Carnegie Mellon University

Software Engineering Institute

# Automating Mismatch Detection and Testing in ML Systems

**NOVEMBER 14, 2022**

Grace Lewis
Principal Researcher and TAS Initiative Lead

# Organizations Struggle with Moving Machine Learning (ML) Components into Production Systems

Challenges include the following:

- ML component performs poorly because model training data is different from production data.
- Large amounts of glue code need to be developed because ML component input/output does not match with system components.
- Production environments and tools are not set up to detect model problems or to collect the necessary data for model troubleshooting and retraining.
- Systems perform poorly (or are unable to deploy) because available computing resources are insufficient to support model inference requirements.
- Organizations acquire ML components that they do not know how to test properly.

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

3

# Different Teams and Workflows

Data Scientists, ML Engineers

Software Engineers

Operations Staff

The development of ML-enabled systems* typically involves three separate activities and workflows.

- model development
- model integration and testing
- model operation

… performed by three different and separate teams

- data science or ML engineering
- software engineering
- operations

… and often no systems context.

\* We define an ML-enabled system (or ML system for short) as a software system that includes one or more machine learning (ML) components.

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

4

# Model Development Workflow

Examples of information required to make better decisions:
- system context
- upstream and downstream components
- production data

**Model Development Environment** *(Data Scientists, ML Engineers)*



Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

5

# Model Integration and Testing Workflow



**System Development Environment** *(Software Engineers)*

ML-Enabled System

ML Component

Output

Test Data

Output

Upstream Components (Test)

Data Pipeline

Trained Model

Downstream Components (Test)

Test Data

Testing Tools

Test Results

Test Cases

Examples of information required to make better decisions:
- system context
- test cases
- test data
- production environment

Legend

3rd Party or Custom Tool

Code

Artifact

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

6

# Model Operations Workflow



Examples of information required to make better decisions:
- system context
- monitoring requirements
- resource requirements

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

7

# ML Mismatch

**ML mismatch** is a problem that occurs in the development, deployment, or operation of an ML-enabled system when different stakeholders—data scientists, ML engineers, software engineers, operations, system owners—make **incorrect assumptions** about systems elements that result in a negative consequence.

As the DoD adopts machine learning to solve mission-critical problems, the inability to detect and avoid ML mismatch creates delays, rework, and failure in the development, deployment, and evolution of ML systems.

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

8

# Previous Work: Descriptors for ML Systems

**Trained Model 36%**

- 11% Evaluation Metrics
- 14% Decisions, Assumptions, Limitations & Constraints
- 8% Versioning
- 17% API/ Specifications
- 2% Data Buffering
- 12% Programming Language/ ML Framework/ Tools/ Libraries
- 14% Model Output Interpretation
- 5% System Configuration Requirements
- 17% Test Cases & Data

**Operational Environment 16%**

- 32% Computing Resources
- 14% Required Model Inference Time
- 54% Runtime Metrics & Data

**Task and Purpose 15%**

- 15% Usage Context
- 18% Task
- 26% Success Criteria
- 12% Data Rights & Policies
- 29% Business Goals

**Raw Data 10%**

- 13% Proxy Data
- 4% Restrictions
- 31% Data Dictionary
- 4% Anonymization
- 48% Metadata

**Development Environment 9%**

- 5% Development & Integration Timelines
- 10% Computing Resources
- 40% Upstream and Downstream System Components
- 45% Programming Language/ ML Framework/ Tools/ Libraries

**Operational Data 8%**

- 16% Data Syntax & Semantics
- 21% Data Sources
- 5% Data Rates
- 21% Data Pipelines
- 37% Data Statistics

**Training Data 6%**

- 62% Data Preparation Pipelines
- 15% Versioning
- 23% Data Statistics

We developed a set of machine-readable descriptors (JSON schema) that define system attributes that need to be specified to avoid mismatch.

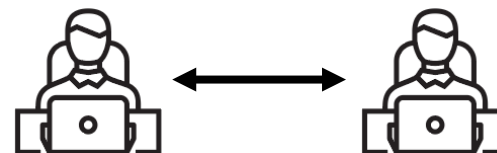- system context
- raw data
- training data
- data pipeline
- trained model
- development environment
- production environment
- production data

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

9

# Project Objectives

Develop a suite of tools to

- automate ML mismatch detection.
- demonstrate extending and using descriptors to support testing production readiness of ML components.
- validate research results with DoD systems and collaborators.
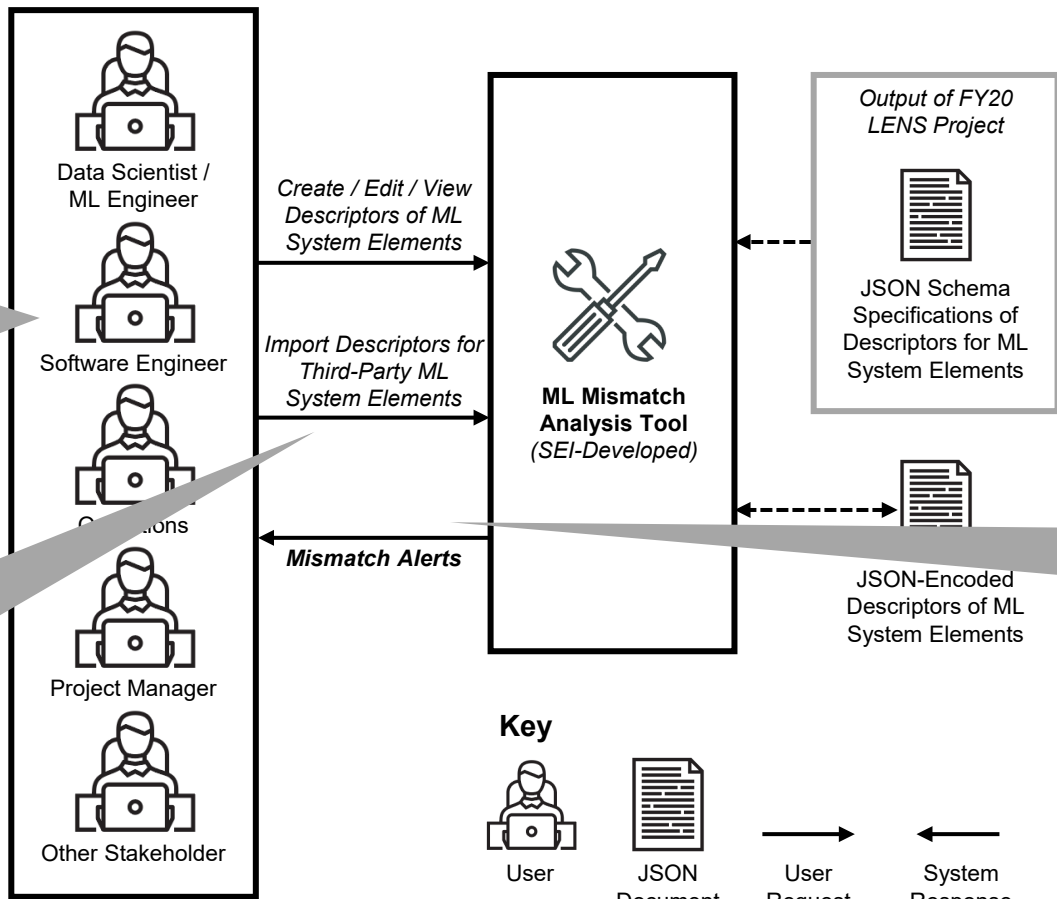    - We are actively looking for project collaborators.

**Collaborators Wanted!**

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

10

# TEC: ML Mismatch Detection Tool

Stakeholders develop and share descriptors for their parts of the system and can view any descriptors.

As system elements are acquired (if applicable), descriptors are imported.

Data Scientist / ML Engineer

Software Engineer

Operations

Project Manager

Other Stakeholder

Create / Edit / View Descriptors of ML System Elements

Import Descriptors for Third-Party ML System Elements

Mismatch Alerts

ML Mismatch Analysis Tool (SEI-Developed)

Output of FY20 LENS Project

JSON Schema Specifications of Descriptors for ML System Elements

JSON-Encoded Descriptors of ML System Elements

Mismatch rules encoded in the tool serve to generate alerts.

**Key**

User

JSON Document

User Request

System Response

Data Flow

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

11

# TEC: ML Mismatch Detection Tool — Demo



Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

12

# ML Component Testing

Testing of ML components is a known challenge, especially for organizations that acquire ML components.

- Results from our ML mismatch study show that one of the top causes for mismatch is lack of information on how to test ML components.

- Results from our study on collaboration challenges in ML systems development show that teams do not exchange information required for testing, leading to problems in production.

- These two findings are consistent with published practitioner studies, which highlight the need for tools and realistic techniques for testing ML-enabled software systems.

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

13

# Testing Production Readiness of ML Components

We define production readiness based on these ML component attributes:

- **Ease of Integration**: ML component inputs and outputs are compatible with upstream and downstream components in the production system.
- **Testability**: The ML component provides either (1) specifications, test cases, or test data that enable testing by software developers or external QA teams or (2) evidence of testing or evidence that teams have followed best practices.
- **Monitorability**: The ML component produces information that monitoring components can use in the production system to detect potential problems.
- **Maintainability**: The ML component defines (or produces) data that teams can use for model retraining and troubleshooting.
- **Quality**: The ML component meets quality requirements.
  - model requirements (e.g., accuracy)
  - system requirements (e.g., inference time, resource consumption)

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

14

# ML Component Testing Assistant

**Extension of descriptors** to support testing for **production readiness** of ML components

Initial testing capabilities targeted at acquisition organizations or QA teams as they perform **acceptance or independent testing** of ML components — where the least amount of understanding and tool support exists

*Output of FY20 LENS Project*

JSON Schema Specifications of Descriptors for ML System Elements

**ML Component Testing Assistant** *(SEI-Developed)*

Test Case Generation

Test Data Generation

ML Model Profiling

Other Testing Assistance

Data Scientist / ML Engineer

Software Engineer

QA

JSON-Encoded Descriptors of ML System Elements

Testing Advice / Recommendations / Results / Descriptor Inputs

ML Component Data

Test Cases and Data

Complete / Execute Test

Profiling and Other Testing Tools

**ML Component (Python)**

Third Party Python Unit Testing Tool (e.g. unittest, mocktest, pytest)

**Legend:**

User

JSON Document

Third Party Tool

ML Component

System Response

Data Flow

Data Flow

User Request

# Implementing and Evaluating Several Approaches to Testing for Production Readiness

| Model Quality | **Completed**: tool for train-test leakage detection |
|---|---|
| **Ease of Integration / Testability** | In progress: black box testing of ML component capabilities |
| **Production Readiness** | In progress: evidence-based process and tool for independent testing of ML component production readiness<br>• mapping practitioner interview data from our two studies to production-readiness attributes<br>• mapping existing SE practices and tools presented as ML testing tools to production-readiness attributes<br>• refining definition of production readiness based on findings and collaborator discussions |

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

16

# Train-Test-Leakage Detection

Trust in models developed by external teams was identified as a collaboration challenge for organizations building ML systems.

- Model developers often report high performance during test and evaluation because models are overfit—inadvertently or intentionally.
- A cause for overfitting is train-test leakage—leaking information about test data into training data.

Developed a static analysis tool that uses data flow analysis and pointer analysis for detection of the following three types of train-test leakage:

- **Overlap**: Test data is directly used as input for training or hyper-parameter tuning.
- **Multi-Test**: Test data is used repeatedly for evaluation.
- **Preprocessing**: Test data and training data are preprocessed together (e.g., normalization, feature selection, vectorization).

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

17

# Looking for Collaborators

As an applied R&D organization, we need collaborators to inform and validate our research through participation in different types of activities.

Near-term, we are looking for

- organizations willing to use and evaluate the automated mismatch detection tool and provide feedback.
- organizations or teams (e.g., DT&E, OT&E) tasked with testing ML components (or ML systems) developed by other organizations to
  - participate in discussions related to definition of production readiness and solution brainstorming.
  - evaluate and use developed processes and tools and provide feedback.

Longer term, we are looking for organizations to integrate SEI-developed approaches into their ML system development and testing processes and tool chains.

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

18

# Team

## SEI Team

Contact us at
info@sei.cmu.edu

**Grace A. Lewis, Ph.D. (PI)**
Principal Researcher and TAS
Initiative Lead

**Rachel Brower-Sinning, Ph.D.**
Machine Learning Research Scientist

**Alex Derr**
Associate Software Engineer

## Collaborators

**Christian Kästner, Ph.D.**
Associate Professor
Carnegie Mellon University
School of Computer Science

**Chenyang Yang**
Ph.D. Student
Carnegie Mellon University
School of Computer Science

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

19

# Publications

ML Mismatch Study: Lewis, G. A., Bellomo, S., & Ozkaya, I. (2021, May). Characterizing and Detecting Mismatch in Machine-Learning-Enabled Systems. In 2021 IEEE/ACM 1st Workshop on AI Engineering-Software Engineering for AI (WAIN) (pp. 133-140). IEEE. https://arxiv.org/abs/2103.14101

Collaboration Challenges Study: Nahar, N., Zhou, S., Lewis, G., & Kästner, C. (2022). Collaboration Challenges in Building ML-Enabled Systems: Communication, Documentation, Engineering, and Process. Proceedings of the 44th International Conference on Software Engineering (ICSE 2022). https://arxiv.org/abs/2110.10234

Distinguished Paper Award

Train-Test Leakage Detection: Yang, C., Brower-Sinning, R., Lewis, G. A., Kästner, C.(2022). Data Leakage in Notebooks: Static Detection and Better Processes. Proceedings of the 37th IEEE/ACM International Conference on Automated Software Engineering (ASE 2022). https://arxiv.org/abs/2209.03345

Automating Mismatch Detection and Testing in ML Systems
©2022

[DISTRIBUTION STATEMENT A] Approved for public release and unlimited distribution.

20