## Natural Language Understanding with Multilingual Neural Models - or the Power of GPU's in

Modern Language Technology



the World's languages

#### Jörg Tiedemann

Department of Digital Humanities University of Helsinki jorg.tiedemann@helsinki.fi

# Ηοι

# How do humans become intelligent?



argue understand analyze learn interpret listen

discuss interact read

# How do machines become intelligent?





# **Democratize NLP ...**







Helsinki-NLP/Opus-MT github.com

22.46 · 21/01/2020 · Twitter for iPhone







💷 elisa 😤

:

Hugging Face @huggingf... · 14h ∨ Let's democratize NLP for all languages! 🌍 🌍 🌍

Today, with v2.9.1, we are releasing 1,008 machine translation models, covering ` of 140 different languages trained by @jorgtiedemann with @marian, ported by @sam\_shleifer. Find your language here: huggingface.co/ models?search=... [1/4]



805

#### You and 8 others







# **Democratize information & Al**





https://peacemachine.net



all the language data I can possibly find

#### my big fat neural network ....



#### my big fat neural network ....



#### my big fat neural network ....





#### Transformer language models (BERT, ELECTRA, ...)

- parameters: 14M (small), 110M (base), 335M (big)
- crawled data: 420 billion words (English), 3 billion (Finnish)
- wikimedia: 4.5 billion words (English), 108 million (Finnish)
- training takes ~ 1-2 weeks (base model)



#### Transformer language models (BERT, ELECTRA, ...)

- parameters: 14M (small), 110M (base), 335M (big)
- crawled data: 420 billion words (English), 3 billion (Finnish)
- wikimedia: 4.5 billion words (English), 108 million (Finnish)
- training takes ~ 1-2 weeks (base model)

Translation model German/English/Finnish/Dutch/Swedish

- 1.46 billion sentence pair in training
- ca. 90 million parameters to be learned
- multi-GPU training on four v100 GPUs
- start to converge after ca 7 days (?) still running ...





Helsingin yliopisto Helsingfors Universitet University of Helsinki



# **Multilingual machine translation**



# Semantic representation learning



(From Sutskever et. alet. al: "Sequence to Sequence Learning with Neural Networks")





Architecture proposed by Cífka and Bojar (2018).

Our implementation in OpenNMT-py (MTM2018)





Architecture proposed by Cífka and Bojar (2018).

Our implementation in OpenNMT-py (MTM2018)

# Implementation of an "Attention Bridge"



Architecture proposed by Cífka and Bojar (2018).

Our implementation in OpenNMT-py (MTM2018)



# Translating image captions (6 languages)





# Natural Language Inference (NLI)

#### **Benchmark for reasoning with language**

A black race car starts up in front of a crowd of people.

contradicts

A man is driving down a lonely road. A soccer game with multiple males playing. entails Some men are playing a sport.



#### **Examples from the SICK dataset**

#### ENTAILMENT, relatedness score = 4.7

The young boys are playing outdoors and the man is smiling nearby The kids are playing outdoors near a man with a smile

#### CONTRADICTION, relatedness score = 3.6

The young boys are playing outdoors and the man is smiling nearby There is no boy playing outdoors and there is no man smiling

#### NEUTRAL, relatedness score = 1.7

A lone biker is jumping in the air A man is jumping into a full pool



## Attention bridge in downstream tasks







visual grounding

- En: A wall divided the city.
- De 1: Eine Wand teilte die Stadt. ×
- De 2: Eine Mauer teilte die Stadt. ✓

 translational grounding





Auto®

## MeMAD: Access to audio-visual content

Car



automatic transcription of videos accessible in many languages

https://memad.eu

think that you saw too many, ovies about the James Bond.

אני חושב שאר

יותר מדי סרטים על

ג'ימס בונד.

# All of this is powered by CSC!

puhti, cPouta, ObjectStorage, allas, ...



# Constant need of billing units ...

OPUS Billing units remaining: 10.

103888 / 1030901

OPUS-MT Billing units remaining:

42.8%

1073760 / 2510000

OPUS-LM

**Billing units remaining:** 

100%

1010000 / 1010000

NLPL Billing units remaining:

45.6%

870941 / 1910000

MeMAD Billing units remaining:

28.8%

463936 / 1610000

CrossNLP Billing units remaining:

0/3820087

Multi-MT Billing units remaining:

0 / 2319155

FoTran Billing units remaining: 5 614261 / 11500000 LingDA Billing units remaining: 49.7% 99364 / 200000

> 22 million BUs used in the projects above, mostly in the last 2 years



## Storage needs keep on growing ...

#### Project scratch



> 20T in total

Other data storage

#### 10T 7.5T ??? T

NLPL (puhti)
OPUS (cPouta)
ObjectStorage, IDA





https://translate.ling.helsinki.fi







- command-line, ssh, bash-scripts, slurm, github, slack
- pyTorch, tensorflow, OpenNMT, MarianNMT, ....
- avoid graphical user interfaces for research & development

bash	tiedeman@puhti-login2:~/research/Opus-MT-train +	tiedeman@puhtinlpl/data/OPUS tiedeman@puhti/Opus-MT-			an@puhti/Opus-MT-train	
File Edit Options Buffers Tools Python I	Help	DGT	ELRA-W0182	ELRA-W0254	ELRA-W0290	
include lib/models/russian.mk		DOGC	ELRA-W0183	ELRA-W0255	ELRA-W0291	
include lib/models/sami.mk		ECB	ELRA-W0184	ELRA-W0256	ELRA-W0292	
include lib/models/wikimedia.mk		EhuHac	ELRA-W0188	ELRA-W0257	ELRA-W0293	
		Elhuyar	ELRA-W0189	ELRA-W0258	ELRA-W0294	
include lib/models/doclevel.mk		ELRA-W0130	ELRA-W0190	ELRA-W0259	ELRA-W0301	
include lib/models/simplify.mk		ELRA-W0131	ELRA-W0191	ELRA-W0260	ELRA-W0305	
		ELRA-W0133	ELRA-W0195	ELRA-W0261	EMEA	
		ELRA-W0134	ELRA-W0196	ELRA-W0262	EUbookshop	
		ELRA-W0135	ELRA-W0213	ELRA-W0263	EUconst	
.PHONY: all		ELRA-W0136	ELRA-W0214	ELRA-W0264	Europarl	
all: \${WORKDIR}/config.mk		ELRA-W0137	ELRA-W0215	ELRA-W0265	Finlex	
\${MAKE} data		ELRA-W0138	ELRA-W0216	ELRA-W0266	fiskmo	
\${MAKE} train		ELRA-W0142	ELRA-W0217	ELRA-W0267	giga-fren	
\${MAKE} eval		ELRA-W0143	ELRA-W0218	ELRA-W0270	GlobalVoices	
\${MAKE} compare		ELRA-W0144	ELRA-W0220	ELRA-W0271	GNOME	
		ELRA-W0145	ELRA-W0221	ELRA-W0272	hrenWaC	
-UU-:F1 Makefile 55% L185 G	it-master (Makefile)	ELRA-W0146	ELRA-W0222	ELRA-W0273	infopankki	
import random		ELRA-W0147	ELRA-W0223	ELRA-W0274	JRC-Acquis	
import re		ELRA-W0148	ELRA-W0226	ELRA-W0275	JW300	
import shutil		ELRA-W0149	ELRA-W0227	ELRA-W0276	KDE4	
from typing import Dict, List, Tuple		ELRA-W0150	ELRA-W0228	ELRA-W0277	KDEdoc	
		ELRA-W0151	ELRA-W0229	ELRA-W0278	komi	
import numpy as np		ELRA-W0152	ELRA-W0235	ELRA-W0279	Makefile	
import torch		ELRA-W0154	ELRA-W0238	ELRA-W0280	MBS	
from torch.nn.utils.rnn import pad_seque	ence	ELRA-W0156	ELRA-W0239	ELRA-W0281	memat	
from torch.utils.data import DataLoader	, Dataset, RandomSampler, SequentialSampl\	ELRA-W0157	ELRA-W0242	ELRA-W0282	MontenegrinSubs	
er		ELRA-W0159	ELRA-W0243	ELRA-W0283	MPC1	
from torch.utils.data.distributed import	t DistributedSampler	ELRA-W0160	ELRA-W0244	ELRA-W0284	MultiParaCrawl	
from tqdm import tqdm, trange		<pre>[tiedeman@p dic freq</pre>	uhti-login2 Op info LICENSE	mono mose	]\$ ls /projappl/n s parsed raw P	
from transformers import (		[tiedeman@puhti-login2 Opus-MT-train]\$ ls /projappl/n				
WEIGHTS_NAME,			latest/ v3/			
AdamW,			[tiedeman@puhti-login2 Opus-MT-train]\$ ls /projappl/m			
BertConfig,			bg.zip da.zip el.zip es.zip fi.zip hu.zip lt.zi			
-UU-:F1 run_language_modeling.py	3% L39 Git-master (Python)	cs.zip de.	zip en.zip e	t.zip fr.z	ip it.zip lv.zi	
		[tiedeman@p	uhti-login2 Op	ous-MT-train	]\$ ]	



# To sum up: What do we need?

#### **Computing resources**

- more GPU's
- shorter job queues (right now on puhti: 2408 jobs waiting)
- storage for data (general and temporary)
- sometimes fast I/O

#### On-going work and plans for the future

- cross-border activities (NLPL, EOSC-nordic)
- better and more efficient access to data
- improved replicability (sharing of data, models, code)





# **Thank You!**

# **Questions and discussions?**

https://blogs.helsinki.fi/language-technology/



# Language Technology in Helsinki

http://blogs.helsinki.fi/language-technology/





FOUND IN TRANSLATION SEMANTICS & MT > 1,000 languages











### https://github.com/Helsinki-NLP





http://opus.nlpl.eu Data collection for MT > 200 languages

**O**PUS





sentimentator



audiovisual data & MT





Methods for Managing Audiovisual Data https://memad.eu



