# Machine Translation

How does Seq2seq behave
with multi-head attention?

# 1.
# Introduction

EUROPEAN
PARLIAMENT

Tao, C., Gao, S., Shang, M., Wu, W., Zhao, D., & Yan, R. (2018, July). Get The Point of My Utterance! Learning Towards Effective Responses with Multi-Head Attention Mechanism. In IJCAI (pp. 4418-4424).

# 1.
# Background

# Bahdanau Attention

- The OG attention
- Single head

## NEURAL MACHINE TRANSLATION
## BY JOINTLY LEARNING TO ALIGN AND TRANSLATE

**Dzmitry Bahdanau**
Jacobs University Bremen, Germany

**KyungHyun Cho**     **Yoshua Bengio**\*
Université de Montréal

### ABSTRACT

Neural machine translation is a recently proposed approach to machine transla-
tion. Unlike the traditional statistical machine translation, the neural machine
translation aims at building a single neural network that can be jointly tuned to
maximize the translation performance. The models proposed recently for neu-
ral machine translation often belong to a family of encoder–decoders and encode
a source sentence into a fixed-length vector from which a decoder generates a
translation. In this paper, we conjecture that the use of a fixed-length vector is a
bottleneck in improving the performance of this basic encoder–decoder architec-
ture, and propose to extend this by allowing a model to automatically (soft-)search
for parts of a source sentence that are relevant to predicting a target word, without
having to form these parts as a hard segment explicitly. With this new approach,
we achieve a translation performance comparable to the existing state-of-the-art
phrase-based system on the task of English-to-French translation. Furthermore,
qualitative analysis reveals that the (soft-)alignments found by the model agree
well with our intuition.

# Multi-Head Attention

- Commonly used for Transformers

- Each head using Scaled Dot-Product Attention

scholarly works.

## Attention Is All You Need

**Ashish Vaswani***
Google Brain
avaswani@google.com

**Noam Shazeer***
Google Brain
noam@google.com

**Niki Parmar***
Google Research
nikip@google.com

**Jakob Uszkoreit***
Google Research
usz@google.com

**Llion Jones***
Google Research
llion@google.com

**Aidan N. Gomez*** †
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser***
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin*** ‡
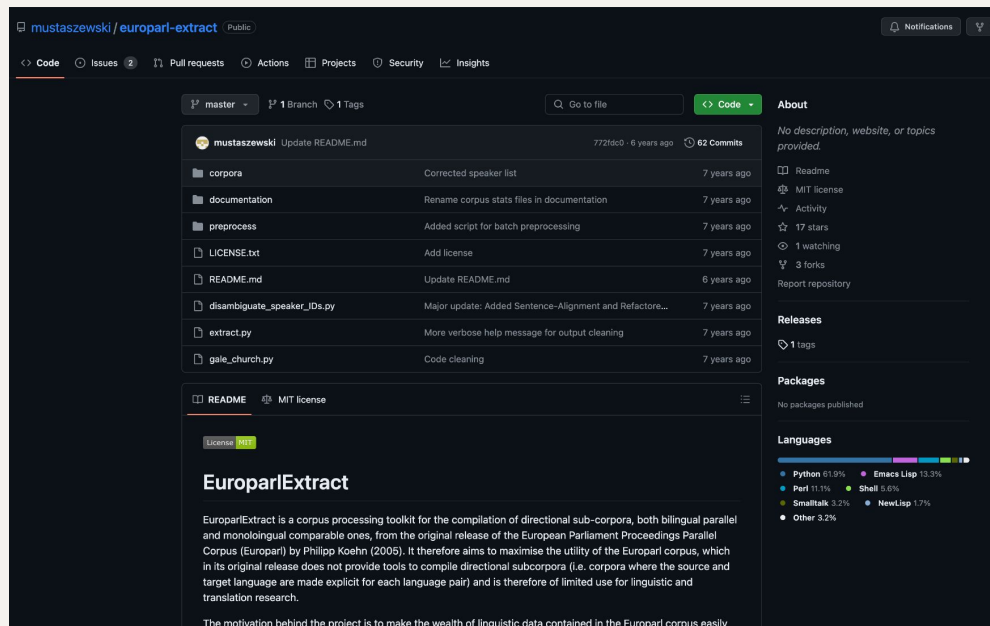illia.polosukhin@gmail.com

### Abstract

The dominant sequence transduction models are based on complex recurrent or convolutional neural networks that include an encoder and a decoder. The best performing models also connect the encoder and decoder through an attention mechanism. We propose a new simple network architecture, the Transformer, based solely on attention mechanisms, dispensing with recurrence and convolutions entirely. Experiments on two machine translation tasks show these models to be superior in quality while being more parallelizable and requiring significantly less time to train. Our model achieves 28.4 BLEU on the WMT 2014 English-to-German translation task, improving over the existing best results, including ensembles, by over 2 BLEU. On the WMT 2014 English-to-French translation task, our model establishes a new single-model state-of-the-art BLEU score of 41.8 after training for 3.5 days on eight GPUs, a small fraction of the training costs of the best models from the literature. We show that the Transformer generalizes well to other tasks by applying it successfully to English constituency parsing both with large and limited training data.
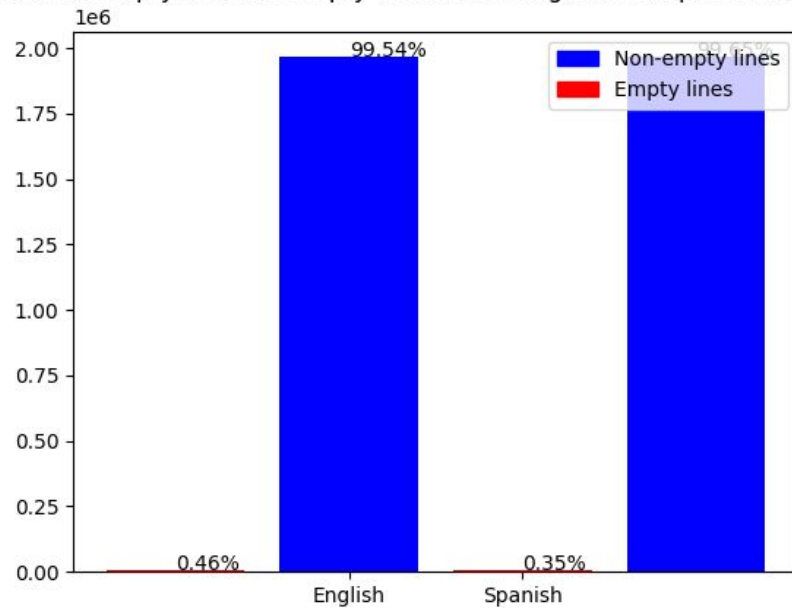
# 2.
# Methods

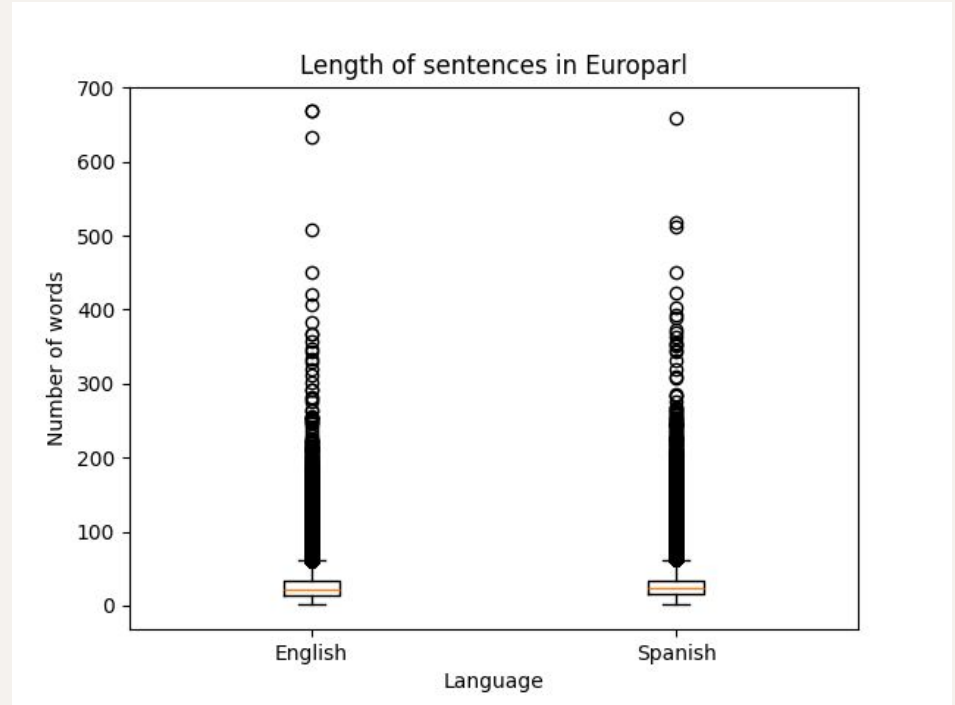# Preprocessing

- Identify speakers

# Preprocessing

- Identify speakers
- Manage empty lines



Number of empty and non-empty lines in the english and spanish translation
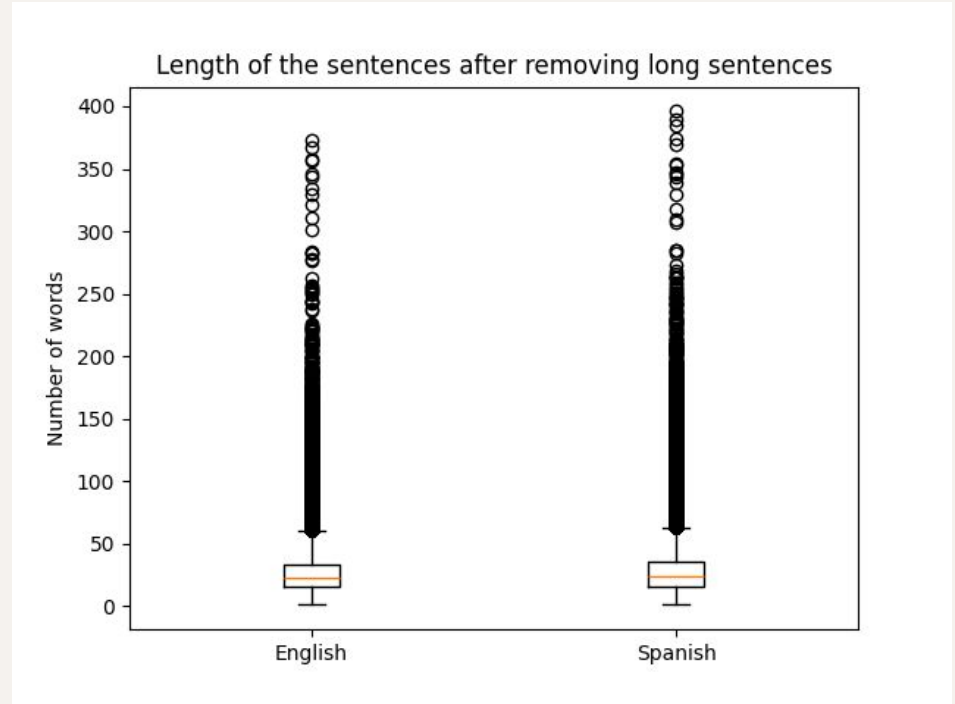
# Preprocessing

- Identify speakers
- Manage empty lines
- Feature engineering

# Preprocessing

- Identify speakers
- Manage empty lines
- Feature engineering



Length of the sentences after removing long sentences

# Sequence to Sequence

Using 'default' hyperparameters

because of long training time

## NLP From Scratch: Translation with a Sequence to Sequence Network and Attention

**Author:** Sean Robertson

This is the third and final tutorial on doing "NLP From Scratch", where we write our own classes and functions to preprocess the data to do our NLP modeling tasks. We hope after you complete this tutorial that you'll proceed to learn how *torchtext* can handle much of this preprocessing for you in the three tutorials immediately following this one.

In this project we will be teaching a neural network to translate from French to English.

```
[KEY: > input, = target, < output]

> il est en train de peindre un tableau .
= he is painting a picture .
< he is painting a picture .

> pourquoi ne pas essayer ce vin delicieux ?
= why not try that delicious wine ?
< why not try that delicious wine ?

> elle n est pas poete mais romanciere .
= she is not a poet but a novelist .
< she not not a poet but a novelist .

> vous etes trop maigre .
= you re too skinny .
< you re all alone .
```

... to varying degrees of success.

This is made possible by the simple but powerful idea of the sequence to sequence network, in which two recurrent neural networks work together to transform one sequence to another. An encoder network condenses an input sequence into a vector, and a decoder network unfolds that vector into a new sequence.

# Evaluation metrics

BLEU score

- Quick to understand
- Similar to how humans evaluate translated sentences
- Widely used

## BLEU: a Method for Automatic Evaluation of Machine Translation

Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu
IBM T. J. Watson Research Center
Yorktown Heights, NY 10598, USA
{papineni,roukos,toddward,weijing}@us.ibm.com

### Abstract

Human evaluations of machine translation are extensive but expensive. Human evaluations can take months to finish and involve human labor that can not be reused. We propose a method of automatic machine translation evaluation that is quick, inexpensive, and language-independent, that correlates highly with human evaluation, and that has little marginal cost per run. We present this method as an automated understudy to skilled human judges which substitutes for them when there is need for quick or frequent evaluations.[1]

the evaluation bottleneck. Developers would benefit from an inexpensive automatic evaluation that is quick, language-independent, and correlates highly with human evaluation. We propose such an evaluation method in this paper.

### 1.2 Viewpoint

How does one measure translation performance? *The closer a machine translation is to a professional human translation, the better it is.* This is the central idea behind our proposal. To judge the quality of a machine translation, one measures its closeness to one or more reference human translations according to a numerical metric. Thus, our MT evaluation system requires two ingredients:

$$Bleu\ (N) = Brevity\ Penalty \cdot Geometric\ Average\ Precision\ Scores\ (N)$$

# 3.
# Results

# Partial results

We can evaluate the results in terms of

- Adequacy: Translation w(t) should adequately reflect the linguistic content of w(s)

- Fluency: Translation w(t) should be fluent text in the target language

- BLEU Score

# Partial results Bahdanau

We can evaluate the results in terms of

- Adequacy: Translation w(t) should adequately reflect the linguistic content of w(s)

- Fluency: Translation w(t) should be fluent text in the target language

- BLEU Score

> It perceives the Union' s sustainable development to be the efficiency of markets, goods, services, capital and employment
= Para ella, el desarrollo duradero de la Unión es la eficacia de los mercados, de los bienes, de los servicios, de los capitales y del trabajo
< la eficacia de los mercados, bienes, servicios, capital y empleo <EOS>

> A huge sum, although it remains within the agreed 20% of the estimate for all European institutions put together
= Una suma astronómica, aunque no sobrepasa el 20% de la estimación acordada para el conjunto de todas las instituciones europeas
< en el acuerdo acordado del 20 % de las previsiones de todas las instituciones europeas que juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos <EOS>

# Partial results Bahdanau

We can evaluate the results in terms of

- Adequacy: Translation w(t) should adequately reflect the linguistic content of w(s)

- Fluency: Translation w(t) should be fluent text in the target language

> It perceives the Union' s sustainable development to be the efficiency of markets, goods, services, capital and employment
= Para ella, el desarrollo duradero de la Unión es la eficacia de los mercados, de los bienes, de los servicios, de los capitales y del trabajo
< la eficacia de los mercados, bienes, servicios, capital y empleo <EOS>

> A huge sum, although it remains within the agreed 20% of the estimate for all European institutions put together
= Una suma astronómica, aunque no sobrepasa el 20% de la estimación acordada para el conjunto de todas las instituciones europeas
< en el acuerdo acordado del 20 % de las previsiones de todas las instituciones europeas que juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos <EOS>

In orange, synonyms, in green coincidences between true translation and model prediction

# Partial results Bahdanau

We can evaluate the results in terms of

- Adequacy: Translation w(t) should adequately reflect the linguistic content of w(s)
- Fluency: Translation w(t) should be fluent text in the target language
- BLEU Score: between 0 and 1 out of 100 pairs

Bahdanau attention:
3.99967114868466e-06

> It perceives the Union' s sustainable development to be the efficiency of markets, goods, services, capital and employment
= Para ella, el desarrollo duradero de la Unión es la eficacia de los mercados, de los bienes, de los servicios, de los capitales y del trabajo
< la eficacia de los mercados, bienes, servicios, capital y empleo <EOS>

> A huge sum, although it remains within the agreed 20% of the estimate for all European institutions put together
= Una suma astronómica, aunque no sobrepasa el 20% de la estimación acordada para el conjunto de todas las instituciones europeas
< en el acuerdo acordado del 20 % de las previsiones de todas las instituciones europeas que juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos juntos <EOS>

In orange, synonyms, in green coincidences between true translation and model prediction

# Full results Multihead & Bahdanau

> The recognition of the uniqueness of each individual and their equal right to the opportunities offered by life is one of the attributes of Europe's united diversity, which is a central element in the union's cultural, political and social integration
= El reconocimiento de la singularidad de cada individuo y su derecho igualitario a disfrutar las oportunidades que ofrece la vida es uno de los atributos de la diversidad conjunta de Europa, que es un elemento central de la integración social, política y cultural de la Unión
< el derecho de la sociedad de la Unión Europea, es decir, una de las características de la diversidad de la Unión Europea, que es un elemento fundamental en la integración de la Unión Europea, la integración de la Unión Europea, la integración de la integración en la esfera de la vida cultural y social de la Unión Europea <EOS>

> The headage payments in Greece for each holding have not risen since 1989-1991
= En mi país, las cuotas sobre el ganado por cada explotación no han aumentado desde 1989-1991
< no siempre que haya sido la luz <EOS>

> Kofi Annan has described the AIDS crisis as being an issue of weapons of mass destruction
= Kofi Annan ha dicho que la crisis del sida es un arma de destrucción masiva
< de la Unión Europea de una especie de destrucción masiva <EOS>

In orange, synonyms, in green coincidences between true translation and model prediction, in yellow original input correctly translated in model prediction

> Therefore, we are now proposing that the food facility and extra appropriations to Palestine, Kosovo and Afghanistan be found by using reserves, for example, the flexibility reserve
= Por lo tanto, ahora proponemos que el mecanismo alimentario y los créditos extras para Palestina, Kosovo y Afganistán se cumplan utilizando las reservas, como por ejemplo, la reserva de flexibilidad
< y los créditos adicionales para Palestina, Kosovo y Afganistán se encuentran utilizando de forma ejemplar por ejemplo, la reserva de flexibilidad de la flexibilidad de flexibilidad <EOS>

> This was brought out in the resolution Parliament adopted in October
= Este asunto fue tratado en la resolución del Parlamento adoptada en octubre
< octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre de octubre aprobado octubre de octubre en octubre de octubre de octubre en octubre de octubre en octubre de octubre <EOS>

In orange, synonyms, in green coincidences between true translation and model prediction, in yellow original input correctly translated in model prediction

BLEU Score:

Bahdanau attention:
3.99967114868466e-06

BLEU Score:

Multi-head attention:
1.7302134354230151e-06

# 4.
# Discussion

# Conclusion & Limitations

- It needs a lot of computational power to train

- Not reliable (the BLEU scores are low)

- Limited context understanding

- Longer sentences are more difficult to translate accurately

# Potential Fixes

- Limit the sentence length
  - **+** The training will be faster, so we will be able to train more epochs
  - **-** We will have to limit the size of the dataset


- Optimize hyperparameters
  - **+** The models would converge faster
  - **-** The process is resource intensive

# Future (potential) work

- Fixing the network

# Future (potential) work

# Thank you!