



Is Twain's 'Awful German Language' really that Awful?

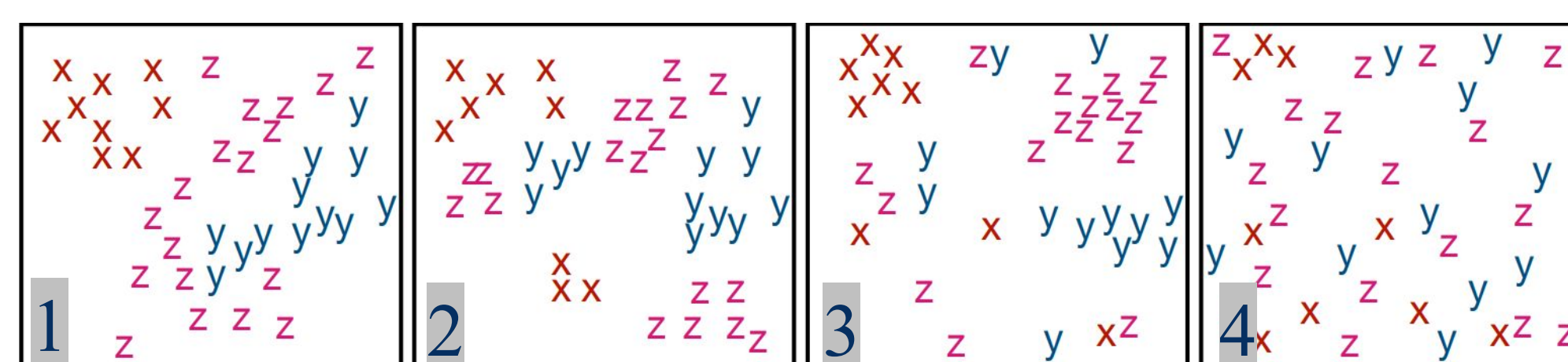
Quantifying the Randomness of Linguistic Gender in 19 Languages

Background and Motivation

- Grammatical Gender:**
 - Linguistic systems where nouns assigned 2 gender categories (e.g masc, fem, neut)
 - Unclear Semantics in gender assignment:**
 - Rarely intuitive
 - But also sometimes clearly organized
 - What's going on?
 - Systematic Issues in Noun Class Literature:**
 - 84% of linguistics literature on English + Indo-European languages (Kidd & Garcia, 2021)
 - Gender findings from Indo-Euro Studies may not generalize.
 - Need more quantitative studies of language
- Examples from Spanish:**
- 'La biblioteca vieja'*
ART.fem library old.fem
'The Old Library'
- BUT*
- 'The Old Library'*
- &*
- Word Vectors*
- Large arrays of numbers that encode meaning. (Bojanowski at al., 2017).
 - Trained w/ lots of written text
 - Words closer together in space closer in meaning.
- 2D Word Embedding of Menu Items**
-

Aims and Hypotheses

- Aim:** Quantify relationship between Gender Assignment + Semantics in as many Language Families as possible
- H1:** Effect Size between Gender and the Gender of semantically-close words is different than 0.
- H2-5:** Gender in Clusters of semantically-close words would follow 1 of 4 patterns:
 - Semantically Pure
 - Pure but outliers assigned systematically
 - Pure but outliers are assigned randomly
 - No different than chance

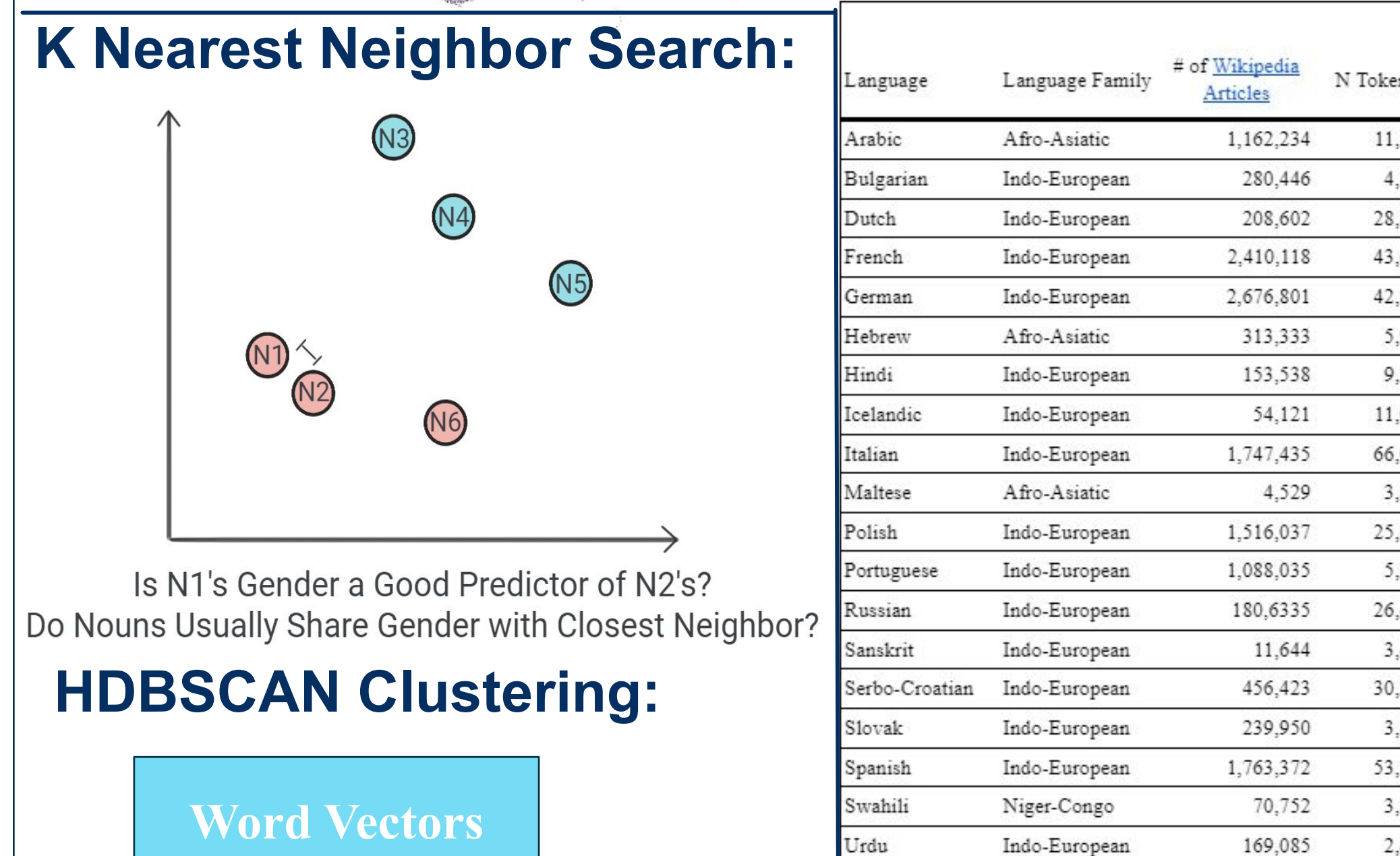
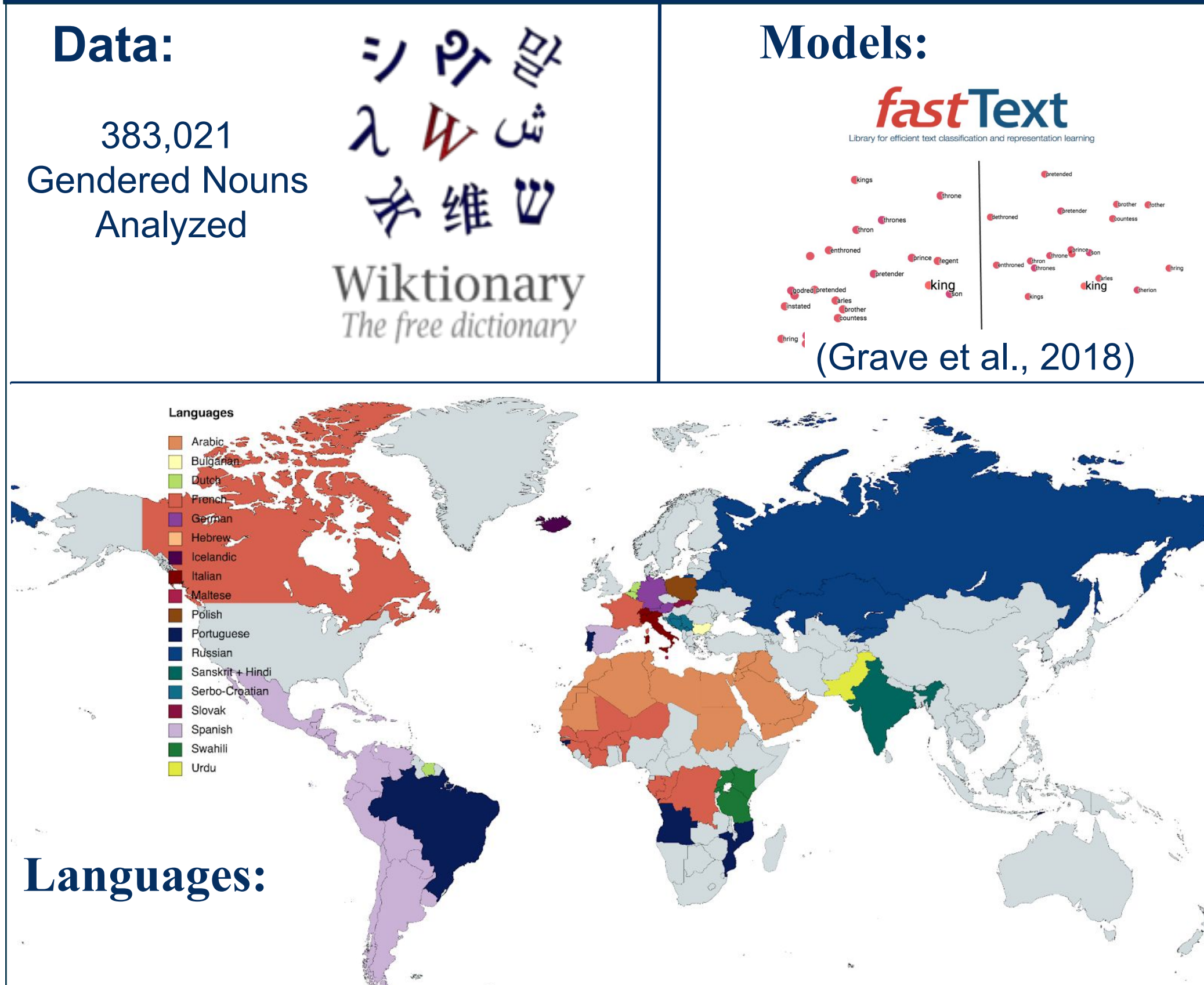


Ethan Amato¹

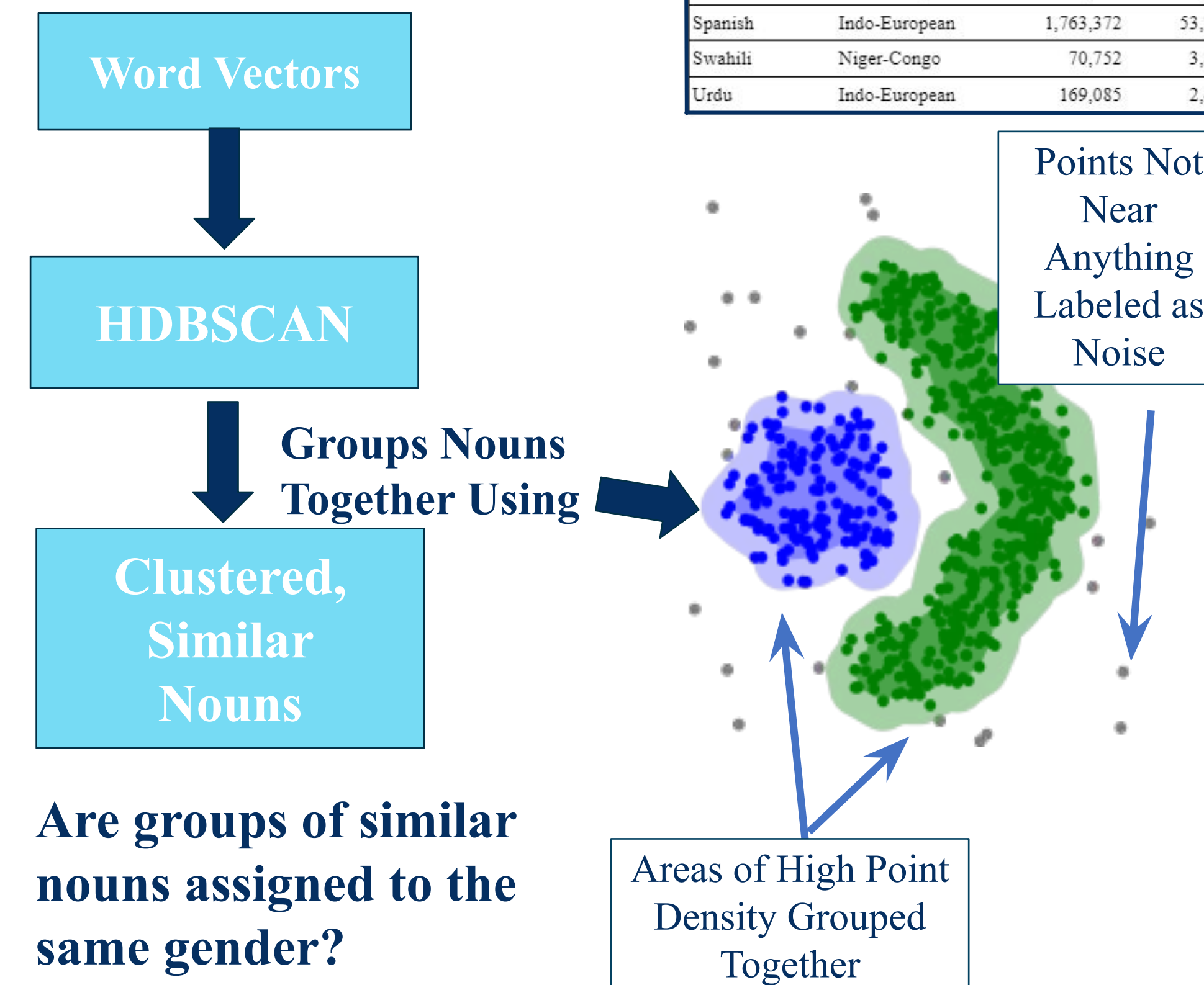
Joshua K. Hartshorne¹

¹Department of Psychology, Boston College

Method

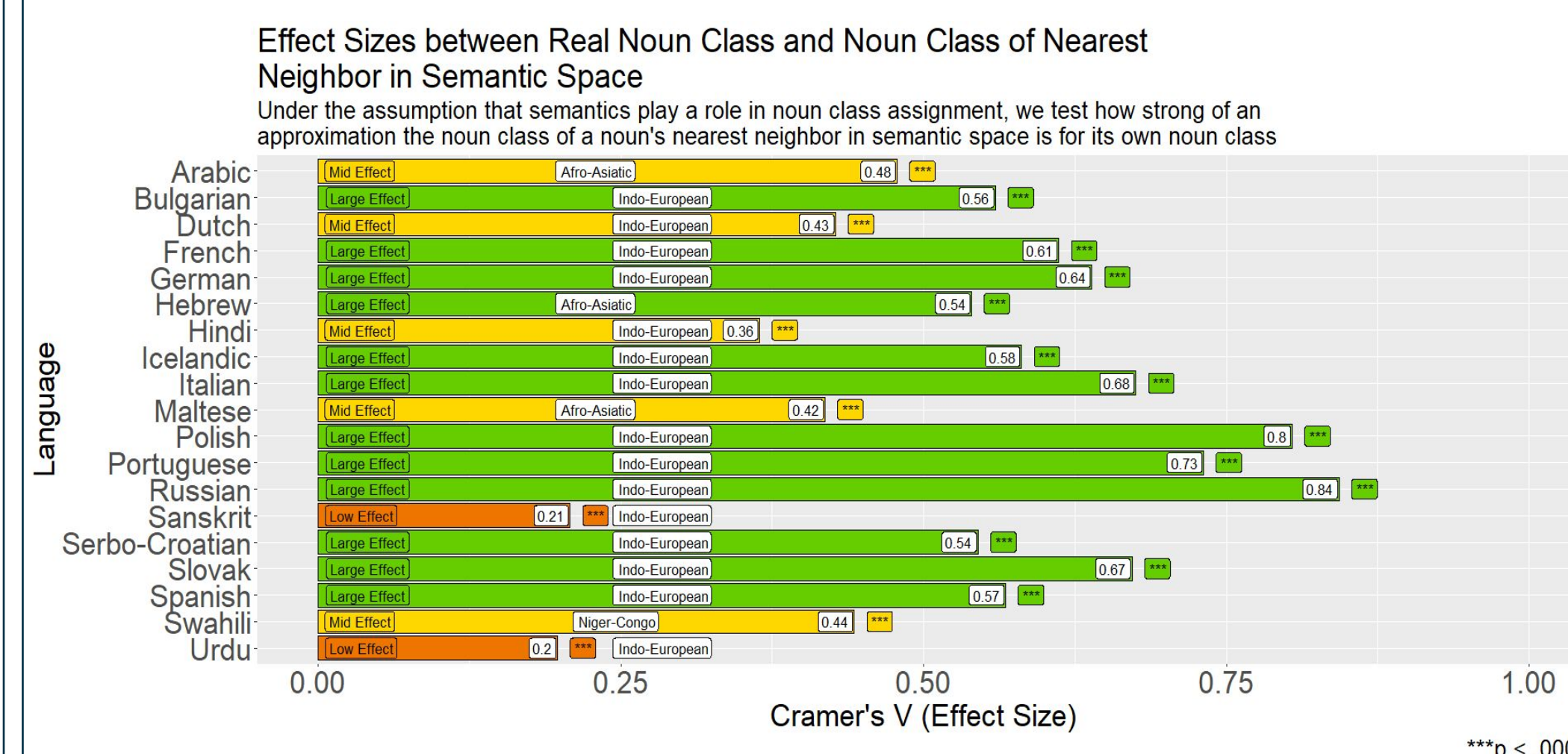


HDBSCAN Clustering:



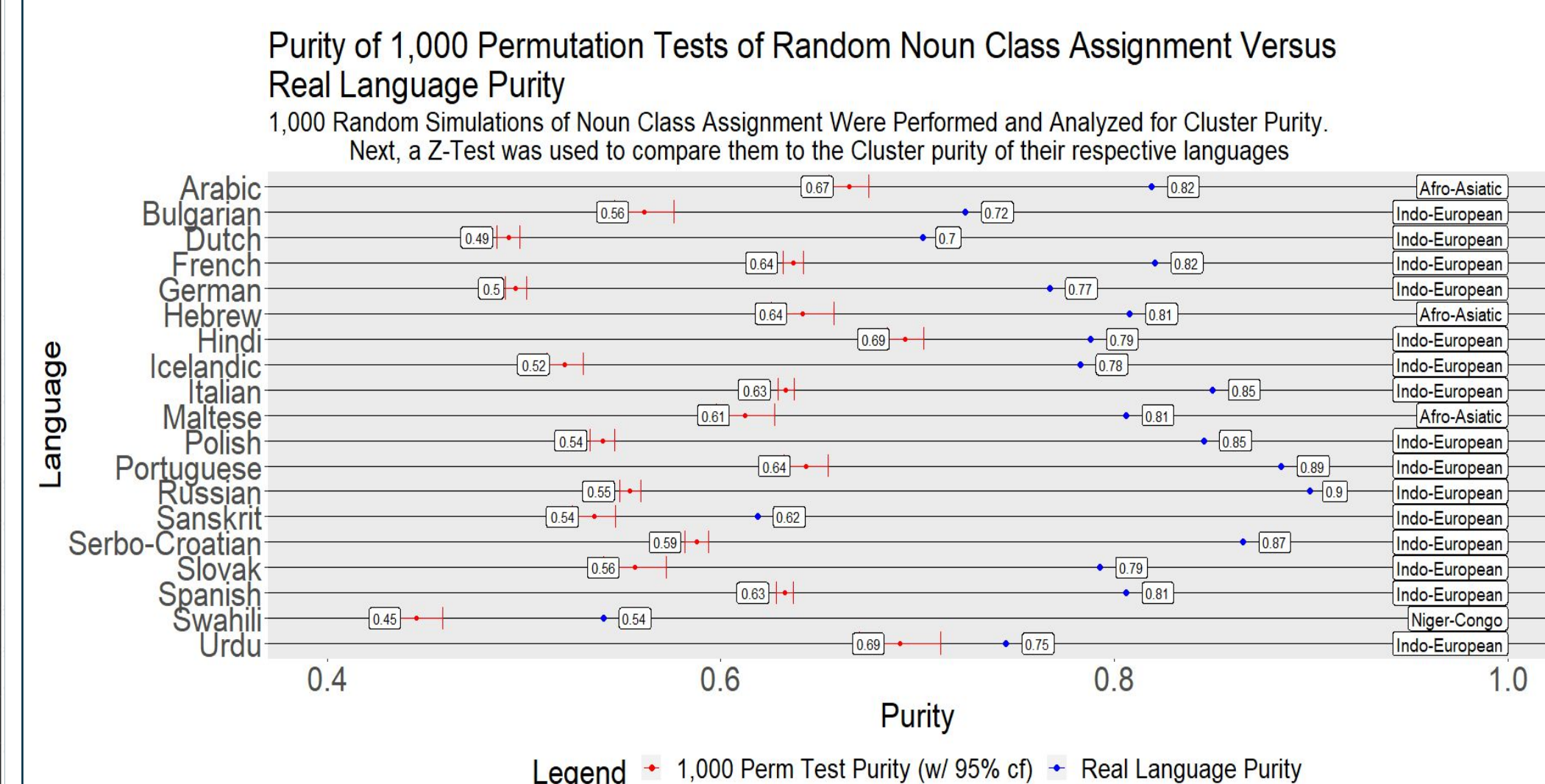
Results

Nearest Neighbor Search

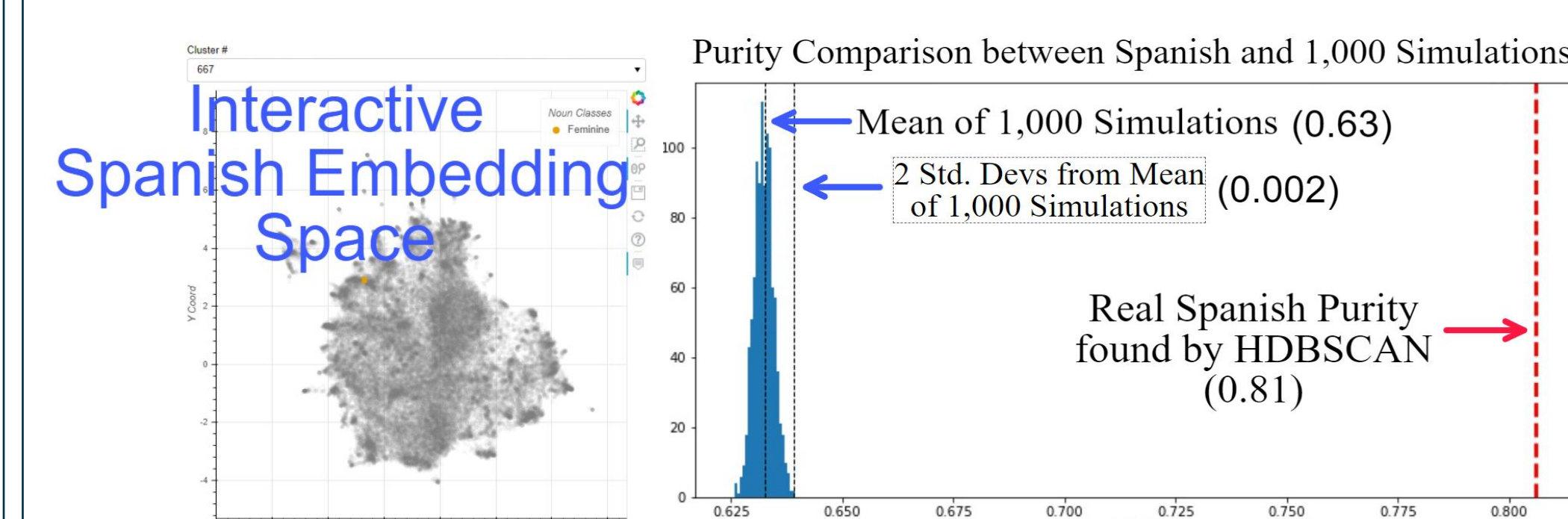


- 17/19 Languages have mid-high effect sizes, but all languages had some degree of systematicity

HDBSCAN:



- All Languages **STRONGLY** outperform 1,000 simulations of random assignment



Conclusions

- Noun Class assignment not completely pure nor completely arbitrary.
- Amount of noise detected by HDBSCAN points towards random assignment of outliers but further testing necessary
- Systematicity between Semantics and Noun class across multiple language families implies cognitive predisposition
- Future testing will include plotting multiple languages onto a common embedding space to see if different languages categorize the same concepts consistently

All Interactive Visualizations uploaded to:
<https://ardianor-ops.github.io/semanticorewebsite/>

References

- Bergen, J. J. (1980). The semantics of gender contrasts in spanish. *Hispania*, 63(1), 48-57. doi:10.2307/340811
- Bojanowski, Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics*, 5, 135-146. https://doi.org/10.1162/tacl_a_00051
- Grave, Bojanowski, P., Gupta, P., Joulin, A., & Mikolov, T. (2018). Learning Word Vectors for 157 Languages.
- Harris, Z. S. (1954). Distributional Structure. *WORD*, 10(2-3), 146-162. doi:10.1080/00437956.1954.11659520
- Kidd, E., & Garcia, R. (2021, September 6). How diverse is child language acquisition?. <https://doi.org/10.31234/osf.io/jpeyq>

Acknowledgements

Thank you to Yujing Huang, Joshua Hartshorne, Wesley Ricketts, Seun Eisape, Aidas Aglinsky, Damian Blasi, and Katie McCauliffe for their endless support and assistance with this project.