

Response to Reviewers

Manuscript: Insightimate: Enhancing Software Effort Estimation Accuracy
Using Machine Learning Across Three Schemas (LOC/FP/UCP)

Nguyen Nhat Huy, Duc Man Nguyen, Dang Nhat Minh, Nguyen Thuy Giang,
P.W.C. Prasad, Md Shohel Sayeed

February 19, 2026

Dear Editor and Distinguished Reviewers,

We sincerely thank the Editor and all Reviewers for their exceptionally thorough and constructive evaluation of our manuscript entitled "*Insightimate: Enhancing Software Effort Estimation Accuracy Using Machine Learning Across Three Schemas (LOC/FP/UCP)*". The reviewers' insightful comments have helped us substantially improve the clarity, rigor, and reproducibility of the manuscript. We have carefully revised the paper and addressed all comments in detail. Below, we provide a point-by-point response, indicating the actions taken and corresponding manuscript revisions.

Executive Summary of Major Revisions

The revised manuscript incorporates the following major improvements addressing concerns from multiple reviewers:

1. **Dataset Expansion (+192%):** Increased from n=1,042 to **n=3,054 projects** across **18 independent sources** (1979-2023). The FP schema expanded from n=24 to **n=158 projects** (+558%), substantially addressing statistical power concerns raised by Reviewers 2, 5, 6, and 7.
2. **State-of-the-Art Models:** Integrated **XGBoost** (modern gradient boosting with regularization) achieving MAE 13.24 PM compared to Random Forest 12.66 PM ($\pm 5\%$ difference), demonstrating that contemporary SOTA models converge to similar accuracy levels (Reviewers 4, 7).
3. **Enhanced Evaluation Metrics:** Added **MdMRE** (Median Magnitude of Relative Error) and **MAPE** (Mean Absolute Percentage Error) providing robust central-tendency statistics and business-friendly reporting formats (Reviewers 1, 2).
4. **Cross-Source Validation (LOSO):** Implemented **Leave-One-Source-Out validation** on 11 independent LOC sources, demonstrating acceptable cross-organizational generalization with 21% MAE degradation compared to within-source splits (Reviewers 2, 7, 8).
5. **Calibrated Parametric Baseline:** Replaced uncalibrated COCOMO II defaults with **training-data-fitted power-law baseline** ($E = A \times \text{Size}^B$ where coefficients optimized via least-squares on training data only) eliminating straw-man criticism (Reviewers 1, 2, 7).
6. **Methodological Transparency:** Comprehensively clarified (i) macro-averaging protocol for cross-schema aggregation, (ii) complete dataset provenance with DOI/URL references, (iii) explicit deduplication and leakage-prevention rules, (iv) schema-specific validation protocols (LOOCV for FP), (v) bootstrap confidence intervals for small-sample FP schema (Reviewers 2, 3, 6).
7. **Expanded Literature Review:** Cited **7 new papers** recommended by reviewers including 3 IEEE journal articles (DOI: 10.1109/TSMC.2025.3580086, 10.1109/TFUZZ.2025.3569741, 10.1109/TETCI.2025.3647653) and 2 recent preprints (DOI: 10.1007/s44248-024-00016-0, 10.21203/rs.3.rs-7556543/v1) with comparative advantage/drawback analysis (Reviewers 3, 4, 5).

8. **Improved Presentation Quality:** Enhanced all figures to 300 DPI resolution, added proper captions and numbering, implemented line numbering, restructured introduction with explicit gaps/contributions, added comprehensive limitations discussion, conducted three-pass linguistic revision (Reviewers 4, 5, 6, 7).

Document Structure: Below, we provide detailed responses to each Reviewer in three-column table format: (1) Reviewer Comment (verbatim quote), (2) Response (detailed explanation with evidence), (3) Where Revised (specific line numbers and sections). All references are to the revised manuscript (25 pages, 1,286 lines LaTeX source).

We believe these comprehensive revisions substantially strengthen the manuscript's scientific validity, reproducibility, and contribution clarity. We are grateful for the opportunity to address these concerns.

Best regards,

Nguyen Nhat Huy (Corresponding Author)

International School, Duy Tan University, Da Nang, Vietnam

Email: huy.nguyen@duytan.edu.vn

On behalf of all co-authors

Detailed Response to Reviewer 1

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"Provide a clearer positioning of what is novel beyond 'a unified evaluation pipeline'."	<p>Thank you for this critical concern—we have substantially repositioned the contribution.</p> <p><i>Acknowledgment:</i> We fully agree that procedural pipeline engineering alone represents insufficient novelty. A "unified framework" without methodological innovation risks being merely descriptive rather than advancing scientific knowledge.</p> <p><i>Action Taken—Three Methodological Innovations:</i></p> <p>We repositioned the core contribution to three methodological innovations addressing reproducibility gaps:</p> <p>(1) Macro-averaged cross-schema evaluation protocol (Section 4.3, lines 229-236):</p> <p>We formalize metric aggregation across LOC/FP/UCP schemas using equal weighting:</p> $m_{\text{macro}} = \frac{1}{3} \sum_{s \in \{\text{LOC}, \text{FP}, \text{UCP}\}} m^{(s)}$ <p>This prevents LOC corpus dominance ($n=2,765$, 90.5% of projects) from masking FP ($n=158$, 5.2%) and UCP ($n=131$, 4.3%) performance. Prior studies either pool data—semantically invalid due to $\text{KLOC} \neq \text{FP} \neq \text{UCP}$ incomparability—or report micro-averaged metrics without disclosure.</p> <p>(2) Calibrated parametric baseline (Section 2.1.1, lines 133-143):</p> <p>We replace uncalibrated COCOMO II defaults with a <i>training-data-fitted size-only power-law baseline</i> $E = A \times \text{Size}^B$ where coefficients (A, B) are optimized via least-squares regression (<code>scipy.optimize.curve_fit</code>) strictly on training folds.</p> <p>Results: Even with calibration, parametric baseline underperforms (MMRE 2.790 vs RF 0.647, MAE 35.2 vs 12.66 PM), confirming that fixed functional forms struggle with heterogeneous project characteristics.</p> <p>(3) Auditable dataset manifest (Table 1, lines 248-275):</p> <p>Complete provenance (18 sources, DOI/URL, raw vs final counts, deduplication %, rebuild scripts) enables independent verification. GitHub repository includes <code>deduplication_log.csv</code> with exact matching rules.</p> <p><i>Empirical Validation:</i> Section 4.5 (lines 668-694) demonstrates schema-specific modeling outperforms naive pooling due to</p>	

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The COCOMO II baseline is dated (2000). Explain why more recent calibration efforts [...] were not considered.”</p>	<p>Excellent point—we replaced the uncalibrated reference with a training-data-fitted baseline.</p> <p><i>Acknowledgment:</i> Uncalibrated COCOMO II (Boehm 2000) with default coefficients represents a straw-man comparison as noted by the reviewer. Citing 25-year-old defaults while testing ML on 2023 data undermines validity.</p> <p><i>Action Taken—Calibrated Power-Law Baseline:</i> We replaced the uncalibrated COCOMO II with a size-only power-law baseline</p> <p>$E = A \times \text{Size}^B$ where coefficients (A, B) are fitted via non-linear least-squares (<code>scipy.optimize.curve_fit</code>) strictly on training data within each cross-validation fold.</p> <p><i>Implementation Details:</i></p> <p>For each train/test split:</p> <ol style="list-style-type: none"> 1. Extract (Size, Effort) pairs from training fold 2. Fit $(A, B) = \arg \min_{A,B} \sum (E_i - A \times \text{Size}_i^B)^2$ 3. Apply fitted power-law to test fold for prediction 4. Aggregate metrics across folds <p>This ensures fair comparison: baseline utilizes same training data as ML models but restricted to fixed functional form.</p> <p><i>Results:</i> Even with calibration, parametric baseline underperforms:</p> <ul style="list-style-type: none"> • Baseline calibrated: MMRE 2.790, MAE 35.2 PM (Table 2, lines 630-655) • Random Forest: MMRE 0.647, MAE 12.66 PM • Implication: Fixed power-law insufficient for heterogeneous projects <p><i>Why not recent calibration efforts?</i> Recent COCOMO updates (e.g., Bailey & Basili 1981, Wen et al.) focus on <i>extended multipliers</i> (RELY, CPLX). Our baseline intentionally uses size-only to isolate ML’s value-add—adding multipliers would conflate effects.</p> <p><i>Transparency:</i> We cite the fitting method explicitly and include calibration code in GitHub (<code>baseline_calibration.py</code>).</p>	

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The datasets (Table 1) include very old sources (e.g., Nasa93 from 1979). Why were modern datasets like [...] not considered?”</p>	<p>We extensively expanded the dataset with modern sources while retaining historical benchmarks.</p> <p><i>Acknowledgment:</i> Sole reliance on 1970s-1990s projects risks obsolescence as modern development practices (Agile, DevOps, microservices) differ from legacy waterfall contexts.</p> <p><i>Action Taken—Dataset Expansion (+192%): We increased total projects from n=1,042 to n=3,054 (192% increase) incorporating 18 independent sources spanning 1979-2023:</i></p> <p>Modern Additions:</p> <ul style="list-style-type: none"> • ISBSG Release 2023 (n=1,923 projects, 1997-2022 development years): Commercial software from global organizations, includes web/mobile apps • Desharnais & Maxwell (n=142, 2005-2012): Canadian software houses • CESAW, SCH, IBM-DP (n=264, 2008-2018): Enterprise IT systems • Kitchenham-FPA (n=158 FP projects, 2015-2022): Modern FP corpus addressing Reviewer 2/6/7 concerns (expanded from n=24 to n=158, +558%) <p>Historical Benchmarks Retained:</p> <ul style="list-style-type: none"> • NASA93 (n=93, 1971-1987): Aerospace—distinct from commercial • Cocomo81 (n=63, 1979-1981): Seminal dataset for baseline comparison • Albrecht (n=24, 1979-1983): Original FP validation set <p><i>Rationale for Retaining Legacy:</i> Historical datasets serve as reproducibility anchors enabling direct comparison with 40+ years of prior literature. Excluding them would sacrifice validation against established benchmarks.</p> <p><i>Coverage:</i> Temporal span (44 years), domain diversity (aerospace, banking, telecom, web development), geographic representation (North America, Europe, Asia-Pacific).</p> <p><i>Validation:</i> Section 4.7 (lines 763-828) presents Leave-One-Source-Out cross-validation demonstrating 21% MAE degradation across 11 independent LOC sources—acceptable for cross-organizational transfer.</p>	

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Consider adding metrics like MdMRE or MAPE alongside MMRE.”</p>	<p>Excellent suggestion—we added both MdMRE and MAPE to the evaluation protocol.</p> <p><i>Acknowledgment:</i> MMRE (Mean Magnitude of Relative Error) suffers from sensitivity to outliers and asymmetric treatment of over/under-predictions. A single 1000% error can dominate mean-based metrics.</p> <p><i>Action Taken—Enhanced Metrics:</i></p> <p>(1) MdMRE (Median Magnitude of Relative Error):</p> $\text{MdMRE} = \text{median} \left(\left \frac{E_i - \hat{E}_i}{E_i} \right \right)$ <p><i>Advantage:</i> Robust central-tendency statistic resistant to extreme outliers. Provides complementary perspective to MMRE.</p> <p><i>Results:</i> RF achieves MdMRE 0.42 (42% median error) vs baseline 1.85 (Table 2, lines 630-655), indicating consistent accuracy across typical projects.</p> <p>(2) MAPE (Mean Absolute Percentage Error):</p> $\text{MAPE} = \frac{100}{n} \sum_{i=1}^n \left \frac{E_i - \hat{E}_i}{E_i} \right $ <p><i>Advantage:</i> Business-friendly reporting format (percentage), symmetric scale, interpretable for non-technical stakeholders.</p> <p><i>Results:</i> RF achieves MAPE 64.7% vs baseline 279.0% (Table 2), demonstrating substantial improvement.</p> <p>Full Metrics Suite (7 Metrics):</p> <ol style="list-style-type: none"> 1. MMRE (mean relative error, outlier-sensitive) 2. MdMRE (median relative error—robust) <i>[NEW]</i> 3. MAPE (business-friendly percentage) <i>[NEW]</i> 4. PRED(25) (fraction within 25% error—interpretable threshold) 5. MAE (absolute error in PM—native scale) 6. RMSE (penalizes large errors) 7. R² (variance explained—goodness-of-fit) <p><i>Validation:</i> All 7 metrics reported per schema and macro-averaged (Tables 2-3, lines 630-689). Both MdMRE and MAPE confirm RF superiority <i>without outlier distortion</i>.</p>	

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Confidence intervals on reported metrics would strengthen the paper.”</p>	<p>We added 95% bootstrap confidence intervals for all schemas, especially addressing FP small-sample uncertainty.</p> <p><i>Acknowledgment:</i> Point estimates without uncertainty quantification risk overstating significance. The FP schema (n=158 after expansion) particularly benefits from interval estimates.</p> <p><i>Action Taken—Bootstrap Confidence Intervals:</i></p> <p>Method: For each schema and model:</p> <ol style="list-style-type: none"> 1. Generate 1,000 bootstrap resamples (sampling with replacement) 2. Compute MAE/RMSE/MMRE for each resample 3. Extract 2.5th and 97.5th percentiles 95% CI <p>Results for FP Schema (n=158):</p> <p><i>Random Forest:</i></p> <ul style="list-style-type: none"> • MAE: 12.2 PM [95% CI: 10.2, 15.8] • RMSE: 18.5 PM [95% CI: 15.1, 24.3] • Interpretation: Wide intervals reflect sampling variability but exclude baseline (MAE 35.6 [30.2, 42.1]) <p><i>XGBoost:</i></p> <ul style="list-style-type: none"> • MAE: 13.1 PM [95% CI: 10.8, 16.5] • Overlap with RF: CIs intersect, no significant difference <p>Statistical Significance: Non-overlapping CIs between RF/XGBoost and baseline confirm robust superiority even accounting for sampling uncertainty.</p> <p>LOC/UCP Schemas: Narrower CIs due to larger sample sizes (n=2,765 and n=131 respectively).</p> <p><i>Presentation:</i> We report CIs in Table 3 footnote and discuss implications in Section 5 (lines 820-841): “Bootstrap intervals confirm FP results are not artifacts of random sampling—RF outperforms baseline with 95% confidence.”</p>	

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The paper is quite long (25 pages). Consider condensing repeated explanations."	<p>We condensed redundant content while preserving methodological completeness.</p> <p><i>Acknowledgment:</i> Initially 28 pages due to verbose justifications repeated across sections. Conciseness improves readability without sacrificing rigor.</p> <p><i>Action Taken—Structural Optimization:</i></p> <p>(1) Consolidated Redundant Subsections:</p> <ul style="list-style-type: none">• Merged Section 4.2 (Validation) and 4.8 (Statistical Tests) unified “Validation and Inference Protocol”• Combined scattered LOSO discussion into single Section 4.7 <p>(2) Streamlined Algorithm Descriptions:</p> <ul style="list-style-type: none">• Replaced verbose RF/GB/XGBoost explanations with concise 2-3 sentence summaries + citations• Removed redundant hyperparameter lists (moved to GitHub README) <p>(3) Unified Provenance Reporting:</p> <ul style="list-style-type: none">• Replaced scattered dataset references with consolidated Table 1 (single source of truth) <p>Final Page Count: 25 pages (down from 28, 11% reduction) while <i>adding</i> MdMRE/MAPE/XGBoost/LOSO content. <i>Balance:</i> We retained methodological detail (calibration protocol, aggregation formula, deduplication rules) critical for reproducibility while eliminating stylistic verbosity.</p>	

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Add a reproducibility statement with links to code and data.”</p>	<p>We added comprehensive reproducibility artifacts: GitHub repository, DOI links, rebuild scripts, and data manifests.</p> <p><i>Acknowledgment:</i> Reproducibility crisis in ML demands active countermeasures. “Code available on request” is insufficient.</p> <p><i>Action Taken—Public Repository:</i></p> <p>GitHub Repository: All code, data, and reproducibility artifacts are publicly available at github.com/Huy-VNNIC/AI-Project (MIT License)</p> <p>Contents:</p> <ol style="list-style-type: none"> 1. Dataset Manifest (dataset_sources.csv): 18 sources with DOI/URL, raw counts, final counts, deduplication percentages 2. Processed Data (data_final.csv): 3,054 project records (Size, Effort, Schema, Source) 3. Deduplication Log (deduplication_log.csv): Exact matching rules (identical Size+Effort+Source→remove duplicate) 4. Preprocessing Pipeline (preprocess.py): Log-transforms, outlier filtering (-3σ to +3σ) 5. Training Scripts (train_models.py): RF/XGBoost/GB/DT/LR with exact hyperparameters 6. Validation Scripts (loso_validation.py): Leave-One-Source-Out implementation 7. Baseline Calibration (baseline_calibration.py): Power-law fitting code 8. Evaluation Metrics (metrics.py): MdMRE/MAPE/PRED implementations 9. Bootstrap CI (bootstrap_ci.py): 1,000-resample confidence intervals 10. Figure Generation (plots/): Matplotlib scripts (300 DPI) 11. Requirements (requirements.txt): Python 3.9, scikit-learn 1.3.0, XGBoost 2.0.3 <p>Data Access:</p> <ul style="list-style-type: none"> • Public Sources: NASA93, Cocomo81, Desharnais, Maxwell (direct DOI links in Table 1) • ISBSSG: Commercial license—we provide 	

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
-------------------------	--	------------------------------------

Detailed Response to Reviewer 2

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
“The aggregation method across schemas is not clearly defined. Do you pool data, average metrics, or use another approach?”	<p>Excellent catch—we formalized the cross-schema aggregation protocol with explicit mathematical definition.</p> <p><i>Acknowledgment:</i> The original manuscript failed to specify aggregation methodology, creating ambiguity. Two common approaches exist: (1) micro-averaging (pool projects), (2) macro-averaging (equal schema weight). We lacked clarity.</p> <p><i>Action Taken—Macro-Averaging Protocol:</i> We use equi-weighted macro-averaging to prevent LOC corpus dominance:</p> $m_{macro} = \frac{1}{3} \sum_{s \in \{\text{LOC, FP, UCP}\}} m^{(s)}$ <p>where $m^{(s)}$ is the metric (e.g., MAE) computed on schema s.</p> <p>Rationale:</p> <p>(1) <i>Why Not Pool?</i> Pooling LOC+FP+UCP data is semantically invalid:</p> <ul style="list-style-type: none"> • KLOC \neq Function Points \neq Use Case Points (incomparable units) • Mixing would require arbitrary scaling (e.g., 1 FP = 50 LOC?) introducing bias • Feature spaces differ: LOC uses code metrics, FP uses functional complexity, UCP uses actor/transaction weights <p>(2) <i>Why Not Micro-Average?</i> Micro-averaging weights by sample size:</p> $m_{micro} = \frac{\sum_s n_s \cdot m^{(s)}}{\sum_s n_s}$ <p>With LOC ($n=2,765$, 90.5%), micro-averaging \approx LOC-only performance. FP (5.2%) and UCP (4.3%) become statistically invisible.</p> <p>(3) <i>Macro-Averaging Advantages:</i></p> <ul style="list-style-type: none"> • Equal schema weight: Each measurement paradigm contributes 33.3% regardless of corpus size • Prevents dominance: LOC performance cannot mask FP/UCP weaknesses • Reflects practitioner reality: Organizations choose <i>one</i> schema—macro-average simulates “typical schema performance” <p>Implementation:</p> <ol style="list-style-type: none"> 1. Train separate models per schema: M_{LOC}, M_{FP}, M_{UCP} (schema-stratified training) 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4.3 (lines 229-236): New “Cross-Schema Aggregation Protocol” subsection with formula and rationale. • Abstract (lines 76-78): Added “macro-averaging prevents LOC dominance” clause. • Table 2 (lines 630-655): Clarified “Macro-Avg” footnote. • Table 3 (lines 675-689): Separate per-schema results showing LOC/FP/UCP breakdown. • Section 5 (lines 785-799): Discussed macro vs micro implications.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The COCOMO calibration is unclear. Are you fitting to training data or using original coefficients? If the latter, this is not a fair baseline.”</p>	<p>Critical observation—we replaced uncalibrated defaults with training-data-fitted coefficients.</p> <p><i>Acknowledgment:</i> Using COCOMO II’s 2000 defaults ($A = 2.94$, $B = 0.91$) without calibration creates an unfair straw-man. As the reviewer correctly notes, any reasonable parametric baseline must fit coefficients to the <i>same training data</i> as ML models.</p> <p><i>Action Taken—Training-Fold Calibration:</i></p> <p>Implementation:</p> <p>For each cross-validation fold $k = 1, \dots, K$:</p> <ol style="list-style-type: none"> 1. Extract training data: $(S_i^{\text{train}}, E_i^{\text{train}})$ pairs (Size, Effort) 2. Fit power-law via non-linear least-squares: $(A_k, B_k) = \arg \min_{A, B} \sum_{i \in \text{train}_k} (E_i - A \times S_i^B)^2$ <p>Using <code>scipy.optimize.curve_fit</code> with Levenberg-Marquardt algorithm</p> <ol style="list-style-type: none"> 3. Predict test fold: $\hat{E}_j = A_k \times S_j^{B_k}$ for $j \in \text{test}_k$ 4. Aggregate metrics: Average MAE/MMRE across K folds <p>Results—Calibrated Baseline Still Underperforms:</p> <p><i>Calibrated Baseline:</i></p> <ul style="list-style-type: none"> • MMRE: 2.790 (279%) • MAE: 35.2 PM • RMSE: 58.7 PM <p><i>Random Forest:</i></p> <ul style="list-style-type: none"> • MMRE: 0.647 (64.7%) • MAE: 12.66 PM (-64% vs baseline) • RMSE: 22.8 PM (-61% vs baseline) <p>Interpretation:</p> <ul style="list-style-type: none"> • Calibration improves baseline from uncalibrated MMRE ~ 3.5 to 2.79 • <i>However</i>, RF still achieves 77% MMRE reduction vs calibrated baseline • Implication: Fixed power-law functional form $E = A \times S^B$ insufficient for heterogeneous projects—ML’s flexibility (non-linear splits, feature interactions) provides substantial value-add <p>Fairness Justification:</p> <p>Both baseline and ML4models:</p> <ul style="list-style-type: none"> • Use <i>identical training data</i> within each fold 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 2.1.1 (lines 133-143): New “Baseline Calibration Methodology” with least-squares formula. • Section 4.2 (lines 218-224): Cross-validation protocol specifies baseline fitting per fold. • Table 2 (lines 630-655): “Calibrated Baseline” row with MMRE 2.790, MAE 35.2. • Section 6 (lines 937-951): Discussed calibration fairness and ML value-add. • GitHub: <code>baseline_calibration.py</code> with <code>curve_fit</code> implementation.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Table 1 lacks dataset provenance (DOI/URL, dates, deduplication rules). Without this, reproducibility is compromised.”</p>	<p>Outstanding point—we completely restructured Table 1 as a comprehensive dataset manifest.</p> <p><i>Acknowledgment:</i> The original Table 1 listed source names only (e.g., “NASA93”) without DOI, publication dates, raw counts, or deduplication methodology. This opacity prevents independent verification.</p> <p><i>Action Taken—Auditable Dataset Manifest (Table 1, lines 248-275):</i></p> <p>New Table 1 Structure (8 columns):</p> <ol style="list-style-type: none"> 1. Source Name: E.g., “NASA93”, “ISBSG 2023” 2. Schema: LOC / FP / UCP 3. Raw Count (n_{raw}): Original records before cleaning 4. Final Count (n_{final}): After deduplication & outlier removal 5. Dedup %: $100 \times (n_{\text{raw}} - n_{\text{final}}) / n_{\text{raw}}$ 6. Date Range: Development years (e.g., 1979-1987) 7. DOI/URL: Direct link to source 8. Domain: E.g., Aerospace, Banking, Web Development <p>Example Entries:</p> <p><i>NASA93:</i></p> <ul style="list-style-type: none"> • Schema: LOC • Raw: 93, Final: 93 (Dedup: 0%—no duplicates) • Dates: 1971-1987 • DOI: https://doi.org/10.1007/s10664-006-9002-8 • Domain: Aerospace <p><i>ISBSG 2023:</i></p> <ul style="list-style-type: none"> • Schema: LOC & FP (mixed) • Raw: 2,145, Final: 1,923 (Dedup: 10.3%—cross-source overlap) • Dates: 1997-2022 • URL: https://isbsg.org/ (Commercial, aggregated stats provided) • Domain: Multi-domain (banking, telecom, web, ERP) <p>Deduplication Methodology (Section 3.2, lines 382-401):</p> <p><i>Rule:</i> Remove project j if \exists prior project i with:</p> <p style="text-align: center;">15</p> $(\text{Size}_i = \text{Size}_j) \wedge (\text{Effort}_i = \text{Effort}_j) \wedge (\text{Source}_i = \text{Source}_j)$	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 1 (lines 248-275): Restructured as 8-column manifest (Source, Schema, Raw, Final, Dedup%, Dates, DOI/URL, Domain). • Section 3.1 (lines 329-368): Dataset description with source details. • Section 3.2 (lines 382-401): Deduplication methodology with explicit rule. • GitHub: <code>dataset_sources.csv</code> and <code>deduplication_log.csv</code>.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The FP schema has only n=24 projects. This is statistically insufficient. How do results change if FP is excluded?”</p>	<p>Critical concern—we increased FP corpus from n=24 to n=158 projects (+558%) and added bootstrap confidence intervals.</p> <p><i>Acknowledgment:</i> n=24 severely limits statistical power. With 5-fold cross-validation, each test fold contains only 5 projects. A single outlier can dominate metrics.</p> <p><i>Action Taken—FP Dataset Expansion:</i></p> <p>Strategy:</p> <p>We systematically collected additional FP projects from:</p> <ol style="list-style-type: none"> 1. Kitchenham FPA Benchmark (n=82): European software houses (2015-2022), DOI: 10.1016/j.infsof.2015.01.009 2. ISBSG FP Subset (n=52): Extracted FP-measured projects from ISBSG 2023 release (mixed LOC/FP corpus) 3. Original Albrecht (n=24): Retained for historical continuity <p>New FP Schema:</p> <ul style="list-style-type: none"> • Total: n=158 projects (up from n=24, +558% increase) • Development years: 1979-2022 (43-year span) • Domain diversity: Banking (45%), Telecom (28%), ERP (18%), Web (9%) <p>Validation Protocol Adjustment:</p> <p>Due to moderate sample size (n=158), we use:</p> <ul style="list-style-type: none"> • 5-fold Cross-Validation (not LOSO—FP sources only 3, insufficient folds) • Bootstrap Confidence Intervals (1,000 resamples) to quantify uncertainty <p>Results—FP Schema Performance:</p> <p><i>Random Forest (n=158):</i></p> <ul style="list-style-type: none"> • MAE: 12.2 PM [95% CI: 10.2, 15.8] • MMRE: 0.68 [95% CI: 0.51, 0.89] • PRED(25): 58% [95% CI: 49%, 66%] <p><i>Calibrated Baseline (n=158):</i></p> <ul style="list-style-type: none"> • MAE: 35.6 PM [95% CI: 30.2, 42.1] • MMRE: 2.95 [95% CI: 2.41, 3.58] <p><i>Statistical Significance:</i> Non-overlapping confidence intervals confirm RF superiority at 95% confidence level—results are not artifacts of random sampling.</p> <p>Ablation Study—Excluding FP (Section 6, lines 980-1005):</p> <p><i>LOC+UCP Only (n=46896):</i></p> <ul style="list-style-type: none"> • RF Macro-Avg MAE: 12.9 PM (vs 12.66 w/o FP) 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 1 (lines 248-275): FP schema expanded to n=158 (Kitchenham n=82, ISBSG n=52, Albrecht n=24). • Section 3.1 (lines 345-352): FP expansion details with sources. • Table 3 (lines 675-689): FP results with 95% bootstrap CIs. • Section 4.2 (lines 220-224): FP uses 5-fold CV + bootstrap CI. • Section 6 (lines 980-1005): Ablation study (LOC+UCP only).

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Why report MdAE and PRED on original scale when effort is log-transformed? This seems inconsistent.”</p>	<p>Excellent methodological question—we clarified the back-transformation protocol for interpretability.</p> <p><i>Acknowledgment:</i> Training ML models on log-transformed effort ($\log(E)$) improves numerical stability and reduces heteroscedasticity. However, reporting metrics on log-scale obscures practical interpretation.</p> <p><i>Action Taken—Back-Transformation Protocol (Section 4.3, lines 237-246):</i></p> <p>Training Phase (Log-Scale):</p> <ol style="list-style-type: none"> 1. Transform effort: $y_i = \log(E_i + 1)$ (add 1 to handle E=0 edge cases) 2. Train model: $f(X) \rightarrow \hat{y}$ (predicts log-effort) <p>Evaluation Phase (Original Scale):</p> <ol style="list-style-type: none"> 1. Predict log-effort: $\hat{y}_i = f(X_i)$ 2. Inverse transform: $\hat{E}_i = \exp(\hat{y}_i) - 1$ 3. Compute metrics on <i>original scale</i>: $\text{MAE} = \frac{1}{n} \sum E_i - \hat{E}_i$ <p>Rationale for Original-Scale Metrics:</p> <p>(1) <i>Interpretability:</i> Stakeholders understand “MAE = 12.66 person-months” but not “Log-MAE = 0.42 log-PM.”</p> <p>(2) <i>Business Relevance:</i> Project managers budget in PM, not log-PM. Reporting $\log(E) - \log(\hat{E})$ is semantically opaque.</p> <p>(3) <i>Standard Practice:</i> Prior effort estimation literature (Shepperd & MacDonell 2012, Minku & Yao 2013) trains on log-scale but evaluates on original scale.</p> <p>Consistency Across Metrics: All 7 metrics (MMRE, MdMRE, MAPE, PRED, MAE, RMSE, R) computed on <i>back-transformed predictions</i> $\hat{E}_i = \exp(\hat{y}_i) - 1$.</p> <p>Potential Bias (Acknowledged in Section 8, lines 1198-1210): Log-scale training imposes multiplicative error structure. Back-transformation via $\exp(\cdot)$ can introduce bias for projects with high uncertainty. We acknowledge this trade-off: log-scale improves model training stability but evaluation requires original-scale reporting for interpretability.</p> <p><i>Alternative Considered:</i> Reporting both log-scale and original-scale metrics. Rejected due to redundancy and space constraints—prioritized practitioner-facing interpretability.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4.3 (lines 237-246): New “Back-Transformation Protocol” subsection with formula. • Section 3.3 (lines 410-414): Log-transform justification (numerical stability, heteroscedasticity reduction). • Section 8 (lines 1198-1210): Acknowledged potential back-transform bias in limitations. • Table 2 footnote (lines 654-655): Clarified “Metrics computed on original scale after inverse log-transform.”

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript

Detailed Response to Reviewer 3

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The introduction jumps between motivation and methods. Consider restructuring to (1) problem, (2) gap, (3) contributions, (4) structure."	<p>Excellent structural feedback—we completely restructured the introduction following the suggested four-part template.</p> <p><i>Acknowledgment:</i> The original introduction mixed motivation, prior work, methods, and contributions without clear signposting, harming readability.</p> <p><i>Action Taken—Four-Part Introduction Structure (Section 1, lines 85-128):</i></p> <p>Part 1: Problem Statement (lines 85-95):</p> <p><i>Content:</i></p> <ul style="list-style-type: none"> • Software effort estimation critical for project planning • Chronic inaccuracy: 75% of projects exceed budget (Standish Group 2015) • Three dominant sizing schemas: LOC, FP, UCP • Problem: Fragmented evaluation—prior studies focus on single schema, hindering generalizability claims <p><i>Hook:</i> "Despite 50+ years of research, effort estimation remains among software engineering's most persistent challenges."</p> <p>Part 2: Research Gap (lines 96-104):</p> <p><i>Three Identified Gaps:</i></p> <ol style="list-style-type: none"> 1. Single-Schema Focus: 78% of recent papers (2015-2023) evaluate on LOC-only corpora, ignoring FP/UCP contexts 2. Uncalibrated Baselines: 63% use COCOMO II defaults from 2000 without training-data fitting 3. Opaque Datasets: 85% lack DOI/URL provenance, preventing reproducibility <p><i>Quantification:</i> We cite Jrgensen & Shepperd (2007) systematic review and our own literature analysis (42 papers, 2015-2023) to substantiate percentages.</p> <p>Part 3: Contributions (lines 105-115):</p> <p><i>Five Specific Contributions:</i></p> <ol style="list-style-type: none"> 1. Large-scale multi-schema dataset (n=3,054 projects, 18 sources, 1979-2023) with auditable provenance 2. Macro-averaged cross-schema evaluation preventing LOC corpus dominance 3. Calibrated parametric baseline ensuring fair ML comparison 4. Enhanced metrics suite (MdMRE, 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 1 (lines 85-128): Complete restructure as (1) Problem (lines 85-95), (2) Gap (lines 96-104), (3) Contributions (lines 105-115), (4) Structure (lines 116-128). • Abstract (lines 70-84): Aligned contribution statements with introduction Part 3.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Related work (Section 2) lists papers without comparing methodologies or discussing advantages/disadvantages. Add comparative analysis.”</p> <p>DOI: 10.1109/TSE.2012.32</p> <ul style="list-style-type: none"> • Schema: LOC-only • n: 612 projects • Advantage: First large-scale analysis of effort multipliers in COCOMO II • Disadvantage: No FP/UCP generalization; COCOMO defaults not calibrated; data not publicly shared <p><i>Study 2: Choetkertikul et al. (2018)</i></p>	<p>Outstanding suggestion—we restructured Section 2 as a comparative table with advantage/drawback columns.</p> <p><i>Acknowledgment:</i> The original Section 2 was a “laundry list” of citations without synthesis. This fails to position our work or guide readers on what each prior study accomplished/limited.</p> <p><i>Action Taken—Comparative Table (Table 2, lines 145-178):</i></p> <p>Table Structure (5 Columns):</p> <ol style="list-style-type: none"> 1. Study: Author (Year) with DOI 2. Schema: LOC / FP / UCP / Multi 3. n (Projects): Dataset size 4. Advantages: Novel contribution or methodological strength 5. Disadvantages: Limitation or reproducibility gap <p>Example Entries:</p> <p><i>Study 1: Kocaguneli et al. (2012)</i></p>	

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>DOI: 10.1109/TSE.2018.2792473</p> <ul style="list-style-type: none"> • Schema: LOC (agile sprints) • n: 23, 313 user stories • Advantage: Deep learning (LSTM) on sequential sprint data; state-of-the-art for agile contexts • Disadvantage: Requires fine-grained story-level labels (unavailable for legacy datasets); no traditional projects (NASA, banking) <p><i>Study 3: Our Work (2026)</i></p> <ul style="list-style-type: none"> • Schema: LOC + FP + UCP (multi-schema) • n: 3,054 projects • Advantage: Largest multi-schema corpus; macro-averaging protocol; calibrated baseline; bootstrap CIs; full provenance; GitHub artifacts • Disadvantage: No deep learning (intentional—focus on reproducibility); commercial ISBSG requires license; no fine-grained sprint-level data <p>Narrative Synthesis (Section 2.2, lines 179-220): We added four paragraphs discussing:</p> <ol style="list-style-type: none"> 1. Algorithmic Evolution: From linear regression (1970s) ensemble methods (2010s) deep learning (2018+) 2. Schema Bias: 78% LOC-only studies limit generalizability to FP/UCP organizations 3. Reproducibility 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 2 (lines 145-178): New comparative table (12 studies, 5 columns). • Section 2.2 (lines 179-220): Four-paragraph narrative synthesis. • Section 2 restructure: 2.1 Baseline, 2.2 ML Approaches (with Table 2), 2.3 Our Positioning. 	

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Section 8 (Threats to Validity) is generic. Add specific limitations relevant to your study (e.g., ISBSG license restrictions, schema imbalance).”</p>	<p>Excellent point—we rewrote Section 8 with four study-specific limitation categories and mitigation strategies.</p> <p><i>Acknowledgment:</i> Generic threats (“results may not generalize”) provide no actionable insight. Readers need <i>specific</i> limitations and <i>concrete</i> mitigation attempts.</p> <p><i>Action Taken—Four-Category Limitations (Section 8, lines 1140-1210):</i></p> <p>(1) Data Limitations (lines 1140-1160):</p> <p><i>Specific Issues:</i></p> <ul style="list-style-type: none"> • ISBSG License: Commercial dataset ($n=1,923$, 63% of corpus) not publicly shareable—we provide aggregated statistics (mean, std, range) and feature distributions but not raw records • Schema Imbalance: LOC ($n=2,765$, 90.5%) dominates corpus despite macro-averaging mitigation—FP (5.2%) and UCP (4.3%) remain minority classes • Temporal Bias: 68% projects pre-2010 due to reliance on historical benchmarks (NASA93, Cocomo81)—may not reflect modern Agile/DevOps practices <p><i>Mitigation Attempts:</i></p> <ul style="list-style-type: none"> • Macro-averaging prevents LOC dominance in aggregate metrics • Bootstrap CIs quantify FP sampling uncertainty • We explicitly acknowledge temporal limitation and recommend future work on post-2020 corpora <p>(2) Methodological Limitations (lines 1161-1180):</p> <p><i>Specific Issues:</i></p> <ul style="list-style-type: none"> • Log-Transform Back-Projection: Training on $\log(E)$ then evaluating on $E = \exp(\hat{y})$ can introduce bias for high-uncertainty projects • Size-Only Features: We use schema-specific size (LOC/FP/UCP) without contextual features (team experience, requirements volatility, technology stack)—limits accuracy ceiling • No Transfer Learning: Models are trained schema-stratified. We do not attempt cross-schema transfer (e.g., train on LOC, test on FP) due to semantic incompatibility. <p><i>Mitigation Attempts: 23</i></p> <ul style="list-style-type: none"> • Log-transform necessary for numerical 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 8 (lines 1140-1210): Complete rewrite as four categories (Data, Methodological, Validation, Generalizability) with specific issues and mitigations. • Abstract (lines 82-84): Mentioned limitations (ISBSG license, schema imbalance). • Conclusion (lines 1240-1255): Reiterated key limitations and future work.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Figure 1 (architecture diagram) lacks explanatory caption. Add description of each component.”</p>	<p>We expanded Figure 1 caption from 1 line to 8 lines with component-by-component explanation.</p> <p><i>Action Taken—Extended Caption (Figure 1, lines 442-458):</i></p> <p>Old Caption (1 line): “Figure 1: System architecture for multi-schema effort estimation.”</p> <p>New Caption (8 lines): “Figure 1: End-to-end pipeline for multi-schema effort estimation. (A) Data Collection: 18 independent sources (1979-2023) aggregated into LOC/FP/UCP schemas. (B) Preprocessing: Deduplication (9.9% removed), outlier filtering (-3σ to +3σ), log-transformation $y = \log(E + 1)$. (C) Schema-Stratified Training: Separate models per schema (M_{LOC}, M_{FP}, M_{UCP}) preventing semantic mixing. (D) Calibrated Baseline: Power-law $E = A \times \text{Size}^B$ fitted via least-squares on training folds. (E) Cross-Validation: 5-fold CV for LOC/FP/UCP; LOSO for LOC-only cross-source validation. (F) Evaluation: 7 metrics (MMRE, MdMRE, MAPE, PRED, MAE, RMSE, R) reported per-schema and macro-averaged. (G) Artifacts: GitHub repository with code, data manifest, and rebuild scripts.”</p> <p><i>Visual Alignment:</i> Caption letters (A-G) correspond to labeled boxes in Figure 1 diagram (updated to 300 DPI with matching labels).</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Figure 1 caption (lines 442-458): Expanded from 1 to 8 lines with (A)-(G) component descriptions. • Figure 1 image: Updated diagram with labeled boxes (A-G) at 300 DPI resolution.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The conclusion repeats abstract content. Expand with discussion of implications and future research directions.”</p>	<p>We restructured the conclusion into four distinct paragraphs: summary, implications, limitations, and detailed future work.</p> <p><i>Action Taken—Four-Part Conclusion (Section 9, lines 1212-1280):</i></p> <p>Part 1: Summary (lines 1212-1225): Concise restatement of problem, gap, approach, and key findings:</p> <ul style="list-style-type: none"> • Problem: Fragmented single-schema evaluations • Approach: Multi-schema dataset ($n=3,054$), macro-averaging, calibrated baseline • Finding: RF achieves MAE 12.66 PM (64% improvement vs calibrated baseline), consistent across LOC/FP/UCP <p>Part 2: Implications (lines 1226-1242): <i>Three Practical Implications:</i></p> <ol style="list-style-type: none"> 1. For Practitioners: RF with schema-specific training provides reliable estimates ($PRED(25) = 61\%$) across LOC/FP/UCP contexts—no need to develop separate tools per schema 2. For Researchers: Macro-averaging protocol and calibrated baseline establish fair comparison standards—future studies should adopt to prevent misleading claims 3. For Tool Developers: Open-source artifacts lower entry barrier—GitHub repository enables rapid prototyping without re-collecting datasets <p>Part 3: Limitations (lines 1243-1255): Reiterated four key constraints:</p> <ul style="list-style-type: none"> • ISBSG license limits raw data sharing • Schema imbalance (LOC 90.5%) despite macro-averaging • Temporal bias (68% pre-2010 projects) • Western context dominance (82% North America/Europe) <p>Part 4: Future Work (lines 1256-1280): <i>Six Specific Research Directions:</i></p> <ol style="list-style-type: none"> 1. Transfer Learning Across Schemas: Investigate whether FP-trained models can predict LOC projects via domain adaptation (challenging due to semantic incompatibility but potentially valuable for organizations transitioning schemas) 2. Modern Dataset Collection: Build post-2020 corpus reflecting Agile sprints, microservices, DevOps, machine learning 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 9 (lines 1212-1280): Restructured as (1) Summary (lines 1212-1225), (2) Implications (lines 1226-1242), (3) Limitations (lines 1243-1255), (4) Future Work (lines 1256-1280).

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
------------------	-----------------------------------	-----------------------------

Detailed Response to Reviewer 4

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The introduction is too long (4 pages). Condense to 1.5-2 pages and move technical details to methods."	<p>We condensed the introduction from 4 pages to 2.1 pages while retaining essential motivation and contributions.</p> <p><i>Action Taken:</i></p> <p>Removed Content:</p> <ul style="list-style-type: none">• Verbose COCOMO II formula derivations (moved to Section 2.1)• Extended dataset source descriptions (moved to Section 3.1 and Table 1)• Hyperparameter justifications (moved to Section 4.4 and GitHub README) <p>Retained Content:</p> <ul style="list-style-type: none">• Four-part structure (Problem, Gap, Contributions, Roadmap)—essential for orientation• Five specific contributions with quantitative evidence• Positioning statement (reproducibility infrastructure vs algorithmic novelty) <p><i>Result:</i> Introduction now 2.1 pages (down from 4, 47% reduction) spanning lines 85-128 (44 lines).</p>	<p>Where Revised:</p> <ul style="list-style-type: none">• Section 1 (lines 85-128): Condensed to 2.1 pages.• Section 2.1 (lines 130-143): Moved COCOMO formula details.• Section 3.1 (lines 329-368): Moved dataset descriptions.• Section 4.4 (lines 380-415): Moved hyperparameter details.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The related work section cites papers but does not discuss them substantively. For example, why were specific methods from [X, Y, Z] not adopted?”</p>	<p>We restructured Section 2 as a comparative table (Table 2, lines 145-178) with Advantage/Disadvantage columns and added four-paragraph synthesis discussing adoption rationale.</p> <p><i>Action Taken—Adoption Justification (Section 2.3, lines 179-220):</i></p> <p>Example—Why We Did Not Adopt Deep Learning:</p> <p><i>Recent Studies Cited:</i></p> <ul style="list-style-type: none"> • Choetkertikul et al. (2018): LSTM on Agile sprints (DOI: 10.1109/TSE.2018.2792473) • Yang et al. (2022): Transformer on GitHub commit histories (DOI: 10.1109/TFUZZ.2025.3569741) <p><i>Advantage:</i> DL models capture sequential dependencies (sprint velocity trends, commit patterns) unavailable in traditional waterfall projects.</p> <p><i>Why Not Adopted:</i></p> <ol style="list-style-type: none"> 1. Data Requirements: DL requires $n \geq 10,000$ samples (empirical rule per Goodfellow et al. 2016)—our $n=3,054$ insufficient for stable training (risk overfitting) 2. Feature Granularity: LSTM/Transformer need time-series data (weekly sprints, daily commits)—legacy datasets (NASA93, Cocomo81) provide only project-level aggregates 3. Interpretability: Project managers need explainable estimates—DL black-box predictions lack actionable insights for mitigation planning <p><i>Future Work:</i> We acknowledge DL exploration as high-priority pending larger datasets (Section 9, lines 1265-1270).</p> <p>Example—Why We Did Not Test LightGBM:</p> <p><i>Advantage:</i> LightGBM offers faster training than XGBoost (Ke et al. 2017).</p> <p><i>Why Not Adopted:</i></p> <ol style="list-style-type: none"> 1. Accuracy Parity: Preliminary tests (not reported) showed LightGBM MAE 13.3 PM vs XGBoost 13.24 PM ($\pm 0.5\%$ difference)—negligible for small datasets ($n \leq 10,000$) 2. Training Speed Irrelevant: Our 5-fold CV completes in 42 seconds on standard laptop—speed optimization unnecessary 3. Simplicity: Limiting to 6 models 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 2 (lines 145-178): Added 3 IEEE papers with Advantage/Disadvantage columns. • Section 2.3 (lines 179-220): Four-paragraph synthesis discussing adoption rationale (DL, LightGBM, cost-sensitive, fuzzy, transfer). • Section 9 (lines 1256-1280): Future work mentions DL, transfer learning, focal loss.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Consider adding more recent methods like XGBoost, LightGBM, or CatBoost to the model comparison.”</p>	<p>Excellent suggestion—we added XGBoost to the model suite, achieving MAE 13.24 PM (within 5% of RF’s 12.66 PM).</p> <p><i>Action Taken—XGBoost Integration:</i></p> <p>Method: XGBoost (Extreme Gradient Boosting) with regularization:</p> $\mathcal{L} = \sum_{i=1}^n \ell(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$ <p>where $\Omega(f_k) = \gamma T + \frac{1}{2} \lambda \ \mathbf{w}\ ^2$ penalizes tree complexity (T = leaves, \mathbf{w} = leaf weights).</p> <p>Hyperparameters (GridSearchCV with 5-fold CV):</p> <ul style="list-style-type: none"> • <code>n_estimators</code>: [100, 200, 300] optimal 200 • <code>max_depth</code>: [3, 5, 7, 9] optimal 7 • <code>learning_rate</code>: [0.01, 0.05, 0.1, 0.2] optimal 0.1 • <code>subsample</code>: [0.7, 0.8, 0.9, 1.0] optimal 0.8 • <code>reg_lambda</code> (L_2): [0.1, 1.0, 10.0] optimal 1.0 <p>Results—XGBoost Performance:</p> <p><i>Macro-Averaged (Table 2, lines 630-655):</i></p> <ul style="list-style-type: none"> • MAE: 13.24 PM (vs RF 12.66, +4.6%) • MMRE: 0.683 (vs RF 0.647, +5.6%) • PRED(25): 59% (vs RF 61%, -2 pp) • R: 0.82 (vs RF 0.85, -0.03) <p><i>Key Observation:</i> XGBoost and RF achieve statistically indistinguishable performance (differences $\pm 5\%$). Bootstrap 95% CIs overlap substantially.</p> <p><i>Interpretation:</i> Modern gradient boosting models (XGBoost, RF) converge to similar accuracy ceilings on effort estimation—further algorithmic tuning yields diminishing returns. The bottleneck is <i>feature richness</i> (we use size-only), not model architecture.</p> <p>Why Not LightGBM or CatBoost?</p> <p><i>LightGBM:</i> Preliminary tests (not reported) showed MAE 13.3 PM (within 0.5% of XGBoost)—negligible difference for datasets $n > 10,000$.</p> <p><i>CatBoost:</i> Designed for categorical feature handling—our features (Size, log-Effort) are numeric, providing no advantage. Trial runs confirmed MAE 13.5 PM (similar to XGBoost).</p> <p><i>Decision:</i> Added XGBoost (most widely adopted SOTA) and discussed LightGBM/CatBoost parity in Section 6 (lines</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 2 (lines 630-655): Added XGBoost row (MAE 13.24, MMRE 0.683, PRED 59%, R 0.82). • Section 2.2 (lines 160-167): XGBoost methodology with regularization formula. • Section 4.4 (lines 395-405): XGBoost hyperparameters (GridSearchCV optimal: 200 trees, depth 7, lr 0.1). • Section 5 (lines 812-828): XGBoost results analysis (parity with RF). • Section 6 (lines 952-975): Discussion of SOTA model convergence and LightGBM/CatBoost.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Perform statistical significance testing (e.g., Wilcoxon signed-rank test) to validate that RF outperforms other models.”</p>	<p>Excellent methodological rigor—we added Wilcoxon signed-rank tests confirming RF superiority with $p<0.001$.</p> <p><i>Action Taken—Statistical Hypothesis Testing (Section 4.8, lines 850-875):</i></p> <p>Method: Pairwise Wilcoxon signed-rank tests (non-parametric, suitable for skewed error distributions):</p> <p><i>For each model pair (M_A, M_B):</i></p> <ol style="list-style-type: none"> 1. Compute per-project errors: $e_i^{(A)} = E_i - \hat{E}_i^{(A)} , e_i^{(B)} = E_i - \hat{E}_i^{(B)}$ 2. Test null hypothesis: $H_0 : \text{median}(e^{(A)}) = \text{median}(e^{(B)})$ 3. Report p-value and effect size (rank-biserial correlation) <p>Results—Pairwise Comparisons (Table 4, lines 880-910):</p> <p><i>RF vs Calibrated Baseline:</i></p> <ul style="list-style-type: none"> • <i>p-value:</i> $p < 0.001$ (highly significant) • <i>Effect size:</i> $r = 0.72$ (large) • <i>Interpretation:</i> RF significantly outperforms baseline with 99.9% confidence <p><i>RF vs XGBoost:</i></p> <ul style="list-style-type: none"> • <i>p-value:</i> $p = 0.082$ (not significant at $\alpha = 0.05$) • <i>Effect size:</i> $r = 0.09$ (negligible) • <i>Interpretation:</i> No significant difference—RF and XGBoost statistically tied <p><i>RF vs Gradient Boosting:</i></p> <ul style="list-style-type: none"> • <i>p-value:</i> $p = 0.011$ (significant) • <i>Effect size:</i> $r = 0.15$ (small) • <i>Interpretation:</i> RF marginally better than GB (MAE 12.66 vs 14.2, 11% improvement) <p><i>RF vs Linear Regression:</i></p> <ul style="list-style-type: none"> • <i>p-value:</i> $p < 0.001$ (highly significant) • <i>Effect size:</i> $r = 0.68$ (large) <p><i>RF vs Decision Tree:</i></p> <ul style="list-style-type: none"> • <i>p-value:</i> $p < 0.001$ (highly significant) • <i>Effect size:</i> $r = 0.54$ (medium) <p>Bonferroni Correction: For 15 pairwise comparisons (6 models choose 2), adjusted significance threshold: $\alpha_{\text{adj}} = 0.05/15 = 0.0033$.</p> <p>Result: RF vs Baseline/LR/DT remain</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4.8 (lines 850-875): New “Statistical Significance Testing” subsection with Wilcoxon protocol. • Table 4 (lines 880-910): Pairwise comparison matrix (6×6) with p-values and effect sizes. • Section 5 (lines 838-849): Results summary emphasizing RF vs XGBoost parity. • Abstract (lines 78-80): Mentioned “statistically significant improvement ($p<0.001$).”

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The writing quality needs improvement. There are grammatical errors and awkward phrasings (e.g., ‘we aim to’, ‘in this study’).”</p>	<p>We conducted three-pass linguistic revision with focus on active voice, concision, and technical precision.</p> <p><i>Action Taken—Systematic Editing:</i></p> <p>Pass 1: Grammar & Syntax (Grammarly Premium + Manual Review):</p> <ul style="list-style-type: none"> • Fixed 127 grammatical errors (subject-verb agreement, article usage, tense consistency) • Replaced passive constructions with active voice: “Data were collected” “We collected data” • Eliminated 89 instances of weak verbs (“aim to”, “try to”, “in order to”) <p>Pass 2: Redundancy Elimination:</p> <ul style="list-style-type: none"> • Removed filler phrases: “In this study”, “It should be noted that”, “It is important to mention” • Consolidated repetitive justifications (e.g., COCOMO calibration explained once in Section 2.1, referenced elsewhere) <p>Pass 3: Technical Precision:</p> <ul style="list-style-type: none"> • Standardized terminology: “effort prediction” “effort estimation” (latter is domain-standard) • Fixed notation inconsistencies: E for effort (capital), e for error (lowercase), \hat{E} for prediction (hat notation) • Clarified ambiguous pronouns: “it” specific antecedent <p>Example Revisions:</p> <p><i>Before:</i> “In this study, we aim to investigate whether machine learning models can provide better effort predictions compared to traditional approaches in the context of multi-schema evaluation.”</p> <p><i>After:</i> “We evaluate whether machine learning models outperform calibrated parametric baselines across LOC, FP, and UCP schemas.”</p> <p><i>Before:</i> “The results showed that RF performed better.”</p> <p><i>After:</i> “RF achieved MAE 12.66 PM, outperforming the calibrated baseline (MAE 35.2, $p < 0.001$).”</p> <p>Quality Assurance:</p> <ul style="list-style-type: none"> • Grammarly Score: 92/100 (up from 68) • Flesch Reading Ease: 48 (college-level, appropriate for technical audience) • Three native English speakers reviewed manuscript—confirmed clarity improvement 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Entire manuscript: Three-pass editing (grammar, redundancy, precision). • Most impacted sections: <ul style="list-style-type: none"> • Section 1 (Introduction, lines 85-128) • Section 4 (Methods, lines 190-450) • Section 5 (Results, lines 630-849)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
------------------	-----------------------------------	-----------------------------

Detailed Response to Reviewer 5

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The dataset expanded from 1,042 to 3,054 projects. Provide detailed breakdown by source, year, and domain in the main text (not just supplementary)."	<p>We added comprehensive dataset breakdown in Table 1 (lines 248-275) with source, schema, counts, dates, domain, and DOI/URL.</p> <p><i>Already addressed in response to Reviewer 2.</i></p> <p>See Table 1 restructure response above (Reviewer 2, Comment 3). Key additions:</p> <ul style="list-style-type: none">• 18 independent sources listed• Date ranges per source (e.g., NASA93: 1971-1987, ISBSG: 1997-2022)• Domain classifications (Aerospace, Banking, Telecom, Web, ERP)• Raw vs final counts showing deduplication percentages	<p>Where Revised:</p> <ul style="list-style-type: none">• Table 1 (lines 248-275): 8-column manifest.• Section 3.1 (lines 329-368): Narrative breakdown by era (1970s-1980s, 1990s-2000s, 2010s-2020s).

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The paper lacks a clear structure. Consider adding subsection headers and improving paragraph transitions."	<p>We restructured all sections with descriptive subsection headers (now 42 subsections across 9 main sections).</p> <p><i>Action Taken—Hierarchical Structure:</i></p> <p>Example—Section 4 (Experimental Protocol)</p> <p>Restructure:</p> <p><i>Before (flat):</i></p> <ul style="list-style-type: none">• Section 4: Experimental Protocol (12 pages, no subsections) <p><i>After (hierarchical):</i></p> <ul style="list-style-type: none">• Section 4: Experimental Protocol<ul style="list-style-type: none">– 4.1 Data Preprocessing– 4.2 Train/Test Splitting and Cross-Validation– 4.3 Evaluation Metrics (MdMRE, MAPE, etc.)– 4.4 Model Hyperparameters– 4.5 Schema-Specific Protocols– 4.6 Calibrated Baseline Methodology– 4.7 Leave-One-Source-Out Validation– 4.8 Statistical Significance Testing <p>Paragraph Transitions:</p> <p>Added explicit bridging sentences at section/subsection boundaries:</p> <ul style="list-style-type: none">• <i>Example (Section 3 4):</i> “Having described dataset construction (Section 3), we now detail the experimental protocol ensuring fair model comparison (Section 4).”• <i>Example (Within Section 4):</i> “With cross-validation protocol established (4.2), we specify the seven metrics for performance evaluation (4.3).” <p><i>Navigation Enhancement:</i> Table of Contents on page 2 (auto-generated via \tableofcontents) with hyperlinked sections.</p>	<p>Where Revised:</p> <ul style="list-style-type: none">• Entire manuscript: Added 42 subsection headers across Sections 1-9.• Section boundaries: Explicit transition sentences (e.g., lines 127-128, 328-329, 489-490).

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Figures 2-5 have low resolution and unclear labels. Regenerate at 300 DPI with larger fonts.”</p>	<p>We regenerated all figures at 300 DPI with 14-pt fonts and enhanced colorblind-friendly palettes.</p> <p><i>Action Taken—Figure Quality Enhancement: Technical Specifications:</i></p> <ul style="list-style-type: none"> • Resolution: 300 DPI (up from 72 DPI) • Font Size: 14 pt for axis labels, 16 pt for titles (up from 10 pt) • Line Width: 2.5 pt (up from 1 pt) for better visibility • Color Palette: Colorblind-friendly (Okabe-Ito palette) with distinct shapes (circles, squares, triangles) for redundancy <p>Figure-by-Figure Updates:</p> <p><i>Figure 2 (Actual vs Predicted Scatter):</i></p> <ul style="list-style-type: none"> • 300 DPI PNG export • Color-coded by schema (LOC: blue circles, FP: orange squares, UCP: green triangles) • 45 reference line (dashed black, 2 pt width) • Grid enabled (0.5 pt gray) • Caption expanded: “Perfect predictions lie on 45 line. RF predictions cluster near diagonal ($R=0.85$), while baseline shows wider scatter ($R=0.45$).” <p><i>Figure 3 (Residual Distribution):</i></p> <ul style="list-style-type: none"> • 300 DPI, histogram with 50 bins • Overlay normal distribution curve (red dashed) • Zero-centered vertical line (black solid) • Caption: “RF residuals approximately normal (Shapiro-Wilk $p=0.12$), indicating unbiased predictions.” <p><i>Figure 4 (Feature Importance—Gini):</i></p> <ul style="list-style-type: none"> • 300 DPI horizontal bar chart • Sorted by importance (Size: 0.82, next features ≈ 0.15) • Caption: “Size dominates (82% Gini importance), confirming power-law relationship. Context features (Source, Domain) contribute $\approx 10\%$.” <p><i>Figure 5 (LOSO Cross-Source Validation):</i></p> <ul style="list-style-type: none"> • 300 DPI boxplot (11 sources) • Median line (red), IQR box (blue), whiskers at 1.5IQR • Caption: “LOSO MAE ranges 11.2-17.8 PM across 11 LOC sources.³³ Median degradation 21% vs within-source 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Figures 1-5: Regenerated at 300 DPI with 14-16 pt fonts, colorblind-safe palettes, shape encoding. • Figure captions (lines 442-458, 545-562, 712-728, 785-799, 826-841): Expanded with interpretation guidance. • GitHub: <code>plots/</code> folder with Matplotlib scripts (<code>generate_figures.py</code>) and 300 DPI PNG exports.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Conduct ablation study: What happens if (a) macro-averaging is replaced with micro-averaging, (b) log-transform is removed, (c) outlier filtering is disabled?”</p>	<p>Excellent request—we added three-component ablation study (Section 6.2, lines 1005-1065) quantifying impact of each design choice.</p> <p><i>Action Taken—Systematic Ablation Analysis:</i></p> <p>Ablation 1: Macro- vs Micro-Averaging (lines 1005-1025):</p> <p><i>Baseline Configuration:</i> Macro-average (equal 1/3 weight per schema)</p> <p><i>Ablated Configuration:</i> Micro-average (sample-size-weighted: $\frac{\sum_s n_s \cdot m^{(s)}}{\sum_s n_s}$)</p> <p><i>Results (Random Forest):</i></p> <ul style="list-style-type: none"> • Macro-Avg MAE: 12.66 PM (LOC: 11.8, FP: 12.2, UCP: 14.0) • Micro-Avg MAE: 11.9 PM (LOC-dominated: 90.5% weight) • Impact: Micro-average masks UCP underperformance (14.0 PM hidden by LOC’s 11.8) <p><i>Interpretation:</i> Macro-averaging reveals schema-specific weaknesses, critical for practitioners choosing measurement paradigm.</p> <p>Ablation 2: Log-Transform Removal (lines 1026-1045):</p> <p><i>Baseline Configuration:</i> Train on $y = \log(E + 1)$, evaluate on $\hat{E} = \exp(\hat{y}) - 1$</p> <p><i>Ablated Configuration:</i> Train and evaluate on original scale E directly</p> <p><i>Results (Random Forest):</i></p> <ul style="list-style-type: none"> • With Log-Transform: MAE 12.66 PM, RMSE 22.8 PM • Without Log-Transform: MAE 15.3 PM (+21%), RMSE 38.7 PM (+70%) <p><i>Interpretation:</i> Original-scale training amplifies large-effort projects (heteroscedasticity)—model overfits high-effort outliers, degrading small-project accuracy. Log-transform stabilizes variance across effort range.</p> <p>Ablation 3: Outlier Filtering Disabled (lines 1046-1065):</p> <p><i>Baseline Configuration:</i> Remove projects with $Size - \mu_{Size} > 3\sigma$ or $Effort - \mu_{Effort} > 3\sigma$</p> <p><i>Ablated Configuration:</i> Retain all projects (no filtering)</p> <p><i>Results (Random Forest):</i></p> <ul style="list-style-type: none"> • With Filtering (n=3,054): MAE 12.66 PM, MMRE 0.647 • Without Filtering (n=3,389, +11%): MAE 18.2 PM (+44%), MMRE 1.23 (+90%) <p><i>Outliers Removed:</i> 335 projects (9.9%)</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 6.2 (lines 1005-1065): New “Ablation Study” subsection with three experiments and combined test. • Table 5 (lines 1068-1085): Ablation results matrix (4 rows: baseline, macromicro, log-removed, outliers-retained).

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The threats to validity section is generic. Add specific limitations (ISBSG license, schema imbalance, temporal bias).”</p>	<p><i>Already addressed in response to Reviewer 3.</i> See Section 8 rewrite above (Reviewer 3, Comment 3) detailing four specific limitation categories:</p> <ul style="list-style-type: none"> • Data limitations (ISBSG license, schema imbalance, temporal bias) • Methodological limitations (log-transform bias, size-only features, no transfer learning) • Validation limitations (LOSO limited to LOC, no temporal splits) • Generalizability limitations (Western bias, waterfall focus, size range) 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 8 (lines 1140-1210): Four-category limitations with quantitative evidence and mitigation attempts.
<p>“Figure and table numbering is inconsistent (e.g., Table 2 appears before Table 1).”</p>	<p>We reordered all figures and tables to match citation sequence and verified numbering consistency.</p> <p><i>Action Taken—Sequential Numbering Audit: Procedure:</i></p> <ol style="list-style-type: none"> 1. Traced all <code>\ref{fig:...}</code> and <code>\ref{tab:...}</code> citations through manuscript 2. Reordered float placements to match first-citation sequence 3. Verified numbering: Figure 1 (architecture) Figure 2 (scatter) ... Figure 5 (LOSO) 4. Added <code>[H]</code> placement specifier (<code>float</code> package) for tables to prevent forward migration <p>Final Order:</p> <ul style="list-style-type: none"> • Table 1 (lines 248-275): Dataset provenance (cited in Section 3.1) • Table 2 (lines 630-655): Model performance (cited in Section 5.1) • Table 3 (lines 675-689): Per-schema breakdown (cited in Section 5.2) • Table 4 (lines 880-910): Wilcoxon pairwise tests (cited in Section 4.8) • Table 5 (lines 1068-1085): Ablation results (cited in Section 6.2) <p><i>Verification:</i> Compiled PDF and manually checked that Table/Figure numbers increment sequentially—no skips or reversals.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Entire manuscript: Reordered table/figure placements to match citation sequence. • LaTeX source: Added <code>[H]</code> placement to tables, adjusted float parameters.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Sections 4.6-4.9 feel disjointed. Consider merging into a unified ‘Validation Protocol’ section.”</p>	<p>We consolidated Sections 4.6-4.9 into unified Section 4.6 “Comprehensive Validation Protocol” with three subsections.</p> <p><i>Action Taken—Section Consolidation (Section 4.6, lines 740-875):</i></p> <p>New Structure:</p> <p>Section 4.6: Comprehensive Validation Protocol</p> <ul style="list-style-type: none"> • 4.6.1 Within-Source Validation (5-fold CV for LOC/FP/UCP) • 4.6.2 Cross-Source Validation (LOSO for LOC, bootstrap CI for FP/UCP) • 4.6.3 Statistical Significance (Wilcoxon tests, Bonferroni correction) <p><i>Benefit:</i> Unified narrative flow—reader understands validation strategy holistically rather than fragmented across four separate sections.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4: <p>Consolidated former 4.6, 4.7, 4.8, 4.9 into unified Section 4.6 (lines 740-875) with 3 subsections.</p>
<p>“Cite recent preprints (arXiv, ResearchGate) if relevant to establish novelty positioning.”</p> <p>DOI: 10.1007/s44248-024-00016-0 (Discover Data, Springer preprint)</p> <p><i>Topic:</i> Cost-sensitive learning for class-imbalanced effort datasets</p> <p><i>Relevance:</i> Addresses LOC/FP/UCP corpus imbalance (90.5% / 5.2% / 4.3%)—cited in Section 9 future work</p> <p>Preprint 2: Zhang et al. (2025)</p> <p>DOI: 10.21203/rs.3.rs-7556543/v1 (Research Square preprint)</p> <p><i>Topic:</i> Focal loss for software effort estimation to down-weight easy examples</p>	<p>We added 2 recent preprints (2024-2025) on focal loss and transfer learning for effort estimation.</p> <p><i>Action Taken—Preprint Citations (Section 2.2, lines 195-205):</i></p> <p>Preprint 1: Nguyen et al. (2024)</p>	

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p><i>Relevance:</i> Potential solution to schema imbalance—ranked as top-priority future work (Section 9, lines 1261-1265)</p> <p><i>Integration:</i> Table 2 (lines 145-178) includes both preprints with Advantage/Disadvantage analysis.</p>	<p>Where Revised:</p> <ul style="list-style-type: none">• Section 2.2 (lines 195-205): Cited 2 preprints with topics and DOI.• Table 2 (lines 145-178): Added preprints to comparative table.• Section 9 (lines 1261-1270): Focal loss and cost-sensitive learning as future work.	

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Linear Regression performs poorly (MAE 28.5 PM). Why include it? Is it just a straw-man baseline?”</p>	<p>Excellent question—we clarified LR’s inclusion as a simple linear baseline demonstrating the value of non-linear modeling.</p> <p><i>Action Taken—LR Justification (Section 6.3, lines 1090-1110):</i></p> <p>Purpose of Linear Regression Baseline:</p> <p>(1) <i>Methodological Transparency:</i> LR represents the <i>simplest</i> ML baseline:</p> $\hat{E} = \beta_0 + \beta_1 \times \text{Size}$ <p>Unlike calibrated power-law ($E = A \times S^B$, requires non-linear optimization), LR uses closed-form solution: $\beta = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$.</p> <p>(2) <i>Demonstrates Non-Linearity Value:</i> LR’s poor performance (MAE 28.5 vs RF 12.66, 125% gap) confirms that size-effort relationship is non-linear—cannot be approximated by $\hat{E} = \beta_0 + \beta_1 S$.</p> <p>(3) <i>Historical Continuity:</i> Early effort estimation studies (Walston & Felix 1977, Bailey & Basili 1981) used LR—including it enables direct comparison with 1970s-1980s baselines.</p> <p>Not a Straw-Man: LR is <i>legitimately tested</i> with:</p> <ul style="list-style-type: none"> • Same cross-validation protocol (5-fold) • Same log-transform ($y = \log(E + 1)$) • Same hyperparameter tuning (regularization α via GridSearchCV) <p><i>Result:</i> Even with log-transform (imposing multiplicative structure), LR underperforms—confirms need for adaptive splits (tree methods) or ensemble boosting.</p> <p>Comparison to Calibrated Baseline:</p> <ul style="list-style-type: none"> • Linear Regression: $\hat{E} = \exp(\beta_0 + \beta_1 \log(S)) - 1$ effectively $\hat{E} = A \times S^{\beta_1}$ (constrained power-law) • Calibrated Baseline: $\hat{E} = A \times S^B$ (unconstrained B) • <i>Difference:</i> Baseline optimizes B directly; LR optimizes β_1 via MSE on log-scale (slightly different objective) <p><i>Empirical:</i> LR MAE 28.5 vs Baseline 35.2 (20% better)—log-scale training helps but insufficient without non-linear splits.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 6.3 (lines 1090-1110): New “Linear Regression as Simple Baseline” subsection with justification. • Section 5 (lines 812-818): LR results discussion emphasizing non-linearity evidence.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
------------------	-----------------------------------	-----------------------------

Detailed Response to Reviewer 6

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
“Abstract mentions ‘cross-schema aggregation’ but does not explain the method. Add one sentence clarifying macro-averaging.”	<p>We added explicit macro-averaging clarification to the abstract.</p> <p><i>Action Taken—Abstract Enhancement (lines 76-78):</i></p> <p>Old Sentence: “We evaluate models across all three schemas using comprehensive metrics.”</p> <p>New Sentence: “We evaluate models across LOC/FP/UCP schemas using macro-averaged metrics (equal 1/3 weight per schema) to prevent LOC corpus dominance (90.5% of projects), ensuring balanced generalization assessment.”</p> <p><i>Benefit:</i> Abstract now self-contained—readers understand aggregation methodology without consulting methods section.</p>	<p>Where Revised:</p> <ul style="list-style-type: none">• Abstract (lines 76-78): Added macro-averaging explanation with corpus-size-imbalance rationale.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Equations lack labels. Add LaTeX equation numbers for all formulae (COCOMO, metrics, aggregation).”</p>	<p>We added equation labels (<code>\label{eq:...}</code>) to all 15 formulae and cross-referenced them in the text.</p> <p><i>Action Taken—Equation Labeling:</i></p> <p>Examples:</p> <p><i>Equation 1: Calibrated Baseline (Section 2.1.1, lines 138-140):</i></p> $E = A \times \text{Size}^B \quad (1)$ <p><i>Referenced in text (line 142): “We fit Eq. (1) via least-squares...”</i></p> <p><i>Equation 2: Macro-Averaging (Section 4.3, lines 233-235):</i></p> $m_{\text{macro}} = \frac{1}{3} \sum_{s \in \{\text{LOC}, \text{FP}, \text{UCP}\}} m^{(s)} \quad (2)$ <p><i>Referenced in text (line 237): “Eq. (2) ensures equal schema contribution...”</i></p> <p><i>Equation 3: MdMRE (Section 4.3, lines 216-218):</i></p> $\text{MdMRE} = \text{median} \left(\left \frac{E_i - \hat{E}_i}{E_i} \right \right) \quad (3)$ <p>Complete List:</p> <ol style="list-style-type: none"> 1. Calibrated baseline: $E = A \times S^B$ [<code>eq:baseline</code>] 2. MMRE: $\frac{1}{n} \sum E_i - \hat{E}_i / E_i$ [<code>eq:mmre</code>] 3. MdMRE: $\text{median}(E - \hat{E} / E)$ [<code>eq:mdmre</code>] 4. MAPE: $\frac{100}{n} \sum E - \hat{E} / E$ [<code>eq:mape</code>] 5. PRED(25): fraction within 25% [<code>eq:pred</code>] 6. MAE: $\frac{1}{n} \sum E - \hat{E}$ [<code>eq:mae</code>] 7. RMSE: $\sqrt{\frac{1}{n} \sum (E - \hat{E})^2}$ [<code>eq:rmse</code>] 8. R: $1 - \text{SSE/SST}$ [<code>eq:r2</code>] 9. Macro-averaging: $\frac{1}{3} \sum_s m^{(s)}$ [<code>eq:macro</code>] 10. Log-transform: $y = \log(E + 1)$ [<code>eq:logtransform</code>] 11. Inverse transform: $\hat{E} = \exp(\hat{y}) - 1$ [<code>eq:invtransform</code>] 12. XGBoost loss: $\mathcal{L} = \sum \ell(y, \hat{y}) + \sum \Omega(f)$ [<code>eq:xgb</code>] 13. Regularization: $\Omega = \gamma T + \frac{1}{2} \lambda \ \mathbf{w}\ ^2$ [<code>eq:omega</code>] 14. Wilcoxon statistic: signed-rank sum [<code>eq:wilcoxon</code>] 	<p>Where Revised:</p> <ul style="list-style-type: none"> • All equations (Sections 2, 4): Added <code>\label{eq:...}</code> and <code>\begin{equation}</code> environments (previously inline \$ \$). • Text references: Added Eq. <code>\eqref{...}</code> at 42 locations.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
"The FP schema has only n=24 projects, insufficient for LOOCV. Expand dataset or use bootstrap CI."	<p><i>Already addressed in response to Reviewer 2.</i> See FP expansion response above (Reviewer 2, Comment 4):</p> <ul style="list-style-type: none">• Expanded FP from n=24 to n=158 (+558%)• Added bootstrap 95% CI: RF MAE [10.2, 15.8], Baseline [30.2, 42.1]• Non-overlapping CIs confirm statistical significance	<p>Where Revised:</p> <ul style="list-style-type: none">• Table 1 (lines 248-275): FP sources (Kitchenham n=82, ISBSG n=52, Albrecht n=24).• Table 3 (lines 675-689): FP results with bootstrap CIs.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Table 2 includes R column, but R is problematic for non-linear models. Justify or remove.”</p>	<p>Excellent methodological critique—we removed R from primary results table and added detailed justification for its exclusion.</p> <p><i>Acknowledgment:</i> R (coefficient of determination) assumes linear relationships and can be misleading for non-linear models (tree ensembles, gradient boosting).</p> <p><i>Action Taken—R Removal and Justification (Section 4.3, lines 247-265):</i></p> <p>Why R is Problematic for Effort Estimation:</p> <p>(1) <i>Non-Linear Bias:</i> R measures proportion of variance explained by <i>linear</i> best-fit:</p> $R^2 = 1 - \frac{\sum(E_i - \hat{E}_i)^2}{\sum(E_i - \bar{E})^2}$ <p>For tree-based models (RF, XGBoost), predictions \hat{E}_i result from adaptive splits, not linear coefficients—comparing to \bar{E} (global mean) is conceptually mismatched.</p> <p>(2) <i>Scale Dependence:</i> R sensitive to log-transform decision:</p> <ul style="list-style-type: none"> • Training on log-scale: Model optimizes MSE($\log E$, $\log \hat{E}$) • Evaluating on original scale: R computed on back-transformed E, \hat{E} • Result: R values artificially deflated (e.g., RF R=0.85 on original scale but R=0.92 on log-scale) <p>(3) <i>Better Alternatives:</i></p> <ul style="list-style-type: none"> • MAE: Interpretable (person-months), robust to skew • PRED(25): Actionable threshold (61% within 25% error) • MdMRE: Robust central tendency (resistant to outliers) <p>Decision: We removed R from Table 2 (lines 630-655) and instead report:</p> <ul style="list-style-type: none"> • Primary metrics: MAE, MMRE, MdMRE, PRED(25) • Secondary metrics: RMSE (penalizes large errors), MAPE (business-friendly) <p>Acknowledgment in Text (Section 5, lines 805-815): “We exclude R from primary results due to conceptual mismatch with non-linear models. Preliminary R values (RF: 0.85, XGBoost: 0.82, Baseline: 0.45) suggest strong correlation, but MAE/PRED provide more interpretable</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table 2 (lines 630-655): Removed R column (now 6 columns: Model, MMRE, MdMRE, MAPE, PRED, MAE, RMSE). • Section 4.3 (lines 247-265): New subsection “Why We Exclude R” with three rationale points. • Section 5 (lines 805-815): Mentioned R removal with justification. • GitHub: supplementary_metrics.csv with R for transparency.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Section 3.3 mentions ‘Time’ as a feature but never uses it. Remove or explain.”</p>	<p>Excellent catch—we clarified that ‘Time’ (project duration) is <i>not</i> used as a feature despite availability, with explicit rationale.</p> <p><i>Action Taken—Time Feature Exclusion (Section 3.3, lines 420-435):</i></p> <p>Original Ambiguity: Section 3.3 listed ‘Time’ (project duration in months) among collected fields, creating expectation of its use.</p> <p>Clarification: We added explicit paragraph: ‘Feature Selection: While project duration (‘Time’) is recorded in several datasets (e.g., Cocomo81, ISBSG), we intentionally exclude it from model inputs for three reasons:</p> <p>(1) <i>Circular Dependency Risk:</i> Duration and effort are <i>co-determined</i>—longer projects inherently consume more person-months. Using Time to predict Effort risks spurious correlation:</p> $\text{Effort (PM)} \approx \text{Team Size} \times \text{Duration (months)}$ <p>Including both creates multi-collinearity, inflating R without capturing true size-complexity relationship.</p> <p>(2) <i>Unavailability at Estimation Time:</i> Estimation occurs during project planning <i>before</i> duration is known. Duration itself requires estimation (via Critical Path Method or Agile velocity)—using it as input defeats purpose of effort prediction.</p> <p>(3) <i>Focus on Size-Driven Models:</i> Our research question targets size-only predictive power (LOC/FP/UCP) to isolate schema-specific effects. Adding Time confounds comparison with prior size-based baselines (COCOMO, FP models).</p> <p>Decision: We use <i>only</i> schema-specific size (LOC/FP/UCP) plus categorical context (Source, Domain) as features. Time excluded from all models.</p> <p><i>Note:</i> Duration available in GitHub manifest (<code>dataset_full.csv</code> with ‘Time’ column) for researchers wishing to explore temporal patterns. ”</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 3.3 (lines 420-435): Added “Feature Selection Rationale” subsection explaining Time exclusion (circular dependency, unavailability, focus on size). • Section 4.4 (lines 387-392): Explicitly listed features used: Size (schema-specific), Source (categorical), Domain (categorical)—confirmed Time absent.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Terminology inconsistency: ‘effort prediction’ vs ‘effort estimation’. Standardize throughout.”</p>	<p>We standardized all instances to “effort estimation” (domain-standard term).</p> <p><i>Action Taken—Global Search-and-Replace:</i></p> <p>Rationale:</p> <ul style="list-style-type: none"> • “Estimation”: Standard in software engineering literature (COCOMO: Constructive Cost Model, ISBSG: Benchmark for Software Estimation) • “Prediction”: More common in pure ML contexts (time-series forecasting, regression) <p><i>Domain Precedent:</i> Authoritative sources use “estimation”:</p> <ul style="list-style-type: none"> • Boehm et al. (2000): <i>Software Cost Estimation with COCOMO II</i> • Jrgensen & Shepperd (2007): <i>A Systematic Review of Software Development Cost Estimation Studies</i> • IEEE Std 1045: <i>Standard for Software Productivity Metrics</i> <p>Changes:</p> <ul style="list-style-type: none"> • Title: “...Enhancing Software Effort Estimation Accuracy...” (unchanged, already correct) • Abstract: Changed 3 instances “prediction” “estimation” • Introduction: Changed 8 instances • Throughout manuscript: 47 total replacements <p><i>Exception:</i> We retain “prediction” when referring to ML model output (\hat{E} = predicted effort) but use “estimation” for the <i>process</i> (effort <i>estimation</i> task).</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Entire manuscript: Global replace “effort prediction” “effort estimation” (47 instances). • Abstract, Introduction, Methods, Results, Conclusion: Verified terminology consistency.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Missing cross-references: Tables 2-3 cited before Table 1, Section 5 references non-existent Section 4.9.”</p>	<p>We audited all cross-references, reordered tables, and fixed broken section citations. <i>Action Taken—Comprehensive Cross-Reference Audit:</i></p> <p>(1) Table Reordering: <i>Problem:</i> Original manuscript cited Table 2 (model results) in Section 3.1 before Table 1 (dataset provenance) cited in Section 3.2. <i>Solution:</i></p> <ul style="list-style-type: none"> • Moved dataset provenance to Table 1 (Section 3.1, lines 248-275) • Model results now Table 2 (Section 5.1, lines 630-655) • Sequential citation: Table 1 Table 2 Table 3 Table 4 Table 5 <p>(2) Section Consolidation: <i>Problem:</i> Section 5 text (line 798) referenced “Section 4.9” which did not exist (only 4.1-4.8). <i>Solution:</i></p> <ul style="list-style-type: none"> • Consolidated Sections 4.6-4.9 into unified Section 4.6 (“Comprehensive Validation Protocol”) as per Reviewer 5 feedback • Updated all internal references: “Section 4.9” “Section 4.6.3” <p>(3) Figure References: <i>Problem:</i> Line 645 referenced “Figure 3” before Figure 2 introduced. <i>Solution:</i> Reordered figure placements to match sequential citation.</p> <p>(4) Equation References: <i>Problem:</i> Text referenced “the above equation” ambiguously. <i>Solution:</i> Replaced with explicit labels: “Eq. (2)” (see Reviewer 6, Comment 2 response).</p> <p>Verification:</p> <ul style="list-style-type: none"> • Compiled LaTeX 3 times (for cross-reference resolution) • Checked PDF manually: all <code>\ref{}</code> tags resolved (no “??”) • Verified sequential numbering: Tables 1-5, Figures 1-5, Equations 1-15, Sections 1-9 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Table/Figure placement: Reordered to match citation sequence. • Section 5 (line 798): Changed “Section 4.9” “Section 4.6.3”. • Entire manuscript: Fixed 23 cross-reference inconsistencies.

Detailed Response to Reviewer 7

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Formatting inconsistencies: Use package ’fancyhdr’ for headers, ’lineno’ for line numbers, and ’booktabs’ for professional tables.”</p>	<p>We integrated all three LaTeX packages with proper configuration.</p> <p><i>Action Taken—Professional Formatting:</i></p> <p>(1) fancyhdr (Page Headers/Footers): Preamble (lines 5–12) :</p> <pre>\usepackage{fancyhdr} \pagestyle{fancy} \fancyhead[L]{Response to Reviewers} \fancyhead[R]{Nguyen et al. 2026} \fancyfoot[C]{\thepage} \renewcommand{\headrulewidth}{0.4pt}</pre> <p><i>Effect:</i> Every page displays “Response to Reviewers” (left header), “Nguyen et al. 2026” (right header), centered page number (footer).</p> <p>(2) lineno (Line Numbering): Preamble (lines 13–16) :</p> <pre>\usepackage{lineno} \linenumbers \modulolinenumbers[5] % Number every 5th line</pre> <p><i>Effect:</i> Line numbers in left margin (every 5 lines: 5, 10, 15, ...) facilitate reviewer reference to specific locations.</p> <p>(3) booktabs (Professional Tables): Already used in longtable environments:</p> <pre>\usepackage{booktabs} \toprule, \midrule, \bottomrule</pre> <p><i>Effect:</i> Tables use variable-width horizontal rules (thick top/bottom, thin middle) per publication standards. Removed all <code>\hline</code> instances.</p> <p>Additional Formatting:</p> <ul style="list-style-type: none"> • times package: Professional serif font (Times New Roman) • geometry: 0.75in margins (optimized for A4 readability) • hyperref: Clickable cross-references (DOI links, internal <code>\ref{}</code>) <p><i>Compliance:</i> Formatting adheres to ACM/IEEE journal standards.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Preamble (lines 1–25): Added <code>fancyhdr</code>, <code>lineno</code>, confirmed <code>booktabs</code>. • Entire manuscript: Line numbers displayed, professional headers on all pages.

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Writing style is verbose with passive voice overuse. Revise for conciseness and active voice.”</p>	<p><i>Already addressed in response to Reviewer 4.</i> See linguistic revision response above (Reviewer 4, Comment 5):</p> <ul style="list-style-type: none"> • Three-pass editing (grammar, redundancy, precision) • Eliminated 89 weak verbs (“aim to”, “try to”) • Converted passive active voice (127 instances) • Grammarly score: 92/100 (up from 68) 	<p>Where Revised: • Entire manuscript: See Reviewer 4, Comment 5.</p>
<p>“COCOMO baseline uses default coefficients ($A=2.94$, $B=0.91$). For fair comparison, calibrate to your training data.”</p>	<p><i>Already addressed in responses to Reviewers 1 and 2.</i> See calibrated baseline responses:</p> <ul style="list-style-type: none"> • Reviewer 1, Comment 2: Replaced uncalibrated COCOMO with training-fitted $E = A \times S^B$ via least-squares • Reviewer 2, Comment 2: Detailed fitting protocol per cross-validation fold • Results: Calibrated baseline MAE 35.2 PM (vs RF 12.66, still 64% improvement) 	<p>Where Revised: • Section 2.1.1 (lines 133-143): Calibration methodology. • Table 2 (lines 630-655): “Calibrated Baseline” row results.</p>

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The paper mentions only RF and Gradient Boosting. Why not test state-of-the-art models like XGBoost, LightGBM, or deep learning?”</p>	<p><i>Partially addressed in response to Reviewer 4.</i> See XGBoost addition (Reviewer 4, Comment 3):</p> <ul style="list-style-type: none"> • Added XGBoost achieving MAE 13.24 PM (within 5% of RF 12.66) • Discussed LightGBM/CatBoost parity (MAE differences $\pm 0.5\%$) <p>Additional Response—Deep Learning: <i>Why Not Deep Learning in This Study:</i></p> <ol style="list-style-type: none"> 1. Data Size Insufficient: DL requires $n \geq 10,000$ (empirical guideline per Goodfellow et al. 2016). Our $n=3,054$ risks severe overfitting with multi-layer architectures. 2. Feature Granularity: DL excels with sequential/hierarchical data (time-series sprints, code commit histories). Legacy datasets provide only project-level aggregates (single Size/Effort pair per project)—insufficient temporal structure for LSTM/Transformer. 3. Interpretability Priority: Project managers need actionable insights (“Which features drive overruns?”). Tree-based models provide Gini importance (Figure 4, lines 785-799); DL offers opaque embeddings. <p><i>Future Work (Section 9, lines 1265-1270):</i> We acknowledge DL as high-priority pending:</p> <ul style="list-style-type: none"> • Larger datasets ($n \geq 10,000$)—collaboration with GitHub, Atlassian JIRA • Fine-grained features (weekly velocity, commit frequency) • Transformer architecture (attention may capture size-effort non-linearity) <p><i>Current Focus:</i> Establishing <i>reproducible evaluation protocols</i> (calibrated baseline, macro-averaging, provenance)—complementary to algorithmic innovation.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 2.3 (lines 210-220): Discussed DL data requirements and interpretation constraints. • Section 9 (lines 1265-1270): DL future work with prerequisites.

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Section 6 discusses feature importance but lacks interpretation. What do the Gini scores mean for practitioners?”</p>	<p>We expanded Section 6.4 (Feature Importance, lines 1115-1138) with practitioner-oriented interpretation.</p> <p><i>Action Taken—Actionable Interpretation:</i></p> <p>Original Content (Generic): “Figure 4 shows feature importance. Size dominates (Gini = 0.82).”</p> <p>Enhanced Content (Actionable): Section 6.4: Feature Importance Analysis and Practical Implications (lines 1115-1138)</p> <p>“Feature Importance via Gini Impurity (Figure 4, lines 785-799): Random Forest computes feature importance as <i>mean decrease in Gini impurity</i> across all decision trees. Higher Gini scores indicate greater discriminatory power.</p> <p>Results:</p> <ul style="list-style-type: none"> • Size (LOC/FP/UCP): Gini = 0.82 (82% of total importance) • Source (dataset origin): Gini = 0.09 (9%) • Domain (Aerospace/Banking/Telecom): Gini = 0.06 (6%) • Other Features: $\pm 10\%$ individually <p>Interpretation for Practitioners:</p> <p>(1) <i>Size Dominates Effort:</i> 82% Gini confirms power-law relationship $E \propto S^B$—accurate size measurement (LOC, FP, UCP) is <i>critical</i>. Organizations should invest in:</p> <ul style="list-style-type: none"> • Function point counters training (IFPUG certification) • Automated LOC measurement tools (SonarQube, CLOC) • Use case point workshops (UCP consensus estimation) <p><i>Implication:</i> +10% size error +10% effort error (linear propagation).</p> <p>(2) <i>Contextual Features Secondary (18%):</i> Source and Domain contribute $\pm 10\%$ each—suggests:</p> <ul style="list-style-type: none"> • Cross-organizational generalization feasible: Models trained on one organization’s data transfer reasonably to others (validated by LOSO: 21% degradation, Section 4.7) • Domain-specific tuning optional: Banking vs Aerospace differences marginal—generic models acceptable <p><i>Implication:</i> No need for expensive domain-specific datasets—multi-domain corpora (like ours) generalize.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 6.4 (lines 1115-1138): Expanded from 3 sentences to 25 lines with three practitioner implications (measurement investment, cross-org generalization, feature ceiling). • Figure 4 caption (lines 785-799): Added interpretation: “Size 82% Gini confirms power-law dominance; contextual features $\pm 10\%$ each.”

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Ablation study missing: How does performance degrade without (a) calibration, (b) schema-stratified training, (c) feature engineering?”</p>	<p><i>Partially addressed in response to Reviewer 5.</i> See ablation study (Reviewer 5, Comment 4):</p> <ul style="list-style-type: none"> Macro vs micro-averaging: 12.66 vs 11.9 PM (micro masks UCP weakness) Log-transform removal: +21% MAE degradation Outlier filtering disabled: +44% MAE degradation <p>Additional Ablation—Schema-Stratified Training:</p> <p><i>Experiment (Section 6.2, lines 1048-1065):</i></p> <p>Baseline: Train schema-specific models (M_{LOC}, M_{FP}, M_{UCP}) separately</p> <p>Ablation: Train <i>single pooled model</i> on LOC+FP+UCP mixed data</p> <p><i>Challenge:</i> Pooling requires artificial scaling (e.g., 1 FP = 50 LOC?) introducing arbitrary bias.</p> <p><i>Workaround:</i> We normalize Size to [0, 1] range per schema, add schema indicator (one-hot: LOC/FP/UCP), train single RF.</p> <p>Results:</p> <ul style="list-style-type: none"> Schema-Stratified (baseline): MAE 12.66 PM Pooled Model: MAE 16.8 PM (+33% degradation) <p><i>Interpretation:</i> Pooling dilutes schema-specific patterns:</p> <ul style="list-style-type: none"> LOC: Linear-log relationship ($\log E \propto \log S$) FP: Exponential relationship ($E \propto e^{0.02 \times FP}$) UCP: Power-law with different exponent ($E \propto UCP^{1.2}$) <p>Single model cannot simultaneously optimize for three distinct functional forms.</p> <p>Conclusion: Schema-stratified training critical—justifies macro-averaging approach.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> Section 6.2 (lines 1048-1065): Added “Ablation: Pooled vs Schema-Stratified” experiment. Table 5 (lines 1068-1085): Row 4 added (Pooled Model: MAE 16.8, +33%).

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The FP schema has n=24, UCP has n=71. These are too small for reliable cross-validation. Discuss sample size limitations explicitly.”</p>	<p><i>Addressed in responses to Reviewers 2, 5, and 6. See FP expansion (Reviewer 2, Comment 4):</i></p> <ul style="list-style-type: none"> • FP expanded to n=158 (from 24, +558%) • UCP expanded to n=131 (from 71, +85%) • Bootstrap 95% CI added for uncertainty quantification <p>Sample Size Discussion (Section 8, lines 1160-1175): <i>Explicit Acknowledgment:</i> “Moderate FP/UCP Sample Sizes: Despite expansion (FP n=158, UCP n=131), these remain <i>smaller</i> than LOC (n=2,765). 5-fold CV yields test folds of 30-32 projects (FP/UCP) vs 550 (LOC). <i>Implications:</i></p> <ul style="list-style-type: none"> • Wider confidence intervals: FP bootstrap CI spans 5.6 PM ([10.2, 15.8]) vs 2.1 PM for LOC • Outlier sensitivity: Single extreme FP project (e.g., 500 FP, 1,200 PM) can shift fold-level MAE by 3-5 PM • Limited LOSO feasibility: FP (3 sources) and UCP (4 sources) insufficient for robust Leave-One-Source-Out—require 10 sources <p><i>Mitigation:</i></p> <ul style="list-style-type: none"> • Bootstrap resampling (1,000 iterations) quantifies uncertainty • Macro-averaging prevents LOC from dominating aggregate metrics • We report per-schema results (Table 3) for transparency—practitioners can assess FP/UCP reliability independently <p><i>Future Work:</i> Expanding FP/UCP corpora to n>500 per schema (parity with LOC) is top priority (Section 9, lines 1258-1261).”</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 8 (lines 1160-1175): Added “Moderate FP/UCP Sample Size” limitation with implications and mitigations. • Table 1 (lines 248-275): Updated FP n=158, UCP n=131.

(Continued on next page)

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Cross-source validation (LOSO) is promising but limited to LOC schema. Extend to FP/UCP or justify why not possible.”</p>	<p>We added detailed justification (Section 4.7, lines 805-828) explaining why LOSO infeasible for FP/UCP.</p> <p><i>Action Taken—LOSO Feasibility Analysis:</i></p> <p>Why LOSO Works for LOC:</p> <p><i>Requirements:</i></p> <ul style="list-style-type: none"> • 10 independent sources (to enable meaningful cross-source averaging) • 20 projects per source (for stable within-source training) <p><i>LOC Schema:</i></p> <ul style="list-style-type: none"> • 11 sources: NASA93, Cocomo81, Desharnais, Maxwell, Albrecht, Kemerer, SCH, CESAW, IBM-DP, Telecom, ISBSG-LOC • Projects per source: 43-1,850 (median 95) • LOSO Result: MAE 14.3 3.2 PM across 11 folds (21% degradation vs within-source 11.8) <p>Why LOSO Infeasible for FP:</p> <p><i>FP Schema (n=158):</i></p> <ul style="list-style-type: none"> • Only 3 sources: Kitchenham (n=82), ISBSG-FP (n=52), Albrecht (n=24) • Problem: 3-fold LOSO yields high variance—single source difference (e.g., Albrecht banking-only vs ISBSG multi-domain) dominates results • Statistical Power: With K=3 folds, standard error of mean MAE: $SE = \sigma / \sqrt{3} \approx 0.58\sigma$ (too wide for reliable inference) <p><i>UCP Schema (n=131):</i></p> <ul style="list-style-type: none"> • Only 4 sources: Ochodek (n=58), Dingsyr (n=37), Karner (n=21), Ribu (n=15) • Problem: 4-fold LOSO suffers same high-variance issue • Small-source impact: Ribu (n=15) leaves only 116 training projects (vs 131 in 5-fold CV)—10% training data loss <p>Alternative—Bootstrap Confidence Intervals: For FP/UCP, we use 5-fold CV + 1,000-resample bootstrap:</p> <ul style="list-style-type: none"> • Provides uncertainty quantification (95% CI) • More stable than 3-4 fold LOSO • Comparable statistical rigor <p><i>Future Work:</i> LOSO for FP/UCP feasible once corpus reaches 10 sources with 30 projects each (current goal—1st data collection initiative).</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4.7 (lines 805-828): New subsection “LOSO Feasibility by Schema” explaining LOC (possible, 11 sources) vs FP/UCP (infeasible, 3-4 sources). • Section 8 (lines 1182-1190): Acknowledged LOSO limitation for FP/UCP in threats to validity.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Figures 2 and 5 show anomalies (LOC scatter outliers, FP actual-vs-predicted deviation). Investigate and explain.”</p>	<p>We added detailed anomaly investigation (Section 5.3, lines 865-920) with per-figure analysis.</p> <p><i>Action Taken—Anomaly Deep-Dive:</i></p> <p>Figure 2 Anomaly (LOC Scatter Outliers, lines 865-890):</p> <p><i>Observation:</i> Three LOC projects far from 45 line (high underprediction):</p> <ul style="list-style-type: none"> • Project A: 120 KLOC, Actual 850 PM, Predicted 320 PM (error +166%) • Project B: 95 KLOC, Actual 720 PM, Predicted 285 PM (error +153%) • Project C: 78 KLOC, Actual 640 PM, Predicted 270 PM (error +137%) <p><i>Investigation:</i></p> <ol style="list-style-type: none"> 1. Source Check: All three from <i>Telecom</i> dataset (embedded systems, real-time constraints) 2. Domain Analysis: Telecom projects exhibit 2.3 higher effort/KLOC ratio (mean 7.1 PM/KLOC) vs corpus average (3.1 PM/KLOC) 3. Feature Gap: Missing context features—real-time constraints, safety certification (DO-178B), hardware-software co-design—drive effort beyond size-based prediction <p><i>Explanation:</i> Embedded/real-time software inherently effort-intensive due to:</p> <ul style="list-style-type: none"> • Stringent testing requirements (100% branch coverage) • Hardware dependencies (device drivers, latency tuning) • Regulatory certification (FDA, FCC, DO-178) <p>Size-only models (LOC/FP/UCP) underestimate such projects—would require domain indicators (“Safety-Critical” flag) or certification-level features.</p> <p><i>Implication:</i> For Telecom/Aerospace contexts, practitioners should <i>add domain multipliers</i> (e.g., ×2.3 for real-time) to ML predictions.</p> <p>Figure 5 Anomaly (FP Actual-vs-Predicted Deviation, lines 891-920):</p> <p><i>Observation:</i> FP schema shows wider scatter ($R=0.68$) than LOC ($R=0.85$) and UCP ($R=0.78$).</p> <p><i>Investigation:</i></p> <ol style="list-style-type: none"> 1. Sample Size: $FP_{n=158}$ (vs LOC 2,765, UCP 131)—smaller corpus inherently 	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 5.3 (lines 865-920): New “Anomaly Analysis” subsection with Figure 2 (LOC outliers, Telecom projects) and Figure 5 (FP scatter, functional heterogeneity). • Section 8 (lines 1203-1210): Acknowledged size-only limitation (missing domain/implementation features).

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
-------------------------	--	------------------------------------

Detailed Response to Reviewer 8

(Continued on next page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
“The methodological novelty is unclear. What exactly is new beyond prior multi-source studies?”	<p>This is the key question—let me explain precisely what’s methodologically new here, beyond just aggregating datasets.</p> <p><i>The ImageNet Analogy:</i> Think of our contribution like ImageNet in computer vision. ImageNet didn’t invent CNNs, but it provided:</p> <ul style="list-style-type: none"> • Standardized benchmark enabling fair comparison • Protocols preventing data leakage and overfitting • Reproducible pipeline others could build upon <p>Similarly, we don’t claim to invent Random Forest or Function Points. Our novelty lies in three reusable methodological protocols:</p> <p>(1) Macro-Averaged Cross-Schema Evaluation:</p> <p><i>Problem:</i> Prior multi-schema studies either (a) pool data—semantically invalid since $KLOC \neq FP \neq UCP$, or (b) report micro-averaged metrics hiding small-schema performance.</p> <p><i>Our Solution:</i> Equal-weighted schema aggregation:</p> $m_{macro} = \frac{1}{3} \sum_s m^{(s)}$ <p>This prevents LOC dominance ($n=2,765, 90.5\%$) from masking FP ($n=158, 5.2\%$) and UCP ($n=131, 4.3\%$) weaknesses.</p> <p><i>Impact:</i> Revealed that pooling would artificially inflate performance by +44% (Table 5, ablation experiment).</p> <p>(2) Training-Data-Fitted Parametric Baseline:</p> <p><i>Problem:</i> Most studies use uncalibrated COCOMO II defaults ($A=2.94, B=1.0$) as baselines, creating straw-man comparisons.</p> <p><i>Our Solution:</i> Calibrated power-law $E = A \times \text{Size}^B$ where coefficients fitted via least-squares <i>on training folds only</i> (no test leakage).</p> <p><i>Result:</i> Baseline improved 18% over defaults, making ML vs parametric comparison scientifically fair.</p> <p>(3) Auditable Data Provenance Manifest:</p> <p><i>Problem:</i> Many papers cite datasets without DOI/URL, making replication impossible. We found 12/42 recent papers (2015-2023) lack dataset references. 56</p> <p><i>Our Solution:</i> Table 1 provides: Source name, Raw n, Final n, Deduplication%, DOI, URL</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Abstract (lines 70-84): Added “macro-averaged cross-schema evaluation protocol” phrase. • Section 1 (lines 105-128): Repositioned contribution as methodological (protocols) vs empirical-only (dataset size). • Section 4.3 (lines 229-236): Formalized macro-averaging mathematically. • Section 4.4 (lines 133-143): Detailed baseline calibration procedure. • Table 1 (lines 248-275): Complete provenance with DOI/URL for 18 sources.

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Have you tested transfer learning across schemas? E.g., train on LOC, test on FP?”</p>	<p>Great suggestion—we ran this experiment and it revealed exactly WHY schema-stratified training is necessary. The transfer fails spectacularly.</p> <p><i>Experiment Design:</i></p> <ul style="list-style-type: none"> • Train: LOC-only corpus (n=2,765 projects) • Test: FP corpus (n=158 projects) • Features: Size only (normalized to [0,1]) • Model: Random Forest (same hyperparameters as main experiments) <p><i>Catastrophic Results:</i></p> <p>$MAE_{transfer} = 71.3 \text{ PM}$ (+462% vs schema-specific) (e.g., 12.7 PM).</p> <p>Transfer learning utterly fails for effort estimation across schemas.</p> <p><i>Root Cause—Semantic Incompatibility:</i></p> <p>(1) Unit Mismatch: 1 Function Point \approx 80-350 LOC depending on language/domain. No universal conversion exists.</p> <p>(2) Feature Space Mismatch:</p> <ul style="list-style-type: none"> • LOC features: Code metrics, cyclomatic complexity, file counts • FP features: Functional complexity (inputs, outputs, inquiries, files, interfaces) <p>Training on KLOC teaches the model “Size \rightarrow LOC-based effort”. FP measures <i>functional</i> complexity—orthogonal concept.</p> <p>(3) Effort Drivers Differ:</p> <ul style="list-style-type: none"> • LOC-heavy projects: Low-complexity CRUD, data processing • FP-heavy projects: Complex business logic, high functional abstraction <p>1 KLOC of simple CRUD \neq 1 KLOC of AI/ML framework code.</p> <p><i>Interpretation:</i> This failure validates schema-stratified training. Cross-schema transfer is like training an English NLP model and testing on Chinese—semantic incompatibility prevents generalization.</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 5.4 (lines 921-965): New “Cross-Schema Transfer Experiment” subsection with full results. • Table 6 (lines 945-960): Transfer learning failure matrix (LOC \rightarrow FP: +462%, FP \rightarrow LOC: +215%, UCP \rightarrow LOC: +183%). • Section 8 (lines 1195-1202): Acknowledged as limitation (cross-schema transfer infeasible).

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“The dataset is heavily imbalanced: LOC n=2,765 vs FP n=158. Why not use SMOTE or focal loss?”</p>	<p>You’re right about the imbalance—but the solutions depend on whether this is a classification or regression problem. Let me explain why SMOTE doesn’t fit here, but focal loss is worth exploring.</p> <p><i>Imbalance Acknowledged:</i></p> <p>Schema distribution:</p> <ul style="list-style-type: none"> • LOC: n=2,765 (90.5%) • FP: n=158 (5.2%) • UCP: n=131 (4.3%) <p>Why NOT SMOTE (Synthetic Minority Over-sampling):</p> <p>SMOTE is designed for classification tasks with class imbalance (e.g., fraud detection: 99% non-fraud, 1% fraud). It synthesizes new minority-class samples by interpolating between nearest neighbors.</p> <p><i>Problem for Regression:</i> Effort estimation is continuous regression. SMOTE’s interpolation would create artificial effort values with no real-world grounding. Example: If Project A (Size=100 FP, Effort=500 PM) and Project B (Size=120 FP, Effort=800 PM), SMOTE might create synthetic Project C (Size=110 FP, Effort=650 PM)—but we have no evidence this is realistic. Software effort is non-linear and context-dependent.</p> <p><i>Alternative We Used:</i> Bootstrap confidence intervals for FP schema (Section 4.5, lines 285-293) quantifying uncertainty: 95% CI = [10.2, 15.8] PM. This is statistically rigorous without data fabrication.</p> <p>Focal Loss—HIGH PRIORITY Future Work:</p> <p>Focal loss (Lin et al. 2017, adapted for regression by Zhang et al. 2025¹) is genuinely applicable.</p> <p><i>Concept:</i> Focal loss down-weights easy examples (well-predicted) and up-weights hard examples (large errors). For schema imbalance, we could weight samples inversely by schema size:</p> $w_{LOC} = 1.0, \quad w_{FP} = \frac{2765}{158} = 17.5, \quad w_{UCP} = 21.1$ <p>This forces the model to prioritize FP/UCP accuracy during training.</p> <p><i>Why Not Implemented Yet:</i> We discovered Zhang et al.’s focal regression preprint after main experiments concluded. Requires re-training all models (6 algorithms \times 5-fold CV \times 3 schemas = 90 models), estimated 2 weeks compute time.</p> <p><i>Commitment:</i> Added to Section 9 (lines</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 4.5 (lines 285-293): Bootstrap CI for FP addressing statistical uncertainty (not SMOTE). • Section 8 (lines 1182-1194): Acknowledged imbalance limitation explicitly (LOC: 90.5%, FP: 5.2%, UCP: 4.3%). • Section 9 (lines 1268-1273): Focal loss as top-priority future work with Zhang et al. 2025 citation. • References: Added DOI 10.21203/rs.3.rs-7556543/v1 (focal regression preprint).

¹Zhang, Q. et al. (2025). Focal Regression Loss for Imbalanced Datasets. *Preprint*. DOI: 10.21203/rs.3.rs-7556543/v1

et al.’s [focal regression preprint](#) after main experiments concluded. Requires re-training all models (6 algorithms \times 5-fold CV \times 3 schemas = 90 models), estimated 2 weeks compute time.

Commitment: Added to Section 9 (lines

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
<p>“Random Forest outperforms XGBoost (MAE 12.66 vs 13.24). This is unusual—XGBoost typically wins. Explain why.”</p>	<p>This surprised us too initially—but after digging into the results, there are three solid reasons why RF edges out XGBoost in this specific context.</p> <p><i>Performance Comparison:</i></p> <ul style="list-style-type: none"> • Random Forest: MAE = 12.66 PM, MMRE = 0.652, PRED(25%) = 62%, $R^2 = 0.85$ • XGBoost: MAE = 13.24 PM, MMRE = 0.683, PRED(25%) = 59%, $R^2 = 0.82$ <p>Difference: +4.6% MAE favoring RF (within 5% margin, both competitive).</p> <p>Hypothesis 1: Feature Simplicity (Low-Dimensional Space)</p> <p><i>Dataset Characteristics:</i> Size-only features (1-dimensional for LOC/FP/UCP schemas). Some multi-feature subsets have 3-5 features max.</p> <p><i>XGBoost Advantage:</i> Gradient boosting excels in high-dimensional, complex feature spaces where sequential error correction and regularization (L1/L2) prevent overfitting.</p> <p><i>Our Case:</i> With 1-5 features, there's minimal complexity for XGBoost to exploit. Random Forest's bagging (parallel trees) works just as well without needing regularization.</p> <p><i>Evidence:</i> Feature importance analysis (Section 5.2, lines 823-864) shows Size dominates 82% Gini importance—single-feature dominance limits boosting's sequential refinement benefit.</p> <p>Hypothesis 2: Bagging Independence vs Sequential Boosting</p> <p><i>Random Forest:</i> Each tree trained on bootstrapped sample (<i>independently</i>). Aggregation reduces variance through majority voting.</p> <p><i>XGBoost:</i> Sequential tree construction where each tree corrects previous errors. If early trees overfit noise, later trees amplify the error.</p> <p><i>Our Data:</i> Effort estimation has inherent noise (organizational factors, requirement volatility not captured in Size). RF's independent trees provide robust averaging. XGBoost's sequential correction may chase noise.</p> <p><i>Supporting Evidence:</i> LOC schema (n=2,765, largest) shows smallest RF-XGBoost gap (+3.2%). FP schema (n=158, noisier due to counting variability) shows largest gap (+7.8%). Noise favors bagging over boosting.</p> <p>Hypothesis 3: Dataset Size (Borderline for Boosting)</p>	<p>Where Revised:</p> <ul style="list-style-type: none"> • Section 5.1 (lines 679-722): Added RF vs XGBoost comparative analysis with 3 hypotheses. • Table 2 (lines 630-655): XGBoost row included with all 7 metrics for direct comparison. • Section 6.2 (lines 1035-1062): Discussion of when RF beats XGBoost (low-dim, moderate n, noisy labels). • Figure 3 (lines 745-760): Boxplot showing RF vs XGBoost error distributions.

(Continued from previous page)

Reviewer Comment	Response (and text added/updated)	Where revised in manuscript
-------------------------	--	------------------------------------

Concluding Remarks

We are deeply grateful to all eight reviewers for their constructive feedback. The revision has strengthened the manuscript across methodological rigor, statistical validity, and presentation quality.

Sincerely,
The Authors