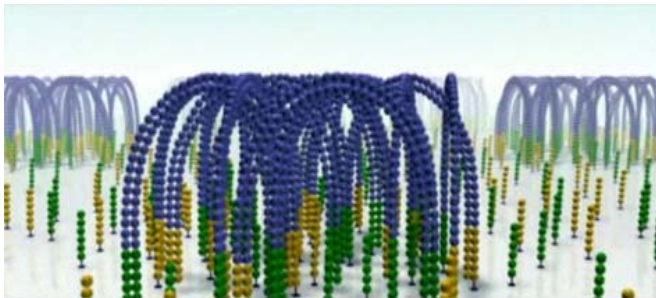# Innovations in Genomic Analysis:
## Downstream analysis of Illumina Sequencing Data

Marco Cappelletti
Product Marketing Manager
Europe

# Agenda
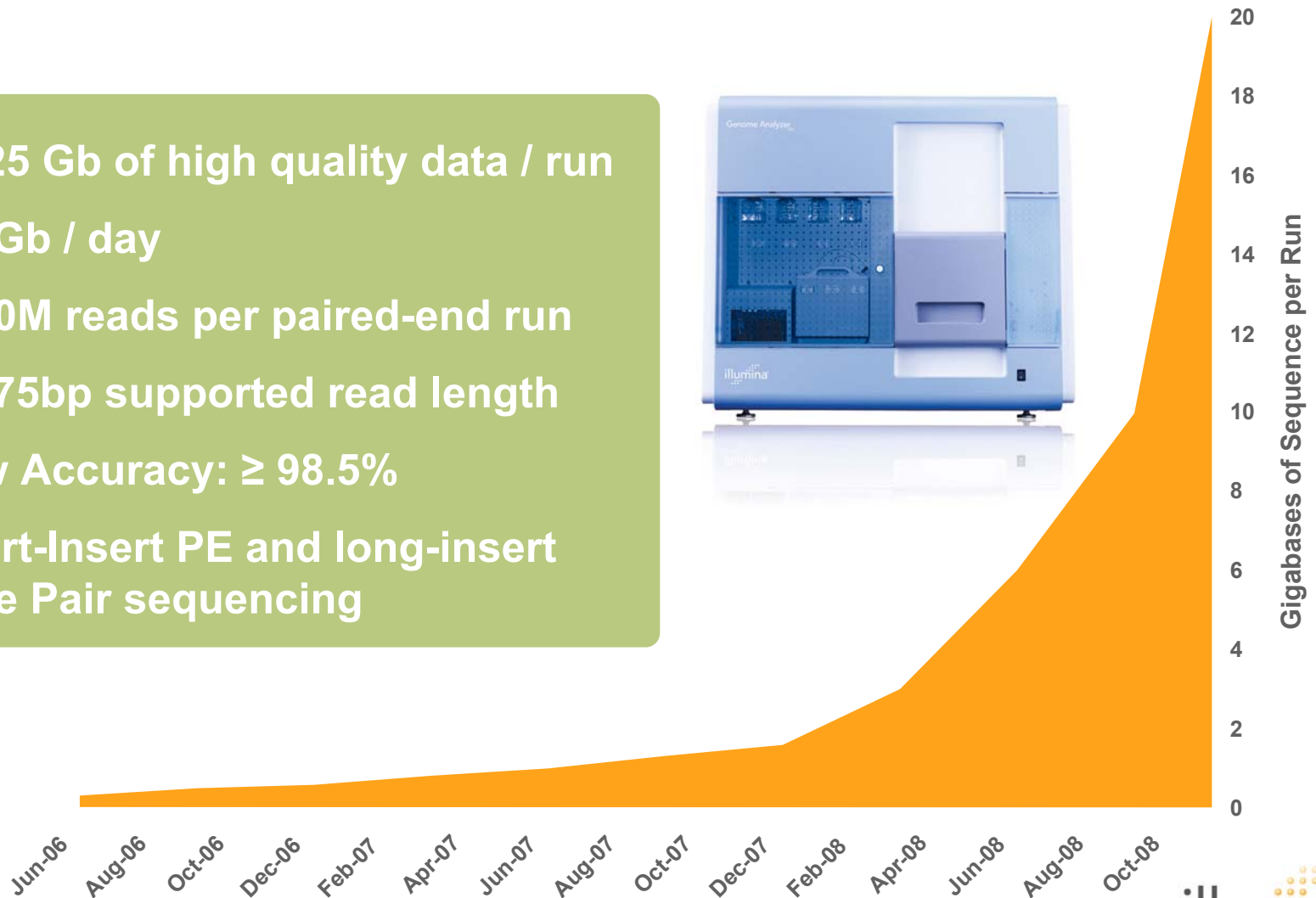
- Genome Analyzer IIx Sequencing Technology

- Applications overview

- GA*II*x SW Improvements
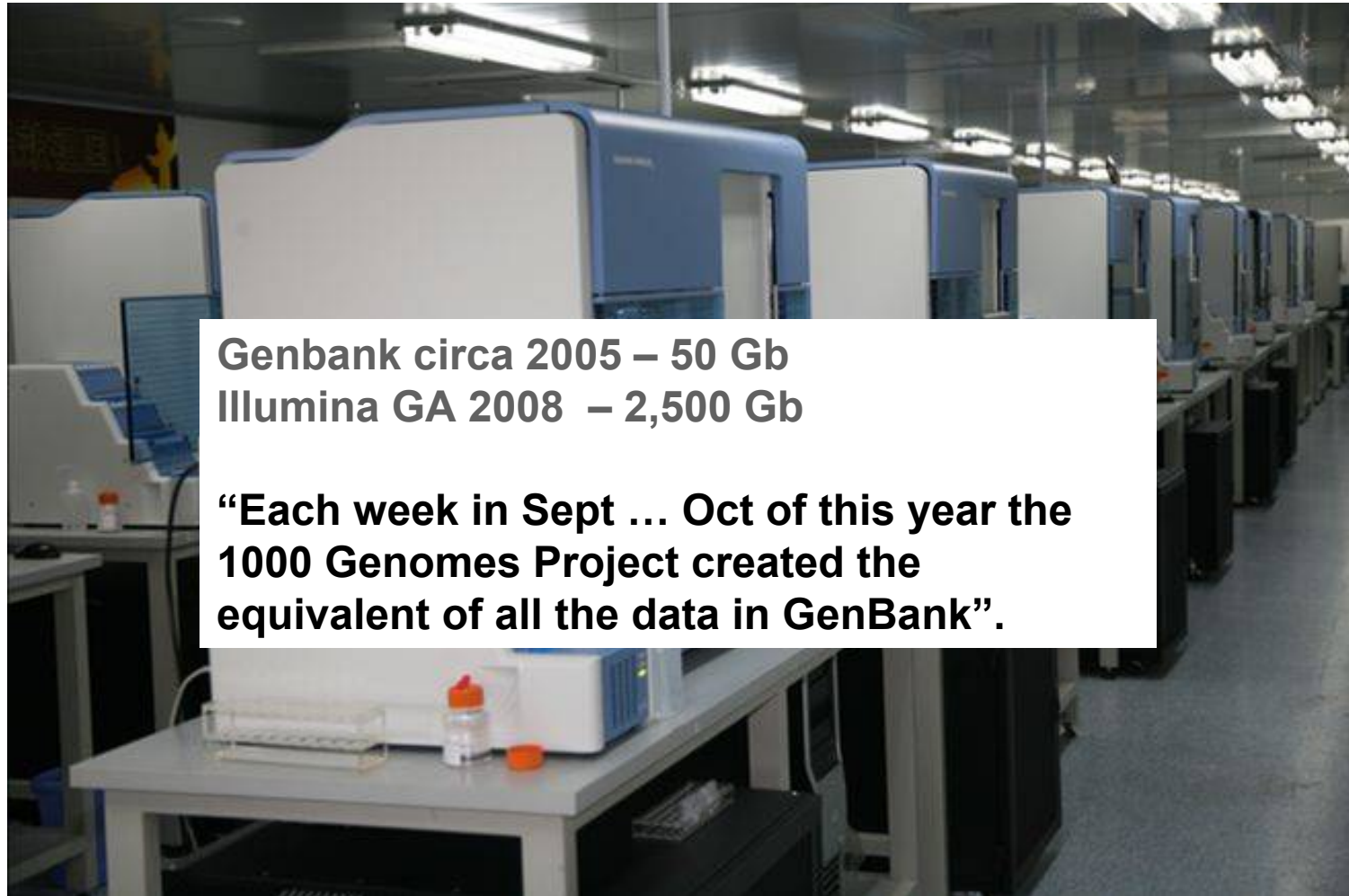
- Downstream Data Analysis

# The Genome Analyzer$_{IIx}$ and Software Advancements
## *65% Increase in data output*

- **20-25 Gb of high quality data / run**
- **2.5 Gb / day**
- **>300M reads per paired-end run**
- **2 x 75bp supported read length**
- **Raw Accuracy: ≥ 98.5%**
- **Short-Insert PE and long-insert Mate Pair sequencing**

Gigabases of Sequence per Run

20
18
16
14
12
10
8
6
4
2
0

Jun-06  Aug-06  Oct-06  Dec-06  Feb-07  Apr-07  Jun-07  Aug-07  Oct-07  Dec-07  Feb-08  Apr-08  Jun-08  Aug-08  Oct-08

illumina®

# Illumina Genome Analyzer: A paradigm shift

Genbank circa 2005 – 50 Gb
Illumina GA 2008  – 2,500 Gb

"Each week in Sept … Oct of this year the 1000 Genomes Project created the equivalent of all the data in GenBank".

illumina®

# Simplest Sequencing Process - 5 to 10 days WF

**1** *Library prep (~ 6 hrs)*

Fragment DNA

⬇

Repair ends / Add A overhang

⬇

Ligate adapters

⬇

Select ligated DNA

**2** *Automated Cluster Generation (~ 5 hrs)*

**Up to 96 samples**

Hybridize to flow cell

⬇

Extend hybridized oligos

⬇

Perform bridge amplification

**3** *Sequencing (~ 4-9 days*)*

**Up to 96 samples**

Perform sequencing on forward strand
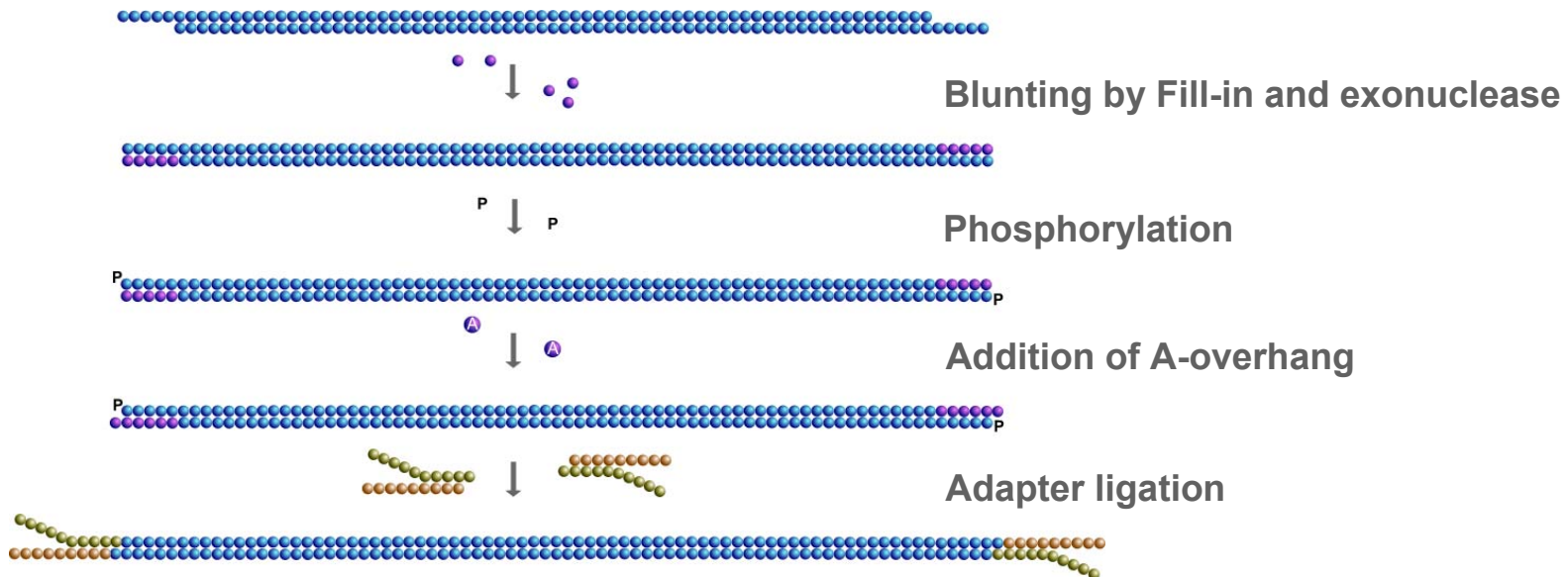
⬇

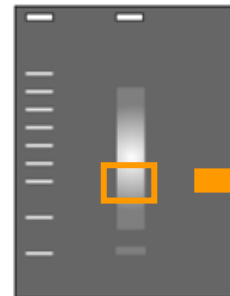Re-generate reverse strand

⬇

Perform sequencing on reverse strand

illumina®

# Flow cell



8 channels

Surface of flow cell coated with a lawn of oligo pairs

5'-PS-TTTTTTTTTAATGATACGGCGACCACCGAGAUCTACAC-3'

5'-PS-TTTTTTTTTCAAGCAGAAGACGGCATACGA

- Contained environment
- No need for clean rooms
- Sequencing performed inside the flow cell

illumina®

# Genomic DNA Library Prep

DNA fragments

**1**

Blunting by Fill-in and exonuclease

Phosphorylation

Addition of A-overhang

Adapter ligation

PCR
6 – 15 cycles

Library

illumina®

# Cluster Generation
## *Cluster station*

- Aspirates DNA samples into flow cell

- Automated amplified clonal clusters



>100M single molecules

>100M single clusters

**Flow cell (clamped into place)**

**DNA libraries**

# Cluster Generation
## *Hybridize Fragment & Extend*

- >150 M single molecules hybridize to the lawn of primers

- Bound molecules are then extended by polymerases

**Adapter sequence**

**3' extension**

illumina®

# Cluster Generation
## *Denature Double-stranded DNA*

- Double-stranded molecule is denatured

- Original template is washed away

- Newly synthesized covalently attached to the flow cell surface

**Newly synthesized strand**

**Original template**

**discard**

illumina®

# Cluster Generation
*Covalently-Bound Spatially Separated Single Molecules*

Single molecules bound to flow cell in a random pattern

illumina®

# Cluster Generation
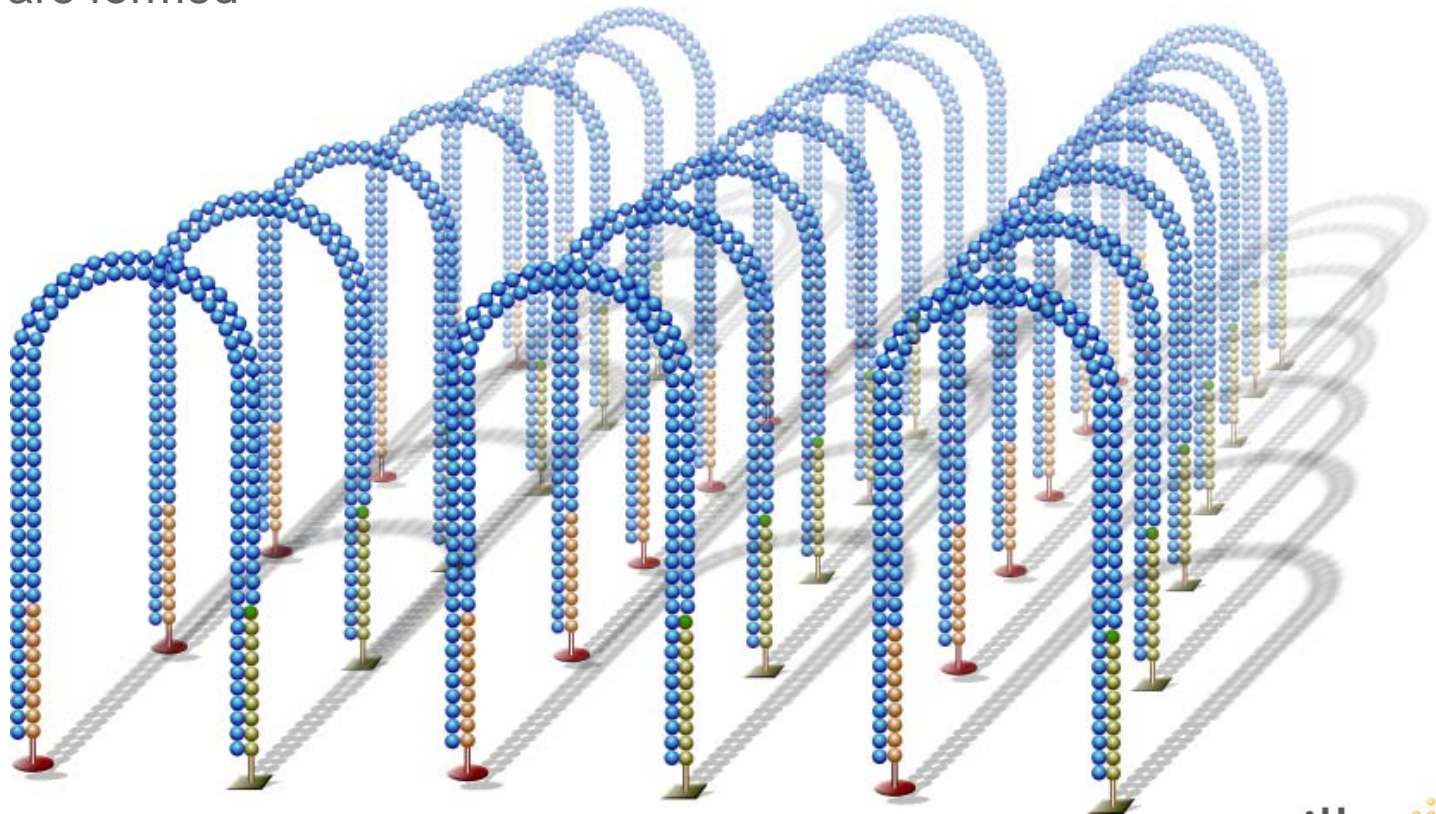## *Bridge Amplification*

- Single-strand flips over to hybridize to adjacent primers to form a bridge
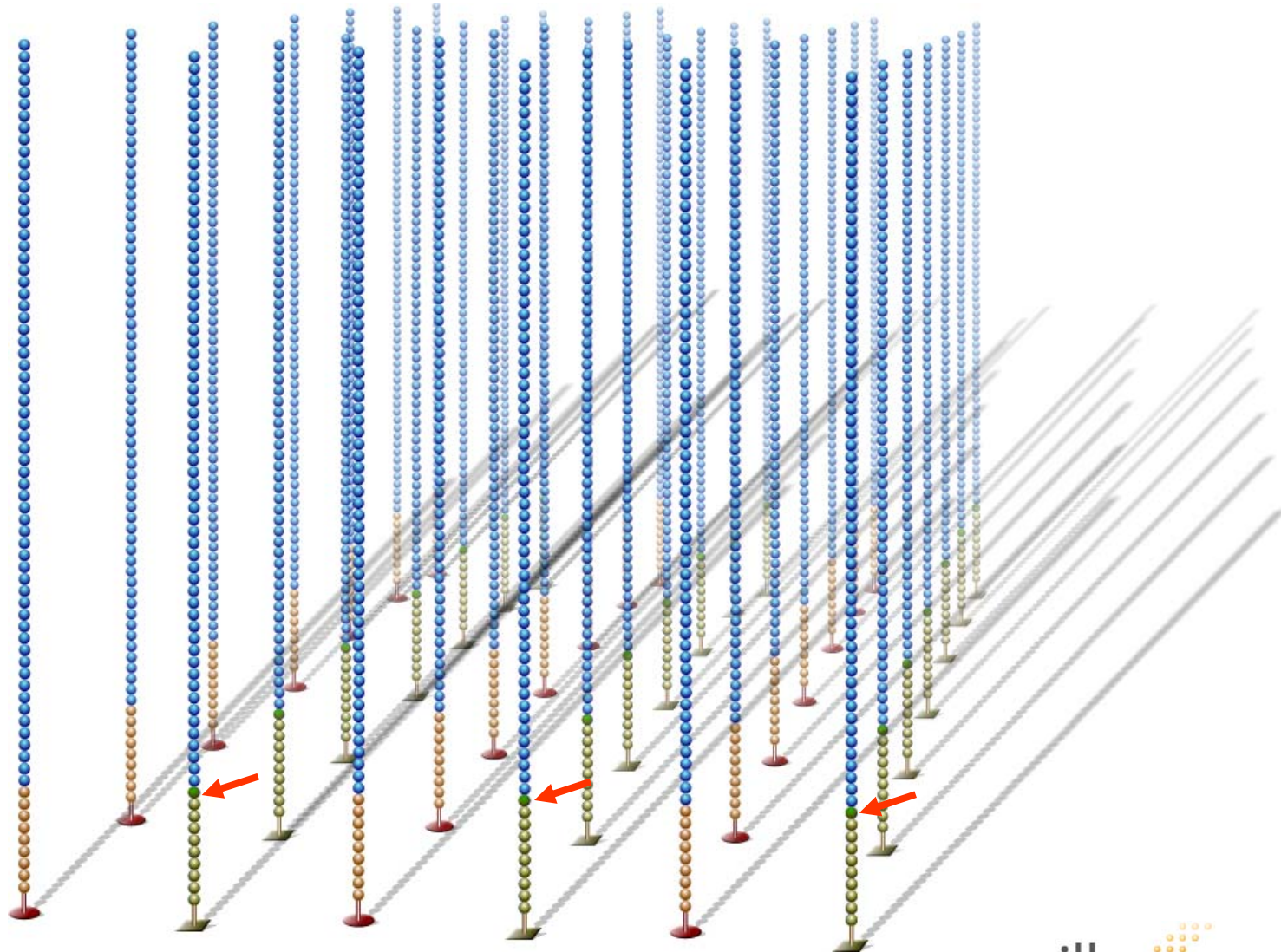
- Hybridized primer is extended by polymerases

illumina®

# Cluster Generation
## *Bridge Amplification*

- Double-stranded bridge is formed

illumina®

# Cluster Generation
## *Bridge Amplification*

- Double-stranded bridge is denatured

- Result: Two copies of covalently bound single-stranded templates

illumina®

# Cluster Generation
## *Bridge Amplification*

- Single-strands flip over to hybridize to adjacent primers to form bridges

- Hybridized primer is extended by polymerase

# Cluster Generation
## *Bridge Amplification*

- Bridge amplification cycle repeated until multiple bridges are formed

illumina®

# Cluster Generation

- dsDNA bridges denatured

- Reverse strands cleaved and washed away

# Cluster Generation

- Leaving a cluster with forward strands only

illumina®

# Genome Analyzer Sequencing reaction

- Sequencing primer is hybridized to adapter sequence

**Sequencing primer**

# Sequencing

**Add 4 Fl-NTP's + Polymerase** — **Incorporated Fl-NTP is imaged** — **Terminator and fluorescent dye are cleaved from the Fl-NTP**

X 36 - 75

illumina®

20

# Sequencing



~1000 copies per cluster

20 um

**>150 Million Clusters Per Flow Cell**

# Broadest range of applications
*Optimized, streamlined and easy-to-use reagent solutions*

**Sample Prep**

**Whole genome**
- Resequencing
- De-novo
- Targeted
- Metagenomics

**Transcriptome**
- RNA-Seq
- DGE
- Small RNA
- miRNA

**Regulation**
- Methylation
- ChIP-Seq

**Automated Cluster Generation**

**Sequencing**

**Epigenomics**

**Transcriptomics**

*de novo* **sequencing**

**Metagenomics**

**Genomics**

**…go where the biology takes you**

illumina®

# Structural Variation Analysis

**Epigenomics**

**Transcriptomics**

*de novo* **sequencing**

**Metagenomics**

**Genomics**

**…to go where the biology takes you**

illumina®

**Reference**

Deletion   Inversion!   Deletion

**Local *de novo* assembly**

| Long Inserts | Deletion? |
| Short Inserts | Inversions? |
| Middle | Normal? |

**Epigenomics**

**Transcriptomics**

*de novo* **sequencing**

**Metagenomics**

**Genomics**

illumina®

# Nature Publications

| Author | Number of Reads | Average Length |
|---|---|---|
| Shi, *et al.* | 389K | 97 bp |
| Dinsdale, *et al.* | 14.6M | 105 bp |
| Mou, *et al.* | 307K | 96 bp |
| Warnecke, *et al.* | 300K | 100 bp |
| Turnbaugh, *et al.* | 1.7M | 93 bp |

illumina®

## Nature Publications

| Paper | Number of Reads | Average Length |
|---|---|---|
| Metatranscriptomics | 389K | 97 bp |
| Biome Profiling | 14.6M | 105 bp |
| Ocean Bacteria | 307K | 96 bp |
| Termite Microbiota | 300K | 100 bp |
| Obesity Microbiome | 1.7M | 93 bp |

illumina®

# Nature Publications

| Reads per run: 150-300 Million | Number of Reads | % of a 150M read run |
|---|---|---|
| | 389K | 0.3% |
| | 14.6M | 9.7% |
| | 307K | 0.2% |
| | 300K | 0.2% |
| | 1.7M | 1.1% |

# The Value of Long, Overlapping Reads

**Epigenomics**

**Transcriptomics**

*de novo* **sequencing**

**Metagenomics**

**Genomics**

illumina®

"
When the read length exceeds a certain threshold, the **read length barrier**, the efficiency reaches nearly 100%, so that the read length indeed does not matter.
"

Chaisson et al., De novo *Fragment Assembly with Short Mate-Paired Reads: Does the Read Length Matter?, Genome Research,* October 22, 2008

illumina®

October, 2008

*Salmonella seftenberg*

| | Illumina |
|---|---|
| N50 contig size: | 139,353 |
| Largest contig: | 395,600 |
| Average contig: | 63,969 |
| Total bases of contigs: | 4.80Mb |
| Coverage of genome: | 99.8% |

**Beijing Genome Institute**

January 9, 2009

- 3GB genome
- Paired 75-base reads
- >95% gene regions
- N50 contig: 300Kb

**Epigenomics**

**Transcriptomics**

*de novo* **sequencing**

**Metagenomics**

**Genomics**

illumina®

**Capabilities**

Length of read

Short-insert Paired Ends

Raw Read Accuracy

**Forward Strand: 126kb**

Exon

Intron

Illumina: 100bp

GAGTGAAGCTCCTGGAGGAACTCAGATGGAGAAGTATCCAGTATGCATCTCGGGGAGAGAGACATTCAGCCTATAATGAATGGAAAAAGGCCCTCTTCAA

AGAACAAAGCACAAGAGTGAAGCTCCTGGAGGAACTCAGATGGAGAAGTATCCAGTATGCATCTCGGGGAGAGAGACATTCAGCCTATAATGAATGGAAAAAGGCCCTCTTCAAGCCTC

| Ensembl | Illumina GA |

**Epigenomics**

**Transcriptomics**

*de novo* sequencing

**Metagenomics**

**Genomics**

illumina®

AT4G25530.1

251 (CpG)

259 (CpG)

273 (CpG)

**% Cytosine methylation**

GATTTGTGGGATACTGAC
CTAAACACCCTATGACTG

illumina®

41

Lister, R., et al., (2008) Highly integrated single-base resolution maps of the epigenome in Arabidopsis. Cell. 133(3), 523-536.

# GA Applications Published 24 months from launch



**Cumulative original papers**

>250 total

| | |
|---|---|
| nature | 42 |
| Science | 14 |

Genome Sequencing, 18%

Transcriptome Analysis, 29%

Data Analysis, 30%

Sequencing Technologies, 2%

Protein-Nucleic Acids Interactions, 16%

Epigenomics, 5%

illumina®

# Genome Analyzer*IIx*
# Software Advancements
*Increased Output, Simplified Computing*

illumina®

# Increased output with reduced computing infrastructure
## *More gigabases of data for fewer gigabytes of computing power!*

## Sequencing Control Software v2.4

- Includes new Real Time Analysis (RTA) feature

- Image extraction and real time base calling on instrument computer

- Shorter time to results
  - Performed simultaneously with sequencing
  - Eliminates need to transfer images and intensities across network
  - Base calls and quality scores within hours of end of run

**Real-Time Analysis**

**Images**      **Intensities**      **Base calls**

illumina®

# New Software Delivers Up to 40% More Data Per Run

- Pipeline 1.4 – Enhanced analysis algorithm
  - Increases yield, improves accuracy
  - Improved cluster delineation
  - More clusters pass filter
  - Lower error rates



- SCS 2.4 – Integrated autofocus
  - Easier to use, removes user error
  - Less variability in focus quality

illumina®

# What is *Real Time Analysis*?

- The RTA module analyzes data as it leaves the Genome Analyzer
  - **Produces base calls, including Phred-like quality scores**
  - **Generates reports, to assess run and library quality**
  - Performs image analysis, generation of cluster intensities

- RTA simplifies the data management process
  - **Eliminates the need to transfer images from computer to computer**
  - Includes optional mechanisms for complete or selective archiving of images
  - Includes optional mechanisms for archiving of intensities

- RTA improves the system performance
  - **Minimizes time to results – base calls and qualities generated within hours of the end of the run**
  - Removes dependencies on network availability
  - Minimizes the time spent analyzing data after the run

illumina®

# Reversing the Trend
## *Simplified computing, Smaller storage needs, Faster analysis*

| Output | Time to Aligned Data | Computing Requirements | Output Files Size |
|--------|----------------------|------------------------|-------------------|
| | | 16Gb Server | 316Gb |
| 65% increase | Up to 30% reduction | 8Gb PC | 17Gb |

illumina®

# FireCrest module: Start From the Spots…



tiff image files

- Clusters identification and assigns intensities to them

- The output is a simple text file



intensity files

**Maximum**

**Threshold**

# From Intensities to Reads

## Intensity Files



## Sequence Files



Transforms intensities into base-calls

# From Reads to Aligned Sequences

sequence files



Gerald



Export.txt

# _export.txt – Flexible Output for Myriad Applications

# CASAVA

**Consensus Assessment of Sequence and Variation**

Aligned Reads



| Consensus assembly |
| --- |
| BINS<br>SNPs<br>Counts |

illumina®

# Consensus Sequence

# Output Interpretation - html

- Project directory/html/Home.html

# Overview of CASAVA Outputs for GenomeStudio™

- Output utilized for DNA experiments
  - Sorted.txt files binned by chromosome and subdivided into 10 megabase bins
  - SNP.txt (1 per chromosome)
  - Run_summary.xml
  - Run.conf file
  - Project.conf file

- Output utilized for RNA experiments
  - Sorted.txt files binned by chromosome and subdivided into 10 megabase bins
  - SNP.txt (1 per chromosome)
  - Exon counts file (1 per chromosome)
  - Gene counts file (1 per chromosome)
  - Splice Junction counts file (1 per chromosome)
  - Run_summary.xml
  - Run.conf file
  - Project.conf file

illumina®

# DNA Sequencing Module

- Direct import of data from Pipeline/CASAVA

- Visualization of SNPs

- Browsing of coverage and consensus reads

- Export of SNP tables

# Exon Counts from RNA Data

## Browser view of exon counts

# A New Look at Alternative Splicing

# Leveraging the GA Informatics Community
## *De Novo Assembly*

- **Velvet – De novo assembly of short reads**
  - Daniel Zerbino and Ewan Birney, EMBL-EBI
  - http://www.ebi.ac.uk/~zerbino/velvet/

- **SSAKE – Assembly of short reads**
  - Group: Rene Warren, et al; British Columbia
  - http://bioinformatics.oxfordjournals.org/cgi/content/full/23/4/500

- **Euler SR – Genomic Assembly**
  - Group: Pavel Pevzner, Mark Chaisson; UC San Diego
  - http://nbcr.sdsc.edu/euler/

# Rapidly Expanding Choice of Open Source Tools
*Genomic Alignment Browsers*

- ## Gbrowse – Genomic Browsing
  - Generic Model Organism Database Project
  - http://www.gmod.org/wlk/index.pho/Gbrowse

- ## UCSC Browser – Genome browsing and comprehensive annotation
  - Generic Model Organism Database Project
  - http://www.genome.ucsc.edu/goldenPath/help/customTrack.html

- ## Anno-J – Genome Annotation and Visualization
  - Computational Systems Biology Center of Excellence
  - http://www.annoj.org/csb_index.shtml

illumina®

# Leveraging the GA Informatics Community
*Alignment and Polymorphism Detection*

- **MAQ – Mapping and Assembly with Quality**
  - Heng Li, Sanger Centre
  - http://maq.sourceforge.net/maq-man.shtml

- **SOAP – Short Oligonucleotide Alignment Program**
  - Ruiqiang Li, Beijing Genomics Institute
  - http://soap.genomics.org.cn/

- **Consed – Alignement and Polymorphism Detection**
  - Green Lab, U. Washington (commercial offering)
  - http://bozeman.mbt.washington.edu/consed/consed.html

# Rapidly Expanding Choice of Open Source Tools

ChIP Sequencing

- ChIP-Seq Peak Finder
  - Barbara Wold, Cal Tech and Rick Meyers, Stanford University
  - http://woldlab.caltech.edu/html/software/

Digital Gene Expression

- Comparative Count Display
  - Alex Lash, NIH
  - ftp://ftp.ncbi.nlm.nih.gov/pub/sage/obsolete/bin/ccd/

- SAGE DGED Tool
  - Cancer Genome Anatomy Project
  - http://cgap.nci.nih.gov/SAGE/SDGED_Wizard?METHOD=SS10,LS108ORG=Hs

illumina®

# Genome Studio:
# Data Analysis Platform for Many Applications

**A single workbench with a growing number of modules**

- **Sequencing**
- **Genotyping (GT)**
- **Gene Expression (GX)**
- **Regulation (M)**
- **ChiP Sequencing (CS)**
- **We are changing the name to reflect a new platform →**
  - **microarray + sequencing**

# Connecting with the larger informatics universe

- Illumina has a 3rd party partnership program – illumina•Connect – launched last year designed to "increase software/hardware ecosystem connecting Illumina to researchers in genomics, genetics and sequencing communities"

  - ~30 vendors and academic partners in the program
  - http://www.illumina.com/pagesnrn.ilmn?ID=229
  - Integrated tools connecting GenomeStudio to several 3rd party apps for microarray and sequencing data analysis

# Delivering on Roadmap Milestones



**15x increase in 2008**

**4-5x increase in 2009**

*2x150 PE*

95G

*2x125 PE*

*Pipeline 1.4*
*SCS 2.4*
*GAIIx*

*2x100 PE*

*SBS Version 3*
*Pipeline 1.3*

55G

35G

25G

15G

Gb per run

100
90
80
70
60
50
40
30
20
10
0

Jan  Feb  Mar  Apr  May  Jun  Jul ----------→ Dec

**2009**

illumina®

# Sequencing Enabled iScan

What is needed?
How does it work?
What can it do?

# iScan Sequencing Module technical Update

- Add on to enable sequencing on the iScan
  - Uses lasers and optics of iScan for imaging of flow cell
  - Fluidics module holding sequencing reagents, pumps and reagent chiller

- Sequencing specs
  - Throughput: ~ 0.5 GB per day, up to 5 GB per run
  - Data density: ~ 32 M clusters per run, 60M Reads
  - Flexible read length: from short single reads to 2x75 bb reads
  - 8 lanes, 1-12 samples per lane

- Applications supported
  - smallRNA, CHiP-Seq & mRNA seq
  - Targeted re-sequencing & re-sequencing of small genomes
  - High Density GWAS, Medium and HD Custom GTP (Array Applications)
  - WG Gene Expression Profiling, Human, Mouse and Rat (Array Applications)

# Thank You

Marco Cappelletti

mcappelletti@illumina.com