

**LOSCHMIDT
LABORATORIES**

Artificial Intelligence in Microbiology

by Stanislav Mazurenko, PhD

mazurenko@mail.muni.cz

- ❑ **Motivation**
- ❑ **Introduction to AI and ML**
- ❑ **Recent applications in Microbiology**



Motivation

Motivation: sequences and chemicals

- ❑ Large volumes of digital data
- ❑ Affordable computing power and storage
- ❑ Complex study objects

GenBank ^{FREE}

Eric W Sayers ✉, Mark Cavanaugh, Karen Clark, James Ostell,
Kim D Pruitt, Ilene Karsch-Mizrachi



GOLD

biobank ^{uk}
Improving the health of future generations

RCSB PDB
PROTEIN DATA BANK

UniProt

BRENDA

FIREPROT ^{DB}

PubChem

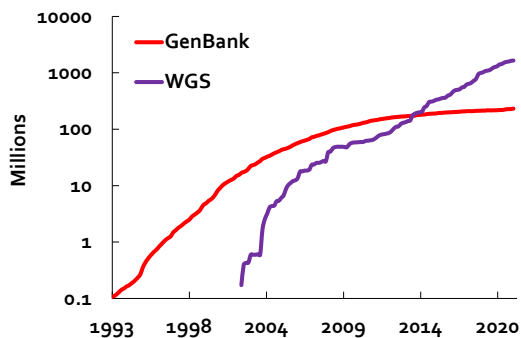
ChEMBL

JOURNAL OF
CHEMICAL
AND MODELING

ZINC 15 – Ligand Discovery for Everyone

Teague Sterling and John J. Irwin*

The total number of sequences



SHARE



UK Biobank Principal Investigator Rory Collins stands amid stored biospecimens from the project's half-million participants. NIGEL HILLER

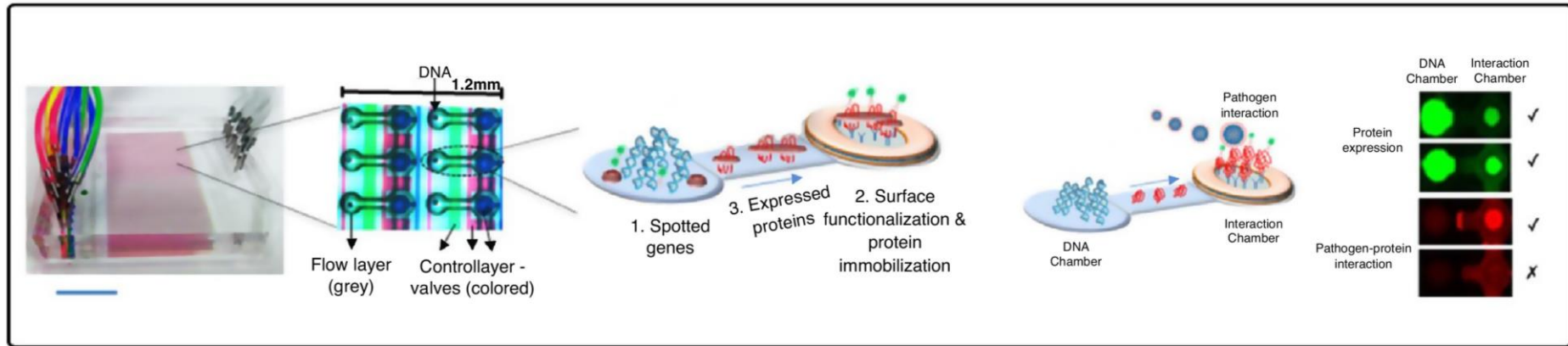
Huge trove of British biodata is unlocking secrets of depression, sexual orientation, and more

By Jocelyn Kaiser, Ann Gibbons | Jan. 3, 2019, 1:20 PM

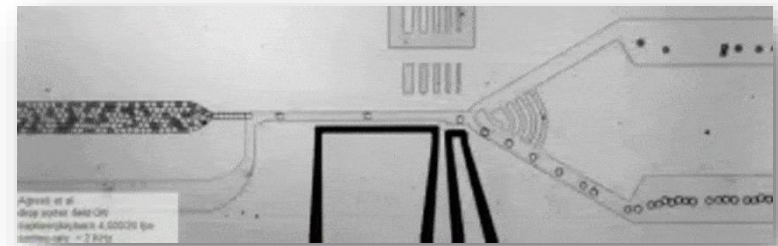
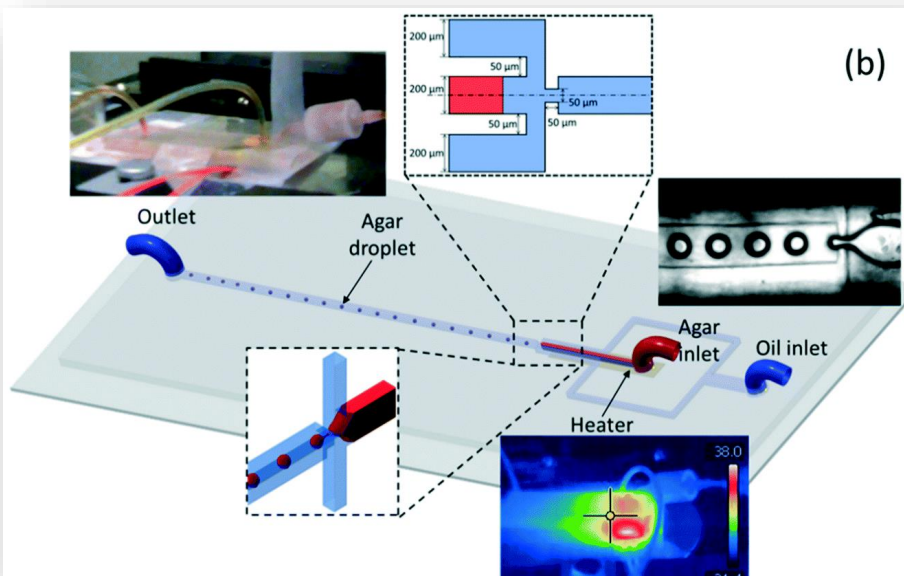
14 Recommended Categories | 253 Origin Categories | 12 Core Categories

Category ID	Description	Items
1014	Brain MRI	922
1005	Cognitive function summary	5
1004	Diet and alcohol summary	321
1002	Early life	13
1007	Education and employment	16
1017	Genomics	30
100113	Geographical and location	13
1015	Heart MRI	39
1019	Linked health outcomes	75
1016	Main abdominal MRI fields likely to be of interest to researchers.	22
1018	Mental health	186
1006	Physical measure summary	66
1001	Primary demographics	8
1003	Self-reported medical conditions	118

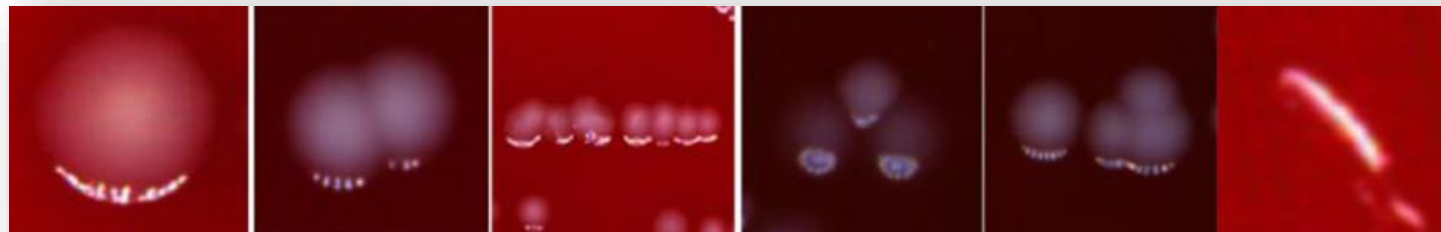
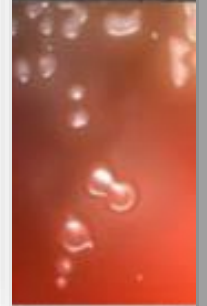
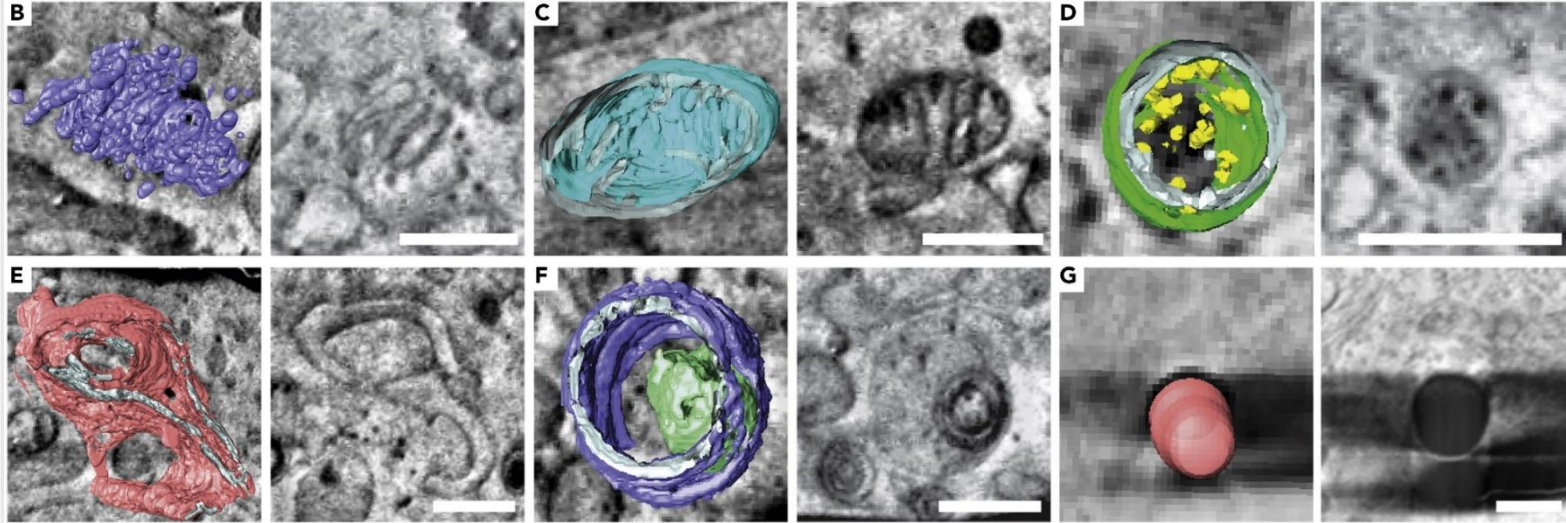
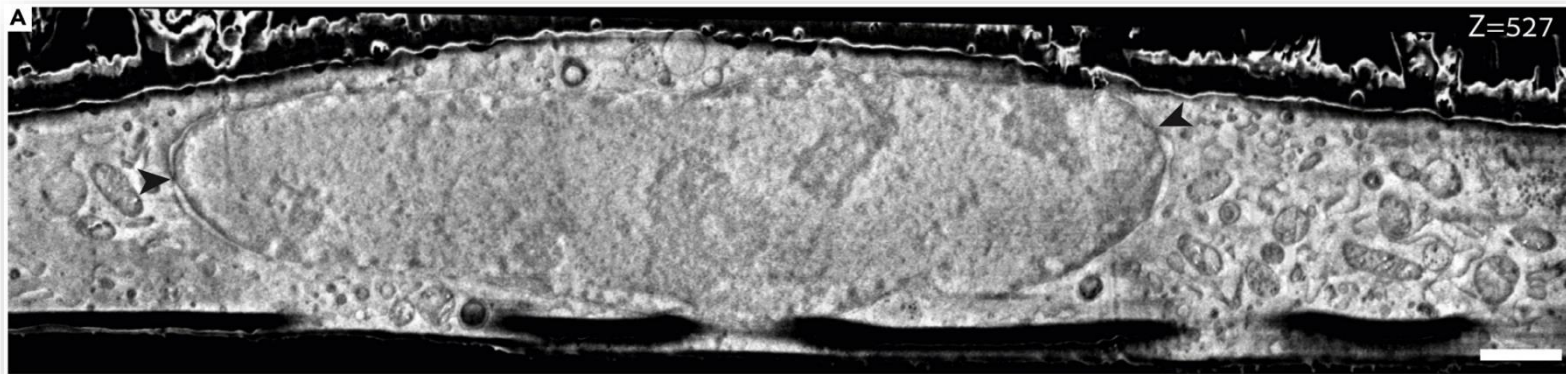
Motivation: big experimental data



Current Opinion in Biotechnology



Motivation: cell imaging

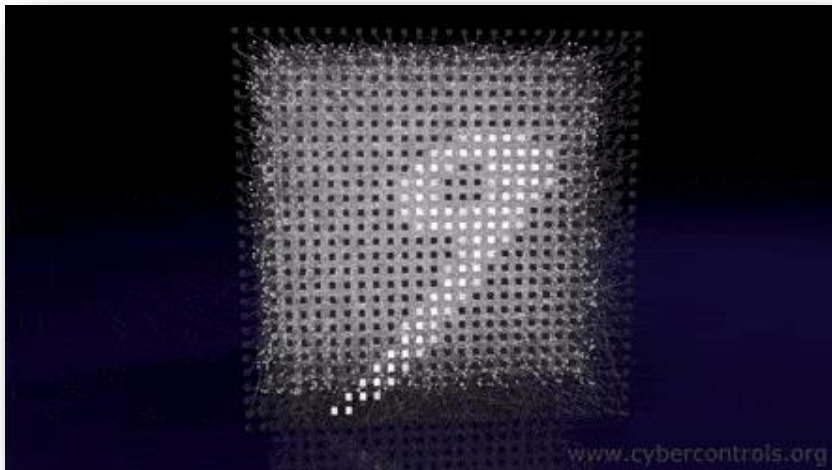
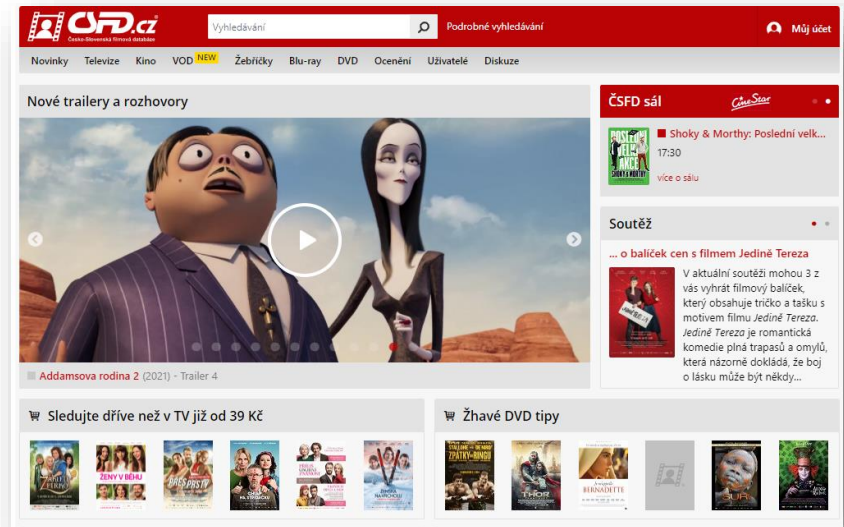


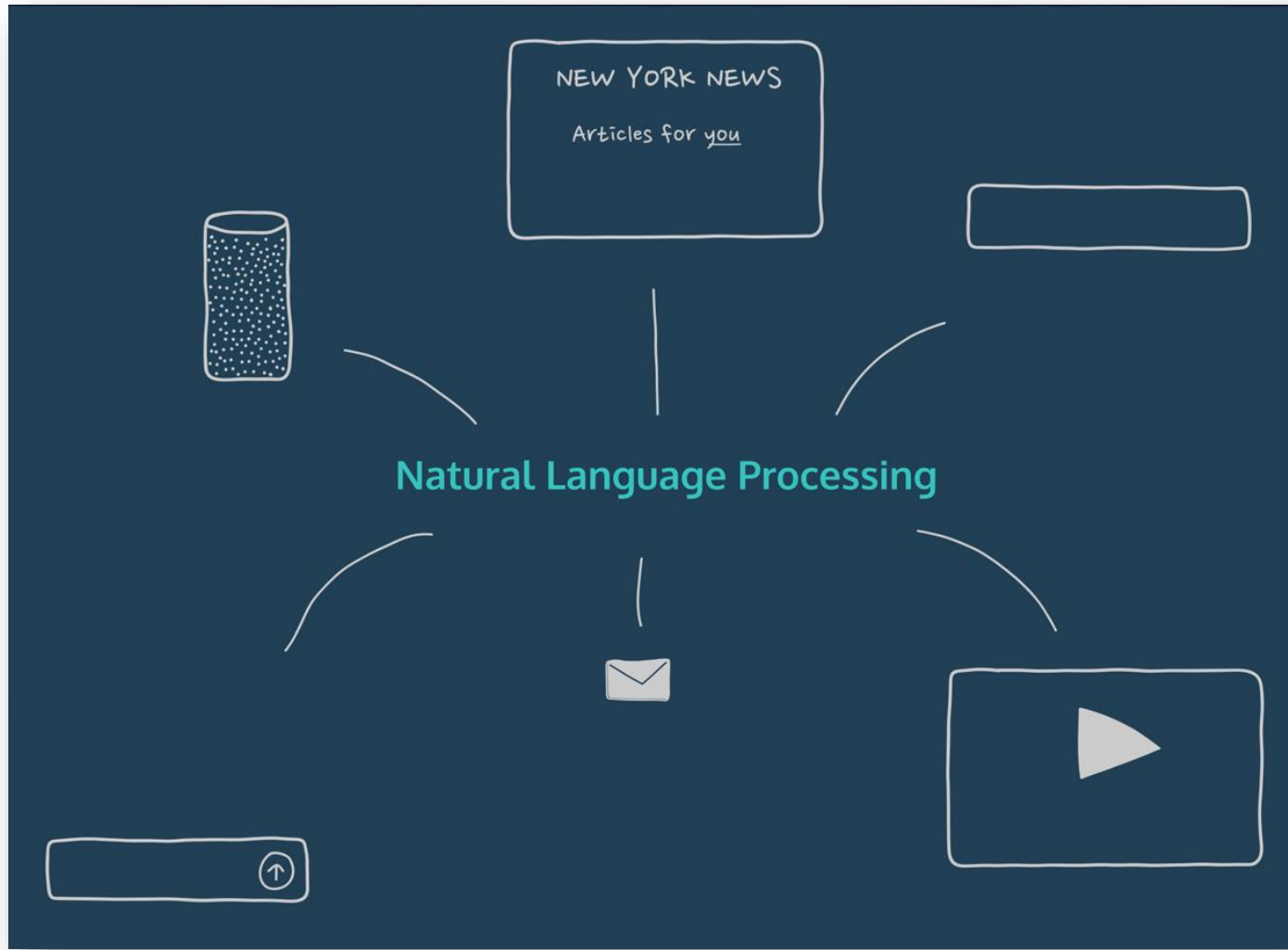


Introduction to AI and ML

Introduction to AI and ML

- Recommendation engines
- Image & speech recognition
- Anomaly detection
- Natural language processing
- Data mining...





Introduction to ML

Faces



Not faces



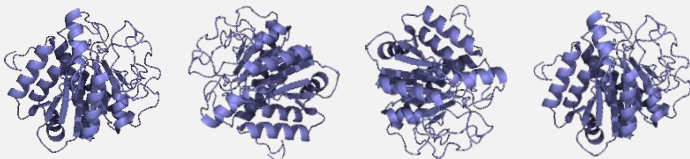
- Historically, people tried to find rules themselves, e.g. detection of particular shapes or color contrasts;
- Often such manual rules are too simplistic to give good results;
- Machine Learning gives the means to generate those rules automatically!



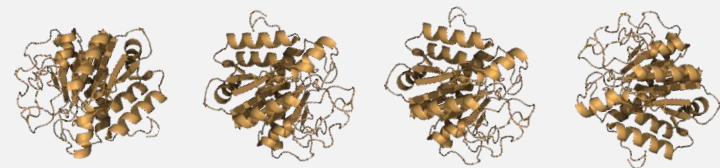
$$F(\text{Image of a woman's face}) = +1$$

$$F(\text{Image of mushrooms}) = -1$$

Non-vaccine candidates:



Vaccine candidates:



Basics of ML: data representation

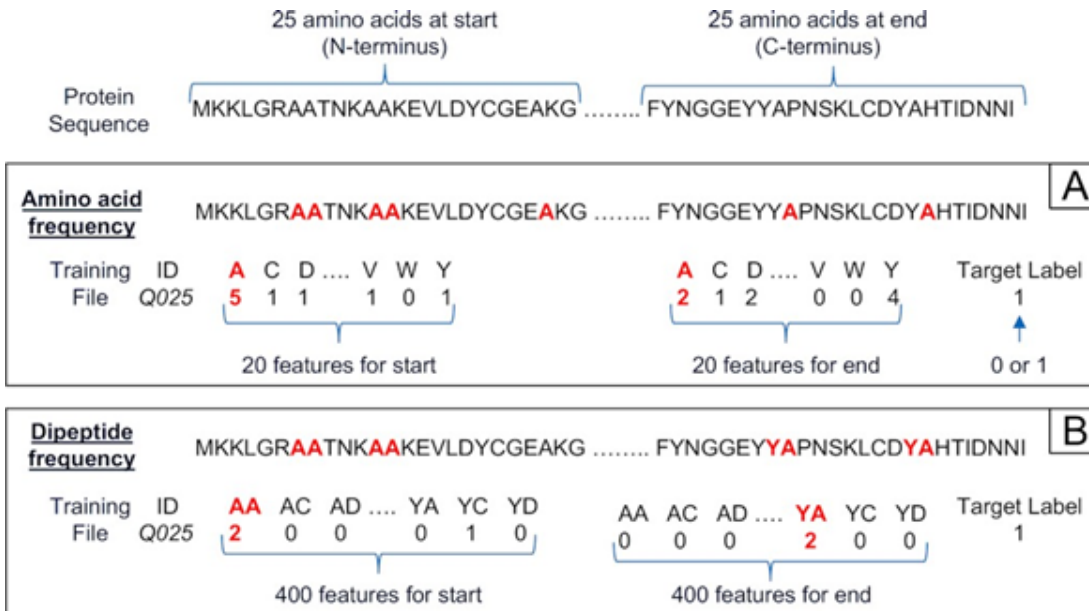
MKKLGRAATNKAAKEVLDYCGEAKG...



Feature vector: (5, 1, 1, -5.67, 0.69, ...)

Examples:

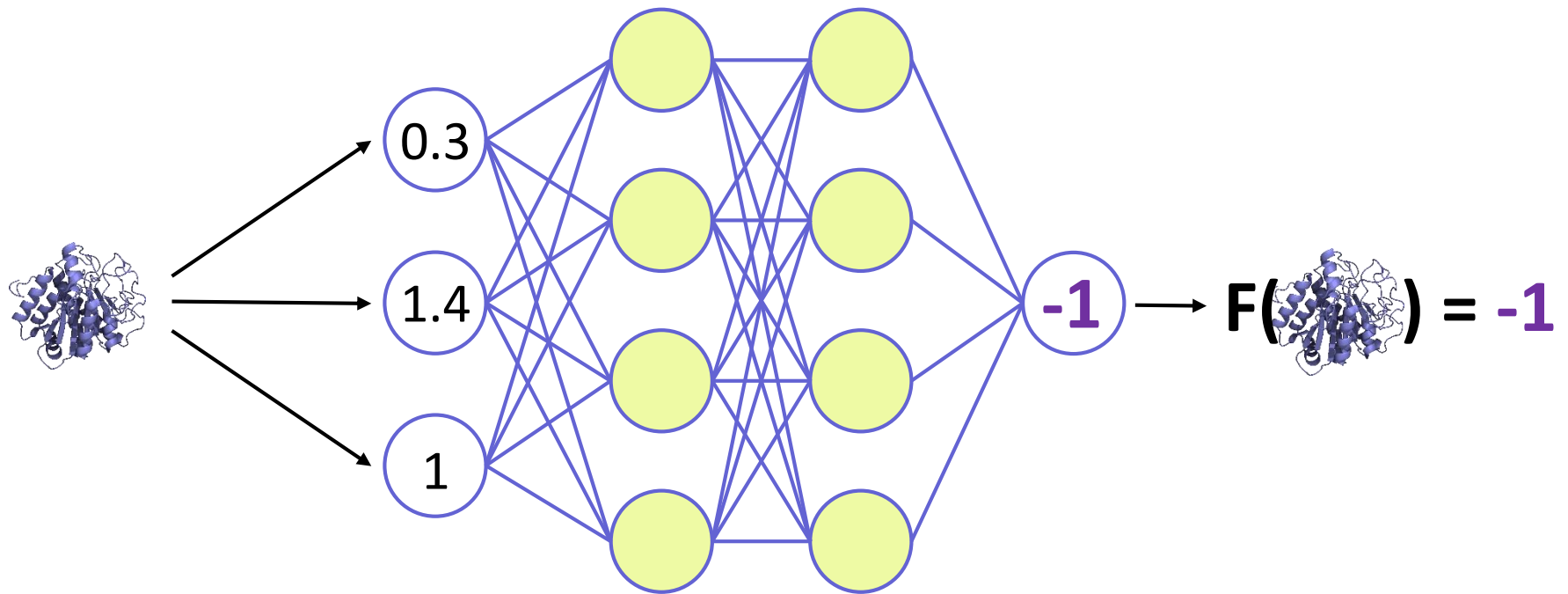
- AA frequency
- AA sequence
- Conservation scores
- Structural elements
- ...



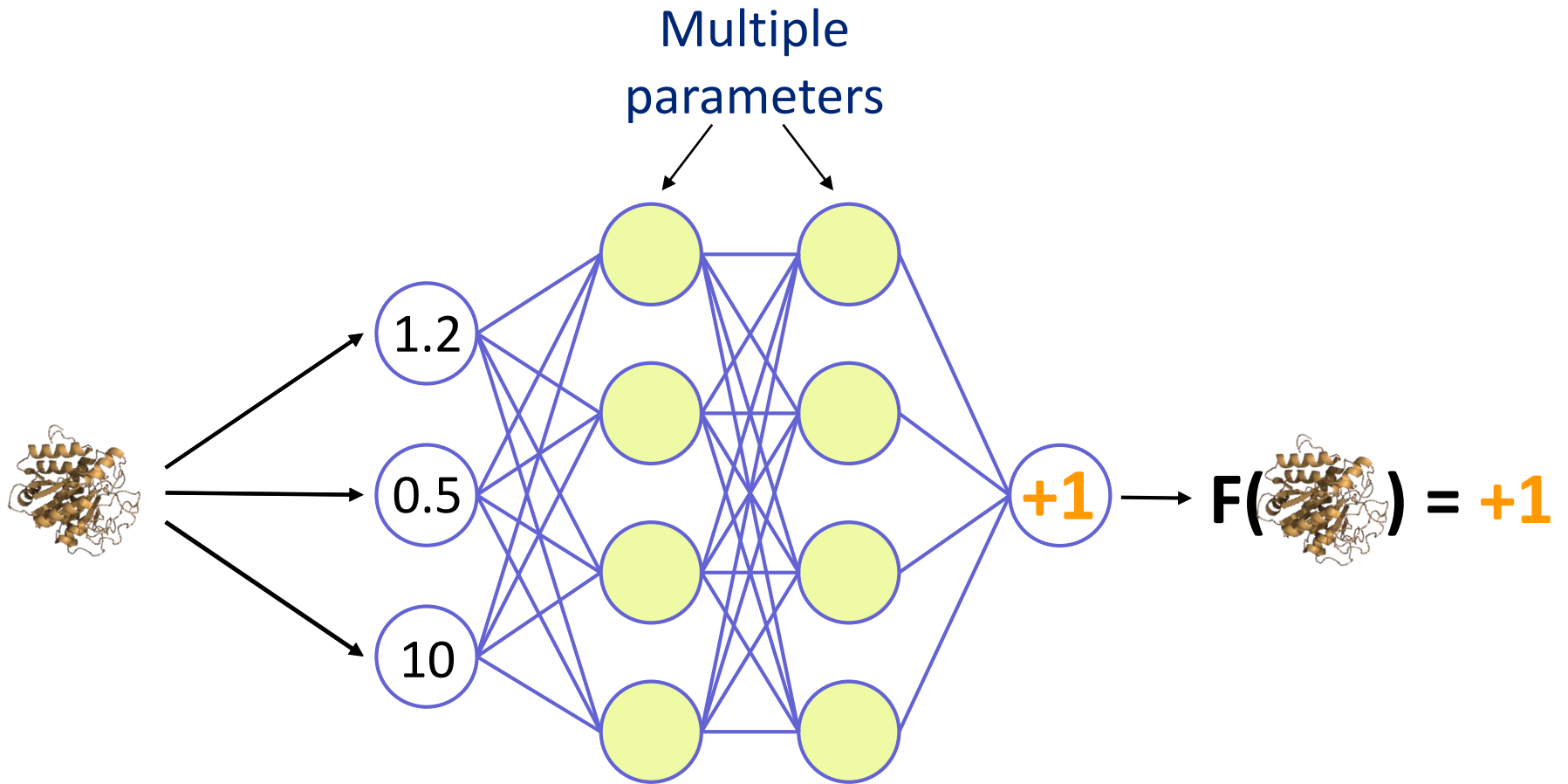
One-hot encoding:

	K	K	L	G	R	A	A	T	...
A	0	0	0	0	0	1	1	0	...
K	1	1	0	0	0	0	0	0	...
L	0	0	1	0	0	0	0	0	...
G	0	0	0	1	0	0	0	0	...
R	0	0	0	0	1	0	0	0	...
T	0	0	0	0	0	0	0	1	...
...									

Basics of ML: training

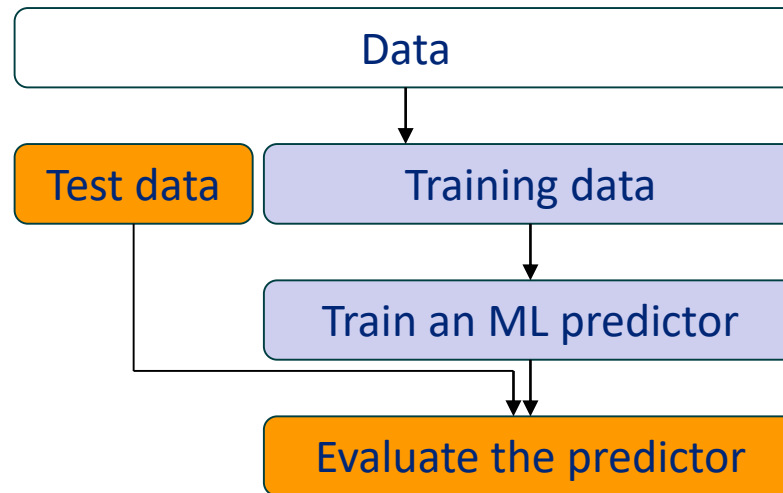


Basics of ML: training

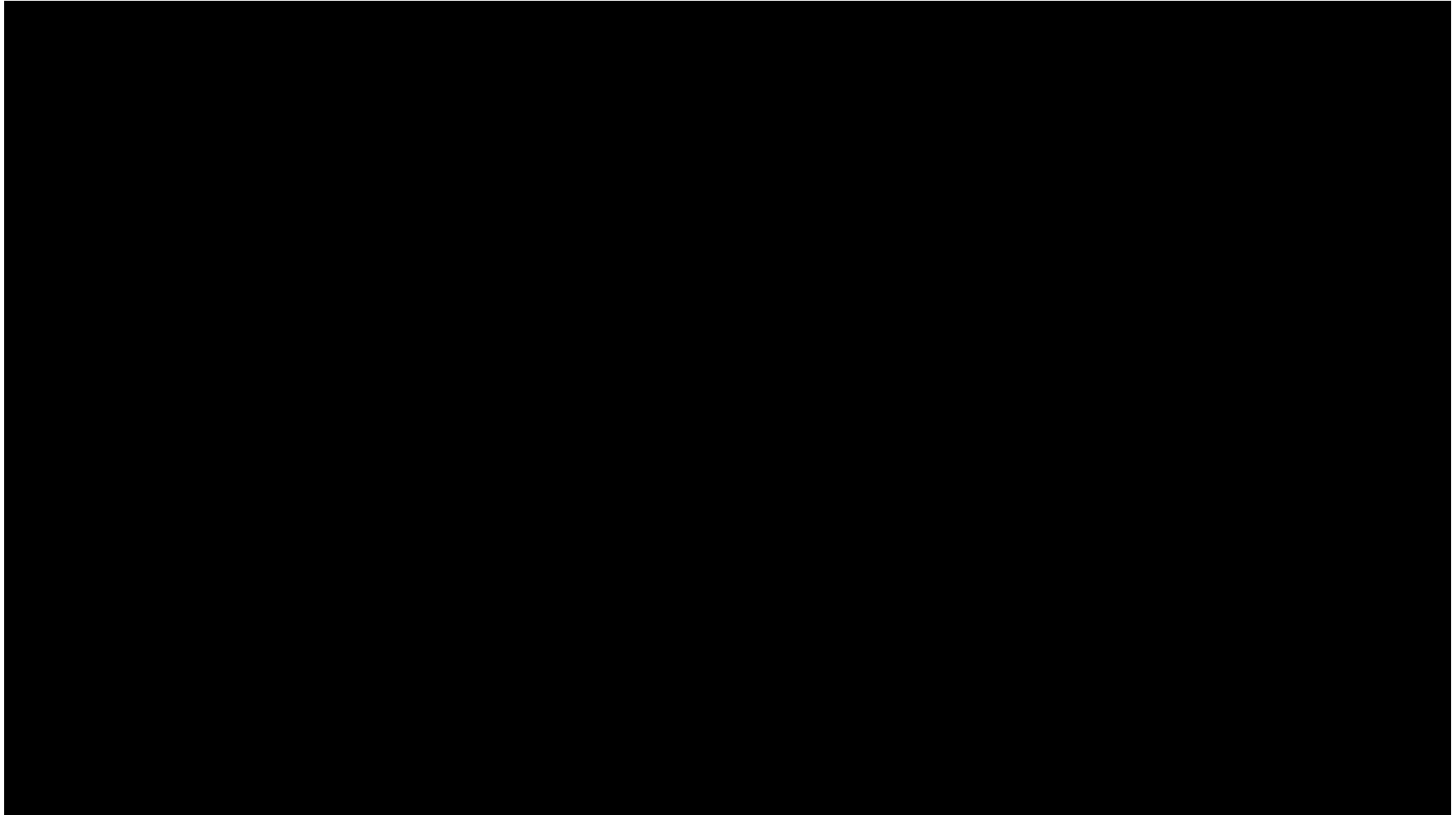


Basics of ML: validation

- The goal of ML is to identify **generalizable patterns** in your training data.
- These patterns must be **valid for future data!**
- Therefore, the core of ML protocol is **to evaluate the predictor on the test data**, hidden from the predictor:



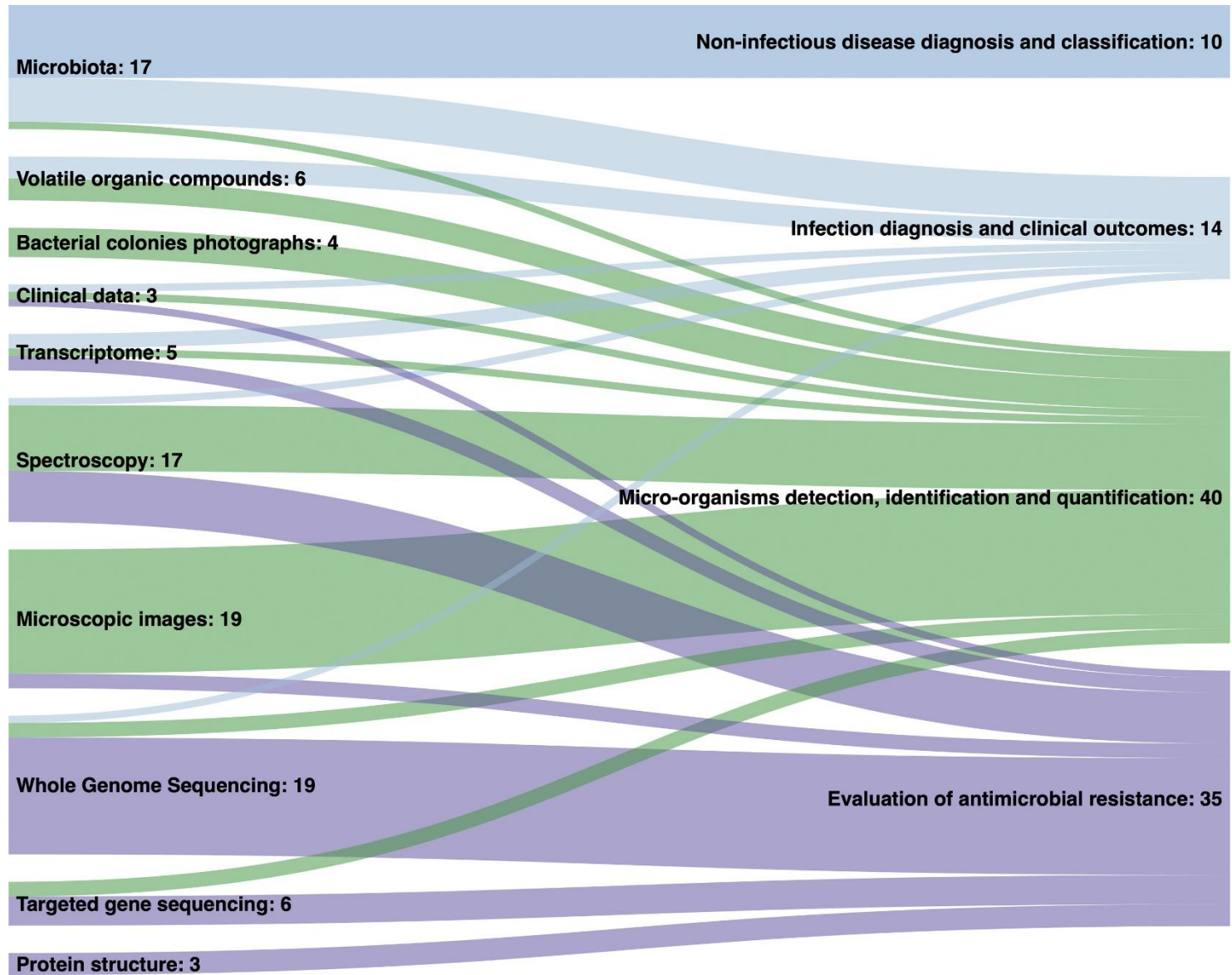
Artificial Neural Networks





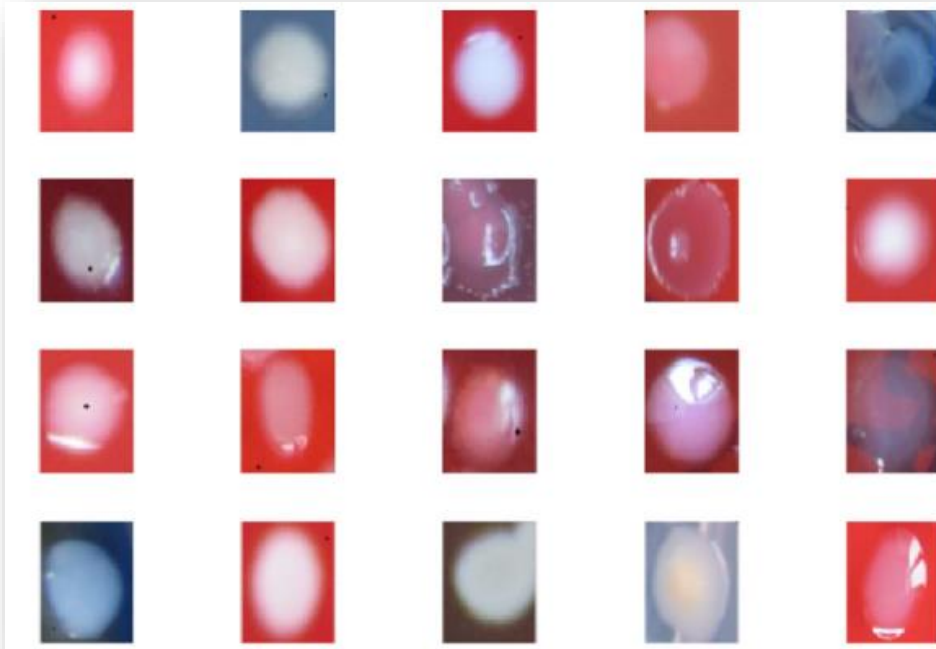
Recent applications

Overview

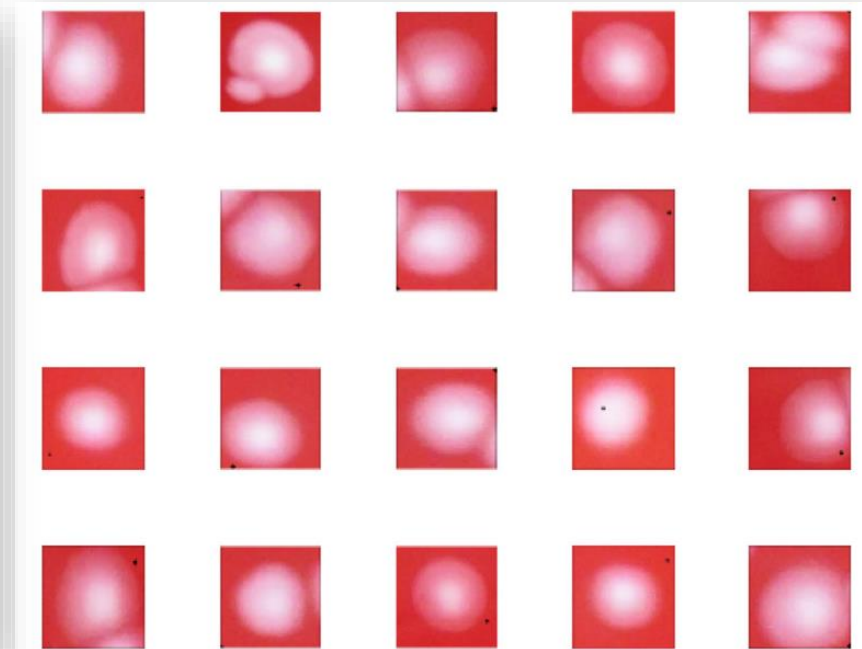


Bacterial colony morphology

Interclass variations

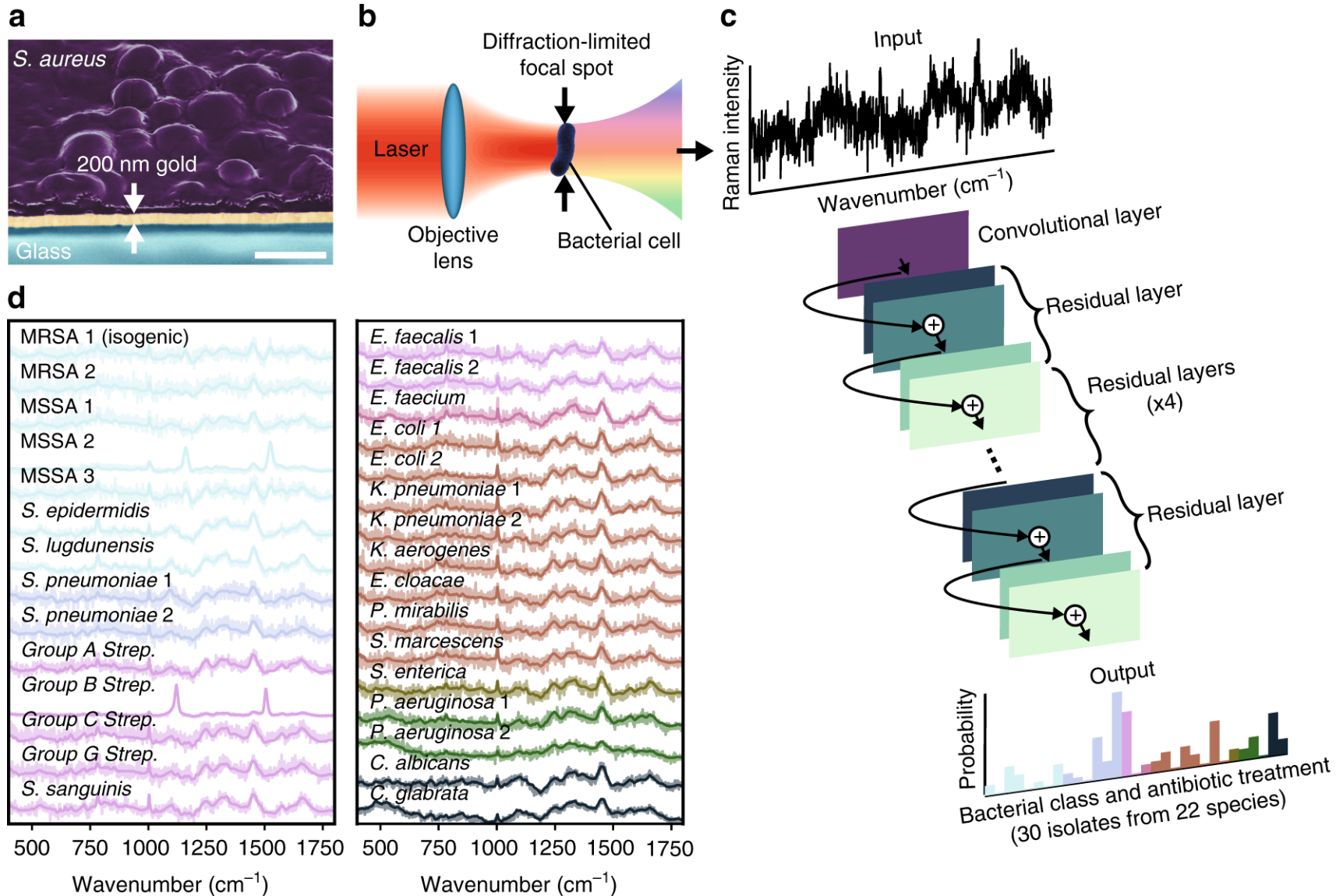


Intraclass variations (*Streptococcus agalactiae*)

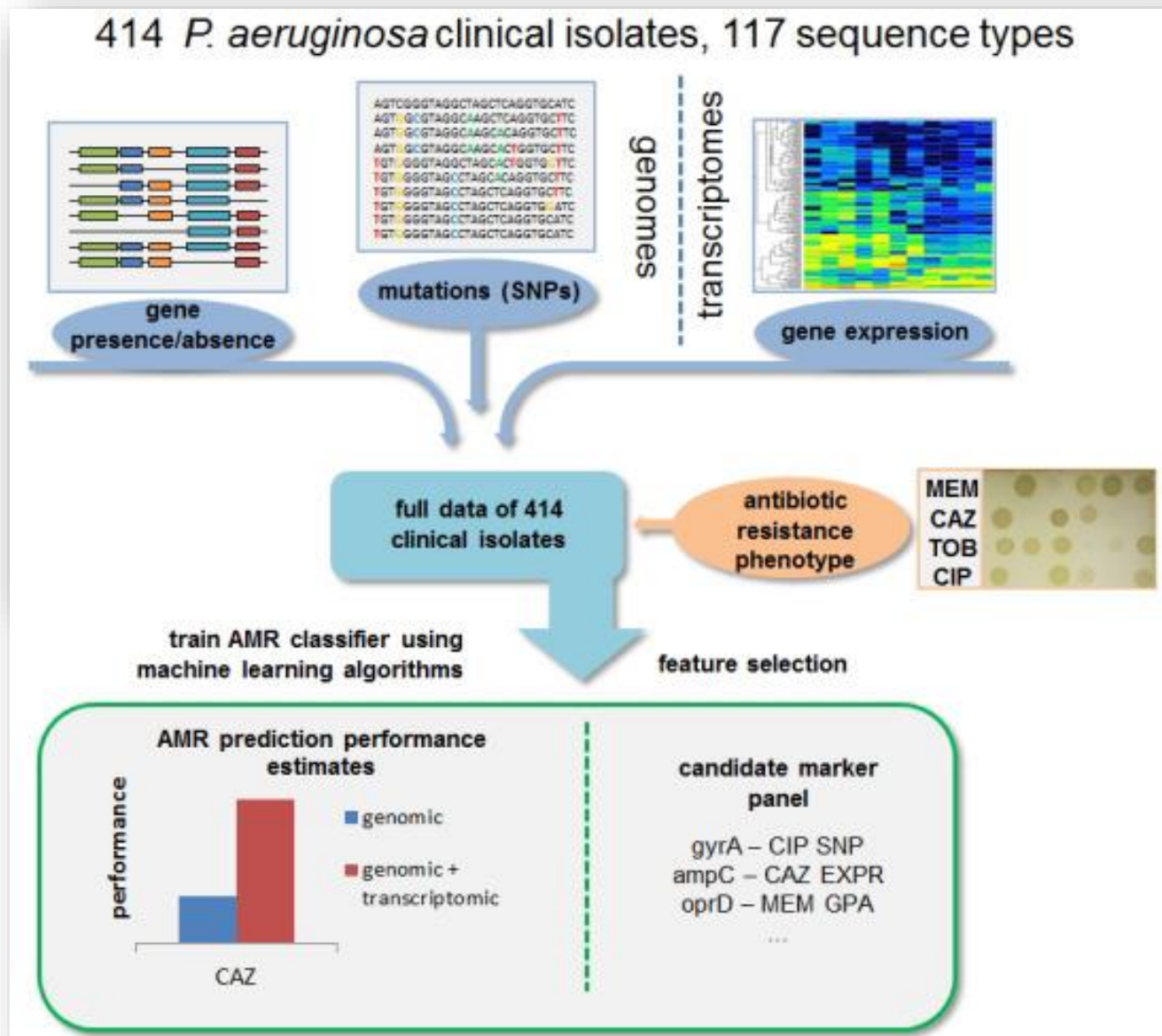


A convolutional neural network was able to discriminate between 18 classes of bacterial colonies.

Identification of pathogens



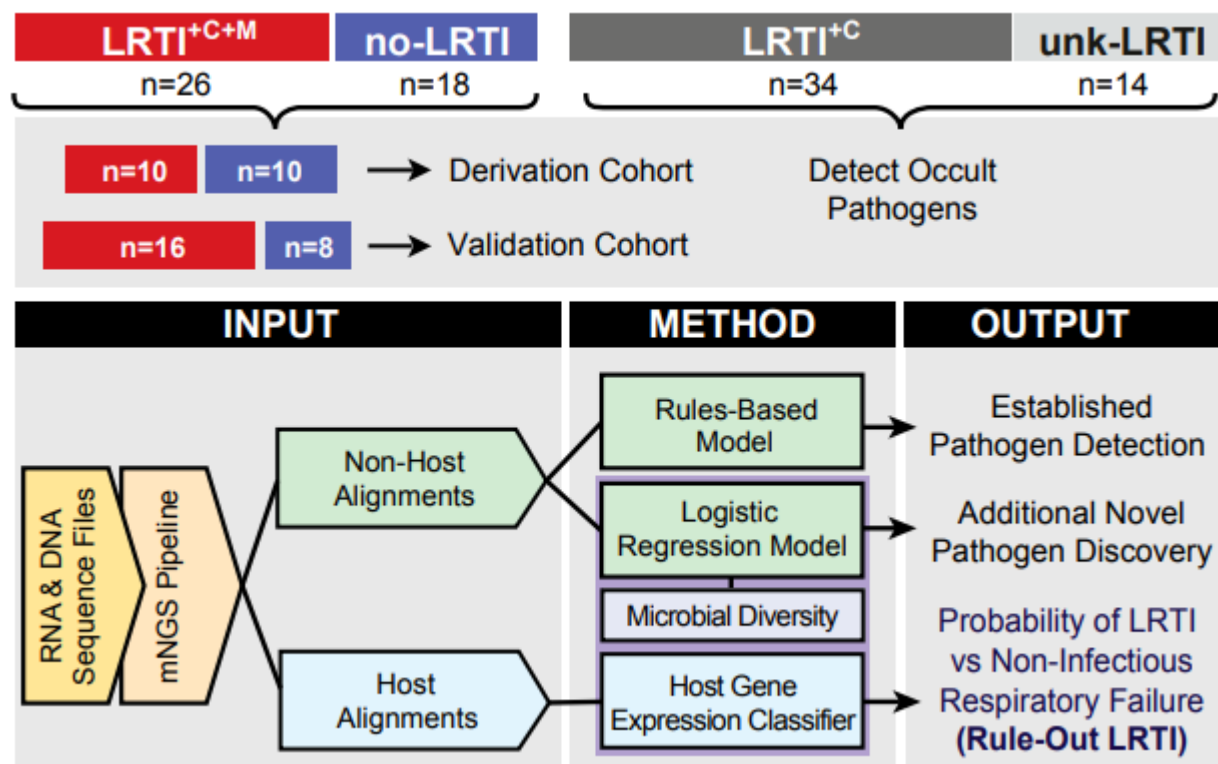
Antimicrobial resistance



Clinical outcomes

Table 1. Demographics and clinical characteristics of study cohort

Cohort characteristics	Cohort overall	LRTI ^{+C+M}	no-LRTI
Patient characteristics			
Total enrolled	92	26	18
Age, y	62 [†]	61	63
Female gender	31 (34%)	6	9
Race			
African American	5 (5%)	0	0
Asian	26 (28%)	10	8
Caucasian	50 (54%)	16	10
Other	11 (12%)	0	0
Hispanic ethnicity	8 (9%)	0	0
Comorbidities and outcomes			
Bacteremia	21 (23%)	10	11
Nonpulmonary infections	29 (32%)	16	13
COPD	12 (13%)	6	6
Diabetes mellitus	6 (7%)	3	3
Congestive heart failure	7 (8%)	4	3
Current smoker	12 (13%)	6	6
Immune suppression	41 (45%)	21	20
Solid-organ transplantation	13 (14%)	7	6
Prior antibiotic use	84 (91%)	42	42
Community acquired pneumonia	42 (46%)	21	21
Hospital acquired pneumonia	13 (14%)	6	7
Ventilator associated pneumonia	3 (3%)	1	2
30-d mortality	18 (20%)	9	9
Clinical metrics			
Max temperature, °C			
Max WBC count, 10 ⁶ cells/μL			
Max heart rate, bpm			
Max respiratory rate, breaths/min			
SIRS criteria, mean			
APACHE III score, mean	97	101	94
Pneumonia severity index, mean	151	148	137



Summary

- ❑ Machine Learning method is a powerful **data-driven alternative** to traditional modelling;
- ❑ One turns data into numbers (**features**) and trains a **generic algorithm** to discriminate between **labels** in the feature space;
- ❑ It is essential to have a separate **test set for evaluation** of the resulting predictor;
- ❑ In Microbiology, **a wide range of tasks** is already solved by Machine Learning.

Bi9680En: AI in Biology, Chemistry, and Bioengineering

- **Období: podzim**
- **Rozsah: přednáška 2 hodiny/týden**
- **Vyučující: Dr. Stanislav Mazurenko**
- **Osnova:**
 - modern bio-challenges: drug design, DNA interpretation, protein engineering
 - types of AI algorithms and workflow for designing predictors
 - clustering algorithms, random forests, artificial neural networks
 - features, databases, and predictors used in applications

