

Retrieval-Augmented Search for Large-Scale Map Collections with ColPali

Jamie Mahowald
 mahowald.jamie@gmail.com
 Personal
 USA

Benjamin Charles Germain Lee
 bcgl@uw.edu
 Information School, University of Washington
 Seattle, Washington, USA

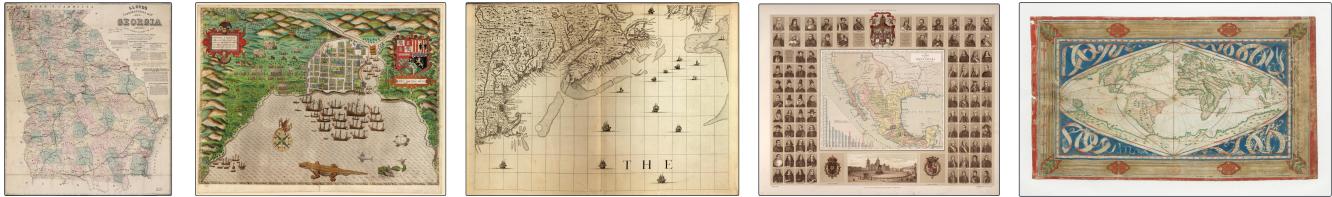


Figure 1: Historical maps from the Library of Congress’s Geography & Maps Division.

Abstract

Multimodal approaches have shown great promise for searching and navigating digital collections held by libraries, archives, and museums. In this paper, we introduce map-ras: a retrieval-augmented search system for historic maps. In addition to introducing our framework, we detail our publicly-hosted demo for searching 101,233 map images held by the Library of Congress. With our system, users can multimodally query the map collection via ColPali, summarize search results using Llama 3.2, and upload their own collections to perform inter-collection search. We articulate potential use cases for archivists, curators, and end-users, as well as future work with our system in both machine learning and the digital humanities. Our demo can be viewed at: <http://www.mapras.com>.

CCS Concepts

- **Information systems → Information retrieval; Digital libraries and archives; Document structure.**

Keywords

Dynamic corpus expansion, search & retrieval, machine learning for visual search, multimodal search, digital libraries, digital humanities

1 Introduction

Maps are key items for institutions that collect and preserve cultural heritage [4]. They reveal historically significant insights into the priorities, capabilities, and limitations of the societies that created them, and so they serve as important subjects of research. The Library of Congress (LOC), the world’s largest map collector, has over 5.5 million maps [18]. Just over 10% are digitized, while a quarter were received before machine-readable cataloging [1].

Searching the LOC’s Geography and Map (G&M) digitized collection amounts to searching through textual metadata or text extracted by optical character recognition (OCR). This method requires thorough and accurate annotation of each searchable document. Moreover, it misses features specific to visual data not captured in the metadata, and it is fixed to the pre-defined hierarchy of

map items. Maps are particularly difficult to characterize because they combine visual and textual elements, posing challenges for traditional, unimodal machine learning solutions.

Recent advances in multimodal machine learning have made this problem more tractable. The contrastive language-image pretraining (CLIP) model class [13] in particular allows a user to compare a text prompt with a pre-computed corpus of image embeddings. By returning the images whose embeddings are most geometrically similar to the prompt embedding (i.e., nearest neighbors search), CLIP can function as a search engine. However, the base CLIP model struggles with text recognition [17], and its hard labeling resists easy fine-tuning on similar-looking training data.

In this work, we apply the ColPali document-retrieval framework [6] to the task of searching maps. Broadly designed as an aid in retrieval-augmented generation (RAG) settings, ColPali has several properties that facilitate this task. Its soft similarity scores can better handle visually similar items, and its multi-vector representations can capture more complex relationships in data. As opposed to models whose embeddings capture only global features, ColPali can attend to specific regions of a document. This framework also makes it easy for users to dynamically upload and add to their corpus, a capability we call retrieval-augmented search (RAS).

In summary, we provide the following contributions:

- (1) We introduce retrieval-augmented search, a variant of RAG, for cultural heritage collections, and we present a demo for 101,233 map images held by the Library of Congress. Our demo can be found at: <http://www.mapras.com>.
- (2) We elaborate on use cases of our map retrieval-augmented search system and detail future directions for our research.
- (3) We release our code at <https://github.com/j-mahowald/mapras>.

2 Related Work

2.1 Library of Congress Maps

As our use case, we use the LOC G&M’s digital collection of 563,698 unique “segments” (unique images) divided among 57,962 “resources”

(composite items). A small sampling can be found in Figure 1. Items contain anywhere from one to 11,981 segments (though most are small, with median 2 segments per item), and metadata is recorded at the item level. These maps span centuries, as well as domains, from geography and topography to historical analysis.

Of the 563,698 images, 439,947 (78% of segments), belong to the collection of Sanborn fire insurance maps, which provide granular data on urban building infrastructure and land use [11]. Because these are easily indexable by metadata, we restrict our tool to the approximately 100,000 non-Sanborn maps in the Library’s digital collection that return valid calls.

The G&M online catalog is broadly searchable via metadata facets including collection, keyword, location, and time period. Each individual catalog entry can be viewed for each map, along with information about the creation and context.



Figure 2: Historical panoramic maps of Seattle, WA (1891), and Santa Fe, NM (1882), showing landmarks, street layouts, and historical buildings.

2.2 ColPali

ColPali offers two base models: the original ColPali built on Google’s PaliGemma (a 3-billion-parameter vision-language model) [3], and ColQwen built on Alibaba’s Qwen2 VLM (2 billion parameters) [6]. Both models use multiple patch embeddings to capture details across different sections of a page. ColQwen generates 768 patch embeddings per page, while ColPali generates 1024, though both use the same embedding dimension of 128. We use ColQwen in this implementation to reduce computational costs.

2.3 Multimodality and Digital Humanities

While map collections are routinely utilized by researchers across disciplines, searching maps oftentimes suffer from fundamental limitations [10]. For example, map metadata is never completely descriptive, and neither is the OCR-able text appearing on a given map. Given the sheer scale of digitized maps available, developing new ways to search these maps is more important than ever.

This paper builds on work in the digital humanities that has demonstrated the value of multimodal search for digital collections held by libraries, archives, and museums. Deemed the “multimodal turn” [15], models such as CLIP [13] have shown great promise for increased discoverability for digital collections ranging from photojournalism collections [8] to lantern slides [14]. This paper extends our previous work surrounding applying CLIP to digitized maps for search & discovery, informed by the perspectives of staff at the Library of Congress [12]. However, as described earlier, CLIP struggles with text recognition, especially with longer passages. For this reason, we have adopted ColPali.

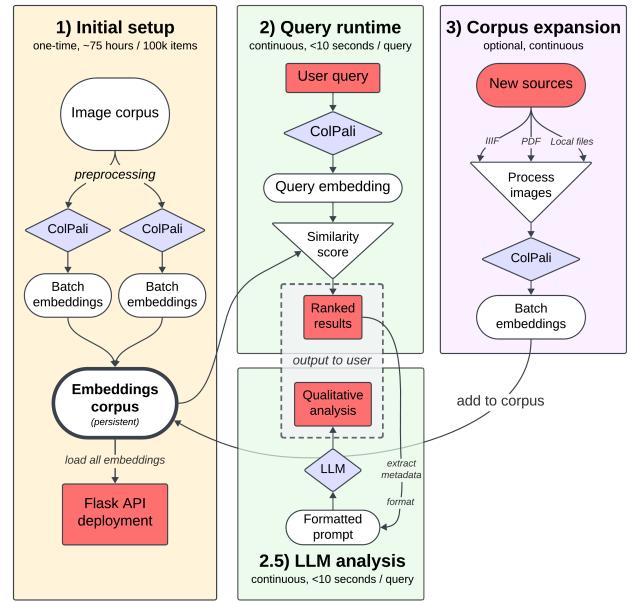


Figure 3: A full flow chart of our 3-stage pipeline. Red shapes indicate direct interface with the user, while blue diamonds are models loaded onto the server. The embeddings corpus is the only persistent object at 28GB.

A number of other works have begun exploring the application of ColPali to digital cultural heritage [5, 9, 16]. In addition to applying ColPali, we utilize Llama 3.2 [7] to generate summaries of search results, in line with RAG-style search.

3 Method

Our fast and reusable retrieval-augmented search engine is motivated by the task of searching map data for cultural institutions:

- (1) First, to ensure accessibility for resource-constrained institutions, the tool should create high-quality embeddings of complex visual data quickly on consumer-grade hardware. The method should be available for practitioners who want to implement it on their own data.
- (2) It should embed several modes into a common latent space to enable comparison between them.
- (3) Because maps have both visual and textual data, it should distinguish between sections of a single document and parse them individually, rather than create a single global embedding.
- (4) It should allow users to easily add new items to the corpus as they come up.

As shown in Figure 3, our method follows three stages: data preparation, embedding generation, and search interface deployment. We describe all three stages below.



analysis: **Main Themes and Time Periods:**
The majority of the results represent maps from various time periods, including 16th-century world maps (Results 1-3), 18th-century world maps (Results 4-6), and 19th-century globe maps (Result 5). The earliest map dates back to 1660, while the most recent one is from 1900.

Subject Matter:
Maps of the world are the dominant type of result, with various authors and cartographers contributing to this category. Some notable contributors include:
 • 16th-century mapmakers like Wit, Frederik de, and Visscher
 • 18th-century cartographers such as Godson, W., Bonne, Rigobert, and Goldthwaite, William M.
 • Modern globe maps by Cartwright (Note: Unfortunately, I couldn't identify the specific creator of this result)

Patterns in Materials:
The types of materials used to create these results are diverse. Some common formats include:
 • Relief maps showing pictorially rendered features
 • World maps with ancillary illustrations and inset maps
 • Globe maps with hachures or other map styles

Key Creators, Locations, or Topics:
The following creators, locations, or topics dominate the results:
 • Wit, Frederik de (mapmaker): Result 1-3
 • Godson, W. (cartographer): Result 2
 • Bonne, Rigobert (cartographer): Result 5
 • Goldthwaite, William M. (cartographer): Result 6

prompt: maps from the Renaissance with illustrations of ships



analysis:

Main Themes:
1. Maps with illustrations of ships from the Renaissance period (1480-1556)
2. Atlases and nautical charts from various time periods
3. Documents related to exploration, trade, and colonization

Time Periods:
 • 15th century: Several maps and documents appear in this period, including those from Henry Popple's "Nautical Atlas of the World" (1519) and Jean André Brémont's "Portolan atlas of the Mediterranean Sea and western Europe" (1670)
 • 16th-18th centuries: Maps and documents from these periods include titles such as "A map of California, from 1597 to 1851" (1597), "Nautical Atlas of the World, Folio 2 Recto, Northeastern Atlantic Ocean and Northern Europe and Folio 2 Verso, Central Atlantic Ocean with the Azores" (1519)
 • Late 19th-early 20th centuries: Documents from this period include maps from Emanuel Miller's "Miller Atlas in the collections of the National Library of France" (1519) and maps from Henry Popple's "A map of California, from 1597 to 1851" (1597)

Subject Matter:
 • Exploration and trade
 • Colonization and territorial claims
 • Cartography and nautical knowledge

Patterns:
 • Many of the results appear to be related to exploration and colonization, with maps and documents showing European powers' interests in various regions.
 • The "Nautical Atlas of the World" series suggests a focus on maritime navigation and cartography.
 • There are also many documents related to California's early history, indicating that this region was an important hub for trade and exploration during this period.

Key Creators, Locations, or Topics:
 • Henry Popple (1519-1556) appears in several results
 • António de Holanda (1480-1556) is mentioned in a few instances
 • Emanuel Miller's "Miller Atlas" is discussed in one of the results
 • The "Nautical Atlas of the World" series, produced by cartographers Pedro Reinel and Jorge Reinel, is highlighted in several results

Figure 4: Search on a non-LOC image results in visually similar images taken from the LOC's corpus. Analysis highlights details on the results taken from metadata.

3.1 Data Preparation

Working from a central CSV file, we download batches of 500 images sequentially, accumulating them in temporary storage and processing them in sub-batches to respect GPU memory constraints. In keeping with cultural heritage norms, we work in the International Image Interoperability Framework (IIIF) throughout our pipeline.

3.2 Embedding Generation

Working in batches, we embed each image in the corpus using the pre-trained ColQwen2 model. On a single NVIDIA T4 Tensor Core GPU, a single embedding processes in 3 to 4 seconds, with its model requiring 5.62GB of GPU memory.¹ Each embedding is put into a dictionary with its unique IIIF identifier, which can be used to reference relevant metadata from the master CSV. Embeddings are saved at the batch level in distributed pickle files.

3.3 Search Interface

Once a base set of embeddings is processed, we deploy a Flask-based REST API that loads all distributed embeddings into memory on startup. The API exposes several endpoints for corpus management and querying, with the primary search endpoint accepting POST requests containing a natural language query string and a parameter k specifying the number of results to return (default: 5).

Query processing follows the late-interaction paradigm: the user's text query is embedded using the same ColQwen2 model, then scored against all document embeddings using the processor's MaxSim computation, which determines the maximum similarity between each query token vector and all document patch vectors,

¹On an AWS g4dn.4xlarge instance, this requires 75 hours of total compute, amounting to approximately \$90 dollars.

Figure 5: The tool allows us to discern features on a map like illustrations of ships that do not appear in the items' metadata.

preserving fine-grained semantic matching while remaining computationally tractable at the scale of the corpus.

In Figures 4 and 5, we show a screenshot of an example search, along with several top results. Results are ranked by similarity score and returned with metadata including the document title, the image itself, the 'loc.gov' API resource, the document type, and numerical score. Response times for queries against our corpus of 25,000+ embedded images average under one second after the initial model load, making the system suitable for interactive exploration.

On a server with a single NVIDIA T4 Tensor Core GPU for AI inference, submitting a query through the 120,000-item embedding corpus takes roughly 6 seconds. Indeed, inference time grows linearly with corpus size, demanding new techniques for small-scale inference once corpus size exceeds a few hundred thousand. We leave this for future work.

3.3.1 Image search. Because vision-language models embed visual and textual data to a common embedding space, the tool can easily accept input images for a reverse-image-search functionality. After images are processed, this follows the same scoring regime as text-to-image search. Notably, the model does not transfer between modes in this case, and the similarity scores between an image query and its result tend to be an order of magnitude higher than between a text query and its result.

3.3.2 Retrieval-augmented Search Framework. Our implementation follows a retrieval-augmented generation (RAG) approach that allows dynamic expansion of the corpus. Users can augment the base set with their own documents, whether from partner institutions' IIIF servers, local digitized materials, or PDF publications, during or between search sessions. These personalized search contexts allow a user to compare their own materials to those of the corpus, and to customize the search engine for their own uses.

More broadly, persistent corpus expansion allows institutions to collaboratively build federated search indices, where multiple organizations contribute embeddings from their collections without centralizing the underlying image files, respecting data sovereignty while enabling cross-institutional discovery.

All additions, regardless of origin, are processed through the identical ColQwen2 embedding pipeline, projected into the same 128-dimensional vector space, and scored using the same late-interaction mechanism. This ensures that user-contributed materials are semantically comparable to the base corpus.

3.3.3 LLM Analysis. Lastly, we include a feature that allows a user to generate qualitative LLM analysis on the results of their search. After a search is completed, the pipeline extracts metadata from the search result, formats it into a prompt, and passes it to a lightweight Llama 3.2-1B, which provides short, digestible information on the major themes, time periods, formats, subjects, and authors of the results. Analysis typically runs for 5-15 seconds, depending on the number of results requested.

4 Use Cases

We envision our retrieval-augmented search system for large-scale map collections to have a number of uses for both collection curators and end-users. On the curatorial side, archivists, curators, catalogers and collection stewards responsible for digitizing collections and providing access to them typically rely on traditional manual methods of cataloging and producing finding aids. This process is time-intensive, especially for larger-scale collections, such as the ones held by the Library of Congress. Put simply, approaching the scale of half a million map images necessitates additional methods to draw connections across a digital collection. With our retrieval-augmented search system, archivists, curators, and catalogers can identify patterns at scale and also interact with the collection in new, multimodal ways. They can test hypotheses about the prevalence of different features during finding aid construction and also make associations that otherwise might not be possible by relying on existing metadata. As digitization efforts continue to grow, and born-digital collections approach orders of magnitude larger, such methods will only become more important.

Similarly, end-users from academic researchers to members of the public face the challenge of how to make sense of large-scale digital collections. Academics searching map collections for specific motifs, aesthetic qualities, or accompanying text might encounter the analogous problem with existing forms of search. Our system enables end-users to define more flexible queries, generate relevant summarizations, and more generally see collections at scales not possible with basic metadata search (i.e., “distant viewing” [2]).

Lastly, we highlight the possibilities for integrating *inter-collection* search. As described in Section 3.3.2, our system enables others to dynamically expand the collection included. The possibilities with this dynamic corpus expansion functionality are manifold. First, archivists and collection stewards can utilize this approach to search *across* collections that are ordinarily siloed. Second, digital humanities scholars can use this approach for exploratory analysis in order to identify similarities and differences between collections. Third, as articulated in Section 3.3.2, we see a broader vision of enabling the multi-institutional construction of federated search indices using these vectorized approaches.

5 Future Work

We conclude by offering a few directions for future work. First, we plan to expand our retrieval-augmented search system into a full vision language model-based RAG system by generating summaries based on what the model finds from the map images directly, not just the map metadata. Second, we will explore finetuning ColPali for specific use with digital cultural heritage collections. Third, we hope to expand our system into a more flexible system for searching a range of digital cultural heritage collections, as well as conduct user evaluations of it. We encourage users to visit <http://www.mapras.com> and test their own prompts, and welcome feedback.

References

- [1] [n. d.]. Research Guides: Cataloging Cartographic Materials: Using the Map Collections — guides.loc.gov. <https://guides.loc.gov/cataloging-cartographic-materials/using-the-map-collections>. [Accessed 27-10-2025].
- [2] Taylor Arnold and Lauren Tilton. 2019. Distant viewing: analyzing large visual corpora. *Digital Scholarship in the Humanities* 34 (03 2019). doi:10.1093/llc/fqz013
- [3] Lucas Beyer, Andreas Steiner, André Susano Pinto, Alexander Kolesnikov, Xiao Wang, Daniel Salz, Maxim Neumann, Ibrahim Alabdulmohsin, Michael Tschanen, Emanuele Bugliarello, Thomas Unterthiner, Daniel Keysers, Skanda Koppara, Fangyu Liu, Adam Grycner, Alexey Gritsenko, Neil Houlsby, Manoj Kumar, Keran Rong, Julian Eisenschlos, Rishabh Kabra, Matthias Bauer, Matko Bošnjak, Xi Chen, Matthias Minderer, Paul Voigtlaender, Ioana Bica, Ivana Balazevic, Joan Puigcerver, Pinelopi Papalampidi, Olivier Henaff, Xi Xiong, Radu Soricut, Jeremiah Harmsen, and Xiaohua Zhai. 2024. PaliGemma: A versatile 3B VLM for transfer. arXiv:2407.07726 [cs.CV] <https://arxiv.org/abs/2407.07726>
- [4] Phaidon Editors and John W. Hessler. 2016. *Map: Exploring the World*. Phaidon Press.
- [5] Nicola Fanelli, Gennaro Vessio, and Giovanna Castellano. 2025. ArtSeek: Deep artwork understanding via multimodal in-context reasoning and late interaction retrieval. arXiv:2507.21917 [cs.CV] <https://arxiv.org/abs/2507.21917>
- [6] Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. 2025. ColPali: Efficient Document Retrieval with Vision Language Models. arXiv:2407.01449 [cs.IR] <https://arxiv.org/abs/2407.01449>
- [7] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhiav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Bin Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaïdis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, Danny Wyatt,

David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Francisco Guzmán, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Govind That-tai, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jack Zhang, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junting Jia, Kalyan Vasuden Alwala, Karthik Prasad, Kartikay Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Kushal Lakhota, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Maria Tsimpoukelli, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Ning Zhang, Olivier Duchenne, Onur Celebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohan Maheswari, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidor, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharann Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vitor Albiero, Vladian Petrovic, Weiwei Chu, Wenhan Xiong, Wenying Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaofang Wang, Xiaoqing Ellen Tan, Xide Xia, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yunling Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aayushi Srivastava, Abha Jain, Adam Kelsey, Adam Shajfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Amos Teo, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Dong, Annie Franco, Anuj Goyal, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharhambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Ce Liu, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Cynthia Gao, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, David Adkins, David Xu, Davide Testuggine, Delin David, Devi Parikh, Diana Liskovich, Dudem Foss, Dingkang Wang, Due Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Eric-Tuan Le, Erik Brinkman, Esteban Arcuate, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Filippos Kokkinos, Firat Ozgenel, Francesca Caggioni, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Grant Herman, Grigory Sizov, Guangyi Zhang, Guna Lakshminarayanan, Hakan Inan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspregon, Hunter Goldman, Hongyuan Zhan, Ibrahim Damlij, Igor Mol'yobog, Igor Tufanov, Ilias Leontiadis, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Janice Lam, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carville, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kiran Jagadeesh, Kun Huang, Kunal Chawla, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madiam Khabsa, Manav Avalani, Manish Bhatt, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Miao Liu, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singh, Nick Egebo, Nicolas Usunier, Nikhil Mehta, Niko-lay Pavlovich Laptev, Ning Dong, Norman Cheng, Oleg Chernoguz, Olivia Hart,

- Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghatham Murthy, Raghu Nayani, Rahul Mitra, Ran-gaprabhu Parthasarathy, Raymond Li, Rebekkah Hogan, Robin Battye, Rocky Wang, Russ Howes, Ruty Rinott, Sachin Mehta, Sachin Siby, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Mahajan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shauna Lindsay, Shauna Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shishir Patil, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Summer Deng, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Koehler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vlad Ionescu, Vlad Poenaru, Vlad Tiborius Mihailescu, Vladimir Ivanov, Wei Li, Wencheng Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaojian Wu, Xiaolan Wang, Xilun Wu, Xinbo Gao, Yaniv Kleinlein, Yanjun Chen, Ye Hu, Ye Qi, Yenda Li, Yelin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu Wang, Yu Zhao, Yuchen Hao, Yundi Qian, Yunlu Li, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, Zhiwei Zhao, and Zhiyu Ma. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] <https://arxiv.org/abs/2407.21783>
- [8] Ying-Hsiang Huang and Benjamin Charles Germain Lee. 2025. Digital Collections Explorer: An Open-Source, Multimodal Viewer for Searching Digital Collections. arXiv:2507.00961 [cs.DL] <https://arxiv.org/abs/2507.00961>
- [9] Haneol Kim, JONGHWAN BAE, Seonwoo Shin, and Sanghun Park. 2025. Fashion-RNA: Interactive Reimagination of Fashion Heritage. In *ICCV 2025 Workshop on Cultural Continuity of Artists*. <https://openreview.net/forum?id=CAKupmBs9w>
- [10] Marta Kuzma and Albina Moscicka. 2020. Evaluation of Metadata Describing Topographic Maps in a National Library. <https://www.proquest.com/working-papers/evaluation-metadata-describing-topographic-maps/docview/2532311981/se-2>
- [11] Geography Library of Congress and Map Division. [n. d.]. Sanborn Maps Collection. <https://www.loc.gov/collections/sanborn-maps/about-this-collection/>
- [12] Jamie Mahowald and Benjamin Charles Germain Lee. 2024. Integrating Visual and Textual Inputs for Searching Large-Scale Map Collections with CLIP. arXiv:2410.01190 [cs.IR] <https://arxiv.org/abs/2410.01190>
- [13] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020 [cs.CV] <https://arxiv.org/abs/2103.00020>
- [14] Thomas Smits and Mike Kestemont. 2021. Towards Multimodal Computational Humanities. Using CLIP to Analyze Late-Nineteenth Century Magic Lantern Slides.. In *CHR*. 149–158.
- [15] Thomas Smits and Melvin Wevers. 2023. A multimodal turn in Digital Humanities. Using contrastive machine learning models to explore, enrich, and analyze digital visual historical collections. *Digital Scholarship in the Humanities* 38, 3 (03 2023), 1267–1280. arXiv:<https://academic.oup.com/dsh/article-pdf/38/3/1267/51309490/fqad008.pdf> doi:10.1093/lhc/fqad008
- [16] Haofeng Wang, Yilin Guo, Zehao Li, Tong Yue, Yizong Wang, Enci Zhang, Rongqun Lin, Feng Gao, Shiqi Wang, and Siwei Ma. 2025. RiverEcho: Real-Time Interactive Digital System for Ancient Yellow River Culture. arXiv:2506.21865 [cs.MM] <https://arxiv.org/abs/2506.21865>
- [17] Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. 2024. Long-CLIP: Unlocking the Long-Text Capability of CLIP. arXiv:2403.15378 [cs.CV] <https://arxiv.org/abs/2403.15378>
- [18] Min Zhang. 2012. From Washington to the world: maps and digital archives at the Library of Congress. *International Journal of Humanities and Arts Computing* 6, 1–2 (2012), 100–110. doi:10.3366/ijac.2012.0041