

XDMoD Value Analytics: Visualizing the Impact of Internal IT Investments on External Funding, Publications, and Collaboration Networks

1 **Olga Scrivner^{*1}, Gagandeep Singh¹, Sara E. Bouchard¹, Scott C. Hutcheson¹, Ben Fulton²,**
2 **Matthew R. Link², Katy Börner¹**

3 ¹Department of Intelligent Systems Engineering, School of Informatics, Computing and Engineering,
4 Indiana University, Bloomington, Indiana, USA

5 ²Pervasive Technology Institute, Indiana University, Bloomington, Indiana, USA

6 *** Correspondence:**

7 Corresponding Author
8 obscrivn@indiana.edu

9 **Keywords: information visualization, scientometrics, impact analysis, grant income, return on**
10 **investment, value analytics, high-performance computing**

11 Abstract

12 Many universities invest substantial resources in the design, deployment, and maintenance of
13 campus-based cyberinfrastructure. To justify the expense, it is important that university
14 administrators and others understand and communicate the value of these internal investments in
15 terms of scholarly impact. This paper introduces two visualizations and their usage in the Value
16 Analytics (VA) module for Open XD Metrics on Demand (XDMoD), which enable analysis of
17 external grant funding income, scholarly publications, and collaboration networks. The VA module
18 was developed by Indiana University's (IU) Research Technologies division, Pervasive Technology
19 Institute, and the [Cyberinfrastructure for Network Science Center](#) (CNS), in conjunction with the
20 University at Buffalo's Center for Computational Research (CCR). It provides diverse visualizations
21 of measures of information technology (IT) usage, external funding, and publications in support of IT
22 strategic decision making. This paper details the data, analysis workflows, and visual mappings used
23 in two VA visualizations that aim to communicate the value of different IT usage in terms of NSF
24 and NIH funding, resulting publications, and associated research collaborations. To illustrate the
25 feasibility of measuring IT values on research, we measured its financial and academic impact from
26 the period between 2012 and 2017 for IU. The financial return on investment (ROI) is measured in
27 terms of IU funding, totaling \$339,013,365 for 885 NIH and NSF projects associated with IT usage,
28 and the academic ROI constitutes 968 publications associated with 83 of these NSF and NIH awards.
29 In addition, the results show that *Medical Specialties*, *Brain Research*, and *Infectious Diseases* are
30 the top three scientific disciplines ranked by the number of publications during the given time period.
31

32 1 Introduction

33 Access to high-performance computing (HPC) systems and advanced cyberinfrastructure generally is
34 critical to advance research in many scholarly fields. Over the last thirty years, supercomputing has
35 expanded from a few monolithic and extremely fast computing systems into a comprehensive “set of
36 organizational practices, technical infrastructure and social norms that collectively provide for the

smooth operation of research and education work at a distance” (Towns et al. 2014). This new form of cyberinfrastructure (Stewart et al. 2017) is also used by a large number of researchers, scholars, and artists, and is more complicated, including high performance computing, storage systems, networks, and visualization systems. This requires a new suite of metrics that can help different stakeholders better understand the value of research cyberinfrastructure. In addition to job- and system-performance monitoring, metrics now available and needed to understand system usage and value include usage modality (e.g., central processing unit [CPU] usage per group or person, number of users, and wait time). All of these metrics are essential to better understand “why users do what they do and how they leverage multiple types and instances of cyberinfrastructure (CI) resources” (Katz et al. 2011). In recent years, there has also been growing interest in measuring the impact of CI on scientific research and publication outcome (e.g., Knepper and Börner 2016, Fulton et al. 2017, Madhavan et al. 2014, among others). Such insights are particularly relevant for campus CI that require significant investment and long-term strategic and financial planning. Visualizations help communicate the value of CI to diverse stakeholders ranging from domain experts to academic deans to financial administrators.

Currently, several CI frameworks enable on-demand rendering of impact metrics (e.g., funding, publications, and citations) that result from using HPC resources. While most monitoring tools for HPC are traditionally “largely passive and local in nature” (Furlani et al. 2013), these new frameworks, such as Deep Insight Anytime, Anywhere (DIA2) (Madhavan et al. 2014) and Extreme Science and Engineering Discovery Environment (XSEDE) Metrics on Demand (XDMoD VA) (Fulton et al. 2017), are open-source, customizable systems with “increased functionality, an improved interface, and high-level charting and analytical tools” (Palmer et al. 2015). DIA2 is a web-based visual analytics system aiming to assess research funding portfolios (Madhavan et al. 2014). Open XDMoD is one of the most widely used software systems in the U.S. and is currently employed by more than 200 institutions to evaluate their HPC usage (Palmer et al. 2015, Fulton et al. 2017). The tool has been developed as an open source software for “metrics, basic accounting, and visualization of CPU (central processing units) and storage usage” at the Center for Computational Research (CCR) of the University at Buffalo (Palmer et al. 2015).

Open XDMoD VA adds a value analytics module to Open XDMoD. The two visualizations presented in this paper contribute to XDMoD VA functionality as follows:

1. Financial and Intellectual Analytics: Analyze grant income and publications by researchers that use Information Technology (IT) and relate that income to use of local IT systems.
2. Co-PI collaboration networks: Analyze research collaborations by researchers that use local IT systems.

In this paper, we detail both types of visual analytics for XDMoD VA and demonstrate how they can help understand the IT impact on academic research. Specifically, we describe the data and visual analytics workflows used for the *Funding and Publication Impact* and the *Co-PI Collaboration Network* visualizations. The *Funding and Publication Impact* visualization uses a Sankey graph to interlink IT usage with funding and publication output. The *Co-PI Collaboration Network* uses NSF and NIH funding data to extract and depict Co-PI collaboration networks together with listings of scholars ranked by the total amount of grants. Both types of visualizations are interactive, supporting overview first, zoom and panning, and details on demand (Shneiderman 1996). The visualizations are rendered using the Web Visualization Framework (WVF) developed by CNS at IU that allows for an effective, highly customizable rendering of interactive visualizations.

82 The remainder of the paper is organized as follows: Section 2 discusses prior work on analyzing and
 83 visualizing the impact of IT resources, particularly on scholarly output. In Section 3, we outline the
 84 data we used, as well as the preprocessing methods needed to render that data useful. Section 4
 85 discusses the general methods applied to render data into visual insights. Section 5 details the results
 86 of using the data and methods to render the *Funding and Publication Impact* and *Co-PI*
 87 *Collaboration Network* visualizations. Key insights and planned developments are discussed in
 88 Section 6.
 89

90 2 Related Work

91 In recent years, there has been growing interest in measuring the impact of high-performance
 92 computing (HPC) on scientific research and publication outcome. Advanced cyberinfrastructure
 93 resources require significant investment, particularly when implemented at the campus level, and
 94 insights about the value of such investments in financial and intellectual terms are essential for the
 95 strategic and financial planning of academic institutions (Fulton et al. 2017). A number of studies
 96 have examined the relationship between these values and CI resources usage. Three studies can be
 97 regarded as a starting point for introducing HPC user metrics: namely, Li et al. (2005), Iosup et al.
 98 (2006), and Lee et al. (2006). In addition to traditional metrics (e.g., job size, system utilization),
 99 these studies included user and group characteristics to analyze system performance. Hart (2011)
 100 extended this idea and highlighted the importance of usage and submission patterns to understand
 101 users and their behavior across HPC resources. Knepper (2011) further investigated the relation
 102 between users, their HPC usage behavior, and their field of science. Particularly, his study examined
 103 PIs, their network affiliation, and the scientific field and allocation size of their research projects
 104 from 2003 to 2011 using TeraGrid, a national (US) computing scientific infrastructure. The results
 105 revealed that PIs constituted 23% (3,334) of the total TeraGrid project users (14,474) and that
 106 molecular biosciences and chemistry are the two scientific fields with the highest number of projects
 107 involving TeraGrid usage, 2,292 and 1,828 respectively. Similarly, Furlani et al. (2012, 2013)
 108 observed that molecular biosciences recently joined physics at the top of the list of sciences with the
 109 greatest CPU usage of Extreme Science and Engineering Discovery Environment (XSEDE), a virtual
 110 system providing digital resources and computing services (Towns et al. 2014). In addition, Furlani's
 111 longitudinal study demonstrated a substantial increase in the number of PIs using XSEDE from ~500
 112 in 2005 to ~1,600 in 2011. The authors noted, however, that often a project PI did not personally
 113 utilize the XSEDE resources, assigning computing tasks to graduate students or postdocs.

114 Other research focuses on various metrics for measuring academic performance with respect to HPC
 115 usage. Apon et al. (2010) suggest measuring the HPC investment in terms of “competitiveness,”
 116 presented as a ranking system. Ranks are calculated using the Top 500 HPC list that reports on the
 117 fastest 500 computers in the world. An institution’s rank is based on their investments in HPC
 118 according to the Top 500 HPC list. In their study, academic performance was characterized by
 119 publication counts and funding awards. Their funding showed that “consistent investments in HPC at
 120 even modest levels are strongly correlated to research competitiveness.” Apon et al.’s study also
 121 presented statistical evidence that NSF research funding and publication counts are good predictors
 122 of academic competitiveness. Subsequently, Knepper and Börner (2016) at the relationship between
 123 the CPU usage of XSEDE resources utilized by PIs and publication records. In addition, they mapped
 124 fields of science into HPC resources, creating a bipartite network. The results demonstrated that
 125 among the 27 top fields, physics, chemistry, astronomy, material research, and meteorology/climate
 126 modeling utilized XSEDE resources the most. They also observed that the type of HPC resources

127 plays an important role, as some systems link to almost all fields of science (e.g., NICS Kraken),
 128 whereas others serve only a small number of fields.

129 Newby et al. (2014) discussed various forms of “return on investment” in studies of research
 130 cyberinfrastructure. Their discussion included a number of different types of value derived from
 131 investment in advanced cyberinfrastructure, including scientific value, workforce development,
 132 economic value, and innovation. More recent work has adopted the financial definition of ROI: “a
 133 ratio that relates income generated … to the resources (or asset base) used to produce that income,”
 134 calculated typically as “income or some other measure of return on investment.” Values greater than
 135 1.0 indicate that return is greater than investment (Kinney and Raiborn 2011). One of the challenges
 136 in measuring ROI in financial terms is that returns may take decades to materialize in some
 137 disciplines (Stewart et al. 2015). Recently, two studies have applied ROI metrics to study the impact
 138 of HPC usage at Indiana University. Thota et al. (2016) compared the annual cost of operating IU’s
 139 Big Red II supercomputer with the funds brought into the university by the researchers that used Big
 140 Red II. The expected annual average cost for Big Red II is around \$15 million dollars. Thota et al.
 141 found that for the year 2013, the total grant income to IU by PIs or Co-PIs who make use of this
 142 supercomputer was more than double the amount of the cost (\$39.8 million). This is suggestive
 143 (although not conclusive) of a favorable financial ROI – only suggestive because the analysis of
 144 Thota et al. was not able to take into account how critical use of the Big Red II supercomputer was to
 145 the grants awarded to users of that system. Similarly suggestive, Fulton et al. (2017) showed a
 146 positive correlation between an increase in HPC usage and IU’s award funding over a period of
 147 years. Stewart et al. (2015) analyzed the ROI for XSEDE (the eXtreme Science and Engineering
 148 Discovery Environment) and argued that the ROI of federal investment in this resource was greater
 149 than 1.

150 **3 Data Acquisition and Preparation**

151 Four data sources are used to analyze and visualize the impact of internal IT usage on external
 152 funding, associated publications, and collaboration networks: IT usage data for faculty, staff, and
 153 students working at an institution; IU award database; NIH and NSF award data for the same
 154 institution together with publications that list these awards in the acknowledgements. These datasets,
 155 as well as their matching, cleaning, and preparation for visualization, are detailed subsequently. The
 156 overall process is illustrated in Figure 1. Please note that all public data and all code is available
 157 online at <http://cns.iu.edu/2017-Value-Analytics.html>.

158 ** Place Figure 1 here

159 **3.1 IT Usage Data**

160 XDMoD HPC resource usage log data is used to extract IT usage information. Among others, the
 161 logs contain five elements: IT system type, IT system name, units used (CPU hours for computing
 162 and Gigabytes for storage), user name, first name, and last name of IT user. An IT system consists of
 163 two types, namely storage and compute node, and several systems within each type.

164 From January 2012 to October 2017, there were two major storage and six computing systems
 165 utilized by (Co-)PIs for NSF and NIH grants at Indiana University.

166 Storage:

- 167 1. Scholarly Data Archive (SDA),
- 168 2. Data Capacitator 2 (DC2)

169 Computing:

- 170 1. Big Red, a supercomputer,
- 171 2. Big Red II, a supercomputer,
- 172 3. Karst, a cluster for serial jobs,
- 173 4. Mason, designed for data-intensive, high-performance computing tasks (Thota et al.
- 174 2016),
- 175 5. Quarry, a computing cluster, and
- 176 6. Carbonate, designed for data-intensive computing, particularly for genome and
- 177 phylogenetic software.

178 Among computing systems, Big Red was decommissioned in 2013 and replaced by Big Red II,
 179 Quarry was decommissioned in January 2015, and Carbonate became available in July 2017 to
 180 replace Mason, scheduled to retire on January 1, 2018.

181 HPC log files do not differentiate between a group and a single user account. Using unique IT user
 182 names and the fields with last and first names, 1,187 instances of group accounts were identified and
 183 removed. The log files were then filtered by the year > 2012-01-01 with the storage and computer
 184 resources > 0. Tables 1 and 2 provide a summary of IT resource usage at IU from January 2012 to
 185 October 2017 with a total of 65,495,233,153 CPU job hours run, 114,893 GB stored, and 4,112
 186 unique users active.

187
 188

189 **Table 1.** Summary of HPC computing jobs at IU by individual users (Jan 2012 - Oct 2017)

| Big Red | Carbonate | Mason | Quarry | Karst | Big Red II | Total CPU-Hours | Users |
|----------------|------------------|--------------|---------------|--------------|-------------------|------------------------|--------------|
| 119 | 134 | 635 | 699 | 1,298 | 1,311 | 65,495,233,153 | 4,197 |

190
 191 **Table 2.** Summary of HPC storage usage at IU by individual users (Jan 2012 - Oct 2017)

| DC2 | SDA | Total Gigabytes | Users |
|------------|------------|------------------------|--------------|
| 252 | 1,687 | 114,893 | 1,939 |

192

193 3.2 Internal IU Award Database

194 IU internal grant data is generated from Kuali Financial Services. This database imports user ID, the
 195 name of award agency, grant ID, and total amount. For this paper, the following query was specified:
 196 a) the start date is 2012-01-01 and b) the grant amount is greater than zero. Out of the total of 28,965
 197 awards between 2012-01-01 and 2018-01-01, 597 grants are from NSF and 2,677 grants are from
 198 NIH agencies. The number of unique (Co-)PI users are 425 for NSF and 690 for NIH.

199 3.3 NIH Grant and Publication Data

200 NIH grant award numbers and linkages to publications that cite award numbers can be downloaded
 201 in bulk using ExPORTER Data Catalog¹ or as a data extract using Research Portfolio Online
 202 Reporting Tools (RePORTER).² The NIH data provides access to both intramural and extramural
 203 NIH-funded research projects from the past 25 years and publications since 1980 (NIH 2017). For
 204 this study, we have used the RePORTER; however, our methodology is applicable to ExPORTER.
 205 When extracting data from the ExPORTER, NIH files are downloaded separately for each year and
 206 then merged, whereas the RePORTER output is already merged. The RePORTER query form also
 207 allows for generating data by means of query elements, such as keywords, organization names,
 208 publications, and others. By combining these elements, the user is able to create highly customized
 209 searches of this NIH funding database. In order to obtain the grant total and the number of
 210 publications for XDMoD, the following query filters were applied: project year (2012-2017),
 211 organizations listed in the NIH lookup (IU Bloomington, IU South Bend, and IUPUI), state (Indiana),
 212 and publication year (2012-2017), as illustrated in Figure 2. The *Agency/Institute/Center* field is kept
 213 with its default set to “admin,” and subprojects are set to be excluded.

214 ** Place Figure 2 here

215 The query was run on August 22, 2017, 9am EST and the query results were exported in CSV format
 216 with relevant fields: namely, *Project Number*, *Contact PI / Project Leader*, *Other PI or Project*
 217 *Leader(s)*, *FY Total Cost by IC*, and *Funding IC*. Publication data was exported with the following
 218 fields: *Core Project Number*, *ISSN*, *Journal*, *PMID*, *PUB Year*, and *Title*. The results comprised 933
 219 grants and 9,838 unique publications that acknowledge funding by these grants.

220 3.4 NSF Grant and Publication

221 NSF grants and associated publications were downloaded using the NSF Award Search Web API.³
 222 The API supports highly customized queries. For this project, a query was run on October 31, 2017,
 223 2pm EST, using the following filters:

- 224 1) *awardeeName*=“Indiana University”,
- 225 2) *startDateStart*=01/01/2012 and
- 226 3) *printFields*=id,publicationResearch,agency,startDate,expDate,fundProgramName,title,piF
 irstName,piLastName,estimatedTotalAmt,coPDPI,primaryProgram,
 awardeeCity,awardeeName.⁴

229 The results comprised 565 unique awards and 245 unique publications. Data was retrieved in JSON
 230 format and converted to a CSV format.

231 3.5 Data Preparation

232 **NIH Award–Publication Linkage.** NIH grants and NIH publications are linked via the project
 233 number. An additional preprocessing step is required for this linkage. The project number from the
 234 grant file is given in a 14-digit format (e.g., 1R01HS022681-01), whereas the publication file is

¹ https://exporter.nih.gov/ExPORTER_Catalog.aspx

² <https://projectreporter.nih.gov/reporter.cfm>

³ <https://www.nsf.gov/developer/>

⁴ <http://api.nsf.gov/services/v1/awards.json?awardeeName=%22Indiana+University%22&offset=26&startDateStart=01/01/2012&printFields=id,publicationResearch,agency,startDate,expDate,fundProgramName,title,piFirstName,piLastName,estimatedTotalAmt,coPDPI,primaryProgram,awardeeCity,awardeeName>

235 assigned an 11-digit format (e.g., R01HS022681). Combining these two files results in 2,046 unique
 236 publication records linked to 293 grants (11-digit format). *Funding and Publication Impact* data then
 237 consist of IT resources, funding agencies, publications, journals, and grant total, merged by (Co-)PI
 238 and project number. The *Co-PI Collaboration Network* includes the names of PIs and Co-PIs, grant
 239 funding total, and the number of grants awarded.

240 ***NSF Award–Publication Linkage.*** For each award, the NSF API data retrieves publications in the
 241 form of a list, which is then split into three fields: authors, publication title, and publication journal.
 242 PIs' first and last name fields were merged, yielding 318 authors, 565 awards, and 245 publications.
 243 Among these awards, 55 awards are associated with publications, and 48 (Co-)PIs are associated with
 244 these 55 awards.

245 ***IT User–IU Award Linkage*** between IT usage data and IU award data was performed using user
 246 IDs. Next, IU awards were linked to NIH and NSF award-publication linkages via project numbers.
 247 The result is a table that links 61 IT users to the very same number of (Co-)PIs with 83 project
 248 awards, and 968 associated publications based on the unique PMID identifier in the case of NIH
 249 award or the unique publication title for NSF awards, as the NSF API data does not provide
 250 publication-unique identifiers. As a result of this merge, the *Funding and Publication Impact* data
 251 include IT resources, funding agencies, publications, journals, and grant total, merged by PI and
 252 project number. The total number of awards associated with IT resources is 657 for NIH and 228 for
 253 NSF awards, totaling \$339,013,365. It should be noted that our main objective is to measure both
 254 financial and academic impacts. We have excluded the IU awards without publications and the
 255 awards for which the (Co-)PI did not use IT resources. As a result, the number of awards and their
 256 publications is lower than the total number of IU awards, thus totaling \$ 21,016,055 for 83 NSF and
 257 NIH awards.

258 ***Data Aggregation*** was performed to determine the number of users per IT resource, the total award
 259 amount per NIH Institute or Center (IC), and the total number of publications per discipline of
 260 science. Table 3 exhibits the total NIH and NSF award amount per IT storage and IT resource for
 261 unique project IDs. NIH ICs identify which Center for Scientific Review (CSR) reviewed the grant
 262 application for a funding decision (NIH Research Portfolio Online Reporting Tools 2017). In
 263 contrast, NSF API does not provide such a field. Table 3 shows the list of NIH IC and NSF together
 264 with the number of awards, publications, and total award amount for the IU dataset.

265 **Table 3.** NIH ICs and NSF funding grants with IU (Co-)PIs that use IT resources

| Funding Agencies | Number of funding awards | Number of publications | Sum of FY Total Cost by IC in \$ |
|------------------|--------------------------|------------------------|----------------------------------|
| NIH-NIGMS | 14 | 156 | 4,057,619 |
| NIH-NHLBI | 6 | 183 | 2,655,533 |
| NIH-NIAMS | 6 | 116 | 1,992,701 |
| NIH-NIAAA | 9 | 61 | 1,523,761 |
| NIH-NIMH | 6 | 37 | 1,423,969 |
| NIH-NCI | 7 | 97 | 1,338,995 |
| NIH-NEI | 3 | 33 | 1,057,979 |

| | | | |
|------------------|----|-----|------------|
| NIH-NIAID | 2 | 10 | 730,446 |
| NIH-NIDDK | 3 | 9 | 674,441 |
| NIH-OD | 2 | 2 | 606,848 |
| NIH-NIA | 4 | 49 | 587,810 |
| NIH-NIBIB | 1 | NA | 500,696 |
| NIH-NLM | 3 | 86 | 364,995 |
| NIH-NICHD | 1 | 56 | 332,956 |
| NIH-NCCIH | 1 | 18 | 299,149 |
| NIH-NIDA | 1 | 5 | 154,000 |
| NIH-NHGRI | 1 | 1 | 5,000 |
| NIH Total | 70 | 919 | 18,306,898 |
| NSF | 13 | 49 | 2,709,157 |
| Total | 83 | 968 | 21,016,055 |

266

267 Publication records were aggregated using the UCSD map of science (Börner et al. 2012), a
 268 classification system that assigns each journal to one or more subdisciplines of science that are
 269 further aggregated into 13 disciplines of science (e.g., mathematics or biology).⁵ It should be noted
 270 that some journal names retrieved from NIH and NSF vary considerably from the UCSD
 271 classification system; several preprocessing steps are necessary, such as lowering cases and
 272 normalizing punctuation. There were 245 cases where publication records did not match the UCSD
 273 map of science dictionary, and 33 publications were associated with more than one discipline with
 274 the same relative association proportion. For these publications, we created two additional categories,
 275 “Unclassified” and “Multidisciplinary.” As a result, each publication is associated with one of the 13
 276 disciplines of science. The number of publications per discipline is given in Table 4.

277

Table 4. Number of publication per discipline associated with the NIH and NSF awards

| Discipline | Number of Publications |
|----------------------|-------------------------------|
| Medical Specialties | 252 |
| Unclassified | 245 |
| Brain Research | 107 |
| Infectious Diseases | 82 |
| Health Professionals | 78 |
| Biotechnology | 59 |
| Chemistry | 44 |
| Multidisciplinary | 33 |
| Social Sciences | 30 |
| Biology | 23 |

⁵ <http://cns.iu.edu/2012-UCSDMap.html>

| | |
|---|------------|
| Math & Physics | 7 |
| Chemical, Mechanical, & Civil Engineering | 5 |
| Electrical Engineering & Computer Science | 2 |
| Earth Sciences | 1 |
| Total | 968 |

Finally, all aggregated data and linkage tables were converted to JSON format as required by the visualization plugin in the XDMoD VA portal. The conversion script is available at <http://cns.iu.edu/2017-Value-Analytics.html>. The JSON format specification for the *Funding and Publication Impact* and *Co-PI Collaboration Network* visualizations is illustrated in Figures 3 and 4.

** Place Figure 3 here

** Place Figure 4 here

4 Methods

4.1 Collaboration Network Extraction

The *Co-PI Collaboration Network* was extracted from NIH grant data as described in Section 3. Co-PIs were calculated by first matching the IU grant database's 'Agency Award Number' to the NIH 'Core Project Number', and then splitting the NIH 'Other PI or Project Leader(s)' field by semicolon and counting the number of grants and total grants for each co-author pair to compute the weight for collaboration edges. For NSF, the Co-PI information was extracted from the field 'coPDPI,' which also includes Co-PI's IDs. These IDs were removed and the field was split by comma.

4.2 Sankey Graphs and Force Network Layout in WVF

Sankey Graphs show the magnitude of flow between nodes in a network as well as the relationship between flows and their transformation (Riehmann et al. 2017). For the *Funding and Publication Impact* graph, we used the D3 Sankey API.⁶ It reads input nodes and weighted links and computes positions using the Gauss–Seidel iterative method (Barrett et al. 1994). First, the horizontal position of the left-most nodes is computed; then, nodes further on right are positioned while minimizing link distance. After all the nodes are positioned, a reverse pass is made from right-to-left, and overlapping nodes are moved to minimize collision. The entire process is repeated several times to optimize the layout.

The final visualization features three types of nodes, namely IT resources on left, funding (e.g., NIH institutes and NSF) in middle, and publication disciplines on right. The height of a bar, or node, is proportional to the maximum of the weighted sum of incoming links and the weighted sum of outgoing links. The nodes are placed in ascending order by their heights.

Force-Directed Graphs are used to display the relationship between objects by calculating the position of each node based on their shared edges. The D3 force-directed graph⁷ applies three primary forces upon the nodes: namely, the sum of all forces, a force between two linked nodes, and a central force using the layout algorithm by Dwyer et al. (2006). The WVF applies *linkStrength* as a

⁶ <https://bost.ocks.org/mike/sankey/>

9

⁷ <https://bl.ocks.org/mbostock/4062045>

309 parameter to calculate node positions and it is constant for all nodes.⁸ In the *Co-PI Collaboration*
 310 *Network* visualization, the nodes are Co-PIs and edges represent their research collaborations. Node
 311 size corresponds to the number of grants and node color denotes the total funding amount in U.S.
 312 dollars per PI. Edge thickness indicates the number of co-authored grants.

313 **The Web Visualization Framework (WVF)** was used to build both visualizations.⁹ WVF is a highly
 314 configurable packaging of several industry-standard web libraries (Angular, D3, HeadJS, Bootstrap,
 315 and many others) that allows its users to quickly build visualization applications. Existing WVF
 316 visualizations support rendering interactive horizontal bar graphs, geospatial maps, network graphs,
 317 bimodal graphs, science maps, and others. The WVF provides lightweight in-browser aggregation
 318 and analysis, but it relies on web services or external data sources to provide primary analyses.
 319 Applications built using WVF plugins allow each visualization within a page to use data from and
 320 interact with other visualization plugin elements, but through loosely coupled data connections. This
 321 allows for both the visualization elements and the aggregation and filtering methods to be replaced or
 322 removed without affecting other elements.

323 5 Results

324 This section explains how the data and methods discussed above are applied to IU institutional data
 325 and what insights were gained. Specifically, we describe the *Funding and Publication Impact* and
 326 *Co-PI Collaboration Network* visualizations, as well as the portal that supports easy access to both.
 327 The visualizations aim to help stakeholders understand the financial ROI measured in terms of total
 328 acquired funding and academic ROI measured by publications associated with these awards.

329 5.1 XDMoD VA Portal

330 The XDMoD portal is an interactive dashboard with an intuitive graphical interface to XDMoD
 331 metrics such as number of jobs, service units charged, CPUs used, or wait time (Furlani et al. 2012).
 332 XDMoD metrics can be broken down by field of science, institution, job size, principal investigator,
 333 and resource. Academic metrics (e.g., publications, citations, and external funding) can be uploaded
 334 by users or incorporated by institutions via Open XDMoD (Fulton et al. 2017). XDMoD Value
 335 Analytics adds new functionality by offering metrics on financial and scientific impact via
 336 visualization plugins, as illustrated in Figure 5. Key features of the VA interface include the ability to
 337 interact and drill-down, allowing users to access additional related information simply by clicking
 338 inside edges and nodes or selecting the desired filters.

339 ** Place Figure 5 here

340 5.2 XDMoD Funding and Publication Impact Visualization

341 The *Funding and Publication Impact* visualization allows users to interactively explore the relations
 342 between IT resource usage (on left in Figure 6), funding awards aggregated by NIH institute and NSF
 343 (in middle), and publications that cite this funding aggregated by scientific discipline (on right).
 344 Sankey graph links take users on an exploratory quest, moving from the IT resources via funding to
 345 papers published in diverse scientific disciplines.

346 ** Place Figure 6 here

⁸ <https://github.com/d3/d3-3.x-api-reference/blob/master/Force-Layout.md>

⁹ <https://github.com/cns-iu/WVF>

347 In this interactive visualization, users are provided with various functionalities on demand, such as
 348 mouse-over and selection; a legend explaining color and size coding can be viewed on demand. For
 349 example, hovering over a particular node will cause that node and all links emanating from it to be
 350 highlighted, whereas hovering over a particular link will highlight that link with the color of the node
 351 from which the link originated, as illustrated in Figure 6. In this example, the user wishes to explore
 352 the connections between IT resource use, grant funding, and the number of publications in the field
 353 of *Brain Research* from *Scientific Discipline* (3). After hovering over one of the links connected to
 354 the *Brain Research* node from the *Scientific Discipline* category on the right, the common link will
 355 connect *Brain Research* publications with the funding agency from the *Funding* category (2) and *IT*
 356 *Resources* (1). Hovering over a link also brings up additional information relevant to the node (e.g.,
 357 number of papers, here 54). To explore other links and nodes, the user can simply double-click on
 358 them to reset. The example in Figure 6 also illustrates how the user may gain insights from the visual
 359 data. In particular, the user's selection reveals that 54 out of 107 papers in *Brain Research*
 360 acknowledge *NIH-NLM* funding and that the grants associated with these papers utilized *Karst*, as a
 361 computing IT resource.

362 **Interpretation:** The visualization shows that in 2012-2017, a total of 114,894 GB and
 363 65,495,118,259 CPU hours were used by externally funded projects that had associated publications.
 364 The grant income to IU by PIs and Co-PIs which use IU HPC resources—for researchers who had
 365 both awards and publications during the period analyzed—is \$21,016,055. A majority of this
 366 funding comes from NIH projects that total \$18,306,898 (87%) with the top three ICs being NIH-
 367 NIGMS, NIH-NHLBI, and NIH-NIAMS. In terms of publications, IT resource usage, via grants,
 368 links to 968 publications. *Brain Research*, *Medical Specialties*, and *Infectious Diseases* have the
 369 largest number of publications.

370 5.3 XDMoD VA Co-PI Network Visualization

371 The *Co-PI Collaboration Network* visualization is shown in Figure 7. It features a force-directed
 372 network layout on the left and a sorted horizontal bar graph on the right. Both visualizations are
 373 coupled so that hovering over an investigator in the network highlights that same investigator in the
 374 bar graph. The node size for each investigator indicates the number of grants received, while the node
 375 color indicates the total award amount; see legend in interactive visualization for details. If two
 376 investigators collaborated (i.e., their names are listed together on a grant), there exists an edge
 377 between them. The thickness of this edge represents the number of times they collaborated together.
 378 To see collaborations in the network, simply hover over the node of a particular investigator. This
 379 will highlight the selected investigator node and all emanating edges leading to other collaborator
 380 nodes in that investigator's network. Similarly, hovering over a bar in the bar graph highlights all
 381 corresponding entries and renders other bars opaque for easy viewing of the selection. The range
 382 filter on the top (1) can be used to increase or decrease the number of node labels in the network
 383 visualization. The plus and minus buttons in the top left can be used to zoom in and out. During
 384 zooming, the legend is automatically updated to ensure that node values and edge thickness remain
 385 accurate.

386 ** Place Figure 7 here

387 **Interpretation:** The visualization helps identify three key elements of the academic and financial
 388 impact: namely, the number of awards, their total dollar amount, and research collaborations.
 389 Collaborations are rendered as a network with nodes representing researchers and edges denoting
 390 their Co-PI relationships. Given the rather short time frame, there are many, relatively small

391 collaboration clusters. Most links are thin, indicating a one-time collaboration; there are few
 392 instances of multiple collaborations denoted by a thicker edge between nodes. A slider (1) filters
 393 labels by the number of grants. The legend (4) provides additional insights on the number of grants,
 394 their total amount, and the number of co-authored grants. The *Total Amount* column on the right (2)
 395 shows researchers sorted by total funding during the years 2012-2017. By selecting a (Co-)PI bar (2),
 396 the collaboration network for that (Co-)PI is highlighted (3).

397 6 Discussion and Outlook

398 The work presented in this paper aims to help researchers, administrators, and funders understand
 399 and communicate the impact of campus cyberinfrastructure investments on scholarly productivity in
 400 terms of funding intake, publication output, and scholarly networks. The visualizations enable
 401 academic institutions to better understand return on investment (ROI) on advanced
 402 cyberinfrastructure for different types of research (e.g., as expressed by NIH ICs and disciplines of
 403 science). As part of the work, we demonstrated different methods for collecting and processing
 404 publicly available data from NIH and NSF official sites and from institutional production systems
 405 that advance the functionality of the XDMoD VA portal.

406 Expanding on the work by Knepper and Börner (2016), we primarily focused on the relationship
 407 between storage and computing resources utilized by (Co-)PIs and associated funding and
 408 publication records. Grant income to the university by (Co-)PIs who used IT resources during the
 409 period analyzed was \$339,013,365 for 885 NIH and NSF projects and grant income from (Co-)PIs
 410 who used IT resources and had both grant awards and publications was \$21,016,055. A total of 968
 411 publications were associated with 83 of these NSF and NIH awards. In addition, the results show that
 412 *Brain Research*, *Medical Specialties*, and *Infectious Diseases* are the top three scientific disciplines
 413 ranked by their publication records during the given time period. Note that only awards associated
 414 with publications and IT resources are displayed; and only funding from two agencies, namely NIH
 415 and NSF, is shown.

416 In the future, we plan to advance the presented work as follows:

- 417 • Institutions will be able to upload not only IT compute cycles and storage usage counts but
 418 also information on storage size and number of compute cycles to provide additional insights
 419 into usage patterns across scientific disciplines.
- 420 • Funding data will be automatically retrieved via NIH RePORTER and NSF APIs, reducing
 421 the amount of manual work involved. Both online resources can be queried periodically to
 422 update award and publication data. Data from other funding agencies might be added as well.
- 423 • Fuzzy matching algorithms will be implemented to increase the number of journals mapped
 424 to scientific disciplines. This will help reduce the number of publications designated as
 425 *Unclassified*.

426 As Fulton et al. (2017) state, “measuring intellectual outcomes is difficult, particularly since the
 427 results of intellectual accomplishments may take years or decades to be fully realized.” To evaluate
 428 data quality and data matching (e.g., by PI name), we are working on a comparison of data retrieved
 429 from NSF/NIH versus data available via IU’s Sponsored Research production databases. Results will
 430 help understand data issues and optimize matching algorithms. Understanding the value of
 431 investment in cyberinfrastructure is challenging, as the impact of such investments has many
 432 dimensions, including intellectual contributions and financial impact. The XDMoD VA modules
 433 facilitate understanding of the role cyberinfrastructure by analyzing a number of metrics and allowing
 434 visualization of the diverse ways in which they impact institutional planning and strategy as well as
 435 the development of human knowledge.

436 **Acknowledgements**

437 We would like to thank Tom Furlani and his team at the University at Buffalo for their leadership in
 438 developing XDMoD; Winona Snapp-Childs, and Robert Henschel of UITS at IU for their expert
 439 input to the design and implementation of XDMoD VA; and Todd Theriault and Craig A. Stewart for
 440 edits and comments on an earlier draft of this paper. This work was partially funded by the National
 441 Science Foundation under grants 1053575 and 1566393, and also supported by the IU Pervasive
 442 Technology Institute. Any opinions, findings, and conclusions or recommendations expressed in this
 443 material are those of the authors and do not necessarily reflect the views of the National Science
 444 Foundation.

445 **Conflict of Interest**

446 *The authors declare that the research was conducted in the absence of any commercial or financial
 447 relationships that could be construed as a potential conflict of interest.*

448 **Author Contributions**

449 Performed data preparation: OS, BF. Analysed data: OS, KB. Designed visualization: GS, SB, KB.
 450 Wrote the paper: OS, KB, GS, SB, SH, MRL.

451 **References**

- 452 Apon, Amy, Stanley Ahalt, Vijay Dantuluri, Constantin Gurdgiev, Moez Limayem, Linh Ngo, and
 453 Michael Stealey. 2010. “High Performance Computing Instrumentation and Research
 454 Productivity in U.S. Universities.” *Journal of Information Technology Impact* 10 (2): 87–98.
 455 https://papers.ssrn.com/sol3/papers.cfm?abstract_id=1679248.
- 456 Barrett, Richard, Michael Berry, Tony F Chan, James Demmel, June M Donato, Jack Dongarra,
 457 Victor Eijkhout, Roldan Pozo, Charles Romine, and Henk Van Der Vorst. 1994. *Templates for
 458 the Solution of Linear Systems: Building Blocks for Iterative Methods*. 2nd ed. Philadelphia, PA:
 459 SIAM. <http://www.netlib.org/templates/templates.pdf>.
- 460 Börner, Katy, Richard Klavans, Michael Patek, Angela M. Zoss, Joseph R. Biberstine, Robert P.
 461 Light, Vincent Larivière, and Kevin W. Boyack. 2012. “Design and Update of a Classification
 462 System: The UCSD Map of Science.” Edited by Neil R. Smalheiser. *PLoS ONE* 7 (7). Public
 463 Library of Science: e39464. doi:10.1371/journal.pone.0039464.
- 464 Dwyer, Tim, Kim Marriott, and Michael Wybrow. 2006. “Integrating Edge Routing into Force-
 465 Directed Layout.” Edited by Kaufmann, Michael and Wagner, Dorothea. In *Graph Drawing*, 8–
 466 19. Berlin: Springer. doi:10.1007/978-3-540-70904-6_3.
- 467 Fulton, Ben, Steven Gallo, Robert Henschel, Tom Yearke, Katy Börner, Robert L. DeLeon, Thomas
 468 R. Furlani, Craig A. Stewart, and Matthew R. Link. 2017. “XDMoD Value Analytics.” In
 469 *Proceedings of the Practice and Experience in Advanced Research Computing 2017 on
 470 Sustainability, Success and Impact - PEARC17*, 1–7. doi:10.1145/3093338.3093358.
- 471 Furlani, Thomas R., Ryan J. Gentner, Abani K. Patra, Gregor von Laszewski, Fugang Wang, Jeffrey
 472 T. Palmer, Nikolay Simakov, et al. 2013. “Using XDMoD to Facilitate XSEDE Operations,
 473 Planning and Analysis.” In *Proceedings of the Conference on Extreme Science and Engineering
 474 Discovery Environment Gateway to Discovery - XSEDE ’13*. New York, NY: ACM Press.
 475 doi:10.1145/2484762.2484763.
- 476 Furlani, Thomas R., Matthew D. Jones, Steven M. Gallo, Andrew E. Bruno, Charng-Da Lu, Amin
 477 Ghadersohi, Ryan J. Gentner, et al. 2013. “Performance Metrics and Auditing Framework Using
 478 Application Kernels for High-Performance Computer Systems.” *Concurrency and*

- 479 *Computation: Practice and Experience* 25 (7): 918–31. doi:10.1002/cpe.2871.
- 480 Furlani, Thomas R., Barry I Schneider, Matthew D Jones, John Towns, David L Hart, Abani K Patra,
 481 Robert L Deleon, et al. 2012. “Data Analytics Driven Cyberinfrastructure Operations, Planning
 482 and Analysis Using XDMoD.” In *Proceedings of the SC12 Conference, Salt Lake City, Utah*.
 483 <https://laszewski.github.io/papers/vonLaszewski-draft-data-analytics-planing.pdf>.
- 484 Hart, David L. 2011. “Measuring TeraGrid: Workload Characterization for a High-Performance
 485 Computing Federation.” *The International Journal of High Performance Computing
 486 Applications* 25 (4): 451–65. doi:10.1177/1094342010394382.
- 487 Iosup, Alexandru, Catalin Dumitrescu, Dick Epema, Hui Li, and Lex Wolters. 2006. “How Are Real
 488 Grids Used? The Analysis of Four Grid Traces and Its Implications.” In *2006 7th IEEE/ACM
 489 International Conference on Grid Computing*, 262–69. Washington, DC: IEEE Computer
 490 Society Press. doi:10.1109/ICGRID.2006.311024.
- 491 Katz, Daniel S., David Hart, Chris Jordan, Amit Majumdar, J.P. Navarro, Warren Smith, John
 492 Towns, Von Welch, and Nancy Wilkins-Diehr. 2011. “Cyberinfrastructure Usage Modalities on
 493 the TeraGrid.” In *Proceedings of the 2011 IEEE International Symposium on Parallel and
 494 Distributed Processing Workshops and Phd Forum*, 932–39. Washington, DC: IEEE Computer
 495 Society Press. doi:10.1109/IPDPS.2011.239.
- 496 Kinney, M.F. and C.A. Raiborn. 2011. Cost Accounting. South-Western, Mason, OH. ISBN: 978-1-
 497 111-97172.
- 498 Knepper, Richard. 2011. “The shape of the TeraGrid.” In *Proceedings of the 2011 TeraGrid
 499 Conference on Extreme Digital Discovery - TG '11*. New York, NY: ACM Press.
- 500 Knepper, Richard, and Katy Börner. 2016. “Comparing the Consumption of CPU Hours with
 501 Scientific Output for the Extreme Science and Engineering Discovery Environment (XSEDE).”
 502 *PloS ONE* 11 (6). doi:10.1371/journal.pone.0157628.
- 503 Lee, Bu-Sung, Ming Tang, Junwei Zhang, Ong Yew Soon, C. Zheng, P. Arzberger, and D.
 504 Abramson. 2006. “Analysis of Jobs in a Multi-Organizational Grid Test-Bed.” In *Proceedings
 505 of the Sixth IEEE International Symposium on Cluster Computing and the Grid (CCGRID'06)*,
 506 59–59. Washington, DC: IEEE Computer Society Press. doi:10.1109/CCGRID.2006.1630950.
- 507 Li, Hui, David Groep, and Lex Wolters. 2005. “Workload Characteristics of a Multi-Cluster
 508 Supercomputer.” In *Lecture Notes in Computer Science*, edited by D Feitelson, L Rudolph, and
 509 U Schwiegelshohn, 176–93. Berlin: Springer. doi:10.1007/11407522_10.
- 510 Madhavan, Krishna, Niklas Elmquist, Mihaela Vorvoreanu, Xin Chen, Yueting Wong, Hanjun Xian,
 511 Zhihua Dong, and Aditya Johri. 2014. “DIA2: Web-Based Cyberinfrastructure for Visual
 512 Analysis of Funding Portfolios.” *IEEE Transactions on Visualization and Computer Graphics*
 513 20 (12): 1823–32. doi:10.1109/TVCG.2014.2346747.
- 514 Newby, Greg, Amy Apon, Nick Berente, Rudolph Eigenmann, Susan Fratkin, David Lifka, and Craig
 515 A. Stewart. 2014. “Return on Investment from Academic Supercomputing.” In *Presentation at
 516 SC14 (Supercomputing 2014)*. <https://scholarworks.iu.edu/dspace/handle/2022/19242>.
- 517 NIH. 2017. “RePORTER User Manual.” Research Portfolio Online Reporting Tools. Retrieved from
 518 https://projectreporter.nih.gov/RePORTER_Manual_files/RePORTERManual.pdf
- 519 Palmer, Jeffrey T, Steven M Gallo, Thomas R. Furlani, Matthew D Jones, Robert L Deleon, Joseph P
 520 White, Nikolay Simakov, Abani K. Patra, Jeanette Sperhac, Thomas Yearke, Ryan Rathsam,
 521 Martins Innus, Cynthia D. Cornelius, James C. Browne, William L. Barth, and Richard T.
 522 Evans. 2015. “Open XDMoD: A Tool for the Comprehensive Management of High-
 523 Performance Computing Resources.” *Computing in Science & Engineering*.
 524 http://www.buffalo.edu/content/dam/www/cer/pdfs/OpenXDMoD_preprint.pdf.
- 525 Riehmann, Patrick, Manfred Hanfler, and Bernd Froehlich. 2017. “Interactive Sankey Diagrams.” In
 526 *IEEE Symposium on Information Visualization*, 2005. INFOVIS 2005., 233–40. IEEE.
 527 Accessed August 24. doi:10.1109/INFVIS.2005.1532152.

- 528 Shneiderman, Ben. 1996. "The Eyes Have It: A Task by Data Type Taxonomy for Information
 529 Visualizations." In *Proceedings of the IEEE Symposium on Visual Language*, 336–43.
 530 Washington, DC: IEEE Computer Society Press. doi:10.1109/VL.1996.545307.
- 531 Stewart, Craig A., Ralph Roskies, Richard Knepper, Richard L. Moore, Justin Whitt, and Timothy
 532 M. Cockerill. 2015. "XSEDE Value Added, Cost Avoidance, and Return on Investment." In
 533 *Proceedings of the 2015 XSEDE Conference on Scientific Advancements Enabled by Enhanced*
 534 *Cyberinfrastructure - XSEDE '15*, 1–8. New York, NY: ACM Press.
 535 doi:10.1145/2792745.2792768.
- 536 Stewart, Craig A., Richard D. Knepper, Matthew R. Link, Marlon Pierce, Eric A. Wernert, and
 537 Nancy Wilkins-Diehr. 2017. Cyberinfrastructure, Science Gateways, Campus Bridging, and
 538 Cloud Computing. In *Encyclopedia of Information Science and Technology*, Third Edition. Vol.
 539 IX. 2014. Hershey, PA. <http://www.igi-global.com>. Available from:
 540 <http://hdl.handle.net/2022/18608>
- 541 Thota, Abhinav S., Ben Fulton, Le Mai Weakley Weakley, Robert Henschel, David Y. Hancock,
 542 Matt Allen, Jenett Tillotson, Matt Link, and Craig A. Stewart. 2016. "A PetaFLOPS
 543 Supercomputer as a Campus Resource." In *Proceedings of the 2016 ACM on SIGUCCS Annual*
 544 *Conference - SIGUCCS '16*, 61–68. New York, NY: ACM Press.
 545 doi:10.1145/2974927.2974956.
- 546 Towns, John, Timothy Cockerill, Maytal Dahan, Ian Foster, Kelly Gaither, Andrew Grimshaw,
 547 Victor Hazlewood, et al. 2014. "XSEDE: Accelerating Scientific Discovery." *Computing in*
 548 *Science and Engineering* 16 (5): 62–74. Washington, DC: IEEE Computer Society Press.
 549 doi:10.1109/MCSE.2014.80.
- 550

551 **Figure Captions**

552 **Figure 1.** Data sources, data linkage, data preparation for XDMoD VA visualizations

553 **Figure 2.** NIH RePORTER online search query

554 **Figure 3.** JSON data schema for the *Funding and Publication Impact* visualization plugin

555 **Figure 4.** JSON data schema for the *Co-PI Collaboration Network* visualization plugin

556 **Figure 5.** XDMoD VA Portal, see interactive version at <http://demo.cns.iu.edu/xdmod-p/portal.html>

557 **Figure 6.** XDMoD VA *Funding and Publication Impact Visualization*, see interactive version at
 558 <http://demo.cns.iu.edu/xdmod-p/impact.html>

559 **Figure 7.** XDMoD VA *Co-PI Network Visualization*, see interactive version at
 560 <http://demo.cns.iu.edu/xdmod-p/co-pi.html>

561