Where are we now?





ATCG.....

Exhausting experiment

Sequencing

Analysis our data

Structure of prepared library



Length of this part is sometimes called "insert size", which totally depends on the length of DNA after shearing















What comes from the sequencer

File suffix is "fq" or "fastq" Ai1_R1 fq.gz **Reads** 8 Ai1_R2.fq.gz AI1_R1 fq **Bases** Ai1.base.pdf 20 Al1_R2.fq composition Ai1.base.png 20 Ai1.quality.pdf **Quality of Bases** Ai1.quality.png 2

Assigned (more professionally called "demultiplexed ") reads according to index close to P7

Expanded file

How reads in xxx_R1.fq looks like



How reads in xxx_R2.fq looks like



What's paired reads



Reads originated from the same molecular before bridge amplification called paired reads

Paired reads in xxx_R1.fq and xxx_R2.fq

xxx_R1.fq

xxx_R2.fq

Paired reads have unique id in their corresponding file and exactly the same name

Compare reads in xxx_R1.fq and xxx_R2.fq

xxx_R1.fq

1	@E00492:247:HFMH3CCXY:8:1101:18436:2909 1:N:0:CTGCAATG
2	CTCAAAAACAACAATAATCTATTATCACAACCAAAGCAACAGCATACATA
3	+
4	FJJJJAJJJJJFJJ<-FJJJFJJJJJJJJJJJJJJJJJJJ
5	@E00492:247:HFMH3CCXY:8:1101:6725:3401 1:N:0:CTGCAATG
6	AATTCAGCGACTTGGTTCACACAATCCTCCATCCCTCCTTTCCTTTACCACGATGTTTGATGATGCTGCAGGTACAGCTGCTGAAGAGACACGAGCTGTGACCAGTAACTA
7	+
8	7JJJAJJJ-FJJJJJJFJJJJAJFJFJA-FFJJFJJJJAAJJF7AJJJJJJJF7A7JJJJJFJJJJAFJJAJFFFJJJ7JFFJJFJJJJJF<7777
9	@E00492:247:HFMH3CCXY:8:1101:8745:3841 1:N:0:CTGCAATG
10	AAATGTAAGGTCAAACTACAACCTTGTCTATGAAGCCGGTGTGCCATTTGCGTAGCTTCCTGTTTTCGGTGGAGAGTTTCTCAGCAGCAGAGCTGCAGCTTGGCATTGTCA
11	+
12	JJJJJJJJJJJJJJJJFJJJFJJJJFJJJJJJJJJJJJ

xxx_R2.fq

1	@E00492:247:HFMH3CCXY:8:1101:18436:2909 2:N:0:CTGCAATG
2	ATAATGGTTTTTGACAAGTTTAGCTTAAAGGAGGAAAATTGCATGTGTTGTGATCTTGTCACTCTCTGTGTCTCCTAGACTCCTGCAGAATTGCAAGGAGGACAGGTTAGA
3	+
4	AJJ7F<- <f<fffj7fj-77aj-af-77-7f<-<<f-7-af7<7aaa7-77aj-a-aa-a-ajf7a<mark>7AAF7J<a7fjaj7a-f< mark=""></a7fjaj7a-f<></f<fffj7fj-77aj-af-77-7f<-<<f-7-af7<7aaa7-77aj-a-aa-a-ajf7a<mark>
5	@E00492:247:HFMH3CCXY:8:1101:6725:3401 2:N:0:CTGCAATG
6	GGTATCTTTAAAACACATTTAACATGCGATACATAACTTATAATTGGTCATGTTTGTCATACAAATTGAATTAGACAGTGAATCAAAGAAAATGTGCTTACACTGACAAGT
7	+
8	JAA7A <fjj<jfj7ff<jfjj<-f-<jj-fjf77f<fjjjfjja-fjjjfj<<jj7jjfff<f-ffa-a7<-f-afajja<ffjfjffff-jjfjfa777-< th=""></fjj<jfj7ff<jfjj<-f-<jj-fjf77f<fjjjfjja-fjjjfj<<jj7jjfff<f-ffa-a7<-f-afajja<ffjfjffff-jjfjfa777-<>
9	@E00492:247:HFMH3CCXY:8:1101:8745:3841 2:N:0:CTGCAATG
10	AATGGTTGCCAAGCAACACAGAGAGAGGAAGTGGCAAGTATGATTGAAGATTGGAGTGATATGAGTGATATTAACATTCAGCTGGGGGGTGATTATTAACATTTAGCAGGATAT
11	+
12	FAF<-7FA<7-<7<<7FFJ7F< <a<aff-7f-a<ff-<a<-ff7ff7ajf77<jf-<af-j-f-afafjjjaf-7<aaaaajafj<fffjfjjjjf<j<<a<<<-a-7< th=""></a<aff-7f-a<ff-<a<-ff7ff7ajf77<jf-<af-j-f-afafjjjaf-7<aaaaajafj<fffjfjjjjf<j<<a<<<-a-7<>

Paired reads are placed at the same line of its corresponding file

Let's start analysis now!!

What's the goal of analysis?

The sequence of loci used to design the baits, we call it "reference"

locus 1

locus 2

locus 3

Intact sequences of corresponding locus of each sample



sample3

flanked by intron or UTR region

What's the goal of analysis?

1	@E00492:247:HFMH3CCXY:8:1101:18436:2909 1:N:0:CTGCAATG
2	CTCAAAAACAACAATAATCTATTATCACAACCAAAGCAACAGCATACATA
3	+
4	FJJJJAJJJJJFJJ<-FJJJFJJJJJJJJJJJJJJJJJJJ
5	@E00492:247:HFMH3CCXY:8:1101:6725:3401 1:N:0:CTGCAATG
6	AATTCAGCGACTTGGTTCACACAATCCTCCATCCCTCCTTCCT
7	+
8	7JJJAJJJ-FJJJJJJFJJJJAJFJFJA-FFJJFJJJAAJJF7AJJJJJJJF7A7JJJJFJJJJAFJJAJFFFJJJ7JFFJJFJJJJJF<7777
9	@E00492:247:HFMH3CCXY:8:1101:8745:3841 1:N:0:CTGCAATG
10	AAATGTAAGGTCAAACTACAACCTTGTCTATGAAGCCGGTGTGCCATTTGCGTAGCTTCCTGTTTTCGGTGGAGAGTTTCTCAGCAGCAGACTGCAGCTTGGCATTGTCA
11	+
12	JJJJJJJJJJJJJJJJFJJJFJJJJFJJJJJJJJJJJJ

Short raw data

There's only one sequence for each sample





3 steps to recover qualified assemblies from raw data

Data preparation

Assembling

Further processing

Data preparation

Demultiplex reads according to inline index

Trim low quality bases and adaptor

Demultiplex reads according to inline index



Reads we got still includes inline index

We need to demultiplex reads according to them, then cut them out

assign reads to its sample

How to demultiplex

Reads before demultiplexing

Paired Reads



Trim low quality bases and adaptor

Why we need to trim low quality bases

A cluster of reads amplified from the same molecular



Trim low quality bases and adaptor

Why we need to trim low quality bases

A cluster of reads amplified from the same molecular

longer the time of

of bases becomes

lower



How to trim low quality bases

E(quality)>=15



ACGGCGTAGGCTGATGATCG

Why we need to trim adaptor?



Why we need to trim adaptor?



How to trim adaptor?



sequence of adaptor and inline index

Reads have been cleaned. Let's start assemble !

What's assemble

Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

Reads 2: TCTGATTCGGTTCGGATCTGGGCAT

Reads 3: ATCTGGGCATGGCGTTCGATGTCGCTAT

3 reads in a sample

What's assemble

Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

Reads 2: TCTGATTCGGTTCGGATCTGGGCAT

Reads 3: ATCTGGGCATGGCGTTCGATGTCGCTAT

Resulting contig:

Contig1 CGGCGGATCTGATGGGATCTGATTCGGTTCGGATCTGGGCAT

"Contig" is the sequence assembled from the reads

What's assemble

Contig1 CGGCGGATCTGATGGGATCTGATTCGGTTCGGATCTGGGCAT

Reads 3: ATCTGGGCATGGCGTTCGATGTCGCTAT

Resulting contig:

CGGCGGATCTGATGGGATCTGATTCGGTTCGGATCTGGGCATGGCGTTCGATGTCGCTAT

Why raw data need to be assembled before various analysis?

150bp

Length of a read

Length of a locus

250bp

Reads are too short to reach the length of the locus

How raw reads magically become sequences of loci of each sample?

Remove PCR duplicates

Parse reads to loci

Assemble parsed reads

Further assemble

Get orthologue assemblies

Remove PCR duplicates

Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

PCR duplicate of Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

PCR duplicate of Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

Reads 2: TCTGATTCGGTTCGGATCTGGGCAT

Resulting contig:

Contig1 CGGCGGATCTGATGGGATCTGATTCGGTTCGGATCTGGGCAT
Remove PCR duplicates

Reads 1: CGGCGGATCTGATGGGATCTGATTCGGTT

Reads 2: TCTGATTCGGTTCGGATCTGGGCAT

Resulting contig:

Contig1 CGGCGGATCTGATGGGATCTGATCGGGTTCGGGATCTGGGCAT

PCR duplicates are redundant for assembly

Remove PCR duplicates



		-
_		
		_

Mixed short reads from lots of loci

Reads of the same color indicate they come from the same loci

I should assembled with which reads?



If reads from different loci assembled together, the resulting contig will be "chimera"

locus 1 locus 2 locus 3

Remember me? I'm the sequence of "reference", used to design the baits



Mixed short reads from several loci

Compare reads with each locus



Mixed short reads from several loci

Reads got different identity with each locus



Mixed short reads from several loci

Select reads with highest identity



Unassigned reads

Most of reads are assigned to different loci. Some reads from nowhere are still unassigned



Assemble parsed reads into longer contigs

Reads:		
Contia:		

Reads:			
Contig: —			

In real case, question is not that easy. We always have loci assigned with more than 2,000 reads

Find overlaps among all reads

Build a graph recording overlaps among all reads

Traverse through the graph to get contigs

Find overlaps among all reads

Align the reads

K-mer

FM-index

Find overlaps among all reads

Align the reads

K-mer

FM-index

Only considerable length of overlap between reads will be kept (25 bp), to guarantee the low probability of accidentally overlap occurs

Build a graph recording overlaps among all reads





Graph of Campylobacter jejuni

Traverse through the graph to get contigs

For each locus, there's only one sequence

Traverse through the all nodes in the graph and each node only pass once

Traverse through the graph to get contigs

For each locus, there's only one sequence

Traverse through the all nodes in the graph and each node only pass once

But this assumption is hard to fulfill



Break the graph into several sub-graph



Discard one of the path



Why we need graph



Which way you should choose







If two path is too diverged (>= 95% identity). The path will be split into several contigs



If two path is too diverged (>= 95% identity). The path will be split into several contigs



Each contig will be aligned to reference. Graph will be reconstructed. The alignment score of each path will be calculated.



Keep the path with higher score.



Why we need it?

Our final aim is to reveal evolutionary history among our enriched species

Orthologues genes are derived from "speciation event". So, the evolutionary history of these genes are identical with the evolutionary history of species

homologs orthologs orthologs paralogs frog^β mouseß frog OL chick OL chick B mouse C.-chain gene hain gene ne duplication early globin gene

What will happen if we use paralog genes to reveal evolutionary history

What will happen if we use paralog genes to reveal evolutionary history



There's various way to find orthologues. The method we used here called reciprocal blast

This method is built on the assumption that orthologous genes have identical or highly related functions and this sharing is greater than for paralogs.

Closest gene between 2 species are potential orthologous gene

Reciprocal blast in general case

(1) sequence of gene which we want to find its orthologous sequence in other organisms(2) genome of these 2 organisms



If a pair of genes in different species are the closest to each other, these 2 genes have "reciprocal best hit"

Reciprocal blast in our pipeline



In our situation, we do not have the genome, only got several contigs only

But, loci of reference and contig have reciprocal best hit, then they are also putative orthologues

Reciprocal blast in our pipeline



Exclude the contig if it does not have reciprocal blast hit with reference loci
Until here we've been reached our first goal

