

# Information Modeling

HL7 WGM

Sept 23, 2020

Robert R. Freimuth

CG Co-Chair

# Agenda

- Why develop an IM?
- Review model
  - Definitional Variant!
- Roadmap

# IM Group: Purpose

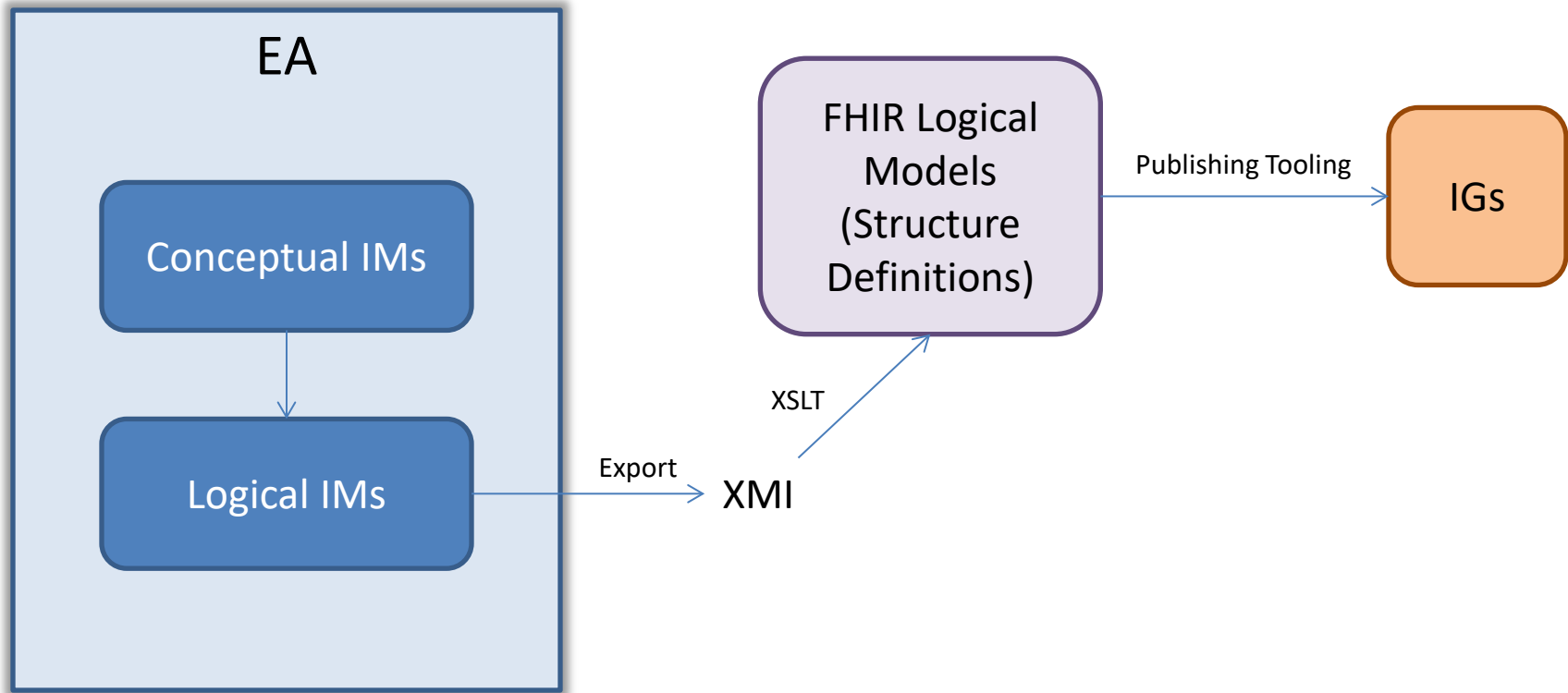
A subgroup of the CG WG will be formed to focus on the **development of an information model to represent the clinical genomics domain**. The creation of this model is consistent with current CG WG project scope statements and it is **a necessary component of the future consolidated standard**, which will include the current Domain Analysis Model and Domain Information Model, as described in the CG WG DAM PSS (Q1/2016).

- A model is needed to unify our specifications

# CG Information Model (IM)

- Technology agnostic
- Conceptual information model
  - Core concepts that are represented in draft standards produced by the WG
- Promote harmonization and inter-artifact consistency among all CG WG standards
- Incorporate relevant concepts requested by other groups working in the CG domain
  - GA4GH
  - ClinGen
  - NCBI
  - etc

# Big Picture

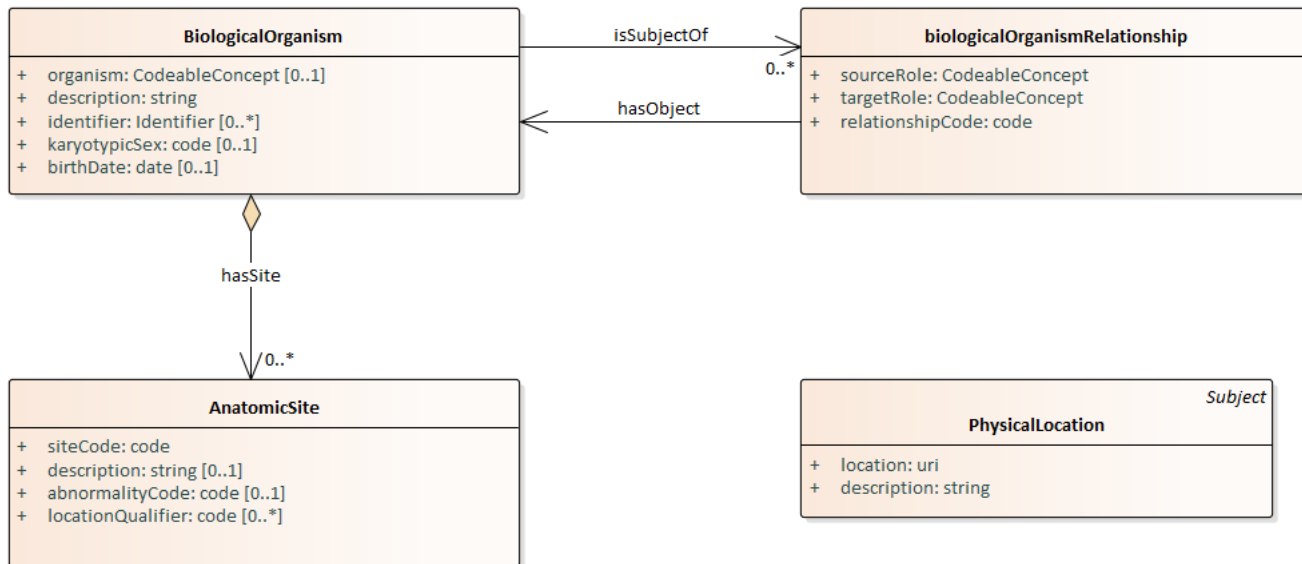


# Things to Keep in Mind

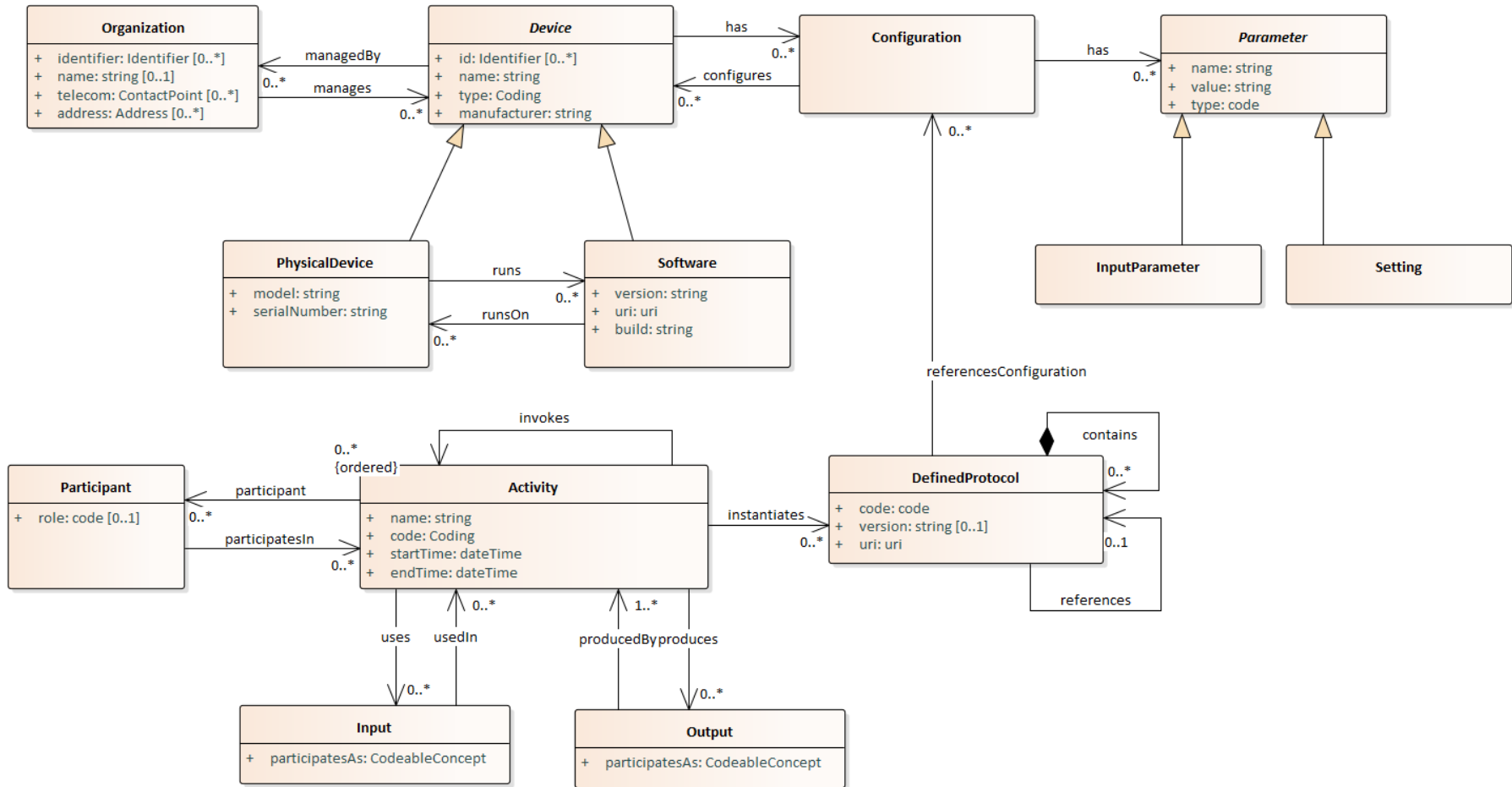
- Domain Analysis Model (DAM)
  - Captures semantics
  - Concepts, relationships
- Intended for SME, not (only) IT
- This is not an implementation

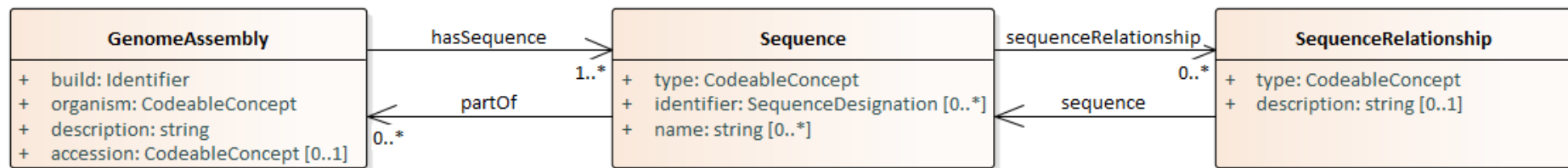
# Things to Keep in Mind

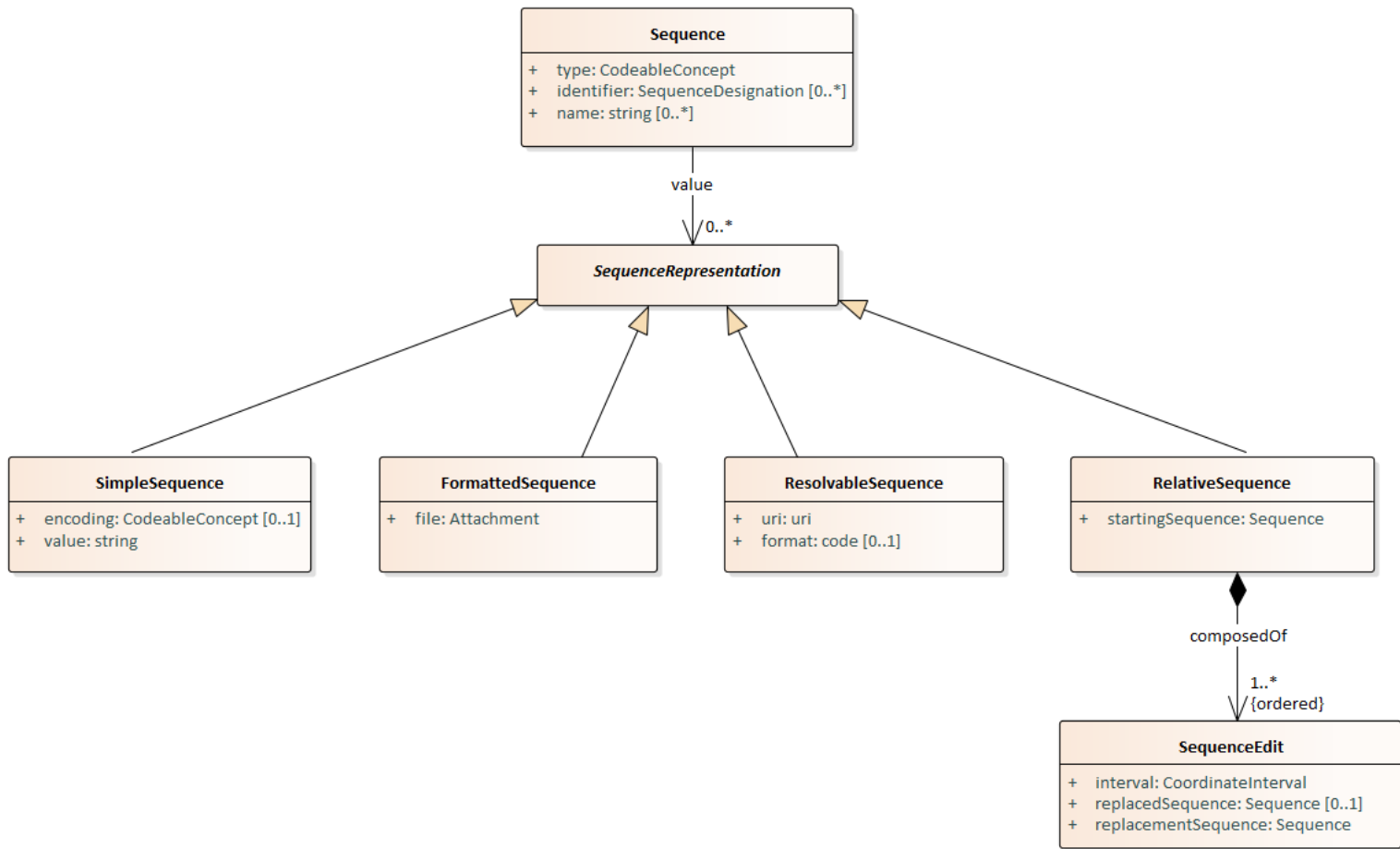
- Builds on existing work
  - Reuse when possible
  - Extend/remodel when necessary
  - Models, ontologies, implementations (FHIR, v2)
- Aligns to other work
  - GA4GH
  - SEPIO/Monarch
  - Specimen DAM
- This is not an implementation (FHIR)

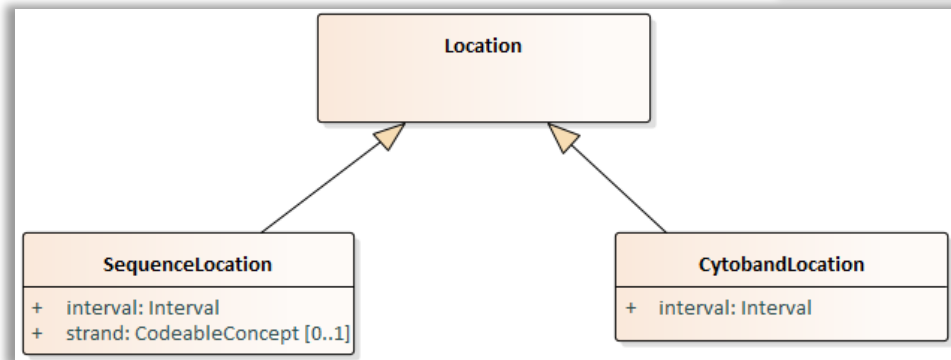
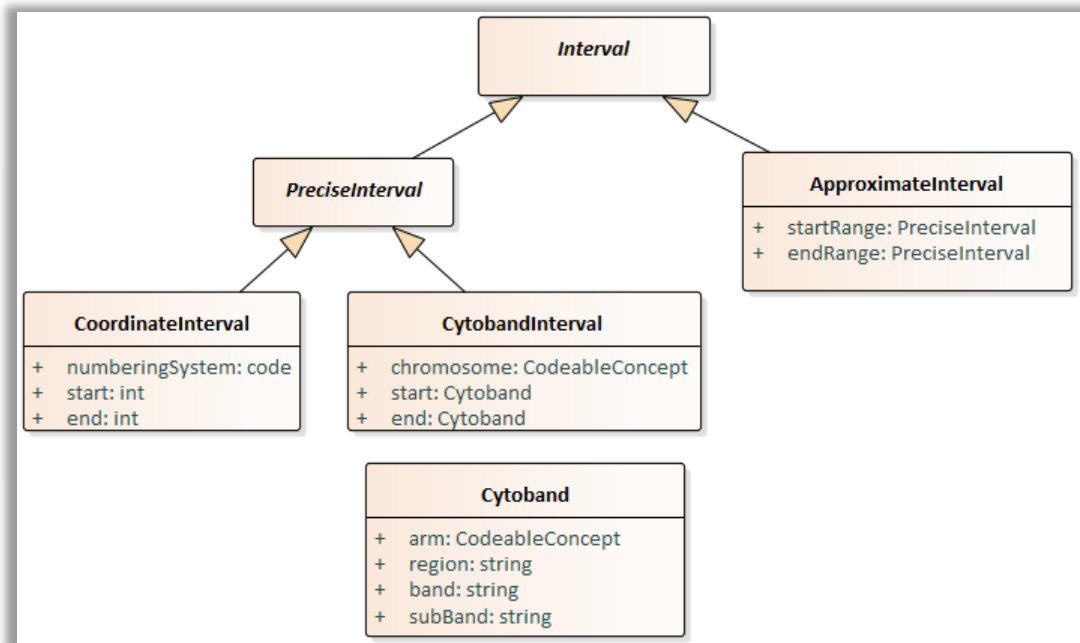
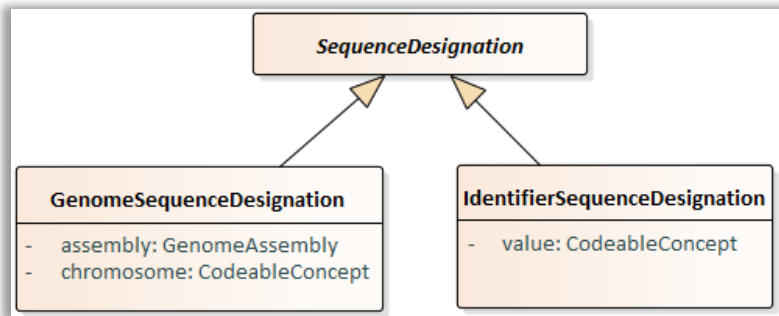


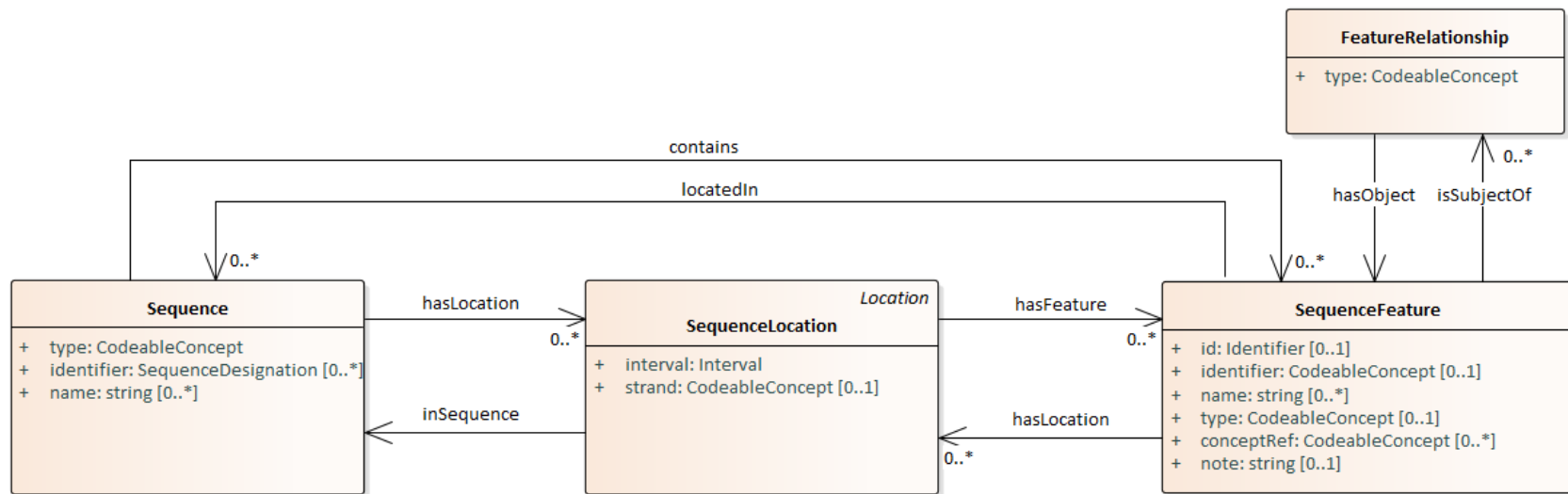


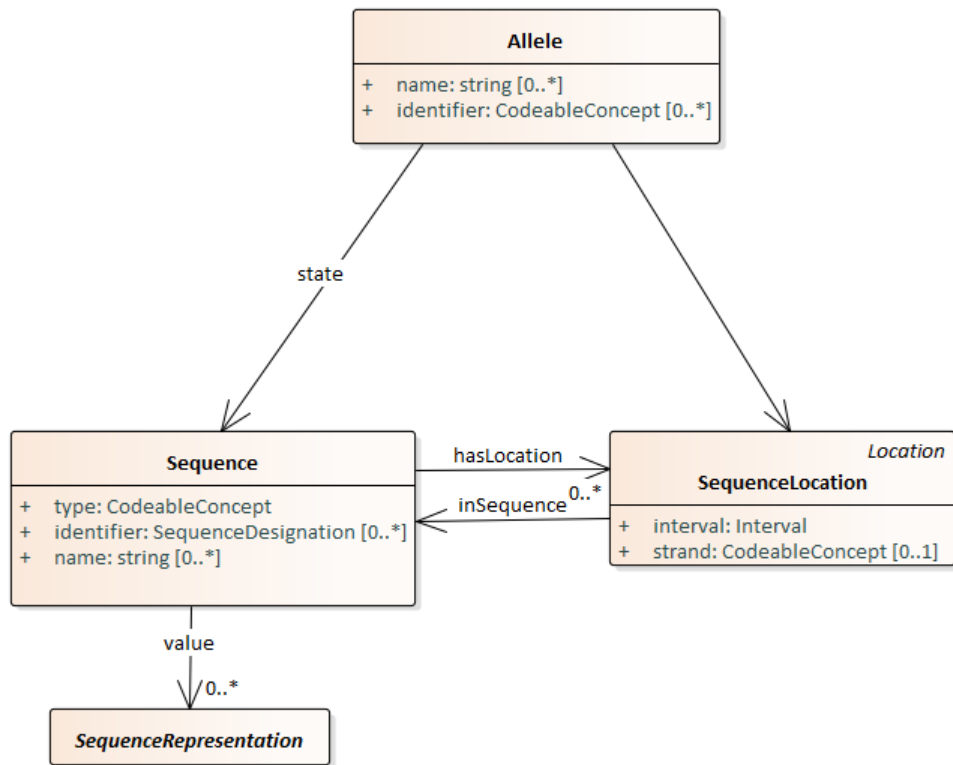






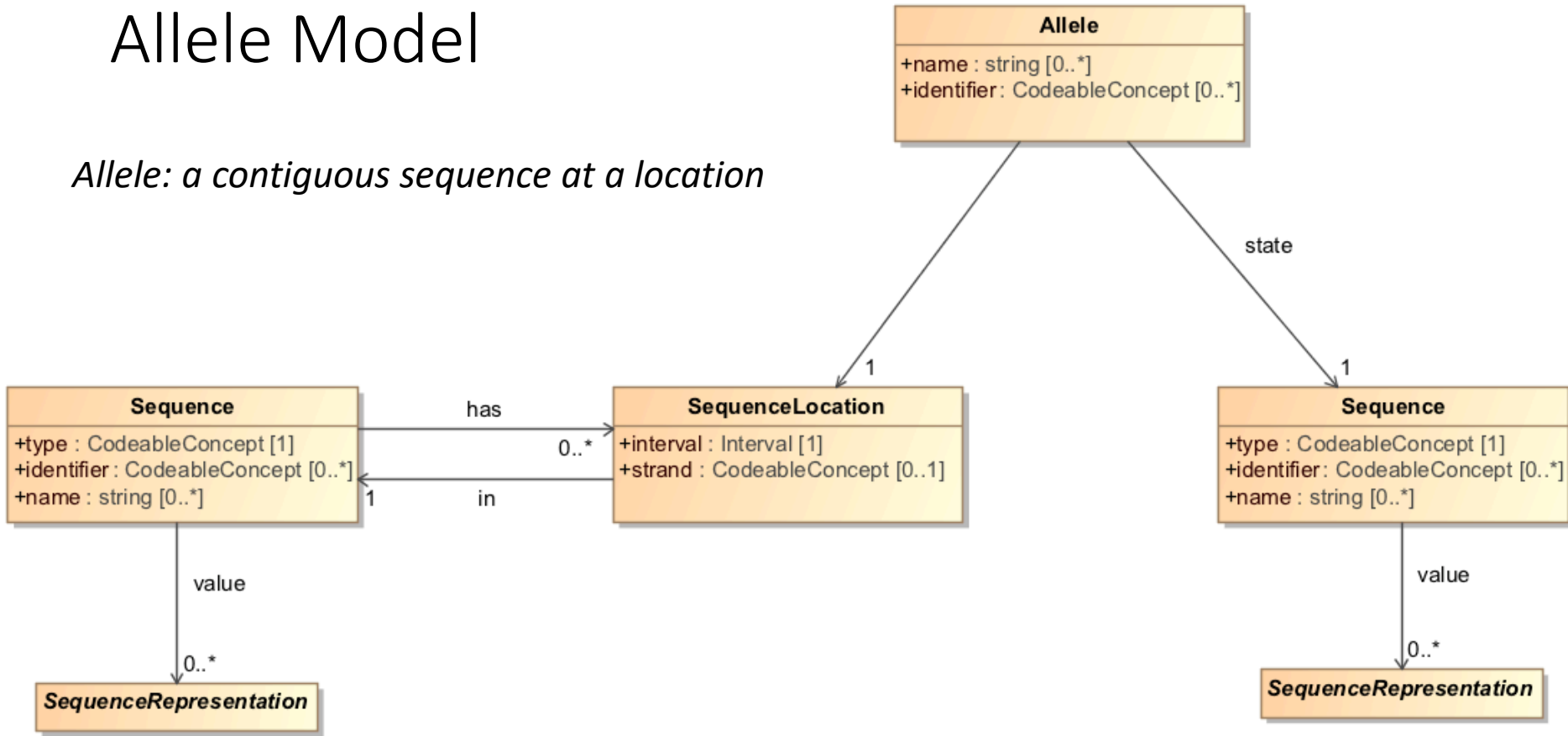






# Allele Model

*Allele: a contiguous sequence at a location*



Where is the allele located?

What was found at this location?

# Example Allele: SNV

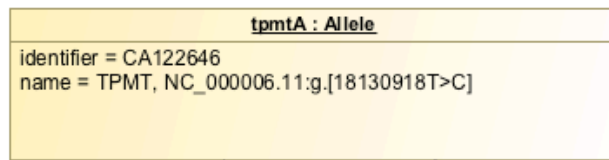
*Allele: a contiguous sequence at a location*

NC\_000006.11:g.18130918T>C

Sequence

Location

State



: state

**NC\_000006.11Seq : Sequence**

identifier = NC\_000006.11  
name = Homo sapiens chromosome 6, GRCh38.p13 Primary Assembly  
type = DNA

: in

: has

: value

**NC\_000006.11Rep : ResolvableSequence**

format = HTML  
uri = [https://www.ncbi.nlm.nih.gov/nucore/NC\\_000006.11/](https://www.ncbi.nlm.nih.gov/nucore/NC_000006.11/)

**seqLoc : SequenceLocation**

interval = coordInterval

**coordInterval : CoordinateInterval**

end = 18130918  
numberingSystem = 1-based Variant/HGVS  
start = 18130918

**stateSequence : Sequence**

type = DNA

: value

**simpleSeq : SimpleSequence**

encoding = IUPAC 1 character genomic  
value = C



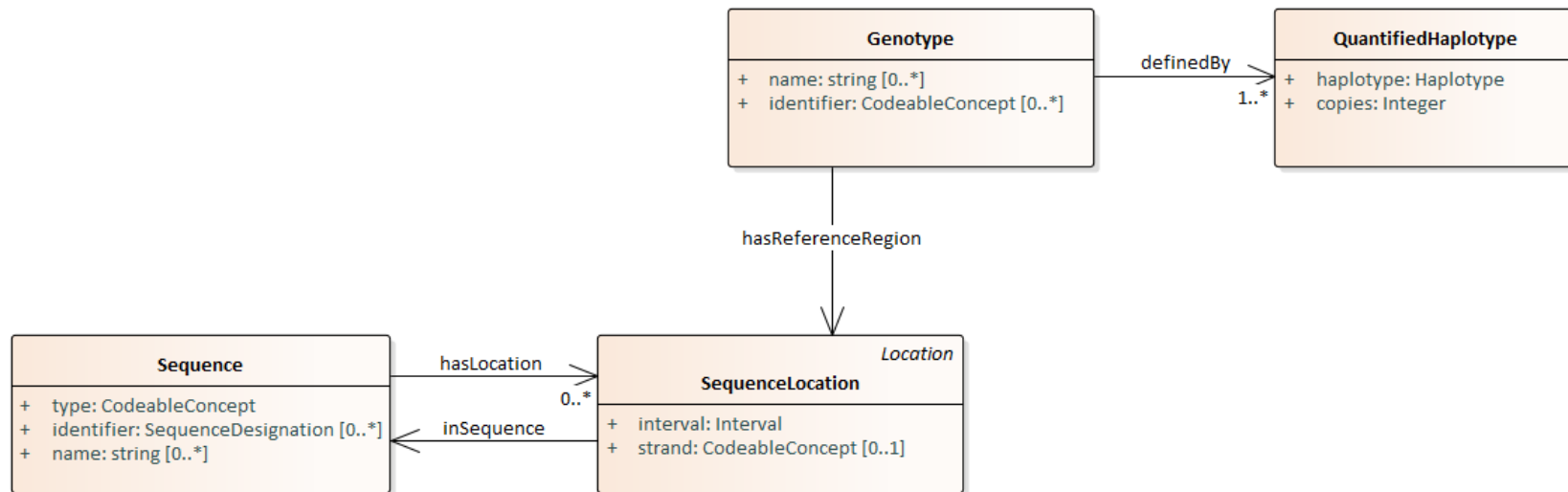
### Haplotype

- + name: string [0..\*]
- + identifier: CodeableConcept [0..\*]
- + allele: Allele [1..\*]
- + location: SequenceLocation [0..1]

### QuantifiedHaplotype

- + haplotype: Haplotype
- + copies: Integer

Draft



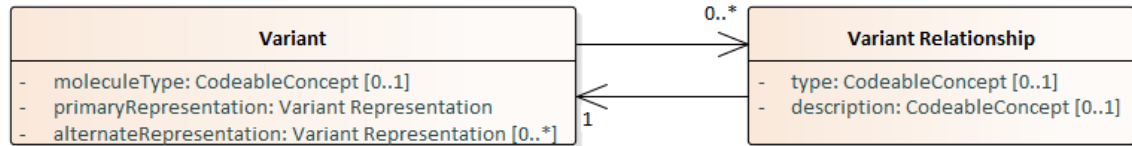
# Definitional Variation

## In scope

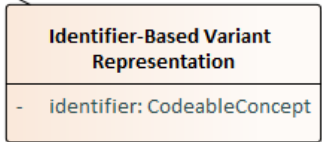
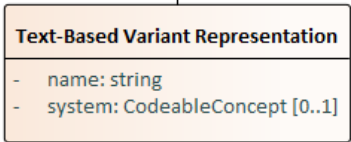
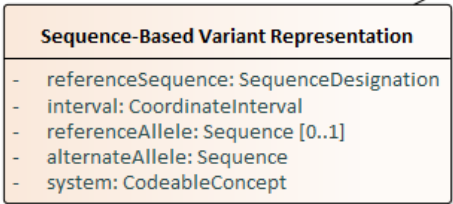
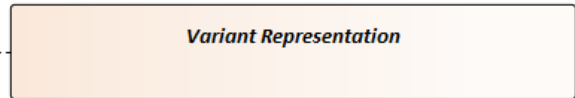
- Genomic
- Known location
- Coordinate-based
- Known change

## Out of scope (for now)

- Non-genomic
- Unknown location
- Non-coordinate based
  - Cytobands
  - Features
- Unknown or inexact change
- Categorical/variation classes
- Structural (unless it can be represented as an "in scope")



Simple and precise variants only

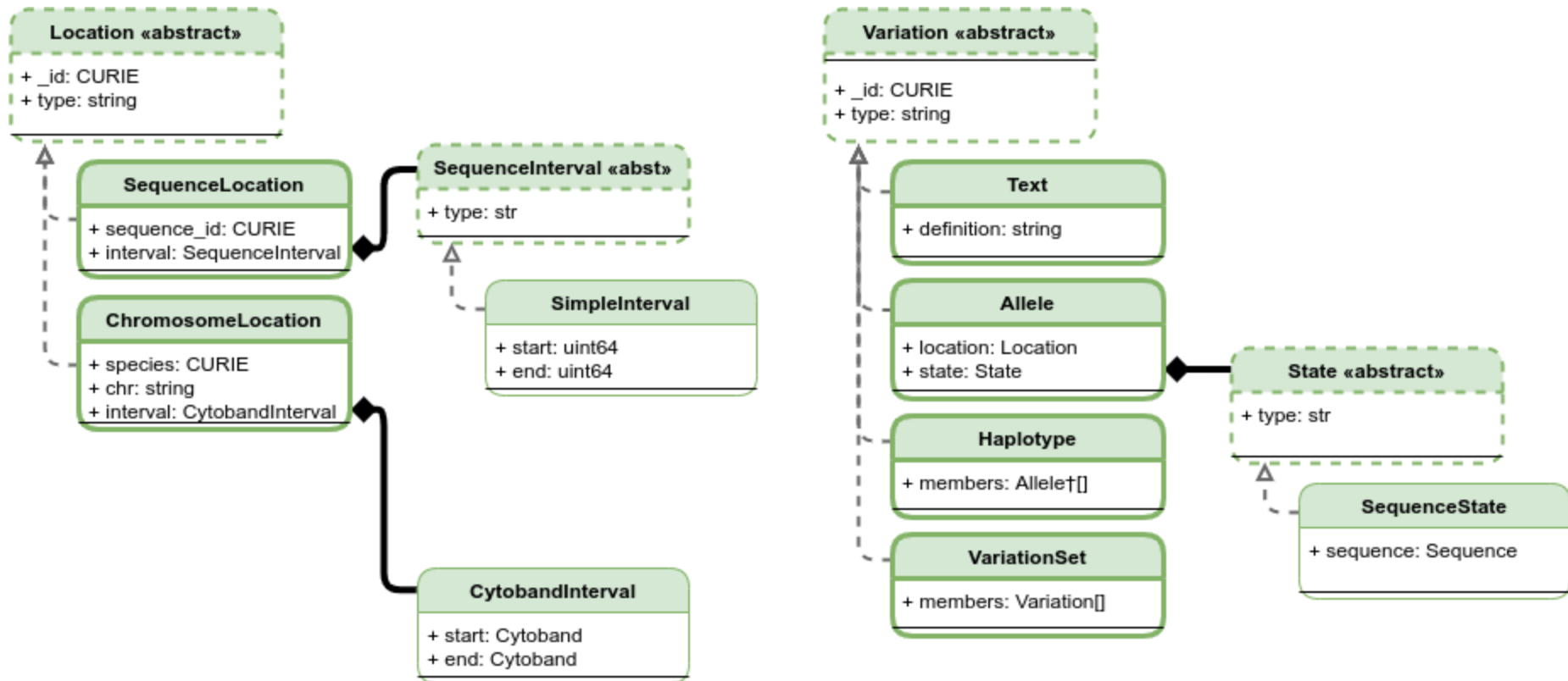


VCF Style  
 HGVS Style  
 SPDI Style  
 gnomAD  
 VRS

HGVS Nomenclature  
 Named (e.g., star alleles)  
 ISCN Nomenclature  
 HLA Nomenclature  
 Legacy names

Examples of codes for use in system

# GA4GH VRS



# Definitional Variant in Context

- Gene
- Disease/condition
- Test (including Panel)
- Observed variant (haplotype, genotype)
- Genomic source class
- Specimen (including subject and metadata)

### Inferred/computed (w/ annot)

- Variant length
- DNA change (HGVS)
- DNA change type
- AA change (HGVS)
- AA change type
- Cytogenomic nomenclature

### Measured data

- Allelic read depth
- Sample allelic frequency
- Array CGH ratio

### Interpretations

- value (e.g., present/absent)
- Allelic state (e.g., heterozygous)
- Copy number
- Chrom CN change type

### Cross-Refs <= *are these definitional?*

- Variation code
- dbSNP ID

### Related/metadata

- Genomic source class (e.g., germline)
- Gene studied (HGNC/NCBI ID) <= *used as Location?*
- Method
- Ref seq assembly <= *captured in IM via GenomeAssembly*

### Other

- Transcript ref seq
- Complex var type (e.g., compound het)

# Prioritizing Next Steps

- Other types of Variation
- Information related to Variation
- Measured, computed, inferred, interpreted information
- Observations of <things>
  
- Develop a logical model for Def Var, Seq, and related classes
  - Identify boundaries between resources
- Transform logical model into FHIR logical models
  - Facilitate review and feedback
- Integrate into IG, propose changes to resource(s)



# Bob's Crystal Ball

- CG Resources (Definitional)

- MolecularSequence

- All types of sequences (DNA, RNA, protein)
    - Representations and relationships
    - GenomeAssembly
    - Annotations and Features

- MolecularVariation

- All types of variation (SNVs, CNVs, structural; alleles; haplotypes; genotypes)
    - Representations and relationships
    - Annotations

- Generic Resources

- Assertion

- Based on SEPIO
    - Profiled for Genomic Implications and inferences related to sequences and variations

