

# **Non-Autoregressive Decoding**

Xiachong Feng

# Outline

- Transformer
- The Importance of Generation Order in Language Modeling *EMNLP18*
- Insertion Transformer:  
Flexible Sequence Generation via Insertion Operations *ICML19*
- Non-Monotonic Sequential Text Generation *ICML19*
- Insertion-based Decoding with automatically Inferred Generation Order
- Levenshtein Transformer
- Conclusion
- Paper List
- Reference

# Transformer

# Transformer

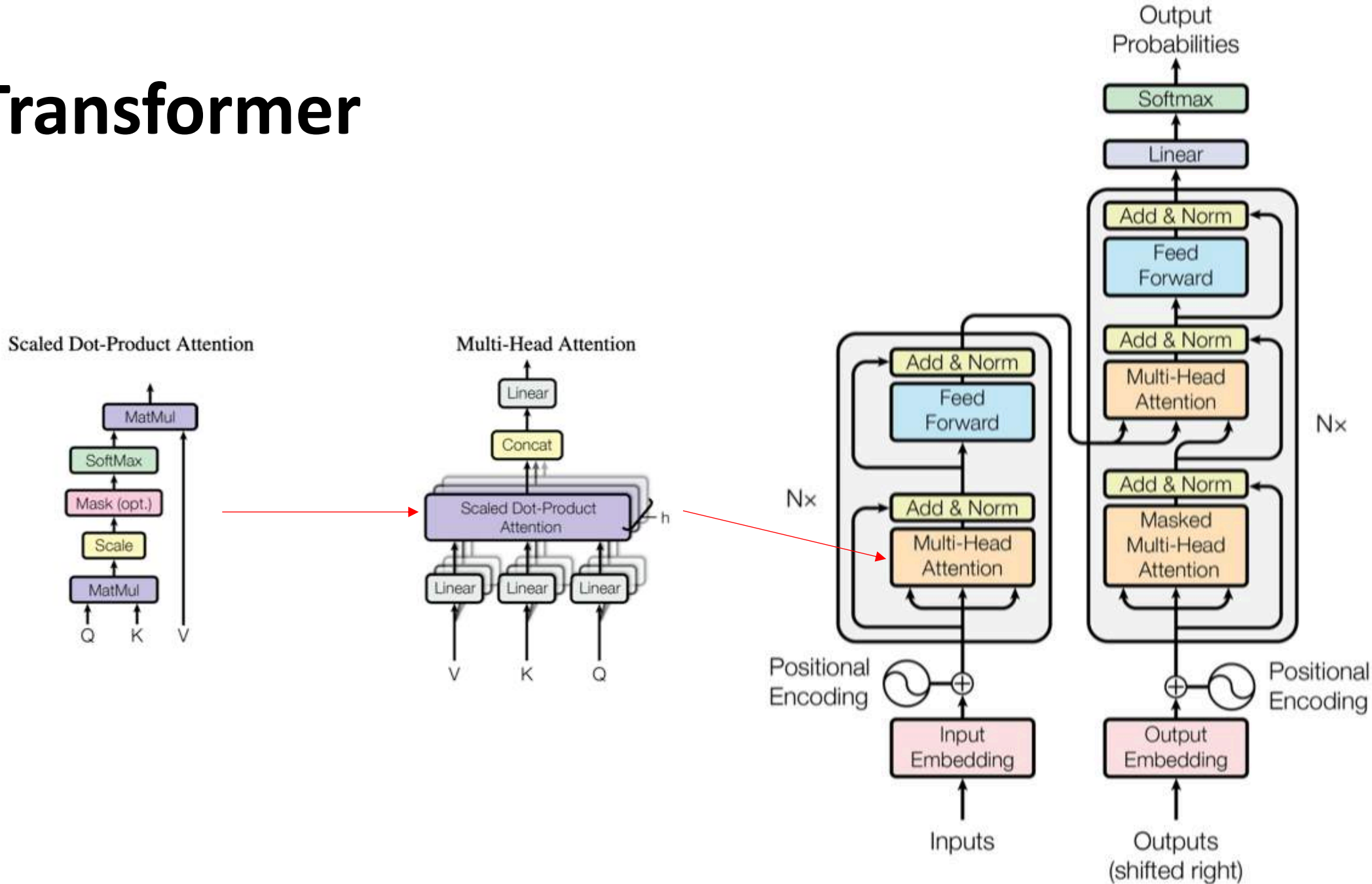
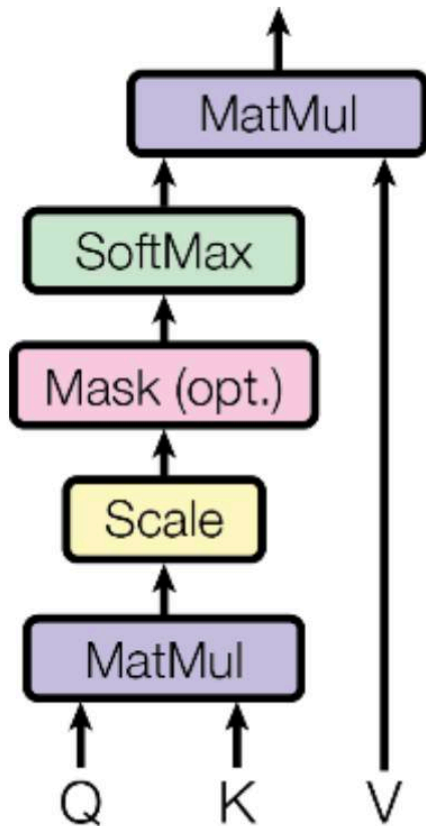


Figure 1: The Transformer - model architecture.

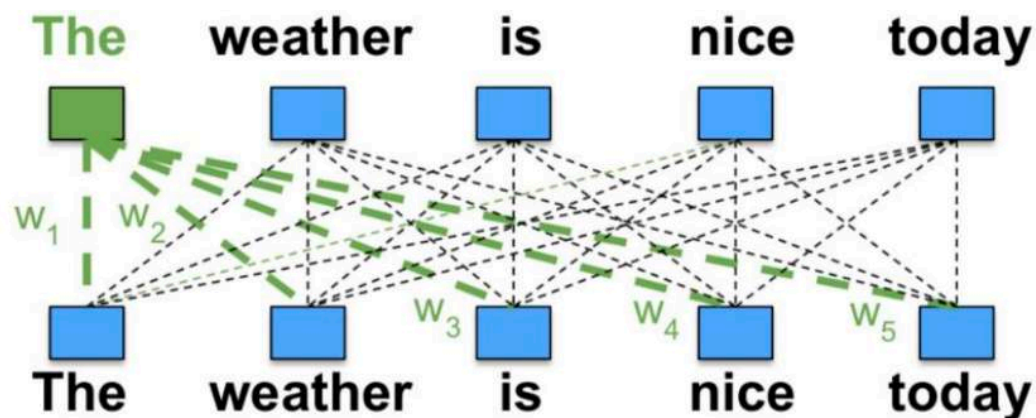
# Scaled Dot-Product Attention



$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

For large values of  $d_k$ , the dot products grow large in magnitude, pushing the softmax function into regions where it has extremely small gradients. To counteract this effect, the dot products are scaled by  $\frac{1}{\sqrt{d_k}}$ .

# Example



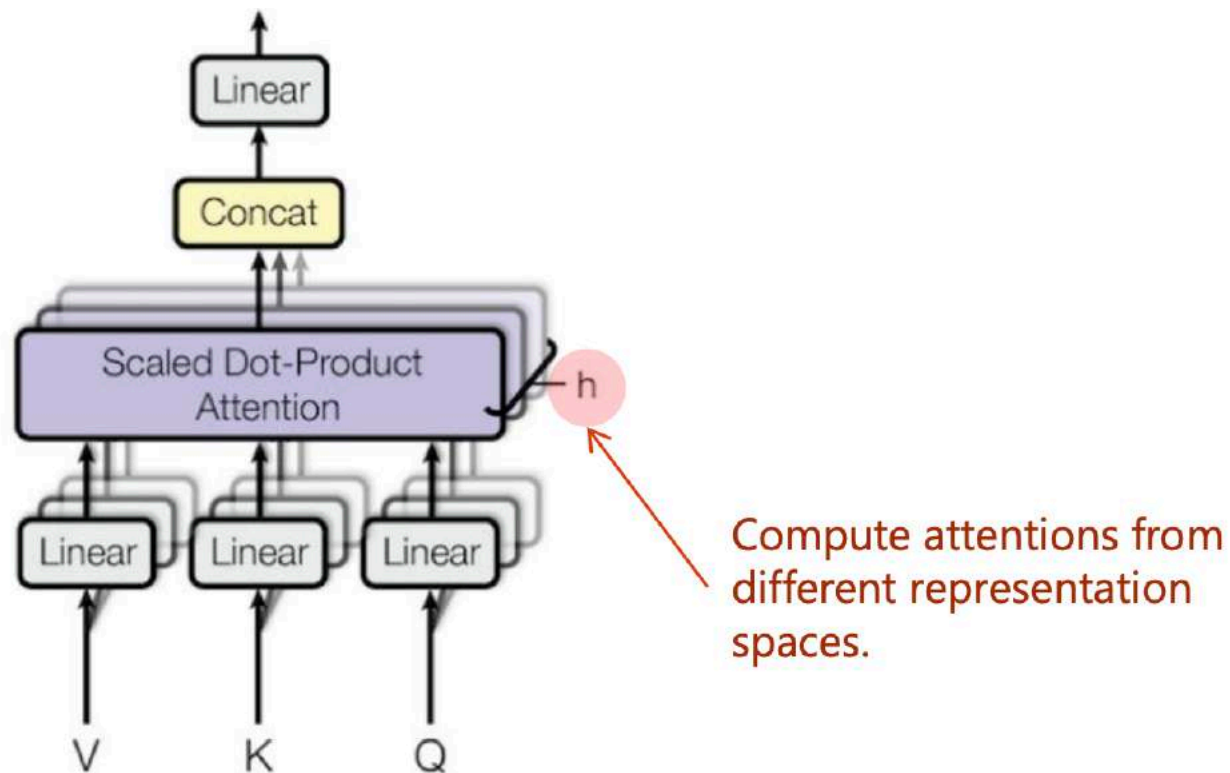
$$w_1, w_2, w_3, w_4, w_5 = \text{softmax} \left( \begin{bmatrix} 0.6 & 0.2 & 0.8 \end{bmatrix} \times \begin{bmatrix} 0.6 & 0.2 & 0.9 & 0.4 & 0.4 \\ 0.2 & 0.3 & 0.1 & 0.1 & 0.1 \\ 0.8 & 0.1 & 0.8 & 0.4 & 0.6 \end{bmatrix} \right)$$

The      The    weather    is    nice    today

$$\begin{bmatrix} 1.8 \\ 2.3 \\ 0.4 \end{bmatrix} = w_1 \times \begin{bmatrix} 0.6 \\ 0.2 \\ 0.8 \end{bmatrix} + w_2 \times \begin{bmatrix} 0.2 \\ 0.3 \\ 0.1 \end{bmatrix} + w_3 \times \begin{bmatrix} 0.9 \\ 0.1 \\ 0.8 \end{bmatrix} + w_4 \times \begin{bmatrix} 0.4 \\ 0.1 \\ 0.4 \end{bmatrix} + w_5 \times \begin{bmatrix} 0.4 \\ 0.1 \\ 0.6 \end{bmatrix}$$

The      The    weather    is    nice    today

# Multi-Head Attention



$$head_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{MultiHead}(Q, K, V) = \text{Concat}(head_1, \dots, head_h)W^O$$

# Transformer

**Layer Normalization** adjusts parameters in each layer to a more stable distribution and accelerates the training of the network.

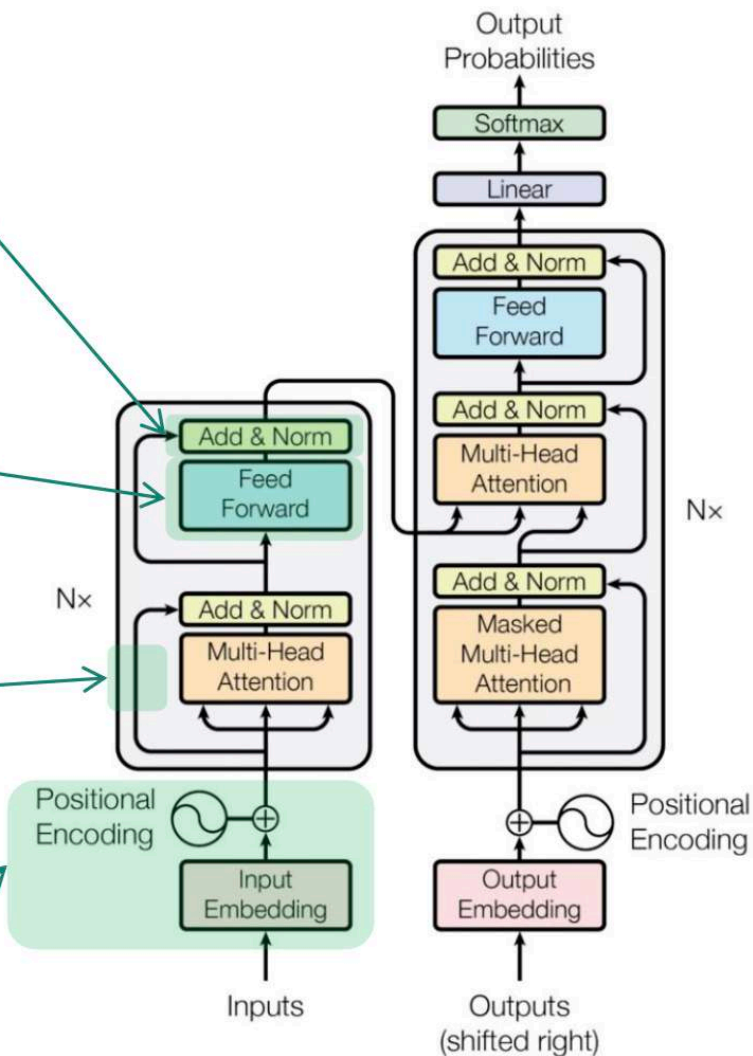
**Position-wise Feed-Forward Networks** is applied to each position:

$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

While the linear transformations are the same across different positions, they use different parameters from layer to layer.

**Residual connections** allow to have deeper models.

**Position encodings** make use of the order of the sequence and they have the same dimension as input embeddings. The two can be summed to tell the neural net that there's an underlying sequential structure.





# **The Importance of Generation Order in Language Modeling**

Nicolas Ford, Daniel Duckworth, Mohammad Norouzi, George E. Dahl

Google Brain

EMNLP18

# Overview

- Linguistic intuition might suggest that we should first generate some abstract representation of what we want to say and then serialize it.
- The best ordering we tried generates function words first and content words last, which cuts against the idea of committing to the general topic of a sentence first and only then deciding exactly how to phrase it.

# Two-pass Language Models

- Produces partially-filled sentence “templates” and then fills in missing tokens
- Partitioning of the vocabulary into a set of first-pass and second-pass tokens to generate sentences.

sentence	common first
" all you need to do if you want the na- tion 's press camped on your doorstep is to say you once had a [UNK] in 1947 , " he noted memorably in his diary . [EOS]	" all you __ to __ if you __ the __ 's __ on __ __ is to __ you __ had a [UNK] in __ , " he __ __ in his __ . [EOS]



# Two-pass Language Models

- Two copies of the Transformer model
  - **Neural language model  $p_1$**  : The first copy just generates the template, so it has no encoder.
  - **Conditional translation model  $p_2$**  : The second copy is a sequence-to-sequence model that translates the template into the complete sentence.

$$p(\mathbf{y}) = p_1(\mathbf{y}^{(1)}) p_2(\mathbf{y}^{(2)} \mid \mathbf{y}^{(1)}) .$$

Sentence  $\rightarrow$  **template**    template  $\rightarrow$  **final**

no encoder

Seq2Seq

# Two-pass Language Models

template

sentence	common first	rare first	function first	content first	odd first
" all you need to do if you want the nation 's press camped on your doorstep is to say you once had a [UNK] in 1947 , " he noted memorably in his diary . [EOS]	" all you __ to __ if you __ the __ 's __ on __ is to __ you __ had a [UNK] in __ , " he __ in his __ . [EOS]	__ __ __ need __ do __ want __ nation __ press camped __ your doorstep __ say __ once __ 1947 __ noted memorably __ diary __ [EOS]	" all you __ to __ if you __ the __ 's __ on your __ is to __ you __ a __ in __ , " he __ in his __ . [EOS]	__ __ __ need __ do __ want __ nation __ press camped __ doorstep __ say __ once had __ [UNK] __ 1947 __ noted memorably __ diary __ [EOS]	" all you need __ __ you __ the nation 's press camped on your doorstep __ say you once had __ __ " __ noted __ his __ . [EOS]
the team announced thursday that the 6-foot-1 , [UNK] starter will remain in detroit through the 2013 season . [EOS]	the __ __ that the __ , [UNK] __ will __ in __ the __ . [EOS]	__ team announced thursday __ 6-foot-1 __ starter __ remain __ detroit through __ 2013 season __ [EOS]	the __ __ that the __ , __ will __ in __ through the __ . [EOS]	__ team announced thursday __ 6-foot-1 __ [UNK] starter __ remain __ detroit __ 2013 season __ [EOS]	the team announced __ the 6-foot-1 __ will remain __ through the 2013 __ . [EOS]
scotland 's next game is a friendly against the czech republic at hampden on 3 march . [EOS]	__ 's __ is a __ the __ at __ on __ . [EOS]	scotland __ next game __ friendly against __ czech republic __ hampden __ 3 march __ [EOS]	__ 's __ is a __ against the __ at __ on __ . [EOS]	scotland __ next game __ friendly __ czech republic __ hampden __ 3 march __ [EOS]	__ 's next game __ the czech republic at hampden on 3 march . [EOS]
of course , millions of additional homeowners did make a big mistake : they took advantage of " liar loans " and other [UNK] deals to buy homes they couldn 't afford . [EOS]	of __ , __ of __ a __ : they __ of " __ " and __ [UNK] __ to __ they __ 't __ . [EOS]	__ course __ millions __ additional homeowners did make __ big mistake __ took advantage __ liar loans __ other __ deals __ buy homes __ couldn __ afford __ [EOS]	of __ , __ of __ a __ : they __ of " __ " and __ to __ they __ . [EOS]	__ course __ millions __ additional homeowners did make __ big mistake __ took advantage __ liar loans __ other [UNK] deals __ buy homes __ couldn 't afford __ [EOS]	of __ __ of additional __ big __ they __ advantage of " liar __ " and other __ deals __ buy homes they couldn __ afford . [EOS]

Table 1: Some example sentences from the dataset and their corresponding templates. The placeholder token is indicated by “\_\_”.

# Results

- It is easier to first decide something about its syntactic structure.
- It is preferable to delay committing to a rare token for as long as possible as all subsequent decisions will then be conditioning on a low-probability event.

Model	Train	Validation	Test
odd first	39.925	45.377	45.196
rare first	38.283	43.293	43.077
content first	38.321	42.564	42.394
common first	36.525	41.018	40.895
function first	36.126	40.246	40.085
baseline	38.668	41.888	41.721
enhanced baseline	35.945	39.845	39.726

Table 2: The perplexities achieved by the best version of each of our models.

# **Insertion Transformer: Flexible Sequence Generation via Insertion Operations**

Mitchell Stern, William Chan, Jamie Kiros, Jakob Uszkoreit

Google Brain, University of California, Berkeley

ICML19

# Insertion Transformer

- $x$  : source canvas (sequence)
- $y$  : target canvas (sequence)
- $\hat{y}_t$  : hypothesis canvas at time  $t$
- $\mathcal{C}$  : content vocabulary (token vocabulary for sequences)
- $l$  : locations  $\in [0, |\hat{y}_t|]$

## Insertion Transformer: Flexible Sequence Generation via Insertion Operations

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

$$p(c, l \mid x, \hat{y}_t) = \text{InsertionTransformer}(x, \hat{y}_t).$$



# Insertion Transformer Model

- **Full Decoder Self-Attention**
  - Remove causal self attention
- **Slot Representations via Concatenated Outputs**
  - Adding special marker tokens at the **beginning** and **end** of the decoder input to extend the sequence length by two.
  - Take the resulting  **$n + 2$**  vectors in the final layer and **concatenate each adjacent pair** to obtain  **$n + 1$**  slot representations.

# Model

$$p(c, l \mid x, \hat{y}_t) = \text{InsertionTransformer}(x, \hat{y}_t).$$

- **Joint content-location distribution**

$$p(c, l) = \text{softmax}(\text{flatten}(HW))$$

*flatten this matrix into a vector*

$H \in \mathbb{R}^{(T+1) \times h}$  matrix of slot representations  
 $W \in \mathbb{R}^{h \times |\mathcal{C}|}$

- **Joint distribution using a conditional factorization**

$$p(c, l) = p(l)p(c|l) = \text{softmax}(Hq) \times \text{softmax}(h_l W)$$

learnable query vector       $l$ -th row of H

# Contextualized Vocabulary Bias

context vector  $g = \text{maxpool}(H) \quad g \in \mathbb{R}^h$

shared bias  $b = gV \quad V \in \mathbb{R}^{h \times |C|}$

$$B = \text{repmat}(b, [T + 1, 1])$$

$$p(c, l) = \text{softmax}(HW + B)$$

Global bias



# Training and Loss Functions

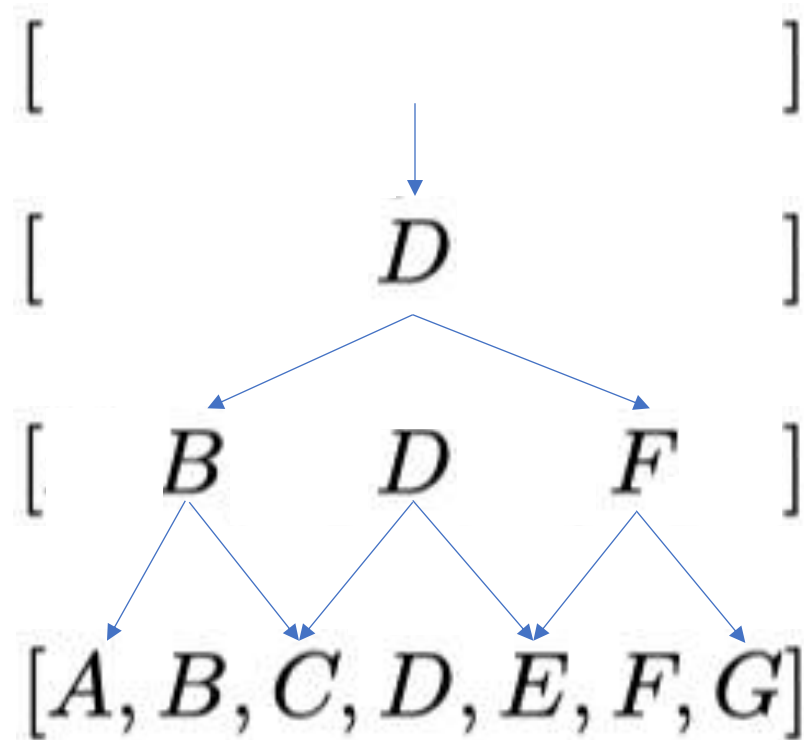
- **Left-to-Right**

- Example :  $(x, y)$
- Sample a length  $k \sim \text{uniform}([0, |y|])$
- Create a new data point  $((x, \hat{y} = (y_1, \dots, y_k)), y_{k+1})$
- Loss : classification loss (negative log-likelihood)
- *Note : only concerns about the last position to insert*

$$\text{loss}(x, \hat{y}) = -\log p(y_{k+1}, k \mid x, \hat{y})$$

# Balanced Binary Tree

- Parallelism



# Balanced Binary Tree

- Example :  $(x, y)$
- Sample a length  $k \sim \text{uniform}([0, |y|])$
- Sample a random subsequence of  $y$  of length  $k$ :  $\hat{y}$ 
  1. Shuffle  $y$
  2. Extract the first  $k$
  3. Reorder

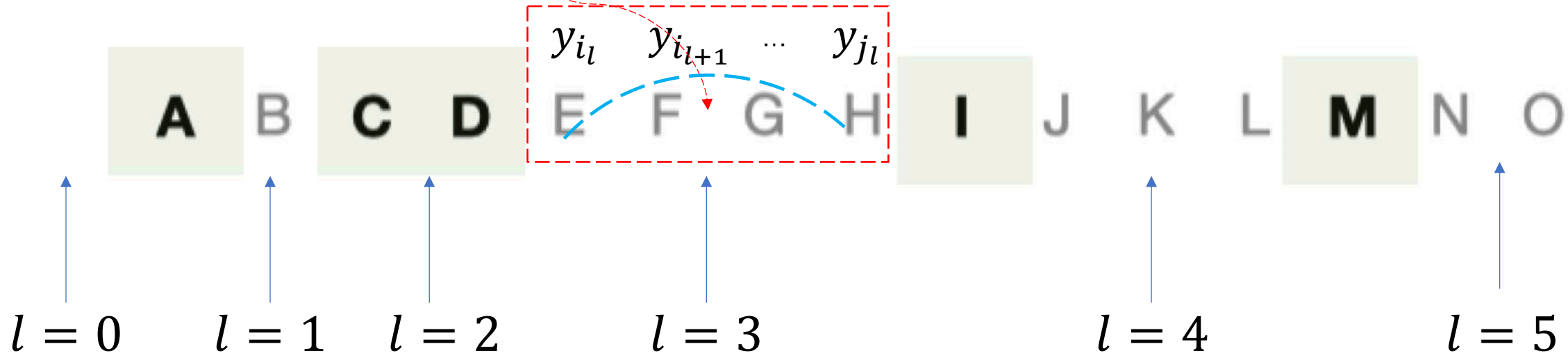
**A** B **C D** E F G H **I** J K L **M** N O

# Soft binary tree loss

$$d_l(i) = \left| \frac{i_l + j_l}{2} - i \right| \quad w_l(i) = \frac{\exp(-d_l(i)/\tau)}{\sum_{i'=i_l}^{j_l} \exp(-d_l(i')/\tau)}.$$

*Distance* ↓  
*w<sub>l</sub>(i)* ↑

span of tokens from the target  
output yet to be produced



$$\text{slot-loss}(x, \hat{y}, l) = \sum_{i=i_l}^{j_l} -\log p(y_i, l \mid x, \hat{y}) \cdot w_l(i).$$

$$\text{loss}(x, \hat{y}) = \frac{1}{k+1} \sum_{l=0}^k \text{slot-loss}(x, \hat{y}, l).$$

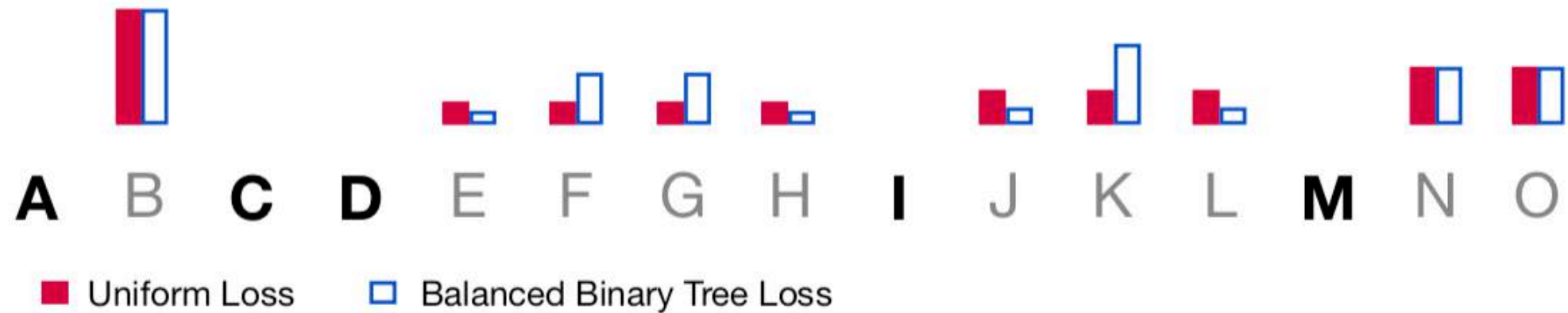
# Uniform

$$w_l(i) = \frac{\exp(-d_l(i)/\tau)}{\sum_{i'=i_l}^{j_l} \exp(-d_l(i')/\tau)}. \quad \tau \rightarrow \infty$$

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



# Balanced binary tree and uniform losses



*Figure 2.* A visualization of the weighting of the per-token negative log-likelihoods in the balanced binary tree and uniform losses. The balanced binary tree loss strongly incentivizes the generation of the center word or center words within each slot.

# Greedy Decoding

- Choose the action with the highest probability

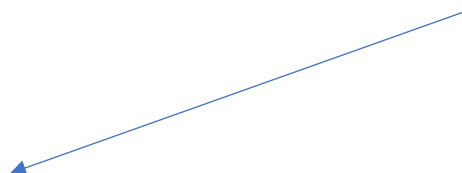
$$(\hat{c}_t, \hat{l}_t) = \operatorname{argmax}_{c, l} p(c, l \mid x, \hat{y}_t).$$

- **sequence finalization**
  - until an end-of-sequence token gets selected
- **slot finalization**
  - restrict the argmax to locations whose maximum-probability decision is not end-of-slot
  - Until the model predicts an end-of-slot token for every location.

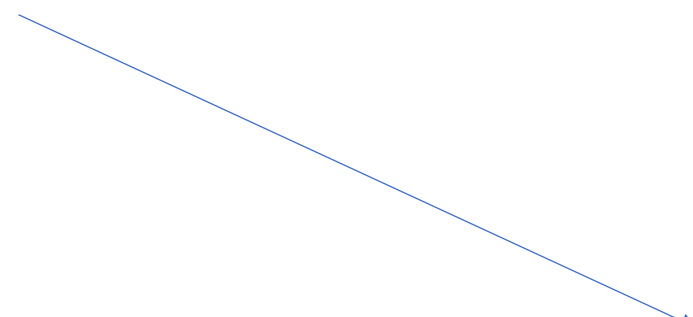
# Parallel Decoding

- For each location  $l$

$$p(c \mid l) : \hat{c}_{l,t} = \operatorname{argmax}_c p(c \mid l, x, \hat{y}_t).$$


$$p(c \mid l) = p(c, l) / p(l) = p(c, l) / \sum_{c'} p(c', l)$$

joint distribution


$$p(c, l) = p(l)p(c|l) = \operatorname{softmax}(Hq) \times \operatorname{softmax}(h_l W)$$

$$p(c, l) = \operatorname{softmax}(\operatorname{flatten}(HW)).$$

factorization

- slot finalization

# Non-Monotonic Sequential Text Generation

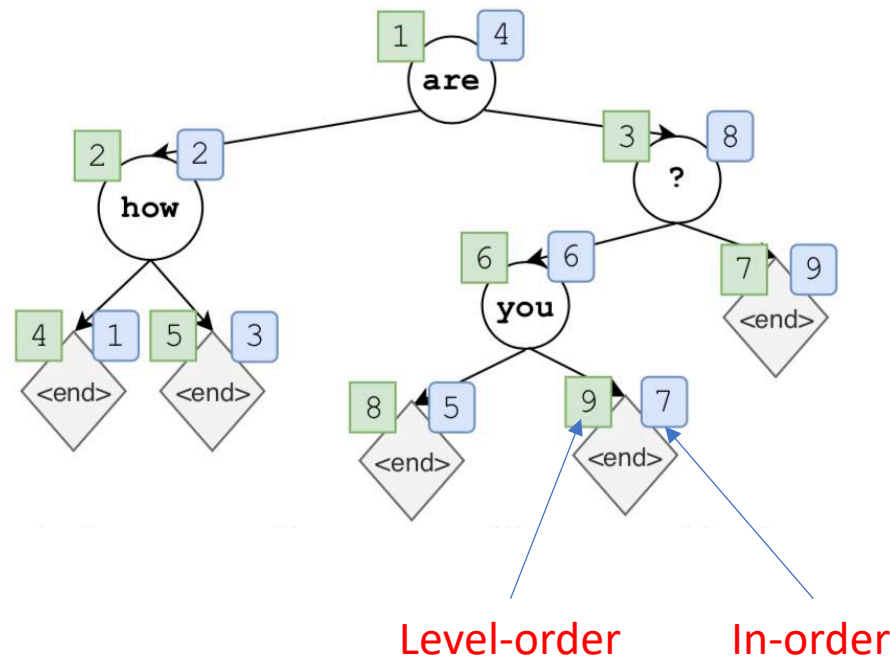
Sean Welleck , Kianté Brantley, Hal Daumé III, Kyunghyun Cho

New York University, University of Maryland, College Park  
Microsoft Research, Facebook AI Research  
CIFAR Azrieli Global Scholar

ICML19

# Overview

- Recursively generating words to its left and then words to its right, yielding a binary tree.
- Learning is framed as imitation learning, including a coaching method which moves from imitating an oracle to reinforcing the policy's own preferences



# Imitation Learning

- Imitation Learning with Recurrent Neural Networks
- Learning to Search Better than Your Teacher *ICML15*
- <https://zhuanlan.zhihu.com/p/25688750>
- <https://blog.csdn.net/WASEFADG/article/details/83651126>
- <https://www.quora.com/What-is-imitation-learning>

# Notation

- Vocabulary  $\tilde{V} = V \cup \{< end >\}$
- State space  $\tilde{V}^*$
- State  $s \in S$  corresponds to a sequence of tokens from  $\tilde{V}$
- Init state: empty sequence  $<>$
- End state:  $< end >$
- Action  $a$  : select an element from vocab and append to the state
- $\tau(t)$ : maps from in-order to level order
- Policy  $\pi(a|s)$

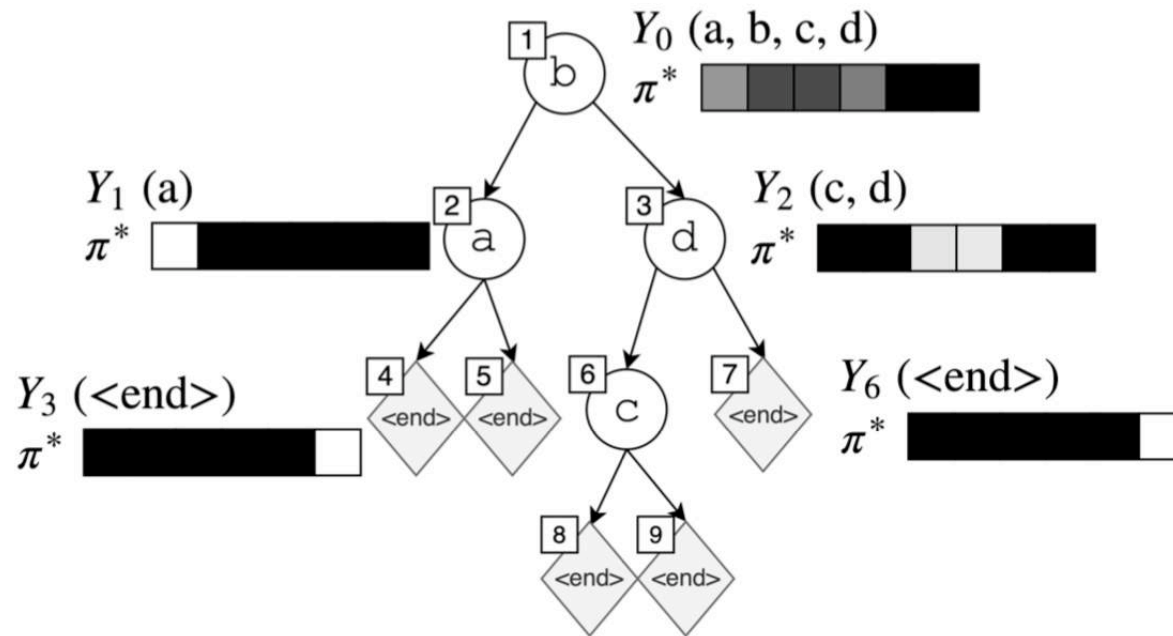
# Challenge

- The sequences  $Y$  alone only tell us what the final output sequences of words should be, but not what tree(s) should be used to get there.

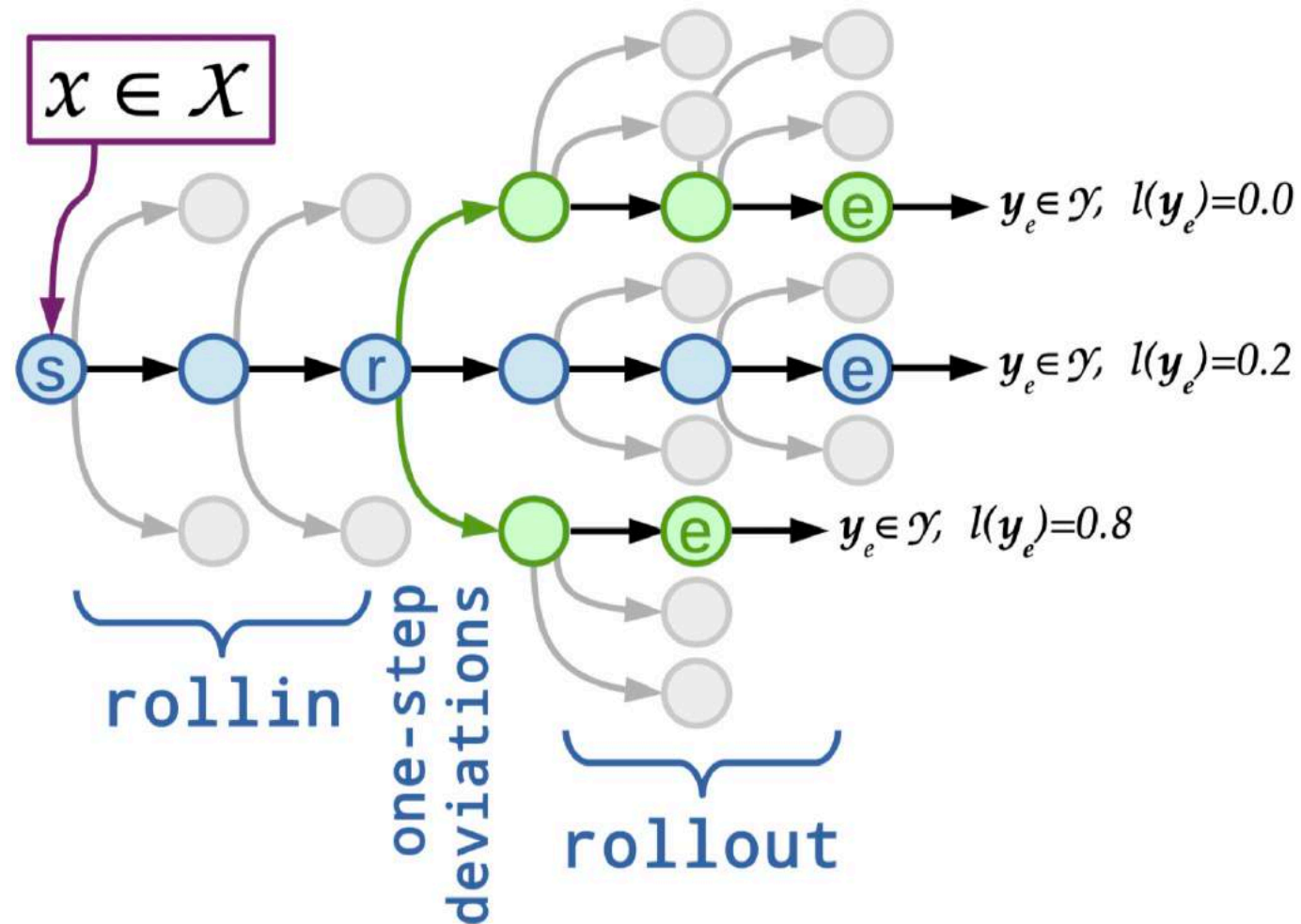


# Imitation Learning

- The first step, an oracle policy's action is to produce any word  $w$  that appears anywhere in  $Y$ .
- All words to the left of  $w$  in  $Y$  are generated recursively on the left (following the same procedure), and all words to the right of  $w$  in  $Y$  are generated recursively on the right.
- The oracle is non-deterministic (many “correct” actions are available at any given time), we inform this oracle policy with the current learned policy, encouraging it to favor actions that are preferred by the current policy.



# Background: Learning to Search



# Loss

- 3  $\mathbb{E}$ 
  - draw states  $s$  according to the state distribution induced by  $\pi^{in}$
  - compute cost-to-go under  $\pi^{out}$ , for all possible actions  $a$  at that state.
- 2  $\mathbb{E}$ 
  - running  $\pi$  for t-many steps
- 1  $\mathbb{E}$ 
  - for one instance

$$\underbrace{\mathbb{E}_{Y \sim D}}_{\text{3}} \underbrace{\mathbb{E}_{t \sim U[2|Y|+1]}}_{\text{2}} \underbrace{\mathbb{E}_{s_t \sim d_{\pi^{in}}^t}}_{\text{1}} \left[ \mathcal{C}(\pi; \pi^{out}, s_t) \right]$$

# Cost Measurement

- when dealing with recurrent neural network policies using a cost function more analogous to a cross-entropy loss can be preferred
- use a KL-divergence type loss, measuring the difference between the action distribution produced by  $\pi$  and the action distribution preferred by  $\pi^{out}$ .
- first sampling one training sequence, running the roll-in policy for  $t$  steps, and computing the KL divergence at that state using  $\pi^*$  (reference *or* oracle )as  $\pi^{out}$ . Learning corresponds to minimizing this KL divergence iteratively with respect to the parameters of  $\pi$ .

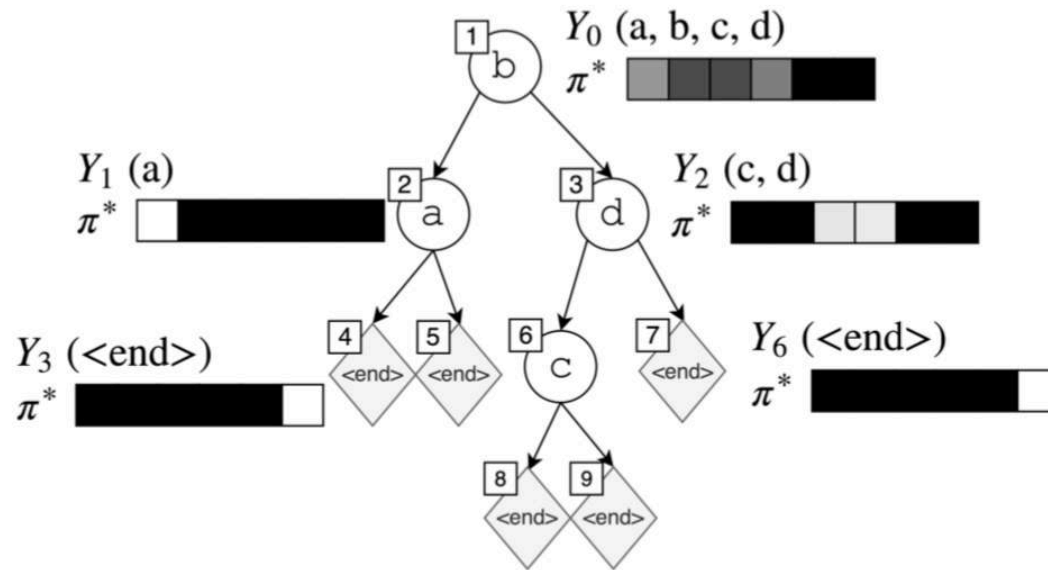
$$\begin{aligned}\mathcal{C}(\pi; \pi^{out}, s) &= D_{KL}(\pi^{out}(\cdot|s) || \pi(\cdot|s)) \\ &= \sum_{a \in \tilde{V}} \pi^{out}(a|s) \log \pi(a|s) + const.\end{aligned}$$

# Roll-In Policies

- In most formal analyses, the roll-in policy is a stochastic mixture of the learned policy  $\pi$  and the oracle policy  $\pi^*$
- *Experimentally*, it has often been found that simply **using the oracle's state distribution** is optimal

roll-out →	Reference	Mixture	Learned
↓ roll-in			
Reference	Inconsistent		
Learned	Not locally opt.	Good	RL

# Oracle Policies



$$\pi^*(a|s_t) = \begin{cases} 1 & \text{if } a = \langle \text{end} \rangle \text{ and } Y_t = \langle \rangle \\ p_a & \text{if } a \in Y_t \\ 0 & \text{otherwise} \end{cases}$$

- Uniform Oracle.  $p_a = 1/n$
- Coaching Oracle
  - preferring actions that are preferred by the current parameterized policy

$$\pi_{\text{coaching}}^*(a|s) \propto \pi_{\text{uniform}}^*(a|s) \pi(a|s)$$

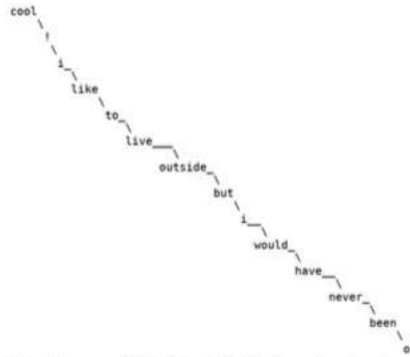
- Annealed Coaching Oracle( $\beta$  from 1 to 0)

$$\pi_{\text{annealed}}^*(a|s) = \beta \pi_{\text{uniform}}^*(a|s) + (1 - \beta) \pi_{\text{coaching}}^*(a|s)$$

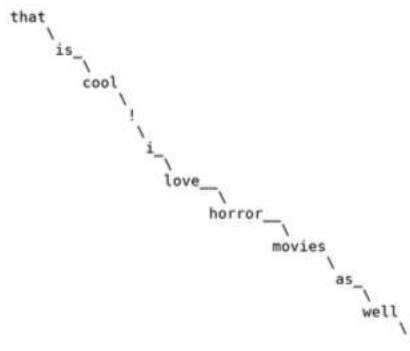
# Word Reordering Examples

Figure 8. Word Reordering Examples. The columns show policies trained with  $\pi_{\text{left-right}}^*$ ,  $\pi_{\text{uniform}}^*$ , and  $\pi_{\text{annealed}}^*$ , respectively.

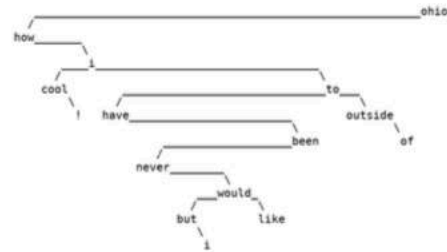
Actual: how cool ! i have never been outside of ohio but i would like to .  
 Predicted: cool ! i like to live outside but i would have never been of ohio .  
 Gen. Order: cool ! i like to live outside but i would have never been of ohio .



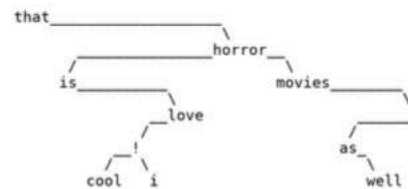
Actual: that is cool ! i love horror movies as well !  
 Predicted: that is cool ! i love horror movies as well !  
 Gen. Order: that is cool ! i love horror movies as well !



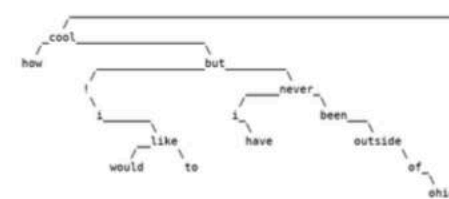
Actual: how cool ! i have never been outside of ohio but i would like to .  
 Predicted: how cool ! i have never but i would like been to outside of ohio .  
 Gen. Order: ohio how . i cool to ! have outside been of never would but like i



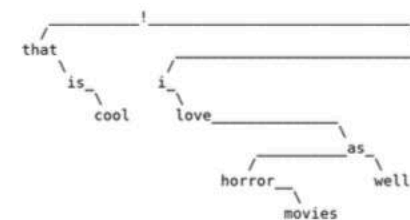
Actual: that is cool ! i love horror movies as well !  
 Predicted: that is cool ! i love horror movies as well !  
 Gen. Order: that horror is movies love ! ! as cool i well



Actual: how cool ! i have never been outside of ohio but i would like to .  
 Predicted: how cool ! i would like to but i have never been outside of ohio .  
 Gen. Order: . cool how but ! never i i been like have outside would to of ohio



Actual: that is cool ! i love horror movies as well !  
 Predicted: that is cool ! i love horror movies as well !  
 Gen. Order: ! that ! is i cool love as horror well movies



# **Insertion-based Decoding with automatically Inferred Generation Order**

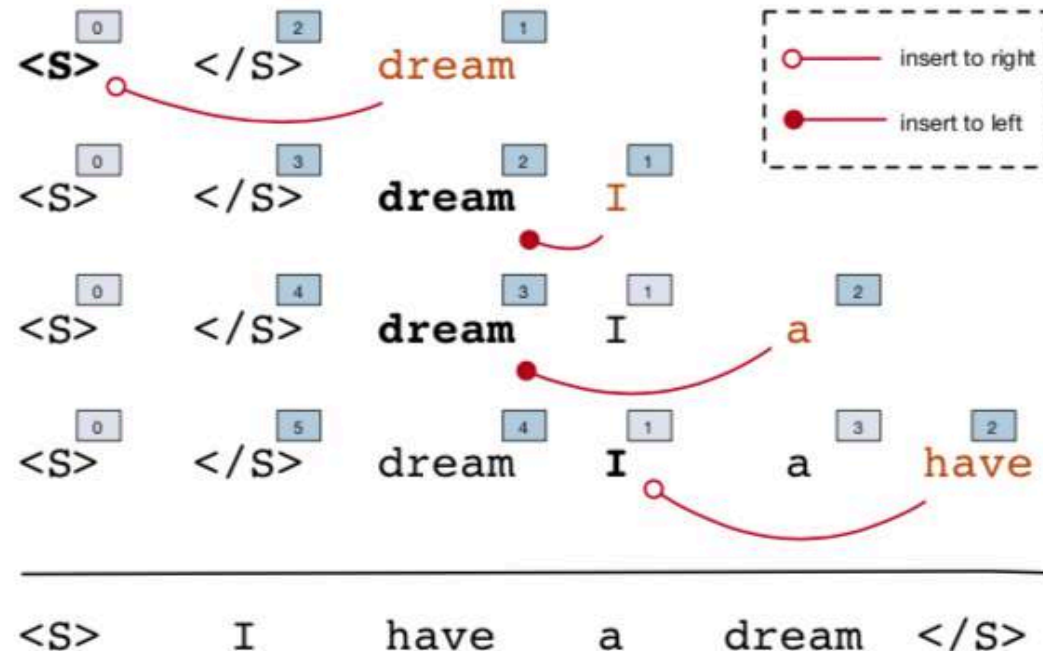
JiataoGu, QiLiu, KyunghyunCho

Facebook AI Research  
New York University



# Motivation

- L2R is not necessarily the optimal option for generating sequences.
- For instance, people sometimes tend to think of central phrases first before building up a whole sentence.



# Orders as Latent Variables

- $P_T$  is the set of all the permutations of  $(1, \dots, T)$
- $\pi = (z_2, z_3, \dots, z_T, z_{T+1}) \in P_T$
- $y_\pi = \{(y_2, z_2), \dots, (y_{T+1}, z_{T+1})\}$ ,  $(y_t, z_t)$  represents the  $t$  –  $th$  generated token and its absolute position
- Two special tokens
  - $(y_0, z_0) = (< s >, 0)$ 、 $(y_1, z_1) = (</s >, T + 1)$
- Object

$$p_\theta(\mathbf{y}|\mathbf{x}) = \sum_{\pi \in \mathcal{P}_T} p_\theta(\mathbf{y}_\pi|\mathbf{x})$$

$$p_\theta(\mathbf{y}_\pi|\mathbf{x}) = p_\theta(\underline{y_{T+2}}|y_{0:T+1}, z_{0:T+1}, x_{1:T'}) \cdot y_{T+2} = < eod >$$

$$\prod_{t=1}^T p_\theta(\underline{y_{t+1}}, \underline{z_{t+1}}|y_{0:t}, z_{0:t}, x_{1:T'})$$

# Relative Representation of Positions

- $r_i^t$ : the relative-position representations of token  $i$  at decode step  $t$
- $r_i^t$  is a vector
- Value : 0, 1, -1  $r_{i,j}^t = \begin{cases} -1 & z_j^t > z_i^t \text{ (left)} \\ 0 & z_j^t = z_i^t \text{ (middle)} \\ 1 & z_j^t < z_i^t \text{ (right)} \end{cases}$
- Matrix  $R^t = [r_0^t, r_1^t, \dots, r_t^t]$  shows the relative-position representations of all the words in the sequence.
- Mapped back to the absolute position  $z_i^t = \sum_{j=0}^t \max(0, r_{i,j}^t)$
- Update

$$R^{t+1} = \left[ \begin{array}{ccc|c} & & & r_{t+1,0}^{t+1} \\ & & & \vdots \\ & R^t & & r_{t+1,t}^{t+1} \\ \hline -r_{t+1,0}^{t+1} & \cdots & -r_{t+1,t}^{t+1} & 0 \end{array} \right]$$

# Insertion-based Decoding

- Given  $y_{0:t}$  and  $r_{0:t}$
- Predict  $y_{t+1}$  and  $r_{t+1}$
- Note : only concerns about the  $y_k$  which has been selected
- $s = -1$  if  $y_{t+1}$  is on the left of  $y_k$ , and  $s = 1$  otherwise.

$$r_{t+1,j} = \begin{cases} s & j = k \\ r_{k,j} & j \neq k \end{cases}, \quad \forall j \in [0, t]$$

# Insertion-based Decoding

---

**Algorithm 1** Insertion-based Decoding

---

**Initialize:**  $\mathbf{y} = (\langle s \rangle, \langle /s \rangle)$ ,  $R = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}$ ,  $t = 1$

**repeat**

    Predict the next word  $y_{t+1}$  based on  $\mathbf{y}$ ,  $R$ .

**if**  $y_{t+1}$  is  $\langle \text{eod} \rangle$  **then**

        break

**end if**

    Choose an existing word  $y_k \in \mathbf{y}$ ;

    Choose the left or right ( $s$ ) of  $y_k$  to insert;

    Obtain the next position  $\mathbf{r}_{t+1}$  with  $k, s$  (Eq. (6)).

    Update  $R$  by appending  $\mathbf{r}_{t+1}$  (Eq. (5)).

    Update  $\mathbf{y}$  by appending  $y_{t+1}$

    Update  $t = t + 1$

**until** Reach the maximum length

Map back to absolute positions  $\boldsymbol{\pi}$  (Eq. (4))

Reorder  $\mathbf{y}$ :  $y_{z_i} = y_i \quad \forall z_i \in \boldsymbol{\pi}, i \in [0, t]$

---

# Transformer-InDiGO

- Relative position-based self-attention

$$e_{i,j} = \frac{(u_i^\top Q) \cdot \left( u_j^\top K + A_{[r_{i,j}+1]} \right)^\top}{\sqrt{d_{\text{model}}}}$$

$$A \in \mathbb{R}^{3 \times d_{\text{model}}}$$

# Transformer-InDIGO

- Word & Position Prediction  $H = (h_0, \dots, h_t)$

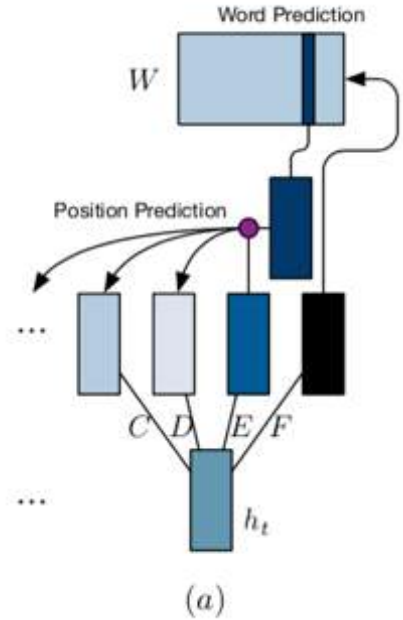
$$p(y_{t+1}, \mathbf{r}_{t+1} | H) = p(y_{t+1} | H) \cdot p(\mathbf{r}_{t+1} | y_{t+1}, H)$$

$$p_{\text{word}}(y | H) = \text{softmax} \left( (\mathbf{h}_t^\top F) \cdot W^\top \right)$$

$$p_{\text{pointer}}(k | y_{t+1}, H) =$$

$$\text{softmax} \left( (\mathbf{h}_t^\top E + W_{[y_{t+1}]}) \cdot \begin{bmatrix} H^\top C \\ H^\top D \end{bmatrix}^\top \right)$$

$$k_{t+1} \in [0, 2t + 1]$$



# Transformer-InDIGO

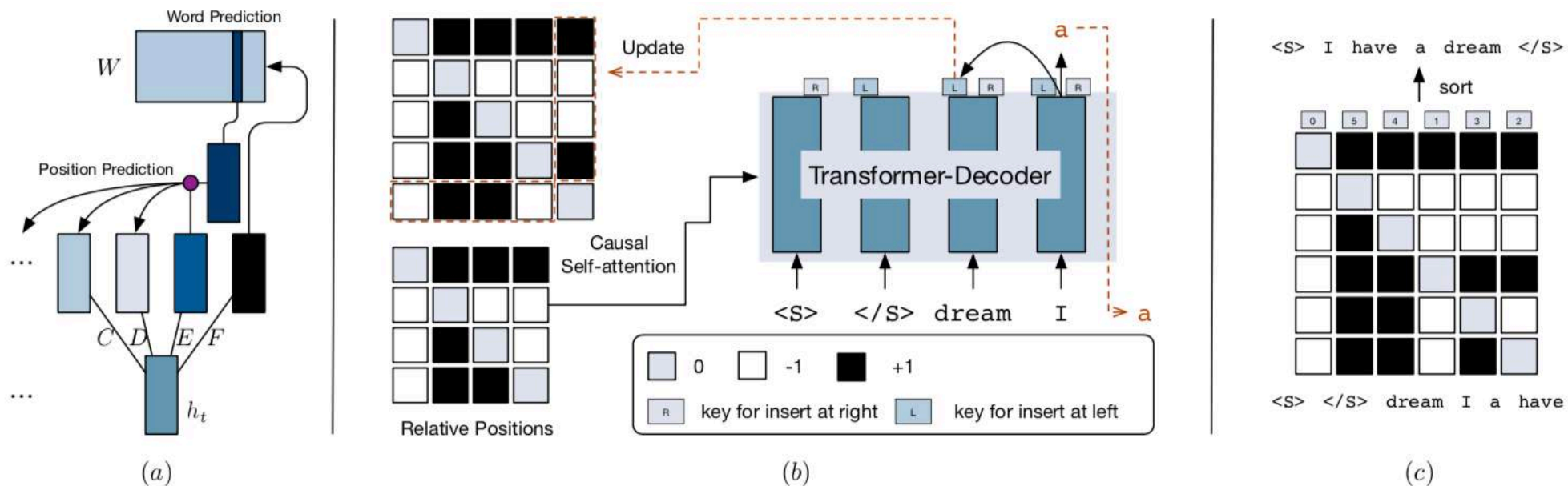


Figure 2: The overall framework of the proposed Transformer-InDIGO which includes (a) the word & position prediction module; (b) the one step decoding with position updating; (c) final decoding output by reordering.



# Learning

- This is intractable since we need to enumerate all of the  $T!$  permutations of tokens.  $p_{\theta}(\mathbf{y}|\mathbf{x}) = \sum_{\pi \in \mathcal{P}_T} p_{\theta}(\mathbf{y}_{\pi}|\mathbf{x})$
- Maximize the evidence lower-bound (ELBO) of the original objective by introducing an approximate posterior distribution of generation orders  $q(\pi|x, y)$ , which provides the probabilities of latent generation orders based on the ground-truth sequences  $x$  and  $y$ :

$$\begin{aligned}\mathcal{L}_{\text{ELBO}} &= \mathbb{E}_{\pi \sim q} \log p_{\theta}(\mathbf{y}_{\pi}|\mathbf{x}) + \mathcal{H}(q) \\ &= \mathbb{E}_{\mathbf{r}_{2:T+1} \sim q} \left( \underbrace{\sum_{t=1}^{T+1} \log p_{\theta}(y_{t+1} | y_{0:t}, \mathbf{r}_{0:t}, x_{1:T'})}_{\text{Word Prediction Loss}} \right. \\ &\quad \left. + \sum_{t=1}^T \underbrace{\log p_{\theta}(\mathbf{r}_{t+1} | y_{0:t+1}, \mathbf{r}_{0:t}, x_{1:T'})}_{\text{Position Prediction Loss}} \right) + \mathcal{H}(q)\end{aligned}$$

# Searched Adaptive Order (SAO)

- beam-search in the space of all the permutations of the target sequence
- Sub-sequence :  $y_{0:t}^{(b)} \in \mathcal{B}$
- Left words :  $y' \in \mathbf{y} \setminus y_{0:t}^{(b)}$
- corresponding position  $\mathbf{r}'$
- select *top-B* sub-sequences as the new set  $\mathcal{B}$  for the next step.

$$\mathcal{L}_{SAO} = \frac{1}{B} \sum_{\pi \in \mathcal{B}} \log p_{\theta}(\mathbf{y}_{\pi} | \mathbf{x})$$

$$q(\boldsymbol{\pi} | \mathbf{x}, \mathbf{y}) = \begin{cases} 1/B & \boldsymbol{\pi} \in \mathcal{B} \\ 0 & \text{otherwise} \end{cases}$$

# Levenshtein Transformer

Jiatao Gu, Changhan Wang, and Jake Zhao (Junbo)

Facebook AI Research  
New York University  
Tigerobo Inc

# Levenshtein Distance

**4 = Levenshtein Distance(Saturday, Sturday)**

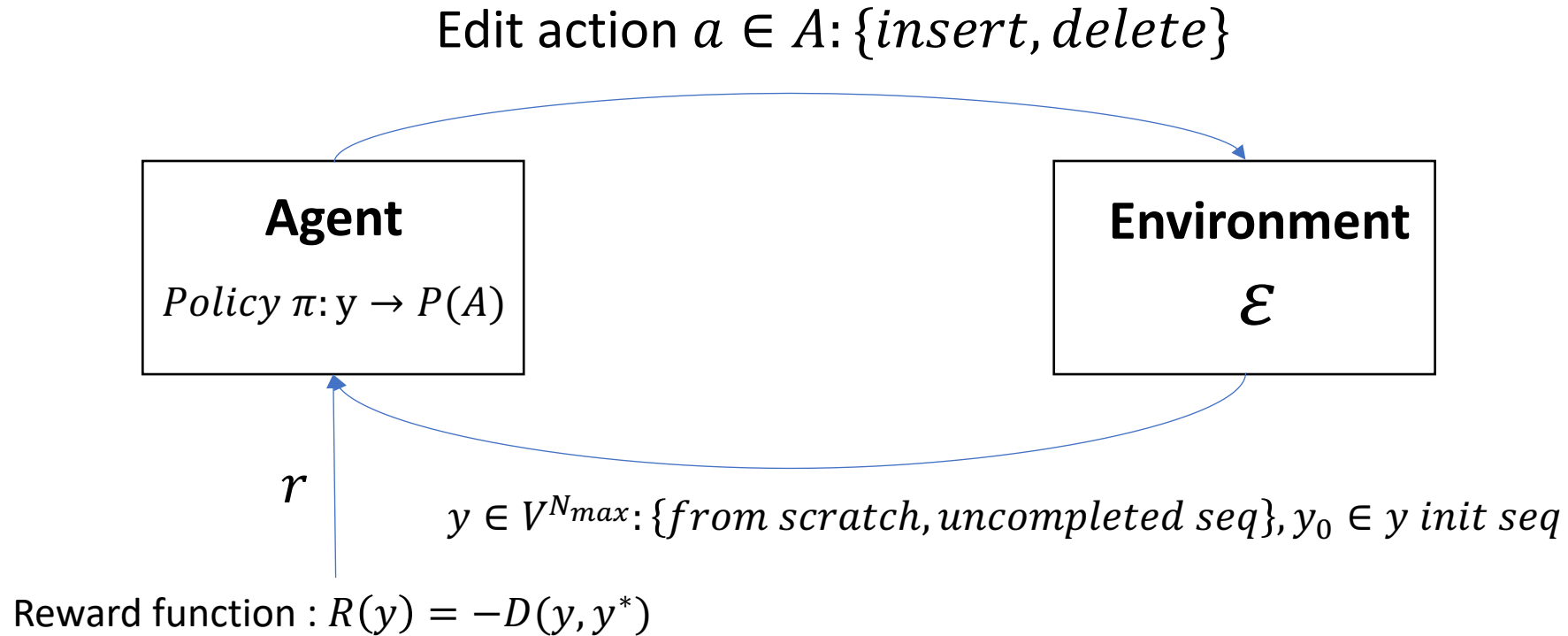
1. Sa**t**urday → Sturday // delete the first a
2. St**u**rday → Surday // delete the first t
3. Sur**r**day → Sunday // replace r with n
4. Sunday → Sundays // add s at the end

# Overview

- Humans can revise, replace, revoke or delete any part of their generated text.
- Atomic operations : insertion and deletion
- Not only generation but also sequence refinement allowing dynamic length changes.
- Partially autoregressive model

# Problem Formulation

- Markov Decision Process (MDP)  $(\mathcal{Y}, \mathcal{A}, \mathcal{E}, \mathcal{R}, y_0)$



# Actions

$$\mathbf{y}^{k+1} = \mathcal{E}(\mathbf{y}^k, \mathbf{a}^{k+1})$$

$$\mathbf{y}^k = (y_1, y_2, \dots, y_n)$$

$\langle s \rangle$

$\langle /s \rangle$

$a \in A: \{insert, delete\}$

# Deletion

- $\pi^{\text{del}}(d|i, \mathbf{y})$  makes a binary decision which is 1 (delete this token) or 0 (keep it)
- Avoid sequence boundary being broken  $\pi^{\text{del}}(0|1, \mathbf{y}) = \pi^{\text{del}}(0|n, \mathbf{y})$



# Insertion

- placeholder prediction and token prediction
- All locations  $(y_i, y_{i+1})$  in  $\mathbf{y}$
- $\pi^{\text{plh}}(p|i, \mathbf{y})$  the possibility of adding one or several placeholders
- $\pi^{\text{tok}}(t|i, \mathbf{y})$  for every placeholder predicted as above, replaces the placeholders with actual tokens in the vocabulary

# Policy combination

- **delete tokens – insert placeholders - replace placeholders** with new tokens
- parallelize the computation within each sub-tasks.

$$\mathbf{a} = \underbrace{\{d_0, \dots, d_n\}}_{\mathbf{d}}; \underbrace{\{p_0, \dots, p_{n-1}\}}_{\mathbf{p}}; \underbrace{\{t_0^1, \dots, t_0^{p_0}, \dots, t_{n-1}^{p_{n-1}}\}}_{\mathbf{t}}$$

$$\pi(\mathbf{a}|\mathbf{y}) = \prod_{d_i \in \mathbf{d}} \pi^{\text{del}}(d_i|i, \mathbf{y}) \cdot \prod_{p_i \in \mathbf{p}} \pi^{\text{plh}}(p_i|i, \mathbf{y}') \cdot \prod_{t_i \in \mathbf{t}} \pi^{\text{tok}}(t_i|i, \mathbf{y}'')$$

# Levenshtein Transformer

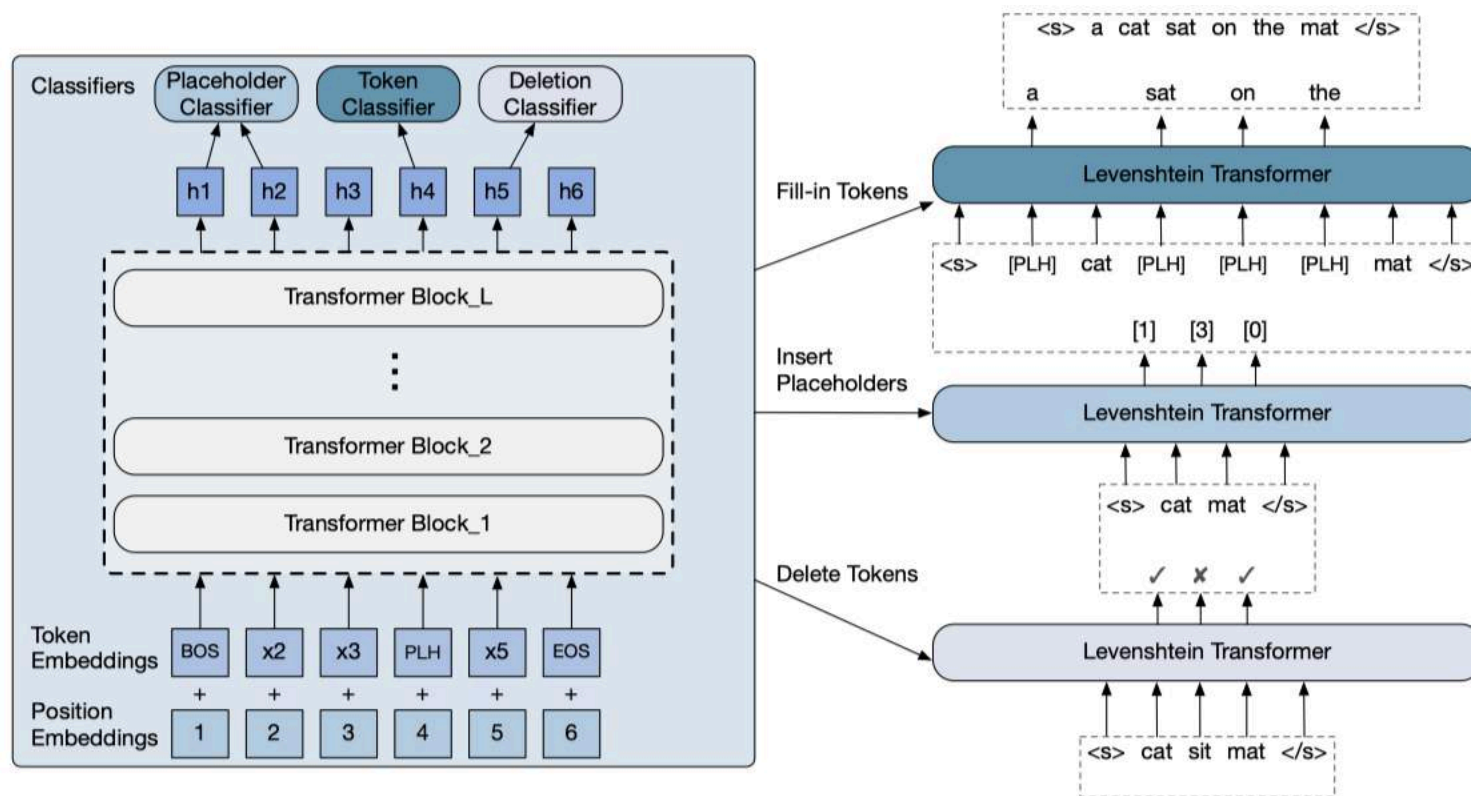


Figure 1: The overall framework of **the decoder of the proposed Levenshtein Transformer**. We show how the same architecture can be applied for three different tasks with specific classifiers. For simplicity, the attention between the encoder outputs is omitted within each Transformer-Block.

# Levenshtein Transformer

- Decoder output :  $(h_0, h_1, \dots, h_n)$ , passed to three policy classifiers

1. **Deletion Classifier:** scans over the input tokens (except for the boundaries) and predict “deleted” (0) or “kept” (1) for each token position

$$\pi_{\theta}^{\text{del}}(d|i, \mathbf{y}) = \text{softmax}(\mathbf{h}_i \cdot A^{\top}), \quad i = 1, \dots, n - 1$$

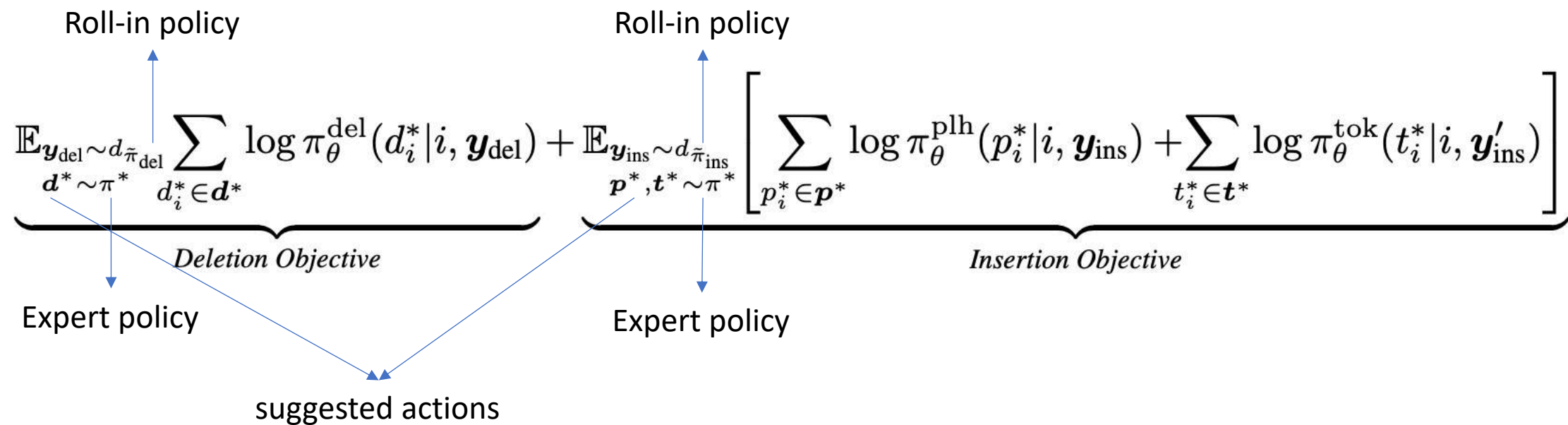
2. **Placeholder Classifier:** predicts the number of tokens to be inserted at every consecutive position pairs

$$\pi_{\theta}^{\text{plh}}(p|i, \mathbf{y}) = \text{softmax}(\text{concat}(\mathbf{h}_i, \mathbf{h}_{i+1}) \cdot B^{\top}), \quad i = 0, \dots, n - 1 \quad B \in \mathbb{R}^{(K_{\max}+1) \times (2d_{\text{model}})}$$

3. **Token Classifier:** fill in tokens replacing all the placeholders.

$$\pi_{\theta}^{\text{tok}}(t|i, \mathbf{y}) = \text{softmax}(\mathbf{h}_i \cdot C^{\top}), \quad \forall y_i = \langle \text{PLH} \rangle$$

# Dual-policy Learning



# Roll-in Policy

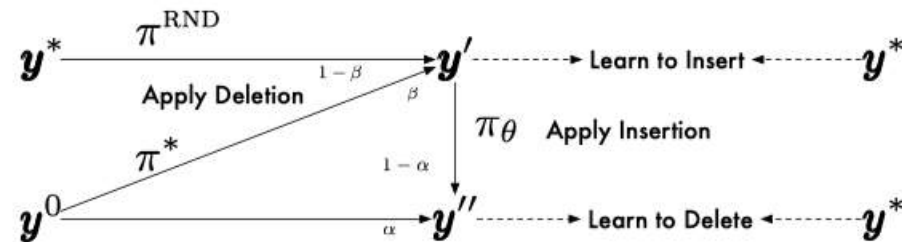


Figure 2: The data-flow of learning.

## • Learning to Delete

$$d_{\tilde{\pi}_{\text{del}}} = \{ \underset{\text{initial input}}{y^0} \text{ if } \underset{u \sim \text{uniform}[0,1]}{u} < \underset{\text{mixture factor}}{\alpha} \text{ else } \underbrace{\mathcal{E}(\mathcal{E}(y', p^*), \tilde{t})}_{\text{output by applying insertion}}, \quad \underset{\text{any sequence ready to insert tokens}}{p^* \sim \pi^*, \tilde{t} \sim \pi_\theta} \}$$

obtained by sampling instead of doing argmax

## • Learning to Insert

$$d_{\tilde{\pi}_{\text{ins}}} = \{ \underbrace{\mathcal{E}(y^0, d^*)}_{\text{deletion output}}, \quad \underset{u \sim \text{uniform}[0,1]}{d^* \sim \pi^*} \text{ if } u < \beta \text{ else } \underbrace{\mathcal{E}(y^*, \tilde{d})}_{\text{random word dropping sequence of the round-truth}}, \quad \tilde{d} \sim \pi^{\text{RND}} \}$$

# Expert Policy

- Oracle:

Levenshtein distance

$$\mathbf{a}^* = \operatorname{argmin}_{\mathbf{a}} \mathcal{D}(\mathbf{y}^*, \mathcal{E}(\mathbf{y}, \mathbf{a}))$$

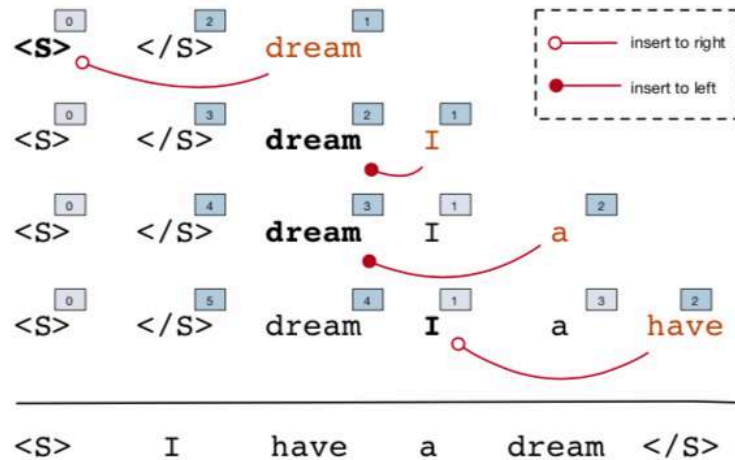
- Teacher Model:
  - first train an autoregressive teacher model using the same datasets and then replace the ground-truth sequence  $\mathbf{y}^*$  by the beam-search result of this teacher-model,  $\mathbf{y}^{AR}$

# Conclusion

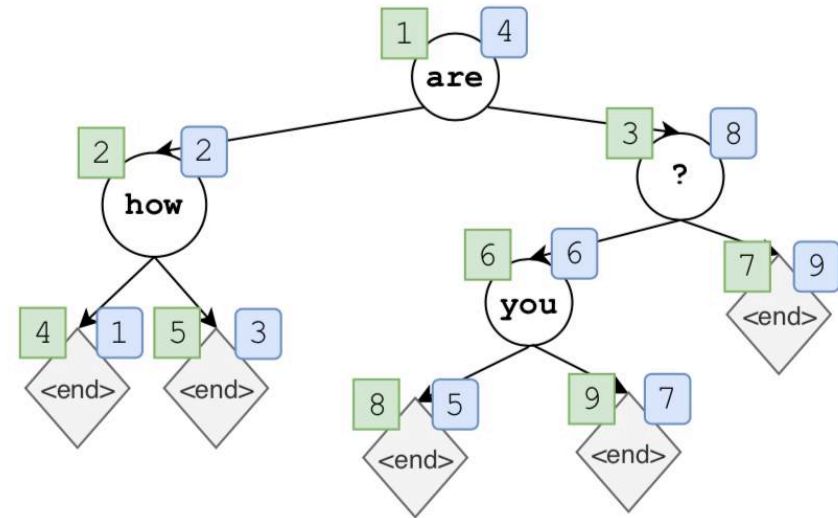
Parallel generation:

$t$	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)

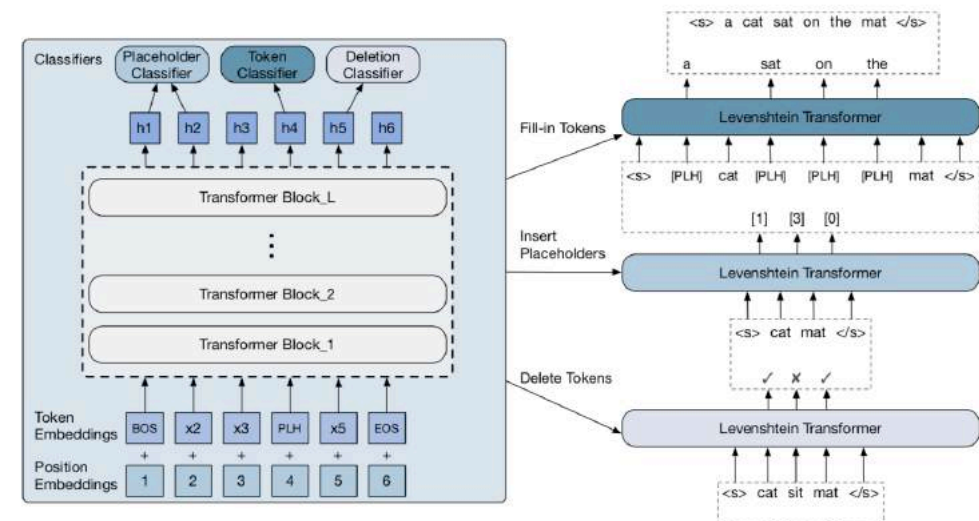
Insertion transformer



InDIGO



Non-Monotonic



Levenshtein



# Paper List

Paper	Conference
Levenshtein Transformer	
Insertion-based Decoding with automatically Inferred Generation Order	
KERMIT: Generative Insertion-Based Modeling for Sequences	
Non-Monotonic Sequential Text Generation	ICML19
Insertion Transformer: Flexible Sequence Generation via Insertion Operations	ICML19
Sequence Generation: From Both Sides to the Middle	IJCAI19
Correct-and-Memorize: Learning to Translate from Interactive Revisions	IJCAI19
Non-Autoregressive Neural Machine Translation	ICLR18
The Importance of Generation Order in Language Modeling	EMNLP18

# Reference

- 香侬读 | 按什么套路生成？基于插入和删除的序列生成方法  
<https://zhuanlan.zhihu.com/p/73417154>
- [https://cips-upload.bj.bcebos.com/ssatt2019/CIPS\\_SSATT\\_2019\\_问答系统\\_唐都钰\\_段楠.pdf](https://cips-upload.bj.bcebos.com/ssatt2019/CIPS_SSATT_2019_问答系统_唐都钰_段楠.pdf)

**Thanks!**