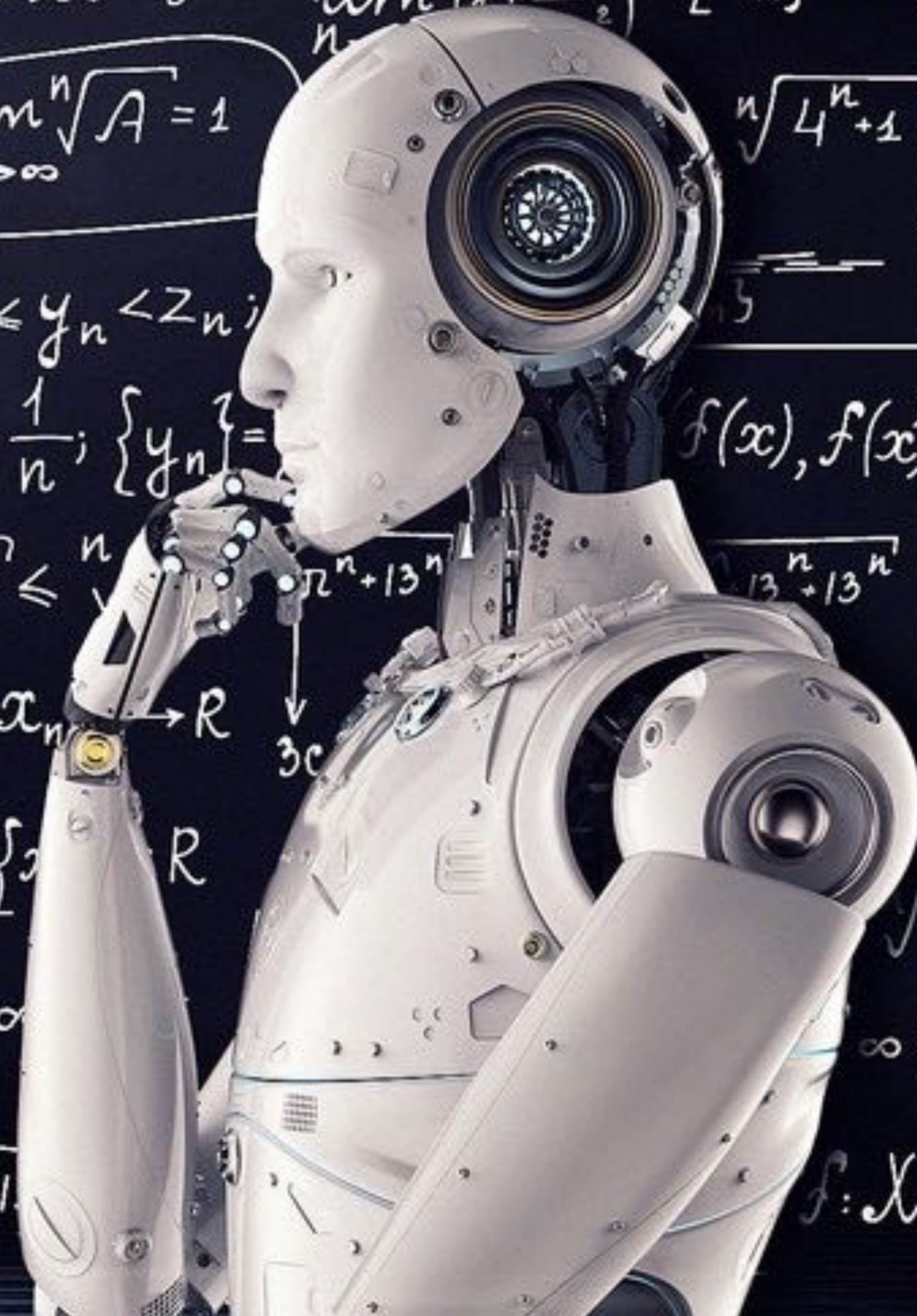


# Deep Learning for Science Opportunities & Challenges

Maxwell Cai, Ph.D. (SURFsara)



Background mathematical content includes:

- Sequences:  $\{x_n\} \subset \mathbb{R}$ ,  $\lim_{n \rightarrow \infty} \frac{n^2 - x}{3}$ ,  $\lim_{n \rightarrow \infty} \sqrt[n]{A} = 1$ ,  $\lim_{n \rightarrow \infty} (1 + \frac{\pi}{n})$
- Graphs: A sine wave graph with labels "lim min" and "lok. min".
- Equations:  $f(x) \Leftrightarrow \exists q \in [0, 1): \forall x, x' \in X$ ,  $x_n - g < \epsilon \quad n \geq n_0: (x_n - g) < \epsilon$ ,  $x_n: \mathbb{N} \rightarrow \mathbb{R}$ ,  $\frac{1}{1 + \frac{1}{n}} = \frac{1}{\frac{n+1}{n}}$ ,  $\{x_n\} + \{y_n\} \stackrel{df}{=} \{x_n + y_n\}; \mathbb{R}$ ,  $\{x_n\} \cdot \{y_n\} \stackrel{df}{=} \{x_n \cdot y_n\}; \mathbb{R}$ ,  $x_n \leq y_n \leq z_n$ ,  $\lim_{n \rightarrow \infty} x_n = g$
- Other:  $\lim_{n \rightarrow \infty} \sqrt[n]{4^n + 1}$ ,  $\lim_{n \rightarrow \infty} \sqrt[n]{0 + 0 + 0 + 13^n} \leq \sqrt[n]{13^n}$ ,  $\lim_{n \rightarrow \infty} \sqrt[n]{13^n} = 13$

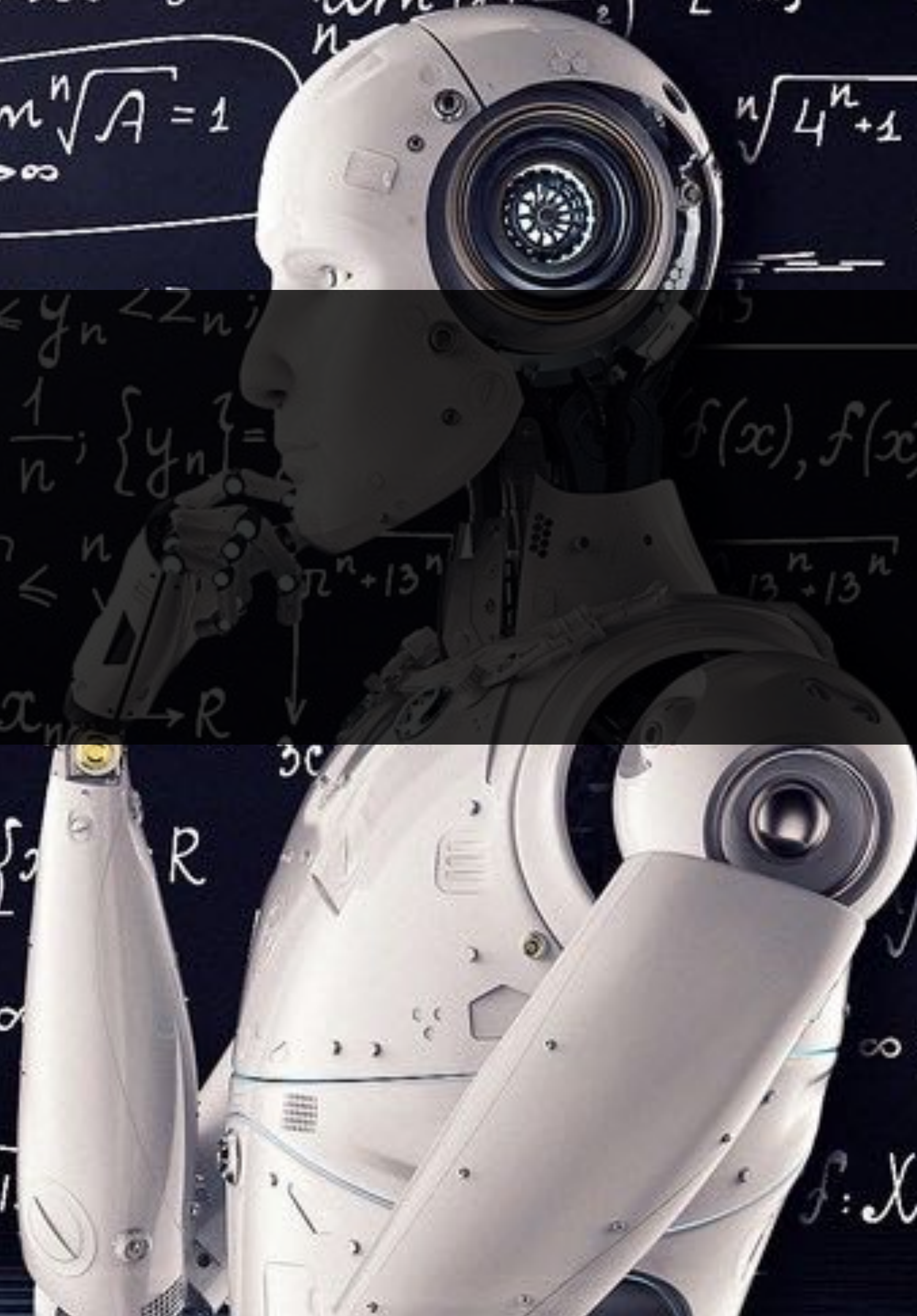


# Deep Learning for Science Opportunities & Challenges

Maxwell Cai, Ph.D. (SURFsara)

## Analytical models

*Elegant & Abstract*

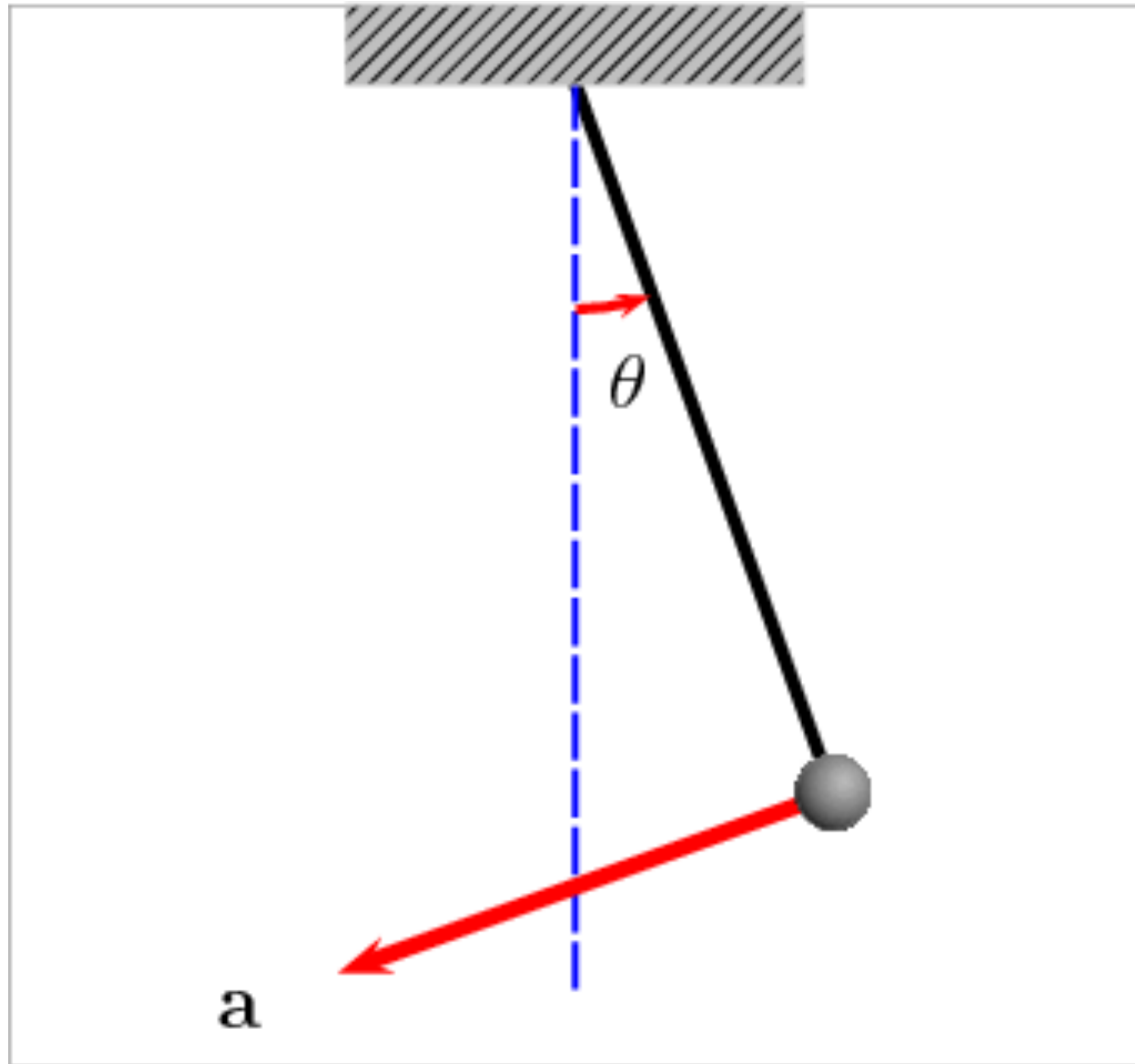


$$\left\{ \frac{1}{n} \right\} = \left\{ \frac{1}{n} \right\}$$
$$x_n: \mathbb{N} \rightarrow \mathbb{R}$$
$$x_n \leq y_n \leq z_n$$
$$\downarrow n \rightarrow \infty \quad \downarrow n \rightarrow \infty$$
$$g \quad g$$

$$\{x_n\} + \{y_n\} \stackrel{\text{df}}{=} \{x_n + y_n\}; \mathbb{R}$$
$$\{x_n\} \cdot \{y_n\} \stackrel{\text{df}}{=} \{x_n \cdot y_n\}; \mathbb{R}$$
$$\lim_{n \rightarrow \infty} \min \quad \lim_{n \rightarrow \infty} \min$$
$$n\sqrt{4} \cdot n\sqrt{13^n} \cdot n\sqrt{13^n}$$
$$\parallel \{x_n\} \cdot \{y_n\} \stackrel{\text{df}}{=} \{x_n \cdot y_n\}; \mathbb{R}$$

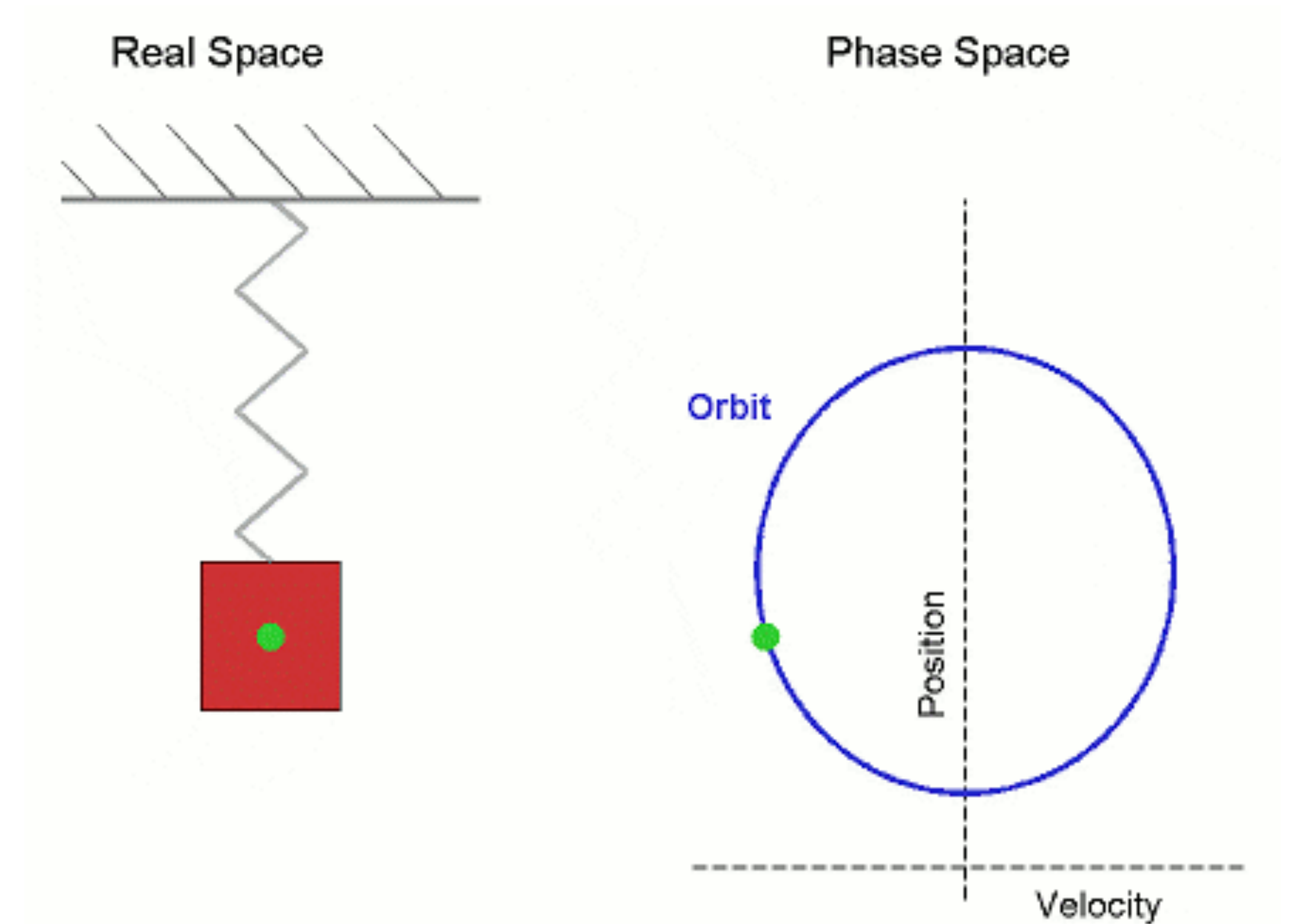


# Analytical model: pendulum & harmonic oscillator

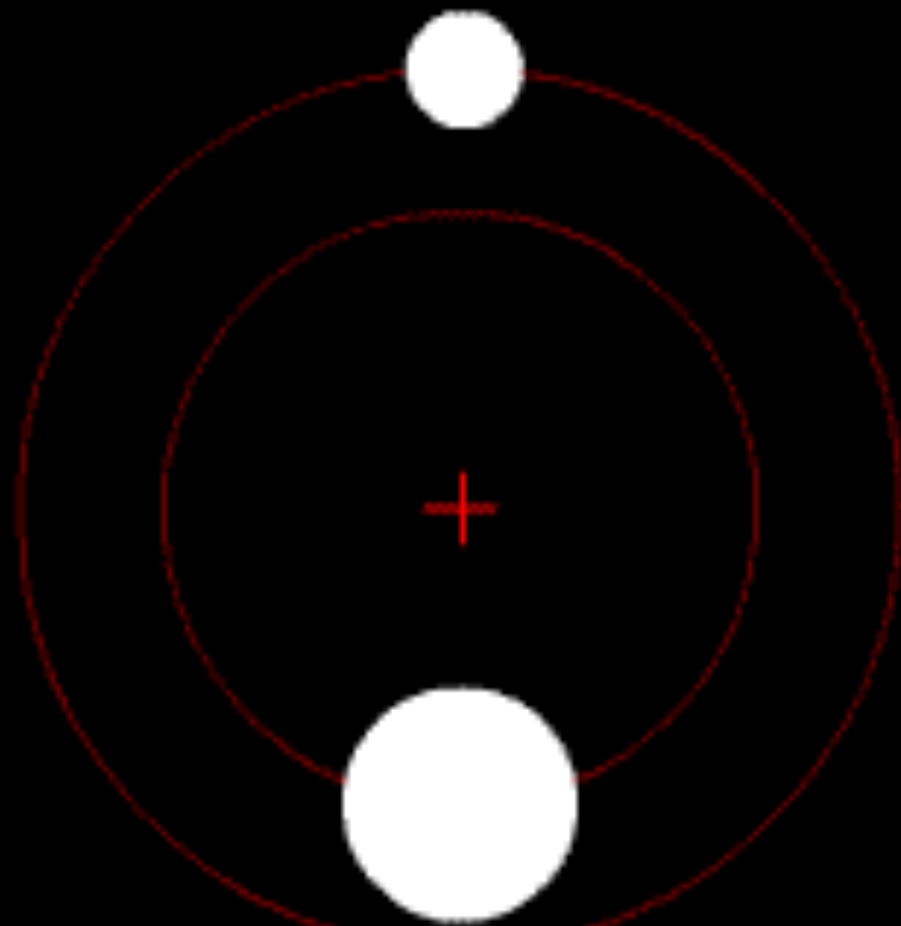


$$\frac{d^2\theta}{dt^2} + \frac{g}{l} \sin \theta = 0$$

$$T_0 \approx 2\pi \sqrt{\frac{l}{g}}$$



# Analytical model: two-body problem



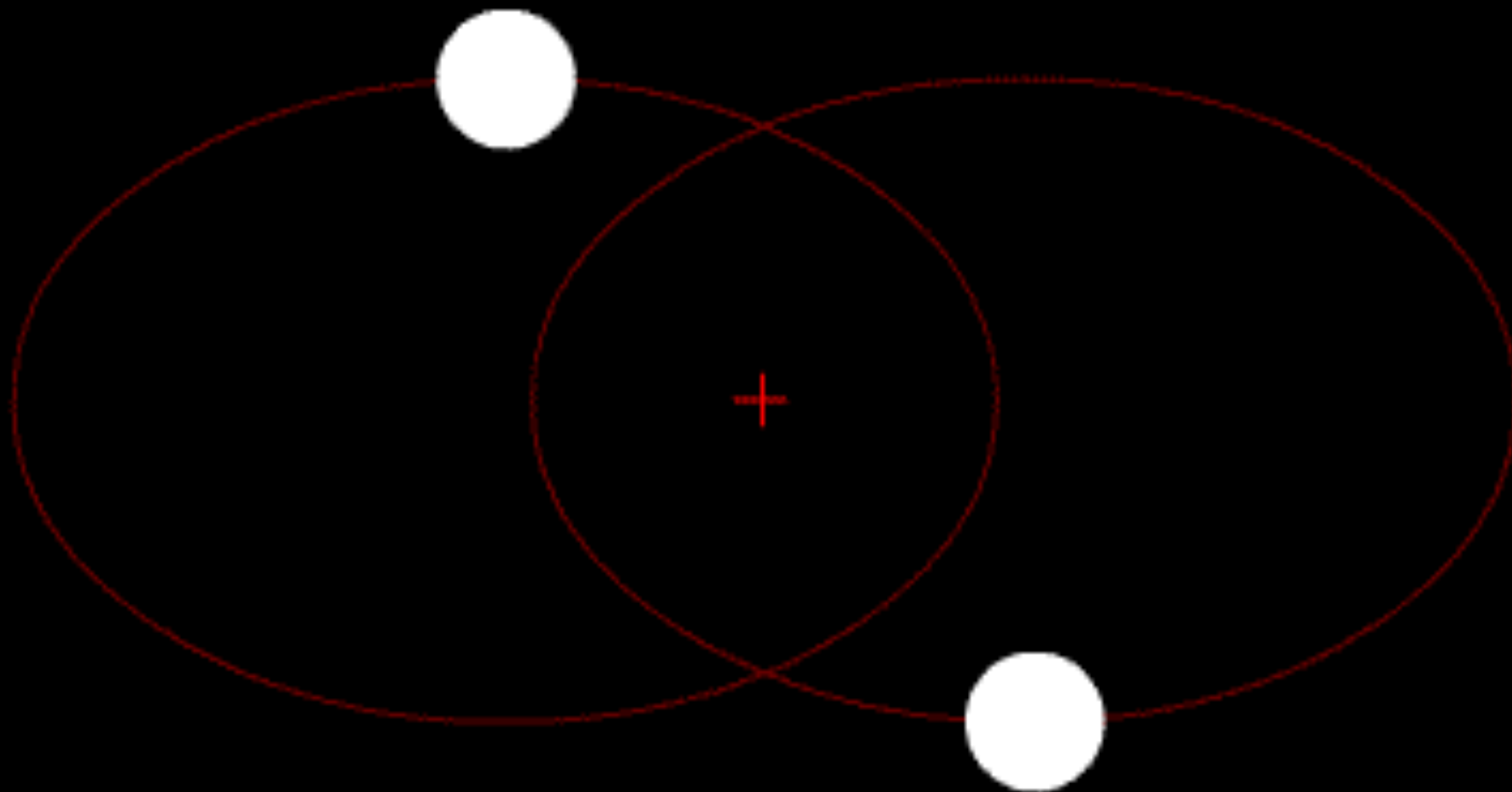
$$\mathbf{F}_{12}(\mathbf{x}_1, \mathbf{x}_2) = m_1 \ddot{\mathbf{x}}_1$$

$$\mathbf{F}_{21}(\mathbf{x}_1, \mathbf{x}_2) = m_2 \ddot{\mathbf{x}}_2$$

$$\ddot{\mathbf{R}} \equiv \frac{m_1 \ddot{\mathbf{x}}_1 + m_2 \ddot{\mathbf{x}}_2}{m_1 + m_2} = 0$$

$$\mu = \frac{1}{\frac{1}{m_1} + \frac{1}{m_2}} = \frac{m_1 m_2}{m_1 + m_2}$$

$$\ddot{\mathbf{r}} = \ddot{\mathbf{x}}_1 - \ddot{\mathbf{x}}_2 = \left( \frac{\mathbf{F}_{12}}{m_1} - \frac{\mathbf{F}_{21}}{m_2} \right) = \left( \frac{1}{m_1} + \frac{1}{m_2} \right) \mathbf{F}_{12}$$



$$\mathbf{x}_1(t) = \mathbf{R}(t) + \frac{m_2}{m_1 + m_2} \mathbf{r}(t)$$

$$\mathbf{L} = \mathbf{r} \times \mathbf{p} = \mathbf{r} \times \mu \frac{d\mathbf{r}}{dt}$$

$$\mathbf{x}_2(t) = \mathbf{R}(t) - \frac{m_1}{m_1 + m_2} \mathbf{r}(t)$$

$$\mathbf{N} = \frac{d\mathbf{L}}{dt} = \dot{\mathbf{r}} \times \mu \dot{\mathbf{r}} + \mathbf{r} \times \mu \ddot{\mathbf{r}}$$

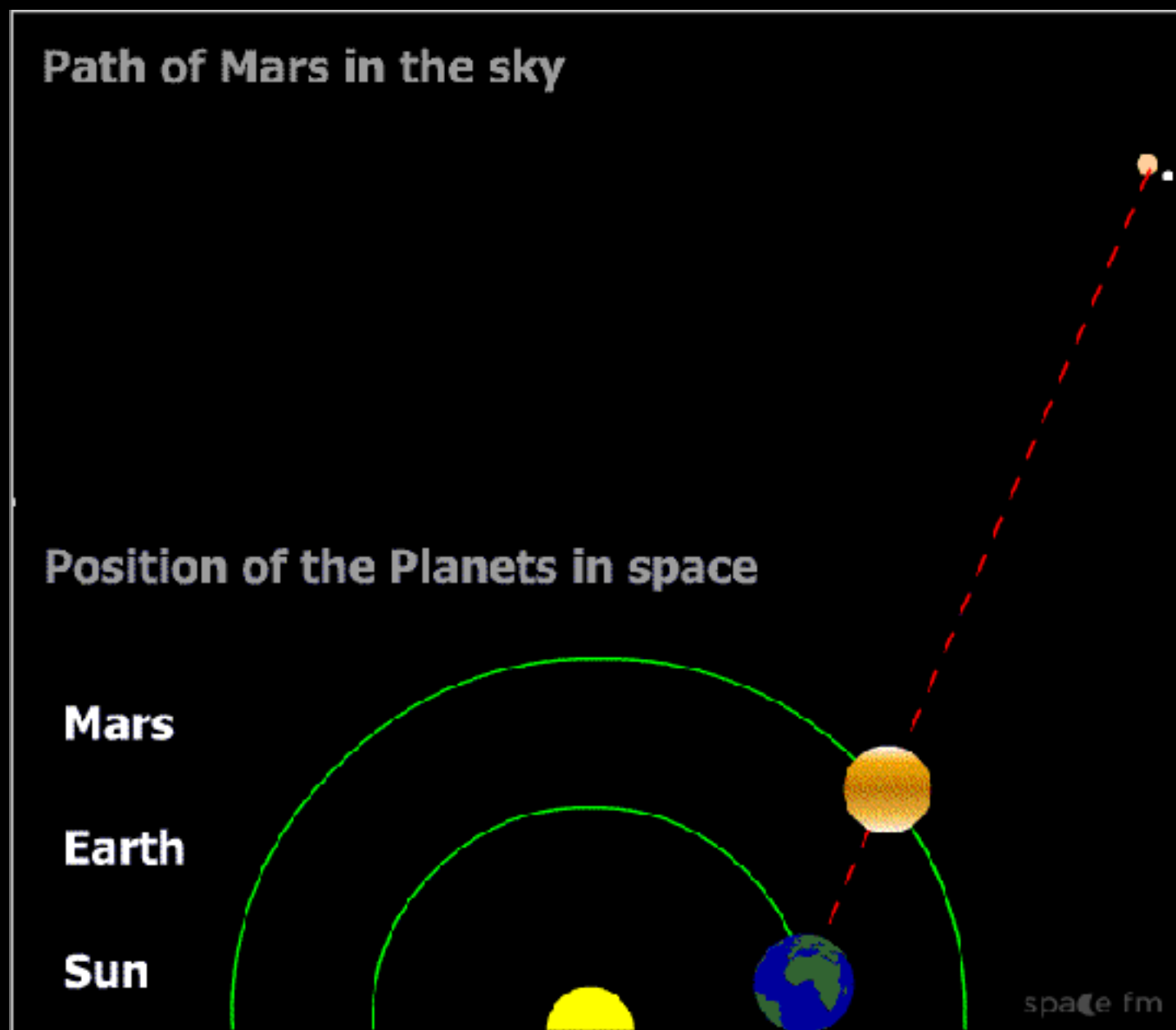
# Analytical models → Numerical models

*Computers are very good at this!*

**But sometimes we don't know how to model...**

*The world is complicated*





$$\frac{a^3}{T^2} = \text{const} = \frac{G(M + m)}{4\pi^2}$$

the rise of empirical models



Tycho Brahe

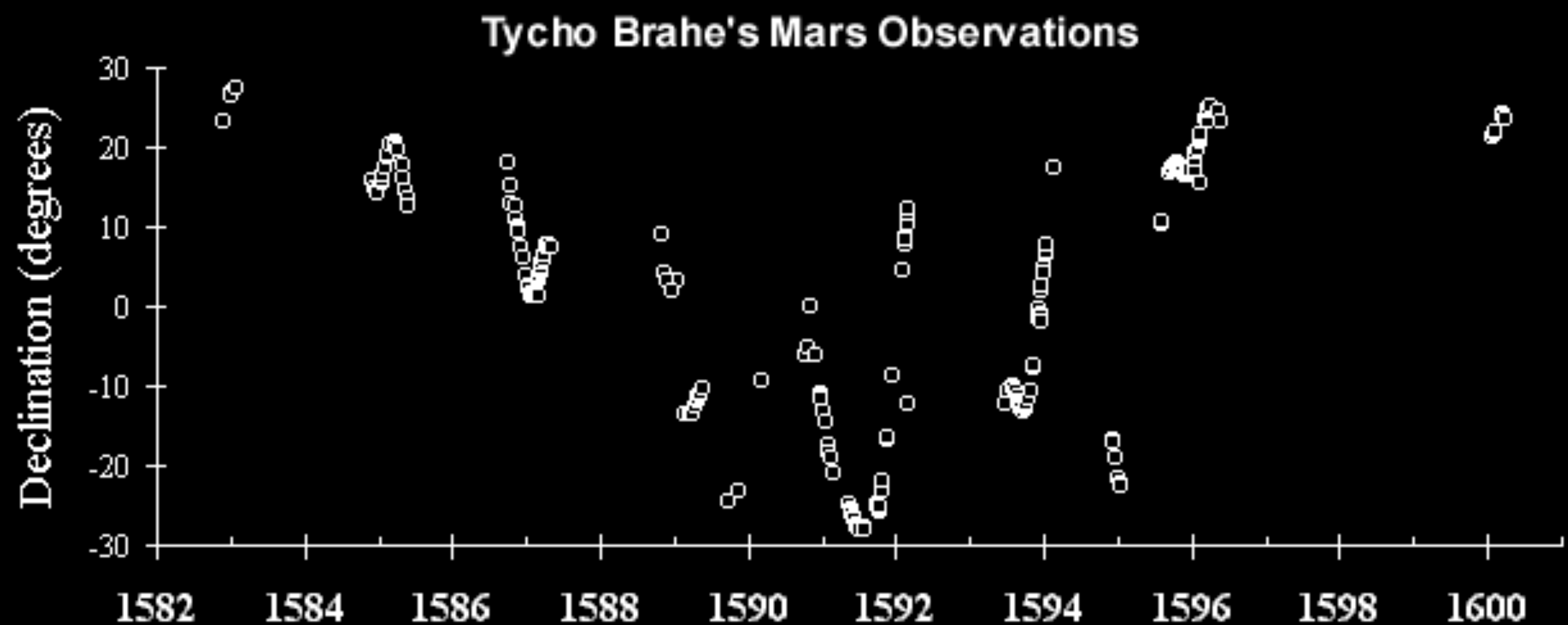
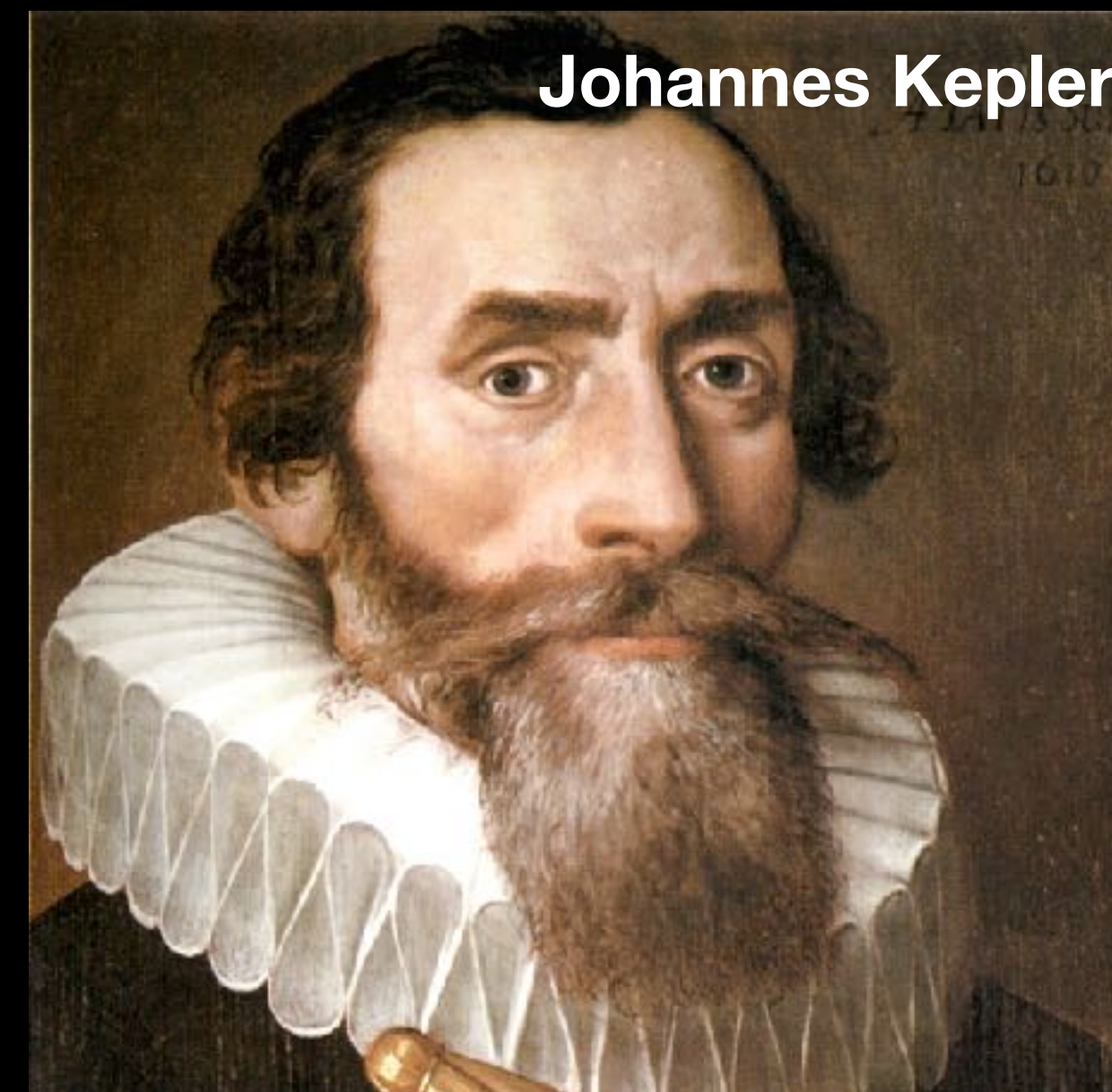


Image Copyright 2000, Wayne Pafko



Johannes Kepler



# The rise of empirical models

## Analytical/numerical models

Driven by **direct knowledge**

**Formulation** required (expensive human efforts)

Based on certain assumptions

Subject to the quality of the **assumptions**

Seek to **emulate** the real-world system

Making predictions can be **expensive**

## Empirical models

Driven by **data**

**Observation** required (can be automated)

No/little assumptions

Subject to the quality of the **data**

Seek to **find patterns** in the data

Making predictions is **straightforward**

Nature, I want to **understand**  
how you think and do, fully...

**Analytical models**

You wish!

**Nature**

Nature, I want to **see** how you  
think and do, more and more...

**Empirical models**

You may.

**Nature**

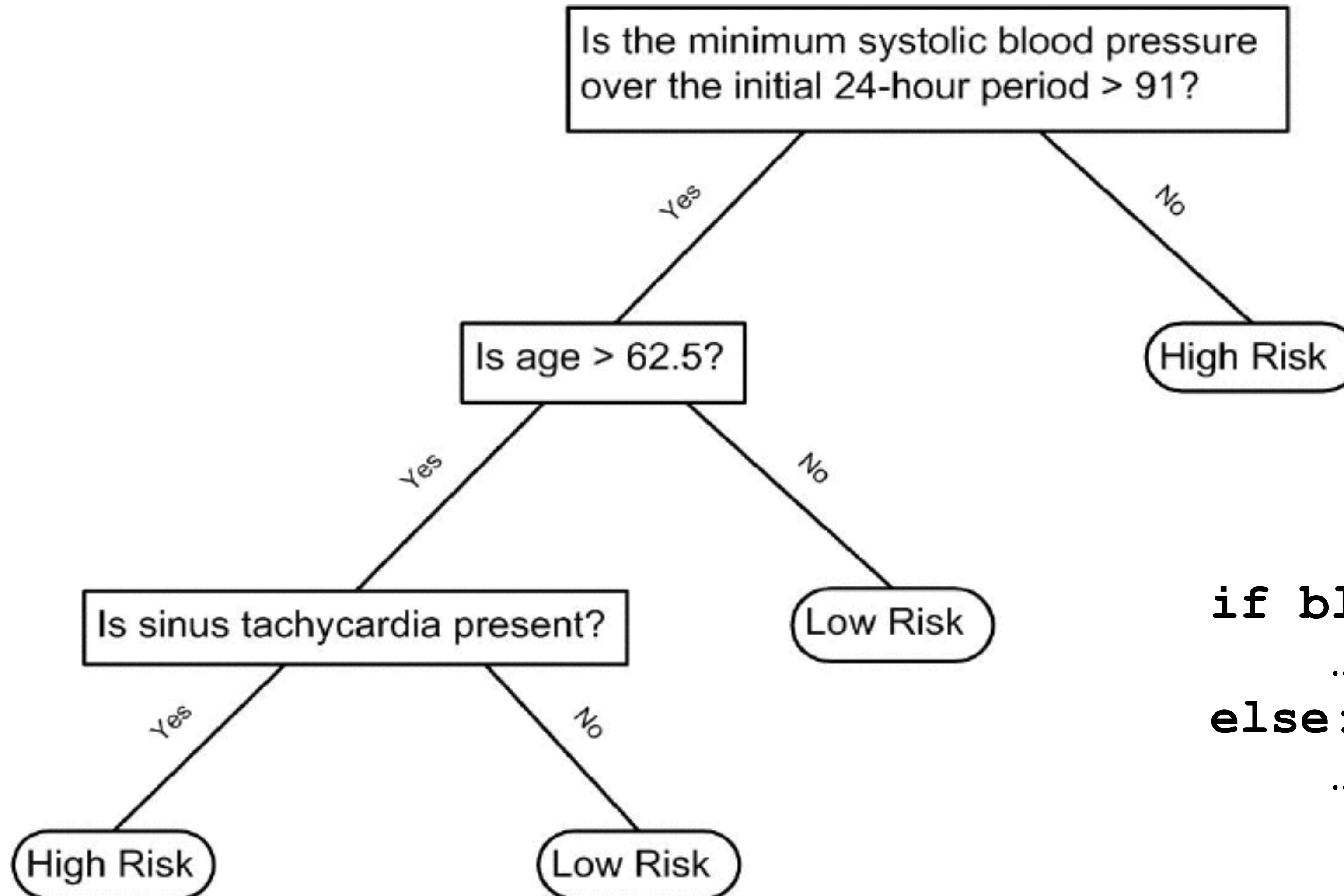


**Complex problem, (huge amount of) complex data**

*New programming paradigm is needed!*



# Decision Tree(s)



```
if blood_pressure > 91:
```

```
...
```

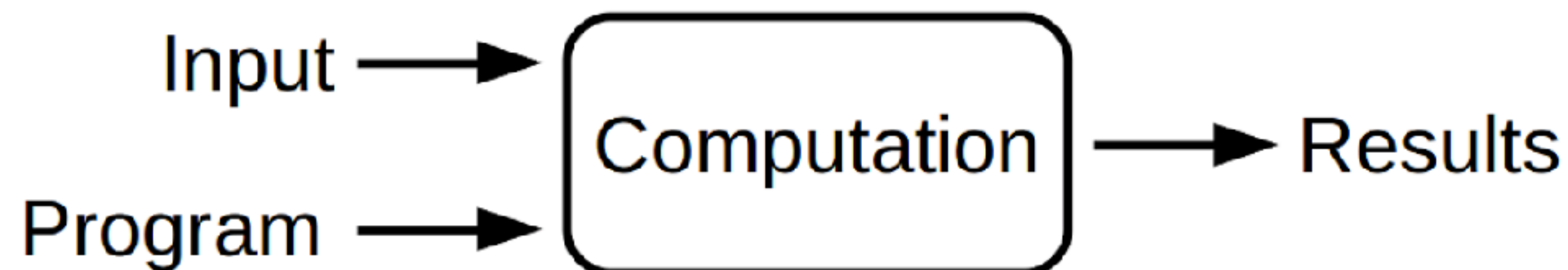
```
else:
```

```
...
```



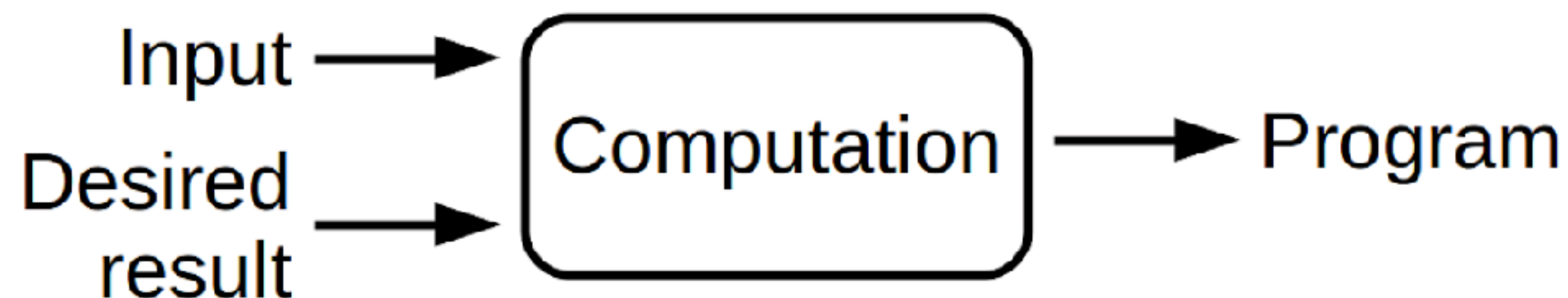
# Machine learning: a new programming paradigm

## Traditional programming



Knowledge base & rules  
Expert systems  
**Human input**

## Machine learning

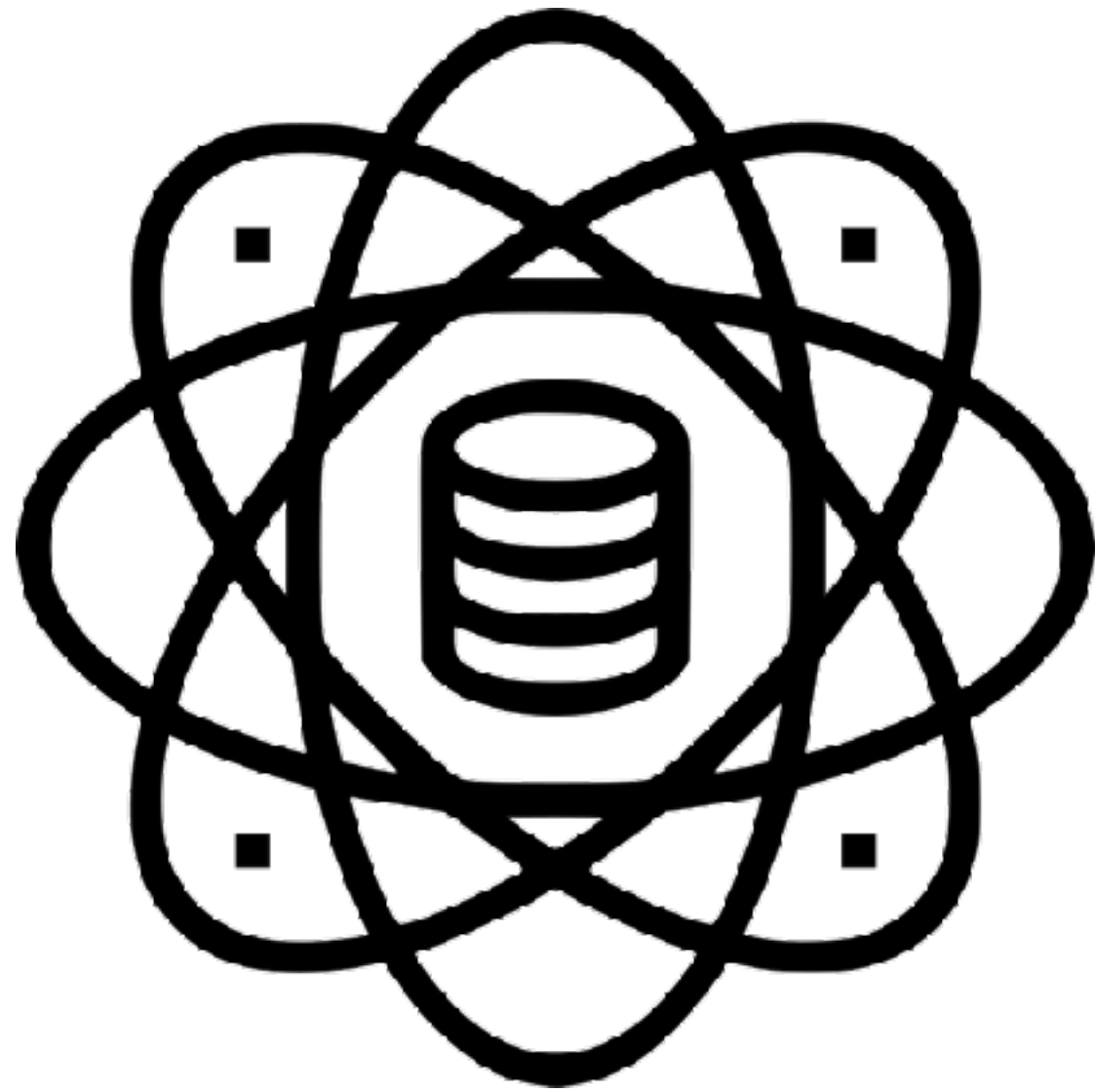


Features and result  
'Decision' system  
**Limited human input**

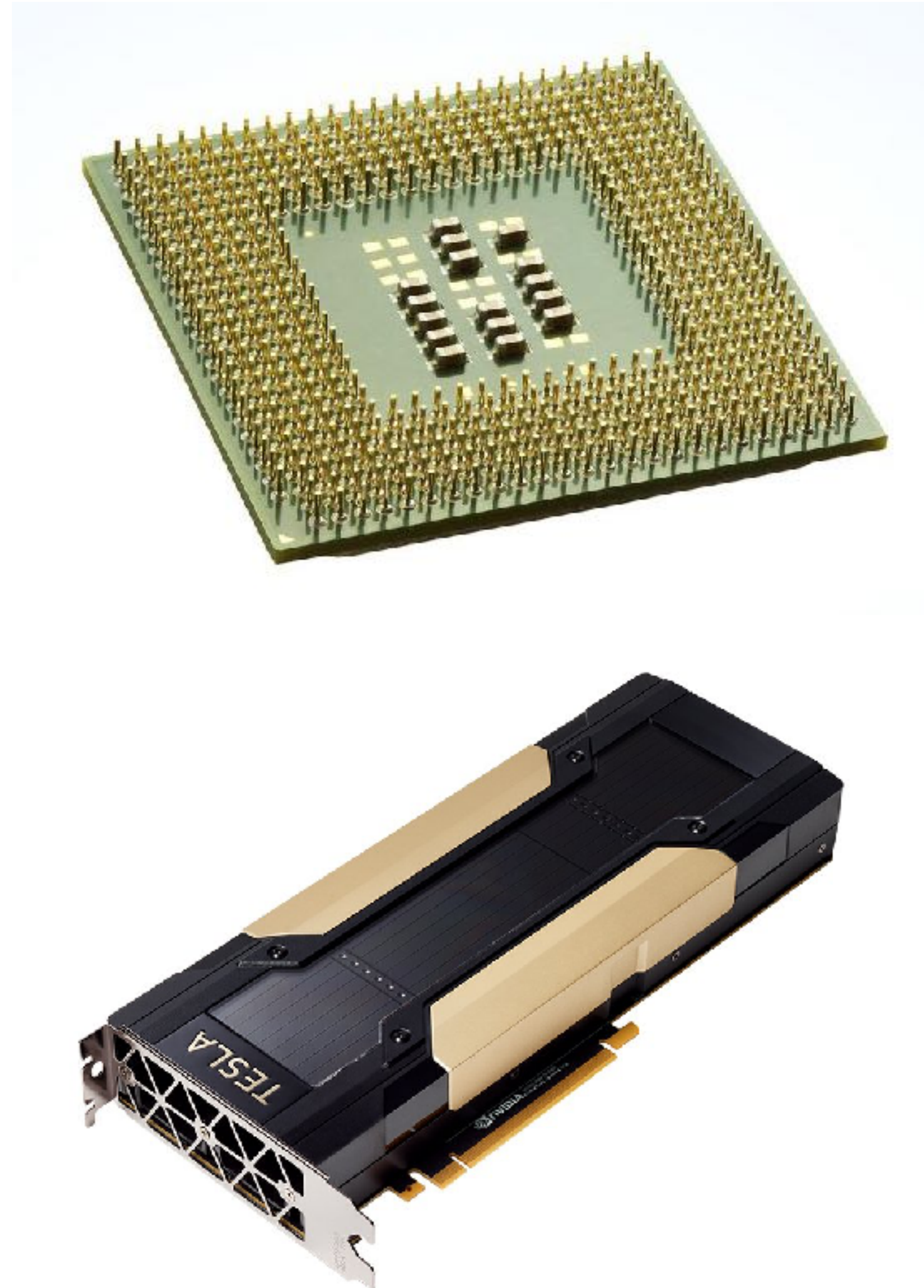


# Three ingredients of AI

(Scientific) data



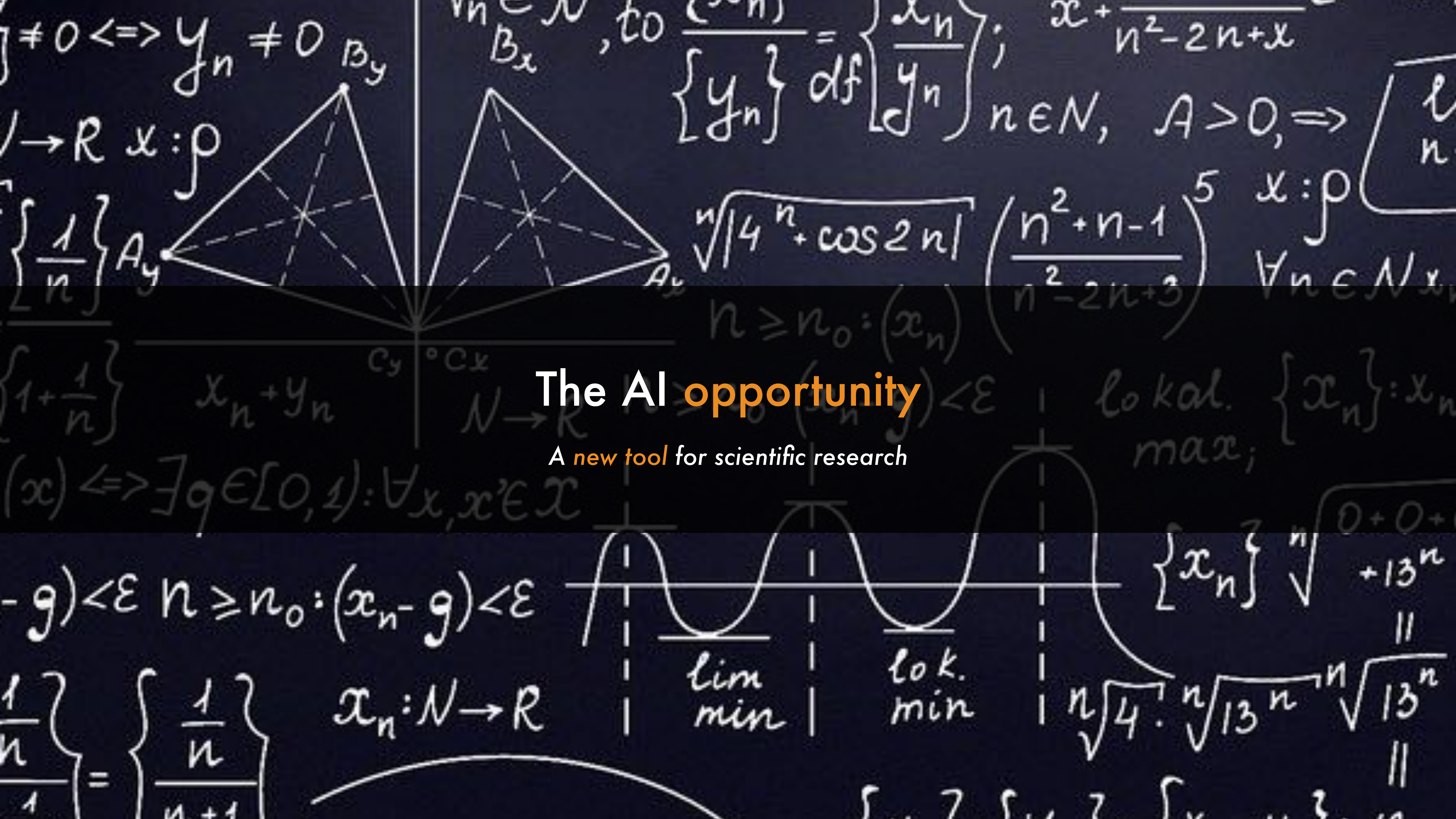
Computational Facilities



Algorithms







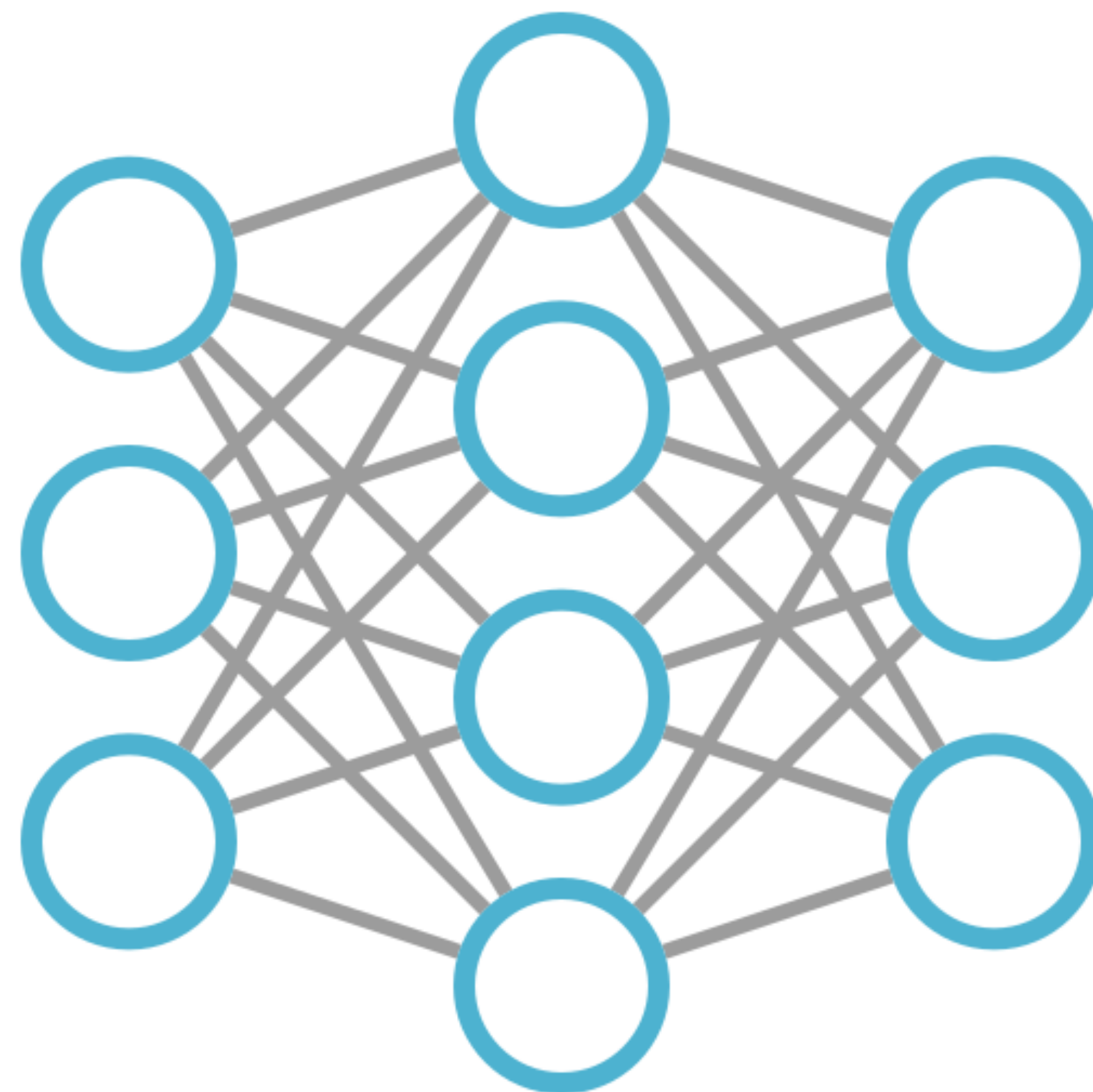
# The AI opportunity

A new tool for scientific research





**Scientific Problem**

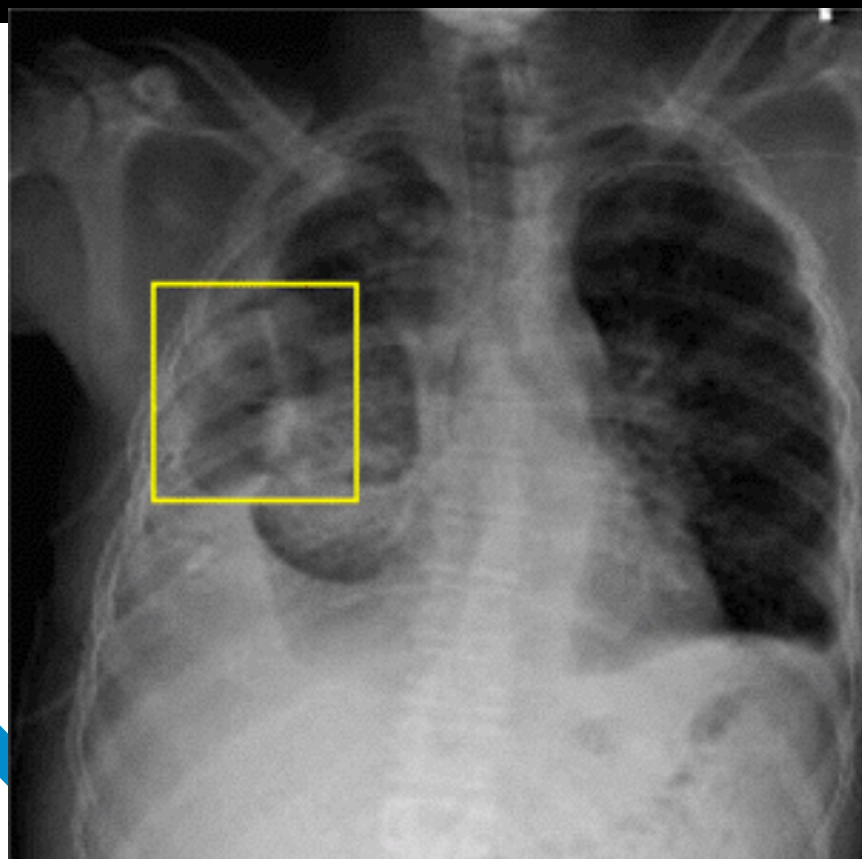


**AI Problem**

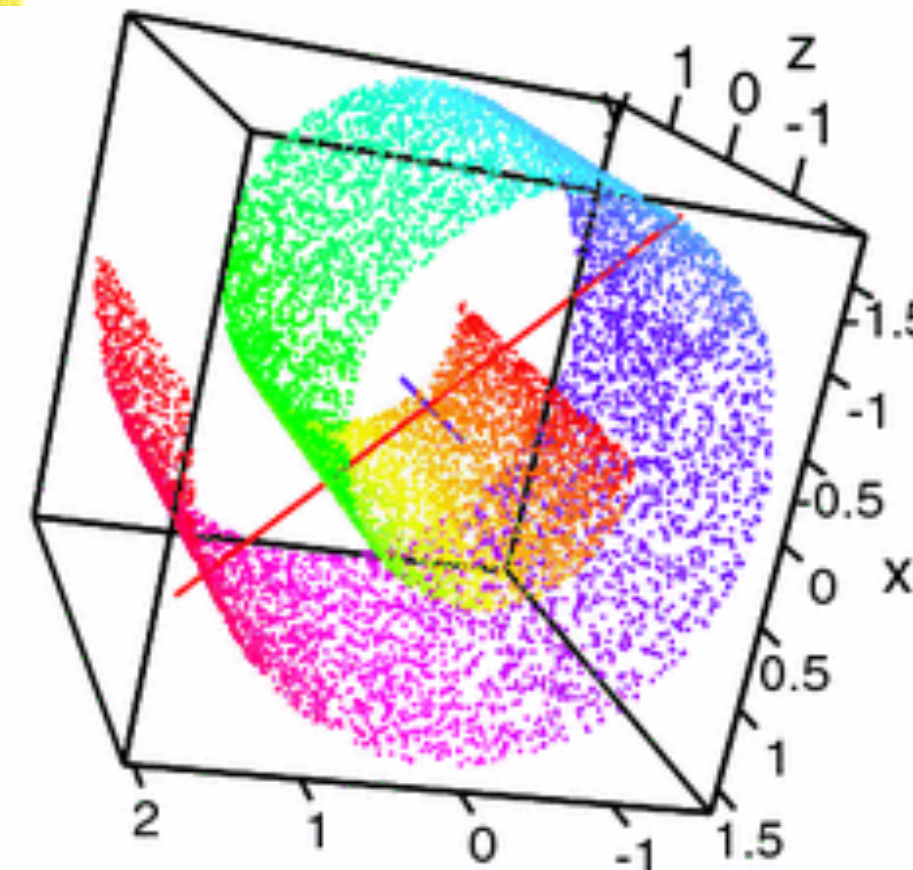


# Scientific Problem → AI Problem

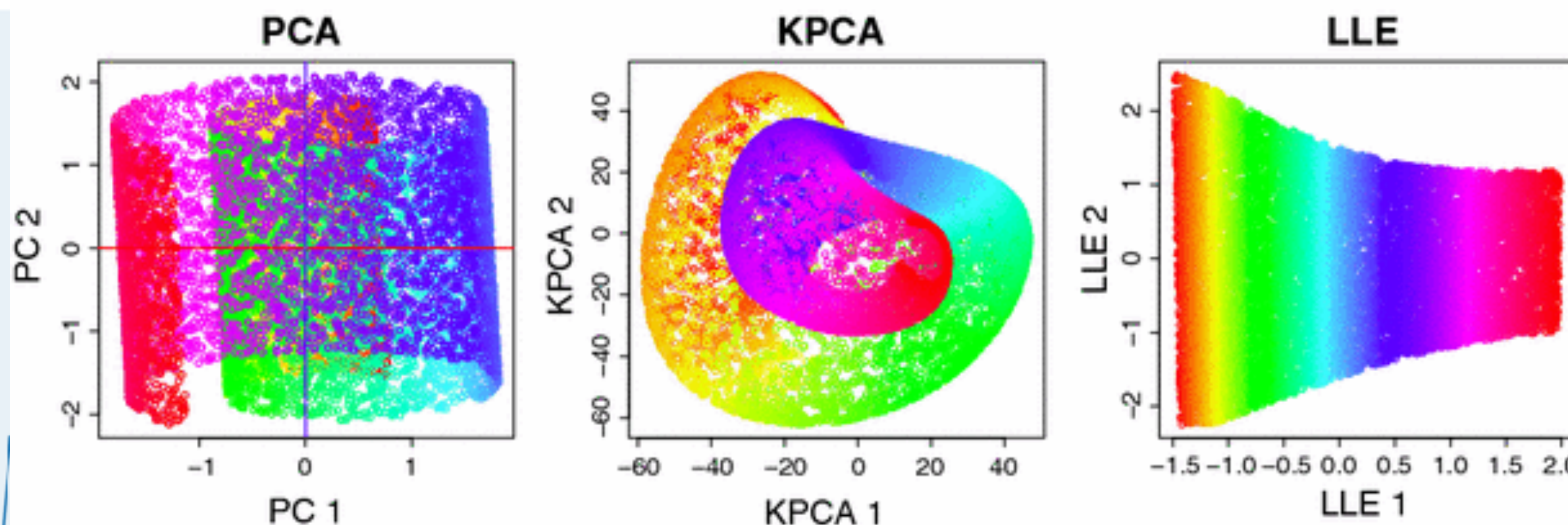
Classification



Dimension reduction

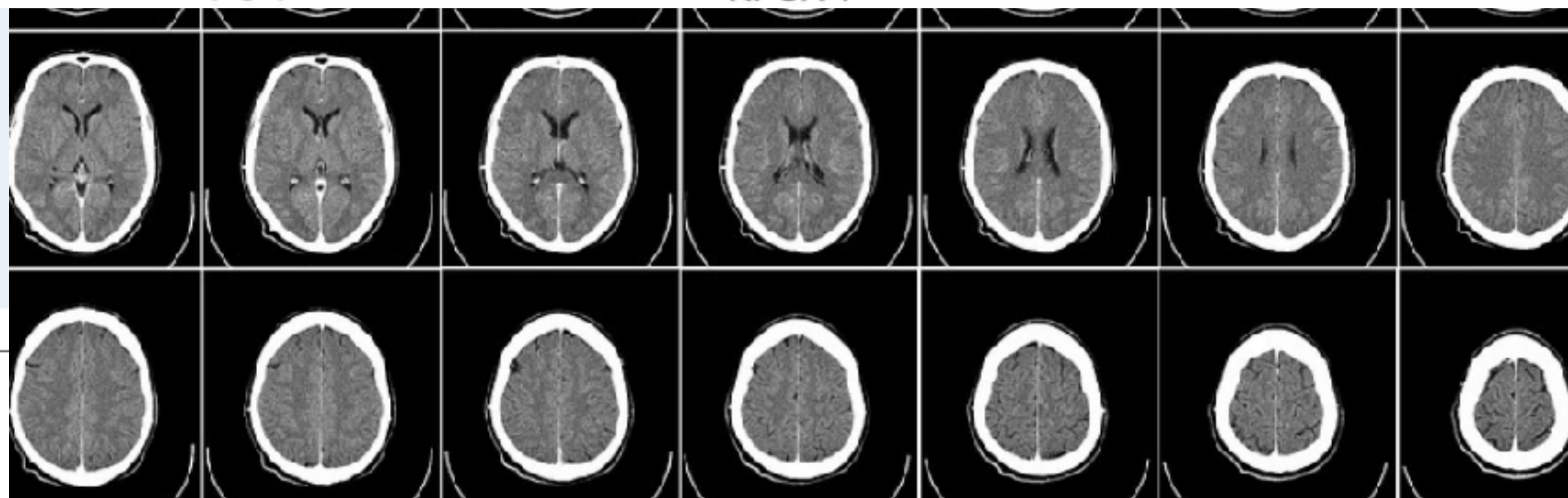
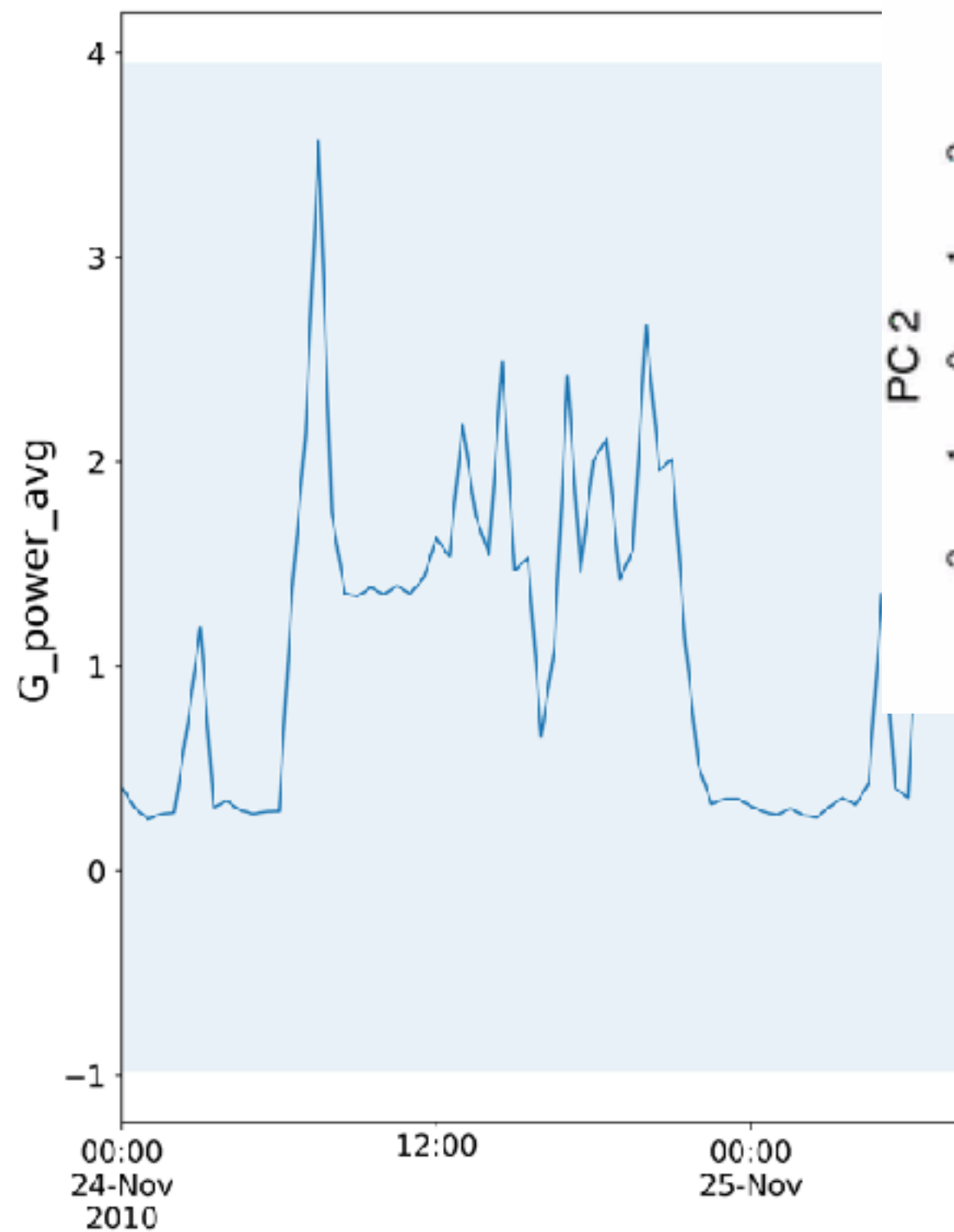


Clustering



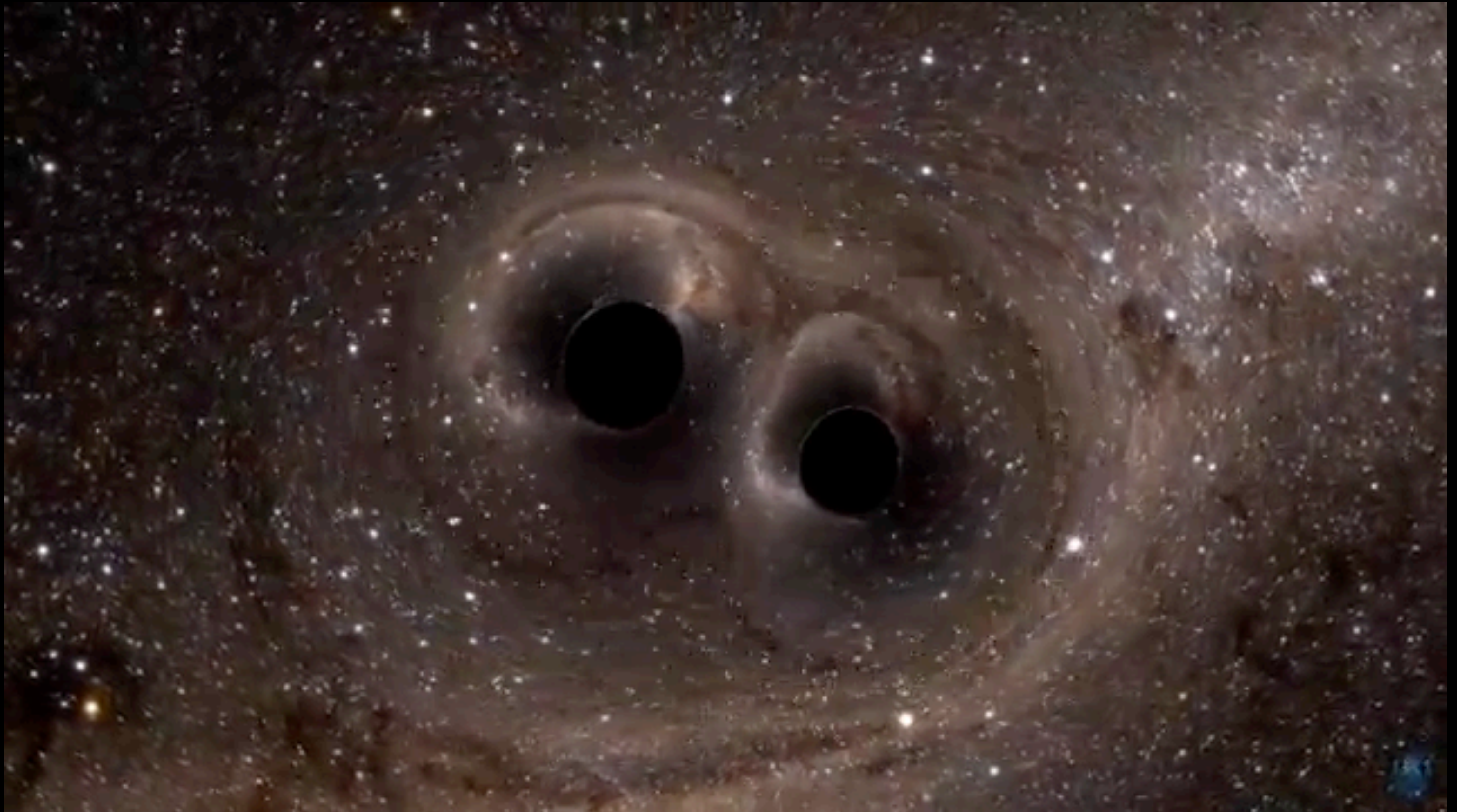
Generation

Regression





# Black hole Mergers → Gravitational Waves

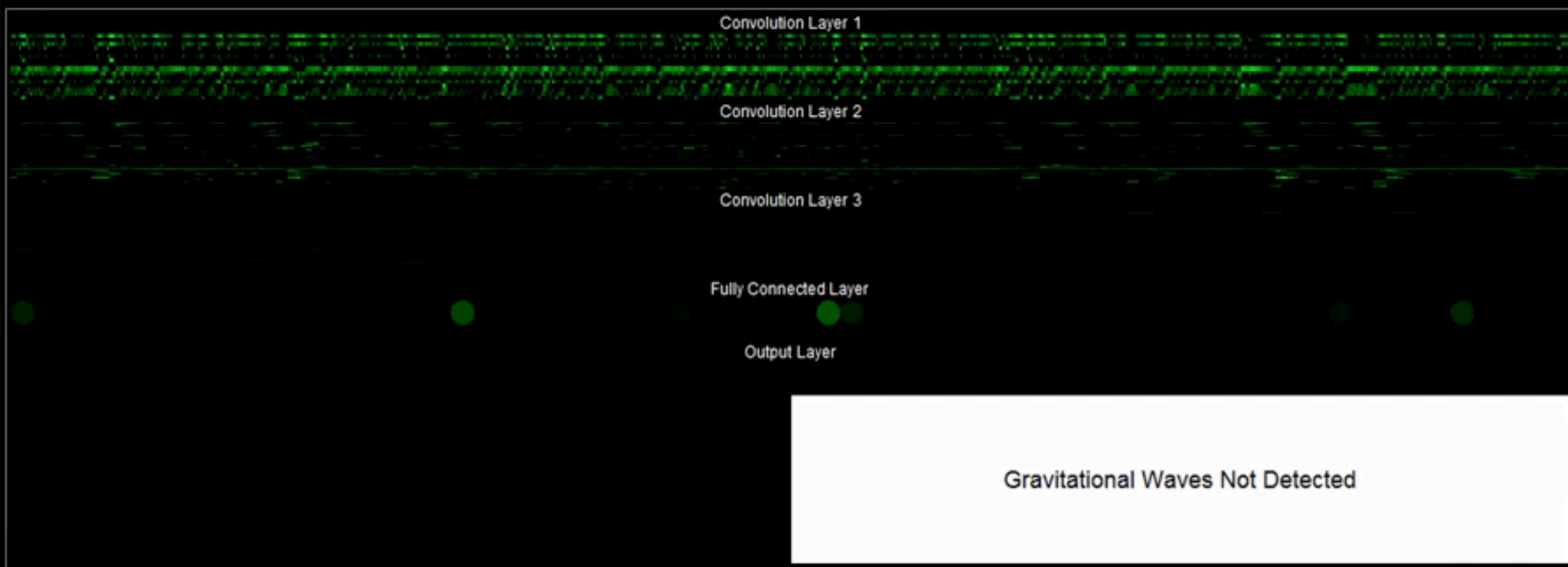
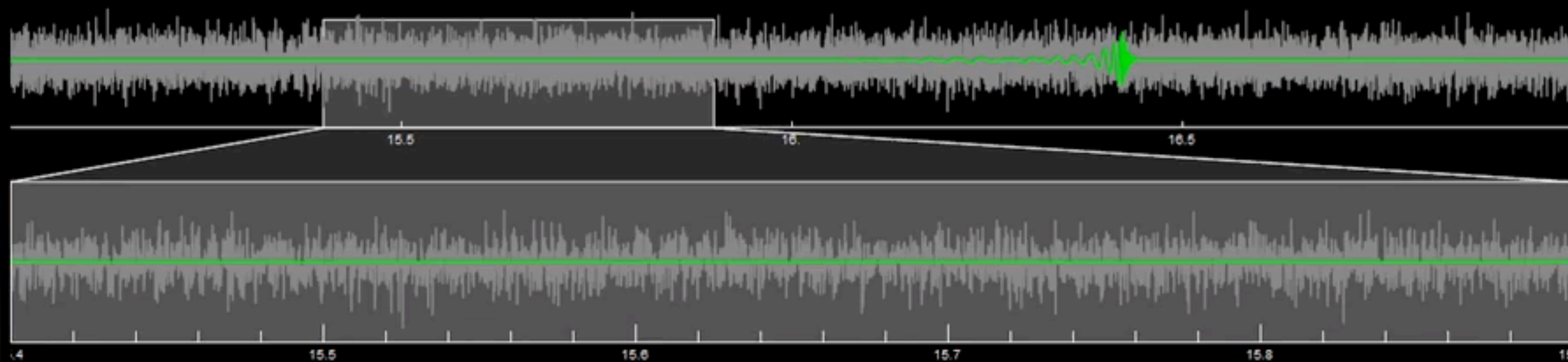
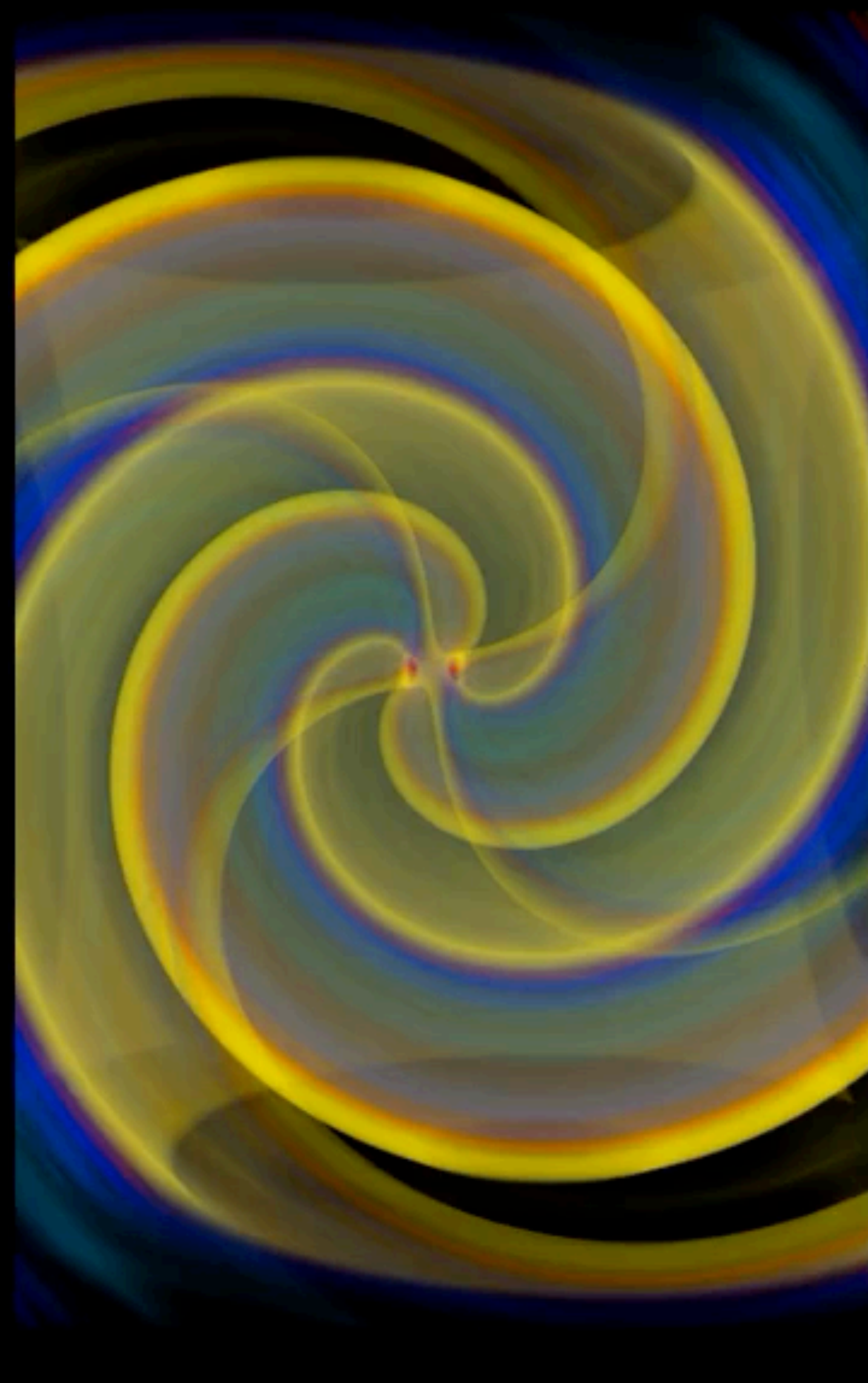


*Credit: European Southern Observatory*



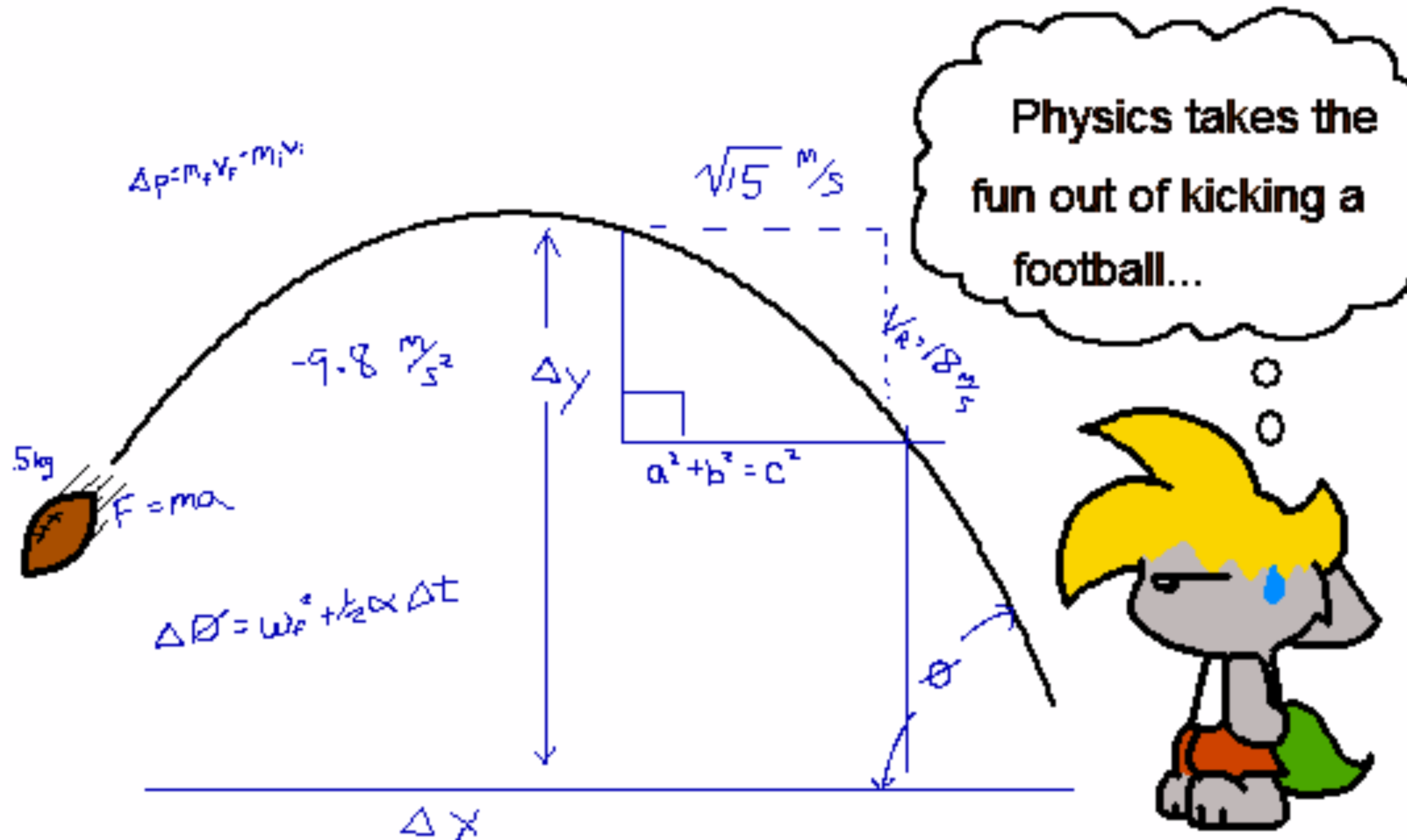
# Detecting Gravitational Waves in Real-Time with Deep Learning

Data from a LIGO Interferometer around the first event (GW150914)

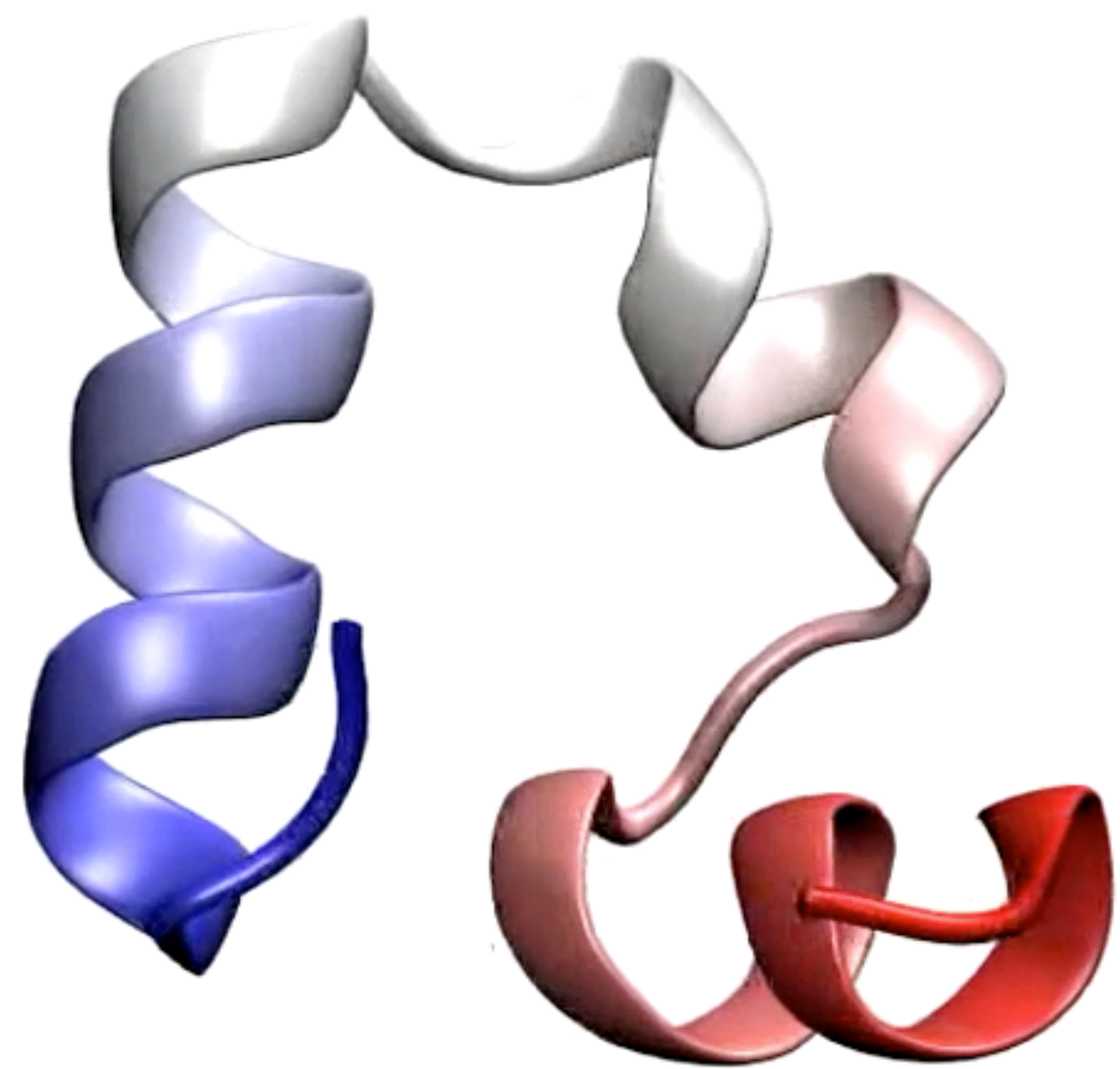




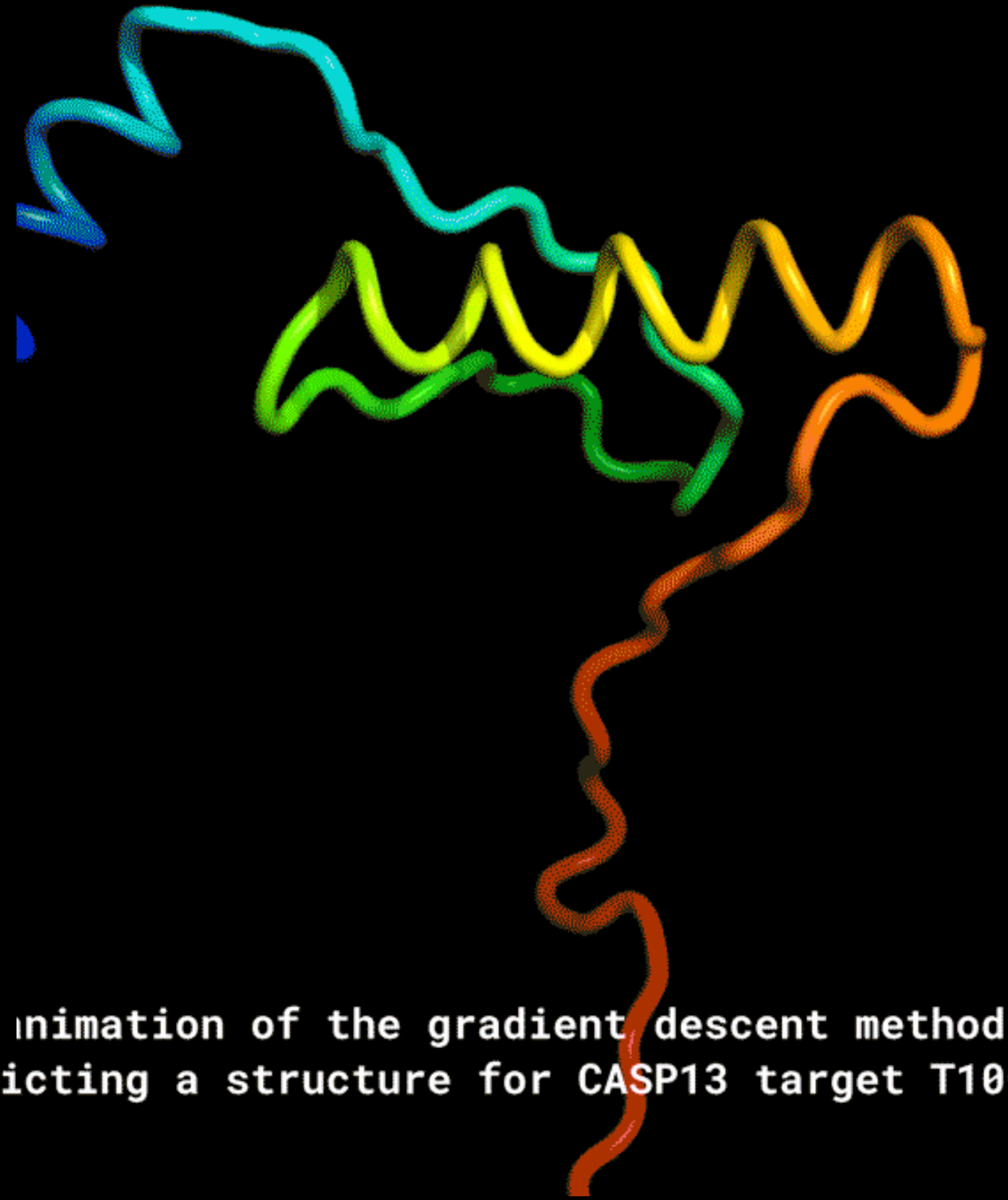
# Let AI to learn the intuition







Animation of the gradient descent method  
predicting a structure for CASP13 target T1008





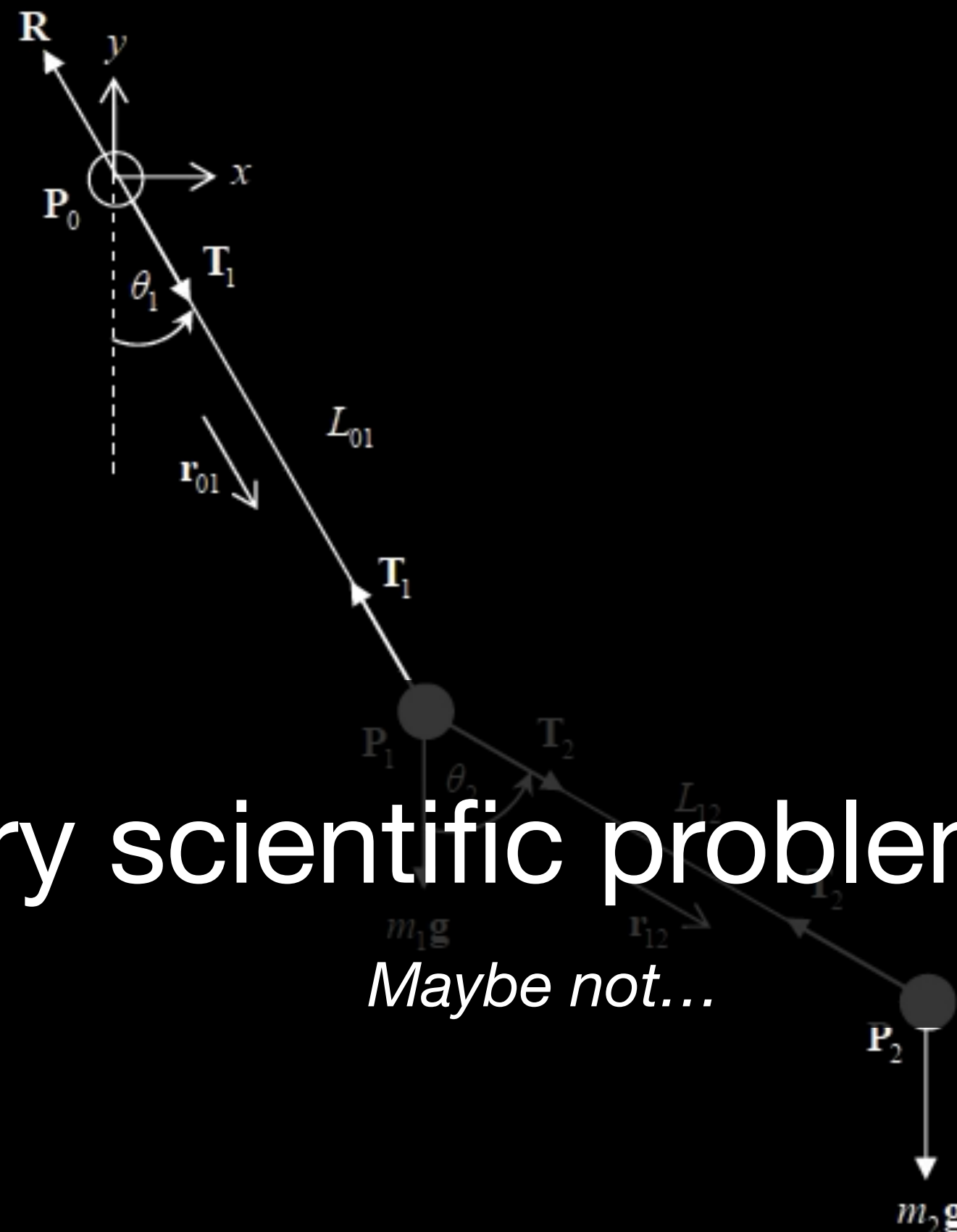
# The AI Challenges

Scientific research is non-trivial; using AI in scientific research is also **non-trivial**





Can we map every scientific problem to an AI problem?



Maybe not...

$$\dot{\theta}_1 = \omega_1$$

$$\dot{\theta}_2 = \omega_2$$

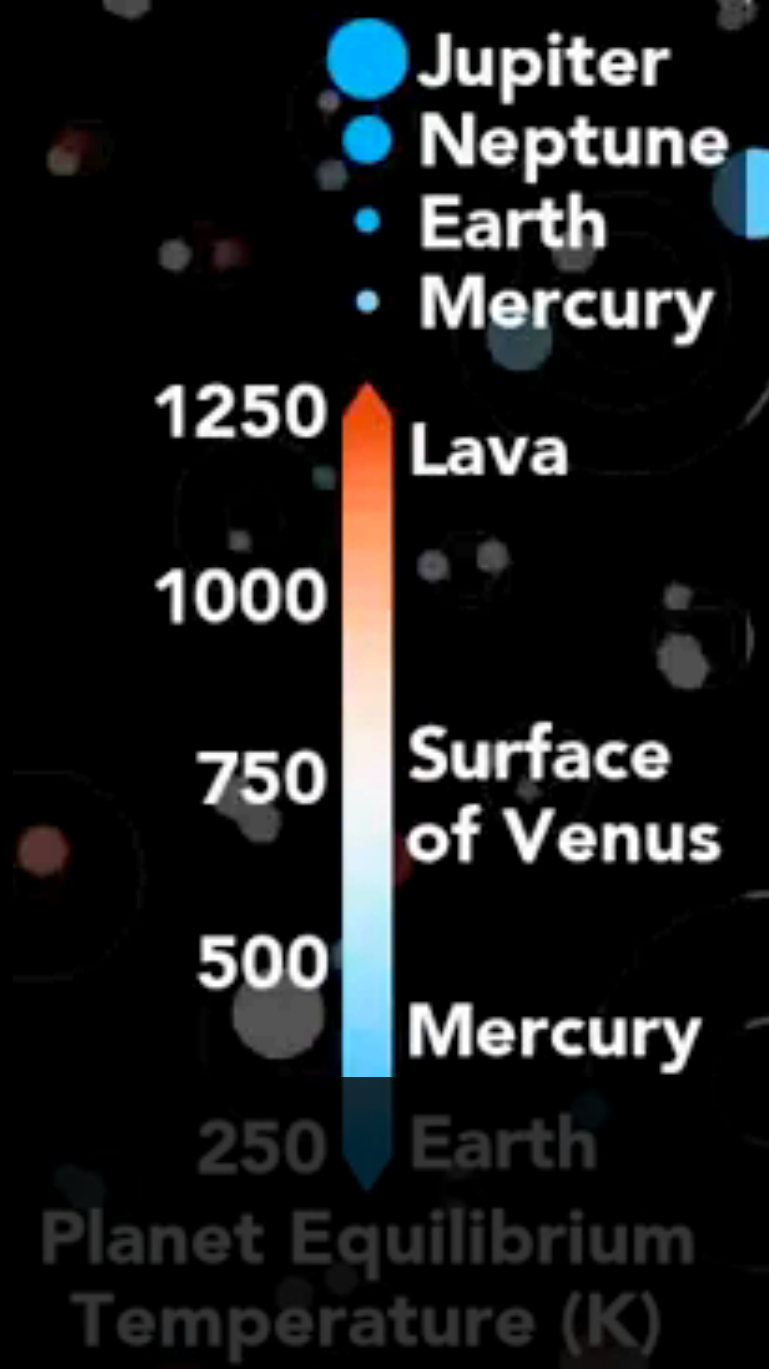
$$\dot{\omega}_1 = \frac{\frac{g}{L_{01}} \left( \sin \theta_2 \cos(\theta_1 - \theta_2) - \left( 1 + \frac{m_1}{m_2} \right) \sin \theta_1 \right) - \omega_2^2 \frac{L_{12}}{L_{01}} \sin(\theta_1 - \theta_2) - \omega_1^2 \sin(\theta_1 - \theta_2) \cos(\theta_1 - \theta_2)}{1 + \frac{m_1}{m_2} - \cos^2(\theta_1 - \theta_2)}$$

$$\dot{\omega}_2 = -\dot{\omega}_1 \frac{L_{01}}{L_{12}} \cos(\theta_1 - \theta_2) + \omega_1^2 \frac{L_{01}}{L_{12}} \sin(\theta_1 - \theta_2) - \frac{g}{L_{12}} \sin \theta_2$$

Mapping



Kepler Orrery IV  
23 Nov 2010  
By Ethan Kruse  
@ethan\_kruse



Sometimes, **outliers** are more interesting!

*But outliers are rare...*

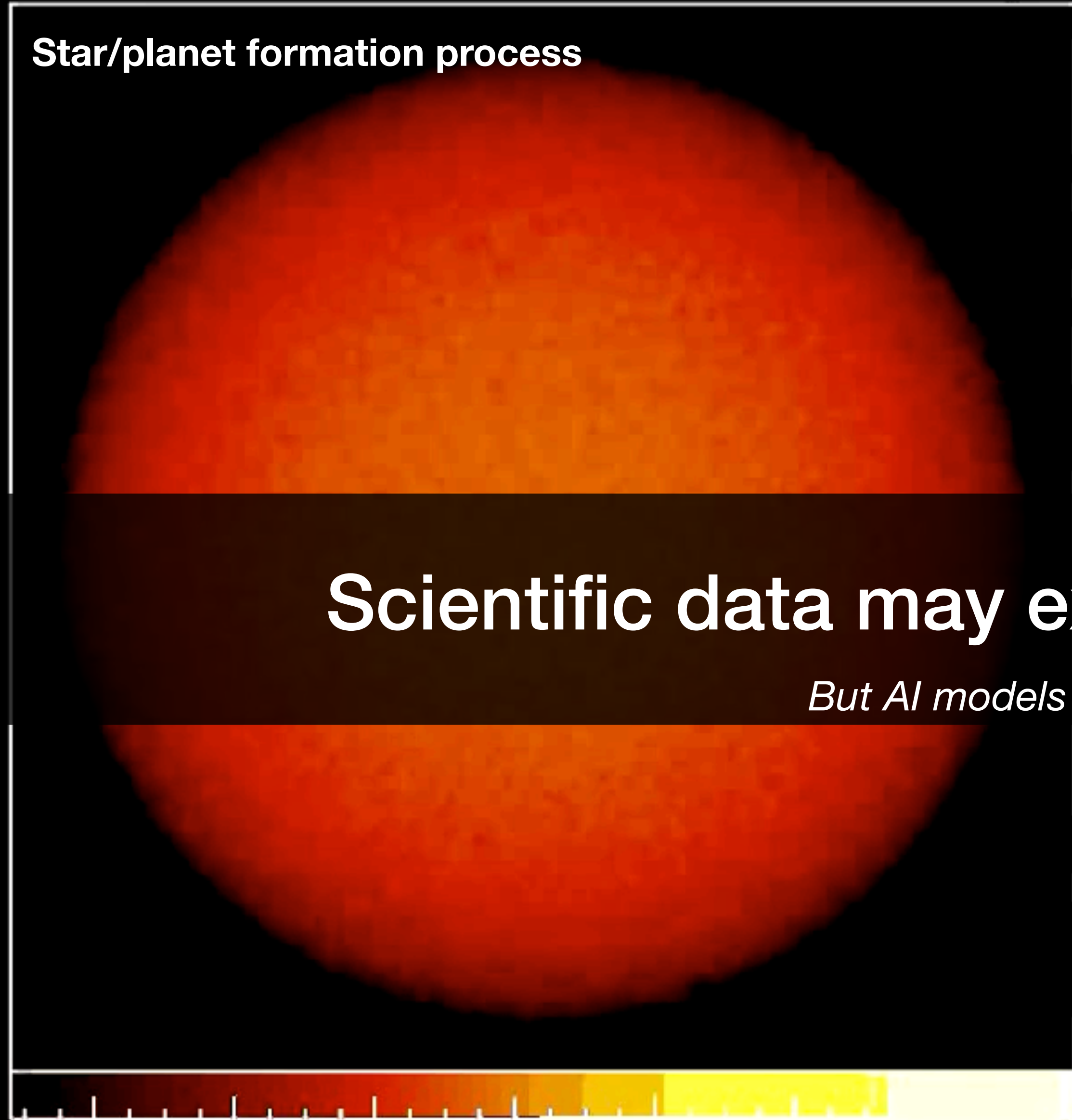
Data complexity



Dimensions: 82500. AU

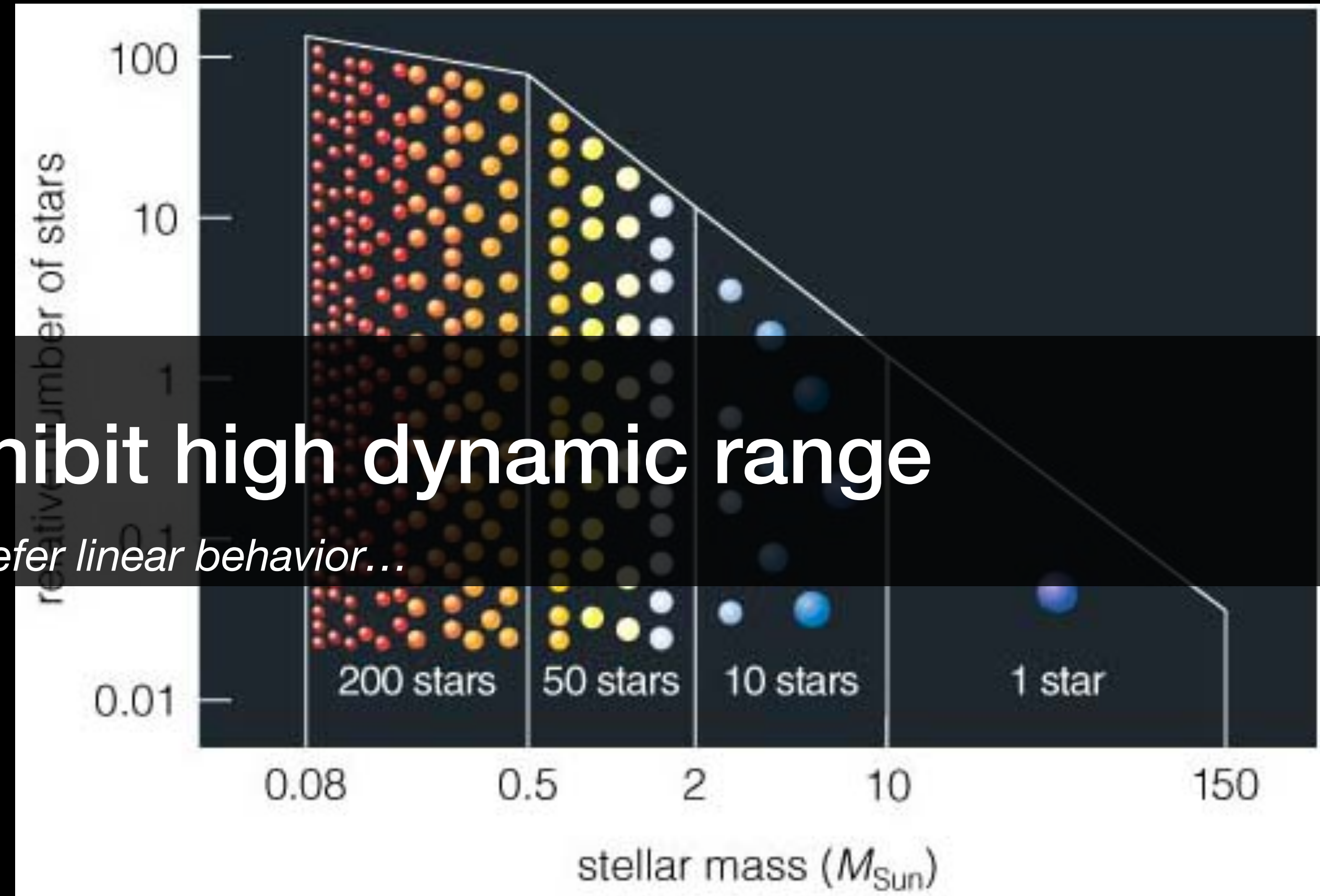
Time: 0. yr

Star/planet formation process



Scientific data may exhibit high dynamic range

*But AI models prefer linear behavior...*



Initial stellar mass

Data complexity

Log Column Density [g/cm²]

Matthew Bate



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

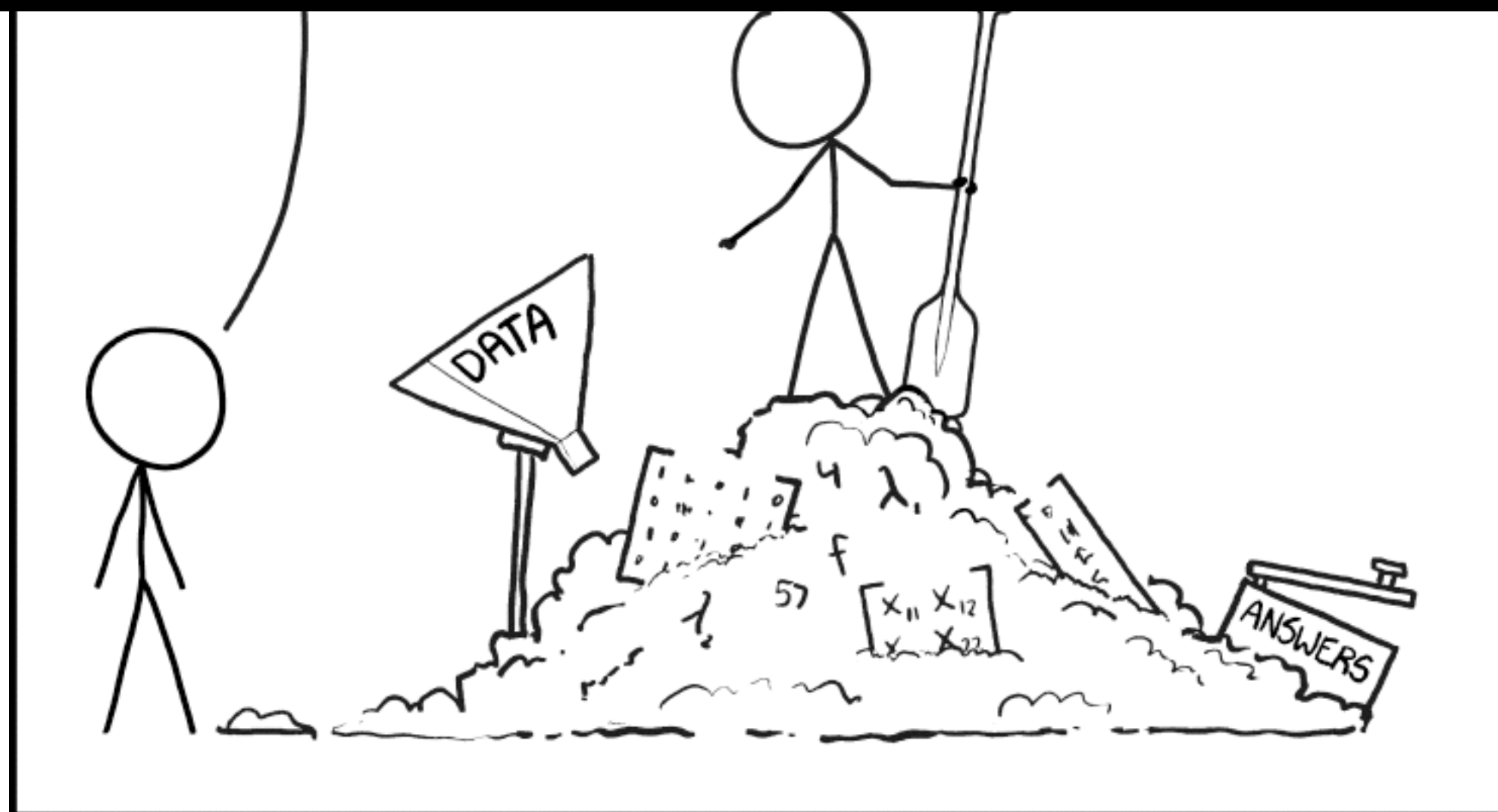
WHAT IF THE ANSWERS ARE WRONG?

**“Please, explain the model!”**, says the scientist.

40000

## Can we simply trust the AI prediction?

Can be dangerous...



Source: xkcd

Number of paper  
interpretable mac

20000

10000

0

98-00 01-03 04-06 07-09 10-12 12-15 15-16 16-17

Interpretability

Image source: Been Kim (Google AI)



# Opportunity

1 Data

2 Compute

2 Algorithms

1 Learning materials

# Challenges

1 Mapping 2

1 Data complexity 2

1 Interpretability 2

1 Community acceptance



1 Domain scientists



2 AI experts



# Conclusions

- ➔ AI is a new programming paradigm
- ➔ AI is a modern version of empirical models
- ➔ AI is not a buzzword; it is a new tool available to the scientific community (after some adaption)
- ➔ AI is not foolproof
- ➔ Leveraging the potentials of AI in scientific research requires joint efforts of domain scientists and AI experts

