



Lars Paulin

New DNA sequencing
technologies

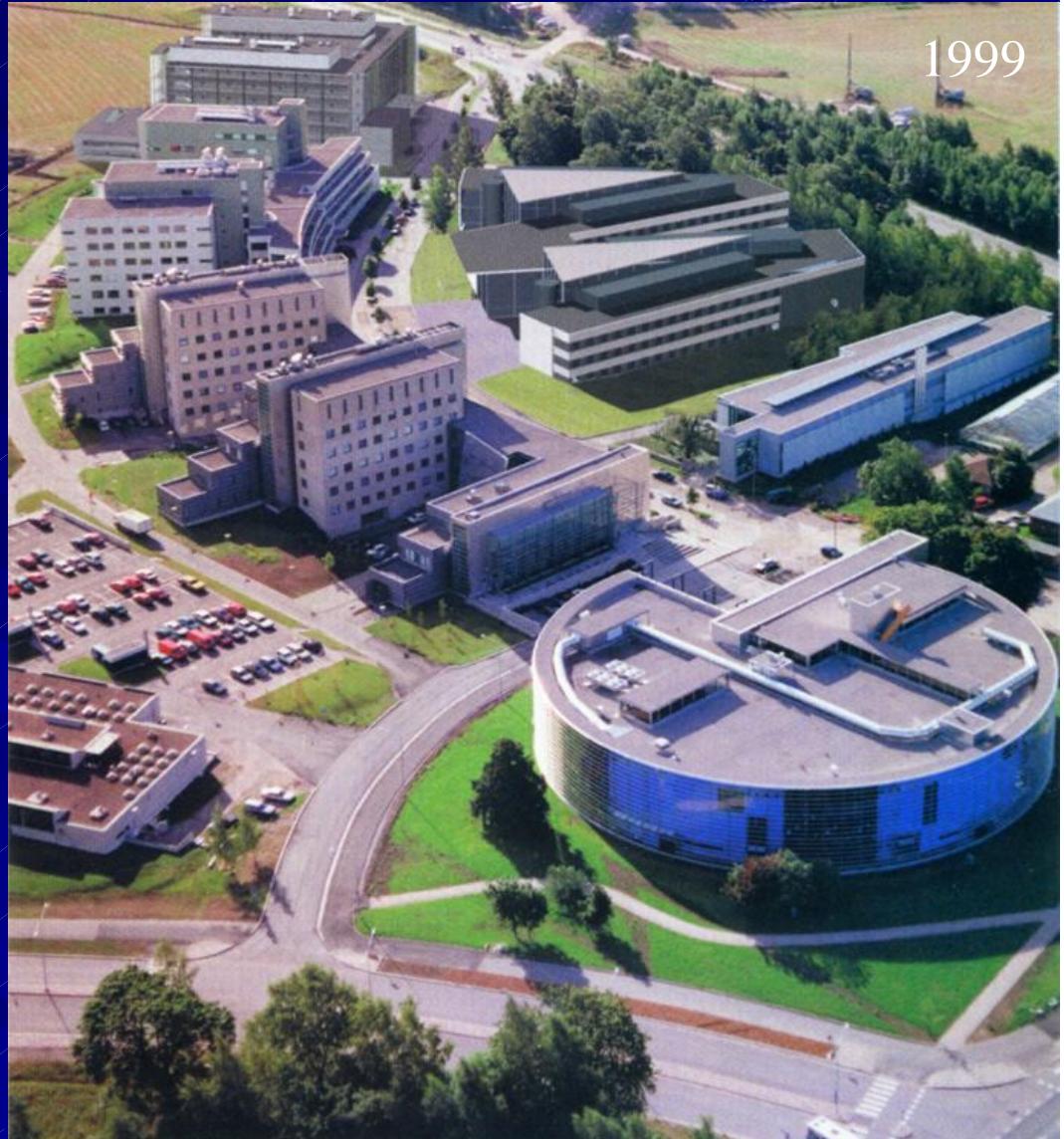
DNA Sequencing and
Genomics

Institute of Biotechnology
University of Helsinki

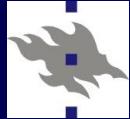
<http://www.biocenter.helsinki.fi/bi/DNA/>

Viikki Science Park

1999



Lars Paulin Institute of Biotechnology University of Helsinki

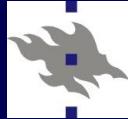


DNA Sequencing Laboratory

Cultivator 2, Viikinkaari 4

- Started in 1990 with DNA Synthesis
- 1991 DNA Sequencing
- 1994 EU Yeast Genome Project
- 1999 - 2000 High-throughput pipeline
- 1999 – 2002 Five EST Sequencing Projects
- 2003 First Microbe Genome Project
 - Move together with Microarray Laboratory to Cultivator 2
- 2006 Genome Sequencer 20

- Core Facility
 - Service DNA sequencing and whole projects
 - Collaborative projects
 - "Research hotel"
 - Develop high-throughput methods
 - Consulting



DNA Sequencing and Microarray Pipe-Line

■ Bacterial growth

- QPix Colony Picker
- QFill 2 Microplate filler
- 4 Multidrop 96 and 384 well
- 4 Titramax shakers

■ PCR

- 1 Tetrad MjResearch [1 x 96, 3 x (2x48)]
- 2 Tetrad MjResearch (4 x 96)
- 1 Eppendorf Mastercycler (384)
- 1 ABI 9700 (2 x 384)
- 4 ABI 9700 (96)

■ Robotics

- Qiagen Biorobot 9600
- Qiagen Biorobot 8000
- Qiagen Biorobot 3000 (2m)
- 2 Tecan Genesis RSP100
 - low volume pipetting
- Tecan Freedom Evo
 - low volume pipetting
- Biomek NX^P, 96-pipettor

■ Sequencers

- ABI 3130 16-Capillary Sequencer
- ABI 3730 48-Capillary Sequencer
- Genome Sequencer 20

■ Centrifuges

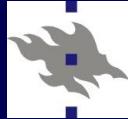
- 2 Eppendorf 5810
- 1 Eppendorf 5804
- 1 Beckman

■ Other

- Tecan SpectraFluor Microplate Reader
- ImageMaster VDS, GE Healthcare
- 8 Servers for analysis

■ Phred, Staden Program, AceDB, ARB etc.

- SUN cluster
- ABI Prism 7000
- LightCycler 480, Roche
- Coulter Counter Z2
- Bioanalyzer
- LIMS, Sapphire



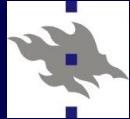
What can we do?

All kinds of service

- Sequencing projects
- Fragment analysis
 - T-RFLP
- SNP detection
- Colony picking
- Plasmid, cosmid, fosmid, BAC purification
- PCR, PCR purification
- rDNA sequencing
- EST, SAGE sequencing
- Whole genome sequencing
- Metagenomic sequencing
- MLST-Sequencing
- etc.

High-throughput

- Picking 3 500 colonies / h
 - petri discs, 22 x 22 cm plate
- Growth ~50 x 96 / day
 - 70 µl - 1 ml
- PCR ~5 000 / day
 - 50 - 100 µl
- PCR purification 12 x 96 / day
 - 2 x 100 µl
- Plasmid, Cosmid, Fosmid, BAC purification
 - 96 / 384 well
- Sequencing 0,5 - 1 Mb / 24h
- Gridding 12 membranes / day
 - 1536 spots in duplicate, 384 controls
- Genome Sequencer FLX



EST Projects 1999 – 2002

Barley

- 44 829 ESTs
- 13 448 clusters
 - 7 062 cluster 1



Gerbera

Genome Res 2005, 15,475-86

- 16 994 ESTs
- 11 266 clusters
 - 8 817 cluster 1



Poplar

Genome Biol, 2005, 6, r101

- 13 845 ESTs
- 8 043 clusters
 - 5 860 cluster 1

Birch

- Wood, Cold, Oxidative
stress, Pathogen
- 73 881 ESTs
 - 20 103 clusters
 - 10 554 cluster 1



Spruce

Plant Mol Biol 2007, 65, 311-328.

- 8 432 ESTs
- 5 817 clusters
 - 4 343 cluster 1

Total

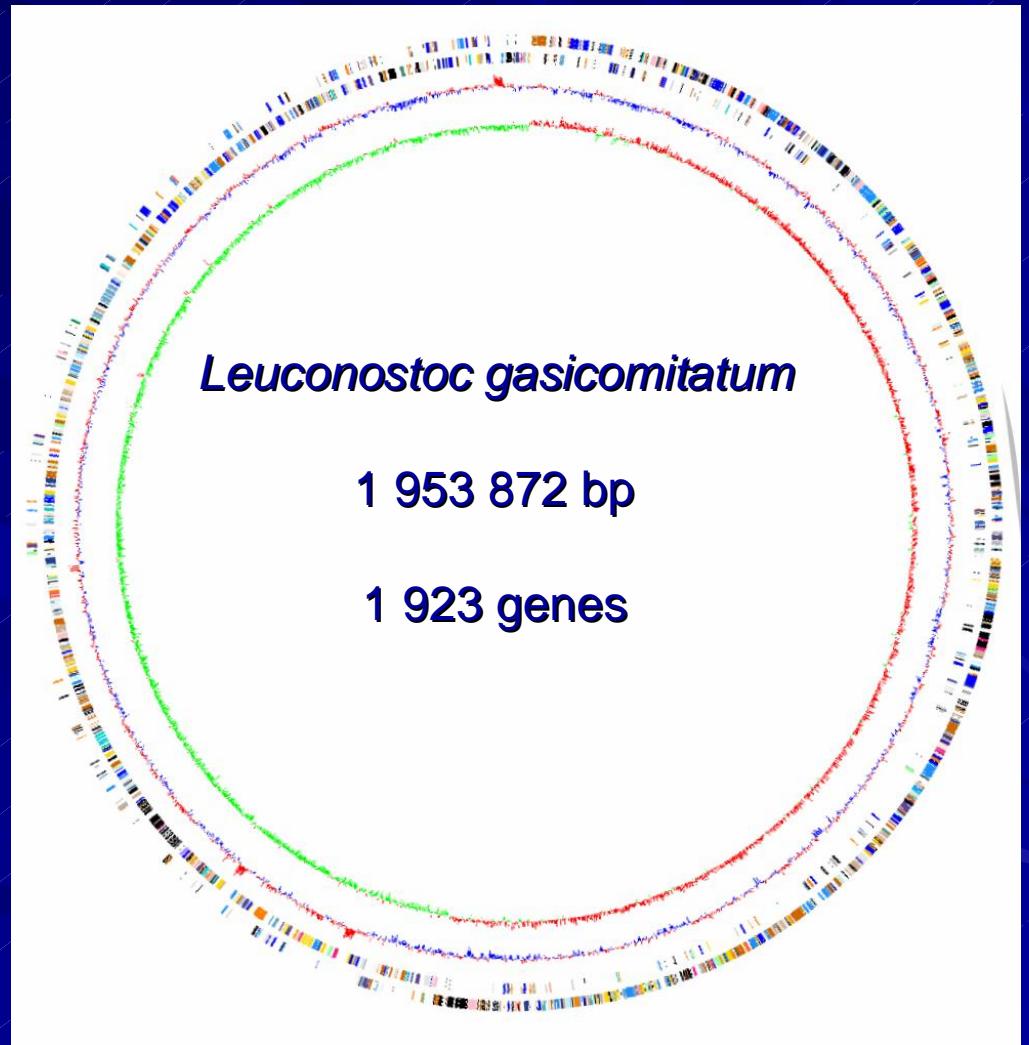
- 157 981 ESTs in Databases
- ~260 000 sequencing reactions

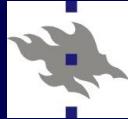


Leuconostoc gasicomitatum genome project

Paulin, Björkroth and Auvinen

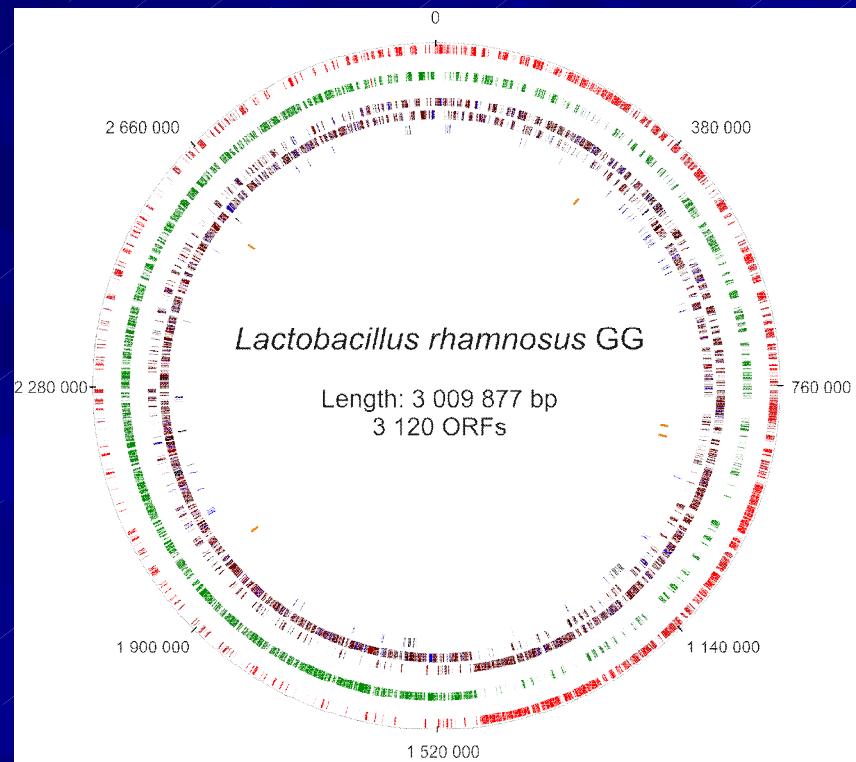
- First described by Johanna Björkroth
 - Appl Environ Microbiol. 2000, 66, 3764-72.
- Food spoiling bacteria
 - can grow at + 6 °C
 - marinaded poultry products
 - fish and meat

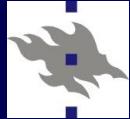




Lactobacillus rhamnosus GG

- Widely used probioite
- Biggest *Lactobacillus* genome so far
 - 3 009 877 bp
 - 3 120 ORFs
 - 2 129 ORFs with assigned functions
 - 991 ORFs without assigned functions
 - 5 rRNA operons
 - G+C % 46,69
 - Annotation done at ERGO
<http://ergo.integratedgenomics.com/ERGO/>
 - Agilent microarrays done
 - The sequence announced at the conference:
SSA GutImpact 3rd Platform meeting
on Foods for Intestinal Health 29.-
31.8.2007, Haikko Manor & Spa,
Finland





Genome Projects

L. gasicomitatum, 1,9 Mb

- Shotgun, fosmids
- Status: finished, annotated

C. pneumoniae, 1,3 Mb

- PCR products
- Status: almost finished

L. rhamnosus, 3 Mb

- Shotgun, fosmids, 454
- 2 Strains, GG and Lc
- Status: 2 finished, 1 annotated

L. oligofermentans, 1,8 Mb

- 454 , fosmids
- Status: closing phase

C. botulinum, 3,8 Mb

- 454, fosmids
- Status: closing phase

H. bizzozeronii, 1,7 Mb

- 454, plasmids
- Status: closing phase

H. salomonis, 1,7 Mb

- 454, plasmids
- Status: closing phase

E. carotovora, 5 Mb

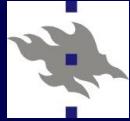
- 454, fosmids
- Status: closing phase

S. islandicus, 2,7 Mb

- shotgun, 454
- Status: closing phase

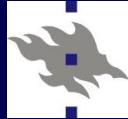
A. brierleyi, 2,7 Mb

- Shotgun, 454
- Status: closing phase



Short History of DNA Sequencing

- 1977
 - Maxam-Gilbert
 - Sanger
- 1986
 - First Automated DNA Sequencer ABI 370 (373)
- 1988
 - Pharmacia ALF
- 1995
 - ABI 377
 - Up to 96 lanes
- 1996
 - First Capillary DNA Sequencer ABI 310
- 1998
 - First 96 Capillary instruments MegaBace, ABI 3700
- 2000
 - ABI 3100, 16 Capillary
- 2002
 - ABI 3730, 48 or 96 Capillary
- 2005
 - Genome Sequencer GS20
- 2006
 - Solexa (Illumina)
- 2007
 - SOLiD



Sanger DNA Sequencing

1. Template

- ssDNA or dsDNA

2. Primer annealing

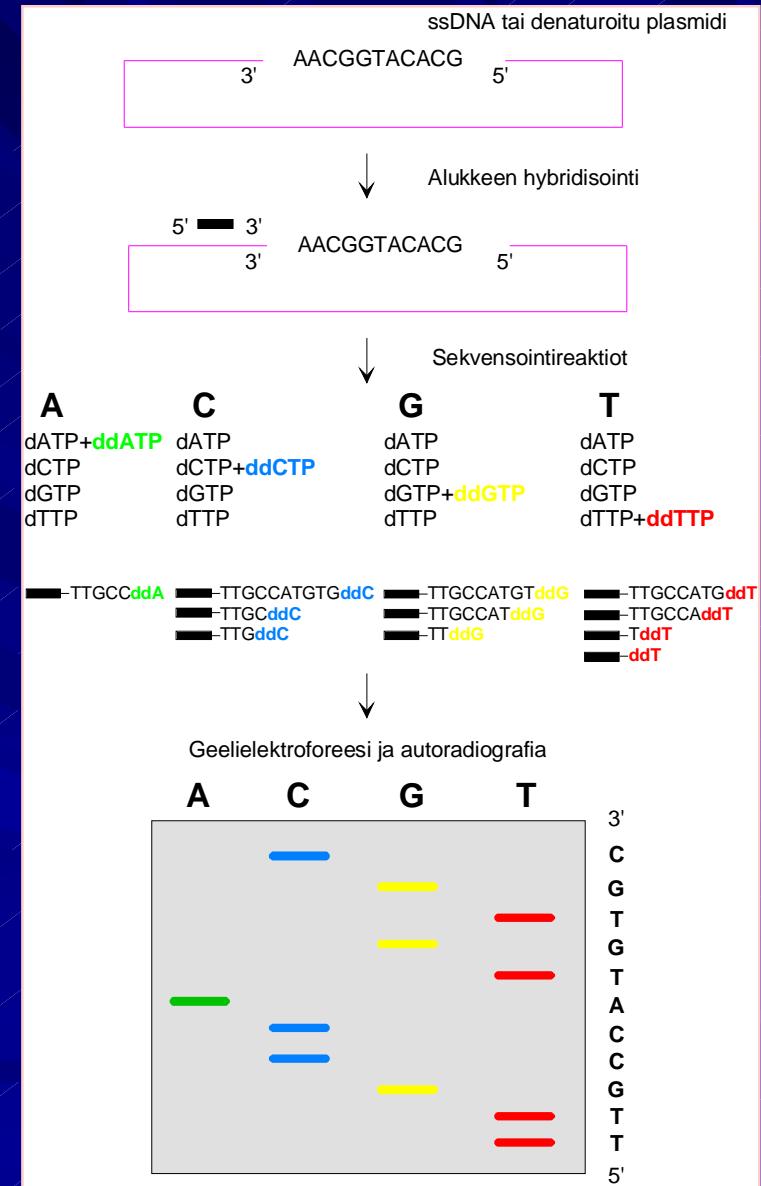
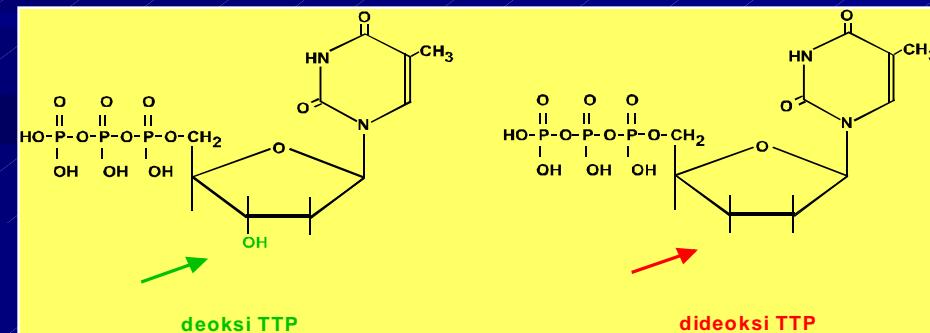
- Sequencing primer

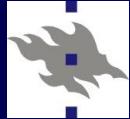
3. Elongation

- DNA polymerase

■ Steps 2 and 3 can be done repeatedly => cycle sequencing

4. Electrophoresis



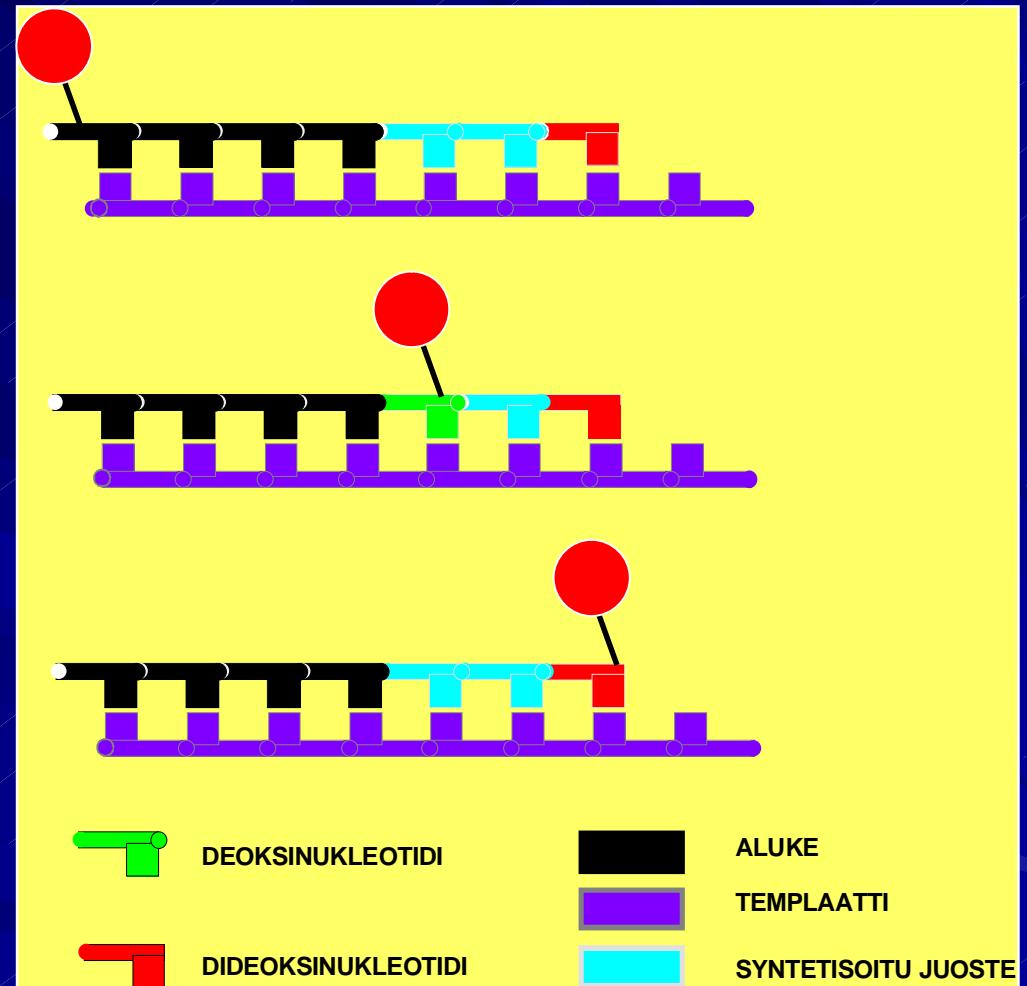


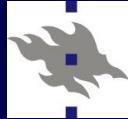
Incorporating Labels

Labelled primers
•1 or 4 labels

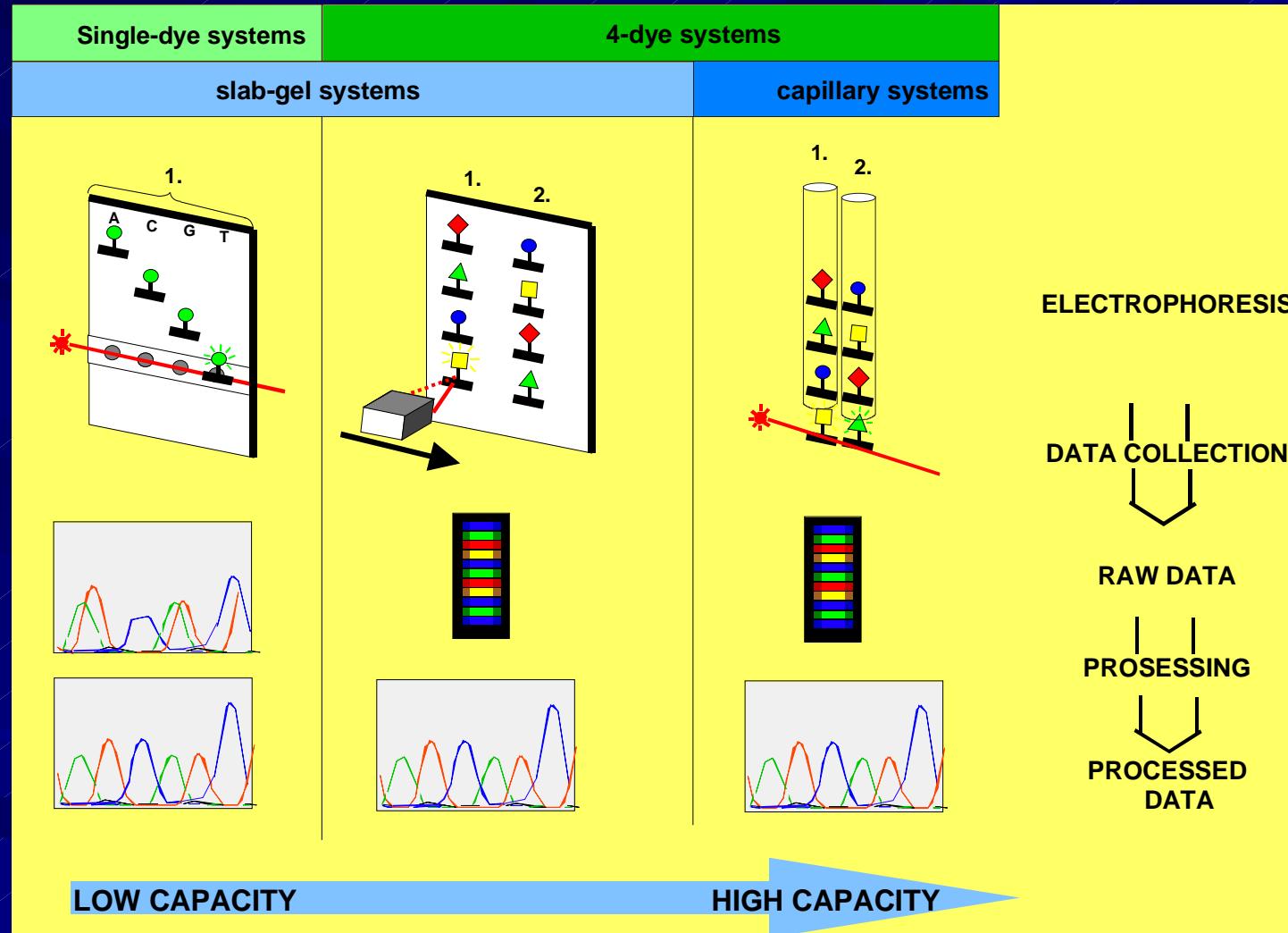
Labelled deoxynucleotides
•1 label

Labelled
dideoxynucleotides
•1 or 4 labels
•BigDye, ET
terminators





Automated DNA Sequencing





Model 3700

d00462_A05_Tas6up_033.ab1

Signal G:172 A:243 T:196 C:173

Page 1 of 2

Version 3.6

DT3700POP5(BD)v3.mob

Tue, Sep 12, 2000 2:37 PM

Basecaller-POP5opt.bcpTas6up

elru

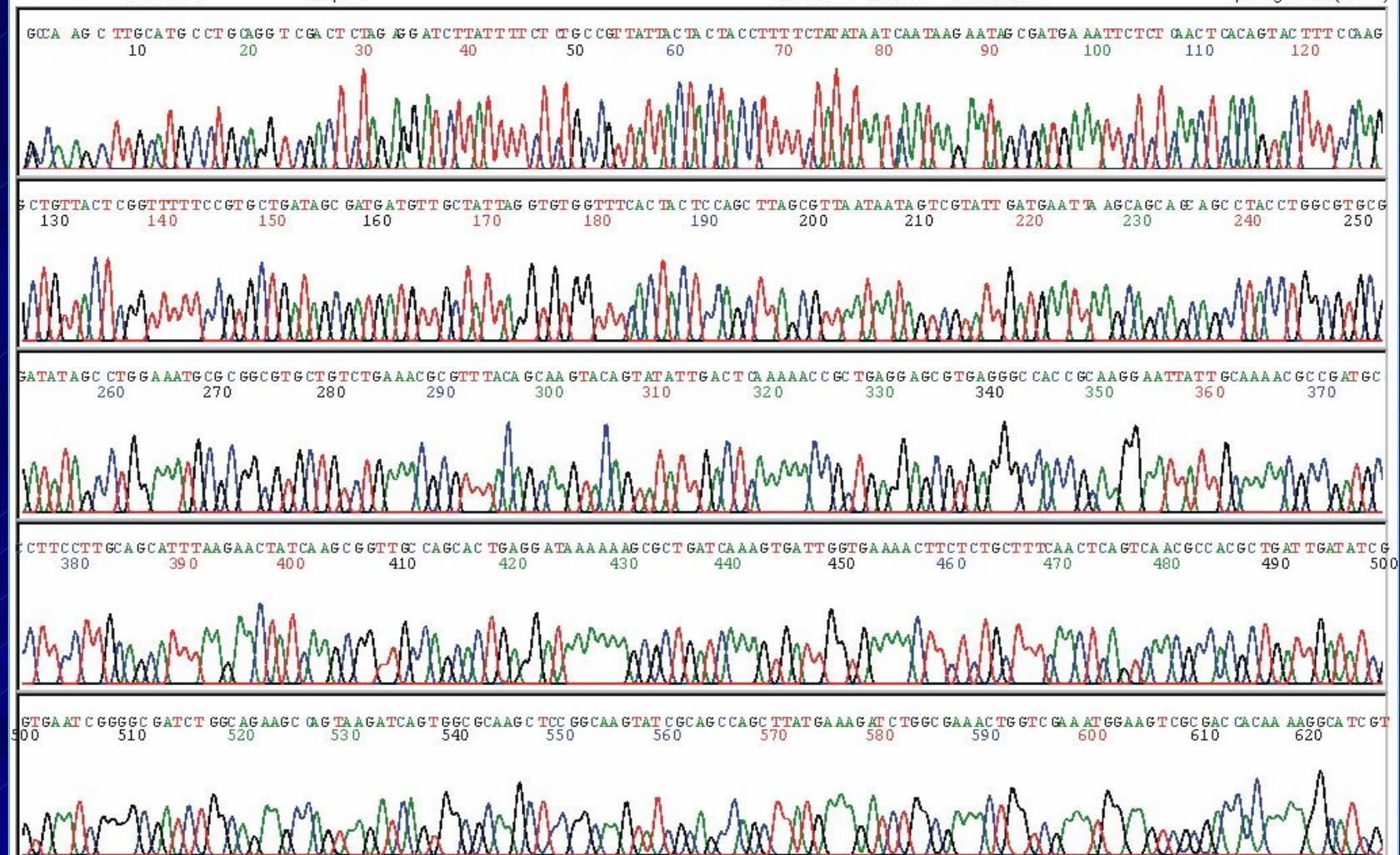
Tue, Sep 12, 2000 1:21 AM

BC 1.1.b.2

Cap 33

Points 2767 to 13845 Pk 1 Loc: 2767

Spacing: 15.52(15.52)





Strategies for Genome Sequencing

■ Shotgun approach

- random sequencing of different sized libraries
- assembly using different software
- closing of gaps using different methods

■ Libraries

- usually made by random shearing of genomic DNA
- 2 kb, 4-6 kb, 10 kb plasmid libraries
- fosmid or cosmid libraries with 30 - 50 kb inserts



Whole Genome Shotgun Sequencing



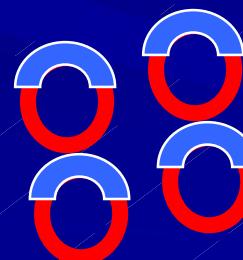
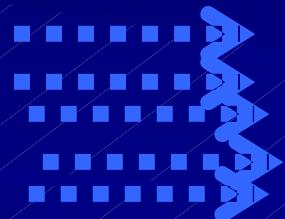
Whole Genome:
~ 3 Mb



Sheared DNA:
~ 2 kb



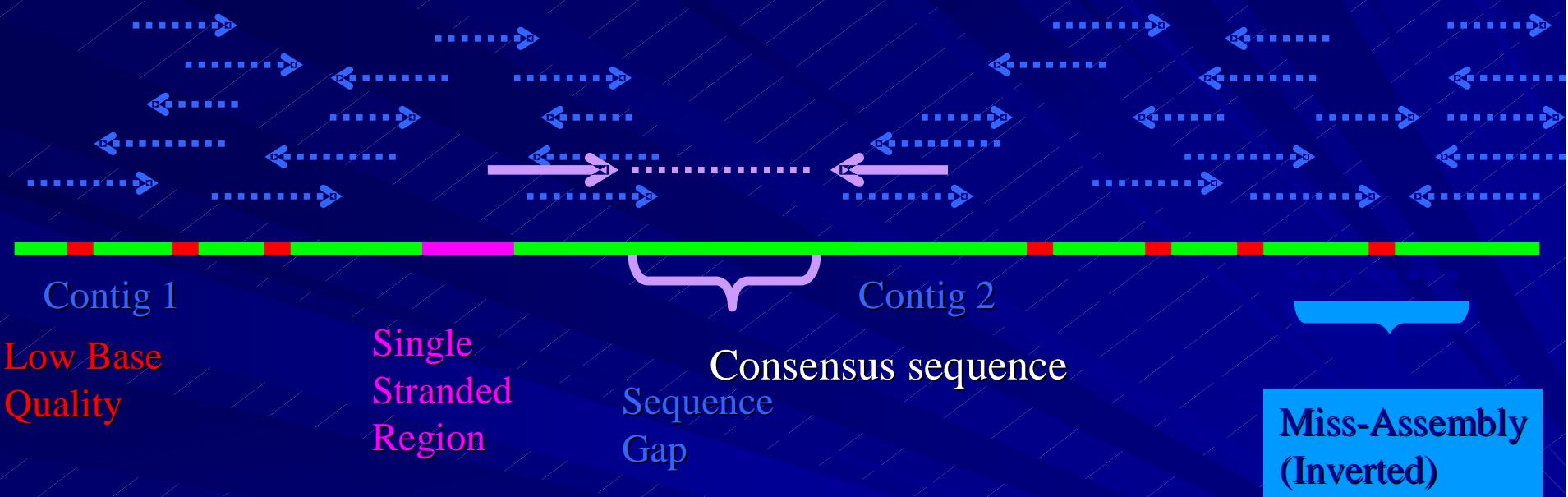
Random
Reads
Both ends



Sequencing
Templates

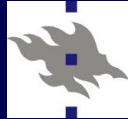


Shotgun Sequencing :ASSEMBLY



- 0.5 -1.0 X (2 reads/kb) - ‘Skimming’
- 3.5 - 4.0 X (~9 reads/kb) - ‘half-shotgun’

- 6.5 - 8.0 X (~18 reads/kb) - ‘pre-finished’
- 10 X (22-24 reads/kb) - ‘deep shotgun’



Phred, Phrap and Staden Package Program

Phred and Phrap

- University of Washington
- Phil Green, <http://www.phrap.org/>

Phred quality score:

$$QV = - 10 * \log_{10}(P_e)$$

where P_e is the probability that the base call is an error.

Phred score	P_e	Accuracy of the base call
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%

Staden Program

- Cambridge, Sanger Center
- Roger Staden,
<http://staden.sourceforge.net/>

■ Trace editing

■ Phrap assembly and Gap4 editing

- display of traces from sequencers
- translations, orfs, RE etc.
- good capacity

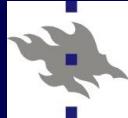


New DNA Sequencing Technology

Parallel Sequencing Technology

- Massive throughput
- Fast sequencing
- No cloning step
- PCR

- Currently three systems ready
 - Genome Sequencer (<http://www.454.com/>,<http://www.roche.com>)
 - 454 Life Sciences, Roche
 - Launched in October 2005
 - Solexa (<http://www.illumina.com>)
 - Illumina
 - Launched 2006
 - SOLiD (<http://www.appliedbiosystems.com>)
 - Applied Biosystems
 - Launched in October 2007



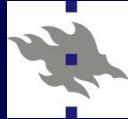
ARTICLES

Genome sequencing in microfabricated high-density picolitre reactors

Marcel Margulies^{1*}, Michael Egholm^{1*}, William E. Altman¹, Said Attiya¹, Joel S. Bader¹, Lisa A. Bemben¹, Jan Berka¹, Michael S. Braverman¹, Yi-Ju Chen¹, Zhoutao Chen¹, Scott B. Dewell¹, Lei Du¹, Joseph M. Fierro¹, Xavier V. Gomes¹, Brian C. Godwin¹, Wen He¹, Scott Helgesen¹, Chun He Ho¹, Gerard P. Irzyk¹, Szilveszter C. Jando¹, Maria L. I. Alenquer¹, Thomas P. Jarvie¹, Kshama B. Jirage¹, Jong-Bum Kim¹, James R. Knight¹, Janna R. Lanza¹, John H. Leamon¹, Steven M. Lefkowitz¹, Ming Lei¹, Jing Li¹, Kenton L. Lohman¹, Hong Lu¹, Vinod B. Makhijani¹, Keith E. McDade¹, Michael P. McKenna¹, Eugene W. Myers², Elizabeth Nickerson¹, John R. Nobile¹, Ramona Plant¹, Bernard P. Puc¹, Michael T. Ronan¹, George T. Roth¹, Gary J. Sarkis¹, Jan Fredrik Simons¹, John W. Simpson¹, Maithreyan Srinivasan¹, Karrie R. Tartaro¹, Alexander Tomasz³, Kari A. Vogt¹, Greg A. Volkmer¹, Shally H. Wang¹, Yong Wang¹, Michael P. Weiner⁴, Pengguang Yu¹, Richard F. Begley¹ & Jonathan M. Rothberg¹

¹454 Life Sciences Corp., 20 Commercial Street, Branford, Connecticut 06405, USA. ²University of California, Berkeley, California 94720, USA. ³Laboratory of Microbiology, The Rockefeller University, New York, New York 10021, USA. ⁴The Rothberg Institute for Childhood Diseases, 530 Whitfield Street, Guilford, Connecticut 06437, USA.

*These authors contributed equally to this work.



Genome Sequencer

(<http://www.454.com/>, <http://www.roche.com>)

■ Genome Sequencer GS20;FLX

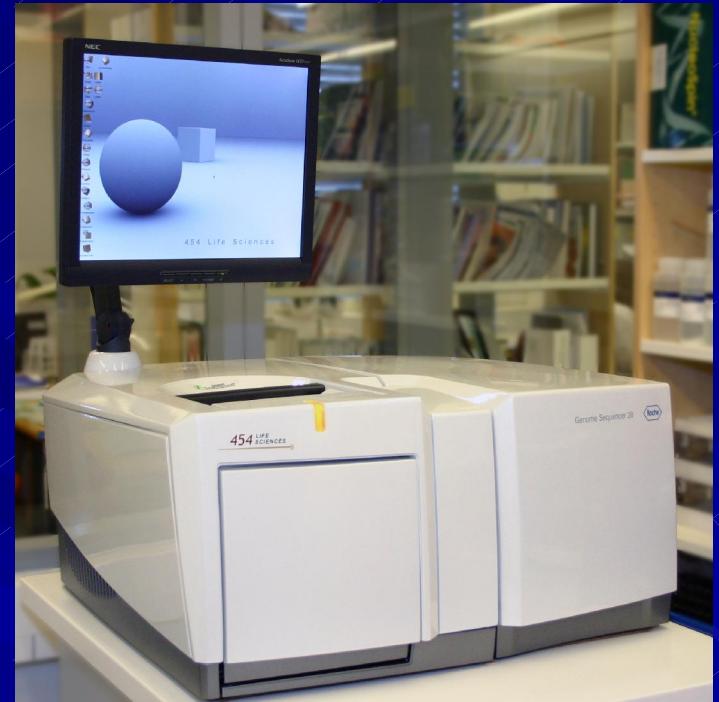
- Manufacturer 454 Life Science
- Marketing Roche

■ Parallel Sequencing

- Shotgun sequencing
 - No plasmid libraries
 - Linkers ligated to fragments
 - Emulsion PCR
 - Picotiter plate, 1 600 000 wells
- Pyrosequencing

(Nyren, P. et al Anal Biochem. 1993, 208,171-5)

- Detection with sensitive CCD camera
- Run time ca. 4,5 h; 7,5 h
- Read lenght 100 -120 bp; 250 – 300 bp
- Raw sequence ca. 25 – 35 Mb/run; 80 – 100 Mb/run





Genome Sequencer GS 20/FLX

DNA Library Preparation

emPCR

Sequencing

emPCR

Sequencing

DNA Library Preparation

1. DNA fragmentation
2. Fragment end polishing
3. Adaptor ligation
4. Library immobilisation
5. Fill-in reaction
6. Single-stranded template DNA (ssDNA) library isolation
7. ssDNA library quality assessment and quantitation

Emulsion PCR Amplification

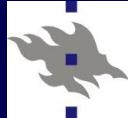
1. Preparation of the live amplification mix
2. ssDNA library capture
3. Emulsification
4. Amplification
5. Bead recovery
6. ssDNA library bead enrichment
7. Sequencing primer annealing

Sequencing/ Genome Sequencer 20 Operation

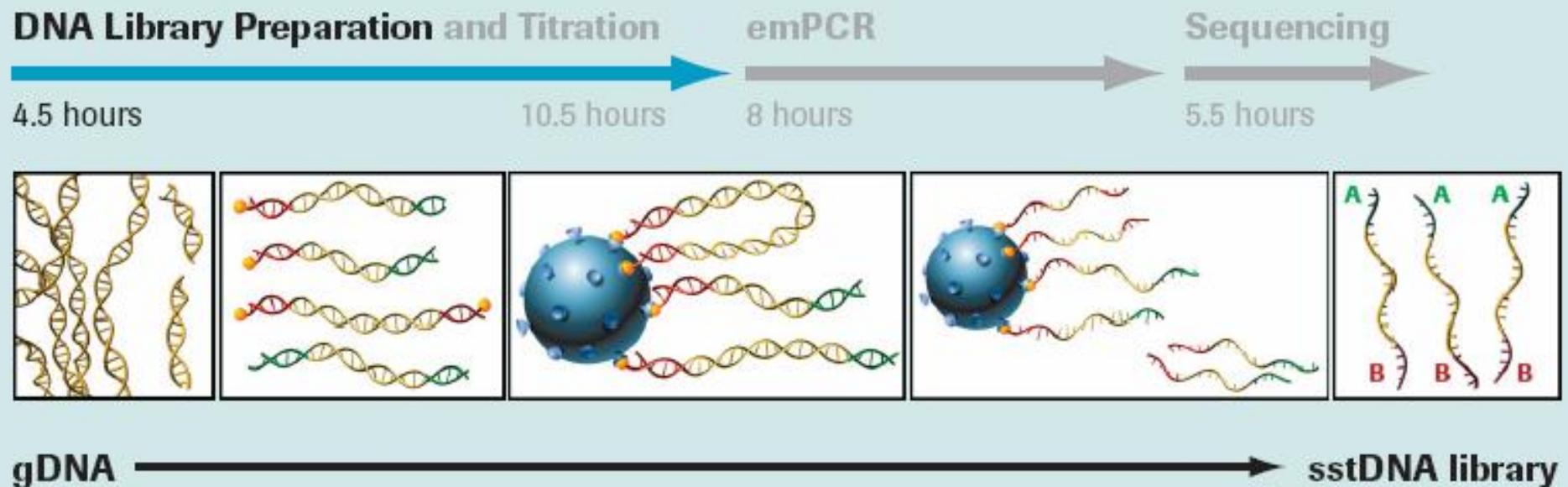
1. The pre-wash Run
2. PicoTiterPlate™ preparation
3. The PREP Run
4. The Sequencing Run

Output

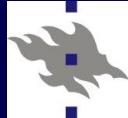
1. FASTA file
2. Assembly
3. Mapping



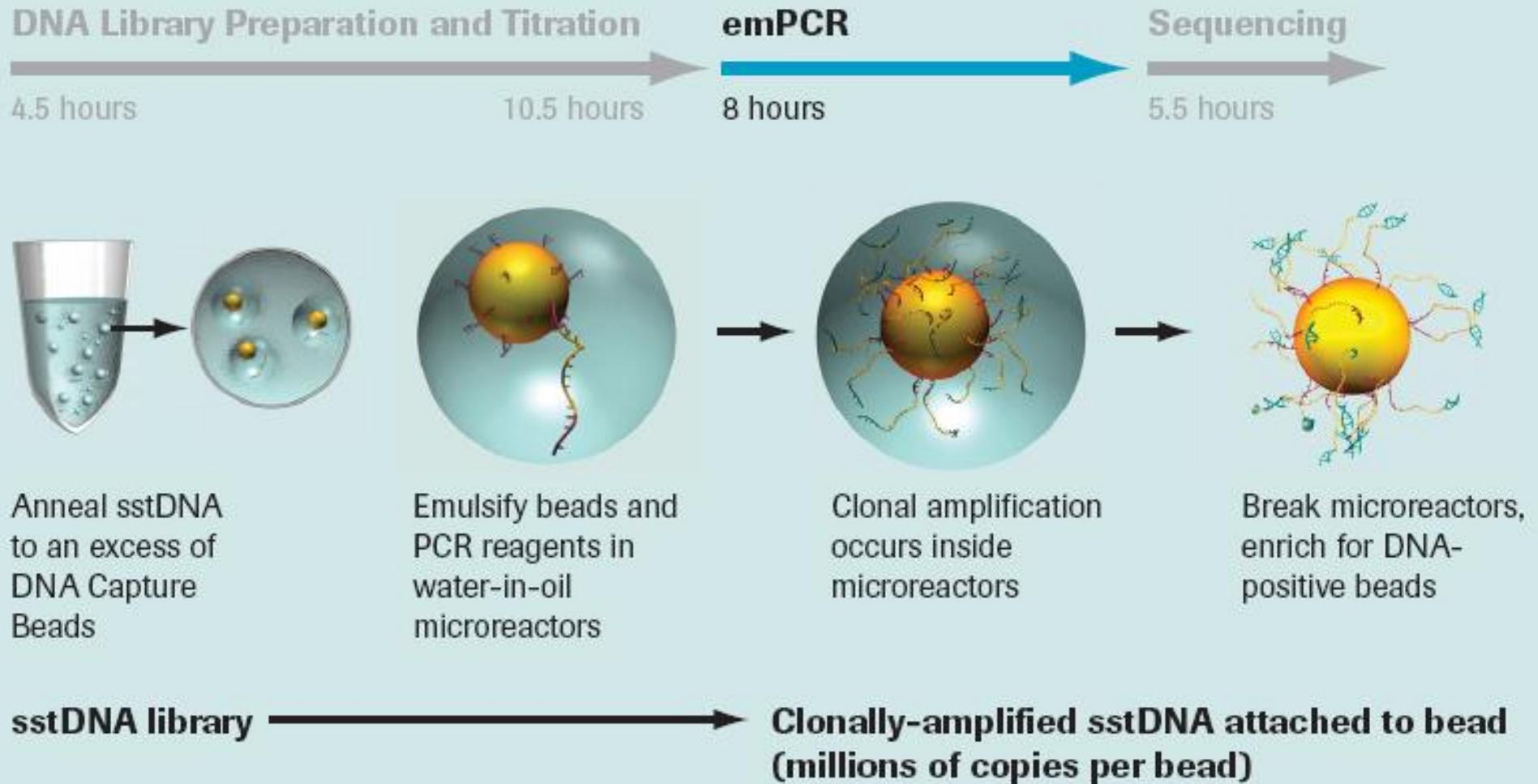
Library preparation

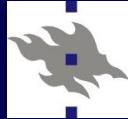


- Genome fragmented by nebulization
- No cloning; no colony picking
- sstDNA library created with adaptors. The adaptors are used as primers, and for binding to beads.
- A/B fragments selected using streptavidin-biotin purification

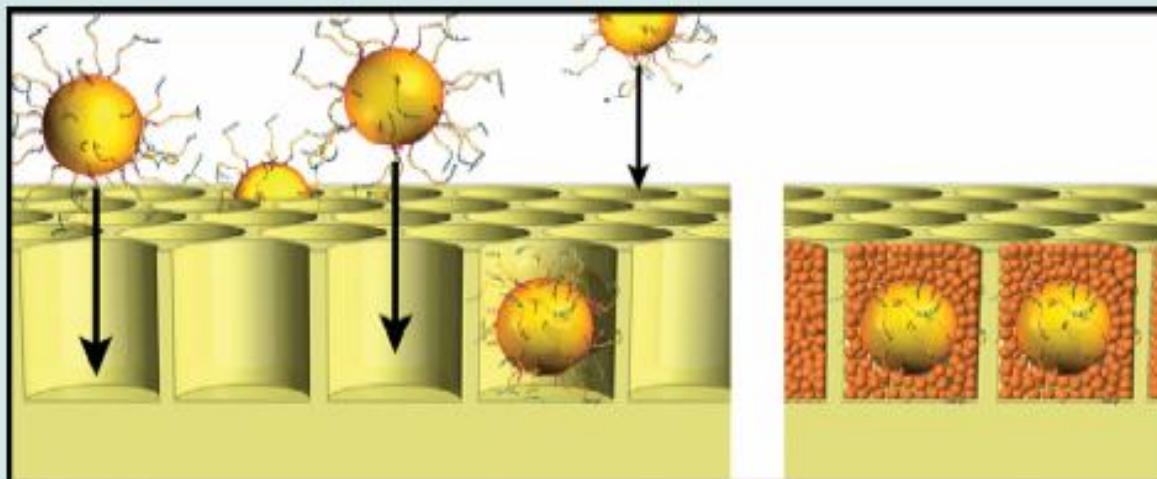
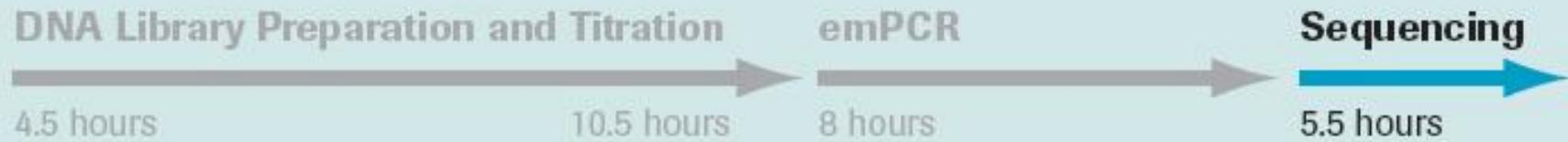


Emulsion PCR



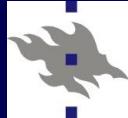


PicoTiterPlate (PTP)

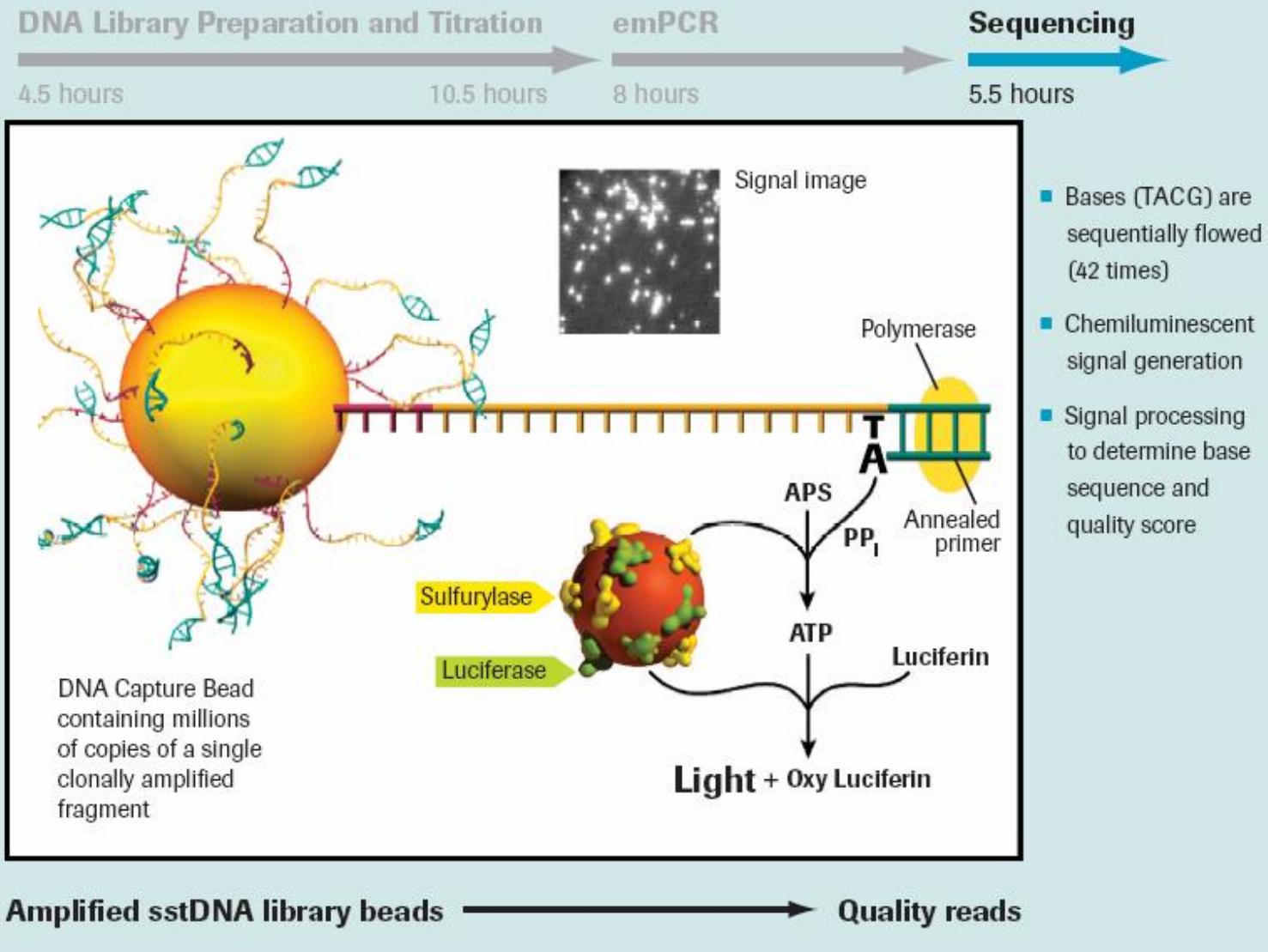


- Well diameter: average of 44 µm
- A single clonally amplified sstDNA bead is deposited per well
- 200,000 reads obtained in parallel on large-format PicoTiterPlate device

Amplified sstDNA library beads → **Quality reads**

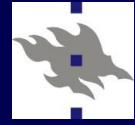


Pyrosequencing

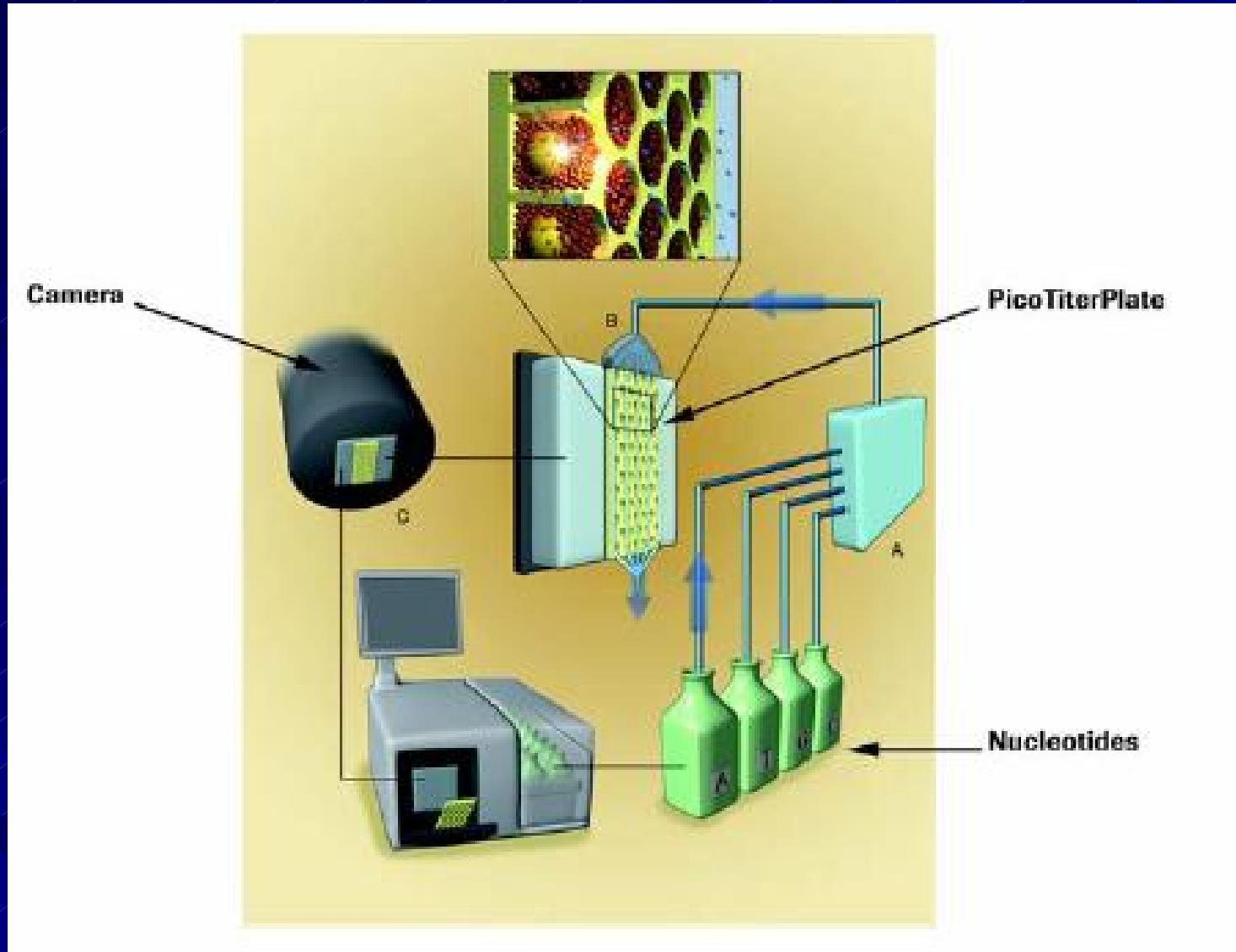


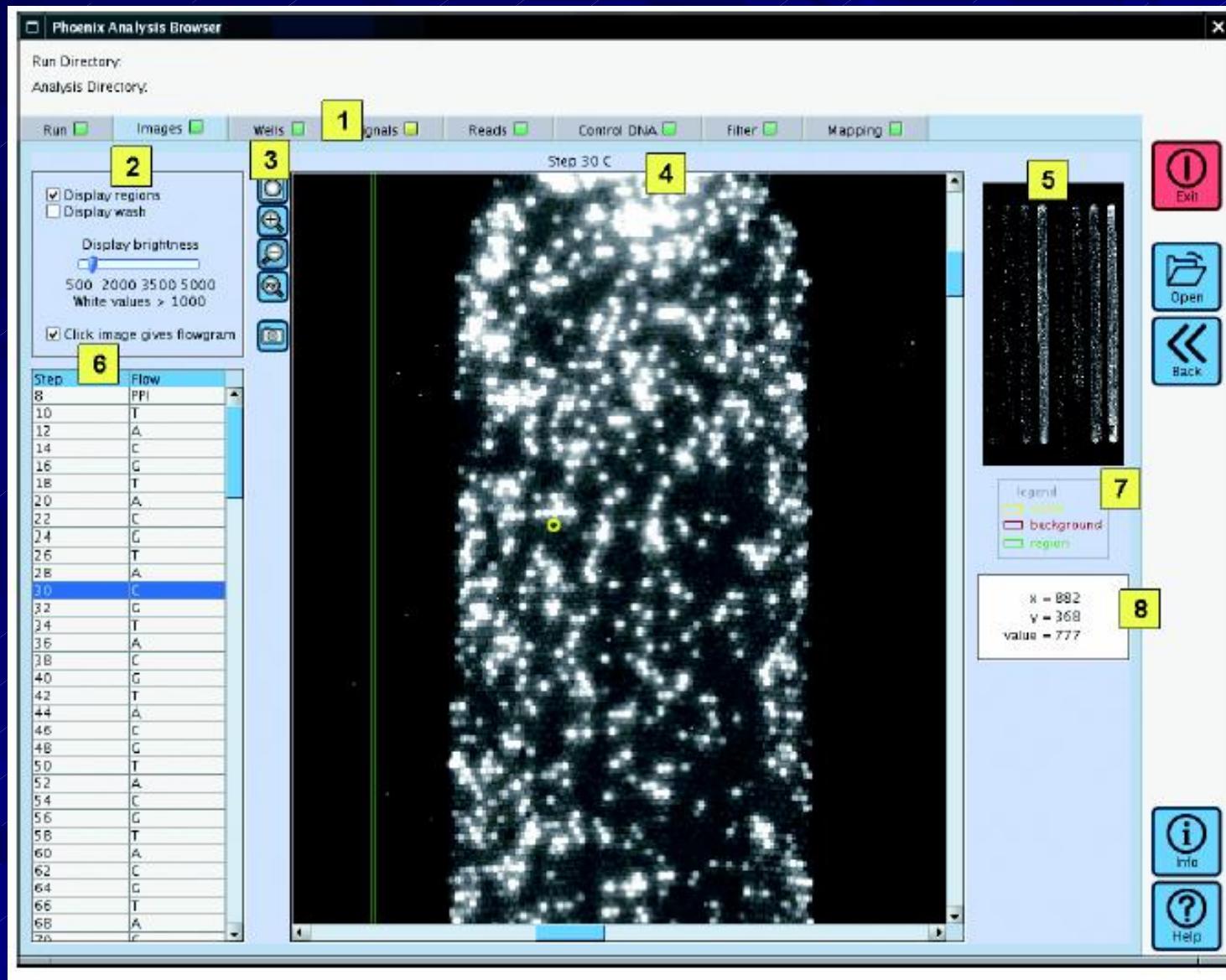
Adaptor Taq TCAG -- CTGA

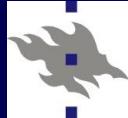
Lars Paulin Institute of Biotechnology University of Helsinki



Genome Sequencer GS20/FLX

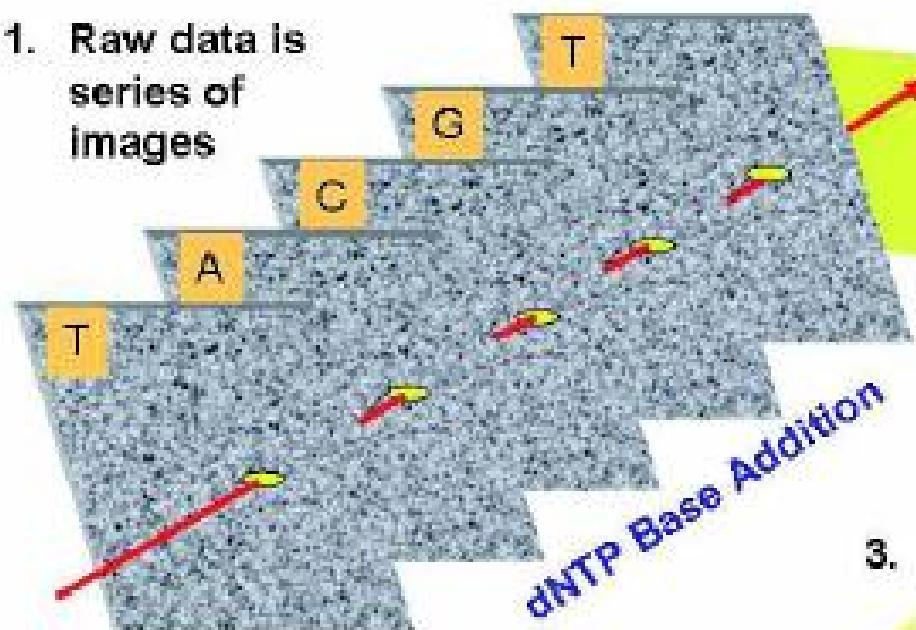






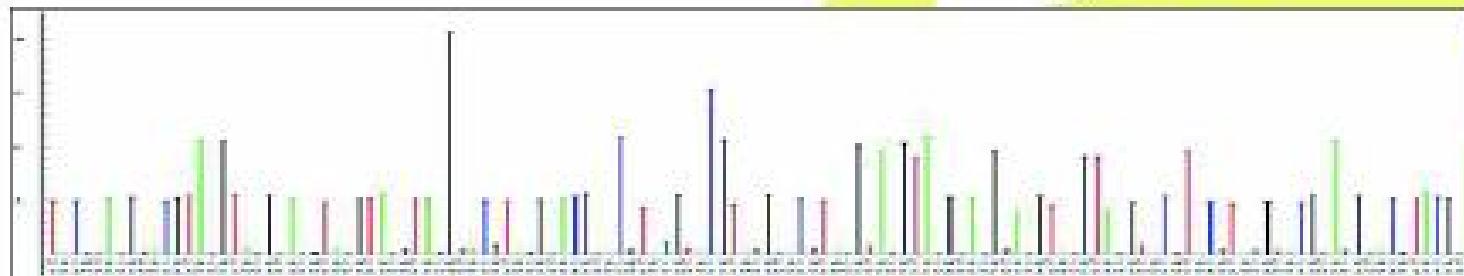
Flowgram

1. Raw data is series of images



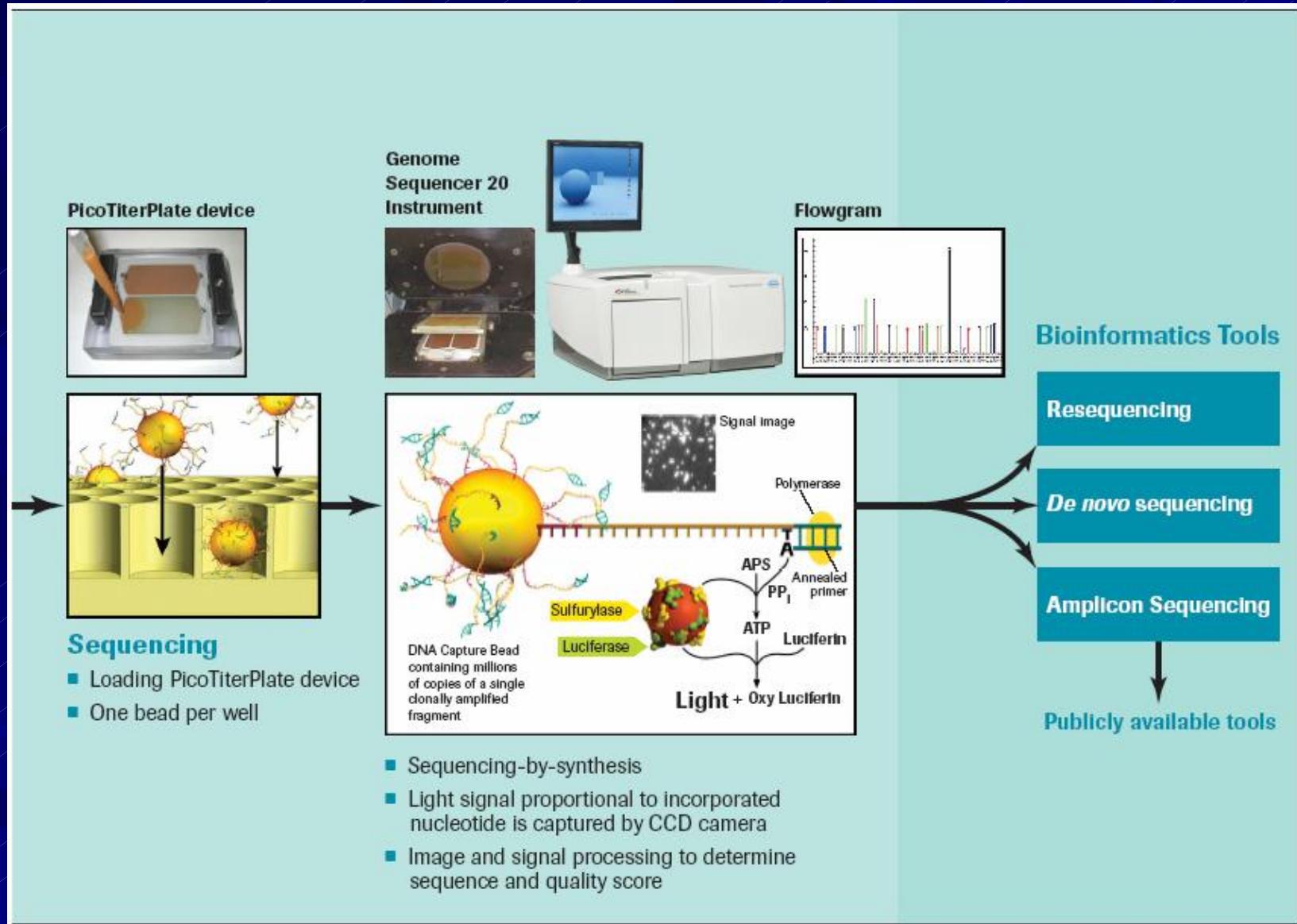
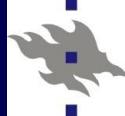
2. Each well's data extracted, quantized and normalized

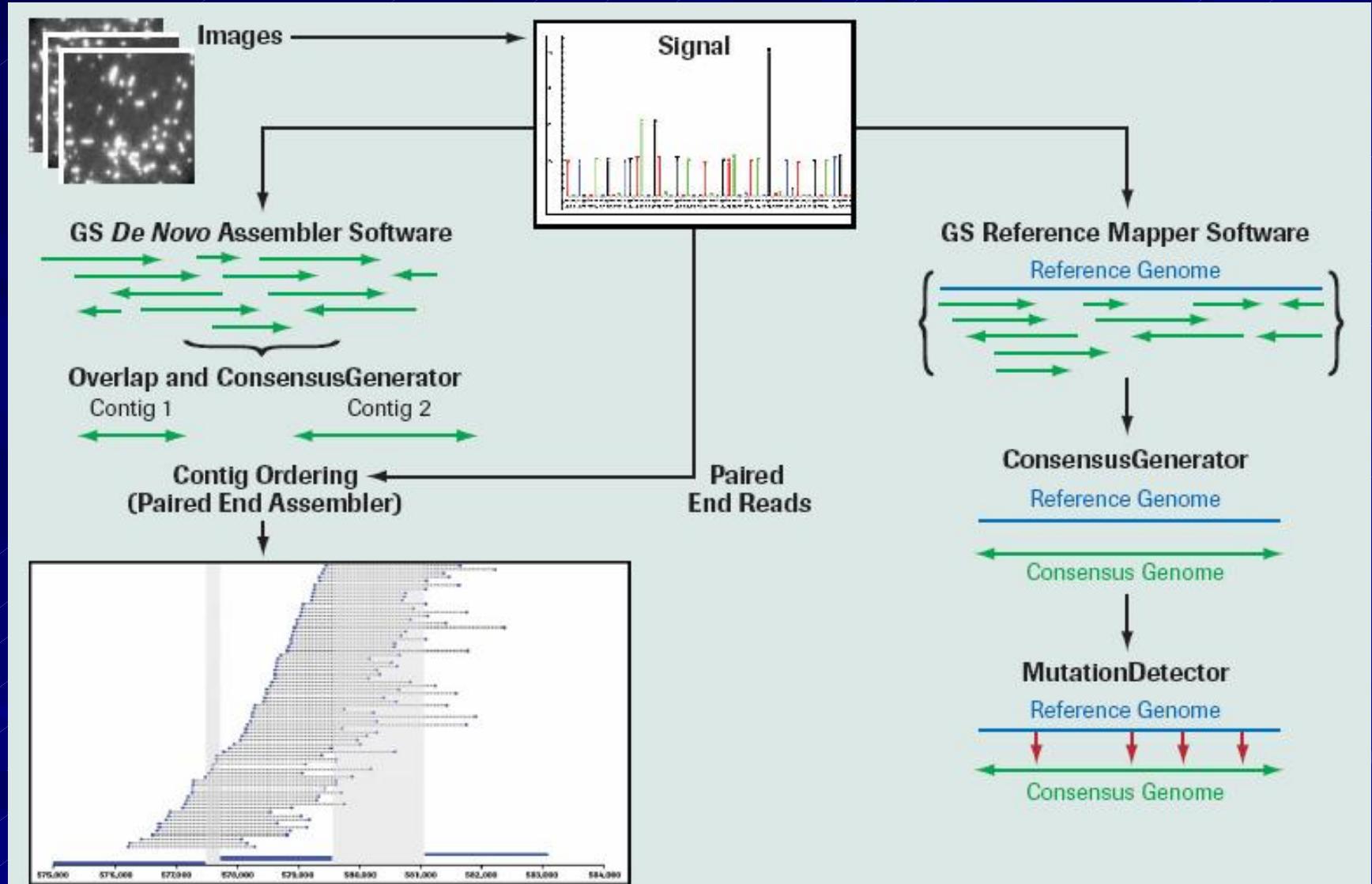
3. Read data converted into "flowgrams"

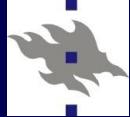


Adaptor Taq TCAG -- CTGA

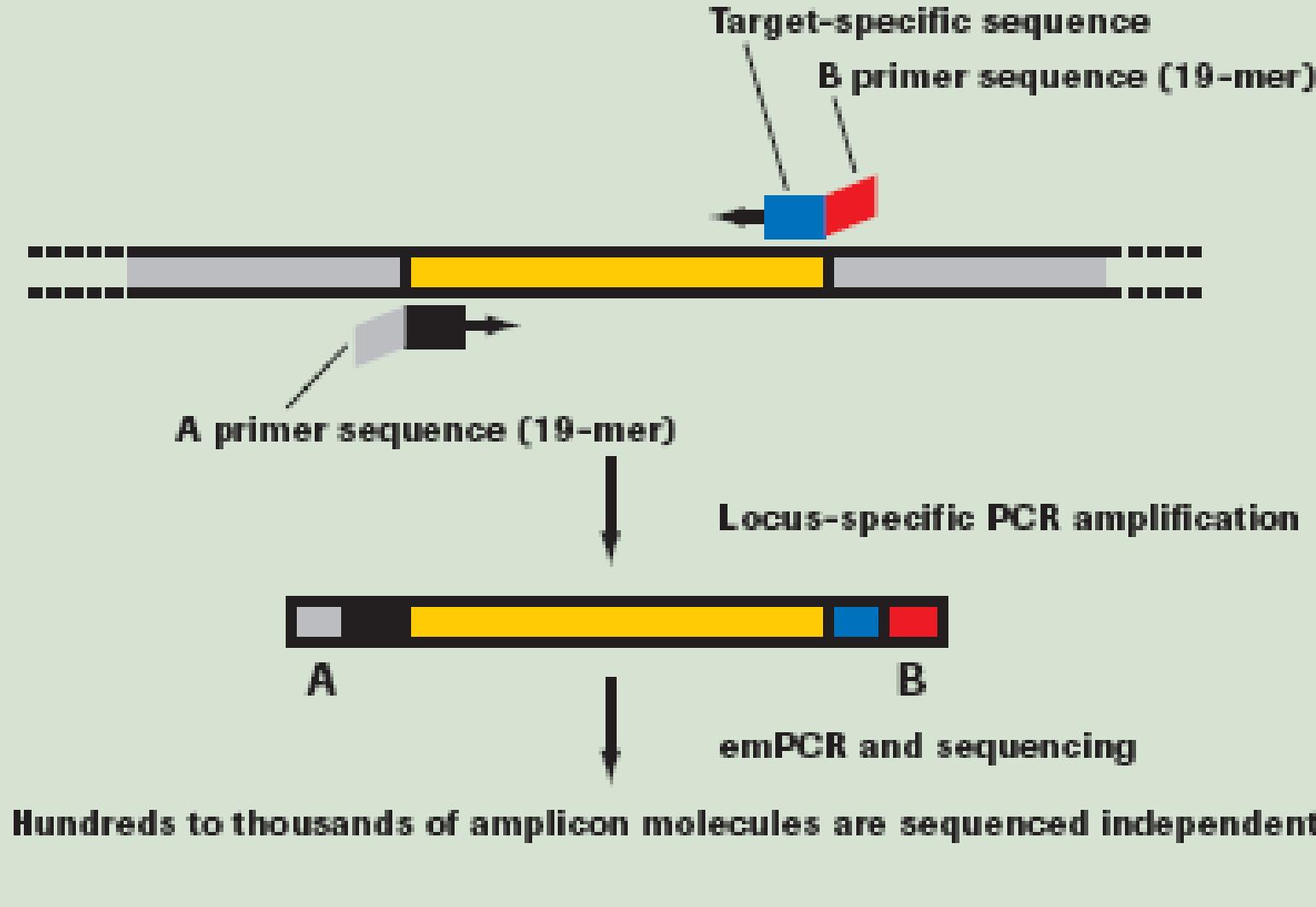
Lars Paulin Institute of Biotechnology University of Helsinki

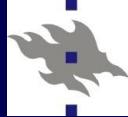




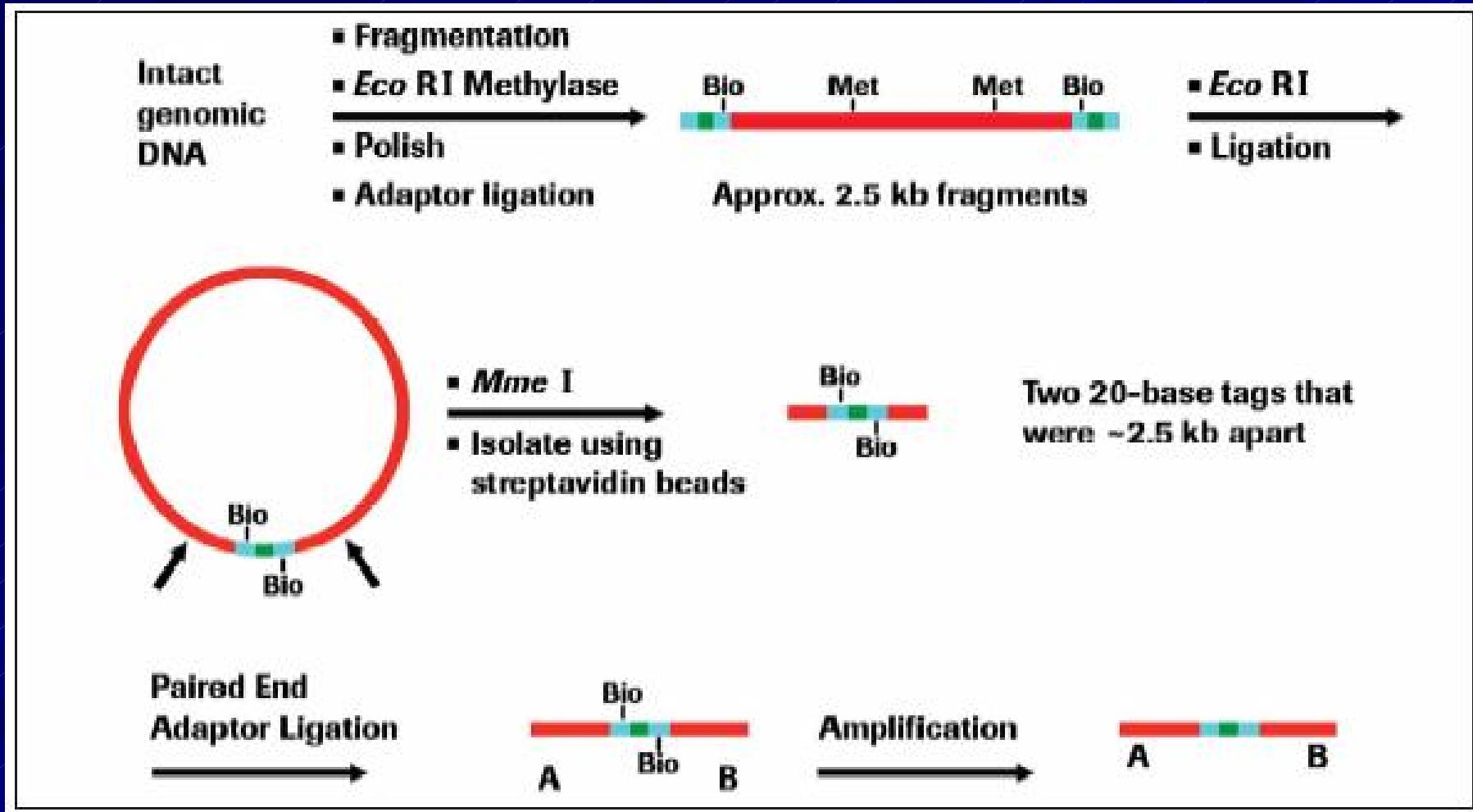


Amplicon sequencing





Paired-end Sequencing



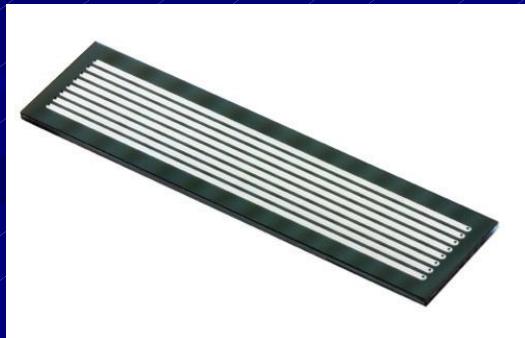


Illumina/Solexa Genome Analyzer

(<http://www.illumina.com>)

■ Clonal Single Molecule Array technology

- Sequencing-by-synthesis technology
- Reversible terminator-based sequencing
 - removable fluorescence
- Flow cell with > 10 million clusters
 - Each cluster ~1,000 copies of template /cm²
- 1–8 samples / run
- 3 laser system (660, 635, and 532 nm)
- Read length 35 - 50 bp, 1- 2 Gb / run
 - Run time 3 – 6 days,



Flow cell

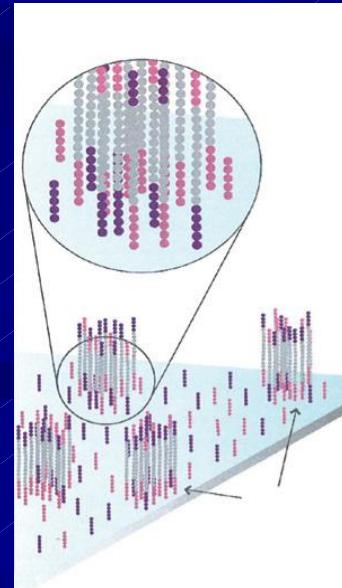
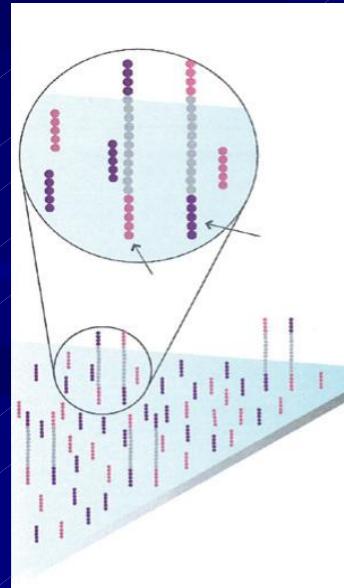
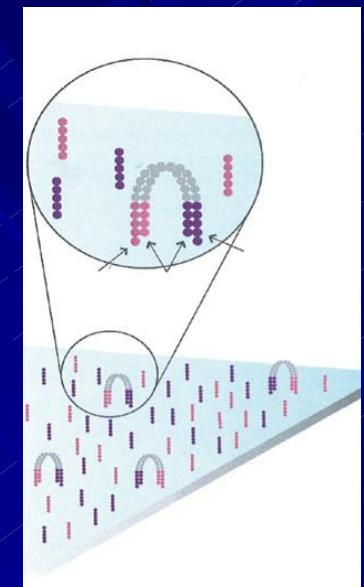
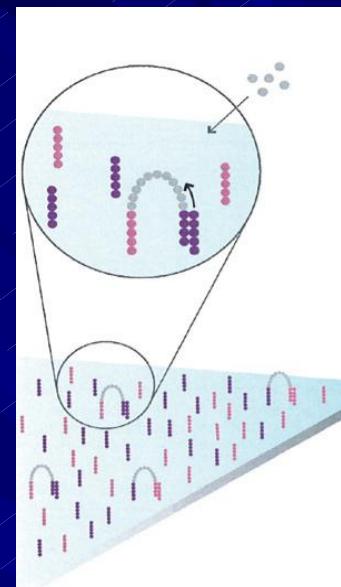
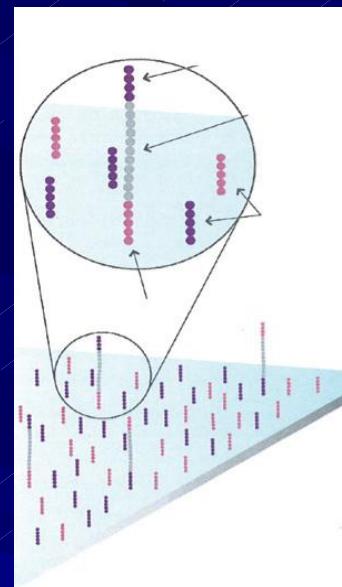
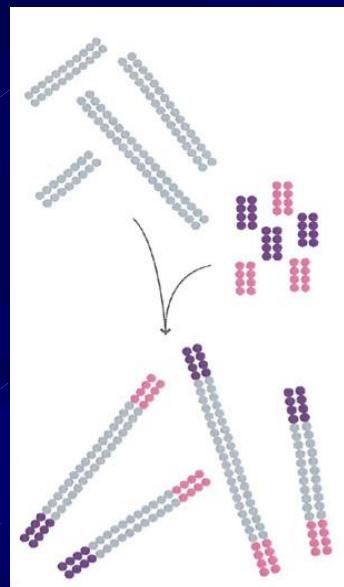


Cluster Station



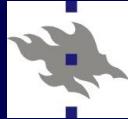


Illumina/Solexa

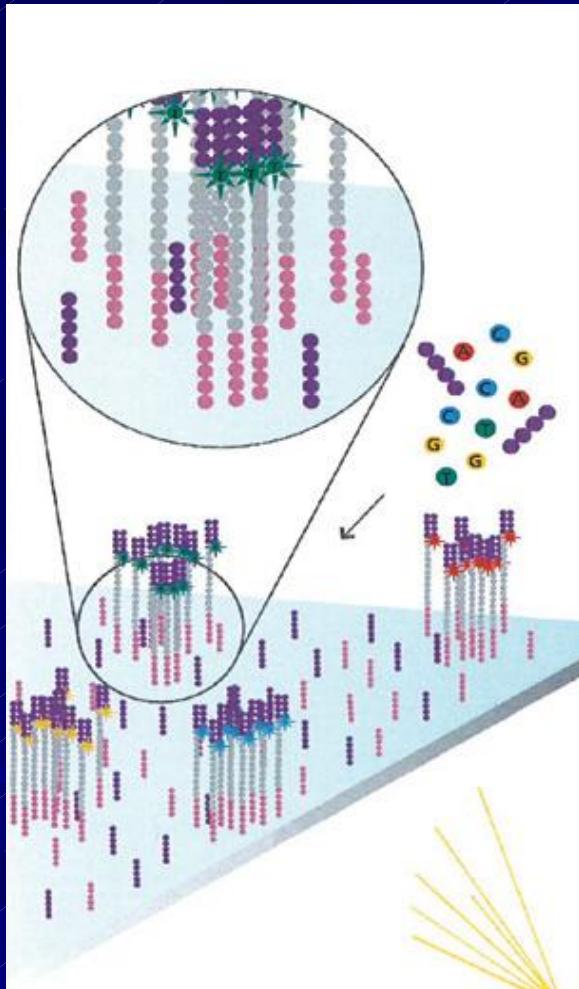


Sample preparation

- 100ng–1 μ g
- Attaching to Flow cell
- Bridging
- PCR
 - Elongation
 - Denaturation
 - Clonal amplification

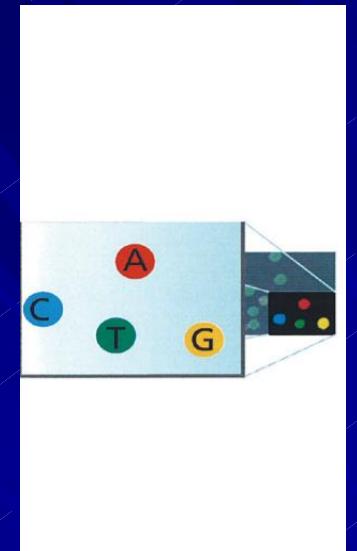
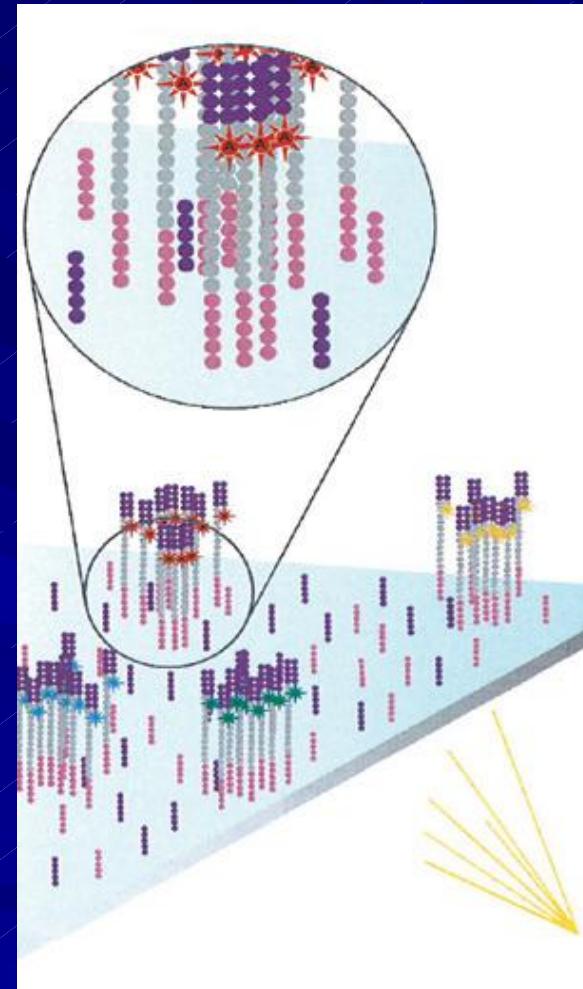
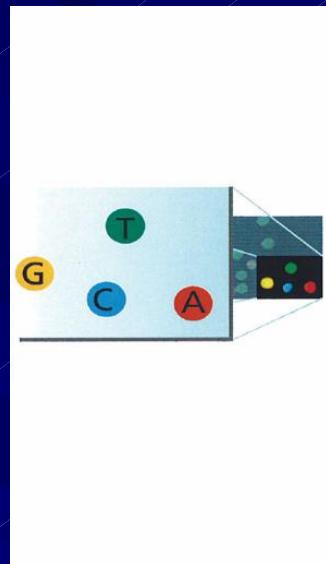


Illumina/Solexa sequencing



Sequencing

- First bases
- Fluorescent reversible terminators
- Detection with laser and CCD camera



Sequencing

- Second bases detected after removal of label and blocking



SOLiD, Applied Biosystems

(<http://www.appliedbiosystems.com>)

■ Sequencing by Ligation

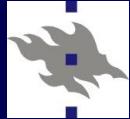
- emPCR
 - Small beads, 1µm
- Attaching to glass slides
- Labelled probes
 - Four colours
 - 2 base encoding system
- Repeated ligation steps
- Detection with 4 Mpixel camera
- Read length 25-30 bp
- 1-2 slides / run
- 1-2 Gb / run
- Run time 5 -10 days

Shendure, J. et.al. Science 2005,
309, 1728-1732

SOLiD



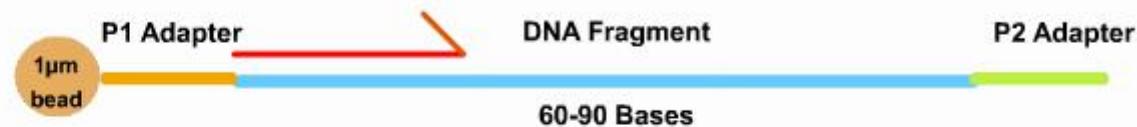
Lars Paulin Institute of Biotechnology University of Helsinki



SOLID

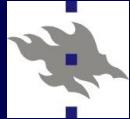
■ Library preparation

Fragment Library (directed resequencing)

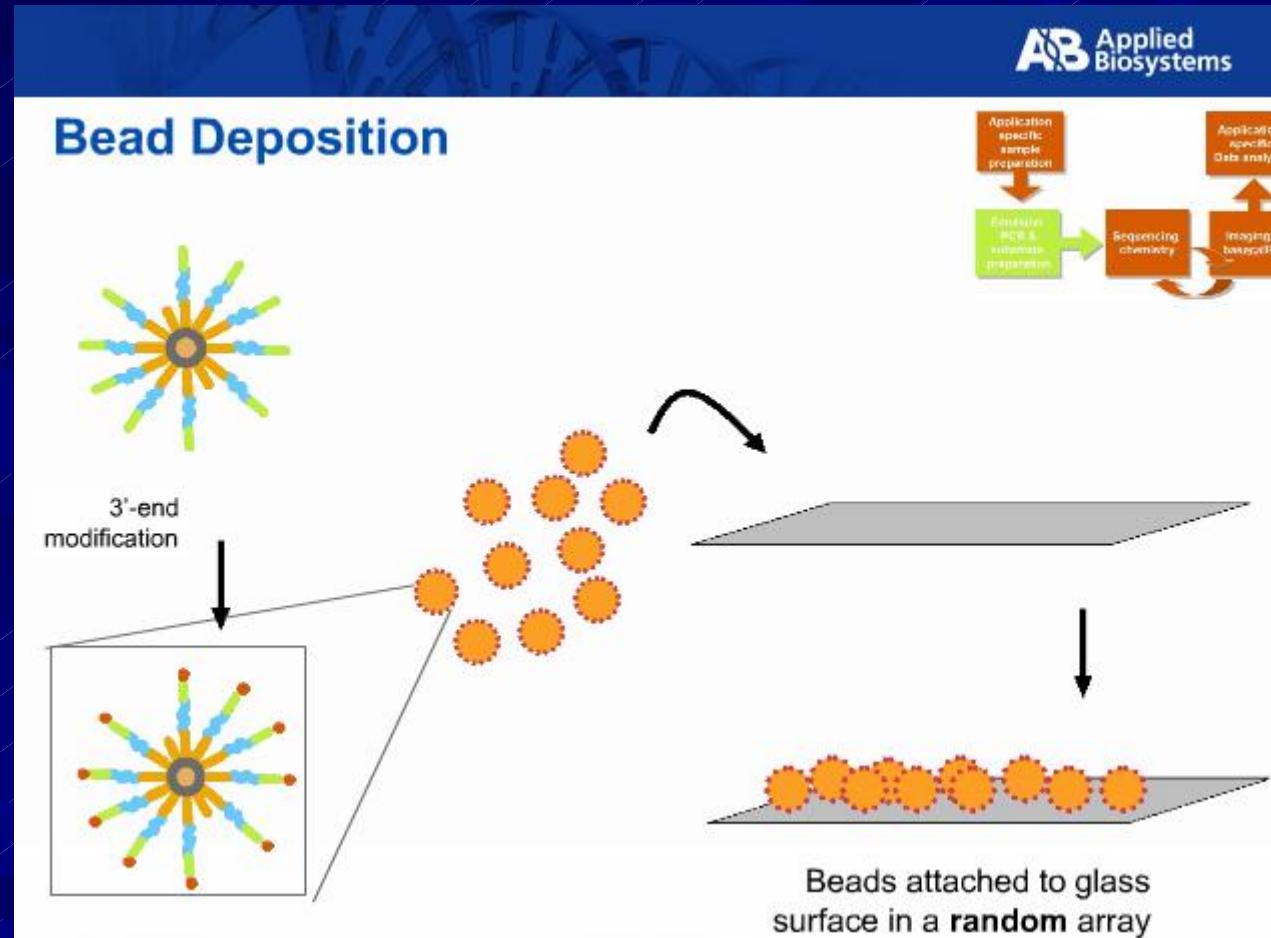


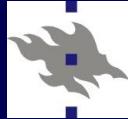
Mate Pair Library (whole genome sequencing)





SOLiD





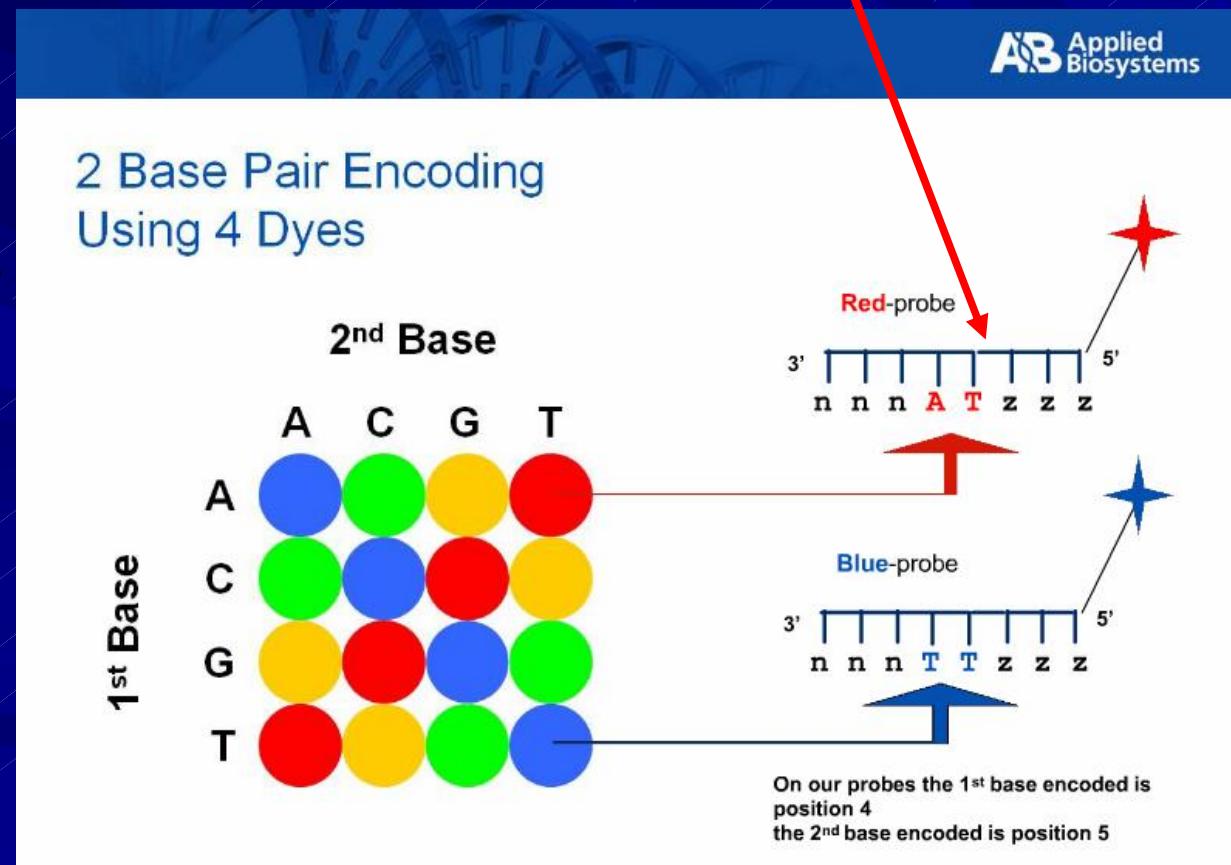
SOLID

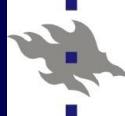
■ Probes

- 1 024 Octamer Probes
- 4 Dyes
- 4 dinucleotides
- 256 probes / dye

N = degenerate bases

Z = universal base

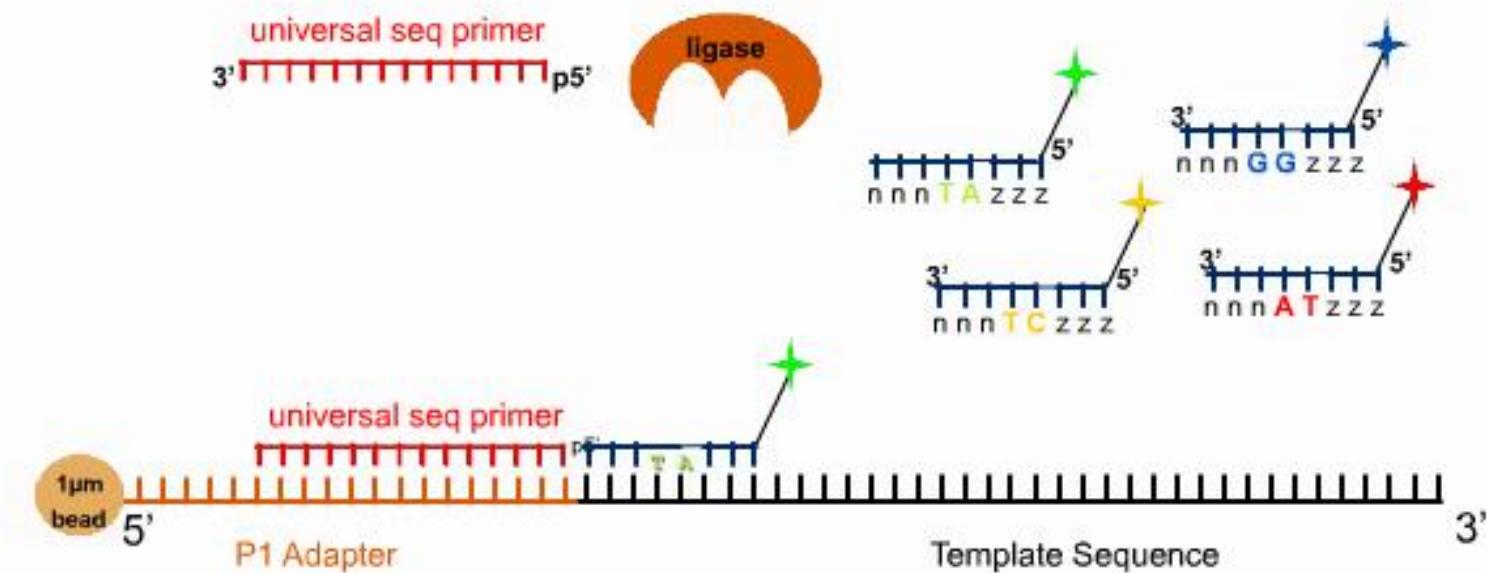
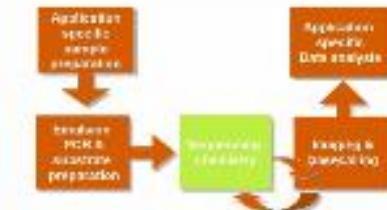


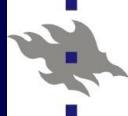


SOLiD

AB Applied Biosystems

SOLiDTM Chemistry System 4-color ligation Ligation reaction

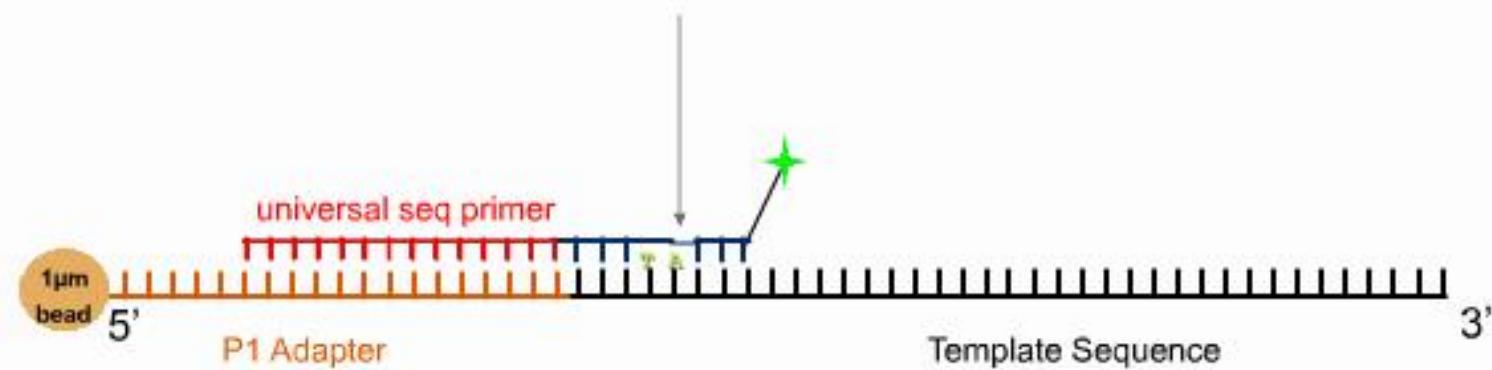
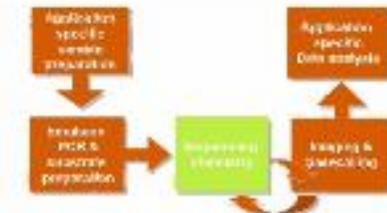


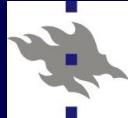


SOLiD

AB Applied Biosystems

SOLiD™ Chemistry System 4-color ligation Cleavage

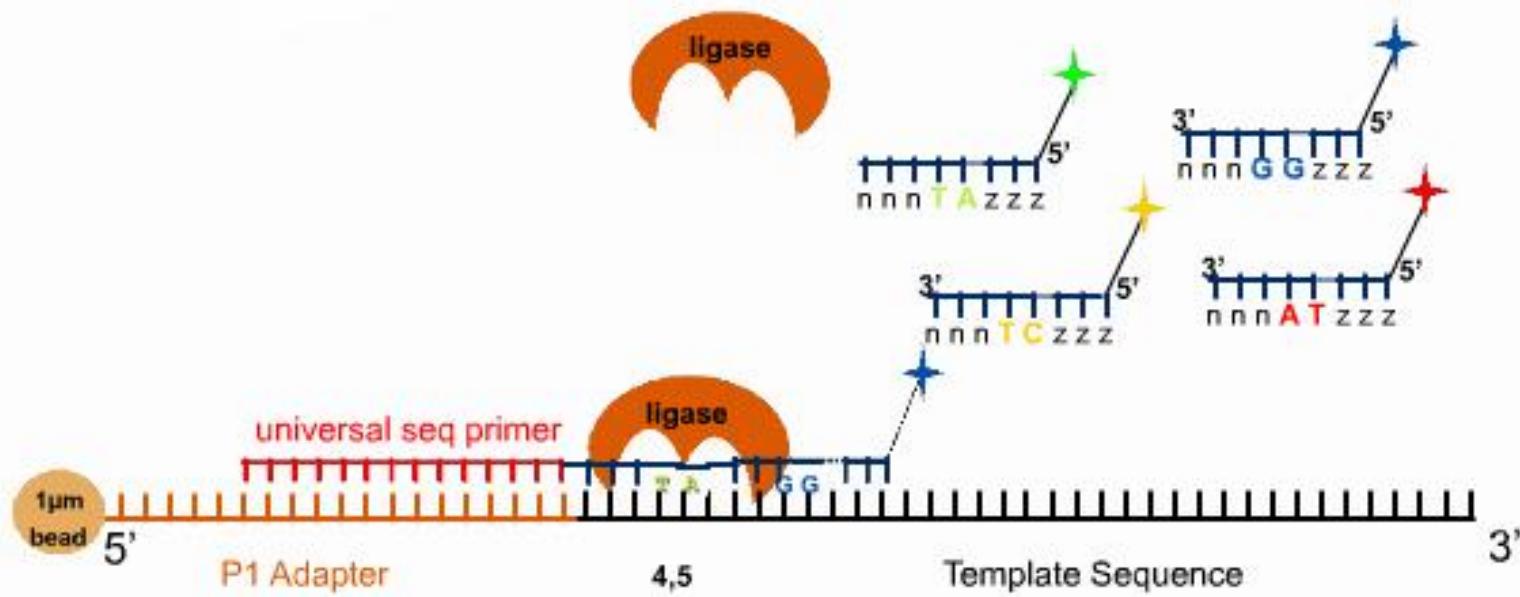
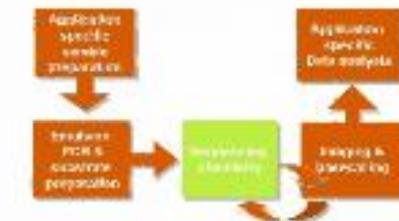




SOLiD

AB Applied Biosystems

SOLiD™ Chemistry System 4-color ligation Ligation (2nd cycle)

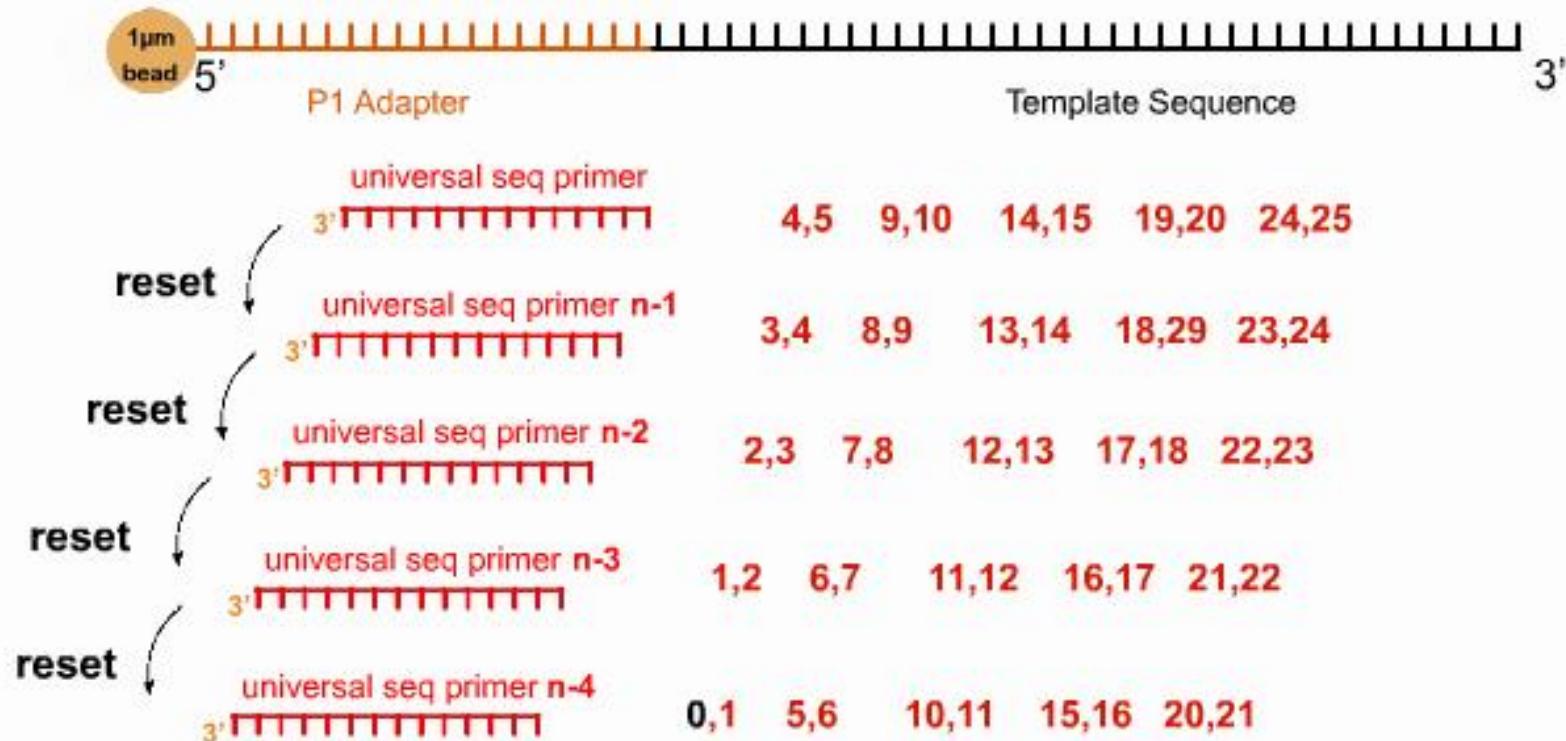


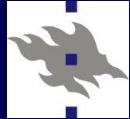


SOLID

AB Applied Biosystems

Sequential rounds of sequencing Multiple cycles per round

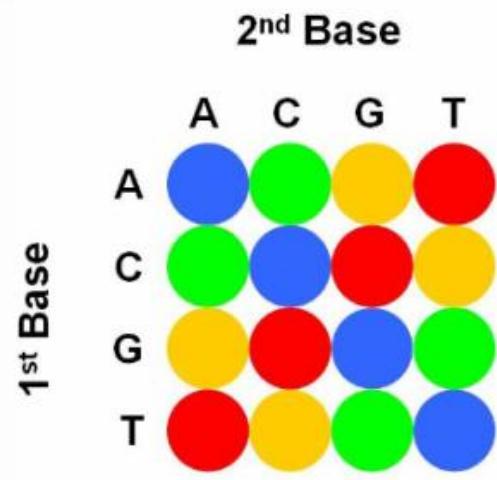
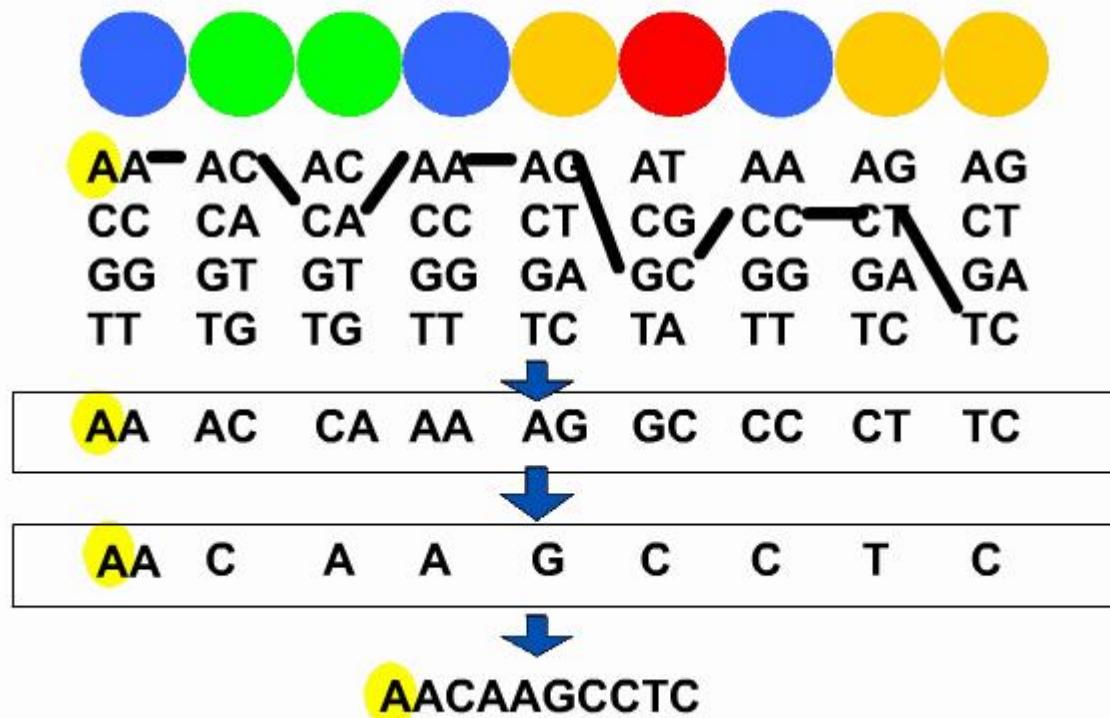




SOLiD

AB Applied Biosystems

Summary of decoding





Applications

■ Whole genome sequencing

- *de novo* sequencing
 - Genome Sequencer FLX

■ Comparative sequencing

- All three systems

■ Metagenomics

- Genome Sequencer FLX

■ Amplicon sequencing

- Mutations / SNP
- All three systems

■ Transcriptome sequencing

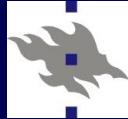
- cDNA
 - All three systems
- Small RNA
 - All three systems

■ ChIP sequencing

- All three systems

■ Methylation sequencing

- All three systems



Other technologies

■ Helicos (www.helicosbio.com)

- Sequencing-by-synthesis
- No PCR amplification
- 25-90 Mb/h

■ VisiGen (www.visigenbio.com)

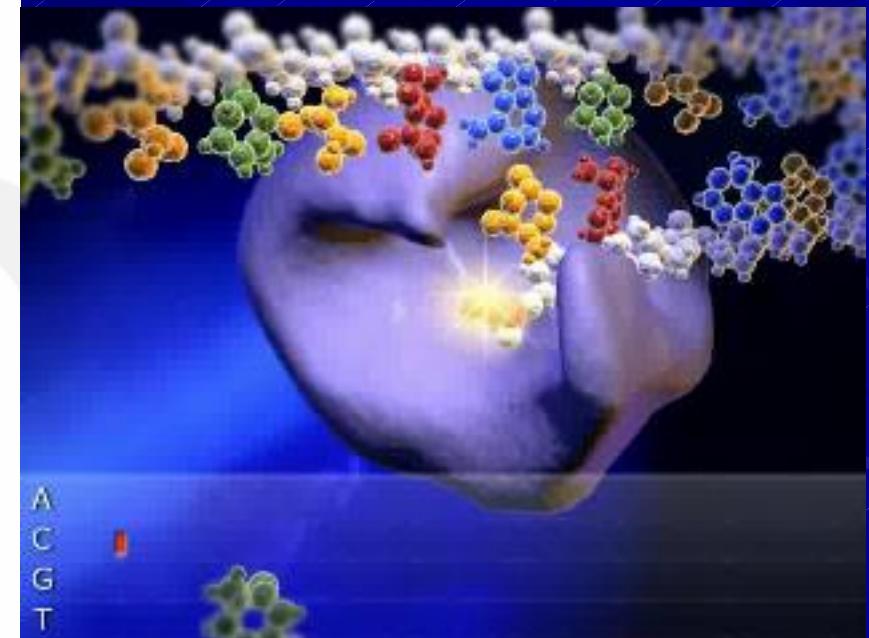
- Real-time detection of DNA synthesis, FRET
- Intact DNA fragments
- 1Mb/sec/machine

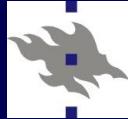
World's First Single Molecule Genetic Analyzer
The HeliScope™ System



Initial throughput is planned to range from 25 to 90 million bases of DNA per hour

- Imaging capacity of the instrument is ~1 Billion bases per hour
- Improvements to the tSMS chemistry and the flow cells will provide customers significant performance gains





<http://genomics.xprize.org/genomics>

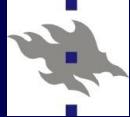
\$10M to the First Team to Sequence 100 Human Genomes in 10 Days

Registered Teams

- 454 Life Sciences (Roche) (www.454.com)
- VisiGen (www.visigenbio.com)
- FfAME (www.ffame.org)
- Reveo (www.reveo.com)
- Base4innovation (www.base4innovation.co.uk/)
- Personal Genome X-Team (PGx), George Church



Lars Paulin Institute of Biotechnology University of Helsinki



Assemblers

- Phrap, Phil Green
 - University of Washington
- TIGR assembler
 - TIGR
- Celera
 - Celera Corporation
- Euler
 - University of California
- Arachne
 - Broad Institute
- CAP3, PCAP
 - Iowa State University
- gsAssembler
 - Roche, 454 Life Science
- Amos
 - University of Maryland
- SHARCGS
 - Genome Res,2007,17,1697
- SAKE
 - Bioinformatics,2007,23,500
- SHRAP
 - PlosONE,2007,2:e484
- New Euler
 - Genome Res,2007,Dec14