

Utilizing Neural MT Engines in Industrial Translation

Toru Shishido

HUMAN SCIENCE

Human Science Co., Ltd.
WWW.SCIENCE.CO.JP



Agenda

- ◆ Evaluation
 - ◆ Raw Output Quality
 - ◆ Throughput PE vs HT
- ◆ Challenges / Best Practices
- ◆ Key Takeaways
- ◆ Q&A

HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Raw Output Quality Evaluation – Details of sample

- ◆ Human assessment
- ◆ Language pair: English-Japanese
- ◆ Translation volume: 3786 words
- ◆ Content type: Software manual

Raw Output Quality Evaluation – Human assessment / how to score

Better



Worse

Meaning and Accuracy

1. Perfectly Understandable

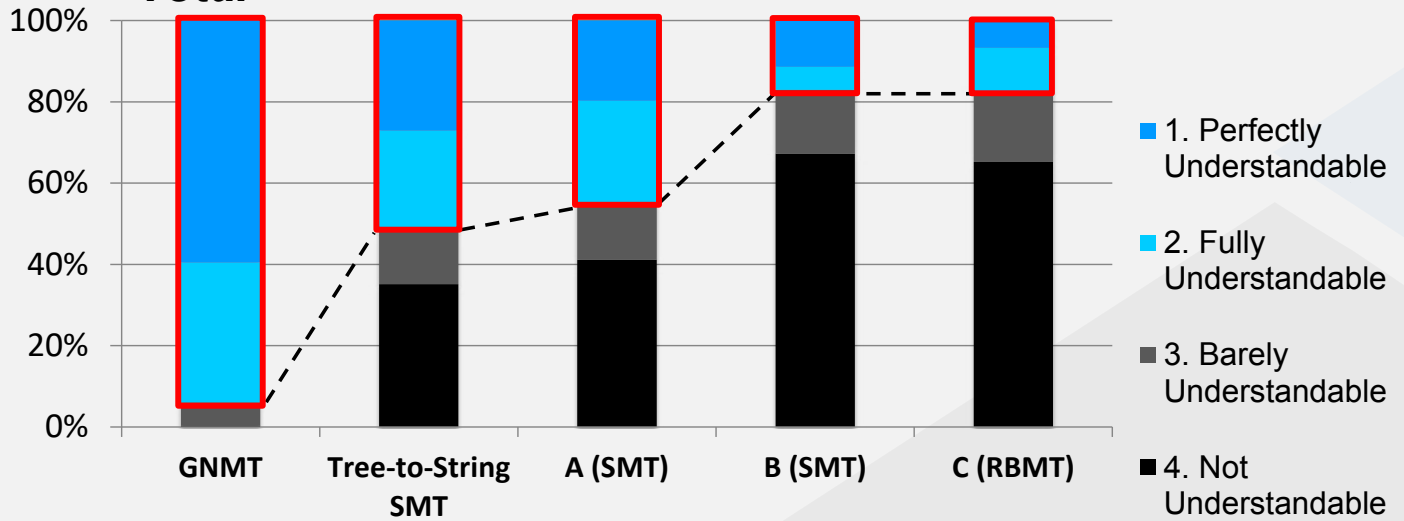
2. Fully Understandable

3. Barely Understandable

4. Not Understandable

Raw Output Quality Evaluation – Results

-Total-

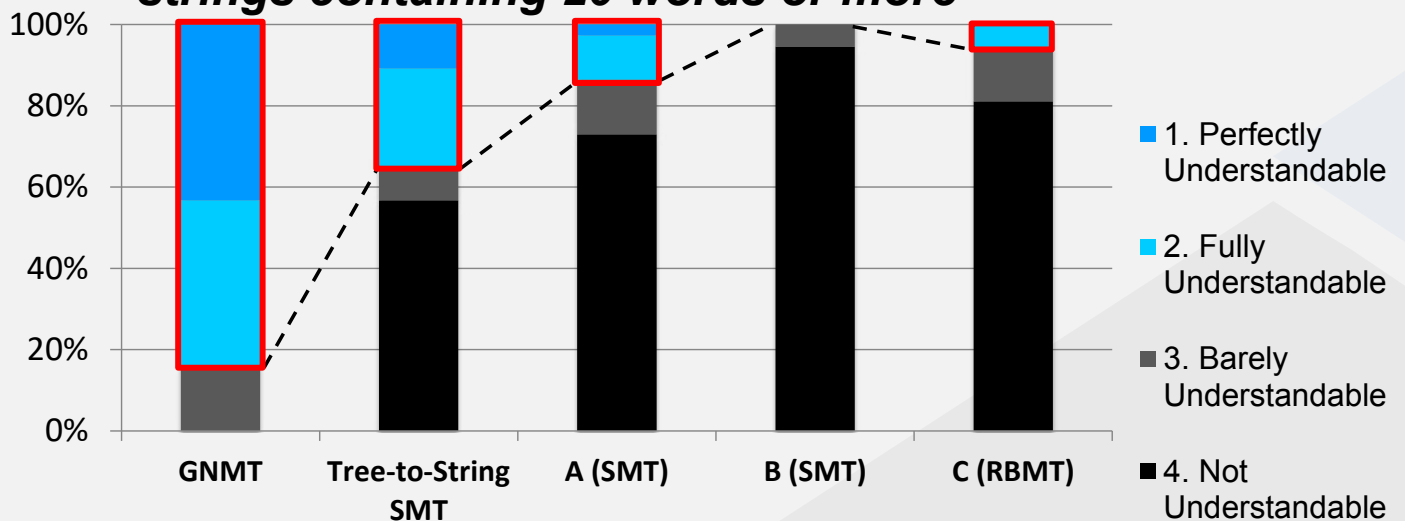


HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Raw Output Quality Evaluation – Results

-strings containing 20 words or more-



HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Raw Output Quality Evaluation – Analysis

- ◆ GNMT is extremely good in translating software manuals
- ◆ The quality of GNMT is high even with long sentences
- ◆ Reasons (speculations):
 - ◆ There are large amount of software manuals on the Internet
 - ◆ Google crawls the Internet for its training corpus
 - ◆ GNMT is like an MT system with a huge translation memory from multiple software vendors



Throughput Evaluation - Details

- ◆ Localization projects of various document types
- ◆ Not in production but completely simulated

(As of July 28, 2017)	Source volume	Number of projects
English-Japanese	49,883 weighted words	36
Japanese-English	10,057 weighted characters	2



Throughput Evaluation - Context levels

- ◆ Translation depends on the information outside the sentence
 - ◆ Other sentences in the document
 - ◆ Basic knowledge of the products / services
 - ◆ Common sense
- LOW: A sentence provides all the information for translation.
- MEDIUM
- HIGH: Information from other sources is needed.



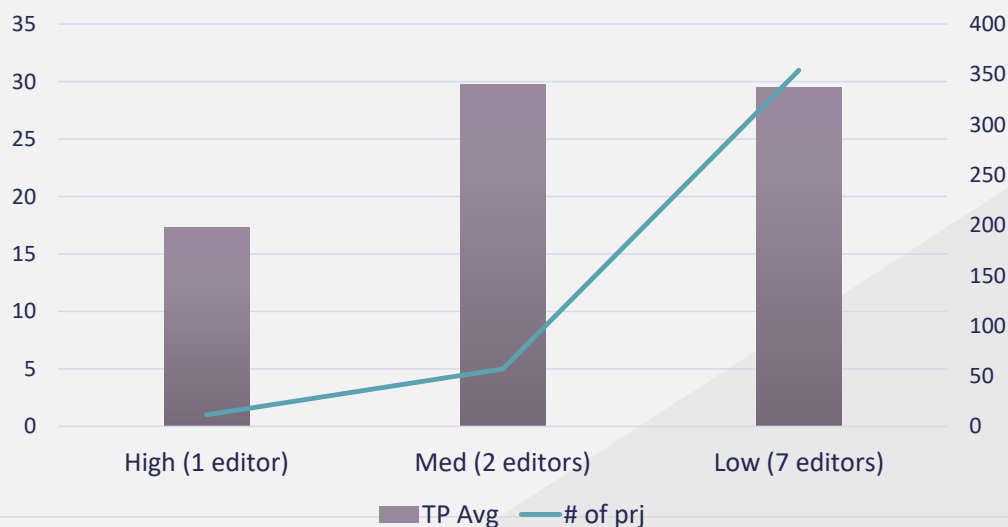
Throughput Evaluation - Results: English-Japanese

Weighted word count	PostEdit time (hr)	Speed (w/hr)	Content type	Context level
2174	11	197.6	Training role play scripts	HIGH
1229	4	307.3	Resource file	Faster than human translation (~250w/hr)
1108	3.5	316.8	FAQ (web services)	
3175	10	317.5	Product information	
1682	4	420.5	FAQ (web services)	LOW
1482	3	494.0	Service description	LOW
1023	2	615.0	Software manual	LOW



Throughput Evaluation - Results : English-Japanese

Throughput average by context level



Throughput Evaluation - Results: Japanese-English

Weighted char count	PostEdit time (hr)	Speed (ch/hr)	Content type	Context level
9352	4	2338.0	Whitepaper	LOW
705	0.33	2350.0	Developer page (UGC)	MEDIUM

Needs to collect more data, but
much faster than manual translation (~500ch/hr)
+ good for UGC



Challenges – Fun Fact

- ◆ There are 24 spelling patterns for the translation of **User Interface**:

ユーザーインターフェース	ユーザーインタフェース	ユーザーインターフェイス	ユーザーインタフェイス
ユーザインターフェース	ユーザインタフェース	ユーザインターフェイス	ユーザインタフェイス
ユーザー▲インターフェース	ユーザー▲インタフェース	ユーザー▲インターフェイス	ユーザー▲インタフェイス
ユーザ▲インターフェース	ユーザ▲インタフェース	ユーザ▲インターフェイス	ユーザ▲インタフェイス
ユーザー・インターフェース	ユーザー・インタフェース	ユーザー・インターフェイス	ユーザー・インタフェイス
ユーザ・インターフェース	ユーザ・インタフェース	ユーザ・インターフェイス	ユーザ・インタフェイス

▲ stands for a single-byte space.

HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Challenges (1) – Following rules in style guides (1)

- ◆ Most of companies have their own style guides and the rules are slightly different, such as spacing rules, brackets, long vowels (*cho-on*), etc.

Spacing rules	Company A	Company B	Company C
Katakana words	User interface ユーザー▲インターフェイス	User interface ユーザインタフェース	User interface ユーザー・インターフェイス
Between single-byte and double-byte characters	From Sept. 19 to 21 9▲月▲19▲日～▲21▲日	From Sept. 19 to 21 9月19日～9月21日	From Sept. 19 to 21 9▲月▲19▲日～21▲日

HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Challenges (1) – Following rules in style guides (2)

- ◆ Most of companies have their own style guides and the rules are slightly different, such as spacing rules, brackets, long vowels (*cho-on*), etc.

	Company A	Company B	Company C
Brackets	Use [] (single-byte) for user interface terms Use 『 』 for book titles and use 「 」 for chapter/section titles	Use [] (double-byte) for user interface terms Use 『 』 for book, chapter and section titles	Use 「 」 (double-byte) for user interface terms Use 『 』 for book, chapter and section titles
Long vowels (<i>cho-on</i>)	User ... ユーザー Printer ... プリンター Programmer ... プログラマー (depends of numbers of syllables)	User ... ユーザー Printer ... プリンター Programmer ... プログラマ	User ... ユーザ Printer ... プリンタ Programmer ... プログラマ

Challenges (2) – Tone (*de-ar* vs *desu-masu* (常体/敬体))

- ◆ There are two major writing styles in Japanese, *de-ar* style vs *desu-masu* style. These styles should be applied appropriately to match the context.

	Source	Raw MT	Post edited
<i>de-ar</i> style (常体)	(This course helps you to:) • Use new services and features from the ABC product to learn about modern technologies.	• ABC 製品の新しいサービスと機能を使用して、最新の技術を学ぶことができます。	• ABC 製品の新しいサービスと機能を使用して、最新の技術について学習する。 (e.g., bullet items)
<i>desu-masu</i> style (敬体)	Use new services and features from the ABC product to learn about modern technologies.	ABC 製品の新しいサービスと機能を使用して、最新の技術を学ぶことができます。	ABC 製品の新しいサービスと機能を使用すると、最新の技術を学ぶことができます。 (e.g., normal texts)

Challenges (3) – Glossary (UI terms / client specific / titles of references)

- ◆ Most companies have UI glossaries and terminologies so the post editors need to apply the appropriate terms.

	Source	Raw MT	Post edited
Example 1	Click on the “Continue” button.	「続行」ボタンをクリックします。	[Continue (続行)] ボタンをクリックします。
Example 2	a getting started guide	スタートガイド	入門ガイド
Example 3	詳細については、「APIを使用した展開」を参照してください。	For details, see “Deployment using API”.	For details, see “Deploying with API”.



Challenges (4) – General terms (Contexts/Inconsistencies)

- ◆ Post editors need to apply the correct translations to context-sensitive terms.
- ◆ Even the translations are correct, they must be consistent.

	Source	Raw MT	Consider when post editing
Example 1	available	利用可能 (<i>able to use</i>) ご利用いただけます (<i>polite “able to use”</i>) あります (<i>exists / be in stock</i>)	Post editors must consider the context of the text since the MT engines do not see the context.
Example 2	question	質問 (<i>an act of asking</i>) 問題 (<i>a problem</i>) 疑わしいこと (<i>a doubt</i>)	
Example 3	server-side	サーバーサイド (<i>server side</i>) サーバー側 (<i>server side</i>)	Both translations are correct, but inconsistent.



Challenges (4) – General terms (new words/buzzwords)

- ◆ Some new words may not be translated correctly sometimes.

	Source	Raw MT	Post edited
deep dive	The XXX Conference is a one-day deep dive into new technology.	XXX Conference は、新たな技術についての 深いダイビング です。 (a recreational diving)	XXX Conference は、新たな技術について考える 1 日間の ディープダイブ (or 分析ワークショップ) です。 (an extensive analysis)
DevOps	DevOps focuses on improving automation.	開発部門 は自動化の改善に重点を置いています。 (Development Dept.)	DevOps では自動化の改善に重点を置きます。

Challenges (5) – Tags / variables

- ◆ In most cases, tags are not properly treated. Also, tags can cause poor translation.

	Source	Raw MT	Post edited
 tag	Cover letter	Coverletter (the tag is omitted)	カバー レター
variable tag	Please ¥{0¥} to try again.	再試行するには ¥ ▲ {0 ▲ ¥} してください。 (unnecessary spaces)	¥{0¥}して、もう一度お試しください。

Challenges and solutions

Issues	Solutions	Can be fixed automatically?
Client-specific style specifications	Apply the rules with regular expression	Some yes, others no
Tone	Check and replace manually in Post Edit	No
Terminology (UI / client-specific / ref mat titles)	Apply some translations from terminology file automatically, and then replace manually in Post Edit (if necessary)	Some yes, others no
Terminology (general/new terms)		
Tags / variables	Delete before MT and insert manually in Post Edit	No

Best Practices

- ◆ Decide the content type to be machine-translated
 - Manuals, user interface, FAQ, UGC, marketing contents
- ◆ Align the final expectations between client and LSP
 - Final quality of translation, TATs, costs, content cycles
- ◆ Then, support and train post editors
 - Appropriate allocation of post editors by content type and final quality expectation, pre-process with SW components, continuous feedback loop

Takeaways - Neural MT for Commercial Use

- ◆ NMT makes the translation hours 1.36x faster and the productivity 1.48x higher (evaluation average)
- ◆ Usable in production both in English-Japanese and Japanese-English pairs in IT localization (incl. UGC)
- ◆ There are issues to be solved manually in Post Edit, but some can be automatically processed with software components

HUMAN SCIENCE
WWW.SCIENCE.CO.JP



Q&A



t-shishido@science.co.jp



www.science.co.jp



+81-3-5321-3111

HUMAN SCIENCE
WWW.SCIENCE.CO.JP

