

Event Detection from Social Media for Epidemic Preparedness

Anonymous ACL submission

Abstract

Social media is an easy-to-access platform providing timely updates about societal trends and events. Discussions regarding epidemic-related events such as infections, symptoms, and locally deployed measures can be crucial for policy making during epidemic outbreaks. In this work, we exploit Event Detection (ED) for extracting and capturing relevant events from social media posts to provide better preparedness for any upcoming epidemic. To facilitate this task, we curate an epidemic event ontology comprising seven generic event types such as *infect*, *symptom*, *prevent*, etc. Using our event ontology and human expert annotation, we construct our epidemic preparedness Twitter dataset SPEED comprising 1,975 tweets and 2,217 event mentions for the COVID-19 pandemic. Experiments reveal that existing ED models and datasets cannot transfer well for our task, highlighting the challenging nature of our dataset. Finally, we provide empirical evidence highlighting the utility and generalizability of our dataset by showing that ED models trained on our COVID-only dataset SPEED, can effectively identify epidemic events and offer timely warnings for three unseen epidemics of Monkeypox, Zika, and Dengue. This generalizability of SPEED lays the foundations for better preparedness against emerging epidemics.¹

1 Introduction

Early epidemic warnings and effective control measures are among the most important tools for policymakers to be prepared against the threat of any epidemic (Collier et al., 2008). World Health Organization (WHO) reports suggest that 65% of the first reports about infectious diseases and outbreaks originate from informal sources and the internet (Heymann et al., 2001). Social media becomes an important information source here, as it's more timely than other alternatives like news and public

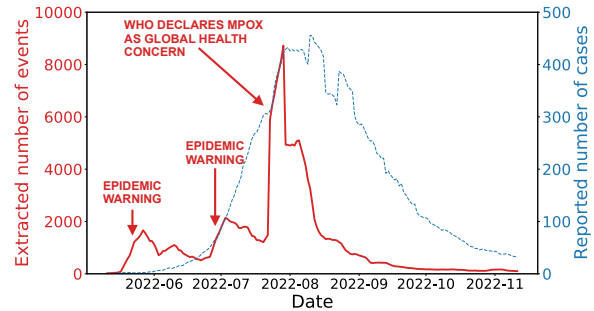


Figure 1: Number of reported Monkeypox cases and the number of extracted events from our trained BERT-QA model from May 11 to Nov 11, 2022. Indicated arrows show how our system can potentially provide early epidemic warnings almost 4-8 weeks before the WHO declared Monkeypox as a pandemic.

health (Lamb et al., 2013), more publicly accessible than clinical notes (Lybarger et al., 2021), and possesses a huge volume of content.² In our work, we explore the social application of information extraction towards building an automated system to efficiently extract epidemic events from social media to provide better epidemic preparedness.

The process of identifying and categorizing significant events based on a pre-defined ontology is a well-established task in NLP, known as *Event Detection* (ED) (Sundheim, 1992; Dodington et al., 2004). However, standard ED datasets mostly focus on general-purpose events for news or Wikipedia domains, and can't be transferred to the epidemic domain (§ 5). Furthermore, most prior epidemiological ED ontologies restrict themselves only to certain diseases or are too fine-grained and specific in nature; while the corresponding datasets majorly focus on news or clinical domains (§ 8). Thus, existing ED ontologies and datasets are not sufficient and models trained on them cannot be readily utilized for extracting events from social

¹Code and data will be released upon acceptance.

²A daily average of 20 million tweets were posted about COVID-19 from May 15 – May 31, 2020.

media for emerging epidemics.

To this end, we construct our own epidemic ED ontology and dataset for social media. Our ontology comprises seven event types - *infect*, *spread*, *symptom*, *prevent*, *cure*, *control*, *death* - chosen based on their relevance for epidemic preparedness, frequency of mentions in social media, and their applicability to various diseases. Our ontology and event definitions are derived from clinical sources (Collier et al., 2008; Babcock et al., 2021) and its sufficiency and coverage are validated by public health experts and quantitative analyses. For our dataset, we choose Twitter as the social media platform and focus on the recent COVID-19 pandemic. Since our task requires domain expertise, we hire six expert annotators to ensure high annotation quality for our dataset. Using our curated ontology and expert annotation, we create our epidemic preparedness dataset **SPEED** (Social Platform based Epidemic Event Detection) comprising 1,975 tweets and 2,217 event mentions. SPEED provides good coverage of events characteristic of any disease and is granular for social media; thus, serves as a valuable ED benchmark for epidemic preparedness from social media.

We benchmark various existing models including four zero-shot models (Shen et al., 2021; Lyu et al., 2021) and two supervised models (Hsu et al., 2022) pre-trained on existing ED datasets of ACE (Doddington et al., 2004) and MAVEN (Wang et al., 2020) on our SPEED benchmark. Experiments reveal that none of the existing models perform well on our dataset mainly owing to the domain-shift and noise in social media as well as unseen epidemic-based event types. Furthermore, training on limited in-domain SPEED data provides significant gains compared to the existing models, highlighting the importance of domain-specific training. Overall, these results reveal how SPEED is a challenging ED dataset.

Tying back to our original motivation of epidemic preparedness, we evaluate the utility and generalizability of our COVID-only dataset SPEED to detect events for any emerging epidemics. More specifically, we evaluate models trained only on SPEED to detect events for three unforeseen epidemics of *Monkeypox*, *Zika*, and *Dengue*. Experiments reveal that SPEED-trained models can successfully detect events for all these epidemics while providing improvements of 29% F1 over zero-shot models and 10% F1 over supervised models trained on small samples of target epidemic data. Further-

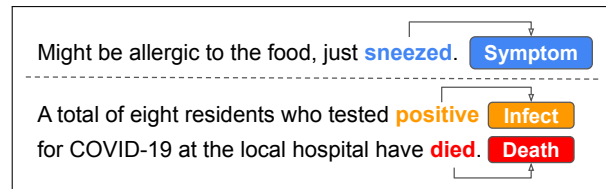


Figure 2: Illustration for the task of Event Detection. Event mentions: Event *symptom* and trigger *sneezed* (1st sentence), Event *infect* and trigger *positive* (2nd sentence), Event *death* and trigger *died* (2nd sentence).

more, by comparing the trends of our extracted events with the actual reported cases, we show that our model can provide early preparedness warnings for the Monkeypox epidemic (Figure 1). These results underscore the strong generalizability and applicability of our dataset SPEED for general epidemic preparedness.

Overall, we make the following contributions: (1) We create an ED ontology and dataset SPEED tailored for predicting epidemic events characteristic of any disease from social media, (2) We show that existing zero-shot models and datasets cannot transfer well to our dataset, highlighting the significance of our dataset, (3) We validate the generalizability of our framework by demonstrating how SPEED-trained ED models using only COVID-tweets can successfully detect events and provide early warnings for three unforeseen epidemics.

2 Task Definition

We employ the task of Event Detection (ED) (Sundheim, 1992; Grishman and Sundheim, 1996) for identifying epidemic events from social media. We define ED based on the ACE 2005 guidelines (Doddington et al., 2004). An **event** is something that happens or describes a change of state and is labeled by a specific **event type**. An **event mention** is the sentence wherein the event is described. Each event mention comprises an **event trigger**, which is the word/phrase that most distinctly highlights the occurrence of the event. **Event Detection** is the task of identifying event triggers from sentences and classifying them into one of the pre-defined event types. The subtask of identifying event triggers is called **Trigger Identification** and classification into event types is **Trigger Classification** (Ahn, 2006). The event types of interest are pre-defined by an **event ontology**. Figure 2 shows examples for three event mentions for the events *symptom*, *infect*, and *death*.

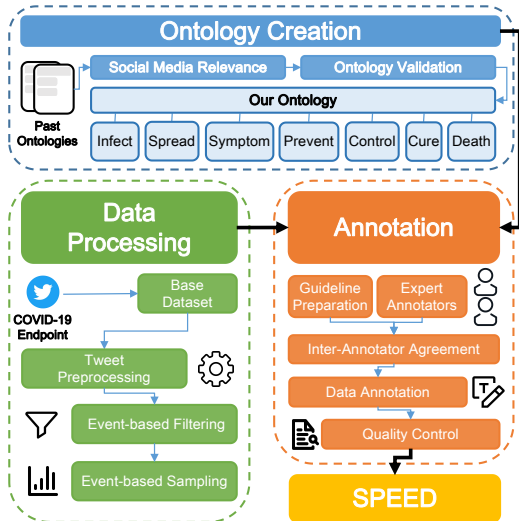


Figure 3: Overview of our dataset creation process.

3 Ontology Creation and Data Collection

We choose social media as our document source since it provides faster and more timely worldly information than other alternatives like news and public health (Lamb et al., 2013) and is more publicly accessible than clinical notes (Lybarger et al., 2021). Owing to its public access and huge content volume, we consider **Twitter**³ as the social media platform and consider the recent **COVID-19 pandemic** as the primary disease for our dataset.

Previous epidemiological ontologies are typically specific to a particular disease, too fine-grained, or cover only a few event types (§ 8 and Table 6) and cannot be readily utilized for ED from social media. Similarly, standard ED datasets don’t comprise epidemiological events and are mostly confined to news or wikipedia domains (§ 8). Owing to these reasons, we create our own event ontology and dataset **SPEED** specific to detecting epidemics from social media. We provide a brief overview of our data creation process in Figure 3 and discuss these steps in more detail below.

3.1 Ontology Creation

Taking inspiration from medical sources like BCEO (Collier et al., 2008), IDO (Babcock et al., 2021), and the ExcavatorCovid (Min et al., 2021a), we curate a wide range of epidemic events while ensuring that they are not biased for specific diseases. We categorize these events into three abstractions of social (events involving larger populations), personal (individual-oriented events), and medical

(medically focused events) types and create our initial ontology comprising 18 event types as reported in Table 19 (§ A.1).

Social Media Relevance To adapt our curated ontology better for social media, we conduct a deeper analysis of the event types based on their frequency and relevance. Majorly, we associate each event type with a certain set of keywords and rank them based on the confidence and frequency of the occurrence of their keywords in social media posts (more details in § A.2). Based on this relevance-based ranking, we merge and discard some event types. Furthermore, we conduct human studies and merge event types to ensure better pairwise distinction.

Ontology Validation and Coverage Drawing upon established epidemiological ontologies serves to guarantee the medical soundness of our ontology. In addition, we assess the sufficiency and comprehensiveness of our ontology and definitions through evaluation by two public health experts. We also quantify our ontology coverage for four diverse diseases by evaluating the percentage of event occurrence in disease-related tweets. We observe a high coverage of 50% for COVID-19, 44% for Monkeypox, 70% for Dengue and 73% for Zika (more details in § A.3), ensuring strong disease coverage of our ontology.

Our final **SPEED** ontology comprises seven major event types that are better suited for social media and cover important aspects of an epidemic. We present our ontology in Table 1 along with event definitions and example event mentions.

3.2 Data Processing

To access a wide range of tweets related to COVID-19, we utilized the Twitter COVID-19 Endpoint released in April 2020. We used a randomized selection of **331 million tweets** between May 15 – May 31 2020, as our base dataset. For preprocessing tweets, we follow Pota et al. (2021): (1) we anonymize personal information like phone numbers, emails, and handles, (2) we normalize any retweets and URLs, (3) we remove emojis and split hashtags, (4) we filter out tweets only in English.

Event-based Filtering Despite COVID-based filtering, most tweets in our base dataset expressed subjective public sentiments, while only 3% comprised mentions adhering to our curated event ontology.⁴ To reduce annotation costs, we further filter

³<https://www.twitter.com/>

⁴Based on keyword-based study conducted on 1,000 tweets

Event Type	Event Definition	Example Event Mention
Infect	The process of a disease/pathogen invading host(s)	Children can also catch COVID-19 ...
Spread	The process of a disease spreading/prevaling massively at a large scale	#COVID-19 CASES RISE TO 85,940 IN INDIA ...
Symptom	Individuals displaying physiological features indicating the abnormality of organisms	(user) (user) Still coughing two months after being infected by this stupid virus ...
Prevent	Individuals trying to prevent the infection of a disease	... wearing mask is the way to prevent COVID-19
Control	Collective efforts trying to impede the spread of epidemic	Social Distancing is our responsibility to reduce spread of COVID-19 ...
Cure	Stopping infection and relieving individuals from infections/symptoms	... recovered corona virus patients cant get it again
Death	End of life of individuals due to infectious disease.	More than 80,000 Americans have died of COVID ...

Table 1: Event ontology comprising seven event types promoting epidemic preparedness along with their definitions and example event mentions. The trigger words are marked in **bold**.

these tweets based on our curated ontology using a simple *sentence embedding* similarity technique. Specifically, we associate each event type with a seed repository of 5-10 diverse tweets. Query tweets are filtered out based on their sentence-level similarity (measured using the BERT sentence embedding model (Reimers and Gurevych, 2019)) with this event-based seed repository. This step filters about 95% tweets from our base dataset significantly reducing the annotation cost.

Event-based Sampling Random sampling of tweets would yield an uneven and COVID-biased distribution of event types for our dataset. We instead perform a uniform sampling - wherein we over-sample tweets linked to less frequent types (e.g. *prevent*) and under-sample the more frequent ones (e.g. *death*). Such an uniform sampling has proven to ensure model robustness (Parekh et al., 2023) - as also validated by our experiments (§ B) - and in turn, would make SPEED generalizable to a wider range of diseases. We sample a total of 1,975 tweets which are utilized for ED annotation.

3.3 Data Annotation

For ED annotation, annotators are tasked with identifying whether a given tweet mentions any of the events outlined in our ontology. If an event is indeed mentioned, annotators are required to identify the specific event trigger. Following the standard ACE dataset (Doddington et al., 2004), we design our annotation guidelines and amend them through several rounds of preliminary annotations to ensure consistency amongst the annotators. Additional details and illustrations of the annotation guidelines and interface are provided in Appendix C.

Annotator Details To ensure high annotation quality and enforce consistency, we choose six experts instead of crowdsourced workers for our annotation. These experts are computer science students studying NLP and are well-versed with the task of ED. They were further trained on our task through multiple loops of annotations and feedback.

Inter-annotator agreement (IAA) We used Fleiss’ Kappa (Fleiss, 1971) for measuring IAA. We conduct two phases of IAA studies: (1) *Guideline Improvement*: Three annotators participated in three annotation rounds with a focus on improving the guidelines. IAA score rose from 0.44 in the first round to 0.59 (70 samples) in the final round. (2) *Agreement Improvement*: All annotators participated in three rounds of annotations. IAA score improved from 0.56 in the first round to a strong 0.65 (50 samples) in the final round.

Quality Control Apart from extensive IAA studies, we deploy two mechanisms to ensure the high annotation quality: (1) *Multi-Annotation*: Each tweet is annotated by two annotators and disagreements are resolved by a third annotator. (2) *Flagging*: Annotators can “flag” ambiguous annotations, which are then resolved and annotated by a third annotator through collective discussion. Both these mechanisms along with a good IAA score ensure that the annotations have high quality.

4 Data Analysis

In this section we present quantitative analyses of our dataset for comparison with other standard ED datasets. Comparison with other epidemiological datasets is discussed in § 8 along with an objective comparison in Table 6.

Dataset	# Event Types	# Sent	# EM	Avg. EM per Event	Domain
ACE	33	18,927	5,055	153.2	News
ERE	38	17,108	7,284	191.7	News
MAVEN	168	49,873	118,732	706.7	Wikipedia
SPEED	7	1,975	2,217	316.7	Social Media

Table 2: Data Statistics for SPEED dataset and comparison with other standard ED datasets. # = “number of”, Avg. = average, Sent = sentences, EM = event mentions.

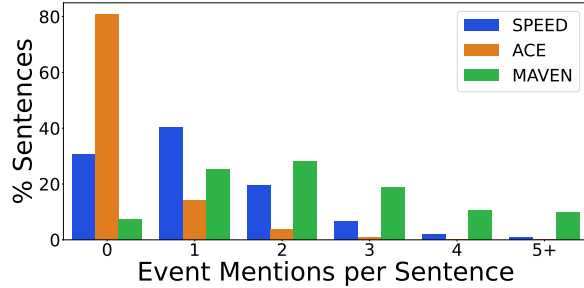


Figure 4: Distribution of number event mentions per sentence. Here % indicates percentage.

Data Statistics Our dataset SPEED comprises seven event types with 2,217 event mentions annotated over 1,975 tweets. We compare SPEED with other standard ED datasets like ACE (Dodgington et al., 2004), ERE (Song et al., 2015), and MAVEN (Wang et al., 2020) in Table 2. Despite the lesser number of sentences and event mentions (since we focus only on 7 event types), SPEED has a reasonable size of 316 average event mentions per event, which is more than the standard ACE and ERE datasets. We also note the differences of the domain of data sources as ACE/ERE focus on News, MAVEN on Wikipedia, while SPEED is based on social media, specifically Twitter data.

Event Mention Density Analysis We compare the distribution of event mentions per sentence with other ED datasets like ACE and MAVEN in Figure 4. We observe that the event density of our dataset is less than MAVEN but better than ACE. This shows that despite having just seven event types, SPEED is a fairly dense dataset.

Trigger Word Analysis We show the diversity of trigger words in SPEED and compare it with other datasets in Table 3. We note that SPEED has a strong average number of triggers per event mention. This demonstrates how SPEED is a diverse and challenging ED dataset.

Dataset	# Unique Triggers	Avg. Triggers per Mention
ACE	1,229	0.24
MAVEN	7,074	0.06
SPEED	555	0.25

Table 3: Comparison of SPEED with ACE and MAVEN in terms of unique trigger words and average number of triggers per event mention. Avg = Average.

5 Transfer Existing Methods

Since existing ED datasets and models are based on general-purpose event ontologies and news/wikipedia domains, they may not transfer well to our social-media-based epidemic detection task. In order to verify this hypothesis, we benchmark the transfer capabilities of these existing methods to our dataset SPEED. For this experimentation, we assume no access to any annotated social media data for epidemic events. We majorly consider the following two families of models:

Zero-shot models do not train on any supervised data and utilize names and definitions of the events for ED. For this, we consider (1) **TE** (Lyu et al., 2021), a pre-trained model that uses event definitions to formulate ED as a textual entailment and question-answering task, (2) **WSD** (Yao et al., 2021) which encodes the contextualized trigger and event definitions jointly and uses a classification head atop for event detection. (3) **TABS** (Li et al., 2022), a model that utilizes two complementary embedding spaces ("mask view" and "token view") to classify examples of new event types. (4) **ETypeClus** (Shen et al., 2021), that extracts salient predicate-object pairs and clusters the embeddings of these pairs in a spherical latent space.

Data transfer models are supervised models pre-trained on other standard ED datasets like ACE (Dodgington et al., 2004) and MAVEN (Wang et al., 2020) and transfer to SPEED in a zero-shot manner. For this, we consider (5) **DEGREE** (Hsu et al., 2022), a generation-based model prompting using natural language templates, (6) **TagPrime** (Hsu et al., 2023), a sequence tagging approach that utilizes priming words to input text to convey more task-specific information.

5.1 Evaluation

We evaluate the above models on the 1,683 tweets from the SPEED dataset. Following previous works (Ahn, 2006), we report the F1-score for the two tasks of Trigger identification (**Tri-I**) and trig-

Model	Tri-I	Tri-C
DATA-TRANSFER		
ACE - TagPrime	0	0
ACE - DEGREE	1.82	1.71
MAVEN - TagPrime	27.65	0
MAVEN - DEGREE	26.72	0
ZERO-SHOT		
TE	9.64	5.54
WSD	17.68	3.65
TABS	3.70	1.61
ETypeClus	17.56	7.66

Table 4: Benchmarking existing zero-shot and data-transfer models on SPEED in terms of Tri-I and Tri-C F1 scores.

ger classification (**Tri-C**) respectively. The results are shown in Table 4. We observe that models pre-trained on the news dataset ACE absolutely fail, while Wikipedia dataset MAVEN pre-training helps to improve Tri-I scores, but still has a nil Tri-C score. The zero-shot models using event definitions perform slightly better, while the best performance is provided by ETypeClus which is an unsupervised clustering model. Overall, **all existing zero-shot and data-transfer models fail to detect epidemic events, mainly owing to the domain shift of social media data and the finer granularity of epidemic events. In turn, this renders SPEED as a challenging ED dataset.**

6 Training with Limited SPEED Data

To improve model performance for SPEED, we conduct experiments trained ED models using limited amounts of in-domain SPEED training data. Majorly, we consider two training paradigms: (1) *Few-shot (FS)*: Models are provided access to n mentions per event (n -shot) for training. We explore 2-shot and 5-shot with three splits of data. (2) *Low Resource (LR)*: Models have access to a limited 100-300 event mentions for training. (Data Statistics in Table 9 in Appendix § D.1).

For training, we consider the following ED models: (1) **DyGIE++** (Wadden et al., 2019), a multi-task classification-based model utilizing local and global context via span graph propagation, (2) **BERT-QA** (Du and Cardie, 2020), a classification model utilizing label semantics by formulating event detection as a question-answering task. We also consider (3) **DEGREE** and (4) **TagPrime** models (as described before in § 5). Other baselines also include (5) **Keyword** (Lejeune et al., 2015),

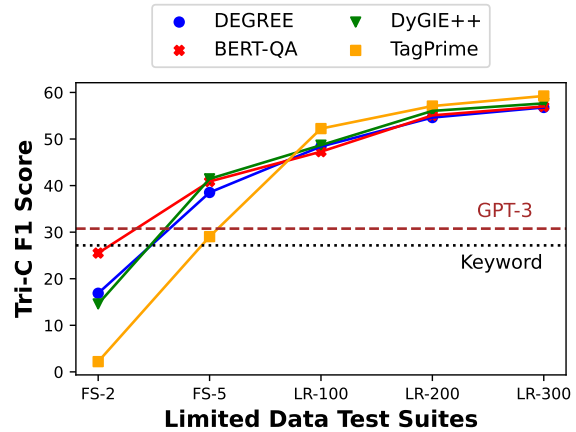


Figure 5: Model performances on the Few-Shot (FS) and Low Resource (LR) test suites in terms of Tri-C F1 scores. Here, LR-XX represents low resource with XX training event mentions and FS-Y represents few-shot with Y training mentions per event.

a popular epidemiological model, that predicts an event if any of the event-specific *curated* keywords are present in the sentence, (6) **GPT-3** (Brown et al., 2020), a large-language model (LLM) baseline using GPT-3.5-turbo as the base model with seven *in-context* examples.

6.1 Evaluation

We follow the same evaluation setup described in § 5.1. Figure 5 presents the model performances for the few-shot and low-resource settings. Majorly, we observe how training on in-domain data can yield performance gains upto 50 F1 points compared to zero-shot and data-transfer methods. We also note that GPT-3 and Keyword baselines are easily outperformed by models trained with just 30 event mentions. Furthermore, these gains are highly consistent for the different ED models. Overall, we note that **small amounts of in-domain training data can provide significant gains in model performance compared to the existing zero-shot and data-transfer models.**

7 Generalization for New Epidemics

Since SPEED focuses solely on COVID-19, its transferability for detecting events for new epidemics remains unknown. To effectively evaluate this generalization, we test if models trained only using our in-domain COVID dataset can detect events for unseen epidemics without any further fine-tuning on the new epidemic data. Specifically, we consider the outbreaks of three diverse diseases

Model	Monkeypox		Zika + Dengue	
	Tri-I	Tri-C	Tri-I	Tri-C
TRANSFER FROM EXISTING DATASETS				
ACE - TagPrime	4.80	0	23.64	0
ACE - DEGREE	12.15	5.14	14.47	0
MAVEN - TagPrime	29.16	0	33.97	0
MAVEN - DEGREE	27.94	0	32.04	0
NO TRAINING + ZERO-SHOT				
TE	16.70	12.11	12.69	9.06
WSD	22.04	4.35	27.93	5.85
ETypeClus	18.31	6.78	13.99	5.33
Keyword	36.40	25.09	25.93	21.69
GPT-3*	42.23	35.33	53.22	14.27
TRAINED FOR TARGET EPIDEMIC				
BERT-QA	59.8	54.08	94.92	80.89
DEGREE	59.58	54.12	86.21	78.76
TagPrime	55.57	49.65	96.67	84.43
DyGIE++	55.83	50.31	73.24	65.65
TRANSFER FROM SPEED				
BERT-QA	67.38	64.17	96.77	81.97
DEGREE	62.95	61.45	88.52	77.69
TagPrime	64.71	61.92	95.24	75.54
DyGIE++	62.76	59.82	91.8	80.34

Table 5: Benchmarking ED models trained on COVID-only SPEED for generalizability to new epidemics of Monkeypox, Zika and Dengue in terms of F1 scores.

of *Monkeypox* (2022), *Zika* (2017), and *Dengue* (2018) as the unseen epidemics.

Experimental Setup For creating datasets, we utilize the Twitter datasets of [Thakur \(2022\)](#) for Monkeypox and [Dias \(2020\)](#) for Zika and Dengue. Using expert annotation for a sample of the tweets, our final evaluation dataset comprises 286 tweets with 398 event mentions for Monkeypox while 300 tweets with 274 event mentions for Zika and Dengue (statistics in § D.3 and § D.4).

For model training, we use a 80-20 split of our COVID-only SPEED dataset to train various ED models (TRANSFER FROM SPEED). For comparison, we benchmark models trained on existing datasets (TRANSFER FROM EXISTING DATASETS) and models requiring no training data (NO TRAINING). As strong baselines, we also consider supervised models trained on a small sample of 300 tweets for the target epidemic (TRAINED ON TARGET EPIDEMIC).

Results We present our results in Table 5. None of the existing data transfer methods or zero-shot methods perform well. Overall, we observe that ED models transferring from SPEED perform the best with model performance ranging from 60-65 F1

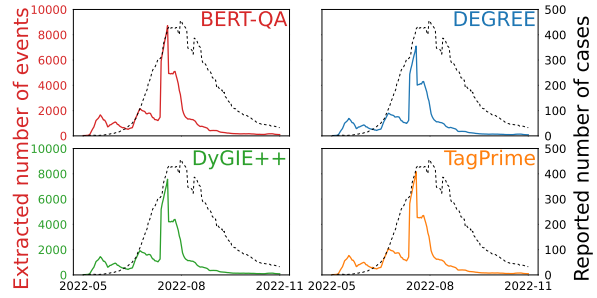


Figure 6: Number of reported Monkeypox cases and the number of extracted events from four trained models from May 11 to Nov 11, 2022.

points, thus **demonstrating the generalizability of our SPEED dataset to new epidemics**. Furthermore, we observe SPEED-trained models even outperform models trained on for Monkeypox by 10 F1 points and are at par for Zika. This outcome is particularly encouraging, as it **demonstrates the resilience of SPEED-trained models, making them highly applicable during the early stages of an unfamiliar epidemic, when minimal to no epidemic-specific data is accessible**.

7.1 Providing Early Epidemic Warnings

As a challenging yet practical evaluation, we evaluate our SPEED-trained models in their ability to provide early warnings for an unknown epidemic. We choose the Monkeypox as the unknown epidemic. We report the number of epidemic events extracted by the BERT-QA trained on SPEED along with the actual number of Monkeypox cases reported in the US⁵ from May 11 to Nov 11, 2022, in Figure 1. As shown by the arrows in the figure, our model could potentially provide two sets of early warnings around May 23 and June 29 before the outbreak reached its peak around July 30. In fact, all our trained ED models are capable of providing these early signals as shown in Figure 6 (further analysis in Appendix F). This robust outcome underscores the **real-time practicality to provide early warnings and broad applicability to unknown epidemics of our SPEED dataset**.

8 Related Work

Event Extraction Datasets Event Extraction (EE) is the task of detecting events (Event Detection) and extracting structured information about specific roles linked to the event (Event Argument

⁵As reported by CDC at <https://www.cdc.gov/poxvirus/mpox/response/2022/mpx-trends.html>

Dataset	Data Source	Sentence Level	Trigger Present	Social Events	Personal Events	Social Media Granular
SPEED (Ours)	Twitter	✓	✓	✓	✓	✓
COVIDKB (Zong et al., 2022)	Twitter	✓	✗	✗	✓	✓
CACT (Lybarger et al., 2021)	Clinical	✗	✗	✗	~	✓
ExcavatorCovid (Min et al., 2021b)	News	✗	✓	✓	✓	✗
BioCaster (Collier et al., 2008)	News	✗	✗	✓	✓	✗
DANIEL (Lejeune et al., 2015)	News	✗	~	✗	~	✓

Table 6: Objective comparison of various epidemiological datasets with our dataset SPEED. We objectify the source of base data (Data Source), the level of annotation granularity (Sentence Level), the presence of trigger information (Trigger Present), the presence of social and personal events in ontology (Social Events and Personal Events), and the suitability of ontology for social media (Social Media Granular). ~ indicates partial presence.

Extraction) from natural text. Earliest works for this task can be dated back to MUC (Sundheim, 1992; Grishman and Sundheim, 1996) and the more standard ACE (Doddington et al., 2004). Over the years, ACE was extended to various datasets like ERE (Song et al., 2015) and TAC KBP (Ellis et al., 2015). Recent progress has been the creation of massive datasets and huge event ontologies with datasets like MAVEN (Wang et al., 2020), RAMS (Ebner et al., 2020), WikiEvents (Li et al., 2021), DocEE (Tong et al., 2022), GENEVA (Parekh et al., 2023) and GLEN (Zhan et al., 2023). These ontologies and datasets cater to general-purpose events and do not comprise epidemiological event types.

Epidemiological Ontologies Earliest works (Lindberg et al., 1993; Rector et al., 1996) defined highly rich taxonomies for describing technical concepts used by biomedical experts. Further developments led to the creation of SNOMED CT (Stearns et al., 2001) and PHSkb (Doyle et al., 2005) that define a list of reportable events used for communication between public health experts. BioCaster (Collier et al., 2008) and PULS (Du et al., 2011) extended ontologies for the news domain. Recent works of NCBI (Dogan et al., 2014), IDO (Babcock et al., 2021) and DO (Schriml et al., 2022) focus on comprehensively organizing human diseases. In light of the recent COVID-19 pandemic, CIDO (He et al., 2020) define a technical taxonomy for coronavirus, while ExcavatorCovid (Min et al., 2021b) automatically extract COVID-19 events and relations between them. Most of these ontologies are too fine-grained or limited to specific events, and can’t be directly used for ED from social media, as also shown in Table 6.

Epidemiological Information Extraction Early works utilized search-engine queries and click-

through rates for predicting influenza trends (Eysenbach, 2006; Ginsberg et al., 2009). Information extraction from Twitter has also been quite successful for predicting influenza trends (Signorini et al., 2011; Lamb et al., 2013; Paul et al., 2014). Over the years, various biomedical monitoring systems have been developed like BioCaster (Collier et al., 2008; Meng et al., 2022), HeathMap (Freifeld et al., 2008), DANIEL (Lejeune et al., 2015), EpiCore (Olsen, 2017). Extensions to support multilingual systems has also been explored (Lejeune et al., 2015; Mutuvi et al., 2020; Sahnoun and Lejeune, 2021). For the COVID-19 pandemic, several frameworks like CACT (Lybarger et al., 2021) and COVIDKB (Zong et al., 2022) were developed for extracting symptoms and infection statistics respectively. Most of these systems focus on the domains of news and clinical notes and use keyword/rule-based or simple BERT-based models, as shown in Table 6. In our work, we explore more recent ED models while focusing specifically on the social media domain.

9 Conclusion and Future Work

In this work, we leverage the framework of Event Detection (ED) to extract epidemic events from social media to promote better epidemic preparedness. To facilitate this, we create our Twitter-based dataset SPEED comprising seven major epidemic event types. Through experimentation, we show how existing datasets and models fail to transfer for our task. Contrastingly, we show how models trained on SPEED can generalize and provide early warnings for unseen emerging epidemics. More broadly, our work demonstrates how event extraction and in general, information extraction can exploit social media to aid policy-making for better epidemic preparedness.

567 Limitations

568 Our work focuses majorly on a single source of
569 social media - Twitter. We haven't explored other
570 social media platforms and how ED would work on
571 those platforms in our work. We leave that for fu-
572 ture work, but are optimistic that our models should
573 be able to generalize across platforms. Secondly,
574 our work mainly only focuses on ED as the pri-
575 mary task, while its sister task Event Argument
576 Extraction (EAE) is not explored. We hope to ex-
577 tend our work for EAE as part of our future work.
578 Finally, we would like to show the generalization
579 of our models on a vast range of diseases. How-
580 ever owing to budget constraints and the lack of
581 publically available Twitter data for other diseases,
582 we couldn't perform such a study. However, we
583 believe showing results on three diseases lays the
584 foundation for generalizability of our model.

585 Ethical Considerations

586 One strong assumption in our work is the avail-
587 ability of internet and social media for discussions
588 about epidemics. Since not everyone has equal ac-
589 cess to these platforms, our dataset, models, and
590 results do not represent the whole world uniformly.
591 Thus, our work can be biased and should be consid-
592 ered with other sources for better representation.

593 Our dataset SPEED is based on actual tweets
594 posted by people all over the world. We attempted
595 our best to anonymize any kind of private informa-
596 tion in the tweets, but we can never be completely
597 thorough, and there might be some private infor-
598 mation embedded still in our dataset. Furthermore,
599 these tweets were sentimental and may possess
600 stark emotional, racial, and political viewpoints
601 and biases. We do not attempt to clean any of such
602 extreme data in our work (as our focus was on
603 ED only) and these biases should be considered if
604 being used for other applications.

605 Since our ED models are trained on SPEED, they
606 may possess some of the social biases embedded
607 in SPEED. Since our work didn't focus on bias
608 mitigation, these models should be used with due
609 consideration.

610 Lastly, we do not claim that our models can
611 be used off-the-shelf for epidemic prediction as
612 it hasn't been thoroughly tested and can have false
613 positives and negatives too. We majorly throw light
614 to show these model capabilities and motivate fu-
615 ture work in this direction. The usage of these
616 systems for practical purposes should be appropri-

ately considered.

References

- David Ahn. 2006. *The stages of event extraction*. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, pages 1–8, Sydney, Australia. Association for Computational Linguistics.
- Shane Babcock, John Beverley, Lindsay G. Cowell, and Barry Smith. 2021. *The infectious disease ontology in the age of COVID-19*. *J. Biomed. Semant.*, 12(1):13.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- YunSeok Choi, Hyojun Kim, and Jee-Hyong Lee. 2022. *TABS: Efficient textual adversarial attack for pre-trained NL code model using semantic beam search*. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5490–5498, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Nigel Collier, Son Doan, Ai Kawazoe, Reiko Matsuda Goodwin, Mike Conway, Yoshio Tateno, Hung Quoc Ngo, Dinh Dien, Asanee Kawtrakul, Koichi Takeuchi, Mika Shigematsu, and Kiyosu Taniguchi. 2008. *Bio-caster: detecting public health rumors with a web-based text mining system*. *Bioinform.*, 24(24):2940–2941.
- Guilherme Dias. 2020. *Tweets dataset on Zika virus*.
- George Doddington, Alexis Mitchell, Mark Przybocki, Lance Ramshaw, Stephanie Strassel, and Ralph Weischedel. 2004. *The automatic content extraction (ACE) program – tasks, data, and evaluation*. In *Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC'04)*, Lisbon, Portugal. European Language Resources Association (ELRA).
- Rezarta Islamaj Dogan, Robert Leaman, and Zhiyong Lu. 2014. *NCBI disease corpus: A resource for disease name recognition and concept normalization*. *J. Biomed. Informatics*, 47:1–10.
- Timothy J. Doyle, Haobo Ma, Samuel L. Groseclose, and Richard S. Hopkins. 2005. *Phskb: A knowledge-base to support notifiable disease surveillance*. *BMC Medical Informatics Decis. Mak.*, 5:27.
- Mian Du, Peter von Etter, Mikhail Kopotev, Mikhail Novikov, Natalia Tarbeeva, and Roman Yangarber. 2011. *Building support tools for russian-language information extraction*. In *Text, Speech and Dialogue - 14th International Conference, TSD 2011, Pilsen, Czech Republic, September 1-5, 2011. Proceedings*, volume 6836 of *Lecture Notes in Computer Science*, pages 380–387. Springer.

672	Xinya Du and Claire Cardie. 2020. Event extraction by answering (almost) natural questions . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 671–683, Online. Association for Computational Linguistics.	728
673		729
674		730
675		731
676		
677	Seth Ebner, Patrick Xia, Ryan Culkin, Kyle Rawlins, and Benjamin Van Durme. 2020. Multi-sentence argument linking . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8057–8077, Online. Association for Computational Linguistics.	732
678		733
679		734
680		735
681		736
682		737
683	Joe Ellis, Jeremy Getman, Dana Fore, Neil Kuster, Zhiyi Song, Ann Bies, and Stephanie M. Strassel. 2015. Overview of linguistic resources for the TAC KBP 2015 evaluations: Methodologies and results . In <i>Proceedings of the 2015 Text Analysis Conference, TAC 2015, Gaithersburg, Maryland, USA, November 16-17, 2015, 2015</i> . NIST.	738
684		739
685		740
686		
687		
688		
689		
690	Gunther Eysenbach. 2006. Infodemiology: Tracking flu-related searches on the web for syndromic surveillance . In <i>AMIA 2006, American Medical Informatics Association Annual Symposium, Washington, DC, USA, November 11-15, 2006</i> . AMIA.	741
691		742
692		743
693		744
694		745
695	Joseph L Fleiss. 1971. Measuring nominal scale agreement among many raters. <i>Psychological bulletin</i> , 76(5):378.	746
696		747
697		
698	Clark C. Freifeld, Kenneth D. Mandl, Ben Y. Reis, and John S. Brownstein. 2008. Model formulation: Healthmap: Global infectious disease monitoring through automated classification and visualization of internet media reports . <i>J. Am. Medical Informatics Assoc.</i> , 15(2):150–157.	748
699		749
700		750
701		751
702		752
703		753
704	Jeremy Ginsberg, Matthew H Mohebbi, Rajan S Patel, Lynnette Brammer, Mark S Smolinski, and Larry Brilliant. 2009. Detecting influenza epidemics using search engine query data. <i>Nature</i> , 457(7232):1012–1014.	754
705		755
706		756
707		757
708		758
709	Ralph Grishman and Beth Sundheim. 1996. Message Understanding Conference- 6: A brief history . In <i>COLING 1996 Volume 1: The 16th International Conference on Computational Linguistics</i> .	759
710		760
711		761
712		762
713	Yongqun He, Hong Yu, Edison Ong, Yang Wang, Yingtong Liu, Anthony Huffman, Hsin-Hui Huang, John Beverley, Asiyah Yu Lin, William D. Duncan, Sivaram Arabandi, Jiangan Xie, Junguk Hur, Xiaolin Yang, Luonan Chen, Gilbert S. Omenn, Brian D. Athey, and Barry Smith. 2020. CIDO: the community-based coronavirus infectious disease ontology . In <i>Proceedings of the 11th International Conference on Biomedical Ontologies (ICBO) joint with the 10th Workshop on Ontologies and Data in Life Sciences (ODLS) and part of the Bolzano Summer of Knowledge (BoSK 2020), Virtual conference hosted in Bolzano, Italy, September 17, 2020</i> , volume 2807 of <i>CEUR Workshop Proceedings</i> , pages 1–10. CEUR-WS.org.	763
714		764
715		765
716		766
717		767
718		768
719		769
720		770
721		771
722		
723		
724		
725		
726		
727		
	David L Heymann, Guénaél R Rodier, et al. 2001. Hot spots in a wired world: Who surveillance of emerging and re-emerging infectious diseases. <i>The Lancet infectious diseases</i> , 1(5):345–353.	772
		773
		774
	I-Hung Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott Miller, Prem Natarajan, Kai-Wei Chang, and Nanyun Peng. 2022. DEGREE: A data-efficient generation-based event extraction model . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1890–1908, Seattle, United States. Association for Computational Linguistics.	775
		776
		777
		778
		779
	I-Hung Hsu, Kuan-Hao Huang, Shuning Zhang, Wenxin Cheng, Premkumar Natarajan, Kai-Wei Chang, and Nanyun Peng. 2023. A simple and unified tagging model with priming for relational structure predictions . In <i>Proceedings of the 61st Conference of the Association for Computational Linguistics</i> , Toronto, Canada. Association for Computational Linguistics.	780
		781
		782
		783
		784
	Alex Lamb, Michael J. Paul, and Mark Dredze. 2013. Separating fact from fear: Tracking flu infections on Twitter . In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 789–795, Atlanta, Georgia. Association for Computational Linguistics.	
	Gaël Lejeune, Romain Brixstel, Antoine Doucet, and Nadine Lucas. 2015. Multilingual event extraction for epidemic detection . <i>Artif. Intell. Medicine</i> , 65(2):131–143.	
	Sha Li, Heng Ji, and Jiawei Han. 2021. Document-level event argument extraction by conditional generation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 894–908, Online. Association for Computational Linguistics.	
	Sha Li, Heng Ji, and Jiawei Han. 2022. Open relation and event type discovery with type abstraction . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 6864–6877, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
	D. A. Lindberg, B. L. Humphreys, and A. T. McCray. 1993. The unified medical language system . <i>Methods of information in medicine</i> , 32(4):281—291.	
	Kevin Lybarger, Mari Ostendorf, Matthew Thompson, and Meliha Yetisgen. 2021. Extracting COVID-19 diagnoses and symptoms from clinical text: A new annotated corpus and neural event extraction framework . <i>J. Biomed. Informatics</i> , 117:103761.	
	Qing Lyu, Hongming Zhang, Elior Sulem, and Dan Roth. 2021. Zero-shot event extraction via transfer learning: Challenges and insights . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International</i>	

785				
786				
787				
788	Zaiqiao Meng, Anya Okhmatovskaia, Maxime Polleri,			
789	Yannan Shen, Guido Powell, Zihao Fu, Iris Ganser,			
790	Meiru Zhang, Nicholas B. King, David L. Buck-			
791	eridge, and Nigel Collier. 2022. Biocaster in 2021:			
792	automatic disease outbreaks detection from global			
793	news media . <i>Bioinform.</i> , 38(18):4446–4448.			
794	Bonan Min, Benjamin Rozonoyer, Haoling Qiu, Alexan-			
795	der Zamanian, Nianwen Xue, and Jessica MacBride.			
796	2021a. Excavatorcovid: Extracting events and rela-			
797	tions from text corpora for temporal and causal			
798	analysis for COVID-19 . In <i>Proceedings of the 2021</i>			
799	<i>Conference on Empirical Methods in Natural Lan-</i>			
800	<i>guage Processing: System Demonstrations, EMNLP</i>			
801	<i>2021, Online and Punta Cana, Dominican Republic,</i>			
802	<i>7-11 November, 2021</i> , pages 63–71. Association for			
803	Computational Linguistics.			
804	Bonan Min, Benjamin Rozonoyer, Haoling Qiu, Alexan-			
805	der Zamanian, Nianwen Xue, and Jessica MacBride.			
806	2021b. ExcavatorCovid: Extracting events and rela-			
807	tions from text corpora for temporal and causal			
808	analysis for COVID-19 . In <i>Proceedings of the 2021</i>			
809	<i>Conference on Empirical Methods in Natural Lan-</i>			
810	<i>guage Processing: System Demonstrations</i> , pages			
811	63–71, Online and Punta Cana, Dominican Republic.			
812	Association for Computational Linguistics.			
813	Stephen Mutuvi, Antoine Doucet, Gaël Lejeune, and			
814	Moses Odeo. 2020. A dataset for multi-lingual epi-			
815	demiological event extraction . In <i>Proceedings of the</i>			
816	<i>Twelfth Language Resources and Evaluation Confer-</i>			
817	<i>ence</i> , pages 4139–4144, Marseille, France. European			
818	Language Resources Association.			
819	Jennifer M Olsen. 2017. Epicore: crowdsourcing health			
820	professionals to verify disease outbreaks. <i>Online</i>			
821	<i>Journal of Public Health Informatics</i> , 9(1).			
822	OpenAI. 2021. ChatGPT: Large-scale language model .			
823	Accessed: June 17, 2023.			
824	Tanmay Parekh, I-Hung Hsu, Kuan-Hao Huang, Kai-			
825	Wei Chang, and Nanyun Peng. 2023. Geneva: Bench-			
826	marking generalizability for event argument extrac-			
827	tion with hundreds of event types and argument roles .			
828	In <i>Proceedings of the 61st Conference of the Associa-</i>			
829	<i>tion for Computational Linguistics</i> , Toronto, Canada.			
830	Association for Computational Linguistics.			
831	Michael J Paul, Mark Dredze, and David Broniatowski.			
832	2014. Twitter improves influenza forecasting. <i>PLoS</i>			
833	<i>currents</i> , 6.			
834	Marco Pota, Mirko Ventura, Hamido Fujita, and			
835	Massimo Esposito. 2021. Multilingual evaluation			
836	of pre-processing for bert-based sentiment analy-			
837	sis of tweets . <i>Expert Systems with Applications</i> ,			
838	181:115119.			
839	Alan L Rector, Jeremy E Rogers, and Pam Pole. 1996.			
840	The galen high level ontology. In <i>Medical Informat-</i>			
841	<i>ics Europe '96</i> , pages 174–178. IOS Press.			
	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:			842
	Sentence embeddings using siamese bert-networks .			843
	In <i>Proceedings of the 2019 Conference on Empirical</i>			844
	<i>Methods in Natural Language Processing</i> . Associa-			845
	tion for Computational Linguistics.			846
	Sihem Sahnoun and Gaël Lejeune. 2021. Multilingual			847
	epidemic event extraction : From simple classifica-			848
	tion methods to open information extraction (OIE)			849
	and ontology . In <i>Proceedings of the International</i>			850
	<i>Conference on Recent Advances in Natural Language</i>			851
	<i>Processing (RANLP 2021)</i> , pages 1227–1233, Held			852
	Online. INCOMA Ltd.			853
	Lynn M. Schriml, James B. Munro, Mike Schor, Dustin			854
	Olley, Carrie McCracken, Victor Felix, J. Allen			855
	Baron, Rebecca C. Jackson, Susan M. Bello, Cyn-			856
	thia Bearer, Richard Lichenstein, Katharine Bisordi,			857
	Nicole Champion, Michelle G. Giglio, and Carol			858
	Greene. 2022. The human disease ontology 2022			859
	update . <i>Nucleic Acids Res.</i> , 50(D1):1255–1261.			860
	Jiaming Shen, Yunyi Zhang, Heng Ji, and Jiawei Han.			861
	2021. Corpus-based open-domain event type induc-			862
	tion . In <i>Proceedings of the 2021 Conference on Em-</i>			863
	<i>pirical Methods in Natural Language Processing</i> ,			864
	pages 5427–5440, Online and Punta Cana, Domini-			865
	can Republic. Association for Computational Lin-			866
	guistics.			867
	Alessio Signorini, Alberto Maria Segre, and Philip M			868
	Polgreen. 2011. The use of twitter to track lev-			869
	els of disease activity and public concern in the us			870
	during the influenza a h1n1 pandemic. <i>PloS one</i> ,			871
	6(5):e19467.			872
	Zhiyi Song, Ann Bies, Stephanie Strassel, Tom Riese,			873
	Justin Mott, Joe Ellis, Jonathan Wright, Seth Kulick,			874
	Neville Ryant, and Xiaoyi Ma. 2015. From light			875
	to rich ERE: Annotation of entities, relations, and			876
	events . In <i>Proceedings of the The 3rd Workshop on</i>			877
	<i>EVENTS: Definition, Detection, Coreference, and</i>			878
	<i>Representation</i> , pages 89–98, Denver, Colorado. As-			879
	sociation for Computational Linguistics.			880
	Michael Q. Stearns, Colin Price, Kent A. Spackman,			881
	and Amy Y. Wang. 2001. SNOMED clinical terms:			882
	overview of the development process and project			883
	status . In <i>AMIA 2001, American Medical Informat-</i>			884
	<i>ics Association Annual Symposium, Washington, DC,</i>			885
	<i>USA, November 3-7, 2001</i> . AMIA.			886
	Beth M. Sundheim. 1992. Overview of the fourth Mes-			887
	sage Understanding Evaluation and Conference . In			888
	<i>Fourth Message Understanding Conference (MUC-</i>			889
	<i>4): Proceedings of a Conference Held in McLean,</i>			890
	<i>Virginia, June 16-18, 1992</i> .			891
	Nirmalya Thakur. 2022. Monkeypox2022tweets: A			892
	large-scale twitter dataset on the 2022 monkeypox			893
	outbreak, findings from analysis of tweets, and open			894
	research questions . <i>Infectious Disease Reports</i> ,			895
	14(6):855–883.			896
	MeiHan Tong, Bin Xu, Shuai Wang, Meihuan Han,			897
	Yixin Cao, Jiangqi Zhu, Siyu Chen, Lei Hou, and			898

899 Juanzi Li. 2022. [DocEE: A large-scale and fine-](#)
900 [grained benchmark for document-level event extrac-](#)
901 [tion](#). In *Proceedings of the 2022 Conference of the*
902 *North American Chapter of the Association for Com-*
903 *putational Linguistics: Human Language Technolo-*
904 *gies*, pages 3970–3982, Seattle, United States. Asso-
905 ciation for Computational Linguistics.

906 David Wadden, Ulme Wennberg, Yi Luan, and Han-
907 naneh Hajishirzi. 2019. [Entity, relation, and event](#)
908 [extraction with contextualized span representations](#).
909 In *Proceedings of the 2019 Conference on Empirical*
910 *Methods in Natural Language Processing and the*
911 *9th International Joint Conference on Natural Lan-*
912 *guage Processing (EMNLP-IJCNLP)*, pages 5784–
913 5789, Hong Kong, China. Association for Computa-
914 tional Linguistics.

915 Xiaozhi Wang, Ziqi Wang, Xu Han, Wangyi Jiang, Rong
916 Han, Zhiyuan Liu, Juanzi Li, Peng Li, Yankai Lin,
917 and Jie Zhou. 2020. [MAVEN: A Massive General](#)
918 [Domain Event Detection Dataset](#). In *Proceedings*
919 *of the 2020 Conference on Empirical Methods in*
920 *Natural Language Processing (EMNLP)*, pages 1652–
921 1671, Online. Association for Computational Linguis-
922 tics.

923 Wenlin Yao, Xiaoman Pan, Lifeng Jin, Jianshu Chen,
924 Dian Yu, and Dong Yu. 2021. [Connect-the-dots:](#)
925 [Bridging semantics between words and definitions](#)
926 [via aligning word sense inventories](#). In *Proceedings*
927 *of the 2021 Conference on Empirical Methods in Nat-*
928 *ural Language Processing*, pages 7741–7751, Online
929 and Punta Cana, Dominican Republic. Association
930 for Computational Linguistics.

931 Qiusi Zhan, Sha Li, Kathryn Conger, Martha Palmer,
932 Heng Ji, and Jiawei Han. 2023. [Glen: General-](#)
933 [purpose event detection for thousands of types](#). *arXiv*
934 *preprint arXiv:2303.09093*.

935 Shi Zong, Ashutosh Baheti, Wei Xu, and Alan Ritter.
936 2022. [Extracting a knowledge base of COVID-19](#)
937 [events from social media](#). In *Proceedings of the 29th*
938 *International Conference on Computational Linguis-*
939 *tics*, pages 3810–3823, Gyeongju, Republic of Korea.
940 International Committee on Computational Linguis-
941 tics.

942 A Ontology Creation - Additional Details 989

943 A.1 Complete ontology 990

944 We present our complete initial event ontology com- 991
945 prising 18 event types organized into 3 abstract cat- 992
946 egories in Table 19. We also describe each event 993
947 type by its definition and also present details about 994
948 the action taken for its role in the final event ontol- 995
949 ogy. 996

950 A.2 Initial analysis of events 997

951 Our initial ontology (§ A.1) was constructed using 998
952 previous ontologies and human knowledge. But 999
953 the suitability of each event type for social media 1000
954 (specifically Twitter) remains unknown. To evalu- 1001
955 ate this suitability, we use frequency and confi- 1002
956 dence as two guiding heuristics and use them for 1003
957 final filtering/merging. We utilize the base Twitter 1004
958 dataset for SPEED for conducting this analysis. We 1005
959 describe each of these heuristics here: 1006

960 **Frequency** To approximately estimate the fre- 1007
961 quency of events, we curate a list of keywords for 1008
962 each event type and count the number of posts con- 1009
963 taining any of these keywords. Keyword curation 1010
964 involves creating a seed list using human expert 1011
965 knowledge and expanding that list using synonyms 1012
966 from external sources like Thesaurus.⁶ We show 1013
967 the results in Figure 7. We observe that most events 1014
968 under the medical abstraction occur much lesser 1015
969 than others. Furthermore, the count variance is 1016
970 large as the most frequent event type *control* is 180 1017
971 times more likely to occur than the least frequent 1018
972 event type *variant*. Since low-frequency events are 1019
973 less likely to be mentioned in a smaller sample of 1020
974 data, we discard or merge such events for our final 1021
975 ontology. 1022

976 **Confidence** For each keyword, we randomly 1023
977 sample a small number of non-duplicate tweets 1024
978 and manually rate the keyword confidence based 1025
979 on the percentage of tweets wherein the semantic 1026
980 meaning of the keyword matches the definition of 1027
981 its event. We mainly categorize this confidence 1028
982 as high, medium, or low.⁷ Take event *control* 1029
983 as an example, it has high confidence keywords 1030
984 such as “quarantine”, “protocol”, and “distancing”; 1031
985 medium confidence keywords such as “restrict”, 1032
986 “postpone”, and “investigate”; low confidence key- 1033
987 words such as “battle”, “limitation”, and “separa- 1034
988 tion”. On the other hand, event *prefigure* does not 1035

⁶<https://www.thesaurus.com/>

⁷We release these keywords as part of our final code. 1036

989 have high confidence keywords, but only medium 990
991 confidence keywords such as “foreshadow” and 992
993 low confidence keywords such as “foretell”. Our 994
995 heuristic suggests that low-confidence keywords 996
997 are more likely to give false positives relative to 998
999 high-confidence ones. Thus, we filter/merge event 1000
1001 types that have a high number of low-confidence 1002
1003 keywords. 1004

1005 Eventually, our final ontology comprises seven 1006
1007 events that are distinguishable, frequent, and have 1008
1009 a low false-positive rate. 1010

1011 A.3 Coverage analysis of ontology 1000

1012 To quantitatively verify the coverage of our ontol- 1001
1013 ogy, we conduct an analysis on four diseases with 1002
1014 very different characteristics - COVID-19, Mon- 1003
1015 keypox, Dengue, and Zika. For each disease, we 1004
1016 randomly sample 300 tweets and then filter them 1005
1017 if they are related to the disease or not. Next, we 1006
1018 annotate the filtered disease-related tweets based 1007
1019 on our ontology and evaluate the proportion of 1008
1020 event occurrences relative to the number of disease- 1009
1021 related tweets. We find that our ontology has high 1010
1022 coverage of 50% for COVID-19, 44% for Monkey- 1011
1023 pox, 70% for Dengue, and 73% for Zika. This in 1012
1024 turn assures that our ontology can be used to de- 1013
1025 tect epidemic events for various different kinds of 1014
1026 diseases. 1015

1016 **Event Type Distribution** As part of our analysis, 1016
1017 we also study our ontology’s event type distribu- 1017
1018 tion for each disease and its correlation with the disease 1018
1019 properties and outbreak stage. We show this event 1019
1020 distribution in Figure 8 for each of the diseases. We 1020
1021 note that distributions for Dengue and Monkeypox 1021
1022 exhibit a strong focus on *spread* and *infect* events. 1022
1023 This makes sense as the data for these diseases was 1023
1024 collected at earlier stages of the outbreak when mit- 1024
1025 igation measures were not being discussed yet. On 1025
1026 the other hand, for COVID-19, the distribution is 1026
1027 vastly dominated by *control* and *death* events. Our 1027
1028 COVID-19 data was collected in May 2020 when 1028
1029 the outbreak had vastly spread in America. Thus 1029
1030 our distribution reflects more notions of lockdowns 1030
1031 and control measures as well reflects the deadly 1031
1032 nature of the disease. 1032

1033 B Uniform Sampling v/s Random 1033

1034 Sampling for Data Selection 1034

1035 Previously Parekh et al. (2023) had shown how uni- 1035
1036 form sampling of data for events can yield more 1036
1037 robust model performance. To validate the same 1037

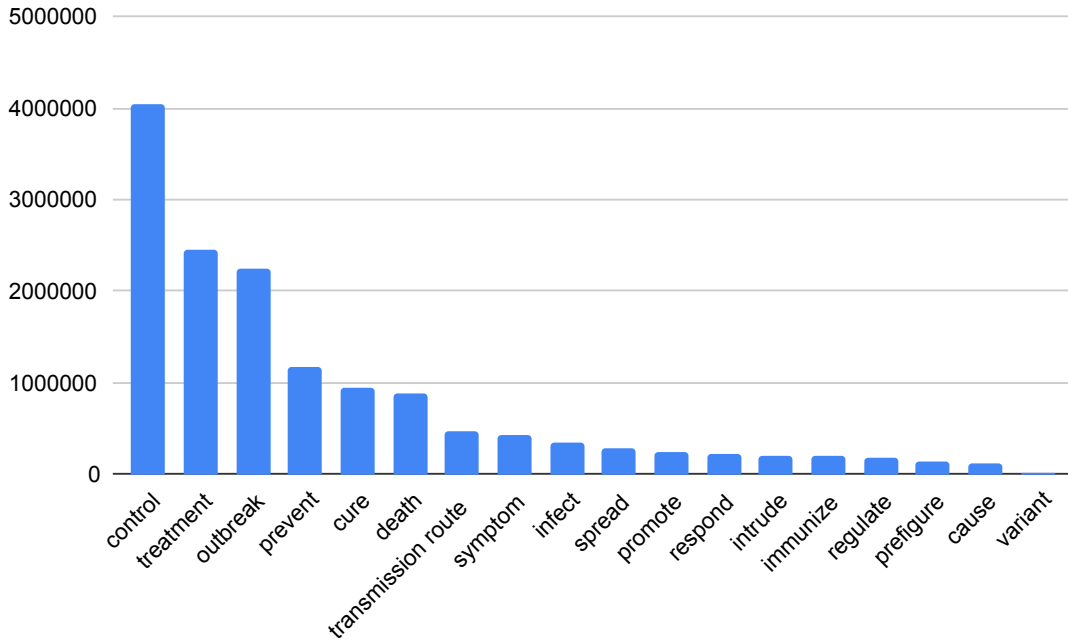


Figure 7: Frequency of occurrence based on keyword search for all event types in the initial complete ontology.

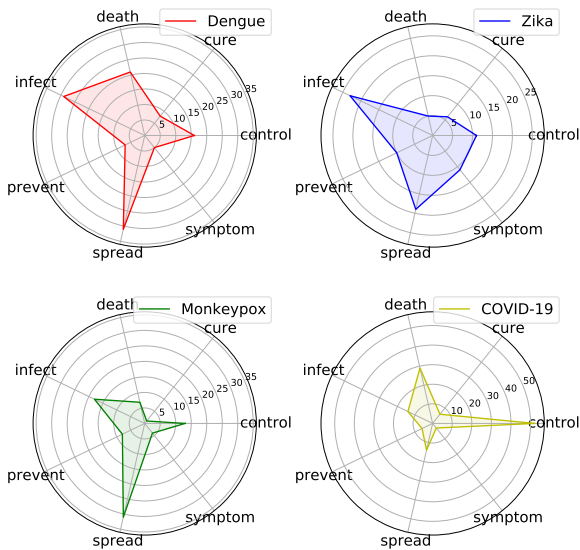


Figure 8: Event type distribution of the disease-related tweets for each disease. Numbers on the axis represent count of mentions for a given event type.

was used for the evaluation of these two sampling techniques.

1046
1047

Model	Tri-I	Tri-C
TRAINED ON UNIFORM DISTRIBUTION		
BERT-QA	58.19	52.30
DEGREE	55.83	52.88
TagPrime	55.48	50.51
DyGIE++	53.22	47.64
Average	55.68	50.83
TRAINED ON RANDOM DISTRIBUTION		
BERT-QA	46.11	43.76
DEGREE	46.11	45.23
TagPrime	25.03	24.15
DyGIE++	51.10	47.35
Average	42.09	40.12

Table 7: Benchmarking ED models trained on uniformly-sampled and randomly-sampled SPEED data on real-distribution based test data of 300 samples.

for our ontology and data, we conduct additional experiments comparing uniform sampling with random sampling. More specifically, we annotate 200 tweets that conform to a ‘real distribution’⁸ based on random sampling and compare the trained models on this data with models trained on 200 tweets of uniform-sampling data. We further annotated 300 tweets based on the ‘real-distribution’ which

⁸Event-based filtering was still applied before sampling.

We present our results in Table 7 averaged over three model runs. We show that in terms of best model performance, uniform sampling is better by 5.5 F1 points compared to random sampling. On average, uniform-sampling trained models outperform the random-sampling trained models by up to 11 points. Both these results prove how despite train-test distribution differences, uniform sampling leads to better training of downstream models.

1048
1049
1050
1051
1052
1053
1054
1055
1056
1057

Generalizability to Other Diseases We also evaluate the models trained on the uniform and random-sampled data for generalizability to other diseases of Monkeypox, Zika, and Dengue. We show the results in Table 8. Clearly, we can see superior generalizability of uniform-sampling trained models as they outperform random-sampling trained models by 37 F1 points for Monkeypox and 28 F1 points for Zika + Dengue. Overall, this result strongly highlights the impact of uniform sampling for robust and generalizable model training.

Model	Monkeypox		Zika + Dengue	
	Tri-I	Tri-C	Tri-I	Tri-C
TRAINED ON UNIFORM SAMPLED DATA				
BERT-QA	56.56	49.30	56.35	46.19
DEGREE	58.35	53.39	58.37	51.27
TagPrime	58.36	53.56	57.05	48.53
DyGIE++	55.73	48.30	56.90	47.10
TRAINED ON REAL SAMPLED DATA				
BERT-QA	9.48	7.97	21.68	20.43
DEGREE	10.76	10.53	19.33	19.00
TagPrime	10.37	8.57	12.78	12.28
DyGIE++	19.59	16.62	26.43	23.40

Table 8: Generalizability benchmarking of ED models trained on 200 samples of uniformly-sampled and randomly-sampled COVID data on other diseases of Monkeypox, Zika, and Dengue.

C Annotation Guidelines and Interface

C.1 Annotation Guidelines

Inspired by [Doddington et al. \(2004\)](#), we develop an extensive set of instructions with tricky cases and examples that have been developed through multiple rounds of expert annotation studies. For our interface, we utilize Amazon Mechanical Turk.⁹ We present the task summary with the major instructions in Figure 14. To reduce ambiguity in trigger selection, we present extensive examples and tricky cases with priority orders as shown in Figure 15. Finally, we also provide a wide range of annotated positive and negative examples as part of the guidelines and show those in Figure 16.

C.2 Annotation Interface

We utilize Amazon Mechanical Turk¹⁰ as the interface for quick annotation. To annotate, annotators can select any word and label it into one of the

⁹<https://www.mturk.com/>

¹⁰<https://www.mturk.com/>

seven pre-defined event types. Event definitions and examples are provided alongside for reference. Each batch (also known as HIT) comprises five tweets for flexibility in annotations. We show the interface and various utilities in Figure 17, 18, and 19 respectively.

D Data Analysis for SPEED

D.1 Benchmarking Test Suites Statistics

We provide the statistics in terms of number of event mentions and tweets for the various benchmarking test suites based on SPEED in Table 9.

	Test Suite	# Mentions	# Tweets
Train	FS-2	14	11
	FS-5	35	24.33
	LR-100	99	67
	LR-200	198	139
	LR-300	306	211
Dev	LR/FS	101	81
Test	All	1,810	1,683

Table 9: Data Statistics for the various benchmarking test suites in terms of number of event mentions and number of tweets. Here, LR-XX represents low resource with XX training event mentions and FS-YY represents few-shot with YY training mentions per event. For FS, we take the average over three different splits of data.

D.2 Event Distribution Analysis

As part of data processing, we attempt to sample tweets in a more uniform distribution between the event types (§ 3.2). In Figure 9, we show the distribution of our dataset in terms of event types. In contrast to tail-ending distributions of other standard datasets like ACE ([Doddington et al., 2004](#)) and MAVEN ([Wang et al., 2020](#)) as shown in Figures 10 and 11 respectively, our distribution of event mentions is more uniform.

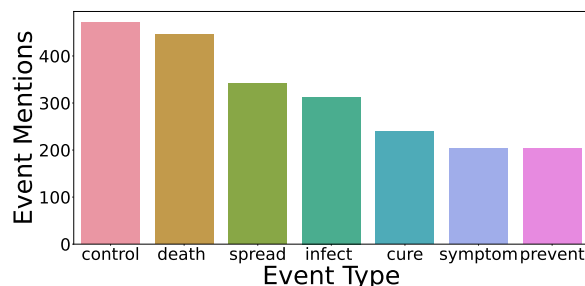


Figure 9: Distribution of event mentions per event type for our dataset SPEED.

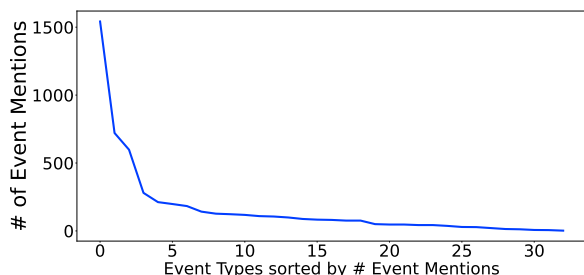


Figure 10: Distribution of event mentions for the event types in the ACE dataset.

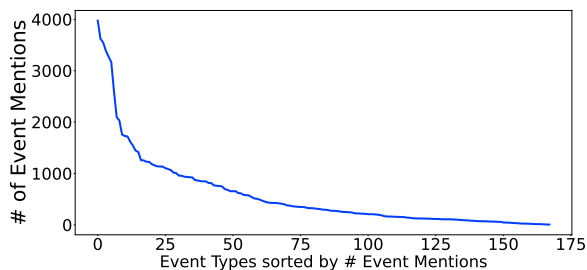


Figure 11: Distribution of event mentions for the event types in the MAVEN dataset.

D.3 Monkeypox Test Data Statistics

We share the data statistics of the evaluation dataset used for Monkeypox in Table 10 split according to each event type. We observe that there is a disparity in distribution across different event types, with *spread* mostly discussed and *cure* and *death* are least discussed.

Event Type	# Event Mentions
infect	78
spread	119
symptom	43
prevent	70
control	62
cure	13
death	13
Total	389

Table 10: Data Statistics for the evaluation dataset for Monkeypox Event Detection categorized by event types.

D.4 Zika + Dengue Test Data Statistics

We share the data statistics of the evaluation dataset used for Zika + Dengue in Table 11 split according to each event type. We observe a more even distribution of event types with more focus on *infect*, *spread*, and *death* well-discussed.

Event Type	# Event Mentions
infect	57
spread	53
symptom	34
prevent	22
control	28
cure	20
death	60
Total	274

Table 11: Data Statistics for the evaluation dataset for Zika+Dengue Event Detection categorized by event types.

E Implementation Details for models

We present the extensive set of hyperparameters and other implementation details about the various ED models we benchmarked in our work.

E.1 BERT-QA

We run our experiments for BERT-QA on an NVIDIA RTX A6000 machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 12.

Pre-trained LM	RoBERTa-Large
Training Batch Size	6
Eval Batch Size	12
Learning Rate	0.001
Weight Decay	0.001
Gradient Clipping	5
Training Epochs	30
Warmup Epochs	5
Max Sequence Length	175
Linear Layer Dropout	0.2

Table 12: Hyperparameter details for BERT_QA model.

E.2 DEGREE

We run our experiments for DEGREE on an NVIDIA RTX A6000 machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 13.

E.3 TagPrime

We run our experiments for TagPrime on an NVIDIA RTX A6000 machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 14.

Pre-trained LM	BART-Large
Training Batch Size	32
Eval Batch Size	32
Learning Rate	0.00001
Weight Decay	0.00001
Gradient Clipping	5
Training Epochs	45
Warmup Epochs	5
Max Sequence Length	250
Max Output Length	20
Negative Samples	15
Beam Size	1

Table 13: Hyperparameter details for DEGREE model.

Pre-trained LM	RoBERTa-Large
Training Batch Size	64
Eval Batch Size	8
Learning Rate	0.001
Weight Decay	0.001
Gradient Clipping	5
Training Epochs	100
Warmup Epochs	5
Max Sequence Length	175
Linear Layer Dropout	0.2

Table 14: Hyperparameter details for TagPrime model.

E.4 DyGIE++

We run our experiments for DyGIE++ on an NVIDIA RTX A6000 machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 15.

E.5 TE

We run our experiments for TE on an NVIDIA 1080Ti machine with support for 8 GPUs. Our hyperparameters are as listed in the original paper (Lyu et al., 2021).

E.6 WSD

We run our experiments for WSD on an NVIDIA A100 machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 16.

E.7 TABS

TABS is an event type induction model, wherein the goal is to discover new event types without a pre-defined event ontology. To adapt this for ED, we follow the end-to-end event discovery setting

Pre-trained LM	RoBERTa-Large
Training Batch Size	6
Eval Batch Size	12
Learning Rate	0.001
Weight Decay	0.001
Gradient Clipping	5
Training Epochs	60
Warmup Epochs	5
Max Sequence Length	200
Linear Layer Dropout	0.4

Table 15: Hyperparameter details for DyGIE++ model.

Pre-trained LM	RoBERTa-Large
Training Batch Size	64
Eval Batch Size	8
Learning Rate	0.00001
Weight Decay	0.01
# Training Epochs	7
Max Sentence Length	128
Max gradient norm	1

Table 16: Hyperparameter details for WSD model.

in (Choi et al., 2022) while making the following modifications: (1) **Dataset Composition:** We utilize ACE (Dodgington et al., 2004) dataset for training and development and our SPEED dataset for testing. Our training data comprises 26 known event types from ACE, the validation set comprises 7 ACE event types, while our test set comprises 7 event types from SPEED. (2) **Candidate Trigger Extraction:** To improve trigger coverage, we extract all nouns and non-auxiliary verbs as candidate trigger mentions. (3) **Evaluation Setup:** Trigger identification (**Tri-I**) F1 score is evaluated using the extracted candidate triggers. For trigger classification (**Tri-C**), we first find the best cluster assignment of the predicted event clusters to the gold event types and then evaluate the F1 score.

We run our experiments for TABS on an NVIDIA RTX 2080 Ti machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 17.

E.8 ETypeClus

For consistency across our evaluations, we follow the re-implementation of the ETypeClus model in (Choi et al., 2022), which consists of the latent space clustering stage of the ETypeClus pipeline and uses the embeddings of trigger mentions to be

Pre-trained LM	BERT-Base
Training Batch Size	8
Eval Batch Size	8
Gradient Accumulation Steps	2
Learning Rate	0.00005
Gradient Clipping	1
# Pretrain Epochs	10
# Training Epochs	30
Consistency Loss Weight	0.2
# Target Unknown Event Types	30

Table 17: Hyperparameter details for TABS model.

the input features. We utilize the contextualized embeddings of the candidate triggers extracted from SPEED for unsupervised training. The candidate trigger extraction process and the evaluation setup are the same as described in § E.7.

We run our experiments for ETypeClus on an NVIDIA RTX 2080 Ti machine with support for 8 GPUs. The major hyperparameters for this model are listed in Table 18.

Pre-trained LM	BERT-Base
Training Batch Size	64
Eval Batch Size	64
Learning Rate	0.0001
Gradient Clipping	1
# Pretrain Epochs	10
# Training Epochs	50
KL Loss Weight	5
Temperature	0.1
# Target Unknown Event Types	30

Table 18: Hyperparameter details for ETypeClus model.

E.9 Keyword

This baseline model basically curates a list of keywords specific to each event and predicts a trigger for a particular event if it matches one of the curated event keywords. Event keywords are curated by expert annotators based on the gold triggers appearing in the SPEED dataset and classified as high confidence, medium confidence, and low confidence based on their occurrence counts and false positive rates (as described in § A.2.¹¹ Although this baseline accesses gold test data, it is meant to be a baseline to provide the upper cap for models of this family.

¹¹We will release the set of keywords with our final code.

This is an event extraction task where the goal is to extract structured events from the text. A structured event contains an event trigger word and an event type.	Task Description
Here are seven events that we are interested in: CONTROL: A CONTROL event are collective efforts trying to impede the spread of a pandemic. INFECT: A INFECT event is the process of a disease or pathogen invading a host or hosts. ... SPREAD: A SPREAD event is the process of a disease spreading or prevailing massively at a large scale.	Ontology and Definitions
Some examples: Input: As the Covid - 19 outbreak spreads at breakneck speed , so does information about the coronavirus . But experts say there ' s a balancing act between sharing findings quickly and taking the time to ensure they ' re scientifically sound . (url) Output: [{"event_type": "SPREAD", "trigger": "spreads"}] Input: signs and symptoms of this phenomenon include fever , rash , abdominal pain , vomiting or diarrhea , along with blood tests showing (url) news headlines & amp ; live updates : A New COVID - 19 Syndrome In Children (url) (url) Output: [{"event_type": "SYMPTOM", "trigger": "symptoms"}] ... Input: We are waiting for the vaccine against the Covid - 19 , when it will be ready ? we need to live in normality . Output: [{"event_type": "PREVENT", "trigger": "vaccine"}]	In-context Examples
Test Sentence: Input: My COVID19 antibodies test came back positive . Crazy . Ive had no symptoms . Please get tested if possible . The more data we have on this the better .	Test Query

Figure 12: Illustration of the prompt used for GPT-3 model. It includes a task description, followed by ontology details of event types and their definitions. Next, we show some in-context examples for each event type and finally, provide the test sentence.

E.10 GPT-3

We use the GPT-3.5 turbo model as the base GPT model. We experiment with ChatGPT (OpenAI, 2021) for tuning our prompts that ensure output consistency. Our final prompt (as shown in Figure 12) comprises a task definition, ontology details, 1 example for each event type, and the final test query. We conducted a looser evaluation for GPT and only match if the predicted trigger text matches the gold trigger text (we didn't check the exact span match basically).

F Predicting Early Warnings for Monkeypox

F.1 Event-wise Analysis

As BERT-QA yields the strongest early warning signal (shown in Figure 6), we conduct an analysis at a more granular level on the contribution of each event type to the early warning signal based on the trained BERT-QA output. We present the results in Figure 13, which leads to the following observations: (1) **Strength of indication varies among event types:** As indicated in Figure 13, event type *infect* and *spread* are strong indicators of the incoming surge in reported cases, while event type *prevent* and *control* can serve as indicators of medium strength. Event type *symptom*, *cure*, and *death* are weak indicators that barely contribute to the early warning signal. (2) **Distribution across**

1237
1238
1239
1240
1241
1242
1243
1244
1245
1246

event types can potentially reveal high-level disease characteristics: We can infer some properties of diseases based on the frequency of mentions about particular events. For example, *death* is less mentioned, which can indicate that *Monkeypox* is less fatal compared to other epidemics like COVID. We would like to mention that these are hypothetical properties based on predictions of our best model (which can be imperfect) and should be taken with a pinch of salt.

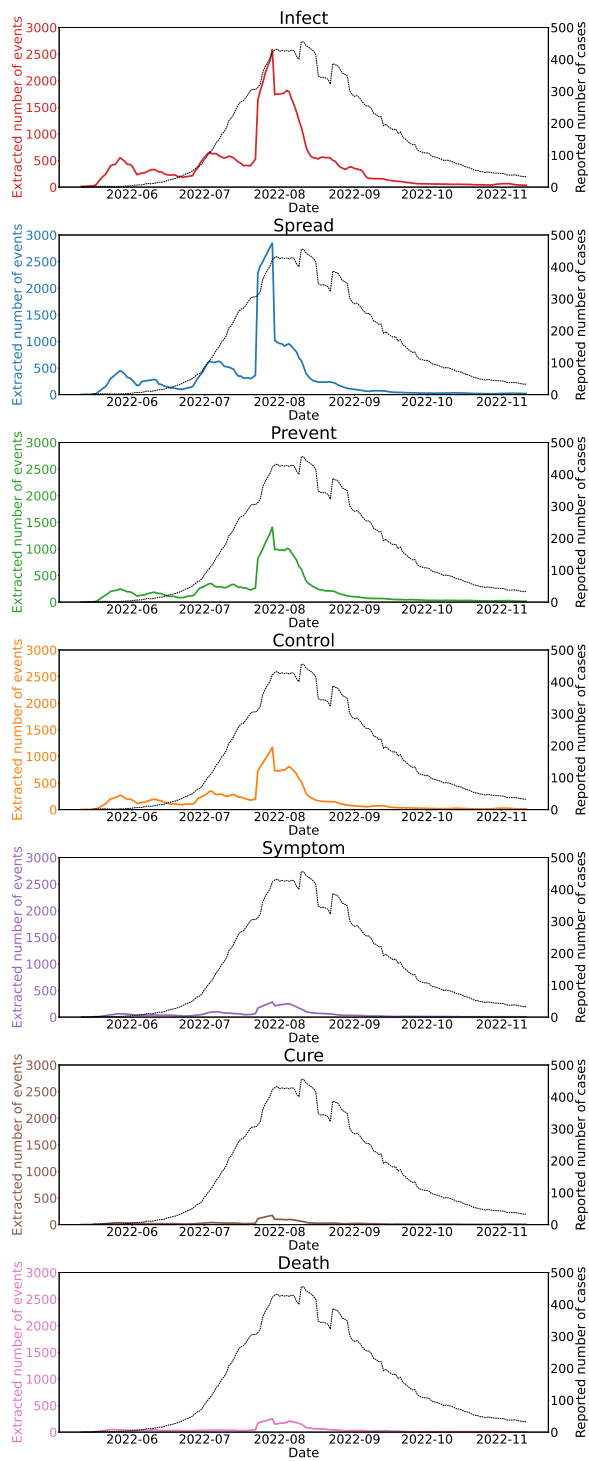


Figure 13: Number of reported Monkeypox cases and the number of extracted events for each SPEED event type from our trained BERT-QA model from XX to XX

An Event is defined as something happens in a sentence. In this task, we are trying to identify whether one or more of the following events exist in a given string: *infect, spread, symptom, prevent, control, cure, and death*. And if an event exist, what is the major **triggering word** that mostly manifest its occurrence.

Event	Definition
infect	The process of a disease/pathogen invading host(s).
spread	The process of a disease spreading/pervailing massively at a large scale.
symptom	Individuals displaying physiological features indicating the abnormality of organisms.
prevent	Individuals trying to prevent the infection of a disease.
control	Collective efforts trying to impede the spread of a pandemic.
cure	Stopping infection and relieving individuals from infections/symptoms.
death	End of life of individuals due to infectious disease.

If there exist any explicit negation of an Event, we say that Event does NOT exist and do not mark it.

Important Notes:

There can be sentences without any events. No need to annotate anything for such sentences.

A trigger word can be linked to one or more events. Choose all possible events in such cases.

Multiple events can be presented in a given sentence. Mark all such events.

The same event can occur multiple times (at different parts) in the same sentence. Mark all occurrences of the event.

You will be able to submit the HIT at the last sentence once you finish annotating all the sentences.

Select "flag" event if you see multiple triggering words or any other tricking situations that needs revisiting, but do not abuse this function.

Figure 14: Task summary and the major annotation guidelines.

Event name	Definition	Action for Final Ontology
SOCIAL SCALE EVENTS		
Prefigure	The signal that precedes the occurrence of a potential epidemic.	Discarded
Outbreak	The process of disease spreading among a certain amount of the population at a massive scale.	Merged into <i>Spread</i>
Spread	The process of disease spreading among a certain amount of the population but at a local scale.	Final Event
Control	Collective efforts trying to impede the spread of a epidemic.	Final Event
Promote	The relationship of a disease driver leading to the breakout of a disease.	Discarded
PERSONAL SCALE EVENTS		
Prevent	Individuals trying to prevent the infection of disease.	Final Event
Infect	The process of a disease/pathogen invading host(s).	Final Event
Symptom	Individuals displaying physiological features indicating the abnormality of organisms.	Final Event
Treatment	The process that a patient is going through with the aim of recovering from symptoms.	Merged into <i>Cure</i>
Cure	Stopping infection and relieving individuals from infections/symptoms.	Final Event
Immunize	The process by which an organism gains immunization against an infectious agent.	Merged into <i>Prevent</i>
Death	End of life of individuals due to infectious disease.	Final Event
MEDICAL SCALE EVENTS		
Cause	The causal relationship of a pathogen and a disease.	Discarded
Variant	An alternation of a disease with genetic code-carrying mutations.	Discarded
Intrude	The process of an infectious agent intruding on its host.	Merged into <i>Infect</i>
Respond	The process of a host responding to an infection.	Discarded
Regulate	The process of suppressing and slowing down the infection of a virus.	Merged into <i>Cure</i>
Transmission route	The process of a pathogen entering another host from a source.	Discarded

Table 19: Complete initial epidemic event ontology comprising 18 event types organized into 3 higher-level abstract categories. We also present details about the event definitions and the action taken for each event type in the final ontology.

Here are more detailed instructions for how to choose the most appropriate triggering word.

Goal: Look for the one word that MOST LIKELY manifests the event's occurrence. You can use the following priority order for annotation:

1. Most of the times, the trigger of the event will be the **main verb** in the sentence.
2. If the verb is ambiguous/vague, the trigger would be a **noun** semantically related to the event.
3. (Rare case) If no such noun exist, the trigger would be any **adjective/adverb** that is related to the event.
4. If still confused, use your best judgement to select the trigger.

In the following illustrations, correct trigger words are marked **blue**.

CASE I : main verb

Example Sentence: "I was coughing and got a fever yesterday and today confirmed I did not get COVID"

Annotation: There are 2 events of symptom

- a. ...**got** a fever...-->Event symptom.
- b. ...was **coughing**... -->Event symptom.
- c. Note 1: "fever" and "COVID" are Not marked as triggering word of the events since the main verbas indicate the event.
Note 2: Here, due to the presence of "and", we have two occurrences of the event symptom.
- d. **Although "get COVID" appears, "not" is the negation emphasizing no infection happens, so event infect does NOT occur**
- e. More examples of main verbs as triggering word:

Example	Event
fight against the pandemic	control
caught a flu	infect
recover from COVID	cure
COVID takes lives	death
prevent infection	prevent
stomach hurts	symptom
number of infection increases	spread

CASE II : nouns

Example Sentence: "Fever, cough, and headache are the most common symptoms of COVID"

Annotation: Here we have 1 event of symptom event:

- a. ...**symptoms** -->Event symptom.
- b. Note: "fever", "cough", and "headache" manifest the symptom event but they are NOT triggering words because "symptom" better manifests the Event.
- c. More examples of nouns as triggering word:

Example	Event
death rate	death
therapy for COVID	cure
infection prevention	prevent
control of spread	control
signs of infection	symptom
spreading of COVID	spread
infection rate	infect

CASE III : adjective

Example Sentence: "I am feverish since 2 days ago"

Annotation: Here we have 1 event of the symptom event

- a. ...**feverish** -->Event symptom.
- b. Note: Here, we do not have a strong verb/noun for marking the trigger. Thus we mark "feverish".
- c. More examples of nouns as triggering word:

Example	Event
get rid of disease	cure
stay cautious against virus	prevent
contagious virus	infect

Figure 15: Guidelines to choose the proper triggering word.

Good Examples

Example 1: "3000+ people are dead due to COVID, so every one please remember to wear a mask and follow the rules to prevent infection and protect our nation from the virus."

Annotation:

- a. **prevent** --> event prevent
- b. **protect** --> event control
- c. **dead**-->event death

Note1: Although "infection" is mentioned, it is prevented, meaning no infection is happening in the sentence, so event infect does NOT exist

Note2: **Do not mark negation of an event.**

Note3: intuitively, people die of COVID must have been infected, but event infect DOES NOT exist here because
An event must be triggered via triggering word and cannot be inferred from another event.

Example 2: "if you ever have a fever, or cough, or have a sore throat, or feel difficult breathing, get tested immediately since you may have been infected."

Annotation:

- a. ...**have** a fever --> event symptom
- b. ...been **infected** --> event infect

Note1: if have more than two parallel phrases triggering an event, only mark the first one instead of all of them.

Note2: event infect has no explicit negation, so event infect exists here.

Bad Examples

Example 1: "Wear a mask"

Wrong annotation:

- a. **wear**-->event prevent

Note1: we may link the action of wearing a mask with pandemic prevention directly, but here it is just an action similar to "read a book" or "eat my lunch".

Note2: if the sentence is instead "wear a mask to prevent COVID." we mark prevent as a triggering word for event prevent instead of "wear"
Look for Events themselves instead of actions/policies related to Events.

Example 2: "Two weeks of quarantine is killing me! May God cure my mind and stop my crazy thoughts."

Wrong annotation:

- a. **killing**-->event death
- b. **cure**--> event cure

Note1: killing does not indicate anybody is dying, and cure does not indicate a therapy against a disease.

Note2: **Do NOT mark hyperbole or rhetorics as Events**

Figure 16: Positive and Negative examples in the annotation guideline.

View instructions

Please read the instructions before attempting the task

Three persons were tested positive for COVID-19 in Karnataka's Dakshina Kannada district. A day after they reached their homes having completed institutional quarantine. (user) (url)

Double click and select the event trigger words present in the above text to in order to annotate, first double click on a word that you think is a potential trigger word for any of the 7 events. Then you would have to choose which event is being triggered by that word

Trigger words and Events:

Submit HIT

Event	Definition	Examples
infect	disease invading host(s). "emphasize infection"	"I have COVID." "High infection rate in U.S. ..."
spread	disease spreading at a large scale. "emphasize dispersion"	"Control the spread of covid." "Infected cases increases ..."
symptom	Individuals displaying abnormal physiological features.	" Symptoms of disease..." "I have a sore throat."
prevent	attempting to avoid infection of a disease. "can be done by individual effort"	"... prevent disease infection." "... protect family from COVID."
control	attempting to control spread of a pandemic. "can NOT be done by individual effort"	"... protect nation from COVID." "... control the spread of flu."
cure	relieving individuals from infections and symptoms.	"This drug can treat smallpox." "I recovered from my illness."
death	End of life of individuals due to infectious disease.	"Stats about Covid death toll..." "The virus kills 50 people..."
flag	Special event used when annotation is ambiguous for some reason.	

Report this HIT | Why Report

Figure 17: Illustration of the default annotation interface on Amazon Mechanical Turk.

amazon mturk
PIPP-Twitter Benchmark B1 (HIT Details) Auto-accept next HIT Requester: Syed Shahriar Hits: 3 Reward: \$0.00 Time Elapsed: 4:01 of 60 Min Return

View instructions Previous Next [Click here to view an exhaustive table](#)

Please read the instructions before attempting the task

Three persons were tested positive for COVID-19 in Karnataka's Dakshina Kannada district a day after they reached their homes having completed institutional quarantine. (user) (url)

select the events that are triggered by "quarantine"

infect spread symptom prevent
 control cure death
 flag

Submit

Double click and select the event trigger words present in the above text to "in order to annotate, first double click on a word that you think is a potential trigger word for any of the 7 events. Then you would have to choose which event is being triggered by that word"

Trigger words and Events:

Submit HIT

Report this HIT | Why Report Return

Event	Definition	Examples
infect	disease invading host(s). "emphasize infection"	"I have COVID." "High infection rate in U.S. ..."
spread	disease spreading at a large scale. "emphasize dispersion"	"Control the spread of covid." "Infected cases increases ..."
symptom	Individuals displaying abnormal physiological features.	"Symptoms of disease..." "I have a sore throat."
prevent	attempting to avoid infection of a disease. "can be done by individual effort"	"...prevent disease infection." "...protect family from COVID."
control	attempting to control spread of a pandemic. "can NOT be done by individual effort"	"...protect nation from COVID." "...control the spread of flu."
cure	relieving individuals from infections and symptoms.	"This drug can treat smallpox." "I recovered from my illness."
death	End of life of individuals due to infectious disease.	"Stats about Covid death toll..." "The virus kills 50 people..."
flag	Special event used when annotation is ambiguous for some reason.	

Figure 18: Illustration of selection of a word within a tweet for annotation in the interface.

amazon mturk
PIPP-Twitter Benchmark B1 (HIT Details) Auto-accept next HIT Requester: Syed Shahriar Hits: 3 Reward: \$0.00 Time Elapsed: 4:40 of 60 Min Return

View instructions Previous Next [Click here to view an exhaustive table](#)

Please read the instructions before attempting the task

Three persons were tested positive for COVID-19 in Karnataka's Dakshina Kannada district a day after they reached their homes having completed institutional quarantine. (user) (url)

Double click and select the event trigger words present in the above text to "in order to annotate, first double click on a word that you think is a potential trigger word for any of the 7 events. Then you would have to choose which event is being triggered by that word"

Trigger words and Events:

Trigger_Word: quarantine
Events: control
Delete Edit

Submit HIT

Report this HIT | Why Report Return

Event	Definition	Examples
infect	disease invading host(s). "emphasize infection"	"I have COVID." "High infection rate in U.S. ..."
spread	disease spreading at a large scale. "emphasize dispersion"	"Control the spread of covid." "Infected cases increases ..."
symptom	Individuals displaying abnormal physiological features.	"Symptoms of disease..." "I have a sore throat."
prevent	attempting to avoid infection of a disease. "can be done by individual effort"	"...prevent disease infection." "...protect family from COVID."
control	attempting to control spread of a pandemic. "can NOT be done by individual effort"	"...protect nation from COVID." "...control the spread of flu."
cure	relieving individuals from infections and symptoms.	"This drug can treat smallpox." "I recovered from my illness."
death	End of life of individuals due to infectious disease.	"Stats about Covid death toll..." "The virus kills 50 people..."
flag	Special event used when annotation is ambiguous for some reason.	

Figure 19: Illustration of the format and options available for a completed annotation in the interface.