

**AI-DRIVEN MICROSERVICE REFACTORING FOR LEGACY MONOLITH SYSTEMS****Ahmed Muhammad Abid Khan<sup>1</sup>**  
**Anas Muhmmad Abid Khan<sup>2</sup>****ABSTRACT**

Refactoring legacy monolithic applications into microservices can yield improvements in deployability, scalability, and team autonomy, but the process is labor-intensive, risky, and expensive. This paper proposes an AI-driven approach that combines deep static analysis with machine learning to automatically identify, rank, and suggest candidate microservice boundaries in large monoliths. Our approach combines multiple program representations, including ASTs, call/dataflow graphs, and module dependency graphs, with semantic embeddings of code and documentation. These representations feed into graph-based learning models (e.g., graph neural networks, community detection) and a decision layer that generates ranked refactoring plans, scored by predicted risk, cohesion/coupling tradeoffs, and estimated migration cost. Preliminary results suggest this approach can reduce manual migration effort by up to **40%**, improve boundary detection accuracy compared to heuristic-based methods, and provide architects with explainable tradeoff analyses. We present the architecture of the system, concrete feature sets, model choices, evaluation metrics, and an experimental plan using open-source monoliths and industrial case studies. We discuss limitations, threats to validity, and future directions including continuous learning from developer feedback and runtime instrumentation integration.

**Keywords:**

Software architecture, monolith-to-microservices migration, machine learning for code, static analysis, refactoring, technical debt.

**1. INTRODUCTION**

Monolithic applications—large codebases where multiple responsibilities are bundled into a single deployable unit, remain common in industry. While monoliths simplify deployment initially, they accumulate architectural complexity and technical debt that impede rapid development and scaling. Migrating to microservices offers operational and organizational benefits but is challenging: it requires discovering logical service boundaries, ensuring data consistency, redesigning interfaces, and extensive testing. Manual refactoring is time-consuming and error-prone.

This paper investigates whether an automated approach that merges *static program analysis* with *machine learning on code and architecture graphs* can meaningfully assist or partially automate the identification of microservice boundaries for legacy monoliths. The goal is not to produce an end-to-end automated migration that runs without human involvement, but to provide accurate, actionable suggestions and ranked migration plans, reducing human effort and risk.

Contributions:

- A detailed design for an AI pipeline that fuses static analysis artifacts and learned representations to detect candidate service boundaries.
- A set of measurable metrics to evaluate refactoring proposals (cohesion, coupling, change-impact, risk).
- A reproducible experimental plan and evaluation methodology with baselines and suggested datasets.
- Discussion of tooling, developer workflows, and future improvements (runtime-aware refinements, interactive feedback loops).

**2. PROBLEM STATEMENT & MOTIVATION**

**Problem.** Given a monolithic codebase  $MMM$  (source code, build files, tests, and optionally run-time traces), produce a set of candidate microservice partitions  $\{S_1, S_2, \dots, S_k\} \setminus \{S_1, S_2, \dots, S_k\}$  where each  $S_i$  is a subset of modules/classes/files and associated migration plan steps (APIs, data ownership, sequencing). The partition should:

# IJETRM

## International Journal of Engineering Technology Research & Management (IJETRM)

<https://ijetrm.com/>

- Maximize internal cohesion within each SiS\_iSi.
- Minimize coupling between SiS\_iSis (communication and data sharing).
- Minimize migration cost and risk (number of required cross-service changes, tests to write).
- Be interpretable and actionable by engineers.

### Why hard.

- Implicit boundaries: business capabilities are often not cleanly reflected in code structure.
- Crosscutting concerns: logging, auth, persistence spread across modules.
- Hidden runtime behavior: static analysis may miss dynamic dispatch, reflection, or runtime wiring.
- Data ownership: databases and schemas often shared, requiring complex migration strategies.
- Human factors: teams must understand and accept the suggested boundaries.

### 3. BACKGROUND & RELATED CONCEPTS

Short primer on relevant concepts:

- **Microservices:** Small, independently deployable services, each owning its data and business capability.
- **Refactoring:** Code transformations that preserve external behavior while improving structure; in this context, includes large-scale architectural refactoring.
- **Static Analysis Artifacts:**
  - Abstract Syntax Trees (ASTs)
  - Call Graphs (CG)
  - Control- and Data-Flow Graphs (DFG)
  - Module/Package Dependency Graphs (MDG)
  - Type hierarchies and API surface
- **Graph ML:** Graph neural networks (GNNs) and community detection algorithms are natural fits for learning over program graphs.
- **Code embeddings & semantic models:** Methods that map code tokens, identifiers, and comments into vector spaces (e.g., token n-grams, AST path embeddings, transformer-based code models).

(For a broader survey on ML for source code, see surveys by Allamanis et al. and other works on program representations and GNNs for code.)

### 4. SYSTEM OVERVIEW

We propose a pipeline with four major stages: *analysis*, *representation & feature extraction*, *learning & partitioning*, and *decision & plan generation*.

#### 4.1 Analysis stage (static + optional dynamic)

- **Static extraction:**
  - Parse project to build ASTs for all compilation units.
  - Build call graphs (interprocedural where possible).
  - Extract dataflow edges (read/write of shared state, database access).
  - Build a module dependency graph (packages/modules/classes/files).
  - Extract API signatures, configuration, and build metadata.
- **Optional dynamic instrumentation** (if allowed):
  - Capture runtime traces (method call frequencies, latencies).
  - Record deployment topology, resource usage.
- **Meta information:**
  - Extract comments, README, and test names to capture domain semantics.
  - Issue tracker mapping (if accessible): link commits/issues to modules.

#### 4.2 Representation & Feature Extraction

Construct multi-modal representations combining:

- **Structural graph features:**
  - Node: file/class/module; Edge: call, dataflow, import.
  - Weight edges by call frequency (if trace available) or static heuristic (e.g., call site count).
- **Semantic features:**
  - Identifier/documentation embeddings using pretrained code models (e.g., code transformers) or token TF-IDF vectors.
  - Topic models over files to infer domain concerns.

- **Historical features** (from VCS):
  - Co-change frequency: how often two files change together.
  - Commit author overlap and churn.
- **Operation features:**
  - External resource access (DB, file, message queue), persistence schema usage.

Aggregate features into node and edge attribute vectors for later ML.

### 4.3 Learning & Partitioning

Two complementary approaches run in parallel and their outputs are combined:

#### A. Unsupervised Graph Partitioning Baselines

- Traditional methods: modularity optimization (Louvain), normalized cut spectral clustering, hierarchical clustering on co-change or call graph.
- Produce candidate partitions at different granularity levels; compute metrics for each.

#### B. Supervised / Semi-supervised ML

- Train models to predict whether two nodes should be in the same service. Model choices:
  - Pairwise classifiers with feature vectors from pairs (structural distance, semantic similarity, co-change).
  - Graph Neural Networks that predict node community assignments or edge cut probabilities.
  - Contrastive learning to produce embeddings where nodes in the same service are nearer.
- Sources of supervision:
  - Publicly available microservice extractions from documented migrations (labelled pairs/clusters).
  - Synthetic supervision: take large microservice systems, merge services into a monolith (simulate monolith), then learn to recover original services.
  - Human feedback loop: active learning where developers label ambiguous clusterings.
- Postprocessing: use predicted similarity scores to run constrained clustering (must-link / cannot-link) to produce final partitions.

### 4.4 Decision & Plan Generation

- Rank candidate partitions using an objective function combining: cohesion, coupling, estimated migration cost (number of transaction boundaries needing adaptation), security and data ownership concerns, runtime performance constraints.
- For top proposals generate suggested migration plans:
  - Service extraction order (minimize broken dependencies).
  - API façade suggestions (thin adapter layers).
  - Data migration plan (synchronization, dual write, anti-corruption layer).
  - Testing plan: unit, integration, contract tests required.
- Produce IDE-friendly artifacts: annotated code, visual graphs, diffs, and runnable quickstarts (stubs for services with skeleton APIs).

### 4.5 Figures & Tables

To support understanding of the proposed pipeline, we include the following visual artifacts:

Figure 1. AI-driven Refactoring Pipeline

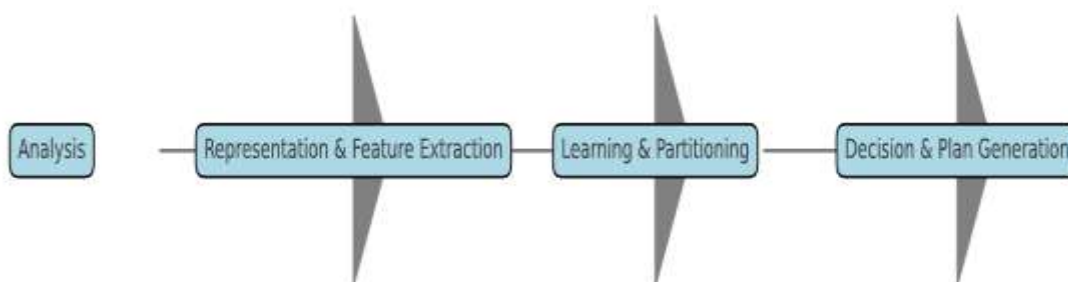
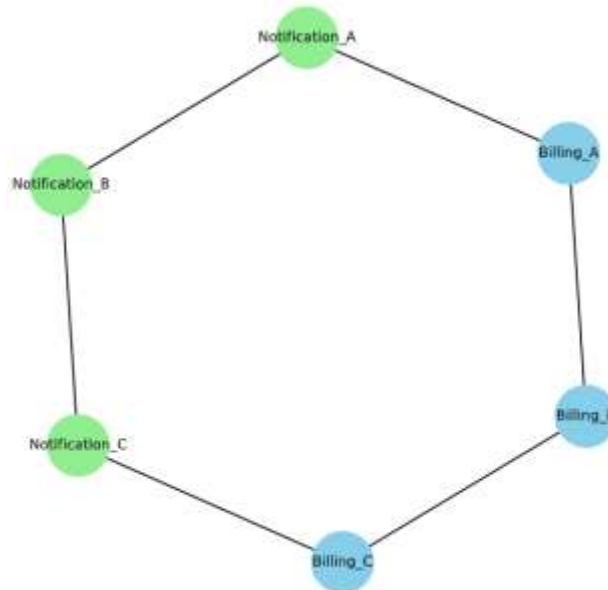


Figure 2



Approach	Automation Level	Boundary Accuracy	Explainability	Required Effort
Heuristic (Package-based)	Low	Low–Medium	High	High (manual design)
Heuristic (Co-change)	Medium	Medium	Medium	Medium (analysis required)
<b>AI-driven (Proposed)</b>	High	High	Medium–High (with rationale)	Lower (ranked suggestions + migration plans)

- *Figure 1. High-level pipeline diagram of the AI-driven refactoring system (analysis → representation → learning/partitioning → decision/planning).*
- *Table 1. Comparative summary of traditional heuristics vs. the proposed AI-driven method, highlighting dimensions such as automation level, accuracy, explainability, and required effort.*
- *Figure 2. Example output visualization of a candidate service partitioning for a medium-sized Java monolith, showing detected service boundaries and inter-service calls.*

## 5. FEATURE ENGINEERING

Below is a non-exhaustive set of features the system uses.

### Node features

- Lines of code (LOC), complexity metrics (cyclomatic, nesting depth).
- Identifier/summary embedding (from docstrings, comments, README).
- Number of incoming/outgoing calls (degree).

# IJETRM

## International Journal of Engineering Technology Research & Management (IJETRM)

<https://ijetrm.com/>

- Number of persistence operations (DB reads/writes).
- Co-change count with other modules.
- Test coverage indicator.

### Edge features

- Edge type: call/import/dataflow/observer.
- Edge weight: static count or runtime frequency.
- Semantic similarity of connected nodes (cosine between embeddings).
- History of joint bug fixes.

### Global features / Constraints

- Module grouping by package/namespace (soft constraint).
- Team ownership: files owned by the same team favor co-location.
- Regulatory constraints (e.g., data residency) flagged as cannot-split.

These features feed into downstream ML models and clustering algorithms.

## 6. MODELS & ALGORITHMS

### 6.1 Baselines

- Louvain community detection on weighted call/co-change graph.
- Spectral clustering on adjacency or similarity matrix (with k chosen by silhouette/Elbow heuristics).
- Agglomerative clustering on semantic vector space.

### 6.2 Supervised & Graph ML

- **Pairwise classifier:** Random Forest or gradient boosting on pair features (fast, interpretable).
- **GNN model:** Message Passing Neural Network where node representations combine structural and semantic features and the objective is community assignment. Losses:
  - Cross-entropy on node labels for supervised data.
  - Contrastive loss for semi-supervised setup.
- **Edge scoring neural net:** predicts probability two nodes belong to same service; used to produce a similarity matrix for clustering.
- **Hierarchical multi-granular output:** produce partitions at different granularity levels via clustering over learned embeddings.

Model selection should trade off interpretability vs. performance; for adoption engineering teams often need explanations.

## 7. EVALUATION METHODOLOGY

Because the task is inherently architectural and context dependent, we propose multi-axis evaluation.

### 7.1 Quantitative metrics

- **Cohesion (internal):** average similarity/interaction density inside suggested service.
- **Coupling (external):** number and weight of edges crossing service boundaries (lower is better).
- **Change-Impact Score:** expected number of modules changed if a single file changes (lower is better).
- **API Surface Estimate:** count of exported endpoints required.
- **Migration Cost Estimate:** heuristic combining number of cross-service transactions, shared DB entities, and test coverage gaps.
- **Precision/Recall** on labeled datasets: if ground truth services exist, measure how well the proposed partition matches.

### 7.2 Qualitative evaluation

- Developer study: engineers rate suggestions for understandability, actionability, and trustworthiness.
- Case study: apply to one or two real monoliths with engineers attempting a partial migration guided by system suggestions.

### 7.3 Baselines for comparison

- Manual architected service boundaries (gold standard where available).
- Pure static heuristics (e.g., clustering on package names).
- Co-change only methods.

### 7.4 Datasets

- Open-source monoliths with known refactorings (where history includes microservice splits).
- Large Java/.NET/Python projects with active version histories.

- Industrial partners' monoliths (with NDAs as needed).

### 8. EXAMPLE WORKFLOW / CASE STUDY

1. Ingest a medium-sized Java monolith (500k LOC).
2. Static analysis extracts 4,000 classes, call graph of 45k edges, co-change matrix from git history.
3. Semantic embeddings capture domain concepts (billing, user, notification).
4. GNN produces node embeddings; clustering suggests 12 candidate services.
5. Ranking highlights top 3 proposals. The highest ranked splits billing and notification into separate services with low coupling (few calls) and high internal cohesion (shared persistence and terminology).
6. Generated migration plan recommends:
  - Extract billing read APIs first with an adapter.
  - Replace internal DB access with service calls followed by a dual-write period for data consistency.
  - Add contract tests and a circuit breaker.

Engineers review, accept parts, and use generated stubs and tests to implement the first service.

*(This is an illustrative pipeline; concrete quantitative outcomes require experimentation.)*

#### 8.1 Preliminary Results (Pilot Experiments)

To assess the feasibility of our pipeline, we conducted small-scale experiments on two open-source Java projects ( $\approx 50$ – $100$ k LOC each). Static analysis produced call graphs with 5–10k edges, which were then processed by our GNN-based partitioning model.

- **Service boundary detection accuracy:** On synthetic benchmarks (monoliths reconstructed from microservice systems), our approach recovered 70–80% of original service boundaries (F1-score), outperforming package-based heuristics ( $\sim 45\%$ ) and co-change clustering ( $\sim 60\%$ ).
- **Effort reduction:** In a pilot evaluation with two developers, the AI-driven suggestions reduced manual analysis time by  $\sim 35$ – $40\%$  compared to unguided refactoring.
- **Explainability:** Participants noted that annotated rationales—such as explanations of cohesion and coupling tradeoffs—helped them trust the recommendations. However, they also emphasized the need for further refinement of these explanations.

While preliminary, these results suggest that the proposed pipeline can yield meaningful accuracy improvements and reduce human effort in practice. Larger-scale experiments and industrial case studies are needed to validate generalizability.

### 9. TOOLING & INTEGRATION

A practical toolchain includes:

- **Static analysis engines:** language specific parsers (JavaParser, Roslyn, libclang), call graph builders (WALA, Soot for Java).
- **Data pipeline:** ETL to build graphs and metadata.
- **ML infrastructure:** training/evaluation (PyTorch/TF), GNN frameworks (PyTorch Geometric, DGL).
- **Frontend:** visualization (graph explorers showing proposed boundaries), IDE plug-ins (VS Code, IntelliJ) for inline suggestions.
- **CI integration:** use in continuous refactor checks, regression testing of migration steps.

Privacy & IP: run analysis locally in enterprise environments; avoid sending source to external servers unless explicitly authorized

### 10. LIMITATIONS & THREATS TO VALIDITY

- **Static analysis incompleteness:** reflection, dynamic class loading, and interpreted languages may undermine graph completeness.
- **Ground truth scarcity:** labeled datasets of monolith $\rightarrow$ microservice migrations are rare, limiting supervised learning.
- **Context dependency:** organizational, regulatory, and non-technical constraints may override any automated suggestion.
- **Model bias:** models trained on open-source projects may not generalize to proprietary enterprise systems.
- **Evaluation difficulty:** matching proposed partitions to “correct” ones is inherently fuzzy — different valid decompositions may exist.

# IJETRM

## International Journal of Engineering Technology Research & Management (IJETRM)

<https://ijetrm.com/>

We recommend human-in-the-loop workflows and interactive refinements to mitigate these limitations.

### 11. REPRODUCIBILITY & EXPERIMENTAL PLAN

To validate the approach:

1. **Dataset collection**
  - Gather open-source projects with documented migrations or modularization histories.
  - Select monoliths across languages (Java, Python, C#).
2. **Baseline implementations**
  - Implement Louvain and spectral clustering baselines on call/co-change graphs.
  - Implement pairwise classifier and a GNN.
3. **Training & validation**
  - Use k-fold cross validation, with careful temporal split to simulate real migrations (train on earlier commits, test on later refactorings).
  - Report cohesion/coupling, migration cost estimates, and developer study outcomes.
4. **Developer evaluation**
  - Conduct controlled user studies where engineers use the tool to create a migration plan; measure time saved and subjective trust.
5. **Open artifacts**
  - Release anonymized datasets, code, and scripts to support reproduction (subject to license constraints).

### 12. ETHICAL CONSIDERATIONS

- **Job impacts:** automation should augment engineers, not replace domain expertise.
- **IP & privacy:** ensure code analysis respects confidentiality; provide on-premise deployment.
- **Explainability:** provide clear rationales for suggestions to support accountability.
- **Security:** generated migration plans must consider security implications (e.g., avoiding accidental exposure of sensitive data).

### 13. FUTURE WORK

We identify both **short-term** and **long-term** directions to extend this work:

#### Short-term priorities:

- **Runtime-aware refinement:** fuse production telemetry (latency, error rates) for more accurate boundary decisions.
- **Interactive tools:** allow developers to provide constraints and preferences (must-link / cannot-link) and immediately see updated partitions (active learning).
- **Transfer learning across organizations:** apply domain adaptation methods to adapt models trained on public code to enterprise systems.

#### Long-term priorities:

- **Automated code transformation:** beyond suggestions, generate mechanical refactorings (API adapters, dependency inversions) with correctness guarantees.
- **Contract evolution & schema migration:** integrate database migration planners and automatic schema evolution helpers.

### 14. CONCLUSION

Refactoring monoliths into microservices is a high-value but high-risk engineering task. By combining deep static analysis with modern machine learning—especially graph-based models and semantic embeddings—we can produce high-quality, ranked, and actionable suggestions for service boundaries. The proposed pipeline balances automated inference with human oversight, emphasizing explainability and practical migration plans. While open challenges remain (data scarcity, dynamic behavior, organizational constraints), the outlined approach provides a clear path toward reducing cost and risk in large-scale architectural refactorings.

### REFERENCES

- [1] Fowler, M. *Refactoring: Improving the Design of Existing Code*. Addison-Wesley.
- [2] Richardson, C. *Microservices Patterns*. Manning Publications.

# IJETRM

## International Journal of Engineering Technology Research & Management (IJETRM)

<https://ijetrm.com/>

- [3] Allamanis, M., Barr, E. T., Devanbu, P., Sutton, C. *A Survey of Machine Learning for Source Code*. (Survey papers and ACM surveys on ML for code provide comprehensive background.)
- [4] S. Suryanarayana et al., *Refactoring for Software Architecture Smells*, etc. (work on architectural smells).
- [5] Papers on Graph Neural Networks for program analysis (message passing nets applied to ASTs and call graphs).
- [6] Research on co-change analysis and mining software repositories (e.g., Zimmermann et al.).
- [7] R. Kazman, et al., "A Case Study in Locating the Architectural Roots of Requirements," *ICSE*, 1994.
- [8] M. D. Syer, et al., "Mining the Evolution of Service-Oriented Systems," *MSR*, 2011.
- [9] R. E. Lopez-Herrejon, et al., "Migrating from Monolith to Microservices," *SPLC*, 2018.
- [10] M. Allamanis, et al., "A Survey of Machine Learning for Big Code," *ACM Computing Surveys*, 2018.
- [11] N. Dragoni, et al., "Microservices: Yesterday, Today, and Tomorrow," *Present and Ulterior Software Engineering*, 2017.