

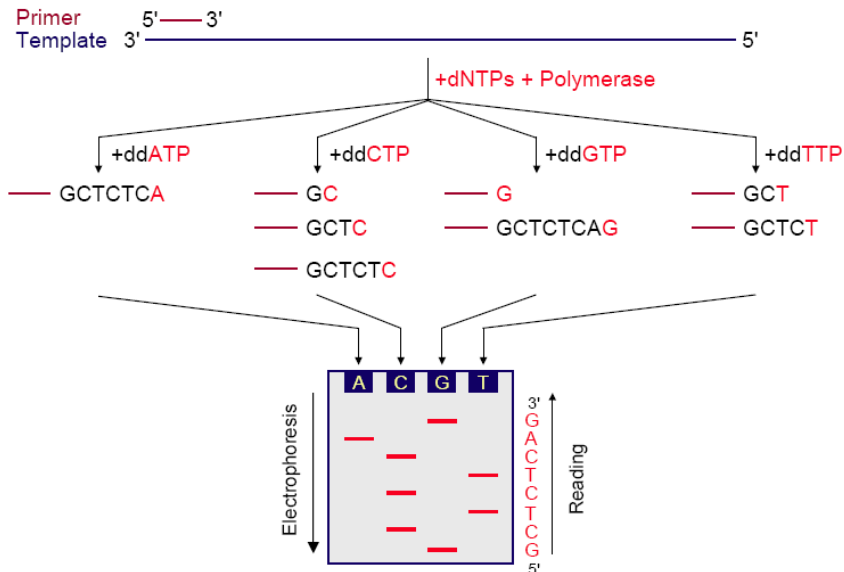
МЕТОДЫ СЕКВЕНИРОВАНИЯ НОВОГО ПОКОЛЕНИЯ

Дмитрий Первушин

19 апреля 2017г.

Сколковский Институт Науки и Технологии
Московский государственный университет им. МВ Ломоносова
Государственный Университет - Высшая Школа Экономики
Centre Regulació Genòmica, Parc Reserca Biomèdica Barcelona
Universitat Pompeu Fabra

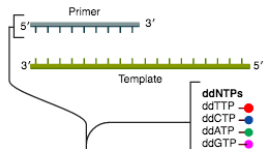
Метод дидезоксинуклеотидов (Нобелевская премия 1980)



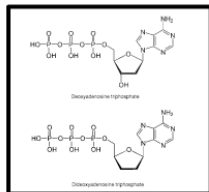
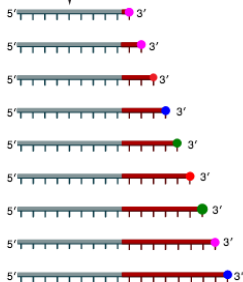
Секвенирование по Сэнгеру

① Reaction mixture

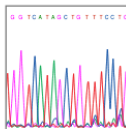
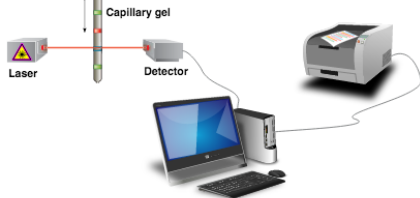
- Primer and DNA template
- DNA polymerase
- ddNTPs with flouochromes
- dNTPs (dATP, dCTP, dGTP, and dTTP)



② Primer elongation and chain termination



③ Capillary gel electrophoresis separation of DNA fragments



Chromatograph

④ Laser detection of flouochromes and computational sequence analysis



APPLICATIONS OF NEXT-GENERATION SEQUENCING

Coming of age: ten years of next-generation sequencing technologies

Sara Goodwin¹, John D. McPherson² and W. Richard McCombie¹

Abstract | Since the completion of the human genome project in 2003, extraordinary progress has been made in genome sequencing technologies, which has led to a decreased cost per megabase and an increase in the number and diversity of sequenced genomes. An astonishing complexity of

Массовое секвенирование коротких фрагментов

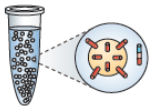
Секвенирование коротких ридов

- Амплификация шаблона : цель – получить много тысяч идентичных молекул ДНК в одной микрообласти пространства для того, чтобы можно было регистрировать сигнал и отличать его от шума
 - Эмульсионная ПЦР
 - Твердофазная генерация кластеров
 - Генерация кластеров в растворе
- Секвенирование : локальное секвенирование амплифицированных шаблонов
 - Секвенирование через лигирование (SBL)
 - Секвенирование через синтез (SBS)
 - Циклическая обратимая терминация (CRT)
 - Однонуклеотидное приращение (SNA)

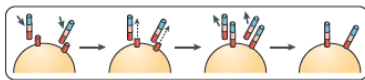
Эмульсионная ПЦР

- К ДНК пришивается адаптер
- К наночастицам пришиваются олигонуклеотиды, комплементарный адаптеру
- С каждой наночастицей связывается только одна молекула ДНК
- Раствор, содержащий реагенты ПЦР, превращается в обращенную эмульсию так, что каждая наночастица попадает в одну мицеллу
- В процессе ПЦР на наночастице образуется колония идентичных молекул
- Эмульсия обращается, частицы осаждаются

■ Emulsion PCR
(454 (Roche), SOLiD (Thermo Fisher), GeneReader (Qiagen), Ion Torrent (Thermo Fisher))



Emulsion
Micelle droplets are loaded with primer, template, dNTPs and polymerase



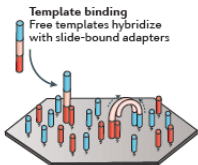
On-bead amplification
Templates hybridize to bead-bound primers and are amplified; after amplification, the complement strand disassociates, leaving bead-bound ssDNA templates



Final product
100–200 million beads with thousands of bound template

Твердофазная генерация кластеров

b Solid-phase bridge amplification (Illumina)



Bridge amplification

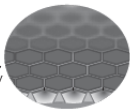
Distal ends of hybridized templates interact with nearby primers where amplification can take place

Cluster generation

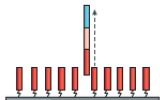
After several rounds of amplification, 100–200 million clonal clusters are formed

Patterned flow cell

Microwells on flow cell direct cluster generation, increasing cluster density

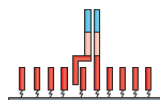


c Solid-phase template walking (SOLiD Wildfire (Thermo Fisher))



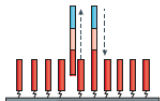
Template binding

Free DNA templates hybridize to bound primers and the second strand is amplified



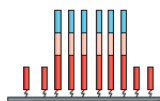
Primer walking

dsDNA is partially denatured, allowing the free end to hybridize to a nearby primer



Template regeneration

Bound template is amplified to regenerate free DNA templates



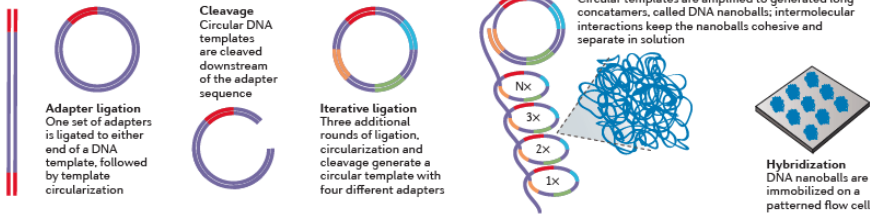
Cluster generation

After several cycles of amplification, clusters on a patterned flow cell are generated

ДНК наношары (Пекинский Институт Биоинформатики)

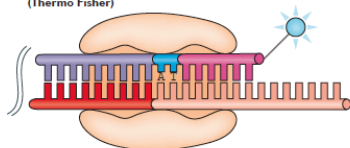
- К ДНК пришивается адаптер и молекула закольцовывается
- Затем ДНК разрезается в другом месте и процесс повторяется несколько раз
- В процессе амплификации катящегося кольца образуются длинные конкатамеры

d In-solution DNA nanoball generation (Complete Genomics (BG1))



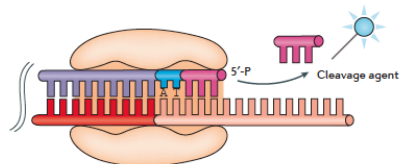
Секвенирование через лигирование (SOLiD)

■ SOLiD (Thermo Fisher)



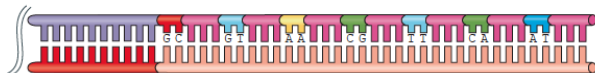
Two-base-encoded probes

Probes with two known bases followed by degenerate or universal bases hybridize to a template; ligase immobilizes the complex and the slide is imaged



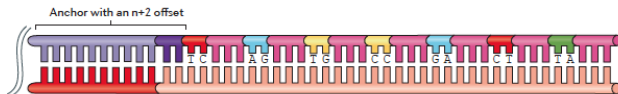
Cleavage

The fluorophore is cleaved from the probe along with several bases, revealing a 5' phosphate



Probe extension

10 rounds of hybridization, ligation, imaging and cleavage identify 2 out of every 5 bases

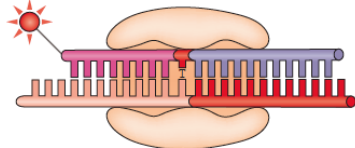


Reset

After a round of probe extension, all probes and anchors are removed and the cycle begins again with an offset anchor

Секвенирование через лигирование (Complete Genomics)

b Complete Genomics (BGI)



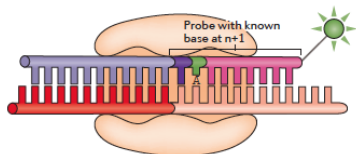
Single-base-encoded probes

A probe with a single known base and degenerate bases hybridizes to a template and is imaged



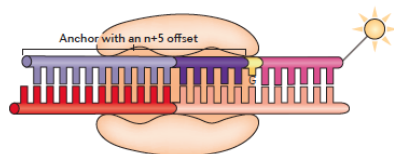
Reset

After each imaging step, both the probe and anchor are removed



Paired-end sequencing

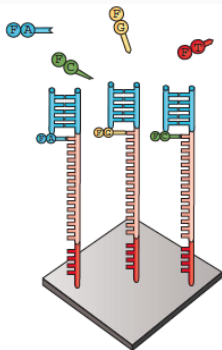
Sequencing is performed for both the left and right sides of the adapter



Offset anchors

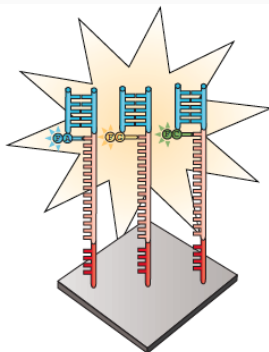
Subsequent rounds of hybridization and ligation use offset anchors to sequence more-distant bases

Секвенирование через синтез CRT (Illumina)



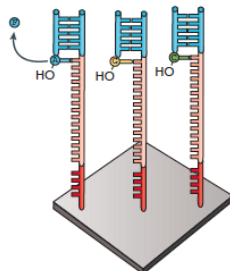
Nucleotide addition

Fluorophore-labelled, terminally blocked nucleotides hybridize to complementary base. Each cluster on a slide can incorporate a different base.



Imaging

Slides are imaged with either two or four laser channels. Each cluster emits a colour corresponding to the base incorporated during this cycle.

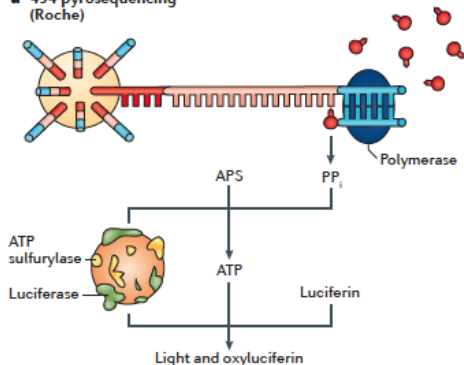


Cleavage

Fluorophores are cleaved and washed from flow cells and the 3'-OH group is regenerated. A new cycle begins with the addition of new nucleotides.

Секвенирование через синтез SNA (Roche 454)

454 pyrosequencing (Roche)

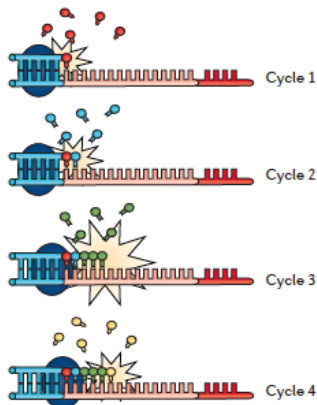


Pyrosequencing

As a base is incorporated, the release of an inorganic pyrophosphate triggers an enzyme cascade, resulting in light

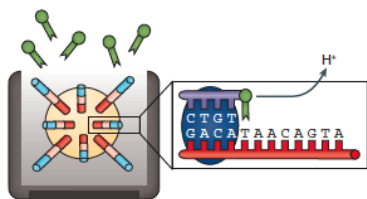
Single nucleotide addition

Only one dNTP species is present during each cycle; multiple identical dNTPs can be incorporated during a cycle, increasing emitted light

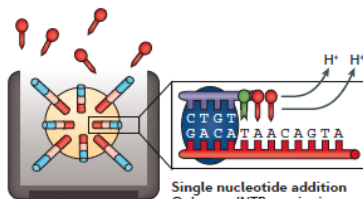
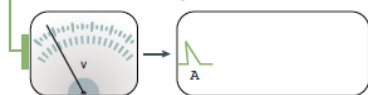


Секвенирование через синтез SNA (Ion Torrent)

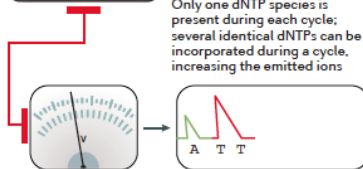
b Ion Torrent (Thermo Fisher)



Semiconductor sequencing
As a base is incorporated, a single H^+ ion is released, which is detected by a CMOS- $ISFET$ sensor



Single nucleotide addition
Only one dNTP species is present during each cycle; several identical dNTPs can be incorporated during a cycle, increasing the emitted ions



- комплементарный металло-оксидный полупроводник (CMOS)
- ион-чувствительный полевой транзистор (ISFET)

Секвенирование одиночных длинных ДНК

Секвенирование одиночных длинных ДНК (PacBio)

Aa Pacific Biosciences

SMRTbell template

Two hairpin adapters allow continuous circular sequencing



ZMW wells

Sites where sequencing takes place

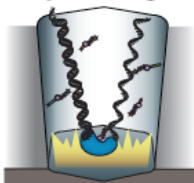


Labelled nucleotides

All four dNTPs are labelled and available for incorporation

Modified polymerase

As a nucleotide is incorporated by the polymerase, a camera records the emitted light



PacBio output

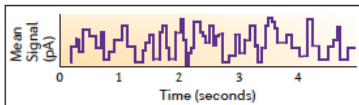
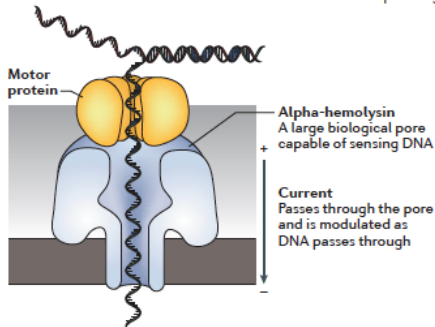
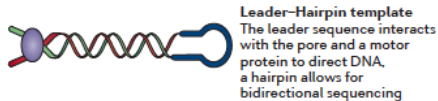
A camera records the changing colours from all ZMWs; each colour change corresponds to one base



- Длинная ДНК лигируется к закольцованным адаптерам
- Модифицированная полимераза прикрепляется к дну микро-колодца ZMT (zero-mode waveguide)
- Полимераза отрезает флуорофор после включения нуклеотида
- Каждая молекула прочитывается многократно благодаря кольцевой структуре

Секвенирование одиночных длинных ДНК (Oxford Nanopore)

Ab Oxford Nanopore Technologies



ONT output (squiggles)
Each current shift as DNA translocates through the pore corresponds to a particular k -mer

- Изменение в проводимости ионного канала зависят от того, какая последовательность находится в поре, что можно интерпретировать как k -мер

Синтетические длинные ряды (Illumina Moleculo)

Ba Illumina

DNA fragment

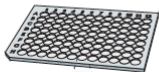
DNA is fragmented and selected to ~10 kb



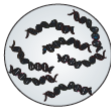
~3,000 molecules per well

Enzymatic cleavage

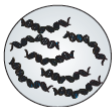
DNA is barcoded and fragmented to ~350 bp



A1



A2



Barcodes

DNA from the same well shares the same barcode

Pooling

DNA from each well is pooled and undergoes a standard library preparation



Sequencing

DNA is sequenced on a standard short-read sequencer

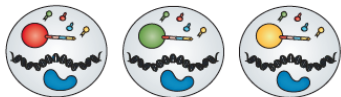
- ДНК фрагментируется до 8–10 килобаз и переносится на подложку с микроколдцами
- Каждый фрагмент разрезается до ~350 нт и к нему пришивается штрихкод, в каждом колдце свой
- ДНК с штрихкодами объединяются и секвенируются по обычному протоколу

Синтетические длинные ряды 10X Genomics

Bb 10X Genomics

Emulsion PCR

Arbitrarily long DNA is mixed with beads loaded with barcoded primers, enzyme and dNTPs



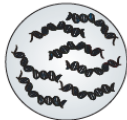
GEMs

Each micelle has 1 barcode out of 750,000



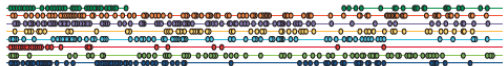
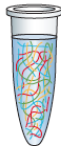
Amplification

Long fragments are amplified such that the product is a barcoded fragment ~350 bp



Pooling

The emulsion is broken and DNA is pooled, then it undergoes a standard library preparation



- Фрагменты ДНК ~100 килобаз помещаются в мицеллы, каждая мицелла содержит адаптер и уникальный штрихкод
- В каждой мицелле короткие фрагменты ДНК амплифицируются и помечаются штрихкодами, каждый из которых соответствует длинному фрагменту

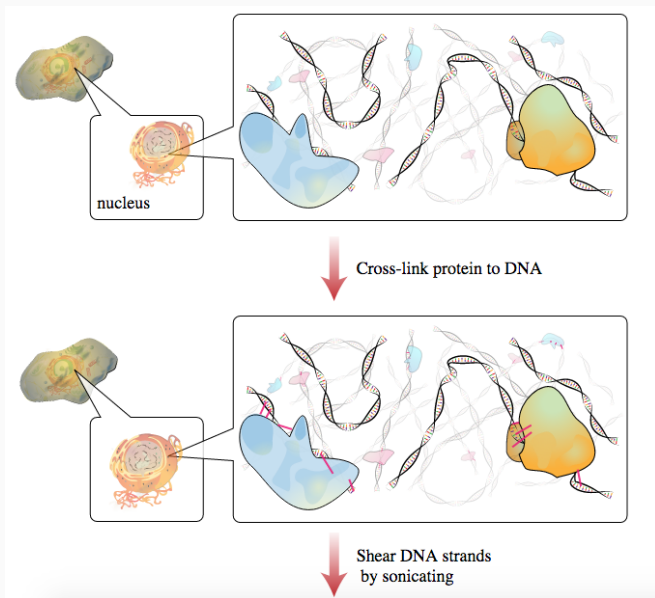
Приложения

Pull-down эксперименты : иммунопреципитация

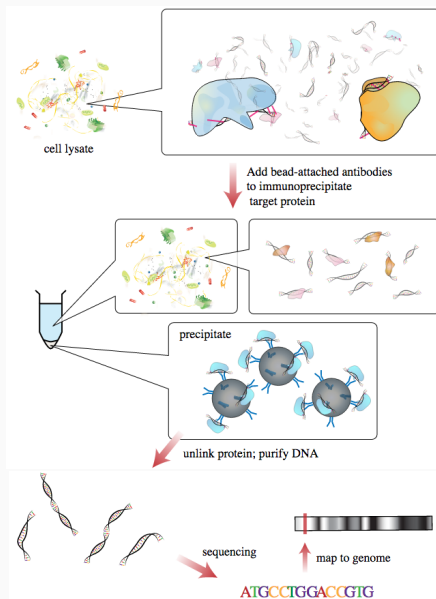
Иммунопреципитация (IP) – техника осаждения антигена путем специфического связывания с антителом

- 1 Иммунопреципитация отдельных белков (IP)
- 2 Иммунопреципитация белковых комплексов (Co-IP)
- 3 Иммунопреципитация хроматина (ChIP)
- 4 Иммунопреципитация рибопротеиновых комплексов (RIP, eCLIP)
- 5 Рибосомный профайлинг
- 6 Бромоеуридиновый профайлинг (Bru-seq)
- 7 Бромоеуридиновый пульсовый профайлинг

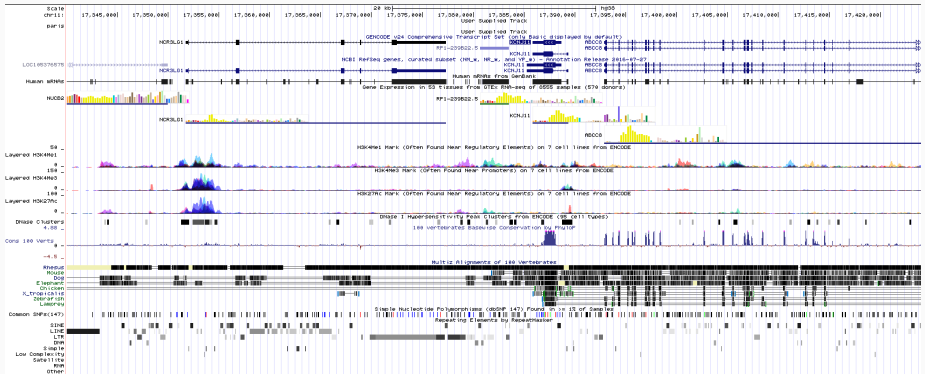
Иммунопреципитация (ChIP)



Иммунопреципитация (ChIP)



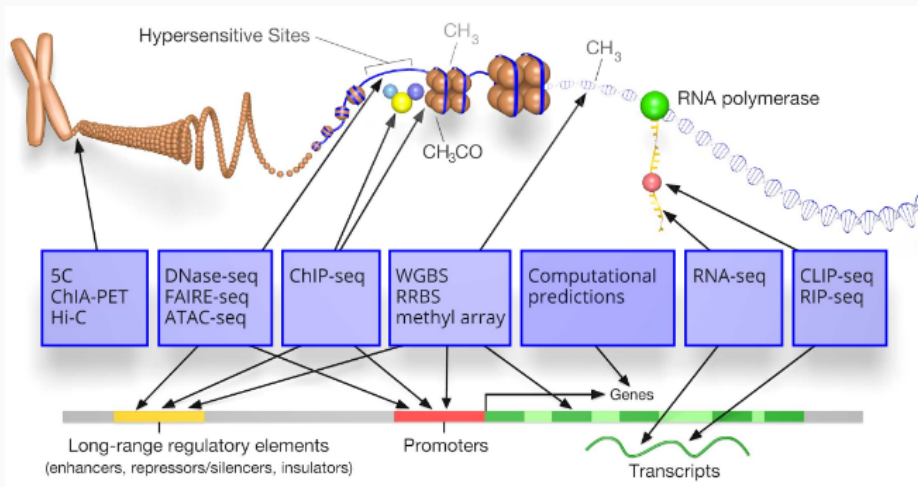
Пример



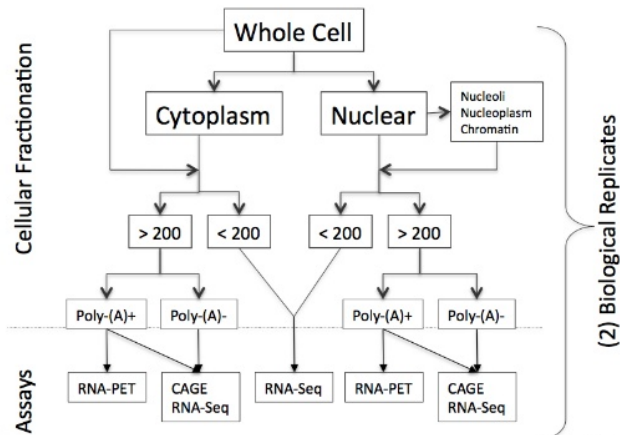
Further applications

- Эксперименты обогащения/деплеции
 - РНКазы R для определения кольцевых РНК
 - Полиаденилированные РНК
 - Деpletion рРНК
- Определение доступности
 - Сайты гиперчувствительности ДНКазы
 - ATAC-seq
- Определение пространственной близости
 - Определение конформации хромосомы (3C)
 - Определение РНК дуплексов (PARIS)

ENCODE consortium (encodeproject.org)



The ENCODE RNA assays (CSHL, CalTech)



23/07/

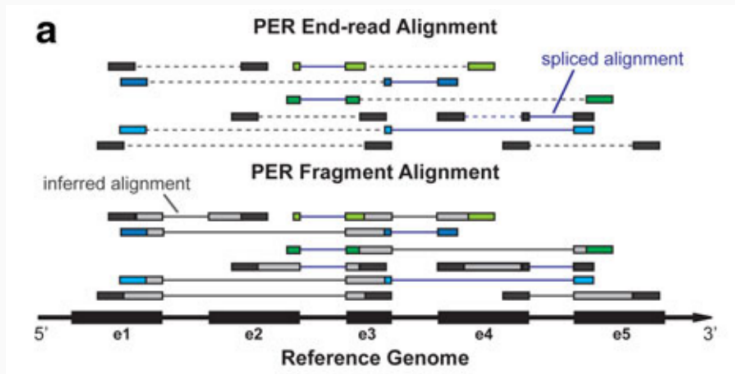
Carrie Davis, CSHL²⁵

Парные прочтения и цепь-специфические прочтения



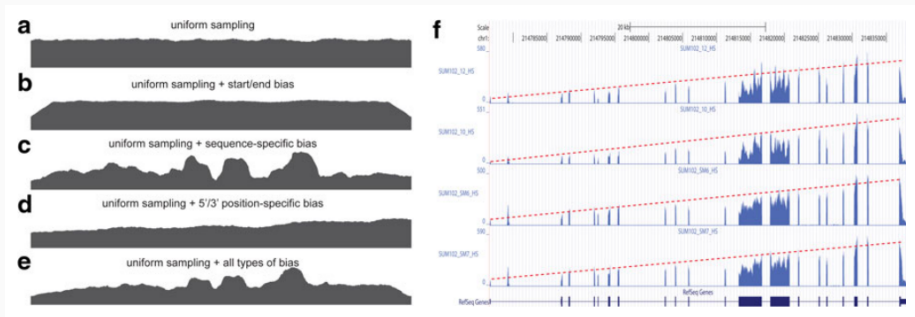
- Если секвенируем одноцепочечную РНК
 - Сначала синтезируем комплементарную цепь (RT)
 - Плюс-цепь синтезируем, но добавляем модифицированный нуклеотид (тиоурин)
 - Двухцепочечный дуплекс лигируем с адаптером
 - Расщепляем тиоурин

- Картирование : имеется большое число коротких прочтений (10^8) и референсный геном
 - Высокая скорость
 - Небольшое число ошибок
 - Неоднозначные картирования
 - Разрывные/нелинейные картирования
- Квантификация : количественная оценка представленности каждого геномного участка
- Peak Calling : локализованный количественный сигнал
- Деконволюция : нужно восстановить представленности и структуру перекрывающихся транскриптов, основываясь на локальной информации о коротких ридах



- Для данной модели транскрипта, какова вероятность получить наблюдаемое распределение ридов?
- Имея наблюдаемое распределение ридов, какая структура транскриптов максимизирует постериорную вероятность видеть то, что мы видим?
- Какая структура транскриптов максимизирует функцию потока в графе?

Неравномерное распределение ридов*

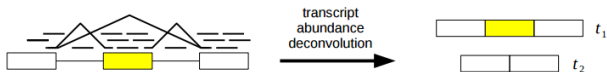


- 3'-смещение из-за прайминга обратной транскрипции
- ГЦ смещение – риды с высоким ГЦ-составом оказываются перепредставлены
- Эффект лаборатории
- Артефакты штрихкодов в новых технологиях Илюмины

* Huang et al, JCB 2013

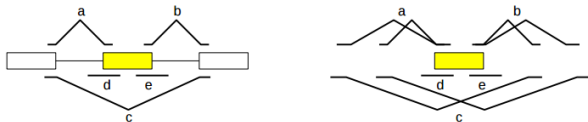
Количественные оценки сплайсинга

1. Transcript-centric



$$psi = \frac{t_1}{t_1 + t_2}$$

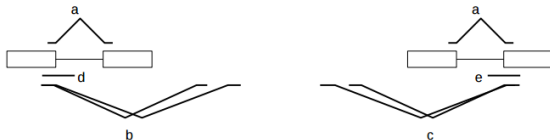
2a. Exon-centric



$$psi = \frac{a+b}{a+b+2c}$$

$$cosi = \frac{a+b+2c}{a+b+2c+d+e}$$

2b. Intron-centric (can be generalized to arbitrary splicing graphs)



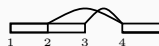
$$psi5 = \frac{a}{a+b}$$

$$psi3 = \frac{a}{a+c}$$

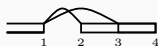
$$cosi5 = \frac{a+b}{a+b+d}$$

$$cosi3 = \frac{a+c}{a+c+e}$$

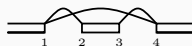
Граф сплайсинга



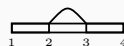
(a)



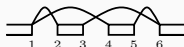
(b)



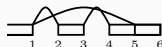
(c)



(d)



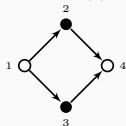
(e)



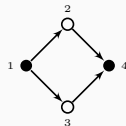
(f)



(g)



(a)



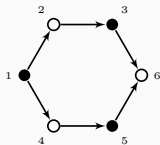
(b)



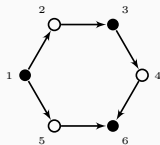
(c)



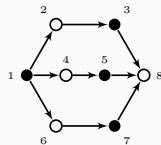
(d)



(e)

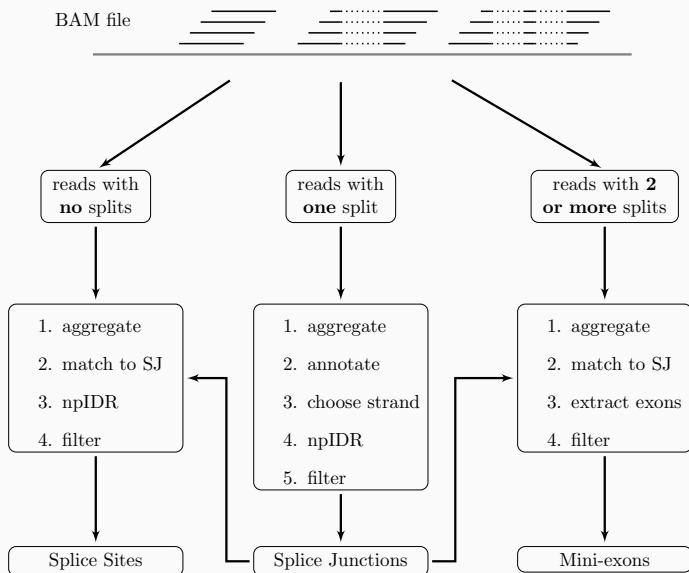


(f)



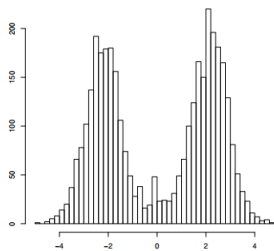
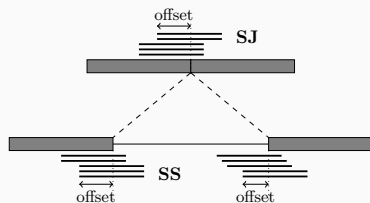
(g)

Ψ is a binomial (multinomial) variable defined for a group of vertex-independent paths

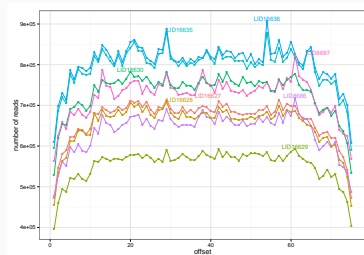


Offset

Offset is the position of the split within the read sequence



$\log_{10}(c_+) - \log_{10}(c_-)$

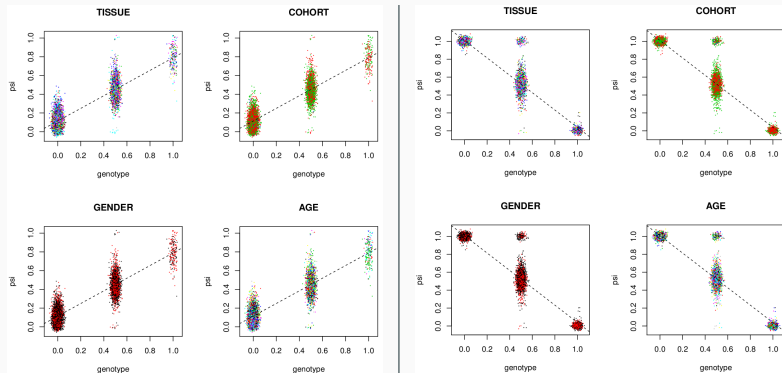


GTEx: splicing QTLs

$$\Psi = b_0 + b_1 * genotype, \quad 0 \leq \Psi \leq 1, \quad 0 \leq genotype \leq 1$$

$$\Delta\Psi = 0.1 \Leftrightarrow b_1 = 0.1$$

$$t = (|b_1| - 0.1) / SE(b_1)$$



Challenge: find sQTLs in unbalanced design

Спасибо за внимание