Занятие №9. Обработка данных секвенирования

Покрытие

- Покрытие (глубина секвенирования) важный параметр методов NGS: кратность прочтения каждого нуклеотида. Для каждой задачи необходимо своё покрытие (обычно устанавливают не менее, чем 30-тикратное покрытие).
- Таким образом, "эффективный" объём данных равен выходу секвенирования, делённому на покрытие.

Оценка необходимого покрытия

 Вероятность того, что нуклеотид не будет определён (Р), исходя из глубины покрытия (с) вычисляется по формуле Ландела–Ватермана:

- Теоретически достаточное покрытие должно позволять определить все нуклеотиды в геноме длиной L (P*L<1).
- Например, для генома человека (L=3*10⁹ п.о.) теоретически достаточно 23-кратного покрытия

Анализ данных секвенирования

- . 1. Очистка "сырых" данных (raw data) (фильтрация ридов по качеству).
- <u>Результат:</u> "примесные" риды удаляются, в остальных обрезаются неточно определённые нуклеотиды
- 2. Сборка генома (слияние ридов для коротких фрагментов) с помощью специальной программы ассемблера.
- <u>Результат:</u> набор длинных фрагментов (контиги) или их упорядоченная последовательность, образующая скэффолд.
- 3. Интерпретация данных (аннотация)
- поиск кодирующих последовательностей и их структурное и

1. Оценка качества ридов: FASTQ – формат записи ридов

@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTC
AACTCACAGTTT

```
!"*((((***+))%%%++)(%%%%).1***-+*"))**55CCF>>>>>CCCCCCC65
```

Каждая последовательность занимает 4 строки:

- первая начинается с @ и содержит название и описание последовательности;
- вторая содержит последовательность (знаки A,G,C,T);
- третья начинается с + и может содержать примечания (технические комментарии секвенирования);
- четвёртая содержит столько же символов, что и вторая, каждый символ указывает вероятность ошибочного определения соответствующего нуклеотида по шкале Phred.

Определение качества ридов по шкале Phred

Каждый символ означает какое-то число (Q) от 0 до 100.
 Вероятность ошибочного определения нуклеотида (Р),
 качество которого оценивается как Q равна:

•
$$P = 10^{-Q/10}$$

"Хорошее" качество при Q≥30 (Р≤0,001=0,1%)

Phred и ASCII:

номера присваивают начиная с 33 символа (!=0) (Phred+33) или с 64 (@=0) (Phred+64)

ASCII control characters				ASCII printable characters					Extended ASCII characters							
00	NULL	(Null character)	32	space	64	@	96	•	128	ç	160	á	192	L	224	Ó
01	SOH	(Start of Header)	33	1	65	A	97	а	129	ü	161	í	193		225	ß
02	STX	(Start of Text)	34		66	В	98	b	130	é	162	ó	194	1.00	226	ô
03	ETX	(End of Text)	35	#	67	C	99	C	131	â	163	ú	195	-	227	ò
04	EOT	(End of Trans.)	36	\$	68	D	100	d	132	ä	164	ñ	196		228	õ
05	ENQ	(Enquiry)	37	%	69	E	101	e	133	à	165	Ñ	197	+	229	õ
06	ACK	(Acknowledgement)	38	&	70	F	102	f	134	à	166	3	198	ã	230	μ
07	BEL	(Bell)	39	·	71	G	103	g	135	ç	167	0	199	Ã	231	þ
08	BS	(Backspace)	40	(72	Н	104	h	136	ê	168	ė	200	Ĩ.	232	Þ
09	HT	(Horizontal Tab)	41)	73	ï	105	ï	137	ë	169	®	201	F	233	Ú
10	LF	(Line feed)	42	*	74	j	106	- 1	138	è	170	7	202	1[234	Û
11	VT	(Vertical Tab)	43	+	75	K	107	k	139	ï	171	1/2	203	75	235	ù
12	FF	(Form feed)	44		76	L	108	î	140	î	172	1/4	204	1	236	
13	CR	(Carriage return)	45		77	M	100	m	141	ì	173	i	205	r =	237	Ý
14	SO	(Shift Out)	46		78	N	110	n	142	Ä	174	**	206	#	238	- 1
15	SI	(Shift In)	47	i	79	O	111	0	143	A	175	>>	207	ר ם	239	
16	DLE	(Data link escape)	48	0	80	P	112	р	144	É	176		208	ð	240	=
17	DC1	(Device control 1)	49	1	81	Q	113	q	145	æ	177		209	Đ	241	±
18	DC2	(Device control 2)	50	2	82	R	114	Г	146	Æ	178		210	Ê	242	÷
19	DC3	(Device control 3)	51	3	83	S	115	S	147	ô	179	=	211	Ë	243	3/4
20	DC4	(Device control 4)	52	4	84	T	116	t	148	ö	180		212	È	244	1
21	NAK.	(Negative acknowl.)	53	5	85	Ü	117	u	149	ò	181	Ā	213	-	245	§
22	SYN	(Synchronous idle)	54	6	86	V	118	v	150	û	182	Â	214	i	245	÷
23	ETB	(End of trans, block)	55	7	87	w	119	w	151	ù	183	À	215	î	247	
24	CAN	(Cancel)	56	8	88	X	120	X	152	ÿ	184	©	216	i	248	0
25	EM	(End of medium)	57	9	89	Ŷ	121		153	Ö	185	4	217	- 1	249	
26	SUB	(Substitute)	58	:	90	Z	122	y z	154	Ü	186	1	218		250	
27	ESC	(Substitute) (Escape)	59		91	[123	{	154	Ø	187	-	219	Г	251	1
28	FS	(Escape) (File separator)	60	<	92	1	124	1	156	£	188	1	220		252	3
29	GS	(Group separator)	61	=	93		125		157	Ø	189	¢	221	1	252	2
30	RS	CONTRACTOR OF THE PROPERTY OF	62	>	94]	126	}	157	×	190	¥	222	-	253	
31	US	(Record separator)	63	?	95	- S	120	~			190		222		255	nhon
127		(Unit separator)	03	- 1	95	<u> </u>			159	f	191	7	223		255	nbsp
121	DEL	(Delete)							-							

Кодировка качества Phred+33

- Качество в формате FASTQ закодировано в ASCII символах
- Существовало несколько стандартов записи качества
- Многим программам важно, чтения во Phred+33 или Phred+64



@M03106:3:000000000-AAWHP:1:1101:17129:1029 1:N:0:84

Базовая таблица кодировки ASCII и ее соответствие качеству прочтения

Q	Шкала	символ	Q	Шкала	символ	Q	Шкала	символ
0	33	1	20	53	5	40	73	I
1	34	cc	21	54	6	41	74	J
2	35	#	22	55	7	42	75	K
3	36	\$	23	56	8	43	76	L
4	31	%	24	57	9	44	77	M
5	38	&	25	58	:	45	78	N
6	39	•	26	59	÷	46	79	0
7	40	(27	60	<	47	80	P
8	41)	28	61	=	48	81	Q
9	42	*	29	62	>	49	82	R
10	43	+	30	63	?	50	23	S
11	44	,	31	64	@	51	84	T
12	45	-	32	65	A	52	85	U
13	46		33	66	В	53	86	V
14	47	1	34	67	C	54	87	W
15	48	0	35	68	D	55	88	X
16	49	1	36	69	Е	56	89	Y
17	50	2	37	70	F	57	90	Z
18	51	3	38	71	G	58	91	[
19	52	4	39	72	Н	59	92	١

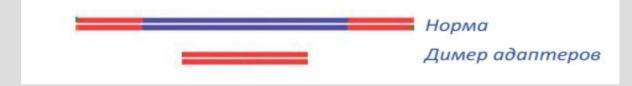
Программа FastQC http://www.bioinformatics.babraham.ac.uk/

Примеры качества по шкалам Phred+33 и Phred+64

Качество по Phred, Q	Символ ASCII (Phred33)	Символ ASCII (Phred64)	Вероятность ошибки	Точность

Источники ошибок в ридах: примеси

- Примеси бывают:
- 1. Артефактные (ошибки секвенирования)
- образование димеров адаптеров



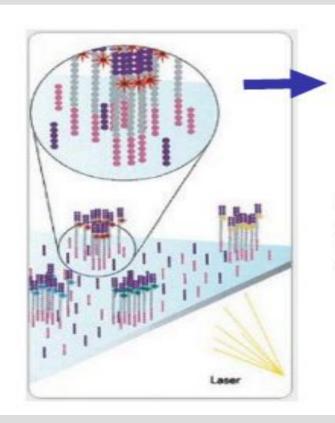
• чтение сквозь – вставки слишком короткие



• 2. Биологические – контаминация

Источники ошибок в ридах: фазировка

 Фрагменты в одном кластере строятся с разной скоростью – секвенатору сложно определить верный нуклеотид.



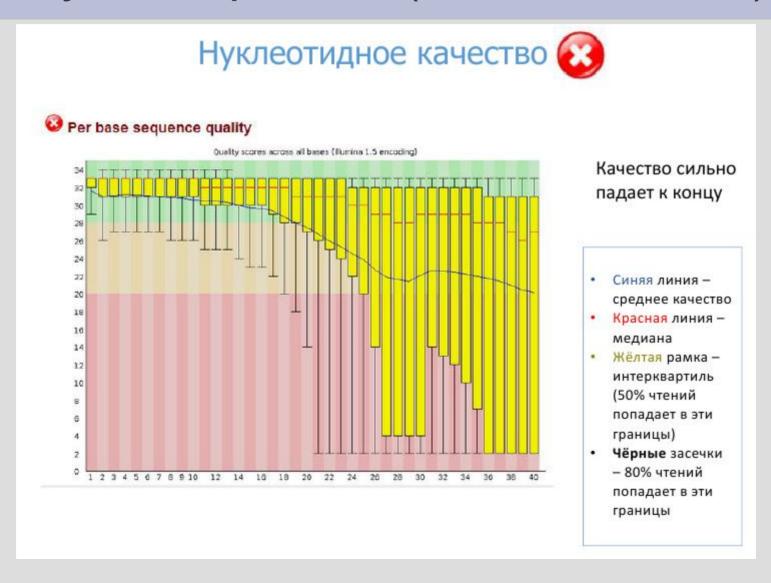
CGCAAGGAAGCTGATTG CGCAAGGAAGCTGATTG CGCAAGGAAGCTGATT CGCAAGGAAGCTGATTGC CGCAAGGAAGCTGATTGC

Чем дольше идёт прогон, тем больше будет накапливаться отстающих и опережающих олигонуклеотидов

Программа FastQC – контроль качества ридов: 1. Среднее нуклеотидное качество – хорошее (все Me>25, все Q1>10)



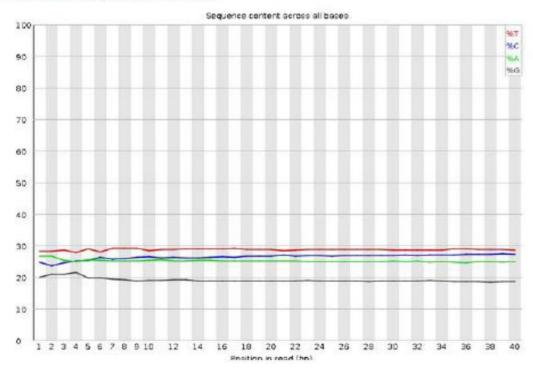
Программа FastQC – контроль качества ридов: 1. Среднее нуклеотидное качество – неудовлетворительное (есть Me<20 или Q1<5)



Программа FastQC – контроль качества ридов: 2. Средний нуклеотидный состав ридов

Нуклеотидный состав

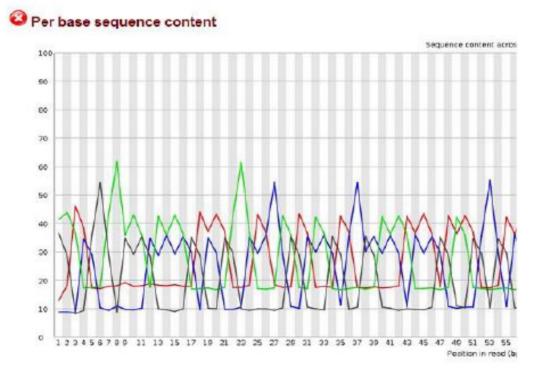
Per base sequence content



- Нуклеотидный состав примерно одинаковый на протяжении всего рида
- В начале чтения небольшие колебания из-за специфичности фрагментации программа обозначает их восклицательным знаком, но это не проблема

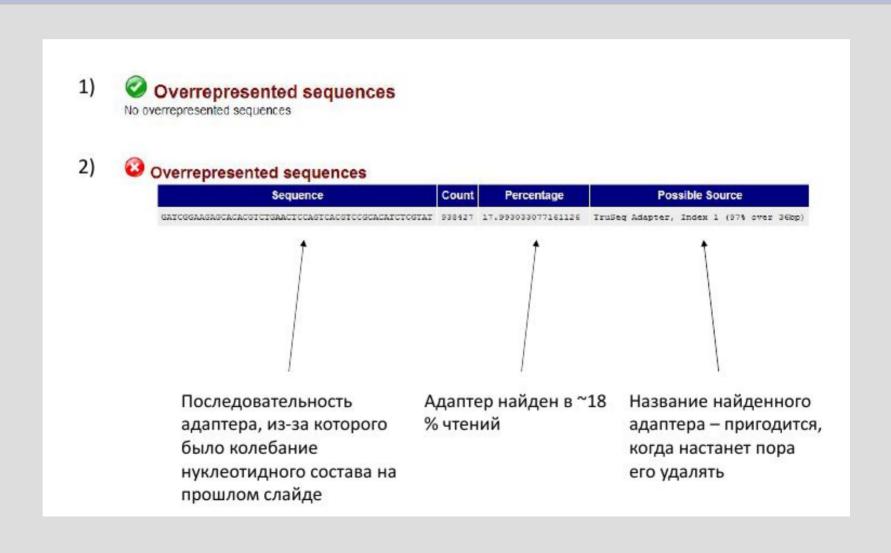
Программа FastQC – контроль качества ридов: 2. Средний нуклеотидный состав ридов

Нуклеотидный состав 🔞



Сильные колебания нуклеотидного состава – скорее всего отсеквенированы димеры адаптеров

Программа FastQC – контроль качества ридов: 3. Чрезмерно представленные последовательности



Очистка "сырых" ридов: тримминг

- 1. Удаление адаптерных последовательностей из ридов
- 2. Отсечение с конца ридов нуклеотидов, качество которых ниже определённого уровня (Q<20 или Q<30)

• Инструмент для тримминга: программа Trimmomatic

Особый этап для метагеномики – Сортировка данных (биннинг)

- 1. Методы, основанные на нуклеотидном составе
- GC-состав
- динуклеотидный состав
- тринуклеотидный состав
- тетрануклеотидный состав

- 2. Методы, основанные на гомологии
- сравнение с базой данных

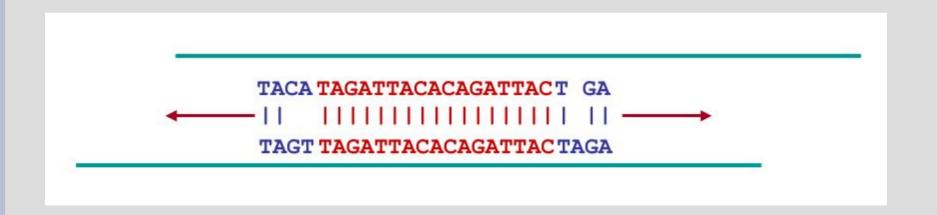
2. Сборка генома (assemby)

- de novo (сборка не секвенированного ранее генома)
- метод OLC (overlap layout concensus) (перекрытие фрагментов) для малого количества длинных фрагментов (Sanger)
- графы де Брёйна для большого количества коротких фрагментов (NGS)
- сборка генома, аналогичного ранее собранному (ресеквенирование) референсному геному (выравнивание на геном, alignment)
- хэш-таблицы

CVMMINCHI IO DODODI O

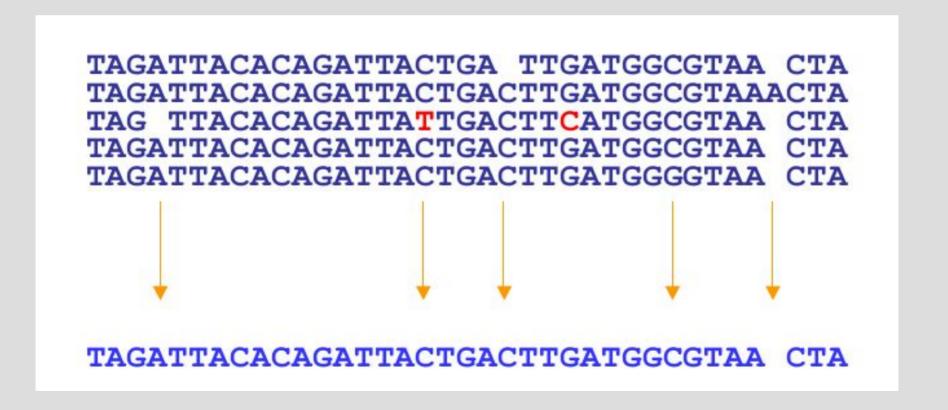
Сборка de novo: Overlap layout consensus: 1

 Поиск пар ридов, имеющих общие k-меры (последовательности длиной k, k=24), смещение двух строк относительно друг друга (выравнивание) до максимального совмещения (≥95% сходства)

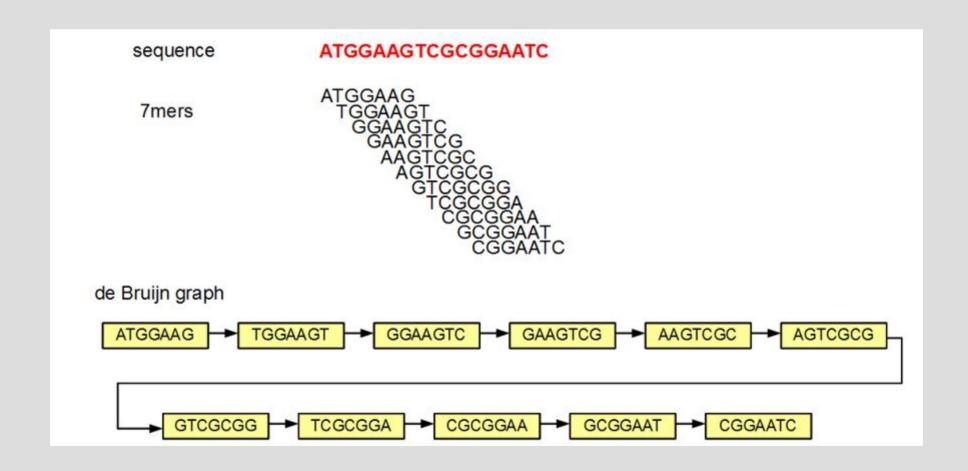


Сборка de novo: Overlap layout consensus: 2

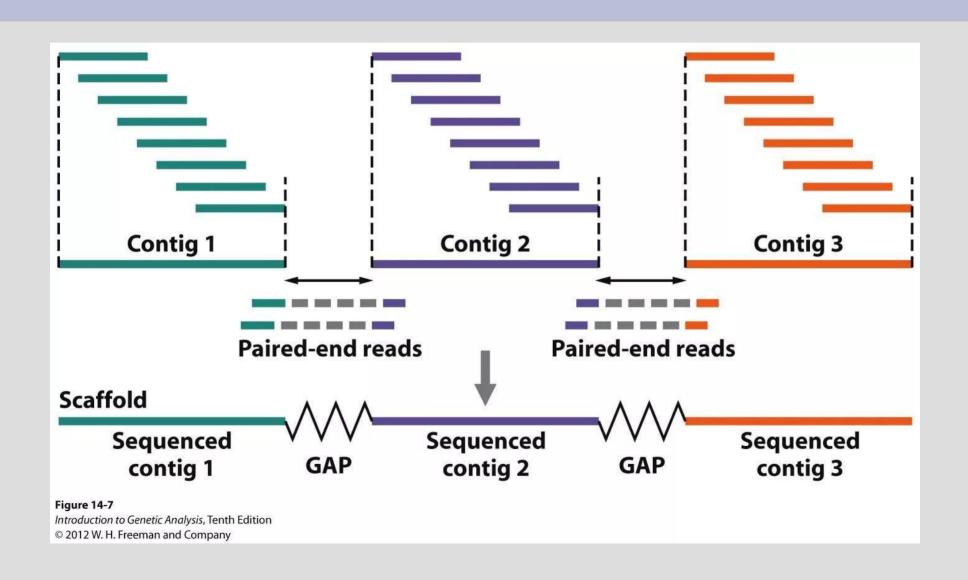
 На базе попарного выравнивания строят множественное выравнивание, корректируют ошибки



Сборка de novo: Графы де Брёйна



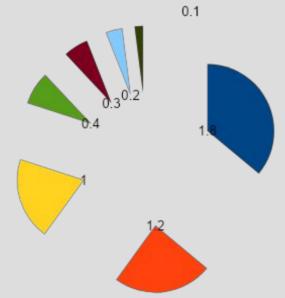
Результат сборки: контиги и скэффолды

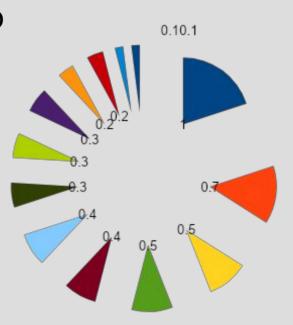


Качество сборки генома

- N50 длина контига, который вместе с остальными контигами большей длины покрывает не менее 50% генома (обычно под геномом понимают суммарную длину всех контигов).
- L50 число контигов не меньших чем N50.

Пример: две сборки генома длиной 5 Мb





Формат представления нуклеотидных последовательностей – FASTA

- >OTU-160-1 Acinetobacter baumannii
- CCTACGGGGGGCTGCAGTGGGGAATATTGGACAATGGGGGGA ACCCTGATCCAGCCATGCCGCGTGTGTGAAGAAGGCCTTATG GTTGTAAAGCACTTTAAGCGAGGAGGAGGCTACTCTAGTTAAT ACCTAGGGATAGTGGACGTTACTCGCAGAATAA
- Каждая последовательность занимает две строки:
- 1). первая строка начинается со знака > и содержит идентификатор (за которым эта последовательность закреплена в некоторой базе данных), через пробел следует опциональное словесное описание;
- 2). вторая строка сама последовательность нуклеотидов.

3. Аннотация

- 1. Поиск белок-кодирующих последовательностей
- на основе гомологии сравнение с уже известными генами
- аннотация ab initio статистический поиск по характерным для белок-кодирующих участков последовательностям (АТС....)
- 2. Поиск других кодирующих последовательностей (гены РНК)

Результаты аннотации

- 1. Структурное описание:
- открытые рамки считывания (ORF) и из расположение
- – структура гена
- кодирующие области
- локализация регуляторных последовательностей
- 2. Функциональное описание
- биохимическая функция белкового продукта
- биологическая функция белка
- экспрессия белка
- участие белка в регуляторных и межбелковых

DOGUMA DOŬATRIAGY