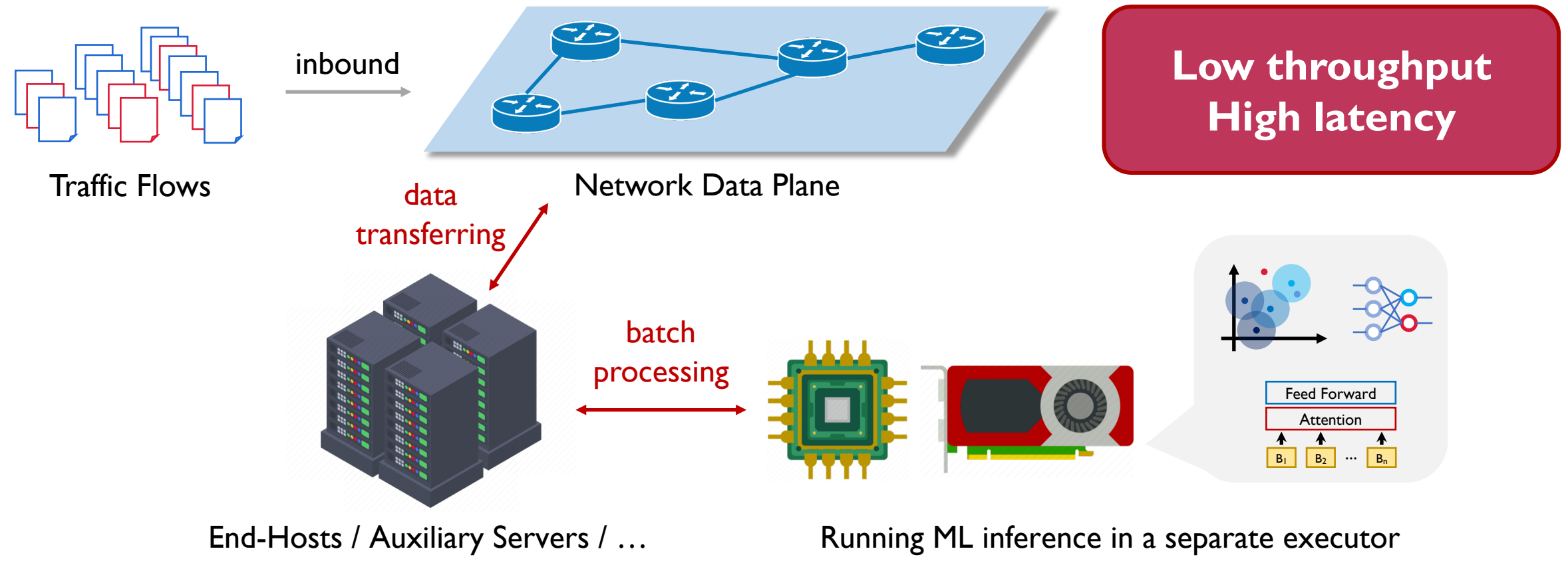# Brain-on-Switch:
# Towards Advanced Intelligent Network Data Plane via NN-Driven Traffic Analysis at Line-Speed

Jinzhu Yan, Haotian Xu, Zhuotao Liu✉, Qi Li, Ke Xu,
Mingwei Xu, Jianping Wu
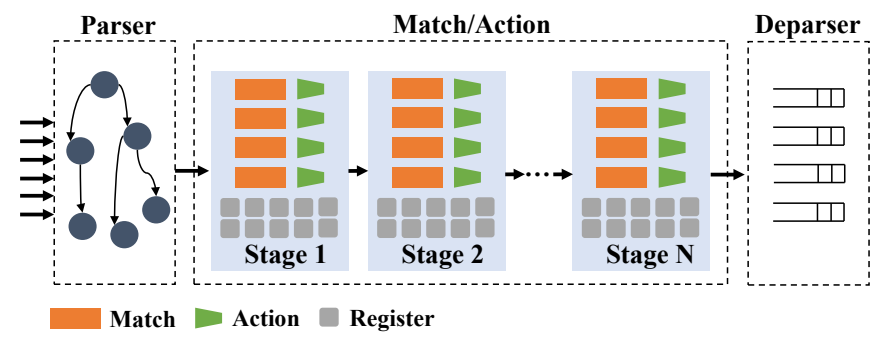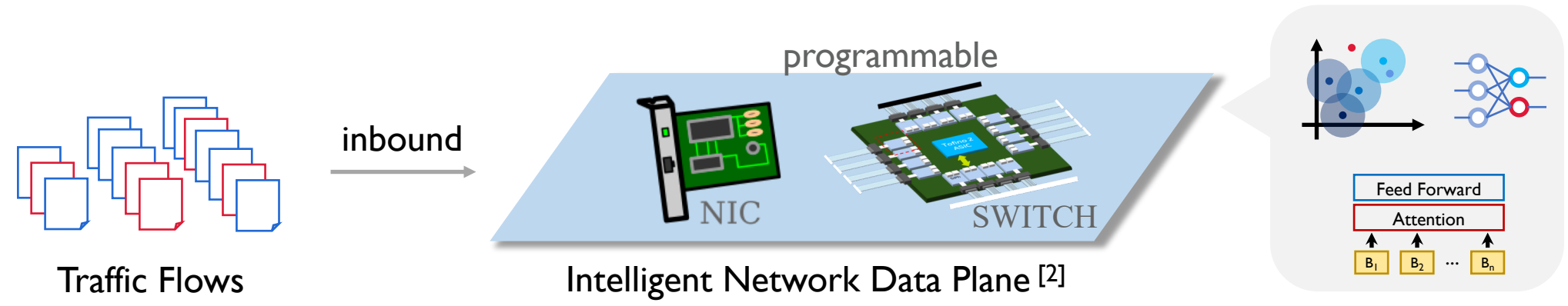
*Tsinghua University*

*April, 2024*

## Bottlenecks of the ML-based traffic analysis on dedicated executor[1]

Traffic Flows

inbound

Network Data Plane

**Low throughput
High latency**

data transferring

batch processing

End-Hosts / Auxiliary Servers / …

Running ML inference in a separate executor

Feed Forward

Attention

$B_1$  $B_2$  …  $B_n$

[1] Giuseppe Siracusano, Salvator Galea, Davide Sanvito, Mohammad Malekzadeh, Gianni Antichi, Paolo Costa, Hamed Haddadi, and Roberto Bifulco.
Re-architecting Traffic Analysis with Neural Network Interface Cards. In USENIX Symposium on Networked Systems Design and Implementation (NSDI), 2022.

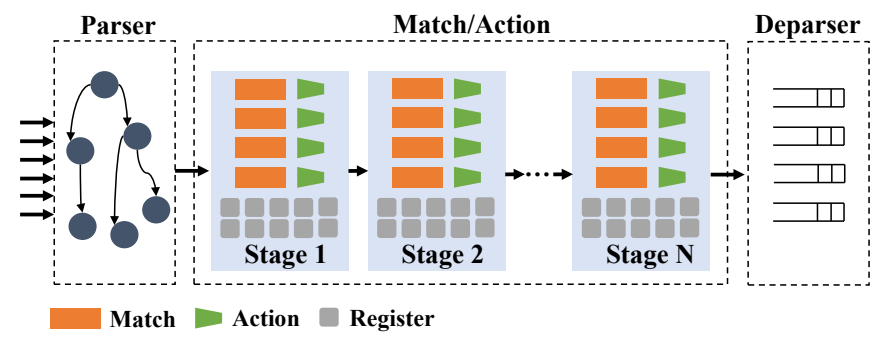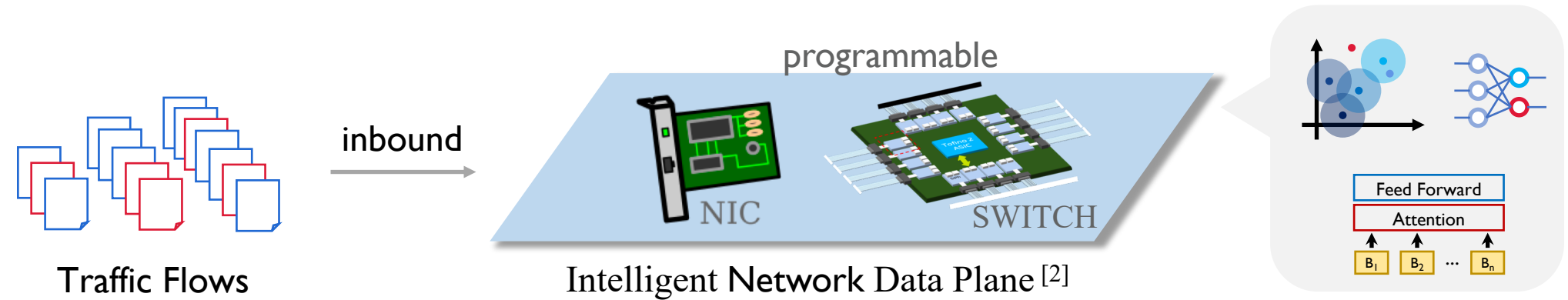## Today's ML-based traffic analysis can be forwarding-native



Traffic Flows

inbound

programmable

NIC

SWITCH

Intelligent Network Data Plane [2]

Feed Forward

Attention

$B_1$ $B_2$ ... $B_n$



Parser

Match/Action

Deparser

Stage 1

Stage 2

Stage N

Match

Action

Register

Protocol-Independent Switch Architecture (PISA)

**Enabling ML inference within network data plane**

1. **Customizable Packet Processing**

2. **Stateful and Persistent Storage**

[2] Guangmeng Zhou, Zhuotao Liu, Chuanpu Fu, Qi Li, and Ke Xu. An Efficient Design of Intelligent Network Data Plane. In USENIX Security Symposium (USENIX Security), 2023.

## Today's ML-based traffic analysis can be forwarding-native



Traffic Flows

inbound

programmable

NIC    SWITCH

Intelligent Network Data Plane [2]



Feed Forward

Attention

$B_1$ $B_2$ ... $B_n$



Parser    Match/Action    Deparser
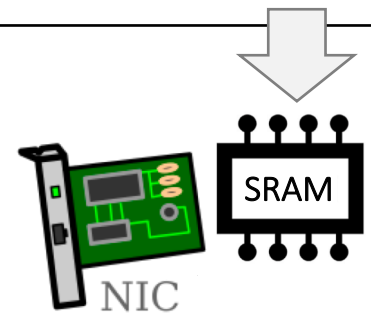
Stage 1    Stage 2    Stage N
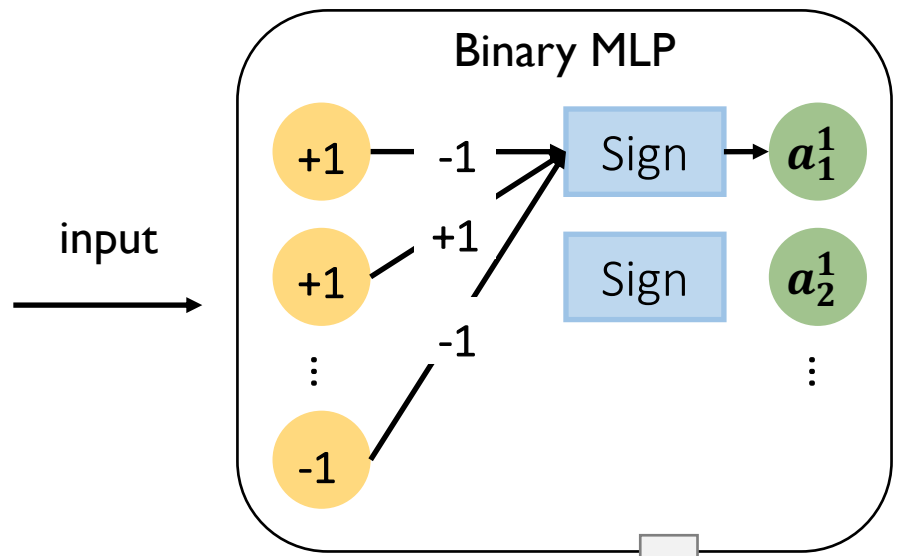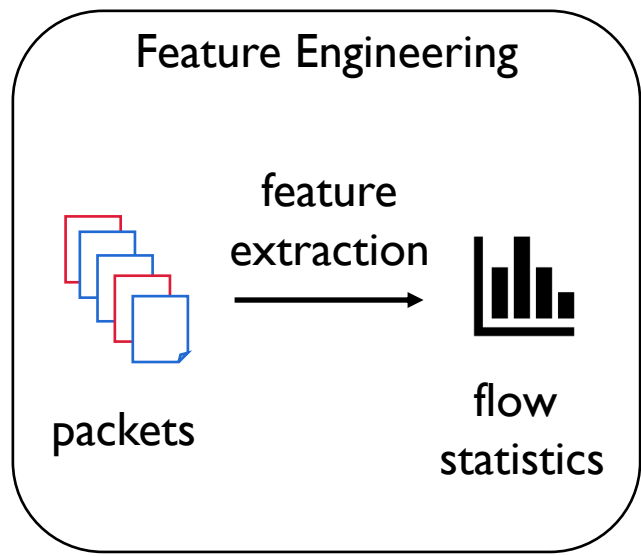
Match    Action    Register

Protocol-Independent Switch Architecture (PISA)

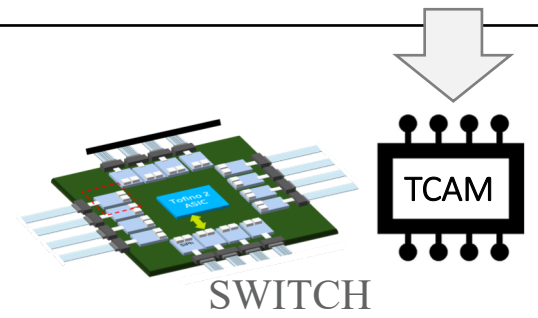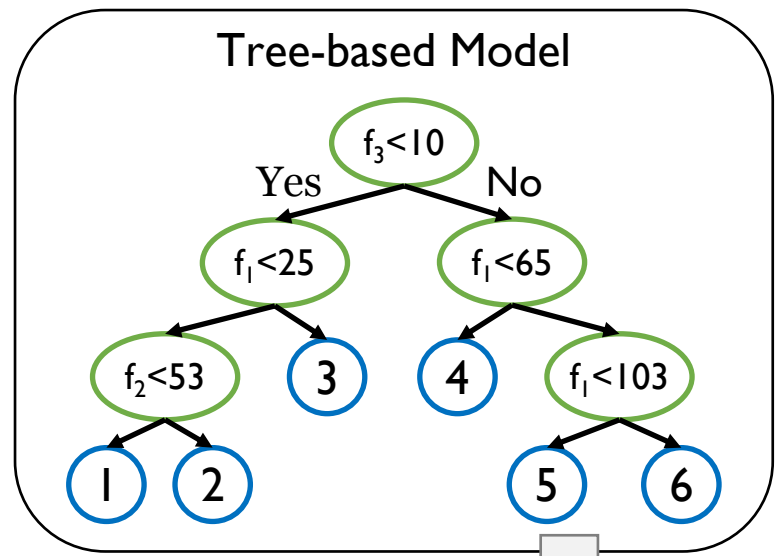**Constraints on ML models**

1. **Computation Constraints (simple OPs, ...)**

2. **Storage Constraints (once register access, ...)**

[2] Guangmeng Zhou, Zhuotao Liu, Chuanpu Fu, Qi Li, and Ke Xu. An Efficient Design of Intelligent Network Data Plane. In USENIX Security Symposium (USENIX Security), 2023.

# Prior traffic analysis art targeting Intelligent Network Data Plane


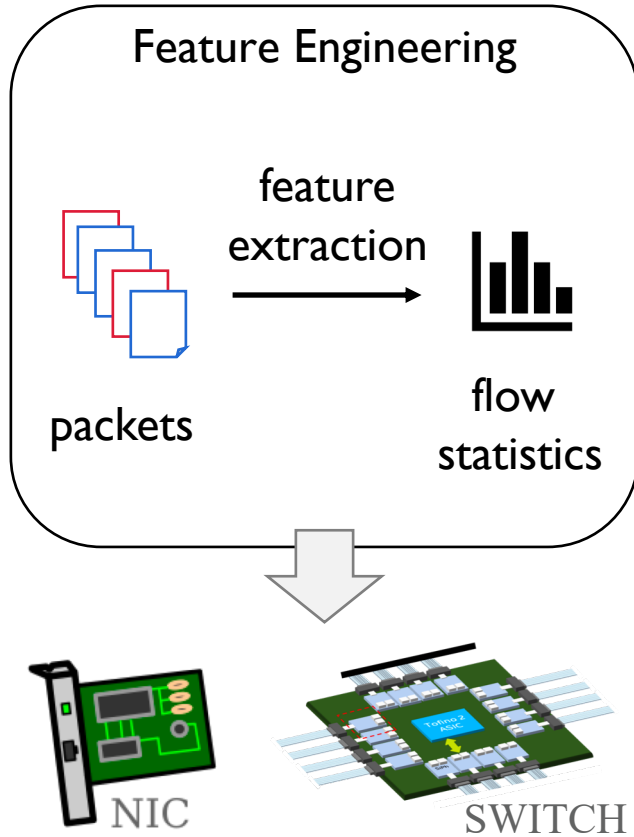
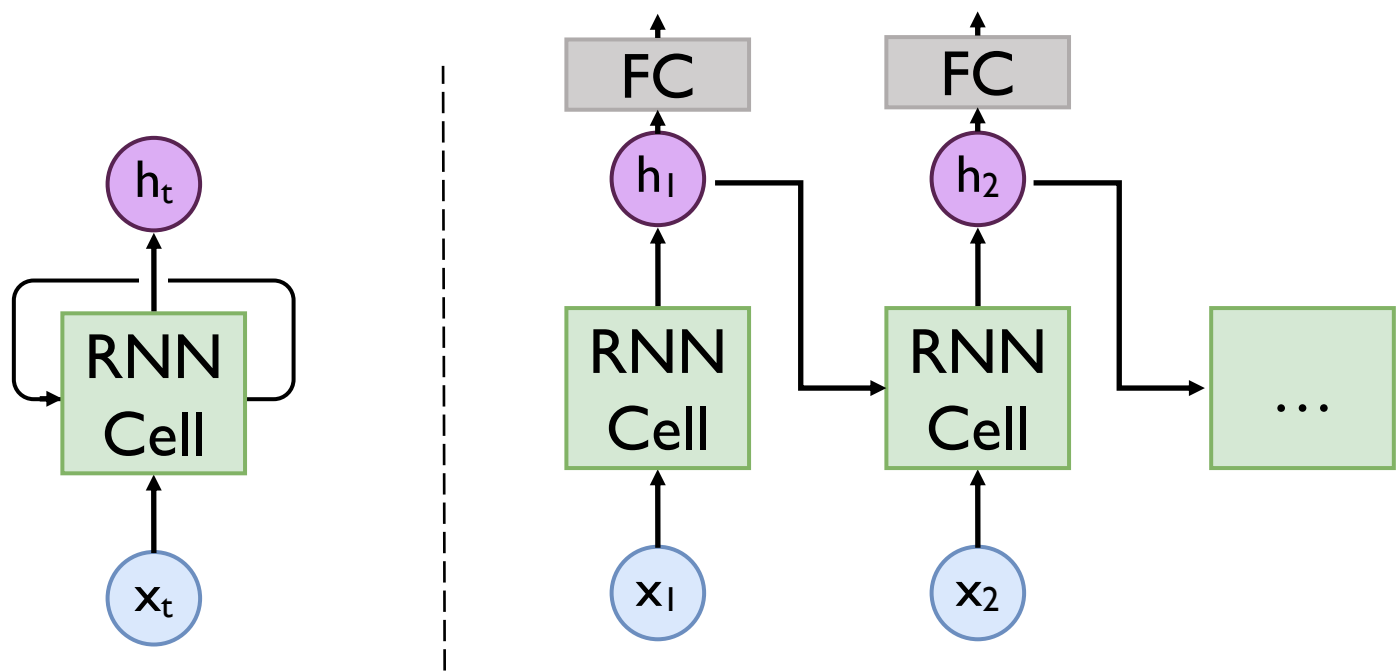Neural Network on the NIC[1]
(NSDI'22)

NetBeacon[2]
(Security'23)

Their models rely on advanced feature engineering to boost accuracy



Fundamental Limitations:

- Critical features are impossible / difficult to compute

- Handling dynamic features as a flow proceeds

- Overheads for computing and storing statistical features

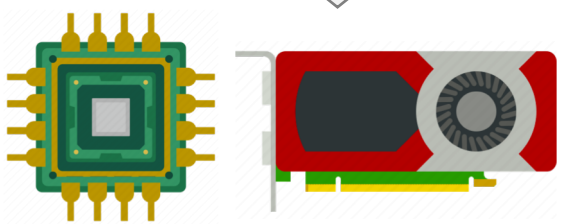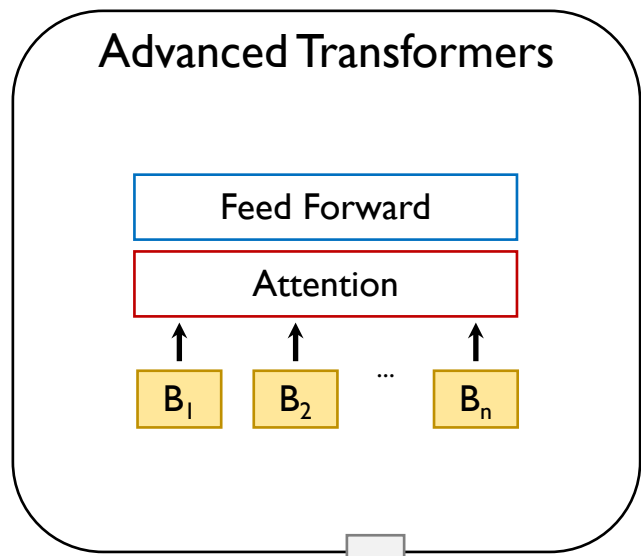- Handcrafted feature engineering and overfitting concerns

## #1 Advance INDP to models that are not limited by the availability of flow features



X: Packet Length, Inter Packet Delay …

- Recurrent computation on raw packet metadata

- Without statistical feature engineering

- Output latest inference result for each packet

# Limited model accuracy on Network Data Plane



**Advanced Transformers**

Feed Forward

Attention

$B_1$ $B_2$ ... $B_n$

Severe Accuracy Degradation

**Binary MLP**

+1 $-1$ Sign $a_1^1$

+1 $+1$ Sign $a_2^1$

$-1$ $-1$

**Tree-based Model**

$f_3 < 10$

Yes       No

$f_1 < 25$       $f_1 < 65$

$f_2 < 53$   3   4   $f_1 < 103$

1   2       5   6

SRAM

NIC

TCAM

SWITCH

Neural Network on the NIC[1]
(NSDI'22)

NetBeacon[2]
(Security'23)

## #2 Complement the on-switch RNN with an off-switch Transformer-based module

inbound

**Traffic Flows**

On-switch: the vase majority of traffic

SWITCH

Co-design

batch processing

Off-switch: boost the overall accuracy

## Challenge 1: implement RNN inference on programmable switch

Recurrent Computation Scheme in RNN

Match-Action Paradigm in PISA



- **Complex calculations** in each RNN time step (multiplications, non-linear functions …)

- **Store and retrieve the hidden states** through RNN time steps

- Simple operations (add, XNOR, shift, …), limited stages

- Each register can only be accessed once, limited storage

Challenge 2: accurately identify the flows for escalation and analyze these flows online



How to classify the vast majority of traffic on-switch and identify the flows with insufficient classification confidence accurately?

How to construct an appropriate system to analyze the escalated flows with a Transformer-based model online?

**BoS** is a hybrid traffic analysis system with the co-design of:

- An on-switch RNN,
- An off-switch Integrated Model Inference System
- A carefully designed flow escalation mechanism



Realize complex RNN computations using a set of novel data plane native operations

Construct an Integrated Model Inference System for fast online traffic analysis with Transformer

Design an escalation mechanism to accurately identify the flows with insufficient confidence from on-switch analysis



Realize complex RNN computations using a set of novel data plane native operations

Construct an Integrated Model Inference System for fast online traffic analysis with Transformer

**Binary RNN**

| length seq | IPD seq |
| Embedding | Embedding |
| *STE* | *STE* |
| FC |
| *STE* |
| GRU | Memory |
| *STE* |
| FC + Softmax |

**Input:**
packet length sequence, Inter-Packet-Delay sequence

**Feature Embedding:**
length embedding || IPD embedding -> Fully Connected

**RNN Cell:**
Gated Recurrent Unit

**Output Layer:**
Fully Connected + Softmax

**Binary RNN**

| length seq | IPD seq |
|:---:|:---:|
| Embedding | Embedding |
| *STE* | *STE* |

| FC |
|:---:|
| *STE* |

| GRU | Memory |
|:---:|:---:|
| *STE* | |

| FC + Softmax |
|:---:|

- **Binary activations & full-precision model weights**
  **Better accuracy than full model binarization**



**Straight-Through Estimator**

**Binary RNN**

| length seq | IPD seq |
|:---:|:---:|

| Embedding | Embedding |
|:---:|:---:|
| *STE* | *STE* |

| FC |
|:---:|
| *STE* |

| GRU | Memory |
|:---:|:---:|
| *STE* | |

| FC + Softmax |
|:---:|

- **Forward propagation based on match-action table lookup**



input-output mapping

| GRU Layer | | |
|:---:|:---:|:---:|
| Input: $x_t$ (6bit) | Input: $h_{t-1}$ (6bit) | Output: $h_t$ (6bit) |
| 000000 | 000000 | 001000 |
| 000000 | 000001 | 000001 |
| ... | ... | ... |
| 111111 | 111110 | 010111 |
| 111111 | 111111 | 111110 |

**Binary RNN**

| length seq | IPD seq |
|---|---|

Embedding **STE** → Embedding **STE**

FC **STE**

GRU **STE** ⟷ Memory

FC + Softmax

- **Expand RNN time steps in serial stages**

① Read the previous hidden state

$ev_t$

② Perform layer forward propagation

GRU **STE** ⟷ Memory

$h_{t-1}$

$h_t$

$h_t$

③ Update the hidden state

☹ ① ② ③ **cannot be realized in one stage**

**Binary RNN**



- **Expand RNN time steps in serial stages**

When a packet arrives, we use the latest $S$ embedding vectors to get an intermediate result.



① A window of $S$ packets

Flow A

packet length

IPD

STE — Straight-through Estimator

Full-precision Layer

Embed
STE

STE
Embed

FC STE

ev

STE
GRU

STE
GRU

STE
GRU

ev₁ ev₂ ●●● evₛ

② Expanded $S$ RNN Time Steps

Output Layer

(0.9, 0.1, 0.0, 0.0, 0.0)

③ An Intermediate Result (A probability vector)

As the flow proceeds, we shift the window by one packet to processing a new segment of embedding vectors repeatedly, which produces many intermediate results.

For the latest packet, we accumulate all previous intermediate results, and select the class with the largest cumulative probability as the final result.



④ Final Result: Class 0

Argmax

(3.0, 1.3, 0.2, 0.1, 0.4)

Accumulate

**STE**    Straight-through Estimator

Full-precision Layer

(0.9, 0.1, 0.0, 0.0, 0.0)
(0.5, 0.2, 0.1, 0.1, 0.1)
...
(0.1, 0.7, 0.0, 0.0, 0.2)

① A window of $S$ packets

packet length

Embed
STE

STE
Embed

IPD

FC    STE

ev

STE
GRU

STE
GRU

STE
GRU

Output Layer
STE

ev$_1$   ev$_2$   ●●●   ev$_S$

Flow A

③ An Intermediate Result (A probability vector)

② Expanded $S$ RNN Time Steps

Embrace advanced models for corner cases with insufficient classification confidence



- Identify the flows with ambiguity from the on-switch analysis by two thresholds

- Forward these flows for escalated analysis using advanced models

① Whether a packet is ambiguous is determined by the Confidence Threshold

② Whether a flow should be escalated is determined by the number of ambiguous packets in the flow, using the Escalation Threshold

Losses for accurately identifying the flows with insufficient confidence



Improve the model's ability to predict the ground-truth class

$$\mathcal{L}_1 = \boxed{-(1 - p_y)^\gamma \log(p_y)} \boxed{- \lambda \sum_{i \neq y} p_i^\gamma \log(1 - p_i)}$$

$$\mathcal{L}_2 = \boxed{-(1 - p_y)^\gamma \log(p_y)} \boxed{- \lambda \, p_{\text{false}}^\gamma \log(1 - p_{\text{false}})}$$

Negate the model's prediction on all non-ground-truth classes / the one with largest probability

24

Enable fast online inference for escalated flows using Transformer-based model



- Non-blocking processing pipeline

- Single-threaded, stateful tasks

Server A
(Packet Generation)

Programmable Switch
(RNN Analysis)

Server B
(IMIS Processing)

## Metrics

- Packet-level macro-accuracy
- SRAM and TCAM consumptions

## Baselines

- Neural Network on the NIC [1]
- NetBeacon [2]

## Tasks

- Encrypted Traffic Classification on VPN
- Botnet Traffic Classification on IoT
- Behavioral Analysis of IoT Devices
- P2P Application Fingerprinting

Table 3: Analysis accuracy for BoS and other two closely related art.

| Methods | BoS | | | NetBeacon [71] (Tree-based Models) | | | N3IC [51] (Binary MLP) | | |
|---|---|---|---|---|---|---|---|---|---|
| Network Load | Low | Normal | High | Low | Normal | High | Low | Normal | High |
| Encrypted Traffic Classification on VPN (ISCXVPN2016) | | | | | | | | | |
| Email | 0.935 / 0.933 | 0.936 / 0.925 | 0.933 / 0.923 | 0.309 / 0.514 | 0.315 / 0.524 | 0.320 / 0.525 | 0.347 / 0.326 | 0.354 / 0.339 | 0.367 / 0.350 |
| Chat | 0.903 / 0.818 | 0.902 / 0.818 | 0.901 / 0.814 | 0.739 / 0.935 | 0.739 / 0.933 | 0.742 / 0.925 | 0.336 / 0.655 | 0.336 / 0.654 | 0.342 / 0.656 |
| Streaming | 0.926 / 0.941 | 0.926 / 0.939 | 0.926 / 0.910 | 0.963 / 0.919 | 0.962 / 0.904 | 0.962 / 0.874 | 0.741 / 0.608 | 0.742 / 0.603 | 0.743 / 0.581 |
| FTP | 0.973 / 0.928 | 0.973 / 0.926 | 0.973 / 0.922 | 0.946 / 0.659 | 0.946 / 0.655 | 0.947 / 0.654 | 0.563 / 0.396 | 0.567 / 0.396 | 0.575 / 0.397 |
| VoIP | 0.968 / 0.958 | 0.968 / 0.958 | 0.968 / 0.957 | 0.938 / 0.882 | 0.939 / 0.881 | 0.939 / 0.882 | 0.883 / 0.783 | 0.884 / 0.782 | 0.886 / 0.787 |
| P2P | 0.905 / 0.927 | 0.903 / 0.928 | 0.876 / 0.930 | 0.810 / 0.959 | 0.798 / 0.959 | 0.778 / 0.960 | 0.578 / 0.739 | 0.577 / 0.742 | 0.565 / 0.748 |
| Macro-F1 | 0.926 | 0.925 | 0.919 | 0.786 | 0.784 | 0.780 | 0.565 | 0.567 | 0.568 |
| Botnet Traffic Classification on IoT (BOTIOT) | | | | | | | | | |
| Data Exfiltration | 0.964 / 0.974 | 0.951 / 0.973 | 0.899 / 0.971 | 0.691 / 0.845 | 0.684 / 0.847 | 0.658 / 0.848 | 0.514 / 0.879 | 0.508 / 0.881 | 0.506 / 0.879 |
| Key Logging | 0.960 / 0.946 | 0.961 / 0.962 | 0.959 / 0.902 | 0.921 / 0.425 | 0.921 / 0.419 | 0.918 / 0.399 | 0.055 / 0.033 | 0.058 / 0.033 | 0.052 / 0.031 |
| OS Scan | 0.996 / 0.996 | 0.995 / 0.989 | 0.995 / 0.966 | 0.838 / 0.963 | 0.841 / 0.963 | 0.844 / 0.945 | 0.831 / 0.693 | 0.830 / 0.677 | 0.831 / 0.672 |
| Service Scan | 0.993 / 0.992 | 0.986 / 0.973 | 0.979 / 0.978 | 0.928 / 0.876 | 0.927 / 0.870 | 0.917 / 0.858 | 0.845 / 0.663 | 0.830 / 0.664 | 0.840 / 0.663 |
| Macro-F1 | 0.978 | 0.974 | 0.955 | 0.785 | 0.782 | 0.769 | 0.547 | 0.542 | 0.541 |
| Behavioral Analysis of IoT Devices (CICIOT2022) | | | | | | | | | |
| Power | 0.926 / 0.887 | 0.924 / 0.882 | 0.921 / 0.882 | 0.819 / 0.726 | 0.820 / 0.724 | 0.817 / 0.724 | 0.639 / 0.750 | 0.640 / 0.750 | 0.640 / 0.748 |
| Idle | 0.922 / 0.943 | 0.921 / 0.942 | 0.918 / 0.941 | 0.810 / 0.938 | 0.808 / 0.938 | 0.806 / 0.936 | 0.618 / 0.640 | 0.620 / 0.642 | 0.622 / 0.646 |
| Interact | 0.934 / 0.946 | 0.934 / 0.948 | 0.934 / 0.943 | 0.871 / 0.786 | 0.873 / 0.786 | 0.872 / 0.784 | 0.651 / 0.504 | 0.655 / 0.506 | 0.661 / 0.510 |
| Macro-F1 | 0.926 | 0.925 | 0.923 | 0.822 | 0.821 | 0.820 | 0.629 | 0.631 | 0.633 |
| P2P Application Fingerprinting (PeerRush) | | | | | | | | | |
| eMule | 0.943 / 0.949 | 0.918 / 0.949 | 0.898 / 0.950 | 0.846 / 0.954 | 0.821 / 0.955 | 0.805 / 0.954 | 0.734 / 0.866 | 0.730 / 0.867 | 0.723 / 0.875 |
| uTorrent | 0.949 / 0.924 | 0.950 / 0.912 | 0.941 / 0.894 | 0.882 / 0.870 | 0.885 / 0.858 | 0.885 / 0.831 | 0.734 / 0.789 | 0.735 / 0.790 | 0.738 / 0.783 |
| Vuze | 0.946 / 0.962 | 0.945 / 0.947 | 0.941 / 0.930 | 0.910 / 0.810 | 0.907 / 0.790 | 0.904 / 0.793 | 0.821 / 0.626 | 0.826 / 0.622 | 0.826 / 0.616 |
| Macro-F1 | 0.945 | 0.937 | 0.925 | 0.877 | 0.866 | 0.858 | 0.755 | 0.755 | 0.752 |

Across 4 tasks, **BoS** achieves an average F1-score improvement of 0.13 and 0.31 than NetBeacon and N3IC.

Table 3: Analysis accuracy for BoS and other two closely related art.

| Methods | BoS | | | NetBeacon [71] (Tree-based Models) | | | N3IC [51] (Binary MLP) | | |
|---|---|---|---|---|---|---|---|---|---|
| Network Load | Low | Normal | High | Low | Normal | High | Low | Normal | High |
| Encrypted Traffic Classification on VPN (ISCXVPN2016) | | | | | | | | | |
| Email | 0.935 / 0.933 | 0.936 / 0.925 | 0.933 / 0.923 | 0.309 / 0.514 | 0.315 / 0.524 | 0.320 / 0.525 | 0.347 / 0.326 | 0.354 / 0.339 | 0.367 / 0.350 |
| Chat | 0.903 / 0.818 | 0.902 / 0.818 | 0.901 / 0.814 | 0.739 / 0.935 | 0.739 / 0.933 | 0.742 / 0.925 | 0.336 / 0.655 | 0.336 / 0.654 | 0.342 / 0.656 |
| Streaming | 0.926 / 0.941 | 0.926 / 0.939 | 0.926 / 0.910 | 0.963 / 0.919 | 0.962 / 0.904 | 0.962 / 0.874 | 0.741 / 0.608 | 0.742 / 0.603 | 0.743 / 0.581 |
| FTP | 0.973 / 0.928 | 0.973 / 0.926 | 0.973 / 0.922 | 0.946 / 0.659 | 0.946 / 0.655 | 0.947 / 0.654 | 0.563 / 0.396 | 0.567 / 0.396 | 0.575 / 0.397 |
| VoIP | 0.968 / 0.958 | 0.968 / 0.958 | 0.968 / 0.957 | 0.938 / 0.882 | 0.939 / 0.881 | 0.939 / 0.882 | 0.883 / 0.783 | 0.884 / 0.782 | 0.886 / 0.787 |
| P2P | 0.905 / 0.927 | 0.903 / 0.928 | 0.876 / 0.930 | 0.810 / 0.959 | 0.798 / 0.959 | 0.778 / 0.960 | 0.578 / 0.739 | 0.577 / 0.742 | 0.565 / 0.748 |
| Macro-F1 | 0.926 | 0.925 | 0.919 | 0.786 | 0.784 | 0.780 | 0.565 | 0.567 | 0.568 |
| Botnet Traffic Classification on IoT (BOTIOT) | | | | | | | | | |
| Data Exfiltration | 0.964 / 0.974 | 0.951 / 0.973 | 0.899 / 0.971 | 0.691 / 0.845 | 0.684 / 0.847 | 0.658 / 0.848 | 0.514 / 0.879 | 0.508 / 0.881 | 0.506 / 0.879 |
| Key Logging | 0.960 / 0.946 | 0.961 / 0.962 | 0.959 / 0.902 | 0.921 / 0.425 | 0.921 / 0.419 | 0.918 / 0.399 | 0.055 / 0.033 | 0.058 / 0.033 | 0.052 / 0.031 |
| OS Scan | 0.996 / 0.996 | 0.995 / 0.989 | 0.995 / 0.966 | 0.838 / 0.963 | 0.841 / 0.963 | 0.844 / 0.945 | 0.831 / 0.693 | 0.830 / 0.677 | 0.831 / 0.672 |
| Service Scan | 0.993 / 0.992 | 0.986 / 0.973 | 0.979 / 0.978 | 0.928 / 0.876 | 0.927 / 0.870 | 0.917 / 0.858 | 0.845 / 0.663 | 0.830 / 0.664 | 0.840 / 0.663 |
| Macro-F1 | 0.978 | 0.974 | 0.955 | 0.785 | 0.782 | 0.769 | 0.547 | 0.542 | 0.541 |
| Behavioral Analysis of IoT Devices (CICIOT2022) | | | | | | | | | |
| Power | 0.926 / 0.887 | 0.924 / 0.882 | 0.921 / 0.882 | 0.819 / 0.726 | 0.820 / 0.724 | 0.817 / 0.724 | 0.639 / 0.750 | 0.640 / 0.750 | 0.640 / 0.748 |
| Idle | 0.922 / 0.943 | 0.921 / 0.942 | 0.918 / 0.941 | 0.810 / 0.938 | 0.808 / 0.938 | 0.806 / 0.936 | 0.618 / 0.640 | 0.620 / 0.642 | 0.622 / 0.646 |
| Interact | 0.934 / 0.946 | 0.934 / 0.948 | 0.934 / 0.943 | 0.871 / 0.786 | 0.873 / 0.786 | 0.872 / 0.784 | 0.651 / 0.504 | 0.655 / 0.506 | 0.661 / 0.510 |
| Macro-F1 | 0.926 | 0.925 | 0.923 | 0.822 | 0.821 | 0.820 | 0.629 | 0.631 | 0.633 |
| P2P Application Fingerprinting (PeerRush) | | | | | | | | | |
| eMule | 0.943 / 0.949 | 0.918 / 0.949 | 0.898 / 0.950 | 0.846 / 0.954 | 0.821 / 0.955 | 0.805 / 0.954 | 0.734 / 0.866 | 0.730 / 0.867 | 0.723 / 0.875 |
| uTorrent | 0.949 / 0.924 | 0.950 / 0.912 | 0.941 / 0.894 | 0.882 / 0.870 | 0.885 / 0.858 | 0.885 / 0.831 | 0.734 / 0.789 | 0.735 / 0.790 | 0.738 / 0.783 |
| Vuze | 0.946 / 0.962 | 0.945 / 0.947 | 0.941 / 0.930 | 0.910 / 0.810 | 0.907 / 0.790 | 0.904 / 0.793 | 0.821 / 0.626 | 0.826 / 0.622 | 0.826 / 0.616 |
| Macro-F1 | 0.945 | 0.937 | 0.925 | 0.877 | 0.866 | 0.858 | 0.755 | 0.755 | 0.752 |

On more challenging tasks with more classes, the improvement is even greater, up to 0.17 and 0.39.

Table 4: Hardware resource utilization.

| Datasets (Tasks) | | ISCXVPN 2016 | BOT IOT | CICIOT 2022 | Peer Rush |
|---|---|---|---|---|---|
| SRAM | Flow Info. (stateful) | 5.21% | 5.21% | 5.21% | 5.21% |
| | EV (stateful) | 3.65% | 3.65% | 3.65% | 3.65% |
| | CPR (stateful) | 5.63% | 3.75% | 2.81% | 2.81% |
| | FE (stateless) | 2.19% | 2.19% | 2.19% | 2.19% |
| | GRU (stateless) | 3.02% | 1.56% | 0.73% | 0.73% |
| | Total⋆ | 23.44% | 20.10% | 18.33% | 18.33% |
| TCAM | Argmax (Total) | 1.74% | 1.04% | 0.69% | 0.69% |

⋆ Including other components not listed, *e.g.,* packet counters for each flow.

- **BoS** uses 23.44%/20.10%/18.33%/18.33% of SRAM in 4 tasks, respectively. (Similar size to NetBeacon)

- **BoS** uses 1.74%/1.04%/0.69%/0.69% of TCAM in 4 tasks, respectively. (20x less than NetBeacon)

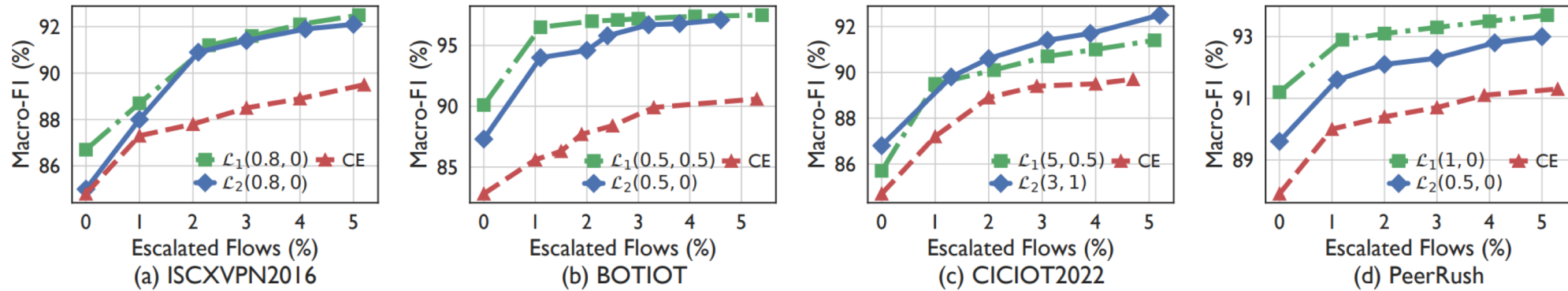- Efficiency of Analysis Escalation & System Performance of IMIS



Figure 9: [Testbed] The trade-off between percentage of escalated flows and the overall accuracy.

- **BoS** effectively accommodates the off-switch analysis model to compensate for on-switch analysis.

- Our losses achieve better trade-off between the amount of escalated flows and the overall accuracy.

- Efficiency of Analysis Escalation & System Performance of IMIS
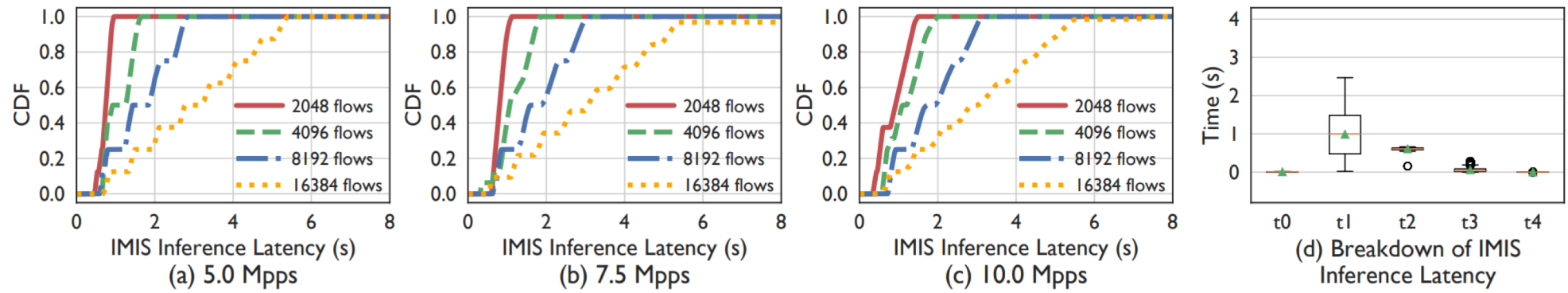


Figure 10: [Testbed] The inference throughput and latency of the off-switch IMIS.

When the number of concurrent flows is below 4096, the maximum end-to-end latency imposed by **IMIS** is less than 2 seconds even for 10.0 Mpps inbound rate (equivalently 41 Gbps as the packet sizes we send are 512 B, and **BoS** typically escalates less than 5% of flows to **IMIS**).

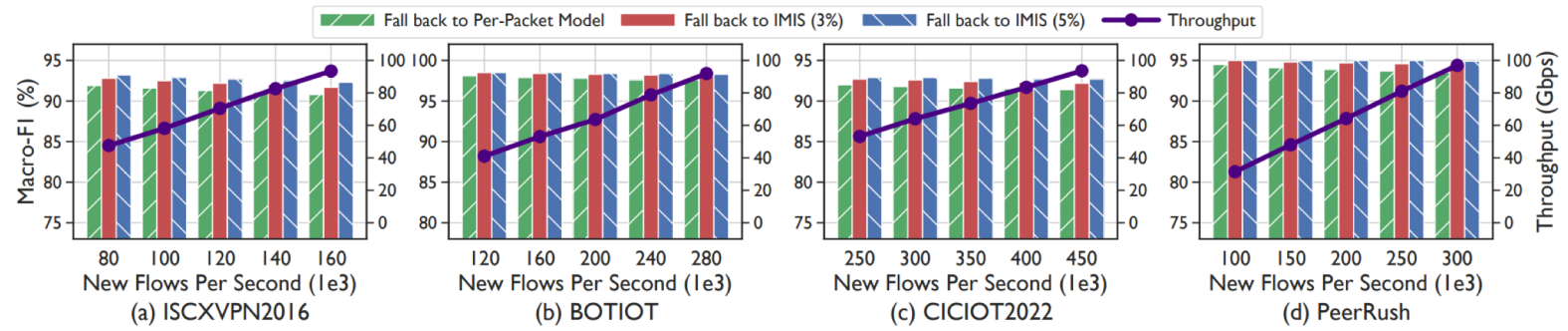- Scaling Test of the Entire System



Figure 11: [Testbed] Scaling test of BoS when we progressively increase the aggregate throughput to 100 Gbps.
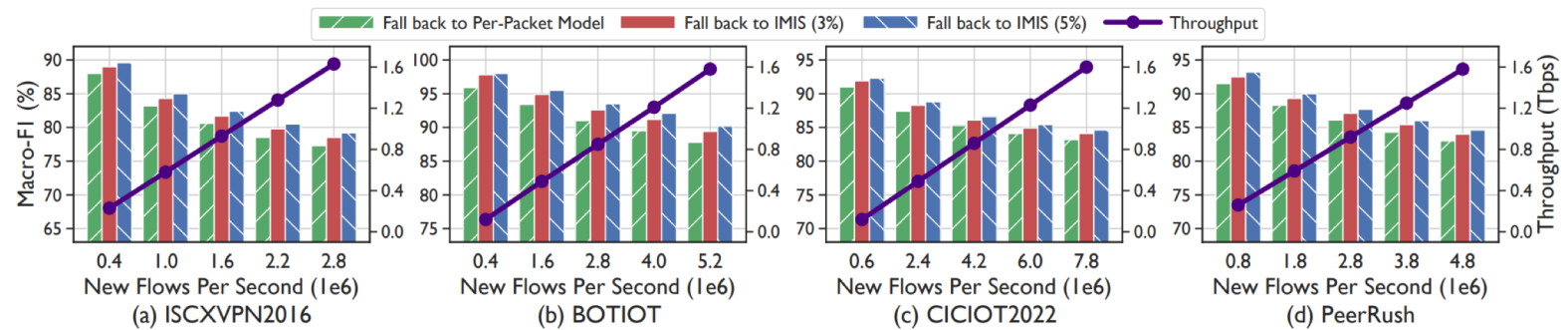


Figure 12: [Simulation] Scaling test of BoS when we progressively increase the aggregate throughput to 1.6 Tbps.

The macro-F1 scores of **BoS** remain nearly identical as the throughput achieves 100Gbps, and reveal a sublinear decline as the throughput achieves 1.6Tbps.

- **BoS** is an online traffic analysis system, which is powered by the co-design of an on-switch RNN, an off-switch Integrated Model Inference System, and a carefully designed flow escalation mechanism.
- As a result, **BoS** can process over 95% of flows with the on-switch RNN accurately, and escalate the remainning ambiguous flows to the off-switch IMIS, outperforming prior works in accuracy, scalability and hardware resource utilization.

Source code: https://github.com/InspiringGroup-Lab/Brain-on-Switch

Homepage of our group: https://inspiringgroup.github.io/

# Brain-on-Switch:
## Towards Advanced Intelligent Network Data Plane via NN-Driven Traffic Analysis at Line-Speed

Thanks! Questions?

yanjz22@mails.tsinghua.edu.cn