

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 (P), исходя из глубины покрытия (c) вычисляется по формуле Ландела–Ватермана:

- $P=e^{-c}$

- Теоретически достаточное покрытие должно позволять определить все нуклеотиды в геноме длиной L ($P*L < 1$).
- Например, для генома человека ($L=3*10^9$ п.о.) теоретически достаточно 23-кратного покрытия

Анализ данных секвенирования

- **1. Очистка “сырых” данных (raw data) (фильтрация ридов по качеству).**
Результат: “примесные” риды удаляются, в остальных обрезаются неточно определённые нуклеотиды
- **2. Сборка генома (слияние ридов для коротких фрагментов) с помощью специальной программы – ассемблера.**
Результат: 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 (P), качество которого оценивается как Q равна:

- $P = 10^{-Q/10}$

- “Хорошее” качество при $Q \geq 30$ ($P \leq 0,001 = 0,1\%$)

Phred и ASCII:

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

ASCII control characters		
00	NULL	(Null character)
01	SOH	(Start of Header)
02	STX	(Start of Text)
03	ETX	(End of Text)
04	EOT	(End of Trans.)
05	ENQ	(Enquiry)
06	ACK	(Acknowledgement)
07	BEL	(Bell)
08	BS	(Backspace)
09	HT	(Horizontal Tab)
10	LF	(Line feed)
11	VT	(Vertical Tab)
12	FF	(Form feed)
13	CR	(Carriage return)
14	SO	(Shift Out)
15	SI	(Shift In)
16	DLE	(Data link escape)
17	DC1	(Device control 1)
18	DC2	(Device control 2)
19	DC3	(Device control 3)
20	DC4	(Device control 4)
21	NAK	(Negative acknowl.)
22	SYN	(Synchronous idle)
23	ETB	(End of trans. block)
24	CAN	(Cancel)
25	EM	(End of medium)
26	SUB	(Substitute)
27	ESC	(Escape)
28	FS	(File separator)
29	GS	(Group separator)
30	RS	(Record separator)
31	US	(Unit separator)
127	DEL	(Delete)

ASCII printable characters		
32	space	
33	!	
34	"	
35	#	
36	\$	
37	%	
38	&	
39	'	
40	(
41)	
42	*	
43	+	
44	,	
45	-	
46	.	
47	/	
48	0	
49	1	
50	2	
51	3	
52	4	
53	5	
54	6	
55	7	
56	8	
57	9	
58	:	
59	;	
60	<	
61	=	
62	>	
63	?	
64	@	
65	A	
66	B	
67	C	
68	D	
69	E	
70	F	
71	G	
72	H	
73	I	
74	J	
75	K	
76	L	
77	M	
78	N	
79	O	
80	P	
81	Q	
82	R	
83	S	
84	T	
85	U	
86	V	
87	W	
88	X	
89	Y	
90	Z	
91	[
92	\	
93]	
94	^	
95	_	
96	`	
97	a	
98	b	
99	c	
100	d	
101	e	
102	f	
103	g	
104	h	
105	i	
106	j	
107	k	
108	l	
109	m	
110	n	
111	o	
112	p	
113	q	
114	r	
115	s	
116	t	
117	u	
118	v	
119	w	
120	x	
121	y	
122	z	
123	{	
124		
125	}	
126	~	

Extended ASCII characters					
128	Ç	160	á	192	Ł
129	ü	161	í	193	ł
130	é	162	ó	194	Ť
131	â	163	ú	195	ť
132	ä	164	ñ	196	—
133	à	165	Ñ	197	+
134	â	166	ª	198	ã
135	ç	167	º	199	Ä
136	ê	168	¿	200	ℒ
137	ë	169	®	201	ℝ
138	è	170	¬	202	ℚ
139	ì	171	½	203	ℙ
140	ï	172	¼	204	℔
141	î	173	ı	205	=
142	Ā	174	«	206	≠
143	Ă	175	»	207	□
144	É	176	⋮	208	∅
145	æ	177	⋮	209	∅
146	Æ	178	⋮	210	Ê
147	ô	179		211	Ë
148	ö	180	†	212	Ě
149	ò	181	‡	213	ı
150	û	182	‡	214	ı̇
151	ù	183	‡	215	ı̈
152	ÿ	184	©	216	ı̇̇
153	Ō	185	≠	217	ı̈̇
154	Ū	186		218	ı̇̈
155	ø	187	¶	219	ı̈̈
156	£	188	¶	220	ı̇̈̇
157	∅	189	¢	221	ı̈̈̇
158	×	190	¥	222	ı̇̈̈
159	f	191	¬	223	ı̈̈̈
				224	Ó
				225	õ
				226	Ô
				227	Õ
				228	ö
				229	Õ
				230	µ
				231	þ
				232	ρ
				233	Ů
				234	Ű
				235	Û
				236	ý
				237	Ÿ
				238	ˉ
				239	˙
				240	≡
				241	±
				242	ˉ
				243	¼
				244	¶
				245	§
				246	÷
				247	˚
				248	˚
				249	˚
				250	˚
				251	˚
				252	˚
				253	˚
				254	■
				255	nbsp

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 ридсах: примеси

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



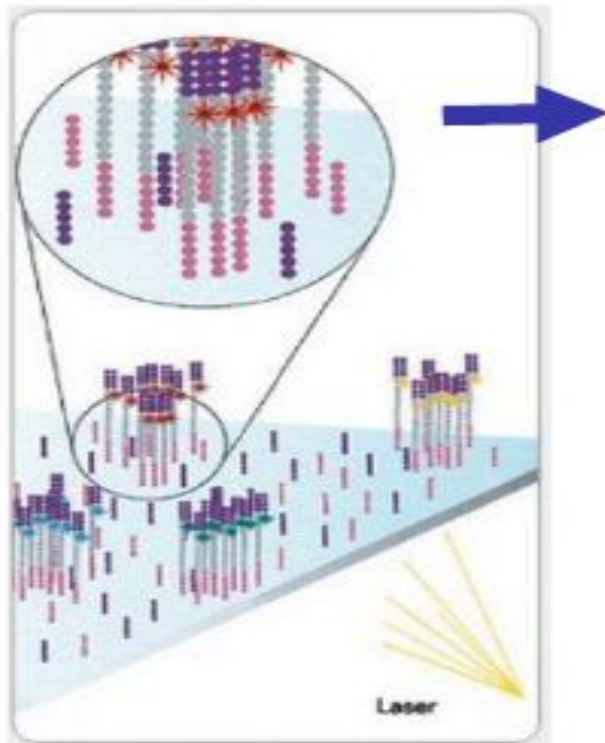
- ЧТЕНИЕ СКВОЗЬ – вставки слишком короткие



- 2. Биологические – контаминация

Источники ошибок в рядах: фазировка

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



CGCAAGGAAGCTGATTG
CGCAAGGAAGCTGATTG
CGCAAGGAAGCTGATTG
CGCAAGGAAGCTGATT
CGCAAGGAAGCTGATTGC
CGCAAGGAAGCTGATTG

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

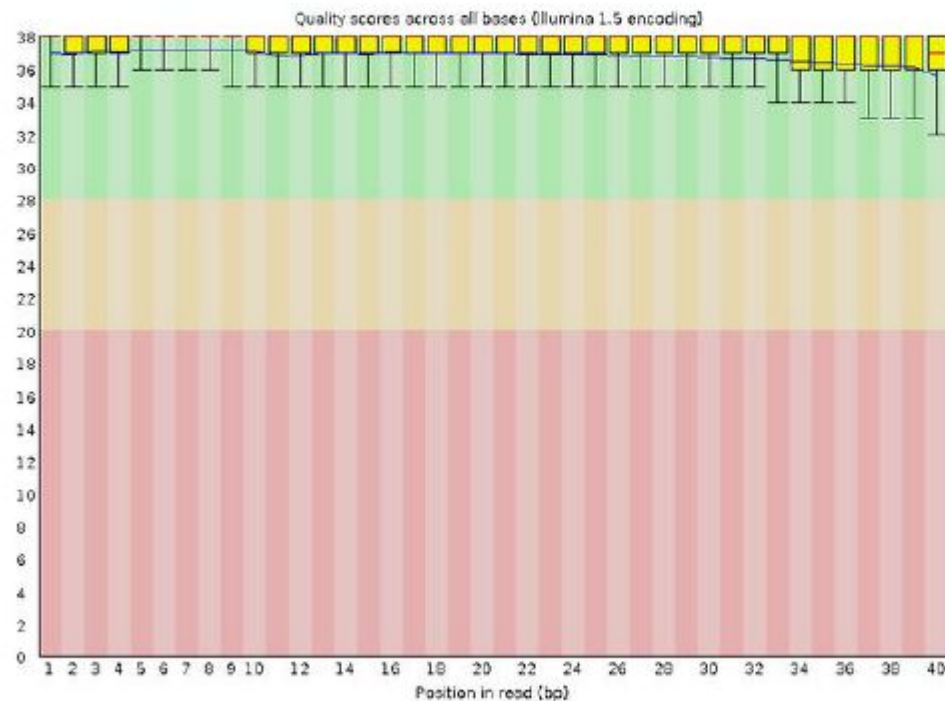
Программа FastQC – контроль качества ридов:

1. Среднее нуклеотидное качество – хорошее (все Me>25, все Q1>10)

Нуклеотидное качество



Per base sequence quality



Качество высокое
на протяжении
всего чтения

- Синяя линия – среднее качество
- Красная линия – медиана
- Жёлтая рамка – интерквартиль (50% чтений попадает в эти границы)
- Чёрные засечки – 80% чтений попадает в эти границы

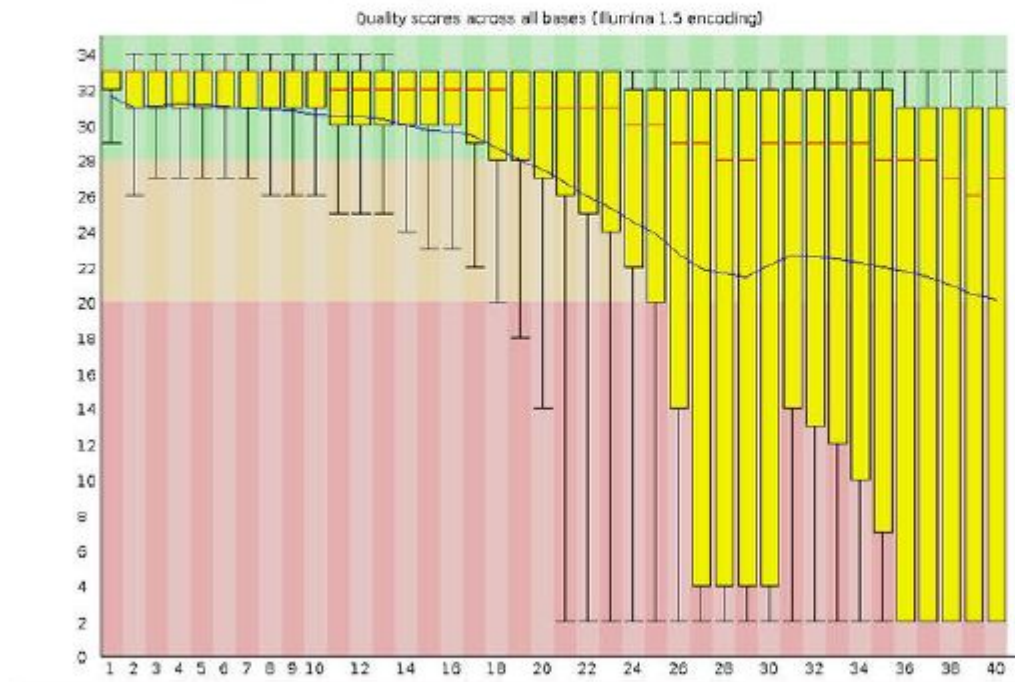
Программа FastQC – контроль качества ридов:

1. Среднее нуклеотидное качество – неудовлетворительное (есть $Me < 20$ или $Q1 < 5$)

Нуклеотидное качество



Per base sequence quality



Качество сильно падает к концу

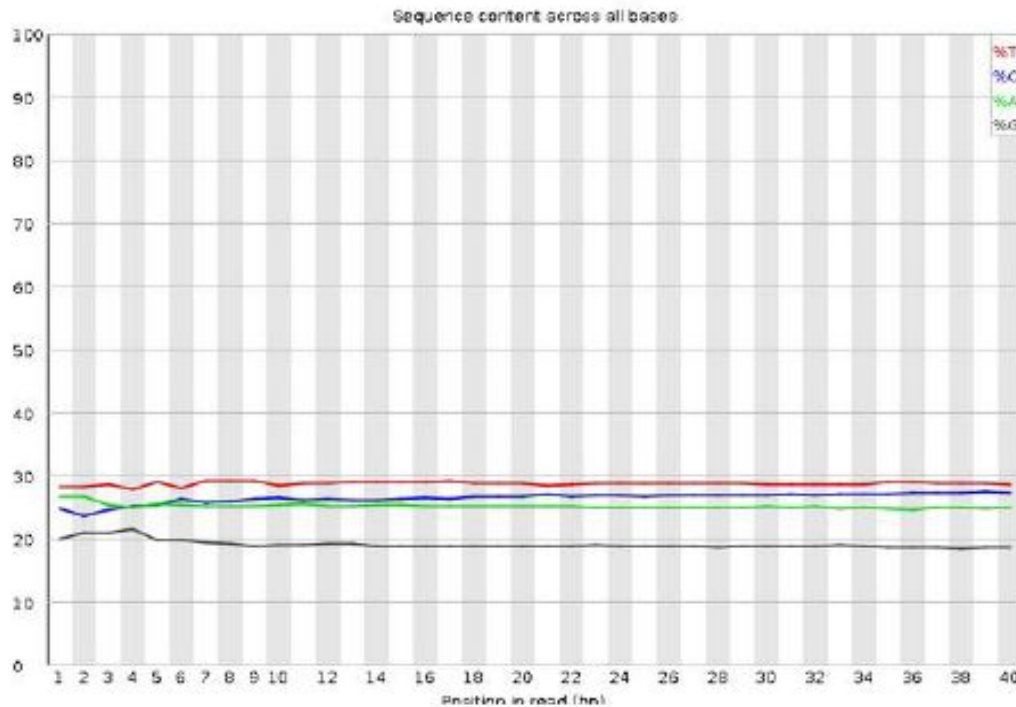
- Синяя линия – среднее качество
- Красная линия – медиана
- Жёлтая рамка – интерквартиль (50% чтений попадает в эти границы)
- Чёрные засечки – 80% чтений попадает в эти границы

Программа FastQC – контроль качества ридов:

2. Средний нуклеотидный состав ридов

Нуклеотидный состав

Per base sequence content



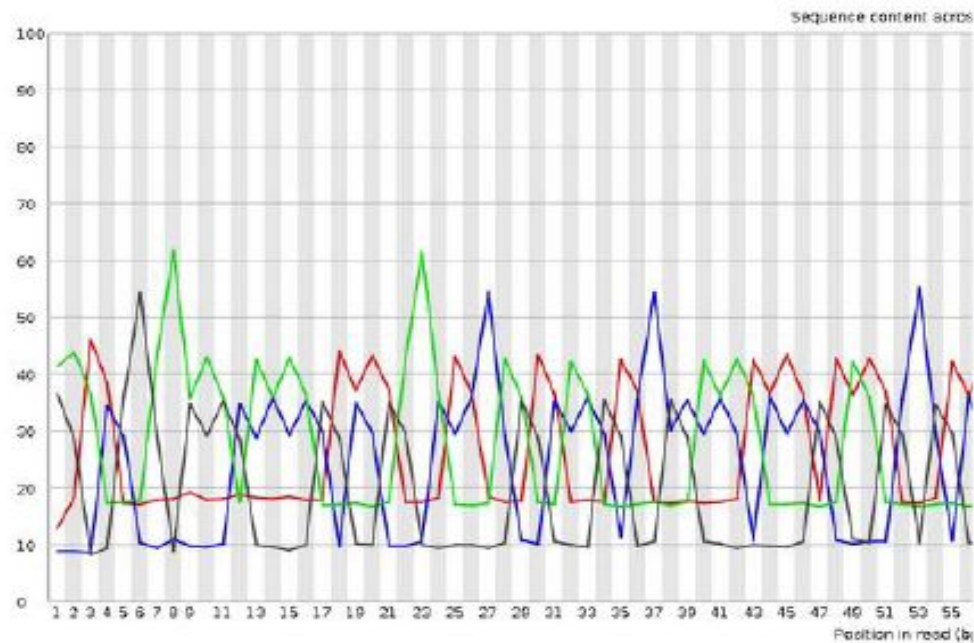
- Нуклеотидный состав примерно одинаковый на протяжении всего рида
- **В начале чтения небольшие колебания** из-за специфичности фрагментации – программа обозначает их восклицательным знаком, но это не проблема

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

Нуклеотидный состав



❌ Per base sequence content



Сильные колебания нуклеотидного состава – скорее всего отсеквенированы димеры адаптеров

Программа FastQC – контроль качества ридов:

3. Чрезмерно представленные последовательности

- 1)  **Overrepresented sequences**
No overrepresented sequences

- 2)  **Overrepresented sequences**

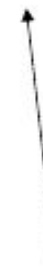
Sequence	Count	Percentage	Possible Source
<code>GATCGGAAAGAGCACACAGCTCTGGAACTCCAGTCAAGSTCCGGCACATCTCGTAT</code>	938427	17.993033077161126	TruSeq Adapter, Index 1 (97% over 36bp)



Последовательность адаптера, из-за которого было колебание нуклеотидного состава на прошлом слайде



Адаптер найден в ~18 % чтений



Название найденного адаптера – пригодится, когда настанет пора его удалять

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 (assembly)

- de novo (сборка не секвенированного ранее генома)
- – метод OLC (overlap layout consensus) (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

суффиксное дерево

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 ($\geq 95\%$ сходства)



Сборка de novo: Overlap layout consensus: 2

- На базе попарного выравнивания строят множественное выравнивание, корректируют ошибки

TAGATTACACAGATTACTGA TTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGCGTAAACTA
TAG TTACACAGATTA**T**TGACTT**C**ATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGCGTAA CTA
TAGATTACACAGATTACTGACTTTGATGGGGTAA CTA



TAGATTACACAGATTACTGACTTTGATGGCGTAA CTA

Сборка de novo: Графы де Брёйна

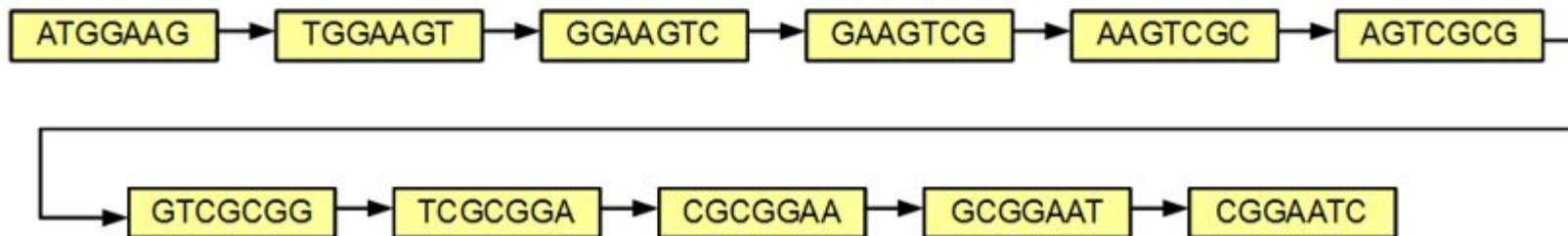
sequence

ATGGAAGTCGCGGAATC

7mers

ATGGAAG
TGGAAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

de Bruijn graph



Результат сборки: контиги и скэффолды

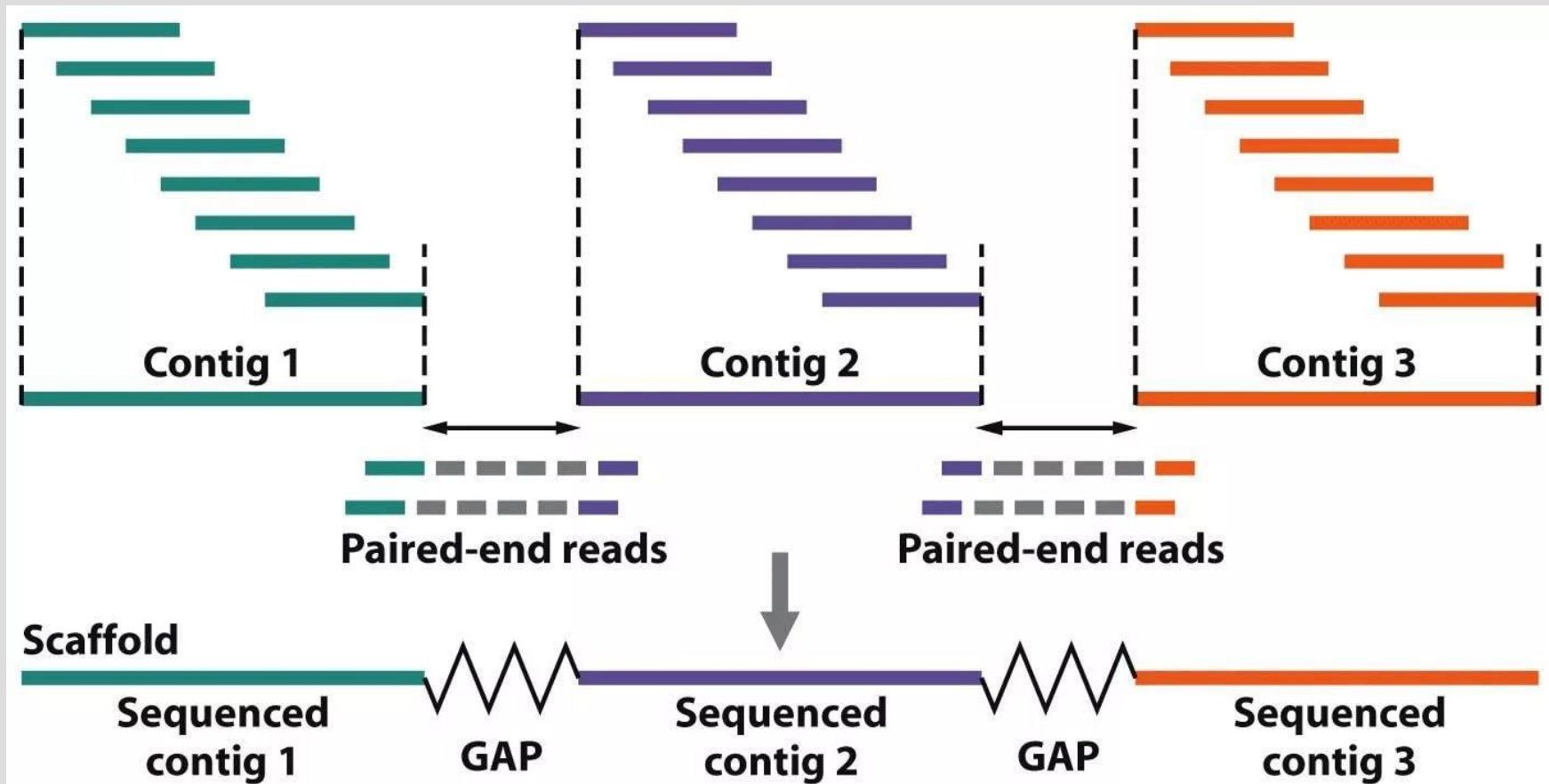


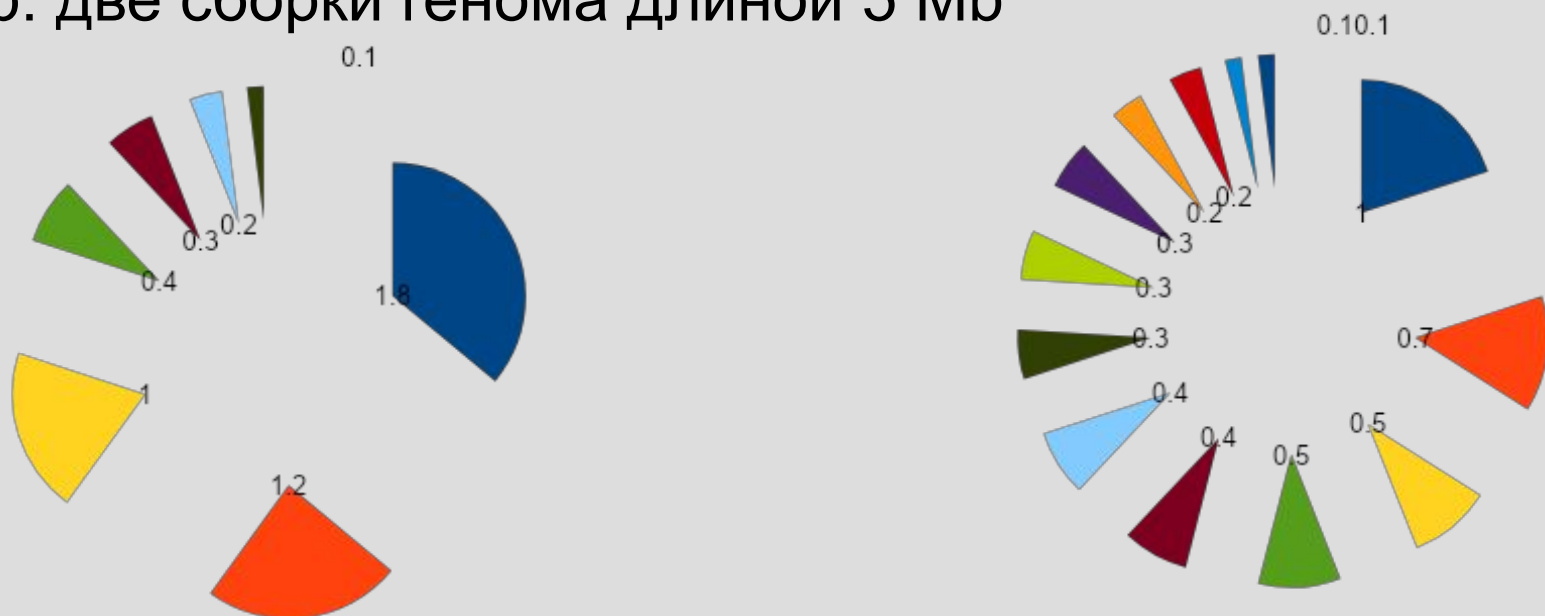
Figure 14-7

Introduction to Genetic Analysis, Tenth Edition

© 2012 W. H. Freeman and Company

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 Mb



Формат представления нуклеотидных последовательностей – FASTA

• >OTU-160-1 *Acinetobacter baumannii*

• CCTACGGGGGGCTGCAGTGGGGGAATATTGGACAATGGGGGGGA
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 (ATG.....)
- 2. Поиск других кодирующих последовательностей (гены РНК)

Результаты аннотации

- 1. Структурное описание:
 - – открытые рамки считывания (ORF) и их 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