# Kidney Metadata and Ontology Design (HuBMAP)

Katy Börner, Leonard Cross, Samuel H. Friedman, Randy Heiland, Bruce Herr II, Paul Macklin, Lisel Record, Ellen Quardokus, James Sluka, Griffin Weber

underlined: most directly involved in ontology development

Intelligent Systems Engineering
Indiana University

May 9, 2019

# Motivation

# Mapping Components (MC):
## Spatial maps of biomolecular data

Given anatomical and molecular data, develop and validate:

**1) Terminologies/Ontologies (Semantics)**

- Reference concepts, e.g., organs, organ parts, cell types, cell states

- Fiduciary concepts: Well-defined landmarks that can be provided by TMCs and used by MC to spatially orient data with respect to 3D structures

**2) 3D Spatial Models interlinked with terminology/ontology**

- Across levels (gross anatomy/organ, tissue, cell level) using hierarchical containment to localize the sample within the body

- Make landmarks visible in 3D models

**3) Interface for semantic and spatial search, filter, review, download of data.**

- Use ontology for query expansion (elastic search), semantic browsing, and as controlled vocabulary (e.g., turning on/off male/female or different cell states).

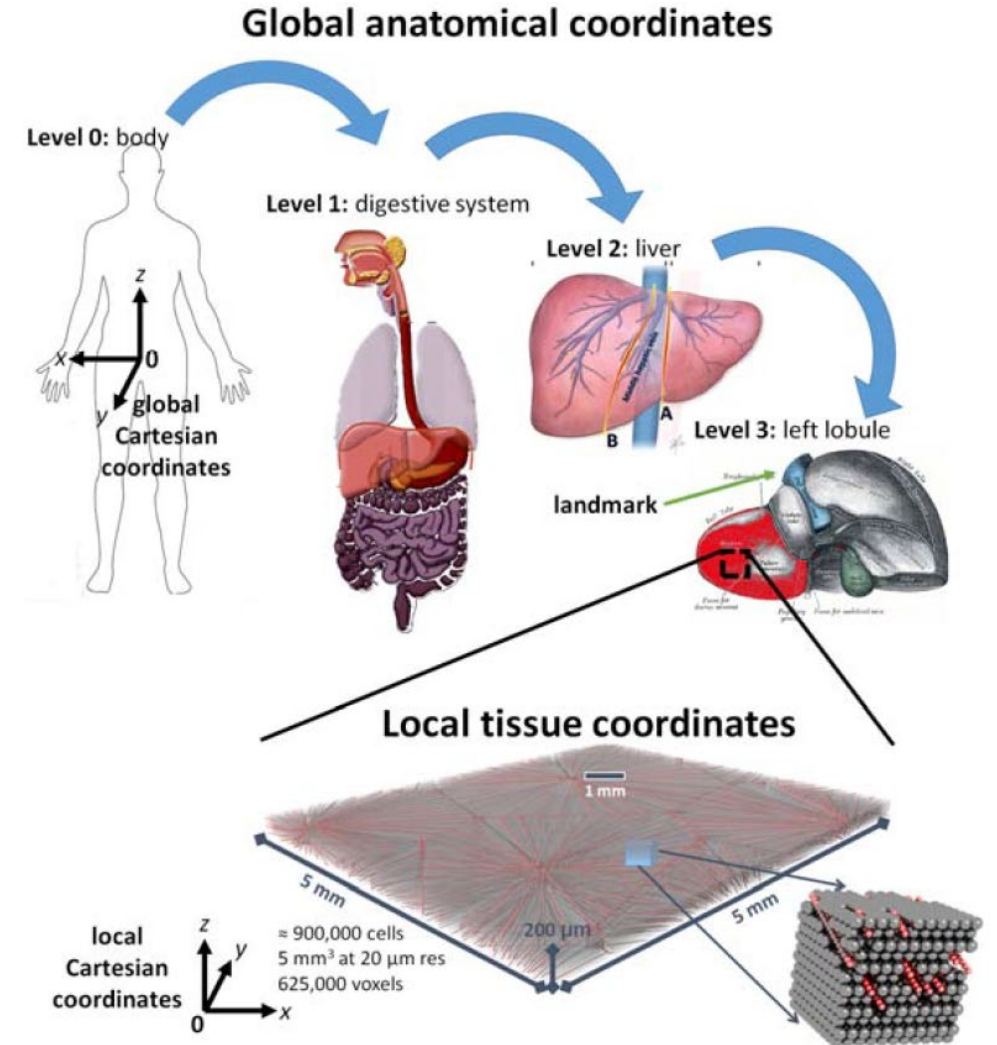- Use 3D models for spatial browsing, confirmation of proper tissue registration, exploring cell context, etc.
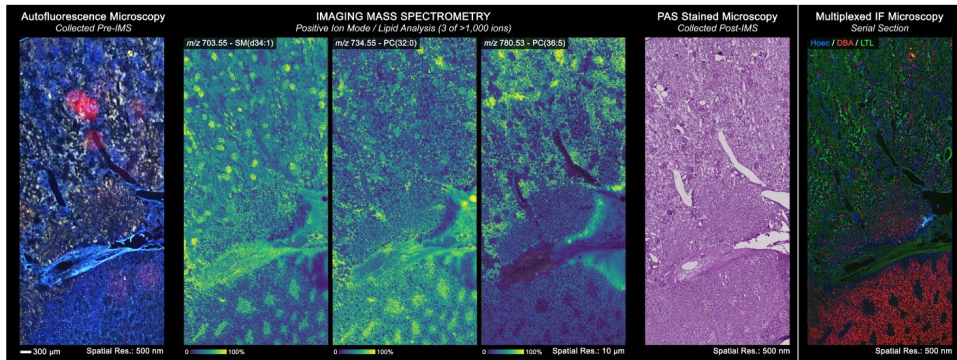


**Fig. 5. CCF concept**, navigating through the global anatomical coordinate system to insert a synthetic tissue sample (from PhysiCell[4]) into the left liver lobe with a local coordinate system.

# We must plan for heterogeneous data

**Kidney: Jeff Spraggins et al., VU**

See data on Globus, BIOMIC_patient-64354



**Heart: Shin Lin, UW**

Year 1: Tissue data for 1-2cm cubed volumes from 9 sites for 1 heart from 1 individual.

## Data Dictionary (115 rows)

| Field # Sort | Field Label Sort | Field Name Sort | Field Units | Field Data T | Lookup Tal | Low Value | High Value | Valid value | IsNullable S | Parent Fiel | Parent Fiel | Can Child b | ReadOnly Sort |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 9 | Donor //ABO: | abo | | char(3) | lkup_abo | | | | TRUE | | | | FALSE |
| 10 | Donor //Date of birth: | dob | | datetime | | | | | TRUE | | | | FALSE |
| 11 | Donor //Gender: | gender | | char(1) | lkup_gender | | | M,F | TRUE | | | | FALSE |
| 12 | Details //Age: | age_in_months | | smallint | | 0 | 1188 | | TRUE | | | FALSE | FALSE |
| 13 | Details //Age Unit: | age_unit | | char(1) | lkup_age_unit | | | M,Y | TRUE | age_in_months | | | TRUE |
| 14 | Details //Height: | hgt_cm | cm | decimal(5, 2) | | 1 | 241.3 | | TRUE | | | | FALSE |
| 15 | Donor hgt_ft // | hgt_ft | ft | int | | 0 | 7 | | TRUE | | | | TRUE |
| 16 | Donor hgt_in // | hgt_in | in | int | | 0 | 11 | | TRUE | | | | TRUE |
| 17 | Details //Weight: | wgt_kg | kg | decimal(7, 4) | | 0.454 | 294.835 | | TRUE | | | | FALSE |
| 18 | Donor wgt_lb // | wgt_lb | lbs | decimal(3, 0) | | 2 | 650 | | TRUE | | | | TRUE |
| 19 | Donor //Ethnicity/race: | race | | bigint | lkup_race_subcat_multi | | | | FALSE | | | | FALSE |
| 30 | Details //History of diabe | hist_diabetes | | smallint | lkup_histdiab_dur | | | | TRUE | | | | FALSE |
| 31 | Donor //History of cance | hist_cancer | | smallint | lkup_histcancer_site | | | | TRUE | | | FALSE | FALSE |
| 32 | Donor History of cancer , | cancer_oth_ostxt | | varchar(50) | | 1 | 50 | | TRUE | hist_cance | 999 | | FALSE |
| 33 | Details //History of hyper | hypertension | | smallint | lkup_histype_dur | | | | TRUE | | | FALSE | FALSE |

## Clinical and Spatial Metadata (21 rows)

## Cell Types, on right

## Cell States (9 rows)

| Cell states | Subset A |
|---|---|
| Proliferating cells | S-phase |
| | G2/M |
| | |
| Cell cyle arrest | G0 |
| | G1/S |
| | G2/M |

| Cell type | Subset A | Subset B | Subset C |
|---|---|---|---|
| Tubular Epithelium | Proximal tubular cells | S1 | |
| | | S2 | |
| | | S3 | |
| | Loop on Henle | Thin descending limg | |
| | | Thin ascending limb | |
| | | Thick limb | medullary |
| | | | cortical |
| | | Macula Densa | |
| | Distal convoluted tubule | | |
| | Connecting segment | | |
| | Collecting duct | Principal cells | |
| | | Intercalated cells | Type A |
| | | | Type B |
| Glomerulus | Epithelium | Visceral | |
| | | Parietal | |
| | Mesangial cells | | |
| Vasculature | Endothelium | Glomerular | |
| | | Peritubular | |
| | | Lymphatic | |
| | Pericytes | | |
| | Juxta Glomerular Cells | | |
| Interstitium | Fibroblasts | Myofibroblasts | |
| | | EPO producing cells | |
| | | Medullary fibroblasts | |
| | Mononuclear cells | Resident macrophages | |
| | | Dendritic cells | |
| | Lymphocytes | T cells | |
| | | B cells | |
| | | NK cells | |

## Cell Types (14)

| endothelial cells | |
|---|---|
| | arterial |
| | capillary |
| | venous |
| | lymphatic |
| cardiomyocytes | |
| | atrial |
| | ventricular |
| | nodal |
| fibroblasts | |
| | fibroblasts |
| | myofibroblasts |
| immune cells | |
| | macrophages |

# Data

**Kidney: Jeff Spraggins et al., VU**

Clinical and Spatial Metadata (21 rows)

| | |
|---|---|
| Sample Number: | 20 |
| Patient Number: | 64354 |
| Procedure ID: | 66598 |
| Date: | 1/30/2019 |
| Age: | 38 |
| Gender: | Female |
| Race: | White |
| Height: | 165.1 cm |
| Weight: | 115.2 kg |
| BMI: | 42.3 |
| Comorbidities: | Obesity |
| Type of Procedure: | Total Nephrectomy |
| Indications for Procedure: | Renal tumor |
| Laterality: | Left |
| Tissue Type: | kidney |
| Dimensions (mm): | L: 19 x W: 13 x H: 7 |
| Anatomical Landmark: | Lower Pole |
| Distance from Tumor: | 7 cm |
| Sample Processing: | Frozen |
| Method of Freezing: | Dry Ice/Isopentane Slurry |
| Embedding Media: | CMC |

**Jeff invited feedback on current fields and format.**
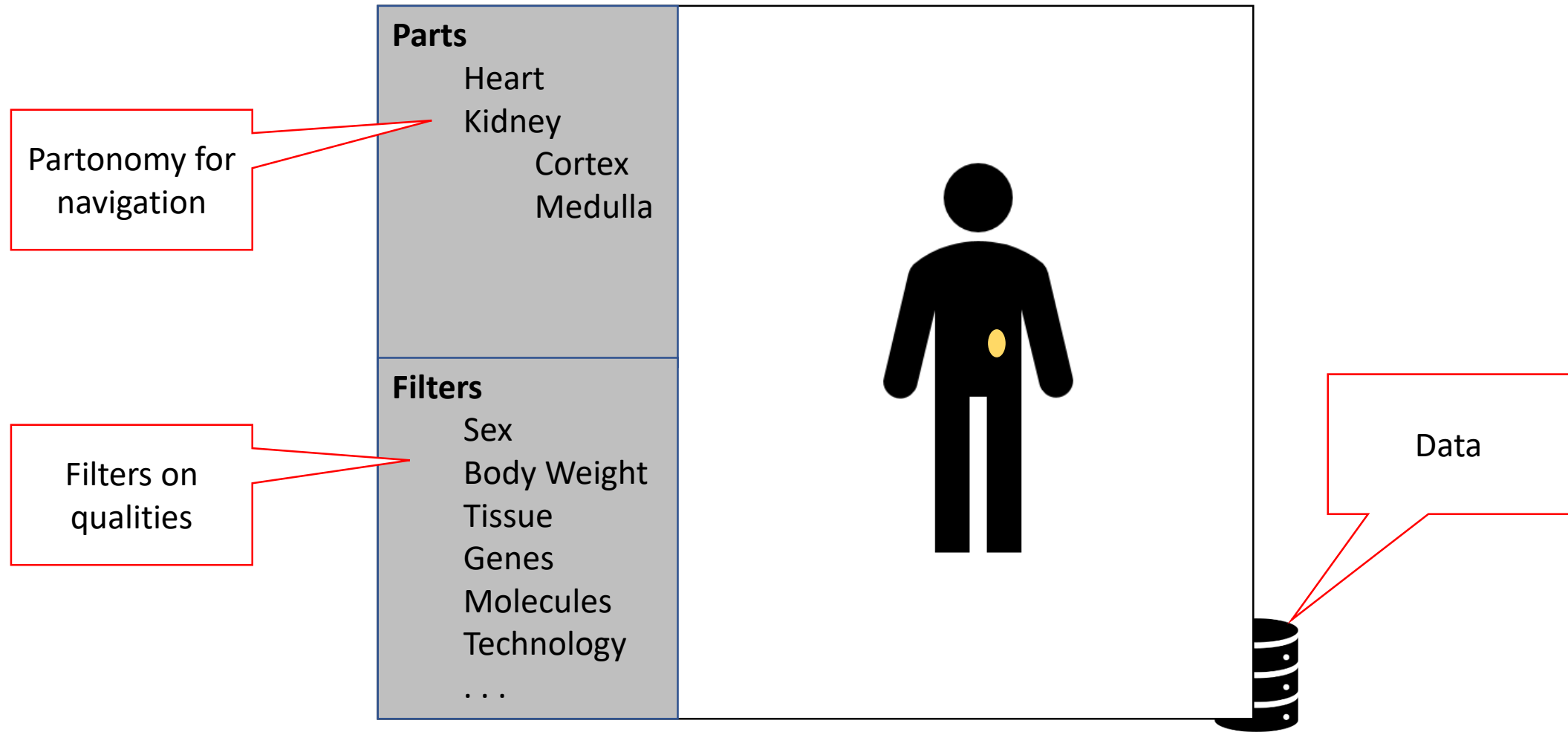
**Heart: Shin Lin, UW**

Data Dictionary (115 rows)

| Field # Sort | Field Label Sort | Field Name Sort | Field Units | Field Data T | Lookup Ta | Low Value | High Value | Valid value |
|---|---|---|---|---|---|---|---|---|
| 9 | Donor //ABO: | abo | | char(3) | lkup_abo | | | |
| 10 | Donor //Date of birth: | dob | | datetime | | | | |
| 11 | Donor //Gender: | gender | | char(1) | lkup_gender | | | M,F |
| 12 | Details //Age: | age_in_months | | smallint | | 0 | 1188 | |
| 13 | Details //Age Unit: | age_unit | | char(1) | lkup_age_unit | | | M,Y |
| 14 | Details //Height: | hgt_cm | cm | decimal(5, 2) | | 1 | 241.3 | |
| 15 | Donor hgt_ft // | hgt_ft | ft | int | | 0 | 7 | |
| 16 | Donor hgt_in // | hgt_in | in | int | | 0 | 11 | |
| 17 | Details //Weight: | wgt_kg | kg | decimal(7, 4) | | 0.454 | 294.835 | |
| 18 | Donor wgt_lb // | wgt_lb | lbs | decimal(3, 0) | | 2 | 650 | |
| 19 | Donor //Ethnicity/race: | race | | bigint | lkup_race_subcat_multi | | | |
| 30 | Details //History of diabe | hist_diabetes | | smallint | lkup_histdiab_dur | | | |
| 31 | Donor //History of cance | hist_cancer | | smallint | lkup_histcancer_site | | | |
| 32 | Donor History of cancer | cancer_oth_ostxt | | varchar(50) | | 1 | 50 | |
| 33 | Details //History of hyper | hypertension | | smallint | lkup_histhype_dur | | | |

Please complete **TMC Landmarks Survey** at
https://goo.gl/forms/x9F8cP1GIzprDxbI2
(complete one survey per organ)

# Goal: Facilitate navigation of multiscale data

**Parts**

    Heart

    Kidney

        Cortex

        Medulla

**Filters**

    Sex

    Body Weight

    Tissue

    Genes

    Molecules

    Technology

    . . .

Partonomy for navigation

Filters on qualities

Data

# Overall CCF Approach

# CCF Ontology: some guiding principles

- **Reuse** existing <u>ontologies</u> and <u>data formats</u> developed for projects similar to HuBMAP to the greatest extent possible
  - GUDMAP / RBK
  - Human Cell Atlas
  - …
- **Reuse** domain-specific ontologies and data formats
  - OME-Tiff (Open Microcopy Community advanced image format)
  - MIAME (Minimum Information About a *Microarray* Experiment)
  - …
- **Leverage** HuBMAP domain expertise!
  - Each TMC is an expert in its organ. Capture this in the organ-specific ontologies.
- Use a **standard Ontology format** and development tools
  - We will use OWL
  - Include test cases in the ontology itself (e.g. both A-box and T-box) for testing, validation and demonstration purposes.
- **Cross-link with existing ontologies** as much as possible
- May need separate partOf (or class/subclass) trees for **simplified navigation** in GUI.

# CCF: Source Ontologies

**Anatomic/Phenotypic**

- Uberon

- Foundational Model of Anatomy (FMA) (has anatomical terms NOT in Uberon)

- Human Phenotype Ontology (HPO)

- Phenotype and Trait Ontology (PATO)

- Organ specific: Kidney Tissue Atlas Ontology (KTAO) and LungMAP

**Tissue/Data Collection**

- Biological Spatial Ontology (BSPO)

- Ontology of Biomedical Investigations (OBI)

- EDAM (Bioinformatics concepts)

**(Sub-)Cellular**

- Cell Ontology (CL)

- Gene Ontology (GO)

- Chemical Entities of Biological Interest (ChEBI)

- RNA Ontology (RNAO)

- Protein Ontology (PR)

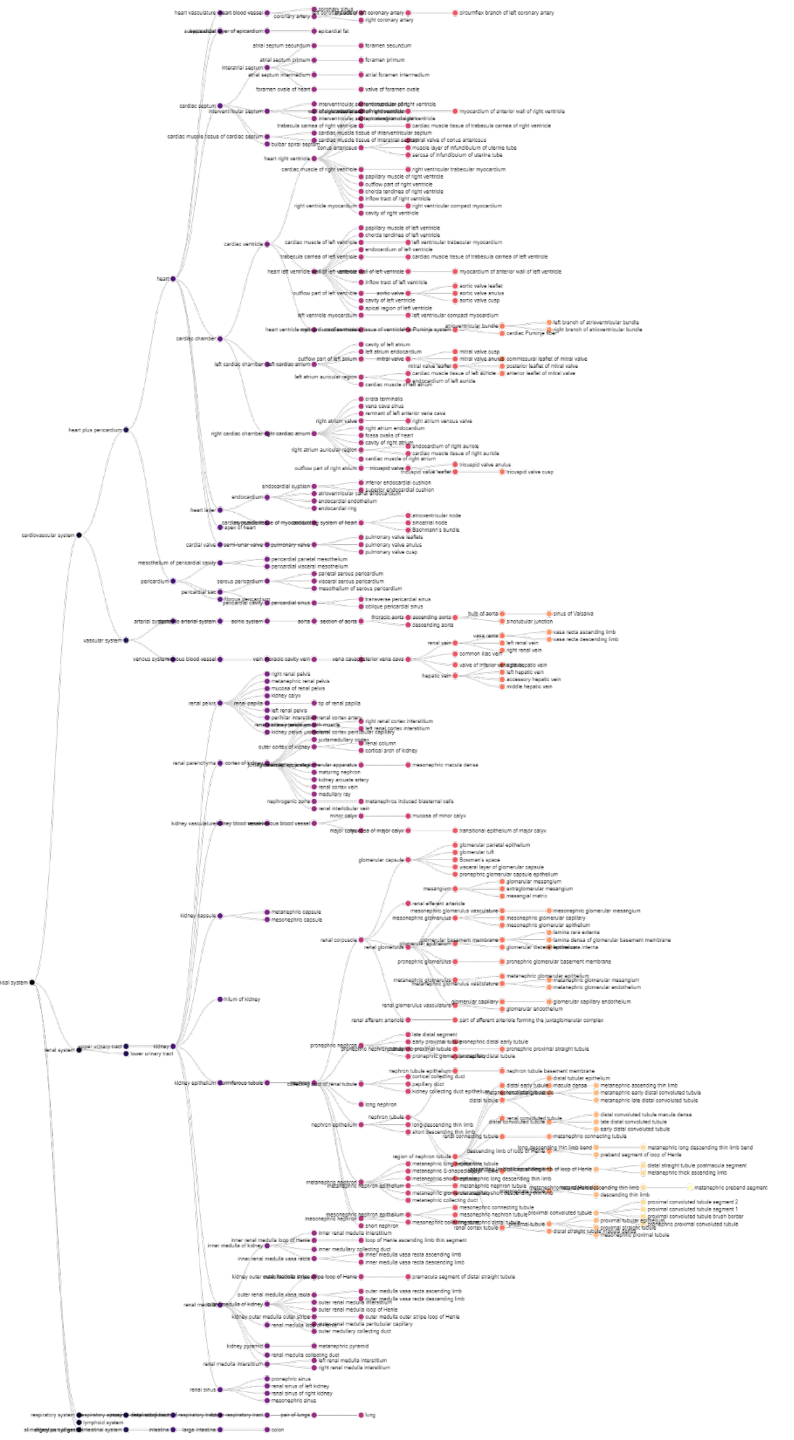- Cell Behavior Ontology (CBO)

**Metadata**

- Basic Formal Ontology (BFO)

- Information Artifact Ontology (IAO)

- Ontology of units of Measure (OM)

- Provenance, Authoring and Versioning ontology (PAV)
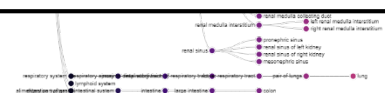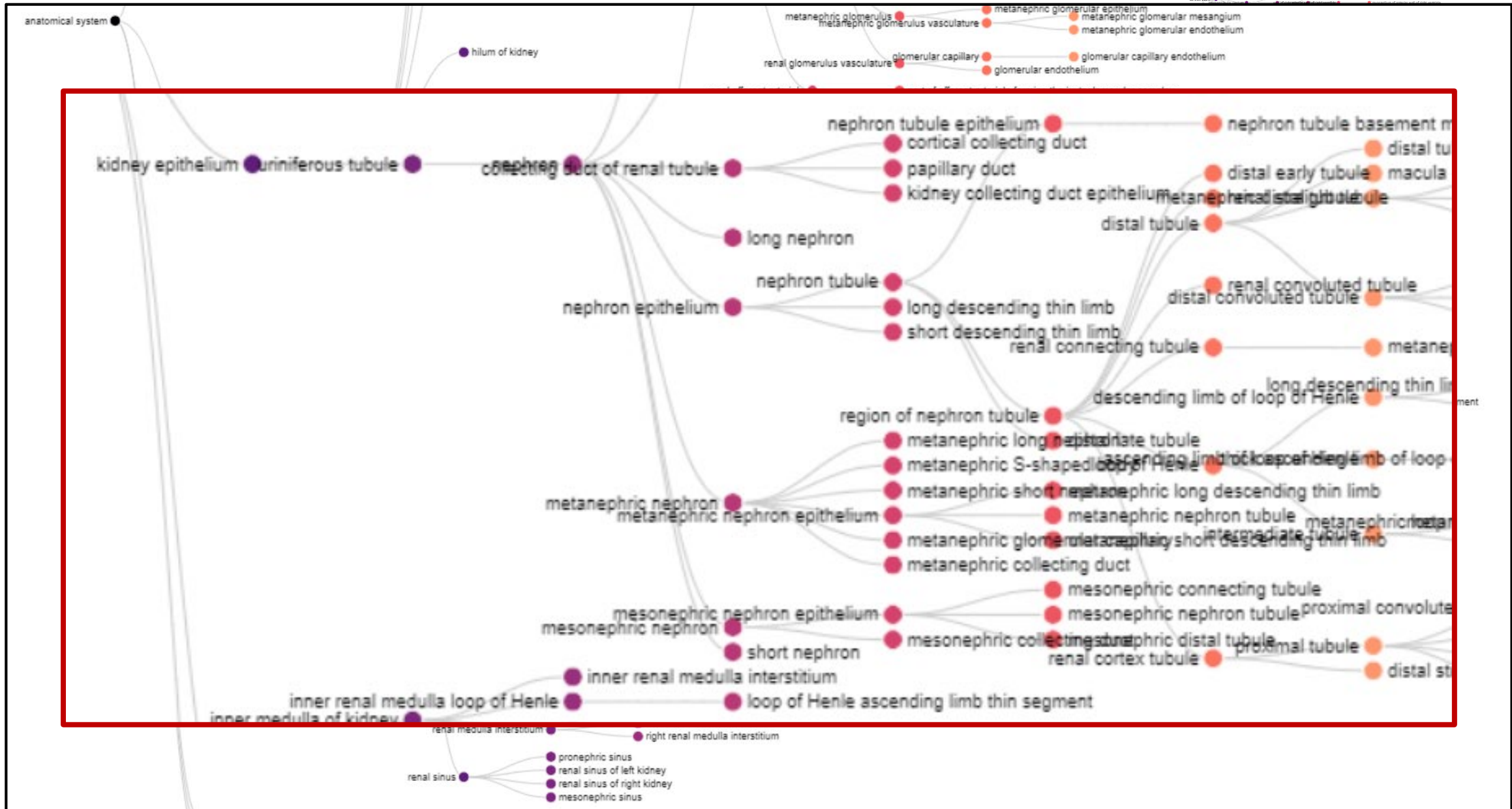
- VIVO (Identifying researchers)

**MeSH and NCI Thesaurus**

# Ontology

# Current CCF Ontology:

- Use Uberon and user-supplied tables of terms to create a SLIM ontology

- Users (initially TMCs) can request missing terms as needed

- "partOf" and other partonomy terms used to help relate concepts
  - Requires domain expertise!
  - Individual TMCs will need to pitch in for their specific organs to refine
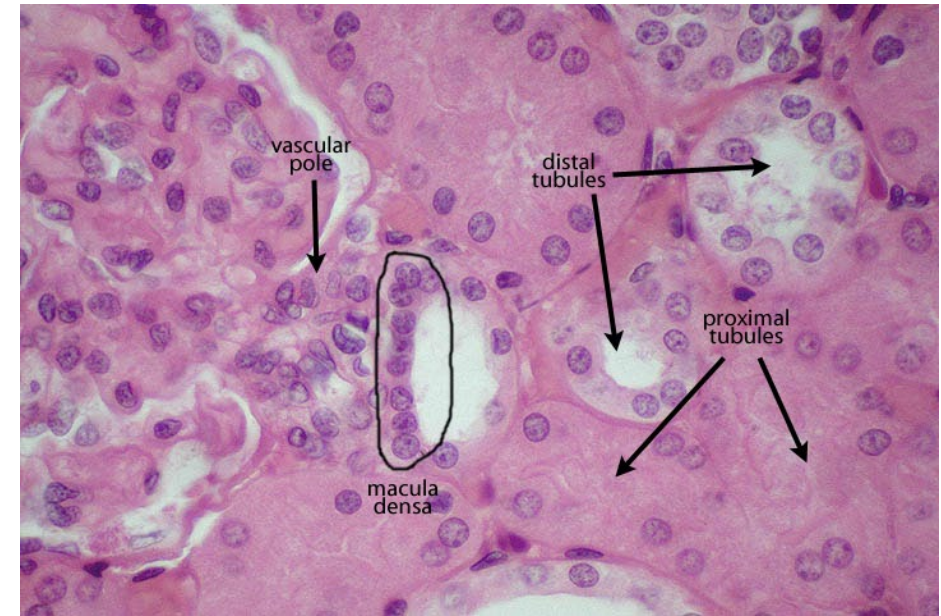
- **Click here to visualize the current CCF ontology**

anatomical system

hilum of kidney

renal system — upper urinary tract — kidney
lower urinary tract

kidney epithelium — uriniferous tubule

metanephric glomerulus — metanephric glomerular epithelium
metanephric glomerulus vasculature — metanephric glomerular mesangium
metanephric glomerular endothelium

renal glomerulus vasculature — glomerular capillary — glomerular capillary endothelium
glomerular endothelium

renal afferent arteriole — part of afferent arteriole forming the juxtaglomerular complex

late distal segment
early proximal tubule — pronephric distal early tubule
pronephric nephron — pronephric nephron tubule — nephric proximal tubule — pronephric proximal straight tubule
pronephric glomerulus capillary — pronephric distal tubule

nephron tubule epithelium — nephron tubule basement membrane
cortical collecting duct — distal tubular epithelium
collecting duct of renal tubule — papillary duct — distal early tubule — macula densa — metanephric ascending thin limb
kidney collecting duct epithelium — metanephric distal tubule — metanephric early distal convoluted tubule
metanephric late distal convoluted tubule

distal tubule

long nephron

distal convoluted tubule macula densa
renal convoluted tubule — late distal convoluted tubule
nephron tubule — distal convoluted tubule — early distal convoluted tubule
nephron epithelium — long descending thin limb
short descending thin limb — renal connecting tubule — metanephric connecting tubule

long descending thin limb bend — metanephric long descending thin limb bend
descending limb of loop of Henle — prebend segment of loop of Henle
region of nephron tubule
metanephric long nephron tubule — distal straight tubule postmacula segment
metanephric S-shaped body — ascending limb of loop of Henle — metanephric thick ascending limb
metanephric short nephron — metanephric long descending thin limb
metanephric nephron — metanephric nephron tubule — metanephric descending thin limb — metanephric prebend segment
metanephric nephron epithelium — intermediate tubule
metanephric glomerulus capillary — metanephric short descending thin limb — descending thin limb
metanephric collecting duct

mesonephric connecting tubule
mesonephric nephron epithelium — mesonephric nephron tubule — proximal convoluted tubule segment 2
mesonephric nephron — mesonephric collecting duct — mesonephric distal tubule — proximal convoluted tubule segment 1
short nephron — proximal convoluted tubule brush border
proximal convoluted tubule — pronephric proximal convoluted tubule
proximal tubular epithelium
proximal tubule — proximal straight tubule
renal cortex tubule — distal straight tubule macula densa — mesonephric proximal tubule
distal straight tubule

inner renal medulla interstitium
inner renal medulla loop of Henle — loop of Henle ascending limb thin segment
inner medulla of kidney — inner medullary collecting duct
inner renal medulla vasa recta — inner medulla vasa recta ascending limb
inner medulla vasa recta descending limb

kidney outer medulla inner stripe loop of Henle — premacula segment of distal straight tubule

outer medulla vasa recta ascending limb
outer renal medulla vasa recta — outer medulla vasa recta descending limb
renal medulla — outer medulla of kidney — outer renal medulla interstitium
kidney outer medulla outer stripe — outer renal medulla loop of Henle
outer medulla outer stripe loop of Henle
renal medulla loop of Henle — outer renal medulla peritubular capillary
outer medullary collecting duct

kidney pyramid — metanephric pyramid

renal medulla collecting duct
renal medulla interstitium — left renal medulla interstitium
right renal medulla interstitium

pronephric sinus
renal sinus of left kidney
renal sinus — renal sinus of right kidney
mesonephric sinus

# Data Formats

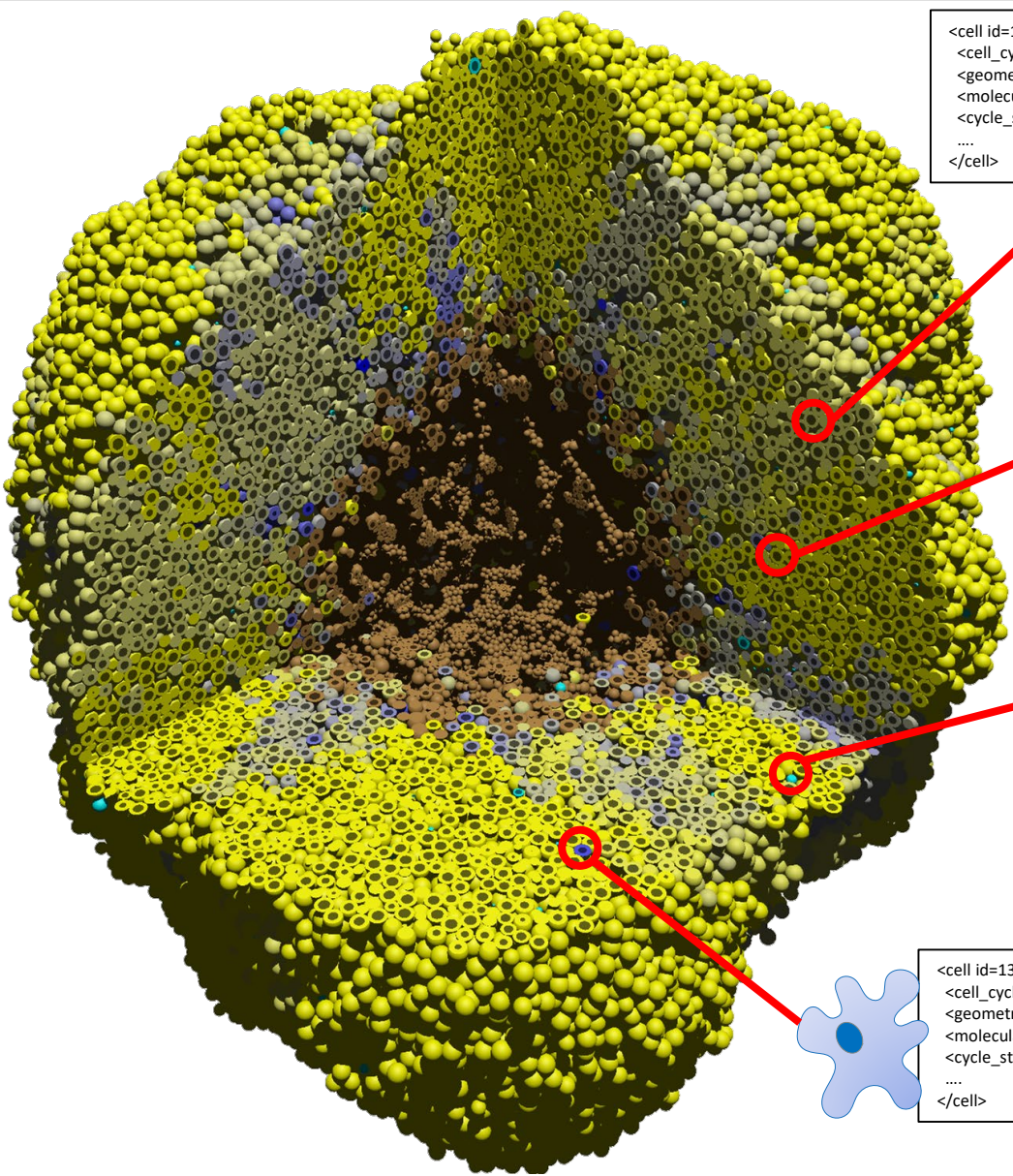# IU CCF Initial (v0.5.0) Image Formats

- **<u>Basic image</u>**: OME-Tiff
  - 2D to 4D data (includes movies)
  - more than three "color channels"
  - More flexible "color" data format (int, float, etc.)

- **<u>Regions of images</u>**: SVG with annotations (aligned with a particular OME-Tiff)

- **<u>Volumetric</u>** (e.g., computed tomography, MR, ultrasound, …)
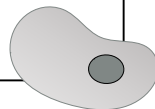  - Data normally represented as volumes or surfaces

# More data



- RNAseq and other OMICS data
  - Challenges in:
    - Data formatting and visual representation of spatial data
    - Extracting *knowledge* from complex, often noisy data
    - Harmonizing data sets from different platforms
    - Validating and benchmarking
- Extracted information
  - Segmented cells and structures
    - SVG overlays
      - e.g., cell apical surfaces, glomerulus podocyte, renal corpuscle, Bowman's space, …
    - Vectorized annotations (e.g., as in MultiCellDS – see below)
  - Multiplexed and massively multichannel imaging (e.g., MALDI)
    - Associate vectors of measurements with segmented structures
    - Additional ontology-driven annotations for the structures
      - e.g., cell type and state by Cell Ontology …
      - Cell morphometric annotations …

# Vectorized annotations of extracted cell features:
## a step from data towards knowledge



```
<cell id=137 type="renal carcinoma">
  <cell_cycle />
  <geometric_properties />
  <molecular_properties />
  <cycle_state/>
  ....
</cell>
```

```
<cell id=137 type="epithelial tubal">
  <cell_cycle />
  <geometric_properties />
  <molecular_properties />
  <cycle_state/>
  ....
</cell>
```

```
<cell id=139 type="podocyte">
  <cell_cycle />
  <geometric_properties />
  <molecular_properties />
  <cycle_state/>
  ....
</cell>
```

```
<cell id=138 type="fibroblast">
  <cell_cycle />
  <geometric_properties />
  <molecular_properties />
  <cycle_state/>
  ....
</cell>
```

```xml
<cellular_information>
  <cell_populations>
    <cell_population type="individual">
      <custom>
        <simplified_data type="matlab" source="BioFVM">
          <filename>output00000540_cells.mat</filename>
        </simplified_data>
        <simplified_data type="matlab" source="PhysiCell">
          <labels>
            <label index="0" size="1">ID</label>
            <label index="1" size="3">position</label>
            <label index="4" size="1">total_volume</label>
            <label index="5" size="1">cell_type</label>
            <label index="6" size="1">cycle_model</label>
            <label index="7" size="1">current_phase</label>
            <label index="8" size="1">elapsed_time_in_phase</label>
            <label index="9" size="1">nuclear_volume</label>
            <label index="10" size="1">cytoplasmic_volume</label>
            <label index="11" size="1">fluid_fraction</label>
            <label index="12" size="1">calcified_fraction</label>
            <label index="13" size="3">orientation</label>
            <label index="16" size="1">polarity</label>
            <label index="17" size="1">migration_speed</label>
            <label index="18" size="3">motility_vector</label>
            <label index="21" size="1">migration_bias</label>
            <label index="22" size="3">motility_bias_direction</label>
            <label index="25" size="1">persistence_time</label>
            <label index="26" size="1">motility_reserved</label>
            <label index="27" size="1">oncoprotein</label>
            <label index="28" size="1">elastic coefficient</label>
            <label index="29" size="1">kill rate</label>
            <label index="30" size="1">attachment lifetime</label>
            <label index="31" size="1">attachment rate</label>
          </labels>
          <filename>output00000540_cells_physicell.mat</filename>
        </simplified_data>
      </custom>
    </cell_population>
  </cell_populations>
</cellular_information>
```

**Bonus 1:** Can represent domain expert knowledge via expert-defined features.

**Bonus 2:** Extracted features could be *directly imported* into computational models.
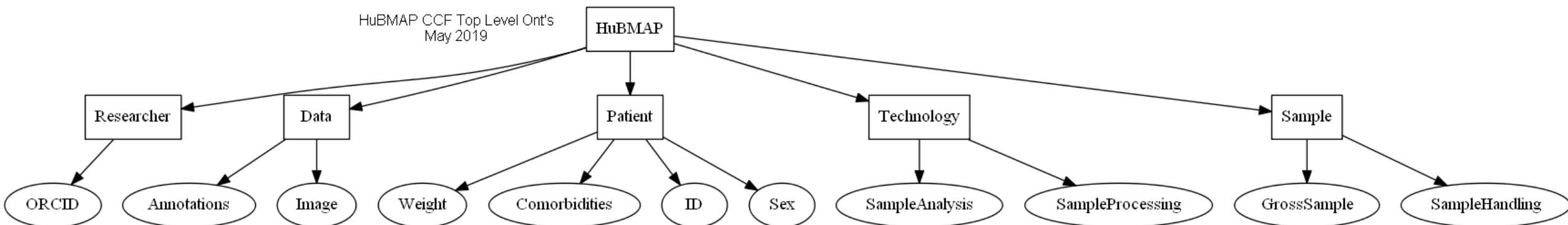
# Metadata

# Some Key Types of Metadata:

- Biological (patient specific data)
  - May include patient differentia such as disease state
- Technological (analysis tools; MS, immunohistochemical, RNAseq, …)
  - Where possible, respect and reuse any technology domain-specific standards, ontologies, etc.
  - Will often include the biological results (e.g., gene expression levels)
- Interpretive
  - Summary of findings, e.g., conversion of gene expression lists into list of highly impacted pathways
  - Algorithms and software used in interpretation and analysis
- Provenance
  - Who processed and analyzed samples
  - Points of contact (to answer questions on the samples and analyses)
  - Can be important for QC and account
- Use metadata
  - Licenses, citation information, ….

# Many data and metadata types and files

- Patient
- Sample
- Analysis Technology
- Results
- Researcher
- . . .

# Example: Metadata for Similar Large-Scale Data & Imaging Projects

**GUDMAP:**



RID: Persistent citable resource identifier
Imaging Data
Genes
Species
Stage
Anatomical Sources
Assay Type
Preparation
Principle Investigator . . .

Kidney Tissue Atlas Ontology from KPMP

**Human cell Atlas, SPARC (informed by BIDS)**



Structured Datasets (like BIDS) provide...

- A convention for organizing data files into folders
- A set of descriptive files that contain information on subjects, experimental information, data set descriptions
- A set of naming conventions for files
- A means to extend the core structure to accommodate most data acquisitions

# Challenge: Inconsistent metadata specifications

**Kidney: Jeff Spraggins et al., VU**

Clinical and Spatial Metadata (21 rows)

| | |
|---|---|
| Sample Number: | 20 |
| Patient Number: | 64354 |
| Procedure ID: | 66598 |
| Date: | 1/30/2019 |
| Age: | 38 |
| Gender: | Female |
| Race: | White |
| Height: | 165.1 cm |
| Weight: | 115.2 kg |
| BMI: | 42.3 |
| Comorbidities: | Obesity |
| Type of Procedure: | Total Nephrectomy |
| Indications for Procedure: | Renal tumor |
| Laterality: | Left |
| Tissue Type: | kidney |
| Dimensions (mm): | L: 19 x W: 13 x H: 7 |
| Anatomical Landmark: | Lower Pole |

**Heart: Shin Lin, UW**

Data Dictionary (115 rows)

| Field # Sort | Field Label Sort | Field Name Sort | Field Units | Field Data T | Lookup Tal | Low Value | High Value | Valid value |
|---|---|---|---|---|---|---|---|---|
| 9 | Donor //ABO: | abo | | char(3) | lkup_abo | | | |
| 10 | Donor //Date of birth: | dob | | datetime | | | | |
| 11 | Donor //Gender: | gender | | char(1) | lkup_gender | | | M,F |
| 12 | Details //Age: | age_in_months | | smallint | | 0 | 1188 | |
| 13 | Details //Age Unit: | age_unit | | char(1) | lkup_age_unit | | | M,Y |
| 14 | Details //Height: | hgt_cm | cm | decimal(5, 2) | | 1 | 241.3 | |
| 15 | Donor hgt_ft // | hgt_ft | ft | int | | 0 | 7 | |
| 16 | Donor hgt_in // | hgt_in | in | int | | 0 | 11 | |
| 17 | Details //Weight: | wgt_kg | kg | decimal(7, 4) | | 0.454 | 294.835 | |
| 18 | Donor wgt_lb // | wgt_lb | lbs | decimal(3, 0) | | 2 | 650 | |
| 19 | Donor //Ethnicity/race: | race | | bigint | lkup_race_subcat_multi | | | |
| 30 | Details //History of diabe | hist_diabetes | | smallint | lkup_histdiab_dur | | | |
| 31 | Donor //History of cance | hist_cancer | | smallint | lkup_histcancer_site | | | |
| 32 | Donor History of cancer | cancer_oth_ostxt | | varchar(50) | | 1 | 50 | |
| 33 | Details //History of hyper | hypertension | | smallint | lkup_histhype_dur | | | |

# AMIS (absolute minimal information solution)

**Figure 2:** Partial AMIS ontology for the CCF UI. White boxes indicate metadata terms included in the Year 1 metadata ontology. Note that the CCF UI ontology only covers the CCF data needed for operation of the UI and that is reflected in the ontology above. A critical issue is that human understandable biological interpretations of the various data sets is required for effective use of the image datasets by end users via the UI.

# IU CCF Initial (v0.5.0) <u>Patient Metadata</u>

| Column Header | Data Type | Comments |
|---|---|---|
| **HuBMAP Sample ID** | string | *assigned by PSE/IEC* |
| **HuBMAP Patient Number** | string | |
| **Source Sample ID** | string | *assigned by clinical unit, deidentified.* |
| **Source Patient Number** | string | |
| **Procedure ID** | string | |
| ~~Procedure Date~~ | ~~formatted date~~ | |
| **Species** | Human (STY:T016) | *required, though always human* |
| **Age** | decimal years | |
| **Sex** | M/F/u (SNOMED extended with "unkown") | |
| **Race** | | *Race, ethnicity, strain* |
| **Height** | meter | |
| **Weight** | killogram | |
| **BMI** | float (calculated localy from height and weight) | |
| **Comorbidities and other clinical classifications** | MEDRA, SNOMED CT or MeSH terms | |
| **Type of Procedure** | MEDRA, SNOMED CT or MeSH terms | |
| **Indications for Procedure** | MEDRA, SNOMED CT or MeSH terms | |
| **Laterality** | MEDRA, SNOMED CT or MeSH terms | |
| **Tissue Type** | MEDRA, SNOMED CT or MeSH terms | |
| **Anatomical Landmark** | MEDRA, SNOMED CT or MeSH terms | |
| **Displacement from Landmark** | Affine transformation matrix | |

# Next?

# Discussion Points

- **<u>Year 2 Plans</u>**
  - Finalize metadata formats with input from the TMCs
  - Unify data formats across TMCs' technologies
  - Work closely with the CCF IU team to ensure data content and formatting is compatible with the needs and goals of the user interface.
  - Work closely with TMCs to insure all anatomical scales are represented with ontology terms

- **<u>Other TMC needs?</u>**
  - Are there unmet needs?
  - Any "must have" features or terms to describe your data?

# Extra reference materials

Global Data Store

CCF Data Store

Data Traceback

Multiple TMC's

PSC (IEC)

User Interface (IU)

End user Interface

| | Data Wrangler |
|---|---|
| | Data Filter (HIPPA, quality control, ...) |

Additional Processing

Tools to assist in sample spatial registration

- Provenance
- Patient
- Sample
- Sample Processing
- Technology (MS, IH, ...)
- Analysis
- Etc.

Propagate needs back to TMC's

- Only the data needed for the GUI

TMC: Tissue Mapping Center
PSC: Pittsburgh Supercomputing center

# What is an ontology?

*An ontology is a <u>particular</u> view of reality that encompasses a defined set of objects, processes and relationships within that reality.*

"Ontological Commitment" →

| Controlled Vocabulary | Hierarchy of Terms (isA) | Full Ontology |
|---|---|---|
| Cell<br>Hepatocyte<br>Leukocyte<br>Organ<br>Heart<br>Liver | 1. Cell<br>   a. Hepatocyte<br>   b. Leukocyte<br>2. Organ<br>   a. Heart<br>   b. Liver | 1. Cell<br>   a. Hepatocyte<br>   b. Leukocyte<br>2. Organ<br>   a. Heart<br>   b. Liver    *partOf* |

adjacentTo, containedIn, derivesFrom, definedBy, participatesIn, contributesTo, downStreamOf, …

# Annotating Images and Images containing identified (not just identifiable) things

In the H&E stained kidney image below what regions (e.g., nephron, tubule, …), cell types, cell states etc. are present? Regions must be annotated, presumably by the TMCs. It is not enough to simply say this image contains cell types X, Y and Z in cell states 1, 2, and 3:



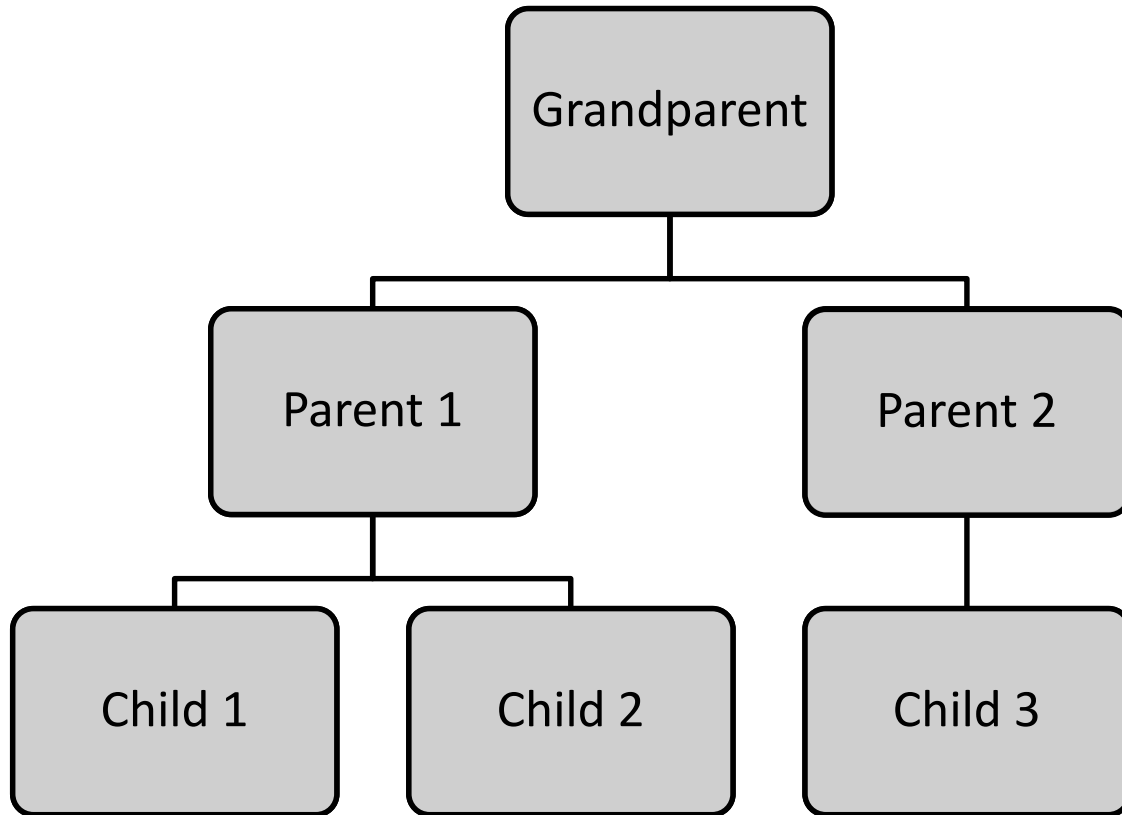This image has been annotated (though the regions aren't very clear):

# Solution: Multiple Slim Ontologies

- Gene Ontology (GO) is a massive ontology. Researchers didn't need the entire ontology, but only a subset of it.
- GO users create "slim" versions of GO:
  - Subsets that contained the terms needed along with the necessary ancestor and children terms
  - Can better illuminate the science of what is going on rather than being overwhelmed by too much information
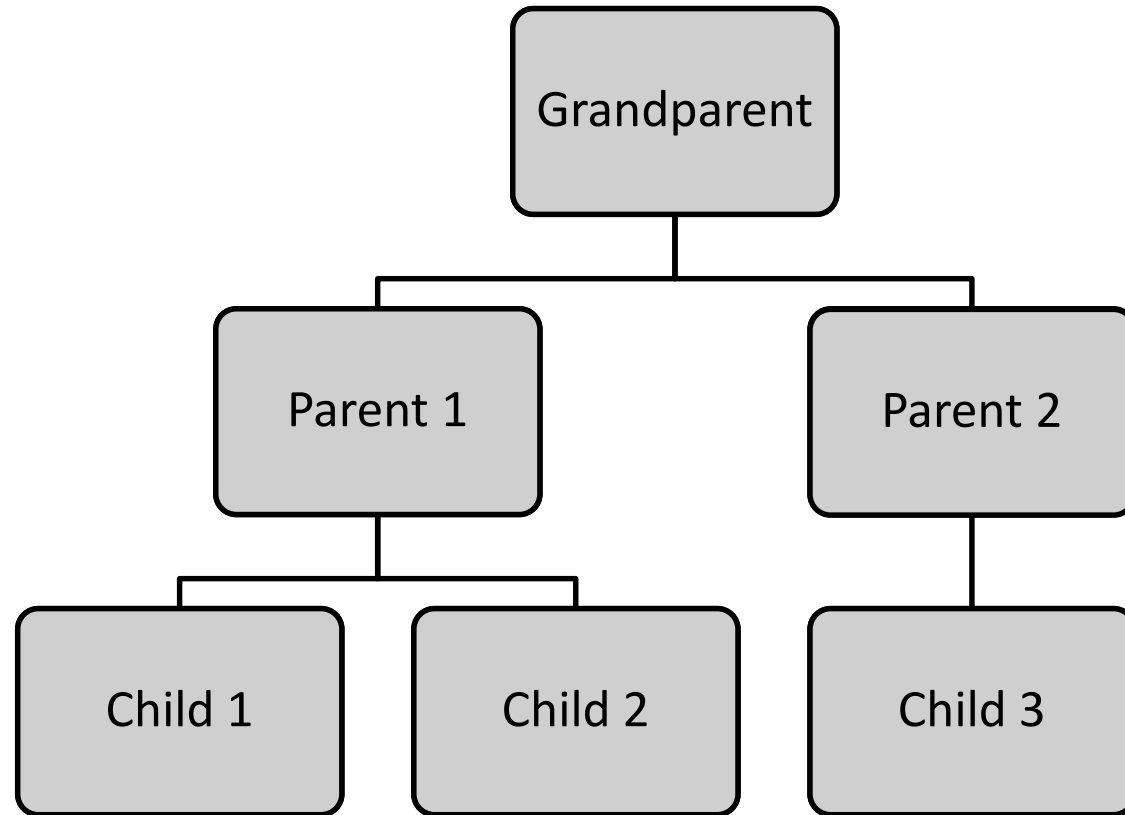- Terms in the main ontology can receive "slim" annotations, indicating their use in the corresponding "slim" ontology

# Multiple Slim Ontologies to Main HuBMAP/CCF Ontology

- Instead of creating a single ontology file, we have software create the ontology based on the needed terms.
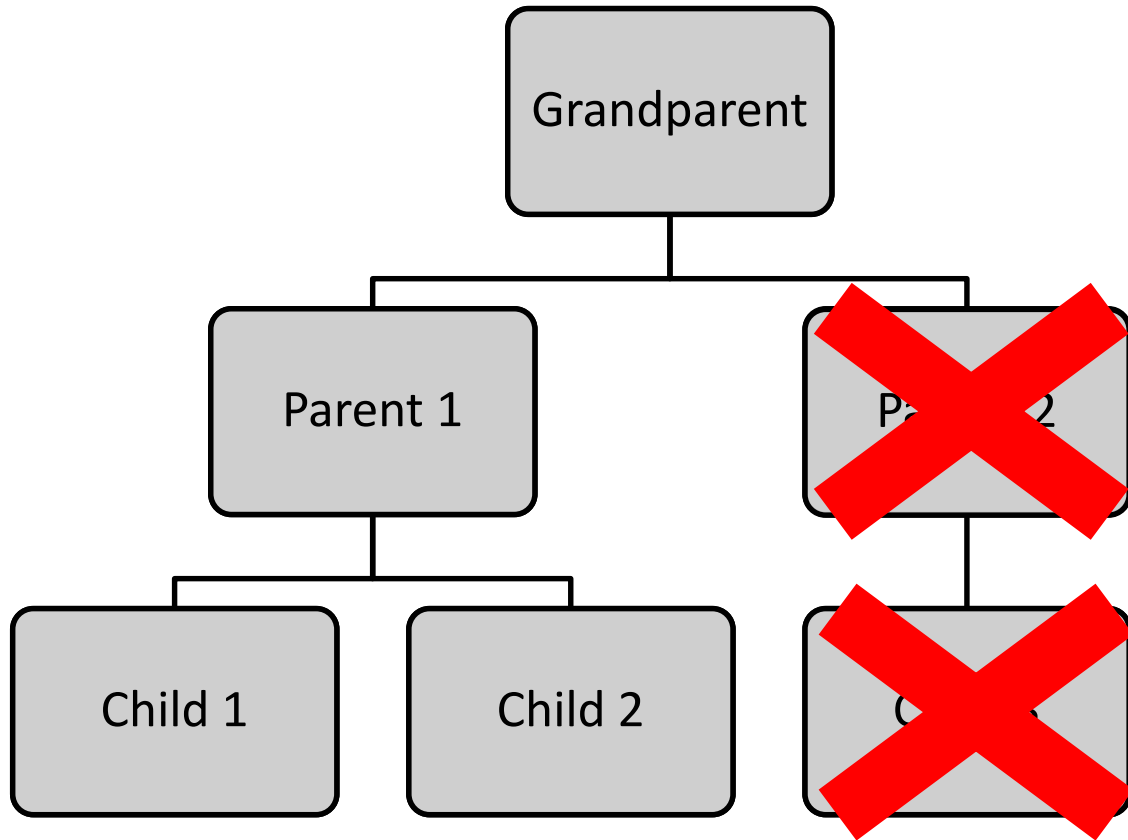- We can ensure that we obtain all appropriate parent/child nodes are included

# Multiple Slim Ontologies to Main HuBMAP/CCF Ontology

- We could specify for the main file we want nodes Parent 1 (from Ontology A) and Parent 2 (from Ontology B) and all their descendent nodes.
- We could specify we want nodes Children 1, 2, and 3 and all of their ancestor nodes.
- We could do a mix and still get the ontology.

# Multiple Slim Ontologies to Main HuBMAP/CCF Ontology

- When viewing the ontology, we can easily eliminate extraneous terms for UI/navigation (e.g. eliminate Parent 2)
- When in a specific part of the body, certain terms could just disappear (e.g. eliminate Parent 2 and descendants)
- With a graph library, very easy to remove nodes or branches.

# Pros/Cons of Multiple Slim Ontologies

## Pros

- Ensures that we have the correct biology in the base ontology.
  - Lots of work already done by using pre-existing ontologies
- Can have slim ontologies for each zoom level, organ, or system
- As the base ontologies continue to update, new information propagates in
- Easy to add additional ontologies through crosswalks
- Elimination of some of the hand editing of ontologies

## Cons

- Need to specify which versions of the input ontologies we are using
- Tracking of slim ontologies could become burdensome
- Lack of hand editing could make slim ontologies harder to use
- Hand editing could be replaced with code, but then additional coding effort is necessary

- Most of these cons are easily surmountable by using GitHub and using a triplestore repository like LungMAP did for its ontology

# Dimensions of CCF Ontology

- Overall anatomical location
  - Heart aorta

- Cell line/type location
  - Endothelial cells

- Chemical location
  - Areas for positive staining for H&E

- Real data will have to be mapped into/annotated with the ontologies

- Some of this data will need to come from the experimental protocols (possible semantic integration with protocols.io ?) and some will likely come from Machine Learning