

NEXT GENERATION SEQENCING **Technologies, Applications & Analysis**

AGTCCGCGAATACAGGCTCGGTAGTCCGCGAATACAGGCTCGGT

Petrus Tang, Ph.D. (鄧致剛)

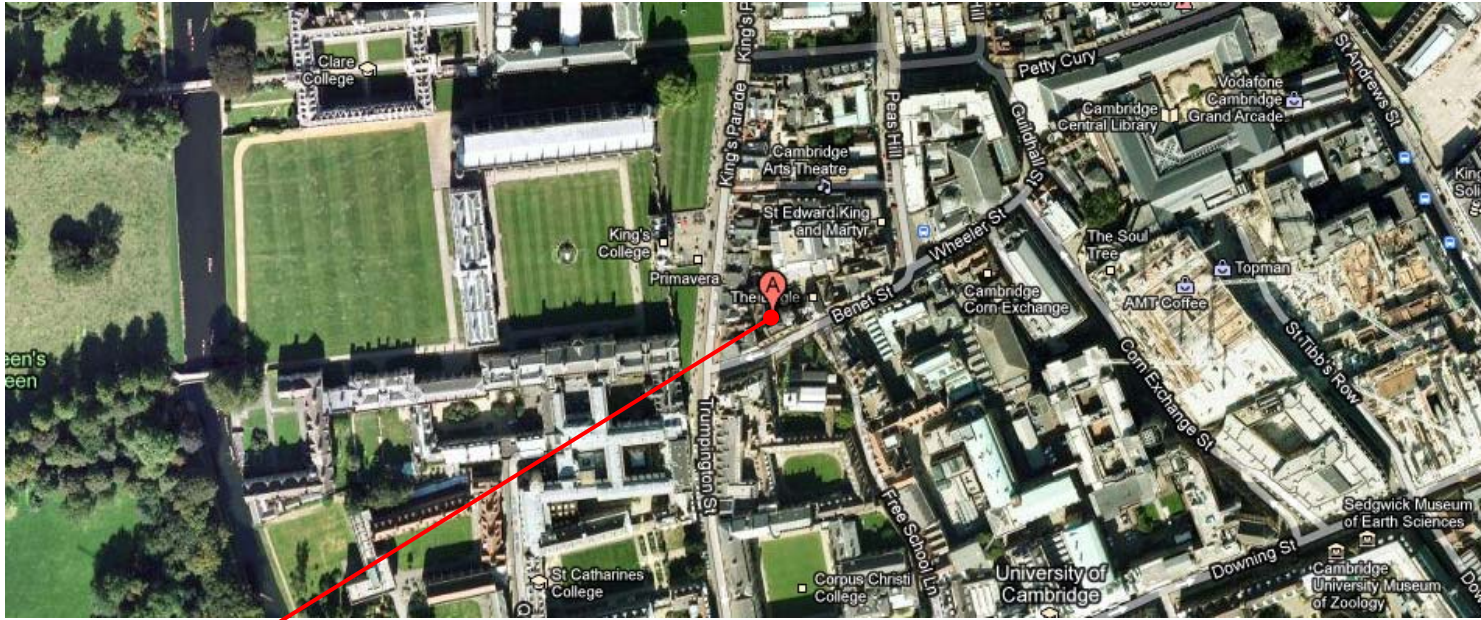
Graduate Institute of Basic Medical Sciences

and

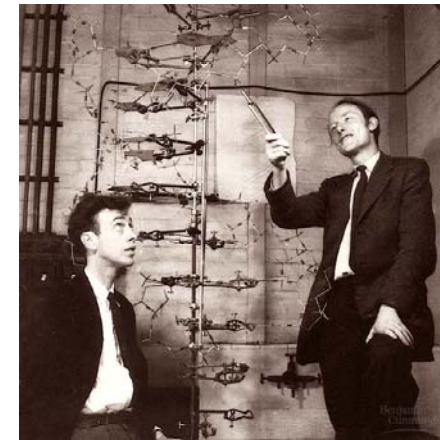
Bioinformatics Center, Chang Gung University.

petang@mail.cgu.edu.tw

<http://petang.cgu.edu.tw>



The Eagle, Cambridge: the place where Francis Crick interrupted patrons' lunchtime on 28 February 1953 to announce that he and James Watson had "discovered the secret of life" after they had come up with their proposal for the structure of DNA





The Nobel Prize in Chemistry 1980

"for his fundamental studies of the biochemistry of nucleic acids, with particular regard to recombinant-DNA"

"for their contributions concerning the determination of base sequences in nucleic acids"



Paul Berg

🕒 1/2 of the prize

USA

Stanford University
Stanford, CA, USA



Walter Gilbert

🕒 1/4 of the prize

USA

Harvard University,
Biological Laboratories
Cambridge, MA, USA

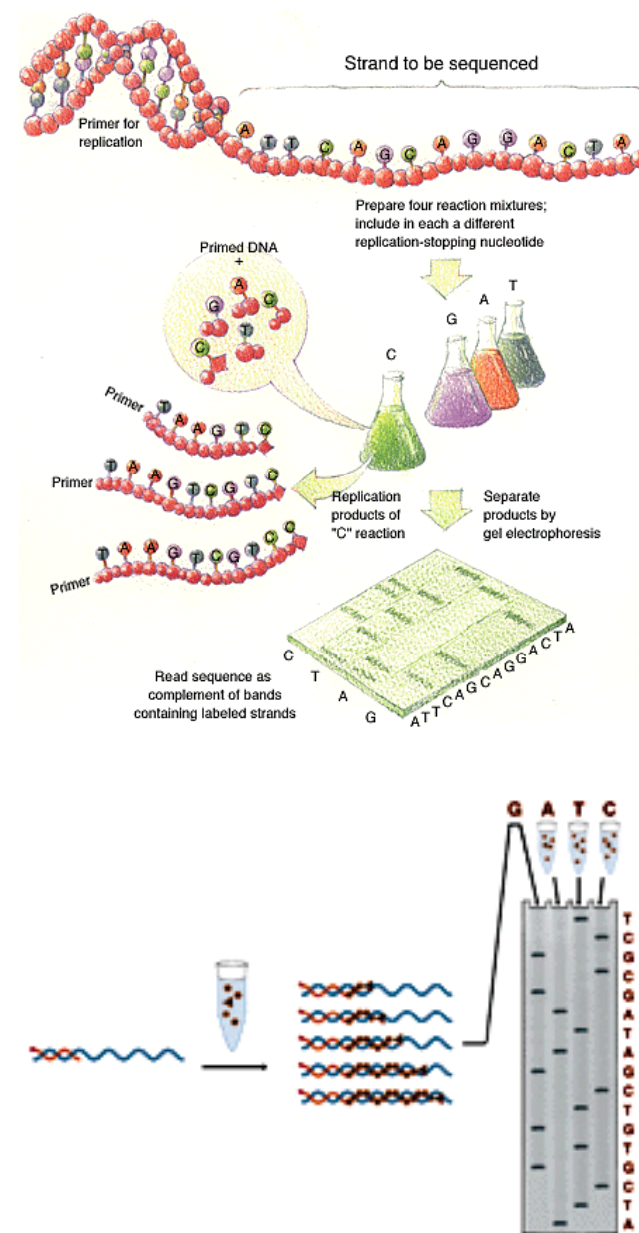


Frederick Sanger

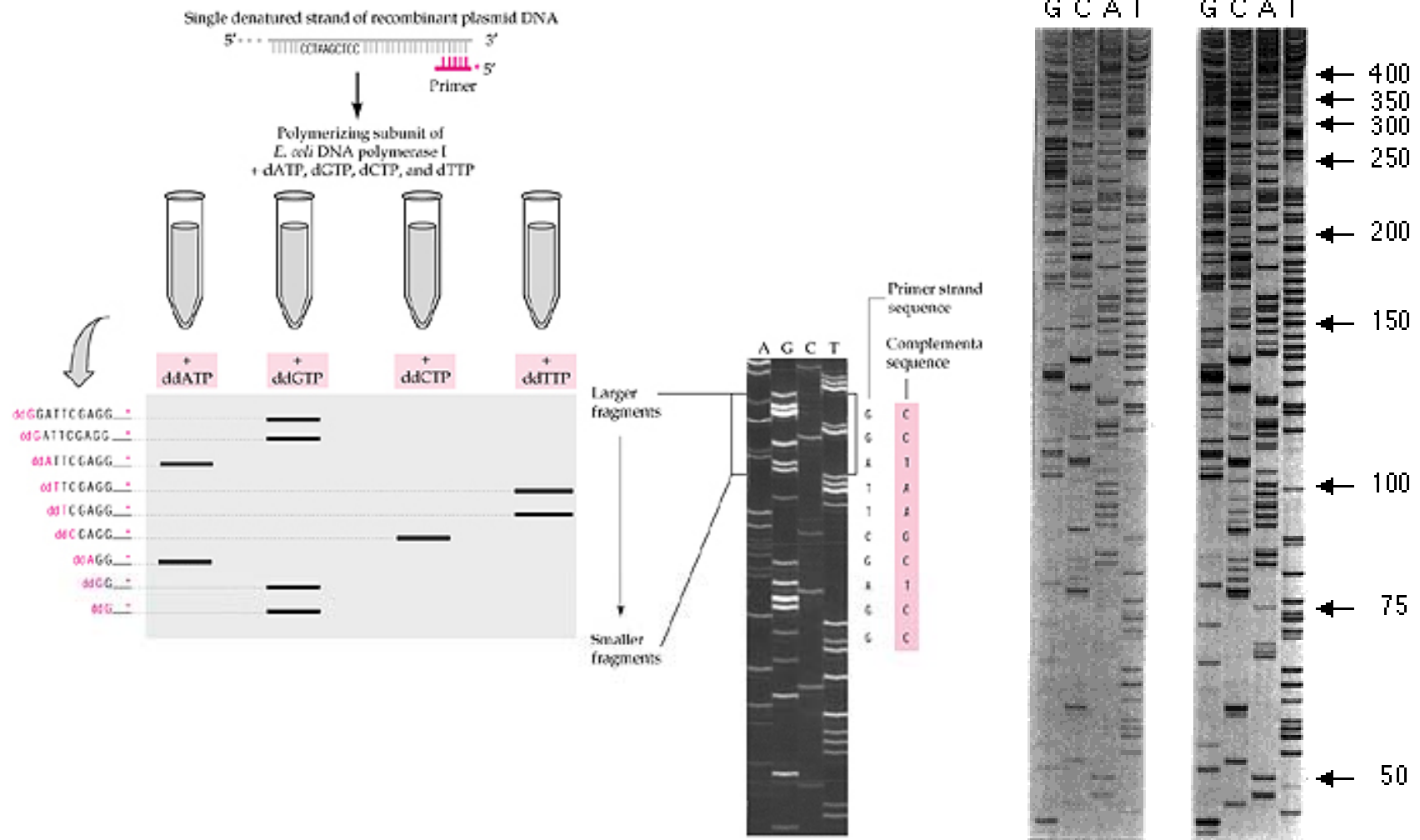
🕒 1/4 of the prize

United Kingdom

MRC Laboratory of
Molecular Biology
Cambridge, United



Sanger Dideoxy Sequencing

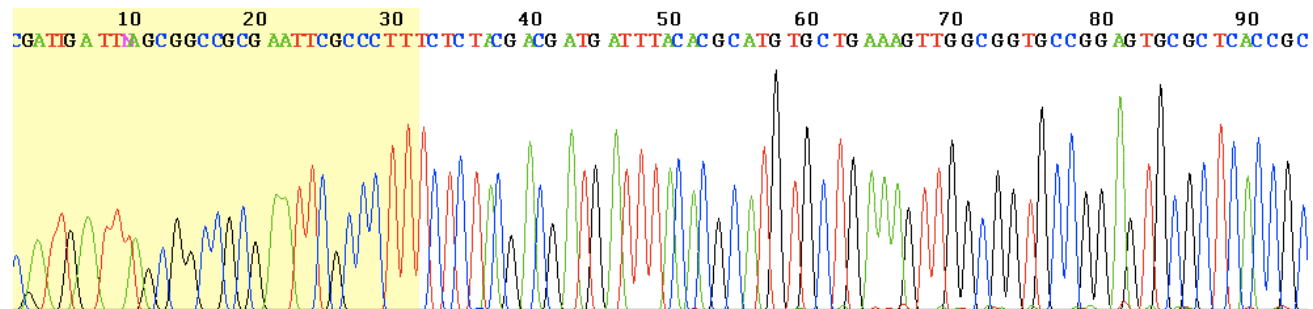


Sanger, F. et al. Nature 24, 687–695 (1977).

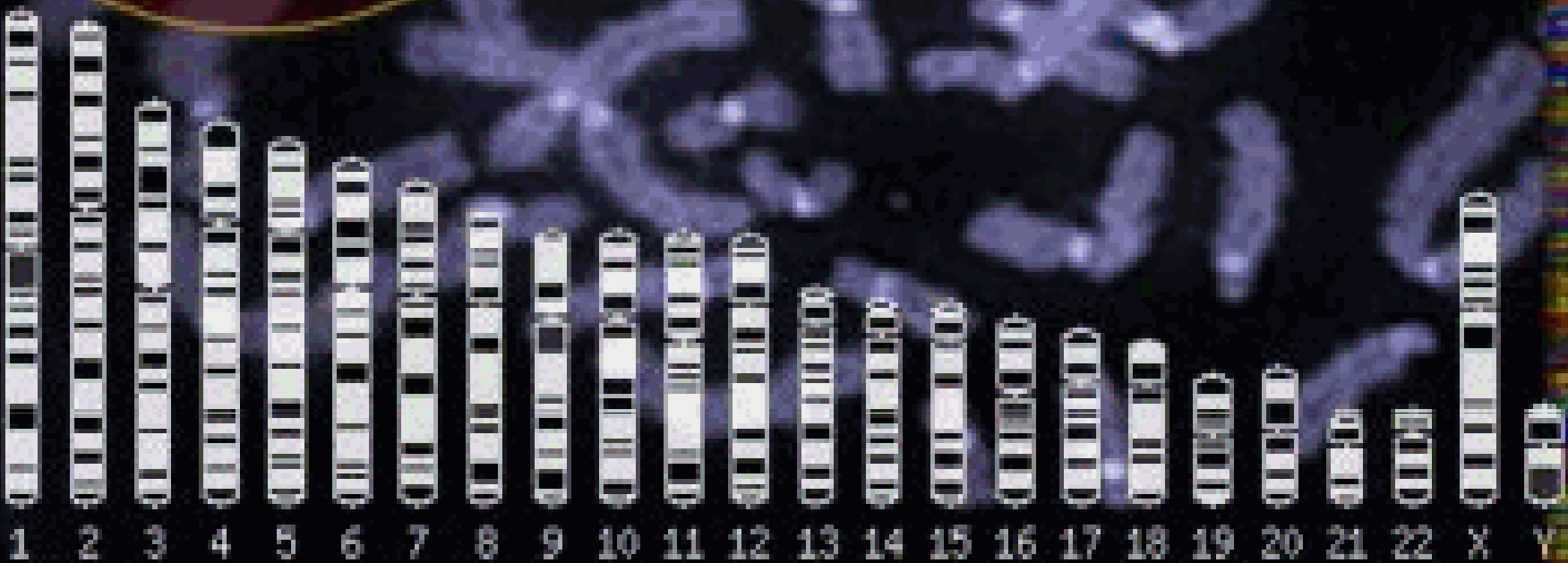
Sanger, F., Nicklen, S. & Coulson, A.R. Proc. Natl. Acad. Sci. USA 74, 5463–5467 (1977).

Basics of the “Old” Technology

- Clone the DNA.
- Generate a ladder of labeled (colored) molecules that are different by 1 nucleotide.
- Separate mixture on some matrix.
- Detect fluorochrome by laser.
- Interpret peaks as string of DNA.
- Strings are 500 to 1,000 letters long
- 1 machine generates 57,000 nucleotides/run
- Assemble all strings into a “whole”.



Human Genome Project

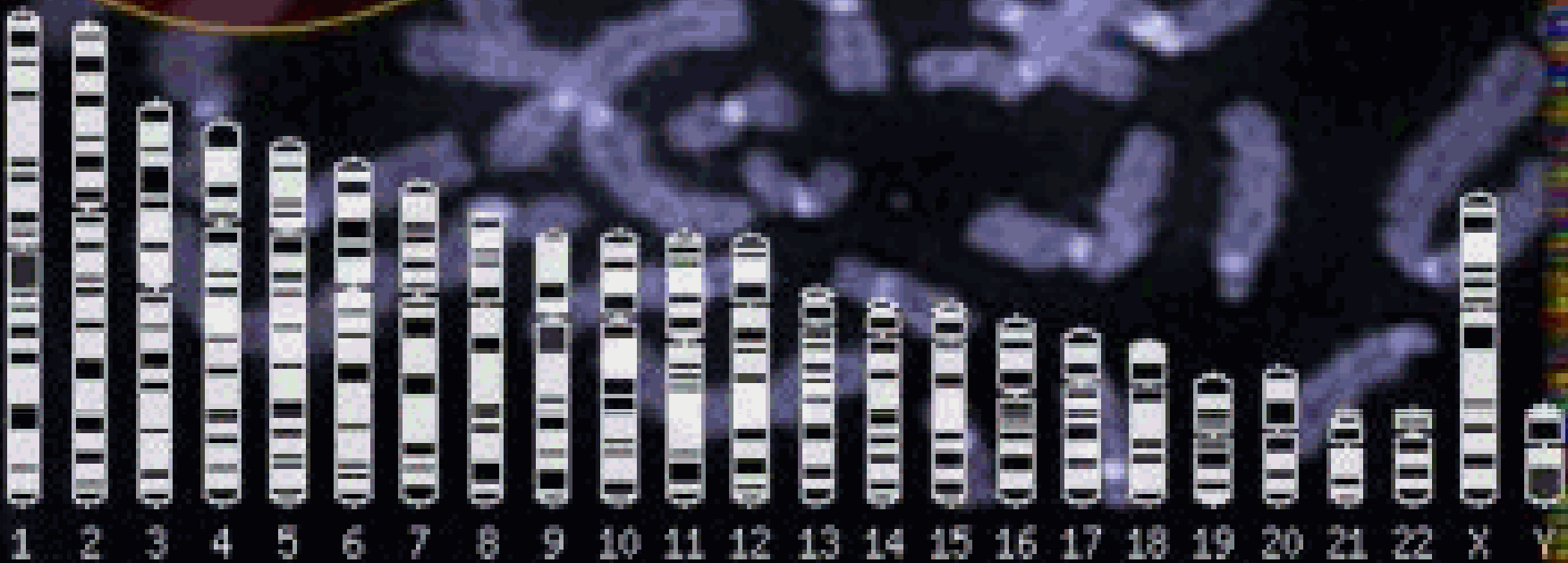


3 billion basepairs

Human Genome Project

3G

3 billion basepairs



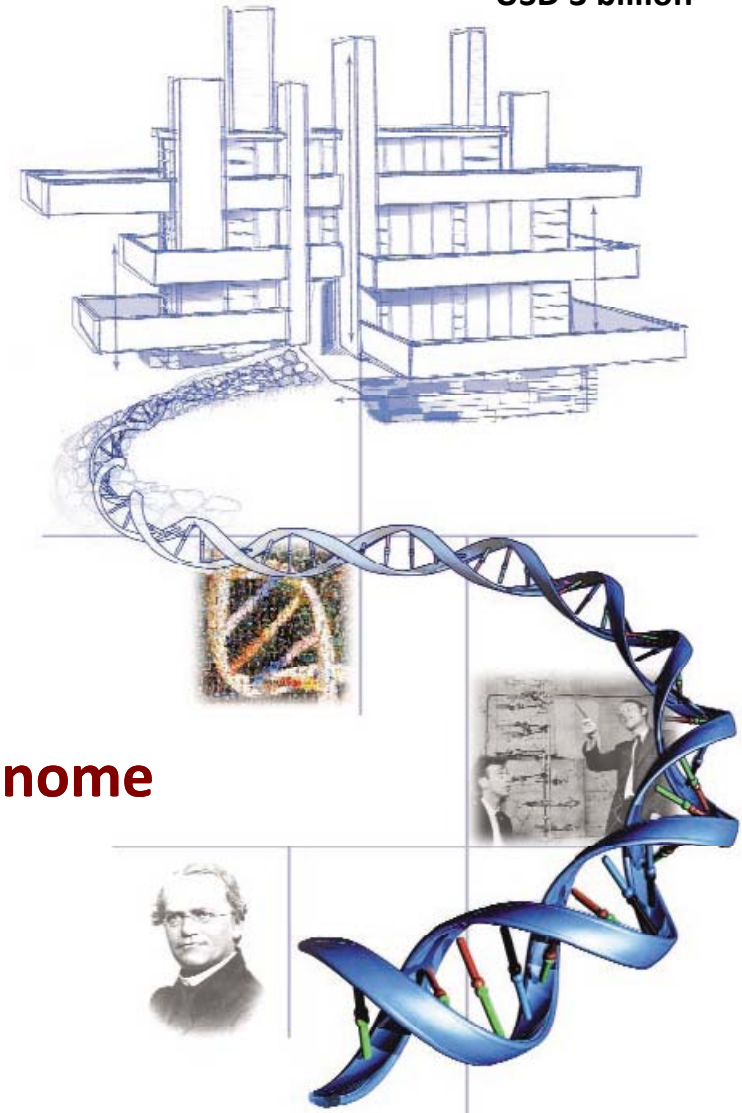
February 2001: Completion of the Draft Human Genome >10 years to finish USD 3 billion



Science, 16 February 2001
Vol. 291, Pages 1145-1434



Nature, 15 February 2001
Vol. 409, Pages 813-960



April 2003: High-Resolution Human Genome

A vision for the future of genomics research

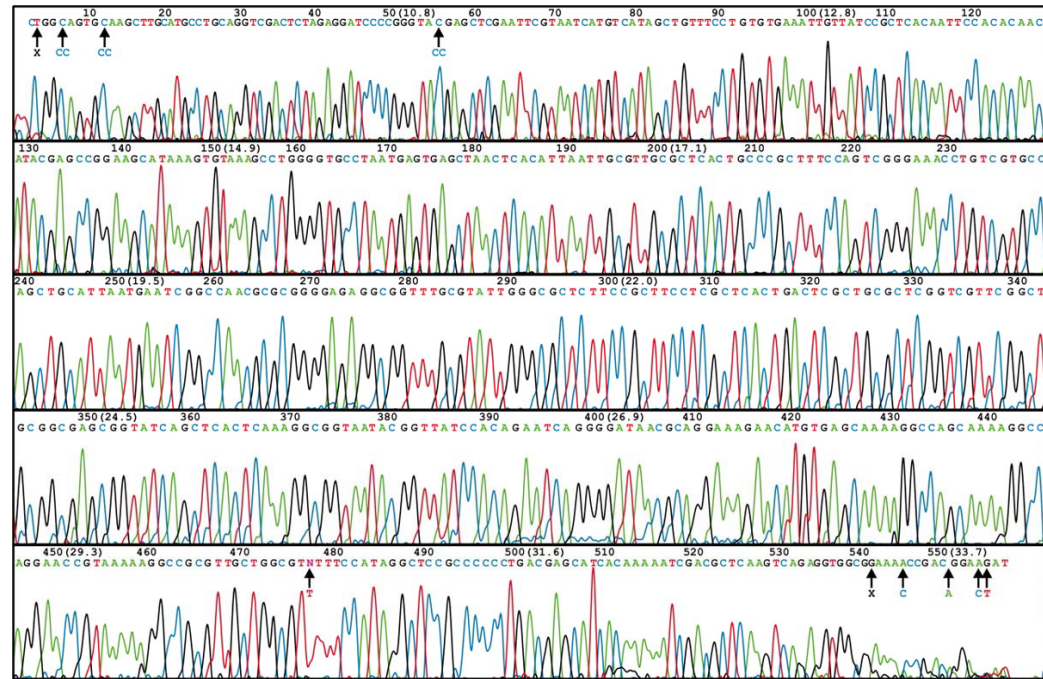
A blueprint for the genomic era.

Francis S. Collins, Eric D. Green,
Alan E. Guttmacher and Mark S.
Guyer on behalf of the US National
Human Genome Research Institute*

50 Years of DNA: *From Double Helix to Health*
A Celebration of the Genome

Nature, 23 April 2003
Vol. 422, Pages 1-13

ABI 3730 XL DNA Sequencer



96/384 DNA sequencing in 2 hrs, approximately 600-1000 readable bps per run.

1-4 MB bps/day

A human genome of 3GB need 750 days to finish 1X coverage





The Sequence of the Human Genome

J. Craig Venter, Mark D. Adams, Eugene W. Myers, Peter W. Li, Richard J. Mural, Granger G. Sutton, Hamilton O. Smith, Mark Yandell, Cheryl A. Evans, Robert A. Hott, Jeannine D. Gocayne, Peter Amanatides, Richard M. Ballwe, Daniel H. Huson, Jennifer Russo Wortman, Qing Zhang, Chinnappa D. Kodira, Xiangqun H. Zheng, Lin Chen, Marian Skupski, Gangadharan Subramanian, Paul D. Thomas, Jinghui Zhang, George L. Gabor Miklos, Catherine Nelson, Samuel Broder, Andrew G. Clark, Joe Nadeau, Victor A. McKusick, Norton Zinder, Arnold J. Levine, Richard J. Roberts, Mel Simon, Carolyn Slayman, Michael Hunkapiller, Randall Bolanos, Arthur Delcher, Ian Dew, Daniel Fasulo, Michael Flanigan, Liliana Florea, Aaron Halpern, Sridhar Hannehalli, Saul Kravitz, Samuel Levy, Clark Mobarry, Knut Reinert, Karin Remington, Jane Abu-Threideh, Ellen Beasley, Kendra Biddick, Vivien Bonazzi, Rhonda Brandon, Michele Cargill, Ishwar Chandramouliswaran, Rosane Charlab, Kabir Chaturvedi, Zuoming Deng, Valentina Di Francesco, Patrick Dunn, Karen Eilbeck, Carlos Evangelista, Andrei E. Gabrielian, Weiniu Gan, Wangmao Ge, Fangcheng Gong, Zhiping Gu, Ping Guan, Thomas J. Heiman, Maureen E. Higgins, Rui-Ru Ji, Zhaoxi Ke, Karen A. Ketchum, Zhongwu Lai, Yiding Lei, Zhenya Li, Jiayin Li, Yong Liang, Xiaoying Lin, Fu Lu, Gennady V. Merkulov, Natalia Milshina, Helen M. Moore, Ashwinikumar K Naik, Vaibhav A. Narayan, Beena Neelam, Deborah Nusskern, Douglas B. Rusch, Steven Salzberg, Wei Shao, Bixiong Shue, Jingtao Sun, Zhen Yuan Wang, Aihui Wang, Xin Wang, Jian Wang, Ming-Hui Wei, Ron Whaley, Chunlin Xiao, Chunhua Yao, Alison Yao, Jane Ye, Ming Zhan, Weiqing Zhang, Hongyu Zhang, Qi Zhao, Liansheng Zheng, Fei Zhong, Wenyang Zhong, Shiaoping C. Zhu, Shaying Zhao, Dennis Gilbert, Suzanna Baumhueter, Gene Spier, Christine Carter, Anibal Cravchik, Trevor Woodage, Feroze Ali, Huijin An, Aderonke Awe, Danita Baldwin, Holly Baden, Mary Barnstead, Ian Barrow, Karen Beeson, Dana Busam, Amy Carver, Angela Center, Ming Lai Cheng, Liz Curry, Steve Danaher, Lionel Davenport, Raymond Desilets, Susanne Dietz, Kristina Dodson, Lisa Doup, Steven Ferreira, Neha Garg, Andras Gluecksmann, Brit Hart, Jason Haynes, Charles Haynes, Cheryl Hainer, Suzanna Hladun, Damon Hostin, Jarrett Houck, Timothy Howland, Chinyere Ibegwam, Jeffery Johnson, Francis Kalush, Lesley Kline, Shashi Koduru, Amy Love, Felecia Mann, David May, Steven McCawley, Tina McIntosh, Mee Moy, Linda Moy, Brian Murphy, Keith Nelson, Cynthia Pfannkoch, Eric Pratts, Vinita Puri, Hina Qureshi, Matthew Reardon, Robert Rodriguez, Yu-Hui Rogers, Deanna Romblad, Bob Ruhfel, Richard Scott, Cynthia Sitter, Michelle Smallwood, Erin Stewart, Renee Strong, Ellen Suh, Reginald Thomas, Ni Ni Tint, Sukyee Tse, Claire Vech, Gary Wang, Jeremy Wetter, Sherita Williams, Monica Williams, Sandra Windsor, Emily Winn-Deen, Keriellen Wolfe, Jayshree Zaveri, Karena Zaveri, Josef F. Abril, Roderic Guigo, Michael J. Campbell, Kimmen V. Sjolander, Brian Kartak, Anish Kejarawal, Huaiyu Mi, Betty Lazareva, Thomas Hatton, Apurva Narechania, Karen Diemer, Anushya Muruganujan, Nan Guo, Shinji Sato, Vineet Bafna, Sorin Istrail, Ross Lippert, Russell Schwartz, Brian Walenz, Shibu Yooseph, David Allen, Anand Basu, James Baxendale, Louis Blick, Marcelo Caminha, John Carnes-Stine, Parris Caulk, Yen-Hui Chiang, My Coyne, Carl Dahlke, Anne Deslattes Mays, Maria Dombroski, Michael Donnelly, Dale Ely, Shiva Esparham, Carl Foster, Harold Gire, Stephen Glanowski, Kenneth Glasser, Anna Glodek, Mark Gorokhov, Ken Graham, Barry Gropman, Michael Harris, Jeremy Heil, Scott Henderson, Jeffrey Hoover, Donald Jennings, Catherine Jordan, John Kasha, Leonid Kagan, Cheryl Kraft, Alexander Levitsky, Mark Lewis, Xiangjun Liu, John Lopez, Daniel Ma, William Majoros, Joe McDaniel, Sean Murphy, Matthew Newman, Trung Nguyen, Ngoc Nguyen, Marc Nodell, Sue Pan, Jim Peck, Marshall Peterson, William Rowe, Robert Sanders, John Scott, Michael Simpson, Thomas Smith, Arlan Sprague, Timothy Stockwell, Russell Turner, Eli Venter, Mei Wang, Meiyuan Wen, David Wu, Mitchell Wu, Ashley Xia, Ali Zandieh, Xiaohong Zhu



Initial sequencing and analysis of the human genome

International Human Genome Sequencing Consortium*

Genome Sequencing Centres (Listed in order of total genomic sequence contributed, with a partial list of personnel. A full list of contributors at each centre is available as Supplementary Information.)

Whitehead Institute for Biomedical Research, Center for Genome Research: Eric S. Lander, Lauren M. Linton, Bruce Birren, Chad Nusbaum, Michael C. Zody, Jennifer Baldwin, Ken Devor, Ken Dewar, Michael Doyle, William Fitzhugh, Roel Funke, Diane Gage, Katrina Harris, Andrew Heaford, John Howland, Lisa Kann, Jessica Lehoczky, Rosie Levine, Paul McEwan, Kevin McKernan, James Meldrum, Jill P. Mesirov, Cher Miranda, William Morris, Jerome Naylor, Christina Raymond, Mark Rosetti, Ralph Santos, Andrew Sheridan, Carrie Sougnez, Nicole Stange-Thomann, Nikola Stojanovic, Aravind Subramanian & Dudley Wyman

The Sanger Centre: Jane Rogers, John Sulston, Rachael Ainscough, Stephan Beck, David Bentley, John Burton, Christopher Clee, Nigel Carter, Alan Coulson, Rebecca Deadman, Pamos Deloukas, Andrew Dunham, Ian Dunham, Richard Durbin, Lisa French, Daren Grafham, Simon Gregory, Tim Hubbard, Sean Humphray, Adrienne Hunt, Matthew Jones, Christine Lloyd, Amanda McMurray, Lucy Matthews, Simon Mercer, Sarah Milne, James C. Mullikin, Andrew Mungall, Robert Plumb, Mark Ross, Ratna Showkeen & Sarah Sims

Washington University Genome Sequencing Center: Robert H. Waterston, Richard K. Wilson, LaDeana W. Hillier, John D. McPherson, Marco A. Marra, Elaine R. Mardis, Lucinda A. Fulton, Asif T. Chinwalla, Kimberlie H. Pepin, Warren R. Gish, Stephanie L. Chissoe, Michael C. Wendt, Kim D. Delehaunty, Tracie L. Miner, Andrew Delehaunty, Jason B. Kramer, Lisa L. Cook, Robert S. Fulton, Douglas L. Johnson, Patrick J. Minx & Sandra W. Clifton

US DOE Joint Genome Institute: Trevor Hawkins, Elbert Branscomb, Paul Predki, Paul Richardson, Sarah Wrenning, Tom Slezak, Norman Doggett, Jan-Fang Cheng, Anne Olsen, Susan Lucas, Christopher Elkin, Edward Uberbacher & Marvin Frazier

Baylor College of Medicine Human Genome Sequencing Center: Richard A. Gibbs, Donna M. Muzny, Steven E. Scherer, John B. Bouck, Erica J. Sodergren, Kim C. Worley, Catherine M. Rives, James H. Gornall, Michael L. Metzker, Susan L. Naylor, Raju S. Kuchertapu, David L. Nelson, & George M. Weinstock

RIKEN Genome Sequencing Center: Yoshiyuki Sakaki, Asao Fujiyama, Masahira Hattori, Tetsushi Yada, Atsushi Toyoda, Takehiko Itoh, Chiharu Kawagoe, Hiromi Watanabe, Yasushi Totoki & Todd Taylor

Genoscope and CNRS UMR-8030: Jean Weissenbach, Roland Heilig, William Saurin, Francois Artiguenave, Philippe Brothier, Thomas Bruls, Eric Pelletier, Catherine Robert & Patrick Wincker

GTC Sequencing Center: Douglas R. Smith, Lynn Doucette-Stamm, Marc Rubinfeld, Keith Weinstock, Hong Mei Lee & JoAnn Dubois

Department of Genome Analysis, Institute of Molecular

Biotechnology: André Rosenthal, Matthias Platzer, Gerald Nyakatura, Stefan Taudien & Andreas Rump

Beijing Genomics Institute/Human Genome Center: Kuaning Yang, Jun Yu, Jian Wang, Guyang Huang & Jun Gu

Multimegabase Sequencing Center, The Institute for Systems Biology: Leroy Hood, Lee Rowen, Anup Madan & Shizen Qin

Stanford Genome Technology Center: Ronald W. Davis, Nancy A. Federspiel, A. Pia Abola & Michael J. Proctor

Stanford Human Genome Center: Richard M. Myers, Jeremy Schmutz, Mark Dickson, Jane Grimwood & David R. Cox

University of Washington Genome Center: Maynard V. Olson, Rajinder Kaul & Christopher Raymond

Department of Molecular Biology, Keio University School of Medicine: Nobuyoshi Shimizu, Kazuhiko Kawasaki & Shinsai Minoshima

University of Texas Southwestern Medical Center at Dallas: Gen A. Evans, Martha Athanasiou & Roger Schultz

University of Oklahoma's Advanced Center for Genome Technology: Bruce A. Roe, Feng Chen & Huaqin Pan

Max Planck Institute for Molecular Genetics: Juliane Ramser, Hans Lehrach & Richard Reinhardt

Cold Spring Harbor Laboratory, Lita Annenberg Hazen Genome Center: W. Richard McCombie, Melissa de la Bastide & Neilay Dedia

GBF—German Research Centre for Biotechnology: Helmut Blöcker, Klaus Homischer & Gabriele Nordisk

* Genome Analysis Group (listed in alphabetical order, also includes individuals listed under other headings): Richa Agarwala, L. Aravind, Jeffrey A. Bailey, Alex Bateman, Serafim Batzoglou, Ewan Birney, Peer Bork, Daniel G. Brown, Christopher B. Burge, Lorenzo Corutti, Hsiu-Chuan Chen, Deanna Church, Michele Clamp, Richard R. Copley, Tobias Doerk, Sean R. Eddy, Evan E. Eichler, Terrence S. Furey, James Galagan, James G. R. Gilbert, Cyrus Harmon, Yoshitake Hayashizaki, David Haussler, Henning Hermjakob, Karsten Hokamp, Wonhee Jung, L. Steven Johnson, Thomas A. Jones, Simon Kashef, Ark Kasparyk, Scot Kennedy, W. James Kent, Paul Kitts, Eugene V. Koonin, Ian Korf, David Kulp, Doron Lancet, Todd M. Lowe, Aoife McLysaght, Tarjei Mikkelser, John V. Moran, Nicola Mulder, Victor J. Pollara, Chris P. Ponting, Greg Schulz, Jörg Schultz, Guy Slater, Arian F. A. Smit, Elia Stupka, Joseph Szustakowski, Danielle Thierry-Mieg, Jean Thierry-Mieg, Lukas Wagner, John Wallis, Raymond Wheeler, Alan Williams, Yut I. Wolf, Kenneth H. Wolfe, Shiaw-Pyng Yang & Ru-Fang Yeh

Scientific management: National Human Genome Research Institute, US National Institutes of Health: Francis Collins, Mark S. Guyer, Jane Peterson, Adam Felsenfeld & Kris A. Wetterstrand; Office of Science, US Department of Energy: Anstoes Patrinos; The Wellcome Trust: Michael J. Morgan

RESEARCH

Open Access

Sequencing and analysis of an Irish human genome

Pin Tong^{1†}, James GD Prendergast^{2†}, Amanda J Lohan¹, Susan M Farrington^{2,3}, Simon Cronin⁴, Nial Friel⁵, Dan G Bradley⁶, Orla Hardiman⁷, Alex Evans⁸, James F Wilson⁹, Brendan Loftus^{1*}

Abstract

Background: Recent studies generating complete human sequences from Asian, African and European subgroups have revealed population-specific variation and disease susceptibility loci. Here, choosing a DNA sample from a population of interest due to its relative geographical isolation and genetic impact on further populations, we extend the above studies through the generation of 11-fold coverage of the first Irish human genome sequence.

Results: Using sequence data from a branch of the European ancestral tree as yet unsequenced, we identify variants that may be specific to this population. Through comparisons with HapMap and previous genetic association studies, we identified novel disease-associated variants, including a novel nonsense variant putatively associated with inflammatory bowel disease. We describe a novel method for improving SNP calling accuracy at low genome coverage using haplotype information. This analysis has implications for future re-sequencing studies and validates the imputation of Irish haplotypes using data from the current Human Genome Diversity Cell Line Panel (HGDP-CEPH). Finally, we identify gene duplication events as constituting significant targets of recent positive selection in the human lineage.

Conclusions: Our findings show that there remains utility in generating whole genome sequences to illustrate both general principles and reveal specific instances of human biology. With increasing access to low cost sequencing we would predict that even armed with the resources of a small research group a number of similar initiatives geared towards answering specific biological questions will emerge.



Whole-genome sequencing and comprehensive variant analysis of a Japanese individual using massively parallel sequencing

Akihiro Fujimoto, Hidewaki Nakagawa, Naoya Hosono, Kaoru Nakano, Tetsuo Abe, Keith A Boroevich, Masao Nagasaki, Rui Yamaguchi, Tetsuo Shibuya, Michiaki Kubo, Satoru Miyano, Yusuke Nakamura & Tatsuhiko Tsunoda

[Affiliations](#) | [Contributions](#) | [Corresponding authors](#)

Nature Genetics 42, 931–936 (2010) | doi:10.1038/ng.691

Received 18 February 2010 | Accepted 10 September 2010 | Published online 24 October 2010

Abstract

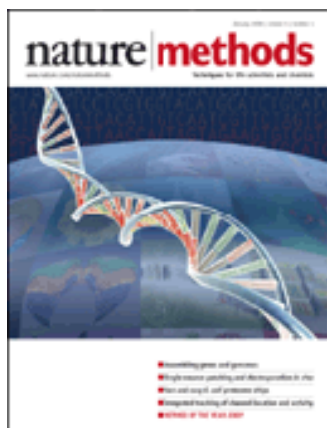
[Abstract](#) • [Author information](#) • [Supplementary information](#)

We report the analysis of a Japanese male using high-throughput sequencing to $\times 40$ coverage. More than 99% of the sequence reads were mapped to the reference human genome. Using a Bayesian decision method, we identified 3,132,608 single nucleotide variations (SNVs). Comparison with six previously reported genomes revealed an excess of singleton nonsense and nonsynonymous SNVs, as well as singleton SNVs in conserved non-coding regions. We also identified 5,319 deletions smaller than 10 kb with high accuracy, in addition to copy number variations and rearrangements. *De novo* assembly of the unmapped sequence reads generated around 3 Mb of novel sequence, which showed high similarity to non-reference human genomes and the human herpesvirus 4 genome. Our analysis suggests that considerable variation remains undiscovered in the human genome and that whole-genome sequencing is an invaluable tool for obtaining a complete understanding of human genetic variation.

Next Generation Sequencing Technology

Massively Parallel Signature Sequencing (MPSS)

NATURE METHODS | VOL.5 NO.1 | JANUARY 2008



Nature Methods' Method of the Year 2007 goes to next-generation sequencing. This series of articles showcase how these novel sequencing methods came into their own in 2007 and the incredible impact they promise to have in a variety of research applications. The Methods to Watch feature provide a glimpse and a wish list for future Methods of the Year.



454 GS FLX

<http://www.454.com/>

illumina® **SOLEXA**

<http://www.illumina.com/pages.ilmn?ID=250>

AB applied biosystems™ **SOLID**

<http://marketing.appliedbiosystems.com/>

Throughput of NGS machines (2007-2009)

| | | | | | | | | |
|---------------------|--------|--------|--------|----------|--------|--------|--------|--------|
| Vendor: | Roche | | | Illumina | | | ABI | |
| Technology: | 454 | | | Solexa | | | SOLiD | |
| Platform: | GS 20 | FLX | Ti | GA | GA II | 1 | 2 | |
| Reads: | 500 k | 500 k | 1 M | 28 M | 80 M | 40 M | 115 M | |
| Fragment | | | | | | | | |
| Read length: | 100 | 200 | 350 | 35 | 50 | 75 | 25 | 35 |
| Run time: | 6 hr | 7 hr | 9 hr | 3 d | 3 d | 4 d | 6 d | 5 d |
| Yield: | 50 Mb | 100 Mb | 400 Mb | 1 Gb | 4 Gb | 6 Gb | 1 Gb | 4 Gb |
| Images: | 11 GB | 13 GB | 27 GB | 500 GB | 1.1 TB | 1.7 TB | 1.8 TB | 2.5 TB |
| PA Disk: | 3 GB | 3 GB | 15 GB | 175 GB | 300 GB | 350 GB | 300 GB | 750 GB |
| PA CPU: | 10 hr | 140 hr | 220 hr | 100 hr | 70 hr | 100 hr | NA | NA |
| SRA: | 500 MB | 1 GB | 4 GB | 30 GB | 50 GB | 75 GB | 100 GB | 140 GB |
| Paired-end | | | | | | | | |
| Read length: | | 200 | | 2×35 | 2×50 | 2×75 | 2×25 | 2×35 |
| Insert: | | 3.5 kb | | 200 b | 200 b | 200 b | 3 kb | 3 kb |
| Run time: | | 7 hr | | 6 d | 6 d | 8 d | 12 d | 10 d |
| Yield: | | 100 Mb | | 2 Gb | 8 Gb | 11 Gb | 2 Gb | 8 Gb |
| Images: | | 13 GB | | 1 TB | 2.2 TB | 3.4 TB | 3.6 TB | 5 TB |
| PA Disk: | | 3 GB | | 350 GB | 500 GB | 600 GB | 600 GB | 1.5 TB |
| PA CPU: | | 140 hr | | 160 hr | 120 hr | 170 hr | NA | NA |
| SRA: | | 1 GB | | 60 GB | 100 GB | 150 GB | 200 GB | 280 GB |

Throughput of NGS machines (2010)

| | | | | | | | | |
|---------------------|--------|--------|--------|----------|--------|--------|--------|--------|
| Vendor: | Roche | | | Illumina | | | ABI | |
| Technology: | 454 | | | Solexa | | | SOLiD | |
| Platform: | GS 20 | FLX | Ti | GA | GA II | 1 | 2 | |
| Reads: | 500 k | 500 k | 1 M | 28 M | 80 M | 40 M | 115 M | |
| Fragment | | | | | | | | |
| Read length: | 100 | 200 | 350 | 35 | 50 | 75 | 25 | 35 |
| Run time: | 6 hr | 7 hr | 9 hr | 3 d | 3 d | 4 d | 6 d | 5 d |
| Yield: | 50 Mb | 100 Mb | 400 Mb | 1 Gb | 4 Gb | 6 Gb | 1 Gb | 4 Gb |
| Images: | 11 GB | 13 GB | 27 GB | 500 GB | 1.1 TB | 1.7 TB | 1.8 TB | 2.5 TB |
| PA Disk: | 3 GB | 3 GB | 15 GB | 175 GB | 300 GB | 350 GB | 300 GB | 750 GB |
| PA CPU: | 10 hr | 140 hr | 220 hr | 100 hr | 70 hr | 100 hr | NA | NA |
| SRA: | 500 MB | 1 GB | 4 GB | 30 GB | 50 GB | 75 GB | 100 GB | 140 GB |
| Paired-end | | | | | | | | |
| Read length: | 200 | | 2x35 | 2x50 | 2x75 | 2x25 | 2x35 | |
| Insert: | 3.5 kb | | 200 b | 200 b | 200 b | 3 kb | 3 kb | |
| Run time: | 7 hr | | 6 d | 6 d | 8 d | 12 d | 10 d | |
| Yield: | 100 Mb | | 2 Gb | 8 Gb | 11 Gb | 2 Gb | 8 Gb | |
| Images: | 13 GB | | 1 TB | 2.2 TB | 3.4 TB | 3.6 TB | 5 TB | |
| PA Disk: | 3 GB | | 350 GB | 500 GB | 600 GB | 600 GB | 1.5 TB | |
| PA CPU: | 140 hr | | 160 hr | 120 hr | 170 hr | NA | NA | |
| SRA: | 1 GB | | 60 GB | 100 GB | 150 GB | 200 GB | 280 GB | |



HiSeq2000 (launched in 2010)

| Read Length | Run Time | Output |
|-------------|-----------|------------|
| 1 x 35 bp | ~1.5 days | 26-35 Gb |
| 2 x 50 bp | ~4 days | 75-100 Gb |
| 2 x 100 bp | ~8 days | 150-200 Gp |

Up to 25GB per day for a 2 x100 bp run



SOLiD 4 (launched in 2009)

NEXT GENERATION SEQUENCING

**Technologies,
Applications &
Analysis**

NGS

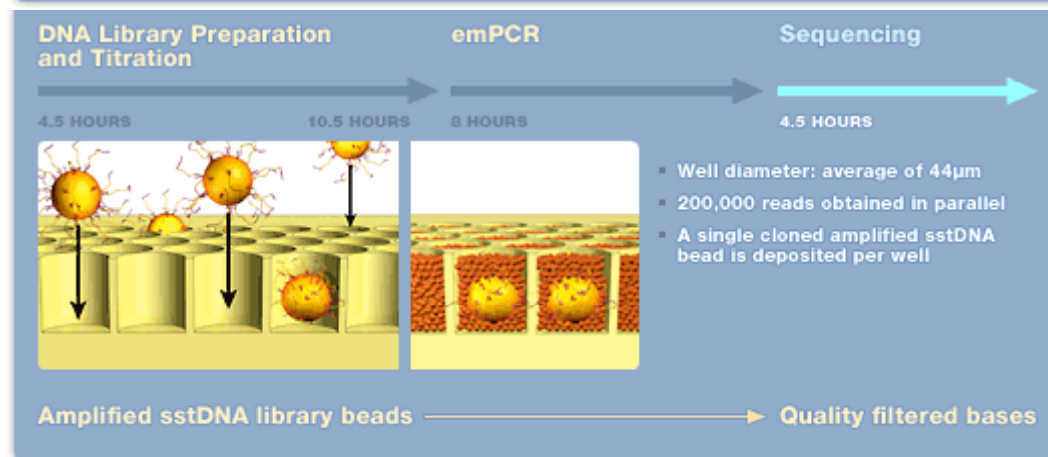
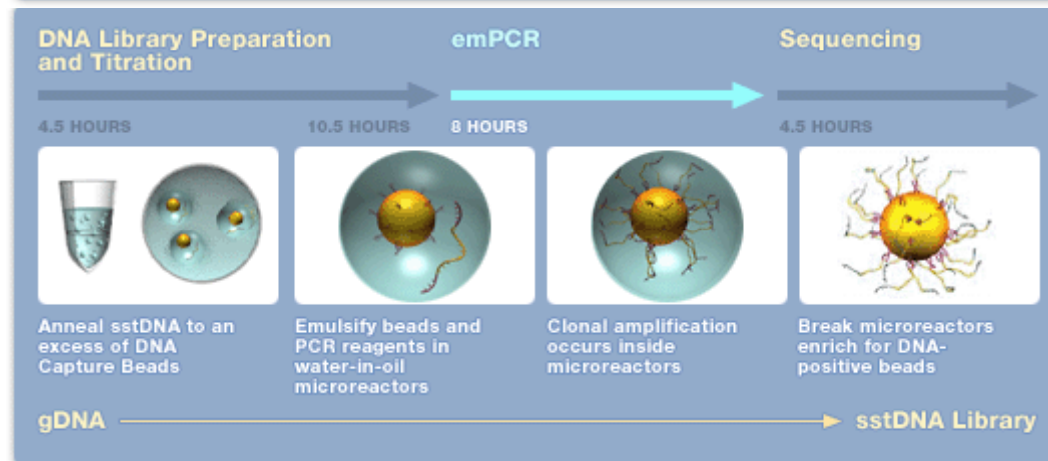
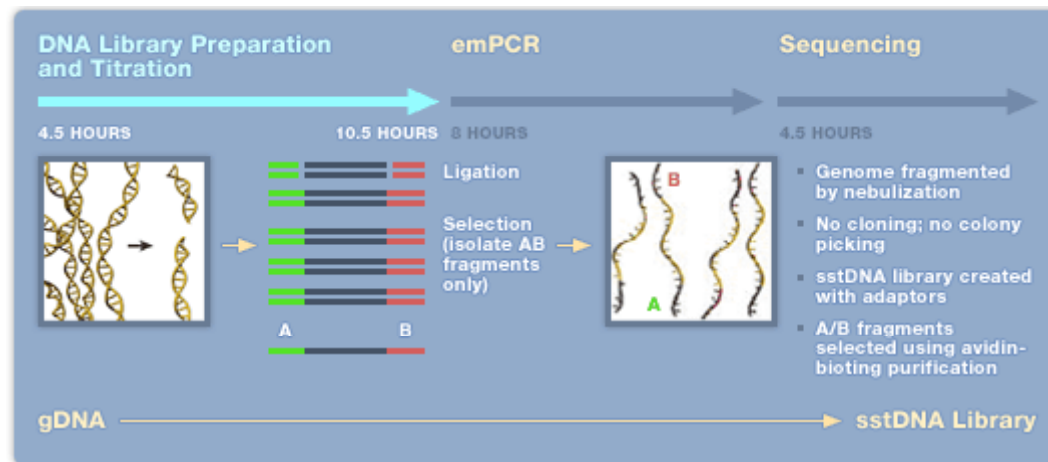
| | Feature generation | Sequencing by synthesis |
|--------|--------------------|--|
| 454 | Emulsion PCR | Polymerase (pyrosequencing) |
| Solexa | Bridge PCR | Polymerase (reversible terminators) |
| SOLiD | Emulsion PCR | Ligase (octamers with two-base encoding) |



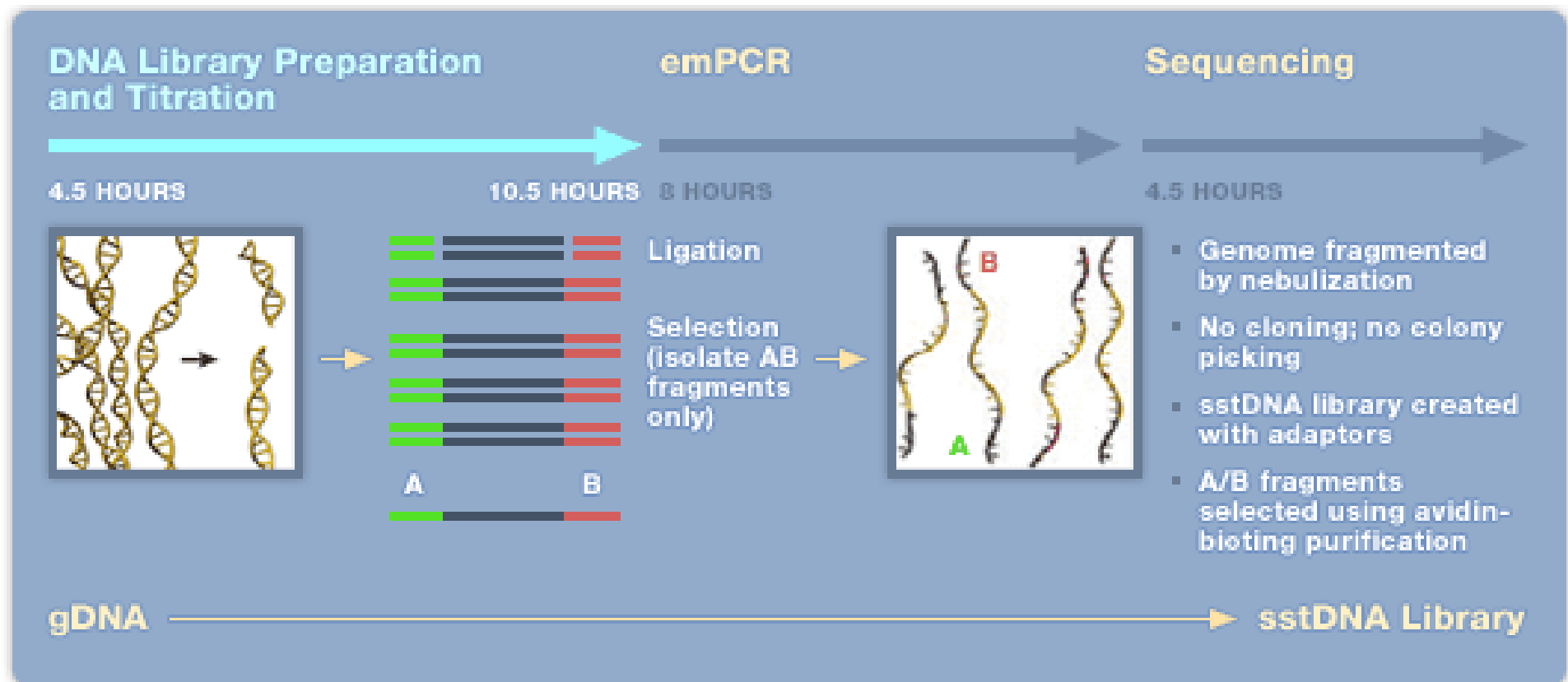
454 SEQUENCING

| | |
|-----------------------|--|
| June 2000 | 454 Life Sciences is founded |
| October 2005 | Release of the Genome Sequencer 20, the first next-generation sequencing system on the market |
| October 2005 | Collaboration agreement signed with Roche Diagnostics |
| December 2005 | 454 Life Sciences Awarded the Wall Street Journal's Gold Medal for Innovation |
| November 2006 | 454 Life Sciences, in collaboration with Svante Paabo, describes in <i>Nature</i> the first million base pairs of the Neanderthal genome and initiates the Neanderthal Genome Project. |
| January 2007 | Release of the Genome Sequencer FLX System |
| March 2007 | Roche Diagnostics completes integration with 454 Life Sciences |
| May 2007 | Complete sequence of Jim Watson published in <i>Nature</i> . First genome to be sequenced for less than \$1 million. |
| November 2007 | Announcement of the 100th peer-reviewed publication enabled by 454 Sequencing |
| June 2008 | 454 Joins the 1000 Genome Project, an international effort to build the most detailed map to date of human genetic variation as a tool for medical research |
| September 2008 | Announcement of the 250th peer-reviewed publication enabled by 454 Sequencing |
| October 2008 | Release of Genome Sequencer FLX Titanium Series reagents, featuring 1 million reads at 400 base pairs in length |

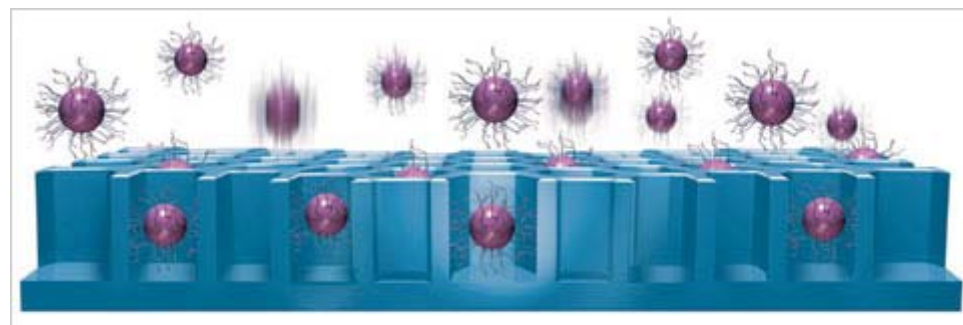
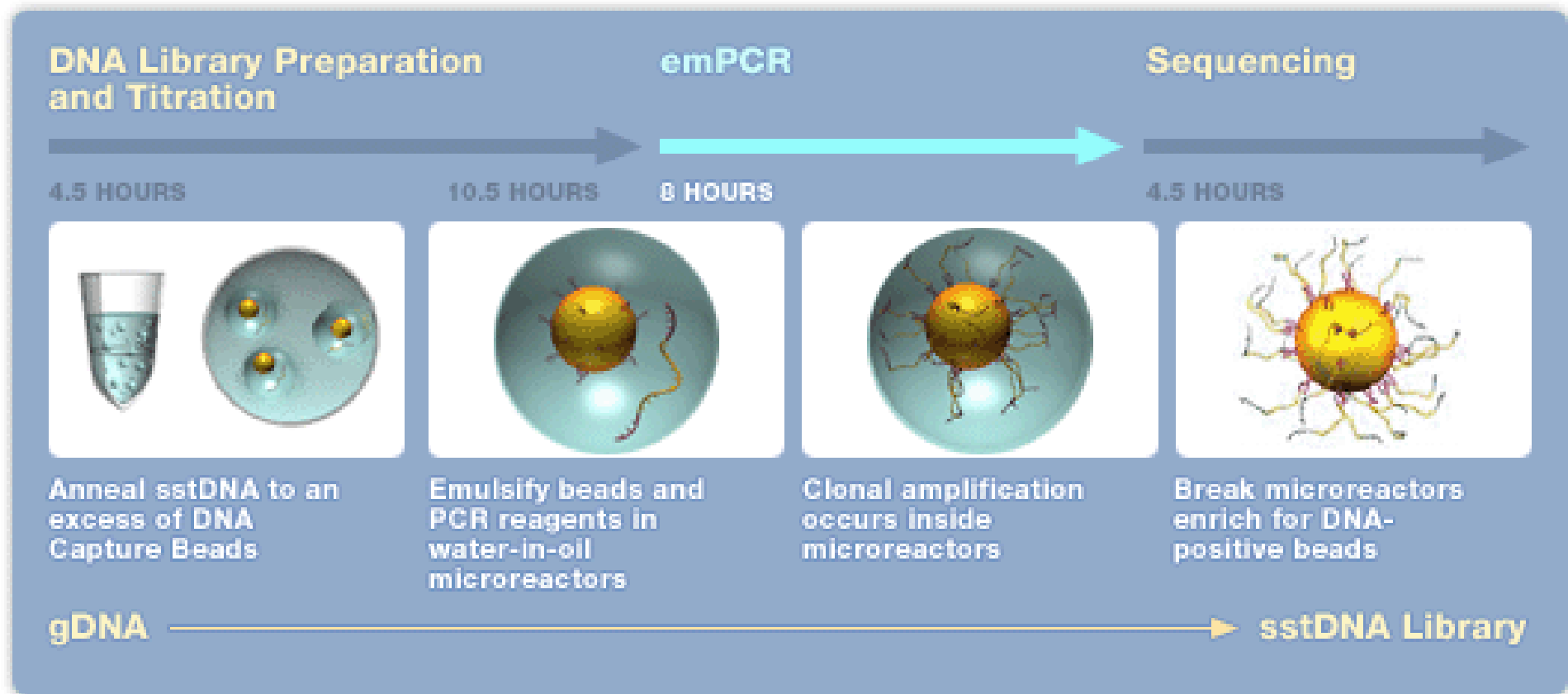




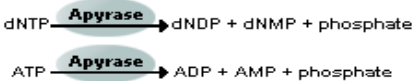
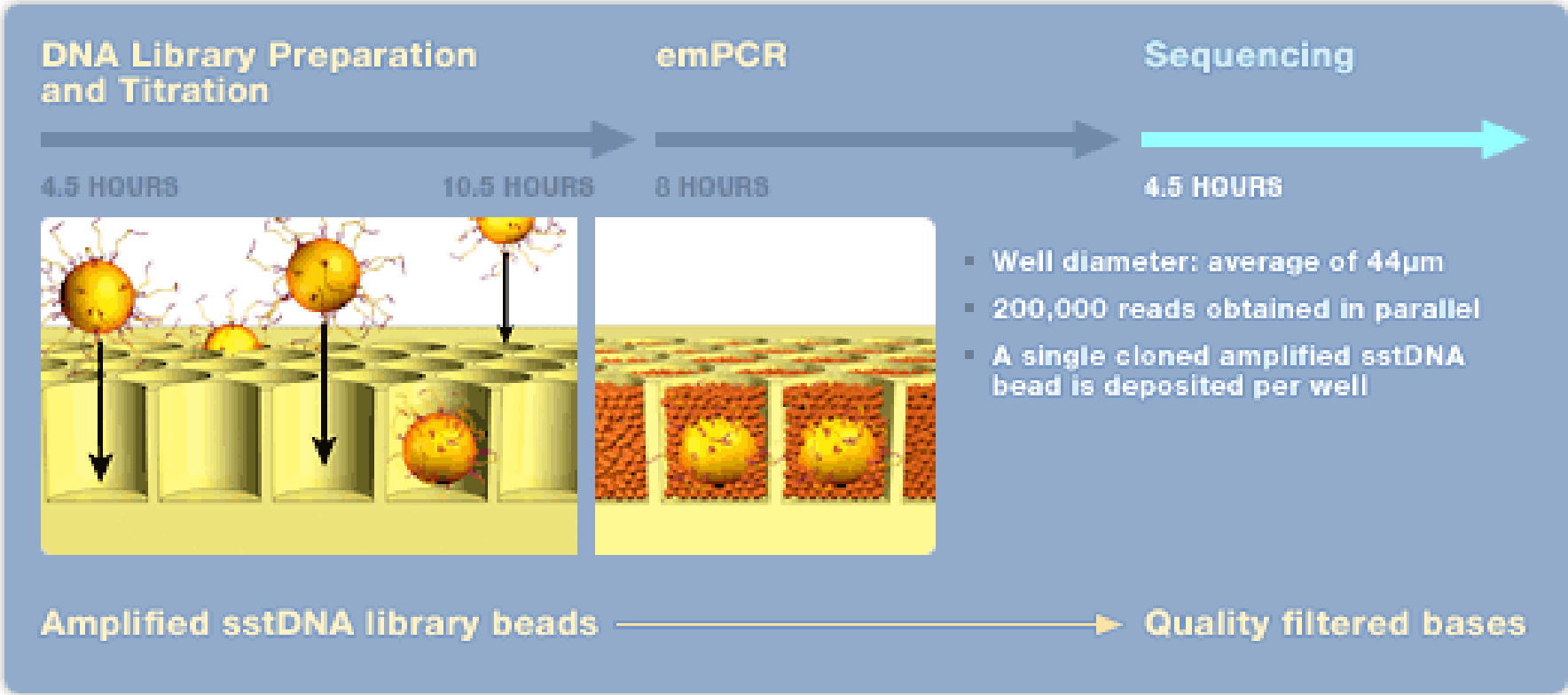
Step 1. DNA Library Preparation



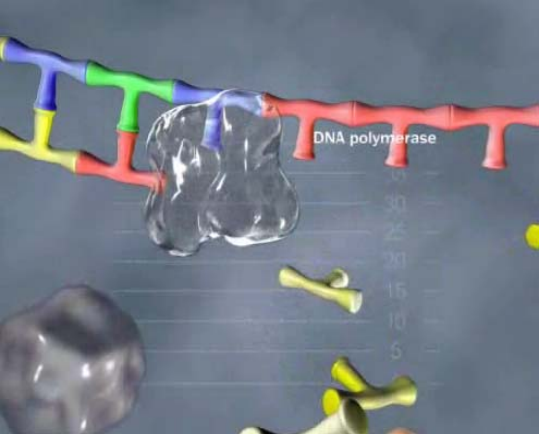
Step 2. Emulsion PCR (emPCR)



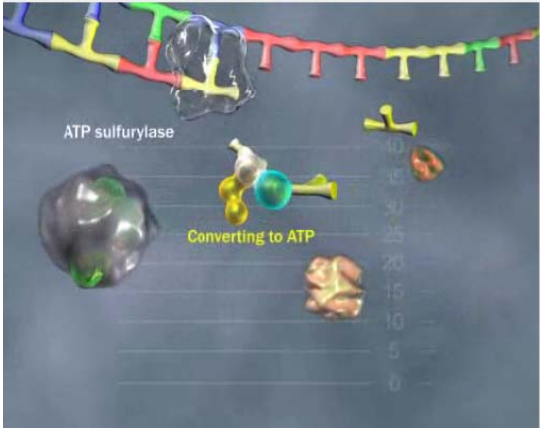
Step 3. Pyrosequencing



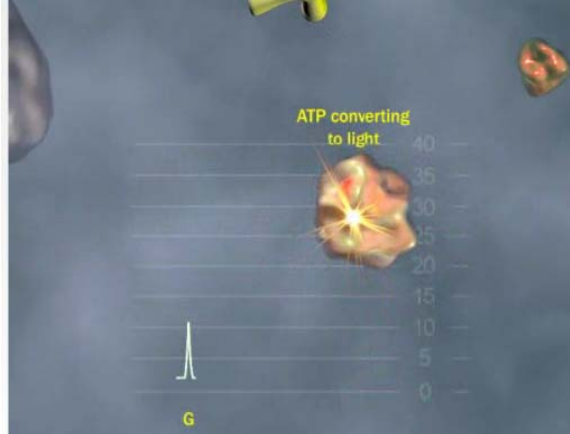
result → read the signal of light



DNA polymerase → add the A.T.C.G

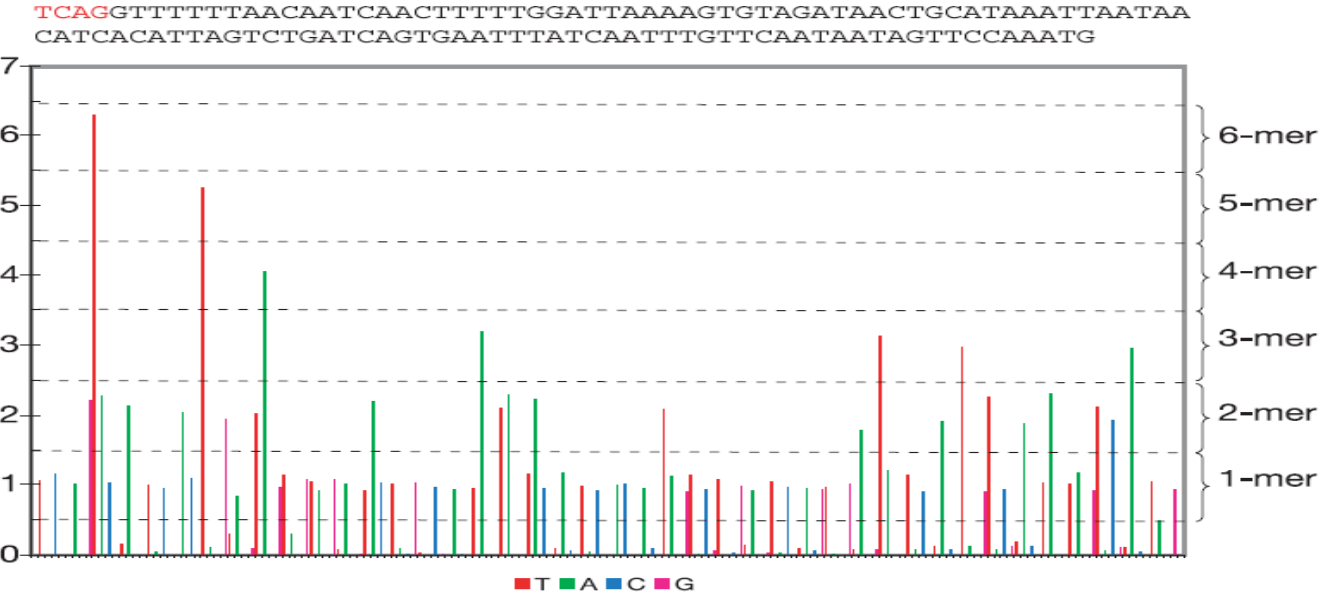


ATP synthase → convert pyrophosphate to ATP



luciferase → react the ATP with luciferin to generate light

apyrase → degrade unincorporated dNTPs and excess ATP



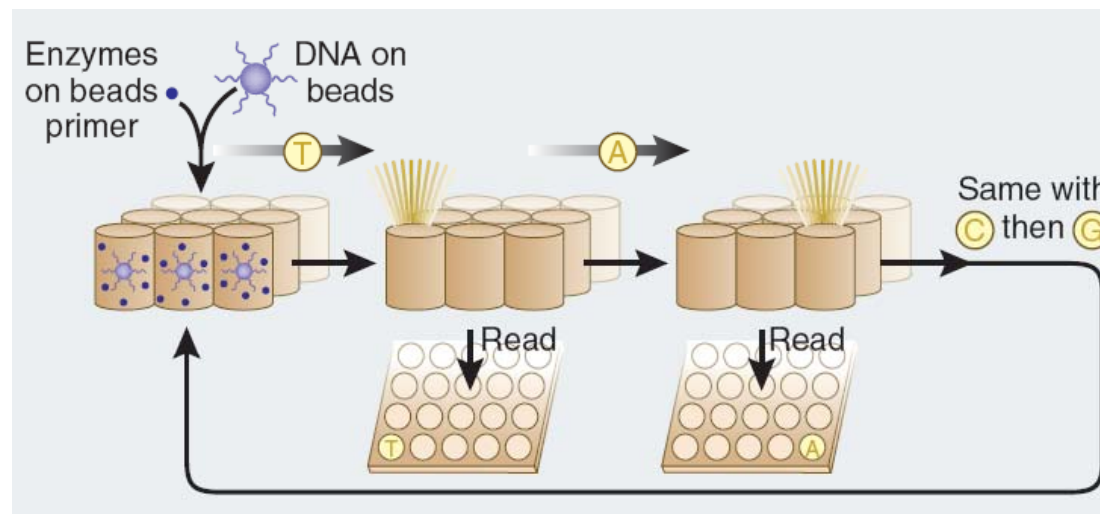
454 technology (Pyrosequencing)

Sample preparation.

Fragments of DNA are ligated to adapters that facilitate their capture on beads (one fragment per bead). A water-in-oil emulsion containing PCR reagents and one bead per droplet is created to amplify each fragment individually in its droplet. After amplification, the emulsion is broken, DNA is denatured and the beads, containing one amplified DNA fragment each, are distributed into the wells of a fiber-optic slide.

Pyrosequencing.

The wells are loaded with sequencing enzymes and primer (complementary to the adapter on the fragment ends), then exposed to a flow of one unlabeled nucleotide at a time, allowing synthesis of the complementary strand of DNA to proceed. When a nucleotide is incorporated, pyrophosphate is released and converted to ATP, which fuels the luciferase-driven conversion of luciferin to oxyluciferin and light. As a result, the well lights up. The read length is between 100 and 150 nucleotides.



First look

Raising the standards in sequencing. What's coming next from 454.

[See the Future of 454 Sequencing](#)

1 2 3 4 5

News

August 5, 2009

454 Life Sciences and Roche NimbleGen Announce Collaboration with Eli Lilly and SeqWright to Sequence Genomic Regions Associated with Psychiatric Disease.

July 29, 2009

Roche and Google.org Start Initiative for Early Discovery of New Diseases

July 16, 2009

Researchers to Decode Antarctic Ice Metagenome with the 454 Sequencing System, to Explore the Effects of Climate Change

[» Read more news](#)

Publications

5 2 5 and counting

> The novel polysaccharide deacetylase homolog Pdi contributes to virulence of the aquatic pathogen *Streptococcus iniae*. *Microbiology*.

> Genomic Diversity and Evolution of *Mycobacterium ulcerans* Revealed by Next-Generation Sequencing. *PLoS Pathogens*.

> Recognition and coupling of A-to-I edited sites are determined by the tertiary structure of the RNA. *Nucleic Acids Research*.

[» Read more publications](#)



Stimulus Funding Grant Support

[» Read more](#)

Contact Us / Support

Get in touch with us directly! Select your area of interest:

Select



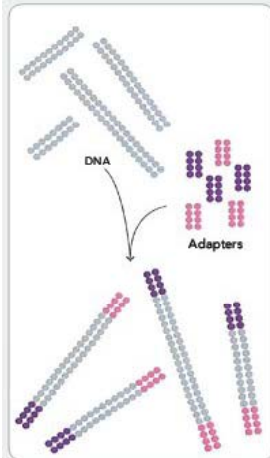
Solexa Genome Analyzer





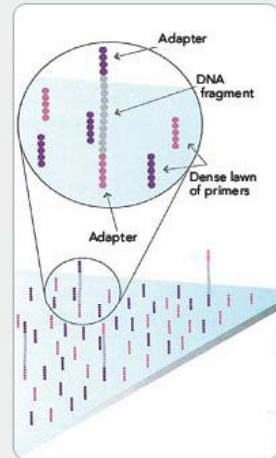
Solexa Genome Analyzer

1. PREPARE GENOMIC DNA SAMPLE



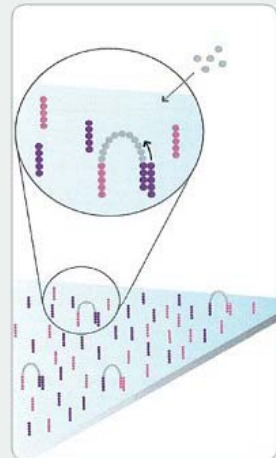
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



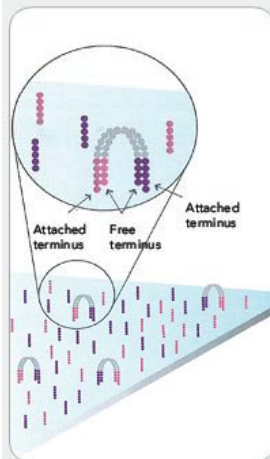
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



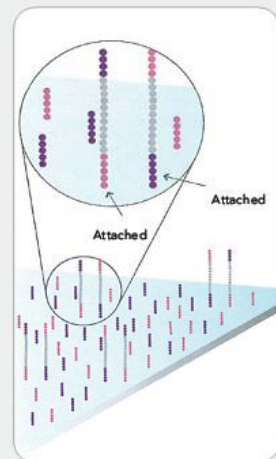
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

4. FRAGMENTS BECOME DOUBLE STRANDED



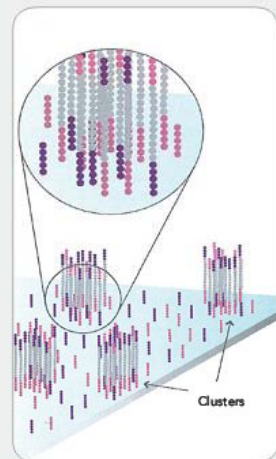
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



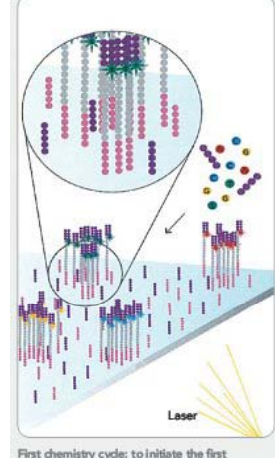
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



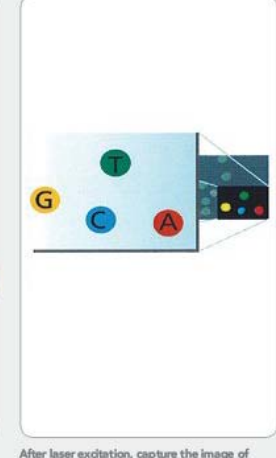
Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

7. DETERMINE FIRST BASE



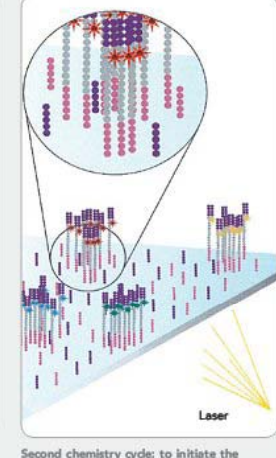
First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

8. IMAGE FIRST BASE



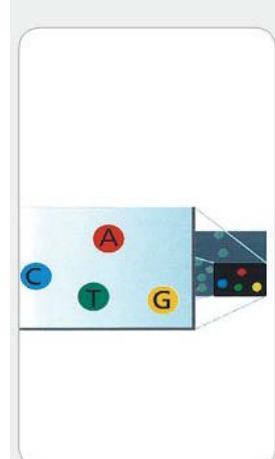
After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

9. DETERMINE SECOND BASE



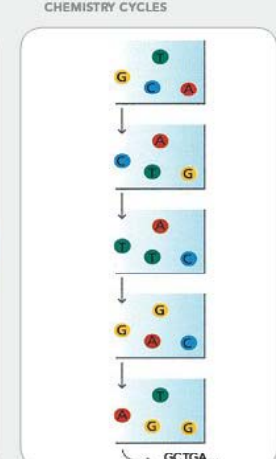
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

10. IMAGE SECOND CHEMISTRY CYCLE



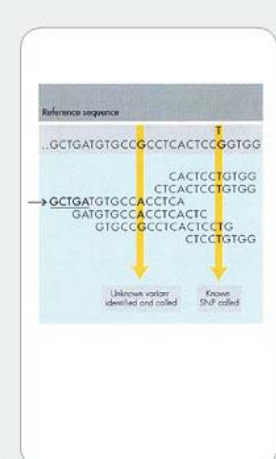
After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA

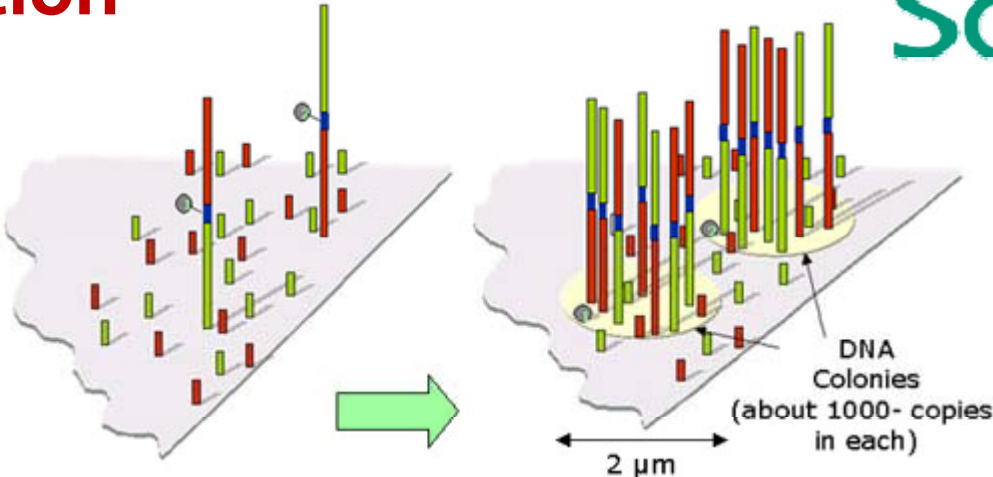


Align data, compare to a reference, and identify sequence differences.

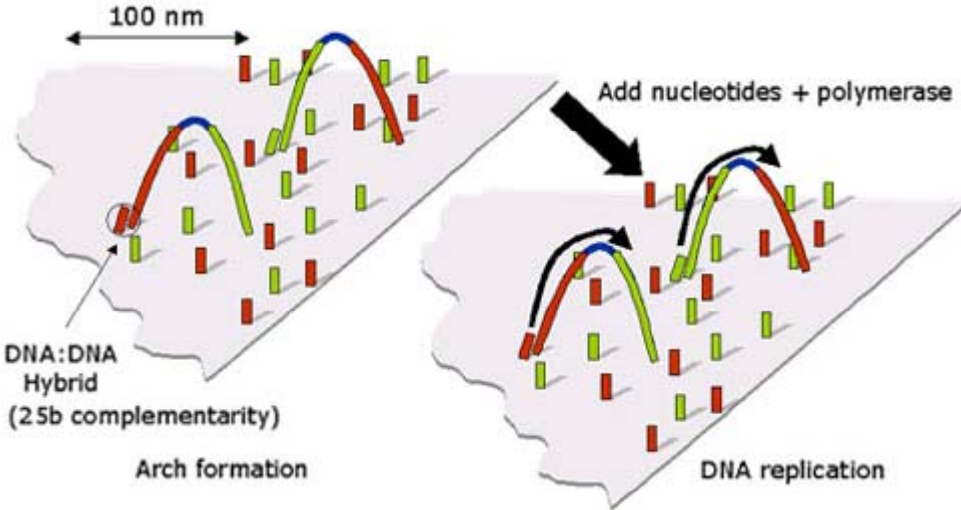
Cluster Generation



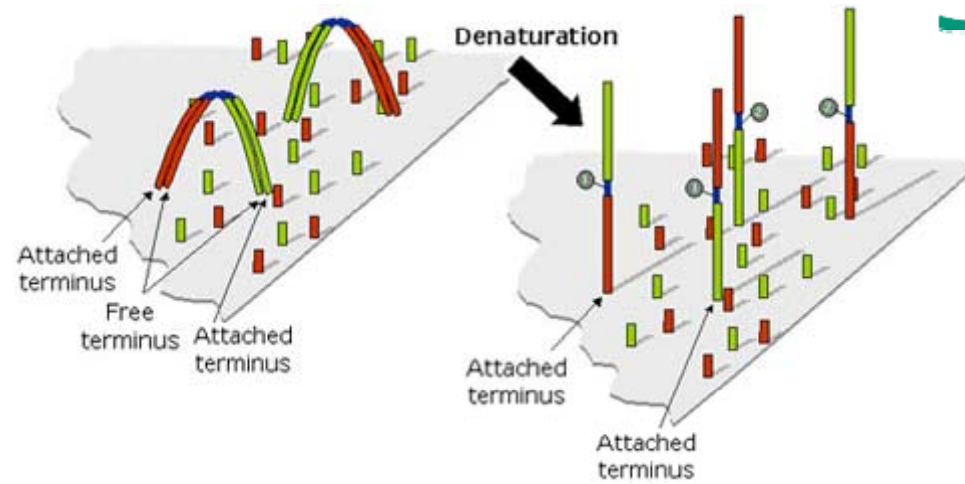
Clusters Array



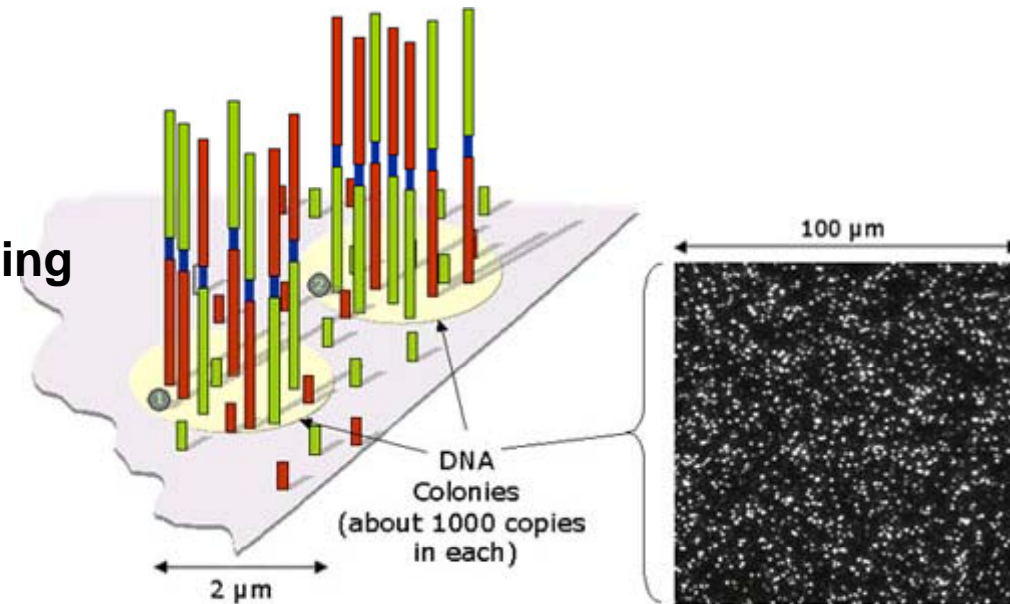
Template Bridging



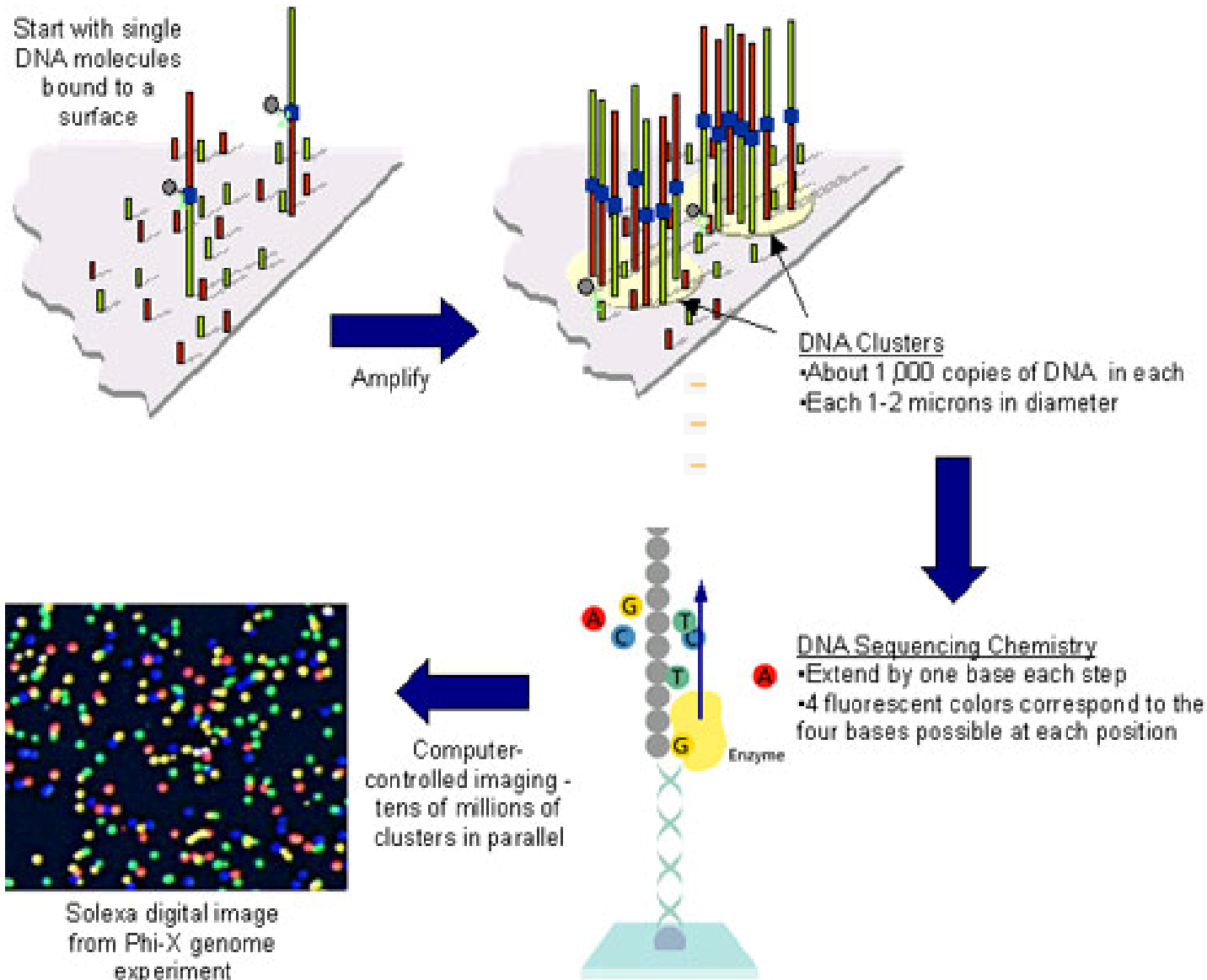
Denaturation

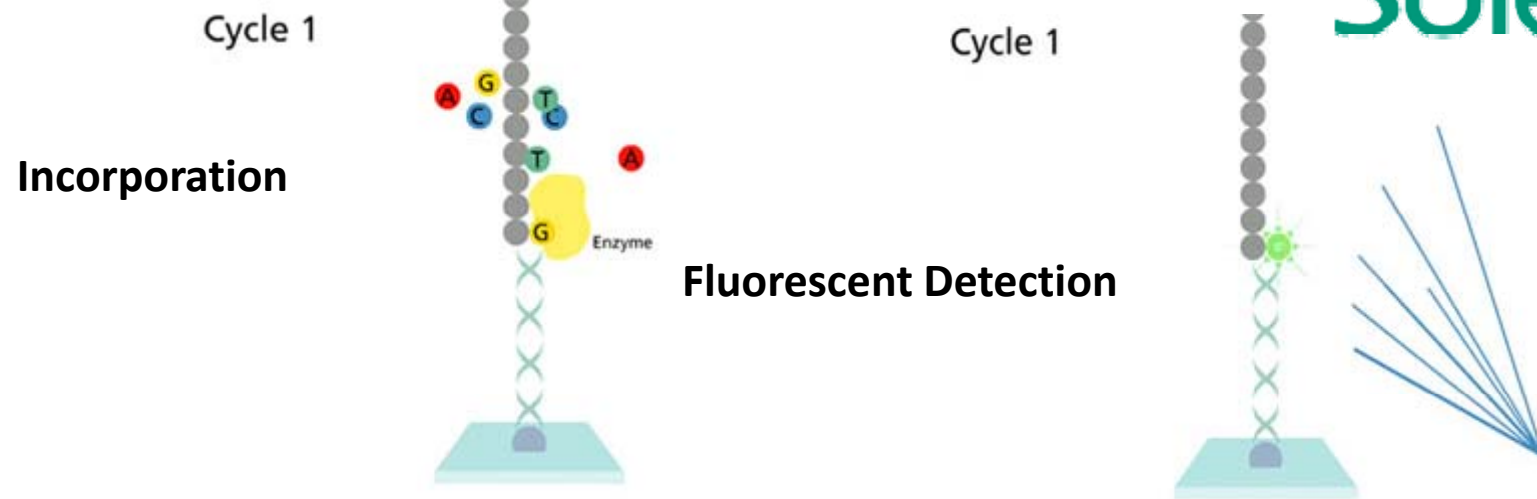


Ready For Sequencing

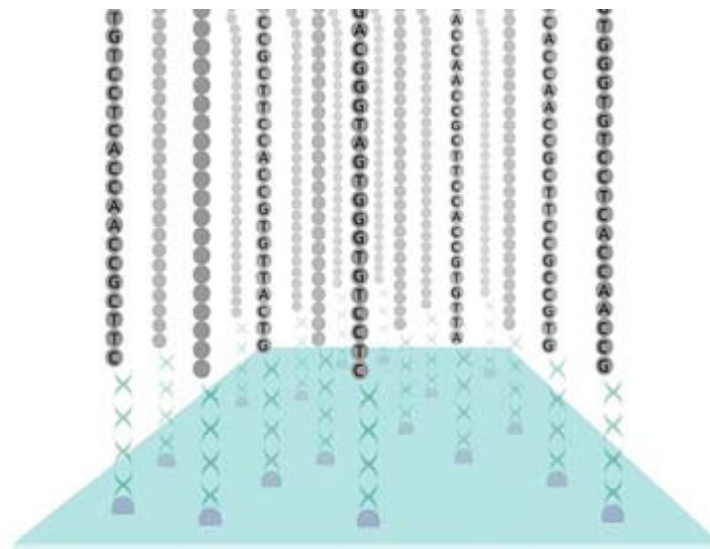


Sequencing-By-Synthesis

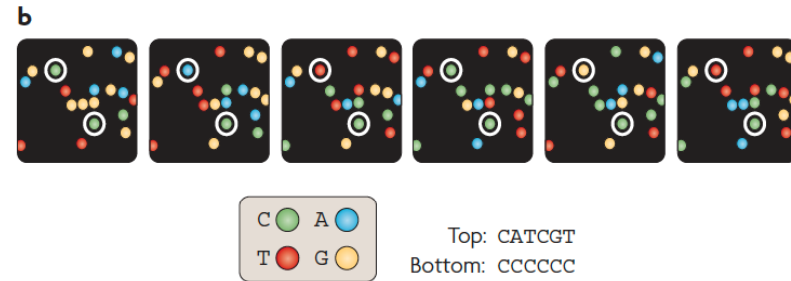
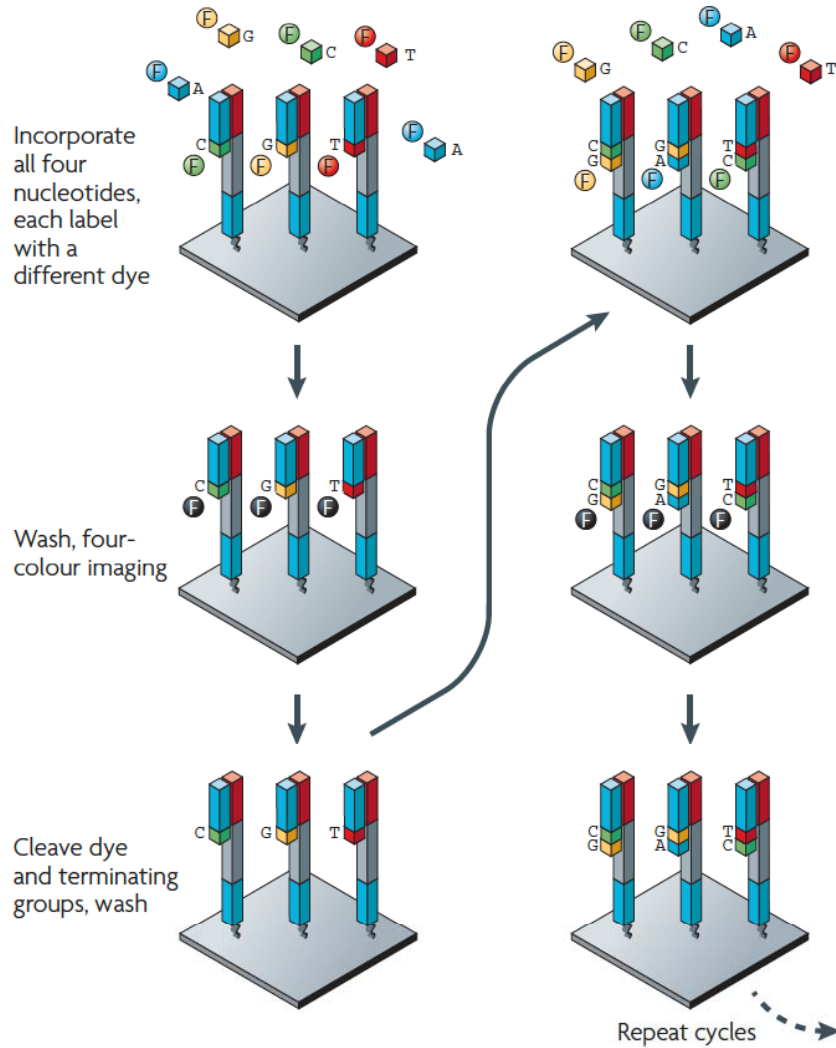




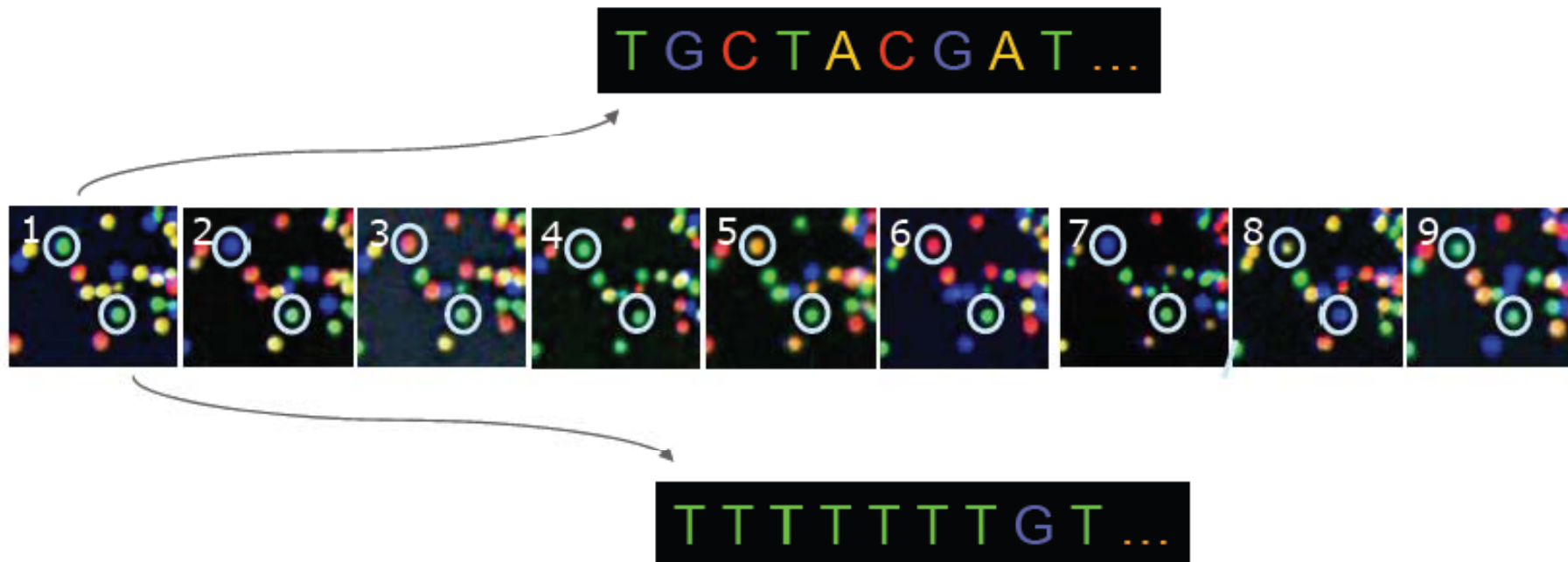
Sequence Generated At Every Site On The Array



a Illumina/Solexa — Reversible terminators



Base calling from raw data



The identity of each base of a cluster is read off from sequential images



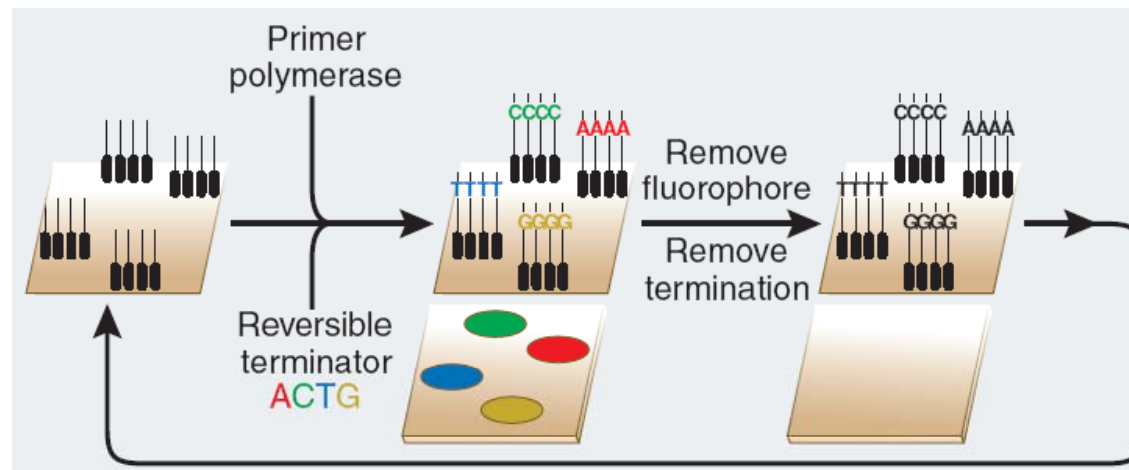
Solexa technology (sequencing-by-synthesis)

Sample preparation.

Fragments of DNA are ligated to end adapters, denatured and bound at one end to a solid surface already coated with a dense layer of the adapters. Each single stranded fragment is immobilized at one end, while its free end 'bends over' and hybridizes to a complementary adapter on the surface, which initiates the synthesis of the complementary strand in the presence of amplification reagents. Multiple cycles of this solid-phase amplification followed by denaturation create clusters of ~1,000 copies of single-stranded DNA molecules distributed randomly on the surface.

Sequencing with reversible terminators.

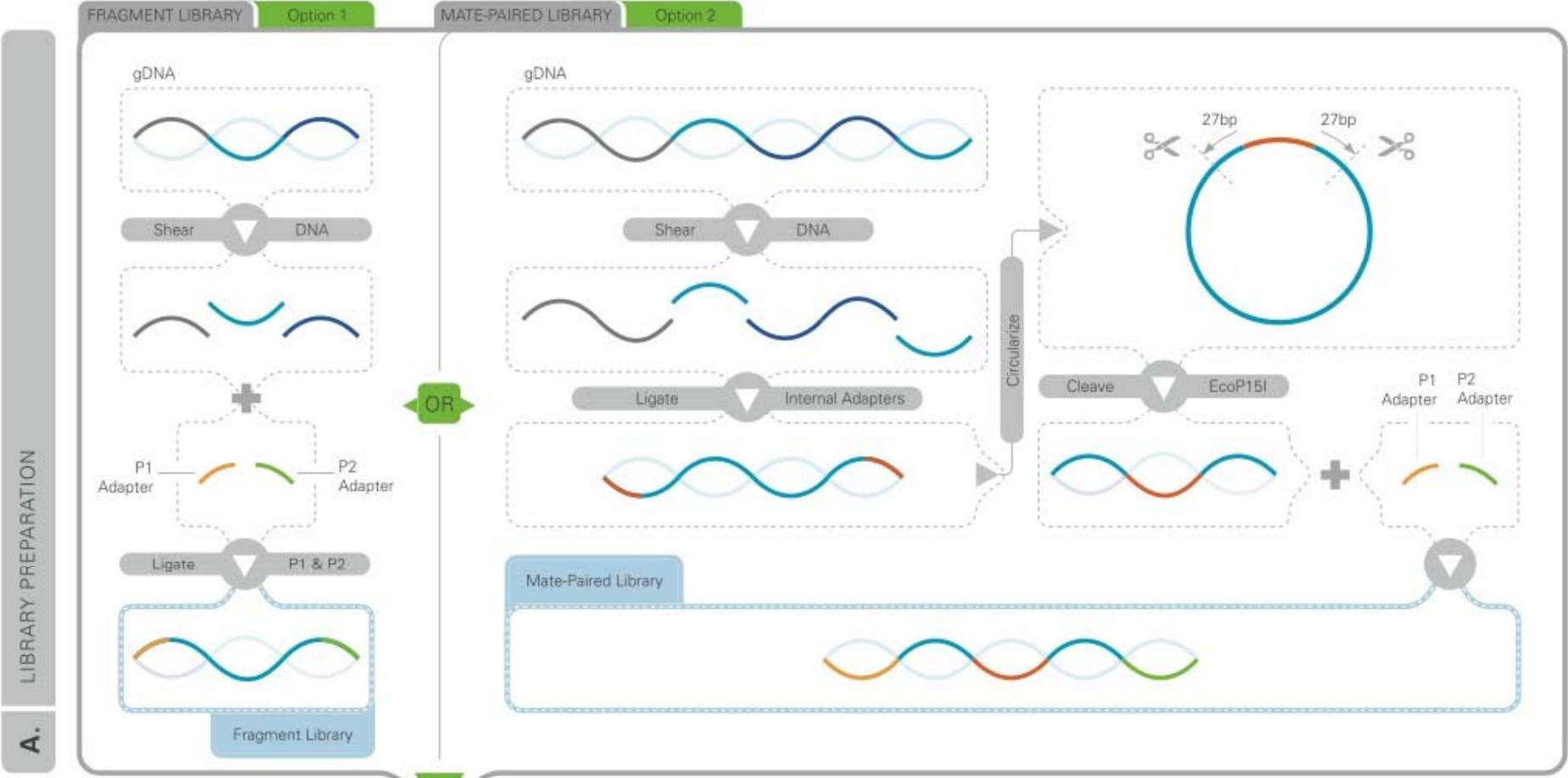
Synthesis reagents, added to the flow cell, consist of primers, DNA polymerase and four differently labeled, reversible terminator nucleotides. After incorporation of a nucleotide, which is identified by its color, the 3' terminator on the base and the fluorophore are removed, and the cycle is repeated for a read length of 30–35 nucleotides.



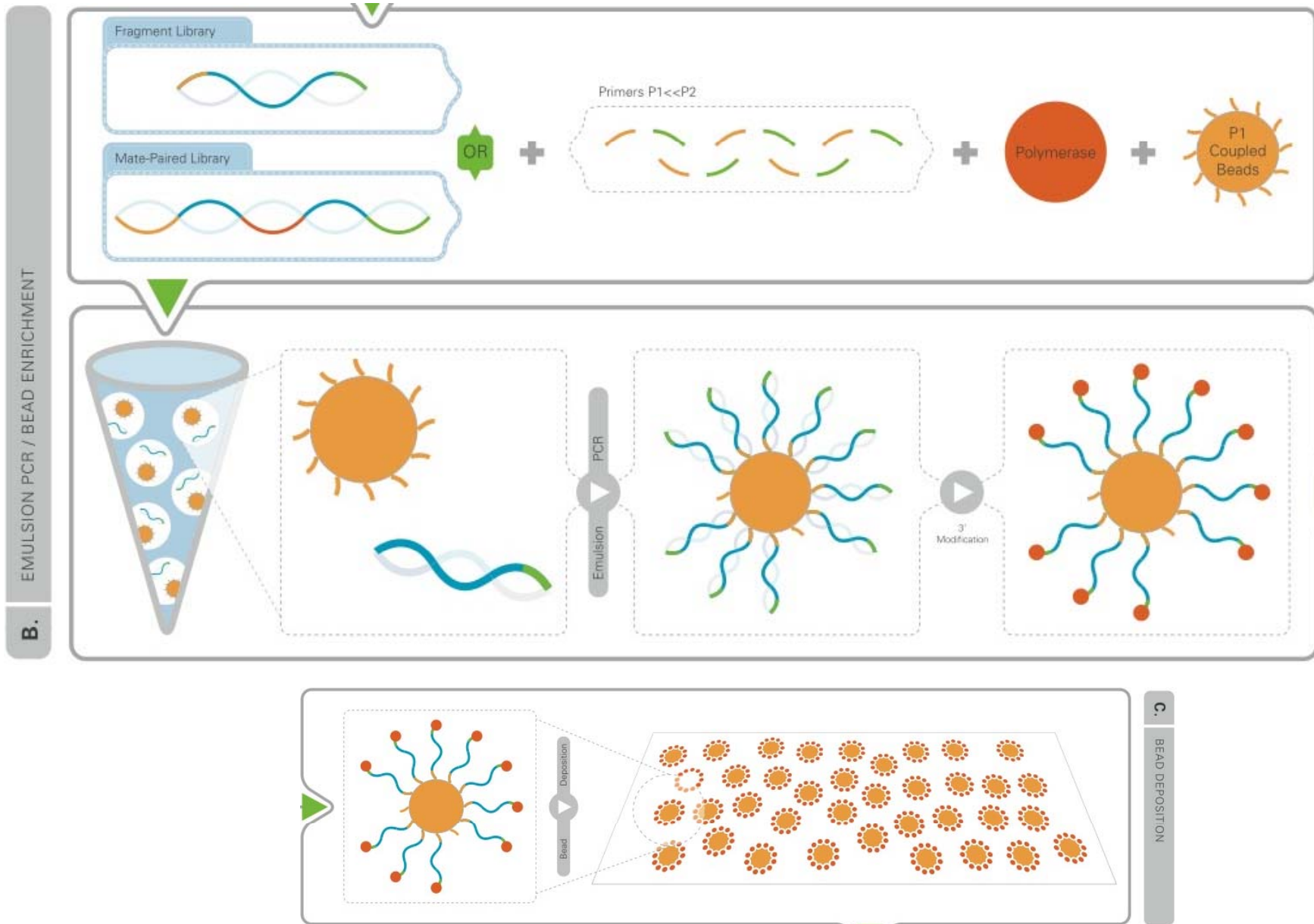
Applied Biosystems SOLiD™ System 2.0



Step 1. Library preparation

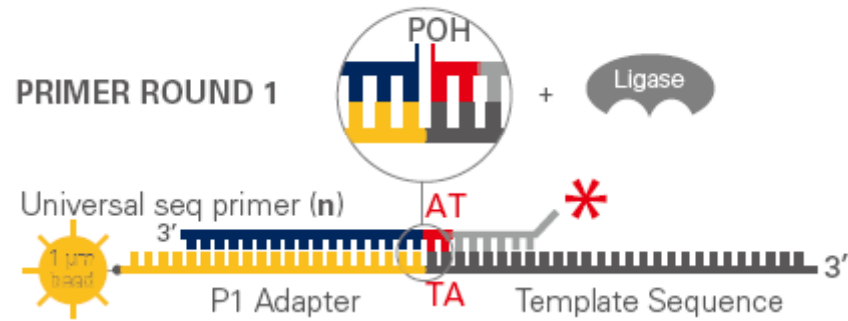


Step 2. Emulsion PCR

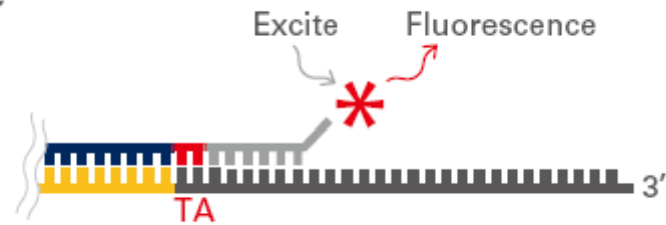


Step 3. Sequencing-by-Ligation

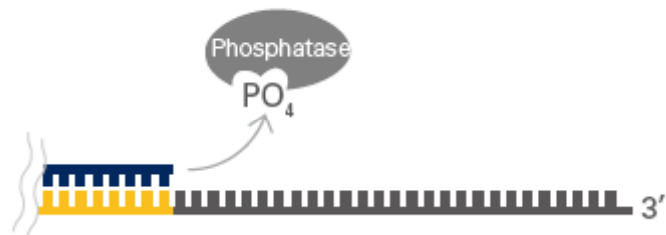
1. Prime and Ligate



2. Image



3. Cap Unextended Strands



4. Cleave off Fluor

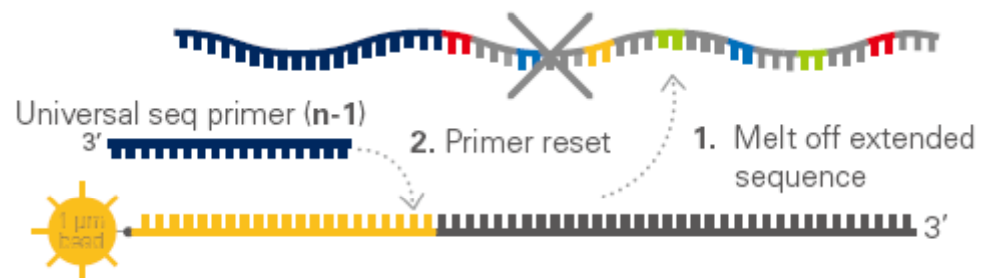


5. Repeat steps 1-4 to Extend Sequence

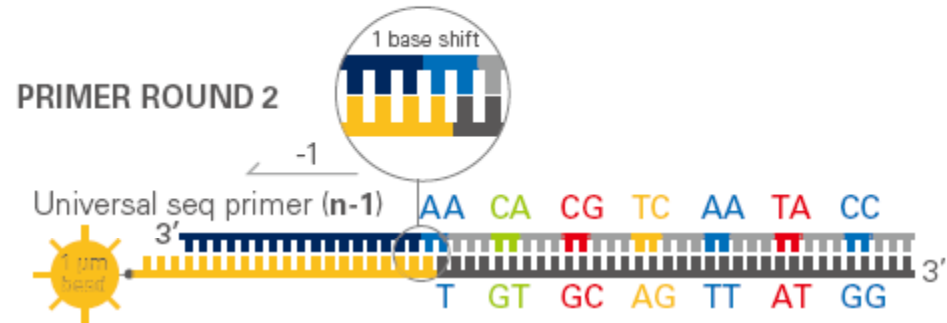
Ligation cycle 1 2 3 4 5 6 7 ... (n cycles)



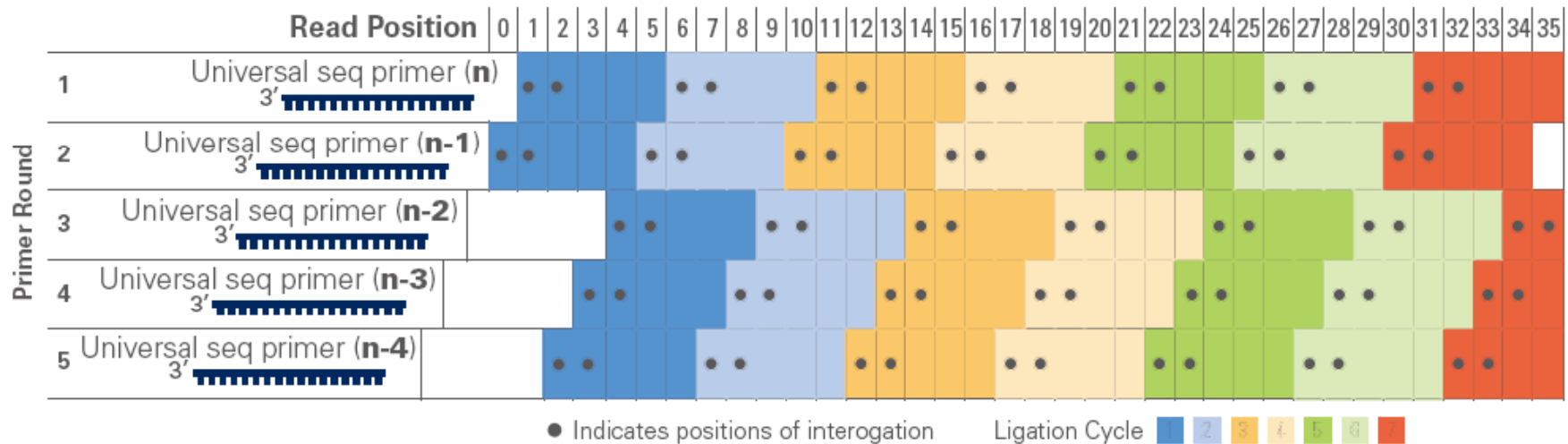
6. Primer Reset



7. Repeat steps 1-5 with new primer



8. Repeat Reset with , n-2, n-3, n-4 primers



Its format is

```
>1_88_1830_R3
```

```
G32113123201300232320
```

>TAG_ID

```
>1_89_1562_R3
```

Color_space

```
G23133131233333101320
```

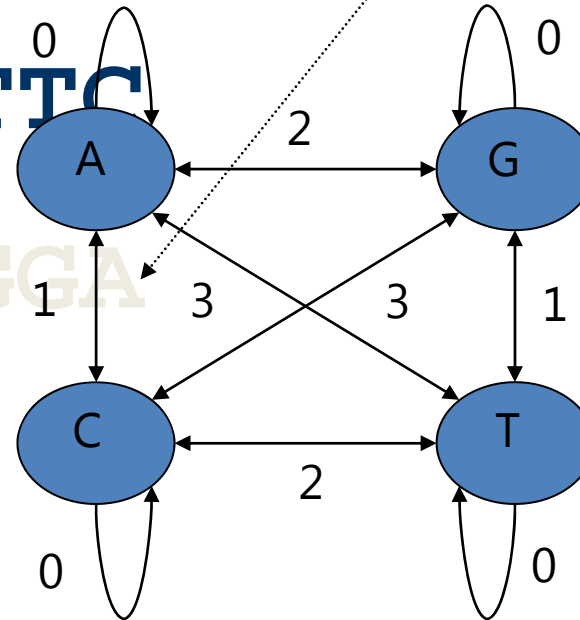

AB SOLiD: Dibase Sequencing

AB SOLiD reads look like this:

T012233102
TGAGCGTTC

T012033102
TGAATAGGA

| | A | C | G | T |
|---|---|---|---|---|
| A | 0 | 1 | 2 | 3 |
| C | 1 | 0 | 3 | 2 |
| G | 2 | 3 | 0 | 1 |
| T | 3 | 2 | 1 | 0 |



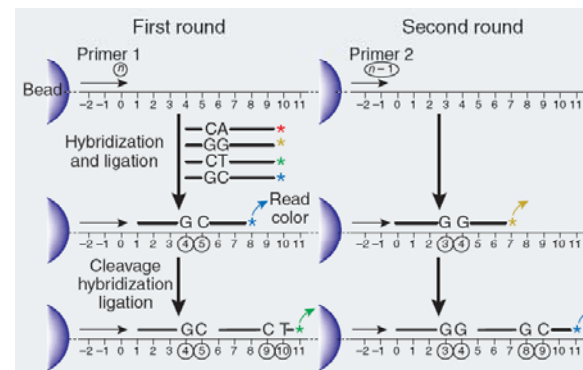
SOLiD technology (sequencing-by-ligation)

Sample preparation.

Fragments of DNA are ligated to adapters and amplified on beads by emulsion PCR. The DNA is denatured and the beads deposited onto a glass slide.

Sequencing by ligation.

A sequencing primer is hybridized to the adapter and its 5' end is available for ligation to an oligonucleotide hybridizing to the adjacent sequence. A mixture of octamer oligonucleotides compete for ligation to the primer (the bases in fourth and fifth position on these oligos are encoded by one of four color labels). After its color has been recorded, the ligated oligonucleotide is cleaved between position 5 and 6, which removes the label, and the cycle of ligation-cleavage is repeated. In the first round, the process determines possible identities of bases in positions 4, 5, 9, 10, 14, 15, etc. The entire process is repeated, offset by one base by using a shorter sequencing primer, to determine positions 3, 4, 8, 9, 13, 14, etc., until the first base in the sequencing primer (position 0) is reached. Since the identity of this base is known, the color is used to decode its neighboring base at position 1, which in turn decodes the base at position 2, etc., until all sequence pairs are identified. The current read length is between 30 and 35 nucleotides



Evolution of Sequencing Technology



Sanger dideoxy-sequencing

ABI 3730XL

Massive parallel sequencing

Roche 454 FLX, Illumina Genome Analyzer, Life Technologies SOLiD

Bead-based em-PCR and sequencing by ligation

Dover Systems' Polonator

Massive parallel sequencing and single molecule sequencing

Pacific Biosciences (single-molecule real-time DNA sequencing (SMRT) technology)

Helicos (true single-molecule-sequencing (tSMS) technology)

VisiGen Biotechnologies (real-time, single-molecule sequencing fluorescence resonance energy transfer (FRET) technology)

Single molecule sequencing and nanopore technology?

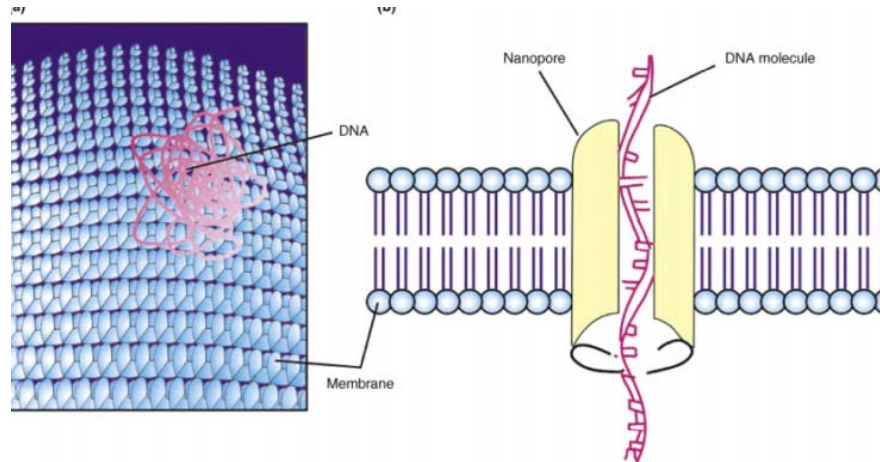
Oxford Nanopore Technologies (label-free, single-molecule sequencing (BASE) technology), Affymetrix, Reveo, Base4innovation, Genome Corp, and Complete Genomics.

Nanopore Sequencing

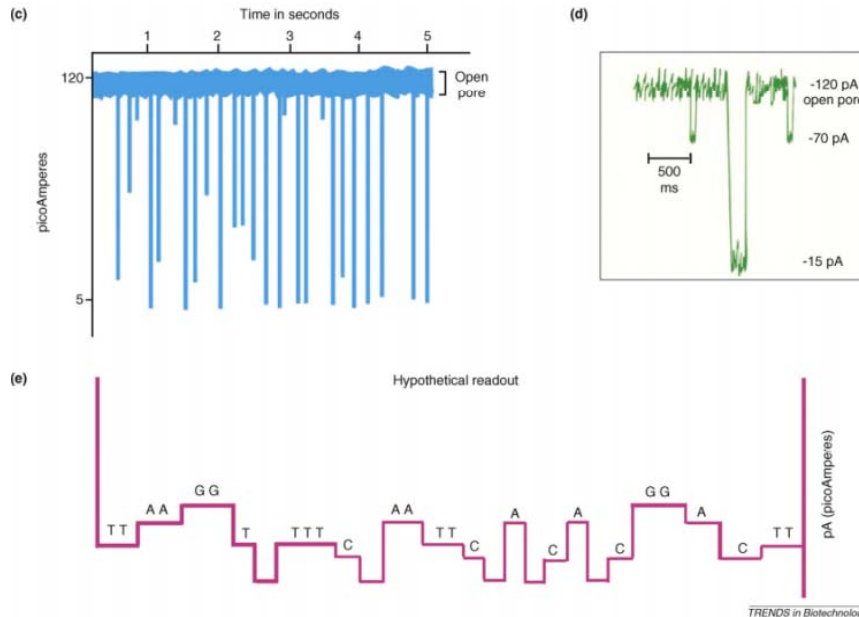
Ion Torrent

Oxford Nanopore

Nucleic acids driven through a nanopore.



Differences in conductance of pore provide readout.



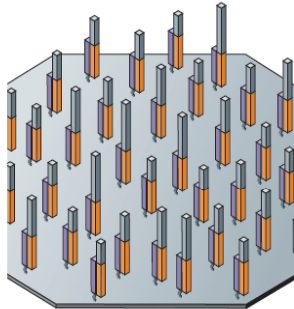
Single Molecule Real-Time Sequencing

Real-time monitoring of PCR activity

Read-out by fluorescence resonance energy transfer between polymerase and nucleotides or

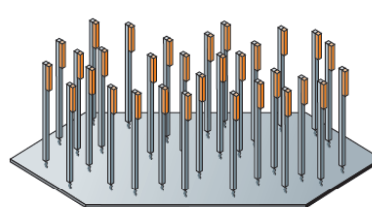
Waveguides allow direct observation of polymerase and fluorescently labeled nucleotides

c Helicos BioSciences: one-pass sequencing
Single molecule: primer immobilized



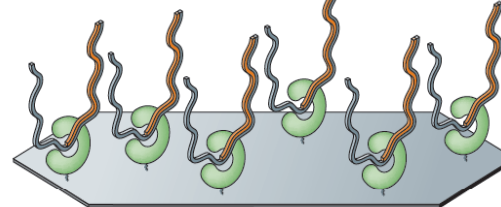
Billions of primed, single-molecule templates

d Helicos BioSciences: two-pass sequencing
Single molecule: template immobilized



Billions of primed, single-molecule templates

e Pacific Biosciences, Life/Visigen, LI-COR Biosciences
Single molecule: polymerase immobilized



Thousands of primed, single-molecule templates

Helicos Biosciences

VisGen Biotechnologies

Pacific Biosciences

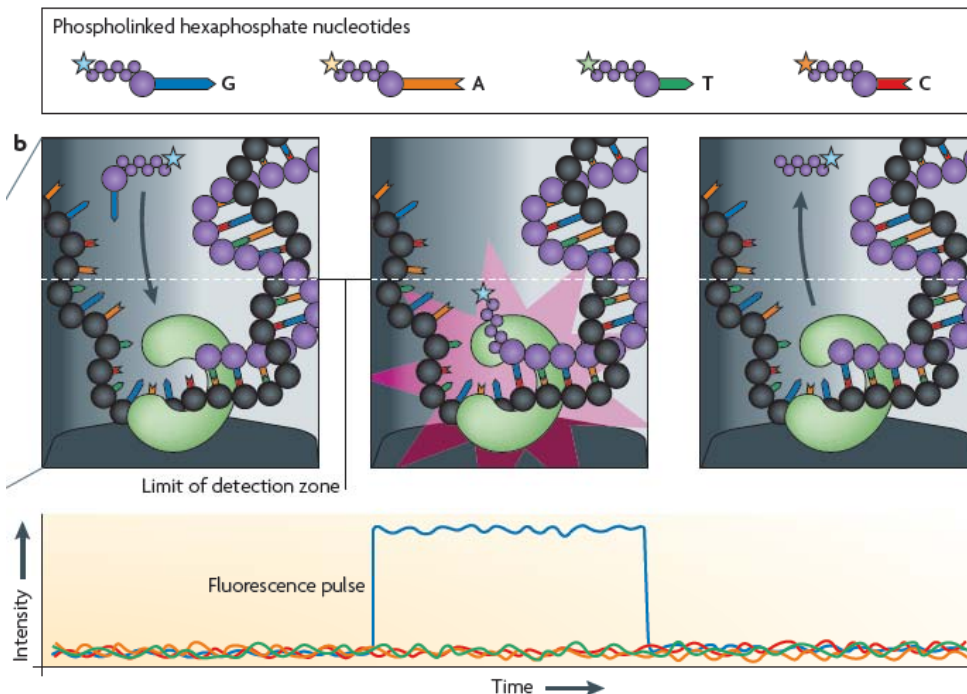
Single Molecule Real-Time (SMRT) Sequencing

www.sciencemag.org SCIENCE VOL 323 2 JANUARY 2009

133

Real-Time DNA Sequencing from Single Polymerase Molecules

John Eid,* Adrian Fehr,* Jeremy Gray,* Khai Luong,* John Lyle,* Geoff Otto,* Paul Peluso,* David Rank,* Primo Baybayan, Brad Bettman, Arkadiusz Bibillo, Keith Bjornson, Bidhan Chaudhuri, Frederick Christians, Ronald Cicero, Sonya Clark, Ravindra Dalal, Alex deWinter, John Dixon, Mathieu Foquet, Alfred Gaertner, Paul Hardenbol, Cheryl Heiner, Kevin Hester, David Holden, Gregory Kearns, Xiangxu Kong, Ronald Kuse, Yves Lacroix, Steven Lin, Paul Lundquist, Congcong Ma, Patrick Marks, Mark Maxham, Devon Murphy, Insil Park, Thang Pham, Michael Phillips, Joy Roy, Robert Sebra, Gene Shen, Jon Sorenson, Austin Tomaney, Kevin Travers, Mark Trulson, John Veceli, Jeffrey Wegener, Dawn Wu, Alicia Yang, Denis Zaccarin, Peter Zhao, Frank Zhong, Jonas Korlach,† Stephen Turner†



Pacific Biosciences — Real-time sequencing

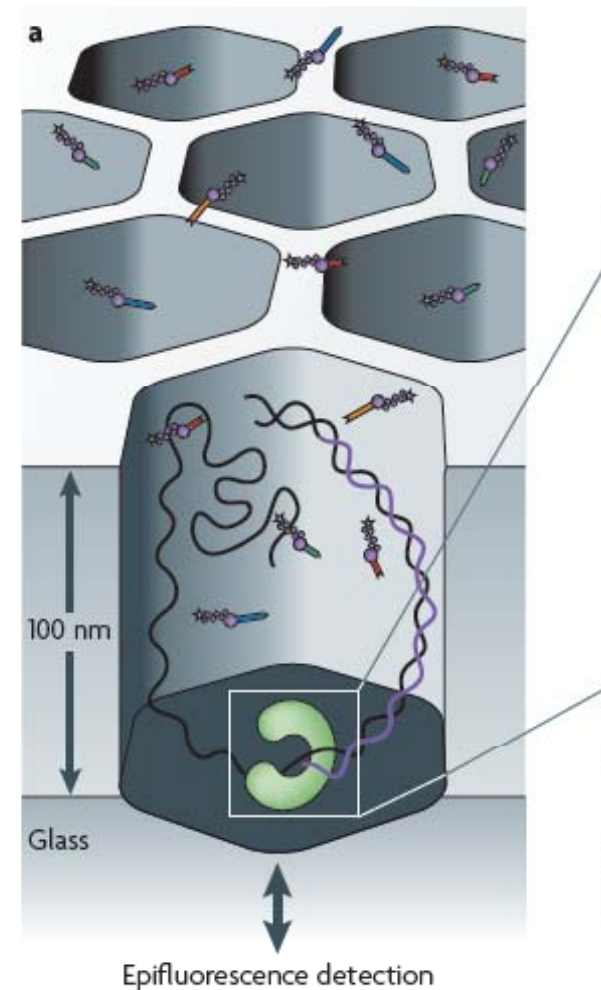




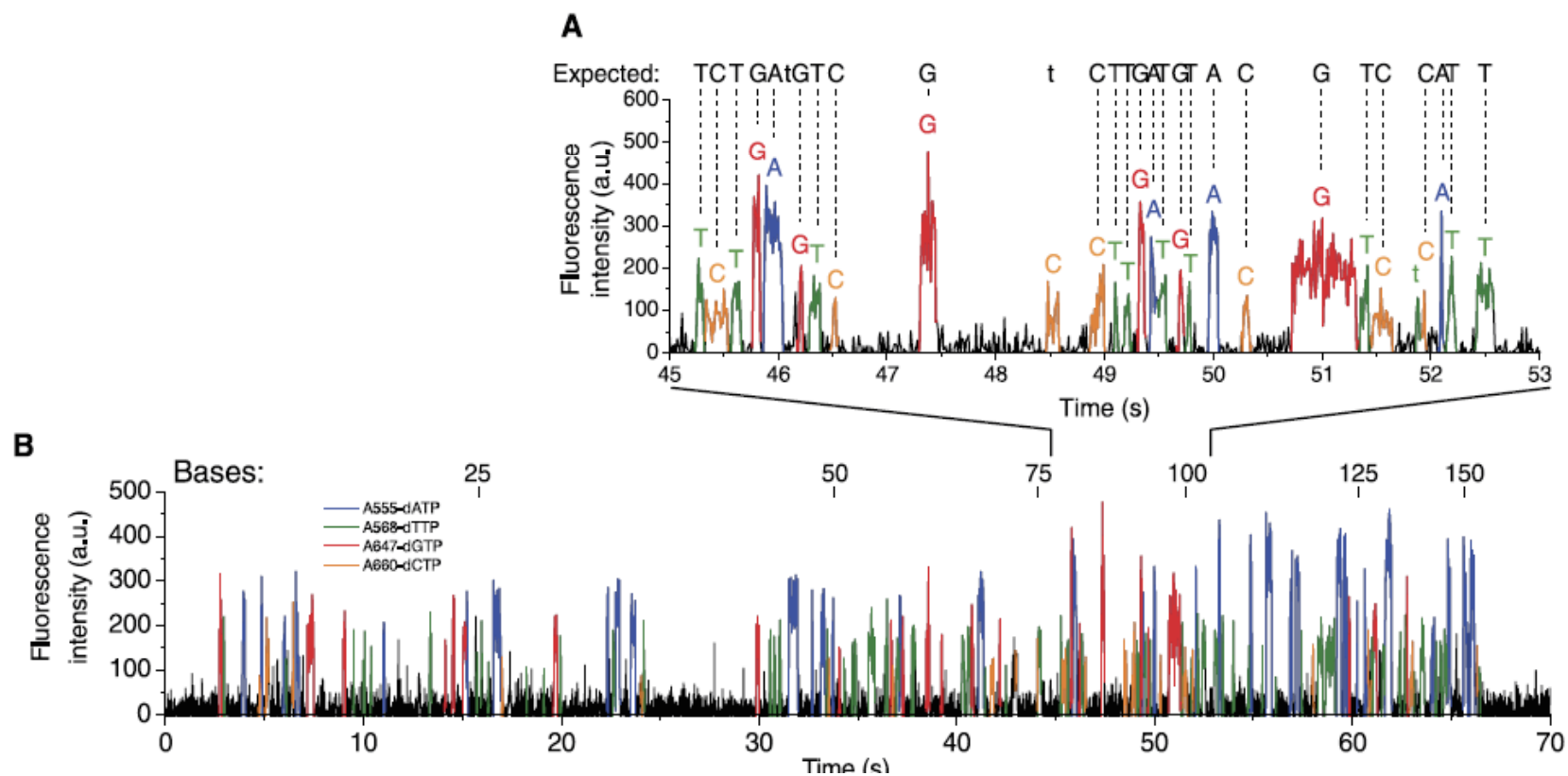
Figure 10. Processive Synthesis with Phospholinked Nucleotides.

Step 1: Fluorescent phospholinked labeled nucleotides are introduced into the ZMW.

Step 2: The base being incorporated is held in the detection volume for tens of milliseconds, producing a bright flash of light.

Step 3: The phosphate chain is cleaved, releasing the attached dye molecule.

Step 4-5: The process repeats.



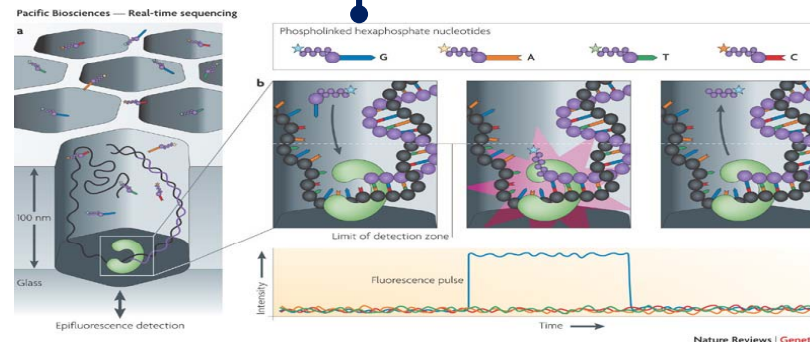
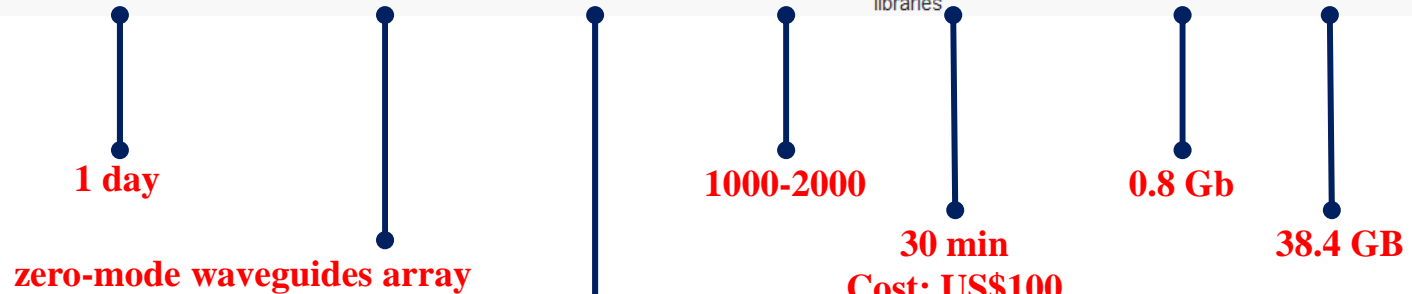
What makes the differences?

easier library preparation, shorter run time, longer read length, higher throughput, lower cost

Summary of second-generation sequencers

| Sequencing platform | Sample requirements | Length of library prep/feature generation (days) | Method of feature generation | Sequencing chemistry | Read length (bases) | Run time | Throughput/run (Gb) | Throughput/day (Gb) |
|---------------------------------|---|--|---|---------------------------|---------------------|--|---------------------|---------------------|
| Roche 454 (FLX-Titanium) | 1 µg for shotgun library, 5 µg for paired end | 3-4 | Bead-based/emulsion PCR | Pyrosequencing | 400-500 | 10 h | 0.4-0.5 | ~1 |
| Illumina Genome Analyzer (GAII) | <1 µg for single or paired-end libraries | 2 | Isothermal 'bridge amplification' on flowcell surface | Reversible terminator SBS | 35-75 | 2 days for 36-cycle single-end run, 4 days for 36-cycle paired-end run | 3-6 | 1.5 |
| ABI SOLiD | <2 µg for shotgun library, 5-20 µg for paired end | 2-4.5 | Bead-based/emulsion PCR | Ligation | 25-75 | 6-7 days for fragment libraries, 8 days for 2 × 25 base paired-end libraries | 10-20 | 1.7-2 |

Pacific Biosciences



Single-molecule real-time DNA sequencing (SMRT)

SMRT™ sequencing sample preparation workflow

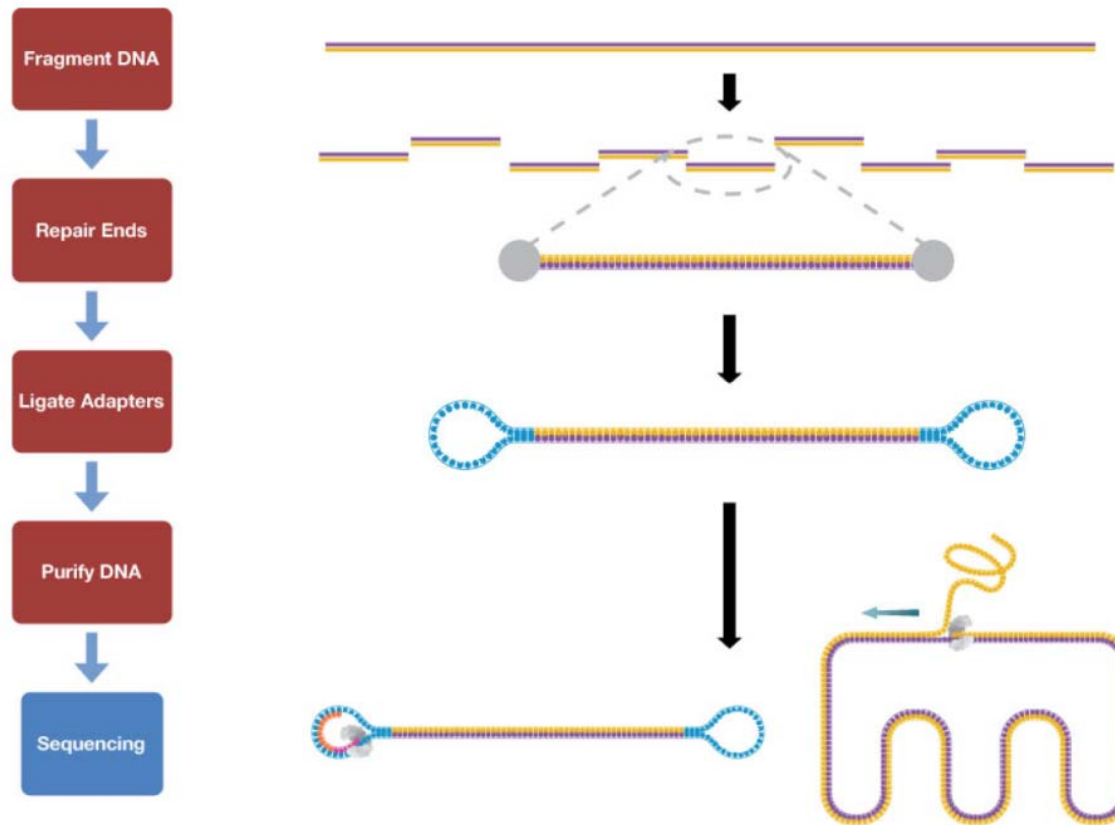
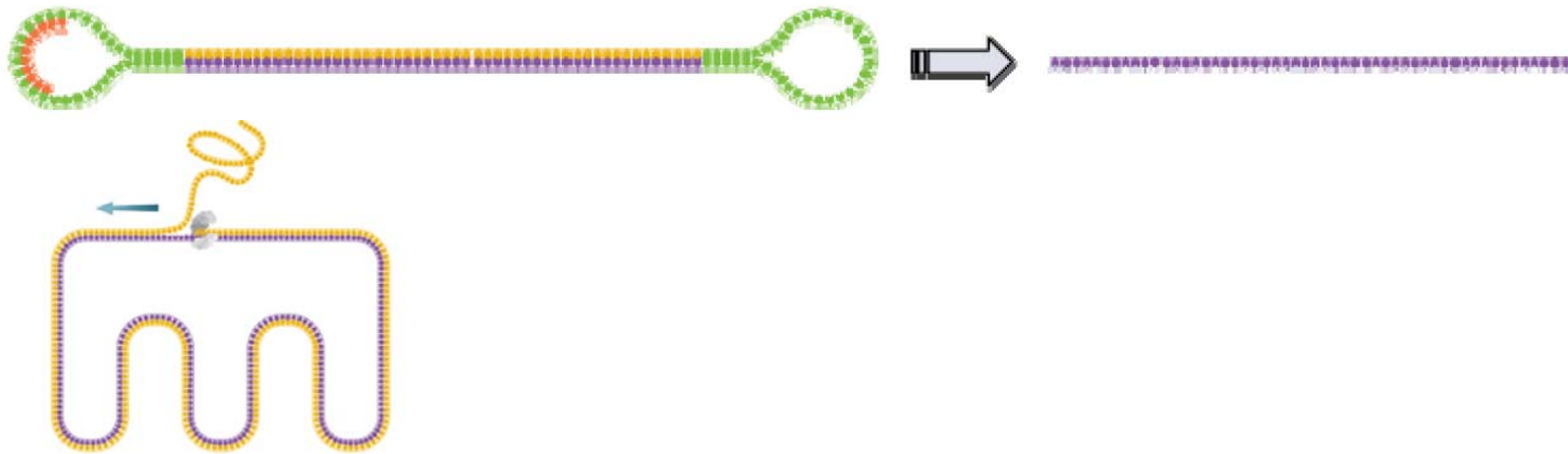


Figure 17. Sample Prep Workflow.

The input sample is first fragmented to the desired size. The ends are repaired and the hairpin structures are ligated to the ends of each fragment. A size selection and purification step selects those fragments with the adaptors attached to both ends. The SMRTbell templates then can go through the sequencing reaction. A strand displacing polymerase enzyme opens the SMRTbell into a circular template and can generate independent reads, both forward and reverse of the same DNA molecule. The quality score increases linearly with the number of times the molecule is sequenced.

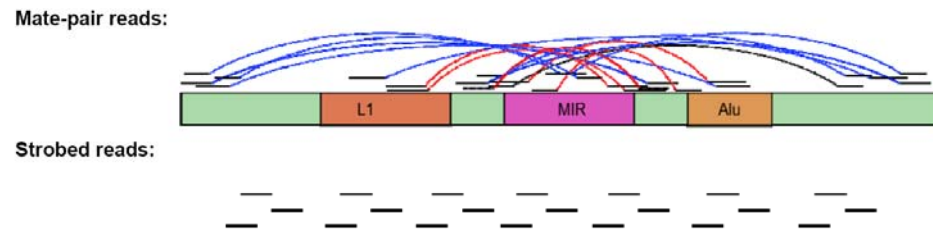
Standard sequencing protocol.

Enzyme processivity enables long readlengths while the speed of synthesis drives fast time to results.



Strobe sequencing protocol.

Strobe sequencing offers greater flexibility and eliminates the need to create multiple libraries of different sizes.



ABI 3730XL

- up to 1100 bases/read
- 96 reads/run
- approx. 1MB/day and machine

First choice for finishing projects; full length cDNA sequencing; single sample sequencing

LONG
READS



Roche GS FLX

- in average 250 bases/read
- up to 400 000 reads/run
- up to 100MB/run/7.5 hours

Optimal for *de novo* sequencing (procaryots & eucaryots) metagenomes, transcriptomes/cDNA, BAC/fosmid pools

EXTREME
SPEED



Illumina Genome Analyzer

- up to 50 bases/read
- up to 60 000 000 reads/run (paired-end)
- up to 2000MB/run/6.5 days
- Sequencing by synthesis

ABI SOLiD

- up to 35 bases/read
- up to 85 000 000 reads/run (paired-end)
- up to 3000MB/run/6 days
- Sequencing by ligation

Highly attractive for resequencing projects of e.g. production strains; small RNAs, SAGE/CAGE and ChIP-Seq; ultra deep mutation/SNPs

HUGE DATA
QUANTITY



VIDEOS

454

SEQUENCING

<http://www.454.com/products-solutions/multimedia-presentations.asp>

Genome Sequencer FLX Multimedia Presentation

Genome Sequencer FLX Standard Series Workflow Presentation

Genome Sequencer FLX Amplicon Sequencing Presentation



http://www.illumina.com/technology/sequencing_technology.ilmn



http://marketing.appliedbiosystems.com/images/Product/Solid_Knowledge/flash/102207/solid.html



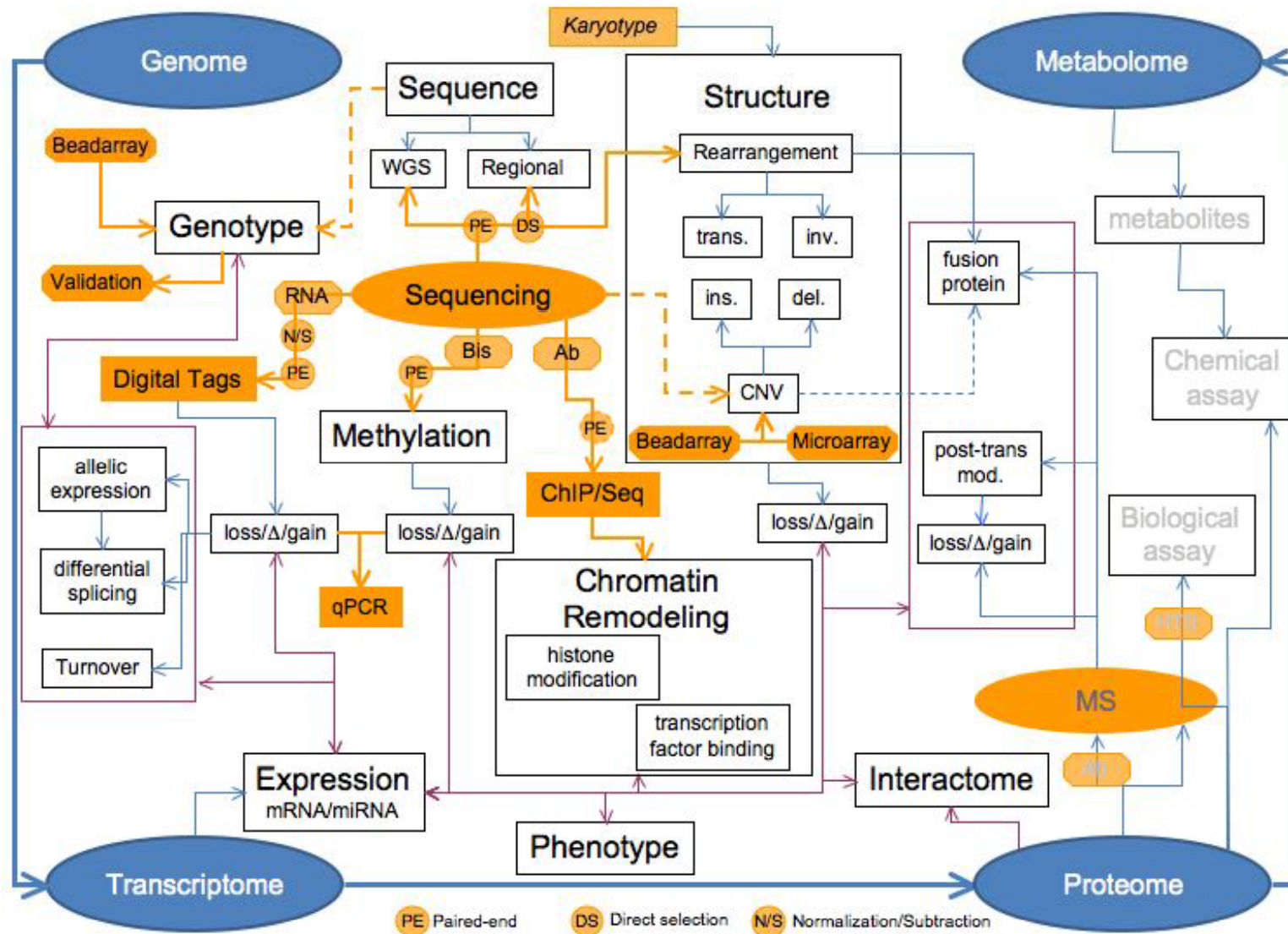
<http://www.helicosbio.com/Technology/TrueSingleMoleculeSequencing/tSMStradeHowItWorks/tabid/162/Default.aspx>



http://visigenbio.com/technology_movie_streaming.html



http://www.pacificbiosciences.com/video_lg.html



ARCHON GENOMICS X PRIZE

Archon X PRIZE for Genomics Teams Media Center Take Action Discover About

"The scaffold has been handed down to us from our ancestors, and through it we are connected to all other life on earth."

- Svante Pääbo

Revolution Through Competition.

▶ TAKE ACTION

ARCHON X PRIZE FOR GENOMICS

- ▶ Introduction
- ▶ [Why Genomics?](#)
- ▶ [PRIZE Overview](#)
- ▶ [Why Whole Genome Sequencing?](#)
- ▶ [The Promise of Personalized Medicine](#)
- ▶ [Frequently Asked Questions](#)
- ▶ [Competition Guidelines \[PDF\]](#)
- ▶ [Register to Compete \[PDF\]](#)

ARCHON X PRIZE FOR GENOMICS

Scientists know that a map of our genome holds boundless potential, ranging from identifying our susceptibility to disease to discovering cures for cancer. But since 1953, when James Watson and Francis Crick concluded that DNA contained the "stuff of life," only a handful of human genomes have been mapped. In fact, it still takes many months and millions of dollars to sequence a single genome.

Understanding our genomes may help delay or even prevent disease. For those suffering from genetic illnesses, personal genetic information can determine which medicines will drive their disease into remission without negative side-effects.

The Archon X PRIZE for Genomics challenges scientists and engineers to create better, cheaper and faster ways to sequence genomes. The knowledge gained by compiling and comparing a library of human genomes will create a new era of preventive and personalized medicine — and transform medical care from reactive to proactive.

The X PRIZE Foundation and scientists the world over dream of the day when we fully understand the human genetic blueprint — enabling us to make informed decisions about our own health and create a brighter future for generations to come.

Please join the X PRIZE Foundation in our challenge to create a breakthrough in genomics that will benefit all of humanity.

[Join the Revolution](#)

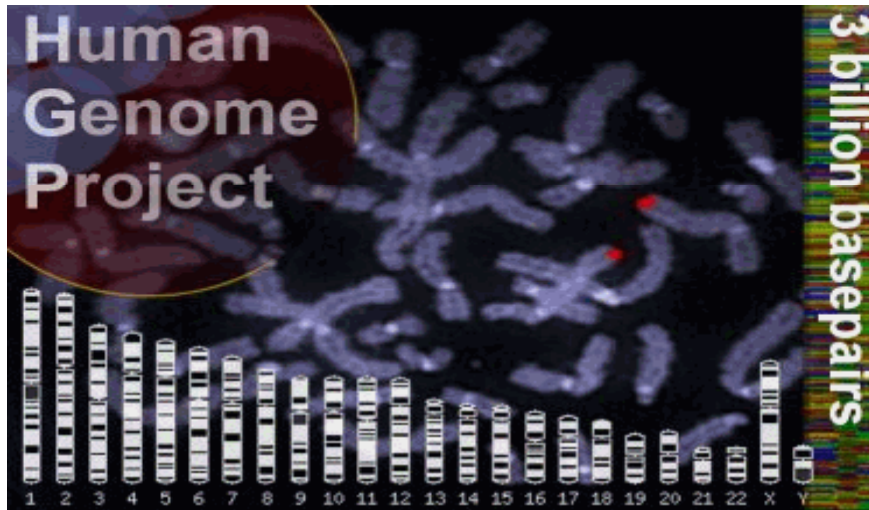
"Many thought that the unbelievable feat in 2001 of deciphering the human genome was the end of the story. Not so! We now firmly believe that it was only the beginning. The true challenge now appears to be in deciphering individual variations within the human genome. It is there that the true importance of genomics resides. Only by unearthing the blueprint details of each and every individual, will we be able to fully fathom health and disease, mental and bodily faculties, and the true secrets of human biology. The Archon X PRIZE for Genomics will greatly enhance scientists' ambition to develop methodologies for individual genomic "blue-printing". I believe the PRIZE will bring about bona-fide successes within 5 years, much earlier than otherwise feasible. It will thereby create the real genomic revolution."

*Prof. Doron Lancet
Head, Crown Human Genome Center
Weizmann Institute of Science
Archon X PRIZE For Genomics
Scientific Advisory Board*



A \$10 MILLION PRIZE
FOR THE FIRST TEAM TO SUCCESSFULLY SEQUENCE
100 HUMAN GENOMES IN 10 DAYS

How much does it cost for a Human Genome?



Human Reference Genome

April, 2003

>10 years to finish
USD 3 billion



James Watson's Personal Genome

June, 2007

1 year
USD 2 million



Craig Venter's Personal Genome

September, 2007

1 year
USD 1 million



YH Genome

November, 2008

1 year
~USD 0.5 million



Personal Genome

When will it available?
How long will it take?
How much?

Sequenced Human Genomes

| Personal Genome | Platform | Genomic template libraries | No. of reads (millions) | Read length (bases) | Base coverage (fold) | Assembly | Genome coverage (%)* | SNVs in millions (alignment tool) | No. of runs | Estimated cost (US\$) |
|--------------------------|---------------------|----------------------------------|-------------------------|---------------------|----------------------|----------------|----------------------|-----------------------------------|-------------|--------------------------|
| J. Craig Venter | Automated Sanger | MP from BACs, fosmids & plasmids | 31.9 | 800 | 7.5 | <i>De novo</i> | N/A | 3.21 | >340,000 | 70,000,000 |
| James D. Watson | Roche/454 | Frag: 500 bp | 93.2 [†] | 250 [§] | 7.4 | Aligned* | 95 ^l | 3.32 (BLAT) | 234 | 1,000,000 [†] |
| Yoruban male (NA18507) | Illumina/Solexa | 93% MP: 200 bp | 3,410 [‡] | 35 | 40.6 | Aligned* | 99.9 | 3.83 (MAQ) | 40 | 250,000 [†] |
| | | 7% MP: 1.8 kb | 271 | 35 | | | | 4.14 (ELAND) | | |
| Han Chinese male | Illumina/Solexa | 66% Frag: 150–250 bp | 1,921 [‡] | 35 | 36 | Aligned* | 99.9 | 3.07 (SOAP) | 35 | 500,000 [†] |
| | | 34% MP: 135 bp & 440 bp | 1,029 | 35 | | | | | | |
| Korean male (AK1) | Illumina/Solexa | 21% Frag: 130 bp & 440 bp | 393 [‡] | 36 | 27.8 | Aligned* | 99.8 | 3.45 (GSNAP) | 30 | 200,000 [†] |
| | | 79% MP: 130 bp, 390 bp & 2.7 kb | 1,156 | 36, 88, 106 | | | | | | |
| Korean male (SJK) | Illumina/Solexa | MP: 100 bp, 200 bp & 300 bp | 1,647 [‡] | 35, 74 | 29.0 | Aligned* | 99.9 | 3.44 (MAQ) | 15 | 250,000 ^{†,‡} |
| Yoruban male (NA18507) | Life/APG | 9% Frag: 100–500 bp | 211 [‡] | 50 | 17.9 | Aligned* | 98.6 | 3.87 (Corona-lite) | 9.5 | 60,000 ^{†,‡} |
| | | 91% MP: 600–3,500 bp | 2,075 [‡] | 25, 50 | | | | | | |
| Stephen R. Quake | Helicos BioSciences | Frag: 100–500 bp | 2,725 [‡] | 32 [§] | 28 | Aligned* | 90 | 2.81 (IndexDP) | 4 | 48,000 [†] |
| AML female | Illumina/Solexa | Frag: 150–200 bp ^{††} | 2,730 ^{††} | 32 | 32.7 | Aligned* | 91 | 3.81 ^{††} (MAQ) | 98 | 1,600,000 ^{†††} |
| | | Frag: 150–200 bp ^{§§} | 1,081 ^{†,§§} | 35 | 13.9 | | 83 | 2.92 ^{§§} (MAQ) | 34 | |
| AML male | Illumina/Solexa | MP: 200–250 bp ^{††} | 1,620 ^{††} | 35 | 23.3 | Aligned* | 98.5 | 3.46 ^{††} (MAQ) | 16.5 | 500,000 ^{†††} |
| | | MP: 200–250 bp ^{§§} | 1,351 ^{†,§§} | 50 | 21.3 | | 97.4 | 3.45 ^{§§} (MAQ) | 13.1 | |
| James R. Lupski CMT male | Life/APG | 16% Frag: 100–500 bp | 238 [‡] | 35 | 29.6 | Aligned* | 99.8 | 3.42 (Corona-lite) | 3 | 75,000 ^{†,††} |
| | | 84% MP: 600–3,500 bp | 1,211 [‡] | 25, 50 | | | | | | |

USD 10000 per Genome in 2015?



HOME

OUR SERVICES

YOUR GENOME

ABOUT US

CONTACT US



Know thyself.

The first personal genomics company to offer complete genome sequencing and analysis services for private individuals.



Our approach



KnomeCOMPLETE™



KnomeSELECT™



Frequently asked questions



Recent news

US\$ XX,000 for a personal genome

Complete Genomics, Inc.
is pleased to announce its
Initial Public Offering
NASDAQ: GNOM
November 11, 2010

[Learn More](#) ▶



Human Genome Sequencing & Analysis Service

Dedicated to **complete human genome sequencing and analysis** provided as an innovative, **end-to-end, outsourced service model**, Complete Genomics enables researchers to conduct **large-scale complete human genome studies**.

By **optimizing our sequencing platform for human DNA**, we are able to achieve **accuracy levels of 99.999%** at a total cost that is significantly less than the total cost of purchasing and using commercially available DNA sequencing instruments.

» [read more](#)

Receive Research-Ready Genomic Data

We offer our customers an end-to-end, outsourced solution that delivers **research-ready genomic data**.

Our **CGA™ Service** provides reliable access to assembled and **annotated sequence data** and our analytical tools enable our customers to **rapidly analyze and compare genomic data**.

» [read more](#)

WEBINAR

Pedigree Genome Sequencing



presented by:
Jared Roach, Ph.D., M.D.,
Institute for Systems Biology

[View](#)

News and Events

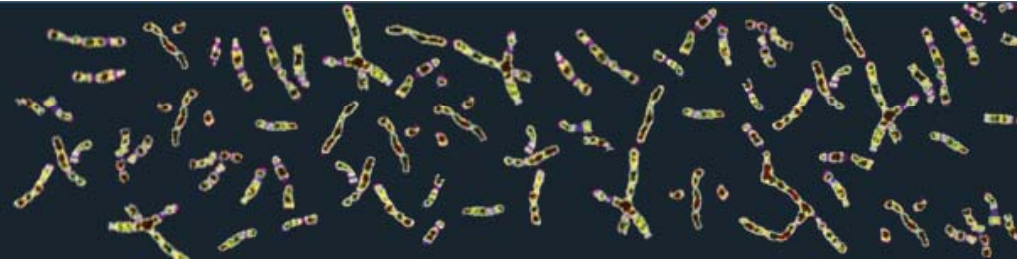
» **Complete Genomics Announces Pricing of Initial Public Offering**

» **Complete Genomics Expands CGA Tools Software Suite**

» **Complete Genomics to Sequence 100 Genomes for NCI Pediatric**

1000 Genomes

A Deep Catalog of Human Genetic Variation



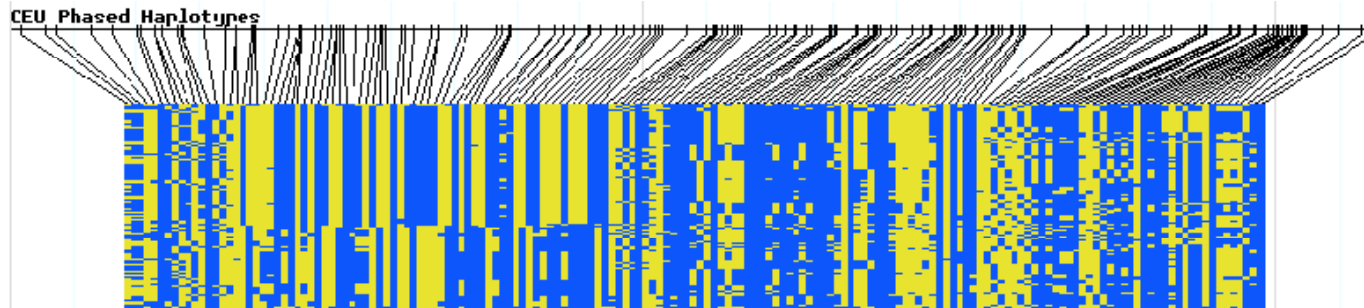
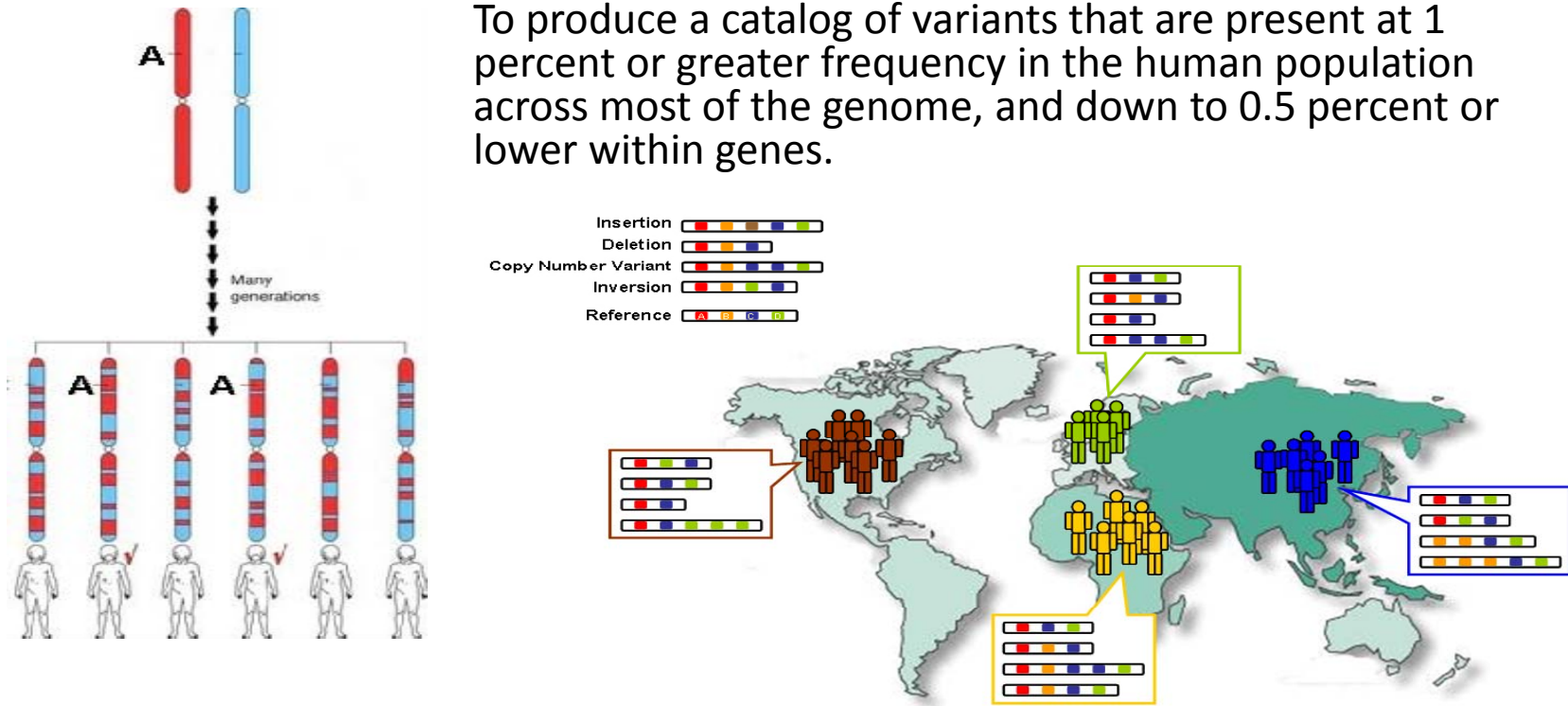
The 1000 Genomes Project is an international research consortium formed to create the most detailed and medically useful picture to date of human genetic variation. The project involves sequencing the genomes of approximately 1200 people from around the world and receives major support from the [Wellcome Trust Sanger Institute](#) in Hinxton, England, the [Beijing Genomics Institute Shenzhen](#) in China and the [National Human Genome Research Institute](#) (NHGRI), part of the [National Institutes of Health](#) (NIH).

Drawing on the expertise of multidisciplinary research teams, the 1000 Genomes Project will develop a new map of the human genome that will provide a view of biomedically relevant DNA variations at a resolution unmatched by current resources. As with other major human genome reference projects, data from the 1000 Genomes Project will be made swiftly available to the worldwide scientific community through freely accessible public databases.

On 4 September 2008, the co-chairs of the analysis group and overall project co-chairs drafted a letter to the NIH Council about 1000 Genomes Project. This letter, available [here](#), reviews the goals, describes the current status, and provide an update on the critical tasks the Analysis Group must accomplish in order to deliver a valuable community resource and achieve the Project's goals.

The Scientific Goals of the 1000 Genomes Project

To produce a catalog of variants that are present at 1 percent or greater frequency in the human population across most of the genome, and down to 0.5 percent or lower within genes.



A map of human genome variation from population-scale sequencing

The 1000 Genomes Project Consortium

Affiliations | Contributions | Corresponding author

Nature 467, 1061–1073 (28 October 2010) | doi:10.1038/nature09534

Received 20 July 2010 | Accepted 30 September 2010 | Published online 27 October 2010

Abstract

Abstract • Introduction • Data generation, alignment and variant discovery • Power to detect variants • Genotype accuracy • Putative functional variants • Application to association studies • Mutation, recombination and natural selection • Discussion • Methods • References • Acknowledgements • Author information • Supplementary information • Comments

The 1000 Genomes Project aims to provide a deep characterization of human genome sequence variation as a foundation for investigating the relationship between genotype and phenotype. Here we present results of the pilot phase of the project, designed to develop and compare different strategies for genome-wide sequencing with high-throughput platforms. We undertook three projects: low-coverage whole-genome sequencing of 179 individuals from four populations; high-coverage sequencing of two mother–father–child trios; and exon-targeted sequencing of 697 individuals from seven populations. We describe the location, allele frequency and local haplotype structure of approximately 15 million single nucleotide polymorphisms, 1 million short insertions and deletions, and 20,000 structural variants, most of which were previously undescribed. We show that, because we have catalogued the vast majority of common variation, over 95% of the currently accessible variants found in any individual are present in this data set. On average, each person is found to carry approximately 250 to 300 loss-of-function variants in annotated genes and 50 to 100 variants previously implicated in inherited disorders. We demonstrate how these results can be used to inform association and functional studies. From the two trios, we directly estimate the rate of *de novo* germline base substitution mutations to be approximately 10^{-8} per base pair per generation. We explore the data with regard to signatures of natural selection, and identify a marked reduction of genetic variation in the neighbourhood of genes, due to selection at linked sites. These methods and public data will support the next phase of human genetic research.

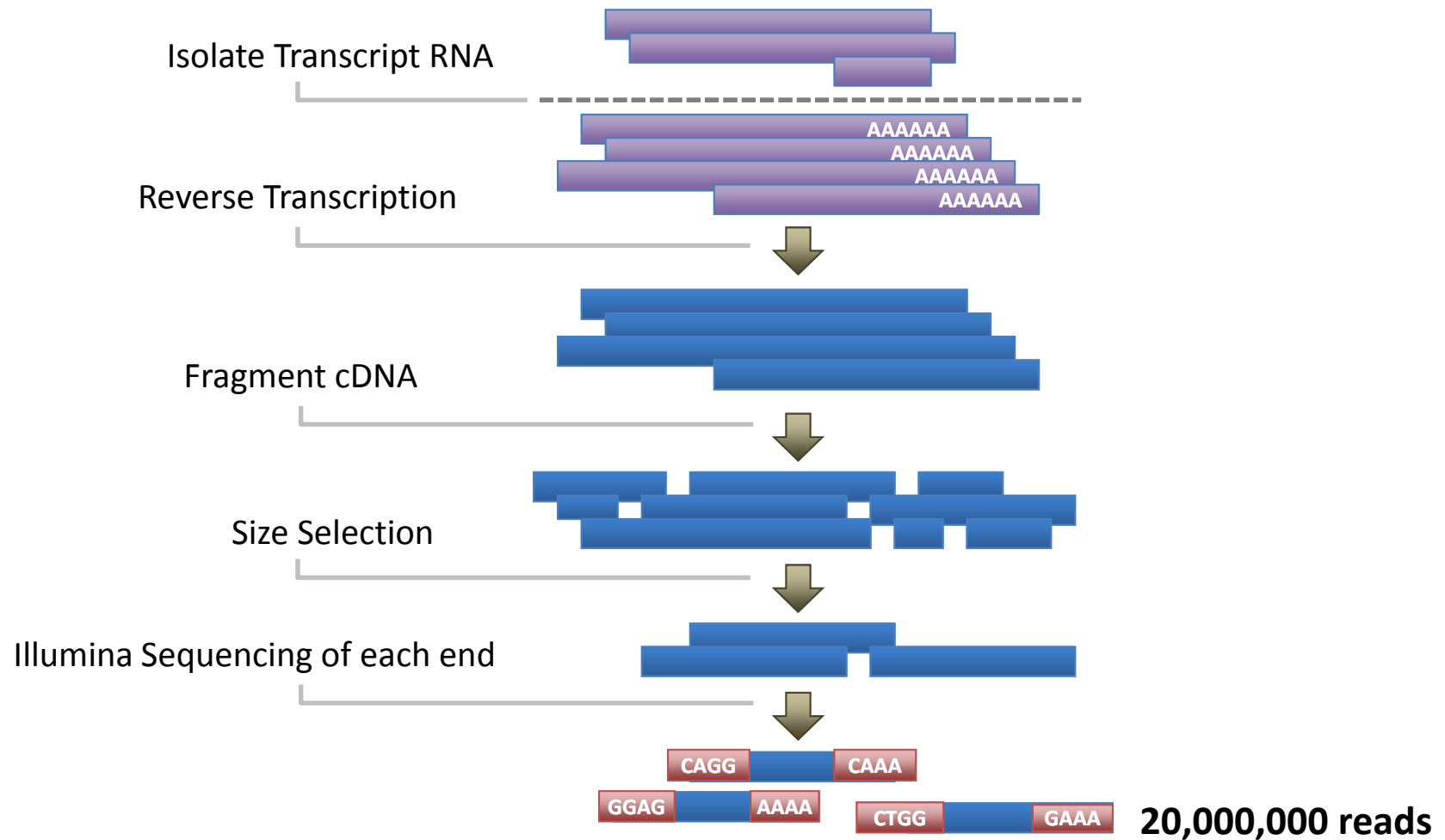
- 日本語要約
- print
- email
- download pdf
- download citation
- order reprints
- rights and permissions
- share/bookmark

| 1000 Genomes Samples | | | | | | | IS- Mar- 10 |
|---|-------------------------------------|------------------------------------|-------------------------------------|-------------|------------|------------|-------------------|
| Population | Status | Cell lines avail (all data approx) | Adult child samples from trio avail | First set | Second set | Third set | Total |
| Utah residents (CEPH) with Northern and Western European ancestry (CEU) | Available - HapMap samples | At Coriell | yes | 100 | | | 100 |
| Toscani in Italia (TSI) | Available - HapMap samples | At Coriell | | 100 | | | 100 |
| British from England and Scotland (GBR) | Available | At Coriell | | 100 | | | 100 |
| Finnish from Finland (FIN) | Available | At Coriell | | 100 | | | 100 |
| Iberian populations in Spain (IBS) | Collecting samples | Jun 2010 | yes | | 100 | | 100 |
| TOTAL European ancestry | | | | 400 | 100 | | 500 |
| Han Chinese in Beijing, China (CHB) | Available - HapMap samples | At Coriell | | 100 | | | 100 |
| Japanese in Tokyo, Japan (JPT) | Available - HapMap samples | At Coriell | | 100 | | | 100 |
| Han Chinese South (CHS) | Available | At Coriell | yes | 100 | | | 100 |
| Chinese Dai in Xishuangbanna (CDX) | Awaiting govt approval | late 2010 | yes | | 100 | | 100 |
| Khmer in Ho Chi Minh City, Vietnam (KHV) | Awaiting govt approval | late 2010 | yes | | 100 | | 100 |
| Chinese in Denver, Colorado (CHD) (sist Derby) | Available - HapMap samples | At Coriell | | 0 | 0 | | 0 |
| TOTAL East Asian ancestry | | | | 300 | 200 | | 500 |
| Yoruba in Ibadan, Nigeria (YRI) | Available - HapMap samples | At Coriell | yes | 100 | | | 100 |
| Luhya in Webuye, Kenya (LWK) | Available - HapMap samples | At Coriell | | 100 | | | 100 |
| Gambian in Western Division, The Gambia (GWD) | Collecting samples | late 2010 | yes | | 100 | | 100 |
| Ghanain in Navrongo, Ghana (GHN) | Final IRS approval expected soon | late 2010 | yes | | 100 | | 100 |
| Malasian in Selayar, Malawi (MAS) | Discussing issues for participation | T | yes | | 100 | | 100 |
| TOTAL West African ancestry | | | | 200 | 300 | | 500 |
| African Ancestry in Southwest US (ASW) | Available - HapMap samples | At Coriell | yes | 81 | | | 81 |
| African American in Jackson, MS (AJM) | IRS approval received, recruitment | late 2010 | yes | | 80 | | 80 |
| African Caribbean in Barbados (ACB) | Collecting samples | late 2010 | yes | | 70 | | 70 |
| Mexican Ancestry in Los Angeles, CA (MXL) | Available - HapMap samples | At Coriell | yes | 70 | | | 70 |
| Puerto Rican in Puerto Rico (PUR) | Available | Jun 2010 | yes | 70 | | | 70 |
| Colombian in Medellin, Colombia (CLM) | Collecting samples | Jun 2010 | yes | | 70 | | 70 |
| Peruvian in Lima, Peru (PEL) | Collecting samples | late 2010 | yes | | 70 | | 70 |
| TOTAL Americas | | | | 281 | 290 | | 500 |
| Ahem in the State of Assam, India | Awaiting govt & IRS approval | 2011 | yes | | | 100 | 100 |
| Kayastha in Calcutta, India | Awaiting govt & IRS approval | 2011 | yes | | | 100 | 100 |
| Raddi in Hyderabad, India | Awaiting govt & IRS approval | 2011 | yes | | | 100 | 100 |
| Marathi in Bombay, India | Awaiting govt & IRS approval | 2011 | yes | | | 100 | 100 |
| Punjab in Lahore, Pakistan | Awaiting IRS approval | late 2010 | yes | | | 100 | 100 |
| TOTAL South Asian ancestry | | | | | | 500 | 500 |
| TOTAL | | | | 1101 | 890 | 500 | 2500 |

Applications on Biomedical Sciences

| Category | Examples of applications |
|---|---|
| Complete genome resequencing | Comprehensive polymorphism and mutation discovery in individual human genomes |
| Reduced representation sequencing | Large-scale polymorphism discovery |
| Targeted genomic resequencing | Targeted polymorphism and mutation discovery |
| Paired end sequencing | Discovery of inherited and acquired structural variation |
| Metagenomic sequencing | Discovery of infectious and commensal flora |
| Transcriptome sequencing | Quantification of gene expression and alternative splicing; transcript annotation; discovery of transcribed SNPs or somatic mutations |
| Small RNA sequencing | microRNA profiling |
| Sequencing of bisulfite-treated DNA | Determining patterns of cytosine methylation in genomic DNA |
| Chromatin immunoprecipitation–sequencing (ChIP-Seq) | Genome-wide mapping of protein-DNA interactions |
| Nuclease fragmentation and sequencing | Nucleosome positioning |
| Molecular barcoding | Multiplex sequencing of samples from multiple individuals |

Transcriptome Sequencing

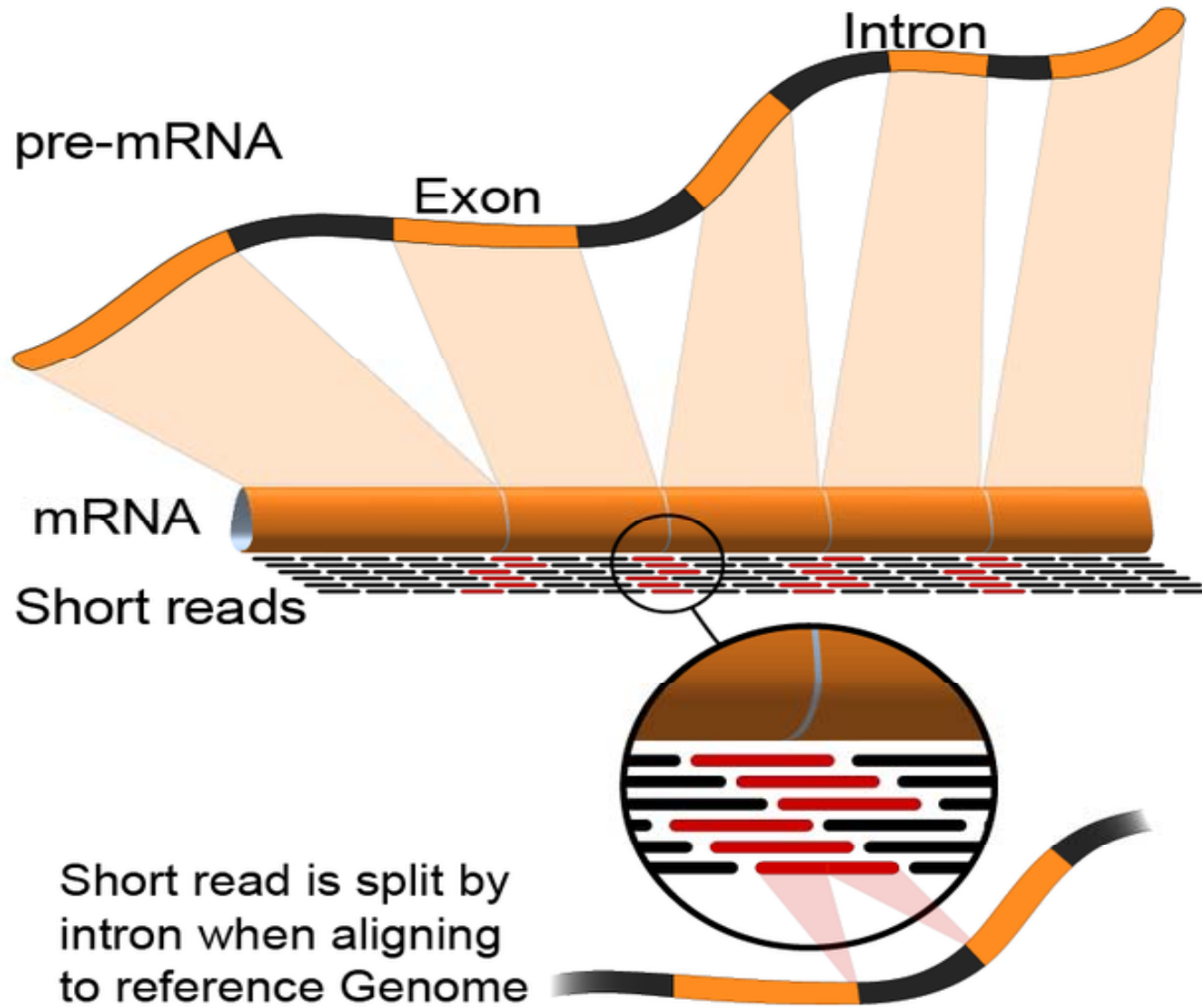


*based on Illumina approach

**strand-specific RNA-seq protocols exist for both Illumina and SOLiD

Slide complements of Andrew McPherson

Transcriptome Sequencing



Analysis Strategies: Reference Sequence Alignment vs De novo Assembly

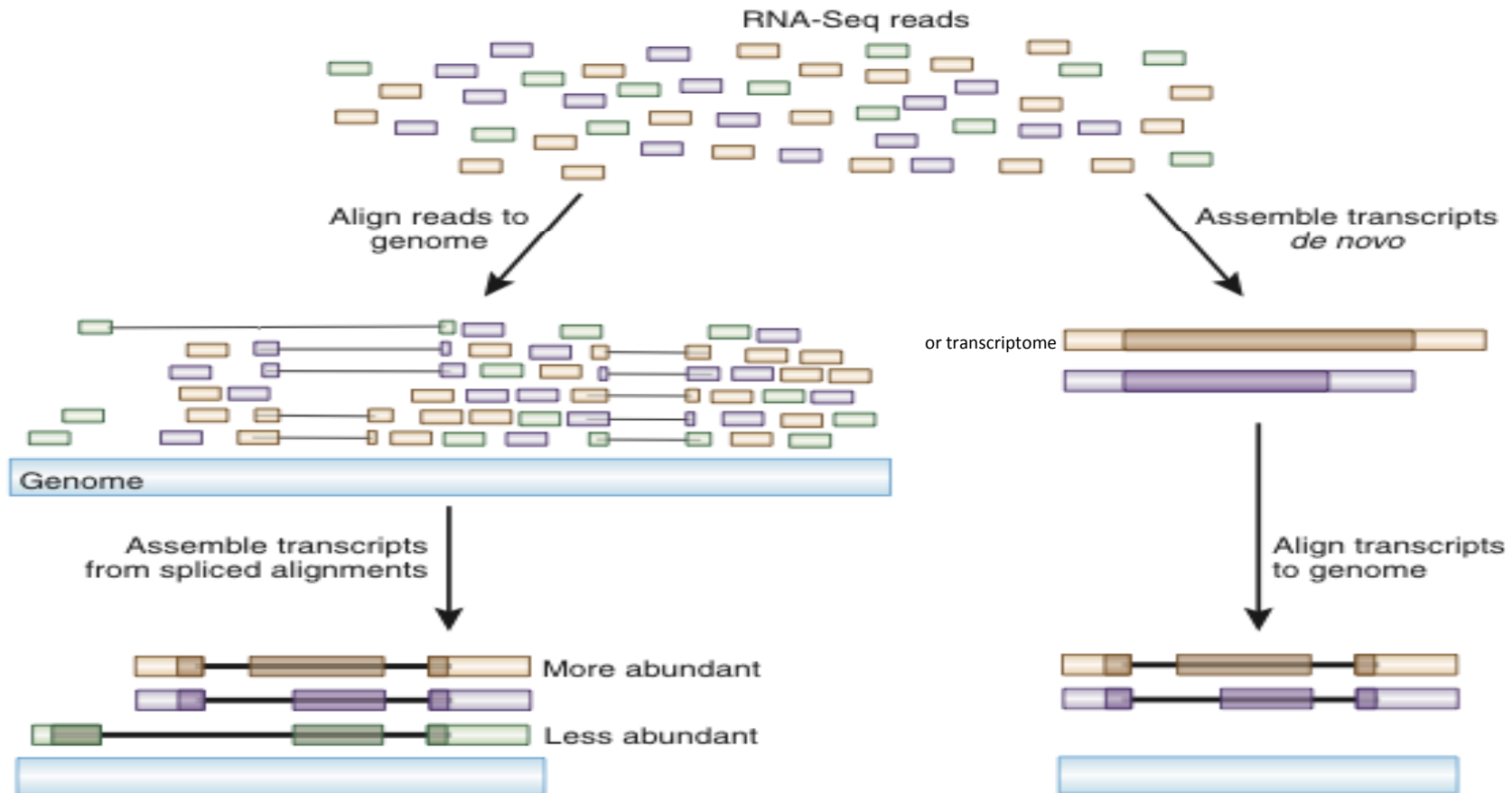


Image from Haas & Zody, 2010

*Assembly is the only option when working with a creature with no genome sequence, alignment of contigs may be to ESTs, cDNAs etc

Drawbacks for each strategy

- **Alignment to genome**

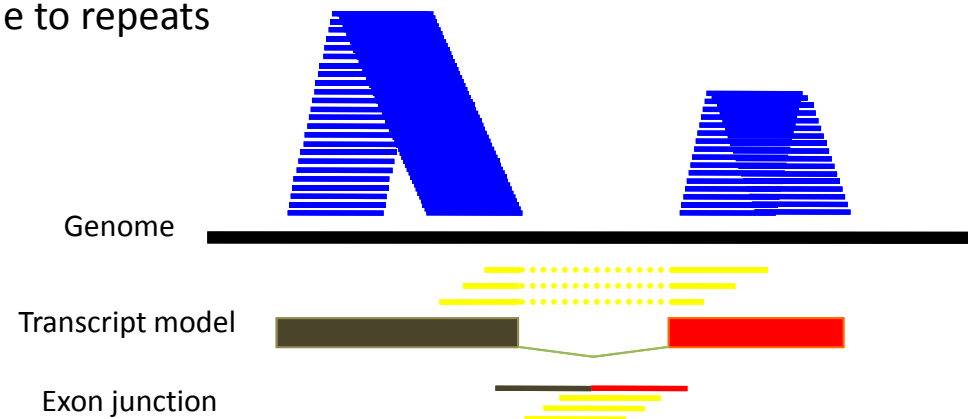
- Computationally expensive
- It is **never** a good idea to simply align RNA-seq data to the genome Need a spliced aligner or a surrogate (such as including exon-exon junction sequences in ‘genome’)

- **Alignment to transcriptome**

- Reads deriving from non-genic structures may be ‘forcibly’ (and erroneously) aligned to genes
 - Incorrect gene expression values
 - False positive SNVs
 - Many other potential problems

- **Assembly**

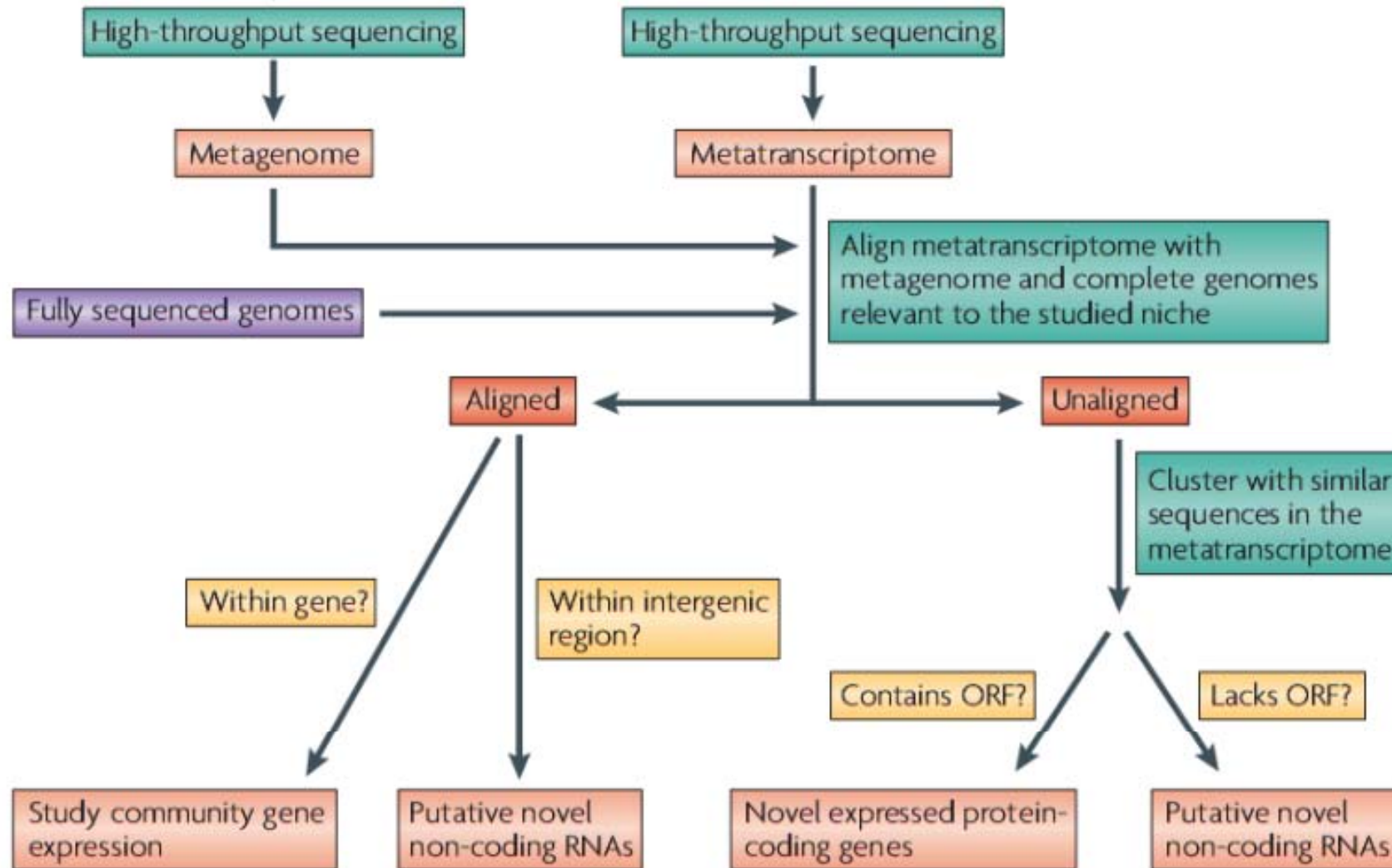
- Low expression = difficult/impossible to assemble
- Misassemblies/fragmented contigs due to repeats
- Requires vast amounts of memory



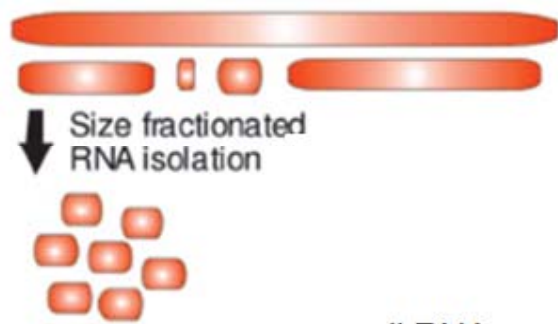
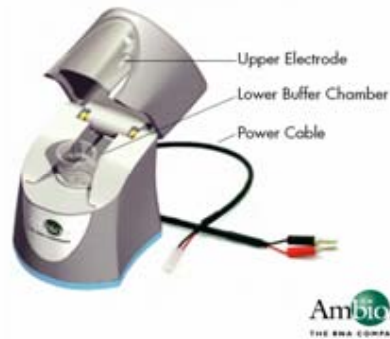
Benefits of each approach

- **Alignment to genome**
 - Allow reads from unannotated loci, introns *et cetera* to align to their correct locations... potential for new biological insights
- **Alignment to transcriptome**
 - Computationally inexpensive
 - Spliced (exon junction) reads map correctly
 - Pairing distance and junction reads may help distinguish individual isoforms (informative/unique regions of transcripts)
- **Assembly**
 - Can provide a more long-range view of transcripts
 - Allows detection of chimeric transcripts and resolution of 'breakpoints'
 - May not necessarily need a genome

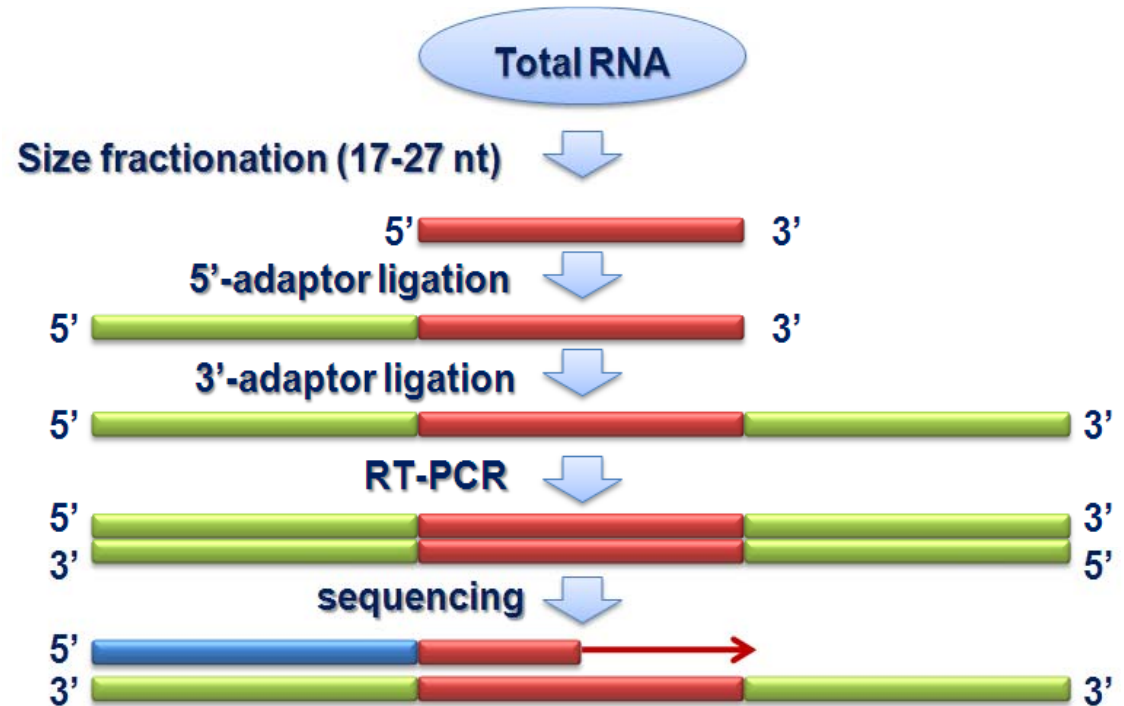
What can you get from Transcriptome Sequencing



Small RNA Sequencing



Illumina SOLEXA sequencing by synthesis



```
@HWI-EAS82_3_FC204V1AAXX:6:1:886:345
AGAGTTCTACAGTCCGGACGATCTCGTATGCCGTC
+HWI-EAS82_3_FC204V1AAXX:6:1:886:345
hhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhhh
```

40,000,000 reads of 35 bp long

Small RNA Sequencing

Workflow

Clean up



Clustering



ncRNA matching
(Rfam V10)



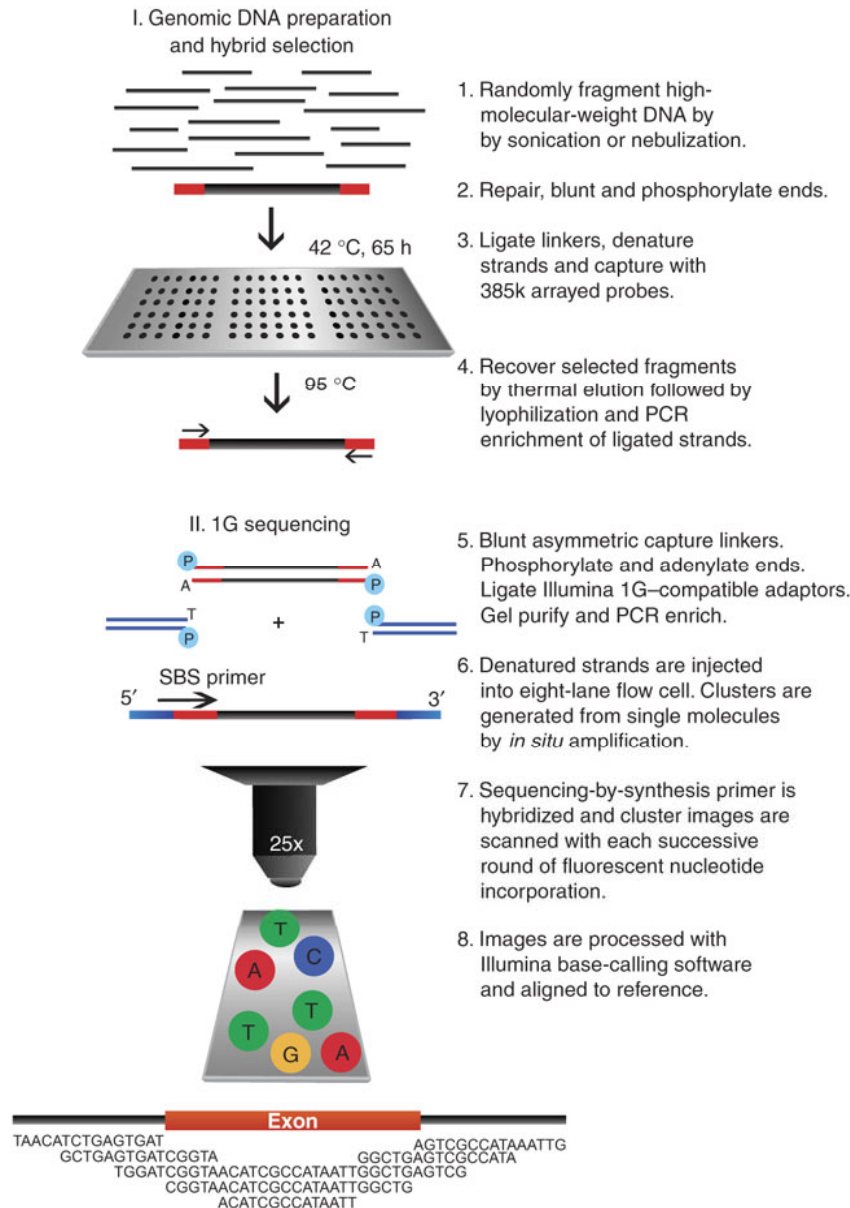
Known miRNA Matching
(miRBase V16)



Comparative miRNAomics



Human Exome Sequencing

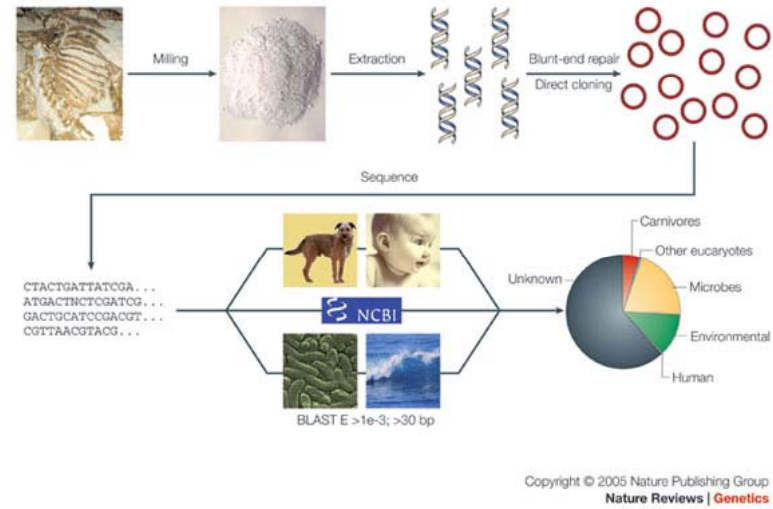
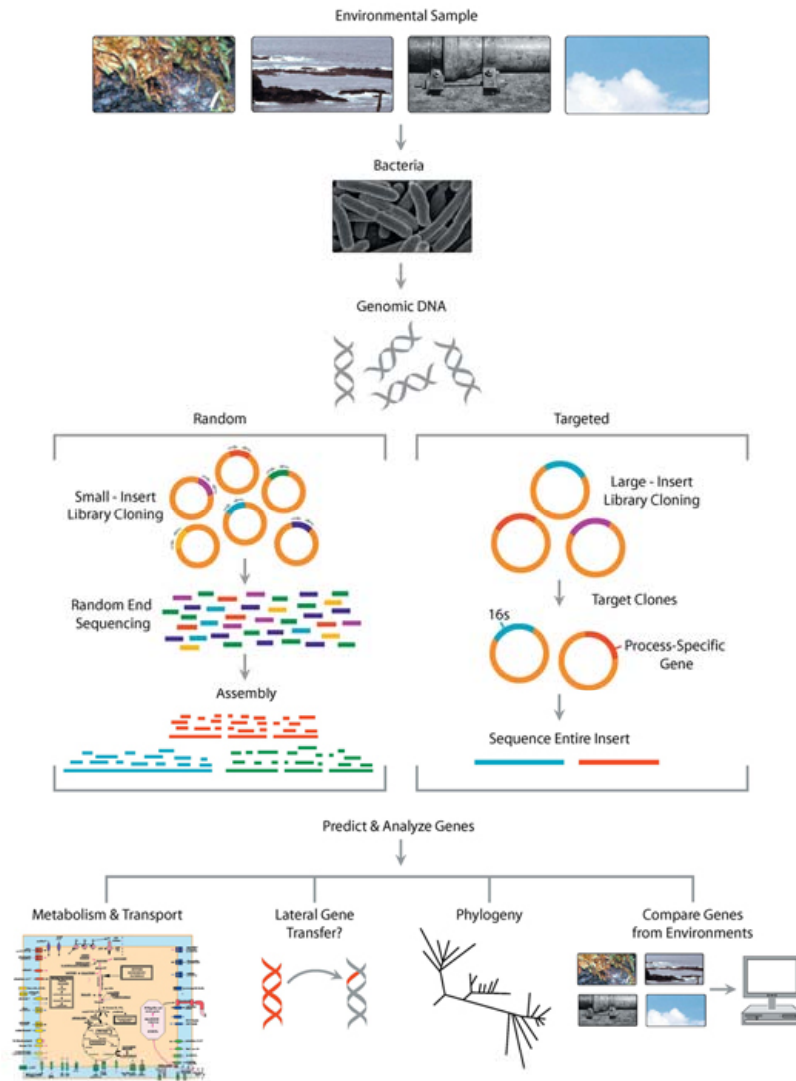


A single Nimblegen 2.1M array is capable of enrichment simultaneously more than 180,000 human coding-exon and 550 miRNA exons available in the CCDS database. The methods significantly improve the efficiency, the cost and throughput of target enrichment compared to conventional PCR-based methods.

The human exome captured library preparation consists of three major processes and lasts seven days, including library preparation for sequence capture, microarray hybridization, elution and library construction for sequencing, using Illumina Genome Analyzer can directly get high quality sequence data for the entire human exome.

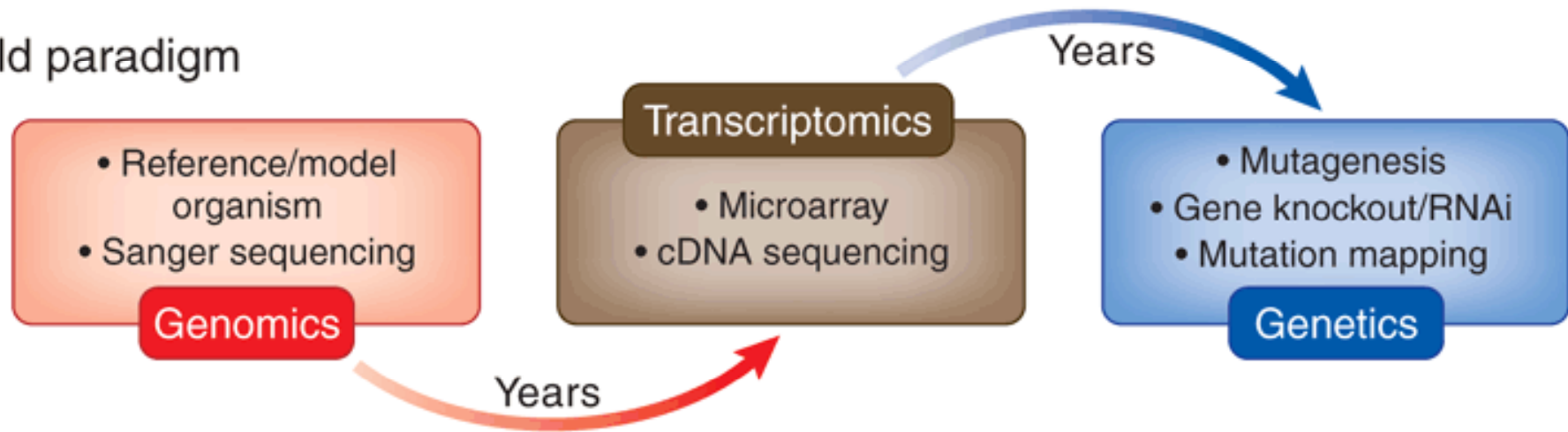
This approach can find the exact genes and mutations causing several complex human diseases, such as cancer, diabetes, obesity and so on.

Metagenomics

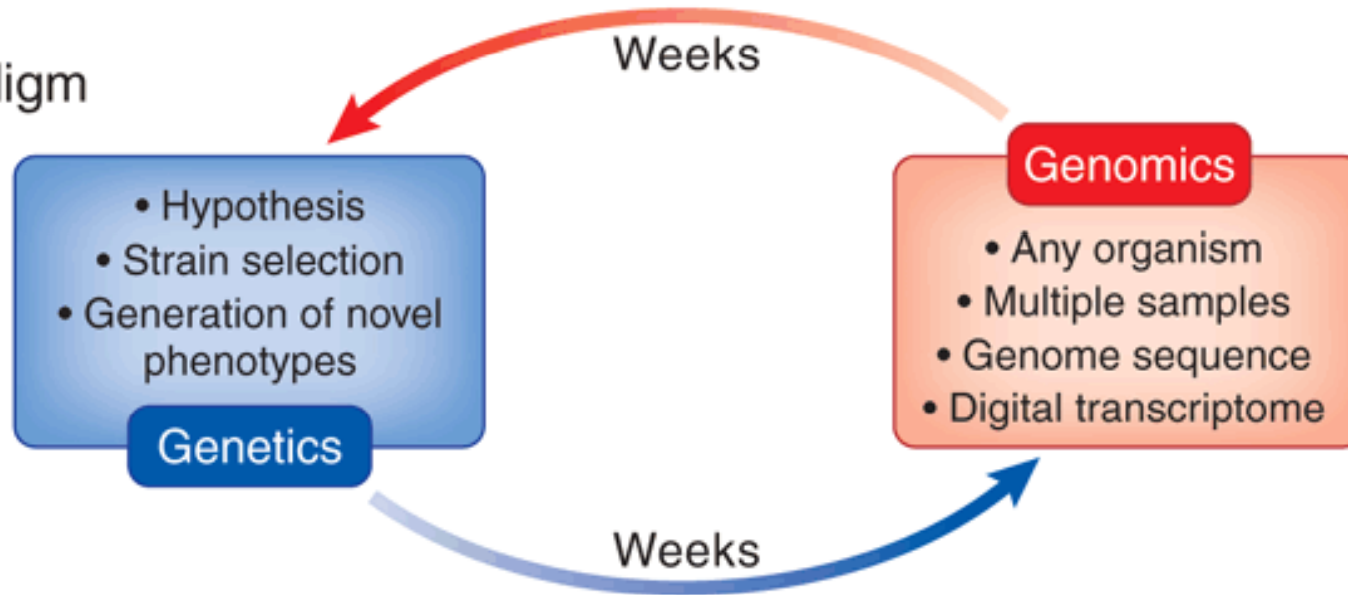


Copyright © 2005 Nature Publishing Group
Nature Reviews | Genetics

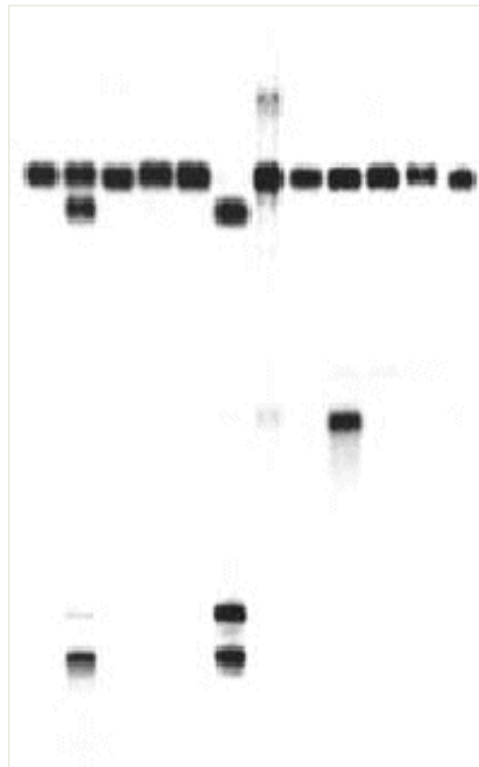
Old paradigm



New paradigm



The “old” biology



The most challenging task for a scientist is to get good data

The “new” biology



The most challenging task for a scientist is to make sense of lots of data

Computing Power

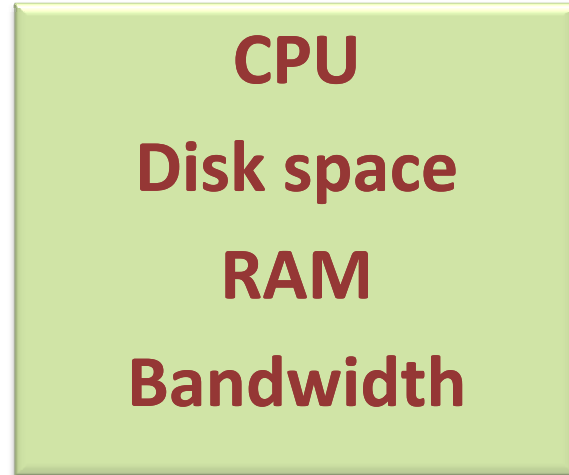
Next gen sequencers generated Giga bp to Tera bp of data



HiSeq2000 (launched in Q1 2010)

| Read Length | Run Time | Output |
|-------------|-----------|------------|
| 1 x 35 bp | ~1.5 days | 26-35 Gb |
| 2 x 50 bp | ~4 days | 75-100 Gb |
| 2 x 100 bp | ~8 days | 150-200 Gp |

Up to 25GB per day for a 2 x100 bp run



CPU

- The speed of your CPU determines how quickly it can process instructions
- Many bioinformatics operations fall into the “embarrassingly parallel” category
- Getting a results faster is as simple as adding more CPUs
- => clusters
- Though clusters still dominate, increasing attention is being paid to servers with large number of cores

For current throughputs, you will need ~8 < dual quad core > nodes per sequencer to handle raw data.

For data analysis, 1-2 nodes is enough for general purpose.

How about use Graphic Processing Unit (GPU) instead?

RAM

- Typical sizing is 2GB of RAM per core
- This works fine for most aligners
 - Certain aligners require a minimum number of cores for optimal efficiency
 - 16-48GB RAM per node
- Assemblers typically need much more RAM
- If you don't have enough RAM, the CPU will need to make use of the disk storage –
- When a computer has run out of RAM it is said to be “swapping”
- Some aligners have a minimum memory limit for optimal performance
- Human de novo requires lots of RAM!

Disk Space

- Sequencers generate a lot of data
- Including quality values and additional files (alignments etc) create multiplier
- Binary formats such as BAM are helping, but there are limited binary formats for basic data
 - Life and Illumina are both moving to binary formats for basic data
- **Debates on what needs to be stored**
- RAID
- Scaling

Bandwidth

Bandwidth of a connection represents the maximum rate of transfer between two points

- E.g. An aligner can process X reads per second on a single CPU at a data rate of Y bytes/sec
 - ~200 million reads in 10 hours
 - Each read 50 bases at 10 bytes per base
 - 2.7 MB/sec
- Design 100TB storage and connect it to a CPU resource
- Design bandwidth to be 10 Mb/sec – plenty of spare bandwidth
- Now we want to complete the job in an hour and get permission to buy 10 more CPUs – great!

BUT

For the 10 CPUs to run at maximum speed, they need to be supplied data at 27MB/sec

Our bandwidth is 10MB/sec

Therefore, no matter how many CPUs we buy, the job will never run faster than ~ 2.5 hours .

How about optic fibers?

<http://petang.cgu.edu.tw>



Molecular Regulation and Bioinformatics Laboratory

Home

MRBLab

Research

CV

Publications

Parasitology

Bioinformatics

The Molecular Regulation & Bioinformatics Laboratory (MRBLab) use bioinformatics approaches to integrate data generated by high-throughput technologies to compare the gene, protein and miRNA expression levels of protozoan as a basis for the development of new chemotherapeutic agents, to elucidate the interactome of pathogen-host and to study the biology of longevity in protozoan.

 Molecular Regulation & Bioinformatics Laboratory

Dept. Parasitology, College of Medicine, Chang Gung University, TAIWAN



PETRUS TANG
Ph.D. (CAM, UK)

Contact Information

Address:

Dept. of Parasitology, College of Medicine,
Chang Gung University,
259 Wenhwa 1st. Road. Kweishan,
Taoyuan 333. Taiwan.

Phone number:

+886-3-2118800#5136

Fax:

+886-3-2118122

E-mail:

petang@mail.cgu.edu.tw

LATEST NEWS

[Announcement]



[Lecture]

● [Bioinformatics Databases](#)

[Publication]

● May2010 - "[DSAP: Deep-Sequencing Small RNA Analysis Pipeline](#)" published in [Nucleic Acids Research](#).

● Mar2010 - "[The Genome of Trichomonas vaginalis](#)" in [Aerobic Parasitic Protozoa: Genomics and Molecular Biology](#)

● Mar2010 - "[Proteomic analysis of the effect of cyanide on Klebsiella oxytoca](#)" published in [Current Microbiology](#).

● Feb2010 - "[Trichomonas vaginalis vast BspA-like gene family: evidence for functional diversity from structural organisation and transcriptomics](#)" published in [BMC Genomics](#)

● Jan2010 - "[Loss of Bikunin in urine as useful marker of bladder carcinoma](#)" published in [J. Urology](#).



THANK YOU