# Ontologies and controlled vocabularies

# Why use ontologies and CVs?

- Very important in all data collection and analysis to manage and share large data sets
- To use the same data labels universally
- To enable quick retrieval of data
- To enable easy comparison of data
- To remove ambiguities

# Ambiguities in names

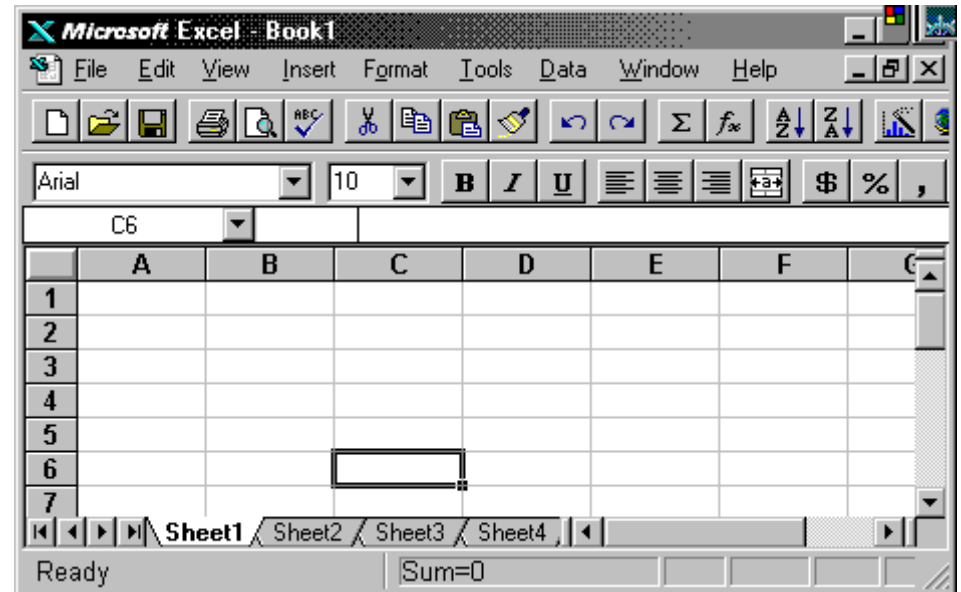- What is a cell?

# Ambiguities in names

- What is a cell?



OR
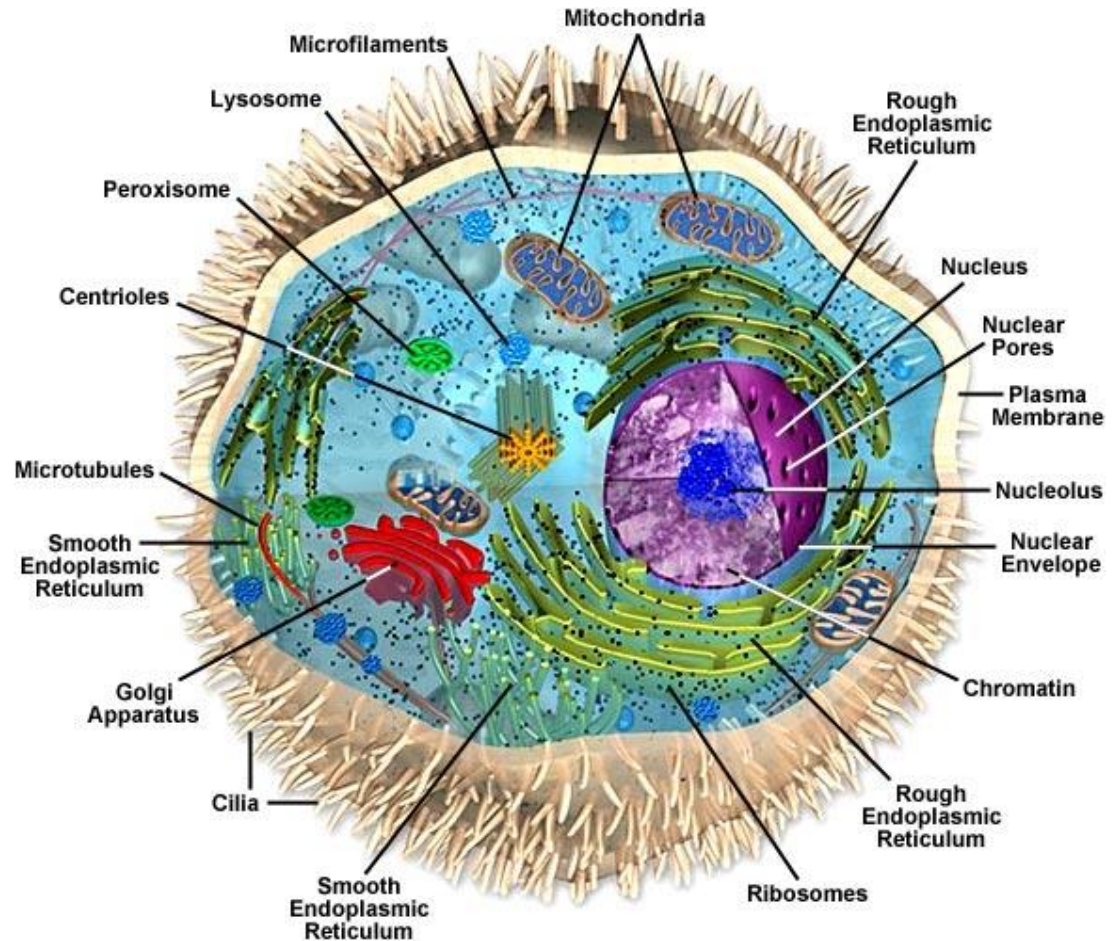
# Ambiguities in names

- What is a cell?

OR

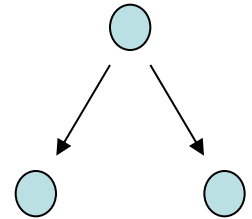# Ambiguities in names

• What is a cell?

# Ambiguities in naming continued

- The same **name** can be used to describe different **concepts**, e.g:
  - Glucose synthesis
  - Glucose biosynthesis
  - Glucose formation
  - Glucose anabolism
  - Gluconeogenesis
- All refer to the process of making glucose
- Makes it difficult to compare the information
- Solution: use **Ontologies** and **Data Standards**

# Ontologies

- An ontology is a formal specification of terms and relationships between them – widely used in biology and boinformatics (e.g. taxonomy)

- The relationships are important and represented as graphs

- Ontology terms should have definitions

- Ontologies are machine-readable

- They are needed for ordering and comparing large data sets

# Open Biomedical Ontologies

http://www.obofoundry.org/

- Central location for accessing well-structured controlled vocabularies and ontologies for use in the biological and medical sciences.

- Provides simple format for ontologies that can encode terms, relationships between terms and definitions of terms including those taken from external ontologies.

# Scope of Open Biomedical Ontologies

- Anatomy
- Animal natural history and life history
- Chemical
- Development
- Ethology
- Evidence codes
- Experimental conditions
- Genomic and proteomic
- Metabolomics
- OBO relationship types
- Phenotype
- Taxonomic classification

# Ontologies of use to H3Africa

- Phenotypes
  - Mammalian phenotype ontology
  - Human phenotype ontology
  - PhenoDB, PhenoTips and PhenomeCentral
  - PhenX Toolkit
  - Symptom ontology
  - Disease Ontology
  - OMIM, SNOMED
- Experiment
  - Experimental factor ontology

# Phenotypes

- …"observable morphological, physiological and behavioural characteristics of an individual in the context of the environment" (MP Ontology)

- Physicians, researchers and clinical laboratories need systems to enable standardized and structured phenotypic data to be collected from patients.

# Phenotype ontology

- http://www.human-phenotype-ontology.org/
- Developed starting from OMIM + literature
- Has 10,000 terms and 50,000 annotations to hereditary di

# Browsing the HPO

# PhenX toolkit

- Well-established, broadly validated measures of phenotypes and exposures relevant to investigators in human genomics, epidemiology, and biomedical research

- Provides detailed protocols, information about the measures

- Provides a toolkit for developing CRF

# PhenX toolkit

- **Domain** –field of research e.g. demographics, anthropometrics, organ systems, complex diseases
- **Collection** –group of measures e.g. population
- **Measure** –way of collecting data from participant
- **Protocol** –standard way to collect & record measure
- **Data Collection Worksheet** –identifies items in a protocol
- **Data Dictionary -**lists each variable included in a protocol with its attributes, including variable names and unique identifiers

# PhenX toolkit example

# Symptom ontology

- http://bioportal.bioontology.org/ontologies/SYMP and
http://symptomontologywiki.igs.umaryland.edu/wiki/index.php/Main_Page

- "A perceived change in function, sensation or appearance reported by a patient indicative of a disease"

# Disease Ontology

- http://do-wiki.nubic.northwestern.edu/do-wiki/index.php/Main_Page and http://www.disease-ontology.org/

- "A disease is a disposition (i) to undergo pathological processes that (ii) exists in an organism because of one or more disorders in that organism."

- Community-driven open source ontology for human diseases

- Link to SNOMED, ICD-9, ICD-10, MeSH, UMLS

# Disease ontology example

# Experimental factor ontology

- http://www.ebi.ac.uk/efo/

- Description of experimental variables

- Combines parts of other ontologies, e.g. anatomy, disease and chemical compounds

- Recommended by EGA

# Experimental ontology example 1

# Experimental ontology example 2

# H3Africa phenotype harmonization WG

- Leverage the significant investment being made in cohorts and genomic analyses

- Multi-site and cross-consortium analyses – more statistically powerful and informative

- Encourage H3Africa to use PhenX standardized phenotype measures for Case Report Forms

- Encouraged to collect a set of 25 essential and 10 discretionary phenotypes

# "Essential" phenotypes

(1) Age & (2) Sex
(3) Country of birth
(4) Current residence
(5) Native language
(6) Ethno-linguistic/tribal affiliation
(7) Country of birth of father and mother
(8) Native language of father and mother
(9) Ethno-linguistic/tribal affiliation of mother and father
(10) Height
(11) Weight
(12) Current medications
(13) Smoking history
(14) Alcohol history

Self-reported personal history (streamlined - yes/no and age at diagnosis):

(15) Hypercholesterolemia
(16) Hypertension
(17) Myocardial infarction
(18) Arrhythmia
(19) Rheumatic fever/rheumatic heart disease
(20) Asthma or reactive airway disease
(21) Stroke
(22) Diabetes
(23) Kidney disease
(24) HIV
(25) Tuberculosis

CBIO
Computational Biology @ UCT

UNIVERSITY OF CAPE TOWN

# "Discretionary" phenotypes

- (1) Schizophrenia
- (2) Cancer (especially cervical cancer)
- (3) Malaria
- (4) Trypanosomiasis
- (5) Substances of abuse history
- (6) Blood pressure measurement
- (7) Urinalysis (albumin, creatinine, protein)

more detailed questions on:

- (8) Diabetes - Type 1 and Type 2
- (9) Kidney disease
  (10) Rheumatic heart disease

Additional phenotypes-measurements being collected by at least 2 of the larger Collaborative Center grants

- Waist circumference
- Hip circumference
- Ultrasound fat measurements
- Family history of stroke, heart attack, hypertension

# Importance of using ontologies

- Share information in a common structure
- Enable reuse of domain knowledge
- Make naming/entities explicit and unambiguous
- Used for searching of data
- For H3Africa it will increase power to do cross-project studies