

Project Management Plan for the Center for Strategic Scientific Initiatives Data Coordinating Center

1. Project Scope

1.1. High level Project Scope

The Center for Strategic Scientific Initiatives (CSSI) is an innovations center within NCI supporting development programs critical to nearly all cancer and biomedical research communities and research programs in new higher risk areas that may question existing paradigms and lead to hypothesis testing. Leidos Biomedical Research (LBR) supports CSSI by providing direct project management support of several defined project efforts and thus can provide timely insight into potential cross-organizational emerging needs. One of the needs identified is an integrated data store that can manage disparate data types generated in the multiple cell or tissue characterization projects currently supported by CSSI. Having a single integrated data store, or Data Coordinating Center (DCC), across all projects will reduce the need to fund a separate data coordinating center (DCC) for each project, as is now the case. Moreover, the availability and maintenance of a common DCC reduces potential project continuity and execution risk when funding cycles for the DCC are not aligned with the data-generation aspects of each project. In addition, developing a uniform way to represent data across CSSI research projects will provide a mechanism that be used by future projects ensure the ability to share their data with the research community. Lastly, but most importantly, integrated management of the data sets across projects will make the data more shareable and therefore more usable by the cancer research community.

To continue to enlarge the general understanding of basic cancer biology, a standard flexible data coordination framework is needed to house multiple subtypes of biologic data generated in support of cancer research by the CSSI. These include, but are not limited to, data from cell lines, animal models, animal/human xenografts, as well as clinically-derived samples. Data of this breadth requires careful attention to the specifications and annotations to make the data useful to cancer researchers. The genomic research community has continually evolved over the last several years to develop a qualified standard that can be used to cross-compare experiments and samples; experience gained in this process can be used to guide the gathering of requirements for a more flexible set of standards that will need to be developed for the CSSI data sets. The Proteomic characterization community is also moving towards adoption and use of standardized data formats for storing and sharing experimental results. Researchers performing physical characterization of biologicals using diverse samples are currently developing methods and formats that will allow the data to be shared with the larger research community. Therefore, projects in progress to house purely genomic data, like the Genomic Data Commons (GDC), funded by the office of Cancer Genomics, may be too inflexible to accommodate the wide-ranging nature of data produced by projects supported by CSSI.

Investigators must be able to accession data into the store, provide a minimal set of metadata that describes the study that was performed, the type of experiment or assay performed, the input and output information whether observation or calculated and the pedigree and/or provenance of the input samples. Information must also be provided about any controls or standards used and the file types to be included in the data store. Information submitted must be validated to ensure that it meets the criteria for the repository on several levels.

Additionally, researchers should be able to access data in the store either for inspection or for download. They should also be able to query across submissions for data with defined characteristics in which they are interested.

1.2. FNLCR Project Scope & Objectives

This project will be comprised of three (3) sequential and interdependent phases. Phase 1 is the analysis phase and its scope is described below. Phase 2 is the first of two software development phases, its output will be a pilot DCC, available to the public, that includes a set of basic functions as well as a subset of CSSI data sets (the proposed data subset is shown in Table 2) Phase 3 is a second software development phase, in this phase additional DCC features will be developed and additional data sets will be included.

This document describes the project management plan for phase 1, the analysis phase. This plan will be continually improved through throughout the course of phase 1 execution leading to a Phase 2 plan. The output of this phase will be a report that documents all of the data, data types, and data formats and data file sizes for the current set of projects originating from the Center for Strategic Scientific Initiatives (CSSI). Phase 1 and phase 2 will focus on CSSI projects (including projects initiated under CSSI) managed by LBR and the larger CSSI portfolio may potentially be added in phase 3. The report will also include a set of preliminary descriptions of how users of the DCC would interact with the system. These use cases, along with the description to the data to be included in the DCC are necessary preconditions not only for the software development phase but also for the development of a framework that outlines the commonalities and differences in the various data sets. This cross-project evaluation of the data is essential because it will be used to determine how the data from the different experiments will need to be stored and related to each other in the CSSI DCC. It should be noted that although phase one is focused on characterization of the data and the potential users, these items will be continued components for the subsequent software development phases.

The first step in the development of the CSSI DCC is to review the breadth of the data developed by completed and anticipated CSSI (and PSON) contracted projects. These data sets, as listed in the table below, will be inventoried and characterized with respect to data types, file types, file sizes and the metadata describing the data and the experiments used to generate the data. This information is essential for the planning of the development of the software system to store the DCC data. Although each individual LBR project manager may have a good understanding of the data types generated for their project, a good understanding of the breadth and gaps of this across all LBR CSSI projects does not currently exist. The delivery of the final report (with the data inventory, use cases and wiki) will be the foundation for future developmental work that will lead to a cohesive data model and data requirements for future research initiatives.

In addition to the technical analysis of the data depth and breadth mentioned above, it is also essential to understand how the user community will want to use the data to support scientific hypothesis evaluation. The team will work with the Stakeholders to determine types of functionality that should be included in a data sharing environment and to develop descriptions of how users would interact with the system, or use cases, for the resource. Additionally, the project team has considerable expertise that will be leveraged to provide the use cases as well as do a preliminary characterization of them. Due to the short timeframe for this phase one evaluation, general users will not be asked to provide their input, with the exception of two experts in the field who can provide high level assistance to this phase. This will be more beneficial in phase 2 when a better idea of the framework has been established. Examples of descriptions would be: as a DCC user I want to be able to download all the experimental data and metadata for a particular experiment, or, as a DCC user I want to search across the genomic data for all the PSON cell lines (or only breast cancer cell lines) for a particular gene (or mutation) of interest. This evaluation phase of the project is the essential first step for the software development required to design the CSSI DCC, because it will outline the best way to store, annotate and distribute the data within the repository. Phases 2 and 3 of this project will occur over the next two fiscal years, however exact schedules cannot be outlined until this evaluation phase is complete. However, it is envisioned that phase 2 and 3 will be approximately one year each.

The combination of the data inventory and the assembly of descriptions of how users would interact with the data is the essential first step in the development of the requirements for the DCC software system including the hardware and storage requirements needed for the DCC.

1.3. FNLCR Key Project Elements

Leidos Biomed PM	ANDREW QUONG andrew.quong@nih.gov 301.360.3401
Leidos Biomed CCB Rep	KATHY TERLESKY Kathy.terlesky@nih.gov 301.846.3916
YT Number	YT14-072
Project ID – Severable funding	001.012.0002.0003
Project Type	Strategic Pilots Incubator
SOW	Design and development of a software system to store, validate and distribute cancer research data produced by projects originating from the Center for Strategic Scientific Initiatives (CSSI).
Contract Amount	\$194,004, Phase 1
Contract Start Date	August 6, 2015
Period of Performance	August 6, 2015 – September 30, 2015
LBR Subcontracts Administrator	Not Applicable
NCI PO (Program Officials)	MICHELLE BERNY-LANG michelle.berny-lang@nih.gov 301.496.1045
COR	JIM CHERRY jim.cherry@nih.gov 301.846.6811 DIANNA KELLY dianna.kelly@nih.gov 301.846.5189
CO	STEPHEN DAVIS stephen.davis2@nih.gov 301.846.1112
NCI Administrative Officer	BRIAN KNESEL brian.knesel@nih.gov 301.402.2705

2. Project Requirements

Table 1. Requirements for the analysis phase of the DCC (through September 30, 2015).

	Requirement	Response
1a	Inventory and characterize current data sets generated by LBR CSSI projects.	Current data and upcoming data will be characterized by file type, file size, data types, and compiled into the final report
1b	Evaluate CPTAC data	CPTAC currently uses a DCC managed by an outside organization. We will coordinate with the NCI CPTAC Program Manager on an appropriate data analysis approach for inclusion in this project..
2	Develop use cases for the DCC to support development of functional requirements for the software development in phase 2	Potential users(data generators, and scientists) will be queried about system functionality and compiled into the final report
3	Purchase hardware to be used as a scalable storage system	The existing data from CSSI and PSON will be used to evaluate the initial hardware needs for the DCC.

In order to develop a plan for the Data Coordinating Center, the data that is currently generated (or has been generated) by CSSI needs to be characterized. Characterization of the data will include evaluating data types, file types, file sizes, and clinical and experimental information associated with the data.

Table 2 includes a list of all of the current LBR yellow tasks associated with CSSI. For this effort it is important to distinguish projects that generate research data, projects that generate biospecimens and biospecimen collection and quality data, and projects that develop algorithms to produce computational data. The scope of this project is focused on research, or wet lab data and these projects are asterisked in the table below.

The Physical Sciences in Oncology Network (PSON) was initiated under CSSI and some of the projects funded by PSON have already been and completed and are therefore proposed to be part of the pilot data set that will be included in Phase 2 of this project. In addition to the PSON data, inclusion of data from Physical characterization of pre-analytical variables in biospecimens, another completed project, is also being proposed as part of the pilot. These data sets contain cell line characterization data and the characterization of individual cells from liquid biopsies and provide a pilot data set representative of the types of disparate data types that will need to be stored as part of the CSSI DCC.

Table 2. LBR Projects associated with CSSI (includes Physical Sciences Oncology Network projects).

Yellow Task #	Center Name	Project status
11-212NS	*Physical Characterization of Parameters in Biospecimens ¹	complete
11-225	*Genomic Analysis of PSON Cell Lines ¹	complete
12-095NS	CPTAC Informatics Tools Development	Project does not generate research data
12-129NS	*HCS of Physical Based Properties in Biospecimens Phase 2	Start up
12-138	*Physical Based Properties of PS-OC Cell Line Panel ¹	complete
12-140NS	CPTAC Tissue Procurement	Biospecimen data
12-158NS	Support for nanosensors pilot project	Out of scope for this effort
13-089NS	CPTAC Tissue Procurement	Biospecimen data
13-118NS	CPTAC Pilot Studies	start up
13-125NS	*Immuno MRM Assay Pilot	start up
13-134NS	*CaSIX	start up
14-083	*Proteomic Characterization of PS-OC Cell Line Panel	start up
14-103NS	CPTAC Computational Omics	Project does not generate research data
14-107NS	*Thrombosis in Cancer Patients	start up
14-115	Nanoinformatics	Out of scope for this effort
14-141	CPTC Web Support (ABCC)	Out of scope for this effort
14-144NS	*CPTAC Immuno-MRM Assays	start up
14-152NS	Prospective Biospecimen Collection Phase 2	Biospecimen data

*Projects that produce research data (wet-lab data)

¹ Project data that will be included as part of the pilot CSSI DCC proposed for phase 2 of the project

3. Technical Approach & Schedule

3.1. Approach: The approach for the analysis phase of the DCC project is outlined in Figure 1.



Figure 1. High level outline of approach to analysis for DCC

3.2. Work Breakdown Structure

1. Phase 1a – Gather Use cases
 - a. Create JIRA or WIKI page to house information
 - i. Project documents
 - ii. Meeting minutes
 - iii. Monthly reports
 - iv. Use cases
 - v. Data inventory
 - b. Query data generating individuals who would input data
 - c. Query scientists that would use data
 - d. Query stakeholders
 - e. Prioritize use cases
 - f. *This is a deliverable that is due September 30, 2015*
2. Phase 1b - Data Characterization Inventory
 - a. Identify data that has already been generated by CSSI projects
 - b. Identify data that will be generated by currently funded projects
 - c. Characterize data
 - i. Data types
 - ii. File types
 - iii. File sizes
 - iv. Clinical information, if applicable
 - v. Experimental information
 - vi. Any other associated information
 - d. *This is a deliverable that is due September 30, 2015*
3. Phase 1c – Initial internal prioritization of use cases
 - a. The internal LBR team will organize the use cases based on categories that indicate their perceived importance to users
 - b. *This is a deliverable that is due September 30, 2015*
4. Hardware purchase
 - a. Evaluate data needs for potential pilot data
 - b. Purchase hardware that is easy to scale and that can be used for future phases
 - c. *This is a deliverable that is due September 30, 2015.*
5. Final Report
 - a. The final report will include the information that was gathered about the CSSI data and use cases.
 - b. The wiki page will contain all necessary documentation and be available for phase 2.
 - c. *This is a deliverable that is due September 30, 2015.*

3.3. Milestone Schedule

Table 3. Milestone and deliverable schedule

Tasks	14-Aug	21-Aug	28-Aug	4-Sep	11-Sep	18-Sep	25-Sep	30-Sep
Phase 1a Collated list of Use Cases						X		Delivery
Use cases								
Data dictionary								
Lessons learned from other DCCs								
Phase 1b Inventory and Characterize the CSSI (and PSON) data sets						X		Delivery
Data in hand								
PSON Genomic Characterization Data								
PSON Physical Characterization Data								
CSSI CTC Data								
Pending data								
PSON Proteomic Characterization Data								
CTC phase II data								
Thrombosis data								
MRM immuno assay data								
CPTAC Data								
Prospective Data								
Phase 1c Preliminary categorization of use cases						X		Delivery
Final Report							X	Delivery
Draft Report								
Review and Revise Report								

The “X” indicates when the information should be completed and provided to the TPM for review and integration into the final report. The Final report will be provided to the team for review by September 25, 2015. The colors indicate each phase of the project and portions will be occurring concurrently.

3.4. Assumptions

1. The LBR CSSI projects that have a wet lab component listed in Table 2 are the only projects with data sources considered for this DCC pilot.
2. Completion of the phases in this project is required prior to development of a full project plan for DCC pilot development.
3. Data from completed projects is available and provided for use.
4. There will be a sufficient number of realistic use cases provided.
5. The ITOG can house the hardware components required for the DCC.

4. Deliverables

- **Data Assessment Deliverable:**
 - Each project under the CSSI umbrella that generates data will be evaluated to determine the types of experimental data that the project will generate, the types and sizes of files that the project will generate, and the information that accompanies the data to describe the data and the experiments (metadata) so that it can be effectively shared to produce a data and metadata inventory. This inventory will be used to make specific

recommendations about the size and type of storage that will be needed to support the CSSI DCC.

- **DCC Use Case Deliverable:**
 - Stakeholders and potential users will be queried in order to gather user descriptions of how they would like to use the data, or use cases. These use cases will be developed iteratively through communication with project stakeholders and subject matter experts in cancer cell biology, molecular cell biology, biospecimen science, genomic and proteomic analysis, metadata analysis, data science and DCC design and development. The use cases delivered at the end of FY15 will need to be prioritized by management before they can be used as input to the planning phase for the software system planned for early in FY16.
- **DCC Hardware:**
 - The hardware required for the analysis of the existing data during the evaluation phase, this hardware will also be able to be used in the subsequent development phase.
- **Meeting minutes**
 - Meeting minutes from major meetings provided within 2 business days of meeting completion
- **Monthly Reports**
 - Summary report of progress to milestones in revised project plan delivered monthly

5. Project Budget

Table 4. Estimated Project Costs

Category	Estimated Cost
Labor and Fringe	\$114,044
<i>Non-Labor Costs</i>	
M&S	
Consultants	\$5,000
Travel	
Hardware	\$75,000
<i>Total Non-Labor Costs</i>	\$80,000
Total Estimated Cost	\$194,044

Table 5. Labor cost composition

Role/Title	LOE	Estimated hours
Project Lead SME (cancer biology, genomics/proteomics)	50%	168
SME (cancer biology, research data standardization, bioinformatics)	30%	104
SME (cancer biology, clinical data, biospecimen metadata)	50%	168
SME (cancer biology, clinical proteomic data, biomedical informatics)	50%	168
Bioinformatics Data Analyst (DCC experience)	50%	168
Systems and Data Architect	20%	68
CSSI TPM	100%	336
Other SMEs and project managers for YT projects	50%	168
Totals	400%	1348

6. Project Team

Table 6. Project team roles and contact information.

Name	Role	Email	Phone
Andrew Quong	Project Lead SME (cancer biology, genomics/proteomics) Core	Andrew.Quong@fnlcr.nih.gov	301.360.3401
Debra Hope	Project Co-Lead SME (cancer biology, research data standardization, research data modeling, research informatics pathology informatics) Core	Deb.Hope@fnlcr.nih.gov	240.276.5777
Mary Anderson	SME (cancer biology, bioinformatics) Core	Mary.Anderson@fnlcr.nih.gov	240.276.5258
Linda Hannick	SME (cancer biology, clinical data, biospecimen metadata)	Linda.Hannick@nih.gov	202.674.1101
Anita Undale	SME (cancer biology, clinical proteomic data, biomedical informatics)	Anita.Undale@fnlcr.nih.gov	240.276.5888
Joan Pontius	SME Bioinformatics metadata (DCC experience) Core	PontiusJ@mail.nih.gov	240.276.6145
Charles Shive	Systems and Data Architect	Charles.Shive@nih.gov	240.276.7616
Corinne Zeitler	TPM Core	Corinne.Zeitler@nih.gov	301.846.6571

*Core members of the project team are highlighted in bold font.

6.1. Team member availability

This project team has reviewed the hourly requirements and will be able to provide their time during FY15. We have determined that there is sufficient depth of expertise and availability across the entire team to cover the required scope for Phase 1 and Phase 2. The greatest risk to personnel availability for this project is

the project-specific demands of their primary project. Because this project emphasizes the linkage points across multiple CSSI projects, we will manage stability of a core team (Project Lead, Project co-Lead, Bioinformatics SMEs, TPM) while also recognizing that input from the SMEs on a rotational basis may be the best means to manage collaboration and project execution.

7. Subcontracting Approach

No subcontractors are currently planned for this project. In the event that this changes during the course of project execution in Phase II, this section will be updated.

8. Communication Plan

A [WIKI page](#) has been created to house collected information and documents that will be available to the project team and customer.

The external users listed in Table 7, who are experts in their field, will be contacted in order to provide f and possible use cases. This will provide important information about lessons learned in the field, as well.

Table 7. External users

Name	Title(s)	Affiliation	Expertise
Cathy H. Wu, Ph.D.	<ul style="list-style-type: none"> Edward G. Jefferson Chair of Bioinformatics & Computational Biology Director, Center for Bioinformatics & Computational Biology (CBCB) Director, Protein Information Resource (PIR) Professor, Computer & Information Sciences Professor, Biological Sciences 	University of Delaware	Bioinformatics and Computational Biology: Biological Text Mining, Biological Ontology, Computational Systems Biology, Protein Structure-Function-Network Analysis, Bioinformatics Cyberinfrastructure
Aydin Tozeren, Ph.D.	<ul style="list-style-type: none"> Distinguished Professor and Director, Center for Integrated Bioinformatics, School of Biomedical Engineering, Science & Health Systems Co-Director of the Greater Philadelphia Bioinformatics Alliance 	Drexel University	Biomarker identification using large scale compilation of microarray data, mathematical modeling of biological systems, cell signaling and adhesion, biomechanics

8.1. Stakeholder Identification and Involvement

Table 8. Stakeholder Identification and Involvement

Name	Title	Affiliations	Type of Communication	Frequency
Rachana Agarwal	Technical Project Manager	Leidos	Phone, email or in person	As Needed
Mary Anderson	TPM II	Leidos	Teleconference	weekly
Michelle Berny-Lang	Project Manager	CSSI	Teleconference	monthly
Braulio Cabral	Director , CBIIT Technical Operations Support	Leidos	Phone or email	As needed
Jim Cherry	Assistant Project Officer	OSO	Teleconference	monthly
Mariam Eljanne	Program Director	DCB (PSON)	Teleconference	monthly
Mike Espey	Program Director	DCB (PSON)	Teleconference	monthly
Sharon Gaheen	TPM II, Genome Data Commons	Leidos	Phone or email	As needed

Name	Title	Affiliations	Type of Communication	Frequency
Dan Gallahan	Deputy Director	DCB	Teleconference	monthly
Emily Greenspan	Program Director	CSSI	Teleconference	monthly
Sean Hanlon	Program Director	DCB (PSON)	Teleconference	monthly
Linda Hannick	TPM II	Leidos	Teleconference	weekly
Debra Hope	Director, Data Science	Leidos	Teleconference	weekly
Mark Jenson	Director, Genomic Data Programs	Leidos	Teleconference	As needed
Dianna Kelly	IT Program Manager	OSO	Teleconference	monthly
Chris Kinsinger	Program Manager	CSSI	Teleconference	monthly
Jerry Lee	Deputy Director	CSSI	Teleconference	monthly
John Otridge	TPM I	Leidos	Teleconference	weekly
Joan Pontius	Bioinformatics Analyst IV	Leidos	Teleconference	weekly
Andrew Quong	Director, Partnership Development Office	Leidos	Teleconference	weekly
Henry Rodriguez	Director	CSSI	Teleconference	monthly
Charles Shive	TPM III, Informatics Architect	Leidos	Teleconference	weekly
Dinah Singer	Director	DCB	Teleconference	monthly
Eric Stahlberg	Bioinformatics Scientist III	Leidos	Teleconference	weekly
Kathy Terlesky	Director, Project Management Operations	Leidos	Teleconference	weekly
Gordon Whiteley	Director Clinical Proteomics	Leidos	Phone, email or in person	As needed
Corinne Zeitler	TPM I	Leidos	Teleconference	weekly

8.2. Teleconferences and meetings

Weekly teleconferences are held with members of the project team to review status of the project and any issues that arise. A monthly meeting with the stakeholders will be held to update everyone on the status. Email and phone calls will be used as necessary to communicate with stakeholders outside of these meetings.

8.3. Stakeholder Project Management Approach

For the purposes of this project, the CSSI Project Manager should interact directly with the Project Lead. The LBR Project Lead may also be asked to interact with the government CORs as needed for the DCC project. Figure 2 highlights the flow of information between the government (blue) and the LBR project team (purple). The project lead, Andrew Quong, will be responsible for conveying information to the government and act as the point of contact for the customer. The project co-lead, Debra Hope, will act as the Leidos Project Team lead, and will facilitate communication between the internal project team since it spans so many different projects and departments.

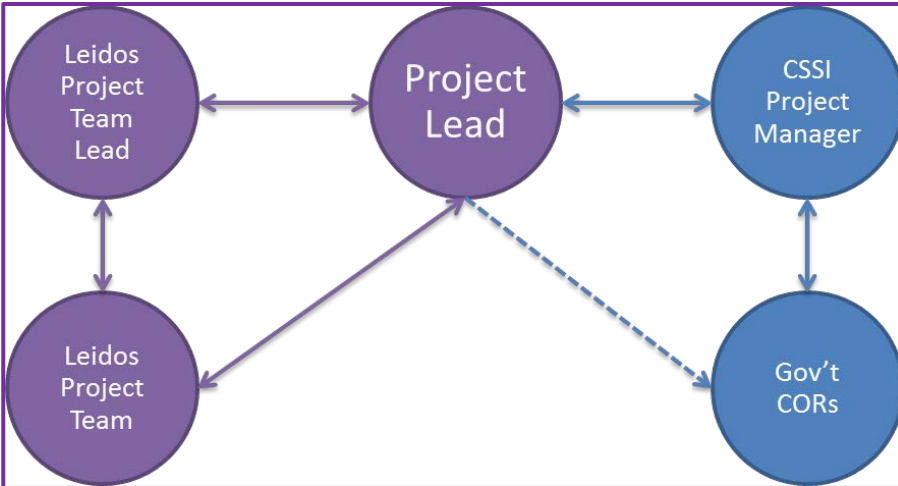


Figure 2. Stakeholder communication plan.

9. Risk Management

9.1. Risk Identification and Risk Management

1. There may not be enough use cases identified. Impact to technical and schedule.
2. It may not be possible to inventory all CSSI data sets due to nonresponse of participating groups. Impact to technical and schedule.
3. ITOG is not able to house the hardware for the pilot. Impact to cost and technical.
4. PSON data may not be available for use in the pilot.

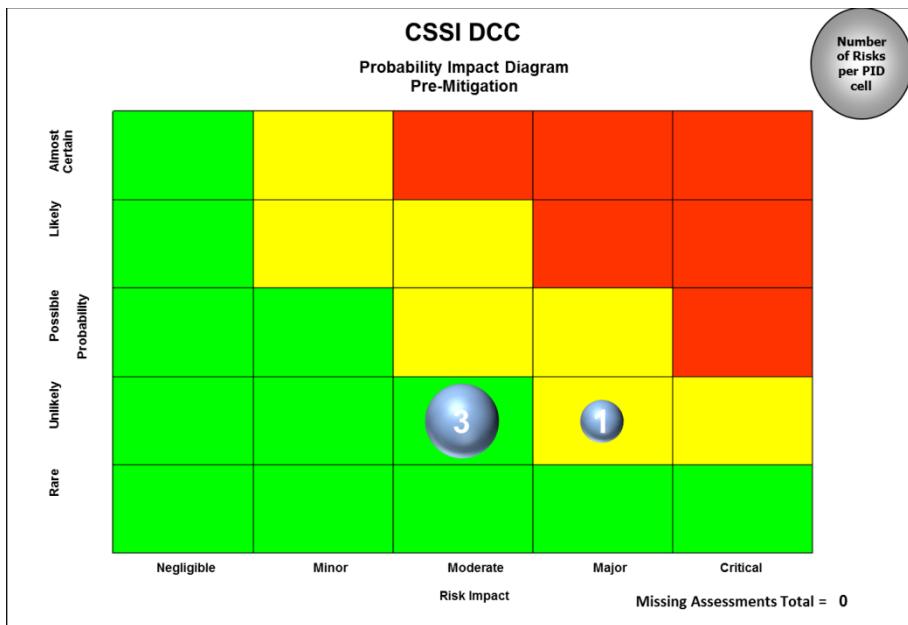


Figure 3. Probability Impact Diagram (PID), Pre-Mitigation

9.2. Risk Register Tools

Table 9. Risk Register

Risk Evaluation Category	Risk Statement	Impact Statement	Risk Owner	WBS Element	Risk Probability	Risk Impact	Risk Urgency	Risk Response	Mitigation Summary
Technical: Scope of Work	Not enough use cases identified	May not be enough information to meaningfully create system requirements in phase 2 or may delay phase 2 to collect more cases			Unlikely	Moderate	Immediate	1-Watch	Will monitor collection of use cases in order to ensure there is a reliable number
Technical: Scope of Work	Unable to inventory all CSSI data sets due to nonresponse of participating groups	May not be enough information to meaningfully create system requirements in phase 2 or may delay phase 2			Unlikely	Moderate	Immediate	1-Watch	Will monitor collection of data to keep on schedule
Technical Baseline Design	ITOG is not able or willing to house hardware for pilot	may impact schedule and cost			Unlikely	Moderate	Immediate	1-Watch	Other options exist to host the hardware, but would like to keep it in house
Technical: Scope of Work	Disparate data cannot be associated within the DCC in a way that is meaningful to users.				Unlikely	Major	Long term	2-Plan Risk Response	Develop a plan to re-evaluate framework or requirements for system

10. Quality Management Plan

The quality metrics are used to assess if the deliverables are meeting project requirements. Quality will be monitored throughout the project. The Project Management Plan and risk register will be updated as needed. Due to the short timeframe, it is imperative that milestones are met and monitored closely. Milestones will be discussed at weekly meetings. The Project Management Plan will be available and, along with the risk register, will be updated as needed.

1. Milestones must be met on time in order to have a final report prepared.
2. Success of the project will be measured by the ability to provide a final report outlining the inventory of CSSI data and characterized use cases, as well as the hardware purchased for the pilot project.

11. Glossary

Accession data – enter new data into the data store

Clinically-derived samples – samples derived from patients

Data coordination framework – a model used to compile the data and meaningfully associate it with related data

Data store - repository of a set of data.

Metadata – the information or set of data that describes the raw or processed data in the data store. This may include, but is not limited to, clinical, experimental or processing information.

Use cases – list of actions or goals that define the interactions between the user and the system

Wet lab – Lab where biological or chemical samples are manipulated.