

He, she, they, they: A first discriminative analysis of third-person pronoun semantics

Dominic Schmitz

Heinrich-Heine-Universität Düsseldorf

In recent years, the use of appropriate third person pronouns has gained increased attention. In English linguistics, this attention overwhelmingly manifests in form of sociolinguistic and syntactic research, which, e.g., investigates the different types of singular *they* (e.g. Conrod, 2022; Han & Moulton, 2022; Konnelly et al., 2020). What is missing, however, is a semantic account of not only third person pronouns, but pronouns in general. The present paper offers a first account of pronoun semantics by example of *he*, *she*, and plural and singular *they*.

While *they* is a third person plural pronoun in its prototypical usage, it has been attested as third person singular pronoun at least since the 15th century (Conrod, 2020). In contemporary English, one can differentiate at least four types of singular *they*: generic indefinite and definite, and specific definite ungendered and gendered. As the latter two types are rather infrequent, the present paper focuses on generic *they*.

To gain insight into the semantics of *he*, *she*, plural and generic *they*, first a corpus of English was created. The corpus consisted of 1000 sentences sampled from COCA (Davies, 2008-), with 100 sentences per pronoun and 700 further sentences. Second, using naive discriminative learning (Baayen et al., 2011), semantic vectors were computed based on the corpus. Semantic vectors capture the semantics of words (Boleda, 2020), but, to the author's knowledge, have not been used for the semantic analysis of pronouns before. As each word form's semantics are represented by a single vector, this leaves us with three vectors for the four pronouns under investigation. As an exhaustive analysis is barely feasible with three vectors, in a third step so-called instance vectors (Lapesa et al., 2018) were computed. To compute instance vectors, for each pronoun attestation the ten preceding and following words' semantic vectors were averaged. Based on the instance vectors, the other words' vectors and the orthographical forms of all words contained in the corpus, a linear discriminative learning network (LDL; Baayen et al., 2019) was implemented in a fourth step. LDL networks map forms onto meaning and vice versa to allow insight into the interrelations of forms and meanings in the mental lexicon. Based on these interrelations, in the fifth and final step, two semantic measures were computed: comprehension quality and semantic activation diversity. The former mirrors how well a word's semantics are comprehended by the network, the latter represents the degree of coactivation by a given word.

Comparing the measures of all four suffixes, it is found that generic *they* is comprehend significantly better than plural *they*. At the same time, generic *they* coactivates entries in the lexicon to the same degree as plural *they* does. When compared to *he* and *she* via cosine similarities, generic *they* shows high similarities. However, generic *they* is more similar to plural *they* than *he* and *she* are. Overall, generic *they* appears to be a singular pronoun with remnants of plurality.

The present findings demonstrate two equally important points. First, distributional semantics can be used for the analysis of pronoun semantics. Second, generic *they* seems to be a fit third person singular pronoun. However, further research is required to investigate in how far its nature is truly is generic.

References

- Baayen, R. H., Chuang, Y.-Y., Shafaei-Bajestan, E., & Blevins, J. P. (2019). The discriminative lexicon: A unified computational model for the lexicon and lexical processing in comprehension and production grounded not in (de)composition but in linear discriminative learning. *Complexity*, 2019, 4895891. <https://doi.org/10.1155/2019/4895891>
- Baayen, R. H., Milin, P., Đurđević, D. F., Hendrix, P., & Marelli, M. (2011). An amorphous model for morphological processing in visual comprehension based on naive discriminative learning. *Psychological Review*, 118(3), 438–481. <https://doi.org/10.1037/a0023851>
- Boleda, G. (2020). Distributional semantics and linguistic theory. *Annual Review of Linguistics*, 6(1), 213–234. <https://doi.org/10.1146/annurev-linguistics-011619-030303>
- Conrod, K. (2020). Pronouns and gender in language. In *The Oxford Handbook of Language and Sexuality*. Oxford University Press. <https://doi.org/10.1093/oxfordhb/9780190212926.013.63>
- Conrod, K. (2022). Abolishing gender on D. *Canadian Journal of Linguistics/Revue Canadienne de Linguistique*, 67(3), 216–241. <https://doi.org/10.1017/cnj.2022.27>
- Davies, M. (2008-). *The Corpus of Contemporary American English (COCA)*. <https://www.english-corpora.org/coca/>
- Han, C. H., & Moulton, K. (2022). Processing bound-variable singular they. *Canadian Journal of Linguistics/Revue Canadienne de Linguistique*, 67(3), 267–301. <https://doi.org/10.1017/CNJ.2022.30>
- Konnelly, L., Cowper, E., Konnelly, L., & Cowper, E. (2020). Gender diversity and morphosyntax: An account of singular they. *Glossa: A Journal of General Linguistics*, 5(1). <https://doi.org/10.5334/GJGL.1000>
- Lapesa, G., Kawaletz, L., Plag, I., Andreou, M., Kisselew, M., & Padó, S. (2018). Disambiguation of newly derived nominalizations in context: A Distributional Semantics approach. *Word Structure*, 11(3), 277–312. <https://doi.org/10.3366/word.2018.0131>