

An Operation Sequence Model for Explainable Neural Machine Translation



UNIVERSITY OF CAMBRIDGE

Felix Stahlberg, Danielle Saunders, Bill Byrne

Department of Engineering

Introduction

- Achieve explainable AI by changing the output representation to explain itself
- For NMT: Output sequence conveys both the translation and an “explanation” in form of alignments
- Potentially works on any Seq2Seq architecture

Operation Sequence NMT (OSNMT)



	Operation	Source sentence	Target sentence (compiled)
		2000 hr の安定動作を確認した	X_1
1	SET_MARKER	2000 hr の安定動作を確認した	$X_2 X_1$
2	2000	2000 hr の安定動作を確認した	X_2 2000 X_1
3	POP_SRC	2000 hr の安定動作を確認した	X_2 2000 X_1
4	hr	2000 hr の安定動作を確認した	X_2 2000 hr X_1
5	POP_SRC	2000 hr の安定動作を確認した	X_2 2000 hr X_1
6	JMP_BWD	2000 hr の安定動作を確認した	X_2 2000 hr X_1
7	SET_MARKER	2000 hr の安定動作を確認した	$X_3 X_2$ 2000 hr X_1
8	of	2000 hr の安定動作を確認した	X_3 of X_2 2000 hr X_1

SET_MARKER: Insert new marker X_i into target sentence

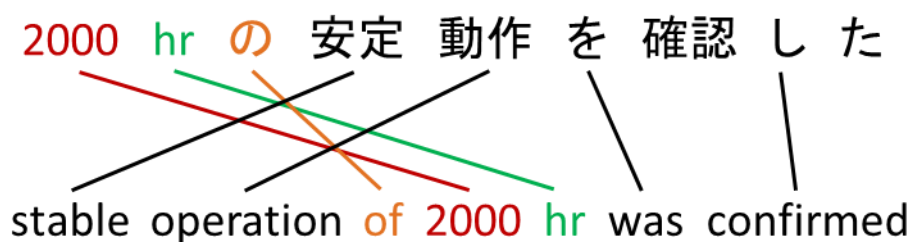
word: Insert word into target sentence

POP_SRC: Move source read head right by 1 token.

JMP_FWD|BWD: Move target write head to next marker.

Source side read head

Target side write head



- Generates target words in source sentence order, interspersed with reordering instructions
- Interpretable behaviour of Transformer attention heads
- OSNMT output can be seen as a parse through a formal (multitext) grammar

Representation	BLEU			
	es-en	pt-en	ja-en	
			dev	test
Plain	37.6	37.5	28.3	28.1
OSNMT	37.1	38.4	28.1	28.8

AER on ja-en against GIZA++ alignments

Method	BLEU	AER
RNN (max attention)	20.2	61%
Transformer (OSNMT)	28.1	20%

