



ROBERT A. AND SANDRA S.  
BORNES JEWISH STUDIES PROGRAM

# INSTITUTE FOR THE STUDY OF CONTEMPORARY ANTISEMITISM

## Detecting Antisemitic Hate Speech & Conspiracy Fantasies – From Raw Data to Smart Detection – July 2025

**Overview:** In this datathon, your team will tackle two interconnected challenges that reflect real-world tasks in hate speech research: the creation of a labeled dataset and the application of machine learning techniques to detect hate speech. Your coordination, documentation, and critical reflection on the results will be assessed and evaluated by experts in this field.

### #1 Challenge: Building and Annotating a Dataset

**Goal:** Create a small but meaningful labeled dataset for hate speech detection.

#### Tasks:

- 1. Data Collection:** Use the *Bright Data*<sup>1</sup> interface to scrape a minimum of at least 100 relevant posts of user-generated content (e.g., tweets, posts, or comments) from X (formerly known as Twitter).
- 2. Sampling and Documentation:** Decide on your scraping focus (e.g., specific hashtags, keywords, seed list or user groups) and document your strategy. Describe how you targeted the issue, the motivations behind your decision, and why this content is relevant and potentially antisemitic. Keep in mind that your dataset should include antisemitic and non-antisemitic content. Machine Learning Models need labeled datasets that include both.
- 3. Annotation:**
  - Apply a definition of antisemitism to annotate your dataset.
  - Apply a standardized annotation form (either use our scheme as prototype, adjust it or create your own scheme based on the type of content you are interested in).<sup>2</sup>
  - We recommend using our online *Annotation Portal*<sup>3</sup> for labeling.
- 4. Deliverables include:**
  - A cleaned and annotated dataset, exported in .csv format, including only posts containing your selected keyword or hashtag.
  - A short dataset report (1–2 pages) that: Defines the labels used during annotation (e.g., antisemitic/non-antisemitic).

---

<sup>1</sup><https://brightdata.com/products/web-scraper/functions>

<sup>2</sup>[https://github.com/AnnotationPortal/DatathonandHackathon.github.io/blob/main/guides/annotation\\_scheme.md](https://github.com/AnnotationPortal/DatathonandHackathon.github.io/blob/main/guides/annotation_scheme.md); see also: <https://arxiv.org/abs/1910.01214>

<sup>3</sup><https://annotate.osome.iu.edu/>

- Describes the keyword/hashtag and time period selected, and why Summarizes the distribution of labeled examples.
- Explains how annotations were conducted and by whom.
- And Includes any challenges or limitations (e.g., platform constraints, ambiguous posts).<sup>4</sup>

## #2 Challenge: Modeling and Evaluation

**Goal:** Use our pre-annotated gold standard dataset to build and evaluate a system to detect antisemitism.

### Tasks:

1. **Dataset Access:** Download the provided annotated datasets:

*Antisemitism on Twitter: A Dataset for Machine Learning and Text Analytics*<sup>5</sup>

*Antisemitism on X: A Dataset Tracking Trends in Counter-Speech and Israel-Related Discourse Before and After October 7*<sup>6</sup>

2. **Modeling:** Use transformer-based models to build a classification system for detecting antisemitic content. This includes not only explicit hate speech but also coded language, conspiratorial narratives, Holocaust distortion.
  - We recommend fine-tuning a model from Hugging Face Transformers. You may choose from general-purpose or task-specific models:
    - [twitter-roberta-base-offensive]<sup>7</sup> — pre-trained on offensive language from Twitter/X
    - [microsoft/mdeberta-v3-base]<sup>8</sup> — strong multilingual classification baseline
    - [bertweet-base]<sup>9</sup> — designed for Twitter-style texts
    - [GroNLP/hateBERT]<sup>10</sup> — pre-trained on hate speech corpora from Reddit
  - Fine-tune your model using tools like Hugging Face’s **Trainer** API or custom PyTorch code.
  - Use a standard train/test split (e.g., 80/20), and optionally a validation set (e.g., 10% of training) for early stopping or tuning.
  - Use a fixed random seed for reproducibility.
3. **Evaluation:** Your submission should include:
  - Standard classification metrics: **precision**, **recall**, **F1-score**, and a **confusion matrix** (on the test set).

---

<sup>4</sup>Our descriptions of the datasets published on Zenodo may be useful as template: <https://zenodo.org/records/14448399>

<sup>5</sup><https://zenodo.org/records/14448399>

<sup>6</sup><https://zenodo.org/records/15025646>

<sup>7</sup><https://huggingface.co/cardiffnlp/twitter-roberta-base-offensive>

<sup>8</sup><https://huggingface.co/microsoft/mdeberta-v3-base>

<sup>9</sup>[https://huggingface.co/docs/transformers/en/model\\_doc/bertweet](https://huggingface.co/docs/transformers/en/model_doc/bertweet)

<sup>10</sup><https://huggingface.co/GroNLP/hateBERT>

- A complete list of **hyperparameters** used (e.g., learning rate, batch size, epochs, weight decay).
- A short **error analysis** describing what types of posts your model struggles with (e.g., satire, rhetorical questions, mixed sentiment or transliterations).
- 3–5 **qualitative examples** of false positives and false negatives, with a brief explanation of why the model may have misclassified them.

**Note:** *While prompt-based classification using large language models (e.g., Chat-GPT, GPT-4) is a popular trend, this task focuses on **training and evaluating reproducible models**.*

### #3 Resources and Structure

- A live tutorial will introduce you to our *Annotation Portal* and *Bright Data*.
- Breakout sessions will follow for Q&A, team coordination, and hands-on work.
- At the end of the datathon competition, each team will submit a report on their workflow and insights, including scripts and documentation of used tools.

### #4 Scoring & Evaluation

Each team can earn a maximum of **100 points**, with an additional **20 bonus points** available for optional advanced tasks.

- **50 points – Task Performance:** Assessed based on the quality of data scraping, annotation consistency, and model performance on the evaluation dataset (e.g., precision, recall, F1-score).
- **50 points – Report and Presentation:** Awarded based on the clarity, structure, and completeness of the submitted report, insights into methodology, annotation rationale, and evidence of teamwork and collaboration.
- **+20 Bonus Points – Optional Advanced Contributions:** Bonus points may be awarded for:
  - Using advanced evaluation metrics such as **Krippendorff’s Alpha** or **Cohen’s Kappa** to report inter-annotator agreement
  - Applying your model to additional **unseen data** beyond the test set
  - Addressing the **social and ethical implications** of antisemitism detection and hate speech classification

The jury includes experts from machine learning, hate speech research, and digital humanities.

**Any Further Questions? *Don't hesitate to reach out!***

For questions or technical support, reach out to the organizing team: Dr. Daniel Miehling (damieh@iu.edu), Prof. Gunther Jikli (gjkeli@iu.edu), and Rachel Kelly (rk18@iu.edu).

**Remember: The competition not only demands curiosity, but also strong teamwork and critical thinking to succeed!**