

Model	Hidden Size	Parameters	RACE (Accuracy)
BERT-large (Devlin et al., 2019)	1024	334M	72.0%
BERT-large (ours)	1024	334M	73.9%
BERT-xtlge (ours)	2048	1270M	54.3%

Table 1: Increasing hidden size of BERT-large leads to worse performance on RACE.

效果
暴力增加隐藏层的数量，翻倍隐藏层
缺点
容量产生梯度消失
计算量非常大

Given the importance of model size, we ask: Is having better NLP models as easy as having larger models? 已经知道模型大小很重要，怎样建造更好的模型

提出新的方案
想办法保持效果又减少参数数量
优化损失函数 introduce a self-supervised loss for sentence-order prediction (SOP)
We establish new state-of-the-art results on the well-known GLUE, SQuAD, and RACE benchmarks for natural language understanding. Specifically, we push the RACE accuracy to 89.4%, the GLUE benchmark to 89.4, and the F1 score of SQuAD 2.0 to 92.2.

分成两部分
By decomposing the large vocabulary embedding matrix into two small matrices, we separate the size of the hidden layers from the size of vocabulary embedding. This separation makes it easier to grow the hidden size without significantly increasing the parameter size of the vocabulary embeddings

SOP primary focuses on inter-sentence coherence and is designed to address the ineffectiveness (Yang et al., 2019; Liu et al., 2019) of the next sentence prediction (NSP) loss proposed in the original BERT.

The backbone of the ALBERT architecture is similar to BERT in that it uses a transformer en -coder (Vaswani et al., 2017) with GELU nonlinearities (Hendrycks & Gimpel, 2016). We follow the BERT notation conventions and denote the vocabulary embedding size as E , the number of encoder layers as L , and the hidden layer size as H .

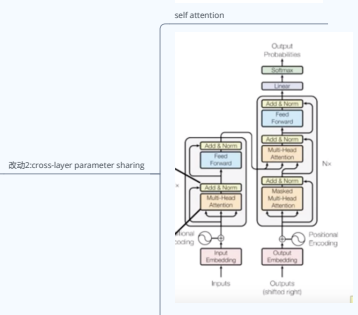
改动1: Factorized embedding parameterization
bert不好的地方: embedding size E is tied with the hidden layer size H , i.e., $E = H$
理由: From a modeling perspective, WordPiece embeddings are meant to learn context-independent representations, whereas hidden-layer embeddings are meant to learn context-dependent representations. 这个 embedding 更多的提供语义，而隐藏层的 embedding 更多地包含上下文的信息
 q^i : query (to match others)
 $q^i = W^q q^i$
 k^i : key (to be matched)
 $k^i = W^k q^i$
 v^i : value to be extracted
 $v^i = W^v q^i$

untying the WordPiece embedding size E from the hidden layer size H allows us to make a more efficient usage of the total model parameters as informed by modeling needs, which dictate that $H \gg E$.

Instead of projecting the one-hot vectors directly into the hidden space of size H , we first project them into a lower dimensional embedding space of size E , and then project it to the hidden space. By using this decomposition, we reduce the embedding parameters from $O(V \times H)$ to $O(V \times E + E \times H)$

THE ELEMENTS OF ALBERT

A little bert



NSP is a binary classification loss for predicting whether two segments appear consecutively in the original text, as follows: positive examples are created by taking consecutive segments from the training corpus; negative examples are created by pairing segments from different documents; positive and negative examples are sampled with equal probability. The NSP objective was designed to improve performance on downstream tasks, such as natural language inference, that require reasoning about the relationship between sentence pairs.

问题: negative sample 不够negative, 两个样本可能都不是同一topic

改进方式: e maintain that inter-sentence modeling is an important aspect of language understanding, but we propose a loss based primarily on coherence. That is, for ALBERT, we use a sentence-order pre-diction (SOP) loss, which avoids topic prediction and instead focuses on modeling inter-sentence coherence. The SOP loss uses as positive examples the same technique as BERT (two consecutive segments from the same document), and as negative examples the same two consecutive segments but with their order swapped.

Model	Parameters	Layers	Hidden	Embedding	Parameter-sharing
BERT	18M	12	768	768	False
large	334M	24	1024	1024	False
xtlge	1270M	24	2048	2048	False
ALBERT	18M	12	768	128	True
large	60M	24	2048	128	True
xtlge	235M	24	4096	128	True

Table 2: The configurations of the main BERT and ALBERT models analyzed in this paper.

Note that for ALBERT-xtlge, we mainly report results on a 12-layer network because a 24-layer network (with the same configuration) obtains similar results but is computationally more expensive

对比对象

downstream下游的结果: we evaluate our models on three popular bench-marks: The General Language Understanding Evaluation (GLUE) benchmark (Wang et al., 2018), two versions of the Stanford Question Answering Dataset (SQuAD; Rajpurkar et al., 2016, 2018), and the Reading Comprehension from Examinations (RACE)

对比结果

Model	E	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
ALBERT	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
large	128	60M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
xtlge	256	1270M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
ALBERT	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
large	128	60M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
xtlge	256	1270M	89.082.9	80.077.2	82.9	91.5	66.7	80.2

F1 and EM

Model	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg	Speedup
BERT	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.77x
large	334M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.0
xtlge	1270M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.0
ALBERT	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.2x
large	60M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.2x
xtlge	235M	89.082.9	80.077.2	82.9	91.5	66.7	80.2	1.2x

对比试验

Model	E	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
ALBERT	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
shared-attention	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
shared-FPN	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
ALBERT	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
shared-attention	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2
shared-FPN	64	18M	89.082.9	80.077.2	82.9	91.5	66.7	80.2

Table 5: The effect of cross-layer parameter-sharing strategies. ALBERT-base configuration.

most of the performance drop appears to come from sharing the FFN-layer parameters, while sharing the attention parameters results in no drop when $E = 128$ (+0.1 on Avg), and a slight drop when $E = 768$ (-0.7 on Avg).

共享参数

SP tasks	Intrinsic Tasks	NSP	SOP	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
NSP	54.5	90.5	52.0	88.481.5	77.274.6	81.6	91.1	62.3	79.2
SOP	54.0	78.9	86.5	89.382.3	80.077.1	82.0	90.3	64.0	80.1

Table 6: The effect of sentence-prediction loss, NSP vs. SOP, on intrinsic and downstream tasks.

loss SOP

Number of layers	Parameters	SQuAD1.1	SQuAD2.0	MNLI	SST-2	RACE	Avg
1	18M	81.025.9	80.050.1	80.4	80.8	60.1	80.5
3	18M	79.809.7	64.461.7	77.7	86.7	54.0	71.2
6	18M	84.478.4	73.871.1	81.2	88.9	60.9	77.2
12	18M	89.083.3	80.077.9	83.3	91.7	66.7	81.5
24	18M	90.583.3	81.079.0	83.3	91.3	68.7	82.1
48	18M	90.083.1	81.079.9	83.4	91.9	66.9	81.4

Table 7: The effect of increasing the number of layers for an ALBERT-large configuration.

网络深度