

ML LABS INTELLIGENCE

# How AI Systems Compound Operational Value

One-off AI wins are useful but fragile. This article shows how to design systems so each deployment improves the economics of the next one.

TECHNICAL

Author Omar Trejo

Date 2026-01-02

ML LABS

[mlabs.com/intel/compounding-effect-of-ai-systems](https://mlabs.com/intel/compounding-effect-of-ai-systems)

---

Most organizations treat AI projects as independent efforts — each starting from scratch with new data collection, model architecture, deployment pipeline, and monitoring. The team learns lessons, but the systems don't.

**Research on the state of AI adoption** (McKinsey, 2024) found that organizations where AI systems share data, infrastructure, and learned representations achieve more value per AI dollar — and the gap widens over time because each shared component reduces marginal cost while increasing quality baseline.

## Why AI Doesn't Compound by Default

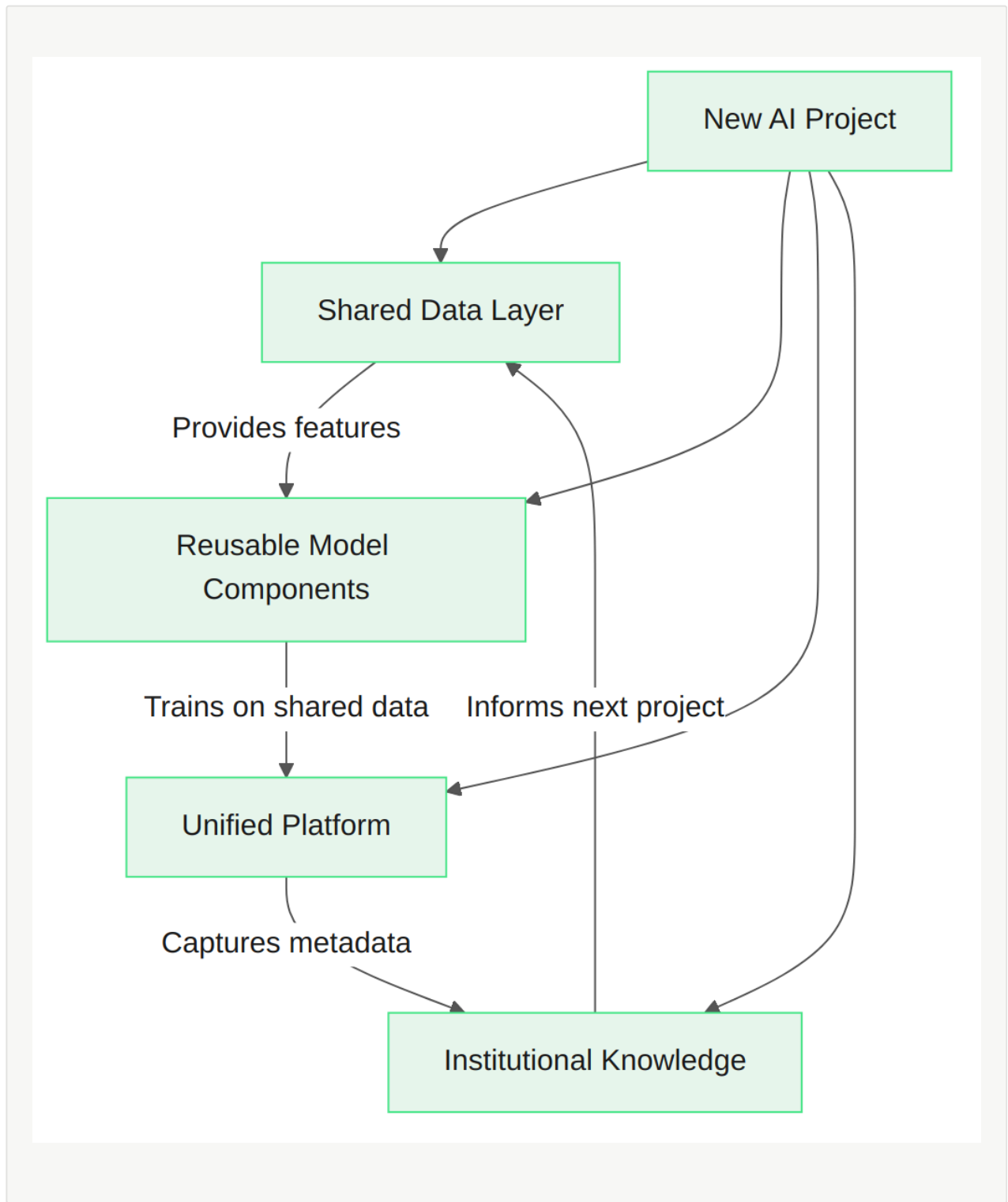
Three forces work against compounding: team incentives reward shipping over contributing to shared infrastructure, architecture fragmentation from each project choosing its own tools, and data silos from overlapping but unshared sources.

// The majority of technical debt in ML systems comes from infrastructure around models, not models themselves — and most of it is duplicated across projects because teams build in isolation.

**Research on ML technical debt** (NeurIPS, 2015) found that model code is a small fraction of real-world ML infrastructure — the surrounding plumbing dominates and is routinely rebuilt. **A shared ML platform initiative** (LinkedIn Engineering, 2019) showed the alternative: shared platforms with reusable feature schemas aimed to double ML engineer effectiveness. By the 10th project on a shared platform, marginal cost approaches a fraction of the original.

## The AI Compounding Flywheel

A compounding AI system has four components that reinforce each other. Each component generates value for all current and future AI projects in the organization.



## Component 1: Shared Data Layer

---

A feature store that every project contributes to and draws from. **One production ML platform** (Del Balso & Hermann, 2017) centralized ~10,000 shared features, eliminating duplicate engineering.

- **Feature registry.** Every feature documented, versioned, discoverable. New projects search before building
- **Feature computation pipeline.** Computed once, shared across all consuming models — no duplicate engineering per project
- **Data quality monitoring.** Automated drift, missing value, and distribution change detection across all features

Start with the 20 most-used features, mandate registry checks for new projects, and target 50%+ feature reuse.

## Component 2: Reusable Model Components

---

Three patterns drive reuse, each compounding across projects.

- **Pre-trained embeddings.** Domain-specific embeddings trained once and shared across projects. **Pre-trained representations** (Stanford NLP, 2014) transfer effectively across tasks
- **Shared encoders.** Text, image, and tabular encoders trained on your domain data become reusable components that transfer to any downstream task
- **Transfer learning infrastructure.** Fine-tuning from a related model is **an order of magnitude faster** (Proceedings of the IEEE, 2021) than training from scratch

## Component 3: Unified Platform

---

A unified platform eliminates per-project setup for experiment tracking, serving, and monitoring. Experiment tracking prevents teams from repeating each other's work. Unified serving handles deployment, scaling, and A/B testing for all models. Shared monitoring covers the entire model fleet with a single on-call rotation. The platform pays for itself after the third project.

## Component 4: Knowledge Base

---

Structured, searchable metadata turning individual project experience into organizational capability: decision logs, performance benchmarks, and failure postmortems. **A systematic lessons-learned database** (NASA, 2024) demonstrates that systematic capture of lessons produces measurable compounding returns even in domains far more complex than ML.

## Measuring Compounding

Three metrics indicate compounding. Flat trend lines mean you're building in isolation regardless of architecture diagrams: **time-to-production per project** (should decrease), **feature reuse rate** (target 50%+), and **cost per project** (should decrease as shared infrastructure absorbs fixed costs).

## Expected Results

Returns accelerate rather than plateau. Each project inherits a running start, so the same team ships more initiatives at higher model quality while maintenance stays low because shared infrastructure is maintained once rather than per-project. The

gap between organizations that compound and those that don't widens with every project shipped.

## The Incentive Problem

The biggest barrier isn't technical — it's organizational. The team building a reusable feature store doesn't get credit when three others benefit. Fix this:

- **Reward reuse.** Track shared component usage and credit builders
- **Reserve capacity.** Allocate 20% of time to shared infrastructure
- **Celebrate leverage.** When reuse accelerates a project, make it visible

## When This Approach Does Not Apply

Compounding requires a minimum threshold — organizations running fewer than 3-4 AI projects concurrently will find shared infrastructure costs more than duplication. The second failure condition is team isolation entrenched by organizational structure; mandating shared components without authority to enforce adoption creates more complexity than independence. Start with the lightest mechanism: a knowledge base where teams record decisions and lessons.

## First Steps

1. **Inventory shared features.** List features across projects. Those used by three or more are first candidates for centralization.

2. **Standardize serving infra.** Pick one model-serving platform. Efficiency gains fund the migration quickly.
3. **Start a knowledge base.** Run a structured retrospective after the next project ships and record it in searchable format.

## Practical Solution Pattern

Every AI initiative should leave behind shared assets — features, model components, deployment patterns, decision logs — that make the next one cheaper and faster. The difference between organizations that compound and those that don't is whether infrastructure is treated as a product that grows or as project scaffolding that gets thrown away.

Most ML engineering cost sits in the surrounding infrastructure, not the models. Without deliberate sharing, that cost resets to zero with every new project. Organizations that invest in reusable foundations create a flywheel: each project is faster, higher quality, and easier to maintain than the last. If you need ongoing technical ownership across several live AI systems, an **AI Engineering Retainer** keeps that flywheel turning.

## References

1. McKinsey & Company. **The State of AI.** *McKinsey & Company*, 2024.
2. Sculley, D., et al. **Hidden Technical Debt in Machine Learning Systems.** *NeurIPS*, 2015.
3. LinkedIn Engineering. **Scaling Machine Learning Productivity at LinkedIn.** *LinkedIn Engineering Blog*, 2019.

4. Uber Engineering. **Michelangelo: Uber's Machine Learning Platform.** *Uber Engineering Blog*, 2017.
5. Pennington, J., Socher, R., & Manning, C. D. **GloVe: Global Vectors for Word Representation.** *Stanford NLP*, 2014.
6. Zhuang, F., et al. **A Comprehensive Survey on Transfer Learning.** *Proceedings of the IEEE*, 2021.
7. NASA. **Lessons Learned Information System.** *NASA*, 2024.



AI SYSTEM 3

Exponential growth in value over time, achieved through advanced AI capabilities and data integration.

SYSTEM 2

Exponential growth in value over time, achieved through advanced AI capabilities and data integration.

EXPONENTIAL VALUE

EXPONENTIAL VALUE

EXPONENTIAL VALUE

ML LABS

Custom AI Systems for High-Value Workflows

[mllabs.com](http://mllabs.com)