

Направления работы Института микробиологии НАН Беларуси в области геномики прокариот

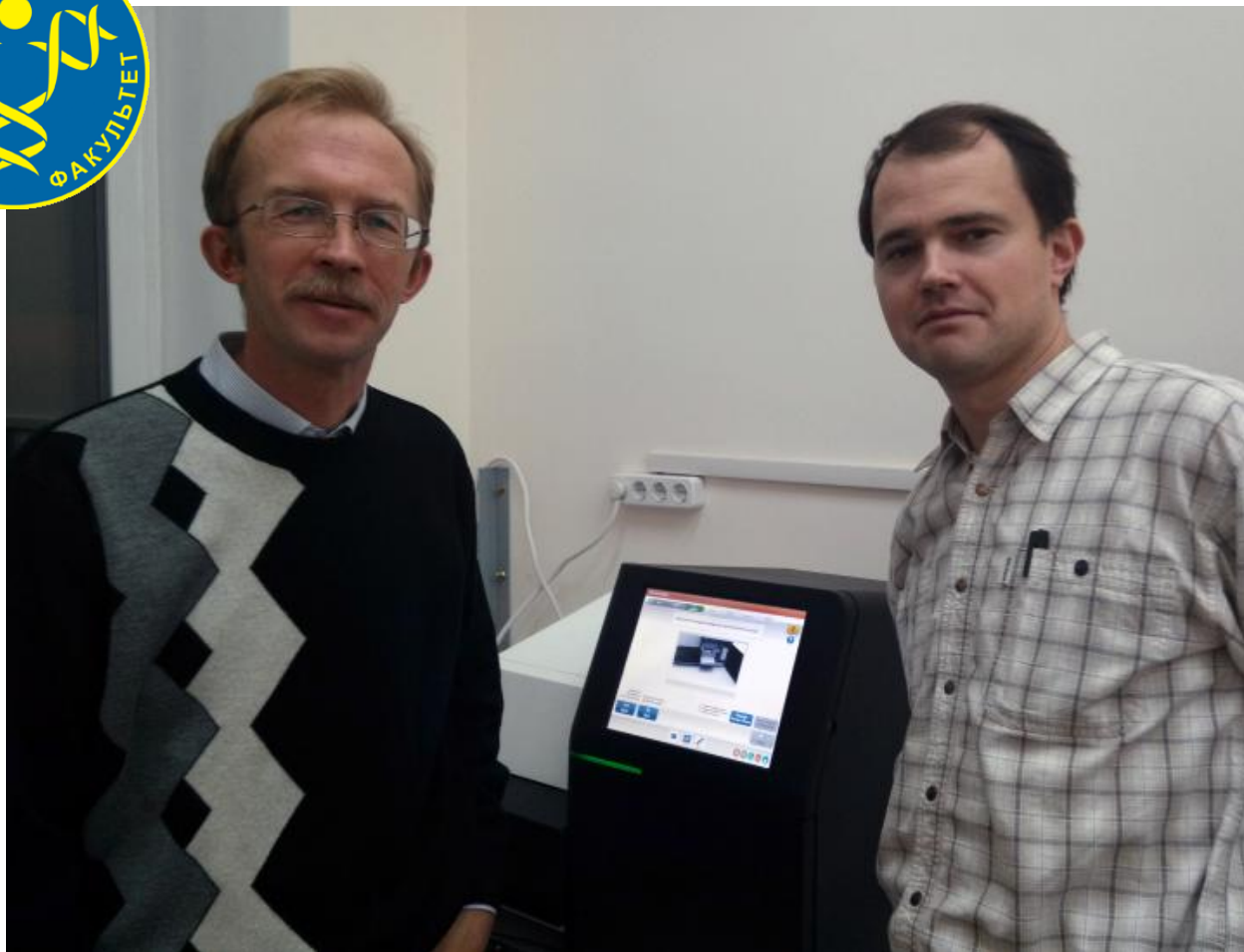


Валентович Л.Н.

Лаборатория «Центр аналитических и генно-инженерных исследований»
Институт микробиологии Национальной академии наук Беларуси

Биологический факультет БГУ

сотрудничает с нами...



Отчёт о первой пятилетке

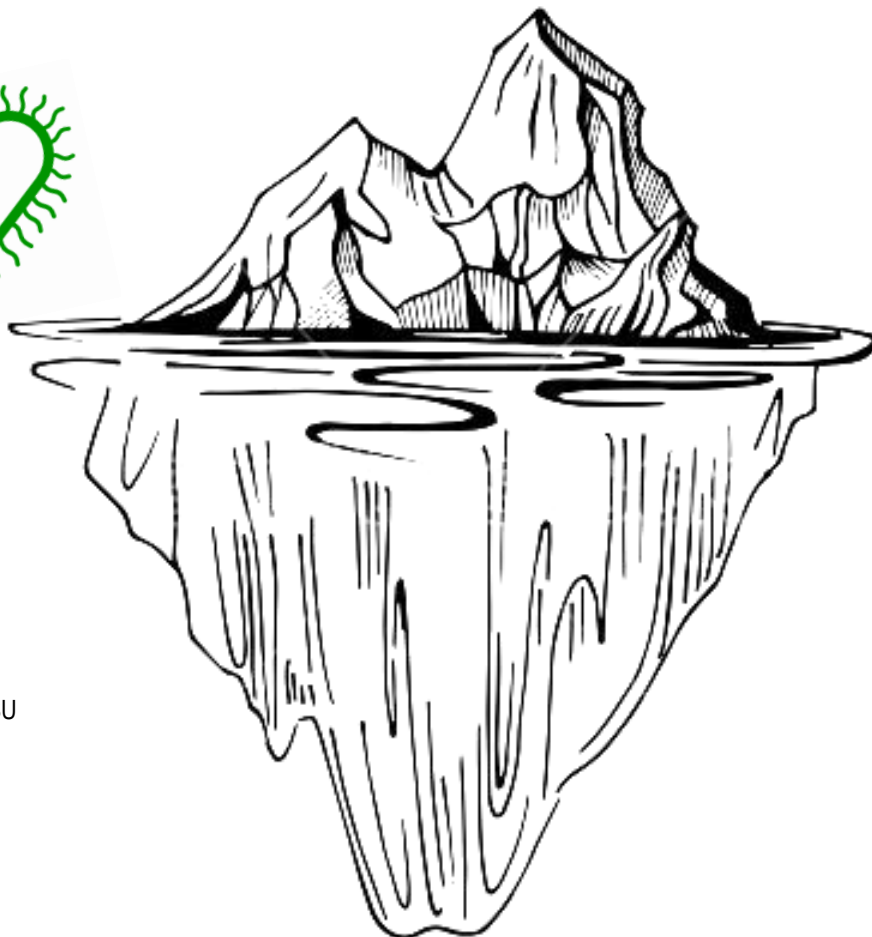
Список проектов по секвенированию (2014-2018 гг.)

• Дрожжей:

- 1) *Glaciozyma antarctica* t3-1a

• Бактерий:

- | | |
|--|---|
| 1) <i>Pectobacterium carotovorum</i> 3-2 | (Genbank: CP024842) |
| 2) <i>Pectobacterium carotovorum</i> 14A | (Genbank: CP034276-CP034278) |
| 3) <i>Pectobacterium atrosepticum</i> 36A | (Genbank: CP024956) |
| 4) <i>Bacillus velezensis</i> BIM B-439D | (Genbank: CP032144) |
| 5) <i>Staphylococcus warneri</i> 22-1 | (Genbank: CP032158-CP032159) |
| 6) <i>Erwinia amylovora</i> E2 | (Genbank: CP024970-CP024971) |
| 7) <i>Erwinia amylovora</i> 1/79 | |
| 8) <i>Pectobacterium carotovorum</i> 25.1 | |
| 9) <i>Bifidobacterium animalis</i> subsp. <i>lactis</i> H1 (BioProject ID: PRJNA263203) | |
| 10) <i>Bifidobacterium animalis</i> subsp. <i>lactis</i> H3 (BioProject ID: PRJNA263288) | |
| 11) <i>Pseudomonas brassicacearum</i> S-1 | |
| 12) <i>Pseudomonas corrugate</i> 3' | |
| 13) <i>Pseudomonas syringae</i> pv. <i>lachrymans</i> 8 | 21) <i>Carnobacterium</i> sp. 11т.7.20.2 |
| 14) <i>Pseudomonas fluorescens</i> BIM B-582 | 22) <i>Porphyrobacter sanguineus</i> gip-4 |
| 15) <i>Pseudomonas guineae</i> 7т.4.20.1 | 23) <i>Bacillus pumilus</i> 63-2-2 |
| 16) <i>Pseudomonas lundensis</i> 2т.2.5.2 | 24) <i>Bacillus pumilus</i> 11-1-1 |
| 17) <i>Pseudomonas</i> sp. 3.C. 1.9 | 25) <i>Bacillus pumilus</i> F6 |
| 18) <i>Leifsonia rubra</i> Хт.6.20.4 | 26) <i>Bacillus pumilus</i> BIM B- 263 |
| 19) <i>Stenotrophomonas maltophilia</i> 92 | 27) <i>Rhodococcus pyridinivorans</i> L5A-BSU |
| 20) <i>Sporosarcina psychrophila</i> 5т.3.20 | 28) <i>Clavibacter michiganensis</i> Н.П. |



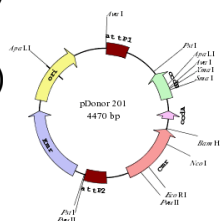
Бактериофагов:

- 1) бактериофаг Pf-10 (Genbank: KP025626)
- 2) бактериофаг фЕа2809 (Genbank: KP037007)



Мегаплазмид:

- 1) pBS72 (Genbank: KX711616)



• Метагеномных образцов (16S рДНК):

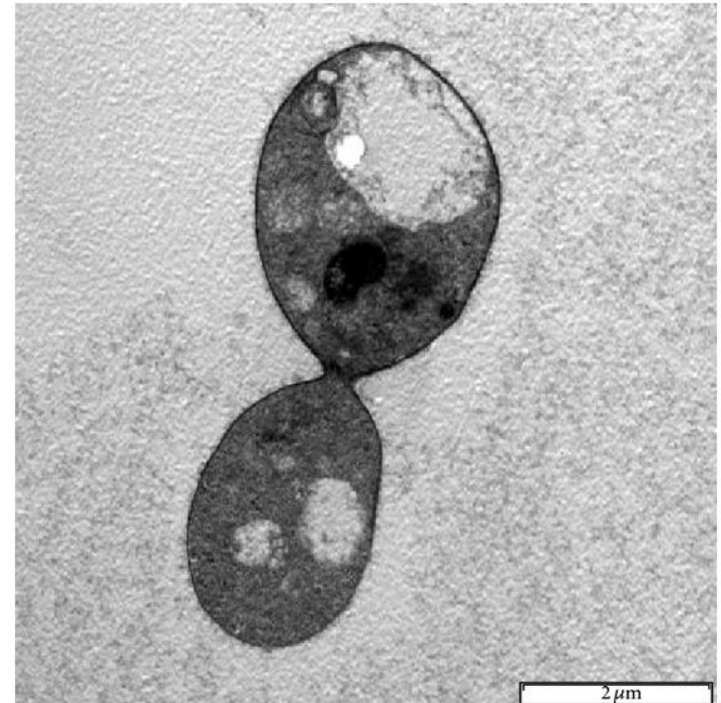
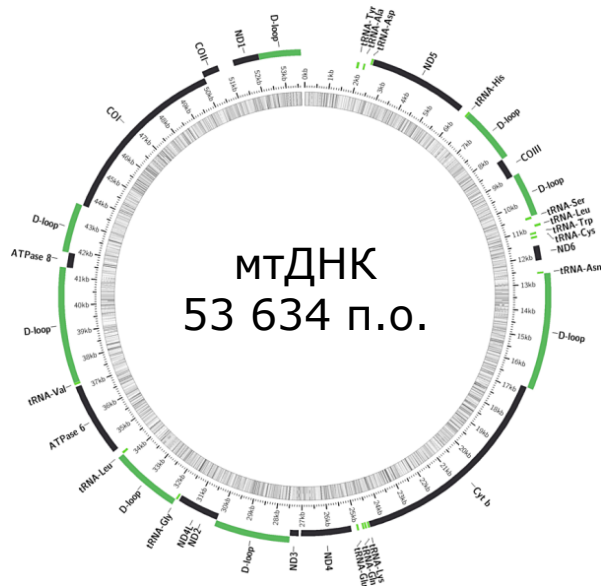
- 1) Озера Нарочь
- 2) Различных биотопов Антарктики
- 3) Микробиома кишечника человека

A map of Antarctica with the continent shaded in light gray. The surrounding oceans are white. Three research stations are marked with red dots and labeled: 'McMurdo' on the eastern coast, 'Palmer' on the western coast, and 'Davis' on the northern coast. A small inset map in the top left corner shows the location of Antarctica within the Southern Hemisphere, with a red box indicating the continent's position.

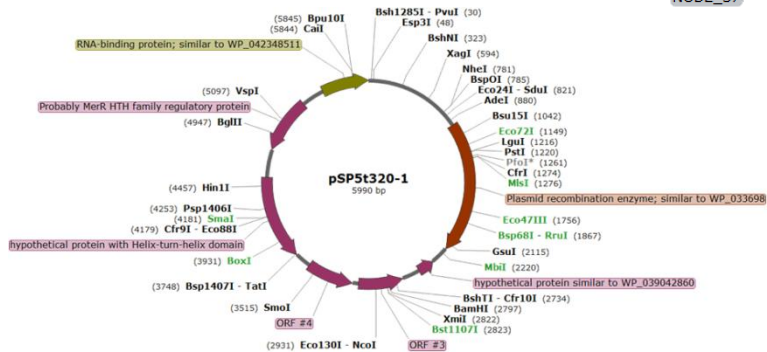
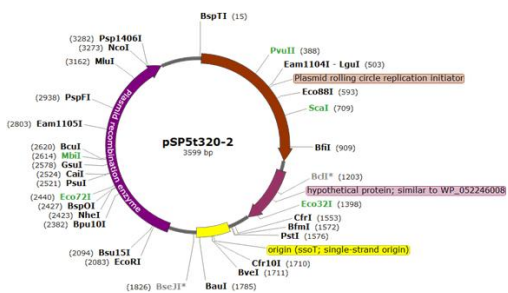
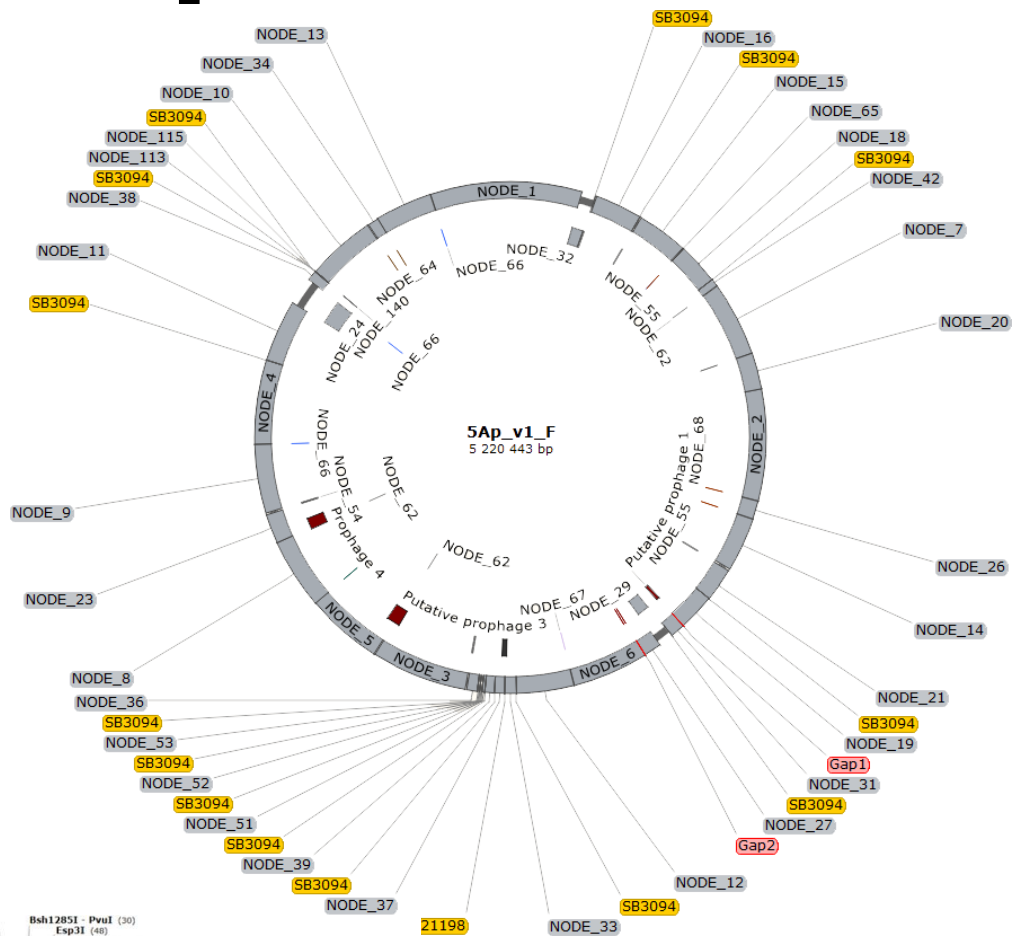
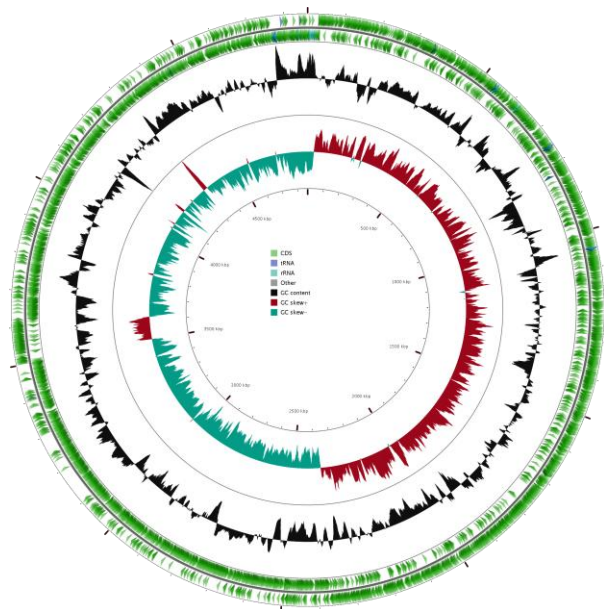
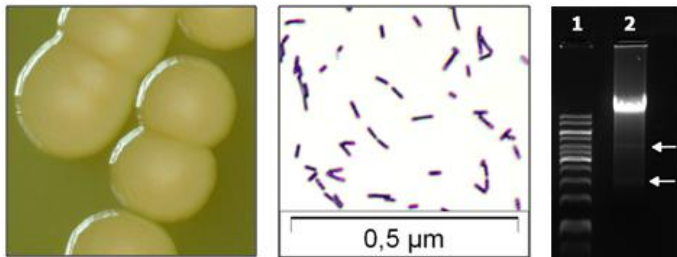
Прочитано нами: $\approx 19\,628\,587$ п.н.

Среднее покрытие: 18

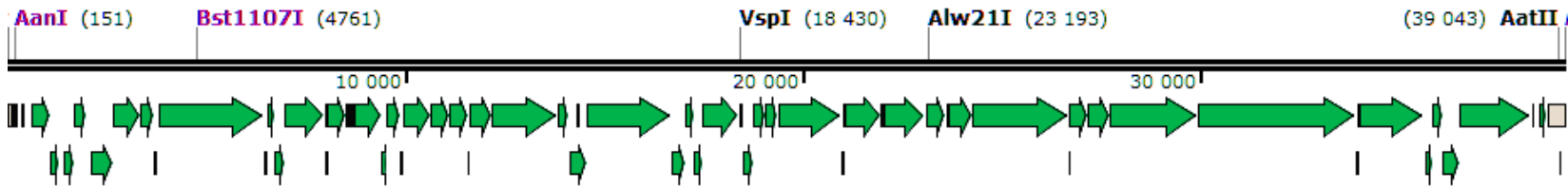
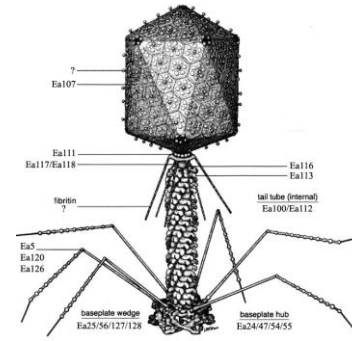
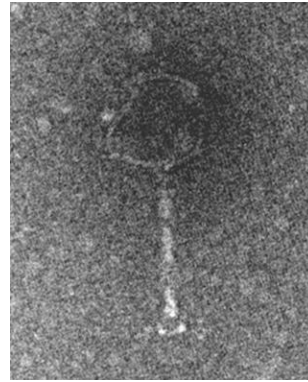
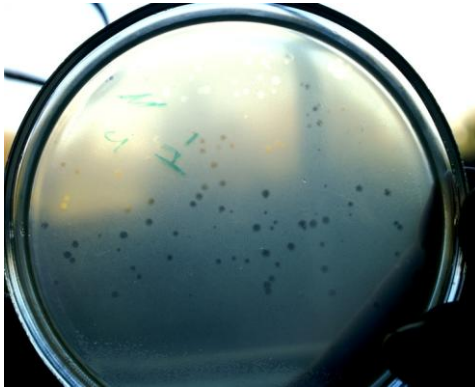
Количество хромосом: ?



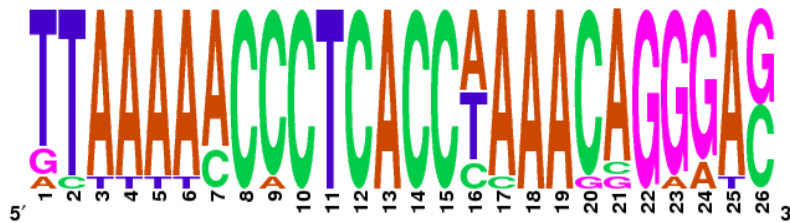
Бактерии



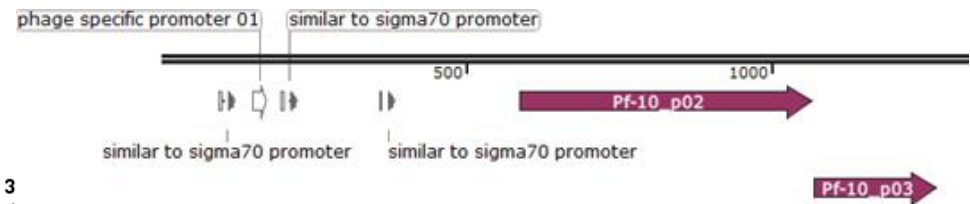
Вирусы и плазмиды



Pf-10
39 167 bp



Фагоспецифический промотор Pf-10



Используемые приборы



Секвенатор **LI-COR 4300**

Определение нуклеотидной последовательности по методу Сэнгера. Размер прочтений – до 1000 пар оснований.



Секвенатор **Illumina MiSeq**

Определение нуклеотидной последовательности с помощью технологии синтеза (SBS) от компании Illumina. Размер прочтения – до 300 пар оснований.

Пробоподготовка:

Nextera XT DNA Library Prep Kit – Illumina

MuSeek Library Preparation Kit, Illumina compatible - Thermo Fisher

Секвенирование: Illumina MiSeq, наборы MiSeq Reagent Kit v3 - 2×300

Исходные данные для работы

Секвенирование по технологии illumina: на выходе миллионы строк в FASTQ формате

```
read :
  filtered :
    control# :
      index_sequence
```

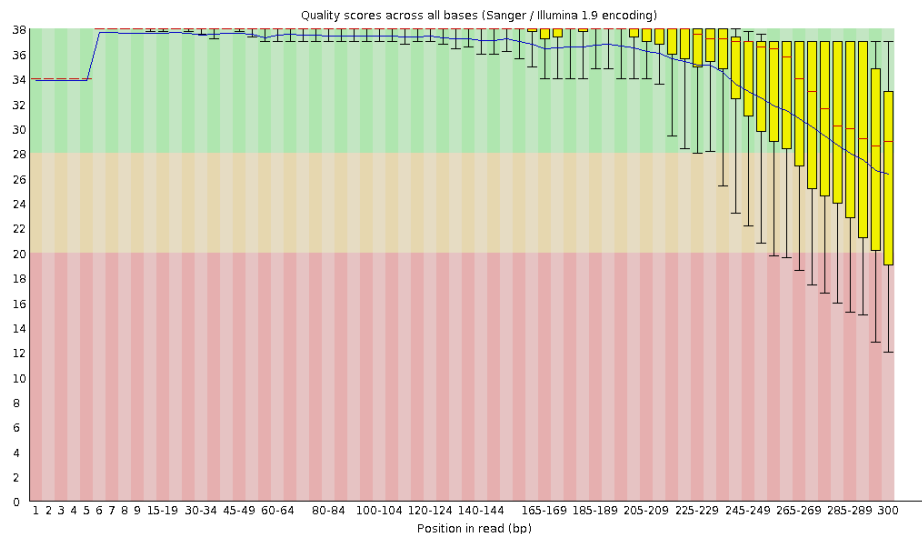
instrument : run#: flowcell_ID :lane: tile : x-pos : y-pos

@HWI-ST992:147:D22HDACXX:8:2313:8106:8355 1:N:0:ATGTCA
CCTCGTTAAGAATGAACGTAGGGAAGCCGTTAGGCATAGTTACGGAGACCCTTGCATAT
+
CCCCFFFFHHHHHJJJJJHIJJJJJJJJJIJJJJJJIGIGIJJJIJJJIJH
@HWI-ST992:147:D22HDACXX:8:2313:8144:8417 1:N:0:ATGTCA
ATAATGAATCTGAAATATAAATTGGAATTATGATTTACAGATGTTCAAACACTTTGAAC
+
CCCCFFFFHHHGFIGGHIJIIJJIEGHJGIHGIIGHIJJGJHGIGFIIGHJEHHIJIJ
@HWI-ST992:147:D22HDACXX:8:2313:8391:8324 1:N:0:ATGTCA
AGTCTATTTTTGAATTTAGTG GTATGACCAGGGCTACTACATTTTCCAATATCTTCAA
+
CC@FFFFFHHHHHFJJJJJHIJFHJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
@HWI-ST992:147:D22HDACXX:8:2313:8483:8332 1:N:0:ATGTCA
AAATTAATTACAGTTATTGTAATTATTATTATATTAGGATGATTCACAAATTGGAAT
+
CCCCFFFFFHGH HHIIIIIJJGGHEIGIIJJJJJJJJJIGIFIJIDGGJJJJJJJJIGHJJII

«Стандартный» путь обработки данных

- Анализ и обработка «сырых» данных (прочтений): **FastQC** и **Trimmomatic**
- Сборка прочтений - получение контигов: **SPAdes** и **A5**
- Фильтрация контигов: **Blast2GO** и **собственные скрипты**
- Анализ оставшихся контигов и сборка их в скаффолды: **тонкая ручная работа в программе Snapgene**
- Секвенирование по Сэнгеру спорных мест. Проверка с помощью ПЦР.
- Проверка сборки – поиск однонуклеотидных вариантов, небольших вставок и делеций: **Pilon**
- Проверка сборки – картирование прочтений на собранный геном: **BWA** и **Bowtie 2**, визуализация: **Tablet**
- Автоматическая аннотация генома: **Rast, Prokka, Phast (Phaster)**
- “Ручная” аннотация генома: **Sigmoid, PHIRE, BLAST** и др.

Анализ и обработка «сырых» данных (прочтений): FastQC и Trimmomatic




<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>



Удаление всех адаптеров и последовательностей с Q меньше 25 и длиной менее 50 нуклеотидов

```
java -jar trimmomatic-0.32.jar PE /media/user/Data/temp/BadH1_S3_L001_R1_001.fastq  
/media/user/Data/temp/BadH1_S3_L001_R2_001.fastq /media/user/Data/temp/lanel_R1_paired.fastq  
/media/user/Data/temp/lanel_R1_unpaired.fastq /media/user/Data/temp/lanel_R2_paired.fastq  
/media/user/Data/temp/lanel_R2_unpaired.fastq ILLUMINACLIP:/home/user/soft/Trimmomatic-  
0.32/adapters/NexteraPE-PE.fa:2:30:10 LEADING:3 TRAILING:3 SLIDINGWINDOW:4:25 MINLEN:50
```



```
ATGAGCACGGCATAGACCCCAAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCG  
TCGGAATTAACAGACAAATCGCTCCAAAGATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGT  
TGAGGTTGTGAGGAAAATGCTCACTGTCCCAAGGAAGATCGGAAGAGCACACGTCTGAACTCCAGT  
CGATGTATCTCGTATGCCGTCTTCTGCTTGAAAAAAGAACACACAACGAAACCGATCGCACG  
CTTTCTCAACGCCATGTACTCTGCGTTGATACCACTGCTTAGATCGGAAGAGCACACGTCTGAA  
CACACGTCTG  
CGTCTTTTGT  
CACCGATGTAT  
AGCACACGTCT  
GTCTGAACTCC  
CGATGTAGCTT  
TATCTCGTATGC  
GGTTCAGGTCTACGAACACTGCCCGGGGAGATCGGAAGAGCACACGTCTGAACTCCAGTAACCGAT  
GATTTACAGCTTTCTTACTGCGTCATCTATATCAGAGAAGATCGGAAGAGCACACGTCTGAACTCCA
```

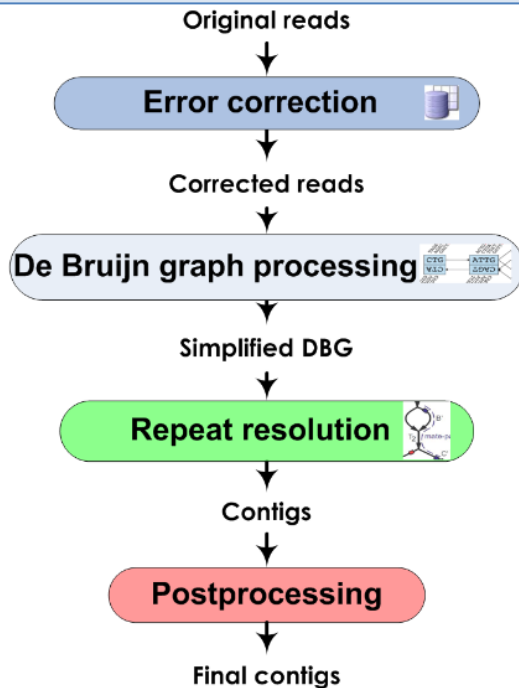
<http://www.usadellab.org/cms/?page=trimmomatic>

Сборка прочтений - получение контигов: SPAdes и A5



SPAdes

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



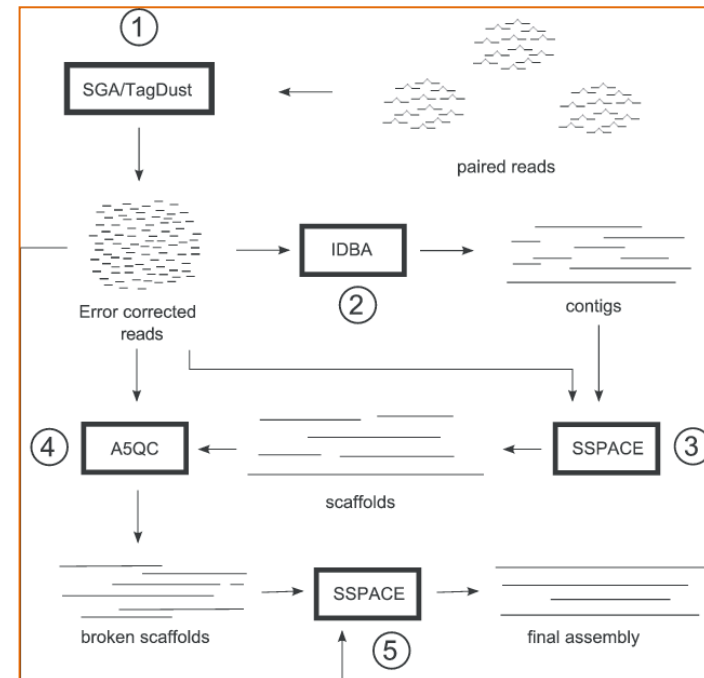
1	2	3
#name	length	coverage
NODE_6	777	14.4252
NODE_4	4666	24.3604
NODE_3	5925	12.647
NODE_2	35102	11.9551
NODE_5	790	62.2327
NODE_1	132734	11.9265



A5-miseq

<https://sourceforge.net/projects/ngopt/>

Read cleaning: **Trimmomatic**;
Error correction: **SGA**;
Contig assembly: **IDBA-UD**;
Crude scaffolding;
Misassembly correction;
Final scaffolding



Фильтрация контигов: Blast2GO и собственные скрипты



<https://www.blast2go.com/>



D:\Work\sequences\NGS\H3\blast2go-h3\blast2go_20140811_0859.dat - Blast2GO V.2.8.0							
File Blast Mapping Annotation Analysis Statistics Select Tools View Support							
GO:0007067,GO:0016021							
<input checked="" type="checkbox"/>	nr	sequence name	seq description	length	#hits	min. eValue	sim mean
<input checked="" type="checkbox"/>	11	NODE_11_length_5513_cov_7.47791_ID_21	streptococcus vt complete genome	5513	3	0,0E0	99.67%
<input checked="" type="checkbox"/>	12	NODE_12_length_3986_cov_51.3921_ID_23	bifidobacterium animalis strain complete g...	3986	3	0,0E0	100.0%
<input checked="" type="checkbox"/>	13	NODE_13_length_3069_cov_35.8838_ID_25	bifidobacterium animalis strain complete g...	3069	3	0,0E0	100.0%
<input checked="" type="checkbox"/>	14	NODE_14_length_2445_cov_84.5207_ID_27	bifidobacterium animalis lactis bi-complet...	2445	3	0,0E0	100.0%
<input checked="" type="checkbox"/>	15	NODE_15_length_1776_cov_58.5992_ID_29	bifidobacterium animalis lactis complete g...	1776	3	0,0E0	99.67%
<input checked="" type="checkbox"/>	16	NODE_16_length_1455_cov_3.61295_ID_31	pectobacterium carotovorum carotovorum c...	1455	3	0,0E0	89.33%
<input checked="" type="checkbox"/>	17	NODE_17_length_1333_cov_132.844_ID_33	bifidobacterium animalis lactis complete g...	1333	3	0,0E0	100.0%
<input checked="" type="checkbox"/>	18	NODE_19_length_1273_cov_0.613438_ID_...	erwinia phage ea9- complete genome	1272	1	0,0E0	98.0%
<input checked="" type="checkbox"/>	19	NODE_21_length_1117_cov_0.919192_ID_...	erwinia phage ea9- complete genome	1117	1	0,0E0	99.0%
<input checked="" type="checkbox"/>	20	NODE_24_length_1029_cov_5.99335_ID_47	bacillus amyloliquefaciens plantarum com...	1029	3	0,0E0	99.0%
<input checked="" type="checkbox"/>	21	NODE_26_length_999_cov_1.79587_ID_51	pseudomonas brassicacearum brassicac...	999	3	0,0E0	94.0%
<input checked="" type="checkbox"/>	22	NODE_27_length_993_cov_13.6871_ID_53	pectobacterium carotovorum carotovorum c...	993	3	0,0E0	90.33%
<input checked="" type="checkbox"/>	23	NODE_28_length_984_cov_4.61494_ID_55	pseudomonas fluorescens complete geno...	984	3	0,0E0	93.33%
<input checked="" type="checkbox"/>	24	NODE_29_length_980_cov_10.4572_ID_57	pectobacterium carotovorum carotovorum c...	980	3	0,0E0	90.67%
<input checked="" type="checkbox"/>	25	NODE_31_length_959_cov_15.5264_ID_61	pectobacterium carotovorum carotovorum c...	959	3	0,0E0	87.67%
<input checked="" type="checkbox"/>	26	NODE_37_length_925_cov_12.312_ID_73	pectobacte				
<input checked="" type="checkbox"/>	27	NODE_39_length_907_cov_6.04231_ID_77	pectobacte				
<input checked="" type="checkbox"/>	28	NODE_41_length_906_cov_0.596919_ID_81	erwinia ph				
<input checked="" type="checkbox"/>	29	NODE_44_length_893_cov_8.81462_ID_87	bacillus ar				
<input checked="" type="checkbox"/>	30	NODE_46_length_886_cov_8.79183_ID_91	---NA---				
<input checked="" type="checkbox"/>	31	NODE_49_length_872_cov_4.0953_ID_97	pectobacte				



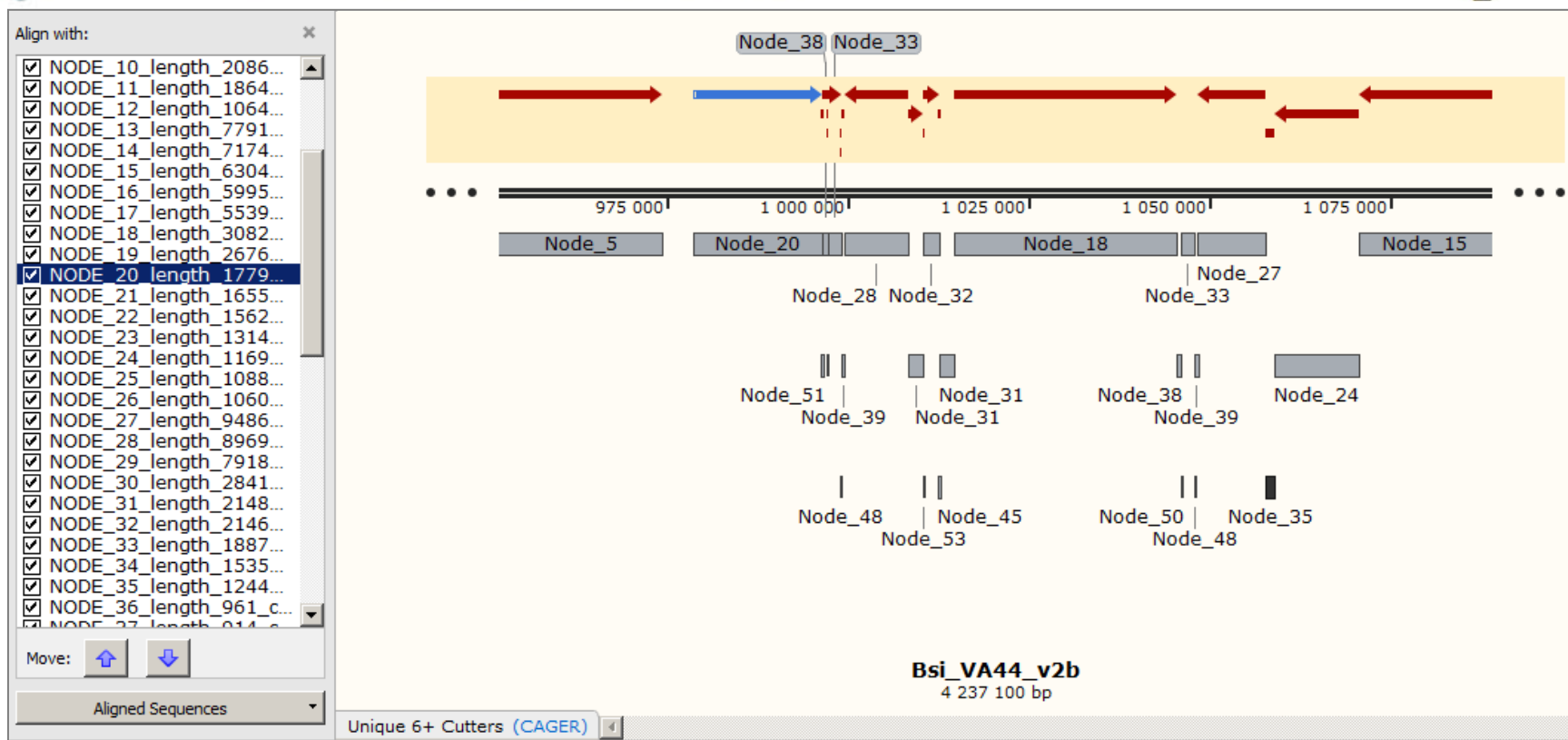
№	Название	Длина	Покрывие	ГЦ	Кэфф.	Аннотация	Начало	Конец
1	NODE_1_length_290866_cov_14.696621	290866	14,696621	61,08	1.0	Pseudomonas corrugata BS3649	-	-
2	NODE_2_length_270694_cov_15.249916	270694	15,249916	59,87	1.0	Pseudomonas corrugata BS3649	-	H_NODE_181 NODE_185_K
3	NODE_3_length_256694_cov_13.726800	256694	13,726800	60,42	0.9	Pseudomonas corrugata BS3649	-	-
4	NODE_4_length_216614_cov_12.766633	216614	12,766633	59,41	0.8	Pseudomonas corrugata RM1-1-4	-	-
5	NODE_5_length_176073_cov_12.056540	176073	12,056540	60,71	0.8	Pseudomonas corrugata RM1-1-4	-	-
6	NODE_6_length_157428_cov_14.060972	157428	14,060972	60,55	0.9	Pseudomonas corrugata RM1-1-4	-	NODE_112_K
7	NODE_7_length_153812_cov_14.955864	153812	14,955864	61,14	1.0	Pseudomonas corrugata BS3649	-	NODE_182_K NODE_184_K

Анализ оставшихся контигов и сборка их в скаффолды: тонкая ручная работа в программе Snapgene

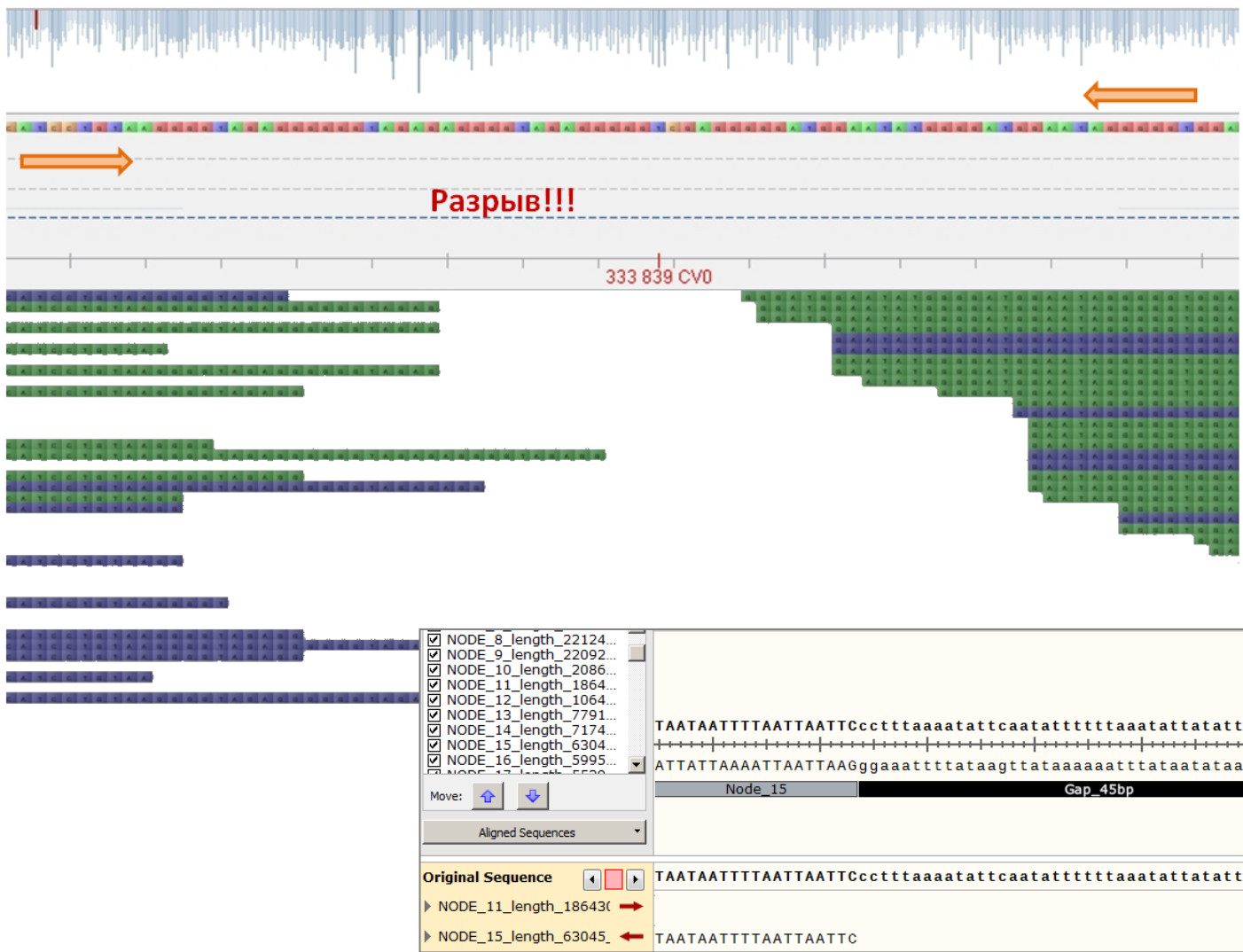


SnapGene®

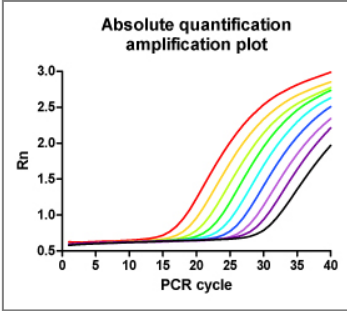
Software for molecular biology



Секвенирование по Сэнгеру спорных мест. Проверка варианта сборки с помощью ПЦР



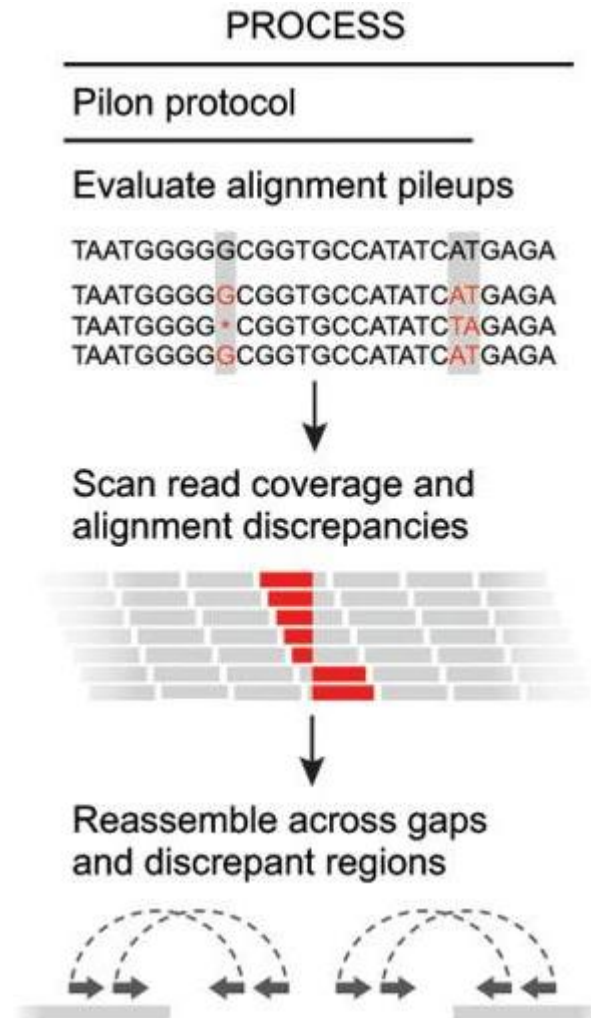
А также проверка количества повторяющихся элементов с помощью количественной ПЦР



Проверка сборки – поиск однонуклеотидных вариантов, небольших вставок и делеций: **Pilon**

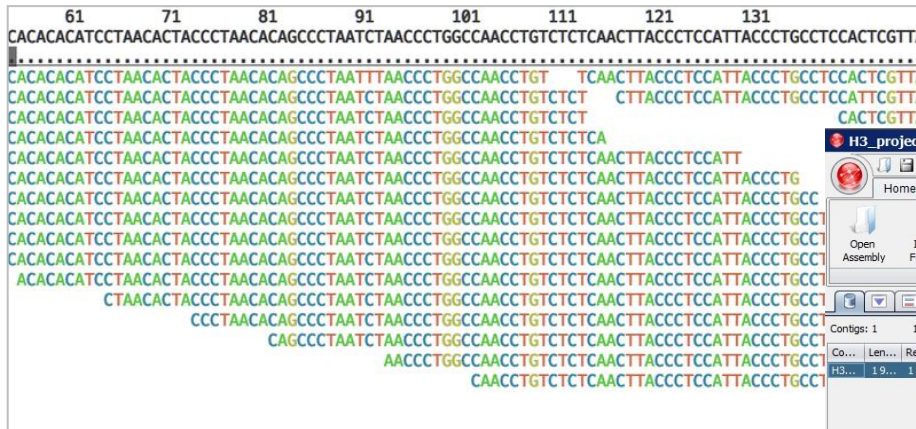


<https://github.com/broadinstitute/pilon/>



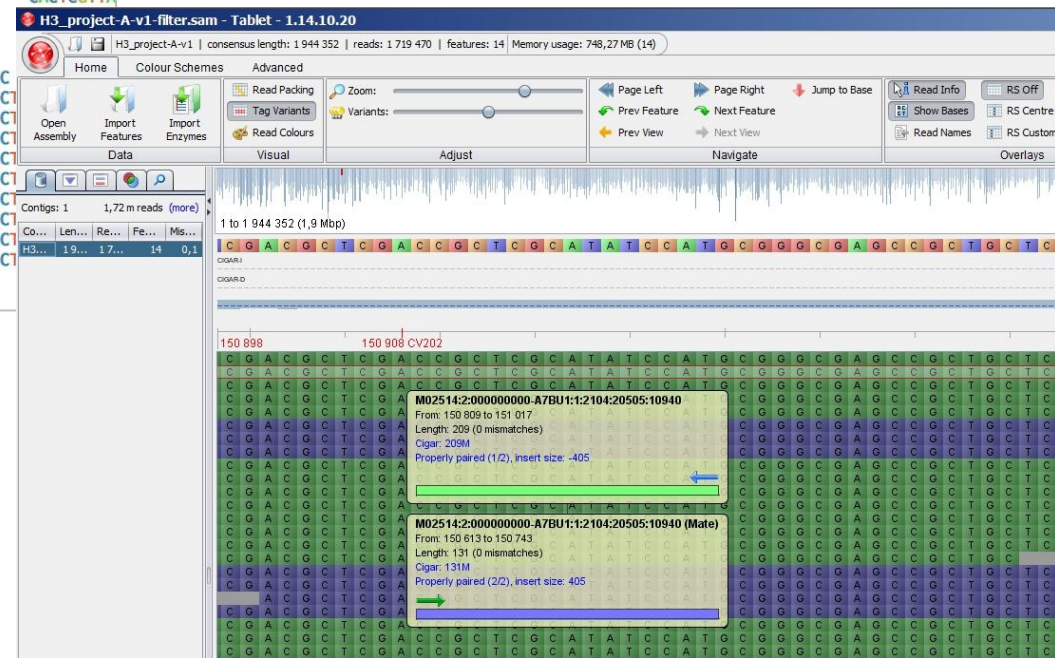
Проверка сборки – картирование прочтений на собранный геном: **BWA** и **Bowtie 2**, визуализация: **Tablet**

<https://ics.hutton.ac.uk/tablet/>



Burrows-Wheeler Aligner

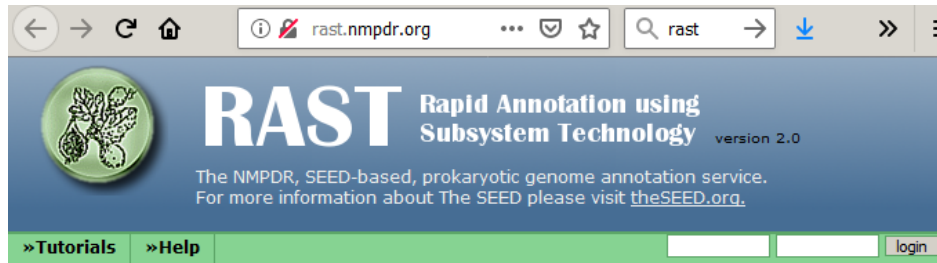
<http://bio-bwa.sourceforge.net/>



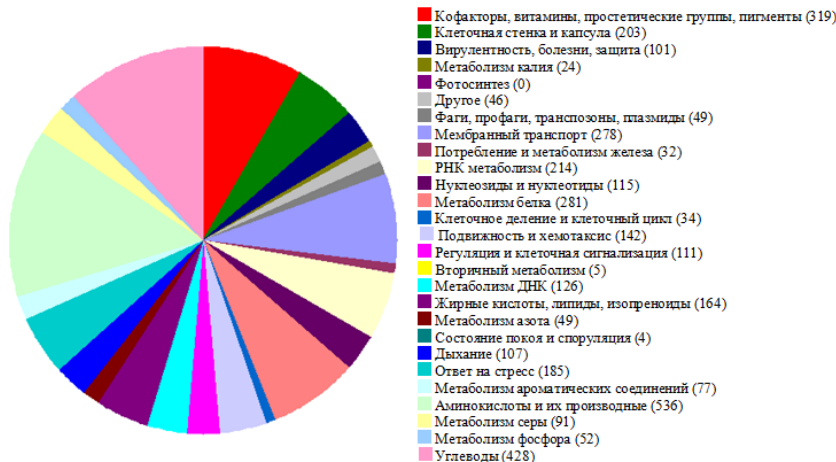
BOW TIE

<http://bowtie-bio.sourceforge.net/bowtie2/index.shtml>

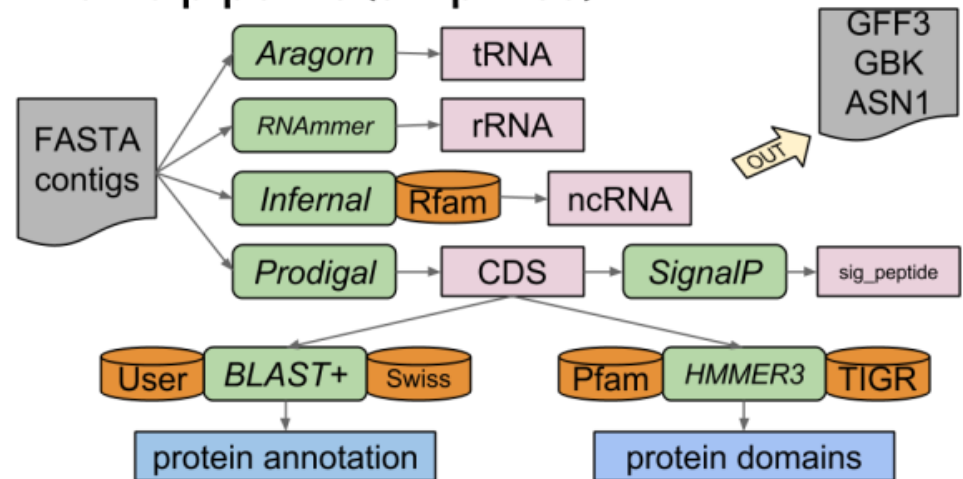
Автоматическая аннотация генома: Rast, Prokka, Phaster и др.



<http://phaster.ca/>

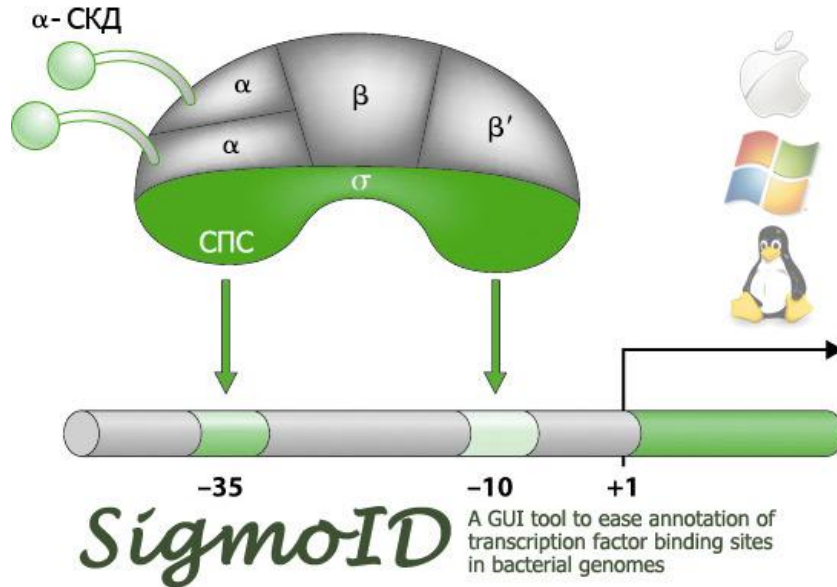


Prokka pipeline (simplified)

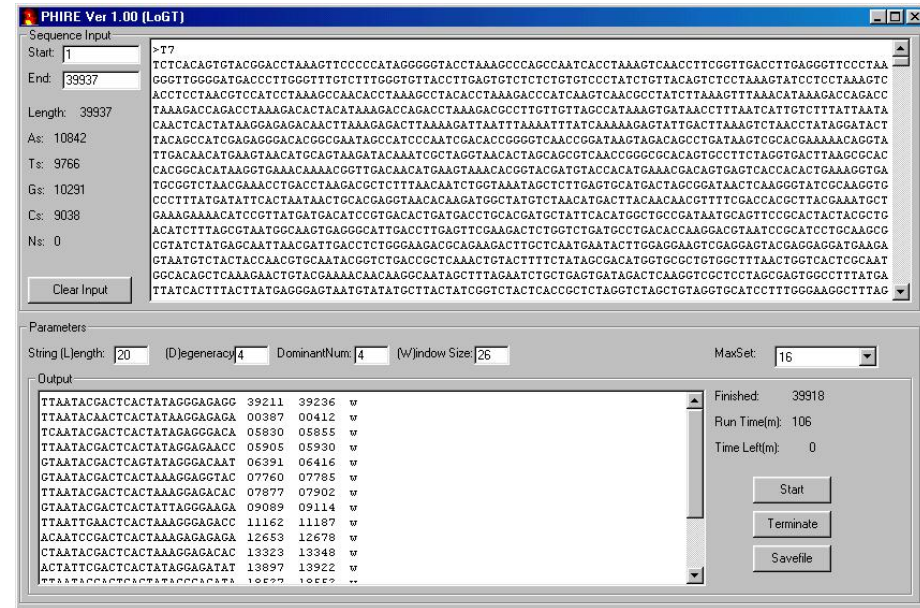


<http://www.vicbioinformatics.com/software.prokka.shtml>

“Ручная” аннотация генома: Sigmoid, PHIRE, BLAST и др.



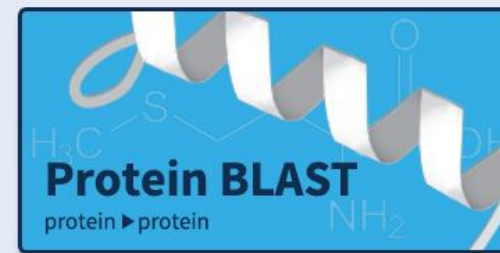
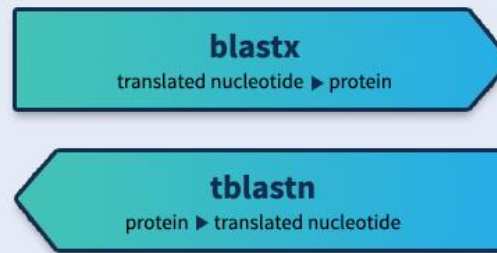
<https://github.com/nikolaichik/SigmoidID>



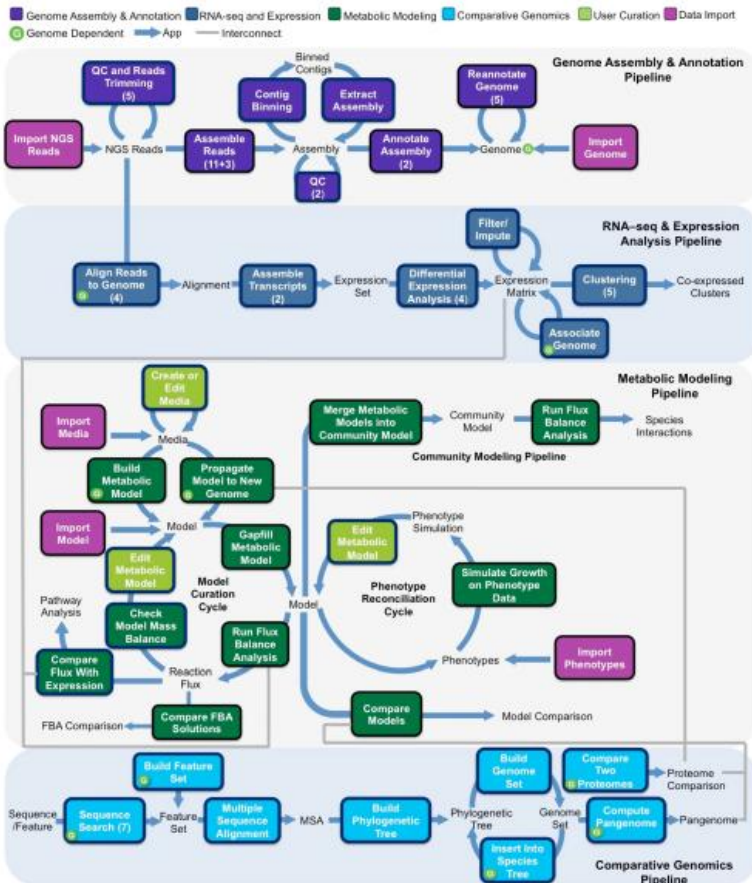
PHIRE (PHage *In silico* Regulatory Elements)

<https://www.biw.kuleuven.be/logit/PHIRE.htm>

Web BLAST



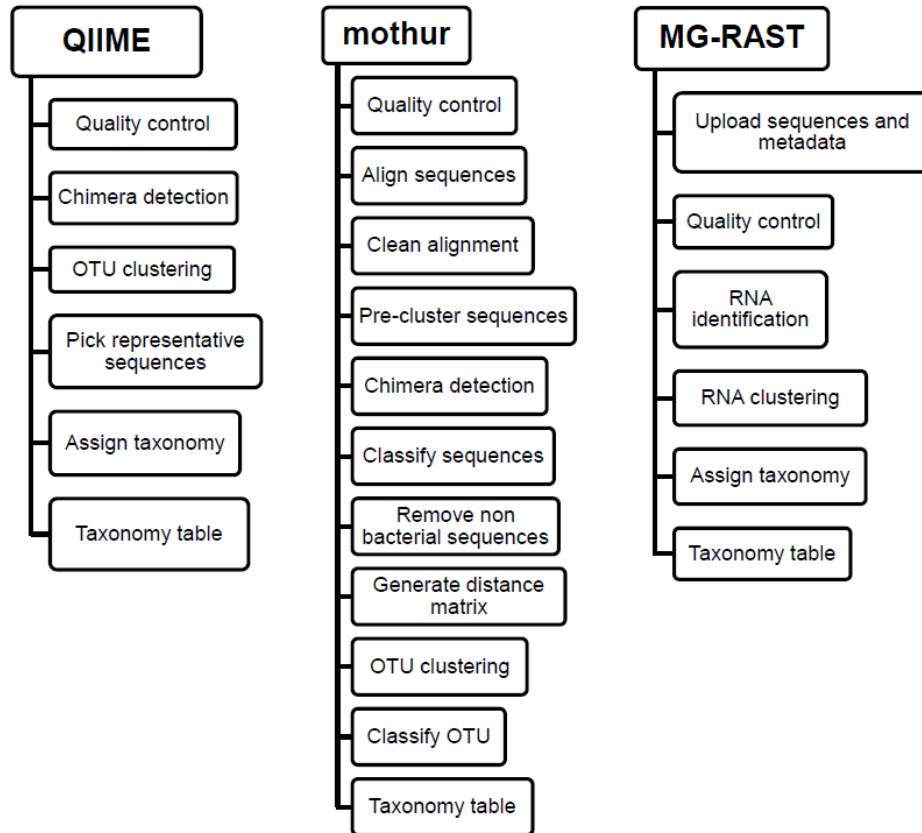
Использование сторонних серверов с предустановленным ПО



--- Common tools ---
--- Tools By Topic ---
Microbiology
Variant calling
Metagenomics
NGS: RNA Analysis
--- Other tools ---
OMERO.biobank

Bioinformatics, Volume 30, Issue 13, 1 July 2014, Pages 1928–1929,
<https://doi.org/10.1093/bioinformatics/btu135>

Работа с фрагментами **16S** рДНК от метагеномных образцов



Б/Д для классификации метагеномных данных

NCBI taxonomy classification
NCBI RefSeq bacterial genomes
NCBI RefSeq viral genomes
NCBI RefSeq GRCh38 human genome
Kraken DB: MiniKraken 4Gb database
CLARK-I DB: RefSeq bacterial+viral genomes
CLARK-I DB: RefSeq viral genomes
DIAMOND DB: UniRef50
DIAMOND DB: UniRef90



ONE CODEX

<https://www.onecodex.com/>



<http://taxonomer.iobio.io/> 19/20

Работа с метагеномными данными

Нарочь
2016



Образец донных
отложений

Получено прочтений: более 4 млн. × 2 (по 300 п.о.)

ОЗУ для сборки с помощью SPAdes: 72 Гб

Собрано контигов: 237 861

Общая длина контигов: 161 633 370 п.о.

Как с ЭТИМ
работать???

Бактерии	86%
Грибы	2,10%
Археи	1,80%
Вирусы	0,04%
Остальное	10%



Спасибо за внимание