

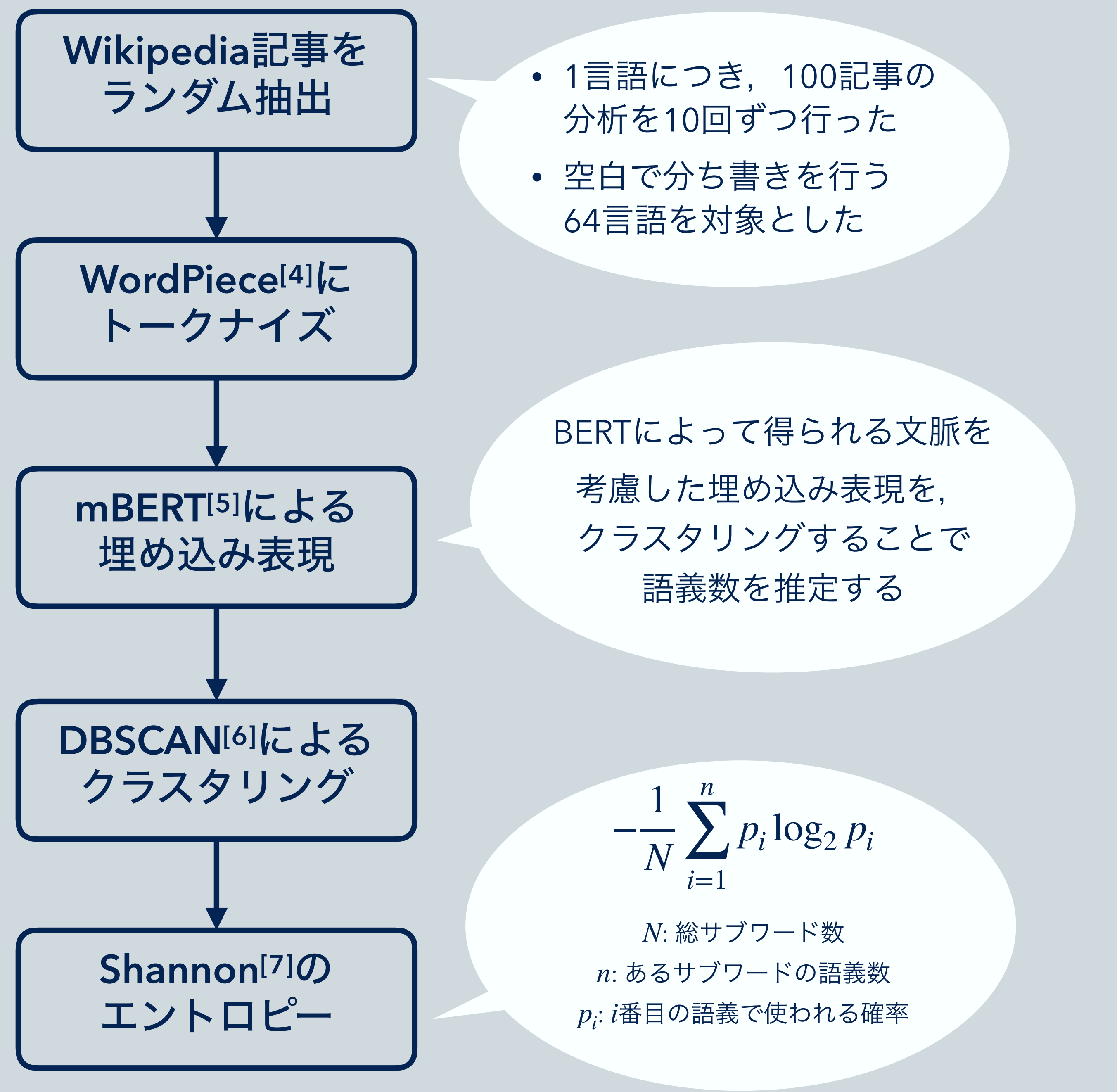
1. 目的

- 形式-意味対応の複雑性(=多義性)を言語の複雑性の1つの側面と捉え、その言語間比較を試みる
- BERTによるサブワードトークンへの分割とその分散表現を利用して、特定の言語に依存しない単位を基準とした計測を試みる

2. 背景

- 言語の等複雑性は20世紀以来、言語学者の間で受容されてきた^{[1][2]}一方で、経験的に信じられてきた一種のドグマであると指摘されている^[3]
- 妥当な言語間比較に関して、総意が得られていない
- これまでの計測手法の対象は、言語の形式的な側面に限られており、意味の側面は考慮されてこなかった

3. 方法論



4. 結果

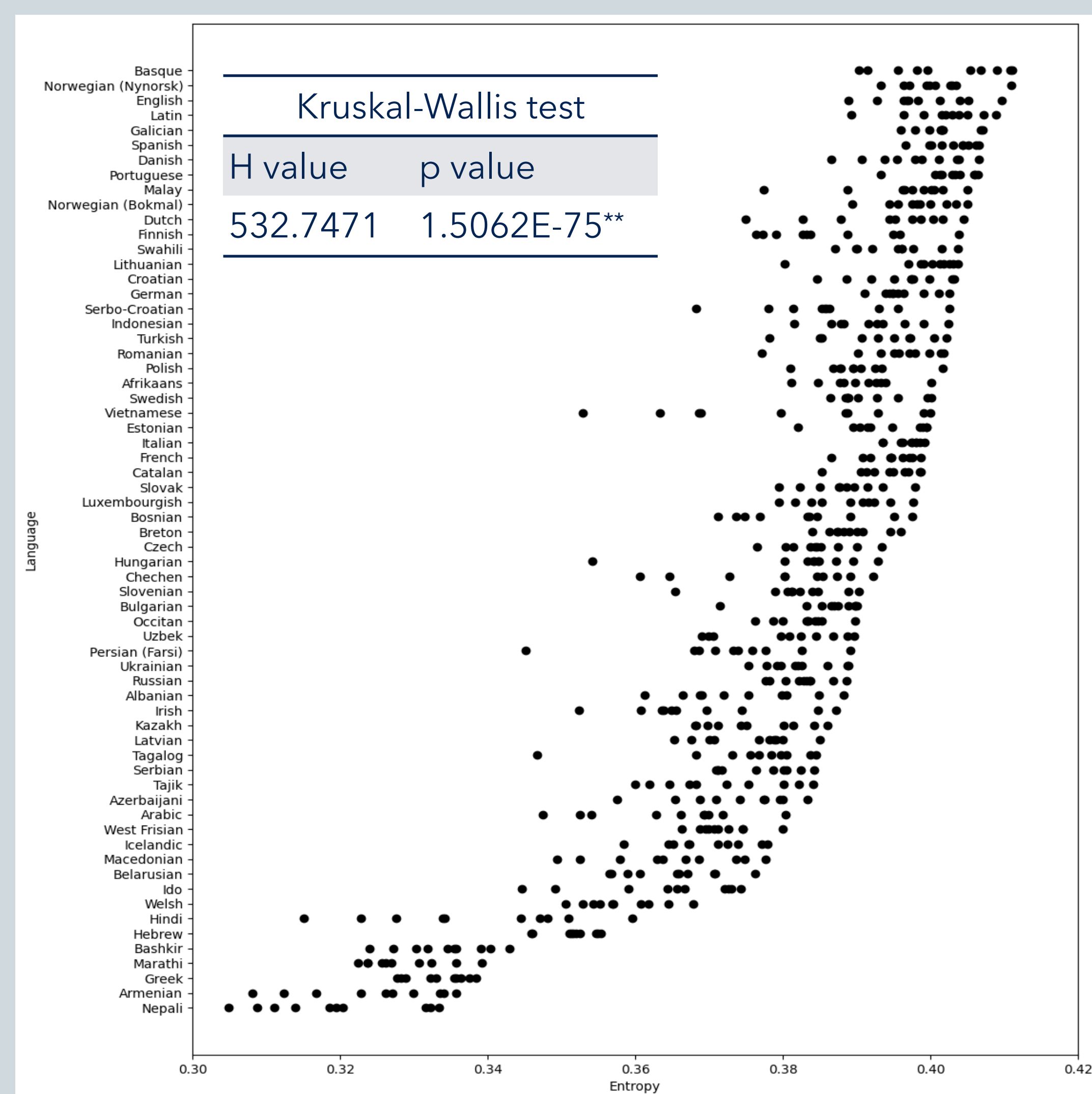


図1. 各言語のエントロピー

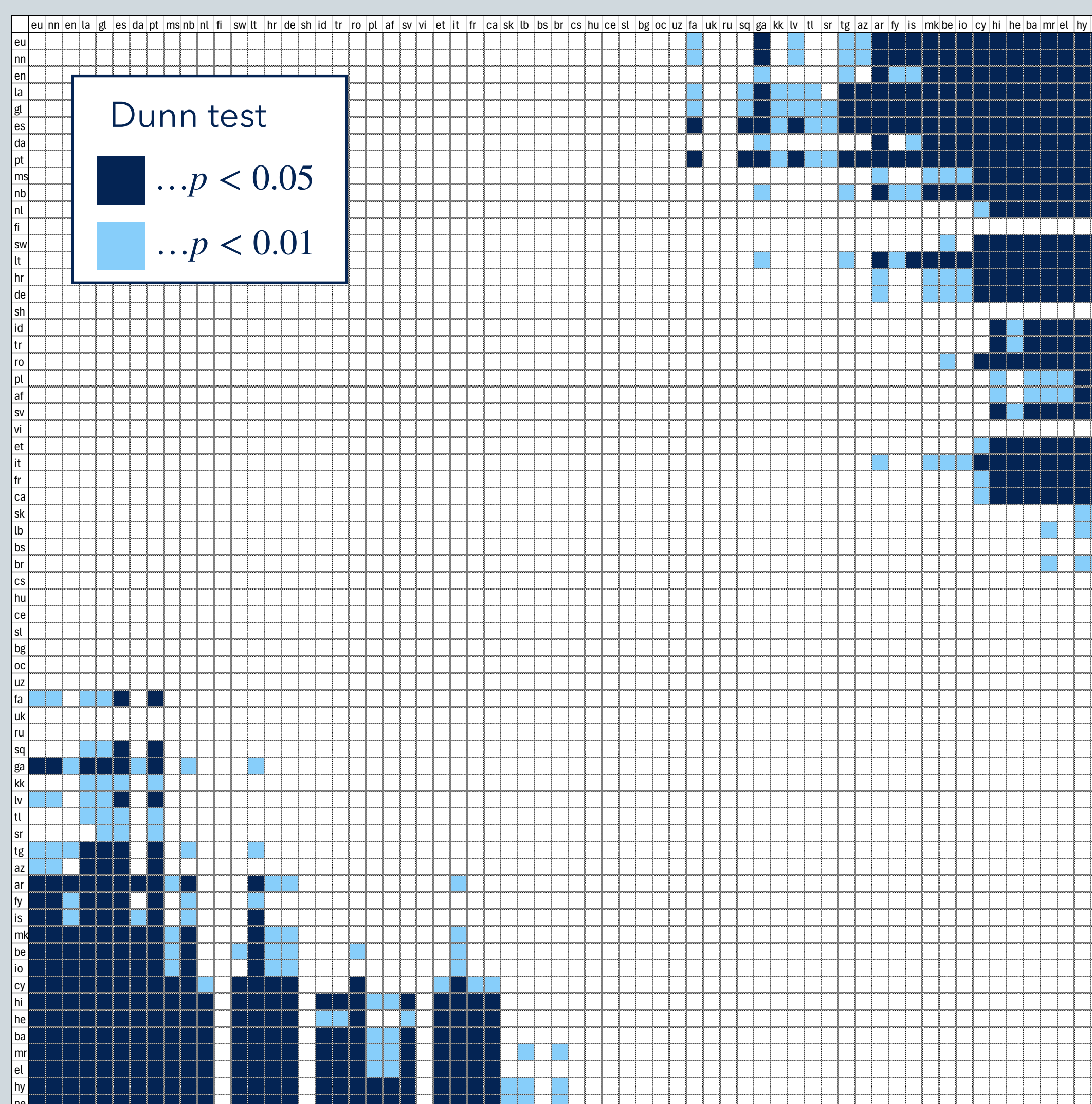


図2. 各言語間の分布の有意差

【図1】

- バスク語が最大
(≈ 0.411 , $\approx 1.329/\text{sw}$)

- ネパール語が最小
(≈ 0.305 , $\approx 1.235/\text{sw}$)

- 狭い範囲に全ての言語が収まっている

- 一方で、Kruskal-Wallis検定では有意差あり

【図2】

- Dunn検定により各言語間の有意差を見ると、上位と下位の言語間に有意差あり

5. 考察

- 「多義性が狭い範囲に収まっている」→等複雑性を支持
 - 「その範囲内で有意差が見られる」→等複雑性を不支持
- これまでの研究が、言語の等複雑性に関して統一的な見解が得られない要因の可能性

6. 結論

- 各言語は狭い多義性の範囲に収まっている事が示唆された
- 同時に、その範囲において各言語間には差異があることも示唆された

【今後の課題】

- 分ち書きを行わない言語への適用可能性を探る
- 語義を離散的に扱っていることの妥当性が自明ではない