# **Microsoft** MPNet: Masked and Permuted Pre-training for Natural Language Understanding

<sup>1</sup>Kaitao Song<sup>\*</sup>, <sup>2</sup>Xu Tan<sup>\*</sup>, <sup>2</sup>Tao Qin, <sup>1</sup>Jianfeng Lu, <sup>2</sup>Tie-Yan Liu <sup>1</sup>Nanjing University of Science and Technology, <sup>2</sup>Microsoft Research

## **Motivation**

- **BERT** (Masked Language Model, **MLM**) allows model to see the full input sentence (*input consistency*), but ignores the dependency between the masked/predicted tokens (output dependency).
- XLNet (Permuted Language Model, PLM) leverages dependency between masked/predicted tokens (*output dependency*) but can only see its preceding tokens rather than the full sentence (*input consistency*).



#### Method

# **MPNet – Inherit advantages of MLM and PLM**

- Autoregressive prediction (avoid the limitation in BERT)
  - Each predicted token condition on previous predicted tokens to ensure *output* dependency
- Position compensation (avoid the limitation in XLNet)
  - Each predicted token can see full position information to ensure *input consistency*



## **Analysis**

<ul> <li>Case Stu</li> </ul>	udy (Factorization)	
Objective	Modeling	
MLM (BERT)	log P(sentence  the task is [M] [M]) + log P(classification  the ta	ask is [M] [M])
PLM (XLNet)	log P(sentence  the task is) + log P(classification  the task	ask is sentence
MPNet	log P(sentence  the task is [M] [M]) + log P(classification  the ta	ask is <mark>sentence</mark>

Comparisons between MPNet and M



# **Experiments**

Setting: Pre-train on English Wikipedia, BooksCorpus, OpenWebText, Stories (160GB) with a batch size of 8192 for 500K Steps. The mask ratio is set as 15%. • GLUE

		MNL	I QNLI	QQP	RTE	SST	MRPC	CoLA	STS	Avg
	Single model on dev set	t								
	BERT (Devlin et al., 2019)	84.5	91.7	91.3	68.6	93.2	87.3	58.9	89.5	83.1
	XLNet (Yang et al., 2019)	86.8	91.7	91.4	74.0	94.7	88.2	60.2	89.5	84.5
	RoBERTa (Liu et al., 2019	9a) 87.6	92.8	91.9	78.7	94.8	90.2	63.6	91.2	86.4
	MPNet	88.5	93.3	91.9	85.8	95.4	91.8	65.0	91.1	87.9
	Single model on test set	t								
	BERT (Devlin et al., 2019)	84.6	90.5	89.2	66.4	93.5	84.8	52.1	87.1	79.9
	ELECTRA (Clark et al., 2	<b>88.5 88.5</b>	93.1	89.5	75.2	96.0	88.1	64.6	91.0	85.8
	MPNet	88.5	93.1	89.9	81.0	96.0	89.1	64.0	90.7	86.5
• SQu	AD									
	Method	SQuAD	v1.1 (dev)	SQ	uADv2	.0 (de	v)   <b>SQ</b> t	ADv2.	0 (test	)
	BERT 2	80.8	/88.5		73.7/7	6.3		73.1/76	.2	
	XLNet [5]	81	.3/-		80.2	/-		-/-		
	RoBERTa [7]	84.6	/91.5		80.5/8	3.7		-/-		
	MPNet	86.9	/92.7		82.7/8	5.7		82.8/85	.8	
• Abla	tion Study									_
	Model Setting			5	QuADv	1.1	SQuADv2	2.0   M	NLI	SST-2
	MPNet				85.0/91	.4	80.5/83.	3   8	6.2	94.0
	- position compensation	on (= PLM)	)		83.0/89	.9	78.5/81.	0   8	5.6	93.4
	- permutation (= MLN	A + output of	dependency	y)   :	84.1/90	.6	79.2/81.	8 8	5.7	93.5
	– permutation & outpu	ut dependen	cy (= MLN	A)	82.0/89	.5	76.8/79.	8 8	5.6	93.3
		. //	••••		/ •	<b>.</b> .	<b>C</b> .			
	Code: ht	tps://	github	<b>0.CO</b>	<u>m/Iv</u>	<u>'licr</u>	<u>osoft</u>	/IVIP	Net	







[M])

MI	LM/PLM (	To	okens and	Position	s)

# Tokens	# Positions
85%	100%
92.5%	92.5%
92.5%	100%