**Hong Ye He (presenter)**, **Zhi Guo Yang, Xiang Ning Chen**
Zhongxing Telecommunication Equipment (ZTE) Corporation

**Session 8:** Security in industrial applications

**Paper S8.1**

ITU KALEIDOSCOPE
ONLINE2020

# 1. Traffic Identification / Classification
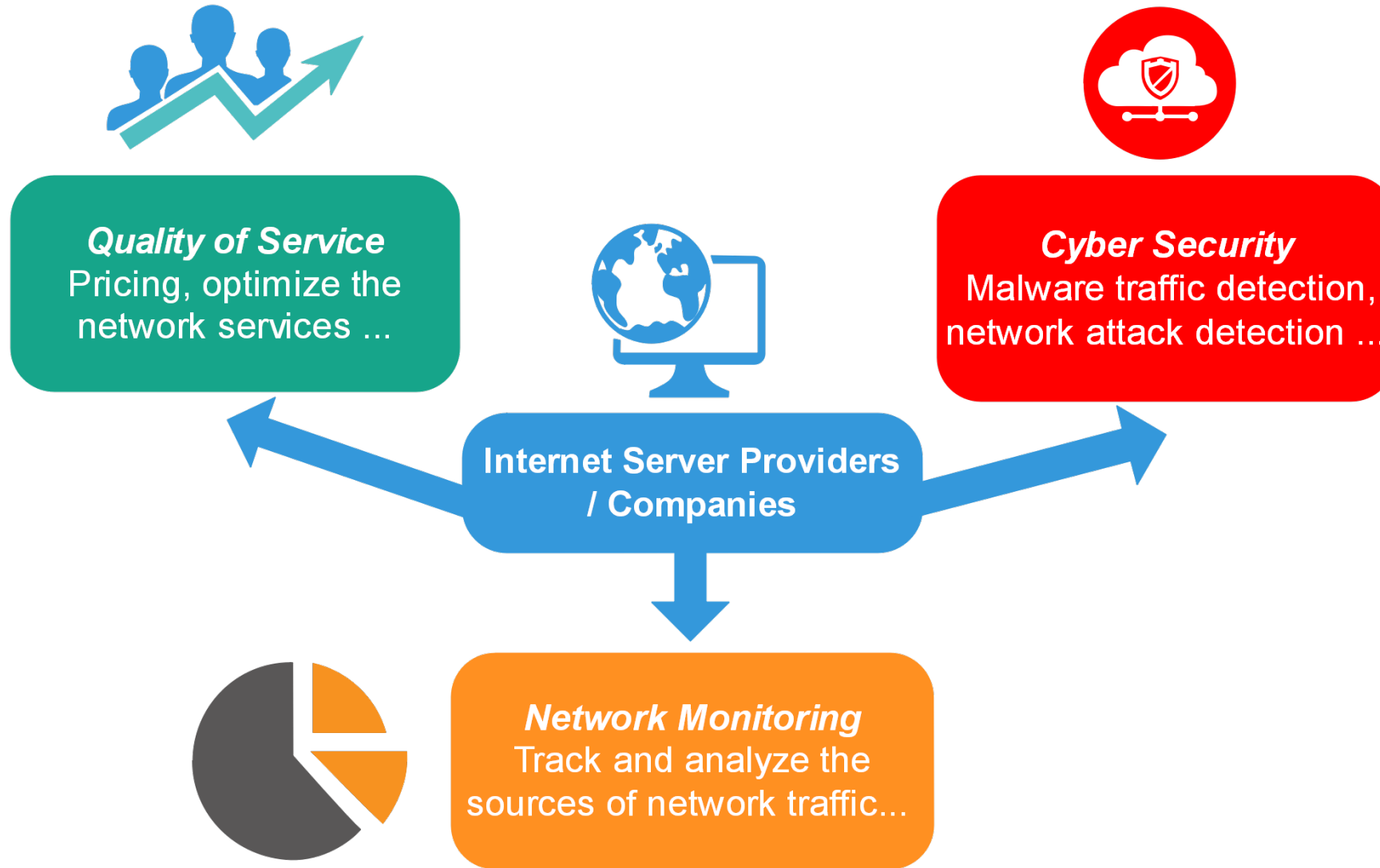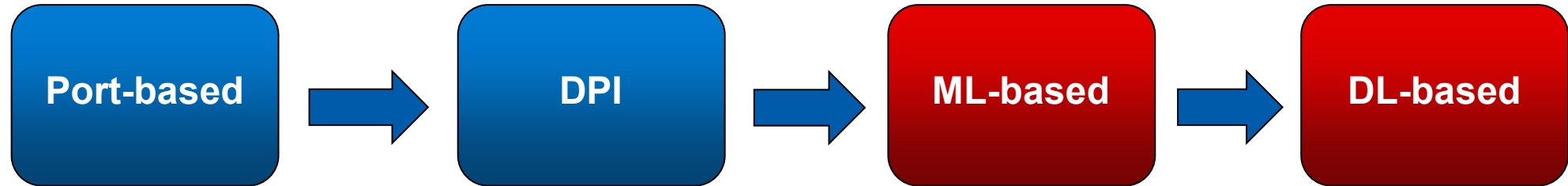
# 2. Traffic Identification - Methods

Port-based → DPI → ML-based → DL-based

- The **port-based** and **deep packet inspection** methods that locate fixed patterns from traffic data.

- Rule-based methods that rely on unencrypted information. Not suitable for encrypted traffic.
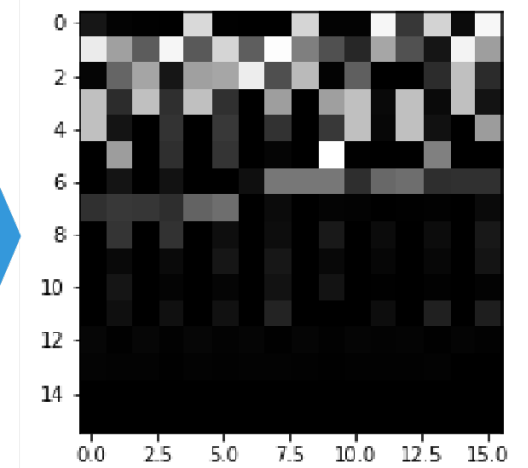
- The **machine learning** methods that extract hand-designed features from traffic.

- The **deep learning** methods that perform representation learning on raw traffic bytes.

- Extract common traffic features. Ideal for encrypted traffic identification.

ITU KALEIDOSCOPE
ONLINE2020

# 3. Deep Learning Based Method - Image Processing

- Current popular DL-based method transform the **raw payload bytes** of traffic packets / flows to grayscale images.

- The purpose is to introduce **image processing** with neural network such as CNN.

- Thus, classification effectiveness is decided by the **representation learning** capacity of the network.
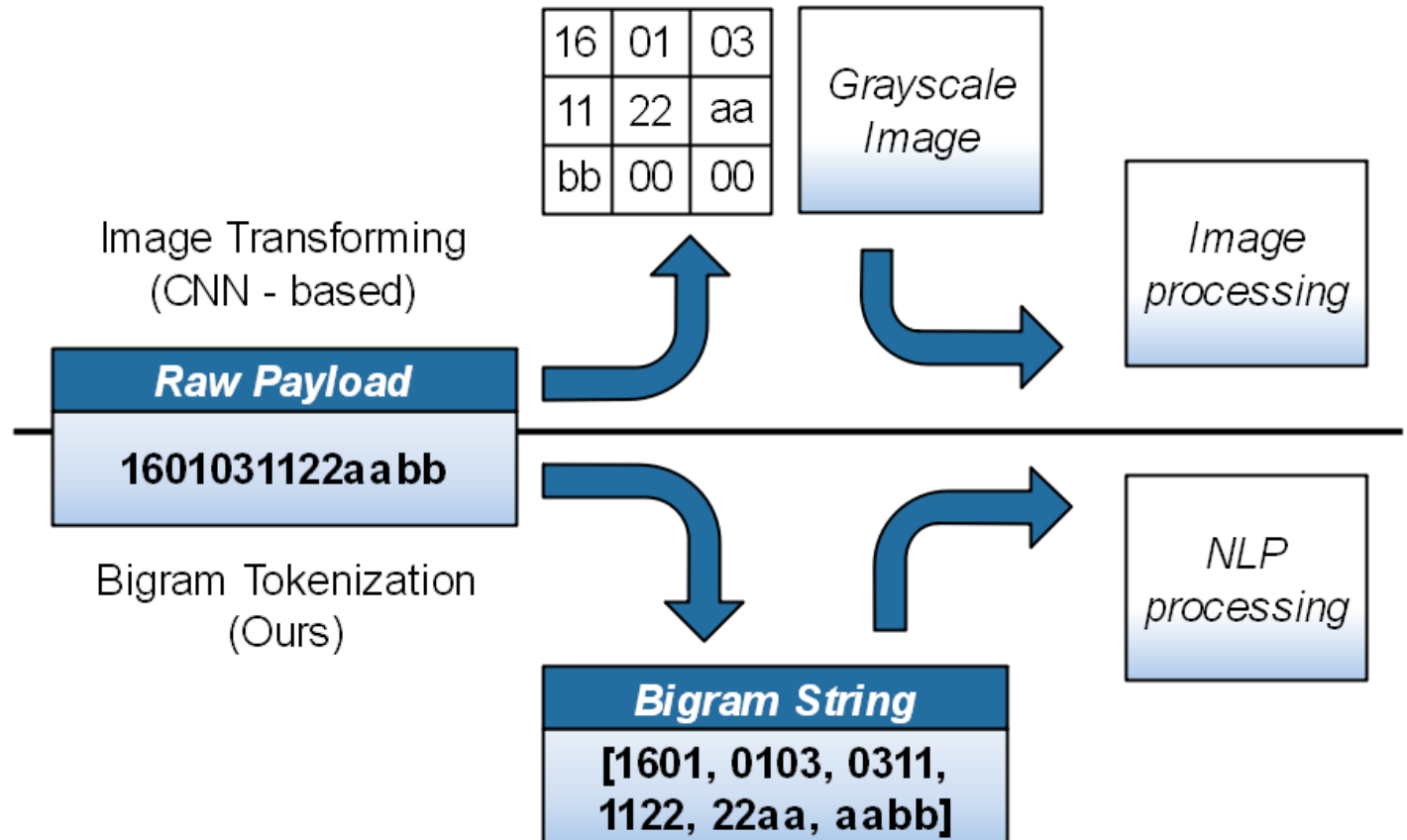
# 4. Introducing NLP Processing

- We perform **bigram tokenization** on encrypted traffic bytes to generate payload bigram strings.

- The traffic identification is transformed to a **NLP classification** task.

- **NLP-related representation learning** can be directly applied to the traffic data.

# 5. Payload Encoding Representation from Transformers (PERT)

**Dynamic Word Embeddbing**

**Encoder from Transformer**



- A Bidirectional Encoder Representations from Transformers (BERT) like structure to apply **NLP representation learning** on raw traffic.

# 6. PERT - Pretraining

- **Language models (LM)** aim to predict words using their contextual inputs.

- LM is originally designed for language generator. But it can also be applied to **initialize NLP encoding network.**

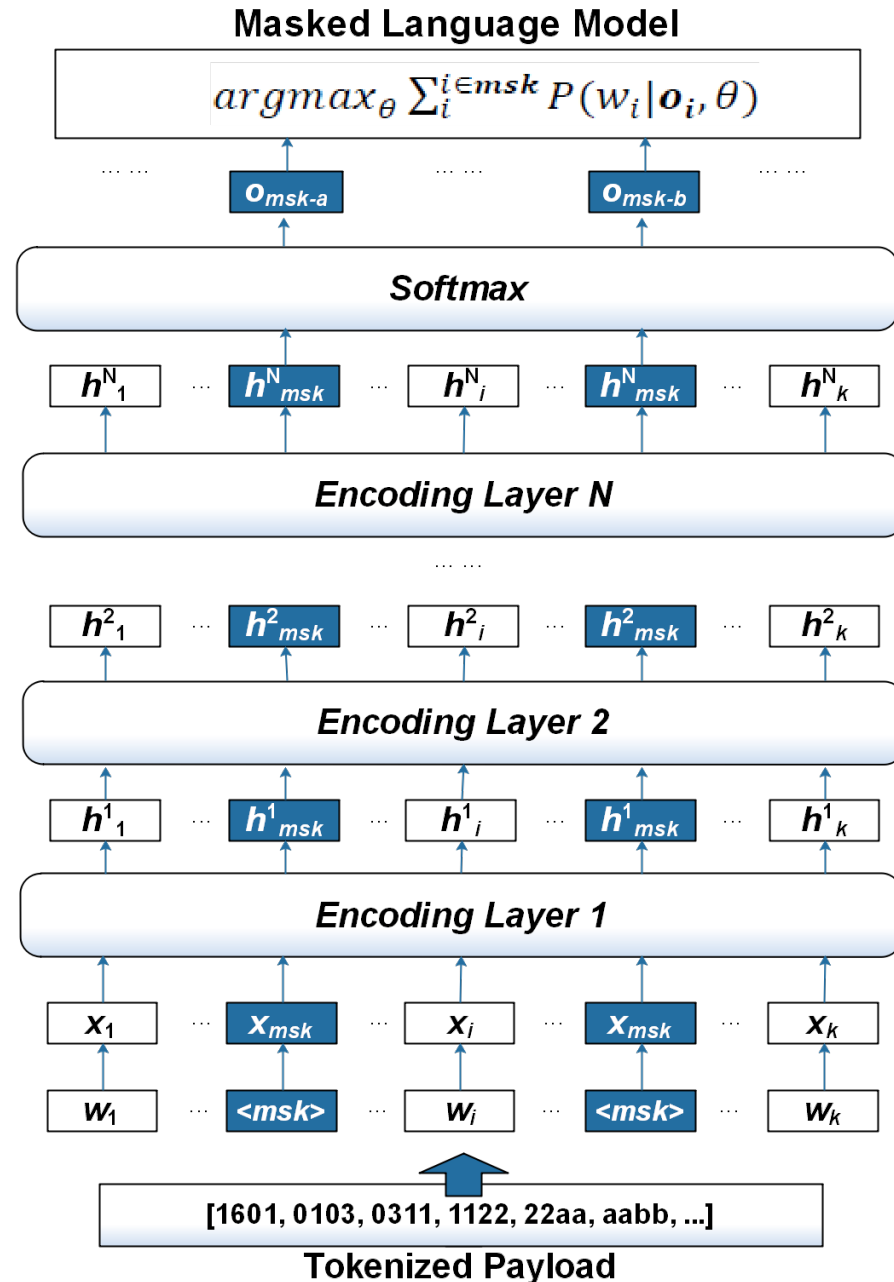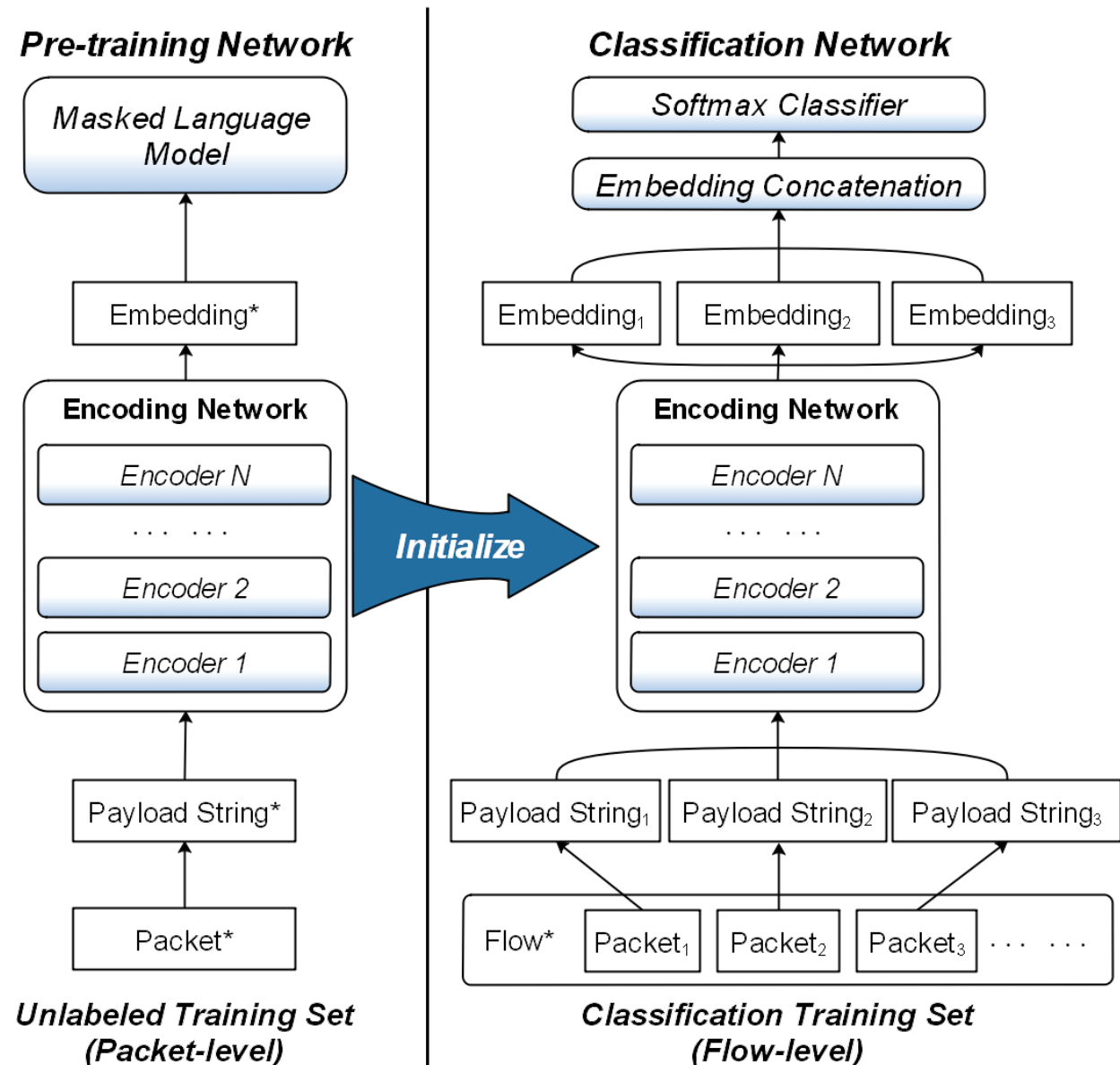- In a BERT-like network, the **masked language model (MLM)** is frequently utilized for initialization.

**Masked Language Model**

$$argmax_\theta \sum_i^{i \in msk} P(w_i | \boldsymbol{o_i}, \theta)$$

$\boldsymbol{o}_{msk\text{-}a}$   $\boldsymbol{o}_{msk\text{-}b}$

*Softmax*

$h^N_1$   $h^N_{msk}$   $h^N_i$   $h^N_{msk}$   $h^N_k$

*Encoding Layer N*

$h^2_1$   $h^2_{msk}$   $h^2_i$   $h^2_{msk}$   $h^2_k$

*Encoding Layer 2*

$h^1_1$   $h^1_{msk}$   $h^1_i$   $h^1_{msk}$   $h^1_k$

*Encoding Layer 1*

$x_1$   $x_{msk}$   $x_i$   $x_{msk}$   $x_k$

$w_1$   *<msk>*   $w_i$   *<msk>*   $w_k$

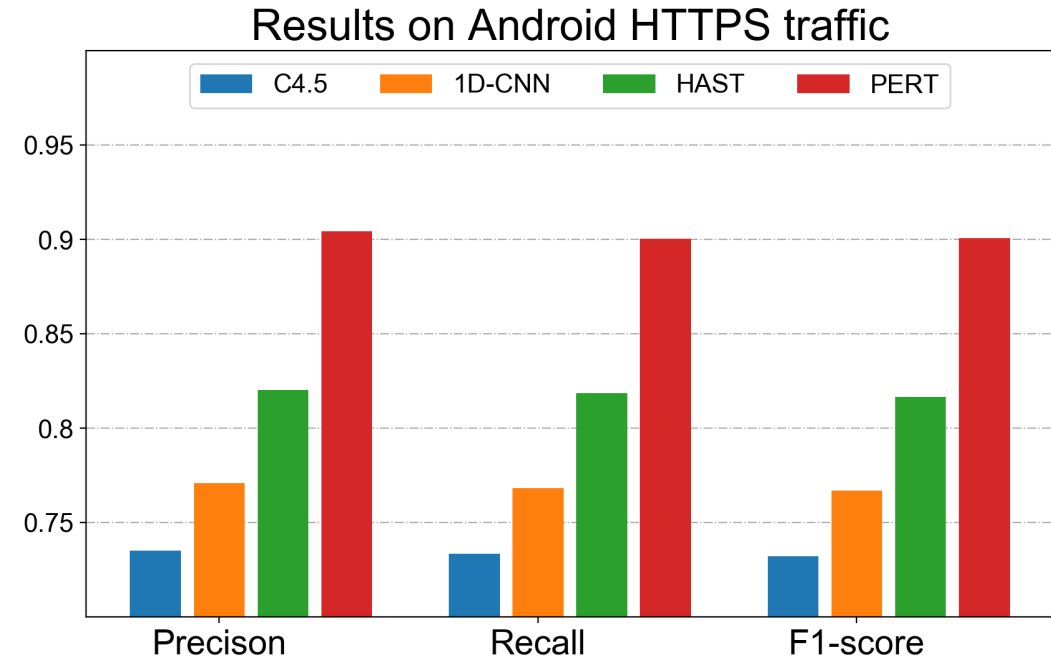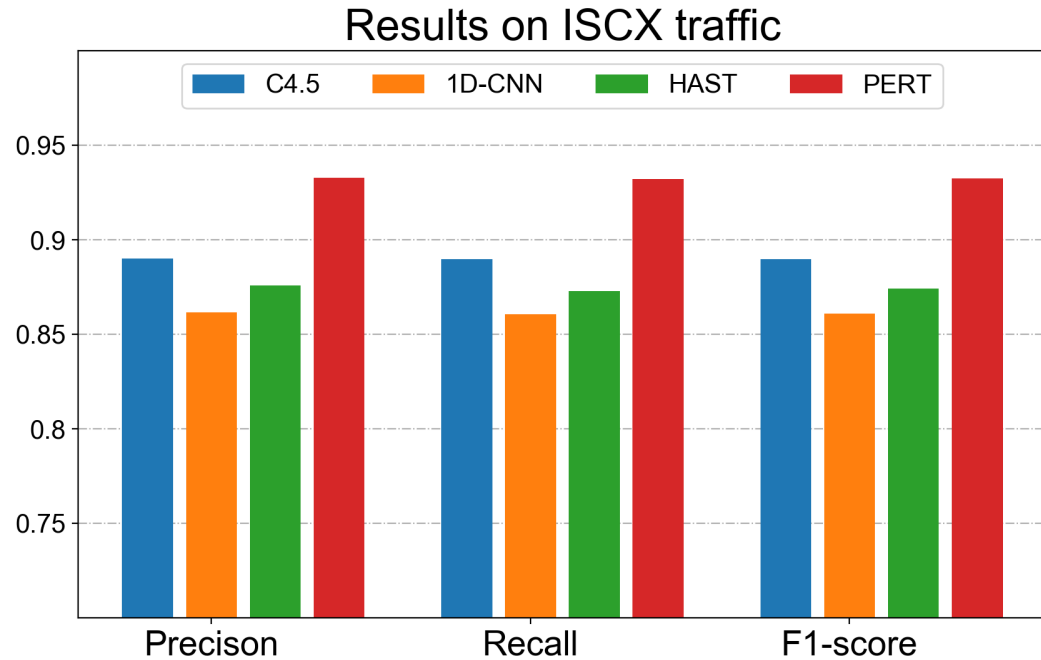[1601, 0103, 0311, 1122, 22aa, aabb, ...]

**Tokenized Payload**

# 7. PERT - Classification

- Encoding network is totally **initialized** by the pre-trained counterpart.

- After applying PERT encoding to the first N packets of a flow, a regular softmax classifier is followed to classify the **concatenated embeddings**.

- Encoding network will be further **fine-turned** during back propagation.

# 8. Encrypted Traffic Identification Experiments



- Classification results on **ISCX** traffic (12 classes) and **Android HTTPS** traffic (100 classes).

ITU KALEIDOSCOPE
ONLINE 2020

# 9. Next Steps

- **More Optimized Encoding Network**

    Follow up the ever-developing BERT research.

- **Flow-level Identification Support**

    Find a better approach to merge the PERT encoded packets.

- **Other NLP Methods**

    Evaluate other NLP methods on tokenized traffic bytes.

# ITU KALEIDOSCOPE
## ONLINE 2020

## Thank you!