**Slide 1**

translations as semantic mirrors

# *"Found in Translation"*

## Emerging Language Spaces Learned From Massively Multilingual Corpora

using massive linguistic diversity

meaning

understanding
speaking

neural machine translation

a thousand languages

*Jörg Tiedemann*
*University of Helsinki*
*jorg.tiedemann@helsinki.fi*

---

**Slide 2**

## Machine Translation (MT)

human translations

meaning

understanding
speaking

dense vector-based representation

encoder
decoder
neural network

source language
target language

Learning Algorithm

---

**Slide 3**

## Multilingual Translation Models

human translations in many languages

meaning

one single model with **shared parameters** across all languages

understanding
speaking

dense vector-based representation

encoder
decoder
neural network

any language
any language

Learning Algorithm

---

**Slide 4**

## Neural Machine Translation

output sentence

$y_1$   $y_2$   $y_3$

Bible translations in > 900 languages

other languages

English

train

256 dimensions

512 dimensions

256 dimensions

attention

$h_1$   $h_2$   $h_3$   $h_4$

256 dimensions

act as trigger to control the influence of certain parameters

256 dimensions

$x_1$   $x_2$   $x_3$   $x_4$

language flags

input sentence

vocabulary: 50,000 sub-word units

## Neural Machine Translation

Bible translations in > 900 languages

other languages ↕ English

train →

The network needs to **compress** information!

Learns to **re-use parameters** for languages with **common properties**

output sentence

$y_1$  $y_2$  $y_3$

attention

$h_1$  $h_2$  $h_3$  $h_4$

$x_1$  $x_2$  $x_3$  $x_4$

language flags

input sentence

vector space of language representations

---

## Emerging Language Space: Batch 0

Legend:
- Italic
- Germanic
- Slavic
- Indo-Iranian
- Celtic
- Albanian
- Baltic
- Greek

language representations as geometric positions

(PCA)

---

## Emerging Language Space: Batch 1

Legend:
- Italic
- Germanic
- Slavic
- Indo-Iranian
- Celtic
- Albanian
- Baltic
- Greek

(PCA)

---

## Emerging Language Space: Batch 2

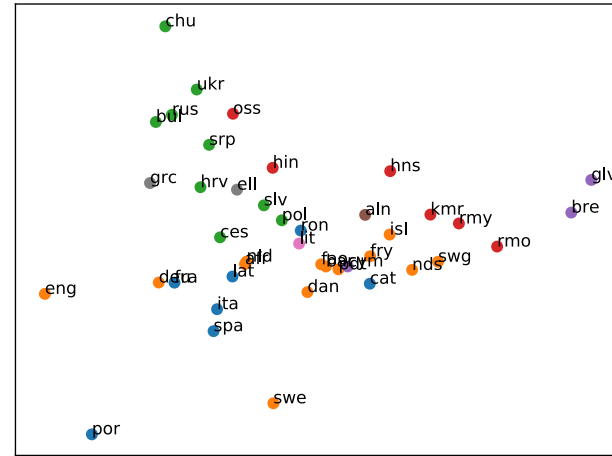Legend:
- Italic
- Germanic
- Slavic
- Indo-Iranian
- Celtic
- Albanian
- Baltic
- Greek
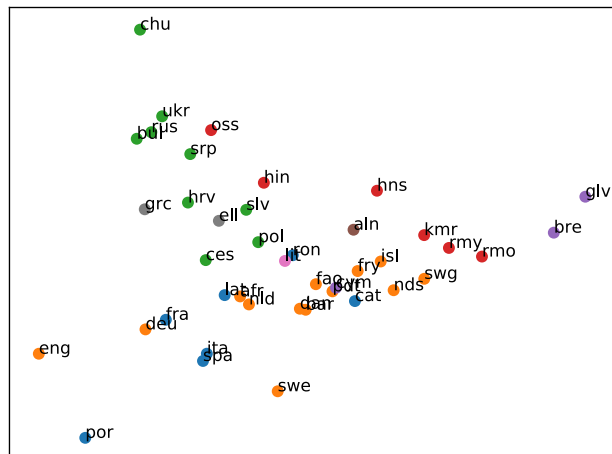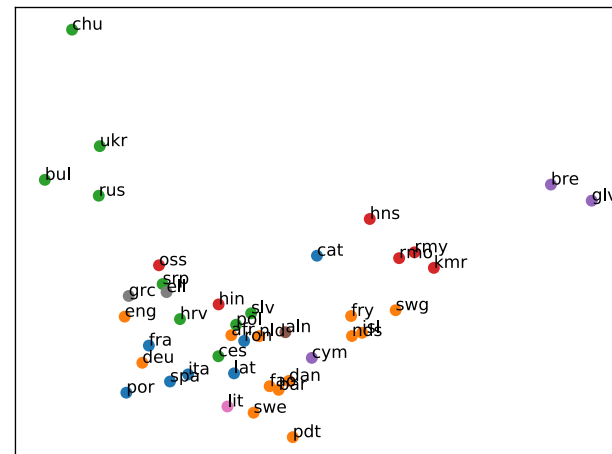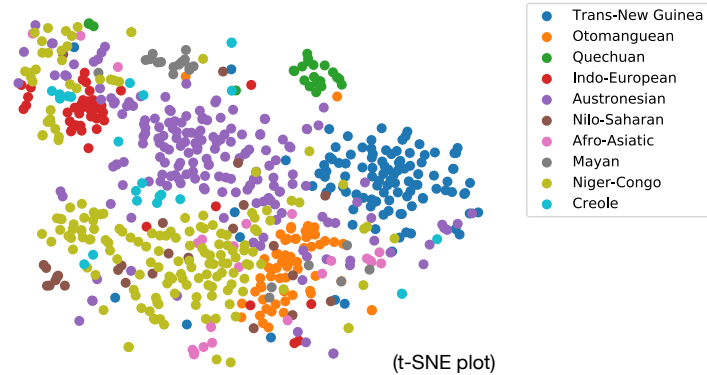
(PCA)

Emerging Language Space: Batch 3

Emerging Language Space: Batch 4

Emerging Language Space: Batch 5

Emerging Language Space: Batch 30

# Language Space of 972 Languages

Rough clusters of language families



Legend:
- Trans-New Guinea
- Otomanguean
- Quechuan
- Indo-European
- Austronesian
- Nilo-Saharan
- Afro-Asiatic
- Mayan
- Niger-Congo
- Creole

(t-SNE plot)

---

# Why Is This Interesting?

Continuous language space
- distances refer to relationship
- languages are not independent discrete units

Completely data-driven approach
- no prior knowledge
- driven by optimizing compression (with translation objective)

Interesting questions for future research
- Can we see specific linguistic properties?
- Combination with other tasks than MT

---

# It Is A Trendy Research Topic

- Academy of Finland project: **Digital language typology: mining from the surface to the core** (Vainio, Toivonen)
- Bjerva, J. and Augenstein, I. (2018) **From Phonology to Syntax: Unsupervised Linguistic Typology at Different Levels with Language Embeddings**, NAACL-HLT 2018.
- Bjerva, J. and Augenstein, I. (2018) **Tracking Typological Traits of Uralic Languages in Distributed Language Representations**, the Fourth International Workshop on Computational Linguistics for Uralic Languages (IWCLUL).
- Chaitanya Malaviya, Graham Neubig, Patrick Littell. **Learning Language Representations for Typology Prediction**, EMNLP 2017.
- Ehsaneddin Asgari and Hinrich Schütze: **Past, Present, Future: A Computational Investigation of the Typology of Tense in 1000 Languages**, EMNLP 2017
- ...

---

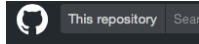# What Is Missing - What Is Next?

Many shortcomings
- data sources are limited and of very narrow domains
- the models are simple and generic
- difficult interpretation of results
- very little interaction with general linguistics

Ideas for the future
- emerging linguistic structures (syntax / semantics)
- diachronic models, different registers, ...
- training for specific phenomena with different objectives