# SEPIO: An Ontology-Based Modeling Framework for Evidence and Provenance Information

Session S27: Open, Expert-curated, Harmonized, and Standardized Precision Oncology Knowledge, Ontologies and APIs

**Matthew Brush**

Oregon Health and Science University
Twitter: #AMIA2018

# The SEPIO Modeling Framework

## The Scientific Evidence and Provenance Information Ontology

- **SEPIO** is the foundation of an **ontology-based modeling framework** for computable representation of **scientific assertions** and the **evidence** and **provenance** supporting them.

- The core SEPIO model is **domain-agnostic** but **extensible** with domain-specific content, and capable of representing any type of assertion and evidence.

- The broader **modeling framework** enables creation of **custom schema** for specific applications called **SEPIO Profiles**.

- **Ontology-based data models** use terms from ontologies to structure data:
  - Mappings from data model **types** and **attributes** to ontology **classes** and **properties**
  - Ontology terms used to build **value sets** for data collection

# Benefits of Ontology-Based Models

1. **Promote understanding and communication** of the domain and the data by providing a unifying conceptual model and terminology

2. **Improve searchability** of data by leveraging relationships in ontology for semantic search and query expansion

3. **Allow algorithmic derivation of new knowledge** based on information encoded in the ontology

   a. Description logic reasoners use computable semantics of ontology to draw novel inferences and perform consistency checks

   b. Semantic similarity algorithms can operate on graphical relationships between ontology terms to infer similarity of entities annotated with them

4. **Facilitate data integration**, **discovery**, and **analysis** by exploiting modern semantic web standards and tools (linked data, RDF, JSON-LD)

# SEPIO Supports Diverse Assertions Across the Basic Science, Translational, and Clinical Spectrum

**Established**

**Monarch Initiative**: genotype-phenotype associations from human and model organism knowledgebases

**ClinGen:** clinical interpretations of genes and variants

**Developing**

**GA4GH**: variant annotations and phenotype associations

**CIViC/VICC**: somatic variant interpretations

**Exploratory**

**HL7 Clinical Genomics**: patient level genomic data

**GloBI/iDigBio:** biodiversity, insect life history assertions

# The SEPIO Worldview

We view the act of making an evidence-based assertion as involving **three universal tasks:**

1. The **identification** of **data** that might be used as evidence for a proposition

2. The **interpretation** of this data as independent **arguments** for or against this proposition

3. The **evaluation** of all arguments towards a final **conclusion** about the truth of the proposition
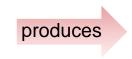
# The SEPIO Core Model

*A central axis of three core classes represent the outputs of each universal task.*

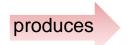## Universal Task

**(3)** **Evaluation** towards a **Conclusion**
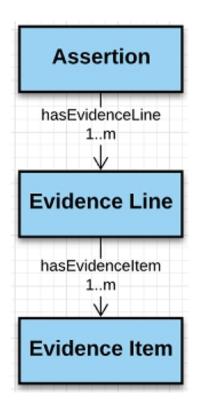
**(2)** **Interpretation** as **Arguments**

**(1)** **Identification** of **Data**

## SEPIO Central Axis



## Example

*Counsyl Genetics' 2015 claim that "BRCA2 c.8023A>G Is pathogenic for Breast Cancer"*

*The argument made by the raw counts and calculated frequency of BRCA2 c.8023A>G, as providing moderate supporting argument for its pathogenicity (ACMG criteria PM2)*

*The 0.0021% population frequency of the BRCA2 c.8023A>G variant in ExAC.*

*Alignment of the model with this universal paradigm is key to its broad applicability.*
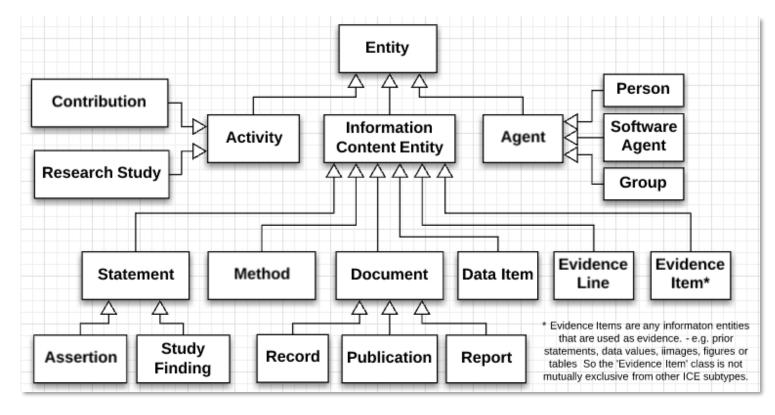
# The SEPIO Modeling Framework

The modeling framework is comprised of four components that enable creation of ontology-based schema for specific domains and applications

1.  **SEPIO Core Ontology:** defines the foundational, domain-agnostic model

2.  **SEPIO Information Model:** provides a view of the ontology specifying how the terms and design patterns it defines can be used to structure data

3.  **SEPIO Profiles:** a mechanism to create custom data models by extending the core model and ontology to provide schema for particular applications

4.  **SEPIO Value Sets:** a model to create ontology-based collections of terms to constrain specific aspects of data collection
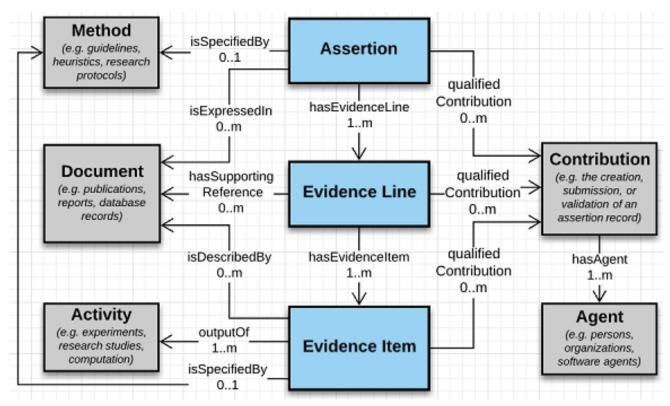
# 1. The SEPIO Core Ontology

- Defines a generic conceptual model on which specific data models are built

- OWL2 ontology, adheres to OBO Foundry Principles

- Reuses terms from community ontologies where possible

- Enables comprehensive mapping of data model **types** and **attributes** to ontology **classes** and **properties**



Hierarchical relationships between the high-level classes/types in the SEPIO core ontology and model

# 1. The SEPIO Core Ontology

- Defines a generic conceptual model on which specific data models are built

- OWL2 ontology, adheres to OBO Foundry Principles

- Reuses terms from community ontologies where possible

- Enables comprehensive mapping of data model **types** and **attributes** to ontology **classes** and **properties**
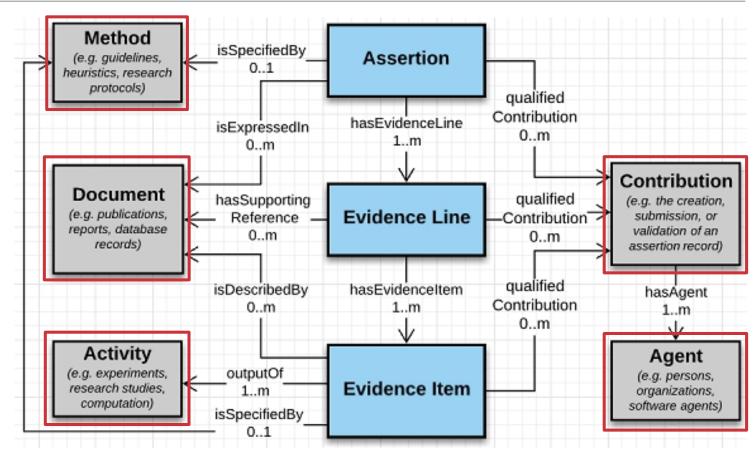


*Relationships between high-level concepts in the core model*

## Provenance

Surrounding the core axis are concepts that let us describe how, when and by whom these core artifacts were created:

- **Contributions** made by **Agents**
- **Activities** they perform to do so, and any **resources** used
- **Methods** that specify activities
- **Documents** that describe them.



*This core model lets us separately track the provenance behind how data is generated, interpreted as evidence, and collectively evaluated to make a final assertion.*
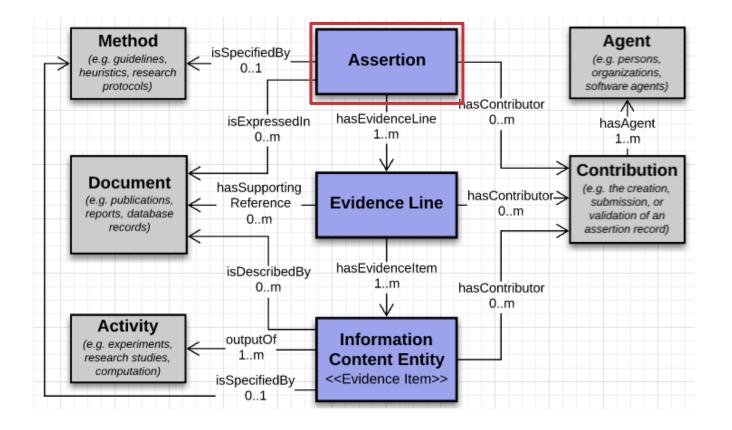
# 1. The SEPIO Core Ontology

## Assertion

**Definition**: an evidence-based statement of purported truth, as made by a particular agent on particular occasion.

**Example**: Counsyl Genetics' 2015 Assertion that the BRCA2 variant c.8023A>G is pathogenic for Breast Cancer.

**Comments**: Assertions put forth a particular 'Proposition' as true. More than one Assertion, made by different agents on different occasions, can put forth the same Proposition.
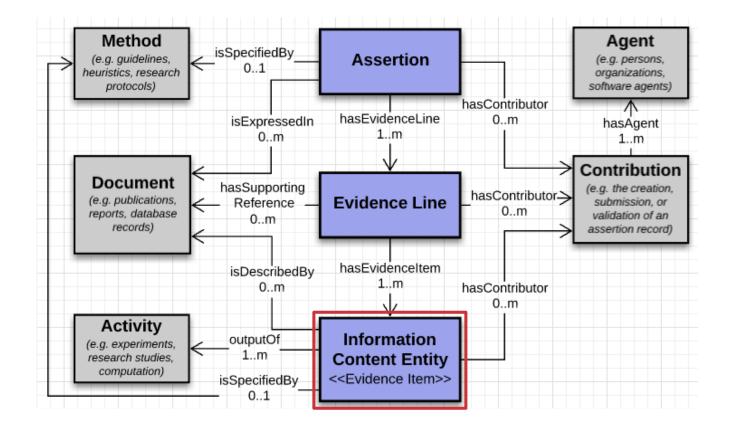
# 1. The SEPIO Core Ontology

## Evidence Item

**Definition**: an individual piece of information that is interpreted to build an argument for or against an Assertion

**Example**: the 0.0021% population frequency of the BRCA2 variant c.8023A>G in healthy NFE individuals in ExAC.

**Comments**: Evidence Items can be primary data, derived statistical calculations, tables/figures depicting these data, statements summarizing the results of a particular study, or prior assertions based on their own lines of evidence.

# 1. The SEPIO Core Ontology

## Evidence Line

**Definition**: independent, meaningful arguments relevant to the validity of an Assertion, that are supported by one or more Evidence Items

**Example**: the argument made for the BRCA2 variant's pathogenicity based on its observed absence in healthy populations.

**Comments**: Representing Evidence Items separately from the 'arguments' they make is an important feature that lets us describe properties of information emerging only through its interpretation as evidence.
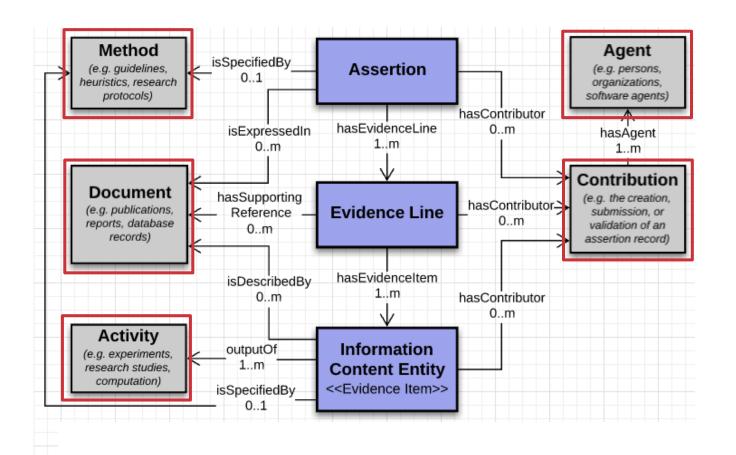
## Provenance

Surrounding the core axis are familiar concepts that let us describe in rich detail who, how, and when these core informational objects were created:

- **Contributions** made to them by **Agents**

- **Activities** they perform to do so

- **Methods** that specify these activities

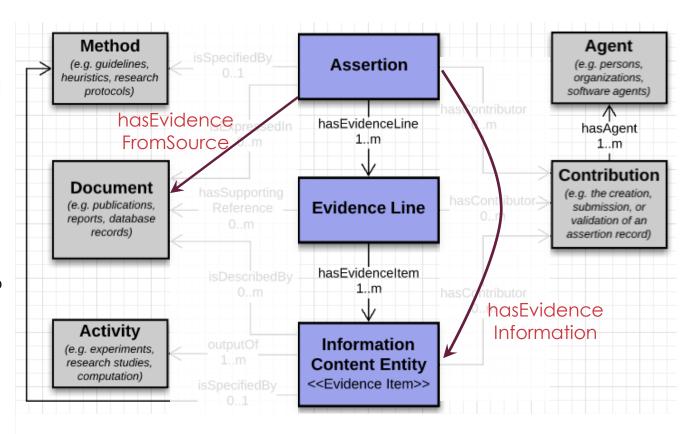- **Documents** created to describe them.

# 1. The SEPIO Core Ontology

## Shortcut Relations

SEPIO provides relationships that can directly link objects normally connected via multiple edges in a fully normalized model, to support concise representation of data with sparser evidence and provenance metadata.

**hasEvidenceFromSource**: connects an Assertion directly to a document that contains info interpreted as evidence to support it
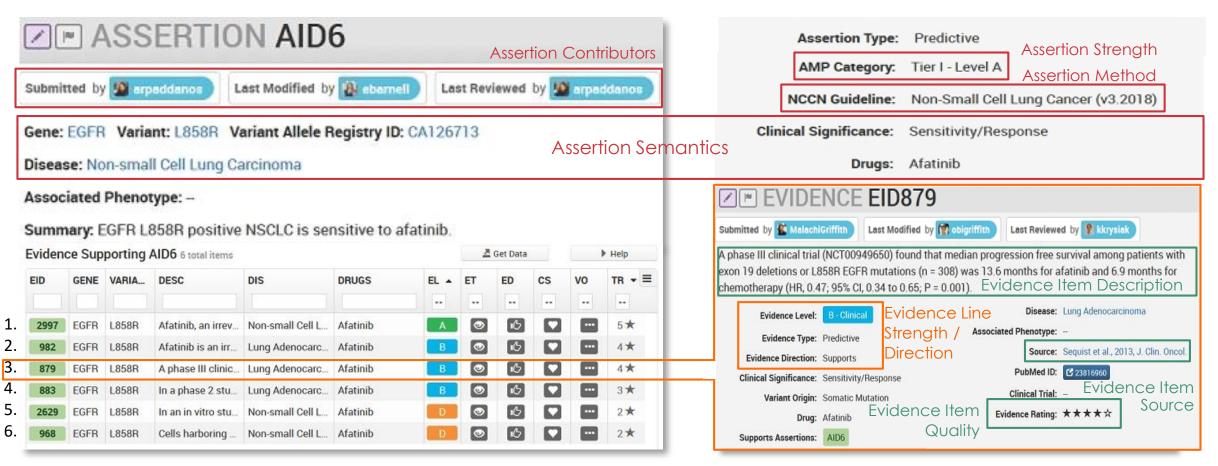
**hasEvideneInformation**: connects an Assertion directly to an Evidence Item

# A CIViC 'Predictive' Variant Interpretation

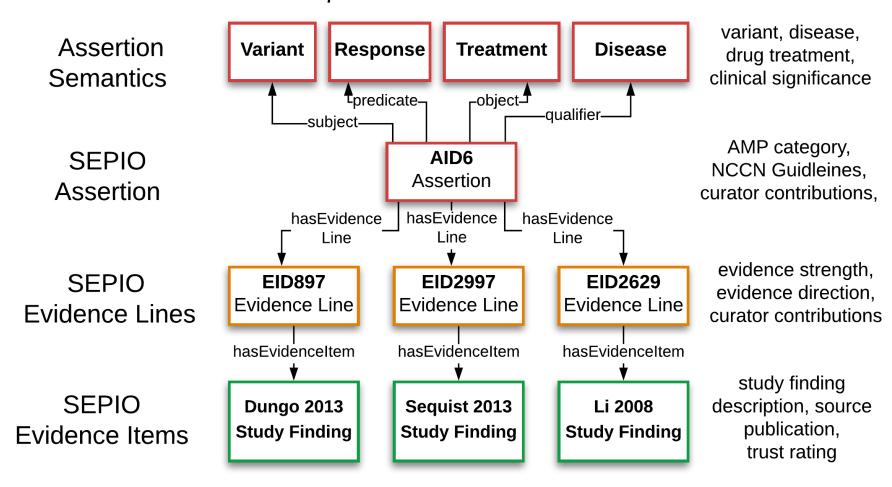**Assertion**: *"EGFR L858R positive NSCLC is sensitive to Afatinib"*



https://civicdb.org/events/assertions/6/summary

# 2. The SEPIO Information Model

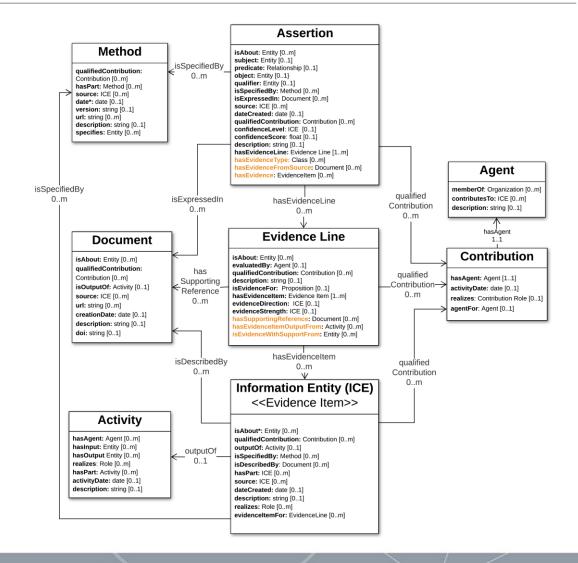- Provides a view of the ontology as a data model

- Specifies all the ways the terms/ design patterns in the ontology can be used to structure data

- Provides a starting point from which to derive a custom schema for a particular application (SEPIO Profiles)

# 3. SEPIO Profiles

Profiles extend/refine the core ontology and information model to generate schema customized for a particular application

Schema typically use formats like JSON-LD that have native support for ontology mapping.
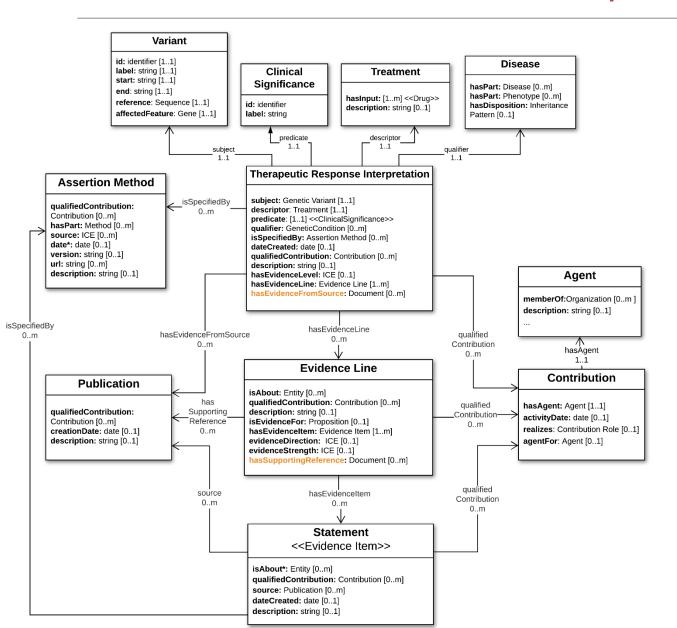
The profiling mechanism allows implementers to:
- Specialize core types for their domain (e.g. Assertion -> Variant Interpretation)
- Apply application-specific cardinalities, data type constraints, and value sets

*SEPIO can support simple or complex profiles, depending on the richness of target data*

# A CIViC 'Somatic Interpretation' Profile



**Therapeutic Response Interpretation**

**subject:** Genetic Variant [1..1]
**descriptor**: Treatment [1..1]
**predicate**: [1..1] <<ClinicalSignificance>>
**qualifier:** GeneticCondition [0..m]
**isSpecifiedBy:** Assertion Method [0..m]
**dateCreated:** date [0..1]
**qualifiedContribution:** Contribution [0..m]
**description:** string [0..1]
**hasEvidenceLevel:** ICE [0..1]
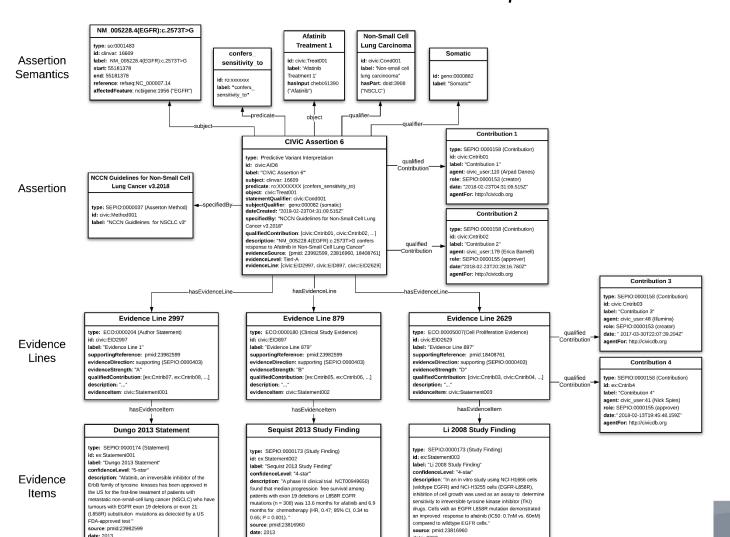**hasEvidenceLine:** Evidence Line [1..m]
**hasEvidenceFromSource:** Document [0..m]

- Assertion semantics
- Assertion strength
- Evidence Line direction and strength
- Evidence Item descriptions, source, quality
- curator contributions to assertions, evidence lines, and evidence items

# A Somatic Interpretation Data Example

**Assertion**: *"EGFR L858R positive NSCLC is sensitive to Afatinib"*

# 4. SEPIO Value Sets

- SEPIO provides **a meta-model** for creating **ontology-based value sets** for particular profiles

- Allow tools to leverage knowledge encoded in ontologies for **improved search and analysis**, and can **make connections to other datasets** using the same ontologies.

**Molecular Consequence Value Set**
- SO:0002012 start lost
- SO:0001578 stop lost
- SO:0001587 stop gained
- SO:0001819 synonymous variant
- SO:0001589 frameshift variant
- SO:0001823 conservative inframe insertion
- SO:0001824 disruptive inframe insertion
- SO:0001825 conservative inframe deletion
- SO:0001826 disruptive inframe deletion
- SO:0001909 frameshift elongation
- SO:0001568 splicing variant
- . . .

# Status and Future Plans for SEPIO
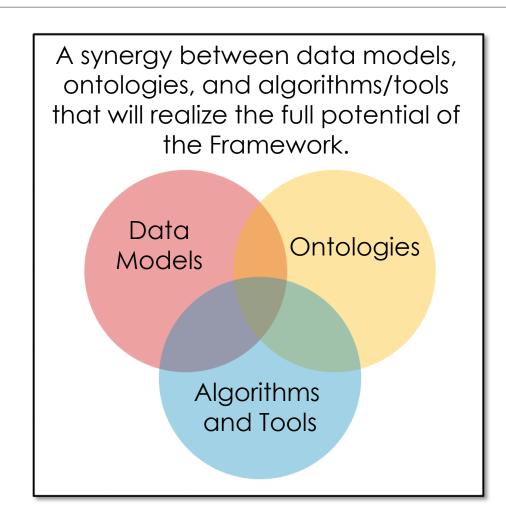
**Short Term:**

- Collect feedback from pilot adopters to improve the model and framework

- Prepare an initial ontology release (~3 mo.)

**Medium Term:**

- Create implementation guide and improved documentation for adopters (3 mo.)

- Develop reference implementation and code to support developers (6 mo.)

**Long Term:**

- Build out a broader infrastructure of supporting ontologies, algorithms, and tools

A synergy between data models, ontologies, and algorithms/tools that will realize the full potential of the Framework.

Data Models

Ontologies

Algorithms and Tools

*Interested adopters - contact us and present your use case. We'd love to work with you!*