# cdisc

## 2022
## US
## INTERCHANGE
### 26-27 OCTOBER | AUSTIN

## Practical Steps for Implementing ML for SDTM Mapping

Presented by Sharon Rossouw, Director, Biostatistics, Bioforum

# Meet the Speaker

## Sharon Rossouw

**Title:** Director, Biostatistics

**Organization:** Bioforum

There is a farm in Africa…. I grew up on a farm in the Zimbabwean bushveld. I completed my schooling and university education in South Africa culminating with a Masters in Biostatistics.
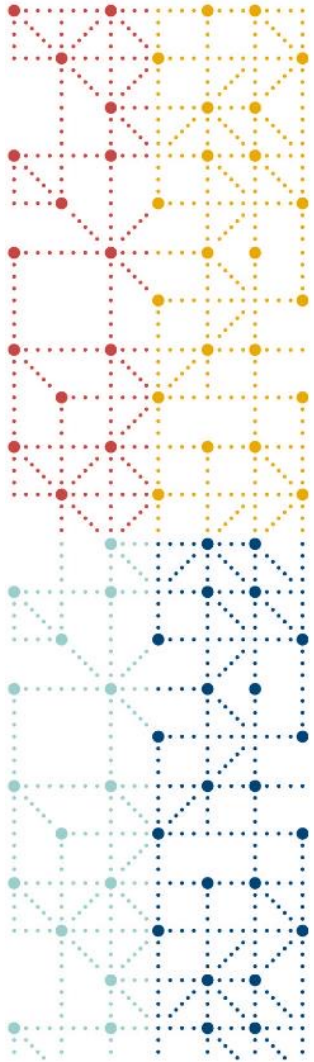
A biostatistician over 25 years of experience providing biostatistical and medical writing services to the pharmaceutical industry and academic institutions.

I am passionate about the training and development of biostatisticians and statistical programmers and have a special interest in process development and implementation.

# Disclaimer and Disclosures

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

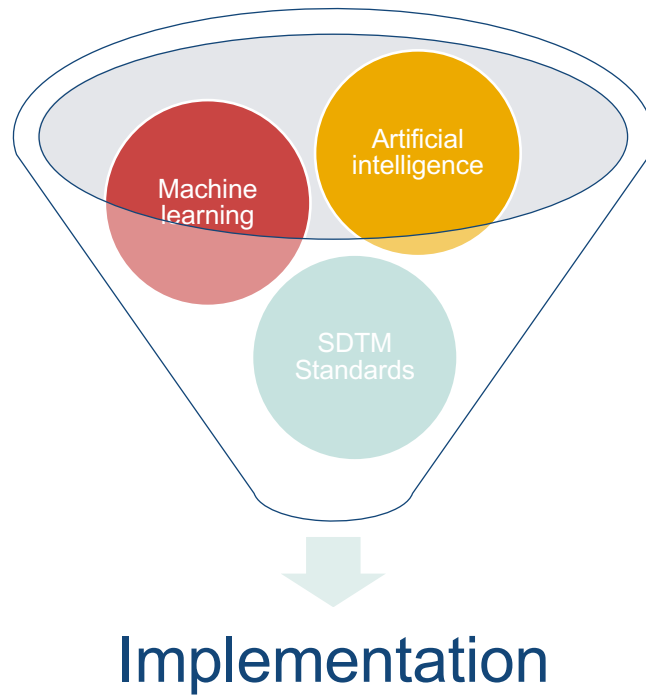- *The author has no real or apparent conflicts of interest to report.*

## Agenda

1. Background
2. Build steps: Develop the models
3. Taking a step back
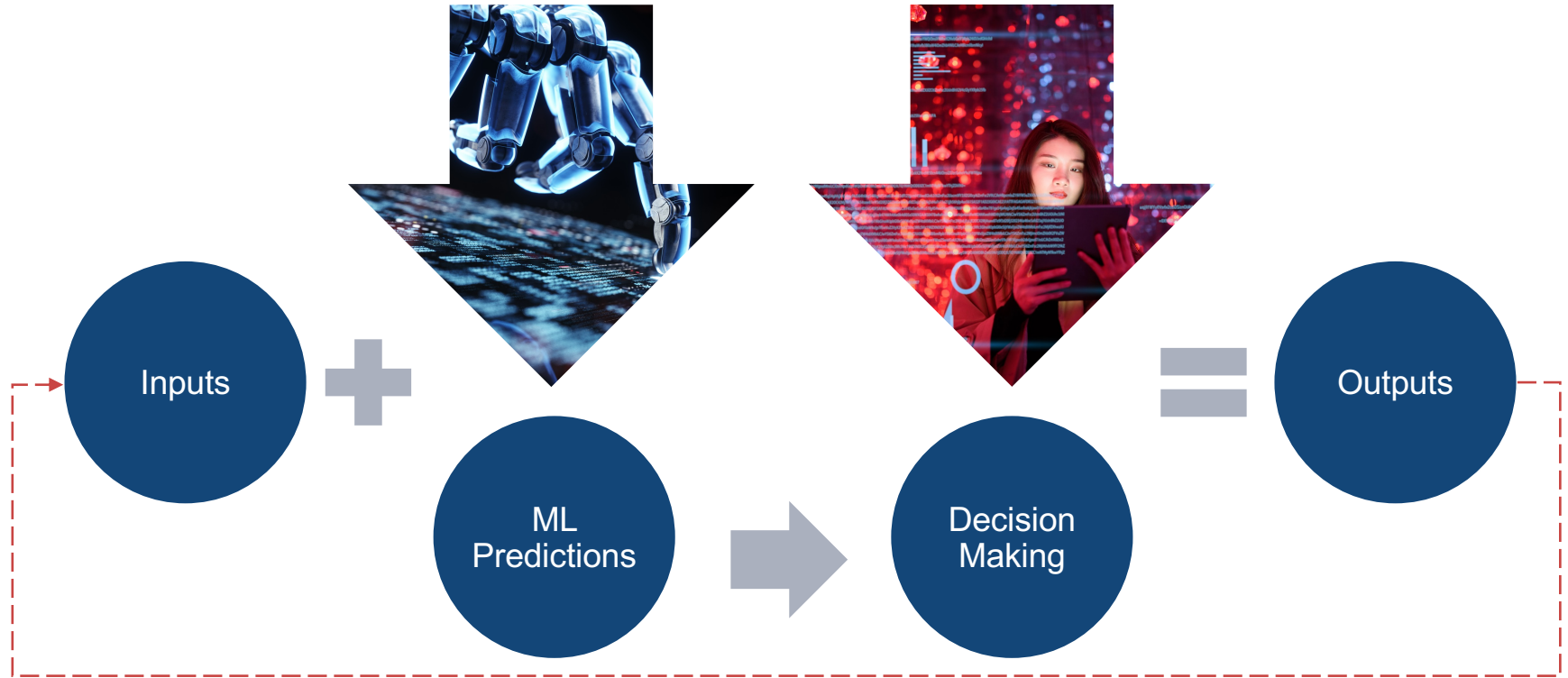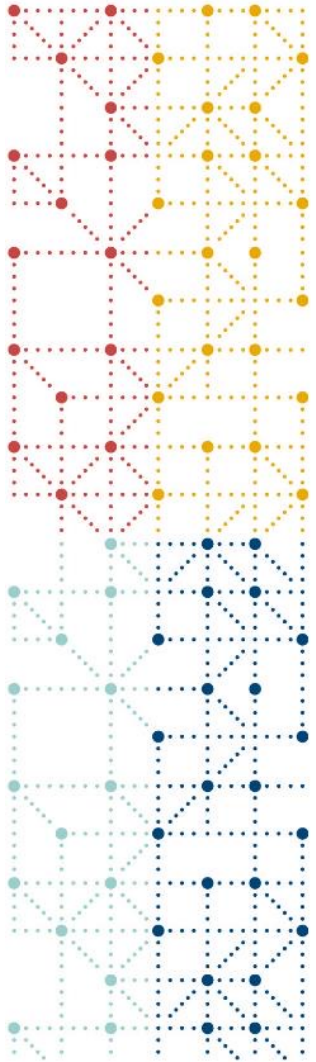4. Refine steps: Applying the models
5. The view from our destination
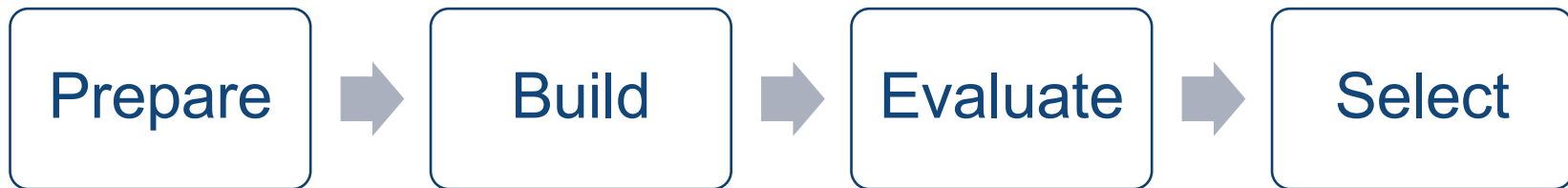
# What are the steps from concept to reality?

Machine learning

Artificial intelligence

SDTM Standards

## Implementation

cdisc

# Where to use ML in the SDTM mapping process?

Inputs + ML Predictions → Decision Making = Outputs

# Build steps: Develop the models

# Build steps: Develop the models

| Prepare | → | Build | → | Evaluate | → | Select |
|---------|---|-------|---|----------|---|--------|

**1st step: Identify the pieces**
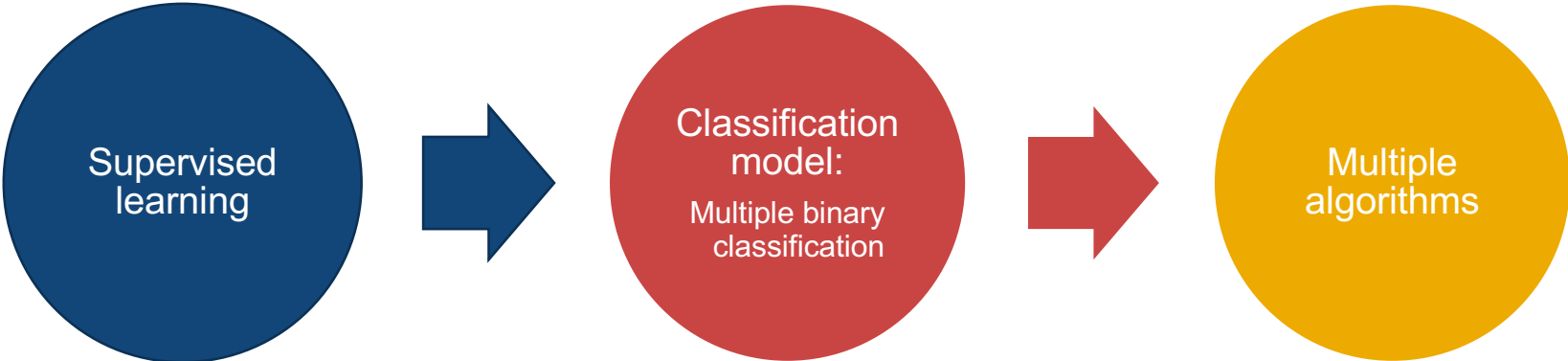- Domain
- Variable
- Controlled terminology
- …

# STEP 1: Prepare

- Training set ➔ Pre-mapped trials

- Raw variables labelled:
  - Domain(s)
  - Variable, etc.

- Raw variable feature extraction:
  - Raw data file characteristics
  - Raw variable metadata
  - Variable values
  - Trial documents
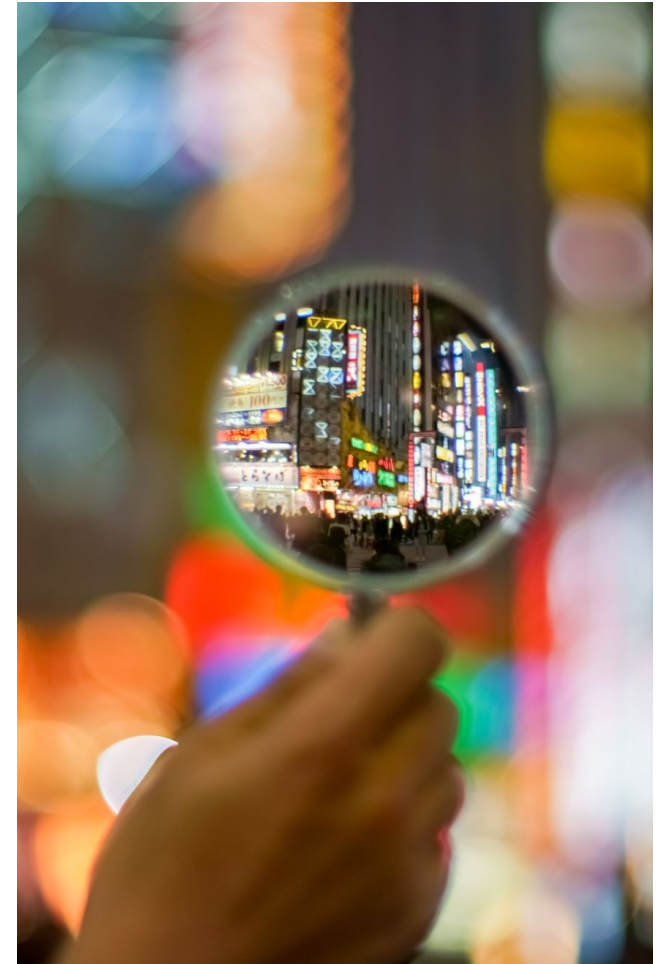
- Training set was tailored to the task

Therapeutic areas

Legacy/ ongoing trials

Training set

Research phases

EDC systems

Trial sponsors

# STEP 2: Build the models

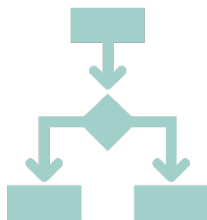Supervised learning → Classification model: Multiple binary classification → Multiple algorithms

# STEP 3: Evaluate the models

- Test model performance
  - Cross validation methodology

- Results
  - Vector of probabilities for each raw variable
  - Probability = Likelihood of mapping to target

- Simple decision rule
  - Select target with highest probability

- Compare
  - Selected target ⇔ Pre-mapped target
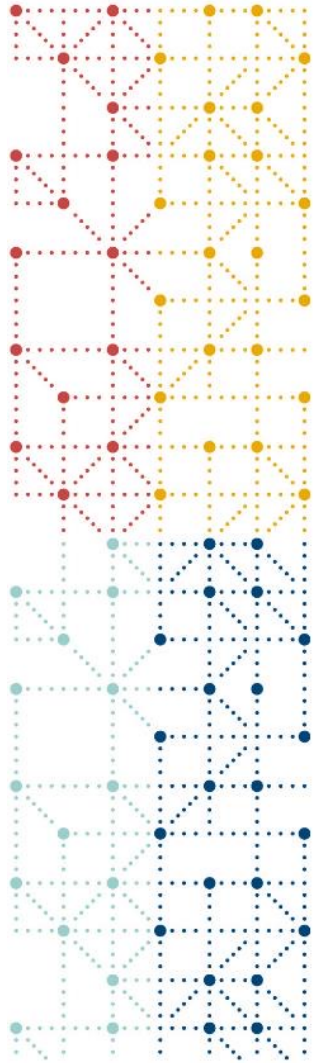
# STEP 4: Select the models

**Domain & Variable mapping**

Random Forest

**Controlled terminology mapping**

Natural Language Processing

- For simplicity, the domain model results are presented
  - Variable model results are briefly mentioned
  - Controlled terminology results are a "forthcoming attraction"

# Where did the first few steps lead?

# Model accuracy

- Confusion Matrix
  ➔ Indication of model quality
- Diagonal frequencies = correct predictions
- Model is mostly correct
- Accuracy:
  - Domain-level ⇨ 71.5% (41 trials)
  - Variable-level ⇨ 83.8% (61 trials)

Columns = Predicted domain

Rows = Pre-mapped domain

| | LB | MH | AE | CM | TS | PR | IE | DS | PC | PE | CO | DM | EC | QS | EG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LB | 3,101 | 3 | 7 | 3 | | 162 | 8 | 11 | 21 | 8 | 38 | 1 | 8 | 3 | |
| MH | 1 | 1,865 | 1 | 30 | 2 | 43 | 7 | 71 | 1 | 2 | 3 | 7 | 16 | 10 | |
| AE | 13 | 25 | 1,688 | 34 | - | - | 13 | 2 | - | 1 | 6 | - | 2 | - | |
| CM | 19 | 109 | 48 | 1,476 | 1 | 18 | 5 | 14 | 10 | 10 | 10 | - | 77 | 3 | 1 |
| TS | 6 | 5 | - | | 1,344 | 1 | 2 | 4 | 48 | 3 | - | 3 | 3 | - | |
| PR | 235 | 53 | 10 | 44 | 12 | 638 | 13 | 34 | 73 | 25 | 11 | 6 | 76 | 14 | |
| IE | 6 | 1 | - | - | 2 | 2 | 1,031 | 60 | - | 1 | - | 1 | 1 | 1 | |
| DS | 16 | 16 | 5 | 6 | 2 | 9 | 35 | 676 | 4 | 4 | 5 | 51 | 26 | 2 | |
| PC | 24 | - | - | - | 38 | 63 | - | 2 | 838 | 4 | 1 | - | 5 | 5 | |
| PE | 4 | 5 | - | 6 | 2 | 8 | - | 1 | 14 | 749 | 15 | - | - | 1 | 2 |
| CO | 46 | 4 | 5 | 6 | 5 | 4 | - | 20 | 95 | - | 601 | 5 | 4 | 1 | 5 |
| DM | 12 | 2 | 4 | 1 | - | 4 | 3 | 34 | 5 | - | 2 | 677 | 7 | - | |
| EC | 64 | 13 | 6 | 59 | 13 | 38 | 4 | 18 | 7 | 4 | - | 9 | 485 | 2 | |
| QS | 5 | 15 | 18 | 16 | 13 | 17 | 10 | 11 | 9 | 3 | 4 | 4 | 2 | 834 | |
| EG | 8 | - | 10 | 8 | 3 | 7 | 1 | - | | - | | - | 2 | 2 | 561 |
| ST | 19 | 1 | 4 | 1 | 3 | 12 | | 9 | | 11 | | 3 | 2 | 5 | |

# A closer look

Most frequently confused:

MH instead of CM

PR instead of LB

LB instead of PR

…
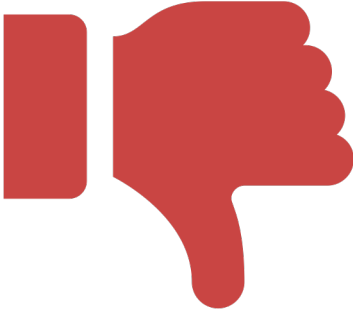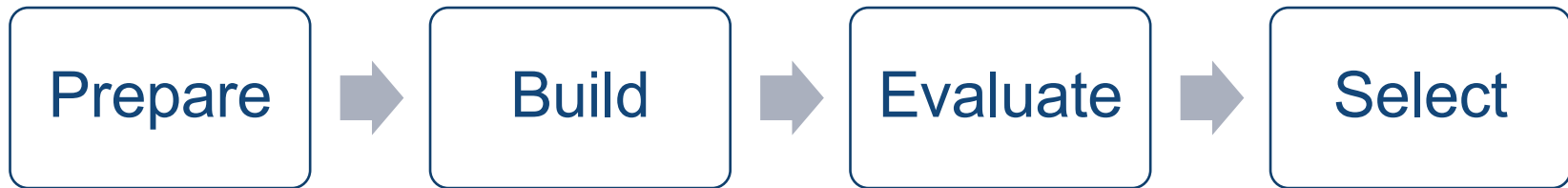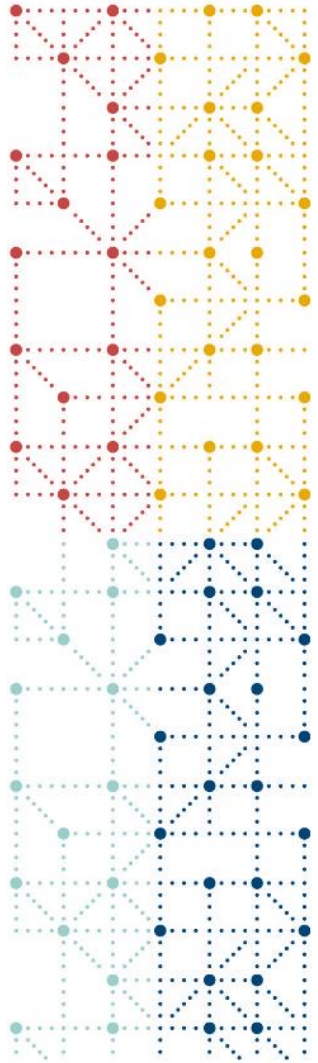
# Taking a step back

# Reviewing the journey

# So what was the next step?

Prepare → Build → Evaluate → Select

# Refine…
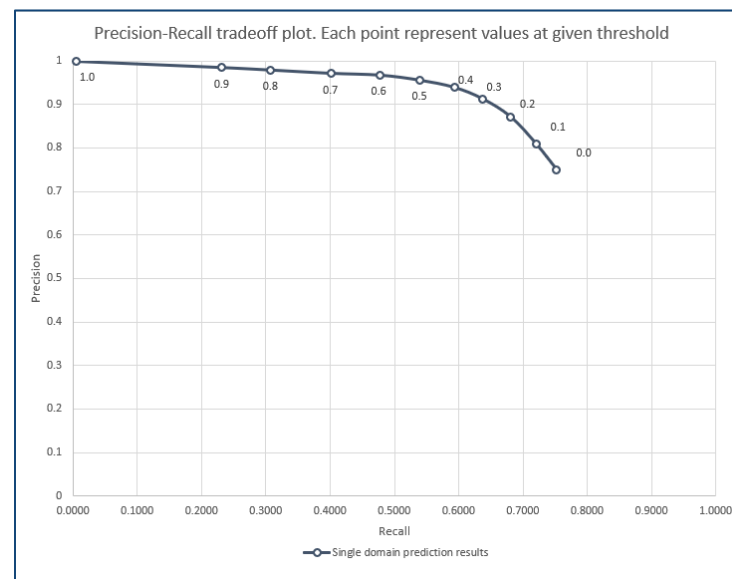
# Refine steps: Domain models

# STEP 1 – Remove predictions with low confidence levels

Remove predictions if confidence ≤ 0.3 ➔ 69.7% of variables retained
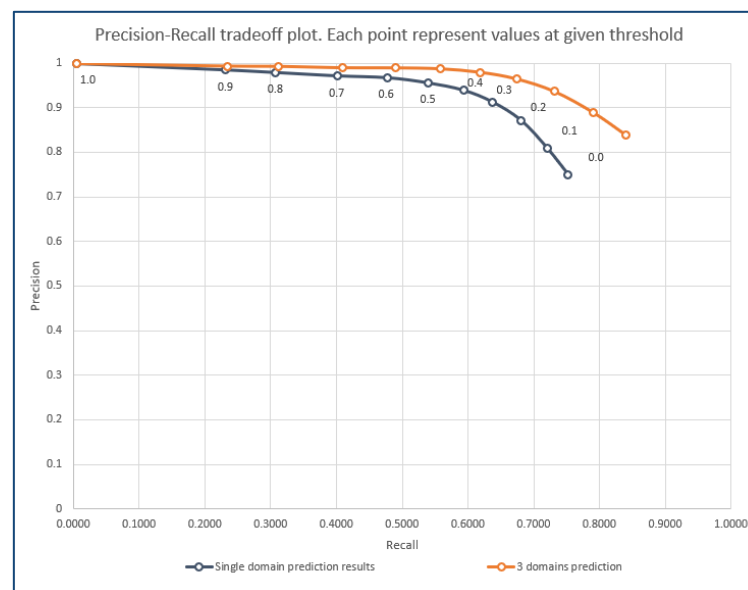- Domain precision: 91.4%, recall of 63.7% (41 trials)

# STEP 1 – Impact

- Now our users were interested…
  - Trust had increased
  - Some trials: nearly error-free

- Investigate mistaken predictions
  - "Correct" recommendation was 2nd/3rd on list
  - How can we adapt the implementation of the models?

# STEP 2 - Provide the top 3 most likely predictions

- Change the decision rule:
  - "Correct" target in the top 3 most likely recommendations

- Using the 0.3 threshold:
  - Precision ➔ 96.6% (STEP 1: 91.4%)
  - Recall ➔ 67.3% (STEP 1: 63.7%)



Precision-Recall tradeoff plot. Each point represent values at given threshold

# STEP 2 – Impact



- Users loved this!
- ⬆ accuracy = ⬆ interest

- Next step:
  - Users had to continuously evaluate 3 options, even when the "perfect" fit was obvious
  - Is there somewhere in between?

# STEP 3 - Dynamic predictions based on cumulative confidence thresholds

- Alternative method investigated:
  - Dynamic cumulative approach

- Tradeoff:
  - Provide a single target
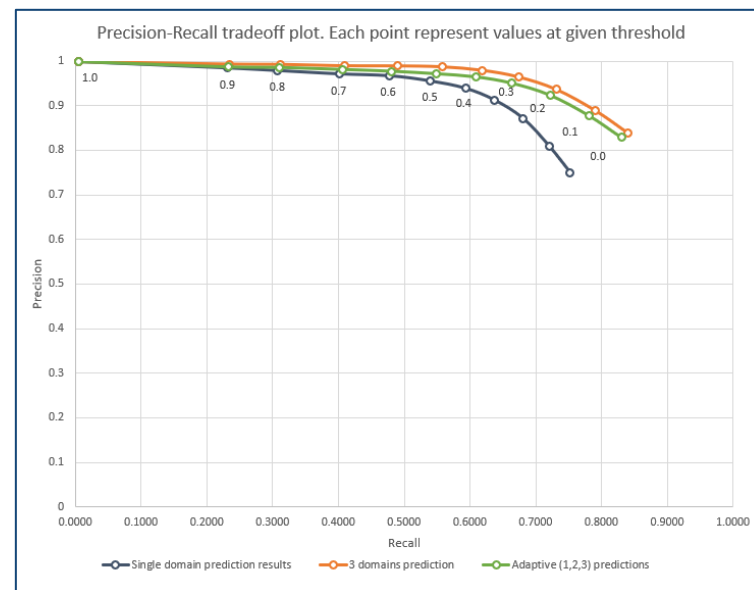  - Need for high confidence

[Top 1] > 0.7 → 1 Prediction

[Top 1+2] > 0.7 → 2 Predictions

[Top 1+2+3] > 0.3 → 3 Predictions



Precision-Recall tradeoff plot. Each point represent values at given threshold

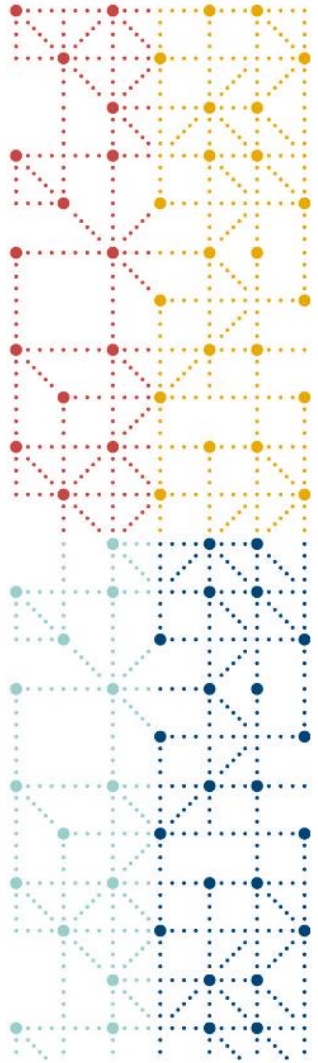Single domain prediction results — 3 domains prediction — Adaptive (1,2,3) predictions

# STEP 3 – Impact

- Maintained high precision
- Single target for about 60% of predictions
  - Other variables: Users could select the prediction from the list
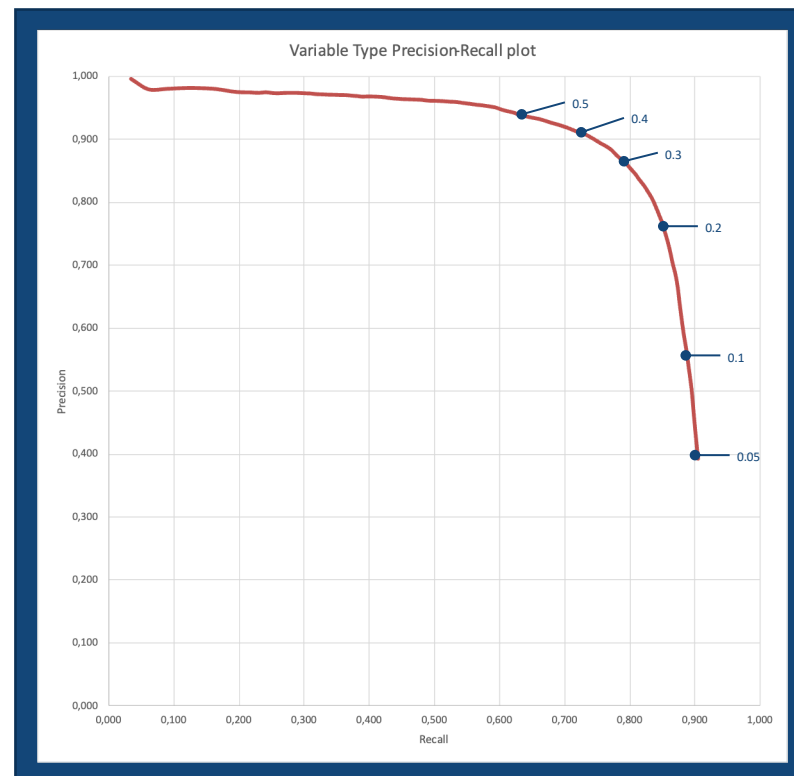- User approved

# Refine steps: Variable models

# STEP 1 – Provide predictions above a static threshold

- Present if likelihood > threshold
  - Up to a maximum of 3 predictions

- Variable-level at 0.3 threshold (61 trials):
  - Precision: 86.9%
  - Recall of 78.6%

- It works to only suggest what you are sure about!

# STEP 2 - Dynamic predictions based on cumulative confidence thresholds
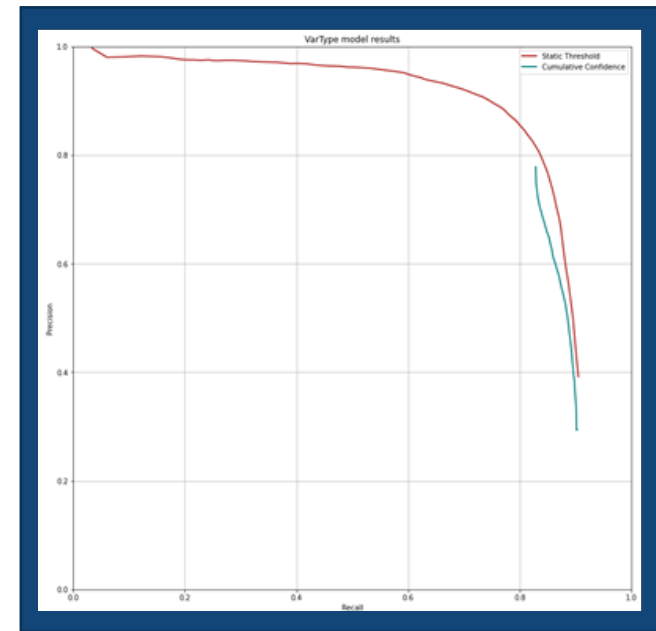
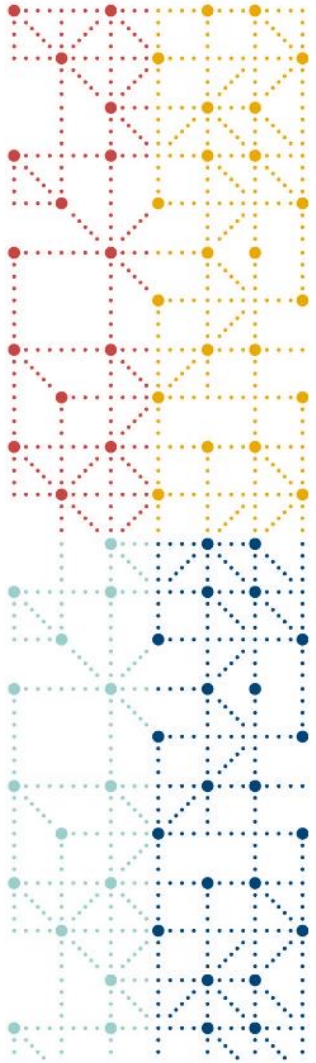[Top 1] > 0.7 → 1 Prediction

[Top 1+2] > 0.7 → 2 Predictions

[Top 1+2+3] > 0.3 → 3 Predictions
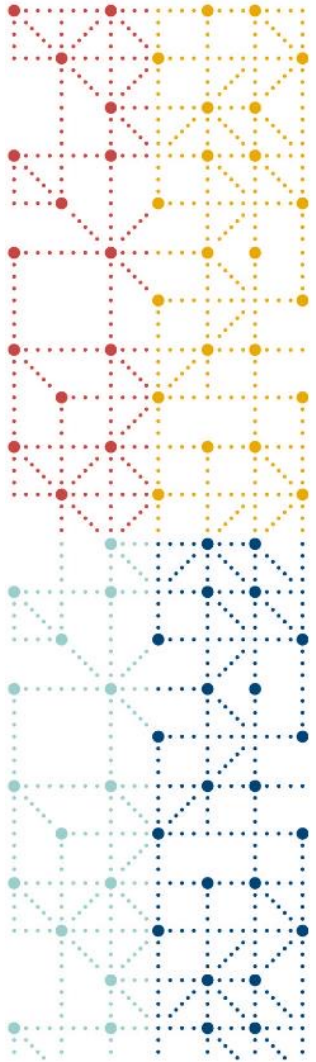
Be practical when implementing ML models!

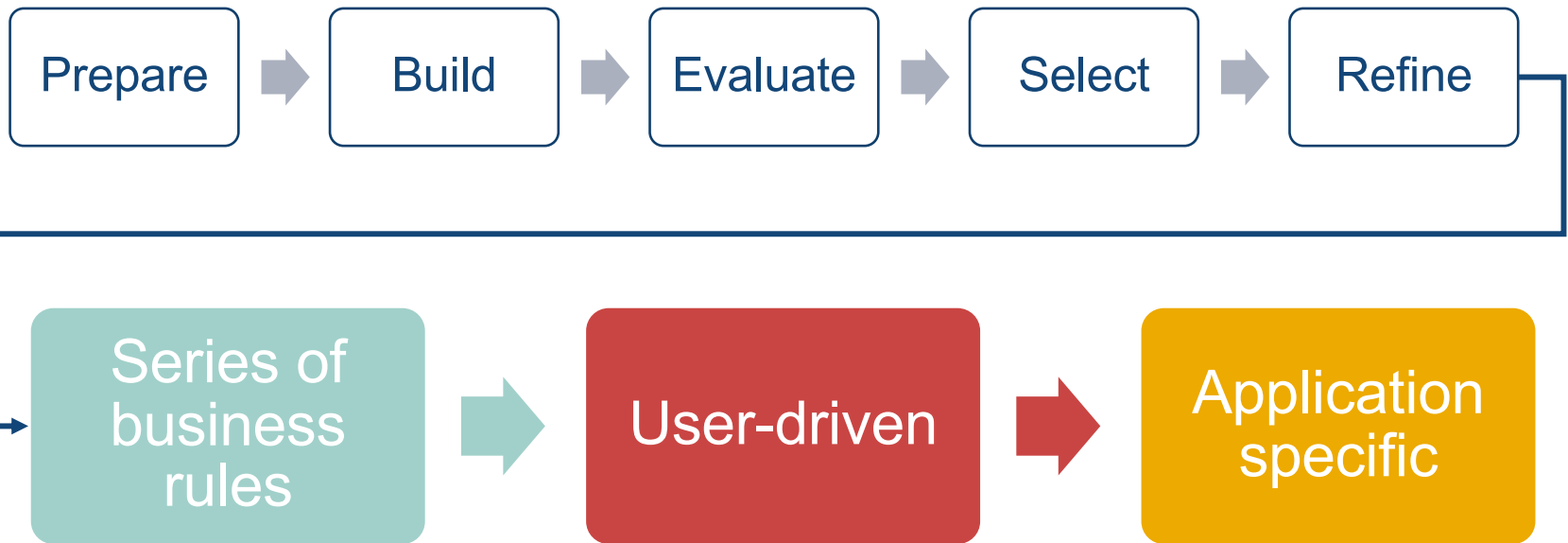# The view from our refined destination
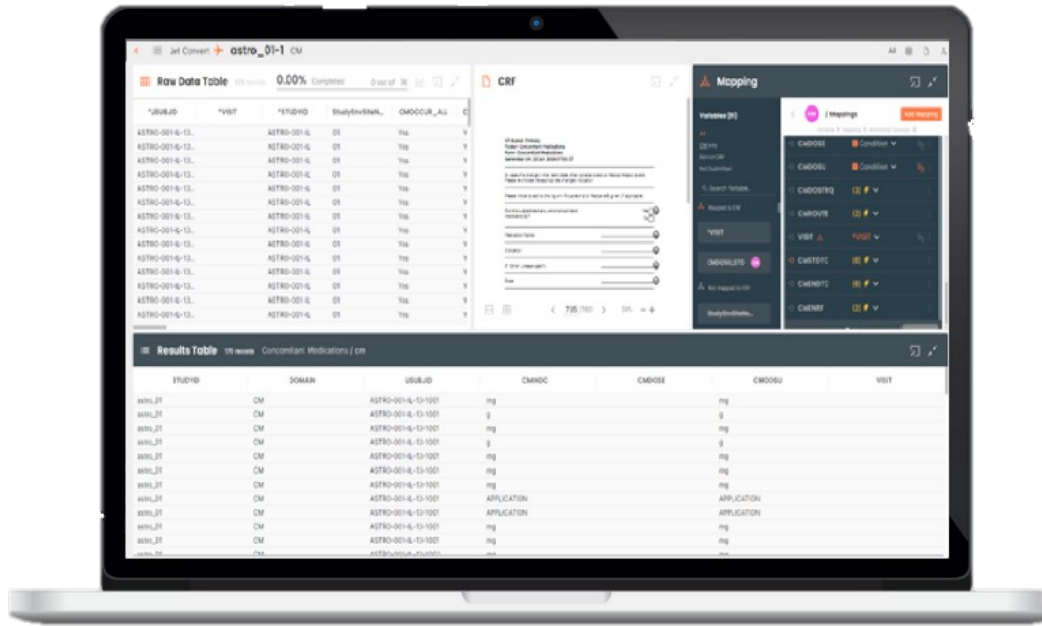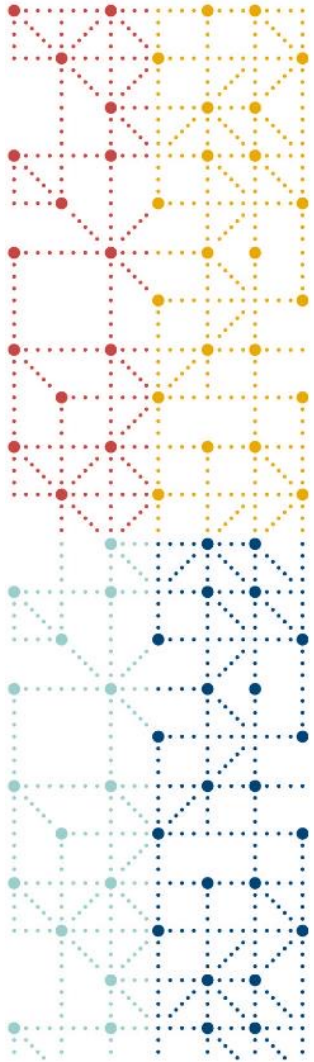
# Where did the journey take us?

# To summarize…

# Practical steps for implementing ML for SDTM Mapping

Prepare → Build → Evaluate → Select → Refine

Series of business rules → User-driven → Application specific

# Using ML in the SDTM mapping process is a reality

# Thank You!



I would like to extend my thanks to the following individuals for assisting in the preparation of this presentation:

- Sergei Merson
- Shahar Cohen
- Lena Hazanov
- Mor Meyerovich
- Eyal Wultz
- Bremer Louw

cdisc