

# AI Safety and the Life Sciences

Gopal P. Sarma MD, PhD

Models, Inference, and Algorithms Initiative

Broad Institute of MIT and Harvard



# Complex Interplay Between AI and the Life Sciences

- **Claim 1:** The life sciences will drive the next decade of creativity in AI
  - Economic and social incentives
  - Rich multi-modal data scaling at high exponential rates
  - Challenging biomedical problems
- **Claim 2:** Integrating AI safety and the life sciences is critical and will require significant transdisciplinary efforts

# Complex Interplay Between AI and the Life Sciences

- **Claim 1:** The life sciences will drive the next decade of creativity in AI
  - Economic and social incentives



# 2019 Global Life Sciences Outlook

Focus and transform | Accelerating change in life sciences

Global health care spending continues to increase dramatically



is projected to reach

**\$10.059 trillion by 2022**

# THE KEY TO SUCCESSFUL DIGITAL TRANSFORMATION IN LIFE SCIENCES



**Deloitte.**



# **TALES OF TRANSFORMATION**

**Regulatory's impact in life sciences  
digital transformation**



# Digital Transformation for Life Sciences

Cisco Healthcare is empowering innovation. Our technologies make it possible to improve customer engagement, reduce discovery and development time, transform the global supply chain, and address changes in the global and regulatory environment, all while enabling security and compliance.



## R&D and Clinical Trials

- Collaboration – video, voice, and virtual meetings
- Cloud and data center
- Data virtualization
- Analytics
- Asset management
- Customer-contact center

## Manufacturing and Supply Chain

- Unified, converged factory network
- Connected machines and edge computing
- Centralized, pervasive wireless
- Analytics
- Services-exchange platform
- Asset management

## Marketing and Sales Engagement

- Collaboration – video, voice, and virtual meetings
- Automated response bots
- Customer-contact center
- Analytics

Security and DNA for Life Sciences



An aerial night view of a city skyline, likely San Francisco, with prominent skyscrapers like the Transamerica Pyramid. The image is overlaid with a blue gradient that is darker at the top and lighter towards the bottom. The text is centered in the upper half of the image.

A GUIDE FOR LIFE SCIENCES COMPANIES

---

# How to Embrace Digital Transformation

[READ NOW](#)



The background of the slide features a teal-to-blue gradient. Overlaid on this are faint, semi-transparent images of laboratory glassware, including a large Erlenmeyer flask on the left, a graduated cylinder in the center, and a round-bottom flask on the right. A faint molecular structure, possibly a hexagonal ring, is also visible in the upper left quadrant.

# Top 3 Digital Transformation Challenges in Life Sciences

[LEARN MORE](#)

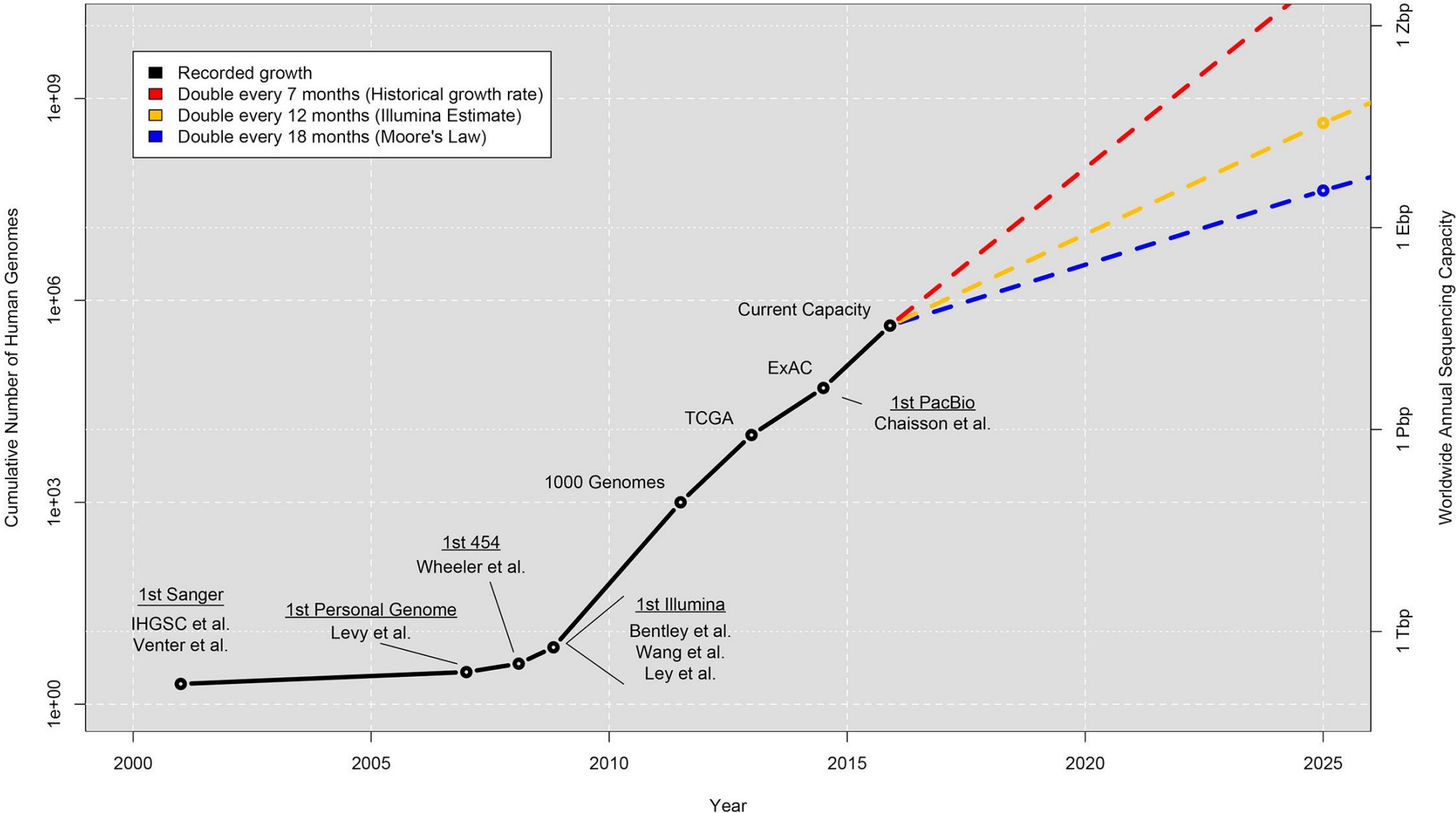
The logo for ADLIB, featuring a stylized green plus sign followed by the word "ADLIB" in white, uppercase letters, with a registered trademark symbol (®) to the right.

# Complex Interplay Between AI and the Life Sciences

- **Claim 1:** The life sciences will drive the next decade of creativity in AI
  - Economic and social incentives
  - Rich multi-modal data scaling at high exponential rates
  - Challenging biomedical problems

# Scaling of DNA Sequencing

Stephens ZD et al. (2015) Big Data: Astronomical or Genomical? PLOS Biology 13(7): e1002195.



# Bio-inspired AI: Uncovering the Genomic Basis of Disease



<https://hail.is>

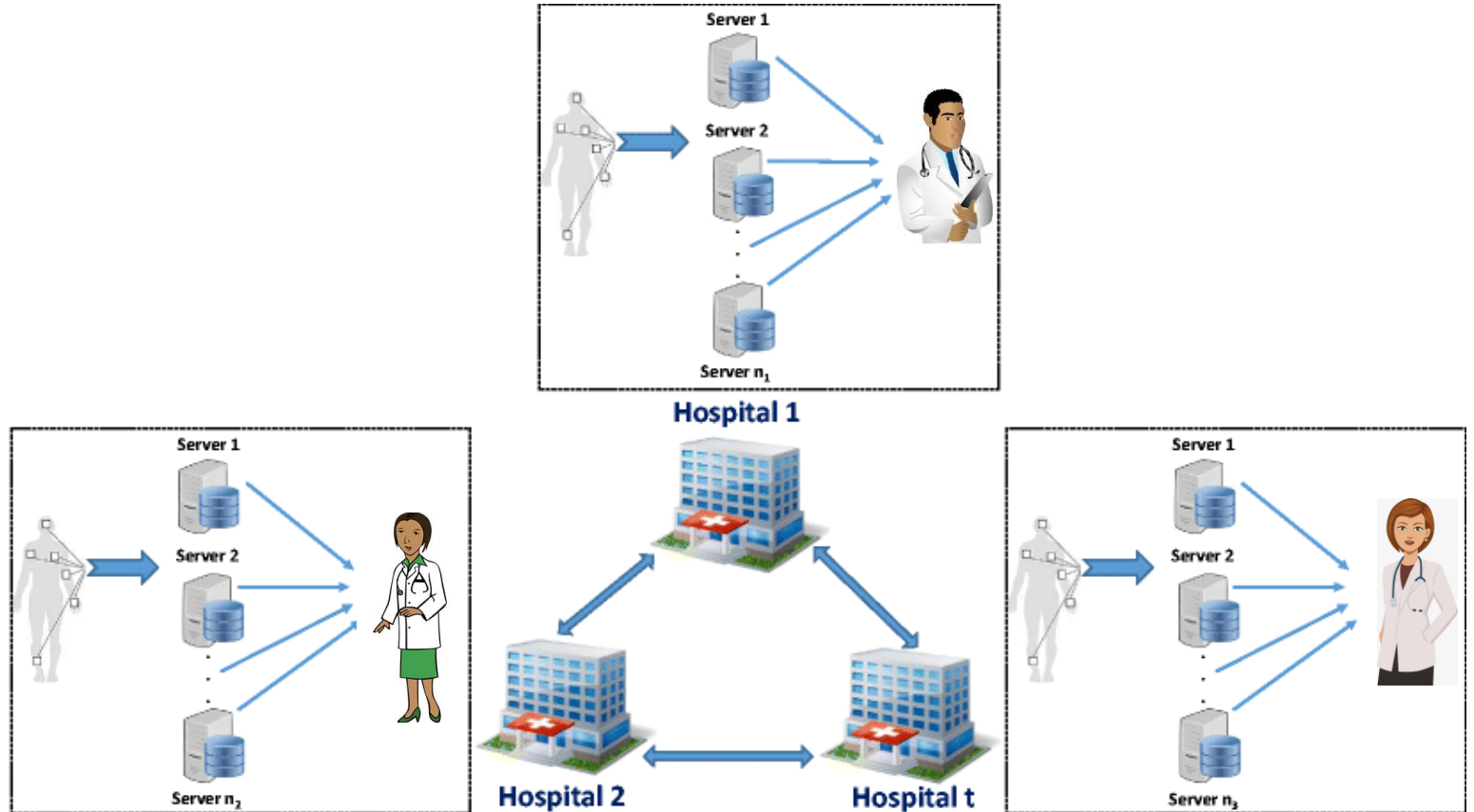
- Open-source, genomic analysis platform
- MatrixTable: first-class support for multi-dimensional structured data, built from the ground-up as a distributed system running in the cloud
- The structure, scale, and analyses of genomic data inspired a new distributed data structure that is broadly useful in biology and beyond



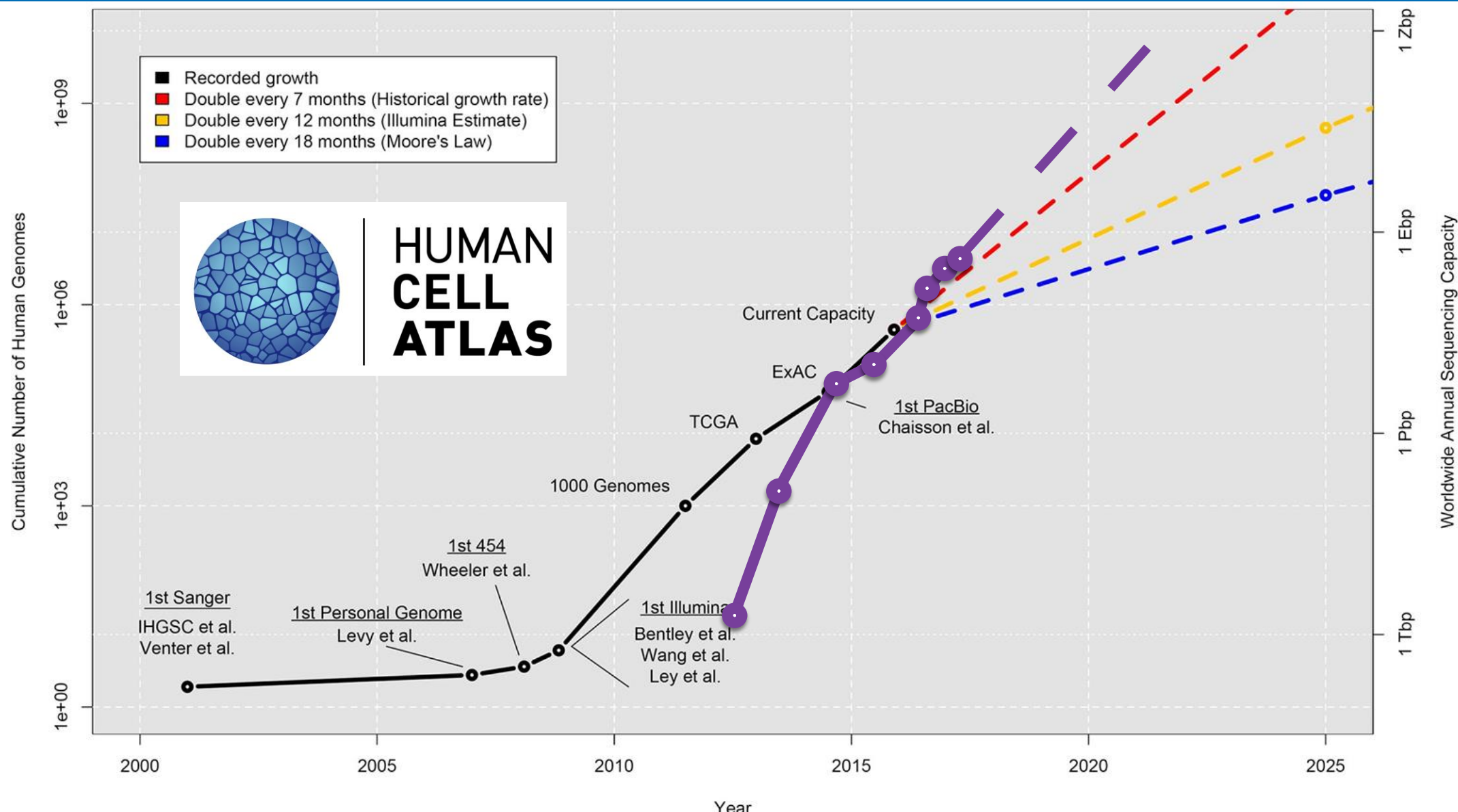
# Bio-Inspired AI: Secure Multi-Party Linear Regression

J. Bloom "Secure multi-party linear regression at plaintext speed." *arXiv preprint* arXiv:1901.09531 (2019).  
B. Berger and H. Cho "Emerging technologies towards enhancing privacy in genomic data sharing" *Genome Biology* (2019)

- Secure multi-party regression at plain text speed
- Addresses issues of data ownership, transmission, and storage at scale



# Scaling of Cells Profiled at the Klarman Cell Observatory







# Bio-Inspired AI: Representation Learning of scRNA Data

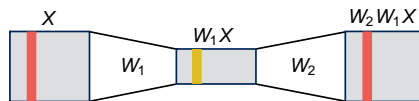


Loss Landscapes of Regularized Linear Autoencoders  
Daniel Kunin<sup>1,2</sup>, Jonathan M. Bloom<sup>1</sup>, Aleksandrina Goeva<sup>1</sup>, Cotton Seed<sup>1</sup>



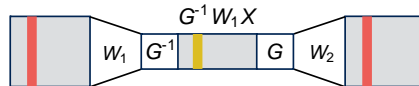
## Background

A linear autoencoder maps  $\mathbb{R}^m \rightarrow \mathbb{R}^k \rightarrow \mathbb{R}^m$ .



$$L(W_1, W_2) = \|X - W_2 W_1 X\|^2$$

LAEs learn the top principal *subspace* but not the principal *directions* or eigenvalues. The optimal latent representation is only defined up to a linear map  $G \in GL_k(\mathbb{R})$ .



LAEs are pseudoinverses at all critical pts.

## Regularization

We prove that  $L^2$ -regularized LAEs are transposes at all critical points and learn the principal directions as the left singular vectors of the decoder. Define  $L_\sigma$  by

$$\|X - W_2 W_1 X\|^2 + \lambda(\|W_1\|^2 + \|W_2\|^2)$$

The minima of  $L_\sigma$  are defined up to an orthogonal map  $O \in O_k(\mathbb{R})$  by

$$W_2 = U_k(I - \lambda^{-1} \Sigma_k^{-2})^{\frac{1}{2}} O = W_1^T$$

where  $X = U \Lambda V^T$  and  $\sigma_1^2 > \dots > \sigma_k^2 > \lambda$ .

$$W_2 W_1 = U_k(I - \lambda^{-1} \Sigma_k^{-2}) U_k^T$$

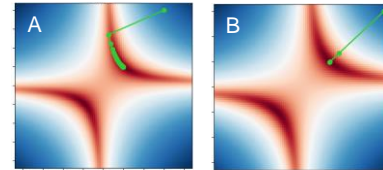
## PCA Algorithms

Hence PCA is a two-step optimization:

1. Train  $L^2$ -regularized LAE on  $X \rightarrow \mathbb{R}^{m \times n}$ .
2. Apply SVD to the decoder  $W_2 \rightarrow \mathbb{R}^{m \times k}$ .

Step 2 is quick. Step 1 options include:

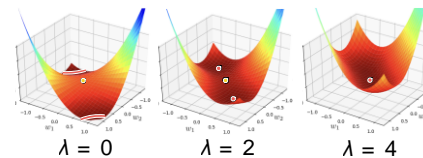
- A. Gradient descent (below).
- B. Solve for  $W_2$ , set  $W_1 = W_2^T$ , iterate.



input  $X \in \mathbb{R}^{m \times n}$ ;  $k \leq m$ ;  $\lambda, \alpha > 0$   
initialize  $W_1, W_2^T \in \mathbb{R}^{k \times m}$   
while not converged  
     $W_1 \leftarrow \alpha(W_2^T(W_2 W_1 - I)X X^T + \lambda W_1)$   
     $W_2 \leftarrow \alpha((W_2 W_1 - I)X X^T W_1^T + \lambda W_2)$   
     $U, \Sigma, V = \text{SVD}(W_2)$   
return  $U, \lambda(I - \Sigma^2)^{-1}$

## Posterior Collapse

Principal directions with eigenvalues below  $\lambda$  collapse as in probabilistic PCA.



Example of collapse for  $X = [2]$ .

## Symmetry and Backprop

$L^2$ -reg LAEs are symmetric at all critical pts.

**Theorem 2.1** (Transpose Theorem). All critical points of  $L_\sigma$  satisfy  $W_1 = W_2^T$ .

*Proof.* Critical points of  $L_\sigma$  satisfy:

$$\begin{aligned} \frac{\partial L_\sigma}{\partial W_1} &= 2W_2^T(W_2 W_1 - I)X X^T + 2\lambda W_1 = 0, \\ \frac{\partial L_\sigma}{\partial W_2} &= 2(W_2 W_1 - I)X X^T W_1^T + 2\lambda W_2 = 0. \end{aligned}$$

We first prove that the matrix

$$C = (I - W_2 W_1)X X^T$$

is positive semi-definite<sup>8</sup>. Rearranging  $\frac{\partial L_\sigma}{\partial W_2} W_2^T$  gives

$$X X^T (W_2 W_1)^T = (W_2 W_1)X X^T (W_2 W_1)^T + \lambda W_2 W_2^T.$$

Both terms on the right are positive semi-definite, so their sum on the left is as well and therefore

$$X X^T (W_2 W_1)^T \succeq (W_2 W_1)X X^T (W_2 W_1)^T.$$

Cancelling  $(W_2 W_1)^T$  via Lemma B.1 gives  $C \succeq 0$ .

We now show the difference  $A = W_1 - W_2^T$  is zero. Rearranging terms using the symmetry of  $C$  gives

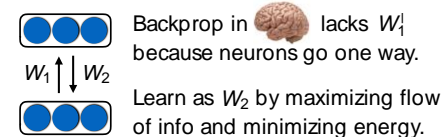
$$0 = \frac{\partial L_\sigma}{\partial W_1} - \frac{\partial L_\sigma}{\partial W_2}^T = 2A(C + \lambda I).$$

Since  $C \succeq 0$  and  $\lambda > 0$  imply  $C + \lambda I \succ 0$ , we conclude from

$$A(C + \lambda I)A^T = 0$$

that  $A = 0$ .  $\square$

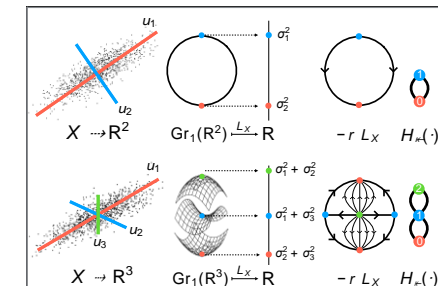
Resolution to weight transport problem:



## Algebraic Topology

We smoothly parameterize the critical manifolds of LAEs with several forms of regularization via one elementary proof.

We factor the loss as a Morse function on the Grassmannian to reveal the dynamics near and between critical manifolds. Morse homology suggests principles and algorithms for deep learning.



**Theorem 4.4** (Curvature Theorem). In local coordinates near any point on the critical manifold indexed by  $I$ , all three losses take the form of a standard degenerate saddle with  $d_I + (k - \ell)(m - \ell)$  descending directions.

- $L$  and  $L_\pi$  have  $k\ell$  flat directions.
- $L_\sigma$  has  $k\ell - \binom{k+\ell}{2}$  flat directions.

The remaining directions are ascending.

**Theorem E.1.**  $L_X$  is an  $\mathbb{R}_2$ -perfect Morse function. Its critical points are the rank- $k$  principal subspaces.

*Proof.* Consider the commutative diagram

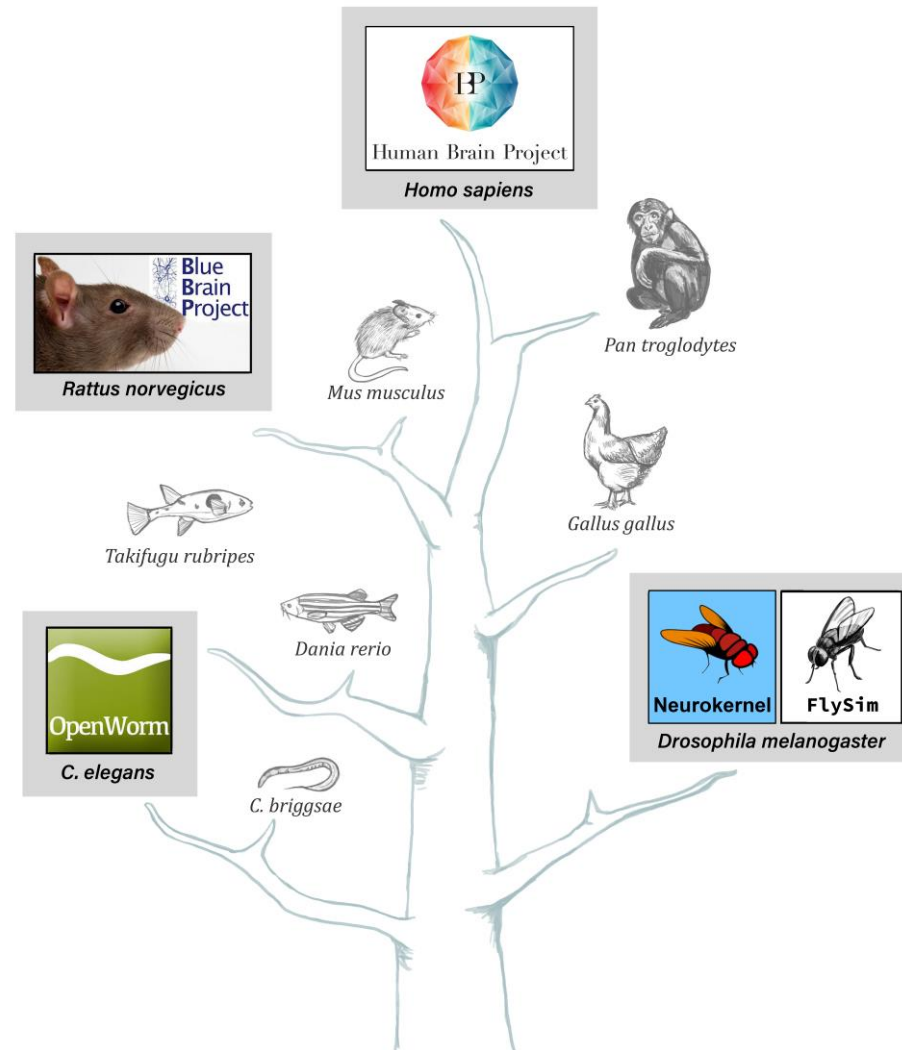
$$\begin{array}{ccc} V_k(\mathbb{R}^m) & \xrightarrow{\pi \circ \text{Im}(OO^T)} & \text{Gr}_k(\mathbb{R}^m) \\ \downarrow \iota \circ O \circ (O^T, O) & & \downarrow L_X \\ \mathbb{R}^{k \times m} \times \mathbb{R}^{m \times k} & \xrightarrow{L} & \mathbb{R} \end{array} \quad (10)$$



# Bio-Inspired AI: Realistic Nervous System Simulations

G. Sarma, A. Safron, and N. Hay, "Integrative Biological Simulation, Neuropsychology, and AI Safety." *AAAI Workshop on AI Safety* (2019)

- Simple organisms show complex behavior that continues to be difficult for modern AI systems.
- Neuronal simulations in virtual environments will allow these biological architectures to be used for AI research.



NEURON

BluePyOpt

NetPyNE

Bionet

Geppetto

ChannelPedia

NeuroMLDB





Emerald Cloud Lab



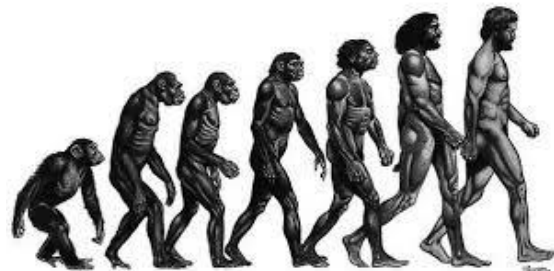
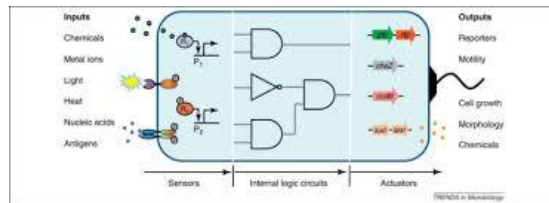
GINKGO  
BIOWORKS  
THE ORGANISM COMPANY



NOVARTIS



Y Combinator



# Complex Interplay Between AI and the Life Sciences

- **Claim 2:** Integrating AI safety and the life sciences is critical and will require significant transdisciplinary efforts

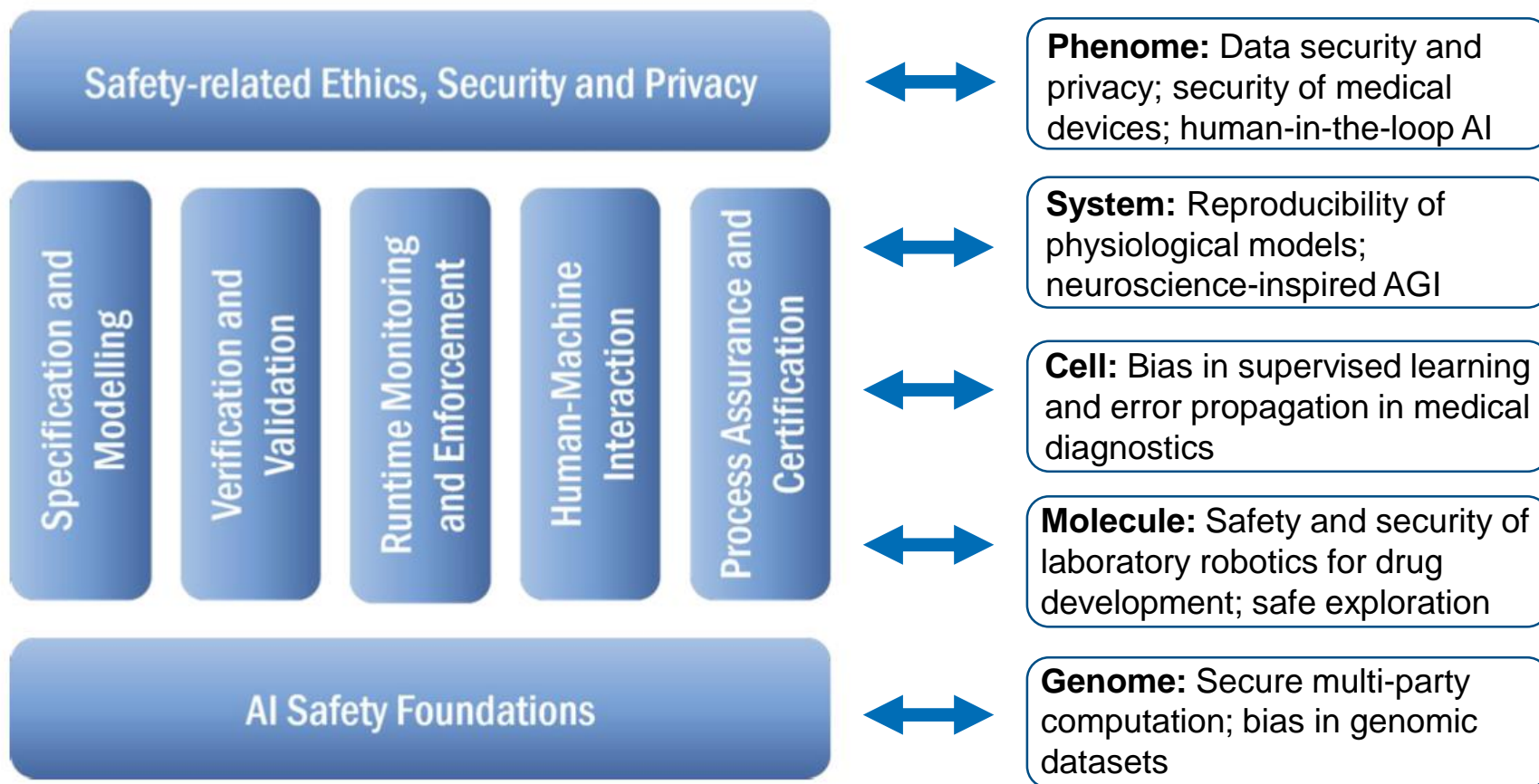
# Complex Interplay Between AI and the Life Sciences

- AI will infuse genomics, cell biology, pharmaceuticals, medical devices, and care delivery
- Nearly all areas of AI safety are relevant to the life sciences
- A **culture of safety** is needed to address a diversity of safety challenges
- Biosecurity and biosafety are models for policy and cultural integration
  - 1975 Asilomar Conference on Recombinant DNA inspired Conference on Beneficial AI in 2017
  - Scientific community actively debating ethics of editing human genomes



# Prioritizing the Life Sciences in the AI Safety Landscape

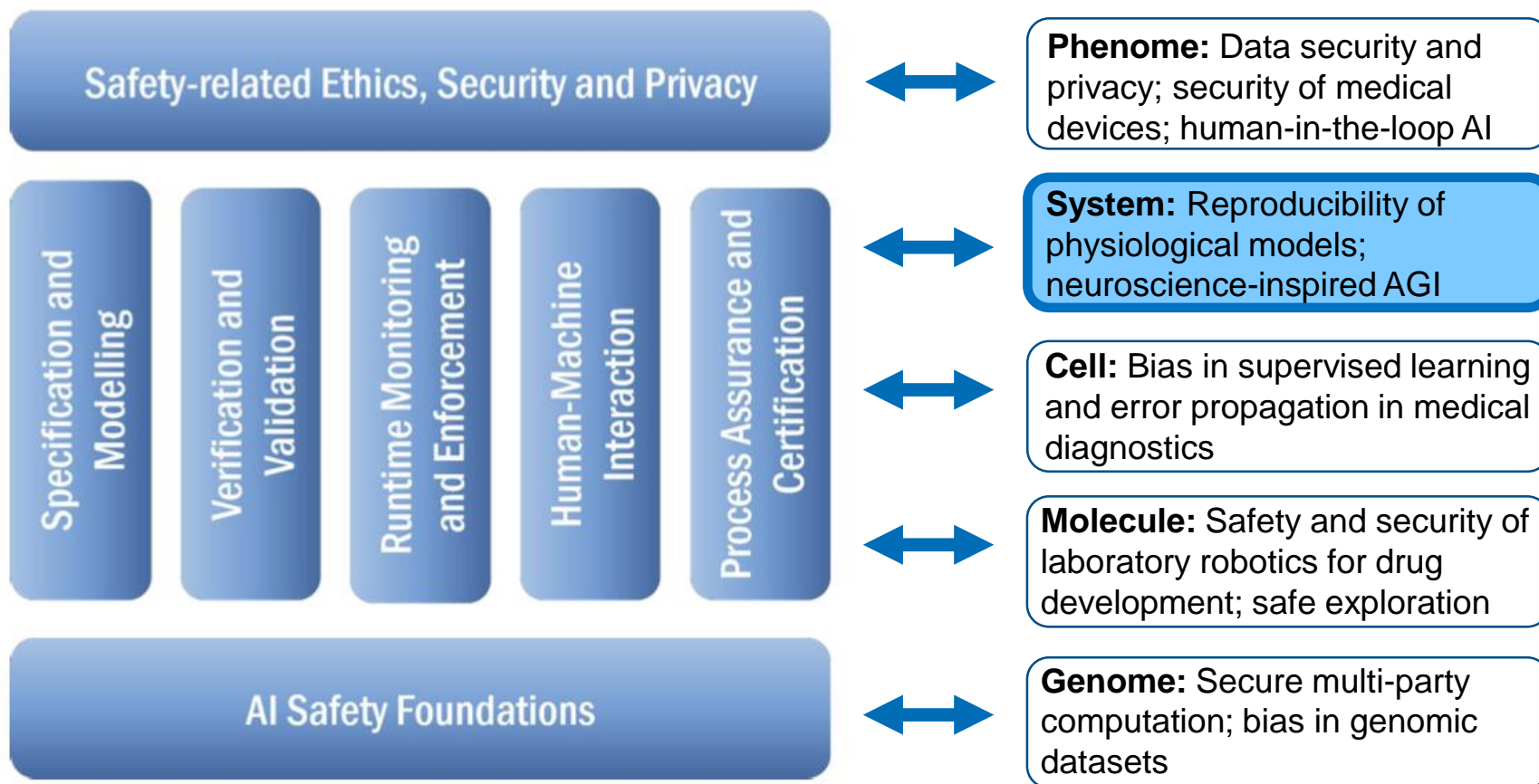
## Life Sciences and AI Safety: A Unique Interplay of Fields





# Prioritizing the Life Sciences in the AI Safety Landscape

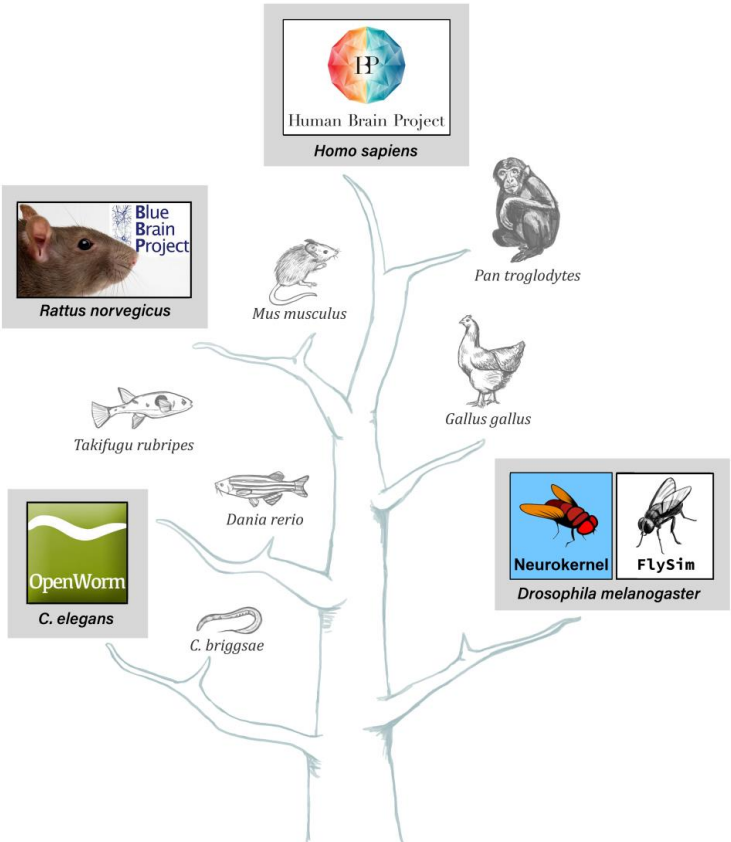
## Life Sciences and AI Safety: A Unique Interplay of Fields





# Prioritizing the Life Sciences in the AI Safety Landscape

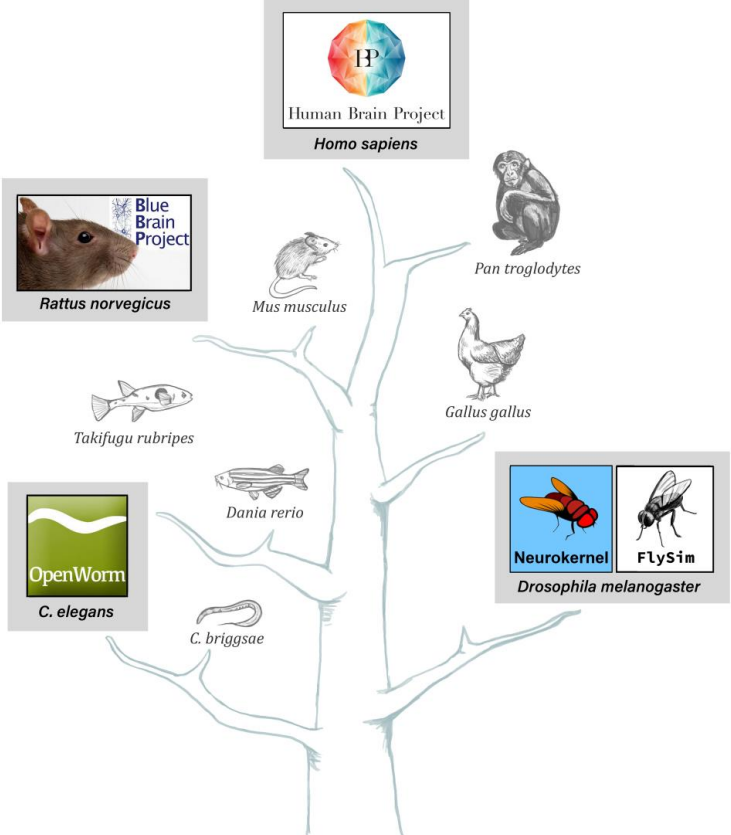
G. Sarma, A. Safron, and N. Hay, "Integrative Biological Simulation, Neuropsychology, and AI Safety." *AAAI Workshop on AI Safety* (2019)  
J. Leike, et al. "AI Safety Gridworlds." *arXiv preprint arXiv:1711.09883* (2017).

Biological Agent and Environment Simulations	Grid Worlds Problems
 <p>A phylogenetic tree diagram illustrating the relationship between various biological agents and their corresponding simulation projects. The tree branches from a common ancestor at the bottom to various species at the top. The species and their associated simulation projects are:</p> <ul style="list-style-type: none"><li><b>Human Brain Project</b> (Homo sapiens)</li><li><b>Blue Brain Project</b> (Rattus norvegicus)</li><li><b>Mus musculus</b></li><li><b>Pan troglodytes</b></li><li><b>Gallus gallus</b></li><li><b>Takifugu rubripes</b></li><li><b>Dania rerio</b></li><li><b>OpenWorm</b> (C. elegans)</li><li><b>C. briggsae</b></li><li><b>Neurokernel</b> (Drosophila melanogaster)</li><li><b>FlySim</b> (Drosophila melanogaster)</li></ul>	<ul style="list-style-type: none"><li>• Avoiding negative side effects</li><li>• Avoiding reward hacking</li><li>• Safe exploration</li><li>• Safe interruptibility</li><li>• Robustness to distributional shift</li></ul>



# Prioritizing the Life Sciences in the AI Safety Landscape

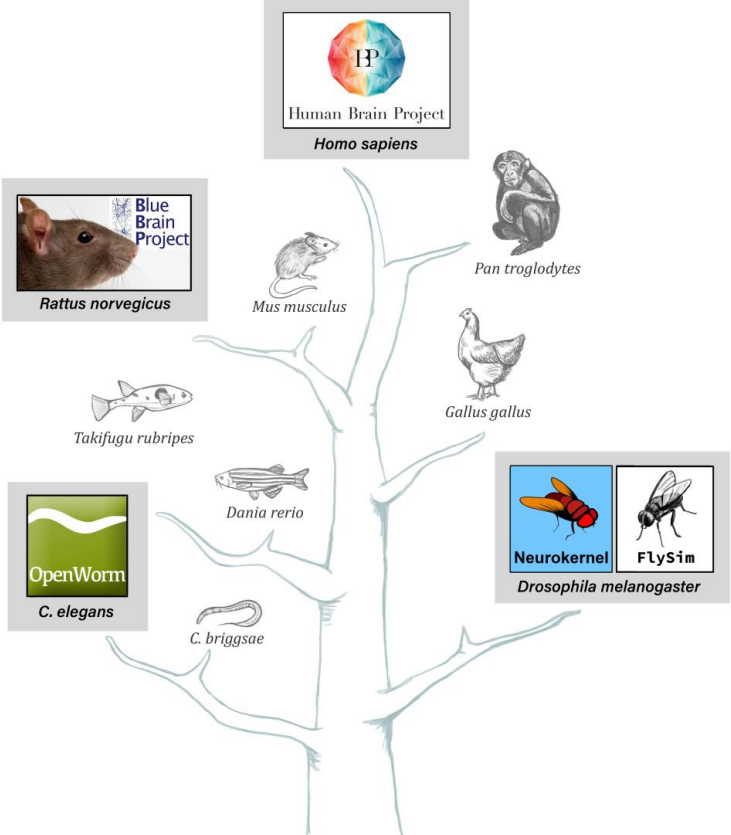
G. Sarma, A. Safron, and N. Hay, "Integrative Biological Simulation, Neuropsychology, and AI Safety." *AAAI Workshop on AI Safety* (2019)  
J. Leike, et al. "AI Safety Gridworlds." *arXiv preprint arXiv:1711.09883* (2017).

Biological Agent and Environment Simulations	Grid Worlds Problems
	<ul style="list-style-type: none"><li>• Avoiding negative side effects</li><li>• <b>Avoiding reward hacking</b></li><li>• Safe exploration</li><li>• Safe interruptibility</li><li>• Robustness to distributional shift</li></ul>



# Prioritizing the Life Sciences in the AI Safety Landscape

G. Sarma, A. Safron, and N. Hay, "Integrative Biological Simulation, Neuropsychology, and AI Safety." *AAAI Workshop on AI Safety* (2019)  
J. Leike, et al. "AI Safety Gridworlds." *arXiv preprint arXiv:1711.09883* (2017).

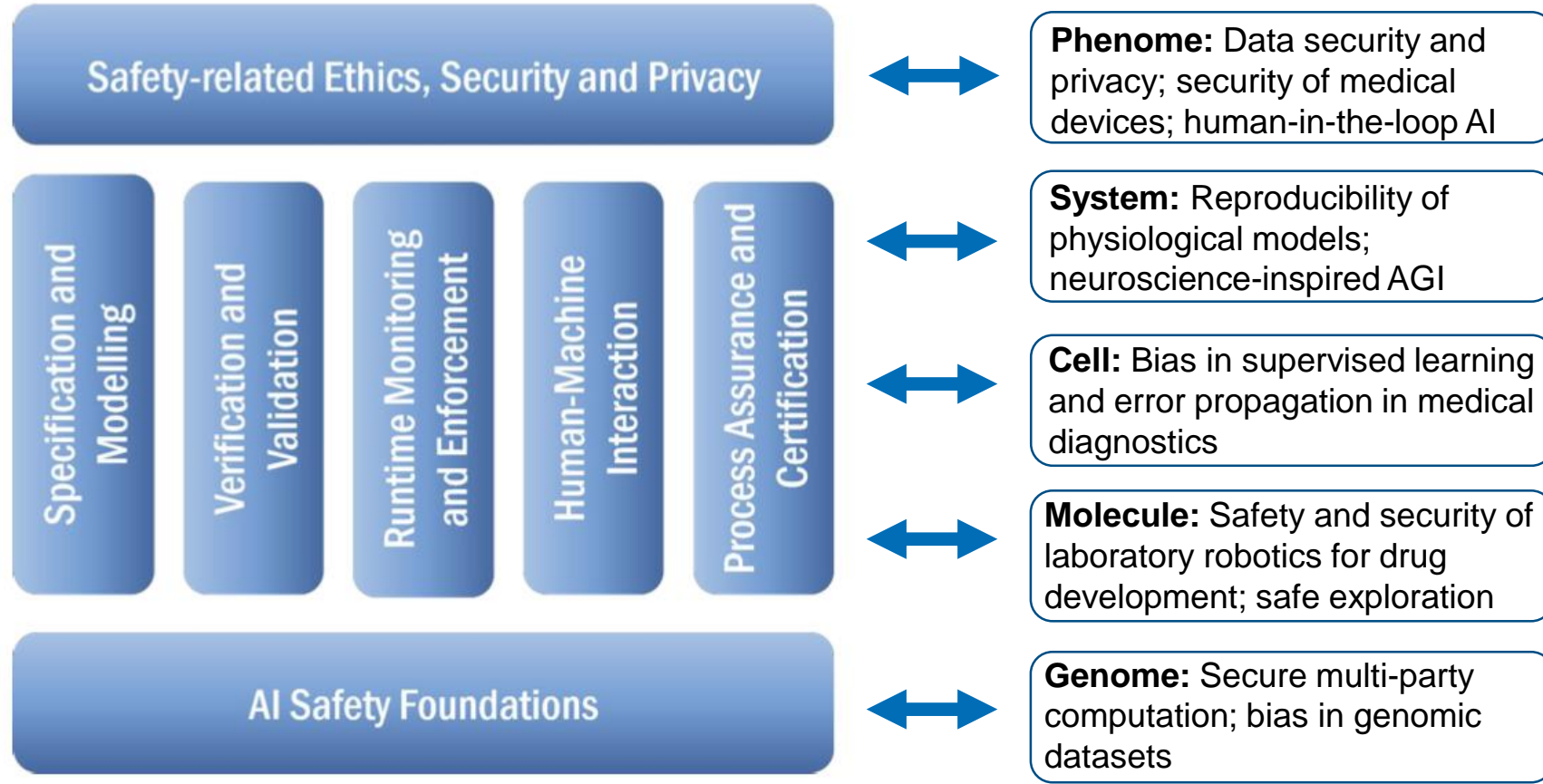
Biological Agent and Environment Simulations	Grid Worlds Problems
	<ul style="list-style-type: none"><li>• Avoiding negative side effects</li><li>• Avoiding reward hacking</li><li>• <b>Safe exploration</b></li><li>• Safe interruptibility</li><li>• Robustness to distributional shift</li></ul>

# Summary

- **Claim 1:** Bio-inspired AI: life sciences will advance AI itself
- **Claim 2:** Bio-inspired AI safety: the life sciences will advance AI safety itself

# Summary

## Life Sciences and AI Safety: A Unique Interplay of Fields





<https://broadinstitute.org/mia>

