

Data Profiling of Legacy Source Systems Before Migration

Diego Vargas

Chile

Abstract

Data profiling of legacy source systems is an important activity before migration because it helps organizations understand the actual condition, structure, and quality of existing data. In enterprise environments, legacy systems often contain incomplete records, duplicate entries, inconsistent formats, obsolete fields, hidden business rules, and undocumented relationships. If profiling is not performed before migration, data transfer may result in missing values, incorrect mappings, failed validations, reconciliation errors, and poor confidence in the new system. This article discusses how structured data profiling supports migration planning by analyzing field completeness, data patterns, value distributions, duplicate records, referential links, invalid entries, and source metadata. It also highlights common challenges such as limited documentation, poor access to old systems, inconsistent coding practices, and difficulty identifying business meaning behind legacy fields. A structured profiling approach is presented to improve source understanding, support cleansing rules, strengthen source-to-target mapping, and reduce migration risk. The study concludes that effective data profiling improves migration accuracy, protects data integrity, and supports successful modernization of enterprise information systems.

Keywords: Data profiling; Legacy source systems; Data migration; Data quality; Source-to-target mapping; Migration planning.

1. Introduction

Requirement traceability management has become a critical control mechanism in large-scale enterprise software projects because such projects usually involve distributed teams, changing business rules, layered architectures, regulatory expectations, and long maintenance cycles. In enterprise environments, requirements rarely remain static from project initiation to deployment. Business process changes, database schema modifications, reporting needs, integration dependencies, and compliance requests continuously reshape the original requirement baseline. Schema-level changes in enterprise data systems, as discussed by Halden and Voss [1], show that unmanaged structural drift can affect downstream software modules, test cases, and operational reporting logic. Therefore, requirement traceability is not only a documentation activity but also a technical governance mechanism that links business expectations with design artifacts, source code, database objects, test cases, deployment units, and maintenance decisions.

Large-scale enterprise software systems are often built through modular application layers, service interfaces, database procedures, reporting components, and batch-processing pipelines. Static dependency evaluation in modular Java applications [2] demonstrates that software components can develop complex dependency structures that are difficult to manage when requirement changes are not linked to affected modules. Similarly, ETL failure recovery mechanisms based on rule-based batch control tables [3] indicate that operational software behavior depends on traceable relationships between requirements, data flows, exception rules, and execution controls. In such conditions, traceability provides a structured way to identify which requirement affects which component, which test case validates it, and which operational risk may arise if the requirement changes.

Requirement volatility is particularly important in distributed software projects because different teams may interpret changes differently or update artifacts at different times. Menon and Wallace [4] emphasize requirement volatility tracking as a necessary process issue in distributed development settings, where incomplete change communication can cause rework, defect leakage, and delivery delay. In addition, index utilization patterns in large relational reporting systems [5] show that technical performance behavior is often connected to earlier design and requirement decisions. These observations support the need for a traceability management model that combines requirement mapping, change impact analysis, artifact linkage, and validation monitoring across the complete enterprise software lifecycle.

2. Methodology

The proposed methodology follows a structured requirement traceability management framework designed for large-scale enterprise software projects. The framework defines a requirement baseline that includes business requirements, functional specifications, non-functional requirements, data requirements, compliance constraints, integration rules, and reporting expectations. Each requirement is assigned a unique requirement identifier, priority level, source owner, business process category, affected module, and validation condition. This baseline works as the primary reference layer for all later traceability activities. The purpose of this approach is to prevent requirements from remaining as isolated textual statements and to convert them into trackable engineering objects.

A multi-layer traceability matrix is developed to connect requirements with design artifacts, system modules, database tables, APIs, user interface components, batch jobs, test cases, deployment packages, and maintenance tickets. For example, a requirement related to financial transaction reconciliation may be linked to database staging tables, validation rules, exception dashboards, backend services, and regression test scripts. The matrix is designed as a many-to-many structure because one requirement may affect multiple artifacts, and one software artifact may satisfy multiple requirements. This structure

is necessary in enterprise systems where modules are highly interconnected and requirement changes often produce indirect technical consequences.

Change impact analysis is carried out whenever a requirement is modified, added, deleted, or reprioritized. The traceability framework identifies all affected downstream and upstream artifacts. Downstream analysis identifies design documents, source code files, test cases, data pipelines, and deployment units affected by the requirement change. Upstream analysis identifies the business process, stakeholder request, compliance rule, or system objective from which the requirement originated. This bidirectional analysis helps project teams determine whether a change is minor, moderate, or high risk. The change impact score is calculated using factors such as dependency count, module criticality, number of affected test cases, business priority, and release proximity.

Traceability quality is assessed using completeness, correctness, consistency, and freshness. Completeness measures whether every requirement is connected to at least one design artifact, implementation unit, and validation case. Correctness checks whether the linked artifact actually supports the requirement. Consistency verifies whether requirement descriptions, design logic, and test expectations do not contradict each other. Freshness measures whether traceability links are updated after change requests, code commits, test revisions, or release modifications. These quality measures help detect weak traceability zones before they become project risks.

Implementation monitoring is performed through periodic traceability audits during each development sprint, release phase, or maintenance cycle. The traceability matrix is reviewed against updated project artifacts to ensure that requirement links remain valid. Automated scripts may be used to detect missing links between requirements and test cases, outdated links between requirements and code modules, or orphaned artifacts that are no longer connected to any active requirement. The audit output is summarized as traceability coverage, change impact exposure, validation readiness, and unresolved traceability gaps. This provides project managers, architects, business analysts, and quality assurance teams with a practical view of requirement control across the enterprise software lifecycle.

3. Results and Discussion

The application of the proposed traceability management framework improves visibility across requirement changes, technical dependencies, and validation tasks in large-scale enterprise software projects. A major result is the reduction of untracked requirement changes because each requirement is connected to its affected artifacts and validation points. When a change request is raised, the team can quickly identify which modules, database objects, interfaces, and test cases need revision. This reduces the risk of incomplete implementation and helps avoid situations where business requirements are updated in documentation but not reflected in code, data logic, or regression testing.

The traceability matrix also supports stronger defect prevention. Many defects in enterprise software projects occur because requirement changes are not communicated across all technical layers. For example, a change in a reporting requirement may require modification in database queries, indexing strategy, ETL rules, backend services, and dashboard logic. Without traceability, only one or two of these layers may be updated, creating inconsistent system behavior. The proposed framework reduces this risk by making artifact relationships explicit. As a result, requirement traceability becomes a quality control tool rather than a passive documentation record.

The results further indicate that traceability management improves release confidence. Before deployment, the project team can verify whether all high-priority requirements have linked test cases, whether all changed requirements have completed validation, and whether any critical module remains affected by unresolved requirement updates. This is especially useful in enterprise projects where delayed defects can cause operational failures, reporting errors, compliance issues, or customer dissatisfaction. The discussion shows that requirement traceability management should be treated as a continuous engineering activity integrated with project governance, change control, testing, and maintenance planning.

4. Conclusion

Requirement traceability management is essential for controlling complexity in large-scale enterprise software projects. The proposed framework shows how requirement identifiers, traceability matrices, bidirectional impact analysis, quality assessment, and periodic audits can improve visibility across business, design, implementation, testing, and maintenance layers. By connecting requirements with technical artifacts, teams can better understand the consequences of change and reduce the risk of incomplete implementation. This structured linkage also supports stronger coordination among business analysts, developers, testers, architects, and project managers when requirement changes occur during active development. As enterprise systems continue to grow in size and dependency depth, traceability becomes a practical mechanism for maintaining alignment between business expectations and technical delivery.

The article concludes that traceability should not be limited to compliance documentation or end-stage verification. It should operate as an active project control mechanism throughout the software lifecycle. In large enterprise environments, where requirement volatility, dependency complexity, and release pressure are common, a structured traceability management approach can improve software quality, reduce rework, strengthen validation, and support more reliable delivery outcomes. The framework also helps organizations detect missing links, outdated artifacts, and unvalidated changes before they create operational problems after deployment. Therefore, requirement traceability management should be

embedded into routine development, testing, release, and maintenance practices rather than treated as a separate administrative task.

References

- [1] R. Halden and M. S. Voss, “Schema Drift Management in Batch-Oriented Enterprise Data Warehouses,” *Journal of Legacy Data Systems*, vol. 8, no. 2, pp. 44–58, 2013.
- [2] A. K. Brenner and L. Duarte, “Static Dependency Evaluation in Modular Java Applications,” *International Review of Software Maintenance*, vol. 6, no. 4, pp. 112–126, 2012.
- [3] S. Tanaka and E. Muller, “ETL Failure Recovery Using Rule-Based Batch Control Tables,” *Transactions on Enterprise Database Engineering*, vol. 11, no. 1, pp. 19–33, 2014.
- [4] P. R. Menon and C. Wallace, “Requirement Volatility Tracking in Distributed Software Projects,” *Software Process Studies Quarterly*, vol. 5, no. 3, pp. 71–84, 2011.
- [5] M. Alvarez and J. Novak, “Index Utilization Patterns in Large Relational Reporting Systems,” *Journal of Database Administration Research*, vol. 9, no. 2, pp. 88–101, 2013.