

Bridging the Provenance Gap across Text, Speech, & Video



Shayne Longpre, Nikhil Singh, Manuel Cherep, Kushagra Tiwary, Joanna Materzynska, William Brannon, Robert Mahari, Naana Obeng-Marnu, Manan Dey, Mohammed Hamdy, Nayan Saxena, Ahmad Mustafa Anis, Emad A. Alghamdi, Vu Minh Chien, Da Yin, Kun Qian, Yizhi Li, Minnie Liang, An Dinh, Shrestha Mohanty, Deividas Mataciunas, Tobin South, Jianguo Zhang, Ariel N. Lee, Campbell S. Lund, Christopher Klamm, Damien Sileo, Diganta Misra, Enrico Shippole, Kevin Klyman, Lester JV Miranda, Niklas Muennighoff, Seonghyeon Ye, Seungone Kim, Vipul Gupta, Vivek Sharma, Xuhui Zhou, Caiming Xiong, Luis Villa, Stella Biderman, Alex Pentland, Sara Hooker, Jad Kabbara

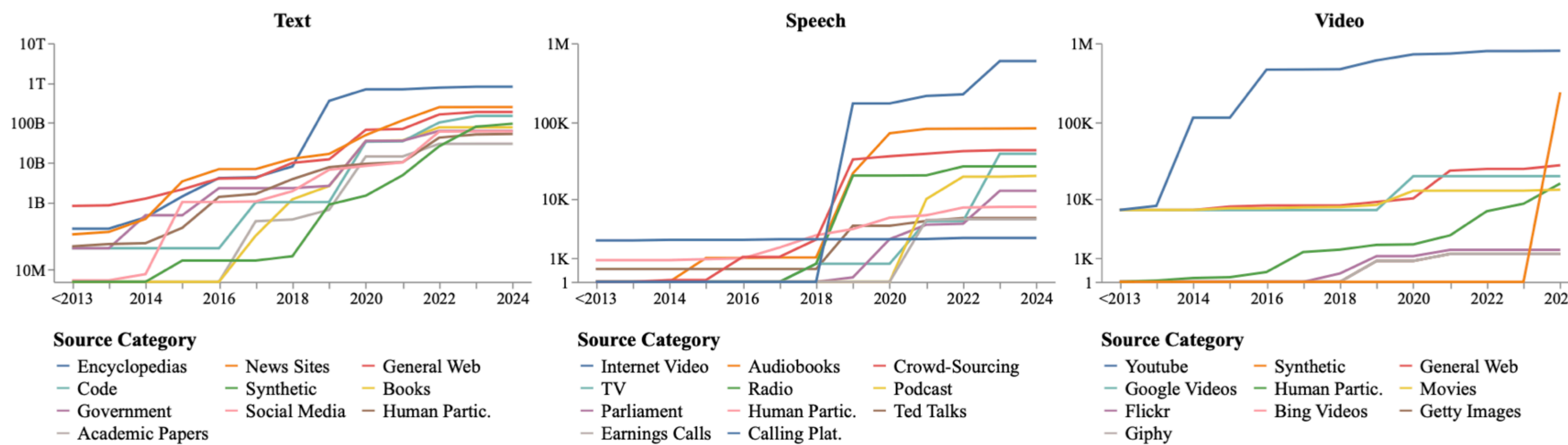
Summary:

- To understand the AI data ecosystem, we audit 3916 of the most popular, public ML datasets for **text**, **speech**, and **video** generation—summarized here
- We trace their **sources**, **license restrictions**, **languages**, and the **geographical origins** of the organizations that published them. (All open sourced!)
- We outline the three major observations on data sourcing, restrictions, and representation.

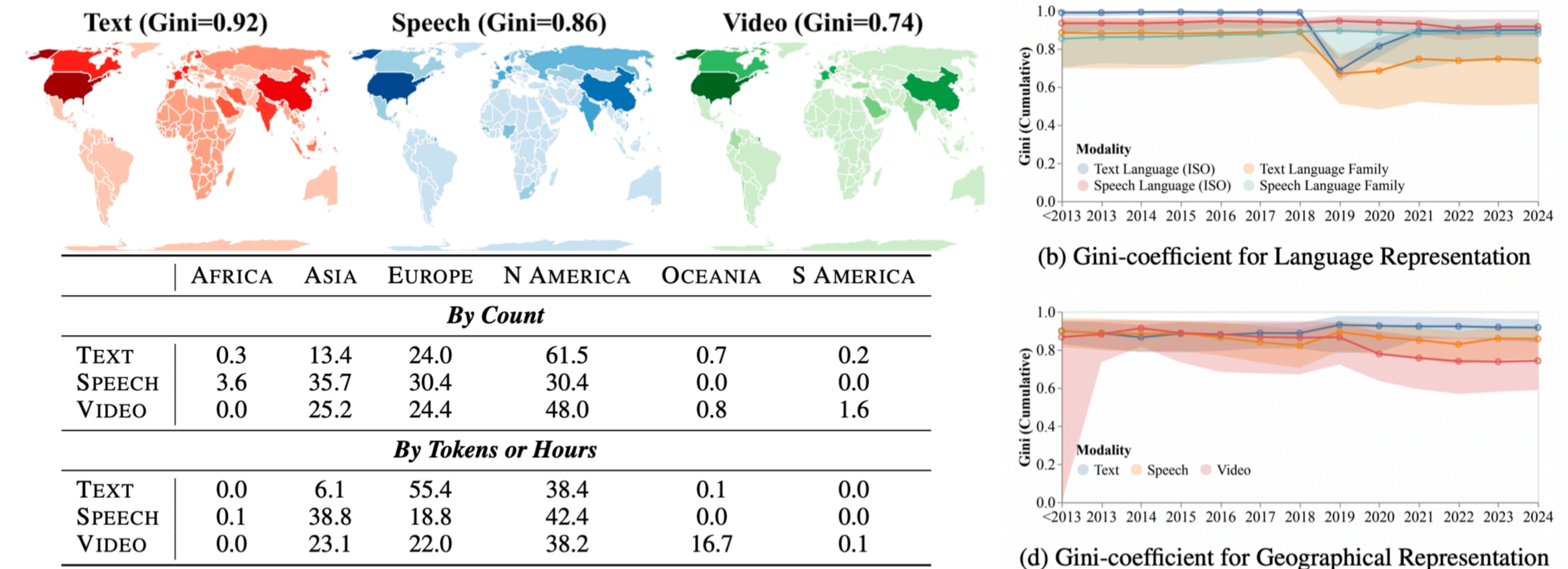
Our audit spans ~4k datasets, from 798 sources, 659 orgs, in 67 countries, & 608 languages.

	DATASETS		SOURCES		CREATOR ORGS		LANGUAGES		TASKS	LICENSES
	#	SIZE	#	DOMAINS	#	COUNTRIES	#	FAMILIES		
TEXT	3717	2.1T	713	23	534	60	502	21	395	50
SPEECH	95	775k	51	16	124	29	260	36	18	19
VIDEO	104	1.13M	44	24	101	23	-	-	33	11
TOTAL	3916	-	798	83	659	67	608	37	443	55

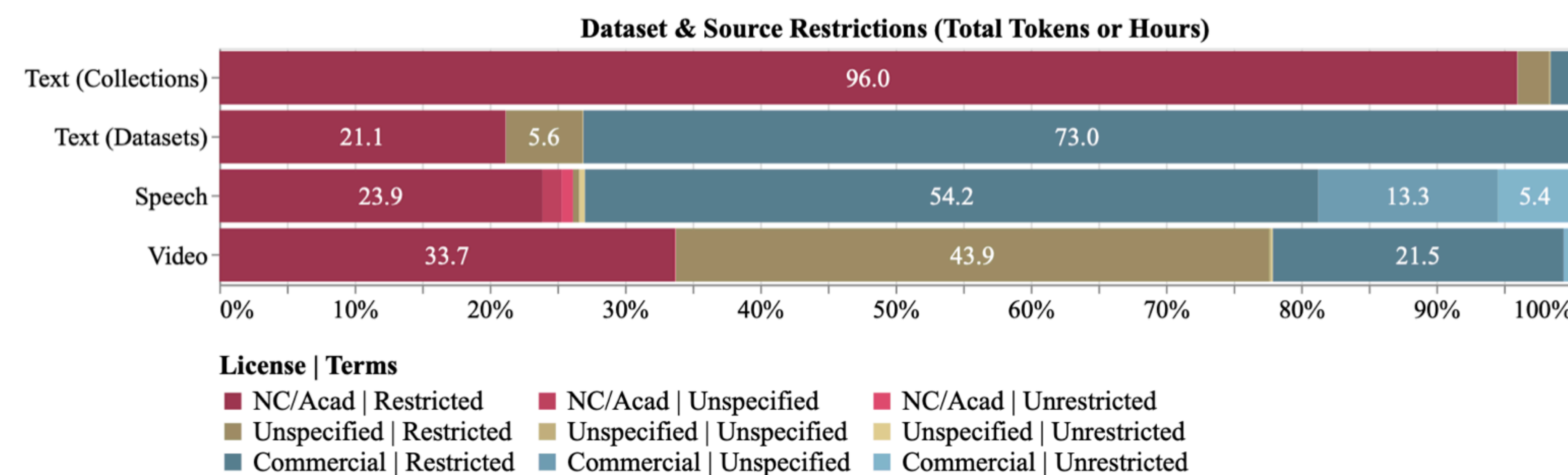
1 Result: Multimodal data is increasingly sourced from the web, social media (YouTube), or synthetically generated



3 Result: Geographical and linguistic representation have not improved for a decade, across the data ecosystem.



2 Result: 20-33% of datasets are non-commercially licensed, but >80% have undocumented source restrictions.



While the amount of data from under-represented creators and languages increases each year, their relative representation remains consistently western-centric, with no significant improvements from >0.7 Gini coefficients.

Check out our paper here:

