

^t This work was mainly performed in NAVER Clova

ST-BERT: Cross-Modal Language Model Pre-Training for End-to-end Spoken Language Understanding

. .

Minjeong Kim^{3*}, Gyuwan Kim^{1,2}, Sang-Woo Lee^{1,2}, Jung-Woo Ha^{1,2} ¹NAVER Clova, ²NAVER AI Lab, ³Upstage AI Research

Embe	d
Moda Embe	ali ed

Experiments

ining				-					
	Model	Full	10%	1%	Model	SmartLights		Snips	
former	Lugosch et al. [2] Wong et al. [2] (PEPT)	98.80%	97.96%	82.78%	Huang and Chen [28]	Far 17 38%	67.98%	-	
†	Wang et al. [5] (BERT) Wang et al. [3] (ERNIE)	98.93% 99.02%	-	-		47.3070	01.90%	09.3370	
s ì	Cho et al. $[5]$ (End (IE)	98.98%	98.12%	83.12%	ST-BERT (Ours)	60.98%	84.65%	96.21%	
	Price [31]	99.3%	-	-	- CM-CLM - text data (pre-training)	54 93%	81.97%	95.07% 95.64%	
	ST-BERT (Ours)	99.50%	99.13%	95.64%	+ DAPT	<u>69.40%</u>	86.91%	96.07%	
esentations	- CM-CLM	99.42%	99.09%	96.51%	Table 2 Smartl ia	hte and Sn	ine Datas	ente	
nbedding	- text data (pre-training)	99.39%	99.04%	89.81%	Table 2. SmartLly	nis and Sh	ips Dalas	0013	
	+ DAPT	99.59%	99.25%	95.83%					
	Table 1. Fluent Speec	Table 1. Fluent Speech Command (FSC) Dataset				(2) Ablation studies			
	1 Implementation				- [- CM-CLM] : Pre-train with CM-MLM only				
Posterior		I. Implementation			- [- text data] : Pre-train with MLM on				
Madal	(1) Pre-training	(1) Pre-training							
IVIODEI	 Pre-training on Librispeech 			 Both pre-training methods are effective The usage of the synthesized voice 					
ech	 Curriculum pre-training 								
	- Initialized from pre-trained BERT-base								
		Due tuelet		T \	affects to the res	sult of Snip	os datas	et	
eform	(2) Domain-Adaptive	(2) Domain-Adaptive Fre-training (DAFT)			(3) Domain-Adaptive Pre-training (DAPT)				
ng	- Continues pre-training with domain-			- Leads to further improvement					
	specific speech	-text pair d	ata, whe	n those	Domorkobly offe		moro po		
	are available				- nemarkably ellective on a more noisy SLU				
	(2) Eine tuning				data (tar field in a	SmartLisg	nt datas	et was	
					recorded from 2	meters av	vay)		
.d	- Fine-tuning on F	-SC, Smar	tlights a	nd	(4) Data shortage sce	enario			
	Snips datasets				- ST_REPT shows	comparat	ivoly ma	rainal	
oneme	- Use synthesized	d audio for	Snips da	ataset				iyilal	
osterior	- Measure perforr	mance twic	ce		performance de	gradation	_		
ddings for the					- Uni-modal pre-t	raining on	speech	suffers	
\sim	2. Results				from a relatively	large perfe	ormance	drop	
	(1) Main regulte				- Leveraging textu	ual informa	ation dur	ing pre-	
		- For all datasets, our model shows			training is critical when limited SLU data are available				
	- FOR all Galasets,								
	: : remarkable resu								

111)							Tra	Transformer		
Subword	†	≜	≜	≜	↑	≜	≜	≜	4	
ings	[CLS]	В	EH	L	L	R	IH	NG	IH	
	+	+	+	+	+	+	+	+	+	
ings	0	1	2	3	4	5	6	7	8	
,	+	+	+	+	+	+	+	+	+	
ings	0	0	0	0	0	0	0	0	0	

-Modal Conditional Language Wodeling (Civi-- Masks out the entire sequence of target modality

- Model needs to predict the masked representation solely conditioned on the source modality - Two types of CM-CLM exist: speech-to-text and text-to-speech CM-CLM - More challenging task than CM-MLM



