

Challenges in building a benchmark of linguistic minimal pairs for low resource languages: The case of Malay and Indonesian

野元 裕樹¹ スリ ブディ レスタリ²
 ファルハン アティラ ビンティ アブドゥル ラザク¹
¹東京外国語大学 ²立命館アジア太平洋大学

ダヴィド ムルヤディ³ 稲垣 和也⁴ 降幡 正志¹
³神田外語大学 ⁴南山大学

【謝辞】本研究はJSPS科研費JP23K25336の助成を受けたものです。

概要

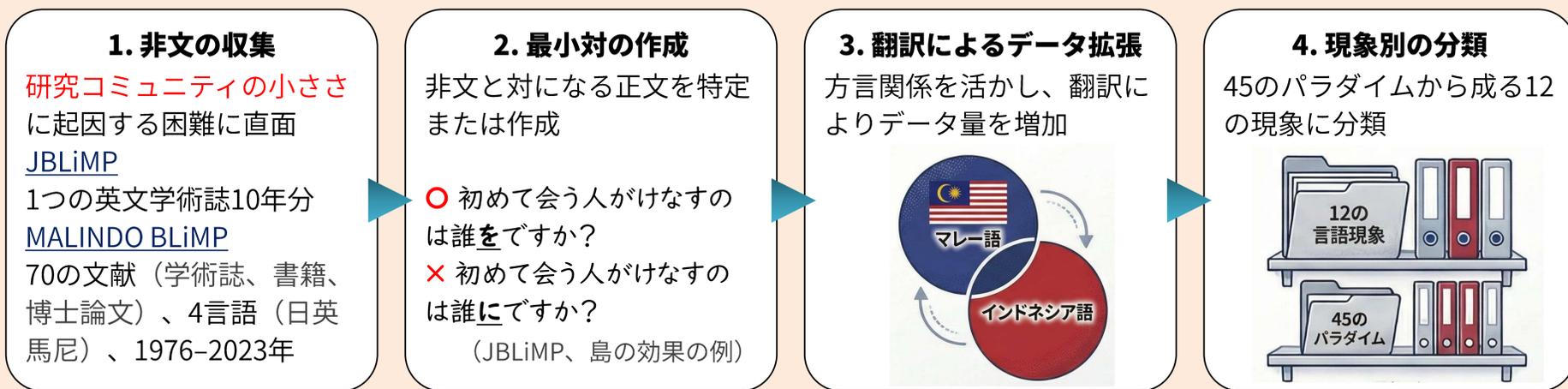
- 現在、マレー語・インドネシア語のBLiMP (Benchmark of Linguistic Minimal Pairs) MALINDO BLiMPを構築中
 - https://github.com/matbahasa/MALINDO_BLiMPにて公開
 - 可能な限り日本語のJBLiMP (Someya & Oseki 2023) に準拠
 - 低資源言語ならではの困難と工夫を報告
- 【参考：BLiMPとは？】
- 言語モデルの文法知識を調べるためのベンチマークの1つ
 - 対象文法現象においてのみ異なる正文・非文のペア（最小対）から成る



本研究のポイント

- 低資源言語でのBLiMP構築では
- 研究コミュニティが小さいため、データの収集・検証が大変
 - 外国人のみの著者による研究が多く、データの信頼性が低い
 - 関連言語間での翻訳によりデータ量不足を補うことが可能

MALINDO BLiMPの構築手法



母語話者に対する容認性判断実験によるコアデータの検証

JBLiMP

- クラウドソーシング「ランサーズ」を利用
- 240名の母語話者、367ペア（22~23ペア/人）
- 2択（どちらがより文法的？）
- 採択基準：
原典と一致 > 50%
- 採択率：90.9% (331/367)

MALINDO BLiMP

- 適切なクラウドソーシング提供会社なし→現地大学の協力+Google Forms
- 308名の話者（母語話者280名）、300ペア（30ペア/人）
- 4択：Aのみ、Bのみ、both、neither
- 採択基準：both + neither (コントラストなし) < 50%
 ∧ same (一致) + both, reverse (逆) + neither ≥ 55%
- 採択率：マレー語 58.0% (174/300)、インドネシア語 63.0% (189/300)

言語学での利用も視野に、高粒度の選択肢



著者の背景の違いが与える影響

JBLiMP

- 大半の原典が現地人（+外国人）による執筆 23/28

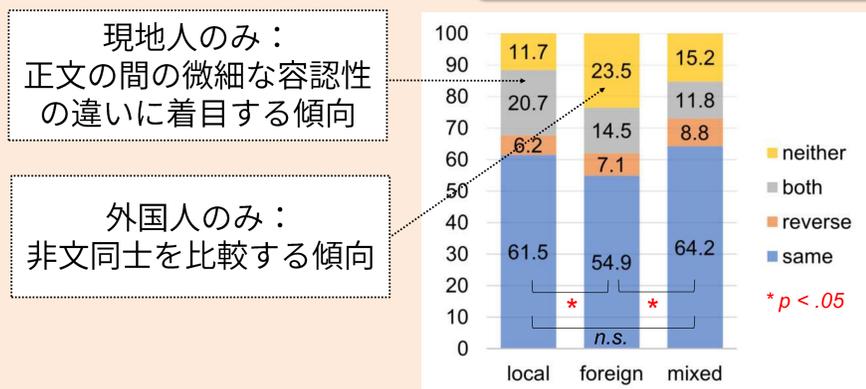
MALINDO BLiMP

- 半数以上の原典が外国人のみによる執筆 39/70
- 現地人のみ 21/70、現地人+外国人 10/70

Q: 外国人が著者の場合、データの信頼性は下がらない？

A: 「現地人のみ」および「現地人+外国人」に比べ、「外国人のみ」は一致率が有意に低くなる

現地研究者との協働は大切！



翻訳が与える影響

Q: 翻訳により拡張したデータの信頼性は下がらない？

A: 信頼性は下がらない

関連言語間での翻訳は有用なデータ拡張手法

