# SHAPE: Shifted Absolute Position Embedding for Transformers

Shun Kiyono<sup>1,2</sup>, Sosuke Kobayashi<sup>2,3</sup>, Jun Suzuki<sup>2,1</sup>, Kentaro Inui<sup>2,1</sup> <sup>1</sup>RIKEN <sup>2</sup>Tohoku University <sup>3</sup>Preferred Networks, Inc.



Kevi

# Overview: building absolute position embedding that can extrapolate and shift-invariant



# **Experiment on Machine Translation (WMT EnDe)**

#### (1) Vanilla


- Standard setting for MT Sanity check of baseline
- performance

#### ③ Interpolate



- Concatenate neighboring sequences
- Evaluate performance on
  Tokens are more infrequent at given position

## **Experimental Result and Observations**

unseen length

2 Extrapolate

Remove sequence

longer than 50 tokens

Dataset	Model	Valid	Test	Speed		
VANILLA	$APE^{\dagger}$	23.61	30.46	x1.00	•	Δ
	$RPE^{\dagger}$	23.67	30.54	x0.91		
	$SHAPE^{\dagger}$	23.63	30.49	x1.01	•	S
EXTRAPOLATE	APE	22.18	29.22	x1.00	•	B
	RPE	22.97	29.86	x0.91		_
	SHAPE	22.96	29.80	x0.99	•	2
INTERPOLATE	APE	31.40	38.23	x1.00		Ы
	RPE*	-	-	-	•	Γ.
	SHAPE	32.50	39.09	x0.99	•	S

SHAPE has no risk of per	rtormance drop
	tpertorm APF

HAPE is as fast as APE while RPE is not

PE is prohibitively slow

HAPE outperforms APE



## **Analysis 1: SHAPE is Shift-invariant**



## **Analysis 2: SHAPE Extrapolates better than APE**



• SHAPE outperforms APE in gray no-training-data region → **better** extrapolation than APE



- SHAPE outperforms APE
- Data augmentation effect?