



Julia Lane

One Sentence Summary: The current approaches to measuring critical technologies can provide robust information about national competencies with a people-centered conceptual framework.

Research Issue

The CHIPS and Science Act requires investments in key technologies to “[grow both curiosity-driven and translational research, ensuring both the creation of new ideas and the ability of those ideas to create new innovations, products, companies, and jobs in the United States...and build new technology hubs across the country, increasing the participation of underrepresented populations and geographies in innovation.](#)” The act also requires NSF’s TIP Directorate to routinely report on the effects of their investments in ten strategic technology areas, including Artificial Intelligence (AI). This requires developing a people-centered framework to understand how ideas flow from research investments to the non-academic sector.

Methods and Data

These requirements pose three challenges. The first challenge is *conceptual*: developing a theory of change. The system should be designed to achieve its objective and goals in a deliberate manner. A “theory of change” for measuring the impact of investments in critical technologies—that is, a causal model or map of how the goals of a program are intended to be achieved—can inform the investment process and provide a framework for its planning and evaluation. This includes articulating the inputs (i.e., available resources to leverage), activities (i.e., actions or work conducted to advance the program), outputs (i.e., the immediate, practical benefits of the program), outcomes (i.e., medium-term results), and longer-term impacts of the overall effort, and how each successively feeds into the next (See the [NAIRR report](#) for an illustrative example). A logical conceptual framework could be derived from the Nobel prize winning work of Paul Romer’s endogenous technological change, which argues that economic growth is generated by the non-rivalrous transmission of ideas – succinctly summarized by Oppenheimer’s famous quote that the best way to transmit knowledge is to wrap it up in a human being. The complementary empirical framework should characterize scientific fields by networks of individuals who work on similar topics, using similar methods and data, and firms grouped into industries based on the clustering of such individuals. It could go beyond using authors on publications, since there is clear evidence that there is attribution bias that excludes many people who have ideas, and do the work, but are not credited. It could include data on grants to be more inclusive and timely.

The second is *technical* and is itself multifaceted. The concepts that underpin cutting-edge AI research and technology shift quickly; cross disciplinary boundaries; and are not consistently identified as AI, especially in applications. The regional focus of TIP investments also poses data collection and confidentiality challenges as it requires levels of granularity in reporting that are not typically achievable with standard federal data sources. Timing is the last technical challenge. To be truly useful, TIP and other stakeholders, including policy makers, require near real-time information. Many “gold standard” data sources that could be used for these tasks (e.g. NCSSES data on doctorate recipients or the Census Bureau’s LEHD) have a two or more year lag making them ineffective for timely reporting and policy making.

The third challenge is in *reporting*. All the new technology areas being emphasized affect firms and workers, and it is essential that reports present metrics for the relevant classes of parties and

communicates findings in a way that is accessible, useful, and responsible. Yet established survey methods, data systems, and classification systems do not identify the companies, jobs, and people affected by NSF investments. As a result, current data systems cannot be used to identify the firms and people who might be affected.

Our [paper](#) for the American Enterprise Institute describes a new, rapidly implementable, conceptual, and empirical approach to tracing how ideas move from investments in research to the marketplace and developing early warning indicators of potential workforce and education impacts. The report proposes a new evidence-based foundation to support US national growth strategies and ensure investments have the greatest chance of success for workers and employers.

Our very preliminary research report developed the first step in the process. The assumption in using simply search strings in bibliometrics is that the distribution of the text will remain similar over time. However, the language of Artificial Intelligence (AI) changes as new methods are developed and others lose popularity. Faced with this problem, the National AI Research Resources Task Force suggested that the best way track AI authors and their work was to leverage the fact that most of the community doing research in AI publishes in AI dedicated conference proceedings. The working framework is, then, that authors in the field tend to continue to produce research within the field and are more constant over time. Such a framework is consistent with the people centered framework that we have outlined. Identifying networks of authors who do research in AI and those applying the latest AI research to problems in other fields can be considered a graph problem.

In technical terms, we would like to find a function f that satisfies the following relationship, $f(G, G') = y$. $G_{ai} = (V, E)$, is a graph with vertices, V , representing documents and authors and edges, E , representing the authorship relationship, $G_p = (V', E')$, is a graph representing a single document that we would like to classify, and $\hat{y} = R^{|V|}$: $0 \leq \hat{y} \leq 1$, is a classification score for each of the nodes in the G_p where values closer to 1 indicate that nodes within G_p are similar to nodes within G_{ai} and values closer to 0 indicate that G_p is not similar nodes within G_{ai} . The seed graph, $G_{ai} = \langle V, E \rangle$ consists of the authors and papers accepted to at least one of the main AI conferences V , and the edges, E , that represent the authorship relationship between the authors and papers. We consider the problem as a node classification problem where each vertex in the graph has a score attribute between 0 and 1 that indicates a confidence for how likely an author is an AI author. G_{ai} is naturally a largely disconnected graph as most authors are not on most papers. All vertices in G_{ai} from an AI conference are assigned a score of 1. In this work we focus on the vertices in V that are authors, and only use authors for inference.

The metadata on research publications is derived from Scopus data. Scopus is Elsevier's abstract and citation database of peer-reviewed literature. Scopus includes data and linkages across 91 million items from over seven thousand publishers, 94 thousand affiliations, and 17 million authors. It is the largest curated abstract and citation database of peer-reviewed literature and provides a useful view on the research landscape. Data and computational resources were provided by the International Center for the Study of Research (ICSR) Lab; a cloud-based computational platform which enables researchers to analyze large, structured datasets. For exploratory projects, replication studies, or when developing new research metrics and indicators, ICSR Lab supports scholarly research by giving access, at no cost, to powerful research metadata and metrics

Insights

We found that there is potential to use this approach to generate sensible results on authors, their institutions, and their funders, although future work will include both enhancing the algorithm and



applying it to grant information. Over the past thirteen years (2010–2022), global research on Artificial Intelligence (AI) has been growing at a rapid pace. Publications in this research area have grown to represent around 2% of all research published in Scopus in 2022, up from less than 1% of all research published in Scopus in 2010. Growth has been particularly high over the past four years; almost half of the 567,740 publications since 2010 were published between 2019 and 2022. The number of sources (e.g., journals, conferences) shows similar trends, although to a lesser degree. The number of researchers publishing has grown similarly. As of 2020, China’s AI research output has been the highest of any country/region. Major funders, in order, are National Natural Science Foundation of China, National Science Foundation, National Key Research and Development Program of China, National Institutes of Health, Fundamental Research Funds for the Central Universities and the Horizon 2020 Framework Programme. The top states, in order, were CA, NY, MA, PA, TX, IL and WA. Of institutions, Carnegie Mellon University ranks first in both years. Stanford University, MIT, UC Berkeley, and the University of Illinois also were top publishers in both years.