

タイトル	WMT20 評価タスクデータにおける自動評価法のメタ評価
著者	越前谷, 博; Echizen ya, Hiroshi; 荒木, 健治; Araki, Kenji
引用	工学研究: 北海学園大学大学院工学研究科紀要(21): 19-32
発行日	2021-12-24

研究論文

WMT20 評価タスクデータにおける自動評価法のメタ評価

越前谷 博* · 荒木 健治**

Meta-evaluation of Automatic Evaluation Metrics by WMT20 Metrics Task Data

Hiroshi Echizen'ya* and Kenji Araki**

Abstract

本報告では機械翻訳分野の国際会議である WMT20 における Metrics Shared Task のデータを用いて行ったメタ評価について述べる。WMT20 の評価タスクでは様々な自動評価法が用いられており、自動評価法の現状を把握するために有効である。メタ評価はシステムレベルとセグメントレベルの2つの指標を用いて行った。また、我々が従来より提案している2つの自動評価法 IMPACT と WE_WPI も含めることで提案手法の他手法に対する位置付けについても検証した。その結果、WE_WPI は to-English (他言語から英語の翻訳) の言語ペアにおいては他手法に比べて比較的上位の評価精度を有していることを確認できた。

1 はじめに

機械翻訳のための自動評価法の研究は従来より盛んに行われており、機械翻訳の変遷と共に多様化している。特に、ここ10年においてはニューラル翻訳 [1, 2, 3, 4] の急速な普及により、文や単語の意味表現に基づく自動評価法の研究が盛んに行われている。自動評価法の研究は1990年代に登場した統計的機械翻訳 [5, 6, 7] により、本格的に研究が始まった。統計的機械翻訳はソースコードが一般公開されており、膨大な対訳コーパスを共有することで容易に稼働できるため、多くの研究者が開発に取り組んだ。その結果、開発サイクルを早めるためにはシステムに対する評価速度が重要となり、自動でシステム訳を評価可能なシステムが不可欠となった。そのような背景のもと人手のように評価の揺れが生ぜず、低コストかつ高速で評価が可能な自動評価法の研究が盛んに行われるようになった。

統計的機械翻訳が主流だった1990年代から2000年代にかけては、表層情報及び WordNet の

ような静的な語彙情報を用いた自動評価法が主流であった。自動評価法のデファクトスタンダードとなっている BLEU [8] は2002年に提案された。BLEU は参照訳と呼ばれる正解訳とシステム訳との間で n-gram 一致率を求めることでシステム訳に対する評価スコアを算出するシステムである。このように自動評価法では参照訳を用いた評価が一般的である。この BLEU がきっかけとなり多くの自動評価法が提案されるようになった。代表的なものとしては TER [9], METEOR [10], RIBES [11] などが挙げられる。我々も表層情報に基づく自動評価法として IMPACT (Intuitive comMon PArts ConTinuum) [12, 13, 14, 15, 16, 17] を提案している。IMPACT は最長共通部分列 [18] を用いて、参照訳とシステム訳間に存在するチャンクを自動的に決定し、そのチャンクの位置と長さに基づいてチャンク列を再帰的に決定することで評価スコアを求める。

2010年代に入ると統計的機械翻訳に取って代わり、ディープラーニングの技術発展に伴いニューラル翻訳の研究が急速に行われるように

* 北海学園大学大学院工学研究科電子情報生命工学専攻

Graduate School of Engineering (Electronics, Information and Life Science Eng.), Hokkai-Gakuen University

** 北海道大学大学院情報科学研究院

Faculty of Information Science and Technology, Hokkaido University

なった。ニューラル翻訳ではそれまでの表層情報や静的な言語資源に依存することなく、単語や文の意味をベクトルで表現することでより高度な言語処理が可能となった。このようなニューラル翻訳の進展に伴い、自動評価法もそれまでの表層情報や静的な言語知識を用いたものだけではなく、ディープラーニングより学習したモデルにより得られる単語分散表現や文ベクトルに基づく手法が提案されるようになった。我々もこのような状況下において単語分散表現に基づく自動評価法として WE_WPI (Word Embedding-based Automatic MT Evaluation using Word Position Information) [19] を提案している。WE_WPI は Earth Mover's Distance (EMD) [20, 21, 22] を言語処理に適応させることでシステム訳に対する評価スコアを得る。EMD を言語処理に適用するためには 3 つのパラメータを定義する必要がある。WE_WPI では特徴量、その特徴量に対する重み、そして、特徴量間における距離計算の 3 つのパラメータにそれぞれ単語分散表現、 $tf \cdot idf$ 、そして、コサイン距離を適用している。

本報告では、提案手法である IMPACT と WE_WPI を含めた様々な自動評価法を用いて行なったメタ評価について述べる。メタ評価には WMT20 (The 2020 Conference on Machine Translation) [23, 24] の Shared Task の一つである Metrics のデータを使用する。国際会議 WMT は毎年開催され、機械翻訳に関する複数の Shared Task を扱っている。評価タスクは Shared Task としては開催初期から取り挙げられており、機械翻訳において非常に重要なタスクとして認識されている。この WMT20 の評価タスクデータによるメタ評価の結果、提案手法 WE_WPI は他言語から英語の翻訳を対象とする to-English においてシステムレベル、セグメントレベルの両方で他手法との比較において評価精度が上位に位置することを確認できた。

2 機械翻訳における自動評価法

WMT20 で使用された自動評価法と我々の 2 つの提案手法 IMPACT, WE_WPI のそれぞれの特徴をまとめたものを表 1 に示す。表 1 の“learned”は推定モデルやベクトルを得るためのモデルを学習により構築する必要があるかどうかを意味し、“yes”はその必要があることを示して

いる。モデルの学習を前提とした自動評価法は学習によりタスクに順応可能なため一般的には評価精度が高くなるが、学習のためのデータや学習時間が別途必要になることがデメリットとなる。また、表 1 の“types”における“src-based”はスコア計算において参照訳を必要としない自動評価法であることを示している。近年では参照訳を準備することのコストを回避するために原文のみを用いる自動評価法の研究も行われている。

2.1 WMT20 における自動評価法

2.1.1 ベースライン

表 1 の sentBLEU, BLEU, TER, chrF, そして、chrF++ は WMT20 ではベースラインとして位置付けられている。sentBLEU は BLEU をセグメントレベルでも評価スコアを算出できるように改良されたものである。TER は挿入、削除、移動、そして、置換の 4 つの編集操作の数に基づく自動評価法である。chrF は参照訳とシステム訳の比較において単語の n-grams ではなく文字の n-grams を用いた自動評価法である。また、chrF++ は単語の unigram と bigram に文字の n-grams を加えた自動評価法である。

2.1.2 BERT-base-L2, BERT-large-L2, mBERT-L2

BERT-base-L2, BERT-large-L2, そして、mBERT-L2 の 3 つの自動評価法はいずれも BERT [25] を fine-tuning することで得られる。その際には、WMT15 から WMT18 の評価タスクデータを用いている。BERT-base-L2 は 12 レイヤーの Transformer アーキテクチャに対して、英語データを用いて pre-trained を行ったものである。mBERT-L2 は 102 言語の Wikipedia データを用いて学習を行なっている。BERT-large-L2 は 24 レイヤーを用いており、WMT20 では to-English のみのメタ評価を行っている。

2.1.3 BLEURT, BLEURT-extended, Yisi-combi, bleurt-Yisi-combi

BLEURT は 2 つのデータを用いて BERT-based の回帰モデルのための学習を行なっている。一つはランダムに取得した 100 万ペアのデータであり、もう一つは WMT15 から WMT19 より取得したデータである。BLEURT-extended

表 1 : WMT20 における自動評価法と提案手法のそれぞれの特徴

metric	features	learned	types
sentBLEU	n-grams		Reference-based
BLEU [8]	n-grams		Reference-based
TER [9]	edit-distance		Reference-based
chrF [26]	character n-grams		Reference-based
chrF++ [27]	character n-grams		Reference-based
parbleu [28]	paraphrases		Reference-based
parachrf++ [28]	paraphrases		Reference-based
paresim [28]	paraphrases	yes	Reference-based
prism [29]	paraphrases		Reference-based
CharacTER [30]	character edit distance		Reference-based
EED [31]	character edit distance		Reference-based
SWSS+METEOR [32]	semantic similarity		Reference-based
MEE [33]	word embeddings		Reference-based
YiSi [35] [36]	contextual word embeddings		Reference-based
BERT-base-L2 [25]	contextual word embeddings	yes	Reference-based
BERT-large-L2 [25]	contextual word embeddings	yes	Reference-based
mBERT-L2 [25]	contextual word embeddings	yes	Reference-based
BLEURT [25]	contextual word embeddings	yes	Reference-based
BLEURT-extended [37]	contextual word embeddings	yes	Reference-based
Yisi-combi [25]	contextual word embeddings	yes	Reference-based
bleurt-combi [25]	contextual word embeddings	yes	Reference-based
COMET [38]	predictor-estimator model	yes	Reference-based
COMET-Rank [38]	predictor-estimator model	yes	Reference-based
COMET-HTER [38]	predictor-estimator model	yes	Reference-based
COMET-2R [38]	predictor-estimator model	yes	Reference-based
COMET-MQM [38]	predictor-estimator model	yes	Reference-based
BAQ, EQ	?	?	Reference-based
COMET-QE [38]	predictor-estimator model	yes	src-based
OpenKiwi-Bert [39]	predictor-estimator model	yes	src-based
OpenKiwi-XLMR [39]	predictor-estimator model	yes	src-based
YiSi-2 [40]	contextual word embeddings		src-based
IMPACT	longest common subsequence		Reference-based
WE_WPI	word embeddings		Reference-based

は WMT15 から WMT19 の人手評価に基づき BERT-based の回帰モデルを学習したものである。Yisi-combi は WMT データで fine-tuned された mBERT モデルに対して YiSi-1 を用いたものである。また、bleurt-combi は複数参照訳を用いて算出されたスコアに対して、YiSi-1, YiSi-2, そして、BLEURT のスコアを組み合わせたものである。

2.1.4 CharacTER

CharacTER は TER を文字レベルで適用した自動評価法である。編集操作は参照訳と一致するまで行われ、システム訳の長さで正規化される。ただし、移動の操作においては単語レベルで行われる。したがって、参照訳と移動の操作後のシステム訳との間のレーベンシュタイン距離は文字レベルで求められることになる。さらに CharacTER

は英語—ロシア語の言語ペアを除き、トークン化されていない参照訳とシステム訳に対して適用される。

2.1.5 COMET

COMET は推定モデルおよび翻訳ランキングモデルを用いて構築される。ニューラルモデルは参照訳とシステム訳をエンコードするために XLM-RoBERTa を用いる。さらに WMT17 から WMT19 における Direct Assessments (DA) に基づいて学習される推定モデルとなっている。COMET-2R は複数参照訳を処理するために学習された COMET の改良版である。COMET-Rank はシステム訳と原文との間の距離及びシステム訳と参照訳との間の距離を最適化するための翻訳ランキングアーキテクチャを用いている。そして、COMET-HTER と COMET-MQM は同じアーキ

テクチャを用いているが、COMET-HTER は Human-mediated Translation Edit Rate (HTER), COMET-MQM は多次元品質推定に基づいている。

2.1.6 EED

EED は “jump” の操作を含めた拡張版の編集距離を character-based で求めることで評価スコアを算出する。“jump” の操作はスペース文字が見つかったタイミングで行われる。

2.1.7 MEE

MEE のスコアはシステム訳と参照訳の単語分散表現の間の類似度を求めることで得られる。単語分散表現による類似度計算に用いられる単語分散表現は Facebook より提供されている fasttext [34] を用いることで取得する。さらに MEE ではスコア計算には表層的な一致、原形の一致 (root match), そして、同義語の一致の 3 つのモジュールを用いている。MEE は既存研究の中で学習済みの fasttext に基づく単語分散表現を用いていることから提案手法 WE_WPI に最も近い自動評価法と言える。しかし、提案手法では文全体の類似度を EMD に基づいて求めている点や EMD に語順情報を導入している点が大きく異なっている。

2.1.8 esim

esim はニューラルモデルに基づく自動評価法である。システム訳と参照訳をベクトルで表現するためにセンテンスアテンションとセンテンスマッチングについてのヒューリスティクスを用いている。

2.1.9 OpenKiwi-Bert, OpenKiwi-XLMR

OpenKiwi-Bert と OpenKiwi-XLMR は品質推定モデルを WMT 評価タスクデータで学習することで実現している。OpenKiwi は WMT20 QE shared task において最先端のモデルとして位置付けられている。

2.1.10 parbleu, parchr++ , paresim

parbleu, parchr++ , そして, paresim は与えられた人手評価と自動生成された言い換えのセットに対して BLEU, charf++ , esim の自動評価法をそれぞれ適用することで評価スコアを算出している。

2.1.11 prism

prism は 39 言語の言語ペアからなるデータを学習して構築された多言語のニューラル翻訳に基づいている。セグメントレベルの評価スコアはシステム訳と参照訳の間のデコーディングにより得られる。また、システムレベルの評価スコアはセグメントレベルのスコアの平均値を用いて得られる。

2.1.12 SWSS+METEOR

SWSS+METEOR は文中の意味的に重要な単語を抽出することで評価スコアを得る。ここでは意味的に重要な単語を定義するために意味表現フレームワークの UCCA を用いる。そのうえで重み付けされた SWSS と METEOR の組み合わせにより評価スコアを算出する。

2.1.13 YiSi-0, YiSi-1, YiSi-2

YiSi は統一的な意味レベルの機械翻訳の品質評価及び利用可能な様々なレベルの言語資源を有する言語に対する推定尺度に基づいた自動評価法である。YiSi-1 は事前学習されたモデルから単語分散表現を抽出し、語彙の意味レベルでの類似度を求める。そのうえでシステム訳と参照訳の間の意味的な類似度を求める。YiSi-2 は参照訳を必要としない自動評価法となっている。事前学習済みの言語モデルから抽出される単語分散表現を多言語へ写像することで意味的な類似度を求め、評価スコアを得ている。YiSi-0 は迅速に評価スコアを算出できるように YiSi-1 を簡素化したものである。ここでは表層情報を用いており、語彙的な類似度を最長共通部分列を文字レベルで適用することで評価スコアを求める。

2.2 提案手法

本節では我々が提案している自動評価法である IMPACT と WE_WPI について述べる。

2.2.1 IMPACT

IMPACT は単語を単位としたチャンクを再帰的に求めることで語順に着目した自動評価法となっている。システム訳と参照訳の間の最長共通部分列を一意に決定するために共通部分の長さとの相対的な位置関係を用いる。一意に決定された共通部分をチャンクと見なし、チャンクの長さに依

存した重みを用いることでチャンクが長く、かつチャンクの出現位置がシステム訳と参照訳間で類似しているほど評価スコアが高くなるようにパラメータにより制御している。また、このチャンクの決定処理を再帰的に行うことで参照訳に対して語順が異なるシステム訳の評価スコアが小さくなるようにパラメータを用いて制御している。

この IMPACT は表層情報のみに基づいているため容易に評価スコアを算出することが可能である。その一方で、表層情報のみに依存するため同義語や活用形による表層上の違いを考慮することができないという問題点もある。

2.2.2 WE_WPI

WE_WPI は単語の意味表現である単語分散表現を Earth Mover's Distance (EMD) に適用した自動評価法である。さらにシステム訳と参照訳間の語順の違いを EMD の距離行列に反映させることで語順についても考慮した評価が可能となっている。具体的には単語分散表現で表された単語間の距離をコサイン距離より求めるが、その距離に単語アライメントにより得られた単語間の位置の相対的なずれを負の重みとして付与することで対応関係にある単語間であっても出現位置が大きく異なる場合には、距離が大きくなる。

また、EMD は輸送問題における最適解を求めるための手法であり、荷物を輸送するための最適な輸送経路を探索する。そのため EMD では荷物を重みとして定義する必要がある。WE_WPI ではその重みについては文レベルの $tf \cdot idf$ を計算することで得ている。また、個々の単語の分散表現は事前学習された fasttext を用いて取得する。そのため、単語分散表現を得るための学習済みのモデルさえ取得できれば容易に評価スコアを得ることが可能である。

3 性能評価実験

3.1 実験データ及び実験方法

本報告では自動評価法のメタ評価を行うために WMT20 の評価タスクデータを用いた。WMT は機械翻訳のみを対象とした、最先端の研究発表が行われる国際会議として広く認知されている。さらに WMT では機械翻訳に関する様々なタスクごとにコンペティションを実施することで機械翻

訳分野の進展を図っている。評価タスクは 2008 年より WMT で取り挙げられており、非常に重要なタスクとして位置付けられている。

WMT20 の評価タスクでは to-English (他言語から英語への翻訳) については 10 言語ペア、out-of-English (英語から他言語への翻訳) については 8 言語ペアのニュースに関するシステム訳と参照訳のデータが提供されている。to-English においてはチェコ語から英語への翻訳 (cs-en)、ドイツ語から英語への翻訳 (de-en)、日本語から英語への翻訳 (ja-en)、ポーランド語から英語への翻訳 (pl-en)、ロシア語から英語への翻訳 (ru-en)、タミル語から英語への翻訳 (ta-en)、中国語から英語への翻訳 (zh-en)、イヌクティトゥット語から英語への翻訳 (iu-en)、クメール語から英語への翻訳 (km-en)、そして、パシュトー語から英語への翻訳 (ps-en) により得られたシステム訳を用いている。out-of-English については英語からチェコ語への翻訳 (en-cs)、英語からドイツ語への翻訳 (en-de)、英語から日本語への翻訳 (en-ja)、英語からポーランド語への翻訳 (en-pl)、英語からロシア語への翻訳 (en-ru)、英語からタミル語への翻訳 (en-ta)、英語から中国語への翻訳 (en-zh)、そして、英語からイヌクティトゥット語への翻訳 (en-iu) により得られたシステム訳を用いている。WMT20 ではこれらの言語ペアのシステム訳に対する参照訳も提供されているため自動評価法を用いて評価スコアを算出することができる。

また、メタ評価ではシステムレベルとセグメントレベルの 2 つの観点からのメタ評価を実施した。具体的には自動評価法が算出した評価スコアと人手評価との間で相関係数を求めた。システムレベルのメタ評価ではピアソンの相関係数を用いた。セグメントレベルのメタ評価ではケンドール τ を用いた。システムレベルのメタ評価においてもケンドール τ を求めているが、ケンドール τ はその特性から少ないデータにおいては複数の自動評価法で同等な値が得られる傾向があり、差別化を図る際には不十分となる。したがって、WMT20 においてはシステムレベルのケンドール τ は参考程度として求めているため、本報告でも同様の扱いとする。

3.2 実験結果

表2と表3にシステムレベルのメタ評価の結果を示す。表2はto-Englishにおけるシステムレベルのピアソンの相関係数、表3はout-of-Englishにおけるシステムレベルのピアソンの相関係数である。また、表4と表5にセグメントレベルのメタ評価の結果を示す。表4はto-Englishにおけるセグメントレベルのケンドール τ 、表5はout-of-Englishにおけるセグメントレベルのケンドール τ である。また、表6と表7にはシステムレベルのケンドール τ を示す。表6はto-English、表7はout-of-Englishとなっている。

表2から表5は各言語ペアにおいて“all”と“-out”の2種類の相関係数が示されている。“all”は全てのシステム訳を用いた相関係数である。それに対して“-out”は“outlier（外れ値）”を意味し、任意の機械翻訳システムのスコアが他の機械翻訳システムのスコアから大きく離れている場合には過度に高い相関係数を導く可能性があるため“all”から取り除いており、その場合の相関係数となっている。また、メタ評価では統計的有意性に基づく相関係数の有意差を示すためにWilliamsの有意検定[41]を使用している。表中の太字の相関係数は同じ言語ペアにおいて他の自動評価法に対して統計的有意差がなかったことを示している。

3.3 考察

我々が提案しているIMPACTとWE_WPIのメタ評価における位置付けについて述べる。IMPACTは表層情報にのみ基づく自動評価法である。そのため、同様のアプローチの自動評価法は表1のfeaturesがn-grams, edit distance, character n-grams, そして, character edit distanceとなっているBLEU, TER, chrF, chrF++, CharacTER, そして, EEDである。また、WE_WPIは事前学習された単語分散表現モデルのfasttextを用いている。したがって、同様のアプローチとしては表1のEEDが該当する。そこで、IMPACTとWE_WPIと類似している自動評価法と比較するために上述した自動評価法について全言語ペアの相関係数の平均を求めた。表8に相関係数の平均を示す。

表8よりシステムレベルのメタ評価においては

表層情報を用いた自動評価法の中ではEEDが最も高い評価精度を示している。IMPACTはそのEEDに比べると低い値であるが上位に位置している。特にout-of-Englishの“-out”はIMPACTが最も高い値を示した。また、単語分散表現を用いる自動評価法においてはWE_WPIはMEEよりも高い評価精度を示した。なお、out-of-Englishの平均値が未記入となっているが、これは相関係数が得られなかった言語ペアがEED, WE_WPIそれぞれに存在していたためである。WE_WPIにおいてはイヌクティット語のfasttextのモデルが取得できなかったためen-iuの評価値が得られず、それに伴い相関係数が得られなかった。

表8のセグメントレベルのメタ評価においては表層情報を用いた自動評価法ではEEDがto-Englishで最も値が高く、out-of-EnglishではchrFが最も高い値を示した。IMPACTは下位に位置し、十分な評価精度を示すことができなかった。この結果からIMPACTは特にセグメントレベルの評価精度が不十分であることが明らかとなった。一方、単語分散表現を用いた自動評価法においては我々が提案するWE_WPIがセグメントレベルにおいてもMEEよりも高い評価精度を示した。WE_WPIは表層情報に基づく自動評価法を含めた場合の比較においてもto-Englishの中で最も高い値を有しており、セグメントレベルの評価においては有効であることが確認できた。

4 まとめ

本報告では、WMT20の評価タスクデータを用いたメタ評価の結果について述べた。メタ評価の結果、提案手法のWE_WPIはto-Englishの言語ペアにおいてWMT20における様々な先行研究に対して比較的高い評価精度が得られることを確認できた。しかし、out-of-Englishにおいては先行研究と比べて上位に位置しているとは言えず不十分であった。今後はその原因について精査し、WE_WPIの評価精度の向上のための改良を進める予定である。

References

- [1] Ilya Sutskever, Oriol Vinyals and Quoc V. Le. 2014. Sequence to Sequence Learning with Neural Networks.

表 2 : to-English におけるシステムレベルのピアソンの相関係数

	cs-en		de-en		ja-en		pl-en		ru-en		ta-en		zh-en		iu-en		km-en		ps-en	
	All	-out	All	-out	All	-out	All	-out	All	-out	All	-out	All	-out	All	-out	All	-out	All	-out
	12	10	12	9	10	7	14	13	11	10	14	12	16	15	11	9	7	7	6	6
BAQ_dyn	-	-	-	-	-	-	-	-	-	-	-	-	-	0.956	0.928	-	-	-	-	-
BAQ_static	-	-	-	-	-	-	-	-	-	-	-	-	-	0.960	0.933	-	-	-	-	-
BERT-base-L2	0.775	0.693	0.997	0.791	0.971	0.789	0.552	0.328	0.919	0.836	0.909	0.746	0.967	0.929	0.704	0.145	0.967	0.967	0.945	0.945
BERT-large-L2	0.784	0.695	0.990	0.800	0.974	0.784	0.520	0.282	0.925	0.843	0.901	0.760	0.962	0.928	0.744	0.211	0.959	0.959	0.950	0.950
BLEU	0.851	0.800	0.985	0.778	0.969	0.826	0.549	0.355	0.884	0.761	0.916	0.807	0.956	0.957	0.569	0.348	0.969	0.969	0.888	0.888
BLEURT	0.792	0.725	0.996	0.770	0.978	0.820	0.591	0.371	0.924	0.844	0.906	0.768	0.966	0.931	0.771	0.320	0.984	0.984	0.955	0.955
BLEURT-extended	0.771	0.668	0.985	0.818	0.961	0.772	0.551	0.298	0.900	0.797	0.897	0.743	0.945	0.931	0.789	0.359	0.985	0.985	0.942	0.942
CharacTER	0.844	0.812	0.998	0.687	0.970	0.895	0.522	0.325	0.927	0.869	0.965	0.880	0.964	0.950	0.763	0.410	0.977	0.977	0.841	0.841
chrF++	0.872	0.806	0.997	0.687	0.968	0.861	0.528	0.312	0.890	0.831	0.951	0.828	0.976	0.954	0.729	0.337	0.978	0.978	0.898	0.898
chrF++	0.867	0.804	0.997	0.699	0.974	0.871	0.538	0.328	0.894	0.833	0.953	0.830	0.975	0.955	0.726	0.392	0.983	0.983	0.900	0.900
COMET	0.783	0.694	0.998	0.773	0.964	0.828	0.591	0.345	0.923	0.836	0.880	0.764	0.952	0.931	0.852	0.605	0.971	0.971	0.941	0.941
COMET-2R	0.777	0.697	0.998	0.772	0.964	0.818	0.584	0.332	0.924	0.843	0.881	0.770	0.949	0.928	0.872	0.644	0.970	0.970	0.949	0.949
COMET-HTER	0.738	0.661	0.995	0.767	0.912	0.702	0.446	0.231	0.867	0.741	0.726	0.595	0.809	0.873	0.770	0.464	0.901	0.901	0.862	0.862
COMET-MQM	0.728	0.612	0.991	0.684	0.906	0.707	0.424	0.222	0.858	0.746	0.767	0.617	0.784	0.862	0.841	0.631	0.914	0.914	0.880	0.880
COMET-QE	0.755	0.622	0.939	0.805	0.892	0.585	0.447	0.218	0.883	0.773	0.795	0.672	0.847	0.887	0.685	0.661	0.896	0.896	0.832	0.832
COMET-Rank	0.705	0.534	0.964	0.757	0.923	0.793	0.483	0.284	0.868	0.732	0.787	0.664	0.877	0.909	0.158	0.214	0.911	0.911	0.855	0.855
EED	0.884	0.838	0.997	0.752	0.974	0.904	0.538	0.299	0.926	0.872	0.958	0.862	0.956	0.932	0.821	0.587	0.990	0.990	0.930	0.930
esim	0.790	0.716	0.998	0.808	0.983	0.822	0.591	0.358	0.928	0.834	0.885	0.801	0.963	0.910	0.807	0.514	0.929	0.929	0.929	0.929
IMPACT	0.848	0.801	0.996	0.799	0.973	0.879	0.536	0.289	0.934	0.870	0.911	0.840	0.952	0.916	0.714	0.587	0.981	0.981	0.910	0.910
mBERT-L2	0.798	0.715	0.995	0.824	0.969	0.811	0.555	0.302	0.908	0.805	0.887	0.740	0.959	0.935	0.837	0.530	0.980	0.980	0.938	0.938
MEE	0.861	0.822	0.995	0.712	0.982	0.900	0.464	0.295	0.927	0.878	0.950	0.835	0.952	0.948	0.771	0.562	0.970	0.970	0.878	0.878
OpenKiwi-Bert	0.726	0.698	0.989	0.741	0.735	0.546	0.355	0.187	0.862	0.695	0.645	0.469	0.625	0.774	-0.126	-0.671	0.751	0.751	0.753	0.753
OpenKiwi-XLMR	0.760	0.680	0.995	0.701	0.931	0.714	0.442	0.171	0.859	0.697	0.792	0.659	0.905	0.899	0.271	-0.577	0.880	0.880	0.865	0.865
parbleu	0.834	0.774	0.986	0.838	0.970	0.833	0.562	0.342	0.877	0.744	0.908	0.801	0.958	0.953	0.624	0.398	0.971	0.971	0.939	0.939
parchr++	0.865	0.810	0.998	0.708	0.974	0.877	0.551	0.347	0.885	0.823	0.942	0.825	0.976	0.956	0.720	0.296	0.985	0.985	0.899	0.899
parsim-1	0.788	0.712	0.998	0.835	0.983	0.819	0.591	0.363	0.926	0.828	0.885	0.797	0.963	0.910	0.800	0.495	0.929	0.929	0.929	0.929
prism	0.818	0.720	0.998	0.775	0.974	0.869	0.502	0.269	0.908	0.839	0.898	0.788	0.957	0.945	0.833	0.616	0.950	0.950	0.966	0.966
sentBLEU	0.844	0.800	0.978	0.786	0.974	0.851	0.502	0.284	0.916	0.833	0.925	0.829	0.948	0.950	0.649	0.469	0.969	0.969	0.888	0.888
SWSS+METEOR	-	-	-	-	0.978	0.919	0.472	0.212	0.925	0.876	0.967	0.862	0.959	0.926	0.766	0.545	0.990	0.990	0.946	0.946
TER	0.845	0.783	0.993	0.766	0.974	0.752	0.586	0.346	0.904	0.829	0.805	0.795	0.956	0.911	0.733	0.616	0.973	0.973	0.935	0.935
WE_WPI	0.838	0.757	0.998	0.743	0.973	0.890	0.573	0.335	0.939	0.861	0.933	0.878	0.965	0.924	0.776	0.654	0.993	0.993	0.922	0.922
YiSi=0	0.876	0.825	0.998	0.786	0.972	0.867	0.453	0.207	0.938	0.874	0.968	0.861	0.956	0.918	0.831	0.563	0.986	0.986	0.932	0.932
YiSi=1	0.832	0.746	0.998	0.783	0.982	0.868	0.543	0.316	0.915	0.833	0.925	0.797	0.961	0.942	0.834	0.590	0.977	0.977	0.953	0.953
YiSi=2	0.764	0.640	0.988	0.404	0.971	0.776	0.437	0.230	0.825	0.814	0.849	0.761	0.964	0.933	0.676	0.371	0.790	0.790	0.942	0.942

表 6 : to-English におけるシステムレベルのケンコーラ

	cs-en	de-en	ja-en	pl-en	ru-en	ta-en	zh-en	iu-en	km-en	ps-en
	12	12	10	14	11	14	16	11	7	6
BAQ_dyn	–	–	–	–	–	–	0.817	–	–	–
BAQ_static	–	–	–	–	–	–	0.867	–	–	–
BERT-base-L2	0.758	0.848	0.822	0.407	0.491	0.604	0.633	0.564	1.000	0.867
BERT-large-L2	0.758	0.848	0.867	0.341	0.564	0.626	0.700	0.527	1.000	0.867
BLEU	0.848	0.697	0.778	0.407	0.455	0.692	0.833	0.309	0.714	0.600
BLEURT	0.758	0.788	0.822	0.407	0.600	0.604	0.650	0.527	1.000	0.867
BLEURT-extended	0.727	0.848	0.778	0.341	0.455	0.582	0.617	0.527	0.905	0.867
CharacTER	0.758	0.758	0.822	0.341	0.745	0.692	0.800	0.527	0.810	0.733
chrF	0.818	0.727	0.822	0.363	0.709	0.714	0.833	0.418	0.619	0.733
chrF++	0.818	0.697	0.778	0.407	0.673	0.714	0.850	0.418	0.619	0.733
COMET	0.727	0.758	0.778	0.407	0.564	0.626	0.733	0.636	1.000	0.867
COMET-2R	0.727	0.788	0.778	0.451	0.527	0.582	0.717	0.600	1.000	0.867
COMET-HTER	0.667	0.788	0.822	0.275	0.491	0.604	0.533	0.564	1.000	0.867
COMET-MQM	0.667	0.727	0.822	0.275	0.455	0.582	0.517	0.636	1.000	1.000
COMET-QE	0.697	0.788	0.778	0.297	0.455	0.516	0.550	0.491	0.905	0.733
COMET-Rank	0.576	0.727	0.822	0.341	0.455	0.626	0.650	0.309	0.810	1.000
EED	0.788	0.727	0.733	0.297	0.782	0.758	0.833	0.636	0.714	0.733
esim	0.727	0.848	0.822	0.451	0.491	0.670	0.717	0.636	1.000	0.867
IMPACT	0.727	0.818	0.689	0.319	0.673	0.802	0.833	0.587	0.714	0.733
mBERT-L2	0.758	0.818	0.822	0.429	0.564	0.604	0.750	0.673	1.000	0.867
MEE	0.758	0.697	0.867	0.363	0.709	0.692	0.783	0.636	0.714	0.733
OpenKiwi-Bert	0.697	0.667	0.733	0.187	0.455	0.429	0.450	–0.055	0.714	0.467
OpenKiwi-XLMR	0.727	0.636	0.822	0.275	0.418	0.560	0.567	0.018	1.000	0.867
parbleu	0.809	0.779	0.778	0.420	0.491	0.685	0.807	0.404	0.714	0.867
parchrf++	0.818	0.727	0.822	0.407	0.709	0.714	0.817	0.491	0.619	0.733
paesim-1	0.727	0.879	0.822	0.451	0.491	0.670	0.700	0.636	1.000	0.867
prism	0.758	0.727	0.867	0.341	0.564	0.648	0.800	0.673	0.714	0.867
sentBLEU	0.788	0.758	0.733	0.297	0.564	0.692	0.850	0.455	0.619	0.600
SWSS+METEOR	–	–	0.822	0.341	0.818	0.736	0.817	0.491	0.714	0.733
TER	0.758	0.788	0.689	0.287	0.600	0.780	0.800	0.514	0.878	0.867
WE_WPI	0.788	0.667	0.778	0.451	0.673	0.736	0.850	0.673	0.714	0.733
YiSi-0	0.758	0.758	0.689	0.231	0.782	0.802	0.833	0.600	0.714	0.733
YiSi-1	0.758	0.758	0.778	0.451	0.564	0.692	0.817	0.673	1.000	0.867
YiSi-2	0.576	0.515	0.778	0.319	0.527	0.582	0.750	0.491	0.810	0.867

Neural Information Processing Systems.

- [2] Minh-Thang Luong, Ilya Sutskever, Quoc V. Le, Oriol Vinyals and Wojciech Zaremba. 2015. Addressing the Rare Word Problem in Neural Machine Translation. Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing. pp.11-19.
- [3] Minh-Thang Luong, Hieu Pham and Christopher D. Manning. 2015. Effective Approaches to Attention-based Neural Machine Translation. Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. pp.1412-1421.
- [4] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia Polosukhin. 2017. Attention Is All You Need.

Proceedings of the 31st Conference on Neural Information Processing Systems (NIPS 2017). pp.6000-6010.

- [5] Peter F. Brown, John Cocke, Stephen A. Della Pietra, Vincent J. Della Pietra, Fredrick Jelinek, John D. Lafferty, Robert L. Mercer and Paul S. Roossin. 1990. A Statistical Approach to Machine Translation. *Computational Linguistics, Vol.16, No.2.* pp.79-85.
- [6] Peter F. Brown, Vincent J. Della Pietra, Stephen A. Della Pietra and Robert L. Mercer. 1993. The Mathematics of Statistical Machine Translation: Parameter Estimation *Computational Linguistics, Vol. 19, No.2.* pp.263-311.
- [7] Richard Zens, Franz Josef Och, and Hermann Ney. 2002. Phrase-Based Statistical Machine Translation. LNAI 2479, pp.18-32. *Springer-Verlag Berlin Heidelberg*

表 7：out-of-English におけるシステムレベルのケンドール τ

	en-cs 12	en-de 14	en-ja 11	en-pl 14	en-ru 9	en-ta 15	en-zh 12	en-iu_full 11	en-iu_news 11
BAQ_dyn	–	–	–	–	–	–	0.697	–	–
BAQ_static	–	–	–	–	–	–	0.788	–	–
BLEU	0.515	0.802	0.818	0.582	0.889	0.829	0.727	0.236	0.236
BLEURT-extended	0.879	0.802	0.782	0.780	0.833	0.771	0.848	0.382	0.345
bleurt-Yisi-combi	–	0.824	–	–	–	–	–	–	–
CharacTER	0.515	0.890	0.782	0.560	0.944	0.771	0.697	0.236	0.345
chrF	0.485	0.868	0.818	0.604	0.889	0.810	0.727	0.345	0.309
chrF++	0.485	0.868	0.782	0.604	0.889	0.829	0.727	0.309	0.309
COMET	0.909	0.846	0.745	0.736	0.722	0.771	0.606	0.382	0.382
COMET-2R	0.909	0.890	0.891	0.714	0.611	0.790	0.606	0.309	0.418
COMET-HTER	0.909	0.802	0.818	0.736	0.667	0.619	0.576	0.491	0.491
COMET-MQM	0.909	0.802	0.818	0.736	0.667	0.619	0.545	0.527	0.455
COMET-QE	0.848	0.802	0.709	0.802	0.667	0.543	0.576	0.600	0.673
COMET-Rank	0.848	0.780	0.782	0.692	0.556	0.524	0.515	0.127	0.345
EED	0.545	0.868	0.782	0.604	0.833	0.867	0.727	0.273	0.273
EQ_dyn	–	–	–	–	–	–	0.727	–	–
EQ_static	–	–	–	–	–	–	0.818	–	–
esim	0.606	0.912	0.855	0.692	0.833	0.752	0.788	0.382	0.455
IMPACT	0.515	0.736	0.818	0.560	0.722	0.867	0.697	0.257	0.236
mBERT-L2	0.788	0.846	0.782	0.736	0.778	0.752	0.909	–	–
MEE	0.576	0.802	–	0.582	0.667	0.829	–	0.273	0.382
OpenKiwi-Bert	0.758	0.780	0.236	0.538	0.722	0.314	0.606	–0.273	0.200
OpenKiwi-XLMR	0.909	0.780	0.818	0.692	0.667	0.657	0.545	0.018	0.200
parbleu	0.504	0.736	0.611	0.633	0.761	0.842	0.718	0.404	0.345
parchrF++	0.515	0.846	0.818	0.670	0.889	–	0.727	–	–
paesim-1	0.667	0.890	0.818	0.692	0.833	0.752	0.818	0.382	0.455
prism	0.818	0.868	0.818	0.670	0.611	0.562	0.576	0.418	0.600
sentBLEU	0.515	0.802	0.855	0.604	0.944	0.867	0.727	0.236	0.273
TER	0.515	0.824	0.018	0.641	0.556	0.752	0.242	0.309	0.309
WE_WPI	0.515	0.780	0.818	0.495	0.833	0.771	0.667	–	–
YiSi-0	0.545	0.846	0.818	0.604	0.944	0.790	0.515	0.236	0.345
YiSi-1	0.606	0.868	0.782	0.626	0.833	0.810	0.758	0.091	0.273
YiSi-2	0.485	0.582	0.527	0.077	0.444	0.886	0.121	0.309	0.455
Yisi-combi	–	0.824	–	–	–	–	–	–	–

表 8：相関係数の平均

	システムレベル				セグメントレベル			
	to-English		out-of-English		to-English		out-of-English	
	All	-out	All	-out	All	-out	All	-out
BLEU	0.765	0.652	0.741	0.505	0.139	0.034	0.302	0.108
CharacTER	0.793	0.667	0.833	0.474	0.173	0.078	0.339	0.172
chrF	0.789	0.651	0.803	0.512	0.183	0.084	0.387	0.211
chrF++	0.791	0.661	0.792	0.515	0.184	0.084	0.383	0.204
EED	0.804	0.698	0.841	0.501	0.187	0.088	0.385	0.203
IMPACT	0.785	0.689	0.814	0.529	0.152	0.049	0.328	0.124
TER	0.777	0.673	0.490	0.473	0.024	–0.090	–0.055	–0.060
MEE	0.787	0.683	–	–	0.114	0.034	–	–
WE_WPI	0.799	0.696	–	–	0.199	0.118	–	–

- [8] K. Papineni, S. Roukos, T. Ward, and Wei-Jing Zhu. 2002. BLEU: a Method for Automatic Evaluation of Machine Translation. Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics. pp.311-318.
- [9] Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. Proceedings of the Association for Machine Translation in the Americas. pp.223-231.
- [10] A. Lavie and A. Agarwal. 2007. Meteor: An Automatic Metric for MT Evaluation with High Levels of Correlation with Human Judgments. Proceedings of the Second Workshop on Statistical Machine Translation. pp.228-231.
- [11] Hideki Isozaki, Tsutomu Hirao, Kevin Duh, Katsuhito Sudoh, and Hajime Tsukada. 2010. Automatic Evaluation of Translation Quality for Distant Language Pairs. Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing. pp.944-952.
- [12] Hiroshi Echizen-ya and Kenji Araki. 2007. Automatic Evaluation of Machine Translation based on Recursive Acquisition of an Intuitive Common Parts Continuum. Proceedings of the Eleventh Machine Translation Summit. pp.151-158.
- [13] Hiroshi Echizen-ya, Terumasa Ehara, Sayori Shimohata, Atsushi Fujii, Masao Utiyama, Mikio Yamamoto, Takehito Utsuro and Noriko Kando. 2009. Meta-Evaluation of Automatic Evaluation Methods for Machine Translation using Patent Translation Data in NTCIR-7. Proceedings of the 3rd Workshop on Patent Translation pp.9-16.
- [14] Hiroshi Echizen-ya and Kenji Araki. 2010. Automatic Evaluation Method for Machine Translation using Noun-Phrase Chunking. Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics. pp.108-117.
- [15] Hiroshi Echizen'ya, Kenji Araki, and Eduard Hovy. Optimization for Efficient Determination of Chunk in Automatic Evaluation for Machine Translation. Proceedings of the 1th International Workshop on Optimization Techniques for Human Language Technology. pp.17-30.
- [16] Hiroshi Echizen'ya, Kenji Araki, and Eduard Hovy. 2013. Automatic Evaluation Metric for Machine Translation that is Independent of Sentence Length. Proceedings of the 9th Recent Advances in Natural Language Processing. pp.230-236.
- [17] Hiroshi Echizen'ya, Kenji Araki, and Eduard Hovy. 2014. Application of Prize based on Sentence Length in Chunk-based Automatic Evaluation of Machine Translation. Proceedings of the Ninth Workshop on Statistical Machine Translation. pp.381-386.
- [18] A. Apostolico and C. Guerra. 1987. The Longest Common Subsequence Problem Revisited. *Algorithmica, Volume 2, issue 1-4*. pp.315-336, Springer.
- [19] Hiroshi Echizen'ya, Kenji Araki, and Eduard Hovy. 2019. Word Embedding-Based Automatic MT Evaluation Metric using Word Position Information. *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp.1874-1883.
- [20] Yossi Rubner, Carlo Tomashi and Leonidas J. Guibas. 1998. A Metric for Distributions with Applications to Image Database. Proceedings of the 1998 IEEE International Conference on Computer Vision. pp.59-66.
- [21] Yossi Rubner, Carlo Tomashi and Leonidas J. Guibas. 2000. The Earth Mover's Distance as a Metric for Image Retrieval. *International Journal of Computer Vision 40 (2), pp.99-121* Kluwer Academic Publishers.
- [22] 柳本豪一, 大松繁. Earth Mover's Distance を用いたテキスト分類. 2007. The 21st Annual Conference of the Japanese Society for Artificial Intelligence. 1G3-4.
- [23] Loic Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. 2020. Findings of the 2020 Conference on Machine Translation (WMT20). Proceedings of the 5th Conference on Machine Translation (WMT). pp.1-55.
- [24] Nitika Mathur, Johnny Tian-Zheng Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020. Results of the WMT20 Metrics Shared Task. Proceedings of the 5th Conference on Machine Translation (WMT). pp.688-725.
- [25] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers). pp.4171-4186.
- [26] Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. Proceedings of the Tenth Workshop on Statistical Machine Translation. pp.392-395.

- [27] Maja Popović. 2017. chrF++: words helping character n-grams. Proceedings of the Second Conference on Machine Translation. pp.612-618.
- [28] Rachel Bawden, Biao Zhang, Andre Tättar, and Matt Post. 2020. ParBLEU: Augmenting metrics with automatic paraphrases for the WMT'20 metrics shared task. Proceedings of the Fifth Conference on Machine Translation (Volume 2: Shared Task Papers). pp.887-894.
- [29] Brian Thompson and Matt Post. 2020. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing. pp.90-121.
- [30] Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. CharacTer: Translation edit rate on character level. Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers. pp.505-510.
- [31] Peter Stanchev, Weiyue Wang, and Hermann Ney. 2019. Eed: Extended edit distance measure for machine translation. Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers). pp.514-520.
- [32] Jin Xu, Yinuo Guo, and Junfeng Hu. 2020. Incorporate semantic structures into machine translation evaluation via ucca. Proceedings of the 5th Conference on Machine Translation (WMT). pp.934-939.
- [33] Manish Shrivastava Ananya Mukherjee, Hema Ala, and Dipti Misra Sharma. 2019. Mee: An automatic metric for evaluation using embeddings for machine translation. Proceedings of the 2020 IEEE 7th International Conference on Data Science and Advanced Analytics. pp.292-299.
- [34] Piotr Bojanowski, Edouard Grave, Armand Joulin and Tomas Mikolov 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics, Volume 5*. pp.135-146.
- [35] Chi-kiu Lo. 2019. YiSi - a unified semantic MT quality evaluation and estimation metric for languages with different levels of available resources. Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers). pp.507-513.
- [36] Chi-kiu Lo. 2020. Extended study on using pre-trained language models and YiSi-1 for machine translation evaluation. Proceedings of the Fifth Conference on Machine Translation: Shared Task Papers. pp.895-902.
- [37] Thibault Sellam, Amy Pu, Hyung Won Chung, Sebastian Gehrmann, Qijun Tan, Markus Freitag, Dipanjan Das, and Ankur Parikh. 2020. Learning to evaluate translation beyond english: Bleurt submissions to the wmt metrics 2020 shared task. Proceedings of the Fifth Conference on Machine Translation. pp.921-927.
- [38] Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. Unbabel's participation in the wmt20 metrics shared task. Proceedings of the Fifth Conference on Machine Translation. pp.911-920.
- [39] Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. 2019. OpenKiwi: An open source framework for quality estimation. Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations. pp.117-122.
- [40] Chi-kiu Lo and Samuel Larkin. 2020. Machine translation reference-less evaluation using YiSi-2 with bilingual mappings of massive multilingual language model. Proceedings of the Fifth Conference on Machine Translation. pp.903-910.
- [41] Evan James Williams. 1959. Regression analysis. wiley.