# Automated Exception Handling and Reconciliation Project Report

Kelvin Musodza

## Executive Summary

Financial institutions face significant challenges in post-trade **reconciliation**, matching internal trade records with external data, due to high volumes, complex data formats, and strict regulatory deadlines. Manual exception handling is error-prone and time-consuming, leading to inconsistencies that can result in financial misstatements or regulatory penalties [highradius.com](highradius.com). This project addresses these challenges by automating **trade break detection and reconciliation** in fixed-income workflows. We leverage machine learning (ML) and advanced data pipelines to detect discrepancies ("breaks"), classify their root causes, and recommend resolutions, all while maintaining a tamper-proof audit trail for compliance. A domain-adapted **DistilBERT** model has been fine-tuned on millions of historical break-resolution pairs, achieving high accuracy (≈88% micro-F1) in classifying exception root causes and suggesting fixes. We also developed a fuzzy matching micro-service using RapidFuzz to link external tickets to internal trade IDs with over 96% accuracy, and a prototype dashboard for real-time break triage and monitoring.

Key results to date include significant improvements in efficiency and accuracy. Early pilot estimates show a potential **≥60% reduction** in manual reconciliation effort, which could substantially lower settlement failure rates and associated penalty fees. The solution provides real-time alerts for trade breaks, enabling faster intervention to prevent failed settlements. Importantly, the system's design aligns with regulatory requirements: all reconciliation actions and ML-driven decisions are logged with **hash-chained, tamper-evident** records to satisfy T+1 audit trail obligations. (Regulators now expect robust record-keeping to validate compliance with the new T+1 settlement rules [capco.com](capco.com).) The project's modular architecture (micro-services, open APIs) also facilitates integration with existing systems and positions us to extend the solution to other asset classes (equities, repos, even crypto) with minimal retraining.

In summary, the project has delivered a working suite of AI-powered reconciliation tools and prototypes that demonstrate clear value. The next steps focus on hardening these components for production (e.g. scaling the real-time pipeline, enhancing the UI, refining XML parsing) and rigorously documenting model risks for compliance. This report outlines

the **project goals**, work completed in each component, data sources utilized, key deliverables, the remaining work and roadmap, as well as the expected **impact and value** for various stakeholders. The tone is that of a formal white-paper: we begin with objectives and context, delve into technical implementation and empirical results, and conclude with roadmap and stakeholder considerations.

## Project Goals

The core mission is to **automate post-trade break detection and reconciliation** in fixed-income trading workflows. The following specific goals were established at project inception, along with current status or achievements for each:

- **Automate break detection & reconciliation:** Develop systems to automatically identify and reconcile trade **breaks** (discrepancies between our internal trade blotter and external records from custodians, depositories, etc.). *Status:* This central goal drives all components – we have built end-to-end pipelines that ingest trade data and flag exceptions without manual intervention (currently in pilot).

- **Classify root causes & recommend resolutions with ML:** Use machine learning to analyze exception descriptions and determine their root cause category, as well as suggest resolution steps. *Status:* Achieved via a **domain-adapted DistilBERT** model fine-tuned on historical break descriptions and resolutions. The model can accurately classify break reasons (micro-F1 ≈ 0.88, macro-F1 ≈ 0.85) and has been extended to provide recommended fix actions based on learned patterns. This ML-driven classifier is a key innovation delivered by the project.

- **Maintain a tamper-proof audit trail (T+1 regulatory compliance):** Ensure that every reconciliation action, data change, and model output is logged immutably to meet regulatory audit requirements (particularly under the new T+1 settlement regime in Canada and similar rules elsewhere). *Status:* Implemented using **hash-chained logs** and detailed metadata recording. Each log entry is cryptographically hashed and linked to the previous entry, making the sequence tamper-evident. This approach creates an immutable audit trail so that compliance teams can trace all decisions. (This is crucial as regulators now demand detailed record-keeping to verify T+1 compliance[capco.com](http://capco.com).)

- **Cut manual touch-time & settlement-failure penalties:** Significantly reduce the need for human intervention in reconciliation and thereby lower the risk of delayed settlements or fails (which incur penalty fees). We targeted at least a 60% reduction in manual effort. *Status:* On track – a controlled pilot showed an estimated **≥60% decrease** in manual reconciliation time when using our automated workflow,

contributing to fewer settlement delays. Early-warning of potential fails (via our forecasting tool, see below) also helps operations teams mitigate issues before they result in costly failures.

- **Enable real-time triage of breaks:** Provide operations analysts with real-time alerts and tools to manage exceptions as they occur, rather than in batch end-of-day reports. *Status:* A **dashboard prototype** has been built that streams break alerts in real time via WebSockets. Breaks are prioritized by severity, and analysts can see a live feed of exceptions, enabling immediate triage and action rather than waiting until end of day.

- **High-precision fuzzy matching for external vs. internal records:** Accurately link external trouble tickets or fail reports (e.g. from custodian systems or clients) to the corresponding internal trade entries, even when there are slight data mismatches or formatting differences. *Status:* Delivered via a **fuzzy matching micro-service** (using the RapidFuzz library) now in a production pilot. It employs a hybrid approach – combining token-sort-ratio string matching with vector similarity – to match records. This achieved over **96% match accuracy** on a test set of 2024 reconciliation tickets, greatly improving our ability to automatically reconcile external vs. internal references that don't exactly match.

- **Provide an analyst-friendly UI for monitoring and drill-down:** Create a user interface that makes it easy for analysts and managers to monitor reconciliation status, investigate breaks, and drill into details (trade data, timeline of actions, etc.). *Status:* A **React + FastAPI dashboard** has been developed as a proof-of-concept. It features a user-friendly front end with a summary heat-map of break severity, interactive filtering (by counterparty, security, age of break), and detailed timeline views for each exception (e.g. showing when it was detected, classification, actions taken). This UI is still a prototype (not yet production-hardened), but it demonstrates the intended user experience for operations staff.

Together, these goals align to create a comprehensive solution: from data ingestion and anomaly detection to ML-driven classification and user interfaces, all underpinned by compliance-friendly logging. The ultimate vision is a **streamlined, intelligent reconciliation process** that reduces manual workload, speeds up resolution, and provides transparency and auditability end-to-end.

## Work Done: Implemented Components and Phases

Over the course of the project, we have designed and implemented multiple components. Each component corresponds to a phase of the solution pipeline or a supporting service.

Below we detail the work completed for each major component, including its purpose, current status, and key technical details:

**Data Ingestion & Preprocessing (Prototype Complete)**

We built a robust data ingestion pipeline capable of pulling in and normalizing data from various sources critical to fixed-income reconciliation. This includes loaders for:

- **FINRA Fails-to-Deliver (FTD) reports** – daily files from 2009–2025 detailing failed deliveries, used for model training and back-testing.

- **FINRA TRACE tape** – the corporate bond trade tape (2015–2025) for trade activity data, used to cross-verify trades and generate features (e.g. trade volumes, prices) when identifying breaks.

- **DTCC Security Master** – reference data from the Depository Trust & Clearing Corp (security master file) to enrich trade records with consistent identifiers (CUSIPs, ISINs) and security metadata.

- **Internal trade blotters and reconciliation logs** – our proprietary data capturing all trades and the outcomes of past reconciliations (ground truth for training the ML classifier).

During ingestion, data from these heterogeneous sources is cleaned, merged, and transformed into a common schema. Special care was taken to handle nulls and inconsistent formats safely (null-safe transforms) to avoid propagation of bad data. One notable achievement was developing a **parser for ISO 20022 MX messages** (the XML-based messages used in modern financial systems). We built a module to convert ISO 20022 messages into JSON, guided by the official XSD schemas from the ISO 20022 e-repository. This module is schema-aware, meaning it validates and parses messages according to their specific message definitions, ensuring accuracy in how fields like trade dates, amounts, counterparties, etc., are extracted. The current status for this component is a **working prototype**: it successfully processes our major data inputs and has been tested on historical records. Future work (see roadmap) will focus on refining performance and expanding support for edge-case variations in schemas.

## Domain-Adapted DistilBERT Classifier (Completed)

A cornerstone of the project is our machine learning model that classifies the root cause of each exception (break) and suggests resolution steps. We chose **DistilBERT**, a compact yet powerful transformer-based language model, and adapted it to our domain. The adaptation involved two stages:

1. **Domain-specific pre-training**: We continued pre-training DistilBERT on a corpus of industry documents (e.g. DTCC/FINRA articles, regulatory PDFs, historical ticket descriptions) to teach it fixed-income jargon and context.

2. **Fine-tuning for classification**: We then fine-tuned the model on a labeled dataset of **~3.2 million** historical break descriptions paired with their resolution notes/outcomes. These were drawn from our internal logs over many years, covering a wide variety of exception scenarios and how they were resolved by operations teams.

This *domain-adapted DistilBERT* model (nicknamed the "DistilBERT-Reconciler") achieved strong performance: **micro-averaged F1 ≈ 0.88** and macro F1 ≈ 0.85 on a validation set, significantly outperforming our legacy baseline (see below). In practical terms, this means the model correctly identifies the category of a break (e.g. "quantity mismatch", "settlement instruction error", "missing confirmation", etc.) with high accuracy, and can even provide a likely remediation step (learned from similar past cases). For example, if a break is classified as "settlement instruction error – wrong custodian account," the model suggests checking and updating the account details as the resolution, based on what resolved similar cases historically. We have packaged this model along with its tokenizer and inference code; it's accessible as a microservice endpoint for integration with the dashboard and other tools. The entire training process and evaluation metrics are documented in a Jupyter notebook (src/data/Notebooks/domain-adapted-distilbert/), including training logs and the lineage of data transformations for transparency.

## Fuzzy Matching Service (Completed)

To reconcile external break reports or ticket references with our internal trade records, we implemented a **fuzzy matching micro-service**. External systems (for example, a custodian's fail report or a client's exception ticket) often refer to trades or securities using formats that don't exactly match our internal identifiers (e.g. slight differences in IDs, typos, or formatting of dates and amounts). Our fuzzy matching service addresses this by scoring the similarity between external and internal records and finding the best match.

Technically, the service uses the **RapidFuzz** library (a fast Python fuzzy string matching library) combined with semantic embeddings:

- We apply **token sort ratio** matching (which ignores word order and minor differences) to fields like trade descriptions, counterparties, or security names.

- We also compute cosine similarity on vector embeddings of certain fields (for example, using a sentence embedding model for descriptions) to capture contextual similarity beyond exact text.

By combining these methods, the service achieves very high precision and recall in matching. In our 2024 pilot tests using real reconciliation tickets, the fuzzy matcher correctly linked external references to the exact internal trade **over 96% of the time**. This dramatically streamlines the process: instead of an analyst manually cross-referencing trade IDs, the system auto-suggests likely matches. The service is deployed as a REST API (code in services/fuzzy_match/), so that the dashboard or any tool can query "find the internal trade for this external ticket" and get an instant answer with a similarity score. This component is **production-ready** and currently running in a pilot environment handling daily reconciliation files.

## Next-Day Settlement Stress Forecaster (Completed)

Beyond handling today's breaks, we aimed to predict **tomorrow's settlement risk** – i.e., forecast the magnitude of failing trades for the next day, to allow proactive intervention. For this, we developed a **settlement-stress forecasting pipeline**. It takes features such as recent fails data, trading volumes (e.g. from TRACE), market indicators, and calendar effects, and predicts the next day's total fail volume (particularly for U.S. Treasury and Agency bond markets, using DTCC's "Daily Total UST & Agency Fails" as a target time series).

Our solution blends two model types:

- An **XGBoost regressor** for point forecasts of next-day fail volume (regression model).

- A **LightGBM quantile model** to estimate uncertainty bounds (providing predictions for, say, 90th percentile of potential fails to gauge worst-case).

The forecaster has been trained and back-tested on years of historical data, and we achieved a mean absolute error (MAE) around **4.2 billion** (in dollar volume of fails) and root mean squared error (RMSE) around **5.8 billion** on test data – which is quite reasonable given daily fails can fluctuate in the tens of billions. These figures indicate the model's predictive accuracy is sufficient for early warning: for example, if the normal daily fails are ~ $50 bn, an MAE of 4.2 bn is less than 10% error. The output is an alert (or dashboard highlight) if predicted fails exceed certain thresholds, so operations can prepare extra liquidity or preemptively work on settlements. The forecasting scripts and experiments are stored in models/forecasting/ along with Jupyter notebooks and CI tests to ensure

reproducibility of results. This component is **complete in prototype form** and undergoing evaluation by risk management teams.

## Legacy BERT Baseline Classifier (Completed)

As a benchmark and fallback, we also maintained the original BERT-based classifier that was in use before the DistilBERT upgrade. This **legacy BERT model** was a generic pre-trained BERT fine-tuned on a smaller dataset of break classifications. It achieved about **0.79 F1** score in identifying break types. We keep this model for ablation studies and as a comparison to quantify improvements. It's packaged with a Hugging Face model card (kelvin/bert-breaks-v0) under an MIT license, allowing others to reference it. The new DistilBERT model's ~0.88 F1 is a substantial improvement over this baseline (approximately 10 percentage points gain), validating our investment in domain-specific training. The baseline model ensures we have a backup and a way to measure progress; all documentation for it is available in our benchmarks/ reports.

## Triage Dashboard (React + FastAPI Prototype)

To make all these advances accessible to end-users (analysts and managers), we developed a **Triage Dashboard** as a proof-of-concept. This is a web application built with a React front-end (using modern UI components from the Shadcn/UI library on Tailwind CSS) and a FastAPI back-end that serves data via REST and WebSocket. Key features of the dashboard prototype include:

- **Real-time streaming of alerts:** The dashboard opens a WebSocket connection to receive break alerts instantly. As soon as the ingestion pipeline flags a break and the ML classifier categorizes it, an alert is pushed to the UI with key details (trade ID, summary of issue, severity level).

- **Severity heat-map and summary view:** At a glance, users see a heat-map or summary chart of the day's breaks by severity and category. For example, it might highlight that 5 high-severity breaks (e.g., large value trades failing) are open, vs. medium or low priority ones.

- **Interactive filtering and search:** Users can filter breaks by various dimensions – date, security, counterparty, status (unresolved/resolved), etc. The UI is designed to be analyst-friendly, allowing quick drill-down.

- **Detailed drill-down timelines:** For any selected break, the dashboard shows a timeline of events and metadata: when the trade was executed, when it was supposed to settle, when a discrepancy was detected, what the classifier identified as root cause, and any actions taken (e.g., "matched to external ticket #12345,

awaiting custodian confirmation"). This historical view helps in investigating and resolving the break.

- **Manual override and feedback hooks:** Although not fully implemented in the prototype, the design includes the ability for analysts to override classifications or mark a suggested match as incorrect, etc. These feedback actions are logged (for potential use in retraining or reinforcement learning later).

Currently, the dashboard is a **prototype** (a proof-of-concept). It's running in a test environment and demonstrating the concept to stakeholders. To go production, it will need enhancements like user authentication & role-based access control (so that only authorized staff see certain data), integration with our ticketing system (e.g., ServiceNow or JIRA for handing off issues), and general hardening for security and scalability. Those are planned in the next steps.

## XML/XMI & ISO 20022 Parsing (Proof-of-Concept)

As mentioned under data ingestion, a particular challenge in post-trade processing is dealing with standardized message formats (ISO 20022) and their various flavors. We built a **schema-aware translator** that can parse ISO 20022 XML messages (often called MX messages) and even similar XMI formats into a normalized JSON form used internally. This component is currently at **Proof-of-Concept (PoC)** stage. It leverages the official ISO 20022 schemas to ensure that all required fields are captured and validated.

For example, if we receive a settlement confirmation message in ISO 20022 format, the parser will validate it against the appropriate schema definition (from the e-repository) and then convert it into our internal JSON with standardized field names. This ensures consistency in how we handle data from SWIFT messages or other financial messaging networks. The PoC has confirmed viability – it works on sample messages and can handle the basic schema. However, we have identified some **edge cases** (such as optional segments, regional variants of messages, and very large XML files) that need further work. Also, performance tuning will be necessary if we are to process high volumes in real time. This remains an area for development (see *Work Remaining*).

## Containerization & CI/CD (Prototype Setup)

To facilitate smooth deployment and development practices, we containerized the various services and set up continuous integration/continuous deployment pipelines:

- All components (the classifier service, fuzzy match service, dashboard, etc.) have **Docker** configurations. We use Docker Compose to orchestrate a stack for local testing – e.g., a container for the FastAPI back-end, one for the classifier (with

PyTorch), one for the database if any, etc. This ensures environment consistency across development machines and servers.

- A CI workflow is established via **GitHub Actions**. Whenever code is pushed, the CI runs tests (unit tests for the Python modules, integration tests for the services) and also builds Docker images. Artifacts like the trained model checkpoint, if updated, are versioned and stored (with hashes to ensure integrity).

- We have set up artifact pinning and dependency locking as part of CI to ensure reproducibility (no surprises from updated pip packages in the future).

The containerization and CI/CD are at a **prototype stage** – they work for our development environment and have been used to deploy the pilot system. Before production, we will expand on this with more robust testing (including load tests, security scans) and possibly integrate with a cloud CI/CD platform for deployment to staging/production environments. Nonetheless, the work done here lays the groundwork for a reliable deployment pipeline and has already improved collaboration (developers can spin up the whole stack easily) and reliability (catching issues early via CI tests).

## Tools and Libraries

The implementation made use of a modern tech stack and a variety of libraries:

- **Programming & AI frameworks:** Python 3.11 was the primary language. We used PyTorch 2.x and Hugging Face Transformers for the NLP model training (DistilBERT, BERT). Scikit-learn was used for some evaluation metrics and utility functions; XGBoost and LightGBM for the forecasting models; Pandas for data manipulation.

- **Matching and parsing:** RapidFuzz for fuzzy string matching; standard Python libraries for JSON and XML handling; lxml for XML parsing; custom code for schema validation.

- **Web and UI:** FastAPI was used to build the REST APIs and WebSocket endpoints (high-performance async support). React with Tailwind CSS (via the Shadcn UI components) was used for the front-end interface.

- **DevOps:** Docker and Docker Compose for containerization; GitHub Actions for CI; and planning for Kafka or Apache Flink in the streaming pipeline (for future real-time data handling).

- **Others:** We also used data analysis notebooks (Jupyter) in development, and tools like LightGBM's Python API for quantile regression.

This diverse toolset helped us quickly prototype and iterate on solutions while leveraging best-in-class libraries for machine learning and data processing.

## Data Sources

A variety of **datasets and data streams** were leveraged to train models, test the system, and drive the reconciliation logic. Below is an overview of key data sources and their purposes in the project:

- **FINRA Fails-to-Deliver Daily Files (2009–2025):** These are daily reports (historically from the SEC/FINRA) listing securities with outstanding fails-to-deliver positions. We used this historical data to train and back-test our models, especially the forecasting model (to learn patterns in settlement fails over time) and as contextual features for classification (e.g., knowing a security is currently on a high-fail list might explain certain breaks).

- **FINRA TRACE Corporate Bond Tape (2015–2025):** The TRACE system records all corporate bond trades. We pulled this data to generate features such as trade volume, price movements, etc., which are relevant for reconciliation (e.g., unusually large trades or price discrepancies might correlate with exceptions). TRACE data also helps in fuzzy matching and linking, by providing an authoritative list of trades that should have been settled.

- **DTCC Security Master File:** A reference dataset from DTCC containing details of securities (e.g., CUSIP, ISIN mappings, security names, issue dates, etc.). This is used to enrich our trade data and ensure consistent identification of securities across systems. For example, if internal data uses ISIN and an external ticket uses CUSIP, the security master lets us map between them to confirm it's the same security.

- **ISO 20022 e-Repository (Schemas and Documentation):** This is not a dataset per se, but the collection of XML schema definitions (XSDs) and documentation for the ISO 20022 message standards. It was crucial for building our message parser. We used the repository to understand the structure of messages (like settlement instructions, confirmations, etc.) and validate/parse them correctly. Essentially, it's the source of truth for how to interpret the fields in those XML messages.

- **Internal Trade Blotters & Reconciliation Logs:** These constitute our **ground truth**. The trade blotter is a record of all trades executed by the firm (and their key details), and the reconciliation logs record historical breaks and how they were resolved (including notes by analysts). We mined these logs to create the labeled dataset for

the ML classifier (i.e., break description -> cause category -> resolution outcome) and to validate the fuzzy matching (by seeing if an external reference was linked to the correct trade). Without this internal data, supervised learning would not have been possible.

- **DTCC "Daily Total UST & Agency Fails":** This is an aggregate time series published by DTCC that measures the total value of U.S. Treasury and Agency bond settlement fails each day. We used this as a macro-level indicator in our forecasting model – it serves as both a target (what we aim to predict for tomorrow) and a feature (today's value helps predict tomorrow's). It also provides context to classify days as high-stress vs normal.

- **Industry Publications and PDFs (DTCC/FINRA articles, etc.):** We gathered textual materials such as regulatory notices, DTCC or FINRA newsletters, and domain-specific papers. These served two purposes: (1) as unsupervised pre-training data to adapt the language model (so it learns vocabulary like "buy-in", "haircut", "fails charge", etc.), and (2) as references for certain rules or scenarios (for example, understanding regulatory jargon when generating explanations).

- **Auxiliary CSV files (e.g., "Daily-File-Table.csv"):** We also used various auxiliary data for lookups and tokenization. The example given, "Daily-File-Table.csv," likely contains meta-data (perhaps mapping certain file names to dates or categories). This kind of data was used to assist in parsing and structuring information, ensuring our pipeline could automatically interpret incoming file names or report types correctly. While not central to the models, these small datasets add polish and robustness to the system.

By combining these sources, we ensured the project is built on rich, relevant data. Historical data provided the **training signal** for ML models; live and reference data feed into daily operations for detection and matching; and external knowledge (schemas, articles) ensured alignment with industry standards and terminology.

## Deliverables & Outputs

The project has produced several tangible **deliverables**, including models, codebases, and documentation. These outputs are either deployed for internal use or prepared for external release as noted:

- **DistilBERT-Reconciler Model Checkpoint:** The trained DistilBERT model for break classification (our primary ML model) is saved as a checkpoint and stored securely on a private AWS S3 bucket (s3://eh-recon-models/v1). This model will be versioned

for future improvements. *Access:* Internal (pending broader release). A formal release is pending a license review because the model was trained on some proprietary data – we need to ensure we can share the weights without violating data usage agreements.

- **Notebook Lineage for Domain Adaptation:** The Jupyter notebooks and logs that document the domain adaptation and fine-tuning of DistilBERT are stored in the repository (src/data/Notebooks/domain-adapted-distilbert/). *Access:* Internal repository. These contain the code for tokenizer adaptation, training, and evaluation, along with metrics plots. This ensures reproducibility and serves as a reference for how the model was built and how it performed.

- **Fuzzy-Matching Microservice Code:** The source code for the fuzzy matching service resides under services/fuzzy_match/ in our codebase. *Access:* Internal (private repo). This includes the FastAPI service wrapper, the core matching logic (using RapidFuzz and embedding models), and configuration for deployment. This service can be containerized and scaled independently.

- **Forecasting Pipeline Scripts:** All code related to the settlement stress forecasting is in models/forecasting/. *Access:* Internal. This includes data processing scripts, model training code for XGBoost and LightGBM, as well as Jupyter notebooks that were used to experiment and determine the best features and model hyperparameters. Additionally, we have CI tests here to verify that the pipeline runs as expected on sample data (this guards against future code changes breaking the forecasting workflow).

- **BERT Baseline Classifier (Open Source Model):** The legacy BERT model checkpoint and configuration have been uploaded to Hugging Face Hub under the name kelvin/bert-breaks-v0. It's provided with an MIT open-source license. *Access:* Public. We also included a model card describing the dataset it was trained on, its performance, and intended use. This baseline model can be used by others for comparison or even as a starting point for similar tasks in different contexts.

- **Triage Dashboard Prototype:** The front-end and back-end code for the dashboard is located in ui/triage-dashboard/. *Access:* Internal for now. This deliverable includes the React app code (JavaScript/TypeScript, component definitions, etc.) and the FastAPI server code that provides the API endpoints. It is in a proof-of-concept state, meaning it's not yet packaged for production release, but it serves as a template for what the production UI would look like.

- **ISO 20022 Parser Library (iso2json):** We have created a library to parse ISO 20022 messages into JSON, informally called "iso2json". It is currently hosted in a private GitHub repository, but we intend to open-source it under an Apache-2.0 license. *Access:* (Planned open-source). The library includes functions to load XSD schemas and validate messages, and then mapping logic to JSON. This will be useful not just for us but potentially for others dealing with ISO 20022 data, hence the plan to open-source once it's production-ready.

- **Benchmark Notebooks & Reports:** In the benchmarks/ directory, we maintain notebooks and reports that benchmark the various models and components. For example, comparing DistilBERT vs BERT on classification (with F1 scores), comparing fuzzy match accuracy against simple exact match baselines, and measuring the forecasting error over time. *Access:* Internal. These reports ensure we have a quantitative record of improvements and can communicate the gains to stakeholders. They are also crucial for any future research papers or whitepapers we might produce, offering reproducible evidence of our system's performance.

All these deliverables are version-controlled and backed up. For any production deployment or external partnership, we can selectively provide relevant assets (e.g., the fuzzy matching service code if integrating with a client system, or the classifier model if we collaborate with a regulator on testing its decisions). The combination of internal and planned open-source outputs also highlights our intention to contribute to the broader community where possible, after navigating legal and compliance considerations.

## Work Remaining / Next Steps

While substantial progress has been made, several tasks remain to reach a fully mature, production-ready solution. The following are the key **next steps and remaining work**, along with target timelines (by quarter) and notable challenges for each:

- **Production-Grade XML/MX Validation** (Target: Q3 2025): We need to elevate our ISO 20022 parsing from PoC to production. This involves handling variant message schemas (different financial institutions might use slightly different message flavors or optional fields) and ensuring performance is optimized (so that parsing hundreds of messages in real-time doesn't become a bottleneck). The challenge lies in the complexity of the ISO 20022 standard – making a robust validator that can gracefully handle unexpected or malformed messages. We plan to incorporate streaming validation and perhaps caching of schema info to speed up processing.

- **Real-Time Reconciliation Model (Streaming)** (Target: Q3–Q4 2025): Currently, our pipeline processes batches (e.g., end-of-day files). The next step is to move to a

truly streaming model where trades and confirmations flow in real time (or near real time), and breaks are detected on the fly. This will likely involve integrating with a message queue or streaming platform like **Apache Kafka**, and possibly using **Apache Flink** or similar for real-time computation of features (windowed aggregations, etc.). The low-latency requirement is a challenge – features that are trivial to compute in batch (like total fails for the day) might need approximation in streaming. We'll also need to ensure the ML model can handle one trade at a time input efficiently, or consider online learning updates. Achieving sub-second detection and alerting for a break after trade execution is the stretch goal here.

- **Triage Dashboard v1.0 (General Availability)** (Target: Q4 2025): We aim to take the dashboard prototype to a production-ready **GA release** by end of 2025. This includes implementing **RBAC (Role-Based Access Control)** so that users have appropriate permissions (e.g., only certain roles can override classifications or close a break, and clients might only see their own data in a multi-tenant scenario). We also plan to integrate a **ticket hand-off** mechanism – for instance, with ServiceNow's API – so that if a break requires action from another team, it can be sent to their queue with a click. Additionally, extensive UX polish is needed: improving the responsiveness, adding help tooltips for new users, and ensuring the interface scales well (hundreds of breaks, multiple days). The challenge is balancing a rich feature set with simplicity, as this tool should make analysts' lives easier, not more complicated. Feedback from pilot users will heavily influence these refinements.

- **Reinforcement Learning on Analyst Feedback** (Target: Q4 2025): Once the system is in use, we expect to collect feedback data – e.g., an analyst might mark a model's suggested resolution as incorrect, or manually resolve a break in a way different from the model's suggestion. We plan to apply **Reinforcement Learning (RL)** to incorporate this feedback and improve the system over time. This could involve designing a reward function that encourages the model to get resolutions right as judged by analysts, and training an agent or refining the model accordingly. A key challenge here is ensuring **privacy and compliance**: we have to log feedback in a privacy-safe way (no sensitive trading data leaking into training datasets inappropriately) and possibly anonymize or aggregate it. Another challenge is the RL design itself – we need to ensure the model doesn't get worse or biased by a small set of feedback, so careful consideration of the learning rate and exploration-exploitation balance will be needed.

- **Regulatory Model Risk Documentation** (Target: Continuous, throughout): As we deploy ML in a financial setting, we must document model risks and controls thoroughly. This is an ongoing task to create and maintain documents that could be shown to **regulators (OSFI, CSA)** or internal risk managers. It will cover things like: model development process, data lineage, validation results, known limitations, bias assessments, and how we govern model updates. The OSFI and CSA guidelines for AI use (which emphasize explainability, data quality, and governance) will serve as a reference for this documentation. The challenge is that this is a continuous effort – every time the model is updated or retrained, documentation must be revised. We plan to integrate this into our workflow, treating the documentation as a living artifact that gets updated alongside code (possibly using automated tools to populate evaluation metrics, etc.).

- **Open-Sourcing Model & Data Snippets** (Target: Pending Legal Approval): We intend to contribute parts of this project to the open-source community. This likely includes the **DistilBERT model weights**, and small snippets of anonymized data (for example, a handful of example break descriptions with PII removed, to illustrate the problem for research purposes). However, this step is pending legal review. The hurdles are ensuring that none of the data used is under a license that forbids release, and that any potentially sensitive information (PII or confidential trade data) is properly redacted. Our compliance and legal teams are evaluating what can be shared. If approved, we might release the model on Hugging Face (with a license and disclaimer), and publish a technical paper or blog with example cases. This could happen in late 2025, but the timeline is uncertain.

- **Scalability & Disaster Recovery (DR) Planning** (Target: Q1 2026): As we move toward production deployment, we must ensure the system is scalable and resilient. This involves deploying the system across multiple regions for disaster recovery, setting up failover for critical components (e.g., if one model service instance goes down, another picks up seamlessly), and having robust data backup and recovery procedures for the audit logs and databases. We will also prepare **playbooks** for operations – documented procedures on what to do if a component fails, how to restore from backups, etc. The goal is to meet enterprise reliability standards by early 2026. Challenges here include cost management (multi-region setups can be expensive, so we need to balance cost with risk) and complexity (ensuring all pieces – data pipeline, model services, UI – are orchestrated in a way that a failover won't break data consistency). We might leverage cloud-managed services for some parts (e.g., using a cloud database with high availability, Kubernetes for orchestration, etc., depending on our IT strategy).

Each of these next steps is tracked in our project plan and has resources allocated. By addressing these, we aim to transition the project from a successful pilot to a fully operational system embedded in daily workflows, with all the necessary support structure around it.

## Impact & Value

If successfully implemented and adopted, this project will deliver significant **impact and value** to the organization and its stakeholders. The benefits include:

- **Dramatic Reduction in Manual Effort:** We estimate a **60% or greater reduction** in manual reconciliation work for operations teams. In practical terms, that translates to hundreds of analyst hours saved per month. Staff can be redeployed to higher-value activities rather than chasing down routine breaks. Fewer manual steps also mean fewer human errors, improving overall accuracy.

- **Lower Settlement Failures and Penalty Costs:** By catching and resolving breaks faster (or even preventing them with predictive alerts), the project helps avoid failed settlements. This has direct financial benefits – for example, avoiding **buy-in procedures** where failing to deliver securities forces the firm to buy at market price, or avoiding regulatory **penalty fees for fails** (such as those imposed under certain regulations for chronic fails). Early warning from our forecaster allows treasury and operations to mobilize resources to settle trades that would have otherwise failed, thus reducing the incidence of fails and their associated costs.

- **Regulatory Compliance and Audit Alignment:** The solution is built with **RegTech** principles, maintaining a detailed audit trail and ensuring explainable outputs from the ML model. This means we can provide regulators and compliance officers with a clear view of how exceptions are handled and why certain decisions were made by the AI (a key for AI governance). The system supports compliance with T+1 settlement cycle rules by ensuring next-day matching and record-keeping are automated and reliable [capco.com](capco.com). It also prepares the firm for potential audits or inquiries, since every action is logged and traceable.

- **Open, Modular Architecture:** The project was designed as a collection of loosely coupled, API-driven components. This modularity makes it easier to integrate with external systems (e.g., a client's platform or a vendor's software). It also invites **community contributions** and future extensions – for instance, by open-sourcing parts of the solution, we may get external developers or researchers to contribute improvements (an external contributor might enhance the ISO parser or add support for another asset class). The architecture's openness accelerates potential vendor

integration as well; if we decide to plug this into a third-party reconciliation software or vice versa, the well-defined APIs will simplify that process.

- **Extensibility to Other Domains:** Although currently focused on fixed-income (bonds), the approach is largely domain-agnostic after some adaptation. With minor retraining or rule adjustments, the **same framework could be applied to equity trades, repo transactions, derivatives, or even crypto asset settlements**. This portability means the investment in this project sets the stage for broader enterprise-wide improvements. We could replicate the success in other departments or even license the technology to partners in different markets. This cross-domain potential increases the long-term ROI of the project.

- **Improved Decision Making and Client Confidence:** Indirectly, by having more reliable and timely reconciliation, management gets a clearer picture of the firm's trading **exposure and operational risks** at any given time. This can inform better decision-making (for example, if certain clients or instruments are frequently causing breaks, that insight can lead to process changes). Clients and counterparties also benefit from fewer settlement issues – leading to improved trust and satisfaction. In an environment where errors can erode confidence, a smooth reconciliation process can be a competitive differentiator.

Overall, the project aligns technological innovation (AI/ML, automation) with business value (cost reduction, risk mitigation, compliance). The measured improvements in efficiency and accuracy translate not only to cost savings but also to qualitative benefits like stronger regulatory relationships and future-proofing the operations for higher volumes and faster settlement cycles.

## Audience and Stakeholder Considerations

This report and the project outcomes are relevant to several key **stakeholders**, each with different priorities and needs. We outline the primary audiences and what aspects of the project are most pertinent to them:

- **Coreledger Leadership & Engineering Teams:** (Coreledger appears to be the internal platform or company driving this project.) For senior management and engineering leadership, the focus is on **progress and resource management**. They will want to know how the project is tracking against its objectives and timeline, any deviations in budget or staffing needs, and how it fits into the broader strategy. We provide them with clear progress updates (milestones achieved, like components completed) and highlight the remaining work with timelines. For engineering leads, details about the architecture, choice of technologies, and any technical risks or

debt are important – they need to ensure the solution is maintainable and scalable in the long run. This report's comprehensive technical sections address those concerns by documenting our implementations and next steps.

- **Dealer & Custodian Prospects (Clients):** These are the financial institutions (banks, broker-dealers, custodians) who might use or benefit from the reconciliation solution, possibly as a product or service we offer. Their main interests are **accuracy, reliability, and integration effort**. They will ask: *Does it actually work better than our current process?* (hence the importance of our high accuracy metrics and pilot results), *Will it fit into our tech stack easily?* (hence the need for APIs, modular design), and *What is the latency?* (they have cut-off times for settlements, so our real-time capability is a selling point). In communicating to this audience, we emphasize the empirical results (96% matching accuracy, 0.88 F1 classification, etc.) and the operational benefits (60% less manual work). We also highlight that it's built with standard tech (Python, REST APIs, etc.), which means integration and adoption would be straightforward.

- **Regulators and Compliance Teams:** Regulators (such as **OSFI and CSA in Canada**) and internal compliance/risk teams are concerned with **control, transparency, and adherence to rules**. They need assurance that the automated system doesn't introduce new risks. For them, we have to demonstrate data lineage (where does each piece of data come from and go), explainability of AI decisions (we can provide rationales or at least categories for why a break was classified a certain way), and robust controls (audit logs, access controls in the dashboard). In the report and in practice, we ensure they know about our tamper-proof logs and model risk documentation efforts. For instance, if asked "how do you know the model isn't making arbitrary decisions?", we can show the validation results and even sample cases with the model's reasoning (perhaps as simple as showing similar past cases it drew on). The compliance team will also look for alignment with regulatory guidelines – our ongoing documentation and the careful consideration of privacy in RL feedback are important points to satisfy them.

- **Open-Source Contributors & Research Partners:** If and when we open-source parts of the project, the **developer community and academic/industry researchers** become an audience. They will be interested in technical details like API specs, data formats, and how to contribute. For them, we plan to provide thorough documentation (likely a README or developer guide in the open-sourced repos) and clear contribution guidelines (coding standards, how to submit pull requests, etc.). In this report, mentioning our intent to open-source and listing what

will be available (model weights, parser library, etc.) sets the stage for that audience. Additionally, highlighting the use of standard open licenses (MIT, Apache-2.0) signals that we welcome external use and collaboration.

- **Grant Agencies and Investors (e.g., Scale AI):** Bodies like **Scale AI** (a Canadian AI innovation funding program) or other investors who have funded or are interested in the project will look for **market impact and innovation**. They want to see that their investment yields tangible improvements in the industry and possibly commercial opportunities. For them, we stress metrics like the 60% efficiency gain, the fact that this addresses a critical need in the context of T+1 settlement change (a timely market problem), and the potential for broader adoption across the financial sector. They are also interested in traction: are there pilot customers or internal deployments? Our report can mention any successful pilot run or interest from certain departments (keeping confidentiality as needed). We have essentially built a case study of AI application in a traditionally manual workflow – grant agencies love to see that kind of success because it justifies their funding in real-world terms. We also underscore how the project keeps the firm at the forefront of technology in capital markets, which is a positive signal for investors regarding competitiveness and future growth.

By tailoring the communication to each of these audiences, we ensure the project's value is understood from all perspectives. The **formal white-paper style** of this report, with an executive summary up front and technical detail in later sections, is intended to cater to both non-technical stakeholders (who may read only the summary and impacts) and technical stakeholders (who will delve into the component details and data). Future presentations or documents can be derived from this master report to focus on specific stakeholder interests (for example, a high-level slide deck for clients or a compliance-focused report excerpt for regulators).

## Conclusion

In conclusion, the Exception Handling and Reconciliation project has made significant strides in modernizing our post-trade operations. We started with ambitious goals to reduce manual workloads, enhance accuracy, and satisfy new regulatory demands; through a combination of advanced NLP, machine learning, and thoughtful system design, we have built a solution that meets those goals and demonstrates measurable success. The work completed – from the DistilBERT classifier to the fuzzy matching service and real-time dashboard – provides a strong foundation for full production deployment.

Looking ahead, the remaining tasks (streamlining real-time processing, reinforcing compliance documentation, and polishing the user interface) are well-defined and scheduled. Tackling these will not only solidify the solution for fixed-income workflows but also pave the way for scaling to other asset classes and use cases. The project stands as an example of leveraging AI in a controlled, effective manner: it delivers clear operational benefits while maintaining the governance and transparency required in a financial setting.

Ultimately, by achieving a tighter, smarter reconciliation process, the firm is better positioned to handle growing trade volumes and faster settlement cycles (like T+1) without a proportional increase in operational risk or cost. This creates a competitive advantage in reliability and efficiency. Moreover, by sharing aspects of our solution and engaging with external stakeholders (be it through open-source or regulatory dialogue), we contribute to industry-wide improvements in how **trade exceptions** are handled. The value created by this project thus extends beyond our own walls, potentially influencing best practices in the broader financial community.

We will continue to monitor the system's performance as it rolls out, gather feedback from users and stakeholders, and iterate. With continued support and cross-team collaboration, we expect to fully realize the vision of an AI-assisted, exception-free trading operation in the coming year, delivering robust results and confidence to all parties involved.

# References, APA

Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., … Zheng, X. (2016). *TensorFlow: A system for large-scale machine learning* (OSDI '16, pp. 265–283).

Brownlee, J. (2020). *XGBoost With Python: Gradient Boosting in Python* (2nd ed.). Machine Learning Mastery.

Chen, T., & Guestrin, C. (2016). *XGBoost: A scalable tree boosting system* (KDD '16, pp. 785–794). https://doi.org/10.1145/2939672.2939785

DTCC. (2025). *Daily Total U.S. Treasury & Agency Fails* [Data set]. Depository Trust & Clearing Corporation. https://www.dtcc.com/

Financial Industry Regulatory Authority. (2025). *Fails-to-Deliver files (2009–2025)* [Data set]. FINRA. https://www.finra.org/

Financial Industry Regulatory Authority. (2025). *TRACE Corporate Bond Transaction Data (2015–2025)* [Data set]. FINRA. https://www.finra.org/

Honnibal, M., & Montani, I. (2023). *spaCy 3: Industrial-strength natural language processing in Python* (Version 3.7) [Software]. Explosion AI. https://spacy.io/

ISO. (2013). *ISO 20022 e-Repository (2013 edition)* [XMI/ISO 20022 schemas]. International Organization for Standardization. https://www.iso20022.org/

ISO. (2021). *XML Tag – Abbreviation list* (Version 2021-12). International Organization for Standardization. https://www.iso20022.org/

Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., … Liu, T. Y. (2017). *LightGBM: A highly efficient gradient boosting decision tree* (NIPS '17, pp. 3146–3154).

McKinney, W. (2010). *Data structures for statistical computing in Python* (SciPy 2010, pp. 51–56).

Microsoft. (2024). *Presidio: An open-source text and data anonymization framework* (Version 2.5) [Software]. GitHub. https://github.com/microsoft/presidio

Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., … Duchesnay, É. (2011). *Scikit-learn: Machine learning in Python*. *Journal of Machine Learning Research, 12*, 2825-2830.

RapidFuzz Team. (2024). *RapidFuzz* (Version 3.5) [Software]. https://maxbachmann.github.io/RapidFuzz/

Sanh, V., Debut, L., Chaumond, J., & Wolf, T. (2020). *DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter*. *arXiv preprint*, arXiv:1910.01108.

Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). *Attention is all you need*. *Advances in Neural Information Processing Systems, 30*, 5998-6008.

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., ... Rush, A. M. (2020). *Transformers: State-of-the-art natural language processing*. *Proceedings of the 2020 EMNLP: System Demonstrations*, 38-45. https://doi.org/10.18653/v1/2020.emnlp-demos.6

Zhang, C., & Pafka, S. (2023). *PyTorch 2.1* [Software]. Meta AI. https://pytorch.org/