# CHARACTERIZING MASSIVE ACTIVATIONS OF ATTENTION MECHANISM IN GRAPH NEURAL NETWORKS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

Graph Neural Networks (GNNs) have become increasingly popular for effectively modeling data with graph structures. Recently, attention mechanisms have been integrated into GNNs to improve their ability to capture complex patterns. This paper presents the first comprehensive study revealing a critical, unexplored consequence of this integration: the emergence of Massive Activations (MAs) within attention layers. We introduce a novel method for detecting and analyzing MAs, focusing on edge features in different graph transformer architectures. Our study assesses various GNN models using benchmark datasets, including ZINC, TOX21, and PROTEINS. Key contributions include (1) establishing the direct link between attention mechanisms and MAs generation in GNNs, (2) developing a robust definition and detection method for MAs based on activation ratio distributions, (3) introducing the Explicit Bias Term (EBT) as a potential countermeasure and exploring it as an adversarial framework to assess models robustness based on the presence or absence of MAs. Our findings highlight the prevalence and impact of attention-induced MAs across different architectures, such as Graph-Transformer, GraphiT, and SAN. The study reveals the complex interplay between attention mechanisms, model architecture, dataset characteristics, and MAs emergence, providing crucial insights for developing more robust and reliable graph models.

033

006

008 009 010

011

013

014

015

016

017

018

019

021

024

025

026

027

028

#### 1 INTRODUCTION

Graph Neural Networks (GNNs) have emerged as a powerful tool for learning representations of graph-structured data, demonstrating remarkable success across various applications such as social network analysis (Min et al., 2021), recommendation systems (Gao et al., 2022) and molecular biology (Zhang et al., 2021). Central to the recent advancements in GNNs is the integration of attention mechanisms, which enable the models to focus on the most relevant parts of the input graph, thereby enhancing their ability to capture intricate patterns and dependencies.

Despite the substantial progress, the phenomenon of Massive Activations (MAs) within attention
 layers has not been thoroughly explored in the context of GNNs. MAs, characterized by exceedingly
 large activation values, can significantly impact the stability and performance of neural networks.
 In particular, understanding and mitigating MAs in GNNs is crucial for ensuring robust and reliable
 model behavior, especially when dealing with complex and large-scale graphs.

In this paper, we aim to bridge this gap by systematically investigating the occurrence and implications of MAs in attention-based GNNs. We focus on edge features in graph transformers, a state-of the-art GNN architecture, and analyze how these features contribute to the emergence of MAs. Our
 study reveals that certain graph structures on edge configurations are more prone to inducing MAs,
 which in turn affects the overall performance and interpretability of the models.

To address these challenges, we propose a novel methodology for detecting and analyzing MAs in
GNNs. Our approach involves a comprehensive evaluation of various GNN architectures, including
GraphTransformer (Dwivedi & Bresson, 2021), GraphiT (Mialon et al., 2021), and SAN (Kreuzer
et al., 2021), across multiple benchmark datasets, like ZINC (Irwin et al., 2012), TOX21 (Mayr
et al., 2016; Huang et al., 2016) and OGBN-PROTEINS (Hu et al., 2020), which differs from their

054 downstream tasks like graph regression, multi-label graph classification, and multi-label node clas-055 sification. We introduce specific criteria for identifying MAs and conduct extensive ablation studies 056 to elucidate the role of edge features in this context.

057 This study represents the first comprehensive investigation of MAs in GNNs, laying the groundwork 058 for future research. Our findings suggest that the scope of MAs analysis can be expanded to include a wider range of architectures and the evaluation of state-of-the-art attack methods, ultimately en-060 hancing our understanding of MAs' influence on GNN performance and robustness. This is crucial 061 for developing more robust and reliable graph transformer models, especially given the increasing 062 popularity and widespread adoption of transformers in various applications today.

Our contributions are threefold:

- We provide the first systematic study on MAs in attention-based GNNs, highlighting their prevalence and impact on model performance.
- We propose a robust detection methodology for MAs, accompanied by detailed experimental protocols and ablation studies.
- We introduce the Explicit Bias Term (EBT) as a potential countermeasure for MAs, and we exploit it in an adversarial framework, called Explicit Bias Attack, to demonstrate the effectiveness of the MAs in compromising GNNs robustness.

Through this work, we aim to shed light on a critical yet understudied aspect of attention-based GNNs, offering valuable insights for the development of more resilient and interpretable graphbased models. 076

077 078

063

064 065

066

067

068

069 070

071

073

074

075

#### **RELATED WORKS** 2

079 080

GNNs have become effective instruments for studying and extracting insights from graph-structured 081 data, with usages spanning fields like fraud detection (Motie & Raahemi, 2023), traffic prediction 082 (Wang et al., 2022) and recommendation systems (Wu et al., 2021). The evolution of GNNs has been 083 marked by significant advancements in their architectures and learning mechanisms, with a recent 084 focus on incorporating attention mechanisms to enhance their expressive power and performance. 085 The introduction of attention in GNNs was largely inspired by the success of transformers in natural language processing (Vaswani et al., 2017). Graph Attention Networks (GATs) (Veličković et al., 2017) were among the first to incorporate self-attention into GNNs, allowing nodes to attend differ-087 880 ently to their neighbors based on learned attention weights. This innovation significantly improved the model's ability to capture complex relationships within graph structures. 089

Building upon the success of GATs, several variants and extensions have been proposed. GraphiT 091 (Mialon et al., 2021) introduced a generalization of transformer architectures to graph-structured 092 data, incorporating positional encodings and leveraging the power of multi-head attention mechanisms. Similarly, the Structure-Aware Network (SAN) (Kreuzer et al., 2021) proposed a novel attention mechanism that explicitly considers the structural properties of graphs, leading to improved 094 performance on various graph-based tasks. 095

- 096 Recent studies on Large Language Models (LLMs) and Vision Transformers (ViTs) have revealed 097 the presence of MAs within their internal states, specifically in the attention layer's output (Xiao 098 et al., 2023; Sun et al., 2024). This phenomenon prompted investigations into the role of these activations in model behavior, performance, and potential vulnerabilities. Similar observations were 099 made in Vision Transformers (ViTs) (Darcet et al., 2023; Dosovitskiy et al., 2020), suggesting that 100 the presence of MAs might be a common feature in transformer-based architectures across different 101 domains. These findings have led to a growing interest in understanding the implications of MAs 102 for model interpretability, robustness, and potential vulnerabilities to adversarial attacks. 103
- 104 The study of internal representations in deep learning models has been a topic of significant interest 105 in the machine learning community. Works such as Bau et al. (2020) have explored the interpretability of neural networks by analyzing activation patterns and their relationships to input features and 106 model decisions. However, the specific phenomenon of MAs in GNNs has remained largely unex-107 plored until now, representing a crucial gap in our understanding of these models.

The intersection of adversarial attacks and GNNs is another relevant area of study that relates to the investigation of MAs. Previous work has explored various attack strategies on graph data, including topology attacks, feature attacks, adversarial training and hybrid approaches (Sun et al., 2022a; Gosch et al., 2024). However, the potential vulnerabilities introduced by MAs represent a novel direction for research in this field. Understanding how MAs might be exploited or manipulated by adversarial inputs could lead to the development of more robust GNN architectures.

However, in the broader context of neural network analysis, techniques for probing and interpreting
model internals have been developed. Methods such as feature visualization (Olah et al., 2017) and
network dissection (Bau et al., 2017) have provided insights into the functions of individual neurons
and layers in convolutional neural networks. Adapting and extending these techniques to analyze
MAs in GNNs could provide valuable insights into their role and impact in possible future works.

Finally, the study of attention mechanisms in various neural network architectures has also yielded insights that may be relevant to understanding MAs in GNNs. Work on attention flow (Abnar & Zuidema, 2020) and attention head importance (Michel et al., 2019) in transformer models has shown that not all attention heads contribute equally to model performance, and some may even be pruned without significant loss of accuracy. These findings raise questions about whether similar patterns might exist in graph transformer models and how they might relate to the presence of MAs.

125 126

127

# 3 TERMINOLOGY OF MASSIVE ACTIVATIONS IN GNNS

Building upon the work on MAs in LLMs (Sun et al., 2024), we extend this investigation to GNNs, focusing specifically on graph transformer architectures. Our study encompasses various models, including GraphTransformer (GT) (Dwivedi & Bresson, 2021), GraphiT (Mialon et al., 2021), and Structure-Aware Network (SAN) (Kreuzer et al., 2021), applied to diverse task datasets such as ZINC, TOX21, and OGBN-PROTEINS (see Appendix A, B, C for details on models' configurations and datasets' composition). This comprehensive approach allows us to examine the generality of MAs across different attention-based GNN architectures.

- 135
- 136 137

## 3.1 CHARACTERIZATION OF MASSIVE ACTIVATIONS

MAs in GNNs refer to specific activation values that exhibit unusually high magnitudes compared to
 the typical activations within a layer. These activations are defined by the following criteria, where
 an activation value is intended to be its absolute value:

Magnitude Threshold: An activation is classified as massive if its value exceeds a predetermined threshold. This threshold is typically set to a value that is significantly higher than the average activation value within the layer, ensuring that only the most extreme activations are considered.

Relative Threshold: In the paper by Sun et al. (2024), MAs were defined as at least 1,000 times larger than the median activation value within the layer. This relative threshold criterion helped differentiate MAs from regular high activations that might occur due to normal variations in the data or model parameters.

 $MAs = \{a \mid a > 100 \text{ and } a > 1000 \times \text{median}(\mathbf{A})\}$ 

- 149 The formal definition was represented as:
- 150
- 151 152

159

161

where  $\mathbf{A}$  represents the set of activation values in a given layer.

However, in contrast to previous studies that employed a fixed relative threshold, our approach adopts a more rigorous method. We estimate MAs by comparing the distributions of activation ratios between a base, untrained model with Xavier weight initializations (Glorot & Bengio, 2010), and a fully trained model. This method ensures a more precise identification of MAs based on empirical data rather than an arbitrary fixed threshold. In this way, the untrained model serves as a reference for identifying unusual activations that emerge during training.

160 3.1.1 DETECTION METHODOLOGY

For both the base and trained models, we detected the MAs following a systematic procedure:



Figure 1: Comparison of MAs for trained vs base models, along all the edges. Activation values have been normalized within each layer by the layer's edge median. Represented ratios have been sorted increasingly for each layer independently.

**Normalization**: We normalized the activation values within each layer, dividing them by the edge 188 median on the layer, to account for variations in scale between different layers and models. This 189 normalization step ensures a consistent basis for comparison. The choice of dividing by the edge 190 median comes from the huge amount of MAs being present, since almost every edge in the layers 191 presenting MAs holds at least one MA, as shown from Figure 1. This is probably caused by the 192 fact that attention is computed between pairs of adjacent nodes only, in contrast to LLMs where it 193 is computed among each pair of tokens, therefore the model tends to spread MAs among almost all 194 the edges to make them "available" to the whole graph. Indeed, Figure 1 indicates that MAs are 195 a common phenomenon across different models and datasets, that they are not confined to specific 196 layers but are distributed throughout the model architecture, and that MAs are an inherent characteristic of the attention-based mechanism in graph transformers and related architectures, not strictly 197 dependent on the choice of the dataset.

Batch Analysis: We analyzed the activations on a batch-by-batch basis, minimizing the batch size, to have suitable isolation between the MAs and to ensure that the detection of MAs is not influenced by outliers in specific samples. For each activation, we computed the ratio of its magnitude to the edge median:

$$ratio(activation) = \frac{abs(activation)}{median(abs(edge\_activations))}$$
(1)

206 207

203 204 205

182

183

185

187

and activations whose ratio exceeds the threshold are flagged as massive. Then, we considered the
 maximum ratio of each batch to detect those containing MAs.

Layer-wise Aggregation: We performed this analysis across multiple layers of the model to identify
 patterns and layers that are more prone to exhibiting MAs. This layer-wise aggregation helps in
 understanding the hierarchical nature of MAs within the model.

Figure 2 reports the analysis results. The batch ratios significantly increase in the trained transformers, concerning base ones, often even overcoming the threshold of 1000 defined by previous works (Sun et al., 2024), showing the presence of MAs in graph transformers, too.



Figure 2: Comparison of MAs on trained against base models, without the use of Explicit Bias Term. Represented ratios have been sorted increasingly for each layer independently.

241

259

260

261

262

264

## 4 METHODOLOGY AND OBSERVATIONS

Focusing on edge features, first, we analyzed the ratio defined in Equation (1), taking the maximum for every batch, across the layers of each selected model and dataset, and visually compared the outcomes to value ranges obtained using the same model in a base state (with its parameters randomly initialized, without training) to verify the appearance of MAs. The graphical comparison, reported in Figure 2, shows ratios over the base range in most of the trained models, representing MAs.

To better characterize MAs, we studied their distribution employing the Kolmogorov-Smirnov statistic (Chakravarti et al., 1967). We found that a gamma distribution well approximates the negative logarithm of the activations' magnitudes, as well as their ratios. Figure 3a shows this approximation for a base model layer. We point out that, according to the existing definition, items on the left of the -3 are MAs.

We then compared the distributions of the log-values between the base and trained models. Figure 3 illustrates this comparison, highlighting a significant shift in the distribution of the trained model compared to the base model. Moreover, this shift underscores the emergence of MAs during the training process, affirming that the threshold around  $-\log(\text{ratio}) = -3$  (e.g., a ratio of 1000 or higher) effectively captures these significant activations, though sometimes it appears to be slightly shifted to the right as in Figure 3c.

- 258 When MAs appear, we have found two possible phenomenons:
  - A lot of massive activation values are added on the left-hand side of the distribution, preventing a good approximation (Figure 3b).
  - A few values appear on the left-hand side of the distribution, as spikes or humps or outof-distribution values, which may or may not deteriorate the approximation, as shown in Figures 3c and 3d.
- For example, histogram in Figure 3a represents the base model with untrained weights (only Xavier initialization). The gamma approximation fits the sample histogram well, with a low Kolmogorov-Smirnov (KS) statistic of 0.020, indicating a very nice fit.
- Figure 3b shows that the distribution of the trained model exhibits a significant shift due to a big hump appearing on the left side, representing extreme activation ratios (MAs). Indeed, the gamma



Figure 3: Activation distributions for base and trained (with MAs) models. In Figure 3d we clearly distinguish a spike on the left of the distribution, corresponding to a ratio of 1000 ( $-\log(ratio) = -3$ ), which identifies the separation between the basic and massive regimes. The approximation pdf is rescaled to match the histogram scale.

approximation does not fit well, with a higher KS statistic of 0.168, indicating a poor match caused by the presence of MAs.

Moreover, in the histogram of Figure 3d the trained model's distribution exhibits a clear spike on the left side at  $-\log(\text{ratio}) = -3$ , corresponding to a ratio of 1000. This separation indicates the distinction between basic and massive activation regimes. The gamma distribution doesn't fit well this time, because of this spike preventing a good approximation, with a KS statistic of 0.027 highlighting the model's shift due to training.

Figure 3c also shows the trained model's distribution, with a noticeable hump on the left side indicating MAs. The gamma approximation fits better than in Figure 3d, with a KS statistic of 0.019, but still indicates the presence of MAs in the trained model, meaning that MAs have been added on the left-hand side of the distribution.

310 4.1 INSIGHTS AND IMPLICATIONS

296 297

298

299

309

311

313

314 315

316

317

318 319

320

321

322

312 From Figures 1 and 2 we can highlight the following points.

#### 1. Dataset Influence:

• The ZINC and OGBN-PROTEINS datasets consistently show higher activation values across all models compared to TOX21, suggesting that the nature of these datasets significantly influences the emergence of MAs. Even though many MAs are emerging form GT on TOX21.

#### 2. Model Architecture:

- Different GNN models exhibit varying levels of MAs. For instance, GraphTransformer and GraphiT tend to show more pronounced MAs than SAN, indicating that model architecture plays a crucial role.
- 3. Impact of Attention Bias:

356

326	L				
327	Dataset	Model	Test loss	Test loss (EBT)	
328			1000		
329	ZINC	GraphTransformer	0.26	0.29	
330		GraphiT	0.13	0.31	
331		SAN	0.18	0.27	
332					
333	TOX21	GraphTransformer	0.25	0.29	
334		GraphiT	0.38	0.32	
335		SAN	0.38	0.31	
336	OGRN PROTEI	NS GraphTransformer	0.13	0.12	
337	OOBIV-I KOTEI	GraphiT	0.13	0.12	
338		SAN	0.14	0.13	
339			0.12	0110	
340			bect that MAs have the function of learned bias, showing the ag bias at the attention layer. This holds for LLMs and Vi II, as shown in Figure 2 where the presence of MAs is affec- the Explicit Bias Term on the attention. Figure 4 and text re intrinsic to the models' functioning, being anti-correlated		
341	• Previous works	suspect that MAs have the			
342	disappear introd	ucing bias at the attention			
343	for our GNNs as	well, as shown in Figure			
344	the introduction	of the Explicit Bias Tern			
345	suggest that MA	As are intrinsic to the mod			
346	the learned bias.				
347					
348	The consistent observation of MAs	MAs in edge features, acro	oss various C	INN models and datasets	
349	to a fundamental characteristic	c of how these models pro	cess relation	nal information.	
350	Inspired by recent advancement	nts in addressing bias insta	ability in LL	Ms (Sun et al., 2024), w	
351	duced an Explicit Bias Term (	EBT) into our graph transf	former mod	els. This bias term is dis	
352	to counteract the emergence of	of MAs by stabilizing the	activation 1	nagnitudes during the a	
353	computation. The EBT is com	puted as follows:			
354		$m{b}_e = m{Q}m{k}m{e}'$			
0					

324 Table 1: Comparison of test loss with and w/o bias for the different models and datasets. In bold the 325 worst performances.

$$e = Qke'$$
 (2)

$$\boldsymbol{b}_v = \operatorname{softmax}(\boldsymbol{A}_e)\boldsymbol{v}',\tag{3}$$

where  $k, e, v \in \mathbb{R}^d$  are the key, edge, and node bias terms (one per each attention head),  $A_e$  is the 357 edge attention output, and d the corresponding hidden dimension.  $b_e$  and  $b_v$  represent the edge and 358 node bias terms and are added to the edge and node attention outputs, respectively. By incorporating 359 EBT into the edge and node attention computations, and adding bias in the linear projections of the 360 attention inputs, we regulated the distribution of activation values, thus mitigating the occurrence of 361 MAs. 362

363 As shown in Figure 4, the introduction of these bias terms significantly reduces the frequency and magnitude of MAs, bringing the activation ratios closer to those observed in the base models. The 364 effect of EBT is evident across all the different datasets. Whether it's ZINC, TOX21, or OGBN-365 PROTEINS, the activation ratios are brought closer to the baseline levels observed in the untrained 366 models. This consistency underscores the general applicability of EBT in various contexts and 367 downstream tasks. Moreover, Figure 4 shows that EBT mitigates MAs across different layers of 368 the models. This is crucial as it indicates that EBT's effect is not limited to specific parts of the 369 network but is extended throughout the entire architecture. For example, GraphTransformer on 370 ZINC without EBT shows MAs frequently exceed  $10^4$ , while when EBT has been applied these 371 ratios are significantly reduced, aligning more closely with the base model's range. 372

Table 1 shows that EBT does not systematically influence the test loss equally across different mod-373 els and datasets. We have considered the test loss metric to keep the approach general, making it 374 extendable to different downstream tasks. This ensures that the proposed method can be applied 375 broadly across various applications of graph transformers. 376

Although the test loss remains relatively unchanged with the introduction of EBT, its presence helps 377 in mitigating the occurrence of MAs, as evidenced by the reduction in extreme activation values



Figure 4: Comparison of MAs on trained against base models, with the use of Explicit Bias Term. Represented ratios have been sorted increasingly for each layer independently.

403

404 405

406

observed in earlier figures. By analyzing these results, it becomes evident that while EBT does not drastically alter the test performance, it plays a crucial role in controlling activation anomalies, thereby contributing to the robustness and reliability of graph transformer models.

In the next section, we will demonstrate how attacking the model with and without MAs can directly impact the robustness of the architectures. This will provide deeper insights into the robustness of graph transformers in the presence of MAs, suggesting their potential pitfalls.

407 408 409

410

426 427 428

## 5 EXPLICIT BIAS ATTACK

411 The study of adversarial attacks on GNNs has become increasingly important as these models are 412 deployed in critical applications. While various attack strategies have been explored (Zügner et al., 413 2018; Sun et al., 2022b), the vulnerability introduced by MAs remains largely unexplored. Un-414 derstanding how MAs can be exploited by adversaries is crucial for developing more robust GNN 415 architectures and their downstream tasks. In this section, we propose the Explicit Bias Attack, a 416 gradient-based method designed to exploit MAs and assess model robustness. Our approach is inspired by gradient ascent attacks previously applied to image classifiers (Goodfellow et al., 2014) 417 418 and adapted for graph data (Dai et al., 2018). By analyzing the effectiveness of gradient ascent attack with and without the presence of EBT and MAs, we aim to provide insights into the role of 419 these activations in model fragility. 420

Therefore, inspired by previous section, we exploited EBT as computed in Equations (2) and (3) to analyze the importance of MAs for a gradient ascent attack at test time, where noise (added to the input feature embedding) is learned to directly maximize the loss function. The effectiveness of an attack is evaluated by comparing the average test loss before and after the attack (i.e., with random and optimized noise, with the same standard deviations, respectively), using a gain defined as

attack gain = 
$$\frac{\text{optimized noise loss} - \text{random noise loss}}{\text{random noise loss}}$$
 (4)

thus a higher gain means a more dangerous attack. We focus on GraphTransformer (GT) with
TOX21 because the presence of MAs in each layer - as shown by Figure 2 - highlights the MAs
effect for the attack, and compare the power of this method with and without the use of the EBT,
which calls off the model's MAs.

Table 2: GT on TOX21 – Comparison of the noise optimization strategy with (no EBT) and without
(EBT) MAs, due to the use of the explicit attention bias. The noise is optimized to maximize the
loss function, and the results are shown for 1000 epochs of noise optimization.

436	Noise dev.	Gain (no EBT, %)	Gain (EBT, %)
437			
438	0.01	1.50	1.41
439	0.03	1.83	1.78
440	0.10	4.73	2.53
441			

442

435

443 Table 2 shows a stable increase of gain when dealing with MAs, using noise with standard devia-444 tion values of 0.01, 0.03, and 0.1 (the input feature embedding has a standard deviation of about 445 0.9) optimized for 1000 epochs on the test set. Table 2 highlights that MAs can be dangerous for 446 the robustness of a model, and potentially exploited by attacks. These results indicate that a gradi-447 ent ascent attack is effective in degrading model performance, especially in the presence of MAs. However, the introduction of explicit bias, consistent with the reduction of MAs, can significantly 448 mitigate the impact of the attack, leading to more robust models. This highlights the importance 449 of considering bias in designing defenses against these types of adversarial attacks, to prevent them 450 from exploiting the presence of MAs. 451

In future work, to enable us to comprehensively assess the correlation between model robustness/fragility and the presence of MAs, we intend to delve deeper into different graph attack configurations while targeting MAs. This will offer a richer understanding of how these vulnerabilities
can be mitigated, in favor of more reliable models.

456 457

458 459

# 6 CONCLUSION AND FUTURE WORK

 This paper presents the first comprehensive study of MAs in attention-based GNNs. We have introduced a novel methodology for detecting and analyzing MAs, focusing on edge features in various graph transformer architectures across multiple benchmark datasets. Our findings reveal that MAs are prevalent across different models and datasets, and demonstrate that they could be effectively leveraged by adversaries to degrade the performance of GNNs.

We showed that the introduction of Explicit Bias Terms (EBT) can effectively mitigate the occurrence of MAs, leading to more stable activation distributions. However, our results also showed that this mitigation does not always translate to improved test performance, highlighting the complex role of MAs in GNNs' behavior.

Furthermore, we introduced the Explicit Bias Attack, a gradient-ascent adversarial framework, that
demonstrates how MAs, if not mitigated by EBT, can expose models to vulnerabilities in their tasks.
This further points out the importance of considering these activations in the context of model robustness.

Future research will expand this analysis to a wider range of architectures and advanced attack methods, further clarifying the influence of MAs on GNN performance and robustness, and potentially
leading to more interpretable and stable graph-based models. Specifically, future research could
explore:

477 478

479

480 481

482

483

484

- **Customized Adversarial MAs**: Developing more adversarial techniques to regulate and attack these activations to enhance model stability and performance, like injecting fake MAs or exploiting state-of-the-art graph attack methods.
- **Downstream-driven MAs**: Leveraging MAs for specific downstream task, investigating how to harness these significant activations to improve models and their interpretability on specific assignments such as link prediction or drug design.
- **Comparative Analysis**: Extending the study to additional models and datasets to generalize the findings further and uncover broader patterns.

486 These insights provide a deeper understanding of the internal mechanisms of attention-based GNNs 487 and highlight the way for improvements in graph learning models. By addressing the challenges and 488 opportunities presented by MAs, we can work towards developing more robust, interpretable, and 489 effective GNN architectures for a wide range of applications.

#### 491 REFERENCES 492

490

493

495

496

497

498

501

- Samira Abnar and Willem Zuidema. Quantifying attention flow in transformers. arXiv preprint arXiv:2005.00928, 2020. 494
  - David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pp. 6541-6549, 2017.
- David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba. 499 Understanding the role of individual units in a deep neural network. Proceedings of the National 500 Academy of Sciences, 117(48):30071–30078, 2020.
- Indra Mohan Chakravarti, Radha Govira Laha, and Jogabrata Roy. Handbook of methods of applied statistics. Wiley Series in Probability and Mathematical Statistics (USA) eng, 1967. 504
- Hanjun Dai, Hui Li, Tian Tian, Xin Huang, Lin Wang, Jun Zhu, and Le Song. Adversarial attack on 505 graph structured data. In International conference on machine learning, pp. 1115–1124. PMLR, 506 2018. 507
- Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need 509 registers. arXiv preprint arXiv:2309.16588, 2023.
- 510 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 511 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An 512 image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint 513 arXiv:2010.11929, 2020. 514
- 515 Vijay Prakash Dwivedi and Xavier Bresson. A generalization of transformer networks to graphs, 2021. URL https://arxiv.org/abs/2012.09699. 516
- 517 Chen Gao, Xiang Wang, Xiangnan He, and Yong Li. Graph neural networks for recommender 518 system. In Proceedings of the fifteenth ACM international conference on web search and data 519 mining, pp. 1623–1625, 2022. 520
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural 521 networks. In Proceedings of the thirteenth international conference on artificial intelligence and 522 statistics, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010. 523
- 524 Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial 525 examples. arXiv preprint arXiv:1412.6572, 2014. 526
- Lukas Gosch, Simon Geisler, Daniel Sturm, Bertrand Charpentier, Daniel Zügner, and Stephan 527 Günnemann. Adversarial training for graph neural networks: Pitfalls, solutions, and new direc-528 tions. Advances in Neural Information Processing Systems, 36, 2024. 529
- 530 Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele Catasta, 531 and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs. Advances in neural information processing systems, 33:22118–22133, 2020. 532
- Ruili Huang, Menghang Xia, Dac-Trung Nguyen, Tongan Zhao, Srilatha Sakamuru, Jinghua Zhao, 534 Sampada A Shahane, Anna Rossoshek, and Anton Simeonov. Tox21challenge to build predictive models of nuclear receptor and stress response pathways as mediated by exposure to environmen-536 tal chemicals and drugs. Frontiers in Environmental Science, 3:85, 2016.
- John J Irwin, Teague Sterling, Michael M Mysinger, Erin S Bolstad, and Ryan G Coleman. Zinc: 538 a free tool to discover chemistry for biology. Journal of chemical information and modeling, 52 (7):1757-1768, 2012.

540 541 542	Devin Kreuzer, Dominique Beaini, Will Hamilton, Vincent Létourneau, and Prudencio Tossou. Re- thinking graph transformers with spectral attention. <i>Advances in Neural Information Processing</i> <i>Systems</i> , 34:21618–21629, 2021.	
543 544	Andreas Mayr, Günter Klambauer, Thomas Unterthiner, and Sepp Hochreiter. Deeptox: toxicity	
545	prediction using deep learning. Frontiers in Environmental Science, 3:80, 2016.	
546	Grégoire Mialon, Dexiong Chen, Margot Selosse, and Julien Mairal. Graphit: Encoding graph	
547 548	structure in transformers. arXiv preprint arXiv:2106.05667, 2021.	
549 550	Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? Advances in neural information processing systems, 32, 2019.	
551	Shangija Min Zhan Gao, Jing Dang, Liang Wang, Va Oin, and Po Eang. Staan, a spatial temporal	
552	graph neural network framework for time-evolving social networks. <i>Knowledge-Based Systems</i> .	
553	214:106746, 2021.	
554	Grand Mathematical Distantic Financial for the data distantic state of the data distantic	
555 556	tematic review. <i>Expert Systems With Applications</i> , pp. 122156, 2023.	
557 558	Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. <i>Distill</i> , 2(11):e7, 2017.	
559	Lichao Sun Yingtong Dou Carl Yang Kai Zhang Ji Wang S Yu Philin Lifang He and Bo Li	
560	Adversarial attack and defense on graph data: A survey. <i>IEEE Transactions on Knowledge and</i>	
561	Data Engineering, 35(8):7693–7711, 2022a.	
563	Lichoo Sun Vingtong Dou Corl Vong Koi Zhong Ji Wong S Vu Dhilin Lifong Ho and Po Li	
564	Adversarial attack and defense on graph data: A survey <i>IEEE Transactions on Knowledge and</i>	
565	Data Engineering, 35(8):7693–7711, 2022b.	
566		
567 568	Mingjie Sun, Xinlei Chen, J. Zico Kolter, and Zhuang Liu. Massive activations in large language models, 2024. URL https://arxiv.org/abs/2402.17762.	
569	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,	
570 571	Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural informa- tion processing systems, 30, 2017.	
572	Peter Veličković Guillem Cucurull Arantya Casanova, Adriana Romero, Pietro Lio, and Voshua	
573 574	Bengio. Graph attention networks. <i>arXiv preprint arXiv:1710.10903</i> , 2017.	
575	Yang Wang, Jin Zheng, Yuqi Du, Cheng Huang, and Ping Li. Traffic-ggnn: predicting traffic flow	
576	via attentional spatial-temporal gated graph neural networks. <i>IEEE Transactions on Intelligent</i>	
577	Transportation Systems, 23(10):18423–18432, 2022.	
578	Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang, and Philip S. Yu. A	
579	comprehensive survey on graph neural networks. IEEE Transactions on Neural Networks and	
580	Learning Systems, 32(1):4–24, 2021. doi: 10.1109/TNNLS.2020.2978386.	
587	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. Efficient streaming	
582	language models with attention sinks. arXiv preprint arXiv:2309.17453, 2023.	
584		
585 586	Xiao-Meng Zhang, Li Liang, Lin Liu, and Ming-Jing Tang. Graph neural networks and their current applications in bioinformatics. <i>Frontiers in genetics</i> , 12:690049, 2021.	
587	Daniel Zügner, Amir Akbarnejad, and Stephan Günnemann. Adversarial attacks on neural networks	
588	for graph data. In Proceedings of the 24th ACM SIGKDD international conference on knowledge	
589	<i>discovery &amp; data mining</i> , pp. 2847–2856, 2018.	
590		
591		
592		
593		

594 595	A DATASET COMPOSITION
596	This section provides additional details on the used datasets throughout the experiments.
597 598 599 600	The ZINC dataset (Irwin et al., 2012) is a benchmark collection for evaluating GNNs in molecular chemistry, where molecules are represented as graphs with atoms as nodes and chemical bonds as edges. Contents include:
601 602 603	• Graphs: The dataset includes over 250,000 molecular graphs. Each molecule is represented by a graph with nodes (atoms) and edges (bonds), incorporating various bond types (e.g., single, double, triple).
604 605	• Node Features: Atoms are described by features that capture their chemical properties, such as atom types, hybridization states, and other atomic attributes.
606 607	• Edge Features: Bonds between atoms are characterized by features representing bond types and additional chemical information.
608 609 610 611	• Task: The primary task is <b>graph regression</b> , where the goal is to predict continuous values associated with each molecule. This often involves predicting molecular properties such as solubility or biological activity.
612 613 614	ZINC is useful for evaluating GNNs' performance in learning molecular representations and predict- ing continuous chemical properties, providing insights into the model's ability to generalize across diverse chemical compounds.
615 616 617 618	The TOX21 dataset (Mayr et al., 2016; Huang et al., 2016) is designed for toxicity prediction and focuses on classifying chemical compounds based on their potential toxicity. It is part of the Tox-icology Data Challenge and features molecular graphs with associated toxicity labels. Contents include:
619 620 621 622	• Graphs: The dataset consists of molecular graphs where nodes represent atoms and edges represent chemical bonds. It includes thousands of molecules with toxicity annotations, and it consists of 7,831 graphs with each graph representing a molecular structure with associated toxicity labels.
623 624	• Node Features: Atoms are encoded with features representing their types, hybridization states, and other chemical properties.
625 626	• Edge Features: Bonds are detailed with features indicating bond types and additional chem- ical attributes.
627 628 629 630	• Task: The main task is <b>multi-label graph classification</b> , where each molecule is classified into multiple toxicity categories. This allows for the prediction of various toxicity endpoints simultaneously.
631 632 633	TOX21 is valuable for assessing GNN models in predicting toxicity from molecular structures, which is crucial for drug discovery and safety evaluation, providing a benchmark for multi-label classification tasks.
634 635 636	The OGBN-PROTEINS dataset, part of the Open Graph Benchmark (OGB) (Hu et al., 2020), fo- cuses on protein function prediction. It contains one large graph representing protein structures, with nodes corresponding to amino acids and edges to their interactions. Contents include:
638 639 640 641	• One Large Graph: OGBN-PROTEINS contains 54,879 nodes and 89,724 edges. These nodes represent amino acids in protein structures, and edges represent interactions or bonds between these amino acids. It includes various protein structures used for functional prediction.
642 643	• Node Features: Amino acids are described by features capturing biochemical properties, such as amino acid type, secondary structure, and other relevant attributes.
644 645	• Edge Features: Edges denote interactions between amino acids and include features reflect- ing the nature of these interactions or spatial relationships.
646 647	• Task: The task is <b>multi-label node classification</b> , where the goal is to predict multiple functional categories for each amino acid node in the protein graph. This involves classifying nodes into various functional classes based on their role in the protein's functionality.

648
 649
 650
 650
 651
 OGBN-PROTEINS is suitable for evaluating GNNs on biological data, specifically in predicting protein functions based on structural information. It provides insights into how well models can handle multi-label node classification tasks in a complex biological context.

#### B MODEL ARCHITECTURE

This section provides additional details on the models' architecture used throughout all the experiments, namely GT (Dwivedi & Bresson, 2021), GraphiT (Mialon et al., 2021) and SAN (Kreuzer et al., 2021). These graph-transformer architectures integrate the principles of both GNNs and transformers, leveraging the strengths of attention mechanisms to capture intricate relationships within graph-structured data. Graph transformers extend the transformer structure, typically used for sequence data, to graphs, operating by embedding nodes and edges into higher-dimensional spaces and then applying multi-head self-attention mechanisms to capture dependencies between nodes.

Mathematically, let  $\mathcal{G} = (V, E)$  be a graph where  $V = \{v_1, ..., v_n\}$  is the set of nodes and  $E \subseteq V \times V$  is the set of edges. Each node  $v_i$  is associated with a feature vector  $x_i \in \mathbb{R}^d$ , and each edge  $(v_i, v_j)$  may have an edge feature  $e_{ij} \in \mathbb{R}^k$ . Therefore, graph transformer models are designed as follows.

666 INPUT EMBEDDING

The initial node features  $X = [x_1, ..., x_n]^T \in \mathbb{R}^{n \times d}$  are typically projected to a higher-dimensional space:

$$\boldsymbol{H}^{(0)} = \boldsymbol{X} \boldsymbol{W}_{in} + \boldsymbol{b}_{in} \tag{5}$$

where  $W_{in} \in \mathbb{R}^{d \times d'}$  is a learnable weight matrix and  $b_{in} \in \mathbb{R}^{d'}$  is a bias vector.

673 POSITIONAL ENCODING

To capture structural information, positional encodings  $P \in \mathbb{R}^{n \times d'}$  are often added:

679

680

681

682

685 686 687

688

693 694

696

697

699 700

670

672

652

653

# $m{H}^{(0)} = m{H}^{(0)} + m{P}$

678 MULTI-HEAD ATTENTION LAYER

The core of a graph transformer is the multi-head attention mechanism. For each attention head i (out of h heads) there are also:

1. Query, Key, and Value Projections:

$$\boldsymbol{Q}_i = \boldsymbol{H}^{(l)} \boldsymbol{W}_i^Q \tag{7}$$

(6)

$$Q_i = H^{(l)} W_i^{(l)}$$

$$K_i = H^{(l)} W_i^K$$
(8)

$$\boldsymbol{V}_i = \boldsymbol{H}^{(l)} \boldsymbol{W}_i^V \tag{9}$$

where  $W_i^Q, W_i^K, W_i^V \in \mathbb{R}^{d' \times d_k}$  are learnable weight matrices, and  $d_k = d'/h$ . 2. Attentions Scores (node features only):

$$\boldsymbol{A}_{i} = \operatorname{softmax}\left(\frac{\boldsymbol{Q}_{i}\boldsymbol{K}_{i}^{T}}{\sqrt{d_{k}}} + \boldsymbol{M}\right), \tag{10}$$

where  $M \in \mathbb{R}^{n \times n}$  is a mask matrix to enforce the graph structure:

$$M_{i,j} = \begin{cases} 0 & \text{if } (v_i, v_j) \in E \text{ or } i = j \\ -\infty & \text{otherwise.} \end{cases}$$
(11)

3. Output of each head:

$$\mathbf{head}_i = \mathbf{A}_i \mathbf{V}_i. \tag{12}$$

4. Concatenation and Projection:

$$\boldsymbol{H}' = \operatorname{Concat}(\mathbf{head}_1, ..., \mathbf{head}_h) \boldsymbol{W}^O, \tag{13}$$

where  $\boldsymbol{W}^{O} \in \mathbb{R}^{d' \times d'}$  is a learnable weight matrix.

683 684

702 703	Feed-Forward Network (FFN)	
704 705 706	Each attention layer is typically followed by a position-wise feed-forward network: $FFN(\boldsymbol{x}) = \max(0, \boldsymbol{x}\boldsymbol{W}_1 + \boldsymbol{b}_1)\boldsymbol{W}_2 + \boldsymbol{b}_2$	(14)
707	where $W_1 \in \mathbb{R}^{d \times d_{ff}}$ , $W_2 \in \mathbb{R}^{d_{ff} \times d}$ , $b_1 \in \mathbb{R}^{d_{ff}}$ , and $b_2 \in \mathbb{R}^d$ are learnable parameters.	
708 709	LAYER NORMALIZATION AND RESIDUAL CONNECTIONS	
710 711 712 712	Each sub-layer (attention and FFN) employs a residual connection followed by layer normaliz $H^{(l+1)} = \text{LayerNorm}(H^{(l)} + \text{Sublayer}(H^{(l)}))$ where Sublayer is either the multi-head attention or the FFN.	ation: (15)
714	Edge Feature Integration	
715 716	GraphTransformer, GraphiT and SAN incorporate edge features:	
717	1. In attention computation:	
719 720	$A_{i,j} =  ext{softmax}\left(rac{oldsymbol{q}_i^Toldsymbol{k}_j + f(oldsymbol{e}_{ij})}{\sqrt{d_k}} ight)$	(16)
721 722	where $f$ is a learnable function (e.g., a small neural network) that projects edge feature	ires.
723 724	2. In value computation: $v_{ij} = V_i + g(e_{ij})$ where <i>q</i> is another learnable function.	(17)
725 726	Global Node	
727 728	Some architectures introduce a global node $v_a$ connected to all other nodes to capture graph	ı-level
729 730	information: $m{h}_g^{(l+1)} = \operatorname{Attention}(m{h}_g^{(l)}, m{H}^{(l)})$	(18)
731 732	OUTPUT LAYER	
733 734	The final layer depends on the task:	
735	• For node classification: $y_{node} = \operatorname{softmax}(H_{node}^{(L)}W_{out} + b_{out})$	
736 737	• For graph classification: $Y_{graph} = MLP(Pool(H^{(L)}))$	
738 739 740	where Pool is a pooling operation (e.g., mean, sum, or attention-based pooling) to switch from node to graph embedding level.	single
741	TRAINING	
742 743 744	The model is typically trained end-to-end using backpropagation to minimize a task-specif function, such as cross-entropy for classification or mean squared error for regression.	c loss
745 746	C KOLMOGOROV-SMIRNOV TEST	
747 748 749 750 751 752	This section provides additional details on the Kolmogorv-Smirnov (KS) test (Chakravarti 1967) used to analyze the distribution of activations. The KS test is a non-parametric te compares the cumulative distribution functions of two samples. It is used to compare a sampl a reference probability distribution (one-sample KS test) or to compare two samples (two-s KS test) with each other.	et al., st that e with ample
753 754 755	In our study, we utilized the KS statistic to compare the distribution of activation values and after training (i.e. base against trained model), identifying Massive Activations (MAs primarily used the one-sample KS test to assess the goodness of fit between our observed acti- distributions and a theoretical gamma distribution.	before). We vation

# 756 C.1 ONE-SAMPLE KOLMOGOROV-SMIRNOV TEST

758 The one-sample KS test can typically be formulated as follows:

760 C.1.1 NULL HYPOTHESIS

759

761

763 764

770

771

778

779

780

782

784

785

786

787 788

789

791

792 793

794

796

797

798

799 800

801

The null hypothesis for the one-sample KS test is:

 $H_0$ : The sample data follows the specified distribution (in our case, a gamma distribution).

765 C.1.2 TEST STATISTIC

The KS statistic  $D_n$  is defined as the supremum of the absolute difference between the empirical cumulative distribution function (ECDF)  $F_n(x)$  of the sample and the cumulative distribution function (CDF) F(x) of the reference distribution:

$$D_n = \sup_{x} |F_n(x) - F(x)|$$
(19)

where  $\sup_x$  denotes the supremum of the set of distances.

774

C.1.3 EMPIRICAL CUMULATIVE DISTRIBUTION FUNCTION

For a given sample  $x_1, x_2, ..., x_n$ , the ECDF is defined as:

 $F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{x_i \le x}$ (20)

where  $\mathbf{1}_{x_i < x}$  is the indicator function, equal to 1 if  $x_i \le x$  and 0 otherwise.

#### 783 C.1.4 CRITICAL VALUES AND P-VALUE

The distribution of the KS test statistic under the null hypothesis can be calculated, which allows us to obtain critical values and p-values. The null hypothesis is rejected if the test statistic  $D_n$  is greater than the critical value at a chosen significance level  $\alpha$ , or equivalently if the p-value is less than  $\alpha$ .

## C.2 APPLICATION TO MAS DETECTION

In our experiments, we used the KS statistic to assess whether the distribution of activation ratios in our GNNs follows a gamma distribution. The process is as follows:

- 1. We computed the activation ratios for each layer of our models, as defined in Equation (1) of the main paper.
- 2. We took the negative logarithm of these ratios to transform the distribution.
- 3. We fit a gamma distribution to this transformed data using maximum likelihood estimation.
- 4. We performed a one-sample KS test to compare our sample data to the fitted gamma distribution.

The KS test statistic provides a measure of the discrepancy between the observed distribution of activation ratios and the theoretical gamma distribution. A lower KS statistic indicates a better fit, suggesting that the activation ratios more closely follow the expected distribution.

- 802 803 804 805
- C.3 INTERPRETATION IN THE CONTEXT OF MAS

Following the described procedure in Section C.2, we employed the KS statistic as quantitative/statistical measure to detect the presence of MAs:

808

• For untrained (base) models, we typically observed low KS statistics, indicating that the activation ratios closely follow a gamma distribution.

010	
810	• For trained models exhibiting MAs, we often saw higher KS statistics. This indicates a
811	departure from the gamma distribution, which we interpret as evidence of MAs.
812	• The magnitude of the KS statistic provided a quantitative measure of how significantly the
813	presence of MAs distorts the expected distribution of activation ratios
814	presence of thirds distorts the expected distribution of derivation ratios.
815	Moreover, we complemented our KS statistic results with visual inspections of the distributions and
816	other analyses as described in the main paper.
817	
818	
819	
820	
821	
822	
823	
023	
024	
020	
826	
827	
828	
829	
830	
831	
832	
833	
834	
835	
836	
837	
838	
839	
840	
841	
842	
843	
844	
845	
846	
847	
848	
849	
850	
851	
852	
853	
954	
855	
956	
050	
050	
858	
859	
860	
861	
862	
863	