

Special Topics: Machine Learning (ML) for Networking

COL867

Holi, 2025

Foundation Model

Tarun Mangla

ML for Networks

Module 1: Case studies of specific network learning tasks

Module 2: Task-agnostic automatic ML pipelines for networks

- Generalized data representation
- **Generalized ML model(s)**

Module 3: Beyond feature engineering and modeling

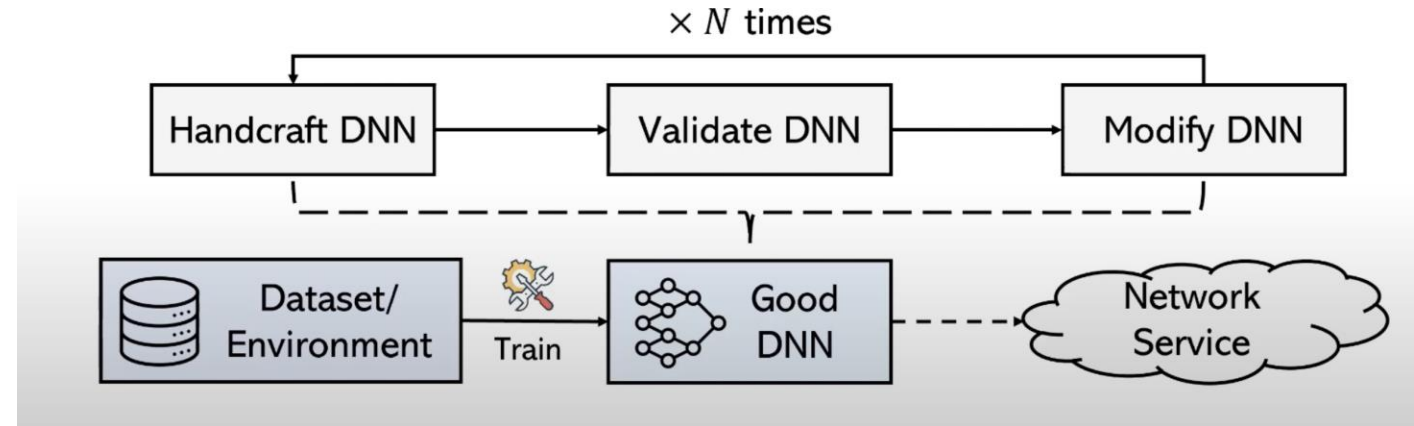
One model to rule them all

- Foundation models: trained on large corpora of unlabeled data using self-supervised learning
- Adapted to different downstream tasks with minimal tuning



Why Foundation Models

- Reduced manual effort in data representation and modeling



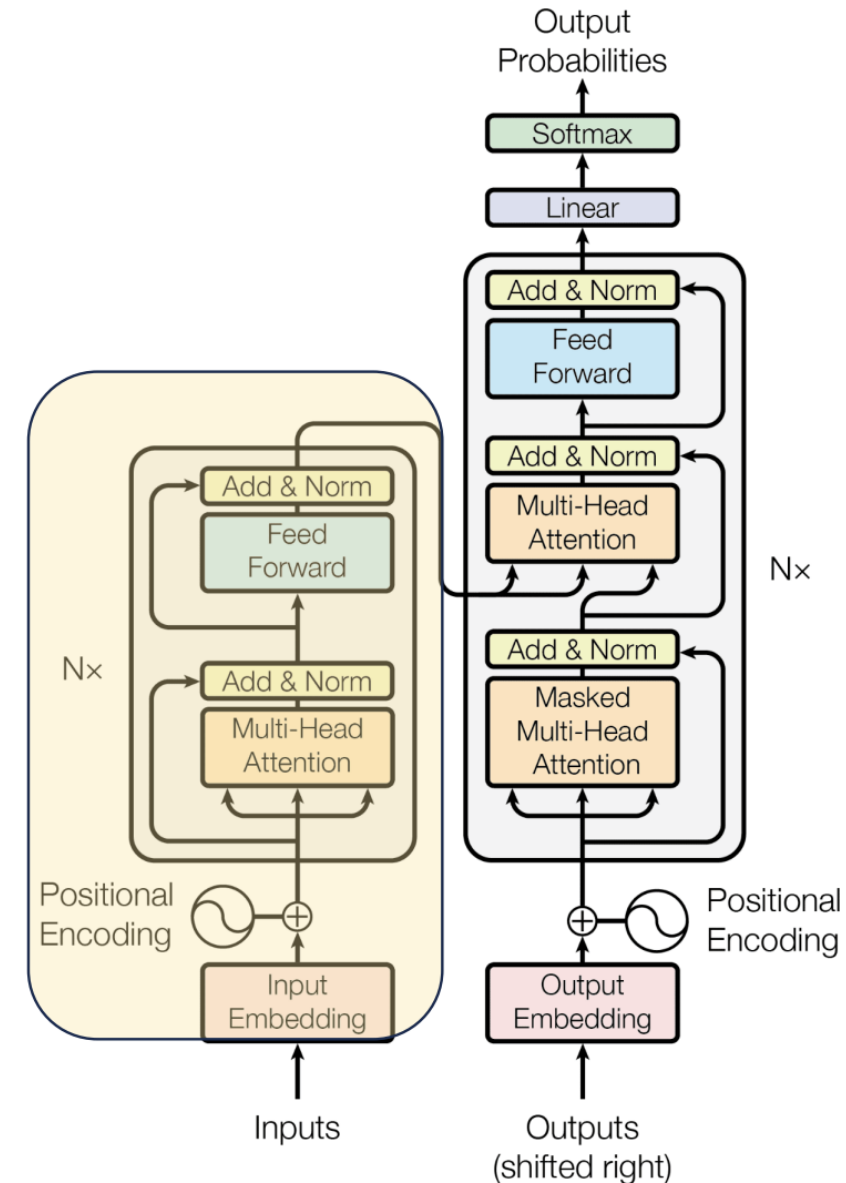
Why Foundation Models

- Reduced manual effort in data representation and modeling
- Requires lesser labeled data
- Emergent abilities



Example from NLP: Bidirectional Encoder Representation (BERT)

- Goal: Build effective language representation that can be applied for a variety of downstream tasks
- Pre-trained on a large corpus
- Uses transformer as the underlying model
- Two important aspects:
 - Tokenization
 - Pre-training task



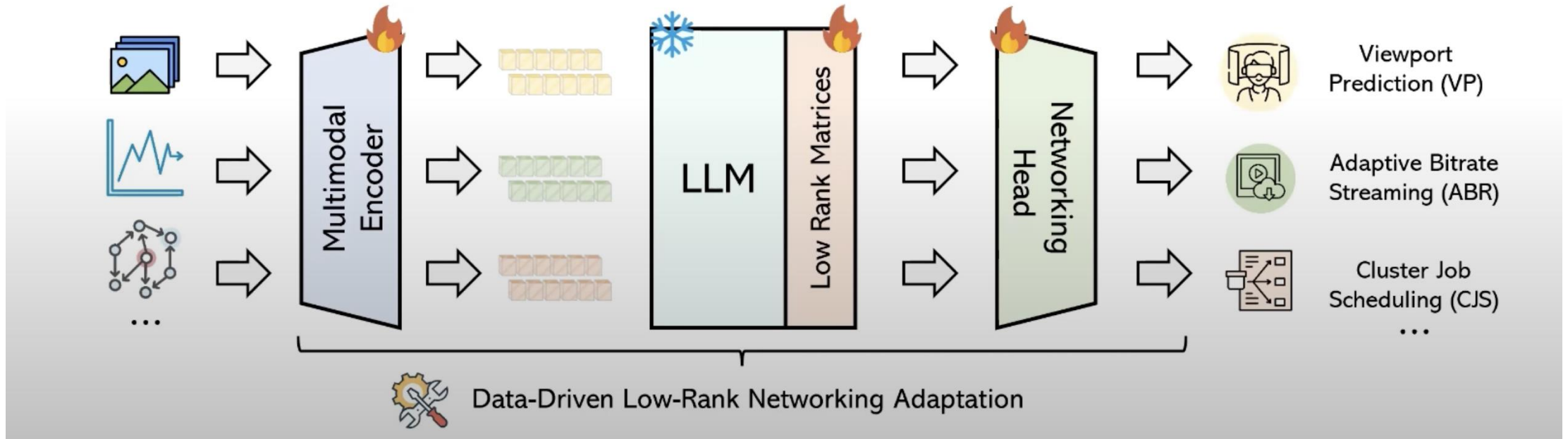
How to build a transformer for network data?

How to Build a Foundation Model for Networks?

Two different paradigms:

- Use a pre-trained large-language model
 - netLLM
- Build a foundation model from scratch*
 - netFound

NetLLM Overview

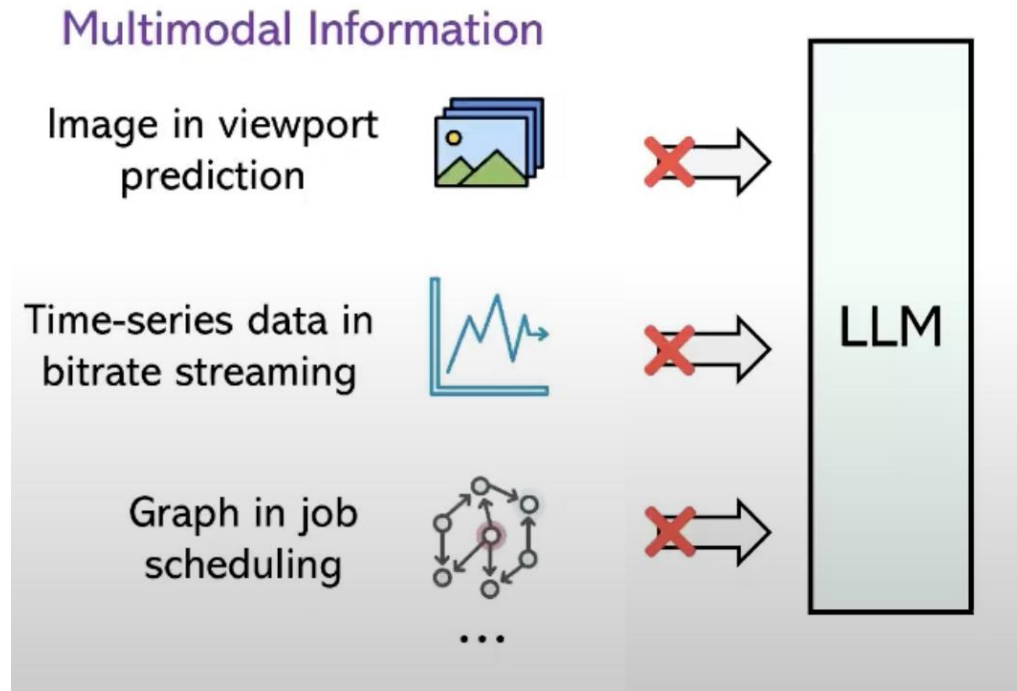


Three major innovations

- Multimodal Encoder
- Networking Head
- Data-driven low-rank networking adaptation

Challenge 1: Network Data is Multimodal

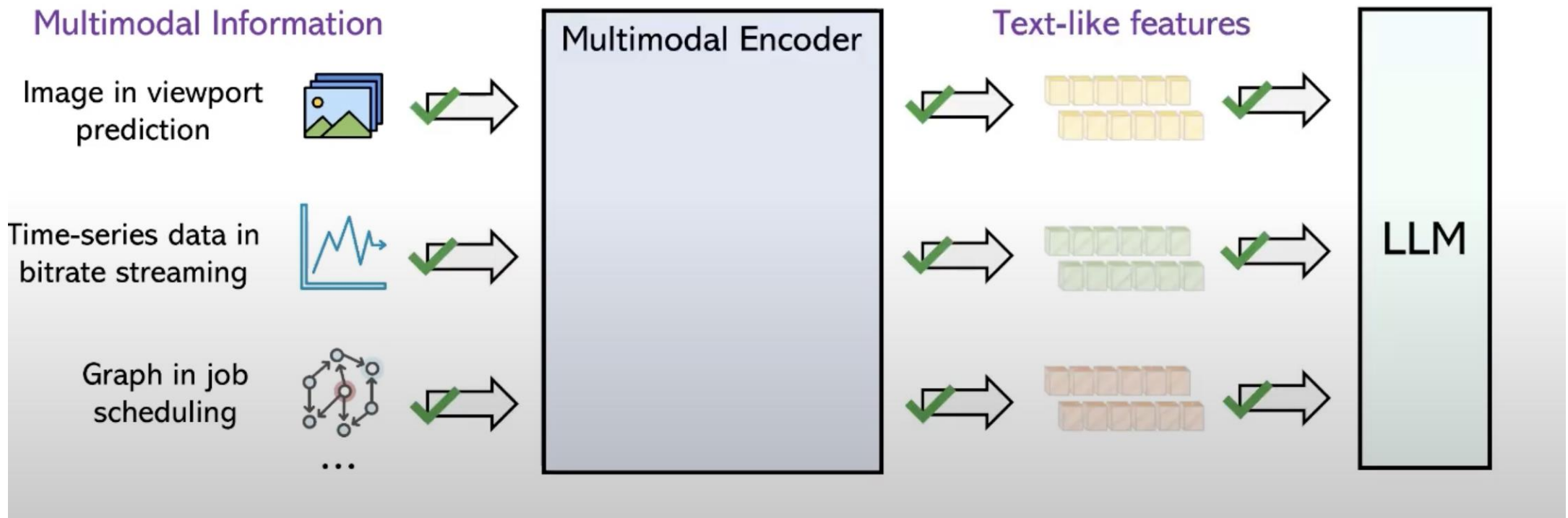
- How to enable the LLM to understand networking information?



Challenge 1: Network Data is Multimodal

- How to enable the LLM to understand networking information?

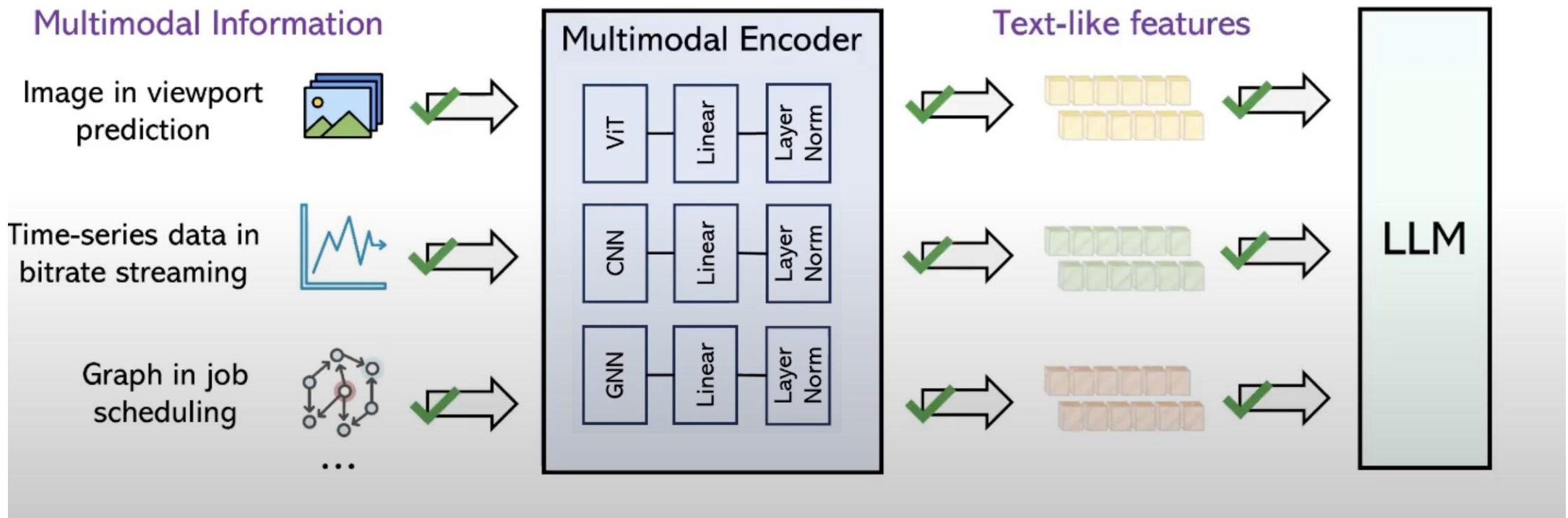
Solution: Project data into the same feature space as texts



Challenge 1: Network Data is Multimodal

- How to enable the LLM to understand networking information?

Solution: Project data into the same feature space as texts

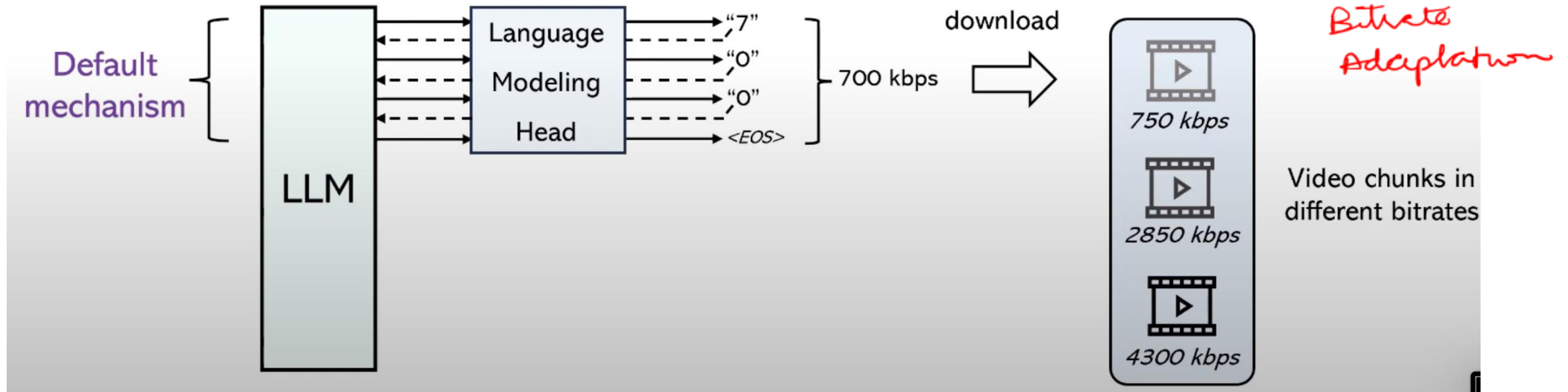


Challenge 2: Generate Output for Network Tasks

Default: Token-based generation with a language modeling output head

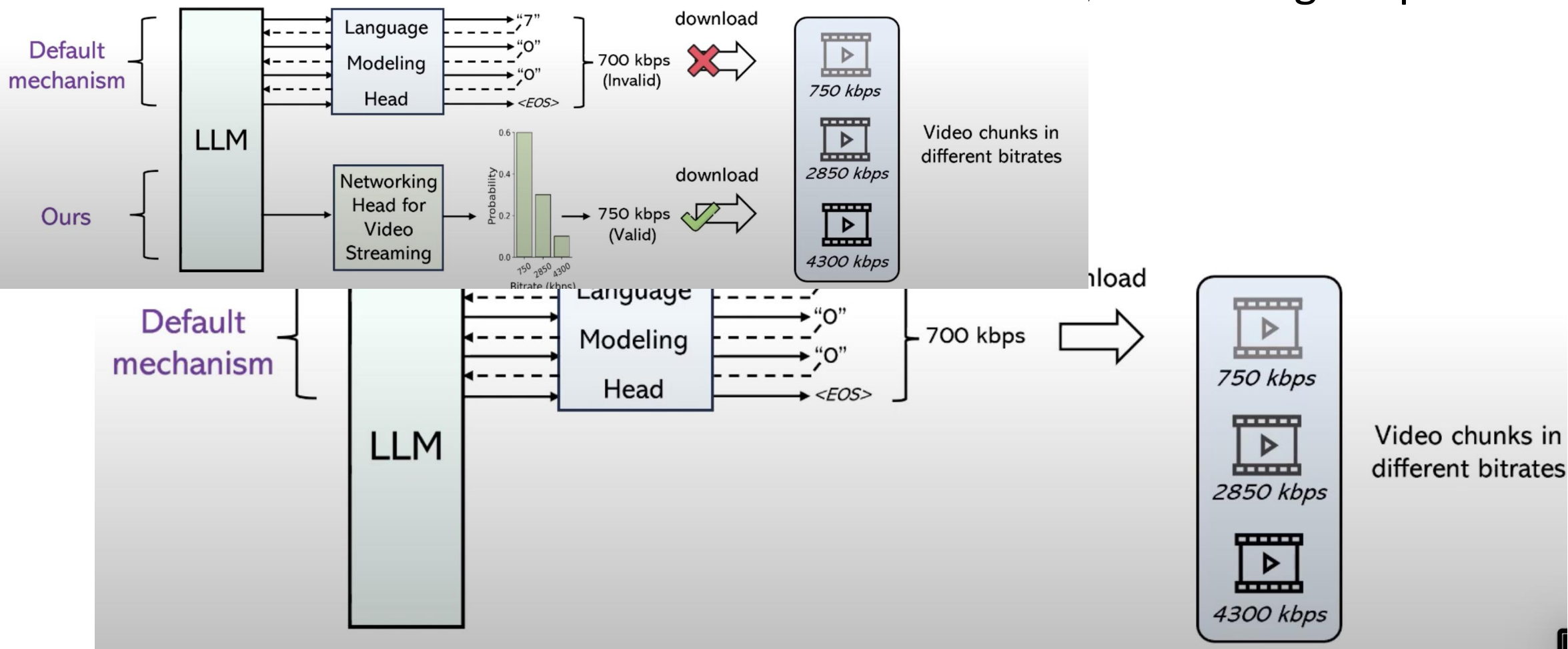
Prediction head

1. High latency, 2. Invalid answers



Challenge 2: Generate Output for Network Tasks

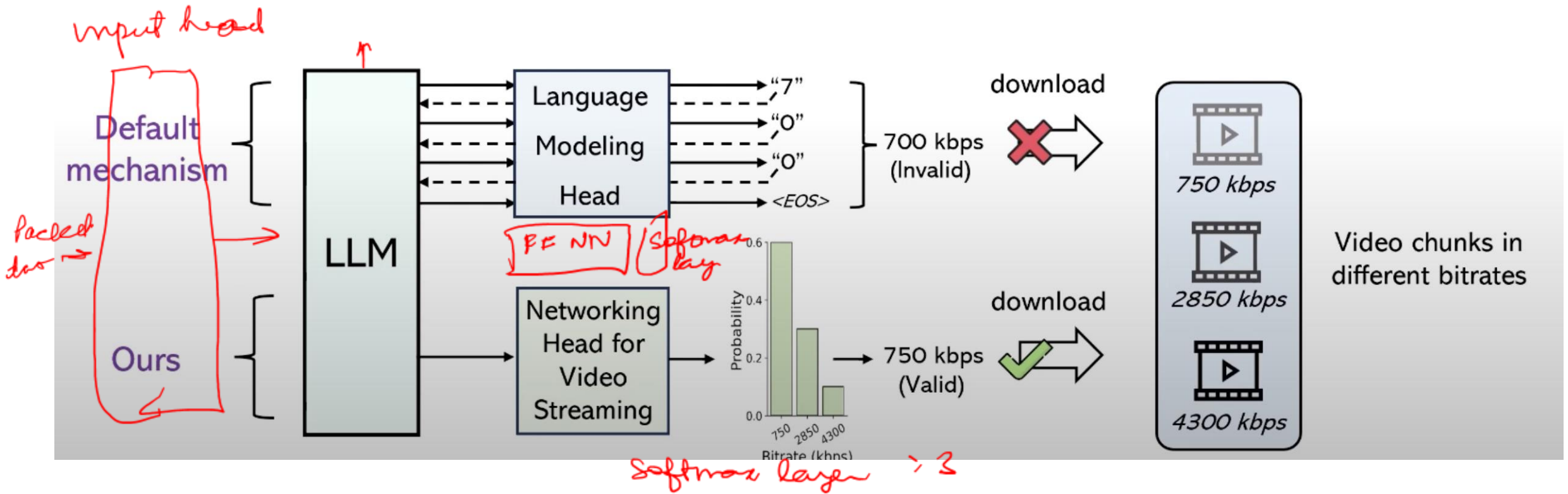
Default: Token-based generation with a language modeling output



Challenge 2: Generate Output for Network Tasks

Default: Token-based generation with a language modeling output head

Solution: Networking head to generate task-specific answers directly



Challenge 3: How to fine-tune the LLMs to learn network knowledge effectively?

- The cost of fine tuning the LLMs are expensive due to large parameter size

- E.g., Llama: 7B – 70.6 B depending on the version, model

$$\phi^l = \phi + \delta$$

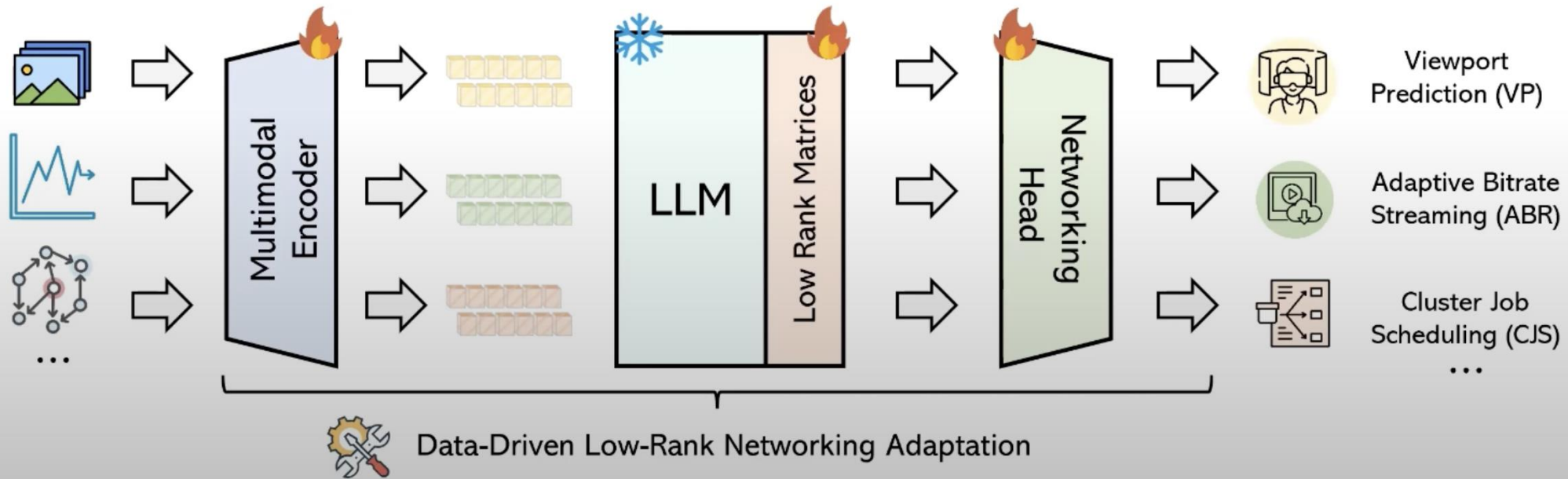
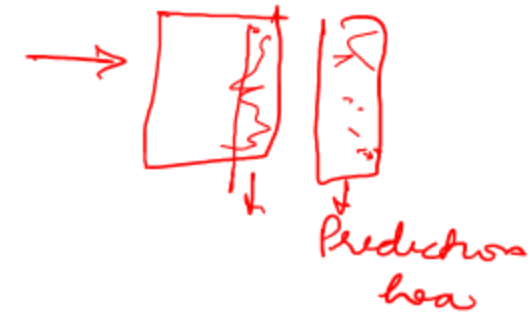
AB

$$= W + W^l \quad d \times k$$

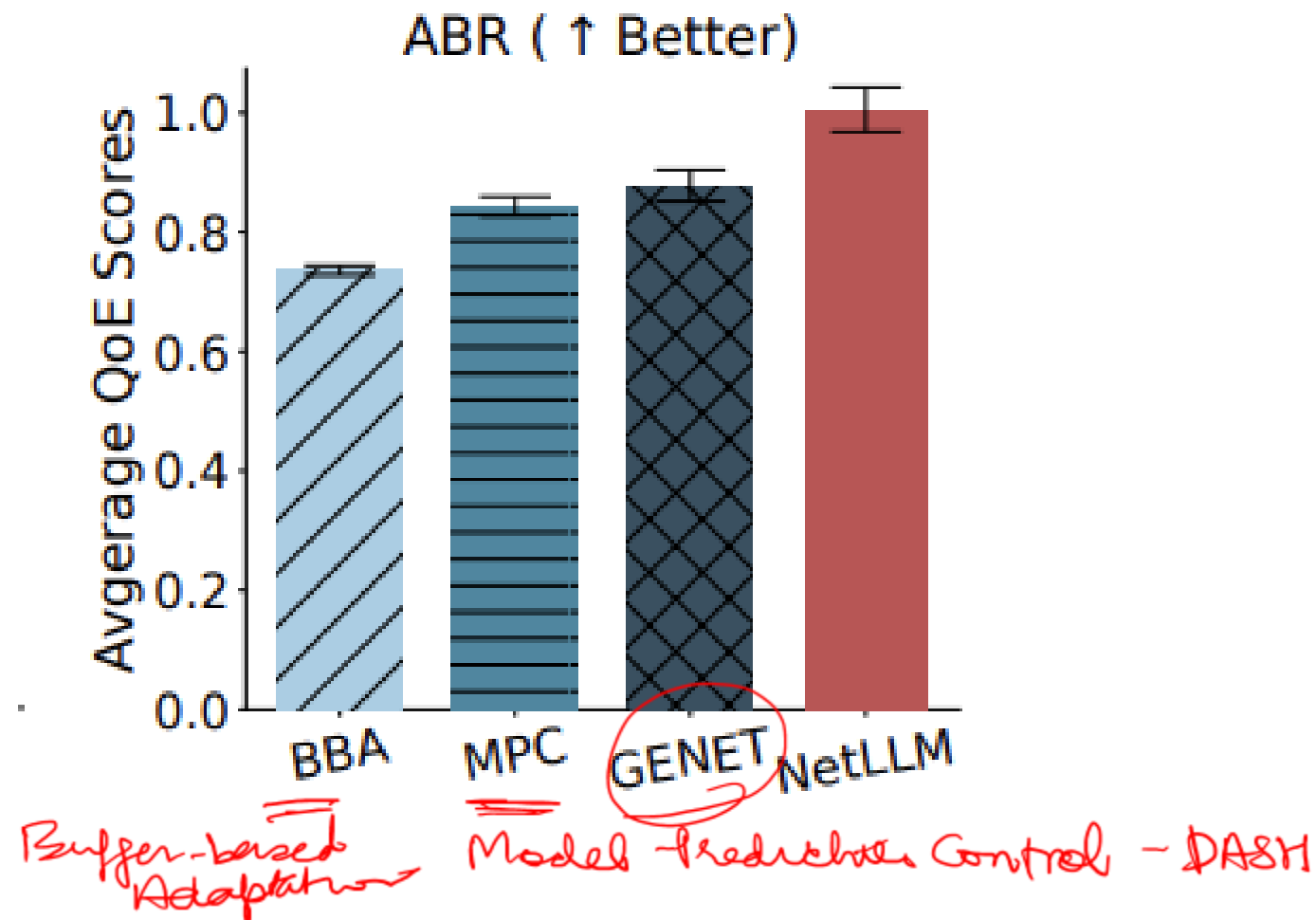
$$+ AB \rightarrow \begin{matrix} d \times r \\ r \times k \end{matrix}$$

- Solution: low-rank networking adaptation

How to train RL methods?



Evaluation: Bitrate Adaptation



Special Topics: Machine Learning (ML) for Networking

COL867

Holi, 2025

Foundation Model

Tarun Mangla

Recap

$$O(LK^2)$$

$$\begin{matrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ h_1 & h_2 \end{matrix}$$

$$O(h_1 h_2) \approx O(k_r) + O(h)$$

Task agnostic \rightarrow generalizes
 \downarrow
 saves effort \rightarrow next

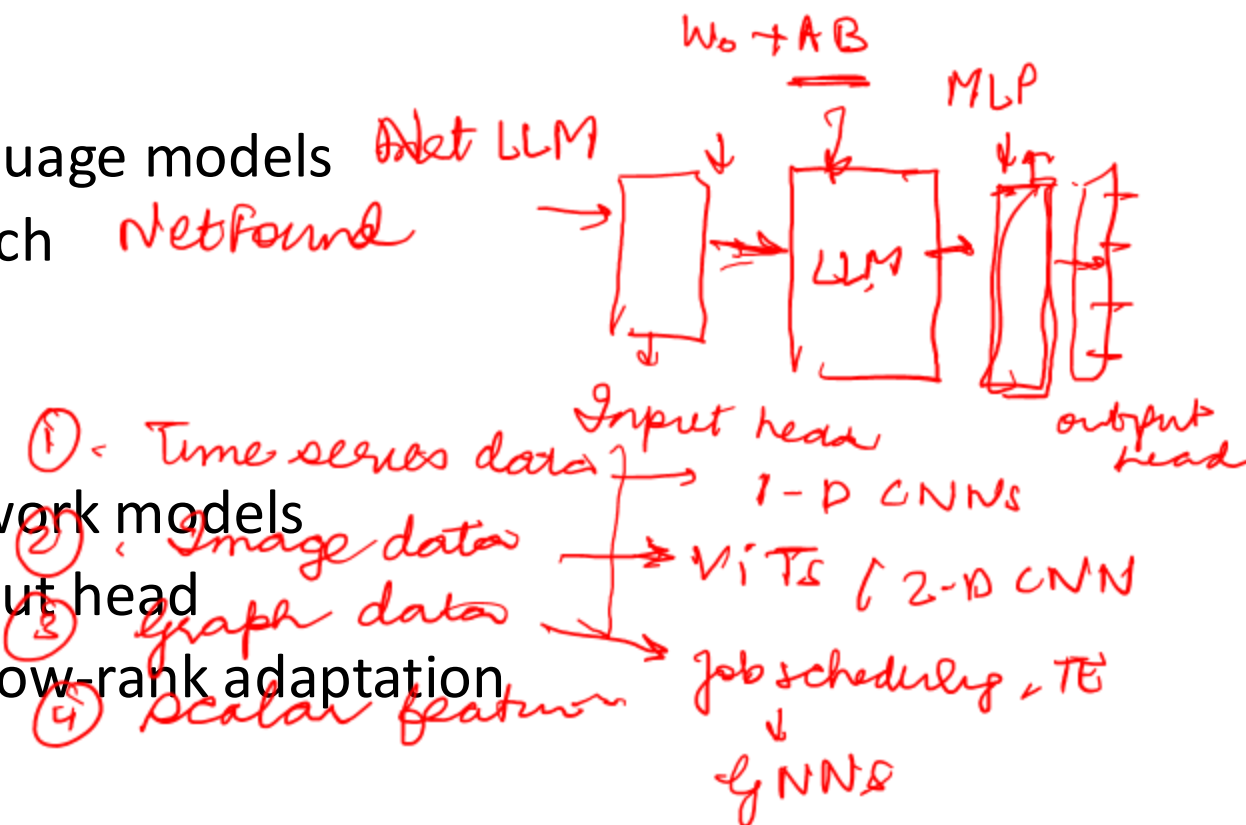
- Foundation model for network data. **Why?**

- Two approaches

- Adapt existing pre-trained large language models
- Build a foundation model from scratch

- NetLLM

- Input:** Use existing deep neural network models
- Output:** Append a task-specific output head
- Training:** Improve fine tuning using low-rank adaptation



How to Build a Foundation Model for Networks?

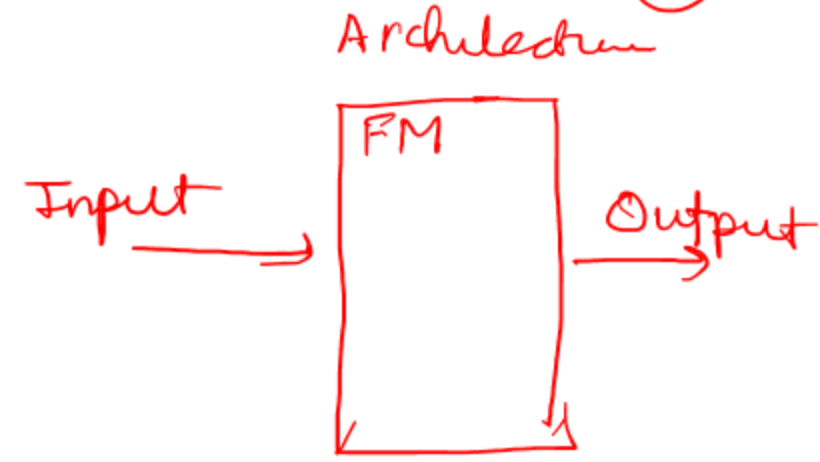
Two different paradigms:

- Use a pre-trained large-language model
 - netLLM
- **Build a foundation model from scratch***
 - netFound

⑥ What should be the ML model? → Architecture

Building NetFound

- ① How to tokenize data?
- ②. What are the tasks → what is the data?
- ③. Incorporate multimodality/
- ④. Latency of the system
- ⑤ Training



④ Training

Given Seq. of Packet 

Input: How to tokenize the data?

- What is the network data?

Packet-level data 

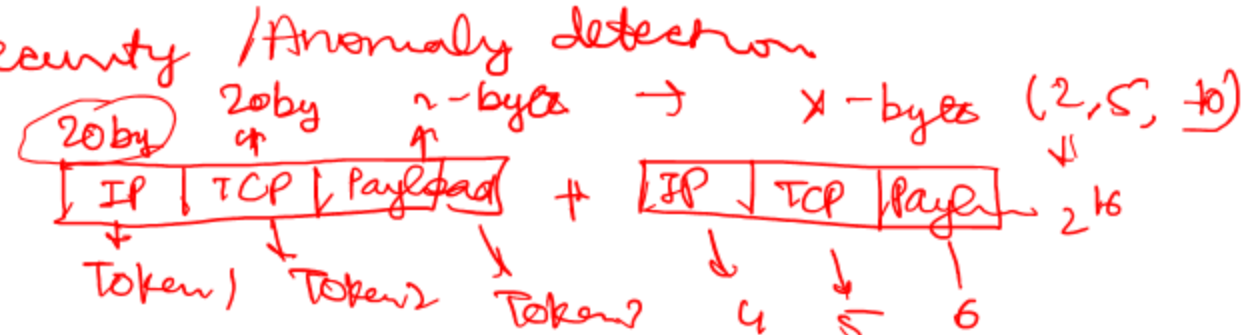
- How to tokenize packet-level data?

① Protocol semantics

② How to handle variable length

- Challenges:

- How to handle both content and statistics (multi-modality)
- How to retain semantic integrity?
- How to handle variable length input?



Missing some information

- : IAT
- : # Pkt per flow
- : Throughput

Preserving Semantic Integrity

- Use a protocol-aware tokenizer

- Consider 2-byte token

- Pad shorter fields

RTP header

TLS header



2 bytes

8-bit field

8 bits

Protocol	Fields				
IPv4	HeaderLen	ToS	TotalLen	Flags	TTL
TCP	Flags	WinSize	SeqNum	AckNum	UrgentPtr
UDP	Length				
ICMP	Type	Code			
Payload	12 bytes				

5
7

6 →

↓
18 tokens

① SNI / DNS
per packet

32-bit → 2 tokens

Content # Metadata

Handling multiple modality

$[0, \dots, T_{max}]$

- Example modalities:

- Temporal details (IAT, pkt sizes)
- Statistical aggregates (throughput)
- Contextual information (downstream/up)

Meta data

t_1 t_2 $t_3 \dots$

1-hot

④ Tokens	CLS-B	MASK	...	0x2f43	CLS-B	0x0004
Position	1	2	...	108	1	2
Direction	1	1	...	1	-1	-1
# bytes in burst	4517	4517	...	4517	54	54
# pkts in burst	6	6	...	6	1	1
IAT	18	18	...	18	25	25
Protocol	17	17	...	17	17	17

Handling Variable Length

- Problem statement
- **Solution:** Data-driven approach to determine the sequence length

→ Seq length can change

↓

T
=

Maxim

Avg

Median

→ Trade-off

Model

- What should be the underlying model? *Transformers*

- Insight: Inherent hierarchy in the network data. Can we leverage it?

Packets → Burst → Flow
↓
Subnet ← Device ← Session

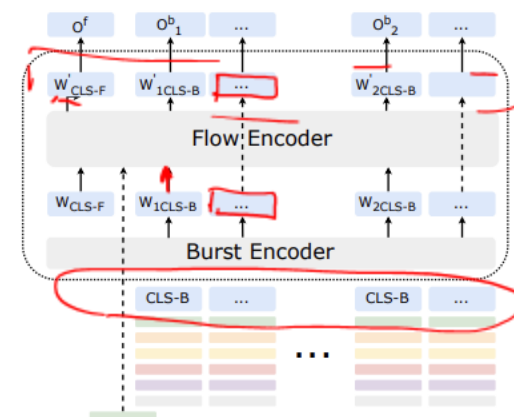
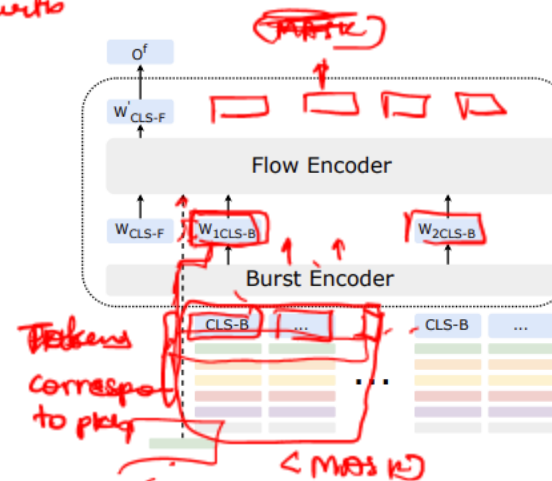
Flow-level classification

Tokenize
FM

F: B₁ B₂ B₃ B₄ ... B_N
↓ Flow encoder
Flow

□ □ □ □ □

Burst: group of pkts with
IAT < δ



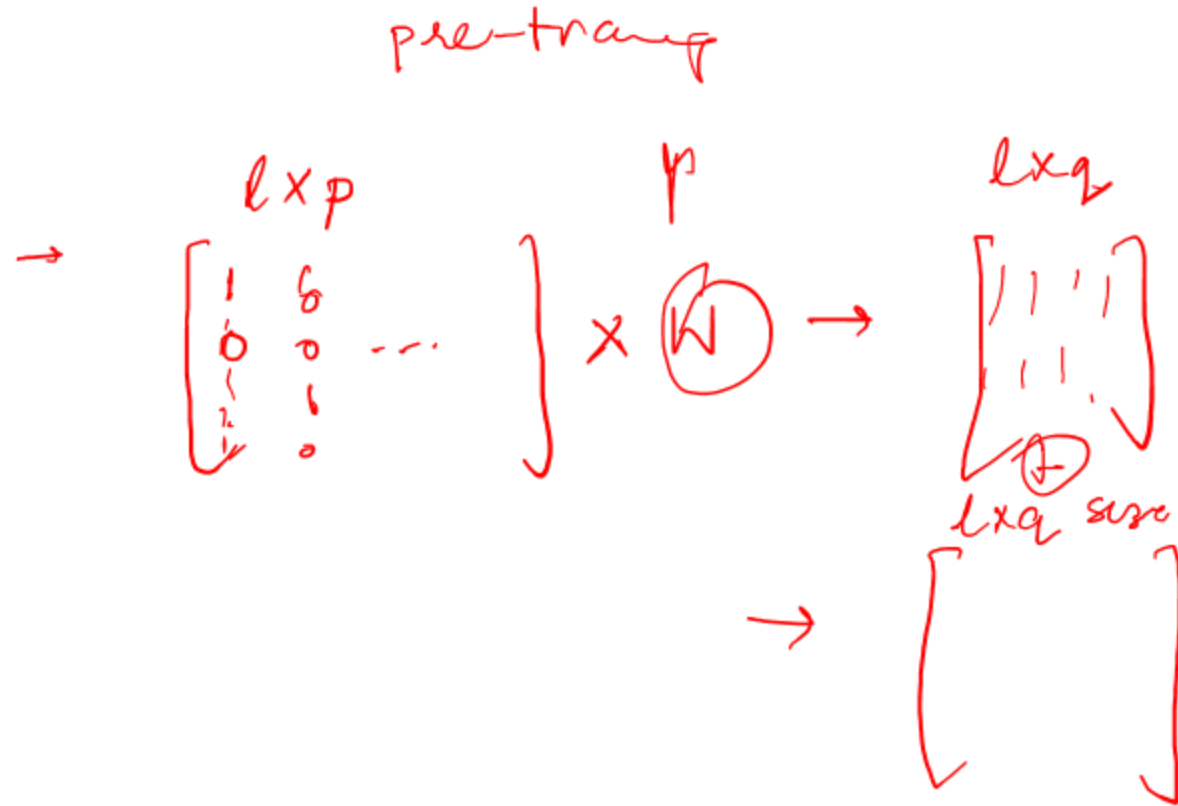
Workflow

- Data-preprocessing
 - Group packets into flows
 - Discard all flows with only 1-2 packets
 - Flow: 6 packets per burst, 12 bursts per flow
- Featurization
 - Extract relevant headers and metadata
- Tokenization
 - Maximum number of tokens per packet (why?): 18

Protocol	Fields				
IPv4	HeaderLen	ToS	TotalLen	Flags	TTL
TCP	Flags	WinSize	SeqNum	AckNum	UrgentPtr
UDP	Length				
ICMP	Type	Code			
Payload	12 bytes				

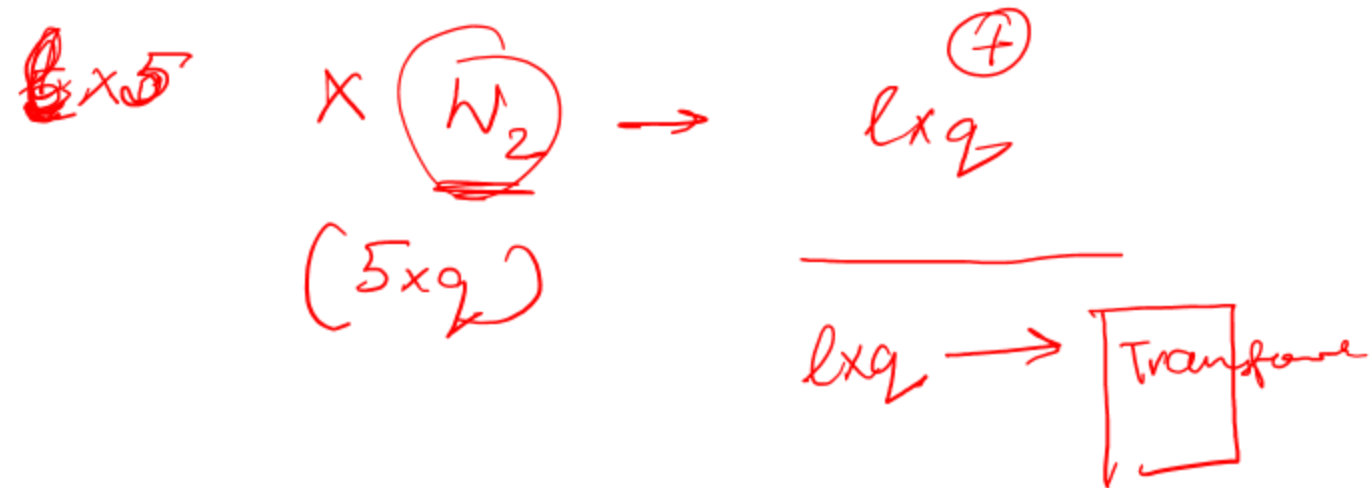
Token Embedding

- Packet field token embedding



- Positional embedding

- Metadata embedding



Training

- Self-supervised pre-training

Mask 30%

- Fine-tuning

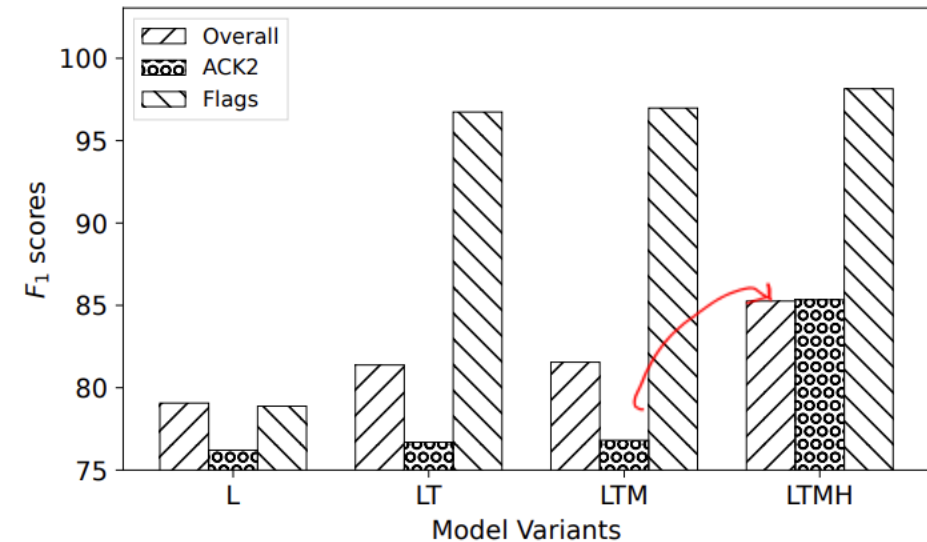
- Add two-layer MLP

• Update weights of pretrained model

Evaluation

Masked Token Prediction

- L: Flat transformer architecture
- LT: Protocol aware tokenizer
- LTM: LT + Metadata
- LTMH: LTM + Hierarchical



Fine-tuning Tasks

Task	Type	Dataset	Curtains (%)	nPrintML (%)	ET-BERT (%)	YaTC (%)	netFound (our) (%)
1	Traffic Classification	Campus dataset	54.53 ± 0.97 $p < 0.001$	87.22 ± 0.12 $p < 0.001$	72.26 ± 0.38 $p < 0.001$	76.54 ± 0.23 $p < 0.001$	96.08 ± 0.04 –
2	Application Fingerprinting	Crossmarkets [59] ($Acc@10$)	20.64 ± 0.13 $p < 0.001$	64.83 ± 0.28 $p = 0.098$	35.62 ± 0.39 $p < 0.001$	58.13 ± 0.89 $p = 0.010$	66.35 ± 0.99 –
3		ISCXVPN-2016 [60]	66.85 ± 2.21 $p = 0.003$	84.10 ± 0.41 $p < 0.001$	77.57 ± 1.20 $p < 0.001$	83.84 ± 0.24 $p < 0.001$	91.02 ± 0.10 –
4	Intrusion Detection	CICIDS2017 [61]	99.75 ± 0.16 $p = 0.082$	99.93 ± 0.01 $p = 0.012$	99.94 ± 0.01 $p = 0.018$	99.92 ± 0.01 $p = 0.005$	99.99 ± 0.01 –
5	HTTP Bruteforce Detection	netUnicorn [5]	96.82 ± 0.22 $p = 0.006$	98.51 ± 0.02 $p < 0.001$	98.63 ± 0.02 $p < 0.001$	98.73 ± 0.10 $p = 0.030$	99.01 ± 0.01 –

Discussion

ML for Networks

Module 1: Case studies of specific network learning tasks

Module 2: Task-agnostic automatic ML pipelines for networks

- Generalized data representation
- Generalized ML model(s)

Module 3: Beyond feature engineering and modeling

Model Architecture

- Want to use transformer

Why Foundation Models for Networking Data?

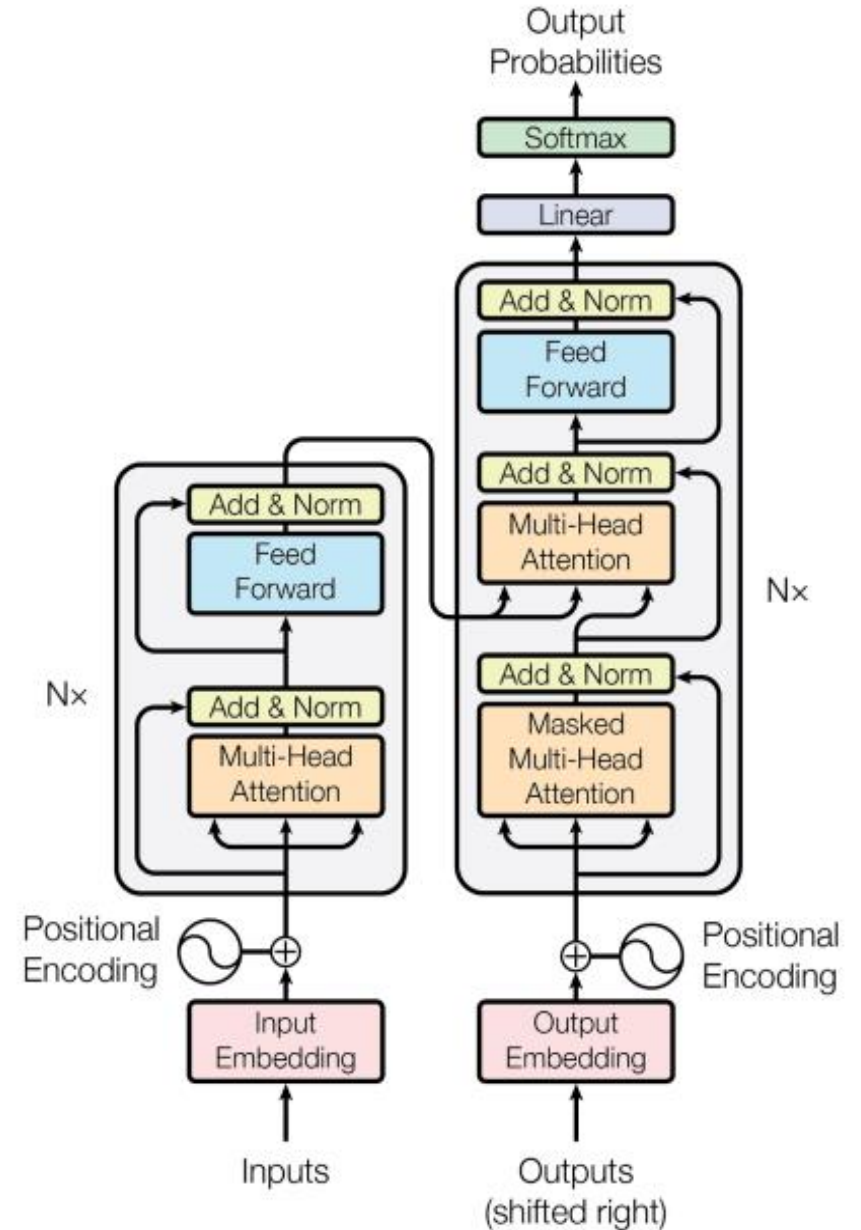
- Network learning approaches consist of classification, reinforcement learning, anomaly detection, generative
 - Foundation models have been successfully applied to these problems
- Abundant unlabeled sequential data
 - Campus networks
 - Data center networks
 - Transit/ISP networks
- Rich semantic content (like text)
 - Well-defined protocols

Transformer Architecture

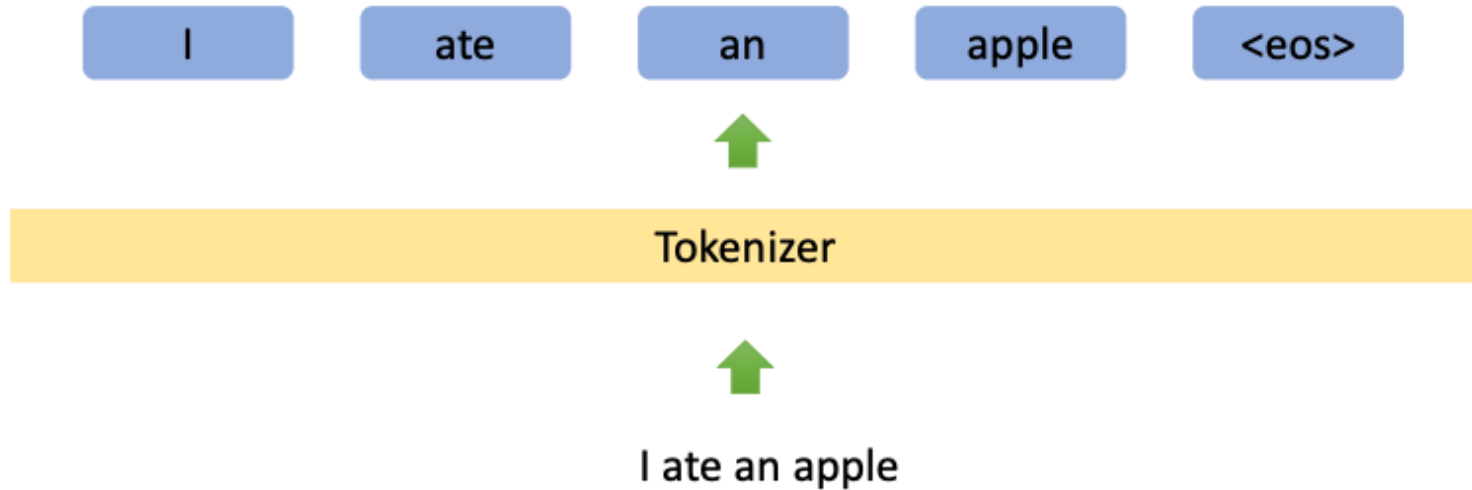
- Tokenization
- Input Embedding
- Positional Encoding
- Residual Connection
- Query
- Key
- Value
- Add & Norm
- Encoder
- Decoder



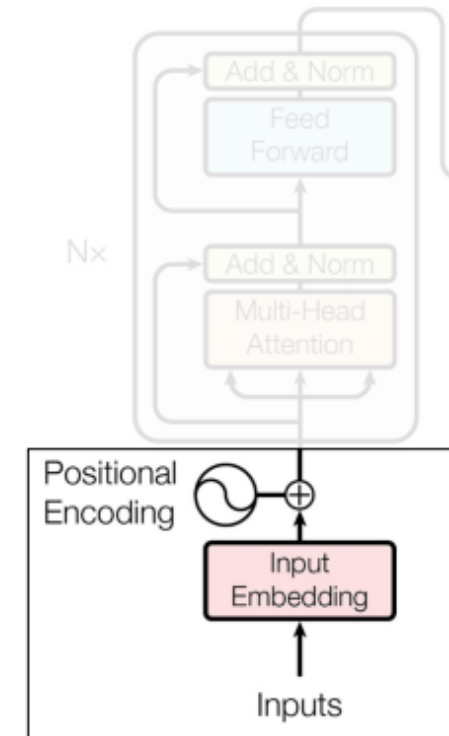
- Attention
- Self-Attention
- Multi-Head Attention
- Encoder Attention
- Probabilities / Logits
- Encoder models
- Decoder only models



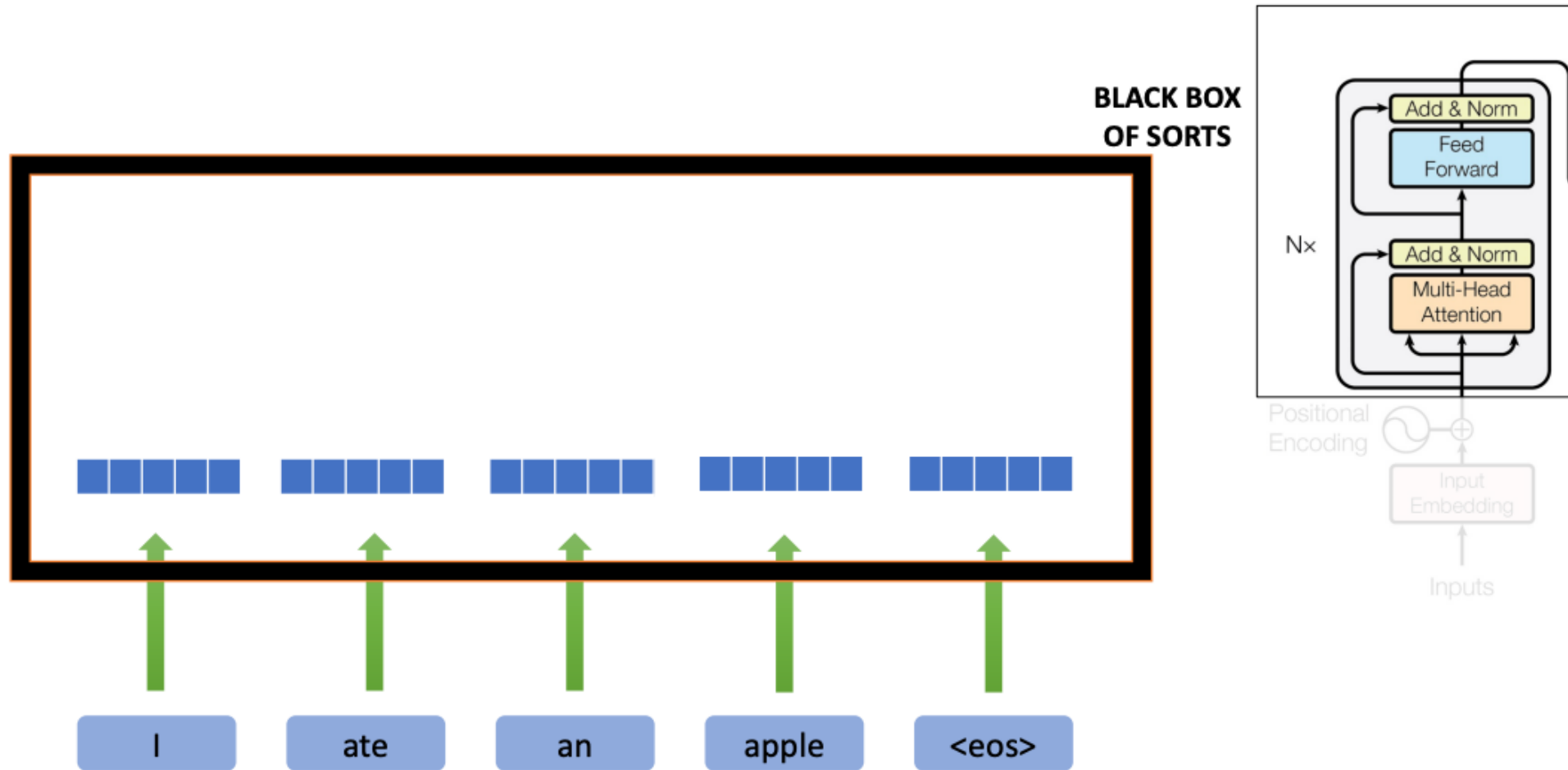
Processing Input



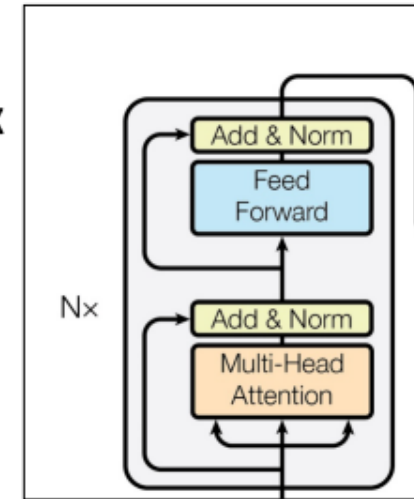
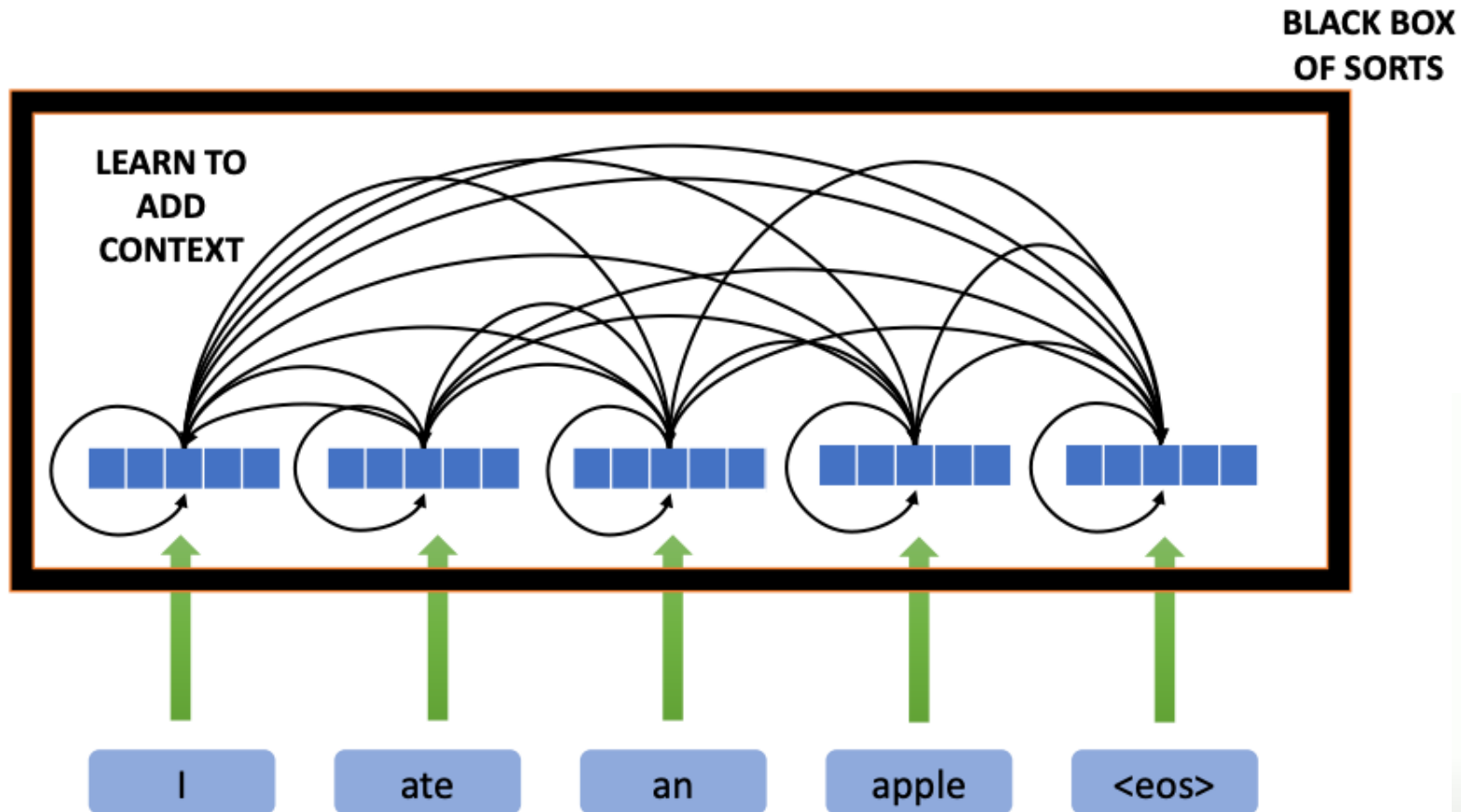
Generate Input Embeddings



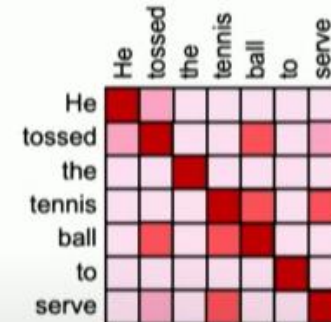
Capturing Context



Capturing Context



Attention weighting: where to attend to!
How similar is the key to the query?

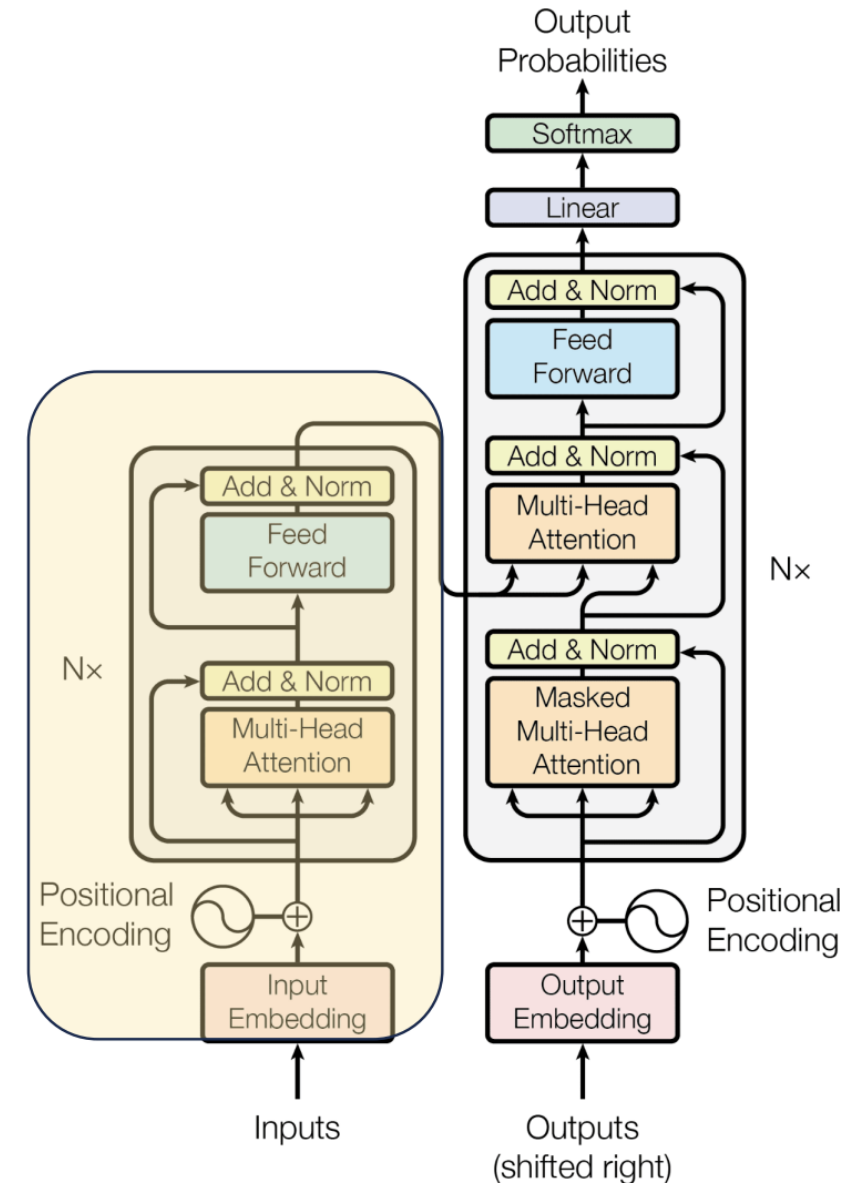
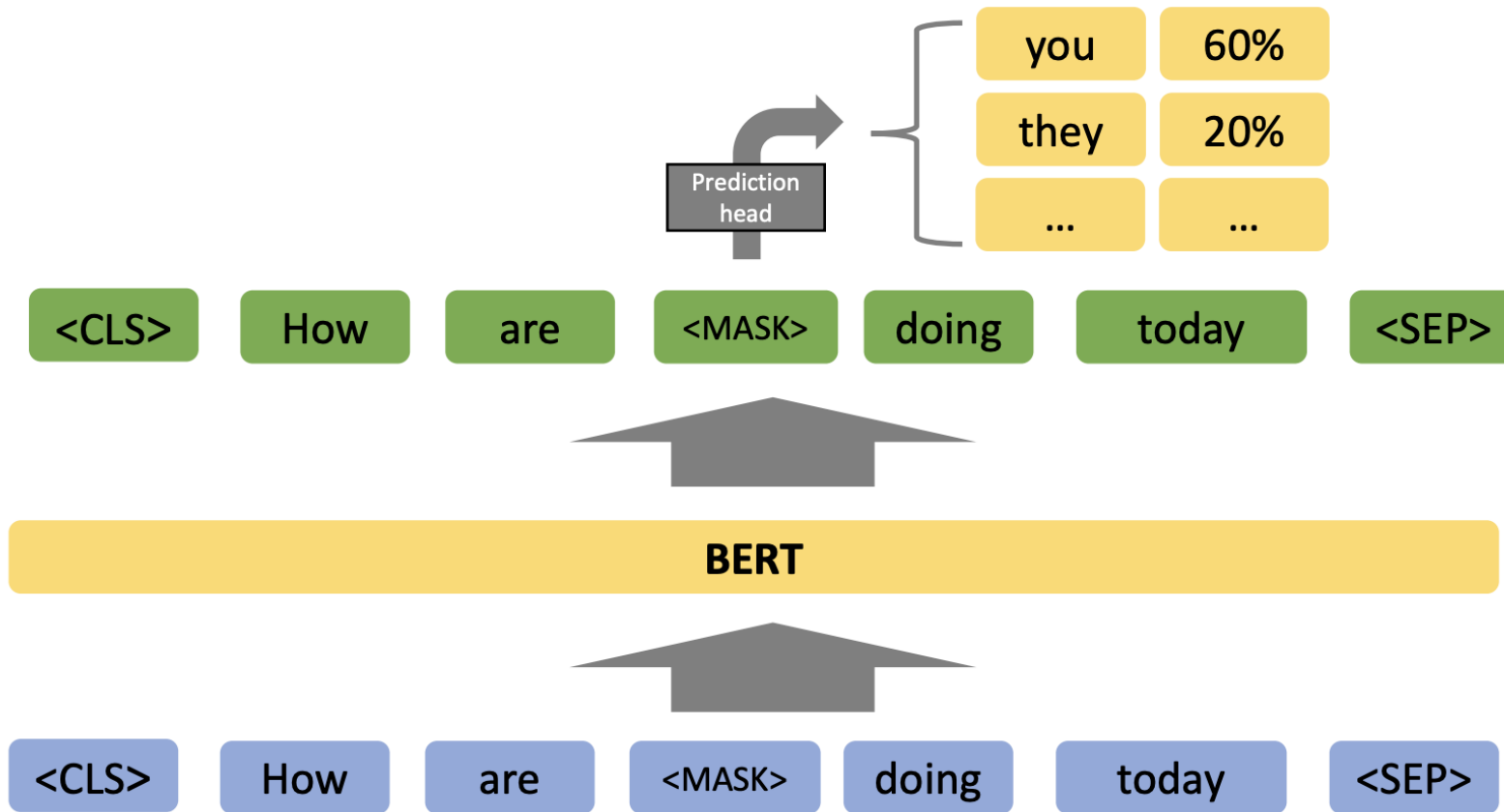


$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right)$$

Attention weighting

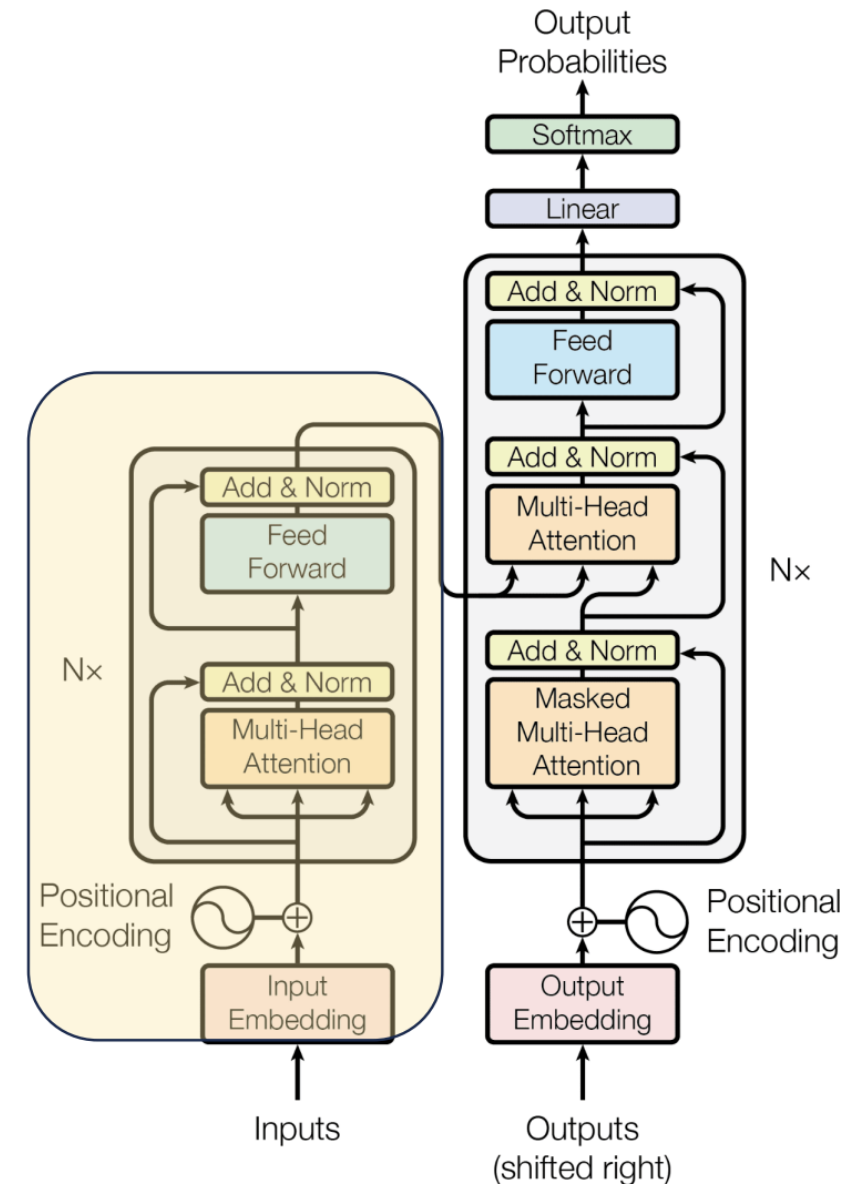
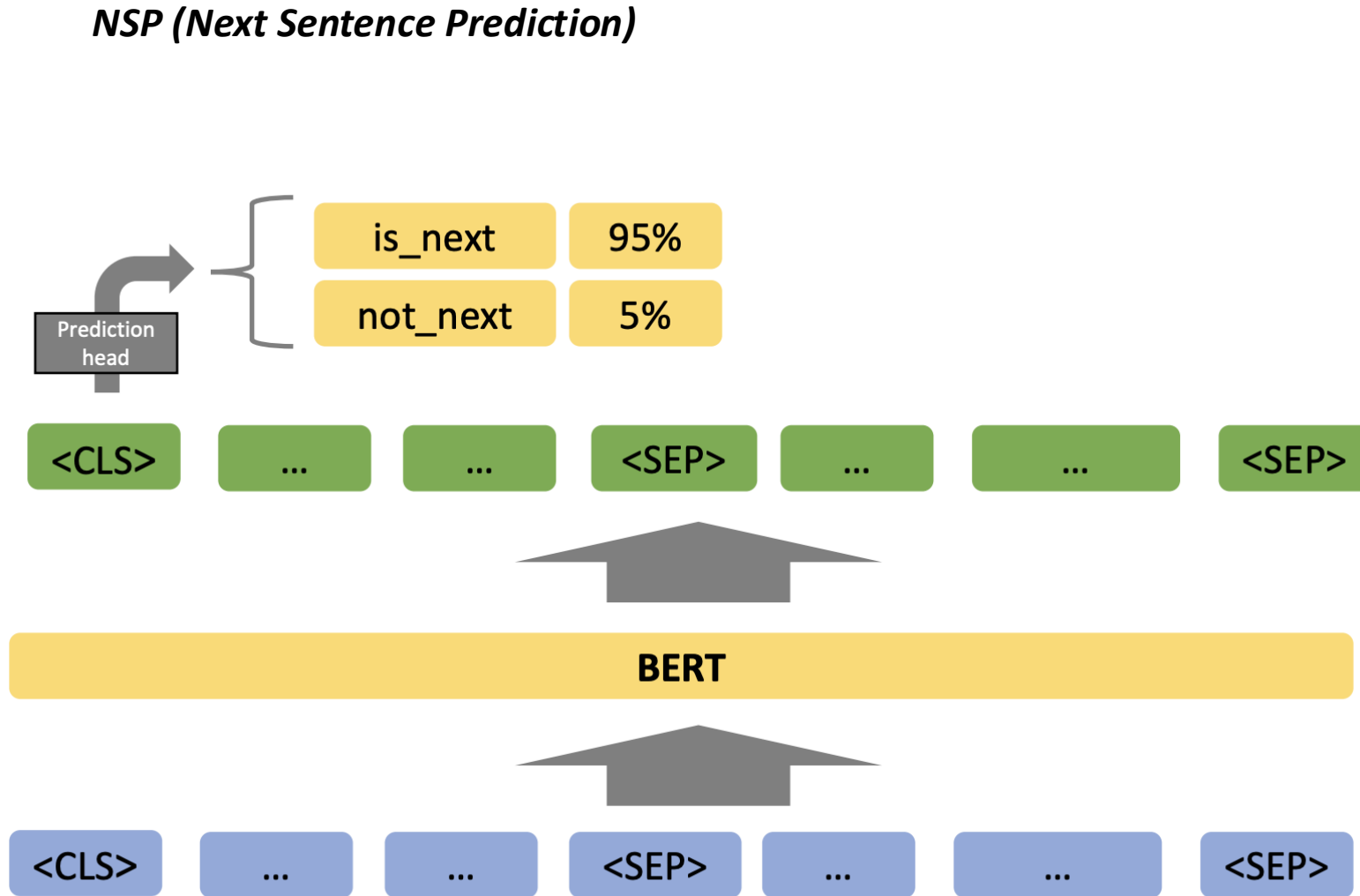
BERT: Bidirectional Encoder Representation

MLM (Masked Language Modeling)



BERT: Bidirectional Encoder Representation

NSP (Next Sentence Prediction)



Why Foundation Models for Networking Data?

- Network learning approaches consist of classification, reinforcement learning, anomaly detection, generative
 - Foundation models have been successfully applied to these problems
- Abundant unlabeled sequential data
 - Campus networks
 - Data center networks
 - Transit/ISP networks
- Rich semantic content (like text)
 - Well-defined protocols

Challenges

- Tokenizer
 - Text: Tokens can be characters or words
 - What is a token for network data?
- Context
 - How do you define context?
 - What are the pre-training tasks?
- **Post mid-term:** Two papers on design of foundation model for networks
 - netFound
 - netLLM

Resources

- CMU Deep Learning lecture: <https://deeplearning.cs.cmu.edu/S24/index.html>
- MIT Deep Learning course: <https://introtodeeplearning.com/>
- Explanation with code: <https://nlp.seas.harvard.edu/2018/04/03/attention.html>
- Jay Alammar, The illustrated transformer: <http://jalammar.github.io/illustrated-transformer/>

Motivation

Network data has a unique context

- **Multi-modal**

- Data from various contexts and perspectives
- Cross-layer interactions
- Network conditions and protocol interactions

- **Hierarchical**

- Packets -> bursts -> flows -> session -> devices -> subnet etc.
- Different learning problems make decisions at different granularity

Can we design a FM keeping in mind this
unique context?

Design Decisions

- What kind of model architecture?
- How to represent network data as input to the model?
- How to (pre-)train the model?

Design Decisions

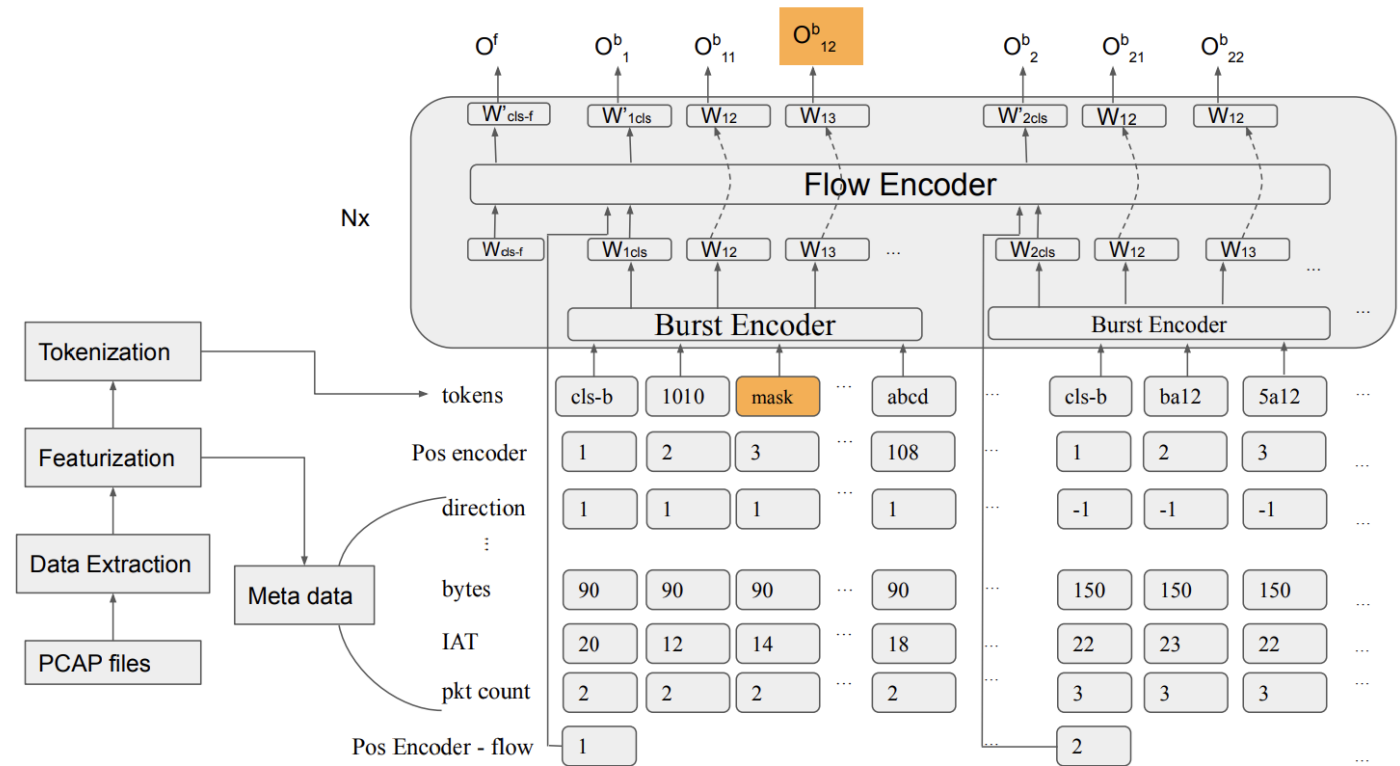
- **What kind of model architecture?**

Transformers

- How to represent network data as input to the model?
- How to (pre-)train the model?

Capturing Multi-modal Inputs

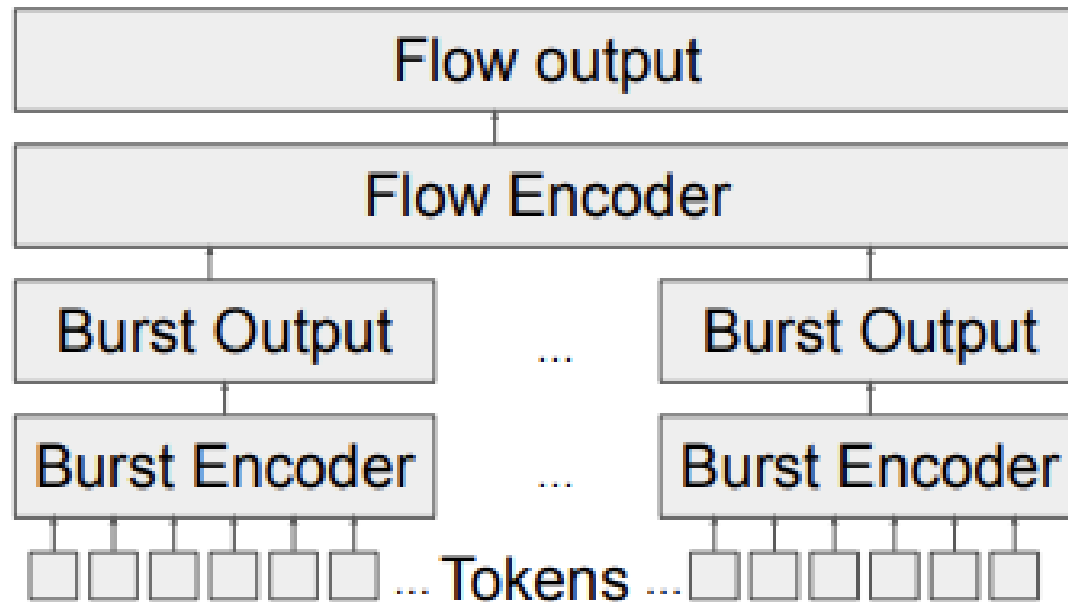
- Raw bytes
 - Similar to ET-BERT, use 2-byte tokens
 - Consider only 12 bytes of payload
- Meta-data
 - Embed other modalities as metadata
 - E.g., direction, time, number of packets in a burst etc.



Capturing Hierarchical Structure

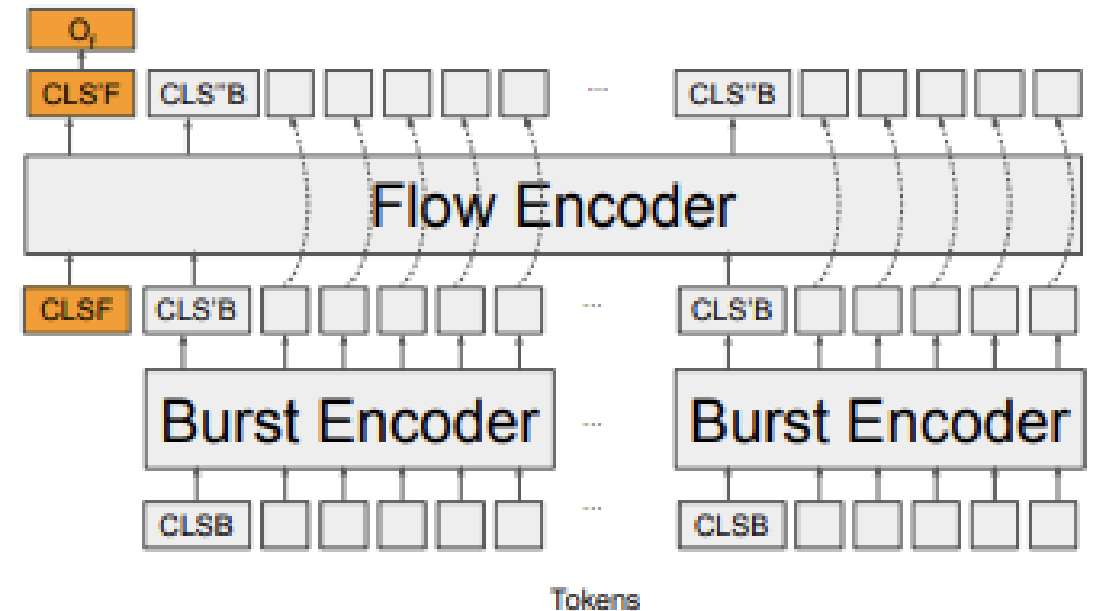
- Hierarchical data: packets -> bursts -> flows etc.
- How to create models to handle hierarchy in data?
- Option 1: Create different models for each level of hierarchy

Capturing Hierarchical Structure



Naive Approach

- Challenging to implement MLM task in this approach



Proposed Approach

- Use skip connections

Design Decisions

- What kind of model architecture?

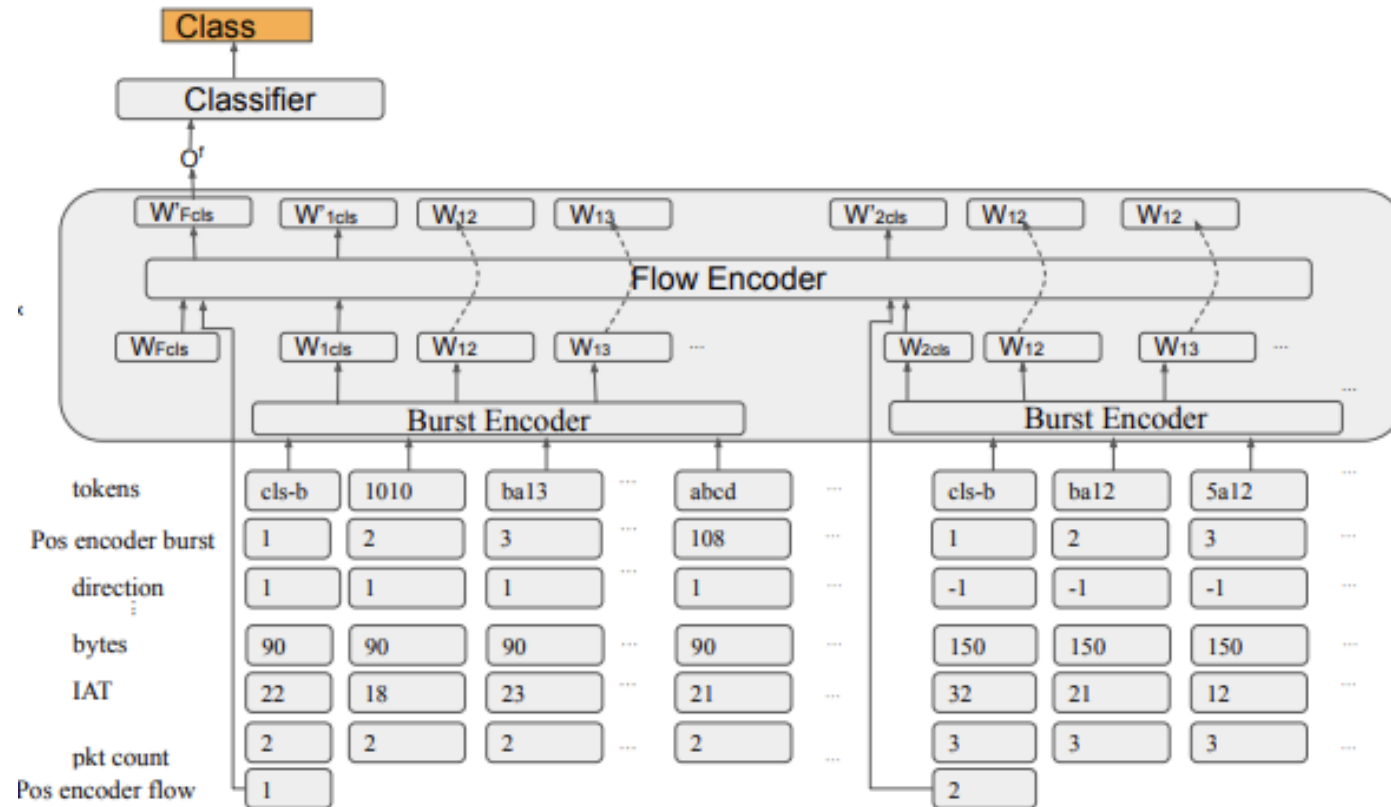
Transformers

- How to represent network data as input to the model?

- **How to (pre-)train the model?**

Masked Language Modeling

Putting It All Together



Workflow

- Tokenization
 - 13 packet fields with a 279-bit vector
 - 18 tokens per-packet
 - Special tokens: [PAD], [CLS-B], [CLS-F], [MASK]
- 6 packets per burst, 12 bursts per flow
- **Pretraining:** Only Mask-language Modeling task

Case Study: Mask Prediction

Pkt #	TCP Flag (Masked)	Burst-2	Length	Window	Seq #	Ack #
1	ACK	CLS	52	506	2726740280	2946828322
2	ACK+PUSH	CLS	735	501	2726740950	2946829950
3	ACK+PUSH	CLS	1339	824	2726752145	2946867610
4	ACK+PUSH	CLS	618	501	2726742967	2946832064
5	ACK	CLS	52	501	2726741633	2946830128
6	ACK	CLS	52	497	2726744816	2946848371

Summary

- Task-agnostic network data representation is useful
- Three approaches:
 - Extract a set of (empirically) effective flow-level features
 - Packet-level bit vector representation
 - Foundation models
- Foundation models are promising. Few open questions
 - What is the best design of an FM for networking?
 - What are these models learning?
 - What are the resource implications of running these models?
 - How do we benchmark these models?
 - ...