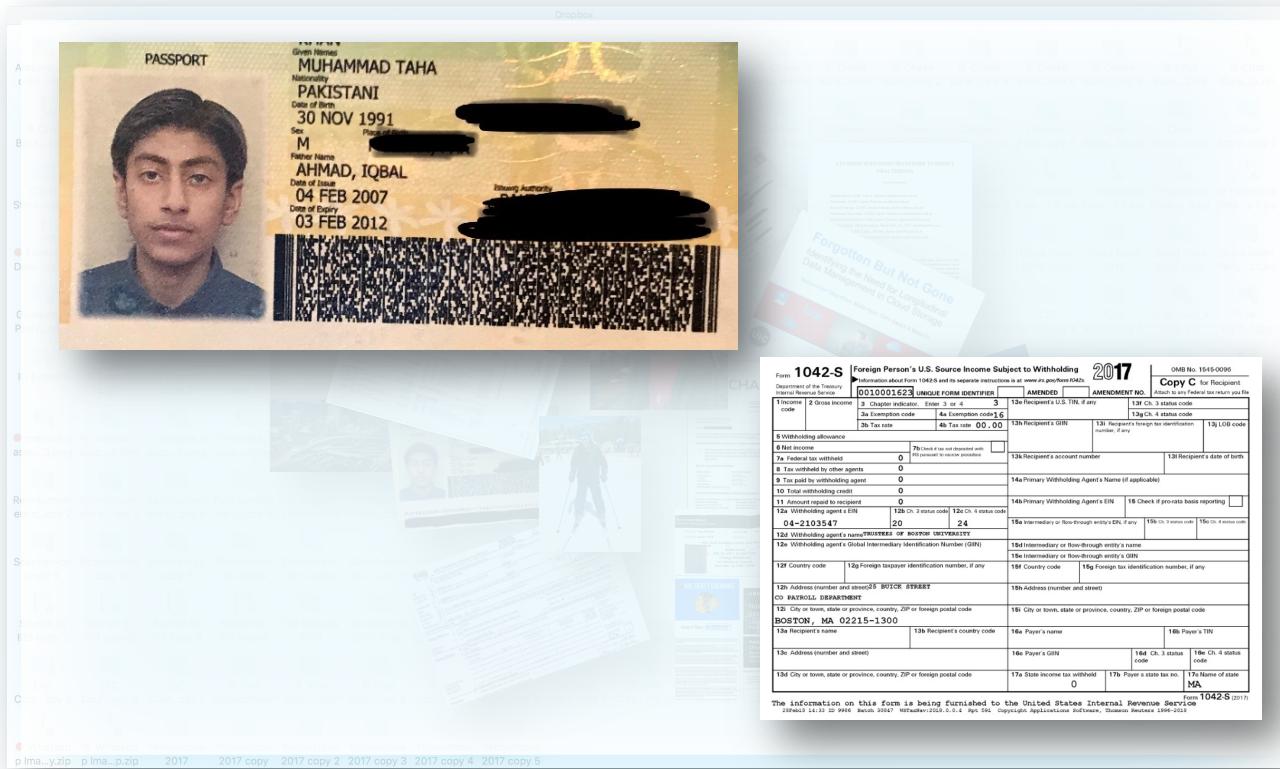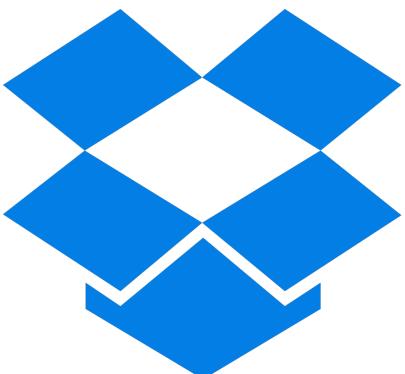# Helping Users Automatically Find and Manage Sensitive, Expendable Files in Cloud Storage

**Mohammad Taha Khan,** Christopher Tran,
Shubham Singh, Dimitri Vasilkov,
Chris Kanich, Blase Ur, Elena Zheleva

# My Personal Dropbox



Taha's Dropbox
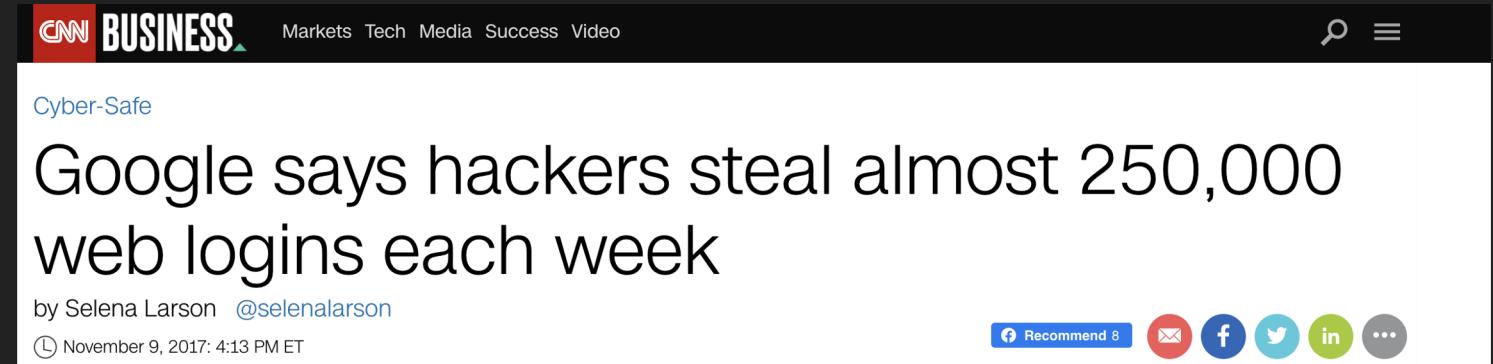User since 2009
8.7GB, 19,245 Files

# The Risk of Breaches



**INDEPENDENT**  News  Voices  Sports

INDY/TECH

## DROPBOX HACK: CLOUD STORAGE COMPANY HACKED, POTENTIALLY REVEALING OVER 60 MILLION PASSWORDS

**CNN BUSINESS**  Markets  Tech  Media  Success  Video

Cyber-Safe

## Google says hackers steal almost 250,000 web logins each week

by Selena Larson  @selenalarson

November 9, 2017: 4:13 PM ET

Recommend 8

Home » Security Bloggers Network » How Social Engineering Tactics Can Crack Multi-factor Authentication

## How Social Engineering Tactics Can Crack Multi-factor Authentication

by Enzoic on April 6, 2021

## Financial Services Experienced 125 Percent Surge in Exposure to Mobile Phishing Attacks in 2020
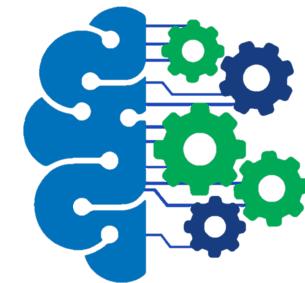
New Lookout Report Finds Increased Attempts to Steal Your Corporate Login Credentials

3

# Goal: Semi-Automated Cloud Management

Machine learning can help users manage large, long-lived archives

Challenges:
- No existing datasets
- File management is subjective and personal
- What features are predictive?

**Contribution:** User-centered design of semi-automated classifiers

# Research Goals and Approach

① ② ③

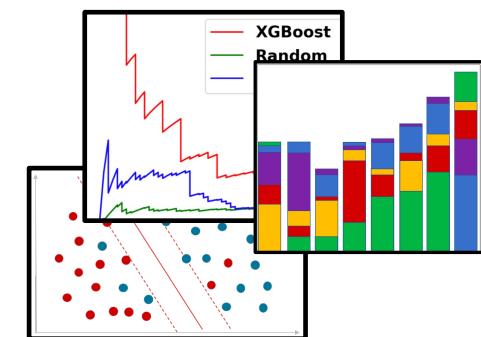**Goals:** Identify relevant file characteristics | Collect features and labels | Develop semi-automated management tools

**Approach:**



**Qualitative Interviews** | **Data-Collection Survey Study** | **Classifier Design and Evaluation**
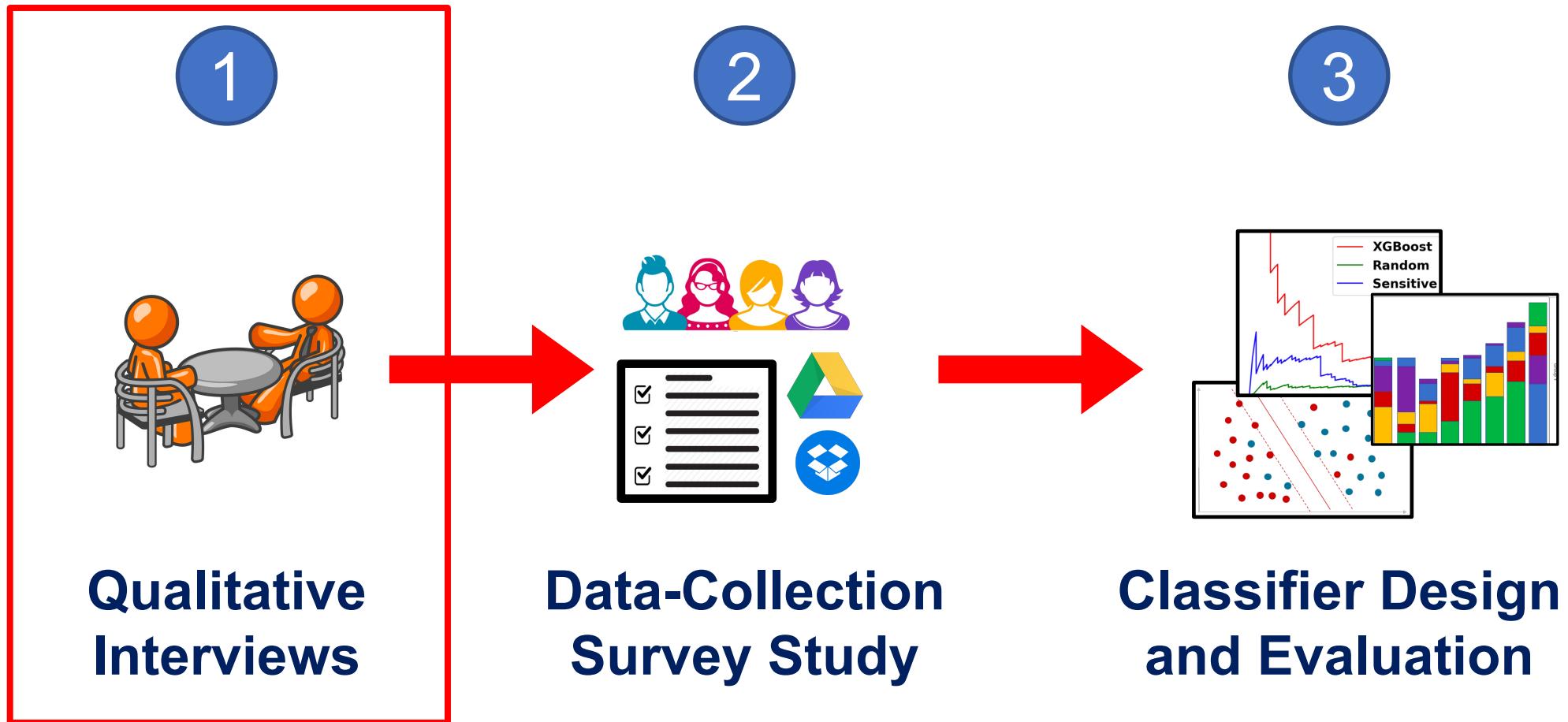
# Sensitivity and Usefulness

Devised metrics of **sensitivity** and **usefulness** to capture subjectivity of file management

Hypothesized file management based on **sensitivity** and **usefulness**

| Sensitivity | Usefulness | Management |
|:---:|:---:|:---:|
| ✗ | ✗ | Delete 🗑 |
| ✓ | ✗ | |
| ✗ | ✓ | Keep as-is |
| ✓ | ✓ | Protect 🔒 |

# Approach



**1** Qualitative Interviews

**2** Data-Collection Survey Study

**3** Classifier Design and Evaluation

# Qualitative Interviews

**Goal:** Understand subjective opinions of file sensitivity and usefulness

Interviewed 17 participants from diverse backgrounds

Explored mental-models of participants

# Characteristics of Potentially *Sensitive* Files

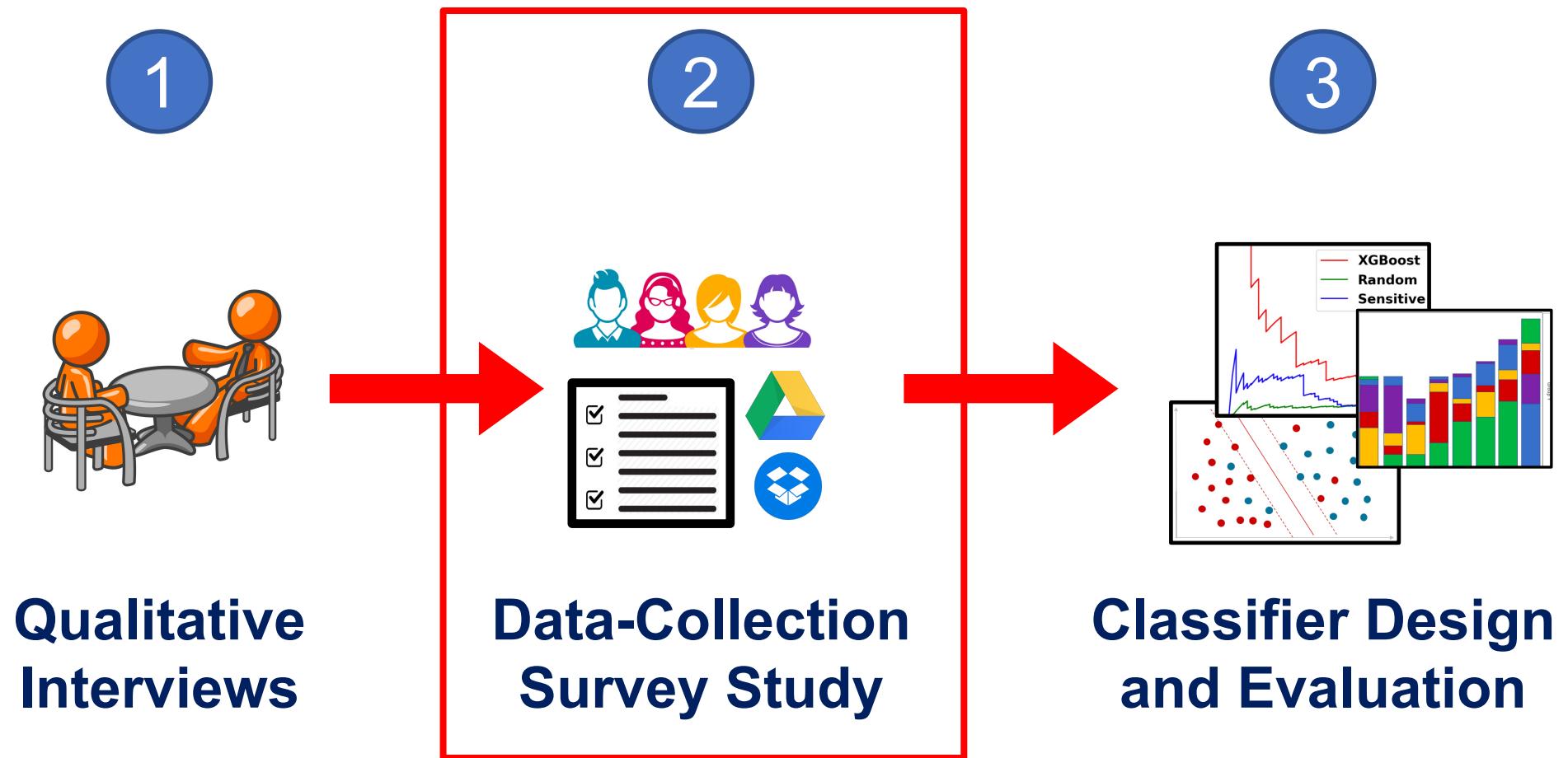| |
|---|
| Personally identifiable or financial content |
| Nude, intimate, or embarrassing content |
| Content concerned with self-presentation |
| Proprietary and confidential information |

# Characteristics of Potentially *Useful* Files

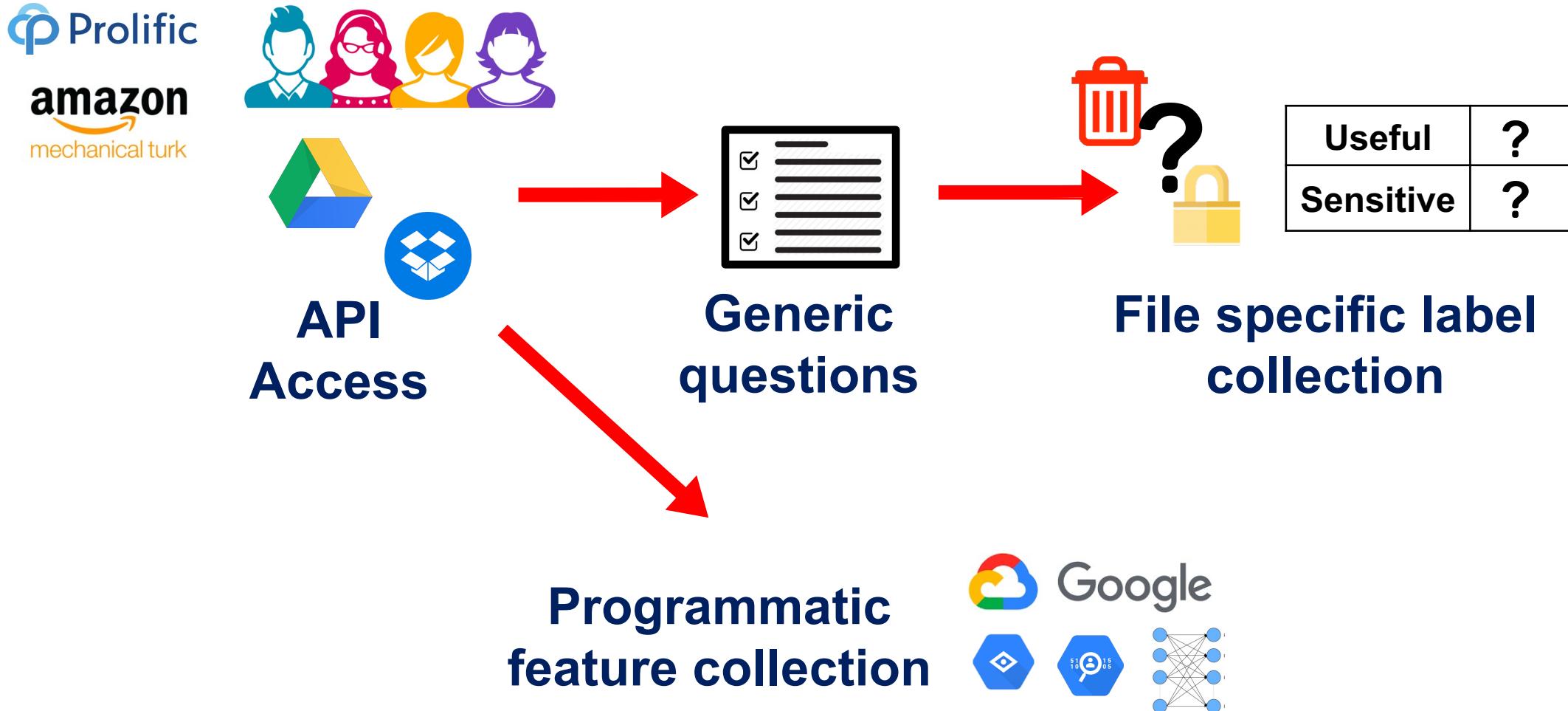| |
|---|
| Files for future reference |
| Files regularly accessed or shared |
| Memories and files with sentimental value |
| Backup archives |

# Approach



**1** Qualitative Interviews

**2** Data-Collection Survey Study

**3** Classifier Design and Evaluation

# Data-Collection Framework



API Access

Generic questions

File specific label collection
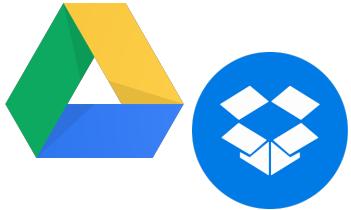
| Useful | ? |
|---|---|
| Sensitive | ? |

Programmatic feature collection

# Features Collected

## Dropbox and GDrive API

- Account age
- File name
- File size
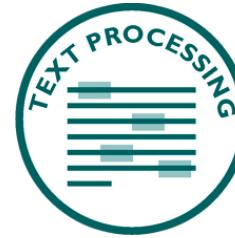- **Access details**
- **Sharing status**
- .
- .
- .

## Google Vision

- Image objects
- **Adult**
- **Racy**
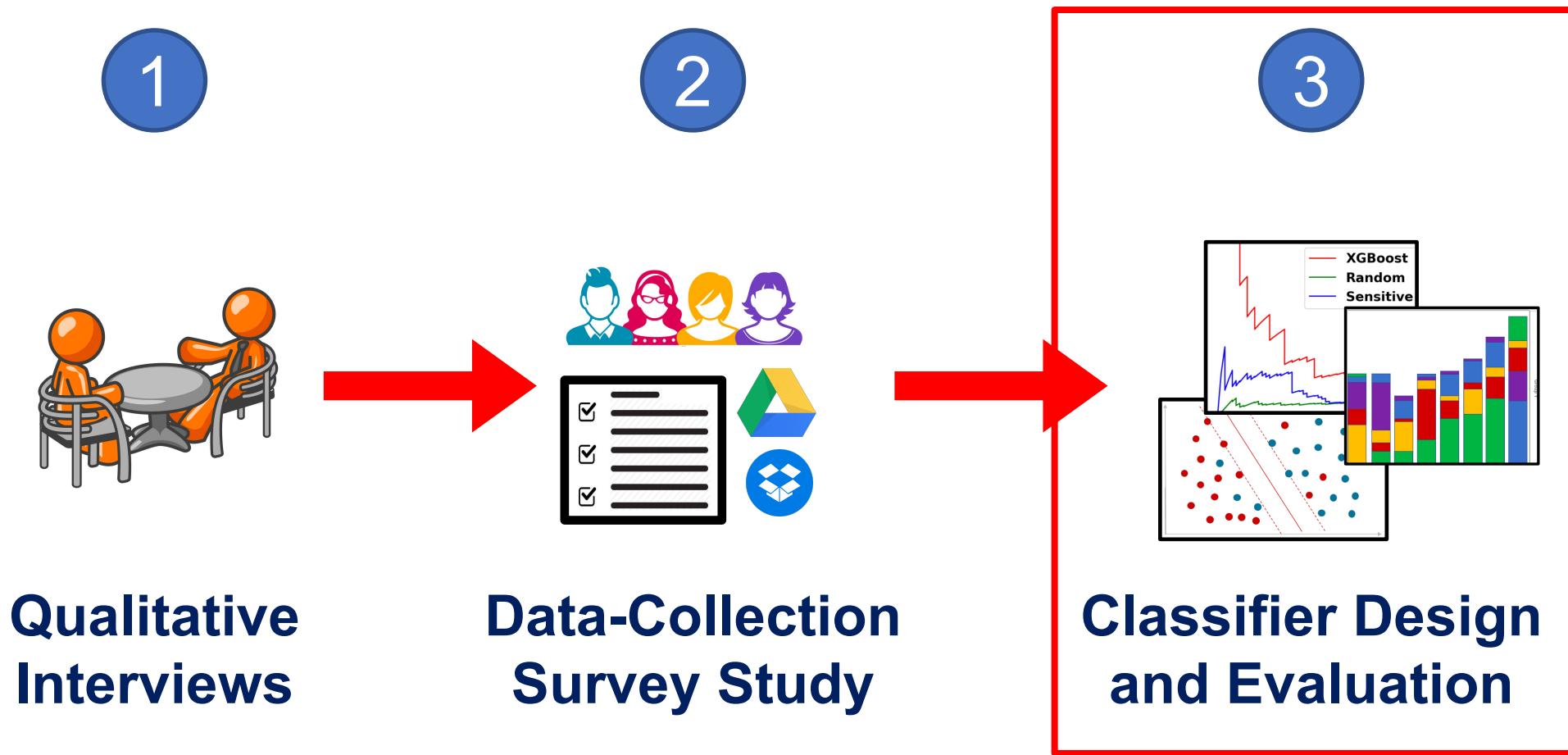- **Violent**
- Spoof
- .
- .

## Google DLP

- **Name**
- **SSN,**
- **Email**
- **License #**
- **Credit card**
- **Bank Info**
- .
- .
- .

## Local text processing

- Doc topics
- Bag of words
- Word2vec
- .
- .
- .

# Approach



**1** Qualitative Interviews

**2** Data-Collection Survey Study

**3** Classifier Design and Evaluation

XGBoost
Random
Sensitive

# Classifier Design and Evaluation

**Goal:** Partially automate file management via machine learning

| Classifier | Prediction Class |
|---|---|
| Sensitivity | Sensitive, Not Sensitive |
| Usefulness | Useful, Not Useful |
| Management | Keep as is, Delete, Protect |

# Performance of the *Sensitivity* Classifier



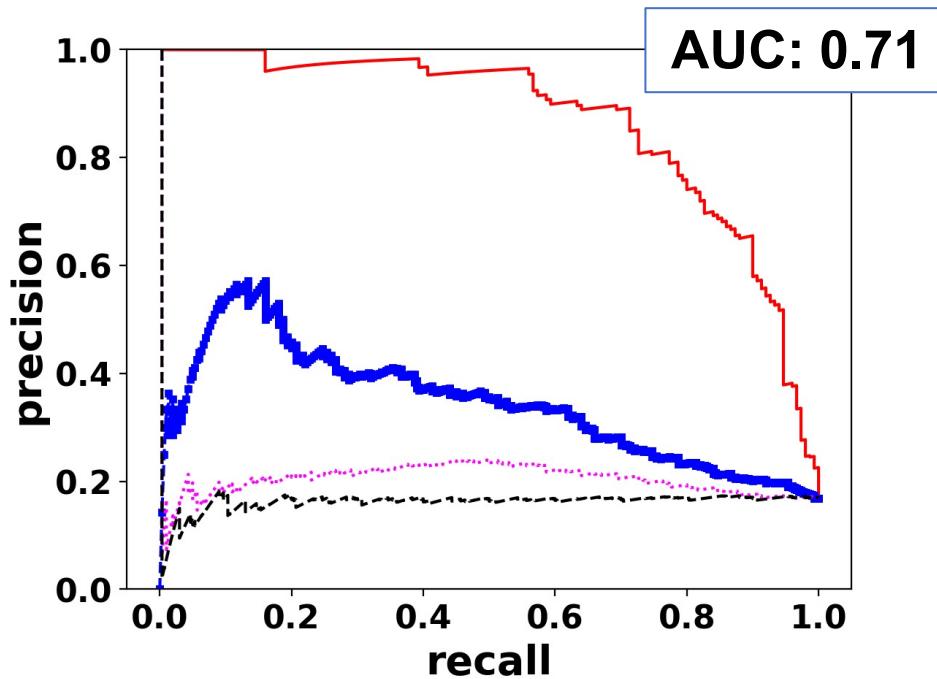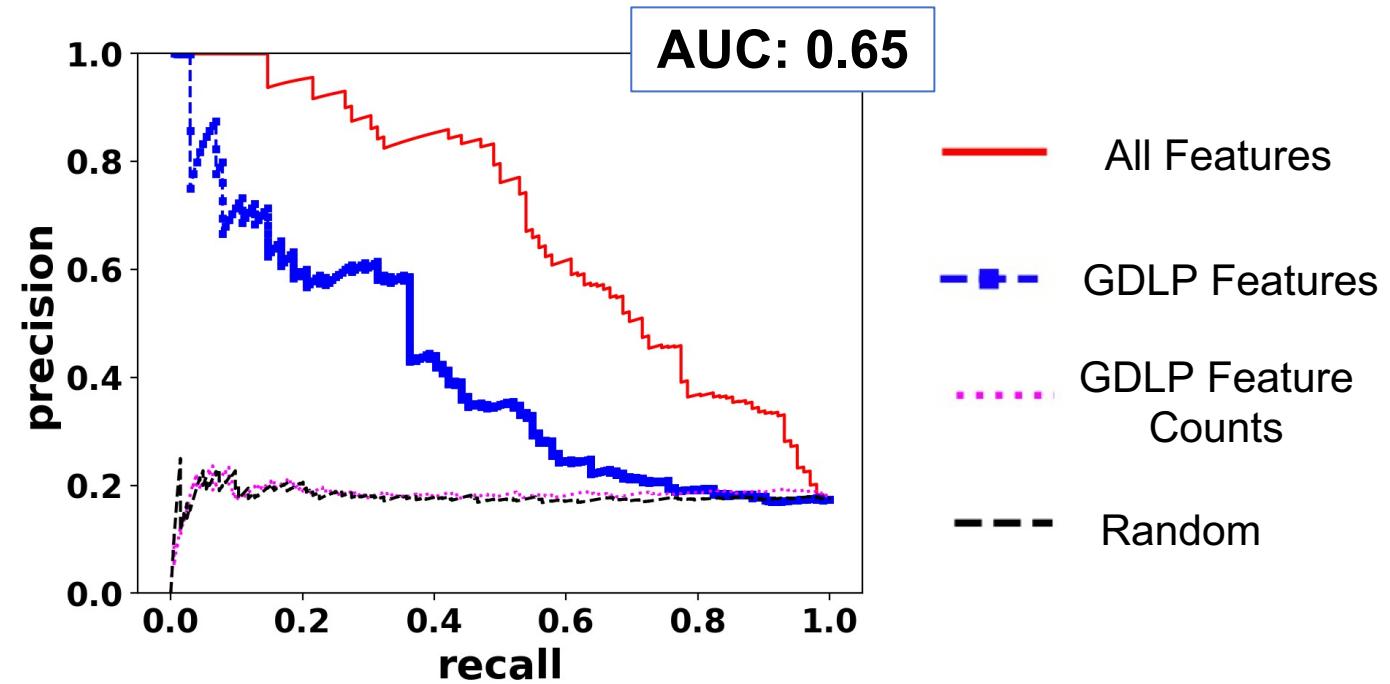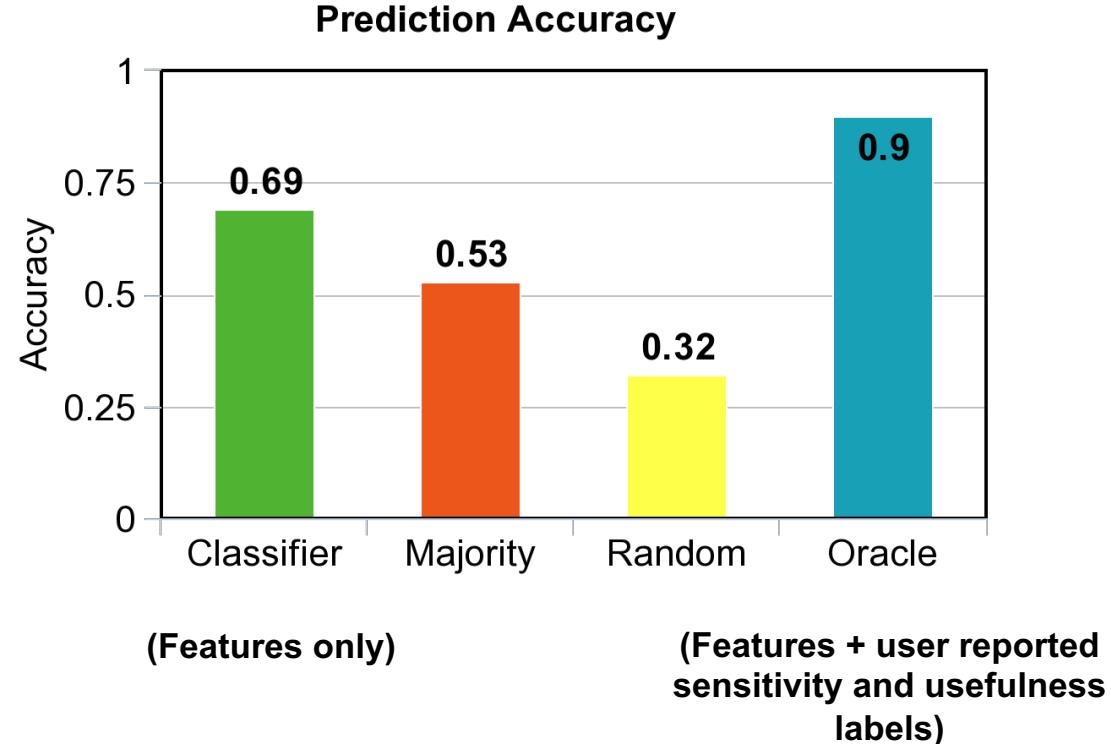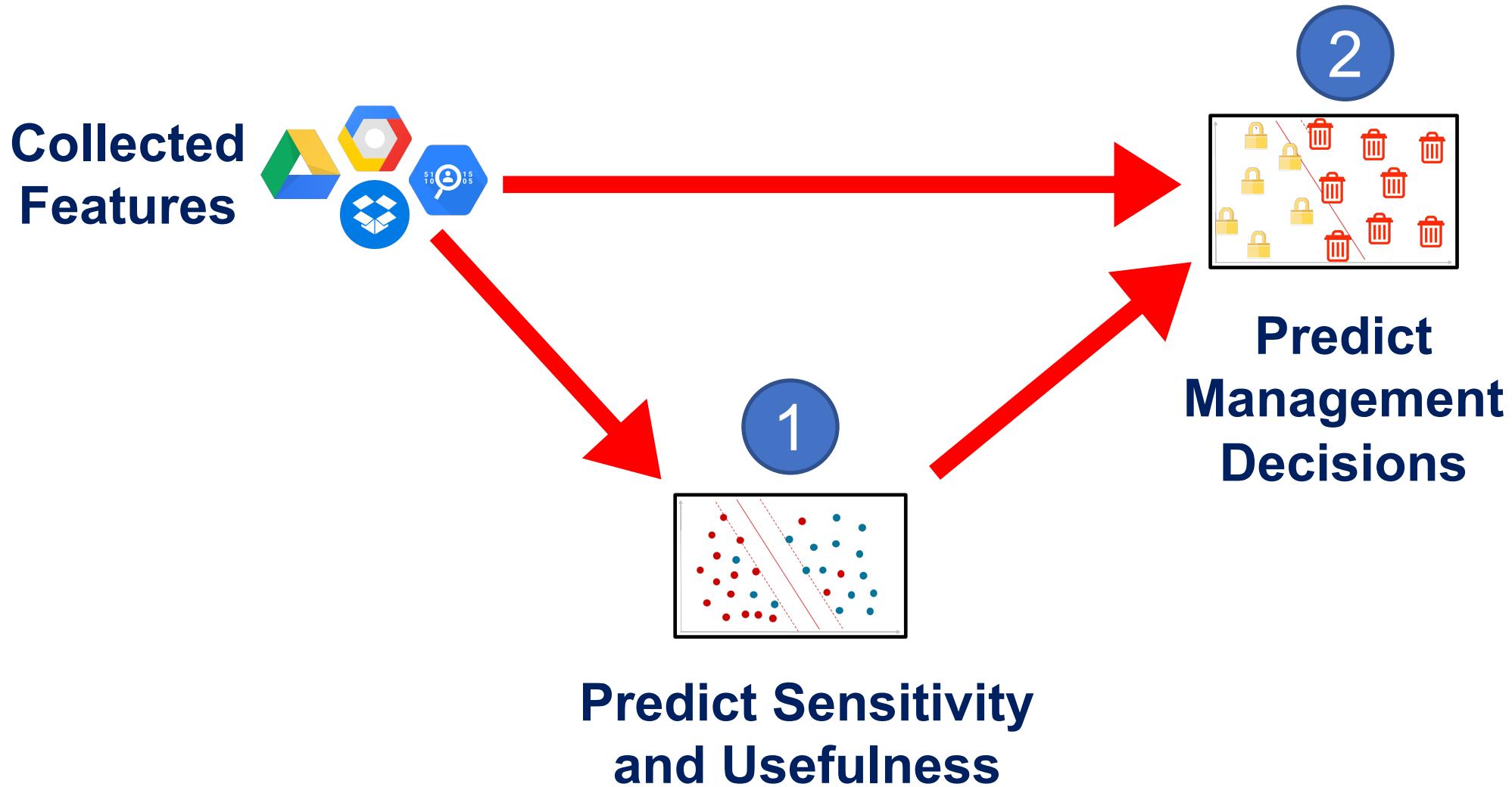Image Sensitivity Classifier

Document Sensitivity Classifier

# Accuracy of the *Management* Classifier

**Prediction Accuracy**



**(Features only)**

**(Features + user reported sensitivity and usefulness labels)**

69% management prediction accuracy by just using collected features

# Two-Step Classification

**Collected Features**

**1** Predict Sensitivity and Usefulness

**2** Predict Management Decisions

# Increased Accuracy of the *Management* Classifier

**Prediction Accuracy**



10% increase in accuracy when training includes sensitivity and usefulness

# Conclusion

Qualitative interviews → characteristics of sensitive and useful files

Quantitative user study → labeled dataset

Predicting file sensitivity and usefulness helps predict file management

Our prototype web app: https://cloudsweeper.app

# Conclusion

- Qualitative interviews → characteristics of sensitive and useful files

- Quantitative user study → labeled dataset

- Predicting file sensitivity and usefulness helps predict file management

- Our prototype web app: https://cloudsweeper.app

**Helping Users Automatically Find and Manage
Sensitive, Expendable Files in Cloud Storage**

**Mohammad Taha Khan** (tkhan@wlu.edu)**,** Christopher Tran, Shubham Singh, Dimitri Vasilkov,
Chris Kanich, Blase Ur, Elena Zheleva