# Overview of WAT2014

Toshiaki Nakazawa    Hideya Mino

Isao Goto    Sadao Kurohashi

Eiichiro Sumita

4th October, 2014 @ WAT2014

# WAT 2014
## The 1st Workshop on Asian Translation

- MT evaluation campaign focusing on Asian languages (Japanese, Chinese and English for this time)

- The first evaluation workshop that uses scientific papers as a domain and Japanese-Chinese as a language pair

- Paragraph-based test set
  - investigate the viability of the context-aware MT

- All the data including test set are OPEN
  - contribute to continuous evolution of MT research by freely distributing the data (like PennTreebank sec. 23)

# Automatic Evaluation in WAT2014

- Prepared an automatic evaluation server
- BLEU, RIBES
- several word segmentation tools

See Evaluation Results:

   http://lotus.kuee.kyoto-u.ac.jp/WAT/evaluation/

Submit Your Translations (need FREE registration):

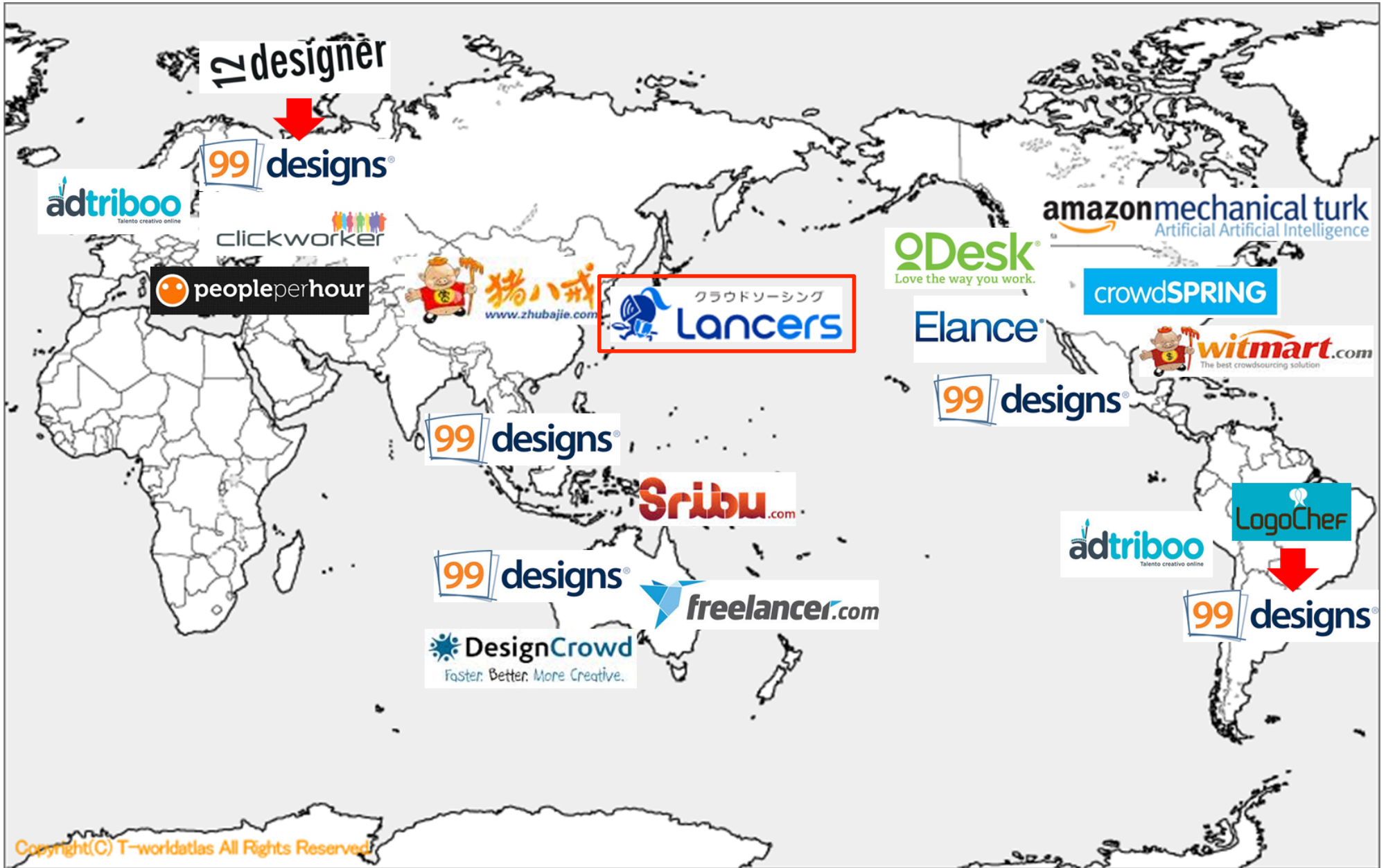   http://lotus.kuee.kyoto-u.ac.jp/WAT/submission/

# Human Evaluation of MT

- Costs a lot of money and time
- Unstable results caused by the different criteria of each evaluators
- Many measures
  - Adequacy/Fluency
  - Ranking
  - Acceptability (NTCIR)
  - Patent Examination Evaluation (NTCIR)

# Human Evaluation in WAT2014

- Costs a lot of money and time
  - using crowdsourcing to reduce them
- Unstable results caused by the different criteria of each evaluators
  - alleviate the divergence by voting
- Many measures
  - HUMAN score

http://www.crowdinfo.jp/2014/02/02/world-crowdsourcing-service/

# Human Evaluation in WAT2014

Phrase-based SMT

- Pairwise evaluation compared to the baseline
  - reduce the number of sentences to be evaluated
  - enable the evaluation of new translation results after the workshop
- 400 sentences selected from the test set by document-based sampling
- Reference translations are not shown

# Sample of the Task



- The order of the baseline and subject outputs are at random

# Pairwise Evaluation by Voting

- 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
  - e.g. MT A vs. MT B

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Worker 1** | A | A | A | A | A | A | Tie | Tie | Tie | B |
| **Worker 2** | A | A | A | Tie | Tie | B | Tie | Tie | B | B |
| **Worker 3** | A | Tie | B | Tie | B | B | Tie | B | B | B |
| **Decision** | A | A | A | A | Tie | B | Tie | B | B | B |

# Comparison to the Baseline

BASELINE-MT-1          SYSTEM1-MT-1
BASELINE-MT-2          SYSTEM1-MT-2
BASELINE-MT-3          SYSTEM1-MT-3
BASELINE-MT-4          SYSTEM1-MT-4
BASELINE-MT-5    VS.   SYSTEM1-MT-5
BASELINE-MT-6          SYSTEM1-MT-6
BASELINE-MT-7          SYSTEM1-MT-7
BASELINE-MT-8          SYSTEM1-MT-8
BASELINE-MT-9          SYSTEM1-MT-9
BASELINE-MT-10         SYSTEM1-MT-10

BASELINE
(Phrase-based SMT)

SYSTEM1

5 wins, 2 losses, 3 ties

# Human Evaluation Score

- Suppose *W* = # of wins, *L* = # of losses and *T* = # of ties, the HUMAN score is

$$HUMAN = 100 \times \frac{W - L}{W + L + T}$$

e.g. sample of the previous page

$$100 \times \frac{5-2}{5+2+3} = 30$$

- Estimate the confidence interval by bootstrap resampling [Koehn, 2004]
  - calculate the human evaluation score on 300 sentences randomly sampled from 400 sentences
  - iterate calculation 1000 times
  - sort the 1000 scores and discard top and bottom 25 scores to get the 95% confidence interval

# Cost of Crowdsourcing

- One judgment by one worker costs <u>5 JPY</u>
- Each sentence requires <u>3 judgments</u>
- We have <u>400 sentences</u> for the human evaluation


- One evaluation of one submission costs

$$5 \times 3 \times 400 = \color{red}{6,000 \text{ JPY}}$$

# OFFICIAL HUMAN EVALUATION RESULTS

# Participants List

| Team ID | J->E | E->J | J->C | C->J |
|---------|------|------|------|------|
| NAIST | ✓ | ✓ | ✓ | ✓ |
| EIWA | ✓ | | | ✓ |
| Kyoto-U | ✓ | ✓ | ✓ | ✓ |
| WEBLIO-EJ1 | | ✓ | | |
| TMU | ✓ | | | |
| BJTUNLP | | | ✓ | |

| Team ID | J->E | E->J | J->C | C->J |
|---------|------|------|------|------|
| NII | ✓ | | | |
| SAS_MT | | ✓ | | ✓ |
| Sense | ✓ | ✓ | ✓ | ✓ |
| NICT | | | ✓ | |
| TOSHIBA | ✓ | | ✓ | |
| WASUIPS | | | ✓ * | ✓ * |

\* Only submitted to the automatic evaluations

Company

Outside Japan

# JE HUMAN score

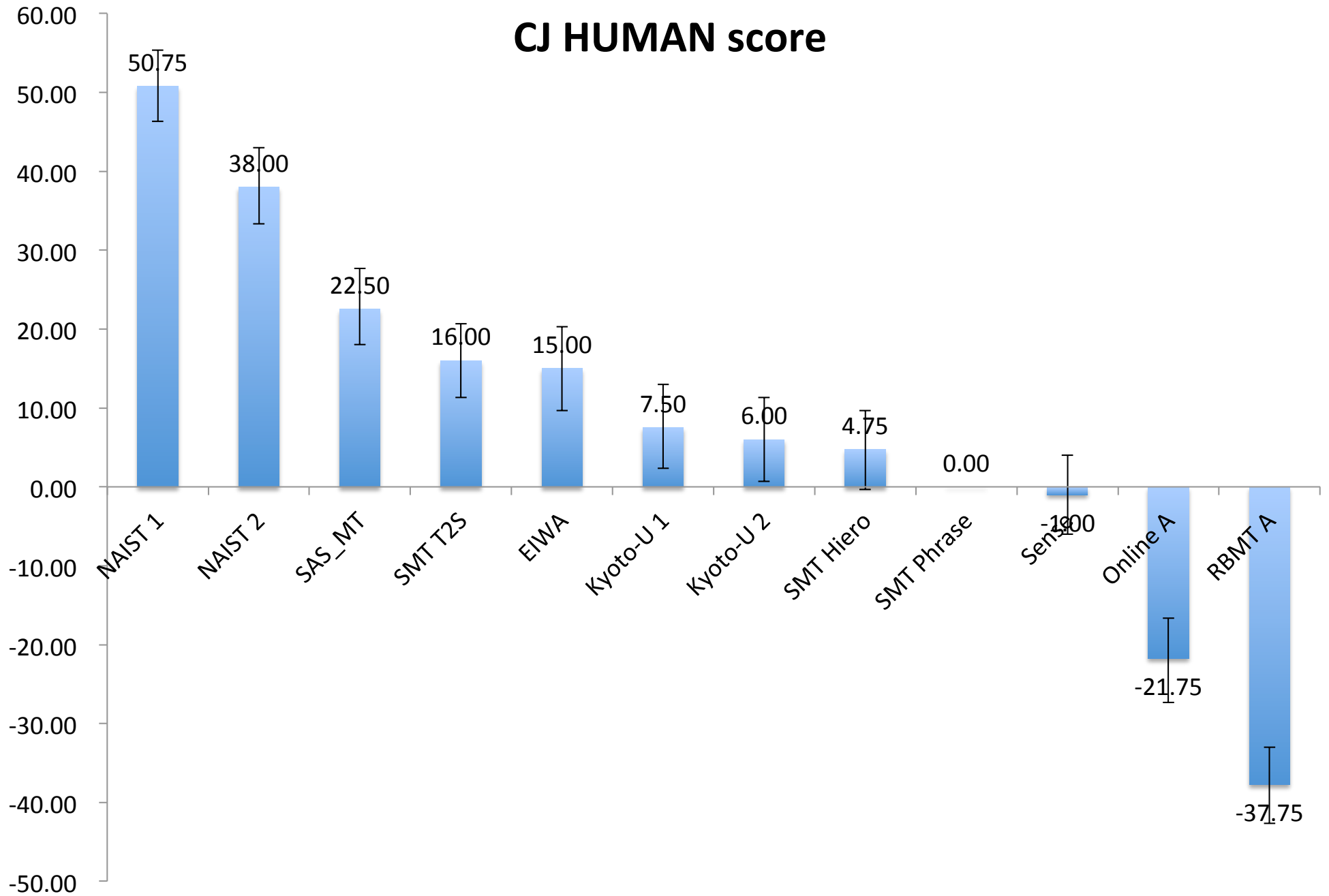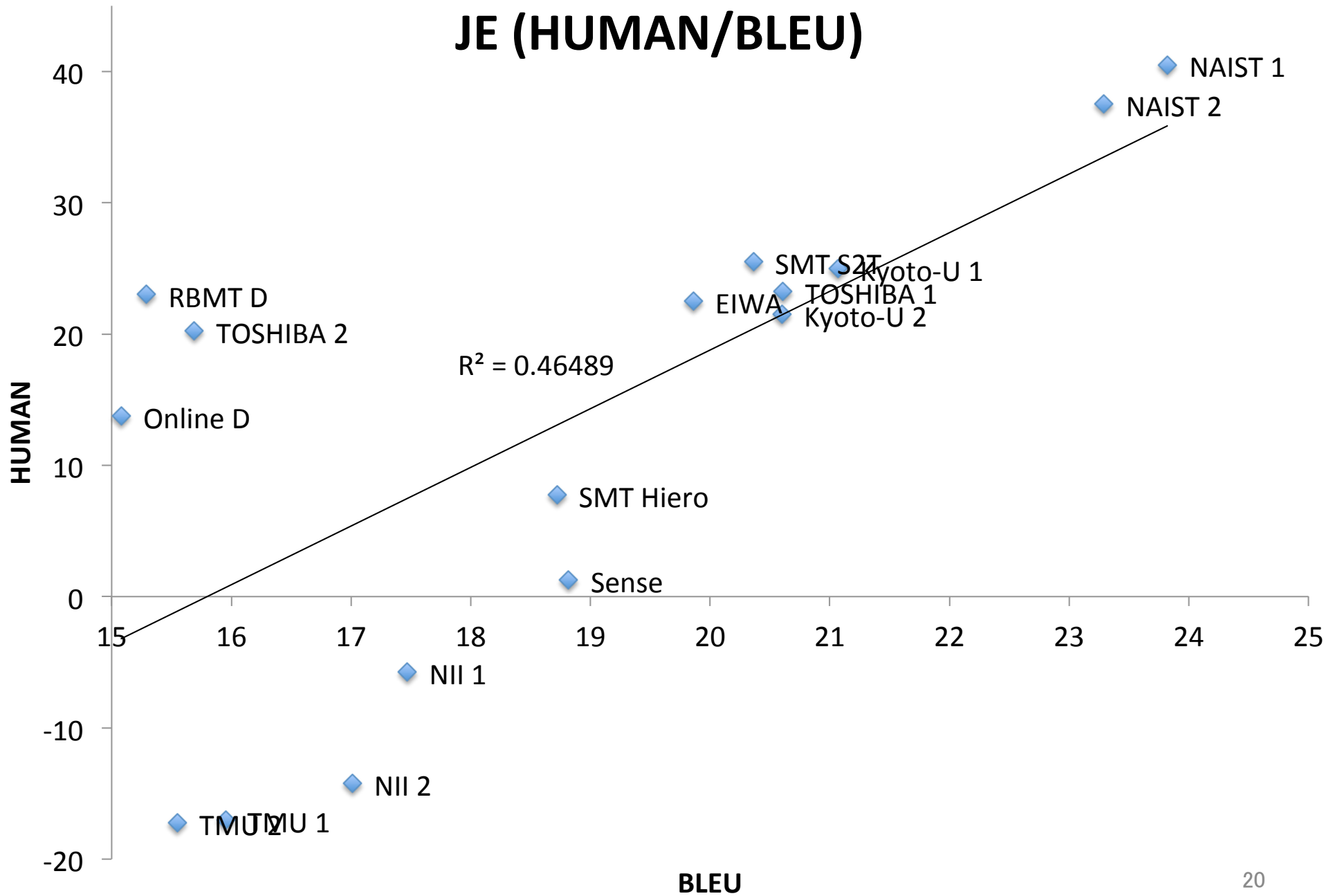| Label | Value |
|---|---|
| NAIST 1 | 40.50 |
| NAIST 2 | 37.50 |
| SMT S2T | 25.50 |
| Kyoto-U 1 | 25.00 |
| TOSHIBA 1 | 23.25 |
| RBMT D | 23.00 |
| EIWA | 22.50 |
| Kyoto-U 2 | 21.25 |
| TOSHIBA 2 | 20.25 |
| Online D | 13.75 |
| SMT Hiero | 7.75 |
| Sense | 1.25 |
| SMT Phrase | 0.00 |
| NII 1 | -5.75 |
| NII 2 | -14.25 |
| TMU 1 | -17.00 |
| TMU 2 | -17.25 |

15

**EJ HUMAN score**
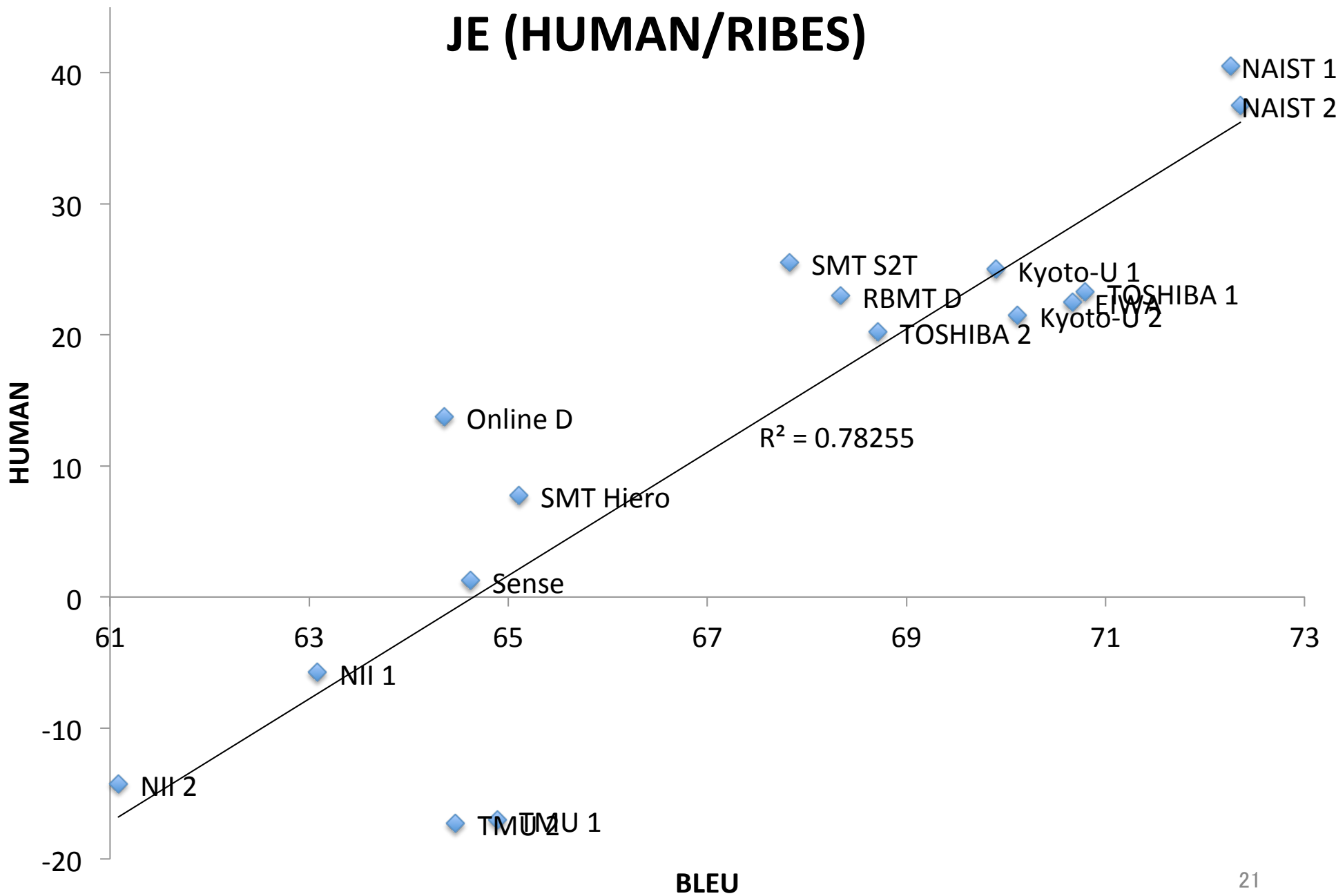
JC HUMAN score

CJ HUMAN score

# CORRELATION BETWEEN BLEU/ RIBES AND HUMAN SCORE
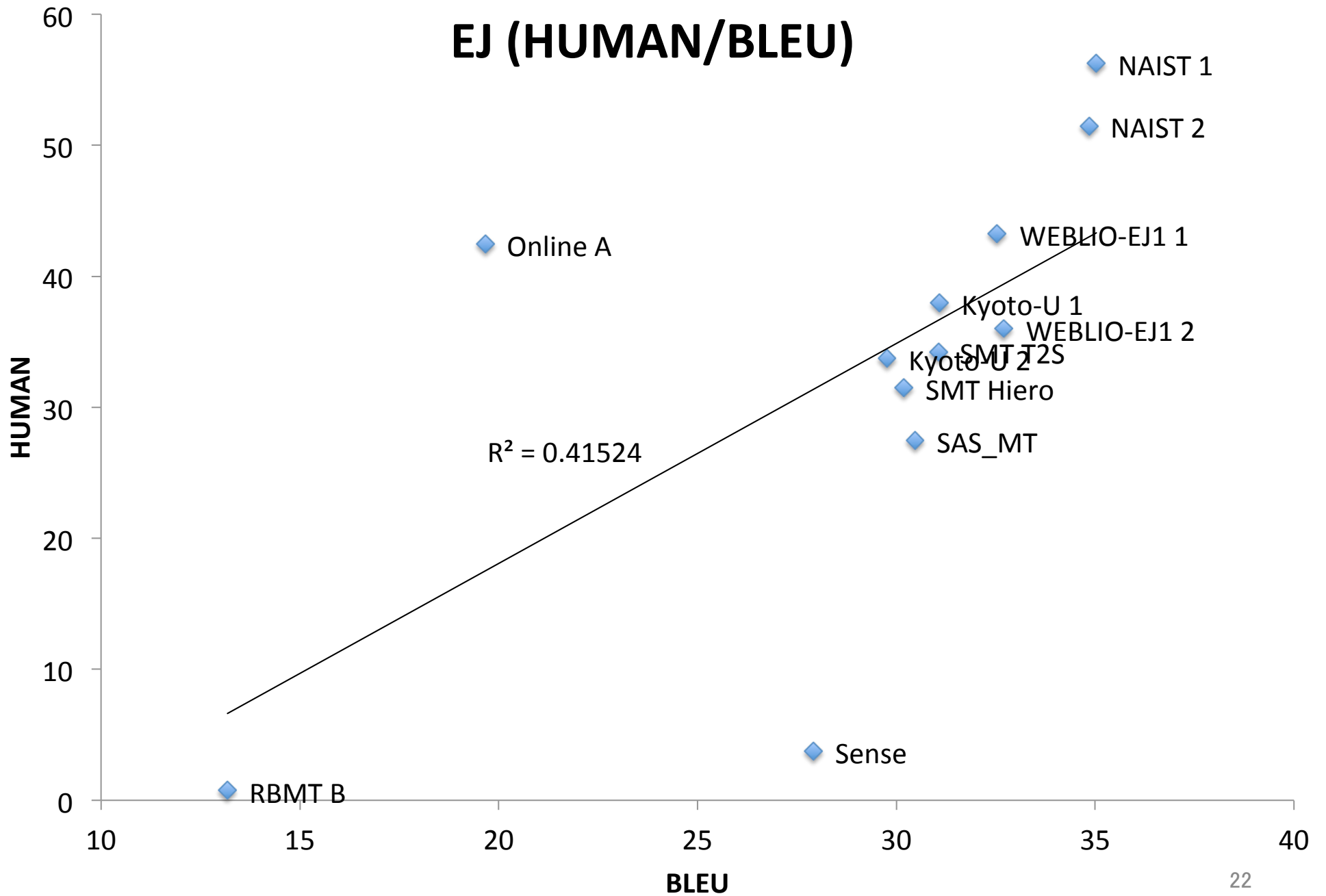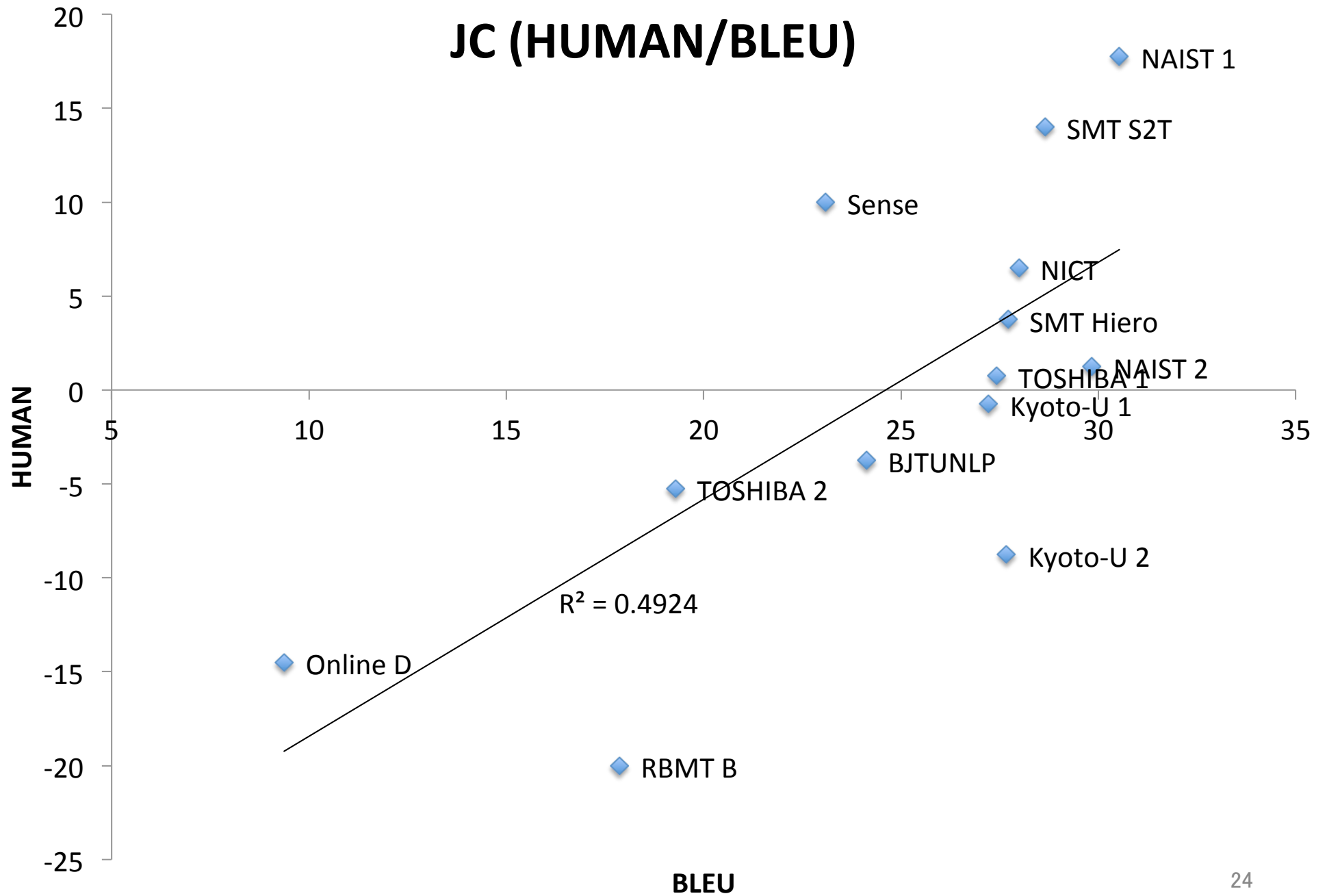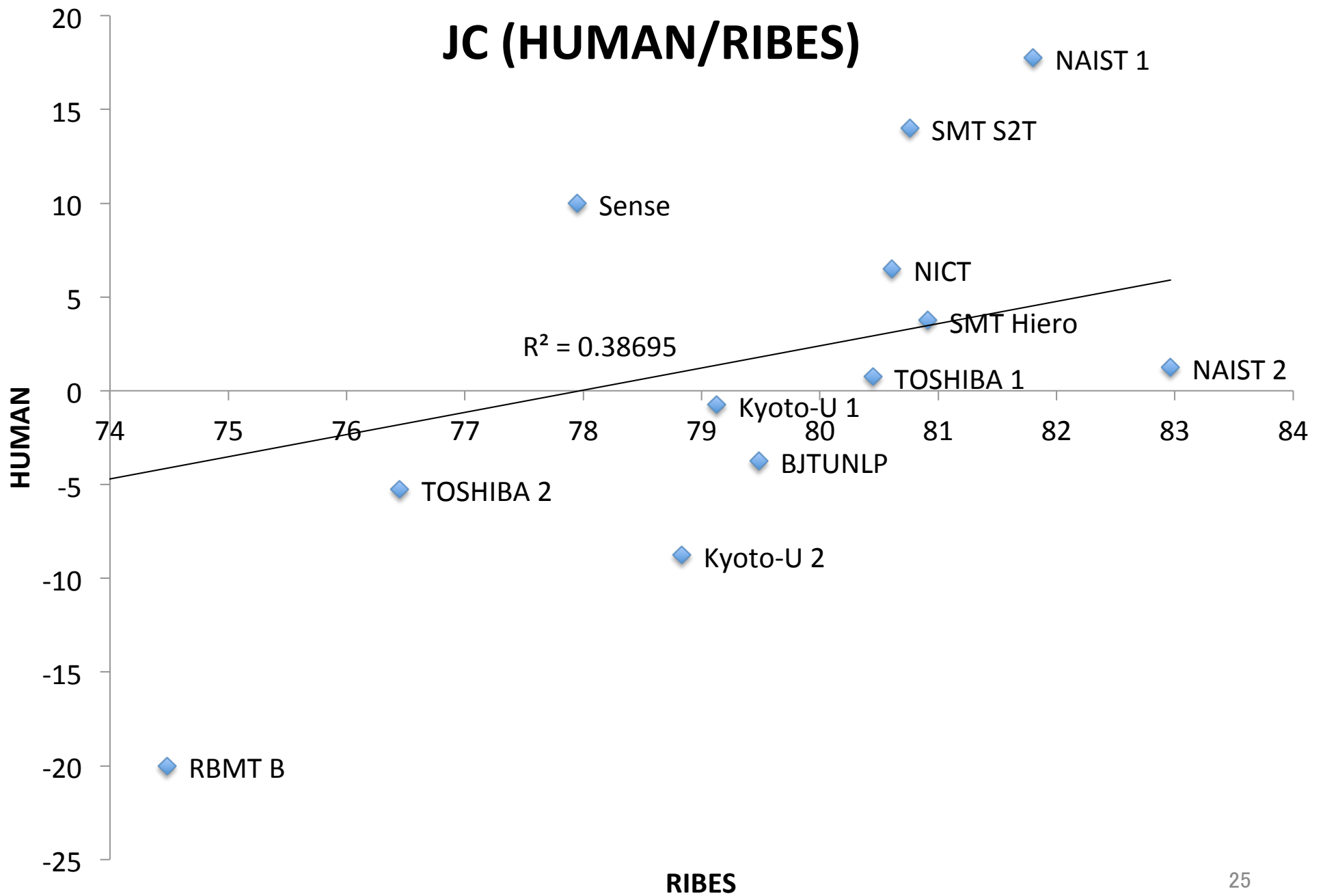
JE (HUMAN/BLEU)

JE (HUMAN/RIBES)

EJ (HUMAN/BLEU)

$R^2 = 0.41524$

EJ (HUMAN/RIBES)

JC (HUMAN/BLEU)

JC (HUMAN/RIBES)

$R^2 = 0.38695$

CJ (HUMAN/BLEU)

CJ (HUMAN/RIBES)

HUMAN

RIBES

$R^2 = 0.70081$

NAIST 1
NAIST 2
SAS_MT
EIWA SMT T2S
Kyoto-U 1
Kyoto-U 2 SMT Hiero
Sense
Online A
RBMT A

# Better Correlation among Corpus-based MT?

- Less correlations between automatic and human evaluations for RBMT and Online

| | All systems | Corpus-based only |
|---|---|---|
| **JE BLEU** | 0.46489 | 0.95098 |
| **JE RIBES** | 0.78255 | 0.83691 |
| **EJ BLEU** | 0.41524 | 0.84418 |
| **EJ RIBES** | 0.75105 | 0.85730 |
| **JC BLEU** | 0.49240 | 0.07937 |
| **JC RIBES** | 0.38695 | 0.10198 |
| **CJ BLEU** | 0.78713 | 0.82592 |
| **CJ RIBES** | 0.70081 | 0.83209 |

* $R^2$ values

# INTER ANNOTATOR AGREEMENT

# Inter Annotator Agreement

| JE | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.162 |
| NAIST 2 | 0.047 |
| SMT S2T | 0.099 |
| Kyoto-U 1 | 0.070 |
| TOSHIBA 1 | 0.098 |
| RBMT D | 0.075 |
| EIWA | 0.083 |
| Kyoto-U 2 | 0.139 |
| TOSHIBA 2 | 0.078 |
| Online D | 0.055 |
| SMT Hiero | 0.119 |
| Sense | 0.245 |
| NII 1 | 0.119 |
| NII 2 | 0.086 |
| TMU 1 | 0.091 |
| TMU 2 | 0.136 |
| Ave. | 0.106 |

| EJ | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.280 |
| NAIST 2 | 0.250 |
| WEBLIO-EJ1 1 | 0.238 |
| Online A | 0.219 |
| Kyoto-U 1 | 0.216 |
| WEBLIO-EJ1 2 | 0.240 |
| SMT T2S | 0.240 |
| Kyoto-U 2 | 0.229 |
| SMT Hiero | 0.277 |
| SAS_MT | 0.248 |
| Sense | 0.395 |
| RBMT B | 0.217 |
| Ave. | 0.254 |

| JC | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.077 |
| SMT S2T | 0.069 |
| Sense | 0.087 |
| NICT | 0.066 |
| SMT Hiero | 0.202 |
| NAIST 2 | 0.093 |
| TOSHIBA 1 | 0.089 |
| Kyoto-U 1 | 0.091 |
| BJTUNLP | 0.198 |
| TOSHIBA 2 | 0.066 |
| Kyoto-U 2 | 0.163 |
| Online D | 0.035 |
| RBMT B | 0.083 |
| Ave. | 0.101 |

| CJ | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.168 |
| NAIST 2 | 0.203 |
| SAS_MT | 0.167 |
| SMT T2S | 0.236 |
| EIWA | 0.175 |
| Kyoto-U 1 | 0.199 |
| Kyoto-U 2 | 0.180 |
| SMT Hiero | 0.274 |
| Sense | 0.228 |
| Online A | 0.239 |
| RBMT A | 0.130 |
| Ave. | 0.200 |

\* Fleiss's Kappa values

# Inter Annotator Agreement

| JE | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.162 |
| NAIST 2 | 0.047 |
| SMT S2T | 0.099 |
| Kyoto-U 1 | 0.070 |
| TOSHIBA 1 | 0.098 |
| RBMT D | 0.075 |
| Online D | 0.055 |
| SMT Hiero | 0.119 |
| Sense | 0.245 |
| NII 1 | 0.119 |
| NII 2 | 0.086 |
| TMU 1 | 0.091 |
| TMU 2 | 0.136 |
| **Ave.** | **0.106** |

| EJ | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.280 |
| NAIST 2 | 0.250 |
| WEBLIO-EJ1 1 | 0.238 |
| Online A | 0.219 |
| Kyoto-U 1 | 0.216 |
| WEBLIO-EJ1 2 | 0.240 |
| SAS_MT | 0.248 |
| Sense | 0.395 |
| RBMT B | 0.217 |
| **Ave.** | **0.254** |

| JC | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.077 |
| SMT S2T | 0.069 |
| Sense | 0.087 |
| NICT | 0.066 |
| SMT Hiero | 0.202 |
| NAIST 2 | 0.093 |
| TOSHIBA 2 | 0.066 |
| Kyoto-U 2 | 0.163 |
| Online D | 0.035 |
| RBMT B | 0.083 |
| **Ave.** | **0.101** |

| CJ | |
|---|---|
| System ID | Kappa |
| NAIST 1 | 0.168 |
| NAIST 2 | 0.203 |
| SAS_MT | 0.167 |
| SMT T2S | 0.236 |
| EIWA | 0.175 |
| Kyoto-U 1 | 0.199 |
| Online A | 0.239 |
| RBMT A | 0.130 |
| **Ave.** | **0.200** |

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

* Fleiss's Kappa values

31

# Case Study

| Submission ID | BLEU | RIBES | HUMAN | Description |
|---|---|---|---|---|
| WEBLIO-EJ1 1 | 32.53 | 0.782 | 43.25 | w/o forest input |
| WEBLIO-EJ1 2 | 32.69 | 0.785 | 36.00 | w/ forest input |

VS.

BASELINE
(WEBLIO-EJ1 1)

WEBLIO-EJ1 2

| HUMAN | Kappa |
|---|---|
| 2.50 ± 4.17 | 0.528 |

EJ HUMAN score

43.25

36.00

WEBLIO-EJ 1 1    WEBLIO-EJ1 2

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

# Conclusion

- 12 participants for the evaluation task
  - including 3 companies and 3 teams outside Japan
- Human evaluation using crowdsourcing
- NAIST team achieved the best results for all the subtasks (congratulations!!)
- Shared the findings of MT for scientific papers
  - http://lotus.kuee.kyoto-u.ac.jp/WAT/papers/papers-2014.html

# Future Perspective

- Automatic evaluation server will keep running even after the workshop
  - promote continuous evolution of MT research
- WAT will be held annually
  - include more languages, domains…
- Let's share your resource!
  - monolingual/bilingual corpora, dictionaries…

# Future Perspective

- 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

- Discuss the importance of both sentence internal/<span style="color:red">external</span> information

# Thank you very much
# for attending WAT2014

# Stay tuned for
# the next WAT workshop!

# Future Perspective

- WAT will be held annually
  - include more languages, domains…
- Let's share your resource!
  - monolingual/bilingual corpora, dictionaries…
- Discuss the importance of both sentence internal/<span style="color:red">external</span> information