

Temporal Betweenness Model (TBM)*

1 Integration of the new node feature: `ptd_feat`

In the multi-layer perceptron (MLP) architecture, an additional scalar feature `ptd_feat` (representing the pass-through degree [1]) is incorporated alongside the primary node feature `src_feat`. The `src_feat` is of dimension $[B, 128]$, where B denotes the batch size, while the `ptd_feat` is a scalar for each node ($[B, 1]$).

1.1 MLP Architecture

In highly imbalanced temporal graphs where only a small fraction of nodes have non-zero betweenness centrality, models must be able to sharply distinguish important from unimportant nodes. We compare two strategies for incorporating **pass-through degree (ptd)** as an additional feature: *Concatenation* and *Feature-wise Linear Modulation (FiLM)* [2].

- **Concatenation (CONCAT)**: Here, the ptd feature is concatenated with the original node feature, treating it as an additional input. This straightforward approach allows the model to learn the importance of ptd during training but may not provide fine-grained control over its influence.
- **FiLM**: FiLM applies feature-wise affine transformations to the node features, conditioned on the ptd feature. This allows for more nuanced modulation, enabling the model to adaptively scale and shift features based on ptd, potentially capturing more complex relationships.

For the FiLM modulation pathway, `ptd_feat` is passed through `ptd_scale_proj` and `ptd_bias_proj`, each projecting $[B, 1] \rightarrow [B, 128]$. These projected vectors serve as the `scale` and `bias` parameters in the FiLM operation: the modulation is computed as `modulated_feat = src_feat * scale + bias`, preserving the original shape $[B, 128]$.

In parallel, the `ptd_feat` is also processed through `ptd_proj`, a nonlinear projection with a ReLU activation and dropout, to produce another $[B, 128]$ vector. This output is then concatenated with the `modulated_feat`, resulting in a combined feature of shape $[B, 256]$.

The concatenated feature is subsequently passed through `input_proj`, which maps it back down to a `final_dim` of 128, maintaining dimensional compatibility for the downstream MLP. Finally, the output of shape $[B, 128]$ is processed by a deep MLP and projected to a single scalar per node via a final linear layer, resulting in an output of shape $[B, 1]$, which is then squeezed to $[B]$ for compatibility with standard loss functions and evaluation.

2 Experiments

Hyper-parameters. The model is conducted with a learning rate of 0.01 for 15 epochs, with the node and time embedding dimension as 128 and number of neighbor samples to 20 (same as TATKC model).

*The code of TBM is available on <https://github.com/filipposchr/TBM>

2.1 Experimental Setup

Datasets. We utilize 13 real datasets from Konect¹ and Network Repository² for evaluation. The statistics are summarized in Table 1. The training dataset consists of 48 real-world datasets from Konect, each ranging from 1000 to 12000 nodes (see Table 1).

Dataset	$ V $	$ E $	$ T $
mathoverflow	24,759	390,414	389,952
facebook-wall	35,817	198,028	194,904
topology	16,564	198,038	32,823
mlwikiquote	43,889	142,340	137,389
plwikiquote	581,646	1,472,273	1,452,278
digg-reply	30,398	87,627	83,943
SMS	44,090	544,607	467,838
retweet-pol	18,469	61,157	60,501
slashdot-reply	51,083	139,789	89,862
wang-amazon	26,911	29,062	2,175
mgwikipedia	220,064	750,811	736,680
tgwiktionary	33,968	81,516	67,065
ltwiktionary	689,678	1,693,277	1,633,334

Table 1: Statistics of the testing datasets.

Evaluation Metrics. We evaluate the model’s effectiveness in terms of top-N% accuracy, Kendall’s tau correlation and Jaccard index.

1. **Top-K Accuracy** is defined as:

$$\text{Top-}N\% = \frac{|\text{Top-}N\% \text{ predicted nodes} \cap \text{Top-}N\% \text{ true nodes}|}{\lceil |V| \times N\% \rceil}$$

where $|V|$ is the number of nodes and $\lceil \cdot \rceil$ denotes the ceiling function. This metric measures the fraction of correctly predicted top-N% nodes compared to the ground truth, and is commonly used to evaluate effectiveness in identifying highly central nodes.

Evaluation Code Behavior:

The function `compute_topk_metrics_all` implements this definition precisely:

- `N = len(preds)` corresponds to $|V|$.
- `topk = max(1, math.ceil(N * (k / 100)))` implements $\lceil |V| \cdot N\% \rceil$.
- `torch.topk(preds, topk)` selects the top-K predicted nodes.
- `torch.topk(labels, topk)` selects the top-K true (ground-truth) nodes.
- The intersection of these sets is computed, and the final metric is:

$$\frac{|\text{predicted top-}K \cap \text{true top-}K|}{K}$$

2. **Weighted Kendall Tau (WKT)** is a rank-based correlation metric that measures the ordinal association between two sequences. In this evaluation, it is used to assess how well the predicted ranking of nodes aligns with the ground truth ranking based on temporal betweenness.

¹Konect is available at <http://konect.cc/>

²Network Repository is available at <https://networkrepository.com/networks.php>

- (a) **Tiebreaking Strategy:** The ranking procedure uses the function `rank_with_id_tiebreak(values)`, with the "ordinal" method, which assigns unique integer ranks based on the order of appearance (row index) when ties occur. This avoids ambiguity from tied ranks, ensuring that every element receives a unique rank, which is required for computing Kendall tau consistently.

- (b) **Computation of WKT** (function `compute_kendall_tau`):

- i. **Standard WKT over all nodes** is computed as:

```
label_ranked_full = rank_with_id_tiebreak(labels)
pred_ranked_full = rank_with_id_tiebreak(preds)
kt_full, _ = weightedtau(pred_ranked_full, label_ranked_full)
```

This version compares the full ranked list of predictions and ground truth across all nodes, regardless of whether the ground truth is zero or not.

- ii. **Filtered WKT on non-zero ground truth values only** is computed as:

```
nonzero_mask = labels > 0
label_filtered = labels[nonzero_mask]
pred_filtered = preds[nonzero_mask]
label_ranked_filt = rank_with_id_tiebreak(label_filtered)
pred_ranked_filt = rank_with_id_tiebreak(pred_filtered)
kt_filt, _ = weightedtau(pred_ranked_filt, label_ranked_filt)
```

This filtered variant restricts the correlation calculation to the subset of nodes with non-zero ground truth centrality. This is important in sparse graphs where most nodes have zero temporal betweenness, which would otherwise dominate the ranking and obscure differences among important nodes.

- (c) The WKT values reported in Table 3 and Table 5 correspond to the *filtered* WKT computed only on nodes with non-zero ground truth (2.a.ii). The standard WKT computed over all nodes is consistently higher.

- (d) **Interpretation:**

- A higher `kt_full` indicates better alignment in global ranking across all nodes.
- A higher `kt_filt` reflects better ranking precision among the truly important (non-zero) nodes.
- In datasets where many nodes have a ground truth value of zero, WKT computed over all nodes can be lower if the model assigns different (noisy) prediction scores to those zero-value nodes (e.g., x , y , z). This introduces noise into the global ranking and reduces Kendall tau. In contrast, WKT computed only over the non-zero ground truth nodes ignores these irrelevant variations and can yield a higher score. This filtered WKT better reflects the model's ranking quality on the meaningful part of the distribution. When the model assigns the same predictions to all zero-labeled nodes, the all-node and non-zero-only WKT scores tend to be closer.

3. **Jaccard Index** is defined as:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

where A and B are two sets. In the context of this model evaluation, A represents the set of predicted important nodes and B the ground truth important nodes.

- (a) **Computation of Jaccard**

The following lines compute the Jaccard index in the `compute_topk_metrics_all` function :

```
pred_full_set = set(torch.topk(preds, len(true_full_set)).indices.tolist())
full_jacc = len(pred_full_set & true_full_set) / len(pred_full_set | true_full_set)
```

This corresponds to the following steps:

- `true_full_set` contains all node indices for which the ground truth value is greater than zero (i.e., nodes deemed important by the centrality ground truth).
- `pred_full_set` is constructed by selecting the same number of top predicted nodes as there are in the true set, ensuring fair comparison.
- The numerator computes the size of the intersection between the predicted and true sets.
- The denominator computes the size of the union of the two sets.

(b) **Interpretation:**

The Jaccard Index here captures the global agreement between the predicted and true sets of important nodes—regardless of ranking—by comparing the overlap of non-zero centrality predictions with ground truth. A higher Jaccard value indicates better alignment between model predictions and actual influential nodes.

2.2 Experimental Result

We evaluate two models: **TSBM** (Temporal Shortest Betweenness Model), which focuses on temporal shortest-path betweenness (sh), and **TSFMBM** (Temporal Shortest-Foremost Betweenness Model), designed for temporal shortest-foremost betweenness (sfm). Specifically, we evaluate the performance of four model variants: **TSBM_C**, **TSBM_{CF}** (see Table 2 and Table 3), and their shortest-foremost equivalents, **TSFMBM_C** and **TSFMBM_{CF}** (see Table 4 and Table 5). Evaluation is conducted using Top-N accuracy, Jaccard index, and Weighted Kendall’s Tau (WKT) metrics across a diverse set of real-world temporal graphs. The **C** variants rely solely on feature concatenation, while the **CF** variants extend this setup with FiLM-based conditioning over the concatenated features.

For comparison, we include WKT scores for the MANTRA model, a recent approximation-based approach for temporal betweenness estimation. The models are trained using the standard `loss_cal` function.

Some WKT scores for the MANTRA [3] model are omitted due to the extremely low connectivity of certain temporal graphs (e.g., 0.0000096 in `tgwiktionary`), where most nodes have zero betweenness (only 235 of 33,968 nodes with non-zero temporal betweenness). In such sparse temporal graphs, reliable approximation would require a very small ε (e.g., $\varepsilon < 0.000009$), making MANTRA impractical or inaccurate. Therefore, WKT is excluded where approximation quality cannot be ensured.

Dataset	Top-1%			Top-5%			Top-10%			Top-20%			Jaccard		
	TSBM _C	TSBM _{CF}	MANTRA	TSBM _C	TSBM _{CF}	MANTRA	TSBM _C	TSBM _{CF}	MANTRA	TSBM _C	TSBM _{CF}	MANTRA	TSBM _C	TSBM _{CF}	MANTRA
mathoverflow	0.86	0.5	0.85	0.82	0.82	0.7	0.76	0.76	0.57	0.8	0.8	0.39	0.95	0.94	0.54
facebook	0.56	0.56	0.85	0.65	0.65	0.83	0.7	0.7	0.79	0.77	0.78	0.67	0.9	0.9	0.69
topology	0.7	0.66	0.68	0.66	0.66	0.65	0.68	0.67	0.55	0.7	0.7	0.39	0.78	0.79	0.44
mlwikiquote	0.76	0.67	0.9	0.84	0.84	0.56	0.8	0.82	0.34	0.63	0.85	0.49	0.77	0.83	0.49
plwikiquote	0.68	0.68	0.49	0.83	0.83	0.04	0.82	0.94	0.08	0.49	0.97	0.48	0.93	0.93	0.47
digg_reply	0.65	0.65	0.88	0.7	0.7	0.79	0.75	0.75	0.6	0.93	0.92	0.35	0.99	0.99	0.67
SMS	0.64	0.64	0.56	0.63	0.64	0.22	0.67	0.67	0.15	0.85	0.99	0.15	0.83	0.86	0.42
rt-pol	0.75	0.75		0.77	0.77		0.91	0.91		0.65	0.87		0.96	0.96	
slashdot	0.81	0.81	0.75	0.82	0.82	0.67	0.84	0.84	0.41	0.85	0.91	0.39	0.9	0.91	0.59
wamazon	0.95	0.95	0.61	0.72	1	0.74	0.87	1	0.87	0.92	1	0.93	1	1	0.64
mgwikipedia	0.84	1		0.18	0.99		0.11	0.99		0.11	0.99		0.95	0.95	
tgwiktionary	0.84	1		0.19	0.99		0.19	0.99		0.23	0.99		0.99	0.99	
ltwiktionary	0.31	1		0.11	0.99		0.13	0.99		0.2	0.99		0.99	0.98	

Table 2: Top-N accuracy and Jaccard similarity for **temporal shortest betweenness (sh)** on a diverse set of real-world temporal graphs. **TSBM_C** refers to the model variant using feature concatenation only, while **TSBM_{CF}** applies FiLM-based conditioning over concatenated features. Bold values highlight the best performance across methods.

Discussion. In highly imbalanced settings, it is possible for the Top- k accuracy to decrease as k increases, contrary to intuition. This counterintuitive behavior typically occurs when the ground truth contains a large number of zero-valued targets, and the model assigns non-zero (noisy) scores to these unimportant nodes. As k increases, the prediction set may include more of these noisy false positives, thereby reducing precision. While the Top-1% predictions often remain focused on genuinely high-centrality nodes, the wider Top-5% or Top-10% sets can become diluted by misranked low-importance nodes.

Dataset	WKT (all nodes)			WKT (non zero)		
	TSBM _C	TSBM _{CF}	MANTRA	TSBM _C	TSBM _{CF}	MANTRA
mathoverflow	0.92	0.94	0.75	0.42	0.87	0.86
facebook	0.91	0.92	0.83	0.86	0.86	0.9
topology	0.91	0.91	0.73	0.84	0.83	0.82
mlwikiquote	0.8	0.95	0.76	0.8	0.79	0.86
plwikiquote	0.86	0.98	0.58	0.83	0.83	0.82
digg reply	0.91	0.97	0.84	0.85	0.85	0.9
SMS	0.89	0.93	0.63	0.84	0.85	0.74
rt-pol	0.9	0.98		0.31	0.87	
slashdot	0.94	0.98	0.83	0.9	0.9	0.86
wamazon	0.94	0.99	0.87	0.86	0.86	0.8
mgwikipedia	0.83	0.99		0.83	0.81	
tgwiktionary	0.59	0.99		0.9	0.83	
ltwiktionary	0.7	0.99		0.92	0.85	

Table 3: Standard (all nodes) and filtered WKT scores for **temporal shortest betweenness (sh)**, computed using the ordinal tie-breaking strategy. The filtered scores consider only nodes with non-zero ground truth values. Results are reported for the TSBM_{CF} variant.

Dataset	Top-1%			Top-5%			Top-10%			Top-20%			Jaccard		
	TSFMBM _C	TSFMBM _{CF}	MANTRA	TSFMBM _C	TSFMBM _{CF}	MANTRA	TSFMBM _C	TSFMBM _{CF}	MANTRA	TSFMBM _C	TSFMBM _{CF}	MANTRA	TSFMBM _C	TSFMBM _{CF}	MANTRA
mathoverflow	0.85	0.58		0.84	0.84		0.79	0.79		0.8	0.8		0.95	0.95	
facebook	0.55	0.55	0.83	0.65	0.65	0.81	0.7	0.7	0.76	0.77	0.77	0.61	0.9	0.9	0.65
topology	0.72	0.68	0.75	0.68	0.68	0.66	0.7	0.7	0.54	0.71	0.71	0.39	0.77	0.8	0.46
mlwikiquote	0.76	0.73	0.90	0.8	0.81	0.59	0.76	0.79	0.41	0.52	0.52	0.63	0.63	0.67	0.46
plwikiquote	0.66	0.66	0.5	0.74	0.74	0.04	0.64	0.64	0.24	0.4	0.37	0.53	0.72	0.73	0.56
digg_reply	0.64	0.64	0.85	0.7	0.7	0.68	0.74	0.74	0.43	0.93	0.93	0.27	0.98	0.98	0.58
SMS	0.65	0.66		0.64	0.64		0.67	0.67		0.85	0.85		0.84	0.86	
rt-pol	0.75	0.75	0.82	0.77	0.77	0.5	0.91	0.91	0.27	0.63	0.62	0.51	0.92	0.92	0.62
slashdot	0.79	0.79	0.74	0.81	0.81	0.61	0.84	0.84	0.37	0.82	0.82	0.37	0.88	0.89	0.56
wamazon	0.96	0.96		0.58	0.55		0.71	0.69		0.84	0.81		0.99	0.99	
mgwikipedia	0.84	0.84		0.17	0.15		0.1	0.12		0.11	0.13		0.76	0.76	
tgwiktionary	0.84	0.84		0.19	0.2		0.18	0.17		0.22	0.25		0.91	0.91	
ltwiktionary	0.11	0.13		0.07	0.03		0.15	0.3		0.24	0.43		0.92	0.92	

Table 4: Top-N accuracy and Jaccard similarity for **temporal shortest-foremost betweenness (sfm)** on a diverse set of real-world temporal graphs. TSFMBM_C refers to the model variant using feature concatenation only, while TSFMBM_{CF} applies FiLM-based conditioning over concatenated features. Bold values highlight the best performance across methods.

Dataset	WKT (all nodes)			WKT (non zero)		
	TSFMBM _C	TSFMBM _{CF}	MANTRA	TSFMBM _C	TSFMBM _{CF}	MANTRA
mathoverflow	0.92	0.9		0.9	0.88	
facebook	0.91	0.91	0.82	0.86	0.86	0.89
topology	0.9	0.9	0.73	0.85	0.84	0.83
mlwikiquote	0.73	0.72	0.80	0.8	0.8	0.83
plwikiquote	0.82	0.72	0.63	0.79	0.79	0.79
digg reply	0.92	0.91	0.8	0.85	0.85	0.86
SMS	0.89	0.89		0.85	0.85	
rt-pol	0.89	0.86	0.81	0.86	0.85	0.86
slashdot	0.95	0.93	0.82	0.9	0.9	0.84
wamazon	0.98	0.9		0.86	0.86	
mgwikipedia	0.8	0.7		0.81	0.8	
tgwiktionary	0.63	0.63		0.87	0.87	
ltwiktionary	0.62	0.66		0.78	0.77	

Table 5: Standard (all nodes) and filtered WKT scores for **temporal shortest foremost-betweenness (sfm)**, computed using the ordinal tie-breaking strategy. The filtered scores consider only nodes with non-zero ground truth values. Results are reported for the TSFMBM_{CF} variant.

This is observed when evaluating the model TSBM_{CF} on the `mathoverflow` dataset, where the Top-1% accuracy reaches 0.50, but the Top-5% drops to 0.82, with Top-10% and Top-20% further declining to 0.76 and 0.80, respectively. The large number of ground-truth-zero nodes in `mathoverflow`, combined with prediction variance among them, leads to reduced performance at broader thresholds despite a strong Top-1% accuracy.

3 Loss function

The default model, **TATKC**, utilizes a pairwise ranking loss function, implemented in `loss_cal` to train the model for approximating TKC ranking scores. This function relies on margin ranking loss, which randomly samples node pairs and encourages the model to preserve correct ordering by scoring higher ground-truth nodes above lower ones.

4 Further on FiLM Architecture

4.1 How FiLM Operates Over Features in the Model

Feature-wise Linear Modulation (FiLM) is a conditioning mechanism that adjusts neural network activations using external signals. In our model, FiLM modulates the primary node representation (`src_feat`) using a graph-theoretic input: the pass-through degree (`ptd_feat`), which measures how frequently a node mediates time-respecting paths.

Instead of processing `src_feat` and `ptd_feat` jointly, FiLM treats `ptd_feat` as a controller. It produces per-feature scaling and shifting parameters applied to `src_feat` via:

$$\text{modulated_feat} = \text{src_feat} \times \gamma(\text{ptd_feat}) + \beta(\text{ptd_feat}),$$

where γ and β are learnable projections of the scalar `ptd_feat` to match the dimension of `src_feat`.

This modulation is applied not just at the input, but throughout deeper MLP layers, allowing the model to adapt the importance of features at multiple transformation stages. Notably, `ptd_feat` itself is not transformed or modulated, preserving its role as a stable, interpretable conditioning signal.

Overall, FiLM enables dynamic reweighting of node features while maintaining the interpretability of the external signal. This leads to more expressive and targeted transformations of `src_feat`, especially valuable in capturing node importance in temporal graphs.

Example: tgwiktionary Dataset. This dataset contains 33,968 nodes, of which only 235 have non-zero temporal betweenness scores. We expect the model to:

- Predict high values for these top 235 nodes;
- Assign uniform low values to all others.

Under FiLM, the model outputs distinct scores for top nodes (e.g., 0.746, 0.551, 0.337), while collapsing all other predictions to a single low value (e.g., -1.41), clearly separating predictions from noise. By contrast, the concatenation-only model produces scattered negative values for irrelevant nodes (e.g., -0.25, -0.51, -0.29), reducing ranking sharpness. Evaluation results confirm this: FiLM achieves $\text{Top-1\%} = 1.0$, while concatenation yields $\text{Top-1\%} = 0.85$. (see Table 2)

4.2 When the CF variant Outperforms or Underperforms C variant

The combination of FiLM and concatenation performs variably depending on dataset characteristics:

- **In dense or structured graphs** (e.g., SMS, slashdot), the CF variant tends to outperform C variant because:
 - FiLM highlights attention-relevant dimensions guided by PTD;

- Concatenation preserves an explicit, interpretable PTD input;
- The two mechanisms complement each other: one reshapes internal representations, the other offers direct supervision.
- **In sparse or skewed graphs** (e.g., mathoverflow), **CF** variant may underperform:
 - FiLM modulation can distort already informative **src_feat** signals;
 - This can hurt top-1% accuracy, where relevant nodes are rare and highly concentrated.

5 Conclusion

FiLM enables sharper separation between important and irrelevant nodes by modulating the feature space with learned scaling factors. This results in:

- Assign diverse and well-separated predictions to top-ranked, high-centrality nodes;
- Produce consistent low outputs for irrelevant or zero-ground-truth nodes;
- Improve ranking quality, both globally (Kendall Tau) and at the top-k level;
- Enhance the interpretability of node importance by clearly separating predictions from noise.

In highly imbalanced temporal graphs—where few nodes exhibit non-zero betweenness—the model must strongly differentiate between the predicted betweenness values and noise. We compared two strategies for incorporating pass-through degree: concatenation and FiLM. In several datasets, FiLM demonstrated better performance by acting as a feature-wise modulation mechanism that conditions node embeddings based on structural importance. This yields more stable predictions for unimportant nodes and better distinction among the top-ranked nodes.

References

- [1] Ruben Becker, Pierluigi Crescenzi, Antonio Cruciani, and Bojana Kodric. Proxying Betweenness Centrality Rankings in Temporal Networks. In Loukas Georgiadis, editor, *21st International Symposium on Experimental Algorithms (SEA 2023)*, volume 265 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 6:1–6:22, Dagstuhl, Germany, 2023. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. doi:10.4230/LIPIcs.SEA.2023.6.
- [2] Marc Brockschmidt. GNN-FiLM: Graph neural networks with feature-wise linear modulation. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 1144–1152. PMLR, 13–18 Jul 2020. URL: <https://proceedings.mlr.press/v119/brockschmidt20a.html>.
- [3] Antonio Cruciani. Mantra: Temporal betweenness centrality approximation through sampling. In Albert Bifet, Jesse Davis, Tomas Krilavičius, Meelis Kull, Eirini Ntoutsi, and Indrė Žliobaitė, editors, *Machine Learning and Knowledge Discovery in Databases. Research Track*, pages 125–143, Cham, 2024. Springer Nature Switzerland.