# AMTA

## MT Quality Evaluations: From Test Environment to Production

ELAINE OCURRAN
Welocalize
October 2015

welocalize
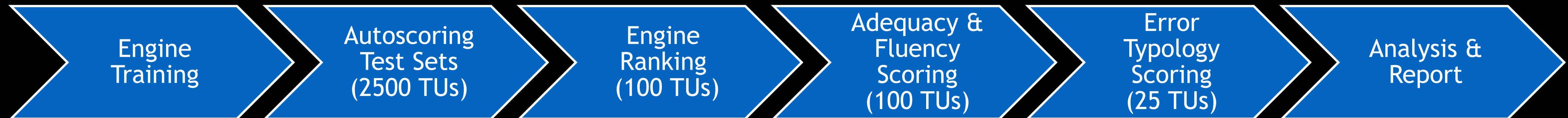doing things differently

# AGENDA

- Our MT evaluation methodologies

- Correlations between automatic scores and human evaluations

- Differences between system autoscores and PE autoscores

- MT evaluations in a production setting

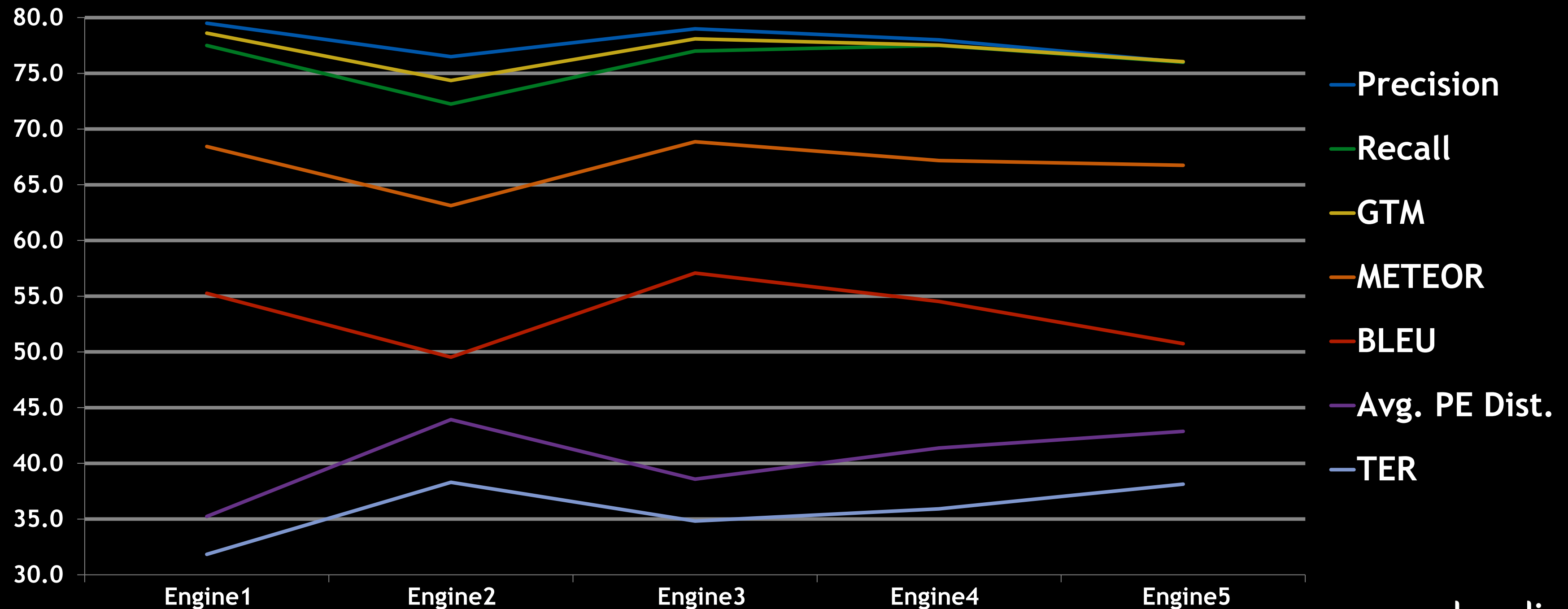- MT evaluations of post-edited files: a case study

welocalize
doing things differently

# OUR EVALUATION METHODS

## A TYPICAL EVALUATION PROCESS PER LOCALE AND PER ENGINE

Engine Training → Autoscoring Test Sets (2500 TUs) → Engine Ranking (100 TUs) → Adequacy & Fluency Scoring (100 TUs) → Error Typology Scoring (25 TUs) → Analysis & Report

welocalize
doing things differently

# OUR EVALUATION METHODS

## AUTOMATIC SCORES GENERATED BY WESCORE



Legend:
- Precision
- Recall
- GTM
- METEOR
- BLEU
- Avg. PE Dist.
- TER

X-axis: Engine1, Engine2, Engine3, Engine4, Engine5

Y-axis: 30.0 – 80.0

welocalize
doing things differently

# OUR EVALUATION METHODS

## HUMAN EVALUATIONS:
## ADEQUACY AND FLUENCY SCORING

**SCORE**

| 5 |
| 4 |
| 3 |
| 2 |
| 1 |

**ACCURACY**

All meaning expressed in the source fragment appears in the translation fragment.

Most of the source fragment meaning is expressed in the translation fragment.

Much of the source fragment meaning is expressed in the translation fragment.

Little of the source fragment meaning is expressed in the translation fragment.

None of the meaning expressed in the source fragment is expressed in the translation fragment.

**FLUENCY**

Native language fluency. No grammar errors, good word choice and syntactic structure. No PE required.

Near native fluency. Few terminology or grammar errors which don't impact the overall understanding of the meaning. Little PE required.

Not very fluent. About half of translation contains errors and requires PE.

Little fluency. Wrong word choice, poor grammar and syntactic structure. A lot of PE required.

No fluency. Absolutely ungrammatical and for the most part doesn't make any sense. Translation has to be re-written from scratch.

welocalize
doing things differently

# OUR EVALUATION METHODS

## HUMAN EVALUATION: ERROR TYPOLOGY

Errors per 25 Segments

Legend: DA, ES, FR, HI-IN, IT, JA, NL, NO, PT-BR, SV

Categories: Mistranslation, Ommission of Info, Addition of Info, Terminology, Text not translated, Untransltable text translated, Word Order, word form agreement, Tense / mood / aspect, Grammar, Wrong spelling, Capitalization, Punctuation, Spacing, Locale Adaptation, Compliance w/ client specs, Style / Register

welocalize
doing things differently

# LESSONS LEARNED

- We always perform autoscoring PLUS human scoring for all our MT evaluations. We have internal thresholds that qualify an engine ready for deployment and it's level of maturity.

- For bake-offs between several engines, we always include engine ranking in addition to our standard scores.

- Productivity tests are valuable during the initial phase of an MT program to build up productivity data for future reference across languages, domains and MT systems.

- Our MT program is now mature and we are able to perform most of our evaluations based on autoscoring PLUS human scoring, and by referencing the productivity data we have collected over a number of years.

welocalize
doing things differently

# NEXT

**Correlations between automatic scores and human evaluations**

welocalize
doing things differently

# CORRELATIONS

## CORRELATIONS BETWEEN AUTOMATIC SCORES AND HUMAN EVALUATIONS

| Pearson's r | Variables | Strength of Correlation | Tests (N) | Locales |
|---|---|---|---|---|
| 0.50576955 | Fluency & METEOR | Strong positive relationship | 150 | 11 |
| 0.50070425 | Fluency & BLEU | Strong positive relationship | 150 | 11 |
| 0.49816365 | Fluency & Recall | Strong positive relationship | 150 | 11 |
| 0.49724893 | Fluency & NIST | Strong positive relationship | 150 | 11 |
| 0.49195687 | Fluency & GTM | Strong positive relationship | 150 | 11 |
| 0.47064566 | Fluency & Precision | Strong negative relationship | 150 | 11 |
| 0.38293518 | Adequacy & NIST | Moderate negative relationship | 150 | 11 |
| 0.31354314 | Adequacy & METEOR | Moderate negative relationship | 150 | 11 |
| 0.2940756 | Adequacy & Recall | Weak positive relationship | 150 | 11 |
| 0.28586852 | Adequacy & GTM | Weak positive relationship | 150 | 11 |
| 0.28386332 | Adequacy & BLEU | Weak positive relationship | 150 | 11 |
| 0.26685854 | Adequacy & Precision | Weak positive relationship | 150 | 11 |
| -0.40270902 | Adequacy & TER | Strong negative relationship | 150 | 11 |
| -0.4788575 | Fluency & PE Distance | Strong negative relationship | 150 | 11 |
| -0.5385275 | Adequacy & PE Distance | Strong negative relationship | 150 | 11 |
| -0.5421933 | Fluency & TER | Strong negative relationship | 150 | 11 |

welocalize
doing things differently

# CORRELATIONS

## THE STRONGEST CORRELATION WAS FOUND BETWEEN FLUENCY AND TER



welocalize
doing things differently

# CORRELATIONS

## THE 2ND STRONGEST CORRELATION WAS FOUND BETWEEN ADEQUACY AND PE DISTANCE



welocalize
doing things differently

# LESSONS LEARNED

- It seems that we cannot rely solely on autoscores as long as the correlation with human judgment is not stronger than the data suggests

- TER and PE Distance show the strongest correlation to both Fluency and Adequacy, and therefor seem closer to human judgment than the other scores.

- Fluency correlates stronger with system autoscores than Adequacy overall.

- PE Distance is the only metric that correlates stronger with Adequacy than Fluency.  PE Distance is also the only character-based metric.

welocalize
doing things differently

# NEXT

**Differences between system
autoscores and post-editing
autoscores**

welocalize
doing things differently

# SYSTEM VS PE AUTOSCORES

## ON AVERAGE, THE POST-EDITING SCORE IS 15 AND 17 POINT HIGHER FOR PE DISTANCE AND BLEU RESPECTIVELY

| Pearson's r | Variables | Strength of Correlation | Tests (N) | Locales |
|---|---|---|---|---|
| 0.832226688 | BLEU (System) & BLEU (PE) | Very strong positive relationship | 57 | 9 |
| 0.832218909 | PE Distance (System) & PE Distance (PE) | Very strong positive relationship | 57 | 9 |



welocalize
doing things differently

# SYSTEM VS PE AUTOSCORES

## CORRELATIONS BETWEEN SYSTEM BLEU AND POST-EDITING BLEU



welocalize
doing things differently

# SYSTEM VS PE AUTOSCORES

## CORRELATIONS BETWEEN SYSTEM PE DISTANCE AND POST-EDITING PE DISTANCE



welocalize
doing things differently

# LESSONS LEARNED

- There is a very high correlation between the MT system autoscores generated during the evaluation phase and the autoscores generated from production using the same engines.

- However, the post-editing autoscores are considerably better than the MT system autoscores by around15%.

- We now differentiate the autoscores in our database as 'System' and 'PE'.

# NEXT

MT evaluations in a production setting

welocalize
doing things differently

# PRODUCTION SETTING

## HOW TO MEASURE POST-EDITING EFFORT

- It is important to monitor the performance of MT and post-editors, especially during the initial launch of a new program

- The use of autoscoring to analyze post-project files is a valuable and cost-effective method to measure the post-editing effort

- They support rate negotiations and can help us to identify over- or under-editing  by post-editors

- TER and PE Distance are useful metrics, with different underlying algorithms

welocalize
doing things differently

# PRODUCTION SETTING

## HOW TO MEASURE POST-EDITING EFFORT

**PE Distance -** lower is better!

- Measures the number of insertions, deletions, substitutions required to transform MT output to the required quality level

- PE Distance values are derived by comparing the post-edited segments with the corresponding machine translation segments

- In our analysis the PE distance applies the Levenshtein algorithm and is **character-based**. This captures morphological post-edits, such as fixing word forms.

welocalize
doing things differently

# PRODUCTION SETTING

## HOW TO MEASURE POST-EDITING EFFORT

**TER - l**ower is better!

- TER stands for  Translation Edit Rate

- It is an error metric for machine translation that measures the number of edits required to change a system output into the post-edited version

- Possible edits include the insertion, deletion, and substitution of single words as well as shifts of word sequences.

- Unlike PE Distance, TER is a **word-based** error metric  and therefor does not capture morphological changes during post-editing.

welocalize
doing things differently

# PRODUCTION SETTING

## TOOLS TO MEASURE POST-EDITING EFFORT

| TOOL | INPUT FILES | OUTPUT REPORT | PROS | CONS |
|------|-------------|---------------|------|------|
| *iOmegaT* | xliff & more | xml | Includes productivity data | Generated in the CAT tool during translation, requires post-editor buy-in |
| *MateCat* | xliff | Excel | Includes productivity data as a built in feature | Generated in the CAT tool during translation, requires post-editor buy-in |
| *Okapi* | xliff | html | Allows us to measure PE distance post-project | Requires access to pre- and post-edited file sets |
| *Post-Edit Compare* | sdlxliff | html | Allows us to measure PE distance post-project | Requires access to pre- and post-edited file sets |
| *Qualitivity* | sdlxliff | Excel | Includes productivity data | Generated in the CAT tool during translation, requires post-editor buy-in |
| *wescore* | tmx | Excel | Allows us to measure PE distance post-project | Proprietary tool, Requires access to pre- and post-edited file sets |

welocalize
doing things differently

# PRODUCTION SETTING

## MATECAT IS A FREE ONLINE CAT TOOL WITH EDITING LOG

welocalize
doing things differently

# PRODUCTION SETTING

## USE POST-EDIT COMPARE TO ANALYSE SDLXLIFF FILES



http://www.translationzone.com/openexchange/app/post-editcompare-495.html

welocalize
doing things differently

# PRODUCTION SETTING

## OKAPI FRAMEWORK TRANSLATION COMPARISON STEP

### Summary

Repartition for Trans1 to Trans2:

| Scores | ED-Scores | | | | FM-Scores | | | |
|--------|-----------|---|-------|---|-----------|---|-------|---|
| | Segments | % | Words | % | Segments | % | Words | % |
| 100 | 139 | 3 | 1414 | 3 | 176 | 4 | 1802 | 4 |
| 90 - 99 | 350 | 8 | 3954 | 8 | 346 | 8 | 3864 | 8 |
| 80 - 89 | 862 | 20 | 9850 | 20 | 674 | 16 | 7659 | 15 |
| 70 - 79 | 971 | 22 | 11137 | 23 | 804 | 19 | 9191 | 19 |
| 60 - 69 | 1078 | 25 | 12423 | 25 | 805 | 19 | 9332 | 19 |
| 50 - 69 | 598 | 14 | 6794 | 14 | 655 | 15 | 7500 | 15 |
| 40 - 59 | 197 | 5 | 2215 | 4 | 392 | 9 | 4479 | 9 |
| 30 - 39 | 33 | 1 | 359 | 1 | 240 | 6 | 2775 | 6 |
| 20 - 29 | 2 | 0 | 22 | 0 | 102 | 2 | 1159 | 2 |
| 10 - 19 | 1 | 0 | 4 | 0 | 36 | 1 | 398 | 1 |
| 0 - 9 | 104 | 2 | 1258 | 3 | 105 | 2 | 1271 | 3 |
| Total | 4335 | 100% | 49430 | 100% | 4335 | 100% | 49430 | 100% |

```
Total Number of Segments:       4335
Total Number of Words:          49430
Average word count per segment: 11.40
Average ED-Score (by segment):  Trans1 to Trans2 = 69.95
Average FM-Score (by segment):  Trans1 to Trans2 = 65.48
Average ED-Score (by word):     Trans1 to Trans2 = 69.76
Average FM-Score (by word):     Trans1 to Trans2 = 65.18
Edit Effort Score:              32.53
```

http://www.opentag.com/okapi/wiki/index.php?title=Translation_Comparison_Step

welocalize
doing things differently

# PRODUCTION SETTING

## QUALITIVITY PLUGIN FOR SDL TRADOS STUDIO



http://www.translationzone.com/openexchange/app/qualitivity-788.html

welocalize
doing things differently

# LESSONS LEARNED

- The use of autoscoring to analyze post-project files is a valuable and cost-effective method to measure the post-editing effort.

- A productivity test requires upfront organization and buy-in from translators.

- It is important to find a tool that works with the given file format and workflow.

- Access to pre- and post-edit versions of projects is required. This is a challenge on some accounts.

- Identification and separation of MT segments from fuzzy segments may be required for some tools.

- Look for consistency across languages and resources. Unusually high or low scores can be a sign of over-editing or under-editing.

welocalize
doing things differently

# NEXT

## MT evaluations of post-edited files: a case study

welocalize
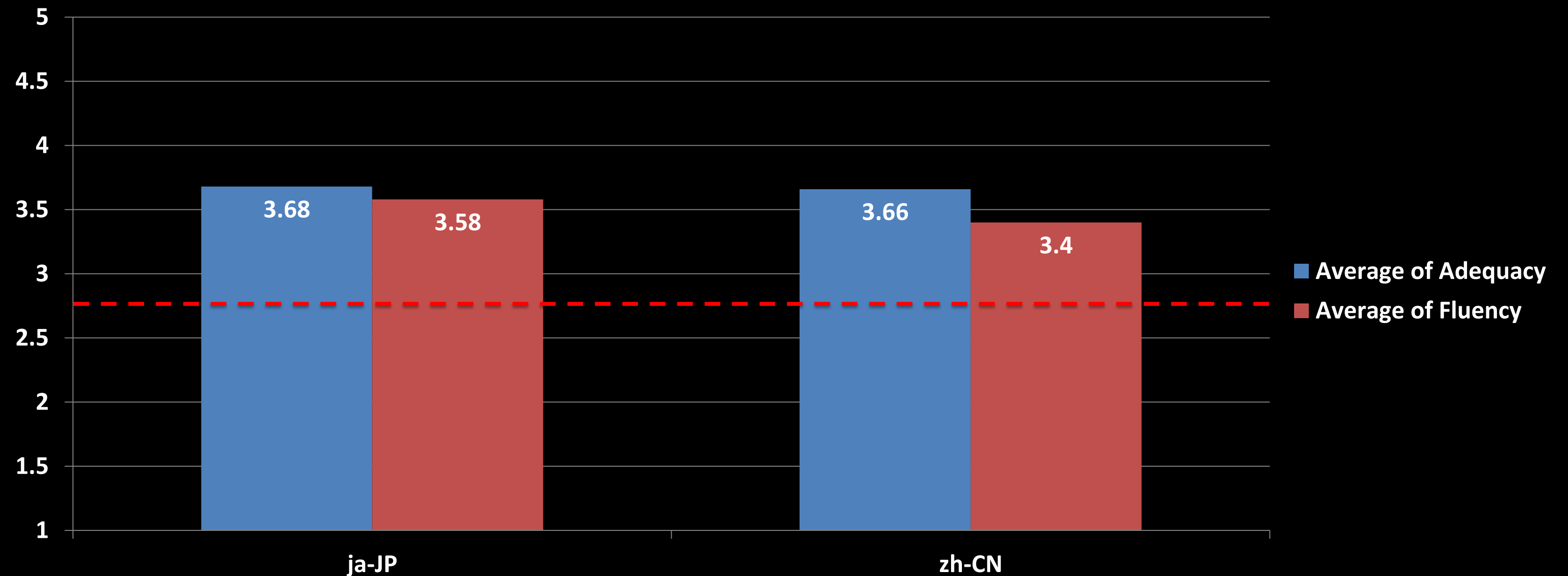doing things differently

# CASE STUDY

## TEST PILOT FOR LIGHT AND FULL POST-EDITING

- Languages: Chinese (Simplified) and Japanese

- The resources are regular translators for this client

- In order to have comparable data, the same resource performed both light and full post-editing tasks of 438 segments
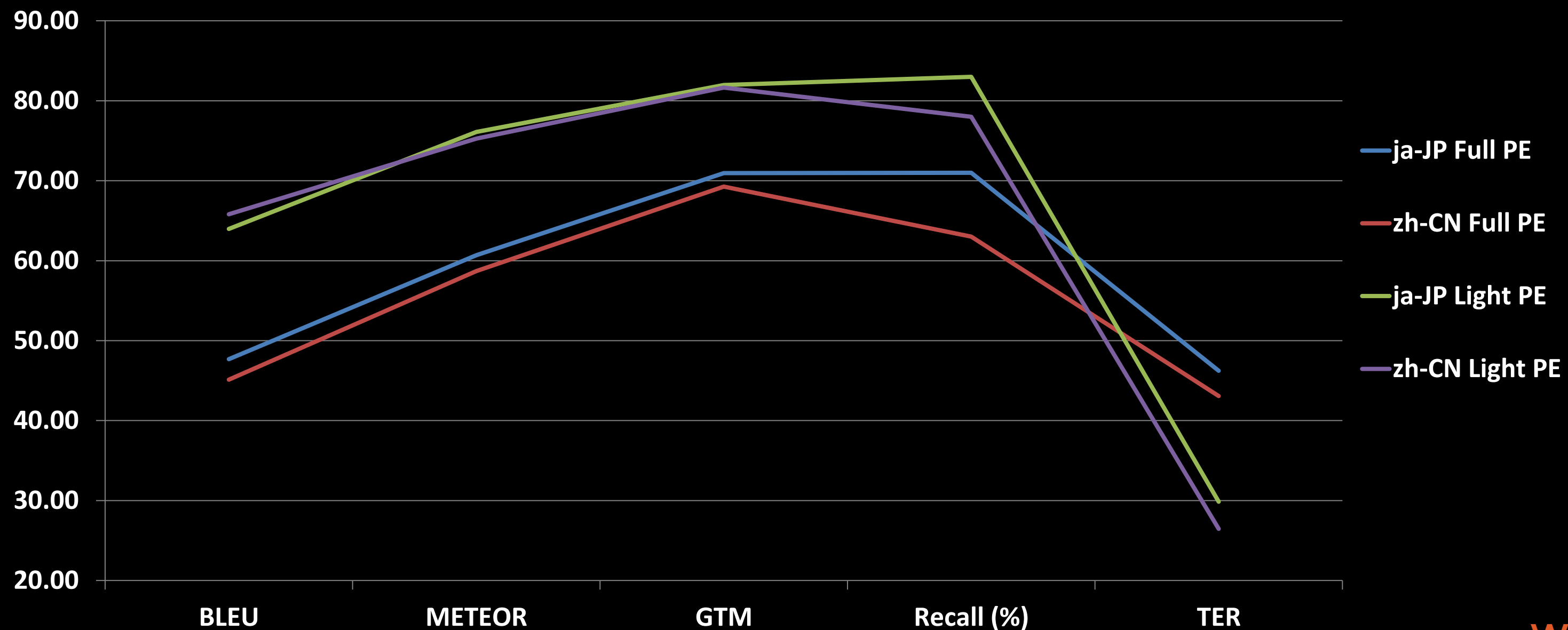
| Full Post-editing | Light Post-editing | LQA of PE Kits | Autoscoring PE Kits | Adequacy & Fluency Scoring | Error Typology Scoring | Analysis & Report |

welocalize
doing things differently

# CASE STUDY: LESSONS

## LESSONS LEARNED

- Using autoscores on post-edited translations can indicate the level of post-editing effort involved for a specific content and MT engine

- The autoscores also illustrate the difference in effort between Light and Full Post-editing, approximately 20 point delta for BLEU and 15 point delta for TER

- The autoscores confirm that the resources have indeed managed to perform two distinct post-editing levels

welocalize
doing things differently