



Retrieving entities from publications in linguistics

Sebastian Nordhoff

2018-09-04, HIRMEOS Workshop, SUB Göttingen

Language Science Press

Linguistics

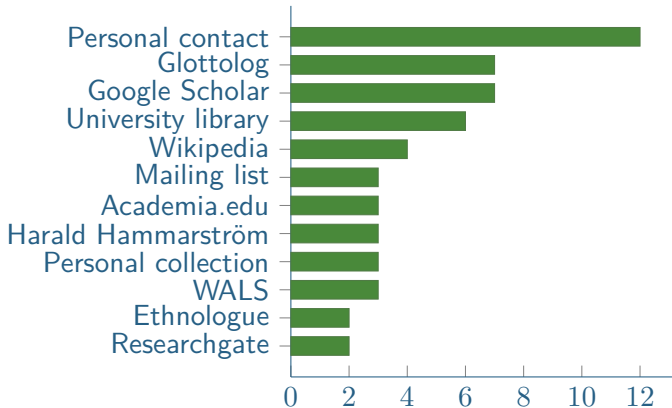
Language Science Press

NERD and linguistics

Testing NERD

- › ca. 25,000 linguists worldwide
- › both monographs and articles
- › rather long publication cycles
- › less output than for instance biology
 - › possibility to keep track
- › less sifting

> For a domain you have little expertise of, how do you find relevant literature?



question asked on list *Linguistic Typology* on 2018-08-29, no predefined answers

n=18, multiple answers possible

```
) start 2014; 75 books; 22 series
```



} Available formats for books

- } pdf
- } tex
- } bib

} Indexes in books

- } Language index
- } Subject index
- } Name index

} Indexes are a discovery tool, similar to NERD.

› A recent book on film subtitles and eyetracking had the following index candidates generated by sketchengine

image composition
eye tracking
speaking direction
typographic identity
fixation duration
audiovisual translation
aesthetic experience
title area
film material
speaker identification
film title
natural focus
text element
image track
information intake
graphical translation

title placement
bottom-centre area
typographic film
german image
split attention
gaze behaviour
reading speed
film identity
typographic film identity
tracking research
first fixation
additional language
narrative text
eye tracking research
visual attention
individual placement



› higher level goals of text and data mining:

› provide better tools for exploration:

Customers who viewed this item also viewed



› automated reasoning:

› gene \longleftrightarrow protein
 › protein \longleftrightarrow disease
 › gene \longleftrightarrow protein \longleftrightarrow disease

› stated goals (Hirneios):

1. enhance discoverability
2. aggregation (word clouds)
3. generate collections
4. highlighting

- › The following knowledge bases can be seen as resources for disambiguation
 - › **authority** (= Name Index)
 - › GND
 - › ORCID
 - › **languoids** [languages, dialects, families] (= Language Index)
 - › Glottolog
 - › **concepts** (= Subject Index)
 - › GOLD
 - › concepticon

} Several platforms have fields for “keywords” (OMP, Zenodo)

Keywords

} But should I really enter strings there?

Alternative names

hhbib_lgcode:

B. 2009)
Meson de Guadalupe
Mixtec-Mixtepec
San Juan Mixtepec
San Juan Mixtepec-Oaxaca

lexvo:

Mixtepec Mixtec [en]

multitree:

Eastern Juxtlahuaca Mixtec
Mixtec, Mixtepec
Mixteco de San Juan Mixtepec
Northern Misteko
Northern Mixteco

René de Saussure and the theory of word formation

Edited by

Stephen R. Anderson

Louis de Saussure

Classics in Linguistics 6



} GND
} ORCID

} A grammar of Komnzo

Language: Anta-Komnzo-Wára-Wéré-Kémä



Classification

- Morehead-Wasur (19)
 - Kanum (4)
 - Morehead-Maró (15)
 - Nambu (8)
 - Tonda (6)
 - Arammba
 - Eastem Tonda (2)
 - Anta-Komnzo-Wára-Wéré-Kémä
 - Anta
 - Kémä
 - Kómnyo
 - Wára
 - Wéré
 - Káncchá
 - Mblale-Rámmo
 - Rema
 - Warta Thuntai
 - Yei

Comments on subclassification

Christian Döhler 2016 :37-42

References

Showing 1 to 10 of 10 entries

Details	Name	Title	Any field	ca	Year	Pages	Doctype	ca	Provider	da
	<input type="text"/>	<input type="text"/>	<input type="text"/>		<input type="text"/>	<input type="text"/>	--any--		--any--	
citation	Christian Döhler 2016	Komnzo: A language of Southern New Guinea	✓		2016	622	grammar		hh	

Glottocode: wara1294 ISO 639-3: wra

Map



show big map

Countries

Links

Alternative names

} cross-linguistic categories don't exist

- › cross-linguistic categories don't exist
- › cross-linguistic categories don't exist

- › cross-linguistic categories don't exist
- › cross-linguistic categories don't exist
- › cross-linguistic categories don't exist

- › cross-linguistic categories don't exist
- › cross-linguistic categories don't exist
- › cross-linguistic categories don't exist
- › something called “dative” in language X cannot be equated with something called “dative” in language Y
 - › General Ontology for Linguistic Description (GOLD) tried and failed

GOLD 2010
[issues](#)
[versions](#)
[xml](#)
[owl/rdf](#)
[gold community](#)
[help](#)
[top](#) [definition](#) [usage](#) [examples](#) [properties](#) [issues](#)

Inallative Case (Concept)

<http://purl.org/linguistics/gold/InallativeCase>
[Thing](#)
[_ Abstract](#)
[_ Linguistic Property](#)
[_ Morphosyntactic Property](#)
[_ Case Property](#)
[_ Inallative Case](#)

Definition:

InallativeCase expresses that something is moving toward the region that is inside the referent of the noun it marks. It has the meaning 'towards in(side)'. Kibrik says that Archi (aqc) possesses a nominal spatial form expressing InallativeCase, namely -aši [Kibrik 1998: 470].

Usage Notes

Examples

[Properties](#)[Values](#)[Definition](#)

User Submitted Issues



- > Is there an “inallative” in *Romanite domum*?
- > Is there an “inallative” in *au foyer*?
- > Is there an “inallative” in *this*?
- > Can you equate the usages in the three examples?
- > take-home-message: it’s complicated, and automated reasoning will not work.

A typology of questions in Northeast Asia and beyond

An ecological perspective

Andreas Hölzl

Studies in Diversity Linguistics



- › *A typology of questions in Northeast Asia and beyond*
- › Book chosen as the most recent publication
- › Variety of countries, languages, ethnic groups, concepts, etc.
- › 546 pages
- › NERD running on local machine

5.8 Mongolic

Table 5.85: Spatial **deictics** in **Mongolian** according to Janhunen (2012b: 131), slightly reduced

	PROX (hearer)	DIST	INT
LOC	naa-n	tzaa-n	xaa-(n)
LOC ABI	naa-n-aas	tzaa-n-aas	xaa-n-aas
LAT	naa-sh	tzaa-sh	xaa-sh
PROI	naa-g.oor	tzaa-g.oor	xaa-g.oor

Table 5.86 shows five of the **interrogatives** that can be found in most modern **Mongolic** languages.

Table 5.86: Five **Proto-Mongolic interrogatives** and their modern representatives

	*ken 'who'	*yaxun 'what'	*alin 'which'	*kejixe 'when'	*kaxana 'where'
Dagur	seng	yoon	aly	sejer	xan
Mongolian	sen	yoon	alyh	sejee	xam
Buryat	sen	yūn	ali	sejee	xamta
Khamnigan Mongol	ken	yeen	ali	kejee	kamta
Ordos	ken	yūn	ali	kejee	kai
Written Oirat	ken	yuu/n	ali	kejee	xamig(h)a
Oirat	ken	yuu/n	al - al-k	kene	xama
Kalmuk	ken	yūn	aly(-k)	keza	xama, aly-d
Shira Yuzhur	ken	yima	aali	kejee	xana
Santa	ken	yang	ali	giczi	khala
Bonan	kang	yang	ane	keet(-)	hala
Kangxi	ko	jo - jai	am(ve)	gadje	yana
Huzhu Mongghul	ken	ya/n	ali	kijet	an-j(i)
Minhe Mangghuer	kan	ya, yang	alyge	kejie	ang(ji)
Monghol	ken	kyan	emah - imas etc.	keja	?

According to Janhunen (2003d: 20) the stem *ke- originally had the meaning 'who' as well as 'what', which is an unlikely scenario from a cross-linguistic point of view. As has been shown by Cysouw (2005), the only place worldwide where this pattern is not **extremely rare** or altogether absent is **South America**.

Proto-Mongolic had two resonances (submorphemes), one in *k- that is still present in most Mongolic languages but changed to x- in **Dagur**, **Buryat** and **Mongolian**, and one in *y- that has survived up to today. Similar changes from *k- to > *x- can be seen in **Turkic**

ALY

Normalized: Upper Aramite language

Domains: Astronomy, Biology, Geography, Sociology

cont: 0.4157



Aramite or Aramite or more specifically Upper Aramite (Upper Aramite), is a **language** spoken in and around **Aramite** (Aramite in Aramite) in the **Southwest**, Australia. The name is sometimes spelled **Aramite** or **Aramite**.

Freebase ID	/m/59288
writing system	Latin script
number of speakers	[exact figure]
instance of	Dialect continuum
UNESCO Atlas of the World's Languages in Danger ID	168

References: [W](#) [B](#)

- › NERD retrieved some pretty specialized concepts
 - › Recall is good
- › NERD also retrieved a lot of irrelevant concepts (“South America”) or lookalikes (business names, radio stations)
- › NERD was rather aggressive and colored whole pages.
- › the system seems to have understood that the book is about linguistics and often selects a linguistic concept. However, sometimes, the concept chosen is off the mark (Australia).
 - › Precision is low.

- › Installation procedure was OK
- › Loading the book in the browser worked out of the box
- › Loading the book in the browser takes several minutes

Questions from a publisher

› in how far does NERD help the readers/authors?

› Exploration/Discovery

› currently, discoverability of content via series, e.g.

Contemporary African Linguistics

› linguists seem to prefer personal/social interaction to automated recommender systems

› Automated reasoning

› limited potential given the fuzzy nature of cross-linguistic concepts

Goals (reprise)

› stated goals (Hirneios):

1. enhance discoverability: better than curated series?
2. aggregation (word clouds): better than index?
3. generate collections: better than curated series?
4. highlighting: is color a value in itself?