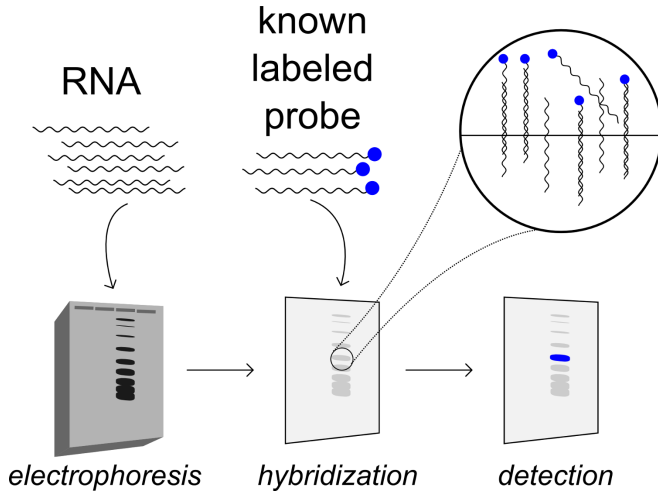# RNA-Seq I
## Measuring Gene Expression

Antonio Marco
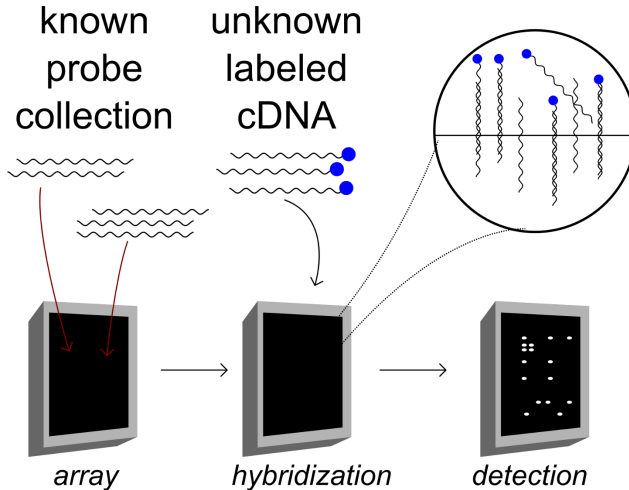
School of Biological Sciences
University of Essex

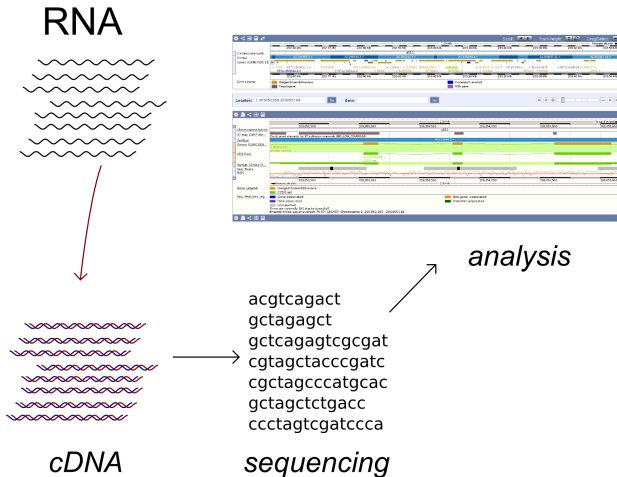6-Jun-17

The 'IS IT THERE?' approach



RNA

known labeled probe
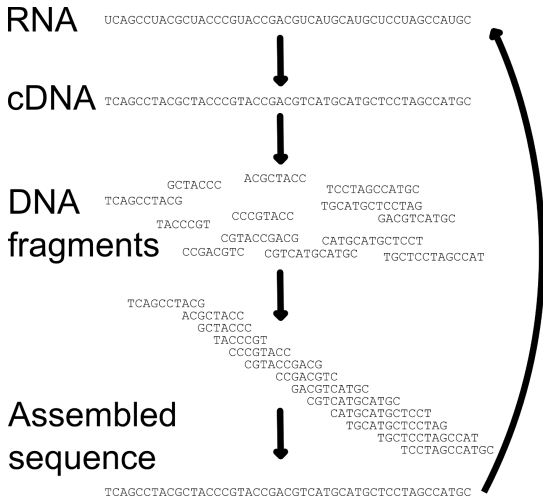
*electrophoresis*  *hybridization*  *detection*

The 'IS ANY OF THESE?' approach



known probe collection

unknown labeled cDNA

*array*    *hybridization*    *detection*

The 'WHAT'S IN THERE?' approach

RNA



analysis

acgtcagact
gctagagct
gctcagagtcgcgat
cgtagctacccgatc
cgctagcccatgcac
gctagctctgacc
ccctagtcgatccca

cDNA    sequencing
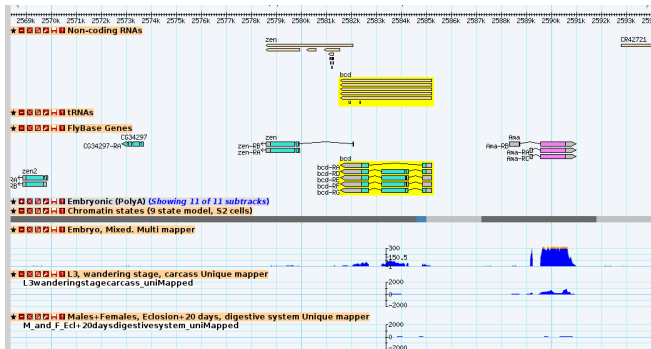
DNA sequencing is done in small fragments

RNA-Seq is a quantitative technique

- Ideally, sequencers always give the actual reads

**ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

**MACHINE READS** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

- Ideally, sequencers always give the actual reads

  **ACTUAL**
  **SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC
  **READS**

- In reality, they often contain errors

  **ACTUAL**
  **SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE** CACC**T**TGCCATC**C**CATCCGATCGCATCGCA**AAC**CA**CT**
  **READS**

- Ideally, sequencers always give the actual reads

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

- In reality, they often contain errors

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACC**T**TGCCATC**C**CATCCGATCGCATCGCA**AAC**CA**CT**

- Good news is, sequencers tell us how confident they are

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACC**T**TGCCATC**C**CATCCGATCGCATCGCA**AAC**CA**CT**

  **MACHINE 'THINKS'**     it is a C     may not be a C     it is a T     can't tell

- Ideally, sequencers always give the actual reads

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC
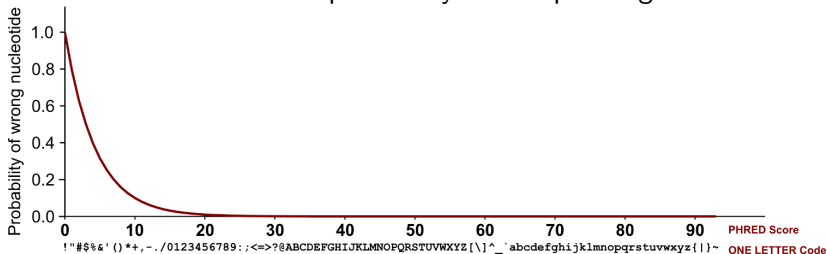
- In reality, they often contain errors

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACC**T**TGCCATC**C**CATCCGATCGCATCGCA**AAC**CA**CT**

- Good news is, sequencers tell us how confident they are

  **ACTUAL SEQUENCE** CACCGTGCCATCGCATCCGATCGCATCGCATCGCATC

  **MACHINE READS** CACC**T**TGCCATC**C**CATCCGATCGCATCGCA**AAC**CA**CT**

  **QUALITY SCORE** ~~~~**?**~~~~~~~**?**~~~~~~~~~~~~~~~~~~**!!!**~~**!!**

- Phred score measures the probability of a sequencing error

- Phred score measures the probability of a sequencing error
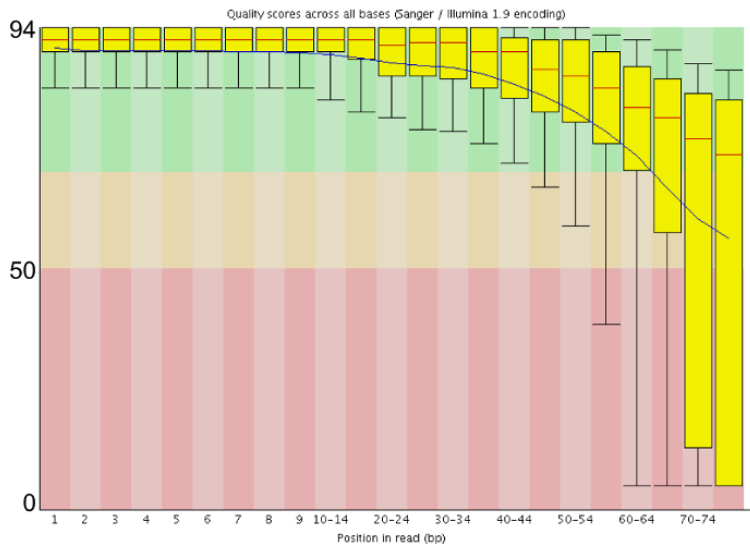


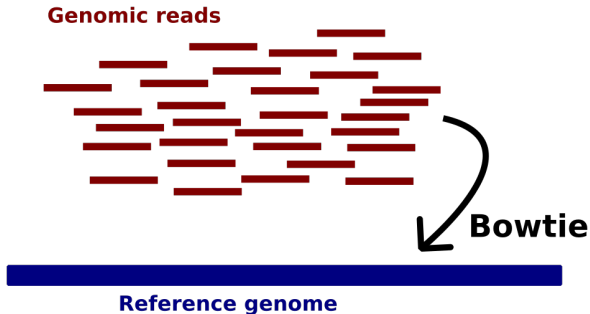- The FASTQ format includes Phred scores in a one-letter code

```
@SEQUENCE_NAME
CATGGCTAGCTGCTAGCTAGCTAGACATTCATCGAAATCGCTAGCCTAGCTACGA
+
!''*((((*∗∗+))%%%++)(%%%).1**∗-+*''))**55CCF>>>>>>C%%%
```
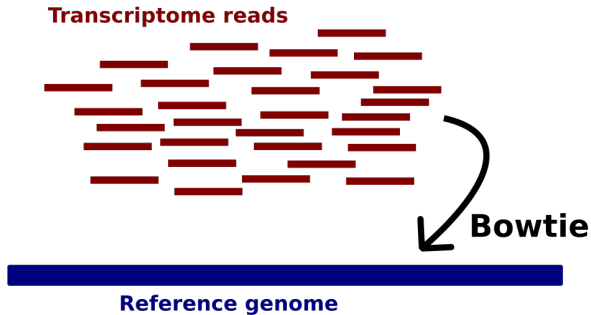
Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Mapping reads across
exon-intron boundaries