# War of Words: Using Large Language Models and Retrieval Augmented Generation to Classify, Counter and Diffuse Hate Speech

**Rohan Leekha, Olga Simek and Charlie Dagli - Massachusetts Institute of Technology Lincoln Laboratory**
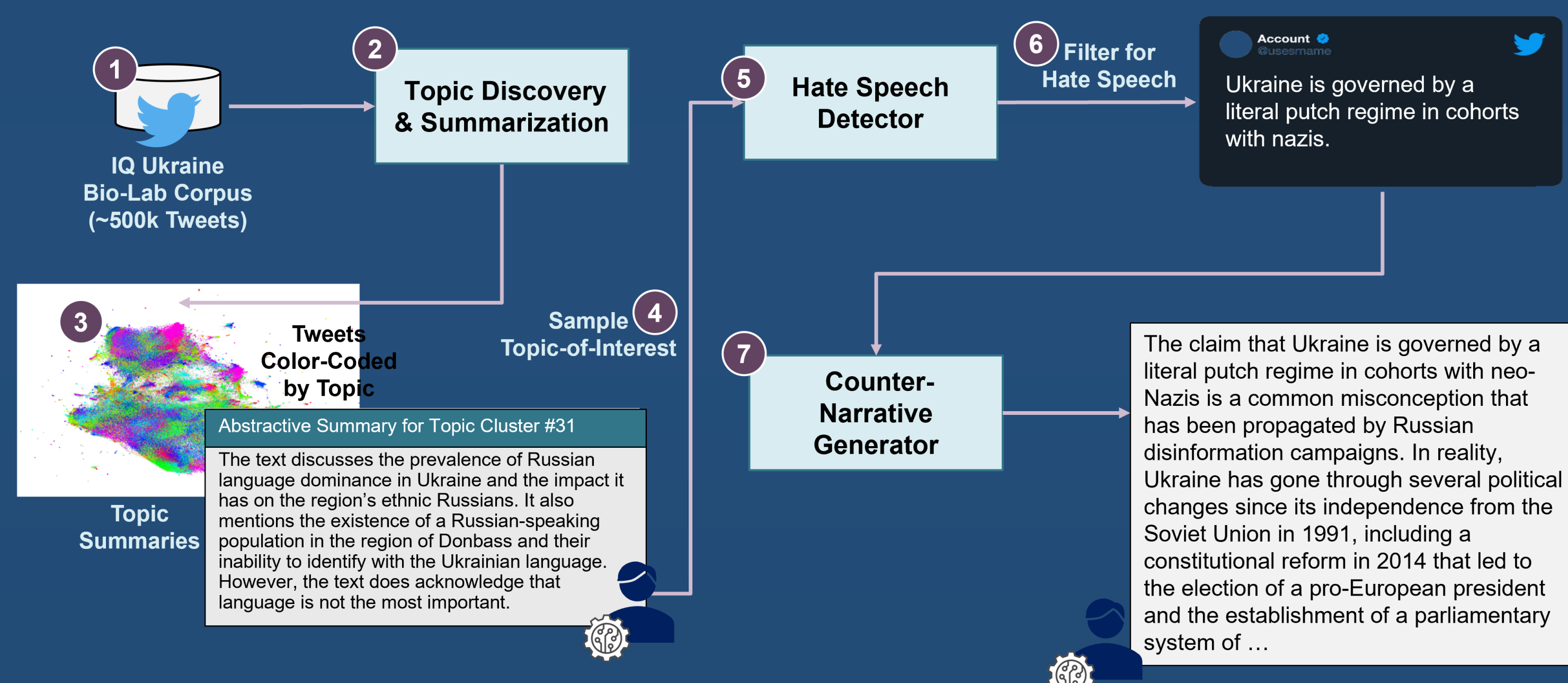
## Introduction

### MOTIVATION

There is a need for a robust and effective defense strategy against malign influence spread by influence operation campaigns perpetrated by foreign state actors and other entities

### PROBLEM

- Social media and, more recently, generative AI make it very easy to generate disinformation and hateful content at scale
- Monitoring via traditional topic/narrative modeling often focuses on low-level content that can be difficult to interpret
- Highly manual workflows typically used for content prioritization and response generation do not scale

## Approach



Counter-Narrative Generation Prototype

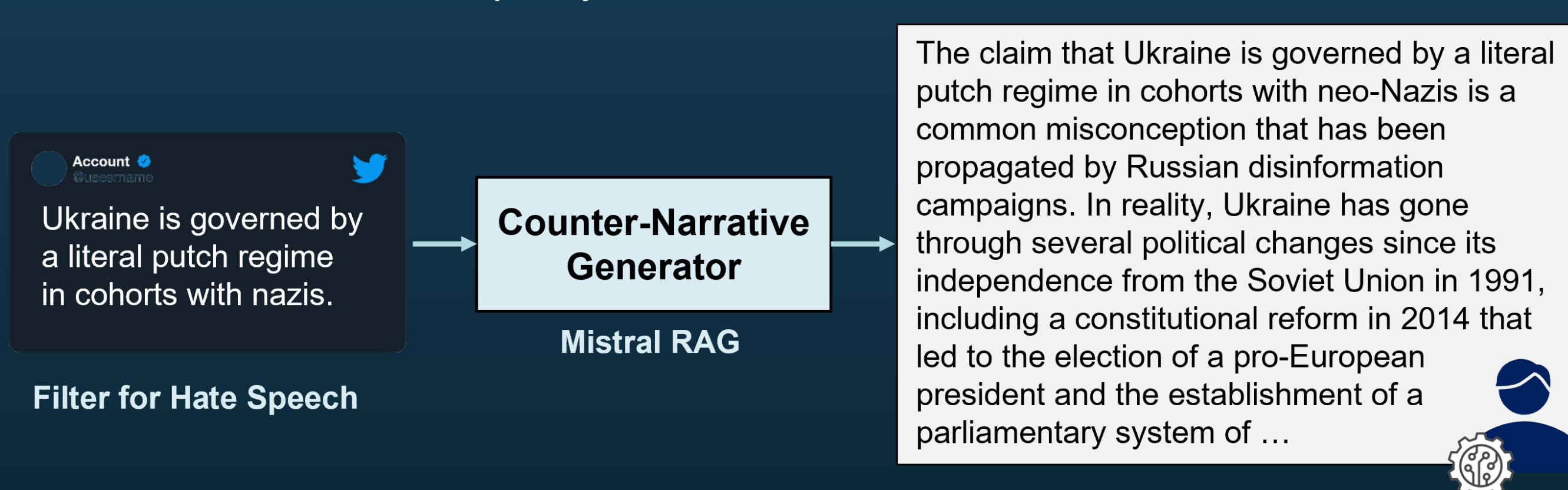### CLASSIFICATION OF SOCIAL MEDIA POSTS CONTAINING INFLAMATORY SPEECH

- Identify extreme content via hate-speech classifier and rank-order tweets in topic clusters
- Zero-shot prompt-based prediction using Mistral Instruct & Twitter's hate speech guidelines
- ~1k hand-labeled evaluation set
- Compare to pre-trained HateBERT1, Fine-Tuned RoBERTa2, prompted LLaMA-7b/LLaMA2-7b



Hate Speech Classification Pipeline

### AUTOMATED RESPONSE TO INFLAMATORY SPEECH

- Retrieval Augmented Generation (RAG) using relevant factual context (Wikipedia, News Articles) to ground auto-generated response (via FAISS over MPNet paragraph embeddings)
- Mistral-7B prompted to make use of factual arguments utilizing the RAG database to generate contextually relevant and correct counter-speech
- Human evaluation on Likert scale of 5 generations for each of 20 random hate speech tweets; 2 raters, 5 quality dimensions



Counter-Narrative Generation Pipeline

## Results

### DATASET AND METRICS FOR EVALUATING HATE SPEECH DETECTION

- Scraped tweets related to Ukraine war and bio-weapons labs during a period leading up to the war
- Time period: between December 2021 and January 2022
- After filtering and removing duplicates, dataset contained ~500k unique tweets

| Model | Precision (⇧) | Recall (⇧) | F1 score (⇧) | Accuracy (⇧) | Time to run in mins (⇩) |
|---|---|---|---|---|---|
| HateBERT Pre-Trained | 0.00 | 0.00 | 0.00 | 0.63 | 117 |
| RoBERTa Fine-Tuned | 0.84 | 0.35 | 0.49 | 0.73 | 105 |
| LLaMA-7B Zero-Shot | 0.38 | **1.00** | 0.54 | 0.38 | 240 |
| LLaMA2-7B Zero-Shot | 0.90 | 0.96 | 0.93 | 0.95 | 102 |
| Mistral-Instruct Prompted | **0.96** | 0.97 | **0.97** | **0.97** | **28** |

Hate Speech Detection Evaluation (BioLab Eval Set)

### METRICS FOR EVALUATING COUNTER-SPEECH GENERATION

- Produced five unique counter-narrative samples for each of 20 randomly selected hateful tweets, resulting in a total of 100 counter-speech samples
- Manually evaluated each counter narrative along 4 dimensions: factuality, relevance, grammaticality and diversity using 1(bad) to 5(good) scale
- Note: one diversity score was assigned for all five counter-narratives responding to a hate tweet
- To ensure an unbiased assessment, two independent raters evaluated the same 100 counter- speech samples
- Inter-rater reliability was quantified using Cohen's Kappa statistic

| Metric | Mean (⇧) | Median (⇧) | Cohen Kappa (⇧) |
|---|---|---|---|
| Factuality | 3.6 | 4 | 0.68 |
| Relevance | 3.8 | 5 | 0.76 |
| Grammaticality | 4.4 | 5 | 0.80 |
| Diversity | 3.7 | 5 | 0.79 |

Counter-Speech Evaluation Metrics

### EXAMPLES OF AUTOMATICALLY GENERATED COUNTER-NARRATIVES



## Conclusion

- Our research demonstrates the effectiveness of zero-shot learning and LLM prompt engineering for nuanced hate speech classification surpassing prior state-of-the-art models
- Our approach effectively combines RAG's information retrieval with LLMs' context processing, overcoming the biases of traditional models and excels in generating coherent and to a large extent relevant and factual counter-narratives

### FUTURE WORK

- Reduce LLM-generated hallucinations (i.e., untruthful information and claims) by expanding RAG database and other approaches
- Develop effective techniques for prompt engineering in order to improve relevance of the generated counter-narratives
- Investigate Langchain's capability to incorporate chat history iterations in order to generate conversational response

LINCOLN LABORATORY
MASSACHUSETTS INSTITUTE OF TECHNOLOGY