

Metafast

High-throughput tool for
metagenome comparison



V. Ulyantsev, S. Kazakov

V. Dubinkina, Tyakht A., Alexeev D.



ITMO UNIVERSITY

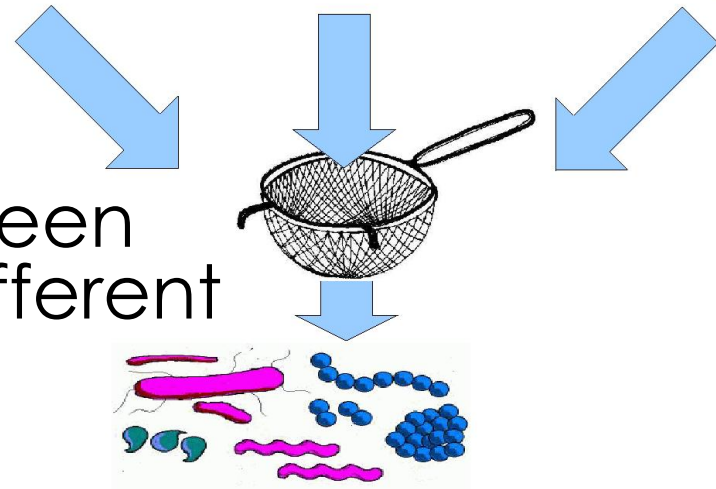
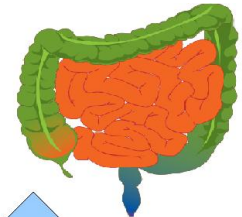
POST



GENOME



Comparative metagenomics

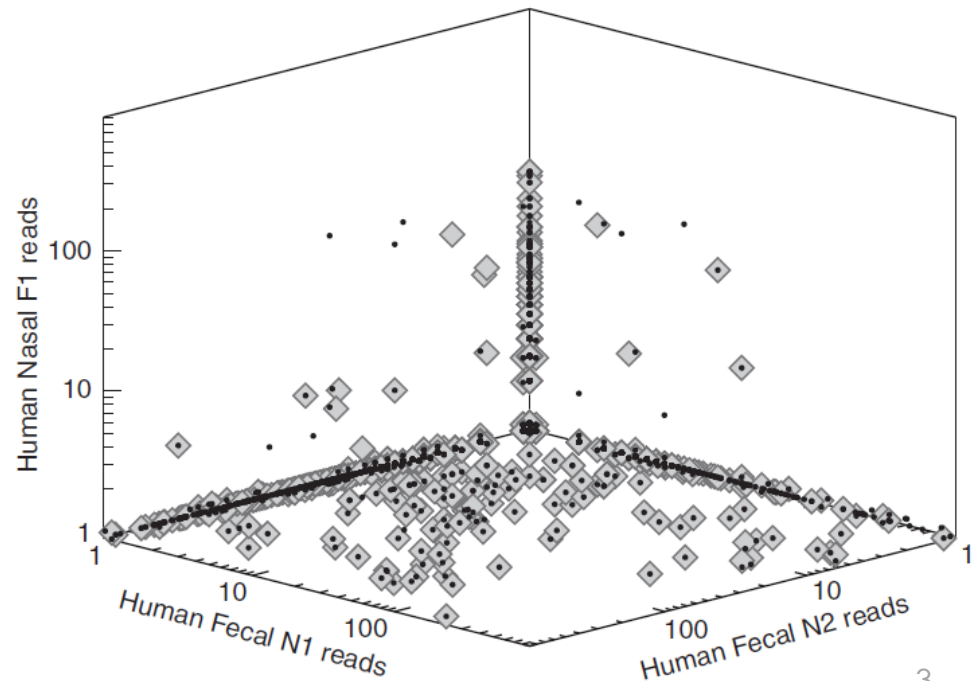


- **Interrelationships** between metagenomes from different samples
 - different biomes
 - different time points
- High level of **unknown** sequences
- Mapping can **limit** the amount of data that can be analyzed

crAss

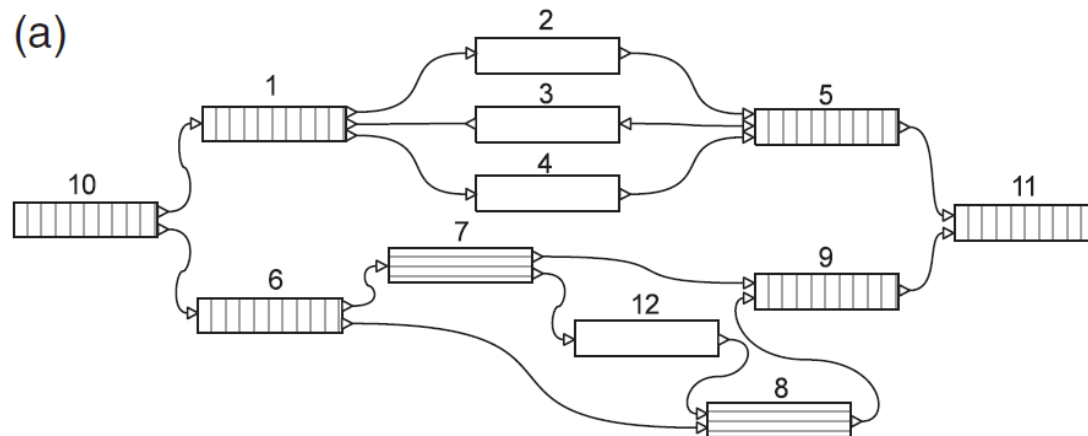
- De novo cross-assembly of all sequence reads
- Number of cross-contigs with reads from both metagenomes

$$d_{i,j} = 1 - \frac{c_{i,j}}{\min\{c_i, c_j\}}$$



MaryGold

- Detect and explore genomic variation between metagenomic sequencing samples
- Detect bubble structures in contig graphs using graph decomposition
- 454 and Illumina data



Challenges

- Reference-based (mapping sequences)
 - High level of unknown sequences
- Assembly-based
 - Slow on large datasets

Solution: fast but low quality “semi-assembly” for every library

New algorithm – MetaFast

Method, based on simple assembly:

A. For each library:

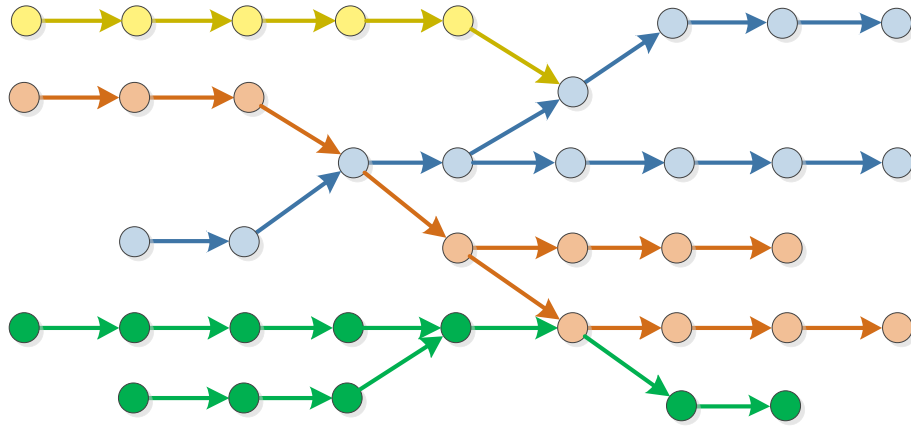
- Construct de Bruijn graph
- Extract simple paths, not contigs

B. For all libraries:

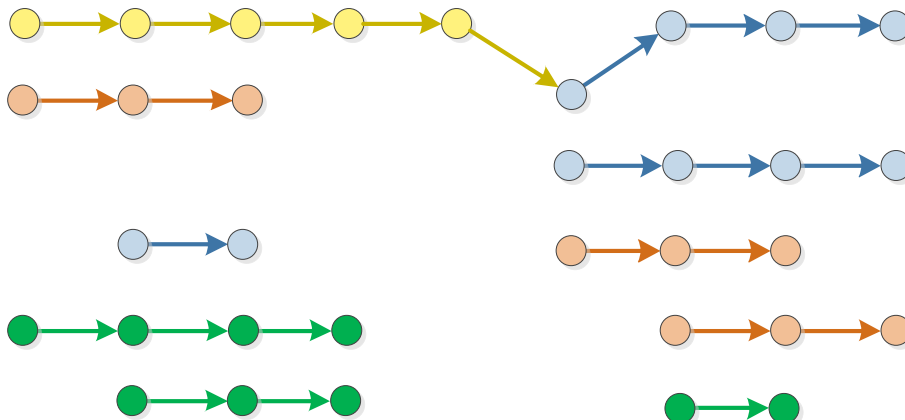
- Construct de Bruijn graph from found paths
- Extract components

C. Calculate characteristic vectors for libraries.

2. Extract simple paths



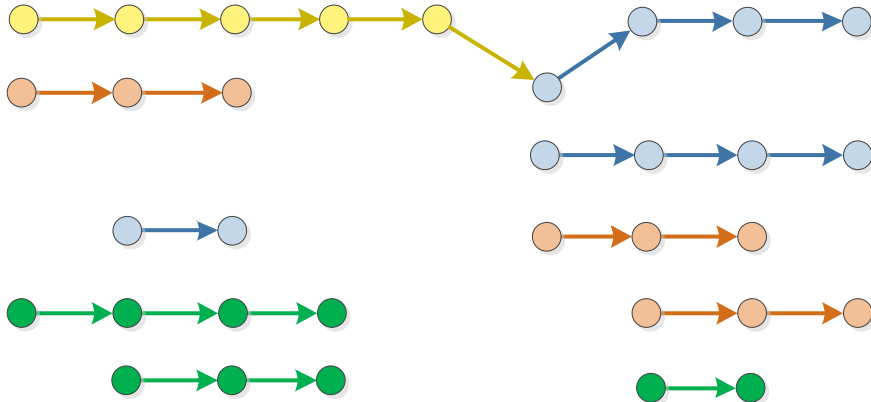
de Bruijn graph



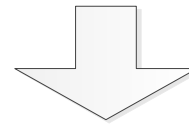
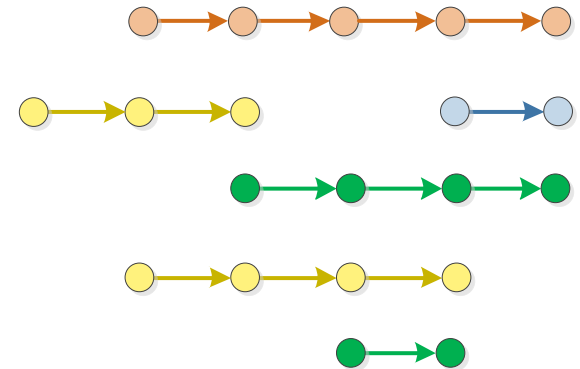
Simple paths

3. Merge paths

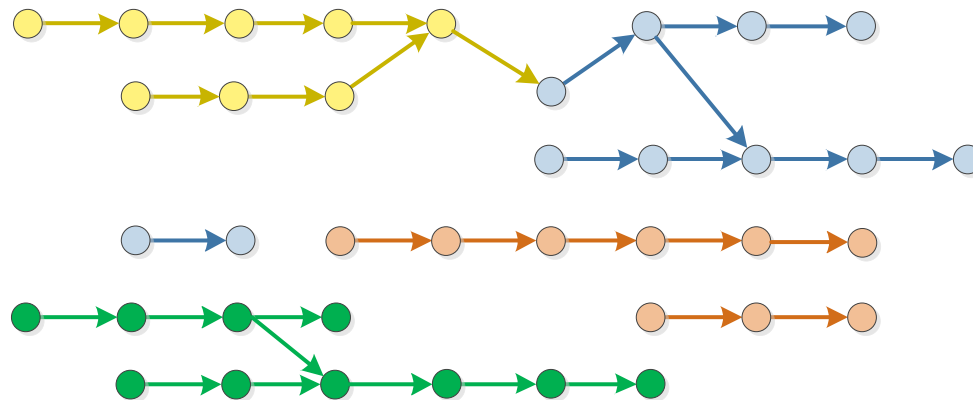
Paths 1



Paths 2

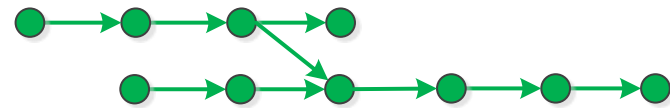


Paths combined in single de Bruijn graph

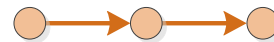
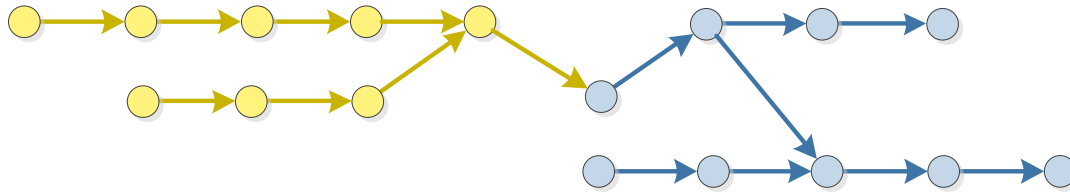


Component

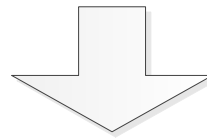
- K-mers set
- $B_1 \leq \text{size} \leq B_2$
- Big component?
 - Iterative algorithm for decomposition



5. Construct characteristic vectors



Components



Library 1

(15, 0, 6)

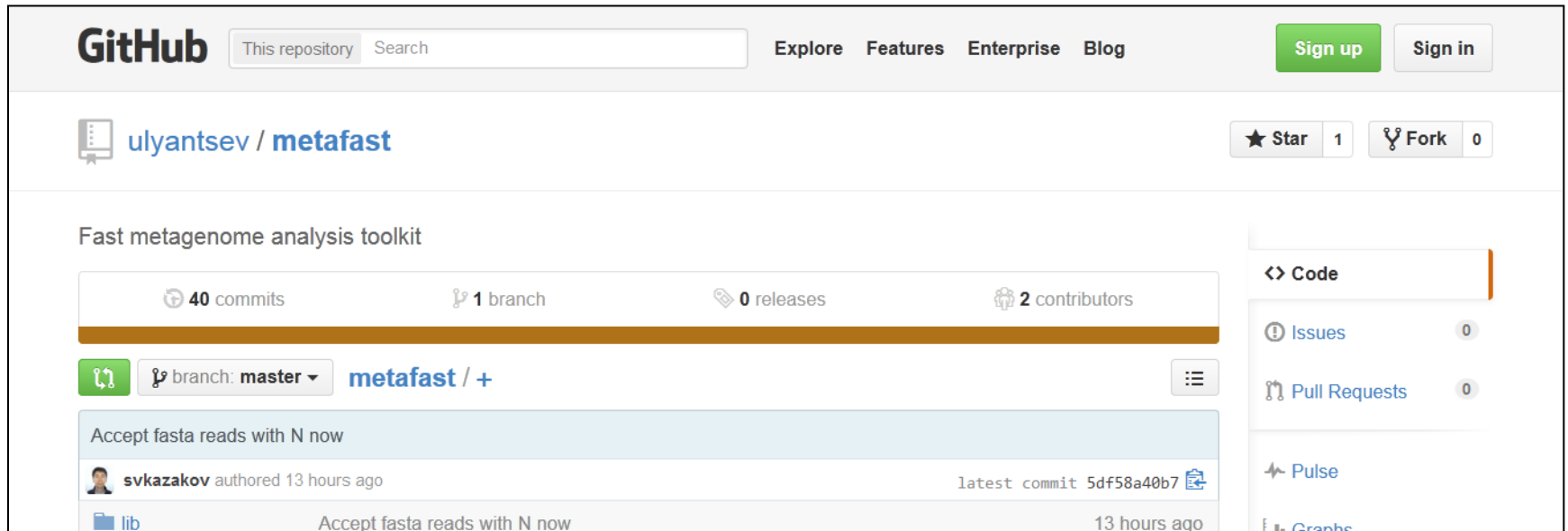
Library 2

(0, 7, 8)

Characteristic
vectors

Implementation

- On *Java*
- Open source project
- <http://github.com/ulyantsev/metafast>

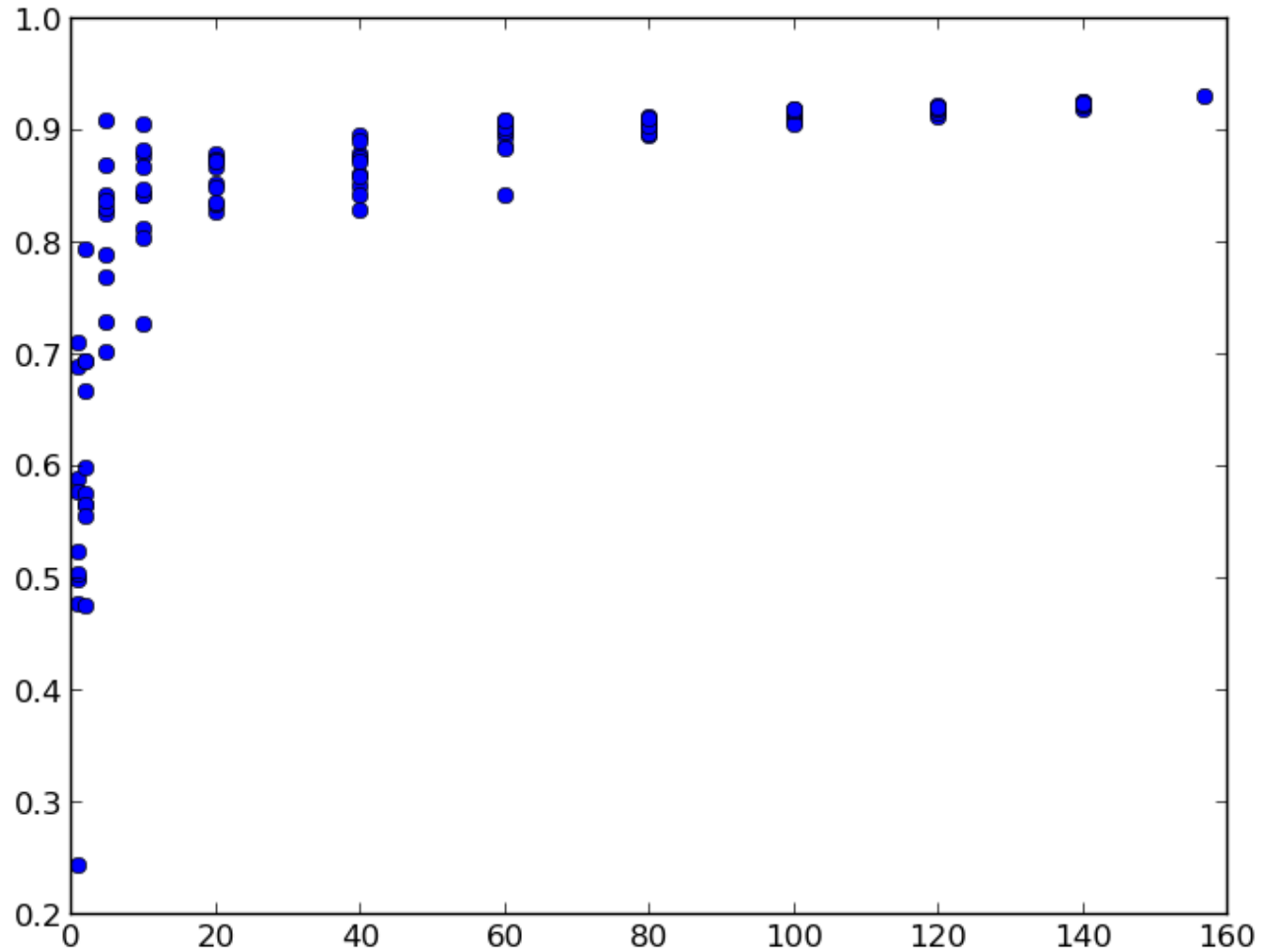


The screenshot shows the GitHub interface for the repository 'ulyantsev / metafast'. The repository is described as a 'Fast metagenome analysis toolkit'. It has 40 commits, 1 branch, 0 releases, and 2 contributors. The current branch is 'master'. A recent commit by 'svkazakov' is shown, titled 'Accept fasta reads with N now', made 13 hours ago. The commit hash is '5df58a40b7'. The repository has 1 star and 0 forks. The right sidebar shows options for 'Code', 'Issues', 'Pull Requests', 'Pulse', and 'Graphs'.

Experiments

- **157** Chinese gut metagenomes (600 Gb)
- **93 %** – correlation between distance matrices based on mapping to known references and our vectors
- About 10 hours cluster time (not full-loaded)

80 libraries is enough



Results & future work

- MetaFast – new approach and cross-platform tool for comparative metagenomics
- Promising initial experiments
- Experiments with simulated data
- New information about existing metagenomes
- Algorithm modifications

Thank you for attention!

<http://github.com/ulyantsev/metafast>

V. Ulyantsev
S. Kazakov

V. Dubinkina
A. Tyakht
D. Alexeev

