



Identifying, cross-referencing, and extracting dark data using GeoDeepDive



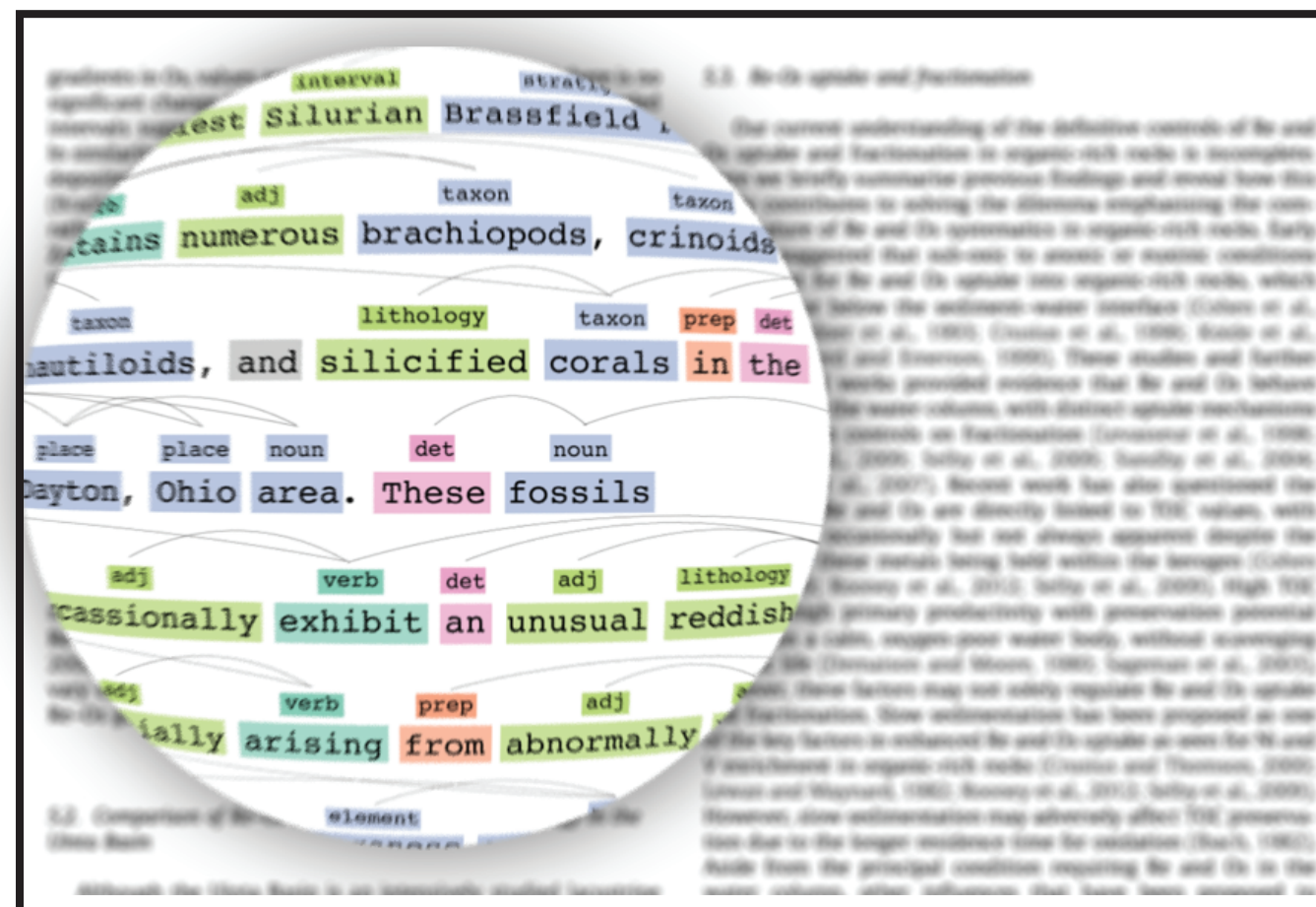
Erika T. Ito, Andrew A. Zaffos, Valerie J. Syverson, Ian A. Ross, and Shanan E. Peters

What is GeoDeepDive?

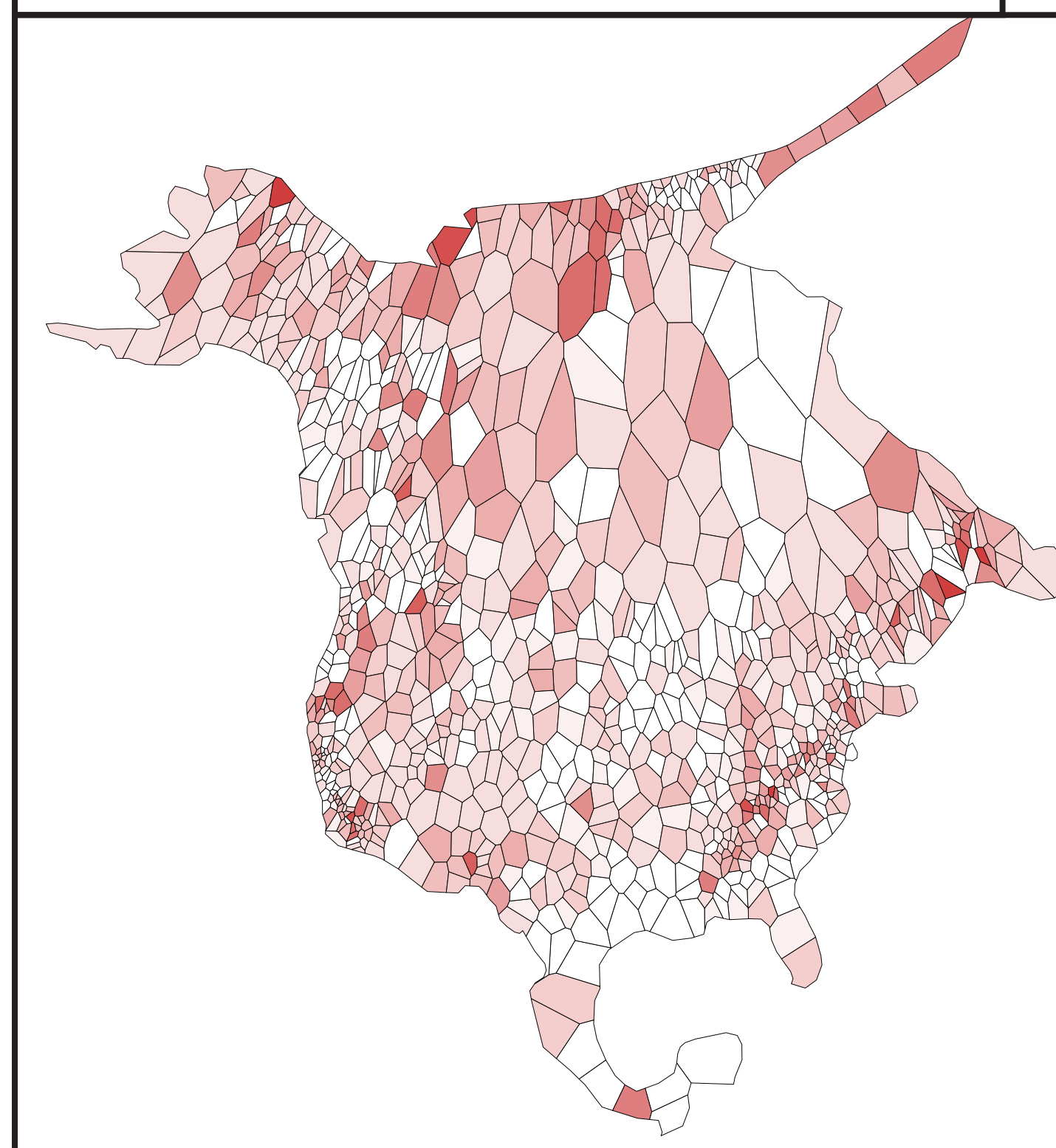
GeoDeepDive is a library of over 3 million machine-readable scientific documents (e.g., journal articles, monographs). It downloads PDFs from partnered publishers (e.g., Elsevier, Wiley, Taylor & Francis), and turns them into machine-readable products using natural language processing, optical character recognition, and other techniques. This processing enables intelligent text-mining and grammatical analysis of scientific documents.

Identifying Dark Data

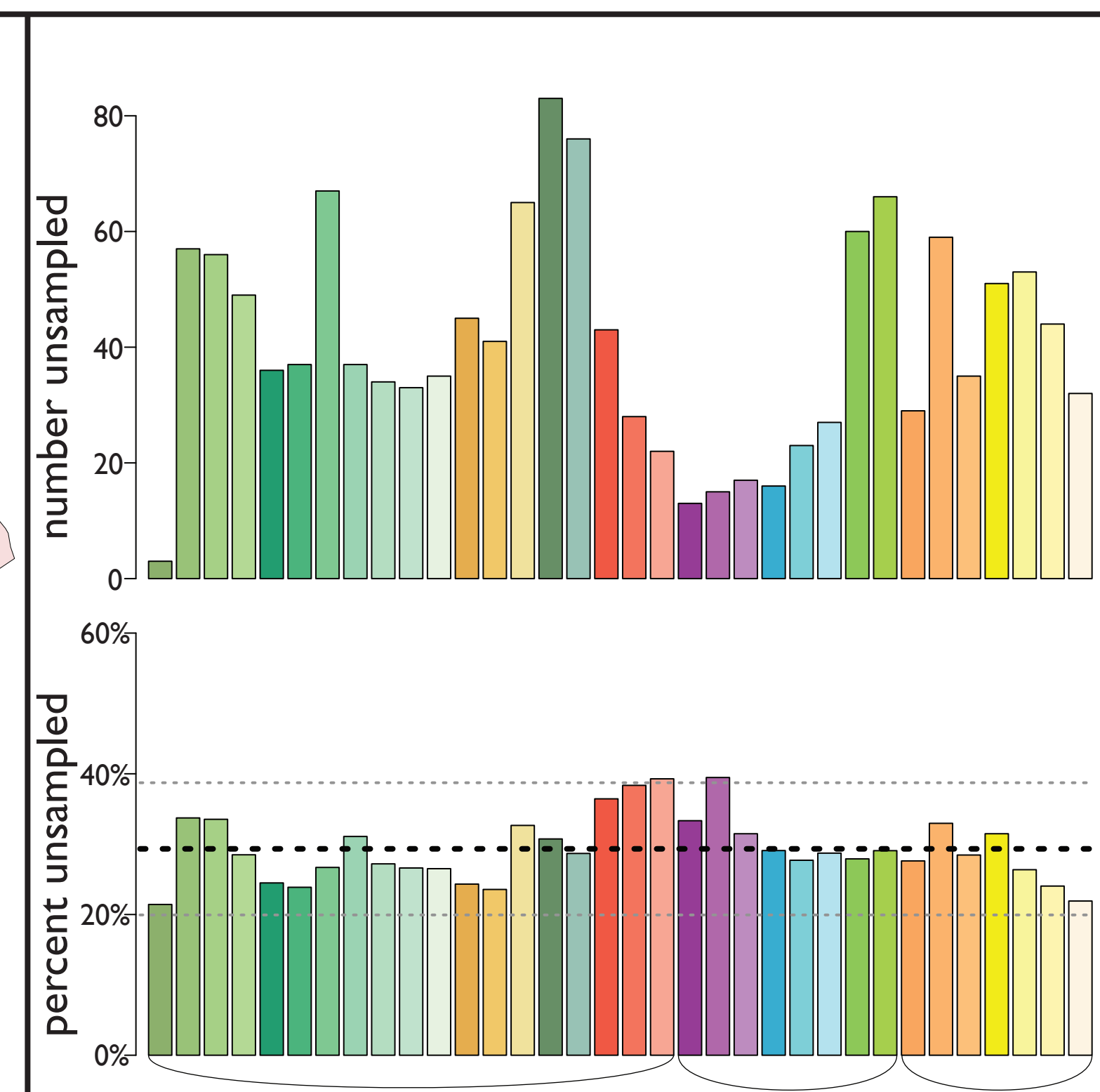
The Paleobiology Database is the largest public repository of global fossil occurrences. However, there are ongoing concerns that certain time-periods, geographic regions, or paleo-environments are systematically understudied. How can we identify understudied areas and test for systematic biases?



1. Identify formal geologic formations not described in the Paleobiology Database.
2. Search documents for the co-occurrence of these formations and keywords indicating they are fossil bearing - e.g., "fossiliferous".
3. Create rules to sift out bad matches, such as using geographic location to verify the correct formation is being referenced.



A map of spatial sampling coverage in the Paleobiology Database relative to the scientific literature.



A barplot of temporal sampling coverage in the Paleobiology Database relative to the scientific literature.

Acknowledgements

Research funded by NSF ICER I343760 and EAR-1150082

GeoDeepDive is made possible by our partnered publishers: Elsevier, Wiley, Taylor & Francis, SEPM, Canadian Science Publishing, PubMed Central, the United States Geological Survey, the Geological Society of America, the Public Library of Science, and the Treatise of Invertebrate Paleontology.

Cross-Referencing Dark Data

```

type: article
title: "New peleciniid wasps (Hymenoptera: Peleciniidae) from Upper Cretaceous Myanmar Amber"
volume: "67"
journal: "Cretaceous Research"
link: [
  {url: "http://www.sciencedirect.com/science/article/pii/S0195667116301409", type: "publisher"},
  {url: "https://data.mendeley.com/datasets/50195667116301409", type: "filepath"},
  {url: "https://doi.org/10.1016/j.cretres.2016.07.003", type: "doi"},
  {url: "https://www.sciencedirect.com/science/article/pii/S0195667116301409", type: "publisher"},
  {url: "https://doi.org/10.1016/j.cretres.2016.07.003", type: "doi"}
]
publisher: "Elsevier"
author: [
  {name: "Guo, Lichao"},
  {name: "Shih, Chungkuei"},
  {name: "Li, Longfeng"},
  {name: "Ren, Dong"}
]
pages: "84-90"
number: "1"
identifier: [
  {type: "doi", id: "10.1016/j.cretres.2016.07.003"}
]
year: "2016"

```

What if we want to link scientific literature to occurrences in online databases like the Paleobiology Database? Once accomplished, we can enhance existing database entries with additional information extracted from the primary literature.

1. Download bibliographic metadata of target articles matching some initial criterion. For example, all articles mentioning a formal taxonomic name.
2. Build a predictive model and extract assumed model predictors from both the article metadata and the Paleobiology Database - e.g., author, publication year, publication name, article title, or taxonomic names.
3. Fit the model to a training set of known correct and false matches to constrain the correct model coefficients.
4. Use the model to predict the likelihood that a fossil occurrence in the Paleobiology Database refers to a specific article in the GeoDeepDive library.

Classic

Burmese amber (CNU coll) (Cretaceous of Myanmar)
Where: Myanmar (26.4° N, 96.7° E; paleocoordinates 12.4° N, 93.8° E)
Coordinates: Myanmar
When: Early-Lower Cretaceous (95.6 - 93.5 Ma)
Environment/Biology: Insectal, amber
Preservation: amber
Collection methods: Repository: Key Lab of Insect Evolution and Environmental Changes, College of Life Sciences, Capital Normal University, Beijing, China
Primary reference: F. Dong, C. K. Shi, and D. Ren. 2015. A new genus of Tanyderidae (Insecta: Diptera) from Myanmar amber, Upper Cretaceous. *Cretaceous Research* 64:299-302. doi:10.1016/j.cretres.2015.07.004
Purpose of describing collection: taxonomic analysis
PaleoDB collection 165681: authorized by Matthew Clapham, entered by Matthew Clapham on 22.01.2015
Creative Commons license: CC BY (attribution)

Taxonomic list
Show authors, comments, and common names
Insecta
Diptera - Tanyderidae - *Diastylus carolinensis*
Phanerozoa - Insecta - Tanyderidae - *Neochorebus sushimae*
Hymenoptera - Ichneumonidae - *Neochorebus longus*, *Calochorebus peruanus*
Hymenoptera - Ichneumonidae - *Calochorebus peruanus*
Hymenoptera - Pelecinidae - *Lagenolepistus*
Hymenoptera - Pelecinidae - *Brachypelecinus euthyntus*, *Abropelecinus tythus*, *Zoropeclecinus periosus*
Hymenoptera - Pelecinidae - *Brachypelecinus euthyntus*
Hymenoptera - Pelecinidae - *Brachypelecinus euthyntus*
Coleoptera - Dermestidae - *Megastoma agrippa*, *Cretodermestes palpatii*, *Attagenus sp.*, *Attagenus securus*
Coleoptera - Carabidae - *Archaeobambolus isabellae*
Coleoptera - Tenebrionidae - *Theobambolus cretaceus*
Coleoptera - Curculionidae - *Brachycolus thysanus*
Diptera - Limoniidae - *Dianopomyia pilosula*
Diptera - Sepsidae - *Sepsisomyia*
Neuroptera - Megaloptera - *Megaloptera reynoldsii*, *Megaloptera reynoldsii*

The first fossil of Lindsaeaceae (Polypodiales) from the Cretaceous amber forest of Myanmar [doi: 10.1016/j.cretres.2016.12.003]

Collembola (Arthropoda, Hexapoda) from the mid Cretaceous of Myanmar (Burma) [doi: 10.1016/j.cretres.2005.07.003]

New beaded lacewings (Neuroptera: Berothidae) from Upper Cretaceous Myanmar amber [doi: 10.1016/j.cretres.2016.08.007]

The earliest Attagenus species (Coleoptera: Dermestidae: Attageninae) from Upper Cretaceous Burmese amber [doi: 10.1016/0022-2011(16)90160-2]

A new earwig (Dermaptera: Pygidicrididae) from the Upper Cretaceous Myanmar amber [doi: 10.1016/j.cretres.2017.02.012]

The first record of Ichneumonidae (Insecta: Hymenoptera) from the Upper Cretaceous of Myanmar [doi:10.1016/j.cretres.2016.11.001]

Suggest papers with similar attributes in GeoDeepDive.

A new genus with a new species, *Brachypelecinus euthyntus* gen. et sp. nov., and two new species, *Abropelecinus tythus* sp. nov. and *Zoropeclecinus periosus* sp. nov., are described and figured from three exquisitely preserved peleciniid wasps in the Upper Cretaceous Myanmar (Burma) amber.

In this paper, we report a new genus and species, *Brachypelecinus euthyntus* gen. et sp. nov., and two new species, *Abropelecinus tythus* sp. nov. and *Zoropeclecinus periosus* sp. nov., based on three Myanmar amber specimens.

Abropelecinus tythus sp. nov., holotype, CNU-HYM-MA2016002.

The forewing of *Abropelecinus tythus* sp. nov. (male, body length 4.4 mm) (Fig. 5B) has hyaline membrane with only two tubular veins present (C and R).

For example, the body lengths of amber species range from 4.4 mm in a male specimen of *Abropelecinus tythus* sp. nov. to 12.7 mm in a female specimen *Zoropeclecinus zigrasi* Engel & Grimaldi, 2013, while the body lengths of compression fossils range from 3.8 mm in a female specimen *Epelecinus minutus* Zhang & Rasnitsyn, 2004 to 50.9 mm in a female specimen *Megapelecinus changi*.

A. Brachypelecinus euthyntus gen. et sp. nov. B. *Abropelecinus tythus* sp. nov. C. *Zoropeclecinus periosus* sp. nov. D. *Zoropeclecinus zigrasi* sp. nov. E. *Pelecinopteron tubuliforme*.

The forewing of *Abropelecinus tythus* sp. nov. with hyaline membrane has only two tubular veins present (C and R), which is the first amber male reported with this type of reduced venation.

Present sentences mentioning each taxon in the text.

Paleobiology Database
revealing the history of life

GeoDeepDive
a digital library for science

(A)

(B)

(C)

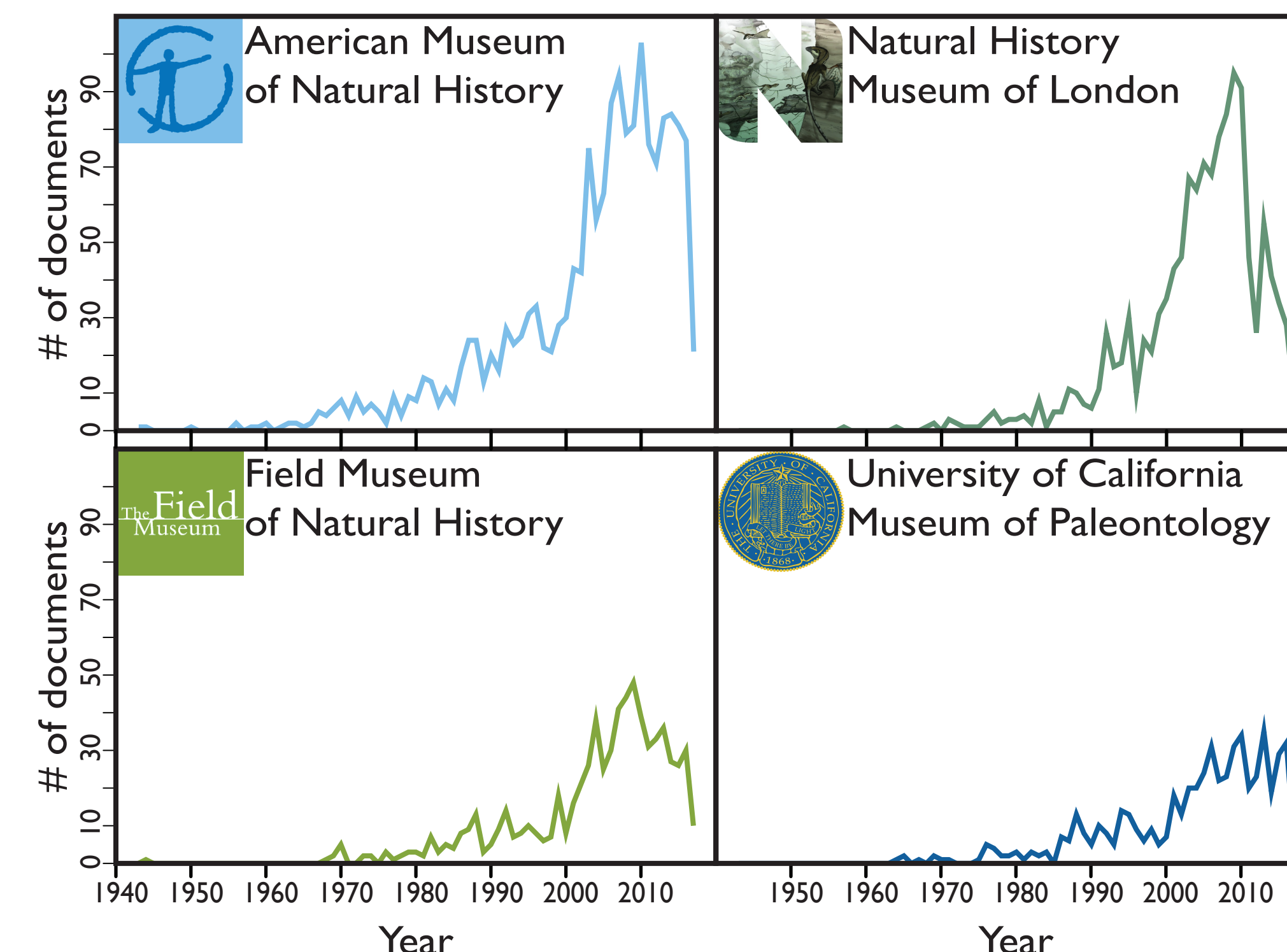
(D)

Display figures mentioning a taxon in the text.

Extracting Dark Data

What if we want to build an entirely new database by extracting information from the literature? We illustrate how to build a database of museum specimens. Such a database allows us to measure the scientific impact of different institutions and specimens.

1. Begin with a list of institution names and codes.
2. Index articles containing codes or institution names.
3. Index articles containing taxonomic names.
4. Find the intersect of institution, taxonomic names, and numeric digits within a sentence and evaluate the probability it is a museum specimen code.



specimen code	#docs
CM 3018	(30)
UCMP 37302	(27)
CM 11338	(26)
FMNH PR2081	(25)
AMNH 7224	(19)
USNM 3529	(19)
USNM 4735	(18)
AMNH 5356	(17)
AMNH 5337	(17)
AMNH 5027	(17)
DMNH 12679	(16)
SMNS 13200	(16)
ROM 1215	(16)
AMNH 5214	(16)
UCMP 77270	(15)
YPM 1883	(15)
BMNH R1111	(15)
YPM 1901	(14)
ROM 873	(14)
AMNH 6810	(14)
AMNH 24450	(14)

Holotype of *Apatosaurus louisae*

Allende meteorite

"Ilex" serrata specimen used on cover of *Cretaceous Research*

Holotype of *Anchisaurus polyzelus*

Complete skull of *Araripesuchus gomesii*

"Sue" (*Tyrannosaurus rex*)

Almost all of the 20 most-cited individual specimens are dinosaurs (listed in green).

Outcomes

1. A machine executed assessment of how thoroughly scientific literature is collated by paleontologists and a high-priority list for future data acquisition.
2. A join table linking different databases or database entries to the primary scientific literature.
3. An evaluation of the scientific impact of different museum specimens and scientific institutions.

www.geodeepdive.org