# Recent Advances and the Future of Neural Machine Translation

**Orhan Firat**[1]     Kyunghyun Cho[2]

[1]Middle East Technical University
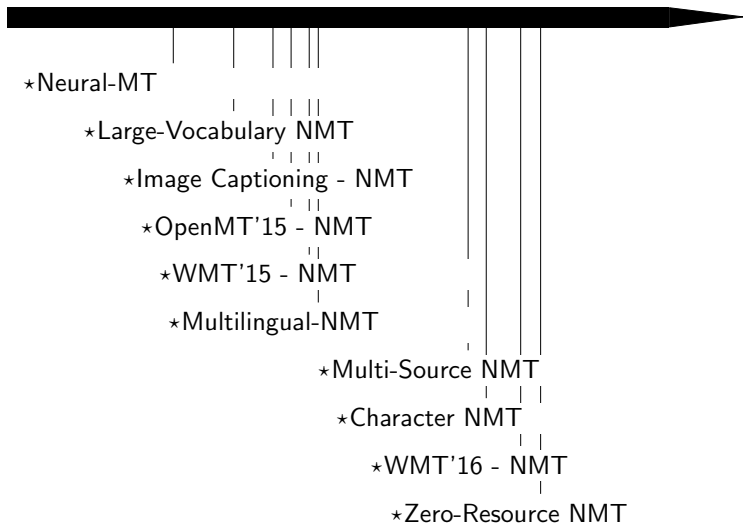[2]New York University

Machine Translation Marathon 2016

The Fog of Progress[1]

and

Artificial *General* Intelligence

---

[1] Hinton video-lectures,https://www.youtube.com/watch?v=ZuvRXGX8cY8

# What is going on?

2014                                                                    2017



⋆Neural-MT

⋆Large-Vocabulary NMT

⋆Image Captioning - NMT

⋆OpenMT'15 - NMT

⋆WMT'15 - NMT

⋆Multilingual-NMT

⋆Multi-Source NMT

⋆Character NMT

⋆WMT'16 - NMT

⋆Zero-Resource NMT
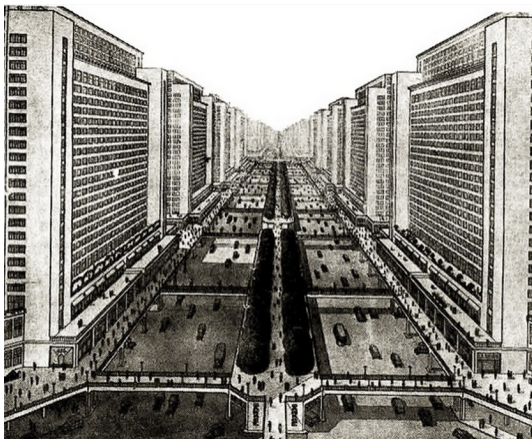
# Warren Weaver-"Translation", 1949

Tall towers analogy:

- Do not shout from tower to tower,
- Go down to the common basement of all towers: *interlingua*

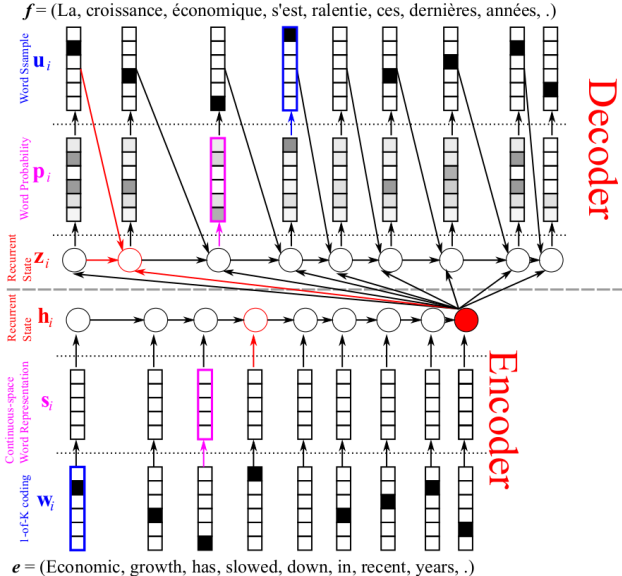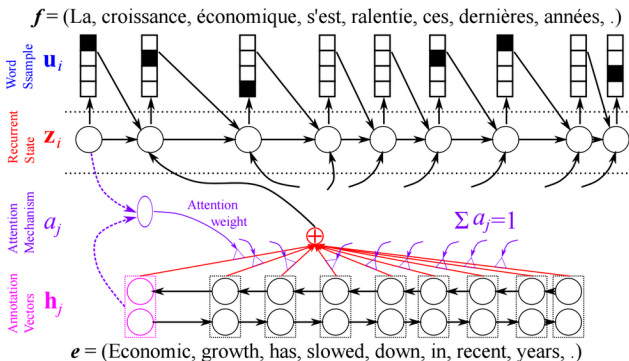# Neural Machine Translation - Encoder Decoder



figure credit, Kyunghyun Cho

# Encoder-Decoder Architecture with Attention

Bahdanau et al.2015



At each timestep in the decoder:

1. Computes a *relevance* score of each annotation
2. Use the weighted sum of the annotations as a *context*

figure credit, Kyunghyun Cho

# Encoder-Decoder Architecture with Attention

Bahdanau et al.'15



$f$ = (La, croissance, économique, s'est, ralentie, ces, dernières, années, .)

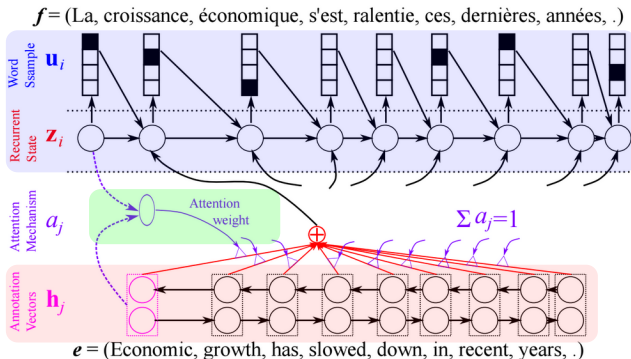$e$ = (Economic, growth, has, slowed, down, in, recent, years, .)
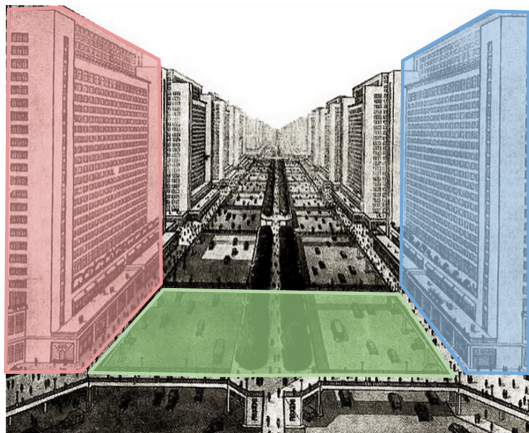
At each timestep in the decoder:

1. Computes a *relevance* score of each annotation
2. Use the weighted sum of the annotations as a *context*

figure credit, Kyunghyun Cho

# Warren Weaver-"Translation", 1949

Tall towers analogy:

- ‣ Red Tower : source language
- ‣ Blue Tower: target language
- ‣ Green Car : alignment function

# Attention-based NMT at work - WMT'16

Seems to be working!



Most of the top-rankers used NMT

# Neural Machine Translation with Finer Tokens

Let's make a poll on pre-processing! 🐥🥷



😫 Why do we use word-level modelling?

- ▸ Words are basic unit of meaning?!
- ▸ Inherent fear of sparsity!
- ▸ Finer granularities → longer sequences

👹 Why can't an NMT system directly learn from the characters?

# Neural Machine Translation with Finer Tokens

Issues with tokenization and segmentation

- ‣ Ineffective way of handling morphological variants:
  'run', 'runs', 'running' and 'runner'
- ‣ How are we doing with compound words?

Issues with treating each and every token separately

- ‣ Fill the vocabulary with similar words
- ‣ Vocabulary size grows linearly w.r.t. the corpus size
- ‣ Rare words, numbers and misspelled words:
  9/11 is a huge contextual information
- ‣ Lose the learning signal of words marked as <UNK>

---

# Granularity in Input and Output Spaces (finer tokens)



Word Level     Sub-Word Level     Character Level

Byte/Unicode Level

[I, really, enjoyed, this, film, .]

[I, real@@, ly, enjoy@@, ed, this, film, .]

**Tokenization**

**BPE-based Segmentation**

**Nothing**

(Sennrich et al.'15, Sennrich et al.'16)

(Costa-jussa et al.'16, Chung et al.'16, Luong et al.'16, **Lee et al.'16**)

[I, _, r, e, a, l, l, y, _, e, n, j, o, y, e, d, _, t, h, i, s, _, f, i, l, m, .]

# Character-Level Decoder without Explicit Segmentation
Chung, Cho and Bengio, ACL'16

Model details,

- ▸ RNNSearch Model
- ▸ Source Side : sub-words (byte pair encoding, BPE)
- ▸ Target Side : either sub-words or characters
- ▸ Three types of decoders:
  1. Sub-word level *base* decoder
  2. Character level *base* decoder
  3. Character level *bi-scale* decoder

Bi-scale decoder:

- ▸ Faster/slower layers for modelling different levels of tokens
- ▸ Use soft gating units for differentiability

# Character-Level Decoder without Explicit Segmentation
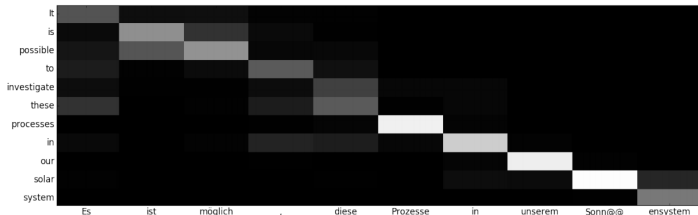
Chung, Cho and Bengio, ACL'16

| | | Src | Trgt | Depth | Attention h1 | Attention h2 | Model | Development Single | Development Ens | Test1 Single | Test1 Ens | Test2 Single | Test2 Ens |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| En-De | (a) | BPE | BPE | 1 | ✓ | | Base | 20.78 | – | 19.98 | – | 21.72 | – |
| | (b) | | | 2 | ✓ | ✓ | Base | $21.26_{20.62}^{21.45}$ | 23.49 | $20.47_{19.30}^{20.88}$ | 23.10 | $22.02_{21.35}^{22.21}$ | 24.83 |
| | (c) | | Char | 2 | | ✓ | Base | $21.57_{20.88}^{21.88}$ | 23.14 | $\mathbf{21.33}_{19.82}^{21.53}$ | 23.11 | $\mathbf{23.45}_{21.72}^{23.91}$ | 25.24 |
| | (d) | | | 2 | ✓ | ✓ | | 20.31 | – | 19.70 | – | 21.30 | – |
| | (e) | | | 2 | | ✓ | Bi-S | $21.29_{21.13}^{21.43}$ | 23.05 | $21.25_{20.62}^{21.47}$ | 23.04 | $23.06_{22.85}^{23.47}$ | 25.44 |
| | (f) | | | 2 | ✓ | | | 20.78 | – | 20.19 | – | 22.26 | – |
| | (g) | | | 2 | ✓ | | | 20.08 | – | 19.39 | – | 20.94 | – |
| | State-of-the-art Non-Neural Approach* | | | | | | | – | | 20.60[1] | | 24.00[2] | |
| En-Cs | (h) | BPE | BPE | 2 | ✓ | ✓ | Base | $16.12_{15.96}^{16.96}$ | 19.21 | $17.16_{16.38}^{17.68}$ | 20.79 | $14.63_{14.26}^{15.09}$ | 17.61 |
| | (i) | | Char | 2 | | ✓ | Base | $17.68_{17.39}^{17.78}$ | 19.52 | $19.25_{18.89}^{19.55}$ | 21.95 | $\mathbf{16.98}_{16.81}^{17.17}$ | 18.92 |
| | (j) | | | 2 | | ✓ | Bi-S | $17.62_{17.43}^{17.93}$ | 19.83 | $\mathbf{19.27}_{19.15}^{19.53}$ | 22.15 | $16.86_{16.68}^{17.10}$ | 18.93 |
| | State-of-the-art Non-Neural Approach* | | | | | | | – | | 21.00[3] | | 18.20[4] | |
| En-Ru | (k) | BPE | BPE | 2 | ✓ | ✓ | Base | $18.56_{18.26}^{18.70}$ | 21.17 | $25.30_{24.95}^{25.40}$ | 29.26 | $19.72_{19.02}^{20.29}$ | 22.96 |
| | (l) | | Char | 2 | | ✓ | Base | $18.56_{18.39}^{18.87}$ | 20.53 | $\mathbf{26.00}_{25.04}^{26.07}$ | 29.37 | $\mathbf{21.10}_{20.14}^{21.24}$ | 23.51 |
| | (m) | | | 2 | | ✓ | Bi-S | $18.30_{17.88}^{18.54}$ | 20.53 | $25.59_{24.57}^{25.76}$ | 29.26 | $20.73_{19.97}^{21.02}$ | 23.75 |
| | State-of-the-art Non-Neural Approach* | | | | | | | – | | 28.70[5] | | 24.30[6] | |
| En-Fi | (n) | BPE | BPE | 2 | ✓ | ✓ | Base | $9.61_{9.24}^{10.02}$ | 11.92 | – | – | $8.97_{8.88}^{9.17}$ | 11.73 |
| | (o) | | Char | 2 | | ✓ | Base | $11.19_{11.09}^{11.55}$ | 13.72 | – | – | $\mathbf{10.93}_{10.11}^{11.56}$ | 13.48 |
| | (p) | | | 2 | | ✓ | Bi-S | $10.73_{10.40}^{11.04}$ | 13.39 | – | – | $10.24_{9.71}^{10.63}$ | 13.32 |
| | State-of-the-art Non-Neural Approach* | | | | | | | – | | – | | 12.70[7] | |

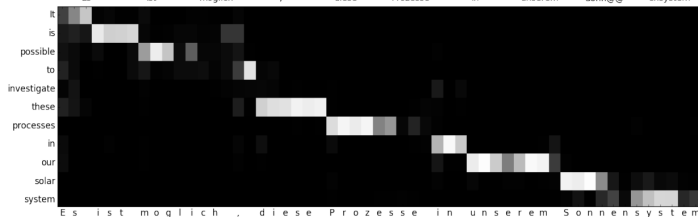# Character-Level Decoder without Explicit Segmentation

Chung, Cho and Bengio, ACL'16



"It is possible to investigate these processes in our solar system"
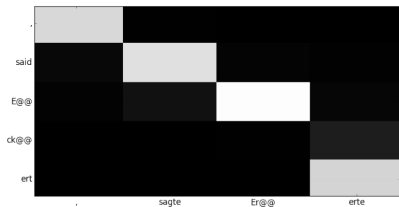
BPE Decoder

Char Decoder

# Character-Level Decoder without Explicit Segmentation
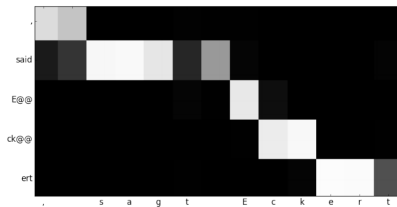
Chung, Cho and Bengio, ACL'16



"said Eckert"

# Neural Machine Translation with Finer Tokens

We are still concerned,

- ‣ Data sparsity problem will last!
    - ‣ but neural nets will less suffer from this issue
      (Bengio et al.,2003)

- ‣ Consequences of increased sequence length!
    - ‣ Capturing long-term dependencies
    - ‣ Will be harder to train
      (but wait we have GRU, LSTM and Attention)
    - ‣ Speed loss, 2-3 times slower

but ...

- ‣ No need to worry about segmentation,
- ‣ Open vocabularies, saves us giant matrices or tricks
- ‣ Naturally embeds multiple languages (**Lee et al.'16**)

---

# Fully Character-Level Multilingual NMT
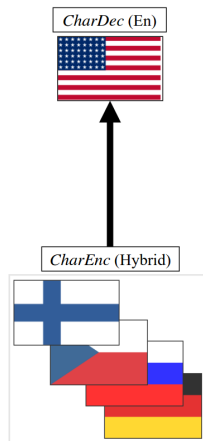
Jason Lee and Kyunghyun Cho, 2016 (in preparation)

Model details,

- RNNSearch model
- Source-Target character level
- CNN+RNN encoder
- *Bi-scale* decoder
- $\{Fi, De, Cs, Ru\} \rightarrow En$

Training,

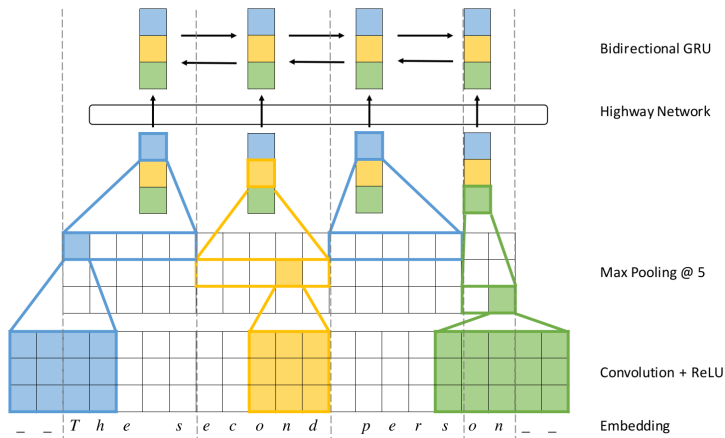- Mix mini-batches
- Use bi-text only



CharDec (En)

CharEnc (Hybrid)

# Fully Character-Level Multilingual NMT

Jason Lee and Kyunghyun Cho, 2016 (in preparation)

Hybrid Character Encoder,

# Fully Character-Level Multilingual NMT

Jason Lee and Kyunghyun Cho, 2016 (in preparation)

Preliminary Results, comparison with *BPE → Char*

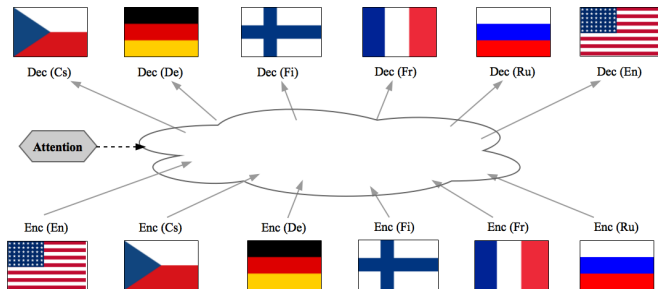| | Model | | | | Valid | Test-1 | Test-2 |
|---|---|---|---|---|---|---|---|
| | bpe2char | | char2char | | | | |
| | single | multi | single | multi | | | |
| De-En | ✓ | | | | 25.64 | 24.59 | 25.27 |
| | | | ✓ | | **26.03** | **25.80** | **25.77** |
| | | ✓ | | | 24.28 | 23.43 | 24.11 |
| | | | | ✓ | 25.45 | 24.27 | 25.06 |
| Cs-En | ✓ | | | | 22.83 | 23.51 | 22.46 |
| | | | ✓ | | 22.85 | 23.38 | 22.03 |
| | | ✓ | | | 22.76 | 23.46 | 21.86 |
| | | | | ✓ | **24.16** | **24.77** | **22.72** |
| Fi-En | ✓ | | | | 14.54 | 13.98 | - |
| | | | ✓ | | 14.18 | 13.10 | - |
| | | ✓ | | | 14.37 | 13.71 | - |
| | | | | ✓ | **15.85** | **15.80** | - |
| Ru-En | ✓ | | | | 21.68 | 26.21 | **22.83** |
| | | | ✓ | | 21.75 | **26.80** | 22.73 |
| | | ✓ | | | 20.91 | 24.59 | 21.93 |
| | | | | ✓ | **22.04** | 25.64 | 22.68 |

# Spoiler - Big Time!

# Spoiler - Big Time!

> ... a better single-pair translation system has never been **the** goal of neural MT ...
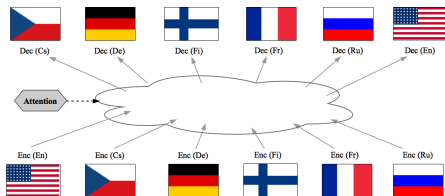>
> *Kyunghyun Cho*

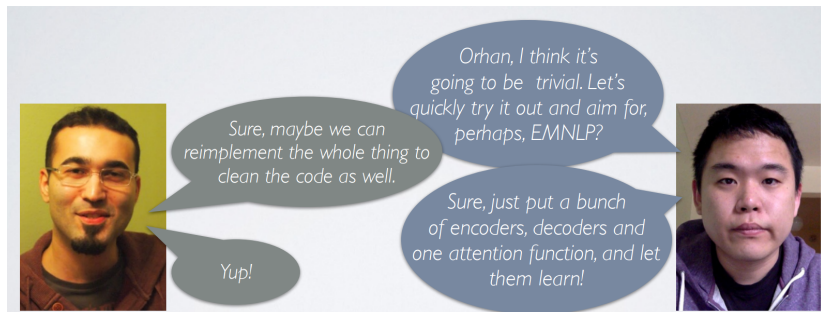# Multi-way, Multilingual Seq2Seq with Attention

# Potential Benefits

1. Positive language transfer across many language pairs/directions
   ‣ Solution to low/zero-resource machine translation
2. # of parameters grows linearly w.r.t. the # of languages
   ‣ as opposed to the quadratic explosion when training many single-pair models.
3. Multi-source translation without requiring any multi-way parallel text
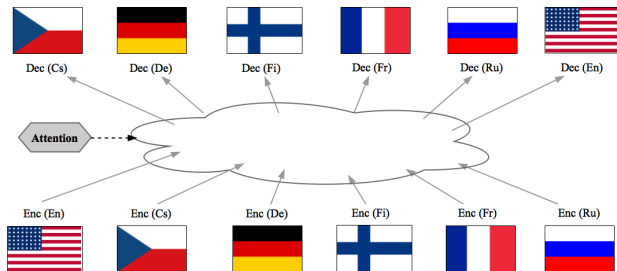   ‣ inspired by but contrary to Zoph & Knight (2016)

‣ Super fun to work on!

Multi-way, Multilingual Neural MT

1. Many-to-many translation
2. One shared attention mechanism
3. No need for multi-way parallel text
4. Scalable in terms of # languages, # sentences
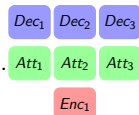5. Extendible to multiple modalities

# Multilingual (Multi-task) Neural Machine Translation
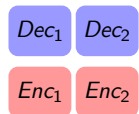
Recent work (chronologically),

- One-to-many (Dong et al., 2015)
  - Each decoder has it's own attention mechanism.
  - Small scale experiments (Europarl).
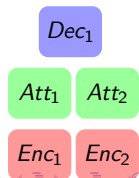  - No support for multiple source sentences.

| $Dec_1$ | $Dec_2$ | $Dec_3$ |
|---|---|---|
| $Att_1$ | $Att_2$ | $Att_3$ |
| $Enc_1$ | | |

- Without attention (Luong et al., 2015)
  - Focus on multi-task learning.
  - No attention, single vector space is shared.
  - Multilinguality is not considered in depth (En↔De).

| $Dec_1$ | $Dec_2$ |
|---|---|
| $Enc_1$ | $Enc_2$ |

- Many-to-one (Zoph and Knight, 2016)
  - Separate attention for each encoder.
  - Necessitates multi-text.
  - Small scale experiments (WMT'15 subset).

| $Dec_1$ | |
|---|---|
| $Att_1$ | $Att_2$ |
| $Enc_1$ | $Enc_2$ |

# Multi-way NMT - Overview

Firat, Cho and Bengio, 2016

| | | Dir | Fr (39m) | | Cs (12m) | | De (4.2m) | | Ru (2.3m) | | Fi (2m) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | → En | En → | → En | En → | → En | En → | → En | En → | → En | En → |
| (a) BLEU | Dev | Single | 27.22 | 26.91 | 21.24 | 15.9 | 24.13 | 20.49 | 21.04 | 18.06 | 13.15 | 9.59 |
| | | Multi | 26.09 | 25.04 | 21.23 | 14.42 | 23.66 | 19.17 | 21.48 | 17.89 | 12.97 | 8.92 |
| | Test | Single | 27.94 | **29.7** | 20.32 | **13.84** | 24 | **21.75** | 22.44 | **19.54** | 12.24 | **9.23** |
| | | Multi | **28.06** | 27.88 | **20.57** | 13.29 | **24.20** | 20.59 | **23.44** | 19.39 | **12.61** | 8.98 |
| (b) LL | Dev | Single | -50.53 | -53.38 | -60.69 | -69.56 | -54.76 | -61.21 | -60.19 | -65.81 | -88.44 | -91.75 |
| | | Multi | -50.6 | -56.55 | -54.46 | -70.76 | -54.14 | -62.34 | -54.09 | -63.75 | -74.84 | -88.02 |
| | Test | Single | -43.34 | **-45.07** | -60.03 | **-64.34** | -57.81 | **-59.55** | -60.65 | -60.29 | -88.66 | -94.23 |
| | | Multi | **-42.22** | -46.29 | **-54.66** | -64.80 | **-53.85** | -60.23 | **-54.49** | **-58.63** | **-71.26** | **-88.09** |

# Warren Weaver-"Translation", 1949

Tall towers analogy:

- Do **NOT** model the individual behaviour of a car,
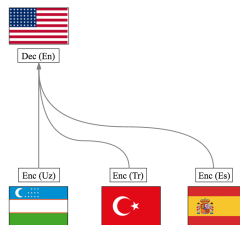- **Model how the highway works!**

# A lot of (if not all) interesting, related questions remain ...

1. What is it good for, other than parameter saving?
2. What if a source sentence is given in multiple languages?
3. What happens with language pairs that are not included during training?
4. How are we going to introduce additional modalities?

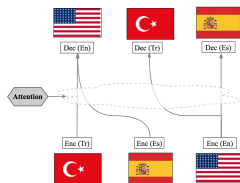# ML-NMT is good for Low-Resource (Firat et al., 2016b)

$[Uz \to En]$

| Added Pairs | Test | |
|---|---|---|
| | Single | Ensemble |
| Uz→En | 42.56 (6.45) | 38.56 (8.81)* |
| + Tr→En | 36.79 (9.34) | 34.49 (11.69)* |
| + Tr→En + Es→En | **35.39 (10.34)** | 33.20 (12.33)* |
| + Tr→En + En→Uz + En→Tr | 36.28 (9.41) | 33.65 (11.30)* |
| MLNMT Ensemble | **31.77 (12.99)**† | |
| Conventional SMT | 32.38 (9.37) | |



$[Tr \to En]$

| Added Pairs | Test | |
|---|---|---|
| | Single | Ensemble |
| Tr→En | 28.58 (17.28) | 24.27 (20.83)* |
| + Es→En | 27.49 (17.75) | 23.94 (20.89)* |
| + Es→En + Fr→En | 26.77 (18.13) | 24.00 (20.90)* |
| + Es→En + En→Tr + En→Es | **26.30 (18.66)** | 24.28 (20.23)* |
| MLNMT Ensemble | **21.78 (22.56)**† | |
| Conventional SMT | 23.42 (18.00) | |

# Where does the improvement coming from?

1. Encoder is shared across *one-to-many* pairs
2. Decoder is shared across *many-to-one* pairs
3. The soft-alignment mechanism is shared across all pairs
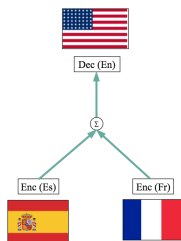
# What can we do more?

1. Share an encoder for multiple, similar source languages
2. Share a decoder for multiple, similar target languages
3. Perhaps, one recurrent net to rule both source and target languages..? (**Lee et al.'16**)

# What if a source sentence is given in multiple languages?
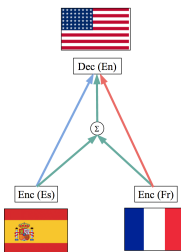
Multi-source neural machine translation

Multi-text given during **training**

1. Train the model for $p(y|x_1, x_2)$
   (Zoph and Knight, 2016)
2. May need to device a merger operation



Multi-text given during **test**

1. Two translation strategies
   (Firat et al., 2016c)
   1.1 Late averaging $p(y|x_1) + p(y|x_2)$
   1.2 Early averaging $p(y|x_1, x_2)$
2. Use existing shared attention for merger operation
   - Simply take the *mean* of representations

# What if a source sentence is given in multiple languages?

Multi-source neural machine translation

Multi-text given during **test**

1. Two translation strategies (Firat et al., 2016c)
    1.1 Late averaging $p(y|z_1) + p(y|x_2)$
    1.2 Early averaging $p(y|z_1, z_2)$
2. Use existing shared attention for merger operation
    - Simply take the *mean* of representations

Single-source translation:

| | Src | Trgt | Multi | | Single | |
|---|---|---|---|---|---|---|
| | | | Dev | Test | Dev | Test |
| (a) | Es | En | 30.73 | 28.32 | 29.74 | 27.48 |
| (b) | Fr | En | 26.93 | 27.93 | 26.00 | 27.21 |
| (c) | En | Es | 30.63 | 28.41 | 31.31 | 28.90 |
| (d) | En | Fr | 22.68 | 23.41 | 22.80 | 24.05 |

Multi-source translation:

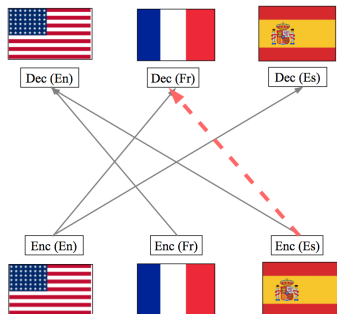| | | Multi | | Single | |
|---|---|---|---|---|---|
| | | Dev | Test | Dev | Test |
| (a) | Early | 31.89 | 31.35 | – | – |
| (b) | Late | 32.04 | 31.57 | 32.00 | 31.46 |
| (c) | E+L | 32.61 | 31.88 | – | – |

# What have we learned from this?

1. No need for multi-way parallel corpora!
2. Because, training with multiple language pairs has encouraged the model to find a common context vector space.
3. Allows us to use simple arithmetic operations in a hopefully flattened manifold.

# What more should we do?

1. Finetuning with multi-way parallel corpus helps, but how far can we go?
2. Larger-scale experiments with more source language pairs.

# Towards Zero-Resource Language Translation

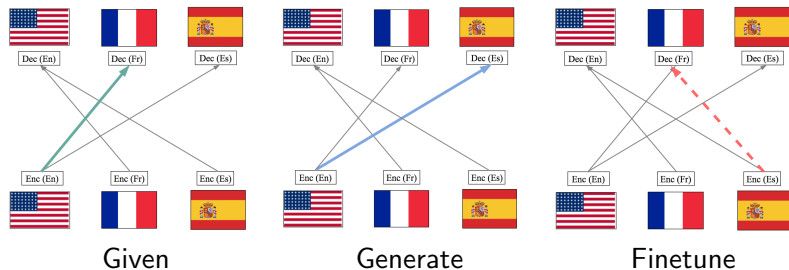Can we translate between a pair without any direct resource?



| | | Pivot | Many-to-1 | | Dev | Test |
|---|---|---|---|---|---|---|
| (a) | ‖ | | | ‖ | < 1 | < 1 |
| (b) | ‖ | √ | | ‖ | 20.64 | 20.4 |

Unfortunately no! Instead $Es \rightarrow En \rightarrow Fr$ is promising

# Towards Zero-Resource Language Translation

- *Es → Fr*: perhaps we can *generate* a pseudo-parallel corpus (Sennrich et al.,2016)
- Still *zero* direct resource



Given        Generate        Finetune

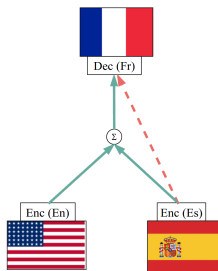# Towards Zero-Resource Language Translation

Firat et al., 2016c

- Generate a pseudo-source ($Es$) given $En - Fr$ parallel corpus
- Finetune $Es - Fr$ using pseudo-parallel corpus

|      | Pseudo Parallel Corpus | | | | True Parallel Corpus | | | |
|------|------|------|------|------|------|------|------|------|
|      | 1k | 10k | 100k | 1m | 1k | 10k | 100k | 1m |
| Dev  | – | – | – | – | – | – | 11.25 | 21.32 |
| Test | – | – | – | – | – | – | 10.43 | 20.35 |
|      | Dev: 20.64, Test: 20.4 | | | | – | | | |
| Dev  | 0.28 | 10.16 | 15.61 | 17.59 | 0.1 | 8.45 | 16.2 | 20.59 |
| Test | 0.47 | 10.14 | 15.41 | 17.61 | 0.12 | 8.18 | 15.8 | 19.97 |

# Towards Zero-Resource Language Translation

Firat et al., 2016c

| Pivot | Many-to-1 | | Pseudo Parallel Corpus | | | | True Parallel Corpus | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | 1k | 10k | 100k | 1m | 1k | 10k | 100k | 1m |
| √ | No Finetuning | | Dev: 20.64, Test: 20.4 | | | | – | | | |
| | | Dev | 0.28 | 10.16 | 15.61 | 17.59 | 0.1 | 8.45 | 16.2 | 20.59 |
| | | Test | 0.47 | 10.14 | 15.41 | 17.61 | 0.12 | 8.18 | 15.8 | 19.97 |
| √ | Early | Dev | 19.42 | 21.08 | 21.7 | 21.81 | 8.89 | 16.89 | 20.77 | 22.08 |
| | | Test | 19.43 | 20.72 | 21.23 | 21.46 | 9.77 | 16.61 | 20.40 | 21.7 |
| √ | Early+ Late | Dev | 20.89 | 20.93 | 21.35 | 21.33 | 14.86 | 18.28 | 20.31 | 21.33 |
| | | Test | 20.5 | 20.71 | 21.06 | 21.19 | 15.42 | 17.95 | 20.16 | 20.9 |



▸ **Multi-source comes to aid!**

  ▸ Teacher (multi-source)
  ▸ Student (zero-resource)

**What have we learned from this?**

1. Don't necessarily need a direct resource.
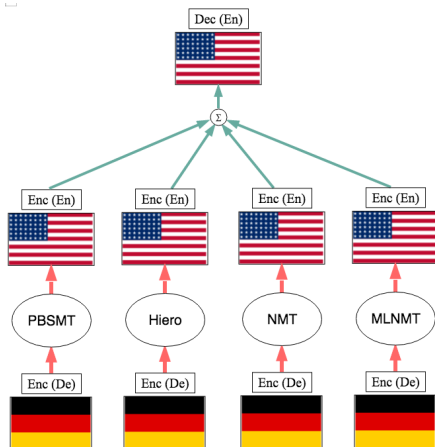2. Multilingual NMT naturally embeds multiple translation strategies.

**What more should we do?**

1. Active learning can bring additional gains.
2. Perhaps, simply more data will constrain the attention mechanism to work with zero-resource pairs automatically.

# Trainable Neural System Combination

Firat, Freitag and Cho - ongoing
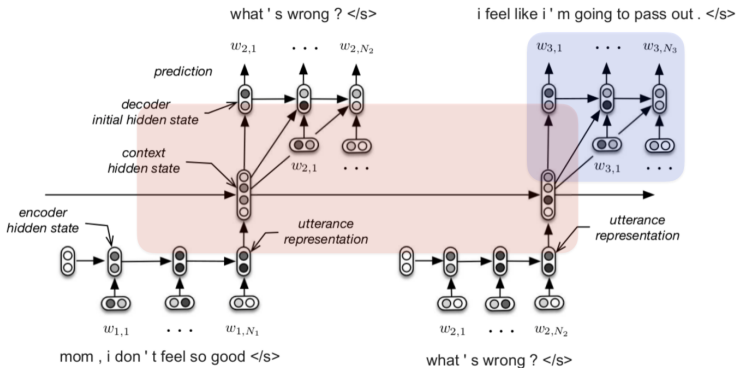
How to combine traditional SMT models with NMT models?

# Larger Context NMT - Going beyond sentences!
Hierarchical Recurrent Encoder-Decoder (Serban et al., 2015, Sordoni et al., 2015)

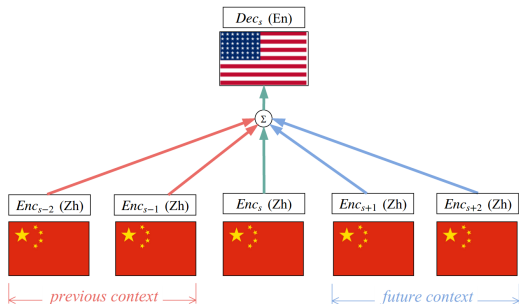For dialogue modelling, capture previous context.

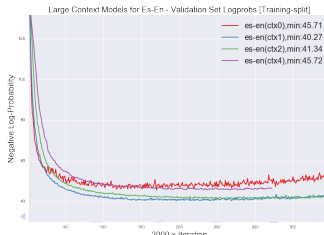**Utterance-level RNN** + **Dialogue-level RNN**

# Larger Context NMT - Going beyond sentences!

Firat, Lauly and Cho - ongoing

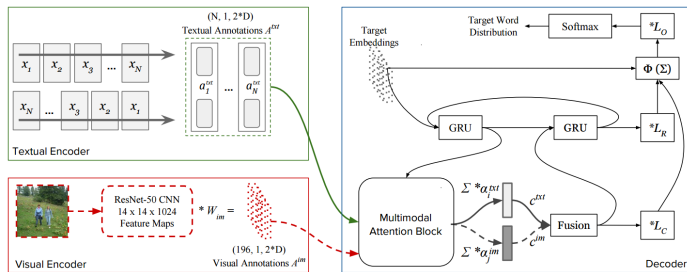Extend context to multiple past/future sentences in the document.



- ▸ New released UN Corpus (Es→En and Zh→En)
- ▸ **Mean** as the merger op
- ▸ Consider context window of size $0, 1, 2$ and $4$

# What about multi-modal NMT?

"Multi-modal Attention for Neural Machine Translation"
Caglayan, Barrault and Bougares, 2016 (actually, just yesterday)



| | | Attention Type | | | Validation Scores | |
|---|---|---|---|---|---|---|
| Model | Fusion | Modality | Decoder | METEOR | BLEU | CIDEr-D |
| NMT | - | - | - | 34.24 (35.59) | 18.64 (21.62) | 58.57 (67.93) |
| IMGTXT | - | - | - | 26.80 | 11.16 | 31.28 |
| MNMT1 | SUM | IND | IND | 33.23 (35.42) | 18.30 (21.24) | 55.45 (65.03) |
| MNMT2 | SUM | IND | DEP | 34.17 (35.48) | 17.70 (20.70) | 53.78 (61.76) |
| MNMT3 | SUM | DEP | IND | 34.38 (35.55) | 18.42 (20.94) | 55.81 (63.37) |
| MNMT4 | SUM | DEP | DEP | 33.67 (34.57) | 17.83 (20.30) | 52.68 (59.63) |
| MNMT5 | CONCAT | IND | IND | 33.31 (34.98) | 17.50 (20.60) | 53.57 (61.46) |
| MNMT6 | CONCAT | IND | DEP | **35.23** (36.79) | 19.30 (22.45) | 60.62 (69.96) |
| MNMT7 | CONCAT | DEP | IND | **35.11** (**37.13**) | **19.72** (**23.24**) | **61.04** (**72.16**) |
| MNMT8 | CONCAT | DEP | DEP | 34.80 (**36.98**) | 19.55 (22.78) | 60.20 (70.20) |

# How far we can extend the existing approaches?

Bigger models, complicated architectures!

RNNs can express/approximate a set of Turing machines,
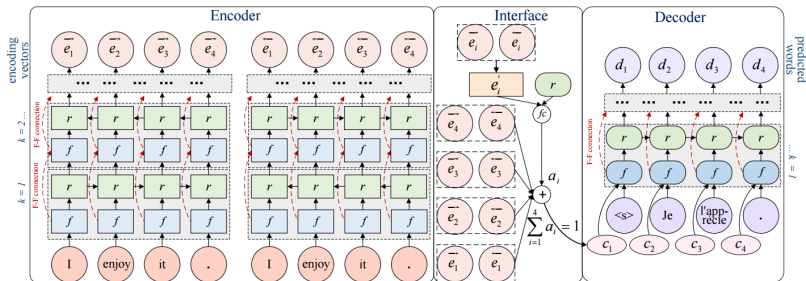
BUT[*] ...

# expressivity ≠ learnability

# How far we can extend the existing approaches?

Fast-Forward Connections for NMT, (Zhou et al., 2016)

Bigger models are harder to train!

- ‣ Deep topology for recurrent networks (16 layers)
- ‣ Performance boost ($+6.2$ BLEU points)
- ‣ Fast-forward connections for gradient flow

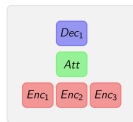# How far we can extend the existing approaches?
## Multi-way Multilingual NMT

Bigger models are harder to train and behave differently!

- ‣ Scheduling the learning process
- ‣ Preventing the unlearning (catastrophic forgetting)
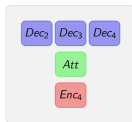- ‣ Layer Normalization (Kiros et al.,2016)

# What are we optimizing?

Explorations on the right objective to be optimized NLL:

- ‣ MERT (Och,2003), MRT for NMT (Shen et al.,2016)
- ‣ Scheduled Sampling (Bengio et al.,2015), Sequence Level Training (Ranzato et al.,2015), Task Loss Estimation (Bahdanau et al.,2015)
- ‣ Actor-Critic (Bahdanau et al.,2016), Reward Augmented ML (Norouzi et al.,2016)
- ‣ Seq2Seq as Beam-Search optimization (Wiseman and Rush, 2016)

*New territory seems to be using new error signals!*

# What Lies Ahead?

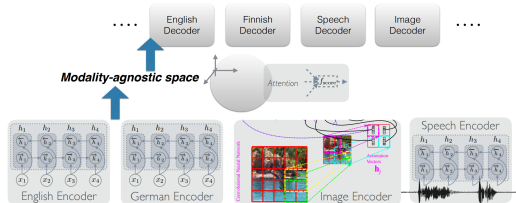Perhaps, we've only scratched the surface!

- ‣ Language barrier, surpassing human level quality.

Revisiting the new territory:

## **Character-level Larger-Context Multilingual** Neural Machine Translation

using,

- ‣ Multiple modalities
- ‣ Better error signals
- ‣ and better GPUs 😎

# One last thing!

Let's remember the game we were playing once more,

# Thank you!

Thanks to: TUBITAK, NSERC, Samsung, IBM, Calcul Quebec, Compute Canada, The Canada Research Chair, CIFAR, NVIDIA, Facebook and Google