# Tracil: AI-Powered Traceability Tool Across CDISC Standards

Kexin Guan, Scientist, Statistical Programming, Merck & Co., Inc.
Junze Zhang, Scientist, Statistical Programming, Merck & Co., Inc.
Anthony Chow, Executive Director, CDISC

# Meet the Speakers

## Kexin Guan

Title: Scientist, Statistical Programming

Organization: Merck & Co., Inc.

Kexin Guan is a Statistical Programmer at Merck & Co. Inc., where she has been part of the Oncology Early Development group since December 2022. She holds a Master's degree in Biostatistics and Bachelor's degree in Applied Mathematics and Statistics.

## Junze Zhang

Title: Scientist, Statistical Programming

Organization: Merck & Co., Inc

Junze Zhang is a Scientist at Merck & Co. Inc., supporting early oncology statistical programming. He earned his Master's in Computer Engineering from NYU and Bachelor's in Computer Science from Oregon State University.
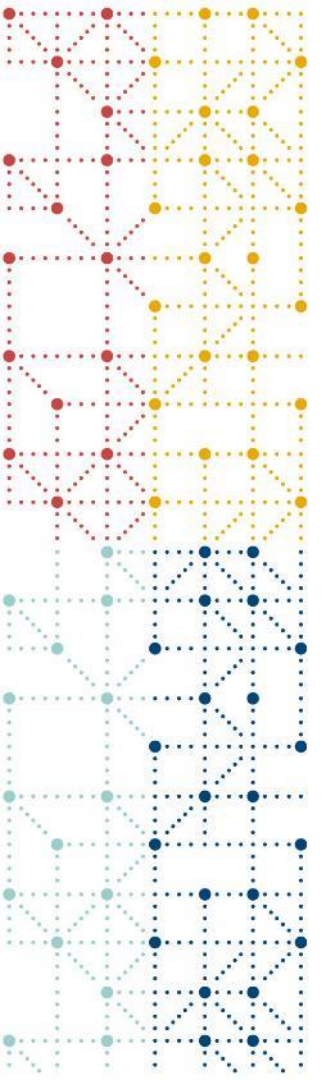
# Disclaimer and Disclosures

- *The views and opinions expressed in this presentation are those of the author(s) and do not necessarily reflect the official policy or position of CDISC.*

- *The views and opinions expressed in this presentation are those of the author(s) and do not necessarily reflect the official policy or position of Merck & Co., Inc.*

# Agenda

1. Background & Motivation
2. Demo
3. App Structure
4. Backend AI Workflow
5. Summary, Lessons Learned & Future Steps

# Background & Motivation

# The Traceability Problem

- Clinical data flow is **complex,** spanning **multiple** systems and **silos**.
- Each layer (Protocol → CRF → SDTM → ADaM → TLFs) adds transformation logic.
- Manual tracing lineage = **heavy review time**.

# Our Project Goals

- **Automate Lineage Inference**
  - Extract relationships from existing metadata without manual mapping.
- **Support CDISC Standards End-to-End**
  - Protocol (USDM) ↔ CRF ↔ SDTM (Define.xml) ↔ ADaM (Define.xml) ↔ TLF (ARS)
- **Provide Explainable AI Results**
  - Every link comes with a natural-language justification.
- **Deliver Human-Friendly Visualization All in One Place**
  - Interactive graphs for regulators, programmers, and statisticians.
- **Explore Data with Natural Language**

# User Flow

## Inputs

- Protocol (USDM / PDF)
- CRF (annotated CRF)
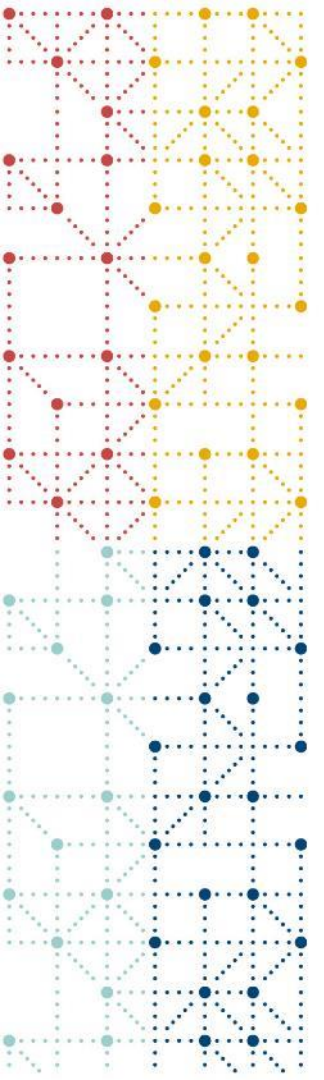- SDTM / ADaM (Define.xml / spec.xlsx)
- TLF (RTF or ARS/ARD)

## Processing

- Preprocessing
- AI model analysis
- Lineage image generation

## Outputs

- Interactive Lineage Graph

# Demo

# Upload & General UI

DEMOS

# ADaM

# Objectives

# Search With Natural Language

# TLF Details with ARS

# App Structure

# App Structure

- Frontend: Javascript with React/Next.js

- Backend: FastAPI

- AI core: Python AI Engine
  - GPT-4o (mini)
  - Text-embedding-3-model from OpenAI

# Backend AI Workflow

# How Tracil's Backend Works

- **Pre-processing**: Python code to convert structured (XML, JSON) and semi-structured (PDF/CRF) inputs into a unified, machine-readable JSON file

- **LLM Reasoning:** passes the inputs to an LLM reasoning layer to infer variable derivations, dependencies, and gaps, returns a JSON graph

- **Post-processing:** standardize the JSON lineage graph so that it can be visualized interactively in the frontend UI

# Pre-Processing: Parsing & Normalization Layer

**Input formats supported:** Specification (.xlsx), define.xml, aCRF, ARS JSON, Protocol PDF, USDM, and TLFs (RTF)

**FastAPI endpoint `/process-files`:**

- Extracts variable metadata, derivation notes, and dataset context
- Normalizes the information extracted into unified schema of names, domain tags, and relationships information

**Why important:**

- Provides the frontend with structured data for visualization and user interaction.

- gives LLM a clean, consistent input so it can reason semantically

# LLM Lineage Builders

**Three Routes by Target**

o Detects what you're tracing: protocol endpoint, ADaM/SDTM variable, table/cell; sends it to the right builder

**Gather the Evidence**

o Collects all supporting pre-processed metadata from the session (aCRF index, protocol text, unified JSON returned by /process-files API, etc.)

o These documents form the "evidence base" for reasoning

**Find What Matters (Chunk + Retrieve)**

o Splits large documents into small readable sections

o Converts each into numerical "embeddings" so the AI can compare meanings

o Selects only the Top K (≈12) most relevant pieces to focus the analysis

**Ask the AI Model**

o Sends the selected context to GPT-4o for reasoning (mini model as fallback).

o The model returns a structured JSON lineage, showing each variable, link, and explanation.

**Parse + Check**

o Cleans the graph and validates to ensure it's proper JSON

# Prompt Design

- **Clear Role Definition:**
  - *System* = "Senior CDISC standards expert"; explicit backtrace & forward instructions; closed node types and canonical IDs (e.g., `ADSL.AGE`, `DM.BRTHDTC`)

- **Structured Thinking:**
  - Each task follows a fixed schema (variable, endpoint, or table) so the AI always knows what format to produce

- **Evidence packing:**
  - The AI reviews only the top relevant document sections, tagged as evidence, before reasoning; pins the target node (e.g., `Target variable: ADSL.TRT01AN`)

- **Explainable Results:**
  - Every link in the lineage includes a short explanation starting with [direct], [reasoned], or [general], and cites where the information came from (define.xml, ARS, CRF, protocol)

- **Direction rules:**
  - hard requirement to emit edges upstream → downstream only; forbid illegal shortcuts (e.g., ADaM → CRF)

- **Resilience:**
  - chat called with `temperature=0.0` + `response_format={"type":"json_object"}`; Built-in checks handle small formatting issues automatically, so results remain valid and reproducible

cdisc

# Post-Processing: Normalize, Validate, Connect

**Normalize the Structure**
- Add tags and standardize names so the graph speaks one language. (e.g., `ADVS.AVAL` → ADaM variable, `VS.VSORRES` → SDTM variable)
- Fix missing labels and edge directions for clear flow (`source` → `result`)

**Validate the Connections**
- Check each link has real evidence; flag gaps if missing
- Remove duplicates/orphans, and add short explanations

**Connect Missing Pieces**
- Re-query files when links are incomplete (e.g., ADaM variable with no SDTM parent)
- Ensure full trace: Protocol → CRF → SDTM → ADaM → TLF

# Conclusion & Future Steps

# Key Takeaways

- Tracil automates lineage across Protocol → CRF → SDTM → ADaM → TLF using AI reasoning

- Converts protocol endpoints, aCRF, define.xml, specifications,  ARS/ARD into a **unified JSON schema**

- Provides **explainable AI outputs** with clear variable relationships

- Makes traceability **easy, fast, interactive, and follows CDISC standards.**

# Limitations

- **Model Accuracy & Stability:** Same input can yield slightly different results due to the nature of LLM

- **Limited Data:** Few open, CDISC-compliant datasets restrict realistic fine-tuning

- **Validation Gap:** Need standardized validation methods to ensure reliable outputs

# Future Steps & Vision

- **Model Expansion:** Test across different LLMs (GPT-5, Gemini, Claude, etc.)

- **Confidence Scoring + User Feedback:** Quantify AI certainty and learn from human corrections

# Thank You!

# Contact Information

- Junze Zhang
  - Junze.zhang@merck.com

- Kexin Guan
  - Kexin.guan@merck.com

- Anthony Chow
  - Achow@cdisc.org

- Tracil GitHub Repository:
  - https://github.com/1mgroot/Tracil