# Addressing bioinformatics challenges in pathogen agnostic detection from metagenomic sequencing

November 19, 2024

Jonathan Allen, Ph.D.
Informatics Scientist
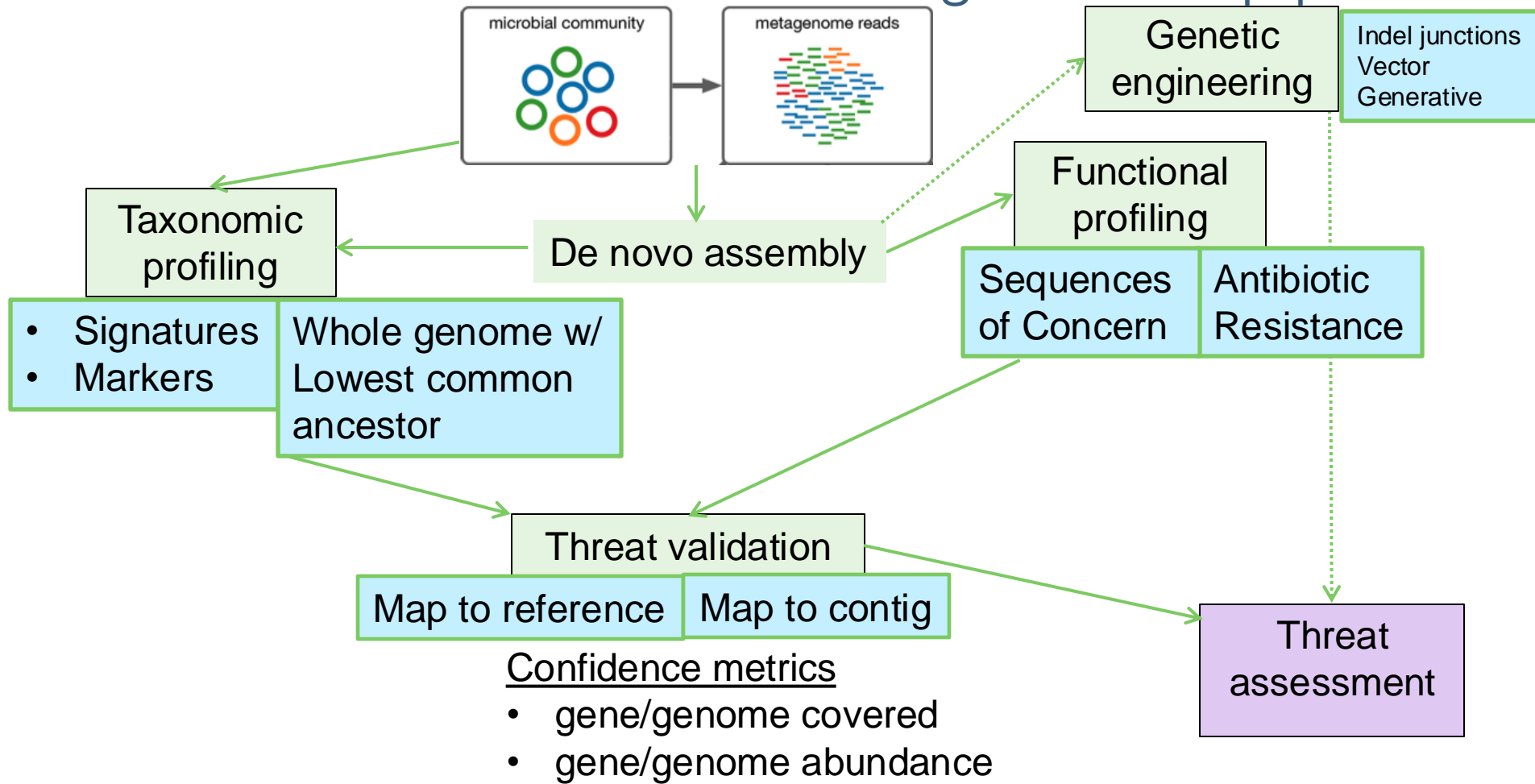
**Lawrence Livermore National Laboratory**

# Pathogen agnostic detection for human health

## Input sample type

- Clinical samples
  - Blood, nasal swabs
- Air filters
  - Dirty (subway, outdoor)
  - Clean(er) (building/space station)
- Wastewater
- Clean rooms

## Organism detection needs

- All human infecting pathogens
- All animal infecting pathogens
- Species emerging from natural background
- Engineered or generative design
- Truly novel – extra terrestrial or otherwise highly divergent

## Kingdom

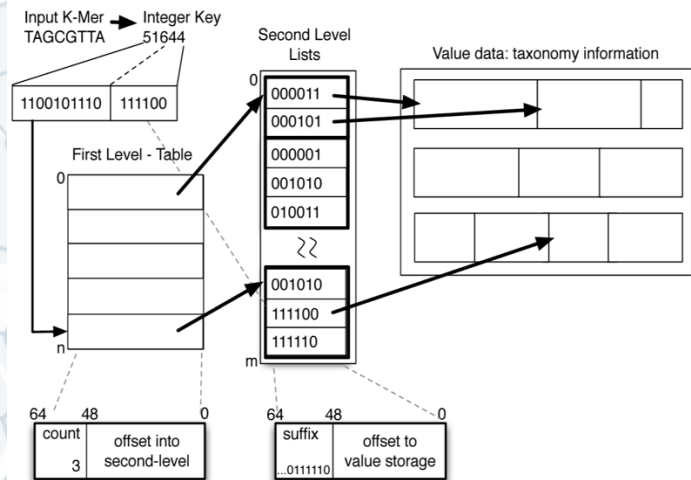- Fungi
- Bacteria
- Microeukaryote
- RNA virus
- DNA virus

# Current state of the art metagenomics pipeline

# Maintaining a reference database for *comprehensive* metagenomic classification

We published 1st scalable k-mer search index 2013 (LMAT)

Ames et al., Bioinformatics 2013



A key innovation was to include draft genomes

- Microbial
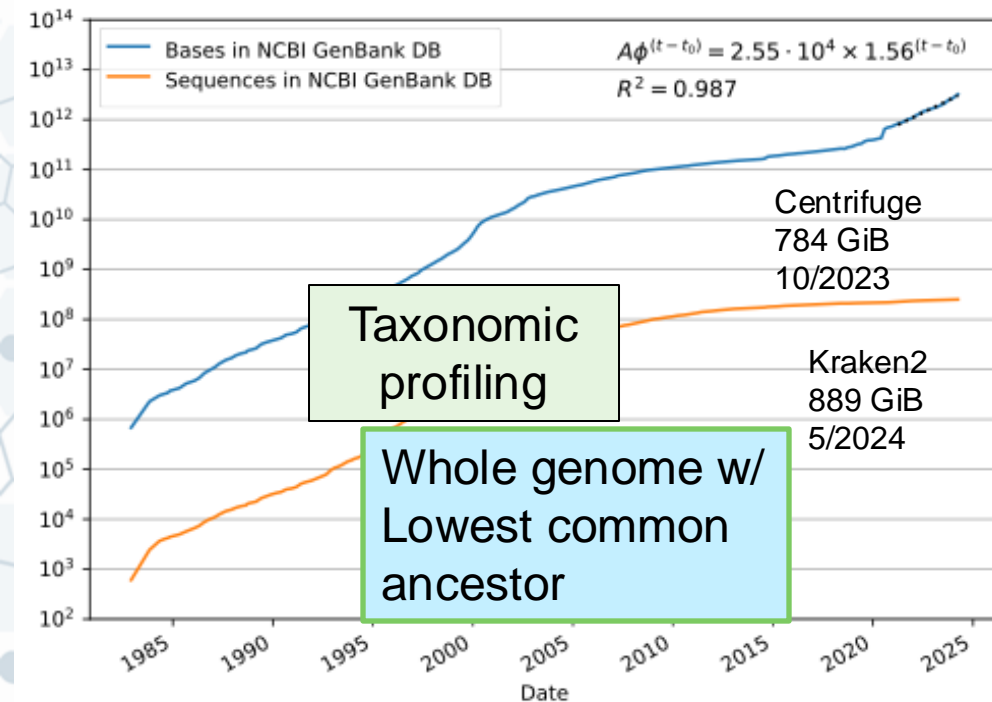- Human

- DB size = 612 GB

We found human sequences in the human microbiome project in screen against 1000 human genomes

## Using populations of human and microbial genomes for organism detection in metagenomes

Sasha K. Ames,[1] Shea N. Gardner,[2] Jose Manuel Marti,[3] Tom R. Slezak,[2] Maya B. Gokhale,[1] and Jonathan E. Allen[2]

[1]Center for Applied Scientific Computing, Lawrence Livermore National Laboratory, Livermore, California 94550, USA; [2]Global Security Computer Applications Division, Lawrence Livermore National Laboratory, Livermore, California 94550, USA; [3]Instituto de Física Corpuscular, CSIC-UVEG, E-46980 Valencia, Spain

# Using the right database for *comprehensive* metagenomic classification is a challenge
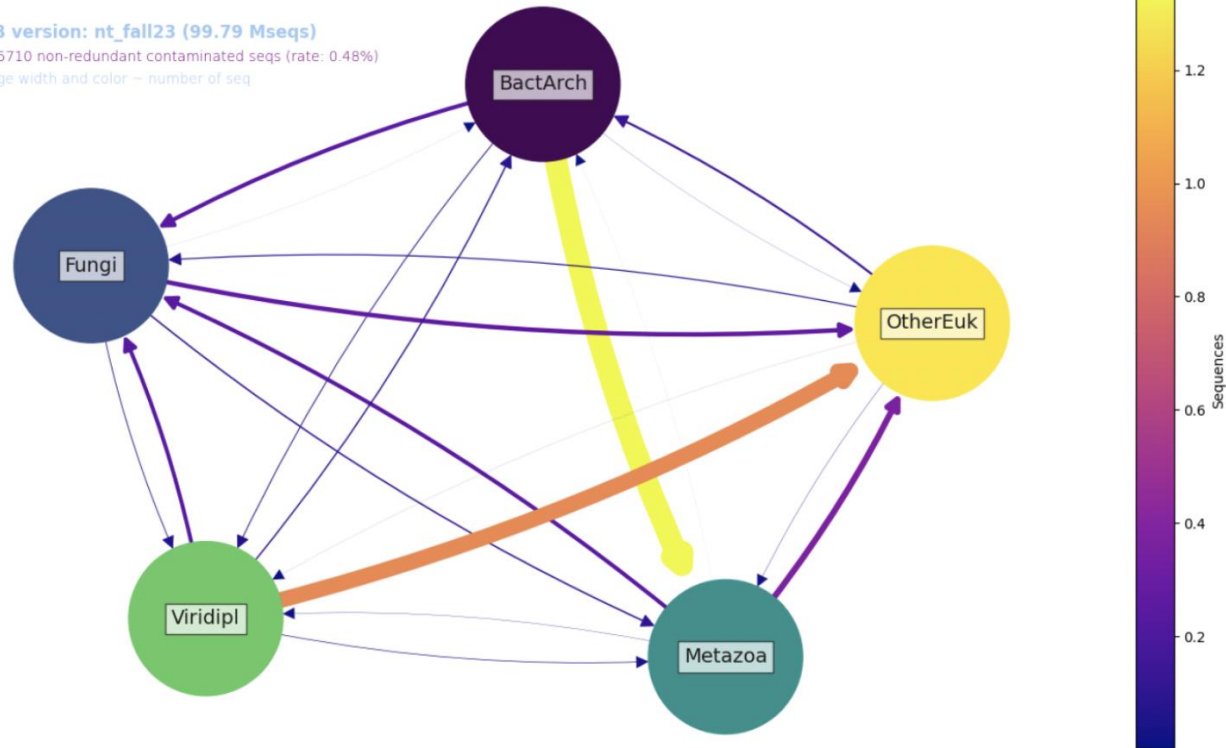


Marti JM, et al. bioRxiv 2024.

CZ ID stopped using full NT database in May 2024 and will use lossy compression

Kraken2's database is not documented for decontamination

We developed the most complete decontaminated DB to date with a published protocol

# Decontamination of reference sequences is a growing computational burden



Now takes 2+ weeks to run on our cluster

Conterminator software

# nt DB content evolution over time

# Two different filtering strategies: by score and abundance

**SCORE**

- "**Destructive**" filter: reads not achieving the threshold quality are removed from the analysis.

- The filter is **inactive by default**.

- Can be set differently for regular and negative control samples.

- Most efficient strategy is running the classifier with a score-permissive value to get a high classification ratio and, once the distribution of scores in known, then "recentrifuge" the samples to the desired threshold.

**ABUNDANCE**

- "**Non-destructive**" filter: reads not achieving the threshold are moved to the parent taxonomic level in a recursive manner until the threshold condition is satisfied.

- The filter is **active by default**.

- Default (automatic) threshold is different per sample: it depends on the log of the total number of reads [relative abundance].

- Can be manually set to a fix value, which can be differently for regular and negative control samples [absolute abundance].

# What is the source of current unknowns?

Waste water sample from contra costa county after human DNA removal

Source of unknowns
- Sequence divergence
- Eukaryotes

Score (avg)

Kickxellaceae

Homalozoon vermiculare

Obtectomera

Acrosorium

Entamoeba

# Unclassified reads present challenges for detection of novel pathogens



Meta2DB: Kok et al, 2024 bioRxiv

Unclassified reads across ~12K human-associated microbiome

# Path forward with traditional genomic search

Make more reference data searchable!!!!

- Metagenomic assembled contigs

- More eukaryotic draft genomic sequences

Pros

- Search algorithms are well understood and relatively scalable with sufficient compute

- Potential for explainability: matches can be assigned to individually sequenced organisms and environments
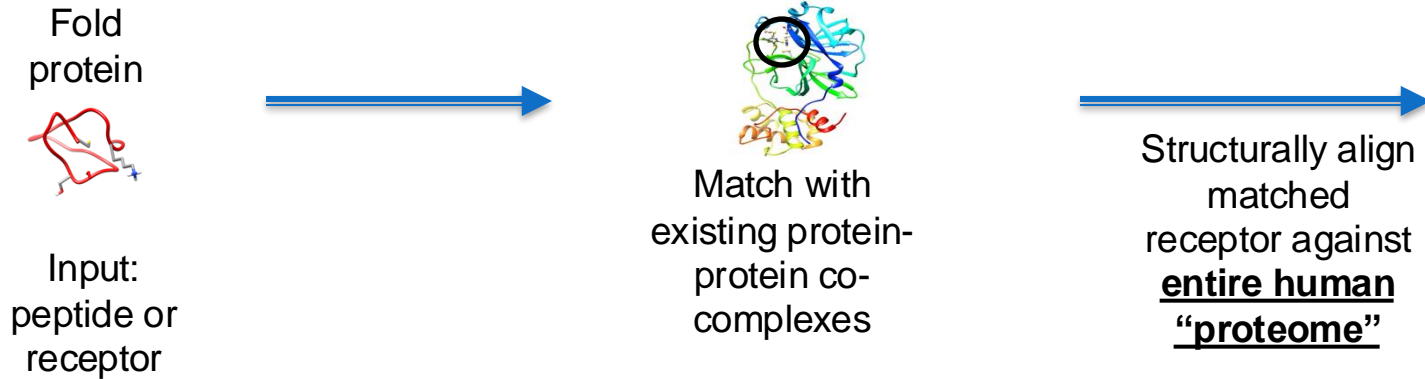
Cons

- Computational cost to build database AND every search is costly

Challenges

- Automated quality control of new reference sequences and incorporating uncertainty into taxonomy calls

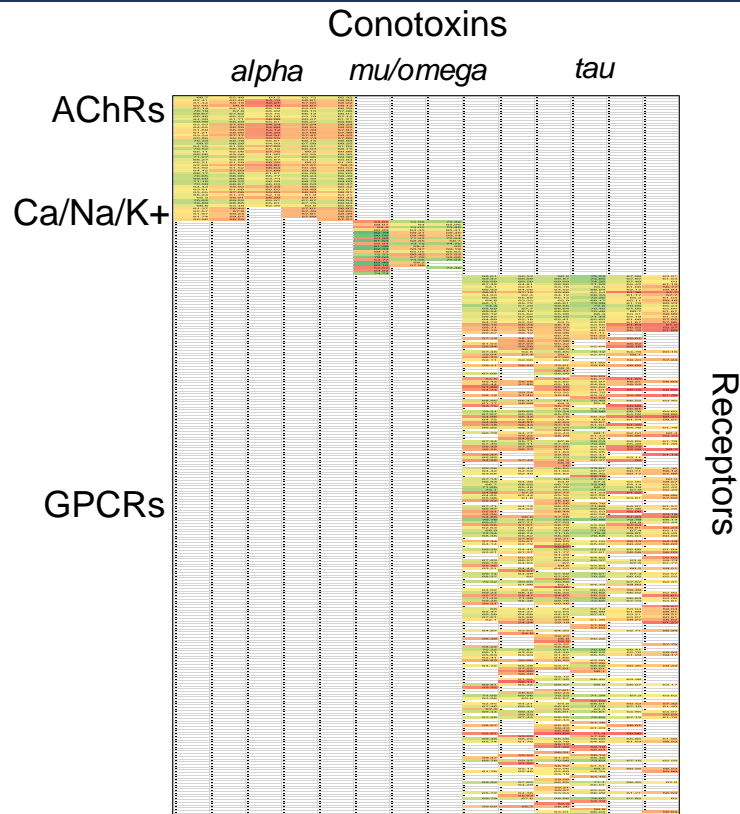# How do we conduct threat assessment for unknowns?

Model pathogen-host interactions
Predict protein-protein interactions
Predict physiologic impact of molecular interactions

Fold
protein



Input:
peptide or
receptor

Match with
existing protein-
protein co-
complexes

Structurally align
matched
receptor against
**entire human
"proteome"**

Use PDBSpheres to structurally align matching templates
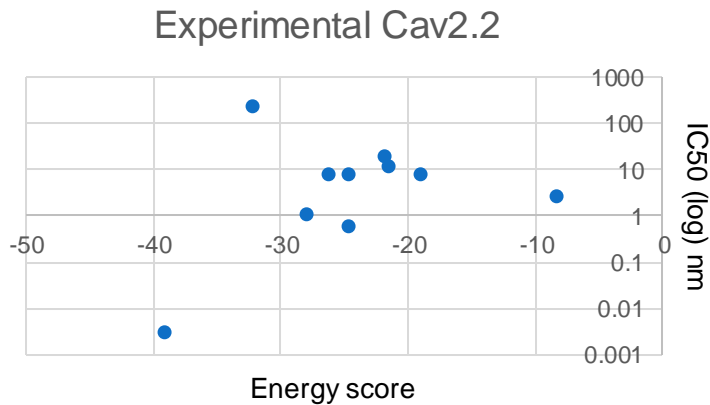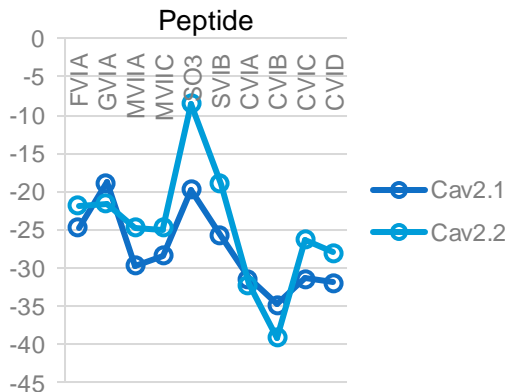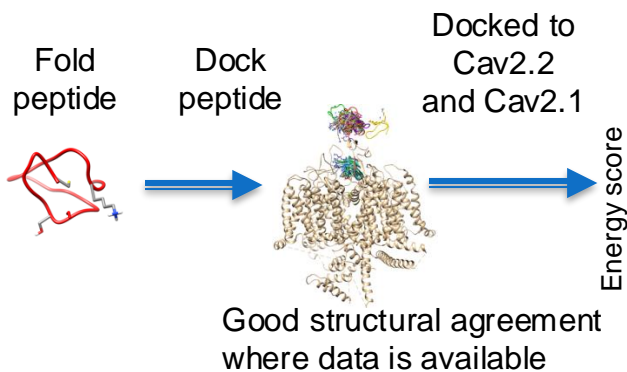Zemla et al., 2023

# Structural alignment identifies receptors with potential functional interference

Selected peptide-protein co-complexes are searched against all human protein pockets with available structures



Conotoxins — *alpha*, *mu/omega*, *tau*

Receptors — AChRs, Ca/Na/K+, GPCRs

# Docking simulations model peptide binding to capture molecular dynamics

Fold peptide → Dock peptide → Docked to Cav2.2 and Cav2.1

Good structural agreement where data is available



We are evaluating docking methods for conotoxins against all human cell receptors with comparisons to experimental data (patch clamp and binding assays)

Current computational cost: ~13 minutes per peptide/receptor pair

# Uncovering microbial dark matter

- Do more with traditional search by searching more reference material than currently done

    - Presents significant engineering challenges: contamination in reference sequences, matching to samples not just reference genomes

- Expand traditional search into 3D structure space and learned feature representations

    - Fundamental limitation to assessing protein function in context of a single protein

- Infer microbial pathogens through improved molecular interaction modeling

    - Mechanisms replication

    - Host cell attachment

    - Microbial transmission

- Major challenges with mapping molecular level mechanisms to physiologic impact:

    - Host to host transmission and disease outcomes

Lawrence Livermore
National Laboratory

# Acknowledgements

Jose Manuel Marti - Recentrifuge
Crystal Jaing - Genomics
Nick Be – Microbiome
Car Reen Kok - Bioinformatics
Ed Lau – MD simulation
Heesung Shim – MD simulation
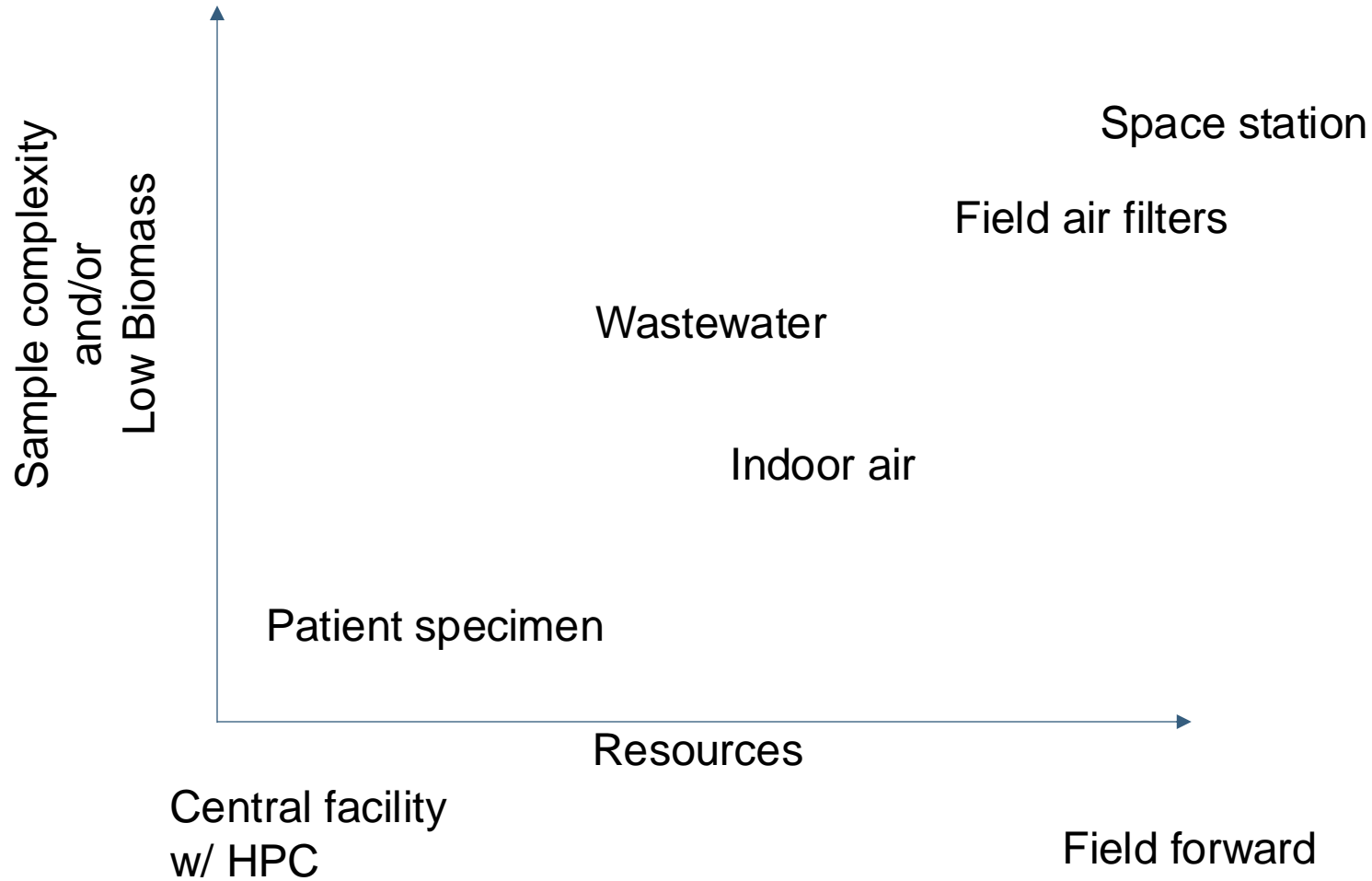
Bioinformatics and Genomics

Data science and Machine Learning

Thank you

# Limitations from operational setting

Example protein language model ESM3

Pretrain model on large corpus of sequence data

- ESM3, AlphaFold3, Chia-1,
- LucaOne/Prot
- ProteInfer
- Evo

Some evidence for improvement in remote homology detection



Lawrence Livermore National Laboratory