•One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: *'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.'*

Warren Weaver to Norbert Wiener, March, 1947

# ORDER YOUR
# KAWHE/COFFEE
## IN MĀORI

**He mōwai māku**  I'll have a flat white
**He pango poto māku**  I'll have a short black
**He pango roa māku**  I'll have a long black
**He rate pīni māku**  I'll have a soy latte
**He kaputino māku**  I'll have a cappuccino
**He rate māku**  I'll have a latte
**He tiakarete wera māku**  I'll have a hot chocolate

**Rahi** Size

**(S) Paku**
**(M) Waenga**
**(L) Nui**

**Kei te pēhea koe?**
How's it going?

**Anei taku kapu mau tonu**
Here is my reusable cup

**Hei kawe atu**
To take away

**Ki konei**
To have here

McCafé

# CLASSIC SOUPS

| | | | Sm. | Lg. |
|---|---|---|---|---|
| 清燉雞湯 | 57. | House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot) | 1.50 | 2.75 |
| 雞飯湯 | 58. | Chicken Rice Soup | 1.85 | 3.25 |
| 雞麵湯 | 59. | Chicken Noodle Soup | 1.85 | 3.25 |
| 廣東雲吞 | 60. | Cantonese Wonton Soup | 1.50 | 2.75 |
| 蕃茄蛋湯 | 61. | Tomato Clear Egg Drop Soup | 1.65 | 2.95 |
| 雲吞湯 | 62. | Regular Wonton Soup | 1.10 | 2.10 |
| 酸辣湯 | 63. 🍲 | Hot & Sour Soup | 1.10 | 2.10 |
| 蛋花湯 | 64. | Egg Drop Soup | 1.10 | 2.10 |
| 雲蛋湯 | 65. | Egg Drop Wonton Mix | 1.10 | 2.10 |
| 豆腐菜湯 | 66. | Tofu Vegetable Soup | NA | 3.50 |
| 雞玉米湯 | 67. | Chicken Corn Cream Soup | NA | 3.50 |
| 蟹肉玉米湯 | 68. | Crab Meat Corn Cream Soup | NA | 3.50 |
| 海鮮湯 | 69. | Seafood Soup | NA | 3.50 |

Egyptian

Greek

- We want a model of $p(\mathbf{e} \mid \mathbf{f})$

**Possible English translation**

**Confusing foreign sentence**

# NOISY CHANNEL MT

p(f|e)

"English"

"Foreign"

p(e)

e

f

channel

*decode*

$$\hat{e} = \arg\max_e p(e|f)$$

$$= \arg\max_e \frac{p(e) \times p(f|e)}{p(f)}$$

$$= \arg\max_e \boxed{p(e)} \times \boxed{p(f|e)}$$

"Language Model"   "Translation Model"

# NOISY CHANNEL DIVISION OF LABOR

- Language model – $p(e)$

  - is the translation fluent, grammatical, and idiomatic?
  - use any model of $p(e)$ – typically an $n$-gram model

- Translation model – $p(f|e)$

  - "reverse" translation probability
  - ensures adequacy of translation

# LANGUAGE MODEL FAILURE



*My legal name is Alexander Perchov.*

# LANGUAGE MODEL FAILURE



*My legal name is Alexander Perchov. But all of my many friends dub me Alex, because that is a more flaccid-to-utter version of my legal name. Mother dubs me Alexi-stop-spleening-me!, because I am always spleening her.*

# LANGUAGE MODEL FAILURE



*My legal name is Alexander Perchov. But all of my many friends dub me Alex, because that is a more flaccid-to-utter version of my legal name. Mother dubs me Alexi-stop-spleening-me!, because I am always spleening her. If you want to know why I am always spleening her, it is because I am always elsewhere with friends, and disseminating so much currency, and performing so many things that can spleen a mother.*

# TRANSLATION MODEL

$p(f|e)$ gives the channel probability – the probability of translating an English sentence into a foreign sentence

$f$ = je voudrais un peu de frommage

$e_1$ = I would like some cheese

$e_2$ = I would like a little of cheese

$e_3$ = There is no train to Barcelona

$p(f|e)$

0.4

0.5

>0.00001

- How do we parameterize *p(f|e)*?

$$p(f|e) = \frac{count(f, e)}{count(e)} \; ?$$

- There are a lot of sentences: this won't generalize to new inputs

# LEXICAL TRANSLATION

How do we translate a word? Look it up in a dictionary!

*Haus: house, home, shell, household*

Multiple translations

Different word senses, different registers, different inflections

*house, home* are common

*shell* is specialized (the Haus of a snail is its shell)

# HOW COMMON IS EACH?

| Translation | Count |
|---|---|
| house | 5000 |
| home | 2000 |
| shell | 100 |
| household | 80 |

# MLE

$$\hat{p}_{\mathrm{MLE}}(e \mid \mathbf{Haus}) = \begin{cases} 0.696 & \text{if } e = \texttt{house} \\ 0.279 & \text{if } e = \texttt{home} \\ 0.014 & \text{if } e = \texttt{shell} \\ 0.011 & \text{if } e = \texttt{household} \\ 0 & \text{otherwise} \end{cases}$$

# LEXICAL TRANSLATION

- Goal: a model *p(e|f,m)*

- where **e** and **f** are complete English and Foreign sentences

$$\mathbf{e} = \langle e_1, e_2, \ldots, e_m \rangle \quad \mathbf{f} = \langle f_1, f_2, \ldots, f_n \rangle$$

# LEXICAL TRANSLATION

Goal: a model $p(e|f,m)$

where **e** and **f** are complete English and Foreign sentences

Lexical translation makes the following *assumptions*:

1. Each word $e_i$ in **e** is generated from exactly one word in **f**

2. Thus, we have a latent *alignment* $a_i$ that indicates which word $e_i$ "came from." Specifically it came from $f_{a_i}$.

3. Given the alignments **a**, translation decisions are conditionally independent of each other and depend *only* on the aligned source word $f_{a_i}$.

- Putting our assumptions together, we have:

$$p(\mathbf{e} \mid \mathbf{f}, m) = \sum_{\mathbf{a} \in [0,n]^m} p(\mathbf{a} \mid \mathbf{f}, m) \times \prod_{i=1}^{m} p(e_i \mid f_{a_i})$$

p(Alignment)     p(Translation | Alignment)

$$p(\mathbf{a} \mid \mathbf{f}, m)$$

- Most of the action for the first 10 years of MT was here. Words weren't the problem. Word *order* was hard.

# ALIGNMENT

- Alignments can be visualized by drawing links between two sentences, and they are represented as vectors of positions:

$$\mathbf{a} = (1, 2, 3, 4)^{\top}$$

- Words may be reordered during translation



$$\mathbf{a} = (3, 4, 2, 1)^{\top}$$

- A source word may not be translated at all



$$\mathbf{a} = (2, 3, 4)^{\top}$$

- Words may be inserted during translation

- E.g. English just does not have an equivalent

- But these words must be explained – we typically assume every source sentence contains a NULL token



$$\mathbf{a} = (1, 2, 3, 0, 4)^\top$$

- A source word may translate into **more than one** target word



$$\mathbf{a} = (1, 2, 3, 4, 4)^\top$$

- More than one source word may **not** translate as a unit in lexical translation



$$\mathbf{a} = ??? \qquad \mathbf{a} = (1, 2, (3, 4)^\top)^\top \ ?$$

# IBM MODEL 1

Simplest possible lexical translation model

Additional assumptions:

  The $m$ alignment decisions are independent

  The alignment distribution for each $\mathbf{a}_i$ is uniform over all source words and NULL

$$\text{for each } i \in [1, 2, \ldots, m]$$
$$a_i \sim \text{Uniform}(0, 1, 2, \ldots, n)$$
$$e_i \sim \text{Categorical}(\boldsymbol{\theta}_{f_{a_i}})$$

Language model says: ☺

Language model says: ☹

# LEARNING LEXICAL TRANSLATION MODELS

How do we learn the parameters *p(e|f)*?

"Chicken and egg" problem:
  If we had the alignments, we could estimate the translation probabilities (MLE estimation)
  If we had the translation probabilities we could find the most likely alignments (greedy)

# IBM 1 - GENERATIVE STORY

We start with an English Sentence $\mathbf{e} = e_1 e_2 \ldots e_n$

1. Choose the length of the Spanish sentence $m$, with uniform probability $\epsilon = \dfrac{1}{M}$, where $M$ is the maximum allowed length of any Spanish sentence in the corpus.
2. Generate an alignment $a_1, \ldots, a_m$ again with uniform probability.
3. Generate Spanish words $f_1, \ldots, f_m$ each with probability $t(f_j \mid e_{a_j})$ or $t(f_j \mid \text{NULL})$



How can we estimate the $t(f \mid e)$ parameters?

# EM ALGORITHM

Pick some random (or uniform) starting parameters

Repeat until bored (~5 iterations for lexical translation models):

1. Using the current parameters, compute "expected" alignments $p(\mathbf{a}_i|\mathbf{e}, \mathbf{f})$ for every target word token in the training data

2. Keep track of the expected number of times $f$ translates into $e$ throughout the whole corpus

3. Keep track of the number of times $f$ is used in the source of any translation

4. Use these estimates in the standard MLE equation to get a better set of parameters

... la maison ... la maison blue ... la fleur ...

... the house ... the blue house ... the flower ...

- Initial step: all alignments equally likely

- Model learns that, e.g., la is often aligned with the

... la maison ... la maison blue ... la fleur ...

... the house ... the blue house ... the flower ...

- After one iteration

- Alignments, e.g., between la and the are more likely

... la maison ... la maison bleu ... la fleur ...

... the house ... the blue house ... the flower ...

- After another iteration

- It becomes apparent that alignments, e.g., between fleur and flower are more likely (pigeon hole principle)

... la maison ... la maison bleu ... la fleur ...

... the house ... the blue house ... the flower ...

$$p(\text{la}|\text{the}) = 0.453$$
$$p(\text{le}|\text{the}) = 0.334$$
$$p(\text{maison}|\text{house}) = 0.876$$
$$p(\text{bleu}|\text{blue}) = 0.563$$

...

• Parameter estimation from the aligned corpus

1. Initialize $t(\,\cdot\,|\,e)$ to uniform: $t(f\,|\,e) = \dfrac{1}{|V_f|}$ where $V_f$ is the Spanish vocabulary, and $e$ is any English word or NULL.

2. E-step: Calculate the *expected* number of times that word $e$ is translated as $f$.
   For each $i, j$ the transition that generates $f_j$ from $e_i$ "competes" with the transitions that generate $f_j$ from the other English words (or NULL). So we update our expected counts $c(f, e)$ as follows:

$$c(f_j, e_i) \leftarrow c(f_j, e_i) + \frac{t(f_j\,|\,e_i)}{t(f_j\,|\,\text{NULL}) + \Sigma_{i'} t(f_j\,|\,e_{i'})} \qquad c(f_j, \text{NULL}) \leftarrow c(f_j, \text{NULL}) + \frac{t(f_j\,|\,\text{NULL})}{t(f_j\,|\,\text{NULL}) + \Sigma_{i'} t(f_j\,|\,e_{i'})}$$

3. M-step: Estimate the model's parameters based on the expected counts.
   Let $t(f\,|\,e) \leftarrow \dfrac{c(f, e)}{\Sigma_f c(f, e)}$ where $e$ is any English word or NULL.

4. Go to step 2.

# CONVERGENCE

das   Haus      das   Buch      ein   Buch

the   house     the   book      a     book

| $e$ | $f$ | initial | 1st it. | 2nd it. | 3rd it. | ... | final |
|-----|-----|---------|---------|---------|---------|-----|-------|
| the | das | 0.25 | 0.5 | 0.6364 | 0.7479 | ... | 1 |
| book | das | 0.25 | 0.25 | 0.1818 | 0.1208 | ... | 0 |
| house | das | 0.25 | 0.25 | 0.1818 | 0.1313 | ... | 0 |
| the | buch | 0.25 | 0.25 | 0.1818 | 0.1208 | ... | 0 |
| book | buch | 0.25 | 0.5 | 0.6364 | 0.7479 | ... | 1 |
| a | buch | 0.25 | 0.25 | 0.1818 | 0.1313 | ... | 0 |
| book | ein | 0.25 | 0.5 | 0.4286 | 0.3466 | ... | 0 |
| a | ein | 0.25 | 0.5 | 0.5714 | 0.6534 | ... | 1 |
| the | haus | 0.25 | 0.5 | 0.4286 | 0.3466 | ... | 0 |
| house | haus | 0.25 | 0.5 | 0.5714 | 0.6534 | ... | 1 |

Phrase-based MT:

Allow multiple words to translate as chunks (including many-to-one)

Introduce another latent variable, the source *segmentation*



Adapted from Koehn (2006)

# EXTENSIONS

Alignment Priors:

Instead of assuming the alignment decisions are uniform, impose (or learn) a prior over alignment grids:



Chahuneau et al. (2013)

# EXTENSIONS

Syntactic structure

Rules of the form:

X之一 →　 one of the X



Chiang (2005), Galley et al. (2006)

# EVALUATION

How do we evaluate translation systems' output?

Central idea: "The closer a machine translation is to a professional human translation, the better it is."

Most commonly used metric is called BLEU

# BLEU: AN EXAMPLE

Candidate 1: *It is a guide to action which ensures that the military always* obey *the commands of the party.*

Reference 1: *It is a guide to action that ensures that the military will forever heed Party commands.*

Reference 2: *It is the guiding principle which guarantees the military forces always being under the command of the Party.*

Reference 3: *It is the practical guide for the army always to heed directions of the party.*

Unigram Precision : 17/18

# ISSUE OF N-GRAM PRECISION

What if some words are over-generated?
  e.g. "the"

An extreme example

  Candidate: *the the the the the the the.*

  Reference 1: *The cat is on the mat.*

  Reference 2: *There is a cat on the mat.*

N-gram Precision: 7/7

**Solution: reference word should be exhausted after it is matched.**

# ISSUE OF N-GRAM PRECISION

What if some words are just dropped?

Another extreme example

Candidate: *the*.

Reference 1: *My mom likes the blue flowers.*

Reference 2: *My mother prefers the blue flowers.*

N-gram Precision: 1/1

**Solution: add a penalty if the candidate is too short.**

# BLEU

Geometric Average

$$\text{BLEU} = (p_1 \cdot p_2 \cdot p_3 \cdot p_4)^{\frac{1}{4}} \max(1,\ e^{1-\frac{r}{c}})$$

Clipped N-gram precisions for N=1, 2, 3, 4

Brevity Penalty

Ranges from 0.0 to 1.0, but usually shown multiplied by 100

An increase of +1.0 BLEU is usually a conference paper

MT systems usually score in the 10s to 30s (40-50s?)

Human translators usually score in the 70s and 80s

# A SHORT SEGUE

Word- and phrase-based ("symbolic") models were cutting edge for decades (up until ~2014)

Such models are still the most widely used in commercial applications

Since 2014 most research on MT has focused on **neural** models

# FULLY NEURAL TRANSLATION

Fully end-to-end RNN-based translation model

Encode the source sentence using one RNN

Generate the target sentence one word at a time using another RNN



Encoder

Decoder

Sutskever et al. (2014)

# ATTENTIONAL MODEL

The encoder-decoder model struggles with long sentences

An RNN is trying to compress an arbitrarily long sentence into a finite-length worth vector

What if we only look at one (or a few) source words when we generate each output word?

Bahdanau et al. (2014)

# THE INTUITION

うち　　の大きな　黒い　犬　が可哀想な　郵便屋　に 噛み ついた　　　。

Our　large　　black　dog　bit　　the　　poor mailman　　.

Bahdanau et al. (2014)

Encoder

Decoder

I    am    astudent    </s>

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Encoder

Decoder

I   am   astudent   </s>

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Encoder

Decoder

softmax

I     am     astudent     </s>

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Context Vector

Encoder

Decoder

I   am   astudent   </s>

Bahdanau et al. (2014)

# THE ATTENTION MODEL

Attention
Model

Context Vector

je

Encoder

I    am    astudent    </s>

Decoder

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Context Vector

Encoder

Decoder

je

je

I    am    astudent    </s>

je

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Encoder

Decoder

je

je

I    am    astudent    </s>

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Context Vector

Encoder

Decoder

je    suis

I    am    astudent    </s>

je

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Context Vector

Encoder

Decoder

je     suis

I     am     astudent     </s>

je     suis

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention Model

Encoder

Context Vector

Decoder

je    suis  étudiant

I    am    astudent    </s>

je    suis  étudiant

Bahdanau et al. (2014)

# THE ATTENTION MODEL



Attention
Model

Context Vector

Encoder

Decoder

je    suis  étudiant </s>

I    am    astudent  </s>

je    suis  étudiant

Bahdanau et al. (2014)

# CONVOLUTIONAL ENCODER-DECODER

Gehring et. al 2017

CNN:

encodes words within a fixed size window

Parallel computation

Shortest path to cover a wider range of words

RNN:

sequentially encode a sentence from left to right

Hard to parallelize

# THE TRANSFORMER

- Idea: Instead of using an RNN to encode the source sentence and the partial target sentence, use self-attention!

Standard RNN Encoder

Self Attention Encoder

word-in-context vector

raw word vector

I    am    astudent    </s>

I    am    astudent    </s>

Vaswani et al. (2017)

# THE TRANSFORMER

Attention Model

Context Vector

Encoder

Decoder

I    am    astudent    </s>

je    suis    étudiant </s>

je    suis    étudiant

Vaswani et al. (2017)

# VISUALIZATION OF ATTENTION WEIGHT

- Self-attention weight can detect long-term dependencies within a sentence, e.g., make … more difficult

Computation is easily parallelizable

Shorter path from each target word to each source word → stronger gradient signals

Empirically stronger translation performance

Empirically trains substantially faster than more serial models

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [17] | 23.75 | | | |
| Deep-Att + PosUnk [37] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [36] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [31] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [37] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [36] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.0** | $2.3 \cdot 10^{19}$ | |

# CURRENT RESEARCH DIRECTIONS ON NEURAL MT

- Incorporation syntax into Neural MT

- Handling of morphologically rich languages

- Optimizing translation quality (instead of corpus probability)

- Multilingual models

# CURRENT RESEARCH DIRECTIONS ON NEURAL MT

- Incorporation syntax into Neural MT

- Handling of morphologically rich languages

- Optimizing translation quality (instead of corpus probability)

- Multilingual models

## Exploring Massively Multilingual, Massive Neural Machine Translation

Friday, October 11, 2019

Posted by Ankur Bapna, Software Engineer and Orhan Firat, Research Scientist, Google Research

# CURRENT RESEARCH DIRECTIONS ON NEURAL MT

- Incorporation syntax into Neural MT

- Handling of morphologically rich languages

- Optimizing translation quality (instead of corpus probability)

- Multilingual models

## Achieving Human Parity on Automatic Chinese to English News Translation

Hany Hassan,* Anthony Aue, Chang Chen, Vishal Chowdhary, Jonathan Clark, Christian Federmann, Xuedong Huang, Marcin Junczys-Dowmunt, William Lewis, Mu Li, Shujie Liu, Tie-Yan Liu, Renqian Luo, Arul Menezes, Tao Qin, Frank Seide, Xu Tan, Fei Tian, Lijun Wu, Shuangzhi Wu, Yingce Xia, Dongdong Zhang, Zhirui Zhang, and Ming Zhou

Microsoft AI & Research

# CURRENT RESEARCH DIRECTIONS ON NEURAL MT

- Incorporation syntax into Neural MT

- Handling of morphologically rich languages

- Optimizing translation quality (instead of corpus probability)

- Multilingual models

- Document-level translation

## Has Machine Translation Achieved Human Parity? A Case for Document-level Evaluation

Samuel Läubli[1]    Rico Sennrich[1,2]    Martin Volk[1]

[1]Institute of Computational Linguistics, University of Zurich
{laeubli,volk}@cl.uzh.ch

# CURRENT RESEARCH DIRECTIONS ON NEURAL MT

- Incorporation syntax into Neural MT

- Handling of morphologically rich languages

- Optimizing translation quality (instead of corpus probability)

- Multilingual models

- Document-level translation

- Domain adaptation and robustness