# What Signals Really Matter for Misinformation Tasks? Evaluating Fake-News Detection and Virality Prediction under Real-World Constraints

**Francesco Paolo Savatteri, Chahan Vidal-Gorène, Florian Cafiero**

Ecole nationale des chartes - PSL

65 rue de Richelieu, 72002, Paris, France

francesco.savatteri@chartes.psl.edu, chahan.vidalgorene@chartes.psl.eu florian.cafiero@chartes.psl.eu

## Abstract

We present an evaluation-driven study of two practical tasks regarding online misinformation: (i) fake-news detection and (ii) virality prediction in the context of operational settings, with the necessity for rapid reaction. Using the EVONS and FakeNewsNet datasets, we compare textual embeddings (RoBERTa; with a control using Mistral) against lightweight numeric features (timing, follower counts, verification, likes) and sequence models (GRU, gating architectures, Transformer encoders). We show that textual content alone is a strong discriminator for fake-news detection, while numeric-only pipelines remain viable when language models are unavailable or compute is constrained. Virality prediction is markedly harder than fake-news detection and is highly sensitive to label construction; in our setup, a median-based "viral" split (<50 likes) is pragmatic but underestimates real-world virality, and time-censoring for engagement features is desirable yet difficult under current API limits. Dimensionality-reduction analyses suggest non-linear structure is more informative for virality than for fake-news detection (t-SNE > PCA on numeric features). Swapping RoBERTa for Mistral embeddings yields only modest deltas, leaving conclusions unchanged. We discuss implications for evaluation design and report reproducibility constraints that realistically affect the field. We release splits and code where possible and provide guidance for metric selection.

**Keywords:** Misinformation detection, Virality prediction, Evaluation study, Textual embeddings, Numeric features, Social media analytics, EVONS, FakeNewsNet, Low-resource settings

## 1. Introduction

The rapid circulation of misleading and fabricated content on social media (Vosoughi et al., 2018; Lazer et al., 2018) has motivated two complementary lines of research: (i) detecting misinformation as early and reliably as possible, and (ii) anticipating the potential spread ("virality") of items before they accumulate engagement (Shu et al., 2017).

Recent work foregrounds multimodal detection—combining text with images and social/temporal signals—yet robust *fusion* and cross-dataset comparability remain open problems (Alam et al., 2022; Comito et al., 2023; Liu et al., 2023; Zeng et al., 2024). In parallel, graph-based methods model propagation and user–source relations with GNNs and are now consolidated by recent surveys (Phan et al., 2023). Beyond veracity, the field increasingly tackles virality prediction and user susceptibility with unified graph, multi-task, and multimodal feature-integration setups (Zhang and Gao, 2024; Jiang et al., 2024). On the data side, newer resources directly pair article content with engagement (e.g., Facebook shares), enabling virality-focused evaluation alongside veracity labels (Krstovski et al., 2022). Finally, the LLM era raises both capabilities and risks for detection pipelines, underscoring the need for stronger evaluation protocols (Chen and Shu, 2024).

While model innovation remains active, many practical deployments hinge on choices that are often treated as secondary in academic reports: how we *define* virality, what *signals* are feasible to collect (textual vs. numeric/metadata), and which *evaluation protocols* best reflect real-world constraints (access limits, time censoring, class imbalance). This paper takes an evaluation-first stance on these design decisions.

*Contribution.* We adopt an *evaluation-first* view and show across two recent datasets that **label design** and **input views** (text-only vs. text+signals) can rival architectural changes, releasing reproducible baselines that make these trade-offs explicit. While our focus is on evaluation design rather than model novelty, we show that lightweight architectures can reach or exceed the performance of recent multimodal and graph-based systems on standard benchmarks such as FAKENEWS-NET (see Section 5).

Concretely, we examine two different tasks: misinformation ("fake news") detection and early virality prediction. What signals and definitions matter most to achieve stable performance under realistic constraints? We compare textual representations (e.g., sentence-level transformer embeddings) with lightweight numeric signals (e.g., user/account statistics, early engagement counts,

and timing features). The two tasks are evaluated on datasets with distinct structures: one based on self-contained news items, the other on temporally ordered sequences of social media posts.

## 2. Tasks and Data

We address two supervised classification tasks across two datasets, yielding four experimental conditions: *(i) fake–news detection* and *(ii) virality prediction* on EVONS (Krstovski et al., 2022) and FAKENEWSNET (Shu et al., 2020).

### 2.1. Tasks

**T1: Fake–news detection.** Given an item (article or tweet series), predict a binary veracity label (*fake* vs. *true*). This task is run on both EVONS articles and FAKENEWSNET tweet series; results are reported with standard classification metrics.

**T2: Virality prediction.** Given an item, predict whether it is *viral* under a dataset–specific definition. On EVONS, virality is defined by the 95th percentile of the Facebook engagement distribution (sum of shares, likes, and comments), i.e., $v_i=1$ if $e_i \geq \tau_{95}$; otherwise $v_i=0$. On FAKENEWSNET, virality uses a median split over total likes per tweet series to ensure class balance.

### 2.2. Datasets

**EVONS.** The dataset comprises $N=92,969$ news articles with fields: title $t_i$, description $c_i$, source $s_i$, and Facebook engagement $e_i$; binary veracity labels are available. Engagement counts come from Facebook Sharing Debugger (or BuzzSumo when a source was blacklisted) and are provided as aggregated totals (no user–level data). Summary statistics show a highly skewed $e_i$ distribution (median 230; max $\sim 4.78$M). For fake–news detection on EVONS, inputs are RoBERTa embeddings of the title and description, concatenated and fed to a small MLP classifier; class imbalance is handled with a weighted loss. For virality prediction on EVONS, the positive class is the top 5% in $e_i$ (95th percentile); oversampling and undersampling were tested early but discarded in favor of weighted loss due to equal or worse performance.

**FAKENEWSNET.** We use the tweet–series setting built from the repository. Data are anonymized via IDs and require API hydration; obtaining dataset's Politifact split was possible through prior academic sharing, while GossipCop could not be recovered under current constraints. Each instance corresponds to a series of tweets referring to the same news article, i.e., a propagation path $S_i=T_{i1},\ldots,T_{in_i}$. Tweets are represented by text embeddings together with numeric features $(\Delta t_{ij}, \textit{followers}ij, \textit{following}ij, \textit{verified}ij, \textit{likes}ij)$. Temporal values are normalized as delays $\Delta t_{ij}=time_{ij} - time_i^{(0)}$ from the first tweet to characterize propagation speed. For fake–news detection on FAKENEWSNET, the dataset is imbalanced (fake ≈30For virality prediction on FAKENEWSNET, total likes per series $L_i=\sum_j \textit{likes}_{ij}$ define virality; labels use a median split.

## 3. Methods

We detail the representations, models, and training setup for both datasets. An overview of the overall experimental pipeline, covering data inputs, feature fusion, model branches, and evaluation flow, is illustrated in Figure 1.

### 3.1. Text and Numeric Representations

**Text embeddings.** All textual fields (article titles and captions on EVONS; tweet texts on FAKENEWSNET) are embedded with RoBERTa into 768-dimensional vectors. Where applicable, per-instance inputs are formed by concatenation (e.g., title+caption) before feeding the classifier. A small sensitivity check compares RoBERTa to Mistral's 1024-d embeddings on one representative model per task/dataset (Section 3.5).

**Numeric features.** On FAKENEWSNET, each tweet contributes the five numeric features $\Delta t, \text{followers}, \text{following}, \text{verified}, \text{likes}$; non-binary features are log-transformed and standardized within each training fold; the binary *verified* flag is used as is. For fusion models, numeric features are projected to 32d and concatenated with text embeddings per tweet. On EVONS, some models additionally incorporate information on article's source.

### 3.2. Models on EVONS

**Fake-news detection (text-only MLP).** A two-layer perceptron takes the concatenated (title, caption) embedding $x \in \mathbb{R}^{1536}$ and produces a single logit $z$ via GELU and dropout; optimization uses binary cross-entropy with positive-class weighting. This model is the same head used wherever we indicate "MLP."

**Virality prediction (architectures).** We evaluate four variants, differing only in the input view and a light fusion mechanism:

1. **MLP (text-only):** identical head as above, on concatenated (title, caption) embeddings.

2. **Source embedding:** a learned embedding for the publisher/source is concatenated to the text vector before the MLP head.

3. **Average engagement feature:** a scalar feature representing the source's average engagement is concatenated to the text vector before the MLP head.

4. **Gated fusion (text ↔ engagement):** text and engagement-derived vectors are fused with a gating unit à la Gated Multimodal Units (elementwise convex combination learned per dimension), then fed to the MLP head.

For the class imbalance in virality, we tested oversampling/undersampling early and retained only *weighted loss*, as data resampling performed equal or worse.

### 3.3. Models on FAKENEWSNET

**Input encoding (series).** Each instance is a series $S_i = (T_{i1}, \ldots, T_{i\ell})$ of tweets. For each $T_{ij}$, we form $x_{ij} \in \mathbb{R}^{800}$ by concatenating a 768d text embedding with a 32d projection of the 5 numeric features; the sequence $(x_{i1}, \ldots, x_{i\ell})$ is the model input.

**Sequence encoders**

1. **BiGRU:** bidirectional GRU (128 per direction) over the sequence; last hidden state → two-layer MLP classifier. Variants replace GRU with **RNN** or **LSTM** (same head).

2. **CNN:** two 1-D conv layers (kernel=3, padding=1, 128 channels) along the sequence, max-pooling over time, then a two-layer perceptron.

3. **Transformer encoder:** one bidirectional encoder layer (8 heads; FFN dim=512) with learned positional embeddings; max-pool over time, then two-layer perceptron.

### 3.4. Training, Model Selection, and Baselines

**Loss and class weighting.** All models use binary cross-entropy; positive-class weights address imbalance (small for fake-news on Evons, larger for virality on Evons).

**Cross-validation and selection.** We use stratified 10-fold cross-validation. Unless otherwise noted, the best epoch per fold is selected by **F1** on the held-out fold; metrics are computed on test splits only (Accuracy, Balanced Accuracy, Precision, Recall, F1, ROC-AUC; reported for the positive class).

**F-$\beta$ analysis (virality).** For virality on Evons, we additionally explore model selection by F-$\beta$ with recall up-weighted relative to precision; we then report the resulting trade-offs (precision drops, recall rises). Baselines are unchanged as they are not trained with F-$\beta$ selection.

**Classical baselines.** On Evons, we include dummy (stratified), logistic regression, and random forest—each fed with the same RoBERTa text embeddings as the MLP—to contextualize the effect of non-linearity.

**Optimization and hyperparameters.** We did not run automatic hyperparameter search; a small manual sweep identified a fixed configuration shared within model families. Evons: LR $1 \times 10^{-4}$, weight decay $0.01$, dropout $0.1$, 50 epochs; FakeNewsNet: LR $8 \times 10^{-5}$, weight decay $0.01$, dropout $0.1$, 100 epochs. Linear LR scheduler; no early stopping.

### 3.5. Ablations and Sensitivity Analyses

**Text vs. numeric ablation (FakeNewsNet).** Using the BiGRU architecture, we remove either the 32d numeric projection (text-only) or the 768d text embedding (numeric-only) to quantify each signal's contribution.

**Sequence length sensitivity (FakeNewsNet).** We vary $\ell \in \{2, 3, 5, 10, 20, 40\}$; virality F1 gains taper beyond $\ell = 5$, which we retain as a practical default for most runs.

**Embedding sensitivity (all settings).** We repeat one run per task/dataset with Mistral embeddings (1024d) instead of RoBERTa. Absolute F1 changes never exceed $0.02$; conclusions are unaffected.

## 4. Results

### 4.1. EVONS

**Fake–news detection.** On EVONS, all models perform strongly, with the text-based MLP approaching ceiling performance. As shown in Table 1, the MLP attains near-perfect scores across metrics (e.g., F1=0.990; ROC–AUC=0.999), while linear and tree baselines are competitive yet clearly behind, indicating a consistent benefit from light non-linearity on top of transformer embeddings.

Table 1: Performance comparison across datasets and tasks (10-fold CV).

| Dataset / Task | Model | Acc | BalAcc | F1 | Prec | Rec | ROC-AUC |
|---|---|---|---|---|---|---|---|
| **EVONS – Fake News Detection** | | | | | | | |
| MLP | | **0.991** | **0.991** | **0.990** | **0.990** | **0.990** | **0.999** |
| Logistic Regression | | 0.972 | 0.971 | 0.969 | 0.971 | 0.967 | 0.996 |
| Random Forest | | 0.932 | 0.931 | 0.926 | 0.925 | 0.927 | 0.983 |
| Dummy (stratified) | | 0.501 | 0.498 | 0.457 | 0.458 | 0.456 | 0.501 |
| **EVONS – Virality Prediction (95th percentile label)** | | | | | | | |
| MLP | | 0.885 | 0.712 | 0.312 | 0.224 | 0.519 | 0.842 |
| Source embedding | | 0.865 | **0.756** | 0.322 | 0.217 | **0.635** | **0.869** |
| Avg. engagement | | 0.869 | 0.751 | 0.322 | 0.218 | 0.621 | 0.867 |
| Gating | | 0.869 | 0.751 | **0.323** | **0.219** | 0.620 | 0.868 |
| Logistic Regression | | 0.761 | 0.783 | 0.252 | 0.150 | 0.807 | 0.866 |
| Random Forest | | 0.855 | 0.699 | 0.266 | 0.178 | 0.526 | 0.811 |
| Dummy (stratified) | | **0.905** | 0.500 | 0.049 | 0.049 | 0.049 | 0.501 |
| **EVONS – Virality Prediction ($F_\beta$, $\beta$ large)** | | | | | | | |
| MLP | | 0.688 | 0.772 | 0.219 | 0.126 | **0.865** | 0.853 |
| Source embedding | | 0.740 | 0.787 | 0.245 | 0.144 | 0.839 | 0.868 |
| Avg. engagement | | 0.728 | 0.784 | 0.239 | 0.139 | 0.847 | 0.867 |
| Gating | | **0.755** | **0.789** | **0.253** | **0.149** | 0.828 | **0.870** |
| **FakeNewsNet – Fake News Detection** | | | | | | | |
| Transformer encoder | | **0.945** | 0.927 | **0.906** | 0.933 | 0.883 | 0.965 |
| GRU | | 0.935 | 0.918 | 0.891 | 0.912 | 0.874 | 0.961 |
| RNN | | 0.941 | 0.926 | 0.901 | 0.919 | 0.886 | 0.963 |
| LSTM | | 0.936 | 0.916 | 0.891 | 0.921 | 0.866 | 0.963 |
| CNN | | 0.928 | 0.904 | 0.876 | 0.912 | 0.846 | 0.962 |
| Logistic Regression | | 0.939 | **0.929** | 0.899 | 0.896 | **0.906** | **0.971** |
| Random Forest | | 0.920 | 0.893 | 0.861 | 0.902 | 0.826 | 0.956 |
| Dummy (stratified) | | 0.578 | 0.499 | 0.300 | 0.300 | 0.300 | 0.499 |
| **FakeNewsNet – Virality Prediction (median split)** | | | | | | | |
| Transformer encoder | | **0.776** | **0.776** | **0.798** | 0.737 | **0.886** | **0.837** |
| GRU | | 0.768 | 0.768 | 0.793 | 0.720 | **0.886** | 0.823 |
| RNN | | 0.762 | 0.762 | 0.787 | 0.719 | 0.871 | 0.818 |
| LSTM | | 0.766 | 0.766 | 0.789 | 0.725 | 0.869 | 0.820 |
| CNN | | 0.770 | 0.770 | 0.790 | **0.731** | 0.864 | 0.830 |
| Random Forest | | 0.743 | 0.743 | 0.769 | 0.701 | 0.853 | 0.808 |
| Logistic Regression | | 0.691 | 0.691 | 0.695 | 0.687 | 0.705 | 0.768 |
| Dummy (stratified) | | 0.500 | 0.500 | 0.000 | 0.000 | 0.000 | 0.500 |

**Virality prediction (95th percentile).** The virality task on EVONS is substantially harder due to the extreme class imbalance and stringent definition (top 5%). Table 1 shows all neural variants clustered around F1 $\approx 0.31$–$0.32$, with a small but consistent edge for gated fusion (F1=0.323). Precision is systematically lower than recall, reflecting the rarity of truly viral items under this labeling.

**Recall-weighted selection.** When models at each epoch are selected based on the highest recall-weighted F-$\beta$ score ($\beta > 1$), they shift toward higher recall at the expense of precision (Table 1). This operating point may be preferable in screening scenarios, though overall F1 drops relative to standard selection.

## 4.2. FakeNewsNet

**Fake–news detection.** On tweet series from FakeNewsNet, sequence models dominate simple baselines; the Transformer encoder yields the best overall performance (Table 1, F1=0.906). GRU/LSTM/RNN and even the CNN are competitive, and logistic regression over text embeddings remains a strong baseline, but lags the Transformer on F1 and ROC–AUC.

**Virality prediction (median split).** With a balanced median-based label, all sequence encoders perform well and close to one another; the Transformer is slightly ahead (Table 1, F1=0.798). Traditional baselines trail the neural models, confirming the value of temporal/contextual information in this

setting.

## 4.3. Ablations and Sensitivity

**Text vs. numeric signals.** GRU ablations in Table 2 show that numeric-only inputs are viable but clearly weaker than text-only, while combining both views recovers most of the performance of the best text-based variant. This supports a practical hierarchy: text is primary, with numeric features providing complementary signal.

**Sequence length.** Increasing the maximum series length $\ell$ barely affects fake-news F1, but correlates strongly with virality F1 (Table 3; $r = -0.029$ vs. $0.965$). Gains taper beyond $\ell = 5$, which we retain as the default for efficiency.

**Embedding sensitivity.** Replacing RoBERTa with Mistral leads to small absolute F1 changes ($\leq 0.02$) and does not alter conclusions across tasks or datasets (Table 3).

## 5. Discussion

**Signals hierarchy.** Across datasets and tasks, the results point to a consistent hierarchy of signals. For *fake–news detection* on EVONS (Table 1) and on tweet series from FakeNewsNet (Table 1), text-based models are very strong: a light nonlinear head on top of transformer embeddings already approaches ceiling on EVONS and clearly outperforms classical baselines on FakeNewsNet. Numeric signals alone are usable but weaker; on FakeNewsNet virality, for example, the numeric-only GRU trails the text-only configuration by a wide margin (Table 2). Overall, *text is primary*, with numeric/meta features acting as complementary cues.

**Task difficulty and label design.** *Virality prediction* is more sensitive to label construction than model choice. Under a stringent top-5% definition on EVONS, all neural variants cluster around F1 $\approx 0.31$–$0.32$ (Table 1), whereas a balanced median split on FakeNewsNet yields substantially higher scores for all sequence encoders (Table 1). This contrast highlights that design choices (tail threshold vs. median) can dominate downstream performance, and that reporting should always pair metrics with an explicit label definition.

**Operating points and screening use-cases.** When selection optimizes a recall-weighted $F_\beta$ (Table 1), models shift to high recall with notable precision loss. This may be appropriate for screening pipelines that prioritize catching rare viral items,

but the trade-off should be stated explicitly. In settings where false positives are costly, the standard F1-selected checkpoints remain preferable.

**Architecture vs. input view.** On FakeNewsNet, sequence encoders (GRU/LSTM/RNN/CNN/Transformer) are all competitive for both tasks, with a slight edge for the Transformer (Tables 1–1). The proximity of their scores suggests that *the choice of input view (text vs. numeric, fusion) and label definition* matters more than incremental architecture tweaks—particularly under realistic constraints.

**Robustness to embedding choice.** Replacing RoBERTa with Mistral yields small absolute changes ($\leq 0.02$ F1) and does not alter conclusions (Table 3). This stability is useful for reproducibility and for deployments where embedding backends may change over time.

**Comparison to prior SOTA.** Other recent work has reported competitive results on the FAKENEWS-NET benchmark. For instance, Jiang et al. (2024) introduce a model integrating domain-specific emotional and semantic features and achieve an F1 of 0.845 on FAKENEWSNET for fake-news detection (Jiang et al., 2024). In comparison, our lightweight Transformer baseline reaches F1 = 0.906 on the same dataset, despite using frozen text embeddings and minimal parameters. This suggests that evaluation-first, resource-efficient architectures can achieve competitive performance without sacrificing transparency or reproducibility, which are often compromised in more complex multimodal systems. However, it should be noted that only original tweets were considered in this study, ignoring retweets. This makes it more difficult to compare results across experiments with different designs.

**Practical recommendations.** (1) Treat virality as a family of tasks indexed by an explicit threshold and time horizon; always report the definition alongside metrics. (2) Use text-only baselines as a strong starting point for fake–news detection, then add light numeric/meta features where available. (3) For early-screening scenarios, consider recall-weighted selection, but document precision trade-offs. (4) Favor simple, well-documented architectures; invest effort in data curation, leakage controls, and clear protocol reporting.

## Ethics and Limitations

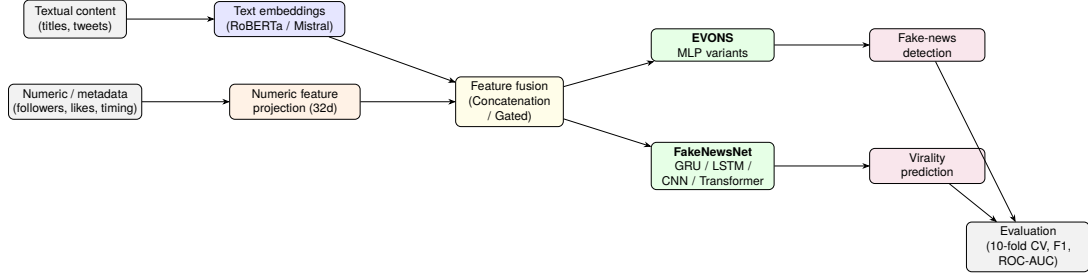**Data provenance and licensing.** We use two publicly known resources: EVONS (news articles

Figure 1: Overall experimental pipeline for fake-news detection and virality prediction across EVONS and FakeNewsNet.

Table 2: Ablation on FakeNewsNet *virality* (GRU; 10-fold CV).

| Input View | Acc | BalAcc | F1 | Prec | Rec | ROC-AUC |
|---|---|---|---|---|---|---|
| All features (text+numeric) | 0.768 | 0.768 | 0.793 | 0.720 | 0.886 | 0.823 |
| Text-only | 0.774 | 0.774 | 0.795 | 0.732 | 0.874 | 0.807 |
| Numeric-only | 0.664 | 0.664 | 0.713 | 0.628 | 0.838 | 0.724 |

with aggregate Facebook engagement counts) and FakeNewsNet (tweet series requiring API hydration). Redistribution follows the original terms of each resource. For FakeNewsNet, we only reference tweet identifiers and metadata permitted by platform policies; any reconstruction must be performed by authorized parties through the platform's API. Evons engagement values are aggregate per-URL counts and do not contain user-level information.

**Privacy and data protection.** No manual collection of personally identifiable information (PII) was performed. We do not store or release user names, free-text bios, or raw social graphs. For FakeNewsNet, we treat author/account indicators as numerical features (e.g., follower counts, verification flags) without retaining raw handles. For Evons, engagement is reported at article level (shares/likes/comments totals), which avoids end-user traceability. When sharing code, we provide feature extractors and documented loaders, not raw user data.

**Construct validity of labels.** Our virality labels reflect pragmatic choices rather than normative definitions. On Evons, "viral" denotes the top 5% of the aggregate engagement distribution; this is intentionally stringent and induces severe class imbalance. On FakeNewsNet, virality is a median split over total likes per series to obtain balanced classes; this under-represents "rare but extreme" events. Fake-news labels come from the underlying resources and inherit their annotation protocols and possible biases.

**Bias, representativeness, and harms.** Both datasets reflect specific sources and collection pipelines. Topical, temporal, and outlet biases may limit generalizability. Models trained for fake-news detection risk over-fitting to source cues or style markers rather than underlying veracity signals. Downstream misuse could include suppressing specific outlets or communities. To mitigate harms, we (i) report class-wise metrics, (ii) compare text vs. numeric signals and simple baselines, and (iii) emphasize transparent labeling choices and clearly stated limitations rather than claims of universal deployment-readiness.

Table 3: Sensitivity analyses for sequence length and embeddings.

| (a) Sequence Length Sensitivity (GRU) | | |
|---|---|---|
| | **Fake–news F1** | **Virality F1** |
| $r(\ell, \mathrm{F1})$ | $-0.029$ | **0.965** |
| Default $\ell$ | 5 (gains taper beyond 5) | |

| (b) Embedding Sensitivity (F1) | | |
|---|---|---|
| **Setting** | **RoBERTa** | **Mistral** |
| EVONS fake–news | 0.990 | **0.992** |
| FakeNewsNet fake–news | 0.891 | **0.906** |
| EVONS virality | 0.323 | **0.337** |
| FakeNewsNet virality | **0.793** | 0.773 |

**Reproducibility and sharing.** We release configuration files and feature-extraction code sufficient to reproduce the reported results, subject to each dataset's license and platform policies. For resources that cannot be redistributed (e.g., hydrated tweets), we document the procedure and required API endpoints so that authorized researchers can re-construct the inputs independently. Scripts and data are available at https://github.com/article-disinfo-lrec/article-

**Dual use and recommended safeguards.**
While early virality prediction could support beneficial moderation and curation workflows, it can also be repurposed for gaming attention or amplifying low-quality content. We recommend that any operational use involve (i) human-in-the-loop review, (ii) auditing for distributional shift, (iii) rate-limiting and red-teaming to prevent adversarial exploitation, and (iv) periodic re-estimation of thresholds and calibration.

**Environmental considerations.** All models are lightweight (frozen embeddings + small heads or shallow sequence encoders) which reduces compute and energy costs.

**Author contributions and conflicts of interest.**
Author 1, 2 and 3 contributed to study design. Author 1 was in charge of implementation, with inputs from authors 2 and 3. Authors 3 was in charge of the final redaction, with inputs from authors 1 and 2.

We declare no financial or personal relationships that could be viewed as potential conflicts of interest related to this work.

## 6. Bibliographical References

Firoj Alam, Stefano Cresci, Tanmoy Chakraborty, Fabrizio Silvestri, Dimiter Dimitrov, Giovanni Da San Martino, Shaden Shaar, Hamed Firooz, and Preslav Nakov. 2022. A survey on multimodal disinformation detection. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6625–6643, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

Canyu Chen and Kai Shu. 2024. Combating misinformation in the age of llms: Opportunities and challenges. *AI Magazine*, 45(3):354–368.

Carmela Comito, Luciano Caroprese, and Ester Zumpano. 2023. Multimodal fake news detection on social media: A survey. *Social Network Analysis and Mining*, 13(1):101.

Wen Jiang, Mingshu Zhang, Xu'an Wang, Wei Bin, Xiong Zhang, Kelan Ren, and Facheng Yan. 2024. A model for detecting fake news by integrating domain-specific emotional and semantic features. *Computers, Materials & Continua*, 80(2).

Kriste Krstovski, Angela Soomin Ryu, and Bruce Kogut. 2022. Evons: A dataset for fake and real news virality analysis and prediction. In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 3589–3596, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.

David M. J. Lazer, Matthew A. Baum, Yochai Benkler, Adam J. Berinsky, Kelly M. Greenhill, Filippo Menczer, Miriam J. Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, Michael Schudson, Steven A. Sloman, Cass R. Sunstein, Emily A. Thorson, Duncan J. Watts, and Jonathan L. Zittrain. 2018. The science of fake news. *Science*, 359(6380):1094–1096.

Peng Liu, Yang Yu, Zhen Li, and Gongshen Liu. 2023. Multi-modal fake news detection via bridging the gap between modals. *Entropy*, 25(4):614.

Huyen Trang Phan, Ngoc Thanh Nguyen, and Dosam Hwang. 2023. Fake news detection: A survey of graph neural network methods. *Applied Soft Computing*, 139:110235.

Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. 2020. Fakenewsnet: A data repository with news content, social context. *Big Data*, 8(3):171–188.

Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. 2017. Fake news detection on social media: A data mining perspective. *arXiv preprint arXiv:1708.01967*.

Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *Science*, 359(6380):1146–1151.

Fengzhu Zeng, Wenqian Li, Wei Gao, and Yan Pang. 2024. Multimodal misinformation detection by learning from synthetic data. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 10467–10484, Miami, Florida, USA. Association for Computational Linguistics.

Xuan Zhang and Wei Gao. 2024. Predicting vulnerable users with graph-based. *Information Processing & Management*, 61(1):103520.