

# *An Ontology-Based Framework for Clinical Research Databases*

Megan Kong<sup>1</sup>, Carl Dahlke<sup>2</sup>, Diane Xiang<sup>1</sup>, David Karp<sup>4</sup>, Richard H. Scheuermann<sup>1,3</sup>

<sup>1</sup>Department of Pathology, <sup>3</sup>Division of Biomedical Informatics, <sup>4</sup>Division of Rheumatology, U.T. Southwestern Medical Center, Dallas, TX, <sup>2</sup> Health Information Systems, Northrop Grumman, Inc., Rockville, MD



- NIH/NIAID/DAIT would like to:
  - maximize the return on the public investment in basic, translational and clinical research
  - allow investigators to more effectively extract meaningful information from the vast amounts of data generated from advanced research technologies
- Bioinformatics Integration Support Contract (BISC) to support all DAIT-funded programs - basic, translational and clinical research
- Immunology Database and Analysis Portal (*ImmPort*) - [www.ImmPort.org](http://www.ImmPort.org)
  - Archive and manage basic and clinical research data
  - Analyze experiment data
  - Integrate with extensive biological knowledge
  - Analyze integrated data
  - Build upon data generated by previous studies

## *ImmPort Overview - System Components*

- **Semi-public web-based database and analysis portal**
  - Multi-level access control
  - Data sharing
- **Data**
  - **Reference data**
    - **Types** - Gene structure, protein function, polymorphisms, metabolic, regulatory, signaling and other networks, protein-protein, gene-gene, host-pathogen interactions
    - **Sources** - NCBI, Uniprot, Swissprot, BIND, Reactome
  - **Experiment data**
    - **Methodologies** – microarrays, SNP genotyping, flow cytometry, ELISA, ELISPOT, qPCR, imaging, etc.
    - **Metadata** (defining how the experiment was performed) - common features of all experiments
    - **Primary results** - from all experiment measurement techniques
    - **Processed results**
      - Interpreted results
      - Analytical metadata
  - **Clinical research data**
    - **Study design** – from clinical protocol
    - **Participant characteristics/phenotype** – from assessments captured in CRF's and laboratory testing
- **Query tools**
  - To support retrieval of reference and experiment data based on specified criteria
  - Pre-defined QBE
  - Customized semantic queries
- **Ontology**
  - Thesaurus function
  - Organize terms and define relationships

- **Analysis tools**
  - **Genetic analysis**
    - LD analysis
    - TagSNP selection
    - Haplotype reconstruction
    - Genotype-phenotype association
  - **Gene expression analysis**
    - Filtering/normalization
    - Clustering
    - Classification
  - **Cell population analysis**
    - Standard FACS statistics
    - Novel population identification based on high dimensional data clustering
  - **Measurement of immune response (e.g., ELISA, ELISPOT)**
    - Statistical analysis of distributions
  - **Biological network analysis**
    - Quantification of topological parameters
    - Module identification
- **Visualization tools**
  - Genome display, including genes, introns, exons, SNPs, tagSNPs, etc.
  - Genetic analysis results
  - Gene expression results
  - Networks, pathways, and molecular interactions
  - Graphing and charting of statistical results

- Population Genetics Analysis Program: Immunity to Vaccines/Infections (PopGen)
- Immune Function and Biodefense in Children, Elderly, and Immunocompromised Populations (Special Pops)
- HLA Region Genetics in Immune-mediated Diseases
- Immune Tolerance Network (ITN)
- Atopic Dermatitis and Vaccinia Network (ADVNI)
- NIAID non-human primate colony

## **Varicella:**

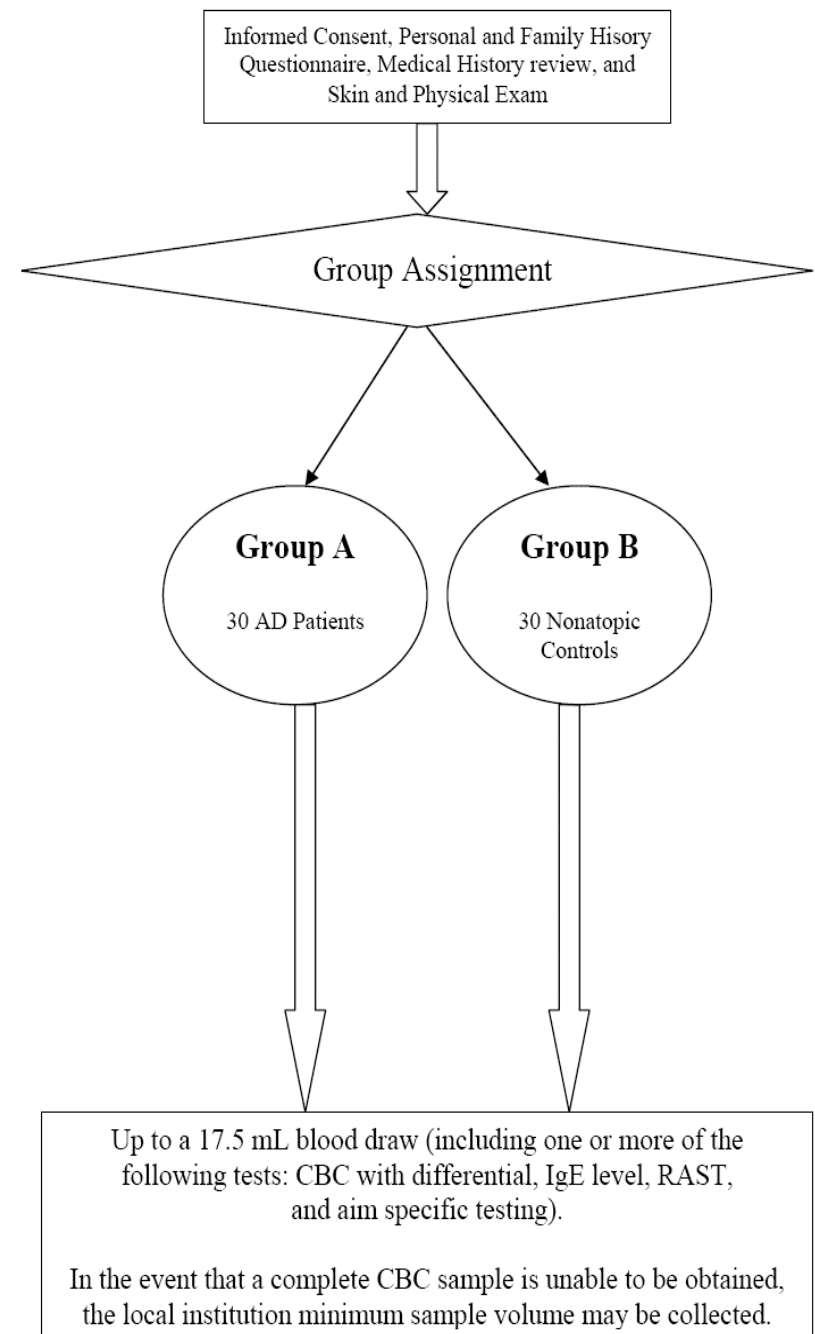
Immune Response to Varicella Vaccination in Subjects with Atopic Dermatitis Compared to Nonatopic Controls

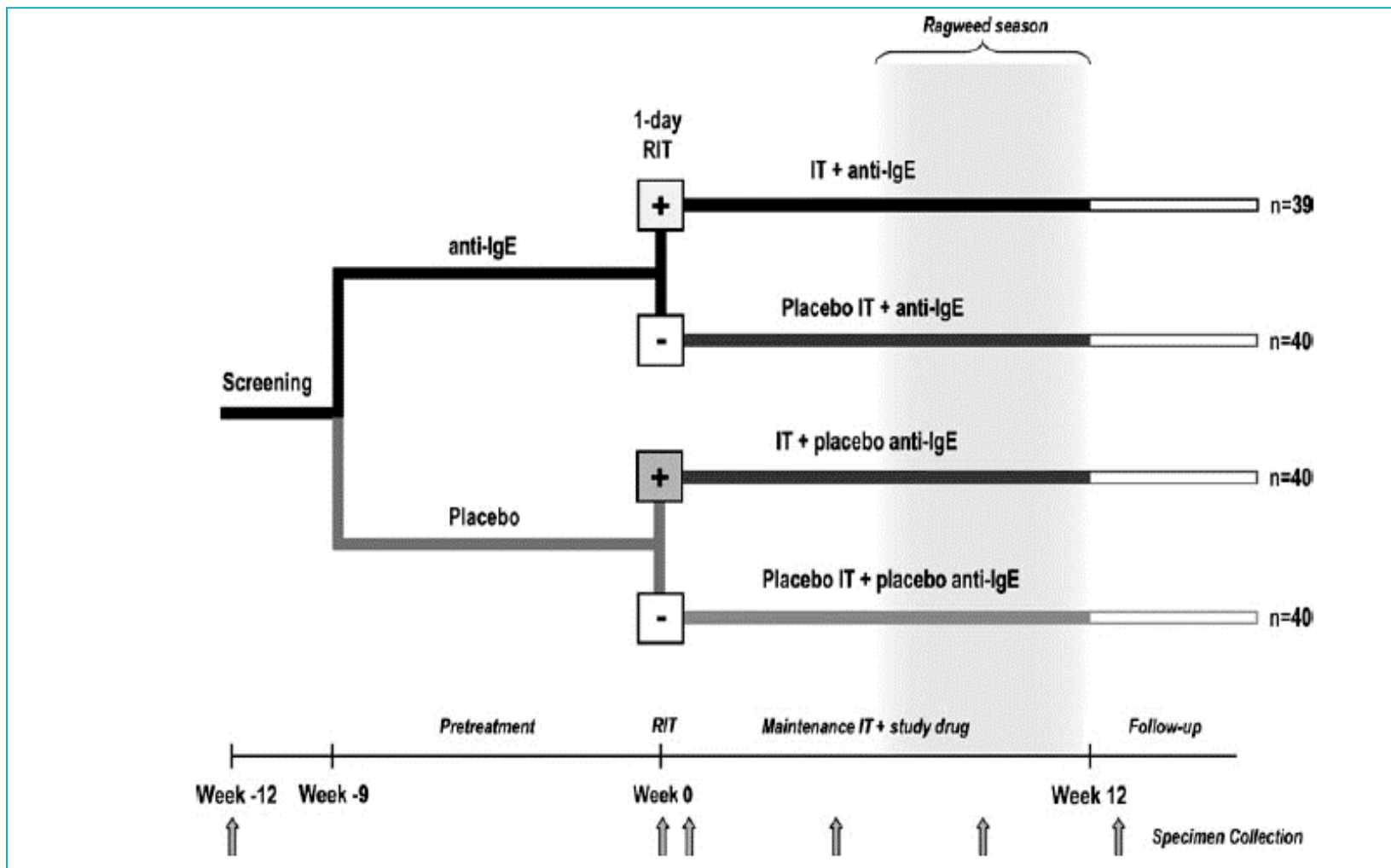
## **Objectives:**

To determine if children with AD have Varicella-specific cell mediated immune (CMI) responses to varicella vaccination that differ from those of nonatopic controls (Th2)

## **Endpoint:**

1. ELISPOT: number of Varicella-specific T-cell
2. ELISA: levels of Varicella-specific antibodies
3. Flow cytometry: perforin generation using 2-color flow cytometry staining for both CD8 and perforin.  
(Perforin levels, which may contribute to defects in CD8<sup>+</sup> cytotoxic T lymphocyte function).







## *Why OBX?*

- Challenges
  - Extensive phenotypic characteristics of interest
  - Complex study designs of treatments, assessments, samplings, etc.
  - Few standards for how they are described or captured
  - Each supported project captures their data in a different format
  - Integration with mechanistic studies of specimens



## *Why OBX?*

- Challenges
  - Extensive phenotypic characteristics of interest
  - Complex study designs of treatments, assessments, samplings, etc.
  - Few standards for how they are described or captured
  - Each supported project captures their data in a different format
  - Integration with mechanistic studies of specimens
- Evaluation of existing clinical data standards
  - BRIDG – designed for supporting the conduct of clinical trials, not a study repository optimized for analysis
  - CDISC STDM – specification for SAS transport files to the FDA
  - CDISC CDASH – models clinical encounters, but not sequence of events; requires use of EPOCHs (contiguous blocks of time)
  - None designed to manage mechanistic studies

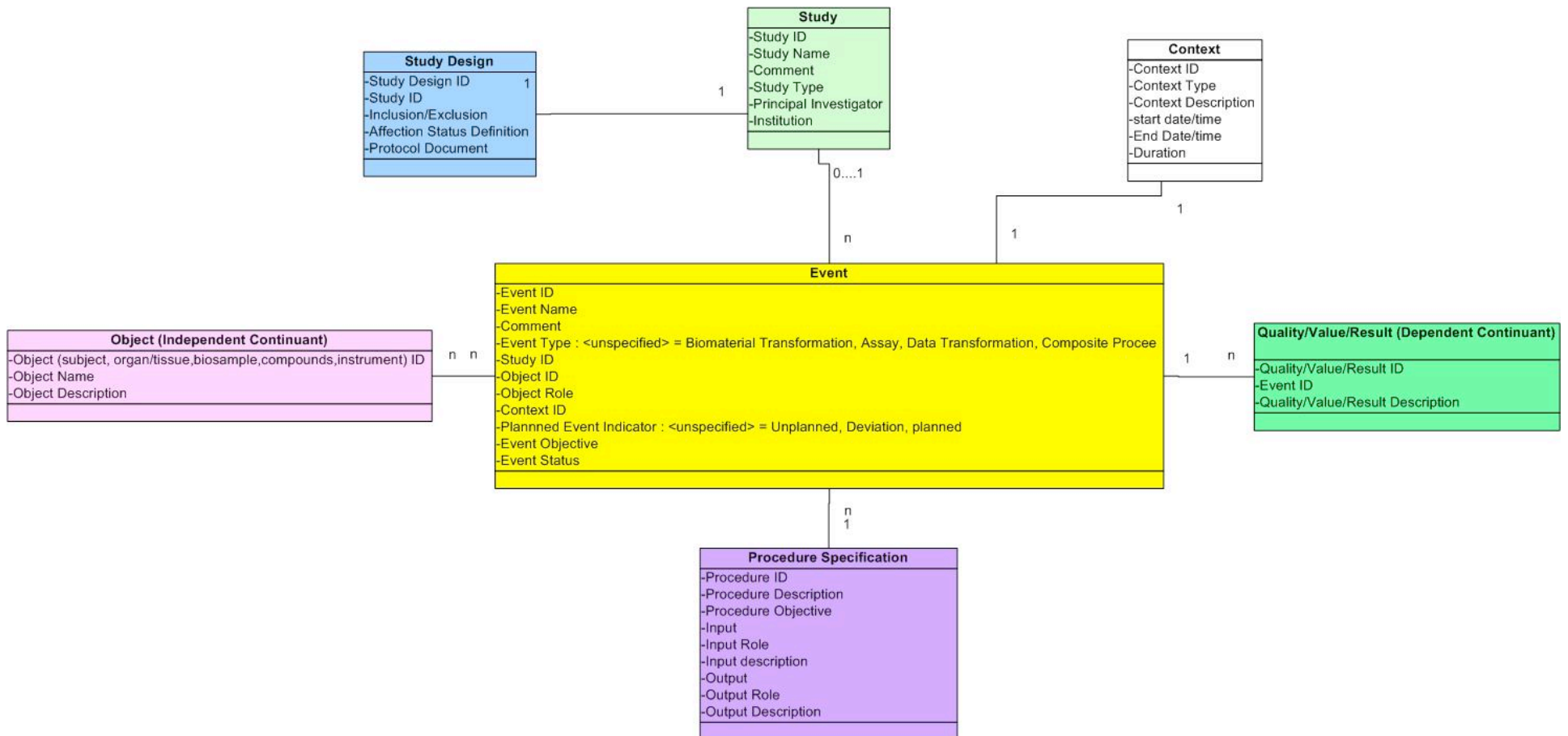
## *Why OBX?*

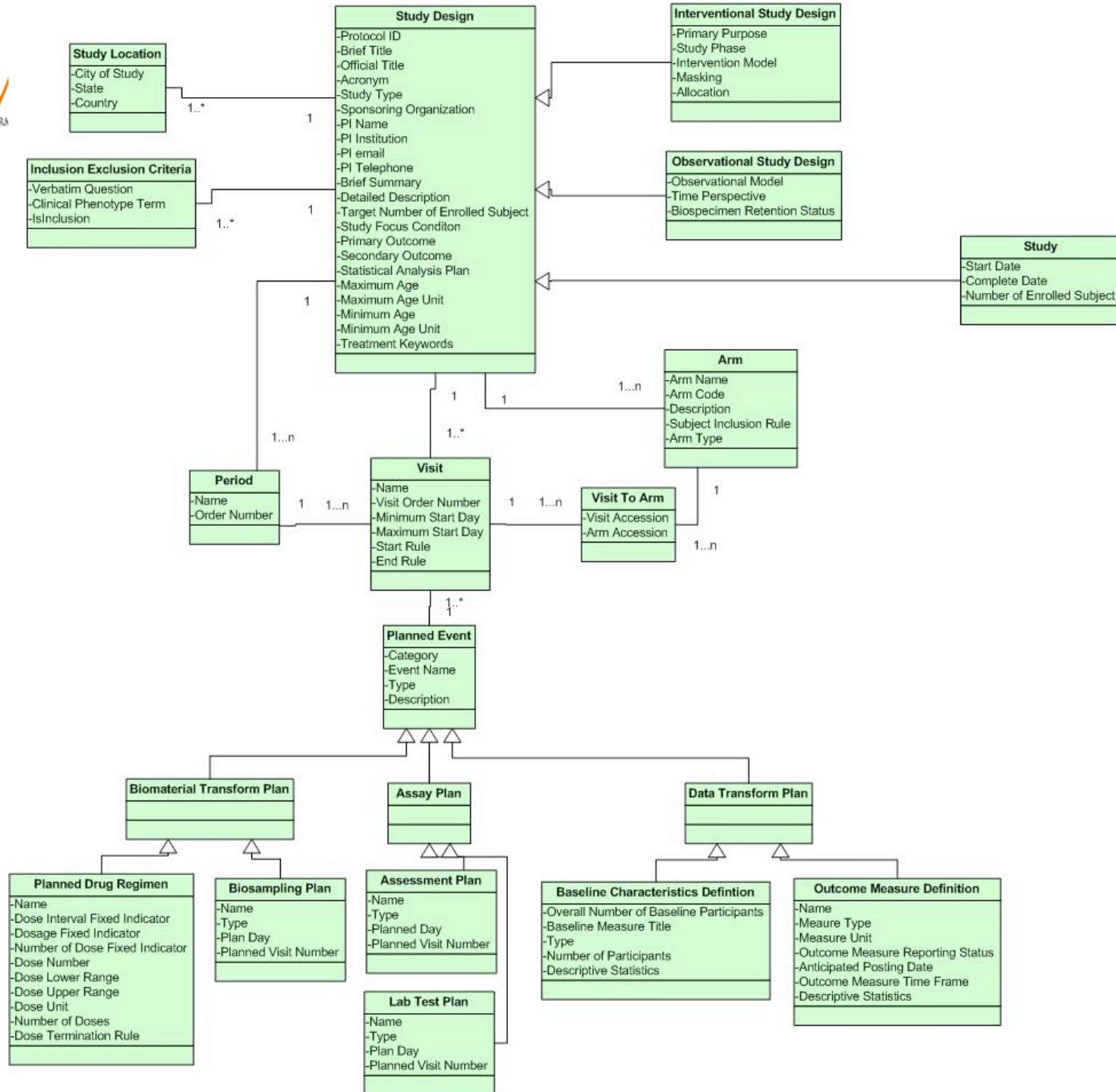
- Challenges
  - Extensive phenotypic characteristics of interest
  - Complex study designs of treatments, assessments, samplings, etc.
  - Few standards for how they are described or captured
  - Each supported project captures their data in a different format
  - Integration with mechanistic studies of specimens
- Evaluation of existing clinical data standards
  - BRIDG – designed for supporting the conduct of clinical trials, not a study repository optimized for analysis
  - CDISC STDM – specification for SAS transport files to the FDA
  - CDISC CDASH – models clinical encounters, but not sequence of events; requires use of EPOCHs (contiguous blocks of time)
  - None designed to manage mechanistic studies
- Ontology-Based eXtensible Data Model (OBX)
  - Logical ontology structure (BFO/OBI) => extensible data model
  - Data element value sets derived from ontology hierarchy

## *Entities of Interest*

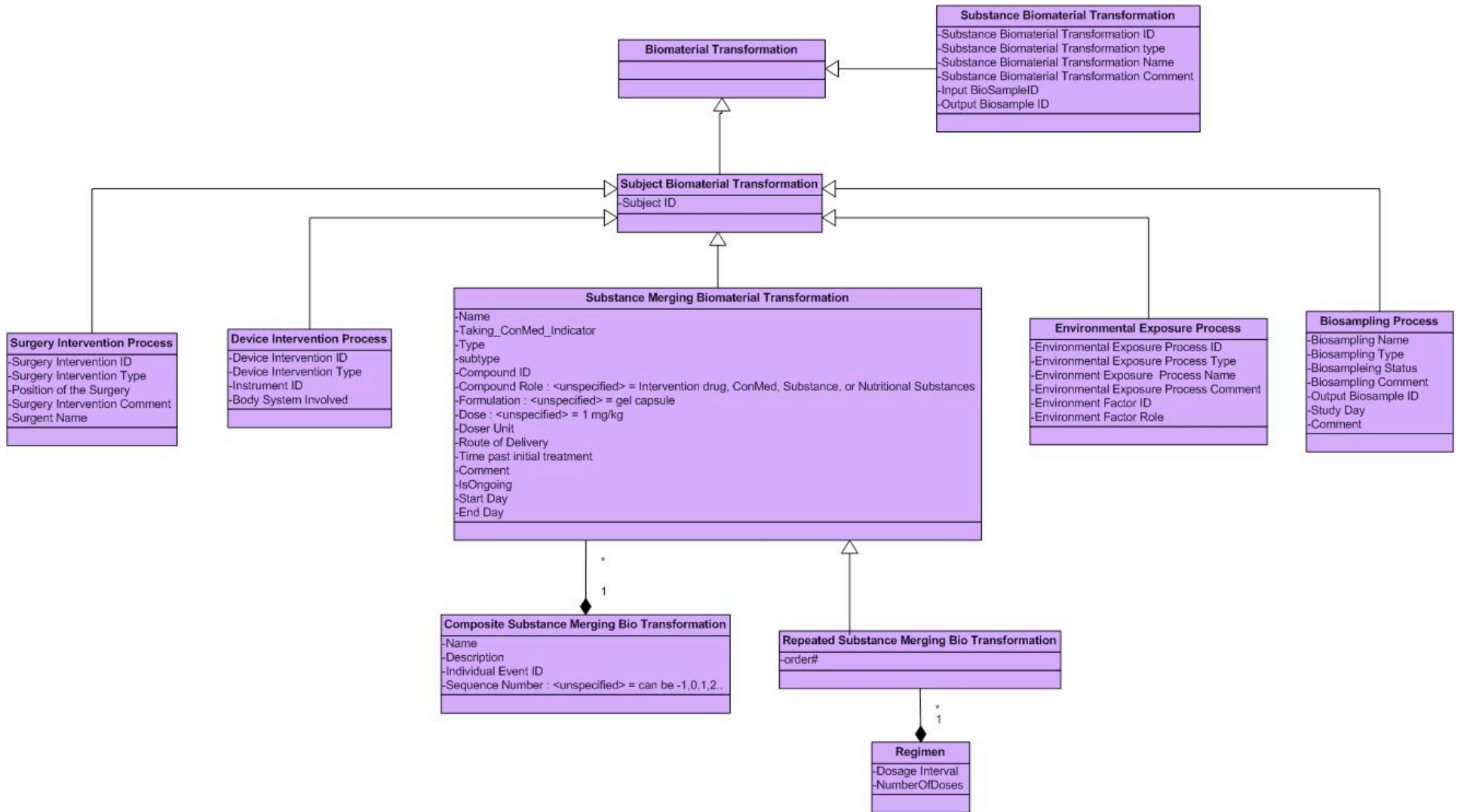
- **Objects**
  - Human subjects
  - Specimens
  - Materials
  - Plans
- **Qualities – clinical phenotype**
  - Malar rash
  - WBC count
  - Weight
  - Blood glucose levels
  - Atopic dermatitis
- **Roles**
  - Principal investigator, study coordinator, study participant
  - Reagent, therapeutic
- **Processes**
  - Recruitment, enrollment, approval
  - Assays, physical exams
  - Blood draw, PBMC isolation
  - Treatment, surgery (interventions)
- **Context**
  - Time
  - Environment

# OBX Framework

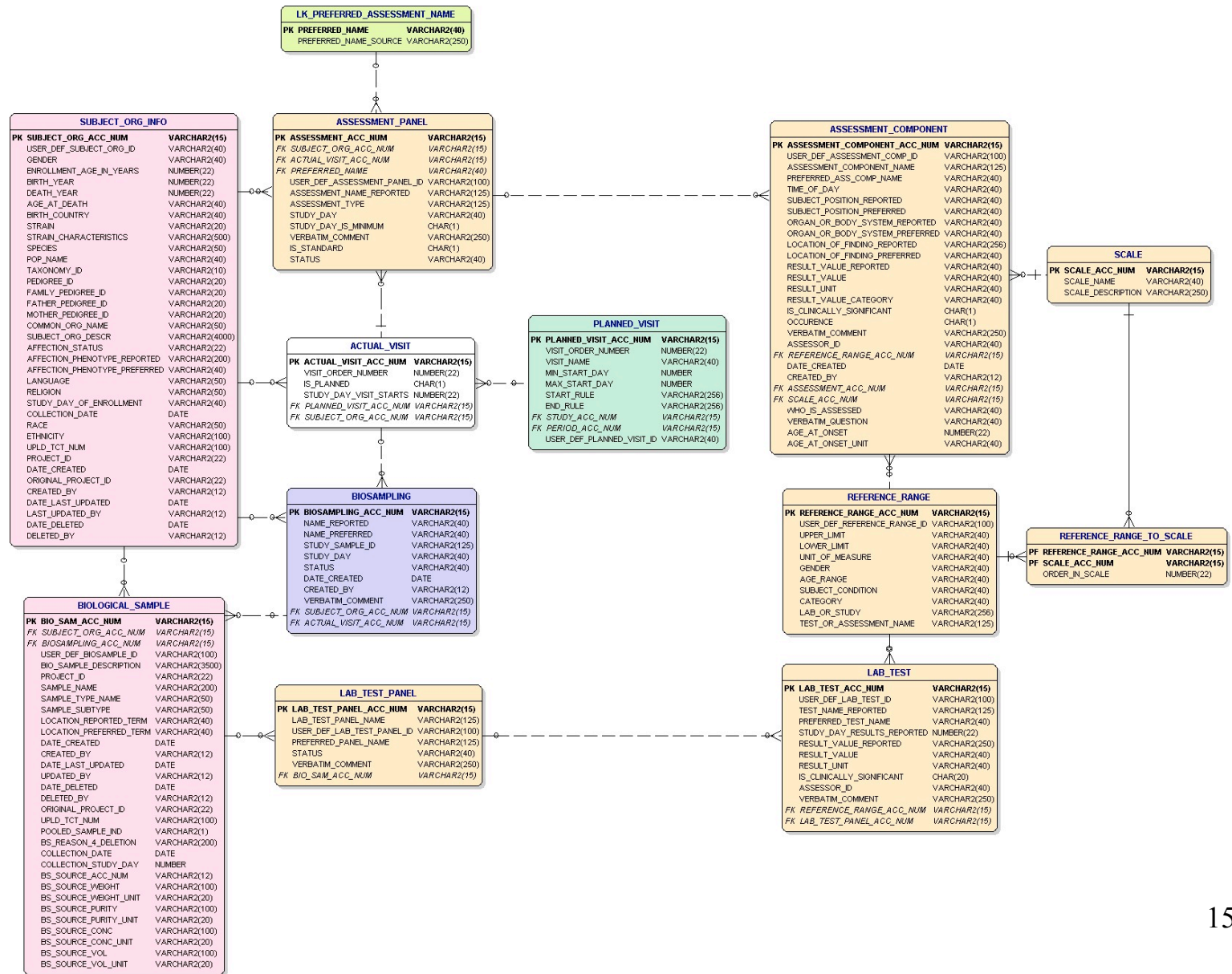




# OBX Biomaterial Transformation



# Assay Database Schema



## *Lab test mapping into OBX*

### Blood Test:

ADVN Fields	ADVN Fields Description	ImmPort Table	ImmPort Attribute
ID	Subject ID	Lab Test	Subject ID
SITE	Study Site	will not capture	
VISITDT	VISIT Date	Lab Test	Visit Date
HGB	Hemoglobin (g/dl):	Lab Test	Type = Blood Test Subtype = Hemoglobin Measurement Test Result Value = HGB Test Result Unit = grams per deciliter Clinical Significance = CSHGB
CSHGB	Hemoglobin Clinical Significant	Lab Test	
HCT	Hematocrit (%):	Lab Test	Type = Blood Test Subtype = Hematocrit Measurement Test Result Value = HCT Test Result Unit = percentage Clinical Significance = CSHCT
CSHCT	Hematocrit Clinical Significant	Lab Test	
WBC	Total WBC (x1000/ $\mu$ l):	Lab Test	Type = Blood Test Subtype = Total White Blood Cell Measurement Test Result Value = HGB Test Result Unit = 1000/ $\mu$ l Clinical Significance = CSHGB
CSWBC	Total WBC Clinical Significant	Lab Test	
NEUT	Neutrophils (x1000/microliter):	Lab Test	Type = Blood Test Subtype = Neutrophils Measurement Test Result Value = NEUT Test Result Unit = 1000/ $\mu$ l Clinical Significance = CSNEUT
CSNEUT	Neutrophils Clinical Significant	Lab Test	
LYMP	Lymphocytes (x1000/ $\mu$ l):	Lab Test	Type = Blood Test Subtype = Lymphocytes Measurement Test Result Value = LYMP Test Result Unit = 1000/ $\mu$ l Clinical Significance = CSLYMP
CSLYMP	Lymphocytes Clinical Significant	Lab Test	
EOSIN	Eosinophils (x1000/ $\mu$ l):	Lab Test	Type = Blood Test Subtype = Eosinophils Measurement Test Result Value = EOSIN Test Result Unit = 1000/ $\mu$ l Clinical Significance = CSEOS
CSEOS	Eosinophils Clinical Significant	Lab Test	



## Load File Examples

Lab_Test_Panel_Acc	Test_Name_Repo	Result_Value	Result_Unit
0101780 blod	Hemoglobin	13.3	g/dl
0315450 blod	Hemoglobin	12.8	g/dl
0315483 blod	Hemoglobin	10.8	g/dl
0315931 blod	Hemoglobin	13.6	g/dl
0316415 blod	Hemoglobin	12.3	g/dl
0101780 blod	Hematocrit	41	%
0315450 blod	Hematocrit	35	%
0315483 blod	Hematocrit	32	%
0315931 blod	Hematocrit	40	%
0316415 blod	Hematocrit	34	%
0101780 blod	Total WBC	12.7	1000/microliter
0315450 blod	Total WBC	5.22	1000/microliter
0315483 blod	Total WBC	4.53	1000/microliter
0315931 blod	Total WBC	8.8	1000/microliter
0316415 blod	Total WBC	11.25	1000/microliter
0101780 blod	Total Neutrophils	3.3	1000/microliter
0315450 blod	Total Neutrophils	0.31	1000/microliter
0315483 blod	Total Neutrophils	1.81	1000/microliter
0315931 blod	Total Neutrophils	1.14	1000/microliter
0316415 blod	Total Neutrophils	3.39	1000/microliter
0101780 blod	Lymphocytes	7.6	1000/microliter
0315450 blod	Lymphocytes	3.92	1000/microliter
0315483 blod	Lymphocytes	1.9	1000/microliter
0315931 blod	Lymphocytes	7.04	1000/microliter
0316415 blod	Lymphocytes	6.22	1000/microliter
0101780 blod	Eosinophils	0.9	1000/microliter
0315450 blod	Eosinophils	0.21	1000/microliter
0315483 blod	Eosinophils	0.18	1000/microliter
0315931 blod	Eosinophils	0	1000/microliter
0316415 blod	Eosinophils	0.98	1000/microliter

Assessment_Comp	Assessment_Co	Location_of_Find	Organ_or_Body	Result_Value
0101780 easi	Erythema	Head	Skin	Moderate
0104097 easi	Erythema	Head	Skin	Moderate
0104302 easi	Erythema	Head	Skin	Mild
0104468 easi	Erythema	Head	Skin	Mild
0106229 easi	Erythema	Head	Skin	Mild
0101780 easi	Erythema	Trun	Skin	Absent
0104097 easi	Erythema	Trun	Skin	Moderate
0104302 easi	Erythema	Trun	Skin	Absent
0104468 easi	Erythema	Trun	Skin	Absent
0106229 easi	Erythema	Trun	Skin	Moderate
0101780 easi	Erythema	Arms	Skin	Mild
0104097 easi	Erythema	Arms	Skin	Absent
0104302 easi	Erythema	Arms	Skin	Mild
0104468 easi	Erythema	Arms	Skin	Mild
0106229 easi	Erythema	Arms	Skin	Mild
0101780 easi	Indurati	Head	Skin	Moderate
0104097 easi	Indurati	Head	Skin	Mild
0104302 easi	Indurati	Head	Skin	Mild
0104468 easi	Indurati	Head	Skin	Mild
0106229 easi	Indurati	Head	Skin	Mild
0101780 easi	Lichenif	Head	Skin	Absent
0104097 easi	Lichenif	Head	Skin	Mild
0104302 easi	Lichenif	Head	Skin	Absent
0104468 easi	Lichenif	Head	Skin	Absent
0106229 easi	Lichenif	Head	Skin	Absent

- OBX is an ontology-informed clinical research data model
- Ontological framework incorporated provides for extensibility and straightforward incorporation of ontology terms as value sets
- Ontological framework also facilitated distributed data model development and guides ongoing mapping strategies
- Successfully loaded two large, distinct clinical studies
- Building the user interfaces for data extraction and analysis
- Refinement through use – mapping, loading, extracting, analysis
- Term submission to OBO Foundry ontologies, especially OBI

## *Acknowledgments*

### *UT Southwestern*

Yu (Max) Qian

Jamie Lee

**Megan Kong**

Jennifer Cai

Jie Huang

Nishanth Marthandan

**Diane Xiang**

Young Bun Kim

Paula Guidry

Eva Sadat

**David Karp**

### *Northrop Grumman*

**Carl Dahlke**

John Campbell

Liz Thompson

**Jeff Wiser**

Mike Attasi

### *OBO Foundry*

**BFO**

**OBI**



Supported by NIH N01AI40076