# Simultaneous Neural Machine Translation

Ashkan Alinejad

Supervisor : Dr. Anoop Sarkar

Simon Fraser University

July 2017

Machine
Translation
RNN Encoder-
Decoder

Simultaneous
Machine
Translation

Neural SMT
Greedy
Decoding
Trainable Agent

**1** Machine Translation
 RNN Encoder-Decoder

**2** Simultaneous Machine Translation

**3** Neural SMT
 Greedy Decoding
 Trainable Agent

# Neural Machine Translation

**RNN structure**

- Encoder

$$h_t = f(x_t, h_{t-1})$$
$$c = q(\{h_1, \ldots, h_{T_x}\})$$

- Decoder

$$p(y) = \prod_{t=1}^{T} p(y_t | \{y_1, \ldots, y_{t-1}\}, c)$$

With an RNN, each conditional probability is modeled as:

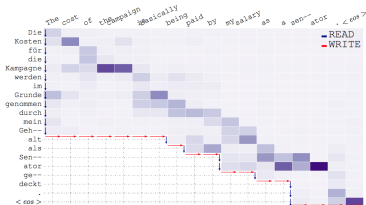$$p(y_t | \{y_1, \ldots, y_{t-1}\}, c) = g(y_{t-1}, s_t, c)$$

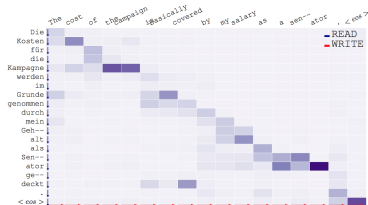# Simultaneous Machine Translation

- Simultaneous Machine Translation is a challenging task of reading from the source language and at the same time, producing the target translation.

- The objective of translation system is defined as a combination of quality and delay.



(a) Simultaneous Neural Machine Translation    (b) Neural Machine Translation

**Previous works**

- Most of the works in this direction are done in the context of speech translation. incoming speech is transcribed and segmented into a translation unit largely based on acoustic and linguistic cues.

- Each of these segments is then translated largely independent from each other

# Neural SMT

Machine
Translation
RNN Encoder-
Decoder

Simultaneous
Machine
Translation

Neural SMT

Greedy
Decoding
Trainable Agent

- Sequentially making two interleaved decisions:
  1. READ
  2. WRITE

- 

$$\text{Input sequence} \quad X = \{x_1, \ldots, x_{T_s}\}$$
$$\text{Decoded Output} \quad Y = \{y_1, \ldots, y_{T_t}\}$$
$$\text{Action sequence} \quad A = \{a_1, \ldots, a_T\}$$

$$T = T_s + T_t$$

• The model structure is an attention-based neural network

$$\text{Encoder} \quad : \quad h_\eta = \phi_{\text{UNI-ENC}}(h_{\eta-1}, x_\eta)$$

$$\text{Decoder} \quad : \quad c_\tau^\eta = \phi_{\text{ATT}}(z_{\tau-1}, y_{\tau-1}, H^\eta)$$

$$z_\tau^\eta = \phi_{\text{DEC}}(z_{\tau-1}, y_{\tau-1}, c_\tau^\eta)$$

$$\text{Output} \quad : \quad p(y|y_{<\tau}, H^\eta) \propto \exp[\phi_{\text{OUT}}(z_\tau^\eta)]$$

$$y_\tau^\eta = \arg \max_y p(y|y_{<\tau}, H^\eta)$$

---

**Algorithm 1** Simultaneous Greedy Decoding

**Require:** $\delta$, $s_0$, Input Pipe $X$, Output Pipe $Y$

 1: Initialize $s \leftarrow s_0$, $C \leftarrow \text{READ}(X, s)$, $C' \leftarrow \{\}$

 2: Initialize the decoder's state $\mathbf{z}_0$ based on $C$

 3: **while true do**

 4:    $\hat{y}_t = \arg\max_{y_t} \log p(y_t | y_{<t}, C)$

 5:    **if** $s \geq T_X$ **then**

 6:        $\text{WRITE}(Y, \hat{y}_t)$, $t \leftarrow t + 1$

 7:    **else**

 8:        $C' \leftarrow \text{READ}(X, \delta)$ if $|C'| = 0$.

 9:        **if** $\Lambda(C, C \cup C')$ **then**

10:           $C \leftarrow C \cup C'$, $s \leftarrow s + \delta$, $C' \leftarrow \{\}$

11:           **continue**

12:        **else**

13:           $\text{WRITE}(Y, \hat{y}_t)$, $t \leftarrow t + 1$

14:        **end if**

15:    **end if**

16:    **if** $\hat{y}_t = \langle \text{eos} \rangle$ **then**

17:        **break**

18:    **end if**

19: **end while**

Machine
Translation
RNN Encoder-
Decoder

Simultaneous
Machine
Translation

Neural SMT
**Greedy
Decoding**
Trainable Agent

- **Wait-If-Worse**

$$\Lambda(C, C \cup C') : (\log p(\hat{y}|\hat{y}_{<t}, C) \\ > \log p(\hat{y}|\hat{y}_{<t}, C \cup C')),$$

where $\hat{y} = \arg \max_y p(y|\hat{y}_{<t}, C)$
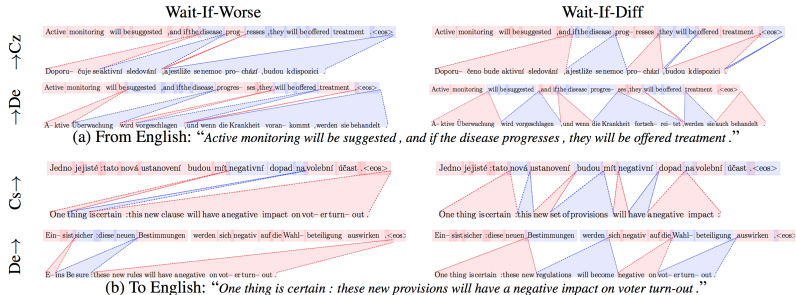
- **Wait-If-Diff**

$$\Lambda(C, C \cup C') : (\hat{y} \neq \hat{y}'),$$

where $\hat{y}' = \arg \max_y \log p(y|\hat{y}_{<t}, C \cup C')$.

Machine
Translation
RNN Encoder-
Decoder

Simultaneous
Machine
Translation

Neural SMT
Greedy
Decoding
Trainable Agent

**Metrics**

- **Quality** The metrics for evaluating quality of the translation is the BLEU score.

- **Delay** $s(t) =$ In each time step for the decoded target symbol $\hat{y}_t$, how many source symbols were required. delay in translation $(T)$:

$$0 < T(X, \hat{Y}) = \frac{1}{|X||\hat{Y}|} \sum_{t=1}^{|\hat{y}|} s(t) \le 1.$$

(a) From English: "*Active monitoring will be suggested , and if the disease progresses , they will be offered treatment* ."

(b) To English: "*One thing is certain : these new provisions will have a negative impact on voter turn-out* ."

|         |      | Cs    | De    | Ru    |
|---------|------|-------|-------|-------|
| En $\uparrow$ | Ours | 15.2  | 19.5  | 17.77 |
|         | $\star$ | 13.84 | 21.75 | 19.54 |
| $\uparrow$ En | Ours | 20.47 | 23.96 | 22.27 |
|         | $\star$ | 20.32 | 24    | 22.44 |

Figure: BLEU scores on the test set (newstest-2015) obtained by the
models used in the paper and ($\star$) from (Firat et al., 2016). Although
our models use a unidirectional recurrent net as an encoder, the
translation qualities are comparable.

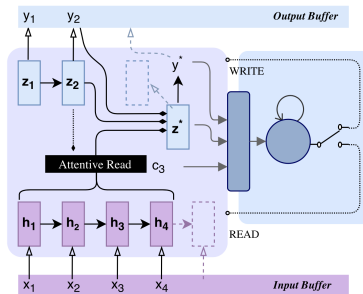**Discussion**

1. They do not have good BLEU score compared to previous works.

2. the waiting criteria proposed in this paper are both manually designed and does not exploit rich information embedded in the hidden representation learned by the recurrent neural networks.

3. The objective of the network is to improve translation quality and do not consider delay during training.

## Trainable Agent

- The idea is to have a separate trainable agent
- The framework can be trained using reinforcement learning and it considers both Quality and Delay during training.

---

**Algorithm 1** Simultaneous Greedy Decoding

---

**Require:** NMT system $\phi$, policy $\pi_\theta$, $\tau_{\text{MAX}}$, input buffer $X$, output buffer $Y$, state buffer $S$.

1: **Init** $x_1 \Leftarrow X$, $h_1 \leftarrow \phi_{\text{ENC}}(x_1)$, $H^1 \leftarrow \{h_1\}$
2:     $z_0 \leftarrow \phi_{\text{INIT}}(H^1)$, $y_0 \leftarrow \langle s \rangle$
3:     $\tau \leftarrow 0$, $\eta \leftarrow 1$
4: **while** $\tau < \tau_{\text{MAX}}$ **do**
5:     $t \leftarrow \tau + \eta$
6:     $y_\tau^\eta, z_\tau^\eta, o_t \leftarrow \phi(z_{\tau-1}, y_{\tau-1}, H^\eta)$
7:     $a_t \sim \pi_\theta(a_t; a_{<t}, o_{<t})$, $S \Leftarrow (o_t, a_t)$
8:     **if** $a_t = \text{READ}$ and $x_\eta \neq \langle/s\rangle$ **then**
9:       $x_{\eta+1} \Leftarrow X$, $h_{\eta+1} \leftarrow \phi_{\text{ENC}}(h_\eta, x_{\eta+1})$
10:      $H^{\eta+1} \leftarrow H^\eta \cup \{h_{\eta+1}\}$, $\eta \leftarrow \eta + 1$
11:      **if** $|Y| = 0$ **then** $z_0 \leftarrow \phi_{\text{INIT}}(H^\eta)$
12:     **else if** $a_t = \text{WRITE}$ **then**
13:      $z_\tau \leftarrow z_\tau^\eta$, $y_\tau \leftarrow y_\tau^\eta$
14:      $Y \Leftarrow y_\tau$, $\tau \leftarrow \tau + 1$
15:      **if** $y_\tau = \langle/s\rangle$ **then break**

---

### Agent

A trainable agent is designed to make decisions
$A = \{a_1, \ldots, a_T\}$, $a_t \in \mathcal{A}$ sequentially based on observations
$O = \{o_1, \ldots, o_T\}$, $o_t \in \mathcal{O}$.

- **Observation:**  $o_{\tau+\eta} = [c_\tau^\eta; z_\tau^\eta; E(y_\tau^\eta)]$

- **Action:**
    - READ: waits to encode the next word
    - WRITE: accepts the candidate and emits it as the prediction

- **Policy:**  a stochastic policy $\pi_\theta$ parameterized by a recurrent neural network

$$s_t = f_\theta(s_{t-1}, o_t),$$
$$\pi_\theta(a_t | a_{<t}, o_{\leq t}) \propto g_\theta(s_t)$$

### Reward Function

At each step the agent will receive a reward signal $r_t$ based on $(o_t,\ a_t)$.

- **Quality** $r_t^Q =$ smoothed BLEU
- **Delay** $r_t^D$
  1. **Average Proportion**
  2. **Consecutive Wait Length**

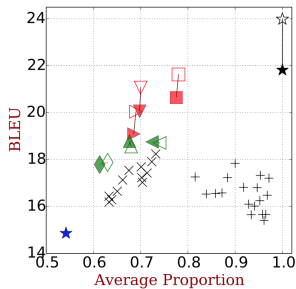The total reward will be computed as:

$$r_t = r_t^Q + r_t^D$$
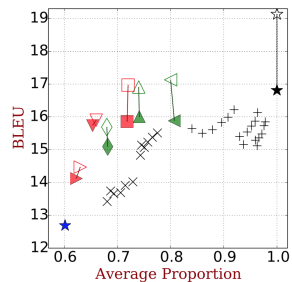
---

**Algorithm 2** Learning with Policy Gradient

---

**Require:** NMT system $\phi$, agent $\theta$, baseline $\varphi$
 1: Pretrain the NMT system $\phi$ using MLE;
 2: Initialize the agent $\theta$;
 3: **while** stopping criterion fails **do**
 4:     Obtain a translation pairs: $\{(X, Y^*)\}$;
 5:     **for** $(Y, S) \sim$ Simultaneous Decoding **do**
 6:         **for** $(o_t, a_t)$ in $S$ **do**
 7:             Compute the quality: $r_t^Q$;
 8:             Compute the delay: $r_t^D$;
 9:             Compute the baseline: $b_\varphi(o_t)$;
10:     Collect the future rewards: $\{R_t\}$;
11:     Perform variance reduction: $\{\tilde{R}_t\}$;
12:     Update: $\theta \leftarrow \theta + \lambda_1 \nabla_\theta [J - \kappa \mathcal{H}(\pi_\theta)]$
13:     Update: $\varphi \leftarrow \varphi - \lambda_2 \nabla_\varphi L$

---

# Results



(d) DE→EN

(c) EN→DE

(◀ ◁: CW=8, ▲△: CW=5, ◆◇: CW=2, ▶ ▷: AP=0.3, ▼▽: AP=0.5, ■□: AP=0.7). For each target, we select the model

**Discussion**

Machine
Translation
RNN Encoder-
Decoder

Simultaneous
Machine
Translation

Neural SMT
Greedy
Decoding
Trainable Agent

📄 Yajie Miao, Mohammad Gowayyed, and Florian Metze.
Eesen: End-to-end speech recognition using deep rnn
models and wfst-based decoding.
07 2015.

📄 Awni Hannun, Carl Case, Jared Casper, Bryan Catanzaro,
Greg Diamos, Erich Elsen, Ryan Prenger, Sanjeev Satheesh,
Shubho Sengupta, Adam Coates, and Andrew Y. Ng.
Deep speech: Scaling up end-to-end speech recognition, 12
2014.

📄 Alex Graves.
*Supervised Sequence Labelling with Recurrent Neural
Networks*.
PhD thesis, Technische Universität München, July 2008.

📄 Vinod Nair and Geoffrey E. Hinton.
Rectified linear units improve restricted boltzmann
machines.
In Johannes Fürnkranz and Thorsten Joachims, editors,

# Thank You !