

Online Neural Machine Translation with 2D Convolutions

Maha Elbayad Laurent Besacier Jakob Verbeek

Under review at EMNLP 2019

July 5, 2019



- Low latency translation.
- 2D convolutional architecture for MT.
- 2D convolutional architecture for low latency MT.
- Deterministic and controllable latency.
- Experimental results.
- Conclusion.

Latency of sequence-to-sequence models

Seq2seq: a conditional language model that assigns probabilities to a sequence of tokens $\mathbf{y} = (y_1, \dots, y_{|\mathbf{y}|})$ given a conditioning sequence $\mathbf{x} = (x_1, \dots, x_{|\mathbf{x}|})$.

$$\log p_{\theta}(\mathbf{y}|\mathbf{x}) = \sum_{t=1}^{|\mathbf{y}|} \log p_{\theta}(y_t|\mathbf{x}, \mathbf{y}_{<t})$$

The model's parameters optimized through MLE:

$$\mathcal{L} = -\log p_{\theta}(\mathbf{y}|\mathbf{x})$$

Latency: stems from $y_t|\mathbf{x}$ as we have to read the full source \mathbf{x} .

Latency of sequence-to-sequence models

Motivation: build a seq2seq model for online / as-you-type translation.

Introduce a hidden variable z_t as the length of the context read by timestep t :

$$\begin{aligned} p_{\theta}(\mathbf{y}|\mathbf{x}) &= \sum_{\mathbf{z}} p_{\theta}(\mathbf{y}, \mathbf{z}|\mathbf{x}) \\ &= \sum_{\mathbf{z}} p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{y}|\mathbf{z}, \mathbf{x}) \\ p_{\theta}(\mathbf{y}|\mathbf{z}, \mathbf{x}) &= \prod_t p_{\theta}(y_t|\mathbf{x}_{\leq z_t}, \mathbf{y}_{<t}) \end{aligned}$$

Latency of sequence-to-sequence models

Two key points:

- 1 Encode $y_t | \mathbf{x}_{\leq z_t}, \mathbf{y}_{<t}$

$$p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x}) = \prod_t p_{\theta}(y_t | \mathbf{x}_{\leq z_t}, \mathbf{y}_{<t})$$

- 2 Model the hidden variables $\mathbf{z} = (z_1, \dots, z_{|\mathbf{y}|})$.

Latency of sequence-to-sequence models

- 1 Encode $y_t | x_{\leq z_t}, y_{<t}$

$$p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x}) = \prod_t p_{\theta}(y_t | \mathbf{x}_{\leq z_t}, \mathbf{y}_{<t})$$

With attention-based encoder-decoder architectures we need to re-score the source for every context size z_t .

With $(h_1, \dots, h_{|y|})$ the decoder states and $(s_1, \dots, s_{|x|})$ the encoder states (both evaluated autoregressively):

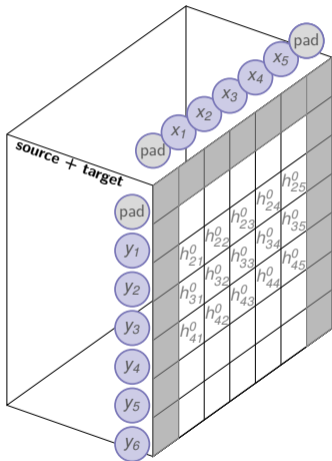
$$\alpha_{z_t} = \text{softmax}(e_i)_{1 \leq i \leq z_t}$$
$$e_i = \text{score}(h_t, s_i), \forall i$$

Latency of sequence-to-sequence models

- 1 Encode $y_t | \mathbf{x}_{\leq z_t}, \mathbf{y}_{<t}$

$$p_{\theta}(\mathbf{y} | \mathbf{z}, \mathbf{x}) = \prod_t p_{\theta}(y_t | \mathbf{x}_{\leq z_t}, \mathbf{y}_{<t})$$

Leverage our 2D convolutional architecture where target encodings with different context sizes are readily available.

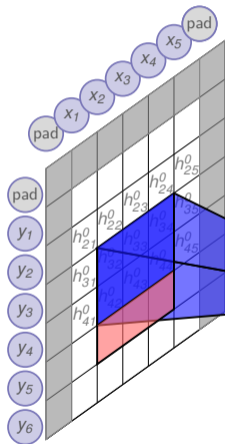


Input grid:

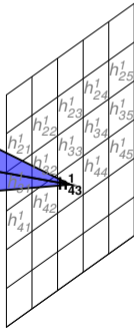
$$\forall i, j \ h_{ij}^0 = \text{concat}(\text{embed}(y_i), \text{embed}(x_j))$$

embed lookup tables for the source and target vocabularies.

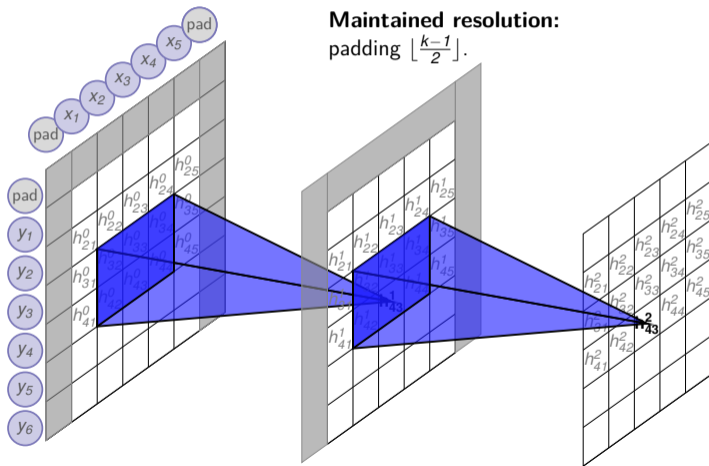
+**Padding**: to maintain the input resolution.

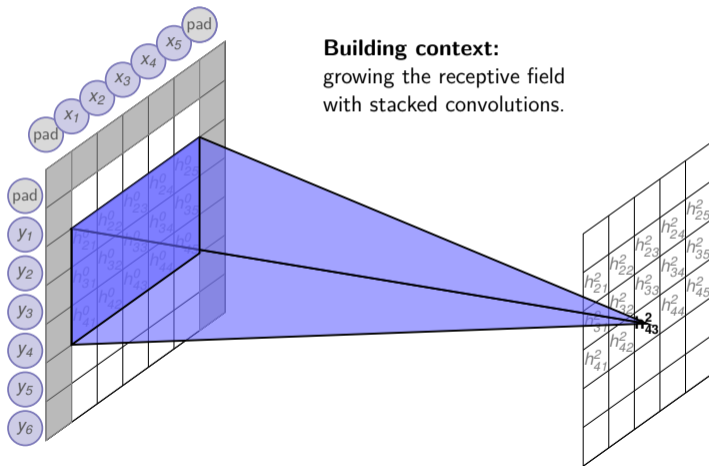


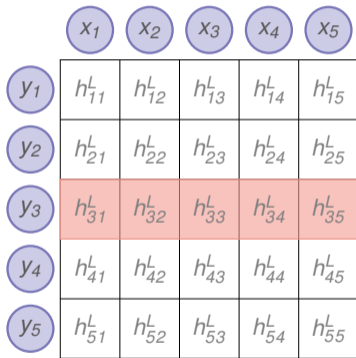
Masked convolution
Effective kernel $k \times (1 + \lfloor \frac{k-1}{2} \rfloor)$



Pervasive attention | 2D convolutional seq2seq model



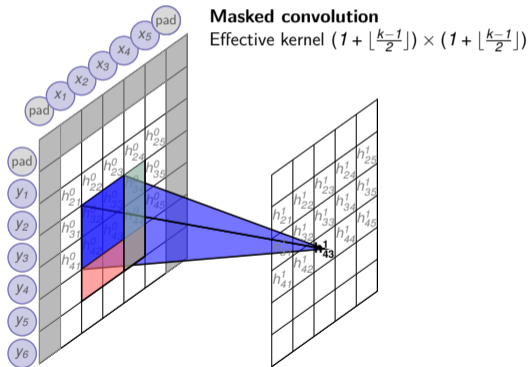




After L convolutional blocks:

- For each timestep t we have x representations.
- To model $y_t|x$ we need a fixed size representation.
- A basic aggregation: $h_t = \text{Max-pool}(h_j)_{1 \leq j \leq |x|}$

- Auto-regressive encoding of the source sequence with masked convolutions.



- Progressive pooling of the hidden states.

Max-pooling the states up to z_{t+1} .

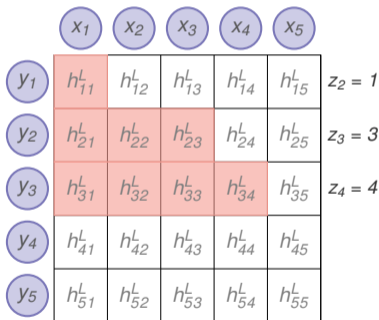
$$h_t = \text{Max-pool}(h_{t, \leq z_{t+1}}).$$

The state is fed to a classifier to predict the next token:

$$p_{\theta}(y_{t+1} | y_{\leq t}, \mathbf{x}_{\leq z_{t+1}}) = \text{softmax}(Wh_t)$$

Less effective alternatives:

- Without pooling: $h_t = h_{t, z_t}$.
- Grid pooling: $h_t = \text{Max-pool}(h_{\leq t, \leq z_t})$.



Two key points:

- 1 Encode $y_t | x_{\leq z_t}, y_{<t}$
- 2 Model the hidden variables $\mathbf{z} = (z_1, \dots, z_{|y|})$.

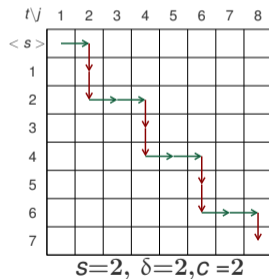
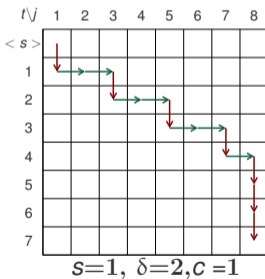
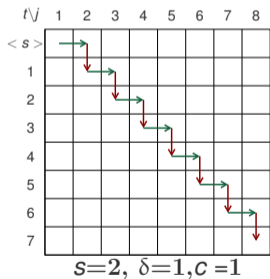
Two modelizations of \mathbf{z} through a '*Controller*':

- Deterministic ($\mathbf{z} \perp \mathbf{x}, \mathbf{y}$) following a pre-defined schedule.
- Dynamic with \mathbf{z} the latent variables of an HMM.

Deterministic z

$$z_t = s + \lfloor \frac{t-1}{c} \rfloor \delta$$

- Shift (s): context size at $t = 1$.
- Stepsize (δ): consecutive reads.
- Catchup (c): consecutive writes.



Two possible options for incorporating the variable \mathbf{z} in the model:

- **At inference:**

With a model trained to generate after reading the full \mathbf{x} (Wait-until-end: WUE), force the model to emit an output following the pre-determined \mathbf{z} .

- **In both training and inference:** Optimize the decoding cross-entropy for the chosen context:

From:

$$\mathcal{L} = - \sum_t \log p(y_t | \mathbf{y}_{<t}, \mathbf{x})$$

To:

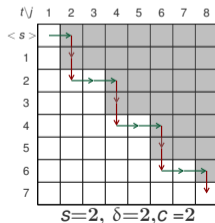
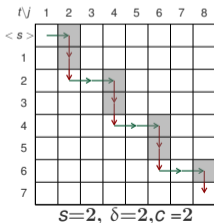
$$\mathcal{L} = - \sum_t \log p(y_t | \mathbf{y}_{<t}, \mathbf{x}_{\leq z_t})$$

- **In both training and inference:** Optimize the decoding cross-entropy for the chosen context:

$$\mathcal{L} = - \sum_t \log p(y_t | \mathbf{y}_{<t}, \mathbf{x}_{\leq z_t}) \quad (\text{Path})$$

We find it beneficial to include larger contexts in the optimization:

$$\mathcal{L} = - \sum_t \sum_{c \geq z_t} \log p(y_t | \mathbf{y}_{<t}, \mathbf{x}_{\leq c}) \quad (\text{Above path})$$



Build an HMM with latent variables \mathbf{z} and observations \mathbf{y} .

- 1 At every position (t, j) make a decision read=1/write=0:

$$\rho_{tj} = \sigma(C^\top h_{tj}),$$

- 2 Estimate the HMM transition probabilities from read/write decisions:

$$M_{tjk} = p_\theta(z_{t+1} = k | z_t = j, \mathbf{y}_{<t}, \mathbf{x}) = \begin{cases} (1 - \rho_{tk}) \prod_{l=j}^{k-1} \rho_{tl}, & \text{if } k \geq j, \\ 0, & \text{otherwise (non-decreasing } z) \end{cases}$$

Build an HMM with latent variables \mathbf{z} and observations \mathbf{y} .

③ EM algorithm:

$$\begin{aligned}\log p_{\theta}(\mathbf{y}|\mathbf{x}) &\geq \log p(\mathbf{y}|\mathbf{x}) - \mathcal{D}_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y}, \mathbf{x})) \\ &\geq \mathcal{H}(q) + \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{z}, \mathbf{y}|\mathbf{x})\end{aligned}$$

Build an HMM with latent variables \mathbf{z} and observations \mathbf{y} .

③ EM algorithm:

$$\begin{aligned}\log p_{\theta}(\mathbf{y}|\mathbf{x}) &\geq \log p(\mathbf{y}|\mathbf{x}) - \overbrace{\mathcal{D}_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y}, \mathbf{x}))}^0 \\ &\geq \mathcal{H}(q) + \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{z}, \mathbf{y}|\mathbf{x})\end{aligned}$$

① E-step: $q := p_{\theta}(\mathbf{z}|\mathbf{y}, \mathbf{x})$. Infer the posterior via Baum-Welch.

Build an HMM with latent variables z and observations y .

③ EM algorithm:

$$\begin{aligned}\log p_{\theta}(\mathbf{y}|\mathbf{x}) &\geq \log p(\mathbf{y}|\mathbf{x}) - \overbrace{\mathcal{D}_{\text{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{y}, \mathbf{x}))}^0 \\ &\geq \underbrace{\mathcal{H}(q) + \sum_{\mathbf{z}} q(\mathbf{z}) \log p_{\theta}(\mathbf{z}, \mathbf{y}|\mathbf{x})}_{Q}\end{aligned}$$

- 1 E-step: $q := p_{\theta}(\mathbf{z}|\mathbf{y}, \mathbf{x})$. Infer the posterior via Baum-Welch.
- 2 M-step: $\theta := \arg \max_{\theta} Q(\theta)$. SGD

$$Q(q, \theta) = \mathcal{W}(q, \theta) + \mathcal{C}(q, \theta).$$

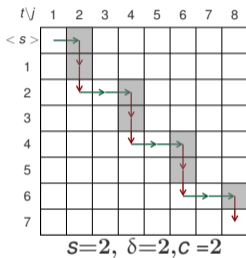
$$\mathcal{W}(q, \theta) = - \sum_{t=1}^{|\mathbf{y}|} \sum_{j=1}^{|\mathbf{x}|} q(z_t = j) \log p_{\theta}(y_t | \mathbf{y}_{<t}, \mathbf{x}_{\leq j}). \quad (\text{writing})$$

$$\mathcal{C}(q, \phi) = \sum_{t=1}^{|\mathbf{y}|-1} \sum_{j=1}^{|\mathbf{x}|} A_{tj} \log \rho_{tj} + B_{tj} \log(1 - \rho_{tj}) \quad (\text{controlling})$$

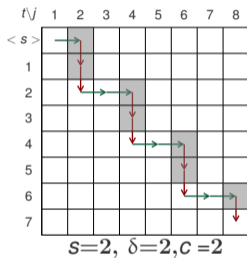
$$A_{tj} = \sum_{l \leq j} \sum_{k > j} q(z_t = l, z_{t+1} = k) \quad ((t, j) \rightarrow (t, j + 1))$$

$$B_{tj} = q(z_{t+1} = j) \quad ((t, j) \downarrow (t + 1, j))$$

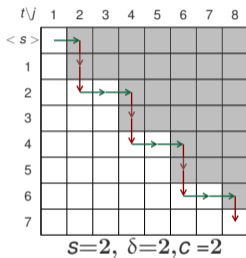
The controller | recap



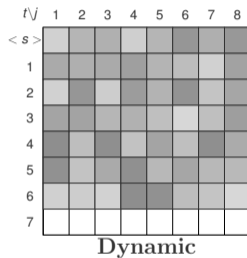
(a) Masked pervasive attention (MPA)



(b) Deterministic - Path



(c) Deterministic - Above path

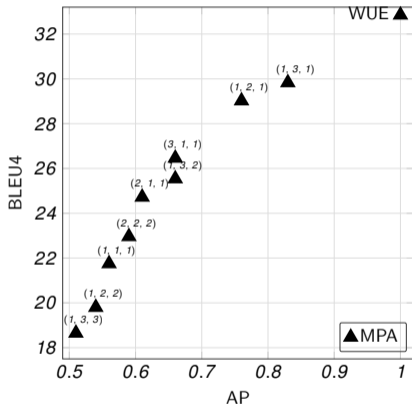


(d) Dynamic

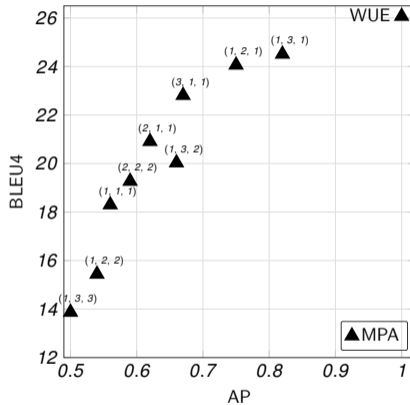
Dataset	Vocab (BPE)	Train	Valid	Test
IWSLT14 De-En/En-De	En(6632), De(8848)	160239	7283	6633
WMT15 Ru-En/En-ru	En(18480), Ru(20000)	806996	8182	3000

- For training, we filter out (\mathbf{x}, \mathbf{y}) if $|\mathbf{x}| > 50$ or $|\mathbf{y}| > 50$.
- Inference with a beam of size 5
- Evaluate the translation quality with case-sensitive tokenized BLEU.
- Decoding latency measured with the average proportion (AP):

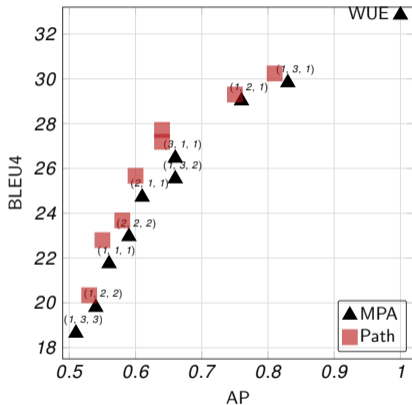
$$AP = \frac{1}{|\mathbf{y}|} \sum_t \frac{z_t}{|\mathbf{x}|}.$$



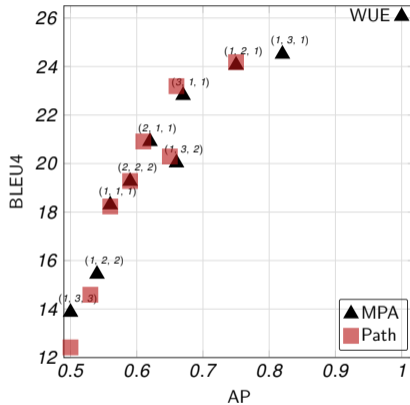
(a) IWSLT14 De-En



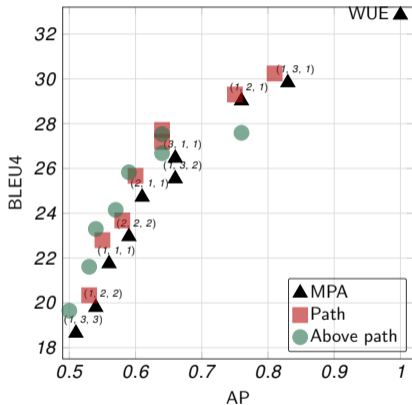
(b) IWSLT14 En-De



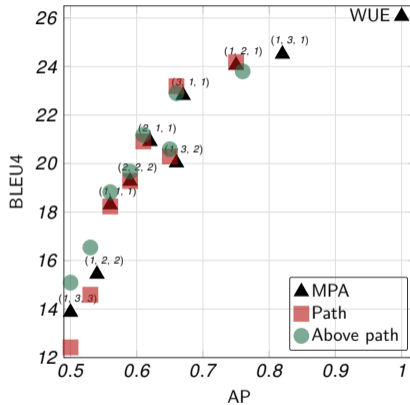
(a) IWSLT14 De-En



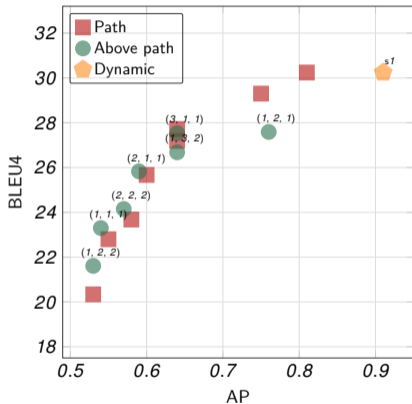
(b) IWSLT14 En-De



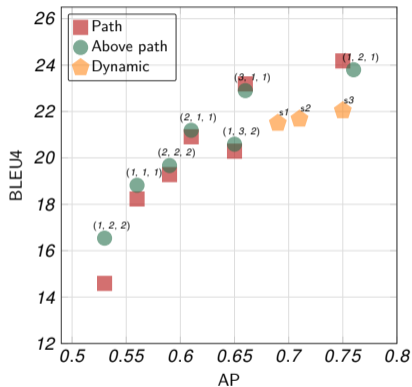
(a) IWSLT14 De-En



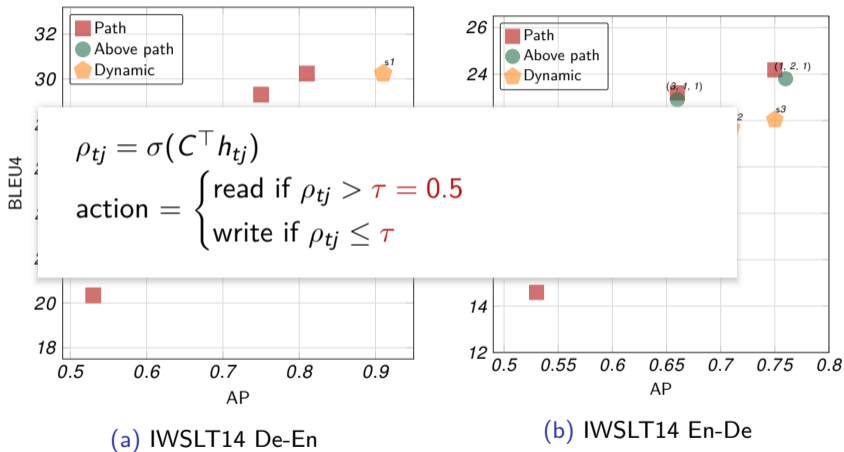
(b) IWSLT14 En-De

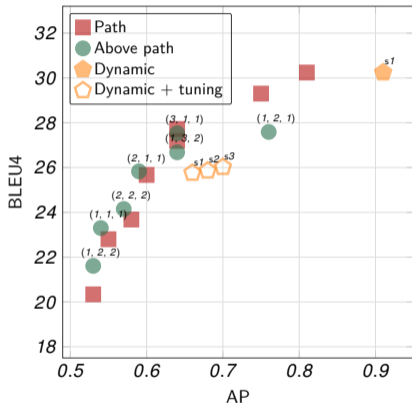


(a) IWSLT14 De-En

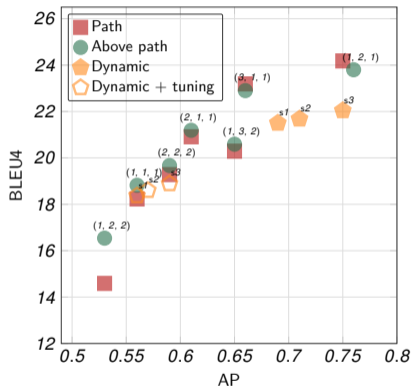


(b) IWSLT14 En-De

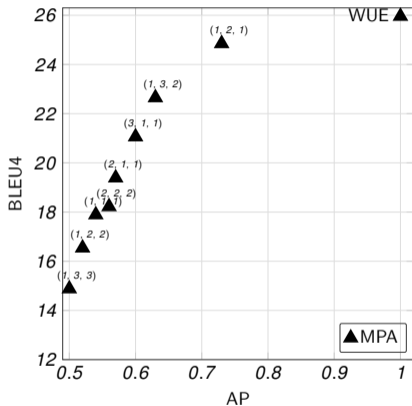




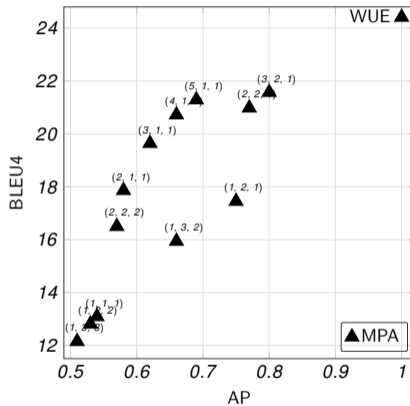
(a) IWSLT14 De-En



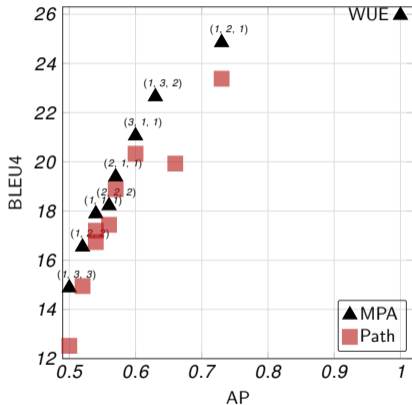
(b) IWSLT14 En-De



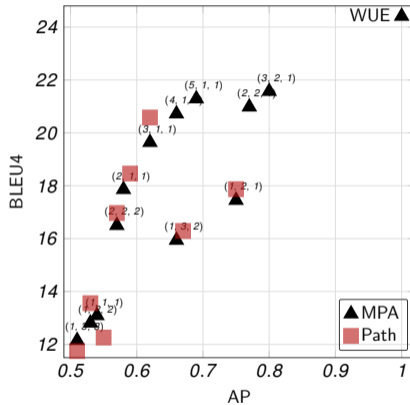
(a) WMT15 Ru-En



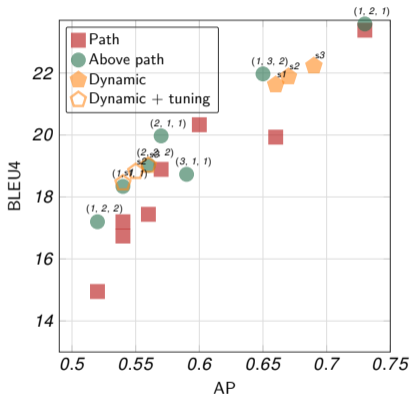
(b) WMT15 En-Ru



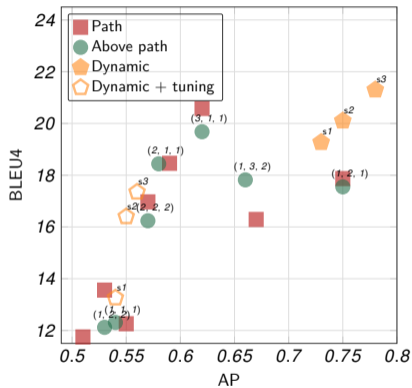
(a) WMT15 Ru-En



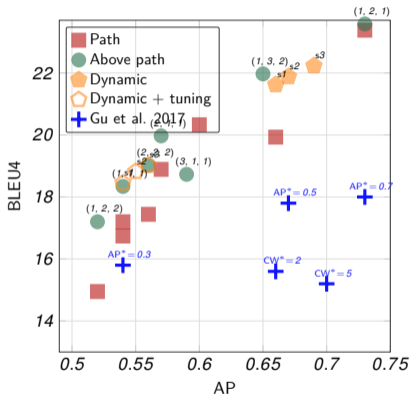
(b) WMT15 En-Ru



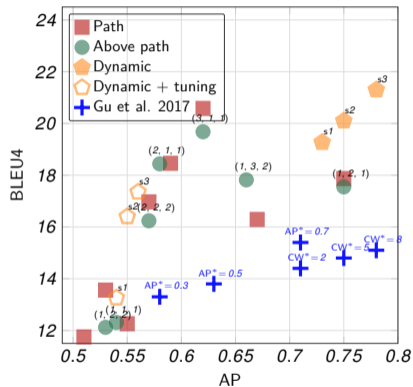
(a) WMT15 Ru-En



(b) WMT15 En-Ru

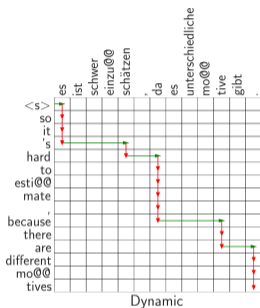
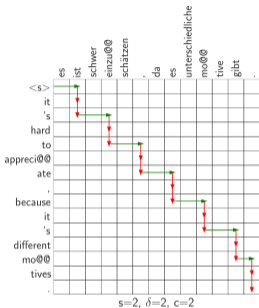
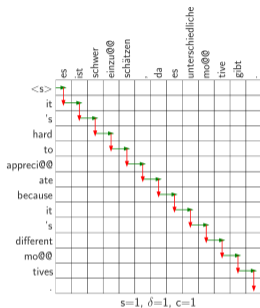
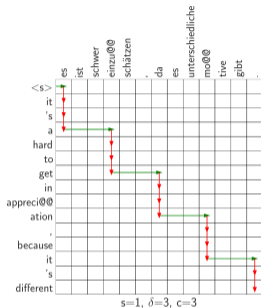


(a) WMT15 Ru-En



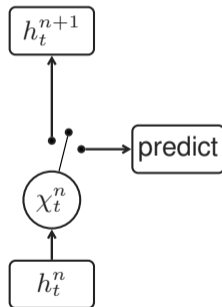
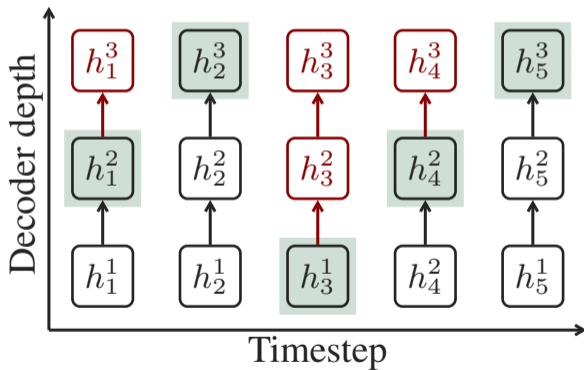
(b) WMT15 En-Ru

Conclusion



- Both deterministic and dynamic models achieve a good BLEU-AP trade-off i.e correct translations with low delay.
- Future work: Analyze whether read/write decisions reveal any meaningful linguistic structure.

facebook Artificial Intelligence



Thank you for your attention.
Questions / Feedback.