

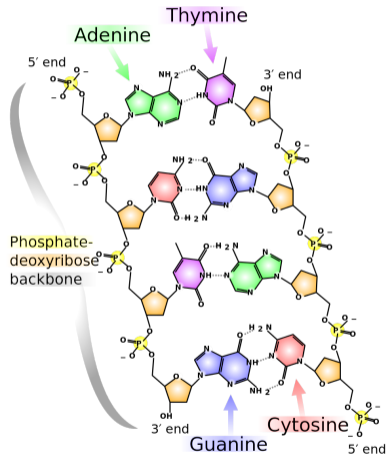
# Introduction to High-Throughput Sequencing and RNA-seq

L Collado-Torres

March 6th, 2013

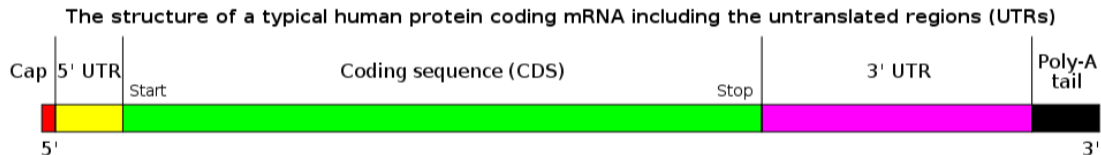
# 1 High Throughput Sequencing

## 2 Sources of variation

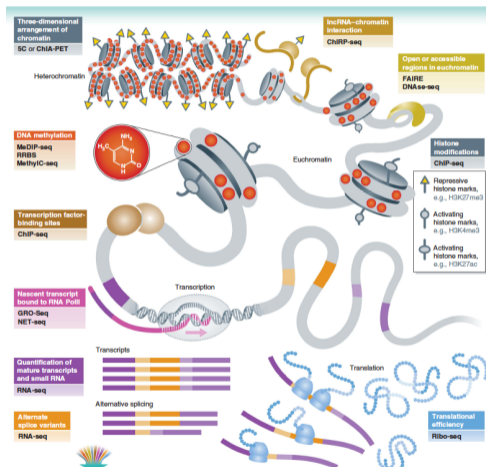
DNA<sup>1</sup>

<sup>1</sup>Wikipedia. DNA. URL: <http://en.wikipedia.org/wiki/DNA> (visited on 03/05/2013).

# Human mRNA<sup>2</sup>



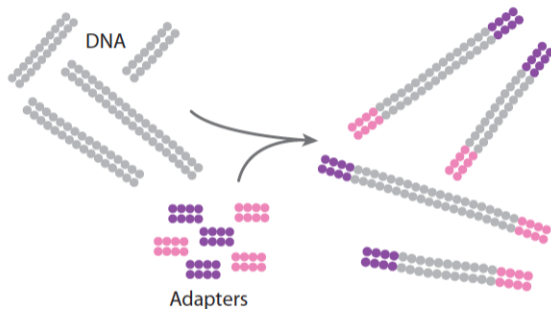
<sup>2</sup>Wikipedia. *Messenger RNA*. URL: [http://en.wikipedia.org/wiki/Messenger\\_RNA](http://en.wikipedia.org/wiki/Messenger_RNA) (visited on 03/05/2013).

Panorama<sup>3</sup>

<sup>3</sup>Wendy Weijia Soon, Manoj Hariharan, and Michael P. Snyder. "High-throughput sequencing for biology and medicine". In: *Molecular Systems Biology* 9.1 (). UR

[http://www.nature.com/msb/journal/v9/n1/fig\\_tab/msb201261\\_F2.html](http://www.nature.com/msb/journal/v9/n1/fig_tab/msb201261_F2.html) (visited on 03/05/2013).

# Prepare DNA<sup>4</sup>



## Prepare genomic DNA sample

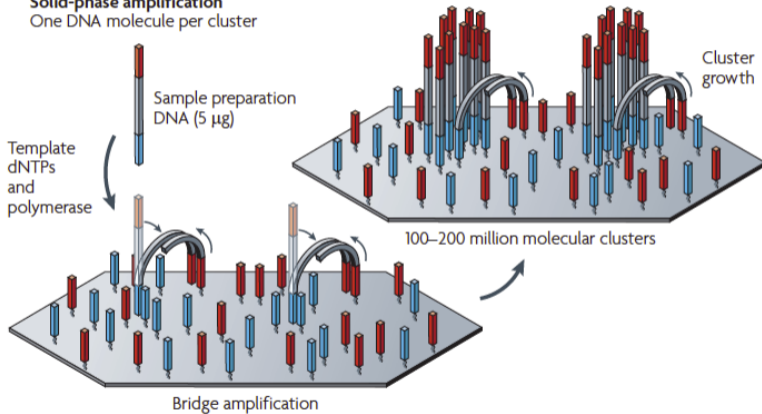
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

<sup>4</sup>Elaine R Mardis. "Next-generation DNA sequencing methods". In: *Annual Review of Genomics and Human Genetics* 9 (2008). PMID: 18576944.

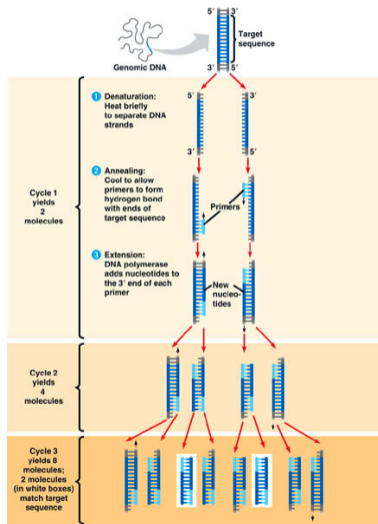
Amplify<sup>5</sup>

**b** Illumina/Solexa  
Solid-phase amplification

One DNA molecule per cluster



<sup>5</sup>Michael L. Metzker. "Sequencing technologies — the next generation". In: *Nat Rev Genet* 11.1 (2010).

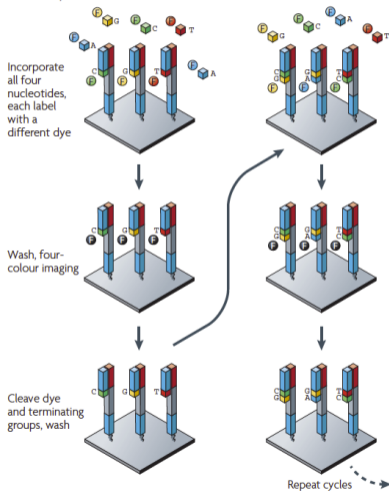
PCR<sup>6</sup>

<sup>6</sup>SCHOOLWORKHELPER. PCR: Uses, Steps, Purpose. URL: <http://schoolworkhelper.net/pcr-uses-steps-purpose/> (visited on 03/05/2013).

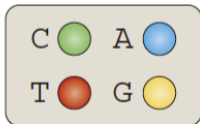
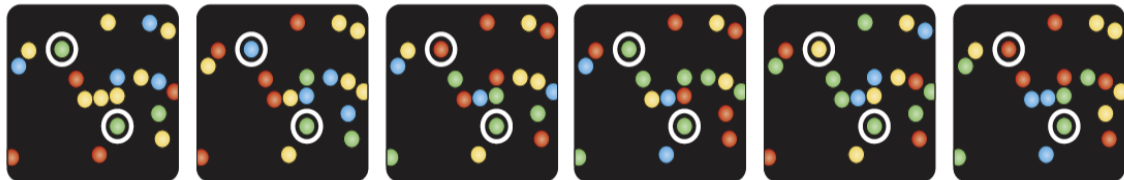


Sequencing by synthesis<sup>7</sup>

## a Illumina/Solexa — Reversible terminators



<sup>7</sup>Michael L. Metzker. "Sequencing technologies — the next generation". In: *Nat Rev Genet* 11.1 (2010).

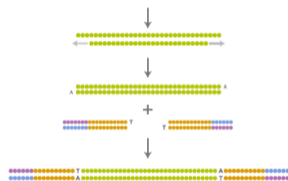
Analyze cluster images<sup>8</sup>

Top: CATCGT  
Bottom: CCCCCC

<sup>8</sup>Michael L. Metzker. "Sequencing technologies — the next generation". In: *Nat Rev Genet* 11.1 (2010).

HiSeq 2000<sup>9</sup>

Figure 3: Next-Generation Sequencing Simplified



Library Preparation  
<6 h (<3 h hands-on)



Cluster Generation  
<4 h (<10 min hands-on)



Sequencing by Synthesis  
1.5-8 days (<10 min hands-on)



RTA v1.7, CASAVA v1.7  
2 days (30 min hands-on)

From simplified sample preparation kits, to automated cluster generation, to streamlined sequencing by synthesis, to complete data analysis, Illumina's HiSeq 2000 sequencing system offers the industry's simplest next-generation sequencing workflow.

<sup>9</sup>Illumina. *HiSeq 2000 Sequencing System*. URL: [http://www.illumina.com/documents/products/datasheets/datasheet\\_hiseq2000.pdf](http://www.illumina.com/documents/products/datasheets/datasheet_hiseq2000.pdf) (visited on 03/05/2011)

# HiSeq 2000

More info on this blog post <http://www.politigenomics.com/2010/01/hiseq-2000.html>

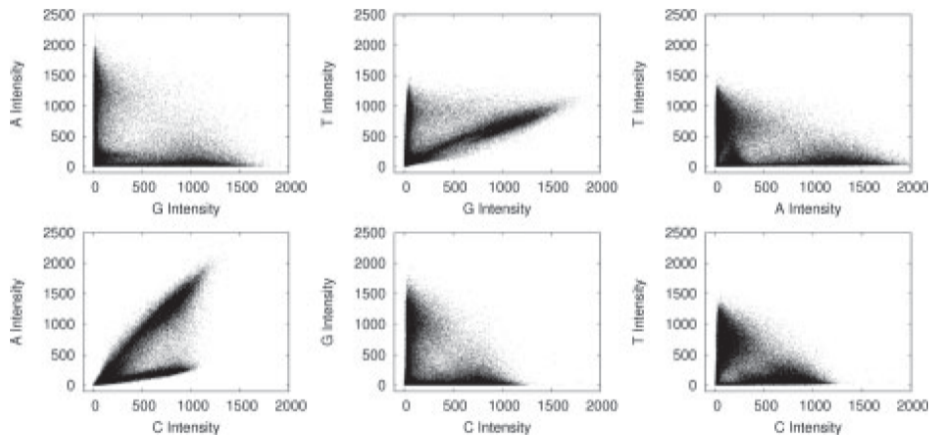
Other 2nd generation sequencers<sup>10</sup>

Sequencer	454 GS FLX	HiSeq 2000	SOLiDv4	Sanger 3730xl
Sequencing mechanism	Pyrosequencing	Sequencing by synthesis	Ligation and two-base coding	Dideoxy chain termination
Read length	700 bp	50SE, 50PE, 101PE	50 + 35 bp or 50 + 50 bp	400~900 bp
Accuracy	99.9%*	98%, (100PE)	99.94% *raw data	99.999%
Reads	1 M	3 G	1200~1400 M	—
Output data/run	0.7 Gb	600 Gb	120 Gb	1.9~84 Kb
Time/run	24 Hours	3~10 Days	7 Days for SE 14 Days for PE	20 Mins~3 Hours
Advantage	Read length, fast	High throughput	Accuracy	High quality, long read length
Disadvantage	Error rate with polybase more than 6, high cost, low throughput	Short read assembly	Short read assembly	High cost low throughput

<sup>10</sup>Lin Liu et al. "Comparison of next-generation sequencing systems". In: *Journal of biomedicine & biotechnology* (2012). PMID: 22829749.

1 High Throughput Sequencing

2 Sources of variation

Cross-talk<sup>11</sup>

<sup>11</sup>Nava Whiteford et al. "Swift: primary data analysis for the Illumina Solexa sequencing platform". In: *Bioinformatics (Oxford, England)* 25.17 (2009). PMID: 19549

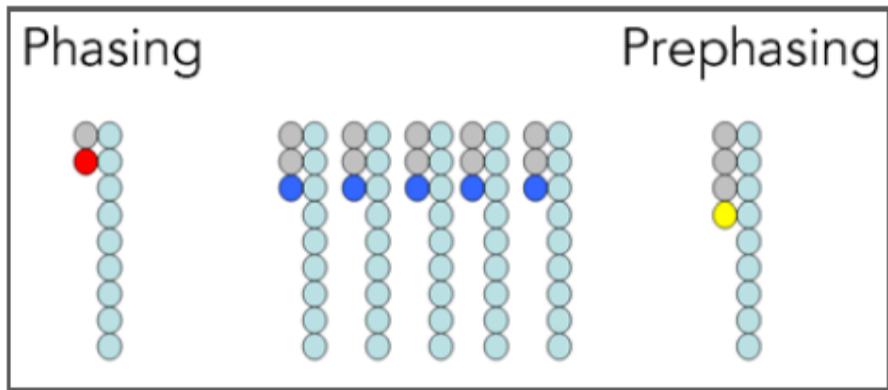
Phasing and pre-phasing<sup>12</sup>

Figure 3 Phasing and Prephasing

<sup>12</sup> Illumina. Pipeline CASAVA User Guide 15003807 ( Pipeline V. 1.4 and Casava V.1.0).

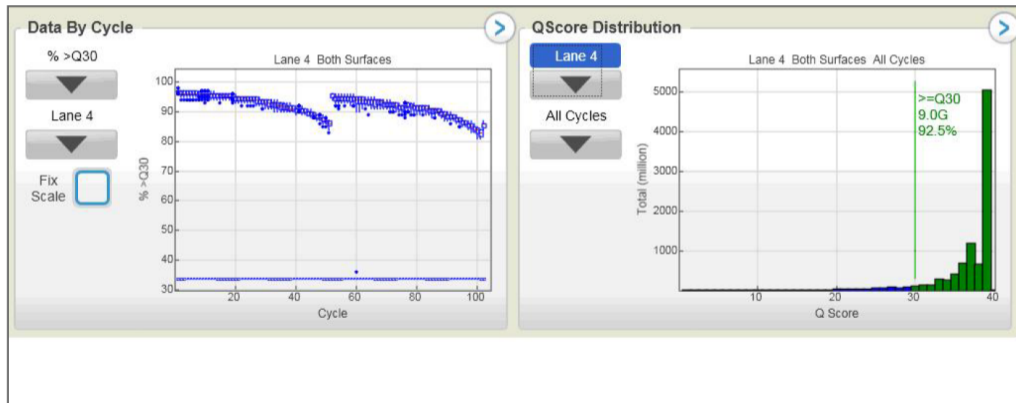


Phasing example<sup>13</sup>

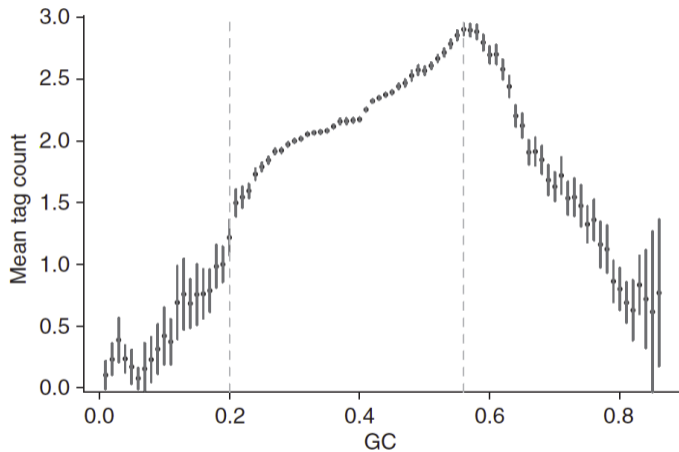
<sup>13</sup>Nava Whiteford et al. "Swift: primary data analysis for the Illumina Solexa sequencing platform". In: *Bioinformatics (Oxford, England)* 25.17 (2009). PMID: 19549

# Sequence quality<sup>14</sup>

Figure 3 SAV Screenshot Showing Excellent Quality Metrics

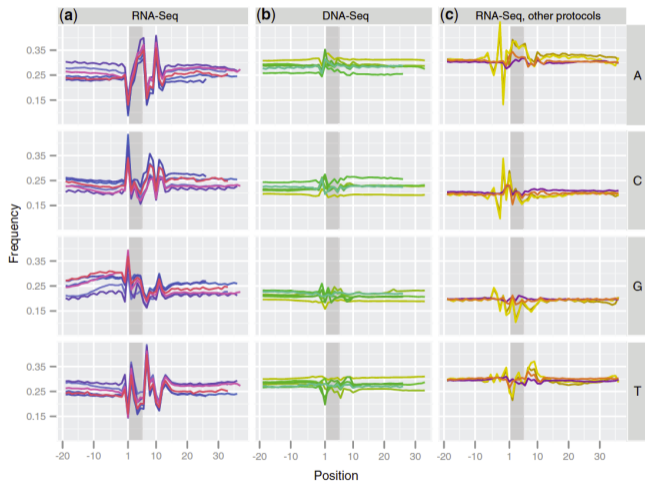


<sup>14</sup> Illumina. CASAVA User Guide (15011196 D). URL: [http://support.illumina.com/downloads/casava\\_user\\_guide\\_15011196.ilmn](http://support.illumina.com/downloads/casava_user_guide_15011196.ilmn) (visited on 03/05/2013).

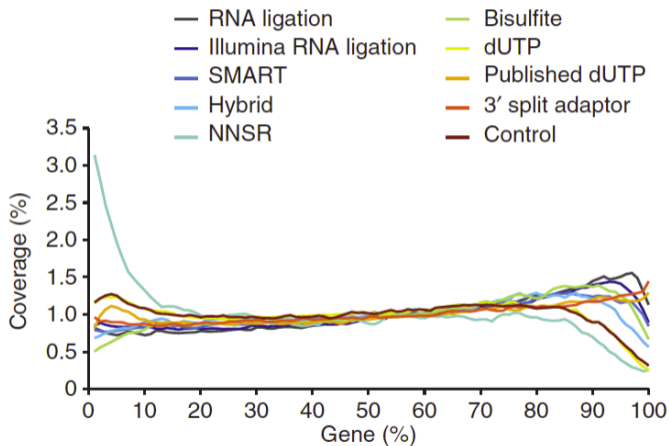
GC bias<sup>15</sup>

<sup>15</sup>Margaret A Taub, Hector Corrada Bravo, and Rafael A Irizarry. "Overcoming bias and systematic errors in next generation sequencing data". In: *Genome Medicine* 2.12 (2010). PMID: 21144010.

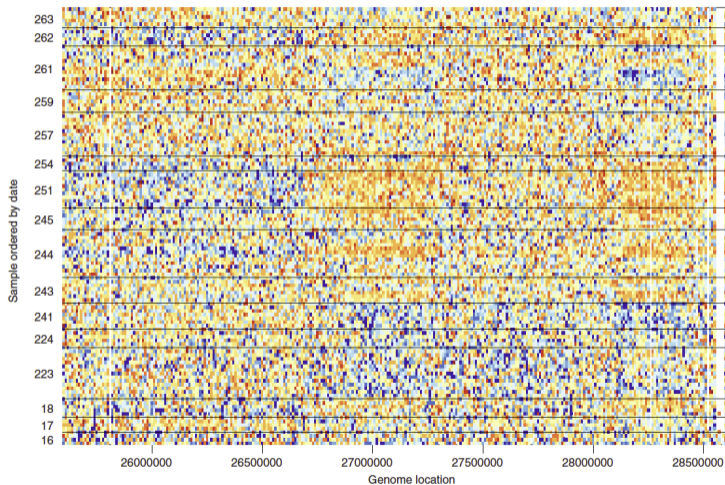
# Random primers bias<sup>16</sup>



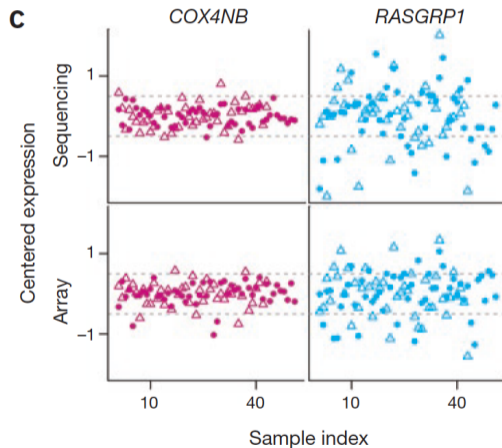
<sup>16</sup>Kasper D Hansen, Steven E Brenner, and Sandrine Dudoit. "Biases in Illumina transcriptome sequencing caused by random hexamer priming". In: *Nucleic Acids Research* (2010). PMID: 20395217.

Library type<sup>17</sup>

<sup>17</sup>Joshua Z Levin et al. "Comprehensive comparative analysis of strand-specific RNA sequencing methods". In: *Nature Methods* 7.9 (2010). PMID: 20711195.

Batch effects<sup>18</sup>

<sup>18</sup>Margaret A Taub, Hector Corrada Bravo, and Rafael A Irizarry. "Overcoming bias and systematic errors in next generation sequencing data". In: *Genome Medicine* 2.12 (2010). PMID: 21144010.

Biological variability<sup>19</sup>

<sup>19</sup>Kasper D Hansen et al. "Sequencing technology does not eliminate biological variability". In: *Nature biotechnology* 29.7 (2011). PMID: 21747377.

## The future

- Further improvements in library preparation
- Single cell sequencing
- Third generation sequencers like Pacific Biosciences

And biostatistical methods =)



# Thanks!

- Google Calendar  
[https://www.google.com/calendar/embed?src=7hprep991i5prd5l5ftksbsfb8%40group.calendar.google.com&ctz=America/New\\_York](https://www.google.com/calendar/embed?src=7hprep991i5prd5l5ftksbsfb8%40group.calendar.google.com&ctz=America/New_York)
- Slides at <http://www.biostat.jhsph.edu/~lcollado/misc/HTSintro.pdf>