

# An LLM-Embedding Semantic Adaptation Network for Post-level Semantic Drift Evaluation

Ning Chen<sup>1</sup> Mingyu Kang<sup>1\*</sup> Jie Li<sup>2</sup> Linyuan Lü<sup>1\*</sup>

<sup>1</sup>University of Science and Technology of China, Hefei, China

<sup>2</sup>Shanghai Artificial Intelligence Laboratory, Shanghai, China

mengtahua@163.com kangmingyu@ustc.edu.cn

lijie.32@outlook.com linyuan.lv@ustc.edu.cn

## Abstract

Evaluating semantic drift is essential for understanding dynamical discourse evolution and opinion formation in online discussions. However, sparse and uneven distributions of event-specific keywords prevent traditional models from capturing post-level semantic drift. Thus, to address this issue, we propose an LLM-embedding Semantic Adaptation Network (LLM-SAN), which is a hybrid semantic drift evaluation model with an LLM-Embedding gated recurrent unit (GRU) module, an LLM-Embedding graph convolutional network (GCN) module and a multi-expert adaptive fusion module. The GRU module is used to extract features from event related posts, and The GCN is used to extract features from temporal graphical topic posts. Then, the features are merged by the multi-expert adaptive fusion module. Finally, this module predicts the future post embedding, and the prediction error is used to evaluate and detect the semantic drift points. Extensive experiments are conducted, and the results show that LLM-SAN achieves the state-of-the-art performance on the semantic drift evaluation task, compared to the other baselines. Ablation experiments are also conducted to show the effectiveness of each module in LLM-SAN. Code is available at: <https://github.com/mengtahua/LLMSAN>.

## 1 Introduction

Semantic drift evaluation is a critical yet challenging task to quantify and measure dynamical semantic changes, as illustrated in Fig. 1. In this task, event-related posts causally evolve and drift with time-varying distribution within public attention and discussion topics (Wu et al., 2024; Wang and Goutte, 2018; Balepur et al., 2023; Kang et al., 2024; He et al., 2024; et al., 2025; Sekar et al., 2026; Frank and Afli, 2025). The vocabulary distributions, sentiment polarity, topic structures and

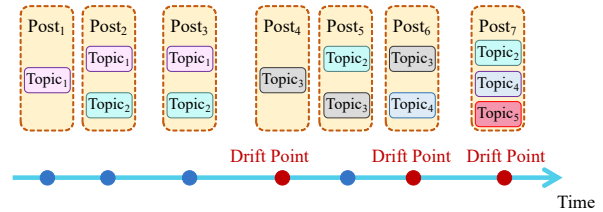


Figure 1: Semantic Drift Evaluation.

contextual semantic associations gradually change during the evolution (Garroppo et al., 2018; Karjus et al., 2020; Gaul and Vincent, 2017; Unankard and Nadee, 2020; Cao et al., 2023; Periti and Tahmasebi, 2024).

There have been studies on event-related content evolution, which can be classified into three types: (i) probabilistic models (Blei et al., 2003; Blei and Lafferty, 2006; Wang and McCallum, 2006) learning topic changes from word-frequency distributions; (ii) incremental models (Wang and Goutte, 2018; Rieger et al., 2022; Patil et al., 2025b) learning dynamic semantic evolution via online detection mechanisms; and (iii) deep-learning models (Wang et al., 2019; Dieng et al., 2019; Rahimi et al., 2024) learning semantic representations through neural encoders.

However, these models all perform at the time level or event level. Thus, they hardly model the rapidly-evolving post content, where the semantic drift is reflected through sparse and unevenly distributed event-specific keywords. This makes them difficult to leverage fine-grained semantic information and perform with low accuracy if event scales vary significantly or salient semantic content is unevenly distributed.

Thus, to address this issue, we propose an LLM-Embedding Semantic Adaptation Network (LLM-SAN) for post-level semantic drift evaluation. In this model, temporal dynamics features are obtained from a gated recurrent unit (GRU) with inputs from an LLM-based encoder, capturing accu-

\* Corresponding authors.

mulated semantic drift across posts. Meanwhile, topical features are obtained from a graph convolutional network (GCN) with inputs from an LLM-based generator, enhancing topic representations and sparse keyword signals. An attention-guided multi-expert mechanism then adaptively fuses these features, which are subsequently used for event-stage segmentation.

Thus, The main contributions of this work are as follows:

1. LLM-SAN model is proposed to evaluate the post-level semantic drift, providing a systematic way to capture fine-grained semantic changes in individual posts.
2. The model uses a GRU with inputs from an LLM-based encoder to capture accumulated post-level semantic drift, a GCN with inputs from an LLM-based generator to enhance topic representations and discontinuous keyword signals, and an attention-guided multi-expert fusion mechanism to adaptively combine multi-scale features, effectively addressing sparse or unevenly distributed event-specific keywords.
3. Experiments demonstrate that LLM-SAN achieves state-of-the-art performance on post-level semantic drift detection. Ablation studies verify the effectiveness of the key modules, and tests on real-world datasets confirm its validity, ultimately facilitating accurate event stage segmentation.

## 2 Related Work

This section introduces related work on semantic drift and content evolution, including probabilistic models, incremental models, and deep-learning models.

**Probabilistic models** describe the documents by using latent topics inferred from word-frequency distributions. Among them, latent dirichlet allocation (LDA) (Blei et al., 2003; Chauhan and Shah, 2021; Goyal and Kashyap, 2022) models document-topic and topic-word distributions under Dirichlet priors. Then, dynamic topic models (Blei and Lafferty, 2006) extend LDA by linking topic parameters across adjacent time slices through state-space assumptions. Further, topics over time (Wang and McCallum, 2006) incorporate temporal priors to capture time-dependent topic distribution changes.

**Incremental models** are designed to track the evolution of semantic in streaming or continuously arriving data. Among them, online LDA (Wang and Goutte, 2018; Fan et al., 2021; Zhou et al., 2023; Balepur et al., 2023) models are proposed to update topic parameters incrementally as new documents arrive sequentially. By contrast, Rolling LDA (Rieger et al., 2022) uses a sliding window to capture the features of continue rolling topics. Moreover, joint dynamic topic model (Zhang and Lauw, 2022) is also proposed to analyze varying topics by the combination of time-aware optimal transport and temporal point process techniques. Graph-based online models (Patil et al., 2025b,a) are proposed to represent topic drift through evolving graph structures.

**Deep-learning models** leverage neural encoders to learn semantic representations. Among them, neural topic models (Wang et al., 2019; Wu et al., 2024; Boutaleb et al., 2024) are proposed to perform topic inference by using neural variational framework. Based on this framework, dynamic embedded topic model (DETM) (Dieng et al., 2019) use word embedding technique to capture the features of topic evolution over time. Then, dynamic structured neural topic model (Miyamoto et al., 2023) is proposed to capture the topic branching and merging features in evolution process by self-attention mechanisms. Further, aligned neural topic model (ANTM) (Rahimi et al., 2024) uses large language model (LLM) to extract time-aware features from evolving topics.

## 3 Methodology

The overall architecture of the proposed LLM-SAN model is presented as shown in Fig. 2.

### 3.1 Temporal Evolution Trend Extraction

Formally, let an event consist of a sequence of temporally ordered posts  $\mathcal{P}_{1:T} = \{p_1, p_2, \dots, p_T\}$ , where each post  $p_t$  is associated with textual content and time step  $t$ .

An LLM-based encoder is applied to generate single post content semantic embeddings:

$$\mathbf{e}_t^{\text{temporal}} = \text{LLM}_{\text{emb}}(p_t), \quad (1)$$

$\mathbf{e}_t^{\text{temporal}} \in \mathbb{R}^{d_1}$  is the semantic embedding of  $p_t$ . The resulting embedding  $\mathbf{e}_t^{\text{temporal}}$  is subsequently input into a GRU module to capture temporal evolution trends and accumulate post-level semantic drift.

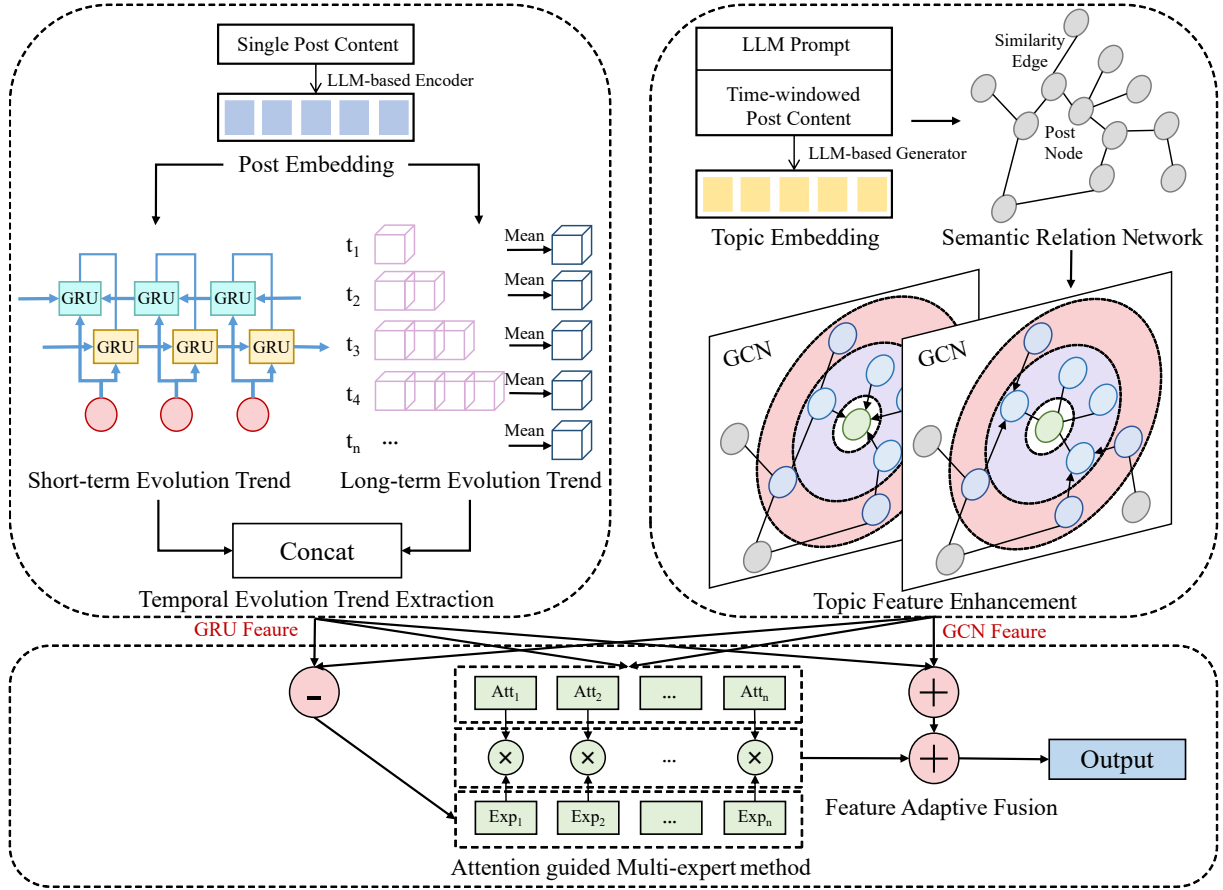


Figure 2: LLM-Embedding Semantic Adaptation Network.

At each time step  $t$ , a two-layer GRU processes the sequence of post embeddings,  $\mathbf{e}_{t-c}^{\text{temporal}}, \dots, \mathbf{e}_{t-1}^{\text{temporal}}$ , from the preceding  $c = 1, 2, \dots, T-1$  time steps, updating its hidden state recurrently to capture short-term evolution trends in event-related posts.

Specifically, given the post embedding at time step  $t-1$ , the GRU hidden state is updated as

$$\mathbf{z}_t = \sigma(\mathbf{W}_z \mathbf{e}_{t-1}^{\text{temporal}} + \mathbf{U}_z \mathbf{h}_{t-1}), \quad (2)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_r \mathbf{e}_{t-1}^{\text{temporal}} + \mathbf{U}_r \mathbf{h}_{t-1}), \quad (3)$$

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h \mathbf{e}_{t-1}^{\text{temporal}} + \mathbf{U}_h (\mathbf{r}_t \odot \mathbf{h}_{t-1})), \quad (4)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t, \quad (5)$$

where  $\mathbf{h}_{t-1} \in \mathbb{R}^{d_h}$  is the GRU hidden state from the previous time step, with hidden dimension  $d_h$ , and the initial hidden state  $\mathbf{h}_0$  is set to a zero vector.  $\mathbf{W}_z \in \mathbb{R}^{d_h \times d_1}$ ,  $\mathbf{W}_r \in \mathbb{R}^{d_h \times d_1}$ , and  $\mathbf{W}_h \in \mathbb{R}^{d_h \times d_1}$  are learnable weight matrices applied to the input embedding  $\mathbf{e}_{t-1}^{\text{temporal}}$ , where  $d_h$  denotes the hidden dimension and  $d_1$  denotes the dimensionality of  $\mathbf{e}_t^{\text{temporal}}$ ;  $\mathbf{U}_z \in \mathbb{R}^{d_h \times d_h}$ ,  $\mathbf{U}_r \in \mathbb{R}^{d_h \times d_h}$ ,

$\mathbf{U}_h \in \mathbb{R}^{d_h \times d_h}$  are learnable weight matrices applied to the previous hidden state, with dimension  $d_h$ ;  $\mathbf{z}_t \in \mathbb{R}^{d_h}$  and  $\mathbf{r}_t \in \mathbb{R}^{d_h}$  are the update and reset gates, controlling information flow;  $\tilde{\mathbf{h}}_t \in \mathbb{R}^{d_h}$  is the candidate hidden state, and  $\mathbf{h}_t \in \mathbb{R}^{d_h}$  is the updated hidden state;  $\sigma(x) = \frac{1}{1+e^{-x}}$  denotes the sigmoid function,  $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$  denotes the hyperbolic tangent function, and  $\odot$  represents element-wise multiplication; Through recursive updates over the previous  $c$  time steps (i.e., using embeddings from time steps  $t-c$  to  $t-1$ ),  $\mathbf{h}_t$  captures short-term evolution trends in event-related posts.

To balance short-term variations with long-term evolution trends, a historical mean embedding is computed over all past posts up to time  $t-1$  as

$$\bar{\mathbf{e}}_t = \frac{1}{t-1} \sum_{i=1}^{t-1} \mathbf{e}_i^{\text{temporal}}, \quad (6)$$

and mapped through a nonlinear transformation:

$$\mathbf{m}_t = \text{ReLU}(\mathbf{W}_m \bar{\mathbf{e}}_t + \mathbf{b}_m), \quad (7)$$

where  $\mathbf{W}_m \in \mathbb{R}^{d_h \times d_1}$  and  $\mathbf{b}_m \in \mathbb{R}^{d_h}$  are learnable weight matrices applied to the historical mean

embedding, and  $\mathbf{m}_t \in \mathbb{R}^{d_h}$  captures the long-term semantic evolution trend accumulated before time step  $t$ .

Then, the short-term representation  $\mathbf{h}_t$  and the long-term trend representation  $\mathbf{m}_t$  are concatenated and projected via a linear layer to obtain the GRU output representation as follows:

$$\mathbf{f}_t^{\text{GRU}} = \mathbf{W}_g[\mathbf{h}_t; \mathbf{m}_t] + \mathbf{b}_g, \quad (8)$$

where  $[\cdot; \cdot]$  denotes vector concatenation, and  $\mathbf{W}_g \in \mathbb{R}^{d_1 \times 2d_h}$  and  $\mathbf{b}_g \in \mathbb{R}^{d_1}$  are learnable weight matrices applied to the concatenated embeddings.  $\mathbf{f}_t^{\text{GRU}} \in \mathbb{R}^{d_1}$  is a representation that jointly preserves long-term stable semantic information while emphasizing short-term semantic shifts, thereby providing a comprehensive characterization of post-level semantic drift.

### 3.2 Topic Feature Enhancement

To encode contextual topic semantics from a preceding time window, a LLM-based generator is applied to the time-windowed post contents  $\mathcal{P}_{t-c:t-1} = \{p_{t-c}, p_{t-c+1}, \dots, p_{t-1}\}$ , where  $c = 1, 2, \dots, T-1$  denotes the length of the temporal window. The LLM-based generator produces a topic-aware embedding:

$$\mathbf{e}_t^{\text{topic}} = \text{LLM}_{\text{gen}}(\mathcal{P}_{t-c:t-1}). \quad (9)$$

where  $\mathbf{e}_t^{\text{topic}} \in \mathbb{R}^{d_2}$ , here  $d_2$  is the dimension of the mean-pooled output of the generator LLM over the window.

The generation process is guided by the following prompt:

**“You are a Twitter user. Based on the following posts published in the previous time period, generate a plausible post for the next time step.”**

A single vector summarizing the window is obtained via mean pooling over the hidden states. The resulting embedding  $\mathbf{e}_t^{\text{topic}}$ ,  $t = c+1, c+2, \dots, T$  is then input into a GCN module to enhance contextual topic representations.

After obtaining the topic-aware embeddings  $\mathbf{e}_t^{\text{topic}}$  from the preceding time window, these embeddings are treated as nodes in a temporal topic graph, and a semantic relation network is constructed to enhance discontinuous keyword signals and topic-level representations using a two-layer GCN.

First, a pairwise similarity matrix  $\mathbf{S}$  is computed among the topic embeddings:

$$\mathbf{S} = \mathbf{E}^{\text{topic}} (\mathbf{E}^{\text{topic}})^{\top}, \quad (10)$$

where  $\mathbf{E}^{\text{topic}} = [\mathbf{e}_{c+1}^{\text{topic}}; \mathbf{e}_{c+2}^{\text{topic}}; \dots; \mathbf{e}_T^{\text{topic}}] \in \mathbb{R}^{(T-c) \times d_2}$  denotes the matrix of topic embeddings corresponding to historical time steps within the preceding window. Edges are formed between node pairs corresponding to the top 5% largest values in the similarity matrix  $\mathbf{S}$ , resulting in an undirected semantic graph  $\mathcal{G}$ . Each node represents a time-specific topic embedding  $\mathbf{e}_t^{\text{topic}}$ , and edges denote semantic similarity between topics across time steps.

Subsequently, a two-layer GCN is applied to aggregate information from semantically related nodes, yielding enhanced topic representations. The update rule at layer  $l+1$  is defined as:

$$\mathbf{H}^{(l+1)} = \sigma \left( \mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{A}} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(l)} \mathbf{W}_p^{(l)} \right), \quad (11)$$

where  $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$  denotes the adjacency matrix  $\mathbf{A} \in \mathbb{R}^{(T-c) \times (T-c)}$  with self-loops  $\mathbf{I} \in \mathbb{R}^{(T-c) \times (T-c)}$ ,  $\mathbf{D} \in \mathbb{R}^{(T-c) \times (T-c)}$  is the degree matrix,  $\mathbf{H}^{(1)} \in \mathbb{R}^{(T-c) \times d_2}$  represents the initial node embeddings,  $\mathbf{H}^{(l)} \in \mathbb{R}^{(T-c) \times d_h}$  for  $l \geq 1$  represents the node embeddings at layer  $l$ ,  $\mathbf{W}_p^{(1)} \in \mathbb{R}^{d_2 \times d_h}$ ,  $\mathbf{W}_p^{(l)} \in \mathbb{R}^{d_h \times d_h}$  for  $l \geq 1$  are learnable weight matrices.

After two graph convolutional layers, the resulting node representations  $\mathbf{H}^{(2)} \in \mathbb{R}^{(T-c) \times d_h}$  are then linearly projected using learnable parameters  $\mathbf{W}_q \in \mathbb{R}^{d_1 \times d_h}$  and  $\mathbf{b}_q \in \mathbb{R}^{d_1}$  to obtain the enhanced topic feature  $\mathbf{f}_t^{\text{GCN}} \in \mathbb{R}^{d_1}$  at time step  $t$ :

$$\mathbf{f}_t^{\text{GCN}} = \mathbf{W}_q \mathbf{H}^{(2)} + \mathbf{b}_q. \quad (12)$$

### 3.3 Feature Adaptive Fusion

After obtaining the GRU feature  $\mathbf{f}_t^{\text{GRU}} \in \mathbb{R}^{d_1}$  and GCN feature  $\mathbf{f}_t^{\text{GCN}} \in \mathbb{R}^{d_1}$ , the two features are first merged through a fixed equal-weight linear fusion as the main predictive signal:

$$\mathbf{f}_t^{\text{base}} = 0.5\mathbf{f}_t^{\text{GRU}} + 0.5\mathbf{f}_t^{\text{GCN}}. \quad (13)$$

To account for conflicts between temporal dynamics and structural semantics, an attention-guided multi-expert method is applied to the feature difference:

$$\mathbf{d}_t = \mathbf{f}_t^{\text{GRU}} - \mathbf{f}_t^{\text{GCN}}. \quad (14)$$

Each of the  $I \in \mathbb{Z}^+$  experts computes a residual projection of this difference:

$$\mathbf{r}_t^{(i)} = \mathbf{W}_e^{(i)} \mathbf{d}_t, \quad i = 1, \dots, I, \quad (15)$$

where  $\mathbf{W}_e^{(i)} \in \mathbb{R}^{d_1 \times d_1}$  is the learnable weight matrix for the  $i$ -th expert.

Attention weights are then computed over the experts to adaptively obtain a weighted residual:

$$\alpha_t = \text{softmax}(\mathbf{W}_b \tanh(\mathbf{W}_a [\mathbf{f}_t^{\text{GRU}}; \mathbf{f}_t^{\text{GCN}}])), \quad (16)$$

$$\mathbf{r}_t^{\text{weighted}} = \sum_{i=1}^I \alpha_t^{(i)} \mathbf{r}_t^{(i)}, \quad (17)$$

where  $\mathbf{W}_a \in \mathbb{R}^{d_1 \times 2d_1}$ ,  $\mathbf{W}_b \in \mathbb{R}^{I \times d_1}$  are learnable parameters and  $\alpha_t \in \mathbb{R}^I$  and  $\mathbf{r}_t^{\text{weighted}} \in \mathbb{R}^{d_1}$ .

Finally, the predicted embedding for the next time step  $t$  is obtained by adding the scaled weighted residual to the base fusion:

$$\mathbf{e}_t^{\text{temporal}} = \mathbf{f}_t^{\text{base}} + \lambda \mathbf{r}_t^{\text{weighted}}, \quad (18)$$

where  $\lambda$  is a small scaling factor controlling the influence of the residual and  $\mathbf{e}_t^{\text{temporal}} \in \mathbb{R}^{d_1}$ .

### 3.4 Implementation of LLM-SAN

In the training process, the loss function is defined as the reconstruction error between the predicted post embeddings and the corresponding ground-truth embeddings:

$$\text{loss} = \frac{1}{T-c} \sum_{t=c+1}^T \|\hat{\mathbf{e}}_t^{\text{temporal}} - \mathbf{e}_t^{\text{temporal}}\|^2, \quad (19)$$

where  $\hat{\mathbf{e}}_t^{\text{temporal}}$  is the estimation of  $\mathbf{e}_t^{\text{temporal}}$ , and  $T-c$  denotes the number of future time steps to be predicted. Note that, we here use samples within time steps  $t = c+1, \dots, T$  to predict that within time steps  $t = 1, \dots, c$ .

The semantic drift score for a post at time step  $t$  is then defined as its prediction error:

$$\epsilon_t = \|\hat{\mathbf{e}}_t^{\text{temporal}} - \mathbf{e}_t^{\text{temporal}}\|^2, \quad (20)$$

where  $\epsilon_t \in \mathbb{R}$  quantifies the deviation between the predicted embedding and the actual embedding.

The overall complexity of LLM-SAN is  $O(T(d_1 d_h + d_h^2) + E d_2 + T d_2^2 + K d_o^2)$  in time,  $O(d_1 d_h + d_h^2 + d_2^2 + K d_o^2)$  in parameter storage, and  $O(T d_h + T d_2 + E + K d_o)$  in activation storage, where  $T$ ,  $E$ ,  $d_1$ ,  $d_h$ ,  $d_2$ ,  $d_o$ , and  $K$  denote the number of time steps, edges, input dimension, GRU hidden dimension, GCN dimension, expert output dimension, and number of experts, respectively.

## 4 Experiment

In this section, comprehensive experiments are conducted to verify the proposed model achieves state-of-the-art semantic drift evaluation performance and enables reasonable event stage segmentation.

### 4.1 Datasets

The experiments are conducted on two benchmark datasets. They are as follows:

**DTELS** (Zhang et al., 2025) is a Chinese news corpus consisting of temporally ordered articles and reference timelines annotated at multiple granularities, the dataset is available at: <https://github.com/chenlong-clock/DTELS-Bench>. Three levels of temporal granularity are provided: coarse-grained (key milestones), medium-grained (important developments), and fine-grained (detailed evolutionary points). Events containing all three granularities and more than 350 documents are retained. Documents are temporally aligned with reference timelines using TF-IDF similarity and summary coverage, and are labeled according to the earliest matched timeline node (3/2/1/0 from coarse to unmatched). The label magnitude indicates the degree of semantic drift. After preprocessing, 18 event datasets are obtained for semantic drift evaluation. See also Appendix A for details.

**Twibot-22** (Feng et al., 2022) is a large-scale Twitter benchmark containing approximately one million users and over 76 million tweets, the dataset is available at: <https://github.com/LuoUndergradXJTU/TwiBot-22>. Event-level datasets are constructed by filtering tweets within predefined time periods using event-specific keywords. Noisy samples are removed, including very short posts, replies, and texts dominated by irrelevant content. Following this procedure, 4 event datasets are extracted for event stage segmentation. See also Appendix B for details.

### 4.2 Parameter Settings and Metrics

Parameter	Value
Number of experts	4
Time window length	12
LLM-based encoder	Qwen3-Embedding-0.6B
Hidden_dim	128
Epochs	3000
$\lambda$	0.1
Learning rate	0.001

Table 1: Hyperparameter settings for LLM-SAN.

Event ID	Rolling-LDA	DETM	ANTM	LLM-SAN (Qwen)	LLM-SAN (Gemma)
1	0.74 ± 0.077	1.13 ± 0.012	1.52 ± 0.565	<b>0.52 ± 0.038</b>	<b>0.52 ± 0.039</b>
2	0.72 ± 0.067	1.06 ± 0.057	1.15 ± 0.465	<b>0.48 ± 0.019</b>	<b>0.47 ± 0.037</b>
3	0.64 ± 0.078	1.27 ± 0.061	1.11 ± 0.512	<b>0.50 ± 0.029</b>	<b>0.48 ± 0.014</b>
4	0.53 ± 0.038	1.20 ± 0.050	1.30 ± 0.543	<b>0.29 ± 0.058</b>	<b>0.29 ± 0.038</b>
5	0.64 ± 0.003	1.17 ± 0.040	1.35 ± 0.241	<b>0.48 ± 0.038</b>	<b>0.43 ± 0.031</b>
6	0.72 ± 0.013	1.27 ± 0.006	1.41 ± 0.168	<b>0.31 ± 0.015</b>	<b>0.30 ± 0.013</b>
7	0.60 ± 0.023	1.14 ± 0.051	1.21 ± 0.279	<b>0.44 ± 0.042</b>	<b>0.43 ± 0.024</b>
8	0.57 ± 0.008	1.30 ± 0.091	1.19 ± 0.118	<b>0.42 ± 0.050</b>	<b>0.50 ± 0.031</b>
9	0.79 ± 0.061	1.34 ± 0.011	1.16 ± 0.530	<b>0.36 ± 0.025</b>	<b>0.36 ± 0.023</b>
10	0.72 ± 0.065	1.20 ± 0.046	1.28 ± 0.260	<b>0.44 ± 0.036</b>	<b>0.41 ± 0.022</b>
11	0.58 ± 0.035	1.24 ± 0.119	1.30 ± 0.348	<b>0.43 ± 0.040</b>	<b>0.41 ± 0.038</b>
12	0.83 ± 0.028	1.34 ± 0.028	1.36 ± 0.468	<b>0.43 ± 0.052</b>	<b>0.47 ± 0.045</b>
13	0.70 ± 0.037	1.22 ± 0.042	1.35 ± 0.228	<b>0.44 ± 0.046</b>	<b>0.48 ± 0.037</b>
14	0.68 ± 0.058	1.22 ± 0.041	1.09 ± 0.301	<b>0.43 ± 0.026</b>	<b>0.38 ± 0.014</b>
15	0.74 ± 0.005	1.27 ± 0.055	1.10 ± 0.352	<b>0.41 ± 0.075</b>	<b>0.41 ± 0.054</b>
16	0.66 ± 0.056	1.17 ± 0.093	1.20 ± 0.146	<b>0.51 ± 0.031</b>	<b>0.52 ± 0.022</b>
17	0.52 ± 0.007	1.14 ± 0.070	1.08 ± 0.373	<b>0.39 ± 0.019</b>	<b>0.43 ± 0.048</b>
18	0.69 ± 0.031	1.19 ± 0.020	1.10 ± 0.435	<b>0.48 ± 0.075</b>	<b>0.45 ± 0.032</b>

Table 2: Semantic Drift Evaluation Performance on DTELS.

The hyperparameter of LLM-SAN are summarized in Table 1. Specifically, the maximum number of stages is set to five. Two types of LLM-based generators were employed: Qwen2.5-3B-Instruct and a Gemma-2-2B-Instruct model fine-tuned for Chinese and German. Each experiment was repeated 5 times using random seeds.

The normalized mean absolute error (MAE) is used to evaluate semantic drift. Given a sequence of predicted semantic drift scores  $\hat{y}_i$  and the corresponding multi-granularity timeline labels  $y_i$  for a sample, then the normalization is applied by

$$\tilde{y}_i = 3 \times \frac{\hat{y}_i - \min(\hat{y})}{\max(\hat{y}) - \min(\hat{y})}, \quad (21)$$

The predicted value is mapped into the range  $[0, 3]$  (Zhong et al., 2022), which is consistent with the granularity labels. The MAE for this sample is then computed as:

$$\text{MAE} = \frac{1}{T - c} \sum_{t=c+1}^T |\tilde{y}_i - y_i| \quad (22)$$

where  $T - c$  is the number of time steps considered. This provides a unified metric for comparing semantic drift prediction across different models.

### 4.3 Baseline Models

To test the effectiveness of LLM-SAN, the following models are considered as baselines:

**Rolling LDA** (Rieger et al., 2022) detects topic changes by monitoring variations in word distributions across consecutive time windows.

**DETM** (Dieng et al., 2019) models semantic evolution using trainable topic embeddings and LSTM-based inference to identify topic change points.

**ANTM** (Rahimi et al., 2024) leverages LLM-generated document embeddings and sliding-window clustering to analyze changes in topic distributions.

All three models are originally designed for coarse-grained, time-slice-level detection. To enable post-level evaluation, they are uniformly adapted: each post is treated as a single time point, text is encoded with Qwen3, and evaluation metrics are adjusted to measure post-level topic change. Additional modifications ensure compatibility with post-level semantic drift analysis and the temporal dynamics captured by LLM-SAN. See also Appendix C for details.

### 4.4 Post-level Semantic Drift Evaluation Performance

To evaluate LLM-SAN on post-level semantic drift prediction, experiments are conducted on the DTELS dataset with post-level labels, and compared with baseline models. As shown in Table 2, LLM-SAN achieves the lowest prediction error, reducing the average error by 34.98% over the second-best method. This demonstrates its effectiveness in capturing semantic drift by: (i) Compared with Rolling LDA, LLM-SAN transforms discrete lexical co-occurrences into continuous sentence-level embeddings via the GRU module, capturing accumulated temporal semantic signals

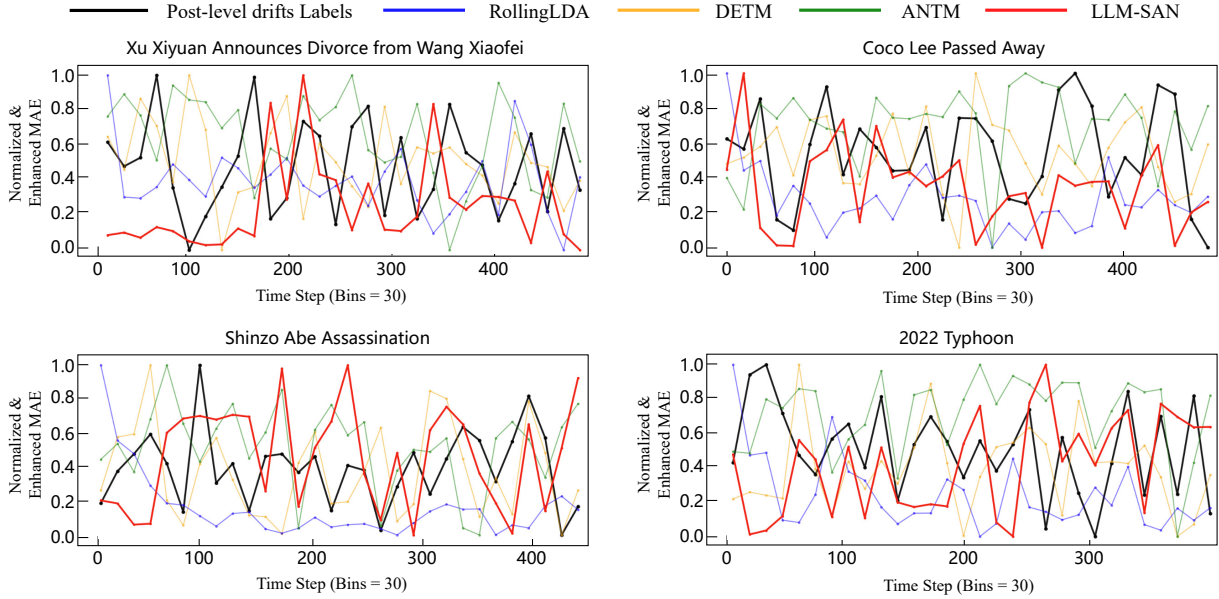


Figure 3: Semantic Drift Trend.

across posts and reducing drift errors. (ii) Compared with DETM, LLM-SAN leverages attention-guided multi-expert fusion to adaptively combine the temporal and topic features at multiple granularities, enhancing sparse keywords and post-level semantic signals, improving drift evaluation accuracy. (iii) Compared with ANTM, LLM-SAN explicitly models semantic dependencies between posts via the GCN module and maintains global semantic continuity through LLM-based encoder and generator features, ensuring robust drift detection even in multi-topic or semantically ambiguous stages.

To further evaluate the methods, the predicted post-level semantic drift values are aggregated into temporal bins of 30 and normalized to  $[0, 1)$ . Subsequently, a Soft-Peak/Valley Enhancement is applied to this normalized curve. This enhancement operates by accentuating local convexity and concavity—conceptually akin to approximating the second derivative—which elevates peaks and deepens valleys. Fig. 3 shows mean drift values for four selected events, where LLM-SAN predictions align most closely with the drift label, capturing both short-term spikes and long-term trends.

#### 4.5 Ablation Study on Feature Modules and Adaptive Fusion

To evaluate the contribution of LLM-SAN components, ablation experiments are conducted on the DTELS dataset. Variants are created by removing the GRU, GCN, or attention-guided multi-expert

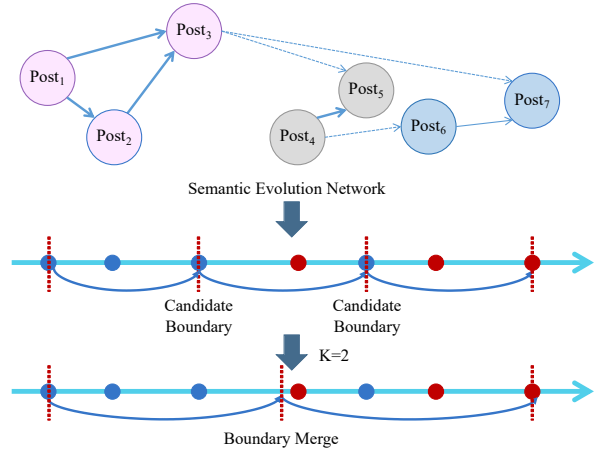


Figure 4: Event Stage Segmentation.

fusion module, and compared with the full model, as shown in Table 3. Several observations can be made: (i) The model generally achieves the lowest errors, indicating that single-dimensional features are insufficient to fully capture the post-level semantic evolution. (ii) Incorporating the attention-guided multi-expert fusion allows the model to adaptively reconcile conflicts between temporal (GRU) and topic (GCN) features, substantially reducing errors.

#### 4.6 Event Stage Segmentation Experiment

As shown in Fig. 4, event stages are inferred from the post-level semantic drift scores produced by the LLM-SAN model. Following established event lifecycle models (Wang et al., 2020), the segmentation identifies stage boundaries that respect temporal

Event ID	Remove GCN	Remove GRU(Qwen)	Remove Expert(Qwen)	LLM-SAN (Qwen)	Remove GRU(Gemma)	Remove Expert(Gemma)	LLM-SAN (Gemma)
1	0.60 ± 0.032	0.58 ± 0.007	0.56 ± 0.017	<b>0.52 ± 0.038</b>	0.57 ± 0.006	0.54 ± 0.013	<b>0.52 ± 0.039</b>
2	0.54 ± 0.026	0.50 ± 0.003	0.49 ± 0.016	<b>0.48 ± 0.019</b>	0.53 ± 0.006	0.50 ± 0.012	<b>0.47 ± 0.037</b>
3	0.59 ± 0.025	0.53 ± 0.003	0.52 ± 0.028	<b>0.50 ± 0.029</b>	0.50 ± 0.004	0.49 ± 0.014	<b>0.48 ± 0.014</b>
4	0.38 ± 0.013	0.33 ± 0.004	0.31 ± 0.006	<b>0.29 ± 0.058</b>	0.33 ± 0.004	0.30 ± 0.015	<b>0.29 ± 0.038</b>
5	0.52 ± 0.034	0.53 ± 0.002	0.53 ± 0.026	<b>0.48 ± 0.038</b>	0.48 ± 0.006	0.47 ± 0.013	<b>0.43 ± 0.031</b>
6	0.39 ± 0.006	0.34 ± 0.004	<b>0.31 ± 0.006</b>	<b>0.31 ± 0.015</b>	0.34 ± 0.006	0.31 ± 0.010	<b>0.30 ± 0.013</b>
7	0.46 ± 0.008	0.45 ± 0.013	<b>0.41 ± 0.02</b>	0.44 ± 0.042	0.42 ± 0.006	<b>0.40 ± 0.011</b>	0.43 ± 0.024
8	0.53 ± 0.010	0.48 ± 0.003	<b>0.42 ± 0.018</b>	<b>0.42 ± 0.050</b>	0.52 ± 0.005	<b>0.48 ± 0.006</b>	0.50 ± 0.031
9	0.46 ± 0.015	0.40 ± 0.004	0.37 ± 0.004	<b>0.36 ± 0.025</b>	0.39 ± 0.008	<b>0.36 ± 0.018</b>	<b>0.36 ± 0.023</b>
10	0.55 ± 0.036	0.47 ± 0.006	0.44 ± 0.012	<b>0.44 ± 0.036</b>	0.44 ± 0.002	0.43 ± 0.027	<b>0.41 ± 0.022</b>
11	0.48 ± 0.008	0.47 ± 0.004	0.44 ± 0.004	<b>0.43 ± 0.040</b>	0.45 ± 0.002	0.44 ± 0.008	<b>0.41 ± 0.038</b>
12	0.52 ± 0.018	0.43 ± 0.005	0.44 ± 0.037	<b>0.43 ± 0.052</b>	0.44 ± 0.008	<b>0.44 ± 0.015</b>	0.47 ± 0.045
13	0.47 ± 0.029	0.43 ± 0.001	<b>0.42 ± 0.010</b>	0.44 ± 0.046	0.48 ± 0.006	<b>0.45 ± 0.016</b>	0.48 ± 0.037
14	0.46 ± 0.034	0.42 ± 0.003	<b>0.42 ± 0.010</b>	0.43 ± 0.026	0.38 ± 0.007	<b>0.37 ± 0.005</b>	0.38 ± 0.014
15	0.43 ± 0.019	<b>0.36 ± 0.004</b>	0.37 ± 0.004	0.41 ± 0.075	0.39 ± 0.008	<b>0.38 ± 0.018</b>	0.41 ± 0.054
16	0.53 ± 0.026	0.53 ± 0.005	0.53 ± 0.027	<b>0.51 ± 0.031</b>	0.55 ± 0.005	0.54 ± 0.028	<b>0.52 ± 0.022</b>
17	<b>0.41 ± 0.012</b>	0.40 ± 0.001	<b>0.38 ± 0.008</b>	0.39 ± 0.019	0.52 ± 0.006	0.45 ± 0.017	0.43 ± 0.048
18	0.66 ± 0.046	0.52 ± 0.008	<b>0.47 ± 0.014</b>	0.48 ± 0.075	0.57 ± 0.012	0.49 ± 0.036	<b>0.45 ± 0.032</b>

Table 3: Ablation Experiment Performance on DTELS.

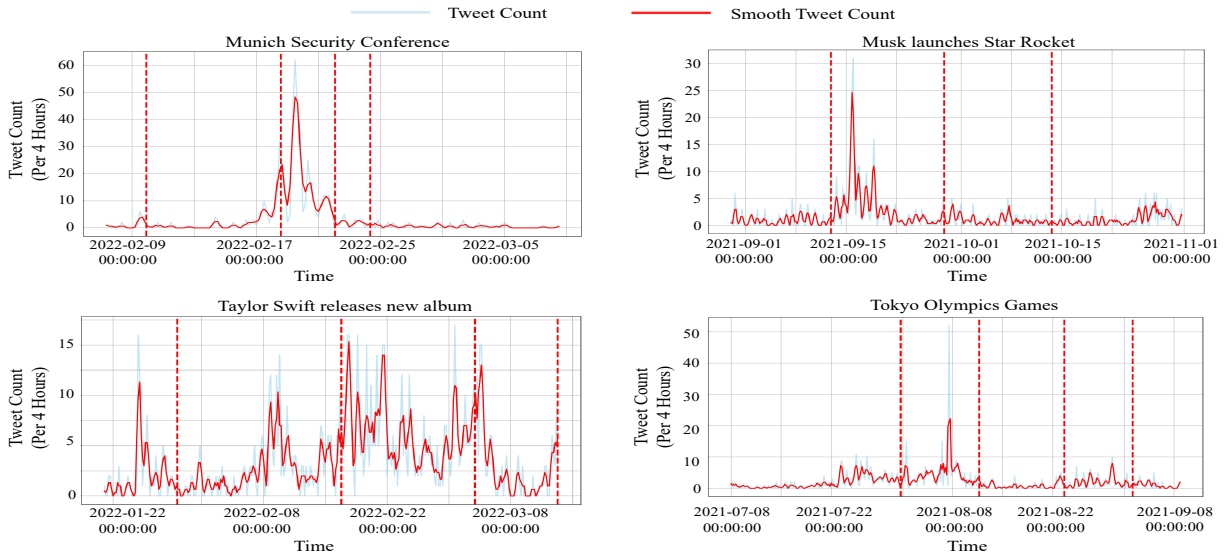


Figure 5: Event Stage Segmentation Experiment.

continuity while capturing major fluctuations in posting activity. The procedure selects timesteps with relatively low semantic drift scores and constructs a semantic evolution network, where each target post is connected to posts contributing to its prediction. Small connected subgraphs are removed to focus on major propagation patterns. Candidate boundaries are then extracted from the latest posts in retained subgraphs, and an iterative merging process ensures that the total number of boundaries does not exceed a predefined limit  $K \in \mathbb{Z}^+$  while maintaining reasonable temporal intervals.

To evaluate the proposed event stage segmentation, experiments are conducted on TwiBot-22 datasets, which offer a large-scale and diverse ba-

sis for analyzing event dynamics. In these experiments, the maximum number of stages is set to  $K = 5$ . Event posts are aggregated in 4-hour windows, smoothed with a three-window moving average, and stage boundaries are visualized as vertical lines on the posting activity curves, as shown in Fig. 5.

Stage boundaries generally occur near small local peaks in low-error regions, marking local maxima or centers of discussion intensity, while adjacent stages are separated by pronounced peaks in post volume, reflecting bursts of activity. These stages do not always correspond to large semantic shifts; instead, they indicate segments where semantic continuity is relatively coherent within each

stage but weaker across stages.

## 5 Conclusion

We propose LLM-SAN for post-level semantic drift evaluation. LLM-SAN integrates a GRU-based module, with inputs from an LLM-based encoder, to capture accumulated temporal semantic signals; a GCN-based module, with inputs from an LLM-based generator, to enhance topic representations and sparse keywords; and an attention-guided multi-expert mechanism to adaptively fuse multi-scale semantic features. Experiments on the DTELS dataset show that LLM-SAN consistently outperforms baseline models in predicting semantic drift, and ablation studies confirm the contribution of each component. Further experiments on the TwiBot-22 dataset demonstrate that LLM-SAN prediction errors can reliably guide event stage segmentation. LLM-SAN captures semantic evolution and can be applied to tasks such as modeling discourse dynamics and analyzing temporal semantic patterns, though its application may require careful adaptation to the specific requirements of different tasks.

## Limitations

Despite its effectiveness in post-level semantic drift evaluation and event stage segmentation, LLM-SAN has some limitations: (i) evaluation is limited to Chinese and English, leaving multilingual analysis for future work; (ii) other signals, such as user interactions, are not explicitly modeled and could further enhance performance; (iii) semantic drift scores, while valuable for analysis, also warrant careful consideration regarding potential misuse—for instance, to steer public discourse or monitor opinions at scale—and should therefore be interpreted and applied with appropriate contextual awareness.

## Acknowledgements

This work was supported by the National Natural Science Foundation of China (T2293771), the China Postdoctoral Science Foundation under Grant Number (2025M771514).

## References

Nishant Balepur, Shivam Agarwal, Karthik Venkat Ramanan, Susik Yoon, Diyi Yang, and Jiawei Han. 2023. DynaMiTE: discovering explosive topic evolutions

with user guidance. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 194–217.

David M Blei and John D Lafferty. 2006. Dynamic topic models. In *International Conference on Machine Learning*, pages 113–120.

David M Blei, Andrew Y Ng, and Michael I Jordan. 2003. Latent dirichlet allocation. *Journal of Machine Learning Research*, 3(Jan):993–1022.

Allaa Boutaleb, Jerome Picault, and Guillaume Grosjean. 2024. BERTrend: neural topic modeling for emerging trends detection. In *Proceedings of the Workshop on the Future of Event Detection*, pages 1–17.

Pengfei Cao, Yupu Hao, Yubo Chen, Kang Liu, Jiexin Xu, Huaijun Li, Xiaojian Jiang, and Jun Zhao. 2023. Event ontology completion with hierarchical structure evolution networks. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 306–320.

Uttam Chauhan and Apurva Shah. 2021. Topic modeling using latent dirichlet allocation: a survey. *ACM Computing Surveys*, 54(7):1–35.

Adji B Dieng, Francisco JR Ruiz, and David M Blei. 2019. The dynamic embedded topic model. *arXiv preprint arXiv:1907.05545*.

Sabbir Mollah et al. 2025. [The telephone game: Evaluating semantic drift in unified models](#). *Preprint*, arXiv:2509.04438.

Wentao Fan, Zhiyan Guo, Nizar Bouguila, and Wenjuan Hou. 2021. Clustering-based online news topic detection and tracking through hierarchical bayesian nonparametric models. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pages 2126–2130.

Shangbin Feng, Zhaoxuan Tan, Herun Wan, Ningnan Wang, Zilong Chen, Binchi Zhang, Qinghua Zheng, Wenqian Zhang, Zhenyu Lei, Shujie Yang, and 1 others. 2022. Twibot-22: towards graph-based twitter bot detection. *Advances in Neural Information Processing Systems*, 35:35254–35269.

Manuel Frank and Haithem Affi. 2025. [Pteb: Towards robust text embedding evaluation via stochastic paraphrasing at evaluation time with llms](#).

Rosario Giuseppe Garroppo, Mohamed Ahmed, Save-rio Niccolini, and Maurizio Dusi. 2018. A vocabulary for growth: topic modeling of content popularity evolution. *IEEE Transactions on Multimedia*, 20(10):2683–2692.

Wolfgang Gaul and Dominique Vincent. 2017. Evaluation of the evolution of relationships between topics over time. *Advances in Data Analysis and Classification*, 11(1):159–178.

- Astha Goyal and Indu Kashyap. 2022. Latent dirichlet allocation-an approach for topic discovery. In *2022 International Conference on Machine Learning, Big Data, Cloud and Parallel Computing*, volume 1, pages 97–102.
- Yanyan He, Mingyu Kang, Duxin Chen, and Wenwu Yu. 2024. Nhce: a neural high-order causal entropy algorithm for disentangling coupling dynamics. *IEEE Transactions on Network Science and Engineering*, 11(6):5930–5942.
- Mingyu Kang, Duxin Chen, Ning Meng, Gang Yan, and Wenwu Yu. 2024. Identifying unique spatial-temporal bayesian network without markov equivalence. *IEEE Transactions on Artificial Intelligence*. DOI: 10.1109/TAI.2024.3483188.
- Andres Karjus, Richard A Blythe, Simon Kirby, and Kenny Smith. 2020. Quantifying the dynamics of topical fluctuations in language. *Language Dynamics and Change*, 10(1):86–125.
- Nozomu Miyamoto, Masaru Isonuma, Sho Takase, Junichiro Mori, and Ichiro Sakata. 2023. Dynamic structured neural topic model with self-attention mechanism. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 5916–5930.
- Megha Ashok Patil, Sunil Kumar, and Sandeep Kumar. 2025a. Gracod: a disruptive graph-aware drift detection algorithm with a gcn-based time-varying module for concept drift detection in short text streams. *International Journal on Digital Libraries*, 26(4):1–26.
- Megha Ashok Patil, Sunil Kumar, and Sandeep Kumar. 2025b. Graphdrift-net: a dynamic graph-based framework for concept drift detection in short unstructured text streams. *The European Physical Journal Plus*, 140(9):884.
- Francesco Periti and Nina Tahmasebi. 2024. A systematic comparison of contextualized word embeddings for lexical semantic change. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4262–4282.
- Hamed Rahimi, Hubert Naacke, Camelia Constantin, and Bernd Amann. 2024. Antm: aligned neural topic models for exploring evolving topics. In *Transactions on Large-Scale Data-and Knowledge-Centered Systems LVI: Special Issue on Data Management-Principles, Technologies, and Applications*, pages 76–97. Springer.
- Jonas Rieger, Kai-Robin Lange, Jonathan Flossdorf, and Carsten Jentsch. 2022. Dynamic change detection in topics based on rolling lidas. In *Text2Story@ ECIR*, pages 5–13.
- Anirudh Sekar, Mrinal Agarwal, Rachel Sharma, Akitsugu Tanaka, Jasmine Zhang, Arjun Damerla, and Kevin Zhu. 2026. Zero-shot embedding drift detection: A lightweight defense against prompt injections in llms.
- Sayan Unankard and Wanvimol Nadee. 2020. Sub-events tracking from social network based on the relationships between topics. In *2020 Joint International Conference on Digital Arts, Media and Technology with ECTI Northern Section Conference on Electrical, Electronics, Computer and Telecommunications Engineering (ECTI DAMT & NCON)*, pages 1–6.
- Kangao Wang, Qizhi Qiu, Mi Wu, and Junan Qiu. 2020. Topic analysis of internet public opinion on natural disasters based on time division. In *2020 3rd International Conference on Advanced Electronic Materials, Computers and Software Engineering*, pages 5–10.
- Rui Wang, Deyu Zhou, and Yulan He. 2019. Atm: adversarial-neural topic model. *Information Processing & Management*, 56(6):102098.
- Xuerui Wang and Andrew McCallum. 2006. Topics over time: a non-markov continuous-time model of topical trends. In *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 424–433.
- Yunli Wang and Cyril Goutte. 2018. Real-time change point detection using on-line topic models. In *Proceedings of the 27th international conference on computational linguistics*, pages 2505–2515.
- Xiaobao Wu, Thong Nguyen, and Anh Tuan Luu. 2024. A survey on neural topic models: methods, applications, and challenges. *Artificial Intelligence Review*, 57(2):18.
- Chenlong Zhang, Tong Zhou, Pengfei Cao, Zhuoran Jin, Yubo Chen, Kang Liu, and Jun Zhao. 2025. Dtels: towards dynamic granularity of timeline summarization. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 2682–2703.
- Delvin Ce Zhang and Hady Lauw. 2022. Dynamic topic models for temporal document networks. In *International Conference on Machine Learning*, pages 26281–26292.
- Ming Zhong, Yang Liu, Suyu Ge, Yuning Mao, Yizhu Jiao, Xingxing Zhang, Yichong Xu, Chenguang Zhu, Michael Zeng, and Jiawei Han. 2022. [Unsupervised multi-granularity summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 4980–4995, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Sulong Zhou, Pengyu Kan, Qunying Huang, and Janet Silbernagel. 2023. A guided latent dirichlet allocation approach to investigate real-time latent topics of twitter data during hurricane laura. *Journal of Information Science*, 49(2):465–479.

## A DTELS Events Details

Event ID	Event Name	Num. of Posts	Time Span
1	Palestine-Israel Military Conflict	422	2019-09-27 00:00:00 – 2024-04-29 07:55:49
2	China Space Station Launch	529	2020-06-23 00:00:00 – 2024-04-30 05:17:22
3	Fukushima Nuclear Waste Discharge	382	2019-10-08 08:18:00 – 2024-04-30 02:38:00
4	Xu Xiyuan Announces Divorce from Wang Xiaofei	482	2021-03-22 14:50:00 – 2024-05-01 00:00:00
5	Evergrande Debt Crisis	312	2020-09-07 11:41:54 – 2024-04-30 02:38:00
6	Lao Rongzhi Case	452	2019-11-27 00:00:00 – 2024-04-30 02:38:00
7	Russia-Ukraine Situation	531	2020-03-11 00:00:00 – 2024-04-30 02:38:00
8	Pelosi Visits Taiwan	446	2020-02-04 00:00:00 – 2024-04-30 04:05:00
9	Tangshan Barbecue Shop Assault	414	2022-05-13 00:00:00 – 2024-04-30 04:52:00
10	2022 Typhoon	386	2020-09-09 23:59:59 – 2038-09-08 00:00:00
11	Shinzo Abe Assassination	430	2020-08-28 00:00:00 – 2024-04-29 23:36:00
12	PLA “Lock Taiwan” Drill	367	2022-07-29 16:18:00 – 2024-04-30 02:38:00
13	Southern Turkey Earthquakes	374	2022-01-08 00:00:00 – 2024-04-30 02:38:00
14	Trump Indictment	345	2020-04-20 11:44:00 – 2024-04-30 02:38:00
15	Myanmar Telecom Fraud	490	2019-12-25 00:00:00 – 2024-04-29 23:35:00
16	Armed Conflict in Sudan Capital	470	2019-10-25 00:00:00 – 2024-04-30 02:38:00
17	Coco Lee Passed Away	455	2020-07-13 00:00:00 – 2024-04-30 02:38:00
18	2023 Typhoon	352	2022-08-01 00:00:00 – 2024-05-01 00:39:05

Tabel 4: DTELS Events Details.

## B Twibot22 Events Details

Event Name	Keywords	Num. of Posts	Time Span
Munich Security Conference	"Munich Security Conference", "MSC"	539	2022-02-03 00:00:00 – 2022-03-20 00:00:00
Elon Musk SpaceX Launch	"SpaceX", "Inspiration4", "civilian spaceflight"	597	2021-08-30 00:00:00 – 2021-10-31 00:00:00
Taylor Swift New Album	"Taylor Swift", "paparazzi", "Ed Sheeran"	1103	2022-01-21 00:00:00 – 2022-03-15 00:00:00
Tokyo Olympics	"Tokyo Olympics", "Olympic Games 2021", "gold medal"	783	2021-07-08 00:00:00 – 2021-09-08 00:00:00

Tabel 5: Twibot22 Events Details.

## C Models Adjusted Details

Model Name	Adjustments for Semantic Drift Evaluation
Rolling LDA	Treat each post as a single time point; use previous 12 posts as a rolling memory window; measure semantic drift with Jensen-Shannon divergence; update model incrementally to adapt to new words and small semantic changes.
DETM	Treat each post individually; encode words with Qwen model; use reparameterization and KL regularization to handle irregular intervals; approximate inference via neural network + reparameterization + mini-batch optimization for efficiency; semantic drift computed as sum of metrics (KL and NLL) between consecutive posts.
ANTM	Treat each post individually; preprocess for Chinese text; encode posts with Qwen; retain unassigned topics as separate; combine puw and jaccard as drift metric; use window of size 12 with overlap 11 to capture drift relative to previous posts.

Tabel 6: Models Adjusted for Semantic Drift Evaluation.