

Overview of WAT2014

Toshiaki Nakazawa
nakazawa@pa.jst.jp

Hideya Mino
hideya.mino@nict.go.jp

Isao Goto
goto.i-es@nhk.or.jp

Sadao Kurohashi
kuro@i.kyoto-u.ac.jp

Eiichiro Sumita
eiichiro.sumita@nict.go.jp

Participants List

Team ID	Organization	J->E	E->J	J->C	C->J
NAIST	Nara Institute of Science and Technology	✓	✓	✓	✓
EIWA	Yamanashi Eiwa College	✓	✓	✓	✓
Kyoto-U	Kyoto University	✓	✓	✓	✓
WEBLIO-EJ1	Weblio, Inc.		✓		
TMU	Tokyo Metropolitan University	✓			
BJTUNLP	Beijing Jiaotong University			✓	
NII	National Institute of Informatics	✓			
SAS_MT	SAS Research and Development Co., Ltd		✓	✓	✓
Sense	Saarland University & Nanyang Technological University	✓	✓	✓	✓
NICT	National Institute of Informatics and Communication Technology			✓	
TOSHIBA	Toshiba Corporation	✓		✓	✓
WASUIPS	Waseda University			✓*	✓*

company outside Japan * Only submitted to the automatic evaluations

HUMAN score

<http://www.lancers.jp>

* Pairwise evaluation compared to the baseline using crowdsourcing
- reduce the cost, time, number of sentences to be evaluated
- enable the evaluation of new translation results after the workshop
* Reference translations are not shown

2つの機械翻訳結果の優劣判断

科学技術論文の英語入力文に対する日本語の機械翻訳結果が2つ表示されています。どちらの翻訳がより正しいかを判断してください。優劣がつかれない場合は、同程度としてください。

入力文: Details of dose rate of "Fugen Power Plant" can be calculated by using DERS software.
翻訳文1: みげん発電所の線量率の詳細はDERSソフトウェアを用いて計算することができます。
翻訳文2: "みげん発電所の線量率"の詳細を用いて計算することができる。"DERS"ソフトウェアである。
○ 1つ目の翻訳の方が良い ○ 2つ目の翻訳の方が良い ○ 同程度

* To guarantee the quality of the evaluation, each pair is evaluated 3 different workers
* The evaluation result of each pair is decided by the voting of 3 judgments

Worker 1	A	A	A	A	A	Tie	Tie	Tie	B
Worker 2	A	A	A	Tie	Tie	B	Tie	Tie	B
Worker 3	A	Tie	B	Tie	B	B	Tie	B	B
Decision	A	A	A	A	Tie	B	B	Tie	B

* 400 sentences are randomly selected (paragraph-based) from the test set

BASELINE (Phrase-based SMT) VS. SYSTEM1

$HUMAN = 100 \times \frac{Wins - Losses}{Wins + Losses + Ties}$

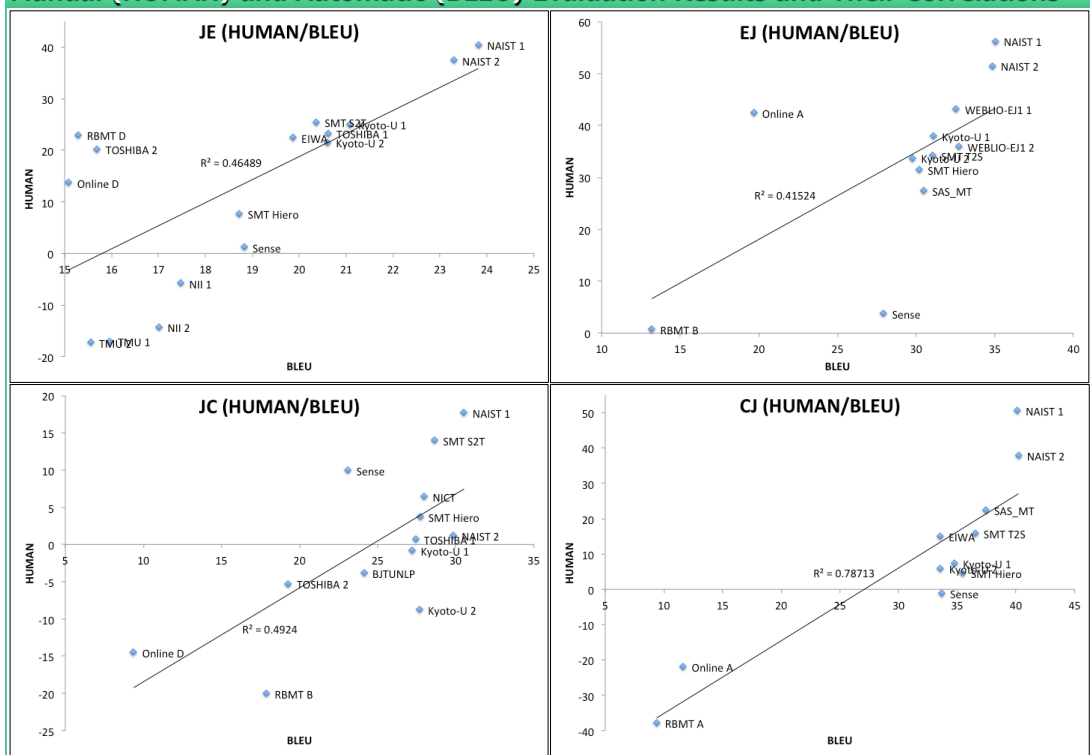
Cost

one judgment by a worker = 5 JPY
×
3 judgments for one decision
×
400 sentences
||
6000 JPY / submission

Time

Depending on the translation direction
On the average, one evaluation finished in a couple of days

Manual (HUMAN) and Automatic (BLEU) Evaluation Results and Their Correlations



Inter Annotator Agreement

JE	EJ	JC	CJ
System ID	System ID	System ID	System ID
Kappa	Kappa	Kappa	Kappa
NAIST 1	NAIST 1	NAIST 1	NAIST 1
NAIST 2	NAIST 2	SMT S2T	NAIST 2
SMT S2T	WEBLIO-EJ1	Sense	SAS_MT
Kyoto-U 1	Online A	NICT	SMT T2S
TOSHIBA 1	Kyoto-U 1	SMT Hiero	EIWA
RBMT D	WEBLIO-EJ1 2	NAIST 2	Kyoto-U 1
EIWA	SMT T2S	TOSHIBA 1	Kyoto-U 2
Kyoto-U 2	Kyoto-U 2	Kyoto-U 1	SMT Hiero
TOSHIBA 2	SMT Hiero	BJTUNLP	Sense
Online D	SAS_MT	TOSHIBA 2	Online A
SMT Hiero	Sense	Kyoto-U 2	RBMT A
Sense	RBMT B	Online D	Online B
NII 1		RBMT B	
NII 2			
TMU 1			
TMU 2			
Ave. 0.106	Ave. 0.254	Ave. 0.101	Ave. 0.200

X->J evaluations are easier than J->X evaluations (because the workers are almost Japanese?)

Case Study

Submission ID	BLEU	RIBES	HUMAN	Forest Input
WEBLIO-EJ1 1	32.53	0.782	43.25	No
WEBLIO-EJ1 2	32.69	0.785	36.00	Yes

BASELINE (WEBLIO-EJ1) VS. WEBLIO-EJ1 2

HUMAN 2.50 ± 4.17

Kappa 0.528

* No significant difference
* Much higher Kappa value
similar outputs can be easily and faithfully judged

Conclusion and Future Perspective

- 12 participants for the evaluation task - including 3 companies and 3 outside Japan
- NAIST team achieved the best results for all the subtasks
- Shared the findings of MT for scientific papers
- Need more investigation to acquire reliable human evaluation results at low cost
- Need to find a better way to compare two systems efficiently and reliably
- Automatic evaluation server will keep running even after the workshop to promote continuous evolution of MT research show: <http://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/> submit: <http://lotus.kuee.kyoto-u.ac.jp/WAT/submission/>
- KEEP ON UPDATING THE STATUS OF YOUR MT ENGINE!!
- WAT will be held annually
- The importance of both sub/super sentential context
- Let's share your resource monolingual/bilingual corpora, dictionaries and so on

Thank you very much for attending WAT2014
Stay tuned for the next WAT workshop!