

# Can Machine Translation Bridge Multilingual Pretraining and Cross-lingual Transfer Learning?



Shaoxiong Ji<sup>1</sup>   Timothee Mickus<sup>1</sup>   Vincent Segonne<sup>2</sup>   Jörg Tiedemann<sup>1</sup>  
<sup>1</sup> University of Helsinki   <sup>2</sup> Universite Grenoble Alpes  
firstname.lastname@helsinki.fi

## Research questions

This paper investigates whether machine translation as a learning objective can improve performances on zero-shot cross-lingual transfer downstream tasks. We attempt to establish whether MT training objectives implicitly foster cross-lingual alignment:

- (i) Do models (re)trained with the MT objective develop cross-lingual representations?
- (ii) Do they generalize well on cross-lingual tasks?
- (iii) Which factors impact their performances?

## Findings

- MT (continued) training objectives do not favor the emergence of cross-lingual alignments more than LM objectives, based on the study on existing publicly available pretrained models.
- We provide evidence from similarity analyses and parameter-level investigations that this is due to separability, which is beneficial in MT but detrimental elsewhere.
- We conclude that MT encourages behavior that is not necessarily compatible with high performances in cross-lingual transfer learning.

## Representation similarity

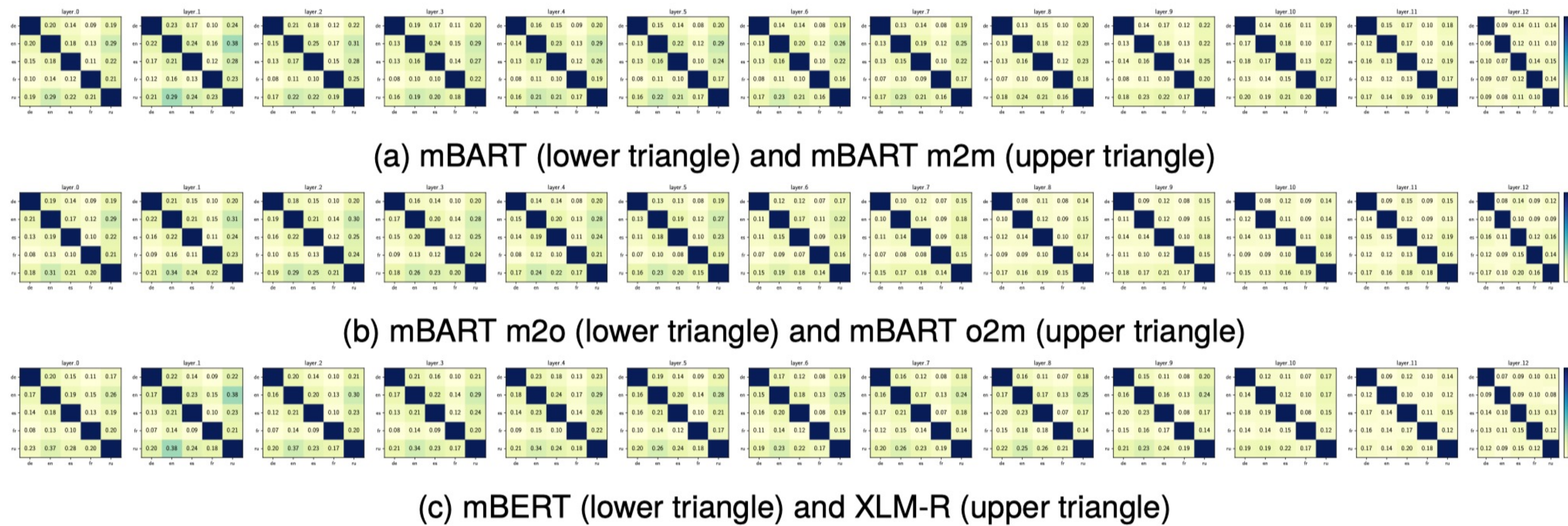


Figure 2: Representational similarity between different languages with representations learned by LMs and MT models

## CP’s effect on scaling

Model	K		Q		V		Out		FC up		FC down	
	$\ \sigma\ $	$d$	$\ \sigma\ $	$d$	$\ \sigma\ $	$d$	$\ \sigma\ $	$d$	$\ \sigma\ $	$d$	$\ \sigma\ $	$d$
mBART	44.76	—	44.85	—	53.73	—	53.45	—	90.25	—	99.63	—
mBART m2m	48.28	4.23	48.29	4.07	55.65	2.73	55.14	3.01	99.28	9.47	107.94	9.63
mBART m2o	48.34	4.23	48.35	4.06	56.19	2.95	55.73	2.99	101.06	11.19	109.71	11.18
mBART o2m	56.13	11.76	56.25	11.74	60.17	7.18	59.32	7.07	116.17	26.34	120.50	22.15

Table 2: SVD scaling effect for mBART and CP models; weight matrices from the 12th layer.

## Quantitative performance

Model	Tasks							
	NC	XNLI	PAWS-X	QAM	QADSM	WPR	NER	POS
mBERT	81.3	65.2	86.6	64.6	63.1	74.4	77.5	76.0
LM XLM-R	<b>82.1</b>	<b>73.5</b>	<b>88.9</b>	<b>67.4</b>	<b>66.9</b>	<b>75.3</b>	<b>78.7</b>	<b>79.7</b>
mBART	82.1	67.6	<b>89.2</b>	<b>67.8</b>	65.5	74.7	77.7	72.7
MT NLLB 600M	76.0	68.3	73.4	61.5	63.9	73.7	54.2	71.4
mBART m2o	80.4	65.9	85.6	63.9	63.9	73.7	61.5	70.8
CP mBART o2m	65.4	48.1	81.7	58.4	62.7	73.2	55.1	55.7
mBART m2m	78.3	60.2	87.2	63.2	62.8	73.7	71.9	69.7

Table 1: Average performance on cross-lingual tasks. We use the base architecture for mBERT and XLM-R. mBART scores are derived from the 12-layer encoder.

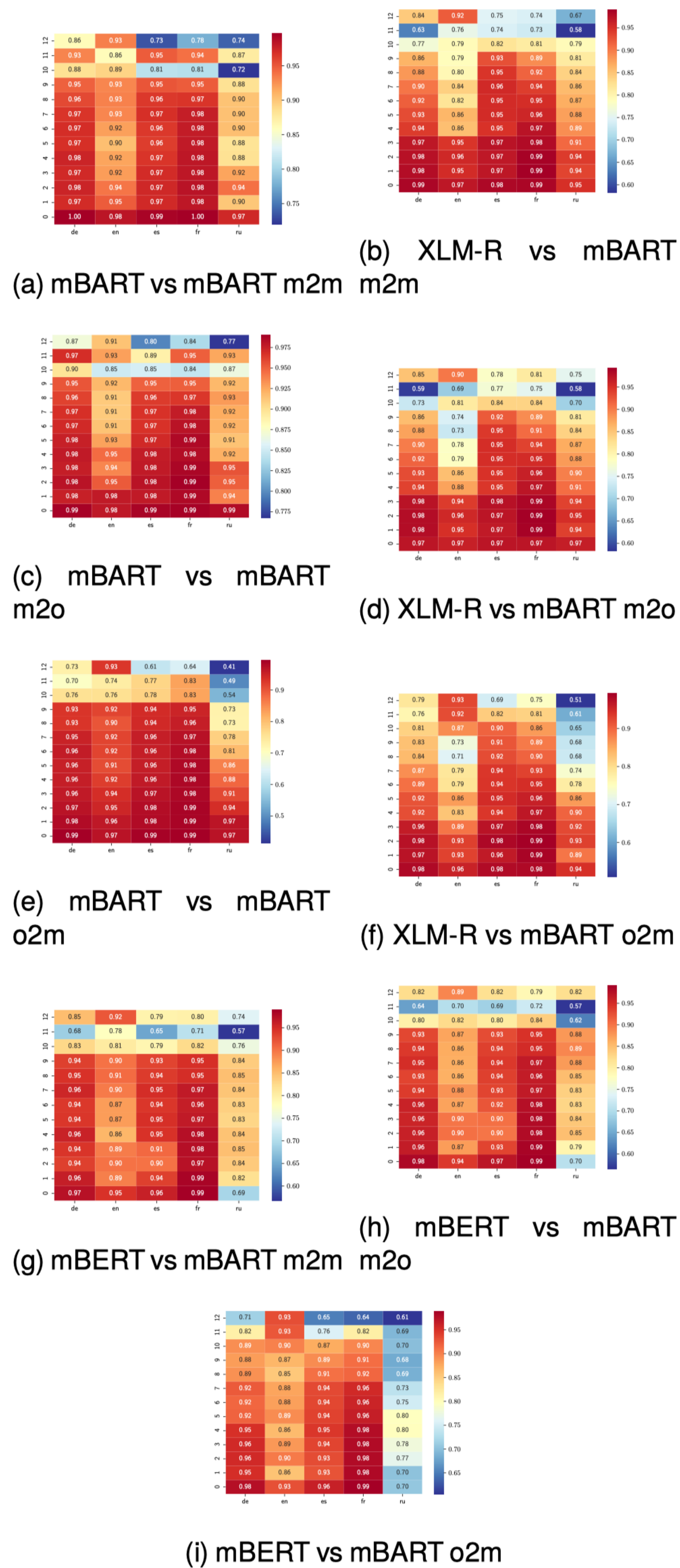


Figure 1: Representational similarity between mBART-based MT models and LMs