



United States  
Environmental Protection  
Agency

# Bioinformatics is for Everyone: Applications To Challenges in Ecotoxicology

Carlie A. LaLone, Ph.D.  
Research Bioinformaticist



Office of Research and Development  
Center for Computational Toxicology and Exposure, Great Lakes Toxicology and Ecology Division

*The views expressed in this presentation are those of the authors  
and do not necessarily reflect the views or policies of the US EPA*

May 5<sup>th</sup>, 2022




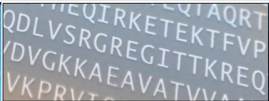



# Bioinformatics

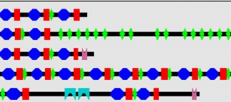
- Combines mathematics, information science, and biology to answer biological questions
- Developing methodology and analysis tools to explore large volumes of biological data
  - Query, extract, store, organize, systematize, annotate, visualize, mine, and interpret complex data
    - Usually pertains to DNA and amino acid sequences

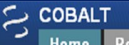
**Let the computers do the work**

**Taxonomy**  
The Taxonomy Database is a curated classification and nomenclature for all of the organisms in the public sequence databases. This currently represents about 10% of the described species of life on the planet.

**Protein**  
The Protein database is a collection of sequences from several sources, including translations from annotated coding regions in GenBank, RefSeq and TPA, as well as records from SwissProt, PIR, PRF, and PDB. Protein sequences are the fundamental determinants of biological structure and function.

**BLAST®**  
Basic Local Alignment Search Tool  
[Home](#) [Recent Results](#) [Saved Strategies](#) [Help](#)

**CDD**  
The Conserved Domain Database is a resource for the annotation of functional units in proteins. Its collection of domain models includes a set curated by NCBI, which utilizes 3D structure to provide insights into sequence/structure/function relationships.

**COBALT**  
Constraint-based Multiple Alignment Tool  
[Home](#) [Recent Results](#) [Help](#)

# AlphaFold Protein Structure Database

Developed by DeepMind and EMBL-EBI

Search for protein, gene, UniProt accession or organism

Examples: [Free fatty acid receptor 2](#) [A1g58602](#) [Q5VSL9](#) [E. coli](#) [Help: AlphaFold DB search help](#)

Feedback on structure: [Contact DeepMind](#)

## DeepFRI Dashboard

DeepFRI is a structure-based protein function prediction



# COFACTOR

Structure-based function predictions

[Enzyme Commission](#) [Gene Ontology](#) [Ligand Binding Site](#)



# I-TASSER

Protein Structure & Function Predictions

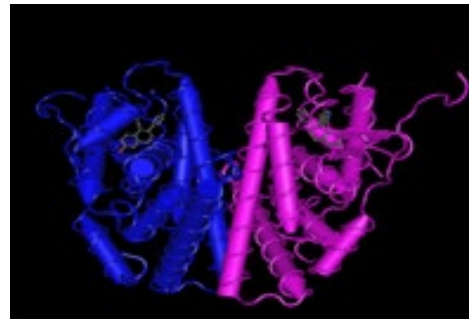
(The server completed predictions for [682524 proteins](#) submitted by [165524 users](#) from [159 countries](#))  
(The [template library](#) was updated on [2022/04/18](#))

# Sequence

MTMTLHTKASGMALLHQIQGNELEPLNRPQLKIPLERPLGE  
VYLDSSKPAVYNYPEGAAYEFNAAAAANAQVYGQTGLPYG  
PGSEAAAFGSNGLGGFPPLNSVSPSPLMLLHPPQLSPFLQ  
PHGQQVPYYLENESGYTVREAGPPAFYRPNSDNRRQGGR  
ERLASTNDKGSMAKESAKETRYCAVCNDYASGYHYGVWSC  
EGCKAFFKRSIQGHNDYMCNATNQCTIDKNRRKSCQACRLR  
KCYEVGMMKGGIRKDRRGGRMLKHKRQRDDGEGRGEVG  
SAGDMRAANLWPSPLMIKRSKNSLALSLTADQMVSAALLA  
EPPILYSEYDTPRPFSEASMMGLLTNLADRELHMINWAKV  
PGFVDLTLDQVHLLECAWLEILMIGLVWRSMHPGKLLFA  
PNLLDRNQGKCEGMVEIFDMLLATSSRFMMNLQGEEF  
VCLKSILLNSGVYFLSSTLKSLEEKDHIHRVLDKITDTLIHLM



# Structure

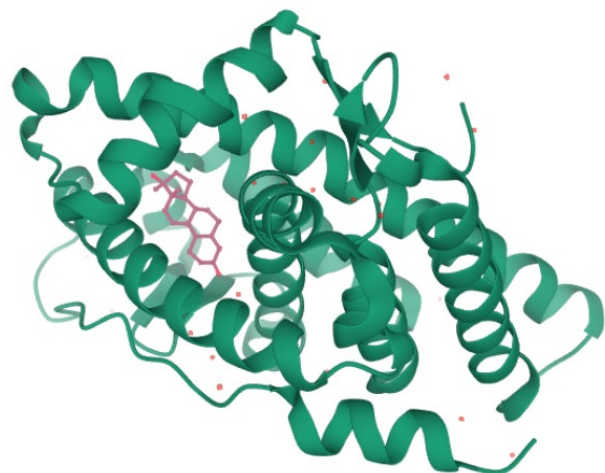


# Function



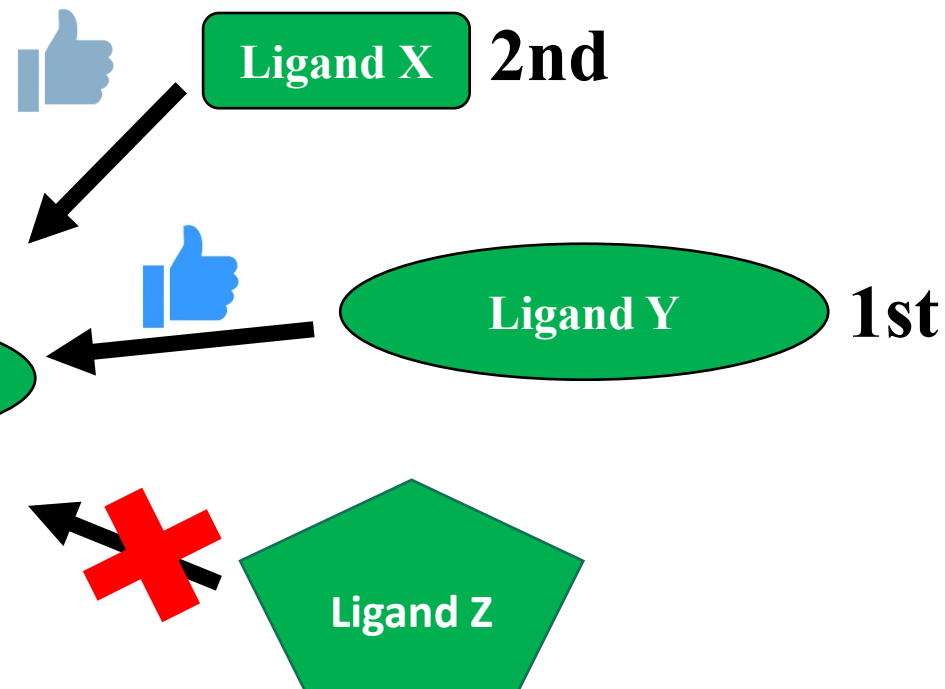
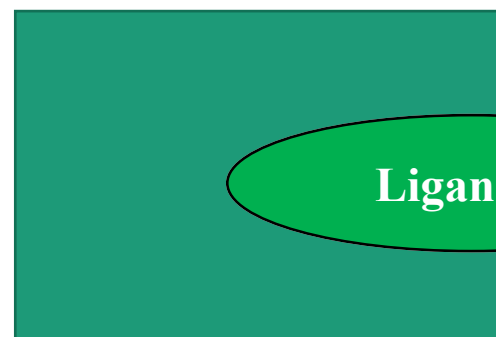
# Bioinformatics

# Advances in Drug Discovery/Development



Structure derived  
from X-ray  
crystallography

Human  
Protein Structure



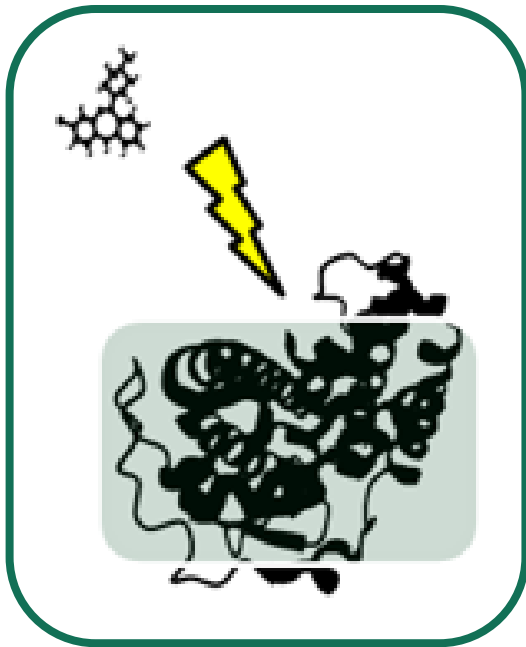
## Bioinformatics Toolbox:

- Molecular modeling
- Molecular docking
- Virtual screening
- Molecular dynamic simulations
- Functional prediction

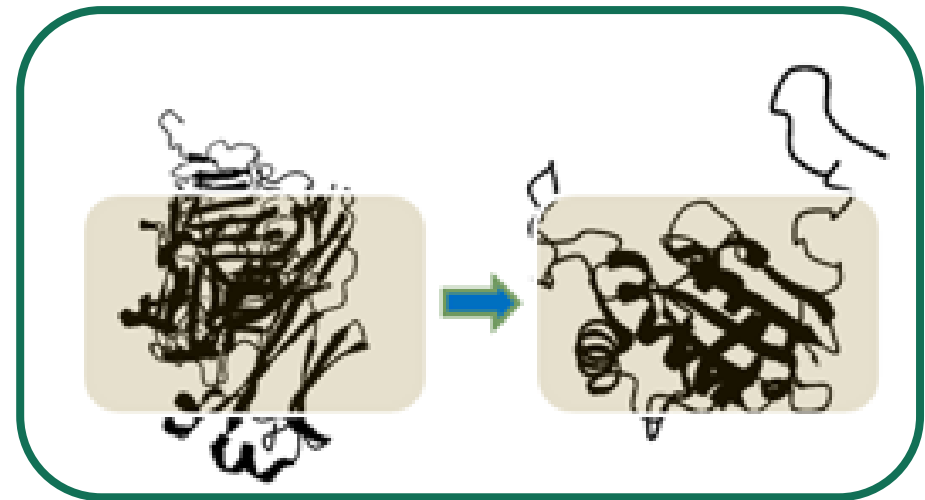


# Numerous bioinformatics approaches

- My interest and focus has been chemical-protein and protein-protein Interactions
  - Utility for cross species extrapolation



Chemical-Protein



Protein-Protein

# Overview Questions

- What species do we rely on for toxicity testing and why?
- Why consider predictive and computational approaches?
- How can bioinformatics help for chemical safety evaluations?
- What tools are available now and moving forward?
- How do we incorporate bioinformatics in decision making?





Chemicals make up the world around us – necessary for our modern society





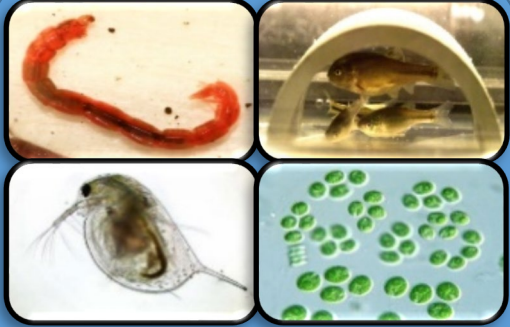


# Toxicity Testing to Understand Chemical Safety

- US EPA Examples:
- *Clean Air Act*
- *Clean Water Act*
- *Resource Recovery Act*
- *Endangered Species Act*
- *Food Quality Protection Act*
- *Endocrine Disruptor Screening Program*
- *Federal Insecticide, Fungicide, and Rodenticide Act*
- *Frank R. Lautenberg Chemical Safety for the 21<sup>st</sup> Century Act*
- *Comprehensive Environmental Response, Compensation, and Liability Act*
- *Guidelines for Deriving Numerical National Water Quality Criteria for the Protection of Aquatic Organisms and Their Uses*

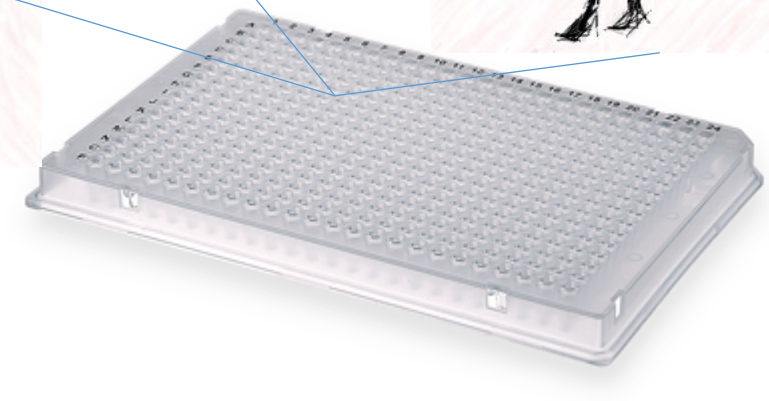
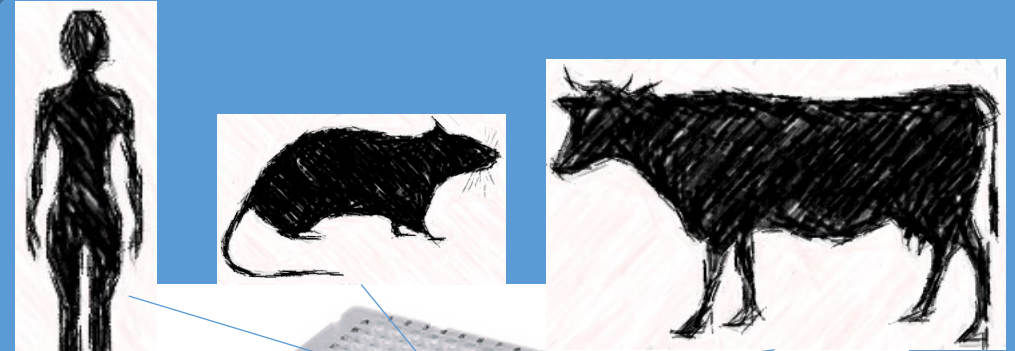


# Need for Advances in Species Extrapolation



High throughput  
transcriptomics

Historic whole organism  
toxicity testing



High-  
throughput  
screening  
assays  
(ToxCast)

Define the taxonomic domain of applicability in AOP development



Use of model organisms as surrogates representing the diversity of species in the environment



**cheap and readily available**



**easy maintenance and good breeding capabilities**

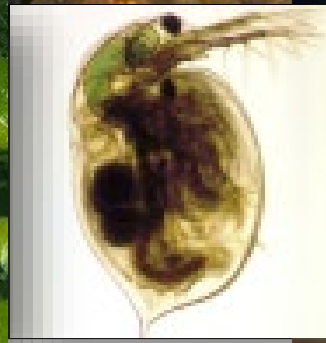


**short lifespans and rapid life cycles**



**ability to control diet and surroundings**

**requires least space and time-consuming care**





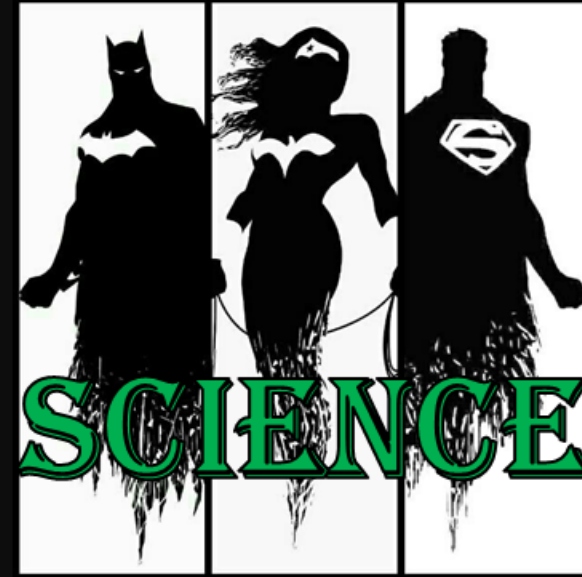
# Species Extrapolation

## What is it?

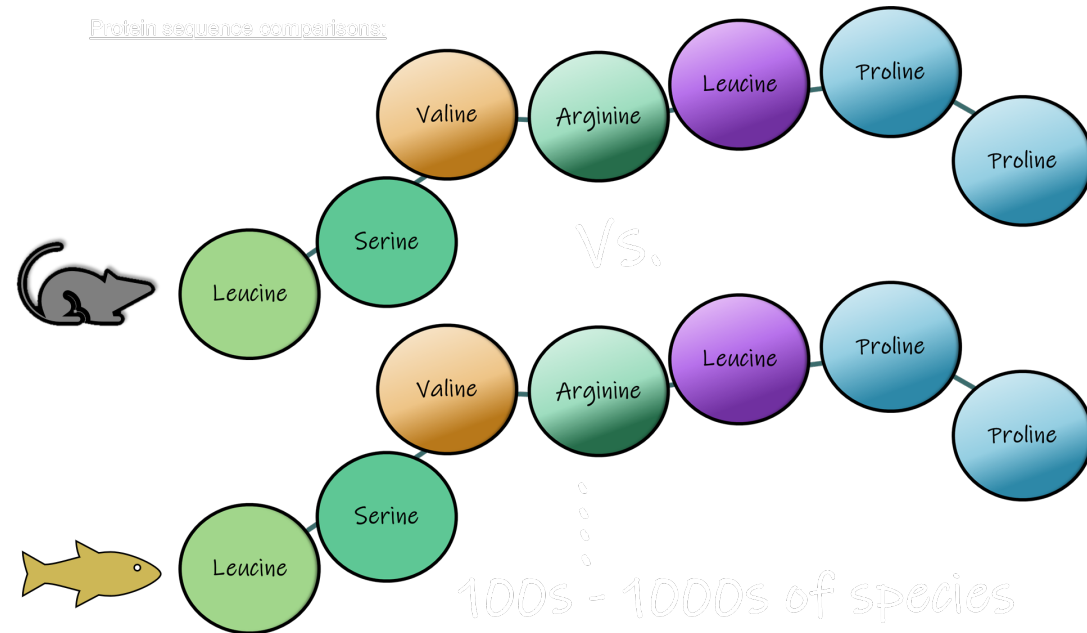
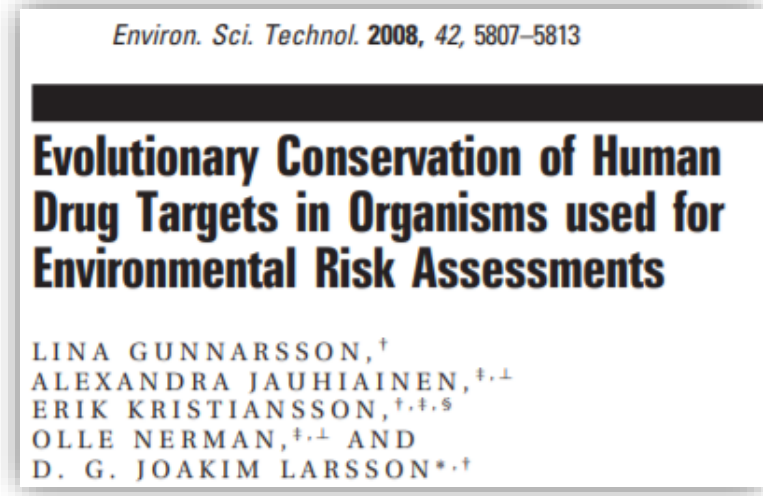
- Using existing **knowledge about one species** to estimate, predict, project, or infer the effect, impact, or **trajectory of another species**
  - For chemical safety typically dealing with toxicity

## Why is it important:

- **Limited or no toxicological data** for the animal or plant species of interest – reliance on surrogate (model organisms)
  - **Impractical to generate new data** for all species
- Testing **resources are limited**
  - International interest to **reduce animal use**
  - Ever-increasing demand to **evaluate more chemicals in a timely** and sometimes expedited manner
- **Sensitivity of species must be estimated** based on scientifically-sound methods of cross-species extrapolation
  - Immense **diversity of species** in the wild
  - Important challenge for species listed under the **Endangered Species Act**



# Bioinformatics for species extrapolation?



- Begin Simple and Advance as the Science Advances
- Consider sequence and structural attributes to understand protein conservation across species



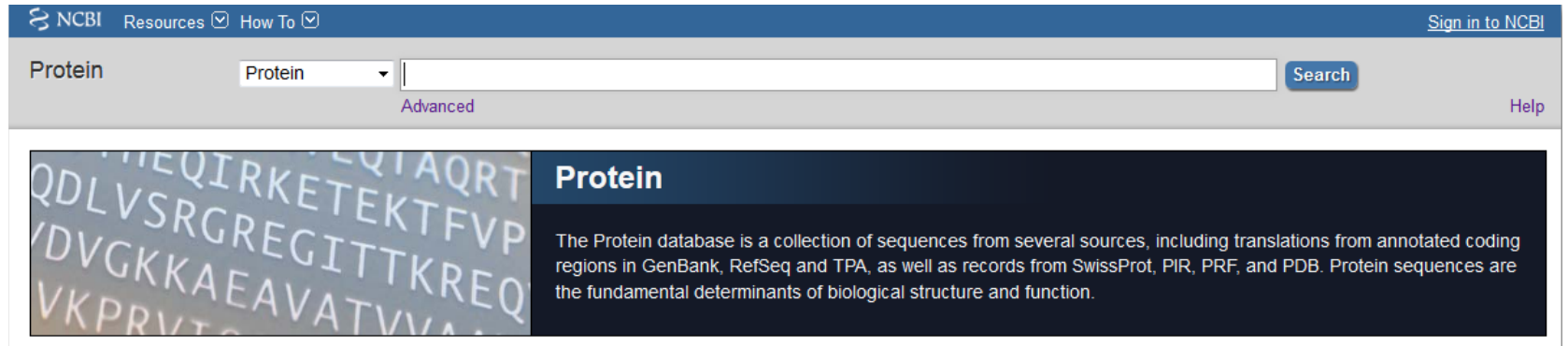

# Where could we begin in understanding species similarities and differences?

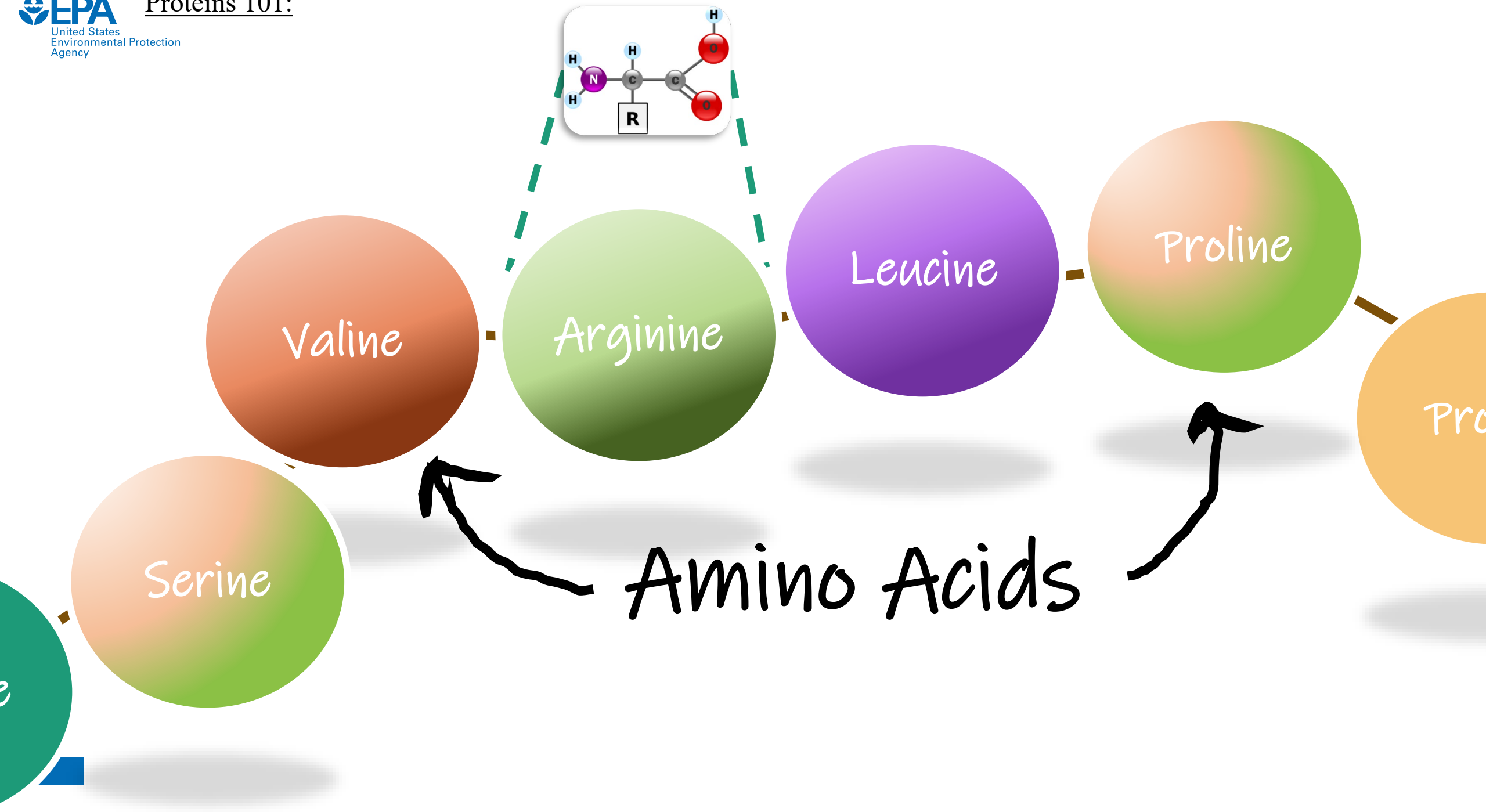
Look for existing, expanding data that does not require the destruction of live organisms

**Sequence and structural data: New tools and technologies have emerged**

- Improved sequencing technologies
- Large databases of sequence data

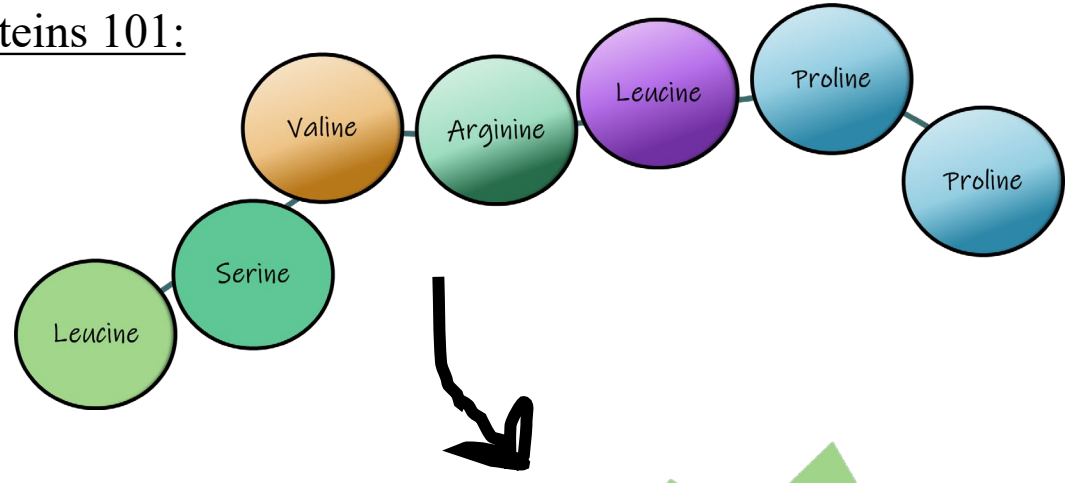
**NCBI: 224,211,842 Proteins representing 117,030 Organisms**



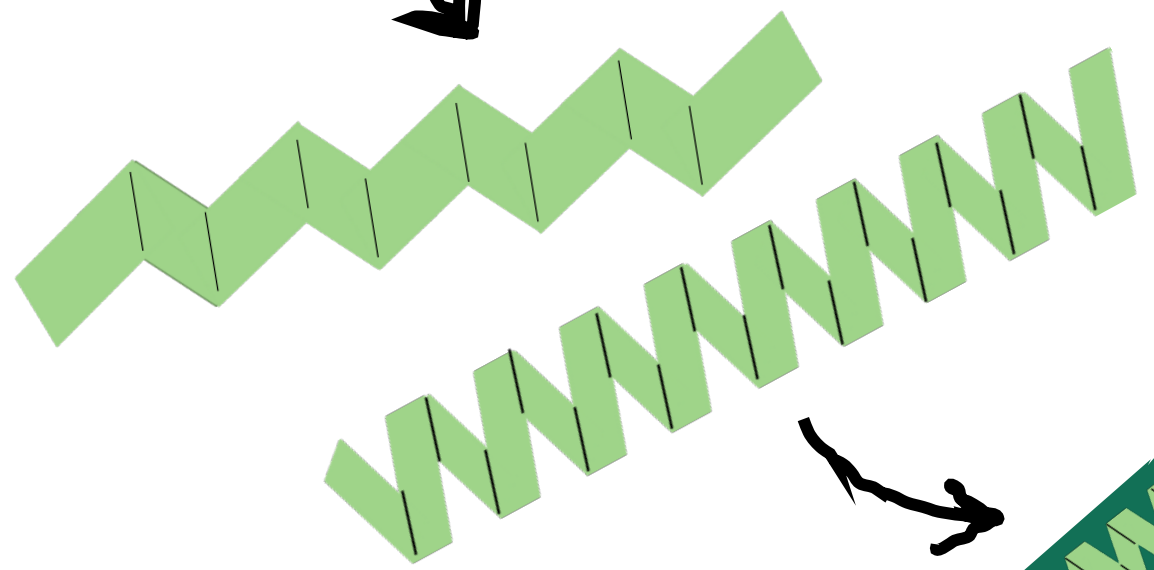




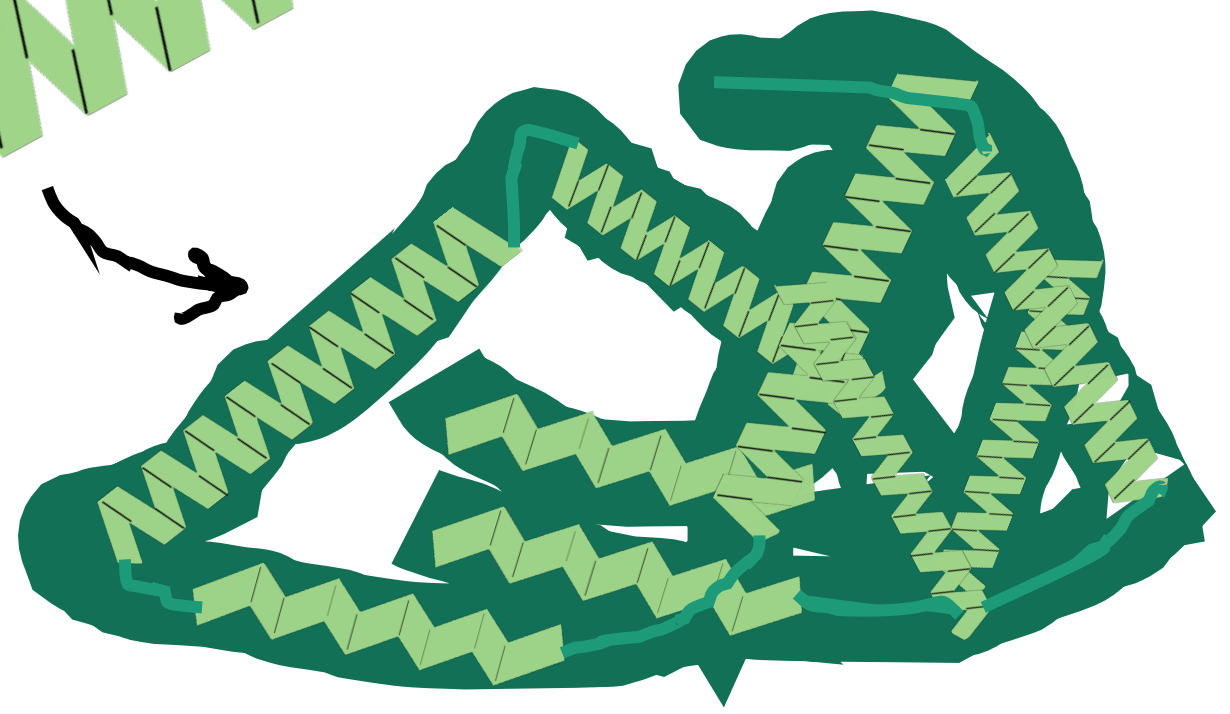
# Proteins 101:



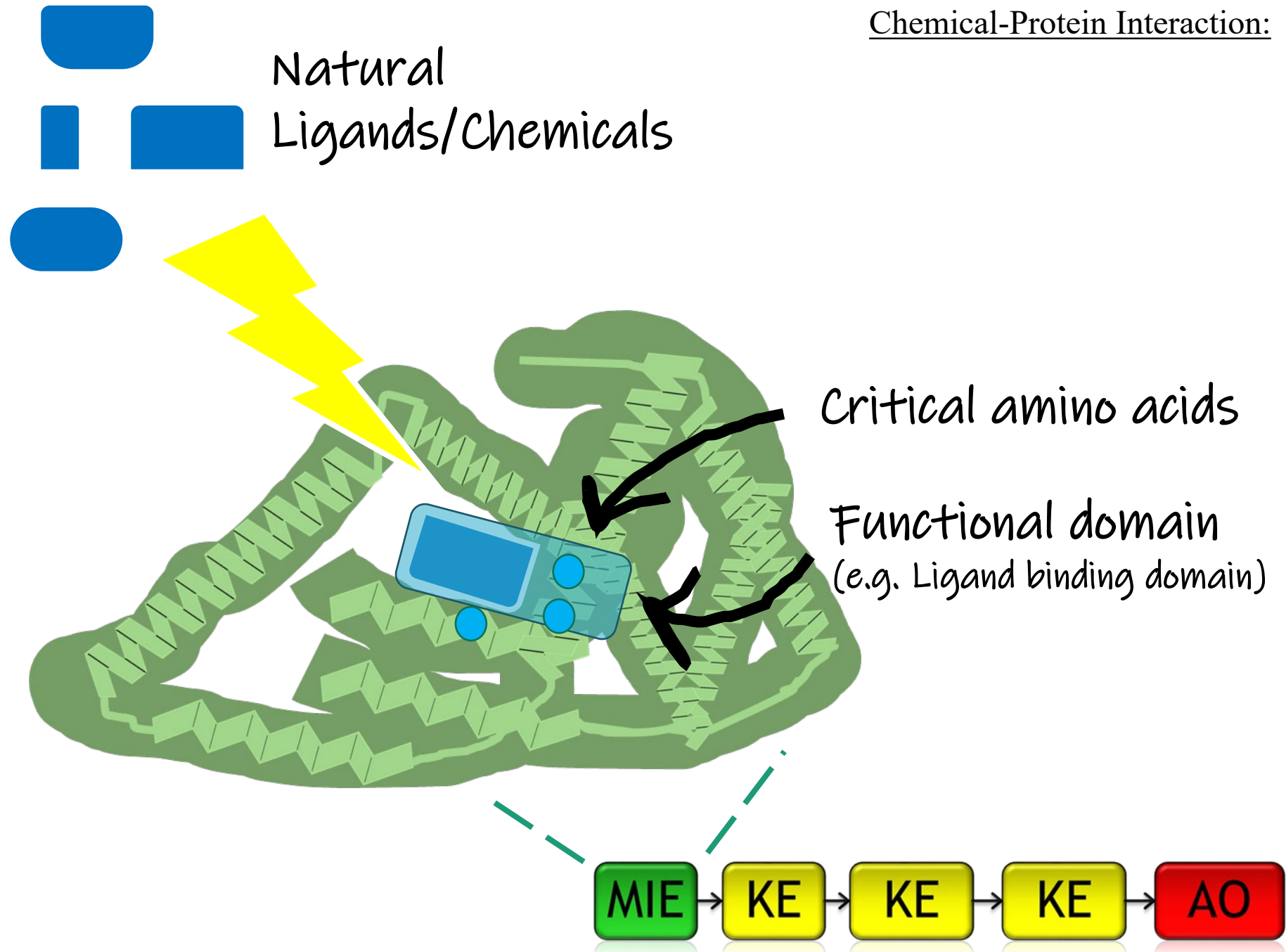
Primary amino acid sequence



Secondary Structure



Tertiary Structure



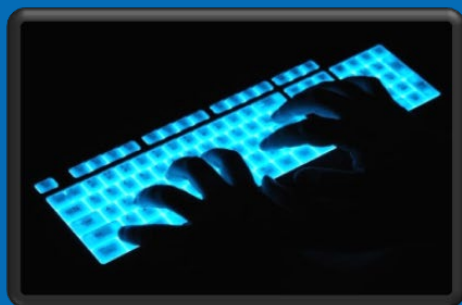
Similarity across species at the molecular level





<https://seqapass.epa.gov/seqapass/>

# Sequence Alignment to Predict Across Species Susceptibility (SeqAPASS)

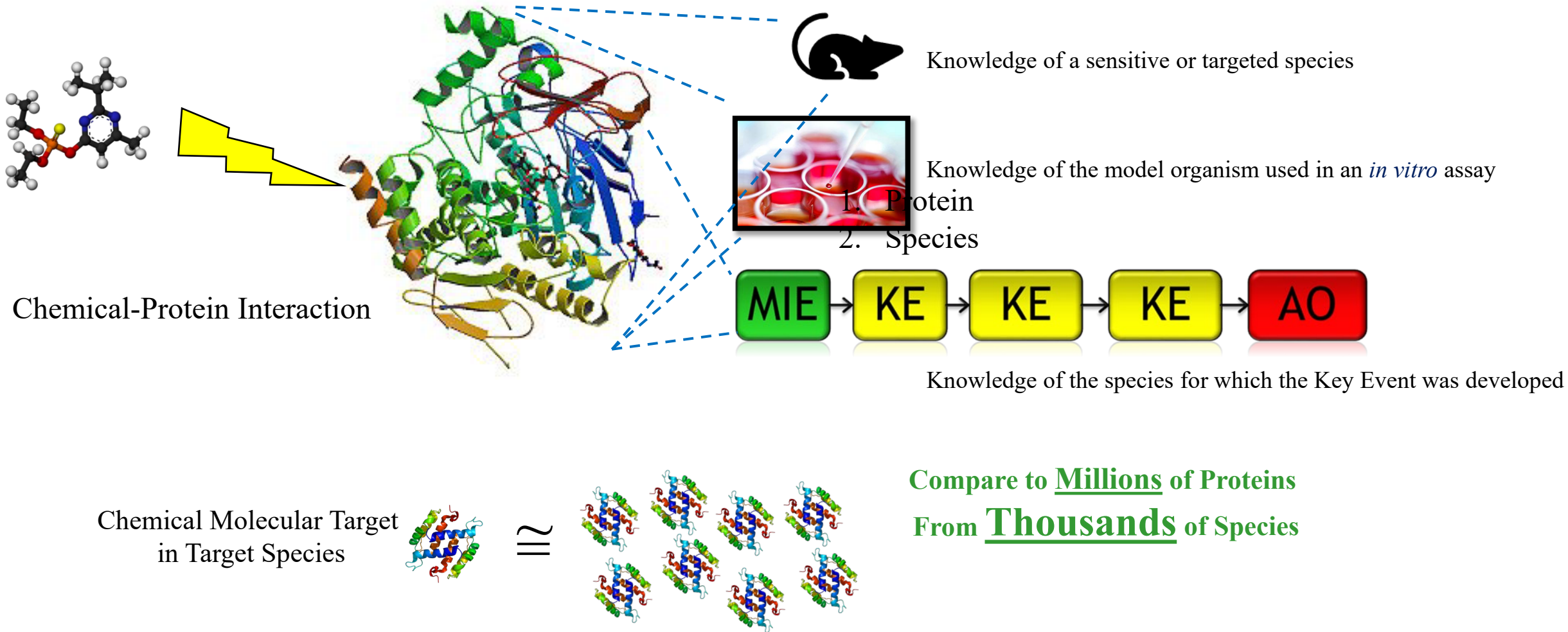


## Sequence Alignment to Predict Across Species Susceptibility (SeqAPASS): A Web-Based Tool for Addressing the Challenges of Cross-Species Extrapolation of Chemical Toxicity

Charlie A. LaLone,<sup>\*,1</sup> Daniel L. Villeneuve,<sup>\*</sup> David Lyons,<sup>†</sup> Henry W. Helgen,<sup>‡</sup>  
Serina L. Robinson,<sup>§,2</sup> Joseph A. Swintek,<sup>¶</sup> Travis W. Saari,<sup>\*</sup> and  
Gerald T. Ankley<sup>\*</sup>



# What information is required for a SeqAPASS query?



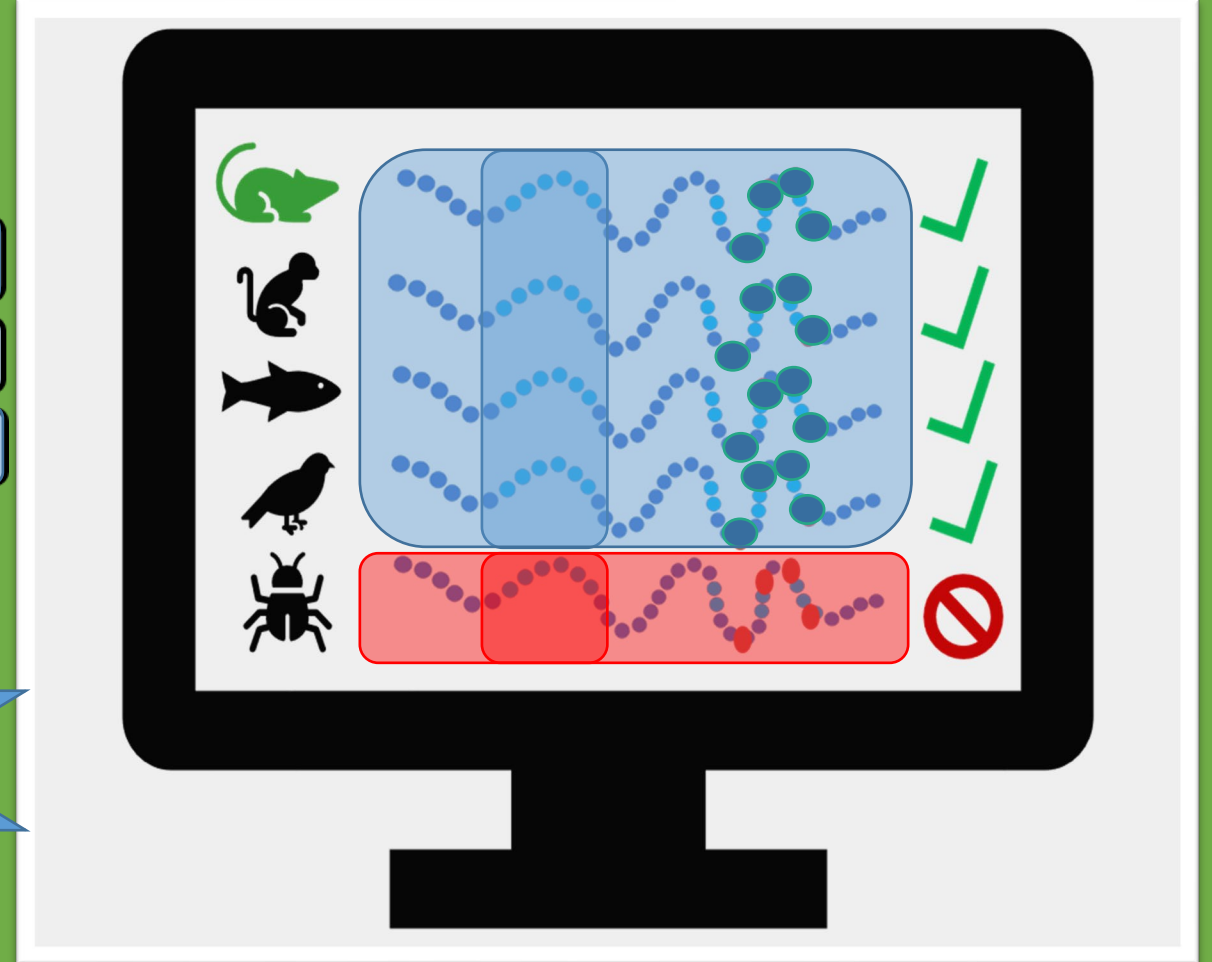
Greater similarity = Greater likelihood that chemical can act on the protein  
Line of Evidence: Predict Potential Chemical Susceptibility Across Species



### Flexible Analysis Based On Available Data

- Level 1** Primary Amino Acid Sequence Alignments
- Level 2** Conserved Functional Domain Alignments
- Level 3** Critical (Close Contact) Amino Acid Conservation

[seqapass.epa.gov/seqapass/](http://seqapass.epa.gov/seqapass/)



## Gather Lines of Evidence Toward Protein Conservation



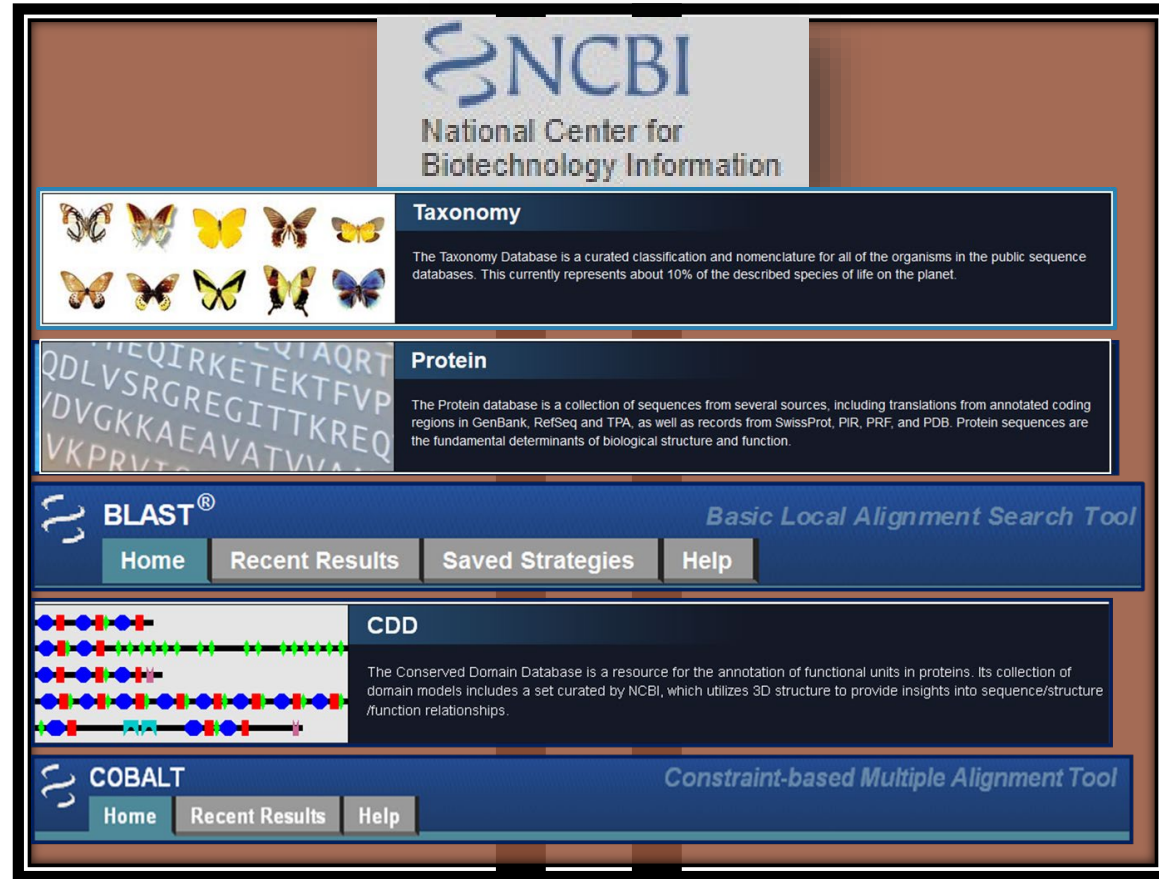


Animation by: Miguel Moravec (EPA CSS) & Andrew Patterson



SeqAPASS

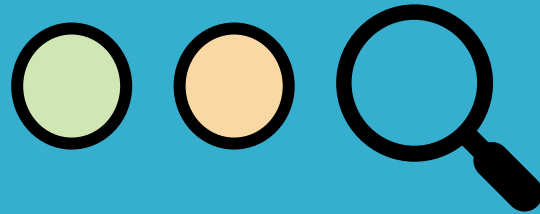




# SeqAPASS

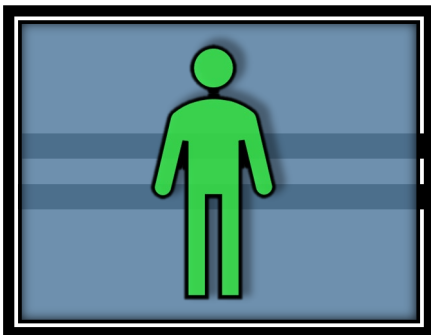
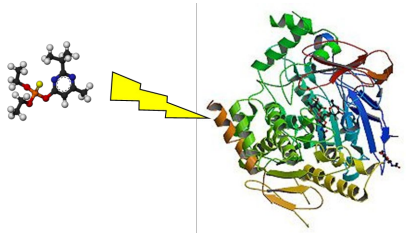
# SeqAPASS

## Level 1



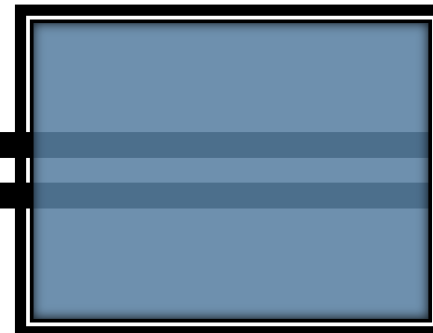
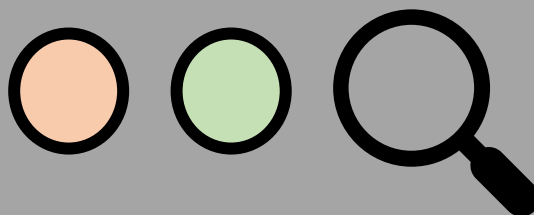
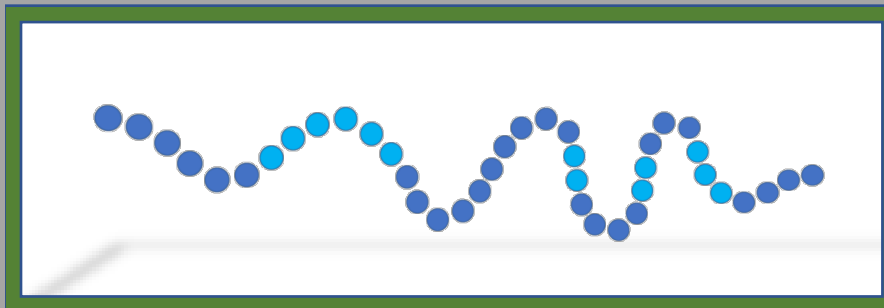


Human Protein Target

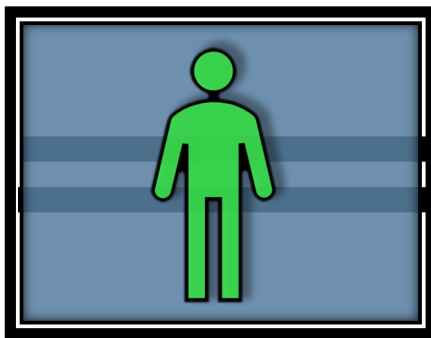
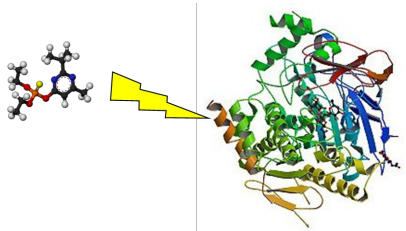


# SeqAPASS

## Level 1

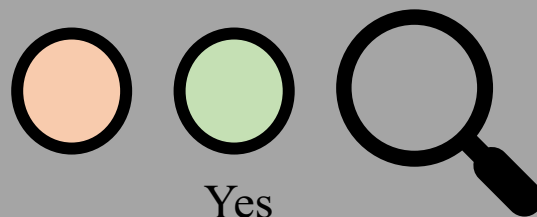
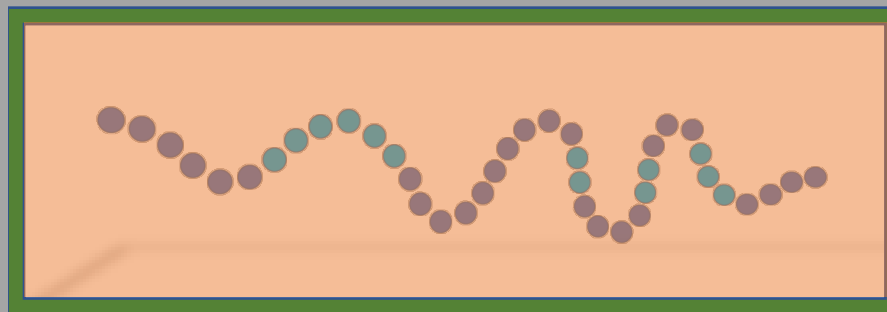


Human Protein Target



# SeqAPASS

## Level 1



Line of Evidence:

Primary amino acid sequence

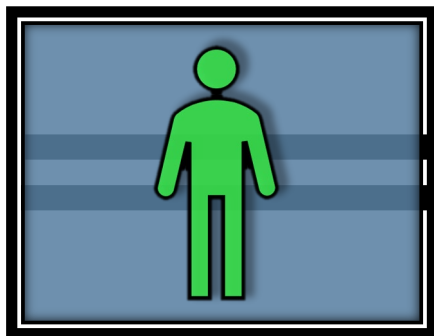
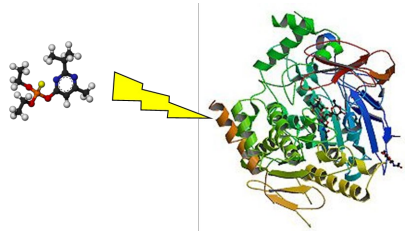
Conserved



Percent similarity

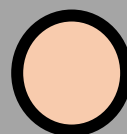
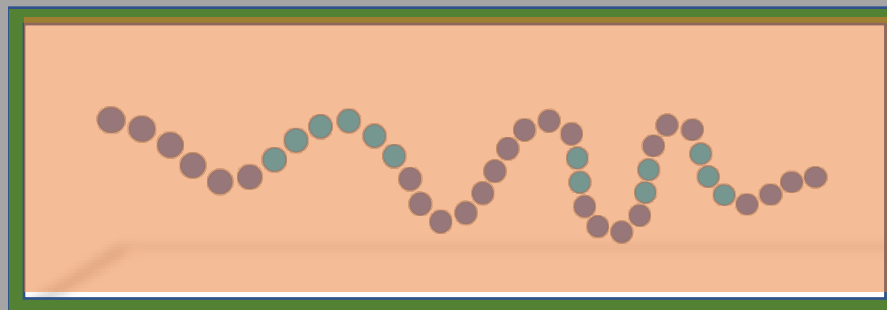


Human Protein Target



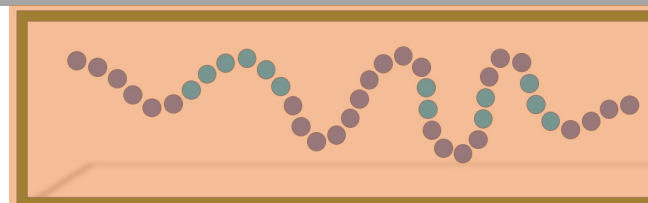
# SeqAPASS

## Level 1

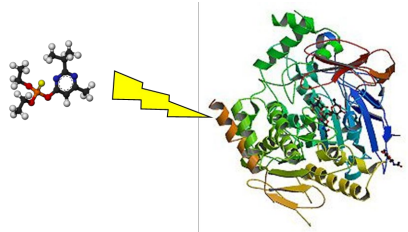


Yes

Line of Evidence:  
Primary amino acid sequence  
Conserved



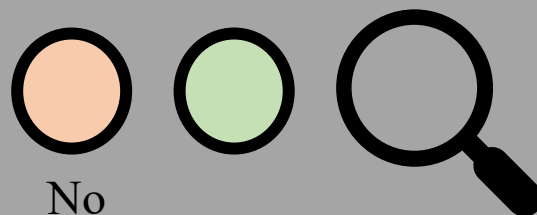
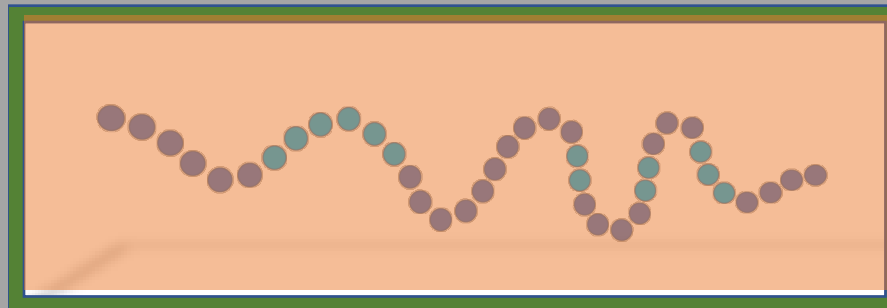
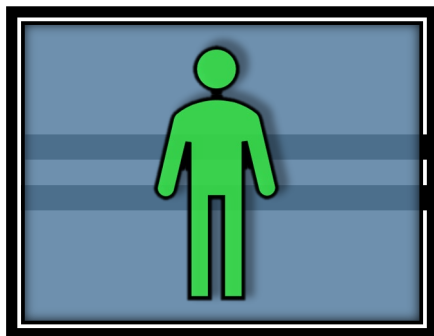
Human Protein Target



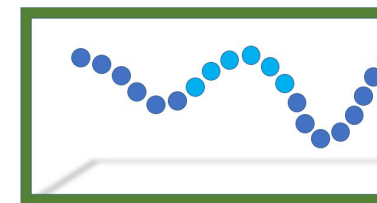
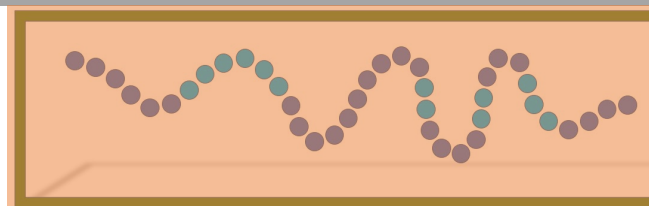
# SeqAPASS

## Level 1

Line of Evidence:  
Primary amino acid sequence  
Conserved



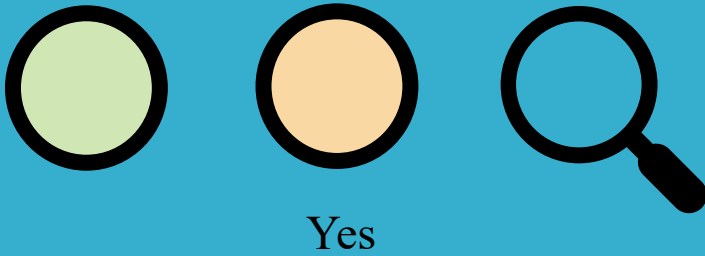
No





# SeqAPASS Level 1

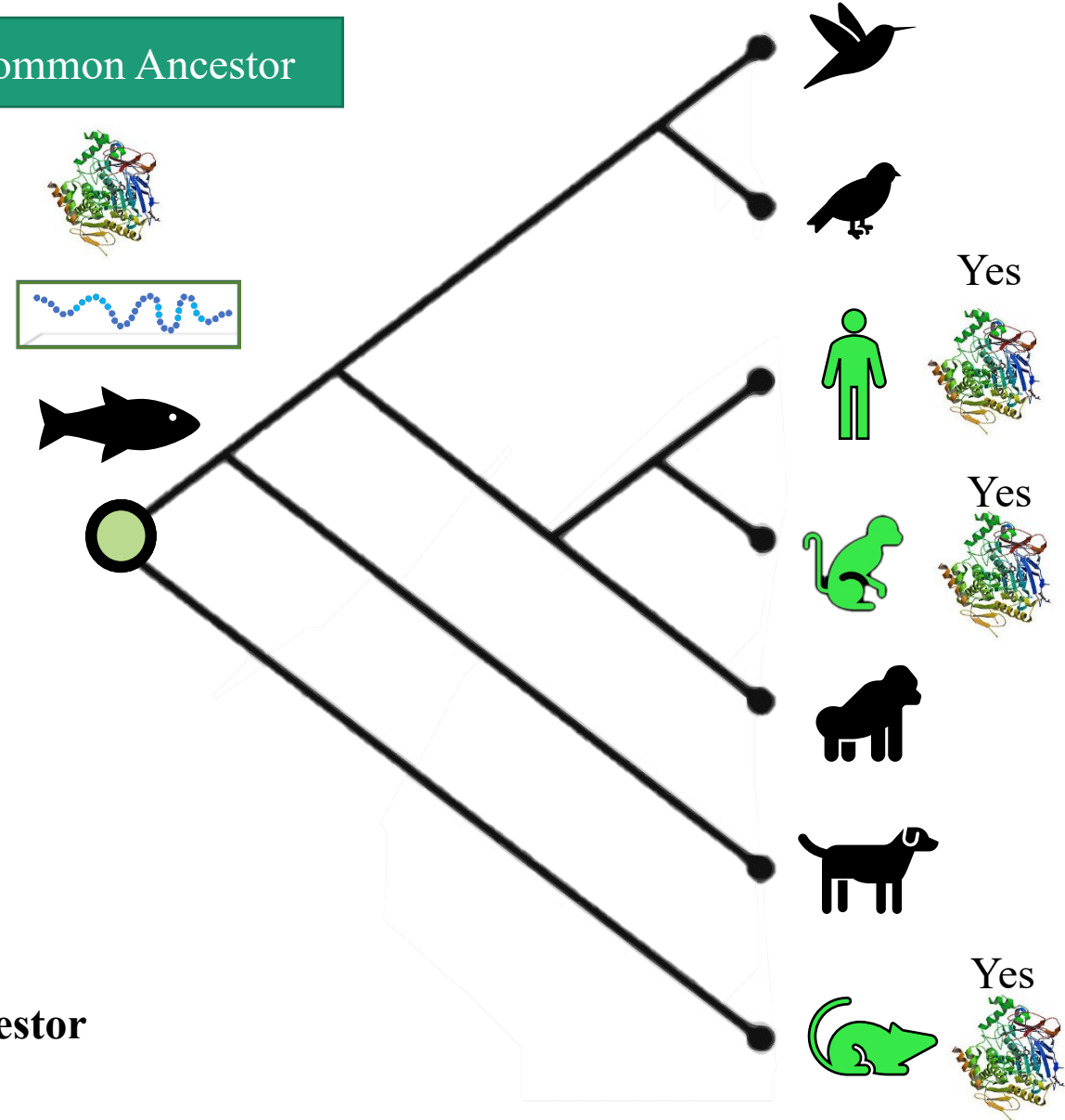
## Ortholog Candidate Identification



Proteins in different species that evolved from a common ancestor

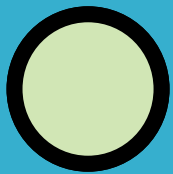
Typically maintain similar function

Common Ancestor

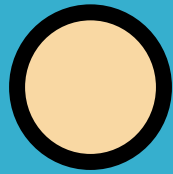


# SeqAPASS Level 1

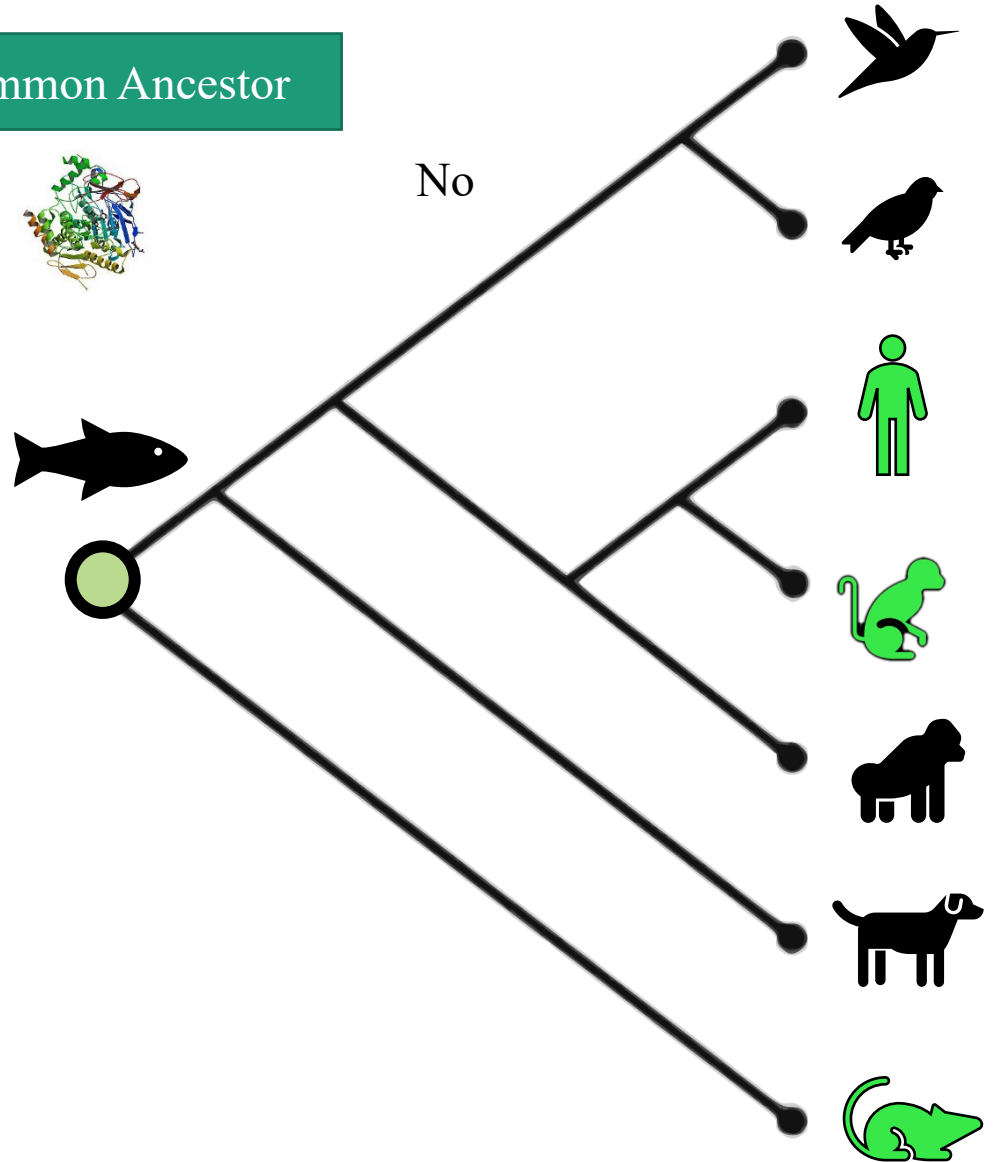
## Ortholog Candidate Identification



No



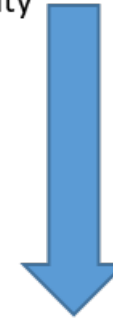
Common Ancestor





# SeqAPASS Level 1

Common Name	Ortholog Candidate	Cut-off	Percent Similarity
Human	Y	33.15	100
Florida manatee	Y	33.15	98.8
Mallard	Y	33.15	82.29
Rock pigeon	Y	33.15	80.93
Green anole	Y	33.15	80.65
Pacific transparent sea squirt	Y	33.15	33.15
Yesso scallop	N	33.15	32.87
Purple sea urchin	N	33.15	26.05
Human whipworm	N	33.15	23.53
Bed bug	N	33.15	21.62



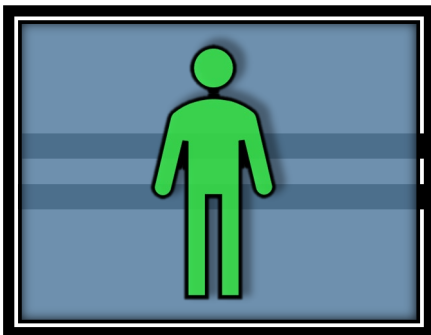
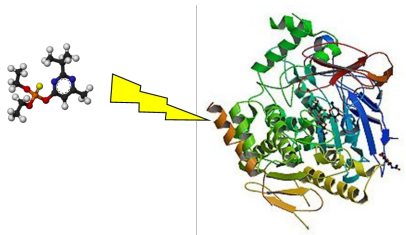
Lowest % Similarity that is still an ortholog

## Example:

Susceptibility Cut-off: Set at 33.15

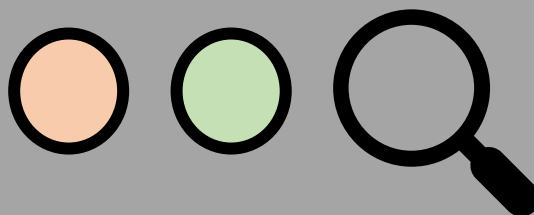
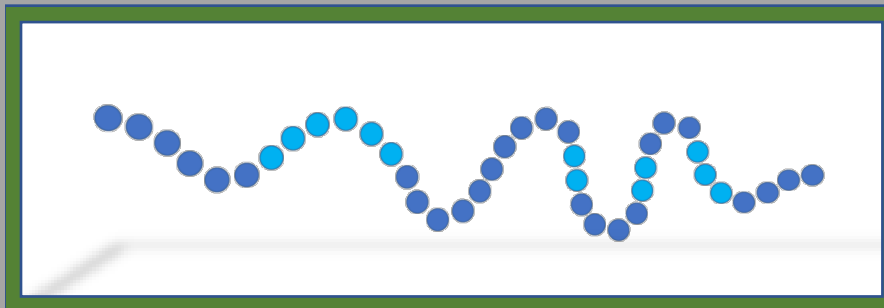
Above cut-off: More likely to be susceptible base on similar **FUNCTION**

Human Protein Target



# SeqAPASS

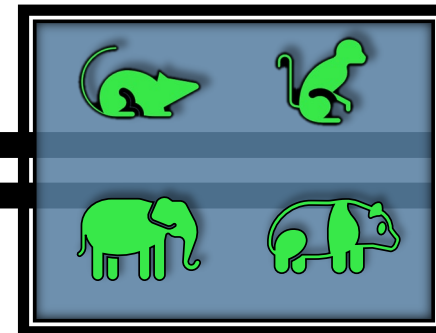
## Level 1



Line of Evidence:

Primary amino acid sequence

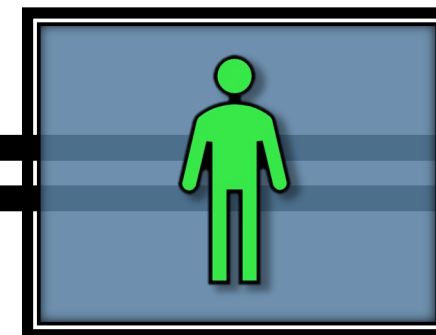
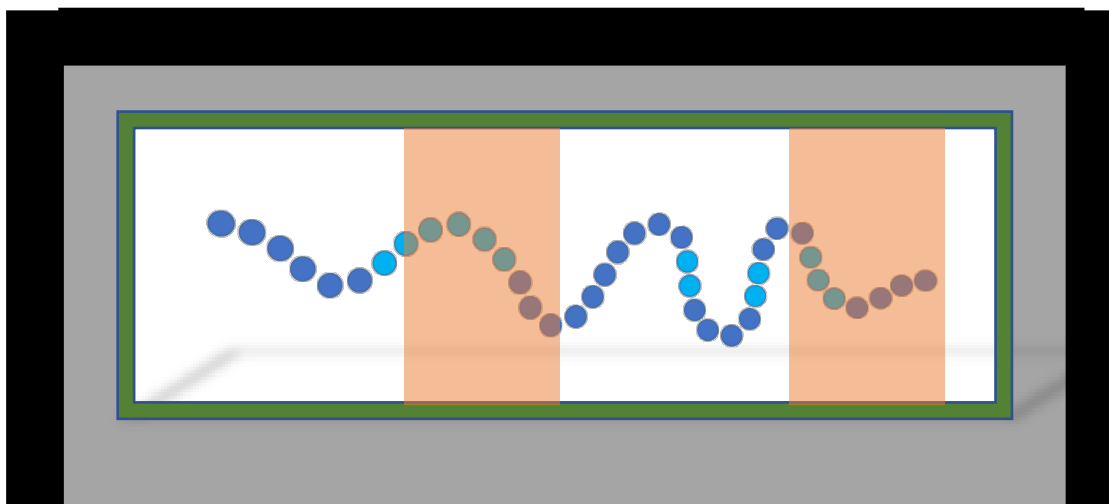
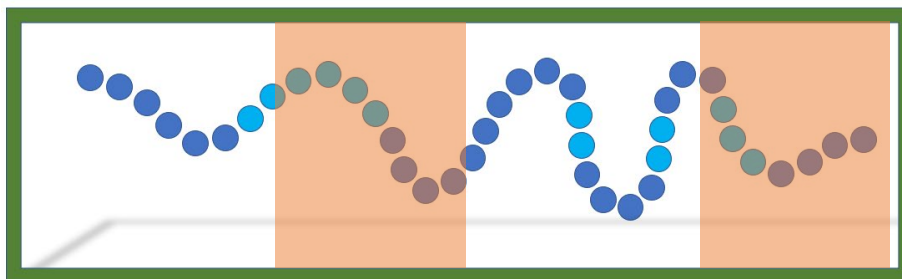
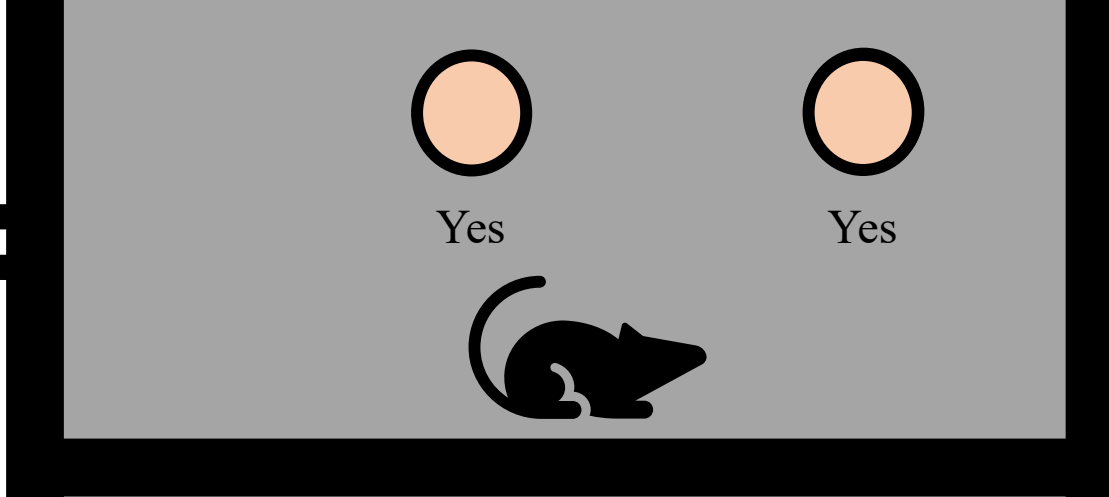
Conserved



Hundreds to Thousands of Species

# Level 2

Line of Evidence:  
Domain  
Conserved

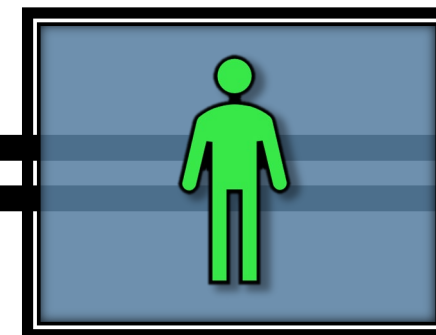
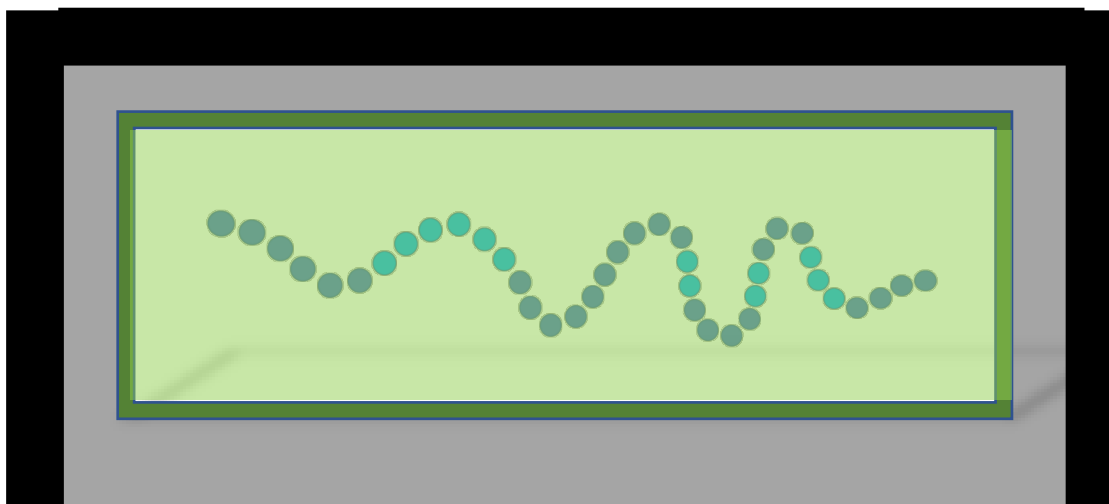
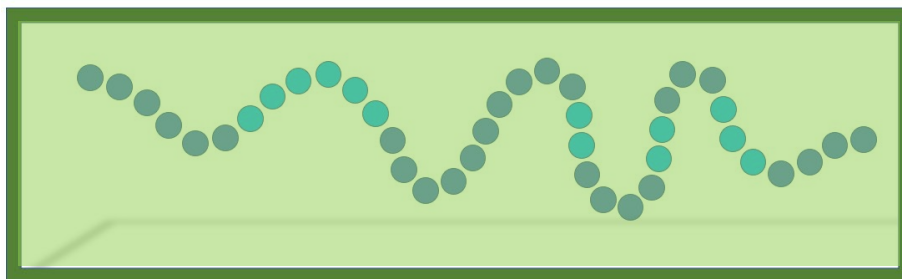
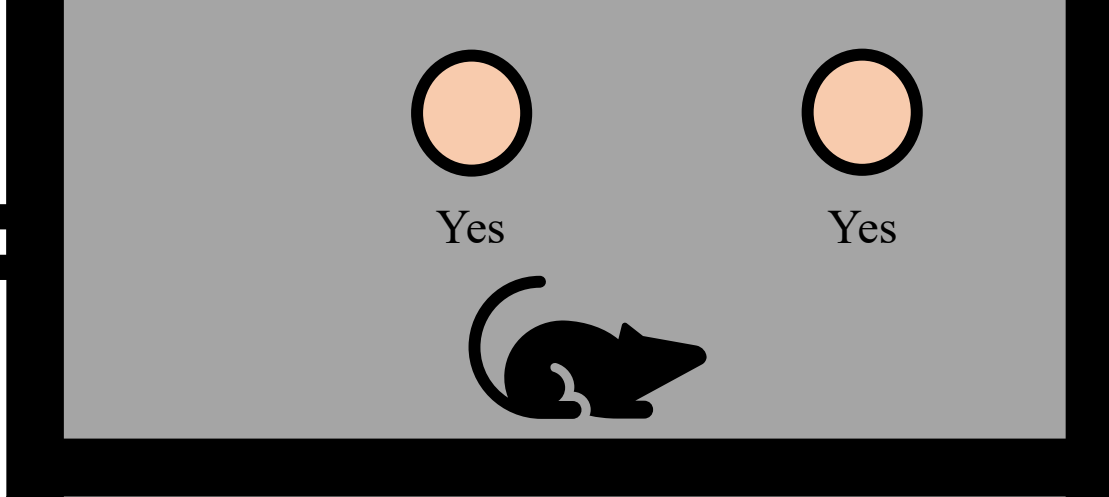


Human Functional Domain(s)





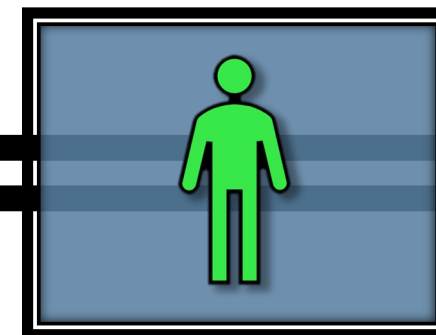
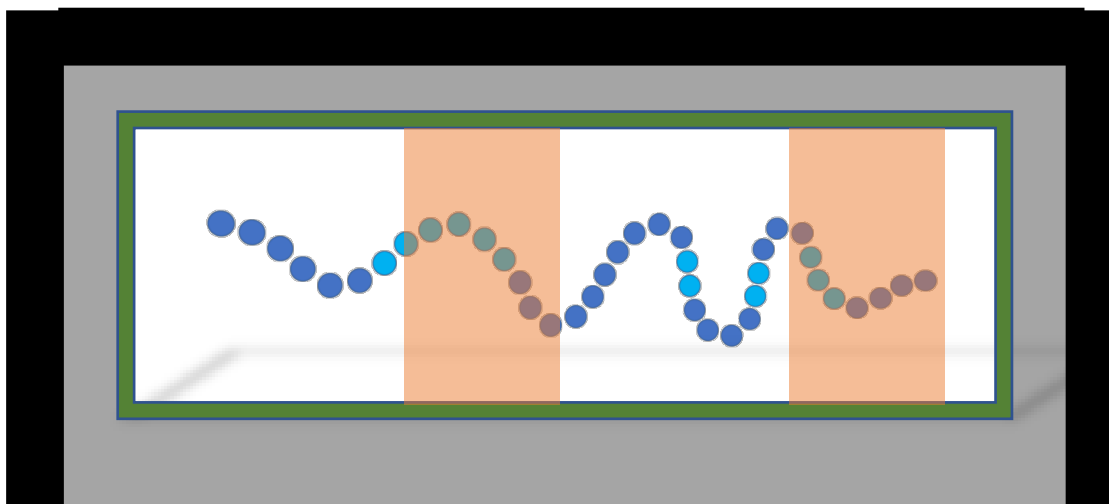
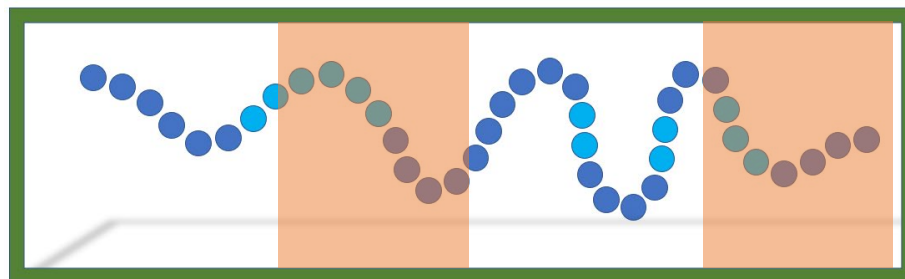
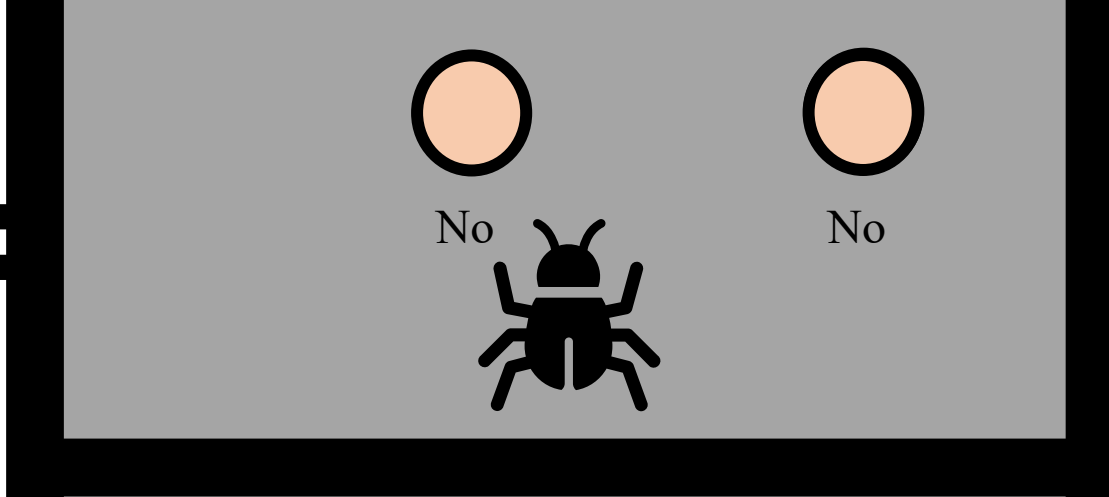
Line of Evidence:  
Domain  
Conserved



Human Functional Domain(s)

# Level 2

Line of Evidence:  
Domain  
Conserved



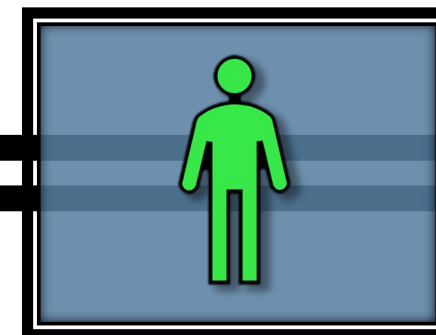
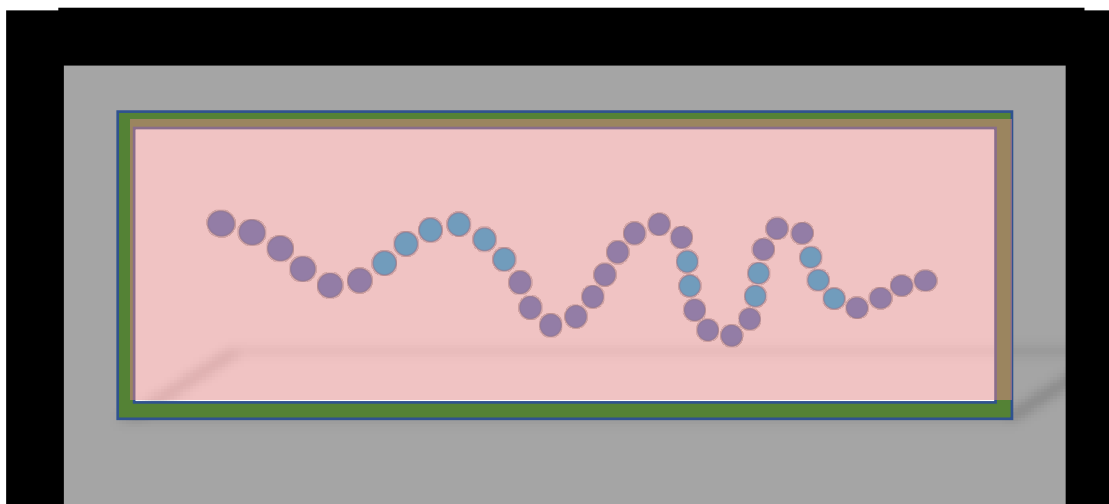
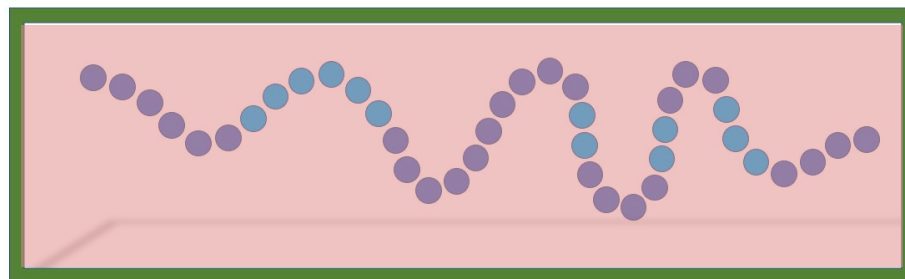
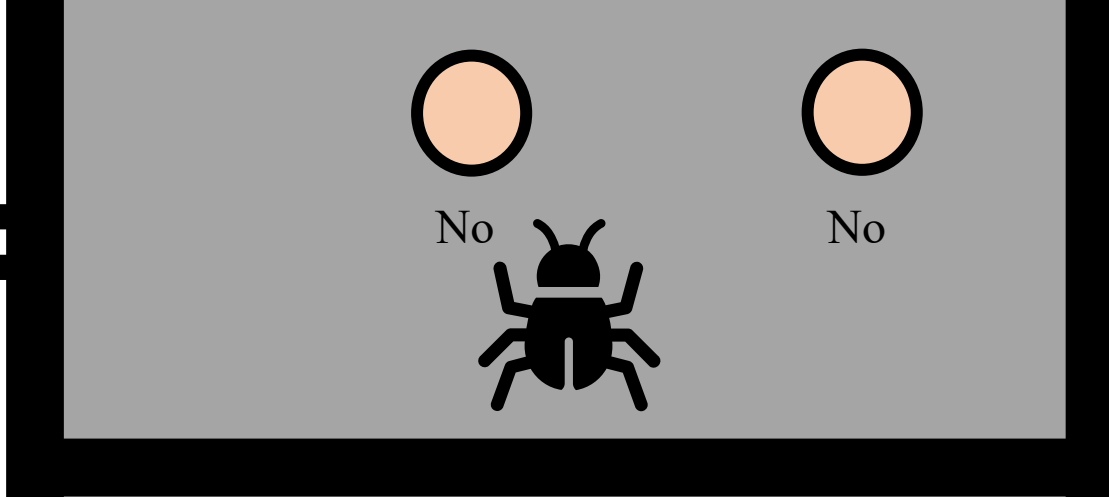
Human Functional Domain(s)

# Level 2

Line of Evidence:

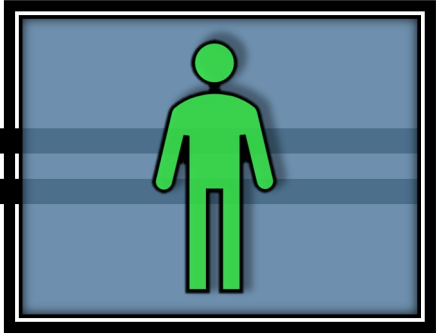
Domain

Not Conserved

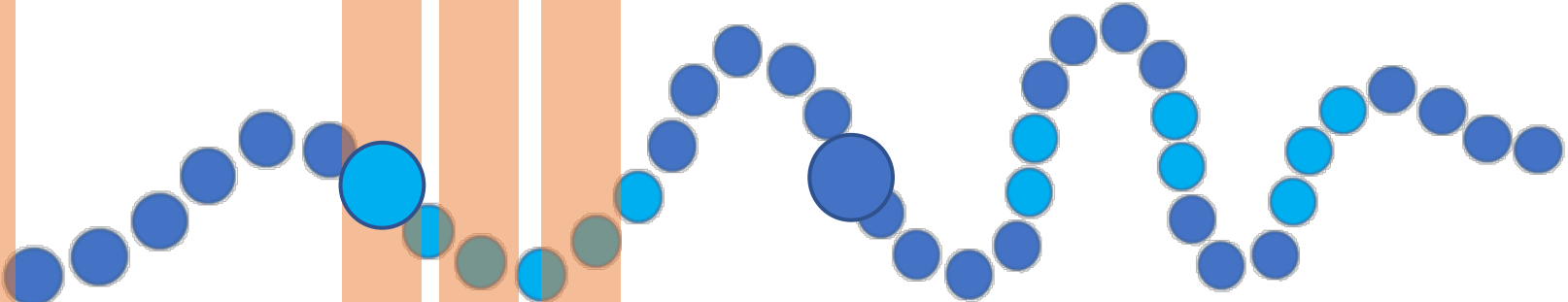
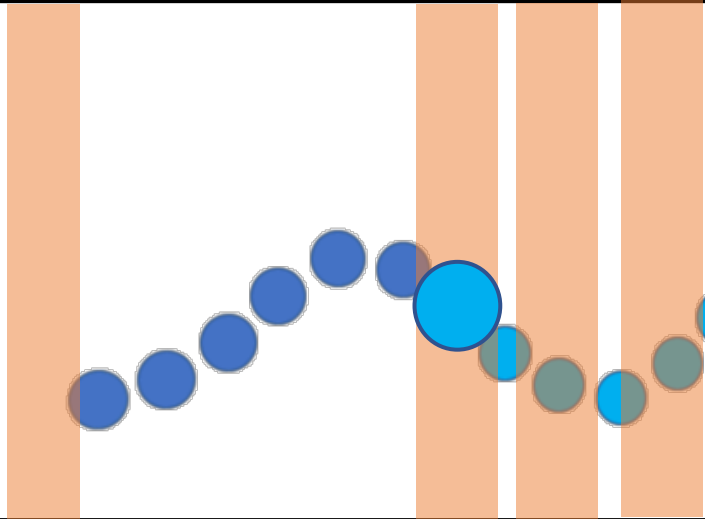
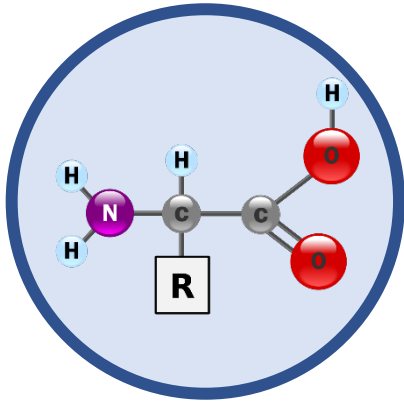
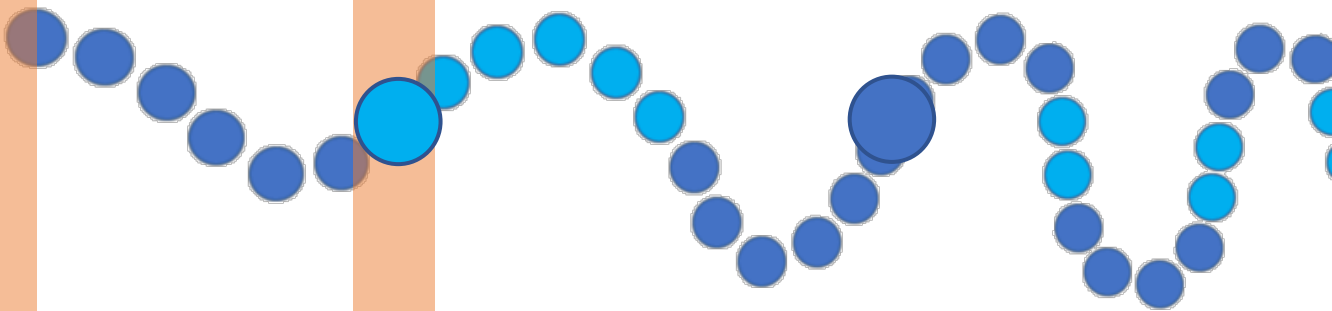
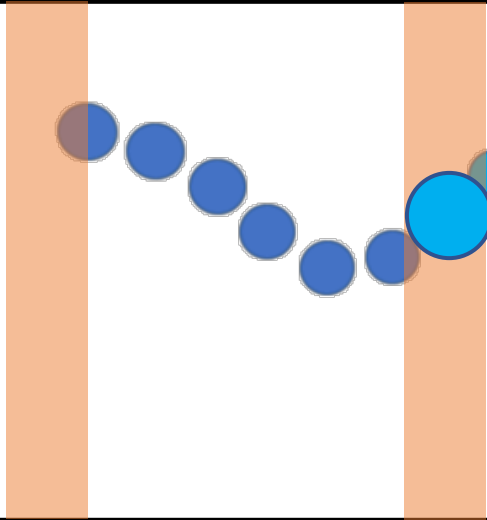
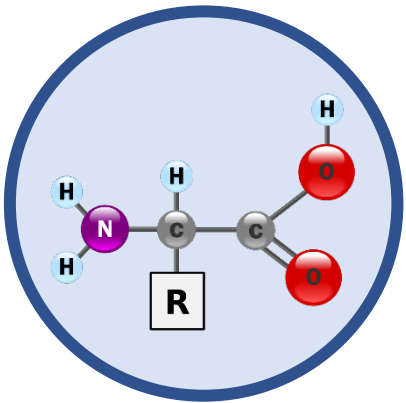


Human Functional Domain(s)

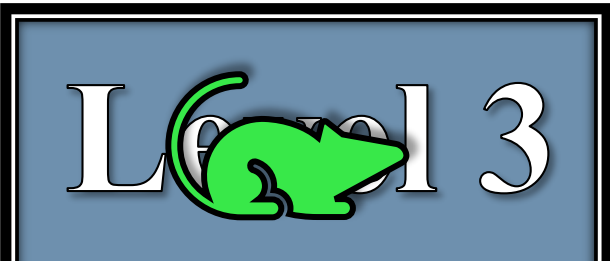


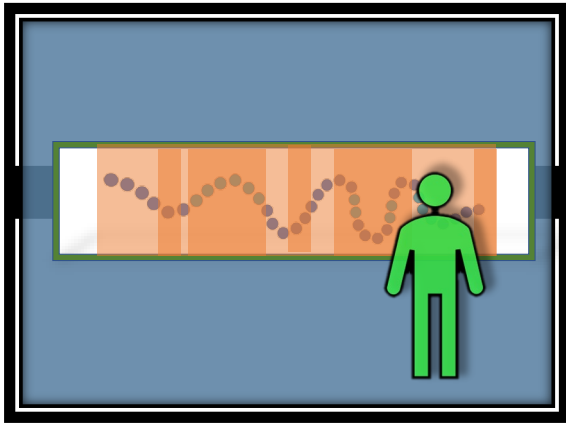


Human Critical Amino Acids



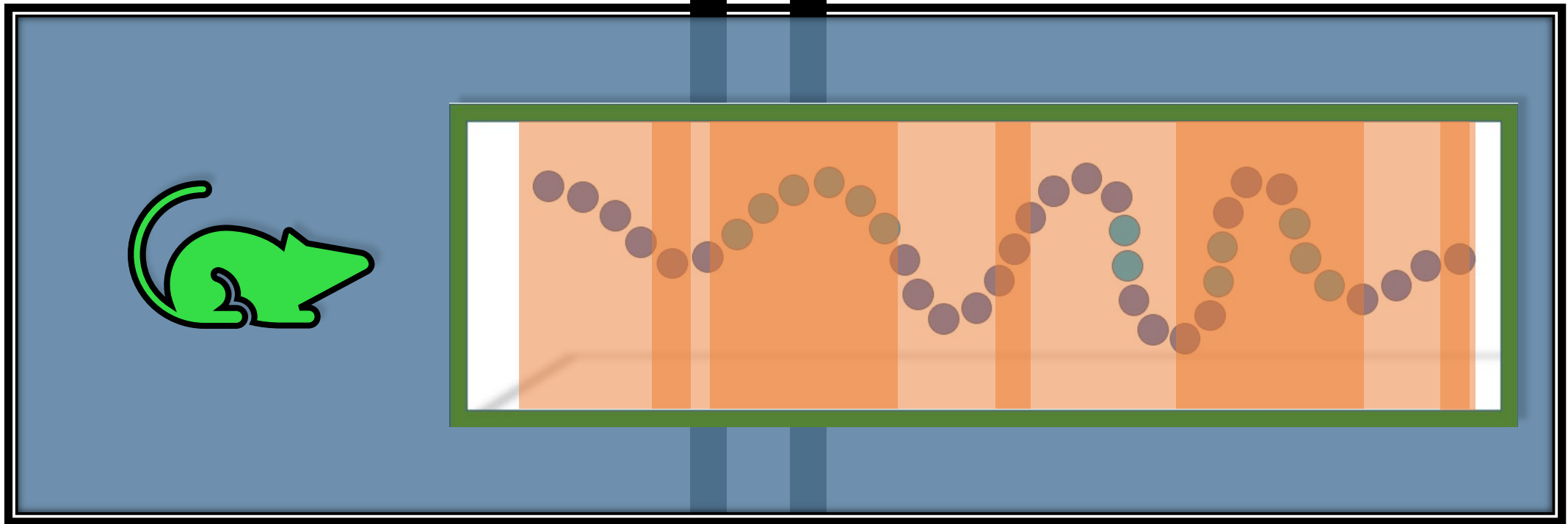
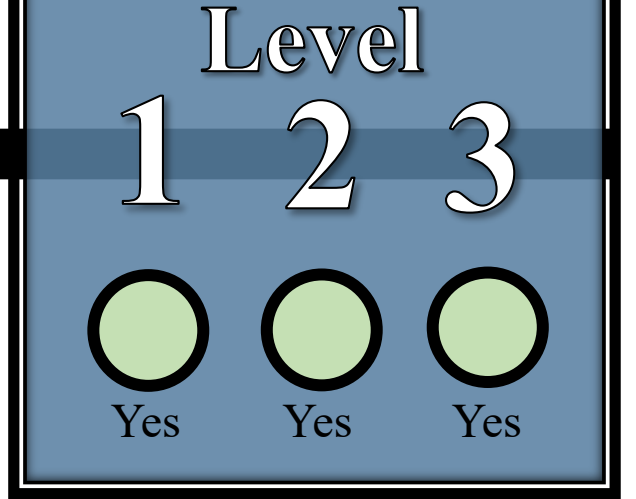
Line of Evidence: Conserved





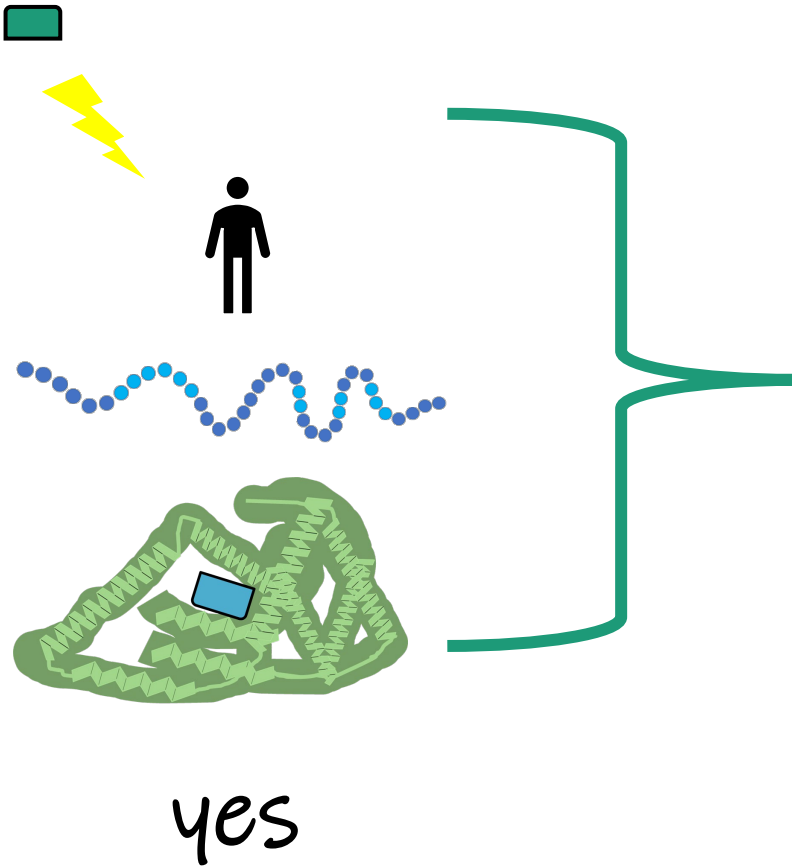
# SeqAPASS











## Summary



**Gather Lines of Evidence for Conservation of Protein Target:**  
**Susceptibility Prediction: Yes or No**

# SeqAPASS Predicts Likelihood of Similar Susceptibility based on Sequence Conservation:



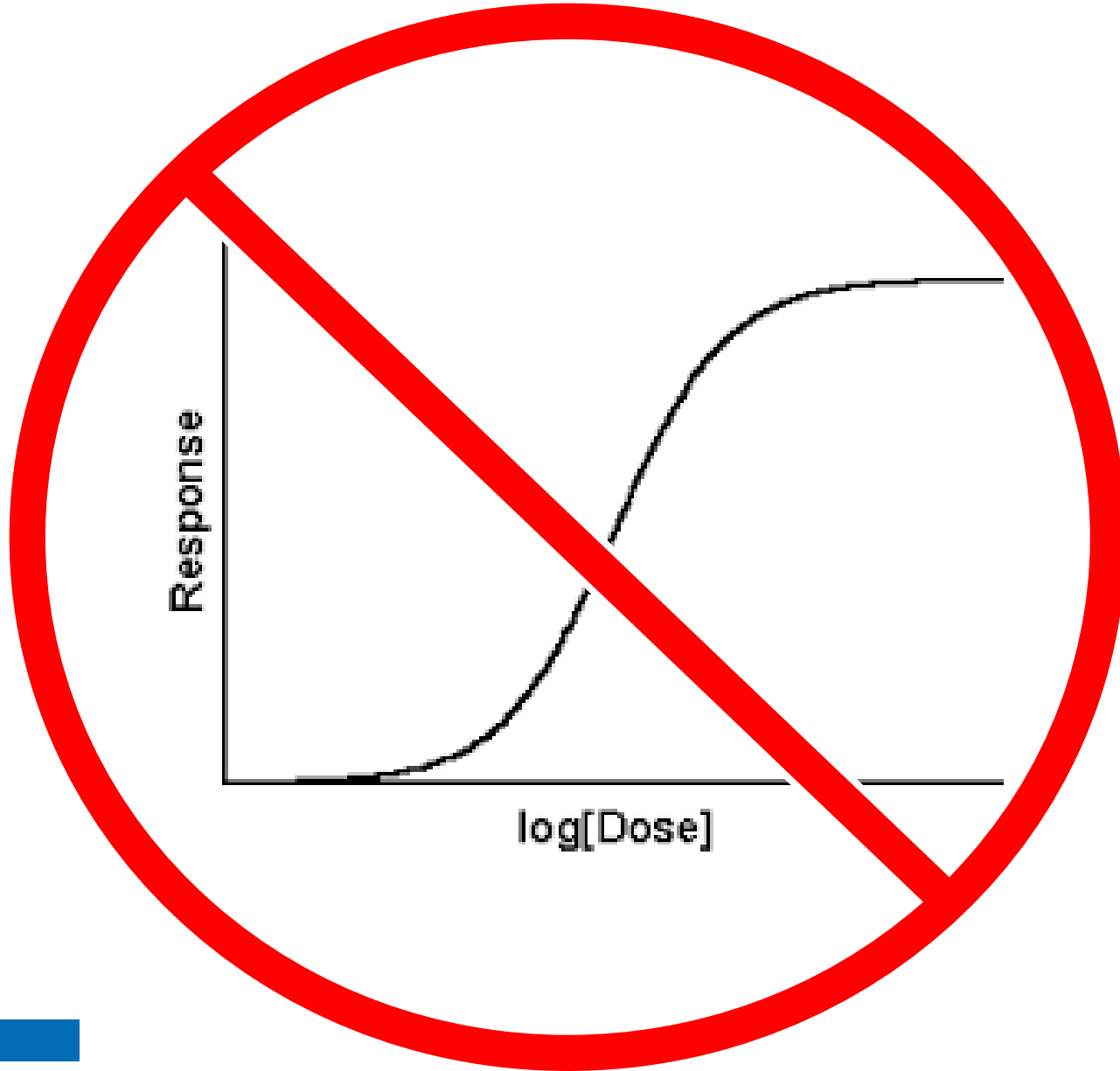
	yes
	yes
	yes
	yes
	yes
	yes
	yes
	no
	yes
	no

Line(s) of evidence indicate

- The protein is conserved
- The protein is NOT conserved



# SeqAPASS DOES NOT predict the degree of sensitivity/susceptibility:



## Factors that make a species sensitive

- Exposure
- Dose
- ADME
- Target receptor availability
- Life stage
- Life history
- etc.
- etc.

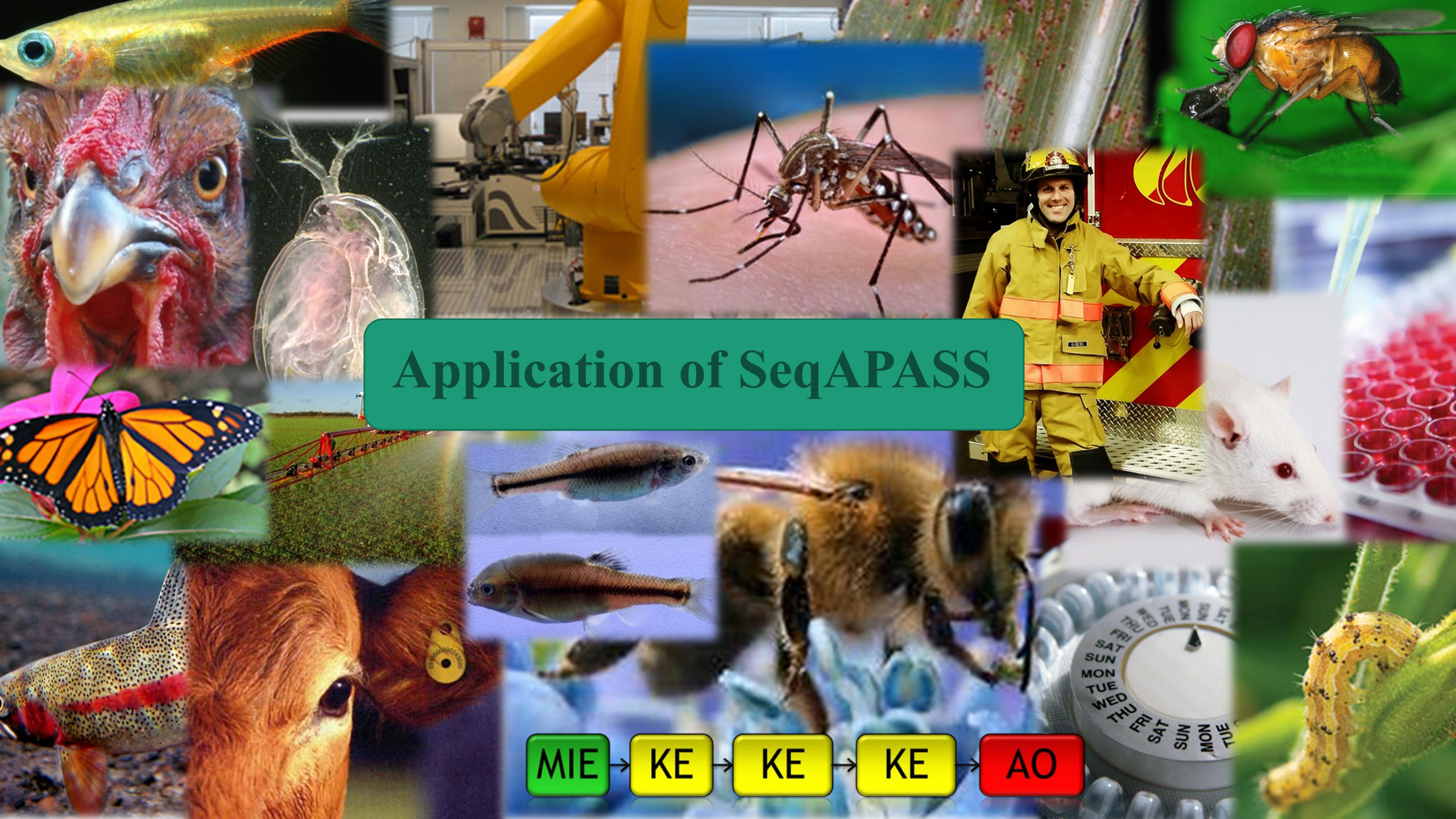




# Strengths of SeqAPASS

- Publicly available to all
- Lines of evidence for conservation for 100s-1000s of species rapidly
- Takes advantage of well-established tools and databases
- Streamlined, consistent, transparent, and published methods
  - Case examples to demonstrate applications
- Guides users to appropriate input
- Evolves as bioinformatics approaches become more user friendly
  - Smart automation or semi-automation





# Application of SeqAPASS



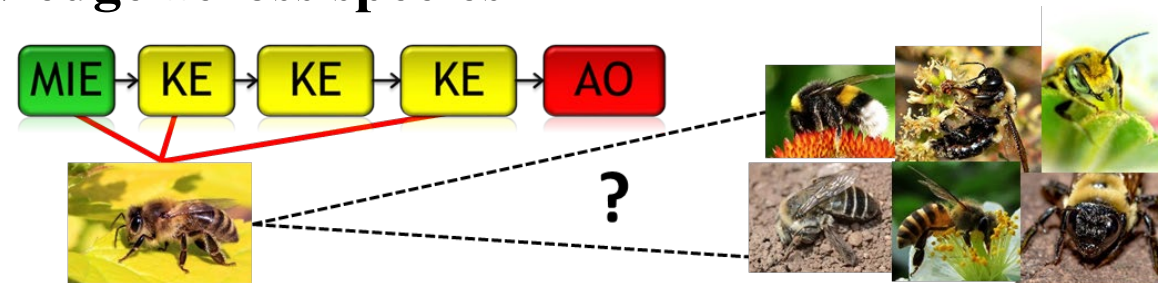


# Applications of Bioinformatics: Case Studies

- **Extrapolate adverse outcome pathway knowledge across species**

- Define the taxonomic domain of applicability

- Apis vs Non-Apis bees



- **Extrapolate high throughput screening data**

- Chemicals that target human estrogen receptor alpha, androgen receptor, steroidogenic enzymes, thyroid axis proteins
  - All ToxCast Assay targets

- **Predict relative intrinsic susceptibility**

- Pesticides
  - Endangered Species Act
  - Derivation of Aquatic Life Criteria

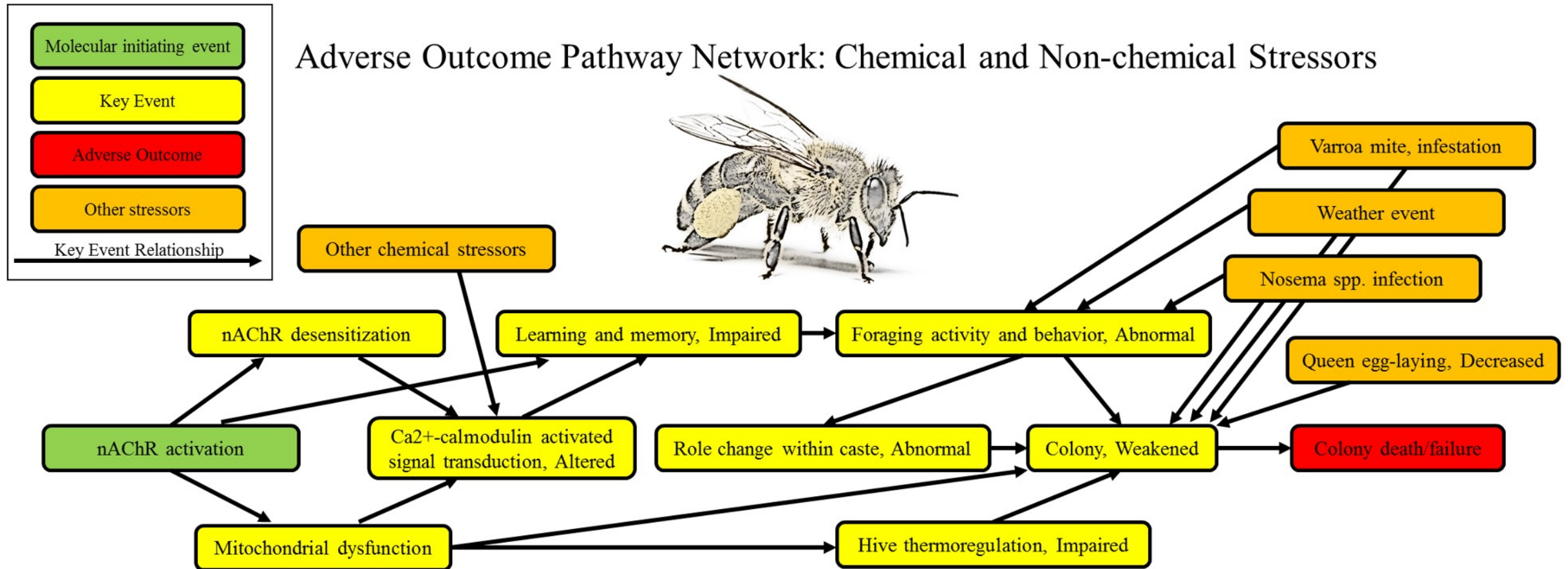
- **Predict chemical bioaccumulation across species**

- Chemicals of concern: PFAS

- **Generate research hypotheses** Strobilurin fungicides

- **Prioritization strategies** Pharmaceuticals





Define Knowledge Gaps

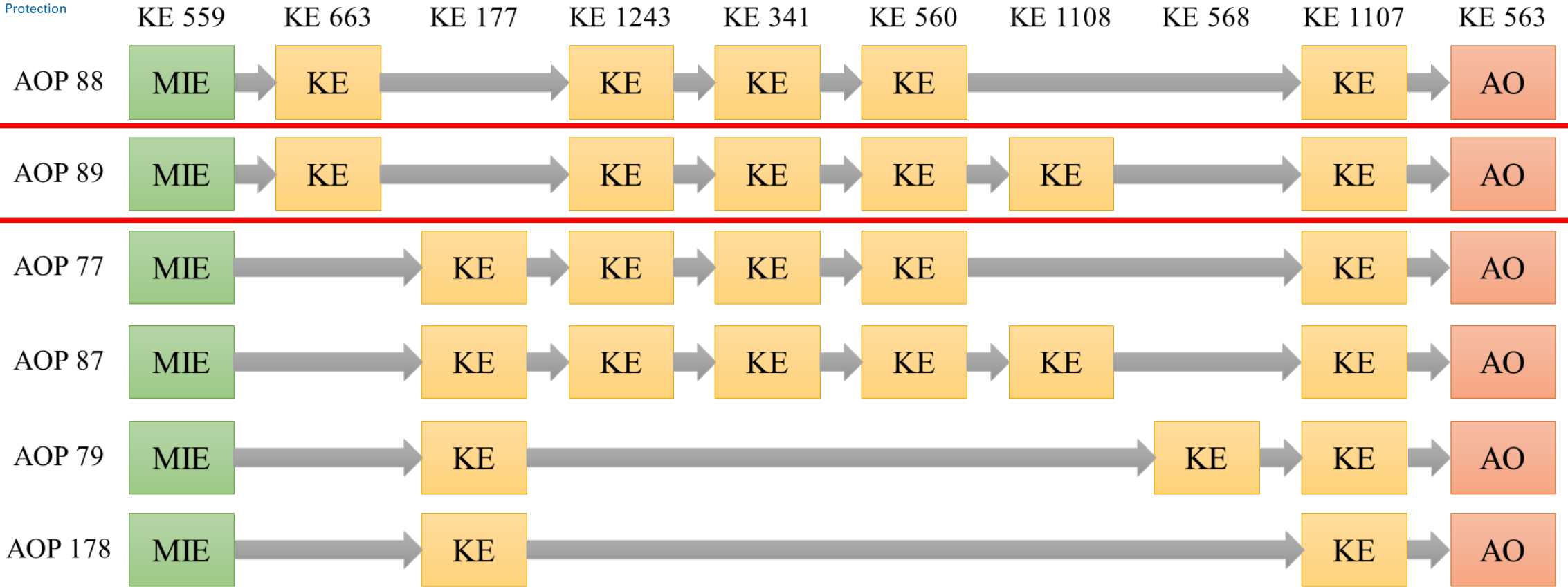
Understand nodes that may be impacted by multiple stressors

Assists in development of mitigation strategies

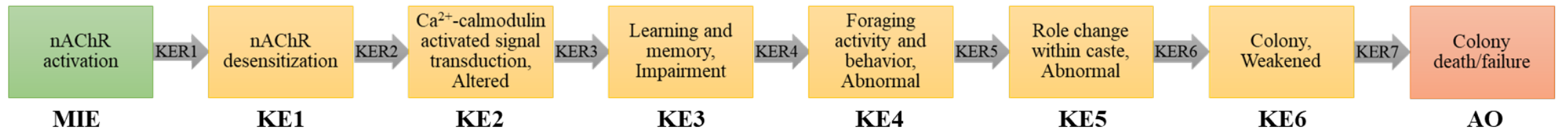
**How to define the taxonomic relevance of an AOP?**

## AOP Wiki Key Event IDs

AOP Wiki AOP IDs



*Modified from LaLone et al., 2017*

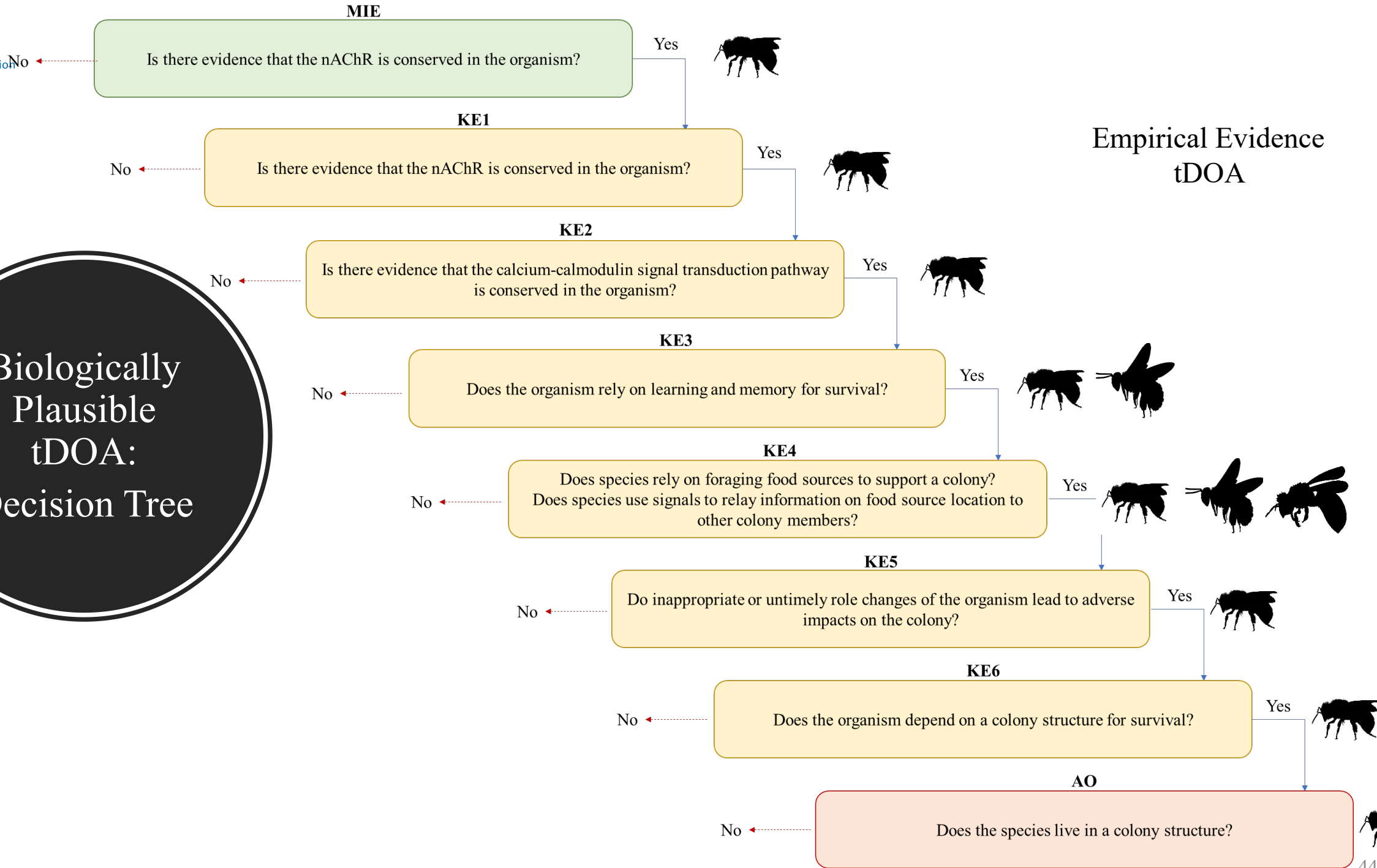


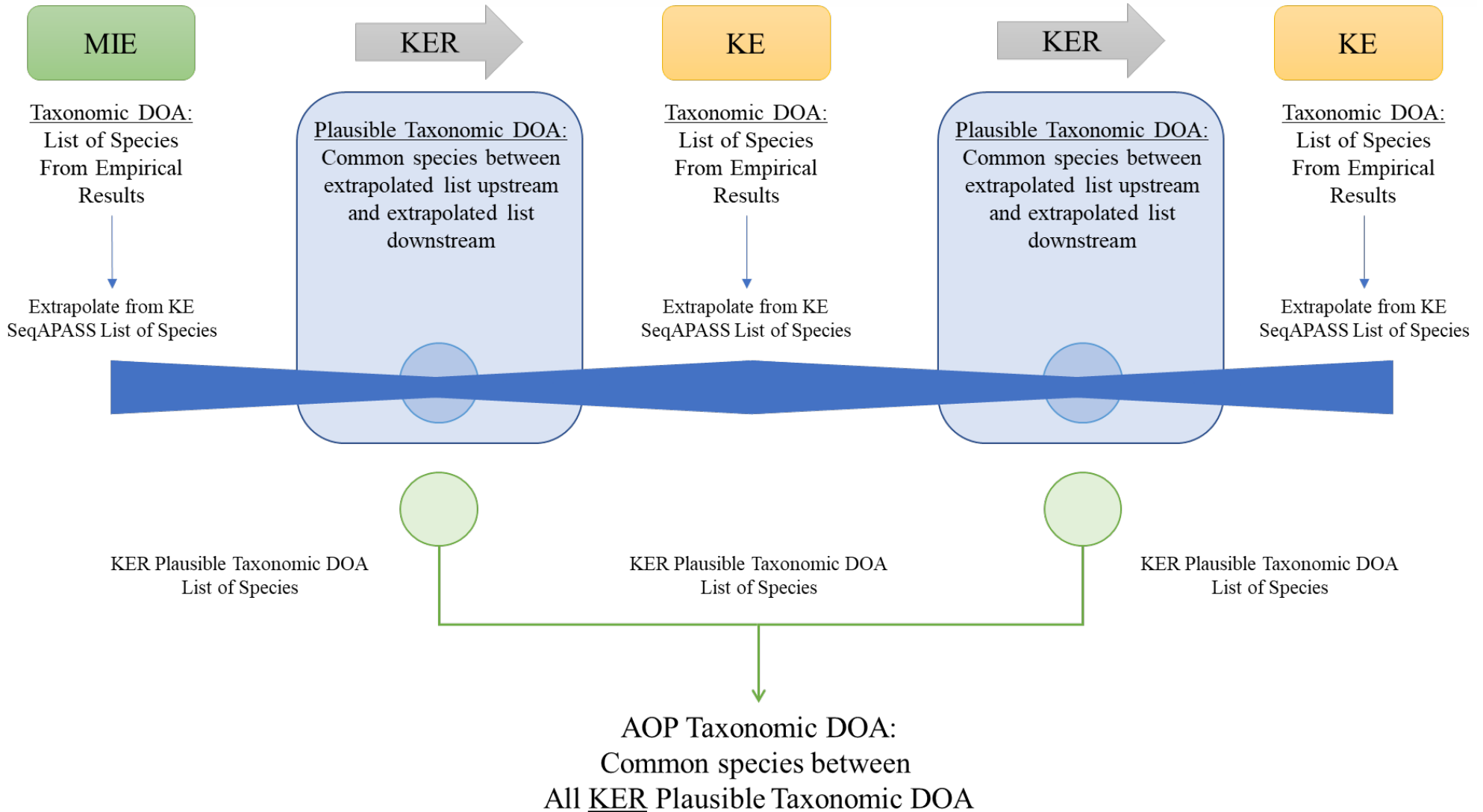
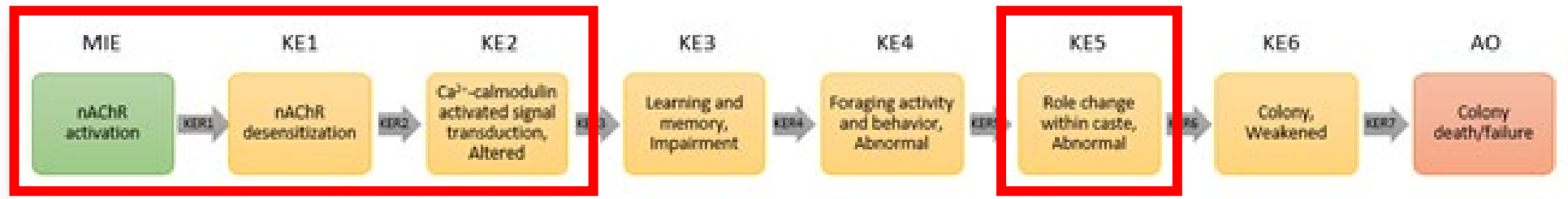
### Case Example

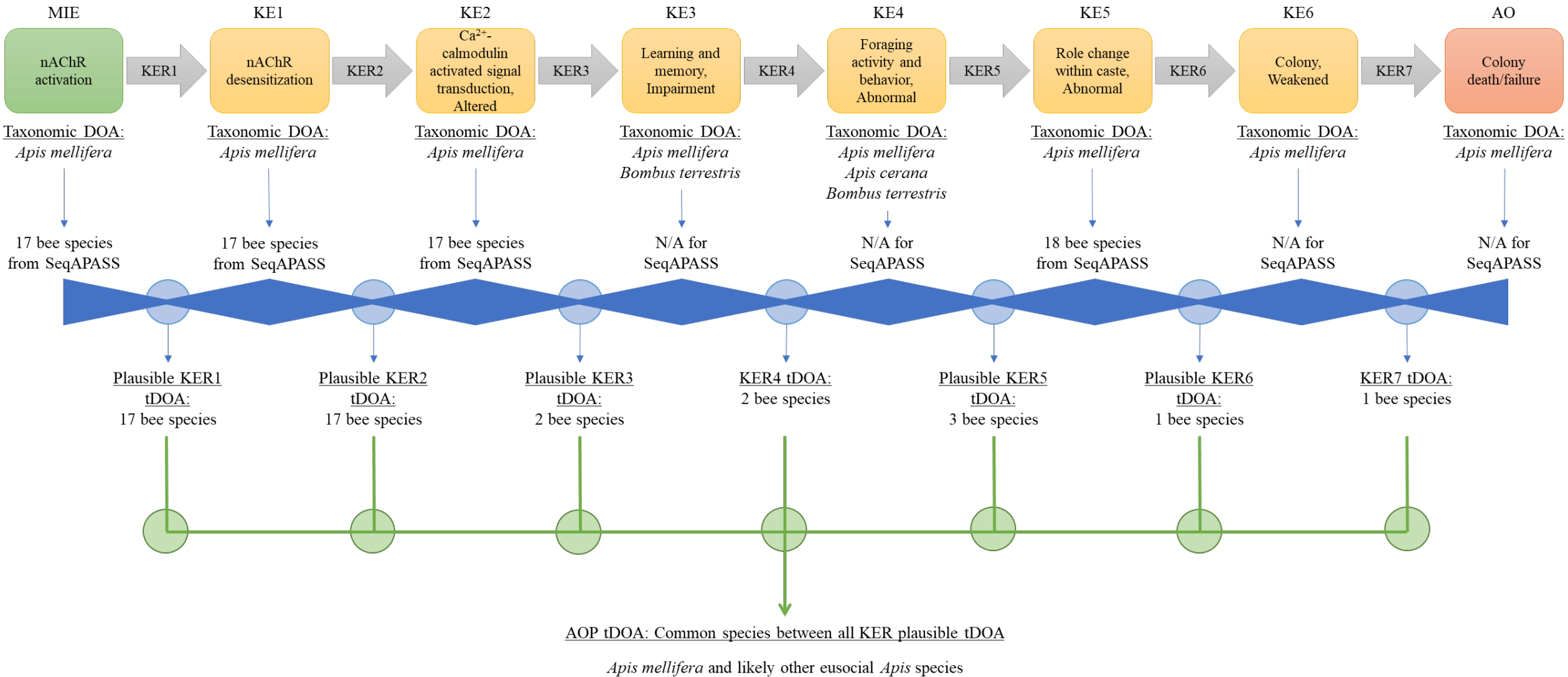
**How to define the taxonomic relevance of an AOP?**

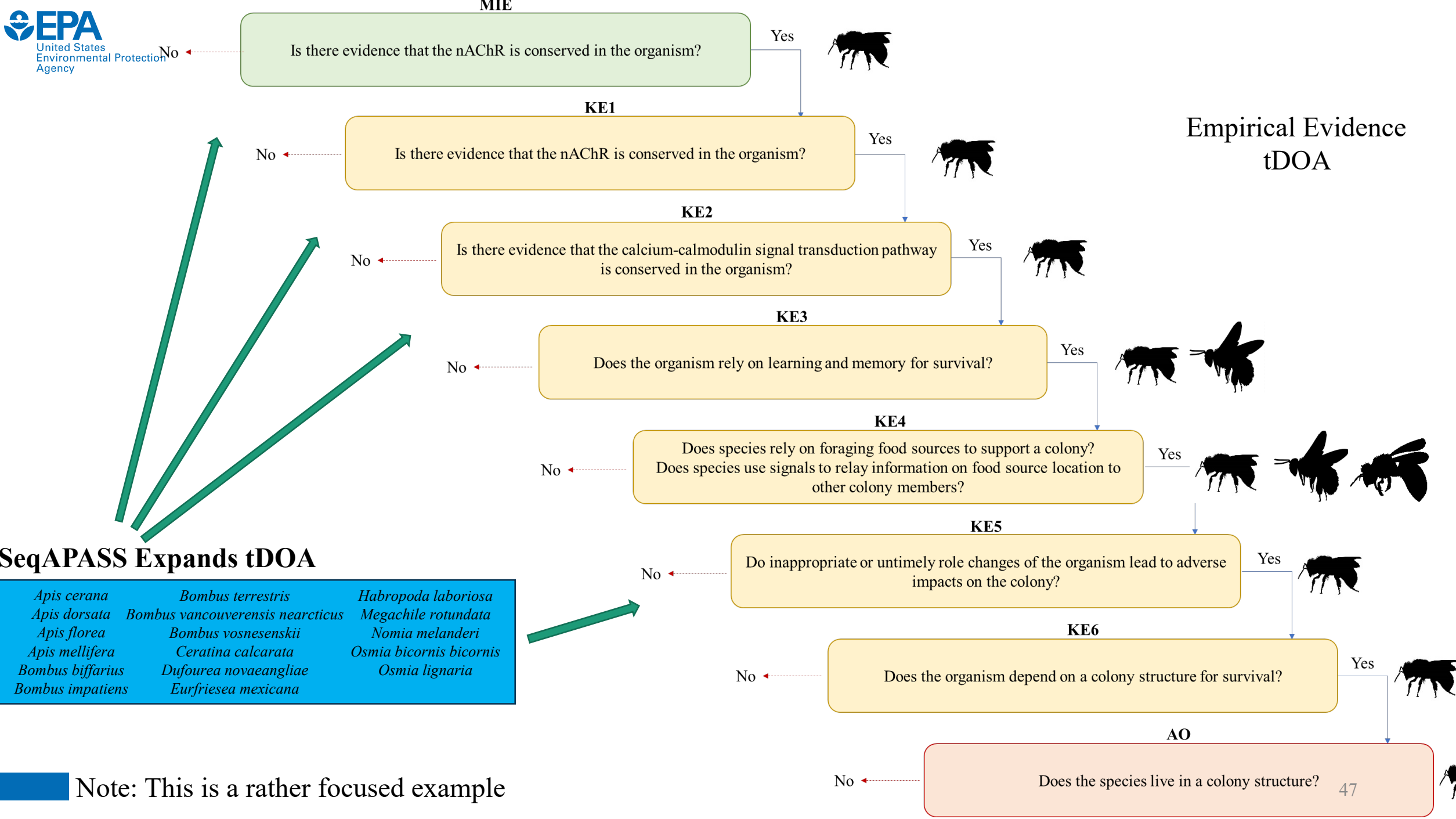


Biologically  
Plausible  
tDOA:  
Decision Tree











Advances in Bioinformatics –  
Future of SeqAPASS

# Always Look Several Steps Ahead



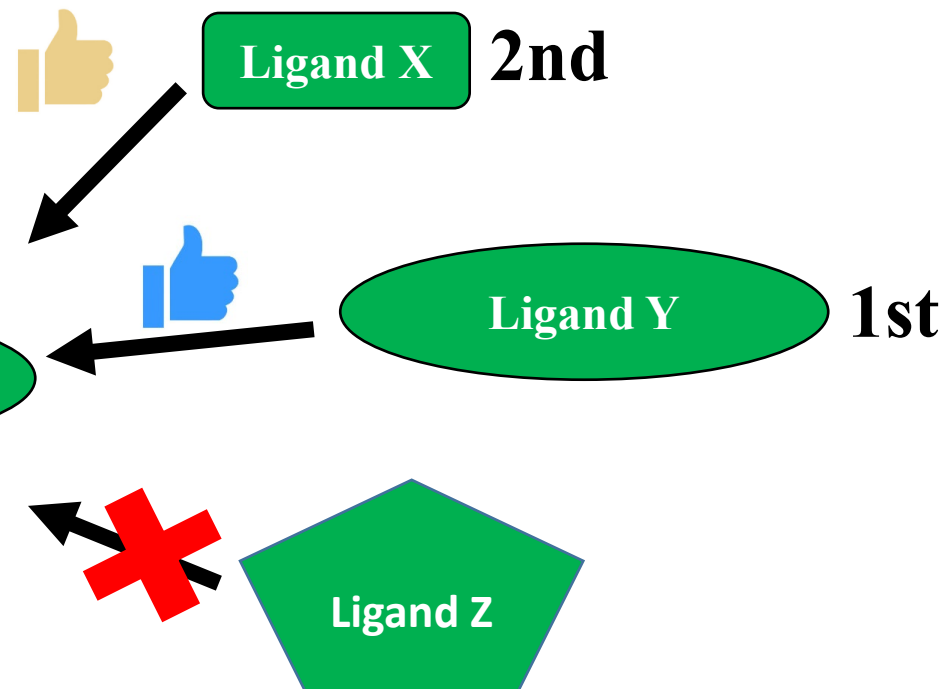
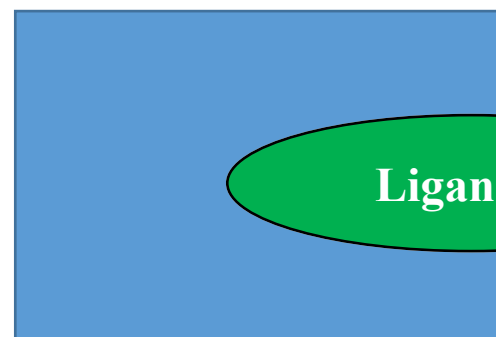
# Advances in Drug Discovery/Development

(COVID-19 has led to advances)



Structure derived  
from X-ray  
crystallography

Human  
Protein Structure



## Bioinformatics Toolbox:

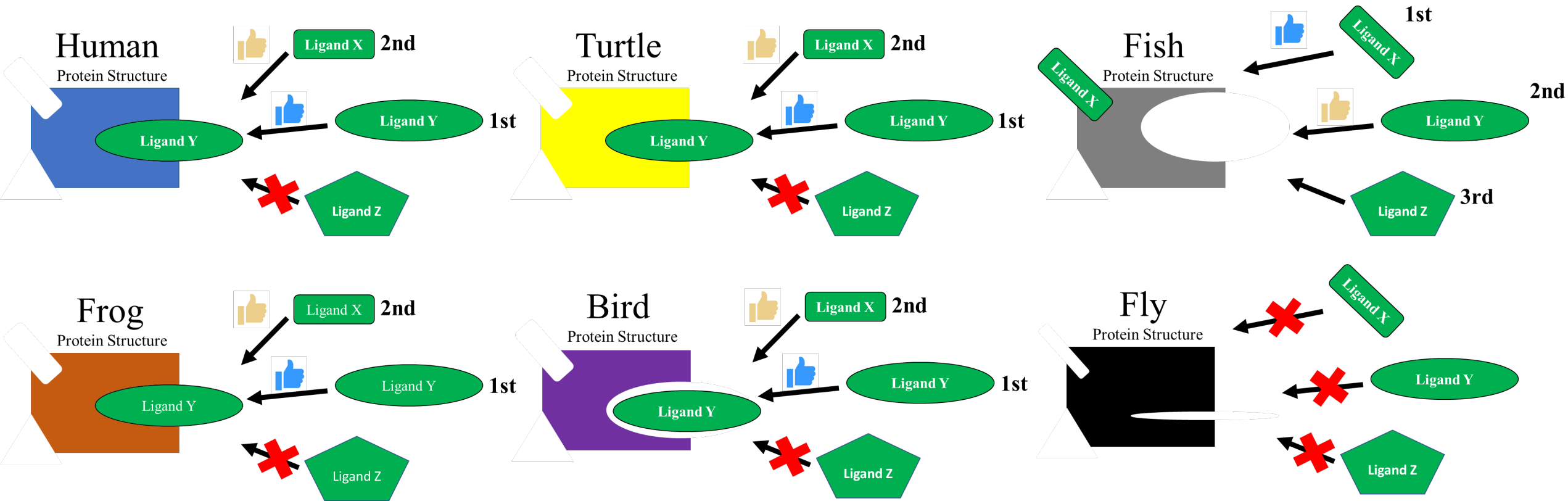
Molecular modeling

Molecular docking

Virtual screening

Molecular dynamic simulations

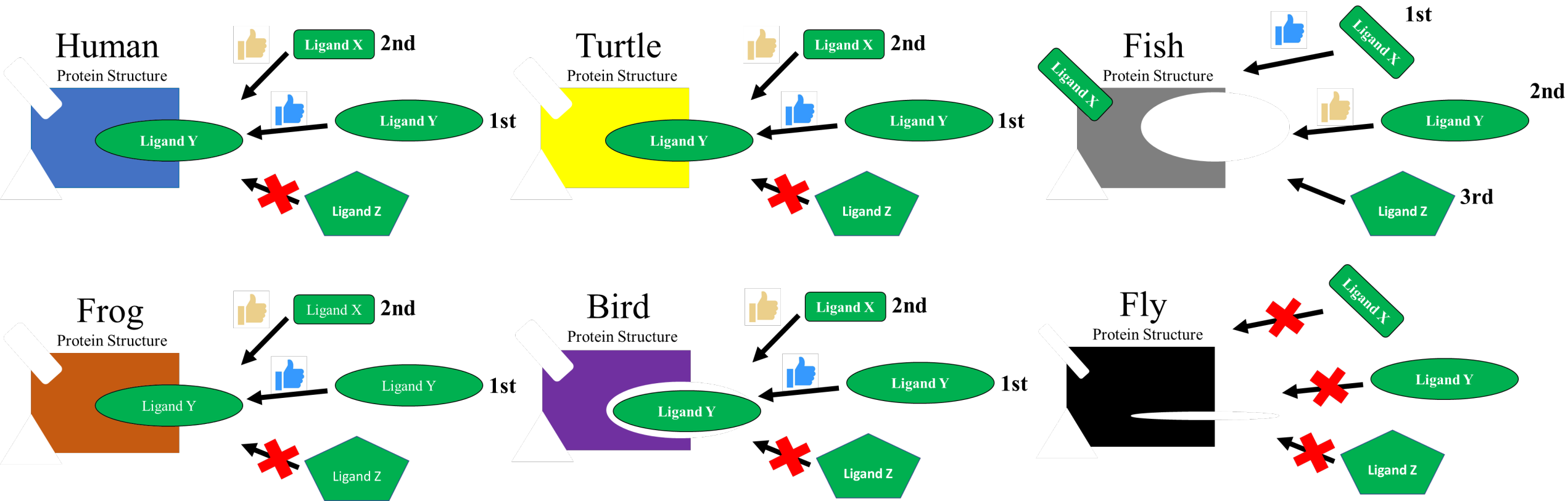
# Application to Species Extrapolation



## Bioinformatics Toolbox:

Molecular modeling  
Molecular docking  
Virtual screening  
Molecular dynamic simulations

# Application to Species Extrapolation



**Bioinformatics Toolbox:**  
Molecular modeling  
Molecular docking  
Virtual screening  
Molecular dynamic simulations

*Thousands/Millions/Billions  
of  
Chemicals*

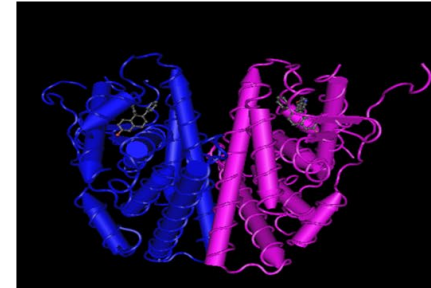


# Sequence

```
MTMTLHTKASGMALLHQIQGNELEPLNRPQLKIPLERPLGE
VYLDSSKPAVYNYPEGAAYEFNAAAAANAQVYGQTGLPYG
PGSEAAAFSGNSLGGFPPLNSVSPSPLMLLHPPQLSPFLQ
PHGQQVPYYLENEPSGYTVREAGPPAFYRPNSDNRRQGGR
ERLASTNDKSGMAMESAKETRYCAVCNDYASGYHYGVWSC
EGCKAFFKRSIQGHNDYMCPTNQCTIDKNRRKSCQACRLR
KCYEVGMMKGIRKDRRGRMLKHKRQRDDGEGRGEVG
SAGDMRAANLWPSPLMIKRSKKNLSLSTADQMVSALLA
EPPILYSEYDTPRPFSEASMMGLLTNLADRELHMINWAKV
PGFVDLTLDQVHLLCAWLEILMIGLVWRSMEHPGKLLFA
PNLLDRNGKQCEVGMVEIFDMLLATSSRFMMNLQGEF
VCLKSILLNSGVYTLSTLSLEEKDHIHRVLDKITDTLIHLM
```



# Structure



**SeqAPASS Results from Level 1**  
Query Sequence FASTA + FASTA from 100s of Aligned Sequences Across Taxa

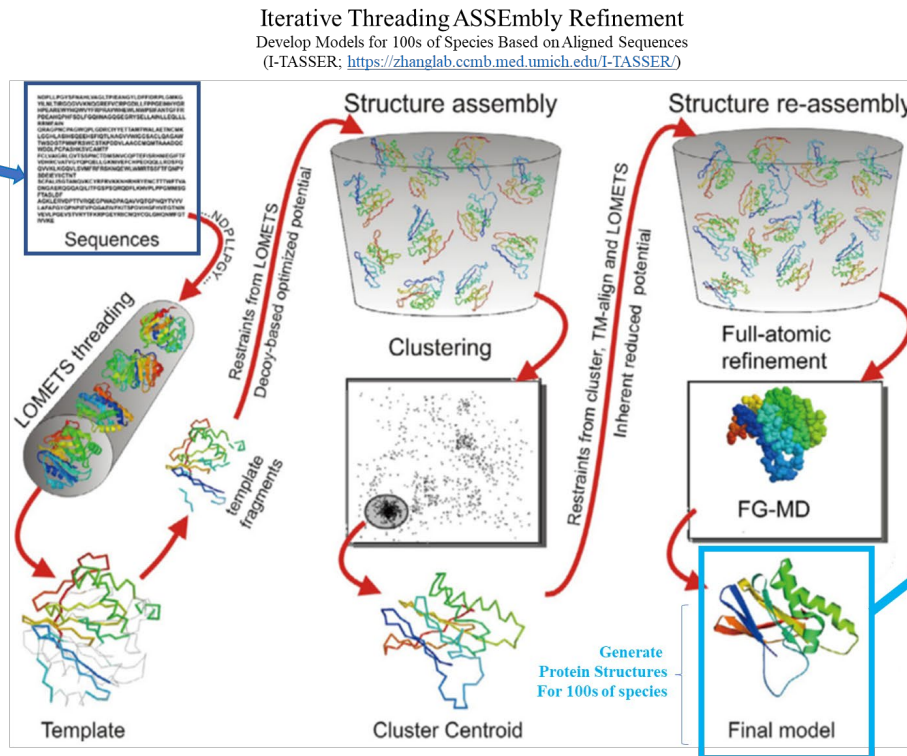
```
>NP_001434.1 Protein X [Homo sapiens]
MSFSGKYQLQSQENFEAFMKAIGLPEELIQKGKDI
KGVSEIVQNGKHKFTITAGSKVIQNEFTVGEECE
LETMTGEKVTVVQLEGDNKLVTFKNIKSVTELN
GDIITNTMTLGDIVFKRISKRI

>NP_787011.1 Protein X [Bos taurus]
MNFSGKYQVQTQENYEFMKAIGLPEELIQKGKDI
KGVSEIVQNGKHKFTITAGSKVIQNEFTVGEECE
MEFMTGEKIKAVVQLEGDNKLVTFKNIKSVTEFN
GDTVSTMTKGDVVKRISKRI

>KFQ76585.1 Protein X [Phoenixcopterus ruber
ruber]
MSFTGKYELQSQENFEAFMKAIGLPEELIQKGKDI
KGVSEIVQNGKHKFTITAGSKVIQNEFTVGEECE
MEFMTGEKIKAVVQLEGDNKLVTFKNIKSVTEFN
GDTVSTMTKGDVVKRISKRI

>NP_001116883.1 Protein X [Xenopus
tropicalis]
MAFAGKYELVHQENFEAFMKAIGLPEELIQKGKDI
KGVSEIVQNGKHKFTITAGSKVIQNEFTVGEECE
LETPTGKVKSVKLEGDNKLVQKAITSTELSG
DTITHVLTNNLVFKRISKRV
```

100s of FASTA

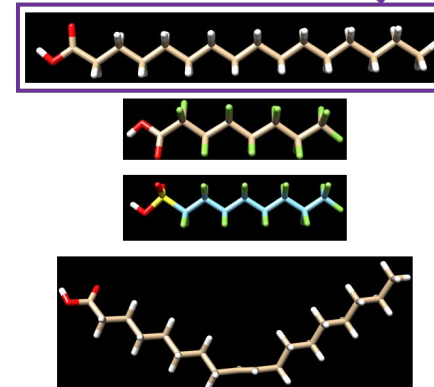


UCSF Chimera  
DockPrep Structures and Minimize Ligands

Protein Structure Models  
From 100s of Species



Ligands of Interest for Docking



AutoDock Vina  
Dock Multiple Ligands to Protein Structures



Collect Predicted Binding Affinity

S	Score	RMSD Lib	RMSD u.b	HBonds (all)	HBond Ligand Atoms	HBond Receptor Atoms
V	-7.1	2.121	2.436	0	0	0
V	-7.0	2.148	6.837	1	1	1
V	-6.9	1.128	2.04	0	0	0
V	-6.9	4.472	7.133	0	0	0
V	-6.7	3.27	7.552	0	0	0
V	-6.7	2.637	3.461	2	2	2
V	-6.6	1.572	3.516	0	0	0
V	-6.6	1.725	3.368	0	0	0

Chimera Model #3.1

REMARK VINA RESULT: -7.1 0.000 0.000

REMARK 15 active torsions:

REMARK status: 'A' for Active; 'I' for Inactive

REMARK 1 A between atoms: C2\_2 and C3\_3

REMARK 2 A between atoms: C3\_3 and C4\_4

REMARK 3 A between atoms: C4\_4 and C5\_5

REMARK 4 A between atoms: C5\_5 and C6\_6

REMARK 5 A between atoms: C6\_6 and C7\_7

REMARK 6 A between atoms: C7\_7 and C8\_8

REMARK 7 A between atoms: C8\_8 and C9\_9

REMARK 8 A between atoms: C10\_10 and C9\_9

REMARK 9 A between atoms: C10\_10 and C11\_11

REMARK 10 A between atoms: C11\_11 and C12\_12

REMARK 11 A between atoms: C12\_12 and C13\_13

REMARK 12 A between atoms: C13\_13 and C14\_14

REMARK 13 A between atoms: C14\_14 and C15\_15

REMARK 14 A between atoms: C15\_15 and C16\_16

REMARK 15 A between atoms: C16\_16 and C17\_17

Government

Industry

# Consortium to Advance Cross Species Extrapolation in Regulation

## Steering Committee:

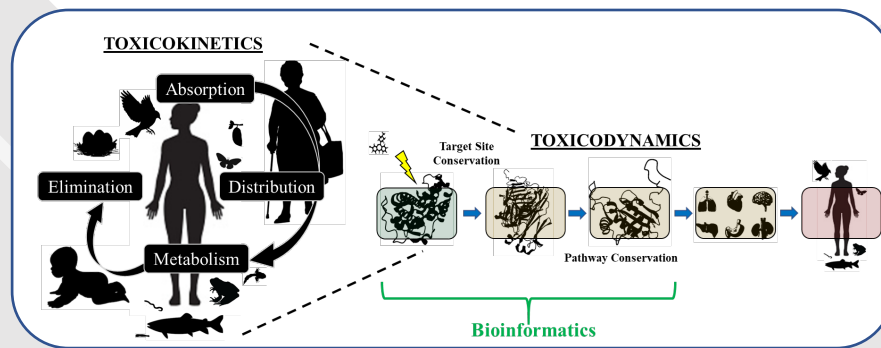
Carlie LaLone (US EPA)  
Geoff Hodges (Unilever)  
Nil Basu (McGill U)  
Steve Edwards (RTI)  
Fiona Sewell (NC3Rs)  
Michelle Embry (HESI)  
Patience Browne (OECD)

1. Define the taxonomic domain of applicability
2. Define the global regulatory landscape/need
3. Develop a bioinformatics toolbox
4. Communicate a shared scientific vision

Interested in Learning more or Joining: Contact [LaLone.Carlie@epa.gov](mailto:LaLone.Carlie@epa.gov) or [Geoff.Hodges@unilever.com](mailto:Geoff.Hodges@unilever.com)

Academia

NGO



# Acknowledgements

## U.S. EPA, ORD

Marissa Jensen (University of Minnesota Duluth)

Sally Mayasich (University of Wisconsin)

Monique Hazemi (ORISE)

Sara Vliet (past ORISE 2021)

Jon Doering (past NRC 2018)

Colin Finnegan (past ORISE 2018)

Donovan Blatz (past ORISE 2021)

## GDIT

Cody Simmons

Audrey Wilkinson

Wilson Menendez

Thomas Transue (past GDIT 2022)

SeqAPASS v6.0 (Released Sept. 2021)



[LaLone.Carlie@epa.gov](mailto:LaLone.Carlie@epa.gov)

<https://seqapass.epa.gov/seqapass/>