

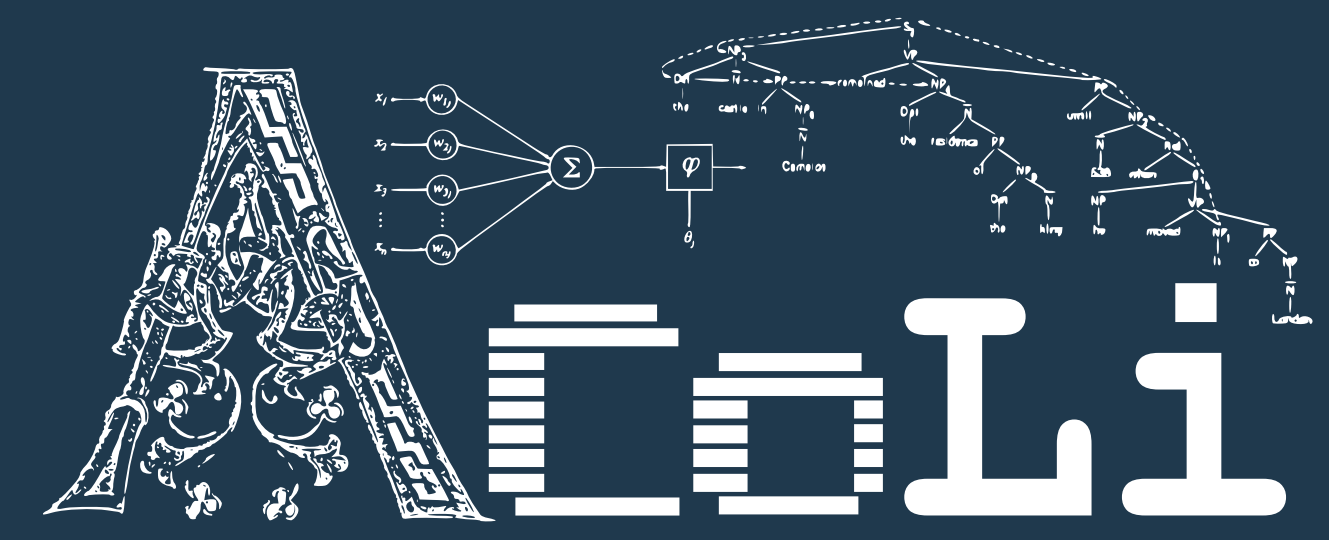
Revisiting the Givenness Hierarchy A Corpus-Based Evaluation



Universität
Augsburg
University

Christian Chiarcos

Applied Computational Linguistics (ACoLi)
University of Augsburg, Germany



Background and Motivation

Human communication is organized around **entities (discourse referents)**, their characteristics and their relations

For a speaker, effective communication usually requires

- to **keep track of entities** (supposed to be) known, identifiable or accessible to the hearer,
- to use communicative devices that allow their **unambiguous identification** by the hearer, and
- to use communicative devices that **indicate their relevance** with respect to the speaker's communicative goals

The **Givenness Hierarchy** (GH, Gundel et al. 1993) is one of several theoretical frameworks to account for this, which the following characteristics

- fine-grained predictions regarding the choice of referring expressions, esp. for demonstratives
- practical relevance (e.g., for language acquisition and human-robot interaction)*
- applied to a broad array of typologically diverse languages* * see references in paper
- comes with concise, cross-linguistically tested annotation guidelines ("coding protocol", 2006)

BUT

- almost no annotated data available
- some GH claims are controversial
 - personal pronoun > demonstrative pronoun ? (against Sgall et al. 1986)
 - demonstrative NP > definite NP ? (against Ariel 1990)
- previous studies were small-scale and usually did not involve statistic significance tests
 - results have not always been significant, especially for demonstratives

We suggest to

- replicate the original findings of Gundel et al. (1990, 1993)
- over corpora with **entity coreference corpora**, not direct annotation for GH
 - the coding protocol defines a decision tree for annotation, *to a large extent* based on (the form) of previous or subsequent mention
 - sufficient in size to expect significant results for less frequent types of referring expressions
 - available (now) for the original set of languages studied by Gundel et al. (1990, 1993)
 - Arabic, English, Chinese, Japanese, Korean, Russian and Spanish

Bootstrapping Givenness from Entity Coreference

Coding protocol (2006):

- for every referring expression, **check the statuses from highest to lowest**, annotate the highest possible status
- for every status, a number of test formulate criteria sufficient for annotation

annotate **in focus** if

- r is subject of the preceding utterance
- r mentioned earlier in same utterance
- r mentioned in both of the two previous utterances
- r is the event of the preceding utterance
- r is an inferred discourse topic

annotate **activated** (if not **in focus** and)

- r is mentioned in the two previous utterances
- r evoked by gesture or gaze
- r is an associated proposition or speech act

annotate **familiar** (if not **activated** or **in focus** and):

- r previously mentioned
- r known from shared background

annotate **unique** (if not **familiar**, etc.):

- expression contains sufficient lexical material to create a unique referent
- r linked via lexical association to activated referent

annotate **referential** (if not **unique**, etc.):

- r mentioned later in discourse
- r linguistically marked for discourse prominence

annotate **type** (if not **referential**, etc.):

- expression encodes interpretable conceptual content

corpus with entity coreference

<= coref + sentence splits + UD *nsubj*
<= coref + sentence splits
<= coref + sentence splits
(only with event coref)
(not available) 60% (3/5 tests covered)

corpus with entity coreference
<= coref + sentence splits
(not applicable: written text)
(not available) 50% (1/2 applicable tests)

corpus with entity coreference
<= coref
(not available) 50% (1/2 tests covered)

corpus with entity coreference
<= approx: ≥ 3 content words
<= possessive pronoun
up to 100% (2/2 tests, approximated)

corpus with entity coreference
<= coref
(this is circular) all non-circular tests

corpus with entity coreference
<= coref (anything subject to coreference annotation)
100% (all tests)

if coreference annotation *fails* to capture a criterion, the bootstrapped status may be lower than the actual status (but never higher) because the hierarchy is implicative, this is imprecise, but **not incorrect**

Empirical Evaluation

10 corpora, 7 languages, 4 coref formats, UD syntax (as provided / spaCy)
up to 11 types of referring expressions* per language

*exactly those defined by Gundel et al. (1990/1993)

	OntoNotes	LitBank	GUM	AnCor	NTC	KoCoNovel	ECMT	RuCor
version	5.0	CU 1.3	CU 1.3	CU 1.3	1.5	—	CU 1.3	CU 1.3
language	ar / en / zh	en	en	es	ja	ko	ko	ru
modality	written	written	written/spoken	written	written	written	written	written
genre	news, web, lit	diverse	diverse	news	news	literature	news	diverse
tokens (K)	325 / 1,750 / 235	190	170	429	1,000	165	439	145

Evaluation for English

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE	total
pron	3966	949	258	7	323	117	4720
dem.prox	105	182	38	72	39	416	416
dem.dist	117	174	42	1	35	34	403
dem.prox N	155	216	156	130	109	216	982
dem.dist N	101	76	19	54	34	90	374
the N	1046	910	1082	1169	1021	2921	8149
a N	691	130	123	731	416	1628	3719
total	5281	2617	1718	2092	2010	5045	16763

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE
pron	0.475+++	0.103+++	-0.074+++	-0.203+++	-0.073+++	-0.319+++
dem.prox	-0.01 n.s.	0.109+++	0 n.s.	-0.053+	0.032+++	-0.059+++
dem.dist	0.003 n.s.	0.125+++	0.007 n.s.	-0.051+++	-0.01 n.s.	-0.062+++
dem.prox N	-0.065+++	0.055+++	0.055+++	0.016(+)	0.003 n.s.	-0.026++
dem.dist N	-0.004 n.s.	0.026++	-0.02+	0.015(+)	-0.007 n.s.	-0.009 n.s.
the N	-0.298+++	-0.07+++	0.125+++	0.089+++	0.051+++	0.177+++
a N	-0.106+++	-0.15+++	-0.101+++	0.134+++	0.008 n.s.	0.189+++

sample results (GUM corpus, Zeldes 2017)

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE
pron	4:0	2:1	0:4	0:3	0:4	0:2
dem.prox	0:1	2:0		0:2	1:1	0:1
dem.dist	0:1	2:0	0:1	0:2	0:1	1:1
dem.prox N	0:4	3:0	3:0	0:1	1:1	0:1
dem.dist N	0:2	2:0	0:1	1:0		
the N	0:4	1:2	4:0	3:0	4:0	2:0
a N	0:4	0:3	0:3	3:0	2:0	2:0

4:0	3:0	2:0, 3:1	1:0, 2:1	1:1	0:1, 1:2	0:2, 1:3	0:3	0:4
in focus	pron	✓	✓	✓	✓	✓	✓	✓
activated	this	✓	✓	✓	✓	✓	✓	✓
	that	✓	✓	✓	✓	✓	✓	✓
	that N	✓	✓	✓	✓	✓	✓	✓
familiar	that N	✗						
unique	the N	✓						
referential	indef. this N	✓						
	n/a	not observable						
type	a N	✓						

aggregate results, 4 corpora: we count **significant** ($p \leq 0.01$) **positive** and **negative** correlations
GH predictions largely confirmed, possible exception *that N* (dem.dist N)

Aggregate Results, 6 non-English Corpora

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE
Ø	es,ja	ja,zh	es,ar	es,ko	es,ko	es,ja
pron	ru,es,ja	ru,ja,zh	es,ar	ru,es,zh	es,ko	es,ja
dem.prox	ko1	ja	es,ja,zh	ko2,zh	es,ko1	ko
dem.med	ko	ja	es,ko,ja	ko1	es,ko2	ko
dem.dist	ko2	zh	zh			es
dem.prox N	ko2	es,ja,zh	ru,es,ko1,ja,zh	zh	ko1 (ja)	es,zh
dem.med N	ko1,ja	ko1	(ja)	es,ko1	es	es,ko1
dem.dist N	zh	zh	zh			ru
the N	es,ar	es	ar	es	ar	es
N	ru,ko,ja,zh	ru,ko,ja,zh	ru,ko2,zh (ja)	ru,ja,zh	ko1	ru,ko,ja,zh
a N	es,zh,ar	es,zh,ar	es,zh,ar	zh,ar	es,zh,ar	

- ✓ GH core assumptions largely confirmed
- ✗ no consistent pattern for dem.dist, dem.med.N, dem.dist
- ✗ **unique**, **referential** and **type** don't seem to be good discriminators

Givenness Hierarchy (Gundel et al. 1993)

- hierarchy of „cognitive statuses“ ranked from highly given (**in focus**) to new, but identifiable in type (**type**)
- implicative hierarchy*
- higher status \subseteq lower statuses
 - e.g. referent r of „*that man*“ is **familiar**
 - r is also **unique**, **referential** and **type**
 - but not **in focus**, or **activated**

- in focus**: referent is the current focus of attention and highly prominent in the local context (\sim *he, she*).
- activated**: referent is present in the local context (\sim *this/that, this man*).
- familiar**: referent is known to both speaker and hearer from prior discourse (\sim *that man*).
- unique**: the referent is uniquely identifiable to hearer and speaker (\sim *the man*).
- referential**: the speaker refers to a specific but possibly unknown entity (\sim *this guy*).
- type** (type identifiable): hearer can identify the category of a referent (\sim *a man*).

in focus > activated > familiar > uniquely identifiable > referential > type identifiable
{it} {this, that, this N} {that N} {the N} {indefinite this N} {a N}

- speakers deviate from lexicalized cognitive statuses to express implicit information
- speakers can use a „lower“ expression to trigger **quantity implicatures**, e.g., when using a definite NP in place of a personal pronoun

- cross-linguistic annotation studies* support the predicted correlations
- for 7 major languages in Gundel (1990/93), for other languages later on

e.g., English
(Gundel 1993)

seems to match the predicted correlations

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE	total
it	214	1					215
HE		1					1
this		15					15
that	1	17					18
this N	1	11					12
that N		10	7				17
the N	30	95	47	108			280
indef. this N					1		1
a N					41	55	96
total	246	150	54	108	42	55	655

re-assessment:
most GH-specific claims are not statistically significant
(for the original statistics)

	IN FOCUS	ACTIVATED	FAMILIAR	UNIQUE	REFERENTIAL	TYPE	
it	0.895++	-0.373++	-0.21++	-0.311++	-0.183+	-0.212++	0.5>+>0.5 n.s.
HE	-0.03 n/a	0.072 n/a	-0.012 n/a	-0.017 n/a	-0.01 n/a	-0.012 n/a	0.5>+>0.5 n.s.
this	-0.119 n.s.	0.281 n/a	-0.046 n/a	-0.068 n/a	-0.04 n/a	-0.046 n/a	n.s. (1/1) n/a
that	-0.111+	0.286 n/a	-0.05 n/a	-0.075 n/a	-0.044 n/a	-0.051 n/a	0.5>+>0.5 n.s.
this N	-0.082(+)	0.224 n/a	-0.041 n/a	-0.061 n/a	-0.036 n/a	-0.041 n/a	0.5>+>0.5 n.s.
that N	-0.127(+)	0.14 n/a	0.195 n/a	-0.073 n/a	-0.043 n/a	-0.049 n/a	
the N	-0.479++	0.227++	0.268++	0.514++	-0.226++	-0.262++	n/a
indef. this N	-0.03 n/a	-0.021 n/a	-0.012 n/a	-0.017 n/a	0.149 n/a	-0.012 n/a	n/a
a N	-0.479++	0.227++	0.268++	0.514++	-0.226++	-0.262++	n.s. not significant

n/a: x' not applicable
n.s.: not significant
(+): marginal, p<0.05
+ : significant, p<0.01
++ : highly significant, p<0.001

Conclusions

GH assumptions largely confirmed
incl. GH-specific predictions for dem.prox, dem.med, dem.prox N

bootstrapping method confirmed
we can now study languages and phenomena **not originally addressed** by Gundel et al.
we can bootstrap training data for **GH tagging**

impulses for GH research

- Does GH really account for the differentiation among demonstratives?
- Should we abandon **unique**?
 - Doesn't differentiate well / cannot be reliably identified (**referential** and **type** can)