

Taxonomy Corpus Analysis: Preliminary Data

Workshop on Taxonomic Disorder, 6/12/2022

Charles H. Pence

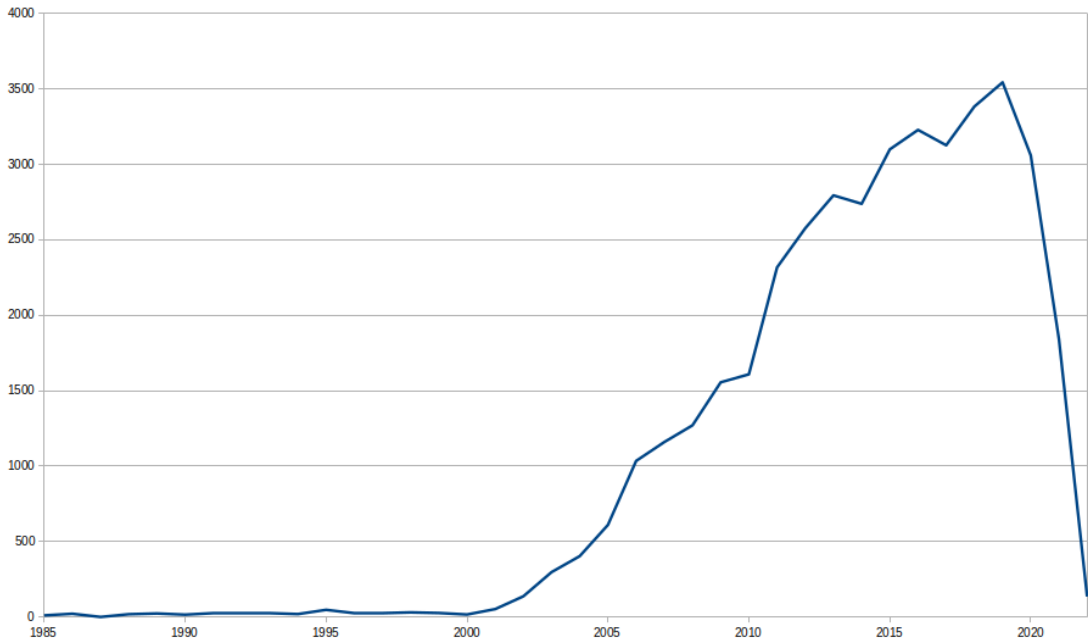
 @pence@scholar.social

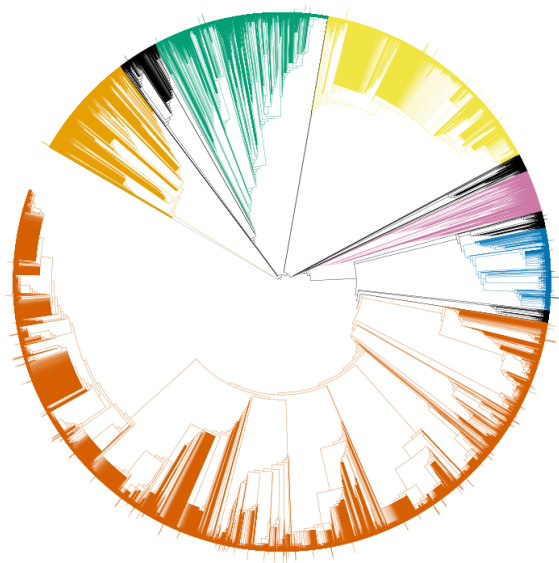
 **UCLouvain**

Institut supérieur de philosophie (ISP)

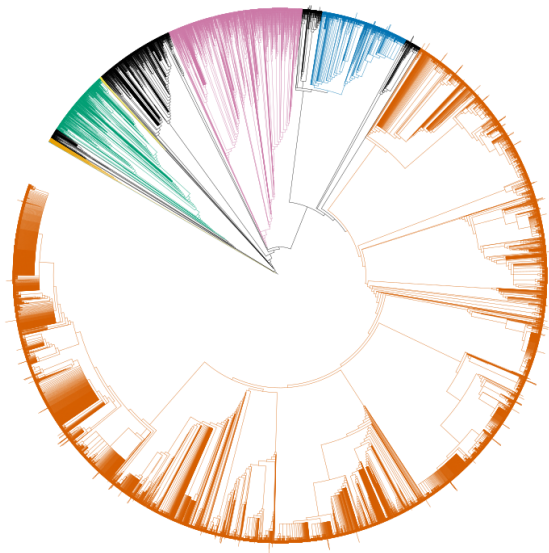


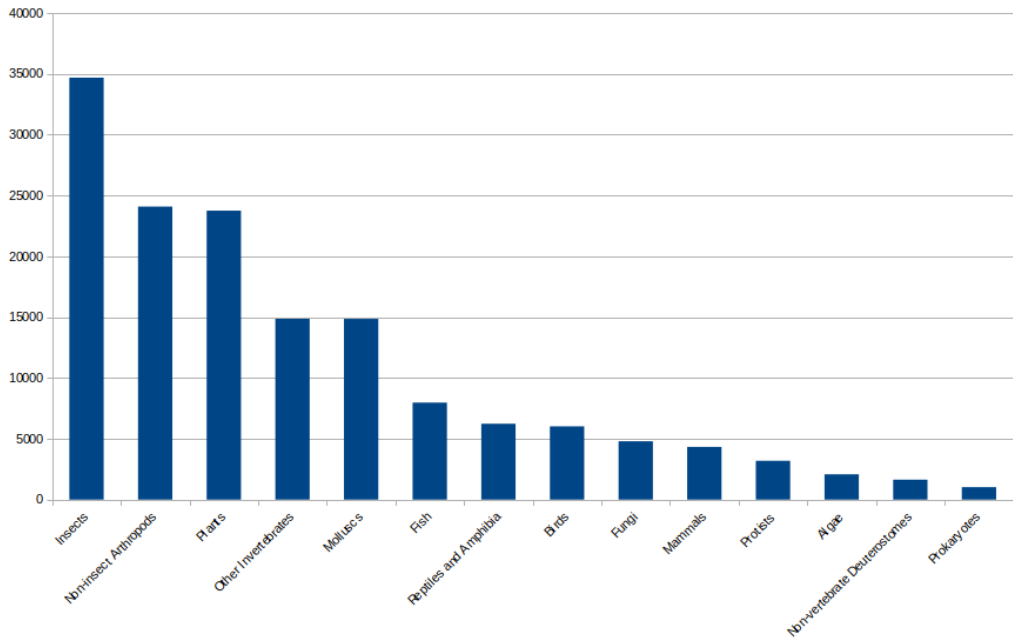
Journal	Publisher	Size
<i>Zootaxa</i>	Magnolia Press	31,348
<i>ZooKeys</i>	Pensoft	4,940
<i>PhytoKeys</i>	Pensoft	820
<i>Journal of Hymenoptera Research</i>	Pensoft	382
<i>MycoKeys</i>	Pensoft	315
<i>Zoosystematics and Evolution</i>	Pensoft	153
<i>Insecta Mundi</i>	Center for Systematic Entomology	1,367
<i>European Journal of Taxonomy</i>	Museum National d'Histoire Naturelle	1,105





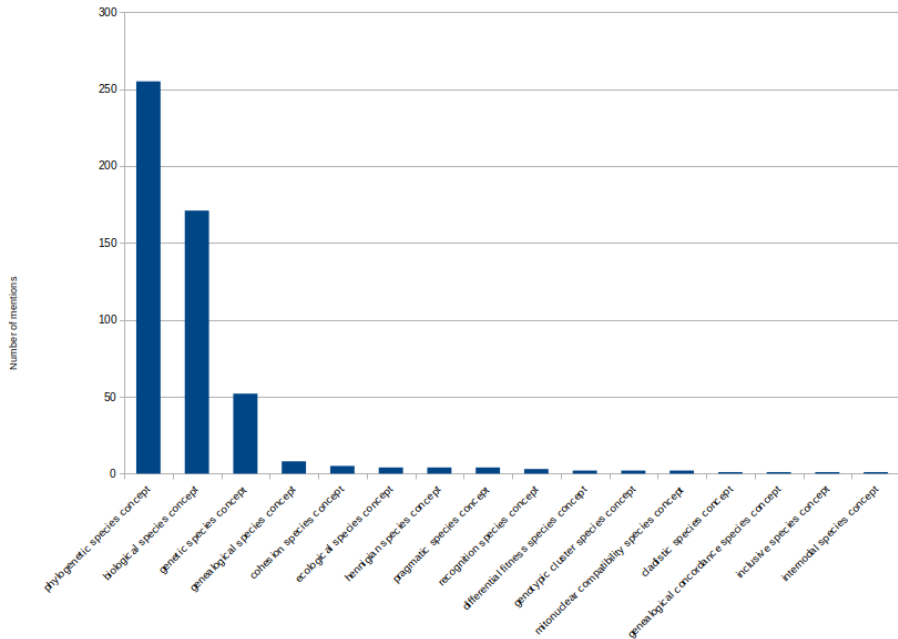
Complete Open Tree of Life





Phylo-Phenetic Species Concept
Phylogenetic Species Concept
Genic Species Concept
Cohesion Species Concept
Genealogical Concordance Species
Concept
Genotypic Cluster Species Concept
Genetic Species Concept
Ecological Species Concept
Recognition Species Concept
Genealogical Species Concept

Biological Species Concept
Differential Fitness Species Concept
Compilospecies Concept
Cladistic Species Concept
Hennigian Species Concept
Internodal Species Concept
Mitonuclear Compatibility Species Concept
Pragmatic Species Concept
Inclusive Species Concept
Biosimilarity Species Concept



Taxa × Species Concepts

The top 5% of non-zero proportions of documents mentioning both the given taxon and the given species concept:

- **Phylogenetic Species Concept:** Protists, Mammals, Birds, Reptiles & Amphibia, Fish
- **Biological Species Concept:** Mammals, Birds

Topic Modeling

Briefly: a kind of unsupervised dimensionality reduction that you can run on a corpus of text. Take documents, normally locations in a 172M-dimensional space (number of word types), and reduce that to 125-D.

(go to live graph:)

<https://cpence.codeberg.page/taxonomy-analyses/>



Interpreting a Topic

Topic 16: popular in mammals

- 0.027*`"colombia"`
- 0.016*`"specie"`
- 0.013*`"type"`
- 0.013*`"peru"`
- 0.010*`"locality"`
- 0.010*`"venezuela"`
- 0.010*`"ecuador"`
- 0.009*`"panama"`
- 0.008*`"distribution"`
- 0.007*`"brazil"`
- 0.007*`"key"`
- 0.006*`"rica"`
- 0.006*`"del"`
- 0.006*`"costa"`
- 0.006*`"genus"`
- 0.006*`"male"`
- 0.006*`"america"`
- 0.006*`"san"`
- 0.006*`"neotropical"`
- 0.005*`"cat"`

Interpreting a Topic

Topic 16: popular in mammals

- 0.027*`"colombia"`
- 0.016*`"specie"`
- 0.013*`"type"`
- 0.013*`"peru"`
- 0.010*`"locality"`
- 0.010*`"venezuela"`
- 0.010*`"ecuador"`
- 0.009*`"panama"`
- 0.008*`"distribution"`
- 0.007*`"brazil"`
- 0.007*`"key"`
- 0.006*`"rica"`
- 0.006*`"del"`
- 0.006*`"costa"`
- 0.006*`"genus"`
- 0.006*`"male"`
- 0.006*`"america"`
- 0.006*`"san"`
- 0.006*`"neotropical"`
- 0.005*`"cat"`

Okay: Central and South American collection sites

Topic 31:

- 0.016* "male"
- 0.016* "genitalia"
- 0.013* "specie"
- 0.009* "female"
- 0.009* "fig"
- 0.008* "brown"
- 0.008* "lepidoptera"
- 0.007* "scale"
- 0.007* "long"
- 0.006* "slide"
- 0.006* "white"
- 0.006* "line"
- 0.006* "new"
- 0.006* "bursae"
- 0.006* "short"
- 0.005* "dark"
- 0.005* "coll"
- 0.005* "forewing"
- 0.005* "holotype"
- 0.005* "leg"

Cautious hypothesis: Lepidopteran anatomy, especially reproductive

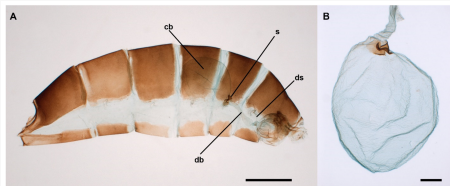
Interpreting a Topic

But wait.

Our lepidopteran reproductive anatomy topic is unusually significant in one group... **in papers that mention molluscs.**

...what?

One Hypothesis



bursa copulatrix, *Leptophobia aripa*



genus *Bursa*, *Bursa granularis*

Some Cool Topics

Topic 9: traditional specimen collection terms

- 0.029* "specie"
- 0.012* "forest"
- 0.012* "habitat"
- 0.010* "area"
- 0.008* "find"
- 0.007* "collect"
- 0.007* "site"
- 0.007* "study"
- 0.007* "record"
- 0.006* "population"
- 0.006* "range"
- 0.006* "high"
- 0.005* "specimen"
- 0.005* "occur"
- 0.005* "know"
- 0.004* "individual"
- 0.004* "region"
- 0.004* "number"
- 0.004* "sample"
- 0.004* "distribution"

Popular in every taxon **except** non-insect arthropods, fish, and fungi.

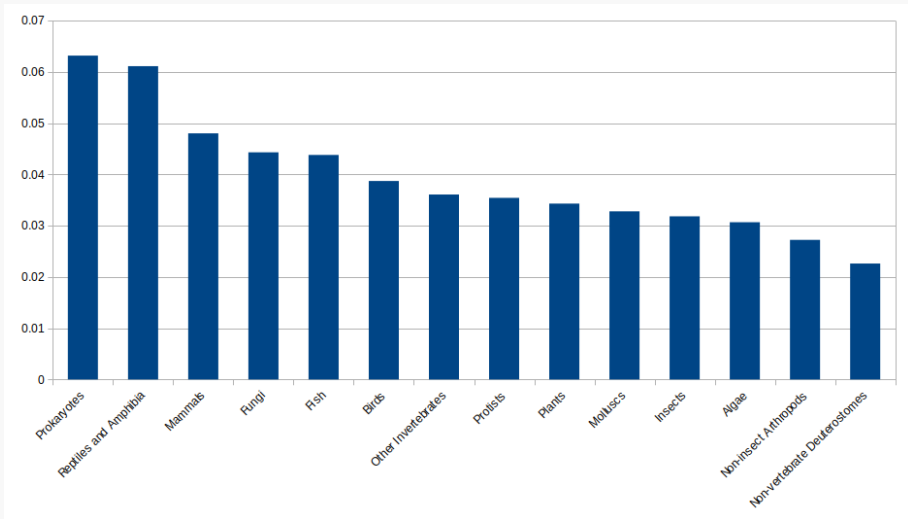
Some Cool Topics

Topic 64: molecular phylogenetics

- 0.021*`"specie"`
- 0.017*`"sequence"`
- 0.016*`"analysis"`
- 0.011*`"molecular"`
- 0.010*`"dna"`
- 0.008*`"phylogenetic"`
- 0.007*`"tree"`
- 0.007*`"clade"`
- 0.007*`"gene"`
- 0.007*`"specimen"`
- 0.007*`"study"`
- 0.007*`"morphological"`
- 0.006*`"support"`
- 0.006*`"group"`
- 0.006*`"genetic"`
- 0.006*`"coi"`
- 0.006*`"datum"`
- 0.006*`"base"`
- 0.005*`"table"`
- 0.005*`"population"`

Among the **top-20 most significant probabilities** in reptiles and amphibia, birds, fish, fungi, and mammals; top-5% in every other group

Molecular Methods



How about disagreement?

Ask the model: what topics are likely to pick words like “disagree,” “disagreement,” “dispute,” “revision,” or “revise”?

For disagreement terms: topics 120 and 43

For revision terms: also 43, as well as 108

What are those topics about?

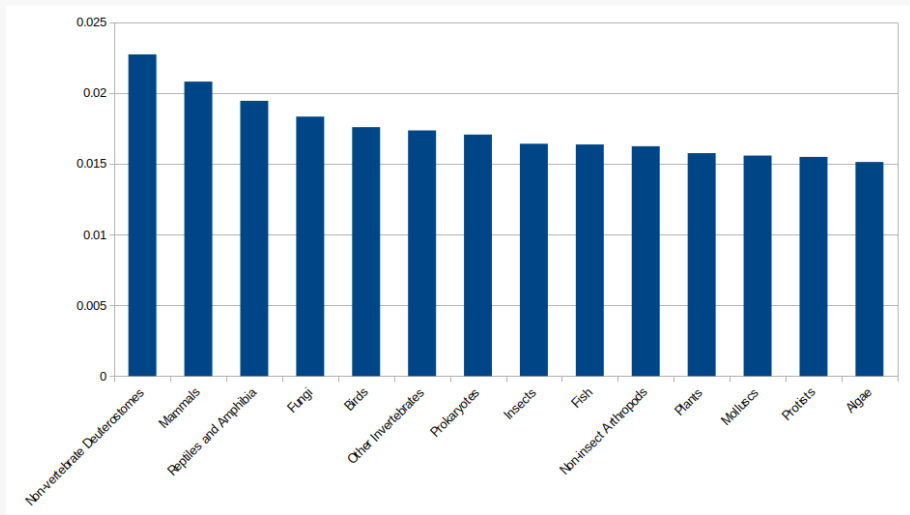
Topic 108 (revision)

- 0.015*"specie"
- 0.009*"limb"
- 0.009*"cladocera"
- 0.007*"alona"
- 0.007*"sinev"
- 0.006*"lake"
- 0.006*"solanum"
- 0.006*"group"
- 0.006*"kotov"
- 0.005*"record"
- 0.005*"van"
- 0.005*"postabdomen"
- 0.005*"portion"
- 0.005*"mol"
- 0.005*"fig"
- 0.005*"distribution"
- 0.005*"genus"
- 0.005*"bor"
- 0.004*"head"
- 0.004*"parthenogenetic"

Topic 120 (disagreement)

- 0.018* "character"
- 0.013* "genera"
- 0.011* "taxon"
- 0.011* "group"
- 0.010* "specie"
- 0.010* "genus"
- 0.009* "phylogenetic"
- 0.008* "include"
- 0.007* "analysis"
- 0.007* "family"
- 0.007* "relationship"
- 0.005* "phylogeny"
- 0.005* "clade"
- 0.005* "morphological"
- 0.005* "classification"
- 0.005* "support"
- 0.005* "press"
- 0.005* "new"
- 0.005* "consider"
- 0.004* "present"

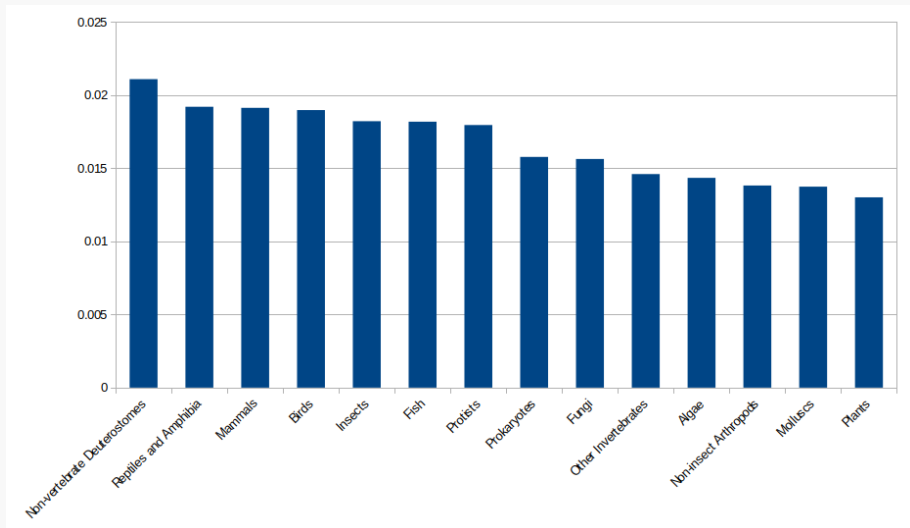
Topic 120 (disagreement — phylogeny)



Topic 43 (disagreement)

- 0.015*"specie"
- 0.011*"name"
- 0.010*"description"
- 0.010*"new"
- 0.008*"publish"
- 0.007*"author"
- 0.007*"nomenclature"
- 0.007*"code"
- 0.007*"publication"
- 0.006*"type"
- 0.006*"article"
- 0.006*"zoological"
- 0.006*"original"
- 0.006*"synonym"
- 0.006*"work"
- 0.006*"list"
- 0.006*"valid"
- 0.005*"international"
- 0.005*"available"
- 0.005*"note"

Topic 43 (disagreement — “new species”)



Coming Soon

Geocoding: how do taxa, topics, and species concepts correlate with mentions of geographic locations?

Help us brainstorm!

charles@charlespence.net

<https://pencelab.be>

 @pence@scholar.social