



THE EXECUTIVE PLAYBOOK FOR PHYSICAL AI DEPLOYMENT IN 2026

Paul Schell, Senior Analyst



EXECUTIVE SUMMARY

CONTENTS

Executive Summary	1
Market Context	2
Successful Use Cases Today	2
Business Needs: Realizing ROI from AI	3
Hybrid Cloud & On-Premises: Paving the Way for Physical AI Adoption in Manufacturing	3
The Production Gap: How to Avoid Deployment Failures.....	4
Successful Deployment of Visual AI: How Roboflow Addresses Industry Challenges	6
Imperatives for Deployment: A 30-60-90 Day Playbook	7
Conclusion	7

- **The biggest barrier to Return on Investment (ROI) is not model quality alone, but the production gap**, where promising industrial pilots fail to scale because of poor data quality, unresolved data pipelines, and weak alignment between business goals, infrastructure, and operating processes.
- **Successful industrial deployments require a systems-engineering approach.** Balancing edge, cloud, and hybrid architectures based on latency, compute, connectivity, and integration needs, while ensuring physical infrastructure (e.g., compute, cameras, lighting, and Programmable Logic Controller (PLC) connectivity) can support increasingly advanced models such as Vision-Language Models (VLMs) and Vision-Language-Action Models (VLAs).
- **Visual Artificial Intelligence (AI) in industrial and manufacturing markets is at an inflection point in 2026**, as the industry shifts from Convolutional Neural Network (CNN)-based models toward transformer architectures that now offer competitive or superior speed and accuracy for real-time deployments, opening the door to more capable and scalable production systems.

- **The proven path to enterprise scale** is a narrow first deployment: a single line, a single use case, and a 90-day window to produce validated ROI, turning one proven success into a repeatable template, rather than a portfolio of stalled pilots.
- **Intuitive end-to-end platforms enable this** by reducing talent bottlenecks by providing interoperability, hardware-agnostic deployment and, ultimately, speed to production.

MARKET CONTEXT

The visual AI market is undergoing a significant evolution. For most of the past decade, CNN architectures have been the de facto standard for real-time visual AI environments due to their speed, simple deployment, and mature ecosystem. Now, transformers are changing the math for AI ROI.

These models have reduced in size and become faster and more capable. At the same time, the required hardware has shrunk in size, cost, and power dissipation, all while becoming capable of running larger models. This combination unlocks countless new use cases in industrial and manufacturing settings, as well as high-throughput cloud environments.

SUCCESSFUL USE CASES TODAY

The discovery of viable, scalable use cases typically follows a framework to ascertain whether: 1) the manufacturing process occurs at a speed and scale that makes it unfeasible for human operators to handle all of the information; 2) there is room for human error or subjectivity to which AI is not prone; and, ultimately, 3) if actions taken on the data can lead to cost or efficiency savings. Today's solutions go beyond basic functionality of the past to deliver:



Inspection: Scanning the condition of equipment or manufactured products for their condition or specifications to detect faults or issues, even before they occur. This includes highly technical manufacturing domains such as semiconductor fabrication to reduce false-reject rates of wafers, leading to significant cost-savings.



Personal Protective Equipment (PPE) Compliance & Zone Safety Monitoring: Recognizing equipment such as hard hats, safety glasses, and high-visibility vests, alongside monitoring “red zones” to prevent collisions between humans and machinery. This can go beyond the traditional deployments by automating Occupational Safety and Health Administration (OSHA) compliance monitoring to reduce workplace injury liabilities and insurance premiums.



Robotic Bin Picking: Guiding robotic arms to identify, pick up, and orient jumbled, overlapping parts from a bulk bin to feed them into an assembly line. This can be implemented in zero-touch continuous assembly lines to guide robotic arms to untangle parts at up to 60 units per minute.

BUSINESS NEEDS: REALIZING ROI FROM AI

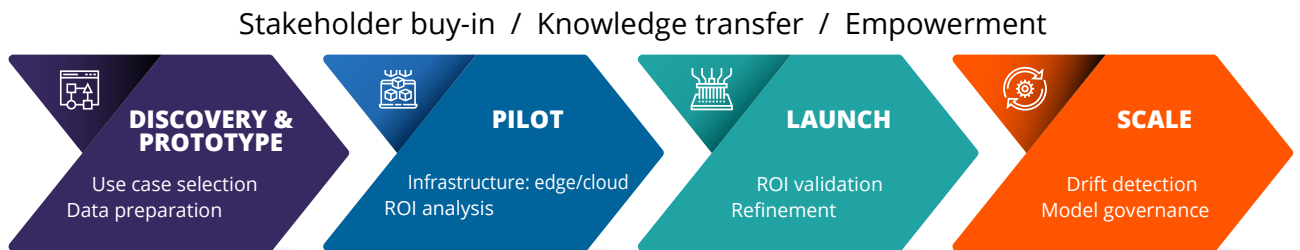
The rapid pace of innovation in vision AI has brought new functionality and efficiency savings—but these benefits are not spread evenly. ABI Research’s annual international Industrial and Manufacturing Survey has identified an increasing number of practitioners reporting skills shortages and complexity as key barriers to scaling transformative AI solutions. Decision makers must prioritize projects of higher value and lower complexity—these are the low-hanging fruit. Regardless of the use case, what matters to enterprises is realizing a return on their investment.

This journey starts with a clear understanding of which pain point or enterprise need should be solved or optimized, down to the specificity of individual manufacturing business units and their most impactful operations. Using this framework, and by identifying the correct, single first-use case, enterprises can uncover high-value vision AI engagements that can tackle tasks that are unproductive or even impossible for human operators—whether due to the volume of inbound data, the mission-criticality of that process, or the acute need for subjectivity and consistency.

These are all areas where vision AI can excel and improve efficiency, product quality, and ultimately, margins.

Figure 1: A Roadmap—Taking Vision AI to Scale

(Source: ABI Research)



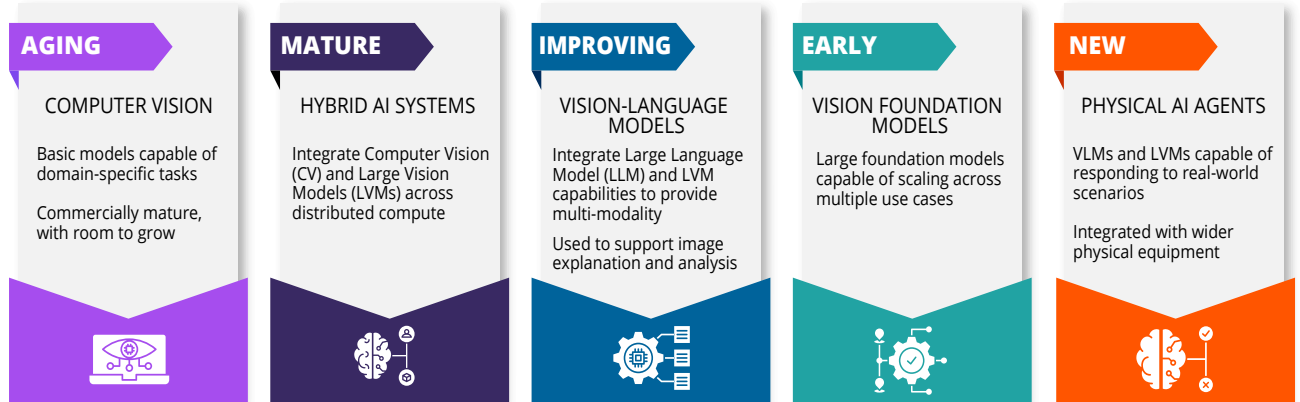
HYBRID CLOUD & ON-PREMISES: PAVING THE WAY FOR PHYSICAL AI ADOPTION IN MANUFACTURING

Vision AI is the perception layer that makes physical AI possible, and a key pillar and consideration for such projects is the compute infrastructure, whether at the edge where the data is created, in the cloud, or across both domains in a hybrid approach. Furthermore, manufacturing environments can impose constraints that neither a pure cloud nor edge architecture satisfy in isolation. Production lines require millisecond-level latency while also benefitting from the reasoning capabilities of larger VLMs for tasks such as root cause analysis, natural language querying of archives, and more.

The transition to more performant transformer-based models previously confined to the cloud is supported by numerous factors such as hardware innovation, performance improvements of smaller models, and use case development. This transition will result in a hybrid AI architecture leveraging several model types for different use cases, forming the basis for the next wave of physical AI—systems where vision models do not simply observe and classify, but also directly control actions in the manufacturing environment.

Figure 2: From Computer Vision to Physical AI Agents

(Source: ABI Research)



THE PRODUCTION GAP: HOW TO AVOID DEPLOYMENT FAILURES

The progress seen in the field of vision AI affects both the hardware and software domains. This includes the expansion of models serving diverse vision AI use cases in manufacturing settings, and the hardware innovations supporting these workloads at the edge (and training models on company-specific data in the cloud). In the context of skills shortages, restricted budgets and a need to see ROI sooner, below are key areas where practitioners fail, and where development platforms that abstract the inherent complexity of vision AI, lower Time to Market (TTM), and enable manufacturers ultimately lead to efficiency gains.

DATA CURATION

Industrial and manufacturing enterprises deploying visual AI continue to run into data challenges. Pilots may succeed when working with curated datasets, but struggle in production environments where data curation and other technical tasks like image labeling slow progress.

These friction points can be directly addressed by platforms using foundation models to accelerate data labeling. Access to open-source datasets to overcome data scarcity, built-in tools to synthetically expand underrepresented classes, and native dataset versioning ensure reproducibility from labeled image to the final deployed model in a production environment.

HARDWARE & SOFTWARE: A HOLISTIC VIEW OF VISION AI

Deploying vision AI, VLMs, and VLA models is more than a software engineering challenge and should be viewed holistically as a systems engineering project. Real ROI is achieved when organizations manage to go beyond AI as a cloud-based software pilot and start treating it as a physical industrial component. By ensuring that edge compute capabilities match model sizes, deployment settings are optimized for data collection, and AI outputs integrate seamlessly with legacy PLCs, industrial leaders can move from AI trials to profitable, at-scale production.

This raises important considerations across the major deployment locations:

- **Edge:** When low latency is mission critical, there are connectivity constraints and data must remain local, making deployment at the edge more suitable.
- **Cloud:** When latency needs are less strict, models require vast computing resources and there is a strong and stable network to maintain the flow of data from the source, so cloud deployments may be advantageous.
- **Hybrid:** Workloads can be partitioned between locations to leverage both realms' strengths, particularly if there is a need for background analysis and iterative training to improve model accuracy with new data, where on-premises hardware may lack the resources.

A platform that standardizes the deployment layer through containerization and unified Application Programming Interfaces (APIs) makes hardware location irrelevant to the software developer. This allows industrial facilities to distribute AI workloads in a hybrid approach, pushing low-latency, mission-critical responses to the edge and keeping heavy reasoning and training in the cloud. This maximizes both performance and ROI, and helps manufacturers where Information Technology (IT)/Operational Technology (OT) network boundaries may be a hurdle. This flexibility is often the difference between a project that scales and one that stalls waiting for infrastructure approvals.

VISION AI KNOW-HOW

Something felt more acutely at the edge is the need for expertise on the ground to deploy, tune, and maintain infrastructure that would otherwise be centralized in a cloud environment. A bottleneck in orchestrating complex hardware-software deployments in industrial settings is the shortage of specialized Computer Vision (CV) engineers and Machine Learning Operations (MLOps) talent.

Managing data pipelines, optimizing models for specific edge devices, and maintaining infrastructure often stalls pilot projects and inflates Total Cost of Ownership (TCO). End-to-end platforms with automated labeling (often powered by VLMs), push-button training, and edge-deployment tools reduce the underlying Machine Learning (ML) complexities. It is imperative that they cover the following functions and capabilities:

- **Low- or No-Code Abstraction:** Platforms that abstract the complexity from MLOps pipelines and enable non-specialist software developers, manufacturing engineers, and other domain experts to develop and maintain CV systems.
- **Automated Edge Compiler:** Quantizing and compiling models for different types of compute hardware at the edge, where performance optimizations are table stakes for vision AI deployments. This must be compatible with common frameworks such as TensorRT for NVIDIA Graphics Processing Units (GPUs).
- **Infrastructure Consolidation:** A united MLOps environment encompassing data curation, model training, and hybrid deployment scenarios—all under one roof.

Platforms with these features can negate the need for a dedicated team of ML PhDs, thereby unblocking deployments and accelerating time-to-ROI, allowing enterprises to develop and scale vision AI solutions with a drastically lower TCO.

SUCCESSFUL DEPLOYMENT OF VISUAL AI: HOW ROBOFLOW ADDRESSES INDUSTRY CHALLENGES

Given the complexities and common snags outlined above, some of the aspects with which industrial organizations typically need assistance, and for which specific solutions exist, are outlined below.

AGILE DATA CURATION, AUTOMATED LABELING AND TRAINING

To successfully deploy visual AI, teams must address the data bottlenecks before anything else. AI-assisted labeling and agile data curation have flipped the traditional workflow to accelerate labeling productivity. Pre-trained backbones such as DINOv2 carry broad knowledge into fine-tuning, meaning production-grade accuracy can be achieved with smaller datasets than CNN-based alternatives. Libraries such as Roboflow's Autodistill use large foundation models to auto-label training data for smaller models and allow human reviewers to correct and approve, rather than annotate from scratch.

In production, conditions shift across sites, and Roboflow addresses this through Roboflow Workflows, with automated pipelines built on a visual drag-and-drop canvas with over 40 pre-built blocks that chain together models, business logic, and external applications. Active learning runs inside Workflows, and the system selectively samples inference images and feeds corrected data into the next training cycle, creating a continuous improvement loop that keeps models current as production conditions change without the need for thousands of lines of custom code.

PLATFORM-AGNOSTIC AND PORTABLE DEPLOYMENT

Hardware agnosticism is important for industrial deployments, especially if it spans the entire vision AI pipeline. This prevents vendor lock-in and allows companies to optimize the size of their physical infrastructure based on their specific ROI and latency needs. Roboflow's platform offers the following key pillars to enable deployment:

- 1. Input Agnosticism:** This means data from an existing overhead Closed-Circuit Television (CCTV) security feed, a drone camera, a smartphone, or a Three-Dimensional (3D) spatial camera can be ingested, regardless of vendor or format.
- 2. Multi-Format Model Compilation:** A challenge in edge deployment is that different silicon chips require different optimizations. It is vital to be able to rapidly export into various optimized formats tailored for specific hardware accelerators, e.g., TensorRT for NVIDIA GPUs, OpenVINO for Intel Central Processing Units (CPUs) and Vision Processing Units (VPUs), and TensorFlow Lite/CoreML for lightweight deployment on mobile devices or Arm-based processors.
- 3. Containerized Edge Deployment:** Managing deployment through an open-source inference engine abstracts away the underlying operating system and hardware. This can then be deployed on various compute platforms.

Roboflow Inference can be used to deploy the same model across an edge device or a cloud instance without rewriting code. The engine automatically selects the best runtime (e.g., TensorRT or ONNX) for the hardware, closing the production gap with software, not manual engineering.

AI CENTER OF EXCELLENCE

Finally, bridging the skills gap is a challenge that is critical for manufacturing businesses. Organizations are addressing this through a Center of Excellence (CoE) model, with a central team that discovers, develops, and vets AI applications to deliver a faster ROI, thereby limiting the number of scattered pilot projects that never reach production. Roboflow offers services to help develop and scale global CoEs for visual AI to curate certified, repeatable projects they know work.

IMPERATIVES FOR DEPLOYMENT: A 30-60-90 DAY PLAYBOOK

Findings indicate that a disciplined 90 days focused on a single industrial vision AI use case, not a plant-wide rollout, are important for a successfully scaled solution further down the line. Even with an enabling developer platform, the identification of successful use cases is still an important part of the journey to successful deployments at scale. Prioritize high-value, low-complexity use cases that can start with models trained on fewer than 1,000 images.

FIRST 30 DAYS

Scope and catalog operations across 3-5 business units (e.g., raw material, manufacturing, Quality Control (QC), distribution) and plot candidates on a value-complexity matrix. Pick a single use case with a clearly visible signal: e.g., misalignment on a line, packaging damage at dispatch, or Stock Keeping Unit (SKU) verification at pick. Define success using Roboflow's V.I.S.U.A.L. framework: Verifiable visual target, measurable Impact in dollars, Scale in images per minute, Usability threshold for accuracy and latency, a specific Action triggered by the model, and a Lifecycle plan with measurable goals. Audit existing cameras and lighting on that one line and engage compliance, legal, and OT security for their approval before proceeding.

MIDDLE 30 DAYS

Prototype on a small dataset. This pilot can be done without thousands of hand-labeled images. Capture no more than 1,000 frames in real production conditions and use AI-assisted labeling: use foundation models like SAM 3 or Grounding DINO via Roboflow to prove viability without heavy data labeling. Train a lightweight custom model and target a realistic 70% to 80% accuracy to prove ROI for most defect and verification tasks. Run it in shadow mode alongside the existing process (no actions triggered) and produce a go/no-go memo with measured performance, projected annual savings, and a confidence range.

FINAL 30 DAYS

Deployment and the iteration loop: Move inference off the workstation and onto a single edge device (e.g., Jetson Orin or Hailo-8) on one line, wired into the PLC so the model triggers an actual reject, alert, or production line-stop. Instrument it to capture edge cases automatically: low-confidence frames, operator overrides, and downstream QC disagreements. Route those back through AI-assisted relabeling and human review, and retrain every 1 to 2 weeks, which can quickly deliver double-digit accuracy gains. End day 90 with a validated ROI figure, a documented runbook, and a go-decision for another production line. This is not an enterprise-wide rollout, but the next deployment line with all of the above lessons baked in.

CONCLUSION

The window for capturing competitive advantage from industrial vision AI is open and the manufacturers pulling ahead in 2026 are those running the right pilot, on a single line, with a clear ROI gate at day 90. The technology has progressed and matured: transformers are deployable at the edge, foundation models have collapsed labeling cycles, and hardware-agnostic platforms have removed most of the integration complexities that stalled the previous generation of projects. What remains is a question of discipline. Pick the use case where the signal is visible, the impact is clearly measurable, and the path from pilot to second line is short and repeatable. Scaled enterprise vision AI programs are built one validated deployment at a time.



Published May 2026

157 Columbus Avenue 4th Floor

New York, NY 10023

+1.516.624.2500

WE EMPOWER TECHNOLOGY INNOVATION AND STRATEGIC IMPLEMENTATION

ABI Research is uniquely positioned at the intersection of end-market companies and technology solution providers, serving as the bridge that seamlessly connects these two segments by driving successful technology implementations and delivering strategies that are proven to attract and retain customers.

©2026 ABI Research. Used by permission. ABI Research is an independent producer of market analysis and insight and this ABI Research product is the result of objective research by ABI Research staff at the time of data collection. The opinions of ABI Research or its analysts on any subject are continually revised based on the most current data available. The information contained herein has been obtained from sources believed to be reliable. ABI Research disclaims all warranties, express or implied, with respect to this research, including any warranties of merchantability or fitness for a particular purpose.