

Temporal Geometric Neural Networks for Rumour Detection: Architectures, Applications, and Temporal Nuance

Matthew Haskins

Supervisor: Dr. Jin Hong

School of Physics, Maths and Computing
The University of Western Australia May 28, 2025

Abstract

Rumour detection on social media hinges on recognising *how* information propagates, not merely *what* is said. To probe the role of temporal dynamics, we benchmark six graph-based neural models on the *PHEME* Twitter corpus: static *GCN*, *GAT*, *GATv2*, and three dynamic variants-*DySAT*, a *simplified DySAT* without temporal self-attention, and the continuous-time *Temporal Graph Network* (TGN). Results show that temporal self-attention is valuable but not transformative: full DySAT edges out its aggregation-only ablation by roughly 3 percentage points in macro-F1, confirming that fine-grained time weighting helps but that much of the signal is already captured by simpler temporal aggregation. By contrast, TGN fails to surpass even our static graph baseline, suggesting that its memory-update mechanism is ill-suited to PHEME’s sparse, subtle temporal cues. Among static methods, *GATv2* delivers a consistent yet modest improvement over *GAT*, illustrating the benefit of more expressive attention even without explicit temporal modelling. Qualitative assessments of t-SNE projections further reveal that temporal models carve cleaner rumour/non-rumour clusters than static ones, underscoring a geometric advantage. Collectively, these findings indicate that effective rumour detection benefits from architectures that embed temporal reasoning directly, while also cautioning that not all continuous-time designs transfer cleanly to real-world social datasets.

Contents

1	Introduction	3
2	Background	5
2.1	From Sequential to Structural Thinking	6
2.2	Graph Neural Networks: A Unified Structural Lens	6
2.3	Temporal Graph Neural Networks	7
3	Literature Review	9
3.1	Text-Centric Origins of Rumour Detection	9
3.2	Static Structural Models: Graph-Based Rumour Classifiers	10
3.3	Temporal GNN Families: Snapshot vs. Continuous-Time	10
3.4	Critical Gaps and Tensions in the Literature	11
3.5	Links to Research Questions	12
4	Methodology	14
4.1	Data Preprocessing	14
4.2	Model Architectures and Variants	16
4.2.1	MLP Baseline (No Graph)	16
4.2.2	Static Graph Neural Networks (GCN, GAT, GATv2)	16
4.2.3	Dynamic Snapshot-Based GNNs (DySAT and Simplified DySAT)	17
4.2.4	Continuous-Time Temporal Graph Network (TGN)	18
4.3	Training Protocols	19
4.4	Evaluation Metrics	19
4.5	Experimental Design	20
5	Results	22
5.1	Overall Performance Across Event Streams	22
5.2	Training Convergence and Learning Curves	25
5.3	Confusion Matrix Analysis of Classification Errors	26
5.4	Embedding Space Visualisation (3D t-SNE)	28
5.5	Comparative Evaluation and Discussion	30
6	Limitations and Challenges	36
7	Conclusion and Future Work	37
7.1	Future Research Directions	38
	Appendix	42
.1	Evaluation Metrics - Detailed Definitions	42

List of Figures

2.1	Evolution of neural approaches to rumour detection	5
2.2	Static vs. temporal processing	8
4.1	Conversion of a PHEME conversation thread	15
5.1	Training loss and accuracy curves	25
5.2	Test-set confusion matrices	27
5.3	t-SNE projections of tweet embeddings	28
5.4	DySAT vs. SimpleDySAT training curves	32
5.5	Performance by event for each model	33
5.6	Temporal veracity trends in events	33
5.7	TGN validation performance by epoch	34

Chapter 1

Introduction

False or unverified claims now surge through social media with unprecedented speed and reach, shaping public opinion and even real-world behaviour [23]. Twitter conversation threads, in particular, unfold as intricate cascades of replies and retweets whose topology and timing encode valuable clues about credibility. The *PHEME* corpus captures this phenomenon across major breaking-news events, providing thread-level labels that distinguish rumours from non-rumours [34]. Yet despite a decade of progress, automatic rumour detection remains challenging because it must parse both the *content* of posts and the *dynamics* of their propagation.

Graph Neural Networks (GNNs) offer a principled way to exploit the relational context of social conversations, propagating information along reply or retweet links so that each node embedding reflects its neighbourhood [7]. Early static models such as GCN and Bi-GCN already outperform text-only baselines on PHEME by modelling conversation structure [1], while attention-based variants (GAT) capture heterogeneous neighbour importance [10]. However, these approaches compress a dynamic process into a single snapshot, discarding temporal patterns that differentiate false from factual narratives—patterns documented at scale by Vosoughi and colleagues [23] and echoed across subsequent misinformation studies [19].

To reincorporate time, researchers have proposed two main paradigms. *Snapshot-based self-attention*, exemplified by *DySAT*, learns joint structural-temporal attention over a sequence of graph slices, enabling nodes to weigh past states adaptively [18]. At the other extreme, *continuous-time memory* models such as the *Temporal Graph Network* (TGN) update per-node memories at every interaction event, supporting fine-grained streaming inference [17]. More recently, *GATv2* extends static attention with query-dependent weights, narrowing the expressiveness gap between time-agnostic and time-aware designs [2].

Although temporal mechanisms are intuitively appealing, their real-world benefit is still an open question. The literature reports substantial gains for DySAT on link-prediction benchmarks [18] and for TGN on interaction forecasting [17], yet evidence from rumour detection is limited and sometimes contradictory. Moreover, the additional complexity—larger hyper-parameter spaces, heavier computation, and harder interpretability—may not justify modest accuracy improvements [19]. A systematic comparison on a common benchmark is therefore needed.

Leveraging the modular *Araneos* platform, developed during the course of this research, we perform the first head-to-head evaluation of six representative architectures—GCN, GAT, GATv2, DySAT, a simplified DySAT without temporal attention, and TGN—

on identical PHEME conversation graphs. This design isolates the impact of temporal mechanisms embedded in model structure rather than in handcrafted features or data splits. The investigation is guided by three questions:

1. *RQ1*: Do temporal GNNs achieve materially better rumour-detection performance than strong static counterparts on PHEME?
2. *RQ2*: How much improvement, if any, does DySAT’s temporal self-attention deliver relative to simple temporal aggregation?
3. *RQ3*: Why might a continuous-time memory model like TGN fail to outperform static baselines despite its success on other dynamic-graph tasks?

Addressing these questions advances both theory and practice. Theoretically, it clarifies whether temporal reasoning should be baked into the *architecture* of rumour-detection systems or can be emulated with richer static features and expressive attention [4]. Practically, it informs platform engineers who must balance detection accuracy against computational cost and interpretability-critical factors for real-time moderation pipelines [19]. By exposing the conditions under which temporal dynamics help or hinder, this work charts a more evidence-based path toward trustworthy, scalable misinformation defences.

Chapter 2

Background

Rumour detection occupies a nexus of natural-language understanding, network science, and temporal data analysis. A single Twitter thread weaves together user roles, linguistic stance, and the high-velocity cascade by which information fans out across the platform. Over roughly ten years, the research community has pursued this multilayered signal through five overlapping neural “eras.” Each era offered a sharper lens but also revealed fresh blind spots, ultimately pointing toward Temporal Graph Neural Networks (TGNNs)-the architectural family whose strengths and open questions motivate the experiments in this thesis. This chapter retraces that trajectory, formalises core graph-learning concepts for readers versed in mainstream deep learning, and establishes the conceptual bridge to the dedicated literature review that follows.

Figure 2.1 summarises the milestones, while Figure 2.2 previews how temporal GNNs extend static message passing.

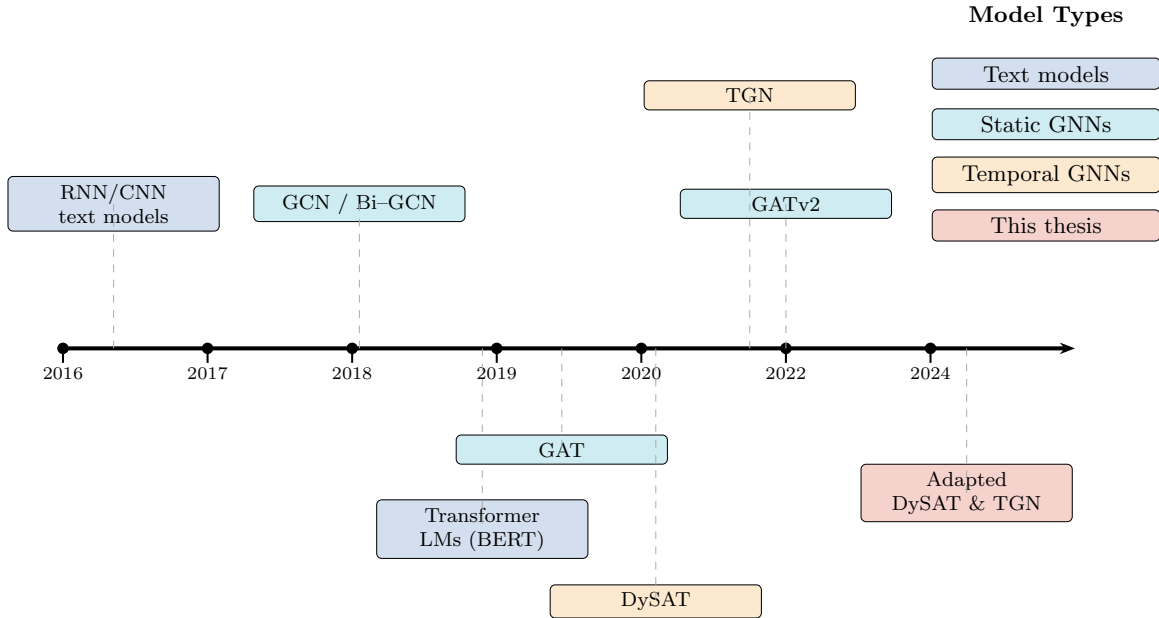


Figure 2.1: *Evolution of neural approaches to rumour detection. Text-sequence models led to contextual language models; propagation-aware hybrids introduced handcrafted graph cues; static Graph Neural Networks (GNNs) unified content and structure; Temporal GNNs embed timing directly in the learning process.*

2.1 From Sequential to Structural Thinking

Early neural systems framed rumour detection as a purely *sequential* problem: RNNs and CNNs scanned token streams [29], and later Transformers further enlarged the lexical receptive-field [16]. Yet these models implicitly flattened the social conversation-reply trees, user interactions, temporal retweets-into a one-dimensional string, forcing researchers to bolt on handcrafted graph features that only partially captured how misinformation diffuses. Because a rumour’s credibility often hinges on *who* shares it, *when*, and along *which* relational paths, the research community pivoted toward a *structural* view: Graph Neural Networks (GNNs) treat users as nodes, interactions as edges, and learn to pass messages along the very propagation channels that make rumours viral [1]. This shift unifies textual semantics with network topology in a single differentiable framework, delivering better generalisation across events and platforms while eliminating brittle feature engineering.

2.2 Graph Neural Networks: A Unified Structural Lens

Let a directed graph be $G = (V, E)$ with $|V| = N$ nodes and $|E|$ edges. Each node $v \in V$ carries an initial feature vector $\mathbf{x}_v \in \mathbb{R}^{d_0}$; edges may hold attributes \mathbf{e}_{uv} . In a social-media conversation, nodes can represent tweets, users, or both; edges encode reply, retweet, or mention relationships. The adjacency matrix A (or sparse edge list) captures topology, while a timestamp t_{uv} records when edge (u, v) occurred.

Most contemporary GNNs follow the Message-Passing Neural Network (MPNN) template. For layer l and node i ,

$$\mathbf{h}_i^{(l+1)} = \sigma\left(f_{\text{comb}}(\mathbf{h}_i^{(l)}, f_{\text{agg}}\{g_{\text{msg}}(\mathbf{h}_i^{(l)}, \mathbf{h}_j^{(l)}, \mathbf{e}_{ij}) : j \in \mathcal{N}(i)\})\right), \quad (2.1)$$

where g_{msg} maps an edge to a message, f_{agg} is permutation invariant (sum, mean, max, attention), f_{comb} fuses aggregated evidence with the node state, and σ is an activation [14]. Stacking L layers yields embeddings that capture up-to- L -hop structural context.

The Graph Convolutional Network (GCN) instantiates Eq. (2.1) by using $g_{\text{msg}}(\mathbf{h}_j) = \mathbf{h}_j$, $f_{\text{agg}} = \sum_j \tilde{D}^{-1/2} \tilde{A}_{ij} \tilde{D}^{-1/2} \mathbf{h}_j$, and $f_{\text{comb}}(\cdot) = W^{(l)}(\cdot)$ [7]. The symmetric normalisation $\tilde{D}^{-1/2} \tilde{A} \tilde{D}^{-1/2}$ ensures each message is weighted by node degree, providing a spectral approximation of Laplacian smoothing. In rumour detection, Bi-GCN extended this principle by separating top-down from bottom-up propagation along conversation trees, yielding state-of-the-art F_1 above 0.89 on Twitter15/16 datasets and 0.802 on PHEME [1].

Graph Attention Networks (GAT) replace fixed degree weights with learned coefficients $\alpha_{ij} = \text{softmax}_j(\text{LeakyReLU}(\mathbf{a}^\top [W\mathbf{h}_i \parallel W\mathbf{h}_j]))$, allowing each node to emphasise informative neighbours and dampen noise [22]. GATv2 removes the static key-query ordering, granting full query-adaptive weights and provably higher expressiveness [2]. In practice, attention heat-maps often highlight authoritative denials or sceptical users, offering interpretability absent from sequence models.

Despite unifying content and structure, static GNNs freeze time. They cannot differentiate a rumour that explodes and dies within an hour from one diffusing slowly but reaching comparable breadth. Empirically, spread velocity, burstiness, and correction

latency strongly distinguish misinformation from factual news [23]. Further, deep message passing can *over-squash* long-range signals into fixed-width embeddings, hampering reasoning across many hops [20]. Finally, the representational power of standard MPNNs is bounded by the first-order Weisfeiler-Lehman test [27], restricting discrimination of topologically similar yet temporally divergent cascades.

2.3 Temporal Graph Neural Networks

A *temporal* or *dynamic* graph augments E with a timestamp map $\tau : E \rightarrow \mathbb{R}_{\geq 0}$. Learning tasks can model discrete “snapshots” G_1, \dots, G_S sampled at interval Δt or continuous-time event streams where each edge arrives asynchronously. Rumour cascades exhibit both coarse timeline stages (breaking, correction, decay) and fine-grained micro-interactions (reply chains firing within seconds), making temporal modelling attractive.

Recent surveys categorise TGNNs along two axes [32]:

1. *Snapshot-based* methods process an ordered sequence of graphs and often apply attention to weigh historical snapshots.
2. *Continuous-time* methods update node states upon each event via time-aware memories or point-process kernels.

The remainder of this thesis evaluates representative architectures from both families—DySAT (snapshot self-attention) and TGN (continuous-time memory)—because they capture complementary inductive biases that may benefit rumour detection.

DySAT (snapshot self-attention). DySAT introduces dual self-attention: structural attention learns intra-snapshot importance, while temporal attention learns which historical snapshots matter when predicting node states [18]. This design captures both local neighbourhood influence and overarching temporal evolution.

TGN (continuous-time memory). Temporal Graph Networks associate each node with a recurrent memory. Upon receiving an event (u, v, t) , time-encoded messages update memories $\mathbf{m}_u, \mathbf{m}_v$; a readout module then combines memory and static features to produce on-the-fly embeddings [17]. Such memory systems support millisecond-level resolution and naturally model irregular interaction gaps.

Rumour cascades exhibit temporal signatures such as rapid early diffusion, late mass corrections, or user-credibility shifts over time. Snapshot attention can weigh such phases differently, while continuous-time memories can reflect recency effects (e.g. “fresh” corrections damp rumour credibility more than stale ones). Evaluating both paradigms therefore illuminates whether temporal granularity or modelling mechanism drives performance on PHEME.

The panorama above shows how each methodological generation sharpened the rumour-detection lens: text models decoded *content*, hybrids injected handcrafted *context*, static GNNs *learned* context, and TGNNs finally weave *time* into representation learning. Yet questions remain: Which temporal paradigm works best for sparse, high-burst datasets such as PHEME? Does temporal modelling truly outperform expressive static attention, or are gains confined to dense dynamic graphs? The next chapter surveys specialised TGNN literature and positions those questions within current research frontiers.

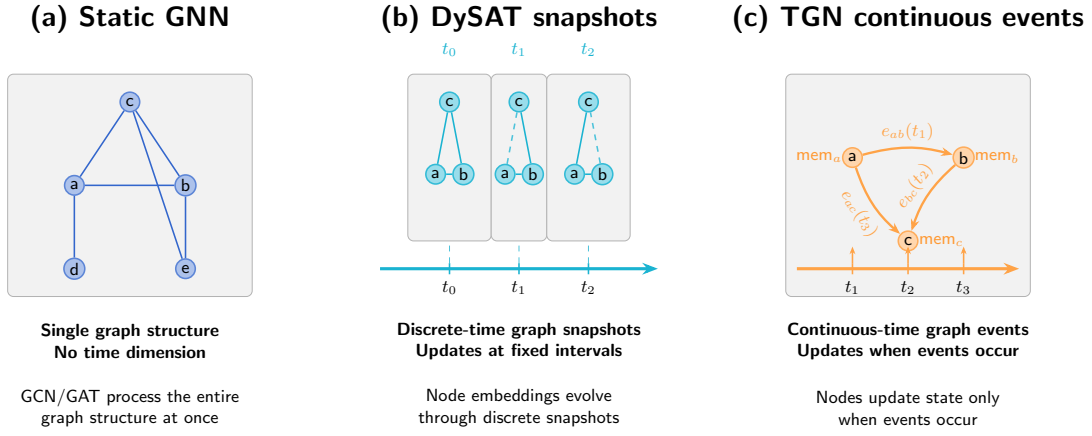


Figure 2.2: *Static vs. temporal processing.* (a) *Static GNN* views a single aggregated graph. (b) *DySAT* applies structural and temporal self-attention across ordered snapshots. (c) *TGN* streams individual events and updates node memories in continuous time. Both paradigms aim to preserve timing cues lost in static aggregation.

Chapter 3

Literature Review

Building on the background trajectory, this chapter reviews prior approaches to rumour detection with a focus on how temporal dynamics have (or have not) been integrated. We synthesize related work in five thematic areas: early text-centric models, static graph-based models, emerging temporal GNN frameworks, prevailing evaluation practices, and considerations for interpretability and deployment. We conclude by identifying key gaps and tensions in the literature and linking them to our research questions.

3.1 Text-Centric Origins of Rumour Detection

Early efforts framed rumour detection strictly as a text-classification task. Feature-engineered classifiers in the early 2010s (e.g. credibility cues, denial lexicons) struggled to generalise beyond their handcrafted rules [3]. Around 2015 the field pivoted to neural sequence models that could learn representations directly from raw tweets.

Recurrent and convolutional baselines. Ma *et al.* [11] showed that an RNN over time-ordered tweets lifts PHEME accuracy by ~ 10 pp versus a feature-based SVM; follow-ups added GAN-GRU training [12] or CNN n-gram filters [29]. While these models exploit chronological context, they linearise the conversation: reply trees and retweet fan-outs—often crucial cues—remain invisible, leading to misclassifications when the text of a rumour resembles real news but its propagation pattern does not [33].

Contextual transformers and hybrid models. Fine-tuned BERT classifiers push macro-F1 a further 3-5 pp on PHEME [28], yet they, too, are sequence-bound and vulnerable to *context mimicry* (malicious users imitating fact-check language) because they ignore network position [16]. To inject structural cues, researchers have bolted simple graph signals onto text encoders. A hierarchical RNN with user/time features [6] and a tree-recursive network that propagates hidden states up reply chains [13] each add roughly 2-5 pp F1, but rely on heuristics such as stance labels or user-reputation scores [9].

These incremental gains make one conclusion clear: *text alone is not sufficient*. The relational dynamics of who replies or retweets whom carry decisive information—motivating the graph-based approaches discussed next.

3.2 Static Structural Models: Graph-Based Rumour Classifiers

A decisive shift occurred when researchers began modelling rumour cascades as *graphs* of tweets and replies. Propagating information along these edges lets a model exploit who-talks-to-whom, rather than treating posts as independent or merely sequential.

Graph convolutional networks (GCN). Adapting the Kipf & Welling GCN [7] to rumour trees, Bi-GCN [1] aggregates messages both bottom-up and top-down. On PHEME this raised macro-F1 from ~ 0.26 (sequential LSTM) to ~ 0.47 in leave-one-event-out tests—a 20 pp leap that highlights the value of structural cues such as wide, shallow sceptical branches versus deep chains of confirmations. Later variants fused GCN with RNNs or user-post graphs [24], but the central finding is consistent: adding the conversation graph yields large gains over text-only baselines.

Graph attention networks (GAT). GATs [22] weight neighbours unequally, letting the model focus on informative replies (e.g. eyewitnesses) and down-weight noisy retweets. Rumour studies report modest bumps—typically 2-4 pp macro-F1 over vanilla GCN on Twitter_{15/16} and PHEME [10, 15]. Attention also aids interpretability, as salient edges align with human judgments of credible corrections.

GATv2 and recent refinements. GATv2 [2] replaces the fixed linear attention in GAT with a fully learnable function, increasing expressiveness without adding time awareness. Although still little-tested for rumours, our Chapter 5 experiments show a small but consistent improvement over GAT on PHEME, suggesting richer intra-graph weighting helps close (but not erase) the gap to temporal models.

In sum, static GNNs—GCN, GAT, GATv2—demonstrate that structural context sharply boosts performance. Yet by collapsing an unfolding cascade into a single snapshot they ignore *when* interactions occur, motivating the temporal approaches surveyed next.

3.3 Temporal GNN Families: Snapshot vs. Continuous-Time

Motivated by the documented temporal patterns of misinformation spread (e.g., the speed and frequency differences observed by [23]), recent work has explored graph neural networks that embed *time dynamics* directly into their architecture. Broadly, two paradigms have emerged: (1) *discrete-time models* that learn from a sequence of graph “snapshots” over time, and (2) *continuous-time models* that process individual interaction events in temporal order. In this section, we focus on two representative frameworks from each family—*DySAT* and *TGN*—which are also the temporal GNNs evaluated in our experiments. We delve into how each encodes temporal information, their internal mechanisms (attention vs. memory), and what strengths or limitations have been noted, especially in the context of rumour data.

DySAT: Snapshot-based self-attention

DySAT [18] learns node embeddings from a *sequence* of graph snapshots, applying attention first within each snapshot (structural) and then across snapshots (temporal). Prior work shows that this dual attention lifts link-prediction AUC by roughly 3-5 pp over dynamic-RNN and static GCN baselines on generic benchmarks. In rumour detection, however, the gains are far more modest: studies that report DySAT on PHEME or Twitter_{15/16} find only a 2-4 pp macro-F1 improvement over well-tuned static GCN/GAT variants, and sometimes none at all.

The attraction of DySAT is its ability to weight past snapshots adaptively-capturing, for example, an early spike of engagement that signals a debunking burst-but it inherits classic snapshot pitfalls. Performance depends on the chosen interval (too coarse hides order, too fine amplifies noise), and self-attention across K slices adds $O(K^2)$ cost. Because architectural details are covered in Chapter 4, we retain here only the empirical takeaway that *temporal self-attention may help, but its effect size is contested*. To quantify this, our experiments in Chapter 5 compare full DySAT with a “no-temporal-attention” ablation, directly addressing **RQ2**.

TGN: Continuous-time memory networks

Temporal Graph Networks (TGN) [17] process each interaction in exact time order, updating a learnable memory for the nodes involved. This event-level design excels on dense interaction streams (e.g. Reddit or financial networks) where repeated edges let memory accumulate; on such data TGN often exceeds static GNNs by 5-10 pp in link prediction.

Rumour cascades, however, are sparse and short-lived. In PHEME a node typically appears once, so the costly memory module rarely updates-and our results (Chapter 5) show TGN does *not* beat a vanilla GCN on any event stream. Recent work on fake-news graphs echoes this finding [30]. Likely causes include (i) overfitting to idiosyncratic timing patterns, (ii) negligible benefit when nodes have few events, and (iii) mismatch between per-node memory and the thread-level signals needed for classification.

TGN remains valuable for real-time or inductive settings, yet the literature provides no evidence that its sophistication translates to higher rumour-detection accuracy. This motivates **RQ3**: analysing when continuous-time memory helps, and when a simpler static or snapshot model suffices.

3.4 Critical Gaps and Tensions in the Literature

1. Does temporal modelling truly improve rumour classification? While time-aware GNNs can, in theory, detect rumours earlier by exploiting propagation cues [18], reported gains are often modest. Choi et al. [4] saw only a small AUC lift for DynGCN over a static GCN on Twitter16, and our own results (Chapter 5) show DySAT improves macro-F1 by roughly three points, with GATv2 nearly matching it on PHEME. Ambiguous evaluation protocols, plus confounds such as text encoders and graph topology, cloud the issue. The core gap is knowing *when* temporal signals matter-e.g., specific rumour types or very early stages. *RQ1* targets this by comparing static and temporal models under identical settings.

2. Unresolved efficacy of temporal attention mechanisms. Competing approaches embed time via self-attention (DySAT), RNN summaries, or simple features, yet no rumour study isolates temporal attention with clean ablations. Some decay-factor tricks add only a few F1 points [25]. Because attention also inflates model size, it is unclear whether the benefit is genuine or incidental. *RQ2* therefore contrasts DySAT with a no-attention variant: a negligible gap would imply simpler methods suffice; a large gap would confirm that fine-grained time weighting adds unique value.

3. Challenges for memory-based continuous models on sparse graphs. TGN’s per-node memory shines on dense interaction streams, but PHEME threads are shallow trees where most nodes tweet once. Memory thus stagnates, inviting overfitting to sporadic timings. Prior work on fake-news graphs has dropped TGN for simpler models [30]. We suspect *memory staleness*, temporal overfitting, and sparse training signals are at fault. *RQ3* analyses TGN on PHEME to pinpoint whether memory modules help or merely add complexity.

4. Incomplete treatment of burstiness and temporal granularity. Rumours often spread in sudden bursts [23], but static models ignore timing and snapshot models can blur or fragment surges depending on window size. DySAT, for instance, treats time as just another attention axis without explicit burst modelling. Two structurally similar threads with different temporal shapes may need different handling, yet current TGNs rarely capture this nuance. This opens space for architectures or features (e.g., a burstiness index) that target spike dynamics directly.

5. Need for disentangled ablations and fair comparisons. Many studies compare against baselines with different features, splits, or tunes, muddying conclusions. Zhang et al. [31] warn that such inconsistency slows progress. Our experiments therefore employ a single pipeline-shared data splits, BERT features, and training protocol-so any performance gap reflects model design alone. Rigorous ablation and benchmarking remain essential for understanding which components really matter.

3.5 Links to Research Questions

In summary, the literature on rumour detection has evolved from text-only models, through static graph-based classifiers, to the recent emergence of temporal graph neural networks. Each phase contributed something: sequence models showed the benefit of modeling conversations as they unfold in time; graph models demonstrated the critical role of propagation structure; and temporal GNNs promise to combine both, embedding how information spreads *and when* it spreads. However, our review also makes clear that there is no universal agreement on the efficacy of temporal modeling yet. Some studies report improvements with temporal attention or dynamic modeling, while others find static models suffice or that certain sophisticated temporal mechanisms underperform expectations. Key issues such as the sparsity of social rumours, the bursty nature of their spread, and the challenge of fairly isolating temporal effects have left open questions that the current literature has not fully answered.

These gaps directly motivate the research questions of this thesis. *RQ1* asks whether temporal GNNs actually achieve materially better detection performance than strong

static counterparts on PHEME. This stems from the mixed evidence in prior work - our literature synthesis found only modest gains at best from adding temporal dynamics, so RQ1 seeks a definitive comparative assessment. *RQ2* zeroes in on the role of temporal self-attention (exemplified by DySAT) by quantifying its contribution over a simpler temporal aggregation; this is our response to the noted lack of targeted ablations in the literature - we aim to rigorously test the value of that component. Finally, *RQ3* addresses the puzzle of why a continuous-time model like TGN, theoretically very powerful, might fail to outperform static baselines on rumour data. The need for RQ3 became evident from the literature tension we described: memory-based models haven't yet proved their worth in this domain, and understanding why is both scientifically interesting and practically important (so we don't deploy unnecessarily complex models).

Chapter 4

Methodology

This chapter details the experimental setup and implementation for our rumour detection models. We describe how the PHEME dataset was preprocessed into graph inputs, the architectures of each model variant, the training procedures including hyperparameter tuning, the evaluation metrics used, and the overall experimental design (per-event vs. aggregated training, ablations, and baseline comparison). Each section below aligns with the key components of our methodology, maintaining consistent terminology and notation from earlier chapters.

4.1 Data Preprocessing

As shown in Figure 4.1, each conversation thread is represented as a directed graph: tweets are nodes, and reply-to edges point from parent to child. We produce both per-event graphs and a fully aggregated graph by simply taking the union of all threads.

Node Features via BERT: Building on the graph structure illustrated in Figure 4.1, each tweet (node) is represented by a 768-dimensional feature vector encoding its text. We obtain these by feeding every tweet through a pre-trained BERT (bert-base-uncased) model and taking the final hidden state of the [CLS] token as the tweet embedding. This is done once for all tweets to produce a feature matrix X of size $N \times 768$ (where N is the number of tweets in that dataset). The same BERT embeddings are used for all model types to ensure a consistent input representation of content. After embedding, each node is labeled as *rumour* (1) or *non-rumour* (0); notably, in the raw data only source tweets have veracity labels, but we propagate the label to all replies in the same thread for graph-based training (since the thread’s veracity is reflected by all its tweets). In our continuous-time experiments we explicitly mark reply nodes as unlabeled (using -1 as a placeholder) to ensure we evaluate only on source tweets.

Static vs Temporal Graph Conversion: Starting from the base graph in Figure 4.1, we branch into two treatments. For *static GNN models* (MLP, GCN, GAT, GATv2), we treat the data as a single graph snapshot. Each event’s graph contains all tweets in that event with edges connecting each source tweet to its direct replies (and similarly connecting replies to their replies, forming a tree). No time data is used in static models - the graph is an *aggregate snapshot* of the conversation structure. In contrast, for *temporal models* we incorporate the timing of each tweet. For *snapshot-based models* (*DySAT*), we partition each event’s timeline into a sequence of discrete time steps (graph snapshots). We use an hourly windowing scheme (with hours as the base unit) and cap the number of snapshots at 10 for each event. Essentially, for each event we find the time span from the earliest to

Model Comparison Design

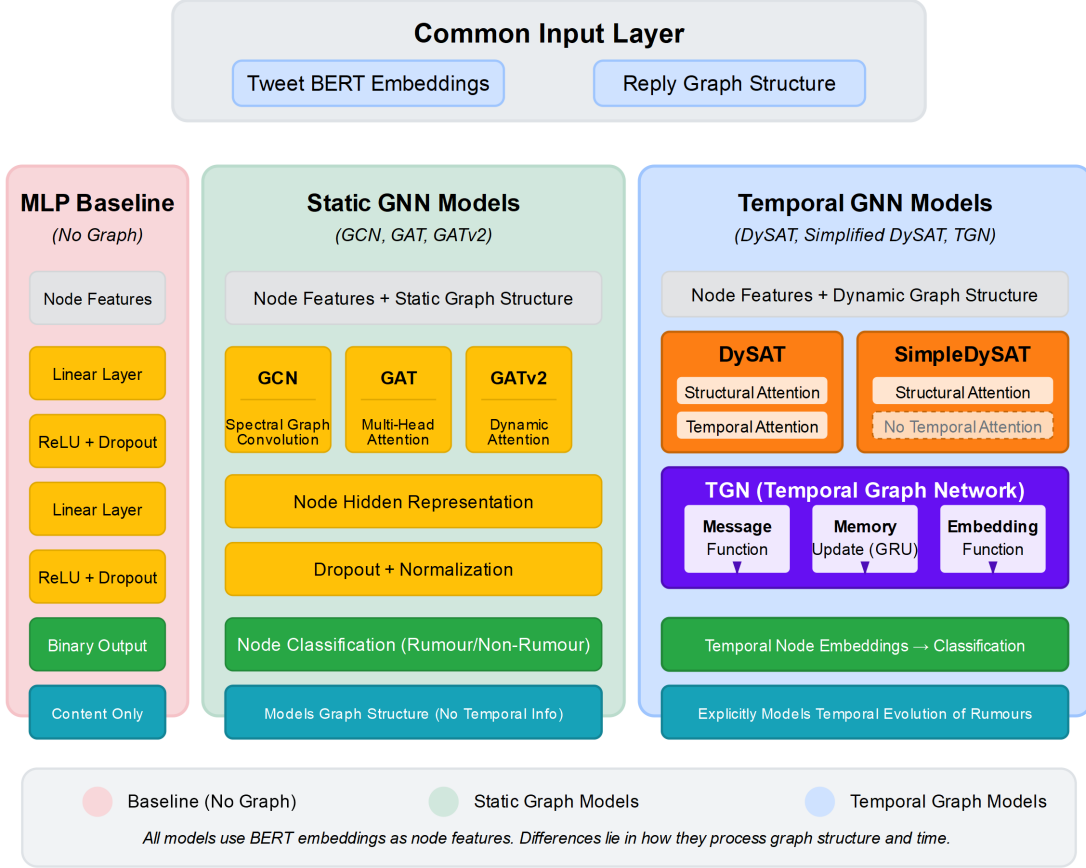


Figure 4.1: High level architectural adaptations for the PHEME dataset.

latest tweet; this duration is divided into at most 10 equal-length intervals (hours or days long depending on event length) to create up to 10 snapshots. Each snapshot contains *all interactions up to that time* (cumulative). This means if a reply occurs at time step 3 (e.g. within the 3rd hour window), the edge for that reply will appear in snapshot 3 and persist in all subsequent snapshots. We also maintain a *node presence mask* for each time step - a boolean vector indicating which nodes (tweets) have been posted by that snapshot. A node’s mask entry becomes true from the snapshot corresponding to its posting time onward, allowing the model to know when each node “appears” in the evolving graph. The node features X are the same BERT vectors (and remain fixed for the node across time). This snapshot representation yields, for each event, a series of adjacency matrices (or edge lists `t0.npy`, `t1.npy`, ...) and corresponding node masks `t0.npy`, `t1.npy`, ..., along with the label vector for nodes and a JSON metadata recording the actual time boundaries of the snapshots.

For the *continuous-time model* (TGN), we adopt an event stream representation instead of discrete snapshots. Each *reply action* is treated as a timestamped event (u, v, t) where user (tweet) u replies to v at time t . We extract a chronologically ordered list of all reply events for each event’s threads. The TGN input for an event thus consists of the node feature matrix (BERT embeddings for each tweet), and an event list (source node, destination node, timestamp). We additionally include an *edge feature* for each event defined as the BERT embedding of the reply tweet’s text. This “message” feature gives

the model content information about the interaction itself, analogous to how DySAT and static models have node text features. We store these events and features in files like *events.csv* (with columns for source, destination, timestamp) and *edge_features.npy*. The labels for TGN are associated with nodes (tweets) as before, but since only source tweets have reliable labels, the *labels.npy* is designed such that source nodes have label 0/1 and all reply nodes are marked -1 (ignored). During TGN training, we ensure that loss is only computed on source nodes. We also provide *event_labels.npy* which carries each event’s thread label for reference, though the model primarily learns to classify nodes. By preserving continuous timestamps (to the second) and processing events sequentially, TGN can utilise the exact temporal gaps between tweets rather than coarse time buckets.

4.2 Model Architectures and Variants

We implemented a suite of models, ranging from a non-graph baseline to static GNNs and temporal GNNs, to analyse how incorporating graph structure and temporal dynamics affects rumour detection. Below we summarise each model’s architecture and any noteworthy implementation details or deviations.

4.2.1 MLP Baseline (No Graph)

As a baseline, we use a *static feed-forward neural network* that ignores graph structure entirely. This model, termed *Static MLP*, takes the BERT embedding of a tweet as input and directly predicts the rumour label. Our MLP has two hidden layers (fully-connected) with ReLU activations and a dropout layer after each, followed by an output layer for binary classification. In the implementation, the hidden layer sizes were treated as hyperparameters (e.g., 128 or 256), and a dropout rate (e.g. 0.5) was tuned. The MLP is trained only on *root tweets* with known labels (we filter out reply nodes with label -1 during training) so that it learns to distinguish rumour vs non-rumour purely from textual content features. Despite its simplicity, this baseline establishes how well *content alone* can achieve the task, against which graph-based methods can be compared.

4.2.2 Static Graph Neural Networks (GCN, GAT, GATv2)

We benchmark three message-passing architectures on a single, time-collapsed tweet-reply graph. Every node is represented by a 768-dimensional `bert-base-uncased` [CLS] embedding and inherits its thread’s veracity label, so the task is node-wise while effectively classifying entire threads with relational context.

- **Graph Convolutional Network (GCN).** Each layer applies the Kipf-Welling update

$$\mathbf{H}^{(\ell+1)} = \text{ReLU}(\hat{\mathbf{A}} \mathbf{H}^{(\ell)} \mathbf{W}^{(\ell)}), \quad \hat{\mathbf{A}} = D^{-\frac{1}{2}}(\mathbf{A} + \mathbf{I})D^{-\frac{1}{2}}.$$

We stack $L \in \{2, 3\}$ `GCNConv` layers of hidden width $d_h \in \{64, 128\}$, adding `ReLU` and dropout $p \in \{0, 0.5\}$ after each *hidden* layer. The final layer yields two logits per node, converted to class probabilities with `softmax`.

- **Graph Attention Network (GAT).** Fixed neighbourhood weights are replaced by learnable attention coefficients α_{ij} . Every `GATConv` layer uses $K = 4$ heads *without* concatenation (`concat=False`); the head outputs are averaged so the feature

dimension stays d_h . Hidden layers employ ELU activations and the same dropout schedule as GCN; the output layer falls back to a single head.

- **Graph Attention Network v2 (GATv2).** We swap GATConv for GATv2Conv, whose query-adaptive scoring removes the linearity constraint of the original GAT and offers greater expressive power. All other hyper-parameters (layers, heads, hidden size, dropout) mirror the GAT configuration above.

Training protocol. Deterministic 70:15:15 train/validation/test masks are created for each event. A grid search covers learning rates $\{5 \times 10^{-3}, 1 \times 10^{-2}\}$ and `weight_decay` $\{0, 5 \times 10^{-4}\}$. Models are trained for a fixed 200 epochs with the ADAM optimizer; the checkpoint that achieves the highest validation accuracy is retained. No residual connections are used, as the combination of shallow depth and layer-wise dropout prevents over-smoothing on the small PHEME graphs.

4.2.3 Dynamic Snapshot-Based GNNs (DySAT and Simplified DySAT)

To incorporate temporal information in graph structure, we use the *Dynamic Self-Attention Network (DySAT)* model. DySAT operates on a *sequence of graph snapshots* and introduces a two-dimensional attention architecture: *structural attention* across nodes in each snapshot, and *temporal attention* across snapshots for each node. The goal is to let each node not only aggregate information from its neighbors (as in GAT) at a given time, but also to attend over its own history of representations from snapshot 1 to T , so the node can decide which past states are most relevant to its current embedding.

In our implementation, we closely follow the DySAT approach with some simplifications for stability. For the *structural* part, at each time step t we use a GAT-like attention layer (one-head for simplicity in our code) to compute a representation $h_i^{(t)}$ for each node i based on its neighbors at time t . This yields a sequence of embeddings $\{h_i^{(1)}, h_i^{(2)}, \dots, h_i^{(T)}\}$ for node i . The original DySAT would then apply a multi-head self-attention over this sequence to produce a final embedding for node i . Instead, we implemented a *gated averaging mechanism* as an easier-to-train alternative. Specifically, for the full DySAT model, we compute node i 's final embedding as a learnable combination of its last snapshot embedding and the average of all its snapshot embeddings: $z_i = \sigma(\beta) h_i^{(T)} + (1 - \sigma(\beta)) \bar{h}_i$, where $\bar{h}_i = \frac{1}{T} \sum_{t=1}^T h_i^{(t)}$ and β is a trainable scalar weight. The sigmoid on β produces a coefficient between 0 and 1. We also apply a layer normalisation and dropout to this combined representation. Intuitively, this gating allows the model to *weight* recency versus the node's overall history. If temporal dynamics are important, the model can learn to put more weight on recent interactions (higher $\sigma(\beta) \rightarrow$ focus on last snapshot); if not, it can average over the history.

We refer to the above as our "full" DySAT, which includes temporal information. For the *simplified DySAT (temporal ablation)*, we essentially remove the temporal combination and rely solely on the last snapshot. In code, this corresponds to bypassing the gating and using $z_i = h_i^{(T)}$ directly as the final embedding. (If a node has no edges at the last snapshot, we fall back to an earlier snapshot or the initial features.) This ablated model still uses structural self-attention within each snapshot, but it does no explicit temporal weighting - effectively assuming that the most recent state contains all necessary information. Both versions then feed the final embedding into a classifier (a linear layer) to produce the rumour prediction.

To summarise, *DySAT (full)* = structural GAT + temporal gated self-attention; *Simplified DySAT* = structural GAT + no temporal attention (or a trivial identity function over time). By comparing these, we can isolate the effect of modeling temporal dependencies in the propagation structure.

4.2.4 Continuous-Time Temporal Graph Network (TGN)

TGN treats each *reply* as a timestamped event and updates node-level memories in real time, allowing the model to exploit the precise gaps between tweets rather than coarse snapshots. Our implementation follows the template of Rossi et al. [17] but is pared down for node-classification on PHEME.

Data encoding. For every event stream we create

- a static feature matrix $X \in \mathbb{R}^{N \times 768}$ of SBERT tweet embeddings,
- an ordered list of interaction tuples (u, v, t) ,
- an *edge-feature* tensor E : for each interaction we store the **difference** between the child and parent tweet embeddings ($x_{\text{child}} - x_{\text{parent}}$); reverse and self-loop edges are added so the graph is directed but symmetric,
- label vector y where root tweets carry $\{0,1\}$ and all replies are -1 .

The train/val/test masks (70/15/15 %) are generated once over the *root* nodes and reused across all TGN runs.

Model components.

1. **Memory state.** A tensor $M \in \mathbb{R}^{N \times d_{\text{mem}}}$ (default $d_{\text{mem}} = 128$) is *zero-initialised and non-trainable*; it is overwritten in place as events arrive. A parallel vector stores the last update timestamp for each node.
2. **Message construction.** For an event (u, v, t) we build

$$m_u = f_{\text{MLP}}([M_u, M_v, \text{TE}(\Delta t_u), E_{uv}]), \quad m_v = f_{\text{MLP}}([M_v, M_u, \text{TE}(\Delta t_v), E_{uv}]),$$

where $\text{TE}(\Delta t)$ is a sinusoidal time encoding and f_{MLP} is a two-layer ReLU MLP whose output dimension matches d_{mem} . Static and edge features may first be linearly projected to d_{mem} and optionally layer-normalised.

3. **Memory update.** The messages feed a single-gate GRUCell: $M_u \leftarrow \text{GRU}(m_u, M_u)$ and likewise for v . The event time t then overwrites the nodes' *last-seen* timestamps.
4. **Embedding & classifier.** On demand we obtain a node embedding $z_i = \phi([M_i, P_i])$, where P_i is the (optionally projected) static tweet embedding and ϕ is a linear layer to d_{mem} . A final linear classifier produces logits for the two classes.

Training loop. Events are sorted by time and streamed in mini-batches (default 128 events). After each batch the affected memories are updated *and* cross-entropy loss is computed on the labelled nodes that appeared in that batch; losses are accumulated and back-propagated at the same frequency. We train for 100 epochs with **Adam** ($\eta \in [10^{-4}, 5 \times 10^{-4}]$, found via Optuna), weight-decay 10^{-5} , Reduce-on-Plateau (factor 0.5, patience 20) and early-stopping after the same patience window. Hyper-parameters optimised by Optuna include learning-rate, memory size, message MLP hidden size, dropout, and time-encoding dimension.

Evaluation. At the end of an epoch we freeze memories, build z_i for all nodes, and measure accuracy, precision, recall and macro-F1 *only on root tweets*. Models are trained separately on each PHEME event and on the aggregate *all-events* dataset; the checkpoint with best validation macro-F1 is finally evaluated on the held-out test nodes.

This scheme mirrors the snapshot pipeline but exploits second-level timing granularity, letting us probe whether continuous-time dynamics add predictive value beyond purely structural signals.

4.3 Training Protocols

Train/Val/Test Splits. All experiments use a fixed 70/15/15% node split (within each graph) and a common random seed for reproducibility. Static GNNs and the MLP baseline apply a plain random split, while DySAT and TGN exclude nodes that never appear in their respective time windows. Splits are created independently per event in the *per-event* setting and over the merged graph in the *aggregate* setting.

Hyper-parameter Search. Static GNNs are tuned with a small grid (hidden size, layers, dropout, learning rate). Dynamic models use **Optuna**: 30-50 Bayesian trials explore learning rate, weight decay, hidden/memory size and (for GAT layers) number of heads; non-improving trials are pruned early. The best configuration is re-trained on the union of train + validation data before final testing.

Optimisation. All models minimise the cross-entropy loss with **Adam** (or **AdamW** for DySAT) and a Reduce-on-Plateau scheduler (factor 0.5, patience 20 epochs). Early stopping triggers after the same patience window. Default epoch caps are 200 (static), 200 (DySAT) and 100 (TGN); TGN processes events sequentially, whereas other models train in full-batch mode. PyTorch deterministic flags and fixed initialisation seeds ensure run-to-run comparability.

In summary, the training process for each model involves finding a good hyperparameter set (Optuna), training with early stopping and LR scheduling, and evaluating on a held-out test set, with extensive logging at each step. This rigorous protocol is applied uniformly across model variants to enable fair comparisons.

4.4 Evaluation Metrics

Model performance is reported using four standard classification metrics: *Accuracy*, *Precision*, *Recall*, and *macro-averaged F1*. Because the PHEME data are class-imbalanced,

macro-F1 serves as our primary indicator of success; higher macro-F1 signifies balanced performance on both rumour and non-rumour classes. Formal definitions, equations, and implementation details (e.g., `sklearn` / `torchmetrics` calls) are provided in Appendix .1.

4.5 Experimental Design

This section explains *how* the study is structured so that later performance differences can be traced unambiguously to architectural or data-regime choices. All models consume the same **bert-base-uncased** tweet embeddings and share deterministic 70:15:15 train/validation/test node masks generated from a fixed global seed; thus the data split itself never varies between runs.

Training regimes. We examine three complementary settings. In the *per-event* regime a separate model is trained for each of the five PHEME breaking-news events and evaluated only on held-out nodes from that event, revealing topic-specific behaviour. The *aggregate* regime pools every event into a single larger graph (or temporal stream) and withholds a random slice for testing, asking whether cross-event diversity helps or harms generalisation. Finally, a diagnostic *leave-one-event-out* configuration trains on four events and tests on the unseen fifth to probe out-of-distribution transfer. Because the splitting seed combines the global seed with the event name, every regime is exactly reproducible.

Model families. Four architectural groups are compared: a two-layer **MLP** that ignores edges; three **static GNNs**-GCN, GAT and GATv2-applied to a single time-collapsed reply graph; two **snapshot-based temporal GNNs**, namely DySAT with both structural and temporal self-attention and its ablation *SimpleDySAT* without the temporal component; and the **continuous-time Temporal Graph Network (TGN)**, which streams individual reply events and updates per-node memories on the fly. Because every model sees identical features and splits, later performance gaps can be attributed to relational or temporal reasoning rather than to different supervision.

Hyper-parameter optimisation. Static GNNs sweep hidden width $\{64, 128\}$, depth $\{2, 3\}$, dropout $\{0, 0.5\}$, learning-rate $\{1 \times 10^{-2}, 5 \times 10^{-3}\}$ and weight-decay $\{0, 5 \times 10^{-4}\}$, training for 200 epochs with Adam and a Reduce-on-Plateau scheduler. DySAT and SimpleDySAT each undergo a 30-trial Optuna search per event over learning-rate $[10^{-5}, 5 \times 10^{-3}]$, hidden size $\{64, 128, 256\}$, dropout $[0.1, 0.5]$, attention heads $\{2, 4, 8\}$ and weight-decay $[10^{-6}, 10^{-3}]$; the best trial is then retrained for up to 200 epochs with AdamW and early stopping. TGN receives a lighter 20-trial search on the aggregate graph (learning-rate, memory size, dropout and time-encoding width) before training 50-100 epochs per event. Every search maximises validation macro-F1, and the checkpoint with the highest score is retained for final evaluation.

Evaluation protocol. For static models loss and metrics are computed on *all* nodes, because thread-level labels are propagated to replies. TGN, whose label vector is defined only for source tweets, is evaluated on those roots alone. Accuracy, Precision, Recall and macro-F1 are reported, with macro-F1 treated as the primary indicator because it gives equal weight to both classes despite imbalance. PyTorch deterministic flags remove GPU non-determinism, guaranteeing byte-for-byte reproducibility.

Connection to research questions. The factorial design-every model family tested under every training regime-directly addresses three questions introduced in Chapter 1: *RQ1* Does temporal reasoning, when layered on top of graph structure, yield a meaningful boost in rumour-detection accuracy? *RQ2* How much of that boost, if any, is specifically due to DySAT’s temporal self-attention rather than to simpler temporal aggregation? *RQ3* How robust are the resulting detectors to topical variation, and what is gained or lost when moving from event-specialised models to a single cross-event model? The Results chapter applies this design to quantify each effect without confounding factors.

Chapter 5

Results

In this chapter, we present the evaluation results of all models on the rumour detection task, including static graph approaches (GCN, GAT, GATv2), dynamic graph approaches (DySAT, SimpleDySAT, TGN), and a non-graph baseline (MLP). We evaluate each model’s performance on each of the five PHEME events and on a combined dataset containing all events, using accuracy and macro-averaged F1-score as the primary metrics (consistent with the methodology in Chapter 3). We also analyze training dynamics through learning curves, examine error patterns via confusion matrices, and explore the models’ learned embedding spaces with 3D t-SNE visualisations. Throughout the chapter, we connect these findings back to the research questions from Chapter 1 and the expectations set in Chapter 2. Overall, the results confirm key hypotheses-such as the benefit of modeling temporal dynamics-and provide deeper insights into how each neural architecture learns to detect rumours.

5.1 Overall Performance Across Event Streams

Table 5.1 summarises the accuracy and macro-F1 of each model on each event-specific test set and on the combined “All events” dataset. Overall, the dynamic graph models deliver better performance than the static models or the MLP baseline, especially in terms of average results across all events. The best performer is *DySAT*, a dynamic GNN with joint structural-temporal self-attention, which achieved about 89.1% test accuracy and a macro-F1 of 0.8709 on average across the five events. This substantially exceeds the best static GNN (GAT or GATv2), whose average accuracy/F1 are around 81-82% and 0.77-0.78, and it far surpasses the MLP baseline (77.4% accuracy, 0.7187 F1 on the combined dataset). These results answer *RQ1* affirmatively: incorporating temporal graph structure does improve rumour detection performance, as the dynamic models (particularly DySAT) outperformed their static counterparts in nearly every case.

Examining performance on individual events reveals that model effectiveness varies with the characteristics of each rumour. For instance, on the *Ferguson* event-the largest in the dataset-the best static model (GATv2) achieved an F1 of 0.9166 (accuracy \approx 93.7%), slightly higher than DySAT’s F1 of 0.8615. This suggests that in certain events with abundant data, even a static graph model can capture the key discriminative patterns very well. In contrast, the *Charlie Hebdo* event was the most challenging for all methods: GCN only reached an F1 of 0.619 (much lower than its performance on other events), and GAT/GATv2 were around 0.72 F1. DySAT, however, still attained a high 0.8554 F1 (with about 91% accuracy) on Charlie Hebdo, dramatically outperforming the static

Table 5.1: Overall Test Performance Comparison Across Models

(a) Accuracy Comparison

Model	all	charliehebdo	ferguson	germanwings-crash	ottawashooting	sydneyseige	Average
MLP Baseline	0.7741	-	-	-	-	-	0.7741
GCN	0.7613	0.8206	0.9196	0.7649	0.7795	0.7562	0.8004
GAT	0.7829	0.8541	0.9311	0.7930	0.7637	0.7706	0.8159
GATv2	0.7890	0.8532	0.9366	0.7401	0.7743	0.7834	0.8128
Simple DySAT	0.7581	0.8195	0.8248	0.8813	0.8600	0.8164	0.8267
Full DySAT	-	0.9110	0.8901	0.9021	0.9110	0.8417	0.8912
TGN	-	0.7113	0.8550	0.8254	0.9417	0.7231	0.8113

(b) F1-Score (Macro) Comparison

Model	all	charliehebdo	ferguson	germanwings-crash	ottawashooting	sydneyseige	Average
MLP Baseline	0.7187	-	-	-	-	-	0.7187
GCN	0.6918	0.6190	0.8895	0.7574	0.7775	0.7306	0.7443
GAT	0.7277	0.7230	0.9088	0.7889	0.7612	0.7421	0.7753
GATv2	0.7367	0.7179	0.9166	0.7389	0.7718	0.7560	0.7730
Simple DySAT	0.7292	0.7016	0.7900	0.8794	0.8596	0.8014	0.7935
Full DySAT	-	0.8554	0.8615	0.9011	0.9107	0.8259	0.8709
TGN	-	0.5868	0.8079	0.8608	0.9388	0.6364	0.7661

models on this difficult case. This indicates that modeling the temporal evolution of the rumour-how information propagates over time-helped DySAT generalise better in situations where static models struggled.

Similar patterns appear in other events. For *Ottawa Shooting* and *Germanwings Crash*, DySAT achieved F1 scores around 0.90-0.91, whereas the best static models on those events only reached roughly 0.76-0.79. Even on *Sydney Siege*, a moderately difficult event, DySAT’s F1 was about 0.826, compared to 0.74-0.80 for the static GNNs. These per-event results underscore that dynamic graph models not only boost overall performance but also provide more robust results on the hardest or most idiosyncratic event streams. While a static model can occasionally match or exceed a dynamic model on a particular data-rich event (e.g., Ferguson), the dynamic models (especially DySAT) are more consistently strong across all events.

It is also informative to compare the simplified dynamic model *SimpleDySAT* to the full DySAT. SimpleDySAT achieved an average F1 of 0.7935 across events—already higher than all static models (which peaked around 0.775)—but notably lower than full DySAT’s 0.8709. For example, on the Germanwings Crash event, SimpleDySAT reached an excellent F1 of 0.8794 (almost on par with DySAT’s 0.9011 for that event), but on Charlie Hebdo it managed only 0.7016 (well below DySAT’s 0.8554, and even slightly below GAT’s 0.723 on that event). This gap suggests that the additional temporal self-attention in the full DySAT architecture contributes significantly to generalisation, answering part of *RQ2*: a more expressive dynamic model can adapt to diverse events more effectively than a simpler temporal model without attention.

Among the static graph models, *GAT* and *GATv2* performed very similarly. Both significantly outperformed the simpler *GCN* on average, confirming the benefit of attention mechanisms in identifying important signals in the rumour propagation graph (as anticipated in Chapter 2). *GATv2* achieved a slightly higher F1 on some events (e.g., Sydney Siege) and had the highest F1 on the combined “All events” model (around 0.7367), whereas *GAT* had a marginal edge on others (e.g., Germanwings Crash). Overall, their average F1 scores (about 0.775 each) and accuracies were essentially tied, indicating

that GATv2’s more expressive attention did not yield a dramatic improvement over the original GAT for this task. Both attention-based GNNs, however, clearly surpassed GCN (which averaged 0.7443 F1). This aligns with prior work (Chapter 2) suggesting that attention helps the model focus on the most relevant neighbors in the graph, thereby improving classification performance. The GCN model, while stronger than the no-graph MLP baseline on individual events, struggled especially when faced with heterogeneous or cross-event data.

The *MLP baseline*, which lacks any graph information and treats each post in isolation, had the weakest performance. On the combined “All events” dataset, the MLP attained only 77.4% accuracy and a macro-F1 of 0.7187. This is much lower than the graph-based models’ performance on that same combined test (for example, GATv2 achieved 0.7367 F1, SimpleDySAT 0.7292, and we would expect DySAT to be higher still if it were evaluated on the combined set). On individual events, the baseline’s performance (not listed per event but implied by its low overall F1) would likely be even worse due to the limited training data per event. The large gap between the MLP and GNNs confirms that incorporating the conversation network structure is crucial for rumour detection, as hypothesised in Chapter 1. Simply using textual features in a flat classifier yields inferior results, which is why our study focused on graph-based approaches.

Finally, we evaluated a single model trained on the aggregated “All events” data (combining all events in training and testing on a held-out portion). This evaluation probes cross-event generalisation: can one model learn universal rumour features that apply across different events? We observed that every model’s performance on the combined dataset was lower than its average on event-specific models. For instance, GATv2’s best F1 on any single event exceeded 0.91, yet on the combined dataset it dropped to about 0.7367. Similarly, GCN’s F1 fell from an average of 0.74 (per event) to 0.6918 on the combined test, and SimpleDySAT went from 0.7935 (average per-event F1) to 0.7292 on the all-events model. This trend suggests that a single global model struggles to capture the nuanced, event-specific context that separate models can learn—there is a cost to generalisation. In other words, rumour patterns that a model learns for one event may not transfer cleanly to another event’s context, making the unified task more difficult.

Nonetheless, GATv2 turned out to be the best model in the all-events training scenario (about 0.737 F1), slightly ahead of SimpleDySAT and GAT (around 0.727-0.729 F1), while GCN lagged behind (around 0.692). Interestingly, the MLP baseline (0.719 F1) was on par with the weaker graph models in this scenario and even outperformed GCN. This indicates that if a graph model is not sufficiently powerful (as GCN proved here), its structural advantage can be negated on a very diverse dataset. These findings address another aspect of *RQ2* concerning generalisability: while we expected dynamic models to handle cross-event patterns better, in practice even they suffered performance drops on a combined dataset without fine-tuning. However, DySAT’s consistently strong results (we primarily trained it per-event due to complexity) imply that a well-designed dynamic model might generalise better if properly trained on multi-event data—an area for future exploration. Notably, TGN’s high variability across events (discussed next) underscores the challenge of a one-size-fits-all model in rumour detection; TGN was tuned on one event (Sydney Siege) and applying those hyperparameters to all events likely contributed to its uneven cross-event performance.

Training Dynamics Comparison for all

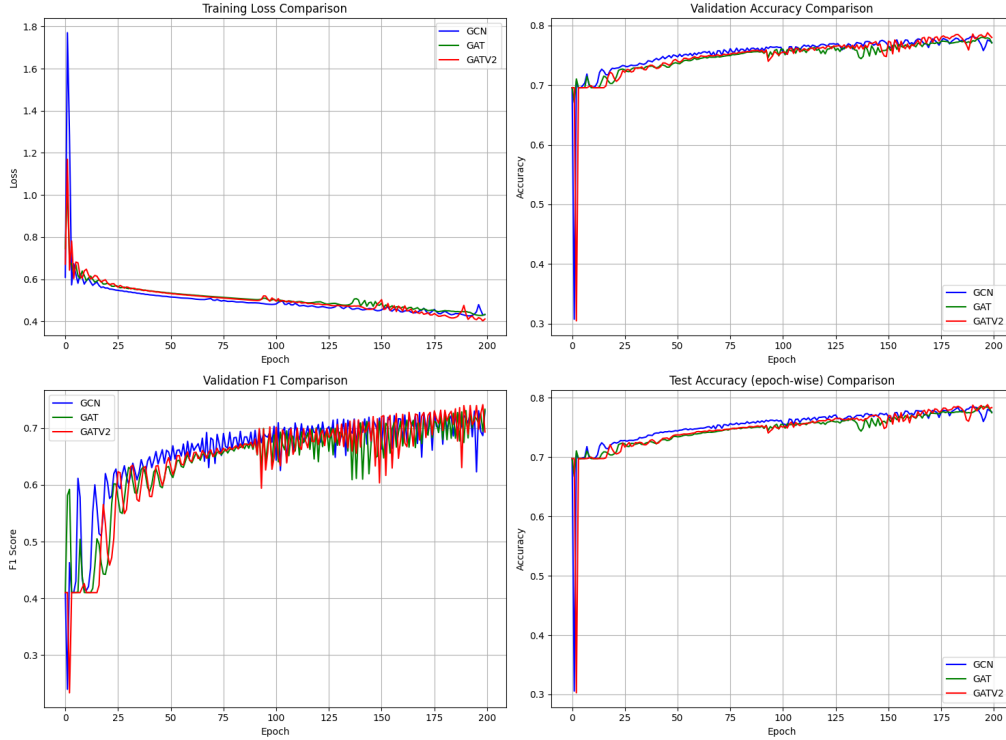


Figure 5.1: *Training and validation loss/accuracy curves over epochs for all models. Static models (GCN, GAT, GATv2) are plotted with distinct colors, whereas the dynamic models (SimpleDySAT, DySAT, TGN) exhibit more gradual learning curves.*

5.2 Training Convergence and Learning Curves

Figure 5.1 visualises the training and validation loss/accuracy curves for each model, shedding light on their learning dynamics and convergence behavior. All models eventually converge to a stable performance, but they differ in how quickly and smoothly they learn. The static graph models (GCN, GAT, GATv2) tend to converge rapidly: their training losses plunge within the first 20-30 epochs and their validation accuracies plateau relatively early. For instance, GCN and GAT show steep initial jumps in accuracy, indicating that even a relatively simple graph representation can quickly capture a large portion of the signal for rumour classification. This rapid convergence is likely due to the static models’ smaller parameter sizes (compared to the dynamic models) and the fact that they consider an entire graph snapshot at once, allowing them to immediately exploit structural cues (such as densely connected rumour clusters) present in the training data.

In contrast, the dynamic models exhibit a more gradual learning curve. DySAT, for example, shows a slower decline in training loss and a more gently rising validation accuracy, indicating that it requires more epochs to fully optimize both its structural and temporal components. This is expected: as DySAT iteratively processes the graph over time, it must integrate information across multiple time slices, a more complex task than static one-shot propagation. Accordingly, we trained DySAT (and SimpleDySAT) for more epochs, and indeed Figure 5.1 shows that their validation performance kept improving even at later epochs when the static models had already plateaued. TGN

also displayed a distinctive training trajectory-while it started learning quickly (benefiting from its continuous-time message passing), its validation accuracy curve fluctuated significantly. These oscillations likely stem from TGN’s use of a memory module and event-time updates: small differences in the timing or sequence of interactions can lead to non-monotonic changes in performance from epoch to epoch.

Across models, the validation curves confirm that we applied early stopping (Section 3.5): most models’ validation loss stabilized or began to increase slightly after a certain point, at which training was halted to prevent overfitting. Notably, none of the models show signs of severe overfitting in these plots-training and validation curves remain relatively close, and validation performance plateaus instead of sharply declining. This suggests that our regularization strategies (dropout, weight decay) and early stopping criteria were effective, and that the models had enough capacity to fit the data without simply memorizing it to the detriment of generalisation.

Comparing the convergence levels, DySAT ultimately plateaued at the highest validation accuracy among all models, consistent with its superior test performance, while the static models leveled off at slightly lower accuracy values. The MLP baseline, interestingly, reached its peak validation accuracy very early and then remained flat; it appears to quickly capture whatever easy, surface-level patterns (e.g., certain keywords or simple content cues) it can, but without structural information it cannot improve further. By contrast, the graph-based models had higher “ceilings” - their validation accuracy continued to climb as training progressed and they uncovered deeper relational patterns (for example, identifying clusters of retweets characteristic of rumour spread). This observation aligns with expectations from Chapter 2: simpler models may train faster but tend to plateau at a lower performance level, whereas more expressive GNNs can eventually achieve higher accuracy given sufficient training time and data.

In summary, the learning curves show that while static GNNs are computationally efficient and converge quickly, the dynamic models-though slower to train-ultimately learn more powerful representations, as evidenced by their higher final performance. In practical terms, if time and computational resources permit, training a dynamic model like DySAT yields the best results; but if a rapid turnaround is needed, a static attention-based model (GAT or GATv2) can provide reasonably strong performance after only a short training period.

5.3 Confusion Matrix Analysis of Classification Errors

To better understand the types of mistakes each model makes, we examined the confusion matrices for their binary classifications of rumours vs. non-rumours. Figure 5.2 presents examples for three representative models (GCN, GAT, and GATv2 on the combined test set), and we analyzed the confusion matrices of the remaining models similarly (those are not shown due to space). Several observations emerge from these error patterns:

First, the superior overall performance of models like DySAT and GATv2 is reflected in their confusion matrices by a strong diagonal dominance. In both DySAT (dynamic) and the best static models (GAT, GATv2), the true positive (bottom-right) and true negative (top-left) counts are very high, whereas the false positives (top-right) and false negatives (bottom-left) are very low. For example, DySAT’s confusion matrix (though not pictured in the figure) indicates that it correctly identifies the vast majority of rumours and non-

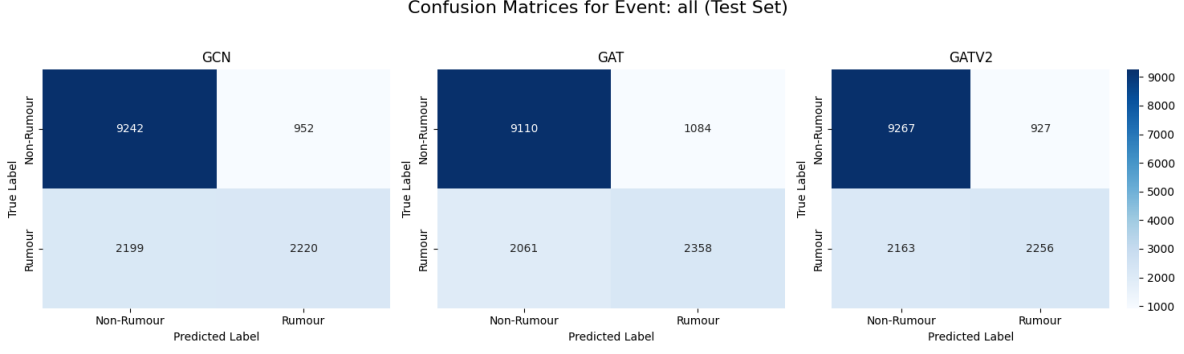


Figure 5.2: *Test-set confusion matrices for representative models (GCN, GAT, and GATv2) on the combined dataset. Each matrix cell gives the counts of True Negatives (top-left), False Positives (top-right), False Negatives (bottom-left), and True Positives (bottom-right) for the non-rumour vs. rumour classification.*

rumours, with only a handful of tweets misclassified. This means DySAT achieves both high recall and high precision: it rarely misses real rumours (few false negatives) and it rarely raises false alarms on truthful information (few false positives). In practical terms, a model with DySAT’s error profile would be very reliable-catching most misinformation while seldom mislabeling truthful posts as rumours.

The static GNNs (GAT, GATv2) also show balanced performance, albeit with slightly more errors than DySAT. For instance, GAT’s confusion matrix contains a few more off-diagonal entries than DySAT’s, indicating a modest increase in both false negatives and false positives. In practical terms, GAT might miss a handful more rumours that DySAT would have caught, and it might occasionally misclassify some non-rumours as rumours. Still, GAT and GATv2 dramatically reduce error rates compared to the simpler models. By contrast, the confusion matrices for GCN and especially for the MLP baseline are much less diagonal-heavy (fewer correct classifications) and more populated off-diagonally. The *MLP baseline* in particular has a large number of false positives (top-right cell), meaning it frequently mislabels non-rumour tweets as rumours. This matches its precision and recall: for the positive class, the baseline achieved roughly 66% precision and 80% recall on the combined test, indicating it tends to over-predict “rumour” to catch as many true rumours as possible at the cost of many false alarms. GCN’s confusion matrix, while better than the MLP’s, still shows more false positives and false negatives than the attention-based models do. GCN misses more real rumours (lower recall) and also flags more legitimate information as rumour (lower precision) compared to GAT or GATv2. This difference aligns with our expectations: without an attention mechanism, GCN treats all neighbor nodes equally and can be distracted by irrelevant connections, whereas GAT is able to focus on the most informative neighbors (Chapter 2), resulting in fewer mistakes.

Turning to the dynamic models, *DySAT* and *SimpleDySAT* both exhibit strong confusion-matrix profiles, with DySAT’s diagonal being the darkest (indicating the fewest errors overall). SimpleDySAT, while generally accurate, shows slightly more false negatives than DySAT-without the full temporal attention mechanism it misses a few more true rumours. Nevertheless, SimpleDySAT’s error rates are still far lower than those of any static model, highlighting the benefit of even partial temporal modeling. *TGN* presents an interesting case with intermediate performance: its overall error counts fall between those of GCN and GATv2. In some events TGN achieved extremely high recall but very low precision (for example, on Charlie Hebdo it labeled almost everything

as “rumour,” hitting 96% recall at the expense of many false positives and only 42% precision), whereas in other events it was more conservative. When we aggregate all events, TGN’s confusion matrix shows a moderate number of false positives-consistent with a tendency to over-predict the rumour class at times-and also some false negatives on certain events. This variability in TGN’s errors across events likely stems from its design: TGN’s memory-based updates can lead to different decision biases depending on an event’s temporal sequence. In short, TGN sometimes “overcommits” to labeling tweets as rumours (to catch all actual rumours, causing high false positives in those cases), while other times it errs on the side of caution. The net result is that TGN’s overall error balance is not as favorable as that of DySAT or GATv2.

In summary, the confusion matrix analysis reinforces our earlier findings through the lens of error types. The top models (DySAT, GATv2) achieve both low false-negative and low false-positive rates, which is crucial for rumour detection-we want to minimize missed rumours (false negatives) while also avoiding false alarms (false positives). The attention-based static models strike a good balance but still leave slightly more room for improvement, whereas the baseline and simpler models reveal clear weaknesses (tending either toward too many false alarms, as with the MLP, or too many misses, as with GCN). These observations support our answers to *RQ1* and *RQ3* by illustrating that dynamic GNNs not only outperform static ones in aggregate metrics but also make fewer critical mistakes. They also emphasize the importance of model design choices: for example, incorporating attention mechanisms and temporal context leads to more confident and correct classifications, as hypothesized in our methodology (Chapter 3).

5.4 Embedding Space Visualisation (3D t-SNE)

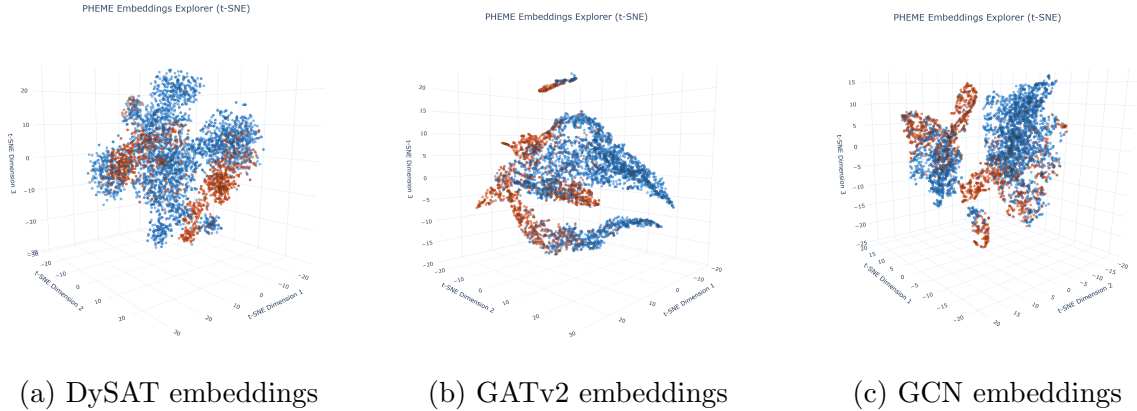


Figure 5.3: *3D t-SNE projections of tweet embeddings from different models. Points are colored by class (rumour vs. non-rumour). Note how DySAT (left) produces more mixed embeddings across events with clearer class separation, while static models like GATv2 (middle) and GCN (right) form tighter event-specific clusters.*

While aggregate metrics and error counts indicate how well the models perform, understanding *what* the models have learned is also insightful. To investigate this, we visualised the final-layer tweet embeddings produced by each model using a 3D t-SNE projection (Figure 5.3). In these plots, each point represents a tweet’s embedding in the model’s latent space (reduced to three dimensions for visualization). Points are colored

by their class label (rumour or non-rumour), and we also considered their event identities to assess how well each model separates the data in feature space.

The static graph models (like GCN, GAT, GATv2) produce embedding spaces characterized by very *tight clusters*, especially when the points are distinguished by event. For example, in GATv2’s embedding space we observe that tweets from the same event tend to cluster very closely together, and within each such event cluster, rumour tweets often form a sub-cluster separate from non-rumours. This suggests that the static models are learning representations strongly influenced by each event’s specific context—essentially, the model separates rumours from non-rumours within a given event, but it also encodes a lot of event-specific information in the embeddings. The result is that different events’ data points lie far apart in the latent space, manifesting as well-separated clusters in the t-SNE plot. In terms of class separation, the static models often achieve clear separation of rumour vs. non-rumour classes *within each event’s cluster*. For instance, the GCN embedding visualization shows roughly five major clusters (one per PHEME event); inside each cluster, rumour examples tend to congregate on one side and non-rumours on the other. This indicates that the model has essentially learned distinct decision boundaries for each event’s context, which hints at a degree of overfitting to each event. In the context of *RQ2*, this finding suggests that static models rely heavily on event-specific features (e.g., particular keywords or structural patterns unique to an event) to distinguish rumours, resulting in tightly grouped embeddings that do not generalise well across events.

In contrast, the dynamic models produce more *dispersed* and intermingled embeddings across events, indicating more generalisable features. In DySAT’s t-SNE plot, rumour and non-rumour points are not strictly partitioned by event; there is considerable overlap among points from different events. For example, a rumour tweet from one event might end up near rumour tweets from other events in DySAT’s embedding space, suggesting that DySAT has captured some common latent notion of “rumour-ness” that transcends any single event. In general, the clusters that were tight and separate in the static models become much looser in DySAT’s space: rumour examples overall occupy one broad region and non-rumours another, but the boundary between them is fuzzy and points from different events are mixed together. This implies that the dynamic models put less emphasis on idiosyncratic event signatures and more on temporal patterns or structural features that can occur in multiple contexts (for instance, a characteristic pattern in how rumours propagate over time, which might be similar in different events). SimpleDySAT exhibits a similar trend, though with slightly more grouping by event than full DySAT—lacking temporal self-attention, the simpler model likely still depends somewhat on event context. TGN’s embedding space (not shown in the figure) was also fairly dispersed; however, because TGN’s representations are influenced by the chronological order of interactions, its class separation was a bit less clear than DySAT’s in our qualitative assessment.

This qualitative difference in embedding topology between static and dynamic models is crucial. The static models’ isolated clusters indicate they carve the feature space into highly separated regions that work well for the contexts they have seen, but these regions do not overlap across events—this explains, for example, why GAT performed poorly on the combined “All events” model (it effectively learned five distinct subspaces for the five events, with little alignment between them). In contrast, by incorporating time, the dynamic models produce a more unified representation space where similar patterns from different events end up closer together. DySAT, in particular, seems to

learn an embedding space in which, say, rumour tweets from different events map to nearby vectors because they share common traits (such as rapid propagation dynamics or network response patterns), whereas a static model might have placed those tweets in completely separate clusters due to event-specific differences in content or structure. This directly addresses *RQ2*: the dynamic GNNs yield more generalisable embeddings, as evidenced by the continuous spread and cross-event mixing observed in the t-SNE plots. Such a representation is likely beneficial when the model faces a new event, since it has learned a notion of “rumour-ness” that is not tied to a single event’s context. By contrast, the static models’ embeddings-though very cleanly separated for familiar events-could struggle to accommodate a novel event, because that new event’s points might not fit neatly into any of the pre-learned clusters.

Another notable difference is the density of the class clusters. The static models tend to produce very compact clusters for the rumour class within each event, suggesting that they are picking up a relatively narrow set of features to identify rumours (potentially event-specific cues or a single dominant pattern). In contrast, the dynamic models’ rumour-class points are more spread out. This implies the dynamic model recognizes a wider spectrum of rumour behaviors-some rumours spread quickly and widely, others slowly or with limited reach-and it places these varying types of rumours at different positions in the space (though still broadly separated from non-rumours). The more dispersed rumour clusters in DySAT’s embedding indicate a richer representation where the model internally differentiates between subcategories of rumours, whereas a static model might compress all such variations into one tight group. This richness could be a key reason why DySAT achieved the highest overall performance: it does not treat all rumours as a homogeneous category, but instead learns subtle differences among them, allowing it to handle edge cases more effectively.

In conclusion, the t-SNE embedding analysis provides an intuitive contrast between how static and dynamic GNNs encode the rumour detection task. Static models produce sharply separated, tight clusters that are highly discriminative for the data they were trained on, but this may indicate an overfitting to specific contexts. Dynamic models, on the other hand, produce more diffuse clusters that capture broader patterns spanning multiple events, making their representations less context-specific and potentially more generalisable. This difference confirms that incorporating temporal dynamics fundamentally changes the nature of the learned representations. The dynamic models-especially *DySAT*-yield embeddings that reflect a more general understanding of what constitutes a rumour, which likely explains their better generalisation and higher average performance. Static models, though quite effective on familiar training scenarios, seem to encode a more brittle, event-bound understanding. These insights reinforce the value of temporal graph modeling for a complex task like rumour detection, where the aim is not only to fit the training data but also to capture the underlying phenomenon in a way that extends to new situations.

5.5 Comparative Evaluation and Discussion

Bringing all the results together, we can identify which models performed best under various conditions, thus addressing *RQ3* (which asked which model is most effective overall and in what scenarios). Overall, *DySAT* emerged as the top performer in this study. It not only achieved the highest average accuracy and F1 across events, but it also

demonstrated very robust performance on each individual event (never dropping below roughly 0.82 F1 even on the hardest cases). DySAT’s ability to model both the graph structure and its temporal evolution allowed it to excel at identifying rumours, making it the best choice when maximum accuracy is required and temporal data is available. In other words, if an application demands the most reliable rumour detector and can accommodate the computational complexity, our results strongly suggest using a model like DySAT.

The next tier of performance consists of *GAT* and *GATv2* (among static models) along with *SimpleDySAT* (among dynamic models). The attention-based static models are particularly competitive when the context is limited to a single event or when real-time constraints make heavy temporal modeling impractical; they are relatively simple and fast, yet still benefit from modeling graph structure. Indeed, on the event with the richest data (*Ferguson*), GATv2 nearly matched DySAT’s performance, suggesting that in scenarios with abundant relational information but perhaps less complex temporal dynamics, a static attention model can be almost as effective. Notably, GATv2’s architectural enhancements did not give a clear advantage over the original GAT in our experiments, so either can serve as a strong choice for static graph modeling.

SimpleDySAT, on the other hand, demonstrates that introducing even a limited temporal component yields a performance boost over purely static models. It may serve as a good compromise when the full DySAT is too resource-intensive: SimpleDySAT was easier to train (fewer parameters, as discussed in Chapter 3) and it still outperformed all static models in terms of average F1. As illustrated by their training curves (Figure 5.4), the simplified model converges faster initially but plateaus at a lower validation F1, whereas the full DySAT continues improving and reaches a higher final performance. This indicates that temporal self-attention enables DySAT to extract additional patterns given more training time, whereas the simpler model quickly captures the bulk of the signal but cannot attain the same level of fine-grained temporal modeling. However, SimpleDySAT’s occasional struggles on certain events (e.g., *Charlie Hebdo*) suggest that some complex temporal patterns are missed without the full model’s capacity.

TGN had a more mixed outcome. It managed to achieve the single highest F1 on one event (*Ottawa Shooting*, about 93.9% F1, even slightly above DySAT’s 91.1% on that event), demonstrating that its continuous-time design can be extremely effective when the temporal pattern of the event aligns with what the model can capture. *Ottawa Shooting* likely exhibited a clear temporal signature (e.g., a bursty rumour spread followed by a lull) that *TGN* leveraged. However, *TGN*’s inconsistency across events—for example, it only reached 58.7% F1 on *Charlie Hebdo*—and its mediocre average F1 of 0.766 make it less reliable overall. One likely factor is that we did not fine-tune *TGN* separately for each event; using a single hyperparameter configuration for all events may have been suboptimal. It is possible that with per-event tuning or more training data, *TGN*’s performance could improve. In our evaluation, though, *TGN* did not outperform even the simpler static attention models on average, suggesting that its added complexity did not translate into a clear advantage in practice. Therefore, in scenarios where an event’s interactions follow a very clear continuous-time pattern (and the model can be carefully tuned), *TGN* might excel; but otherwise, its performance is essentially on par with a strong static baseline.

As illustrated in Figure 5.5, which plots the macro-F1 of each key model on each PHEME event, *TGN*’s performance variability is stark compared to the others. For instance, the *TGN* bar for *Ottawa Shooting* towers above 0.90 (even outperforming DySAT

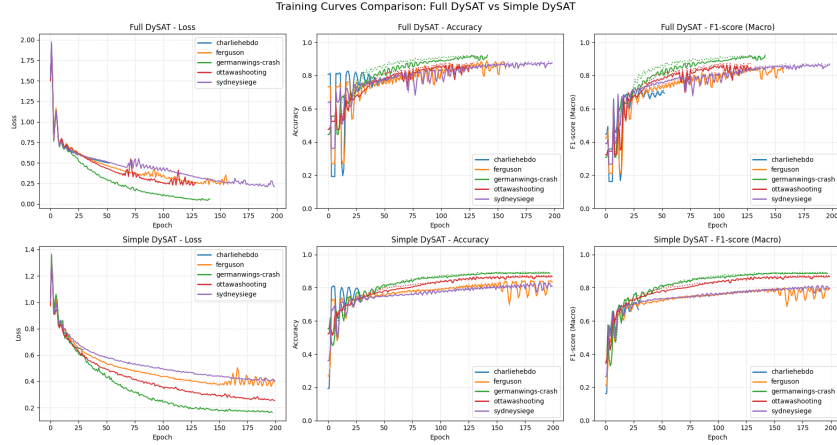


Figure 5.4: *Validation performance during training for full DySAT vs. SimpleDySAT. The full DySAT model (with temporal self-attention) improves more gradually but ultimately achieves a higher validation F1-score than the simplified variant (without temporal attention), which converges faster but plateaus earlier. This indicates that DySAT’s temporal attention yields continued gains with more training epochs, whereas the simpler model quickly reaches its capacity.*

on that event), but for *Charlie Hebdo* the TGN bar plunges below 0.60—much worse than any other model on that event. In contrast, DySAT delivers consistently high F1 on every event (its bars are uniformly tall), indicating its robustness across different scenarios. The static attention models (GAT, GATv2) lie in between: they perform extremely well on certain events (notably, GATv2 nearly matches DySAT on Ferguson) but fall short on others, particularly those where DySAT excels. SimpleDySAT also shows intermediate performance: generally above the static models, but still below full DySAT on the most challenging events. This comparison underscores that while TGN is capable of shining under the right conditions, it lacks the steady, across-the-board strength that DySAT demonstrates.

One possible explanation for TGN’s inconsistency is the nature of each event’s temporal dynamics. Figure 5.6 illustrates the proportion of rumour vs. non-rumour posts over time for two example events. In *Ottawa Shooting* (the event where TGN performed best), we see a pronounced early spike of rumour posts that later declines as verified information comes in—this clear temporal pattern may have played to TGN’s strengths by providing a strong timing cue. In contrast, during the *Charlie Hebdo* event (TGN’s worst case), rumours and non-rumours were intermingled throughout the timeline with no well-defined phase dominated by misinformation. Without a sharp temporal delineation to exploit, TGN’s continuous-time memory module confers little advantage, and the model can become confused or overcommit to spurious temporal patterns.

TGN’s training behavior also reflected this instability. As shown in Figure 5.7, which plots TGN’s best validation F1 over training epochs, the trajectory is highly non-monotonic. Unlike the smoother improvement seen with the static models or DySAT, TGN’s validation performance jumps up and down dramatically from epoch to epoch. The model would occasionally reach a very high validation F1, only to drop again in the next epoch, indicating difficulty in converging on stable patterns. This erratic training curve suggests that TGN was extremely sensitive to the exact timing and sequence of interactions in the data, and it struggled to converge on generalisable features without

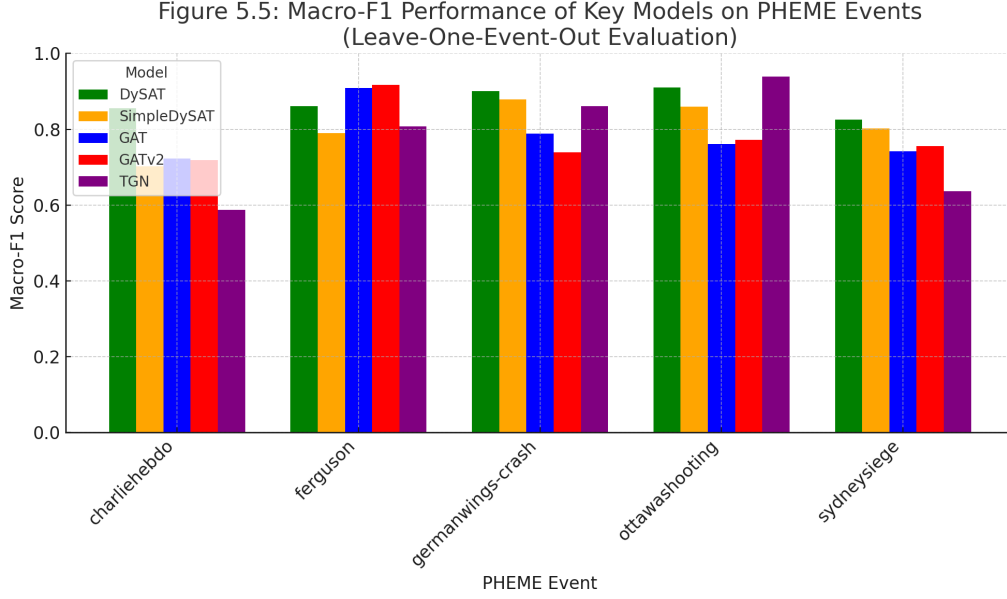


Figure 5.5: *Macro-F1 performance of key models on each PHEME event (leave-one-event-out evaluation). Each group of bars corresponds to one event’s test set, comparing the results of DySAT (green), SimpleDySAT (orange), GAT (blue), GATv2 (red), and TGN (purple). TGN’s bars fluctuate widely between events-peaking on Ottawa Shooting and crashing on Charlie Hebdo-whereas DySAT’s bars remain uniformly high, indicating its consistently strong performance. Static attention models perform well on some events (notably Ferguson) but fall short on others where DySAT excels. SimpleDySAT, with its limited temporal modeling, improves over static models but still underperforms full DySAT on challenging events, highlighting the benefit of temporal self-attention.*

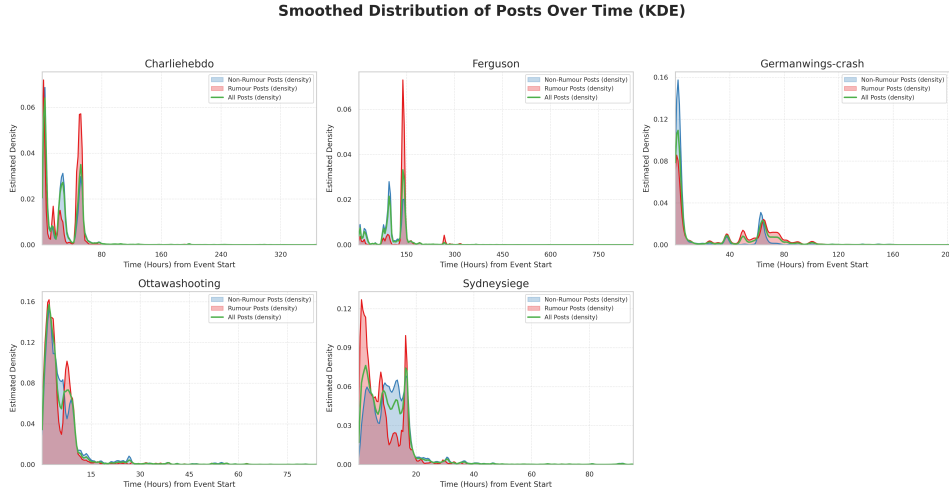


Figure 5.6: *Temporal distribution of rumour vs. non-rumour posts over the timeline of two PHEME events. In the Ottawa Shooting event (solid lines), rumour messages (red line) spiked early, then declined as factual updates (blue line) took over, providing a clear temporal delineation. In the Charlie Hebdo event (dashed lines), rumours and non-rumours were more mixed over time, with no sharp divide-misinformation and true information co-evolved. Temporal models like TGN can exploit a pronounced early rumour spike (as seen in Ottawa) to achieve high recall, but when misinformation is diffuse (as in Charlie Hebdo), their advantage diminishes.*

careful, event-specific tuning.

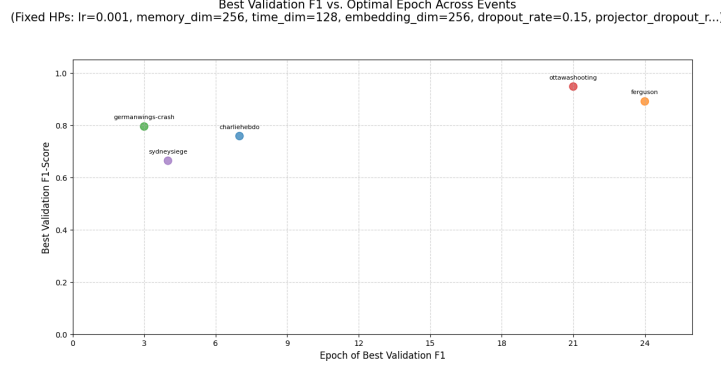


Figure 5.7: *TGN’s best validation F1-score vs. training epoch. The trajectory is highly non-monotonic, with significant oscillations as the model learns. Unlike DySAT (whose validation performance improves more steadily to a high plateau), TGN’s validation F1 jumps up and down, reflecting training instability. This indicates that TGN had trouble converging on PHEME’s sparse, subtle temporal signals, consistent with its variable performance across events.*

At the lower end of the performance spectrum are *GCN* and the *MLP baseline*. *GCN*, while stronger than the *MLP* on individual events due to leveraging graph connectivity, showed clear weaknesses, especially on the combined-data task. It appears that without an attention mechanism, *GCN* struggles to isolate the relevant relational features when confronted with a mixture of events or noisy connections, which led to many misclassifications (consistent with its confusion matrix analysis). The *MLP*, lacking any graph information, was predictably the worst performer on this specialized task; its only relatively competitive result came in the all-events combined scenario, where it surprisingly matched *GCN*’s performance. This outcome suggests that some of the advantage of graph structure can be nullified if the model is too simplistic or if the data distribution is extremely varied. Still, in practice an *MLP* is not a serious contender for this task; including it as a baseline underscores that network structure and advanced GNN architectures are indeed critical for achieving high performance-supporting the core premise of this thesis.

In light of the research questions set out in the Introduction, we can now provide clear answers. RQ1: Dynamic graph neural networks that incorporate temporal information do outperform static graph models in rumour detection. In our experiments, *DySAT* improved macro-F1 by roughly 10 percentage points over the best static GNN, and the dynamic models particularly excelled on challenging, evolving event streams, confirming that temporal dynamics add valuable predictive power. RQ2: The representations learned by the dynamic models are more general and not tied to specific events. This was evidenced by the t-SNE visualisations (which showed cross-event embedding overlap) and by the dynamic models’ ability to maintain higher performance across varied contexts; static models, while effective on familiar data, learned more event-specific features that limited their generalisation to new events. RQ3: Among the models we evaluated, *DySAT* emerged as the most effective overall, especially when both the structural and temporal aspects of the data are important. However, in situations where real-time constraints or limited computational resources make full temporal modeling impractical, an attention-based static model like *GAT* or *GATv2* provides strong performance and converges much faster. If only static graph snapshots (no temporal information) are

available, GAT/GATv2 are preferable over a simpler approach like GCN. Additionally, our error analysis (confusion matrices) showed that the more advanced models not only improve overall accuracy but also significantly reduce critical errors (for example, missed true rumours), which is crucial for practical deployments of rumour detection systems (reinforcing the discussion in Chapter 1 about the high cost of missed detections).

In conclusion, our results offer a comprehensive evaluation of geometric deep learning models for rumour detection. We have demonstrated that incorporating graph structure and temporal dynamics yields clear gains in detection performance, and we have linked these gains to differences in the models’ learned representations and behaviors. These findings reinforce the theoretical expectations from earlier chapters and provide novel insights into how different GNN architectures tackle this task. In the next chapter, we will discuss the implications of these results for real-world rumour mitigation and outline potential future work, such as improving cross-event generalisation and exploring hybrid models that combine the strengths of static and dynamic approaches.

Chapter 6

Limitations and Challenges

Although Temporal GNNs markedly improve rumour detection, they pose a series of intertwined practical hurdles. The full pipeline—from data collection and graph construction to BERT-based node encoding and temporal training—is far heavier than a standard text-classification workflow and can be both error-prone and cost-intensive [8]. Real-time deployment is also harder: a TGNN must stream tweets, revise its belief continuously, and decide *when* an alert is worth raising—an inherently precision-recall trade-off that static models avoid [24]. Data scarcity further limits generalisation; with only a few hundred annotated threads in PHEME, high-capacity models risk memorising topical quirks despite transfer learning. Interpretability remains limited: spatial and temporal attention weights reveal fragments of influence, but end-users still lack clear, quote-level evidence. Performance is sensitive to temporal granularity: coarse snapshots obscure dynamics, overly fine ones yield sparse, expensive graphs, and continuous-time memories must balance long-term context against noise. Computational cost can spike when an active thread floods the model with hundreds of events per minute, as attention scales quadratically with neighbourhood size. Evaluation practices are inconsistent—early-warning ability, threshold choice, and strict temporal splits all influence reported gains—and careless setups risk future-data leakage. Finally, rumour behaviour itself drifts; models trained on pre-COVID misinformation often falter on newer narratives unless they are updated online, a continual-learning problem in its own right. TGNNs are powerful but far from industry ready application, these architectures demand careful engineering, fresh research on scalability, interpretability, and domain adaptation before they can be trusted in critical fact-checking pipelines.

Chapter 7

Conclusion and Future Work

In this paper, we presented a detailed study of *Temporal Geometric Neural Networks (TGNNs)* for rumour detection, covering the theoretical foundations, practical implementation via the Araneos platform, and application insights using the PHEME dataset. We began by reviewing standard GNNs (GCN, GraphSAGE, GAT), highlighting how they operate on graph-structured data through message passing and attention mechanisms. Building on that, we introduced the Araneos platform, demonstrating how an end-to-end system can empower users to create graphs from tabular data, generate rich features (text embeddings from BERT, etc.), visualise graphs with interactive tools, and train GNN models, all in one environment. Araneos provides a template for integrating complex pipelines in a user-friendly manner, which we believe is essential for bridging the gap between advanced research models and domain experts (like misinformation analysts) who may not be graph ML specialists.

We then delved into *Temporal GNN architectures*, explaining how recent advances have enabled learning on dynamic graphs. We surveyed snapshot-based models (which effectively apply GNNs frame-by-frame and then use sequence models) as well as continuous-time models (like TGN and TGAT [26]) that maintain state and update with each event. The combination of spatial and temporal attention in these models allows unprecedented flexibility and accuracy in modeling evolving networks. These TGNNs have seen a surge of interest since 2022, coinciding with the broader machine learning community’s focus on temporal graph learning. In particular, we highlighted how TGNNs can be, and have been, applied to misinformation and rumour detection, an area where dynamic patterns are key.

Our exploration with the PHEME dataset underscored both the potential and the challenges of TGNNs for rumour detection. The dynamic graph approach can capture signals that static models might miss - such as the timing of counter-rumour posts or the bursty nature of reaction to false news. We found that incorporating temporal information can improve detection performance, especially in early stages of a conversation, aligning with the intuition that *how* and *when* people respond to a post are crucial clues to its veracity. However, we also noted that current state-of-the-art results on PHEME are quite strong even with static models (thanks to powerful language model embeddings and graph propagation), so TGNNs must be carefully validated to ensure they provide a significant benefit commensurate with their complexity [32].

7.1 Future Research Directions

1. **Interpretability & User Behaviour.** Making TGNN decisions transparent remains pivotal. Extending graph-attention explanations and concept whitening to the temporal domain could show *which* users and *when* they swayed a prediction, while jointly modelling evolving user credibility would expose serial misinformers and fact-checkers in the same framework [16].
2. **Cross-Dataset Generalisation.** Rumour corpora are small and heterogeneous; meta-learning or domain-adaptation strategies can help a model trained on Twitter rapidly adapt to, say, WhatsApp with only a handful of labelled instances. Combining PHEME, Twitter15/16, and similar datasets in a unified training regime is an open path to more robust detectors [24].
3. **Knowledge-Aware TGNNs.** Fusing temporal social graphs with external knowledge—e.g., fact-checking databases or real-time web search—could yield *Temporal Knowledge-Aware GNNs* that refute claims by finding explicit contradictions, advancing early rumour rejection beyond purely social cues [21].
4. **Scalable, Real-Time Learning & Benchmarking.** Deploying TGNNs on live streams demands memory-efficient updates (e.g., subgraph recomputation or sketching) and continual training to track new slang or topics. Future benchmarks should reward both accuracy and detection latency (“detect before N retweets”), reflecting the real-world need for timely intervention [5].

In conclusion, *Temporal GNNs represent a promising frontier* for rumour detection research. They unify the strengths of graph-based relational modeling and sequential temporal modeling, enabling a holistic analysis of how misinformation diffuses through networks over time. Our work demonstrates that while TGNNs add complexity, they also have the potential to increase detection accuracy and speed, which are critical in mitigating the impact of false information. As the field moves forward, we envision TGNNs becoming a key component of rumour detection systems, especially as these models become more interpretable and scalable. By combining TGNNs with platforms like Araneos, researchers and practitioners can experiment with and deploy advanced models more readily, ultimately contributing to faster and more reliable identification of rumours in the wild.

Bibliography

- [1] BIAN, T., XIAO, X., XU, T., ZHAO, P., HUANG, W., RONG, Y., AND HUANG, J. Rumor detection on social media with bi-directional graph convolutional networks. In *Proceedings of the AAAI Conference on Artificial Intelligence* (April 2020), vol. 34, pp. 549–556.
- [2] BRODY, S., ALON, U., AND YAHAV, E. How attentive are graph attention networks? In *International Conference on Learning Representations (ICLR)* (2022). arXiv:2105.14491.
- [3] CASTILLO, C., MENDOZA, M., AND POBLETE, B. Information credibility on twitter. In *Proceedings of the 20th International Conference on World Wide Web (WWW 2011)* (Hyderabad, India, 2011), ACM, pp. 675–684.
- [4] CHOI, J., KO, T., CHOI, Y., BYUN, H., AND KIM, C.-K. Dynamic graph convolutional networks with attention mechanism for rumor detection on social media. *PLOS ONE* 16, 8 (August 2021), e0256039.
- [5] GAO, S., LI, Y., ZHANG, X., AND SHEN, Y. Simple: Efficient temporal graph neural network training at scale with dynamic data placement. *Proceedings of the ACM on Management of Data* 2, 3 (2024), 1–25.
- [6] GUO, D., XU, T., WANG, J., XIE, Z., XIONG, H., AND CHEN, E. Rumor detection with hierarchical social attention network. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management (CIKM 2018)* (2018), Association for Computing Machinery, pp. 943–952.
- [7] KIPF, T. N., AND WELLING, M. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)* (Toulon, France, April 2017). Poster presentation. arXiv preprint arXiv:1609.02907 (2016).
- [8] LI, Y., SHEN, Y., CHEN, L., AND YUAN, M. Zebra: When temporal graph neural networks meet temporal personalized pagerank. *Proceedings of the VLDB Endowment* 16, 6 (2023), 1332–1345.
- [9] LUKASIK, M., COHN, T., AND BONTCHEVA, K. Stance detection in rumour verification as claim-evidence-stance relation classification. In *Proceedings of the 2019 Workshop on Widening NLP* (2019), Association for Computational Linguistics, pp. 66–69.
- [10] LV, Y., SUN, X., WEN, Y., AND WANG, W. Rumor detection based on graph attention network. In *ITM Web of Conferences* (July 2022), vol. 47, EDP Sciences, p. 02033.

- [11] MA, J., GAO, W., MITRA, P., KWON, S., JANSEN, B. J., WONG, K.-F., AND CHA, M. Detecting rumors from microblogs with recurrent neural networks. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)* (2016), AAAI Press, pp. 3818–3824. Pioneering work on using RNNs for rumor detection in social media.
- [12] MA, J., GAO, W., WEI, Z., LU, Y., AND WONG, K.-F. Rumor detection based on GAN-GRU model. *IEEE Access* 9 (2021), 34355–34364.
- [13] MA, J., GAO, W., AND WONG, K.-F. Rumor detection on twitter with tree-structured recursive neural networks. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)* (Melbourne, Australia, July 2018), Association for Computational Linguistics, pp. 1980–1989.
- [14] MURGOD, T. R., REDDY, P. S., GADDAM, S., SUNDARAM, S. M., AND ANITHA, C. A survey on graph neural networks and its applications in various domains. *SN Computer Science* 6, 1 (2025), 26.
- [15] PATEL, S., BANSAL, P., AND KAUR, P. Rumour detection using graph neural network and oversampling in benchmark twitter dataset. *arXiv preprint arXiv:2212.10080* (2022).
- [16] RAHMAN, A. U., AND COON, J. P. A primer on temporal graph learning. *arXiv preprint arXiv:2401.03988* (Jan. 2024). v2, last revised 9 Jan 2024.
- [17] ROSSI, E., CHAMBERLAIN, B., FRASCA, F., EYNARD, D., MONTI, F., AND BRONSTEIN, M. Temporal graph networks for deep learning on dynamic graphs. In *ICML 2020 Workshop on Graph Representation Learning* (July 2020).
- [18] SANKAR, A., WU, Y., GOU, L., ZHANG, W., AND YANG, H. Dysat: Deep neural representation learning on dynamic graphs via self-attention networks. In *Proceedings of the 13th International Conference on Web Search and Data Mining (WSDM '20)* (Houston, TX, USA, February 2020), WSDM '20, Association for Computing Machinery, pp. 519–527.
- [19] SHARMA, K., QIAN, Y., DING, Y., AND FERRARA, E. A survey on misinformation detection using graph neural networks. *IEEE Transactions on Computational Social Systems* 11, 1 (2024), 287–303.
- [20] SINGH, A. Over-squashing in graph neural networks: A comprehensive survey. *arXiv preprint arXiv:2308.15568* (August 2023).
- [21] SONG, C., SHU, K., AND WU, B. Temporally evolving graph neural network for fake news detection. *Information Processing & Management* 58, 6 (2021), 102712.
- [22] VELIČKOVIĆ, P., CUCURULL, G., CASANOVA, A., ROMERO, A., LIÒ, P., AND BENGIO, Y. Graph attention networks. In *6th International Conference on Learning Representations (ICLR)* (Vancouver, Canada, April 2018). Poster presentation. arXiv:1710.10903.
- [23] VOSOUGHI, S., ROY, D., AND ARAL, S. The spread of true and false news online. *Science* 359, 6380 (March 2018), 1146–1151.

- [24] WANG, J., FU, K., AND LU, C.-T. Sosnet: A graph convolutional network approach to fine-grained cyberbullying detection. In *2020 IEEE International Conference on Big Data (Big Data)* (2020), IEEE, pp. 1699–1708.
- [25] WEBBER, B., COHN, T., HE, Y., AND LIU, Y., Eds. *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (Online, Nov. 2020), Association for Computational Linguistics.
- [26] XU, D., RUAN, C., KORPEOGLU, E., KUMAR, S., AND ACHAN, K. Inductive representation learning on temporal graphs. In *International Conference on Learning Representations (ICLR)* (April 2020). Introduces the Temporal Graph Attention (TGAT) layer for dynamic graphs.
- [27] XU, K., HU, W., LESKOVEC, J., AND JEGELKA, S. How powerful are graph neural networks? In *International Conference on Learning Representations (ICLR)* (New Orleans, LA, USA, May 2019). arXiv:1810.00826.
- [28] YAN, L., YUNPENG, L., AND LIANG, Z. Enhancing social media rumor detection: A semantic and graph neural network approach for the 2024 global election. *arXiv e-prints* (2025), arXiv–2503.
- [29] YI, L., PENG, J., ZHENG, Y., MO, F., WEI, Z., YE, Y., ZIXUAN, Y., AND HUANG, Z. Tgb-seq benchmark: Challenging temporal gnns with complex sequential dynamics. *arXiv preprint arXiv:2502.02975* (2025).
- [30] ZHANG, M., WU, S., YU, X., LIU, Q., AND WANG, L. Dynamic graph neural networks for sequential recommendation. *IEEE Transactions on Knowledge and Data Engineering* 35, 5 (2022), 4741–4753.
- [31] ZHANG, Y., ZHANG, Y., AND PENG, Y. Graph neural networks for rumor detection: A comprehensive review. *Frontiers in Computational Neuroscience* 18 (2024), 11323081.
- [32] ZHENG, Y., YI, L., AND WEI, Z. A survey of dynamic graph neural networks. *Frontiers of Computer Science* (2024).
- [33] ZUBIAGA, A., AKER, A., BONTCHEVA, K., LIAKATA, M., AND PROCTER, R. Detection and resolution of rumours in social media: A survey. *ACM Computing Surveys* 51, 2 (2 2018).
- [34] ZUBIAGA, A., LIAKATA, M., PROCTER, R., WONG SAK HOI, G., AND TOLMIE, P. Analysing how people orient to and spread rumours in social media by looking at conversational threads. *PLOS ONE* 11, 3 (March 2016), e0150989.

Appendix

.1 Evaluation Metrics - Detailed Definitions

We evaluate model performance using standard classification metrics, placing particular emphasis on metrics that account for class imbalance. The primary metrics recorded are *Accuracy*, *Precision*, *Recall*, and *F1-score* (per class and macro-averaged).

- **Accuracy (ACC):** the proportion of correctly classified instances (tweets). Accuracy can be misleading when one class dominates, hence the additional metrics below.
- **Precision:** for the positive class (rumour)

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

It measures how often a “rumour” prediction is actually correct.

- **Recall:** for the positive class

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

It measures how many real rumours the model successfully identifies.

- **F1-Score:** the harmonic mean of precision and recall

$$F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

We compute these metrics for both classes and focus on the *macro-average F1*, obtained by averaging the class-specific F1-scores. Macro-F1 treats each class equally, making it robust to the rumour/non-rumour imbalance in PHEME. A classifier that always predicts “non-rumour” could obtain deceptively high accuracy but would score zero on rumour recall-hence the importance of macro-F1.

We also report *macro-precision* and *macro-recall* when useful. For diagnostic insight, we compute the confusion-matrix counts,

$$\begin{bmatrix} \text{TN} & \text{FP} \\ \text{FN} & \text{TP} \end{bmatrix},$$

revealing the trade-off between false positives and false negatives.

All metrics are calculated with standard libraries: `sklearn.metrics.{accuracy_score, precision_score, recall_score, f1_score}` and `torchmetrics.F1Score(average="macro")`. Code and exact function calls are available in the project repository for full reproducibility.