

Language Transfer for Early Warning of Epidemics from Social Media



Mattias Appelgren¹, Patrick Schrempf^{1,2}, Matúš Falis¹, Satoshi Ikeda¹, Alison Q O'Neil^{1,3}

¹ Canon Medical Research Europe Ltd. ² University of St Andrews ³ University of Edinburgh

Purpose

- Analyse statements on social media to identify red flag medical symptoms, allowing **early detection** of the spread of disease.
- Explore how data in one language might be used to build models for a different language, using Japanese as our target language.

Experiments

Using BERT variants that are trained on different corpora (original BERT [1], mBERT [2] and jBERT [3]) we perform three experiments:

- Zero-shot transfer** with multilingual pre-training
- Training on **machine translated data**
- Mixing translated data** with original data

For translation we use Google and Amazon Translate [5,6].

Data

We use the MedWeb ("Medical Natural Language Processing for Web Document") dataset [2] which consists of pseudo-tweets in three different languages: **Japanese, English and Chinese**.

Dataset	#Pseudo-Tweets	Mean #labels per example	Influenza	Diarrhoea	Hayfever	Cough	Headache	Fever	Runny nose	Cold	#Examples with no labels
Training	1,920	0.997	106	182	163	227	251	345	375	265	530
Test	640	0.933	24	64	46	80	77	93	123	90	195

Table 1: MedWeb dataset overview statistics.

Pseudo-tweet

	Pseudo-tweet	Labels
(ja)	風邪を引くと全身がだるくなる。	Cold
(en)	The cold makes my whole body weak.	
(zh)	一感冒就身酸无力。	
(ja)	アトピーと花粉症が重なってつらい	Hay fever & Runny nose
(en)	It's really bad. My eczema and allergies are acting up at the same time.	
(zh)	敏症加花粉症, 受死了。	
(ja)	今日インフルの手術じゃないただの注射なのにビビる	No labels
(en)	I'm so scared of today's flu shot, and it's not even surgery or anything.	
(zh)	今天只打不做流感手, 但是害怕。	

Table 2: Example pseudo-tweet triplets.

Results

Model	Source	Train	Test	Exact Match Accuracy	F1 macro
Baselines					
Majority class classifier	-	-	-	0.305	-
Random classifier	-	-	-	0.130 (0.012)	0.118 (0.007)
Iso et al. [1]	-	EN	EN	0.795	-
Iso et al. [1]	-	JA	JA	0.825	-
Iso et al. [1]	-	ZH	ZH	0.809	-
BERT	-	EN	EN	0.847 (0.003)	0.884 (0.004)
jBERT	-	JA	JA	0.843 (0.012)	0.880 (0.006)
mBERT	-	ZH	ZH	0.835 (0.004)	0.876 (0.006)
Zero-shot transfer					
mBERT	-	EN	JA	0.305 (0.001)	-
mBERT	-	ZH	JA	0.507 (0.007)	0.484 (0.032)
Machine translation					
mBERT	EN	TJA	JA	0.740 (0.011)	0.740 (0.012)
mBERT	ZH	TJA	JA	0.774 (0.008)	0.821 (0.010)
mBERT	EN	TJA (x2)	JA	0.754 (0.009)	0.758 (0.034)
mBERT	ZH	TJA (x2)	JA	0.804 (0.004)	0.849 (0.098)

Table 3: Overall results, given as mean (standard deviation) of 5 runs, for different training/test data pairs using English (EN), Japanese (JA), Chinese (ZH) and translated Japanese (TJA).

Figure 1 (right):

t-SNE plot of max-pooled output of mBERT final layer (before fine tuning). 20 sentence triplets are linked to show the mapping between languages.

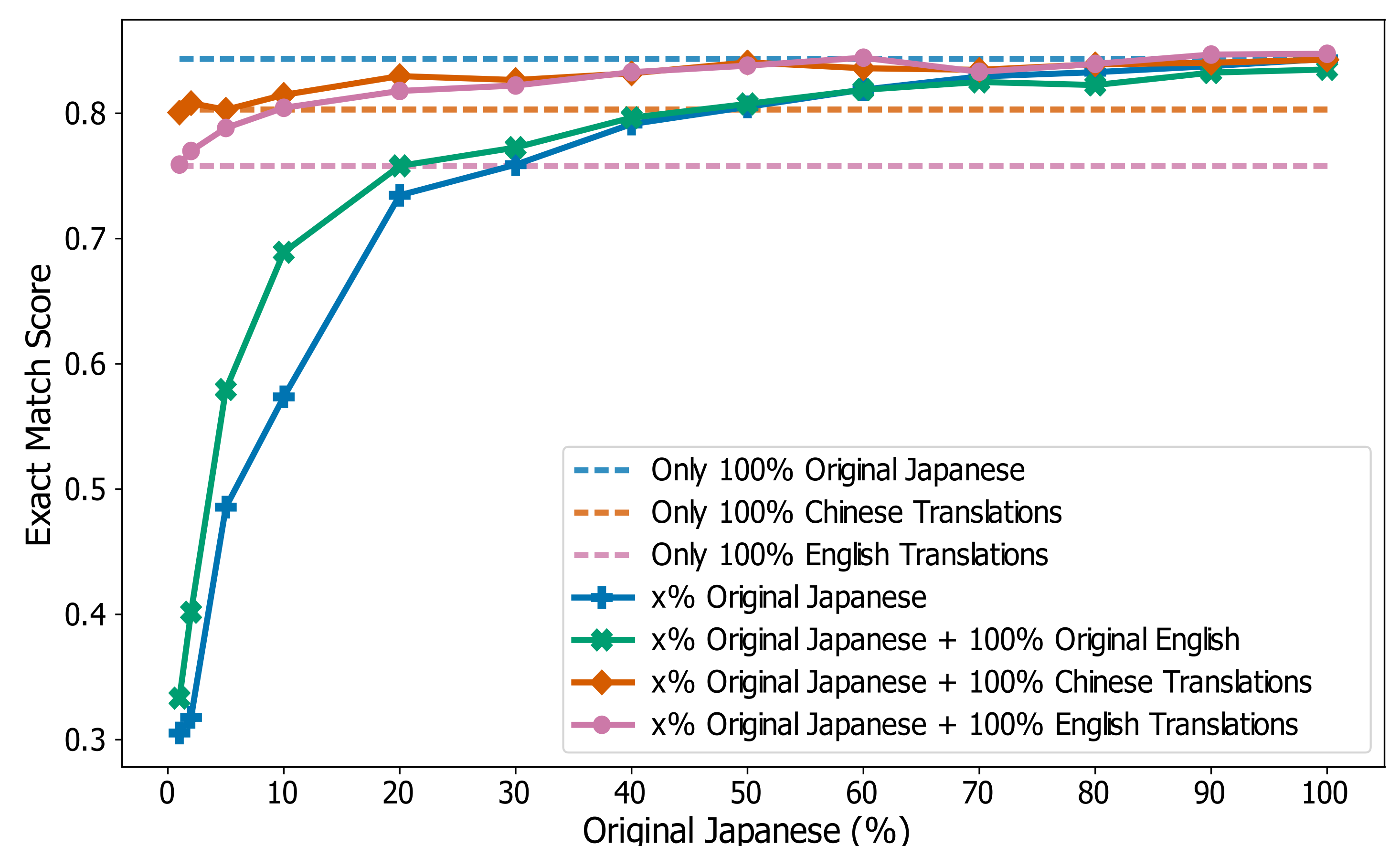
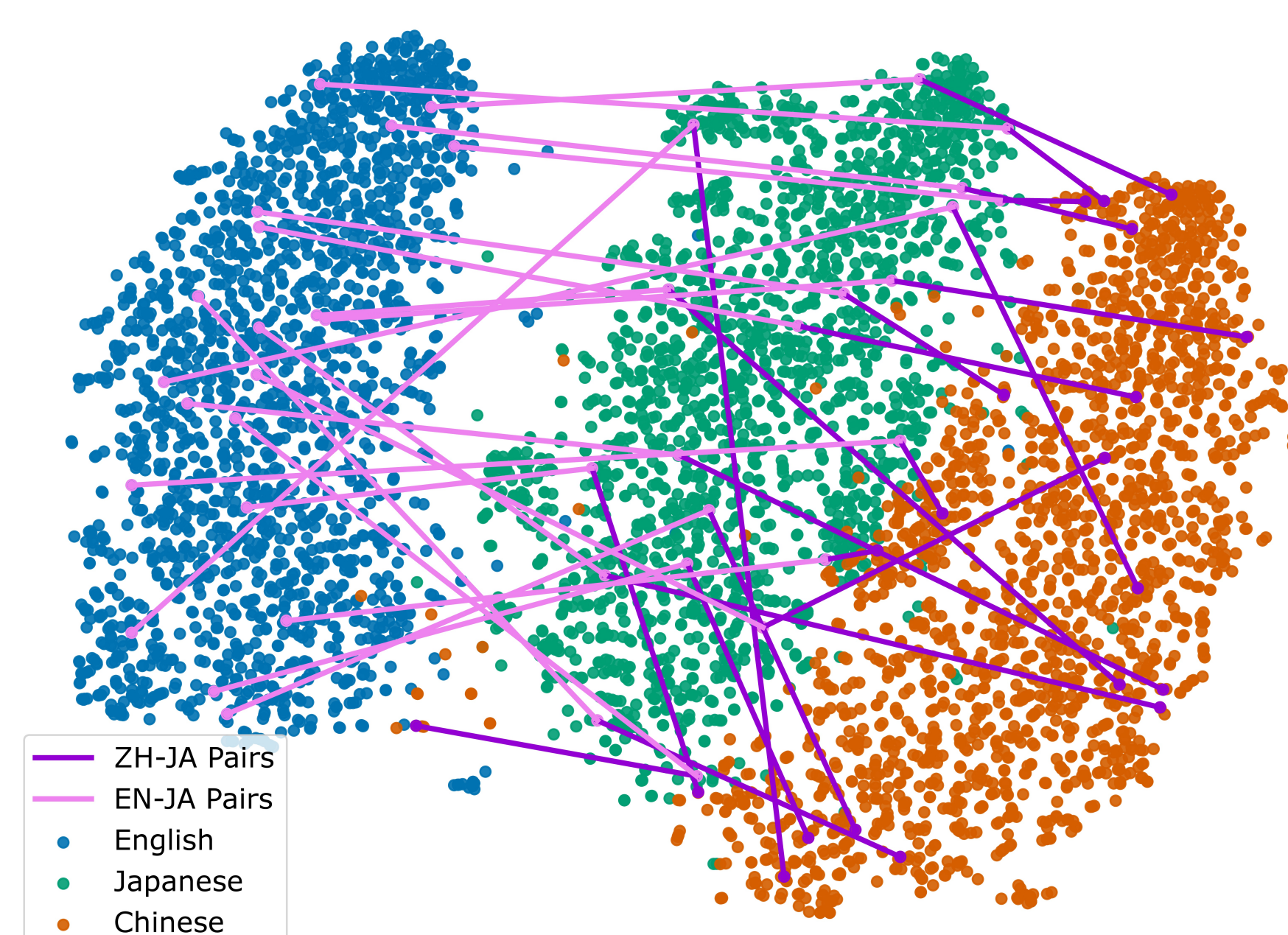


Figure 2: Exact match accuracy when training on different proportions of the original Japanese training set, with or without either the original English data or the translated data.

- Zero-shot transfer using mBERT performs poorly. t-SNE inspection (Fig. 1) shows different language representations are not shared.
- Training on machine translations gives promising performance, which is improved by adding small amounts of original target data (Fig. 2).
- Machine translations are sometimes reasonable but not consistent with labels: 風邪 = cold (the illness) is translated into 寒さ = cold (low temperature).

Conclusions

- The **choice of source language** impacts the performance, with ZH-JA being a better language pair than EN-JA.
- Training on machine translated data** shows promise, especially when used with small amounts of target language data.

References

- Iso et al., NTCIR13 MedWeb Task: Multi-label classification of Tweets using an Ensemble of Neural Networks., NTCIR, 2017
- Devlin et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding., NAACL, 2019
- Kikuta et al., BERT pre-trained model trained on Japanese Wikipedia articles.
- mBERT, <https://github.com/google-research/bert/blob/master/multilingual.md>
- Google Translate, <https://cloud.google.com/translate/>
- AWS Translate, <https://aws.amazon.com/translate/>