

25 июня 2018г



Новые горизонты алгоритмической биоинформатики

Алла Львовна Лapidус

Professor, Department of Cytology and Histology, SPbSU
Deputy director, Centre for Algorithmic Biotechnology, SPbSU

О чем



1. Центр алгоритмической биотехнологии - немного сухой статистики
2. Что такое биоинформатика?
3. Наука в Центре алгоритмической биотехнологии
4. Образование

О чем



1. Центр алгоритмической биотехнологии - немного сухой статистики
2. Что такое биоинформатика?
3. Наука в Центре алгоритмической биотехнологии
4. Образование



В СПбГУ с конца 2014 года - в рамках открытого конкурса СПбГУ на создание исследовательских лабораторий под руководством ведущих ученых



«Центр алгоритмической биотехнологии» (ЦАБ)

Руководитель лаборатории – П.А. Певзнер, профессор Калифорнийского университета в Сан-Диего (University of California San Diego, USA), к.ф-м.н.

Заместитель руководителя – А.Л.Лapidус, профессор кафедры цитологии и гистологии СПбГУ, к.б.н.

«Центр алгоритмической биотехнологии» (ЦАБ)

Численность лаборатории

2015 год - 14 человек, из которых 2 кандидата наук

2018 год - 20 человек, из которых 4 кандидата наук + 1 работа принята к защите на получение степени СПбГУ

Средний возраст - 27 лет

«Центр алгоритмической биотехнологии» (ЦАБ)

Научные партнеры (**Collaborators**)

- ✓ **Yale University**
 - ✓ **Scripps Institution, UCSD**
 - ✓ **University of British Columbia**
 - ✓ **The Wellcome Trust Sanger Institute**
 - ✓ **University of Hong Kong**
 - ✓ **John Hopkins University**
 - ✓ **UCSD**
 - ✓ **JCVI**
 - ✓ **Joint Genome Institute (JGI, LBNL)**
 - ✓ **EMBL-EBI**
 - ✓ **Brain and Mind Research Institute,
(Weill Cornell Medicine, USA)**
 - ✓ **Institute for Agricultural and Forest Systems
in the Mediterranean (Italy)**
- Juno Therapeutics***
- Genentech***
- AstraZeneca Oncology***

«Центр алгоритмической биотехнологии» (ЦАБ)

Научные партнеры (Collaborators)

St Petersburg:

- ✓ А.С. Готов, СПбГУ, ИТБМ
- ✓ Центр Добржанского, СПбГУ
- ✓ А.Н. Суворов, Институт экспериментальной медицины
- ✓ ИТМО

Moscow

Dmitry Chudakov and Mikhail Shugay, Institute of Bioorganic Chemistry

Novosibirsk

Институт цитологии и генетики СО РАН

«Центр алгоритмической биотехнологии»

(ЦАБ)

2015-2018 год

Публикации – 53:

Nature family – 5; PNAS – 3; Genome Biology – 2; Genome research – 2


Доклады на международных конференциях – 62

Преимущественно приглашенные


РИД

Государственную регистрацию прошли 5 программных продуктов
и подготовлены к регистрации еще 2

Количество цитирований публикаций в системе Google Scholar

10.  **Alla Lapidus** Цитируется: 22 464
St. Petersburg State University, Санкт Петербургский государственный университет
Подтвержден адрес электронной почты в домене spbu.ru
bioinformatics genomics MOOC

Все позиции с 1-ой по 9-ую и с 11 по 20-ую занимают мужчины...

20.  **Alexey Gurevich** Цитируется: 5 229
St. Petersburg State University, Center for Algorithmic Biotechnology
Подтвержден адрес электронной почты в домене spbu.ru
Bioinformatics Computational Biology Genome Assembly Metabolomics

Премии за высокую цитируемость



О чем

1. Центр алгоритмической биотехнологии - немного сухой статистики
2. Что такое биоинформатика?
3. Наука в Центре алгоритмической биотехнологии ,,,"
4. Образование

Term Bioinformatics

Термин **Биоинформатика** был введен в 1970-ом году голландскими биологами Паулиеной Хогвег и Беном Хеспере, изучавшими свойства сложных динамических биологических систем

Что такое биоинформатика?

Компьютеры в биологии

Программирование

Математика

BIG data

Анализ данных

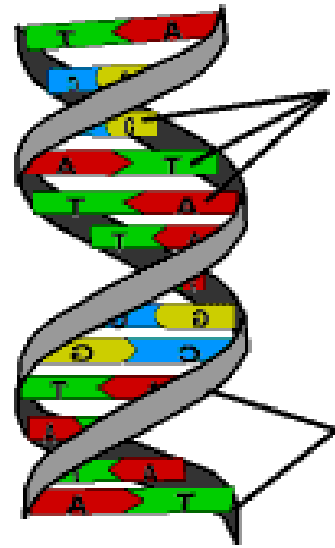
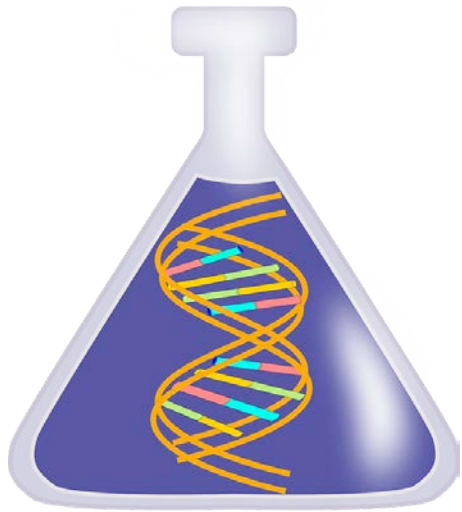
Осмысление данных

Применение

Биолог: ДНК

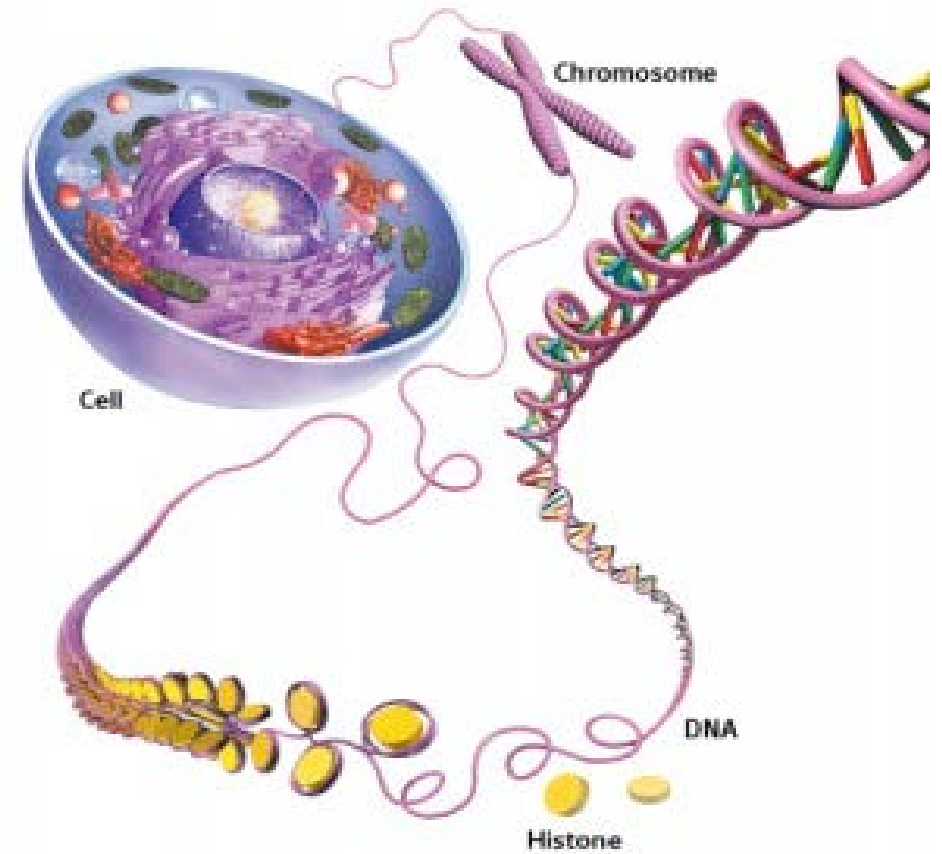


www.shutterstock.com · 54489088

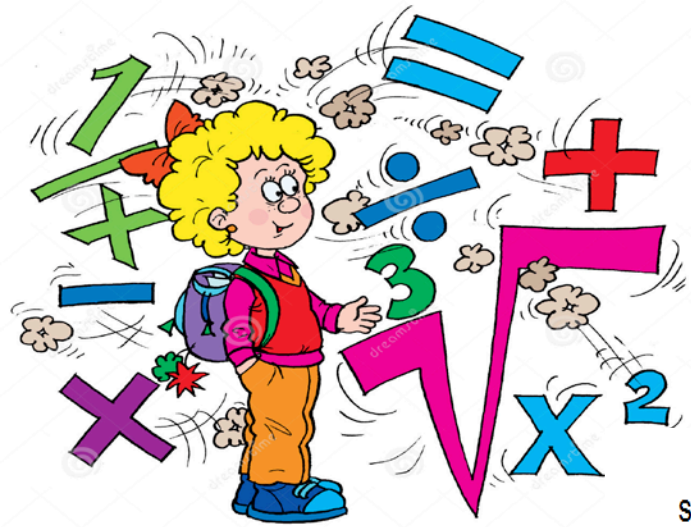


bases -
varying
part of
molecule

side units -
do not
vary



Математик: ДНК



A T G C

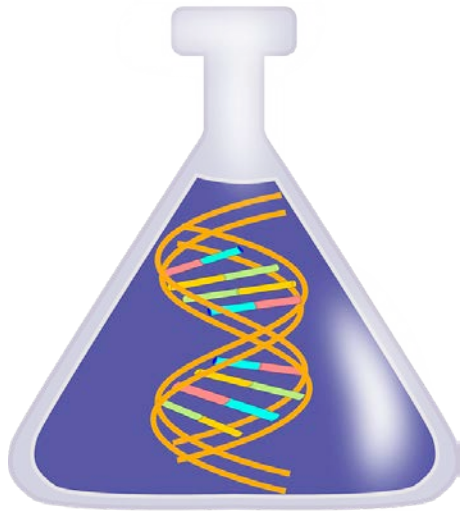
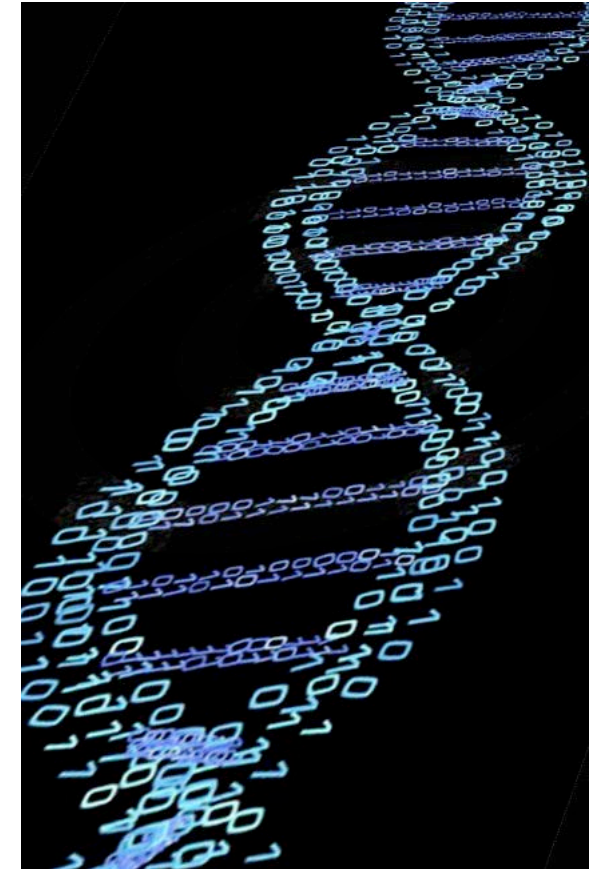
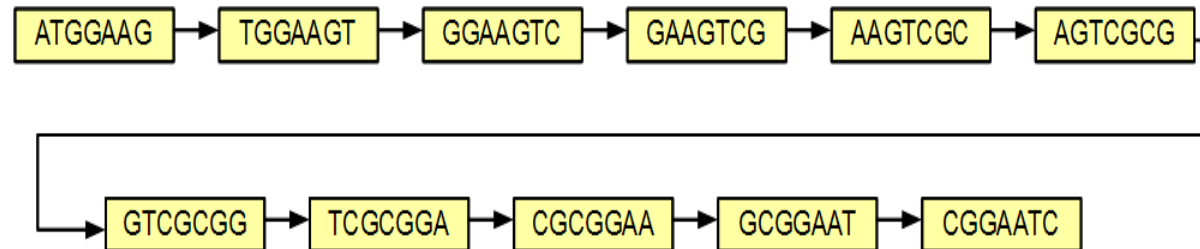
sequence

ATGGAAGTCGCGGAATC

7mers

ATGGAAG
TGGAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

de Bruijn graph



Кто такой биоинформатик?

Biologist

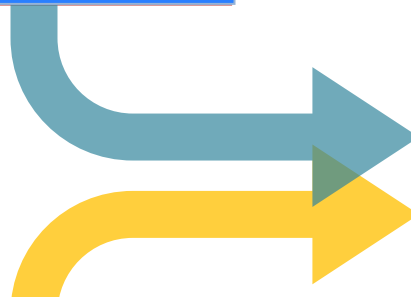
collects molecular data:
DNA, RNA & Protein sequences,
gene expression, etc.

Computer scientist

(+Mathematicians, Statisticians, etc.)
Develops tools, pieces of software and
algorithms to store and analyze data.

Bioinformatician

Asks and answers biological questions
by analyzing molecular data



Покопались в данных и



1,003 reference genomes of bacterial and archaeal isolates expand coverage of the tree of life

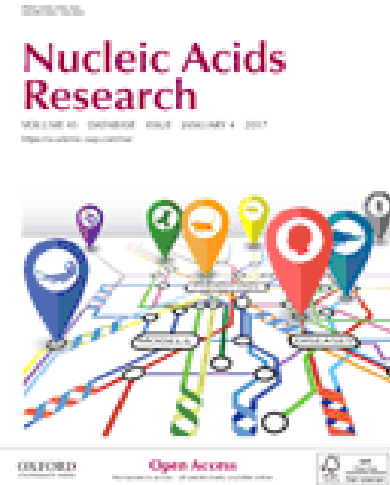
Nature Biotechnology 35, 676–683 (2017) doi:10.1038/nbt.3886



Evolutionary drivers of thermoadaptation in enzyme catalysis.

614 protein families with currently unknown structures

Science 20 Jan 2017: Vol. 355, Issue 6322, pp. 289-294
DOI: 10.1126/science.aah3717



IMG-ABC: new features for bacterial secondary metabolism analysis and targeted biosynthetic gene cluster discovery in thousands of microbial genomes

Nucleic Acids Res (2017) 45 (D1): D560-D565. DOI: <https://doi.org/10.1093/nar/gkw1103>

Genomic encyclopedia of bacteria and archaea: sequencing a **myriad of type strains.**

PLoS Biol. 2014 Aug 5;12(8):e1001920. doi: 10.1371/journal.pbio.1001920.



Биоинформатика – это НАУКА

1-й закон биоинформатики

GARBAGE IN



GARBAGE OUT



О чем



1. Центр алгоритмической биотехнологии - немного сухой статистики
2. Что такое биоинформатика?
3. Наука в Центре алгоритмической биотехнологии
4. Образование

«Центр алгоритмической биотехнологии» (ЦАБ)

Основные научные интересы – алгоритмическая биоинформатика

в области:

геномной сборки

метагеномики, транскриптомики, иммуногеномики,

поиска новых антибиотиков

Визитная карточка лаборатории №1



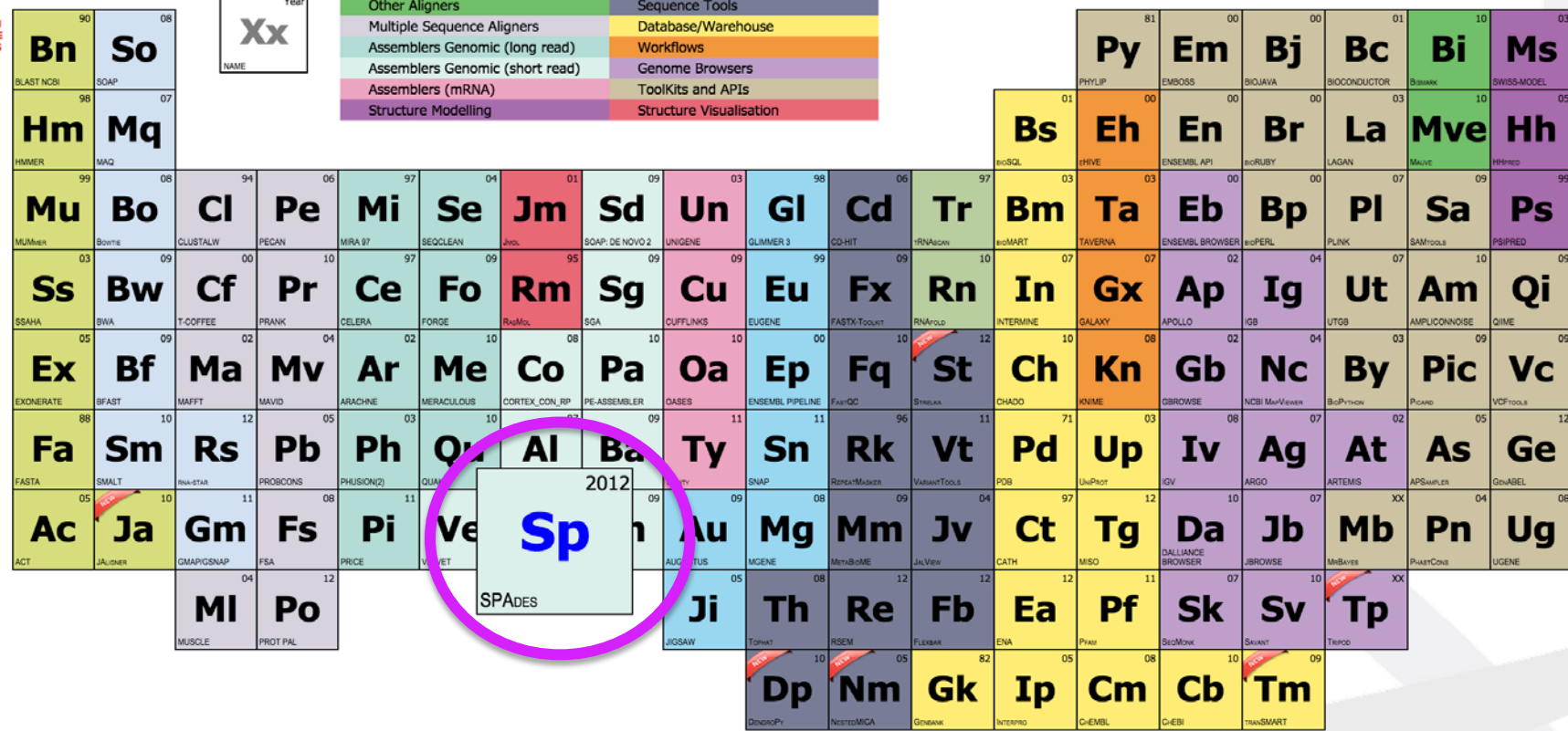
The Elements of Bioinformatics

Search by name:
 Filter by year: 1970 2014

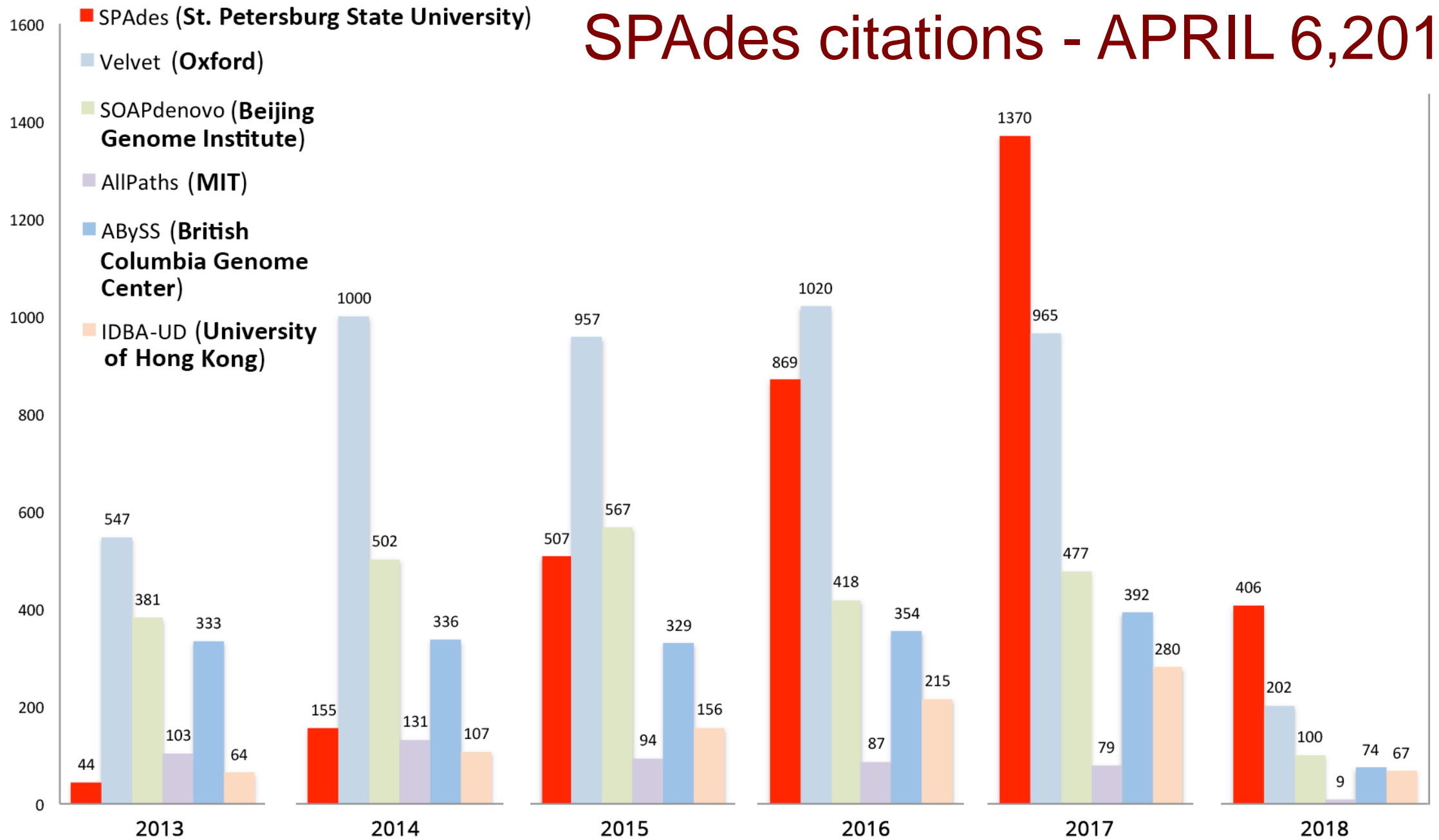
KEY TO TOOL TYPE (stripes indicate new update)

Aligners (pairwise)	Gene Prediction (mRNA)
Aligners (short read)	Gene Prediction (ncRNA)
Other Aligners	Sequence Tools
Multiple Sequence Aligners	Database/Warehouse
Assemblers Genomic (long read)	Workflows
Assemblers Genomic (short read)	Genome Browsers
Assemblers (mRNA)	ToolKits and APIs
Structure Modelling	Structure Visualisation

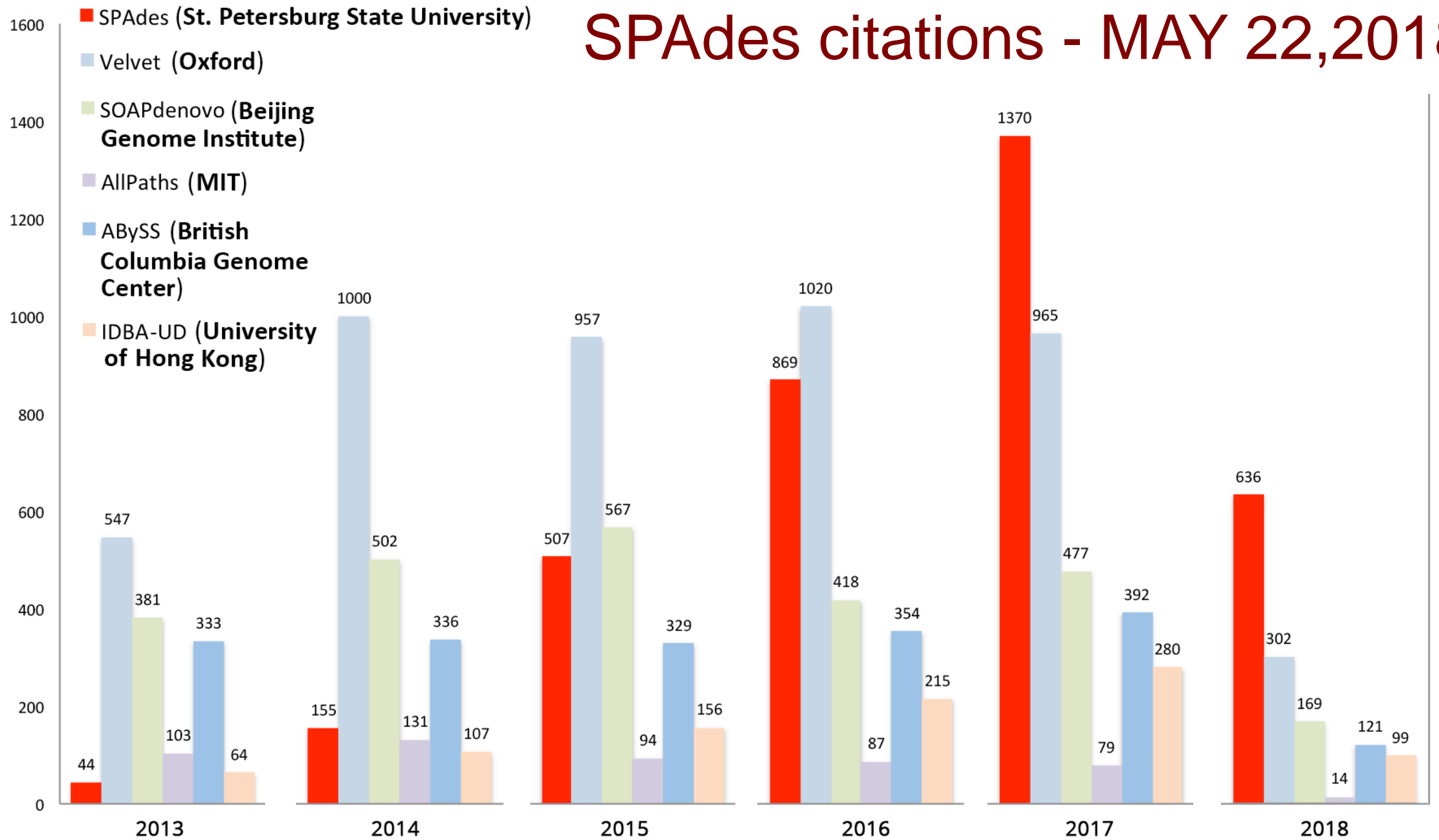
OPEN SOURCE TOOLS



SPAdes citations - APRIL 6, 2018



SPAdes citations - MAY 22, 2018



2. DNA fragmentation

1. DNA isolation

3. Library preparation

4. Sequencing run

6. Genome reconstruction (assembly)



5. Sequencing reads

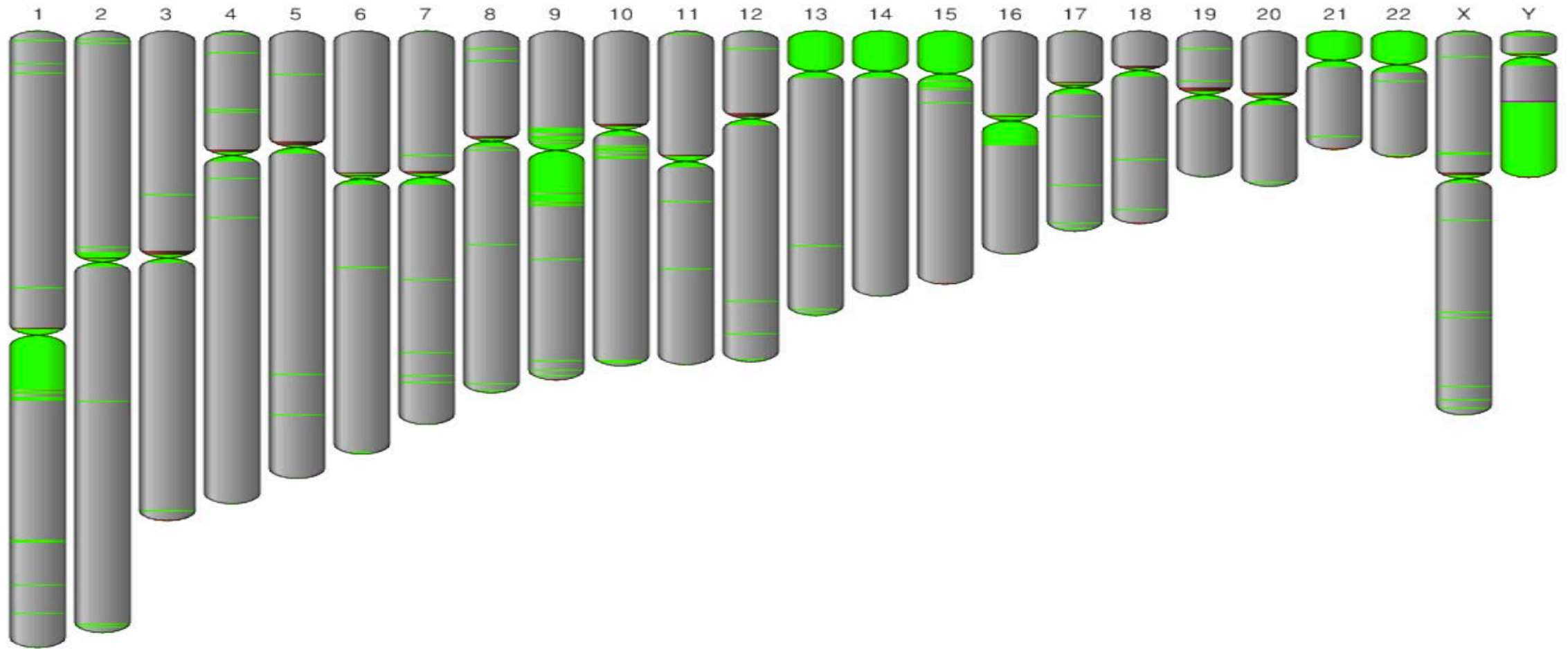
Why genome assembly and why it is hard...



Why genome assembly and why it is hard...



Human genome



■ Gaps

Human genome Project - not the genome! - is finished



Визитная карточка лаборатории №2



QUAST - Genome assembly evaluation tool

S.aureus dataset, reference-based evaluation, contig alignment viewer:



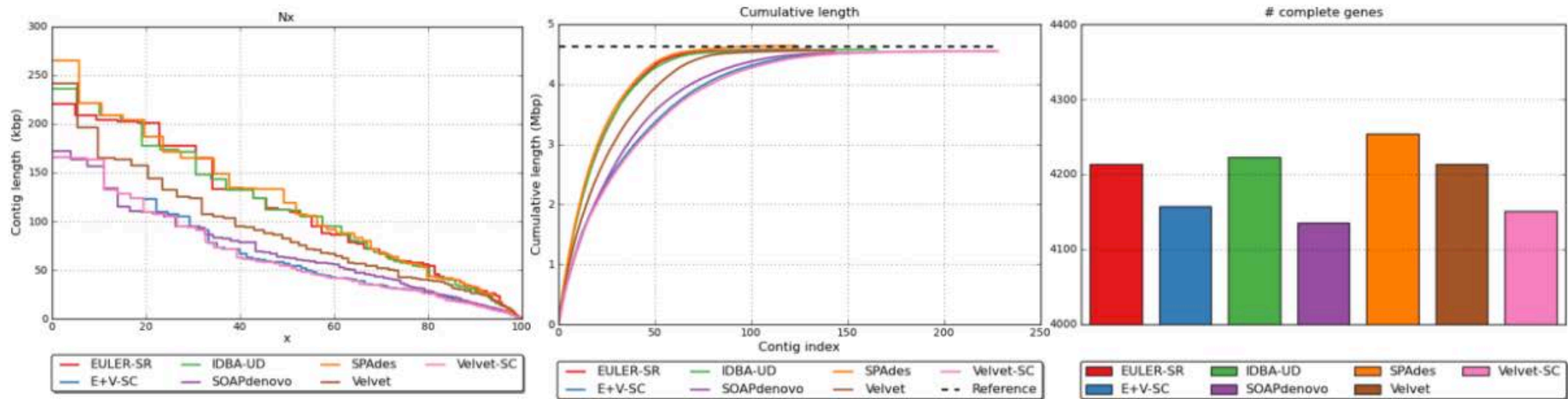


Визитная карточка лаборатории №2

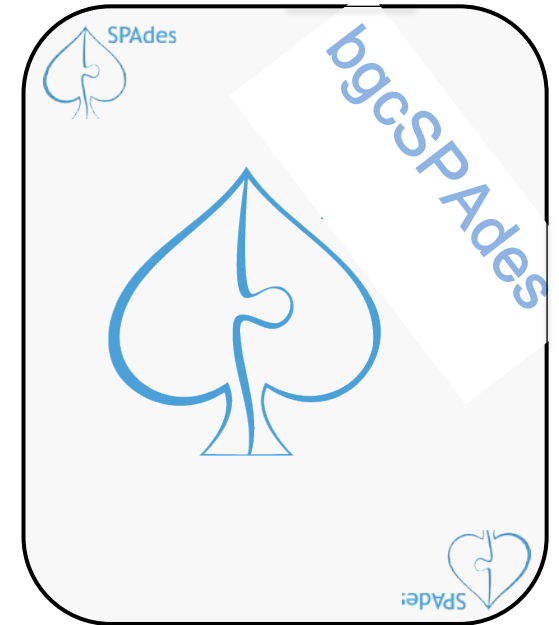
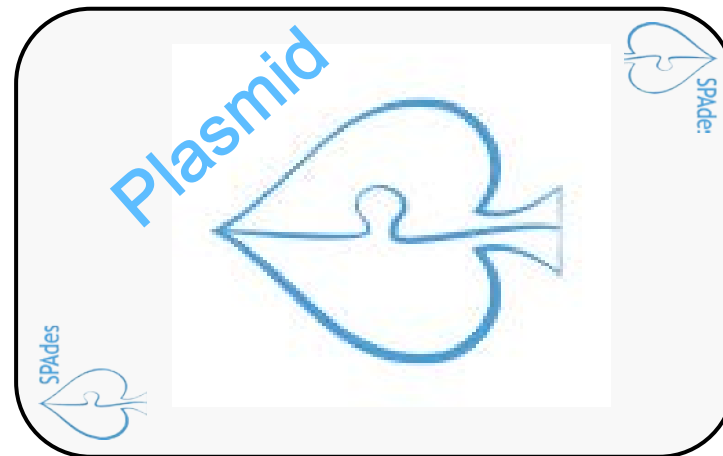
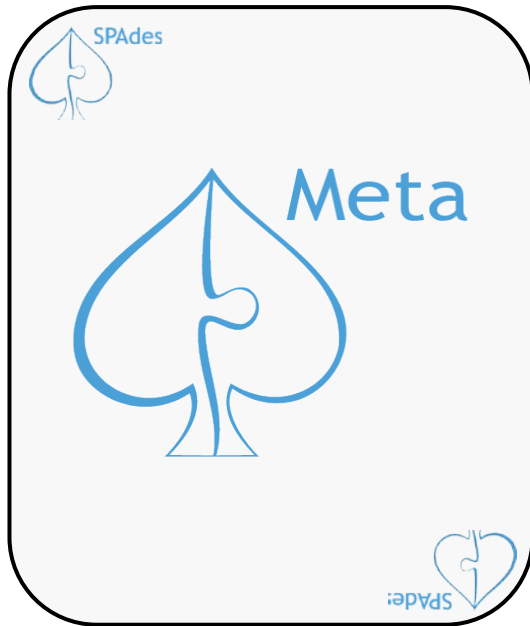
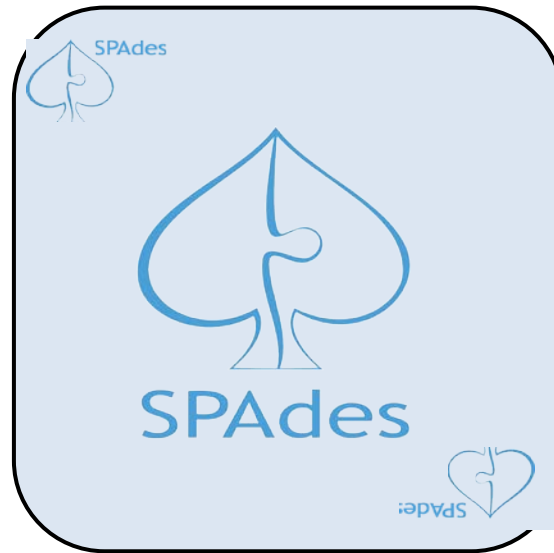
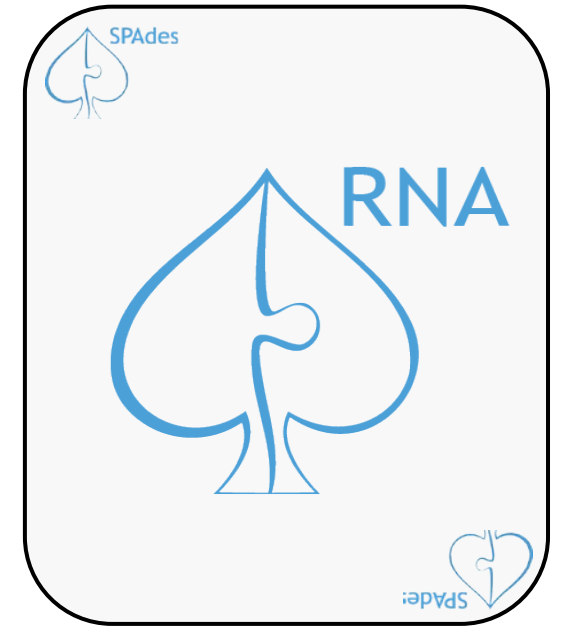
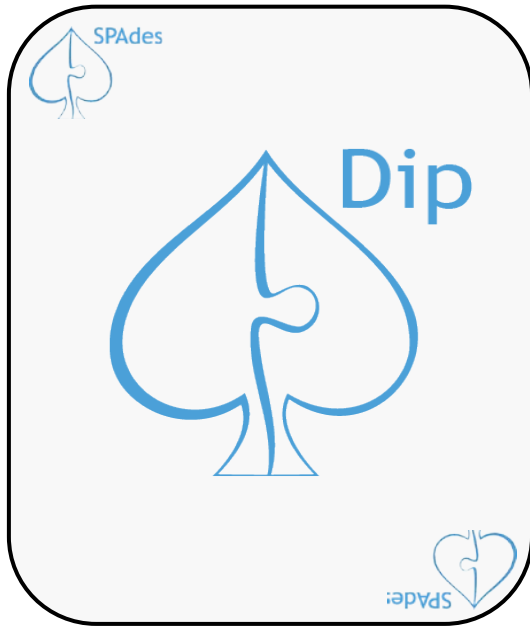


QUAST - Genome assembly evaluation tool

Examples of QUAST output



SPAdes Toolbox



Etc...

Antibiotics Discovery

Genome Assembly

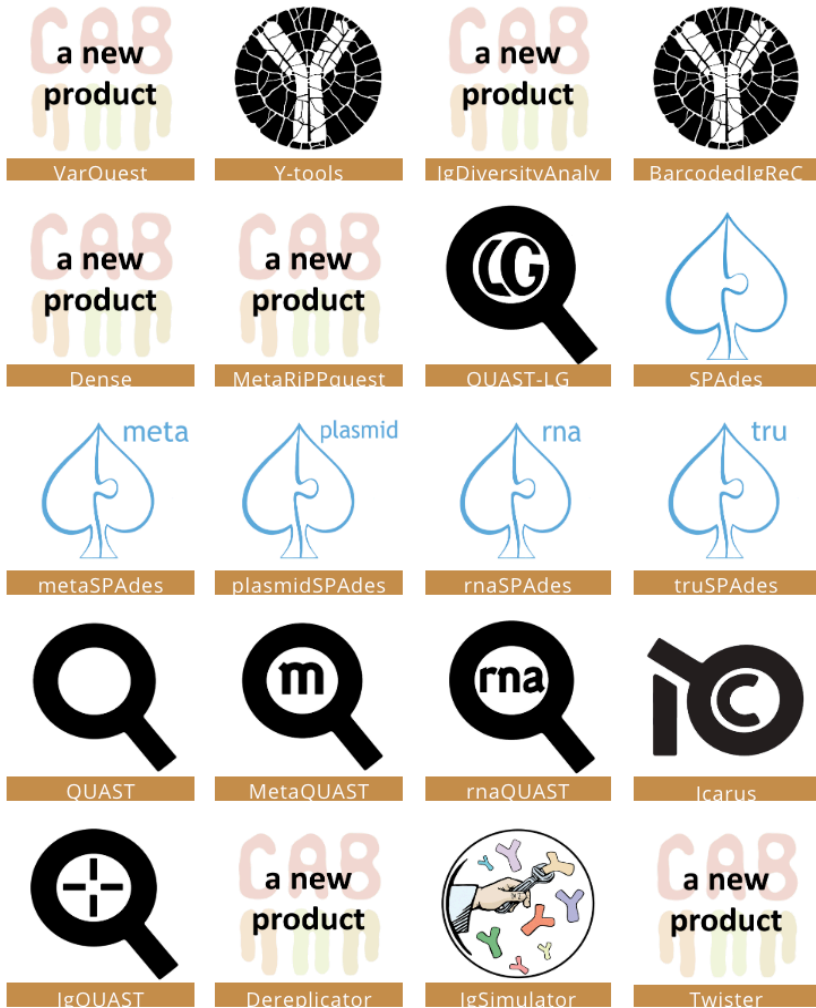
Transcriptomics

TSLR Analysis

Metagenomics

Immunoinformatics

Proteomics



Antibiotics discovery

Immunoinformatics

Education

<http://cab.spbu.ru/our-software/>

Antibiotics discovery

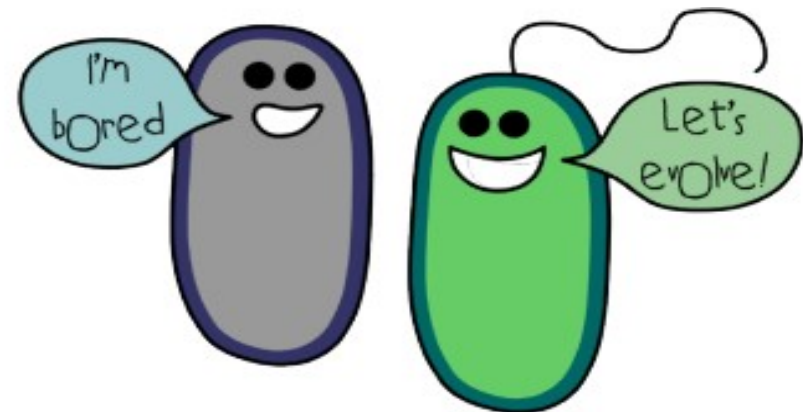
- Antibiotics — natural substances that kill bacteria or suppress their reproduction
- People use antibiotics more than they should (even for viruses)
- Some bacteria may survive



Antibiotics discovery

They evolve and become resistant

- Drug resistance is a major concern worldwide!
- People need new-generation drugs



how multi-drug resistant
bacteria came to be....

Antibiotics discovery challenge

How to discover potential antibiotic among millions of proteins?



Metabolo-genomics

the emerging approach that combines **genomic** and **proteomic** data to identify biologically active organic compounds that could be used as a natural source for novel antibiotics and other drugs



Metabolo-genomics

- DNA sequencing technologies in **genomics**
- high-throughput mass spectrometry-based analysis that includes database search of **mass spectra** against already known metabolites
- **algorithms** and **bioinformatics** tools to expedite research

Metabolo-genomics

- DNA sequencing technologies in **genomics**
- high-throughput mass spectrometry-based analysis that includes database search of **mass spectra** against already known metabolites
- **algorithms** and **bioinformatics** tools to expedite research

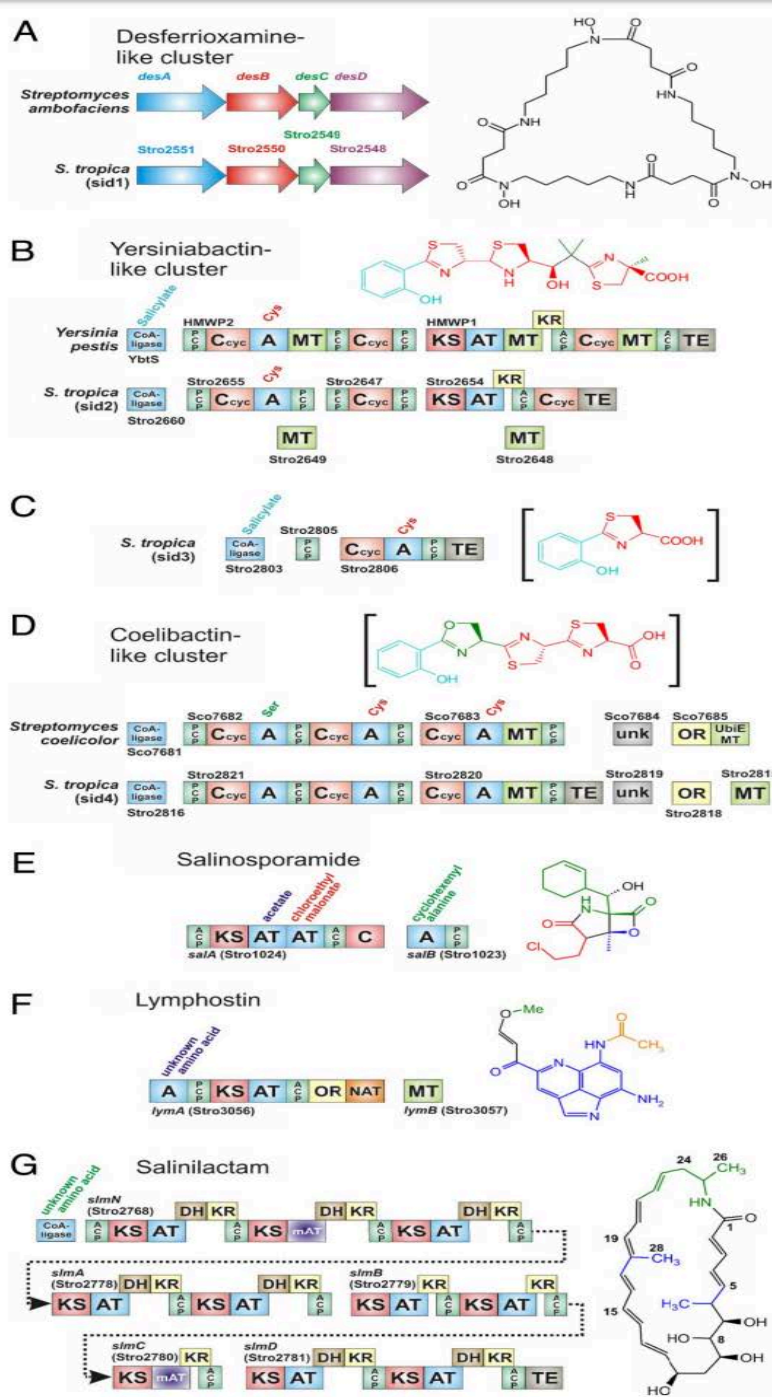


Fig. 2. Selected genes from *S. tropica* modular biosynthetic enzyme systems

Genome sequencing reveals complex secondary metabolome in the marine actinomycete *Salinispora tropica*

Daniel W. Udway*, Lisa Zeigler*, Ratnakar N. Asolkar*, Vasanth Singan†, Alla Lapidus†, William Fenical*, Paul R. Jensen*, and Bradley S. Moore*^{‡§}

*Scripps Institution of Oceanography and †Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California at San Diego, La Jolla, CA 92093-0204; and ‡Department of Energy, Joint Genome Institute–Lawrence Berkeley National Laboratory, Walnut Creek, CA 94598

Актиномицеты - богатейший источник вторичных метаболитов, на которые приходится более половины всех антибиотиков, обнаруженных на сегодняшний день





Streptomyces coelicolor

“Antibiotics factory”

29 biosynthetic gene clusters, including 3 NRPS
(nonribosomal peptide synthetase)

- Coelichelin: 3 A-domains Assembled with SPAdes in 1 contig in minutes
- Nogalamycin: 2 A-domains Assembled with SPAdes in 1 contig in minutes
- Calcium-dependent antibiotic: 12 A-domains - But these were not





Streptomyces coelicolor

“Antibiotics factory”

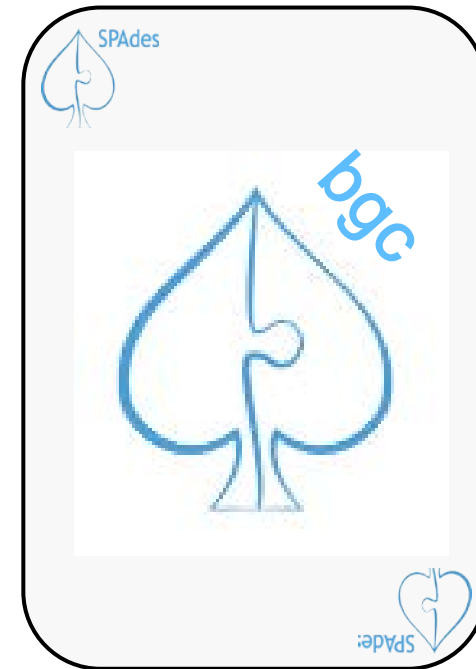
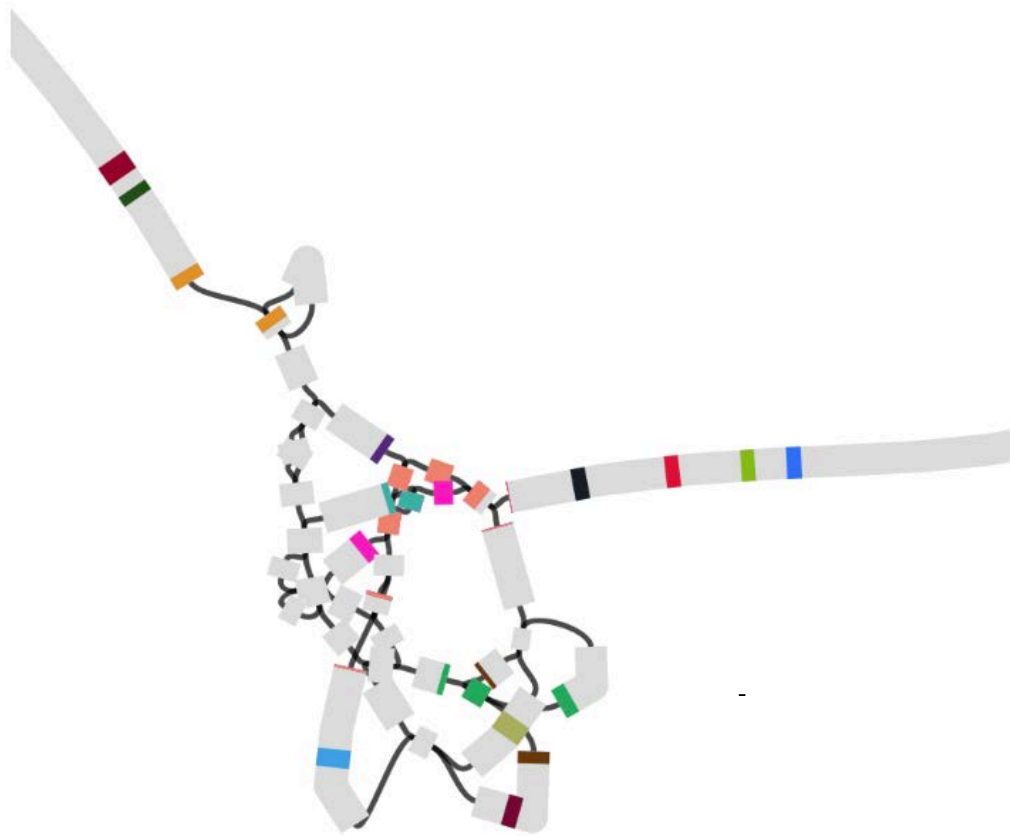
29 biosynthetic gene clusters, including 3 NRPS
(nonribosomal peptide synthetase)

- Coelichelin: 3 A-domains Assembled with SPAdes in 1 contig in minutes
- Nogalamycin: 2 A-domains Assembled with SPAdes in 1 contig in minutes
- Calcium-dependent antibiotic: 12 A-domains -> need special tool to deal with complex structures to assemble these domains!

New tool to solve difficult repeat structures

bgcSPAdes

converts complex subgraphs like this



Minutes vs 8 months!

Metabolo-genomics

- DNA sequencing technologies in **genomics**
- high-throughput mass spectrometry-based analysis that includes database search of **mass spectra** against already known metabolites
- **algorithms** and **bioinformatics** tools to expedite research

Global Natural Products Social (GNPS) molecular network



Please Login to Use Workflows

The Future of Natural Products Research and Mass Spectrometry



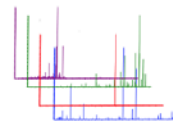
[Tweet](#) [Share](#)

Data Analysis



The [Data Analysis](#) portal will allow you to organize and visualize your mass spectrometry data. Leveraging the molecular networking techniques, there are additional tools to aid in understanding the unknowns in your sample. Check out the [documentation](#) and live [demo](#). Further, a separate [dereplication workflow](#) is provided as a standalone workflow.

Create Public MassIVE Datasets



[Submit](#) your own data to be made public MassIVE datasets. These MassIVE datasets must be prefixed with GNPS to be visible to other GNPS users. Take advantage of [continuous identification](#) to learn more about your dataset after publication automatically. New hits to the community curated libraries and related datasets are reported. [Documentation](#)

GNPS содержит более миллиарда масс-спектров биологически активных соединений, собранных во всем мире

=> необходимы алгоритмы поиска Натуральных пептидных продуктов (PNP) в этой массе данных

<https://gnps.ucsd.edu/ProteoSAFe/static/gnps-splash.jsp>

Antibiotics discovery challenge

How to discover potential antibiotic among millions of proteins?

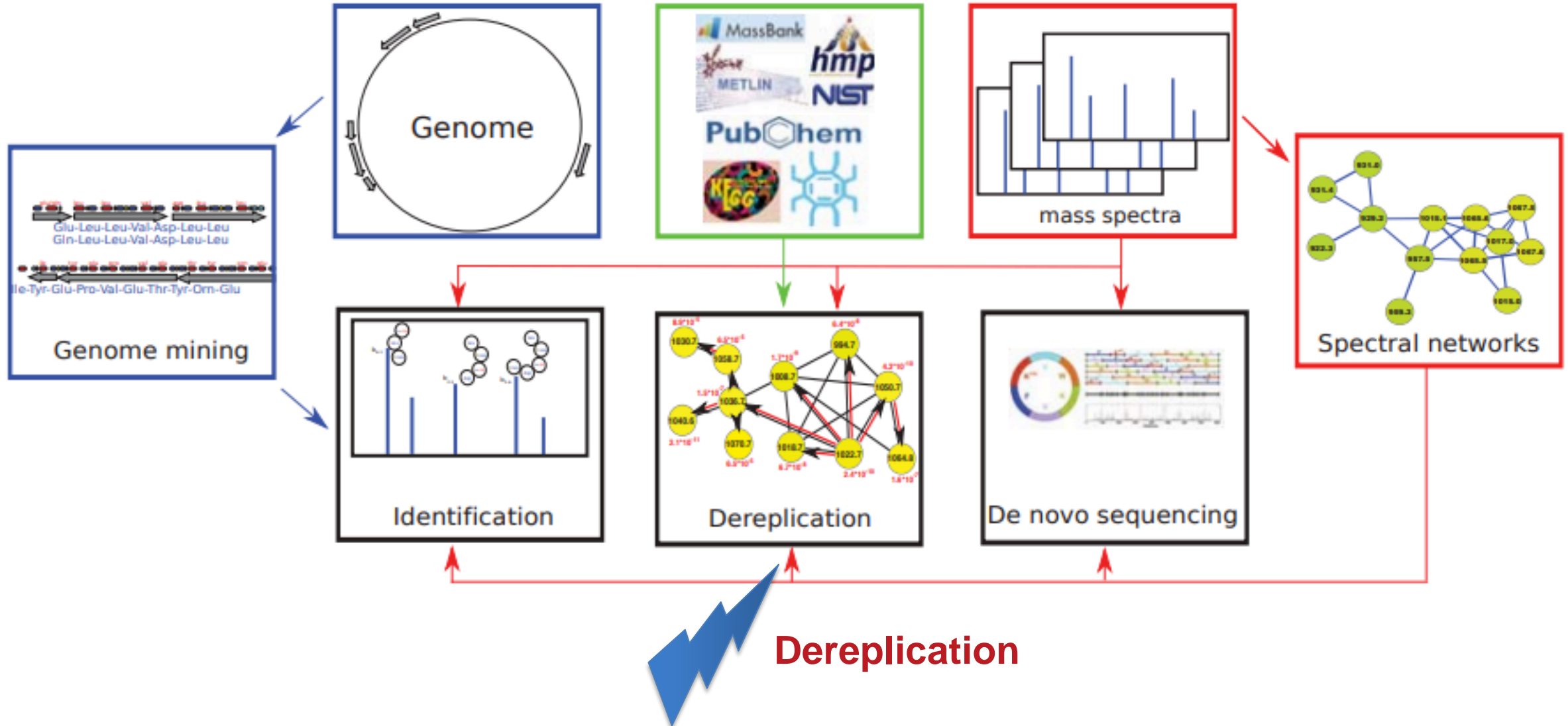


PNP discovery challenge

How to discover new PNP among millions of proteins?

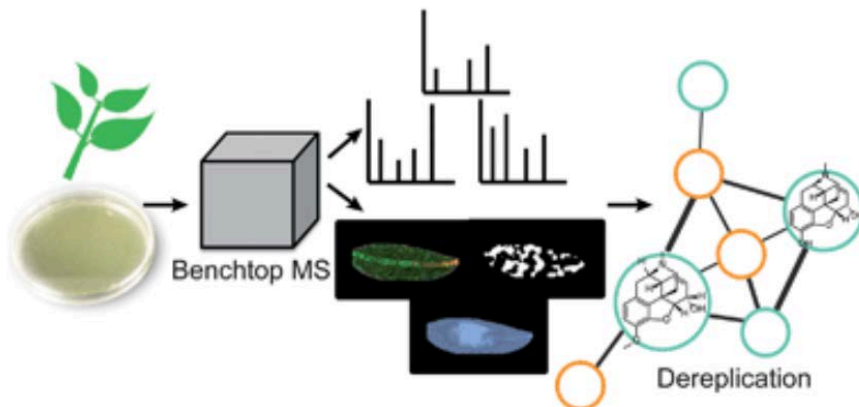


Metabolo-genomics



Dereplication

- тестирование образцов смесей для распознавания и исключения из рассмотрения тех активных веществ, которые уже были ранее охарактеризованы



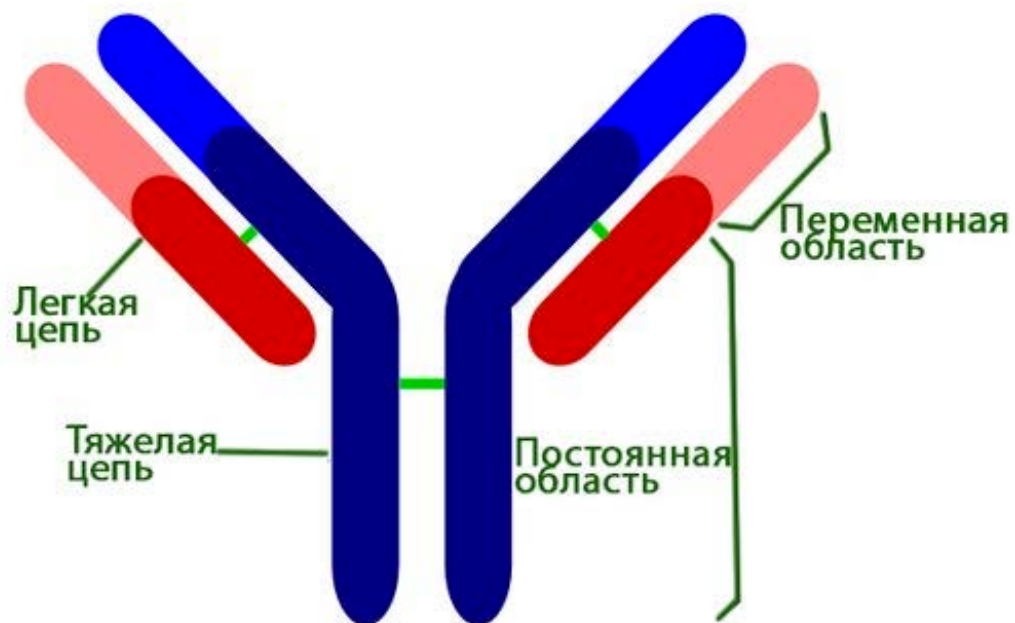
Dereplicator

- a computational tool for identification of known natural products from LC/MS-MS data.

- ✓ использует базу данных **химических** структур
- ✓ генерирует in-silico **масс-спектры** соединений, предсказывая, как они фрагментируются во время масс-спектрометрии
- ✓ объединяет их с **экспериментальными LC / MS-MS** для обнаружения сходства
- ✓ конвертирует оценки подобия в статистическую значимость



Immunoinformatics at CAB

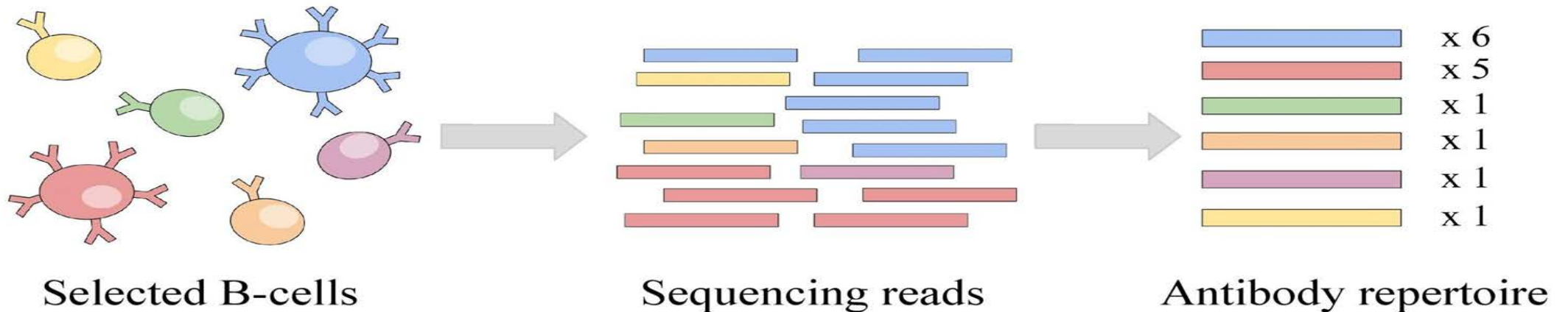


Антитело (англ. antibody) — белок (иммуноглобулин), синтезируемый В-лимфоцитами в организме человека и животного в ответ на попадание в него чужеродного вещества и обладающий специфическим сродством к этому веществу.

Антитела продуцируются иммунной системой для нейтрализации антигенов (чужеродных молекул)

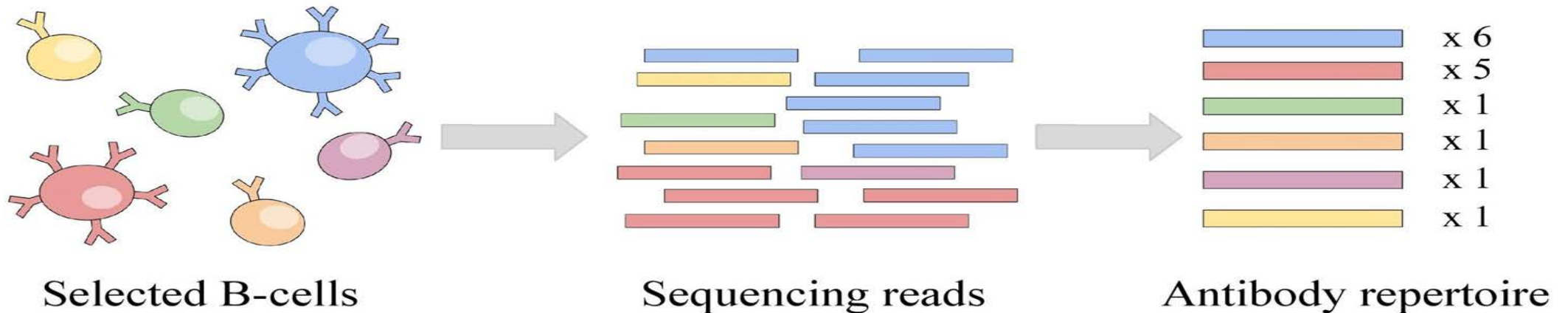
Репертуар антител

Репертуар антител — суммарное количество типов антител, характеризующих функциональное состояние иммунной системы индивидуума в данный момент

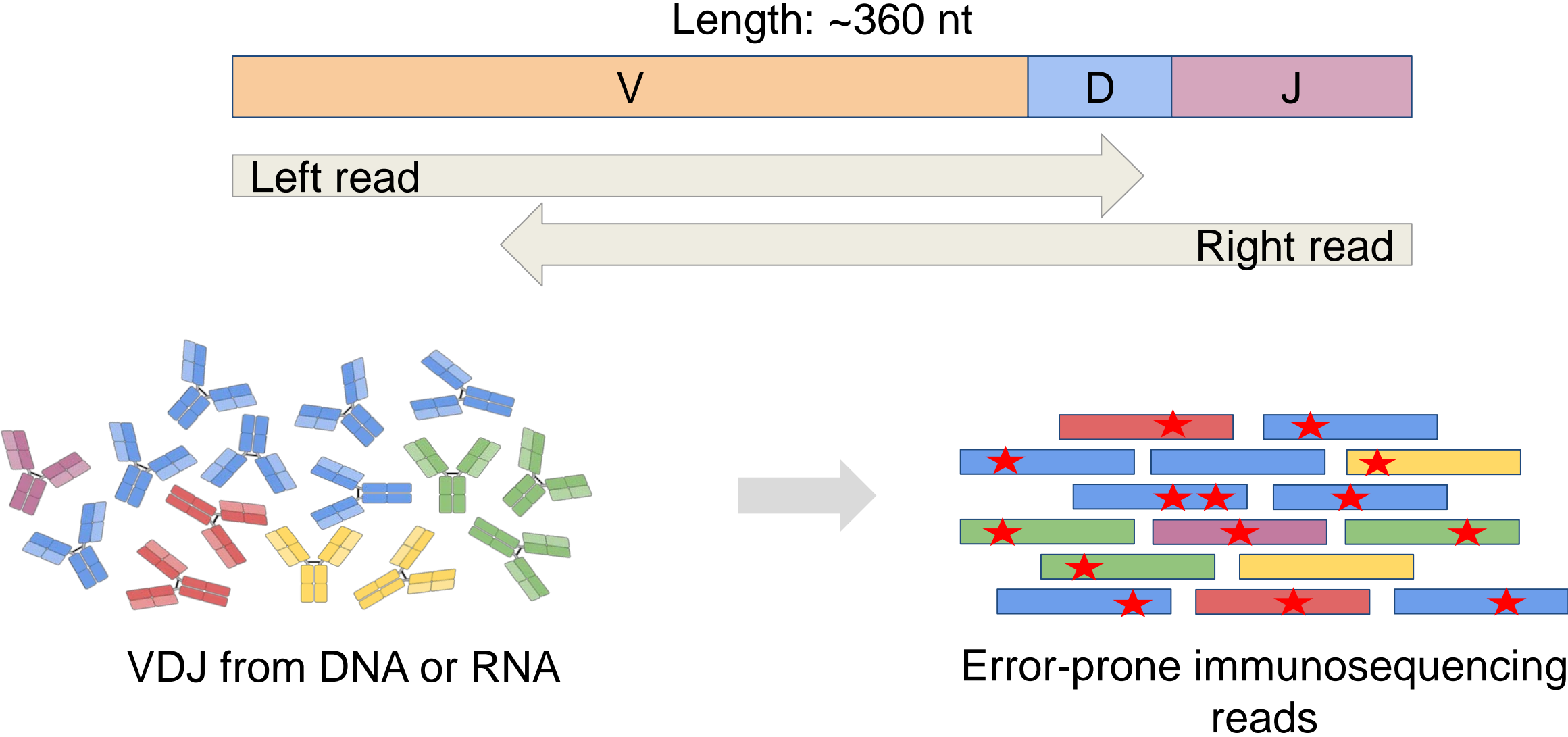


Репертуар антител

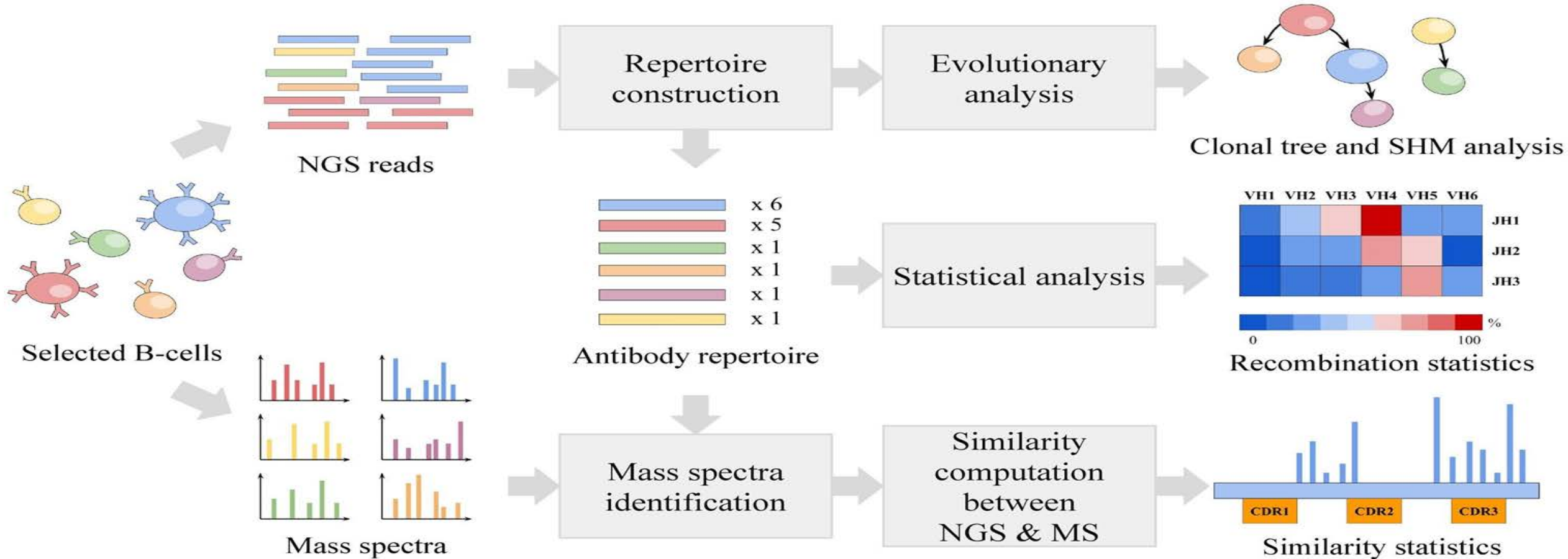
Построение репертуара антител из данных сиквенса ДНК (NGS) и его анализ являются важными шагами в разработке лекарственных препаратов на основе антител и клинических исследований



Antibody repertoire sequencing (Rep-seq)



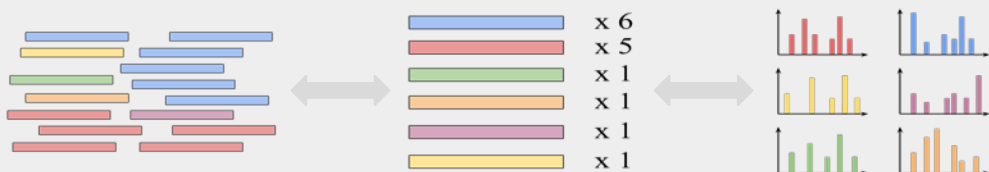
Immunoinformatic analysis



Immunosequencing data opens new bioinformatics challenges

Immunoinformatics at CAB

IgRepertoireConstructor



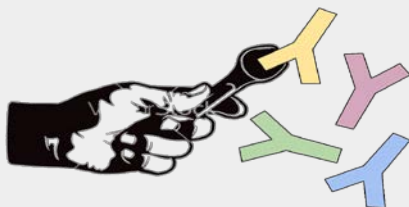
tool for construction of antibody repertoire using mass spectra

IgAnalyzer & IgQUAST



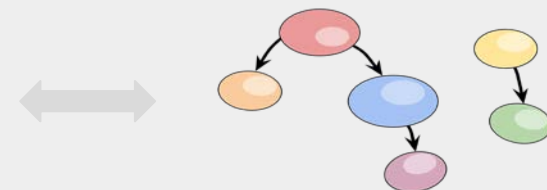
tools for “quality” assessment of antibody repertoire

IgSimulator



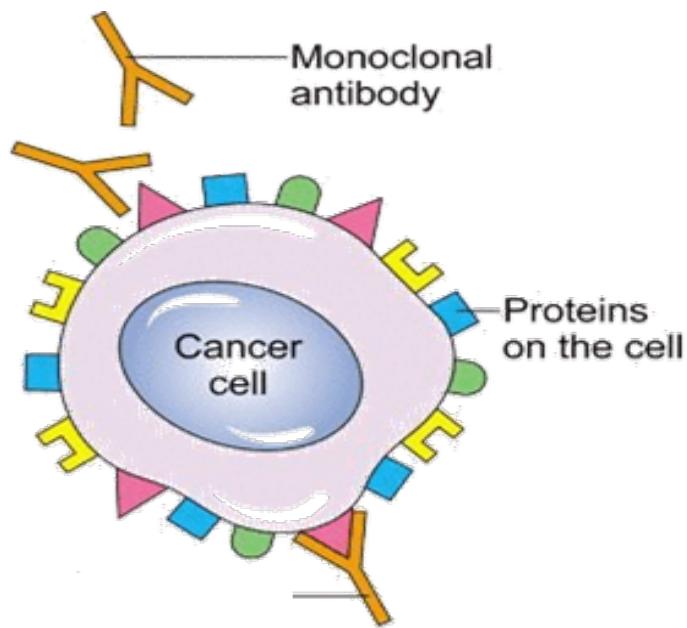
versatile immunosequencing simulator

AntEvolvo



tool for construction of clonal trees and analysis of somatic hypermutation (SHM)

Биомедицинские применения



Cancer treatment based on monoclonal antibody therapy



Studying of autoimmune disorders and immunodeficiency



Immunization analysis and monitoring of treatment

Review

Bioinformatic and Statistical Analysis of Adaptive Immune Repertoires

Victor Greiff,¹ Enkelejda Miho,¹ Ulrike Menzel,¹ and Sai T. Reddy^{1,*}

REVIEW

Open Access



Practical guidelines for B-cell receptor repertoire sequencing analysis

Gur Yaari^{1*} and Steven H. Kleinstein^{2,3*}

Sequencing and bioinformatics analysis of adaptive immune repertoire are becoming routine in clinical studies



HHS Public Access

Author manuscript

Nat Rev Rheumatol. Author manuscript; available in PMC 2015 April 01.

Published in final edited form as:

Nat Rev Rheumatol. 2015 March ; 11(3): 171–182. doi:10.1038/nrrheum.2014.220.

Sequencing the functional antibody repertoire—diagnostic and therapeutic discovery

William H. Robinson

Division of Immunology and Rheumatology, CCSR 4135, 269 Campus Drive, Stanford, CA 94305, USA. wrobins@stanford.edu

REVIEW

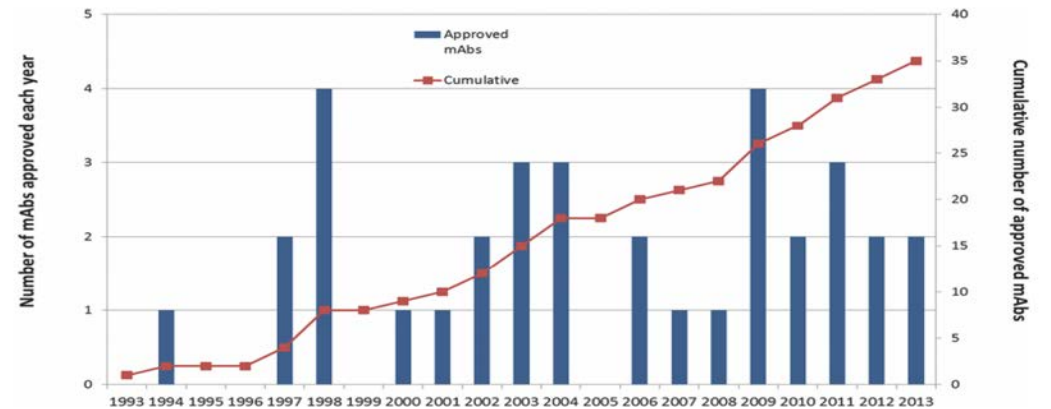


The promise and challenge of high-throughput sequencing of the antibody repertoire

George Georgiou^{1–4}, Gregory C Ippolito^{3,4}, John Beausang^{5,6}, Christian E Busse⁷, Hedda Wardemann⁷ & Stephen R Quake^{5,6,8,9}

Antibody drugs

In 2017, global sale revenue for antibody drugs was \$90 billion



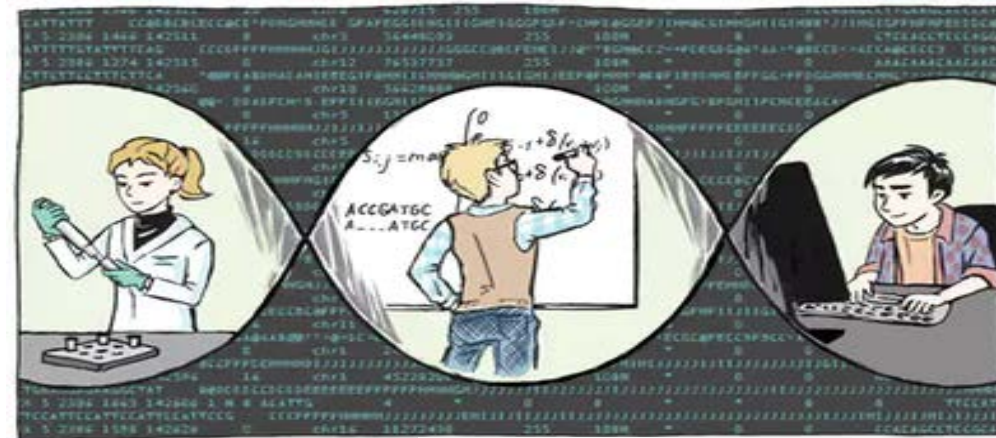
О чем



1. Центр алгоритмической биотехнологии - немного сухой статистики
2. Что такое биоинформатика?
3. Наука в Центре алгоритмической биотехнологии
4. Образование

Bioinformatics education: **online**

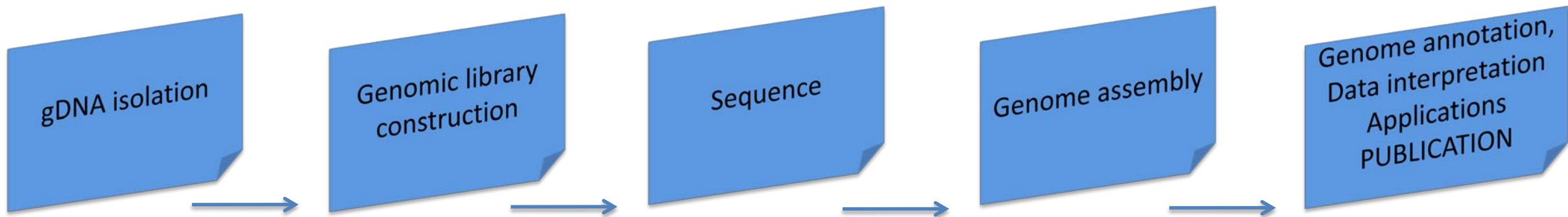
- **Introduction to bioinformatics for biologists**
 - Coursera (2014)
- **Introduction to bioinformatics: Metagenomics (2016)**
 - Coursera, Opened Education Platform, Stepik
- **Introduction into Linux for biologists (2015)**
 - Stepik



[\(Introduction to Bioinformatics\)](#)

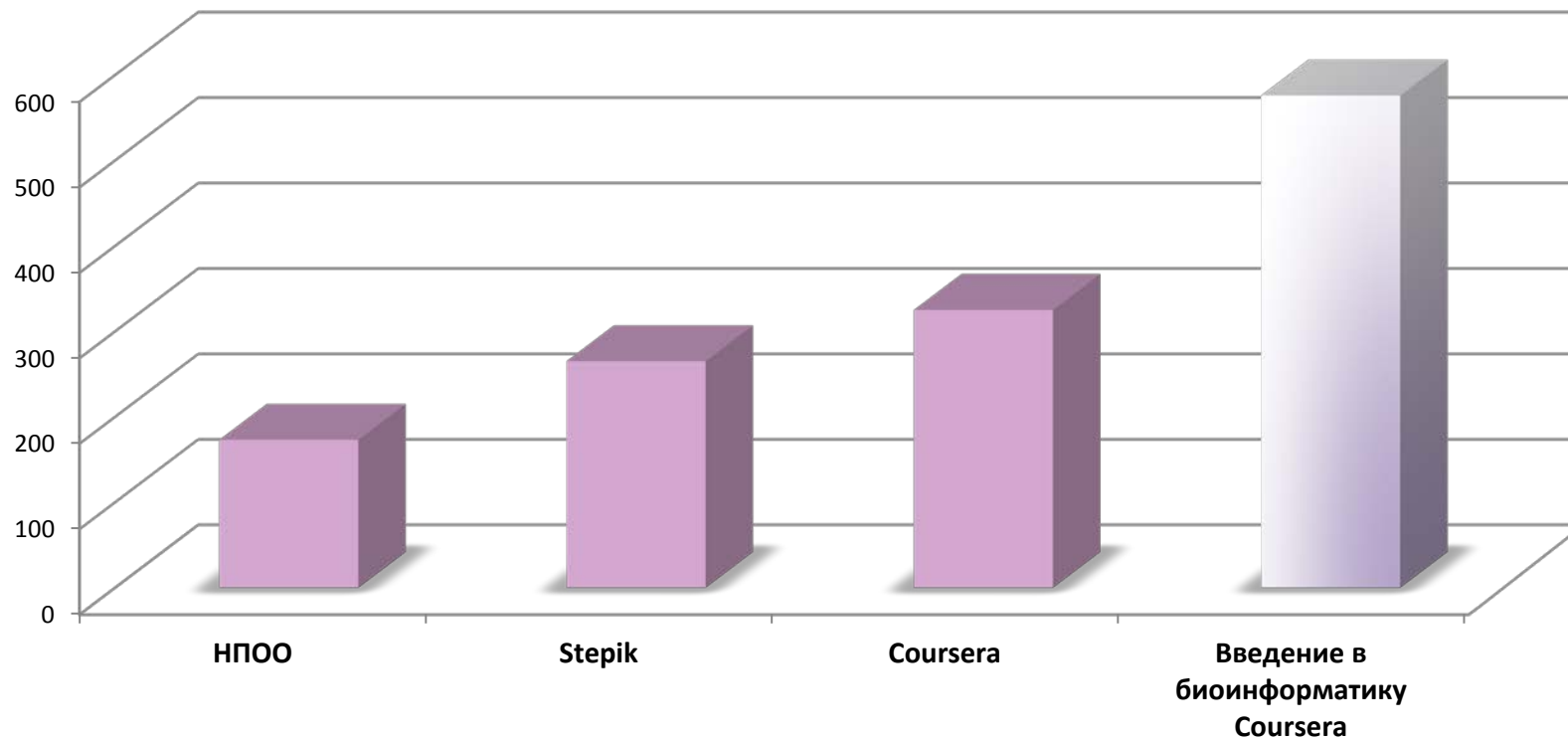
How they differ from other bioinformatics courses on Coursera (and not only Coursera)

- На русском языке
- New - “along the project” – wet lab + algorithms
- Practical



“Введение в биоинформатку” &

“Введение в биоинформатку: МЕТАГЕНОМИКА”



TOTAL: 1339 new students/month

Bioinformatics education: offline

- Bioinformatics for biologists
 - Introduction to bioinformatics (1 semester main course, MS)
 - Two year MS program “Bioinformatics” (**will start September 2018**)
- Algorithms, statistics and sequencing data analysis for computer scientists
- Genomics workshops (Russia, Finland, Italy)
- Summer schools for undergraduate students
- Summer internships at CAB

Bioinformatics education: offline





[HOME](#)

[REGISTRATION](#)

[SUBMIT ABSTRACT](#)

[WORKSHOP](#)

[TRAVEL](#)

[VENUE](#)

[COMMITTEES](#)

BIOINFORMATICS: FROM ALGORITHMS TO APPLICATIONS

July 16-19, 2018

Saint Petersburg, Russia

[REGISTER NOW](#)

<http://biata2018.spbu.ru/>

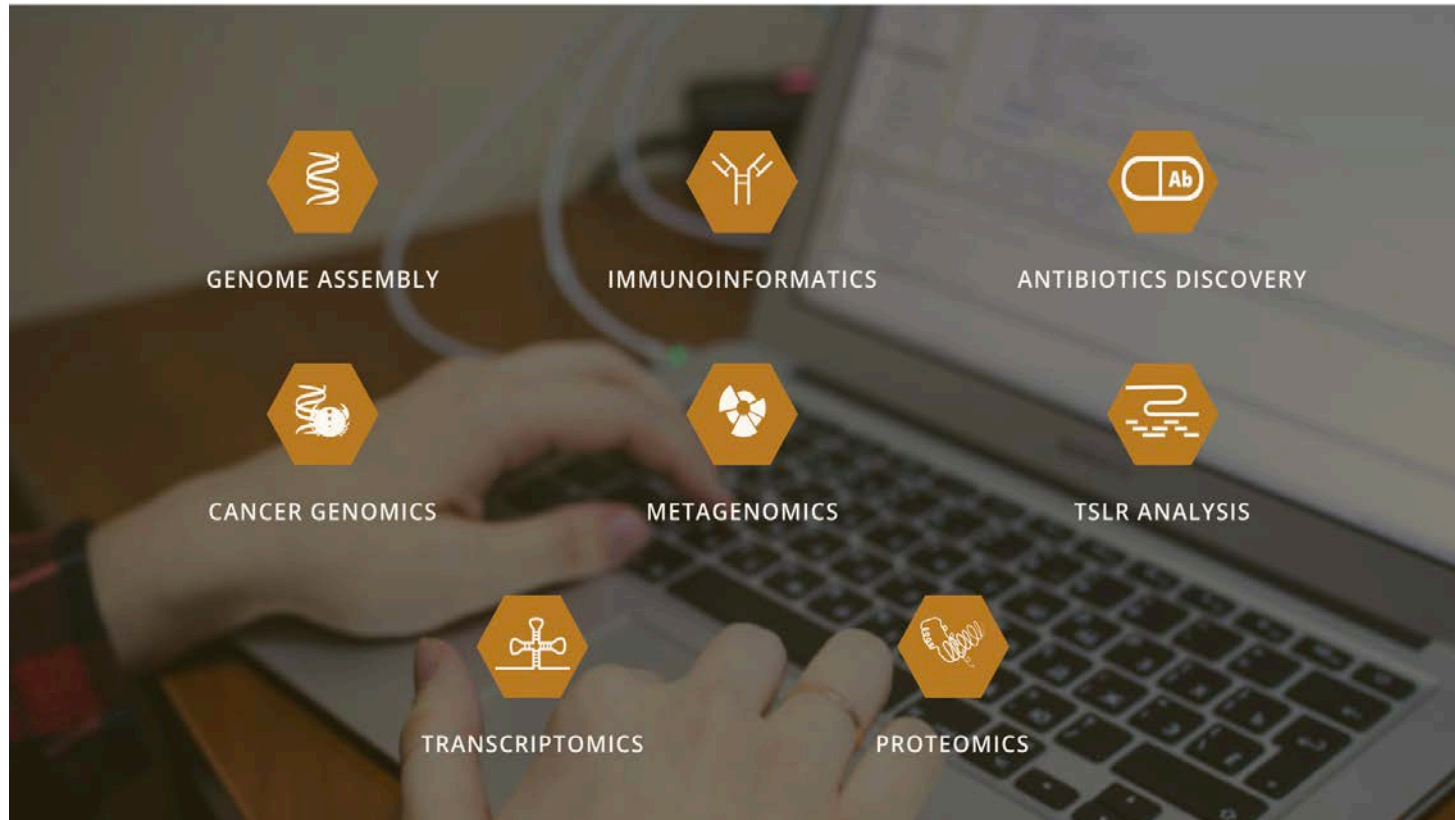
Acknowledgments



**Center for
Algorithmic
Biotechnology**

St. Petersburg State University

HOME ABOUT US RESEARCH SOFTWARE MEMBERS CONTACTS



(grant #14-50-00069)

<http://cab.spbu.ru/>