

AI Engineering Effectiveness: A Predictive Framework for Measuring AI-Assisted Software Delivery Velocity

Kunal Kumar

Thales Group

Abstract: *This paper presents the AI Engineering Effectiveness Framework (AEEF), a structured model for assessing and predicting AI-assisted software delivery performance. The framework argues that delivery velocity is influenced by multiple interacting dimensions rather than AI capability alone. Seven factor groups are incorporated: Human, Technical, Process, Organizational, AI, Domain, and Environmental factors. The model uses weighted scoring to estimate delivery effectiveness during early adoption stages and evolves toward regression-based prediction as historical data become available. The framework provides a practical mechanism for identifying delivery bottlenecks, comparing teams, and improving decision-making through explainable metrics. The approach supports evidence-based AI adoption and offers a pathway from anecdotal productivity claims to measurable delivery outcomes.*

Keywords: AI engineering effectiveness, software delivery velocity, AI-assisted development, engineering productivity, software engineering metrics

1. Introduction

AI-assisted software development is frequently discussed through anecdotal evidence, isolated productivity reports, or tool-centered benchmarks. However, software delivery outcomes are shaped by a broader set of engineering conditions, including technical quality, process maturity, organizational constraints, and the availability of reliable delivery environments. A tool that performs well in a controlled benchmark may still generate inconsistent value in production if the surrounding engineering system is unstable.

This paper proposes the AI Engineering Effectiveness Framework (AEEF) as a structured method for evaluating AI-assisted software delivery in real organizational settings. The framework is based on the view that engineering velocity is influenced by multiple interacting factors and that AI should therefore be treated as one part of a larger delivery system rather than as a standalone explanation for productivity improvement.

The paper has three objectives. First, it defines a factor-based framework for assessing AI-assisted delivery readiness. Second, it introduces a weighted scoring model that organizations can use when historical data are limited. Third, it outlines a path from heuristic estimation toward regression-based prediction as delivery data mature.

2. Literature Survey

Recent work on engineering productivity measurement emphasizes that developer effectiveness cannot be represented adequately by a single output metric. Holistic frameworks such as SPACE highlight the importance of combining performance, communication, collaboration, and efficiency of flow when evaluating software engineering productivity. This multidimensional perspective is directly relevant to AI-assisted delivery because gains from AI depend not only on coding speed but also on trust, review quality, system reliability, and workflow continuity.

Similarly, predictive analytics literature shows that future outcomes can be estimated more reliably when historical data are combined with structured statistical models. Regression methods remain widely used for estimation and forecasting tasks because they provide interpretable relationships between input variables and observed outcomes. These ideas support the proposed AEEF progression from weighted assessment to data-driven calibration.

Existing industry discussion on AI-assisted engineering often focuses on coding acceleration, automated test generation, or documentation support. While useful, these observations frequently remain tool-specific and do not fully account for surrounding delivery conditions. The present framework extends this conversation by combining AI factors with software engineering, organizational, and environmental variables in a single predictive structure.

3. Problem Definition

Organizations adopting AI for software delivery often begin by asking which tool provides the greatest productivity improvement. This question is incomplete because productivity is influenced not only by AI capability, but also by requirements stability, technical debt, architecture maturity, test automation, review delays, interruptions, dependency complexity, and decision latency.

A second challenge is the lack of a repeatable method for connecting AI usage with measurable delivery outcomes. Teams often report perceived time savings or faster implementation, but these claims may not translate directly into reduced cycle time, increased throughput, or improved predictability. Without a structured framework, it becomes difficult to compare teams, identify bottlenecks, or estimate expected gains from AI adoption.

The problem addressed in this paper is therefore the absence of a practical, explainable, and scalable model for measuring

AI-assisted delivery effectiveness across real engineering environments.

4. Methodology and Framework Design

4.1 Conceptual Model

The AI Engineering Effectiveness Framework models delivery velocity as a function of seven categories:

- Human factors (H)
- Technical factors (T)
- Process factors (P)
- Organizational factors (O)
- AI factors (A)
- Domain factors (D)
- Environmental factors (E)

The conceptual relationship may be written as:

$$Velocity = f(H, T, P, O, A, D, E)$$

This structure positions AI as a major contributor, but not the sole determinant, of delivery performance.

4.2 Factor Definitions

The factor groups are defined as follows:

- Human factors: experience, AI workflow skill, team stability, and ability to validate AI output.
- Technical factors: codebase complexity, technical debt, architecture maturity, test automation, tooling quality, and environment stability.
- Process factors: requirements stability, planning quality, review process, deployment process, and SDLC maturity.
- Organizational factors: meeting load, interruptions, decision latency, dependency complexity, and leadership support.
- AI factors: prompt quality, context quality, acceptance rate, AI coverage, cost efficiency, and tool maturity.
- Domain factors: business complexity, regulatory complexity, integration complexity, and data complexity.
- Environmental factors: infrastructure stability, build efficiency, test environment availability, and incident load.

These dimensions collectively describe whether a team is positioned to convert AI capability into reliable software delivery improvement.

4.3 Score Normalization

The proposed method begins by collecting raw observations for each team, story, or epic across the seven factor groups. These observations are then converted into normalized scores on a 0 to 100 scale.

A practical normalization approach is:

- Direct percentage measures such as build success rate or automation coverage may be used directly after appropriate capping.

- Negative indicators such as review delay or incident frequency may be inverted so that higher scores always indicate more favorable conditions.
- Qualitative observations such as prompt maturity or leadership support may be converted through a rubric with clearly defined score bands.

An example of score normalization is shown below.

Raw observation	Interpretation	Normalized score
Build success rate = 92%	High delivery stability	92
Average review turnaround = 48 hours	Moderate process friction	60
Prompt maturity = defined and reusable	Strong prompt discipline	80

4.4 Weighted Scoring

In the early phase, AEEF is computed using assigned weights that reflect organizational priorities. An example starting configuration is shown below.

Factor	Weight
Human	0.20
Technical	0.15
Process	0.15
Organizational	0.10
AI	0.25
Domain	0.05
Environmental	0.10

The composite score is calculated as:

$$AEEF = 0.20H + 0.15T + 0.15P + 0.10O + 0.25A + 0.05D + 0.10E$$

A worked example is shown using the following scores: (H = 80), (T = 70), (P = 75), (O = 60), (A = 85), (D = 65), and (E = 80). The resulting AEEF value is 76.25.

4.5 Velocity Gain Estimation

In its initial form, AEEF serves as a predictor of expected velocity gain using the following relationship:

$$\hat{V}G = 1 + \alpha \times \frac{AEEF}{100} \times M$$

where α is a scaling coefficient and M is a work-mode multiplier. The multiplier reflects the reality that AI gains vary by task type. Documentation and unit-test generation are generally more compressible than architecture-heavy design or legacy refactoring.

A representative work-mode scale may be defined as follows:

Work type	Multiplier M
Documentation	0.90
Simple API development	0.80
Unit test generation	0.75
Existing code enhancement	0.60
Legacy refactoring	0.50
Distributed workflow	0.45
Architecture-heavy work	0.35

5. Results and Discussion

The main result of the framework is not a fixed benchmark, but an explainable measurement system that helps organizations estimate and interpret AI-assisted delivery performance. The AEEF score summarizes multiple determinants of engineering effectiveness into a single metric while preserving visibility into the underlying factor groups.

This allows several practical forms of analysis:

- Teams with similar AI tooling may still differ significantly in expected delivery gain because of process quality, architecture maturity, or environmental stability.
- Low measured benefit from AI does not necessarily imply weak model performance; it may indicate bottlenecks in requirements, review workflows, or infrastructure.
- Product and engineering leadership can use the factor scores to target interventions more precisely, such as improving context quality, reducing dependency delays, or strengthening test automation.

The framework is also useful because it supports progressive maturity. In the early phase, weighted scoring offers immediate operational value. In later stages, prediction quality can improve as organizations compare expected and actual outcomes and calibrate the model using empirical data.

6. Statistical Evolution

Once sufficient delivery data have been collected, the framework can transition from assumed weighting to learned estimation. A regression-based form may be written as:

$$VG = \beta_0 + \beta_1H + \beta_2T + \beta_3P + \beta_4O + \beta_5A + \beta_6D + \beta_7E + \beta_8M + \epsilon$$

Where β_0 through β_8 represents the learned coefficients for each factor and ϵ represents the error term

8. Summary

This formulation enables the model to learn from historical delivery outcomes rather than relying entirely on expert judgment. It also supports quantitative evaluation of which factor groups contribute most strongly to realized velocity gain.

Model quality may be assessed using metrics such as R², residual error patterns, and agreement between predicted and observed outcomes. These measures help determine whether the model is sufficiently reliable for planning, comparison, and governance use.

7. Benefits, Limitations, and Future Scope

The framework provides several practical benefits:

- It offers a structured way to evaluate AI adoption beyond tool-specific anecdotes.
- It supports explainable measurement of software delivery readiness.
- It helps product and engineering teams identify the real sources of delivery friction.
- It creates a path from qualitative assessment to predictive modeling.

At the same time, some limitations should be recognized:

- Certain factors such as prompt quality, context quality, and human capability may be difficult to score consistently without clear rubrics.
- Some variables may be correlated, which can complicate regression interpretation.
- The definition of velocity gain must be carefully standardized to avoid inconsistent use across teams.

Future work should focus on empirical validation, factor standardization, cross-team benchmarking, and quality-adjusted productivity measures. Further extensions may also examine interaction effects among variables and domain-specific calibration models.

AEEF END-TO-END EXAMPLE: FROM TRAINING DATA TO PREDICTION
Objective: Predict Velocity Gain (VG) using AEEF factors

STEP 1: CAPTURE & SCORE FACTORS (0-100)
7 AEEF Factor Scores

Human	Skills, experience, AI fluency, team stability	80
Technical	Code quality, test automation, tech debt, architecture	70
Process	SOLC maturity, planning, reviews, requirement clarity	75
Organization	Meeting load, decision latency, dependencies, communication	60
AI	Prompt quality, context quality, acceptance rates, AI coverage	85
Domain	Business complexity, integration complexity, regulatory complexity	65
Environment	Infra stability, build reliability, test env availability	80

STEP 2: CALCULATE AEEF

Weights (Initial):
w_H 0.20, w_T 0.15, w_P 0.15, w_O 0.10, w_A 0.25, w_D 0.05, w_E 0.10, Sum = 1.00

AEEF = 0.20H + 0.15T + 0.15P + 0.10O + 0.25A + 0.05D + 0.10E

Example Calculation (for a project):
AEEF = 0.20(80) + 0.15(70) + 0.15(75) + 0.10(60) + 0.25(85) + 0.05(65) + 0.10(80) = 76.25

STEP 3: TRAINING DATASET (Historical Projects)
Used to learn the relationship between AEEF and Actual Velocity Gain (VG)

Project ID	H	T	P	O	A	D	E	AEEF	Actual VG
P1	65	60	70	50	60	50	70	59.00	1.20
P2	80	70	75	60	85	65	80	73.25	1.70
P3	55	50	60	40	50	40	60	50.25	1.10
P4	90	80	85	70	90	70	85	81.75	2.20
P5	70	65	70	55	75	60	75	66.00	1.45
P6	85	75	80	65	85	65	85	76.75	1.90
P7	45	40	50	35	40	35	50	41.00	1.00
P8	95	85	90	75	95	75	90	87.75	2.35
P9	60	55	65	45	60	50	65	55.75	1.25
P10	75	70	75	60	80	65	80	71.25	1.65

STEP 4: TRAIN THE MODEL (Simple Linear Regression)
Model: VG = $\beta_0 + \beta_1 \times$ AEEF

Using Least Squares Regression the model learns the best β_0, β_1 that minimize the total error.

From the above data, we get:
 β_0 (Intercept) = 0.335
 β_1 (slope) = 0.0207

Trained Model:
VG = 0.335 + 0.0207 × AEEF

$R^2 = 0.92$ (Model explains 92% of the variation in VG)

STEP 5: PREDICTION FOR A NEW PROJECT
(Not part of training data)

New Project Scores: H: 70, T: 65, P: 70, O: 50, A: 80, D: 60, E: 75

Calculate AEEF for New Project:
AEEF = 0.20(70) + 0.15(65) + 0.15(70) + 0.10(50) + 0.25(80) + 0.05(60) + 0.10(75) = 66.75

Predict Velocity Gain (VG):
VG = 0.335 + 0.0207 × AEEF
= 0.335 + 0.0207 × 66.75
= 1.716
= 1.72x Velocity Gain

Interpretation: This project is expected to deliver work about 1.72 times faster with AI compared to without AI.

Convert VG to Effort (if needed):
If Baseline Effort (without AI) = 100 hours
Expected AI-Enabled Effort = 100 / 1.72 = 58.1 hours
Effort Reduction = 41.9%

STEP 6: MODEL EVALUATION (On Training Data)

Project ID	AEEF (X)	Actual VG (Y)	Predicted VG ($\hat{Y} = \beta_0 + \beta_1 X$)	Error (Y - \hat{Y})	Abs Error	Sq Error
P1	59.00	1.20	1.558	-0.358	0.358	0.127
P2	73.25	1.70	1.855	-0.155	0.155	0.024
P3	50.25	1.10	1.376	-0.276	0.276	0.076
...

$R^2 = 0.92$ (Good Fit) | MAE = 0.11 | RMSE = 0.14

STEP 7: AS MODEL EVOLVES

1. Add more projects to training data.
2. Recalculate AEEF and Actual VG.
3. Retrain the model → Updated β_0, β_1 .
4. Model becomes more accurate over time.

END-TO-END SUMMARY FLOW

Capture Raw Data (7 Factors) → Calculate AEEF (0-100) → Train Model (Regression) → Predict VG for New Projects → Use for Planning & Decisions.

- Step 1: Define the 7 AEET dimensions (Human, Technical, Process, Organization, AI, Domain, Environment).
 - Step 2: Define measurable sub-factors for each AEET dimension.
 - Step 3: Collect raw engineering, AI, process, and delivery metrics from projects.
 - Step 4: Normalize all collected metrics to a common 0–100 scoring scale.
 - Step 5: Calculate H, T, P, O, A, D, and E scores using weighted sub-factor aggregation.
 - Step 6: Calculate the overall AEEF score using weighted factor scores.
 - Step 7: Capture actual velocity gains from completed historical projects.
 - Step 8: Build a training dataset containing factor scores, AEEF, and actual velocity gain.
 - Step 9: Train a statistical model to learn the relationship between AEET factors and velocity gain.
 - Step 10: Validate model accuracy using historical project outcomes.
 - Step 11: Calculate H, T, P, O, A, D, and E for a new project.
 - Step 12: Calculate AEEF for the new project.
 - Step 13: Apply the trained model to predict expected velocity gain.
 - Step 14: Convert predicted velocity gain into effort, cost, and schedule estimates.
 - Step 15: Compare predictions with actual outcomes after project completion.
 - Step 16: Continuously retrain and refine the model as more project data becomes available.
- [4] IBM. 6 Ways to Enhance Developer Productivity with- and Beyond- AI. Available from: <https://www.ibm.com/think/insights/developer-productivity>
- [5] ACM Digital Library. Predictive Models in Software Engineering. Available from: <https://dl.acm.org/doi/10.1145/3503509>

9. Conclusion

The AI Engineering Effectiveness Framework provides a practical and explainable method for evaluating AI-assisted software delivery within the broader engineering system. By integrating Human, Technical, Process, Organizational, AI, Domain, and Environmental factors, the framework moves beyond anecdotal productivity claims and toward evidence-based assessment.

The weighted scoring model offers immediate usability in early adoption stages, while the transition to regression-based calibration supports increasing predictive accuracy as data maturity improves. The framework therefore provides a credible basis for engineering decision-making, product planning, and enterprise AI adoption strategy.

References

- [1] IBM. What is predictive analytics? IBM. Available from: <https://www.ibm.com/think/topics/predictive-analytics>
- [2] Microsoft. Developer experience. Microsoft. Available from: <https://developer.microsoft.com/en-us/developer-experience>
- [3] Octopus Deploy. The SPACE Metrics: A Holistic Measure of DevOps Productivity. Available from: <https://octopus.com/devops/metrics/space-framework/>