# Datathon and Machine Learning Competition on Antisemitism Workshop 2 – Creating a Discourse Dataset from X

Presentation available at:
https://github.com/damieh1/datathon_2025/blob/main/Datathon_Workshop-Session-2.pdf

# What happened so far?

**Workshop 1 Recap:**

- Prof. Jikeli introduced patterns of antisemitism online ▪ ▪ ▪ ➤

- Focus on dynamics before and after October 7



**Conceptual and Theoretical Background:**

- Reach out to him with questions about content/context ▪ ▪ ▪ ➤ *Missed it? His slides are available in the GitHub repo*

**Team formation:**

- Most groups should be in place by now

- Team Communication & Responsibilities

# Today's Session

1. Bright Data Introduction

2. Hands-on Session for Challenge #1

   • Work in your teams to collect, annotate, and prepare data

3. Q&A → Conceptual: Prof. Jikeli | Content: Dr. Miehling | ML: Prof. Cavar

INDIANA UNIVERSITY

# Meet the Instructors



Workshop #1

**Prof Günther Jikeli**

[gjikeli@iu.edu](mailto:gjikeli@iu.edu)



Workshop #2

**Dr Daniel Miehling**

[damieh@iu.edu](mailto:damieh@iu.edu)



Workshop #3

**Prof Damir Cavar**

[dcavar@iu.edu](mailto:dcavar@iu.edu)

# Workshop Schedule Overview

| Date | Focus | | Description |
|------|-------|---|-------------|
| **July 13** | Workshop 1 – Kickoff | ✅ | Input on antisemitism online & team formation |
| **July 20** | Workshop 2 – Challenge #1 | | Scraping, annotating, exporting discourse data (Hands-On Session) |
| **July 27** | Workshop 3 – Challenge #2 | | Automated Content Detection: The Basics |
| **August 5** | Final Submission Deadline | | Submit both challenge deliverables and documentation |

# Bright Data



https://www.youtube.com/watch?v=AGaiVApKfmc

# Tutorial Challenge #1

**Team Setup, Roles & Prerequisites**

Each team should have:

- **Data Manager** → uses Bright Data for scraping
- **Portal Manager** → sets up the annotation project
- **Annotation Team** → annotates and reviews the content
- Discuss roles & workflow internally
- Don't forget to choose a **team name or tag**!

**Access to Key Links**

- Challenge Description (PDF)
- Annotation Portal
- GitHub Overview
- Colab Script for Preprocessing

INDIANA UNIVERSITY

# Agenda & Objectives

**Agenda**

1. Working with X's Advanced Search

2. Scraping with Bright Data

3. Pre-Processing the Data

4. Annotations Portal Walkthrough

**Objectives: What You'll Learn Today:**

- How the Annotation Portal works

- How to approach the ML challenge

- Where to find tools and datasets

- How to succeed as a team

# Prerequisites & Setup

Before You Start:

A computer with internet access

A X account

A Gmail account

Access to Google Colab (https://colab.research.google.com/)

Stop & Do Now:

Register on the Annotation Portal: https://annotate.osome.iu.edu

Check Github: https://github.com/AnnotationPortal/DatathonandHackathon.github.io/blob/main/README.md

Read Challenge Description: https://github.com/damieh1/datathon_2025/blob/main/Datathon_Challenge.pdf

INDIANA UNIVERSITY

# Challenge #1

📍 **Subtasks include:**

1. Define your scraping focus (hashtags, user groups, topics) and document your rationale and potential biases.

2. Use the Bright Data interface to scrape at least 100 relevant user-generated posts from X.com.

3. Annotate your data using a structured definition of antisemitism and hate speech.

4. Prepare a X/Twitter dataset, and include a dataset report with label definitions, distribution information, and annotation rationale.

# Earn Bonus Points

📍 **Deliverables for Challenge #1**
- Adapting and implementing an existing definition of antisemitism
- Reporting how the data was scraped and which guidelines were used to classify and annotate the data in a standardized way

Gain **+10 bonus points** by evaluating the consistency of your team's annotations using an inter-annotator agreement (IAA) metric.
This means:
- Having at least two annotators label a shared subset of the data
- Calculating a formal agreement score, such as:
  - **Cohen's Kappa** (for binary or pairwise categorical annotation)
  - **Krippendorff's Alpha** (especially for multi-class or missing data)
Clearly report:
- Which subset was double-annotated
- Your score and a brief interpretation (e.g., "moderate agreement," "high agreement")

# Working With X's Advanced Search Function

Top of the pop-up menu

Bottom of the pop-up menu



1. Go to: → X (Twitter) → https://x.com/search-advanced

2. Specify dates, e.g., May 8, 2024.

3. Click "Search."

4. Select posts with a minimum of 200 views.

5. Go to user profiles and copy URLs to a spreadsheet.

6. Goal: Find a wide range of users who engage in online discourse.

# Working With Bright Data



1. Go to: Bright Data → Web Scrapes → X (Twitter) → Posts → Discover by URL → https://brightdata.com/cp/scrapers/no_code

2. Click: Add Inputs → https://x.com/RandomXUser

3. Specify Number of Posts → max. 250 per User

4. Start Collecting → Runes the Query

5. Download Output as .CSV

# Working With Google Colab



1. Go to: Google Colab → https://colab.research.google.com

2. Upload: Bright Data Output .CSV

3. Parse Data
   → Run Code on Colab

4. Prepare Data for Annotation Portal
   → Run Code on Colab

5. Download compatible .CSV Output

Click this →  Open in Colab

INDIANA UNIVERSITY
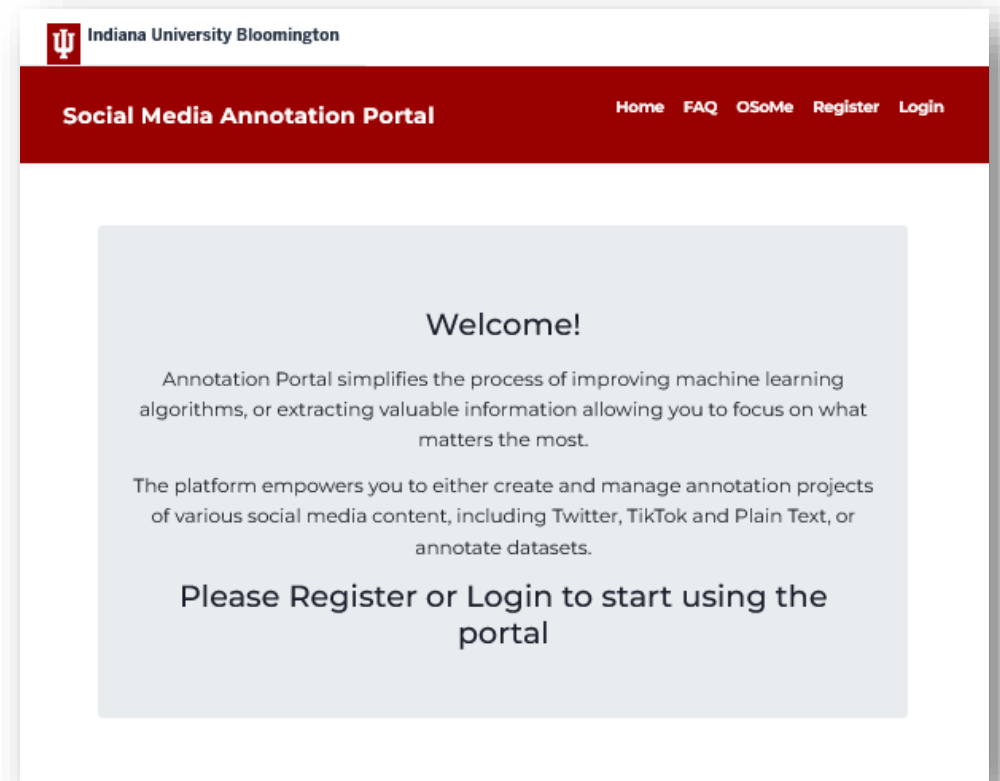
# Annotation Portal Walkthrough

**Use the Annotation Portal:**
1. Create New Project
2. Create a Sample
3. Create Annotation Scheme and Questions
4. Important! Do not start annotating before the schema has been fully created.
5. Export when done



https://annotate.osome.iu.edu/

# What's next?

📝 **Work with your Team on Challenge #1**

📅 **Important Dates:**

- Workshop 1 – July 13: Kick-Off & Introduction (Team Assignment & Communication) ✅
- Workshop 2 – July 20: Hands-On Session (Data collection, preprocessing & annotation) ✅
- Workshop 3 – July 27: Introduction to Automated Detection (ML modeling & evaluation)
- Final submission deadline: August 5.

📁 **All materials available at:**

- [https://github.com/damieh1/datathon_2025](https://github.com/damieh1/datathon_2025)

💬 **Q&A – We'll now open the floor for questions!**

# Thanks for your attention

**Dr. Daniel Miehling**

-----------------------------------------------

Visiting Assistant Professor
Borns Jewish Studies Program

-----------------------------------------------

Samerian Foundation Research Fellow
Institute for the Study of
Contemporary Antisemitism (ISCA)

-----------------------------------------------

Indiana University Bloomington
damieh@iu.edu