# Deep learning approaches for deciphering composition and functional roles of the ocean microbiome
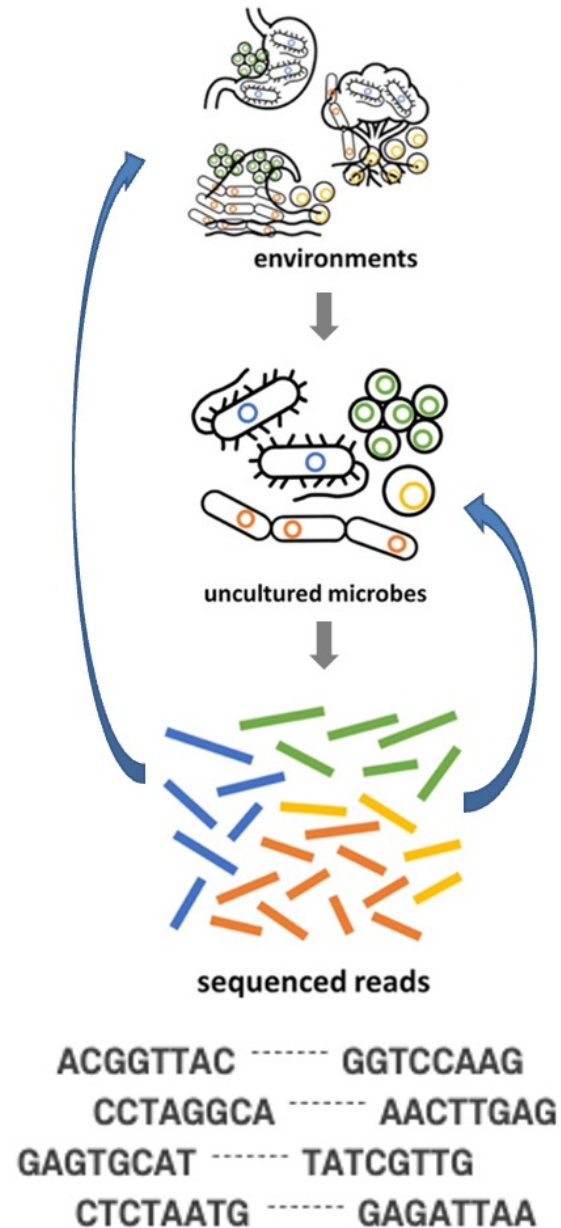
Mina Rho
Hanyang University

# Introduction to the marine microbiome

# Nutrient cycling in marine microbes



An example of energy production in *Phototrophicus methaneseepsis* ZRK3

# Close relationship with human as one health framework



https://nih.go.kr/

nt. Microbiol.(2016)   Front Cell Infect Microbiol 2023
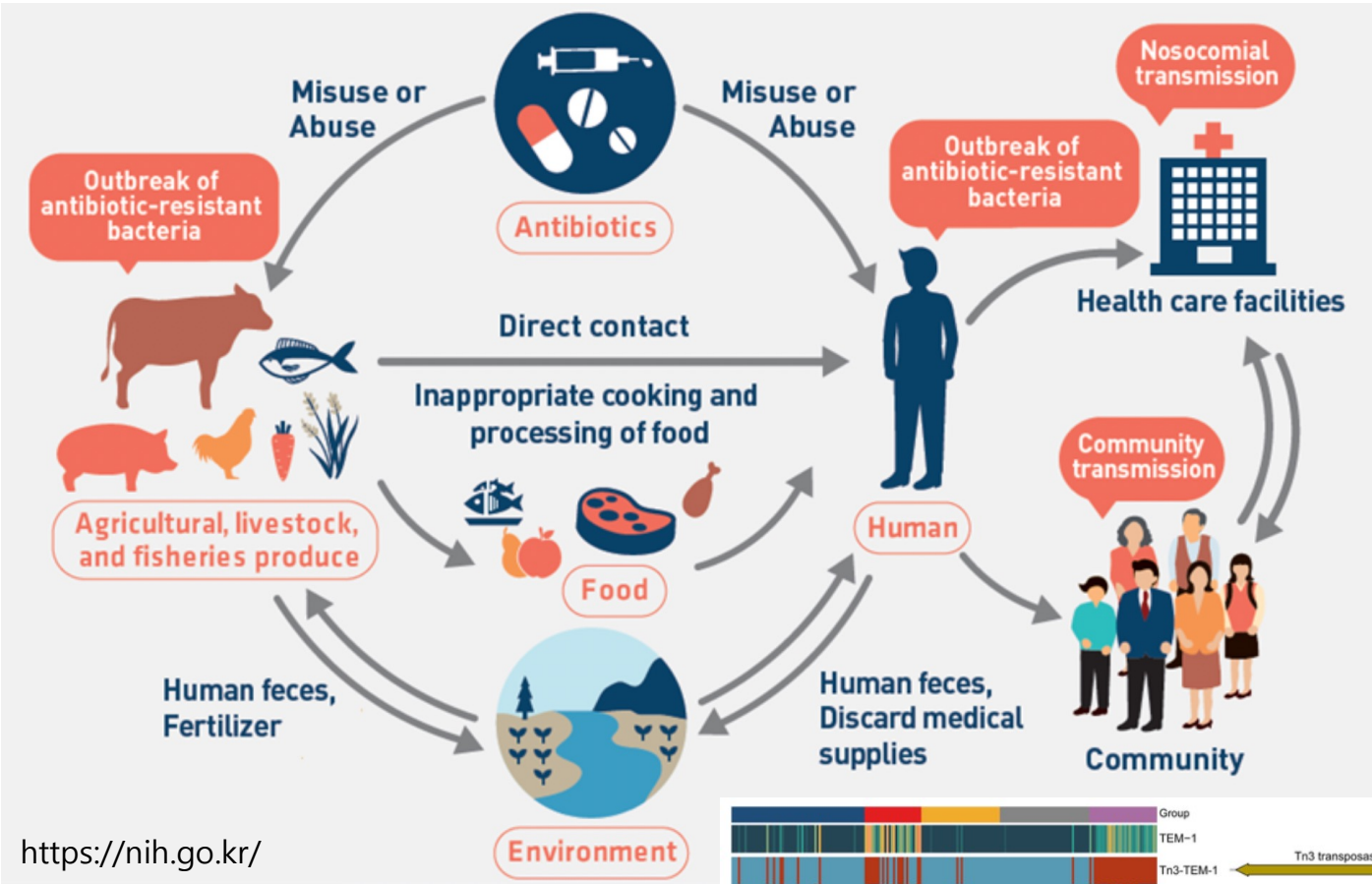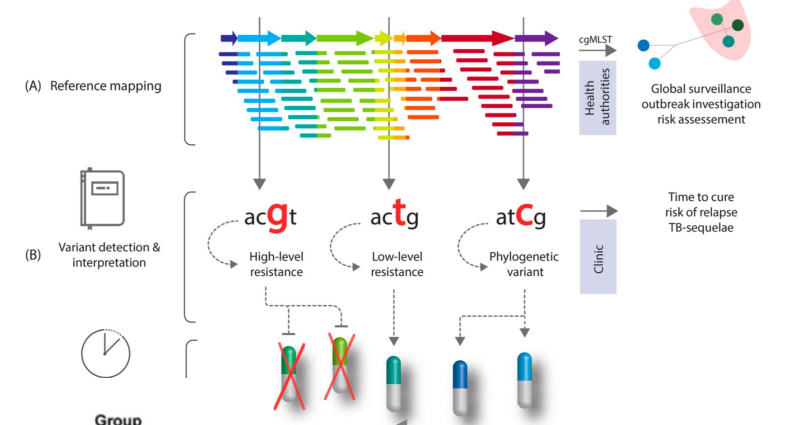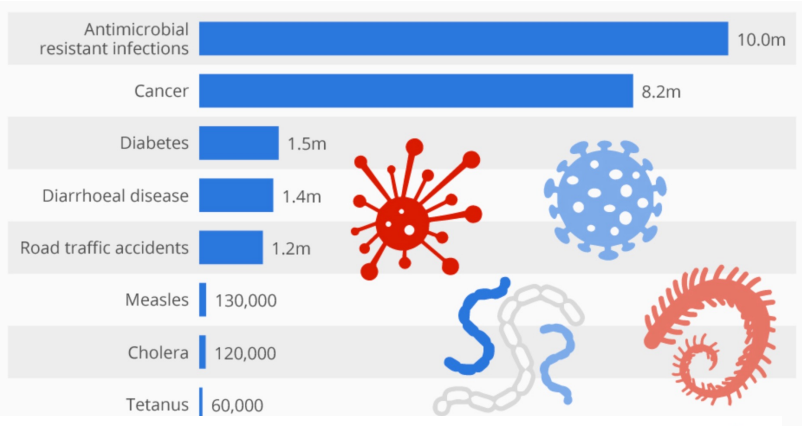
# AIs for biomedical problems

# Large language models

# Large language models to biological sequence data

- It consists of a series of words

  **Genome / Gene / Protein** ----- **Document / Sentence**

  **Nucleotides / Amino acids** ----------- **Words**

- Local relationship between words

  **Motifs** ------------- **Idioms / word order**

- Remote relationship between words

  **Co-evolution** -------------- **Context**

# Self-supervised learning to biological sequencing data

- Using sequence fragments of 150 -250 bps, the sequence should be assigned (classified) to one of the taxa at a certain taxonomy level

10,119 viral genomes
    2,293  eukaryotic DNA viruses
    2,733  eukaryotic RNA viruses

5,529 bacterial genomes
    (one genome for each species)

Eukaryote-related viruses were retained
    (Based on ICTV annotation)

The length of reads was set to 151 and 251
Insert size from 300 to 800

# VIBE: Taxonomy classification from sequencing read sequence

## Order-level classifier for RNA viruses

- SOTA methods:
  - CHEER / Skip-Gram + parallel CNN model
  - Kraken2 / $k$-mer homology method
- Test data:
  - Read-level validation set
  - Genome-level validation set

# VIBE: Taxonomy classification from sequencing read sequence



- **Betacoronavirus**
- **Alphacoronavirus**
- **Gammacoronavirus**
- **Deltacoronavirus**
- **SRR14403295**

- The ViBE was re-trained without the SARS-CoV-2 reference genome

- COVID-19 samples were tested with the model

- (A) 4-mer frequency, (C) embedded vector by fine-tuned model

Four different hemoglobin protein sequences

```
CAA37898.1    -----------MSTLEGRGFTE--EQEALVVKSWSAMKPNAGELGLKFFLKIFEIA
P68871.2      -----------MVHLTPEEKSA-------VTALWG-KV-NVDEVGGEALGRLLVVY
CAA77743.1    MHSSIVLATVLFVAIASASKTRELCMKSLEHAKVG-TSKEAKQDGIDLYKHMFEHY
AAA29796.1    MHSSIVLATVLFVAIASASKTRELCMKSLEHAKVG-TSKEAKQDGIDLYKHMFEHY
                :    :       :          .        :. : * .    :::

CAA37898.1    KLFSFLKDSNVPL--ERNPKLKSHAMSVFLMTCESAVQLRKAGKVTVRESSLKKLGASHF    105
P68871.2      RFFESFGDLSTPDAVMGNPKVKAHGKKVLG-AFS-------DGL----AHLDNLKGTFAT    88
CAA77743.1    KYFKHRENY-TPADVQKDPFFIKQGQNILL-ACHVLCATY-DDR----ETFDAYVGELMA    112
AAA29796.1    KYFKHRENY-TPADVQKDPFFIKQGQNILL-ACHVLCATY-DDR----ETFDAYVGELMA    112
                : *.     :   .*        :* .  :. .::    :            .              *

CAA37898.1    KHGVAD-------EHFEVTKFALLETIKEAVPETWSPEMKNAWGEAYDKLVAAIKLEMKP    158
P68871.2      LSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANALAHK---    145
CAA77743.1    RHE--RDHVKIPNDVWNHFWEHFIEFLG--SKTTLDEPTKHAWQEIGKEFSHEISHHGRH    168
AAA29796.1    RHE--RDHVKVPNDVWNHFWEHFIEFLG--SKTTLDEPTKHAWQEIGKEFSHEISHHGRH    168
                : :.       ::  :              : *: :       .    : .
```

- There is no hard sequence similarity threshold for "safe" function prediction

- Sequences that are more than 30-40% identical are considered as having the same or a very similar function

# His-Me finger endonucleases

## His-Me finger nuclease

- Conserved Histidine residue (His)
- Catalytic metal ion (Me)
- Finger-like structure (finger)

## Example

- Cas9 enzyme
  in the CRISPR-Cas9 genome editing technology

Capturing conserved residues and secondary structure
is essential to classify His-Me finger nuclease

**Classes in SCOPe**

| | |
|---|---|
| a: | All alpha proteins |
| b: | All beta proteins |
| c: | Alpha and beta proteins (a/b) |
| d: | Alpha and beta proteins (a+b) |
| e: | Multi-domain proteins (alpha and beta) |
| f: | Membrane and cell surface proteins and peptides |
| g: | Small proteins |



A    B

Capturing **conserved residues** and **secondary structure** is essential to classify His-Me finger nuclease

Pooling: light attention
- **1D-CNN + max pooling**: capturing secondary structure
- **Convolution x attention**: capturing conserved residues

Using sequence-level representation,

**binary classification** for detecting His-Me finger nuclease

and **secondary structure class classification**

were performed <u>simultaneously.</u>

# FuncPred: Function prediction from protein sequences



Aggregation method employing both convolution and attention outperformed other pooling methods

Aggregation method successfully focused on both
conserved residues and conserved structures

# FuncPred: Function prediction from protein sequences

His-Me finger proteins in the SCOPe database are remote homologous against the Pfam database

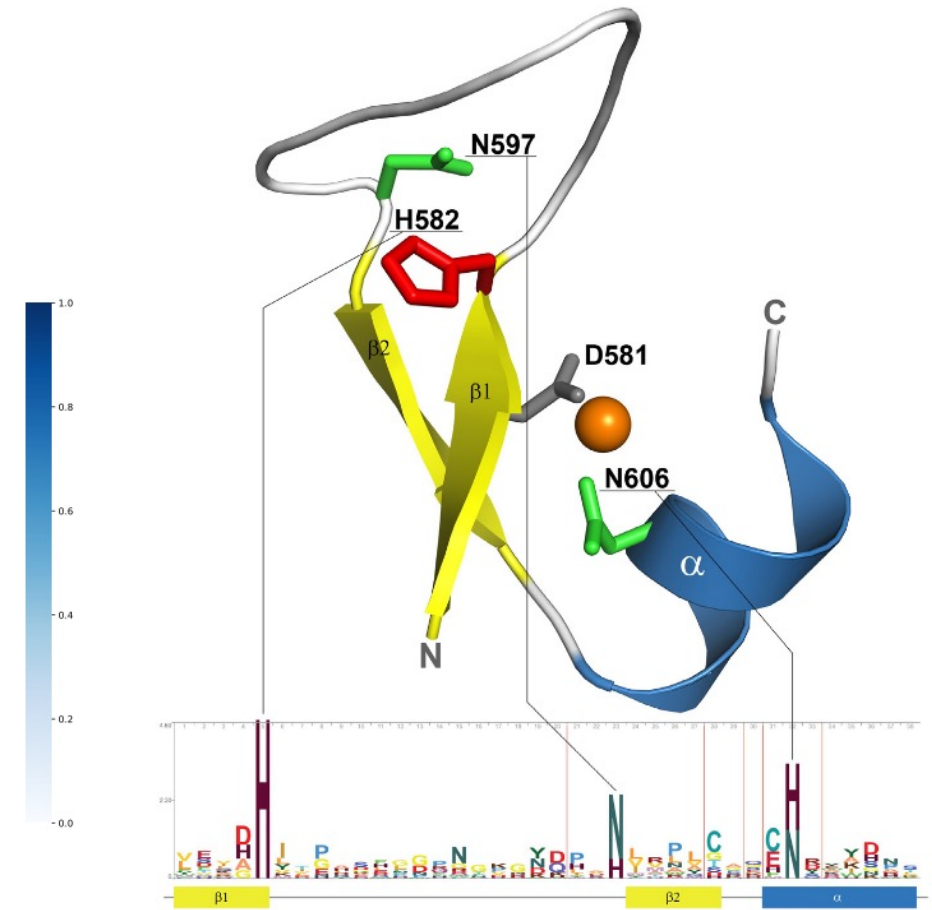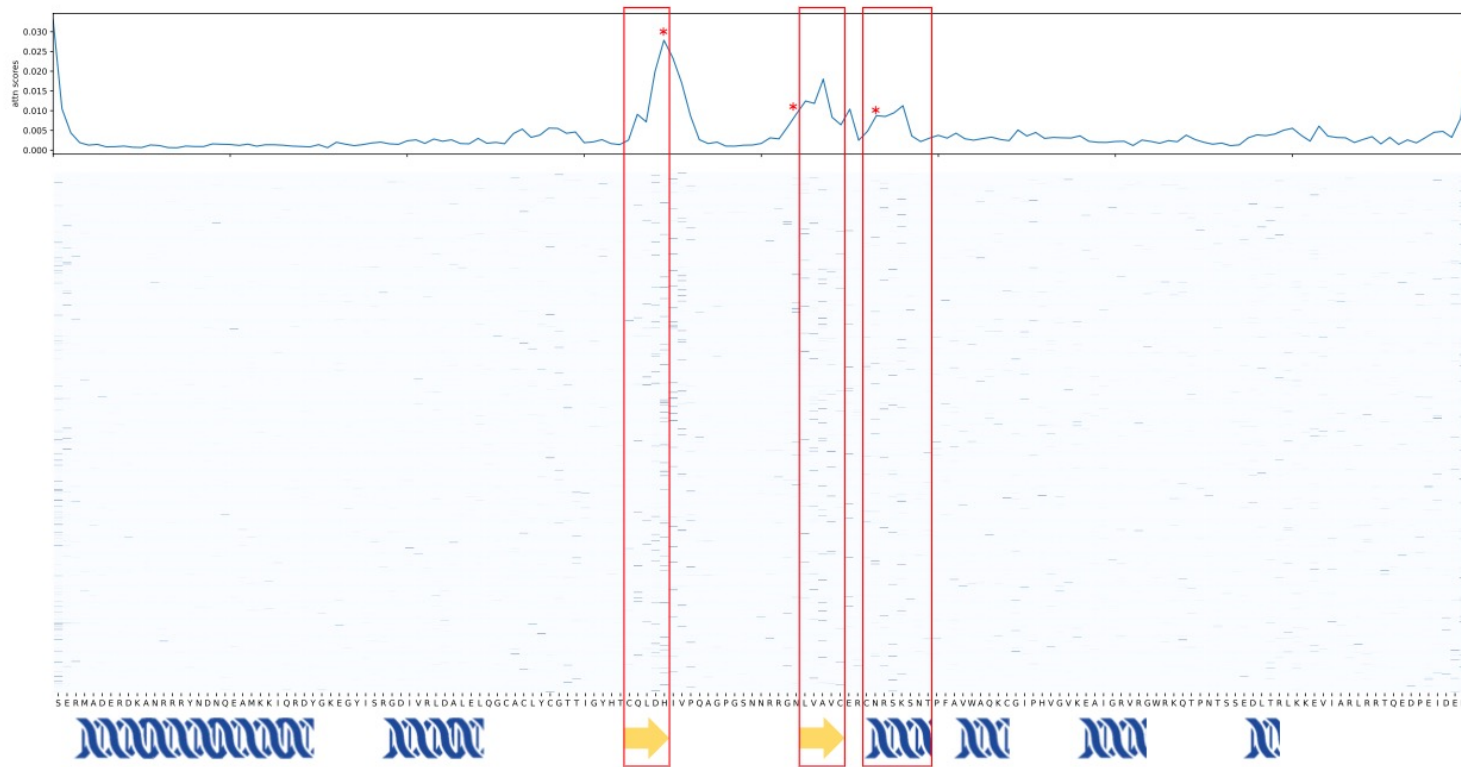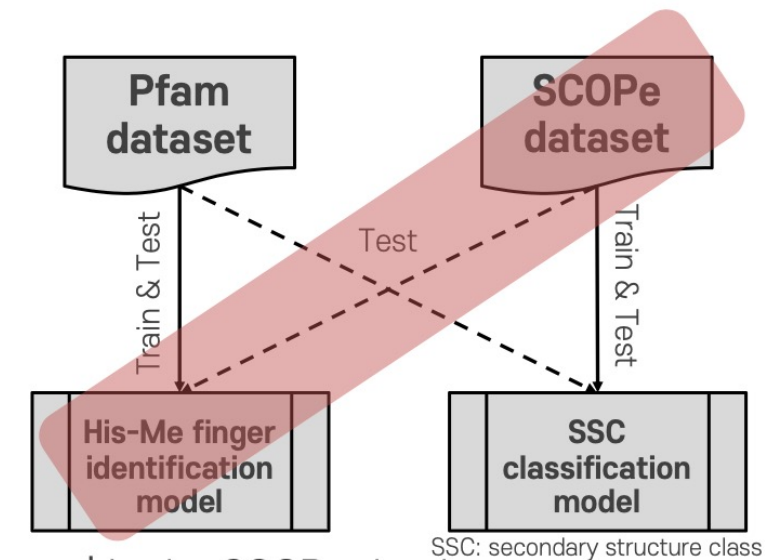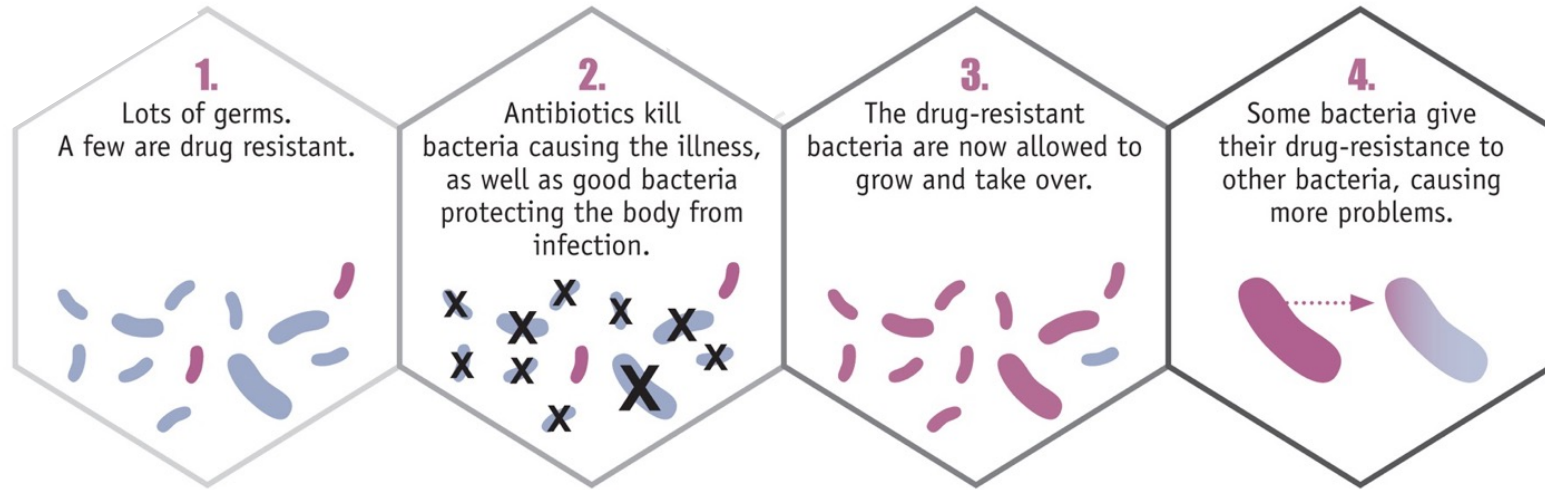| query ID | subject ID | percent identity | query coverage | alignment length | query length | query start | query end | subject start | subject end | e-value |
|---|---|---|---|---|---|---|---|---|---|---|
| d1zm8a_ | NUCA_NOSS1/43-263 | 99.55 | 92.5% | 221 | 239 | 8 | 228 | 1 | 221 | 6.42E-164 |
| d4e3ya_ | A0A240BSN3_SERFI/24-245 | 94.12 | 92.1% | 221 | 240 | 1 | 221 | 2 | 222 | 2.27E-161 |
| d1u3em1 | B6V2J8_BPSP1/53-98 | 100.00 | 43.8% | 46 | 105 | 53 | 98 | 1 | 46 | 1.59E-29 |
| d1a73a_ | A0A1E3PUR8_LIPST/1-90 | 46.15 | 24.1% | 39 | 162 | 88 | 126 | 16 | 54 | 4.E-06 |
| d1e7la2 | A0A2H4YFJ9_9CAUD/1-97 | 90.72 | 94.2% | 97 | 103 | 1 | 97 | 1 | 97 | 1.15E-63 |
| d2pu3a_ | Q5E7R4_ALIF1/26-232 | 91.30 | 100.0% | 207 | 207 | 1 | 207 | 1 | 207 | 3.01E-146 |
| d1v0da_ | A0A0H2UHU4_RAT/123-344 | 96.85 | 90.6% | 222 | 245 | 19 | 240 | 1 | 222 | 7.69E-164 |
| d4ogca2 | F9PLJ4_9ACTO/560-606 | 97.87 | 29.2% | 47 | 161 | 54 | 100 | 1 | 47 | 1.01E-29 |
| d5axwa2 | K9B5K9_9STAP/65-120 | 82.14 | 38.1% | 56 | 147 | 51 | 106 | 1 | 56 | 4.05E-29 |
| d4oo8a2 | CAS9_STRP1/821-872 | 98.08 | 39.1% | 52 | 133 | 47 | 98 | 1 | 52 | 1.52E-30 |
| d6w0va_ | A0A1I6EC40_9FIRM/69-117 | 38.71 | 24.8% | 31 | 125 | 94 | 124 | 15 | 45 | 2.E-04 |
| d3qsvd_ | A0A671YF54_SPAAU/33-134 | 98.02 | 80.8% | 101 | 125 | 25 | 125 | 1 | 101 | 5.23E-73 |

☐ Remote homologous proteins

☐ Highly homologous proteins, which was not reported as His-Me finger nuclease



SSC: secondary structure class

1. Lots of germs. A few are drug resistant.

2. Antibiotics kill bacteria causing the illness, as well as good bacteria protecting the body from infection.

3. The drug-resistant bacteria are now allowed to grow and take over.

4. Some bacteria give their drug-resistance to other bacteria, causing more problems.

10124

16889

☐ Non-pathogen
☐ Pathogen

j

IS

Integrase    Transposase

<5kb    <5kb

MGE1    ARG    MGE2

Integrase/Transposase

0  1000 2000 3000 4000

Number of MGEs

- Antibiotic resistance genes are often located on plasmids or transposons and can be transferred from cell to cell by conjugation, transformation, or transduction

- Resistome is a collection of all the antibiotic resistance genes and related elements in bacteria.

# Take-home message

- Deep learning discovers multiple representations with different levels of abstractions for genomes and chemical compounds

- Generalization and specialization can be achieved by self-supervised learning in the pre-training process and task-specific learning in the fine-tuning process

- Large language model can be used to understand genomic characteristics and interactions of various taxa and functions

# Acknowledgements