A Multilingual Neural Machine Translation Model for COVID-19



Alexandre Bérard Zae Myung Kim Vassilina Nikoulina Eunjeong Lucy Park Matthias Gallé Kweon Woo Jung



NAVER LABS EUROPE NAVER PAPAGO







Motivation

Help translate COVID-19-related text, by creating a **public, biomedical-domain, multilingual** machine translation model





Motivation

1. Help translate COVID-19 guidelines,* news articles and social media reactions

* We recommend expert verification for critical translations

State-of-the-art public NMT systems are only available for a couple of languages.

None is adapted to the biomedical domain.





Sebastian Seung @SebastianSeung

1/ SAVE LIVES by translating Korean→English! 75 page playbook for @KoreaCDC's fight against **#Covid19**. Let's translate into English by Monday morning. No time to lose! Please share the link. bit.ly/CovidPlaybook

9:15 PM · Mar 27, 2020 · Twitter Web App

1.3K Retweets 93 Quote Tweets 2.4K Likes



TAUS @T21Century · Apr 6

a result of the TAUS Corona Crisis Project, the first batch of corona corpora are available in four language pairs. Everybody can download it for free! Take a look now: bit.ly/347eLrn





Motivation

1. Help translate COVID-19 guidelines,* news articles and social media reactions

2. Apply existing NLP tools to other languages than English

Examples of English-centric NLP tools for COVID-19:

- Rapidly Deploying a Neural Search Engine for the COVID-19 Open \bullet Research Dataset: Preliminary Thoughts and Lessons Learned
- Document Classification for COVID-19 Literature \bullet
- What Are People Asking About COVID-19? A Question \bullet Classification Dataset
- Measuring Emotions in the COVID-19 Real World Worry Dataset \bullet
- NLP-based Feature Extraction for the Detection of COVID-19 Misinformation Videos on YouTube





1/ SAVE LIVES by translating Korean→English! 75 page playbook for @KoreaCDC's fight against **#Covid19**. Let's translate into English by Monday morning. No time to lose! Please share the link. bit.ly/CovidPlaybook

9:15 PM · Mar 27, 2020 · Twitter Web App

1.3K Retweets 93 Quote Tweets 2.4K Likes

TAUS @T21Century · Apr 6









Challenges

In-domain data availability

Multilinguality



Compatibility with an existing framework

Evaluation



Related Work

- **TAUS Corona Crisis corpora** \bullet
- SYSTRAN Corona Crisis models \bullet
- TICO-19 resources \bullet
- NLP COVID-19 Workshop \bullet

Publications:

- Facilitating Access to Multilingual COVID-19 Information \bullet via Neural Machine Translation, Way et al. (2020)
- A System for Worldwide COVID-19 Information lacksquareAggregation, Aizawa et al. (2020)



TAUS @T21Century · Apr 6

As a result of the **TAUS Corona Crisis** Project, the first batch of **corona corpora** are available in four language pairs. Everybody can download it for free! Take a look now: bit.ly/347eLm





CONTENTS

- 2. Domain Adaptation
- 3. Our Model
- 4. Evaluation and examples

1. Introduction to Neural Machine Translation



1. Introduction to Neural Machine Translation





Transformer (Vaswani et al., 2017) is now the standard in MT

It is also the architecture of BERT (Devlin et al., 2018), the basis of most current NLP work



































Figures from "The Illustrated Transformer" blog post

















1.2 Training





Training data: thousands/millions of (source, target) pairs of translated sentences



1.2 Training

for epoch in range(N):
for source, target in data:
 probs = model(source, target[1:]) # forward pass (model outputs)
 loss = cross_entropy(probs, target) # forward pass (loss)
 optimizer.backward(loss) # backward pass: compute gradients
 optimizer.step() # update the model weights

In practice, **source** and **target** are tensors corresponding to batches of multiple sentences High parallelism is achieved thanks to matrix multiplications on GPUs





Training data filtering: language identification & length-based filtering 1.



Numbers from *Khayrallah and Koehn (2018)*



Language id	Length filtering	Copies	Hallucinations
No	No	12%	5%
Yes	No	0%	3%
Yes	Yes	0%	1%

Table from *Berard et al. (2019)*



- **1. Training data filtering:** language identification & length-based filtering
- **Unicode normalization** 2.







- **1. Training data filtering:** language identification & length-based filtering
- Unicode normalization
- **3.** BPE (Byte Pair Encoding)



Credits to Chris Manning's "NLP with Deep Learning" lecture







- **1. Training data filtering:** language identification & length-based filtering
- 2. Unicode normalization
- 3. BPE (Byte Pair Encoding)

l, o, w, e, r, **BPE vocabulary:** n, w, s, t, i, d, es, est, lo







1.4 Evaluation

How to measure progress in MT? How to compare two systems?

1. Automatic metrics \rightarrow BLEU, METEOR, chrF, etc.

> Parallel corpus

- 2. Human evaluation
- \rightarrow Relative ranking
- \rightarrow Direct assessment









- One single model for multiple source languages



Joint multilingual BPE model & shared vocabulary and embeddings





Choose target language with source-side tag: <2FR>, <2IT>, etc.







But most training data is paired with English...























→ "Monolingual Adapters for Zero-Shot Neural Machine translation", Philip et al. (2020)



BLEU differences with bilingual models, from high to low-resource languages



Any-to-English (BLEU Δ)

Graphs from Google's "Massively Multilingual NMT in the Wild" paper



English-to-Any (BLEU Λ)



1.6 MT Framework: fairseq

Ļ	<u>pytorch/fairseq</u>			
	python	pytorch	artificial-in	
	☆ 9.9k	Python	MIT license	

Fast training, thanks to:

- 1. Multi-GPU and delayed updates
- Large batches: increased learning rate and faster convergence 2.
- Mixed-precision training 3.



ce-to-Sequence Toolkit written in Python.

ntelligence

Updated 1 hour ago 82 issues need help





2. Domain Adaptation



2.1 Domain shift







2.2 Domain-specific issues

- Lexical ambiguity \rightarrow carte in French (means card, map, or menu)
- Unknown words or expressions \rightarrow "noix de Saint-Jacques" \bullet
- Wrong level of politeness $\rightarrow tu / vous$ in French
- Robustness to noise \rightarrow typos, missing punctuation, etc. \bullet
- Others \rightarrow dialect, punctuation norms, etc.





2.3 Examples of domains



INPUT: 방탄소년단 보라해 **GENERIC:** BTS Borahae **DOMAIN-ADAPTED:** BTS, I purple you.





INPUT: 여자친구 노래 들어보자

GENERIC: Let's listen to your girlfriend's song.

DOMAIN-ADAPTED: Let's listen to GFRIEND's song.


2.3 Examples of domains

INPUT: 노래 틀고 놀자 그래 좋은 것 같아 노래 틀어줘

GENERIC: Let's play some music. Yeah.

Sounds good. Play some music.

DOMAIN-ADAPTED: Let's play some music.

I think it's good. Play some music.





INPUT: 그렇지x2 멋있어x2 GENERIC: That's right. x2 It's cool. x2 DOMAIN-ADAPTED: That's right. x2 Cool. x2



2.3 Examples of domains





Meilleur rapport qualité/prix asiatique de Liège! Grande carte et tout est très bon. Préférences pr: le poulet croquant aigre douce, poulet curry du chef, cuisses de poulet farcies et le canard laqué!

Vote positif 1

Vote négatif

Best value for money in Asia from Liège! Large card and everything is very good. Preferences: sweet sour crunchy chicken, chief curry chicken, stuffed chicken legs and lacquered duck!

Vote positif 1

Vote négatif

-> "Machine Translation of Restaurant Reviews: New Corpus for Domain Adaptation and Robustness", Berard et al. (2019)





Mai 10, 2014

Mai 10, 2014



2.4 What is domain adaptation?







2.5 Fine-tuning

Continue training the generic model on the in-domain data only

Advantages:

- Quick and easy \bullet
- Best in-domain performance

Drawbacks:

- One model per domain





Catastrophic forgetting \rightarrow quality degradation on the other domains



2.6 Domain tags

Train a model from scratch with data from multiple domains, with control tags indicating the domain

Advantages:

- Easy to implement lacksquare
- Multi-domain models with easy control at inference time

Drawbacks:

Very costly \rightarrow need to re-train the model from scratch each time a new domain comes in









3. Our Model

3.1 Training data

Generic: ParaCrawl, United Nations, Europarl, AI Hub, OpenSubtitles, WikiMatrix, IWSLT (TED Talks), etc.

Biomedical: Medline, TAUS, UMLS, etc.

Language	Total	Generic	Back- translated	Biomedical
French	124.0	119.7		4.3
Spanish	88.9	86.2		2.7
German	84.5	81.6		2.9
Italian	43.9	43.1		0.8
Korean	13.4	5.5	7.8	0.1
Total	354.8	336.2	7.8	10.8

Numbers in millions of sentence pairs (with **English** as target)





3.2 Pre-processing

Joint BPE with SentencePiece \rightarrow partially shared source/target vocabulary & • embeddings



target_embed = source_embed[:38000] *# share parameters*





3.2 Pre-processing

- Joint BPE with SentencePiece \rightarrow partially shared source/target vocabulary & \bullet embeddings
- **Inline casing:** robustness to capitalization



 \rightarrow "NAVER LABS Europe's Systems for the WMT19 Machine Translation Robustness Task", Berard et al. (2019)





Transformer Big with:

- *Wide* 6-layer encoder of size 8192 \bullet
- Shallow 3-layer decoder \bullet

Total parameters: 254M Checkpoint size: 2.7GiB

















- Up-sampling biomedical data by × 2
- <medical> tag before each source biomedical sentence

Language	Generic	Back- translated	Biomedical	Total
French	119.7		4.3 × 2	128.3
Spanish	86.2		2.7 × 2	91.6
German	81.6		2.9 × 2	87.4
Italian	43.1		0.8 × 2	44.7
Korean	5.5	7.8	0.1 × 2	13.5
Total	336.2	7.8	10.8	365.5



2 biomedical sentence



- Up-sampling biomedical data by × 2
- <medical> tag before each source biomedical sentence
- Per-language sampling with temperature 5

Language	Generic	Back- translated	Biomedical	Total	Frequency	Sampling probability
French	119.7		4.3 × 2	128.3	0.35	0.23
Spanish	86.2		2.7 × 2	91.6	0.25	0.22
German	81.6		2.9 × 2	87.4	0.24	0.21
Italian	43.1		0.8 × 2	44.7	0.12	0.19
Korean	5.5	7.8	0.1 × 2	13.5	0.04	0.15
Total	336.2	7.8	10.8	365.5	1.0	1.0



2 biomedical sentence rature 5



3.4 On-the-fly preprocessing





Our dynamic pipeline





3.4 On-the-fly preprocessing

Advantages:

- Less user-side pre-processing hassle lacksquare
- Easy to randomly sample from multiple corpora \bullet
- Easy to apply stochastic pre-processing (e.g., BPE dropout) \bullet

Implementation issues:

- How to efficiently pre-process and batch to avoid creating a bottleneck \bullet
- How to shuffle





3.4 On-the-fly preprocessing







3.5 How to use the model

Attps://github.com/naver/covid19-nmt

1. Install fairseq

git clone https://github.com/pytorch/fairseq cd fairseq python3 -m venv env python3 -m pip install sentencepiece



- source env/bin/activate python3 -m pip install --editable .



3.5 How to use the model

Attps://github.com/naver/covid19-nmt

2. Download the model

> covid19-nmt/Covid19/checkpoint_best.pt



git clone https://github.com/naver/covid19-nmt.git cat covid19-nmt/Covid19/checkpoint best.pt.part* \



3.5 How to use the model

Attps://github.com/naver/covid19-nmt

3. Good to go!

- fairseq-interactive covid19-nmt/Covid19 --user-dir covid19-nmt/Covid19 \ --path covid19-nmt/Covid19/checkpoint best.pt --medical \ -s src -t en --bpe covid19 --buffer-size 1000 --max-tokens 8000 --fp16 \ < INPUT | grep "^D" | cut -f3 > OUTPUT

 \rightarrow "A Multilingual Neural Machine Translation Model for Biomedical Data", Berard et al. (2020)





4. Evaluation and examples





4.1 New Korean-English test set

Human translation of handpicked English sentences into Korean

Covers: safety guidelines, gov. briefings, clinical tests, biomedical exp., etc.

Source	Sentence pairs
arXiv	500
KCDC	258
All	758

arXiv \rightarrow Abstracts of biomedical papers on SARS-CoV-2 and COVID-19 **KCDC** \rightarrow Official guidelines and reports from Korea Centers for Disease Control





4.2 Generic evaluation

Language	Model	News
	Ours	41.0
French	NLE@WMT19	40.2
	OPUS-MT	38.0
	Ours	41.3
German	FAIR@WMT19	41.0
	OPUS-MT	39.5
Spanich	Ours	36.6
Spanish	OPUS-MT	35.2





Naver Labs Europe's Systems for the WMT19 **Machine Translation Robustness Task**

Alexandre Bérard Ioan Calapodescu **Claude Roux**

Naver Labs Europe first.last@naverlabs.com

Facebook FAIR's WMT19 News Translation Task Submission

Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, Sergey Edunov Facebook AI Research, Menlo Park, CA & New York, NY.



4.2 Generic evaluation

Language	Model	News	TED Talks
	Ours	41.0	41.1
French	NLE@WMT19	40.2	-
	OPUS-MT	38.0	38.9
	Ours	41.3	31.6
German	FAIR@WMT19	41.0	32.0
	OPUS-MT	39.5	30.4
Spanich	Ours	36.6	48.8
Spanish	OPUS-MT	35.2	47.3
Italian	Ours	-	42.2
Italiali	OPUS-MT	-	40.4
Koroon	Ours	-	21.3
Korean	OPUS-MT	-	17.8





Naver Labs Europe's Systems for the WMT19 **Machine Translation Robustness Task**

Alexandre Bérard Ioan Calapodescu **Claude Roux**

Naver Labs Europe first.last@naverlabs.com

Facebook FAIR's WMT19 News Translation Task Submission

Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, Sergey Edunov Facebook AI Research, Menlo Park, CA & New York, NY.



4.3 Biomedical-domain evaluation

Participation in the WMT20 Biomedical task Official BLEU results (Medline 2020):

Language	Our model	(Next) best m
French	43.1	44.1
German	34.1	34.8
Spanish	50.6	46.4
Italian	42.5	42.0

Remarks:

- Single model, not an ensemble
- Multilingual
- Not tuned on Medline \rightarrow also good on other domains

-> "NAVER LABS Europe's Participation in the Robustness, Chat, and Biomedical Tasks at WMT 2020", Berard et al. (2019)





BLEU on our new test set:

Language	Our model	OPUS-MT
Korean	37.4	20.0



SOURCE: 결과는 ACE2 발현이 코점막 상피세포에서 나타난다는 것을 보여줍니다.

OPUS-MT: The result is that **ACE2 spins out of** the **epidural upper blood** cell.



- **REFERENCE:** The result indicates that the **ACE2 expression appears** in **nasal epithelial** cells.
- **OUR MODEL:** The results show that **the expression of ACE2 occurs** in the **nasal epithelial** cells.



SOURCE: 근육통 또는 관절통(OR=0.28; 95%CI=0.10~0.66), 미각장애(OR=0.28;



- 95%CI=0.05~0.92) 및 후각소실증(OR=0.23; 95%CI=0.04~0.75)은 보호 인자였습니다.
- **REFERENCE:** Myalgia or artralgia (OR=0.28; 95%CI=0.10 to 0.66), dysgeusia (OR=0.28;
- 95%CI=0.05 to 0.92) and anosmia (OR=0.23; 95%CI=0.04 to 0.75) were protective factors.
- **OUR MODEL:** Muscle or joint pain (OR=0.28; 95%CI= $0.10^{\circ}0.66$), taste disturbance (OR=0.28;
- 95%CI=0.05~0.92) and olfactory loss (OR=0.23; 95%CI=0.04~0.75) were protective factors.
- **OPUS-MT:** A muscle or joint pain (OR) 0.89% CI (0.10 to 0.66), a taste disorder (. 28; 95% CI = (0.05.92) and a olfactory osmosis (OR = 0.23; 95%Cl = 0.04.0.0.75) were protective factors.



사례 및 잠재적 사례)과 초과 사망률은 거의 차이가 나지 않습니다. (confirmed and possible cases) and excess mortality. (confirmed and potential cases) and the excess mortality rate for coronavirus-19. excess mortality rate is very little different.



- SOURCE: 벨기에의 경우, 코로나바이러스감염증-19에 대한 보고된 사망률(확인된
- **REFERENCE:** Belgium has virtually no discrepancy between COVID-19 reported mortality
- **OUR MODEL:** In Belgium, there is **little difference** between the reported mortality rate
- **OPUS-MT:** In Belgium, the reported mortality rate for **Coronar virus infection -- 19 --** and the



나타냈습니다.

REFERENCE: Conversely, for **under-expressed** genes, **pathways** indicated **repression of lymphocyte differentiation** and T cell activation.

OUR MODEL: Conversely, in the case of the **hypoplastic** gene, the **pathway** showed

Iymphocyte differentiation inhibition and T cell activation.

OPUS-MT: On the other hand, in the case of a **low-end** gene, the **path** showed **lymph nodes** and T cells activated.



SOURCE: 반대로, 저발현 유전자의 경우 경로는 림프구 분화 억제 및 T 세포 활성화를



PCR에 의한 SARS-CoV-2 RNA는 음성으로 나타났다.

SARS-CoV-2 RNA by RT-PCR resulted negative.

while SARS-CoV-2 RNA by RT-PCR was negative.

CoV-2 RNA in RT-PCR is negative.



- SOURCE: CSF는 68%의 경우에 hyperproteinorrachia 및 / 또는 다구증을 보인 반면, RT-
- **REFERENCE:** CSF showed **hyperproteinorrachia** and/or **pleocytosis** in 68% of cases whereas
- **OUR MODEL:** CSF showed **hyperproteinorrhea** and/or **multiple sclerosis** in 68% of cases,
- **OPUS-MT:** While CSF shows **hyperproteinocia** and / or **multiplexia** in 68% of cases, SARS-



- New languages: Chinese and Russian (available in TAUS Corona Crisis) corpora)
- More training data (689M)
- New domain tags: <medical>, <medline>, <ted>, <food>, <subtitles>, <religious>, <wiki>









BLEU scores by domain Averages over all languages



BLEU scores by domain Chinese-English





BLEU scores by domain **Russian-English**







- 6 epochs = 13 days





BLEU difference with or without tag Averages over all languages









Conclusion





Conclusion

Contributions:

- Multilingual multi-domain model, compatible with *fairseq*
- New Korean-English COVID-19 test set lacksquare
- Participation in WMT20 Biomedical Task \rightarrow top performance in 2/4 languages
- \bullet • Version 2.0 of this model with more languages and more domains

Future work:

- Add more languages
- Share a one-to-many model
- Share domain-specific adapter layers \bullet






Thank You



