

# Social Learning of Knowledge Representations in Wikipedia

Peeter Tinit<sup>1</sup>, Stefan Hartmann<sup>2</sup>

1-Tallinn University; 2-University of Bamberg

## Introduction

- When we present knowledge about the world, it is useful to represent similar things in similar ways. This helps an experienced reader (conventions) and may grasp something about reality (iconicity).
- This can be a difficult task: humans are notoriously bad at maintaining clean categories of information<sup>1</sup>.
- In this project, we study the degree of information standardization in Wikipedia.

- Are Wikipedias becoming more structured?
- Is that structure becoming conventional?
- What mechanisms support it?

An exploratory corpus-based social learning study.

## Methods

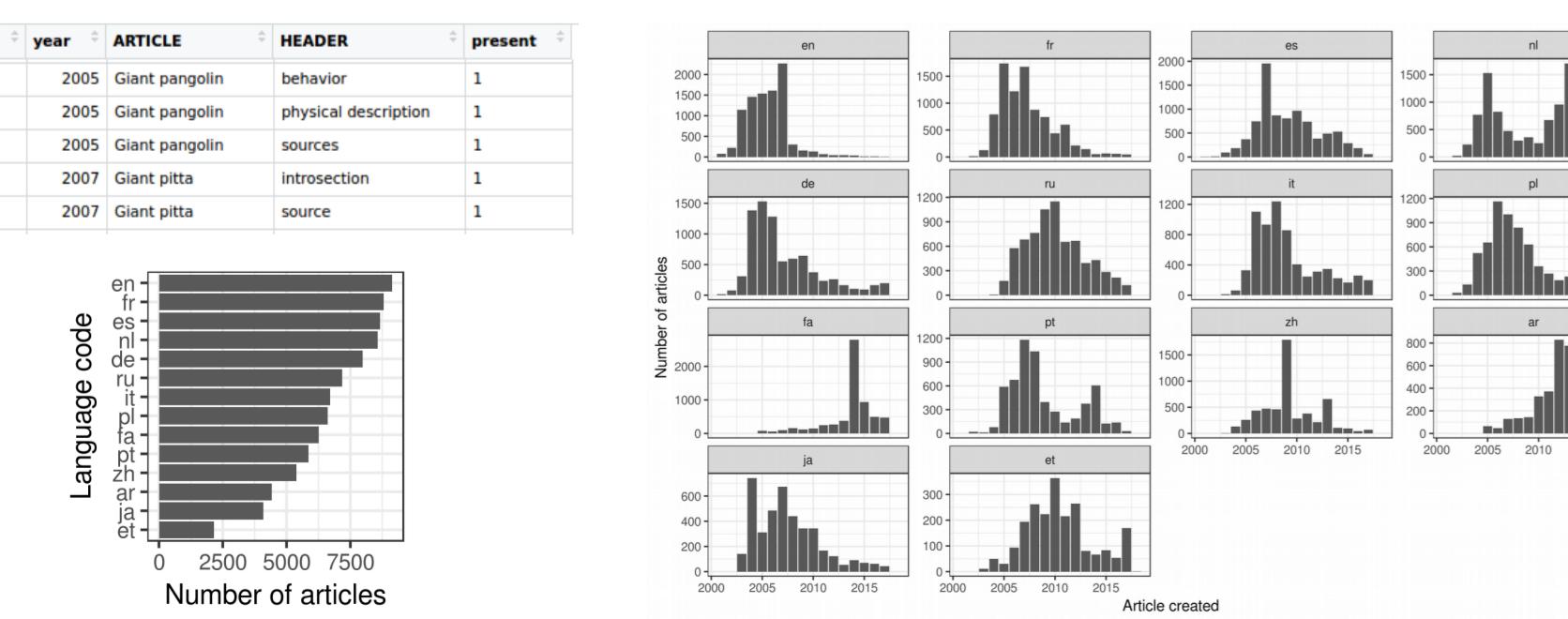
### Object

- Wikipedia - collaborative encyclopedia.
- ~50M users, 200k active writers.
- Ideally, one phenomenon – one article.
- Information distributed into subsections based on aspects of the phenomenon.

### Data

- Corpus of Wikipedia subsection headers
- Article histories on 9189 bio. species, described in at least 20 Wikipedias.
- 13 most active Wikipedias + Estonian
- Presence or absence of headers in time

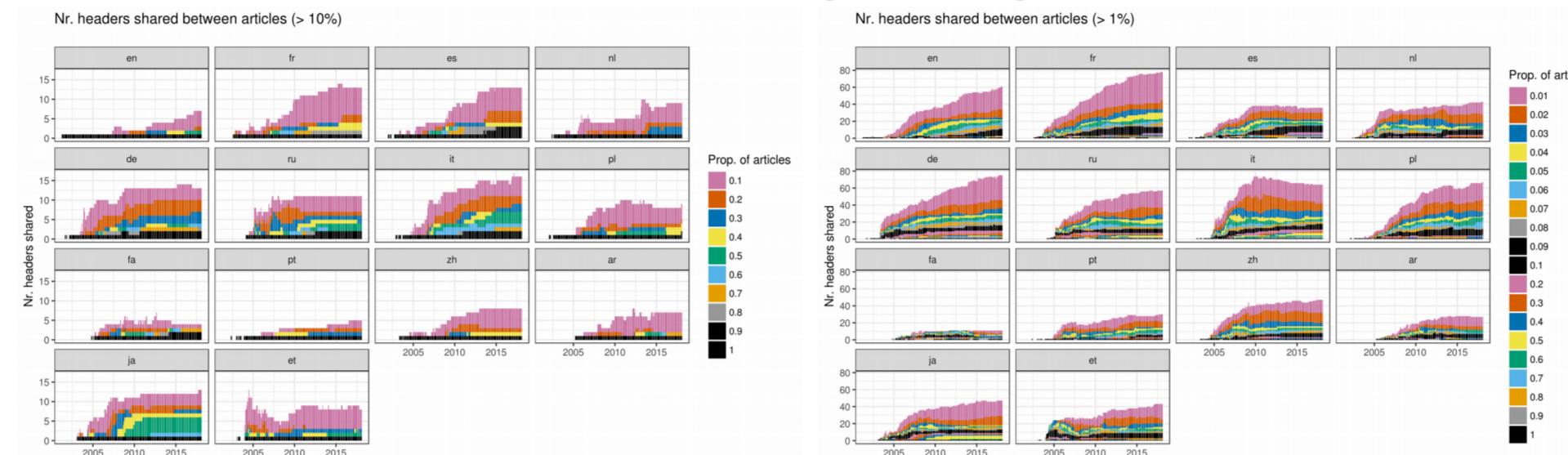
Contents [hide]	
1	Taxonomy
1.1	Classification
1.2	Etyymology
1.3	Subspecies
2	Description
2.1	Pathology
2.2	Genomics
3	Ecology
3.1	Diet
3.2	Predators
4	Behavior
4.1	Reproduction
5	Uses and human interaction
5.1	Early references
5.2	Western discovery
5.3	Panda diplomacy
5.4	Biofuel
5.5	Conservation
5.6	In zoos
5.7	Population chart
5.8	Reference in medicine
5.9	In cryptozoology
6	See also
7	References
8	External links



## Increase in structure

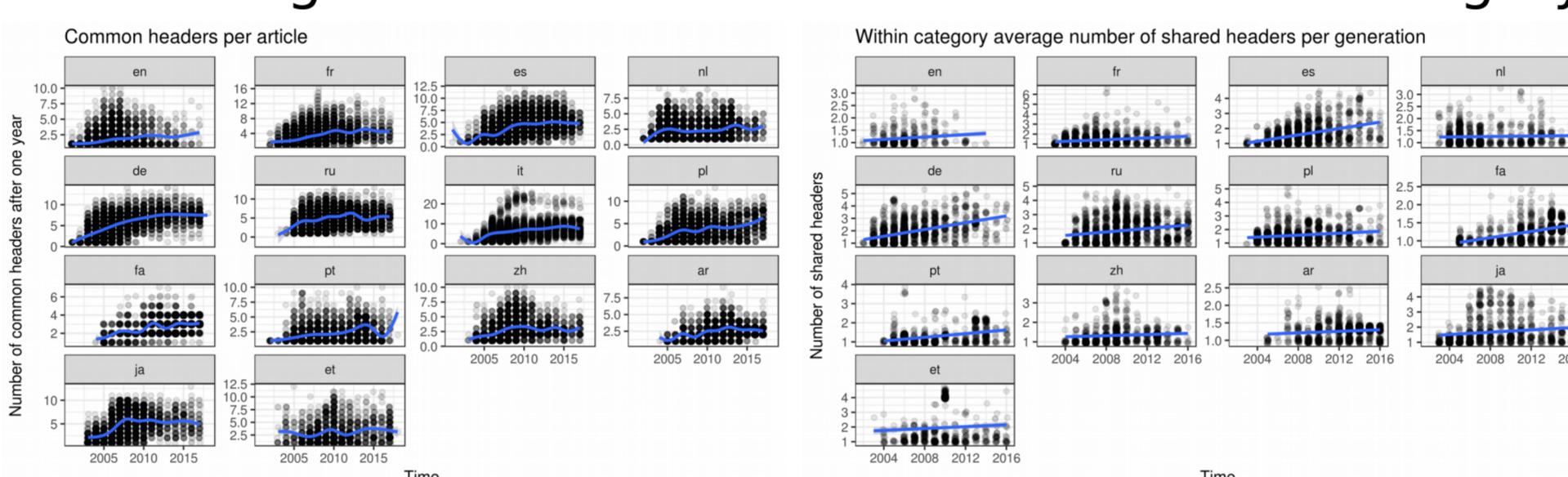
Standardization of headers:

- Over time, the number of headers shared between the articles in our set is growing.



"Generations" of articles from each year

- Popular headers (> 1% articles) after 365 days
- Number of headers per article
- Average number of headers shared within category

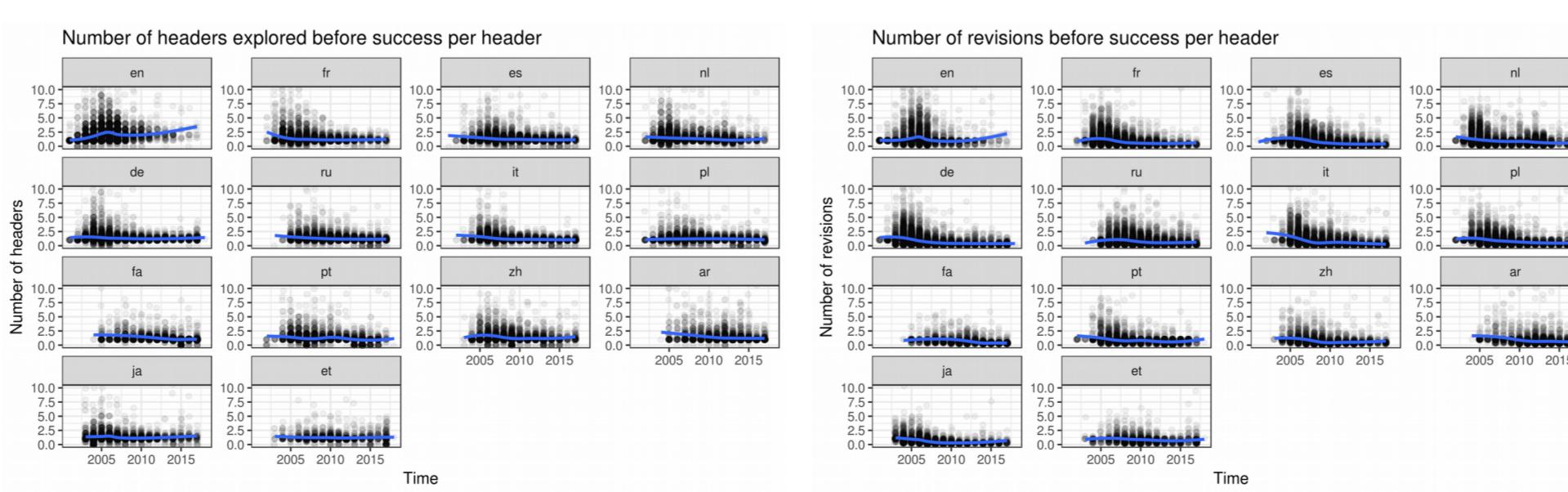


Are these conventions?

- We can check if particular headers in new articles are based on earlier prevalence in the set.
- Different sets can be made: e.g. accumulated frequencies, frequencies weighted by pageviews, frequencies boosted by prestige, additions the "generation" before, etc.
- This can be done within categories (here: *mammals*, *birds*, *snakes* etc) or across the dataset. See more in top right.

## Search space

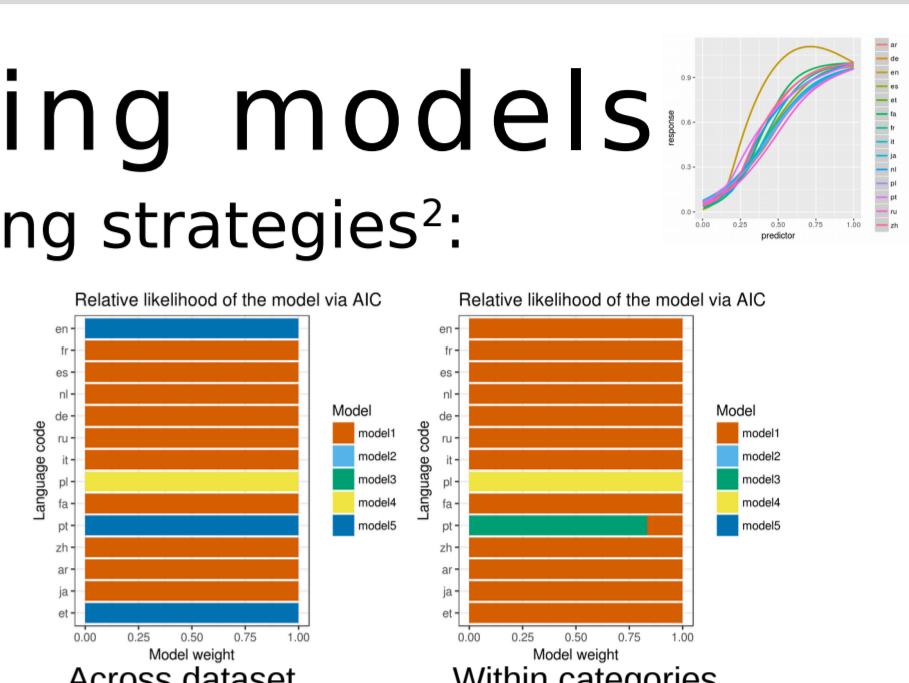
Are newer articles quicker to find the "right" headers?



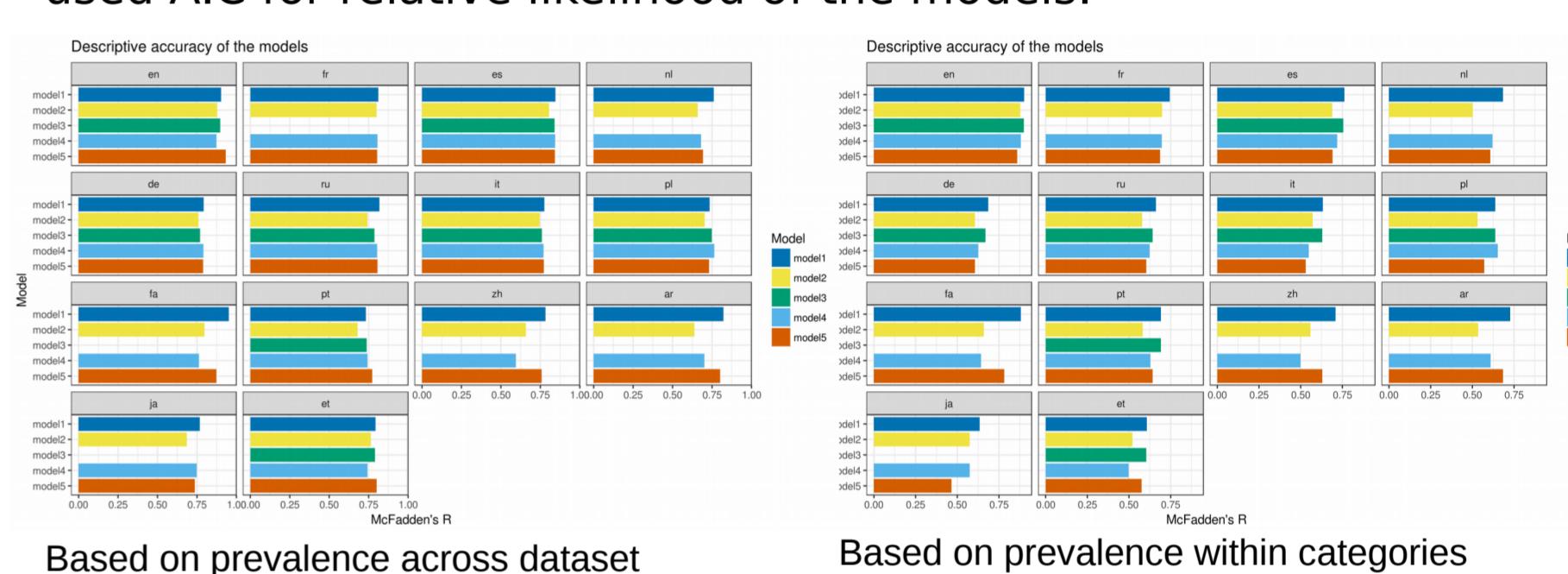
## Social learning models

We contrasted plausible learning strategies<sup>2</sup>:

- Simple replication
- Weighted by popularity
- Weighted by prestige
- Recent additions
- Average



We fit logistic regression models with different learning sets, and used AIC for relative likelihood of the models.



## Discussion

- Wikipedias are getting more structured, however with notable language differences.
- Digital data offer rich resources to study social learning mechanisms<sup>3</sup>, however often they may be hard to distinguish.
- Here, simple replication seems fit best, but others are close by.

Future aims:

- Better ways to distinguish learning trajectories when they make similar predictions.
- Explanations for differences between Wikipedias: Community size? Cultural differences?
- A closer look at translational equivalence of articles.

## References

- Latour, B. (1993). *We have never been modern*. Harvard University Press
- Kendal, R. L., Boogert N.J., Rendell, L., Laland, K.N., Webster, M., Jones P.L. (In press). Social Learning Strategies: Bridge-Building between Fields. *Trends in Cognitive Sciences*.
- Acerbi, A. (2016). A Cultural Evolution Approach to Digital Media. *Front. Hum. Neurosci.*