

BERTology

Jie Tang

Department of Computer Science and Technology

Tsinghua University

May 20, 2020

Overview



2 Pre-BERT Era





5 Conclusion



Overview



2 Pre-BERT Era

3 BERT



5 Conclusion

Background

- Language model: given an sequence of length n, it assigns a probability $p(x_1, x_2, x_3, \ldots, x_n)$ to the whole sequence. The probability of standard LM can be decomposed as $p(x_1, x_2, \ldots, x_n) = \prod_{t=1}^n p(x_t | x_1, \ldots, x_{t-1}).$
- Sequence to sequence (seq2seq) learning: given an input sequence $x_1, x_2, ..., x_m$ and an output sequence $y_1, y_2, ..., y_n$, the objective of seq2seq learning is to maximize the likelihood $p(y_n, y_{n-1}, ..., y_1 | x_1, x_2, ..., x_m)$. Common seq2seq methods decompose this objective as $p(y_n, y_{n-1}, ..., y_1 | x_1, x_2, ..., x_m) = \prod_{t=1}^n p(y_t | y_{t-1}, ..., y_1; x_1, x_2, ..., x_m)$.

Pre-BERT Era

- Semi-supervised Sequence Learning
- context2vec: Learning Generic Context Embedding with Bidirectional LSTM
- Pre-trained seq2seq: Unsupervised Pretraining for Sequence to Sequence Learning
- ELMo: Deep contextualized word representations
- OpenAl GPT: Improving Language Understanding by Generative Pre-Training

Semi-supervised Sequence Learning¹

- This paper presents two approaches that use unlabeled data to improve sequence learning with recurrent networks.
- One is to predict what comes next in a sequence, which is a conventional language model.
- The other is to use a sequence autoencoder, which reads the input sequence into a vector and predicts the input sequence again.
- These two methods can be used as a "pretraining" step for a later supervised sequence learning algorithm.

¹Dai, Andrew M., and Quoc V. Le. "Semi-supervised sequence learning." Advances in neural information processing systems. 2015.

Semi-supervised Sequence Learning (cont.)

• The sequence autoencoder is inspired by seq2seq, except that it is an unsupervised learning model. The objective is to reconstruct the input sequence itself.



• The weights obtained from the sequence autoencoder can be used as an initialization of downstream LSTM networks.

$context2vec^1$

 context2vec is an unsupervised model for efficiently learning a generic context embedding function from large corpora, using bidirectional LSTM. The architecture is based on word2vec's CBOW but replaces its context modeling with LSTM.



¹Melamud, Oren, Jacob Goldberger, and Ido Dagan. "context2vec: Learning generic context embedding with bidirectional lstm." CoNLL 2016.

Pre-trained Seq2seq¹

- This work presents a general unsupervised learning method to improve the accuracy of seq2seq models.
- The weights of a seq2seq model are initialized with the pretrained weights of two language models and then fine-tuned with labeled data.
- All parameters in a shaded box are pretrained from language models.



¹Ramachandran, Prajit, Peter J. Liu, and Quoc V. Le. "Unsupervised pretraining for sequence to sequence learning." arXiv preprint arXiv:1611.02683 (2016).

ELMo^1

- ELMo introduces a new type of deep contextualized word presentation, which are functions of the internal states of a deep bidirectional language model (biLM).
- ELMo uses the concatenation of independently trained left-to-right and right-to-left LSTMs to generate features for downstream tasks.



¹Peters, Matthew E., et al. "Deep contextualized word representations." arXiv preprint arXiv:1802.05365 (2018).

ELMo (cont.)

• ELMo formulation jointly maximizes the log likelihood of the forward and backward directions:

$$\frac{\sum_{k=1}^{N} \left(\log p\left(t_{k} | t_{1}, \dots, t_{k-1}; \Theta_{x}, \vec{\Theta}_{LSTM}, \Theta_{s} \right) + \log p\left(t_{k} | t_{k+1}, \dots, t_{N}; \Theta_{x}, \overleftarrow{\Theta}_{LSTM}, \Theta_{s} \right) \right)$$
(1)

where Θ_x and Θ_s are the parameters of token representation and softmax layer.

ELMo Results

TASK	PREVIOUS SOTA		OUR BASELINE	ELMO + BASELINE	INCREASE (ABSOLUTE/ RELATIVE)
SQuAD	Liu et al. (2017)	84.4	81.1	85.8	4.7 / 24.9%
SNLI	Chen et al. (2017)	88.6	88.0	88.7 ± 0.17	0.7 / 5.8%
SRL	He et al. (2017)	81.7	81.4	84.6	3.2 / 17.2%
Coref	Lee et al. (2017)	67.2	67.2	70.4	3.2 / 9.8%
NER	Peters et al. (2017)	91.93 ± 0.19	90.15	92.22 ± 0.10	2.06 / 21%
SST-5	McCann et al. (2017)	53.7	51.4	54.7 ± 0.5	3.3 / 6.8%

OpenAl GPT^1

- OpenAI GPT uses generative pre-training of a language model on a diverse corpus of unlabeled text, followed by discriminative fine-tuning on each specific task.
- OpenAI GPT uses a left-to-right Transformer.

	Text Task Prediction Classifier	Classification	Start	Text	Extract	+ Transform	ner 🔸 L	.ine	ar	
	Layer Norm	Entailment	Start	Premise	Delim	Hypothesis	Extract	- [Transformer	+ Linear
	Feed Forward	Similarity	Start	Text 1	Delim	Text 2	Extract	- [Transformer	++ Linear
12x —	Layer Norm		Start	Text 2	Delim	Text 1	Extract	-	Transformer	۲Ÿ
	Masked Multi Self Attention		Start	Context	Delim	Answer 1	Extract	-	Transformer	+ Linear
l		Multiple Choice	Start	Context	Delim	Answer 2	Extract	ŀ	Transformer	Linear
	Text & Position Embed		Start	Context	Delim	Answer N	Extract	ŀ	Transformer	+ Linear

¹Radford, Alec, et al. "Improving language understanding by generative pre-training." URL https://s3-us-west-2. amazonaws. com/openai-assets/researchcovers/languageunsupervised/language understanding paper. pdf (2018).

OpenAl (cont.)

• GPT uses a standard language modeling for pre-training

$$\sum_{i} \log P\left(u_{i}|u_{i-k},\dots,u_{i-1};\Theta\right)$$
(2)

where P is modeled using a multi-layer Transformer

$$h_{0} = UW_{e} + W_{p}$$

$$h_{l} = \text{transformer_block} (h_{l-1}) \forall i \in [1, n]$$

$$P(u) = \text{softmax} \left(h_{n}W_{e}^{T}\right)$$
(3)

where W_e is the token embedding and W_p is the position embedding.

Results on GLUE

Method	Classif	ication	Seman	GLUE		
	CoLA (mc)	SST2 (acc)	MRPC (F1)	STSB (pc)	QQP (F1)	
Sparse byte mLSTM [16]	-	93.2	-	-	-	-
TF-KLD [23]	-	-	86.0	-	-	-
ECNU (mixed ensemble) [60]	-	-	-	<u>81.0</u>	-	-
Single-task BiLSTM + ELMo + Attn [64] Multi-task BiLSTM + ELMo + Attn [64]	$\frac{35.0}{18.9}$	90.2 91.6	80.2 83.5	55.5 72.8	<u>66.1</u> 63.3	64.8 68.9
Finetuned Transformer LM (ours)	45.4	91.3	82.3	82.0	70.3	72.8

Overview



2 Pre-BERT Era





5 Conclusion

BERT^1

- BERT is designed to pre-train deep bidirectional representations from unlabeled text by jointly conditioning on both left and right context in all layers.
- The pre-trained BERT model can be fine-tuned with just one additional output layer to create state-of-the-art models for a wide range of tasks. range of tasks



¹Devlin, Jacob, et al. "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805 (2018).

BERT (cont.)

- Overall pre-training and fine-tuning procedures for BERT.
- BERT uses two unsupervised tasks: masked language model (MLM) and next sentence prediction (NSP).



BERT Results

Results on GLUE

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.8	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	87.4	91.3	45.4	80.0	82.3	56.0	75.1
BERTBASE	84.6/83.4	71.2	90.5	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	92.7	94.9	60.5	86.5	89.3	70.1	82.1

Results on SQuAD

System	D	ev	Te	st	System	D	ev	Те	st
	EM	F1	EM	FI	bystem	EM	E1	EM	D 1
Top Leaderboard System	s (Dec	10th,	2018)			ENI	гі	ENI	гі
Human	-	-	82.3	91.2	Top Leaderboard Systems	(Dec	10th	2018)	
#1 Ensemble - nlnet	-	-	86.0	91.7	Human	06.2	200.0	2010)	00.4
#2 Ensemble - QANet	-	-	84.5	90.5	Human	80.5	89.0	80.9	89
D.I.V.I					#1 Single - MIR-MRC (F-Net)	-	-	74.8	78.0
Publishe	ed				#2 Single - nlnet	-	-	74.2	77
BiDAF+ELMo (Single)	-	85.6	-	85.8	#2 Shigie - hiller	-	-	/4.2	//
R.M. Reader (Ensemble)	81.2	87.9	82.3	88.5	Publishe	d			
Ours					unet (Ensemble)	-	-	71.4	74.9
BERT _{BASE} (Single)	80.8	88.5	-	-	SLOA (Simple)			71 4	74
BERTLARGE (Single)	84.1	90.9	-	-	SLQA+ (Single)	-		/1.4	/4.4
BERTLARGE (Ensemble)	85.8	91.8	-	-	0				
BERTLARGE (Sgl.+TriviaOA)	84.2	91.1	85.1	91.8	Ours				
BERTLARGE (Ens.+TriviaQA)	86.2	92.2	87.4	93.2	BERT _{LARGE} (Single)	78.7	81.9	80.0	83.
DERTEARDE (EIIST TITTING)	0012		0/11						

(a) SQuAD 1.1

(b) SQuAD 2.0

Overview



2 Pre-BERT Era

3 BERT



5 Conclusion

Post-BERT Era

- RoBERTa (2019)
- XLNet (2019)
- ERNIE (Tsinghua) (2019)
- ERINE (Baidu) (2019)
- ALBERT (2019)
- ELECTRA (2020)

$RoBERTa^1$

- RoBERTa: A Robustly Optimized BERT Pretraining Approach
- RoBERTa finds that BERT is under-trained due to hyperparameters and training set.
- BERT can be improved by:
 - longer inputs (only full length sentence)
 - larger batch size and learning rate
 - larger dataset
 - No next sentence prediction
 - Dynamic masking

¹Liu, Yinhan, et al. "Roberta: A robustly optimized bert pretraining approach." arXiv preprint arXiv:1907.11692 (2019).

RoBERTa Results

• Results on SQuAD

Madal	SQu/	D 1.1	SQuAD 2.0			
Model	EM	F1	EM	F1		
Single models	on dev	, w/o da	ita augm	entation		
BERTLARGE	84.1	90.9	79.0	81.8		
XLNet _{LARGE}	89.0	94.5	86.1	88.8		
RoBERTa	88.9	94.6	86.5	89.4		
Single models	on test	t (as of .	July 25, 2	2019)		
XLNet LARGE			86.3†	89.1^{+}		
RoBERTa			86.8	89.8		
XLNet + SG-	Net Ve	rifier	87.0^{\dagger}	89.9 †		

• Results on GLUE

	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
Single-task si	ngle models	on dev								
BERTLARGE	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet LARGE	89.8/-	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa	90.2/90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	91.3	-
Ensembles on	test (from l	eaderboa	rd as of	July 25,	2019)					
ALICE	88.2/87.9	95.7	90.7	83.5	95.2	92.6	68.6	91.1	80.8	86.3
MT-DNN	87.9/87.4	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2/89.8	98.6	90.3	86.3	96.8	93.0	67.8	91.6	90.4	88.4
RoBERTa	90.8/90.2	98.9	90.2	88.2	96.7	92.3	67.8	92.2	89.0	88.5

XLNet^1

- We first review and compare the conventional auto-regressive language modeling and BERT for language pretraining.
- The objective of auto-regressive language modeling:

$$\max_{\theta} \log p_{\theta}(\mathbf{x}) = \sum_{t=1}^{T} \log p_{\theta}\left(x_{t} | \mathbf{x}_{< t}\right) = \sum_{t=1}^{T} \log \frac{\exp\left(h_{\theta}\left(\mathbf{x}_{1:t-1}\right)^{\top} e\left(x_{t}\right)\right)}{\sum_{x'} \exp\left(h_{\theta}\left(\mathbf{x}_{1:t-1}\right)^{\top} e\left(x'\right)\right)}$$
(4)

• The objective of BERT is to reconstruct masked tokens $\overline{\mathbf{x}}$ from a corrupted version $\hat{\mathbf{x}}:$

$$\max_{\theta} \log p_{\theta}(\overline{\mathbf{x}}|\hat{\mathbf{x}}) \approx \sum_{t=1}^{T} m_{t} \log p_{\theta}\left(x_{t}|\hat{\mathbf{x}}\right) = \sum_{t=1}^{T} m_{t} \log \frac{\exp\left(H_{\theta}(\hat{\mathbf{x}})_{t}^{\top} e\left(x_{t}\right)\right)}{\sum_{x'} \exp\left(H_{\theta}(\hat{\mathbf{x}})_{t}^{\top} e\left(x'\right)\right)}$$
(5)

¹Yang, Zhilin, et al. "Xlnet: Generalized autoregressive pretraining for language understanding." Advances in neural information processing systems. 2019.

XLNet (cont.)

 Instead of using a fixed forward or backward factorization order as in conventional Auto-Regressive models, XLNet maximizes the expected log likelihood of a sequence w.r.t. all possible permutations of the factorization order.

$$\max_{\theta} \quad \mathbb{E}_{\mathbf{z} \sim \mathcal{Z}_T} \left[\sum_{t=1}^T \log p_{\theta} \left(x_{z_t} | \mathbf{x}_{\mathbf{z}_{< t}} \right) \right]$$
(6)

where \mathcal{Z}_T is the set of all possible permutations of the length-T index sequence.



XLNet Results

• Results on SQuAD

SQuAD2.0	EM	F1	SQuAD1.1	EM	F1
Dev set results (single mod	lel)			
BERT [10]	78.98	81.77	BERT† [10]	84.1	90.9
RoBERTa [21]	86.5	89.4	RoBERTa [21]	88.9	94.6
XLNet	87.9	90.6	XLNet	89.7	95.1
Test set results o	n leaderba	oard (singl	e model, as of Dec	: 14, 2019)	
BERT [10]	80.005	83.061	BERT [10]	85.083	91.835
RoBERTa [21]	86.820	89.795	BERT* [10]	87.433	93.294
XLNet	87.926	90.689	XLNet	89.898 ‡	95.080 ‡

• Results on GLUE

Model	MNLI	QNLI	QQP	RTE	SST-2	MRPC	CoLA	STS-B	WNLI
Single-task single	e models on de	v?v							
BERT [2]	86.6/-	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-
RoBERTa [21]	90.2/90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	-
XLNet	90.8/90.8	94.9	92.3	85.9	97.0	90.8	69.0	92.5	-
Multi-task ensem	bles on test (fi	om leade	rboard as	s of Oct 2	28, 2019)				
MT-DNN* [20]	87.9/87.4	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0
RoBERTa [*] [21]	90.8/90.2	98.9	90.2	88.2	96.7	92.3	67.8	92.2	89.0
XLNet*	90.9/90.9 [†]	99.0 †	90.4 †	88.5	97.1 [†]	92.9	70.2	93.0	92.5

$\mathsf{ERNIE}\;(\mathsf{Tsinghua})^1$

- ERNIE: Enhanced Language Representation with Informative Entities
- ERNIE utilizes both large-scale textual corpora and knowledge graphs to train an enhanced language representation model, which can take full advantage of lexical, syntactic, and knowledge information simultaneously.



¹Zhang, Zhengyan, et al. "ERNIE: Enhanced language representation with informative entities." arXiv preprint arXiv:1905.07129 (2019).

$\mathsf{ERNIE} (\mathsf{Baidu})^1$

- ERNIE: Enhanced Representation through Knowledge Integration
- ERNIE is designed to learn language representation enhanced by knowledge masking strategies, which includes entity-level masking and phrase-level masking.



¹Sun, Yu, et al. "Ernie: Enhanced representation through knowledge integration." arXiv preprint arXiv:1904.09223 (2019).

ALBERT^1

- ALBERT uses two parameter-reduction techniques to lower memory consumption and increase the training speed of BERT.
- Factorize embedding parameterization: reduce the embedding parameters from $O(V \times H)$ to $O(V \times E + E \times H)$ where $H \gg E$.
- Cross-layer parameters sharing: a default decision for ALBERT is to share all parameters across layers.
- Sentence-order prediction replaces NSP

Mod	lel	Parameters	Layers	Hidden	Embedding	Parameter-sharing
	base	108M	12	768	768	False
BERT	large	334M	24	1024	1024	False
	base	12M	12	768	128	True
AL DEDT	large	18M	24	1024	128	True
ALDERI	xlarge	60M	24	2048	128	True
	xxlarge	235M	12	4096	128	True

¹Lan, Zhenzhong, et al. "Albert: A lite bert for self-supervised learning of language representations." arXiv preprint arXiv:1909.11942 (2019).

ALBERT Results

• Results on GLUE

Models	MNLI	QNLI	QQP	RTE	SST	MRPC	CoLA	STS	WNLI	Avg
Single-task single	models on	dev								
BERT-large	86.6	92.3	91.3	70.4	93.2	88.0	60.6	90.0	-	-
XLNet-large	89.8	93.9	91.8	83.8	95.6	89.2	63.6	91.8	-	-
RoBERTa-large	90.2	94.7	92.2	86.6	96.4	90.9	68.0	92.4	-	-
ALBERT (1M)	90.4	95.2	92.0	88.1	96.8	90.2	68.7	92.7	-	-
ALBERT (1.5M)	90.8	95.3	92.2	89.2	96.9	90.9	71.4	93.0	-	-
Ensembles on test	(from lead	lerboard	as of Sep	ot. 16, 20	019)					
ALICE	88.2	95.7	90.7	83.5	95.2	92.6	69.2	91.1	80.8	87.0
MT-DNN	87.9	96.0	89.9	86.3	96.5	92.7	68.4	91.1	89.0	87.6
XLNet	90.2	98.6	90.3	86.3	96.8	93.0	67.8	91.6	90.4	88.4
RoBERTa	90.8	98.9	90.2	88.2	96.7	92.3	67.8	92.2	89.0	88.5
Adv-RoBERTa	91.1	98.8	90.3	88.7	96.8	93.1	68.0	92.4	89.0	88.8
ALBERT	91.3	99.2	90.5	89.2	97.1	93.4	69.1	92.5	91.8	89.4

• Results on SQuAD and RACE

Models	SQuAD1.1 dev	SQuAD2.0 dev	SQuAD2.0 test	RACE test (Middle/High)
Single model (from leaderboy	ard as of Sept. 23,	2019)		
BERT-large	90.9/84.1	81.8/79.0	89.1/86.3	72.0 (76.6/70.1)
XLNet	94.5/89.0	88.8/86.1	89.1/86.3	81.8 (85.5/80.2)
RoBERTa	94.6/88.9	89.4/86.5	89.8/86.8	83.2 (86.5/81.3)
UPM	-	-	89.9/87.2	- 1
XLNet + SG-Net Verifier++	-	-	90.1/87.2	-
ALBERT (1M)	94.8/89.2	89.9/87.2	-	86.0 (88.2/85.1)
ALBERT (1.5M)	94.8/89.3	90.2/87.4	90.9/88.1	86.5 (89.0/85.5)
Ensembles (from leaderboard	d as of Sept. 23, 2	019)		
BERT-large	92.2/86.2	-	-	-
XLNet + SG-Net Verifier	-	-	90.7/88.2	-
UPM	-	-	90.7/88.2	
XLNet + DAAF + Verifier	-	-	90.9/88.6	-
DCMN+	-	-	-	84.1 (88.5/82.3)
ALBERT	95.5/90.1	91.4/88.9	92.2/89.7	89.4 (91.2/88.6)

ELECTRA¹

• This paper proposes a replaced token detection task. This approach replaces some tokens with plausible alternatives sampled from a generator network, and then trains a discriminative model to predict whether each token was replaced by a generator sample or not.

$$\mathcal{L}_{\mathrm{MLM}}\left(\boldsymbol{x}, \theta_{G}\right) = \mathbb{E}\left(\sum_{i \in \boldsymbol{m}} -\log p_{G}\left(x_{i} | \boldsymbol{x}^{\mathrm{masked}}\right)\right)$$
(7)

$$\min_{\theta_{G},\theta_{D}} \sum_{\boldsymbol{x} \in \mathcal{X}} \mathcal{L}_{\text{MLM}}\left(\boldsymbol{x},\theta_{G}\right) + \lambda \mathcal{L}_{\text{Disc}}\left(\boldsymbol{x},\theta_{D}\right)$$
(8)



¹Clark, Kevin, et al. "Electra: Pre-training text encoders as discriminators rather than generators." arXiv preprint arXiv:2003.10555 (2020).

ELECTRA FLOPS

 ELECTRA substantially outperforms MLM-based methods such as BERT and XLNet given the same model size, data, and computation.



ELECTRA Results

• Results on GLUE

Model	Train FLOPs	CoLA	SST	MRPC	STS	QQP	MNLI	QNLI	RTE	WNLI	Avg.*	Score
BERT	1.9e20 (0.06x)	60.5	94.9	85.4	86.5	89.3	86.7	92.7	70.1	65.1	79.8	80.5
RoBERTa	3.2e21 (1.02x)	67.8	96.7	89.8	91.9	90.2	90.8	95.4	88.2	89.0	88.1	88.1
ALBERT	3.1e22 (10x)	69.1	97.1	91.2	92.0	90.5	91.3	-	89.2	91.8	89.0	_
XLNet	3.9e21 (1.26x)	70.2	97.1	90.5	92.6	90.4	90.9	-	88.5	92.5	89.1	-
ELECTRA	3.1e21 (1x)	71.7	97.1	90.7	92.5	90.8	91.3	95.8	89.8	92.5	89.5	89.4

• Results on SQuAD

Model	Train FLOPs	Params	SQuAI	D 1.1 dev	SQuAD 2.0 dev		SQuA	D 2.0 test
			EM	FI	EM	FI	EM	FI
BERT-Base	6.4e19 (0.09x)	110M	80.8	88.5	-	-	-	-
BERT	1.9e20 (0.27x)	335M	84.1	90.9	79.0	81.8	80.0	83.0
SpanBERT	7.1e20 (1x)	335M	88.8	94.6	85.7	88.7	85.7	88.7
XLNet-Base	6.6e19 (0.09x)	117M	81.3	-	78.5	-	-	-
XLNet	3.9e21 (5.4x)	360M	89.7	95.1	87.9	90.6	87.9	90.7
RoBERTa-100K	6.4e20 (0.90x)	356M	-	94.0	-	87.7	-	-
RoBERTa-500K	3.2e21 (4.5x)	356M	88.9	94.6	86.5	89.4	86.8	89.8
ALBERT	3.1e22 (44x)	235M	89.3	94.8	87.4	90.2	88.1	90.9
BERT (ours)	7.1e20 (1x)	335M	88.0	93.7	84.7	87.5	-	-
ELECTRA-Base	6.4e19 (0.09x)	110M	84.5	90.8	80.5	83.3	-	-
ELECTRA-400K	7.1e20 (1x)	335M	88.7	94.2	86.9	89.6	-	-
ELECTRA-1.75M	3.1e21 (4.4x)	335M	89.7	94.9	88.0	90.6	88.7	91.4

Overview

1 Background

2 Pre-BERT Era

3 BERT





Conclusion

- Background
 - language model
 - sequence to sequence model
- Pre-BERT Era
 - semi-supervised sequence learning
 - ontext2vec
 - pre-trained seq2seq
 - ELMo
 - OpenAI GPT
- BERT
- Post-BERT Era
 - RoBERTa
 - XLNet
 - ERNIE (Tsinghua & Baidu)
 - ALBERT
 - ELECTRA

Thanks.

HP: http://keg.cs.tsinghua.edu.cn/jietang/ Email: jietang@tsinghua.edu.cn