

COBOL to Java and Newspapers Still Get Delivered

Alessandro De Marco
The New York Times Company
New York, NY, USA
alex.demarco@nytimes.com

Valentin Iancu
Modern Systems International Ltd.
Bucharest, Romania
viancu@modernsystems.com

Ira Asinofsky
The New York Times Company (Retired)
New York, NY, USA
irabevasi@aol.com

Abstract—This paper is an experience report on migrating an American newspaper company’s business-critical IBM mainframe application to Linux servers by automatically translating the application’s source code from COBOL to Java and converting the mainframe data store from VSAM KSDS files to an Oracle relational database. The mainframe application had supported daily home delivery of the newspaper since 1979. It was in need of modernization in order to increase interoperability and enable future convergence with newer enterprise systems as well as to reduce operating costs. Testing the modernized application proved to be the most vexing area of work. This paper explains the process that was employed to test functional equivalence between the legacy and modernized applications, the main testing challenges, and lessons learned after having operated and maintained the modernized application in production over the last eight months. The goal of delivering a functionally equivalent system was achieved, but problems remained to be solved related to new feature development, business domain knowledge transfer, and recruiting new software engineers to work on the modernized application.

Index Terms—Software testing, Mainframe, COBOL, Java, Software application migration, Code translation

I. INTRODUCTION

Since 1979, this American news media company relied on a software application running on its mainframe as the core IT system supporting daily home delivery of its newspaper. The application had grown to more than two million lines of COBOL code implementing billing, customer account maintenance, delivery routing, and other business-critical functionality. As a legacy system, it “represent[ed] years of accumulated experience and knowledge” [1], while it “significantly resist[ed] modification and evolution” [2]. It was also very expensive to operate in comparison to more modern systems at the company.

An attempt to redevelop the home delivery application between 2006 and 2009 failed. In 2015, with mounting pressure to quickly lower costs, a different modernization approach was selected. Instead of redeveloping the application from scratch, the aim was to migrate the application off the mainframe and onto Linux servers by using a code and data translation approach. This approach promised to deliver an application that would be functionally equivalent, cheaper to operate, and easier to integrate with in comparison to the original. An evaluation of alternate approaches determined that a second attempt at redeveloping the application would have been much more expensive, and rehosting [3] would have continued to lock-up data in proprietary technology.

A vendor provided the technology to convert the code and data [4]. Based on an early *proof-of-concept* trial of the

vendor’s technology, it became clear that even though the translated software could work as expected it would likely be more costly to maintain and enhance than *handcrafted* Java software. It would also require knowledge of COBOL programming idioms and mainframe concepts that most Java software engineers would not possess. Despite the disadvantages, at an estimated cost of less than a tenth of the 2006-2009 redevelopment initiative and with a projected timeline of just one year, senior IT management opted to move forward.

As the project team leaders, we report on what went well, what did not, and lessons learned. Even though the modernized application went live a year later than originally planned, this is a success story.

There are known challenges inherent in legacy code translation [5]. However, in our project, testing the application proved to be the most time-consuming, difficult, and underestimated area of work. Unexpected testing obstacles caused significant delays. Development of an elaborate testing process was required in order to test functional equivalence between the legacy application and the modernized application, as the terms are defined in [6], while supporting some functionality changes and feature enhancements along the way.

The main contributions of this paper are:

- 1) Our process to test functional equivalence between the legacy and modernized applications.
- 2) The obstacles that we encountered while testing the modernized application.
- 3) Issues encountered in production due to gaps in the testing process and other lessons learned.

In II, we describe our migration methodology. In III, we explain our testing process. In IV, we cover the main challenges encountered in the testing process, production issues, and lessons learned. We conclude in V.

II. MIGRATION METHODOLOGY

A. An 8-Step Process

We employed an 8-Step process to migrate our application off the mainframe and onto Linux servers. The steps are listed below. For a detailed explanation of the steps, refer to [4].

- Step 1: Collect Inventory in the Legacy System
- Step 2: Break down into Work Packets
- Step 3: Database Remodelling
- Step 4: Data Migration
- Step 5: Code / JCL Conversion
- Step 6: Testing
- Step 7: User Acceptance Testing
- Step 8: Cutover

Work Packets are defined as one or more component groups of the application that could be translated and tested together. For each Work Packet, we executed Steps 3-6. Once all Work Packets had been converted and tested, we executed Steps 7 and 8.

B. Technology Conversion Mapping Summary

A high-level technology conversion mapping is provided in the table below. A detailed explanation of the vendor’s translation technology (i.e. first 4 rows) is available in [4].

Technology	Legacy	Modernized
Programming Language	COBOL	Java
Database	VSAM KSDS files	Relational Database (Oracle)
Batch Jobs	JCL	Spring Batch XML (JSR-352)
User Interface Screens	BMS Maps	JSF, HTML/CSS, Javascript; accessible via a Web Browser
Security and Access Control	RACF	Spring Security and Active Directory
Middleware	CICS	Jetty, Apache-CXF, JPA/Eclipselink
Reporting	QMF and DB2	Jasper Reports and Oracle
Encryption	Megacryption	Java Cryptography Extension
Screen Automation / Macros	IBM HATS	Newly developed Web Services
Batch Process Scheduler	CA7	Control-M
Monitoring and Alerting	RMF, SMF, Omegamon	NewRelic, nagios, Sumologic
Development and Deployment	Changeman	git, gradle, jenkins, puppet, ansible

C. Framework Support For Utilities and Middleware Services

In addition to translating the code and data, the vendor provided a runtime framework that implemented many mainframe services and utilities. This allowed the translated code to run on the modernized platform while continuing to interface with its environment in a similar way to how it did on the mainframe. This approach was also used in [7].

D. New Component Development

When legacy application dependencies were not supported by the runtime framework (e.g. REXX, GVEXPORT), if an off-the-shelf software package was not available as a substitute, or if it made more sense to make use of capabilities of the new environment (e.g. Oracle database backups and restore points; file system snapshots), then replacement components were developed. In addition, since automatically converting the CA7 batch schedule for Control-M was unsuccessful, we redeveloped the batch job schedule for Control-M.

III. TESTING PROCESS

To assure that the modernized application would be *functionally equivalent* to the legacy application, we developed the testing process described next. By functional equivalence we mean that the modernized application would produce the same output as the legacy application given the same input. As the translation from COBOL to Java preserved business rules and much of the internal component hierarchy of the legacy

application, we were able to test functional equivalence at the component group level first, and gradually build up to testing functional equivalence of the whole application.

A. Black-Box Testing of Component Groups

Component groups assembled related components that would be runnable and testable together as a “black-box” through existing externally accessible interfaces, such as SOAP Web Services, User Interface screens, database tables, and files. Given the same inputs, the output of legacy component groups were compared to the output of their modernized counterparts. When they matched, the test was considered to have passed. When they did not, root-cause analysis was performed to find the source of the mismatch.

B. Test Environments

A new test region was created on the mainframe to support the testing process. It served as the *source of truth* for expected behavior.

QA analysts and developers installed a Linux virtual machine and a database on their individual computers. This allowed everyone to test modernized components locally.

Development, Staging, and Production environments were provisioned in a private datacenter. A shared Oracle database and a Control-M batch scheduler were configured in each. These environments let us build and test infrastructure and configuration automation code, the full batch schedule, and infrastructure-related performance improvements, since these could not be tested on individual developer machines. Static test data was used in the Development environment, and dynamic *current day* test data was used in the Staging environment. To generate current day test data, the batch process ran every day in the mainframe test region and output files were then transferred over. This enabled testing of scenarios that depended on the day of week or month, and it validated that one day’s batch output would be processed correctly the next day.

C. Stages In the Testing Process

The testing process, Step 6 of II-A, was broken down into stages, with each stage progressively increasing in scope and level of difficulty in isolating the root-cause of test failures. The stages are summarized in the table below.

Stage 1: Pre-Delivery	
D:	Tests done by the translation technology vendor prior to delivering translated code.
E:	Individual developer machines.
S:	One batch job, one set of screens, or a Web Service; automated testing of 10 batch jobs for regression testing.
Stage 2: Data Migration Validation	
D:	Tests done by the QA analysts to make sure that data loaded into the database matched VSAM KSDS files.
E:	Individual QA analyst machines, and Development.
S:	All migrated data.

Stage 3: Component Group	
D:	Tests done by QA analysts to make sure that a component group worked as expected.
E:	Individual QA analyst machines, and, for Web Services and UI screens, both Development and Staging.
S:	One batch job, one or more related UI screens, one SOAP Web Service.

Stage 4: Batch Process (Operating on Static Test Data)	
D:	Automated tests for the end to end batch process. Also acted as a batch regression test system.
E:	Development
S:	The full batch, but configured to run the same day, everyday in order to simplify root-cause analysis.

Stage 5: Batch Process (Operating on Dynamic Test Data)	
D:	Automated tests for the end to end batch process.
E:	Mainframe test region and Staging, running in parallel.
S:	The full batch, but configured to compare current day output between the legacy and modernized processes.

Stage 6: Batch Process (Operating on Full Production Data)	
D:	Automated execution of the batch process in production to benchmark performance.
E:	Production (prior to cutover).
S:	The full batch, operating on production data.

Stage 7: System Integration	
D:	Tests done by QA analysts in collaboration with other teams/systems within the organization.
E:	Staging
S:	The whole enterprise system, with new transactions entered via client systems, processed by the batch, and flowing to downstream consumers via reports and file feeds.

D = Description, E = Environment, S = Scope

IV. DISCUSSION

We discuss the main obstacles encountered in the Testing Process defined in III, types of production issues that were not caught in testing, and lessons learned.

A. Testing Obstacles

Testing accounted for approximately 70% to 80% of the time spent on the project. The project was completed a year later than expected due primarily to these obstacles. We have grouped the obstacles into three areas.

1) Initial State of the Legacy Application:

a) *Lack of Tests:* The vast majority of legacy application components did not have a high enough level of test coverage to codify the required information for functional equivalence testing. As a consequence, much time was spent trying to analyze inputs, outputs, and sometimes even the internal behavior of components when tests failed and root-cause analysis had to be performed. This was especially time-consuming and difficult for the more complex batch jobs (i.e. many man-months of work).

b) *Lack of Batch Automation:* The CA7 batch process scheduler had never been installed in the existing mainframe test region. In order to test functional equivalence of the batch process end-to-end, Control-M was installed and configured in

a new test region as noted in III-B. This work was unplanned and required a great amount of effort.

c) *Obsolete Code:* In operation for more than 35 years, the legacy application had accumulated a fair amount of obsolete code [8]. Due to a lack of adequate maintenance over the years, we spent time identifying this code and removing it to reduce the amount of translation and testing work to do.

d) *Interfaces with Other Systems:* The legacy application generated more than 3500 data file feeds and reports for downstream consumers daily. We did not know all the consumers, many transfers used insecure protocols, and connections were configured in a variety of different ways. Discovering the consumers, upgrading to secure protocols, and harmonizing configuration management caused significant delays.

2) Data Formats:

a) *EBCDIC Files That Contained Computational Fields:* Many files on the mainframe contained a mix of display fields and computation fields. The mainframe environment has appropriate tools for working with these file formats, but the Linux platform does not. Files of this kind had to be transferred over because they were required for functional equivalence testing of batch jobs. When transferring these files off the mainframe, the display fields needed to be converted from EBCDIC to ASCII, but the computational fields (e.g. COMP3) had to remain binary-encoded. ETL software was developed for this which required additional development work, became a processing bottleneck, and was a source of errors. Handling the large volume of files and many versions of each further complicated the matter.

b) *File Processing Tools On Linux:* A lack of tools to inspect, modify, and compare files with fixed block or variable block structures and computational fields on the Linux platform necessitated the development of new tools. This work was not part of the original plan.

3) Batch Performance:

a) *Multi-layout VSAM KSDS Files:* VSAM KSDS files with multiple record layouts (i.e. REDEFINES) could be translated to one or many tables. We obtained better performance when mapping to one table, but improved maintainability (i.e. fewer null-value columns) when mapping each layout to its own table. Since reading records in an indexed VSAM file often involved moving a pointer to an indexed location given a key, and then iteratively moving the pointer forward or backward to the next or previous record until some condition was met, when translated to one relational database table using a result set cursor as the pointer, this usually performed well. When translated to 97 tables, the worst case we encountered, this was painfully slow because open cursors had to be maintained on result sets from each of the 97 tables, and then data access logic decided which cursor would have the next or previous record.

b) *In-Memory VSAM Cache:* When testing some batch jobs, we encountered performance problems that could not be solved by database remodelling as described previously. The technology vendor developed an in-memory representation of VSAM KSDS files that acted like a cache. Encapsulated within

data access middleware and tuned with external configuration settings, no application code needed to be changed. Once all necessary data had been loaded into memory, operations were performed against the in-memory data structure, and at the end of processing, changes were written back to the database. By performing these operations in-memory, network latency and database I/O bottlenecks were eliminated, but we could not run other batch jobs at the same time if they depended on the same database tables.

B. Types of Issues That The Testing Process Failed to Uncover

Despite the heavy investment in testing, there were gaps in the testing process. After completing User Acceptance Testing and Cutover steps of II-A, we encountered problems in production due to the gaps. We have grouped them by type with examples.

a) *Unexpected User Input:* On the first day of operation, a subscriber contacted our call center to complain about non-delivery of the newspaper over the prior 147 days. Complaints over such a long time period are extremely uncommon. The algorithm to compute the credit due the subscriber involved a date computation that overflowed a field that supported a maximum of 99 days. This overflow caused a runtime exception and a critical batch job failed. When attempting to reproduce this on the mainframe, we observed that the date value was simply truncated causing the calculation to be incorrect, but the job did not fail.

b) *Concurrency Issues:* On another day, Control-M unfortunately scheduled two batch jobs to run at the same time. They both wrote to the same output file which corrupted the data in that file. The corrupted data was later loaded into an in-memory VSAM cache (see IV-A3b) which caused a critical batch job to fail.

c) *Inefficient Processing Idioms:* Several batch jobs involved maintenance of data. On the mainframe, data maintenance was often done by transferring VSAM data to files which would then be pruned and sorted using file processing tools. The cleaned up files would then be reloaded into VSAM and reindexed. The translation of these batch jobs resulted in operations (i.e. millions of DELETES and INSERTs) that locked up our database for hours on two different occasions causing online transaction processing outages. These maintenance jobs were redeveloped to make use of more efficient in-database maintenance operations or decommissioned altogether.

C. Lessons Learned

a) *Application Understanding:* As noted in a 2013 Gartner Survey cited in [8], most modernization projects are delivered later than originally planned “as a direct result of poor legacy application understanding”. The translation approach did not eliminate the need for us to develop a deep level of application understanding in critical parts of the system. We could have developed this knowledge while improving the initial state of the application as a precursor project to this one. This would likely have alleviated the problems noted in

IV-A1 and informed planning, estimation, and management of the project to follow.

b) *Project Management:* We attempted to break work down into smaller tasks and measure time taken to complete the tasks in order to forecast when similar future work would be completed, but this did not work well. Later Work Packets were more batch-oriented than earlier ones, so testing obstacles related to batch processing (IV-A2 and IV-A3) caused early delivery forecasts to be way off. This problem was discovered too late to remedy, and may have been avoided had we had a deeper understanding of the application when structuring the Work Packets at the outset.

V. CONCLUSION

We employed a code and data translation approach and an 8-step methodology to migrate a legacy application off the mainframe to Linux servers. We achieved the goal of delivering a modernized application that is functionally equivalent to the original. Functional equivalence was demonstrated via an elaborate testing process. It was delivered later than expected due to unexpected testing obstacles, but still at much lower cost than a prior redevelopment project that failed. Despite a few production issues, the modernized application has proven to be quite stable in production over the last eight months and newspapers still get delivered daily.

On the other hand, new feature development remains challenging. It requires knowledge of COBOL programming idioms and mainframe concepts that Java software engineers on the team do not possess. Mainframe COBOL developers on the team know the idioms and concepts, but are not yet proficient in Java. We are cross-training developers and hope that this will help with knowledge transfer. We also note that it has been difficult to recruit new Java developers due to a lack of interest in working with the translated code.

REFERENCES

- [1] K. Bennett, “Legacy systems: Coping with success,” *IEEE software*, vol. 12, no. 1, pp. 19–23, 1995.
- [2] M. L. Brodie and M. Stonebraker, *Legacy Information Systems Migration: Gateways, Interfaces, and the Incremental Approach*. Morgan Kaufmann Publishers Inc., 1995.
- [3] Micro Focus, “Micro Focus Enterprise Server Overview.” [Online]. Available: <https://www.microfocus.com/products/enterprise-suite/enterprise-server/>
- [4] Modern Systems International Ltd., formerly BluePhoenix Solutions, “Mainframe Migration COBOL to Java,” 2014. [Online]. Available: <https://www.platformmodernization.org/bluephoenix/Lists/ResearchPapers/Attachments/2>
- [5] A. A. Terekhov and C. Verhoef, “The realities of language conversions,” *IEEE Software*, vol. 17, no. 6, pp. 111–124, 2000.
- [6] R. Khadka, P. Shrestha, B. Klein, A. Saeidi, J. Hage, S. Jansen, E. van Dis, and M. Bruntink, “Does software modernization deliver what it aimed for? a post modernization analysis of five software modernization case studies,” in *Software Maintenance and Evolution (ICSME), 2015 IEEE International Conference on*. IEEE, 2015, pp. 477–486.
- [7] H. M. Sneed and K. Erdoes, “Migrating AS400-COBOL to Java: a report from the field,” in *Software Maintenance and Reengineering (CSMR), 2013 17th European Conference on*. IEEE, 2013, pp. 231–240.
- [8] Modern Systems International Ltd., “Three Essentials for Reducing Risk When Migrating from COBOL to Java,” 2017. [Online]. Available: <http://modernsystems.com/essentials-cobol-to-java/>