# Towards a Scientific Impact Measuring Framework for Large Computing Facilities - a Case Study on XSEDE

Fugang Wang
Indiana University
2719 10th Street
Bloomington, Indiana, U.S.A.

Gregor von Laszewski[*]
Indiana University
2719 10th Street
Bloomington, Indiana, U.S.A.
laszewski@gmail.com

Geoffrey C. Fox
Indiana University
2719 10th Street
Bloomington, Indiana, U.S.A.

Thomas R. Furlani
Center for Computational
Research
University at Buffalo, SUNY
701 Ellicott Street
Buffalo, New York, 14203

Robert L. DeLeon
Center for Computational
Research
University at Buffalo, SUNY
701 Ellicott Street
Buffalo, New York, 14203

Steven M. Gallo
Center for Computational
Research
University at Buffalo, SUNY
701 Ellicott Street
Buffalo, New York, 14203

## ABSTRACT

We present a framework that (a) integrates publication and citation data retrieval, (b) allows scientific impact metrics generation at different aggregation levels, and (c) provides correlation analysis of impact metrics based on publication and citation data with resource allocation for a computing facility. Furthermore, we use this framework to conduct a scientific impact metrics evaluation of XSEDE. We carry out an extensive statistical analysis correlating XSEDE allocation size to the impact metrics aggregated by project and field of science. This analysis not only helps to provide an indication of XSEDE's scientific impact, but also provides insight regarding maximizing the return on investment in terms of allocation by taking into account the field of science or project based impact metrics. The findings from this analysis can be utilized by the XSEDE resource allocation committee to help assess and identify projects with higher scientific impact. It can also help provide metrics regarding the return on investment for XSEDE resources, or other institutional or campus resources for which an analysis of impact based on publications is important.

## Categories and Subject Descriptors

H.4 [**Information Systems Applications**]: Miscellaneous; D.2.8 [**Software Engineering**]: Metrics—*complexity measures, performance measures*

## General Terms

Theory, Measurement

---

[*]Corresponding Author.

## Keywords

Scientific impact, bibliometric, h-index, Technology Audit Service, XDMoD, XSEDE

## 1. INTRODUCTION

It is a well-known fact that many science and engineering innovations and discoveries are increasingly dependent on access to high performance computing resources. For many researchers, this demand is met by large-scale compute resources that cannot typically be supported by any single research group. Accordingly, dedicated large-scale computing facilities play an important role in scientific research, in which resources are shared among groups of researchers, while the facilities themselves are managed by dedicated staff. Indeed, the National Science Foundation and the Department of Energy have supported such facilities for many years. One such facility is the Extreme Science and Discovery Environment (XSEDE) [10]. XSEDE allocates resources to approved projects, which represent a substantial financial investment by NSF. Thus, justification for their use is warranted and questions regarding the scientific impact of these resources naturally arise, including:

1. Is there a way to measure the impact that such facilities provide to scientific research?

2. Is there a correlation between the size of a given allocation and the scientific impact of an individual user, a given project, or a field of science?

3. When evaluating a proposal request, what is the criteria to judge whether the proposal has the potential to lead to impactful research, and how does one obtain metrics to substantiate this?

To answer these questions, first we need a process to quantify the scientific outcome for the individual researchers. Secondly, we need to define and generate metrics to measure the scientific impact for individual researchers and higher level aggregated entities. Finally, we correlate the impact metrics to the consumed resources, to provide insight on how the computing facility benefits and impacts the science conducted utilizing its resources. In this paper, we present a

framework that addresses these questions. It is important to point out that measuring scientific impact can be quite controversial and that the presented results do not necessarily represent an absolute measure of the impact of a scientific project, but rather the results we present represent one of many factors that together define the scientific impact.

Furthermore, while we have restricted our analysis of scientific impact as it relates to XSEDE, the work presented here has general applicability to not only HPC resources, but to other resources including campus based HPC centers, beam lines, and other expensive equipment.

In particular, we focus our effort to identify impact based on scientific publications as the base unit of the research productivity and obtain data, as well as derive various metrics based on publication data to measure the impact of individual users, projects, Field of Science (FOS), and XSEDE itself as a whole.

In the following sections we briefly review related work (Section 2), then present our designed framework (Section 3) and implementation details (Section 4). Then, we discuss results and their impact (Section 5). Finally, we outline ongoing activities, indicate future plans (Section 6), and provide a summary (Section 7).

## 2. RELATED WORK

Our choice of using publication as the basic unit to measure the scientific impact is supported by the fact that bibliometrics based criteria is one of the de-facto standards to measure the impact of research. For example, publication derived metrics are broadly used in faculty recruit/promotion, and institutional rankings [19]. While usage based metrics are proposed by some [12, 14, 13], citation based metrics are a widely accepted measure. For instance, nanoHub uses publication and citation derived metrics to measure the impact of their project [7].

In addition to the intuitive measures like number of publications and number of citations, h-index [17] and g-index [15] are two other popular metrics. The publication count is often related to a measure of the productivity, while citation counts are often related to the quality, or impact of the work published. As h-index and g-index calculate the metric by combining this data, they measure both the productivity and the quality, thus providing a general measure for impact.

There are existing tools to measure the metrics for individual users, e.g. Scholarometer [18] and Publish or Perish [8]. These could be potentially leveraged to analyze a relatively small group of users, e.g., the work [11] showing TeraGrid's impact based on limited data from one resource allocation meeting consisting of only 112 selected PIs. However, neither of the tools provides a scalable solution to the large community we are concerned with here, namely the over 20,000 users who have utilized TeraGrid/XSEDE resources. While more formal publication based metrics, either based on citation or usage, are still the most widely employed criteria,

Other non publication based metrics have been proposed by altmetrics [1] while also considering measures for datasets and code; as well as mentioning of a snippet of work via social networking. We acknowledge these efforts willl be useful as they correlate other usage, however at this time we still lack standards and a well-established way to objectively derive scientific impact from the multitude of data sources.

Furthermore, often we do not have reliable data available.

## 3. SYSTEM DESIGN

We have designed a software framework to support measuring scientific impact via a publication and citation based approach. The framework is based on distributed set of services. The service-oriented system consists of components for (a) publication and citation data retrieval (e.g., from NSF award database, Google Scholar, and ISI Web of Science), (b) parsing and processing while correlating data from various databases and services, such as the XSEDE central database (XDcDB), which stores all usage data for jobs run on XSEDE resources, and (c) the Partnerships Online Proposal System (POPS) database, which stores publication and grant funding information for PI's applying for XSEDE allocations. The system also includes components for metrics generation and an analysis system for different aggregation levels (users, projects, organization, Field of Science), as well as a presentation layer using a lightweight portal in addition to exposing some data via RESTful API.

Fig 1 shows the layered system architecture, with an emphasis on the relationships between related components especially those integrating with databases. On the core **App** layer we have the database mining and publication mashup components. The database mining component generates XSEDE user specific publication data as well as, user, project, and Field of Science (FOS) views. The publication mashup component aggregates the publication data mined from the previous component, in addition to those from XDcDB, and from other available external services. It also retrieves citation data for each publication from external services. Another essential task of this component is to generate metrics for users, projects, and FOS in which the POPS database is involved to get proposal and project data. This data will be stored into the mashup database which can then be integrated into the XDMoD [16] system at our partner site at University of Buffalo. We also expose some data and analysis results via RESTful API and a portal as denoted on the **Service/GUI** layer. The **Data** layer illustrates the databases involved. The **External Resource/Services** layer lists the third party resource and services that we are currently using or have experimented or plan to investigate.

To conduct the analysis the general workflow includes obtaining the publication data for each XSEDE user, and then retrieving the citation data for each publication. Hence, the data is originally collected per user and per publication basis. As part of processing the data we are aggregating it based on organization, XSEDE project/account, and FOS. By correlating the data (for example the Service Units (SU) awarded by XSEDE) our intention is to identify if the analysis may reveal patterns and trends of how XSEDE can impact the sciences and possibly helps to achieve a better measure of return on investment (ROI) for NSF. While we are using the system to analyze the scientific impact of XSEDE, the framework itself is flexible enough so that it could be easily adapted to other similar systems for impact measure and analyses.

## 4. IMPLEMENTATION

We have implemented the system following best practices and leveraging popular tools and frameworks. The core system is developed in Python and various libraries are utilized
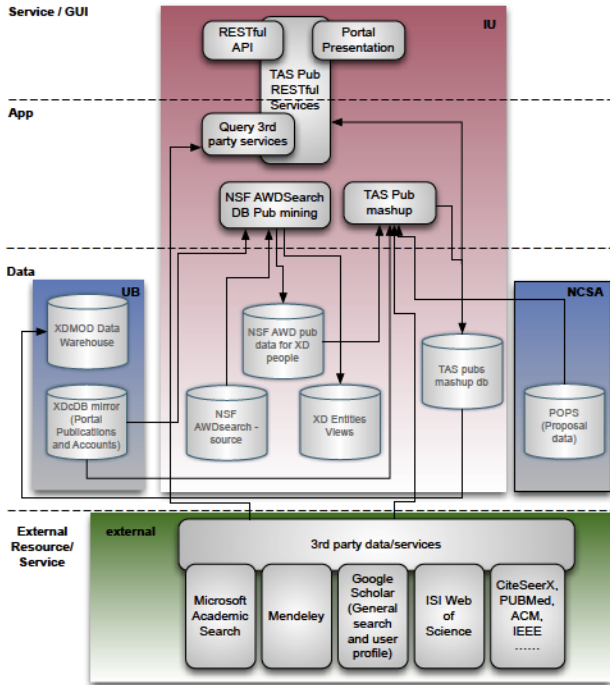
Figure 1: The Architecture of the Framework

to help interact with various databases, developing service and web frameworks, and so on. The portal/web tier follows the Ajax approach which provides better experience to users while viewing the presented data and chart.

Publication and citation data retrieval was a complex but essential part of our study, so we provide details of this process next.

## 4.1 Publication Data Acquisition



(a) Number of Publications  (b) Distribution of number
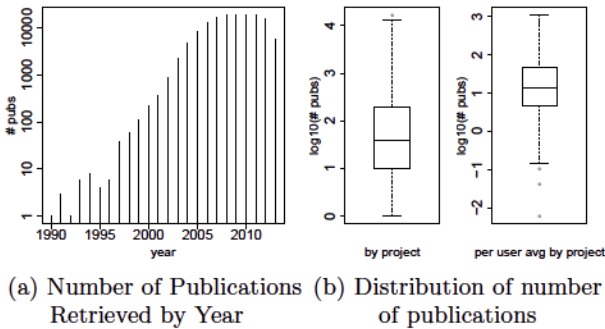Retrieved by Year              of publications

Figure 2: Distribution of the Publication Data

Given the size of the XSEDE user database, which as of Jan 2014 was over 20,000 users, we needed to employ an automated approach to obtain publication data on behalf of each user. Publication citation data are available via subscribed resources such as ISI Web of Science [4] or open access such as Google Scholar [2], Microsoft Academic Search [6], and Mendeley [5], however they unfortunately usually do not provide unlimited access, making automated publication retrieval impractical for some services. Another

approach is to obtain the publication data directly from the users. This is desirable since user curated data tends to be more accurate in comparison to automated publication mining. Additionally, it can provide extra information regarding a publication's association with the system, e.g., to which project a given publication is associated with. XSEDE provides such functionality via the user portal. However, this framework was just recently introduced. Thus, the collected data is small and insufficient for our in-depth analyses. The vetting and gathering of data by users through Web forms is also conducted by other projects, such as the nanoHub citation analysis [7]. Our framework supports pluggable data sources that allow for the mining of databases and/or accessing 3rd party service APIs for publication data. In this study, we focus only on two of the data sources - the user submitted publication data via the XSEDE portal, and the extensive NSF award database for automated mining. The former source has user curated data with project affiliation information, and thus in principle it gives a measure of *direct* impact of XSEDE. However, it has limited data entries as of date. On the other hand, the NSF award database contains an extensive compilation of publications that can be automatically mined to pull out all publications for a given XSEDE user. While we cannot directly correlate the publications obtained in this way with XSEDE resources (since the NSF database contains all NSF related publications for a given user regardless of whether the publication was associated with XSEDE use), it does nonetheless provide a measure on a general or *indirect* impact of XSEDE. As a given XSEDE user is affiliated with accounts/projects, and the projects are part of one or more FOS, we can thus tag a publication as being related to the projects and a FOS based on these indirect correlations. Although not ideal, it provides analytic capabilities to analyze an *indirect* impact. Based on this technique, we have been able to obtain over 142,000 publication entries for over 20,000 XSEDE users as of Jan 2014. This by itself is a substantial accomplishment as we know of no other database that has this level of detail that can be correlated to researchers participating in XSEDE. To provide a quick overview of the data analyzed we refer to Figure 2 showing the yearly distribution of the publications (histogram in (a)), and the distribution of number of the publications by project (left boxplot in (b)) and by per user for each project (right boxplot in (b)).

## 4.2 Citation Data Retrieval

While the publication data from the user curated data might be more ideal, we need to conduct an automated search to identify the subsequent citations of the publication recovered from the NSF award database to help provide an indication of the quality of the research.

Due to the size of this publication data (over 142,000 publications), the only realistic way to accomplish this is with an automated process. We have used Google Scholar and ISI Web of Science as sources for the citation data. In order to compare the two methods of obtaining citations, we explored Google Scholar and ISI data for a subset of the publication data, and did a comparison of the results. While a similar comparison has been attempted [20], it was restricted to a very small sample size - 2 people and about 100 publications. In comparison, our study included 33,861 publications and 1,462 users; moreover an author has been identified to have an XSEDE account.

(a) Citation counts comparison for a subset of our publication data

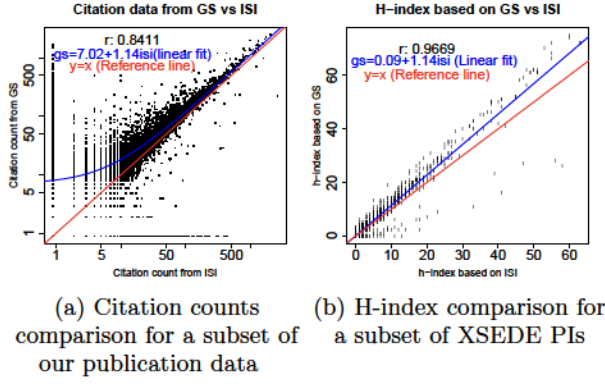(b) H-index comparison for a subset of XSEDE PIs

Figure 3: Comparison of metrics derived from GS vs ISI

The result of this activity is depicted in Figure 3. Part (a) shows the correlation of the citation data from Google Scholar (GS) with the ISI Web of Science (ISI). And part (b) shows the h-index derived from Google Scholar citation data correlating to that calculated from ISI citation data. In either case a high positive correlation is observed. The Pearson correlation coefficients (r) are 0.84 and 0.97 respectively. The very strong correlation of the h-index values are mostly due to the fact that one of the two factors determining the h-index, the number of publications, stay the same for a particular user.

Based on this study, while being aware of the limitations, we were able to use the ISI citation data to get very similar measures for most of the data especially for the h-index metric. This is especially useful if we consider that we do have issues to retrieve a complete citation data set from Google scholar for each of our relevant users and publications based on restricted access rights. Thus, the following analyses are only using citations from ISI.

## 5. RESULTS AND ANALYSES

The previous section described the method used to extract publication and citation data for XSEDE users. With this data now in hand, we discuss the metrics derived from it with the goal of providing a measure of scientific impact. We also conduct analyses to determine if a correlation exists between the data and various categories such as the *field of science*.

## 5.1 Direct impact of XSEDE

By using the user vetted submitted publications only, we were able to show the *direct* scientific impact of XSEDE. As of Jan 27, 2014, there are currently 837 publications registered, involving 882 XSEDE users as authors, 220 organizations, 331 XSEDE projects, and a total of 11,258 citations to date. Please note that these values are based on continuously growing publication data. So far, not all users have contributed their publications to XSEDE projects, or they have not uploaded them to the portal yet. We are working with the XSEDE portal team to provide a publication discovery service in which a user would be presented a list of publications to curate. This is expected to ease the publication acquisition process so more user curated publication data will be available for our future analyses.

Based on the currently available data, we calculated a series of metrics aggregated by user, organization, project,

and FOS. For each entity, we include the number of publications (as header *# of Pubs*), number of citations (as header *Cited by*), h-index and g-index. We also include the $m$ factor of h-index which indicates the slope of the h-index over the years spanned by the publications. This can be used to compare the efficiency between peers if they have the same h-index. Another metric we compute is i10-index [3] which was first introduced in Google Scholar to measure the publication count of publications receiving over 10 citations each. For all the metrics excluding the $m$ factor for h-index, we also compute a *recent* version which was computed using only the publications published from the last 5 years. The time limited comparison helps to compare the peers based on recent work by eliminating effects from older publications.
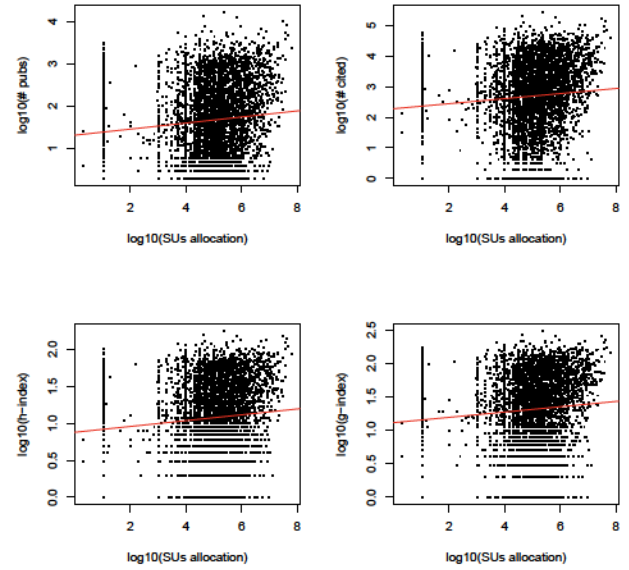
## 5.2 Project metrics vs SUs allocation



Figure 4: Impact Metrics (number of publications, number of citations, h-index, g-index) vs SUs for all projects

|  | Correlation with SUs allocated | r (Pearson's) | df | p-value |
|---|---|---|---|---|
| All Projects | # pubs | 0.242 | 6278 | < 2.2e-16 |
|  | # cites | 0.243 |  | < 2.2e-16 |
|  | h-index | 0.228 |  | < 2.2e-16 |
|  | g-index | 0.220 |  | < 2.2e-16 |
| Research Projects | # pubs | 0.381 | 1677 | < 2.2e-16 |
|  | # cites | 0.377 |  | < 2.2e-16 |
|  | h-index | 0.319 |  | < 2.2e-16 |
|  | g-index | 0.305 |  | < 2.2e-16 |
| Campus Champion Projects | # pubs | 0.335 | 86 | 0.001 |
|  | # cites | 0.315 |  | 0.003 |
|  | h-index | 0.344 |  | 0.001 |
|  | g-index | 0.325 |  | 0.002 |
| Startup Projects | # pubs | 0.025 | 3944 | 0.118 |
|  | # cites | 0.027 |  | 0.091 |
|  | h-index | 0.031 |  | 0.048 |
|  | g-index | 0.035 |  | 0.029 |

Table 1: Correlation between SUs allocated vs the impact metrics for each project

Figure 4 shows the correlation analysis of impact metrics (number of publications, number of citations, h-index, and g-index) versus XSEDE resource allocation (number
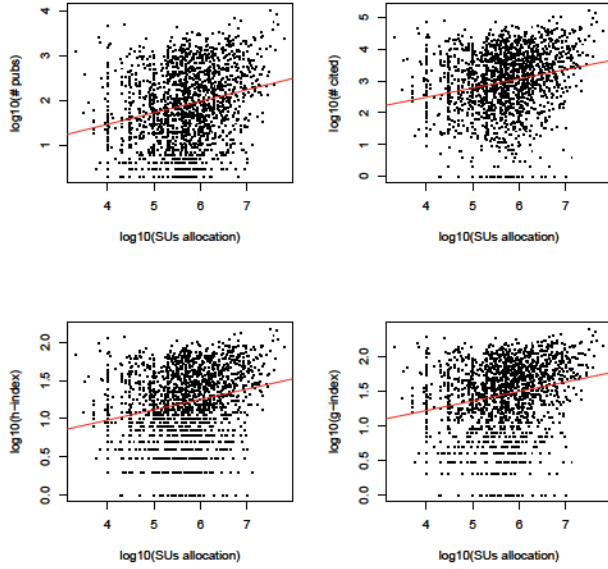
Figure 5: Impact Metrics (number of publications, number of citations, h-index, g-index) vs SUs for research projects



Figure 6: Impact Metrics (number of publications, number of citations, h-index, g-index) vs SUs for FOS's

of SU's) for an individual project (research, start-up, campus champion, etc). Previous work showed a stronger correlation between the citation and SUs [11] using a much smaller sample size taken from a specific XSEDE resource allocation meeting. However, we observed a weaker correlation, if any. When categorizing the projects based on the types (research, startup, campus champion, etc.), it shows a slightly stronger correlation, although still not as strong in correlation to each category other than for the startup projects/allocations. Figure 5 shows the analysis for research projects only. Table 1 lists the correlation coefficient values as well as the p-values showing the significance of the test. Please note in Figure 4 and 5 we included a regression line showing the upper trends of the correlation, i.e., higher SUs allocation correlating to higher impact metrics, but not suggesting a linear relationship. This correlation analysis does not show causality especially since the funding and impact are expected to be related in a feedback loop.

## 5.3   Metrics vs SUs allocation on FOS level

While for individual projects we do not observe strong correlations between impact metrics and the resource allocations, Figure 6 shows stronger positive correlation on the FOS level (132 FOS involved). The Pearson correlation coefficients (r) are 0.704, 0.712, 0.651, 0.648 respectively for the four impact metrics - number of publications, number of citations, h-index and g-index. This result is statistically significant and show a very strong relationship between these variables.

The stronger correlations are most likely caused by the effect of the different sizes of the FOS's. However, this does not diminish the conclusion of the analysis that shows how XSEDE impacts science from different disciplines, e.g., by approving more projects and granting more allocations for certain FOS's.
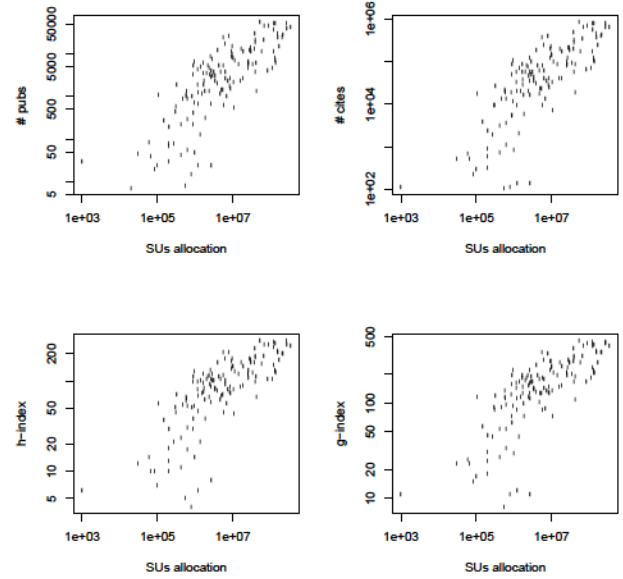
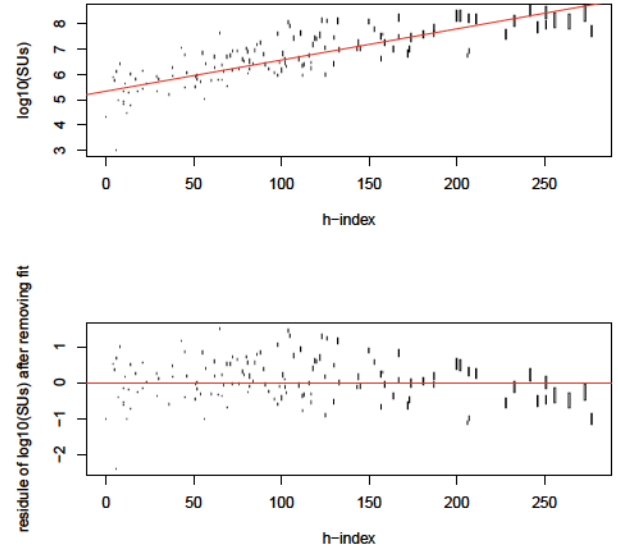Figure 7 shows the SUs allocated (transformed in logarith-



Figure 7: SUs vs h-index for each FOS with trend (above) and residual analysis (bottom)

mic scale) vs the h-index produced for each FOS, while the circle size is proportional to the size (number of projects) of the FOS. It also shows that after removing the fitted trend, we can see a divergence of the SUs received, from the expected SUs trend to produce the given impact judging by h-index. This could imply that certain FOS's are more efficiently (requiring less than expected resources) producing a given impact while some others require more than expected SUs to produce the same impact. An interactive version of this plot is available via the portal interface [9].
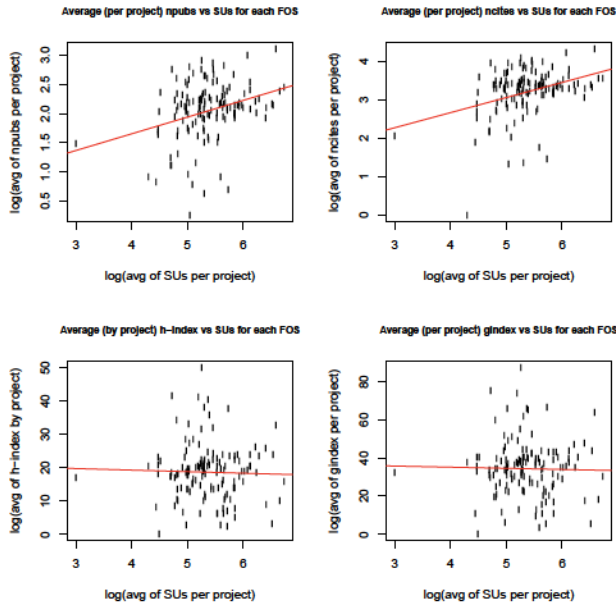
Figure 8: Impact Metrics (number of publications, number of citations, h-index, g-index) vs SUs for FOS (avg by project)



Figure 9: Correlation coefficient (r) of impact metrics vs SUs on project level for each FOS

| | Correlation with average SUs allocated | r (Pearson's) | df | p-value |
|---|---|---|---|---|
| Average per project for each FOS | # pubs | 0.221 | | 0.010 |
| | # cites | 0.222 | 132 | 0.010 |
| | h-index | -0.043 | | 0.620 |
| | g-index | -0.035 | | 0.688 |

Table 2: Correlation between average SUs allocated vs the average impact metrics (by projects) for each FOS

As we see, the size of FOS significantly effects the impact as well as the allocations (for h-index as in Figure ??). We can eliminate this effect by comparing the average values within each FOS by dividing the number of projects, as shown in Figure 8, while Table 2 has the values. It shows the weak correlation of per project based metrics vs SUs for the number of publications and citations, which is actually not significantly different than the result presented in Table 1. We did not observe any correlation between allocation and h-index or g-index. This is probably caused by the fact that these two metrics do not work well when being averaged as they are not cumulative or additive values.

However, as shown in Figure 9, within each FOS, the project level metrics vs SUs correlations are typically higher especially for large size FOS's. With increasing size of the FOS (n=10, n=50, and n=100 are denoted as vertical lines), the correlation appears positively higher and more significant. Figure 10 shows the distribution of correlation coefficients (r) between number of publications and allocations for each project within the same FOS, while grouped by size of FOS (number of projects). Note the general trend that the extremes and ranges are narrowing, and the medians are increasing (above 0.4 for groups of FOS with more than 50 projects), along with the increase of the FOS size. This suggests that for the majority of the FOS, impact metrics for a project do have a positive correlation with SUs allocated to the project. By investigating the individual data points, we
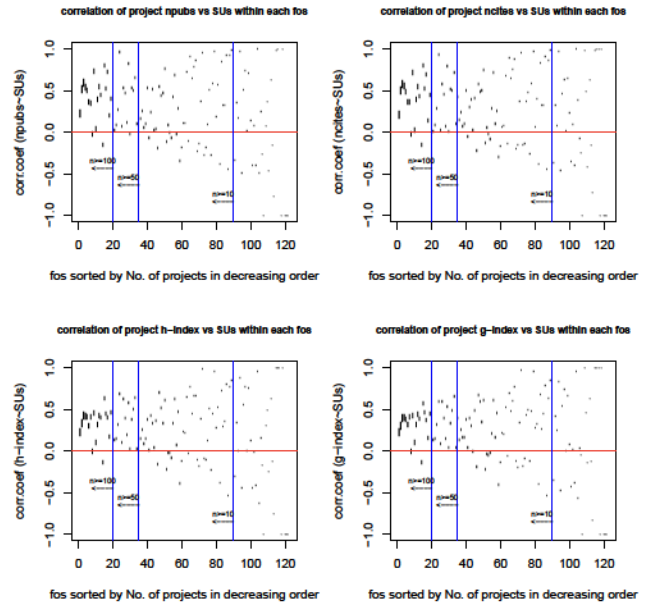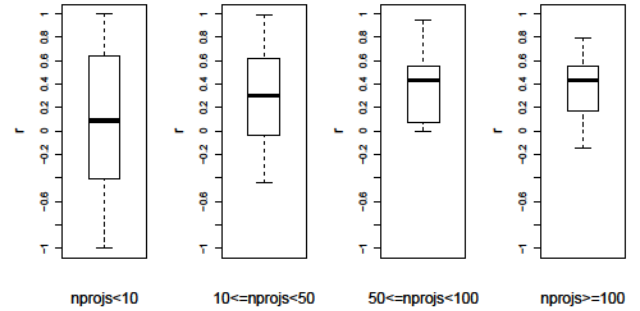


Figure 10: Distribution of r grouped by size of FOS

would be able to find in which FOS this correlation appears much stronger, and in which they are weak. This could be potentially used during resource allocation to help determine which projects should be preferred when resources are limited but demands are high.

## 5.4 Scientific Impact Produced per SU Allocation Unit

The publications database acquired from the NSF awards database includes all publications from XSEDE users rather than just those relevant to XSEDE. As such, these publications present only an indirect measure of the scientific impact of XSEDE, diluted by the presence of many publications that are not related to the XSEDE resources. A more ideal, and direct measurement of XSEDE's scientific impact is obtained from the user curated publication database. We have shown in the previous section that these scientific impact metrics can be used to measure the scientific impact of XSEDE in general, as well as comparing individual users, projects, and FOS with their peers. As these metrics are
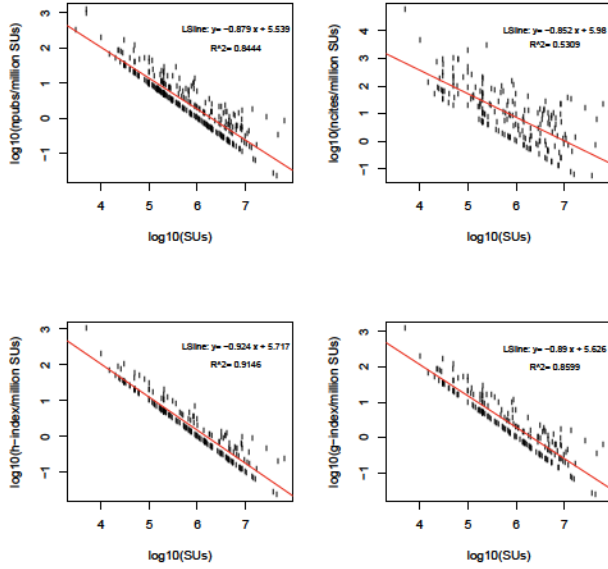
Figure 11: Four different measures of scientific impact per SU allocated. Note that the Y-axis gives scientific impact scaled by SU allocation. Therefore, these plots indicate that as the allocation size grows there is a diminishing scientific impact per SU allocated

obtained from the publications that are tagged as results from an XSEDE project, we also can do an analysis of the scientific impact produced per SU allocation unit.

We have calculated the scientific impact for those involved projects (302 out of more than 6,000 in total) based on the direct metrics obtained earlier and SUs allocated to them (in million SUs). Figure 11 shows a series of four log-log plots in which four different scientific impact metrics for each project are scaled by the SU allocation then plotted against the total allocation. Previously we have demonstrated the positive correlation of scientific impact metrics and the resource allocation of projects within each FOS. Figure 11 suggests that based on these metrics of scientific impact, that is number of papers, citations, h-index, and g-index scaled by SU's, sponsoring a larger number of smaller scale projects could actually produce a higher scientific impact than a smaller number of very large projects. In other words, we cannot expect a project that received double the amount of SUs of what another project did to produce double the impact, as measured by number of publications, citation counts, h-index, and g-index.

Unfortunately, to date, the number of user curated publications is still too small. With more such data available over time, we anticipate to repeat our analysis in order to derive the scientific impact of XSEDE and demonstrate the relationship between XSEDE funded allocations and a variety of scientific impact metrics.

## 6. ONGOING AND FUTURE WORK

This paper does not address the researchers name ambiguity issue, which deserves dedicated research. The root cause of this issue is that the metadata of the publications simply does not include enough information to distinguish similar names that can be uniquely associated to XSEDE

user names. This is not a problem specific to our study but for the automated bibliometrics analysis in general. In the future, we plan to tackle the problem based on other available data such as field of science, organization, funding data, co-author relationship etc. while conducting machine learning techniques as well as adopting social network based analyses.

A useful compromise is to let users curate their publication list. We will include processes assisting in the curration of date into the workflow of vetting the papers. One pathway we currently pursue is to work with the XSEDE portal team while providing the publication data we have collected as a publication discovery service, in the hope to provide more convenient way for users to quickly populate the vetted publications library.

We have also started another similar activity, in which we are attempting to extract and parse the publication data from past TeraGrid/XSEDE quarterly reports. This data, while not curated on per user basis, does have project level association information, and thus, can serve quite well for most of our analyses.

As for the resource allocation, we currently only considered the Service Units (SUs), or cpu-hours, as this is the dominant factor thus far to measure resource allocation in XSEDE. With the increasingly importance and bigger needs of storage allocations from big-data applications, and Virtual Machine (VM) based allocations for those interested into cloud computing, we will need to put these also into the equation to cover more forms of resources in addition to SUs.

Finally, we are conducting social networking related analyses among publications, users, projects, FOS's, etc. based on citation and co-authorship relations. Mining social networking media such as Twitter and Facebook is also planned to obtain usage data, among other altmetrics, to compliment the publication-based scientific impact studies.

## 7. CONCLUSION

This paper presents a framework to facilitate the measuring of scientific impact and evaluation of ROI for large computing facilities. We have used this framework to conduct an evaluation of scientific impact of XSEDE by deriving various metrics and carrying out extensive statistical analyses. The major accomplishments include:

1. We have devised a process to obtain and manage publication and citation data from various sources for a given group of people. We have followed this workflow to obtain over 142,000 publications as well as the citation count data for over 20,000 XSEDE users.

2. Based on the consolidated relevant bibliometrics data various scientific impact metrics are derived for users and other aggregated levels such as projects and field of science.

3. The results are presented via a lightweight portal, and are also exposed via database integration or RESTful services to other portals, including the XDMoD portal and the XSEDE portal. For example, we expose the publication data via RESTful service API to the XSEDE portal team as a publication discovery service. This will help facilitate the identification and curation of XSEDE enabled publications by XSEDE users.

4. Statistical analyses were carried out correlating the impact metrics with projects/proposals, field of science, and allocation data to help provide metrics that can be used to quantify the impact of XSEDE resources on scientific research. These analyses do show a positive correlation between XSEDE funded allocations and various scientific impact metrics. When more XSEDE directly related data are available we expect to provide much more insight into the scientific impact of the XSEDE program.

5. We have conducted preliminary analyses on scientific impact produced per SU allocation unit based on user curated publication data with a limited sample size. This may provide a way to measure the ROI of XSEDE. We will conduct similar analyses when having more user curated publications to further solidify the results.

It is obvious that continuous work is crucial to conduct longitudinal tracking of the data and deal with the issues that XSEDE has so far provided limited amount of publication data. Important is to note that this work has pioneered the workflow and the analysis capability on how to achieve the data gathering. This framework can be reused by various groups enabling different services as part of XSEDE to assist the portal and auditing teams. Moreover, this framework and its service oriented model makes it possible to expand its usage beyond those targeted by XSEDE resources. It could be employed within other organizations such as Department of Energy (DOE) or even a department of a university. Those that would like to consult with us on such specializations, can contact us for further details.

## 8. ACKNOWLEDGMENTS

## 9. REFERENCES

[1] altmetrics. Web Page. URL: http://altmetrics.org/manifesto/.

[2] Google Scholar. Web Page. URL: http://scholar.google.com/.

[3] i-10 index | google scholar citations open to all. URL: http://googlescholar.blogspot.com/2011/11/google-scholar-citations-open-to-all.html.

[4] ISI Web of Science. Web Page. URL: http://wokinfo.com/.

[5] Mendeley. Web Page. URL: http://www.mendeley.com/.

[6] Microsoft Academic Search. Web Page. URL: http://academic.research.microsoft.com/.

[7] nanoHUB.org - Citations. Web Page. URL: https://nanohub.org/citations.

[8] Publish or Perish. Web Page. URL: http://www.harzing.com/pop.htm.

[9] Tas scientific impact metrics and analysis dev/testing portal. URL: http://fgdev.pti.indiana.edu:8088/xdportalpub/.

[10] XSEDE. Web Page. URL: https://www.xsede.org/.

[11] J. Bollen, G. Fox, and P. R. Singhal. How and where the TeraGrid supercomputing infrastructure benefits science. *Journal of Informetrics*, 5(1):114–121, 2011.

[12] J. Bollen, M. A. Rodriguez, and H. Van de Sompel. MESUR: Usage-based Metrics of Scholarly Impact. In *Proceedings of the 7th ACM/IEEE-CS Joint Conference on Digital Libraries*, JCDL '07, pages 474–474, New York, NY, USA, 2007. ACM. URL: http://doi.acm.org/10.1145/1255175.1255273, doi:10.1145/1255175.1255273.

[13] J. Bollen, H. Van de Sompel, A. Hagberg, and R. Chute. A principal component analysis of 39 scientific impact measures. *PloS one*, 4(6):e6022, 2009.

[14] J. Bollen, H. Van de Sompel, and M. A. Rodriguez. Towards Usage-based Impact Metrics: First Results from the Mesur Project. In *Proceedings of the 8th ACM/IEEE-CS Joint Conference on Digital Libraries*, JCDL '08, pages 231–240, New York, NY, USA, 2008. ACM. URL: http://doi.acm.org/10.1145/1378889.1378928, doi:10.1145/1378889.1378928.

[15] L. Egghe. Theory and practise of the g-index. *Scientometrics*, 69(1):131–152, 2006.

[16] T. R. Furlani, B. L. Schneider, M. D. Jones, J. Towns, D. L. Hart, S. M. Gallo, R. L. DeLeon, C.-D. Lu, A. Ghadersohi, R. J. Gentner, A. K. Patra, G. von Laszewski, F. Wang, J. T. Palmer, and N. Simakov. Using xdmod to facilitate xsede operations, planning and analysis. In *Proceedings of the Conference on Extreme Science and Engineering Discovery Environment: Gateway to Discovery*, XSEDE '13, pages 46:1–46:8, New York, NY, USA, 2013. ACM. URL: http://doi.acm.org/10.1145/2484762.2484763, doi:10.1145/2484762.2484763.

[17] J. E. Hirsch. An index to quantify an individual's scientific research output. *Proceedings of the National academy of Sciences of the United States of America*, 102(46):16569–16572, 2005.

[18] J. Kaur, D. T. Hoang, X. Sun, L. Possamai, M. JafariAsbagh, S. Patil, and F. Menczer. Scholarometer: A social framework for analyzing impact across disciplines. *PloS one*, 7(9):e43235, 2012.

[19] P. Thomas and D. Watkins. Institutional research rankings via bibliometric analysis and direct peer review: A comparative case study with policy implications. *Scientometrics*, 41(3):335–355, 1998.

[20] K. Yang and L. I. Meho. Citation analysis: a comparison of Google Scholar, Scopus, and Web of Science. *Proceedings of the American Society for Information Science and Technology*, 43(1):1–15, 2006.