# Virtual AskQC Office Hours

Data and algorithms and bibs, oh my!

OCLC Metadata Quality

April 2023

OCLC

# Housekeeping

This session is being recorded

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Housekeeping

This session is being recorded

All session recordings, slides, and notes are available at oc.lc/askqc

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

# Housekeeping

This session is being recorded

All session recordings, slides, and notes are available at oc.lc/askqc

**Enter questions in chat to "Everyone" at any time during the presentation**

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC
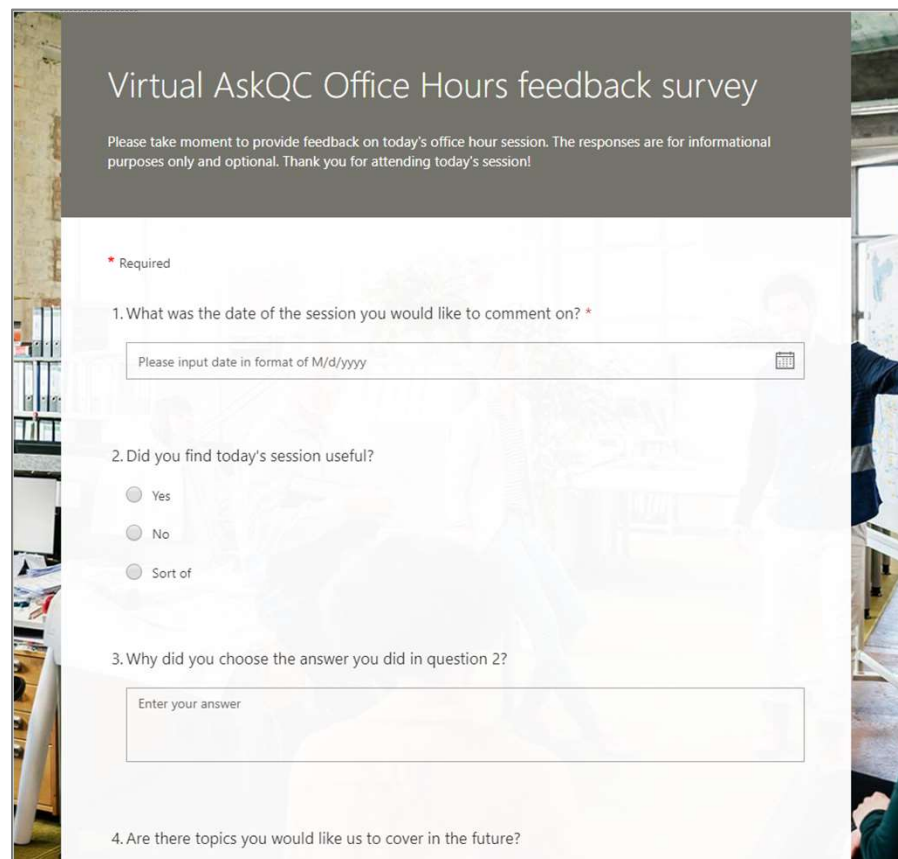
# Housekeeping

This session is being recorded

All session recordings, slides, and notes are available at oc.lc/askqc

Enter questions in chat to "Everyone" at any time during the presentation.

**After the session, you will be directed to a brief, optional survey**

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# On the call today

**Hayley Moreno**
Senior Data Analyst

**Laura Ramsey**
Senior Metadata Operations Manager

**Nathan Putnam**
Director, Data Quality & Governance

**Jenny Toves**
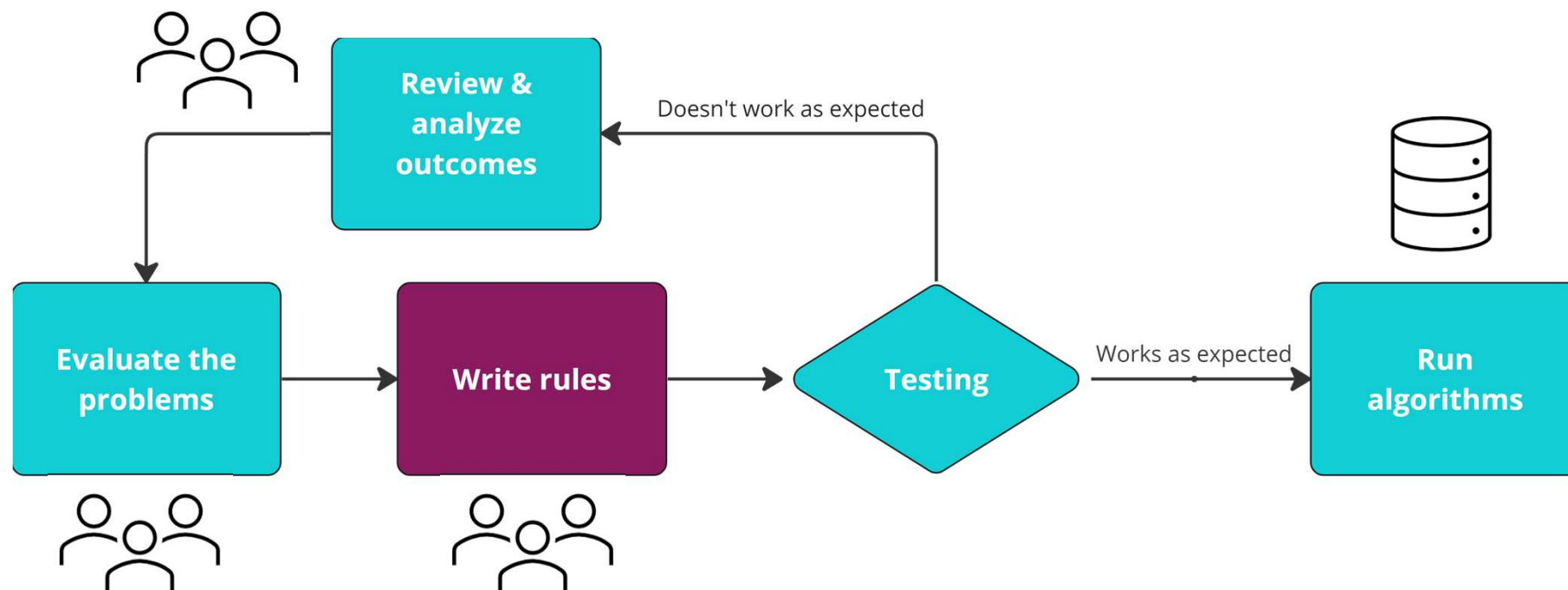Senior Technical Manager

**Becky Dean**
Lead Data Analyst

**Robin Six**
Senior Data Analyst

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Data and algorithms and bibs, oh my!

**Laura Ramsey**
Senior Metadata
Operations Manager

**Nathan Putnam**
Director,
Data Quality & Governance

OCLC

# Rules-based algorithms



Review & analyze outcomes

Doesn't work as expected

Evaluate the problems

Write rules

Testing

Works as expected

Run algorithms

Adapted from: Géron, Aurélien. Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow. Second edition. O'Reilly Media. Kindle Edition.

# Rules-based algorithms



Candidate duplicate records

# Rules-based algorithms

**Rule 1: if X then Y**

**Rule 2: if A then B**

**Rule 3: if M then N**

**Rule …: if … then …**

**Rule n: if $n_1$ then $n_2$**

Lorem ipsum dolor sit amet, consectetur adipisci elit, sed eiusmod tempor incidunt ut labore et dolore magna aliqua. Ut enim ad minim veniam, quis nostrum exercitationem ullam corporis suscipit laboriosam, nisi ut aliquid ex ea commodi consequatur. Quis aute iure reprehenderit in voluptate velit esse cillum dolore eu fugiat nulla pariatur. Excepteur sint obcaecat cupiditat non proident, sunt in culpa qui officia

Candidate duplicate records ➡ Matching & merging rules

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Rules-based algorithms



Candidate duplicate records ➡ Matching & merging rules ➡ **De-duplicated record**

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Where does this leave us?

**Duplicate Detection & Resolution**

- 1 million duplicates removed

**Member Merge Program**

- 20,000 duplicates removed

**OCLC Quality Staff Activities**

- 12,000 duplicates removed

Over **1 million duplicate** records removed between July 2022 and March 2023.

Source: OCLC, MQ Accumulated Stats | Data as of March 2023.

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Where does this leave us?

**Duplicate Detection & Resolution**

- 1 million duplicates removed

**Member Merge Program**

- 20,000 duplicates removed

**OCLC Quality Staff Activities**

- 12,000 duplicates removed

Over **1 million duplicate** records removed between July 2022 and March 2023.

*…but we all know there are more out there.*

Source: OCLC, MQ Accumulated Stats | Data as of March 2023.

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Machine Learning

The science (and art) of programming computers to learn from data



Review & analyze outcomes

Doesn't work as expected

Evaluate the problems

Write rules

Testing

Works as expected

Run algorithms

Data

Adapted from: Géron, Aurélien. Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow. Second edition. O'Reilly Media. Kindle Edition.

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Machine Learning



Training set

OCLC

# Machine Learning



Training set       "Learns" differences

OCLC

# Machine Learning



Training set        "Learns" differences        New data

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Machine Learning



Training set     "Learns" differences     New data     Results

OCLC

# Machine Learning projects

- OCLC Data Labelling Project

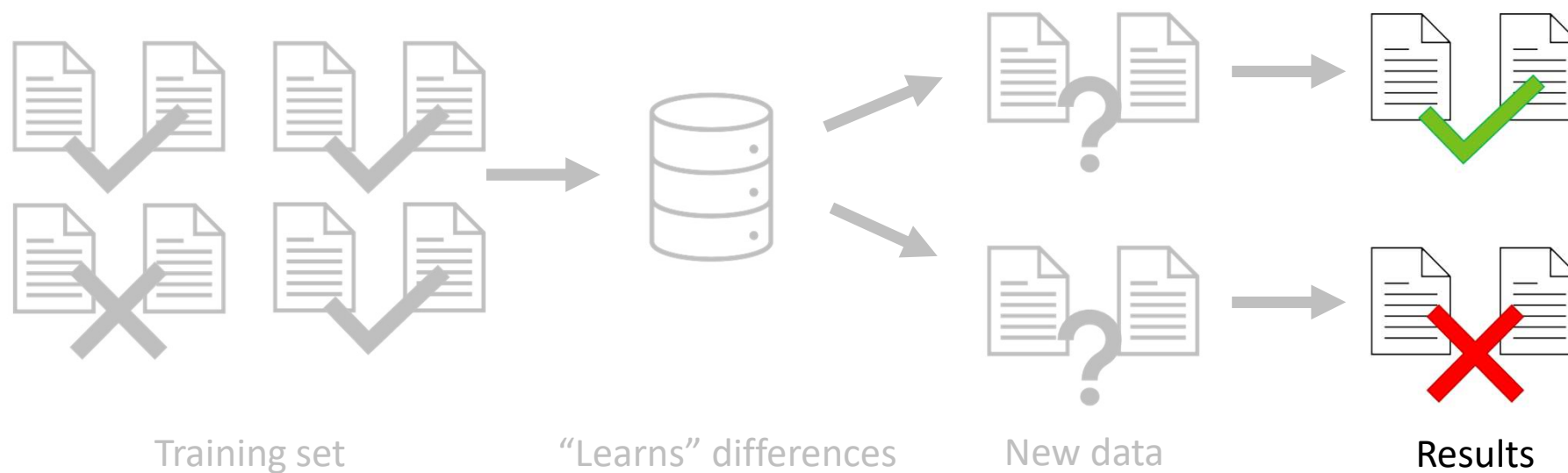  *Objective*: *Tap into the expertise of the cataloging community to validate and enhance our model's understanding of duplicate records, which ultimately improves the quality of WorldCat for the entire cooperative and library community*

- A sample of 60,000 sets for the tool was generated using MinHash

- The labelled data should be representative of WorldCat in terms of language of item, language of cataloging, material type, etc.

- For more information about  MinHash: https://en.wikipedia.org/wiki/MinHash

OCLC

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

# Labelling project stats

## Comparison outcomes

- 40% → the pair IS a duplicate
- 60% → the pair IS NOT a duplicate

## Top 5 fields used for decisions

- 245 – Title statement
- 300 – Physical description
- 260 – Publication, etc. statement
- 100 – Main entry - personal name
- 264 – Publication, etc. statement

**Daily contribution counts**



Source: OCLC, Data Science| Data as of April 2023.

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Labelling project stats

## Language of cataloging



Source: OCLC, Data Science| Data as of April 2023.

## Item format

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Progress so far

- Group 1—Exact duplicates
  - Primary match keys are normalized and compared


- Group 2—Vendor duplicates
  - Records from identified vendors with looser criteria


- Group 3—Model duplicates
  - Labelled data used to train model

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Considerations

- Like rule-based instructions, ML isn't a catch-all solution
  - Models provide a degree of accuracy
- Accepting the degree of uncertainty
  - Is cleaning up millions of duplicate records worth a thousand or so incorrect merges?
- Shift in thinking
  - A rule for every situation is not possible, especially complex data sets
  - Use ML to fill in gaps
  - Experts still involved in training set and evaluation

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# On the call today

**Hayley Moreno**
Senior Data Analyst

**Laura Ramsey**
Senior Metadata
Operations Manager

**Nathan Putnam**
Director, Data Quality &
Governance

**Jenny Toves**
Senior Technical Manager

**Becky Dean**
Lead Data Analyst

**Robin Six**
Senior Data Analyst

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC

# Thank you!

**May Virtual AskQC Office Hours**
What's in a name? Descriptive access points overview

---

Tuesday, 9 May at 10:00 AM Eastern
Thursday, 18 May at 4:00 PM Eastern

---

**Registration and session links available at oc.lc/askqc**

Send cataloging policy questions at any time to askqc@oclc.org



Photo by Eric Rothermel on Unsplash

Virtual AskQC Office Hours: Data and algorithms and bibs, oh my!

OCLC