

JOINT INSTITUTE
交大密西根学院

ELMA: Early-Exit Offloading for Embedded Question Answering Applications

Capstone Design

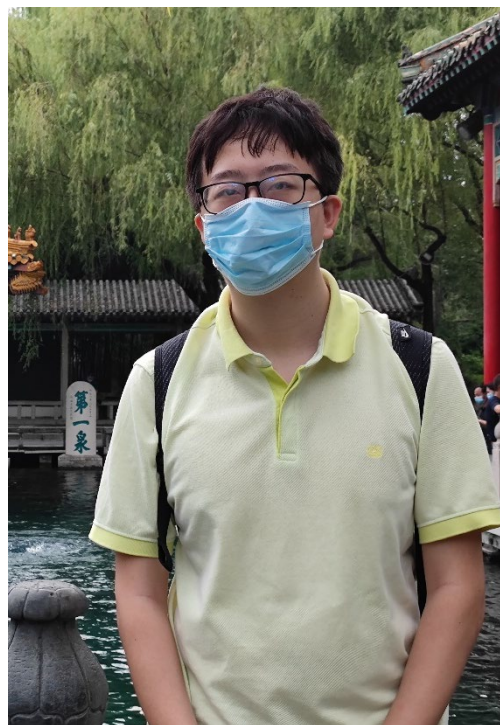
Dec 15, 2021

Group:	3
Instructor:	Prof. An Zou
Sponsor:	UM-SJTU Joint Institute
Group Member:	Yihua Liu, Shuocheng Chen, Yiming Ju

Team Members



JOINT INSTITUTE
交大密西根学院



Yihua Liu

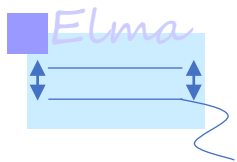


Shuocheng Chen



Yiming Ju

- Introduction
- Design Specifications
- Concept Generation & Selection
- Design Description
- Implementation & Validation
- Discussion & Conclusion



JOINT INSTITUTE
交大密西根学院

1. Introduction

Introduction

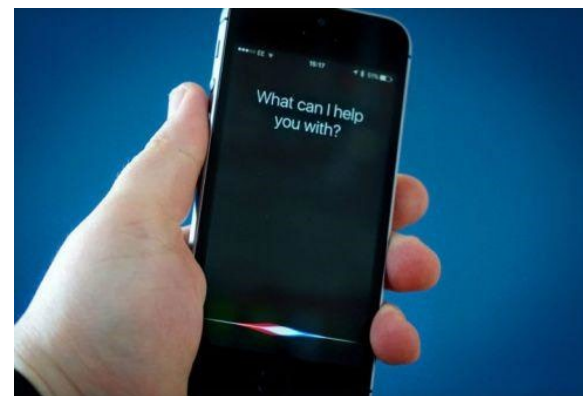


JOINT INSTITUTE
交大密西根学院

NLP



Embedded System

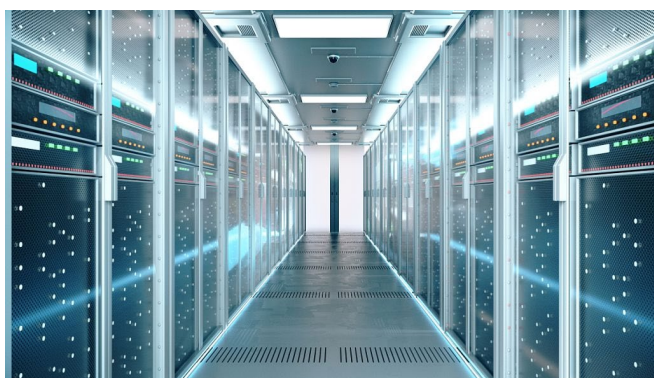


16000 words / s

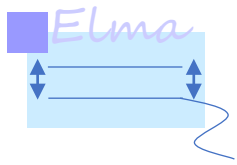
0.2 s

Large computation

High performance



Cloud



JOINT INSTITUTE
交大密西根学院

2. Design Specifications



Design Specifications



Customers Requirements

Functionality

Question answering on embedded systems

Similar accuracy

Efficiency

Network latency

Faster prediction

Security

Personal sensitive data

A safer method

Cost

Need network connection

Work without network



Design Specifications



Engineering Specifications

Functionality

Train on the cloud and predict locally

Offloading

Efficiency

The whole model for prediction is large

Early exit

Security

Encryption and decryption consume resources

Simplify

Cost

Storage and computation on the cloud

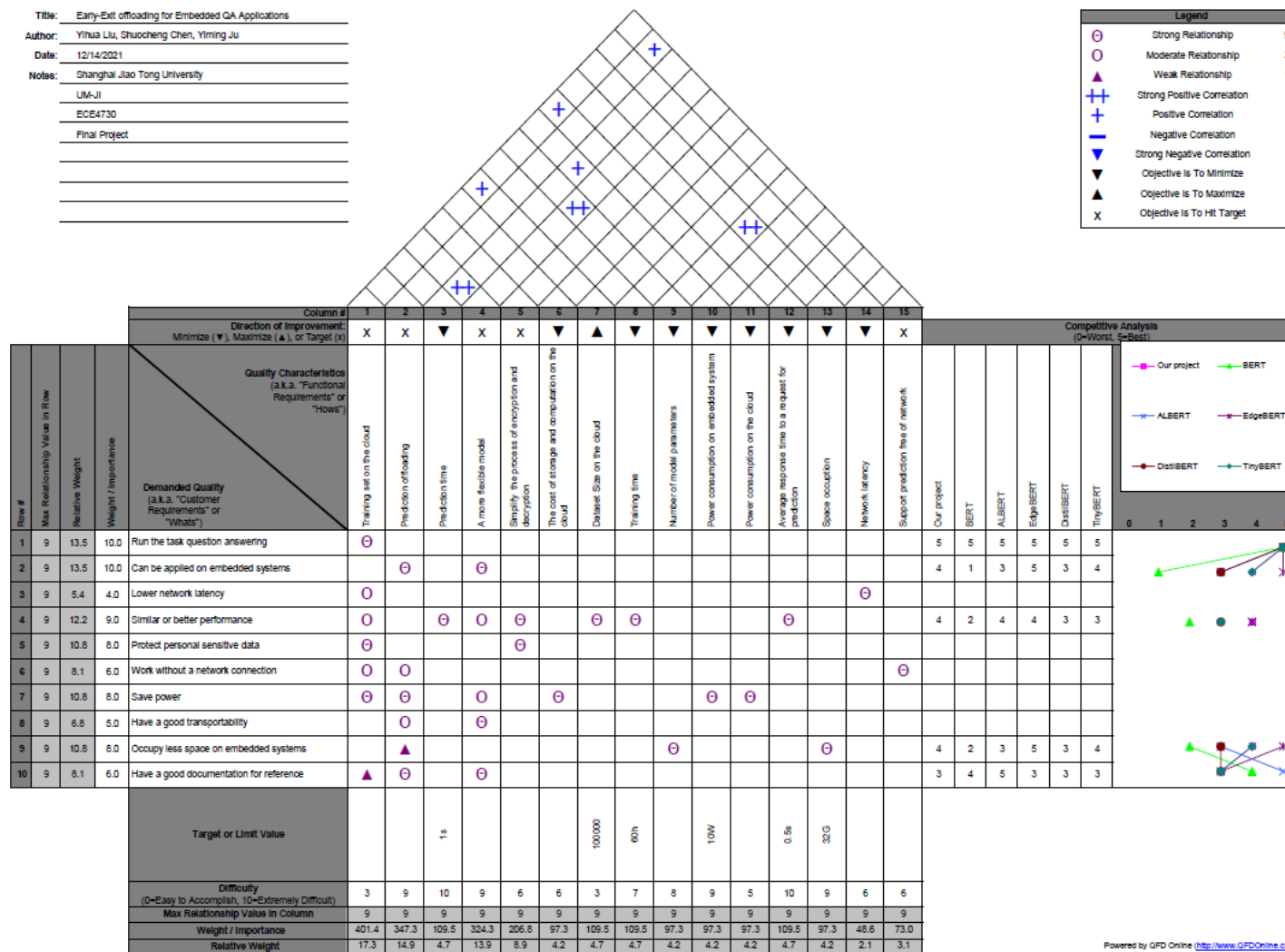
Work remotely

Design Specifications



JOINT INSTITUTE
交大密西根学院

Title: Early-Exit offloading for Embedded QA Applications
Author: Yihua Liu, Shuocheng Chen, Yiming Ju
Date: 12/14/2021
Notes: Shanghai Jiao Tong University
UM-JI
ECE4730
Final Project



3. Concept Generation & Selection



Concept Generation & Selection



JOINT INSTITUTE
交大密西根学院

Dataset selection

SQuAD2.0
The Stanford Question Answering Dataset

Environment Selection



Develop Kit Selection



Concept Generation & Selection



BERT vs. ALBERT Tensorflow vs. PyTorch

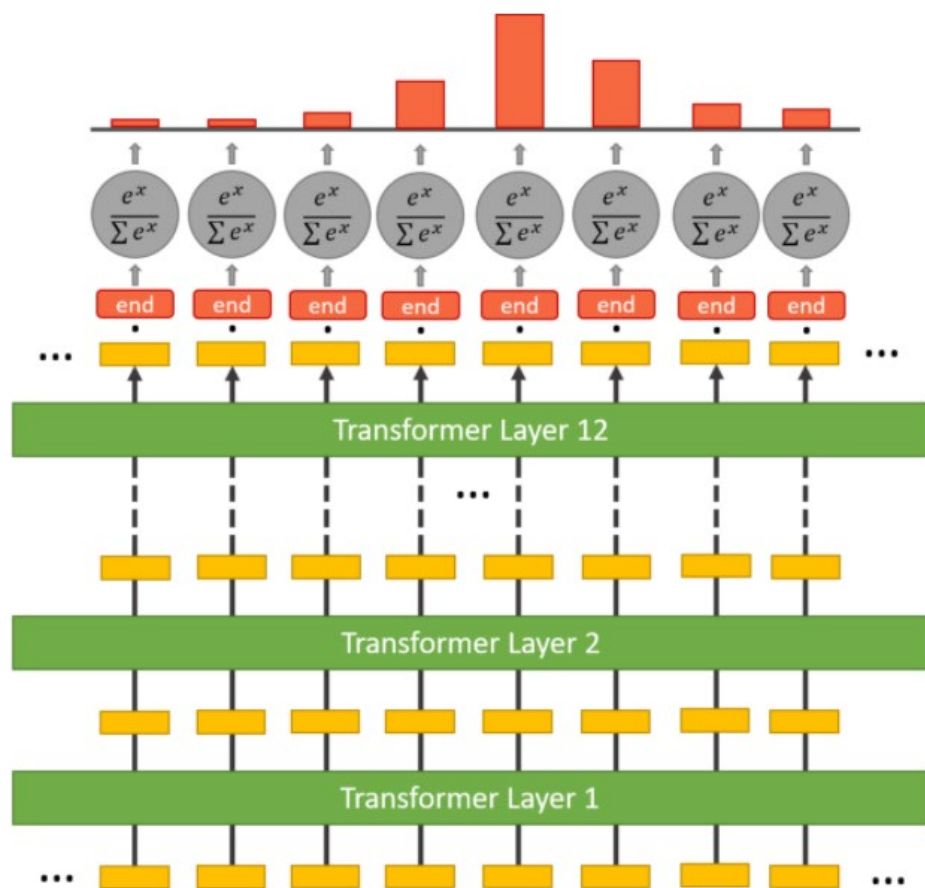
Criterion	Weight(%)	BERT				ALBERT			
		Tensorflow		PyTorch		Tensorflow		PyTorch	
		Score	Rating	Score	Rating	Score	Rating	Score	Rating
Performance	0.35	8	2.8	8	2.8	9	3.15	9	3.15
Space occupation	0.35	4	1.4	4	1.4	7	2.45	7	2.45
Flexibility	0.2	3	0.6	8	1.6	3	0.6	8	1.6
Derivative	0.1	7	0.7	5	0.5	7	0.7	5	0.5
Total			5.5		6.3		6.9		7.7

ALBERT + Pytorch

Concept Generation & Selection



ALBERT



Parameter Sharing



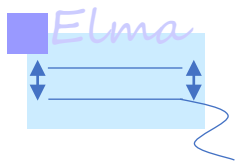
Save memory space

Embedding Factorization

Sentence Order Prediction



Better performance



JOINT INSTITUTE
交大密西根学院

4. Design Description

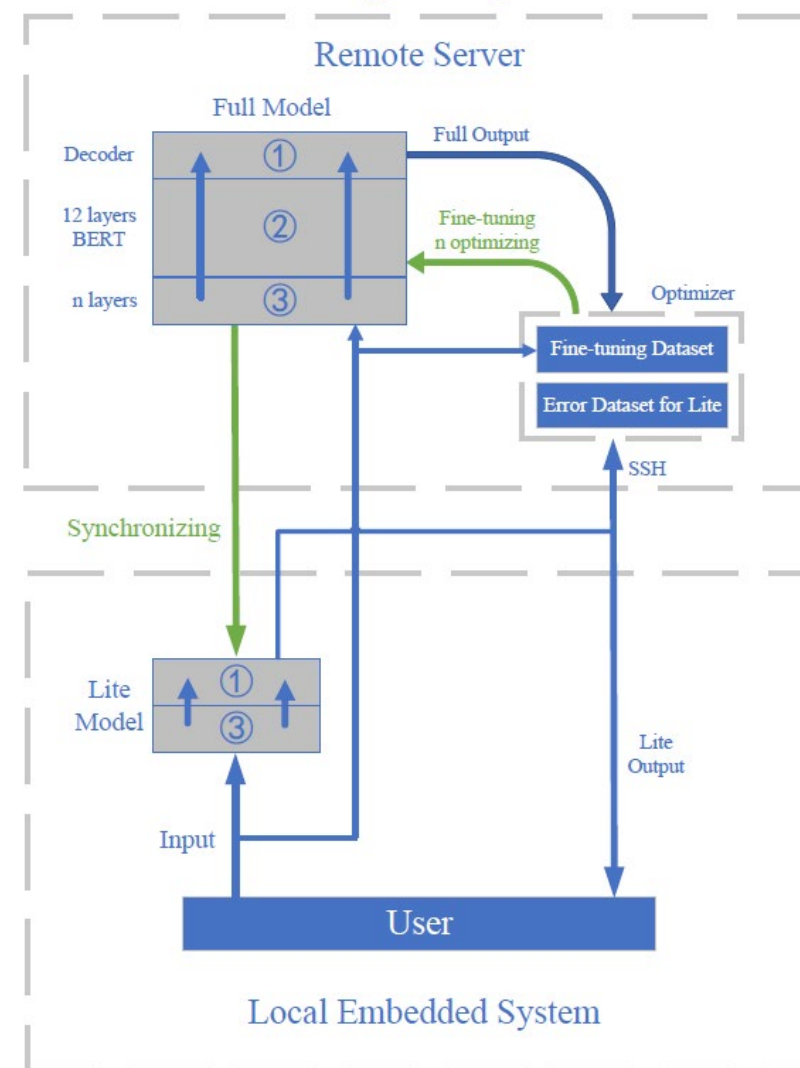
Design Description



Remote Server: Full model

Embedded System: Lite model

Concept Diagram



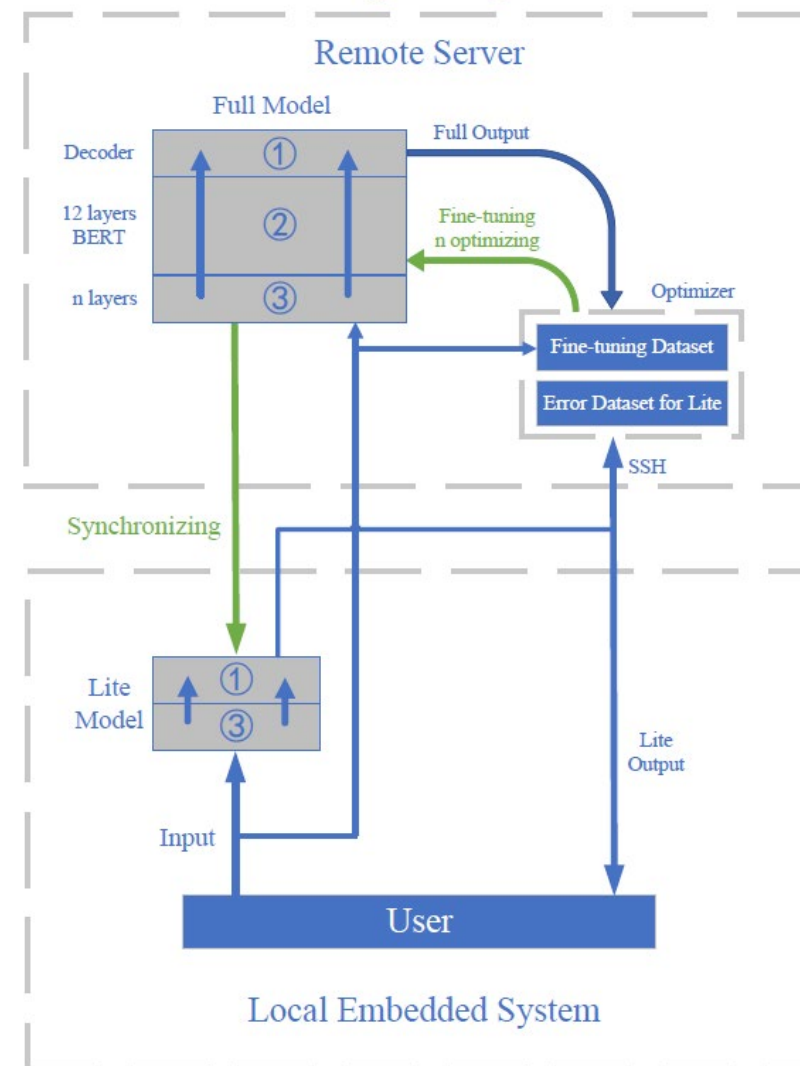
Design Description



Remote Server: Full model

- 12 layers ALBERT
- Optimizer
- Fine-tuning

Concept Diagram



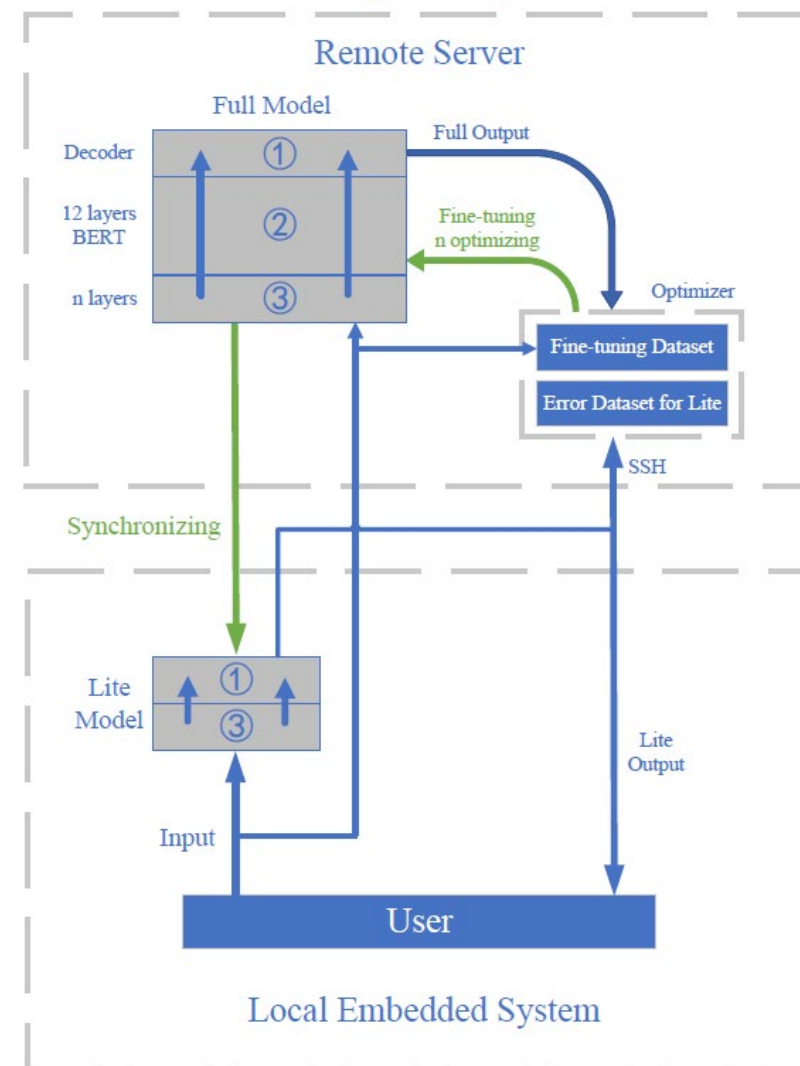
Design Description



Embedded System: Lite model

- Early-Exit
- Synchronization

Concept Diagram



Early Exit



Input: x_n (transformer output), E_T (threshold of entropy)

Output: z_N (output of last transformer) or unsolvable

For $n = 1 \dots N$ **do**

$z_n = f_{\text{exit}}(x_n)$

$y_n = \text{softmax}(z_n)$

$e_n = \text{entropy}(y_n)$

if $e < E_T$ **then**

 return $\text{argmax}(y_n)$

End if

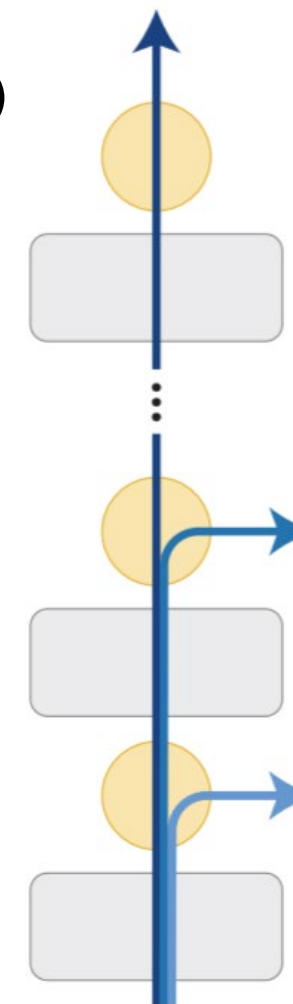
End for

If the question is solvable

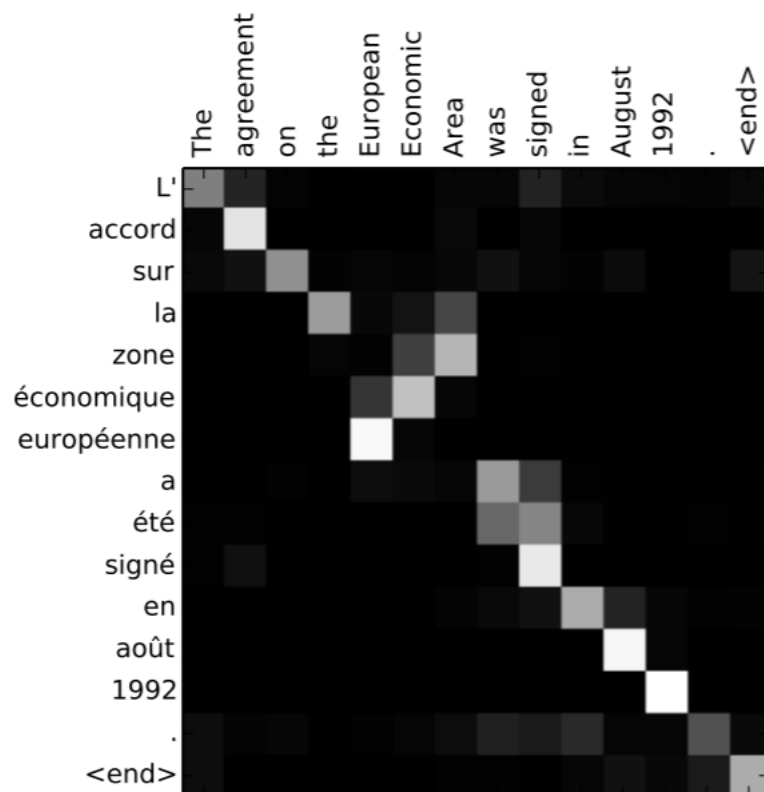
 return z_N

Else

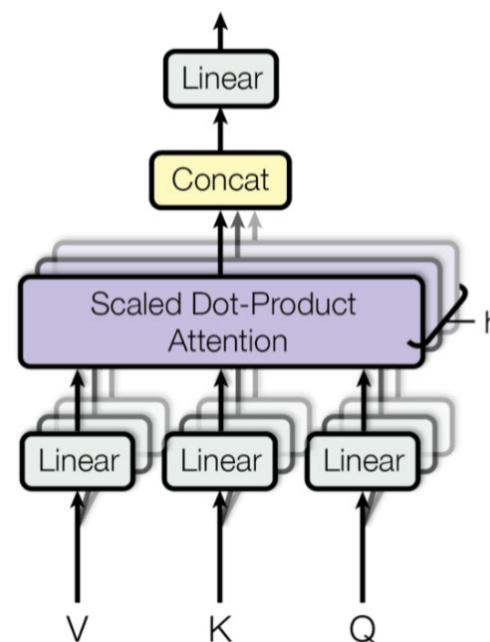
 return unsolvable question to the server



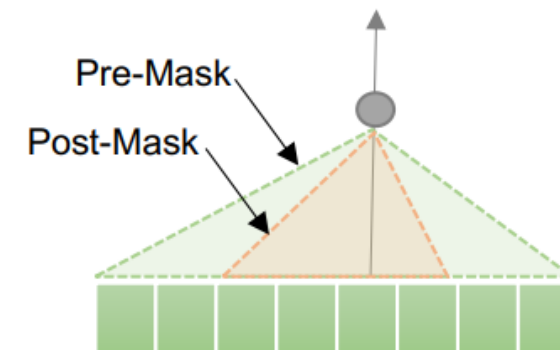
Adaptive Attention



Attention Mechanism



Multi-Head Attention



Adaptive Attention Span

Fp16 Quantization



Floating Point defined by IEEE 754

Numerical form:

$$V_{10} = (-1)^S * M * 2^E$$

Sign bit S determines whether number is negative or positive

Significand (mantissa) M usually a fractional value in range [1.0,2.0)

Exponent E weights value by a $(-/+)$ power of two

Analogous to scientific notation

Addition

$$\begin{aligned} (\pm s_1 \times b^{e_1}) + (\pm s_2 \times b^{e_2}) &= (\pm s_1 \times b^{e_1}) + (\pm s_2 / b^{e_1 - e_2}) \times b^{e_1} \\ &= (\pm s_1 \pm s_2 / b^{e_1 - e_2}) \times b^{e_1} = \pm s \times b^e \end{aligned}$$

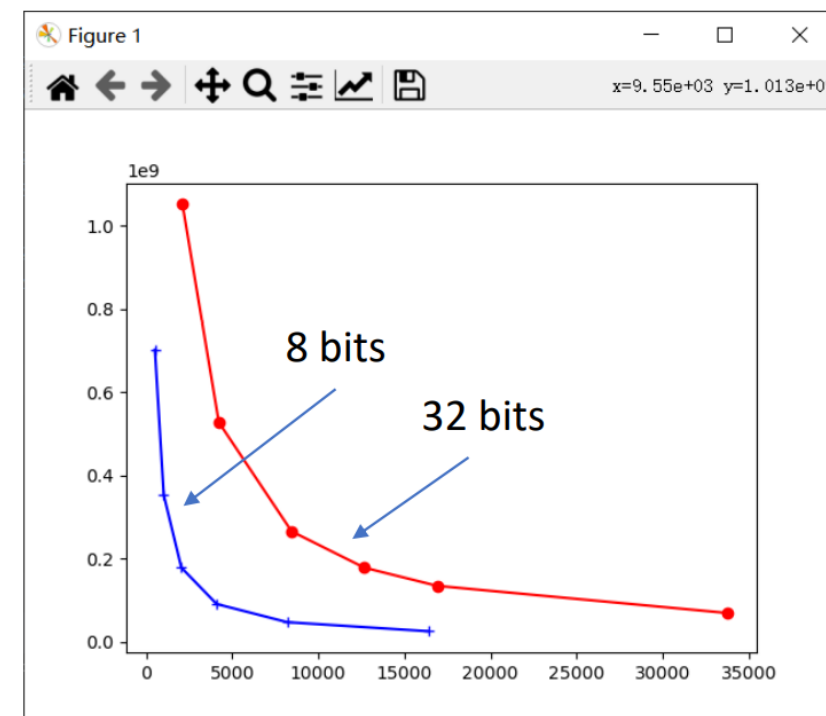
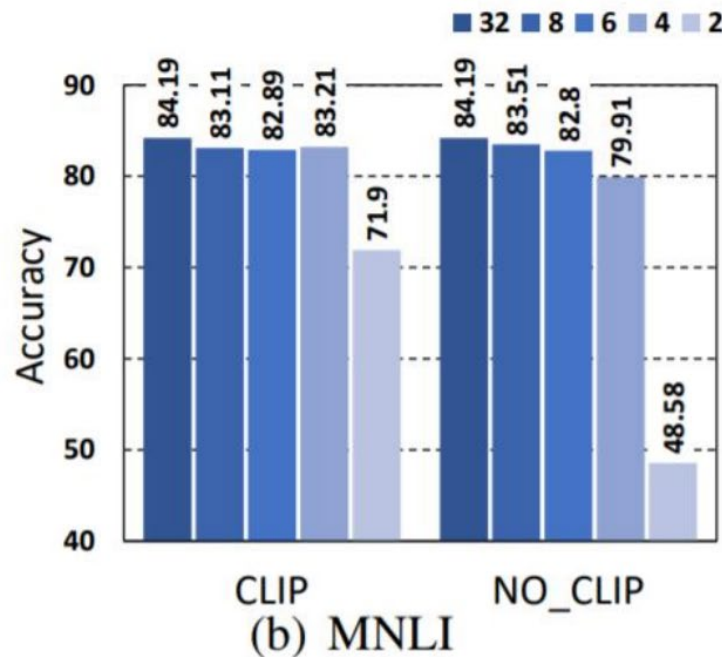
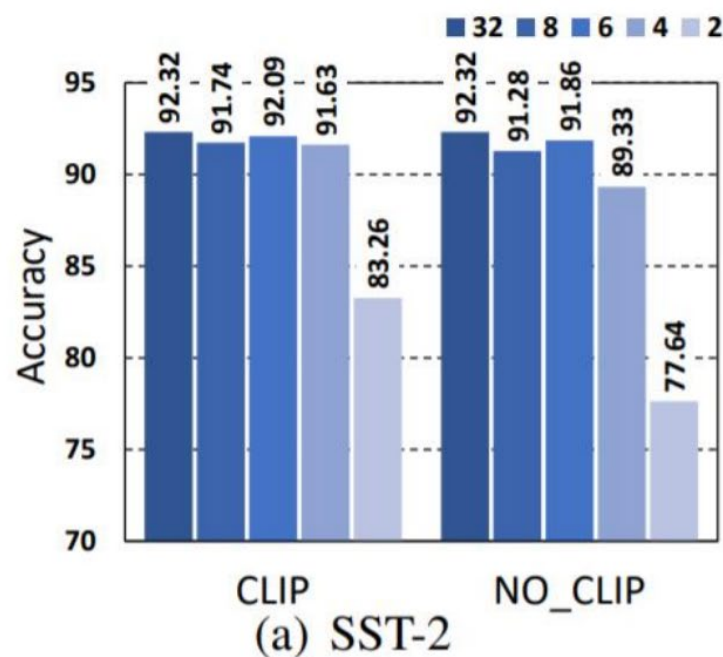
Multiplication

$$(\pm s_1 \times b^{e_1}) \times (\pm s_2 \times b^{e_2}) = (\pm s_1 \times s_2) \times b^{e_1 + e_2}$$

Division

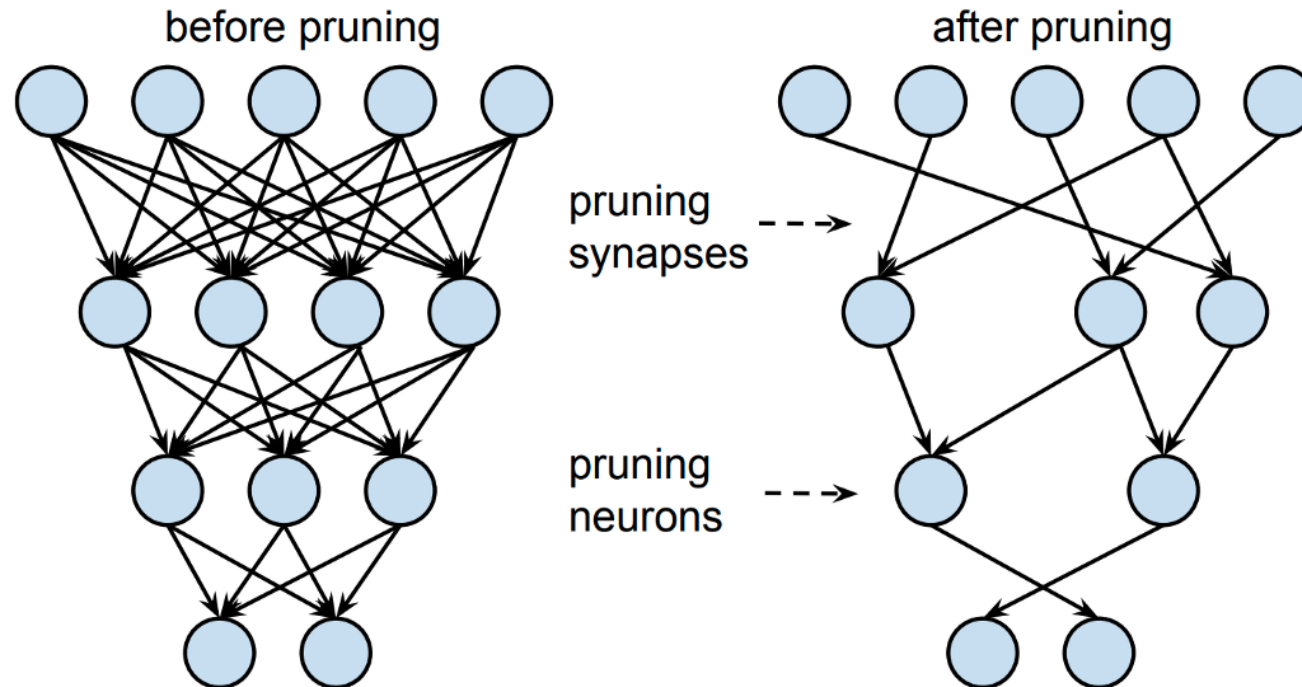
$$(\pm s_1 \times b^{e_1}) / (\pm s_2 \times b^{e_2}) = (\pm s_1 / s_2) \times b^{e_1 - e_2}$$

Fp16 Quantization





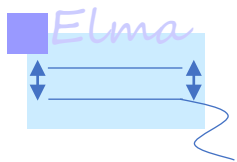
Network Pruning



Method:

Movement Pruning

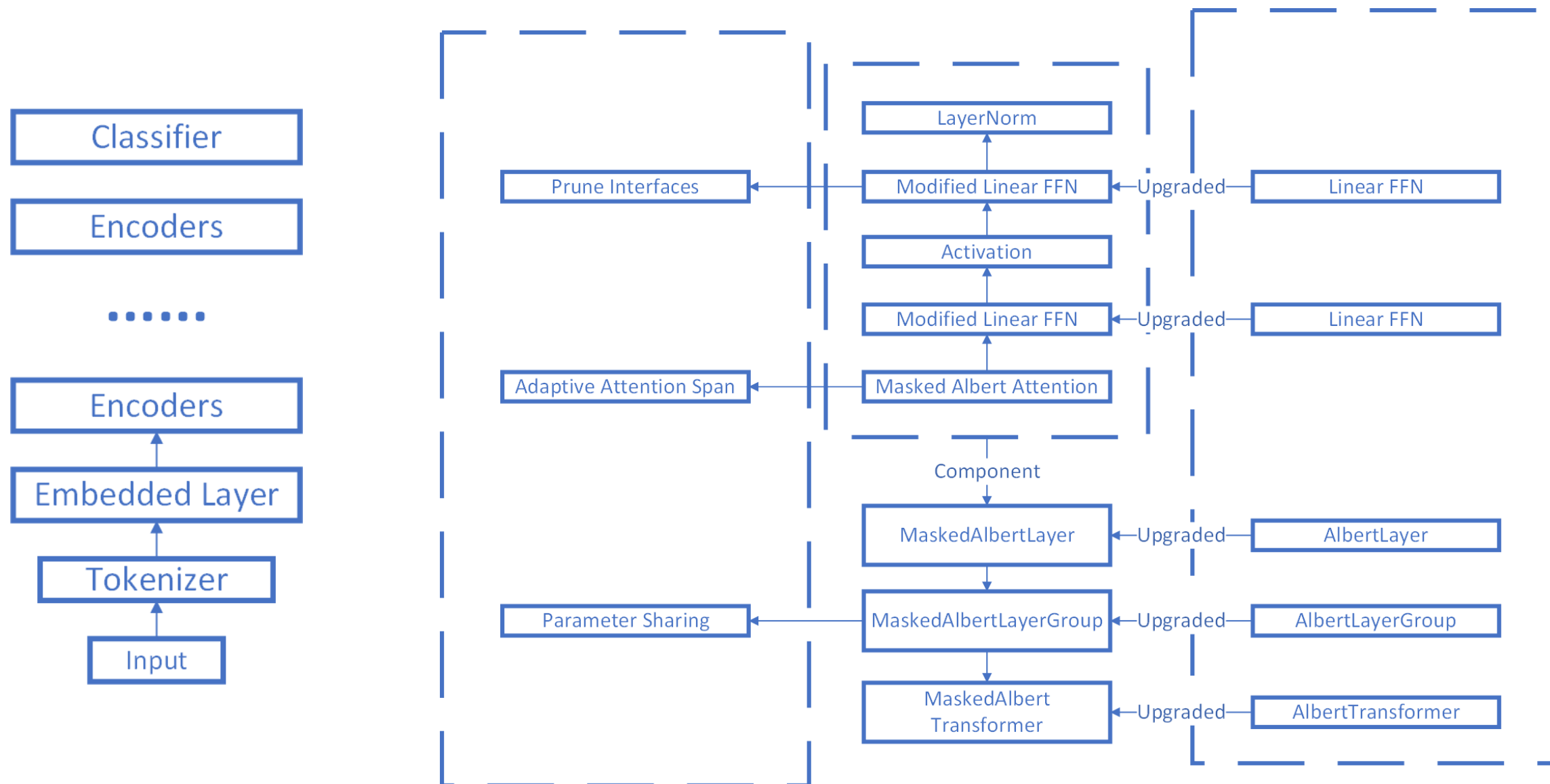
Magnitude Pruning



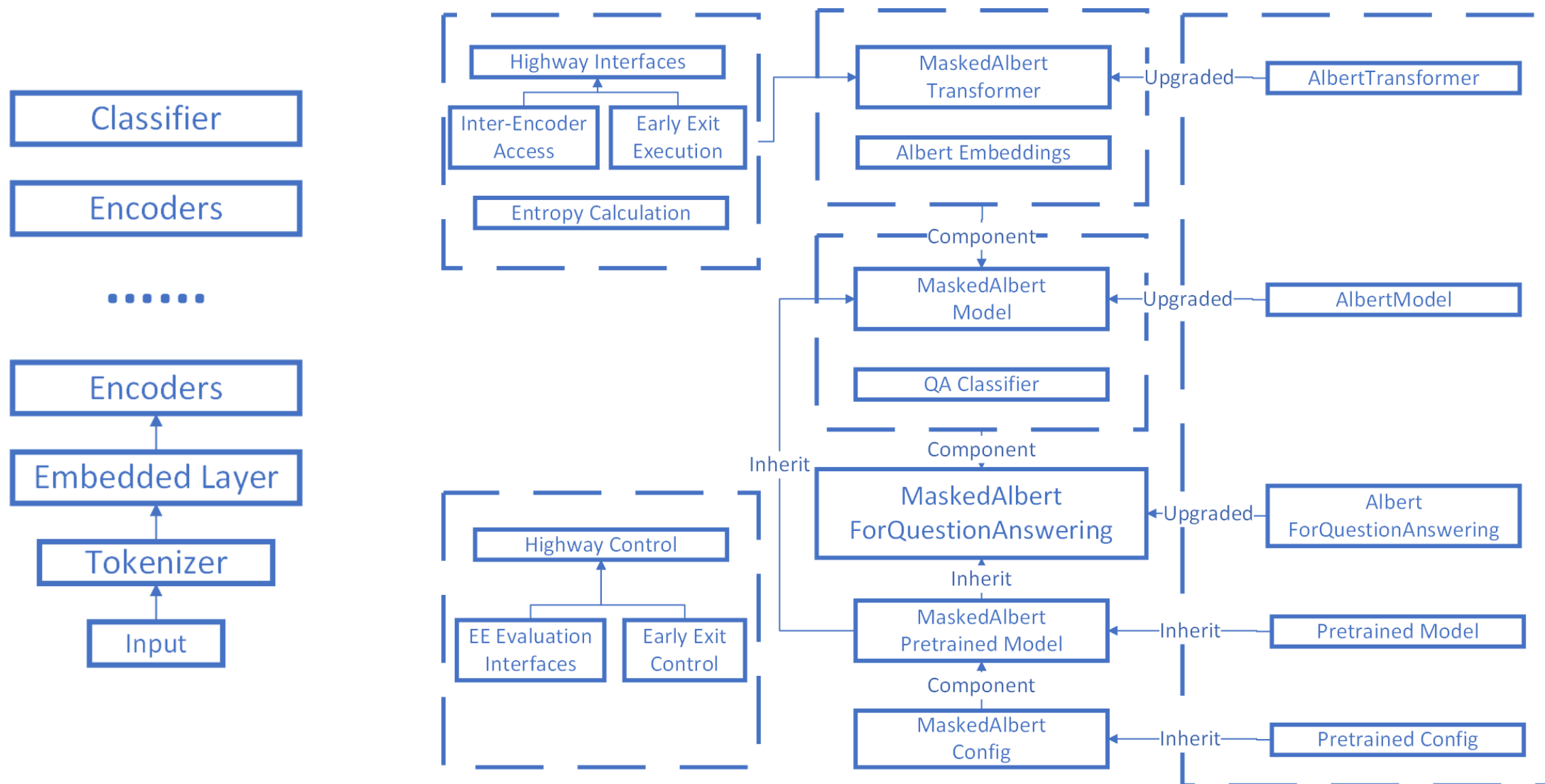
JOINT INSTITUTE
交大密西根学院

5. Implementation & Validation

Implementation



Implementation



FP16

```
if fp16:
    with amp.scale_loss(loss, optimizer) as scaled_loss:
        scaled_loss.backward()
else:
    loss.backward()
```

Magnitude Prune

```
mask = MagnitudeBinarizer.apply(inputs=tensor, threshold=threshold)
pruned_model[name] = tensor * mask
```

Validation



$\text{exact} = (\text{norm}(\text{orig answer}) == \text{norm}(\text{pred answer}))$

$\text{f1} = (2 * \text{precision} * \text{recall}) / (\text{precision} + \text{recall})$

$\text{precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}, \text{recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$

63 Aug 03, 2020	AMBERT-H (single model) ByteDance	76.710	79.659
63 Aug 03, 2020	AMBERT-S (single model) ByteDance	76.563	79.776
64 Jan 05, 2019	synss (single model) bert_finetune	76.055	79.329
65 May 21, 2021	mgrc single model	75.344	78.381
65 Apr 05, 2021	BERT-Base-L (single model) Anonymous	75.457	78.232
66 Dec 18, 2018	ARSG-BERT (single model) TRINITY RESEARCH LABS, Active.ai https://active.ai	74.746	78.227
66 Aug 29, 2020	BERT-Base-V (single model) Anonymous	75.073	77.805
66 Nov 05, 2018	MIR-MRC(F-Net) (single model) Kangwon National University, Natural Language Processing Lab. & ForceWin, KP Lab.	74.791	77.988
67 Aug 06, 2020	BERT-Base-DT (single model) Anonymous	74.769	77.706
68 Dec 03, 2020	BERT-Base-V2 single model	74.656	77.404

Metrics	Values
HasAns_exact	65.72199730094466
HasAns_f1	71.02531166019085
HasAns_total	5928
NoAns_exact	84.44070647603027
NoAns_f1	84.44070647603027
NoAns_total	5945
exact	75.09475280047165
f1	77.74261328405724
Total	11873

6. Discussion & Conclusion

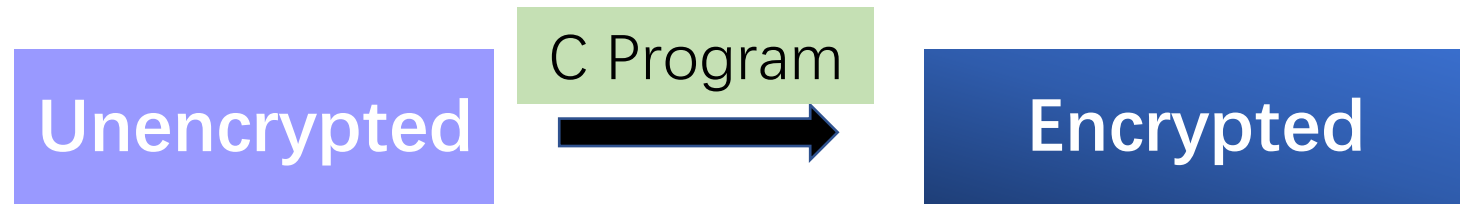


Discussion

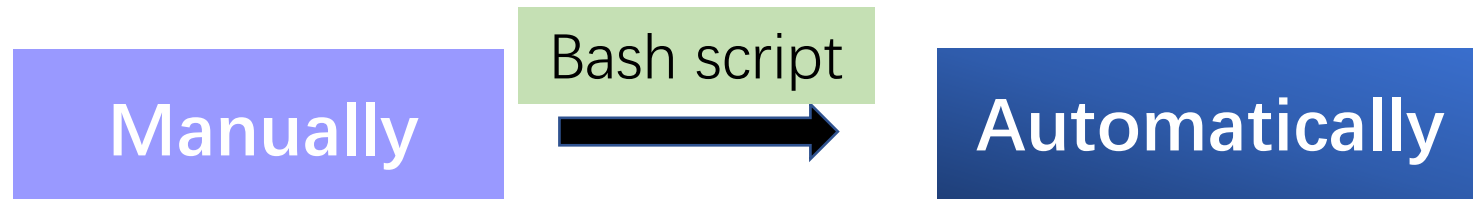


Functionality

- Encryption: Localized encryption



- Synchronization: Network communication automation



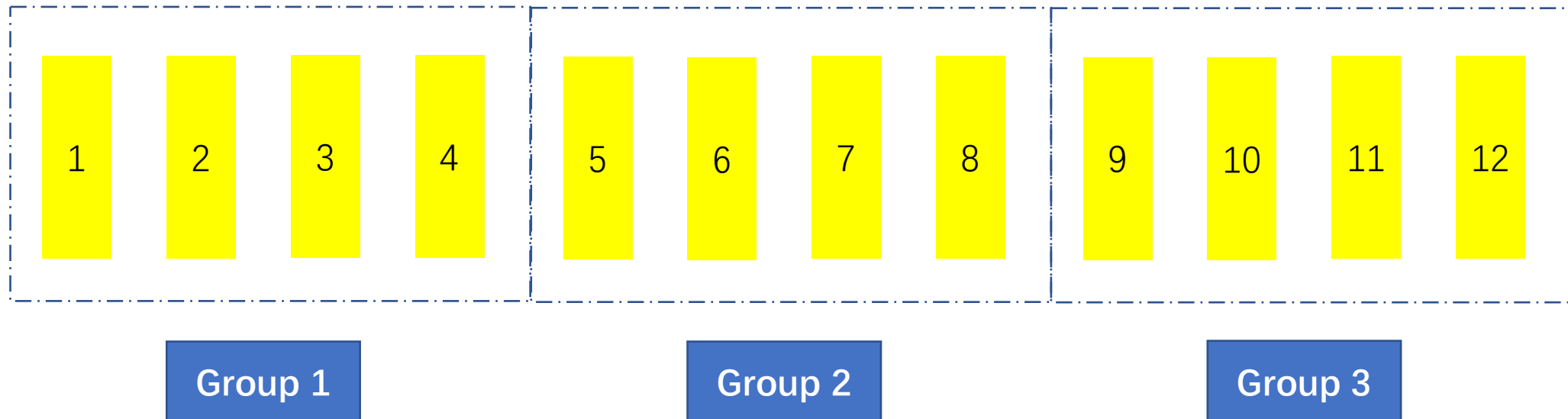


Discussion

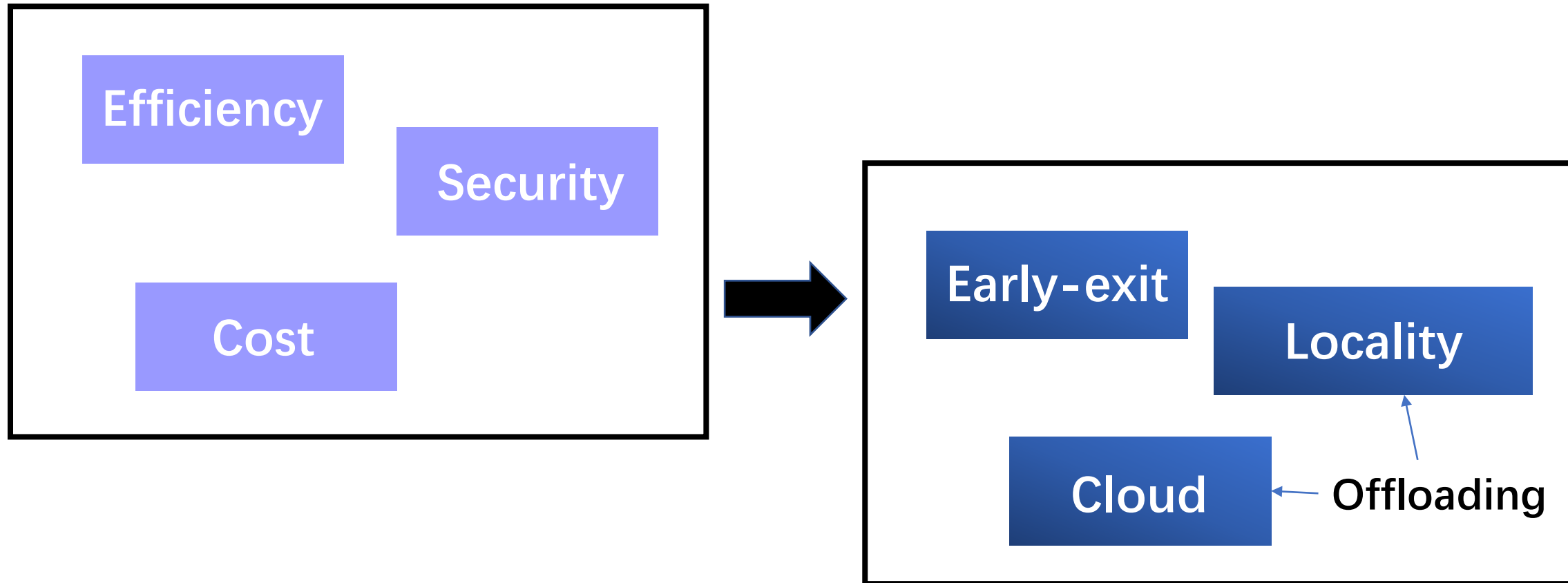


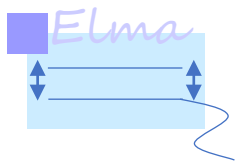
Optimization

- Distillation rate
- Choice of entropy
- Introduce a true error dataset
- Parameter sharing within each of the layer groups



Conclusion





JOINT INSTITUTE
交大密西根学院

Thanks!

ELMA: Early-Exit Offloading for Embedded Question Answering Applications

Group:	3
Instructor:	Prof. An Zou
Sponsor:	UM-SJTU Joint Institute
Group Member:	Yihua Liu, Shuocheng Chen, Yiming Ju