

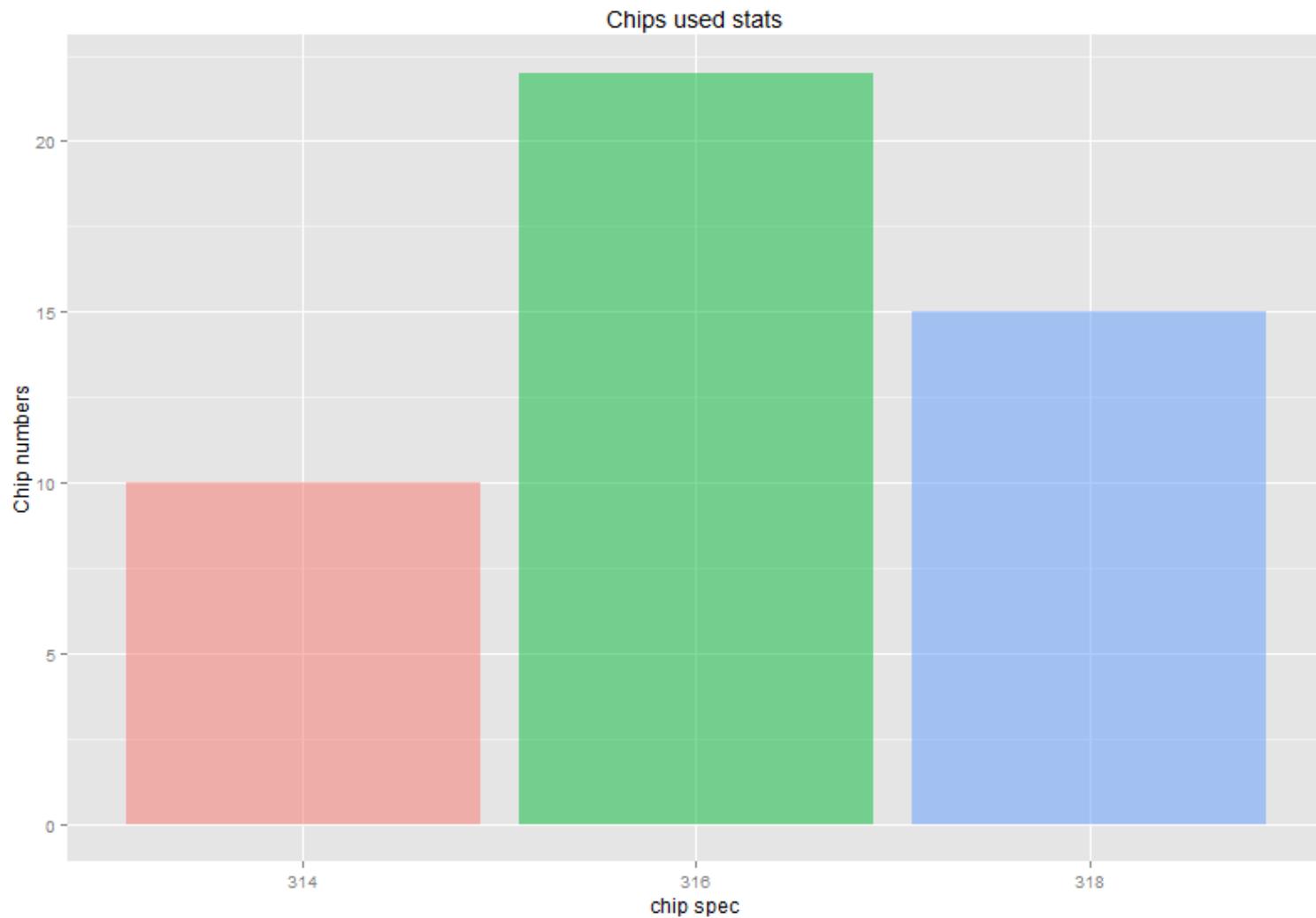
Наши

~~Мои~~ первые 100 проектов на
IonTorrent - рассказ
биоинформатика

Алексеев Дмитрий
зав. лаб. Биоинформатики НИИ ФХМ

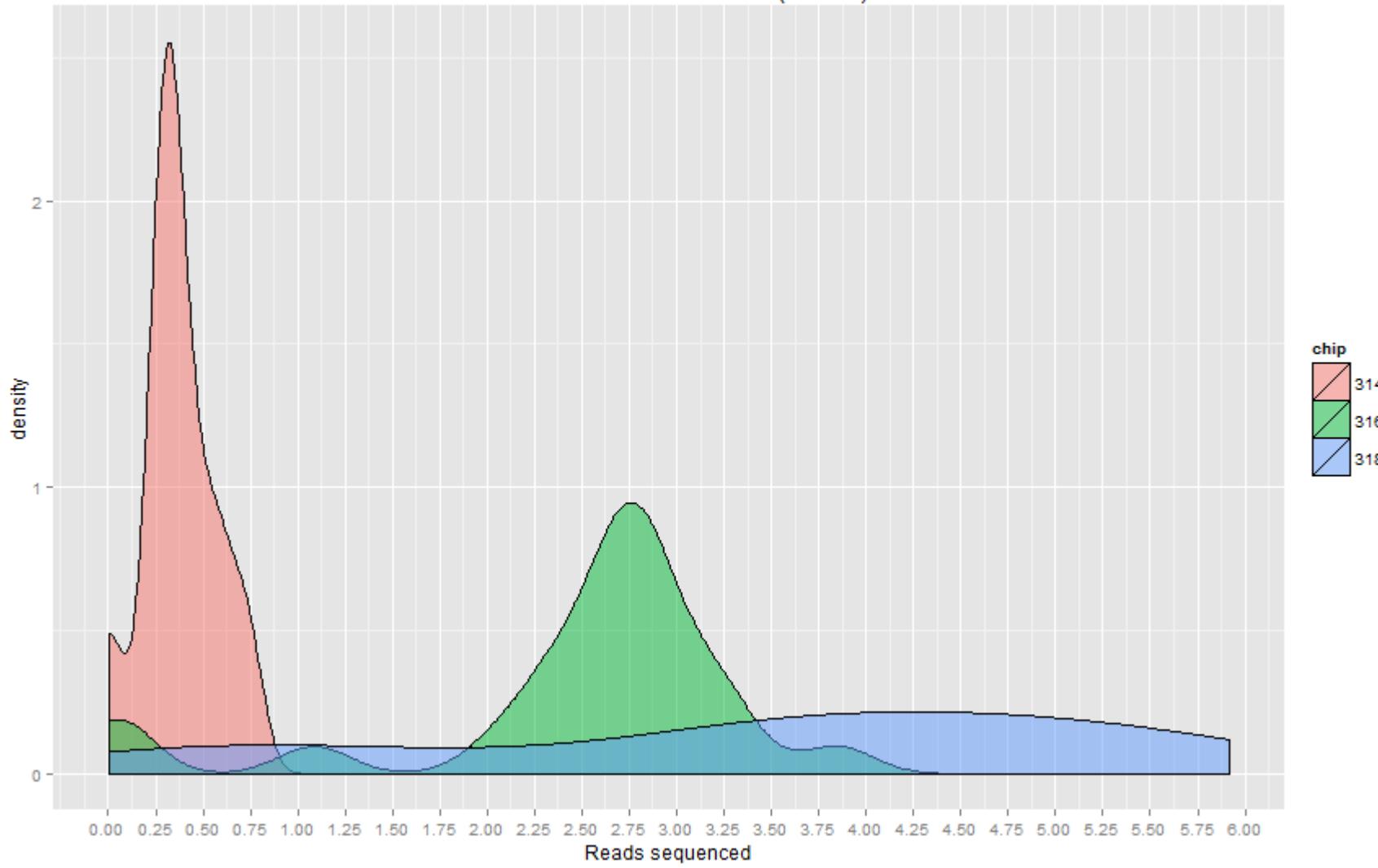


Статистика – 54 чипа

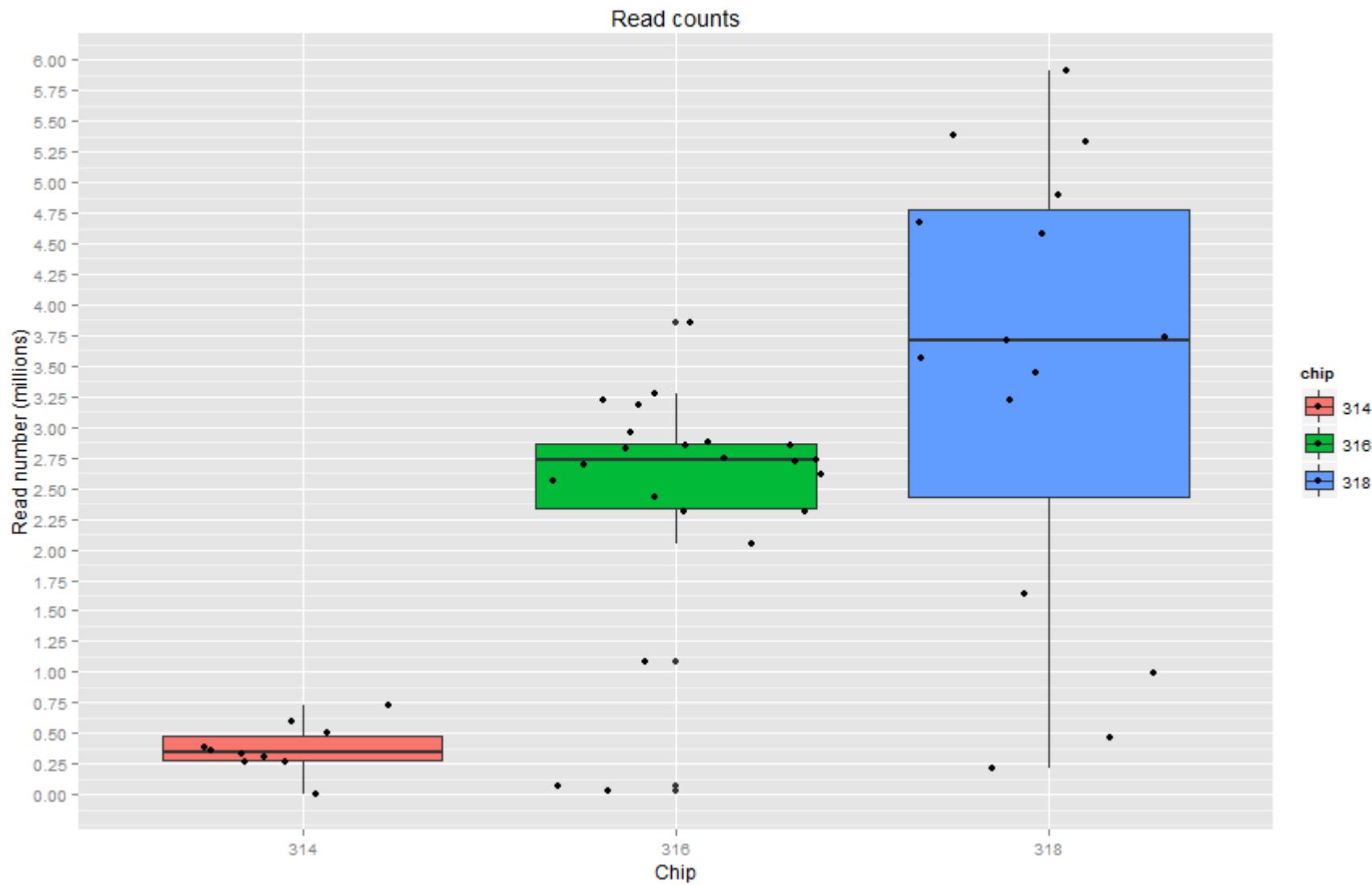


Кол-во ридов

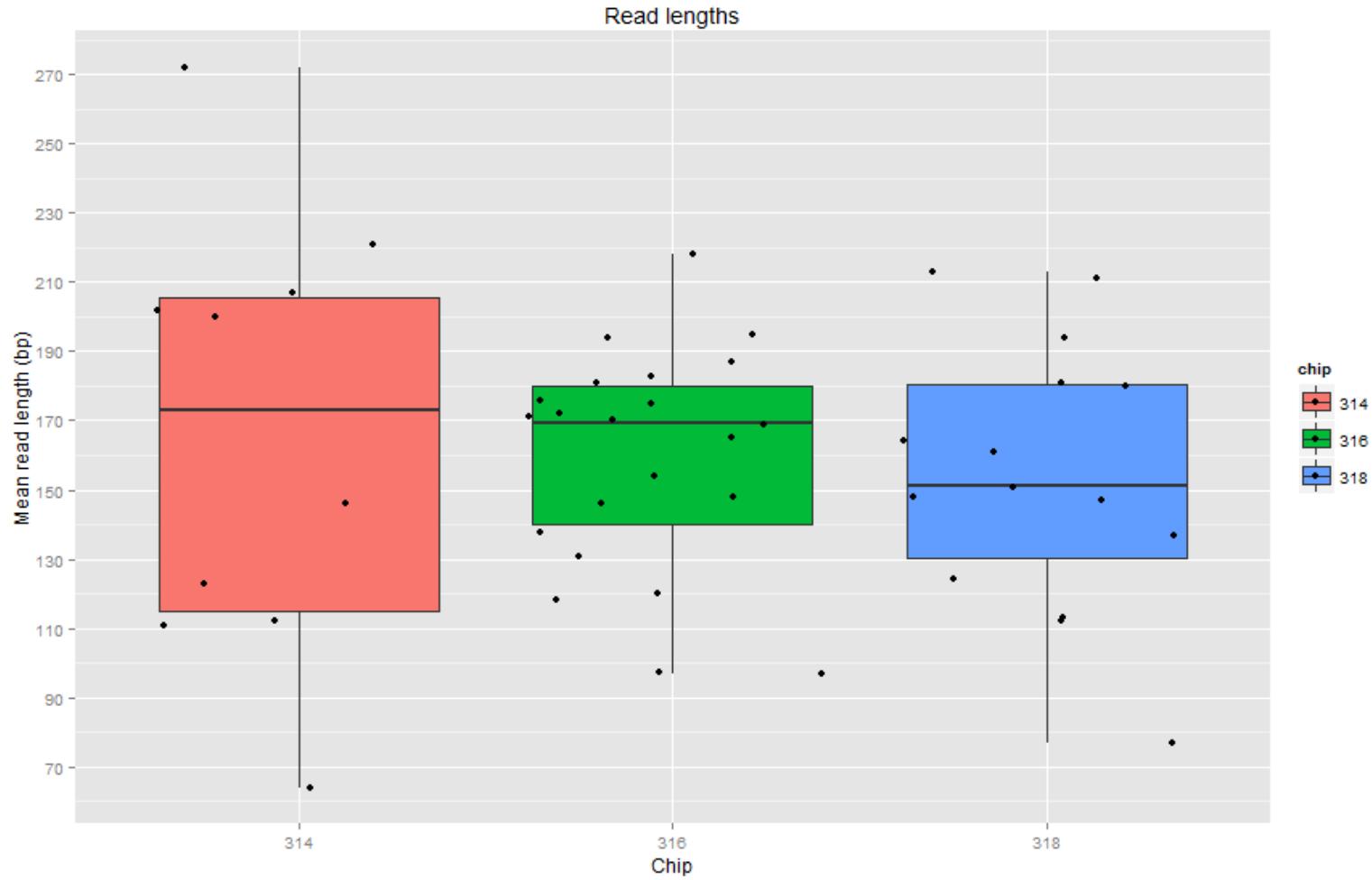
Distribution of reads numbers (millions)



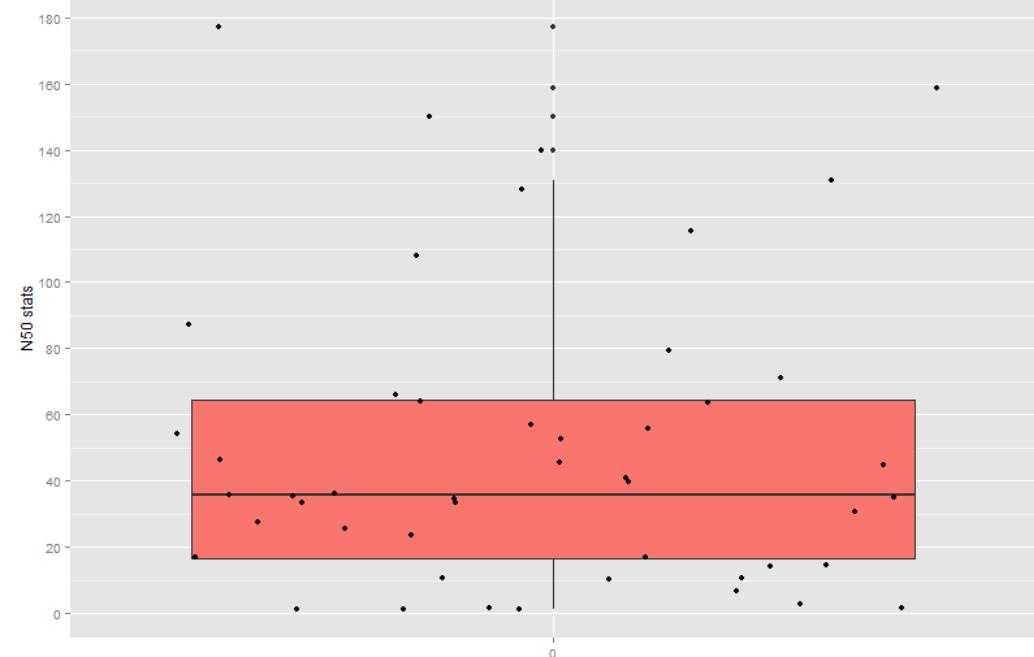
Кол-во ридов



Длины ридов



N50 in Kbp

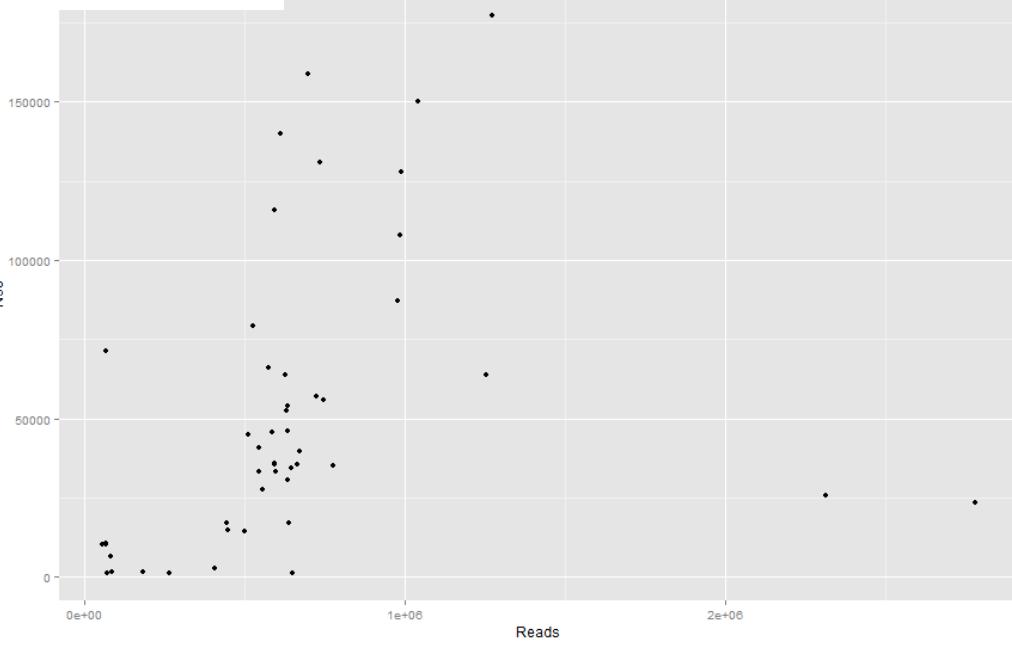


N50

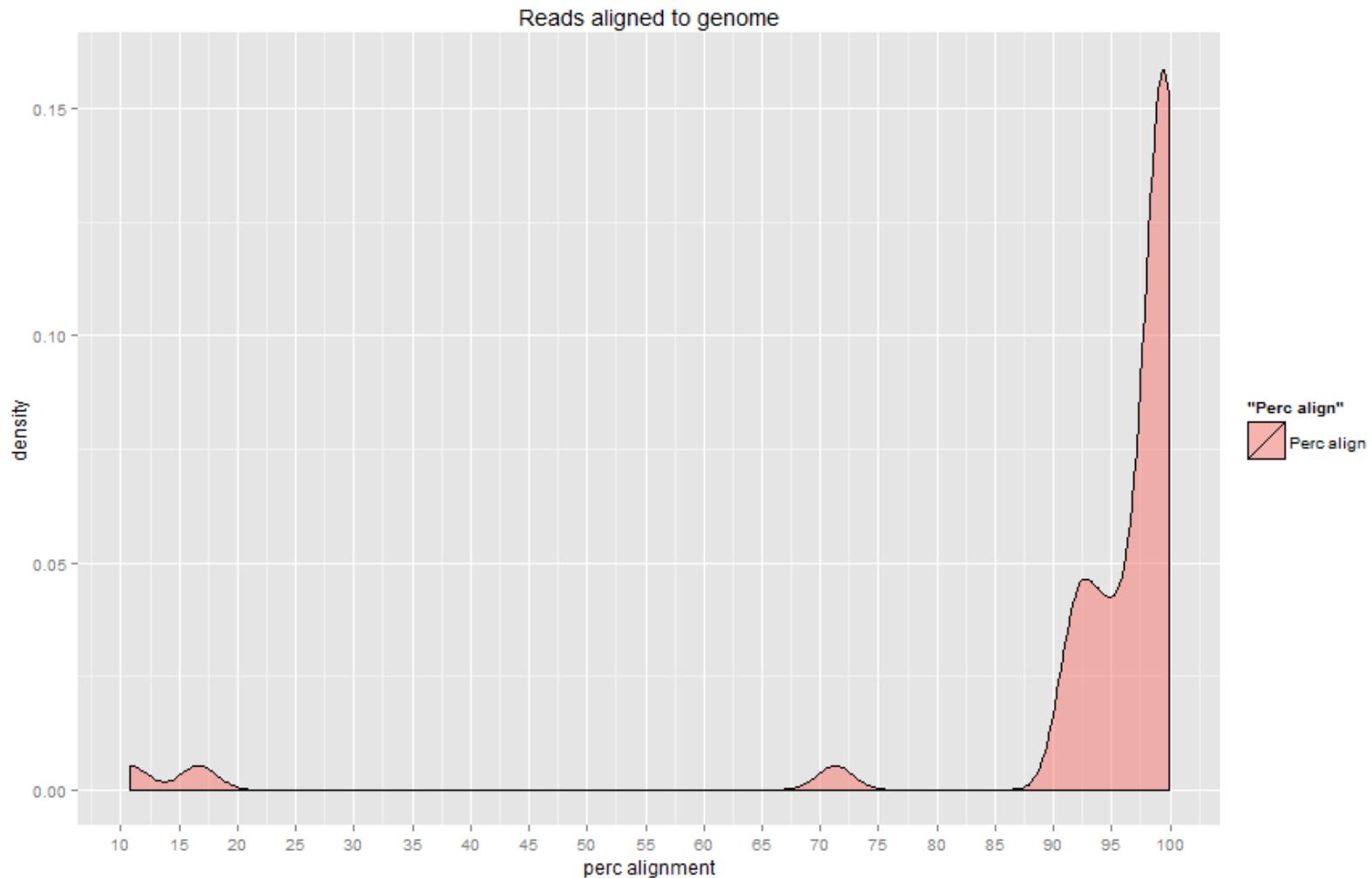
"N50"
N50



N50



Риды картировавшиеся на сборку (mira)

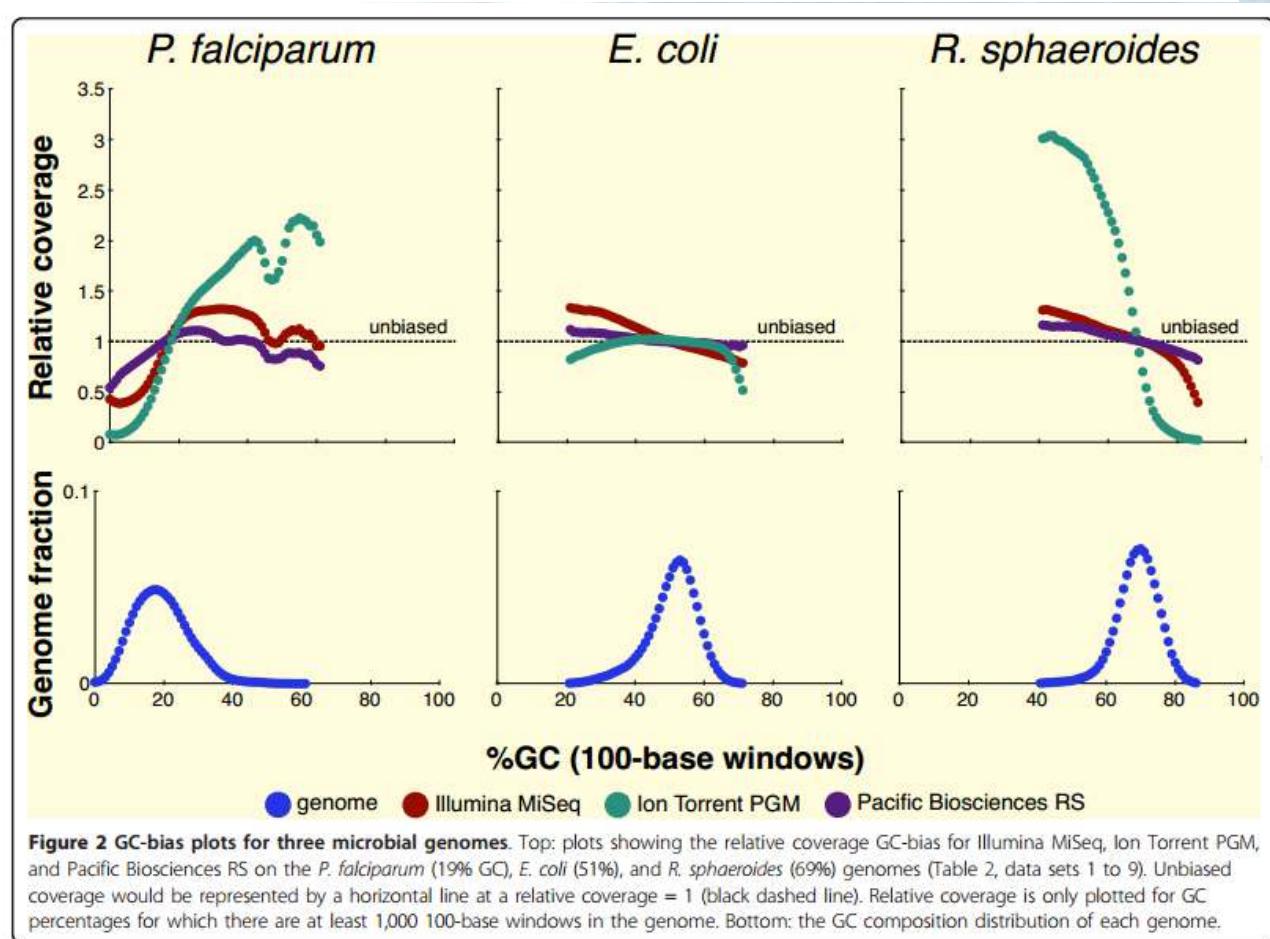


RESEARCH

Open Access

Characterizing and measuring bias in sequence data

Michael G Ross*, Carsten Russ, Maura Costello, Andrew Hollinger, Niall J Lennon, Ryan Hegarty, Chad Nusbaum and David B Jaffe



Парные фрагменты - SSPACE

- Позволяет удлинять и соединять (скаффолдинг) контиги на основе одной и более библиотек парных ридов

BIOINFORMATICS APPLICATIONS NOTE

Vol. 27 no. 4 2011, pages 578–579
doi:10.1093/bioinformatics/btq683

Genome analysis

Advance Access publication December 12, 2010

Scaffolding pre-assembled contigs using SSPACE

Marten Boetzer^{1,2}, Christiaan V. Henkel³, Hans J. Jansen³, Derek Butler¹ and Walter Pirovano^{1,*}

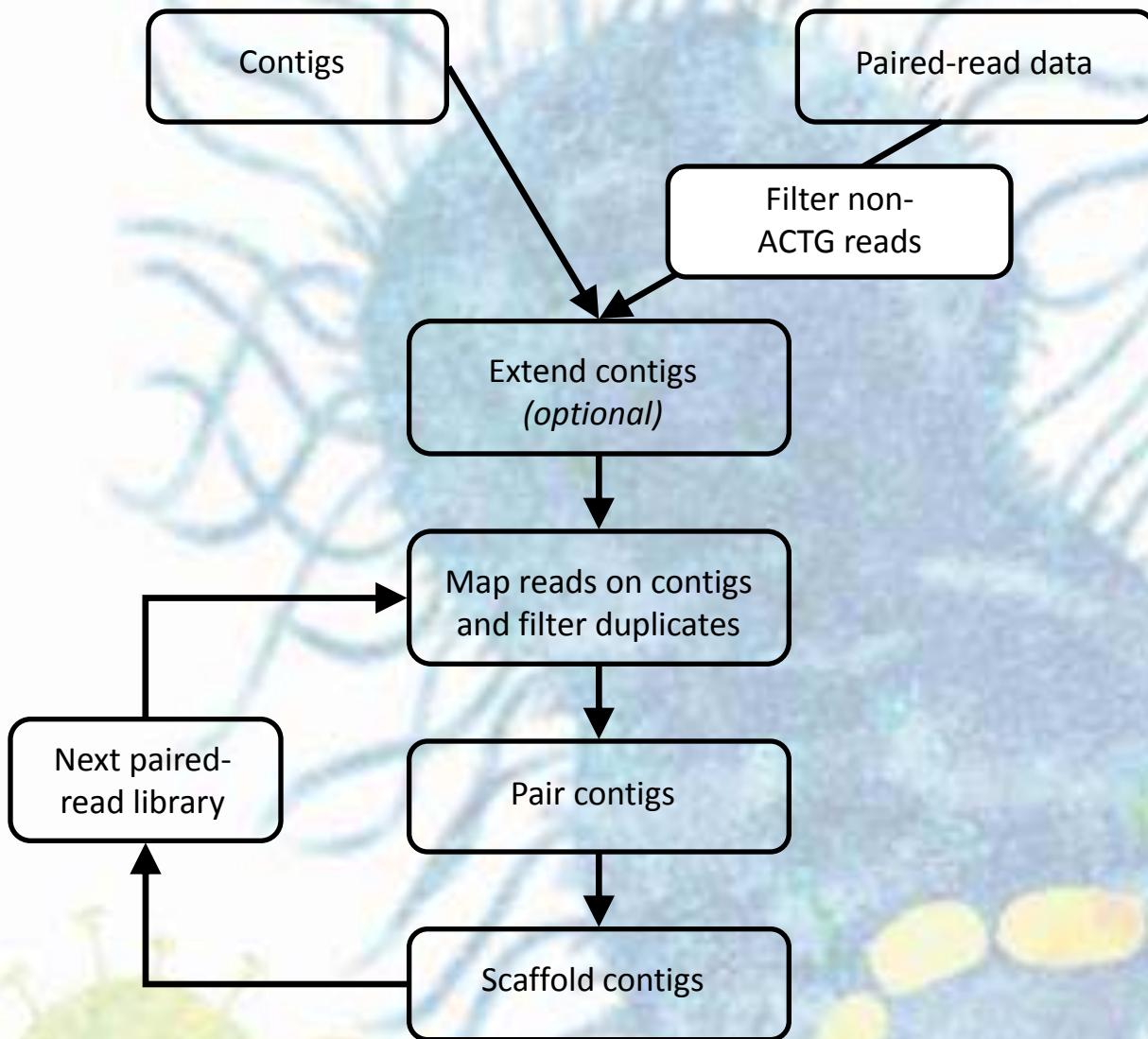
¹BaseClear B.V., Einsteinweg 5, 2333 CC Leiden, ²Leiden Institute for Advanced Computer Science, Leiden University, Niels Bohrweg 11, 2333 CA Leiden and ³ZF-screens B.V., Niels Bohrweg 11, 2333 CA Leiden, The Netherlands
Associate Editor: John Quackenbush

ABSTRACT

Summary: De novo assembly tools play a main role in genome sequencing projects. One of the major challenges in de novo assembly is the scaffolding of pre-assembled contig sequences.

To date only a few programs are able to scaffold pre-assembled contig sequences. A commonly used tool is Bambus (Pop *et al.*, 2004), although it was not designed for the current generation of

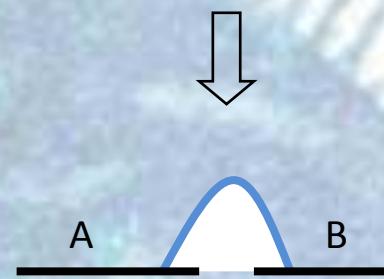
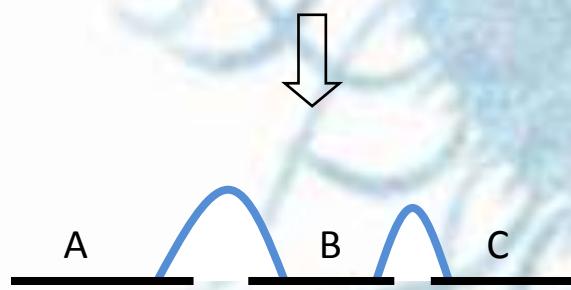
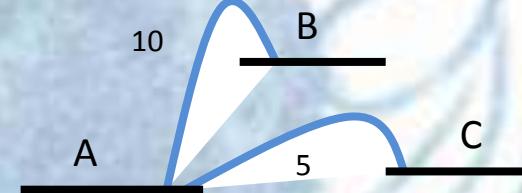
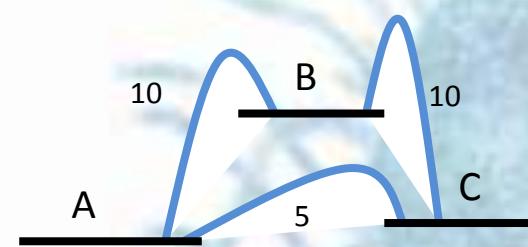
SSPACE algorithm



SSPACE algorithm

- Вероятные пары контигов определяются на основе картирования парных ридов на сборку
- Картирование считается успешным, если расстояние между ридами удовлетворяет указанному при запуске диапазону
- Скаффолдинг начинается с наиболее длинного контига. Если число связей между ним и другим контигом превышает порог, то они объединяются в скаффолд (итерационный процесс)
- В случае наличия альтернативных путей ($A \leftrightarrow B$, $A \leftrightarrow C$) вычисляется отношение числа связей между ними. Если оно меньше, чем порог (0.7) выбирается лучший путь. Если больше – скаффолдинг прерывается

Альтернативные пути скаффолдинга



(a)

(b)

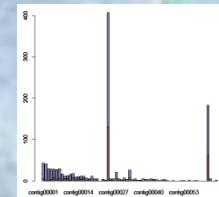
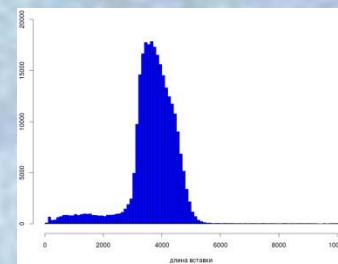
Представление скаффолдов: .dot файл, dotty



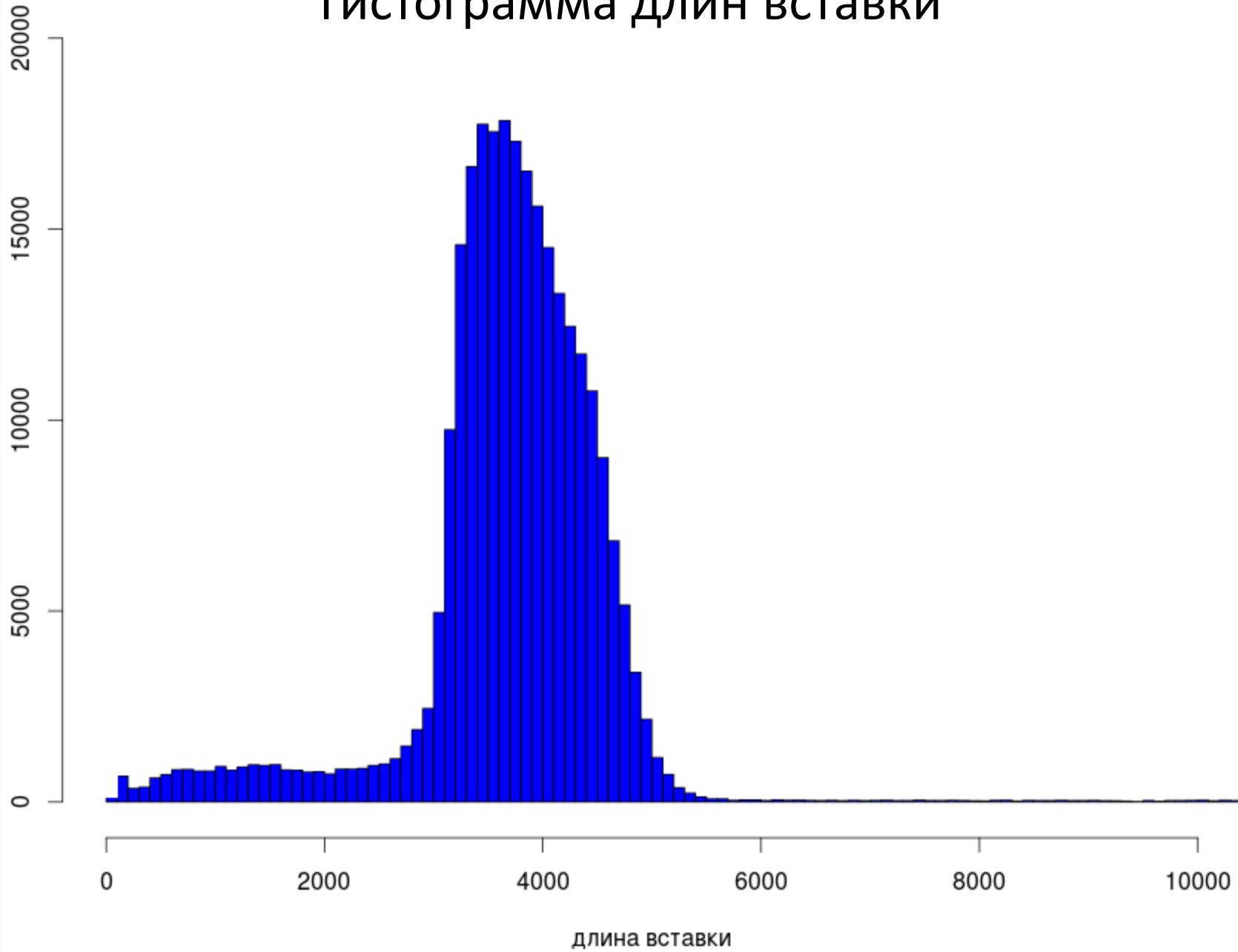
Пример статистики скаффолдинга

- READING READS r2b1: Total inserted pairs = 437,615
- READING READS r3b1: Total inserted pairs = 460,399
- Inserted contig file;
 - Total number of contigs = 135
N50 = 41,390
- After scaffolding r2b1:
 - Total number of scaffolds = 51
 - N50 = 143,152
- After scaffolding r3b1:
 - Total number of scaffolds = 25
 - N50 = 321,140

Пайплайн для сбора статистики

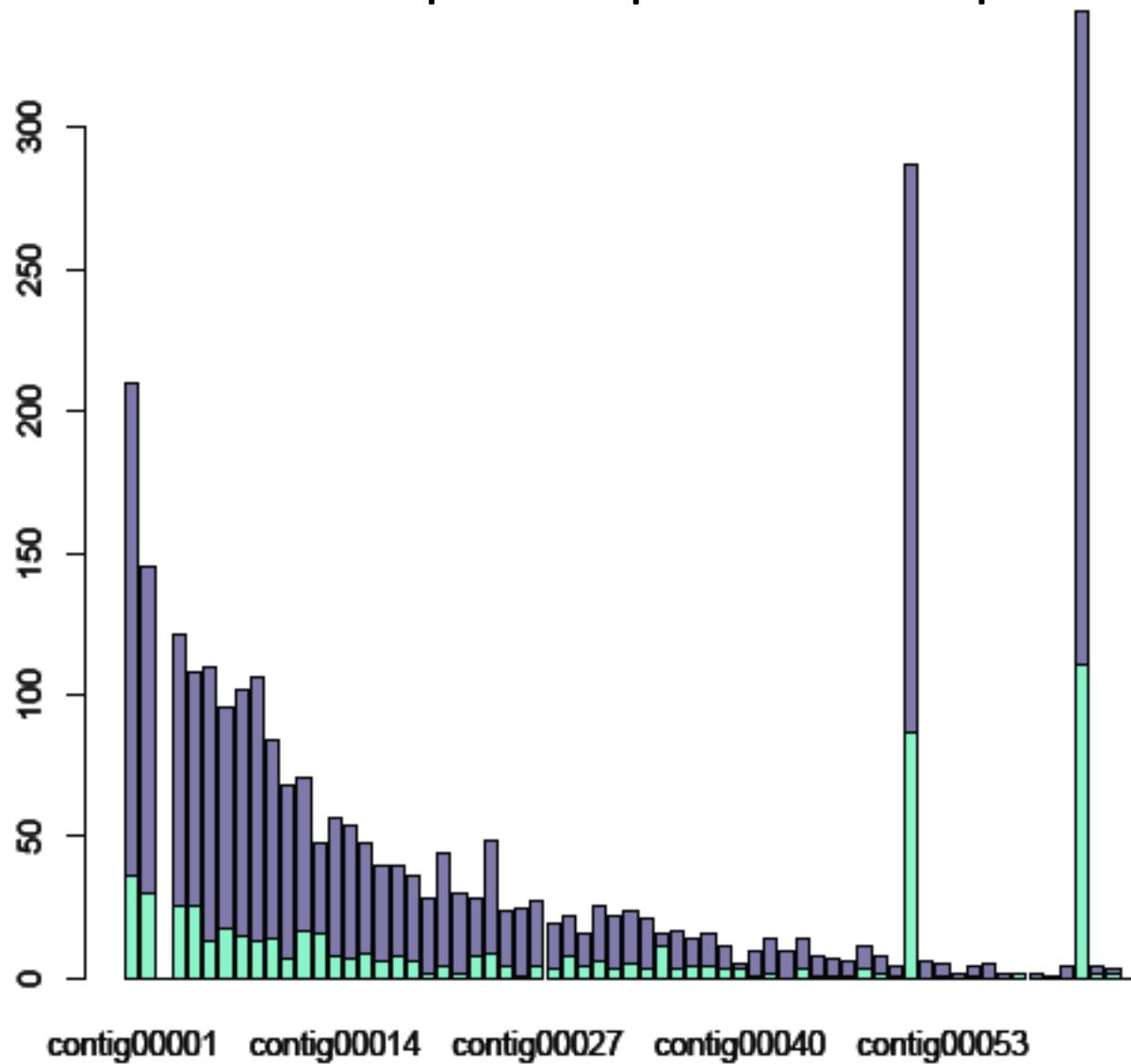


Гистограмма длин вставки

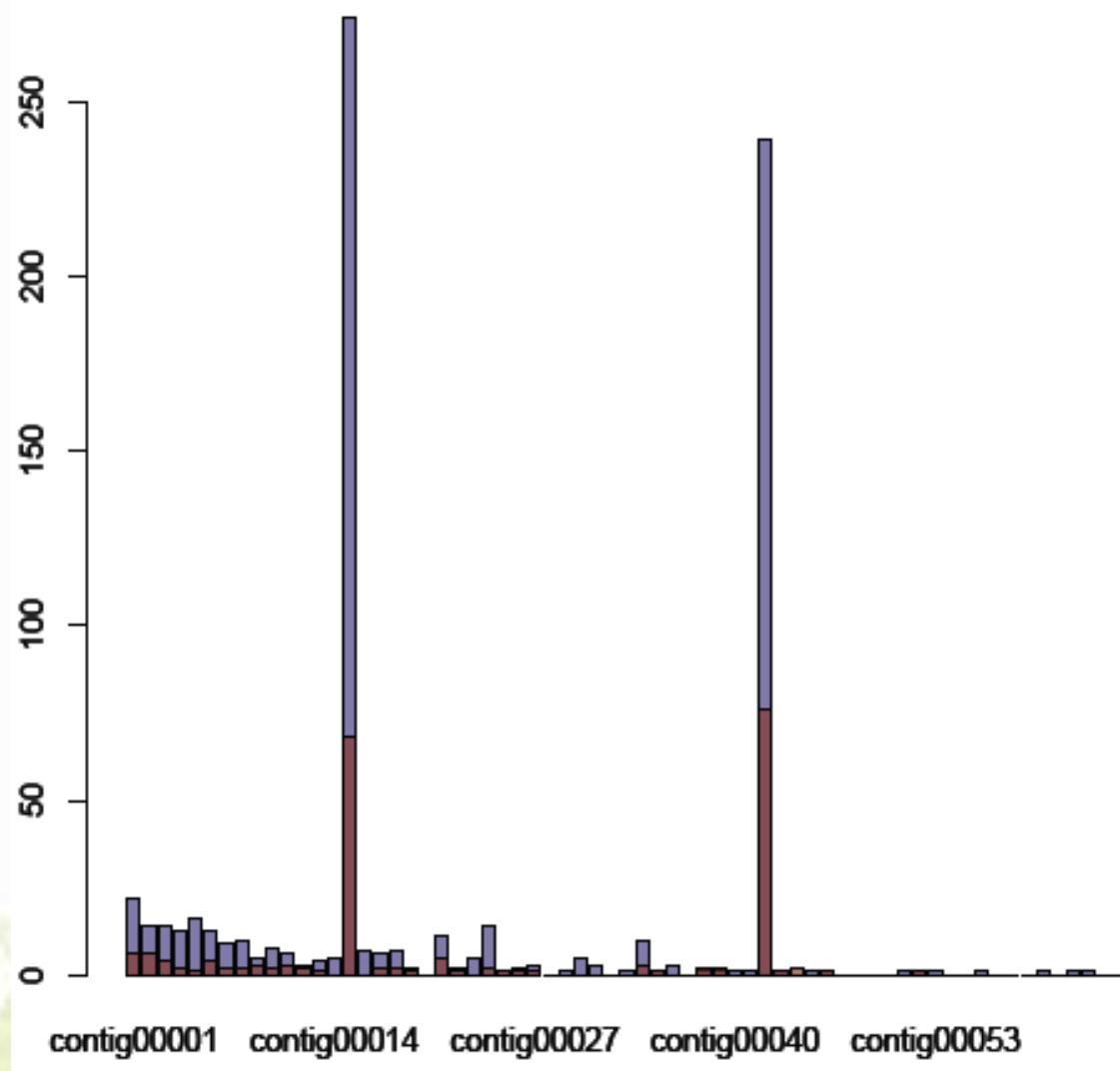


Contig 3: определить пары просто

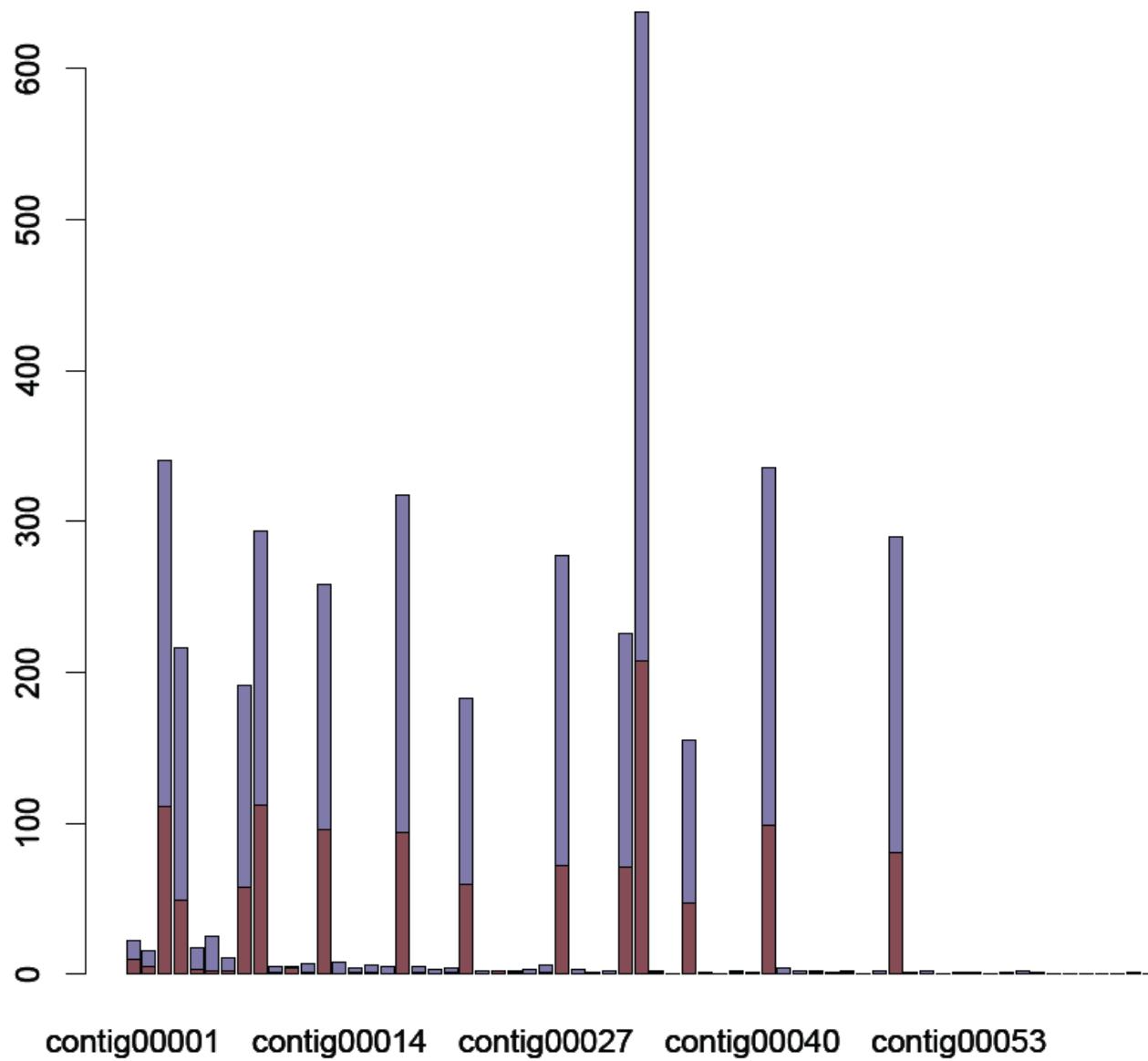
Биологические повторы - хорошая воспроизводимость



Contig 37: определить пары просто



Contig 62: определить пары сложно



Скаффолдинг сборок 454 и Ion Torrent

- *Pseudomonas stutzeri*, 454+Ion Torrent
 - Total number of contigs = 91
 - N50 = 130,512
 - Total number of scaffolds = 5
 - N50 = 2,709,945
 - Max scaffold size: 2,709,945
- *Pseudomonas stutzeri*, Ion Torrent
 - Total number of contigs = 369
 - N50 = 47,929
 - Total number of scaffolds = 74
 - N50 = 377,361
 - Max scaffold size: 1,618,709

Скаффолдинг сборок 454 и Ion Torrent

- *Enterococcus faecium*, 454
 - Total number of contigs = 70
 - N50 = 95,885
 - Total number of scaffolds = 15
 - N50 = 441,245
 - Max scaffold size: 1,575,801
- *Enterococcus faecium*, Ion Torrent
 - Total number of contigs = 156
 - N50 = 100,114
 - Total number of scaffolds = 50
 - N50 = 1,960,091
 - Max scaffold size: 1,960,091

M.tuberculosis
454 FLX
149 контигов
N50 = 101797

Библ парных
фрагментов
(2-3 тыс + 4-5 тыс)
(под разными
баркодами)

49 scaffolds
N50 = 494662

M.tuberculosis
454 FLX
165 контигов
N50 = 88489

Библ парных
фрагментов
(2-3 тыс + 4-5 тыс)
(под разными
баркодами)

35 scaffolds
N50 = 738326



Тема:

Расшифровка механизмов формирования детерминант патогенности стрептококками группы *viridans* на основании сравнения полноразмерных нуклеотидных последовательностей их геномов

Проблема:

Прояснение молекулярных механизмов формирования патогенного потенциала клинически значимыми микроорганизмами рода *Streptococcus*, в частности *S. pneumoniae* (пневмококка), вызывающего у человека пневмонию, сепсис, менингит

Задачи:

- Геномное секвенирование клинических изолятов зеленящих (*viridans*) стрептококков, отличных по видовой принадлежности (*S. mitis* и *S. pneumoniae*) и клиническим проявлениям со стороны пациента (инвазивные и неинвазивные *S. pneumoniae*)
- Инвентаризация и систематизация известных клеточных компонентов и структур (ферментов, протеинов клеточной стенки, капсулы, пилей и др.), обычно рассматриваемых в качестве факторов патогенности и вирулентности, на основании анализа полноразмерных нуклеотидных последовательностей геномов изолятов зеленящих стрептококков подгруппы *mitis*.
- Выявление новых факторов, которые могут играть роль в адаптации и инвазии микроорганизмов рода *Streptococcus* в макроорганизме, методами сравнительной геномики.
- Выявление структурных и функциональных особенностей генома, определяющих выбор между паразитизмом и комменсализмом на примере сравнения геномов *S. pneumoniae* и *S. mitis*.

Sequencing Platform	Sequencing Depth	Assembly Method	Large Contigs
Ion Torrent PGM™	35x	GS De Novo Assembler v. 2.8	115
Ion Torrent PGM™	34x	GS De Novo Assembler v. 2.8	103
Ion Torrent PGM™	40x	GS De Novo Assembler v. 2.8	107
Ion Torrent PGM™	31x	GS De Novo Assembler v. 2.8	124
Ion Torrent PGM™	20x	GS De Novo Assembler v. 2.8	226
Ion Torrent PGM™	26x	GS De Novo Assembler v. 2.8	180
Ion Torrent PGM™	19x	GS De Novo Assembler v. 2.8	208
Ion Torrent PGM™	55x	GS De Novo Assembler v. 2.8	88
Ion Torrent PGM™	75x	GS De Novo Assembler v. 2.8	60
Ion Torrent PGM™	48x	GS De Novo Assembler v. 2.8	110
Ion Torrent PGM™	74x	GS De Novo Assembler v. 2.8	37
Ion Torrent PGM™	53x	GS De Novo Assembler v. 2.8	67
Ion Torrent PGM™	77x	GS De Novo Assembler v. 2.8	34

К настоящему моменту прочитаны геномы 8 штаммов пневмококков (4 инвазивных, 4 неинвазивных), 5 штаммов *S. mitis* и одного штамма псевдопневмококка. Все геномы можно найти в базе данных NCBI

Вид	Штамм	Accession в NCBI
<i>Streptococcus pneumoniae</i>	357	PRJNA201317
<i>Streptococcus pneumoniae</i>	2009	PRJNA201318
<i>Streptococcus pneumoniae</i>	801	PRJNA201319
<i>Streptococcus pneumoniae</i>	845	PRJNA201320
<i>Streptococcus pneumoniae</i>	1488	PRJNA201321
<i>Streptococcus pneumoniae</i>	1542	PRJNA201322
<i>Streptococcus pneumoniae</i>	3051	PRJNA201323
<i>Streptococcus pneumoniae</i>	1779	PRJNA206047
<i>Streptococcus mitis</i>	11/5	PRJNA201324
<i>Streptococcus mitis</i>	13/39	PRJNA201325
<i>Streptococcus mitis</i>	17/34	PRJNA206048
<i>Streptococcus mitis</i>	18/56	PRJNA206049
<i>Streptococcus mitis</i>	29/42	PRJNA206050
<i>Streptococcus pseudopneumoniae</i>	G42	в процессе оформления

Первичная характеристика штаммов по геномным данным: «виртуальное серотипирование» *Streptococcus pneumoniae*

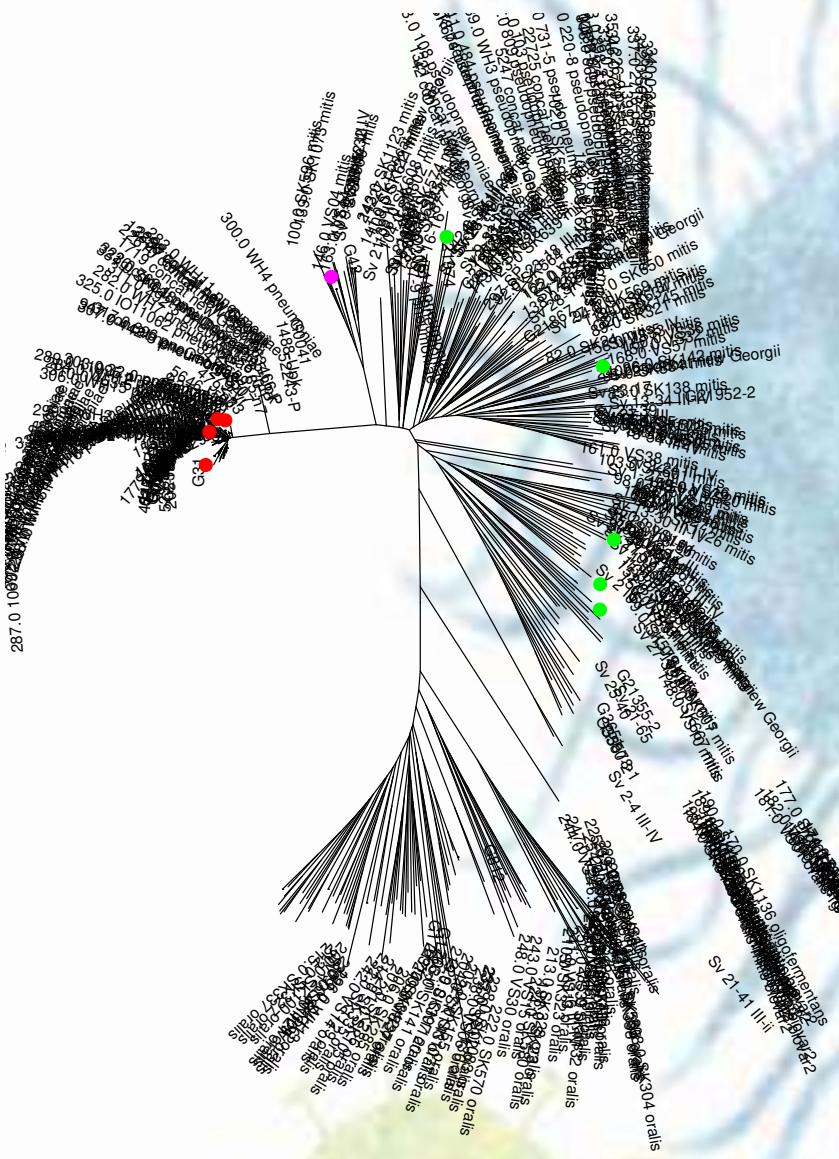
Рекомендации CDC по определению 40 серотипов с помощью сероспецифичной мультипраймерной ПЦР



S38					
	A	B	C	D	E
1	query id	subject id	% identity	alignmentlength	mismatches
2	ser_2	Sp_SPAR27g 39586194gbALCC01000004.1	100.00	289	0
3	ser_4	Sp_SPAR4g 193804931gbAE005672.3	100.00	430	0
4	ser_4	Sp_SPAR4g 395870110gbALCE01000001.1	100.00	430	0
5	ser_4	Sp_2082170g 395615788gbALBO01000003.1	100.00	430	0
	ser_4	Sp_2081074g 395611663gbALBM01000005.1	100.00	430	0
7	ser_6A_6B_6C_6D	Sp_SPAR55g 395877268gbALCF01000001.1	99.20	250	2
8	ser_6A_6B_6C_6D	Sp_SPAR55g 395884239gbALCK01000002.1	98.40	250	4
9	ser_6A_6B_6C_6D	Sp_GA60190g 395888325gbALCL01000003.1	98.00	250	5
10	ser_6A_6B_6C_6D	Sp_GA60132g 395909931gbALCV01000001.1	98.00	250	5
11	ser_6A_6B_6C_6D	Sp_GA60080g 395903898gbALCR01000010.1	98.00	250	5
12	ser_6A_6B_6C_6D	Sp_07AR0125g 328694797gbAFBY01000240.1	98.00	250	5
13	ser_6C_6D	Sp_GA60190g 395888325gbALCL01000003.1	100.00	727	0
14	ser_6C_6D	Sp_07AR0125g 328694797gbAFBY01000240.1	100.00	727	0
15	ser_6C_6D	Sp_GA60132g 395909931gbALCV01000001.1	99.86	727	1
16	ser_6C_6D	Sp_GA60080g 395903898gbALCR01000010.1	99.86	727	1
17	ser_7C_7B_70	Sp_SPAR95g 395872337gbALCD01000001.1	99.62	260	0
18	ser_7F_7A	Sp_2070109g 395578353gbALBA01000026.1	100.00	599	0
19	ser_7F_7A	Sp_2070108g 395581726gbALAZ01000006.1	100.00	599	0
20	ser_8	Sp_2071247g 395605107gbALBK01000004.1	100.00	201	0
21	ser_8	Sp_2081685g 395611008gbALBN01000001.1	99.50	201	0
22	ser_9V_9A	Sp_GA56113g 395896137gbALCP01000002.1	100.00	816	0
23	ser_9V_9A	Sp_GA17301g 395882932gbALCI01000001.1	100.00	816	0
24	ser_9V_9A	Sp_2070531g 3955902975gbALBD01000005.1	100.00	816	0
25	ser_9V_9A	Sp_2070425g 395590755gbALBC01000004.1	100.00	816	0
26	ser_10A	Sp_2080076g 395597163gbALBH01000002.1	100.00	628	0
27	ser_10A	Sp_2070035g 395575797gbALAY01000004.1	100.00	628	0
28	ser_11A_11D	Sp_AP200g 1306408173gCP02121.1	100.00	463	0
29	ser_12E_12A_44_46	Sp_CDCC288-049 158031652gbARGE01000001.1	99.73	376	1

- 8 наших геномов *S. pneumoniae*
 - 64 генома из базы данных NCBI

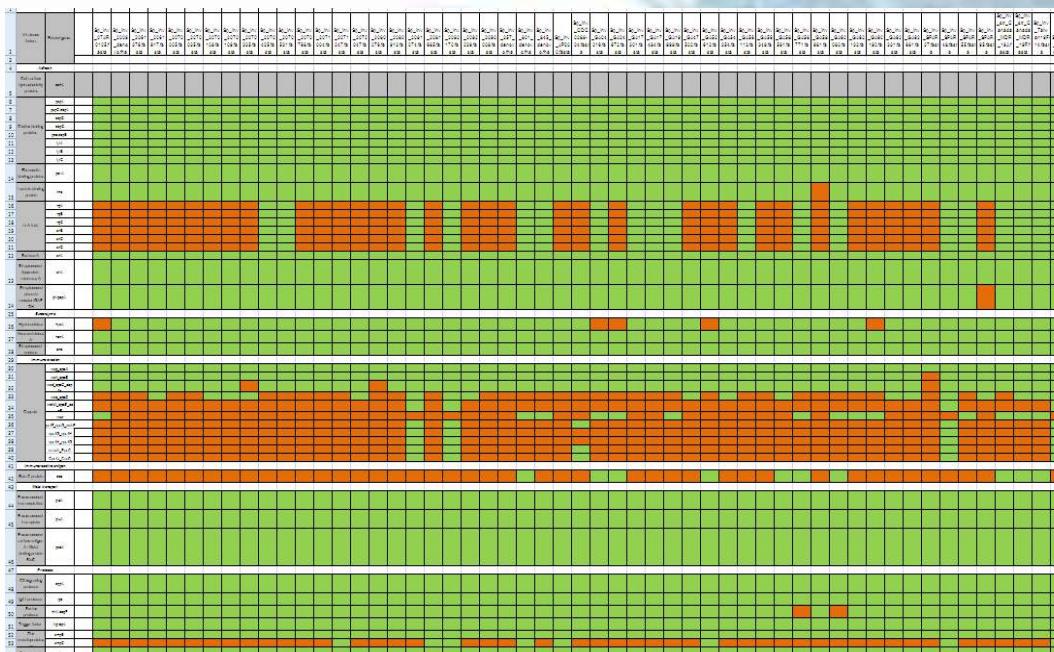
Первичная характеристика штаммов по геномным данным: филогенетика



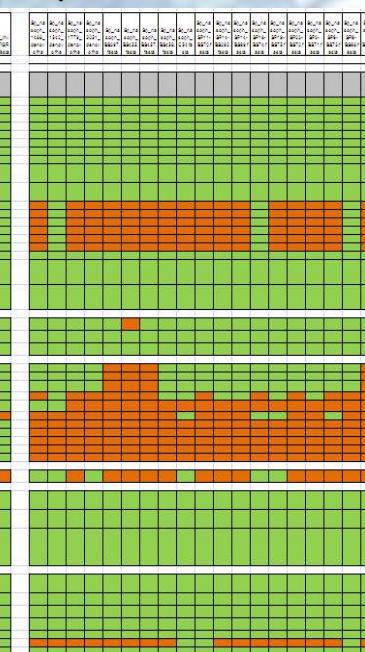
- 427 штаммов из базы данных MLSA *S. viridans*
- 64 генома из базы данных NCBI
- 14 геномов из нашей коллекции

Сравнительный анализ факторов вирулентности штаммов подгруппы mitis группы зеленящих стрептококков

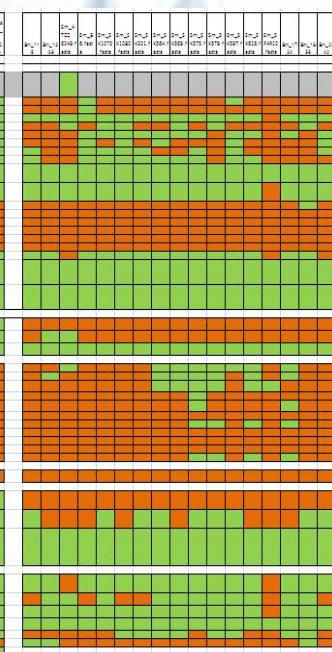
S. pneumoniae (инвазивные)



S. pneumoniae
(неинвазивные)



S. mitis



Структура капсулного оперона представителей группы зеленящих стрептококков

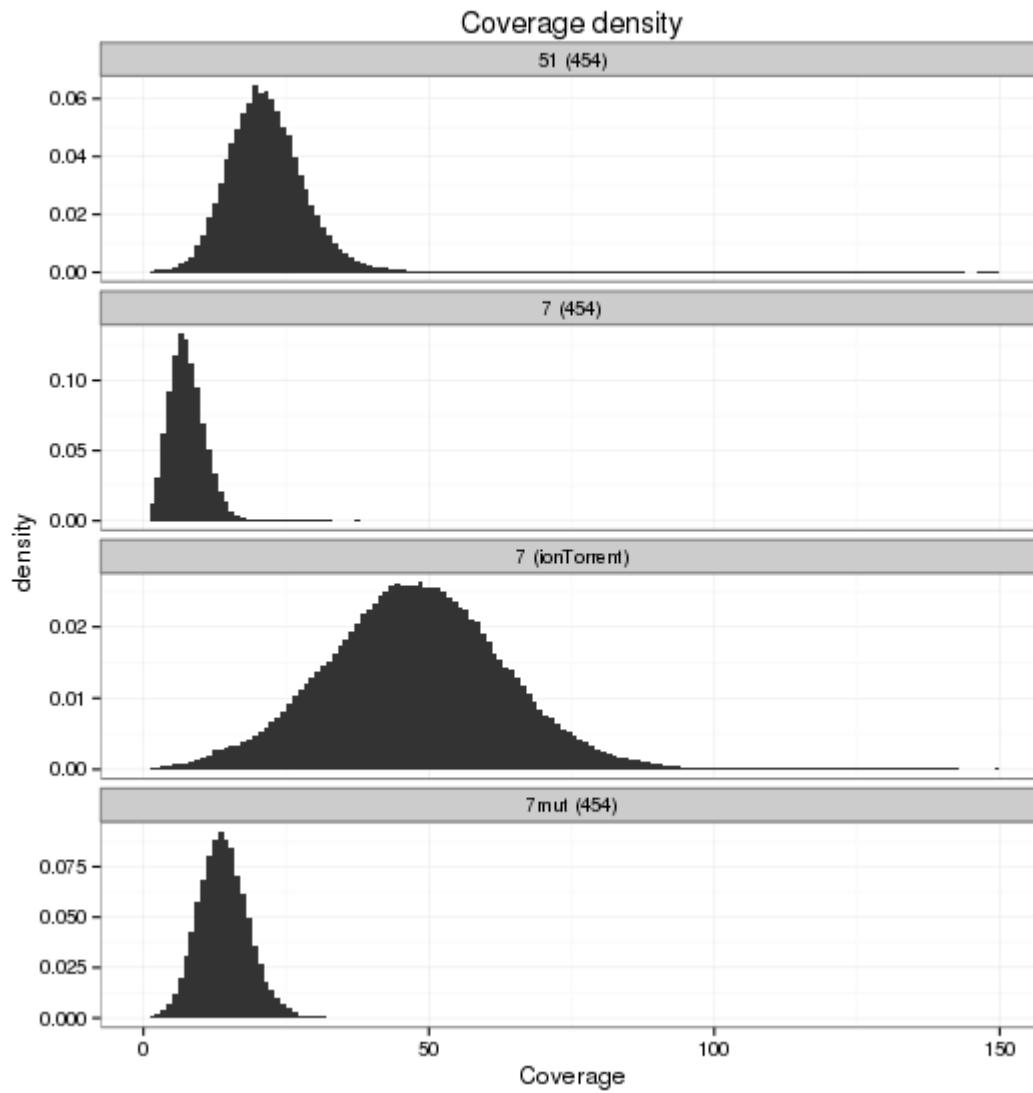
Капсула – один из основных факторов вирулентности пневмококка



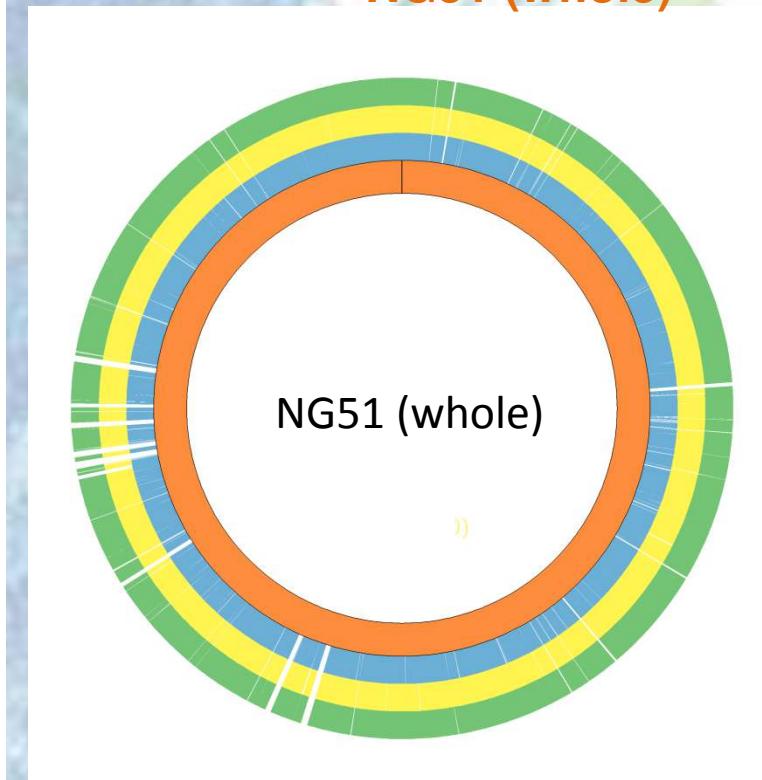
Neisseria gonorrhoeae
Клинические изоляты

	PEN	TET	CIPR	SPEC	CRO	Azi
NG3						
NG3mut						
NG19						
NG51						

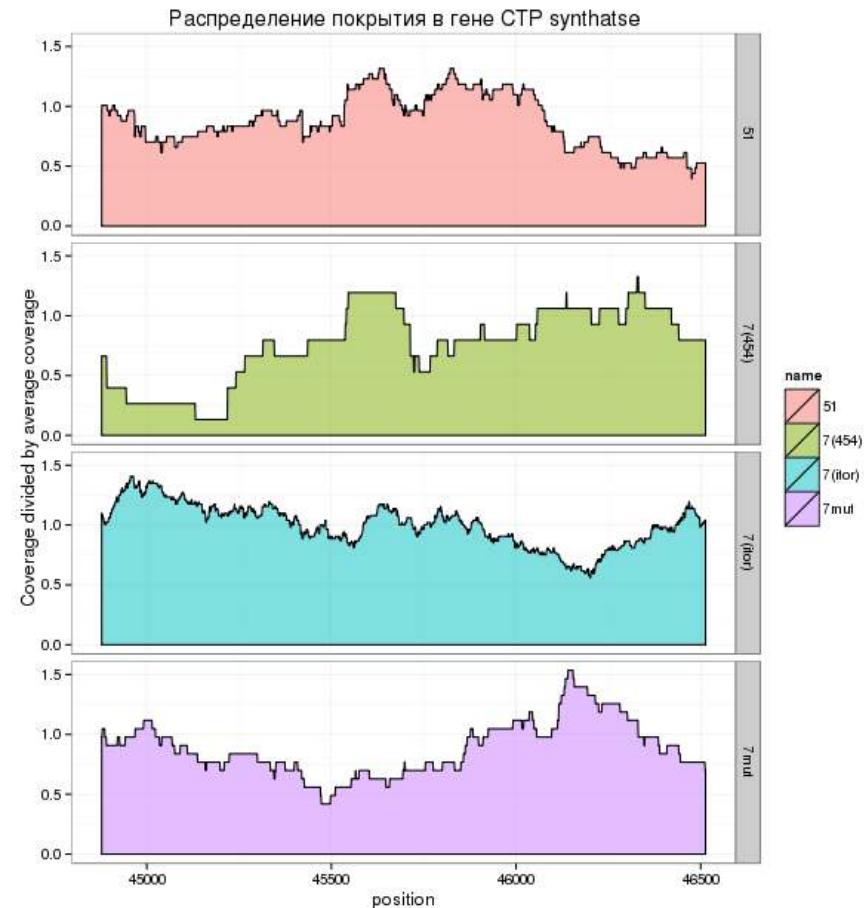
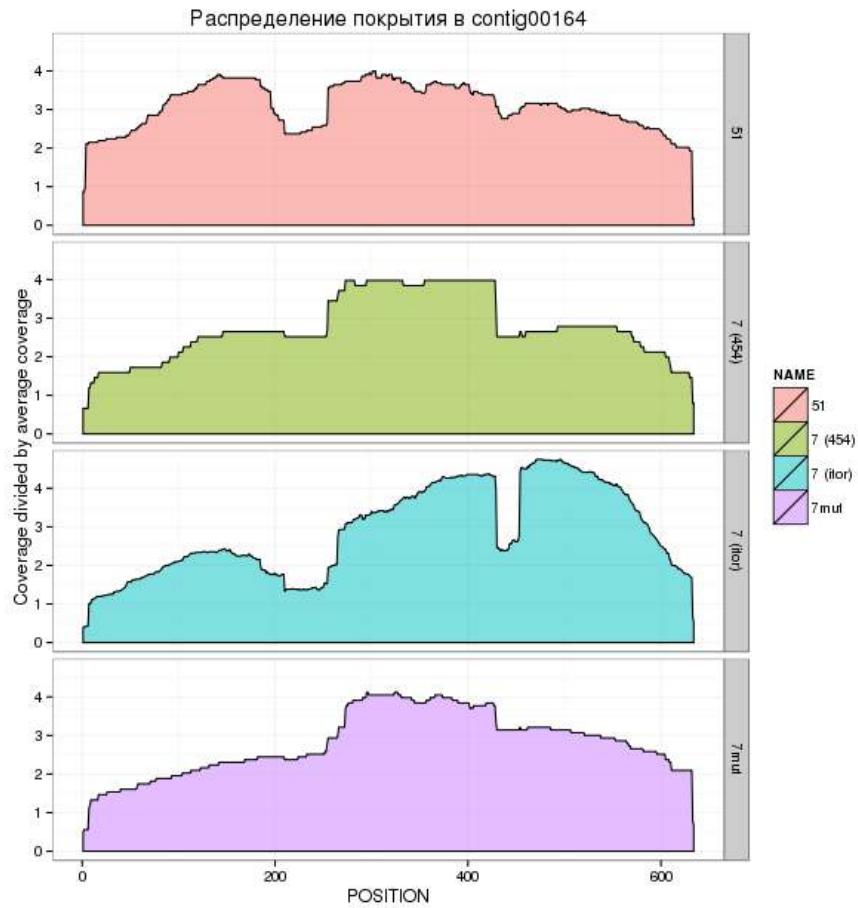
Neisseria gonorrhoeae



NG7mut (454)
NG7 (ionTorrent)
NG7 (454)
NG51 (whole)



Neisseria gonorrhoeae



Neisseria gonorrhoeae

Внутригеномные повторы

- DNA uptake sequences 5'-ATGCCGTCTGAA-3'

Штамм	Количество повторов	
	DUS10	DUS12
<i>N. gonorrhoeae</i> NG3	1849	1419
<i>N. gonorrhoeae</i> NG3mut	1905	1451
<i>N. gonorrhoeae</i> NG19	1930	1479
<i>N. gonorrhoeae</i> NG51	1913	1462
<i>N. gonorrhoeae</i> FA1090	1965	1522
<i>N. gonorrhoeae</i> NCCP11945	1966	1520
<i>N. gonorrhoeae</i> 35/02	1915	1475



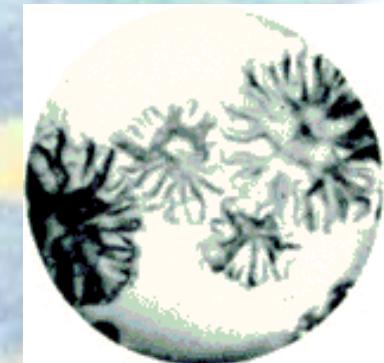
Сборка de-novo и анализ генома *Pseudomonas stutzeri KOS6*

Объект исследования: *Pseudomonas stutzeri*



- Грамотрицательные бактерии, которые широко распространены в природе
- Впервые были извлечены из спинномозговой жидкости. Условно-патогенные микроорганизмы
- Обитают в очень разнообразных средах: почва, толща морской воды и осадки, сточные воды
- Обладают очень широким спектром метаболических функций
- Модельный организм для изучения денитрификации
- Способности организма усваивать широкий спектр субстратов может сделать его применимым в биоремедиации и очистки сточных вод.
- Высокая частота естественной трансформации
- Использовались итальянским микробиологом Giancarlo Ranalli для реставрации старинных фресок

Pseudomonas stutzeri KOS6

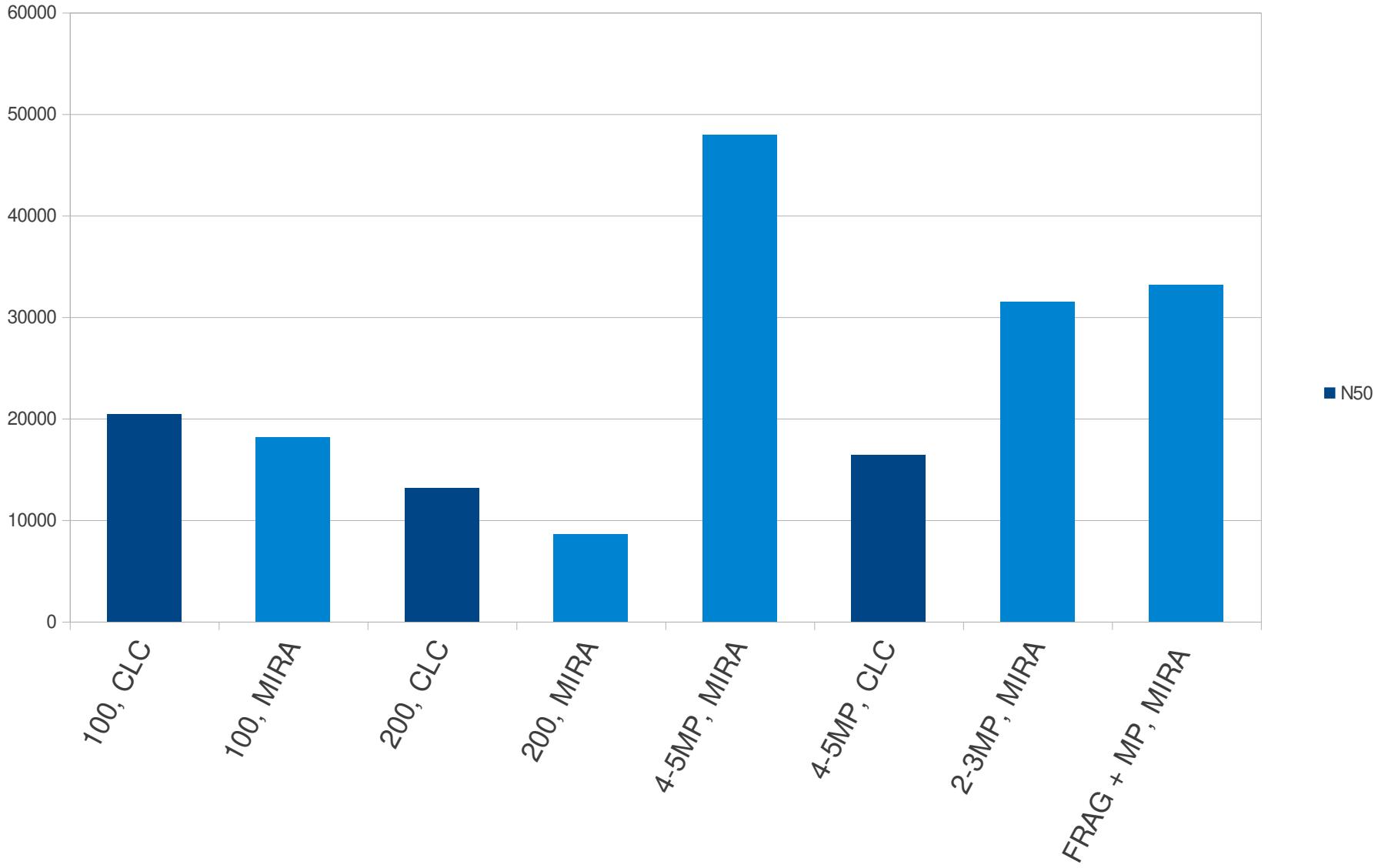


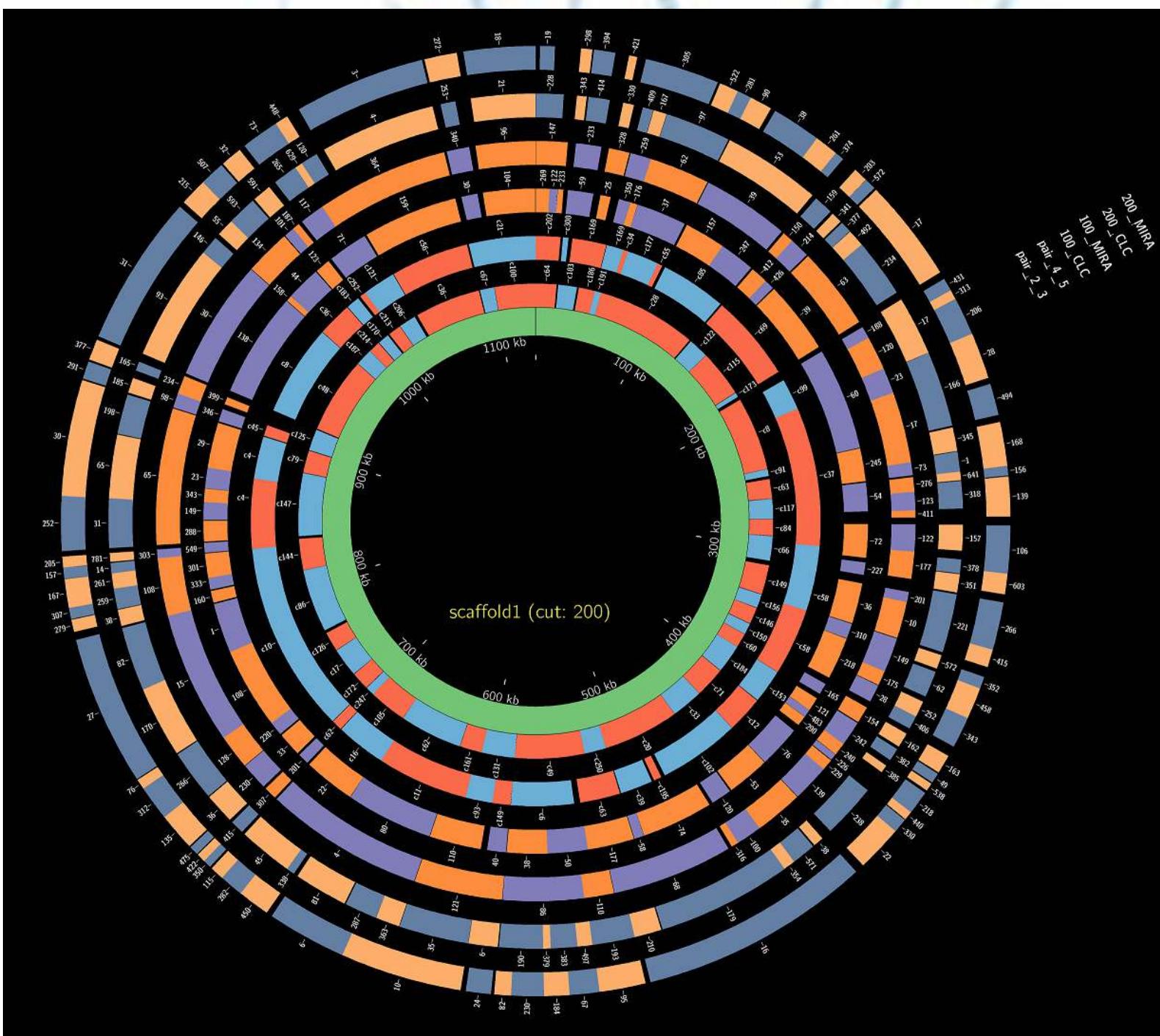
Результаты секвенирования

Ion Torrent PGM

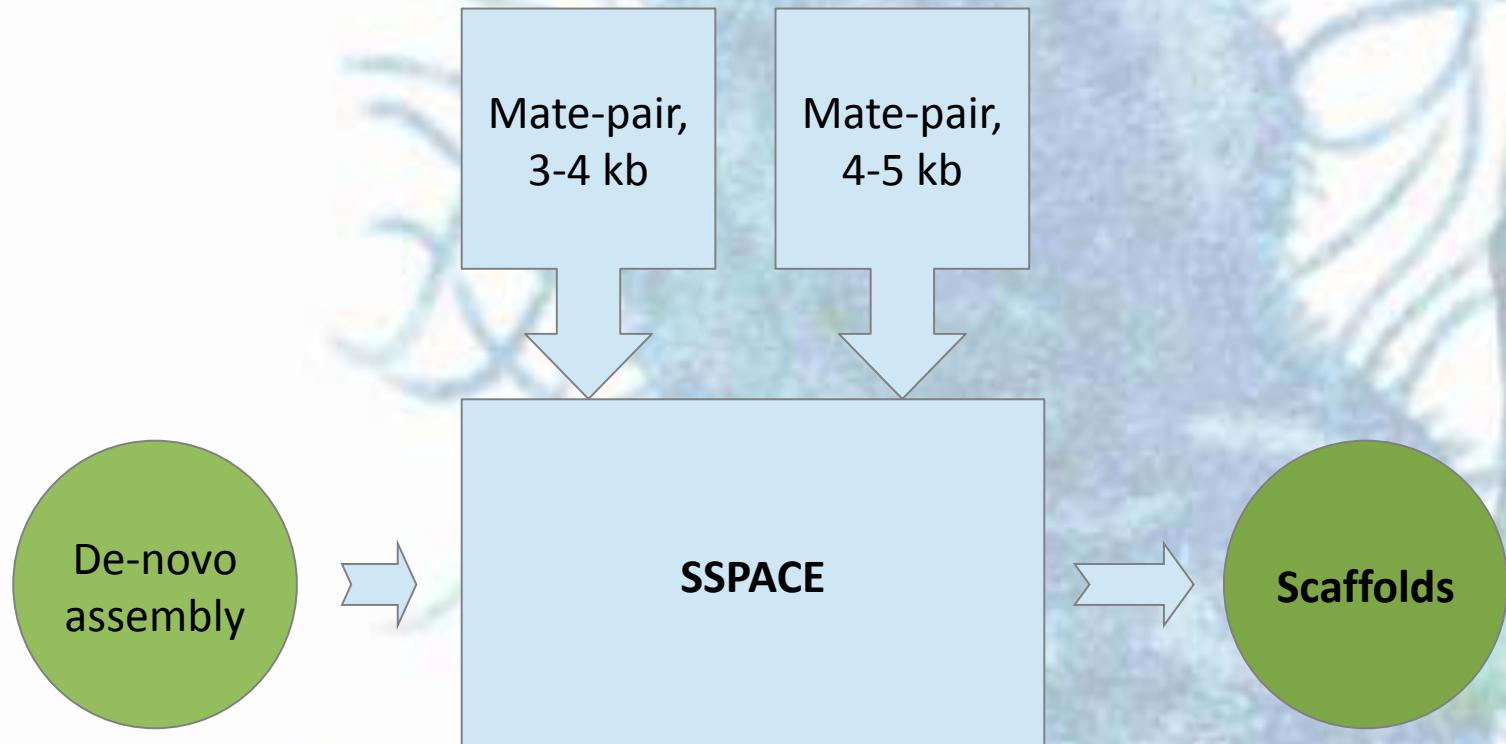
- *Две фрагментные библиотеки:*
 - 100 пн (покрытие 57x)
 - 200 пн (покрытие 135x)
- *Две парные (mate-pair) библиотеки:*
 - вставки 2.5 кб (покрытие 149x, 73% парных)
 - вставки 4.5 кб (покрытие 142x, 87% парных)

Сборка de-novo, N50





Скаффолдинг



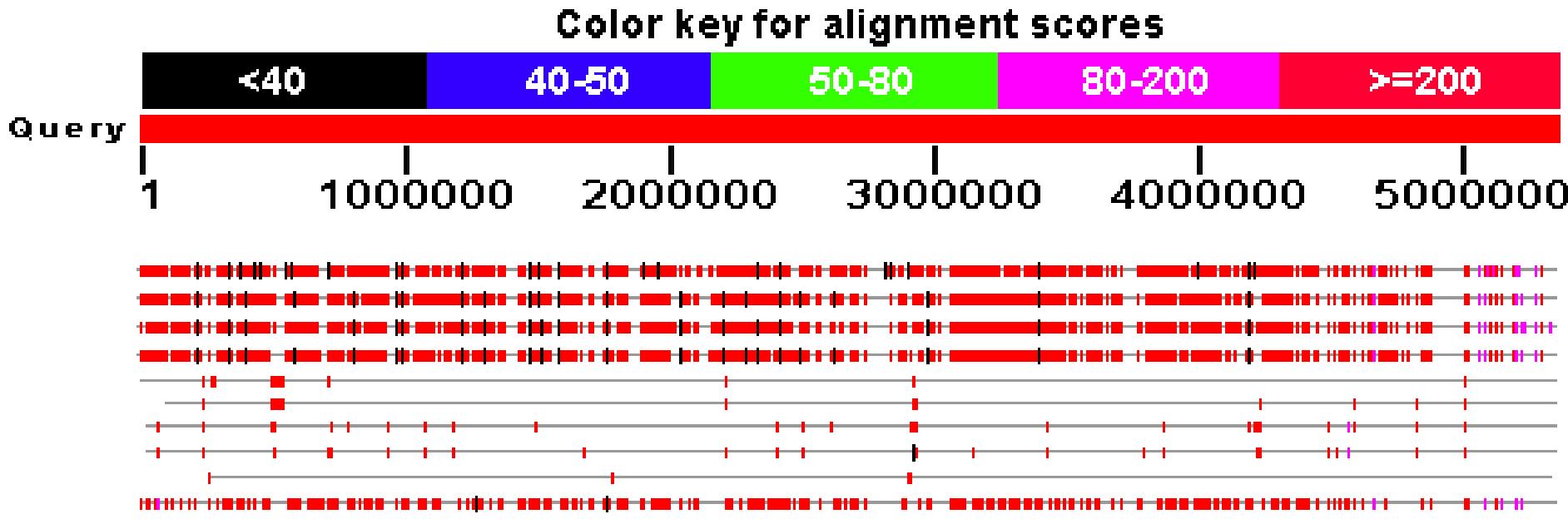
Скаффолдинг

- Mate pair 4-5 кб:
 - 74 скффолдов
- Mate pair 2-3 кб
 - 90 скффолдов
- Fragment + mate pair («ручная», Vector NTI)
 - 26 скффолдов

Анализ генома



Whole-genome BLAST



Горизонтальный перенос

- *Tolumonas auensis*

Бактерия производящая толуол. Обитает в том числе в бескислородных озерных отложениях

Толуол (от исп. Tolu, толуанский бальзам) — метилбензол, бесцветная жидкость с характерным запахом, относится к аренам. Продукт каталитического риформинга бензиновых фракций нефти.

- *Parvibaculum lavamentivorans DS-1*

Был извлечен из очистных сооружений (Германия). Участвует в переработке линейных алкилбензолсульфонатов (в сульфофинилкарбоксилаты)

115 генов

76 – гипотетические (неизвестные)

Известные штаммы *P. stutzeri*

KOS6

шлам

CCUG_29243

Нафталин дегр, загр. м. ос.

ATCC_14405

морск

DSM_4166

ризосф

ATCC_17588

клинич

A1501

ризосф

DSM_10701

Почва бтрансформ прир в-в

SDM_LAC

Утилиз лактат

T13

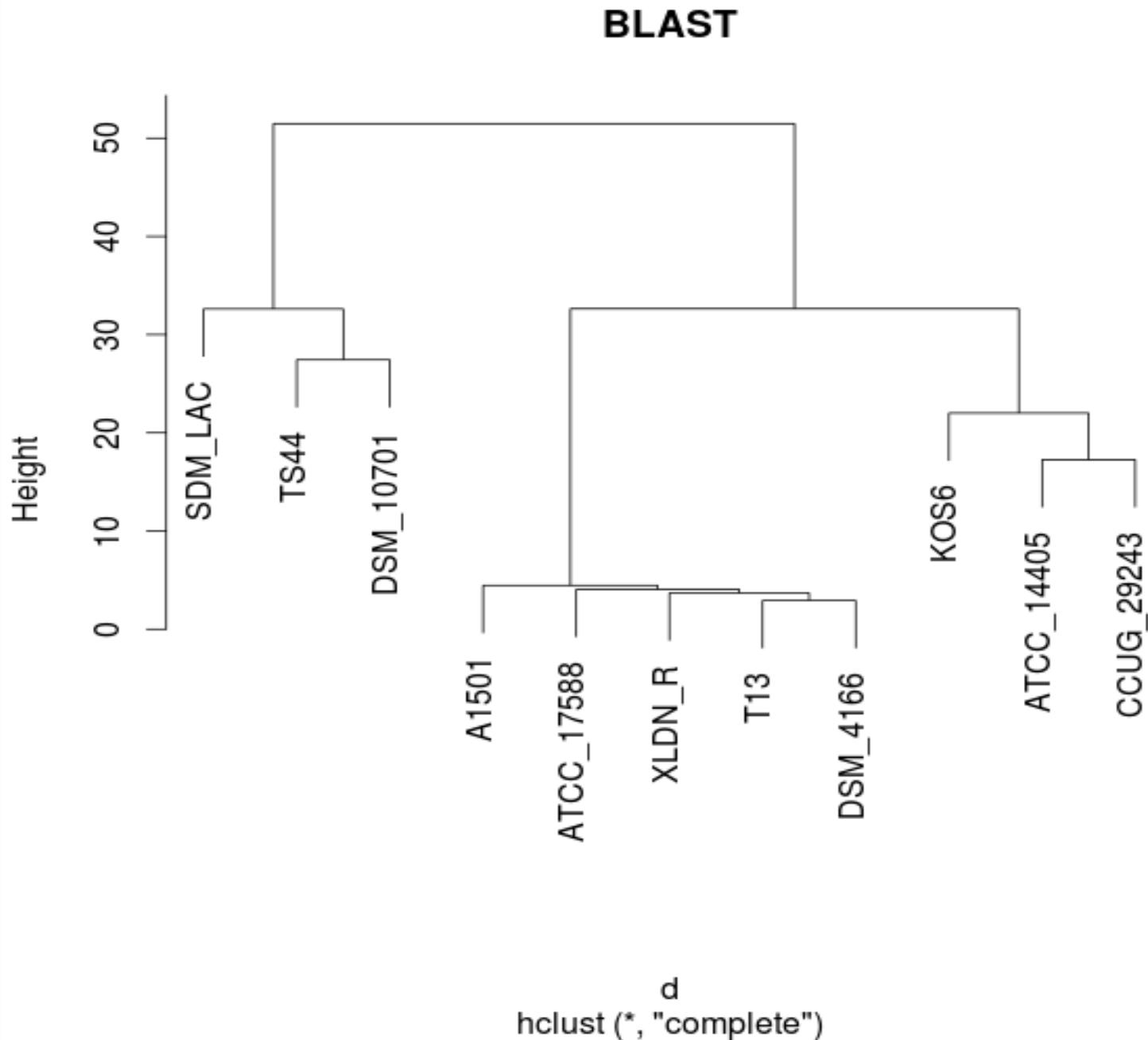
Активный ил

TS44 -

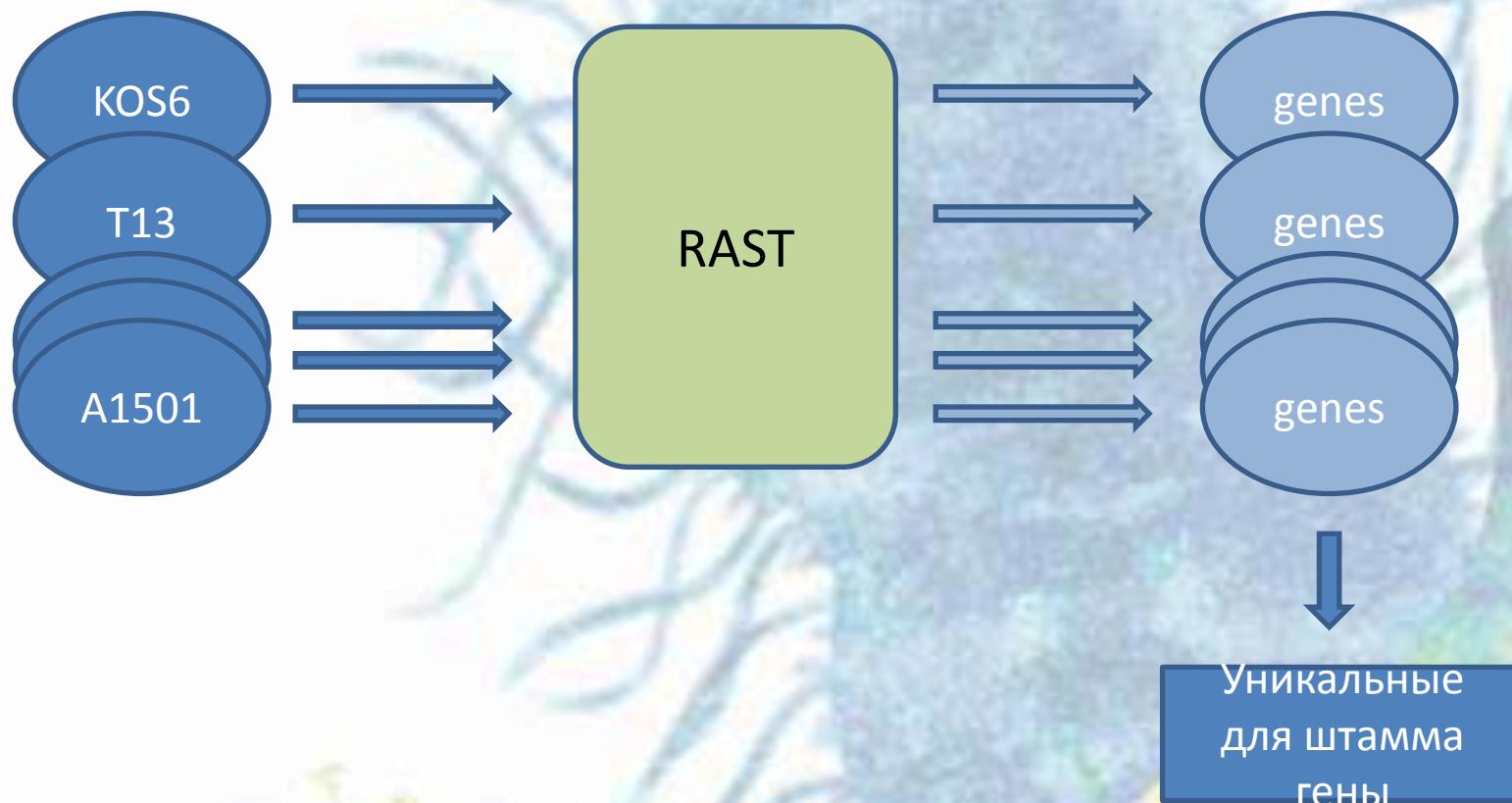
Загрязн почва трансф арсенита

XLDN_R -

Дегр крезола

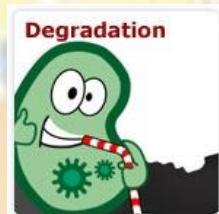


Как найти «уникальные» гены



- [1] "Enoyl-CoA hydratase [branched-chain amino acid degradation] (EC 4.2.1.17)"
- [2] "Aldehyde dehydrogenase A (EC 1.2.1.22)"
- [3] "Sialic acid transporter (permease) NanT"
- [4] "Acyl hydratase"
- [5] "FIG016502: iron uptake protein" - не уникальный...
- [6] "RecD-like DNA helicase Atu2026"
- [7] "RecD-like DNA helicase YrrC"
- [8] "VgrG-3 protein" - не уникальный...
- [9] "TRAP-type C4-dicarboxylate transport system, small permease component"
- [10] "TRAP dicarboxylate transporter, DctM subunit, unknown substrate 5"
- [11] "TRAP dicarboxylate transporter, DctQ subunit, unknown substrate 5"
- [12] "TRAP transporter solute receptor, unknown substrate 5"
- [13] "2,4-dihydroxyhept-2-ene-1,7-dioic acid aldolase (EC 4.1.2.-)"
- [14] "2-hydroxyhepta-2,4-diene-1,7-dioate isomerase (EC 5.3.3.-)"
- [15] "2-oxo-hepta-3-ene-1,7-dioic acid hydratase (EC 4.2.-.-)"
- [16] "3,4-dihydroxyphenylacetate 2,3-dioxygenase (EC 1.13.11.15)"
- [17] "5-carboxymethyl-2-hydroxymuconate delta-isomerase (EC 5.3.3.10)"
- [18] "5-carboxymethyl-2-hydroxymuconate semialdehyde dehydrogenase (EC 1.2.1.60)"
- [19] "5-carboxymethyl-2-oxo-hex-3-ene-1,7-dioate decarboxylase (EC 4.1.1.68)"
- [20] "Homoprotocatechuate degradative operon repressor"
- [21] "Transcriptional activator of 4-hydroxyphenylacetate 3-monooxygenase operon, XylS/AraC family"
- [22] "gliding motility protein MgIA"
- [23] "Cyn operon transcriptional activator"
- [24] "Phage capsid scaffolding protein"
- [25] "Phage tail fiber protein"
- [26] "Phage major tail tube protein"
- [27] "Phage tail completion protein"
- [28] "Phage tail length tape-measure protein"
- [29] "Phage tail protein"
- [30] "Phage tail sheath monomer"
- [31] "ParD protein (antitoxin to ParE)"
- [32] "Vanillate O-demethylase oxidoreductase (EC 1.14.13.-)"
- [33] "Carbon monoxide dehydrogenase large chain (EC 1.2.99.2)"
- [34] "Carbon monoxide dehydrogenase medium chain (EC 1.2.99.2)"
- [35] "Carbon monoxide dehydrogenase small chain (EC 1.2.99.2)"
- [37] "transcriptional regulator, Crp/Fnr family"

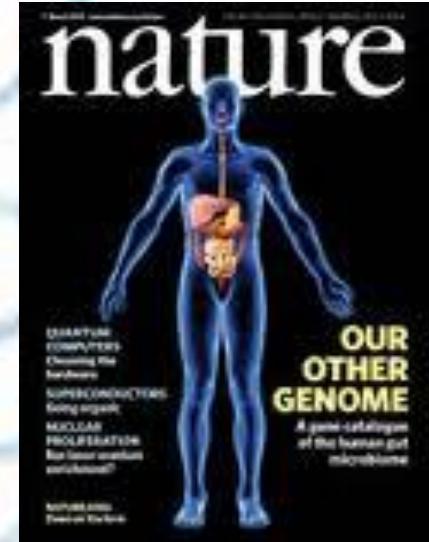
		Денитрификация, Азотфиксация nor, nir	nif	Деградация нафталена, nah	VI type
KOS6	шлам	+	+	+	+
CCUG_29243	Нафталин дегр, загр. морск. осадки.	+	-	+	+
ATCC_14405	Морская вода	+	-	-	-
DSM_4166	Ризосфера	+	+	-	-
ATCC_17588	Клинический	+	-	-	-
A1501	Ризосфера	+	+	-	-
DSM_10701	Почва	+	-	-	-
SDM_LAC	Утилизатор лактата	+	-	-	-
T13	Активный ил	+	-	-	-
Tolumonas_auensis	озер. отложения	+	+	-	+
TS44	Загрязненная почва трансф. арсенита	+	-	-	+
XLDN_R	Деград-я крезола	+	-	-	-



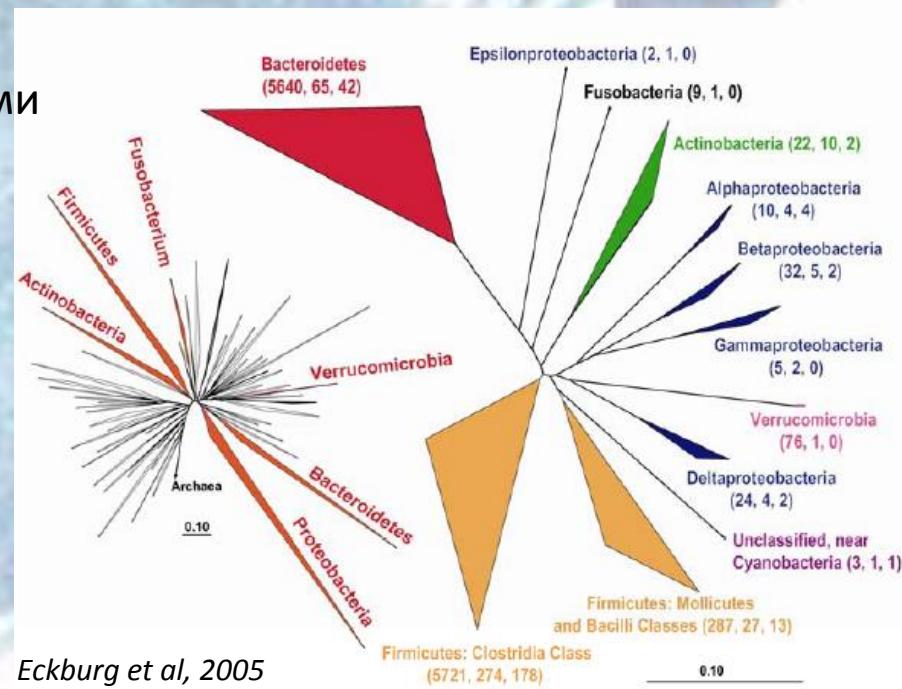
Выводы

- Произведена сборка de-novo и скаффолдинг *P. stutzeri KOS6*
 - Лучшие показатели: mate-pair, MIRA
- Присутствует 100 кб вставка, гомологичная участкам геномов бактерий, вероятно, обитающих в сходных условиях
- Комбинация способностей к денитрификации, азотфиксации, деградации нафталена является уникальной среди других известных штаммов *P. stutzeri*

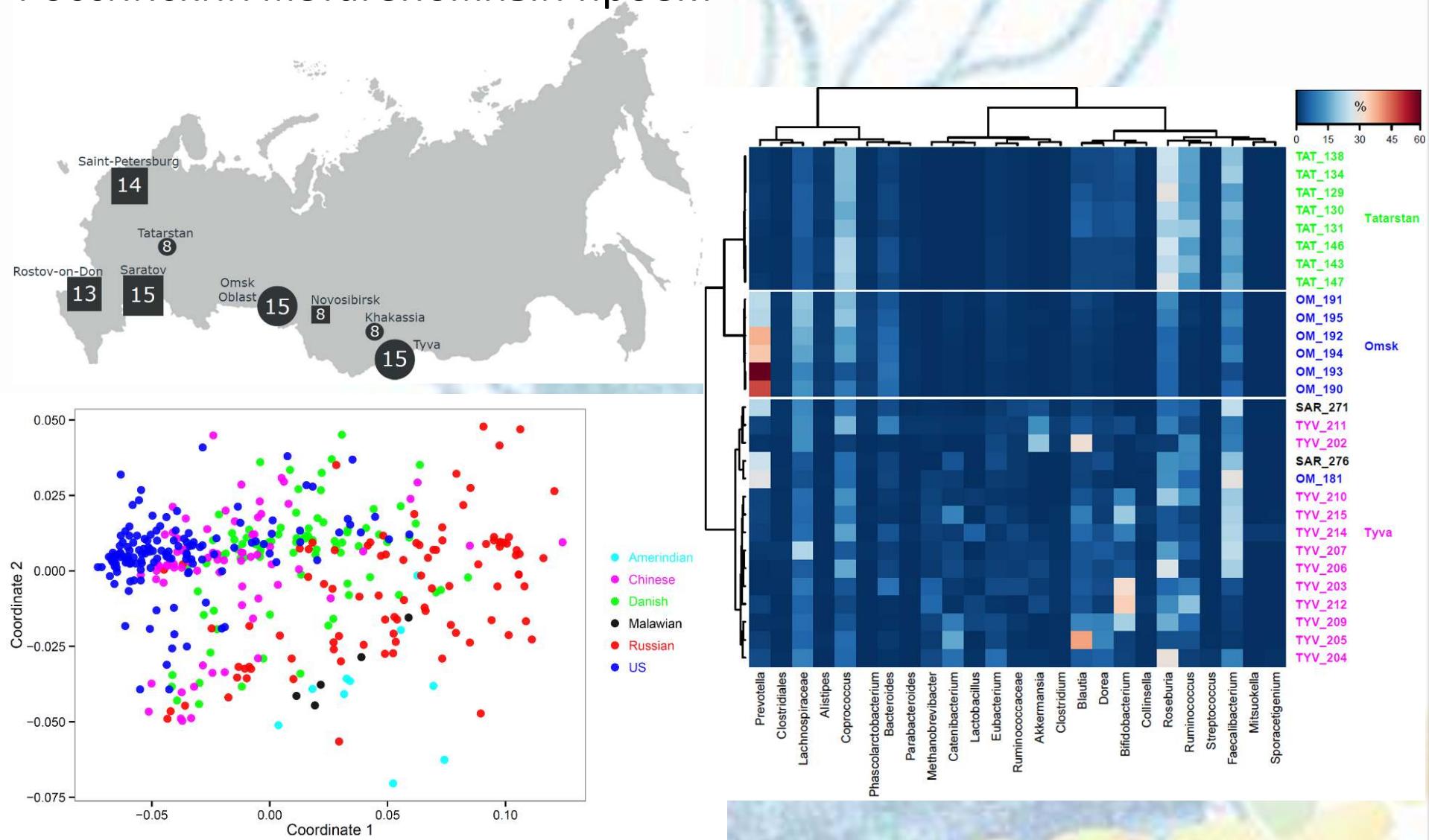
Микробиота кишечника человека



- Метагеном
- 300-1000 видов бактерий
- Число клеток $\sim 10^{15}$: на 2 порядка раз больше, чем клеток человека
- Число генов в метагеноме $\sim 10^6$: на 2 порядка больше, чем генов человека
- Функции:
 - Метаболизм
 - Предотвращение колонизации патогенами
 - Иммунитет человека
 - Защита от воспалительных заболеваний
 - Связь с высшей нервной деятельностью

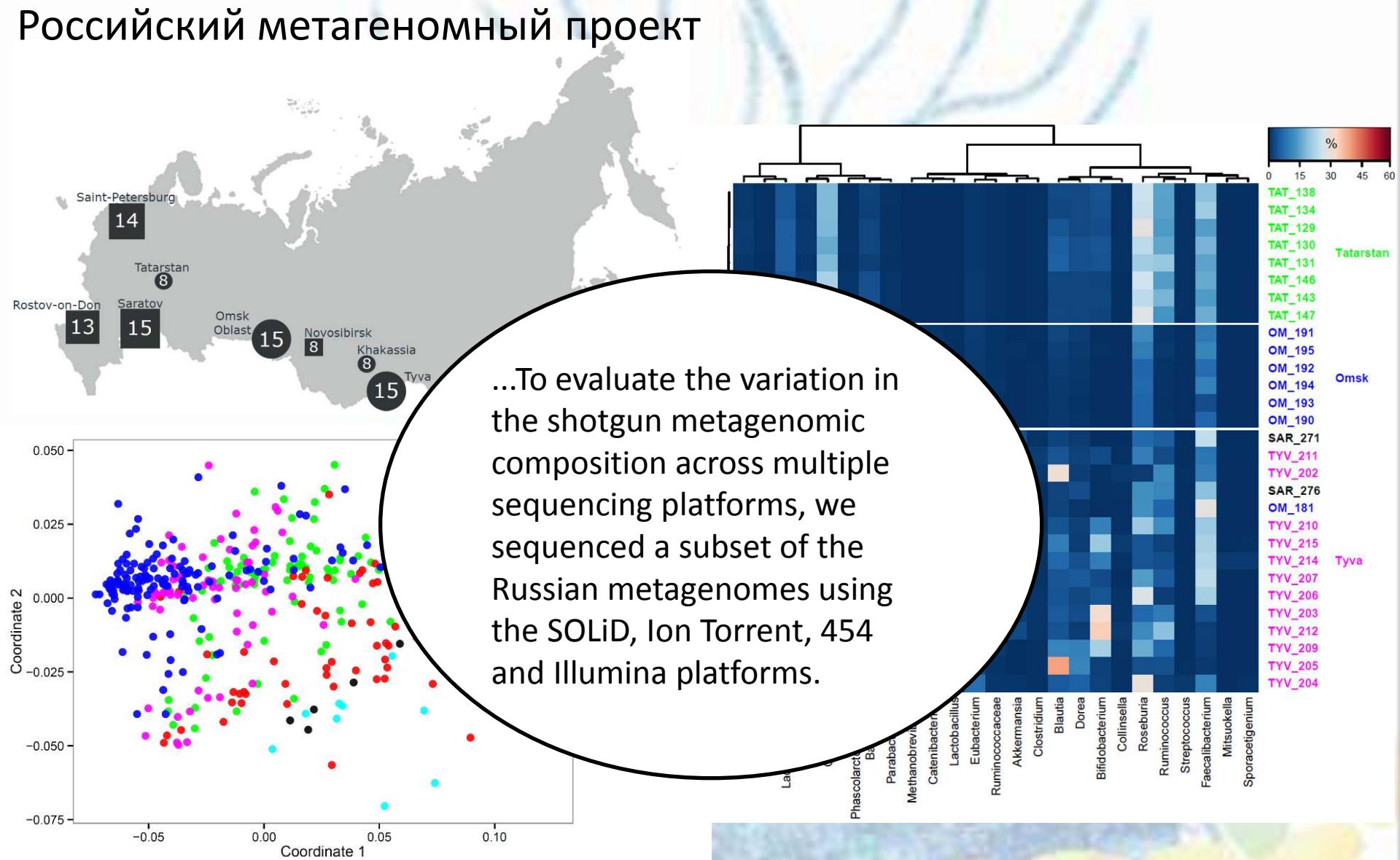


Российский метагеномный проект



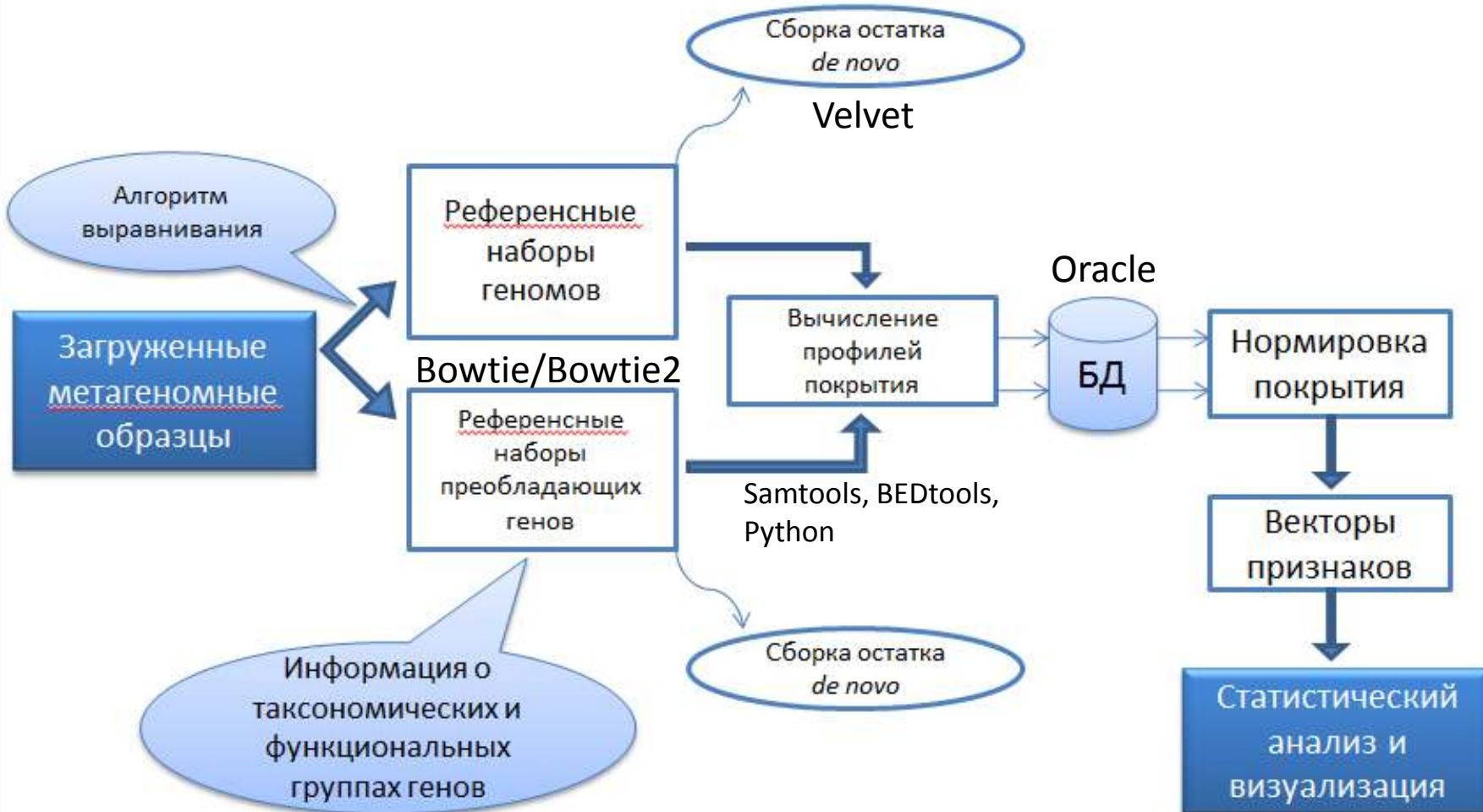
Tyakht, A. et al. Novel human gut microbiota community structures in urban and rural populations in Russia. *Nature Communications*, 2013 (accepted).

Российский метагеномный проект



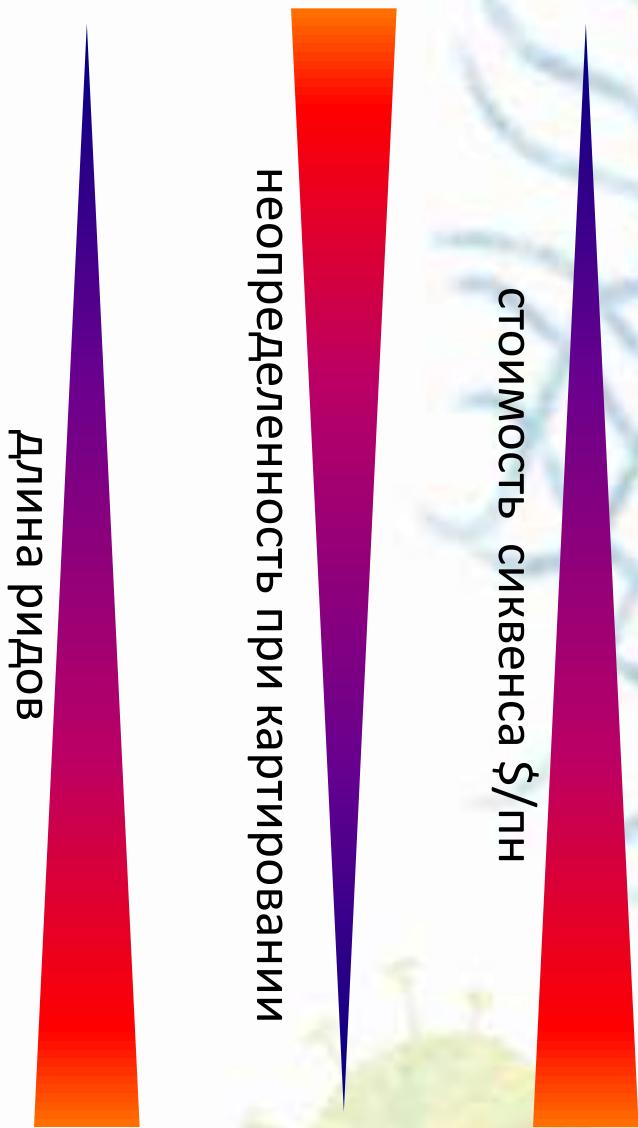
Tyakht, A. et al. Novel human gut microbiota community structures in urban and rural populations in Russia. *Nature Communications*, 2013 (accepted).

Высокопроизводительный конвейер для анализа метагеномных данных

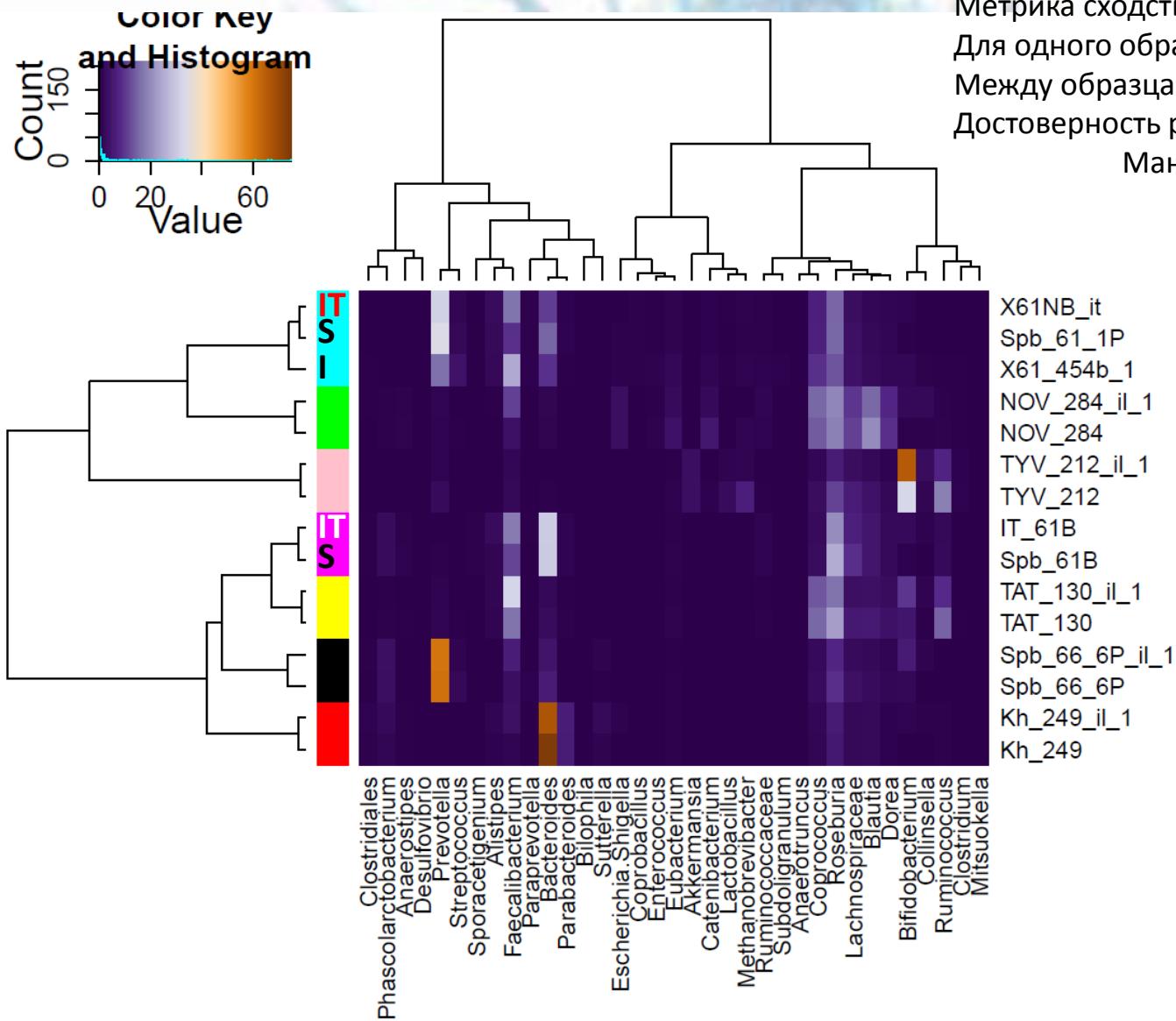


R

Выбор платформы секвенирования для метагеномного исследования



Сравнение состава микробного сообщества, полученного на разных платформах



Признаки: родовой состав

Метрика сходства: корреляция Спирмена

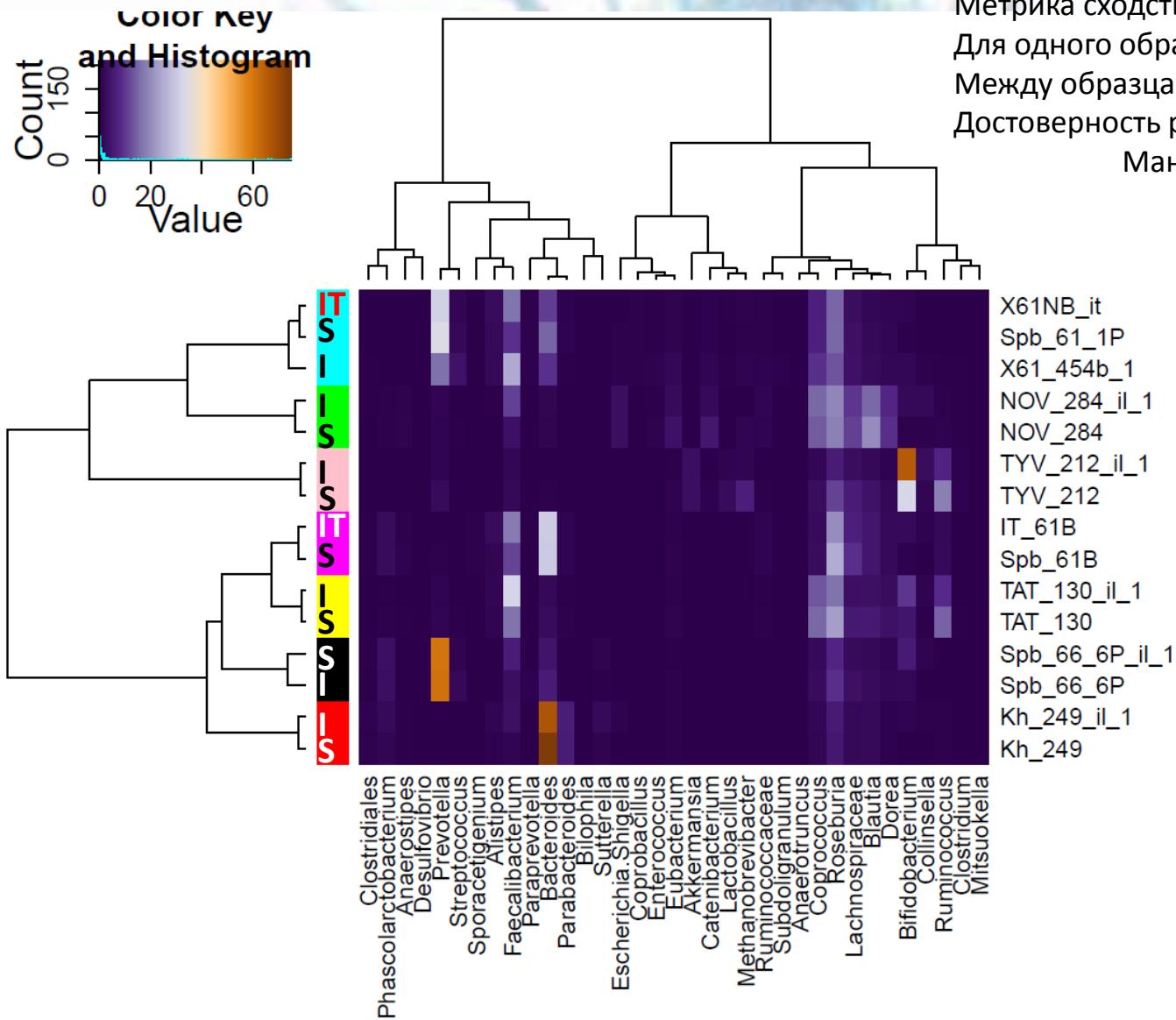
Для одного образца: $0,95 \pm 0,03$

Между образцами: $0,74 \pm 0,10$

Достоверность различия: тесты

Манна-Уитни, $P < 0,01$

Сравнение состава микробного сообщества, полученного на разных платформах



Признаки: родовой состав

Метрика сходства: корреляция Спирмена

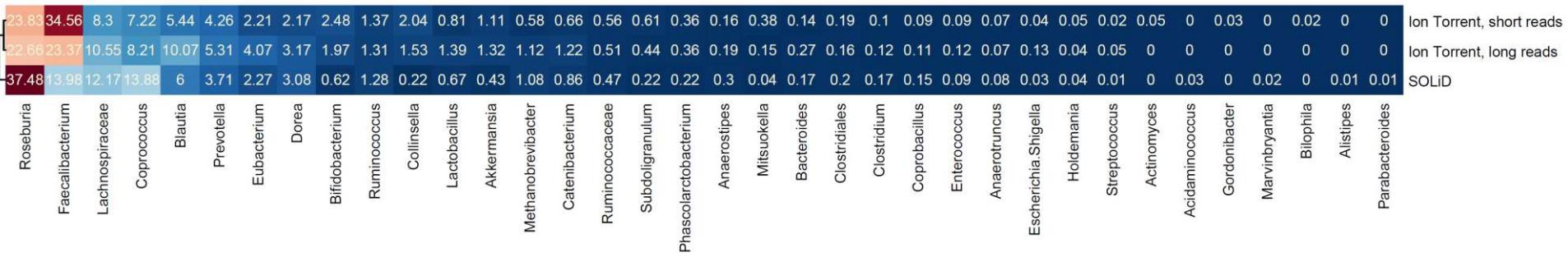
Для одного образца: $0,95 \pm 0,03$

Между образцами: $0,74 \pm 0,10$

Достоверность различия: тесты

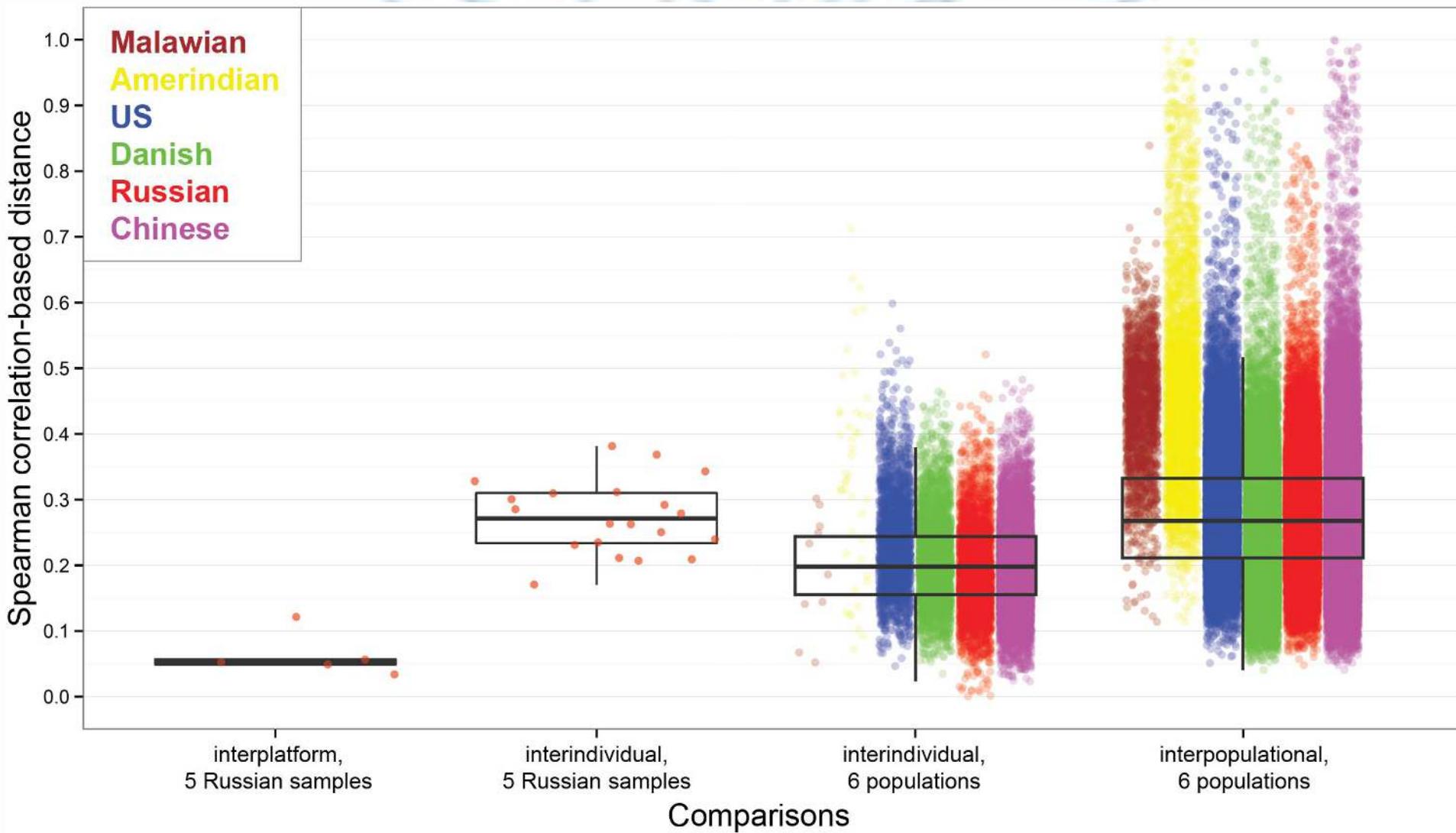
Манна-Уитни, $P < 0,01$

Разные длины ридов на Ion Torrent

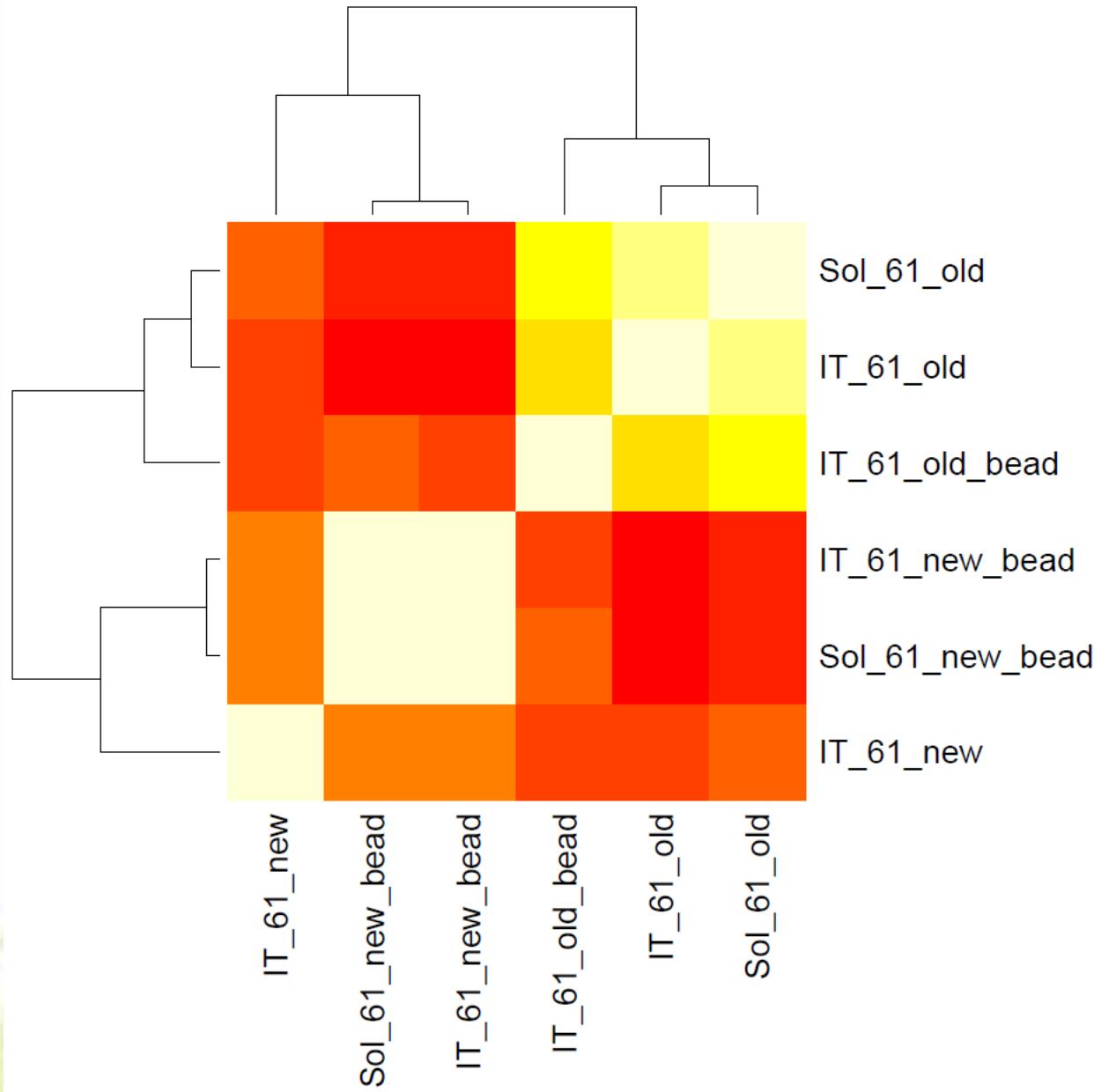


- Ion Torrent
 - 120 пн
 - 250 пн
- Solid – 50 пн
- Корреляция: 0.95 ± 0.02 s.d.

Фоновая корреляция

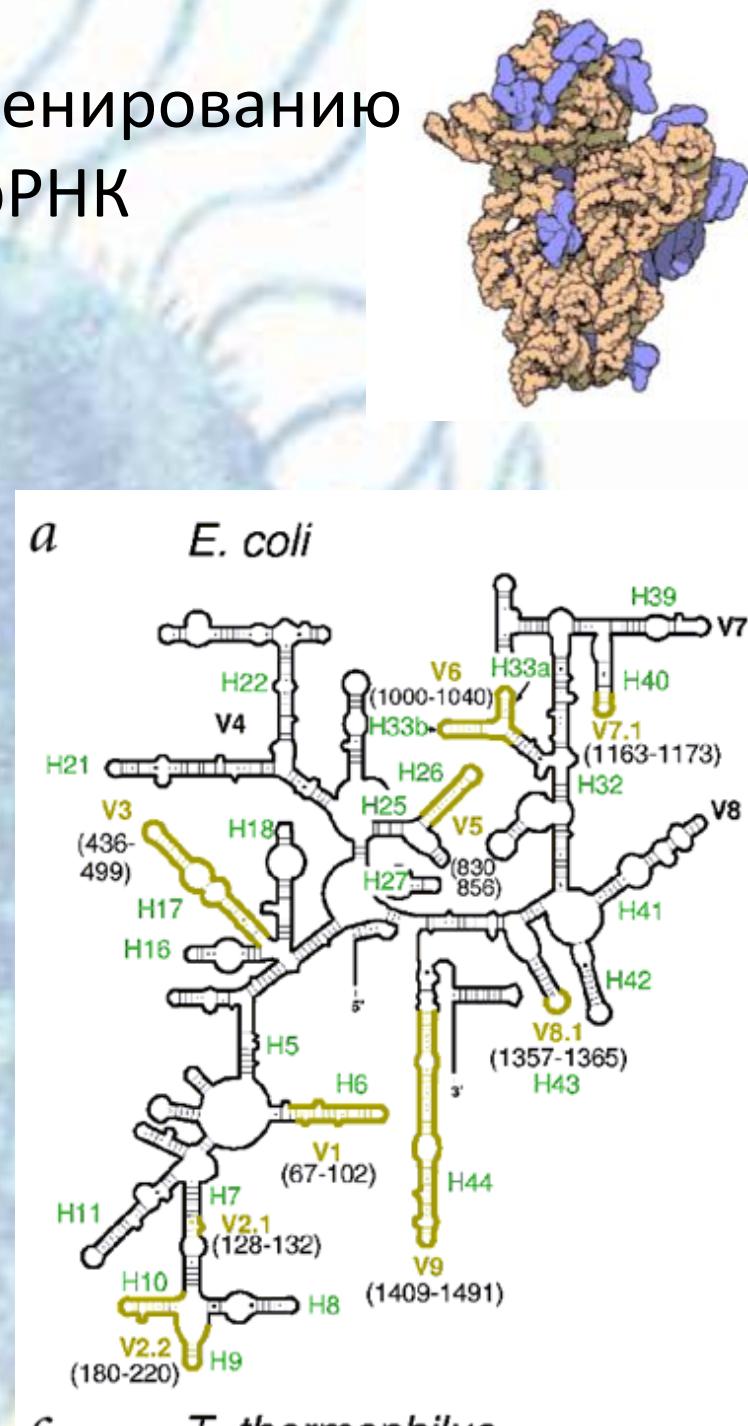


Факторы вариабельности микробиотного профилирования: временная вариация > метод пробоподготовки > платформа



Анализ микробного состава по секвенированию последовательностей 16S рРНК

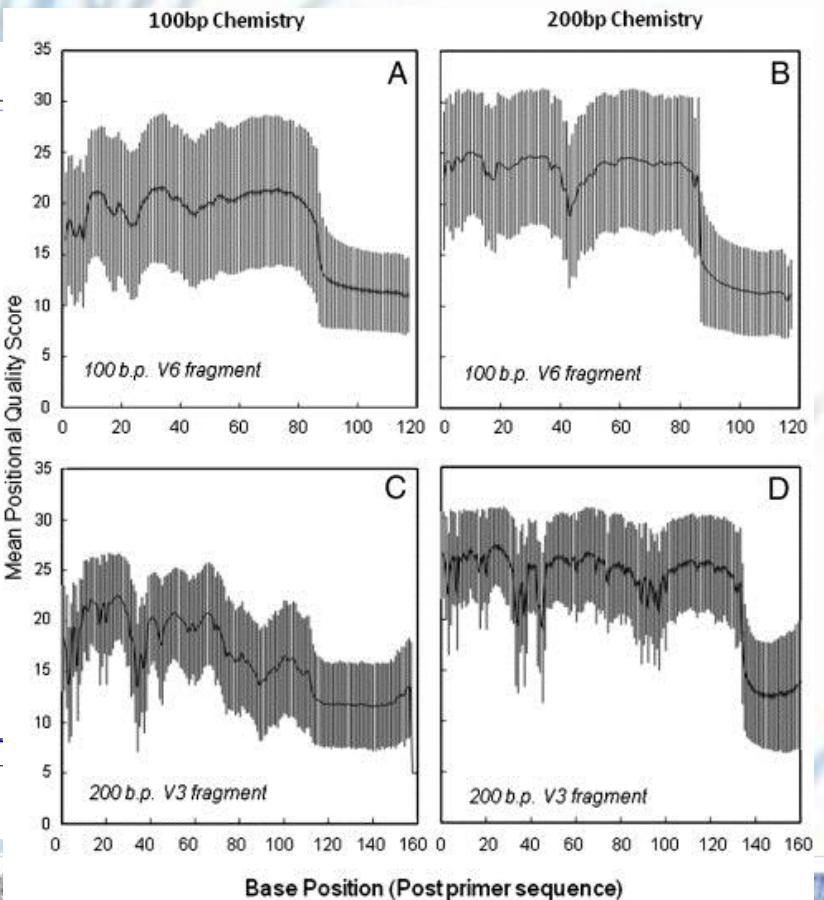
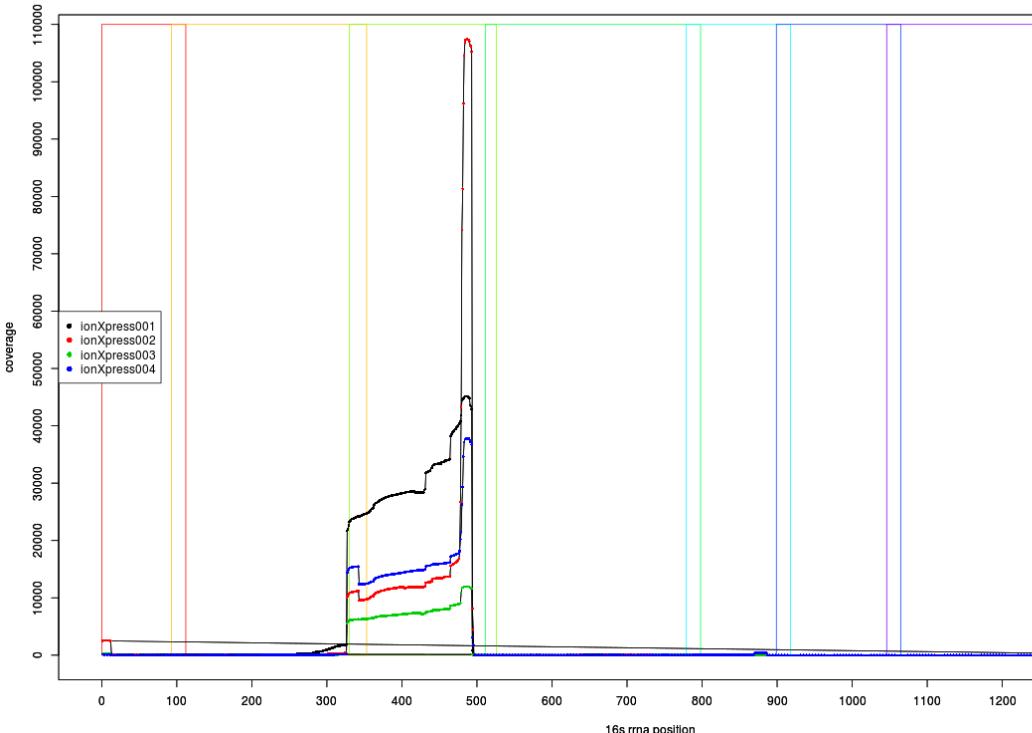
- Внутривидовое сходство сиквенса 98-99%
- Выделяется из тотального ДНК с помощью универсальных праймеров
- Секвенируется целиком либо вариабельные области
 - Для микробиоты кишечника: стандартный протокол Human Microbiome Project (<http://www.hmpdacc.org>) для V13, V35
- Для эукариот – можно использовать 18S рРНК.



Преимущества 16S формата

- Эффективность, поточный анализ
- Секвенирование из малого количества бактериальной ДНК (биоптат тканей человека, мокрота)
- Детекция бактериального заражения даже при низкой глубине покрытия

Ion torrent 16S



Volume 91, Issue 1, October 2012, Pages 80–88



Microbial 16S rRNA Ion Tag and community metagenome sequencing using the Ion Torrent (PGM) Platform

Andrew S. Whiteley^{a, b}, Sasha Jenkins^{b, c}, Ian Waite^{b, c}, Nina Kresanje^{d, e}, Hugh Payne^f, Bruce Mullan^f, Richard Allcock^{d, e}, Anthony O'Donnell^{b, c}

Пример: Влияние среды на микробиоту. Работники вредного производства (n=13), QIIME+R

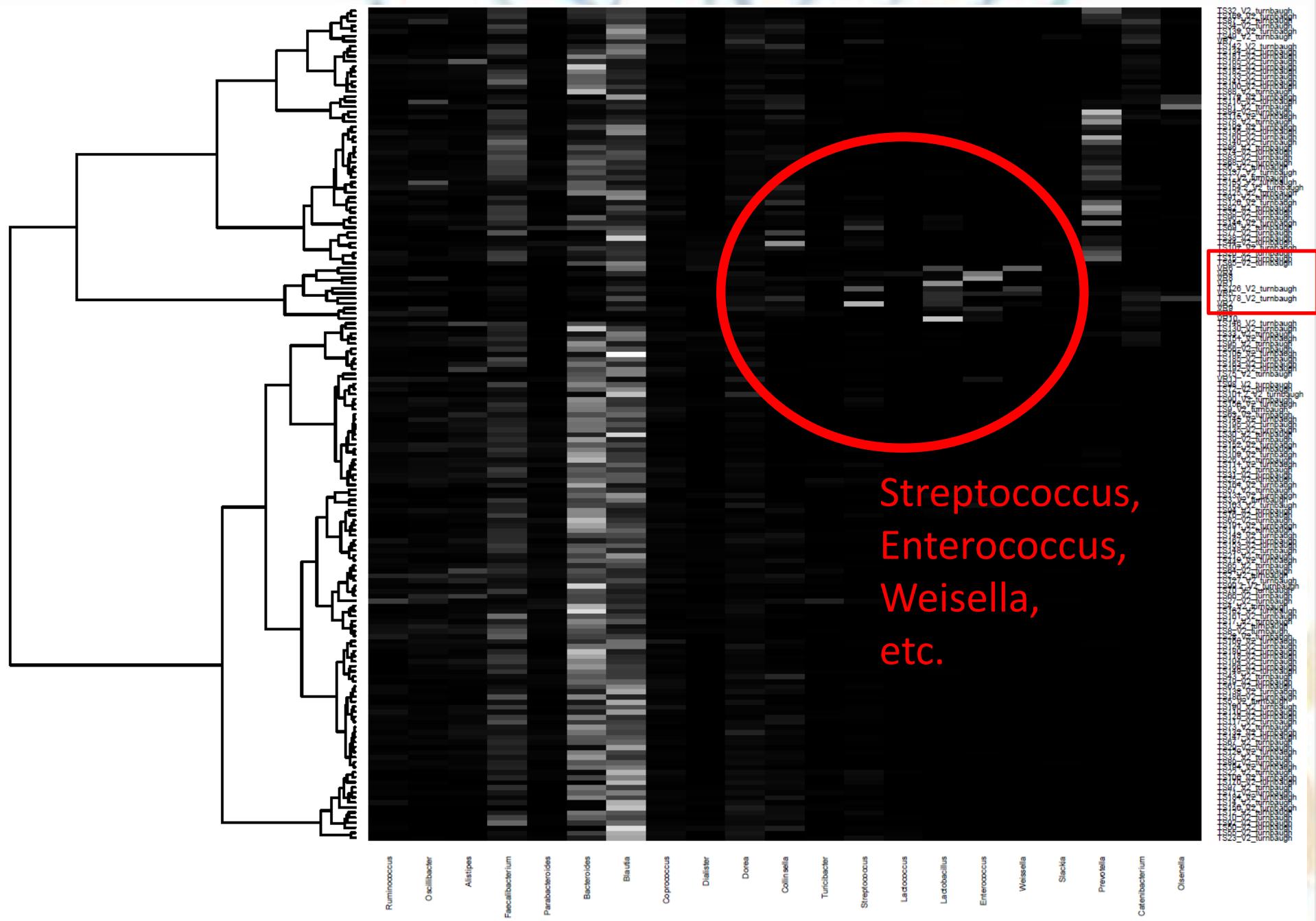


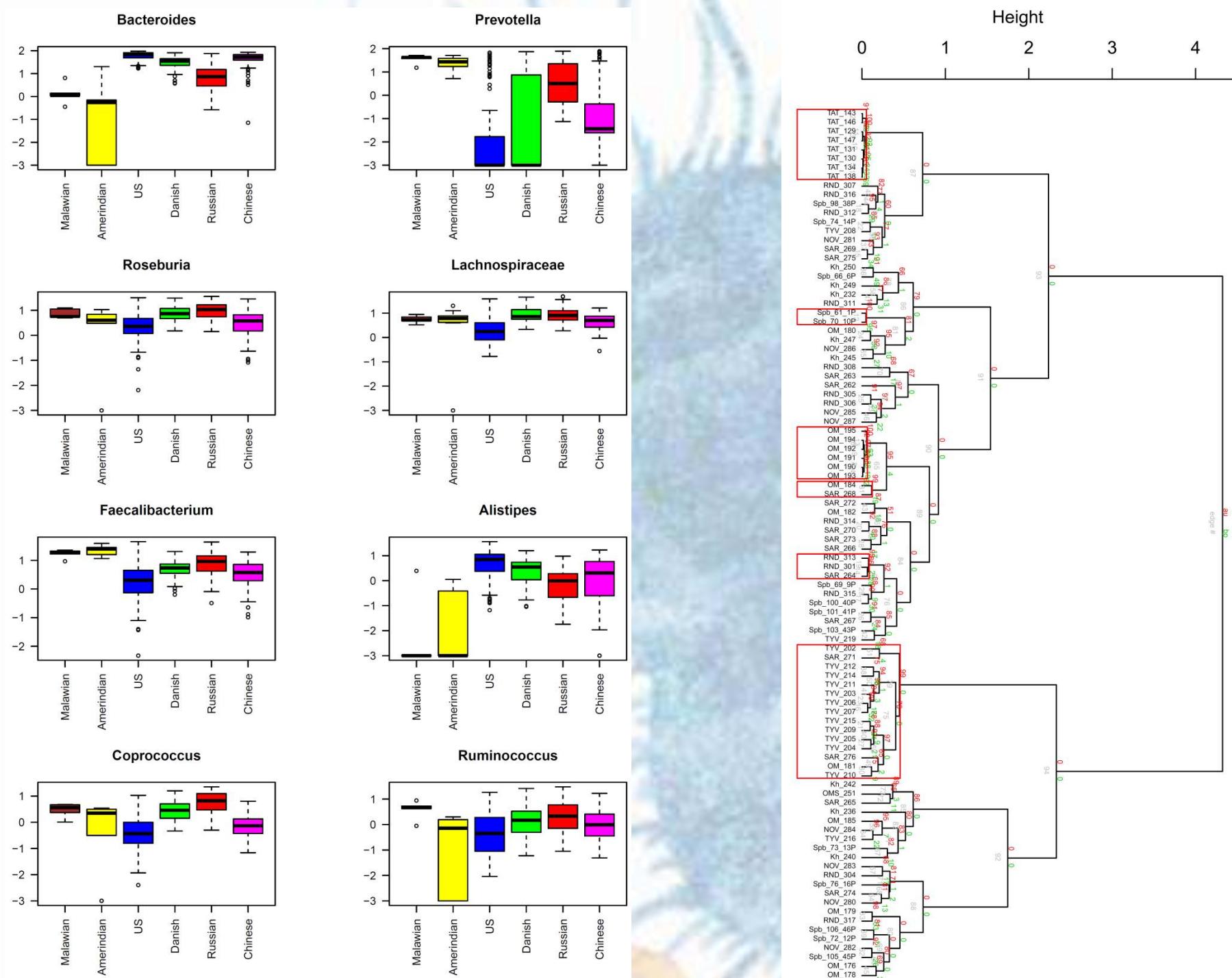
Consensus Lineage

k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Erysipelotrichi;o_Erysipelotrichales;f_[Coprobacillaceae];g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Dorea;s_formicigenerans
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Veillonellaceae;g_Dialister;s_invisus
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Streptococcaceae;g_Lactococcus;s_garvieae
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_[Ruminococcus];s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Coprococcus;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g_Collinsella
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g_Oscillospira;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactobacillus;s_ruminis
k_Bacteria;p_Firmicutes;c_Erysipelotrichi;o_Erysipelotrichales;f_Erysipelotrichaceae;g_[Eubacterium];s_biforme
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Clostridiaceae;g_Clostridium
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Blautia;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactobacillus;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Peptostreptococcaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae
k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Dorea;s__
k_Bacteria;p_Firmicutes;c_Erysipelotrichi;o_Erysipelotrichales;f_[Coprobacillaceae];g_Catenibacterium;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Blautia;s_obeum
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Enterococcaceae;g_Enterococcus;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae
k_Bacteria;p_Tenericutes;c_Mollicutes;o_RF39;f__;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Coprococcus;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Lactobacillaceae;g_Lactobacillus;s_brevis
k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g_Adlercreutzia;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Clostridiaceae;g_Clostridium
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Peptostreptococcaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Streptococcaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Blautia;s__
k_Bacteria;p_Bacteroidetes;c_Bacteroidia;o_Bacteroidales;f_Bacteroidaceae;g_Bacteroides;s__
k_Bacteria;p_Proteobacteria;c_Gammaproteobacteria;o_Enterobacteriales;f_Enterobacteriaceae;g_Escherichia;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_[Ruminococcus];s_gnavus
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Leuconostocaceae;g_Weissella;s__
k_Bacteria;p_Firmicutes;c_Erysipelotrichi;o_Erysipelotrichales;f_Erysipelotrichaceae;g_[Eubacterium];s_cylindroides
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_Streptococcaceae;g_Streptococcus;s_luteiae
k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Coriobacteriales;f_Coriobacteriaceae;g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Ruminococcaceae
k_Bacteria;p_Firmicutes;c_Bacilli;o_Lactobacillales;f_g__;s__
k_Bacteria;p_Firmicutes;c_Clostridia;o_Clostridiales;f_Lachnospiraceae;g_Sphaerotilus

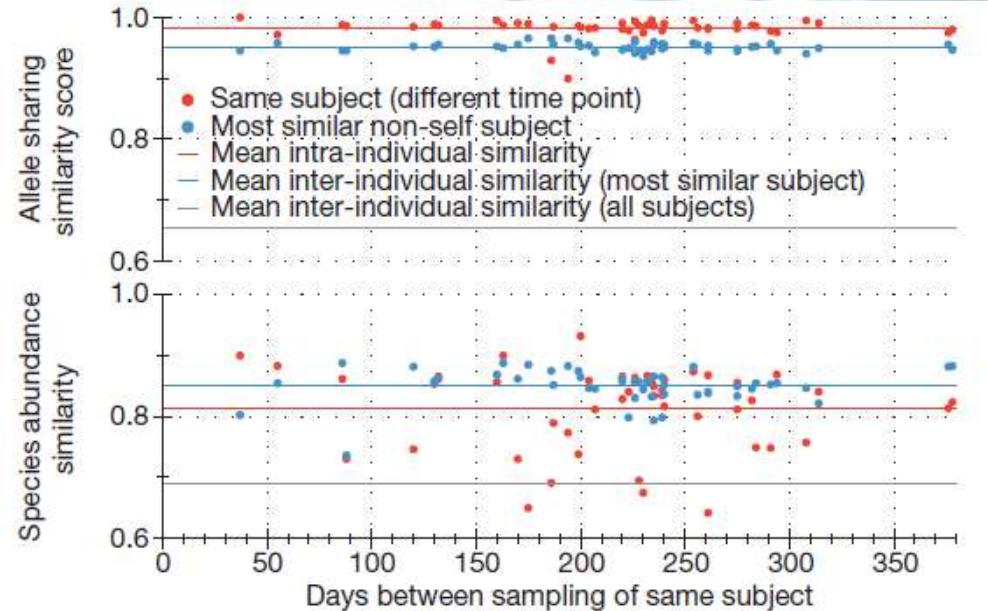
#OTU ID	VR9	VR4	VR3	VR8	VR10	VR11	VR2	VR5	VR7	VR1	VR6
65	173	694	559	15		99	102				
892	8	258	7	4	4	570	53	60	784		
1188	271	231	67	13	11	549	13	78	1151	15	2343
1524		1784	1	1		6	2				1689
1920	192		33		1	66	19	4	1081	75	
2349	725	38	147	359	33	152	242	55	166	294	
2467	16	334	22	2		35	34	25	1453	62	2
2663	1519	9121	737	12588	654	126	239	448	37	107	
3181	954		373	761	54		1171	75	544	208	
3250	876	30	953	2319	60	1256		115	3		
3363	1005		275		6		1880	290		4772	65
3452	366	1	639	81	9		890	32	1656	6706	
3632		1	1854								
3669	56	4280		78	2461		12	195		147	17
3858	394	725	57	58	4	42	237	106	22	262	238
3968		58		21	125		269	3		3844	
4202	4	41	1117	3643	507		167	39		1929	241
4266		402	19	257	9631		3434	267	4	1	6322
4270		975				59				327	
4910	279			52		89	9	2	2304	9	
5222	315		183		276		1268	245	2780		
5348	101	1879	204	334	43	86	233	71	1250	408	24
5464	558	3398		20	6	1977	2	2		1	
5474	1483		13	70	5	15		4		9	
6363		1	3374							2	
6384	183	7	103	1	4	109	5	6	1216	23	55
6450	217	258	141	1	1	2089	5	75	382	19	310
6460		252	2	78	27		14	252			3605
7153	47	751		5		46				521	
8232	3	551	6917	119	590		250	32	1	236	98
8542		1139	353	208		1	8	8		97	
8895		190		4	1	8	3511	1	195	121	38
8900	752	231	130	120	15	323	1592	89	2156	1367	12152
9055		81		3		3898	1	75			
9331	878	39	3	1	2	1049	27	20		45	5
9979	1095	1414	56	836	24	1972	1978	7	2450	1145	32
10493					83	3	843	424	4		10301
10892	1608					26		482			
11432				2	10	1	4558	40		4	
11492					122	348	106	202	41	1090	
11702		153	484	623	267	239	1427	108	28	594	
11707		1		79	1628	1	656	63	50		2
11751		270	5057	200	270	207	210	1050	204	700	

Пакет R: сравнение с внешней контрольной группой (Turnbaugh et al, 2010, 154 образца)





SNP variations in the metagenome



Temporal stability of genomic variation patterns.

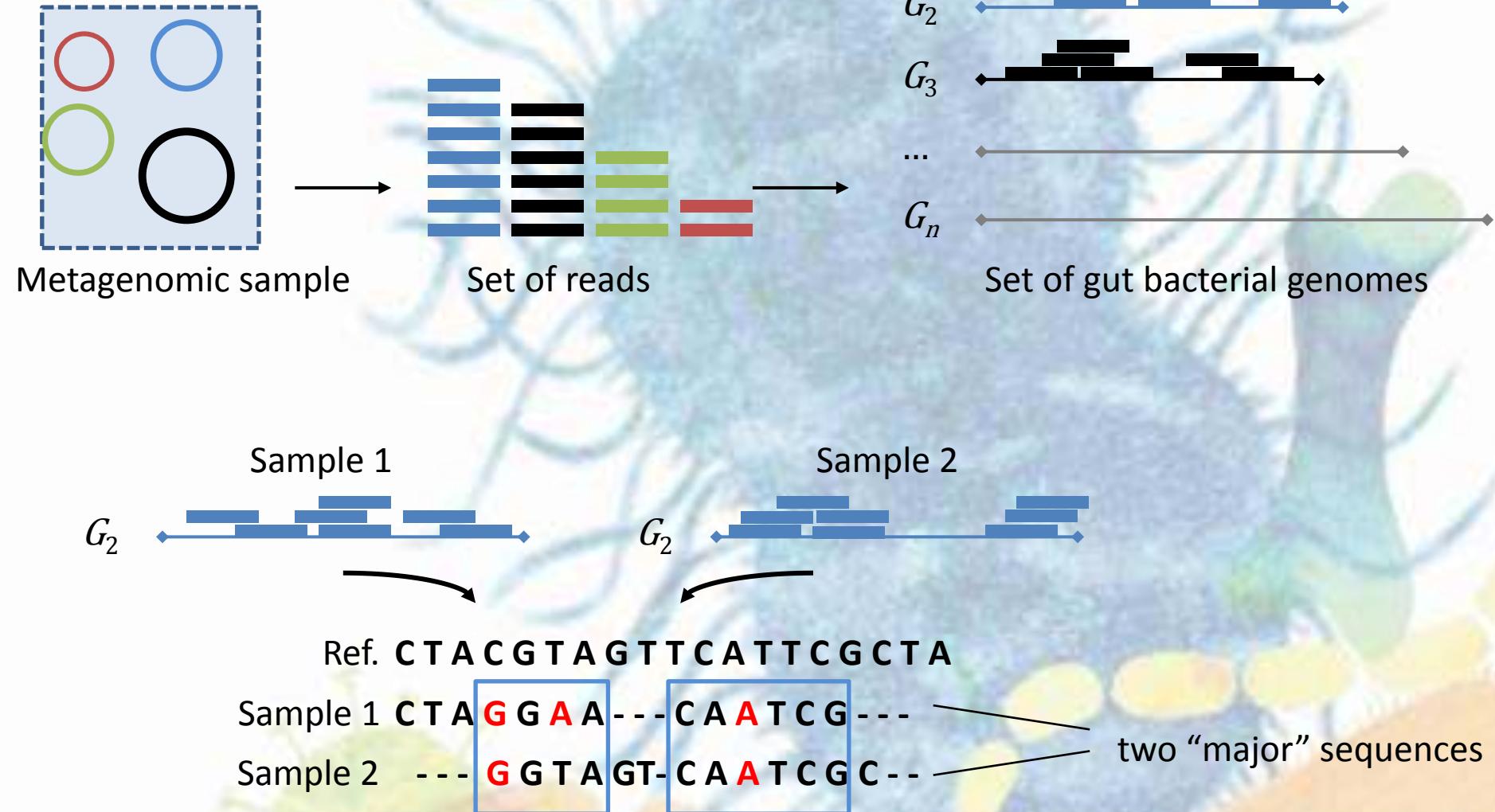
Data: 139 (USA), 110 (Europe)
metagenomic samples

Distance measure: fixation index and
nucleotide diversity

Results:

- Individual-specific variation patterns are stable over time
- Absence of clear continental stratification apart from small number of species.

Our method of study of metagenome genetic variety



Sources of nucleotide differences in metagenomes

- True SNPs
- Mapped reads from other bacterial genome
- Errors of sequencing

So the main problem was proper choice of coverage, allele sharing and allele support thresholds to generate “major” sequence for each sample. In our opinion, distance between such major sequences is more robust characteristic than fixation index.



Data

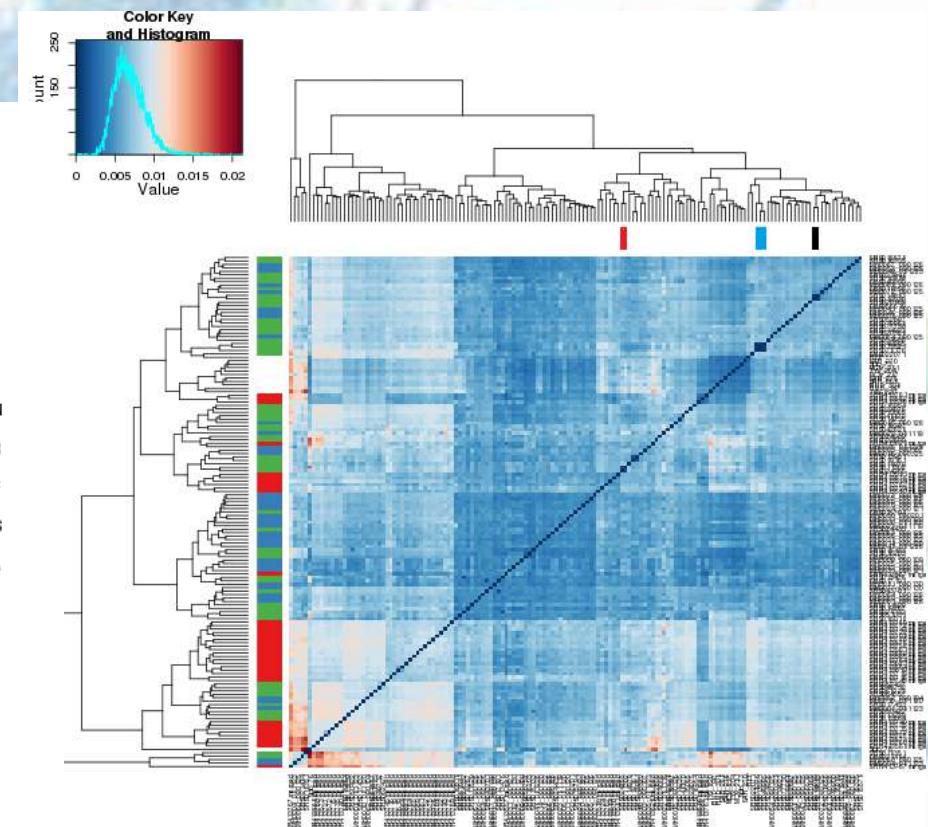
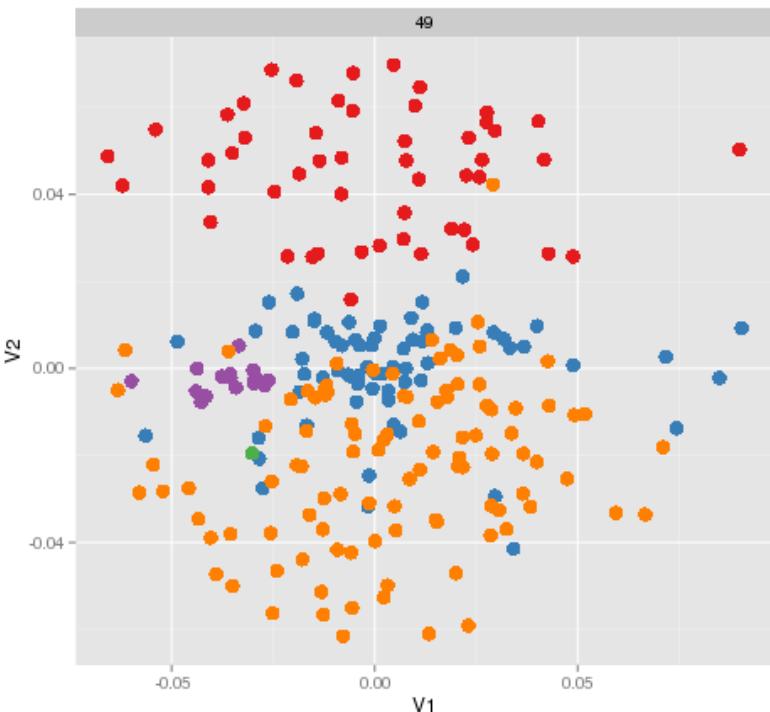
Host country	Source of the data	Number of samples	Number of individuals	Sequencing platform
USA	HMP	138	50 (sampled once) 41 (sampled twice) 2 (sampled three times)	Illumina
Denmark	MetaHit	84	84	Illumina
China	BGI-Shenzhen	100	50 (type 2 diabet) 44 (healthy)	Illumina
Russia	RIPCM	20	20	SOLiD

Total: 342 gut metagenome samples of individuals from 4 countries.

Metagenome reads have been mapped on the catalog of 444 bacterial genomes

We have chosen 93 bacterial genomes that have abundance level greater than 1% at least in 50 individuals.

Result



Example: *Faecalibacterium prausnitzii* A2-165

Consistency of results

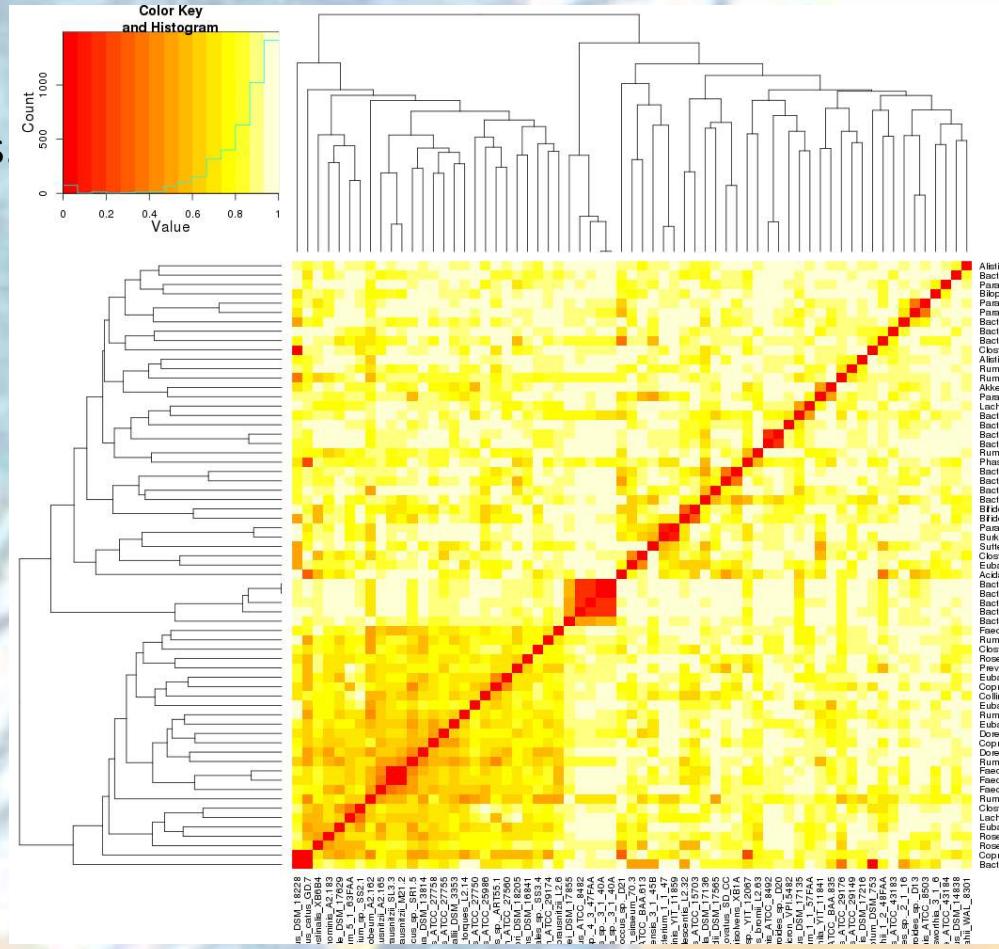
On the heatmap we can see dense clusters of different strains of the same bacterial species:

- *Faecalibacterium prausnitzii* (4 strains)
 - *Bifidobacterium adolescentis* (2 strains)
 - *Bacteroides vulgatus* (2 strains.)
- ...and some others

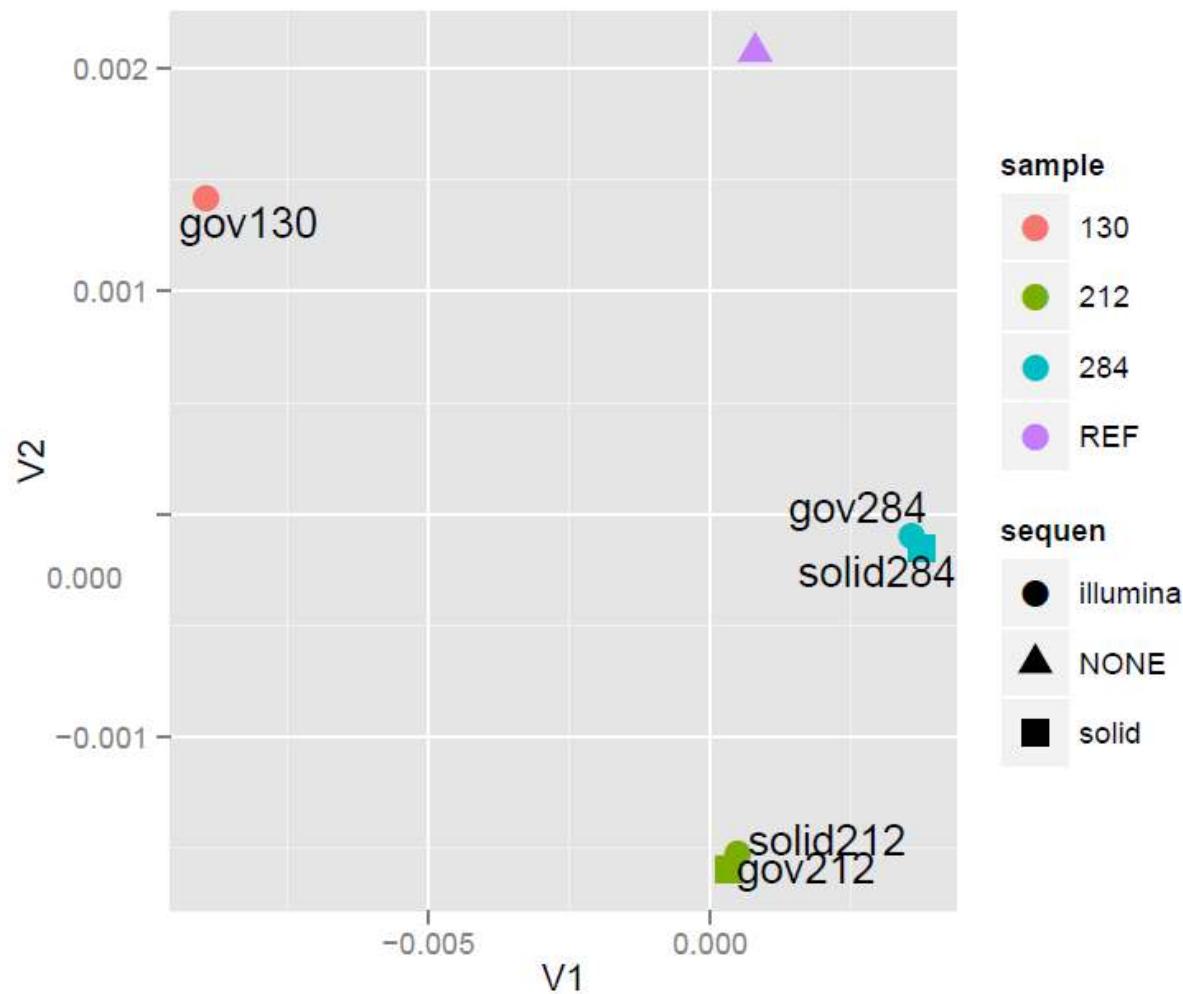
Similarity measure:

$$d(M_i, M_j) = 1 - |\text{cor}(M_i, M_j)|$$

$\text{cor}(M_i, M_j)$ - Mantel test Spearman correlation between distance matrices



Dependence on sequencing platform



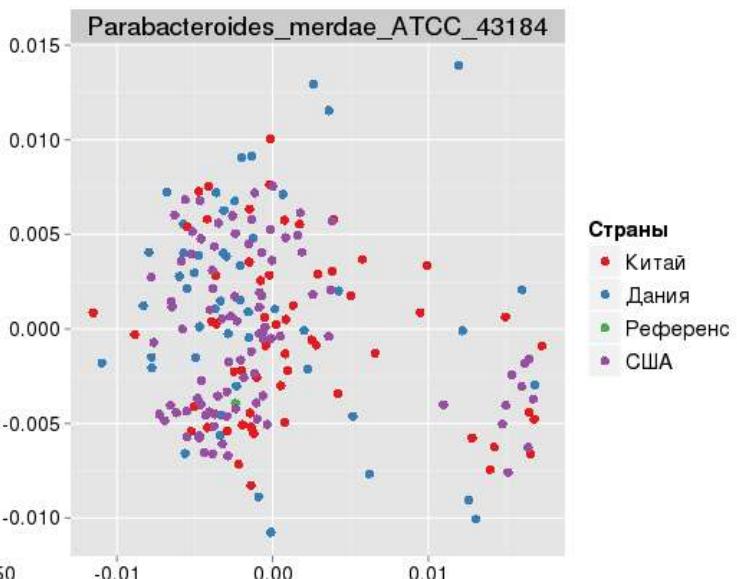
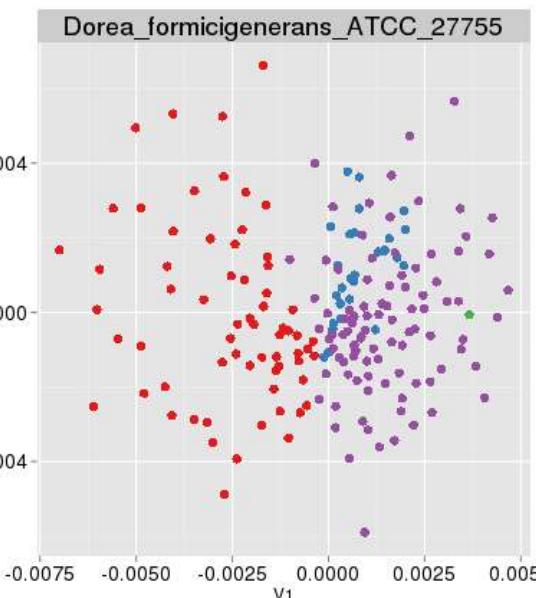
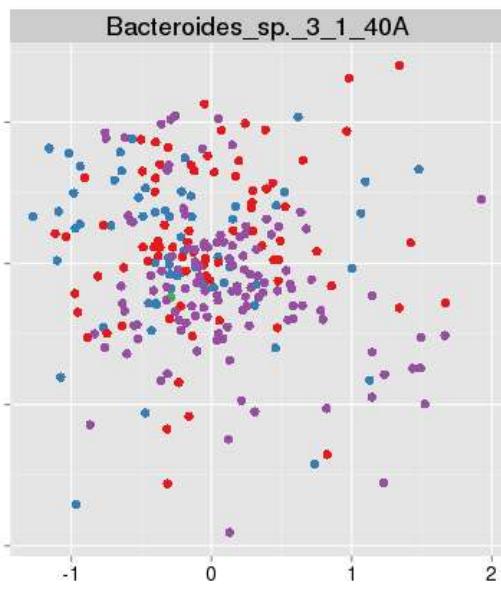
We've sequenced 5 metagenomic samples on different platforms.

Stratification patterns of different bacterial species

Absence of any stratification

Geographic factor

Non-geographic factor

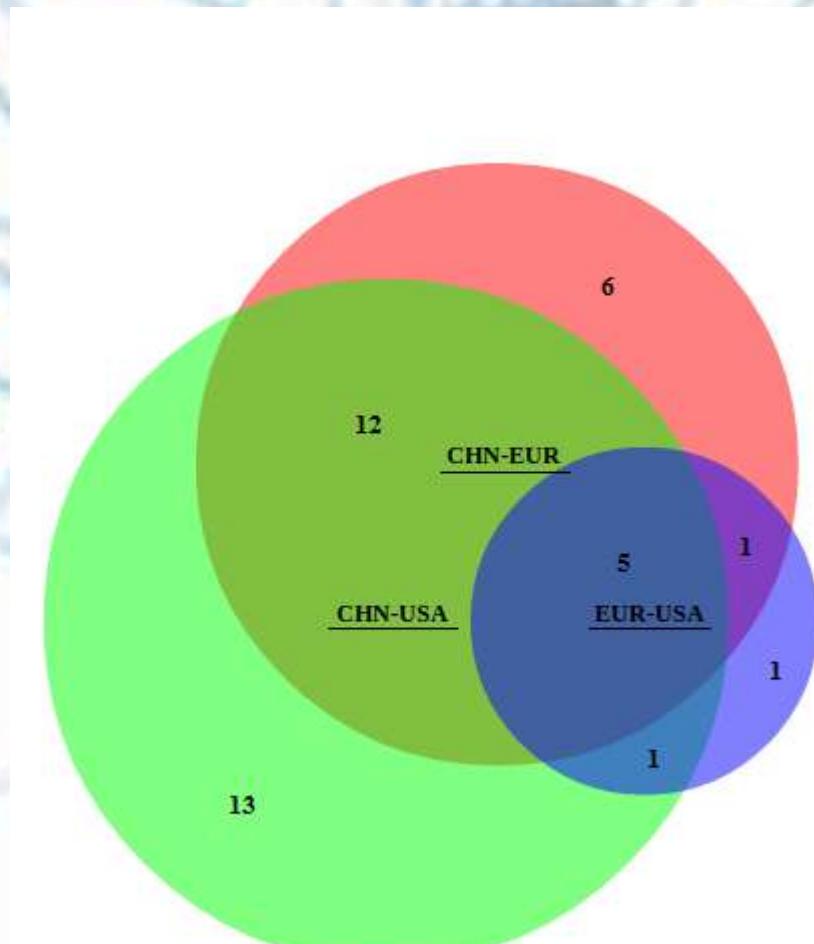


Страны

- Китай
- Дания
- Референс
- США

MDS-projections of distance matrices of different bacterial species.

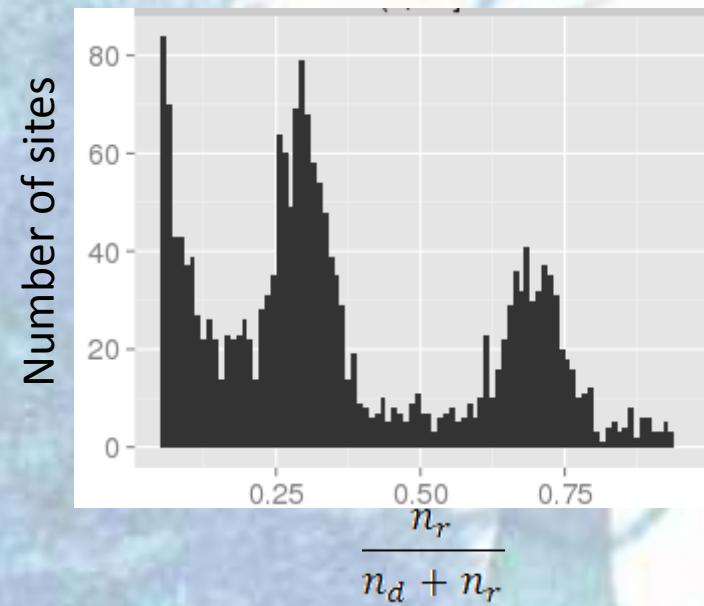
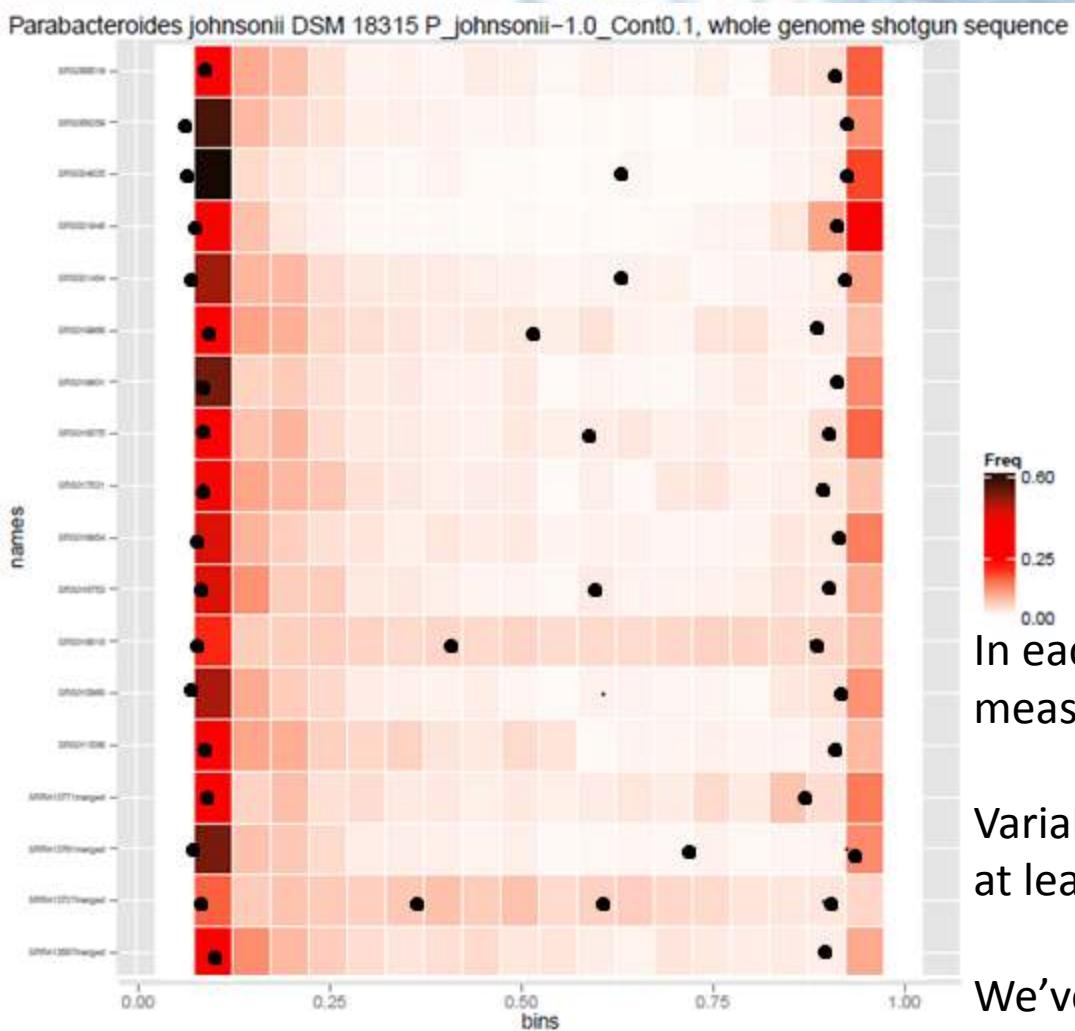
Geographic stratification in different “directions”



Circle means the number of bacterial species with significant difference in Wilcoxon test in particular “direction” (inter-country)
Total number of tested bacterial species: 65.

“Heterogeneity” of variable sites in samples

Is it a manifestation of co-existence subpopulation of same bacterial species in the same organism?



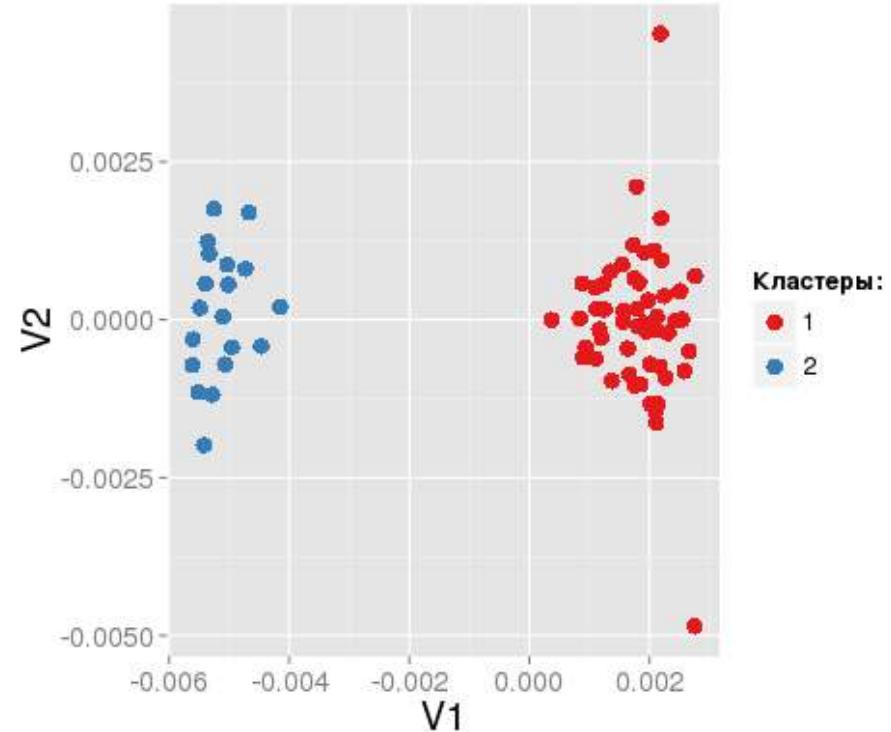
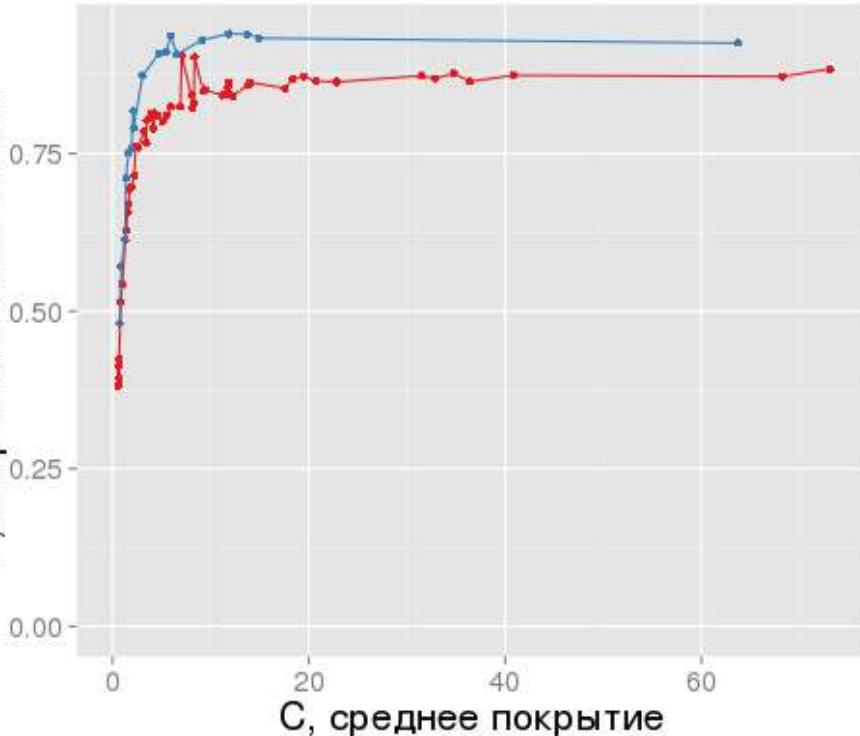
In each variable site of all samples we have measured “recessive” allele fraction.

Variable site means that it differs from reference at least N of M samples.

We've used different thresholds: {N, M}

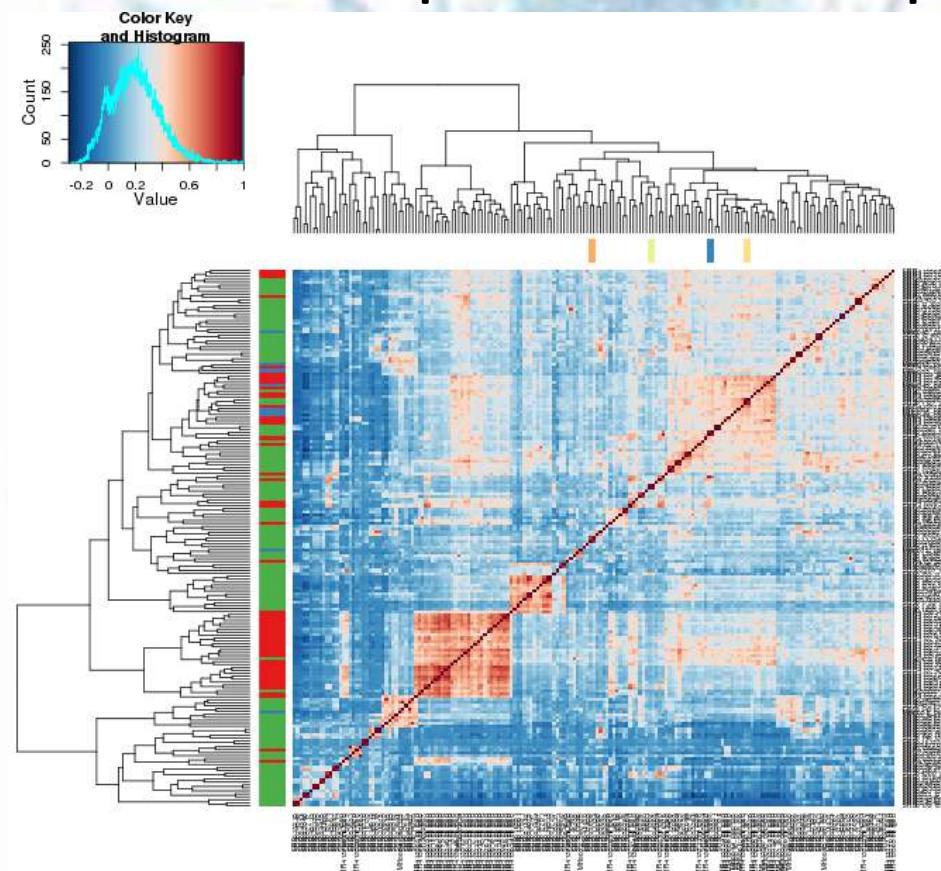
Polymorphism-level difference is connected with gene composition difference between samples

F, покрытая часть генома



$$F \approx \beta(1 - e^{-c}) = \beta \left(1 - \exp\left(\frac{NL}{G}\right)\right)$$

Correlation between coverage profiles of different samples of same species



We have chosen samples with coverage ≥ 5 . Then

Разработка метода картирования
сайтов инициации транскрипции с
помощью технологии секвенирования
второго поколения на примере
Mycoplasma gallisepticum

Фисунов Г.Ю., Горбачёв А.Ю.,
Семашко Т.А., Мазин П.В.

Зачем искать сайты инициации транскрипции

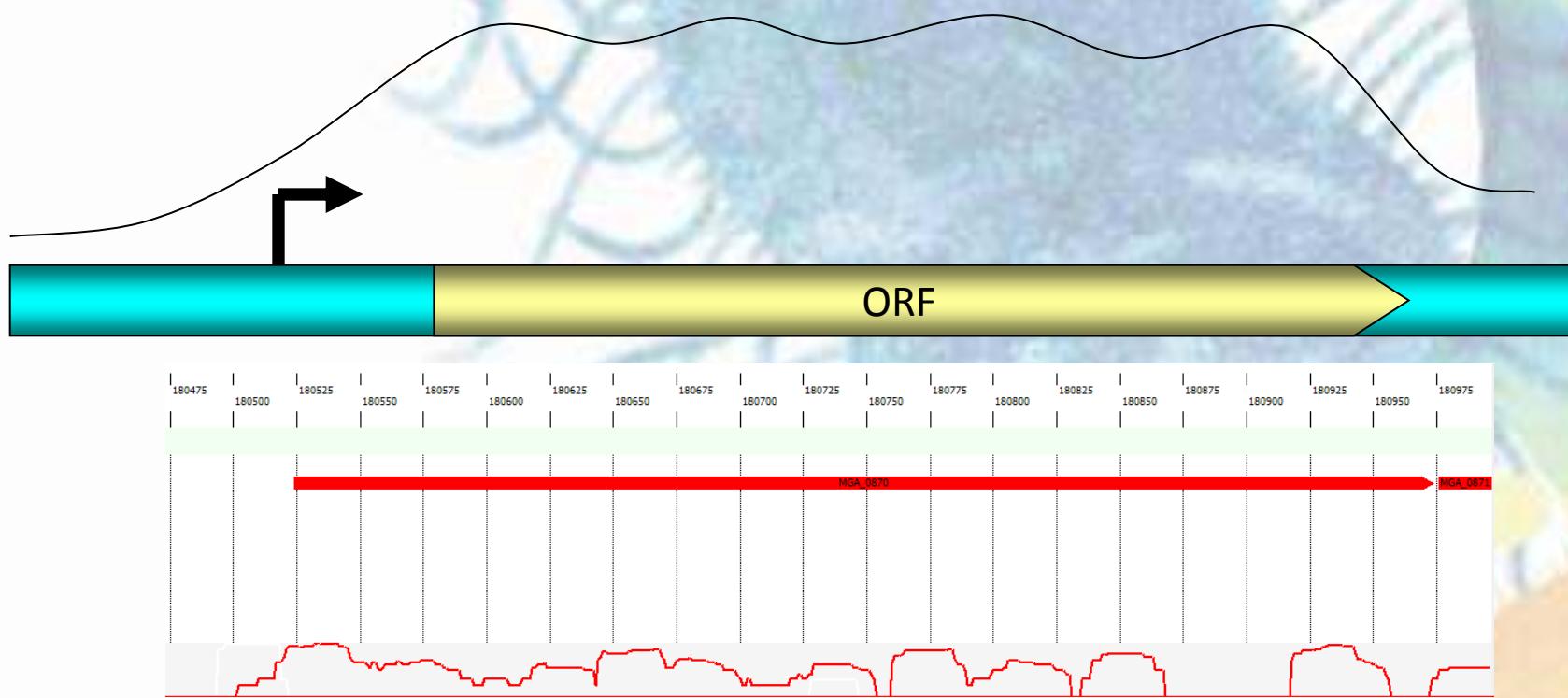
- Поиск сайтов связывания транскрипционных факторов

Технологическая проблема

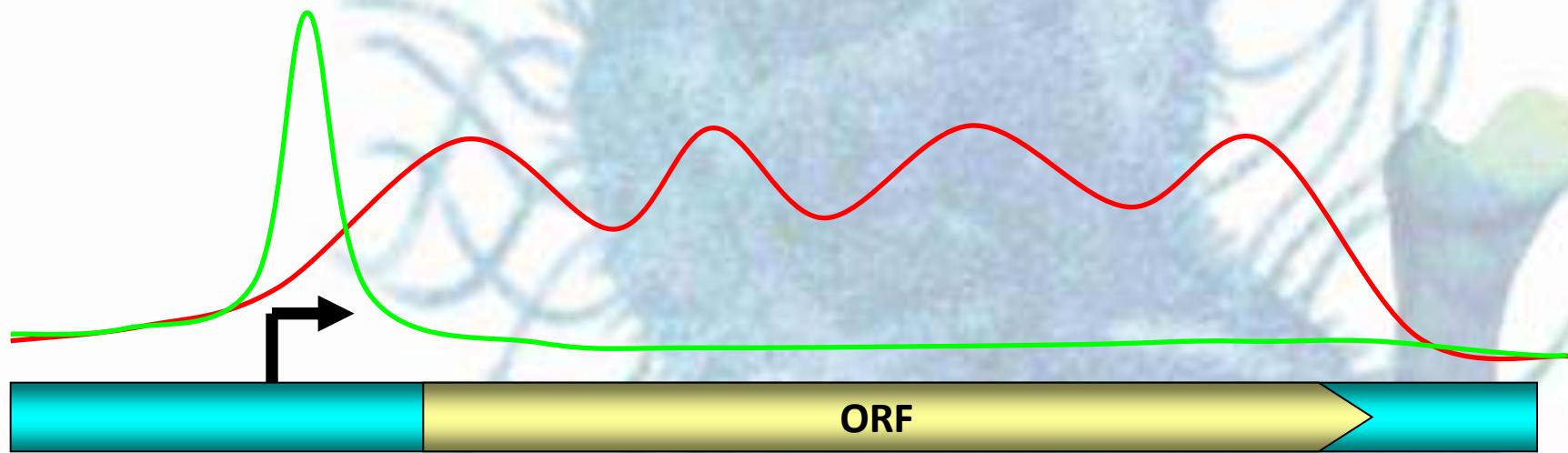


мРНК после фрагментации

Покрытие при стандартном приготовлении библиотеки



Идея метода



— Покрытие при стандартном приготовлении библиотеки

— Покрытие 5'-обогащённой библиотеки

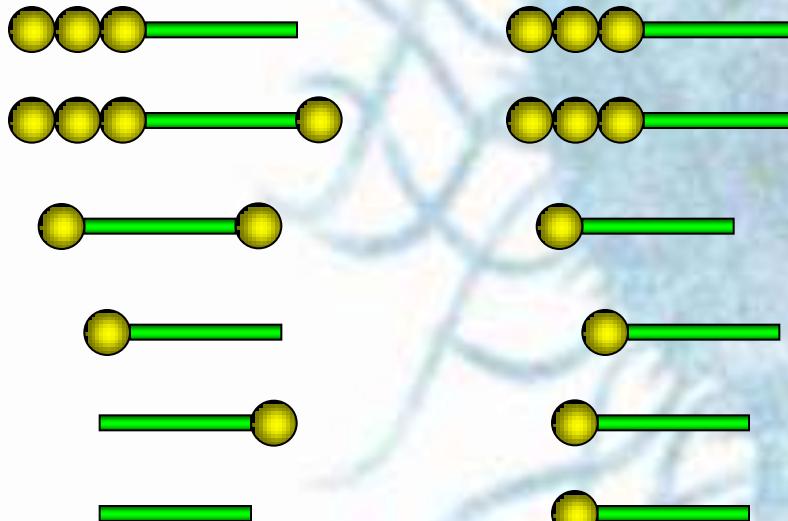
Шаг 1 - фрагментация



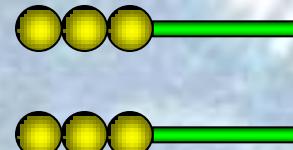
● = Фосфат

Шаг 2 – обогащение 5`-последовательностями

Починка концов

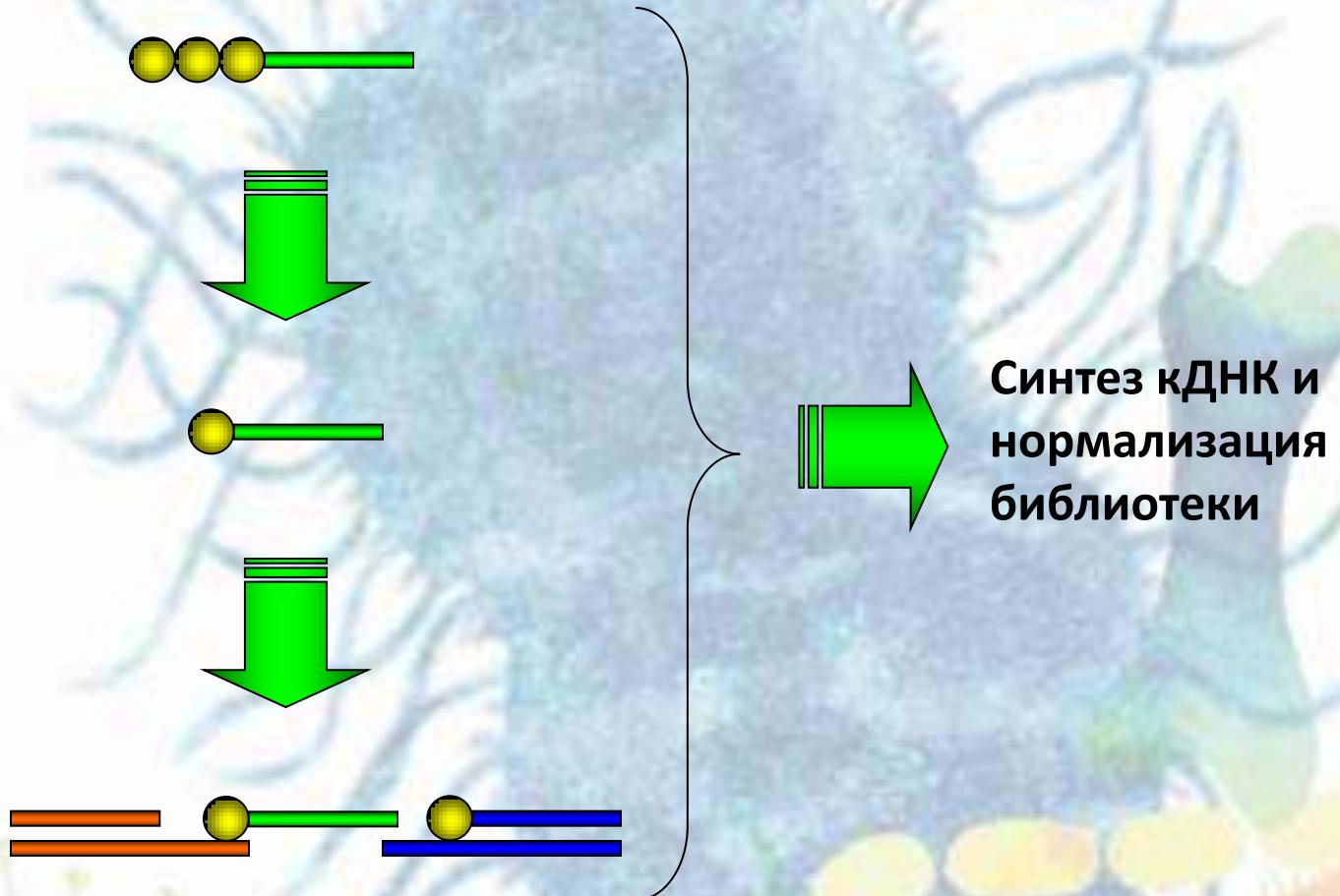


Обработка 5`-фосфат-зависимой нуклеазой



● = Фосфат

Шаг 3 – снятие защитной группы и лигирование адаптеров



○ = Фосфат

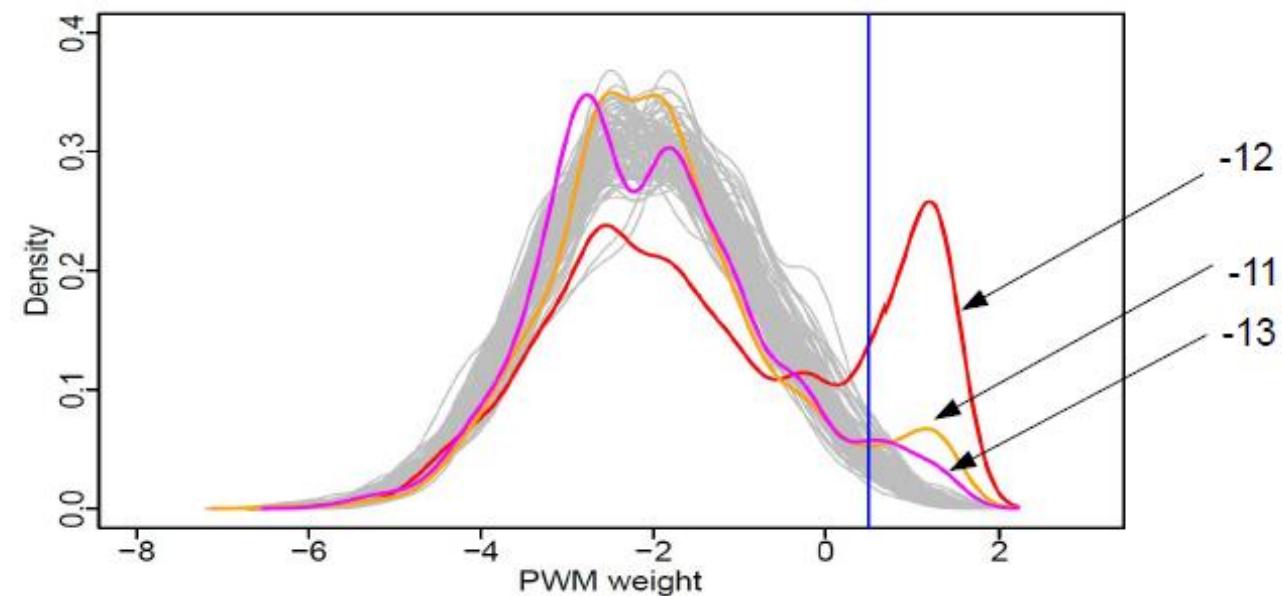


Шаг 4 - секвенирование

Шаг 5 – анализ результатов

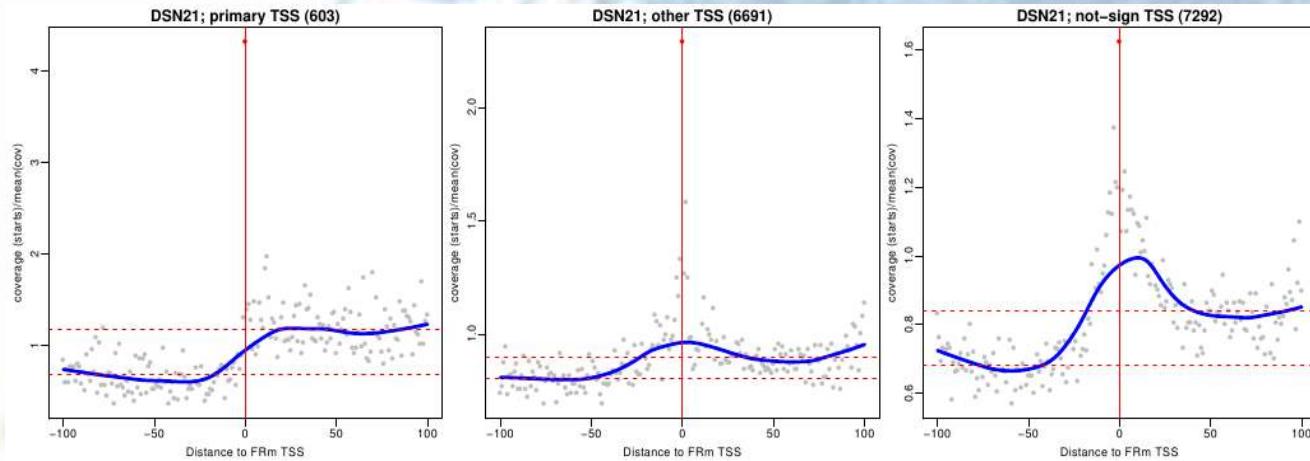
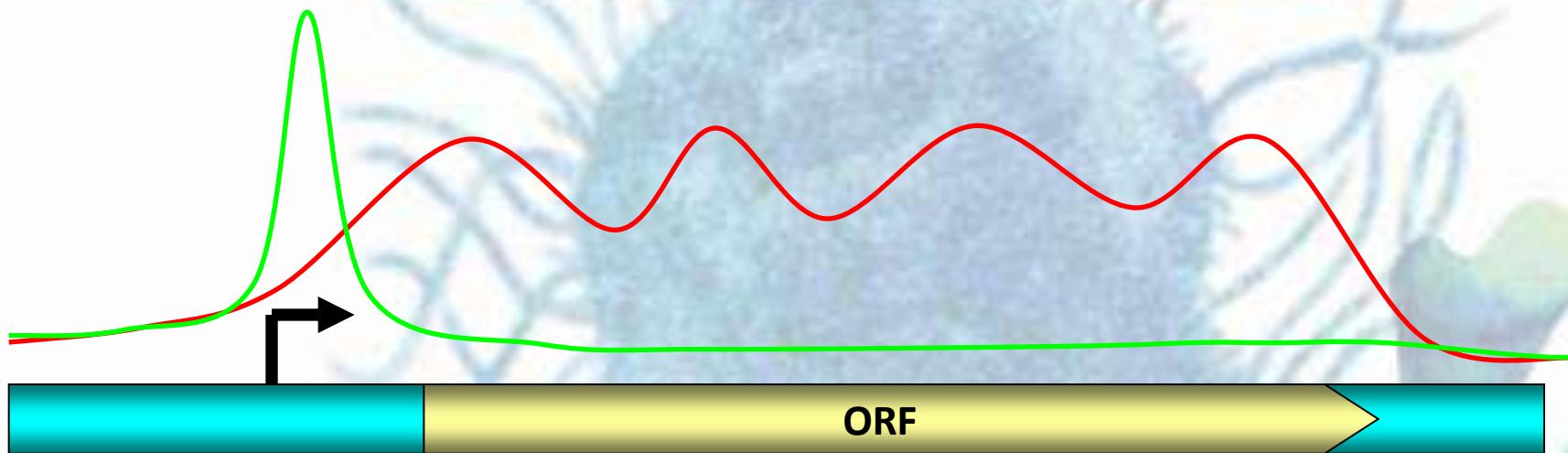


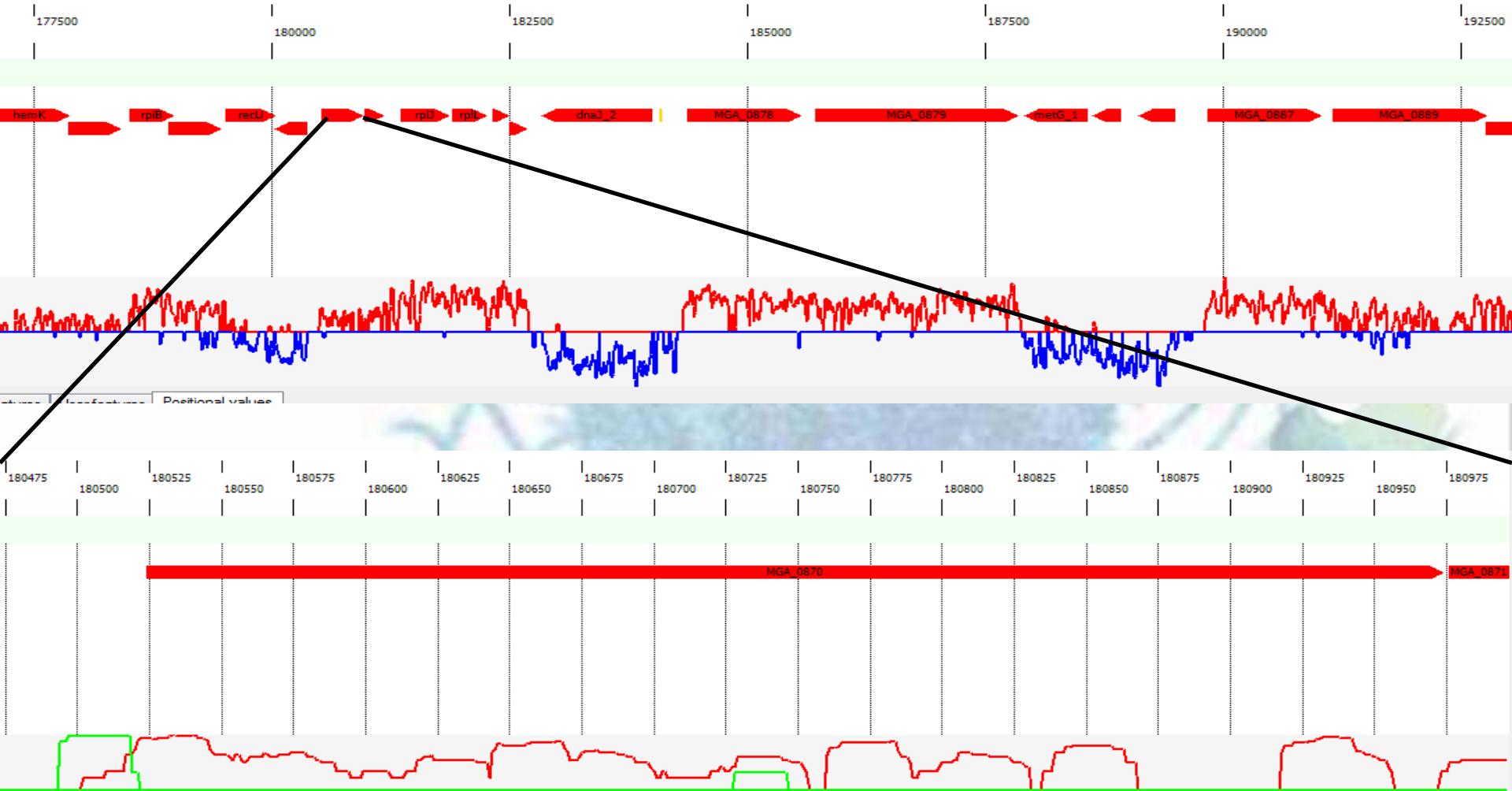
Веса ТАТА-боксов,
расположенных в
разных позициях
относительно старта



Шаг 5 – анализ результатов

Проверка по наличию ступеньки в покрытии



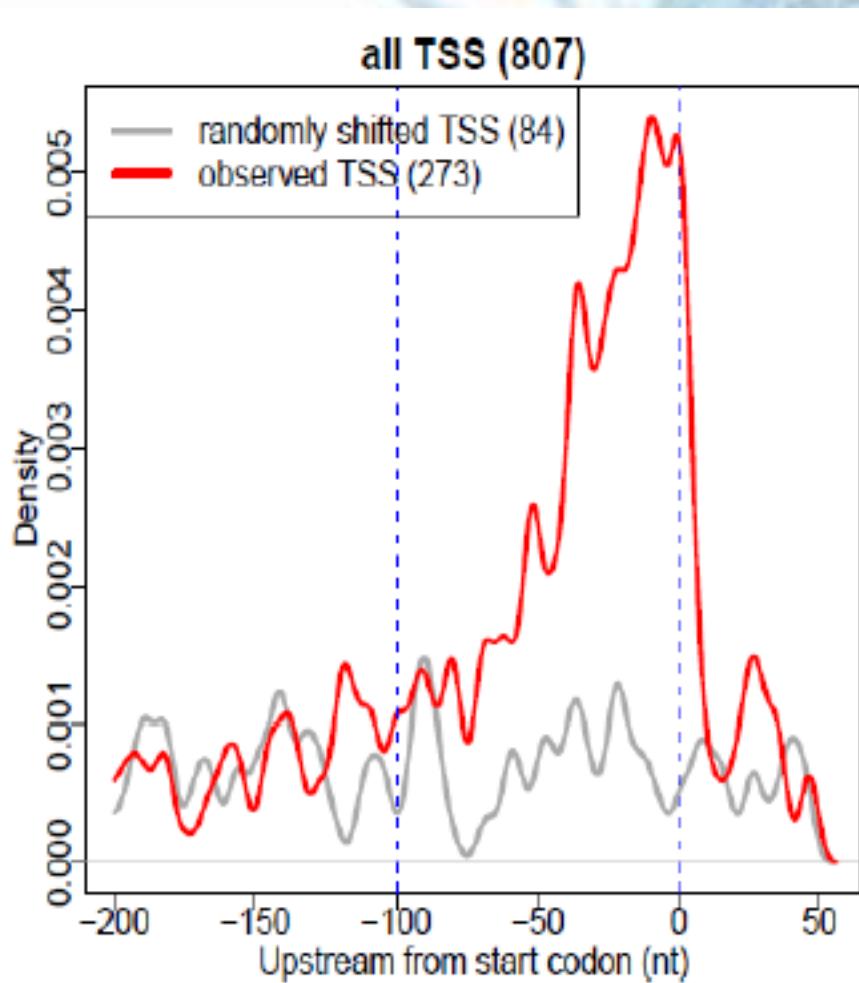


Покрытие по + цепи

Покрытие по - цепи

Покрытие по + цепи с обогащением 5'- концами

Результаты – как соотносятся старты транскрипции (TSS) и ORF (длина 5`-UTR)

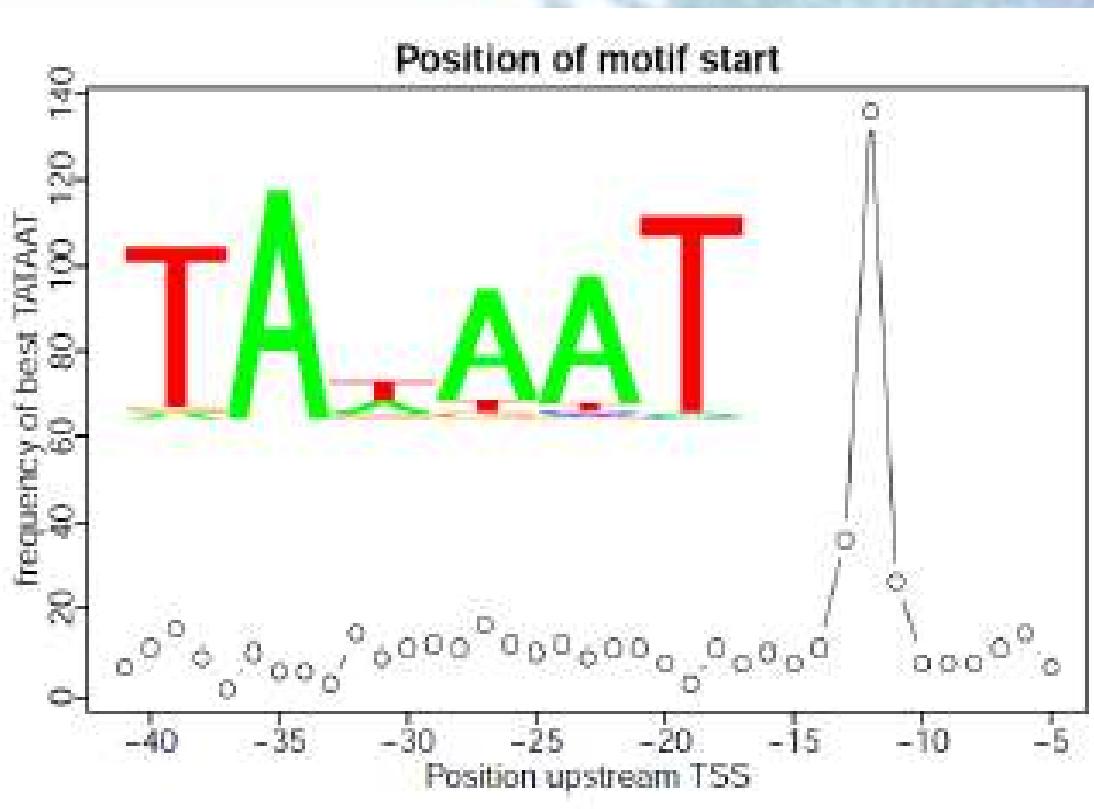


TSS в среднем
располагается от 0 до -50
bp от старта рамки

Результаты

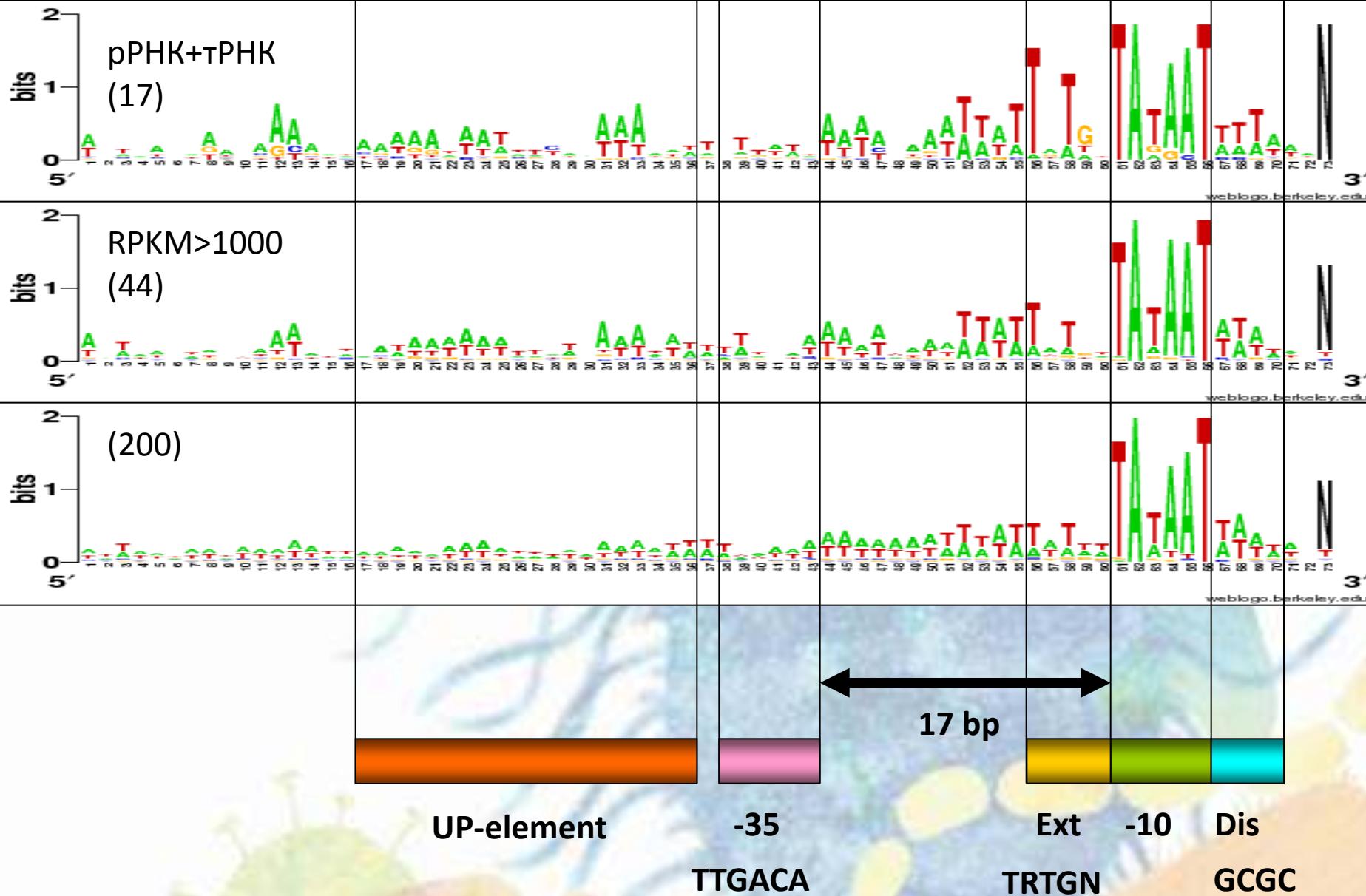
Всего идентифицировано 807 TSS

200 TSS имеют ТАТА-бокс, стоят вблизи рамки и имеют ступеньку в покрытии



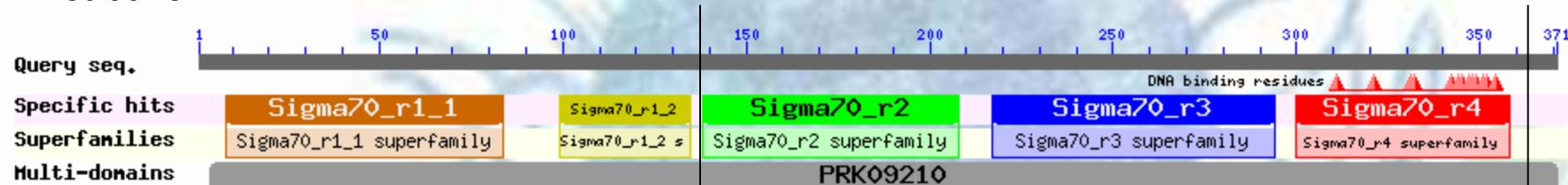
ТАТА-бокс стоит на позиции -12 от старта транскрипции

Структура микоплазменного промотора



Строение микоплазменного сигма-фактора

B. subtilis



M. gallisepticum



Выводы

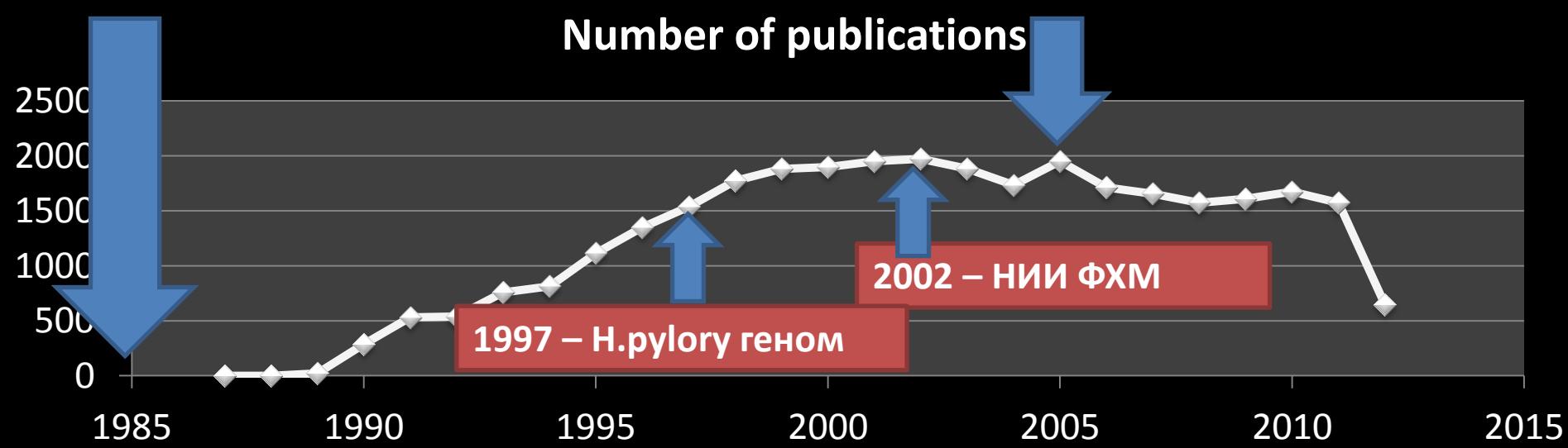
- Метод работоспособен
- Ключевые элементы микоплазменного промотора – ТАТА-бокс и Ext-элемент
 - Консенсус - WWWTTRTGNTAWAATWWW
- Консервативный -35 элемент отсутствует
- Длина 5`-UTR порядка 50 bp



Открытие бактерии

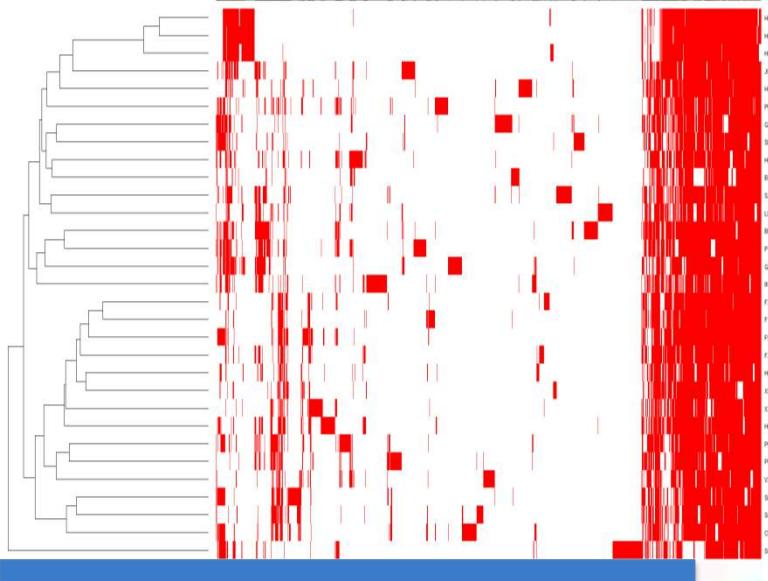
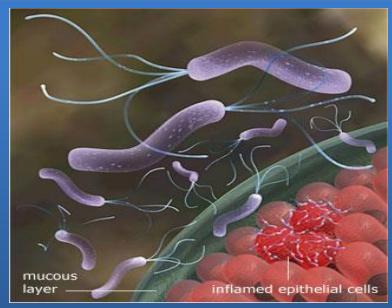
J Robin Warren, Barry Marshall

Нобелевская премия

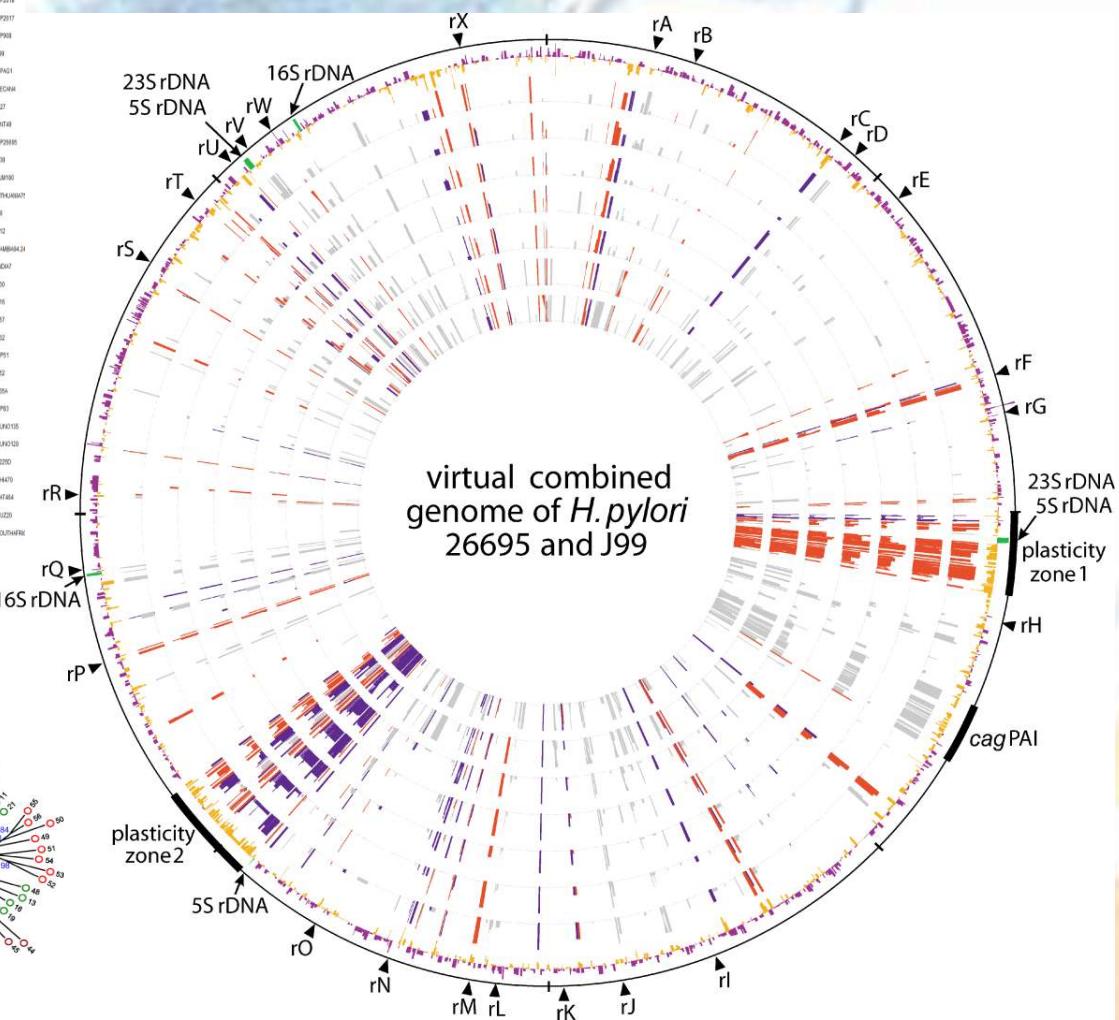
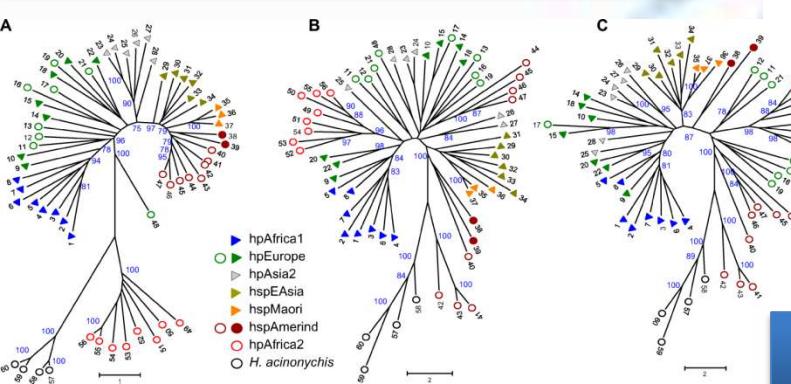


Micro and macro heterogeneity

Helicobacter pylori

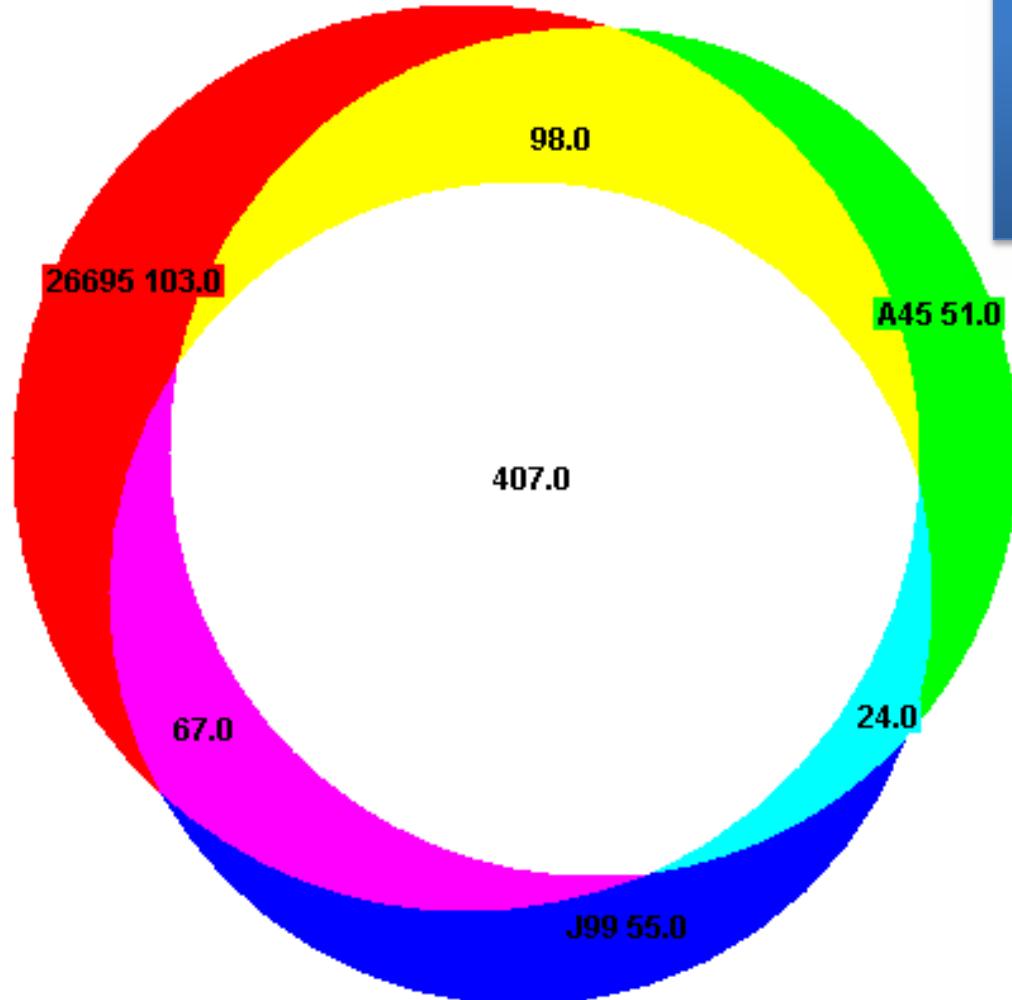


3186 genes in 31 strain
having 1051 shared gene

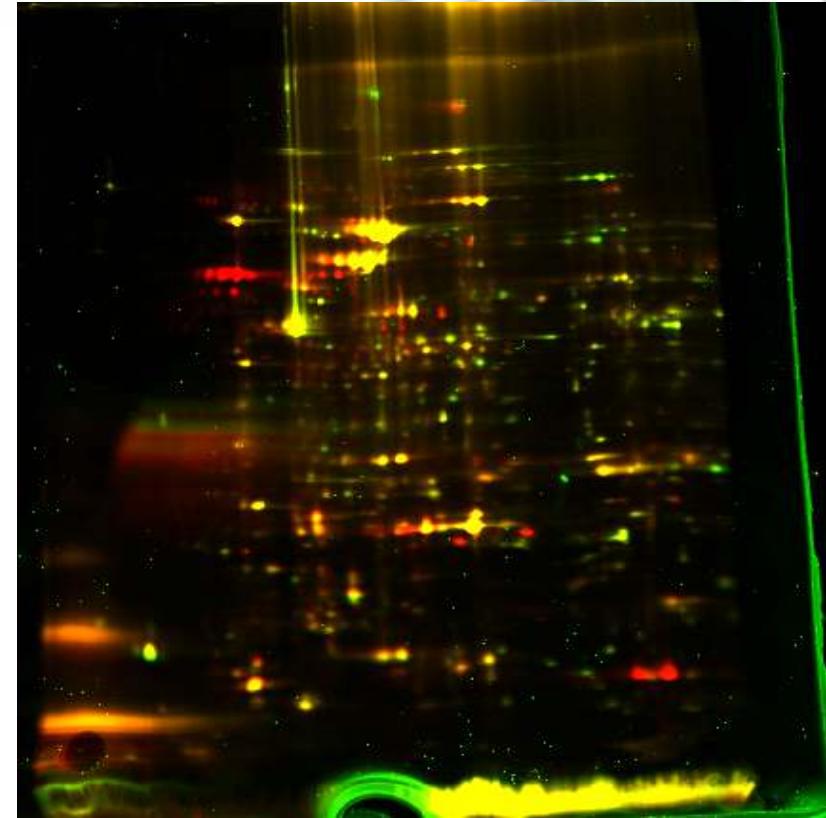


Phylogeny based on different genes

Протеомные карты *Helicobacter pylori* – первый взгляд

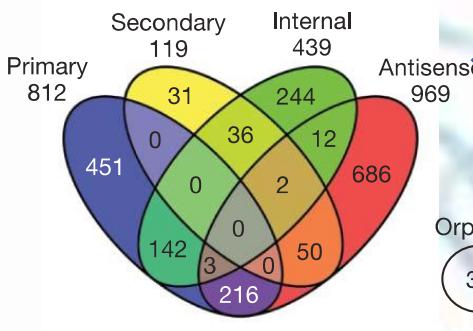
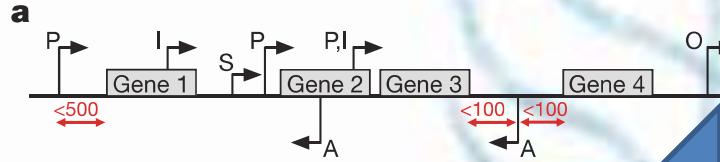


26695 – 675 белков
J99 – 553 белков
A45 – 580 белков





Транскриптомное профилирование бактерий

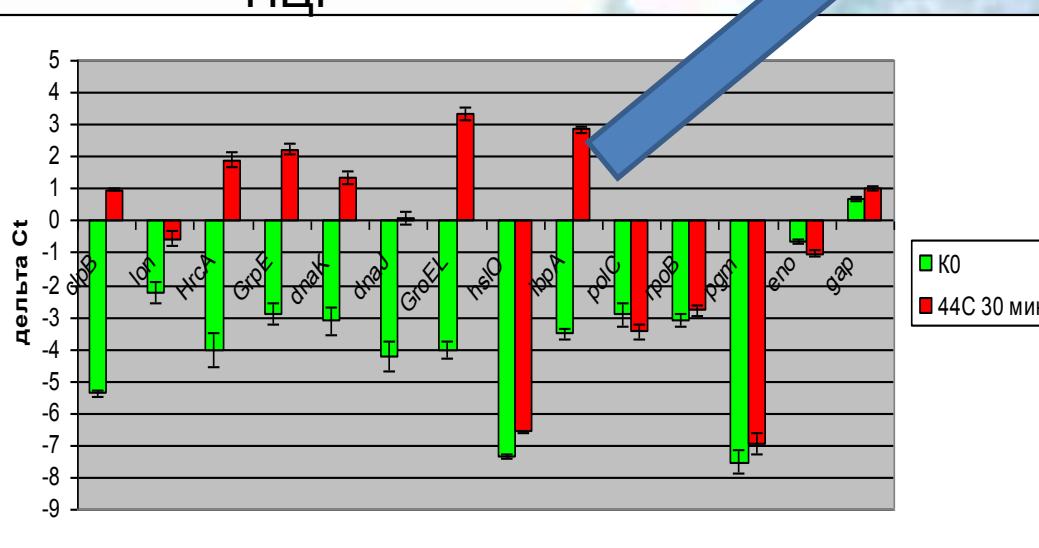


Более половины РНК
- Антисмысловое



Полное
транскрипционное
профилирование

ПЦР



Метилирование

Helicobacter pylori - Phasevariation



a45	J99	26695
ACGT	ACGT	ACGT
	ACNGT	ACNGT
ATTAAT	ATTAAT	ATTAAT
CATG	CATG	CATG
CCATC		
CCGG	CCGG	CCGG
	CCNNGG	
CCTC	CCTC	CCTC
CCTTC		CCTTC
	CGWCG	
		CTGCAG
		GAAGA
GANTC	GANTC	GANTC
GATC	GATC	GATC
GCGC	GCGC	GCGC
GGCC		
GRRG	GRRG	GRRG
GTAC	GTAC	GTAC
GTNNAC	GTNNAC	GTNNAC
	GTSAC	
TCGA	TCGA	TCGA
TCNGA		
TCNNGA	TCNNGA	TCNNGA
TGCA		TGCA

Although we cannot exclude another as yet undescribed role for these modH phase variation in *H. pylori* biology, we have confirmed phasevarion mediated epigenetic mechanism of gene expression does operate in *H. pylori* - **Phasevarion Mediated Epigenetic Gene Regulation in *Helicobacter pylori* – Plos One – 12.2011**

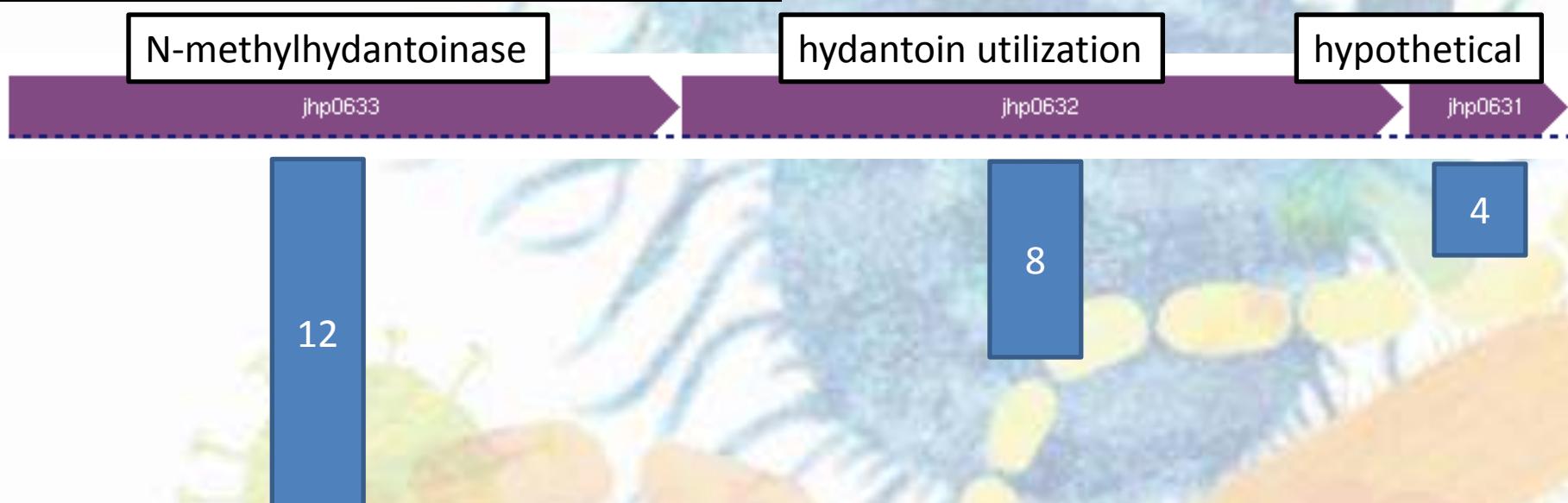
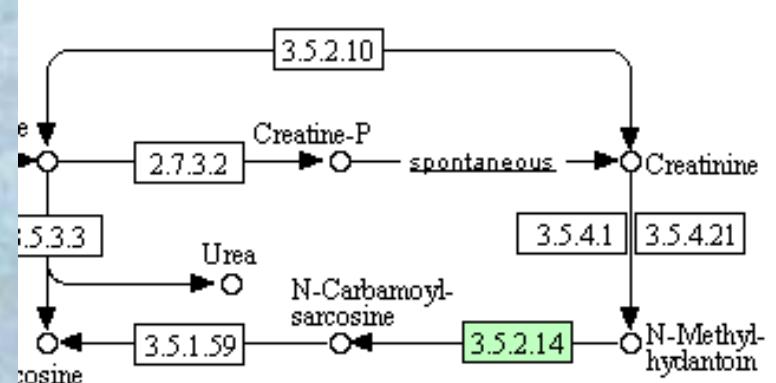
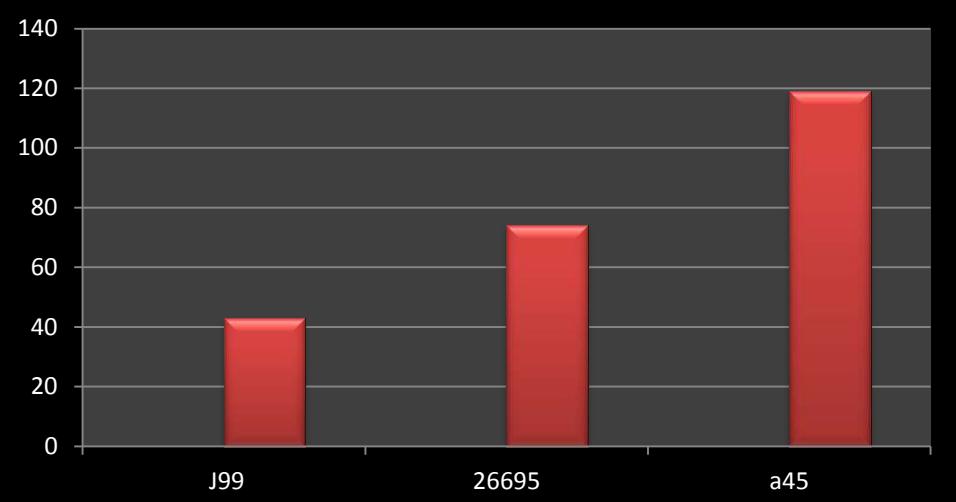


Метилирование

Helicobacter pylori - Phasevariation

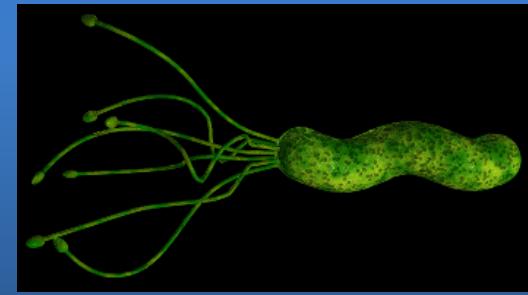


Motive TGCA – 158 protein

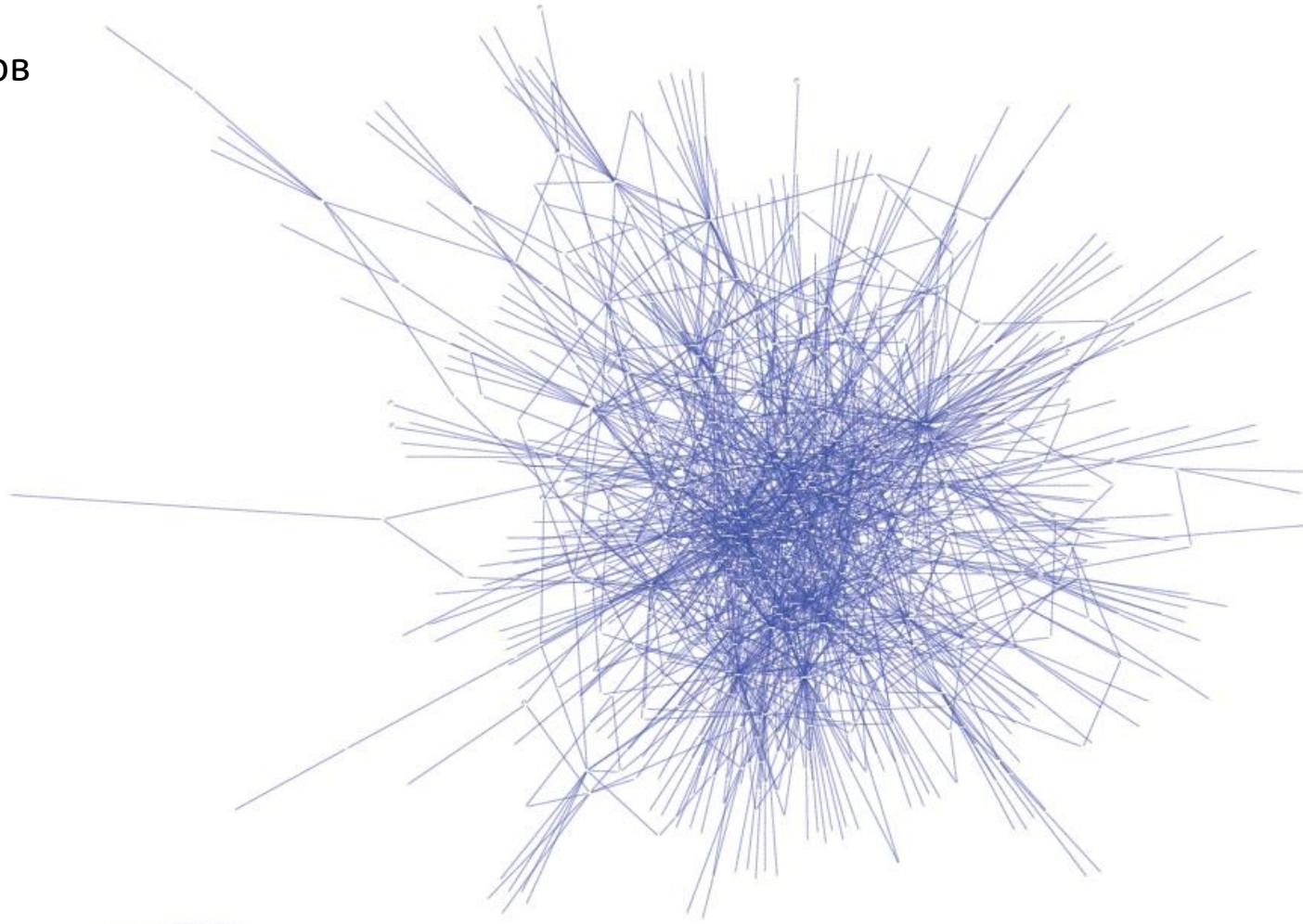


Интерактом

Helicobacter pylori



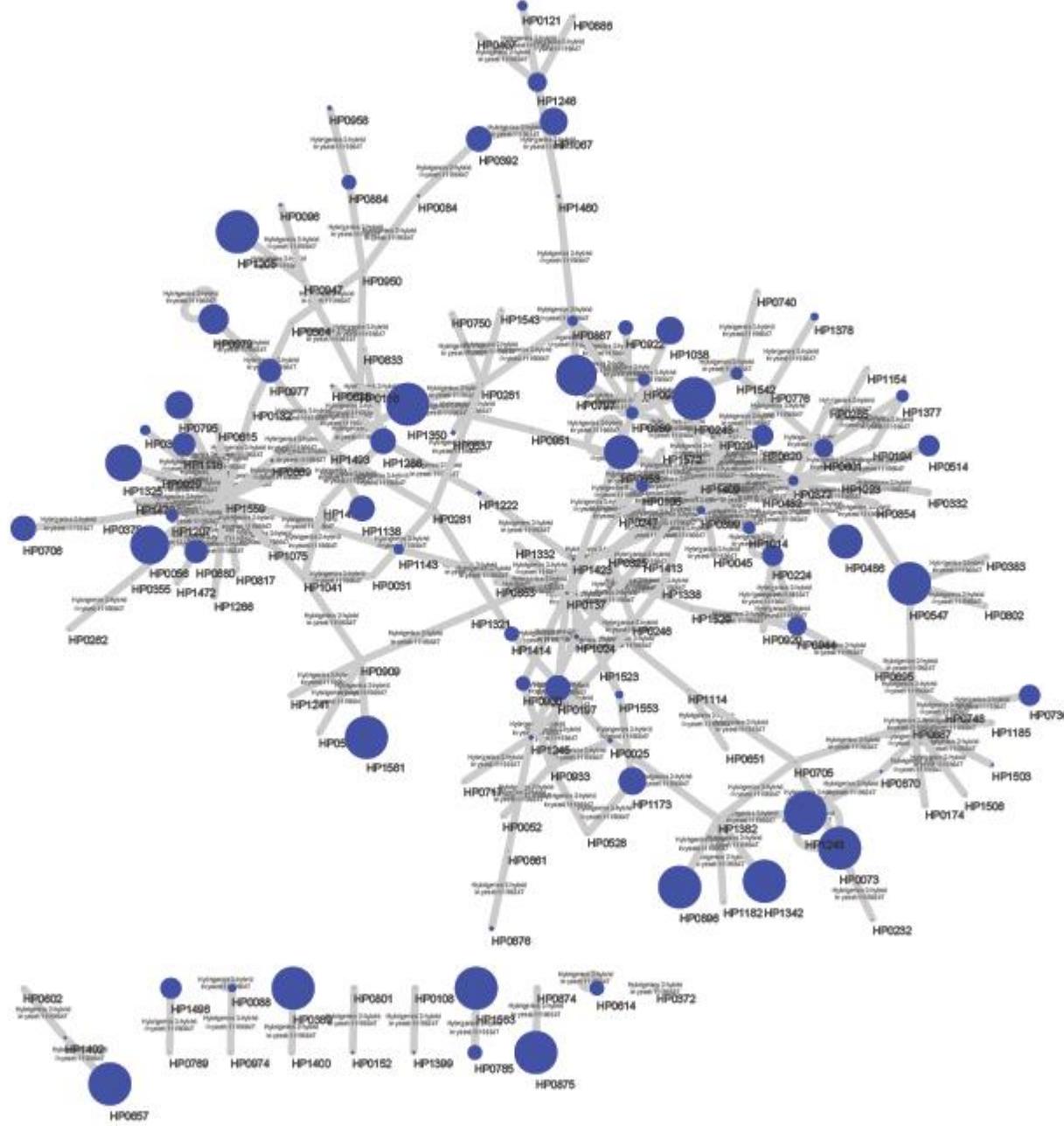
1200 белков



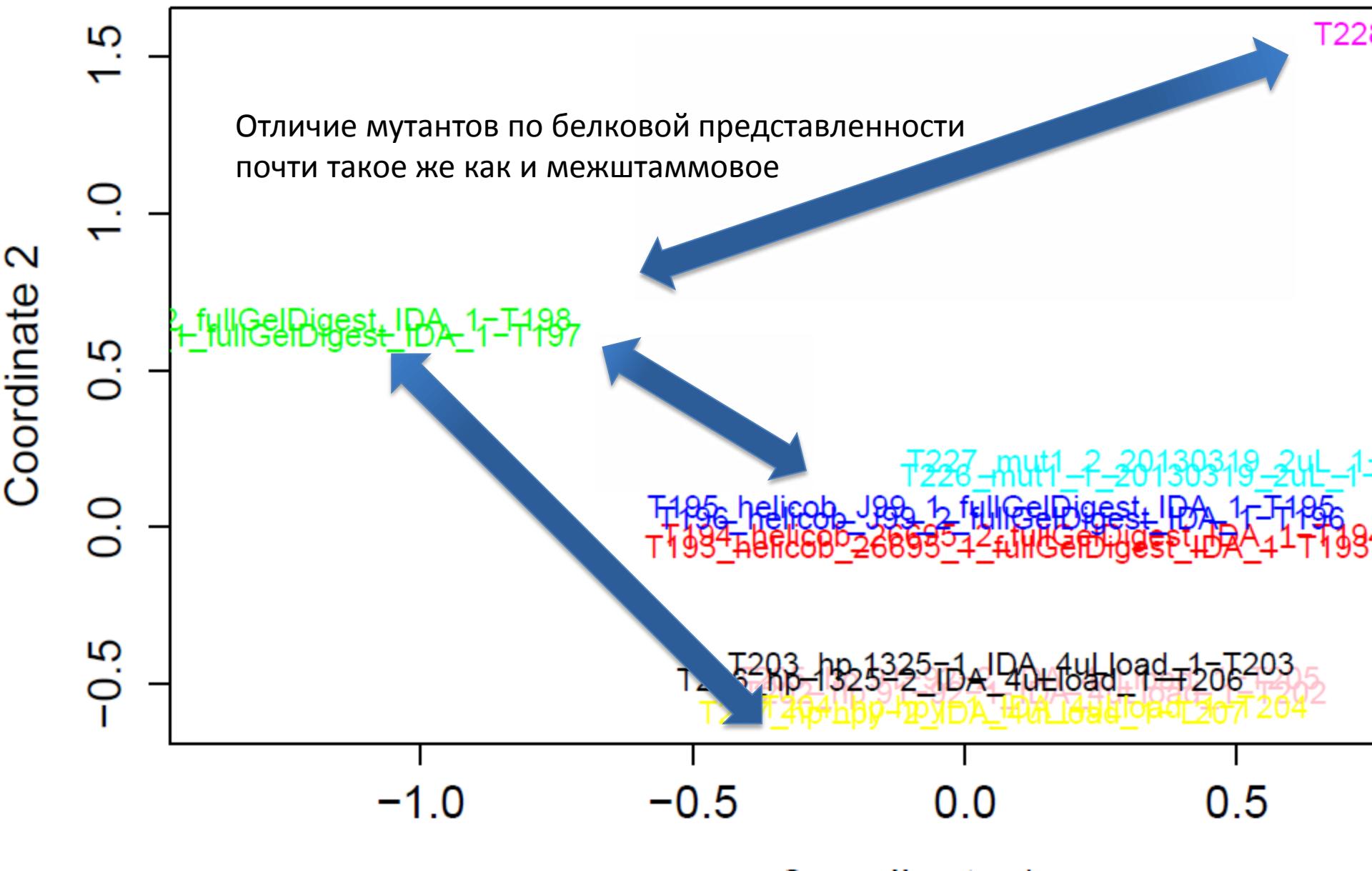
26695

A45

J99



Nonmetric MDS



Благодарю за внимание

