



# NGS Data Analysis

#### METHODS AND PROTOCOLS







## Sequencing Technology









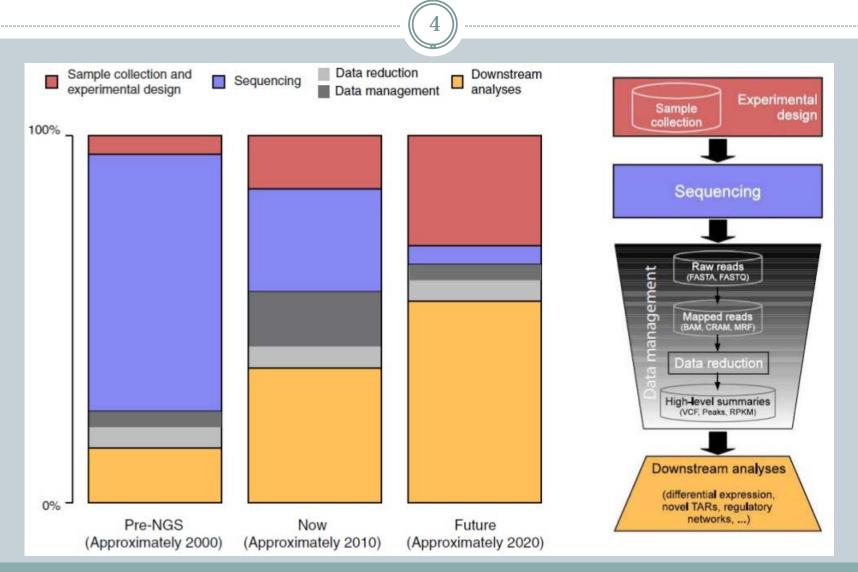
## Changes and Timing past decade







#### Overview of costs (past, present and near future)

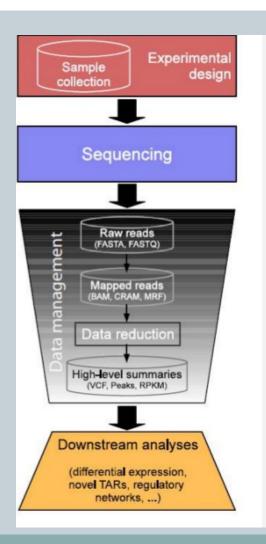






## Steps in sequencing experiments





#### Data analysis

Raw machine reads... What's next?

#### Preprocessing (machine/technology)

- adaptors, indexes, conversions,...
- machine/technology dependent

#### Reads with associated qualities (universal)

- FASTQ
- QC check

#### Depending on application (general applicable)

- 'de novo' assembly of genome (bacterial genomes,...)
- Mapping to a reference genome → mapped reads
  - SAM/BAM/...

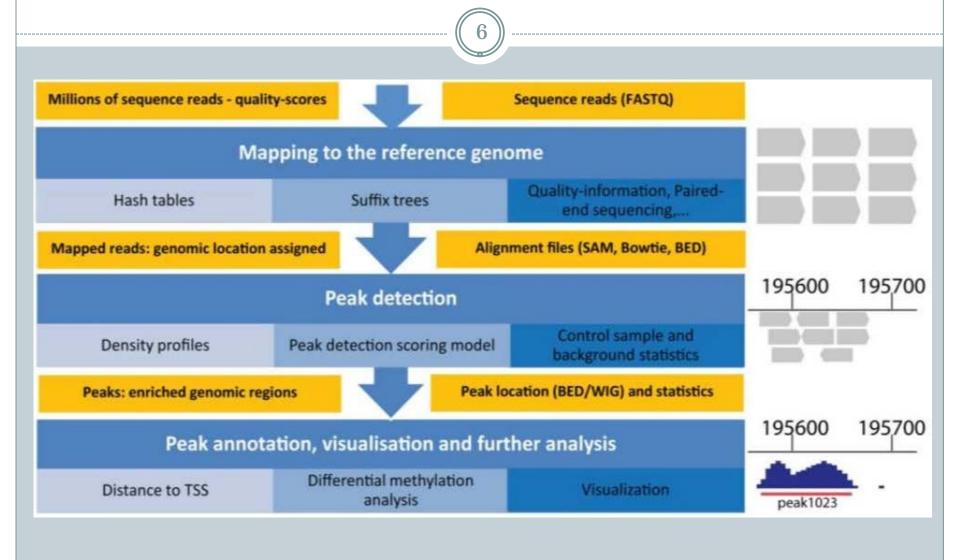
#### High-level analysis (specific for application)

- SNP calling
- Peak calling





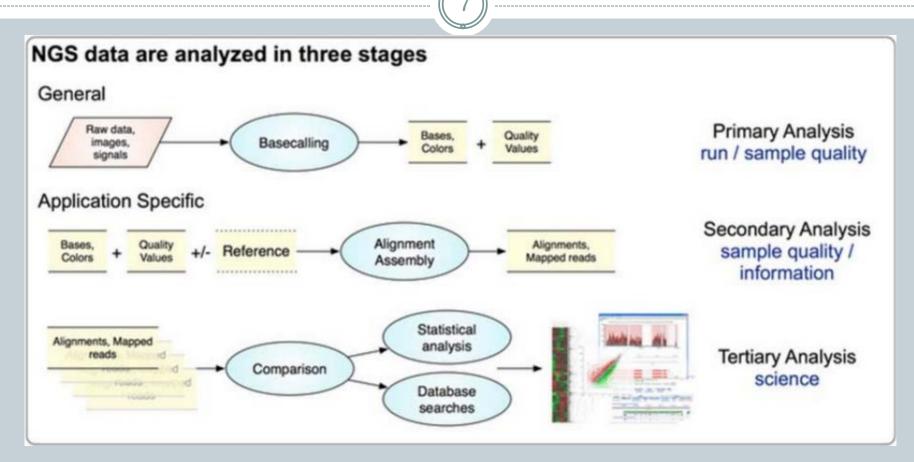
### NGS analysis workflow







## The three stages of NGS data analysis



We will focus mostly on the first two...





#### NGS Applications are **sequencing** applications



- Whole Genome Sequencing
- Gene Regulation
- Epigenetic Changes
- Metagenomics
- Paleogenomics
- Transcriptome Analysis
- Resequencing
- ....







## Why QC and preprocessing

9)

- Sequencer output
  - Reads + quality
- Natural questions
  - Is the quality of my sequenced data ok?
  - If something is wrong, can I fix it?
- Problem: HUGE files

```
@HWI-EAS225:3:1:2:854#0/1
GGGGGGAAGTCGGCAAAATAGATCCGTAACTTCGGG
+HWI-EAS225:3:1:2:854#0/1
a`abbbbabaabbababbb^`[aaa`_N]b^ab^``a
@HWI-EAS225:3:1:2:1595#0/1
GGGAAGATCTCAAAAACAGAAGTAAAACATCGAACG
+HWI-EAS225:3:1:2:1595#0/1
a`abbbababbbabbbbbbbbbbbbbabb`aaababab\aa_`
```





## **Sequencing Data Formats**



#### Raw sequence reads:

Represent the sequence ~ FASTA

>SEQUENCE\_IDENTIFIER
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT

Extension: represent the quality, per base ~ FASTQ – Q for quality

```
@SEQUENCE_IDENTIFIER
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>CCCCCCC65
```

- OK, the strange signs at the last line indicate the quality at the corresponding base...
   But what's the decoding scheme? (Nerd alert ahead !!)
- We want to represent quality scores ~ Phred scores
- Q= -10 log P (with P being the chance of a base called in error)

Phred quality scores are logarithmically linked to error probabilities						
Phred Quality Score	Probability of incorrect base call	Base call accuracy				
20	1 in 100	99 %				
30	1 in 1000	99.9 %				
40	1 in 10000	99.99 %				





## Quality before content



```
@SEQUENCE_IDENTIFIER
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>CCCCCCC65

Example of the identifier line for Illumina data (non-multiplexed):
#@machine_id:lane:tile:x:y:multiplex:pair
@HWUSI-EAS100R:6:73:941:1973#0/1
```

- Phred + 33 → Sanger
- Illumina 1.3 + → Phred +64
- Illumina 1.5 + → Phred +64
- Illumina 1.8 + → Phred +33
- Solid → Sanger

Check your instument + version → FastQC will give you a hint which scoring scheme is probably used

Extensions: FASTQ / FQ





## What is quality?



Genome Res. 1998 Mar;8(3):175-85.

Base-calling of automated sequencer traces using phred. I. Accuracy assessment.

Ewing B1, Hillier L, Wendl MC, Green P.

Author information

#### Abstract

The availability of massive amounts of DNA sequence information has begun to revolutionize the practice of biology. As a result, current large-scale sequencing output, while impressive, is not adequate to keep pace with growing demand and, in particular, is far short of what will be required to obtain the 3-billion-base human genome sequence by the target date of 2005. To reach this goal, improved automation will be essential, and it is particularly important that human involvement in sequence data processing be significantly reduced or eliminated. Progress in this respect will require both improved accuracy of the data processing software and reliable accuracy measures to reduce the need for human involvement in error correction and make human review more efficient. Here, we describe one step toward that goal: a base-calling program for automated sequencer traces, phred, with improved accuracy, phred appears to be the first base-calling program to achieve a lower error rate than the ABI software, averaging 40%-50% fewer errors in the data sets examined independent of position in read, machine running conditions, or sequencing chemistry.

PMID: 9521921 [PubMed - indexed for MEDLINE] Free full text

RESEARCH

Base-Calling of Automated Sequencer Traces Using *Phred.* I. Accuracy Assessment

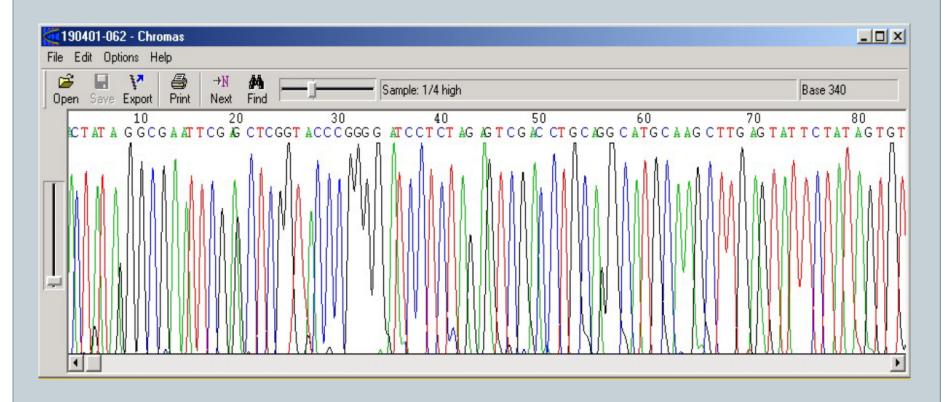
Brent Ewing, <sup>1</sup> LaDeana Hillier, <sup>2</sup> Michael C. Wendl, <sup>2</sup> and Phil Green <sup>1,3</sup>





# Trace File (high quality)



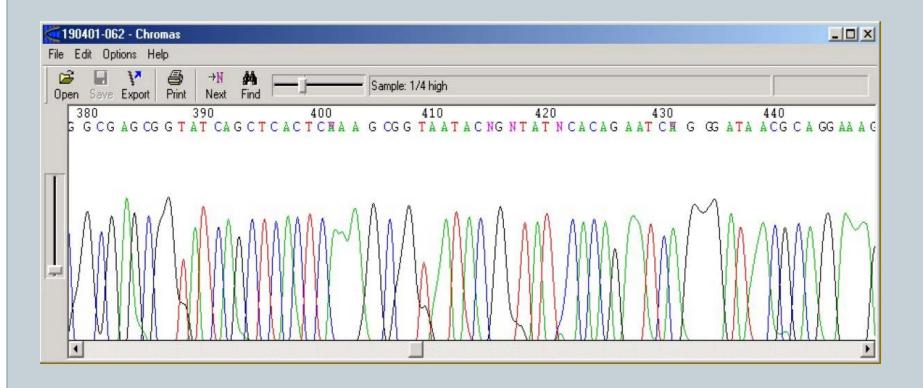






### Trace File (Medium Quality)



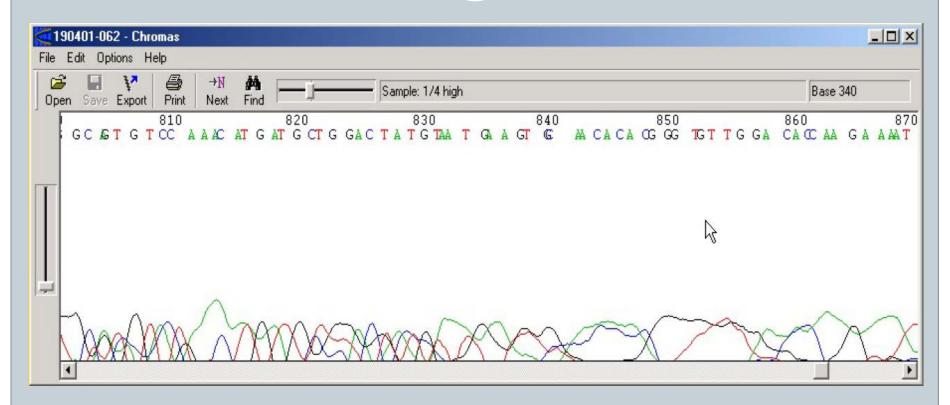






## Trace File (Low Quality)





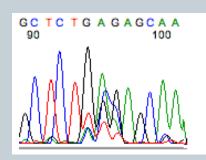


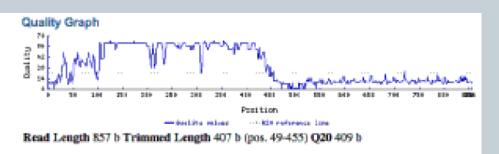


## **Phred Quality Scores**



- Phred is a program that assigns a quality score to each base in a sequence. These scores can then be used to trim bad data from the reads, and to determine how good an overlap actually is
- Phred scores are logarithmically related to the probability of an error:
  - o a score of 10 means 10% error probability,
  - 20 means a 1% chance,
  - o 30 means a 0.1 chance, etc
- A score of 30 is usually considered the minimum acceptable score.









#### **FASTQ File Format**



- Each read is represented by four lines:
- @ followed by read ID
- 2. Sequence
- 3. + optionally followed by repeated read ID
- 4. Quality line
  - Same length as sequence
  - Each character encodes the quality of the respective base

```
@HWI-EAS225:3:1:2:854#0/1
GGGGGGAAGTCGGCAAAATAGATCCGTAACTTCGGG
+HWI-EAS225:3:1:2:854#0/1
a`abbbbabaabbababb^`[aaa`_N]b^ab^``a
@HWI-EAS225:3:1:2:1595#0/1
GGGAAGATCTCAAAAACAGAAGTAAAACATCGAACG
+HWI-EAS225:3:1:2:1595#0/1
a`abbbababbbabbbbbbbbbbbbabb`aaababab\aa_`
```





#### **FASTQC**

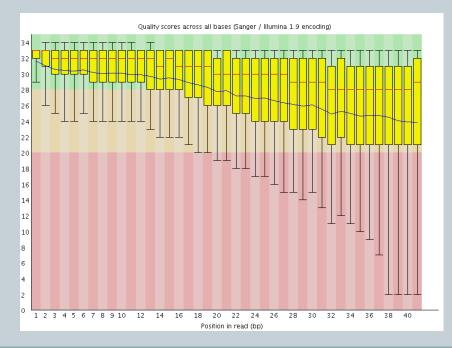
18

 As the name implies, FastQC is way to quickly see some summary statistics to check the quality of your NGS run.

It runs both as a GUI (requires Java) and as a command line

program.

- O Provides several statistics:
  - **Per Sequence Quality**
  - Per sequence quality scores
  - Per base sequence and GC content
  - **▼** Per Sequence GC Content
  - × etc..







## Trimming



- Knowing quality → Act accordingly
- Adapter trimming
  - May increase mapping rates
  - Absolutely essential for small RNA
     Probably Improves de novo assemblies
- Quality trimming
  - May increase mapping rates
  - May also lead to loss of information
- Lots of software:
  - Cutadapt, Trim Galore!, PRINSEQ, etc.





## **Mapped Reads**



- Mapping: "align" these raw reads to a reference genome
  - Single-end or paired-end data?
  - O How would you align a short read to the reference?
- Old-school: Smith-Waterman, BLAST, BLAT,...
- Now: mapping tools for short reads that use intelligent indexing and allow mismatches

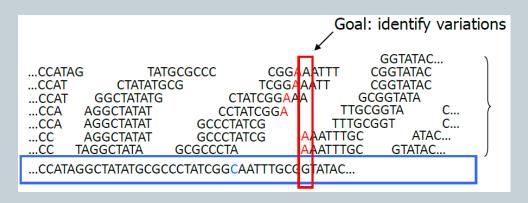




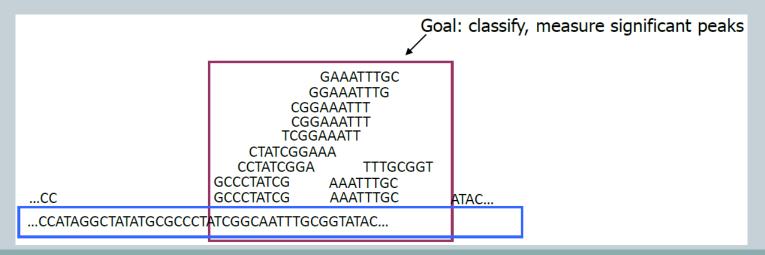
### Short read applications

21)

Genotyping



RNA-Seq, ChIP-Seq, Methyl-Seq,...







## ... always a problem

22

#### Finding the alignment is a computational bottleneck

```
GGTATAC...
...CCATAG
                 TATGCGCCC
                                   CGGAAATTT
                                                  CGGTATAC
...CCAT
                                                  CGGTATAC
             CTATATGCG
                                 TCGGAAATT
...CCAT
         GGCTATATG
                             CTATCGGAAA
                                                GCGGTATA
                                              TTGCGGTA
...CCA
        AGGCTATAT
                           CCTATCGGA
                                            TTTGCGGT
...CCA
        AGGCTATAT
                        GCCCTATCG
                                       AAATTTGC
                                                        ATAC...
...CC
        AGGCTATAT
                        GCCCTATCG
      TAGGCTATA
                                       AAATTTGC
                                                     GTATAC...
...CC
                     GCGCCCTA
...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...
```

GAAATTTGC
GGAAATTTG
CGGAAATTT
CGGAAATTT
TCGGAAATT
CTATCGGAAA
CCTATCGGA
CCTATCGGA
GCCCTATCG AAATTTGC
GCCCTATCG AAATTTGC
GCCCTATCG AAATTTGC

...CCATAGGCTATATGCGCCCTATCGGCAATTTGCGGTATAC...

...CC





## Defining the question



- Given a reference and a set of reads, report at least one "good" local alignment for each read, if one exists
  - Approximate answer to question: where in genome did read originate
- What is "good"? For now we concentrate on:
- Fewer mismatches = better
- Failing to align a low-quality base is better than failing to align a high-quality base

```
...TGATCATA...
GATCAA

better than ...TGATCATA...
GAGAAT

...TGATATTA...
FILL
GATCATA

better than ...TGATCATA...
GATCATA

GATCATA

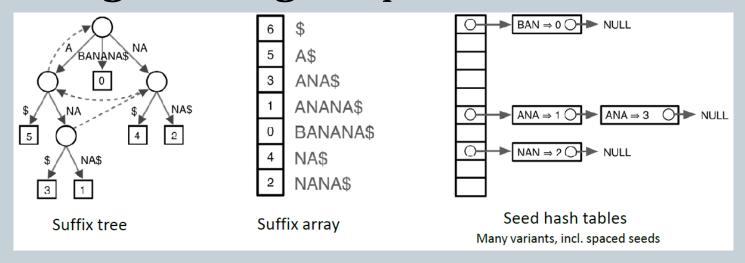
contact co
```





## A few technical aspects (geeky stuff)

- **24**
- Genomes and reads are too large for direct approaches like dynamic programming
- Indexing/Hashing is required



Choice of index is key to performance

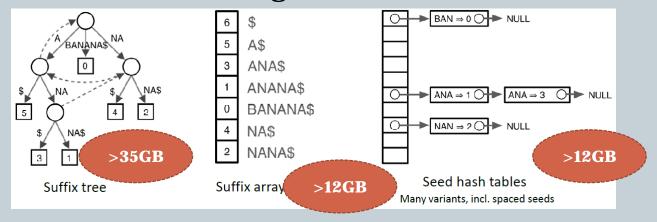




## A few technical aspects (geeky stuff)

(25)

Genome indices can be big. For human:



- Large indices necessitate painful compromises
- 1. Require big-memory machines
- Use secondary storage
- Build new index each run
- 4. Subindex and do multiple passses

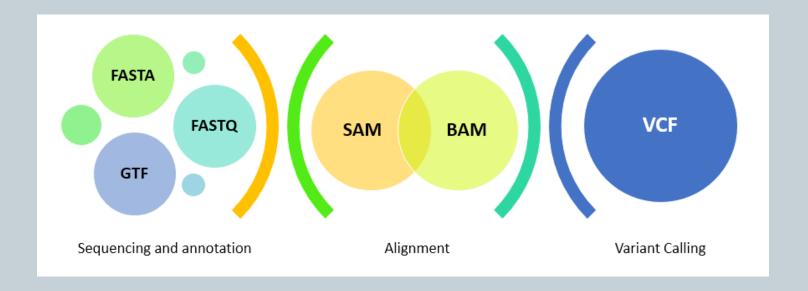




#### Interlude

26)

#### (not only) NGS File Formats







## The Sequence Alignment/Map Format

(27)

- Generic alignment format
- Supports short and long reads
- Supports different sequencing platforms
- Flexible in style, compact in size, computationally efficient to access

- SAM File Format
  - BAM is the binary version of the SAM file; not human readable but indexed for fast access for other tools / visualization / ...





#### **SAM Fields**



```
DESCRIPTION OF THE 11 FIELDS IN THE ALIGNMENT SECTION
# QNAME: template name
#FLAG
#RNAME: reference name
# POS: mapping position
#MAPQ: mapping quality
#CIGAR: CIGAR string
#RNEXT: reference name of the mate/next fragment
#PNEXT: position of the mate/next fragment
#TLEN: observed template length
#SEQ: fragment sequence
#QUAL: ASCII of Phred-scale base quality+33
#Headers
@HD VN:1.3 SO:coordinate
@SO SN:ref LN:45
#Alignment block
r001 163 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5H6M * 0 0 AGCTAA * NM:i:1
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
```







- Browser Extensible Data (location / annotation / scores).
  - o used for mapping / annotation / peak locations
  - o extension: bigBED (binary)

```
FIELDS USED:
# chr
# start
# end
# name
# score
# strand

track name=pairedReads description="Clone Paired Reads" useScore=1
#chr start end name score strand
chr22 1000 5000 cloneA 960 +
chr22 2000 6000 cloneB 900 -
```

- BEDGraph files (location, combined with score)
  - used to represent peak scores

```
track type=bedGraph name="BedGraph Format" description="BedGraph format" visibility=full color=200,100,0 altColor=0,100,200 priority=20 #chr start end score chr19 59302000 59302300 -1.0 chr19 59302300 59302600 -0.75 chr19 59302600 59302900 -0.50
```







- WIG files (location / annotation / scores): wiggle
  - used for visualization or to summarize data, in most cases count data or normalized count data (RPKM)
  - o extension: BigWig binary versions, often used in GEO for

ChIP-seq peaks









#### General Feature Format

- used for annotation of genetic / genomic features, such as all coding genes in Ensembl
- often used in downstream analysis to assign annotation to regions/peaks/....

```
FIELDS USED:
# segname (the name of the sequence)
# source (the program that generated this feature)
# feature (the name of this type of feature - for example: exon)
# start (the starting position of the feature in the sequence)
# end (the ending position of the feature)
# score (a score between 0 and 1000)
# strand (valid entries include '+', '-', or '.')
# frame (if the feature is a coding exon, frame should be a number between
0-2 that represents the reading frame of the first base. If the feature is
not a coding exon, the value should be '.'.)
# group (all lines with the same group are linked together into a single
item)
track name=regulatory description="TeleGene(tm) Regulatory Regions"
#chr source feature start
                                  end scores tr fr group
chr22 TeleGene enhancer 1000000 1001000 500 + . touch1
chr22 TeleGene promoter 1010000 1010100 900 + . touch1
chr22 TeleGene promoter 1020000 1020000 800 - . touch2
```





32

#### Variant Call Format

used for SNP representation

```
##fileformat=VCFv4.1
##fileDate=20090805
##source=mylmputationProgramV3.1
##reference=file:///seg/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sapiens".taxonomy=x>
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP.Number=1.Type=Integer.Description="Total Depth">
##INFO=<ID=AF, Number=A, Type=Float, Description="Allele Frequency">
##INFO=<ID=AA, Number=1, Type=String, Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2, Number=0, Type=Flag, Description="HapMap2 membership">
##FILTER=<ID=q10, Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP, Number=1, Type=Integer, Description="Read Depth">
##FORMAT=<ID=HQ.Number=2,Type=Integer,Description="Haplotype Quality">
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT NA00001 NA00002 NA00003
20 14370 rs6054257 G A 29 PASS NS=3;DP=14;AF=0.5;DB;H2 GT;GQ;DP;HQ 0|0;48:1:51,51 1|0;48:8:51,51 1/1:43:5:...
20 17330 . T A 3 q10 NS=3;DP=11;AF=0.017 GT:GQ:DP:HQ 0[0:49:3:58,50 0]1:3:5:65,3 0/0:41:3
20 1110696 rs6040355 A G,T 67 PASS NS=2;DP=10;AF=0.333,0.667;AA=T;DB GT;GQ;DP;HQ 1|2;21;6;23,27 2|1;2:0;18,2 2/2;35;4
20 1230237. T. 47 PASS NS=3;DP=13;AA=T GT:GQ:DP:HQ 0|0:54:7:56,60 0|0:48:4:51,51 0/0:61:2
20 1234567 microsat1 GTC G.GTCT 50 PASS NS=3:DP=9:AA=G GT:GQ:DP 0/1:35:40/2:17:2 1/1:40:3
```





# aaaand back to the story

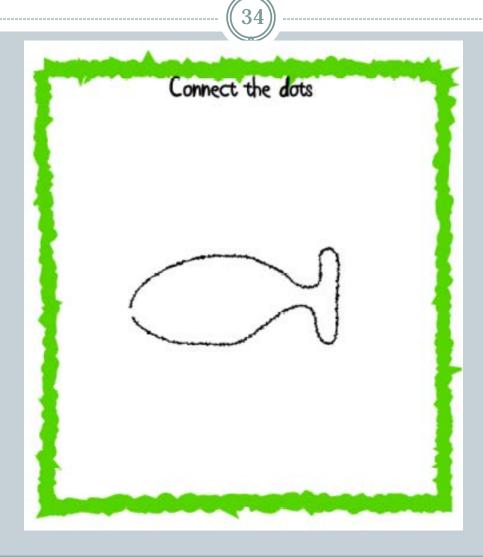
33)







# Assembly simplified



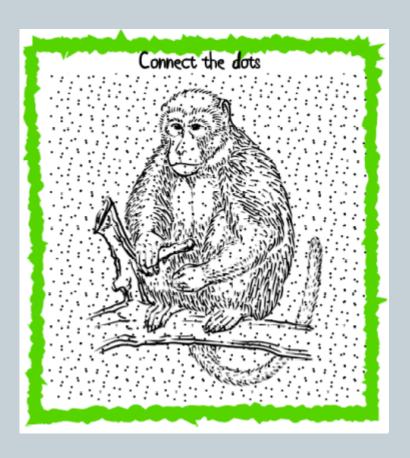


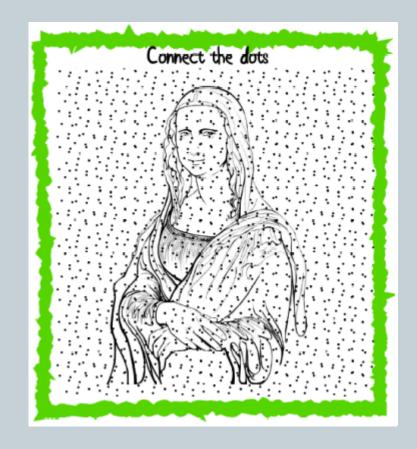


## Assembly simplified

35)

Impossible to assemble manually



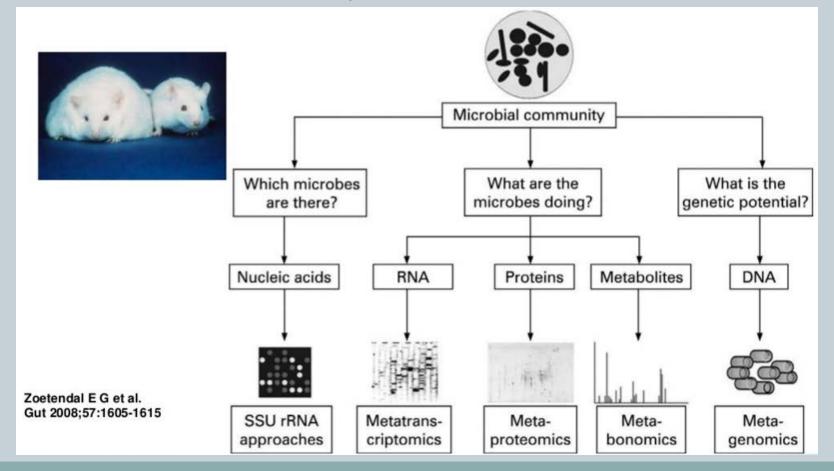






#### Metagenomics

(and other community based "omics")

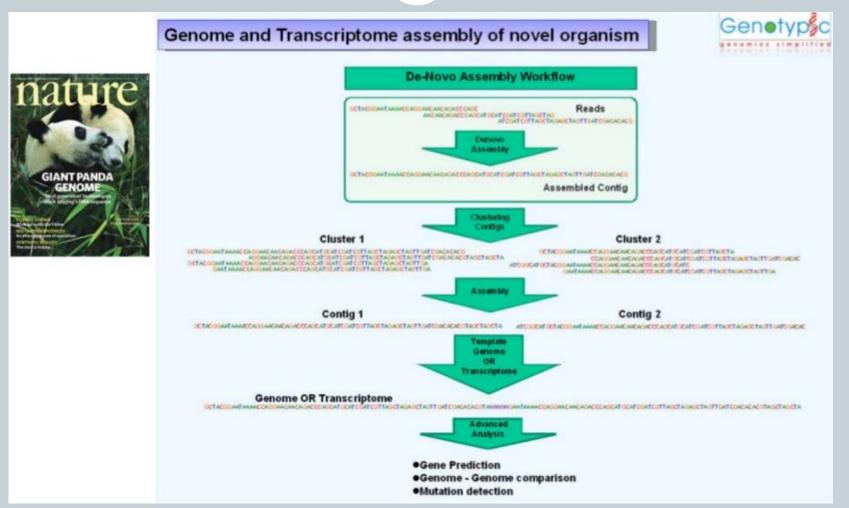






## De-novo sequencing









#### **BowTie**

(38)

- BowTie is the most commonly used aligner
- Employs an indexing algorithm that can trade flexibility between memory usage and running time
- For Human data (NCBI 36.3) on an 2.4 GHz AMD Opteron:

Physical memory Target	Actual peak memory footprint	Wall clock time
16 GB	14.4 GB	4h:36m
8 GB	5.84 GB	5h:05m
4 GB	3.39 GB	7h:40m
2 GB	1.39 GB	21h:30m





### **TopHat**

(39)

- TopHat is one of many applications for aligning short sequence reads to a reference genome.
- It uses the BOWTIE aligner internally.
- Genome alignments from TopHat were saved as BAM files, the binary version of SAM (samtools.sourceforge.net/).
- Other alternatives are BWA, MAQ, OLego, Stampy, Novoalign, etc





#### We've aligned the data. Then what?

40

Depending on the target study.

Gene	Treatment 1				Treatment 2		
1	14	18	10	47	13	24	
2	10	3	15	1	11	5	
3	1	0	10	80	21	34	
4	0	0	0	0	2	0	
5	4	3	3	5	33	29	
•	•	•	•	•	•	•	
•	•	•	•	•	•	•	
•	•	•	•	•	•	•	
53256	47	29	11	71	278	339	

Total 22910173 30701031 18897029 20546299 28491272 27082148





## Differential Expression

41

 To determine if gene 1 is DE, we would like to know whether the proportion of reads aligning to gene 1 tends to be different for experimental units that received treatment 1 than for experimental units that received treatment 2

```
14 out of 22910173 47 out of 20546299
```





#### **Cufflinks**



- CuffLinks is a program that assembles aligned RNA-Seq reads into transcripts, estimates their abundances, and tests for differential expression and regulation transcriptome-wide.
- CuffDiff is a program within CuffLinks that compares transcript abundance between samples





## Putting it all together

43

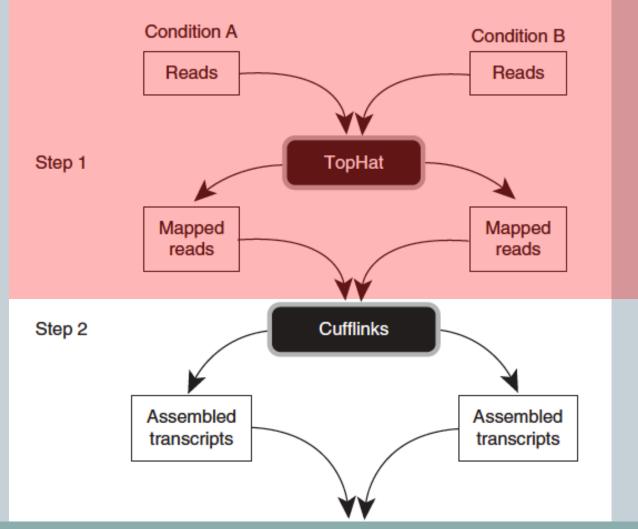
- There are several protocols, tools and (now) platforms that are specific to NGS data analysis:
  - Tuxedo protocol
  - Galaxy platform
  - Chipster platform





#### Tuxedo





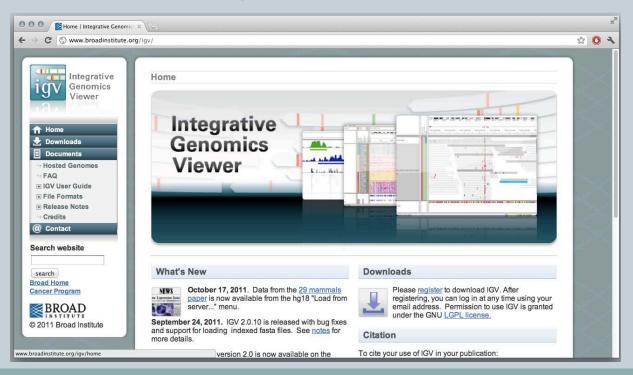




## In closing: Visualization is key



- The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets.
  - It supports a wide variety of data types, including array-based and nextgeneration sequence data, and genomic annotations.



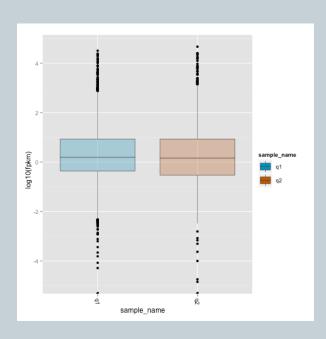


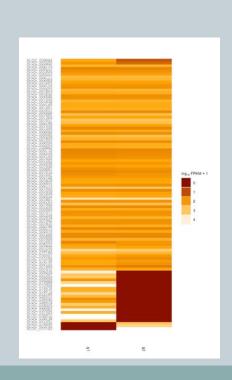


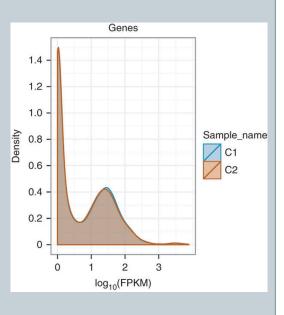
#### CummeRBund

46

#### Downstream Analysis







INA3,

atory La

(47)

# Thank you for your patience!