

# PheMT: A Phenomenon-wise Dataset for Machine Translation Robustness on User-Generated Contents

Ryo Fujii<sup>1</sup>, Masato Mita<sup>2,1</sup>, Kaori Abe<sup>1</sup>, Kazuaki Hanawa<sup>2,1</sup>, Makoto Morishita<sup>3,1</sup>, Jun Suzuki<sup>1,2</sup>, Kentaro Inui<sup>1,2</sup>  
 1. Tohoku University 2. RIKEN 3. NTT Communication Science Laboratories

## Summary

- A new dataset for evaluating the robustness of Japanese-to-English MT systems on UGC
- Provide focused evaluation on four linguistic phenomena with the idea of contrastive datasets
- Evaluated the effect of the phenomena with both in-house and widely used off-the-shelf systems
- Discovered a unique preprocessing method towards improving the performance on *Variant*

## Background

- UGC are prevailing in our real-life communication
  - e.g., social media, blog posts, user reviews
- A shared task on machine translation robustness <sup>[1]</sup>

More attention towards handling UGC to promote cross-cultural communication



The performance of current MT systems on UGC is still far behind

## Q. Why is it difficult to translate UGC ?

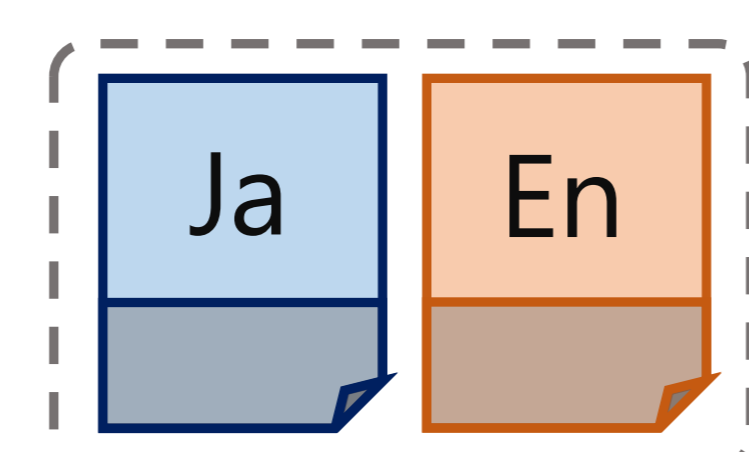
Still not clear...

**We need a solid basis for more detailed analysis !**

[1] Li et al. (2019), Findings of the first shared task on machine translation robustness.  
 [2] Michel and Neubig (2018), MTNT: A Testbed for Machine Translation of Noisy Text.

## Creating phenomenon-wise dataset

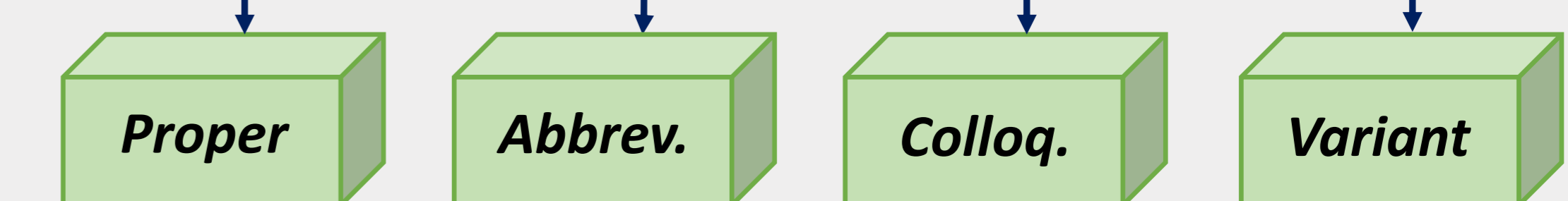
MTNT corpus <sup>[2]</sup>



Label *	Examples { Orig. / Norm. }
<i>Proper Noun</i>	安倍首相 ( <i>abeshushō</i> , meaning PM Abe)
<i>Abbreviated Noun</i>	{ アプデ / アップデート } ( <i>apude</i> , meaning update)
<i>Colloquial Expression</i>	{ かなちい / かなしい } ( <i>kanachii / kanashii</i> , meaning sad)
<i>Variant</i>	{ アリガトウ / ありがとう } ( <i>arigatou</i> , meaning thank you)

\* Please refer to the paper for the definition

Step1: Annotating phenomena labels



Step2: Extracting targeted expressions / alignments

Ja: 地味な **アプデ** (*apude*, meaning update, abbreviated) だが  
 En: That's a plain **update** though

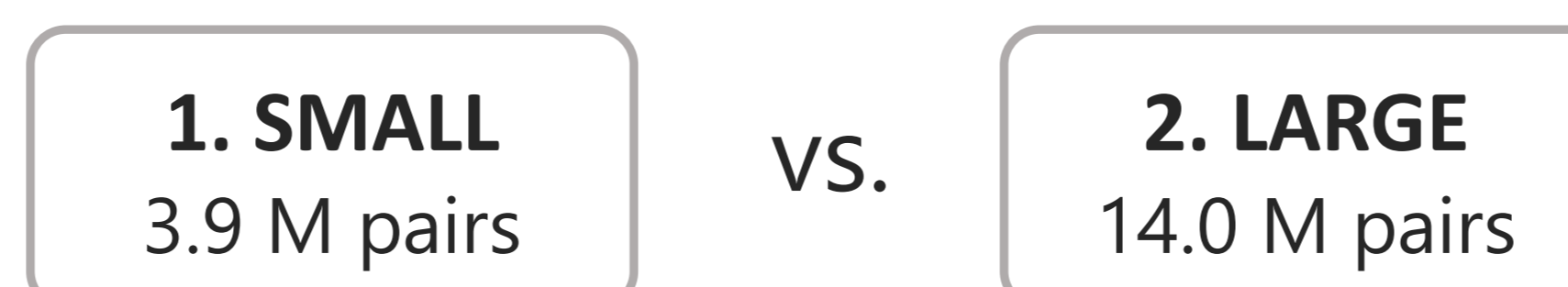
Step3: Normalizing the expressions (except *Proper Noun*)

Orig.: 地味な **アプデ** (*apude*) だが  
 Norm.: 地味な **アップデート** (*update*, canonical) だが

## Translation models

- The five **in-house models** :

Q1. Effect of training data size ?

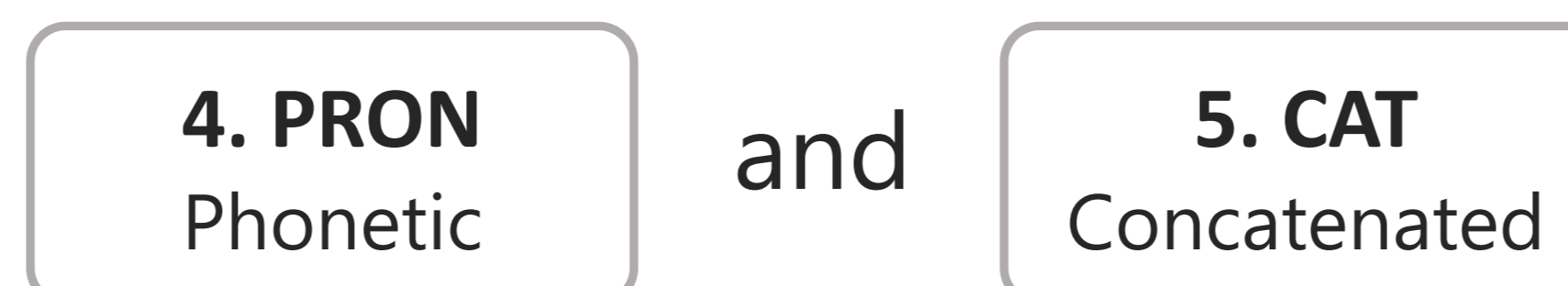


Q2. Effect of tokenization ?



Q3. Susceptible to local improvement ?

Trained on a fully-pronunciation based corpus to absorb symbolic differences in *Variant*



- **Off-the-shelf systems** : Google, DeepL

## Phenomenon-wise evaluation

Our robustness measure :

**The difference of arbitrary metrics** for (Orig. / Norm.) input

Translation accuracy with extracted alignment (raw acc. only for *Proper*)

	SMALL	LARGE	CHAR	PRON	CAT	Google	DeepL
<i>Proper</i>	(34.3)	(49.7)	(47.1)	(43.2)	(48.0)	(55.2)	(50.5)
<i>Abbrev.</i>	+6.4 (24.1 / 30.5)	-0.6 (33.6 / 33.0)	+0.6 (34.2 / 34.8)	+1.1 (30.2 / 31.3)	-1.2 (34.2 / 33.0)	-4.3 (41.1 / 36.8)	-1.2 (39.1 / 37.9)
<i>Colloq.</i>	+5.8 (18.0 / 23.8)	+9.9 (14.5 / 24.4)	+4.1 (17.4 / 21.5)	+21.5 (8.7 / 30.2)	+16.9 (15.7 / 32.6)	+7.0 (19.2 / 26.2)	+5.8 (22.7 / 28.5)
<i>Variant</i>	+19.5 (15.5 / 35.0)	+25.2 (13.6 / 38.8)	+20.4 (13.6 / 34.0)	+10.7 (25.2 / 35.9)	+8.8 (26.2 / 35.0)	+14.6 (23.3 / 37.9)	+16.6 (18.4 / 35.0)

- A1. High coverage with larger training data was effective for nouns, while not for UGC-specific phenomena
- A2. Char-based tokenization worked well with *Colloq.*, which share most of the characters with their canonical forms
- A3. Our dataset could detect the improvement against *Variant*, which was proven to be more problematic to current systems