

# **ALBERT**

**: A LITE BERT FOR SELF-SUPERVISED LEARNING  
OF LANGUAGE REPRESENTATIONS**

# INTRODUCTION

- Memory Limitation
- Training Time
- Memory Degradation

# ELEMENTS OF ALBERT

- Factorized Embedding Parameterization
- Cross-layer Parameter Sharing
- Sentence Order Prediction

# Factorized Embedding Parameterization

- $E$  : vocabulary embedding size
- $H$  : hidden size
- $H \gg E$

# Cross-layer Parameter Sharing

- Share all parameters accross layers

# Inter-sentence Coherent Loss

- Use a sentence-order prediction(SOP) loss
  - positive example
    - : same as BERT (two consecutive segments from the same document)
  - negative example
    - : the same two consecutive segments but with their order swapped

# EXPERIMENT

## Factorized Embedding Parameterization

Model	$E$	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
ALBERT base not-shared	64	87M	89.9/82.9	80.1/77.8	82.9	91.5	66.7	81.3
	128	89M	89.9/82.8	80.3/77.3	83.7	91.5	67.9	81.7
	256	93M	90.2/83.2	80.3/77.4	84.1	91.9	67.3	81.8
	768	108M	90.4/83.2	80.4/77.6	84.5	92.8	68.2	82.3
ALBERT base all-shared	64	10M	88.7/81.4	77.5/74.8	80.8	89.4	63.5	79.0
	128	12M	89.3/82.3	80.0/77.1	81.6	90.3	64.0	80.1
	256	16M	88.8/81.5	79.1/76.3	81.5	90.3	63.4	79.6
	768	31M	88.6/81.5	79.2/76.6	82.0	90.6	63.3	79.8

Table 3: The effect of vocabulary embedding size on the performance of ALBERT-base.

# EXPERIMENT Sentence Order Prediction

SP tasks	Intrinsic Tasks			Downstream Tasks					Avg
	MLM	NSP	SOP	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	
None	54.9	52.4	53.3	88.6/81.5	78.1/75.3	81.5	89.9	61.7	79.0
NSP	54.5	90.5	52.0	88.4/81.5	77.2/74.6	81.6	<b>91.1</b>	62.3	79.2
SOP	54.0	78.9	86.5	<b>89.3/82.3</b>	<b>80.0/77.1</b>	<b>82.0</b>	90.3	<b>64.0</b>	<b>80.1</b>

Table 5: The effect of sentence-prediction loss, NSP vs. SOP, on intrinsic and downstream tasks.



# EXPERIMENT

## Cross-layer Parameter Sharing

	Model	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
ALBERT base $E=768$	all-shared	31M	88.6/81.5	79.2/76.6	82.0	90.6	63.3	79.8
	shared-attention	83M	89.9/82.7	80.0/77.2	84.0	91.4	67.7	81.6
	shared-FFN	57M	89.2/82.1	78.2/75.4	81.5	90.8	62.6	79.5
	not-shared	108M	90.4/83.2	80.4/77.6	84.5	92.8	68.2	82.3
ALBERT base $E=128$	all-shared	12M	89.3/82.3	80.0/77.1	82.0	90.3	64.0	80.1
	shared-attention	64M	89.9/82.8	80.7/77.9	83.4	91.9	67.6	81.7
	shared-FFN	38M	88.9/81.6	78.6/75.6	82.3	91.7	64.4	80.2
	not-shared	89M	89.9/82.8	80.3/77.3	83.2	91.5	67.9	81.6