

機械翻訳に対するプリエディットのタスク分類に関する考察

株式会社アスカコーポレーション
大阪大学大学院 情報科学研究科 マルチメディア工学専攻
早川威士

- プリエディット (pre-editing) ¹
 - MTに掛ける前の原文を修正する
 - 原文をMTが処理しやすい文章にすることで、MT出力を改善できるだろうという発想
 - 単言語の作業
- ポストエディットとの比較
 - ポストエディットは原文と一致した情報にすればいいので、正確性が保証できていればいい
 - プリエディットは何をしたらいいかの基準がない
 - 原文を修正してもMT出力が改善されるかどうかはわからない
 - 何もしなくてもMT出力には問題がないかもしれない

1. Johnson, R., & Whitelock, P. (2003). Machine translation as an expert task. Readings in Machine Translation, 233.

- 日→英翻訳におけるリエディット作業の有効性
 - リエディットはMTのエラーを未然に防ぐことができるか
 - またその有効性は定量的に表現できるか
- リエディットが必要な日本語原文の問題箇所の特定制
 - リエディットが必要な／不要な箇所を原文の特性にもとづいて分類することで、事前に知ることができるか
- リエディットに求められる能力の特定制
 - リエディットを行う上で言語的スキルはどの程度求められるか

問題箇所とタスク

- MTに依存する問題
 - 人間には理解できるが機械には理解困難
- 自然言語に帰属する問題
 - 言語間に普遍的に存在する問題
 - 言語の運用に起因する問題

タスク分類

- プリエディットのタスクを11パターンに分類
 - MT依存的な問題：e1～e4
 - 言語的な問題：e5～e8
 - 言語運用の問題：e9～e11

No	修正内容	発生頻度 (文当たり)
e1	主語や目的語の脱落を補う	23%
e2	指示語を具体化する	12%
e3	複文を分割する/入れ子構造を外す	19%
e4	文構造を明確化する	28%
e5	固有名詞・新語の事前置換	19%
e6	専門用語の事前置換	13%
e7	日本固有表現の書き換え	8%
e8	婉曲表現の書き換え	11%
e9	誤字、脱字の修正	8%
e10	英語文法に沿った書き換え	9%
e11	翻訳不要語句の削除	8%

- 使用した例文

- ウェブサイトからパラレルコーパスを取得可能な100文対
 - 専門性の高い文として医学・医薬品関連の文章を使用
- 日本語文に対してMTを適用（プリエディットの対象文）

- 作業者

- 言語スキルの有無による差を測定できるように作業者を設定
 - 医薬分野の翻訳専門家：2名＋非翻訳者：2名

- MTシステム

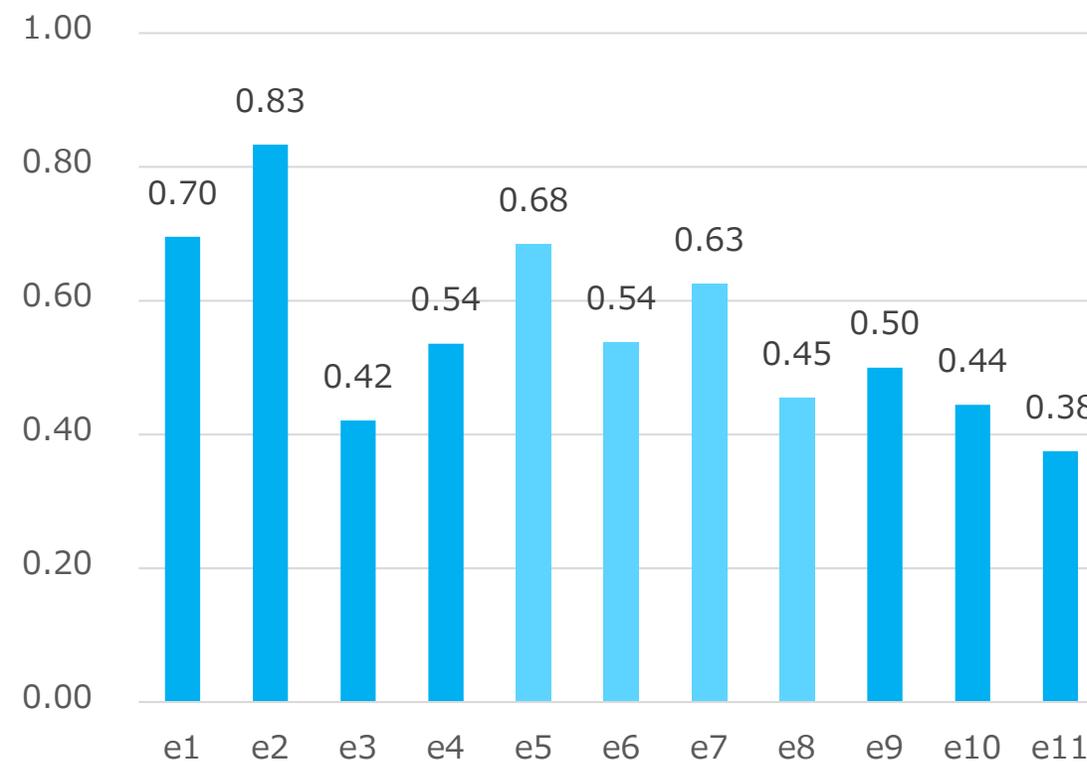
- 情報通信研究機構（NICT）みんなの自動翻訳 @TexTra
 - 汎用NMT：2019年8月時点のモデルを適用

- **プリエディット対象把握の正確性**
 - 原文の問題点を把握し、適切にエディットできているか
 - 原文に存在する問題点のうち、プリエディットされた割合を測定（プリエディット率）
- **プリエディットによるMTエラーの防止**
 - プリエディットによって実際にMTのエラーが防止されたか（エラー防止率）
 - 問題箇所への対応を以下の4つに分類し、エラーとプリエディットの関連を測定
 - ① プリエディットを行い、エラーを防止できた
 - ② プリエディットしたが、エラーになった
 - ③ プリエディットしなかったが、エラーにならなかった
 - ④ プリエディットを行えず、エラーになった

結果：ベースライン（MTの生出力）のエラー

- 原文の問題箇所が実際にMTによってエラーになった割合
 - 割合の範囲は38%～83%
 - MT依存的な問題は実際に出力のエラーになりやすい

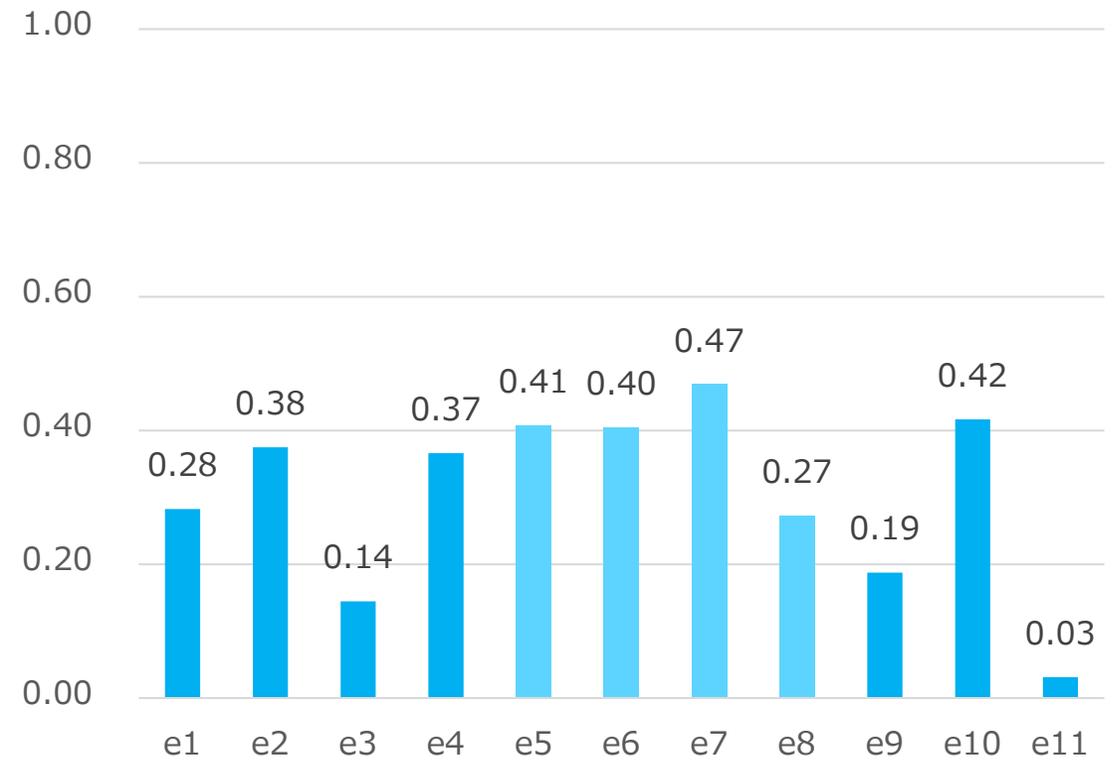
ベースラインのエラー率



結果：プリエディットによるエラー防止

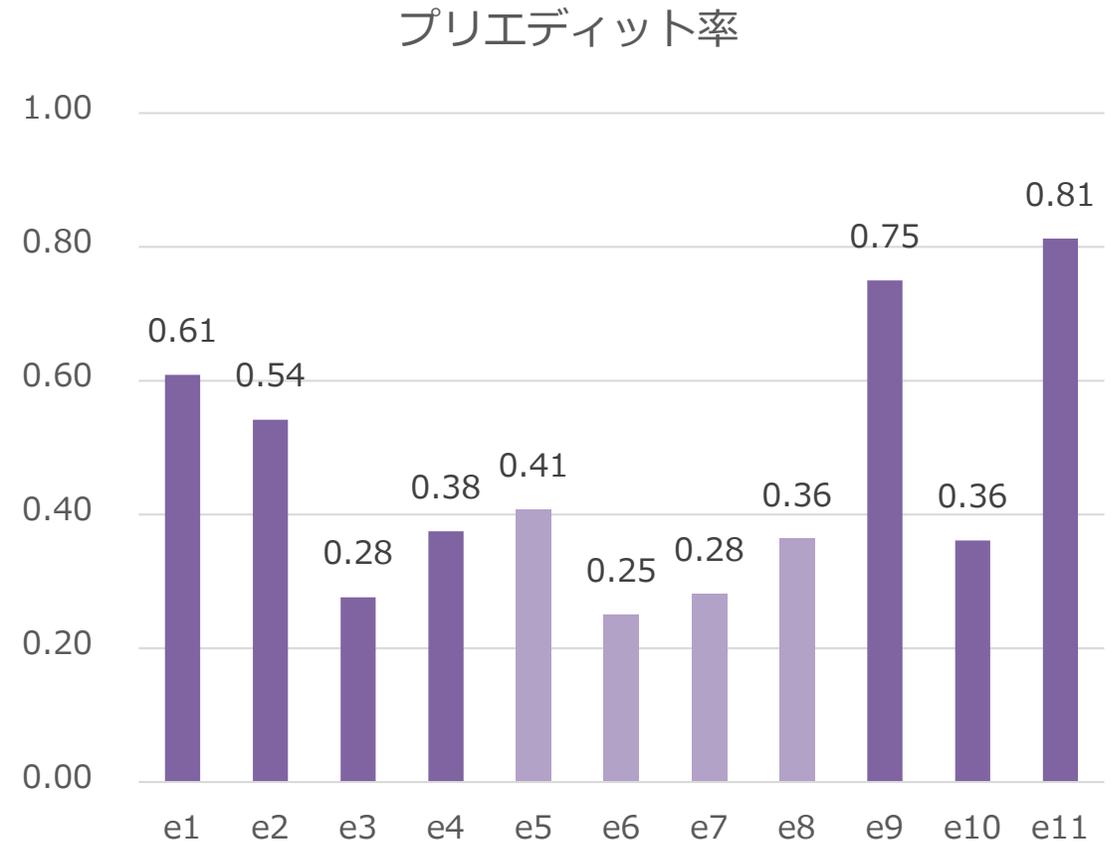
- 原文の問題箇所に対してプリエディットを行った後のエラー率
 - 実際にエディットしたかどうかを問わず、出力がエラーになった割合
 - 複文、入れ子構造（e3）と翻訳不要箇所（e11）のエラーはよく抑制されている

プリエディット後のエラー率



結果：プリエディット率

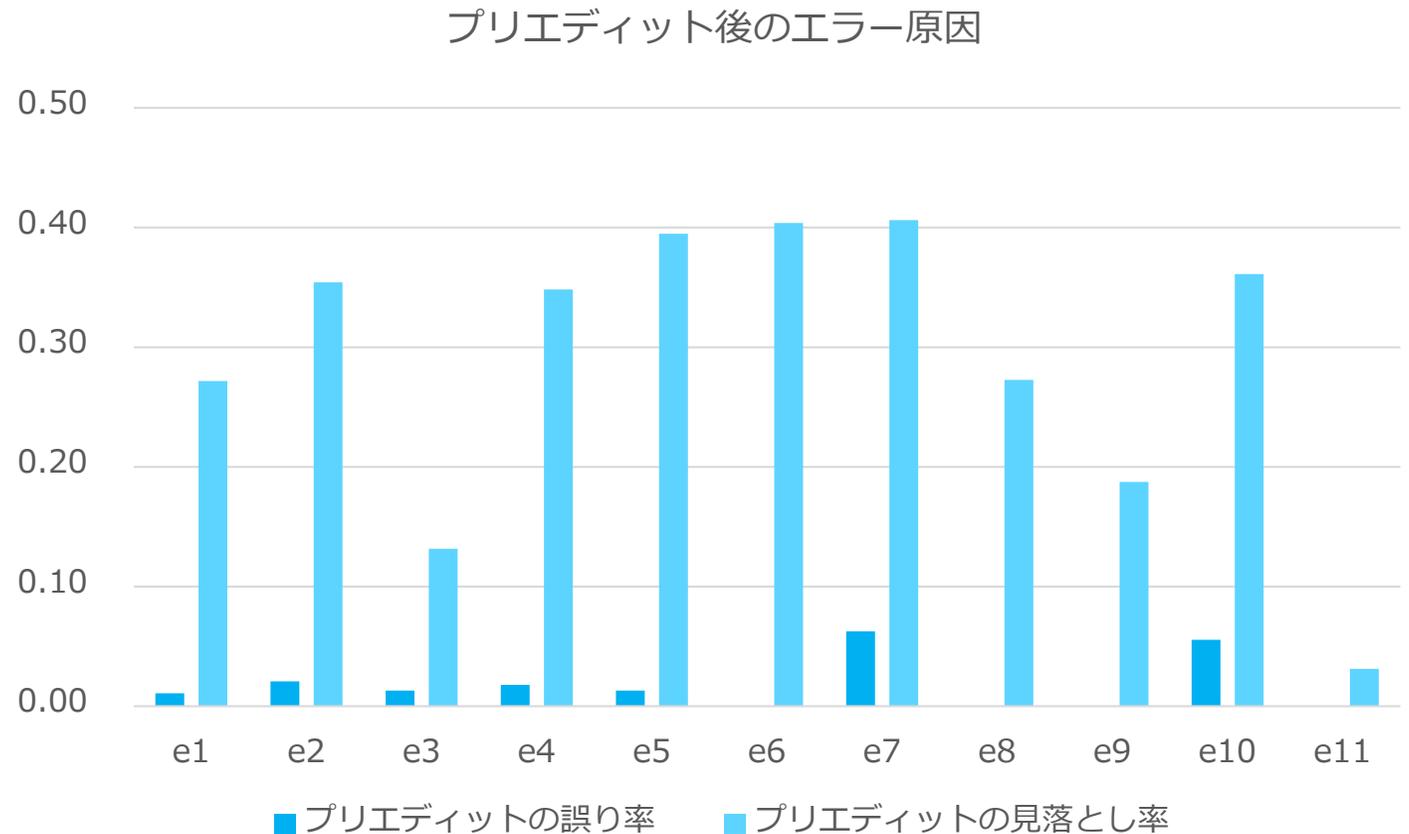
- 問題箇所に対してプリエディットを実施できた割合
 - 問題の検出力と編集作業の難易度を反映
 - 最低値：25%
(e2, 専門用語の事前置換)
最高値：81%
(e11, 翻訳不要語句の削除)
 - タスクの内容によって作業難易度の幅が大きい



結果：プリエディット後のエラー原因

• プリエディットの誤りと見落とし

- プリエディットの誤り
= エディットをしたが
エラーになった
- プリエディット見落とし
= エディットすべき
箇所の見落とし
- プリエディットの誤り率
は非常に低く、見落とし
の方が起こりやすい



結果：言語的スキルの貢献度

- 翻訳専門家と非専門家のエラー防止率

	e1	e2	e3	e4	e5	e6	e7	e8	e9	e10	e11
翻訳者	0.70	0.50	0.29	0.59	0.63	0.35	0.25	0.41	0.81	0.39	0.75
非翻訳者	0.50	0.50	0.26	0.16	0.16	0.15	0.19	0.32	0.69	0.22	0.88

- タスクによっては (e4, e5) 専門スキルによって防止率に差がある
- 一方で、それほど差がないタスク (e9, e11) もある

- **プリエディット作業の有効性**
 - 原文に問題があると30%以上の割合で訳文にもエラーが起こるため、事前の対処が必要
 - プリエディット後の誤りはエディットの不適切さではなく、ほぼ見落としが原因
 - プリエディットはエラー防止にある程度有効な対策となる
- **潜在的な原文の問題の把握**
 - MTのエラーが起こりやすい原文の問題点は分類可能である
 - NMTは言語運用上の問題に対して比較的頑健であるが、文脈判断 (e1, e2) などMT (言語処理) の課題に対してはエラー頻度が高かった
 - 文脈判断 (e1, e2)、表記上の問題 (e9, e11) についてはプリエディットによる介入効果が高かった

- プリエディット作業に必要な能力

- 作業者を言語スキルを持つ翻訳者とそうでない非翻訳者に層別化することで、タスクへの適性を視覚化することができた
- 文構造の明確化（e4）や固有名詞の認識（e5）は翻訳の仕事で常に意識されていることであり、その経験の差がプリエディット率に反映されている
- 一方で、誤字脱字の修正（e9）、翻訳不要語句の削除（e11）などのタスクは言語スキルを問わず達成可能であった
- エラーの原因となり得、かつプリエディット作業も困難なタスクについては別のアプローチの検討も必要

ご清聴ありがとうございました



コトバも
医療技術だと考える