

What is the goal of life?

Readings for today

- Vlastelica M. (2019). Learning Theory: Empirical Risk Minimization. Towards Data Science. (Blog post)
- PAC learning. Wikipedia.
- Valiant, L. G. (2009). Evolvability. Journal of the ACM (JACM), 56(1), 1-21.

Topics

- Entropic view of life
- Structure of a learnable problem

Entropic view of life

Let's talk about life

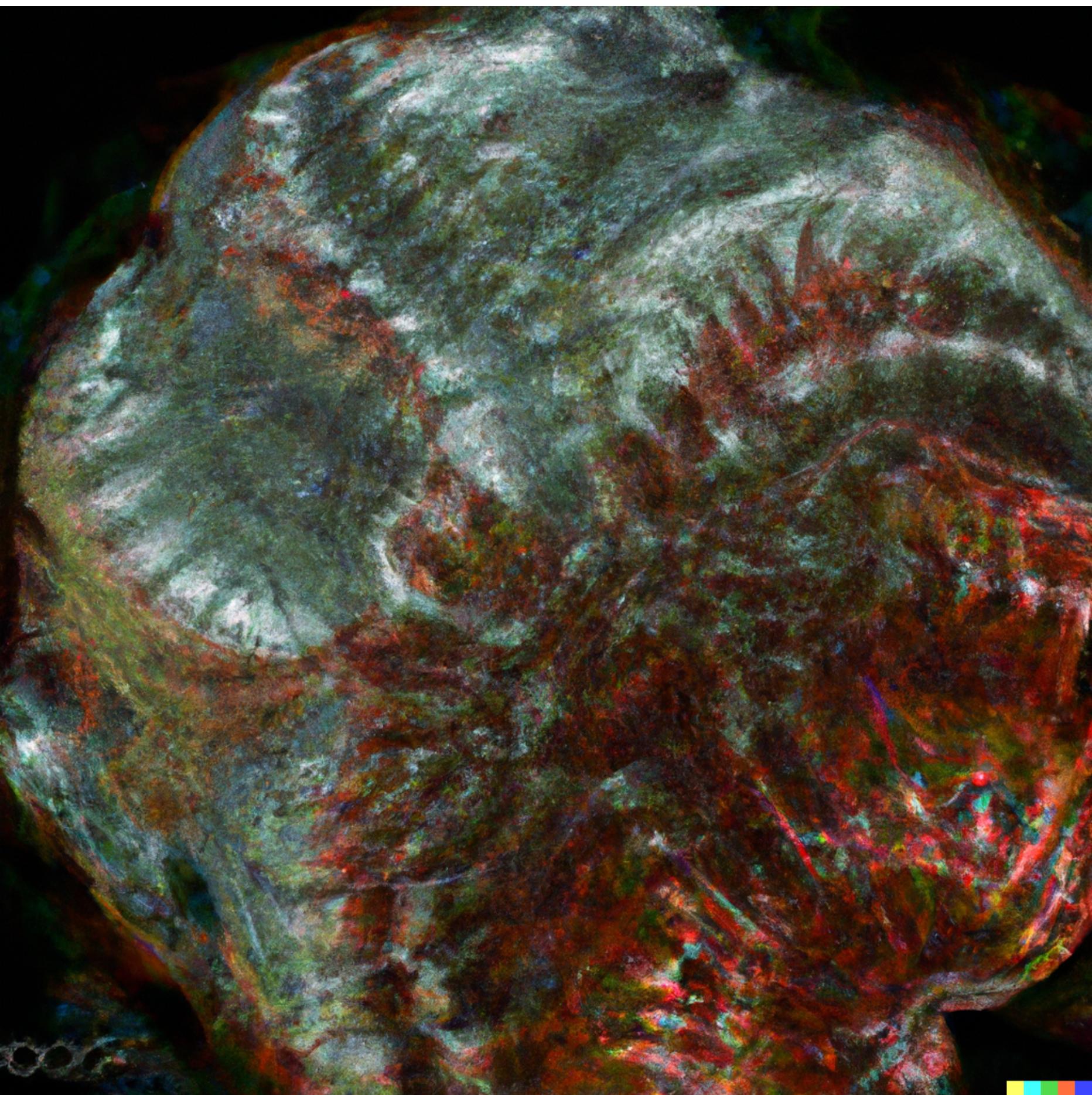
Question: What is life?

Answer: It's complicated

- Cellular life
- Autocatalytic molecules
- Generation of proteins
- Reproduction
- Adaptation

Two basic properties:

- **Metabolism:** maintain their own structure by continually resynthesizing their own elements.
- **Replication:** Produce copies of themselves in the immediate vicinity.



A better question

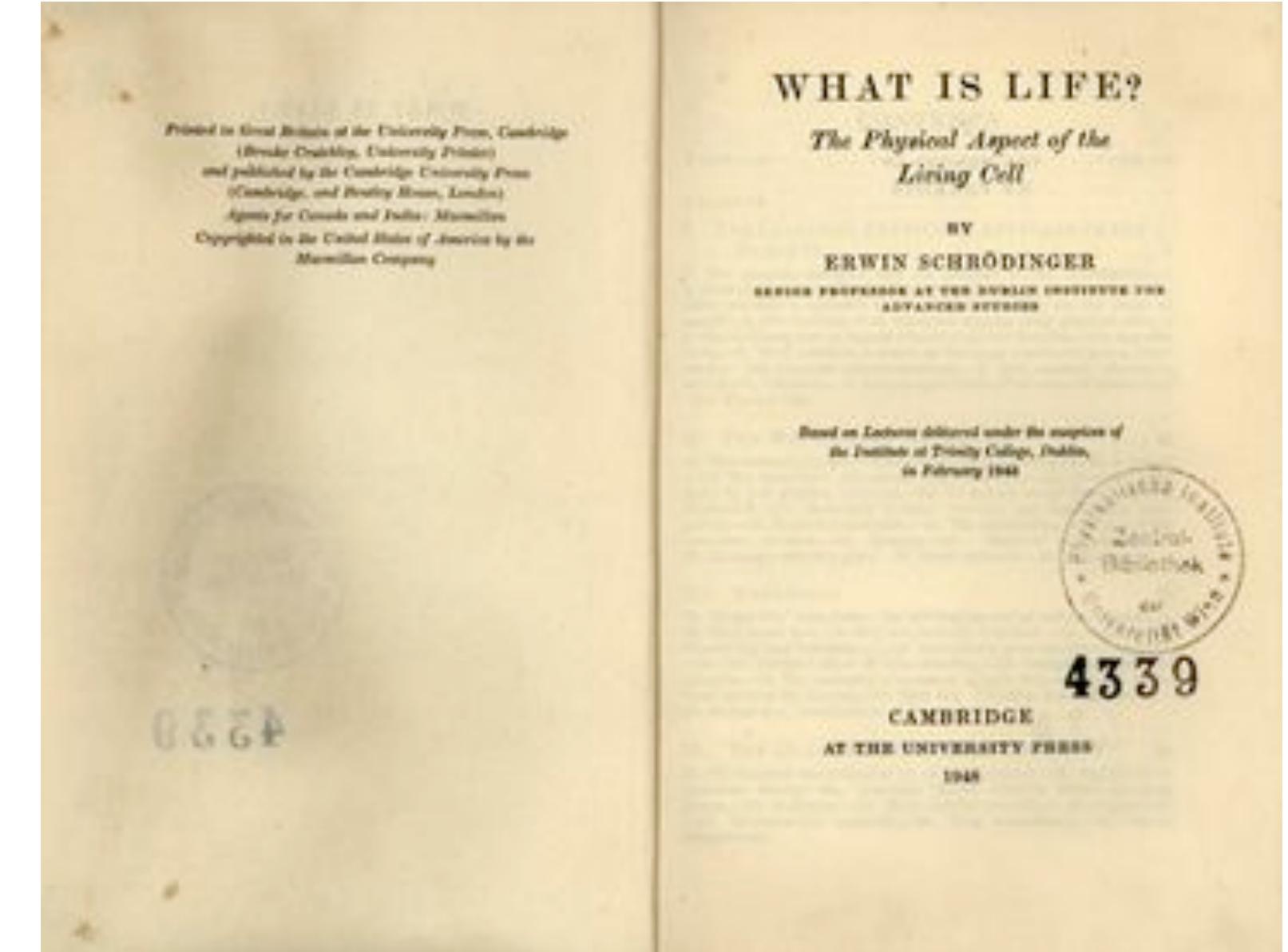
Question: What does life do?

Answer: “Order from order” - E. Schrödinger

- Life controls/reduces entropy of its environment.
 - Manipulates local entropy so as to maintain molecular stability.

Entropy: Number of possible states that matter can take on.

$$S = -k_B \sum_i p_i \ln(p_i)$$

Stochastic model of evolution

Goal: Maximize fitness of organism

- Random mutations (introduce change)
- Natural selection (select best fit)

Problem: If the process is a random search, then there are not enough seconds since the start of the universe to produce life as complex as what we have currently on earth.



NP-hard problem

P vs. NP complexity

P: Can be solved in polynomial time.

Example: *Long multiplication* algorithm.

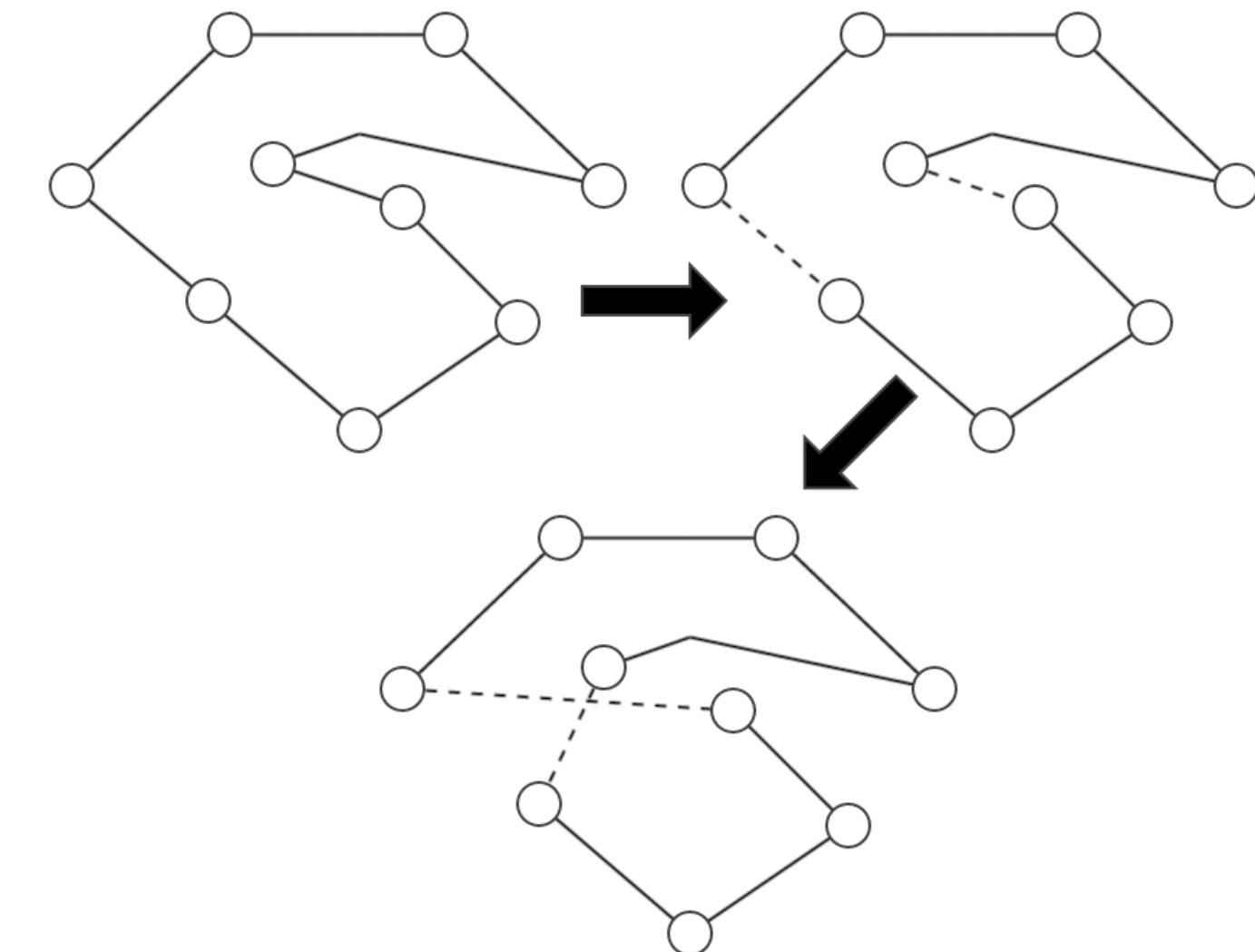
Multiplying 2 n -bit integers requires n^2 steps (i.e. $O(n^2)$)

$$\begin{array}{r} & 23958233 \\ \times & 5830 \\ \hline 00000000 & (= 23,958,233 \times 0) \\ 71874699 & (= 23,958,233 \times 30) \\ 191665864 & (= 23,958,233 \times 800) \\ + 119791165 & (= 23,958,233 \times 5,000) \\ \hline 139676498390 & (= 139,676,498,390) \end{array}$$

NP: Can be solved in nondeterministic polynomial time.

Example: *Traveling salesman problem*

Euclidean solution: $O(n(\log n)^{(O(c\sqrt{d}))^{d-1}})$



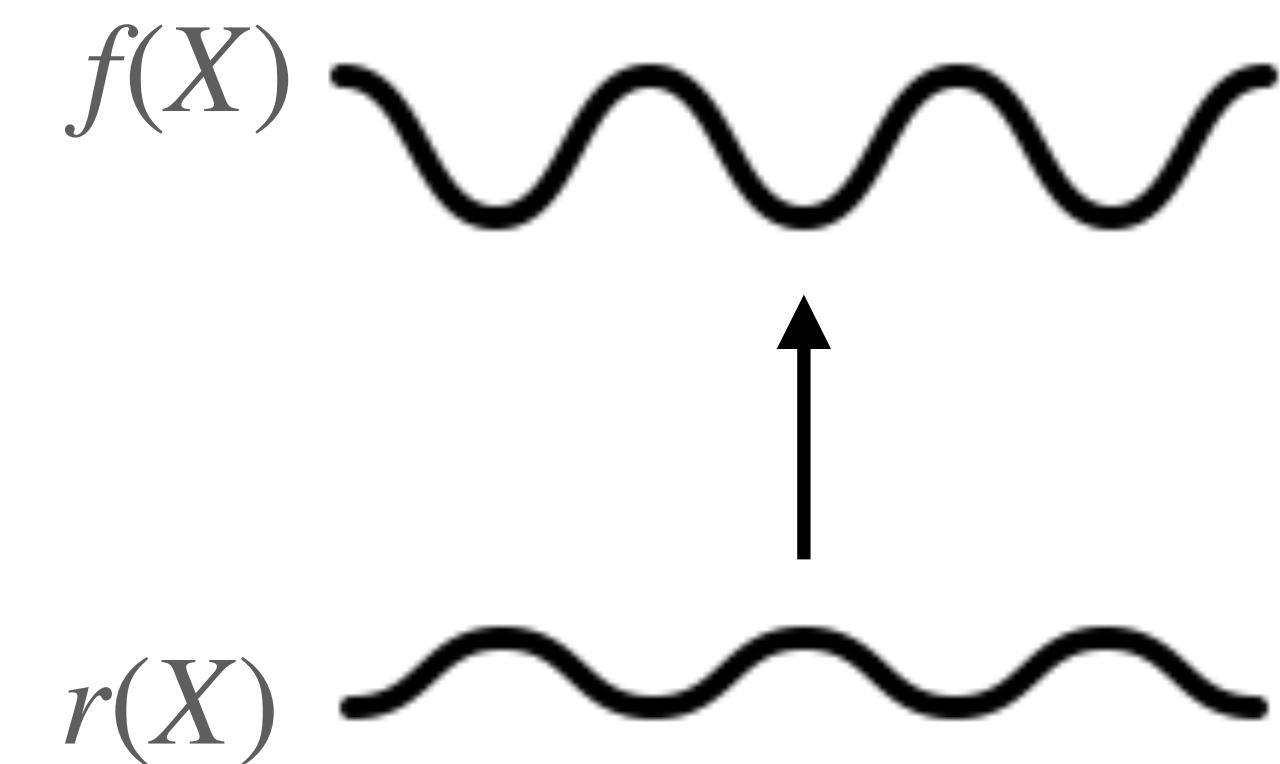
Evolvability

Question: What if life is a structured learning problem?

Process: 1. The goal of cellular life is to produce proteins

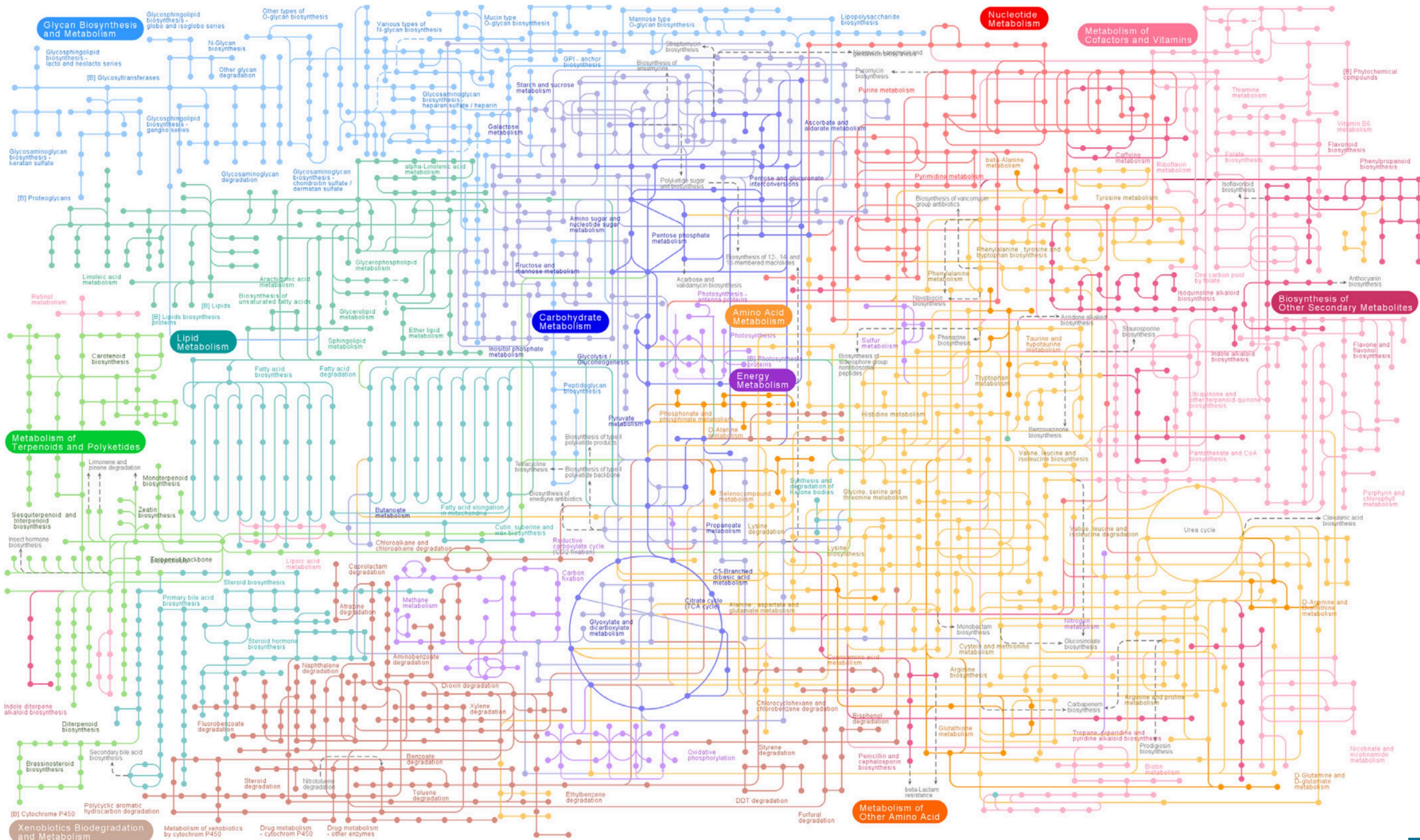
2. Functional process:
- x_i : single DNA segment
 - f : ideal DNA \rightarrow protein distribution
 - r : real DNA \rightarrow protein distribution
 - D_n : distribution over x

3. Goal: Learn the coding (DNA) and translation (RNA) change r to meet the ideal f .



$$\text{Perf}_f(r, D_n) = \sum_{x \in X_n} f(x)r(x)D_n(x)$$

How $r(x)$ is implemented in a cell



Structure of a learnable problem

What is learnable?

Learnability theory: Can the true h be learned?

Determine whether the true h is able to be learned given a particular data scenario (e.g., particular $P_{(X,Y)}$)?

Computational theory: Can the solution for the true h be computed?

P : Problem can be *computed* in polynomial time.

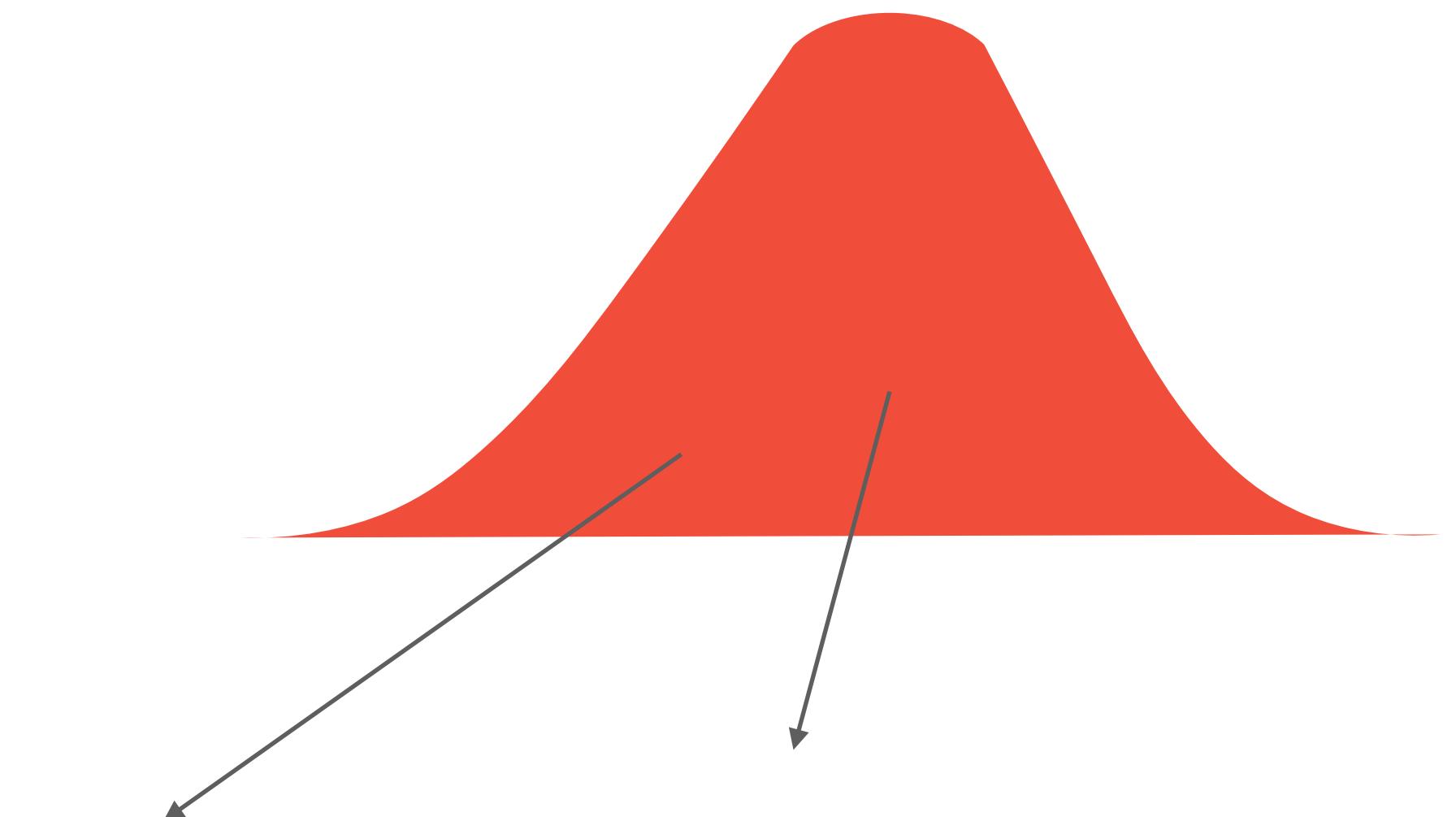
We can only work with
“ P hard” problems

NP : Problem requires non-deterministic polynomial time.

Empirical risk minimization

Expected Risk

$$\begin{aligned} E_{\text{risk}}(h, n, P) &= \int_{\underbrace{(\mathbf{X}, \mathbf{Y})_n}_{\text{train}}} \underbrace{R(h)}_{\text{risk}} \underbrace{dP_{(\mathbf{X}, \mathbf{Y})_n}}_{\text{distribution}} \\ &= \int_{\underbrace{(\mathbf{X}, \mathbf{Y})_n}_{\text{test}}} \int_{\underbrace{(\mathbf{X}, \mathbf{Y})}_{\text{train}}} \underbrace{\ell(h(X), Y)}_{\text{loss function}} \underbrace{dP_{X, Y}}_{\text{distribution}} \underbrace{dP_{(\mathbf{X}, \mathbf{Y})_n}}_{\text{distribution}} \end{aligned}$$



Assumption: Both the training and test data come from the same distribution.

Empirical risk minimization

$$h_{true} \rightarrow h_{best}$$

Generalization: Given a fixed **training data set**,
find the model that *best* predicts
future **unseen (test) data set**.

$$E_{\text{risk}}(h, n, P) = \underbrace{\int_{(\mathbf{X}, \mathbf{Y})_n} \underbrace{\int_{(\mathbf{X}, \mathbf{Y})}_{\text{training}}}_{\text{test}}}_{\text{loss function}} \underbrace{\ell(h(X), Y)}_{dP_{X,Y}} \underbrace{dP_{(X,Y)_n}}_{\text{distribution distribution}}$$

Probably Approximately Correct (PAC) Learning

Q: How do you get “good enough” learning so as to be useful?

PAC learning requires a learner to:

1. Approximate the true h
2. Be computationally feasible (P problem)

Approximately: A hypothesis $h \in H$ is approximately correct if its error over the training data $P_{(X,Y)}$ is bounded by some ϵ , with $0 \leq \epsilon \leq \frac{1}{2}$.

Probably: The h is probably approximately correct at a generalization error rate of δ , if its prediction accuracy is $1 - \delta$, with $0 \leq \delta \leq \frac{1}{2}$.

Probably Approximately Correct (PAC) Learning

Bounding the sample size:

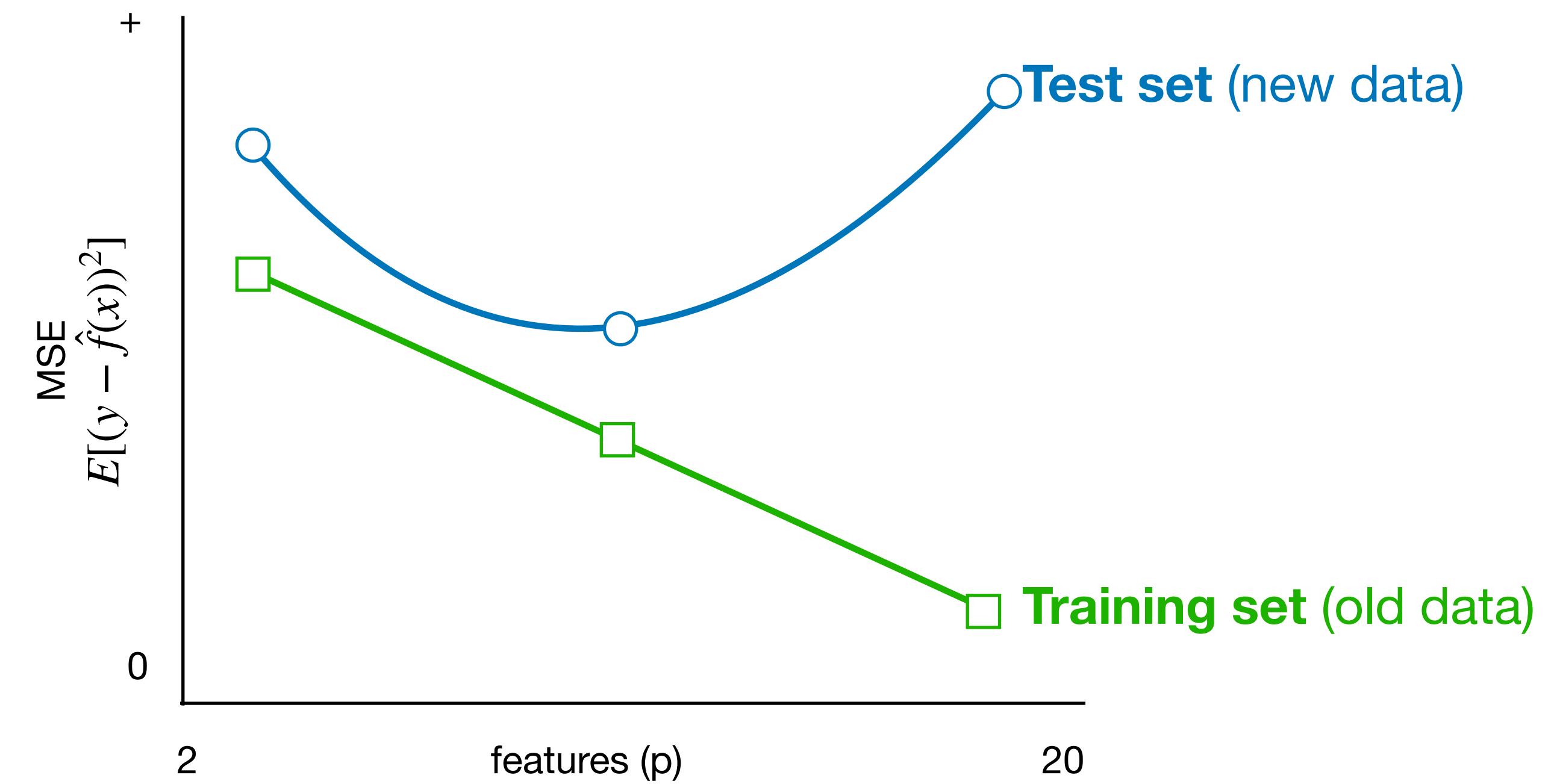
- as $\downarrow \epsilon$ & $\downarrow \delta$, or $\uparrow H$, then $\uparrow n$

$$n \geq \frac{1}{\epsilon} (\ln |H| + \ln \frac{1}{\delta})$$

e.g., complexity of model

The bias-variance tradeoff:

- $\delta \rightarrow$ Test error
- $\epsilon \rightarrow$ Training error



Evolvability

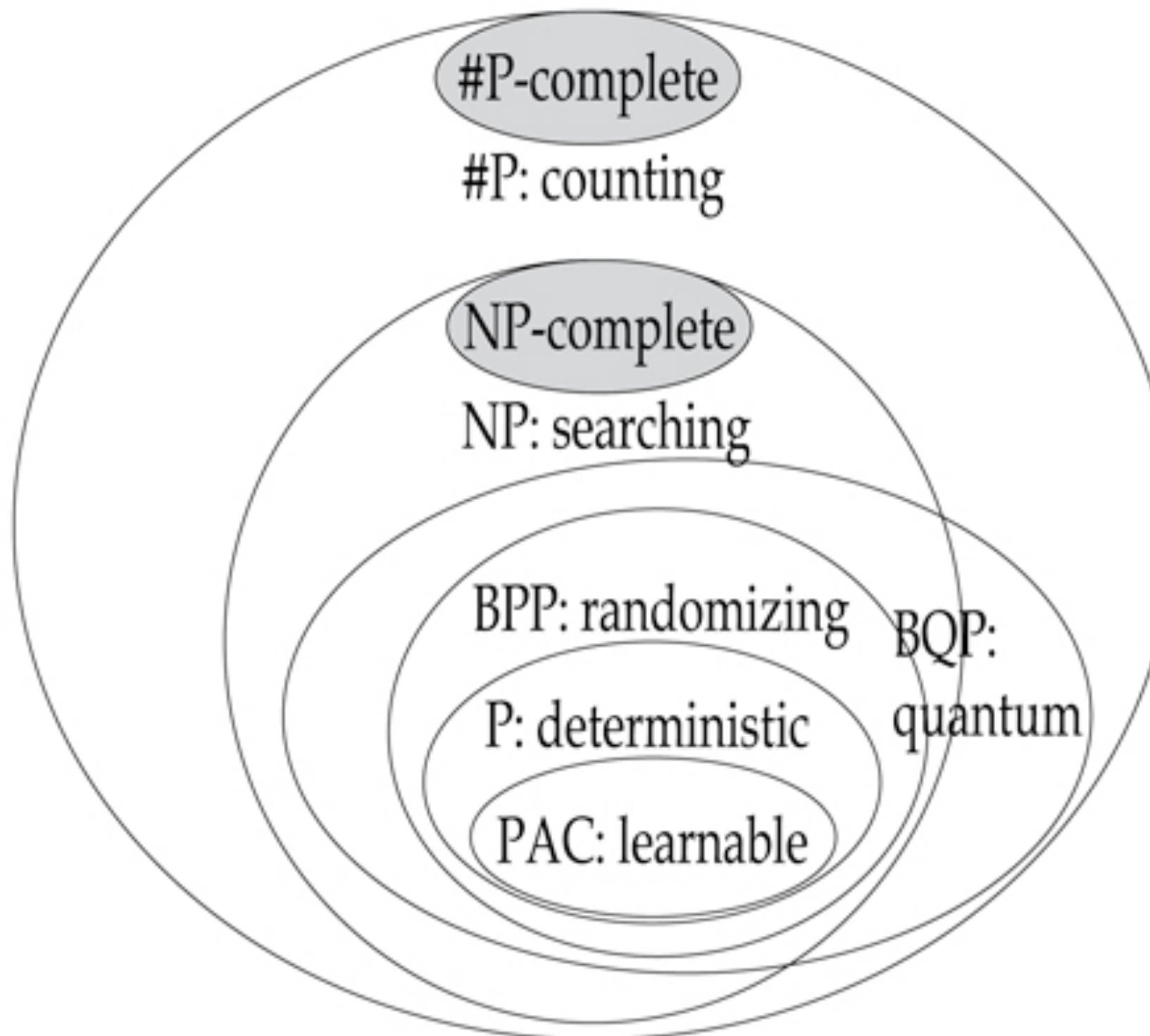
Definition: Let F be a class of functions, R be a class of representations (e.g., $r(x)$), and D a class of distributions on X . Then F is *evolvable* by R over D if there exists a set of *polynomial* translation, sample, selection, and tolerance functions such that for all n variables and all $\epsilon > 0$, for all ideal functions $f \in F_n$ and representations $r_0 \in R_n$, with probability at least $1 - \epsilon$

$$\text{Perf}(f, r_{g(n, \frac{1}{\epsilon})}) \geq 1 - \epsilon$$

where the sizes of neighbors $N(r)$ for $r \in R_n$ are at most $p(n, \frac{1}{\epsilon})$, samples size is $s(n, \frac{1}{\epsilon})$, tolerance is $t(\frac{1}{n}, \epsilon)$, and generation size is $g(n, \frac{1}{\epsilon})$.

Learnable evolutionary functions

Levels of complexity



Parity function: A Boolean function that has value true if and only if an odd number of its arguments have value true

$$f(x) = x_1 \oplus x_2 \oplus \dots \oplus x_n$$

“exclusive or”

Complexity: NP-hard

Con(dis)junction Function: A Boolean function that has value true if all of its arguments have value true (or at least one of its arguments are true for disjunctions).

Complexity: P-hard

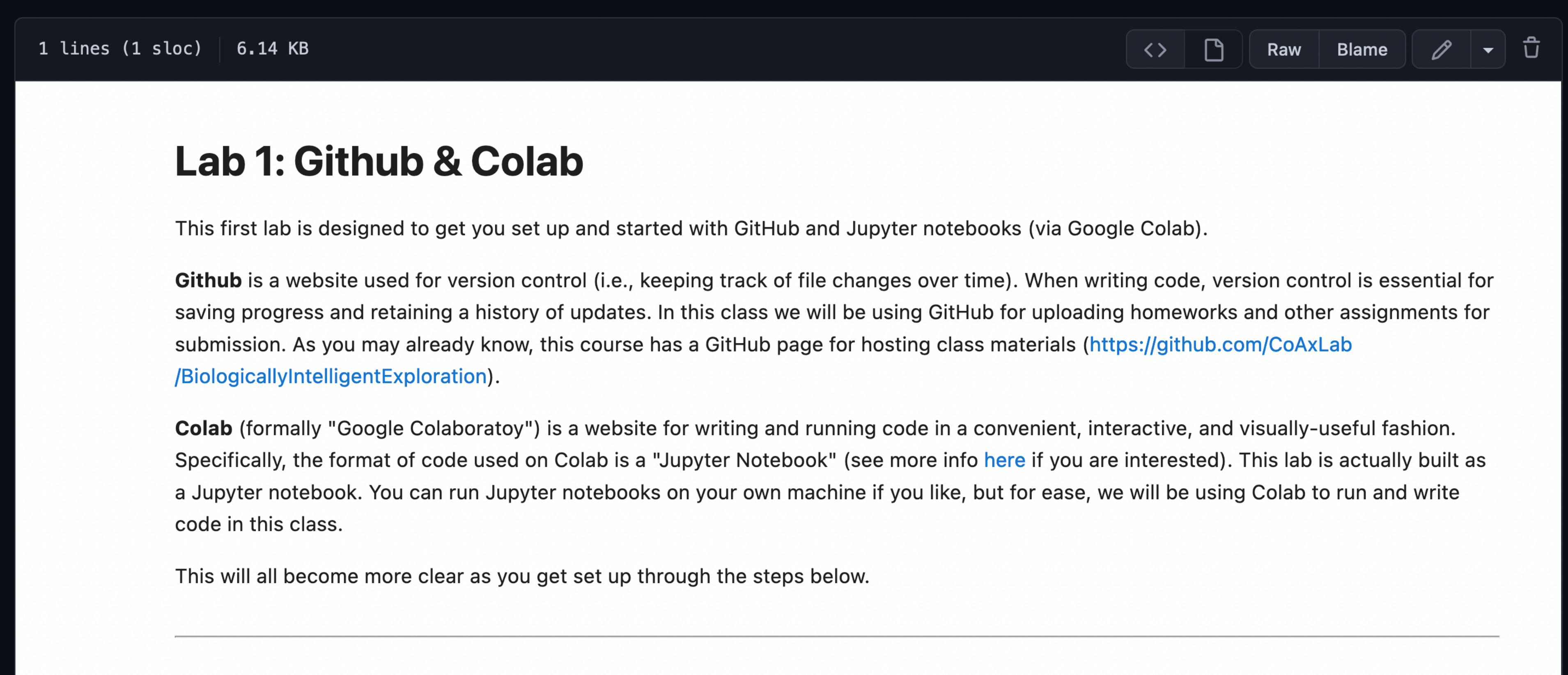
How we think the DNA to protein translation happens

Take home message

- Life as a random search problem is NP-hard and thus unlearnable.
- Life as a structured learning problem is P-hard and learnable in a reasonable (and quick) amount of time.
- Thus, life is a learning process.

Lab 1: Github and Colab

URL: <https://github.com/CoAxLab/BiologicallyIntelligentExploration/tree/main/Labs>



The screenshot shows a GitHub file viewer with a dark theme. At the top, it displays "1 lines (1 sloc) | 6.14 KB". To the right are standard GitHub file operations: a diff icon, a copy icon, "Raw" (which is currently selected), "Blame", a pen icon for editing, a dropdown menu, and a trash bin icon.

Lab 1: Github & Colab

This first lab is designed to get you set up and started with GitHub and Jupyter notebooks (via Google Colab).

GitHub is a website used for version control (i.e., keeping track of file changes over time). When writing code, version control is essential for saving progress and retaining a history of updates. In this class we will be using GitHub for uploading homeworks and other assignments for submission. As you may already know, this course has a GitHub page for hosting class materials (<https://github.com/CoAxLab/BiologicallyIntelligentExploration>).

Colab (formally "Google Colaboratory") is a website for writing and running code in a convenient, interactive, and visually-useful fashion. Specifically, the format of code used on Colab is a "Jupyter Notebook" (see more info [here](#) if you are interested). This lab is actually built as a Jupyter notebook. You can run Jupyter notebooks on your own machine if you like, but for ease, we will be using Colab to run and write code in this class.

This will all become more clear as you get set up through the steps below.
