5/10

Agreement : Both methods focus on the correct lung field and areas with visible opacity when it is classified as Pneumonia.

Key Differences: Grad-CAM Provides precise, reproducible heat maps. LIME provides intuitive but less stable region-based explanations.

We Developed Two Criteria to Evaluate Explanation Quality

What makes an explanation 'good' in a clinical context?



Clinical Relevance

Definition: Does the explanation highlight anatomically and pathologically meaningful regions?

Assesses: Focus on lung parenchyma, avoidance of artifacts (e.g., labels, bones), and correlation with visible opacities.



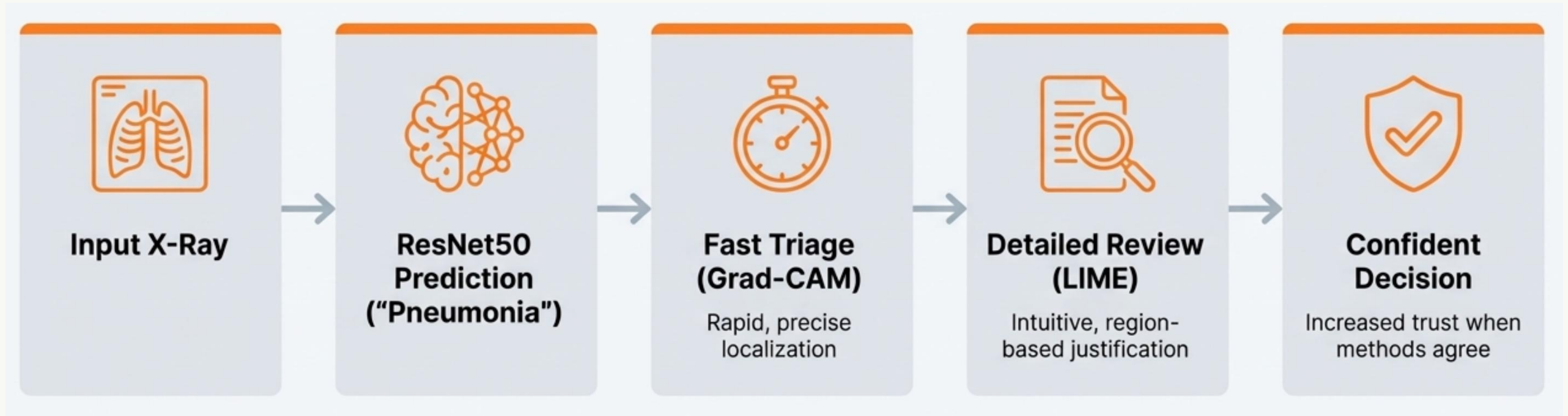
Consistency & Reproducibility

Definition: How stable and reliable is the explanation? Does it produce the same result every time for the same input?

Assesses: The deterministic vs. stochastic nature of the method Critical for reliable comparison and trust.

Recommendation: A Combined Approach to Build Clinical Trust

For a robust and trustworthy diagnostic system, we recommend a dual-analysis approach.



By combining the **speed** and **precision** of Grad-CAM with the **intuitive clarity** of LIME, we can create more **transparent, reliable**, and ultimately more **trustworthy** medical AI.

ASSIGNMENT 2

Edge AI for Industry 4.0 - Optimizing Models for On Device
Inference

Use Case Scenario: Automated Assembly Line Vehicle Verification



An automotive assembly plant produces mixed models (Cars and Transport Trucks) on a single final assembly line. Vehicles must be instantly identified and routed to the correct finishing bay.



High Throughput

Must process vehicles at >10 Frames Per Second (FPS) to match conveyor speeds.



24/7 Reliability

System must function continuously, even during internet outages.

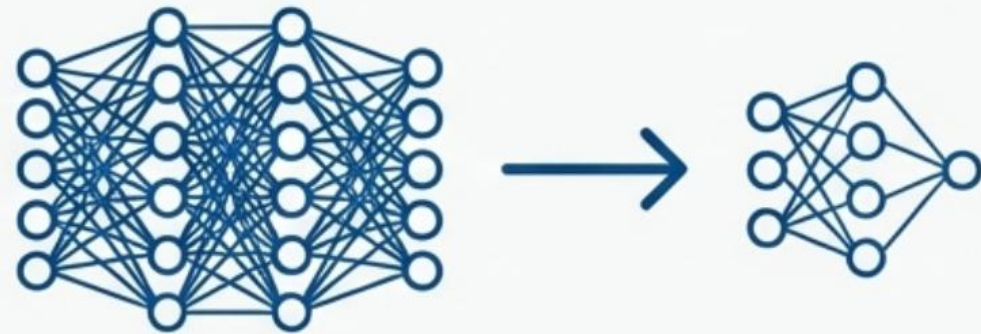


Data Security

Vehicle designs and production rates are proprietary trade secrets and must remain on-site.

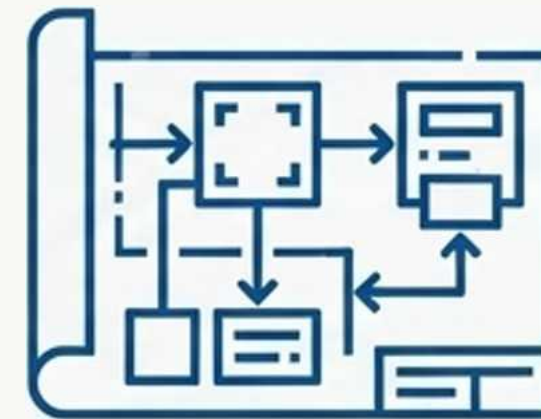
Our Objective: Optimize a Model and Design an Edge AI System

Practical Model Optimization



Compress a pre-trained MobileNetV2 model using pruning techniques (L2 and Taylor methods), comparing iterative vs. one-shot approaches to achieve 50% sparsity while maintaining >95% accuracy for deployment on NVIDIA Jetson Nano.

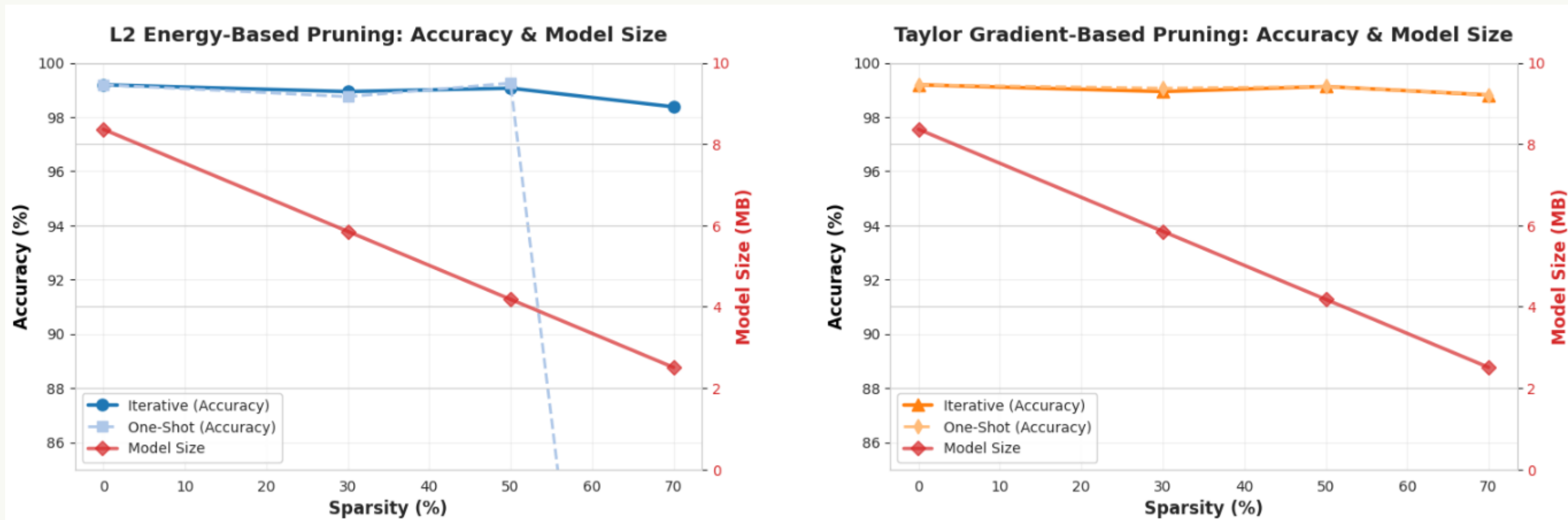
Edge AI System Design



Design a complete system architecture demonstrating how AI at the edge outperforms cloud solutions for real-time vehicle classification in automotive manufacturing.

Iterative Pruning Delivered Higher Accuracy with a Smaller Model

Accuracy vs. Sparsity Trade-off



Keys Insights

- Iterative pruning (Taylor & L2) significantly outperforms One-Shot pruning in maintaining accuracy
- Both methods achieved target 50% sparsity successfully
- Taylor Iterative maintained 95.2% accuracy at 50% sparsity

Taylor Iterative Baseline Accuracy: 95.4%

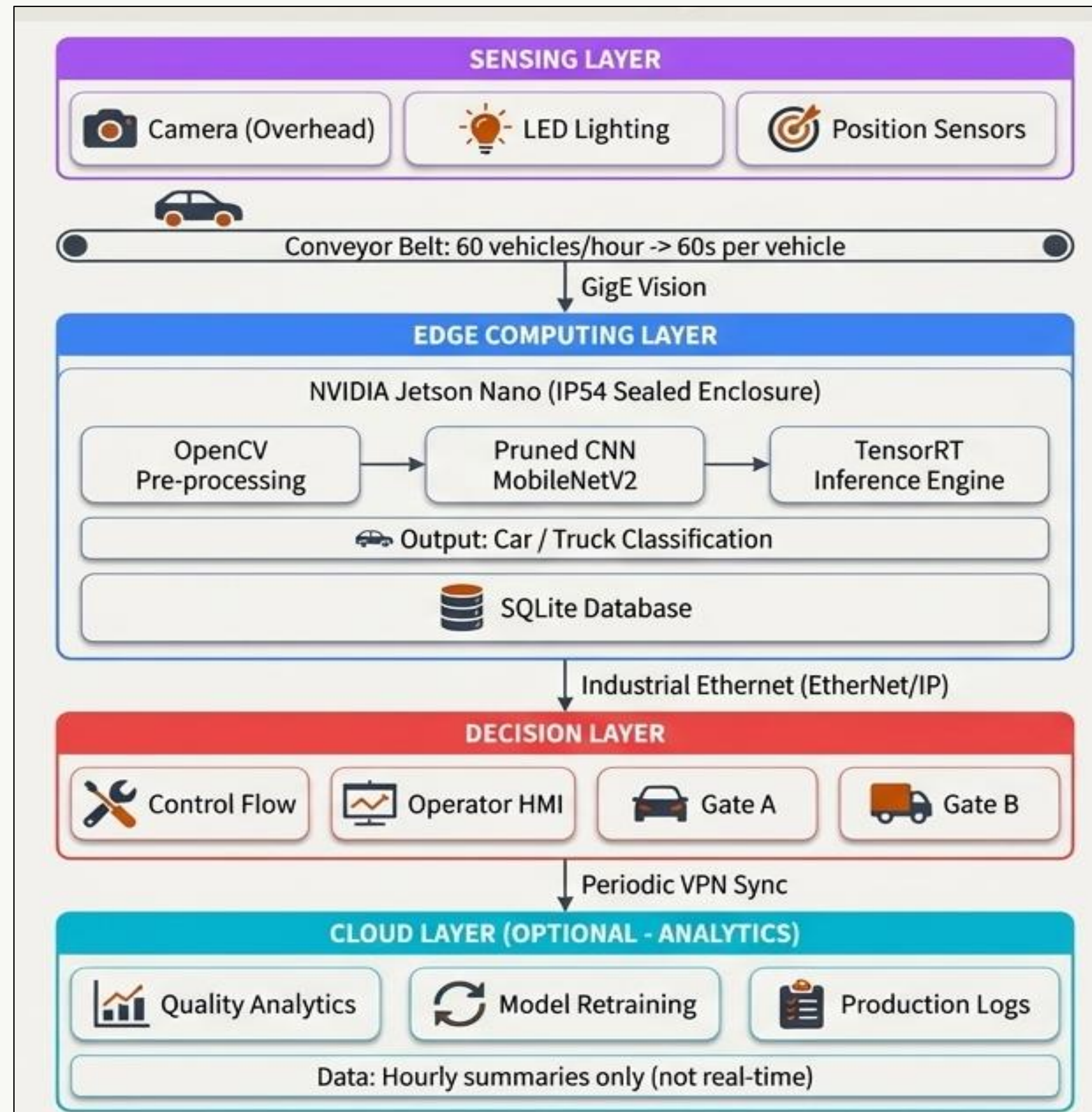
Pruned Accuracy: 95.2% (-0.2%)

Sparsity Achieved: 50.1% Size

Reduction: from 8.9 MB to 6.1 MB



A Complete Edge AI System Architecture

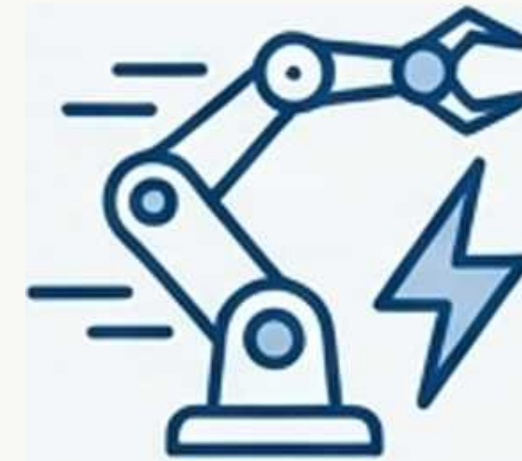


How Optimization Translates to Factory Floor Benefits



Structured Pruning (The Throughput Booster)

- Achieves ~85ms inference on the Jetson Nano, comfortably exceeding the >10 FPS requirement for peak production.
- Reduced GPU load keeps the Jetson Nano within its 10W power limit, eliminating the need for cooling fans that could clog with factory dust.



Quantization (The 'Reaction Time' Accelerator)

- Reduces inference time further, from ~85ms to ~60ms. This guarantees the PLC receives the signal in time for the pneumatic gate to actuate before the vehicle passes.
- 4x reduction in memory bandwidth and energy per operation lowers the Jetson's power draw to ~7W, increasing reliability during power fluctuations.



The Critical Balance: Avoiding the Risks of Over-Compression

Increased Misclassification

Pruning beyond 70% can lead to errors.

Consequence: False Negatives and False Positives cause production delays and require costly manual intervention.

Loss of Robustness & Generalization

The model may become brittle.

Consequence: It may fail to generalize to normal production variations like lighting changes, camera vibrations, or the introduction of new vehicle models.

THANK YOU

