Synergizing Unsupervised Episode Detection with LLMs for Large-Scale News Events

Anonymous ACL submission

Abstract

State-of-the-art automatic event detection struggles with interpretability and adaptability to evolving large-scale key events, unlike episodic structures, which excel in these areas. Often 004 overlooked, episodes represent cohesive clusters of core entities (e.g., "protesters", "police") 007 performing actions at a specific time and location. Each key event can be represented as a partially ordered sequence of episodes. This paper introduces a novel task, episode detection, which identifies episodes within a news corpus of key event articles. Detecting episodes 012 poses unique challenges, as they lack explicit temporal or locational markers and cannot be 015 merged using semantic similarity alone. While large language models (LLMs) can aid with these reasoning difficulties, they suffer with 017 long contexts typical of news corpora. To address these challenges, we introduce **EpiMine**, an unsupervised framework that identifies a key event's candidate episodes by leveraging natural episodic partitions in articles, estimated through shifts in discriminative term combinations. These candidate episodes are more cohesive and representative of true episodes, synergizing with LLMs to better interpret and refine them into final episodes. We apply EpiMine to 027 our three diverse, real-world event datasets annotated at the episode level, where it achieves a 59.2% average improvement across all metrics compared to baselines. 031

1 Introduction

Given the saturation of real-time news accessible at our fingertips, reading and processing a key event's critical information has become an increasingly daunting challenge. Consequently, recent work on automatic textual event detection has attempted to integrate the manner in which humans neurologically perceive/store events into textual event detection methods. Specifically, neuroscientists studying event representations in human memory



Figure 1: Example event structure hierarchy. Given a key event node's corpus, detect its episode children and their respective relevant text segments.

find that events are stored in a top-to-bottom hierarchy, as demonstrated in Figure 1. The deeper the hierarchical event level, the more fine-grained its corresponding text granularity (Zhang et al., 2022): we consider a theme as corpus-level (all articles discussing the 2019 Hong Kong Protests), key event as document-level (an article typically discusses a full one to two day key event), episode as segmentlevel, and atomic action as sentence or phrase-level.

043

044

045

051

054

057

060

061

062

063

064

065

066

067

068

069

071

072

Furthermore, neurological research (Baldassano et al., 2017; Khemlani et al., 2015) indicates that events are encoded into memory as *episodic structures*. Representing events as discrete episodes helps us piece together a *coherent and concise narrative* by focusing on *meaningful* clusters of actions, reactions, and developments, rather than examining each in isolation or as a whole. Despite its strengths, **existing automatic event extraction works fail to consider the episode-level**.

For instance, key event detection focuses on identifying "a set of thematically coherent documents" for each key event (Zhang et al., 2022; Liu et al., 2023), but manually parsing large clusters of articles is inefficient and lacks interpretability. Timeline summarization (Steen and Markert, 2019; Li et al., 2021a; Gholipour Ghalandari and Ifrim, 2020; Chen et al., 2023) addresses this by providing dates and compact summaries, yet it suits historical themes better than evolving key events that require finer granularity. Event chain mining (Jiao et al., 2023) takes a fine-granularity approach by

170

171

172

125

126

127

identifying temporally ordered atomic actions, butits phrase-level granularity is often too fine andpractically redundant for large-scale events (e.g., inFigure 1, the actions all describe the same episode).To bridge this gap, we propose the novel task ofepisode detection to pave the way for a more effective event representation.

073

074

079

081

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119 120

121

122

123

124

Episode detection aims to detect episodes from a news corpus containing key event articles. An episode can be described as a cohesive cluster of potentially diverse subjects performing actions at a certain time and location, occurring as part of a larger sequence of episodes under a specific key event. We introduce EpiMine, which detects meaningful episodic events and their corresponding text segments in a large key event corpus, all without any level of human supervision or labeled training data. EpiMine consists of: (1) episode indicative term mining, (2) episode partitioning, (3) LLM-enhanced episode estimation, and (4) episode-segment classification. Collectively, they tackle the unique challenges of episode detection, detailed below:

Challenge 1: Episodes are not timestamped. Key event detection partitions a thematic corpus into document-level clusters by heavily relying on explicit temporal features, like publication dates, being associated with the key event articles (Zhang et al., 2022). However, this assumption fails at the episode-level, where there is no guarantee to have a distinct timestamp associated with each text segment that discusses a new episode. Fortunately, we can take advantage of the idea that journalists naturally partition news articles by sequentially discussing distinct episodes:

Example: An article likely completes its discussion of the episode A, *protesters storming the Legislaive Council*, before episode B, "*protesters vandalized the Legislative Chamber*" (Figure 3).
Hence, to partition articles into distinct episode segments, EpiMine must identify whether two consecutive segments are discussing the same or different episodes– bringing us to our next challenge.

Challenge 2: Episodes contain semantically diverse actions. Each episode features a *set of unique atomic actions*, which can help determine if two segments discuss the same episode. However, for clustering actions, existing methods (Jiao et al., 2023) rely heavily on semantic similarity. This is not realistic for episode-segment clustering:

Example: "protesters spray-painted slogans" and "they unfurled the colonial-era flag" will fall under

the same episode, but are semantically different and unlikely to be clustered.

Alternatively, we can identify salient terms unique to the same episode (episode A: "barriers" and "shoved"; episode B: "*defaced*" and "*walls*"), by exploiting corpus-level signals. For example, if "*defaced*" and "*walls*" are frequently mentioned together across the corpus (or their respective synonyms) and not with other terms, then they are a *discriminative co-occurrence*. When terms between two segments discriminatively co-occur, this indicates the same episode is being discussed. Conversely, if a sufficient shift in term combinations occurs, then a different episode is being discussed.

Challenge 3: Articles often do not feature all episodes. Real-time news reporting often provides an incomplete coverage of multi-day events, with individual articles potentially omitting or partially addressing key episodes. Consequently, while LLMs could assist with the first two challenges, requiring multiple articles hinders their use given their long context limitations (Li et al., 2024; Liu et al., 2024). To address this challenge, EpiMine seeks to select a minimal set of articles that maximizes both the quantity and quality of event episodes. It then merges any article partitions across these articles which likely discuss the same episode and synergizes with an LLM to provide a more fluent interpretation of the candidate episodes, accounting for the episode's core entity, actions, object, location, and time period. This allows EpiMine to finally map the remaining nonsalient article segments to these episodes, pruning any candidates which are not sufficiently supported by the remaining articles. We summarize our core contributions:

- Episode detection: *novel* task to detect episodes & their segments from a key-event corpus.
- **EpiMine**, an unsupervised episode detection method which introduces discriminative term co-occurrence and episode partitioning.
- **Three novel datasets**, reflecting a diverse set of real-world themes and thirty global key events (no key event corpus exists for this task).
- EpiMine outperforms all baselines by, on average, a **59.2% increase across all metrics**.

Reproducibility: We provide our dataset and source code¹ to facilitate further studies.

¹anonymous.4open.science/r/epimine-8782

175

176

177

178

179

181

182

184

186

187

188

190

193

194

195

196

198

199

202

207

208

209

211

212

213

214

215

216

217

218

219

222

2 Related Works

2.1 Event Extraction

Event extraction has been widely studied, focusing on event detection (Liu et al., 2018a; Du and Cardie, 2020; Li et al., 2021b; Lu et al., 2021; Qi et al., 2022; Jiao et al., 2022), event relation extraction (Han et al., 2019; Wang et al., 2020; Ahmad et al., 2021), and salient event identification (Liu et al., 2018b; Jindal et al., 2020; Wilmot and Keller, 2021). Recent work has also addressed event process understanding (Zhang et al., 2020; Chen et al., 2020), though these often rely on expensive expert annotations. Some studies have introduced unsupervised methods to address annotation challenges (Weber et al., 2018; Li et al., 2020). Some overlapping work exists in topic discovery, where (Yoon et al., 2023) proposes unsupervised stream-based story discovery- computing article embeddings based on their shared temporal themes. Recently, large language models have demonstrated powerful general and event extraction-specific reasoning abilities (Pai et al., 2024; Gao et al., 2024).

However, traditional and LLM-driven methods either, (1) focus on phrase/sentence-level events (analogous to actions in Figure 1), or (2) require human-curated event ontologies, often overlooking interpretable, yet meaningful granularities and open-domain texts, which go beyond pre-defined event types. While unsupervised granular event extraction has been explored (Zhang et al., 2022; Jiao et al., 2023) at the document and phrase-level, episode detection is a more interpretable granularity that remains a largely unexplored, yet vital area.

2.2 Timeline Summarization

Timeline summarization (TLS) identifies key dates and concise descriptions for major events. Early methods were extractive, focusing on ranking events for thematic timelines (Nguyen et al., 2014) or using submodular frameworks to model temporal dimensions (Martschat and Markert, 2018). Abstractive methods later emerged, such as sentence clustering and multi-sentence compression (Steen and Markert, 2019). More recent approaches are graph-based, such as event-graph representations for salient sub-graph compression (Li et al., 2021a) and heterogeneous GATs for redundancy reduction (You et al., 2022). While they effectively summarize key events as high-level timelines, they focus on historical themes. Episode-level timelines for ongoing news remain underexplored.

3 Methodology

To tackle episode detection, we propose a novel unsupervised framework, EpiMine. As shown in Figure 2, EpiMine consists of the following four core components: (1) episode indicative term mining, which identifies combinations of salient terms likely to discriminatively co-occur within an episode and not across episodes; (2) episode partitioning, which partitions each article into approximate isolated episodes based on consecutive shifts in the term co-occurrence distribution, (3) LLMenhanced candidate episode estimation, which clusters the top partitions into candidate episodes and utilizes LLM-based reasoning to produce fluent and meaningful episodes, and (4) episode-segment classification, which maps confident segments to their respective episode clusters.

3.1 Preliminaries

3.1.1 Problem Definition

Definition 1 (Episode). An episode E_i is one of a partially ordered sequence of subevents, $\{E_1, \ldots, E_i, \ldots, E_k\}$, of a key (major) event E, where typically $2 \le k \le 20$, and E_i does not overlap with E_j if $i \ne j$. Actions in the episode E_i can be either semantically similar or diverse, but typically have relatively tight time, location, and/or thematic (entities, actions, objects) proximity.

EpiMine aims to extract episodes from a news corpus, where an episode occurs as a significant component of a larger group of episodes that fall under a specific key event. For instance, in Figure 1, without knowing Episode #1, "Protesters stormed the Legislative Council Complex", readers would not fully understand Episode #3, "Police dispersed protesters". Hence, episodes help us understand the overall key event and are especially useful for events that are currently evolving, where finer context is required for sufficiently understanding them.

Definition 2 (Episode Detection). Given a corpus \mathcal{D} about one key event, where each document $d \in \mathcal{D}$ is a news article, the task is to obtain a set of text segment clusters $\mathcal{E} = \{E_1, E_2, \ldots, E_k\}$. Each episode cluster $E_i \subset \mathcal{S} = \{s_1^1, s_2^1, \ldots, s_{|d|}^{|\mathcal{D}|}\}$, where \mathcal{S} contains all the text segments identified in each document $d \in \mathcal{D}$, and every two clusters do not have overlapping text segments (i.e., $E_i \cap E_j = \emptyset$ for $i \neq j$).

It is important to note that k, the number of episodes, is not known in advance and oftentimes,

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

267

269

270

271



Figure 2: We detail the overall framework of EpiMine.

a news article segment may discuss either episodes of a different key event (e.g., an episode with similar aspects that occurred in a different historical key event) or multiple episodes of the current key event. Nonetheless, our goal is to detect the most relevant episodes to the current key event at hand and consequently mine the most distinctive text segments for each of these (hence our constraint of $E_i \cap E_j = \emptyset$ for $i \neq j$).

272

273

276

277

279

281

284

290

292

3.2 Episode Indicative Term Mining

Hong Kong protesters **broke** through the final **barrier** separating them and the interior of the city's Legislative Council Complex. Many of the **glass panels** lining the building's exterior had already been <u>smashed</u> in the afternoon, but it took another hour or so before <u>protesters</u> were able to **breach** the **metal shutters**.

Within minutes, <u>protesters</u> began **spray-painting slogans** onto the **corridor walls** and **vandalizing** the **portraits** of <u>Legislative Council</u> presidents. They also **draped** the **flag** of colonial-era <u>Hong Kong</u> at the **podium** of the legislative **chamber**.

Figure 3: Natural partition between two episodes in a key event article. An episode's discriminative terms are **bolded**; salient non-discriminative terms are <u>underlined</u>.

Lacking supervision, our goal is to identify *potential* candidates for episodes. Episodes are often described in relation to each other and usually lack timestamps or locations consistently mentioned within their segments. For example, the phrase "police dispersed protesters" may not have a precise timestamp because it is a *response* to "protesters stormed the Legislative Council Complex," and some journalists may consider the implicit ordering adequate. Additionally, the same episode can be described using different entities and actions– journalists may *report different perspectives*. For example, both "protesters shoved against the barricades" and "the police used pepper-spray on the protesters" describe the episode "protesters stormed the Legislative Council Complex". However, they are semantically different, focused on different core entities and actions. Thus, we cannot depend on a consistent subject-action-object triple or an explicit time/location mapped to each episode in the article. 293

294

295

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

329

To circumvent this challenge, we exploit the idea that journalists naturally partition news articles according to episodes, forming episode fragments. For example, as shown in Figure 3, an article will likely complete its discussion of episode #1, "Protesters stormed the Legislative Council Complex" (red), before fully shifting to discussing episode #2, "protesters vandalized the Legislative Chamber" (blue). Across these episode fragments, certain salient terms are featured (e.g., protesters, legislative, vandalizing, podium). We adapt the idea of event salience from (Jiao et al., 2023) specifically for the task of episode detection, allowing us to identify terms which are (1) distinct and significant to understanding a given key event, and (2) frequently found in a key event's segments and infrequently in other background/general articles. Thus, we identify a set of salient terms for each segment within the corpus (details in Appendix E).

Discriminative Co-occurrence. In Figure 3, we can see that the first episode fragment, and Episode #1 in general, features a combination of similar terms, such as "protesters", "barrier", and "breach". Likewise, the second episode may include a combination of terms similar to "protesters", "spraypainting", and/or "flag". We note that despite some journalists choosing to only describe the protesters spray-painting, while others focus on the protesters

422

375

draping the colonial-era flag, we must be able to recognize that their respective salient terms are likely to co-occur within the same episode.

However, we make a novel distinction between a co-occurrence and a discriminative cooccurrence. A salient term a (e.g., "protesters") 335 may often co-occur with a salient term b ("spray-336 painting") within an episode. However, if a also frequently co-occurs with many other terms in various episodes ("protesters broke"), a and its co-occurrences are less useful for distinguishing 340 episodes. Thus, (a, b) is not a discriminative co-341 occurrence. 342

> **Definition 3 (Discriminative Co-occurrence).** A pair of terms (a, b) discriminatively co-occur if (1) they frequently appear together in episode E_i , and (2) neither a nor b appear as frequently with other terms w in other episodes $E_{\notin i}$.

348

351

358

364

367

371

We compute the discriminative occurrence d between salient term pair (a, b) using the following:

$$\mathbf{d}(a,b) = \log\left(\frac{freq(a,b)}{\max(\bar{f}_a,\bar{f}_b)}\right) \times \log\left(\frac{|T|}{\max(|F_a|,|F_b|)}\right),$$

where $\bar{f}_a = \frac{1}{|T|} \sum_{\forall w_i \in T} freq(a,w_i)$, and
 $F_a = \{freq(a,w_i) > 1 \ \forall \ w_i \in T\}$
(1)

The first log term ensures that the pair's cooccurrence (freq(a, b)) is statistically significant $(\geq$ the max of a and b's mean vocabulary-wide co-occurrence respectively). The second log term ensures the pair is a discriminative match, penalizing cases where a or b frequently co-occurs with a large portion of the salient term set T. For example, co-occurrences with "protesters" are not discriminative because "protesters" is a core entity in all episodes and thus frequently co-occurs with many terms in T. In contrast, ("slogans", "flags") is a discriminative co-occurrence since both terms frequently appear together in segments discussing episode #2 and rarely co-occur with other terms $w_i \in T$. If a and b are the same term or close synonyms (determined by statistically significant semantic similarity), they have maximum co-occurrence. By leveraging multiple articles in a large key event corpus, we have sufficient statistical support to ensure our output reflects the average realistic reporting of the key event and its episodes.

3.3 **Episode Partitioning**

With the ability to identify discriminative cooccurrences, we can use a key transitive property 374

articles, where not all combinations of an episode's 376 discriminative terms explicitly co-occur: If (a, b) and (b, c) are both discriminative cooccurrences, then (a, c) is *also likely* to be a discriminative co-occurrence. To illustrate this, we have the following text segment excerpts of a news article (the salient and discriminative terms are *italicized*): 1. Protesters defaced the Hong Kong emblem, spray-painted slogans, and unfurled the flag. 2. The portrait of LegCo president was defaced. 3. A slogan on the wall reads: "The government forced us to revolt". 4. Police said at least 13 people had been arrested on suspicion of involvement in the prodemocracy protest. We can naturally see that segments 1-3 all discuss the "protesters vandalized the Legislative Chamber" episode, while segment 4 discusses the "police dispersed protesters" episode. We can systematically replicate this partitioning process by considering the discriminative co-occurrence score between all pairwise combinations of terms from segments (i - 1) and (i). If the average discriminative co-occurrence and static semantic similarity between each term a from (i - 1) and b from (i) is statistically significant ($\geq \mu_d - \sigma_d$) for that specific article d (e.g., notably (slogans, defaced) for segments 1-3), we hypothesize that the same episode is being discussed and *merge* them into one episode fragment. If not (e.g., (slogans, arrested) for segments 3-4), this indicates that a different episode is being discussed, and we partition them into two

to resolve episode co-references within and across

3.4 LLM-Enhanced Episode Estimation

are provided in Appendix G.

episode fragments. Further implementation details

LLMs demonstrate strong event-specific reasoning at the phrase or sentence level (Pai et al., 2024; Gao et al., 2024), but they struggle with understanding long contexts (Li et al., 2024; Liu et al., 2024). This limitation hampers their ability to process all episode fragments for detecting episodes. Additionally, noisy retrieval significantly affects reasoning performance (Shen et al., 2024). To address these challenges, we propose a synergistic approach that enhances in-context episode reasoning by reducing the number of required fragments while improving

their cohesiveness and quality. We first identify
the set of articles that maximizes the *quantity and quality* of potential episodes, where each article is
ranked by multiplying two metrics:

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

- 1. *Quality of episode fragments*: A top article should primarily consist of episode fragments containing salient terms that discriminatively co-occur. This reduces the rank of general fragments which summarize/analyze the event. We average each episode fragment's mean inner-discriminative co-occurrence (across all pairwise combinations of its salient terms).
 - Quantity of episode fragments: A top article should ideally contain all ground-truth episodes. Therefore, we take the log of the number of episode fragments in the article.

After ranking all articles, we select the top $\delta\%$ and resolve potential co-references to the same episode across these top articles. We apply agglomerative clustering (Murtagh and Contreras, 2012) to the top episode fragments using a pre-computed distance matrix. The distance between two fragments (inversed) is calculated using the same discriminative and static semantic similarity score used in Section 3.3). Clusters with a statistically insignificant number of episode fragments are pruned.

Finally, we provide episode fragment clusters as a more interpretable context for the LLM to resolve two challenges: (1) missing time and location stamps in fragments, and (2) semantic inconsistencies within clusters. The LLM summarizes each cluster by identifying its core attributes– entities, actions, objects, location, and time. It then outputs the *episode attributes, relevant keywords* for extraction, and the top *extracted text segments* (prompt & example in Appendix I).

3.5 Episode-Segment Classification

With these core summaries of the episode clusters, we obtain a generalized description of each candidate episode. For each candidate, we encode its LLM-based core attributes and extracted segments to compute a simple episode representation. We use these to assign an episode and confidence score to each encoded input segment. Further details on our encoding process are provided in Appendix H.

Episode-Segment Confidence Estimation. Directly mapping a text segment to its top episode based on cosine similarity risks misclassifying episode-irrelevant segments or those discussing multiple episodes (e.g., a journalist's summary). To avoid classifying such segments and ensure nonoverlapping episode clusters (as discussed in Section 3.1.1), we must determine the confidence of a segment discussing a single episode. 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

We compute segment s_i 's cosine similarity to its top two episodes $(e_i^0 \text{ and } e_i^1)$. A larger gap $(e_i^0 - e_i^1)$ reflects greater confidence in classifying s_i to e_i^0 . Each gap is normalized by the sum of all segment-episode gaps across the corpus, ensuring confidence is relative to the key event:

$$s_{i,\text{confidence}} = \frac{e_i^0 - e_i^1}{\sum_{l=1}^{|\mathcal{S}|} (e_l^0 - e_l^1)}$$
(2)

Segments with statistically significant confidence in their top episode are assigned to their respective episode clusters E_i . Episodes with no assigned segments are pruned, yielding the **final detected episodes and clusters**, \mathcal{E} .

4 Experiments

For implementing **EpiMine**, we use the following hyperparameters across all datasets: $\delta = 25\%$, $sim_thresh = 0.75$. We also use Claude-2.1 as our base LLM (A). All other hyperparameters are set to their respective default values. We provide all experimental settings in Appendix A.

Table 1: Statistics of our collected datasets. The numbers are averaged per key event.

Theme	# docs	# episodes	# segments
Terrorism/Attacks	32.2	5.9	290.3
Natural Disasters	36.2	7.4	324.6
Political Events	70.2	7.5	667.7

4.1 Datasets

We conduct our experiments on three novel thematic, real-world news corpora selected from Wikipedia² over the last decade. For each theme, we manually collect approximately 10 key events composed of multiple articles and ensure that distinct *episodes* exist in each. The articles are obtained from the Wikipage references of each key event– filtered with constraints in time, language, and relevance. Further details on the criteria/process, each theme, and corresponding key events are in Appendices D and K. We also segment each article to match our setting (App. F).

²https://en.wikipedia.org/wiki/

	Terr	orism (5.36	eps)	Natural Disasters (7.4 eps)			Politics (7.5 eps)		
Methods	5-prec	5-recall	5-F1	5-prec	5-recall	5-F1	5-prec	5-recall	5-F1
EMiner	8.64	0.25	0.48	10.37	0.19	0.37	8.66	0.16	0.32
K-means	21.23	21.23	21.23	27.85	28.47	$\frac{28.14}{22.00}$	16.04	16.04	16.04
K-means + A	28.58	14.04	18.26	37.40	16.58		27.18	17.36	<u>18.25</u>
EvMine	23.03	15.02	17.45	28.15	8.02	12.25	5.36	4.00	4.58
EvMine + A	37.88	15.70	21.33	43.56	13.22	19.40	32.73	12.98	17.28
EpiMine (A)	71.21	<u>22.07</u>	<u>32.43</u>	70.98	<u>28.46</u>	34.53	62.67	<u>21.54</u>	<u>29.23</u>
- No Confidence (A)	<u>61.97</u>	30.19	38.45	<u>43.66</u>	20.78	27.76	<u>60.29</u>	27.73	24.77
- No LLM	37.73	21.62	24.77	37.19	14.78	17.52	<u>30.64</u>	23.51	19.06

Table 3: Compares top-5 salient terms which (1) have the highest cosine-sim (CS) and (2) discriminative cooccurrence (DC), with the given keyword.

Keyword	CS	DC		
broke	stormed, ransacked,	glass, doors, metal, build-		
	dashed, occupied, rushed	ing, teargas		
slogans	spray, placards, painted, defaced, pictures	reads, wall, damage, started, portraits, spray		

4.2 Baselines

509

510

511

513

514

515

517

518

519

520

522

523

526

528

530

532

536

We compare against the following methods using the evaluation metrics (App. J): (1) Kmeans (Likas et al., 2003): given the # of groundtruth episodes, it clusters segments using ST (Reimers and Gurevych, 2019) embeddings; (2) **EvMine** (Zhang et al., 2022): a document-level unsupervised key event detection method adapted to segment level for episode detection; (3) EMiner (Jiao et al., 2023): unsupervised event chain miner that clusters atomic actions, adapted to episodes; (4) No Confidence: an ablation that uses max cosine-similarity instead of confidence from Equation 2; (5) No LLM: an ablation that uses estimated episode clusters from Section 3.4 to compute our episode representations directly. We also integrate A into K-means and EvMine using our same prompt (Appendix I). All baseline and ablation details are in Appendix B.

4.3 Overall Results & Analysis

In Table 2, EpiMine shows an average **80.8%** increase in 5-precision, a **34.0%** increase in 5-recall, and a **62.8%** increase in 5-F1 over all baselines. Notably, despite both K-means and K-means + A being *given the ground-truth number of episodes*, they are **significantly outperformed by EpiMine** (both the base model and no confidence ablation). Additionally, EvMine and EMiner, originally de-

signed for key event and atomic action levels of event granularity, **fail to address the unique challenges of episode detection**. We further analyze our results through extensive quantitative and qualitative studies, including a detailed case study on the "2019 Hong Kong Legislative Protest" (as shown in Figure 1), leading to the following takeaways: 537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

1. LLMs require effective episode fragment clusters for synergistic episode estimation. As shown in Table 4, LLMs without any initial clusters as guidance (\mathbb{A} , GPT-4³) fail to detect high-quality episodes, miss most ground-truth episodes, and include irrelevant atomic actions (e.g., "Brian Leung pulls off mask"). Similarly, using low-quality baseline clusters results in poor performance. EvMine detects episodes that all reflect the same event, "Protesters vandalized the Legislative Chamber". While K-means produces more distinct episodes, it does not capture the most critical, gold episodes. In contrast, EpiMine's episodes are both distinct and meaningful, attributed to its cluster quality (quantitatively confirmed by EpiMine-No LLM's competitive performance).

EpiMine's clusters also elicit the LLM to identify more meaningful temporal information. Unlike most baseline episodes which have "July 1, 2019" as the time attribute, EpiMine's episodes feature more descriptive temporal cues: "after breaking in", "after midnight", "in a news conference at 4 am on July 2". Moreover, EpiMine's "incorrect" episode #4 is a significant sub-event of the key event discussed by many articles. This *strongly demonstrates the impact of EpiMine's candidate episode clusters* as input into the LLM; LLMs **alone cannot perform quality episode detection**.

³https://chat.openai.com/

Table 4: Gold and detected episodes (a maximum of five are included for brevity) for the "2019 Hong Kong Legislative Protests" key event. We specify the gold/detected episode attributes for each episode cluster in the following semicolon-separated format: core entity; action; object; time; location. "Not detected" denotes that no more episodes were generated by the model. We note the number of detected episodes beside the model name.

Model	Episode #1	Episode #2	Episode #3	Episode #4	Episode #5
Gold (5 eps)	Activists; headed; towards the Legislative Council Complex; 1 July 2019; Hong Kong	Protesters; stormed; the Legisla- tive Council Complex; around 9:00 pm; Hong Kong;	Protesters; damaged/defaced; portraits, furniture, emblem, etc.; 1 July 2019; Legislative Council Complex	Police; started using; tear gas to disperse protesters; 12:05 am 2 July; around the Legislative Council com- plex	Police; arrested; individuals in connection with the inci- dent; between 3 July and 5 July; Hong Kong
K-means + A (4 eps)	Protesters; storm and vandal- ize; Legislative Council build- ing; July 1, 2019; Legisla- tive Council complex in Ad- miralty, Hong Kong	Hong Kong government; con- demns; protesters storming leg- islative building; July 1, 2019; Hong Kong	Hong Kong protesters; ex- press; demands for freedom and democracy; July 1, 2019; Hong Kong Legislative Council	Hong Kong police; adopt; more restrained tactics; July 1, 2019; Hong Kong Leg- islative Council	Not detected
EvMine + A (4 eps)	Protesters; vandalize; Hong Kong legislative building; July 1, 2019; Hong Kong legisla- tive building	Protesters; occupy and vandal- ize; Hong Kong legislative cham- ber; July 1, 2019; Hong Kong legislative building	Protesters; spray paint; slogans and demands; July 1, 2019; Hong Kong legislative building	Protesters; deface; Hong Kong emblem; July 1, 2019; Hong Kong legisla- tive building	Not detected
Claude (3 eps)	Protesters; storm; Hong Kong's Legislative Council; July 1, 2019; Hong Kong's Legislative Council building	Police; retreat and avoid con- frontation; protesters storming Hong Kong's Legislative Coun- cil; July 1, 2019; Hong Kong's Legislative Council building	Brian Leung Kai-ping; pulls off mask and reads protesters' de- mands; inside Hong Kong's Leg- islative Council; July 1, 2019; Legislative Council chamber	Not detected	Not detected
GPT-4 (2 eps)	Hong Kong protesters; storm Legislative Council; govern- ment and police; July 1, 2019; Legislative Council Complex, Hong Kong	Hong Kong citizens; march against extradition bill; "Carrie Lam and Chinese government; June 2019; Various locations in Hong Kong	Not detected	Not detected	
EpiMine (7 eps)	protesters; broke into and oc- cupied; Hong Kong's legisla- tive building; July 1, 2019; Hong Kong	protesters; vandalized; the leg- islative building; after breaking in; Hong Kong	police; fired tear gas at; protesters; after midnight on July 1; outside the legislative building	Carrie Lam; condemned; the protesters' actions; in a news conference at 4am on July 2; Hong Kong	police; began making ar- rests of; protesters involved; in the days after; Hong Kong

2. The strengths of discriminative cooccurrence complement those of cosine similarity. Table 3 illustrates the qualitative strengths of our novel discriminative co-occurrence metric. Both cosine similarity (CS) and discriminative cooccurrence (DC) offer different, complementary strengths. CS identifies similar words that play a similar role or are synonyms within an episode (e.g., "broke", "ransacked"), while DC identifies the key surrounding actions and objects that cooccur within the same episode (e.g., "slogans", "wall"). This is quantitatively supported by our ablations (Table 5; Appendix C), which show a significant decrease in the quality of our top articles (partitioned into episodes) without our salience and discriminative co-occurrence measures.

573

574

575

576

578

580

582

586

587

588

589

590

592

596

3. Fragment ranking identifies top articles. In Figure 4, we conduct a sensitivity analysis of the top articles chosen to estimate candidate episodes (Section 3.4). We compare the gold episodes contained in the set of the top $\delta\%$ articles as we vary δ . By ranking the articles based on their likelihood of containing both high-quality and numerous episodes, we find that *EpiMine's top article selection covers the vast majority of episodes* by $\delta = 25\%$ and more comprehensively around $\delta = 45\%$. This is significant as it helps *minimize both the noise and the amount of data* needed to accurately detect all episodes.



Figure 4: Percentage of key event's gold episodes captured in the $\delta\%$ top articles chosen during the candidate episode estimation. Results averaged across themes.

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

5 Conclusion

In this work, we proposed EpiMine, a novel, unsupervised episode detection method for large-scale news events. EpiMine performs (1) episode indicative term mining- identifying combinations of salient terms that are likely to discriminatively cooccur within an episode and not across episodes, (2) episode partitioning, which partitions each article into approximate isolated episodes, (3) LLMenhanced episode estimation, which clusters the top partitions into candidate episodes and synergizes with an LLM to produce fluent and meaningful episodes, and (4) episode-segment classification, which maps confident segments to their respective episode clusters. EpiMine significantly outperforms all baselines on the vast majority of key events, as shown through extensive quantitative and qualitative analysis.

629

630

635

647

650

653

656

664

665

6 Limitations & Future Work

While EpiMine serves as an intuitive, unsupervised framework which demonstrates a more interpretable granularity for event analysis (episodes), it contains a few limitations that form the foundation for future, impactful research areas.

We note that our confidence metric influences EpiMine to be more conservative in its episodesegment classification. The "No-Confidence" ablation in Table 2 of Appendix C shows that using confidence (Equation 2) improves EpiMine's precision but reduces recall. This conservative approach results in EpiMine being cautious when assigning segments to episode clusters, excluding segments with insufficient confidence. This indicates that our method relies on both our clustering method and confidence scoring for precise episode detection. EpiMine users can determine if they prefer a confidence or no-confidence method depending on their use-case (higher precision versus higher recall). Nonetheless, both versions still outperform all baselines.

We also note that the key event theme also has an impact on EpiMine's performance. Specifically, natural disaster episodes are typically sequential and semantically distinct: disaster begins \rightarrow warning \rightarrow evacuation \rightarrow damage/deaths \rightarrow relief. As K-means is uniquely given k, the number of episodes, and relies on semantic similarity, it performs well with distinct episodes. However, we still see that its reliance on surface-level semantics leads to lower precision.

Finally, in the ablation studies shown in Table 5 of Appendix C, the politics dataset does show slight improvements in precision and F1 when *discriminative co-occurrence is replaced*, due to more term overlap across episodes, resulting in less distinct, sequential episodes.

Further work towards the temporal analysis of episodes within articles can be explored, as well as extending our work to primarily multilingual news settings with low resources.

7 Ethics Statement

Based on our current methodology and results, we do not expect any significant ethical concerns, given that subtasks like episode detection within the news event extraction and analysis is a standard problem domain across data mining applications. Furthermore, having the method rely on zero supervision helps as a barrier to any user-inputted biases.

However, one minor factor to take into account is any hidden biases that exist within the large language models used as a result of any potentially biased data that they were trained on. We used these pre-trained language models for refining the fluency of the detected episode clusters and did not observe any concerning results, as it is a low-risk consideration for the domains that we studied. 669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

709

710

711

712

713

714

715

716

717

718

719

720

721

722

References

- Wasi Uddin Ahmad, Nanyun Peng, and Kai-Wei Chang. 2021. GATE: graph attention transformer encoder for cross-lingual relation and event extraction. In Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021, pages 12462–12470. AAAI Press.
- Christopher Baldassano, Janice Chen, Asieh Zadbood, Jonathan W Pillow, Uri Hasson, and Kenneth A Norman. 2017. Discovering event structure in continuous narrative perception and memory. *Neuron*, 95(3):709– 721.
- Muhao Chen, Hongming Zhang, Haoyu Wang, and Dan Roth. 2020. What are you trying to do? semantic typing of event processes. In *Proceedings* of the 24th Conference on Computational Natural Language Learning, CoNLL 2020, Online, November 19-20, 2020, pages 531–542. Association for Computational Linguistics.
- Xiuying Chen, Mingzhe Li, Shen Gao, Zhangming Chan, Dongyan Zhao, Xin Gao, Xiangliang Zhang, and Rui Yan. 2023. Follow the timeline! generating an abstractive and extractive timeline summary in chronological order. *ACM Transactions on Information Systems*, 41(1):1–30.
- Freddy Y. Y. Choi. 2000. Advances in domain independent linear text segmentation. In *1st Meeting of the North American Chapter of the Association for Computational Linguistics*.
- Xinya Du and Claire Cardie. 2020. Event extraction by answering (almost) natural questions. In *Proceedings* of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 671–683. Association for Computational Linguistics.
- Jun Gao, Huan Zhao, Wei Wang, Changlong Yu, and Ruifeng Xu. 2024. Eventrl: Enhancing event extraction with outcome supervision for large language models. *arXiv preprint arXiv:2402.11430*.
- Demian Gholipour Ghalandari and Georgiana Ifrim. 2020. Examining the state-of-the-art in news timeline summarization. In *Proceedings of the 58th Annual*

835

836

780

Meeting of the Association for Computational Linguistics, pages 1322–1334, Online. Association for Computational Linguistics.

723

724

727

733

734

735

736

737

740

741

742

743

744

745

746

747

748

749

750

751

758

759

763

769

770

774

775

776

777

778

779

- Rujun Han, I-Hung Hsu, Mu Yang, Aram Galstyan, Ralph M. Weischedel, and Nanyun Peng. 2019. Deep structured neural network for event temporal relation extraction. In Proceedings of the 23rd Conference on Computational Natural Language Learning, CoNLL 2019, Hong Kong, China, November 3-4, 2019, pages 666–106. Association for Computational Linguistics.
- Yizhu Jiao, Sha Li, Yiqing Xie, Ming Zhong, Heng Ji, and Jiawei Han. 2022. Open-vocabulary argument role prediction for event extraction. In *Findings of the Association for Computational Linguistics: EMNLP* 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022, pages 5404–5418. Association for Computational Linguistics.
- Yizhu Jiao, Ming Zhong, Jiaming Shen, Yunyi Zhang, Chao Zhang, and Jiawei Han. 2023. Unsupervised event chain mining from multiple documents. In *Proceedings of the ACM Web Conference 2023*, pages 1948–1959.
- Disha Jindal, Daniel Deutsch, and Dan Roth. 2020. Is killed more significant than fled? A contextual model for salient event detection. In *Proceedings of the* 28th International Conference on Computational Linguistics, COLING 2020, Barcelona, Spain (Online), December 8-13, 2020, pages 114–124. International Committee on Computational Linguistics.
- Priyanka Kargupta, Tanay Komarlu, Susik Yoon, Xuan Wang, and Jiawei Han. 2023. MEGClass: Extremely weakly supervised text classification via mutuallyenhancing text granularities. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 10543–10558, Singapore. Association for Computational Linguistics.
- Sangeet S Khemlani, Anthony M Harrison, and J Gregory Trafton. 2015. Episodes, events, and models. *Frontiers in human neuroscience*, 9:590.
- Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2024. LooGLE: Can long-context language models understand long contexts? In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16304–16333, Bangkok, Thailand. Association for Computational Linguistics.
- Manling Li, Tengfei Ma, Mo Yu, Lingfei Wu, Tian Gao, Heng Ji, and Kathleen McKeown. 2021a. Timeline summarization based on event graph compression via time-aware optimal transport. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 6443–6456, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Manling Li, Qi Zeng, Ying Lin, Kyunghyun Cho, Heng Ji, Jonathan May, Nathanael Chambers, and Clare R. Voss. 2020. Connecting the dots: Event graph

schema induction with path language modeling. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 684–695. Association for Computational Linguistics.

- Sha Li, Heng Ji, and Jiawei Han. 2021b. Documentlevel event argument extraction by conditional generation. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2021, Online, June 6-11, 2021, pages 894–908. Association for Computational Linguistics.
- Aristidis Likas, Nikos Vlassis, and Jakob J Verbeek. 2003. The global k-means clustering algorithm. *Pattern recognition*, 36(2):451–461.
- Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. *Transactions of the Association for Computational Linguistics*, 12:157–173.
- Xiao Liu, Zhunchen Luo, and Heyan Huang. 2018a. Jointly multiple events extraction via attention-based graph information aggregation. In *Proceedings of the* 2018 Conference on Empirical Methods in Natural Language Processing, pages 1247–1256, Brussels, Belgium. Association for Computational Linguistics.
- Zheng Liu, Yu Zhang, Yimeng Li, and Chaomurilige. 2023. Key news event detection and event context using graphic convolution, clustering, and summarizing methods. *Applied Sciences*, 13(9):5510.
- Zhengzhong Liu, Chenyan Xiong, Teruko Mitamura, and Eduard H. Hovy. 2018b. Automatic event salience identification. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 -November 4, 2018, pages 1226–1236. Association for Computational Linguistics.
- Yaojie Lu, Hongyu Lin, Jin Xu, Xianpei Han, Jialong Tang, Annan Li, Le Sun, Meng Liao, and Shaoyi Chen. 2021. Text2Event: Controllable sequence-tostructure generation for end-to-end event extraction. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 2795–2806, Online. Association for Computational Linguistics.
- Sebastian Martschat and Katja Markert. 2018. A temporally sensitive submodularity framework for timeline summarization. In *Proceedings of the 22nd Conference on Computational Natural Language Learning*, pages 230–240, Brussels, Belgium. Association for Computational Linguistics.
- Yu Meng, Jiaxin Huang, Guangyuan Wang, Zihan Wang, Chao Zhang, Yu Zhang, and Jiawei Han. 2020. Discriminative topic mining via category-name guided

943

944

945

946

text embedding. In *Proceedings of The Web Conference 2020*, pages 2121–2132.

837

838

841

847

848

851

852

853

857 858

870

871

872

877

878

882

884

887

- Fionn Murtagh and Pedro Contreras. 2012. Algorithms for hierarchical clustering: an overview. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 2(1):86–97.
- Kiem-Hieu Nguyen, Xavier Tannier, and Véronique Moriceau. 2014. Ranking multidocument event descriptions for building thematic timelines. In COL-ING 2014, the 25th International Conference on Computational Linguistic, pages 1208–1217.
- Liu Pai, Wenyang Gao, Wenjie Dong, Lin Ai, Ziwei Gong, Songfang Huang, Li Zongsheng, Ehsan Hoque, Julia Hirschberg, and Yue Zhang. 2024. A survey on open information extraction from rule-based model to large language model. *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9586–9608.
 - Zheng Qi, Elior Sulem, Haoyu Wang, Xiaodong Yu, and Dan Roth. 2022. Capturing the content of a document through complex event identification. In *Proceedings of the 11th Joint Conference on Lexical and Computational Semantics*, pages 331–340, Seattle, Washington. Association for Computational Linguistics.
 - Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence embeddings using Siamese BERTnetworks. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3982–3992, Hong Kong, China. Association for Computational Linguistics.
 - Xiaoyu Shen, Rexhina Blloshmi, Dawei Zhu, Jiahuan Pei, and Wei Zhang. 2024. Assessing "implicit" retrieval robustness of large language models. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8988– 9003, Miami, Florida, USA. Association for Computational Linguistics.
 - Julius Steen and Katja Markert. 2019. Abstractive timeline summarization. In Proceedings of the 2nd Workshop on New Frontiers in Summarization, pages 21– 31, Hong Kong, China. Association for Computational Linguistics.
 - Haoyu Wang, Muhao Chen, Hongming Zhang, and Dan Roth. 2020. Joint constrained learning for eventevent relation extraction. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 696–706. Association for Computational Linguistics.
 - Zihan Wang, Dheeraj Mekala, and Jingbo Shang. 2021. X-class: Text classification with extremely weak supervision. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for*

Computational Linguistics: Human Language Technologies, pages 3043–3053, Online. Association for Computational Linguistics.

- Noah Weber, Leena Shekhar, Niranjan Balasubramanian, and Nate Chambers. 2018. Hierarchical quantized representations for script generation. In Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018, pages 3783– 3792. Association for Computational Linguistics.
- David Wilmot and Frank Keller. 2021. Memory and knowledge augmented language models for inferring salience in long-form stories. In *Proceedings of the* 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, pages 851–865. Association for Computational Linguistics.
- Susik Yoon, Dongha Lee, Yunyi Zhang, and Jiawei Han. 2023. Unsupervised story discovery from continuous news streams via scalable thematic embedding. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 802–811.
- Jingyi You, Dongyuan Li, Hidetaka Kamigaito, Kotaro Funakoshi, and Manabu Okumura. 2022. Joint learning-based heterogeneous graph attention network for timeline summarization. In *Proceedings* of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 4091–4104.
- Hongming Zhang, Muhao Chen, Haoyu Wang, Yangqiu Song, and Dan Roth. 2020. Analogous process structure induction for sub-event sequence prediction. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020, pages 1541– 1550. Association for Computational Linguistics.
- Yunyi Zhang, Fang Guo, Jiaming Shen, and Jiawei Han. 2022. Unsupervised key event detection from massive text corpora. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '22, page 2535–2544, New York, NY, USA. Association for Computing Machinery.

A Experimental Settings

For implementing **EpiMine**, we use the following hyperparameters across all datasets: $\delta = 25\%$, $sim_thresh = 0.75$. All other hyperparameters are set to their respective default values. We provide all experimental settings in Appendix A.To determine statistical significance, we check for $\geq \mu - \sigma$. For our word representations, we use bert-base-uncased. For our sentence representations, we use all-mpnet-base-v2. We choose

	Terrorism			Na	Natural Disasters			Politics		
Ablations	5-prec	5-recall	5-F1	5-prec	5-recall	5-F1	5-prec	5-recall	5-F1	
EpiMine-Top	0.2292	0.2435	0.2144	0.3817	0.2232	0.2450	0.1051^\dagger	0.2233	0.1201^{\dagger}	
TF-IDF	0.0985	0.1403	0.1059	0.3284^{\dagger}	0.1919^{\dagger}	0.2221^{\dagger}	0.0907	0.1908	0.0916	
No DC	0.1968^{\dagger}	0.1752^{\dagger}	0.1707^{\dagger}	0.2520	0.1546	0.1785	0.1126	0.2108^{\dagger}	0.1299	

Table 5: Ablation studies conducted on top 25% of article episode clusters (Section 3.4).

947to use Claude-2.14 for fluent candidate episode948estimation due to its strong structured JSON/XML949input and output formatting abilities. However, this950proprietary model can be replaced with any open-951source model as EpiMine is model-agnostic. We952use only one NVIDIA GeForce GTX 1080 for all953experiments; for non-API models, we utilize two954NVIDIA-RTX A6000s.

B Baselines

955

956

957

961

962

963

965

966

967

968

969

970

971

974

975

976

978

979

982

983

We compare against the following methods using the evaluation metrics specified in Appendix J.

- **K-means** (Likas et al., 2003): No. of groundtruth episodes is given; clusters segments based on semantic similarity of ST (Reimers and Gurevych, 2019) embeddings.
- EvMine (Zhang et al., 2022): Unsupervised framework for key event detection that leverages peak phrases and detects communities using event-indicative features. We extend the original document-level method to the segment level for episode detection.
- EMiner (Jiao et al., 2023): Unsupervised event chain mining that performs atomic action clustering. For episode detection, we map its final output, a list of events, back to the original sentences from which each event was extracted, treating these sentences as segments. To retrieve more episode-associated segments, we use ST (Reimers and Gurevych, 2019) to select the k most similar segments to each cluster sentence.

We also include the following full and partial ablations of EpiMine (clusters segments from all articles vs. top $\delta\%$ articles, respectively):

• No Confidence: A full ablation, where all input segments are classified based on the episode with max cosine similarity instead of using the confidence score from Equation 2.

• No LLM: We take the estimated episode clusters from Section 3.4 that normally would have been inputted into the LLM, and instead use them to compute our episode representations directly. These representations are used for our classification step (Section 3.5), run on the full dataset with confidence.

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1001

1004

1005

1006

1007

1008

1010

1011

1012

- **EpiMine-Top**: A partial ablation which directly outputs the intermediate episode clusters formed based on the top articles identified in Section 3.4 without inputting them into the LLM-based episode estimation step.
- **TF-IDF**: A partial ablation which replaces the salience and synonym expansion step (Section 3.2) with TF-IDF.
- **No DC**: A partial ablation which replaces the discriminative co-occurrence score (Equation 1) with raw pair frequency.

C Ablation Studies

In Table 5, we note a significant decrease in the quality of our top articles (partitioned into episodes) without our salience and discriminative co-occurrence measures. The politics dataset does show slight improvements in precision and F1 when discriminative co-occurrence is replaced, due to more term overlap across episodes, resulting in less distinct, sequential episodes

D Key Event Corpus Dataset Construction

Given that our task is novel and no large-scale key 1013 event-specific news corpus is available for this task 1014 where the key events are guaranteed to contain 1015 distinguishable episodes, we briefly discuss how 1016 we collect the input corpus from online news data. 1017 Given our set of key events (as listed in Section 1018 4.1), we first scrape the external reference list from 1019 their corresponding Wikipedia page and select the 1020 news articles that have been published within two 1021

⁴claude.ai/

months of given key event's start date (e.g., all arti-1022 cles selected for "January 6 2021 Capitol Attack" 1023 would have been published between November 6-1024 March 6). This is important as we want to prior-1025 itize the news articles which focus on describing the episodes of the key event and their correspond-1027 ing aspects as opposed to primarily opinions or 1028 analyses. This allows us to motivate our task as 1029 one critical for currently evolving key events which 1030 required a more fine-grained episodic timeline. Fur-1031 thermore, it is consequently unlikely for a single article to cover all of the episodes and exclusively 1033 episodes under a key event. Despite this being 1034 more challenging, it is acceptable as the goal of 1035 our task is to extract only the key event-related 1036 episodes, which must be substantiated by multiple documents in either case. 1038

1039

1040

1041

1042

1044

1045

1046

1048

1050

1051

1052

1053

1054

1055

1057

1058

1059

1062

1063

1064

1065

1066

1067

1068

1069

1071

During the collection process, we targeted selecting a diverse set of key events topics within a theme. For instance, we attempted to cover every type of "natural disaster", including tornados, wildfires, and etc. When selecting key events, we leave out those with less than 20 hyperlinks in the Wikipage. Table 1 summarizes the statistics for these datasets. We also construct a background news corpus of approximately 4,000 long news articles using the New York Times corpus for topic categorization (Meng et al., 2020).

We list all of our selective themes and their corresponding key events included in the dataset that we construct:

• Terrorism and attacks: 2021 Atlanta spa shootings; 2014 Montgomery County Shootings; 2021 Indianapolis FedEx shooting; 2022 Cincinnati FBI field office attack; 2019 Jersey City shooting; 2019 Naval Air Station Pensacola shooting; 2022 Greenwood Park Mall shooting, 2018 Capital Gazette shooting; 2021 Collierville Kroger shooting; 2019 Kyoto Animation arson attack

• Natural disasters: 2023 Tornado outbreak sequence; 2023 Hawaii Wildfires; 2021 Western Kentucky tornado; 2017 Mocoa landslide; 2010 Haiti earthquake; 2021 Henan floods; 2019 Nyonoksa radiation accident; 2022 North American winter storm; 2011 Fukushima nuclear accident

• Political Events: 2020 Kyrgyz Revolution, 2019 Storming of the Hong Kong Legislative Council Complex; 2019 Siege of the Hong Kong Polytechnic University; 2017 Zimbabwean coup; 2018 Italian government formation; 2021 Jan-
uary 6 United States Capitol attack; 2018 Thai1073Cave Rescue Operation; 2018 Armenian Revolu-
tion; 2017 Lebanon–Saudi Arabia dispute; 20131074Tunisian political crisis1076

1077

1078

1079

1080

1082

1083

1084

1085

1086

1087

1088

1090

1092

1104

E Identifying Salient Terms for Episode Detection

Definition 4 (Salience). A term is <u>salient</u> if it is (1) distinct and significant to understanding a given key event, as well as (2) frequently found in a key event's segments and infrequently in other background/general articles.

We define the salience score of a term w_i within segment s as the following function, where $freq(w_i)$ is the number of key event segments that w_i is contained in, N_{bg} is the number of news articles in the background corpus we construct (using general New York Times articles), and $bgf(w_i)$ is the number of background articles that w_i is present in.

Salience
$$(w_i) = (1 + \log^2 (freq(w_i)))$$

 $\times \log \left(\frac{N_{bg}}{bgf(w_i)}\right)$ (3)

Stop words and infrequent terms $(freq(w_i) < 5)$ are assigned a salience score of -1. A key event's 1094 set of salient terms T is comprised of the terms with 1095 a salience score above the mean salience across 1096 the entire vocabulary. In the case of infrequent 1097 synonyms used by a journalist as a stylistic choice 1098 (e.g., "demonstrations", "marches"), we expand 1099 T with terms that are similar (cosine-similarity) 1100 to their static word representations (average of its 1101 contextualized word embeddings across entire key 1102 event corpus). 1103

F Key Event News Article Pre-Processing

Given that the expected output for the episode de-1105 tection task is a *cluster of text segments*, we first 1106 must segment each key event news article. We 1107 would like to ideally preserve both the primary 1108 aspects (e.g., core entities and their actions) and 1109 peripheral aspects (e.g., reactions to a core entity's 1110 action) relevant to that episode, which may be help-1111 ful for cross-document episode co-reference res-1112 olution. In order to do this, we utilize the text 1113 segmentation method, C99 (Choi, 2000). Further-1114 more, in order to assist with the cohesiveness of the 1115 segment, we employ entity co-reference resolution 1116

1168

1198

1199

1200

1201

1202

1203

1204

1205

1117before performing segmentation, which assists with1118retaining the context across text segments ("They1119surrounded the legislative building [...]" \rightarrow "The1120protesters surrounded the legislative building [...]").1121Our core methodology is given these text segments1122(in their raw form, without co-references resolved)1123and their source articles as the primary inputs.

G Additional Details for Episode Partitioning

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

We note that for determining semantic similarity between the terms of two segments, we use both (1) the average cosine similarity between all unordered pairs of terms between segment (i-1) and (i), and (2) the cosine similarity between the average of static term representations in (i - 1) and the average of static term representations in (i). Furthermore, we filter out any non-salient segments before episode partitioning to avoid any influence of noisy segments (e.g., journalist's analysis, summary statements, historical comparisons, and other generic noise) on the quality of our episode fragments.

Finally, following (Wang et al., 2021; Kargupta et al., 2023), we take the harmonic mean of all pairwise discriminative co-occurrence scores instead of a simple average. This allows us to prioritize the more salient *and* discriminative terms when determining the episode partitions. For instance, if "protesters" consistently occurs throughout the majority of episodes and thus has a low average discriminative co-occurrence, then it is not as informative for episode partitioning.

H Episode & Segment Representations

We compute a candidate episode representation by encoding both its LLM core attributes and extracted segments using SentenceTransformers (ST) (Reimers and Gurevych, 2019). Following extremely weakly supervised text classification works, such as (Wang et al., 2021; Kargupta et al., 2023), we take the harmonic mean of these representations, as the latter extracted segments are likely not as significant as the earlier extractions and core attributes. We similarly encode all input *segments* with the same ST model.

I Claude-2 Prompt & Example for Candidate Episode Estimation

Prompt. We use the following prompt for estimating fluent candidate episodes from our input episode fragment clusters. We denote k as the number of episode fragment clusters outputted after clustering the top article episode fragments in Section 3.4.

Task: You are a key news event analyzer 1169 that is aiming to detect episodes (a 1170 representative subevent that reflects a 1171 critical sequence of actions performed by 1172 a subject at a certain and/or location) 1173 based on text segments from different 1174 news articles. Given the above groups of 1175 article segments, predict at least 2 and 1176 at most {k} potential episodes of the 1177 Some groups may fall under key event. 1178 the same episode. Output your answer 1179 inside the tags <answer></answer> as a 1180 JSON object where each item is also a 1181 JSON with the key "title" with the value 1182 containing the [subject, action, object, 1183 time, location] of the episode, a key 1184 "keywords" with the string value being a 1185 list of 5-10 associated keywords unique 1186 to that specific episode, and a final key 1187 "example sentences" with a value being 1188 a list of 2-5 extracted sentences from 1189 the input segment groups. Feel free to 1190 output less than {k} episodes if you feel 1191 that any are redundant (could fit under 1192 an existing candidate episode). The 1193 title, keywords, and example sentences 1194 of a predicted episode should not be able 1195 to be placed under another different 1196 predicted episode. 1197

Example. Below, we provide an example of EpiMine's candidate episode estimation step (Section 3.4). Specifically, the LLM identifies the core attributes (subject, action, object, time, location) of each unique cluster, relevant keywords, and top extracted text segments given the input clusters:

'title': ['Protesters', 'storm and 1206 vandalize', 'Hong Kong's Legislative 'July Council building', 1, 2019', 1208 'Legislative Council building in 1209 Admiralty, Hong Kong'] 1210 'keywords': ['vandalism, graffiti, 1211 violence, escalation, ransacking'] 1212 'example_sentences': ['Hundreds of 1213 anti-extradition bill protesters finally 1214 broke into the legislature after many 1215 hours of attacking the public entrance 1216 1217and ransacked the building, including1218displaying the colonial Hong Kong flag1219in the chamber.', 'Slogans on the wall1220read: "Murderous regime", and "There are1221no rioters only a tyrannical regime."]

J Evaluation Metrics

Following a recent work on key event detection (Zhang et al., 2022), we adapt the k-prec, k-recall, and k-F1 to quantitatively evaluate the episode detection results. We use these metrics to evaluate how the model output matches the ground truth episodes using the top-k segments within each detected episode. Formally, suppose there are N ground truth episodes $\mathcal{G} = \{G_1, G_2, \ldots, G_N\}$, each of which is a set of text segments related to its corresponding episode. $\mathcal{E} = \{E_1, E_2, \ldots, E_K\}$ are the model predicted episodes, each of which is a ranked list of segments, and $E_{j,k}$ means the top-k segments within E_j . Then, the k-metrics are defined as follows:

$$\begin{aligned} \mathbf{k}\text{-}\mathrm{prec} &= \frac{\sum_{G_i \in \mathcal{G}} \mathbbm{1}(\exists E_j \in \mathcal{E}, E_{j,k} \cap G_i \geq \frac{k}{2})}{\sum_{E_j \in \mathcal{E}} \mathbbm{1}(|E_j| \geq k)} \\ \mathbf{k}\text{-}\mathrm{recall} &= \frac{\sum_{G_i \in \mathcal{G}} \mathbbm{1}(\exists E_j \in \mathcal{E}, E_{j,k} \cap G_i \geq \frac{k}{2})}{N} \\ \mathbf{k}\text{-}\mathrm{F1} &= \frac{2 \cdot \mathbf{k}\text{-}\mathrm{prec} \cdot \mathbf{k}\text{-}\mathrm{recall}}{\mathbf{k}\text{-}\mathrm{prec} + \mathbf{k}\text{-}\mathrm{recall}} \end{aligned}$$

1240

1241 1242

1244

1245 1246

1222

1223 1224

1225

1226

1228

1229

1230

1231 1232

1233

1234

1235 1236

1237

1238

1239

K Claude-2 Prompt for Dataset Annotation

We automatically annotate our dataset using Claude-2.1 using the prompt below (before an additional human-verification stage):

You are а news event analyzer that 1247 labels text segments of a news article 1248 with their matching event episode 1249 I will give you several description. 1250 text segments, and several episodes of 1251 a key event in tuples. We define an 1252 1253 episode as the following: an episode is a set of thematically coherent text 1254 segments discussing a particular set of 1255 core entities performing actions for or 1256 1257 towards an object(s) at a certain time and/or location during a real-world key 1258 The entities, actions, objects, event. 1259 time, and location can all be considered aspects of an episode. 1261

	1262
[one-shot demonstration & format	1263
specification]	1264
Please help classify the text segments	1265
under different episodes (the output	1266
value for each segment should be an	1267
integer key of each episode). If you	1268
think a text segment cannot be used	1269
to describe any episodes, please use	1270
"X" in the output to indicate the lack	1271
of an episode tuple number for that	1272
segment. If a text segment is very	1273
general, does not describe the key event	1274
at hand, or can be matched to multiple	1275
episodes, then please use a "M" in the	1276
output to indicate the multiple episode	1277
mapping for that segment. There should	1278
be a value assigned to each of the	1279
<pre>len(segments) segments (segment_0,,</pre>	1280
<pre>segment_len(segments)-1).</pre>	1281