

AI LAB 2023

1.1 - INTRODUCTION

FRANCESCA M. BUFFA



LAB STRUCTURE

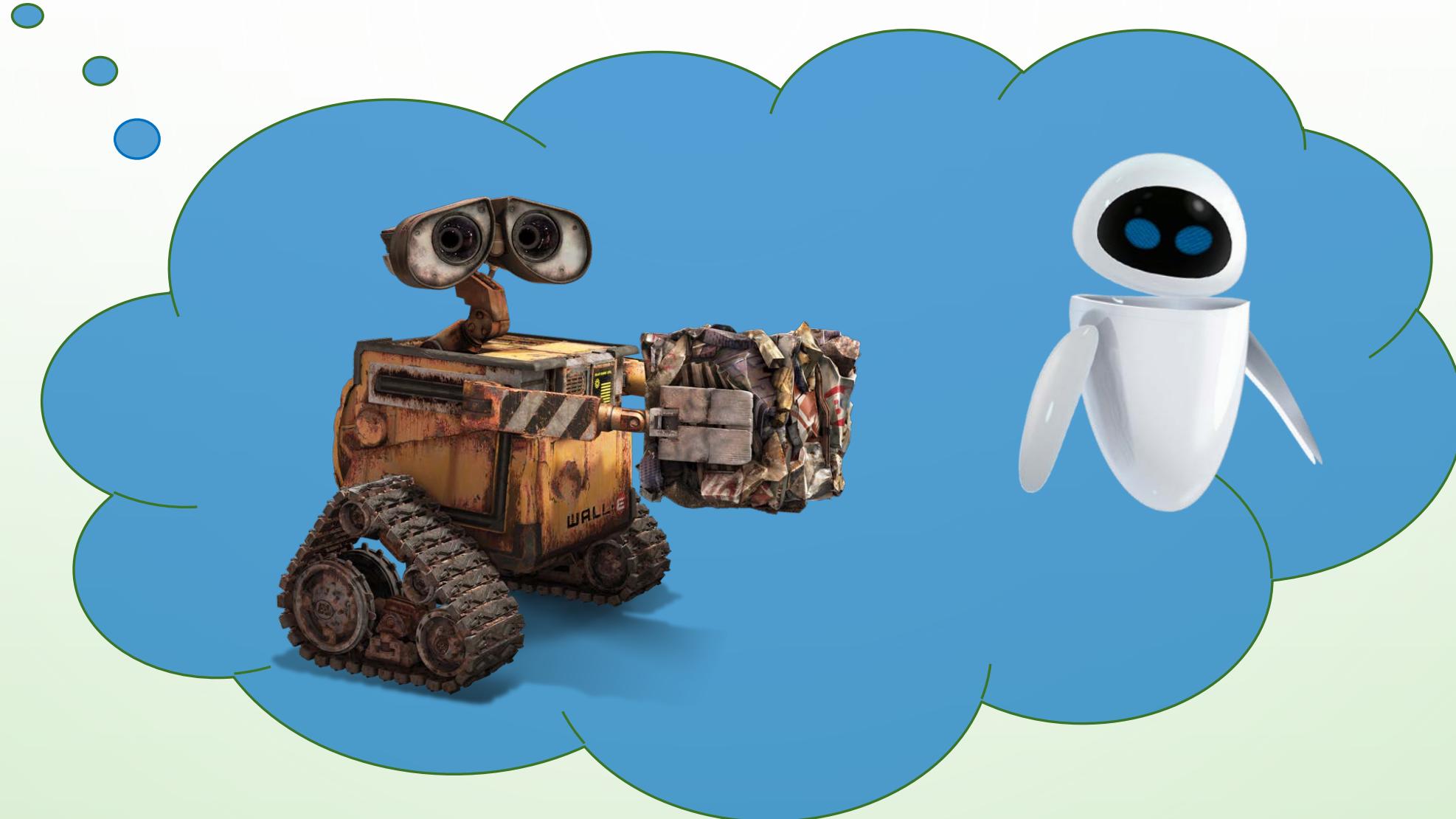
- INTRODUCTION
- THE DATA
- THE AI-LAB CHALLENGE – PART 1
- PART 1 - SHARING AND DISCUSSION

- UNSUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 2
- PART 2 - SHARING AND DISCUSSION

- SUPERVISED LEARNING EXAMPLES
- THE AI-LAB CHALLENGE – PART 3
- LARGE PROJECTS AND DATABASES

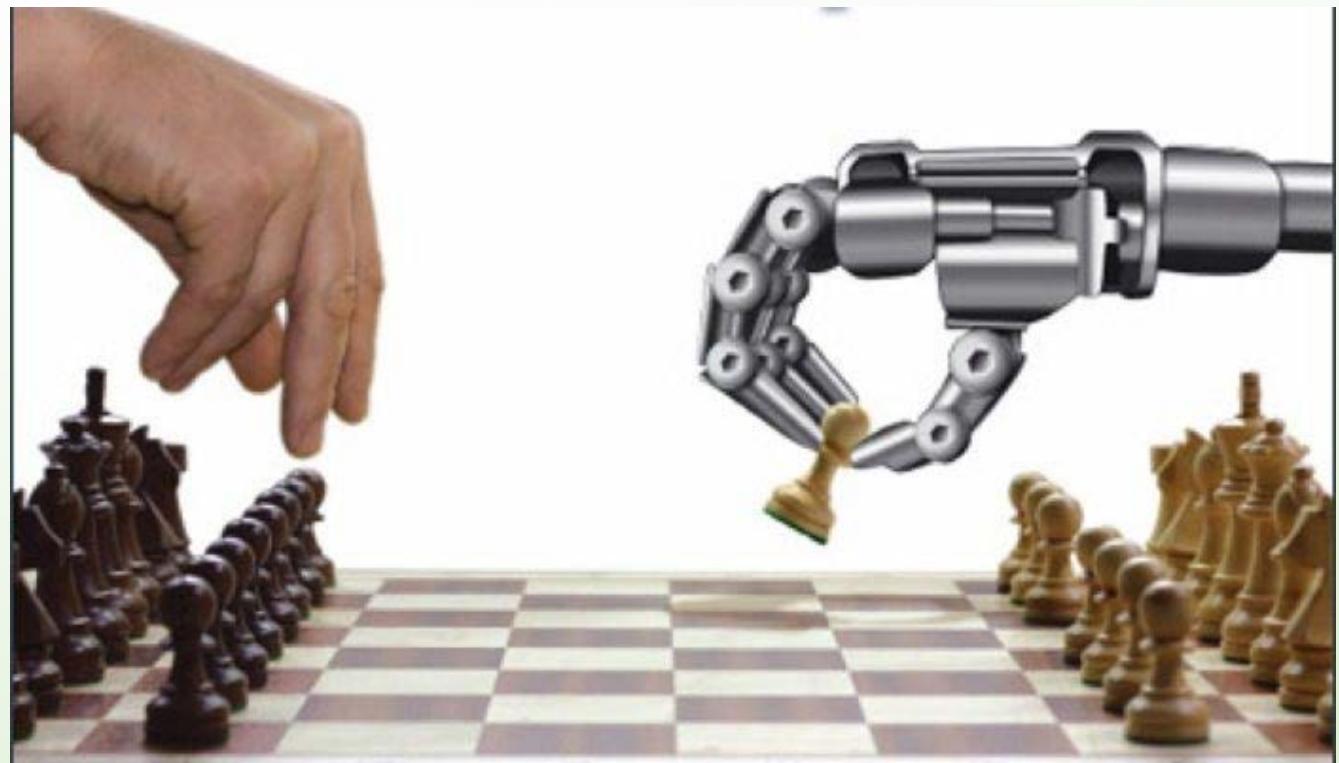
- THE AI-LAB CHALLENGE PARTS 1-3, SHARING AND DISCUSSION
- DATA INTERPRETATION
- DISCUSS AND PREPARE WORKSHOP PRESENTATIONS

AI

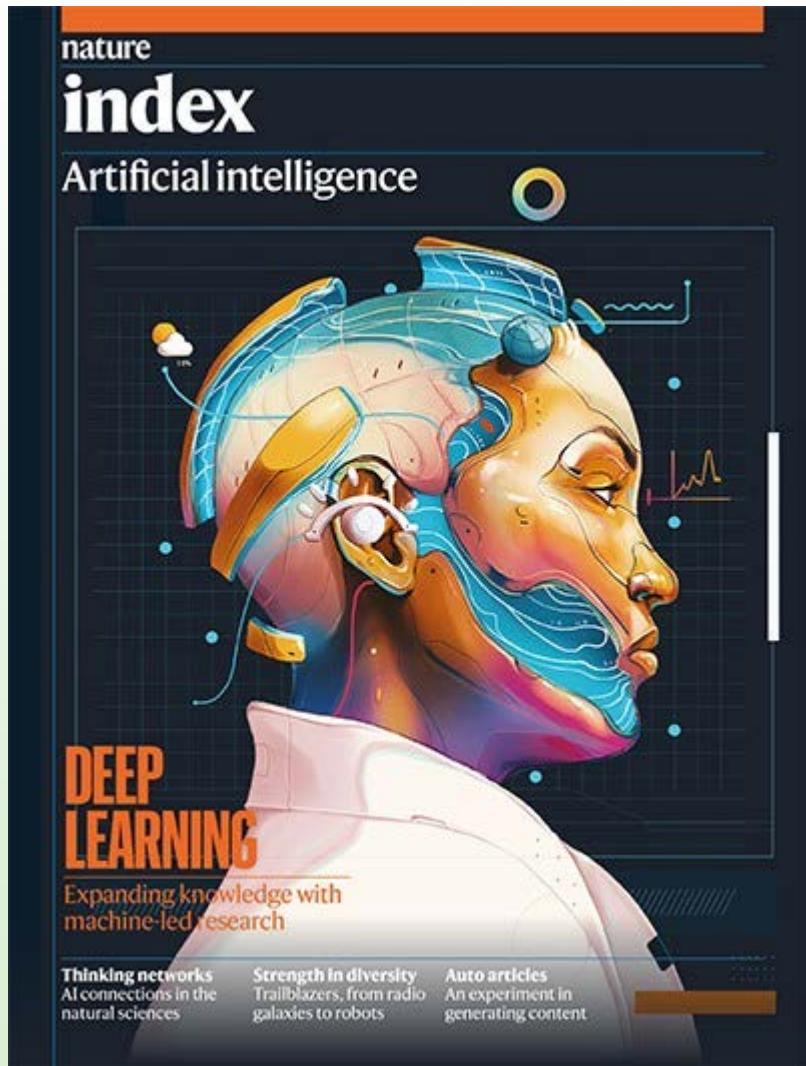


“NARROW” AI

DEVELOPED FOR ONE SPECIFIC TASK

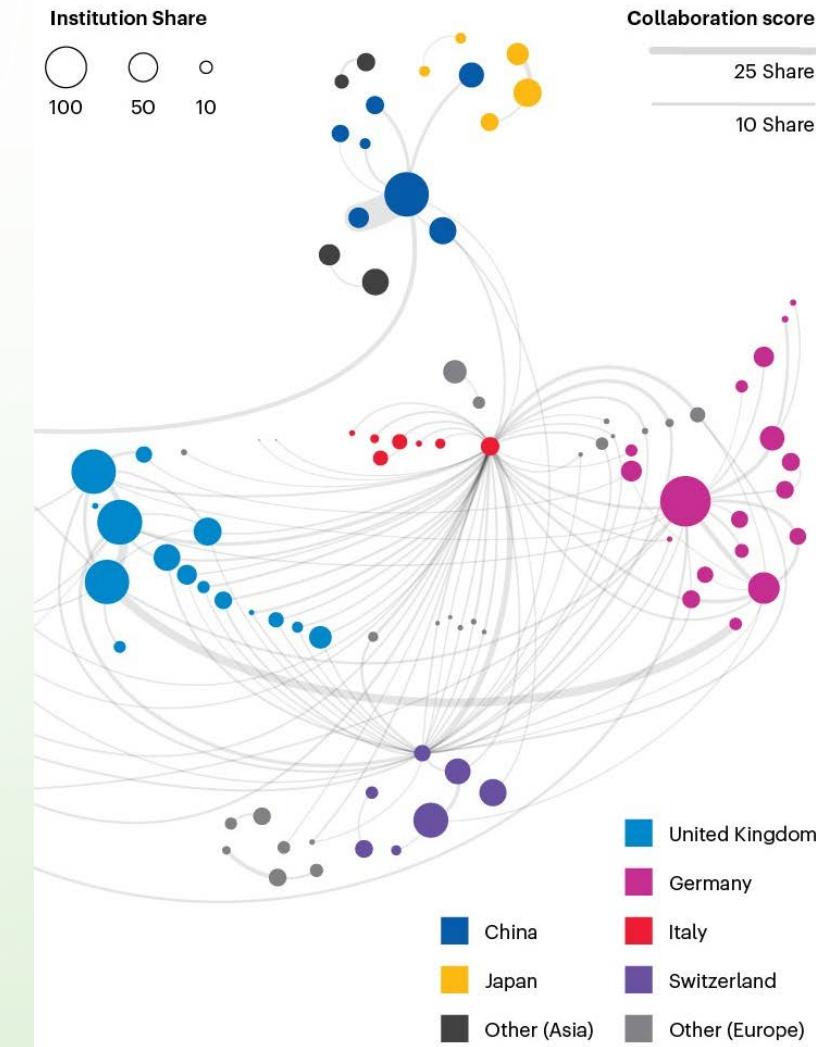


AI: THE INTERNATIONAL LANDSCAPE



<https://www.nature.com/articles/d41586-020-03409-8>

Top 200 collaborations (lines) among 146 institutions (circles) between 2015 and 2019. Circles sized according to institution's share in AI. The thickness of the lines corresponds to a particular collaboration's total share between two institutions.



<https://www.nature.com/collections/hchdibjchj>

TASKS THAT ARE HARDER FOR AI:

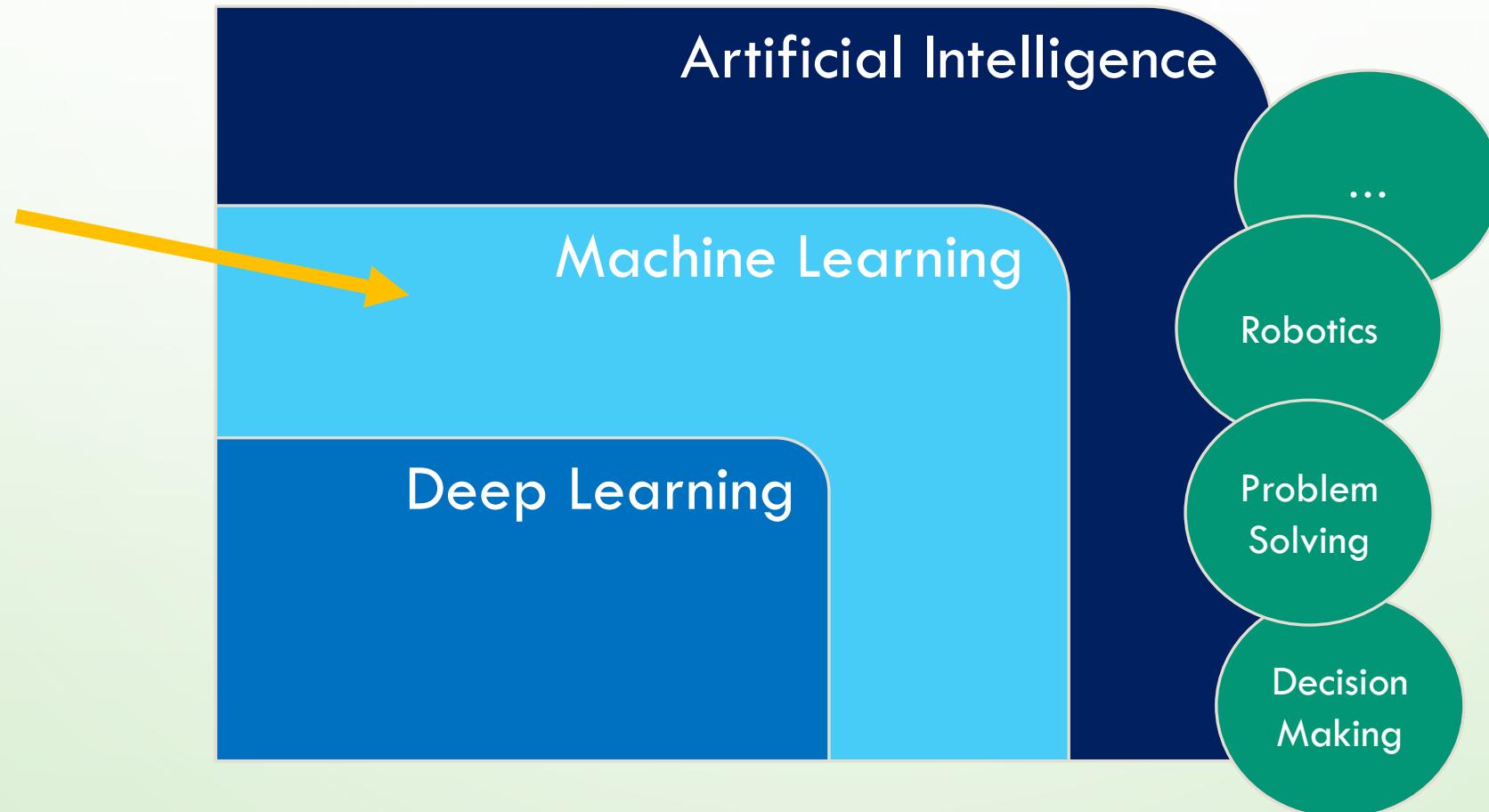
- UNSTRUCTURED DATA
 - NEED STRUCTURED DATA TO OPERATE EFFECTIVELY (EG. NUMERICAL DATA IN A DATABASE).
 - IMAGES, VIDEOS, AND TEXT DOCUMENTS ARE HARDER
- AMBIGUITY AND UNCERTAINTY
 - INCOMPLETE, LIMITED OR CONFLICTING INFORMATION CAN MAKE IT DIFFICULT FOR AI TO MAKE ACCURATE PREDICTIONS
- CREATIVITY AND INTUITION
 - NEW IDEAS AND SOLUTIONS ARE POSSIBLE BUT LIMITED AS AI RELIES ON PRE-EXISTING DATA AND PATTERNS
- CONTEXTUAL UNDERSTANDING
 - SOCIAL, CULTURAL, AND HISTORICAL FACTORS THAT SHAPE HUMAN BEHAVIOUR
 - HARD TO INTERPRET HUMAN LANGUAGE, BEHAVIOUR, AND INTERACTIONS.
- LEARNING FROM LIMITED DATA
 - NEED LARGE AMOUNTS OF DATA TO TRAIN AI MODELS
 - PERFORMANCE NOT EASY TO IMPROVE WHEN DATA ARE LIMITED OR INCOMPLETE
 - INACCURATE OR BIASED PREDICTIONS

TASKS WHERE AI CAN OUTPERFORM HUMANS:

- COMPLEX CALCULATIONS
 - AI IS FASTER AND MORE ACCURATELY THAN HUMANS
- DATA ANALYSIS
 - AI CAN PROCESS LARGE AMOUNTS OF DATA QUICKLY AND ACCURATELY
- PATTERN RECOGNITION
 - AI CAN IDENTIFY PATTERNS IN DATA OR IMAGES THAT MAY BE DIFFICULT FOR HUMANS TO DISCERN, OR TOO RARE
- REPETITIVE TASKS
 - ASSEMBLY DATA ENTRY WITHOUT GETTING BORED OR MAKING MISTAKES
- OBJECT TRACKING
 - TRACK OBJECTS AND MOVEMENTS ACROSS LARGE AREAS, SECURITY SURVEILLANCE OR TRAFFIC MONITORING
- PLAYING GAMES
 - ABLE TO BEAT HUMAN PLAYERS AT COMPLEX GAMES SUCH AS CHESS, GO, AND POKER.

OVERALL, TASKS THAT REQUIRE: PROCESSING LARGE AMOUNTS OF DATA, PERFORMING COMPLEX CALCULATIONS OR SIMULATIONS, IDENTIFYING PATTERNS RARE OR DIFFICULT TO DETECT PATTERNS.

AI, MACHINE LEARNING



APPLICATION OF MACHINE LEARNING

What do you want the machine learning system to do?

I want to see if there are natural clusters or dimensions in the data I have about different situations.

I want to learn what actions to take in different situations.

Do you want the ML system to be active or passive?

ACTIVE

The system's own actions will affect the situations it sees in the future.

PASSIVE

The system will learn from data I give it.

Yes

Could a knowledgeable human decide what actions to take based on the data you have about the situation?

No

Do you have access to data that describes a lot of examples of situations and appropriate actions for each situation?

Yes

Could there be patterns in these situations that humans haven't recognized before?

No

Yes

Will the system be able to gather a lot of data by trying sequences of actions in many different situations and seeing the results?

No

Yes

UNSUPERVISED
LEARNING MAY BE APPROPRIATE

clustering
anomaly detection

SUPERVISED
LEARNING MAY BE APPROPRIATE

neural nets
support vector machines
regression
recommender systems

MACHINE LEARNING IS NOT USEFUL

REINFORCEMENT
LEARNING MAY BE APPROPRIATE

TYPICAL STEPS IN MACHINE LEARNING APPLICATIONS

- UNDERSTAND THE APPLICATION CONTEXT
- ASK A QUESTION
- GENERATE OR COLLECT THE DATA
- LOOK AT THE DATA
- UNDERSTAND THE DATA
- CLEAN AND PROCESS THE DATA
- SELECT ALGORITHM OR ALGORITHMS
- PLAN THE TRAINING APPROACH
- TRAIN AND VALIDATE YOUR MODEL
- FINALIZE YOUR MODEL
- TEST YOUR MODEL BY PREDICTING ON NEW DATA

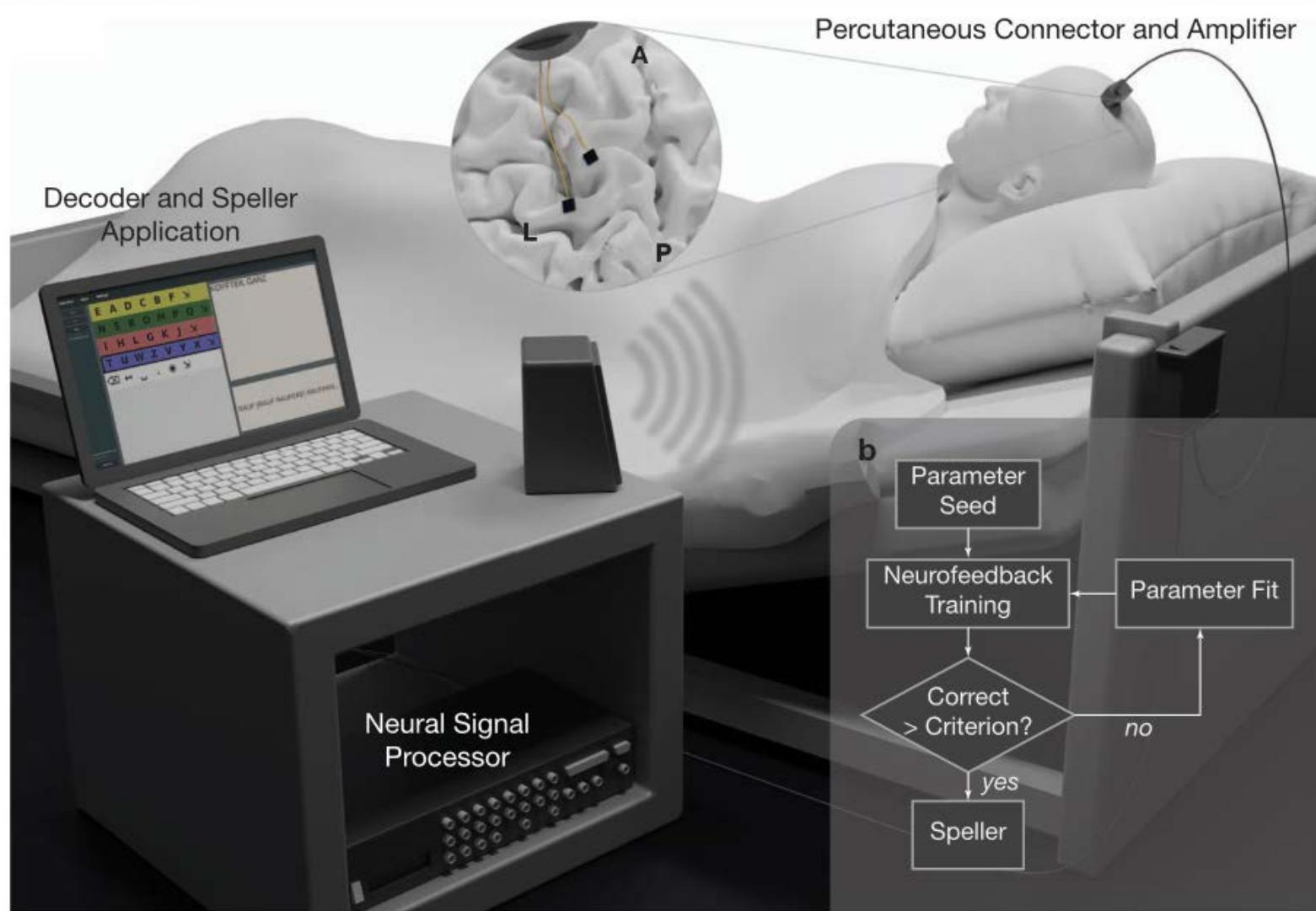
UNDERSTAND THE CONTEXT OF APPLICATION

- YOUR CONTEXT WILL BE LIFE SCIENCES AND HEALTH
- WHY DID WE CHOOSE THIS CONTEXT?
- INTRODUCTION TO THE BASIC CONCEPTS
- EXAMPLES OF PREVIOUS AI APPLICATIONS

HEALTH MAJOR APPLICATION FOR AI

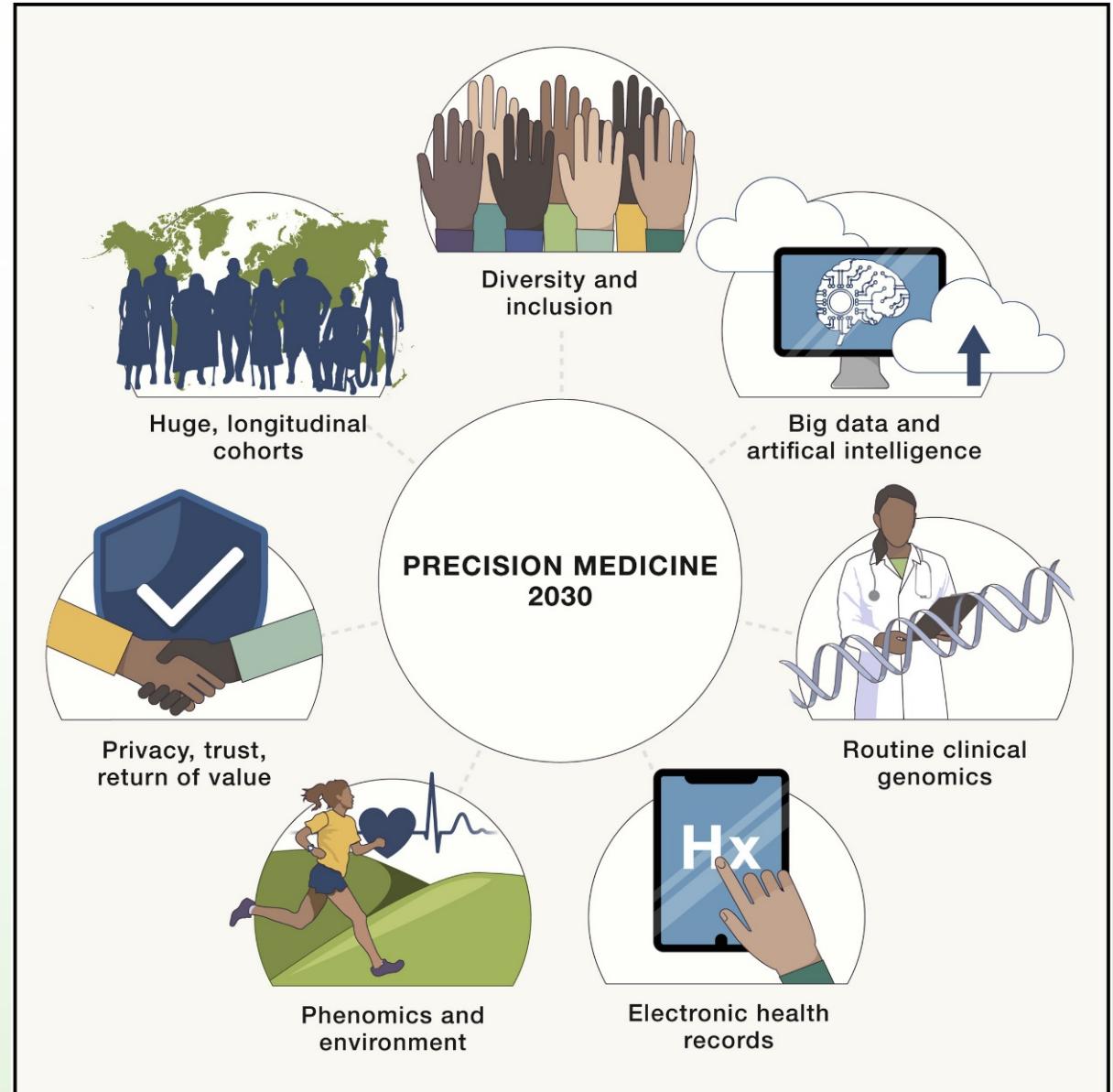
APPLICATION	POTENTIAL ANNUAL VALUE BY 2026	KEY DRIVERS FOR ADOPTION
Robot-assisted surgery	\$40B	Technological advances in robotic solutions for more types of surgery
Virtual nursing assistants	20	Increasing pressure caused by medical labor shortage
Administrative workflow	18	Easier integration with existing technology infrastructure
Fraud detection	17	Need to address increasingly complex service and payment fraud attempts
Dosage error reduction	16	Prevalence of medical errors, which leads to tangible penalties
Connected machines	14	Proliferation of connected machines/devices
Clinical trial participation	13	Patent cliff; plethora of data; outcomes-driven approach
Preliminary diagnosis	5	Interoperability/data architecture to enhance accuracy
Automated image diagnosis	3	Storage capacity; greater trust in AI technology
Cybersecurity	2	Increase in breaches; pressure to protect health data

SPELLING INTERFACE USING INTRACORTICAL SIGNALS ENABLED VIA AUDITORY NEUROFEEDBACK TRAINING



AI IN HEALTH

- New technologies enable to obtain and store large amount of biological and medical data.
- Both the quantity and type of data collected has changed dramatically, and is very diverse.
- To achieve impact on health we need to integrate and transform such data in information.
- This requires collaboration between disciplines: medicine, biology, genomics, engineering, ethics, law, big data, maths, stats, AI.



16 February 2001

Science

Vol. 291 No. 5507
Pages 1145–1434 \$9

THE HUMAN GENOME



AMERICAN ASSOCIATION FOR THE ADVANCEMENT OF SCIENCE

15 February 2001

nature

www.nature.com

the human genome

Nuclear fission
Five-dimensional energy landscapes

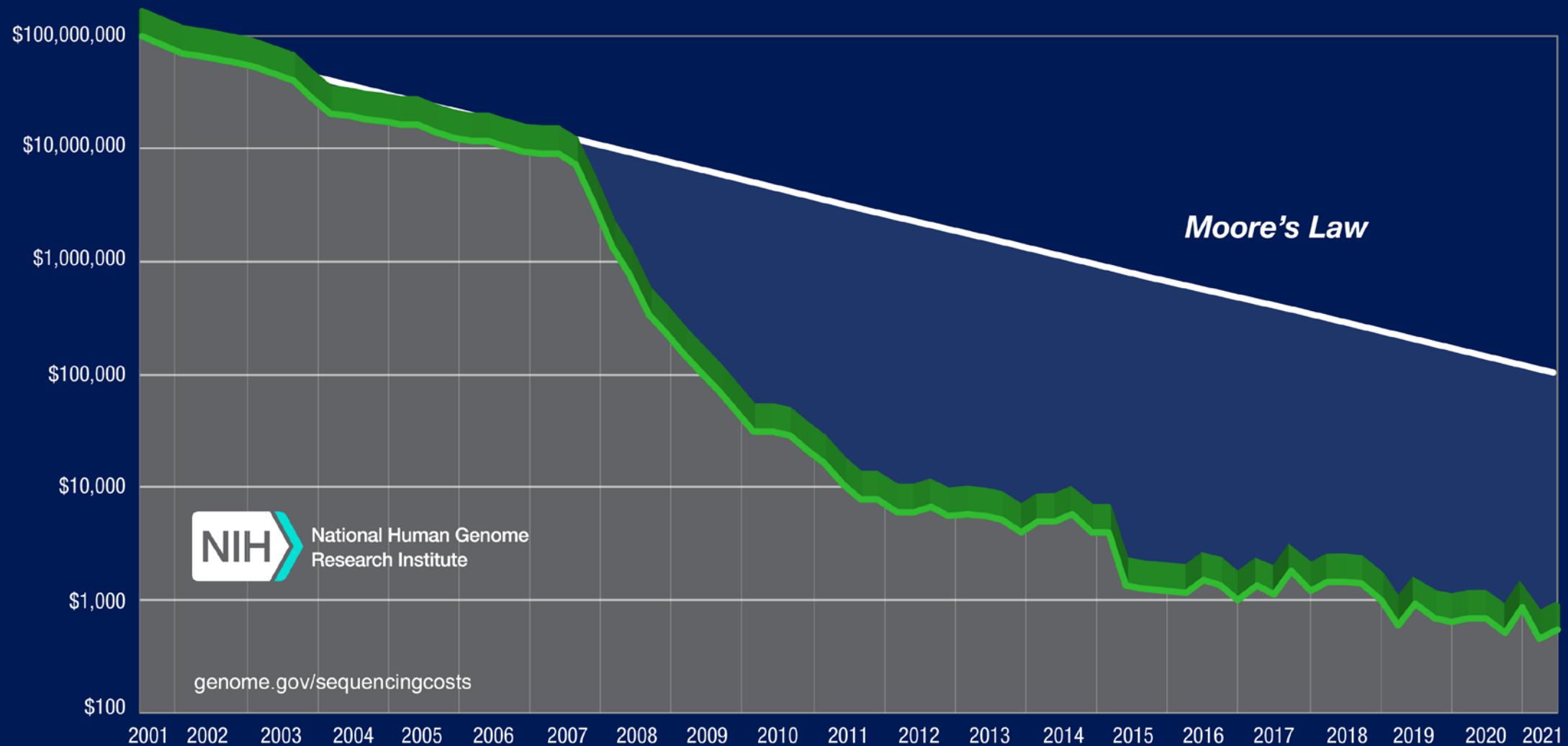
Seafloor spreading
The view from under the Arctic ice

Career prospects
Sequence creates new opportunities

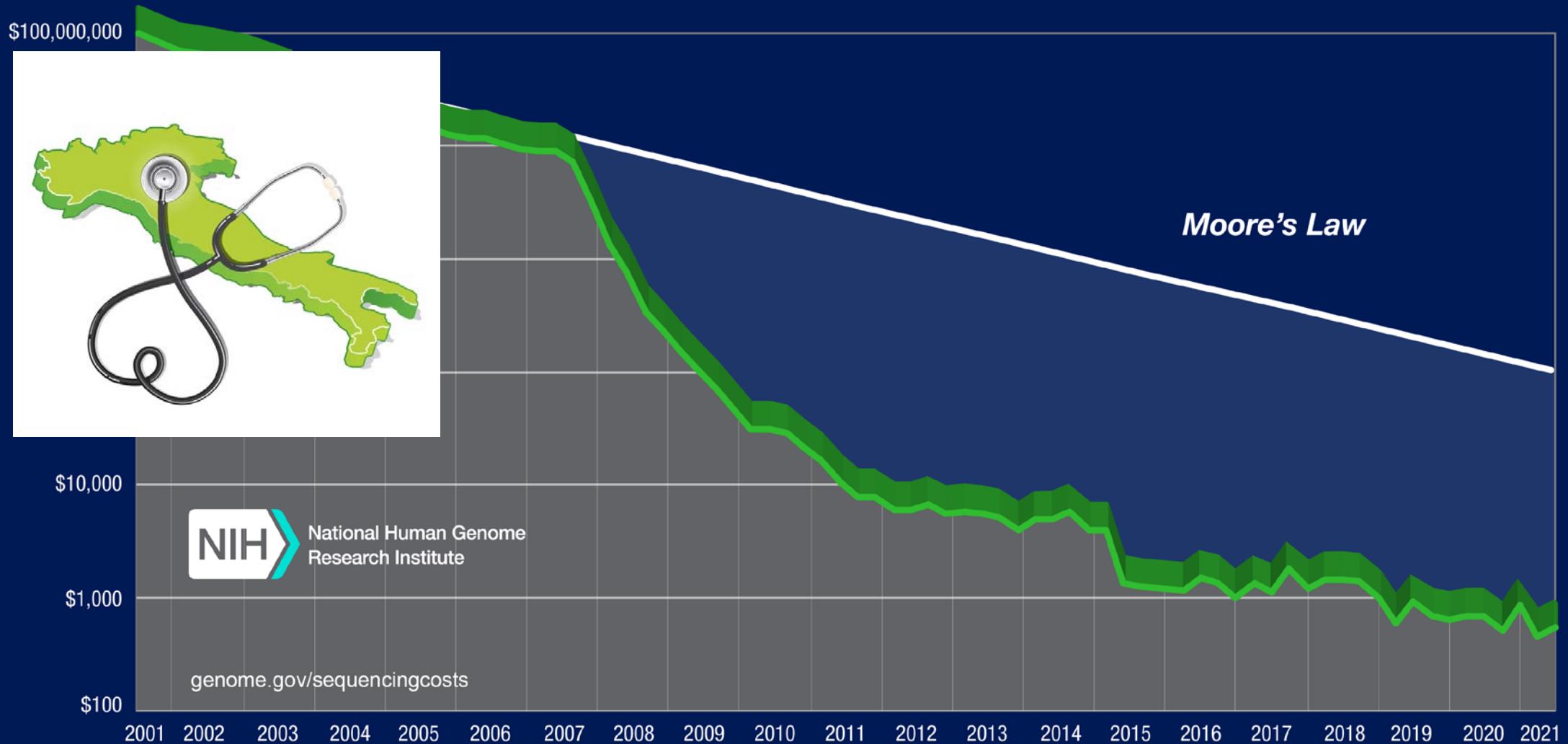
naturejobs

genomics special

Cost per Human Genome



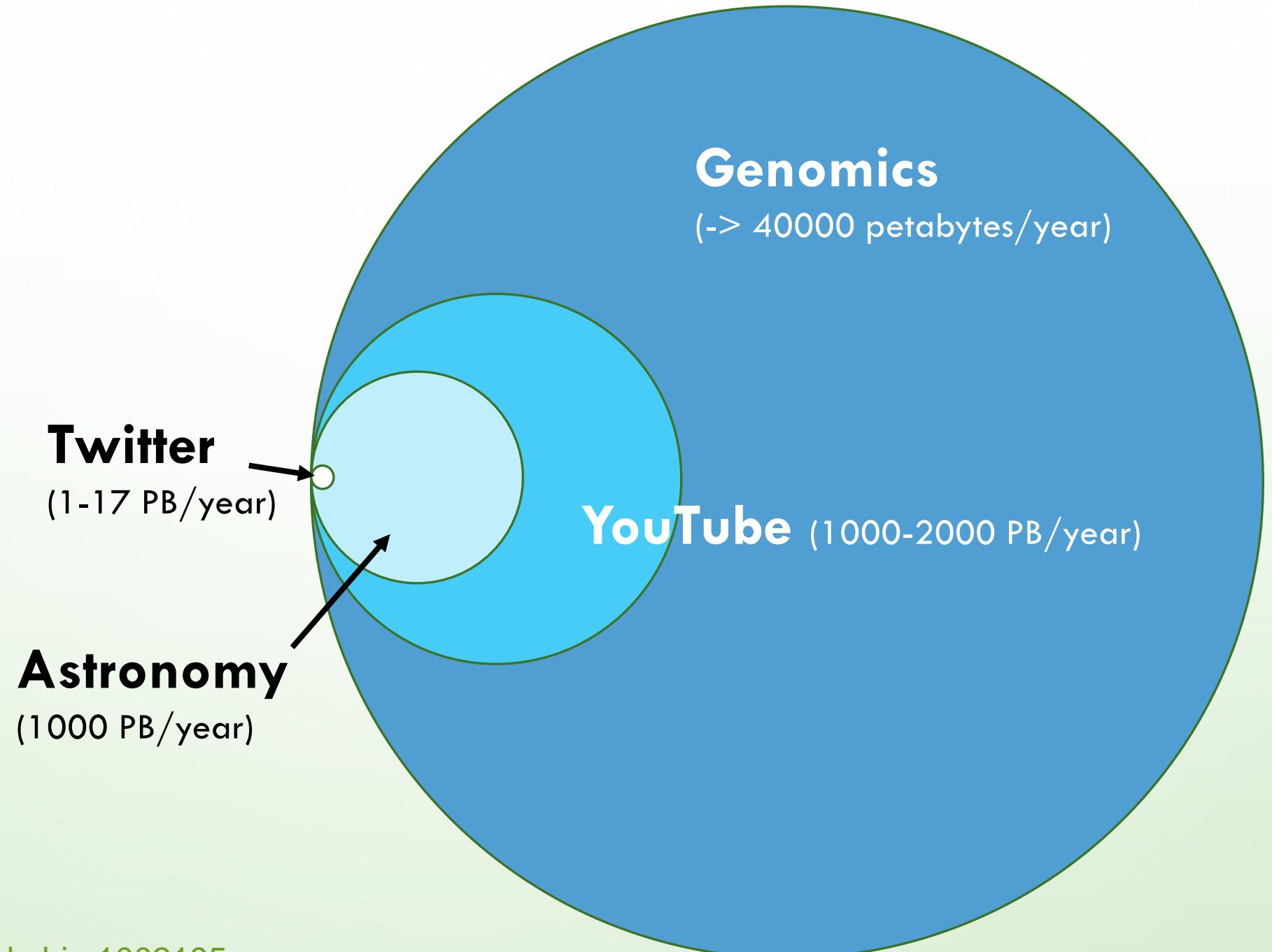
Cost per Human Genome



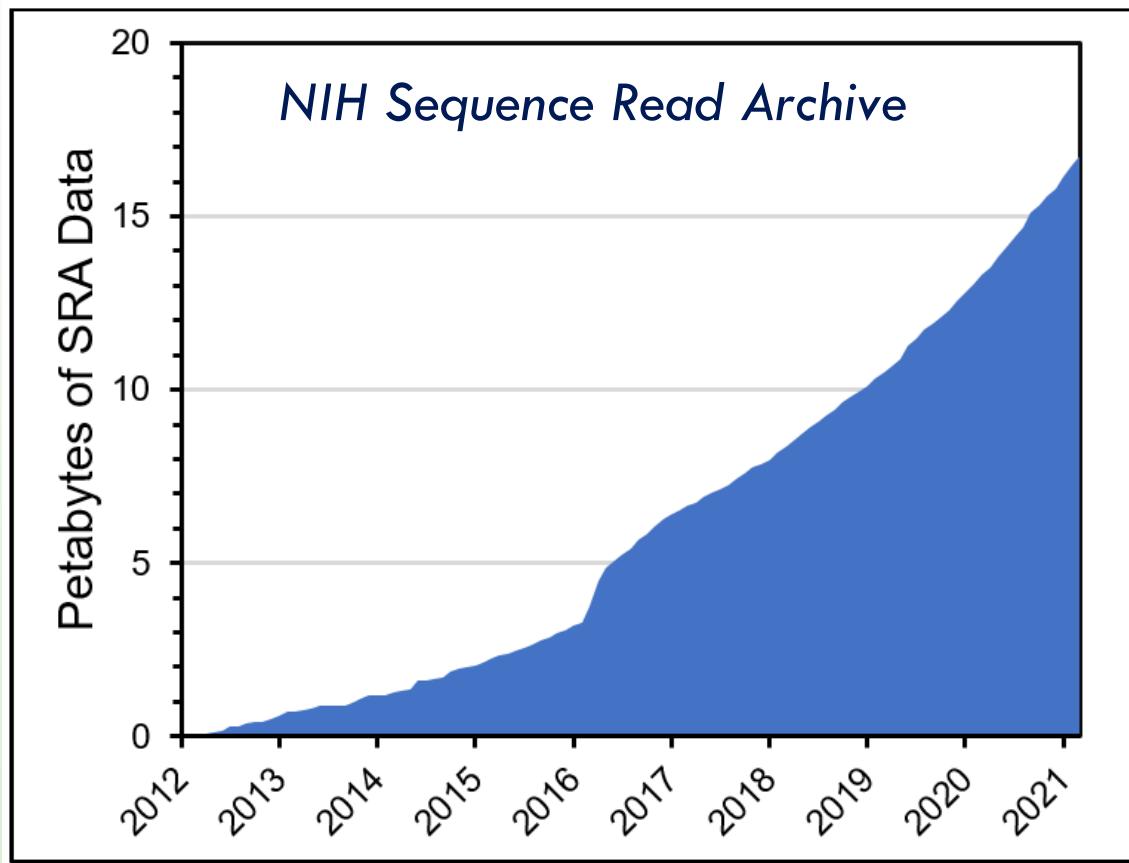
Cost per Human Genome



Big data



TACKLING PETABYTE SCALE SEQUENCE SEARCH CHALLENGES



Dr. Susan Gregurick, NIH Associate Director for Data Science and ODSS Director, said: “We all share a common problem and a need to develop, enhance, and implement methods that streamline data access, search or findability, and ultimately data reuse.”





Your genes

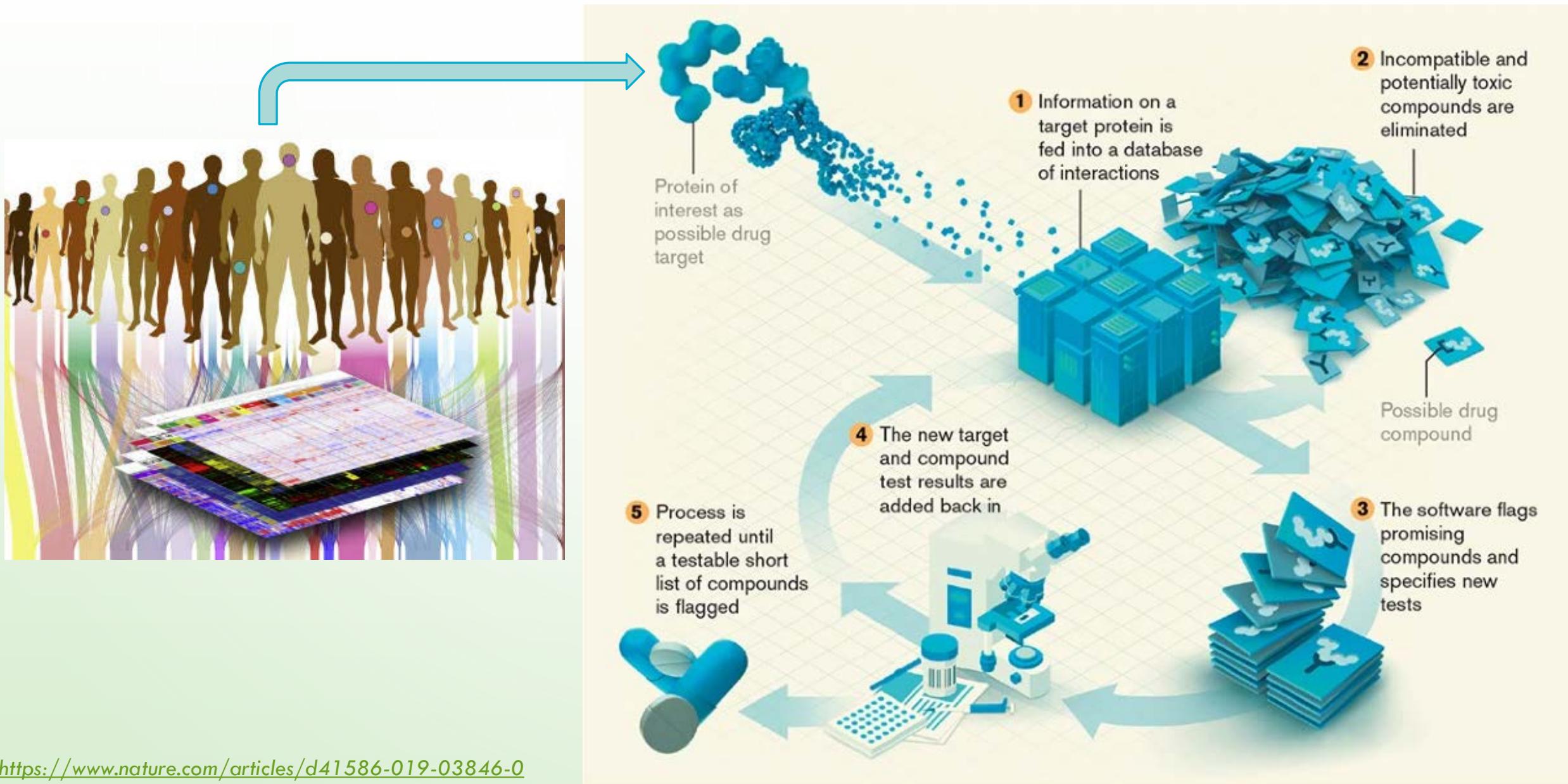


Your
environment
& lifestyle

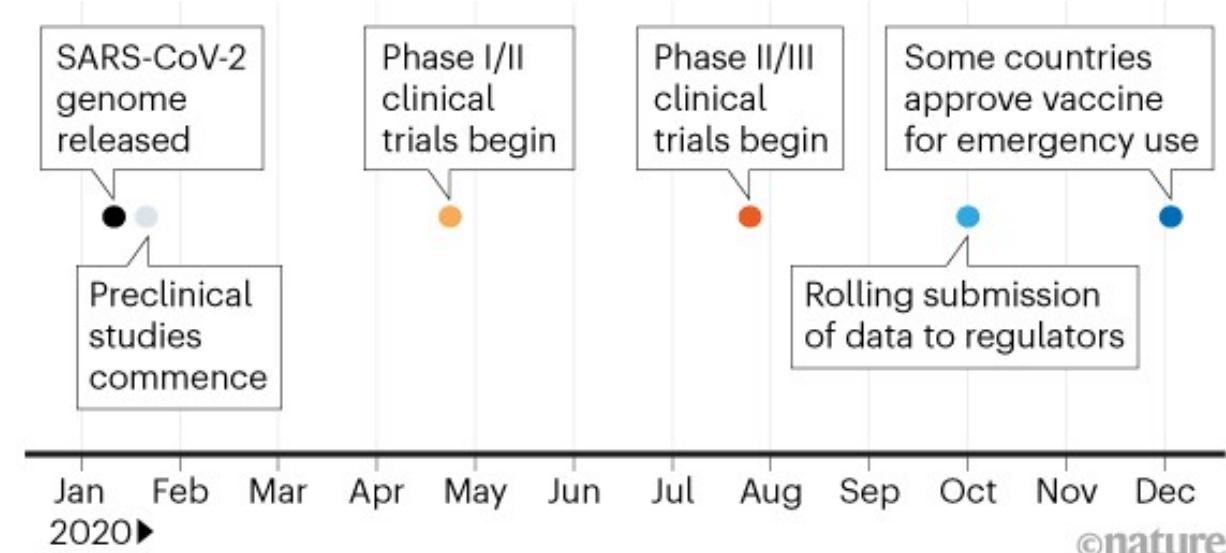
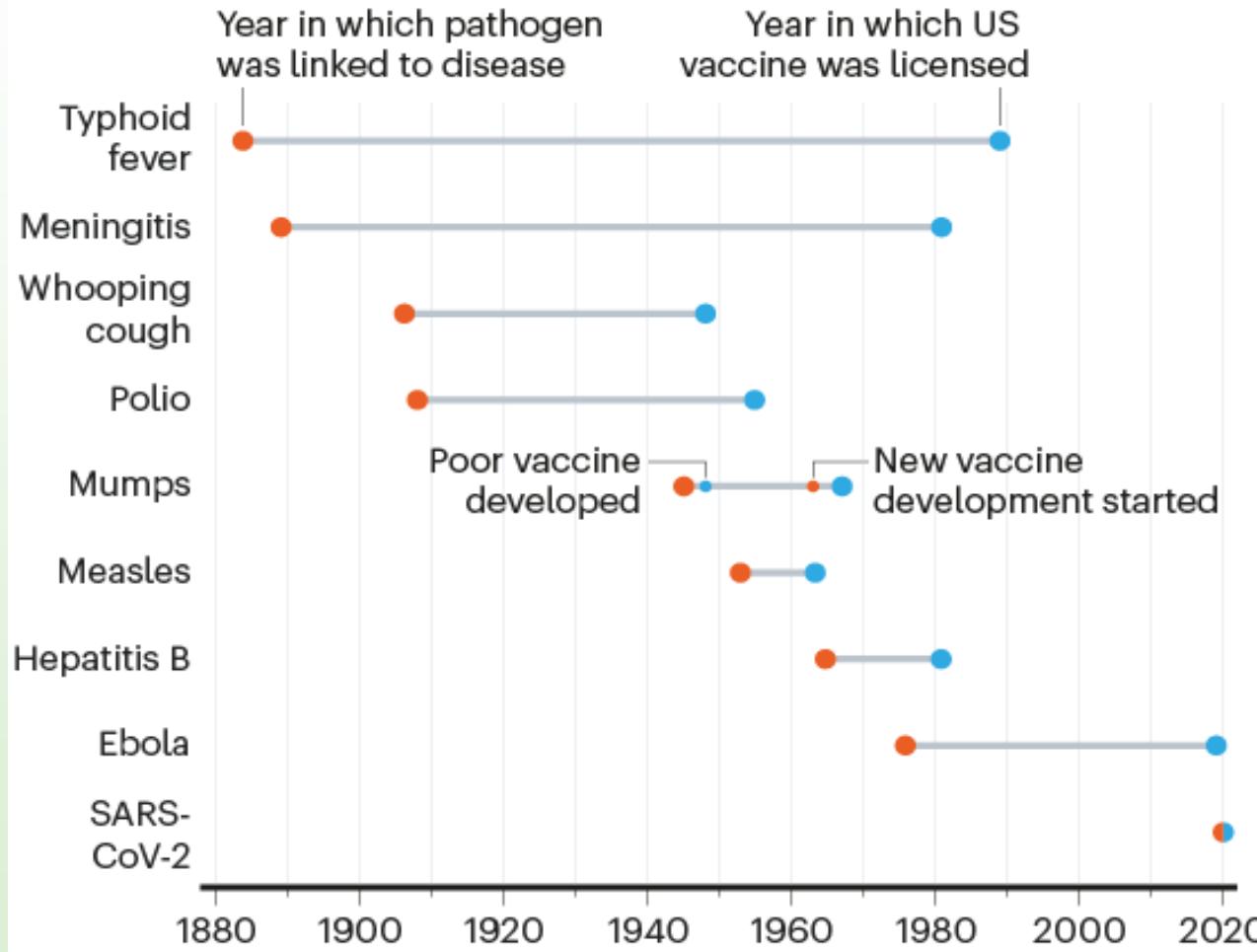


YOU!

SPEEDING UP THE SEARCH FOR DRUGS WITH AI



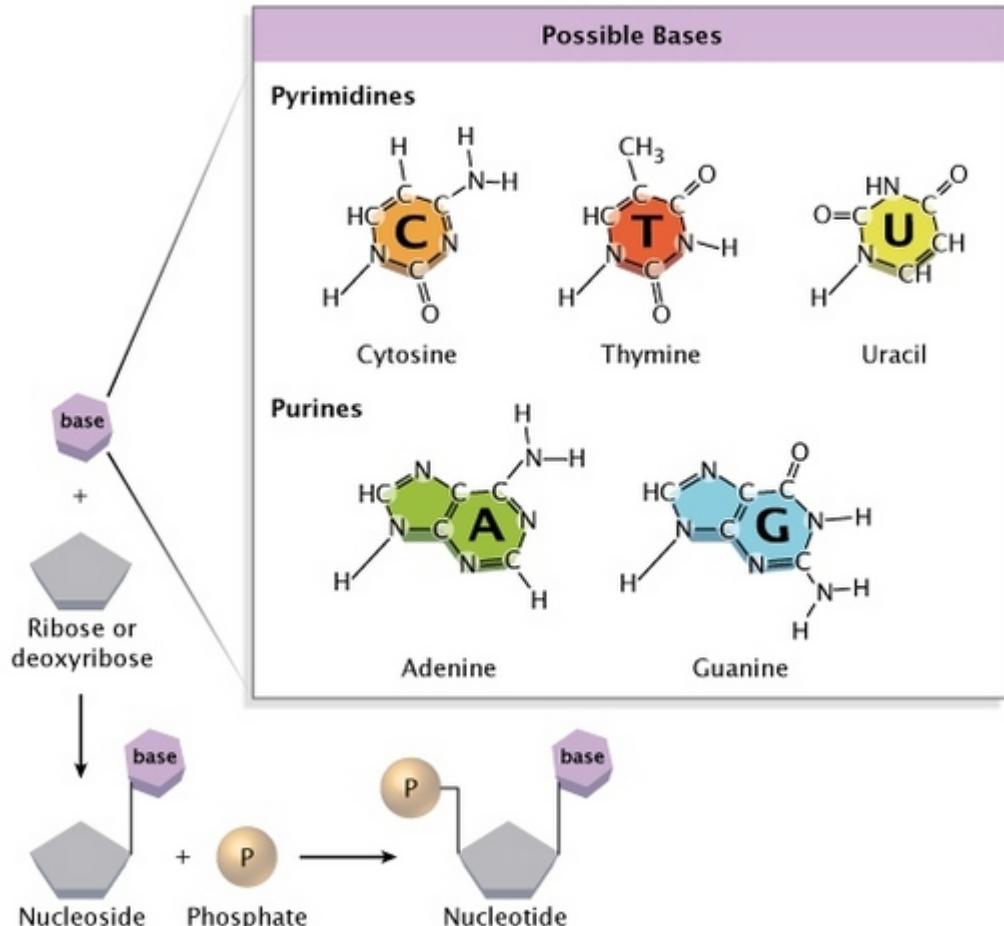
AI and vaccine development



Dave Johnson, PhD in Information Physics, chief data and AI officer at Moderna: “We’ve seen how this digital infrastructure and how these algorithms can really help push things forward”.

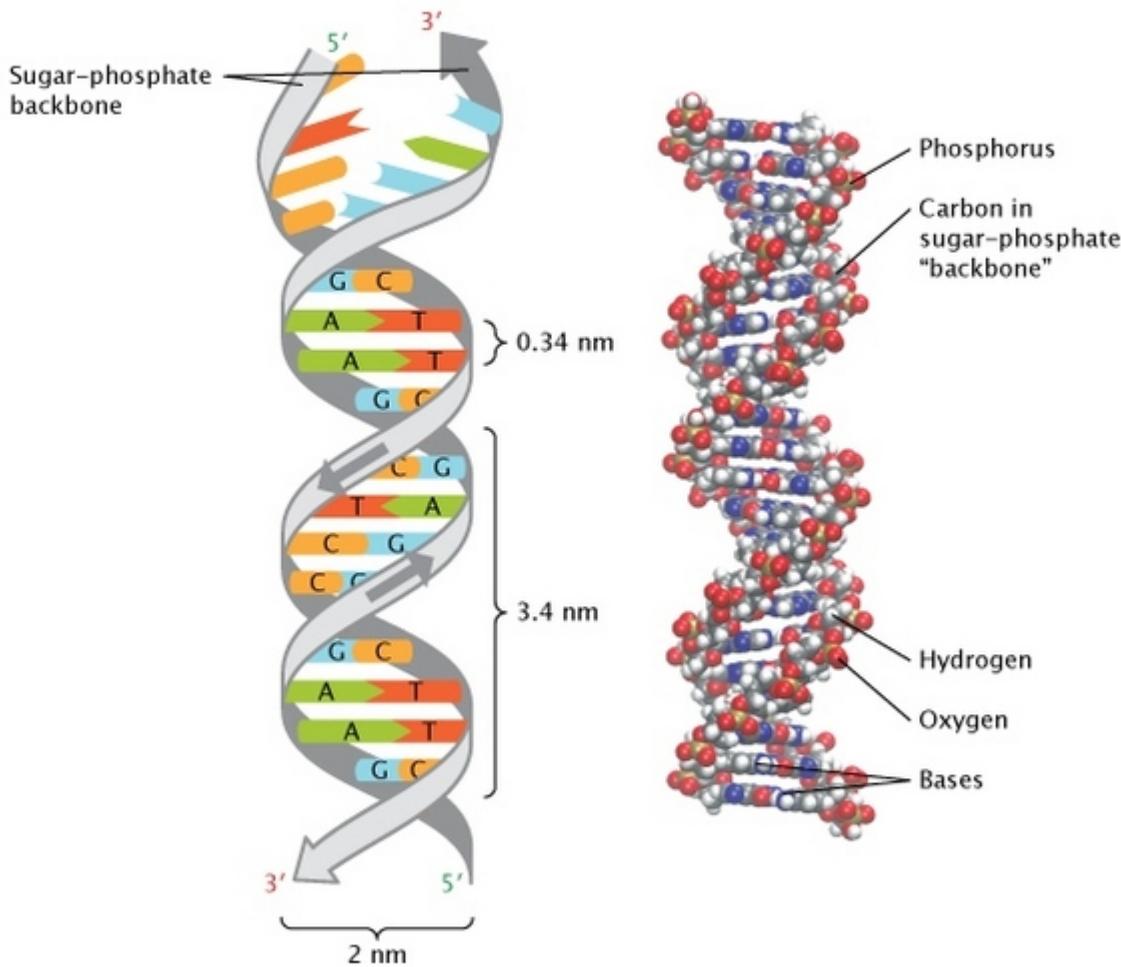
<https://sloanreview.mit.edu/audio/ai-and-the-covid-19-vaccine-modernas-dave-johnson/>

DNA/RNA STRING REPRESENTATION



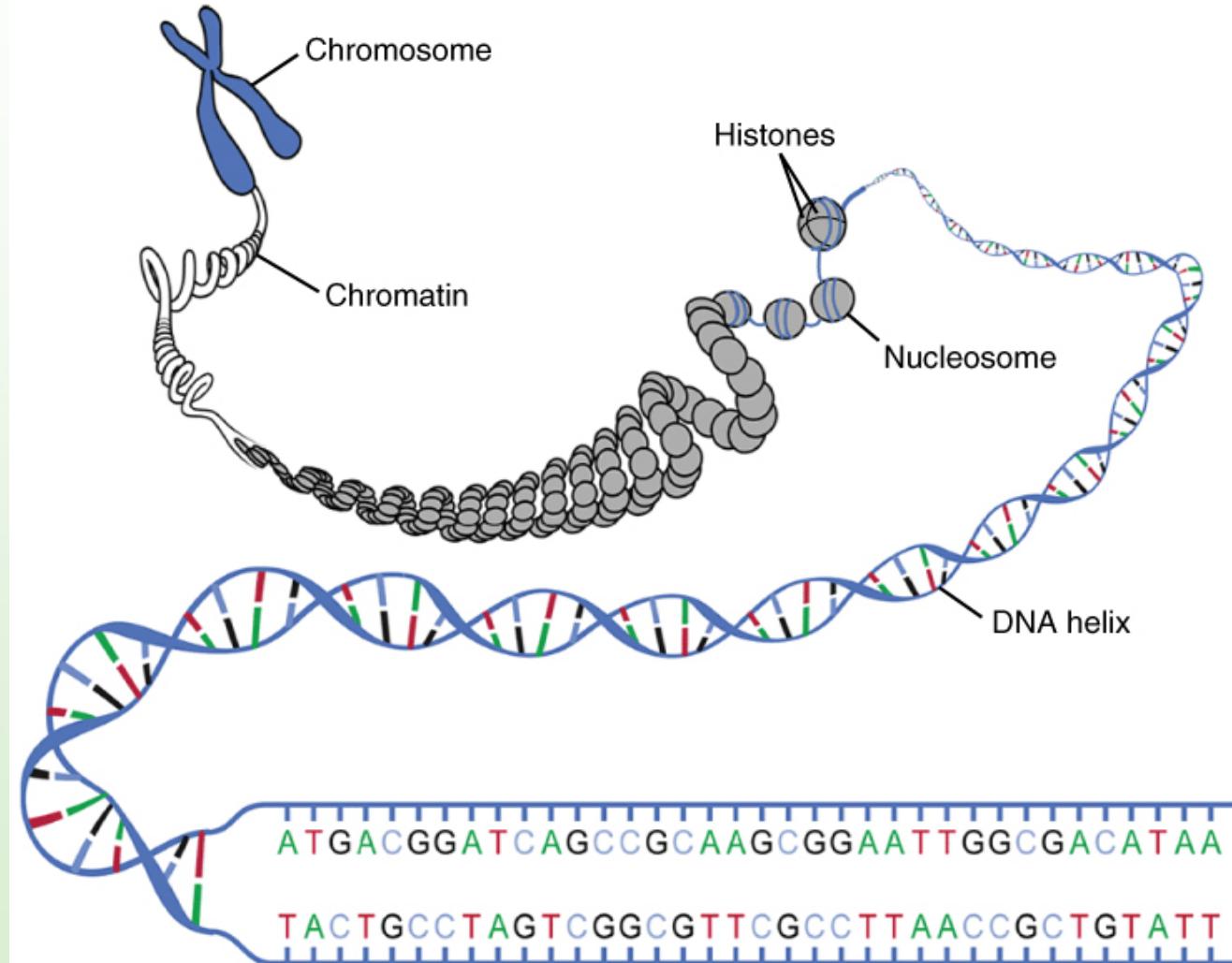
- A (ADENINE)
- T (THYMINE)
- C (CYTOSINE)
- G (GUANINE)
- ORGANIZED IN A DOUBLE HELIX OF PAIRED BASES: A-T, C-G
- IN RNA T -> U (URACIL)

DNA/RNA STRING REPRESENTATION

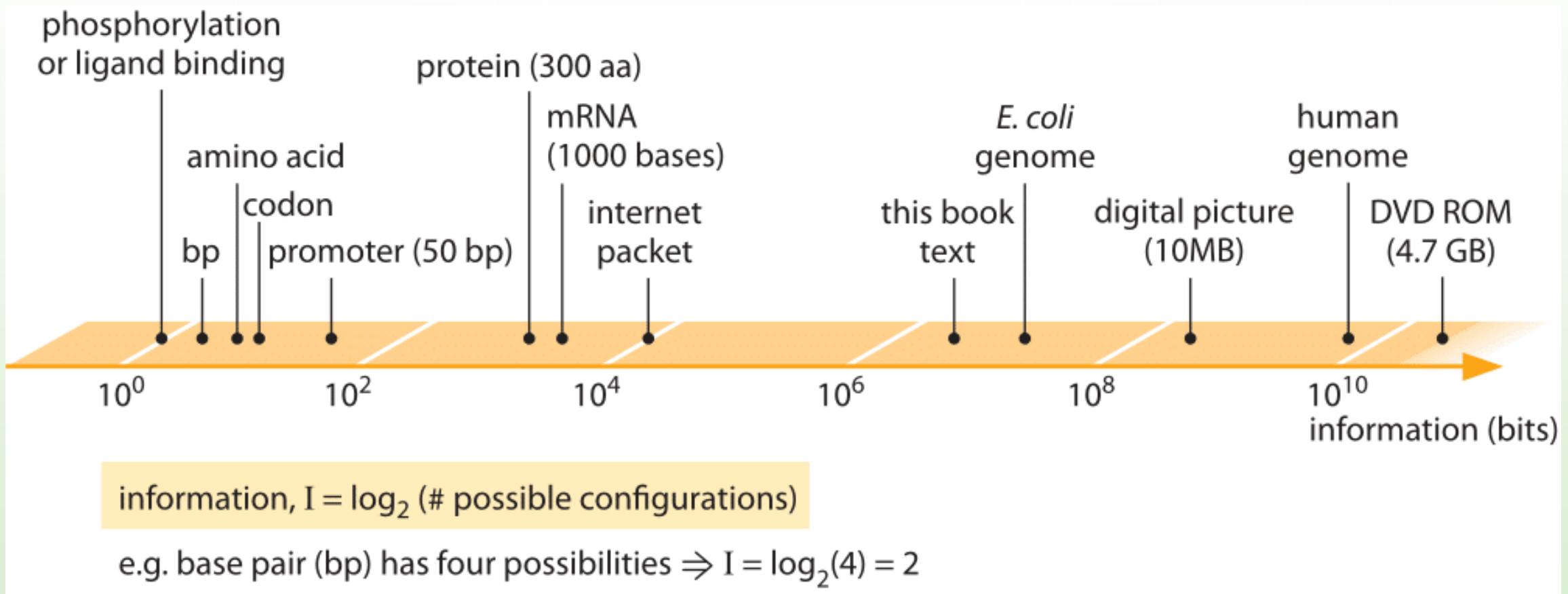


- A (ADENINE)
- T (THYMINE)
- C (CYTOSINE)
- G (GUANINE)
- ORGANIZED IN A DOUBLE HELIX OF PAIRED BASES: A-T, C-G
- IN RNA T -> U (URACIL)

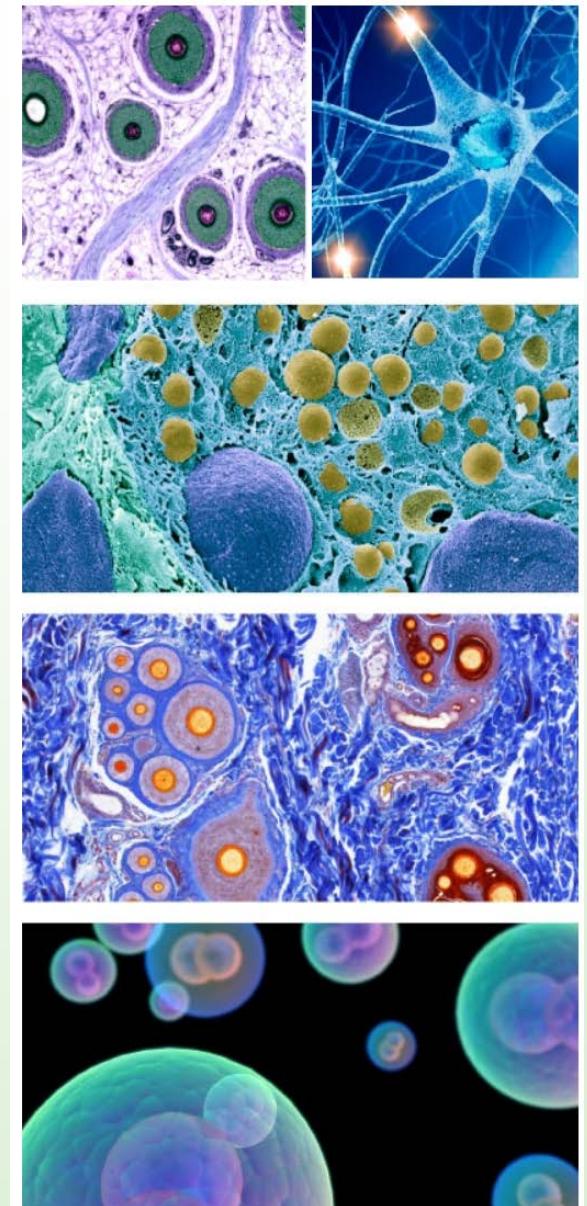
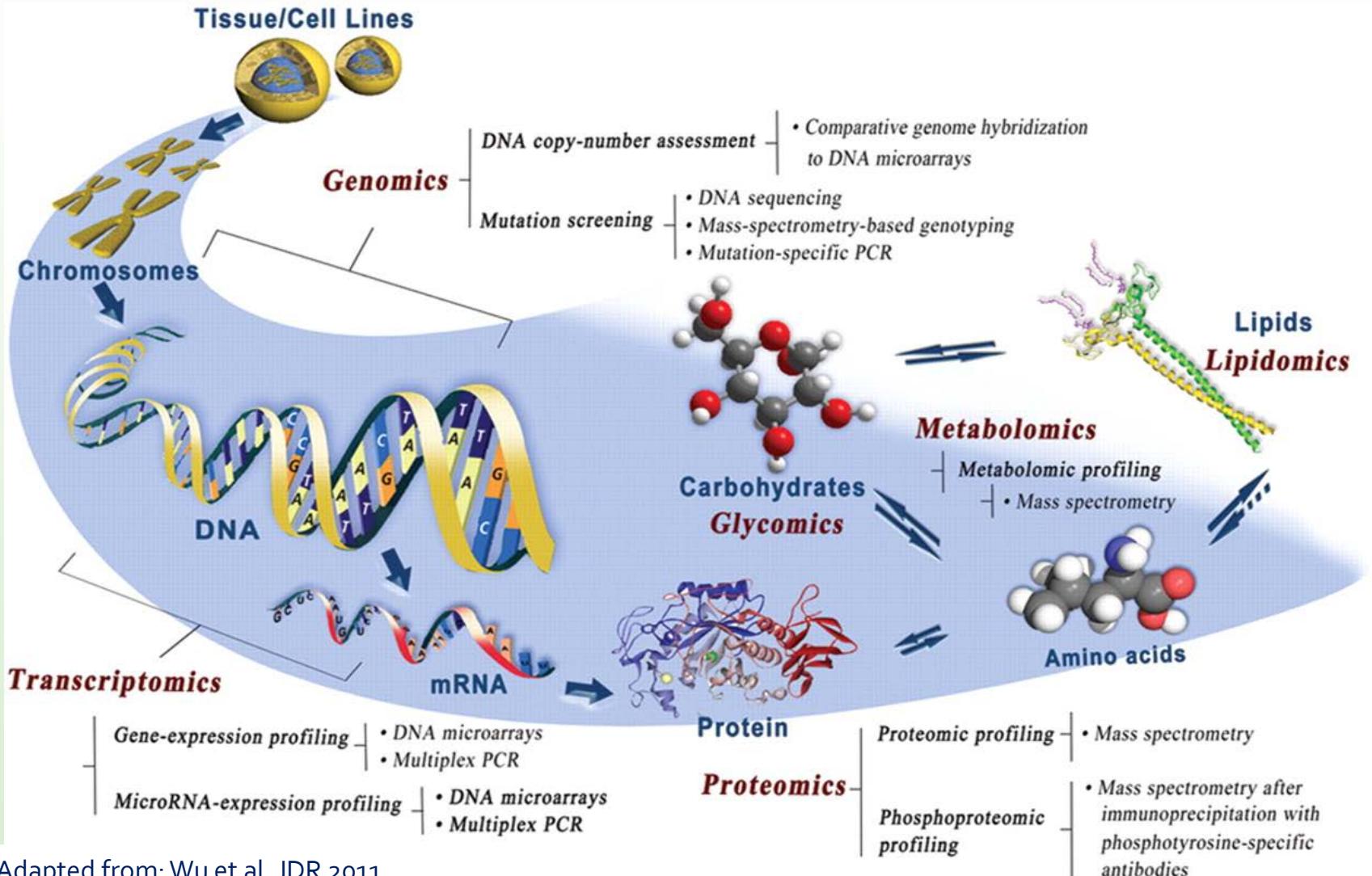
THE DIGITAL CODE OF DNA



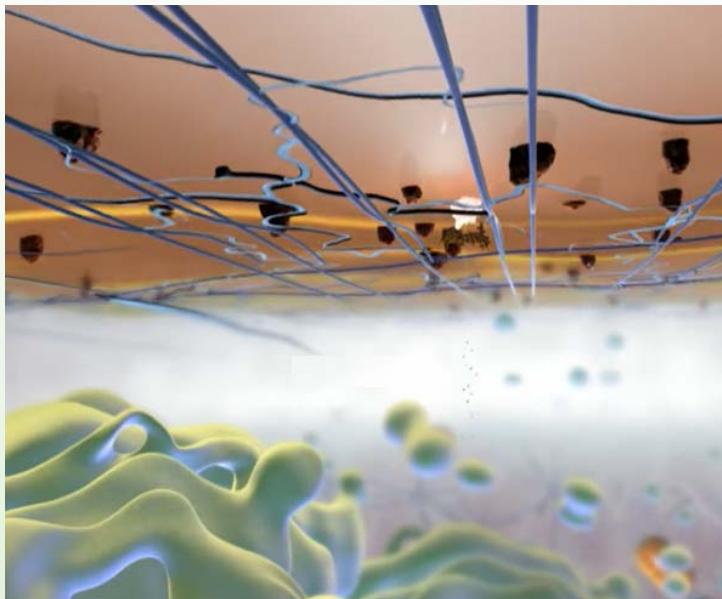
INFORMATION CONTENT: BIOLOGICAL ENTITIES VS STORAGE DEVICES



From genotype to phenotype

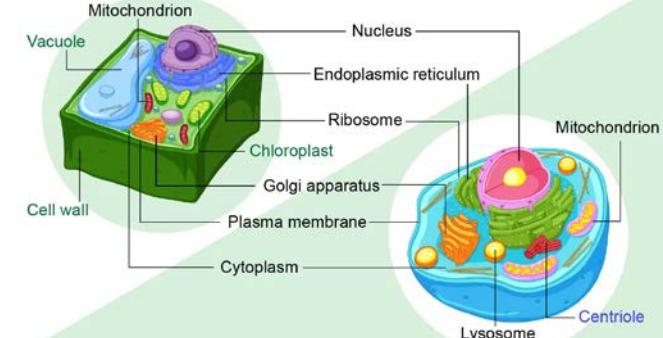


THE CELL



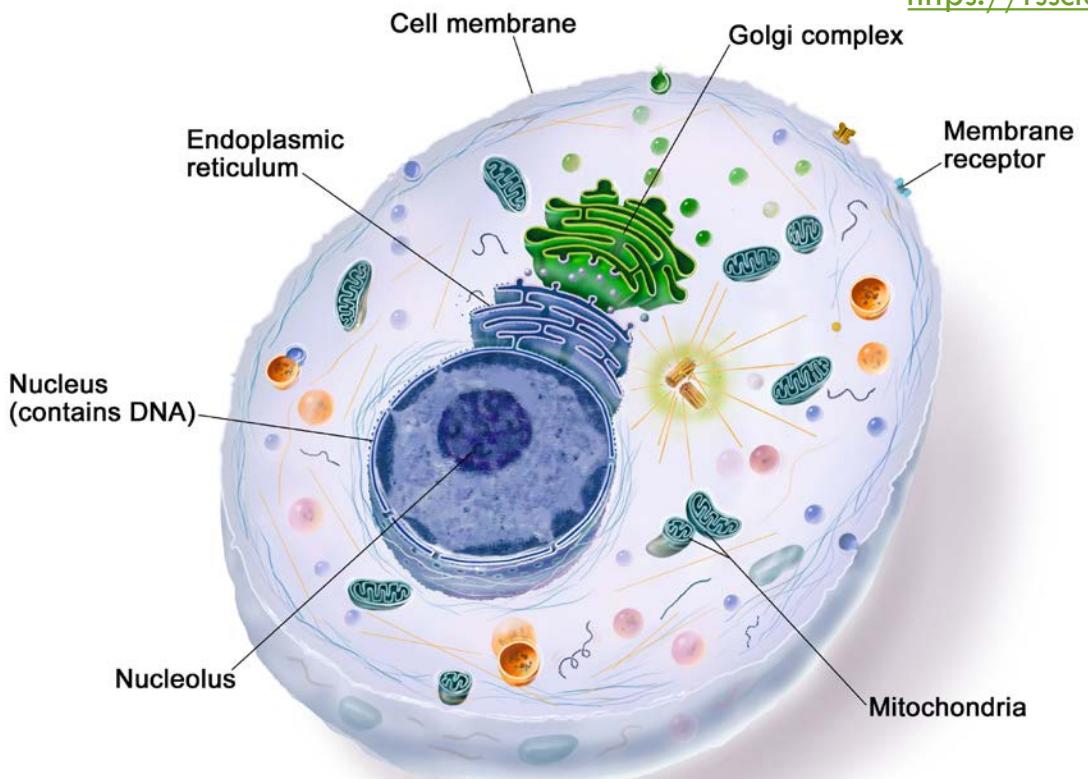
<https://www.youtube.com/watch?v=2KQbVr9kFO0>

PLANT CELL



ANIMAL CELL

<https://rsscience.com/about-me/>



<https://www.cancer.gov/publications/dictionaries/cancer-terms/def/cell>

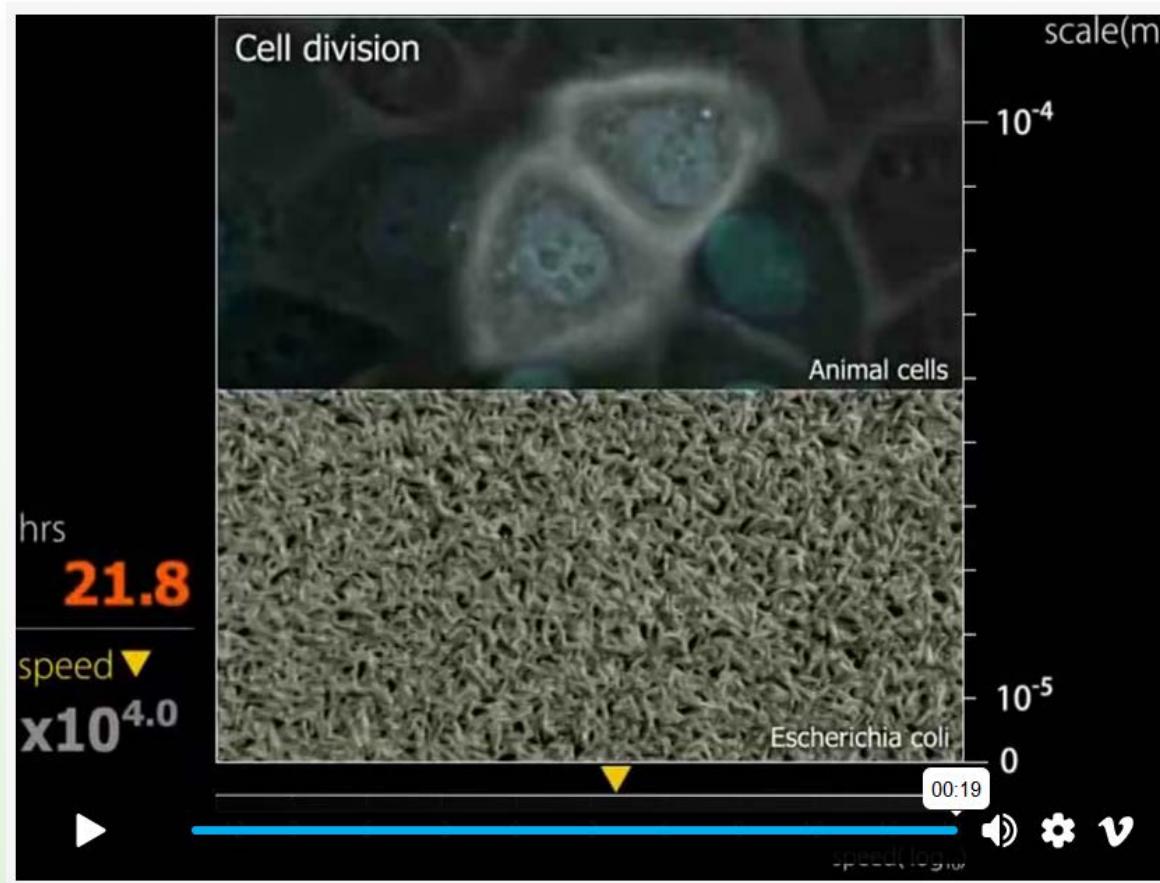
© 2014 Terese Winslow LLC
U.S. Govt. has certain rights

GLOSSARY FOR VIDEO

<https://www.genome.gov/genetics-glossary>

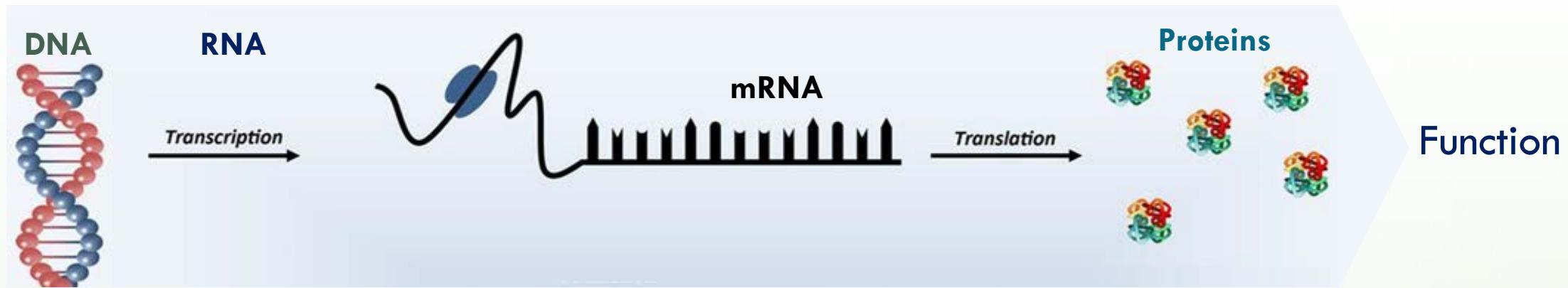
- PLASMA MEMBRANE: MEMBRANE FOUND IN ALL CELLS THAT SEPARATES THE INTERIOR OF THE CELL FROM THE OUTSIDE ENVIRONMENT
- EXTRACELLULAR MATRIX: LARGE NETWORK OF PROTEINS AND OTHER MOLECULES THAT SURROUND, SUPPORT, AND GIVE STRUCTURE TO CELLS AND TISSUES IN THE BODY, COMMUNICATION AND MOVEMENT
- GOLGI: ORGANELLE THAT PROCESSES AND PACKAGES PROTEINS AND MOLECULES, ESPECIALLY THOSE TO BE EXPORTED FROM THE CELL
- LYSOSOME: ORGANELLE CONTAINING ENZYMES CAPABLE OF BREAKING DOWN PROTEINS, NUCLEIC ACIDS, CARBOHYDRATES, LIPIDS
- CENTRIOLES: ORGANELLES LOCATED IN THE CYTOPLASM NEAR THE NUCLEAR ENVELOPE. CENTRIOLES PLAY A ROLE IN ORGANIZING MICROTUBULES THAT SERVE AS THE CELL'S SKELETAL SYSTEM.
- MICROTUBULES: MAJOR COMPONENTS OF THE CYTOSKELETON. INVOLVED IN MITOSIS, CELL MOTILITY, INTRACELLULAR TRANSPORT, AND MAINTENANCE OF CELL SHAPE.

CELL DIVISION

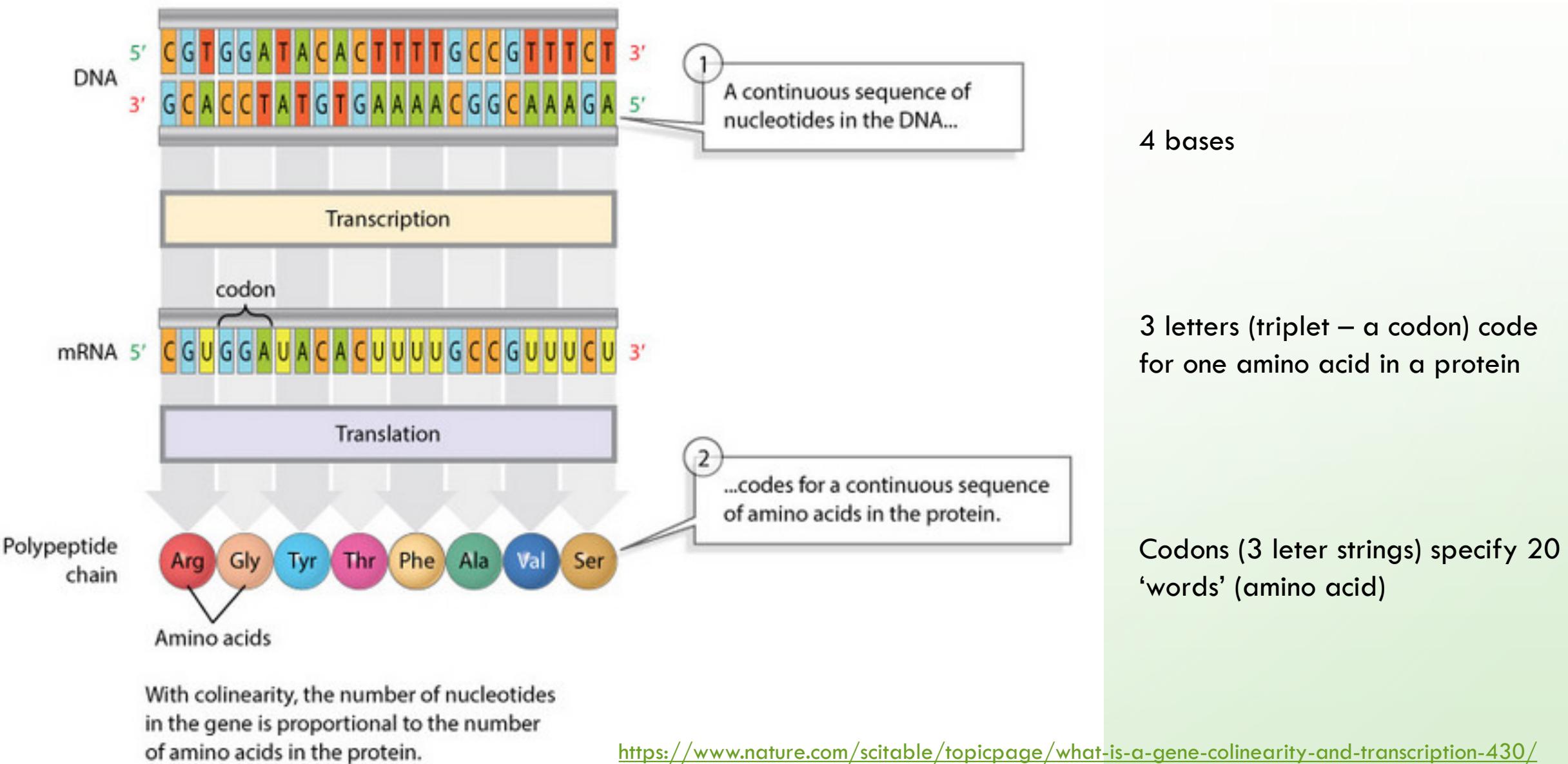


https://vimeo.com/89112986?embedded=true&source=vimeo_logo&owner=9445085

DNA, RNA AND PROTEINS



FROM DNA TO FUNCTION



$4^3 = 64$ possible ways (order important, repetition allowed) to code for 20 amino acids and stop codons

Second nucleotide				Third nucleotide
	U	C	A	
U	UUU Phe UUC UUA Leu UUG	UCU Ser UCC UCA UCG	UAU Tyr UAC UAA STOP UAG STOP	UGU Cys UGC UGA STOP UGG Trp
C	CUU CUC Leu CUA CUG	CCU Pro CCC CCA CCG	CAU His CAC CAA Gln CAG	CGU CGC Arg CGA CGG
A	AUU AUC Ile AUA AUG Met	ACU ACC ACA Thr ACG	AAU Asn AAC AAA Lys AAG	AGU Ser AGC AGA Arg AGG
G	GUU GUC Val GUA GUG	GCU GCC GCA Ala GCG	GAU Asp GAC GAA Glu GAG	GGU GGC Gly GGA GGG

- Alanine
- Arginine
- Asparagine
- Aspartic Acid
- Cysteine
- Glutamic acid
- Glutamine
- Glycine
- Histidine
- Isoleucine
- Leucine
- Lysine
- Methionine
- Phenylalanine
- Proline
- Serine
- Threonine
- Tryptophan
- Tyrosine
- Valine

A SINGLE BASE CHANGE CAN CREATE DEVASTATING GENETIC DISORDERS

e.g. Sickle-Cell Anemia

In sickle-cell anemia, the gene for the beta chain of the hemoglobin protein (the oxygen-carrying protein that makes blood red) is mutated.

Beta hemoglobin (beta globin) is a single chain of 147 amino acids.

One single-base mutation causes the sixth amino acid in the chain to be valine, rather than glutamic acid. This changes the red blood cell shape and ability to carry oxygen, with huge consequences for the individual.

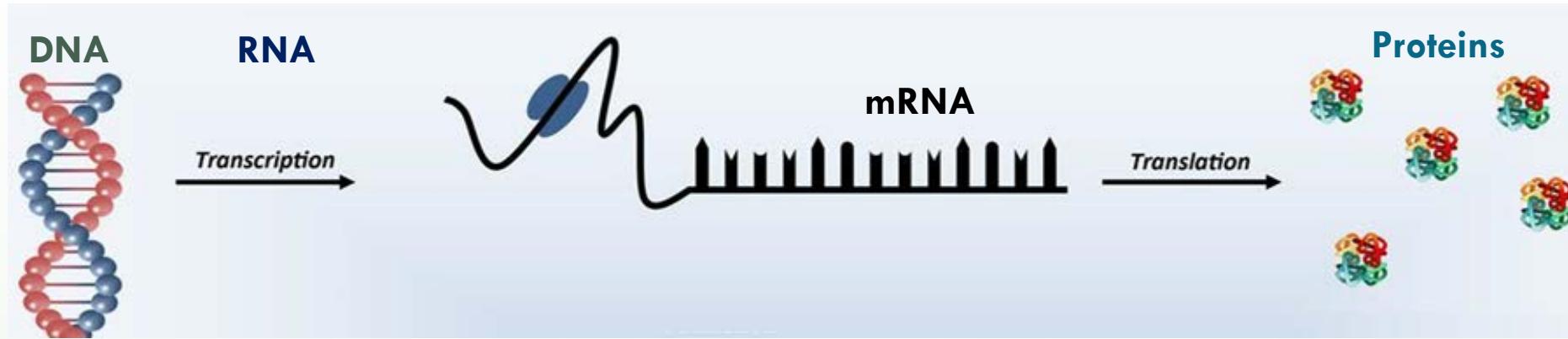
Wild-Type

ATG	GTG	CAC	CTG	ACT	CCT	GAG	GAG	AAG	TCT	GCC	GTT	ACT
Start	Val	His	Leu	Thr	Pro	Glu	Glu	Lys	Ser	Ala	Val	Thr

Mutant

ATG	GTG	CAC	CTG	ACT	CCT	GTG	GAG	AAG	TCT	GCC	GTT	ACT
Start	Val	His	Leu	Thr	Pro	Val	Glu	Lys	Ser	Ala	Val	Thr

DNA, RNA AND PROTEINS

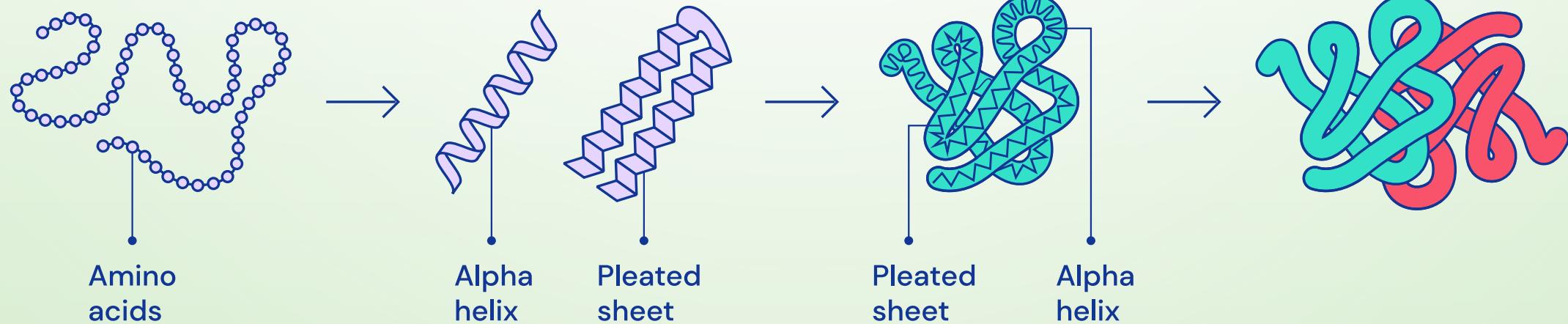


Every protein is made up of a sequence of amino acids bonded together

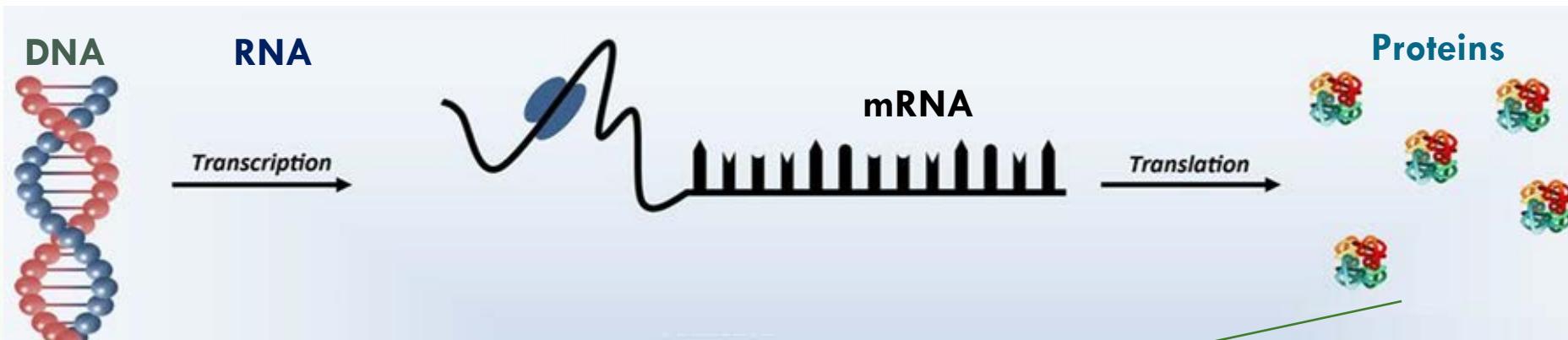
These amino acids interact locally to form shapes like helices and sheets

These shapes fold up on larger scales to form the full three-dimensional protein structure

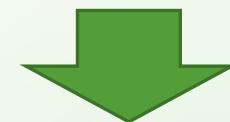
Proteins can interact with other proteins, performing functions such as signalling and transcribing DNA



DNA, RNA AND PROTEINS



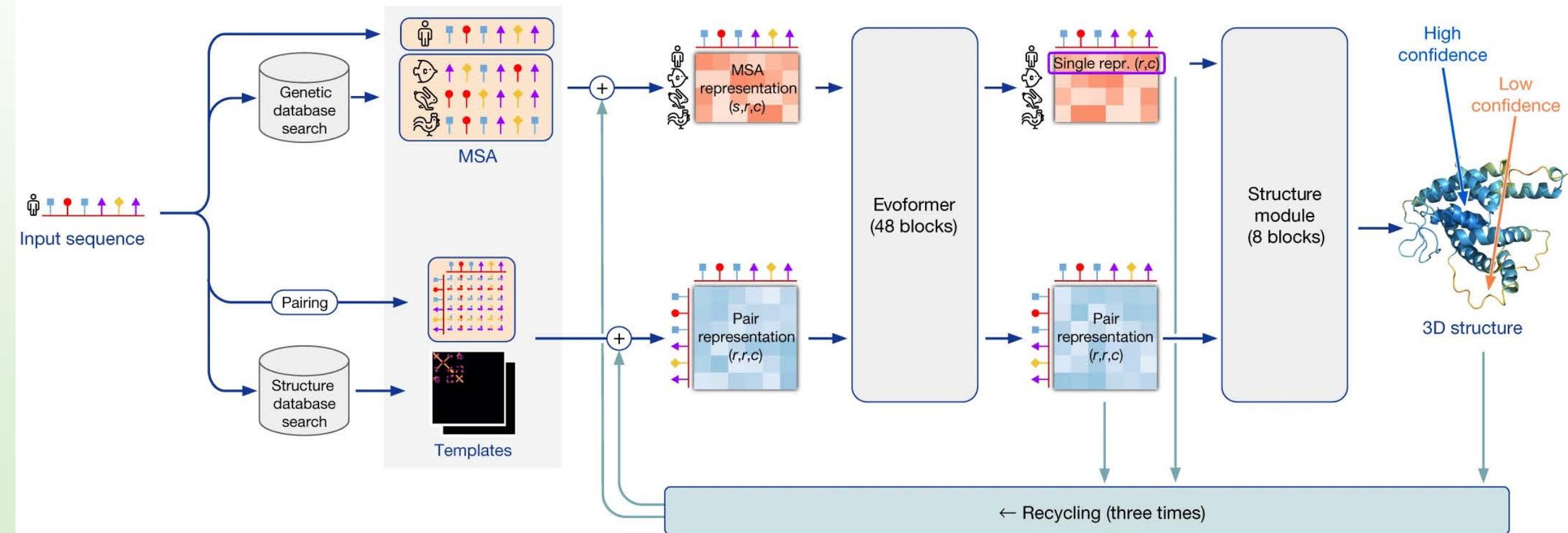
Protein 3D structure



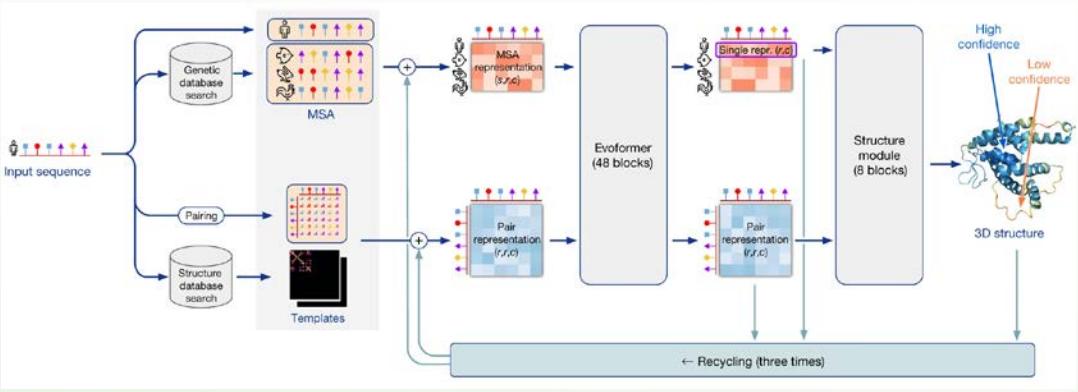
Protein function

USING AI TO PREDICT PROTEIN STRUCTURE

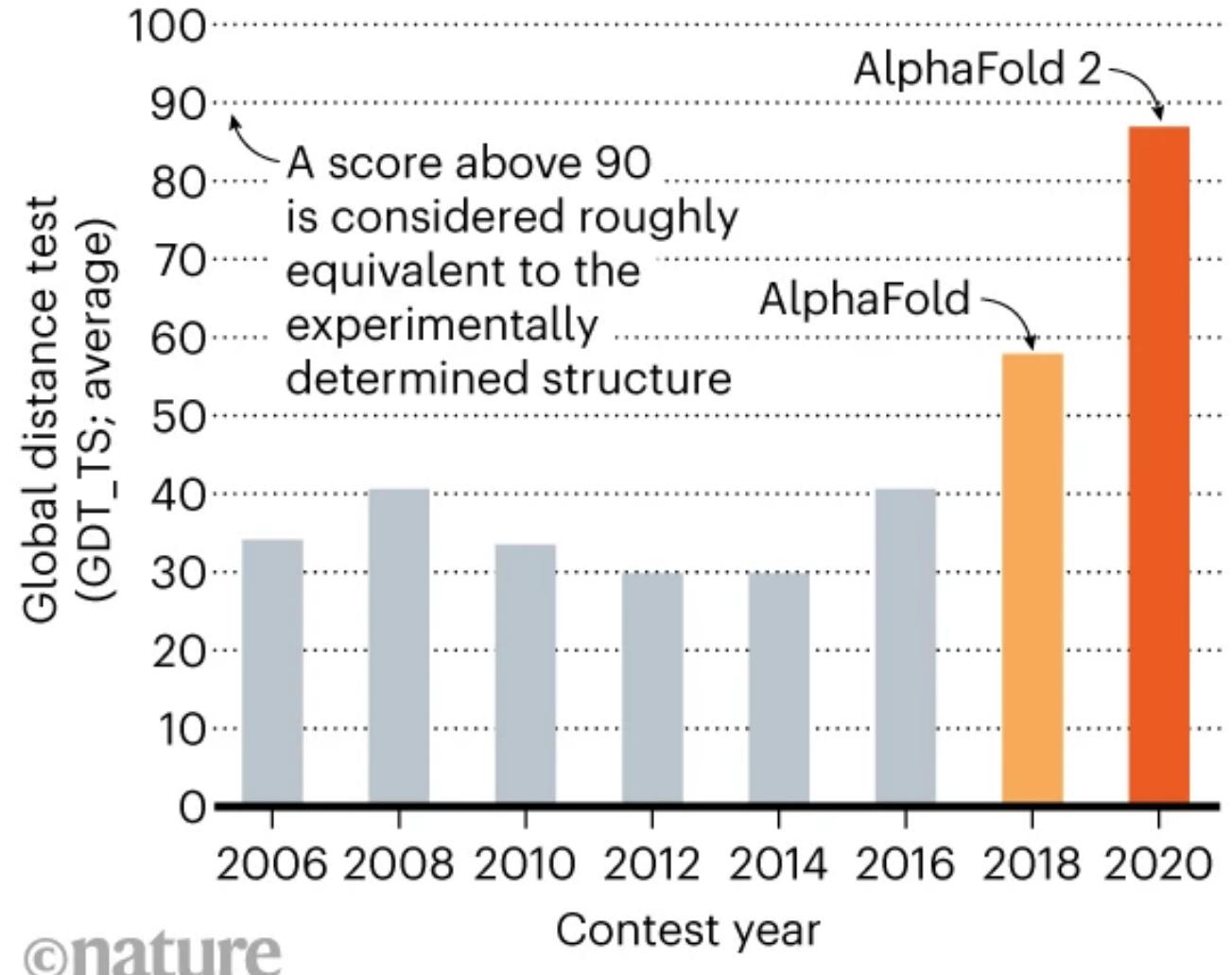
Can we determine a protein's 3D shape from its amino-acid sequence?



ALPHAFOLD: USING AI TO PREDICT PROTEIN STRUCTURE



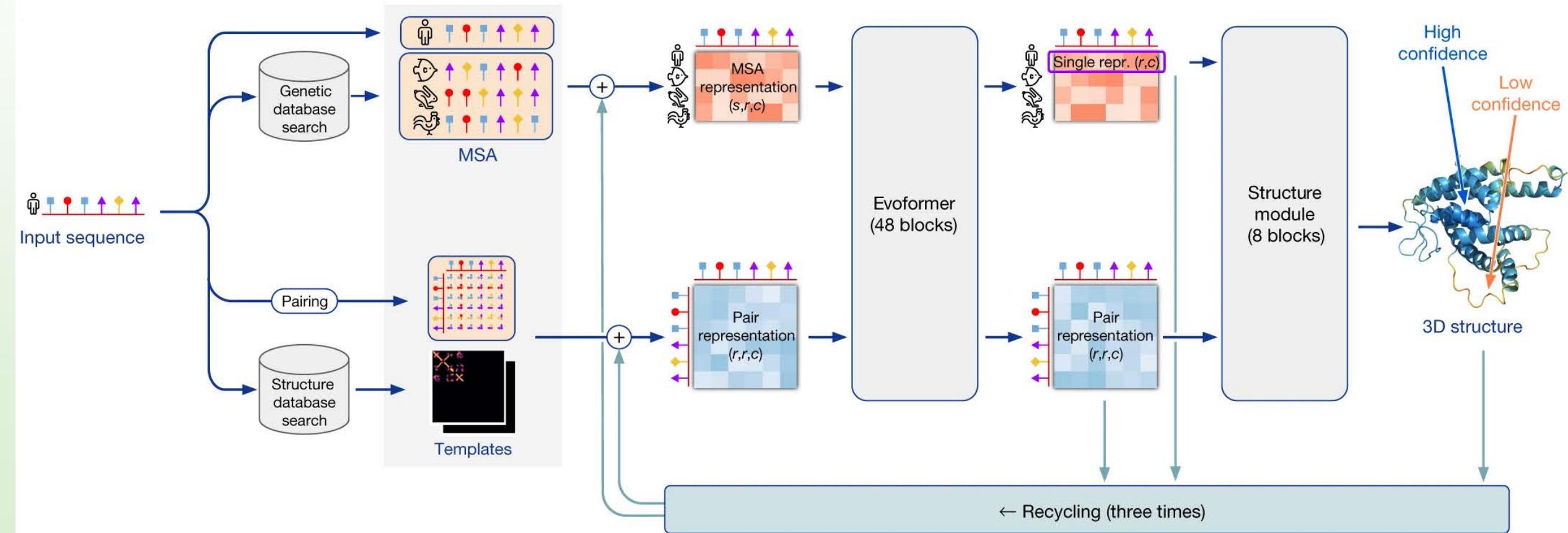
<https://alphafold.ebi.ac.uk/>



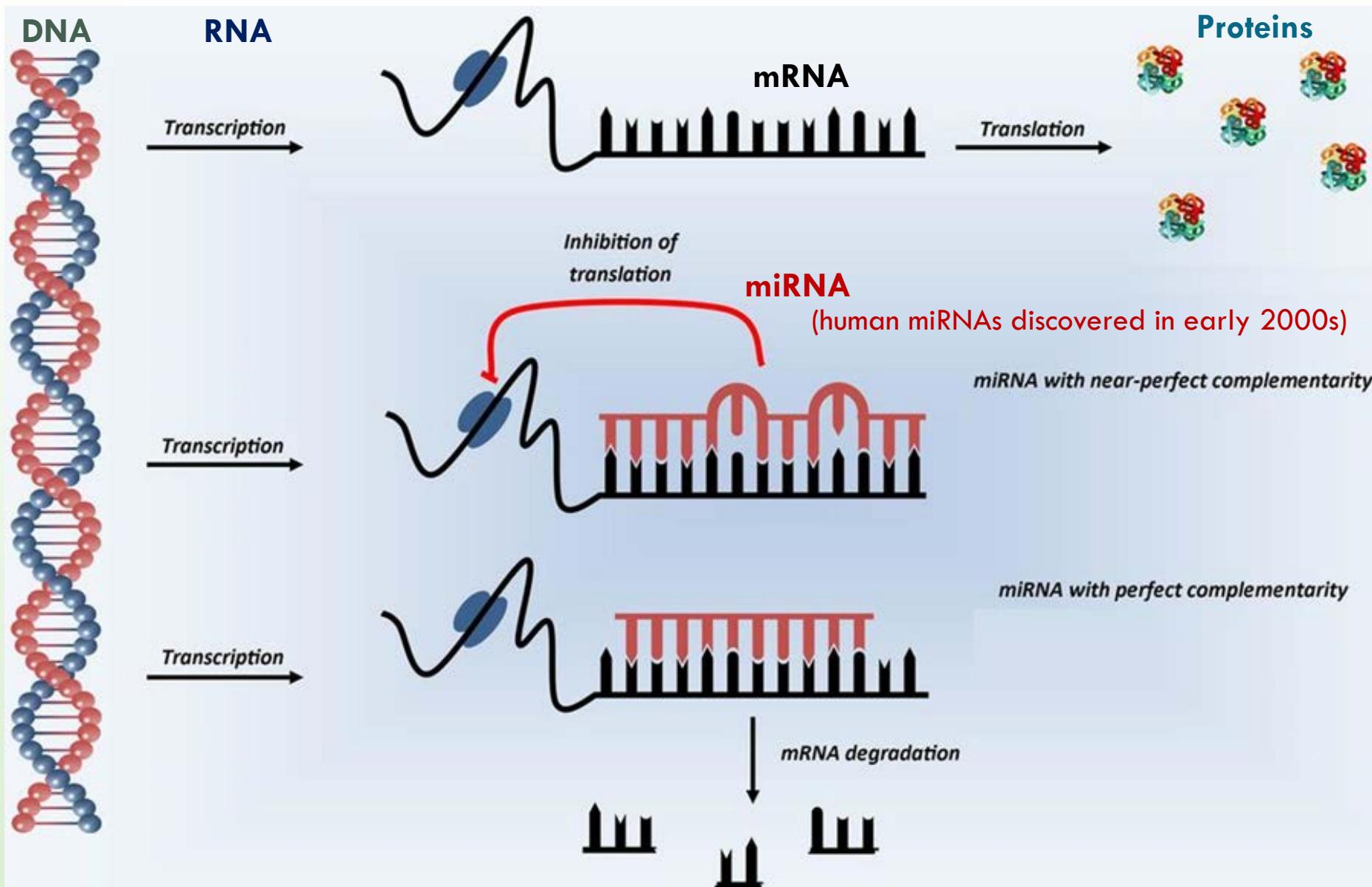
©nature

USING AI TO PREDICT PROTEIN STRUCTURE

Can we determine a protein's 3D shape from its amino-acid sequence?

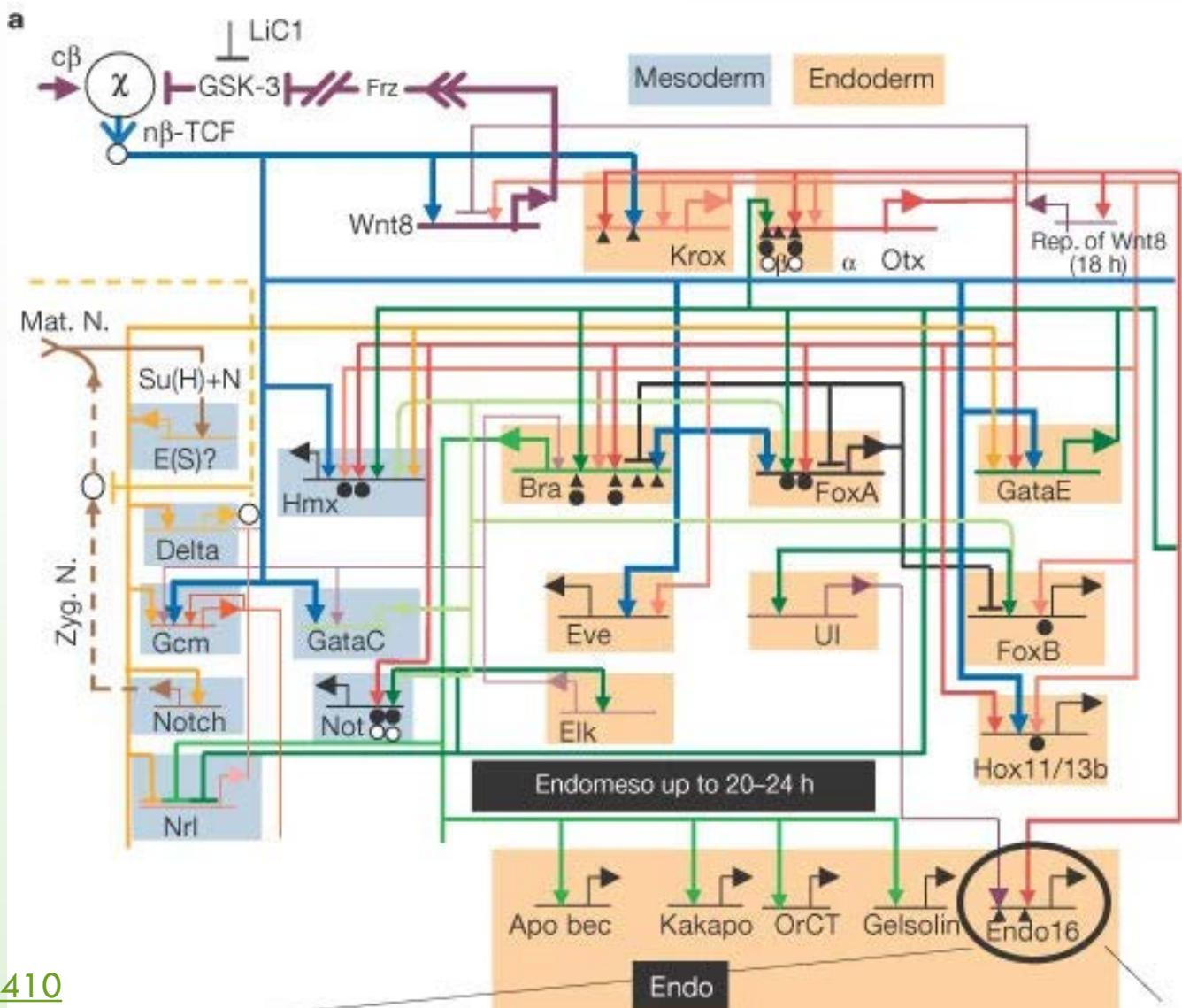


DNA, RNA AND PROTEINS



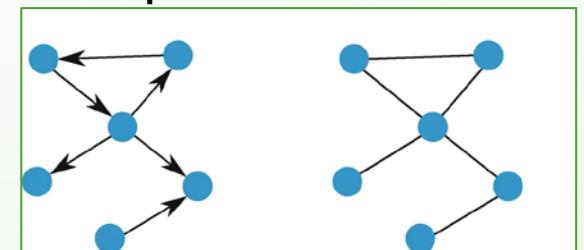
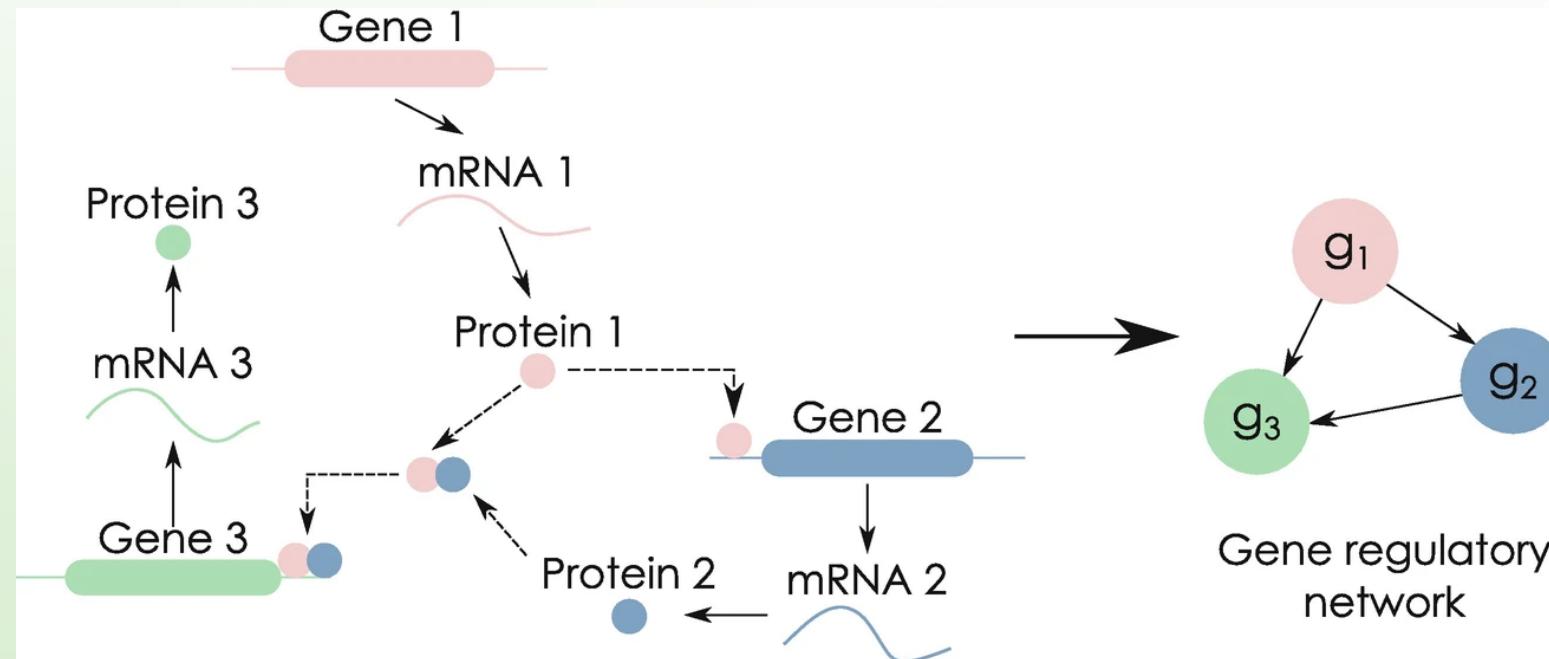
CODING NETWORK OF GENES

Gene regulatory networks that specify the behaviour of the genes can be represented computationally using methods from STEM



REGULATION OF GENE EXPRESSION IN EUKARYOTES

At any given time, a complex set of interactions between genes, RNA molecules and proteins determine which genes are activated, and the amount of protein or RNA product.

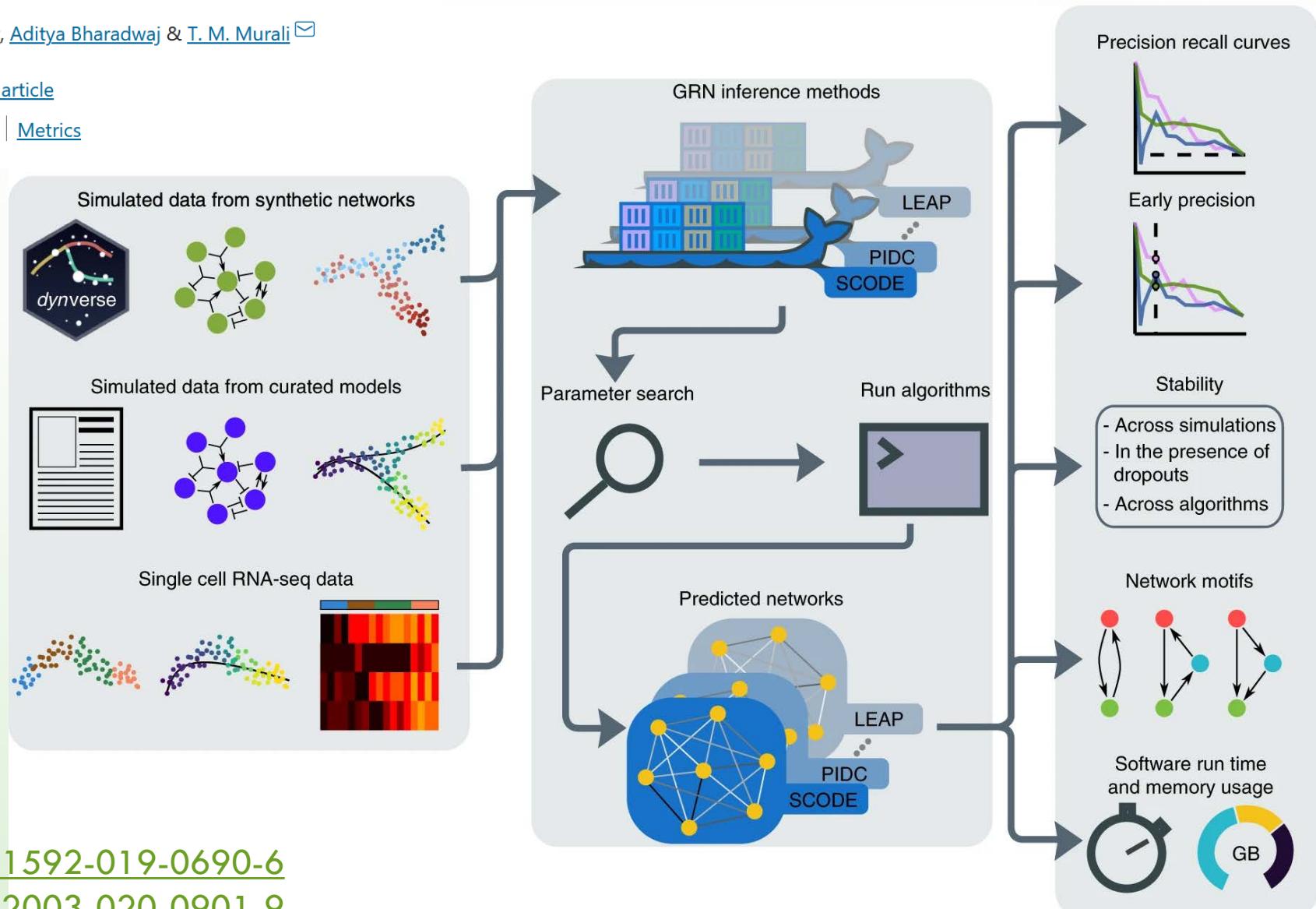


Benchmarking algorithms for gene regulatory network inference from single-cell transcriptomic data

Aditya Pratapa, Amogh P. Jalihal, Jeffrey N. Law, Aditya Bharadwaj & T. M. Murali 

[Nature Methods](#) 17, 147–154 (2020) | [Cite this article](#)

34k Accesses | 168 Citations | 64 Altmetric | [Metrics](#)



<https://www.nature.com/articles/s41592-019-0690-6>

<https://www.nature.com/articles/s42003-020-0901-9>