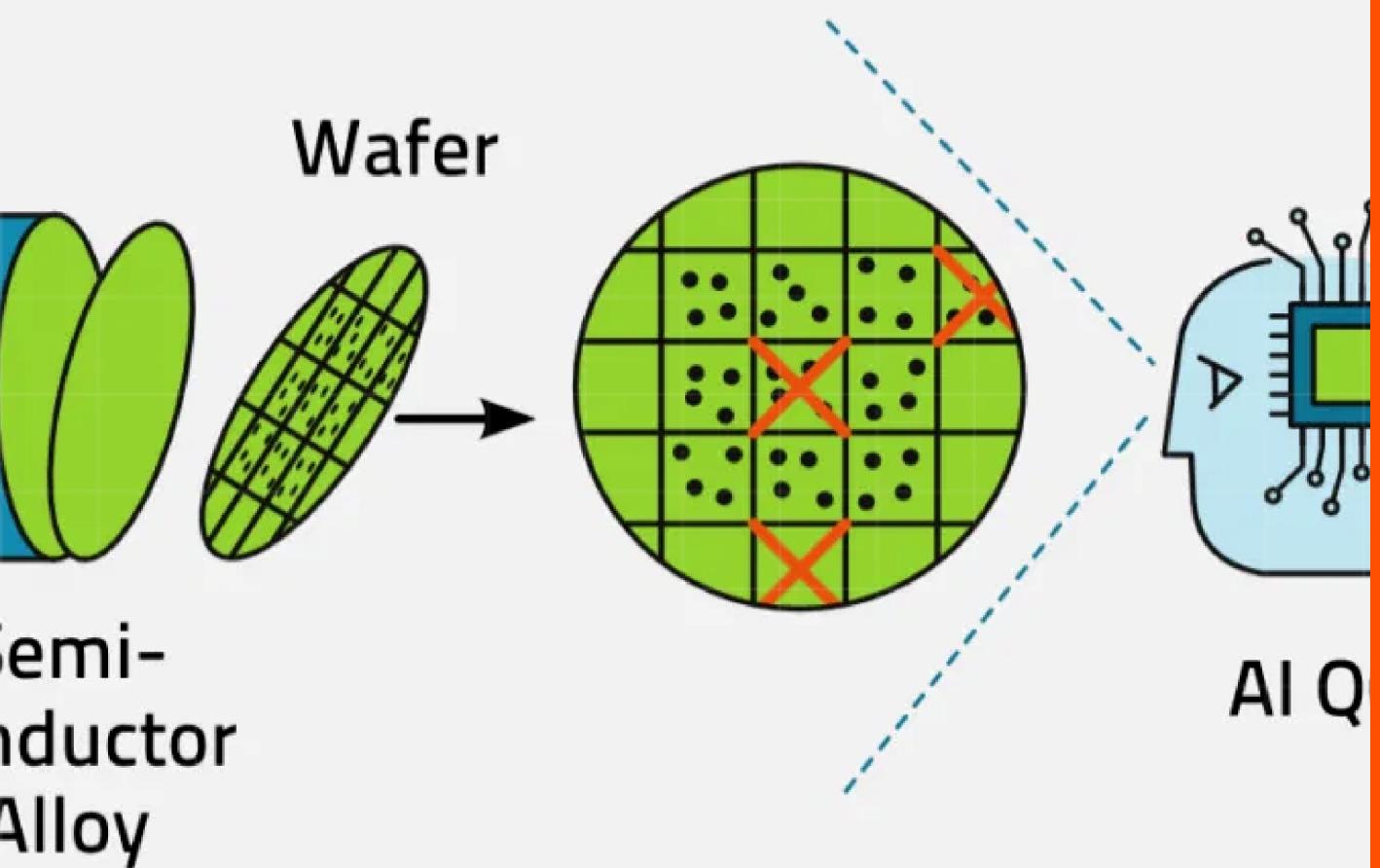


Semiconductor QC via AI



AI-DRIVEN AUTOMATED WAFER DEFECT ANALYSIS SYSTEM

Maximize Yield · Minimize Cost · Accelerate Root Cause Analysis

99%+ Defect Detection Accuracy

1000+ Wafers/Hour Throughput

100% Automated Continuous Learning

MANUFACTURING CHALLENGES IN WAFER DEFECT ANALYSIS

Current Industry Bottlenecks



High Labor Costs

Expert engineers manually review defect maps for every wafer. Time-consuming and error-prone process for thousands of wafers daily.

~\$150K / Year Avg Cost



Slow Analysis Speed

Traditional methods require hours to complete root cause analysis, causing significant production delays and capacity loss.

2-4 Hours / Analysis Cycle



Inconsistent Identification

Subjective classification by different engineers leads to high data inconsistency, impacting process improvement reliability.

High Data Inconsistency

BUSINESS IMPACT

⬇ Yield Loss (\$500K-\$1M/Year)

⌚ Capacity Drop 15-20%

✖ Inaccurate RCA

END-TO-END AI AUTOMATION SOLUTION

⚡ Real-Time Automated Analysis

Complete defect detection, classification, and RCA in <5 seconds after wafer entry. Zero human intervention required.

🌀 99%+ Accuracy

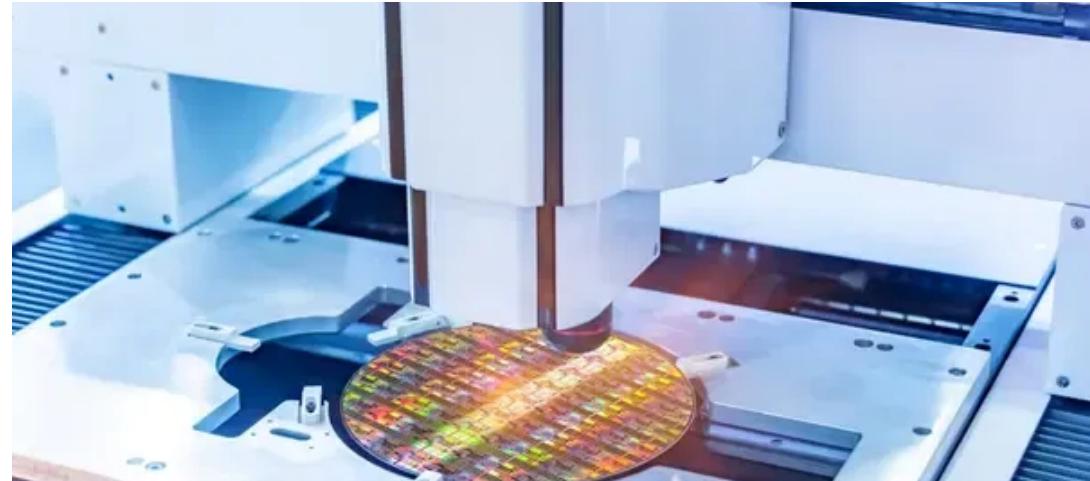
Utilizing state-of-the-art deep learning models (YOLOv10 + Vision Transformer), achieving 99.1% accuracy on WM-811K benchmark.

⌚ Continuous Automated Learning

System automatically collects expert feedback and uses GAN to generate synthetic training data for weekly auto-retraining.

≡ Scalable Deployment

Kubernetes-based microservices architecture supports horizontal scaling, handling 1000+ wafers/hour throughput.



EXPECTED ROI

50%

LABOR REDUCTION

2-3%

YIELD IMPROVEMENT

<5 Sec

ANALYSIS TIME

12-18 Mo

PAYBACK PERIOD

MICROSERVICES ARCHITECTURE ENABLES HIGH AVAILABILITY AND PERFORMANCE

APPLICATION SERVICES LAYER

API Gateway

Kong / Nginx

Data Ingestion

Python FastAPI

Preprocessing

CUDA / OpenCV

Inference Service

NVIDIA Triton

GAN Service

PyTorch StyleGAN2

MLOps Pipeline

Kubeflow / MLflow

INFRASTRUCTURE & DATA LAYER

 Kubernetes Cluster

 2x NVIDIA H100

 PostgreSQL + MinIO

 RabbitMQ

<50ms

INFERENCE LATENCY

1000+

WAFERS / HOUR

99.5%

SYSTEM UPTIME

95%

GPU UTILIZATION

STATE-OF-THE-ART DEEP LEARNING MODEL STACK

YOLOv10-Medium

98.5% Accuracy

8-12ms Latency

Real-time detection with precise bounding box localization for all defect types.

DeiT-Tiny

96.2% Accuracy

5-8ms Latency

Vision Transformer architecture. Data-efficient and handles class imbalance exceptionally well.

ResNet50 + CBAM

96.96% Accuracy

Attention Mechanism

Focuses on relevant defect regions for confidence score refinement.

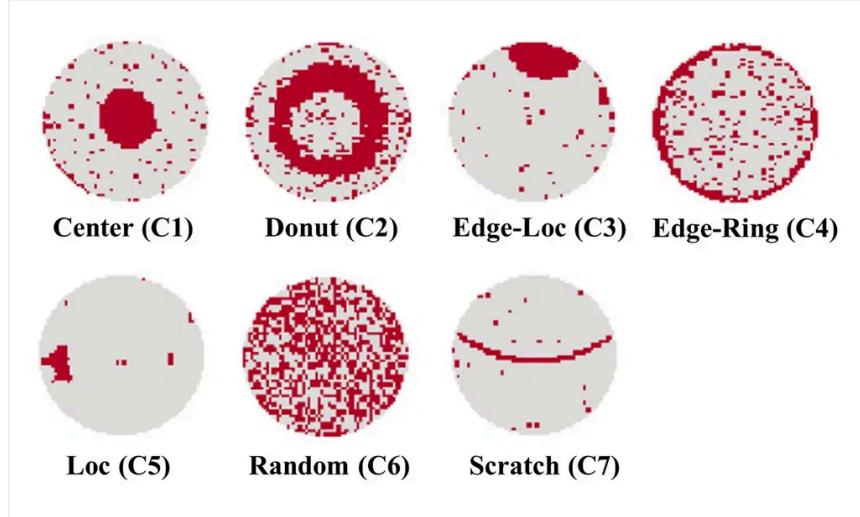
TensorRT Optimization

Speedup: **2.5x**

Memory Reduction: **40%**

Power Saving: **30%**

DEFECT DETECTION



CLASSIFICATION

ENSEMBLE VOTING

ENSEMBLE PERFORMANCE

Combined Accuracy

99.1%

Total Latency

<50ms

Model Size

~200MB

GAN-BASED SYNTHETIC DATA GENERATION

Solving Data Scarcity with Generative AI

GENERATIVE METHODS

StyleGAN2 Architecture

Generates high-fidelity, physically plausible synthetic defect images to expand the training dataset.

Conditional GAN (cGAN)

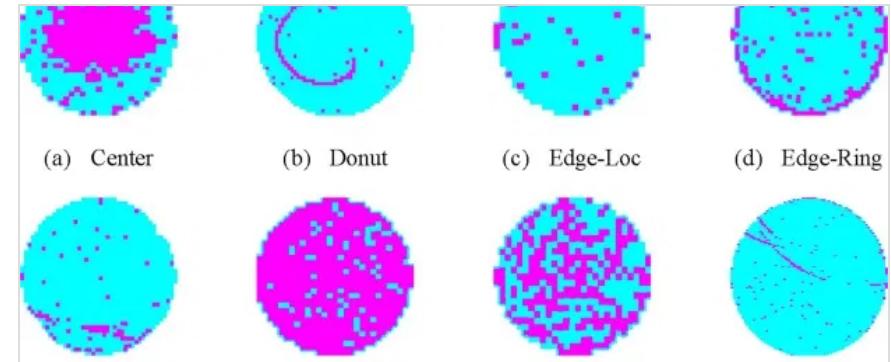
Targeted generation of specific defect types to augment minority classes and balance the dataset.

CycleGAN

Unpaired image-to-image translation enabling temporal data augmentation and domain adaptation.

AUGMENTATION STRATEGIES

- ✓ Geometric Transformations
- ✓ Defect-Aware Augmentation
- ✓ Mixup / CutMix
- ✓ SMOTE for Spatial Data



PERFORMANCE IMPACT

10x

Training Data Expansion

+3%

Accuracy Gain
(96% → 99%)

+25%

Minority Class F1 Score

< 5

FID Score
(High Quality)

CONTINUOUS LEARNING AND MODEL IMPROVEMENT

Automated MLOps Pipeline Enables Continuous Evolution

AUTOMATED LEARNING CYCLE



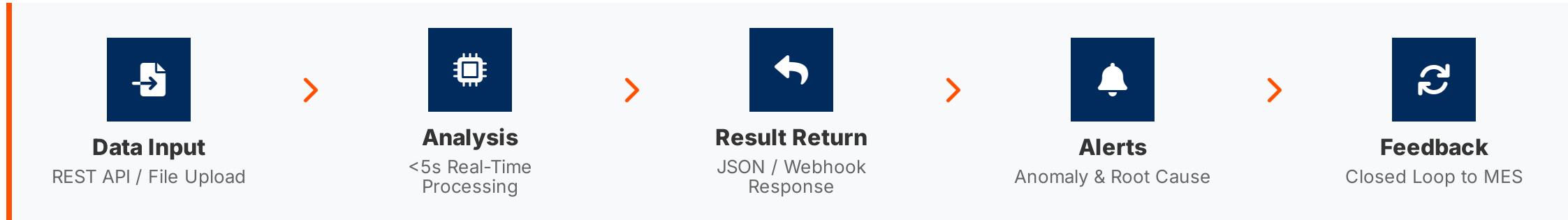
Performance Monitoring

Real-time Accuracy Target	$\geq 98\%$
False Alarm Rate	< 5%
Inference Latency (P99)	< 50ms
Model Drift Detection	Automated

Expected Benefits

- 👉 **Continuous Improvement:** Expect 0.5-1% model performance gain per month.
- ⚡ **Rapid Adaptation:** Time to recognize new defect types reduced from weeks to days.
- ✅ **High ROI on Feedback:** Every expert correction directly contributes to model accuracy.

REAL-TIME INFERENCE AND MES INTEGRATION

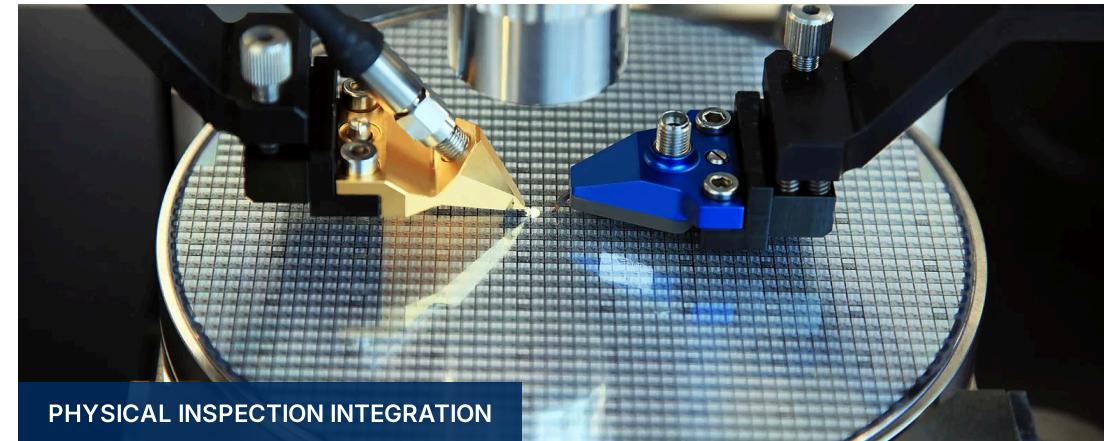


API Specification

- RESTful Design, OpenAPI Compliant
- JWT Token Authentication
- High Throughput: 1000+ Requests/Min

Alert Mechanism

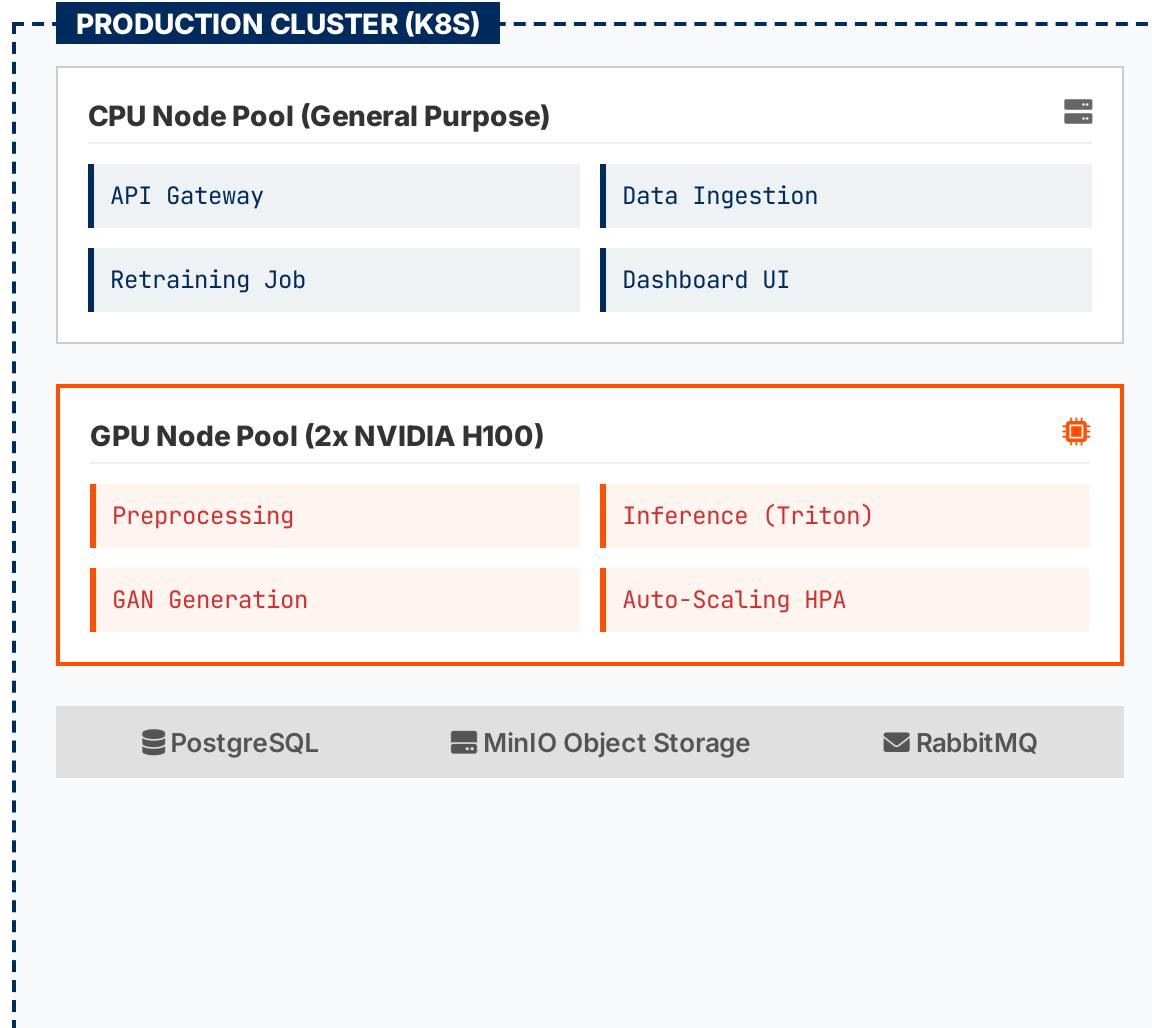
- Real-Time Anomaly Detection
- Automatic Root Cause Identification
- Configurable Thresholds & Escalation



PHYSICAL INSPECTION INTEGRATION

KUBERNETES-NATIVE ARCHITECTURE

Cloud-Native Deployment for Unlimited Scalability



Auto-Scaling Strategy

↔ Horizontal Scaling

Auto-scale Pods based on CPU/GPU util & queue depth.

↑ Vertical Scaling

Dynamic provisioning of GPU nodes for peak loads.

❤ Self-Healing

Auto-restart failed Pods; <30s Recovery Time (RTO).

1000+

WAFERS / HOUR

99.5%

AVAILABILITY

+500

WPH PER H100

95%

GPU UTILIZATION

Cost Optimization

Tiered storage lifecycle management and on-demand node scaling prevent over-provisioning.

8-MONTH PHASED DELIVERY ROADMAP

Structured Implementation Plan with Clear Milestones



QUANTIFIABLE SUCCESS METRICS AND KPIS

Measurable Business and Technical Outcomes

TECHNICAL PERFORMANCE

METRIC	TARGET	BASELINE	IMPROVEMENT
Defect Accuracy	≥ 99%	85-90%	+9-14%
False Alarm Rate	< 5%	10-15%	-5-10%
Inference Latency	≤ 50ms	2-4 Hours	99.9%
Throughput	≥ 1000 WPH	100-200	5-10x
Availability	≥ 99.5%	95%	+4.5%

BUSINESS IMPACT

MANUAL REVIEW

-50%

Save \$75K-\$150K labor cost annually

YIELD IMPROVEMENT

2-3%

Generate \$1M-\$1.5M additional revenue

DELAY REDUCTION

80%

Improve capacity utilization by 15-20%

PAYBACK PERIOD

12-18 Mo

Profitability begins in Year 2

USER SATISFACTION TARGETS

4.5/5 Engineer Satisfaction

4.5/5 MES Integration

4.5/5 Management ROI

RISK MANAGEMENT AND MITIGATION STRATEGIES

Proactive Identification and Control

TECHNICAL

Model Performance Drift

MITIGATION

Continuous learning pipeline, A/B testing in shadow mode, and automatic rollback mechanisms.

Integration Complexity

MITIGATION

Standardized REST APIs, pre-built adapter tools, and dedicated integration support team.

OPERATIONAL

User Adoption

MITIGATION

Comprehensive training programs, phased rollout strategy, and user feedback loops.

System Failure

MITIGATION

Redundant deployment ($N+1$), automatic failover, and real-time monitoring alerts.

BUSINESS

Implementation Delay

MITIGATION

Agile development methodology, strict milestone checkpoints, and risk monitoring.

Market Changes

MITIGATION

Modular architecture allows rapid adaptation to new manufacturing processes.

RISK MONITORING



Weekly Risk Assessments



Real-time KPI Dashboard



Emergency Response Plans

COMPETITIVE ADVANTAGES AND DIFFERENTIATION

Technical Leadership

Superior performance metrics compared to traditional AOI and competitor solutions.

99.1% Accuracy <50ms Latency Auto-Learning

Cost Effectiveness

Rapid ROI through significant labor reduction and yield improvement.

12-18 Mo Payback \$1.5M Savings/Year

Easy Integration

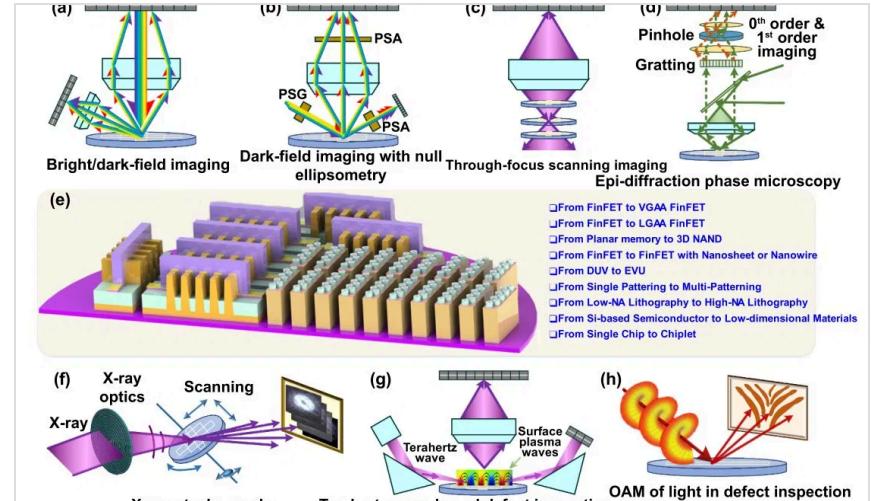
Seamless deployment into existing manufacturing environments without workflow disruption.

Standard REST API Multi-MES Support

Long-Term Value

Future-proof architecture that evolves with your manufacturing process.

Automated MLOps Continuous Improvement



VENDOR ADVANTAGES

- ✓ Professional AI/ML Team
- ✓ 24/7 Technical Support
- ✓ Proven Industry Track Record
- ✓ Regular Updates & Upgrades

INVESTMENT DECISION AND NEXT STEPS

INITIAL INVESTMENT
\$500K-800K
Dev & Deployment

ANNUAL OPEX
\$200K-300K
GPU & Maintenance

YEAR 1 ROI
\$1M-1.5M
Yield & Labor Savings

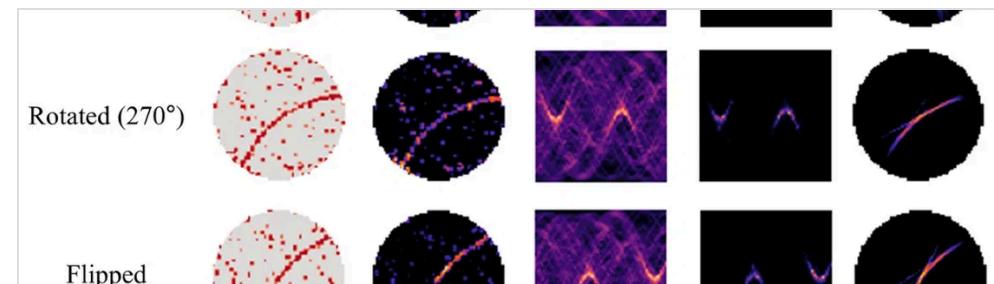
3-YEAR BENEFIT
\$3M-4.5M
Cumulative Value

DECISION CRITERIA MET

- ✓ Technical Feasibility Verified
- ✓ Risk Mitigation Complete
- ✓ Business Value > 100% ROI
- ✓ Timeline Reasonable (8 Mo)

RECOMMENDED NEXT STEPS

- 1 Executive Committee Approval This Month
- 2 Project Team Assembly Month 1
- 3 Pilot Deployment (Line 1) Month 6
- 4 Full Production Rollout Month 8



REQUIRED RESOURCES

- AI/ML Development Team (6-8 FTE)
- GPU Infrastructure (2x NVIDIA H100)
- Manufacturing Engineering Support (2-3 FTE)

RECOMMENDED ACTION: APPROVE PROJECT INITIATION