

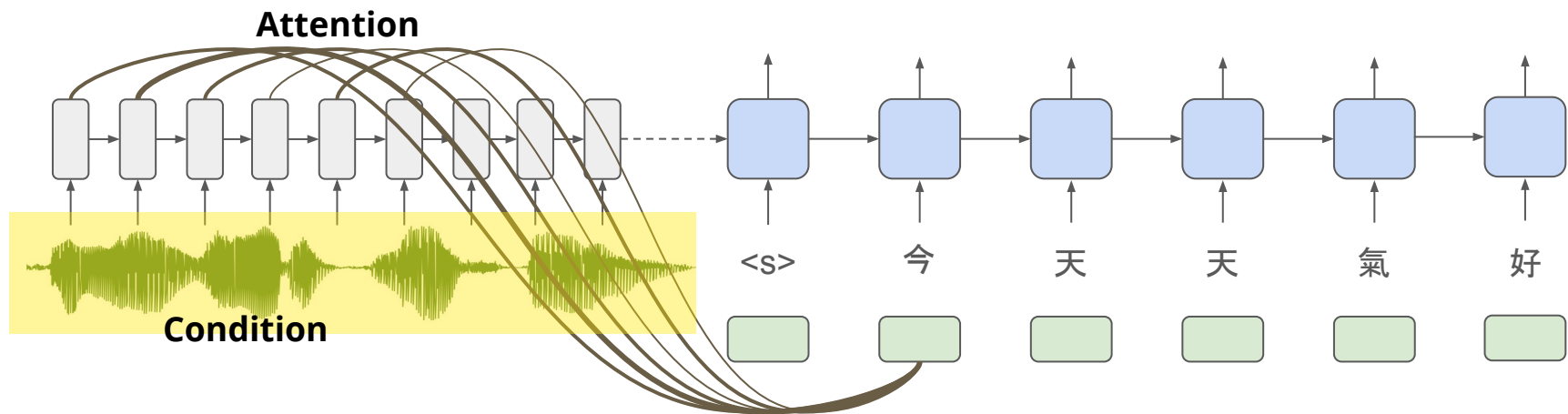
# Non-Autoregressive Sequence Generation

Yung-Sung Chuang 莊永松

<https://voidism.github.io/home/>

2020/05/13

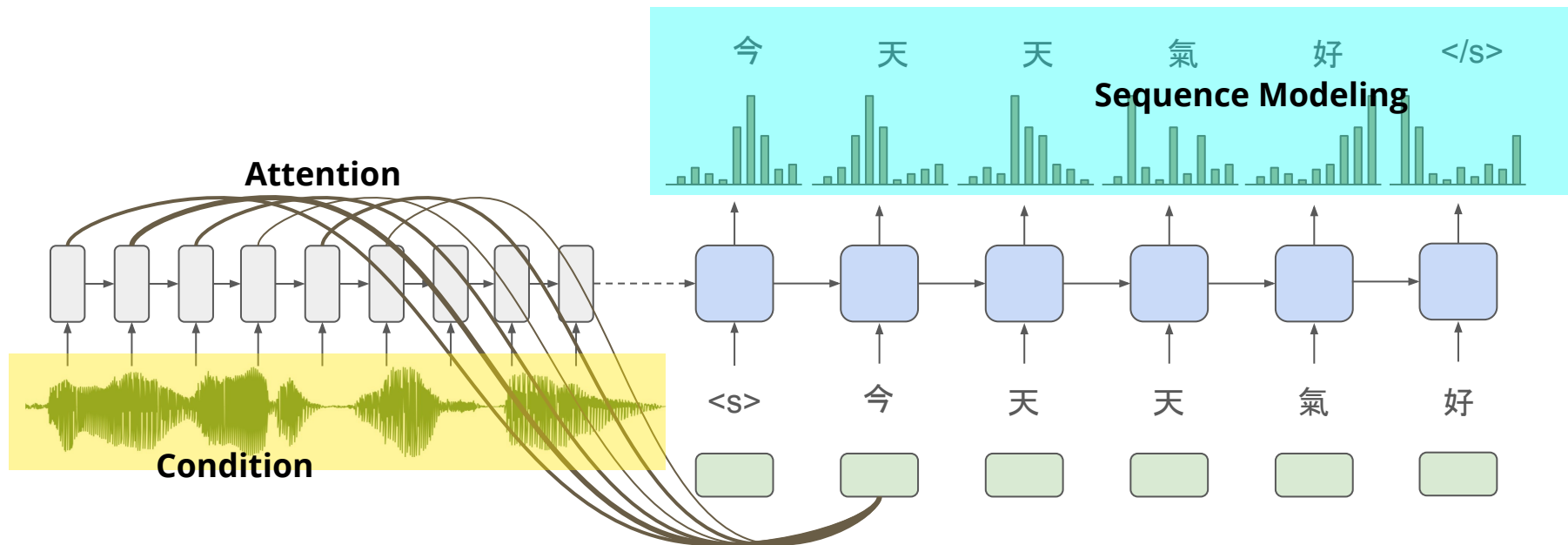
# What We have done in HW1...



=> Conditional Sequence Generation

# What We have done in HW1...

$$p(Y|X) = \prod_{t=1}^T p(y_t | y_{<t}, X)$$



=> Conditional Sequence Generation

# More Conditional Sequence Generation



A woman is throwing a frisbee in a park.

**Image Caption Generation**

Hello, World!



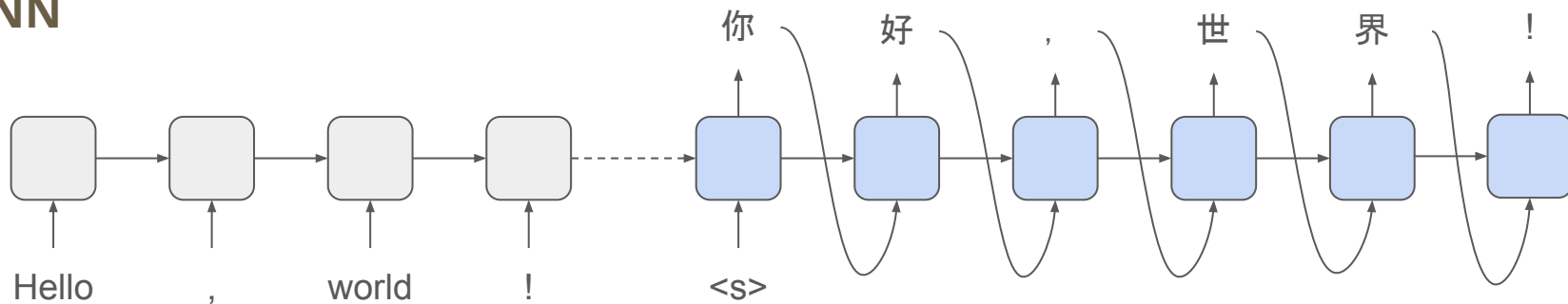
你好，世界！

**Machine Translation**

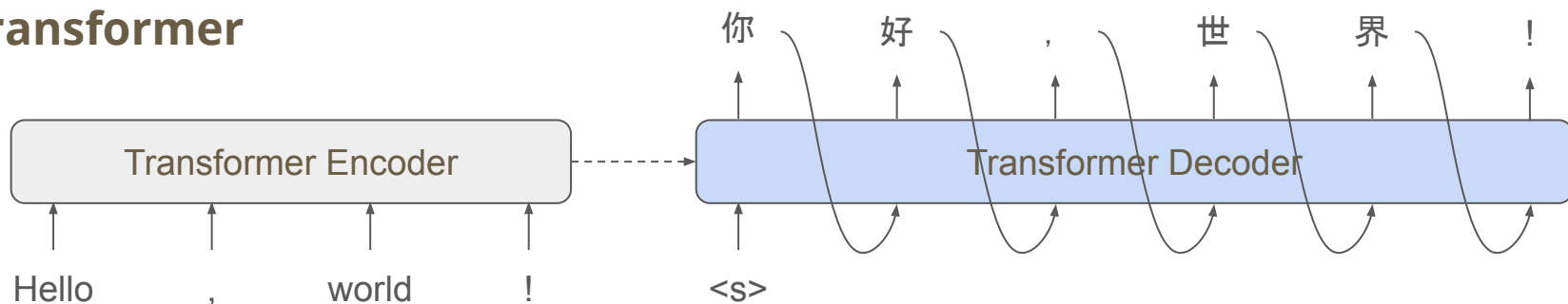
# Autoregressive model (inference stage)

$$p(Y|X) = \prod_{t=1}^T p(y_t | y_{<t}, X)$$

## RNN



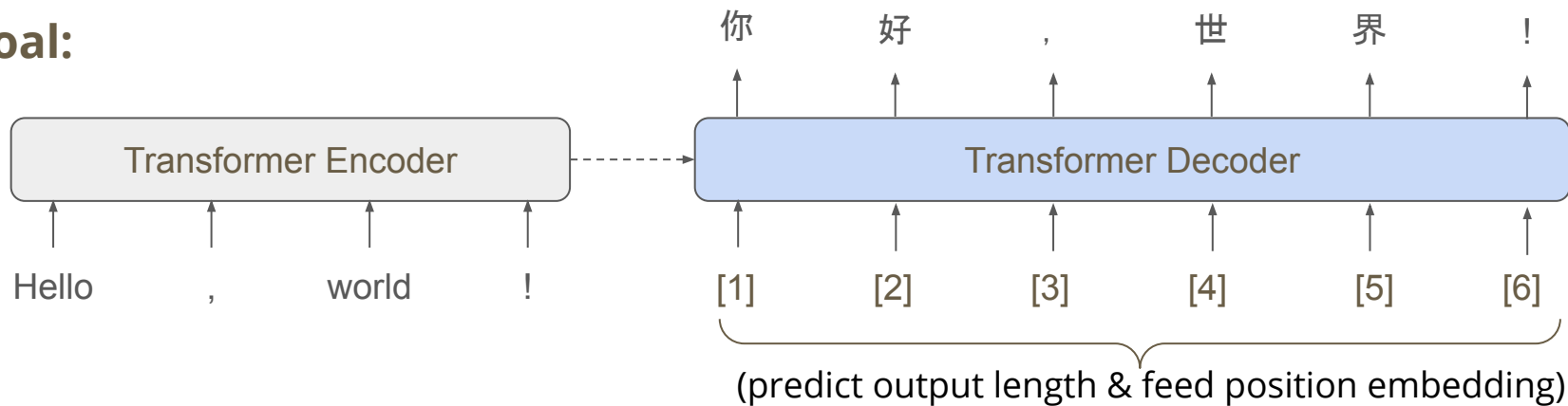
## Transformer



# Non-autoregressive model (mostly by Transformer)

$$p(Y|X) = \prod_{t=1}^T p(y_t | X)$$

**Goal:**



# What's the problem?

Output would be average of images



70%



30%



model learned

## Text-to-Image

a dog is running



a bird is flying



- Traditional supervised approach

c<sup>1</sup>: a dog is running



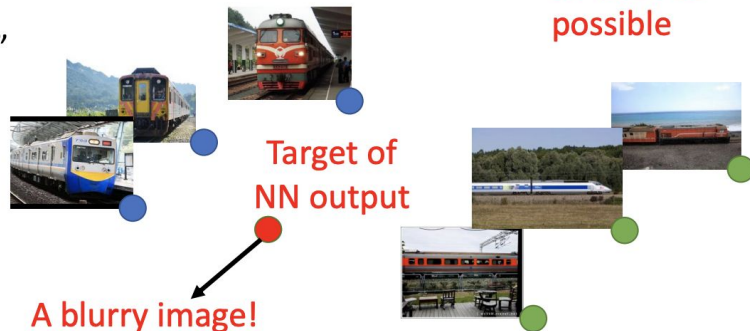
NN

Image

as close as possible

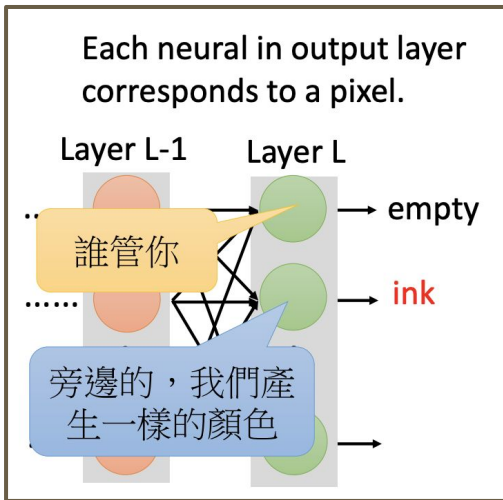


Text: "train"



# What's the problem?

No dependency on output structure



No Latent Variable



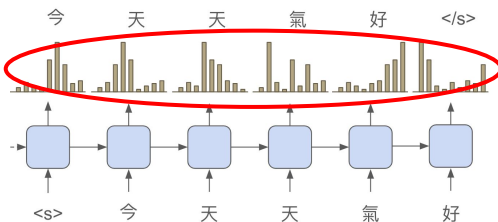
70%



30%

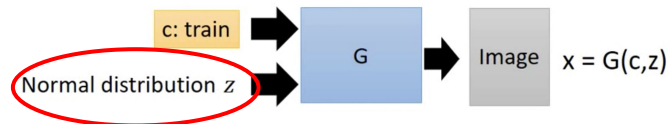
Solutions:

Autoregressive model  
ex: PixelRNN



GAN

Conditional GAN



# What's the problem?

哈囉

50%

你好

50%

Transformer Encoder

Hello

!

Transformer Decoder

[1]

[2]

[3]

嗨

!

哈

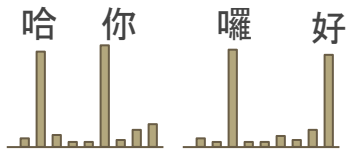
囉

!

你

好

!



你囉




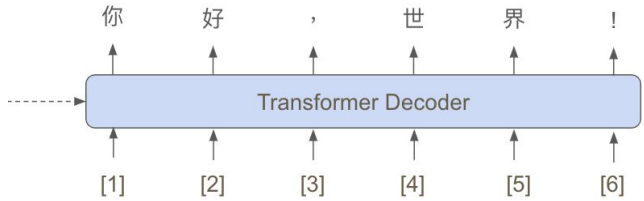

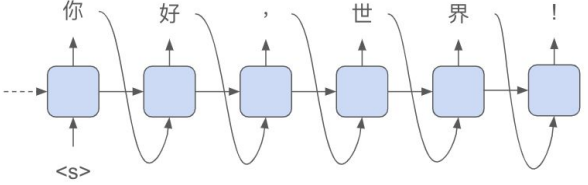


哈好



Hello

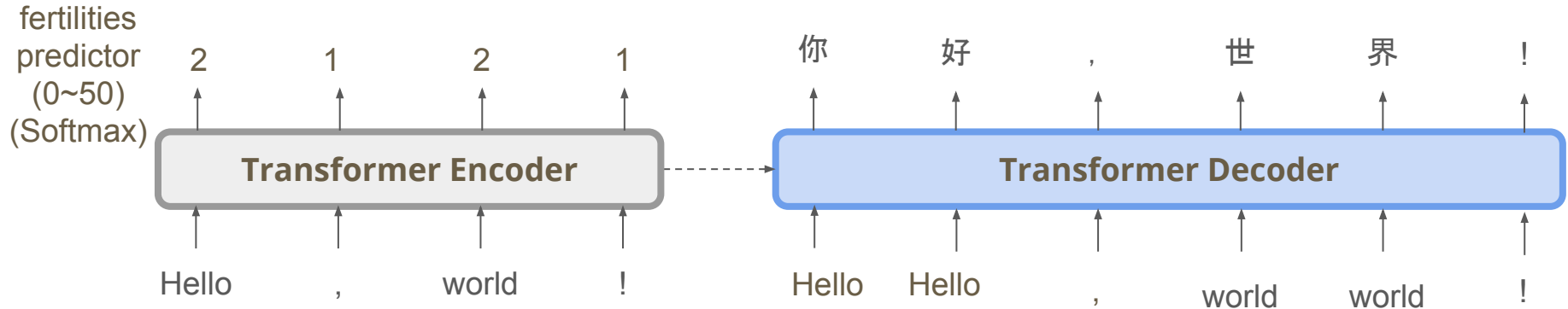
Transformer Decoder

=> multi-modality problem!

generate target / methods	Image	Text
<b>naive approach</b>	Deconvolution layer + L2 loss  Average of images (bad)	non-autoregressive decoder 
<b>autoregressive</b>	ex: PixelRNN, VQVAE-2 	autoregressive decoder 
<b>Generative Adversarial Networks w/ non-auto model</b>		

# Vanilla NAT (Non-Autoregressive Translation)

- Predict **fertility** as latent variable & Copy input words
- Represents sentence-level “plan” before writing Y



# Fertility

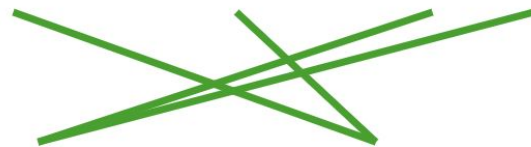
## Fixed targets

1. Labels comes from external aligner
2. Observing attention weights in auto-regressive models

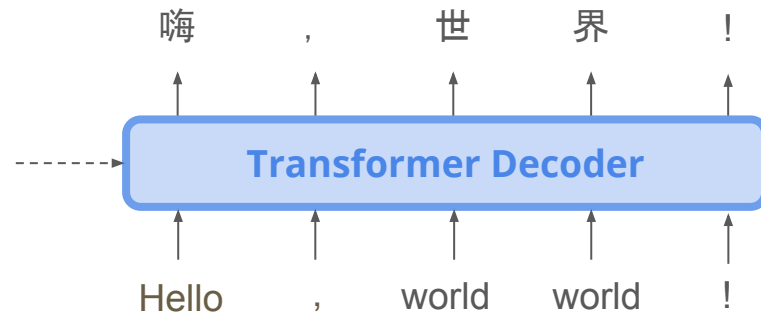
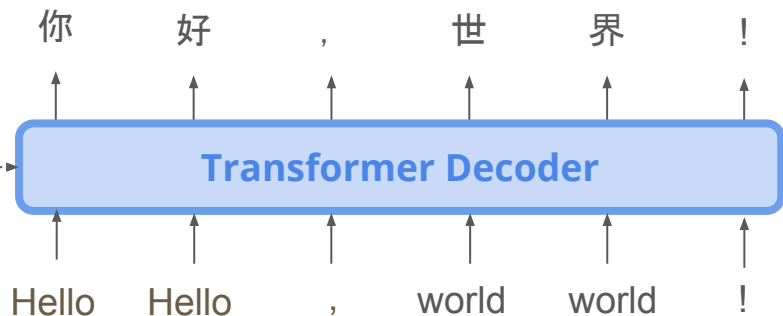
## Fine-tune (after NAT model converges)

1. Updating fertility classifier with REINFORCE

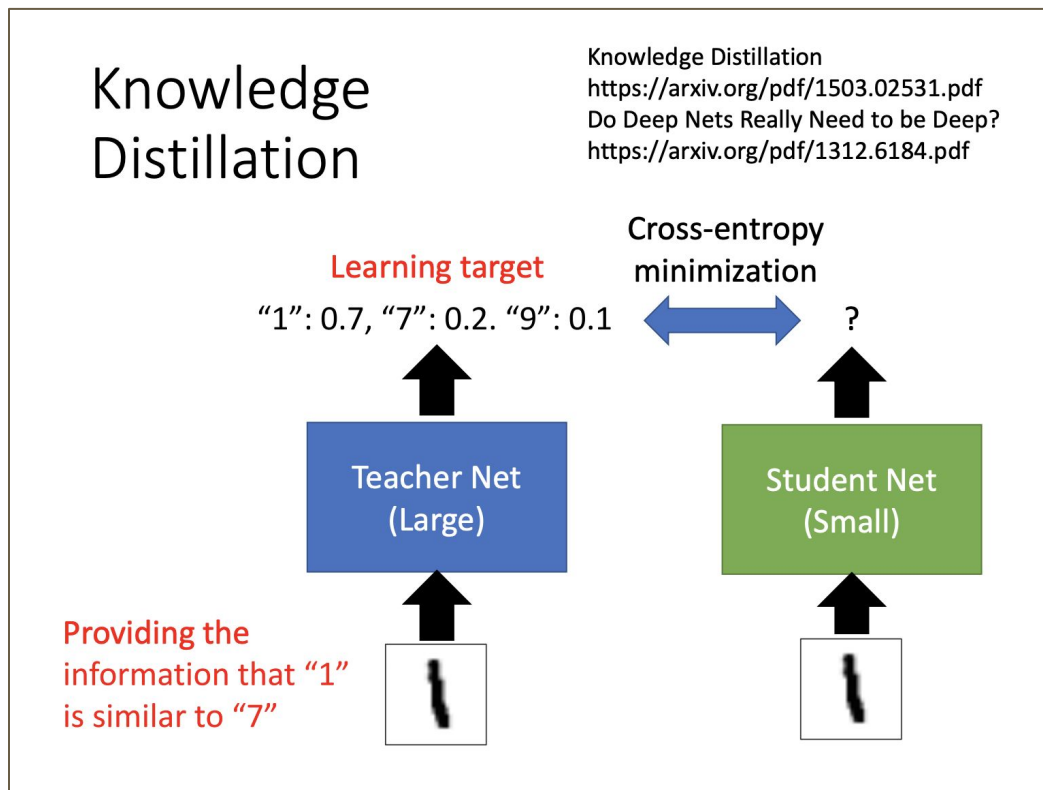
Word alignments are fun



Lustige Wortalignierungen

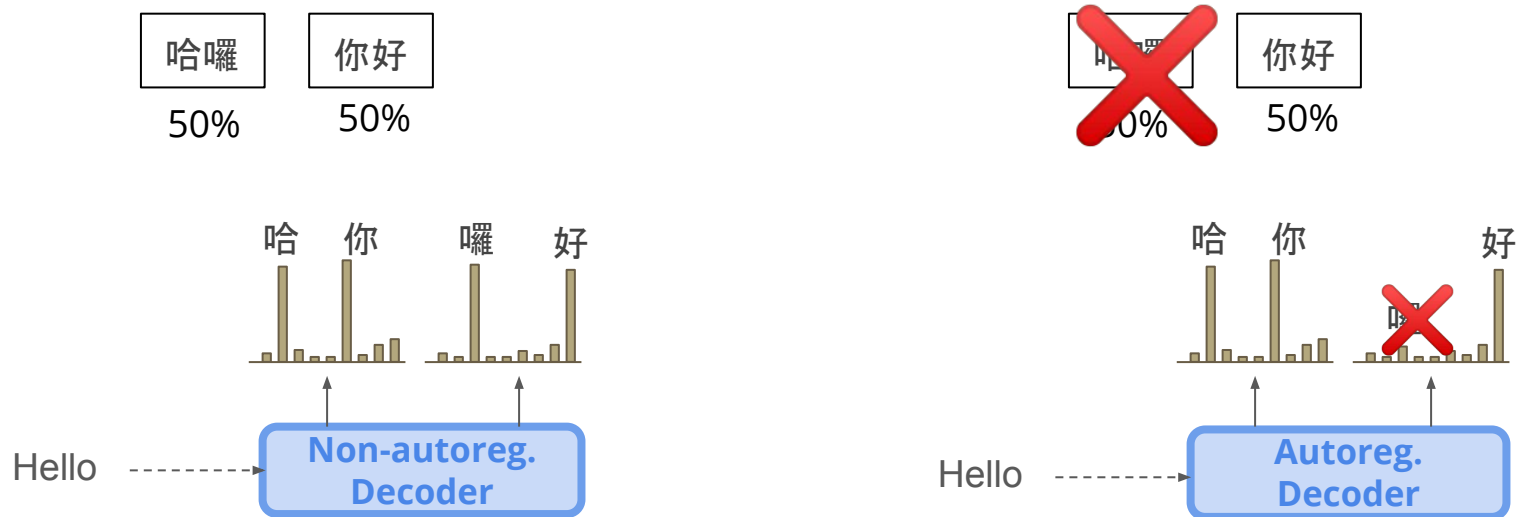


# Sequence-level knowledge distillation



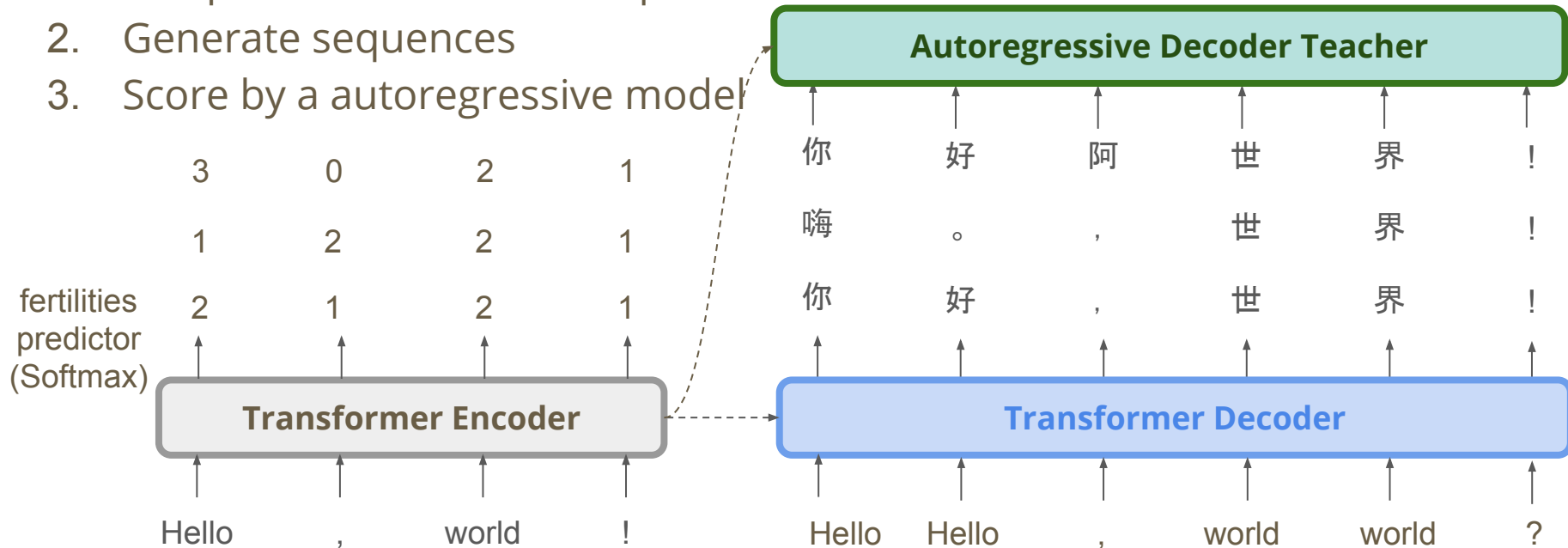
# Sequence-level knowledge distillation

- Teacher: Autoregressive model, Student: Non-autoregressive model
- Construct new corpus by autoregressive teacher model
- Teacher's greedy decode output as student's training target



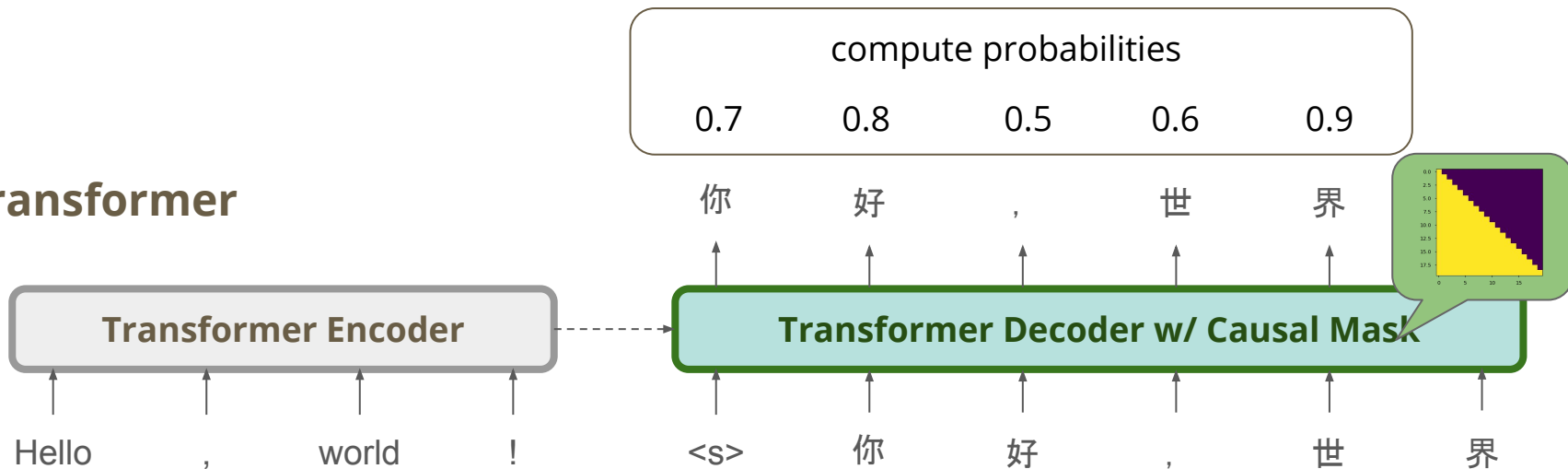
# Noisy Parallel Decoding (NPD)

1. Sample several fertile sequences
2. Generate sequences
3. Score by an autoregressive model



# Autoregressive model w/ teacher forcing

Transformer



# Experiments

Models	WMT14		WMT16		IWSLT16		
	En→De	De→En	En→Ro	Ro→En	En→De	Latency / Speedup	
NAT	17.35	20.62	26.22	27.83	25.20	39 ms	15.6×
NAT (+FT)	17.69	21.47	27.29	29.06	26.52	39 ms	15.6×
NAT (+FT + NPD $s = 10$ )	18.66	22.41	29.02	30.76	27.44	79 ms	7.68×
NAT (+FT + NPD $s = 100$ )	19.17	23.20	29.79	<b>31.44</b>	28.16	257 ms	2.36×
Autoregressive ( $b = 1$ )	22.71	26.39	31.35	31.03	28.89	408 ms	1.49×
Autoregressive ( $b = 4$ )	23.45	27.02	31.91	31.76	29.70	607 ms	1.00×

# Experiments

Models	WMT14		WMT16		IWSLT16		
	En→De	De→En	En→Ro	Ro→En	En→De	Latency / Speedup	
NAT	17.35	20.62	26.22	27.83	25.20	39 ms	15.6×
NAT (+FT)	17.69	21.47	27.29	29.06	26.52	39 ms	15.6×
NAT (+FT + NPD $s = 10$ )	18.66	22.41	29.02	30.76	27.44	79 ms	7.68×
NAT (+FT + NPD $s = 100$ )	19.17	23.20	29.79	<b>31.44</b>	28.16	257 ms	2.36×
Autoregressive ( $b = 1$ )	22.71	26.39	31.35	31.03	28.89	408 ms	1.49×
Autoregressive ( $b = 4$ )	23.45	27.02	31.91	31.76	29.70	607 ms	1.00×

# Experiments

Models	WMT14		WMT16		IWSLT16		
	En→De	De→En	En→Ro	Ro→En	En→De	Latency / Speedup	
NAT	17.35	20.62	26.22	27.83	25.20	39 ms	15.6×
NAT (+FT)	17.69	21.47	27.29	29.06	26.52	39 ms	15.6×
NAT (+FT + NPD $s = 10$ )	18.66	22.41	29.02	30.76	27.44	79 ms	7.68×
NAT (+FT + NPD $s = 100$ )	19.17	23.20	29.79	<b>31.44</b>	28.16	257 ms	2.36×
Autoregressive ( $b = 1$ )	22.71	26.39	31.35	31.03	28.89	408 ms	1.49×
Autoregressive ( $b = 4$ )	23.45	27.02	31.91	31.76	29.70	607 ms	1.00×

# Experiments

Models	WMT14		WMT16		IWSLT16		
	En→De	De→En	En→Ro	Ro→En	En→De	Latency / Speedup	
NAT	17.35	20.62	26.22	27.83	25.20	39 ms	15.6×
NAT (+FT)	17.69	21.47	27.29	29.06	26.52	39 ms	15.6×
NAT (+FT + NPD $s = 10$ )	18.66	22.41	29.02	30.76	27.44	79 ms	7.68×
NAT (+FT + NPD $s = 100$ )	19.17	23.20	29.79	<b>31.44</b>	28.16	257 ms	2.36×
Autoregressive ( $b = 1$ )	22.71	26.39	31.35	31.03	28.89	408 ms	1.49×
Autoregressive ( $b = 4$ )	23.45	27.02	31.91	31.76	29.70	607 ms	1.00×

# Experiments

Models	WMT14		WMT16		IWSLT16		
	En→De	De→En	En→Ro	Ro→En	En→De	Latency / Speedup	
NAT	17.35	20.62	26.22	27.83	25.20	39 ms	15.6×
NAT (+FT)	17.69	21.47	27.29	29.06	26.52	39 ms	15.6×
NAT (+FT + NPD $s = 10$ )	18.66	22.41	29.02	30.76	27.44	79 ms	7.68×
NAT (+FT + NPD $s = 100$ )	19.17	23.20	29.79	<b>31.44</b>	28.16	257 ms	2.36×
Autoregressive ( $b = 1$ )	22.71	26.39	31.35	31.03	28.89	408 ms	1.49×
Autoregressive ( $b = 4$ )	23.45	27.02	31.91	31.76	29.70	607 ms	1.00×

uniform: 看對應長度, 四捨五入直接 Copy inputs

# Ablation

Distillation		Decoder Inputs			Fine-tuning			BLEU	BLEU (T)
$b=1$	$b=4$	+uniform	+fertility	+PosAtt	$+\mathcal{L}_{KD}$	$+\mathcal{L}_{BP}$	$+\mathcal{L}_{RL}$		
				✓				$\approx 2$	
		✓		✓				16.51	
			✓	✓				18.87	
✓		✓		✓				20.72	
	✓	✓		✓				21.12	
✓			✓					24.02	43.91
✓			✓	✓				25.20	45.41
✓		✓		✓	✓	✓		22.44	
✓			✓	✓			✓	×	×
✓			✓	✓		✓		×	×
✓			✓	✓	✓	✓		25.76	46.11
✓			✓	✓	✓	✓	✓	<b>26.52</b>	<b>47.38</b>

uniform: 看對應長度, 四捨五入直接 Copy inputs

# Ablation

Distillation		Decoder Inputs			Fine-tuning			BLEU	BLEU (T)
$b=1$	$b=4$	+uniform	+fertility	+PosAtt	$+\mathcal{L}_{KD}$	$+\mathcal{L}_{BP}$	$+\mathcal{L}_{RL}$		
				✓				$\approx 2$	
		✓		✓				16.51	
			✓	✓				18.87	
✓		✓		✓				20.72	
	✓	✓		✓				21.12	
✓			✓					24.02	43.91
✓			✓	✓				25.20	45.41
✓		✓		✓	✓	✓		22.44	
✓			✓	✓			✓	×	×
✓			✓	✓		✓		×	×
✓			✓	✓	✓	✓		25.76	46.11
✓			✓	✓	✓	✓	✓	<b>26.52</b>	<b>47.38</b>

uniform: 看對應長度, 四捨五入直接 Copy inputs

# Ablation

Distillation		Decoder Inputs			Fine-tuning			BLEU	BLEU (T)
$b=1$	$b=4$	+uniform	+fertility	+PosAtt	$+\mathcal{L}_{KD}$	$+\mathcal{L}_{BP}$	$+\mathcal{L}_{RL}$		
				✓				$\approx 2$	
		✓		✓				16.51	
			✓	✓				18.87	
✓		✓		✓				20.72	
	✓	✓		✓				21.12	
✓			✓					24.02	43.91
✓			✓	✓				25.20	45.41
✓		✓		✓	✓	✓		22.44	
✓			✓	✓			✓	×	×
✓			✓	✓		✓		×	×
✓			✓	✓	✓	✓		25.76	46.11
✓			✓	✓	✓	✓	✓	<b>26.52</b>	<b>47.38</b>

uniform: 看對應長度, 四捨五入直接 Copy inputs

# Ablation

Distillation		Decoder Inputs			Fine-tuning			BLEU	BLEU (T)
$b=1$	$b=4$	+uniform	+fertility	+PosAtt	$+\mathcal{L}_{KD}$	$+\mathcal{L}_{BP}$	$+\mathcal{L}_{RL}$		
				✓				$\approx 2$	
		✓		✓				16.51	
			✓	✓				18.87	
✓		✓		✓				20.72	
	✓	✓		✓				21.12	
✓			✓					24.02	43.91
✓			✓	✓				25.20	45.41
✓		✓		✓	✓	✓		22.44	
✓			✓	✓			✓	×	×
✓			✓	✓		✓		×	×
✓			✓	✓	✓	✓		25.76	46.11
✓			✓	✓	✓	✓	✓	26.52	47.38

# Evolution of NAT

## 1. Vanilla NAT

- Gu et al., ICLR'18

## 2. Iterative Refinement

- iNAT, Lee et al., EMNLP'18
- Mask-Predict, Ghazvininejad et al., EMNLP'19

## 3. Insertion-based

- Insertion Transformer, Stern et al., ICML'19
- KERMIT, Chan et al., arXiv'19

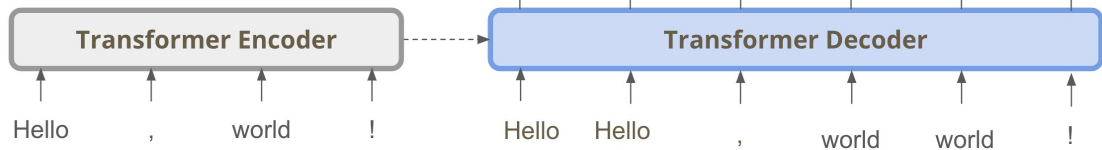
## 4. Insert+Delete

- Levenshtein Transformer, Gu et al., NeurIPS'19

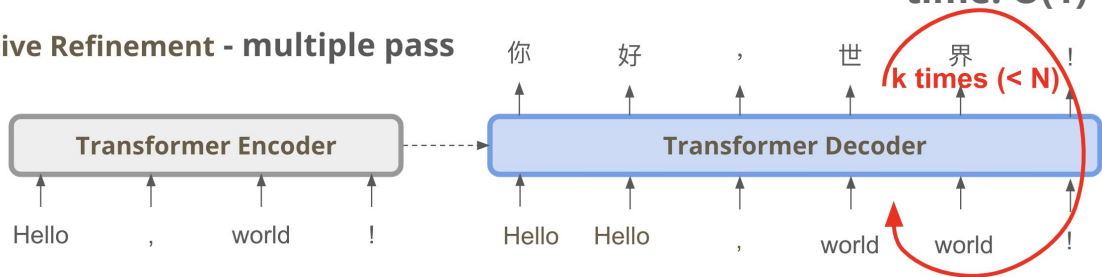
## 5. CTC-based

- E2E NAT w/ CTC, Jindrich Libovick et al., EMNLP'20
- NAT w/ Latent Alignments, Saharia et al., arXiv'20

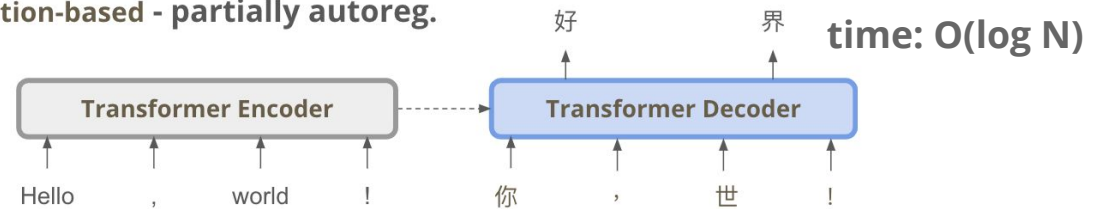
Vanilla NAT - one pass



Iterative Refinement - multiple pass



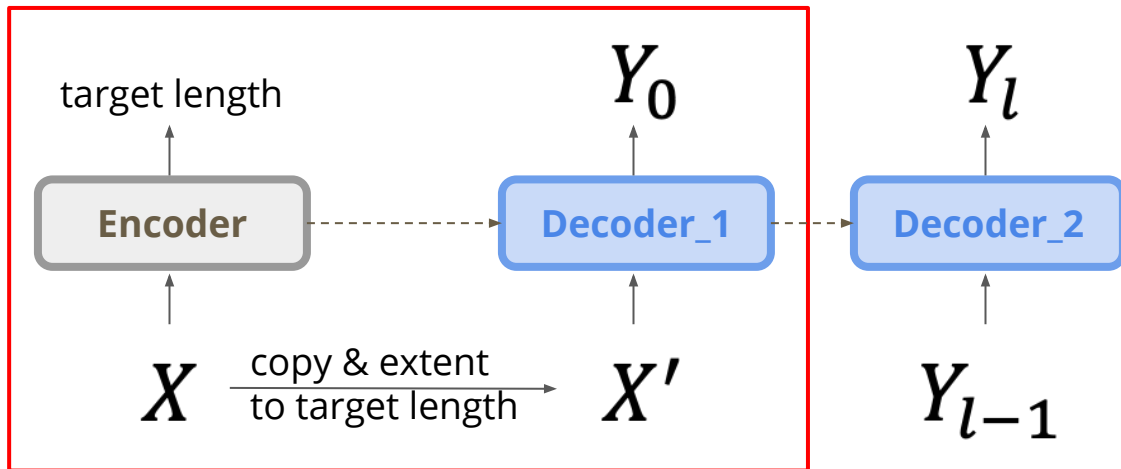
Insertion-based - partially autoreg.



$[\ ] \rightarrow [D] \rightarrow [B, D, F] \rightarrow [A, B, C, D, E, F, G]$

# NAT with Iterative Refinement

$$X \rightarrow Y_0 \rightarrow Y_1 \rightarrow Y_2 \rightarrow \dots \rightarrow Y_L = Y$$



## Corruption Process

- (1) replace  $y_{t+1}$  with  $y_t$
- (2) replace  $y_t$  with a random token
- (3) swap  $y_t$  and  $y_{t+1}$ .

## Stochastically mix cost function

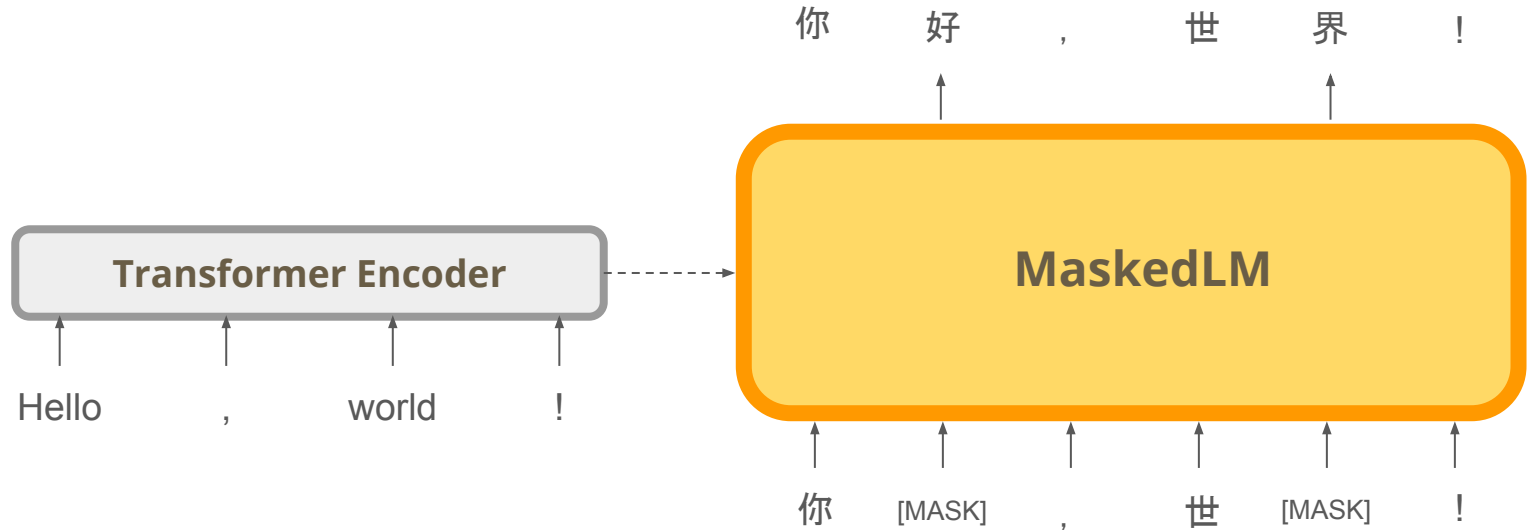
$$J(\theta) = - \sum_{l=0}^{L+1} \left( \alpha_l \sum_{t=1}^T \log p_{\theta}(y_t^* | \hat{Y}^{l-1}, X) \right. \\ \left. + (1 - \alpha_l) \sum_{t=1}^T \log p_{\theta}(y_t^* | \tilde{Y}, X) \right), \quad (4)$$

# Results

		IWSLT'16 En-De				WMT'16 En-Ro				WMT'14 En-De				MS COCO		
		En→	De→	GPU	CPU	En→	Ro→	GPU	CPU	En→	De→	GPU	CPU	BLEU	GPU	CPU
AR	$b = 1$	28.64	34.11	70.3	32.2	31.93	31.55	55.6	15.7	23.77	28.15	54.0	15.8	23.47	4.3	2.1
	$b = 4$	28.98	34.81	63.8	14.6	32.40	32.06	43.3	7.3	24.57	28.47	44.9	7.0	24.78	3.6	1.0
NAT	FT	26.52	–	–	–	27.29	29.06	–	–	17.69	21.47	–	–	–	–	–
	FT+NPD	28.16	–	–	–	29.79	31.44	–	–	19.17	23.30	–	–	–	–	–
Our Model	$i_{\text{dec}} = 1$	22.20	27.68	573.0	213.2	24.45	25.73	694.2	98.6	13.91	16.77	511.4	83.3	20.12	17.1	8.9
	$i_{\text{dec}} = 2$	24.82	30.23	423.8	110.9	27.10	28.15	332.7	62.8	16.95	20.39	393.6	49.6	20.88	12.0	5.7
	$i_{\text{dec}} = 5$	26.58	31.85	189.7	52.8	28.86	29.72	194.4	29.0	20.26	23.86	139.7	23.1	21.12	6.2	2.8
	$i_{\text{dec}} = 10$	27.11	32.31	98.8	24.1	29.32	30.19	93.1	14.8	21.61	25.48	90.4	12.3	21.24	2.0	1.2
	Adaptive	27.01	32.43	125.9	29.3	29.66	30.30	118.3	16.5	21.54	25.43	107.2	20.3	21.12	10.8	4.8

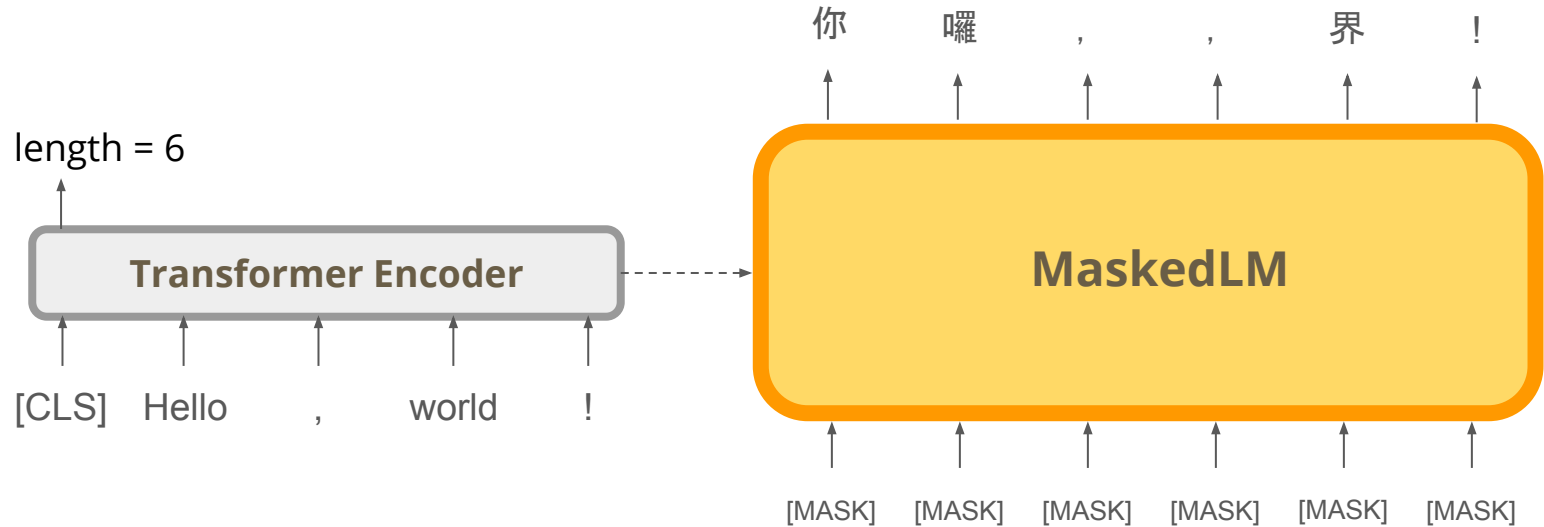
# Mask-Predict

## Conditional Masked Language Models (CMLM):



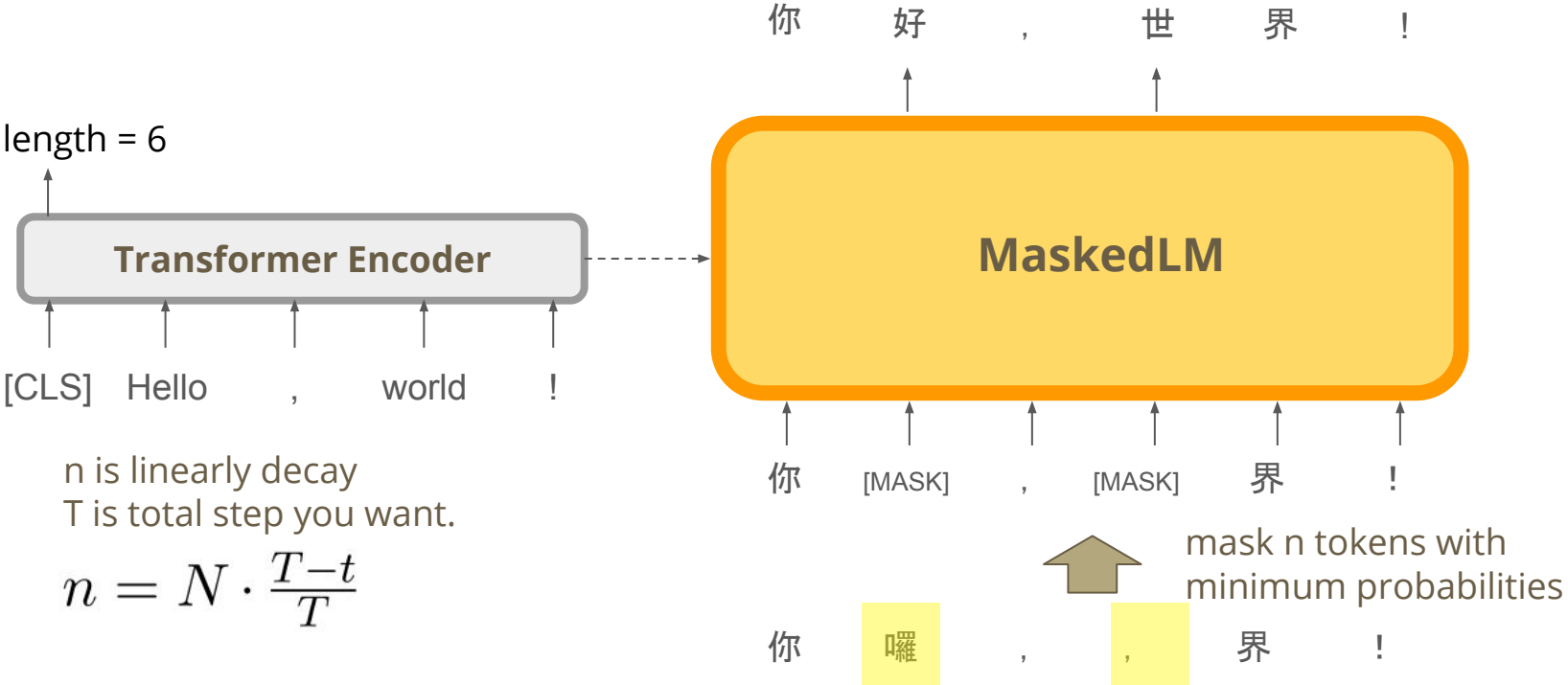
# Mask-Predict

$t = 0$



# Mask-Predict

t = 1



# Mask-Predict

**Example:** (target length is fixed)

---

*src* Der Abzug der franzsischen Kampftruppen wurde am 20. November abgeschlossen .

---

$t = 0$  The **departure of the French combat completed completed on** 20 November .

$t = 1$  The **departure** of French combat troops was **completed** on **20 November** .

$t = 2$  The withdrawal of French combat troops was completed on November 20th .

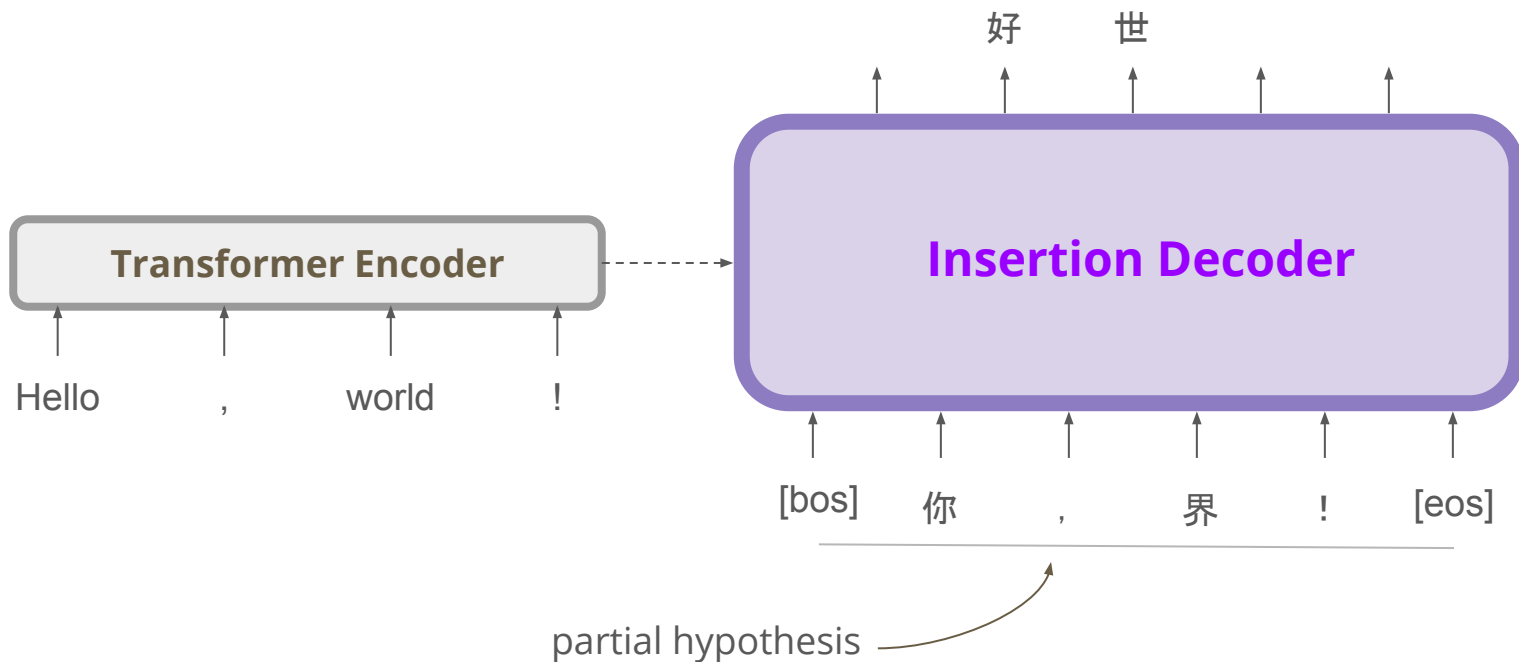
---

# Results

Model	Dimensions (Model/Hidden)	Iterations	WMT'14		WMT'16	
			EN-DE	DE-EN	EN-RO	RO-EN
NAT w/ Fertility (Gu et al., 2018)	512/512	1	19.17	23.20	29.79	31.44
CTC Loss (Libovický and Helcl, 2018)	512/4096	1	17.68	19.80	19.93	24.71
Iterative Refinement (Lee et al., 2018)	512/512	1	13.91	16.77	24.45	25.73
	512/512	10	21.61	25.48	29.32	30.19
(Dynamic #Iterations)	512/512	?	21.54	25.43	29.66	30.30
<i>Small CMLM with Mask-Predict</i>	512/512	1	15.06	19.26	20.12	20.36
	512/512	4	<b>24.17</b>	<b>28.55</b>	<b>30.00</b>	30.43
	512/512	10	<b>25.51</b>	<b>29.47</b>	<b>31.65</b>	<b>32.27</b>
<i>Base CMLM with Mask-Predict</i>	512/2048	1	18.05	21.83	27.32	28.20
	512/2048	4	<b>25.94</b>	<b>29.90</b>	<b>32.53</b>	<b>33.23</b>
	512/2048	10	<b>27.03</b>	<b>30.53</b>	<b>33.08</b>	<b>33.31</b>
Base Transformer (Vaswani et al., 2017)	512/2048	<i>N</i>	27.30	— —	— —	— —
Base Transformer (Our Implementation)	512/2048	<i>N</i>	27.74	31.09	34.28	33.99
Base Transformer (+Distillation)	512/2048	<i>N</i>	27.86	31.07	— —	— —
Large Transformer (Vaswani et al., 2017)	1024/4096	<i>N</i>	28.40	— —	— —	— —
Large Transformer (Our Implementation)	1024/4096	<i>N</i>	28.60	31.71	— —	— —

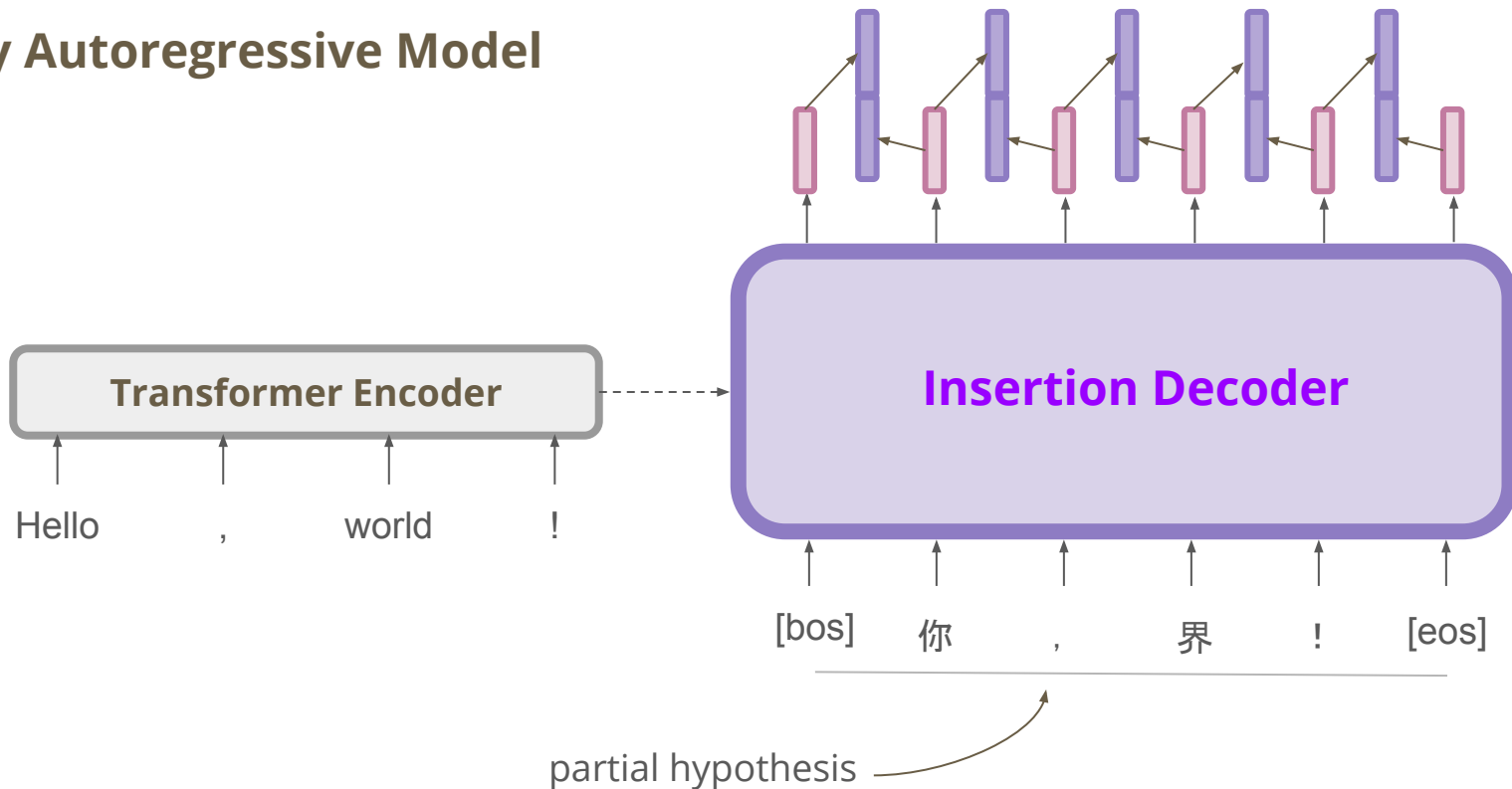
# Insertion Transformer

## Partially Autoregressive Model



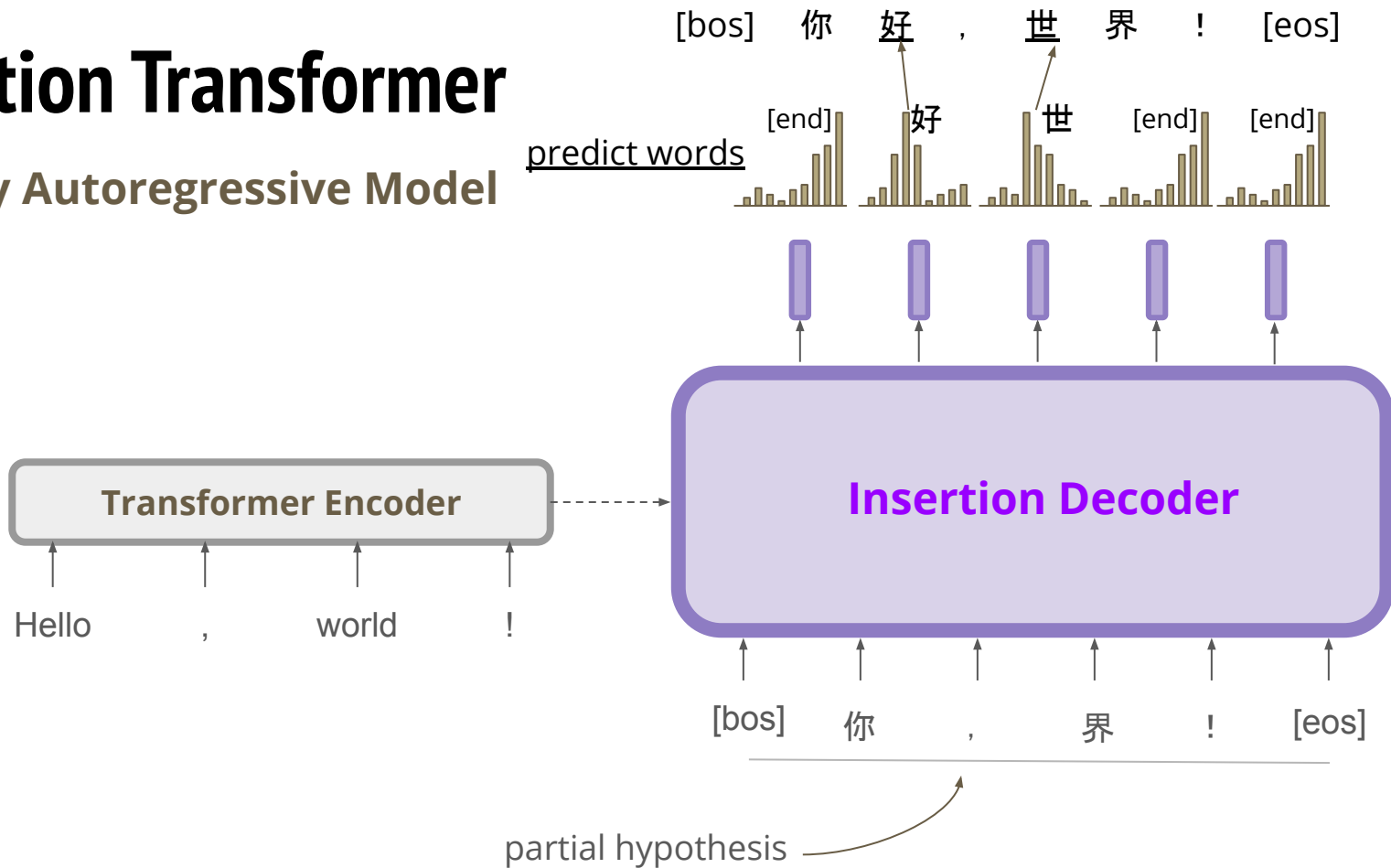
# Insertion Transformer

## Partially Autoregressive Model



# Insertion Transformer

## Partially Autoregressive Model



# Training Examples

origin:  $y_1$   $y_2$   $y_3$   $y_4$   $y_5$   $y_6$   $y_7$   $y_8$   $y_9$   $y_{10}$

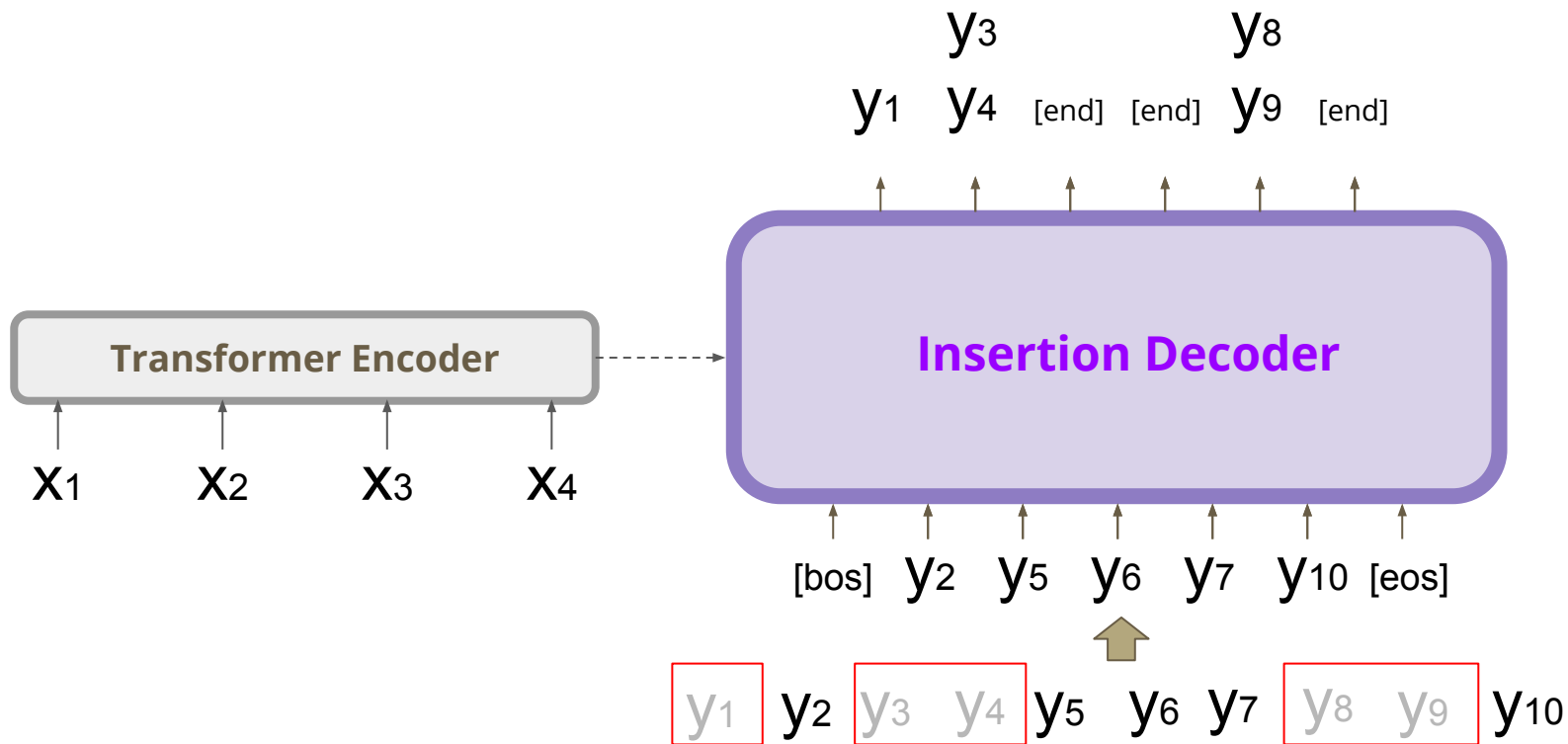
shuffled:  $y_7$   $y_{10}$   $y_2$   $y_6$   $y_5$   $y_8$   $y_9$   $y_4$   $y_3$   $y_1$

random sample  $k$ :  $y_7$   $y_{10}$   $y_2$   $y_6$   $y_5$   $y_8$   $y_9$   $y_4$   $y_3$   $y_1$   
 $k \sim \text{Uniform}([0, |y|])$

select first  $k$  words:  $y_1$   $y_2$   $y_3$   $y_4$   $y_5$   $y_6$   $y_7$   $y_8$   $y_9$   $y_{10}$

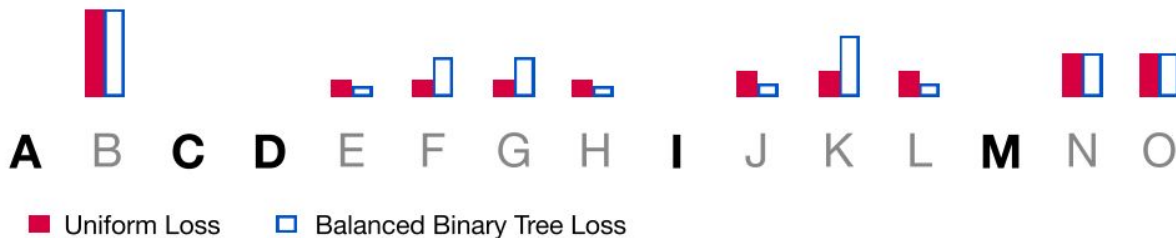
get slots to insert:  $y_1$   $y_2$   $y_3$   $y_4$   $y_5$   $y_6$   $y_7$   $y_8$   $y_9$   $y_{10}$

# Training Examples



$[\ ] \rightarrow [D] \rightarrow [B, D, F] \rightarrow [A, B, C, D, E, F, G]$

# Multiple target words to predict?



**Uniform Policy**  $\text{slot-loss}(x, \hat{y}, l) = \frac{1}{j_l - i_l + 1} \sum_{i=i_l}^{j_l} -\log p(y_i, l | x, \hat{y})$ .

## Balanced Binary Tree Policy

$$w_l(i) = \frac{\exp(-d_l(i)/\tau)}{\sum_{i'=i_l}^{j_l} \exp(-d_l(i')/\tau)} \Rightarrow \text{slot-loss}(x, \hat{y}, l) = \sum_{i=i_l}^{j_l} -\log p(y_i, l | x, \hat{y}) \cdot w_l(i).$$

softmax weighting policy

Middle word has priority!!!

# Example

**Input:** But on the other side of the state, that is not the impression many people have of their former governor.

**Output:** Aber auf der anderen Seite des Staates ist das nicht der Eindruck, den viele von ihrem ehemaligen Gouverneur haben.

**Parallel decode (binary tree loss):**

Aber\_ auf\_ der\_ anderen\_ Seite\_ des\_ Staates\_ ist\_ das\_ nicht\_ der\_ Eindruck\_, \_ den\_ viele\_ von\_ ihrem\_ ehemaligen\_ Gouverneur\_ haben\_ ..

Aber\_ auf\_ der\_ anderen\_ Seite\_ des\_ Staates\_ ist\_ das\_ nicht\_ der\_ Eindruck\_ , \_ den\_ viele\_ von\_ ihrem\_ ehemaligen\_ Gouverneur\_ haben\_ ..

Aber\_ auf\_ der\_ anderen\_ Seite\_ des\_ Staates\_ ist\_ das\_ nicht\_ der\_ Eindruck\_ , \_ den\_ viele\_ von\_ ihrem\_ ehemaligen\_ Gouverneur\_ haben\_ ..

Aber\_ auf\_ der\_ anderen\_ Seite\_ des\_ Staates\_ ist\_ das\_ nicht\_ der\_ Eindruck\_ , \_ den\_ viele\_ von\_ ihrem\_ ehemaligen\_ Gouverneur\_ haben\_ ..

Aber\_ auf\_ der\_ anderen\_ Seite\_ des\_ Staates\_ ist\_ das\_ nicht\_ der\_ Eindruck\_ , \_ den\_ viele\_ von\_ ihrem\_ ehemaligen\_ Gouverneur\_ haben\_ ..

**Input:** They want to create a post on the college's equal opportunities committee to ensure that their opinions can be aired freely.

**Output:** Sie wollen einen Posten im Ausschuss für Chancengleichheit des Kollegiums einrichten, um sicherzustellen, dass ihre Meinungen frei zur Sprache gebracht werden können.

**Parallel decode (uniform loss):**

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ .. um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ .. dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

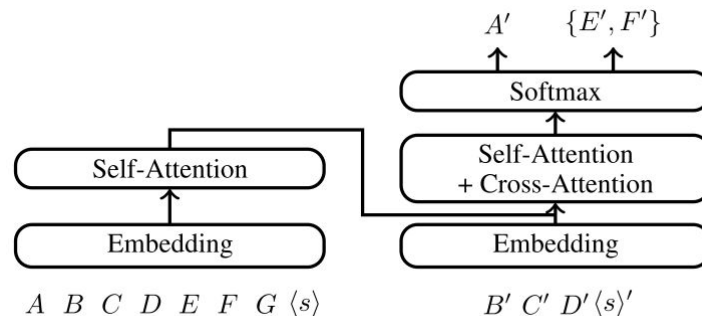
Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

Sie\_ wollen\_ einen\_ Posten\_ im\_ Ausschuss\_ für\_ Chancengleichheit\_ des\_ Koll\_egi\_ ums\_ einrichten\_ , \_ um\_ sicherzustellen\_ , \_ dass\_ ihre\_ Mein\_ ungen\_ frei\_ zur\_ Sprache\_ gebracht\_ werden\_ können\_ ..

# Decode stage



Serial generation:

$t$	Canvas	Insertion
0	[]	(ate, 0)
1	[ <u>ate</u> ]	(together, 1)
2	[ate, <u>together</u> ]	(friends, 0)
3	[ <u>friends</u> , ate, together]	(three, 0)
4	[ <u>three</u> , friends, ate, together]	(lunch, 3)
5	[three, friends, ate, <u>lunch</u> , together]	(⟨EOS⟩, 5)

$n$  operations

Parallel generation:

$t$	Canvas	Insertions
0	[]	(ate, 0)
1	[ <u>ate</u> ]	(friends, 0), (together, 1)
2	[ <u>friends</u> , ate, <u>together</u> ]	(three, 0), (lunch, 2)
3	[ <u>three</u> , friends, ate, <u>lunch</u> , together]	(⟨EOS⟩, 5)

$\lfloor \log_2 n \rfloor + 1$  operations

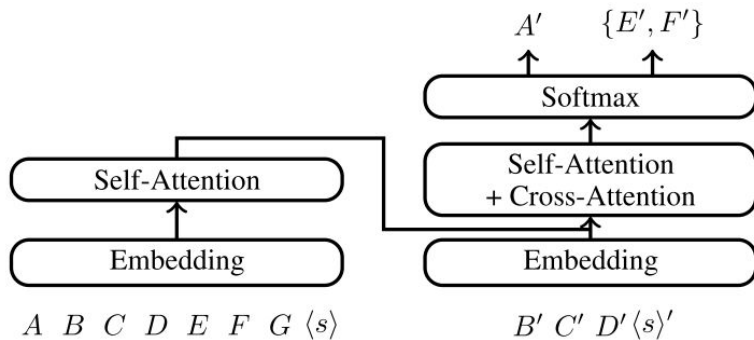
# Comparison

Model	BLEU	Iterations
Autoregressive Left-to-Right Transformer (Vaswani et al., 2017)	27.3	$n$
Semi-Autoregressive Left-to-Right SAT (Wang et al., 2018)	24.83	$n/6$
Blockwise Parallel (Stern et al., 2018)	27.40	$\approx n/5$
Non-Autoregressive NAT (Gu et al., 2018)	17.69	1
Iterative Refinement (Lee et al., 2018)	21.61	10
Our Approach (Greedy)		
Insertion Transformer + Left-to-Right	23.94	$n$
Insertion Transformer + Binary Tree	27.29	$n$
Insertion Transformer + Uniform	27.12	$n$
Our Approach (Parallel)		
Insertion Transformer + Binary Tree	27.41	$\approx \log_2 n$
Insertion Transformer + Uniform	26.72	$\approx \log_2 n$

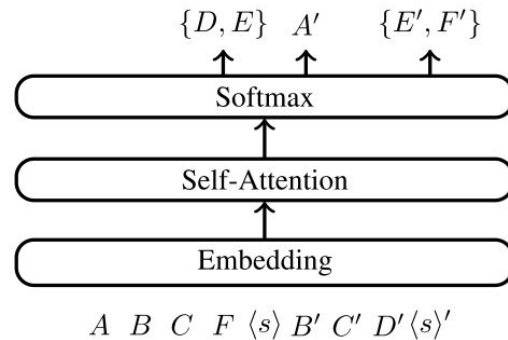


From  $P(y|x)$  to  $P(x, y)$ ,  $P(x)$ ,  $P(y)$ ,  $P(x|y)$ ,  $P(y|x)$

# KERMIT: Kontextuell Encoder Representations Made by Insertion Transformations



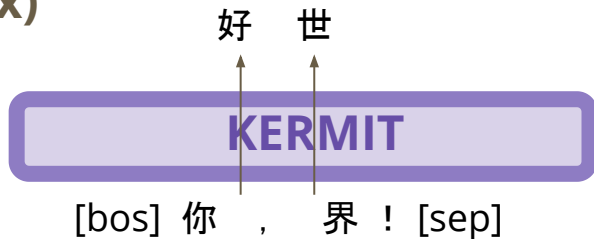
(c) Insertion Transformer



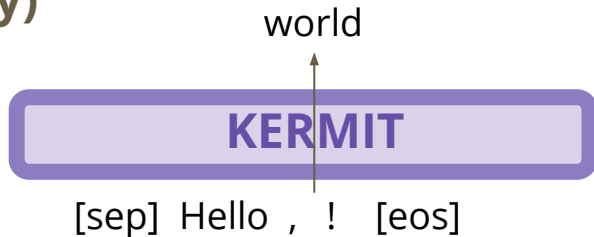
(d) KERMIT

# Advantages

$P(x)$



$P(y)$



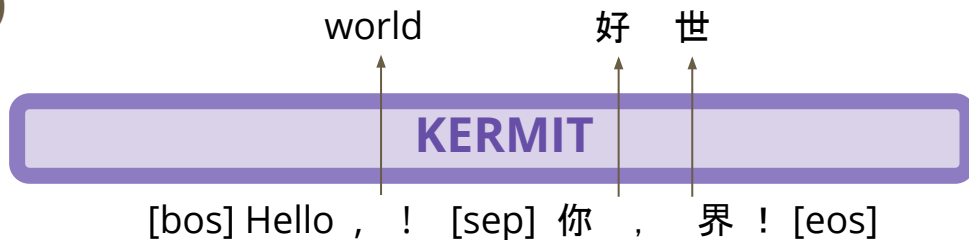
$P(y|x)$



$P(x|y)$



$P(x, y)$



# Results

Model	$\leftrightarrow$	En $\rightarrow$ De	De $\rightarrow$ En	Iterations
Autoregressive				
Transformer (Vaswani et al., 2017)	$\times$	27.3		$n$
Transformer (Our Implementation)	$\times$	27.8	31.2	$n$
Non-Autoregressive				
NAT (Gu et al., 2018)	$\times$	17.7	21.5	1
Iterative Refinement (Lee et al., 2018)	$\times$	21.6	25.5	10
Blockwise Parallel (Stern et al., 2018)	$\times$	27.4		$\approx n/5$
Insertion Transformer (Stern et al., 2019)	$\times$	27.4		$\approx \log_2 n \ll 10$
KERMIT				
Unidirectional ( $p(y   x)$ or $p(x   y)$ )	$\times$	27.8	30.7	$\approx \log_2 n \ll 10$
Bidirectional ( $p(y   x)$ and $p(x   y)$ )	$\checkmark$	27.2	27.6	$\approx \log_2 n \ll 10$
Joint ( $p(x, y)$ )	$\checkmark$	25.6	27.4	$\approx \log_2 n \ll 10$
+ Marginal Refining ( $p(x)$ and $p(y)$ )	$\checkmark$	25.8	28.6	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Unidirectional Finetuning	$\times$	28.7	31.4	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Bidirectional Finetuning	$\checkmark$	28.1	28.6	$\approx \log_2 n \ll 10$

# Results

Model	$\leftrightarrow$	En $\rightarrow$ De	De $\rightarrow$ En	Iterations
Autoregressive				
Transformer (Vaswani et al., 2017)	$\times$	27.3		$n$
Transformer (Our Implementation)	$\times$	27.8	31.2	$n$
Non-Autoregressive				
NAT (Gu et al., 2018)	$\times$	17.7	21.5	1
Iterative Refinement (Lee et al., 2018)	$\times$	21.6	25.5	10
Blockwise Parallel (Stern et al., 2018)	$\times$	27.4		$\approx n/5$
Insertion Transformer (Stern et al., 2019)	$\times$	27.4		$\approx \log_2 n \ll 10$
KERMIT				
Unidirectional ( $p(y   x)$ or $p(x   y)$ )	$\times$	27.8	30.7	$\approx \log_2 n \ll 10$
Bidirectional ( $p(y   x)$ and $p(x   y)$ )	$\checkmark$	27.2	27.6	$\approx \log_2 n \ll 10$
Joint ( $p(x, y)$ )	$\checkmark$	25.6	27.4	$\approx \log_2 n \ll 10$
+ Marginal Refining ( $p(x)$ and $p(y)$ )	$\checkmark$	25.8	28.6	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Unidirectional Finetuning	$\times$	28.7	31.4	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Bidirectional Finetuning	$\checkmark$	28.1	28.6	$\approx \log_2 n \ll 10$

# Results

Model	$\leftrightarrow$	En $\rightarrow$ De	De $\rightarrow$ En	Iterations
Autoregressive				
Transformer (Vaswani et al., 2017)	$\times$	27.3		$n$
Transformer (Our Implementation)	$\times$	27.8	31.2	$n$
Non-Autoregressive				
NAT (Gu et al., 2018)	$\times$	17.7	21.5	1
Iterative Refinement (Lee et al., 2018)	$\times$	21.6	25.5	10
Blockwise Parallel (Stern et al., 2018)	$\times$	27.4		$\approx n/5$
Insertion Transformer (Stern et al., 2019)	$\times$	27.4		$\approx \log_2 n \ll 10$
KERMIT				
Unidirectional ( $p(y   x)$ or $p(x   y)$ )	$\times$	27.8	30.7	$\approx \log_2 n \ll 10$
Bidirectional ( $p(y   x)$ and $p(x   y)$ )	$\checkmark$	27.2	27.6	$\approx \log_2 n \ll 10$
Joint ( $p(x, y)$ )	$\checkmark$	25.6	27.4	$\approx \log_2 n \ll 10$
+ Marginal Refining ( $p(x)$ and $p(y)$ )	$\checkmark$	25.8	28.6	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Unidirectional Finetuning	$\times$	28.7	31.4	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Bidirectional Finetuning	$\checkmark$	28.1	28.6	$\approx \log_2 n \ll 10$

# Results

Model	$\leftrightarrow$	En $\rightarrow$ De	De $\rightarrow$ En	Iterations
Autoregressive				
Transformer (Vaswani et al., 2017)	$\times$	27.3		$n$
Transformer (Our Implementation)	$\times$	27.8	31.2	$n$
Non-Autoregressive				
NAT (Gu et al., 2018)	$\times$	17.7	21.5	1
Iterative Refinement (Lee et al., 2018)	$\times$	21.6	25.5	10
Blockwise Parallel (Stern et al., 2018)	$\times$	27.4		$\approx n/5$
Insertion Transformer (Stern et al., 2019)	$\times$	27.4		$\approx \log_2 n \ll 10$
KERMIT				
Unidirectional ( $p(y   x)$ or $p(x   y)$ )	$\times$	27.8	30.7	$\approx \log_2 n \ll 10$
Bidirectional ( $p(y   x)$ and $p(x   y)$ )	$\checkmark$	27.2	27.6	$\approx \log_2 n \ll 10$
Joint ( $p(x, y)$ )	$\checkmark$	25.6	27.4	$\approx \log_2 n \ll 10$
+ Marginal Refining ( $p(x)$ and $p(y)$ )	$\checkmark$	25.8	28.6	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Unidirectional Finetuning	$\times$	28.7	31.4	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Bidirectional Finetuning	$\checkmark$	28.1	28.6	$\approx \log_2 n \ll 10$

# Results

Model	$\leftrightarrow$	En $\rightarrow$ De	De $\rightarrow$ En	Iterations
Autoregressive				
Transformer (Vaswani et al., 2017)	$\times$	27.3		$n$
Transformer (Our Implementation)	$\times$	27.8	31.2	$n$
Non-Autoregressive				
NAT (Gu et al., 2018)	$\times$	17.7	21.5	1
Iterative Refinement (Lee et al., 2018)	$\times$	21.6	25.5	10
Blockwise Parallel (Stern et al., 2018)	$\times$	27.4		$\approx n/5$
Insertion Transformer (Stern et al., 2019)	$\times$	27.4		$\approx \log_2 n \ll 10$
KERMIT				
Unidirectional ( $p(y   x)$ or $p(x   y)$ )	$\times$	27.8	30.7	$\approx \log_2 n \ll 10$
Bidirectional ( $p(y   x)$ and $p(x   y)$ )	$\checkmark$	27.2	27.6	$\approx \log_2 n \ll 10$
Joint ( $p(x, y)$ )	$\checkmark$	25.6	27.4	$\approx \log_2 n \ll 10$
+ Marginal Refining ( $p(x)$ and $p(y)$ )	$\checkmark$	25.8	28.6	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Unidirectional Finetuning	$\times$	28.7	31.4	$\approx \log_2 n \ll 10$
$\hookrightarrow$ Bidirectional Finetuning	$\checkmark$	28.1	28.6	$\approx \log_2 n \ll 10$

# GLUE

Model	Generative?	CoLA	SST-2	MRPC	STS-B	QQP	MNLI-(m/mm)	QNLI	RTE	WNLI	AX	Score
GPT (Radford et al., 2018)	✓	45.4	91.3	82.3/75.7	82.0/80.0	70.3/88.5	82.1/81.4	87.4	56.0	53.4	29.8	72.8
BERT (Devlin et al., 2019)	✗	60.5	94.9	89.3/85.4	87.6/86.5	72.1/89.3	86.7/85.9	92.7	70.1	65.1	39.6	80.5
KERMIT	✓	60.0	94.2	88.6/84.3	86.6/85.6	71.7/89.0	85.6/85.2	92.0	68.4	65.1	37.6	79.8

# Zero-shot ClozeQA

Model	Exact Match	F1
GPT-2	10.9	16.6
+ Oracle Length	12.2	18.2
BERT	12.3	18.9
+ Oracle Length	16.2	23.1
KERMIT	<b>20.9</b>	<b>30.3</b>

Table 4: SQuAD zero-shot cloze question answering.

Plymouth has a post-war shopping area in the city centre with substantial pedestrianisation. At the west end of the zone inside a grade II listed building is the Pannier Market that was completed in 1959 – pannier meaning "basket" from French, so it translates as "basket market". In terms of retail floorspace, Plymouth is ranked in the top five in the South West, and 29th nationally. Plymouth was one of the first ten British cities to trial the new Business Improvement District initiative. The Tinside Pool is situated at the foot of the Hoe and became a grade II listed building in 1998 before being restored to its 1930s look for £3.4 million. **What notable location was named a grade II listed building in 1998? \_\_\_ was named a grade II listed building in 1998.**

Model	Answer
GPT-2	→ "Plymouth"
+ Oracle Length	→ "A listed building"
BERT	→ "plymouth"
+ Oracle Length	→ " : the pool"
KERMIT	→ <b>"the tinside pool"</b>
Correct	→ "Tinside Pool"

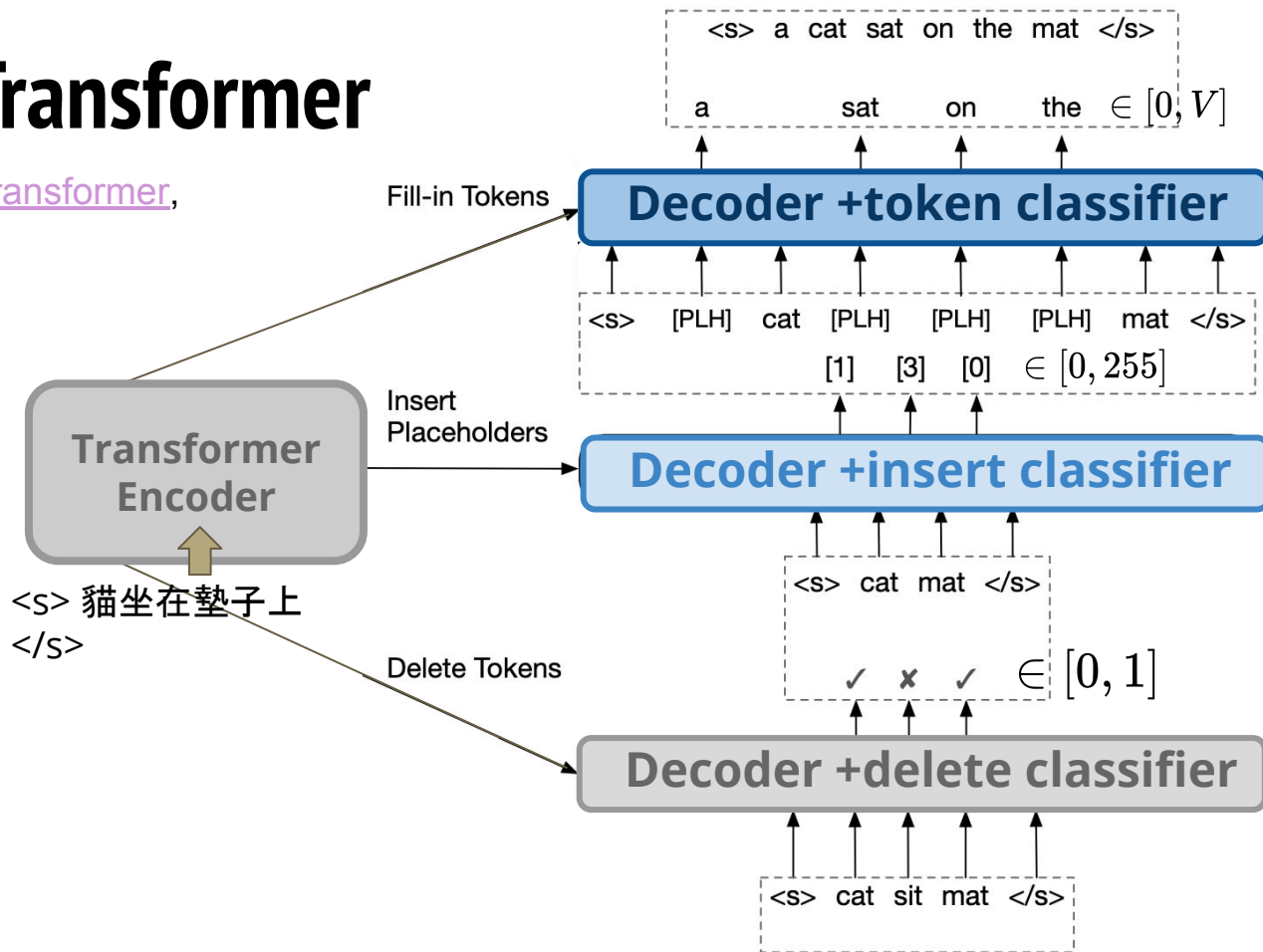
# Multilingual KERMIT: It's Not Easy Being Generative

*KERMIT can be trained by feeding it  $N$  way parallel-data, bilingual data, or monolingual data. At inference, KERMIT can generate translations for a particular target language, or up to  $N - 1$  languages in parallel.*



# Levenshtein Transformer

[1905.11006] Levenshtein Transformer,  
Gu et al., NeurIPS 2019.



How to train?

Imitation Learning

# Learn from expert: Levenshtein Distance Algorithm

```
>>> import Levenshtein
>>> Levenshtein.distance("ABCEFGHJJ", "ABCDEFGHI")
3
>>> Levenshtein.editops("ABCEFGHJJ", "ABCDEFGHI")
[('insert', 3, 3), ('delete', 7, 8), ('replace', 8, 8)]
```

*In this paper we decompose 'replace' into delete+insert.*

# Imitation Learning

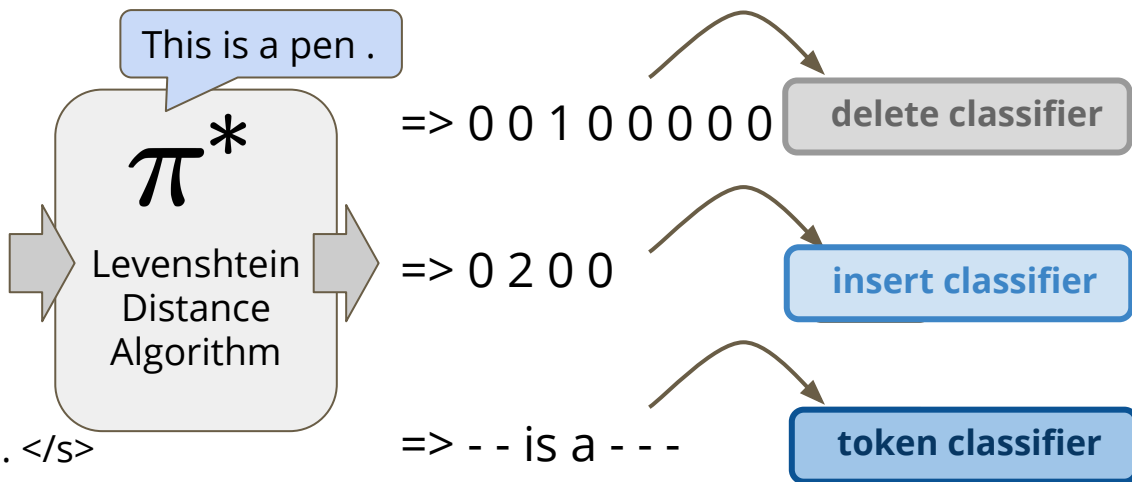
$y_{del}$  sentences that need to delete words

$y_{ins}$  sentences that need to insert words

$y_{del}$  => <s> This is is a pen . </s>

$y_{ins}$  => <s> This pen . </s>

$y'_{ins}$  => <s> This [PLH] [PLH] pen . </s>





# Results

	Dataset	Metric	Transformer		Levenshtein Transformer	
			greedy	beam4	oracle	teacher
Quality $\uparrow$	Ro-En	BLEU	31.67	32.30	<b>33.02</b>	—
	En-De	BLEU	26.02	26.56	24.43	<b>26.67</b>
	En-Ja	BLEU	42.86	<b>43.68</b>	42.36	43.17
		ROUGE-1	34.91	35.19	35.57	<b>36.08</b>
	Gigaword	ROUGE-2	17.05	17.58	17.11	<b>18.33</b>
		ROUGE-L	32.66	32.98	33.55	<b>33.81</b>
Speed $\downarrow$	Ro-En	Latency (ms) / $I_{\text{DEC}}$	326 / 27.1	349 / 27.1	<b>97 / 2.19</b>	—
	En-De	Latency (ms) / $I_{\text{DEC}}$	343 / 28.1	369 / 28.1	126 / 2.88	<b>92 / 2.05</b>
	En-Ja	Latency (ms) / $I_{\text{DEC}}$	261 / 22.6	306 / 22.6	112 / 2.61	<b>106 / 1.97</b>
	Gigaword	Latency (ms) / $I_{\text{DEC}}$	116 / 10.1	149 / 10.1	98 / 2.32	<b>84 / 1.73</b>

Table 1: Generation quality (BLEU  $\uparrow$ , ROUGE-1/2/L  $\uparrow$ ) and latency (ms  $\downarrow$ ) as well as the average number of decoder iterations ( $I_{\text{DEC}}$ ) on the standard test sets for LevT and the autoregressive baseline (with both greedy and beam-search outputs). We show the results of LevT trained from both oracle and the autoregressive teacher model.

# Results

	Dataset	Metric	Transformer		Levenshtein	Transformer
			greedy	beam4	oracle	teacher
Quality $\uparrow$	Ro-En	BLEU	31.67	32.30	<b>33.02</b>	—
	En-De	BLEU	26.02	26.56	24.43	<b>26.67</b>
	En-Ja	BLEU	42.86	<b>43.68</b>	42.36	43.17
		ROUGE-1	34.91	35.19	35.57	<b>36.08</b>
	Gigaword	ROUGE-2	17.05	17.58	17.11	<b>18.33</b>
		ROUGE-L	32.66	32.98	33.55	<b>33.81</b>
Speed $\downarrow$	Ro-En	Latency (ms) / $I_{\text{DEC}}$	326 / 27.1	349 / 27.1	<b>97 / 2.19</b>	—
	En-De	Latency (ms) / $I_{\text{DEC}}$	343 / 28.1	369 / 28.1	126 / 2.88	<b>92 / 2.05</b>
	En-Ja	Latency (ms) / $I_{\text{DEC}}$	261 / 22.6	306 / 22.6	112 / 2.61	<b>106 / 1.97</b>
	Gigaword	Latency (ms) / $I_{\text{DEC}}$	116 / 10.1	149 / 10.1	98 / 2.32	<b>84 / 1.73</b>

Table 1: Generation quality (BLEU  $\uparrow$ , ROUGE-1/2/L  $\uparrow$ ) and latency (ms  $\downarrow$ ) as well as the average number of decoder iterations ( $I_{\text{DEC}}$ ) on the standard test sets for LevT and the autoregressive baseline (with both greedy and beam-search outputs). We show the results of LevT trained from both oracle and the autoregressive teacher model.

# Results

	Dataset	Metric	Transformer		Levenshtein Transformer	
			greedy	beam4	oracle	teacher
Quality $\uparrow$	Ro-En	BLEU	31.67	32.30	<b>33.02</b>	—
	En-De	BLEU	26.02	26.56	24.43	<b>26.67</b>
	En-Ja	BLEU	42.86	<b>43.68</b>	42.36	43.17
		ROUGE-1	34.91	35.19	35.57	<b>36.08</b>
	Gigaword	ROUGE-2	17.05	17.58	17.11	<b>18.33</b>
		ROUGE-L	32.66	32.98	33.55	<b>33.81</b>
Speed $\downarrow$	Ro-En	Latency (ms) / $I_{\text{DEC}}$	326 / 27.1	349 / 27.1	<b>97 / 2.19</b>	—
	En-De	Latency (ms) / $I_{\text{DEC}}$	343 / 28.1	369 / 28.1	126 / 2.88	<b>92 / 2.05</b>
	En-Ja	Latency (ms) / $I_{\text{DEC}}$	261 / 22.6	306 / 22.6	112 / 2.61	<b>106 / 1.97</b>
	Gigaword	Latency (ms) / $I_{\text{DEC}}$	116 / 10.1	149 / 10.1	98 / 2.32	<b>84 / 1.73</b>

Table 1: Generation quality (BLEU  $\uparrow$ , ROUGE-1/2/L  $\uparrow$ ) and latency (ms  $\downarrow$ ) as well as the average number of decoder iterations ( $I_{\text{DEC}}$ ) on the standard test sets for LevT and the autoregressive baseline (with both greedy and beam-search outputs). We show the results of LevT trained from both oracle and the autoregressive teacher model.

# Results

	Dataset	Metric	Transformer		Levenshtein	Transformer
			greedy	beam4	oracle	teacher
Quality $\uparrow$	Ro-En	BLEU	31.67	32.30	<b>33.02</b>	—
	En-De	BLEU	26.02	26.56	24.43	<b>26.67</b>
	En-Ja	BLEU	42.86	<b>43.68</b>	42.36	43.17
		ROUGE-1	34.91	35.19	35.57	<b>36.08</b>
	Gigaword	ROUGE-2	17.05	17.58	17.11	<b>18.33</b>
		ROUGE-L	32.66	32.98	33.55	<b>33.81</b>
Speed $\downarrow$	Ro-En	Latency (ms) / $I_{\text{DEC}}$	326 / 27.1	349 / 27.1	<b>97 / 2.19</b>	—
	En-De	Latency (ms) / $I_{\text{DEC}}$	343 / 28.1	369 / 28.1	126 / 2.88	<b>92 / 2.05</b>
	En-Ja	Latency (ms) / $I_{\text{DEC}}$	261 / 22.6	306 / 22.6	112 / 2.61	<b>106 / 1.97</b>
	Gigaword	Latency (ms) / $I_{\text{DEC}}$	116 / 10.1	149 / 10.1	98 / 2.32	<b>84 / 1.73</b>

Table 1: Generation quality (BLEU  $\uparrow$ , ROUGE-1/2/L  $\uparrow$ ) and latency (ms  $\downarrow$ ) as well as the average number of decoder iterations ( $I_{\text{DEC}}$ ) on the standard test sets for LevT and the autoregressive baseline (with both greedy and beam-search outputs). We show the results of LevT trained from both oracle and the autoregressive teacher model.

# Results

	Dataset	Metric	Transformer		Levenshtein Transformer	
			greedy	beam4	oracle	teacher
Quality $\uparrow$	Ro-En	BLEU	31.67	32.30	<b>33.02</b>	—
	En-De	BLEU	26.02	26.56	24.43	<b>26.67</b>
	En-Ja	BLEU	42.86	<b>43.68</b>	42.36	43.17
		ROUGE-1	34.91	35.19	35.57	<b>36.08</b>
	Gigaword	ROUGE-2	17.05	17.58	17.11	<b>18.33</b>
		ROUGE-L	32.66	32.98	33.55	<b>33.81</b>
Speed $\downarrow$	Ro-En	Latency (ms) / $I_{\text{DEC}}$	326 / 27.1	349 / 27.1	<b>97 / 2.19</b>	—
	En-De	Latency (ms) / $I_{\text{DEC}}$	343 / 28.1	369 / 28.1	126 / 2.88	<b>92 / 2.05</b>
	En-Ja	Latency (ms) / $I_{\text{DEC}}$	261 / 22.6	306 / 22.6	112 / 2.61	<b>106 / 1.97</b>
	Gigaword	Latency (ms) / $I_{\text{DEC}}$	116 / 10.1	149 / 10.1	98 / 2.32	<b>84 / 1.73</b>

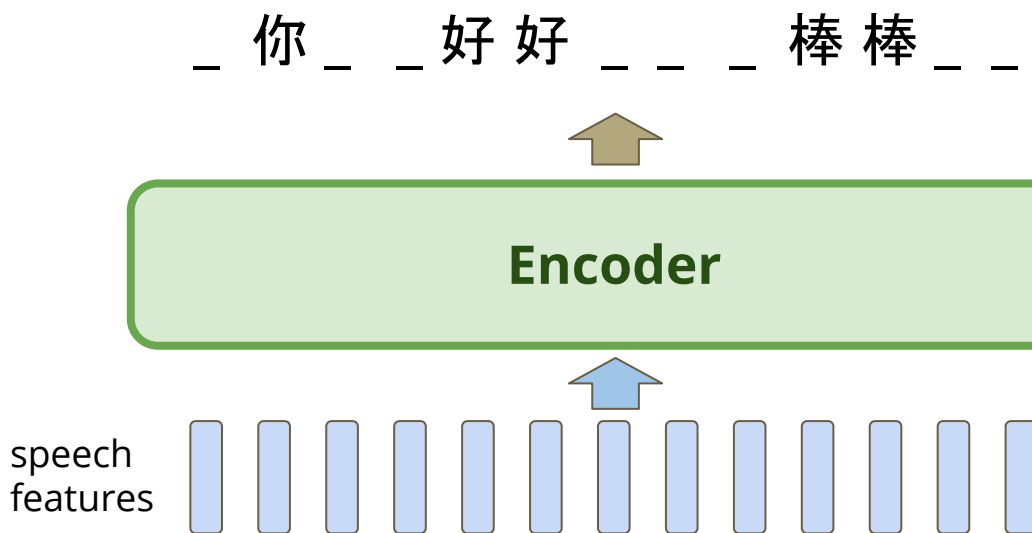
Table 1: Generation quality (BLEU  $\uparrow$ , ROUGE-1/2/L  $\uparrow$ ) and latency (ms  $\downarrow$ ) as well as the average number of decoder iterations ( $I_{\text{DEC}}$ ) on the standard test sets for LevT and the autoregressive baseline (with both greedy and beam-search outputs). We show the results of LevT trained from both oracle and the autoregressive teacher model.

# CTC

- also non-autoregressive model (for speech recognition)
- speech data doesn't have serious multi-modality problem!

shortcoming:

1. Falls behind LAS
2. Can't be refined



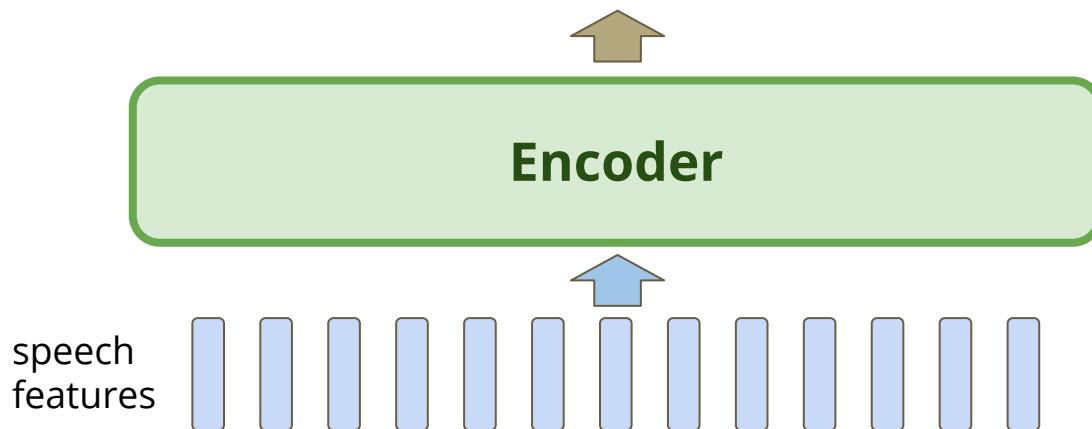
# CTC

- also non-autoregressive model (for speech recognition)
- speech data doesn't have serious multi-modality problem!

\_ 你 \_ \_ 好好 \_ \_ \_ 棒 棒 \_ **棒**

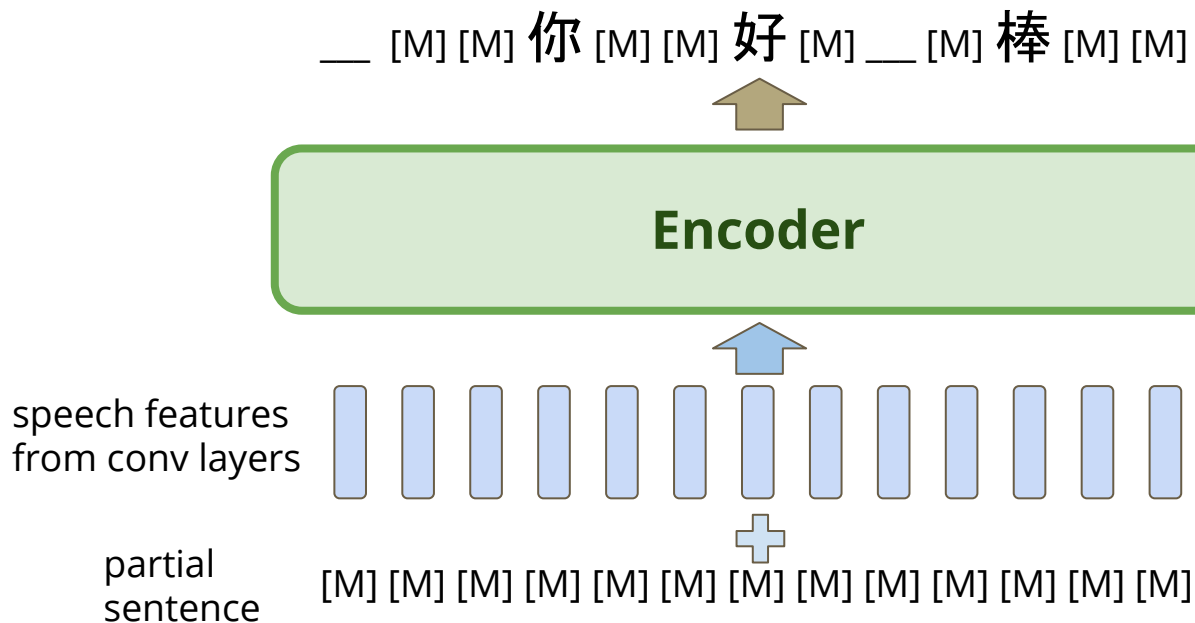
shortcoming:

1. Falls behind LAS
2. Can't be refined



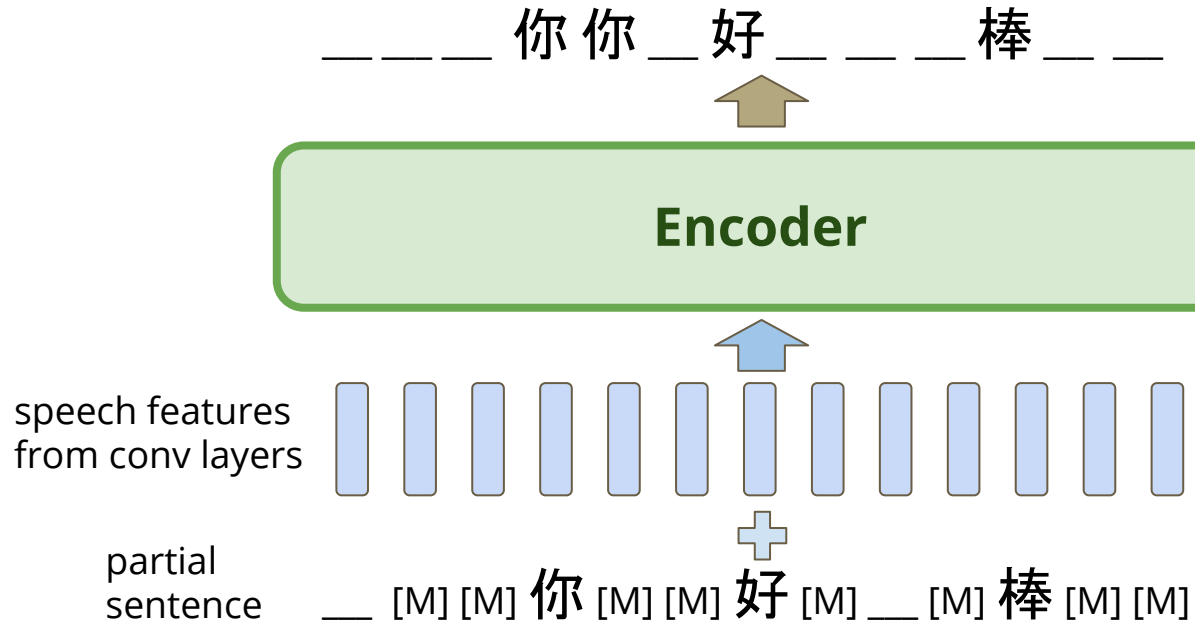
# Imputer (CTC+Mask-Predict)

t = 0



# Imputer (CTC+Mask-Predict)

t = 1

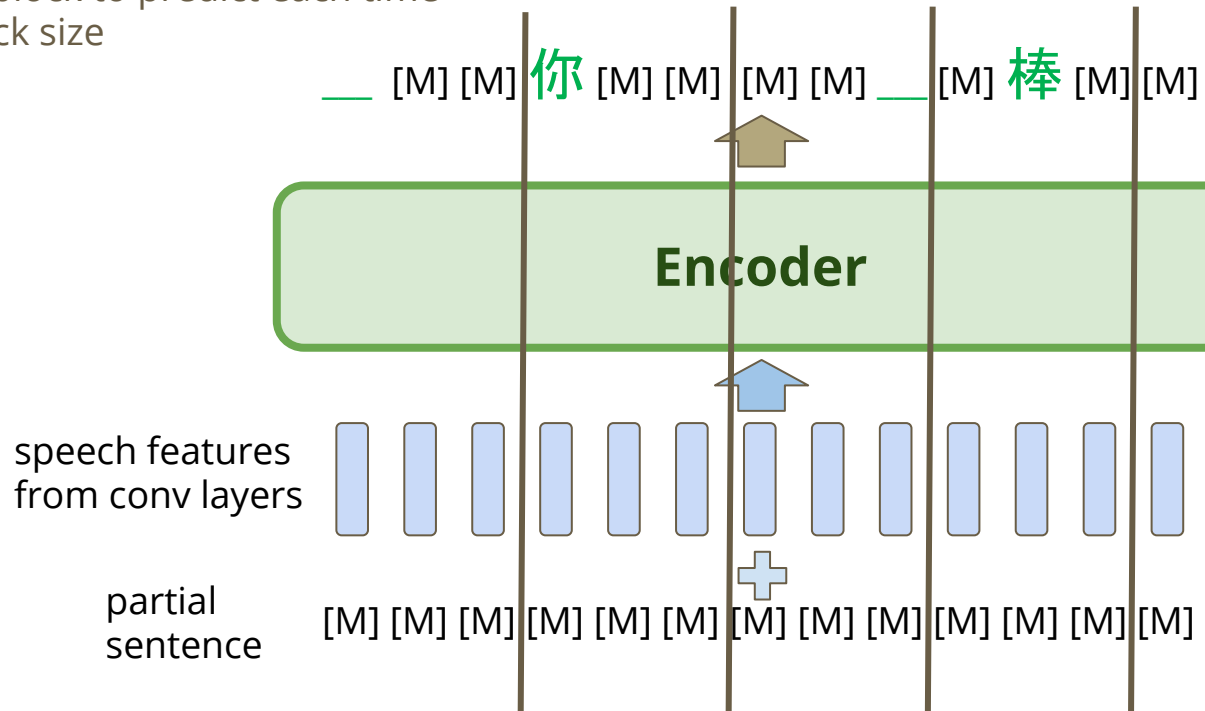


# Block Decoding

ex: Decode Block Size = 3

pick one position in a block to predict each time

=> max iteration = block size



# Example

b |

Alignment

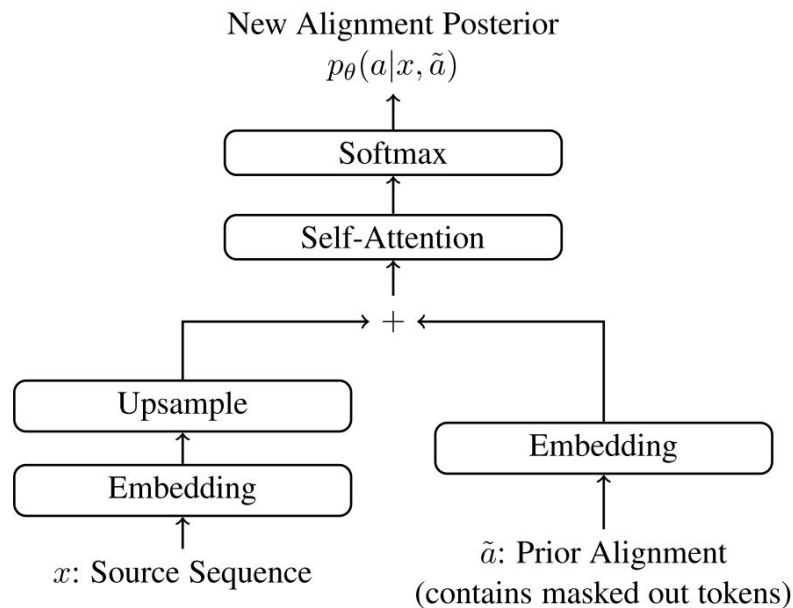
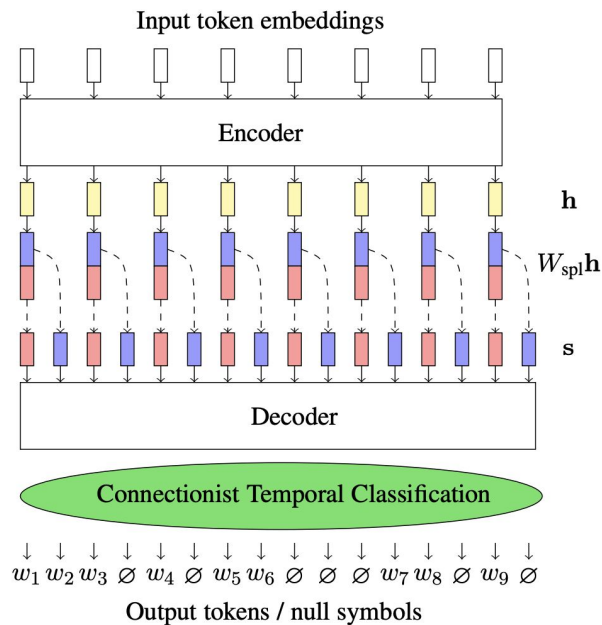
0	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅	∅
1	∅	∅	∅	[HAVE	∅	∅	∅	∅	∅	[L	∅	∅	∅	∅	∅	∅	∅	∅	Y	∅	∅	∅	∅	∅	∅	∅	∅
2	∅	∅	∅	[HAVE	∅	∅	[TO	∅	∅	[L	∅	∅	IVE	∅	∅	∅	∅	∅	Y	∅	∅	[SE	∅	∅	∅	∅	∅
3	∅	∅	-	[HAVE	∅	∅	[TO	∅	-	[L	∅	∅	IVE	∅	∅	∅	[M	∅	Y	∅	∅	[SE	∅	-	∅	∅	∅
4	∅	∅	-	[HAVE	∅	∅	[TO	-	-	[L	∅	∅	IVE	-	∅	∅	[M	∅	Y	∅	∅	[SE	-	-	∅	∅	∅
5	-	∅	-	[HAVE	∅	∅	[TO	-	-	[L	∅	-	IVE	-	∅	∅	[M	-	Y	∅	∅	[SE	-	-	-	∅	∅
6	-	∅	-	[HAVE	∅	-	[TO	-	-	[L	∅	-	IVE	-	[WITH	∅	∅	[M	-	Y	∅	-	[SE	-	-	-	∅
7	-	-	-	[HAVE	∅	-	[TO	-	-	[L	-	-	IVE	-	[WITH	∅	-	[M	-	Y	∅	-	[SE	-	-	-	∅
8	-	-	-	[HAVE	-	-	[TO	-	-	[L	-	-	IVE	-	[WITH	-	-	[M	-	Y	-	-	[SE	-	-	-	-

# Results

Table 1. Wall Street Journal Character Error Rate (CER) and Word Error Rate (WER).

Model	CER	WER	Iterations
seq2seq			
Bahdanau et al. (2016a)	6.4	18.6	n
Bahdanau et al. (2016b)	5.9	18.0	n
Chorowski & Jaitly (2017)	-	10.6	n
Zhang et al. (2017)	-	10.5	n
Chan et al. (2017)	-	9.6	n
Kim et al. (2017)	7.4	-	n
Serdyuk et al. (2018)	6.2	-	n
Tjandra et al. (2018)	6.1	-	n
Sabour et al. (2019)	3.1	9.3	n
CTC			
Graves & Jaitly (2014)	8.4	27.3	1
Liu et al. (2017)	-	16.7	1
CTC (Our Work)	5.6	15.2	1
Imputer (IM)	6.2	16.5	8
Imputer (DP)	4.9	12.7	8

# Applying CTC/Imputer on text generation



[1811.04719] [End-to-End Non-Autoregressive Neural Machine Translation with Connectionist Temporal Classification](#), Libovický et al., EMNLP 2018

[2004.07437] [Non-Autoregressive Machine Translation with Latent Alignments](#), Chan et al., arXiv 2020

# Results

compared with single step models

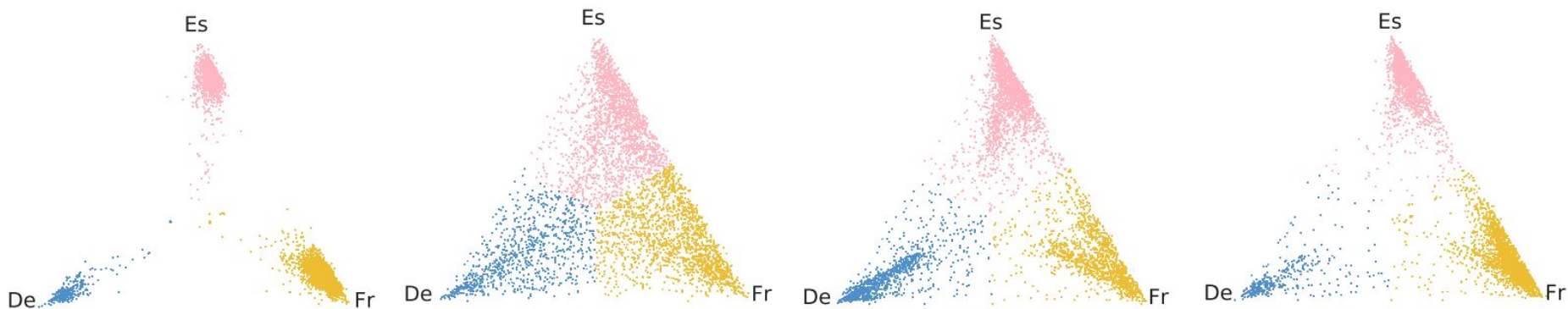
Method	Iterations	WMT'14		WMT'16	
		En→De	De→En	En→Ro	Ro→En
<i>Non-Autoregressive</i>					
Iterative Refinement (Lee et al., 2018)	1	13.9	16.7	24.5	25.7
NAT with Fertility (Gu et al., 2018)	1	17.7	21.5	27.3	29.1
CTC <sup>†</sup> (Libovicky and Helcl, 2018)	1	17.7	19.8	19.9	24.7
Mask-Predict (Ghazvininejad et al., 2019)	1	18.0	19.3	27.3	28.2
SMART (Ghazvininejad et al., 2020b)	1	18.6	23.8	-	-
Auxiliary Regularization (Wang et al., 2019)	1	20.7	24.8	-	-
Bag-of-ngrams Loss (Shao et al., 2020)	1	20.9	24.6	28.3	29.3
Hint-based Training (Li et al., 2019)	1	21.1	25.2	-	-
FlowSeq (Ma et al., 2019)	1	21.5	26.2	29.3	30.4
Bigram CRF (Sun et al., 2019)	1	23.4	27.2	-	-
AXE CMLM (Ghazvininejad et al., 2020a)	1	23.5	27.9	30.8	31.5
<i>Our Work</i>					
CTC <sup>†</sup>	1	25.7	28.1	32.2	31.6
Imputer	1	<b>25.8</b>	<b>28.4</b>	<b>32.3</b>	<b>31.7</b>

# Results

compared with multiple refine steps models

Method	Iterations	WMT'14		WMT'16	
		En→De	De→En	En→Ro	Ro→En
<i>Autoregressive</i>					
Base Transformer <sup>2</sup>	$n$	27.8	31.2	34.3	34.0
<i>Non-Autoregressive</i>					
Insertion Transformer (Stern et al., 2019)	$\approx \log_2 n$	27.4	-	-	-
KERMIT (Chan et al., 2019b)	$\approx \log_2 n$	27.8	30.7	-	-
Iterative Refinement (Lee et al., 2018)	10	21.6	25.5	29.3	30.2
Mask-Predict (Ghazvininejad et al., 2019)	4	25.9	29.9	32.5	33.2
	10	27.0	30.5	33.1	33.3
SMART (Ghazvininejad et al., 2020b)	4	27.0	30.9	-	-
	10	27.7	<b>31.3</b>	-	-
<i>Our Work</i>					
Imputer	2	27.5	30.2	33.7	33.4
	4	28.0	31.0	34.3	34.0
	8	<b>28.2</b>	<b>31.3</b>	<b>34.4</b>	<b>34.1</b>

# Knowledge Distillation in NAT



(a) AT Baseline

(b) NAT Baseline

(c) NAT Random Select

(d) NAT Distill

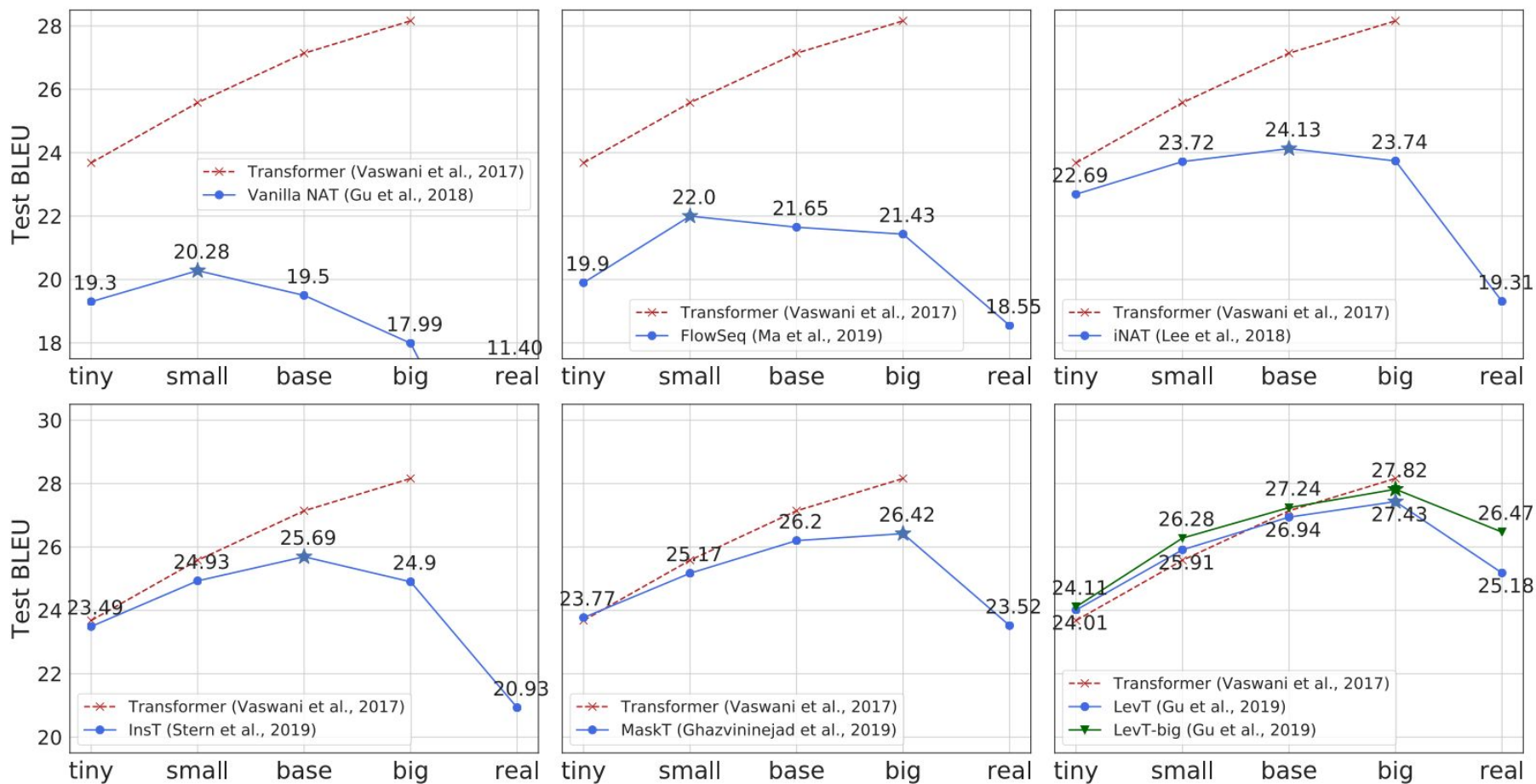


Figure 4: The performance of NAT models of varying capacity trained on both the real and the distilled data from tiny, small, base and big AT models.

# Reference



**Jiatao Gu**

@Facebook AI Research

<https://jiataogu.me/>

1. **Vanilla NAT** (Gu et al., 2018)  
*Non-Autoregressive Neural Machine Translation, ICLR 2018*
2. **Levenshtein Transformer** (LevT, Gu et al., 2019)  
*Levenshtein Transformer, NIPS 2019*
3. **Knowledge Distillation in NAT** (Zhou et al., 2019)  
*Understanding Knowledge Distillation in Non-autoregressive Machine Translation, ICLR 2020*

# Reference



William Chan @Google Brain  
& Mitchell Stern

<http://williamchan.ca/>

4. **Insertion Transformer** (InsT, Stern et al., 2019)  
*Insertion Transformer: Flexible Sequence Generation via Insertion Operations, ICML 2019*
5. **KERMIT** (Chan et al., 2019)  
*KERMIT: Generative Insertion-Based Modeling for Sequences, Preprint*
6. **Multilingual KERMIT** (Chan et al., 2019)  
*Multilingual KERMIT: It's Not Easy Being Generative, PGR@NeurIPS 2019*
7. **Imputer** (Chan et al., 2020)  
*Imputer: Sequence Modelling via Imputation and Dynamic Programming, Preprint*
8. **NAT with Latent Alignments** (Saharia et al., 2020)  
*Non-Autoregressive Machine Translation with Latent Alignments, Preprint*

*(The first author of "Listen, Attend, Spell")*

# Reference

## Others

9. **NAT with Iterative Refinement** (Lee et al., 2018)  
*Deterministic Non-Autoregressive Neural Sequence Modeling by Iterative Refinement, EMNLP 2018*
10. **MaskPredict** (Ghazvininejad et al., 2019)  
*Mask-Predict: Parallel Decoding of Conditional Masked Language Models, EMNLP 2019*
11. **CTC** (Graves et al., 2006)  
*Connectionist Temporal Classification : Labelling Unsegmented Sequence Data with Recurrent Neural Networks, ICML 2006*
12. **NAT with CTC** (Libovický et al., 2018)  
*End-to-End Non-Autoregressive Neural Machine Translation with Connectionist Temporal Classification, EMNLP 2018*