Finding the where of scholarly publishing: automating scientific article geolocation

BRAZIL

Jason W. Karl, Ph.D. Connor Bryant College of Natural Resources University of Idaho



PER

MEXICO



ALGERIA

NIGER

MAL

What happens if you want to search for research from a specific place?



<u>https://journalmap.org</u> - Geographic-based literature searching



SEARCH COLLECTIONS PUBLISHERS HELP API DASHBOARD





Challenge: Extracting location information from published studies

>90% of ecology/agriculture studies report a study location, but use countless ways to do it*

<50% of ecology/ studies report geographic
coordinates for their study area</pre>

tats within the Craters of the Moon National Preserve in southern Idaho, USA (Fig. 1). The Craters of the Moon area was created during at least eight volcanic eruptive periods that blanketed portions of the Snake River Plain with basaltic lava flows between 15,000 and 2,000 years ago (Owen 2008). Within the area of the lava flows, some areas were not covered by the lava and formed kipukas-islands of vegetation surrounded by lava. Two of these kipukas, Little Park (3,150 ha) and Laidlaw Park (35,040 ha) were sampled for this study.

2. Study area

For this study we considered the Bureau of Land Management's (BLM) 97 308-ba Wildborse grazing allotment in southern Idaho (Fig. 1, 43.028°N, 113.864°W). The majority of the study area is in public ownership with the Bureau of Land Management (BLM) being the largest single land steward-managing approximately 93,317 ha (95.8%) of the study area. Approximately 1305 ha (1.3%) of the study area is in private ownership, and 2843 ha (2.9%) managed by the state of Idaho. The dominant land use in the study area is cattle and sheep grazing.

STUDY AREA

This study was conducted in the West Central Sage-grouse planning area in west-central Idaho, which included parts of Washington, Adams, Gem, and Payette counties (centered at 44°26'N, 116°38'W) The 374,700 ha planning area was established to conserve a small and isolated population of sage-grouse that was considered the most likely to be extirpated within the state (Idaho Sage-grouse Advisory Committee 2006). Exact population numbers are unknown but it was estimated that the population was significantly

Study area	Site	Geographic coordinates	Average elevation, m	Average annual precipitation,1 cm	Photo acquisition dates
Idaho	1	Lat 43°12'35.9'N,	1 227	28	27 August 2008
		long 116°44'15.6'W			
	2	Lat 43"6"32.1"N,	1 629	69	27 August 2008
		long 116°46'33.8'W			
	3	Lat 43°3'58.1"N,	2 082	89	27 August 2008
		long 116°45'23.2"W			
Nevada	1	Lat 36'21'45.4"N,	366	15	17 March 2009
		long 114°25'1.9'W			
	2	Lat 36°22'12.5"N,	470	16	17 March 2009
		long 114°26'50.9'W			
	3	Lat 35°17'50.7"N,	1 601	24	17 March 2009
		long 115°33'14.0'W			
New Mexico	1	Lat 32°34'37.0"N,	1 230	26	23 August 2008
		long 106°0'42.0'W			
	2	Lat 32°29'38.8"N,	1 641	48	23 August 2008
		long 105°40'59.7"W			
	3	Lat 32°22'33.8"N,	1 440	37	23 August 2008
		long 105°39'25.6'W			

*Examples from papers I have published. I'm part of the problem!

Article Geotagging ## Regular expression of parsing bounding box coordinates in form of ## Latitude #1 to Latitude #2, Longitude #1 to Longitude #2

test arr.push line



- Locations can be extracted from article texts*)? ### prepend stuff not necessary
- Pattern matching for coordinate values^{#1 direction (before the coord)}
- Place-name recognition is challenging the coord value)
 - Studies have many irrelevant place names

(([L|l]atitude|[L|l]at(\.)?)\s*)? ### more fluff
(?<lat2dir>N|S|[Nn]orth|[Ss]outh)?\s* ### Latitude #2 direction (before the coord)

- The Challenge: develop approaches to identify and extract locations from journal articles based on place names.
 - Irrelevant place names \s*(,|;|to|and|,\sand),?\s*((between|from)\s*)? ### Latitude\Longitude separator
 - Uncommon/non-standard place names ((IL))ongitude(s)?[[L]]ong(\.)?]\s*)?((between|from)\s*)? ### more fluff
 - Things that lookalike place names Longitude #1 direction (before the coord value)
 •?E.g.1deGanada.Goose[0-9]{1,10})?`?\s*(?<lon1min>[0-9]{1,2}(\.[0-9]{1,10})?('|'|'))?\s*(?<lon1sec>[0-9]{1,2}(\.[0-9]{1,10})?(''|''|''))?\s* ### Longitude #1 (\q<lon1dir>)?\s* ### Longitude #1 direction (after the coord value)

Performance of Natural Language Processing Tools for Article Geotagging





Godfrey, Kenyon, and Karl (in prep)

*Based on comparison to confirmed locations of 120 PLoS One articles.

Evaluating LLM Tools for Article Geotagging



- 3 LLM tools
 - ChatGPT 3.5 (OpenAI)
 - Bard (Google)
 - Claude 2 (Anthropic)
- 120 PLoS One articles
 Confirmed locations
- Prompt development
 - Uploaded PDFs
- Measured indicators
 - # articles returning locations
 - Average distance to confirmed location
 - Repeatability of LLM results

Can we use modern AI tools like ChatGPT to automate article geotagging?



Confirmed locations of 120 PLoS One articles used in Karl (2018) paper

LLM Article Geotagging Results



	% Articles Returning a Location	% Articles Accurately Located	Average Distance from Conf. Location	PDF Effects of Previ 4.86 MB Is the study described in this document specific to a location? Yes, the study described in this paper is specific to a location in the Central Pre-Pyrenees region of Spain. Here are the key points that indicate it is localized:		
Bard	85%	64%	272km	 The abstract mentions the research was conducted in oak forests in the Central Pre-Pyrenees, Spain. The study area section provides details on the location, geology, climate, and vegetation specifics to the Central Pre-Pyrenees region. Figure 1 shows a map with the locations of the 10 formular Central Pre-Pyrenees area. 		
ChatGPT	70%	69%	174km	Where was the study described in this document conducted? Output results as CSV code (a CSV) with name of study area and geographic coordinates for a centroid (Latitude, Longitude) and bounding box (Min Lat, Min Long, Max Lat, Max Long).		
Claude	95%	95%	33km	Study Area,Latitude,Longitude,Min Lat,Min Long,Max Lat,Max Long Central Pre-Pyrenees, Spain,42.2,-0.2,42.1,-0.3,42.3,-0.1		
				Example prompts and responses from Claude 2		

LLM Geotagging Repeatability



	Avg Bounding Width	Avg Bounding Height	% Returning Location	Avg # Success (out of 10)
ChatGPT 3.5	48.72km	147.56km	82%	7.5 (range 1-10)
Bard	639.73km	1800.00km	91%	10
Claude 2 Pro	7.67km	49.73km	100%	10



- LLMs are generative AI tools
 - May produce different results to the same prompt
- Locations requested 10 times for a set of 10 articles
- Generally, variability was minor for repeat requests
- Claude 2 performed the best



Finding the where of scholarly publishing: automating scientific article geolocation

BRAZIL

Jason W. Karl, Ph.D. Connor Bryant College of Natural Resources University of Idaho



PER

MEXICO



ALGERIA

NIGER

MAL