

A Quantitative Analysis of Linguistic Metadata

Kavon Hooshlar
University of Hawai'i at Mānoa

ICLDC4, 2/27/15

This talk

Show some pretty graphs

Explain lots of horrible assumptions

See if you think this is meaningful by the end if it

Inspiration

How does language size relate to endangerment?

How does endangerment relate to documentation?

Size and Endangerment

“Lots of small endangered languages”

“They don’t correlate because of outliers”

Endangerment and Documentation

“There are so many languages going extinct that aren’t documented”

“Endangerment alone can’t be a reason to document a language”

An example



Paul Newman

SOAS Department Seminar

Oct. 15, 2013

**The Law of Unintended
Consequences: How the
Endangered Languages
Movement Undermines Field
Linguistics as a Scientific
Enterprise.**

Check the data

Each topic implies knowledge of correlations

Number of Speakers vs Endangerment

Endangerment vs Documentation

Modern linguistic resources claim to give us access to this data, so lets actually check...

The data

Goal: a set of metadata for every language

- Number of speakers
- Endangerment level
- Metric of completed linguistic research

Completed Linguistic Research?

ELCat Documentation Index

Raw scores given for Boasian Trilogy

(Grammars, Dictionaries/Lexicon, Texts/Corpora)

My goal is to follow this model

using previously compiled data

Data -> Information

Linguists collect metadata on individual languages during their work

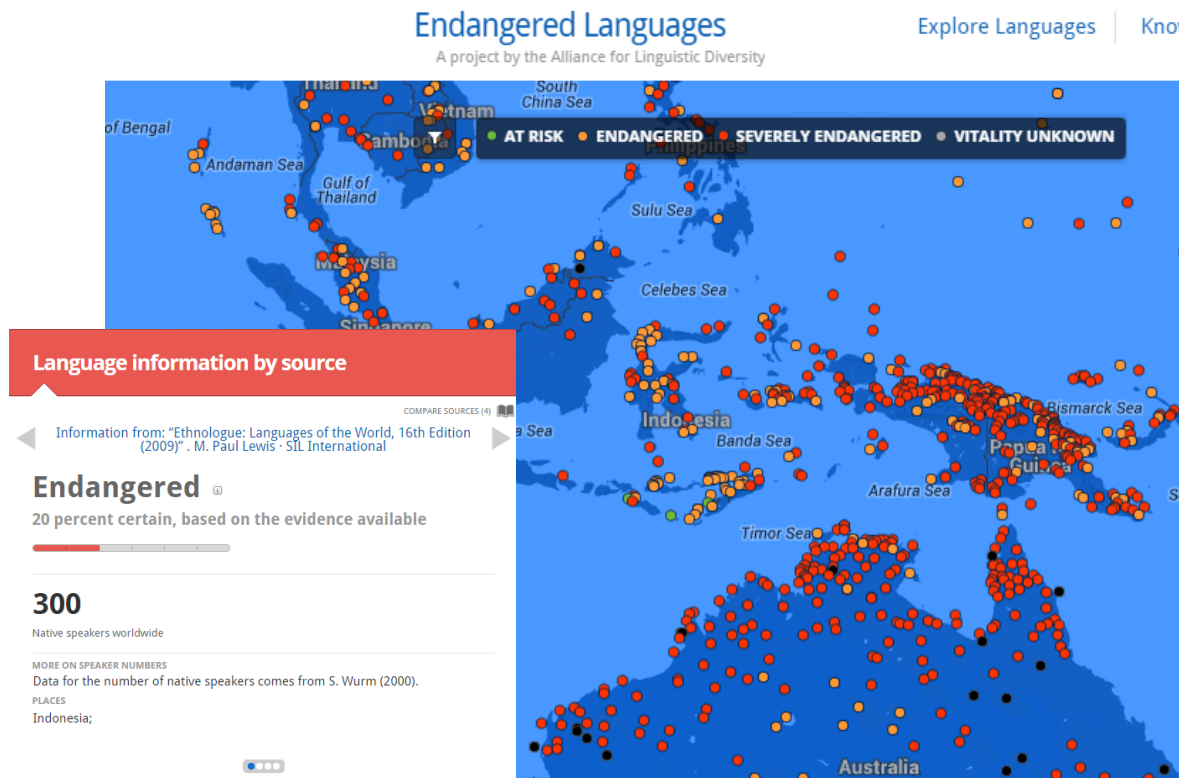
Other linguistics compile all of this metadata

Still other linguists (me!) process the data

Assumption: My analysis can only be as good as the work of those who came before me



Potential data sources: ELCat



endangeredlanguages.com


Pro: Info quality

Con: endangered languages only

Potential data sources: OLAC

Participating Archives • OLAC • Delivered by the Penn Libraries
Printer-Friendly Page

OLAC Language Resource Catalog

Search for language resources 

AB1-043

Title: AB1-043

Link to the object: <http://hdl.handle.net/10125/29556>

Online: Yes

Archive: [Kaipuleohone](#) (see archive description)

Contributor: Pius, Anna (speaker)
Berez, Andrea L. (researcher)
Berez, Andrea L. (depositor)
Gabriel, Waim (participant)

Date: 2013-07-09

Description: Anna Pius talks about her son's Karim Leg event.
Region: Pex Village, Asaro District, Eastern Highlands Province, Papua New Guinea

Content language: Kuman

Subject language: Kuman

Language family: Trans-New Guinea
Papuan

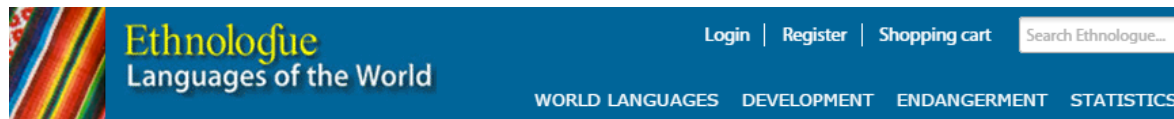
Find Related Information:

- ☒ Archive: Kaipuleohone
- ☐ Online: Yes
- ☐ Subject language: Kuman
- ☐ Language family: Papuan
- ☐ Language family: Trans-New Guinea
- ☐ Geographic region: Pacific
- ☒ Country: Papua New Guinea
- ☐ Linguistic type: Primary text
- ☐ DCMI type: sound
- ☐ Content language: Kuman
- ☐ Date: 2000 and later
- ☐ Date: 2010 - 2019

Pro: Archive info

Con: Unorganized,
unmaintained

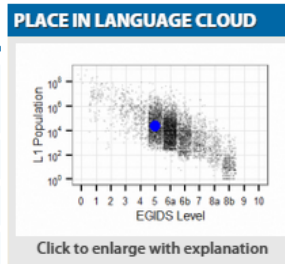
Potential data sources: Ethnologue



Anal

Print

LANGUAGE		FEEDBACK
A language of <u>India</u>		
ISO 639-3	<u>anm</u>	
Alternate Names	Namfau	
Population	23,200 in India (2001 census). Population total all countries: 23,250.	
Location	Southeast Manipur state: Chandel district, Chandel, Chakpikarong, and engnoupal subdivisions, on Chakpi river banks. Possibly in Bangladesh.	
Language Maps	<u>India, Map 5</u>	
Language Status	5 (Developing).	
Classification	<u>Sino-Tibetan, Tibeto-Burman, Sal, Kuki-Chin-Naga, Kuki-Chin, Northern</u>	
Dialects	Laizo, Mulsom. Reportedly most similar to Lamgang [<u>lmk</u>] (Kuki Naga).	



ETHNOLOGUE PRODUCTS

Languages of India
An Ethnologue Country Report,
282 pp., \$24.95
\$24.95
[Add to cart](#)

Pro: All languages

Con:

Resource data

Potential data sources: Glottolog

Glottolog Home Languages Families L-Search **References** R-Search About

Name / glottocode / iso

References

Showing 1 to 100 of 194,261 entries

← Previous 1 2 3 4 5 Next →

Details	Name	Title	ca	Year	Pages	Doctype	ca	Provider
	<input type="text"/>	<input type="text"/>		<input type="text"/>	<input type="text"/>	--any--		--any--
more	Resmi P. 2010	Process of convergence: a comparative case study of the language of kudumbis and konkanis of Kerala		2010	125	specificfeature , wordlist , sociol		hh
more	L'Unité Bozo 1980	Guide de Transcription et Lexique Bozo [Tieyaxo sawananbaana yee a xarabuye]		1980	130	grammarsketch , dictionary		hh
more	Brohez 1905, 1905	Ethnographie Katangaise: Population et Colonisation (Ethnographie : Les Balubas)		1905	44	ethnographic		hh
more	Faridanona 1977	Rantimbôlana diksionera Tsimihety		1977		dictionary		hh
more	ILV 1966	Hablemos Español y Huasteco		1966	140	wordlist		hh

Pro: All languages
Resource data
“One man job”

Con: Archive data
“One man job”

Hammarström, Harald. "Automatic annotation of bibliographical references with target language." *Proceedings of the Workshop on Multi-source Multilingual Information Extraction and Summarization*. Association for Computational Linguistics, 2008.

LANGUAGE

FEEDBACK

A language of Papua New Guinea

ISO 639-3

gim

Alternate Names

Labogai

Population

22,500 (Wurm and Hattori 1981).

Location

Eastern Highlands Province, Okapa district; small area in adjacent Chimbu Province.

Language Maps

Papua New Guinea, Map 9Papua New Guinea, Map 10

Language Status

5 (Developing).

Classification

Trans-New Guinea, Kainantu-Goroka, Gorokan, Fore

Dialects

East Gimi, West Gimi (Gouno).

Language Development

Literacy rate in L1: 5%–15%. Literacy rate in L2: 5%–15%. NT: 1994.

References

Showing 1 to 100 of 194,261 entries

Details▲	Name	Title	ca	Year	Pages	Doctype	ca	Provider
	<input type="text" value="Search"/>	<input type="text" value="Search"/>		<input type="text" value="Search"/>	<input type="text" value="Search"/>	--any-- ▼		--any-- ▼
more	Resmi P. 2010	Process of convergence: a comparative case study of the language of kudumbis and konkanis of Kerala		2010	125	specificfeature, wordlist, socling		hh
more	L'Unité Bozo 1980	Guide de Transcription et Lexique Bozo [Tieyaxo sawananbaana yee a xarabuye]		1980	130	grammarsketch, dictionary		hh
more	Brohez 1905, 1905	Ethnographie Katangaise: Population et Colonisation (Ethnographie : Les Balubas)		1905	44	ethnographic		hh
more	Faridanona 1977	Rantimbôlanja diksionera Tsimihety		1977		dictionary		hh
more	ILV 1966	Hablemos Español y Huasteco		1966	140	wordlist		hh

My “completed research” metric

A percentage of completion for each language: 0-100%

Grammar 0-33%

Each grammar rated
for page numbers
and age

Diminishing returns
on multiple
grammars

Dictionary 0-33%

Total number of
pages

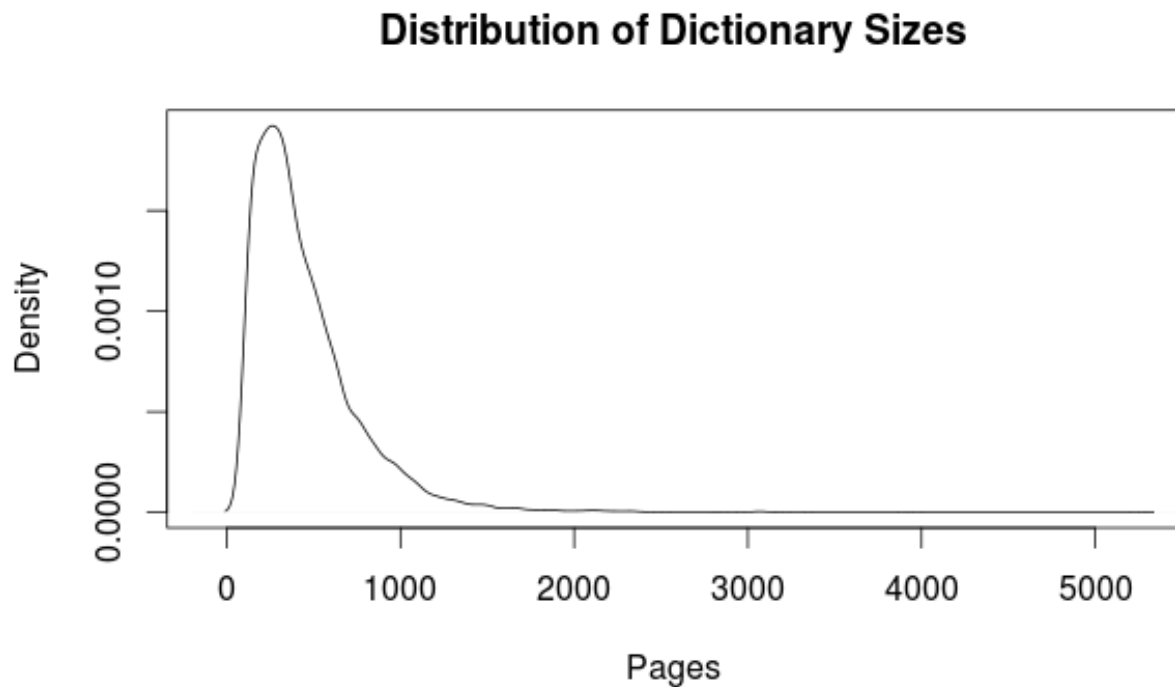
250pp is 33%

Texts 0-33%

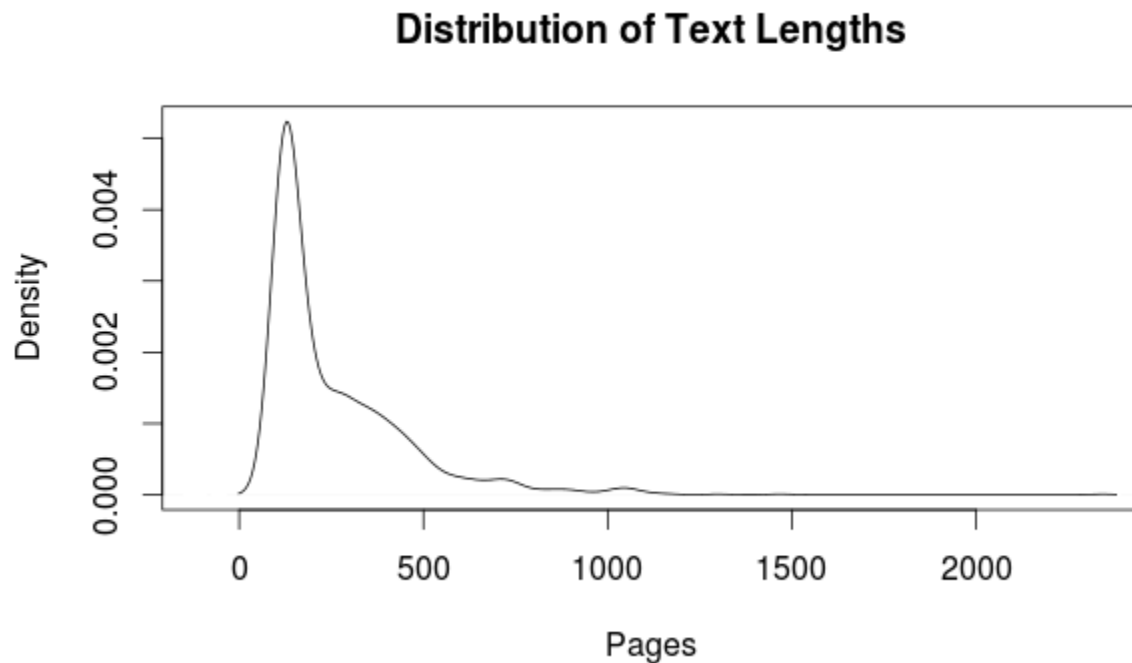
Total number of
pages

100pp is 33%

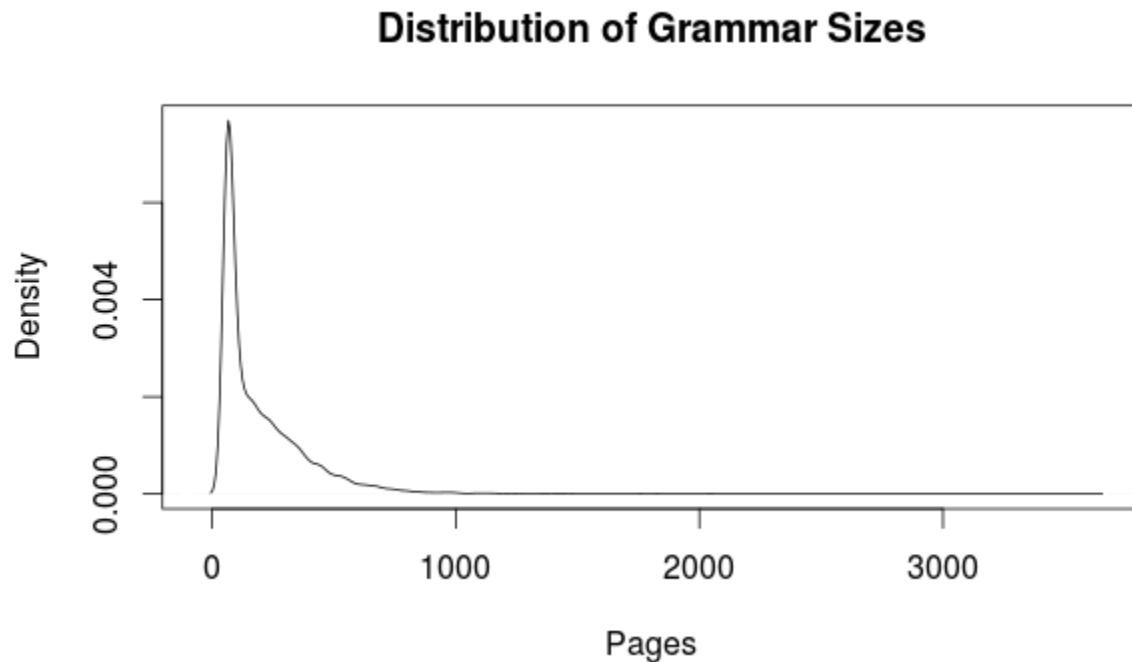
Dictionary sizes in the data



Text lengths in the data



Grammar sizes in the data



Issues

Reliance on pages numbers

Texts vs archived resources

The multiple label problem

(affecting up to 650 languages)

Other assumptions



Speaker numbers: total of L1 speakers

Only using data that has an associated ISO code

Some data created through algorithms

Data pulled in February, doesn't include items in press

Repetition:

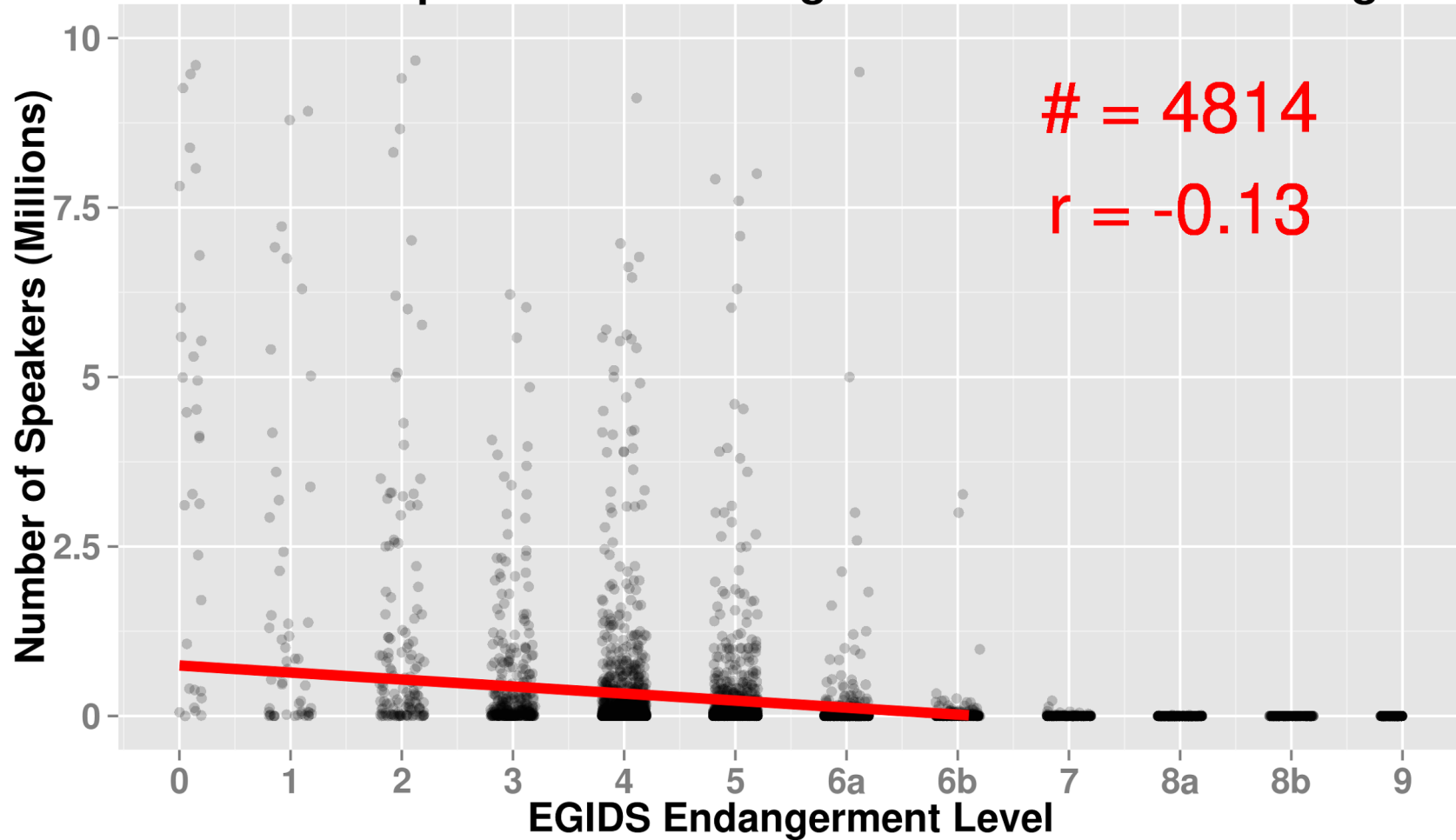
Information is only as good as the data it's based on

One final thought...

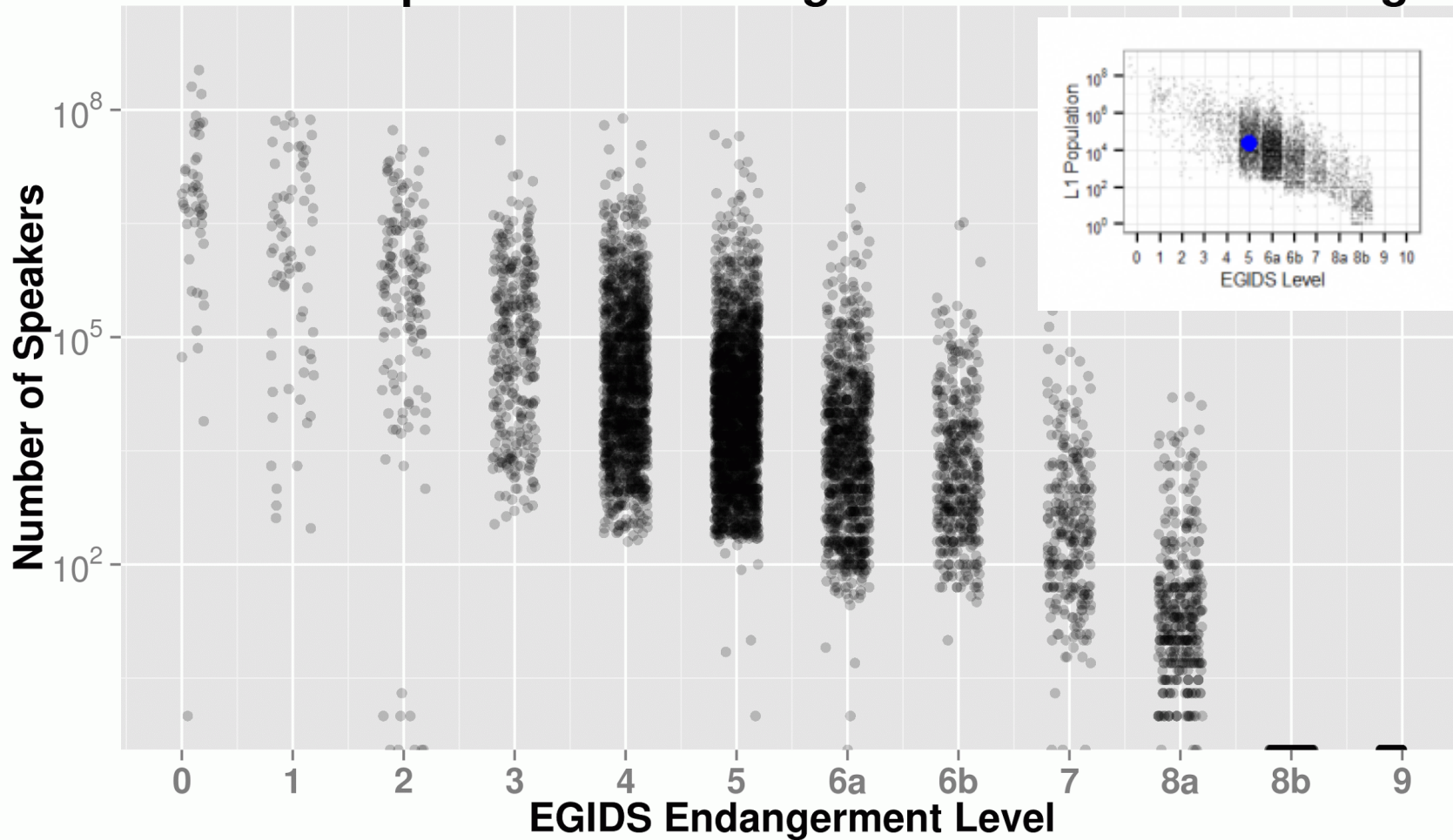
General strategy: sacrifice *some* amount of accuracy to make the analysis feasible

General goal: look for correlations and patterns among large sample sizes so that variance and outliers are acceptable

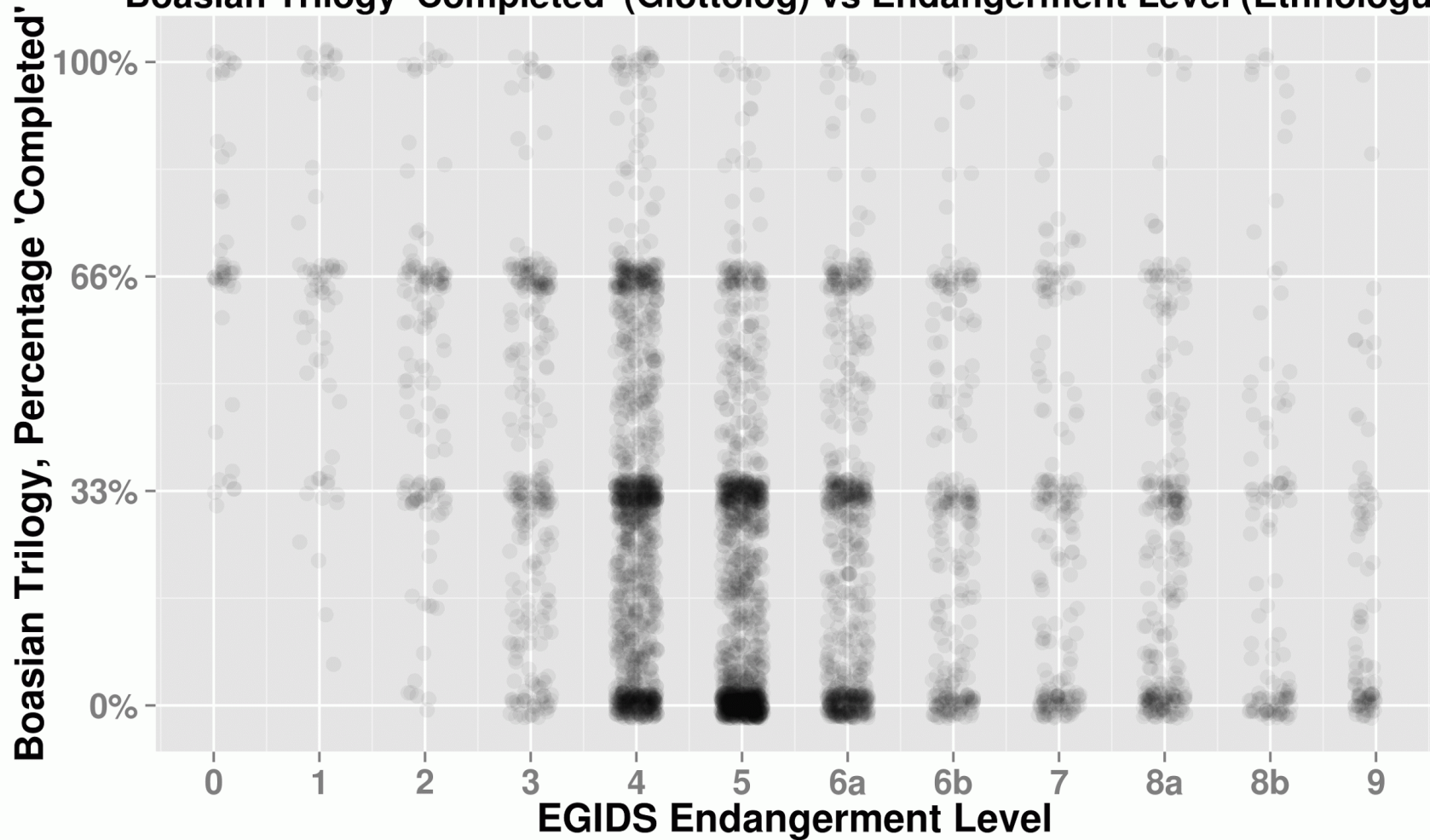
Number of Speakers vs Endangerment Level on Ethnologue



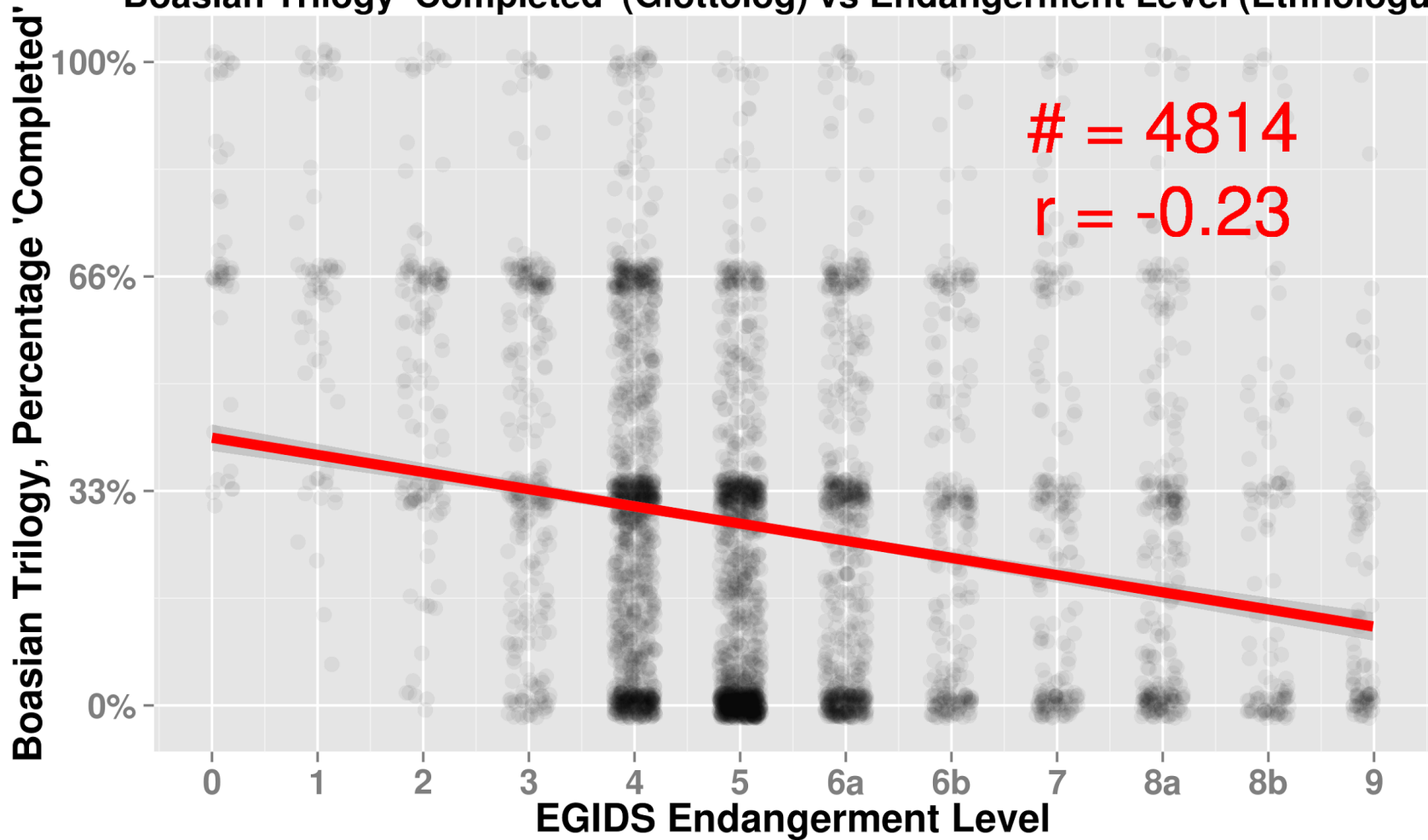
Number of Speakers vs Endangerment Level on Ethnologue



Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



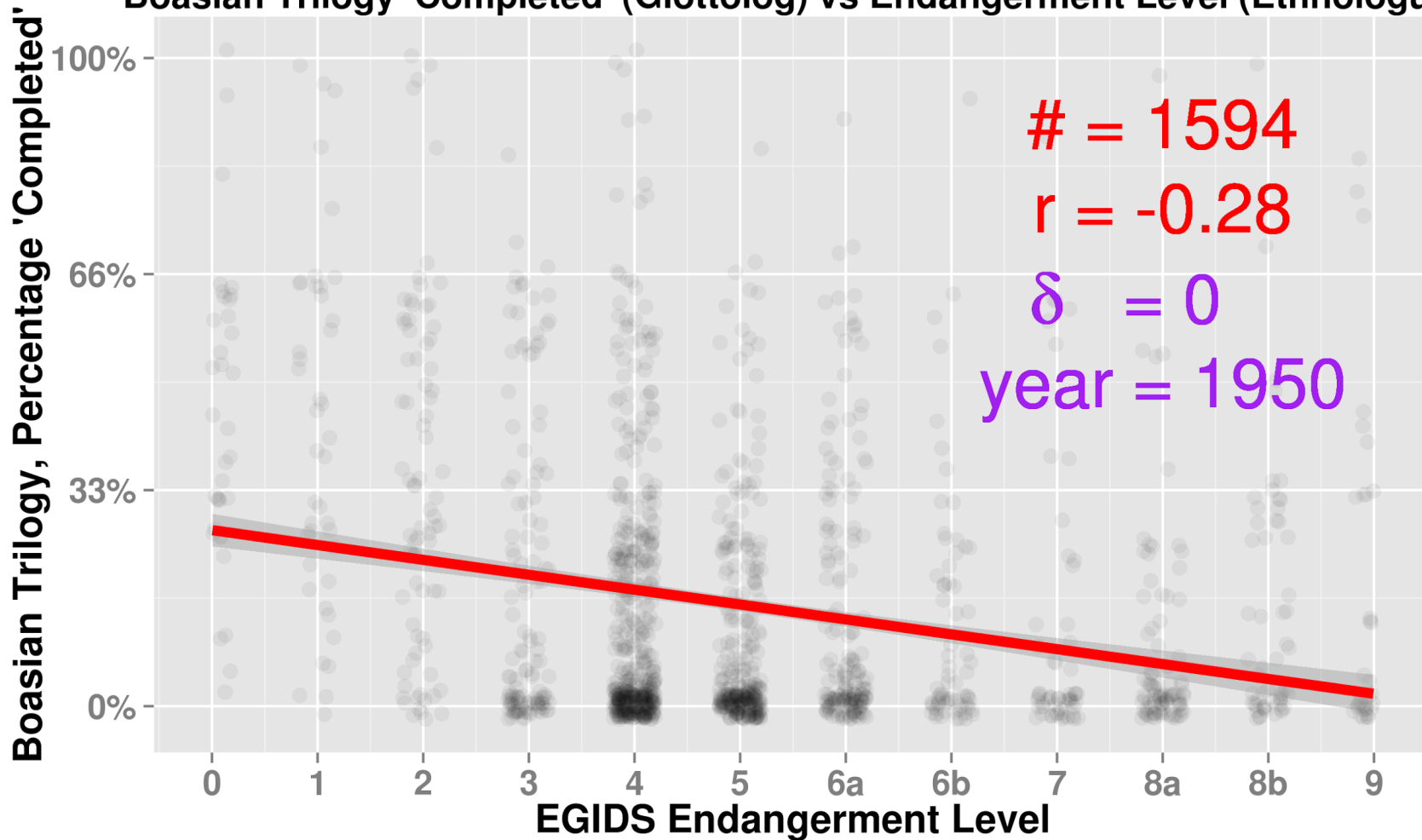
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



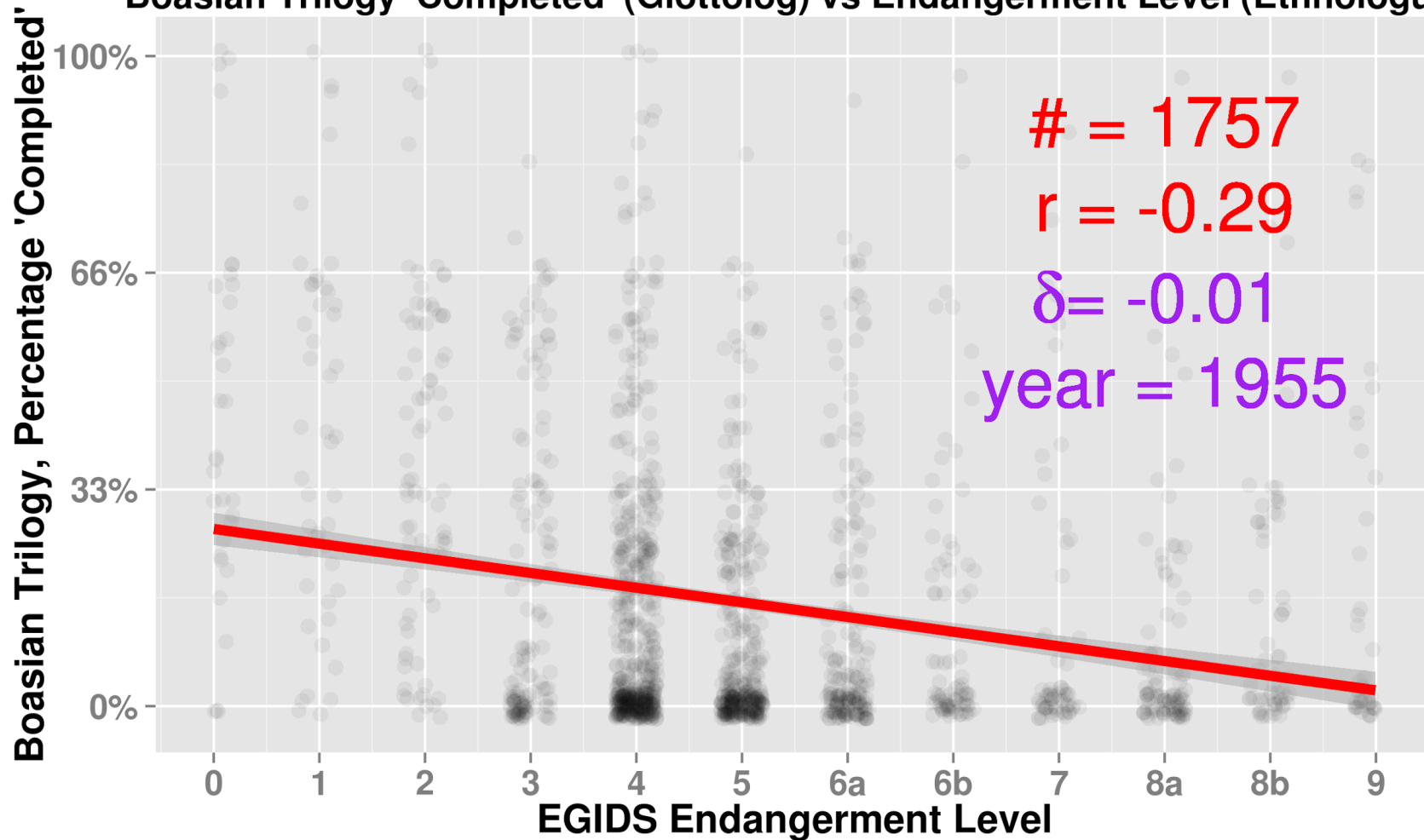
What next?

time travel, of course

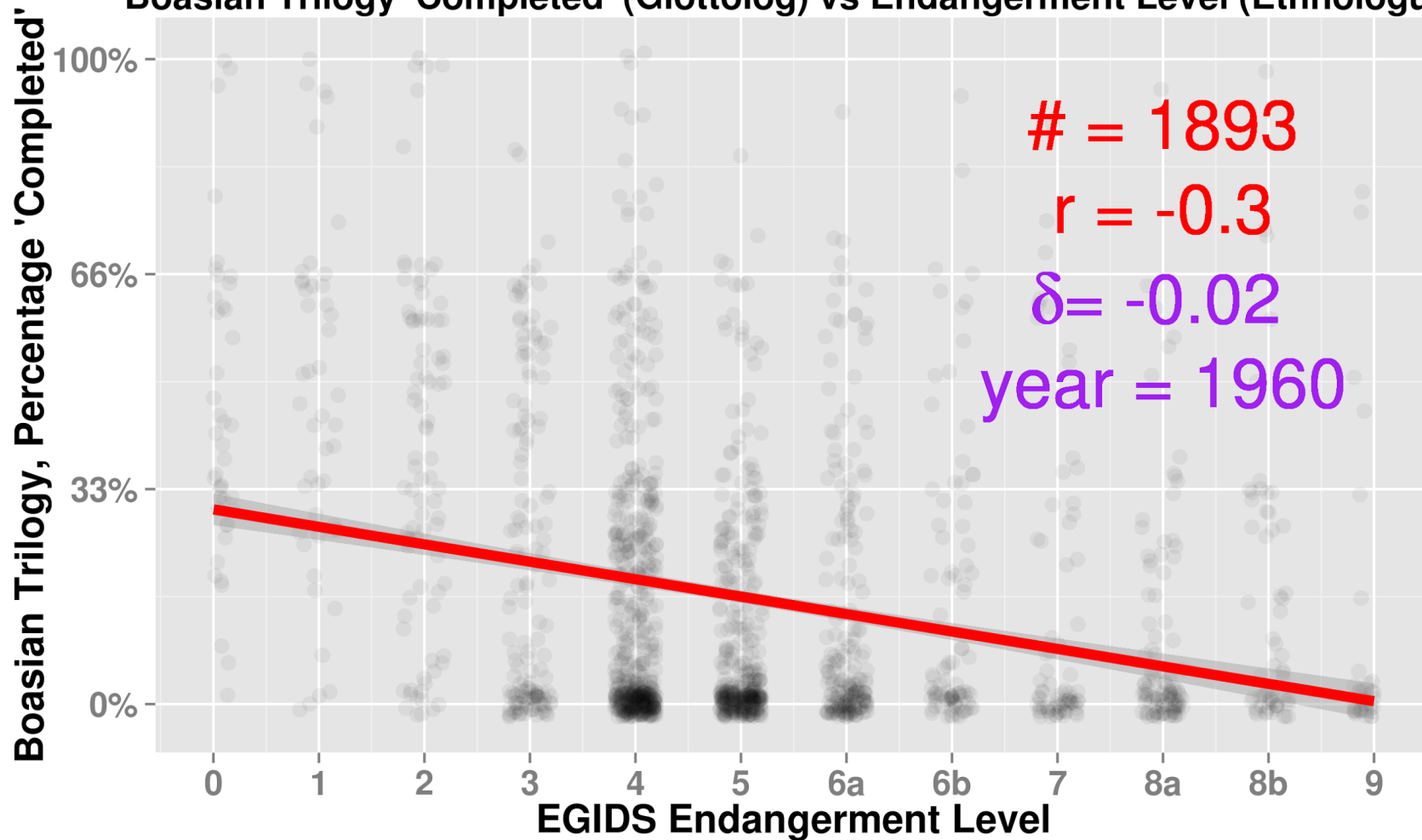
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



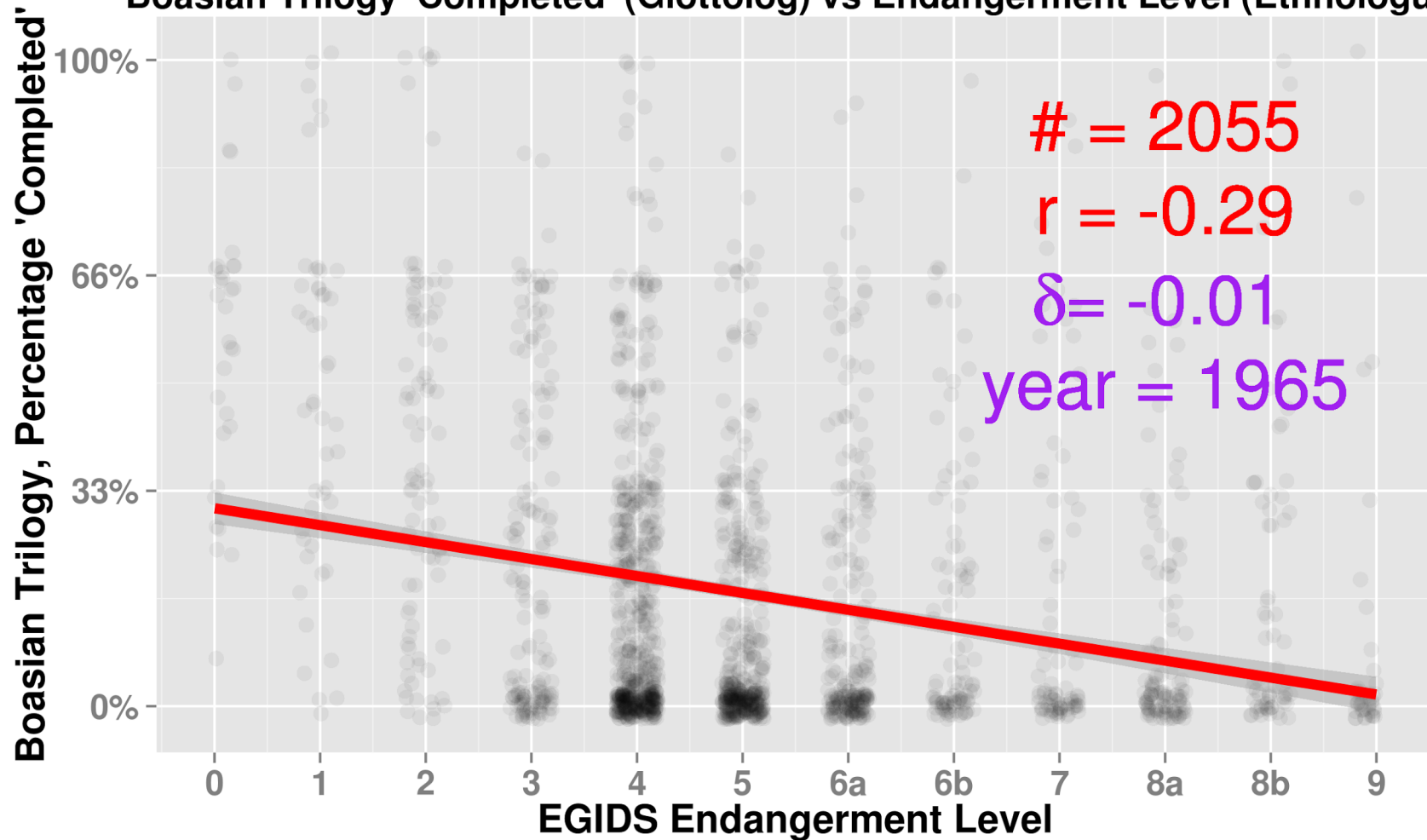
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



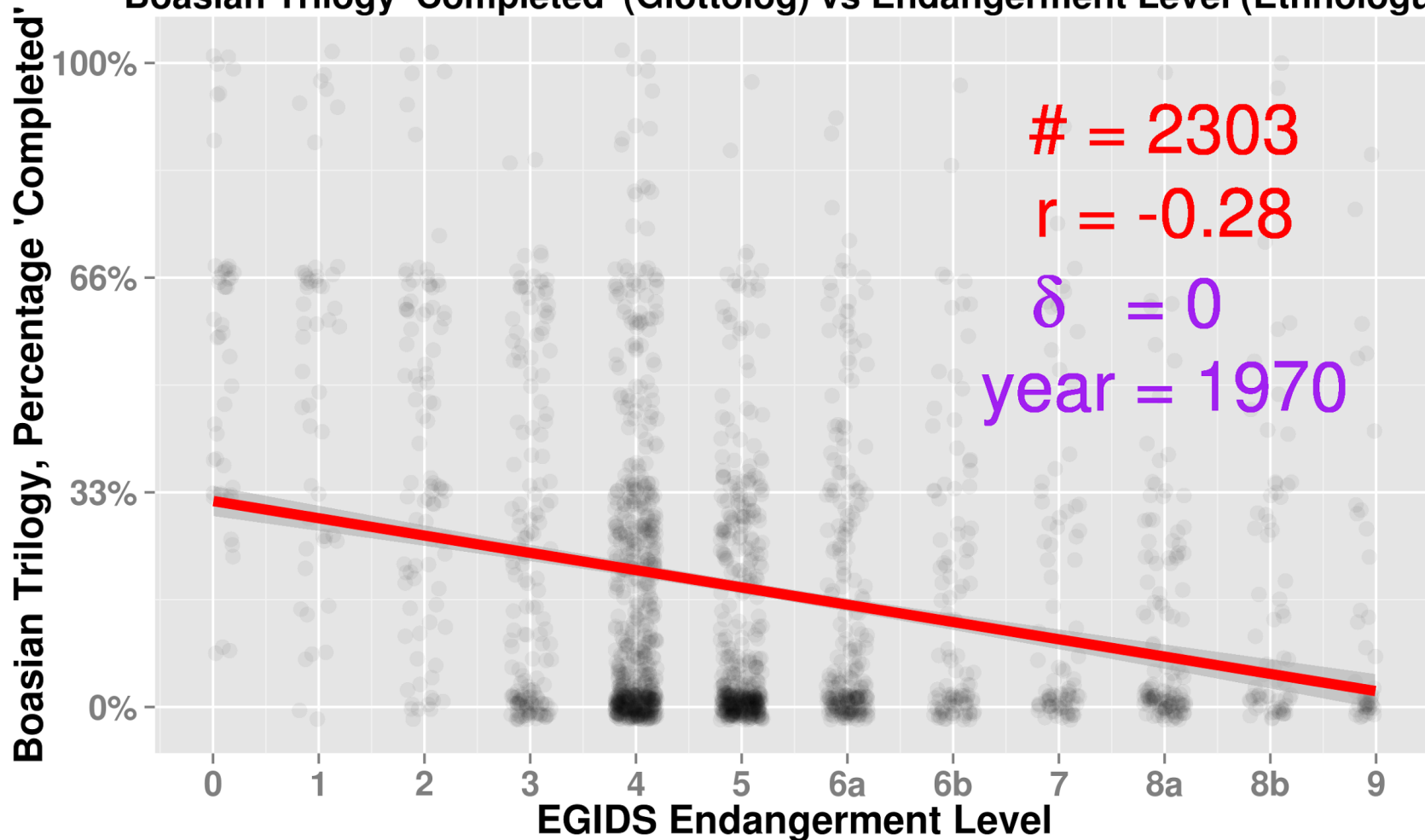
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



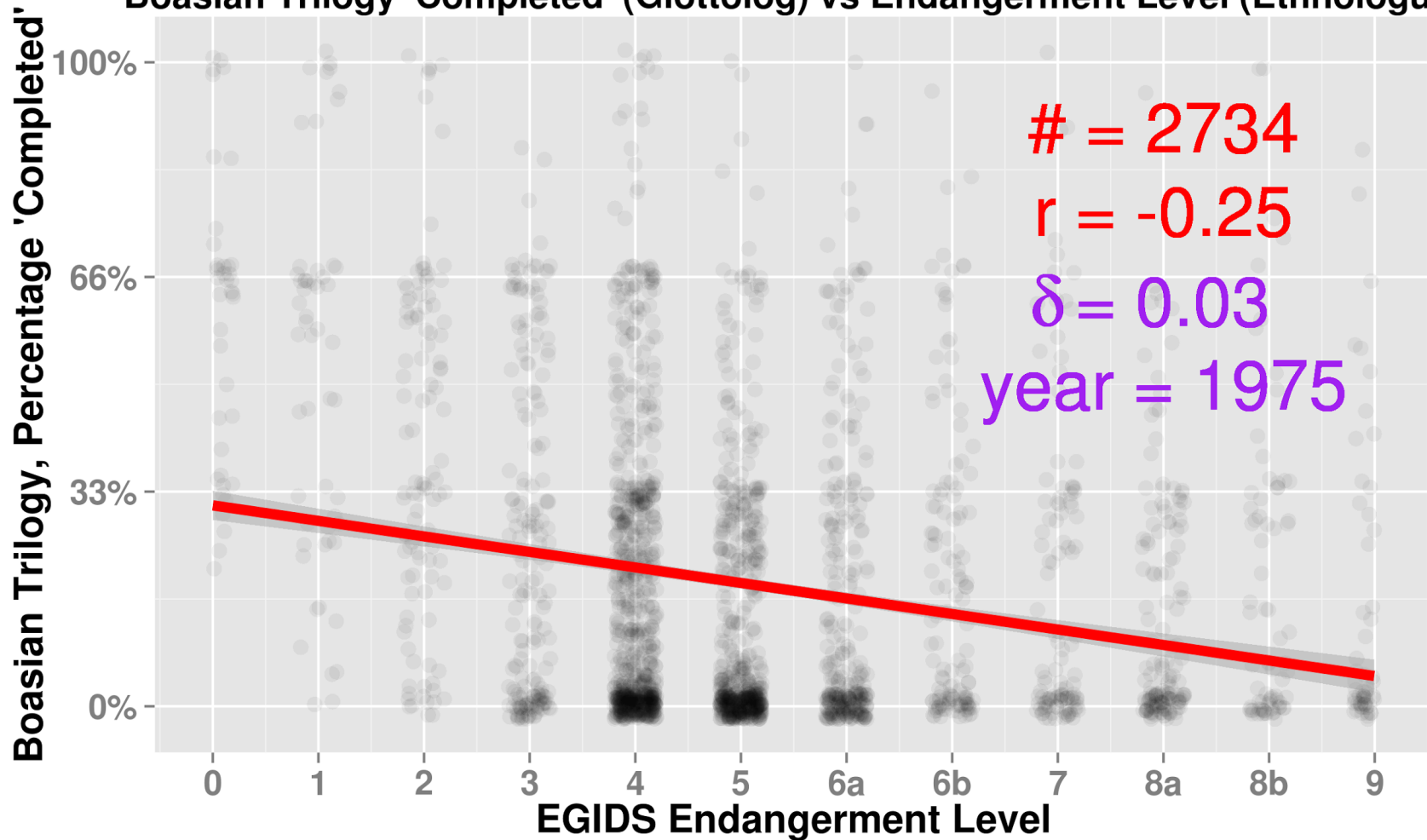
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



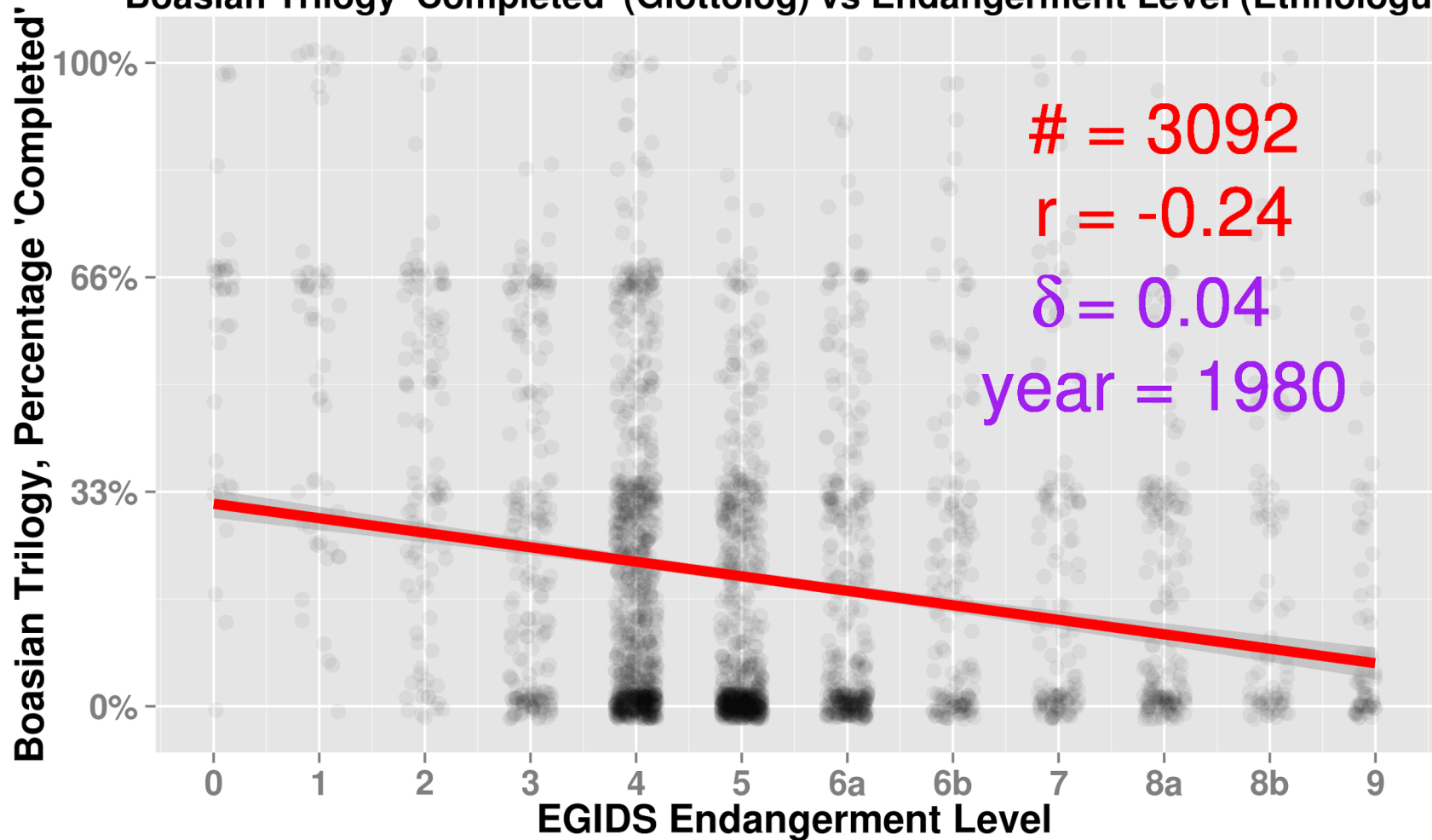
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



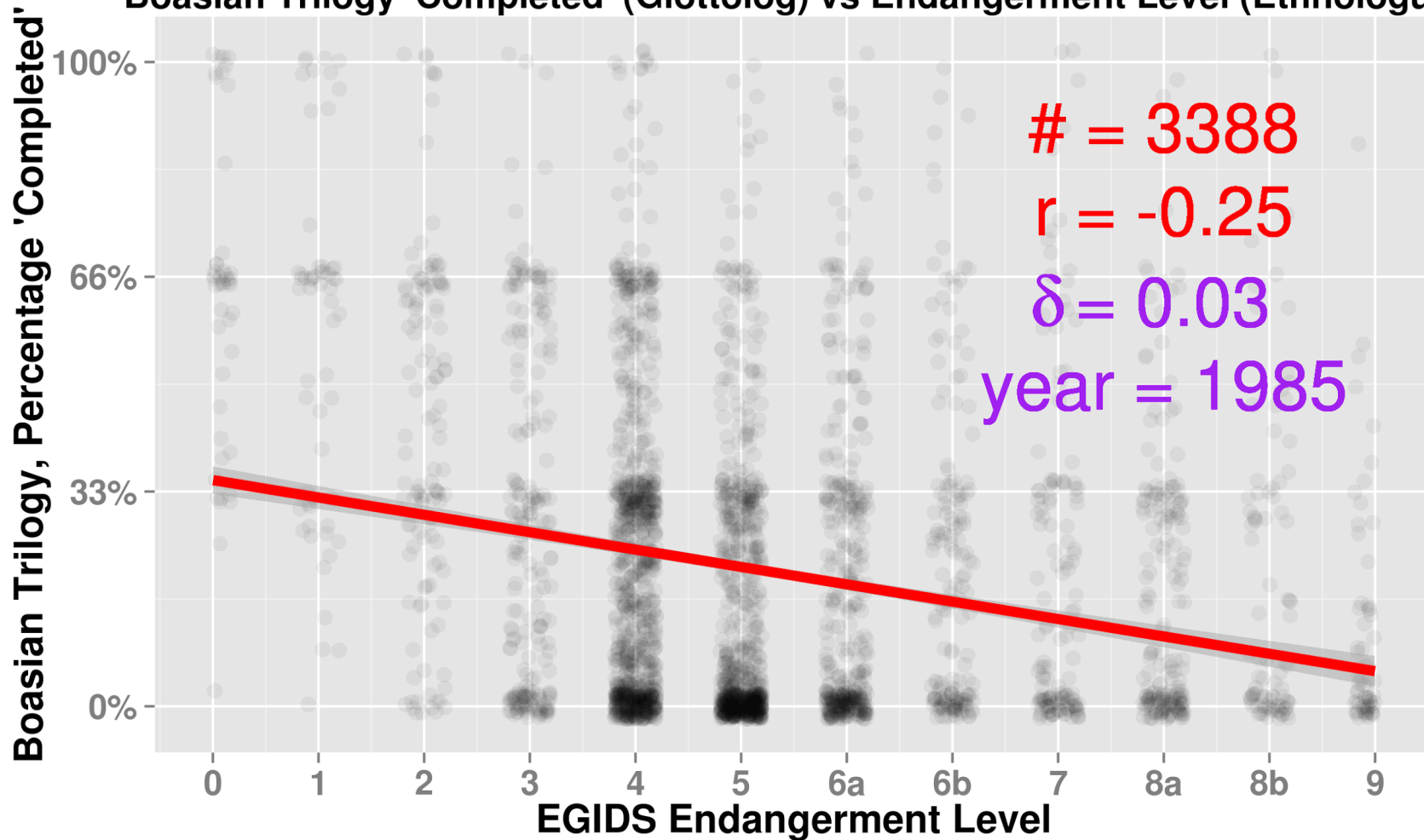
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



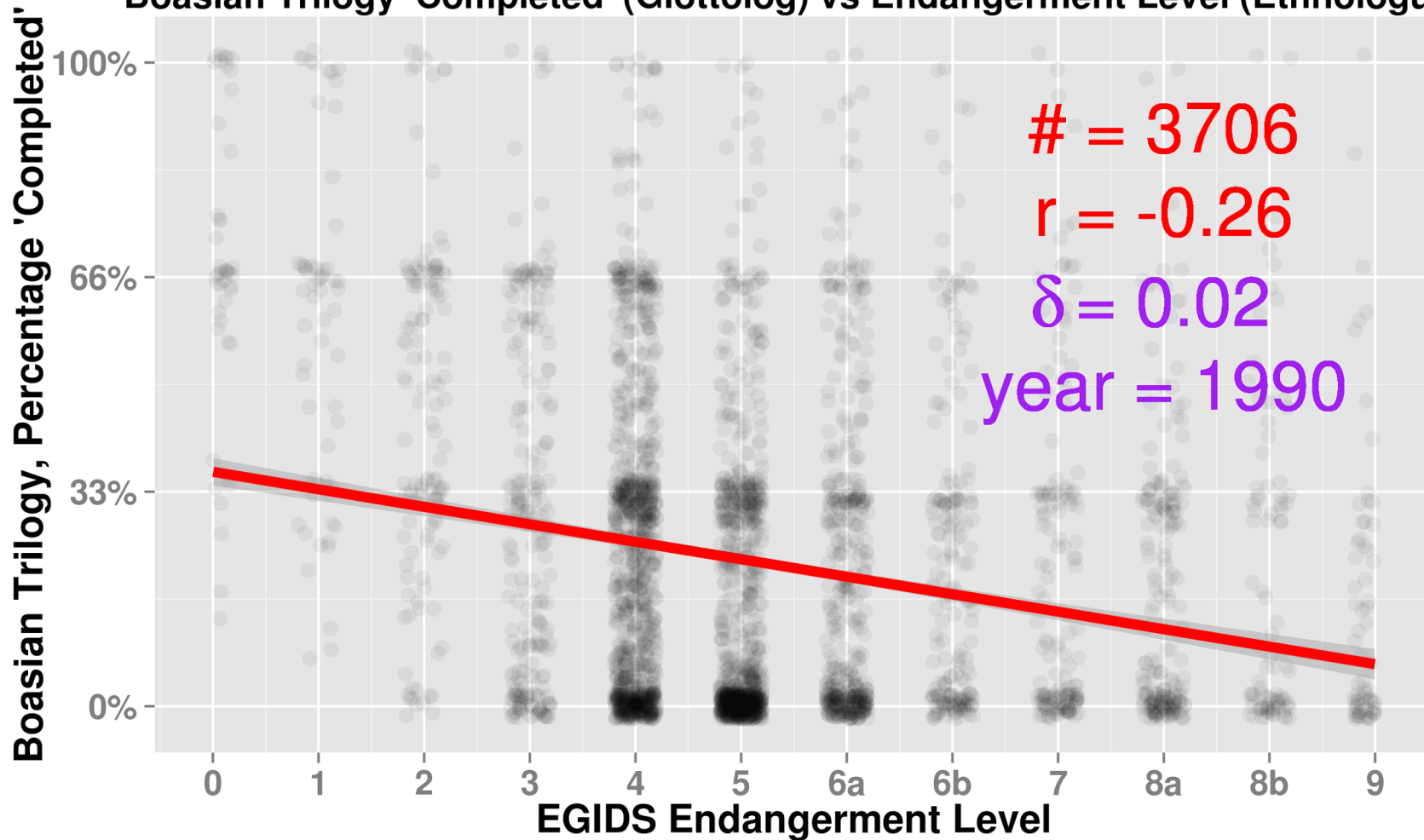
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



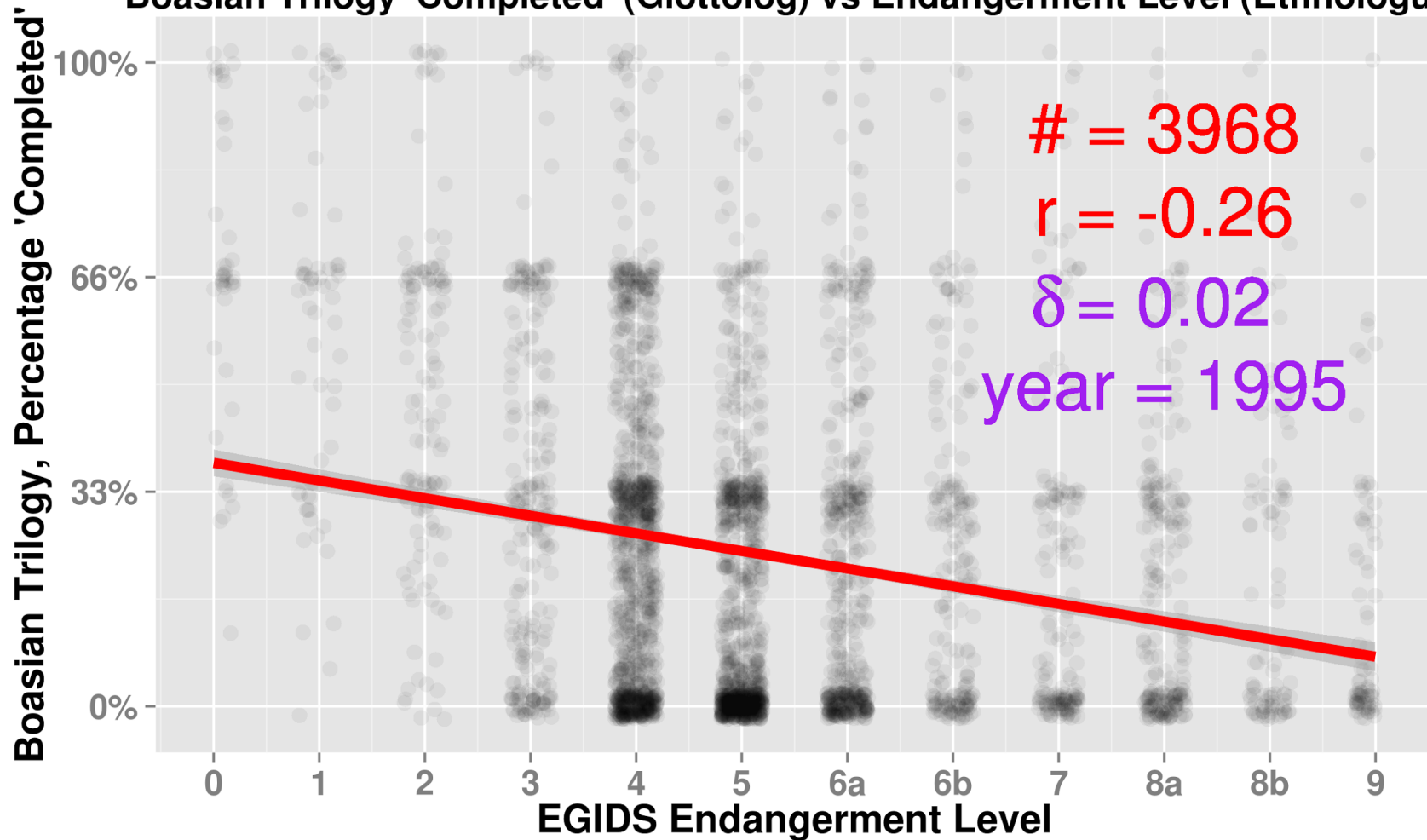
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



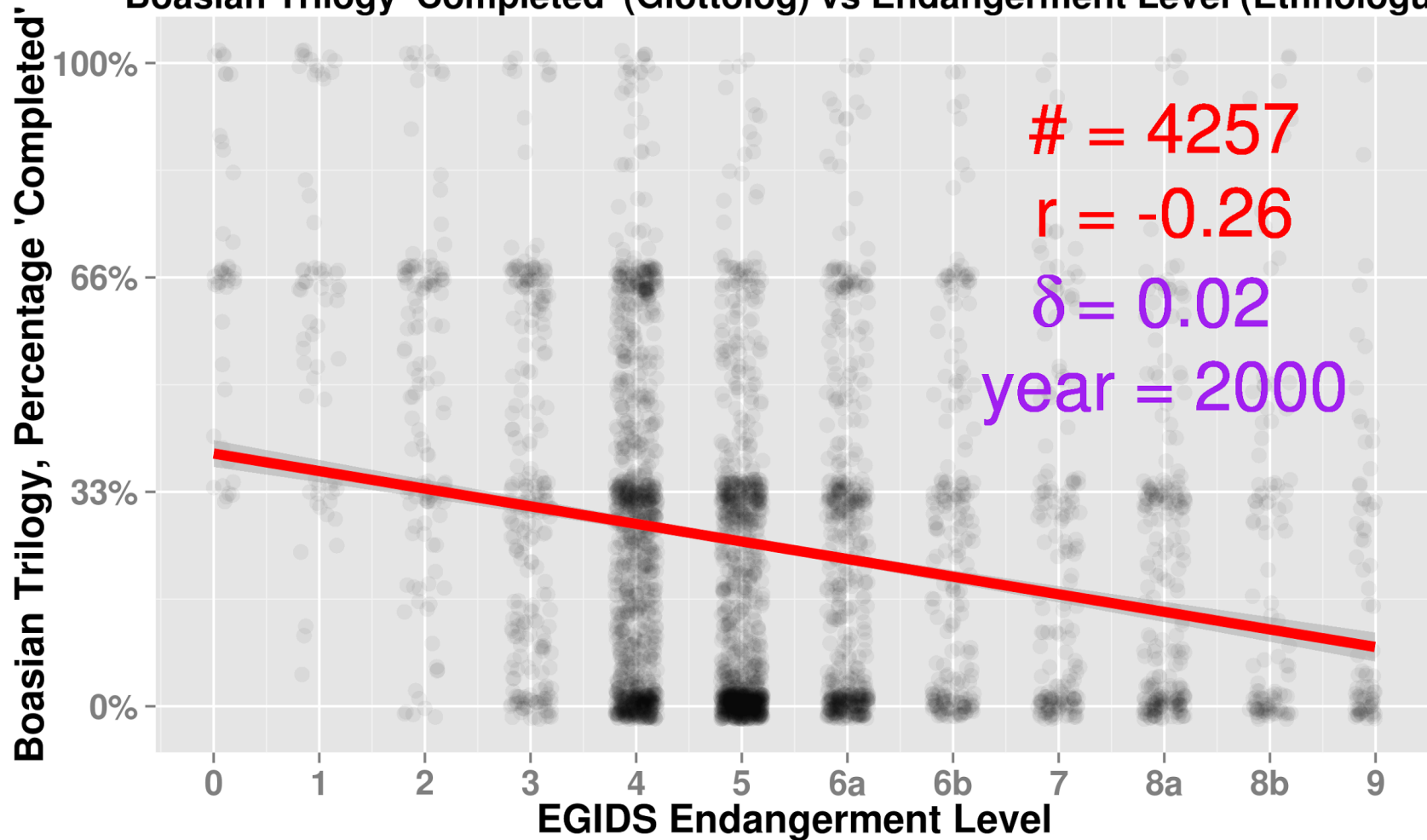
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



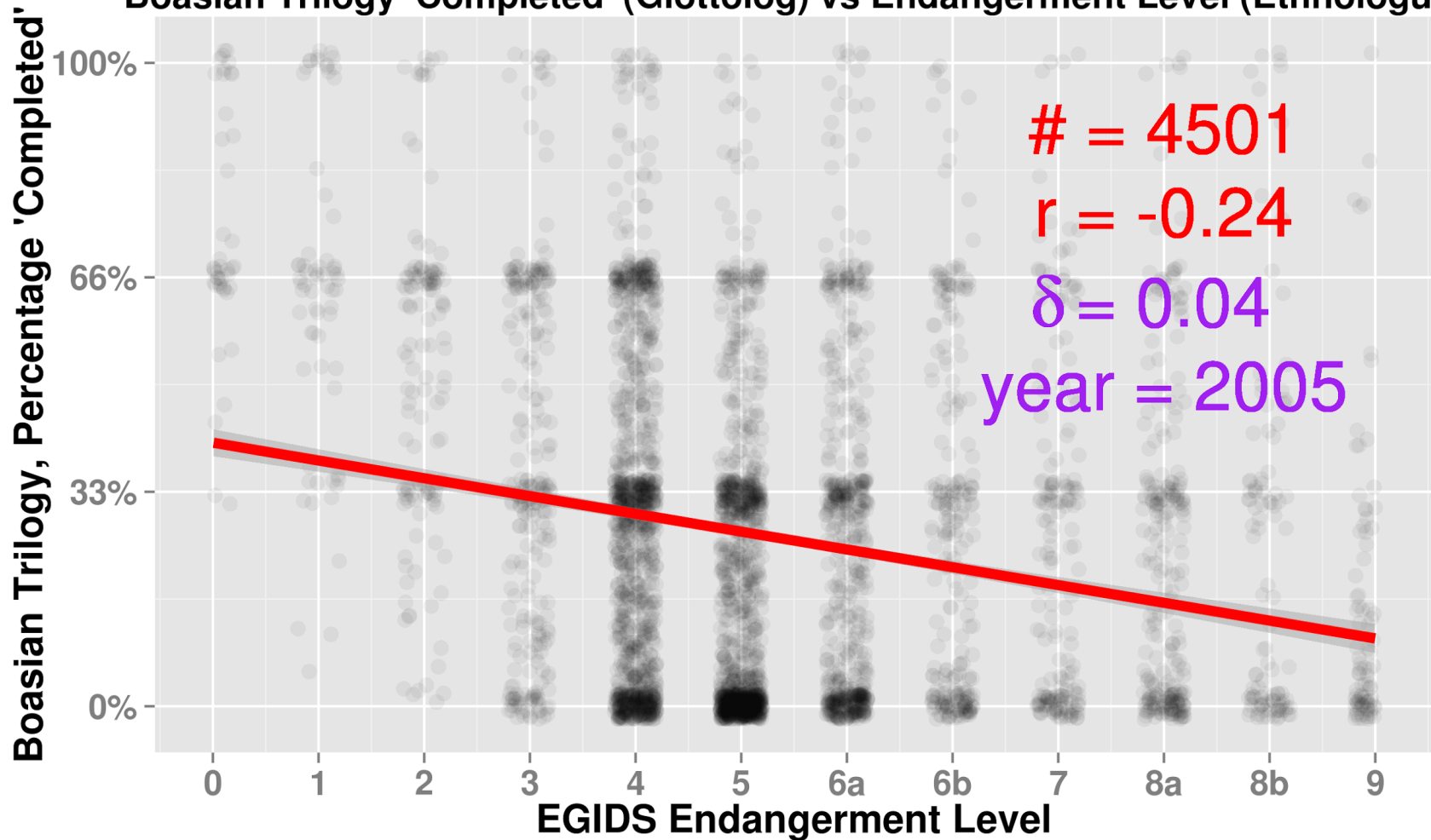
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



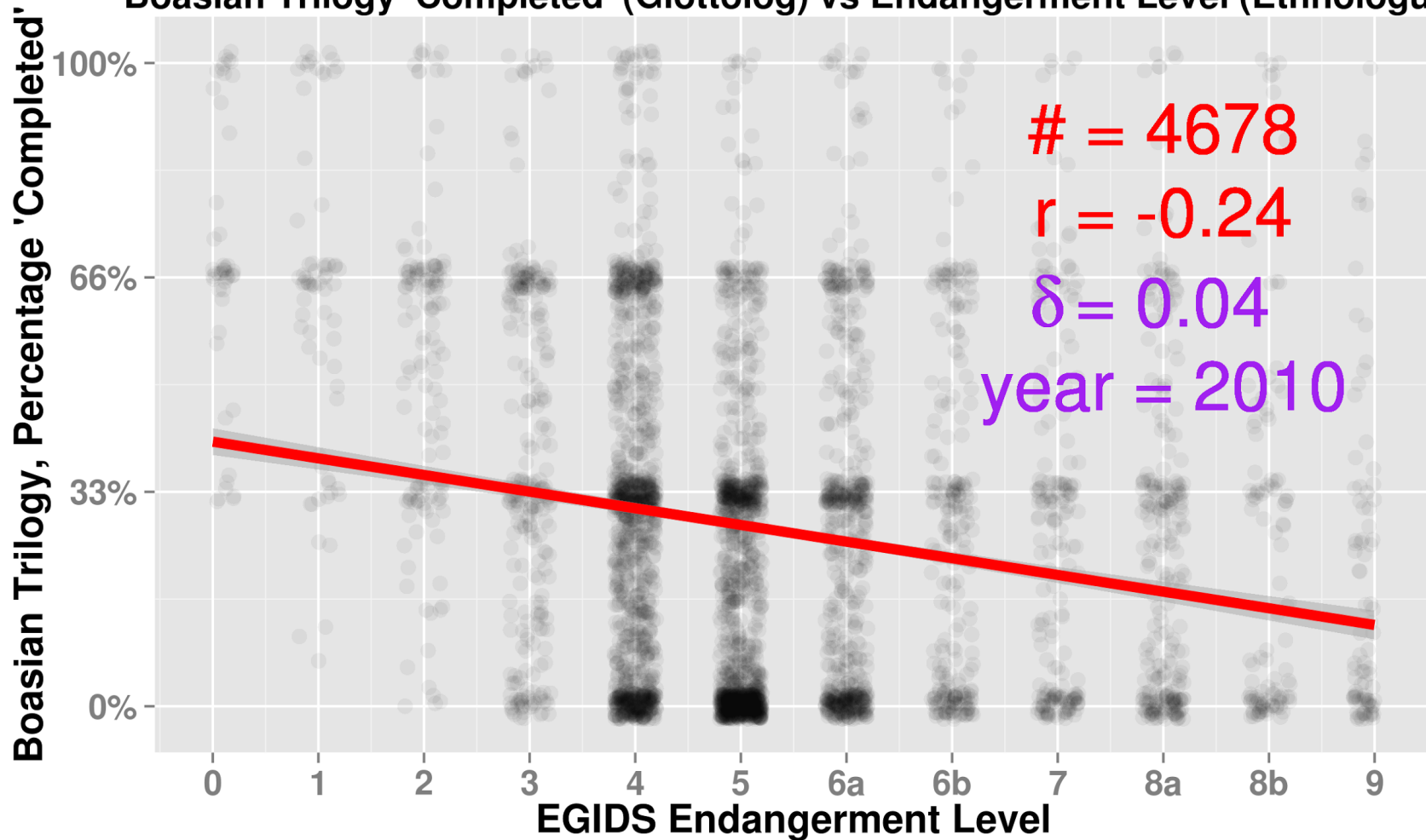
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



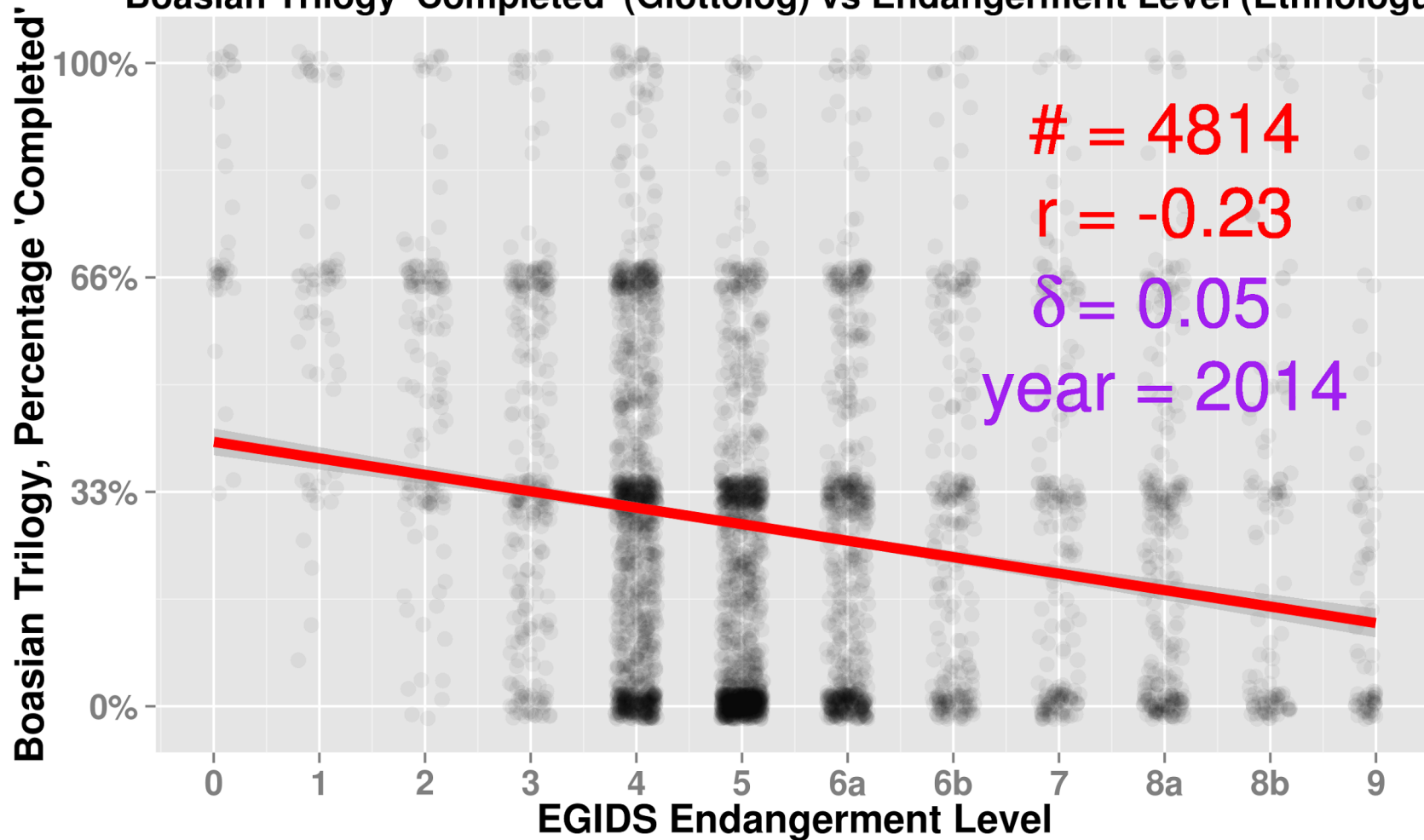
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



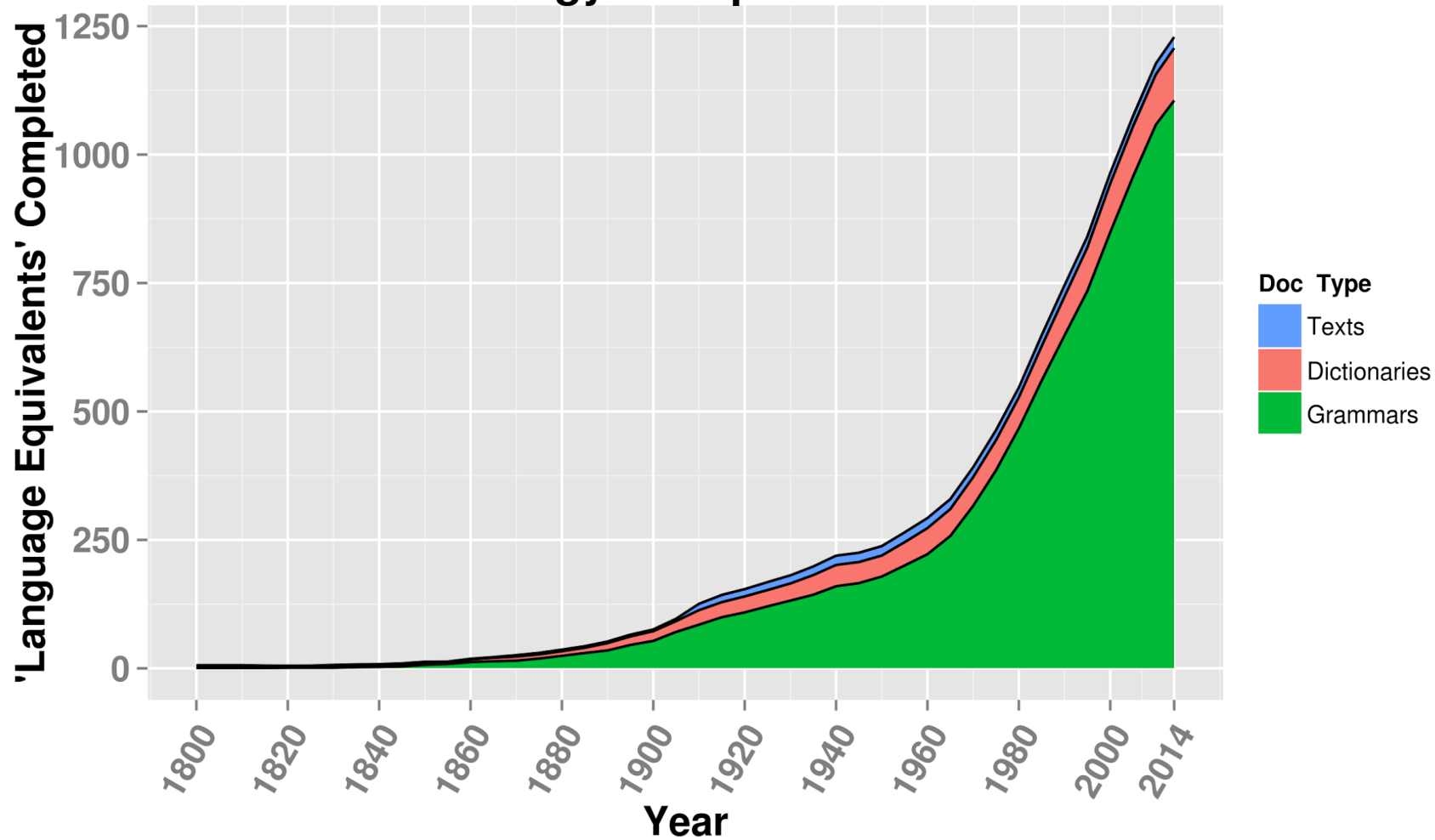
Boasian Trilogy 'Completed' (Glottolog) vs Endangerment Level (Ethnologue)



Finally,

“Boasian completion” over time

Boasian Trilogy 'Completion' over time



Conclusions

I didn't assume I could answer any new questions,
my goal was to use data to verify what we qualitatively assume

Seeing the magnitudes of known phenomena is valuable

Lots of room for improving and expanding my methods

Lets move from a "one-man" to crowdsourced model

mahalo nui loa!