



The 22nd IEEE International Conference on  
Mobile Ad-Hoc and Smart Systems  
(MASS 2025)

# *An Interpretable Multi-Modal Transformer-Based Intrusion Detection System Utilizing Log Messages and PCAP Files*

*Nadia Niknami<sup>1</sup>, Vahid Mahzoon<sup>2</sup>, Rajorshi Biswas<sup>3</sup>, Slobadan Vucetic<sup>2</sup>, and Jie Wu<sup>2</sup>*

*<sup>1</sup> Villanova University*

*<sup>2</sup> Temple University*

*<sup>3</sup> Penn State University(Berks)*

# Outline



Background & Motivation



Key Inside



Idea Overview



Evaluation & Conclusion

# Problem Background

- Challenges in traditional IDS:
  - Signature-based limitations,
  - Unseen attacks,
  - Poor interpretability
- Logs capture system-level events and context.
- PCAP files provide detailed network-level information.
- By fusing them, we get a more holistic understanding of network behavior

# Signature-based vs. Deep Learning vs. Ours

- Signature-based → misses zero-day
- CNN/LSTM-based → weak interpretability, ignores semantics
- Transformer fusion (ours) → interpretable and multi-modal

Prior deep models improved accuracy but still face three key issues:

- **Single modality** – using only packets or only logs.
- **Limited interpretability** – no way to explain detections.
- **Data inefficiency** – models need extensive labeled data

# Key Insight: Transformers Unify Modalities

Logs tell us *what happened*. PCAP shows *how it happened*.

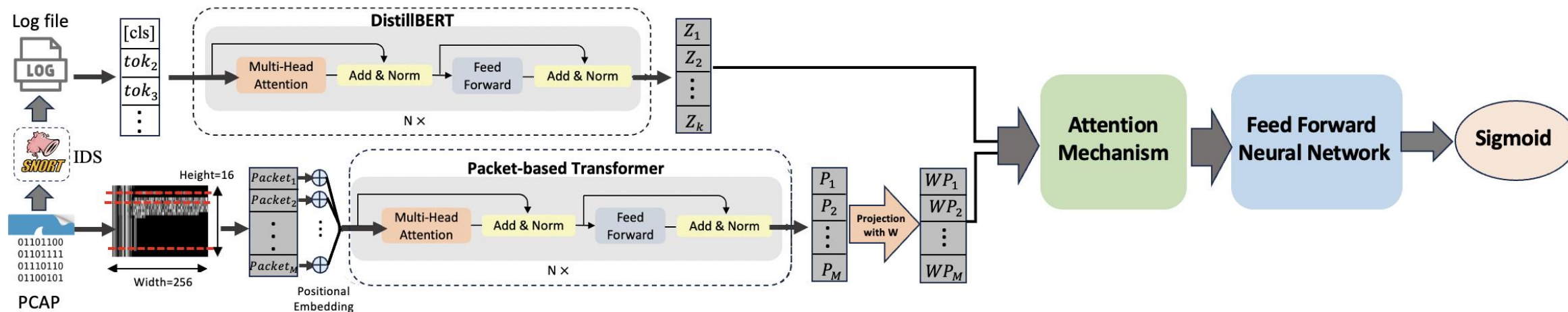
Our idea is to fuse them through a **multi-modal model** so that each modality compensates for the weaknesses of the other

- The **Transformer architecture** was originally designed for text, but later expanded to vision and multimodal tasks.
- We realized this flexibility makes Transformers ideal for merging **packet-level spatial patterns** and **log-level semantic patterns** into one unified framework

# Contributions

- A **multi-modal Transformer framework** for IDS.
- Integration of **DistilBERT** for semantic log embeddings and a **packet-based Transformer** for traffic data.
- An **attention mechanism** that makes the model interpretable by highlighting important packets or log entries.
- Extensive evaluation on CICIDS-2018, demonstrating superior accuracy and generalization to *zero-day attacks*

# Idea Overview



- A **DistilBERT** for logs
- A **packet-based Transformer** for PCAP

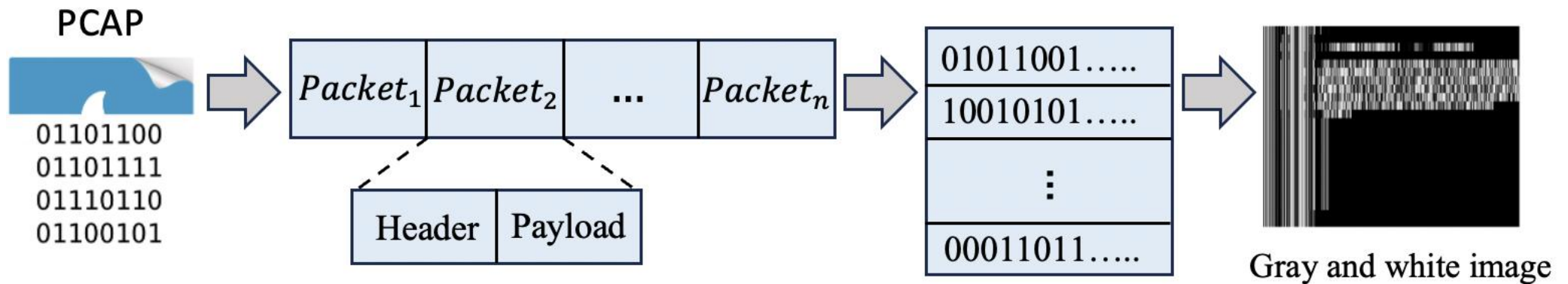
The outputs are fused with an **attention module** that learns which modality matters more for each flow

Three goals in bold:

1. Accurate detection (even zero-day)
2. Interpretable results
3. Efficient computation

# Data Preprocessing

- Log tokenization
- PCAP-to-image transformation



# Attention & Interpretability

- The attention mechanism assigns importance weights to each log token and packet.
- It helps the model decide which elements are most indicative of an attack.  
This not only boosts performance but makes decisions explainable.

# Evaluation

TABLE I: Performance metrics of different methods

Method	Individual Classes Accuracy				Overall		
	DDoS	BruteForce	DoS	Bot	Accuracy	Recall	F1-score
<b>DistilBERT</b>	0.7484	0.8119	0.7555	0.6975	0.7534	0.5073	0.6731
<b>CNN-Packet</b>	0.8662	0.9000	0.8950	0.8307	0.8784	0.8458	0.8745
<b>Packet-based Transformer</b>	0.9047	0.9250	0.9175	0.8537	0.9122	0.8776	0.8986
<b>TransIDS</b>	<b>0.9365</b>	<b>0.9700</b>	<b>0.9642</b>	<b>0.8975</b>	<b>0.9424</b>	<b>0.8948</b>	<b>0.9389</b>

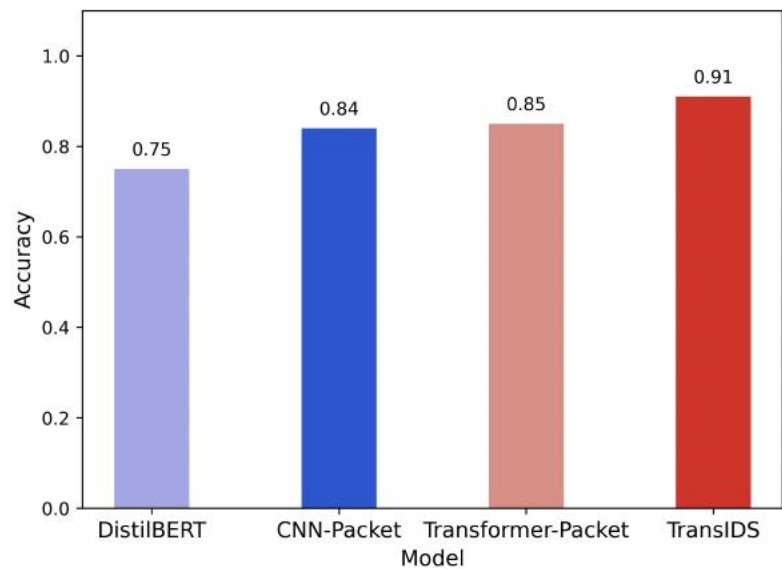


Fig. 4: Multi-class classification.

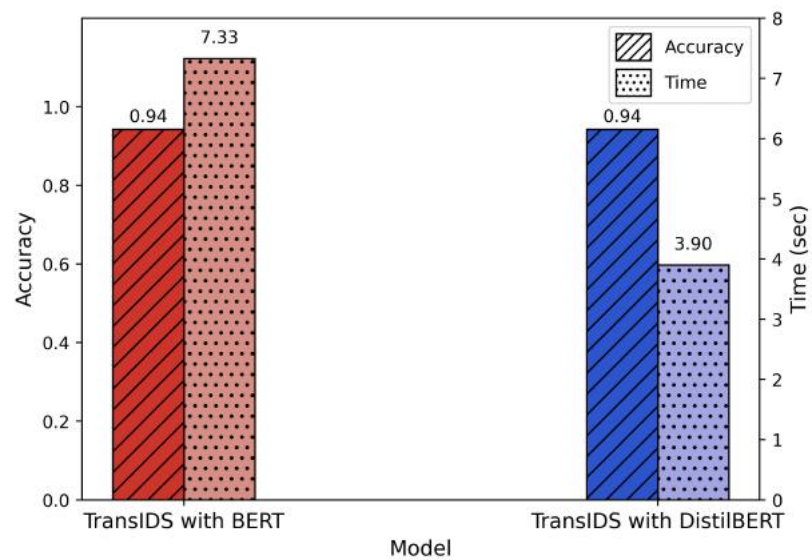


Fig. 5: BERT vs DistilBERT TransIDS.

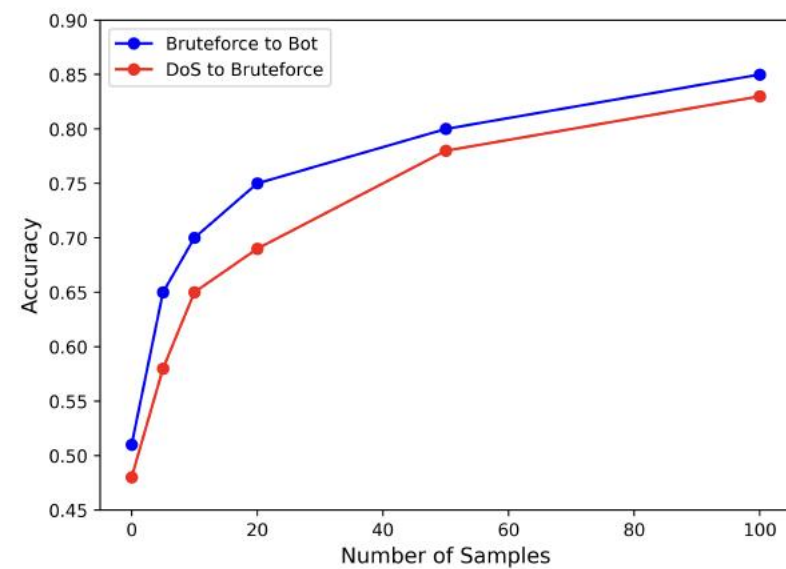
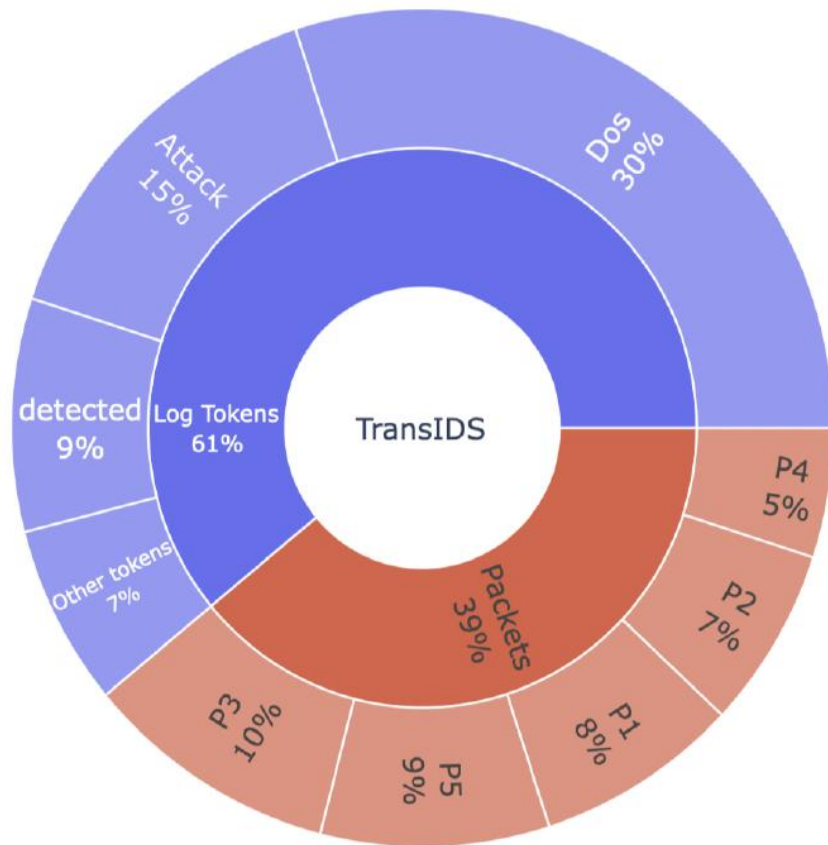
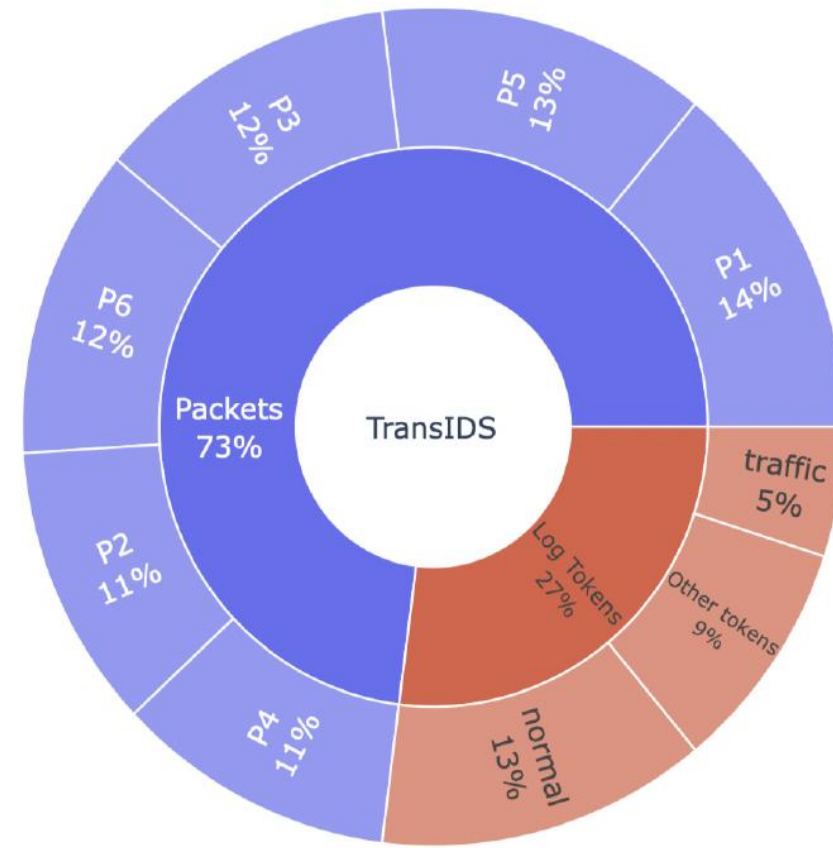


Fig. 6: Few-shot learning.



(a) Flow 1



(b) Flow 2

*(a) Flow 1 consists of 5 packets. It is detected by the log-based (signature-based IDS) as “Possible SYN DoS, Possible DoS attack detected”.*

*(b) Flow 2 contains 6 packets. It was not detected by the log-based system and identified as “It is a normal traffic”.*

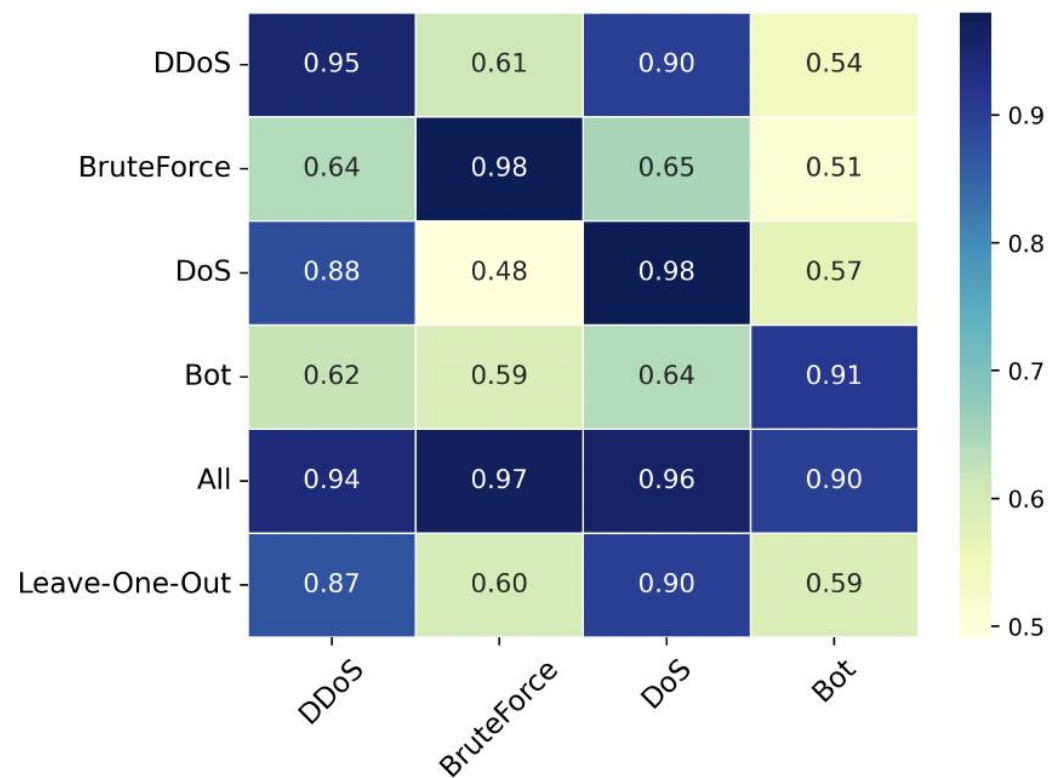


Fig. 8: Results on Zero-day Attacks

# Conclusion & Future Work

- In summary, **TransIDS** fuses PCAP and log data via Transformers to deliver interpretable, multi-modal intrusion detection. It outperforms prior deep models while maintaining transparency.
- Future directions include scaling to larger datasets and applying the framework to *real-time streaming logs* and *cross-domain attacks*