



HPC pielietojums bioinformātikā

Edgars Liepa
2024



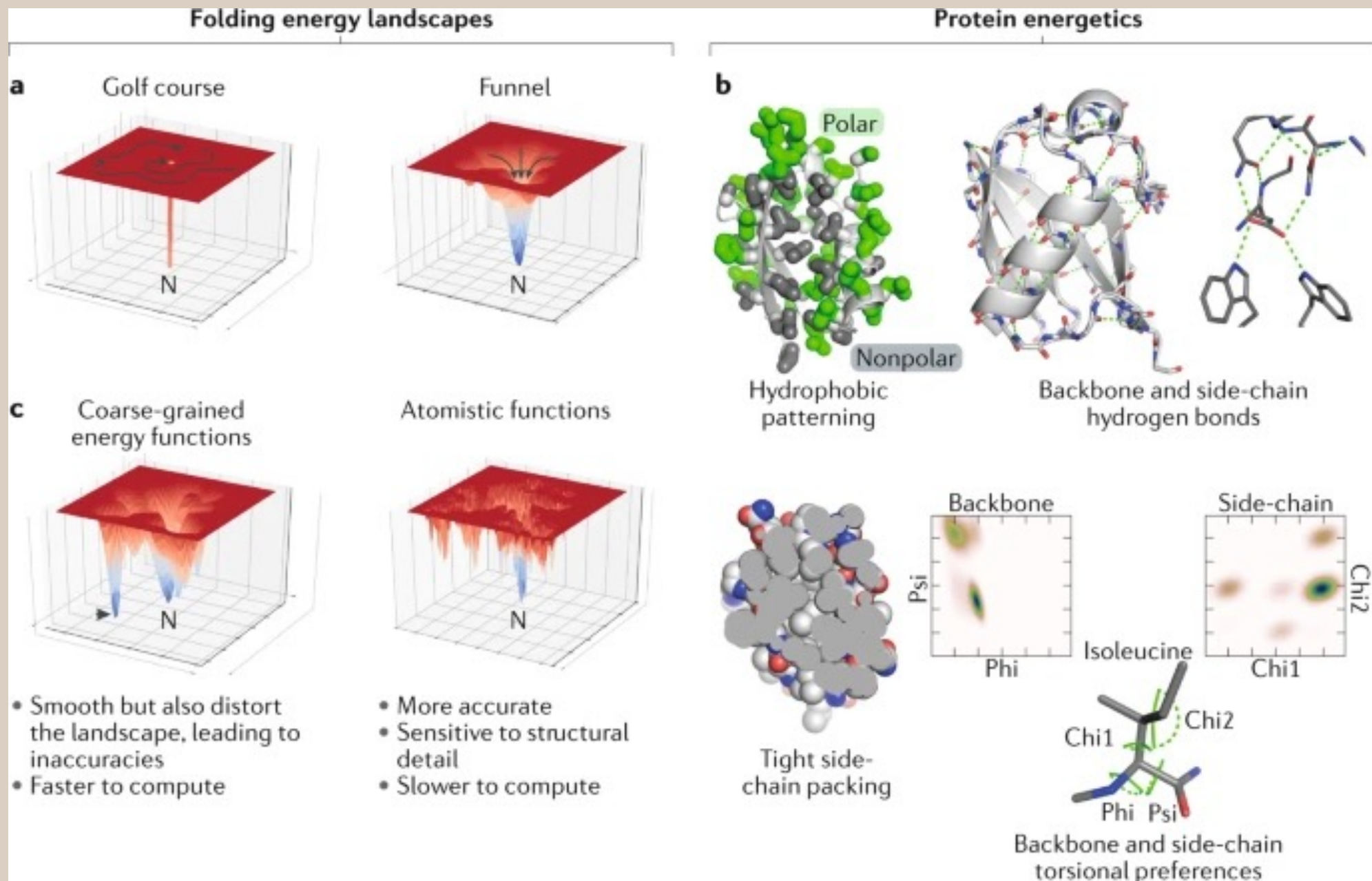
Par mani

- Bioinformātikis LV BMC
 - D. Fridmaņa grupa
- PhD students LU BF

Web: edgarsliepa.lv

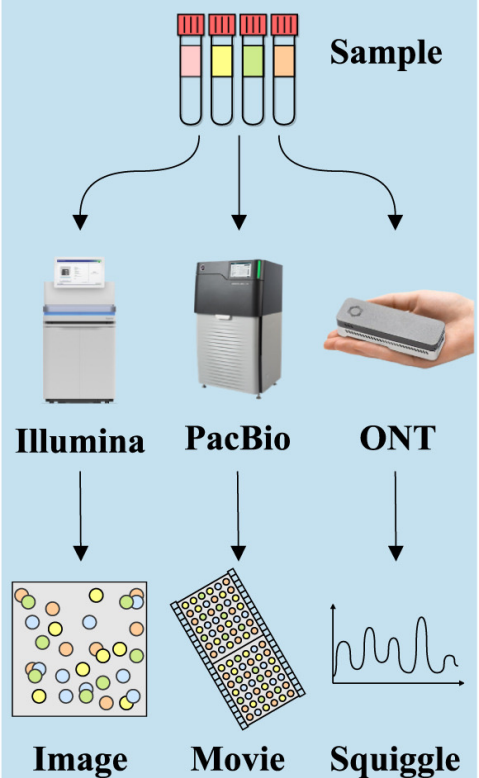
LinkedIn: [linkedin.com/in/edgars-liepa/](https://www.linkedin.com/in/edgars-liepa/)

Github: github.com/EdgarsLiepa

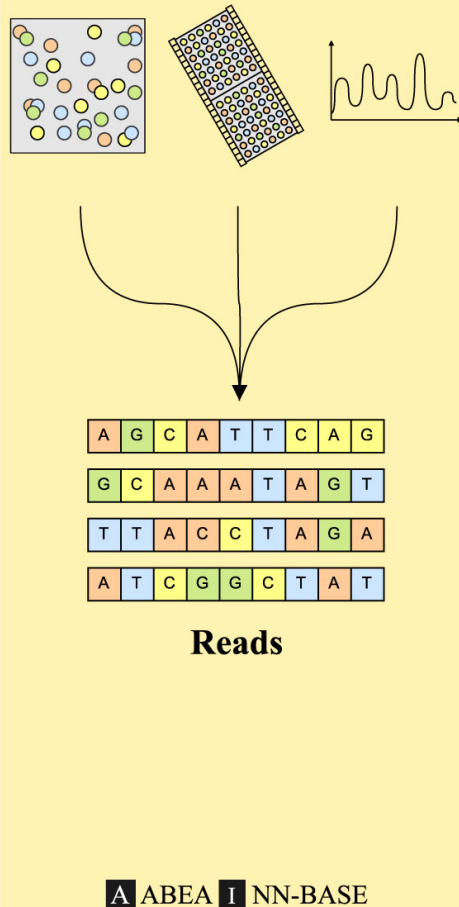


McMaster, Benjamin, et al. "Can AlphaFold's breakthrough in protein structure help decode the fundamental principles of adaptive cellular immunity?." *Nature Methods* (2024): 1-11.

1. Sequencing



2. Basecalling



3.a. Genome Resequencing

1. Read Mapping

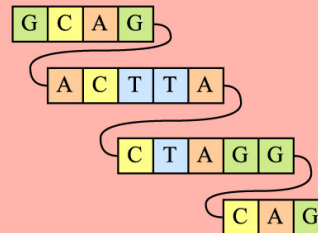
1. Seed

Index	
Seed	Position
A T C	1, 10
A T G	2, 4, 5

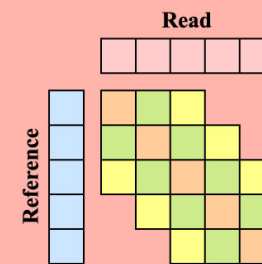
Reference
=?
Read

G FMI

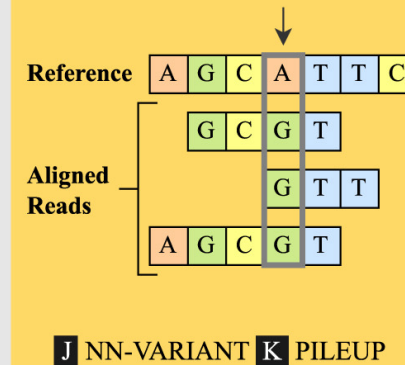
2. Chain



3. Extend



2. Variant Calling



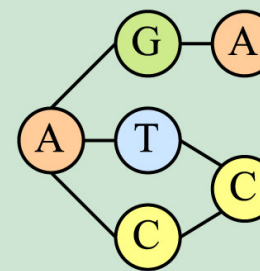
3.b. Genome Assembly (De-Novo Assembly)

1. K-mer Counting

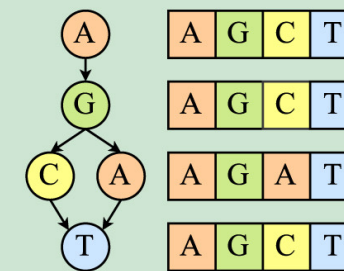
k-mer	#
A T A	120
A T T	9
A T C	412
A T G	986

H KMER-CNT

2. De-Bruijn Graph Construction

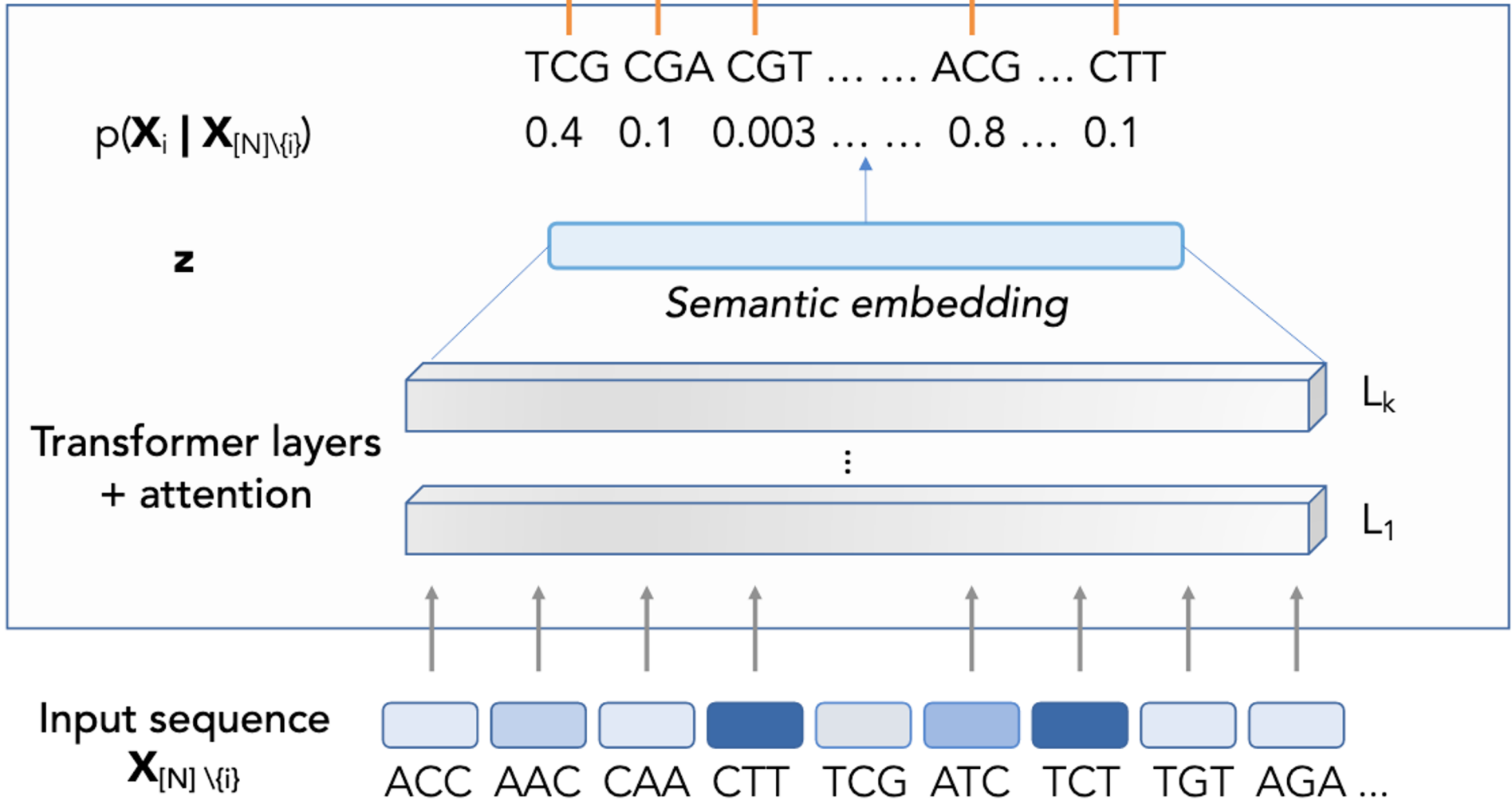


3. Multiple Sequence Alignment





GenSLM



$$p(\mathbf{X}_i | \mathbf{X}_{[N] \setminus \{i\}})$$

\mathbf{z}

Semantic embedding

Transformer layers
+ attention

L_k

⋮

L_1

Input sequence

$$\mathbf{X}_{[N] \setminus \{i\}}$$

ACC AAC CAA CTT TCG ATC TCT TGT AGA ...

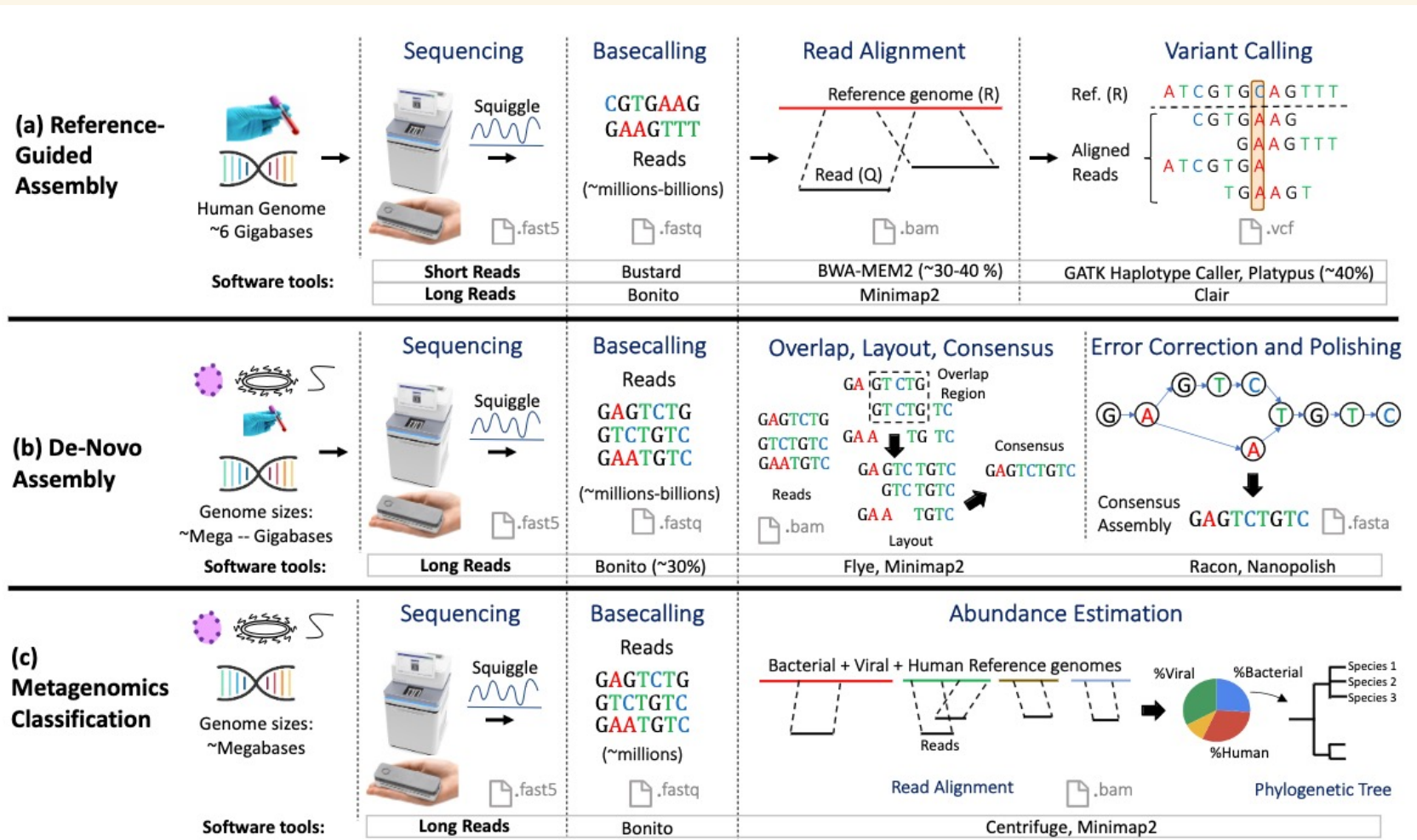
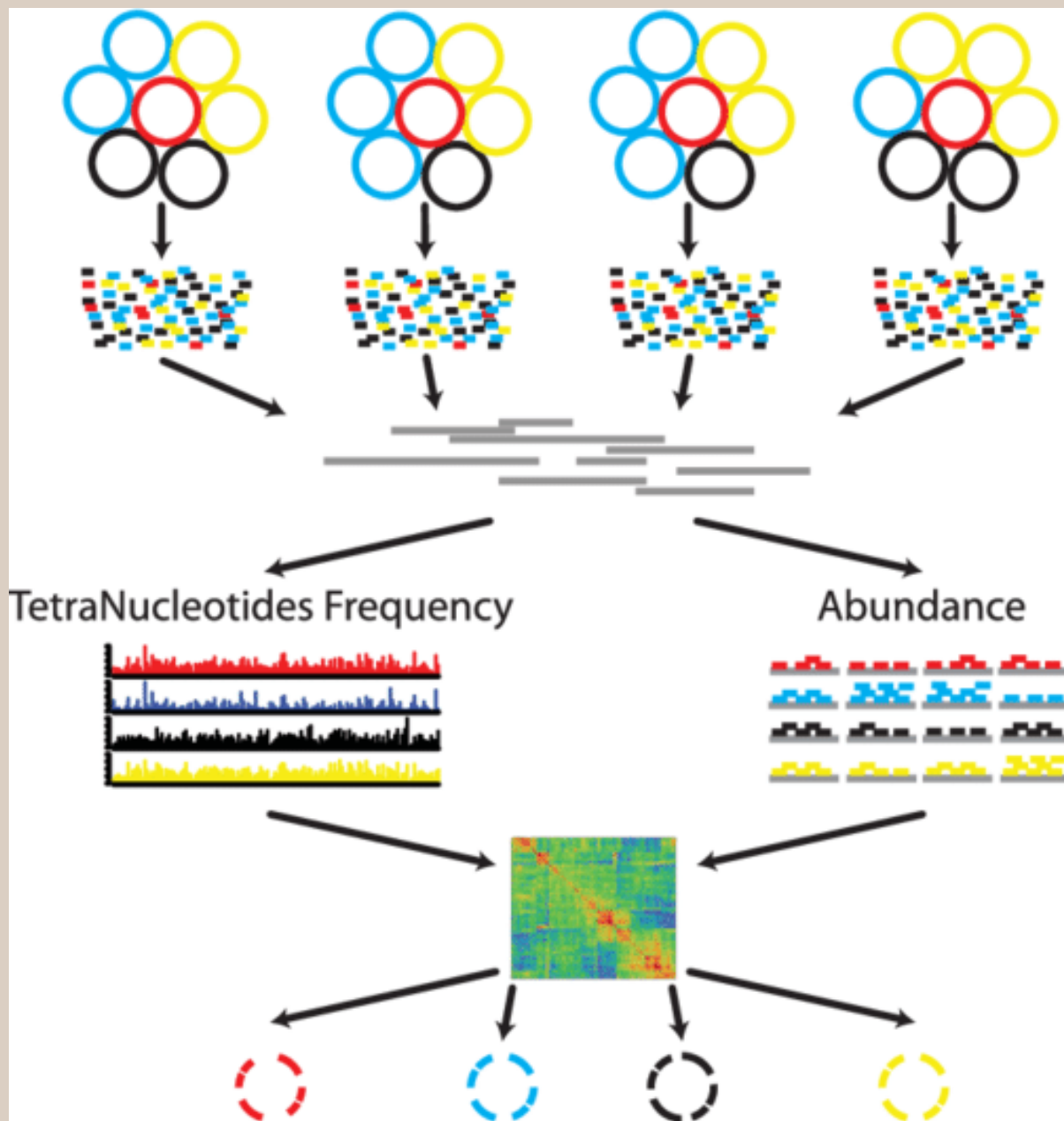


Fig. 1. Common workflows in genomics

Subramaniyan, Arun, et al. "Genomicsbench: A benchmark suite for genomics." *2021 IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*. IEEE, 2021.



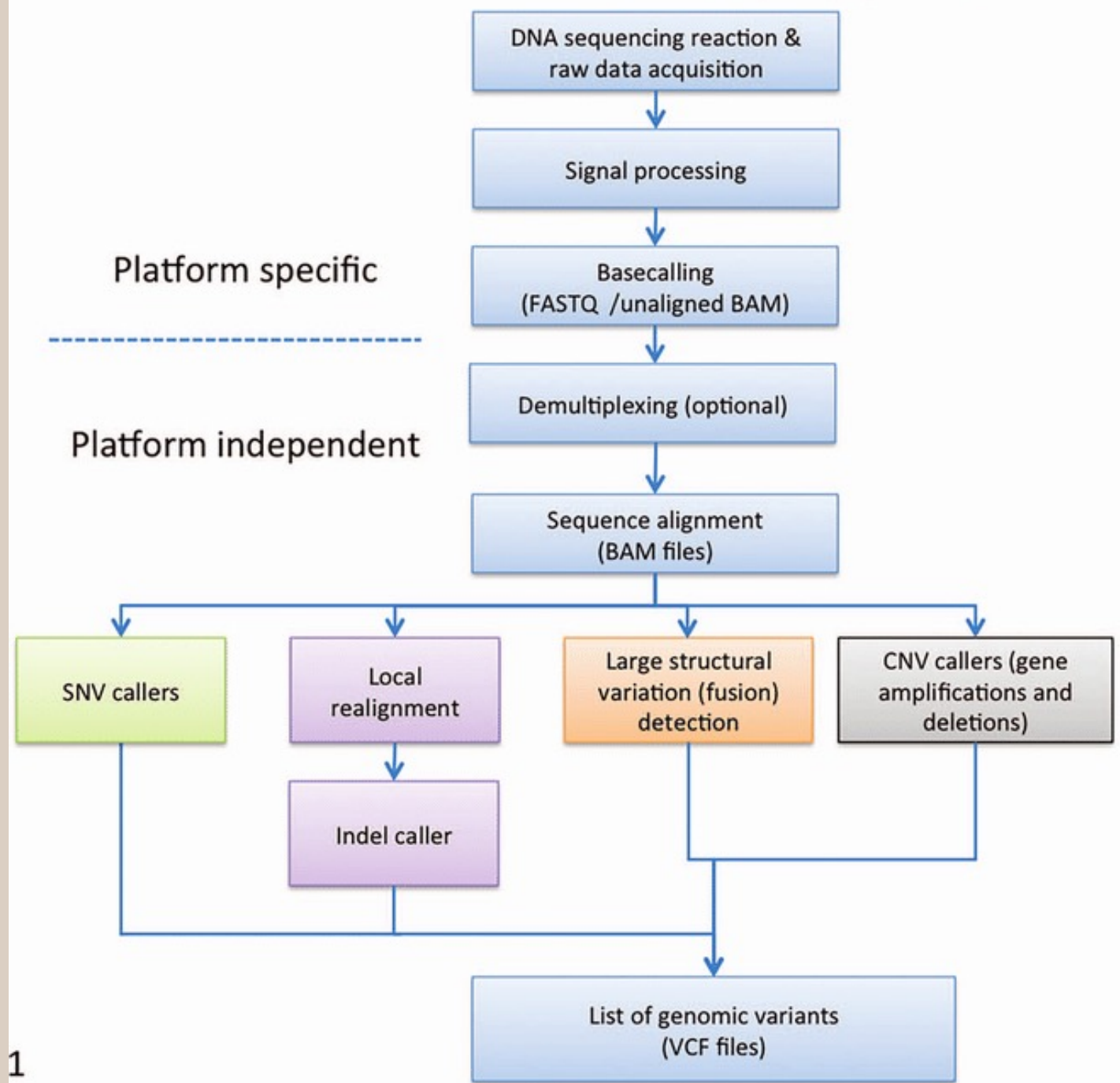
Preprocessing

- 1 Samples from multiple sites or times
- 2 Metagenome libraries
- 3 Initial de-novo assembly using the combined library

MetaBAT

- 4 Calculate TNF for each contig
- 5 Calculate Abundance per library for each contig
- 6 Calculate the pairwise distance matrix using pre-trained probabilistic models
- 7 Forming genome bins iteratively

NGS data analysis pipeline



File size range

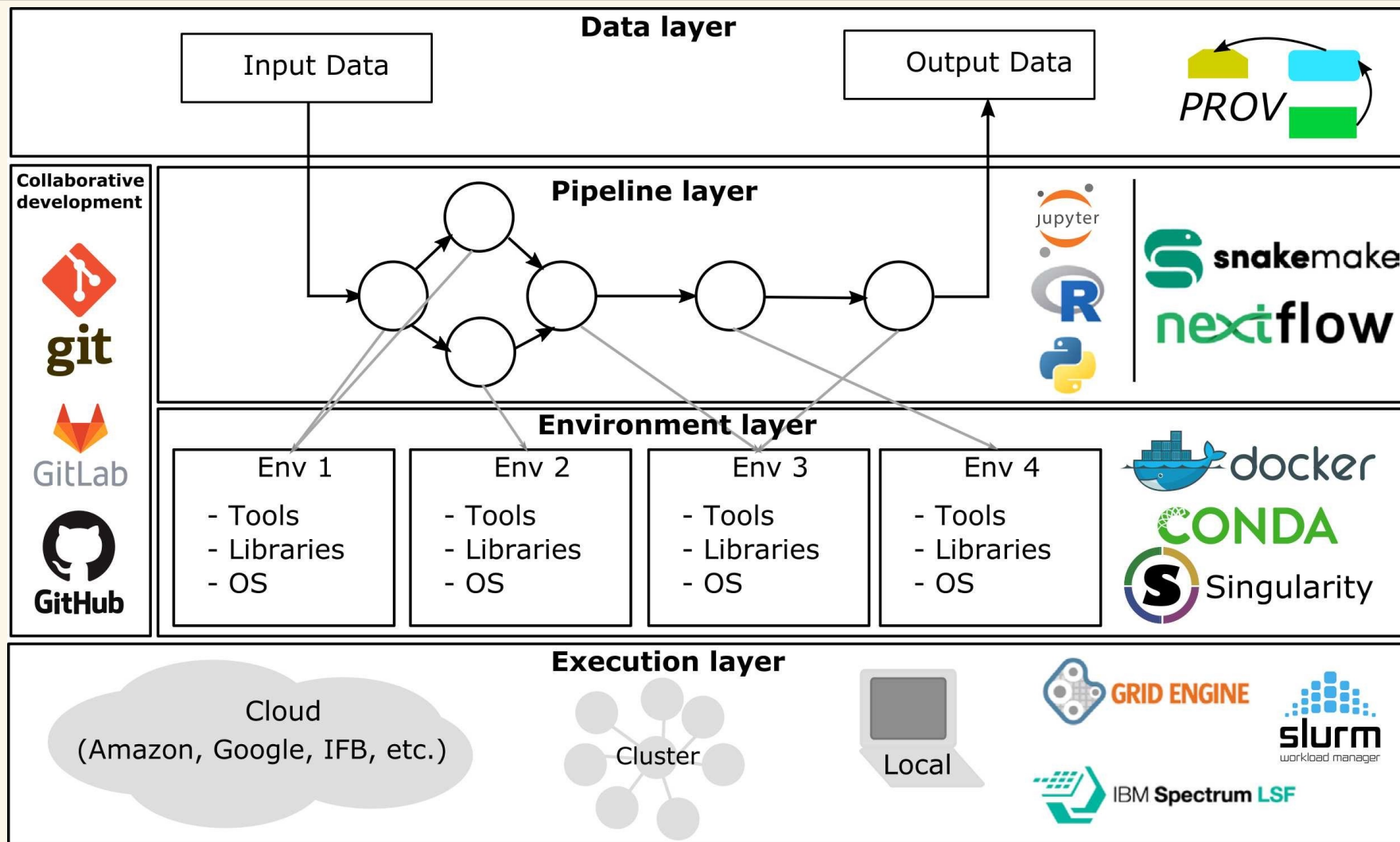
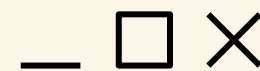
Terabytes

Gigabytes

Megabytes / Kilobytes

Roy, Somak, et al. "Next-generation sequencing informatics: challenges and strategies for implementation in a clinical environment." *Archives of pathology & laboratory medicine* 140.9 (2016): 958-975.

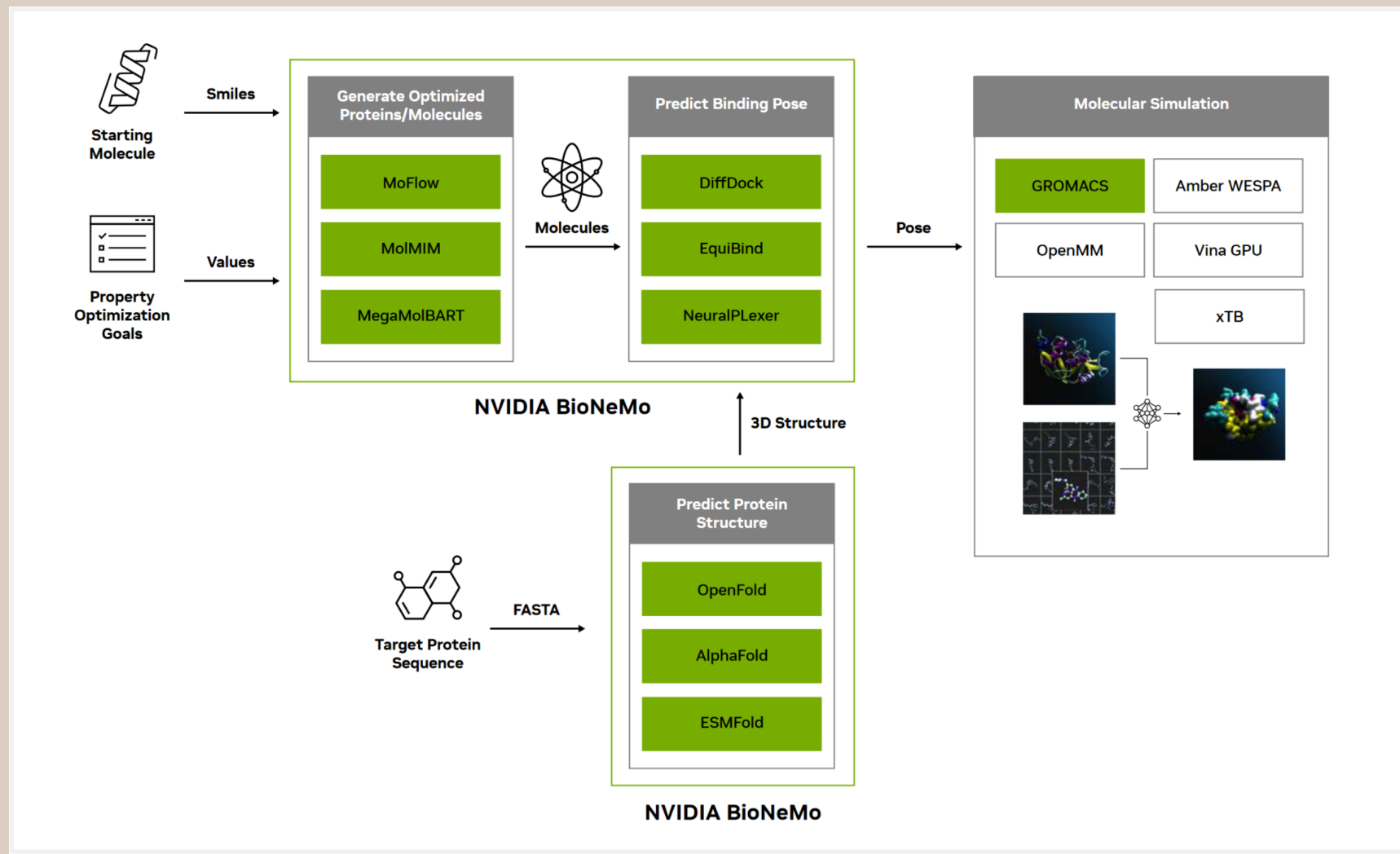
Zinātniskās darba plūsmas pārvaldības sistēma



Djaffardjy, Marine, et al. "Developing and reusing bioinformatics data analysis pipelines using scientific workflow systems." *Computational and Structural Biotechnology Journal* 21 (2023): 2075-2085.

Small Molecule Discovery

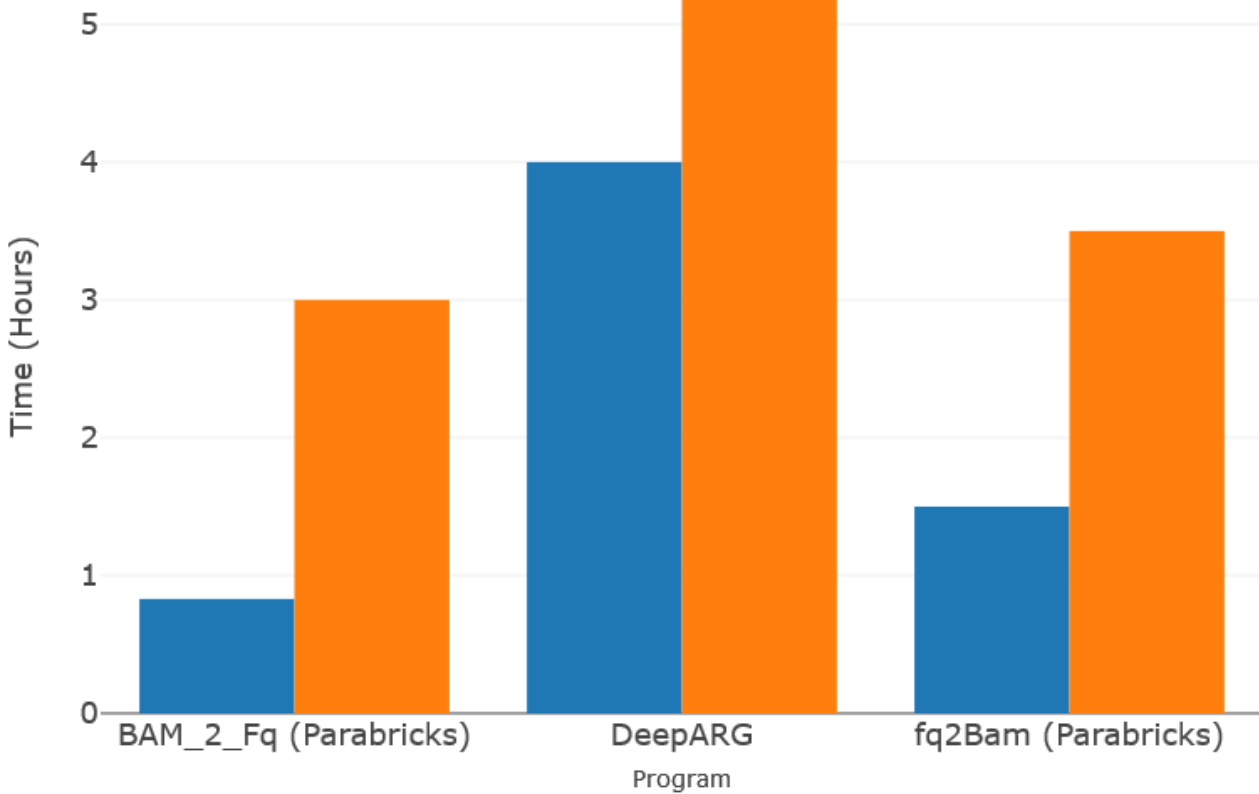
Discover drugs faster with small-molecule virtual screening accelerated by molecular generative AI models



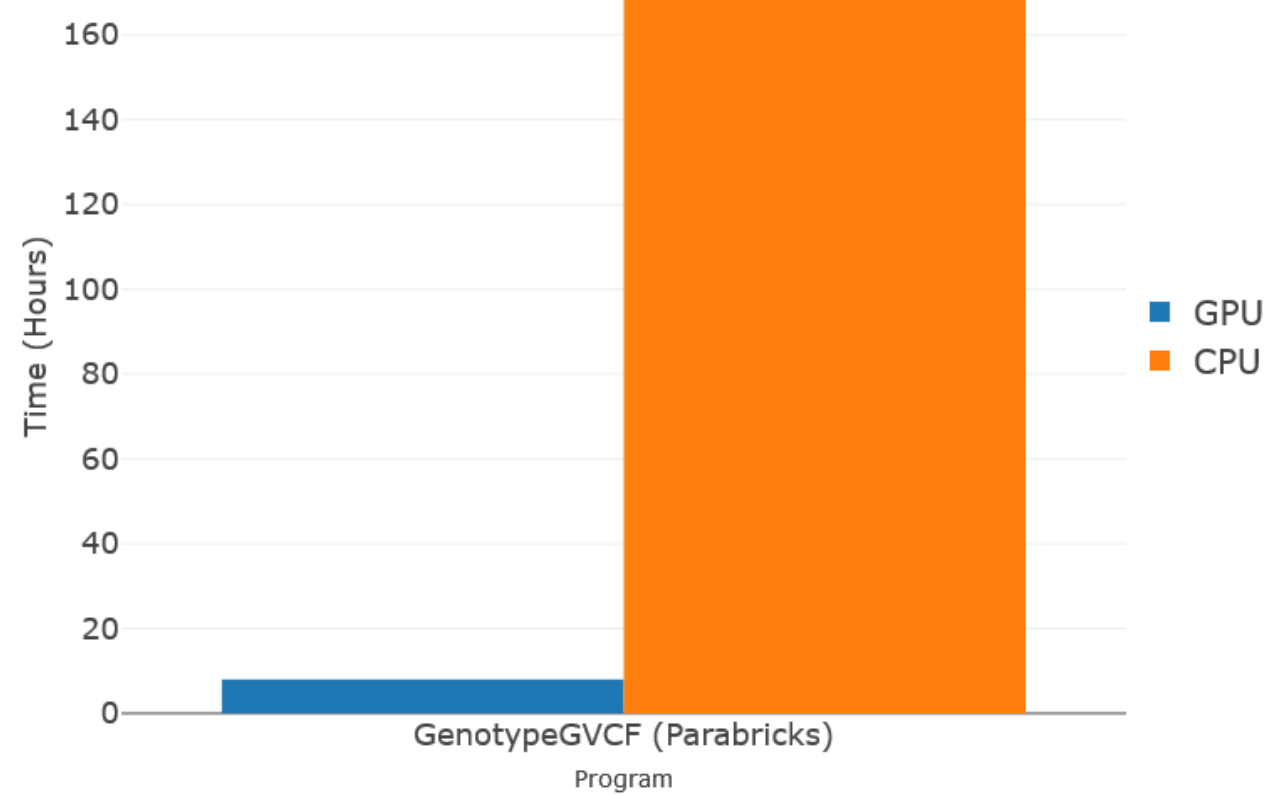
Attēls: <https://resources.nvidia.com/en-us-hc-biopharma/hc-solution-overview-5>

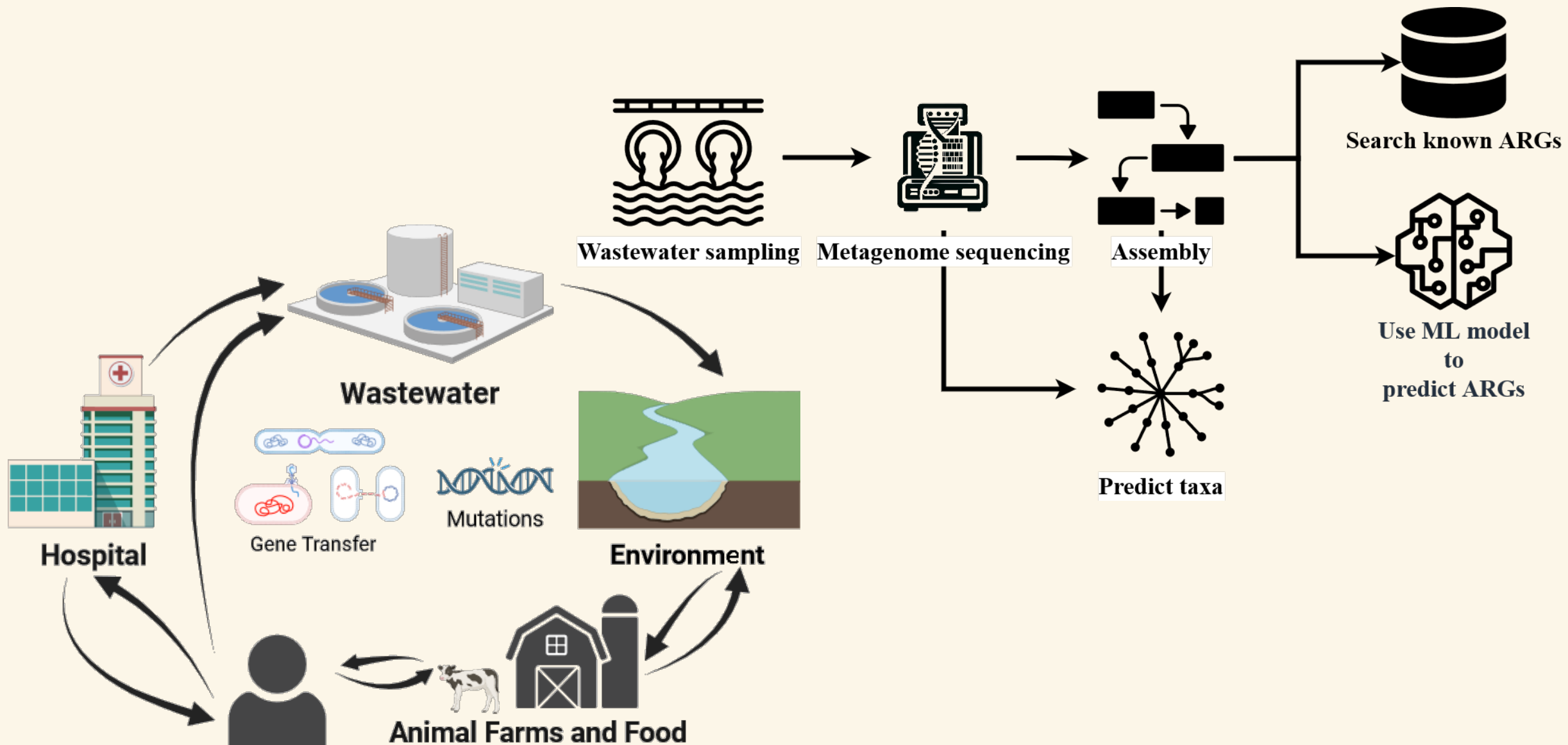
GPU runtime comparison on HPC

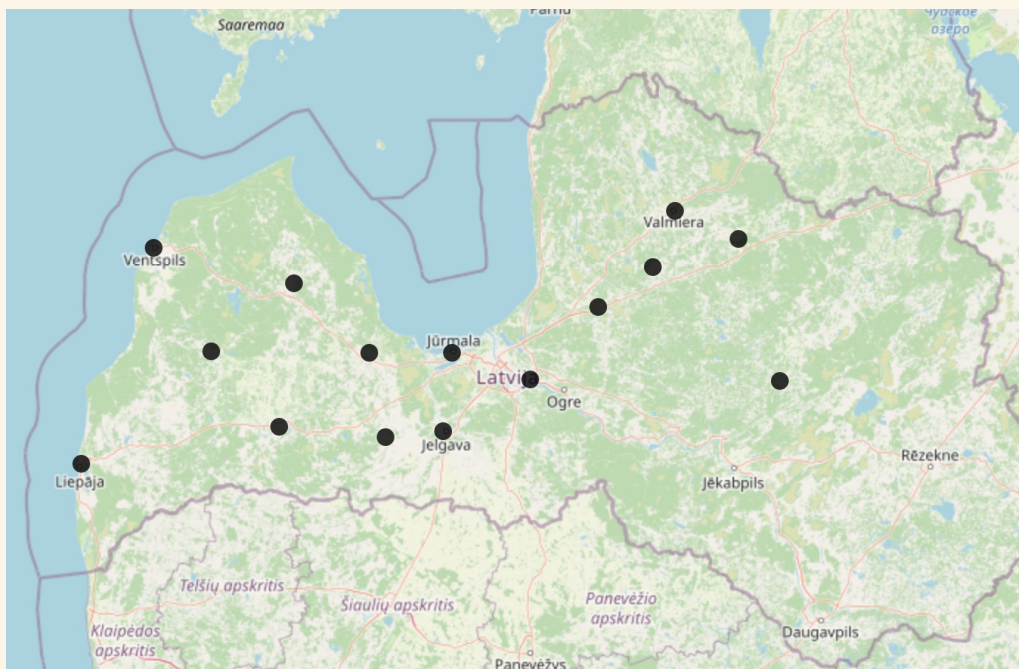
Execution Time for metagenome seq. reads



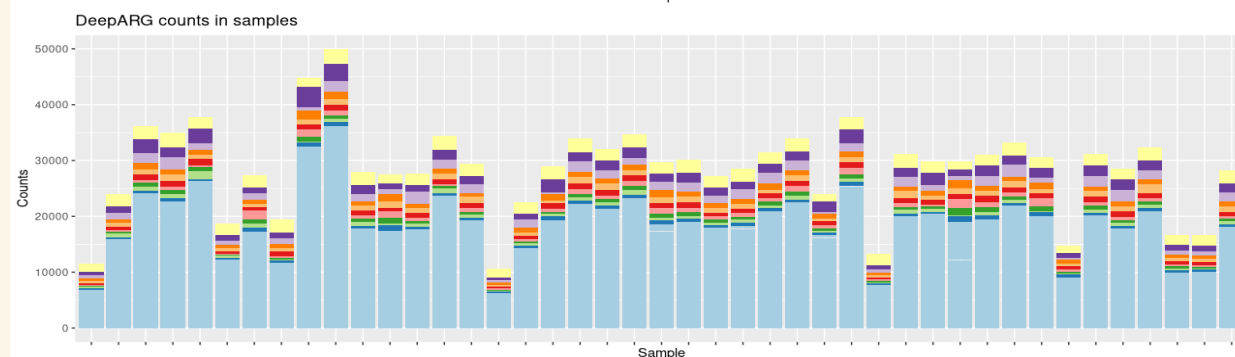
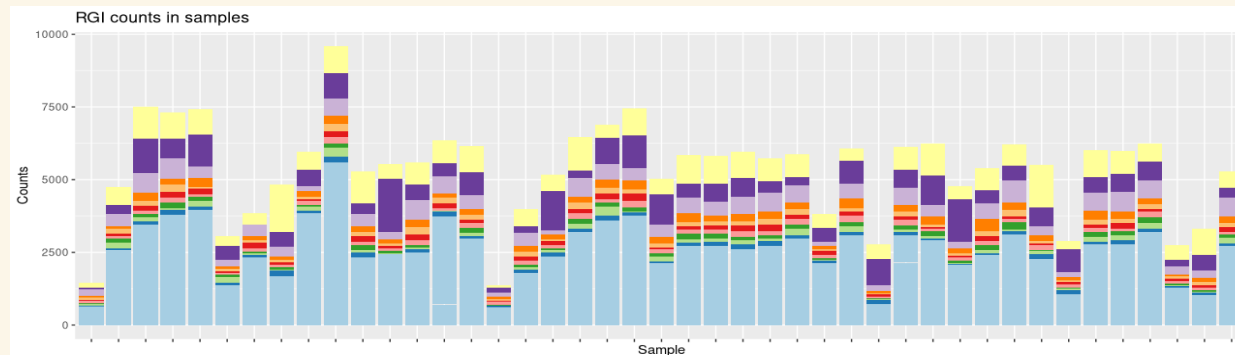
Execution Time for Genotyping full genome



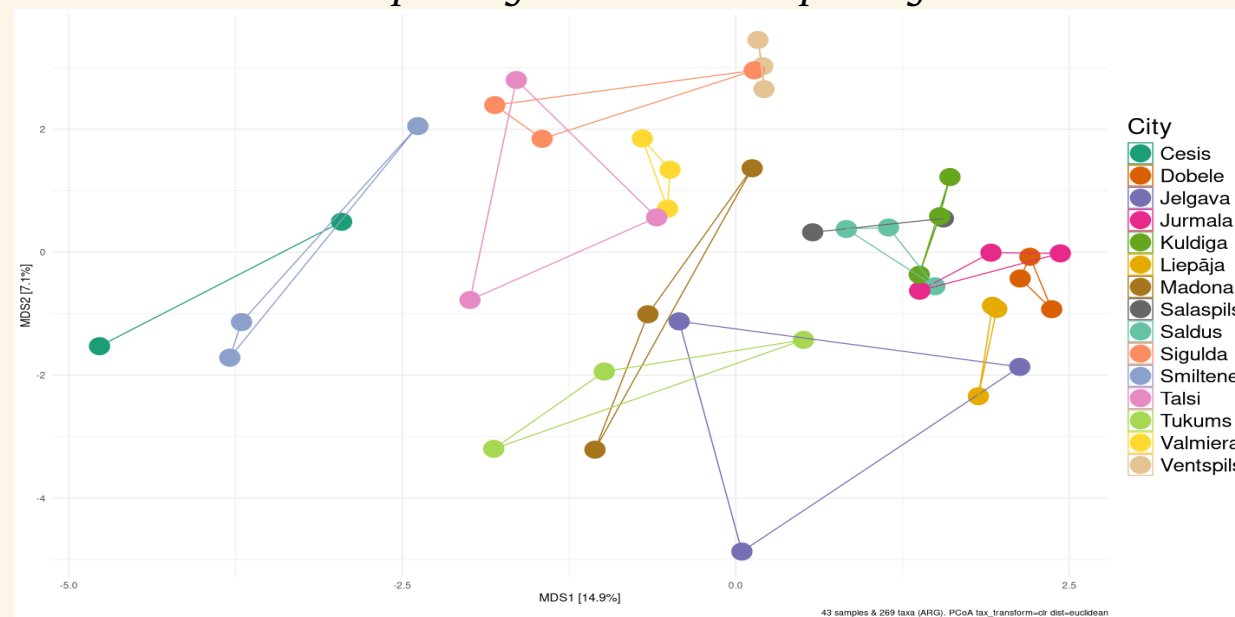




Att.: Paraugu ievākšanas vietas no notekūdeņu stacijām Latvijas pilsētās



Att.: 10 biežāk sastopamo gēnu daudzums paraugos



Att.: Rezistences gēnu beta diversitātes grafiks pa pilsētām

Metagenoma sekvencēšana atklāj antibiotiku rezistentu gēnu (ARG) un baktēriju sastāvu slimnīcas vidē.

Karbapenemāzi producējošas *Klebsiella pneumoniae* izolāti no:

- 6 pacientiem intensīvās terapijas nodālās
- 10 izlietnēm intensīvās terapijas nodālās
- 2 notekūdeņu paraugiem.

Metagenoma paraugi no:

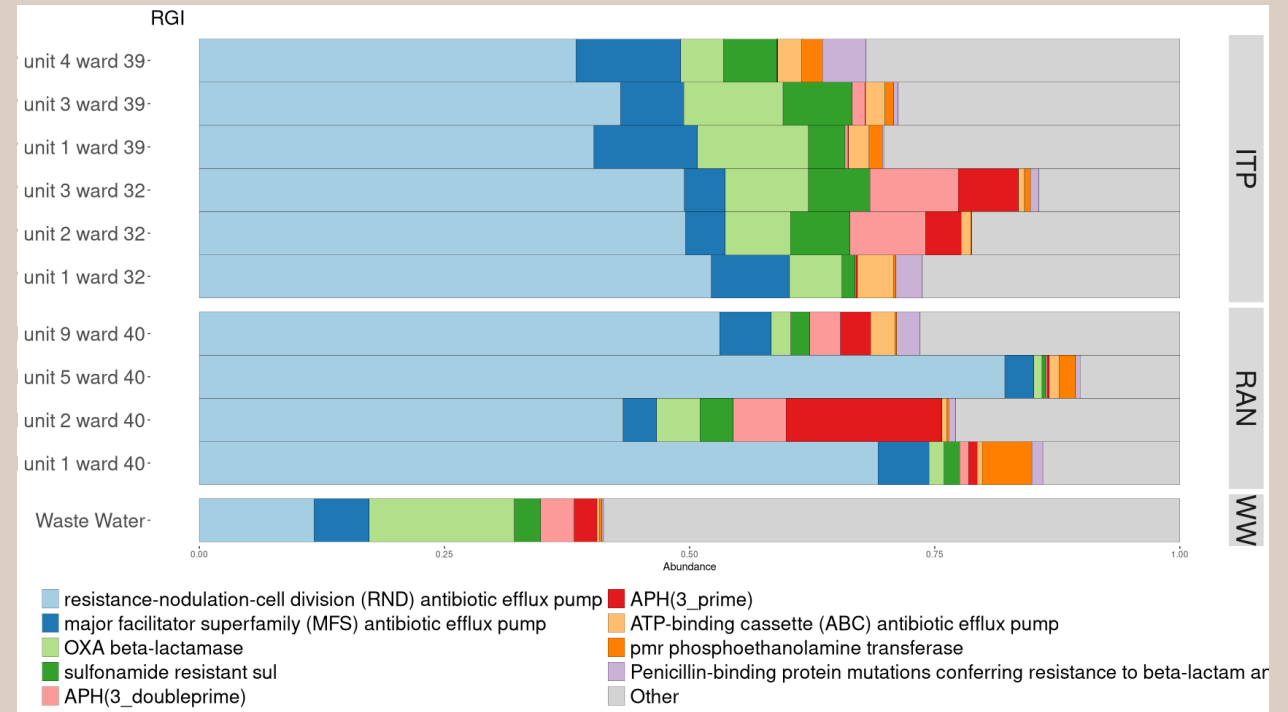
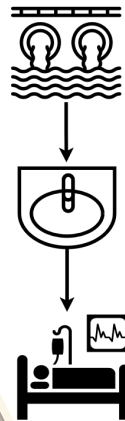
- 10 izlietnēm divās intensīvās terapijas nodālās
- Slimnīcas ēkas notekūdeņu paraugi.

- Kopā identificēti **118 dažādi rezistences gēni**

- **94% no gēniem** pacientu izolātos atrodami arī metagenoma datos.

- No **66 gēniem** pacientu izolātos **24** tika atrasti visos paraugos, **32 vismaz trīs paraugos**.

- **15 karbapenemāzes rezistenci gēni** no kuriem **6 gēni** tika atrasti visos izolētajos paraugos.



Attēls: Biežāk sastopamās antibiotiku rezistences gēnu grupas slimnīcu nodaļās un notekūdeņos

- 1) Noteiktie ARG parāda **visu baktēriju aizsardzības mehānismu kopumu**
- 2) Pacientos atrastajās baktēriju pilno genomu sekvencēs tika atrasti **tie paši gēnu veidi**, kas izlietnēs un notekūdeņos.
- 3) Pacientos atrastajās baktērijās novērojama līdzīgu rezistences mehānismu uzkrāšanās.
- 4) AI rīku pielietošana uzrāda salīdzināmus rezultātus,



Projekta Nr.: VPP-EM-BIOMEDICĪNA-2022/1-0001



Paldies par
uzmanību