

Genomic Data Analysis using Python

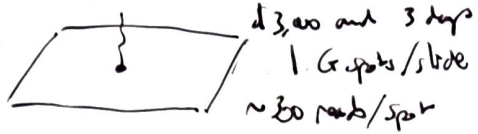
C:\Users\lobsta\AppData\Local\Programs\Python\Python37\python.exe -m venv .env

Genome \rightarrow 23 pairs of chromosomes \sim 3 Gb
 approx 1/1000 nucleotide different for any two copies of the human genome
 \rightarrow polymorphism

Sequencing Primer (5') \rightarrow Sequence \rightarrow 3' (ddNTPs)
 Sanger: detectable (hydrolyt group) nucleotides \rightarrow radioactivity
 \rightarrow gel

Capillary: fluorescent dyes and laser reads

NBS: microarray slide
 \rightarrow "pictures"



Overlap detection

Seq1 = $(x_1, x_2, \dots, x_{300})$

Seq2 = $(y_1, y_2, \dots, y_{300})$

	x_1, x_2, \dots, x_{300}	with matrix
y_1	1 1 0	
y_2	0 0 1	
\vdots		
y_{300}	1 0 0	

$\sim 10^5$ entries



10^7 reads \rightarrow 1-fold coverage

random assembly, reads excess coverage \rightarrow 100-fold coverage

$(10^9 \times 10^9 \text{ fragments})$ $10^9 \times 10^9 \times 10^5 = 10^{23}$ operations \sim 100 z flops
 8 bits letter = 26 bits 4 letters \leftarrow 1 byte
 1 Gb for entire genome
 to save days on faster machine

$4^{10} = 2^{20} \approx 1 \text{ million possible 10-mers of 4 letters}$

AAAAAAAAAA
AAAAAA AAC

↓
table
v-index
||
hash table

Seq 1, pos 2
Seq 1
Seq 1, pos 1



→ 300 different 10-mers for each 300 b read

If two sequences overlap, there should be multiple sequences in the hash table that match both sequences. The two sequences do not have to match all of the same sequences in the hash table because there will be regions that do not overlap within the 300 nucleotide read.

This will reduce the computer 1 million-fold, by aligning only analog sequences. 10-mers are also tradeoff, with 3-mers requiring bigger hash table leading to better certainty, though, small errors can be minimized using 10-mers.

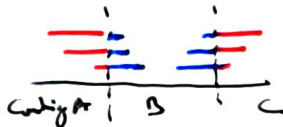
Transposable elements = repeats ~15% of genome

300 b SINES → ~ million copies randomly distributed

300 b LINES → 60,000
(by interspersed elements)

We can identify repetitive sequence in the hash table, looking at unexpected frequencies, to remove from the overlap analysis → gaps in genome.

300b 300b 300b = jumping fragments / clones



jumping clone
mRNA → cDNA ↔ genome

If the mRNA is transcribed then the region of the genome is more likely to be significant to changes in regions that were not transcribed → disease

Sequence Analysis using Python

seqChrums-mm3.fa.zip : July 2017 NCBI 37/mm3

```
> chr6
XXXXXXXXXXXX...
> chr11
XXXXXXXXXX...
```

genome-euro.ucs.c.edu

mm3-sel-chrums-knownGene.txt

name	chrom	start	txStart	txEnd	cdsStart	cdsEnd	cdsStart	cdsEnd	exonStart	exonEnds	proteinId	✓	display
UC003aaw.1	chr6	+3638518	3262019	3238718	xxx	3	xxx,xxx,xxx	xxx,xxx,xxx	xxx,xxx,xxx	xxx,xxx,xxx	uc003aaw.1		

Generated table format

table generated
"A gene prediction."

string	name;	"Name of gene"
string	chrom;	
chr[1]	strand;	
uint	txStart;	
uint	txEnd;	
uint	cdsStart;	
uint	cdsEnd;	
uint	exonStart;	
uint[exonCount]	exonStarts;	
uint[exonCount]	exonEnds;	

Schemas for knownGene:

Generated table +

string	proteinId;	"UniProt display ID, UniProt accession number, RefSeq protein ID"
string	alignId;	"unique identifier (GENCODE + RefSeq ID)"