

# On potential statistical universals of grammar in discourse: Evidence from Multi-CAST

**Geoffrey Haig,<sup>1</sup> Nils N. Schiborr,<sup>1</sup> Stefan Schnell<sup>1,2</sup>**

<sup>1</sup> *University of Bamberg*, <sup>2</sup> *Centre of Excellence for the Dynamics of Language*

DGfS2020, Hamburg

*Workshop Corpus-based typology:*

*Spoken language from a cross-linguistic perspective*

4–6<sup>th</sup> March 2020

# Overview

1. Corpus-based typology with Multi-CAST
2. Cross-linguistic uniformity in the distribution of full ('lexical') expressions
3. Light Human Subjects
4. The subject/object person asymmetry
5. Conclusions

(1) Corpus-based typology with Multi-CAST

# Traditional research in discourse & grammar

- roots in the functionalist tradition: Chafe, Givón, Prince, Du Bois, among many others
- ‘grammar’ shaped and constrained by demands of successful communication, rather than an autonomous module
- explicitly cross-linguistic, empirical perspective
- remains hugely influential, e.g. in Cognitive Grammar, Grammaticalization

# Corpus-based typology

- couples the functionalist tradition with digital corpora, and methodologies from corpus linguistics and variationist sociolinguistics (Schnell & Barth 2018)
- complements grammar-based, or “data-reduction” typology (Wälchli 2009)
- bottom-up, data-driven, probabilistic rather than categorial generalizations
- attends to variation, attends to context
- In Multi-CAST: focus on spoken language, monologic, indigenous content, sample breadth rather than corpus breadth

# Multi-CAST

*Multilingual Corpus of Annotated Spoken Texts*

Multi-CAST

[Annotations](#)

[The corpora](#) ↓

[Research](#)

[Contribute](#)

[People](#)

[More](#) ↓

[Contact](#)

**Multi-CAST**, the *Multilingual Corpus of Annotated Spoken Texts*, is a collection of annotated texts from a typologically diverse section of languages.

- ◆ multiple levels of parallel [annotations](#), time-aligned with audio recordings,
- ◆ including comparative morphosyntactic annotations for [cross-corpus typological research](#)
- ◆ chiefly monologic, natural narrative texts from [twelve languages](#), encompassing roughly 21 500 clause units
- ◆ available in multiple file formats, including as EAF files for the [linguistic annotation software ELAN](#), as XML and TSV files, and via the [multicastR](#) package for R
- ◆ freely accessible under a [CC-BY 4.0 licence](#)

## Getting started with Multi-CAST

<b>collection overview (!)</b>	<a href="#">PDF</a>	355 KB	v2.1	20/01/12	<a href="#">archive</a>
<i>research context</i>	<a href="#">PDF</a>	337 KB	v1.1	18/05/25	<a href="#">archive</a>
<i>full collection</i>	<a href="#">XML</a>	39 MB	2001	20/01/12	<a href="#">archive</a>
	<a href="#">TSV</a>	6.9 MB	2001	20/01/12	<a href="#">archive</a>
<i>full metadata</i>	<a href="#">TSV</a>	4 KB	2001	20/01/12	<a href="#">archive</a>

## Citing Multi-CAST

Haig, Geoffrey & Schnell, Stefan (eds.), *Multi-CAST: Multilingual corpus of annotated spoken texts*. ([multicast.aspra.uni-bamberg.de/](http://multicast.aspra.uni-bamberg.de/)) (date accessed)

(2) The cross-linguistic uniformity in the use of full (“lexical”) expressions

# Lexical versus reduced (Kibrik 2011) forms of referring expressions

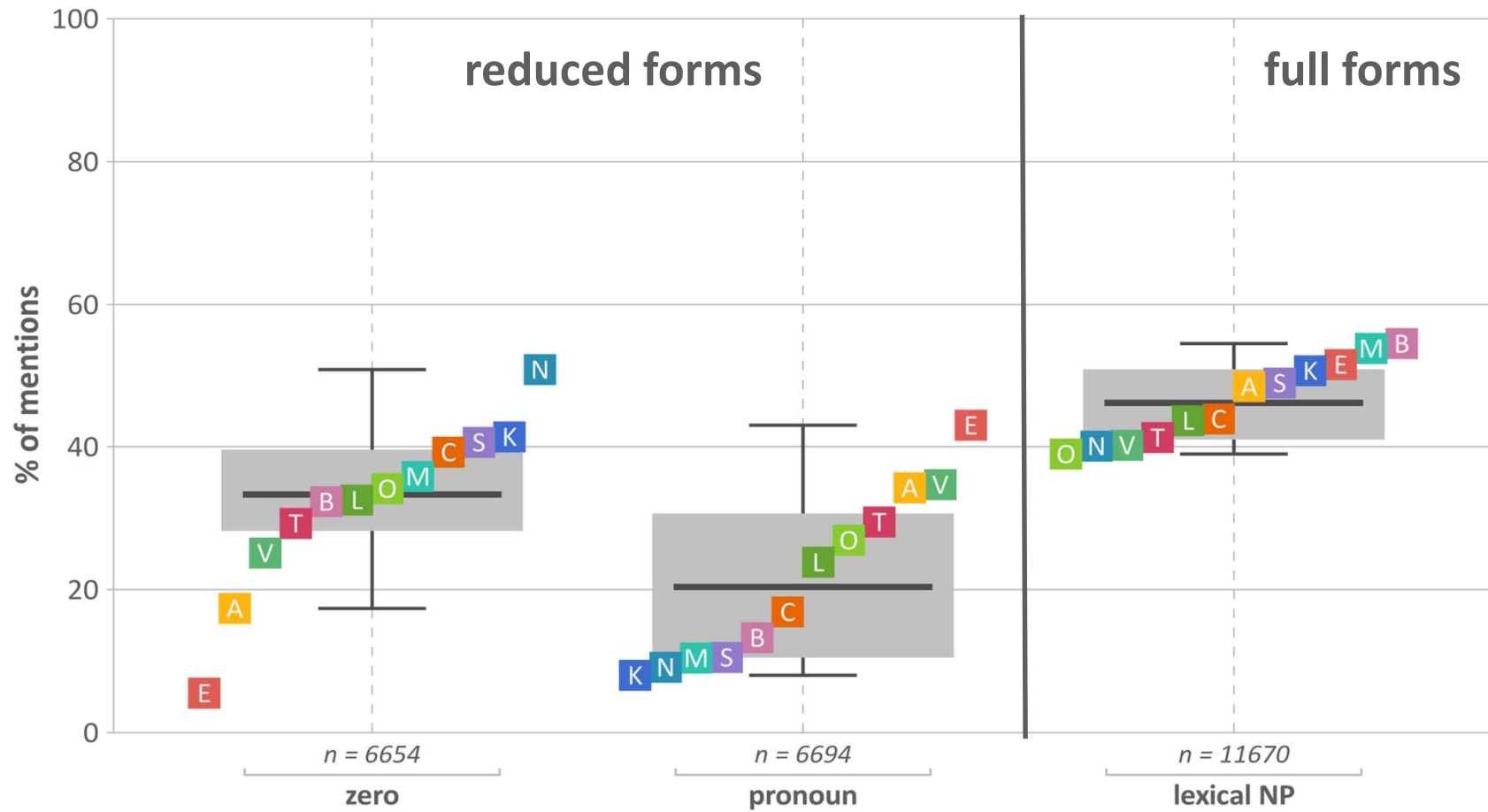
## lexical 'full' expressions / lexical NPs:

*A new syntax professor / Amanda / that woman / the supervisor ...*

## reduced / 'light' forms:

*she / her / Ø*





- |                        |                           |                        |                  |
|------------------------|---------------------------|------------------------|------------------|
| <b>A</b> Arta          | <b>M</b> Mandarin         | <b>S</b> Sanzhi Dargwa | <b>O</b> Tondano |
| <b>C</b> Cypriot Greek | <b>N</b> Nafsan           | <b>B</b> Tabasaran     | <b>L</b> Tulil   |
| <b>E</b> English       | <b>K</b> Northern Kurdish | <b>T</b> Teop          | <b>V</b> Vera'a  |

# Uniformity of lexical expressions: interim summary

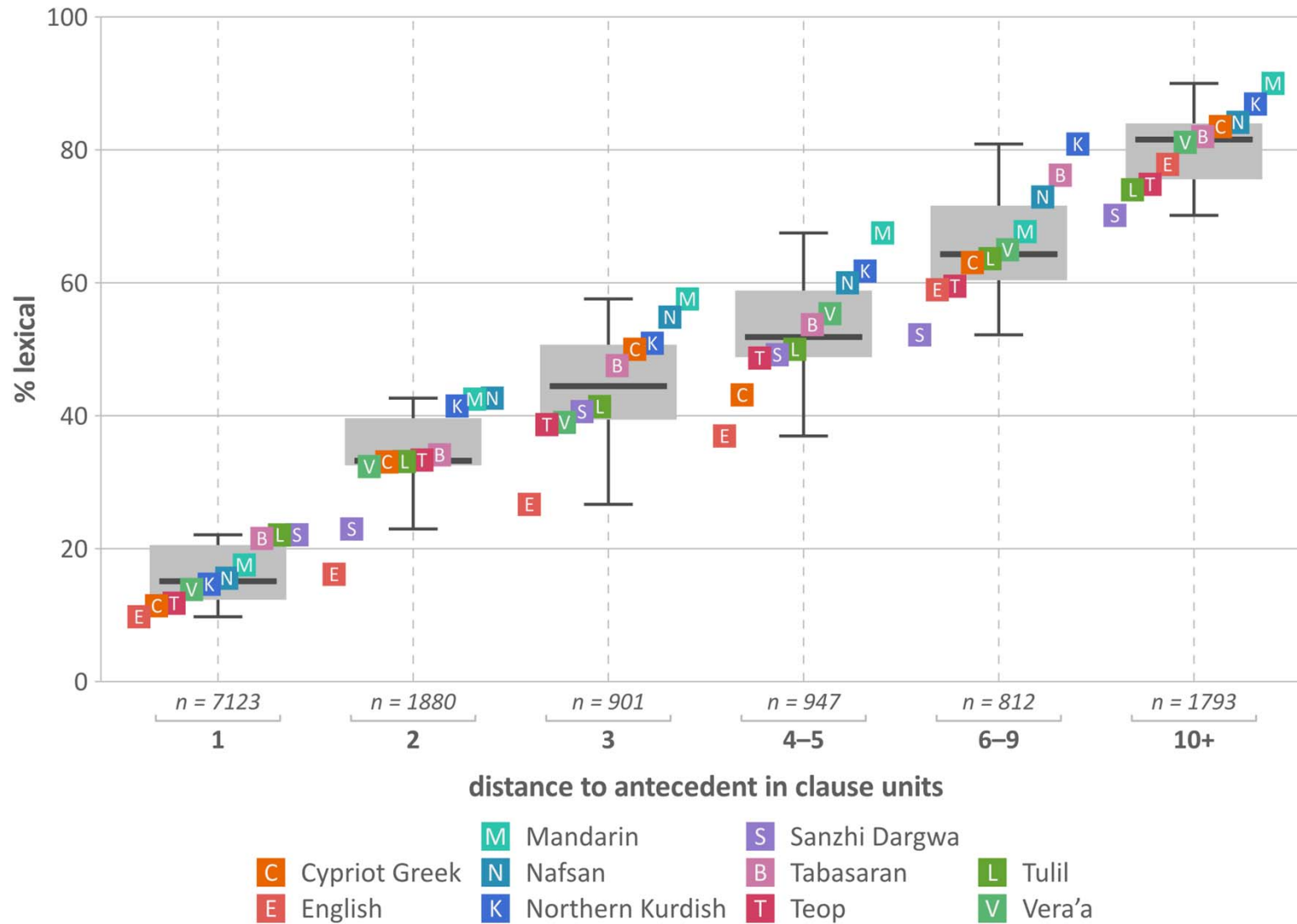
- discourse is carried by a **relatively uniform bedrock of lexical expressions** (40–60%), regardless of language
- the locus of cross-linguistic variability is **the respective contributions of zero and pronouns among the reduced expressions** (Schiborr, in prep., Schnell & Barth 2018)

# Explanations for uniformity of lexical expression

- lexical forms are used with similar rates across languages because their use is largely determined by **the same factors**
- most powerful factor: **anaphoric distance**
- the pronoun vs. zero choice is tempered by language-specific inherited historical accidents of morphosyntax, not treated in today's presentation, e.g.
  - presence of agreement morphology
  - informativity of pronouns (gender, number etc.)
  - differing effects with subjects and objects (Schnell & Barth 2018, resub.; Schwenter 2006, 2014)

(**exception**: same-subject clause sequence contexts favour zero subjects across all languages (Torres Cacoulos & Travis 2019, Vollmer 2019, Schiborr, in prep.)

# Anaphoric distance and lexical expression (Schiborr, in prep.)



# Uniformity of lexical expressions: theoretical implications

- suggests a re-evaluation of the view that informativeness of discourse is language specific (i.e. that some languages are apparently ‘less explicit’, rely more on ‘pragmatic inference’, typologies of ‘pragmatic vs. syntactic’, ‘hot vs. cold’ languages; Stoll & Bickel 2009, Huang 2000)
- e.g. Mandarin: actually among the highest levels of lexical expressions in our sample (Vollmer 2019)
- little evidence for an across-the-board impact of ‘accessibility’ dictating zero vs. pronoun, and lexical vs. reduced (Schiborr, in prep.)

## Uniform rates of lexical NPs: candidate universal

- in spontaneous unplanned discourse, between 40–60% of referring expressions are lexical NPs, regardless of language

(3) Light Human Subjects:  
The skewed distribution of new referents in syntax

# Light human subjects

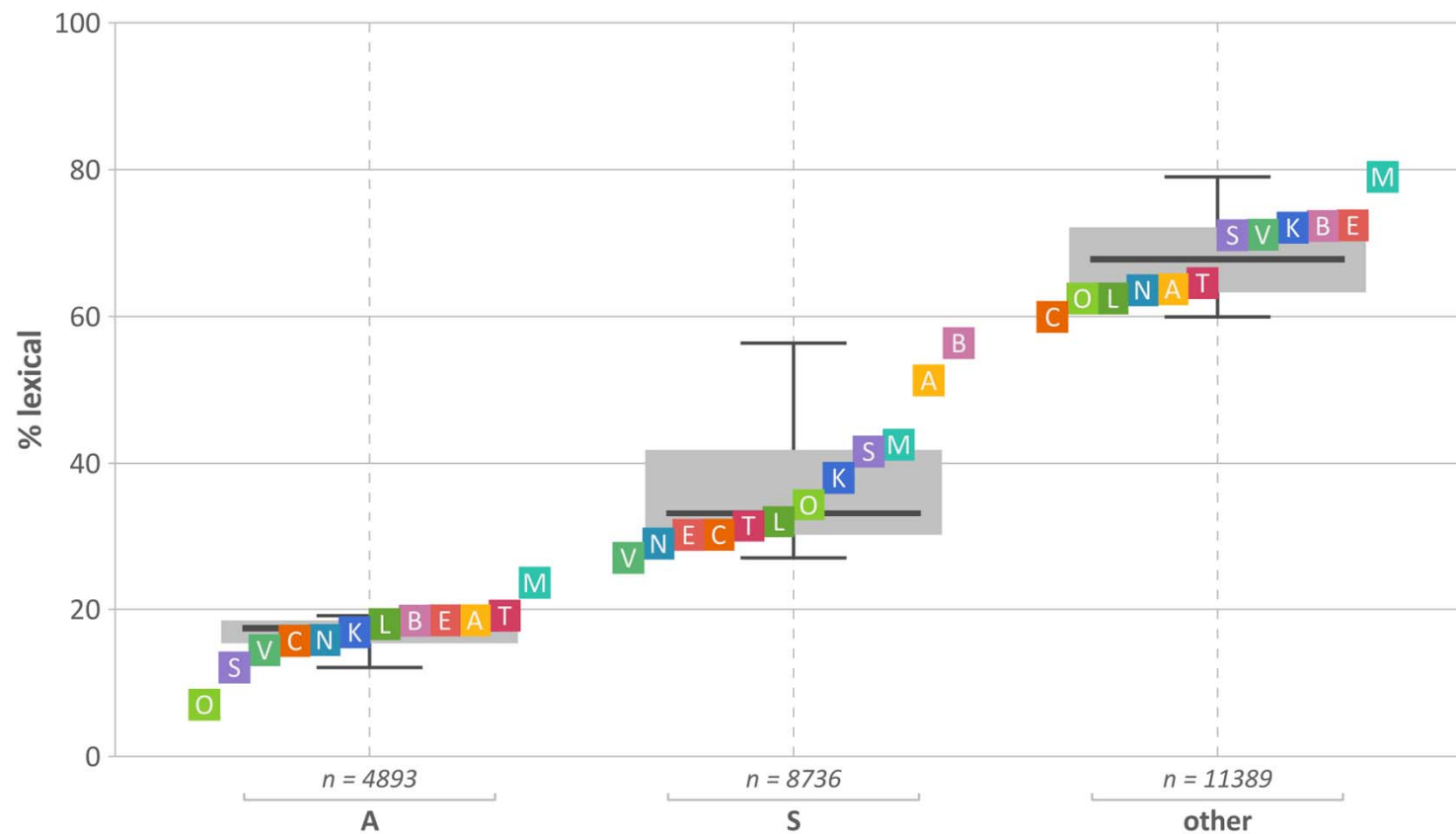
- original observation by Du Bois 1987: **Avoid new/lexical A**
- **lexical** referential forms (with new referents) vs. **reduced** referential forms (pronouns, zero) are not evenly distributed across syntactic functions (Du Bois 1987, 2003, 2017)
- **transitive subjects (A)** apparently particularly favour reduced as opposed to lexical forms
- transitive (A) and intransitive subjects (S) apparently differ in this respect, with S clustering with P (objects)
- Du Bois' explanation for Avoid lexical A is related to information management in discourse, e.g. avoidance of more than one new referent per clause



# Light human subjects

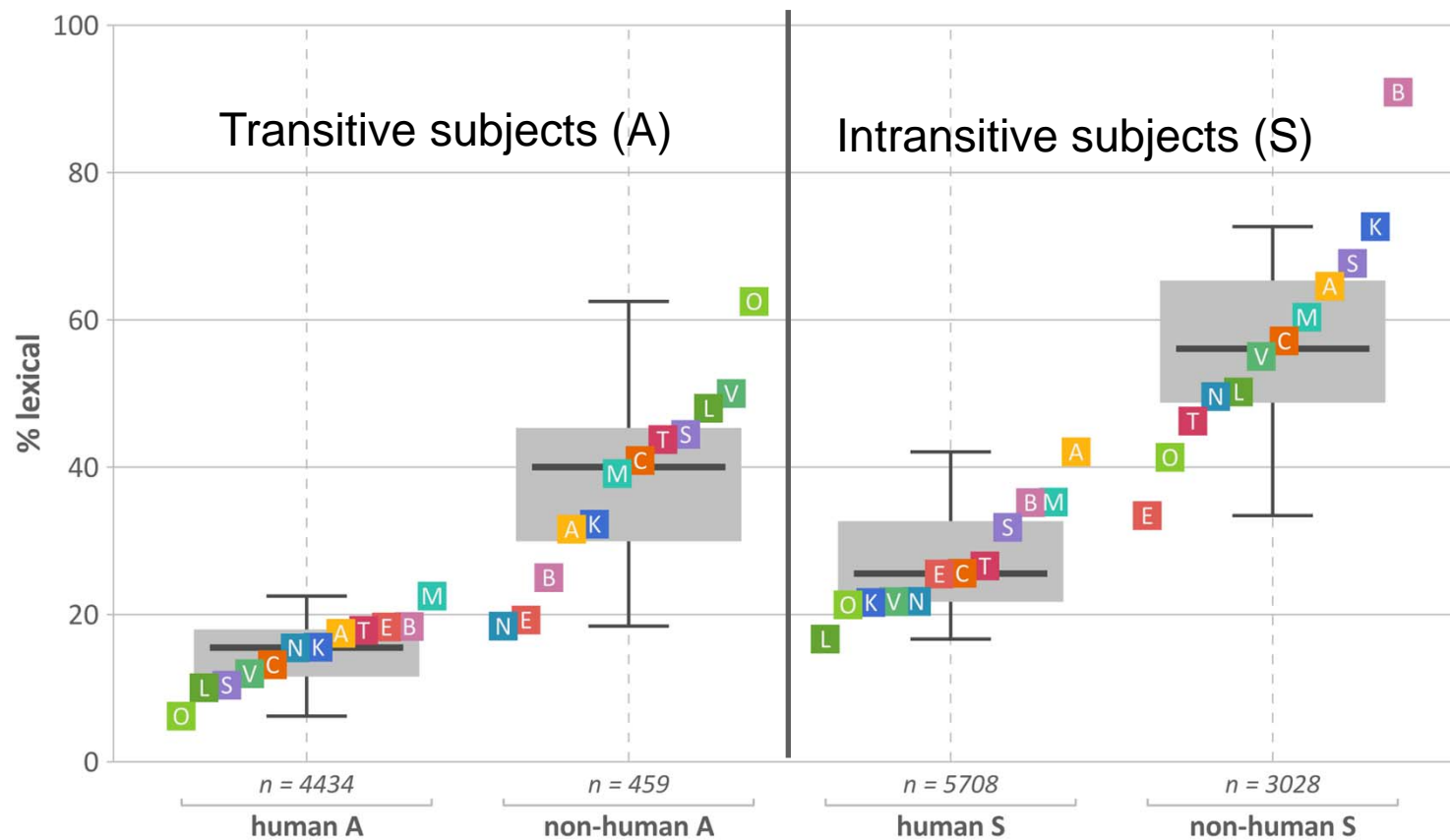
- these claims have been questioned on empirical and conceptual grounds (e.g. Payne 1987, Kärkkäinen 1996, Haspelmath 2003, Everett 2009, Haig & Schnell 2016):
- no clustering of S and P; S and A closer than predicted
- role of information management overestimated; animacy accounts for most of the variation
- data from Multi-CAST ...

# Distribution of lexical arguments: A, S, and other



- |   |  |   |   |
|---|--|---|---|
| <span style="background-color: yellow; border: 1px solid black; padding: 2px;">A</span> Arta          | <span style="background-color: cyan; border: 1px solid black; padding: 2px;">M</span> Mandarin         | <span style="background-color: purple; border: 1px solid black; padding: 2px;">S</span> Sanzhi Dargwa | <span style="background-color: lightgreen; border: 1px solid black; padding: 2px;">O</span> Tondano |
| <span style="background-color: orange; border: 1px solid black; padding: 2px;">C</span> Cypriot Greek | <span style="background-color: blue; border: 1px solid black; padding: 2px;">N</span> Nafsan           | <span style="background-color: pink; border: 1px solid black; padding: 2px;">B</span> Tabasaran       | <span style="background-color: lightgreen; border: 1px solid black; padding: 2px;">L</span> Tulil   |
| <span style="background-color: red; border: 1px solid black; padding: 2px;">E</span> English          | <span style="background-color: blue; border: 1px solid black; padding: 2px;">K</span> Northern Kurdish | <span style="background-color: red; border: 1px solid black; padding: 2px;">T</span> Teop             | <span style="background-color: green; border: 1px solid black; padding: 2px;">V</span> Vera'a       |

# Human vs. non-human A and S



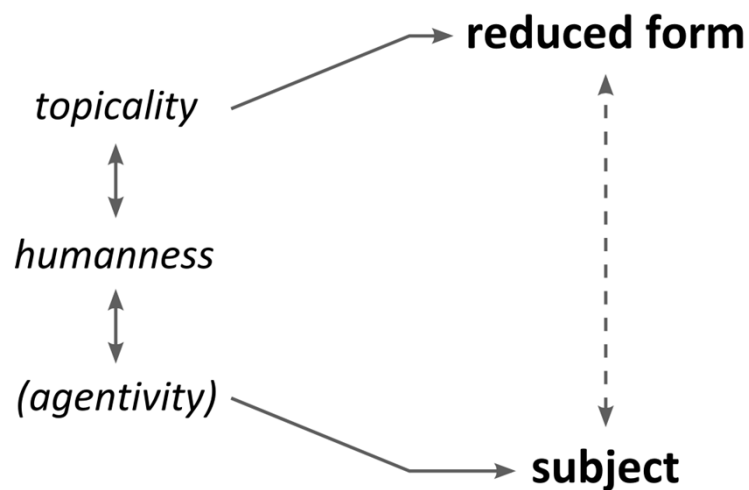
- |   |               |   |                  |   |               |   |         |
|---|---------------|---|------------------|---|---------------|---|---------|
| A | Arta          | M | Mandarin         | S | Sanzhi Dargwa | O | Tondano |
| C | Cypriot Greek | N | Nafsan           | B | Tabasaran     | L | Tulil   |
| E | English       | K | Northern Kurdish | T | Teop          | V | Vera'a  |

# Light human subjects: interim summary

- the impact of transitivity (A vs. S) has been overrated
- the relevant generalization couples 'humanness' with 'subject' (S or A)
- not a question of 'constraints on information management in discourse', but a more general strategy reflecting cognitive prominence of human, topical entities

# Explanations

- rather than a direct link of syntactic role and transitivity ('A') with information status, a more general concern with **human actors** drives the distribution:



- the significant factor is the **pragmatic and semantic prominence of human referents**

## Light human subjects: candidate universal

- in spontaneous unplanned discourse, human subjects are generally (>75%) reduced

(4) The person asymmetry across subjects and objects

# Subject/object asymmetry: main finding

- asymmetry between subjects and objects wrt. to various parameters regularly noted (e.g. Haig 2018; Schnell & Barth, resubm.; Dalrymple & Nikolaeva 2011)
- objects exhibit more complex patterns of pronominalization, with language-specific factor weightings (Schnell & Barth 2018)
- the single most robust difference appears to be robust regularities in the distribution of person values in transitive clauses



# Subject/object asymmetry: main finding

- the person value of transitive subjects is determined by **content and genre**:

**conversational  
narratives**

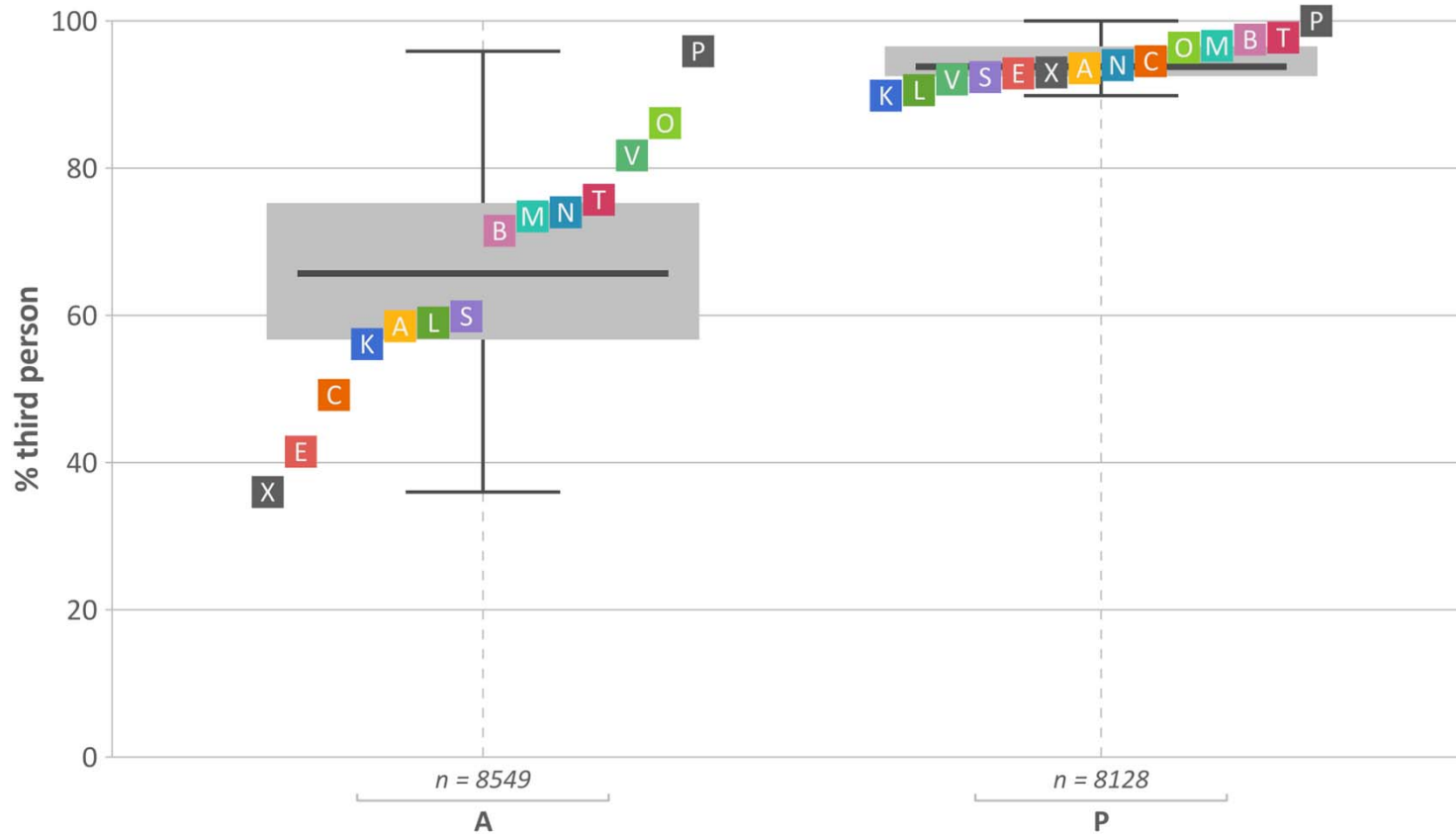
→ high levels of 1<sup>st</sup>/2<sup>nd</sup> person, low levels of 3<sup>rd</sup> person  
→ low 1<sup>st</sup>/2<sup>nd</sup> person, high 3<sup>rd</sup> person

- the person value of objects is **impervious to content and genre**:

**all genres:**

→ overwhelmingly 3<sup>rd</sup> person

# Person values, subjects vs. objects



- |                           |                   |                        |                  |
|---------------------------|-------------------|------------------------|------------------|
| <b>A</b> Arta             | <b>P</b> Persian  | <b>X</b> English (SB)  | <b>O</b> Tondano |
| <b>C</b> Cypriot Greek    | <b>M</b> Mandarin | <b>S</b> Sanzhi Dargwa | <b>L</b> Tulil   |
| <b>E</b> English (MC)     | <b>N</b> Nafsan   | <b>B</b> Tabasaran     | <b>V</b> Vera'a  |
| <b>K</b> Northern Kurdish | <b>T</b> Teop     |                        |                  |

## The person asymmetry: candidate universal

- in spontaneous unplanned discourse, objects are overwhelmingly (> 90%) third person, regardless of content and genre
- the person values of transitive subjects, on the other hand, are dependent on content and genre

## (5) Summary: candidate universals

Spontaneous unplanned discourse appears to comply with the following quantitative universals:

- **Light Human Subjects**  
the majority (> 75%) of human subjects are reduced in form (pronominal, zero)
- **Uniform rates of lexical expression**  
between 40–60% of referring expressions are lexical NPs;  
the respective rates of pronoun and zero, on the other hand, are subject to cross-linguistic variability
- **The subject/object asymmetry in person values**  
at least 90% of all objects are third person;  
there is no comparable constant rate for subject person values

# Thanks!



University of Bamberg



ARC CENTRE OF EXCELLENCE FOR  
THE DYNAMICS OF LANGUAGE



DOBES



Deutsche  
Forschungsgemeinschaft



Australian Government  
Australian Research Council



VolkswagenStiftung

# References

- Dalrymple**, Mary & **Nikolaeva**, Irina. 2011. *Objects and information structure*. Cambridge: Cambridge University Press.
- Du Bois**, John. 1987. The discourse basis of ergativity. *Language* 63(4), 805–855.
- Du Bois**, John. 2003. Discourse and grammar. In Tomasello, Michael (ed.), *The new psychology of language: Cognitive and functional approaches to language structure*, 47–88. Mahwah, NJ: Erlbaum.
- Du Bois**, John. 2017. Ergativity in discourse and grammar. In Coon, Jessica & Massam, Diane & Travis, Lisa D. (eds.), *The Oxford handbook of ergativity*, 23–57. Oxford: Oxford University Press.
- Everett**, Caleb. 2009. A reconsideration of the motivations for preferred argument structure. *Studies in Language* 33(1), 1–24.
- Haig**, Geoffrey. 2018. The grammaticalization of object pronouns: Why differential object indexing is an attractor state. *Linguistics* 56(4), 781–818. (DOI: 10.1515/ling-2018-0011)
- Haig**, Geoffrey & **Schnell**, Stefan. 2015. *Multi-CAST: The Multilingual Corpus of Annotated Spoken Texts*. (multicast.aspra.uni-bamberg.de)
- Haig**, Geoffrey & **Schnell**, Stefan. 2016. The discourse basis of ergativity revisited. *Language* 92(3), 591–618. (DOI: 10.1353/lan.2016.0049)
- Haspelmath**, Martin. 2003. *Ditransitive constructions in the world's languages*. Handout, March 2003, University of California, Berkeley.
- Huang**, Yan. 2000. *Anaphora: A cross-linguistic study*. Oxford: Oxford University Press.
- Kärkkäinen**, Elise. 1996. Preferred argument structure and subject role in American English conversational discourse. *Journal of Pragmatics* 25(5), 675–701.
- Payne**, Doris L. 1987. Information structuring in Papago narrative discourse. *Language* 63(4), 783–804.
- Schiborr**, Nils N. In preparation. *Lexical anaphora: A corpus-based typological study of referential choice*. PhD dissertation, University of Bamberg.

# References

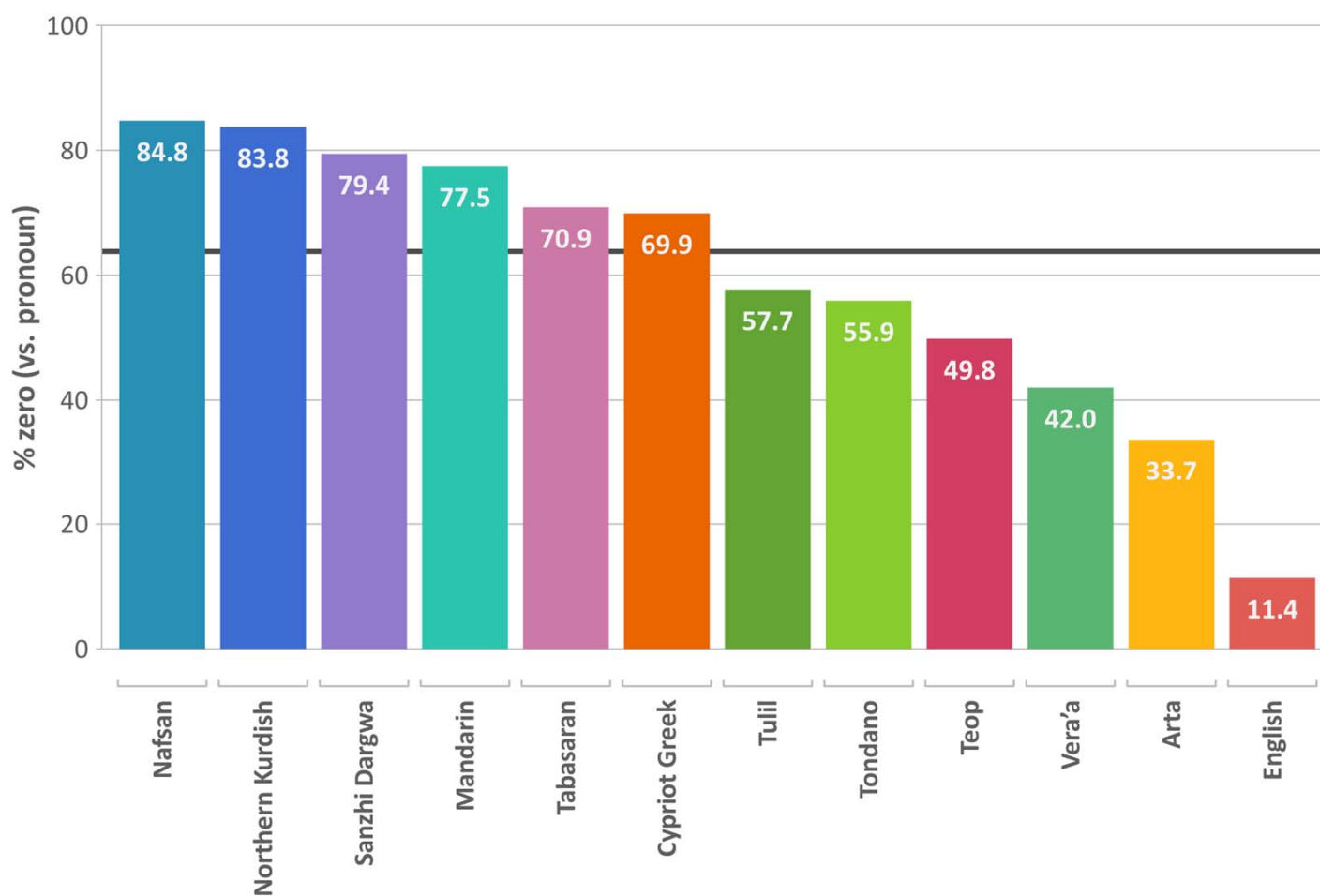
- Schnell**, Stefan & **Barth**, Danielle. 2018. Discourse motivations for pronominal and zero objects across genres in Vera'a. *Language Variation and Change* 30(1), 51–81. (DOI: 10.1017/S0954394518000054)
- Schnell**, Stefan & **Barth**, Danielle. Resubmitted. Towards subject--predicate agreement in Vera'a. Submitted to *Language Variation and Change*.
- Schwenter**, Scott. 2006. Null objects across South America. In Face, Timothy L. & Klee, Carol A. (eds), *Selected proceedings of the 8th Hispanic Linguistics Symposium*, 23–36. Somerville, MA: Cascadilla Proceedings Project.
- Schwenter**, Scott. 2014. Two kinds of object marking in Portuguese and Spanish. In Amaral, Patrícia & Carvalho, Ana M. (eds.), *Portuguese-Spanish interfaces: Diachrony, synchrony, and contact*, 237–260. Amsterdam: John Benjamins.
- Stoll**, Sabine & **Bickel**, Balthasar. 2009. How deep are differences in referential density? In Guo, Jiansheng & Lieven, Elena & Budwig, Nancy & Ervin-Tripp, Susan & Nakamura, Keiko & Özçaliskan, Seyda (eds.), *Crosslinguistic approaches to the psychology of language*, 543–555. London: Psychology Press.
- Torres Cacoulios**, Rena & **Travis**, Catherine E. 2019. Variationist typology: Shared probabilistic constraints across (non-)null subject languages. *Linguistics* 57(3), 653–692.
- Vollmer**, Maria C. 2019. *How radical is pro-drop in Mandarin? A quantitative corpus study on referential choice in Mandarin Chinese*. MA thesis, University of Bamberg.

From here on: spare slides



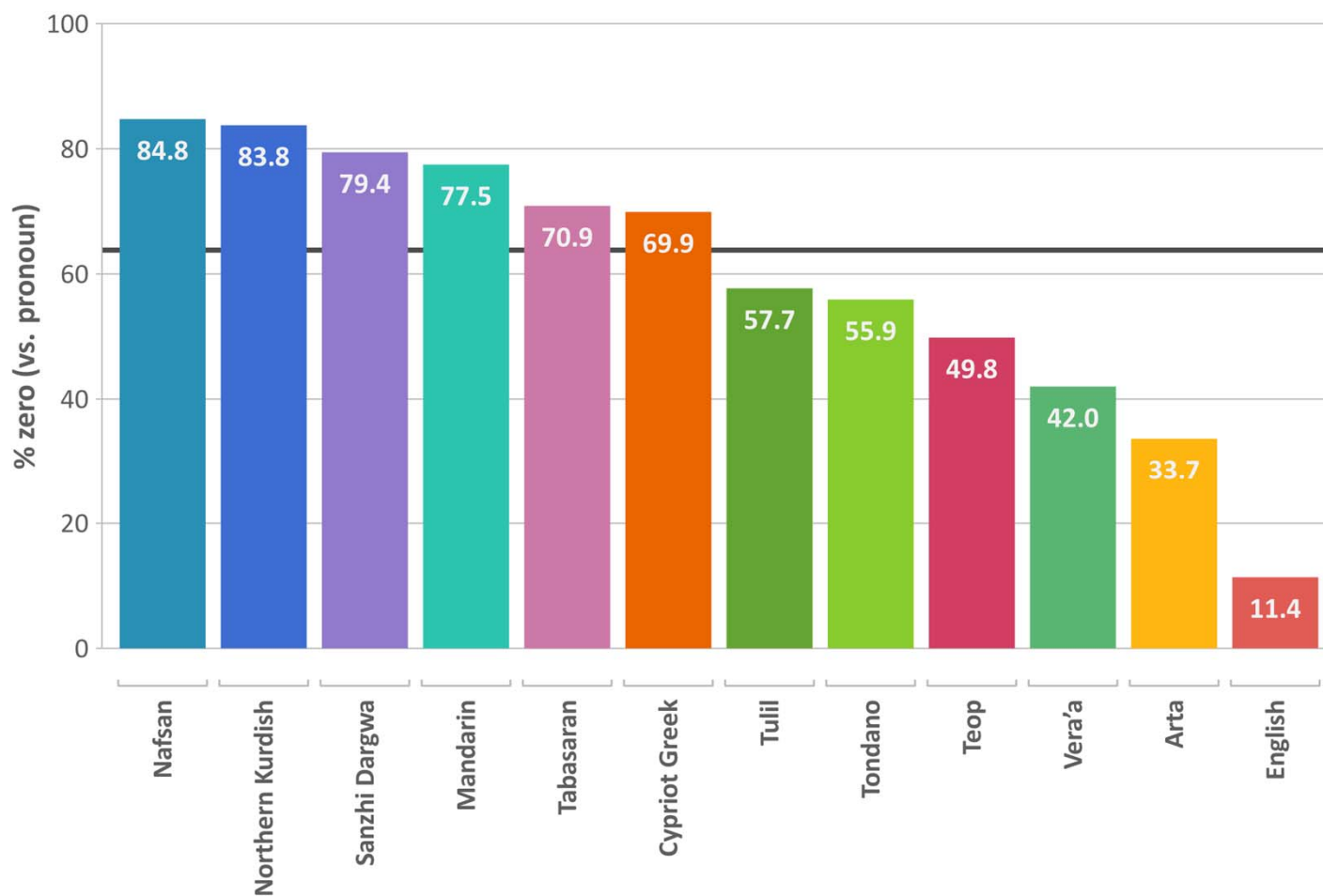
# Typological variation in pronoun use

- use of non-lexical forms (pronouns and zero) cross-linguistically variable
- expected under traditional conceptions of pro-drop / referential density



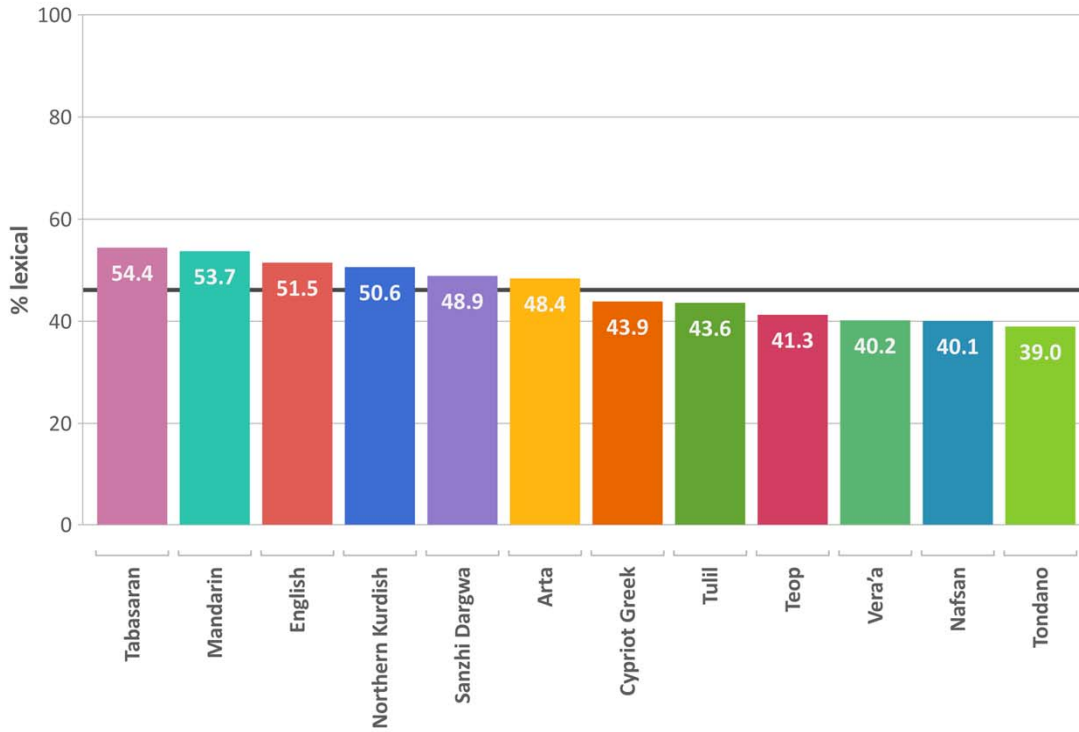
# Possible factors and trends

- differences between subjects and objects (Haig 2018; Schnell & Barth 2018; Dalrymple & Nikolaeva 2011); objects more sensitive to prominence-related factors
- same-subject contexts might make a difference  
agreement might make a difference

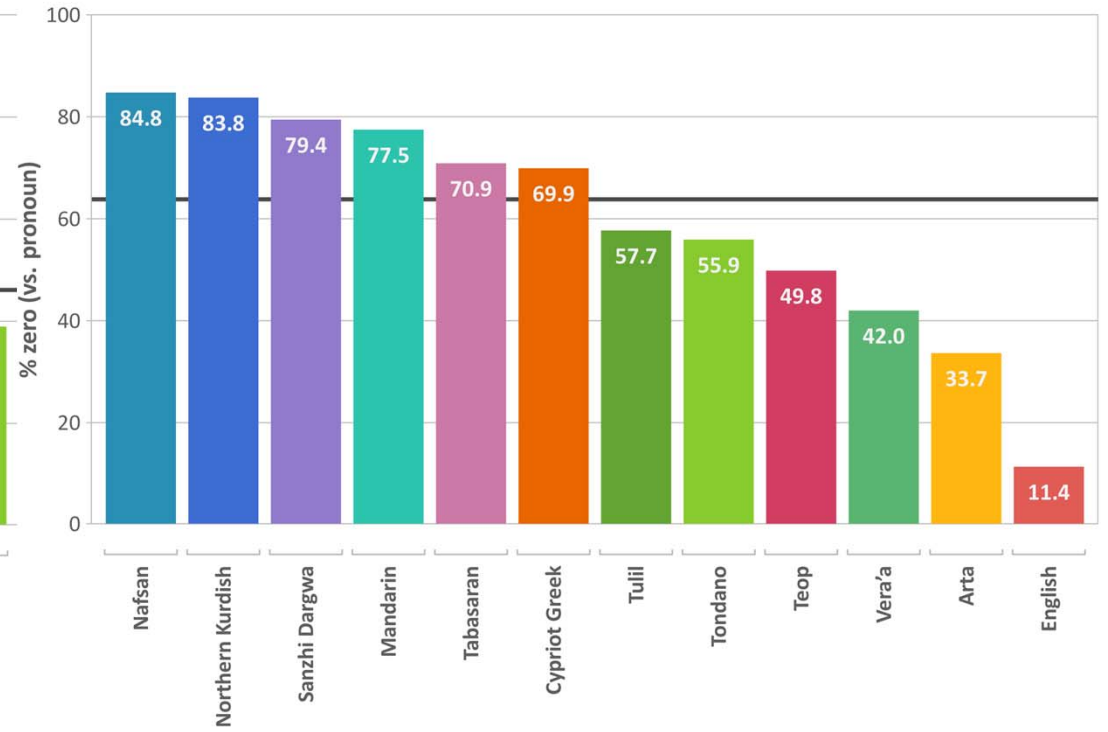


# Uniformity of lexicality vs. the variance of pronoun/zero

## % lexical forms



## % zero among reduced forms



# Complementarity in subjects and conditions on objects

- differences between subjects and objects
- **subjects:**
  - complementarity between clause-level and agreement paradigm exponence in subjects (see Schnell & Barth (resub) on subjects in Vera'a; Fuß 2005; Rosenqvist 2008);
  - connection to ongoing processes of grammaticalization
- **objects:**
  - might show typological categorization between
    - (A) either rigorously pronominal (English, Tuli)
    - (B) or conditioned (humanness, discourse topicality, person),  
see Schnell & Barth (2018), Schwenter (2006, 2014)

# Pronouns vs zeroes: interim summary

- tremendous diversity in use of pronouns vs zeroes across languages
- possible universal trends: complementarity in subjects and conditions on objects

## Universal IIIa (subjects):

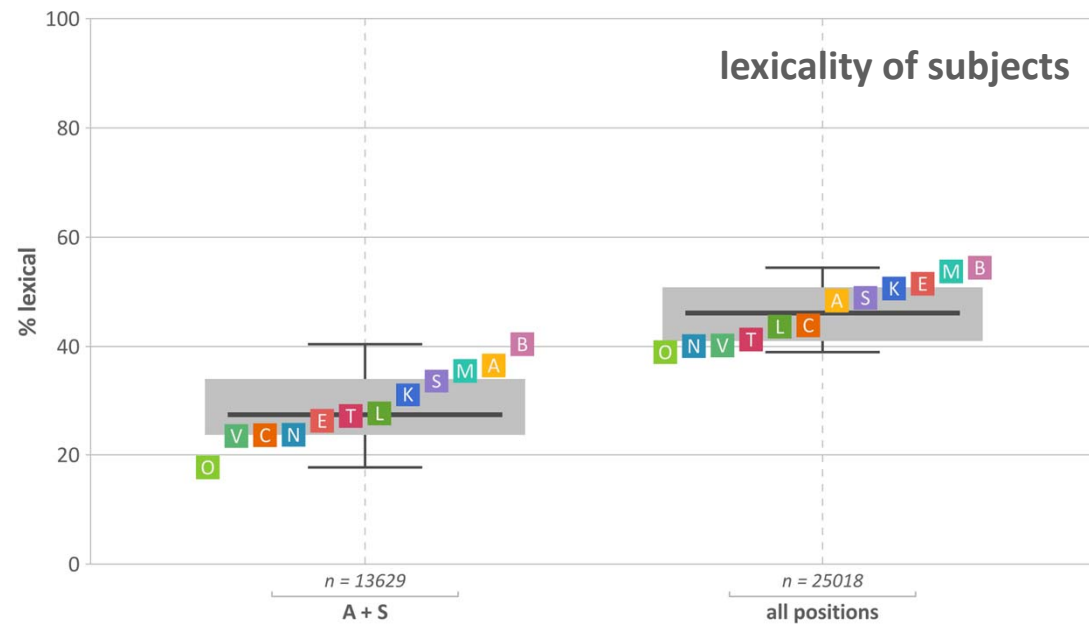
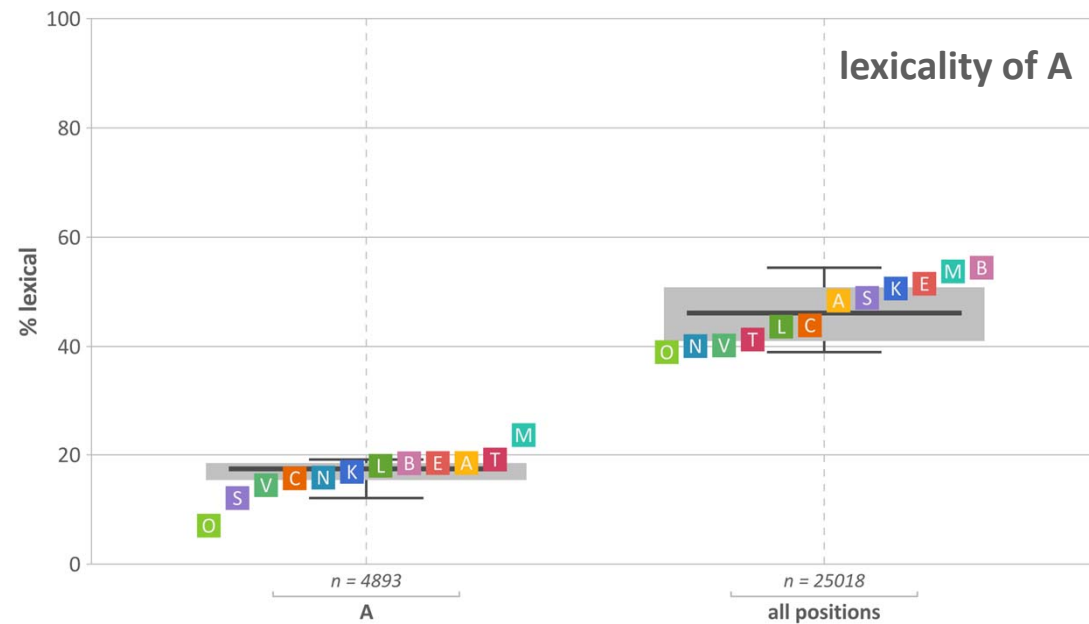
- **Observational**: subjects tend towards complementarity in pers/num marking.
- **Explanatory**: (diachronic, preferences in change) grammaticalization of agreement  
(functional-adaptive) economy of marking

## Universal IIIb (objects):

- **Observational**: if pronouns not categorical, then mark prominence/salience (or the unexpected, Aissen 2003)
- **Explanatory**: if available, use some form to marking (might be in agreement, Haig 2018)

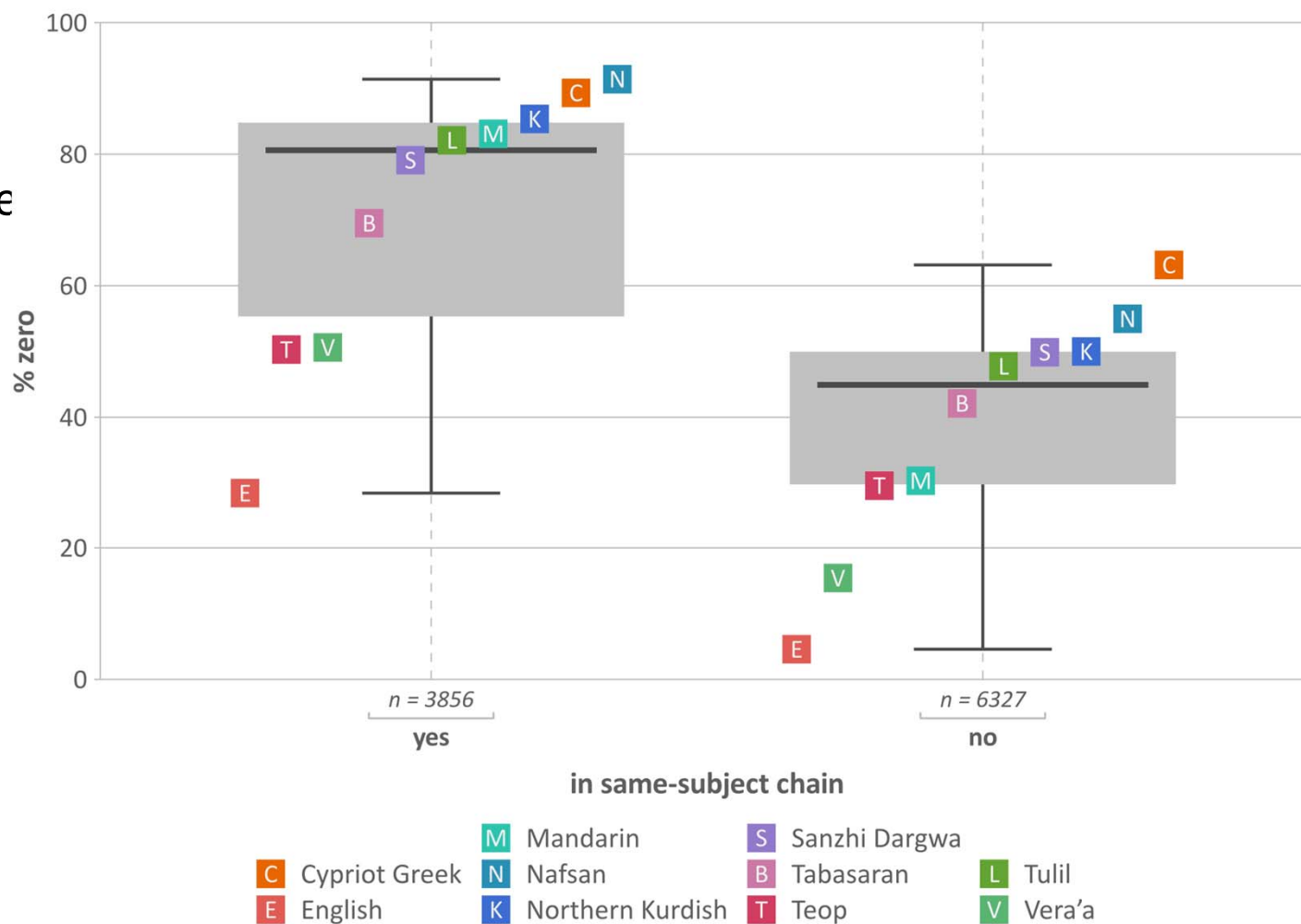
# Light subjects: a robust finding

- Empirically robust: low lexicality / newness of A
- in fact: **light subjects** (Chafe 1987, 1994)



# What drives the choice between pronoun and zero?

- for Accessibility Theory
- (Ariel 1990; Givon 1983),
- motivations are essentially the same as those pertaining to choice of lexical NPs
- (antecedent relations, prominence, etc)
- others stress the functional similarity of both types of expression
- repeated observation: prevalence of zero in same-subject clausal chains



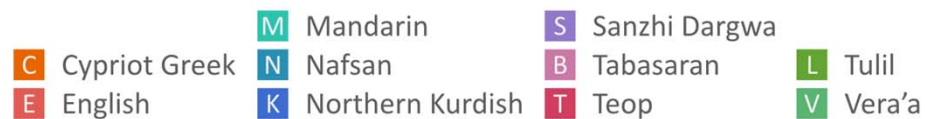
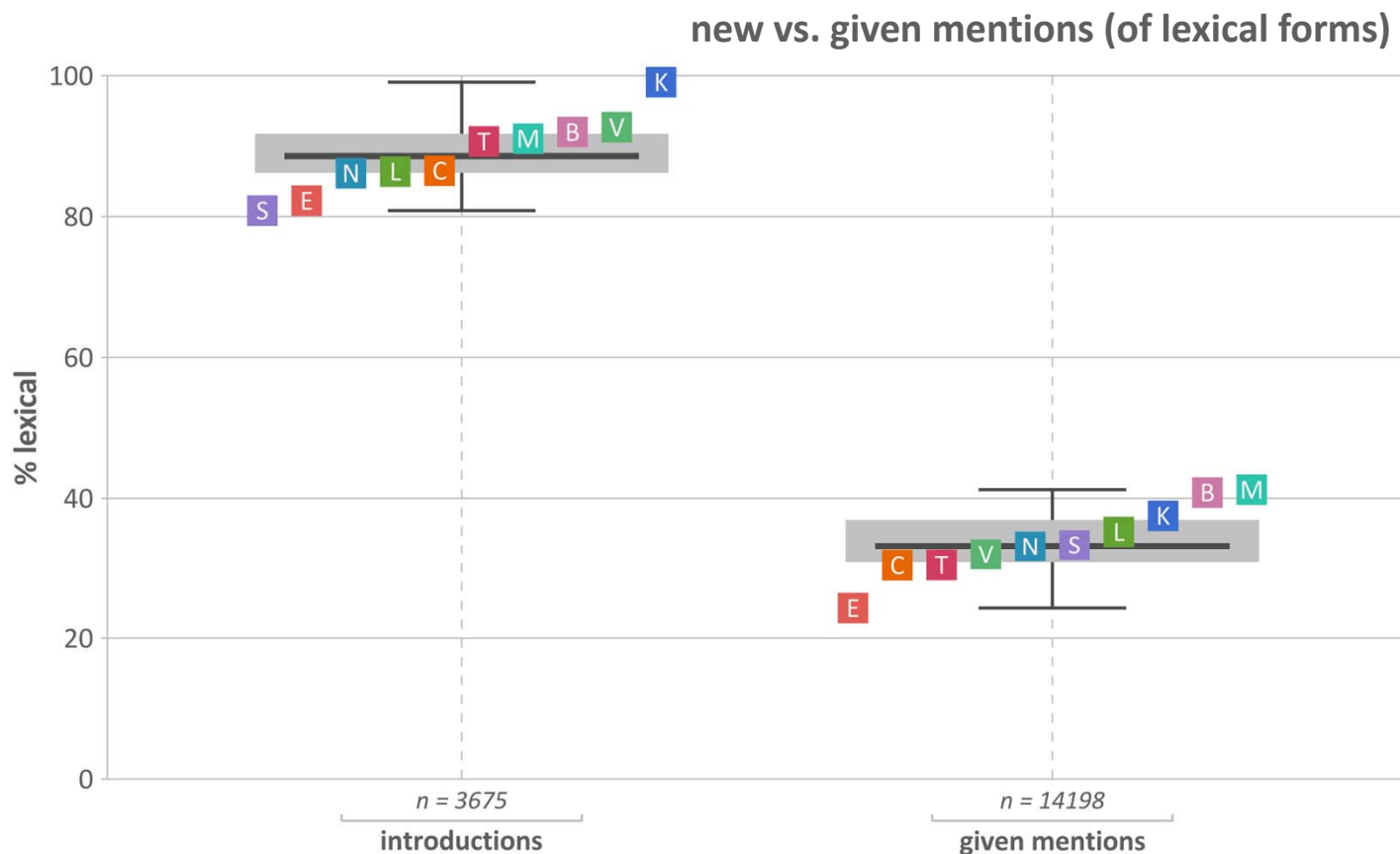
# What drives the cross-linguistic variation?

- Various accounts in terms of holistic properties of languages:
  - pro-drop parameter in various versions (**Roberts & xxx 2011**)
  - related to agreement
  - radical pro-drop
  - referential density (Bickel 2003)
- Complex interaction of holistic properties and specific patterns of variation
- (Torres Cacoulos & Travis 2019)



# Newness correlates with lexical form

- nearly all new referents receive lexical form of expression
- contrasts with occasional observations of non-lexical introductions (Stoll & Bickel 2009)
- most lexical NPs have given referents



# Light human subjects

- Original observation: **lexical** referential forms (with new referents) vs. **reduced** referential forms (pronouns, zero) are not evenly distributed across syntactic functions (Du Bois 1987, 2003, 2017)
- Transitive subjects (A) apparently particularly favour lexical as opposed to reduced forms
- Constraints: Avoid new/lexical A
- Dubious on both empirical and conceptual grounds (Haig & Schnell 2016; Everett 2009; etc)

