

## 4.2 - Life Sciences

Current Working Draft

Semantic Infrastructure 2.0 needs to address metadata and terminology related requirements from the life sciences domain. This will enable interoperability both between different sub-domains within life sciences, and between life sciences and other domains in caBIG® such as clinical trials and electronic health records.

While life sciences will leverage common semantic functionalities such as Enterprise Conformance and Compliance Framework (ECCF) registry, modeling, forms, behavioral semantics, terminology and value sets, there are aspects specific to the life sciences domain that need to be addressed, including but not limited to:

- High change semantic environment where novel concepts which have not been previously characterized need to be described
- Computational or analytical workflow compositions processing raw data to derive knowledge
- Description of statistical processes related to computational processes to achieve the above
- Working with the platform, enabling semantic description of raw data potentially of large volume (for example next-gen sequencing, imaging data)
- Support of provenance to trace data acquisition and data ownership and also to achieve reproducibility of analytical results
- The caBIO ECCF service specification project on molecular and pathway annotation services from the Integrated Cancer Research (ICR) Workspace

This section highlights some key use cases that depend on data semantics. These use cases provide a representative set to capture the requirements of the life sciences domain. A comprehensive set of all life sciences use-cases can be found at on the [ICRi WG GForge wiki archive](#). This section includes the following:

- [Discovering a biomarker](#)
- [Finding biomaterial to validate a biomarker](#)
- [Extending the use of a biomarker](#)
- [Exploring predictive power of gene expression in breast cancer metastasis](#)
- [Oncologists in formulating ideas for new clinical studies](#)
- [Multi-Center Ancillary Study in the context of a Consortium Clinical Trial \(extension from Enterprise Use Cases\)](#)
- [Overlay of protein array data on the regulatory pathways with links to patient and cell culture data.](#)
- [Animal model use case](#)
- [An outside researcher requests access to a consortium's Prostate SPORes Federated Biorepositories, eleven instances of caTissue Suite independently maintained and managed](#)
- [High throughput sequencing using DNA sequencing to exhaustively identify tumor associated mutations](#)
  - [Variant: Version A](#)
  - [Variant: Version B](#)
  - [Variant: Version C](#)
  - [Nanoparticle Delivery System](#)
  - [Nanoparticle Delivery System B](#)
  - [Scenario 13](#)
  - [Scenario 13 B](#)
  - [Scenario 14](#)
  - [Scenario 14 B](#)
  - [Scenario 14 C](#)
- [NanoParticle Ontology](#)



### Note

However, part of the Infrastructure Inception activities include Prototyping [Orchestrations and/or Choreographies](#) (including Life Science workflows) as well as outreach to communities to address other major use cases and requirements. The life sciences communities are engaged with the Roadmap inception efforts now on the following:

- [Lymphoma "workflow" use cases and requirements](#)
- [Dynamic Extensions use cases and requirements](#)
- [Imagining Use cases and requirements](#)

Refer to the pages listed for the use cases and requirements gathering activities. These will be moved to relevant sections in the Roadmaps when mature as use cases, requirements and resulting architecture design.

### Discovering a biomarker

A scientist is trying to identify a new genetic biomarker for HER2/neu negative stage I breast cancer patients. Using a caGrid-aware client, the scientist queries for HER2/neu negative tissue specimens of Stage I breast cancer patients at LCCC (University of North Carolina Lineberger Comprehensive Cancer Center-NC Cancer Hospital) that also have corresponding microarray experiments. Analysis of the microarray experiments identify genes that are significantly over-expressed and under-expressed in a number of cases. The scientist decides that these results are significant, and related literature suggest a hypothesis that gene A may serve as a biomarker in HER2/neu negative Stage I breast cancer. To validate this hypothesis in a significant number of cases the scientist needs a larger data set, so the scientist queries for all the HER2/neu negative specimens of Stage I breast cancer patients with corresponding microarray data and also for appropriate control data from other cancer centers. After retrieving the microarray experiments the scientist analyzes the data for over-expression of genes A.

### Finding biomaterial to validate a biomarker

In scenario 1, the scientist has validated a biomarker based on available microarray experiments provided by various cancer centers. Now, the scientist would like to request biomaterial in the form of formalin-fixed, paraffin embedded tissue specimens from patients with the appropriate clinical outcomes. The scientist would like to validate the genetic biomarker in a different series of cases, this time using a different technique such as immunohistochemistry. The scientist queries for the presence of appropriate tissue using a caGrid-aware client and for the appropriate contact information of the person(s) responsible for the tissue repository. The scientist contacts the person(s) to begin the protocol for retrieving biomaterials.

### **Extending the use of a biomarker**

The scientist would like to check if genes A could also be used as biomarker for other types of cancer. The flow of events will be similar to Scenario 1 with the exception that the specimen query will not be restricted to Stage I breast cancer patients.

### **Exploring predictive power of gene expression in breast cancer metastasis**

The scientist would like to explore if gene expression patterns can predict how breast cancer will metastasize. The scientist queries all the specimens of breast cancer patients from other cancer centers where their metastasis sites are liver, bone and brain. The scientist then retrieves the corresponding microarray experiments for these specimens. The scientist analyzes the microarray experiments to explore for a correlation between expression profiles and metastasis sites.

### **Oncologists in formulating ideas for new clinical studies**

The oncologist often wants to first find out the answers to questions such as: How many patients have been seen at our institution with disease x? How does that compare with other institutions? What is the average survival of patients with disease x? How is it different if they are treated with drug x or y? How many patients are there with disease x and TNM stage y at diagnosis? How many patients with disease x relapse after treatment y? This use case is about enabling oncologists to ask these exploratory questions of their clinical databases as well as those at other institutions accessible on caGrid.

### **Multi-Center Ancillary Study in the context of a Consortium Clinical Trial (extension from Enterprise Use Cases)**

The following are the steps in a translational research scenario.

- Within a consortium of cooperating institutions an investigator conducts a search across the clinical data repositories to investigate the feasibility of a potential clinical research idea.
- Within the consortium, the research question is circulated to gauge interest.
- Members of the consortium discuss the research question and approve it as viable.
- The research question is formalized by the coordinating center into a clinical research ancillary protocol for validation of a biomarker as predictive of tumor shrinkage in the context of treatment using an investigational agent and posts to the consortium for consideration (for example, Do patients with a particular marker respond better to treatment with the agent?).
- Consortium member sites choose to join the protocol and agree to accrue patients on to it, collect bio samples from each participant and ship a defined set of bio samples to Central Pathology.
- Participating consortium sites each submit common consent forms, case report forms, and boilerplate Material Transfer Agreements (MTAs) to the appropriate local regulatory offices.
- The protocol meta data, case report forms, standard operating procedures, MTA documents, and other items as needed are finalized and disseminated to each participating site.
- Participants are screened on the basis of eligibility by study coordinator at each site.
- Patients are accrued (by physician or patient self referral) by local staff onto the protocol at each site, and the accrual event is reported to the coordinating center.
- Bio samples are collected and relevant clinical annotations including tumor measurements are collected at the appropriate time points as indicated in the protocol (these are for calculating the primary end point, tumor shrinkage).
- Follow-up appointments are scheduled as specified in the protocol.
- Bio samples are periodically sent to Central Pathology.
- Central Pathology re-labels the samples to hide the source and identity.
- Central Pathology sends out batches of collated bio samples to each of the participating biomarker assay labs.
- A basic scientist at biomarker lab submits the result of biomarker assays.
- Patients are followed for three years from primary treatment date. An annual follow-up visit occurs and a blood sample is taken. Additional clinical annotations are collected.
- The trial closes, and all the data are made accessible to the statisticians.
- A statistician communicates the clinical significance of and evidence for biomarker response prediction.
- A clinical researcher, basic scientist and statistician write a scientific paper reporting the results.
- Data are made available according to funding agencies requirements.

## Overlay of protein array data on the regulatory pathways with links to patient and cell culture data.

A clinical research scientist wants to be able to predict the efficacy of tyrosine kinase inhibitors as cancer chemotherapeutic agents. The fact that many oncogenes are tyrosine kinases would predict that such agents should be effective, but several have been synthesized and tested in clinical trials, and the results have been disappointing in the extreme, with more cases of tumor growth stimulation than inhibition. The clinician hypothesizes that these unexpected effects are the result of regulatory feedback loops.

To test this hypothesis, he requires software tools for modeling regulatory pathways. In addition, he needs to determine the state of such pathways in different patients by measuring the state of phosphorylation of the elements (proteins) of these pathways using reverse phase protein arrays. Because the consequences of treating the wrong patient with the wrong agent are so severe, the response of the tumor to the inhibitors will be tested in vitro, on cell cultures established from tumor biopsies. However, biospecimens and data from those patients who participated in clinical trials of these reagents before their ineffectiveness was appreciated is also available.

Outputs measured on these cultures and biospecimens will include growth rate (determined by flow cytometry or visual counting of cells at different time points, extent of cell death determined similarly, photomicrographs, reports of microscopic observations by trained investigators, rate of DNA synthesis measured by radioisotope or fluorescent labeled precursor uptake and incorporation, and staining with various immune reagents followed by high throughput robotic microscopy and automated image analysis. To develop an understanding that will result in giving the correct drugs to the correct patients, data from the protein arrays will be overlaid on the regulatory pathways and linked to patient and cell culture data.

## Animal model use case

The following are the steps describing the scenario.

- A bench scientist chooses candidate glioblastoma genes using human Genome-Wide Association Studies (GWAS), for example, The Cancer Genome Atlas (TCGA).
- The scientist also uses pathway analysis to postulate how multiple "hits" may be involved in tumorigenesis, to direct design of genetically altered mice.
- The scientist uses targeted gene transfer to deliver mutated genes to inbred mice.
- Using inbred mice provides uniform genetic background in which researcher can also investigate mutated candidate genes in conjunction with other gene knockouts.
- The scientist finds mutated gene x expressed in gene y knockout mouse results in glioblastoma development that parallels human pathology.
- The scientist validates that the model reacts in similar way to current therapeutic treatments.
- The scientist uses a mouse model to test new therapeutic treatments, including combinations of drugs chosen to inhibit multiple pathways.
- Clinical scientists use mouse model results to design clinical trials to treat glioblastoma, incorporating genomic information on patients.

## An outside researcher requests access to a consortium's Prostate SPOREs Federated Biorepositories, eleven instances of caTissue Suite independently maintained and managed

- A Research Fellow at a University has been working at identifying SNPs that might be related to aggressive forms of Prostate Cancer. The Fellow has narrowed the search down to 21 SNPs, and discusses the results with the mentor.
- The mentor just returned from a Bio repository presentation where the mentor learned of the Consortium's federated biobanks built on caTissue, and he mentions that this might be a valuable resource that could aid the research. The mentor suggests that the Fellow contact a colleague at Memorial Sloan Kettering Cancer Center (MSKCC), who is a member of the Consortium.
- The Fellow drafts an email briefly explaining the research and sends it to the MSKCC cancer center member. The Fellow asks about the possibility of searching across the Consortium's federated biobanks for cases that have X and Y and at least 3 years of outcome data. The goal is to collect enough tissue to construct a Tissue Microarray.
- The MSKCC member responds and directs the Fellow to the consortium hub web site where there are details on the policies and procedures for requesting an account to be able to submit a query that would search each of the 11 instances of caTissue Suite. He agrees to be her sponsor.
- The Fellow completes an online form that requests and abstract of the research, the name of the Fellow's institution, the non-profit status of that institution, the name of the sponsoring Consortium member and that person's institution, and a required checkbox indicating that The Fellow has read and agrees to the terms of use.
- The Oversight Committee (OC) of the Consortium Biorepositories has a standing regularly scheduled telephone conference call during which the committee reviews requests to query the federated biorepositories. After each regular call a new set of primary reviewers are elected who will be responsible for thoroughly reading new requests and presenting them to the other members of the OC for a vote.
- The OC has authored a set of appropriate use policy documents against which all requests are measured. The current primary reviewers read the Fellow's application and report to the other members of the OC.
- The OC votes to approve the Fellow's request.
- The Fellow is notified of the OC's decision, and is supplied with an account to the MSKCC instance of caTissue, since this application was sponsored by a member at MSKCC.
- The Fellow logs into caTissue Suite at MSKCC as a researcher, and formulates the parameters of the query. The Fellow submits the query and after a period of time sees a results set that span 8 of the 11 instances of caTissue Suite at the Consortium sites.
- The Fellow uses this information to request tissue from four institutions to build the tissue microarray (TMA).

## High throughput sequencing using DNA sequencing to exhaustively identify tumor associated mutations

This is a basic research use case that easily becomes translational when the output of this use case is used, for example, to identify targets for biomarker studies or drug candidates for clinical trials.

### Variant: Version A

Version A is "Sequencing of selected genes via Maxim Gilbert Capillary ("First Generation") sequencing." *Nature*. 2008 Sep 4 - Epub ahead of print

1. Develop a list of 2000 to 3000 genes thought to be likely targets for cancer causing mutations.
2. As a preliminary (lower cost) test, pick the most promising 600 genes from this list.
3. Develop a gene model for each of these genes.
4. Hand modify that gene model, for example, to merge small exons into a single amplicon.
5. Design primers for PCR amplification for each of these genes.
6. Order Primers for each exon of each of the genes.
7. Test Primers.
8. In parallel with steps 1-7, identify matched pairs of tumor samples and normal tissue from the same individual for the tumors of interest.
9. Have pathologists confirm that the tumor samples are what they claim to be and that they consist of a high percentage of tumor tissue.
10. Make DNA from the tumor samples, confirming for each tumor that quantity and quality of the DNA are adequate.
11. PCR amplify each of the genes.
12. Sequence each of the exons of each of the genes for each tumor and normal pair of DNA samples.
13. Find all the differences between the tumor sequence and normal sequence.
14. Confirm that these differences are real using custom arrays, the sequenome (Mass Spec) technology and biotage or both. (A biotage is pyrosequencing-based technology directed specifically at looking for SNP-like changes.)
15. Identify changes that are seen at a higher frequency than what would occur by chance.
16. Relate the genes in which these changes are seen to known signaling pathways.

Existing tools for each step are:

1) None; a completely manual process. 2) None; a completely manual process. 3) Data is uploaded from the UCSC Genome Browser to Genboree which has modules for all of the required tasks. 4) Same as 3. 5) Primer3 embedded into a local pipeline developed at the HGSC that keeps primers away from repeats and SNPs. Gaps where this pipeline is unable to create primers are filled in by hand. 6) Manual process. 7) Manual process. 8) It is not known how this was done by the HGSC, but caTissue and similar products can be used here. 9) Manual process. The pathology imaging initiative of Tissue Banks and Pathology Tools (TBPT) might fit in here. 10) Manual process. 11) Manual process. Could a Laboratory Information Management System (LIMS) help here? 12) Software provided as part of the ABI sequencer. 13) Combination of custom, ad-hoc software and manual processes. 14) Manual process. 15) Combination of custom, ad-hoc software and manual processes. 16) Manual process. This *should not be a manual process*, but almost always is, or it is of low quality.)

### Variant: Version B

Version B. As above, except globally sequence all genes. *Science* 321: 1807-1812 (2008). Delete steps 1 and 2 and replace step 3 with: 3) Develop a gene model for each of the genes in the Human genome.

### Variant: Version C

Version C. Whole genome sequencing using second generation sequencers. *Hypothetical*.

1. Identify matched pairs of tumor samples and normal tissue from the same individual for the tumors of interest.
2. Have pathologists confirm that the tumor samples are what they claim to be and that they consist of a high percentage of tumor tissue.
3. Make DNA from the tumor samples, confirming for each tumor that the quantity and quality of the DNA are adequate.
4. Sequence each of the sample pairs to the required fold coverage (7.5 to 35-fold, depending on the technology and read length).
5. Map the individual reads to the canonical human genome sequence.
6. Find all the differences between the tumor sequence and normal sequence.
7. Confirm that these differences are real using custom arrays, the sequenome (Mass Spec) technology or biotage or both. (Biotage is a pyrosequencing-based technology directed specifically at looking for SNP-like changes).
8. Identify changes that are seen at a higher frequency than what would occur by chance.
9. Relate the genes in which these changes are seen to known signaling pathways.

Existing tools for each step are:

1) caTissue or similar product. 2) caTissue or similar product pathology imaging tools to be developed by TBPT. 3) caTissue or similar product. 4) Combination of custom, ad-hoc software and manual processes. 5) Proprietary, platform-dependent software, a wide variety of non-caBIG-compatible software packages: Solexa Mapper, Mosaic, 454 Mapper, Velvet Mapper, Solid Mapper (uses a non-standard sequence representation model), Mac. 6) Combination of custom, ad-hoc software and manual processes. 7) Manual process. 8) Combination of custom, ad-hoc software and manual processes. 9) Manual process. This *should not be a manual process*, but almost always is, or it is of low quality.)

## Nanoparticle Delivery System

This is a scenario based on finding a nanoparticle delivery system to target a drug which in its free form causes significant side effects. Sorafenib is a Raf kinase inhibitor that disrupts the key Ras/Raf/MEK/ERK cellular pathway that is up-regulated in renal cell carcinoma, glioblastoma multiforme (GBM), and stomach cancer. The drug has significant side effects and a scientist hypothesizes that nanoparticle-assisted targeted delivery of the drug will reduce the required dosing and its side effects.

A scientist interested in targeting this drug to GBM does research on possible nanoparticle-delivery systems that have the following properties:

- Biocompatibility
- Sufficiently long intravascular half-life to allow for repeated passage through and interactions with the activated endothelium

- The ability to have ligands and proteins conjugated on the surface in multivalent configuration to increase the affinity and avidity of interactions with endothelial receptors
- The ability to have functional groups for high-affinity surface metal chelation or radio-labeling for imaging
- The ability to encapsulate drugs
- The capability to have both imaging and therapeutic agents loaded on the same vehicle

Furthermore, the scientist looks for information on nanoparticles that could potentially target the GBM. Integrin-targeted nanoparticles are identified. Synthesis involves ultraviolet (UV) cross-linking of an  $\alpha v \beta 3$ -integrin-targeting ligand attached to diacytyle phospholipids and a cationic lipid. These are sonicated to form polymerized vesicles and the  $\alpha v \beta 3$ -targeted NP can serve as a scaffold for the attachment of therapeutic agents for imaging and therapy.

The physical characteristics have been determined. These include size, zeta potential, and the relevant IC50. In a cell adhesion assay, the 10 of 19 effect of multivalency on IC50 is also measured. Selectivity was also demonstrated in a receptor-binding assay and it is also shown that the  $\alpha v \beta 3$ -targeted NP is not rapidly cleared from the target tissue. Previous studies have shown this particle to be highly stable, to have no measurable toxicity and to specifically target tumor associated vasculature in GBM when conjugated to GFP. Furthermore the particle has been used as an imaging agent when conjugated with Gd3+ or Iodine2+. The  $\alpha v \beta 3$ -targeted NPsorafenib is synthesized. Sorafenib absorption characteristics are available and the concentration of the drug in the system is determined via spectroscopy methods. Other physical properties are characterized.

## Nanoparticle Delivery System B

This is the Nanoparticle Delivery System scenario extended. The scientist investigates what data sets are available for in vivo use of the drug. A breast cancer xenograph subcutaneous model is found and cell lines from this system are also available. However, toxicity data for the drug in animal models are not publicly available. The scientist contacts the drug manufacturer and begins in vitro testing. PK/PD in vitro tests, including drug uptake, toxicity and effectiveness, are performed in the model system cell lines, and related and control cell lines by comparing the effects of drug alone, nanoparticle alone, and the combination. Next is in vivo testing with three established animal tumor models. The drug alone, nanoparticle alone, and the combination are administered and tumor size (and other parameters) are monitored. Finally efficacy, dosing, and side effects of the current dosing protocol are compared with targeted nanoparticle delivery of sorafenib.

### Scenario 13

This is a scenario based on in vitro profiling of nanomaterial activity. A scientist has created a library of surface-modified nanoparticles with potential as in vivo imaging agents. The scientist would like to use an in vitro approach to gain insight on potential toxicity of these nanoparticles, and exclude those that might be problematic prior to using costly and time-intensive in vivo methods. The mode of administration is considered in selecting a variety of cell types to use in the in vitro assays. Cell cultures are started. Each nanoparticle is added to cultures of each cell type at multiple biologically-relevant concentrations. Multiple cell-based activity assays are used to test each combination of nanoparticle type and cell type, resulting in each nanoparticle being tested in all conditions. Hierarchical clustering algorithms are used to group the nanoparticles based on their activity profiles. Class predictions can be made and verified. Understanding of structure-activity relationships increases, and in vivo correlations among nanoparticles can be tested, and compared with in vitro correlations.

### Scenario 13 B

This is Scenario 13 extended. How can an investigator use the dataset described above (and others created in similar ways) to make choices about nanoparticle design to optimize the chance that it would have a favorable in vivo activity? A scientist wants to maximize the circulating half-life of a nanomaterial. One material that has a long half-life is known and the scientist wonders if other nanomaterial compositions have similarly long half-lives.

The scientist would like to look at all available datasets, to see which nanomaterials act similarly to the known agent with a 11 of 19 long half-life. The scientist first queries across cancer center datasets to identify other nanoparticles with the best half-life. Initially, those data sets that use the same experimental protocol and a similar or better half-life are retrieved and compared. Next, the scientist wishes to broaden the search to include data sets that do not explicitly measure half-life, but a common set of cell-based assays. The data sets are normalized and combined. Hierarchical clustering algorithms are used to group the nanoparticles based on their activity profiles across the various cell-based assays. The scientist queries for nanoparticles that cluster closest to the starting nanoparticle with a long half-life, based on their behavior in the cell-based assays. The scientist then tests the hypothesis that the cluster neighbors will also have long half-lives in vivo.

### Scenario 14

This is a scenario based on identifying in vivo imaging probes using in vitro cell binding data. The scientist in the previous scenario would like to increase the imaging potential of candidate nanoparticles by modifying them and looking for cell type-specific binding capabilities.

The scientist submits a protocol to the institutional review board (IRB) and begins work upon approval. Libraries of surface-modified nanoparticles with appropriate pharmacokinetic and toxicity profiles are selected and screened for cell binding in vitro using cell cultures of "background" and "target" cell types or classes. The apparent concentration of binding or uptake of each nanoparticle to the different cell classes is measured. Metrics for differential binding to target versus background cells are calculated, and statistical significance is calculated by permutation. (These calculations employ analysis modules available through GenePattern.

To validate the increased specificity for binding target cells, those that provide the best discrimination are further tested ex vivo. Under IRB approval, anatomically intact human tissue specimens containing target and background cells are collected. The tissues are incubated with nanoparticles and evaluated for nanoparticle localization using microscopy. Further validation is conducted in vivo using an animal model. Animals are injected with the nanoparticle and another tissue specific probe and intravital microscopy is used to determine the extent of co-localization. The scientist contacts the tech transfer office to pursue next steps.

### Scenario 14 B

This is Scenario 14 extended, Customizing cell lines to identify nanoparticle probes. Varying the cell lines chosen for the study can help to generate analogous datasets. A scientist wants to find a nanoparticle that targets cancer cells bearing a specific oncogene mutation. Cell assays are performed in multiple cell lines that either do (target) or do not (background) bear this oncogene mutation. The data are analyzed as above to find particles that discriminate between the presence and absence of the mutation. The scientist then tries to validate these probes using independent tumor samples, or in mice genetically engineered to bear tumors that either do or do not express the mutation under study.

## Scenario 14 C

This is Scenario 14 extended, Analyzing existing datasets to identify nanoparticle probes. When many nanoparticles have been screened for their uptake in many different cell lines across many cancer centers, a scientist imports all the datasets that involve nanoparticle binding or uptake to cells. The cell lines are reclassified into target or background cells based on a set of criteria (such as tissue type or presence or absence of an oncogene mutation) and an analogous analysis is performed to identify nanoparticles that exhibit differential binding and uptake to different classes of cell lines.

### **NanoParticle Ontology**

This is a scenario based on evaluating and enriching the NanoParticle Ontology (NPO), an ontology which is being developed at Washington University in St. Louis to serve as a reference source of controlled vocabularies and terminologies in cancer nanotechnology research. Concepts in the NPO have their instances in the data represented in a database or in literature. In a database, these instances include field names, field entries, or both for the data model. The NPO represents the knowledge supporting unambiguous annotation and semantic interpretation of data in a database or in the literature. To expedite the development of the NPO, object models must be developed to capture the concepts and inter-concept relationships from the literature. Minimum information standards should provide guidelines for developing these object models, so the minimum information is also captured for representation in the NPO.

Nanotechnology is being applied to clinical therapeutics, but this use case could be extended to development of any specialized therapeutics. There are various pre-existing databases holding experimental data that need to be accessible across the entire community to facilitate rational nanomaterial design. Two strategies are being employed. The first is to establish semantic interoperability by finding areas of semantic overlap in the current database models based on controlled vocabularies (NCI Thesaurus, NCI Metathesaurus, Nanoparticle Ontology). The second is to develop a data submission standard based on the extension of standardized models (Biomedical Research Integrated Domain Group (BRIDG), Life Sciences Domain Analysis Model (LS-DAM)) where extensions are supported by controlled vocabularies. New vocabulary is needed to support both of these strategies. New concepts are curated in the controlled vocabularies as appropriate and term definitions are reviewed by the community.