

Progress Toward Cancer Data Ecosystems

Robert L. Grossman, PhD

Abstract: One of the recommendations of the Cancer Moonshot Blue Ribbon Panel report from 2016 was the creation of a national cancer data ecosystem. We review some of the approaches for building cancer data ecosystems and some of the progress that has been made. A data commons is the collocation of data with cloud computing infrastructure and commonly used software services, tools, and applications for managing, integrating, analyzing, and sharing data to create an interoperable resource for the research community. We discuss data commons and their potential role in cancer data ecosystems and, in particular, how multiple data commons can interoperate to form part of the foundation for a cancer data ecosystem.

Key Words: Cancer data ecosystems, clinical informatics, cloud computing, data commons, data ecosystems, data sharing, genomic databases

(*Cancer J* 2018;24: 122–126)

Accelerating Research Through Data Sharing

With the new instruments producing ever larger amounts of data relevant to biology, medicine, and health care¹ and with the maturing of data science,² there are increasing opportunities to use data science to make discoveries. Sharing cancer data can provide the data required to make basic discoveries about the mechanisms of cancer,³ accelerate the discovery of biomarkers,⁴ inform the repurposing of drugs,⁵ and aggregate data to provide the evidence required to discover genomic alterations associated with rare subtypes of cancer or other rare diseases.⁶ An overview of the importance of sharing molecular data in cancer can be found in Siu et al.⁷ A succinct target for data sharing in cancer was expressed by former Vice President Biden with the phrase “The goal [...] is to accelerate progress in cancer research and treatment enough to do in 5 years what would have taken 10 [years].”⁸

A fundamental tension in biomedical research is the trade-off between (i) ethical, legal, and regulatory requirements to obtain consent from patients and to protect patients and their data and (ii) the right of patients to benefit from research⁹ (Fig. 1). There are always potential risks when different data sets and large data sets are brought together.^{10,11} For example, through what are sometimes called side-channel attacks,¹² it is becoming increasingly easy to uncover the likely identities of some individuals when large deidentified data sets are made available. An example by Gymrek et al.¹³ shows how genotype data, general geolocation information (say, associated with the hospital that provided the biospecimen), and a local white pages can be used to create a short list of individuals who may have potentially contributed data to the 1000 Genomes Project.

One of the 10 recommendations of the Cancer Moonshot Blue Ribbon Panel¹⁴ was to build a national cancer data ecosystem “...that will enable all participants across the cancer research and care continuum to contribute, access, combine, and analyze diverse data that will enable new discoveries and lead to lowering

the burden of cancer in our country.”¹⁴ The Blue Ribbon Panel report lists “allow both public and private information resources to be readily discovered and connected through the use of a common information architecture”¹⁴ as one of the goals of a cancer data ecosystem.

In this review, we address 2 questions related to information architectures for supporting a cancer data ecosystem:

- (1) What are some of the systems that can be used for sharing cancer genomic and clinical data with security and compliance in order to gain access to a critical mass of data that can shorten the time required to make discoveries that benefit cancer patients?
- (2) How can we build an ecosystem of interoperating systems for sharing data and applications that access them so that research scientists can most effectively use the data that are being shared?

Data Clouds and Data Commons

Data commons usually rely on cloud computing infrastructure, which is based on the virtualization¹⁵ of computing, storage, networking, and other services required for computing. Virtualization provides the required computing infrastructure using software simulation rather than physical hardware. This provides 2 main advantages: first, different simulated hardware configurations can use the same underlying physical hardware; second, the underlying physical hardware can be used more efficiently. With this approach, storage and computing resources can be provisioned elastically as needed by users, even at large scale.¹⁶ Cloud computing infrastructure is available from open source software, such as OpenStack,¹⁷ which can be run by an organization locally or on products from third-party cloud service providers, such as Amazon's AWS, Google's GCP, or Microsoft's Azure.

For the purposes here, we can define a *data commons* as a system that collocates data with cloud computing infrastructure and commonly used software services, tools, and applications for managing, analyzing, and sharing data to create an interoperable resource for the research community.¹⁸

Some of the most important services for working with cancer data include (1) authentication and authorization services, (2) digital ID services, (3) metadata services, and (4) services for launching virtual machines, containers, notebooks, or similar infrastructure for computing over data.

Additional services that are sometimes supported include (5) using a common data dictionary and integrating data with respect to a common data model, (6) processing the integrated data with respect to a common set of bioinformatics pipelines, and (7) exposing the integrated and harmonized data using an application programming interface (API). Notice that services (1) to (3) when exposed via an API, as is standard, is enough to make data *findable* and *accessible* as usually defined,¹⁹ whereas (5) is enough to make data *interoperable*.¹⁹

The term *data commons* is beginning to be used to refer to a system with capabilities (1) to (7), and the term *data cloud* to refer to a system when properties (1) to (4) are supported.

One of the important uses of data clouds is to support the large-scale processing of NGS data. Researchers can define their

From the University of Chicago, Chicago, IL.

Reprints: Robert L. Grossman, PhD, University of Chicago, 900 E 57th Street, KCBD 10142, Chicago IL 60637. E-mail: robert.grossman@uchicago.edu.

The authors have disclosed that they have no significant relationships with, or financial interest in, any commercial companies pertaining to this article.

Copyright © 2018 Wolters Kluwer Health, Inc. All rights reserved.

ISSN: 1528-9117

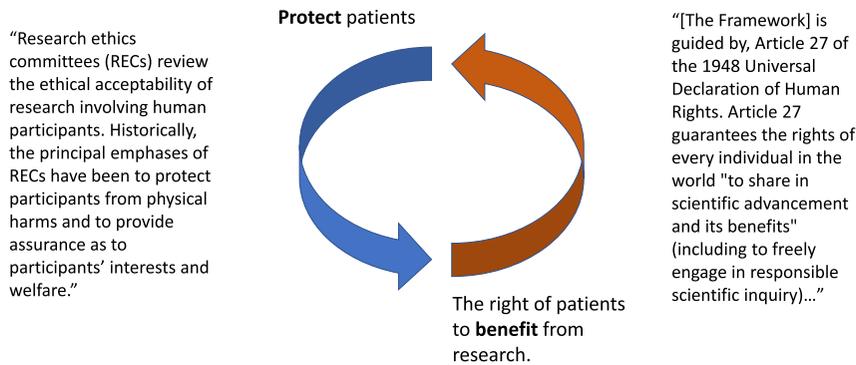


FIGURE 1. Both quotes in the figure are from Knoppers.⁹ One of the challenges with building data ecosystems for cancer is balancing the ethical requirements to protect patients and their data from the right that all individuals have to benefit from scientific advances.⁹ Today, data ecosystems integrate cancer data that not only can accelerate research and improve outcomes for cancer patients, but also can lead to potential exposures of data. Balancing these 2 requirements is one of the challenges of building a data ecosystem for cancer.

own bioinformatics workflows, work with members of an analysis working group (AWG) to develop workflows for the working group, or use workflows developed and optimized by third-party vendors.

Data clouds that support user-defined workflows to support individual researchers or AWGs include the University of Chicago's Bionimbus Protected Data Cloud,²⁰ OICR's Cancer Collaboratory,²¹ the Galaxy Cloud,^{22,23} Seven Bridges Cancer Genomics Cloud,²⁴ Broad's FireCloud,²⁵ and the Institute for Systems Biology Cancer Genomics Cloud.²⁶ The latter 3 were part of the National Cancer Institute (NCI) Cloud Pilots (now Cloud Resources).²⁷ A good analysis of using multiple clouds to support an AWG is reported by Yung et al.²⁸

Data clouds that support running vendor-supplied bioinformatics pipelines (sometimes called *NGS as a service*) include the DNAnexus platform^{29,30} and Globus Genomics.³¹

The NCI Genomic Data Commons (GDC) and the BloodPAC Data Commons³² are examples of a data commons satisfying requirements (1) to (7).^{33,34} Data commons usually have portals for exploring data and downloading data, as well as submitting data, and APIs so that third-party software services, applications, notebooks,³⁵ and cloud-based collaborative workspaces can access data managed by the commons.

Pediatric Data Commons

Because of the much small number of patient data available within the pediatric cancer research community, data sharing is essential for many discoveries.³⁶ As an example, the International Neuroblastoma Risk Group collects clinical data and genotype data about neuroblastoma tumors from patients from around the world to create a database and associated data commons.³⁷ Given the relatively small number of neuroblastoma tumors each year, without such an effort there would not be enough data to reliably classify neuroblastoma tumors.³⁸

Other examples of data commons supporting the pediatric cancer community include Cavatica (www.cavatica.org) and the Kids First Data Resource (d3b.center/kidsfirst/).

Standards

Standards for data commons and data ecosystems are immature and still emerging. Efforts worth mentioning include the Global Alliance for Genomics and Health,³⁹ which is establishing standards for sharing clinical and genomic data within frameworks for data governance, security, privacy, and patient benefits. The FORCE11 FAIR principles¹⁹ describe guidelines so that data can be findable, accessible, interoperable, and reusable (FAIR).

A challenge with efforts to establish standards for data commons is that the standards can get so far ahead of any scalable implementations that they become irrelevant. In response to this, 3 large-scale projects for sharing clinical and genomic data began to collaborate on APIs, while competing on building platforms. The 3 groups are some of the developers behind the (i) NCI GDC and NCI Cancer Research Data Commons (NCRDC), (ii) the NIH All of Us Research Program (<https://allofus.nih.gov/>), and (iii) the Human Cell Atlas.⁴⁰ The collaboration is called the Commons Alliance or Data Biosphere. It is important to note that it is a commitment of a group of the core software developers supporting these projects that are working together to collaborate on standardized APIs and not, at this time, a commitment of the organizations themselves. Note that a similar approach was taken by the NCI Cloud Pilots (“collaborate on standardized APIs and compete on implementations”).²⁷

Cancer Data Ecosystem—Initial Steps

Currently, one way that cancer data ecosystems are emerging is the ability for several data commons to interoperate and to share a common ecosystem of applications by each using a core set of common services. For example, the NCI GDC,³³ the BloodPAC Data Commons for liquid biopsies,³² and the Data Catalog for the Kids First Data Resource, plus several additional data commons, can all interoperate because they are built with the same open source software (Gen3) and are built around the following core services that are all accessible through open APIs:

- (1) authentication and authorization services,
- (2) digital ID services, and
- (3) metadata services.

In addition, they all use the NCI Thesaurus⁴¹ and are based in part on the GDC Data Model and use the GDC API.³⁴ With the digital ID services, data can be in any private or public cloud, located from its metadata (3) using any convenient searching and indexing services, and accessed by a client through an API, as long as the user is authenticated and authorized to access the data using (1).

More information about the Gen3 software stack for building data commons and data ecosystems can be found in Gen3 Project Team.⁴²

The NCI Cancer Data Ecosystem

Perhaps the most mature cancer data ecosystem has been developed by the NCI and currently consists of the NCI GDC and the NCI Cloud Resources.²⁷ It is called the NCRDC. The Cloud

Resources²⁴ in the NCRDC use AWS and the GCP so that users can compute over data from the GDC and other resources and contain different applications for analyzing data. Importantly, the Cloud Resources also support workspaces. Additional commons are being added to the NCRDC including a proteomics and imaging commons, with additional commons and resources in planning. The NCRDC Framework Services can be used so that the different data commons and other NCRDC resources can share authentication, authorization, ID, and metadata services.

More specifically, the NCRDC data ecosystem is based on the following framework services:

- (1) authentication (AuthN) and authorization (AuthZ) services,
- (2) metadata validation using the NCI Thesaurus⁴¹ and related services,
- (3) an extensible data model based on the GDC data model,³⁴
- (4) APIs for containers and workflows based on the Global Alliance for Genomics and Health WES and TES standards, and
- (5) services to set up and access collaborative workspaces.

In particular, this enables FAIR resources to be built and for the data commons and resources to interoperate.

Currently, public clouds typically charge egress charges for moving data outside the cloud boundary. Network peering is the idea that 2 Internet service providers exchange network traffic at no charge, which is in contrast to the charges imposed when an Internet service provider charges its customers for transmitting network traffic.⁴³ Data peering between the large Internet service providers was an important driver for the growth of the Internet. An important goal for building a data ecosystem is to support data peering between data commons so that research data can be accessed from one data commons at no cost, even when it is stored in another data commons or cloud.¹⁸ Because public clouds make money when their users compute over data, supporting some level of data peering should be quite feasible.

There is no agreed-to definition of a data ecosystem at this time, but at the minimum, a data ecosystem (as opposed to a data commons) should support:

- (1) authentication and authorization services so that a community of researchers can access an ecosystem of data and applications with a common (research) identity;

- (2) the ability for multiple data commons to interoperate, preferably through data peering;
- (3) shared data models to simplify the ability for third-party applications to access data from multiple data commons; and
- (4) a collection of applications that are powered by APIs that are FAIR compliant.

A properly designed and operated data commons is not limited by the type of data, the amount of data, or the type of applications that can be supported, but usually by the priorities, interests, and funds of the sponsors and operators of the data commons. For this reason, one of the most important roles of a data ecosystem is to support multiple data commons that can be interoperable and a rich variety of third-party applications over them.

Architectures for Data Ecosystems to Support Cancer Research

There is no consensus system architecture at this time for the key components in a data ecosystem to support cancer research. Below, we discuss 4 options that are being used.

Commons With an Open API

With this approach, a data commons exposes an open API, which can be used by large-scale cloud resources. This was the approach that emerged during 2014–2016 with the NCI GDC³³ and NCI cancer clouds developed by Broad,²⁵ Seven Bridges Genomics,²⁴ and the Institute for Systems Biology.²⁶ This architecture is enabled by the GDC API,³⁴ which is used both by the GDC itself for all its applications, by the 3 cancer clouds, and by an increasing number of third-party applications.

Walled Garden

Another approach is to develop a data ecosystem within a proprietary system. Examples include DNAnexus²⁹ and Seven Bridges Genomics.²⁴ These types of systems may expose APIs for users and may be part of a broader ecosystem but are not themselves open source.

Narrow Middle Design

An important design philosophy when designing complex systems is sometimes called the *end-to-end principle*, “[which

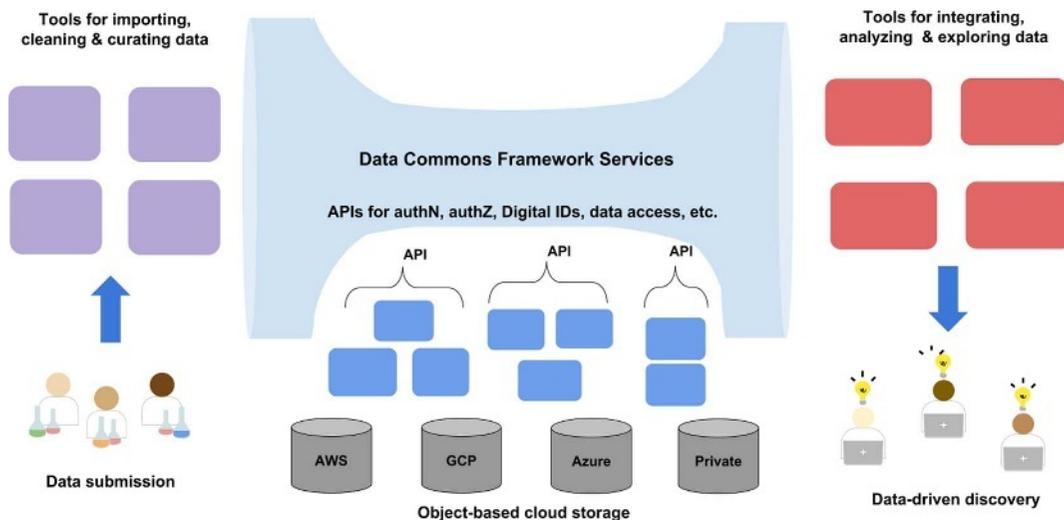


FIGURE 2. This figure illustrates the “narrow middle” architecture for data ecosystems that uses a core set of framework services that can support a rich collection of rapidly evolving applications and services for ingesting and curating data (at one end) and a rich collection of rapidly evolving applications and services for exploring and analyzing data (at the other end).

TABLE 1. Some of the Different Federation Architectures That Have Been Used for Building Data Ecosystems for Cancer

Design	Example
Single system with open APIs supporting third-party applications	NCI GDC with third-party applications, notebooks, and workspaces ³⁴
Federation of multiple independent interoperating systems with APIs supporting third-party applications	The Matchmaker Exchange Platform for Rare Disease Gene Discovery ⁶
Federation of multiple independent data commons with common services and common core data models, with open APIs supporting third-party applications	BloodPAC Data Commons ³² and other Gen3 data commons ⁴² and applications that can access data from them
Federation of multiple independent interoperating systems over a set of common services and common core data models, with open APIs supporting third-party applications	NCRDC, including GDC, Cloud Resources, Proteomics Databases ²⁷
Federation of multiple independent systems connected via APIs	Kids First Data Resource (d3b.center/kidsfirst)

argues] for moving function upward in a layered system, closer to the application that uses the function.”⁴⁴ The underlying architecture for the Internet was based in part on this principle, with the end result that the protocols such as tcp and http in the middle were relatively simple and changed slowly, in contrast to the Internet applications at the “ends,” which were much more complicated and evolved rapidly.⁴⁴ In our context of data ecosystems for cancer, we might term this the *narrow middle design*, in which the minimal core services necessary to support a data ecosystem are kept few, small, and simple, with innovation and rich applications at the ends: for getting data into the ecosystem (at the “input end”) and for exploring, integrating, harmonizing, and analyzing data in the ecosystem (at the “output” end). Notice that the data commons framework services in the National Cancer Research Data Commons are an example of this approach (Fig. 2).

Modular Components and Services

Another approach being explored for building data commons and ecosystem is to support enough core modular services and components that are FAIR compliant¹⁹ that a complete commons or data ecosystem can be built. For example, the NIH Data Commons Pilot Phase Consortium is taking this approach and developing a system by integrating a suite of components and services, including persistent identifier services, metadata services, data indexing services, data search services, services to support workspace, and services for working with complex phenotype and genotype data (commonfund.nih.gov/commons).

Once multiple systems have been developed, it is common to federate them⁴⁵ in order to create a large ecosystem. There are multiple ways that this can be done,⁴⁵ and Table 1 describes some of them.

Cancer Data Ecosystems—Next Steps

Moving from the current state of cancer data ecosystems requires making 4 transitions:

First, the various current systems for sharing cancer data, including the NCI GDC, NCI Cloud Resources, NCBI’s ClinVar,⁴⁶ AACR’s Project GENIE,⁴⁷ ASCO’s CancerLinQ,⁴⁸ and so on, must themselves interoperate (Fig. 3). These systems are organized into different “lanes,” with different regulatory environments and different incentives for sharing and operated by different entities with different sustainability models, which makes interoperating more challenging.

Second, patient portals and related infrastructure supporting patient-partnered research must be integrated into all aspects of the ecosystem. This technology is still maturing, and which systems will emerge as the dominant ones supported by large groups of patients is still not at all clear.

Third, changes must be made to make it easier to get data into the ecosystem and to curate the data that are part of it. This effort requires both better tools and a larger better-trained workforce of bioinformaticists and data scientists.

Fourth, sustainability is both a major challenge and a major reason that current systems do not interoperate. Today’s systems

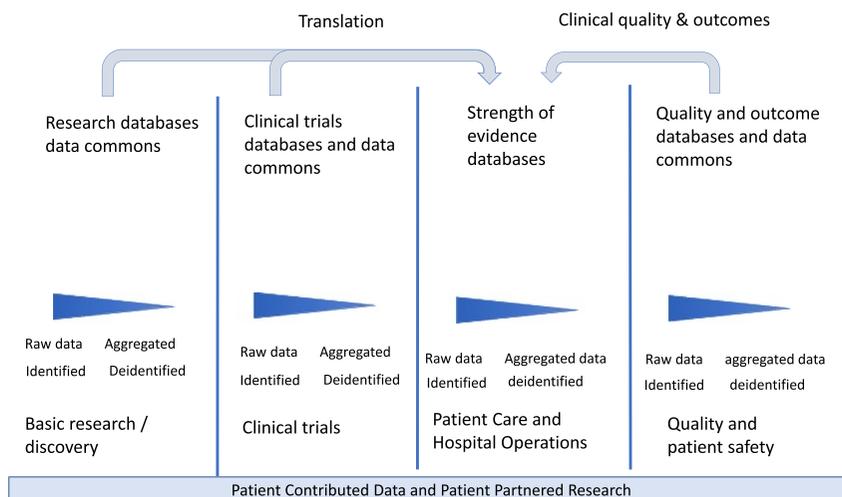


FIGURE 3. The cancer data ecosystem includes data from (1) research, (2) clinical trials, (3) patient care and hospital operations, and (4) data quality and outcomes research. An important new source of data will be patient-contributed data.

often limit access to dues-paying members or sign commercial agreements that limit how data can be shared. Funders, including federal agencies and foundations, must build data sharing and support for data sharing ecosystems into how they support research, or we risk not having the critical mass of data available that are needed to improve patient outcomes. Payers must also insist on data sharing so that the required evidence is available to inform standard of care instead of relying on black box systems that do not share the underlying evidence that is used for decision support.

REFERENCES

- Stephens ZD, Lee SY, Faghri F, et al. Big data: astronomical or genetical? *PLoS Biol.* 2015;13:e1002195.
- Donoho D, ed. *50 Years of Data Science*. Based on a presentation at the Tukey Centennial Workshop; 2015; Princeton, NJ.
- Weinstein JN, Collisson EA, Mills GB, et al. The Cancer Genome Atlas Pan-Cancer Analysis Project. *Nat Genet.* 2013;45:1113.
- Khleif SN, Doroshow JH, Hait WN. AACR-FDA-NCI Cancer Biomarkers Collaborative consensus report: advancing the use of biomarkers in cancer drug development. *Clin Cancer Res.* 2010;16:3299–3318.
- Li J, Zheng S, Chen B, et al. A survey of current trends in computational drug repositioning. *Brief Bioinform.* 2015;17:2–12.
- Philippakis AA, Azzariti DR, Beltran S, et al. The Matchmaker Exchange: a platform for rare disease gene discovery. *Hum Mutat.* 2015;36:915–921.
- Siu LL, Lawler M, Haussler D, et al. Facilitating a culture of responsible and effective sharing of cancer genome data. *Nat Med.* 2016;22:464.
- Ferber D. Biden blasts science denialists, calls to dramatically speed up fight against cancer. *Science Magazine* February 19, 2018.
- Knoppers BM. Framework for responsible sharing of genomic and health-related data. *HUGO J.* 2014;8:3.
- Karp DR, Carlin S, Cook-Deegan R, et al. Ethical and practical issues associated with aggregating databases. *PLoS Med.* 2008;5:e190.
- Caulfield T, McGuire AL, Cho M, et al. Research ethics recommendations for whole-genome research: consensus statement. *PLoS Biol.* 2008;6:e73.
- Erllich Y, Narayanan A. Routes for breaching and protecting genetic privacy. *Nat Rev Genet.* 2014;15:409–421.
- Gymrek M, McGuire AL, Golan D, et al. Identifying personal genomes by surname inference. *Science.* 2013;339:321–324.
- Panel BR. Cancer Moonshot Blue Ribbon Panel Report. 2016 [cited 2018]. Available at: <https://www.cancer.gov/research/key-initiatives/moonshot-cancer-initiative>. Accessed April 19, 2018.
- Scarfone K, Souppaya M, Hoffman P. Guide to security for full virtualization technologies. National Institute of Standards and Technology Special Publication. 2011;800–125.
- Mell P, Grance T. The NIST definition of cloud computing (draft): recommendations of the National Institute of Standards and Technology. National Institute of Standards and Technology; 2011.
- Pepple K. Deploying Openstack. Champaign: O'Reilly Media, Inc., 2011.
- Grossman RL, Heath A, Murphy M, et al. A case for data commons: toward data science as a service. *Comput Sci Eng.* 2016;18:10–20.
- Wilkinson MD, Dumontier M, Aalbersberg IJ, et al. The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data.* 2016;3:160018.
- Heath AP, Greenway M, Powell R, et al. Bionimbus: a cloud for managing, analyzing and sharing large genomics datasets. *J Am Med Inform Assoc.* 2014;21:969–975.
- Yung CK, Mihaiescu GL, Tiernay B, et al. *The Cancer Genome Collaboratory*. Philadelphia: AACR; 2017.
- Afgan E, Baker D, Coraor N, et al. eds. Galaxy CloudMan: delivering cloud compute clusters. *BMC Bioinformatics.* 2010;11(suppl 12):S4.
- Afgan E, Baker D, Coraor N, et al. Harnessing cloud computing with Galaxy Cloud. *Nat Biotechnol.* 2011;29:972.
- Lau JW, Lehnert E, Sethi A, et al. The Cancer Genomics Cloud: collaborative, reproducible, and democratized—a new paradigm in large-scale computational research. *Cancer Res.* 2017;77:e3–e6.
- Birger C, Hanna M, Salinas E, et al. FireCloud, a scalable cloud-based platform for collaborative genome analysis: strategies for reducing and controlling costs. *bioRxiv.* 2017: 209494.
- Reynolds SM, Miller M, Lee P, et al. The ISB Cancer Genomics Cloud: a flexible cloud-based platform for cancer genomics research. *Cancer Res.* 2017;77:e7–e10.
- Hinkson IV, Davidsen TM, Klemm JD, et al. A comprehensive infrastructure for big data in cancer research: accelerating cancer research and precision medicine. *Front Cell Dev Biol.* 2017;5:83.
- Yung CK, O'Connor BD, Yakneen S, et al. Large-scale uniform analysis of cancer whole genomes in multiple computing environments. *bioRxiv.* 2017: 161638.
- Reid JG, Carroll A, Veeraraghavan N, et al. Launching genomics into the cloud: deployment of Mercury, a next generation sequence analysis pipeline. *BMC Bioinformatics.* 2014;15:30.
- Shringarpure SS, Carroll A, Francisco M, et al. Inexpensive and highly reproducible cloud-based variant calling of 2,535 human genomes. *PLoS One.* 2015;10:e0129277.
- Madduri RK, Sulakhe D, Lacinski L, et al. Experiences building Globus Genomics: a next-generation sequencing analysis service using Galaxy, Globus, and Amazon Web services. *Concurr Comput.* 2014;26:2266–2279.
- Grossman RL, Abel B, Angiuoli S, et al. Collaborating to compete: Blood Profiling Atlas in Cancer (BloodPAC) consortium. *Clin Pharmacol Ther.* 2017;101:589–592.
- Grossman RL, Heath AP, Ferretti V, et al. Toward a shared vision for cancer genomic data. *N Engl J Med.* 2016;375:1109–1112.
- Wilson S, Fitzsimons M, Ferguson M, et al. Developing cancer informatics applications and tools using the NCI Genomic Data Commons API. *Cancer Res.* 2017;77:e15–e18.
- Kluyver T, Ragan-Kelley B, Pérez F, et al., eds. Jupyter Notebooks—a publishing format for reproducible computational workflows. *Proceedings of the 20th International Conference on Electronic Publishing*, Göttingen, 2016; 87–90.
- Volchenboum SL, Cox SM, Heath A, et al. eds. Data commons to support pediatric cancer research. American Society of Clinical Oncology Educational Book. Presented at the American Society of Clinical Oncology Meeting; Chicago, 2017.
- Volchenboum SL, Cox SM, Heath A, et al. Data commons to support pediatric cancer research. *Am Soc Clin Oncol Educ Book.* 2017;37: 746–752.
- Pinto NR, Applebaum MA, Volchenboum SL, et al. Advances in risk classification and treatment strategies for neuroblastoma. *J Clin Oncol.* 2015; 33:3008–3017.
- Lawler M, Siu LL, Rehm HL, et al. All the world's a stage: facilitating discovery science and improved cancer care through the Global Alliance for Genomics and Health. *Cancer Discov.* 2015;5:1133–1136.
- Rozenblatt-Rosen O, Stubbington MJT, Regev A, et al. The Human Cell Atlas: from vision to reality. *Nature.* 2017;550:451–453.
- Sioutos N, de Coronado S, Haber MW, et al. NCI Thesaurus: a semantic model integrating cancer-related clinical and molecular information. *J Biomed Inform.* 2007;40:30–43.
- Gen3 Project Team. Gen3 Documentation. 2018 [cited March 20, 2018]. Available at: <https://uc-cdis.github.io/gen3-user-doc/>. Accessed April 13, 2018.
- Huston G, ed. Interconnection, peering, and settlements. *proc INET*; 1999.
- Saltzer JH, Reed DP, Clark DD. End-to-end arguments in system design. *ACM Transactions on Computer Systems (TOCS).* 1984;2:277–288.
- Busse S, Kutsche R-D, Leser U, et al. Federated information systems: concepts, terminology and architectures. *Forschungsberichte Fachbereichs Informatik.* 1999;99:1–38.
- Landrum MJ, Lee JM, Riley GR, et al. ClinVar: public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* 2013;42:D980–D985.
- Consortium APG. AACR Project GENIE: powering precision medicine through an international consortium. *Cancer Discov.* 2017;7:818–831.
- Schilsky RL, Michels DL, Kearbey AH, et al. Building a rapid learning health care system for oncology: the regulatory framework of CancerLinQ. *J Clin Oncol.* 2014;32:2373–2379.