Ontologies to Unify **Genomics** and Phenomics Across **Species**

Anne E Thessen

annethessen@gmail.com









On behalf of:

Melissa Haendel Chris Mungall Chris Chute Damian Smedley **David Osumi-Sutherland** Peter Robinson **Tudor Groza** Sebastian Köhler Ada Hamosh Monica Muñoz-Torres Jules Jacobsen Kent Shefchek Nomi Harris Tim Putman Anne Thessen Harshad Hedae Justin Reese Kevin Schaper Lauren Chan Matt Brush Nicolas Matentzoglu Nicole Vasilevsky Sabrina Toro Seth Carbon Sierra Moxon Tudor Groza

Funded by: **NIH Office of Director:** 1R240D011883; NIH-UDP: HHSN268201300036C, HHSN268201400093P; NCINCI/Leidos #15X143, BD2K U54HG007990-S2 (Haussler) & BD2K PA-15-144-U01 (Kesselman)

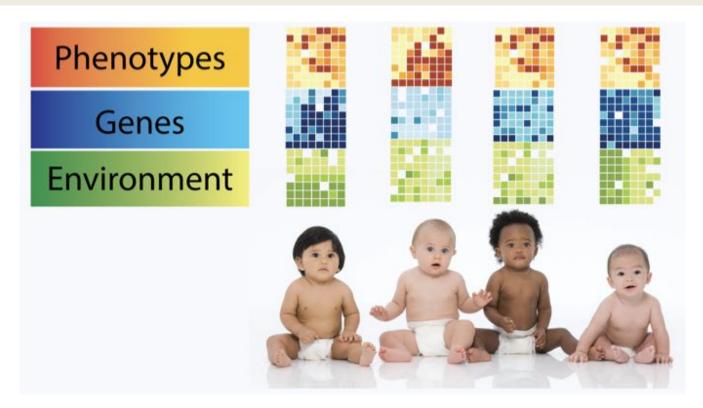
These slides are at bit.ly/NAScompanion

The genome is sequenced, but....



...we still don't know very much about what it does

The genome is sequenced, but....



...we don't understand environmental effects

1a) Individual person as biological subject

1b) Individual person as patient

Phenotypes

Genes

Environment

4) Many individuals in the EHR



2) A disease as reference knowledge (derived from many patient's presentation)

3) Cross species comparison for diagnosis

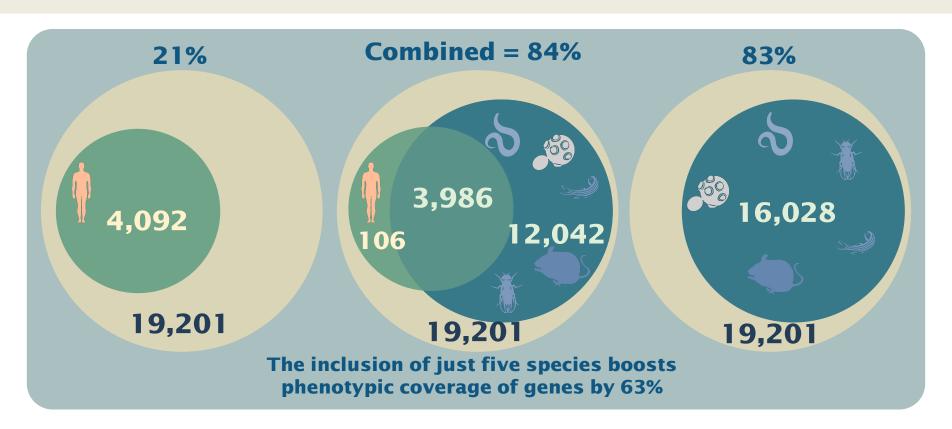






More species, more coverage





Other species aren't just relevant, each has unique phenotypes



The dog's retina has area centralis (analogous to the human macula) & fovea-like region, similar to humans; useful to study naturally occurring cone diseases



• Aged cats are natural models of Alzheimer's Disease: they form Abeta oligomers, neurofibrillary tangles, and have neuronal loss



Naked Mole Rats don't get cancer



Armadillos are a natural host of *M. leprae*, the mycobacterium that causes leprosy (only one besides humans)



 Tree shrews' glioblastomas are morphologically & genetically similar to humans (& more similar than mouse models)



Great pond snails are models of inflammation-mediated memory dysfunction, and show evidence of spontaneous neural tissue regeneration after injury



Silkworms are a model for uric acid metabolism. Decreases in plasma uric acid are correlated with clinical progression of Parkinson's Disease

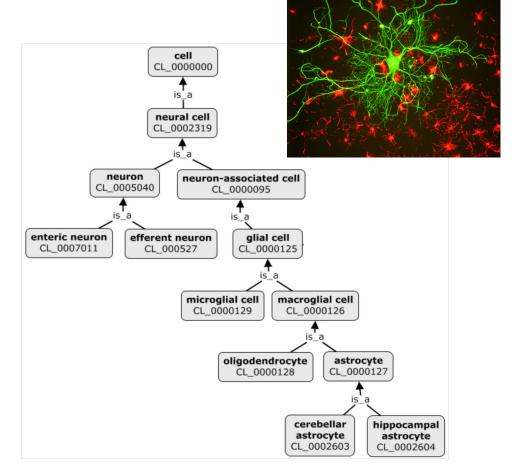
What is an ontology?

DEFINITION:

A formal, computational representation of knowledge in a particular domain.

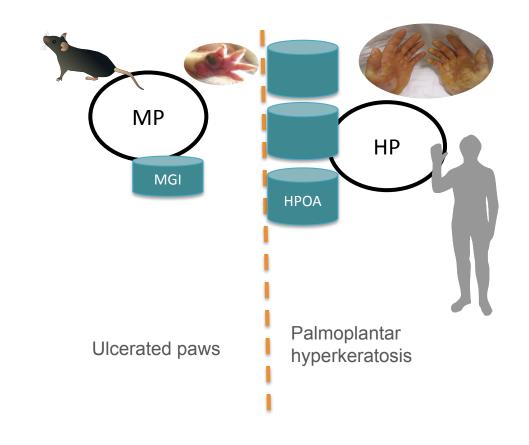
KEY FEATURES:

- Terms are defined
- Semantics relationships between terms are defined, allowing logical inference and sophisticated data queries
- Terms are arranged in a hierarchy
- Expressed in a knowledge representation language such as RDFS, OBO, or OWL

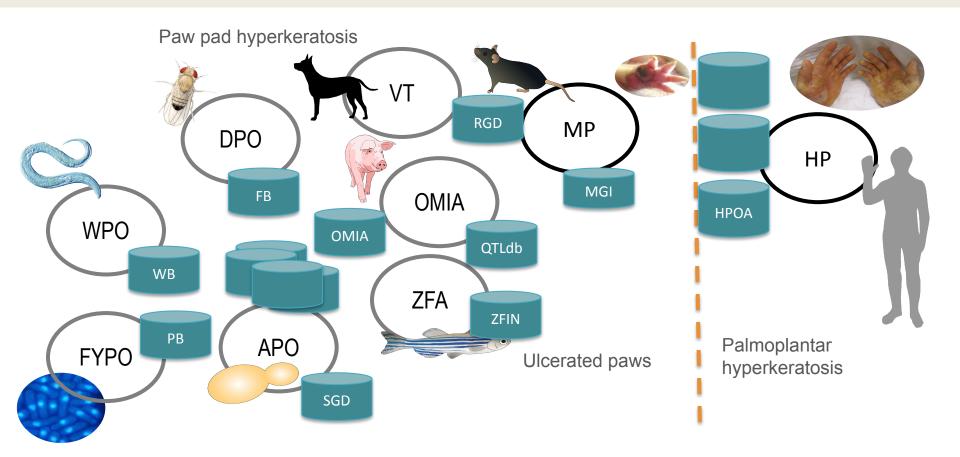


bit.ly/ontology101

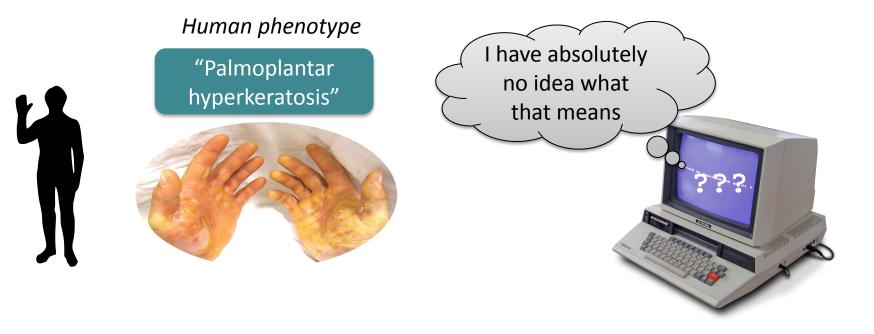
Challenge: Each data source uses its own vocabulary



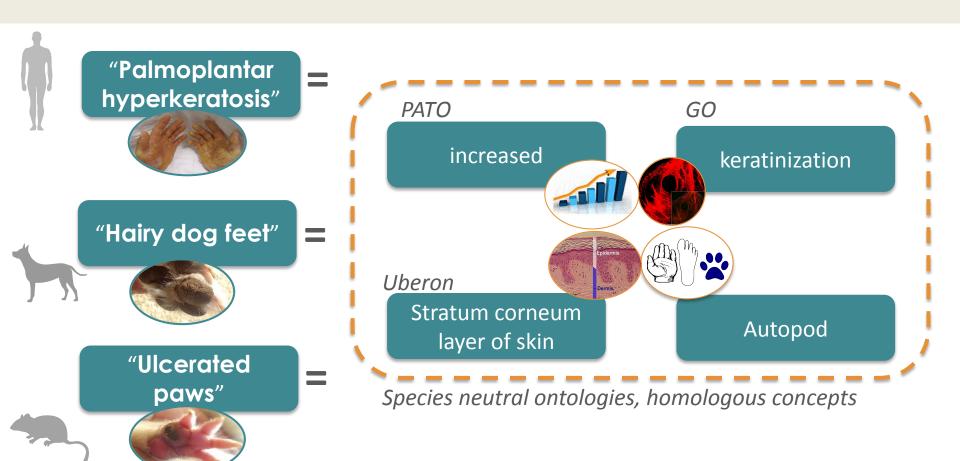
Challenge: Each data source uses its own vocabulary



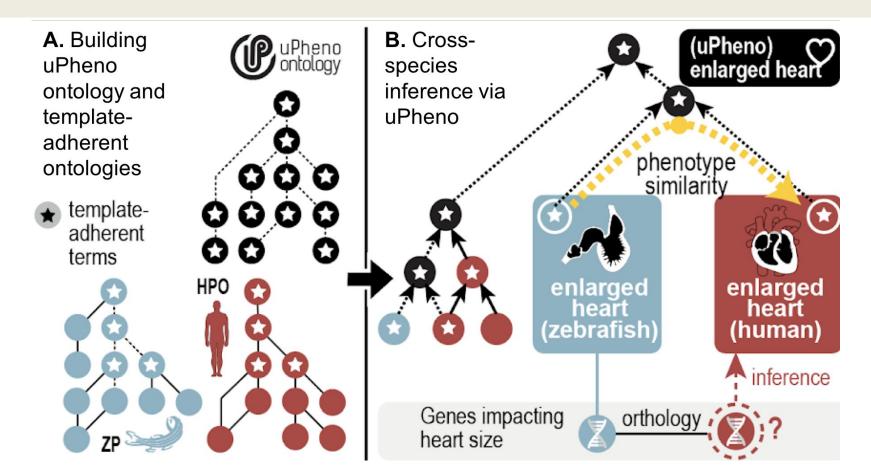
Can we help machines understand phenotypes?



Decomposition of complex concepts allows interoperability

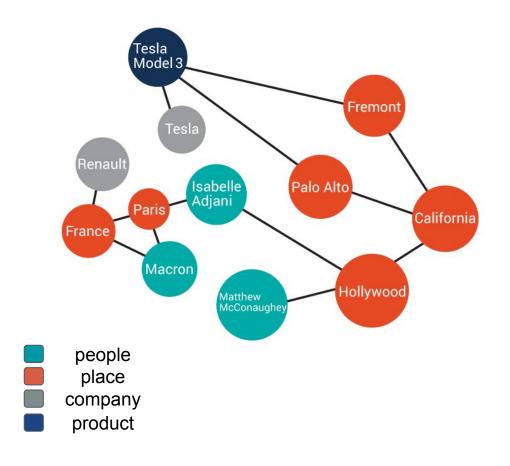


Why do we care about phenotype integration?



Ontology or Knowledge Graph?

- Two different ways to represent the same data
- Require slightly different modeling techniques
- Ontologies are community-driven schema that moves very slowly
- Knowledge graphs grow more quickly and can be more responsive
- More tools and services available for knowledge graphs

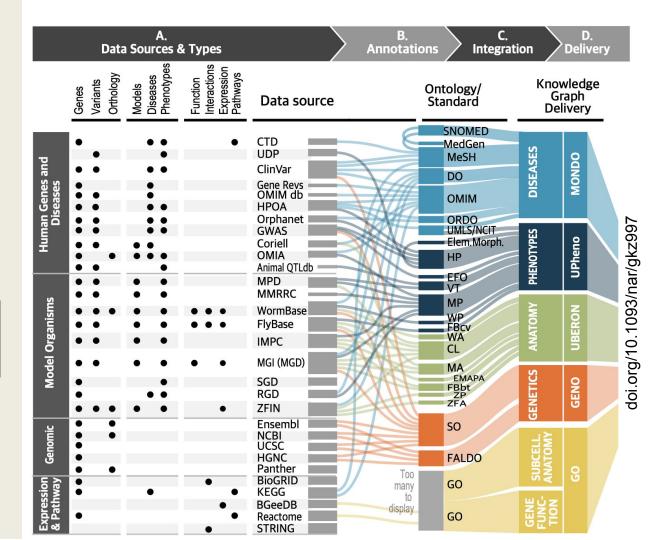


How did we do this?

- monarchinitiative.org
- Integrator of cross species genotype-phenotype data
- OWL and DOS-DP
- Uses OBO Foundry ontologies
- obofoundry.org







Use Case: Rare disease diagnostics

- Fuzzy matching
- Collections of phenotypes and genes

Legend



Perfect Match



Fuzzy Match



No Match

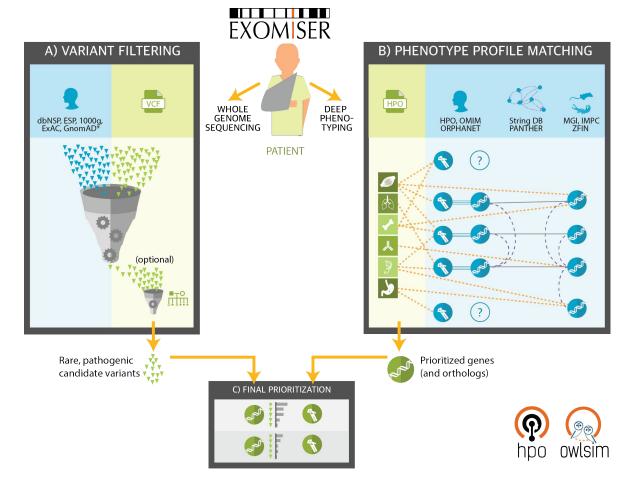






Use Case: Variant prioritization for clinical diagnostics

- Phenotype profile includes abnormally low bone mineral density
- Find all human orthologs
 of genes for which the
 presence of variants
 (alleles) is correlated with
 a phenotypic effect on
 bone mineral density.

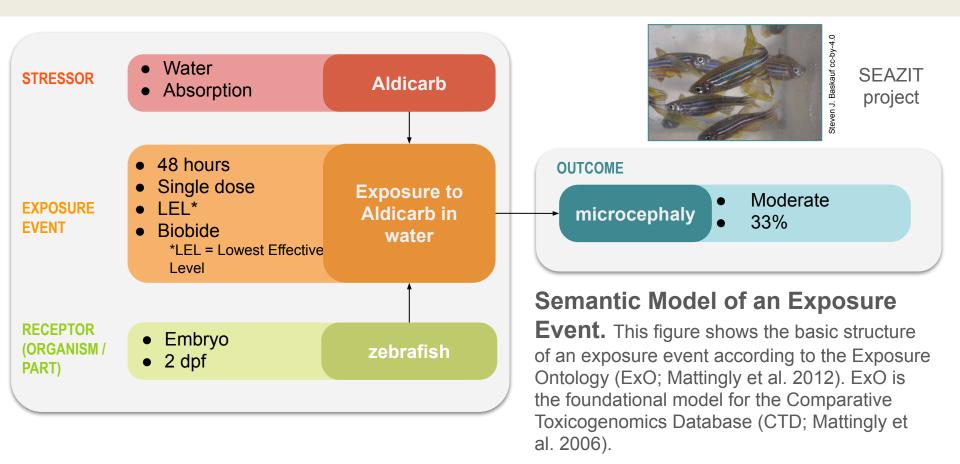




Exposures: The missing piece

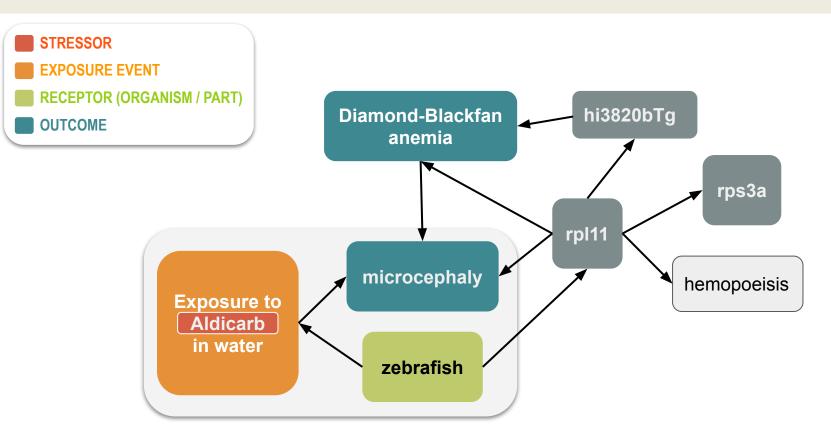


Exposure Ontology make environmental data computable



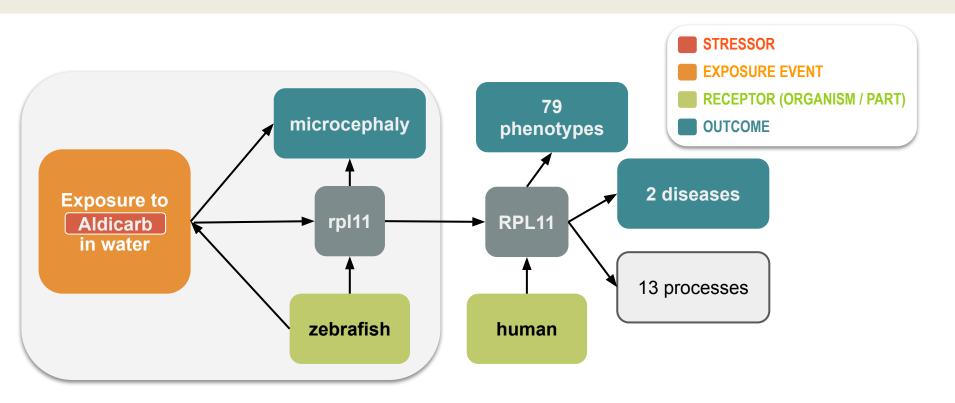
Integration with the Monarch Knowledge Graph





Integration with the Monarch Knowledge Graph

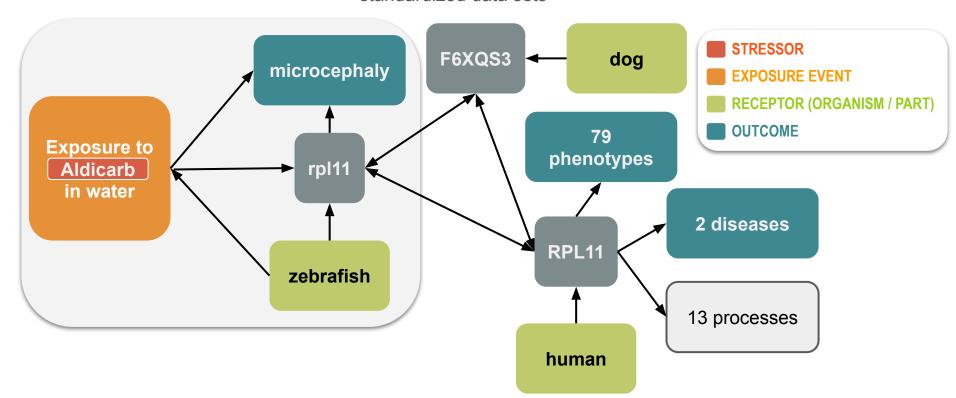




Integration with the Monarch Knowledge Graph



Not many connections for dog - need standardized data sets



We have a machine-readable language for describing some exposures



CHEBI:6651

$$=$$
 CCOC($=$ 0)CC(SP($=$ S)(OC)OC)C($=$ 0)OCC

CheBI is a chemical ontology

But others are harder to define



Image: Zol87 CC by/nc

Summary

- Ontologies and knowledge graphs are an effective tool for integrating biological data across species and terminologies
- Integrated data sets can be used for inference and semantic similarity analysis
- Data standards are required (Dog Phenotype Ontology)
- Future work will include development of environmental exposure data models
- Companion animals share many human exposures can add power to studies

Want to help?

- annethessen@gmail.com
- Join the team (we and our partners are hiring trainees and staff)
 - 21772 Instructor Senior Instructor Scientific / Medical Writer
 - 23591 Postdoctoral Fellow
 - 21903 Consortium Program Assistant
 - Professional Research Assistant
 - Postdoctoral Fellow (Veterinary)
- Give us feedback on the Monarch tools and data. (T shirts and swag for user experience interviews) Inquire at info@tislab.org
- Promote responsible licensing of data reusabledata.org

The Monarch Initiative

www.monarchinitiative.org

annethessen@gmail.com



PDs: Melissa Haendel, Chris Mungall, Peter Robinson

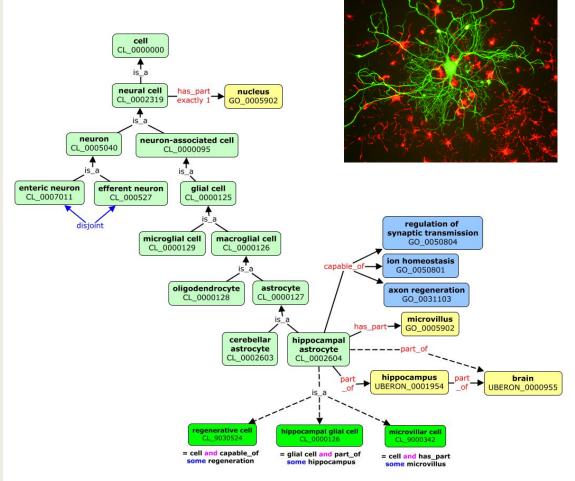
What is an ontology?

DEFINITION:

A formal, computational representation of knowledge in a particular domain.

KEY FEATURES:

- Terms are defined
- Semantics relationships between terms are defined, allowing logical inference and sophisticated data queries
- Terms are arranged in a hierarchy
- Expressed in a knowledge representation language such as RDFS, OBO, or OWL



bit.ly/ontology101

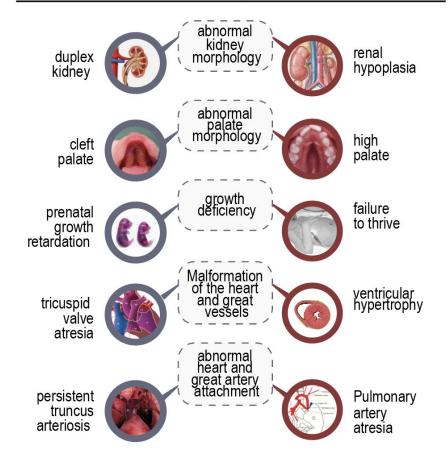
Ontologies as a tool for data integration

- Across species
- MAYBE REMOVE



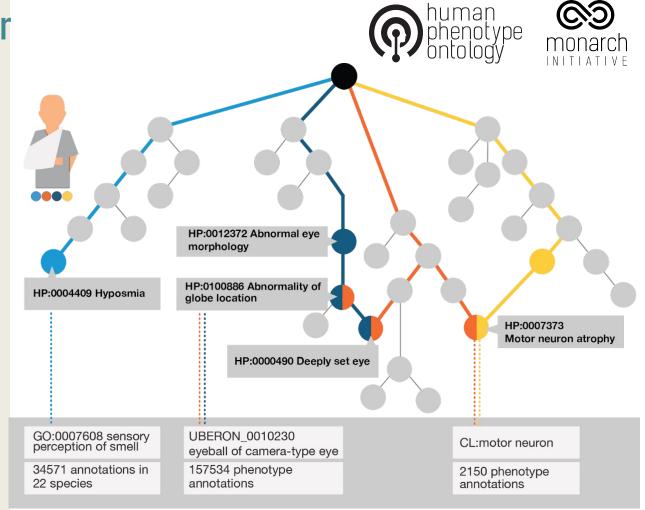






Ontologies as a tool for data integration

- Across species
- Granularity mismatch
- MAYBE REMOVE
 BECAUSE WE HAVE 8

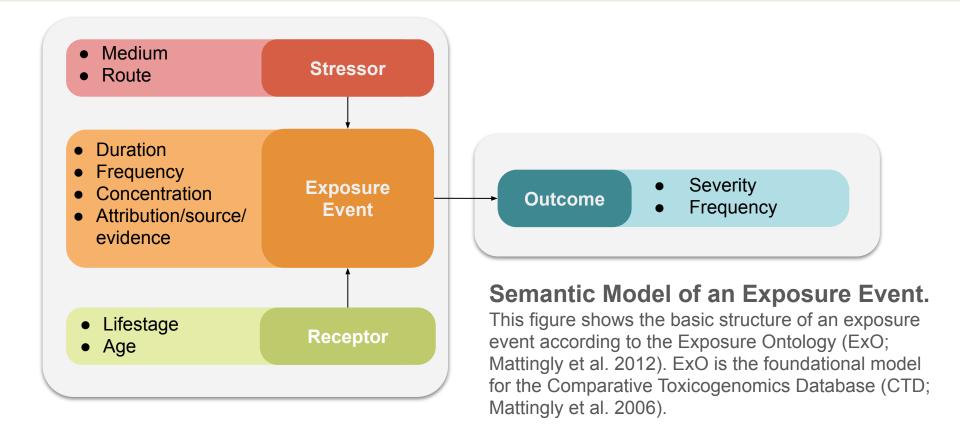


Exposures: The missing piece

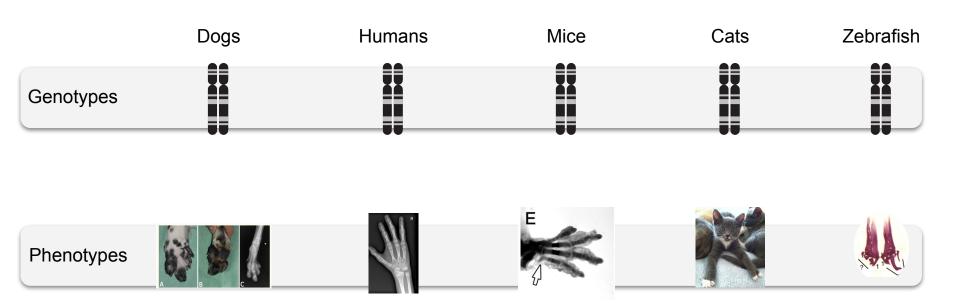




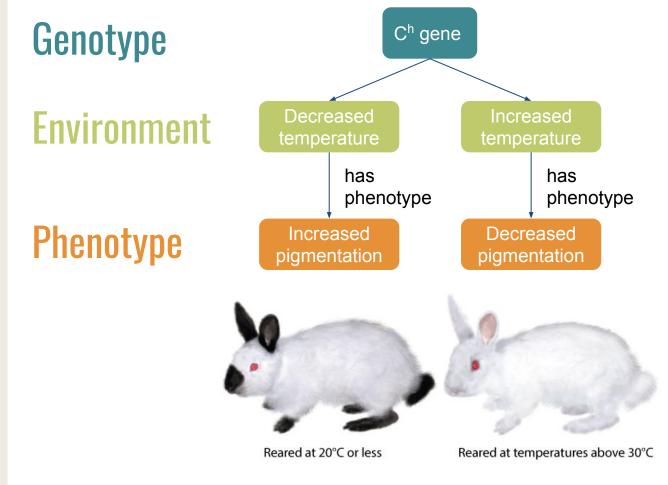
Environmental Exposures Modeling with ExO



Why do we care about phenotype integration?



Modeling Genotype + Environment = Phenotype



Lobo, I. (2008) Environmental influences on gene expression. Nature Education 1(1):39

- 1. What are the strategies for standardizing, sharing, and aggregating health records and relevant metadata across species? What are the best practices for collection, storage, and analysis of biosamples to assess exposures (e.g., biorepository resources, DNA susceptibility, DNA methylation, microbiome, etc.)?
- 2. What are the obstacles (e.g., scientific, infrastructure, ethical, and financial) to using companion animals as biomonitors, and what are potential solutions?
- 3. What are the best practices for data integration of human and companion animal data? Have you identified best practices or lessons learned in efforts to link datasets?
- 4. Are there opportunities to promote data sharing and collaboration to advance research on the role of companion animals as sentinels for predicting environmental exposure effects? If so, what are they?5. What incentives are needed to encourage data sharing and collaboration? How can disincentives for data sharing be
- 5. What incentives are needed to encourage data sharing and collaboration? How can disincentives for data sharing be overcome and the use of incentives become more standard? What infrastructure and resources are needed? Which groups need to be engaged?
- 6. What are the implications for expanded, systematic collection of this data? Could this help to identify environmental hazards and provide an early warning system for public health interventions?