10
 20
 30
 40
 50
 60
 7

 CGATIG A TTIAGC GGC CGCG AATTCGC CC T T TC TC TACG ACG ATG AT TTAC AC G C ATG TG C TG AAAG TTO

Harnessing Sequencing Technologies for Phylogenetics

Workshop in Applied Phylogenetics March 12, 2017

Joanna C. Chiu jcchiu@ucdavis.edu Department of Entomology and Nematology UC Davis

Outline

(1) Sequencing technologies : From Sanger sequencing to Oxford Nanopore Minion

(2) Types of NGS data for phylogenetics

Thanks to Brian Moore and Jonathan Eisen for ideas on this presentation!

GENOMICS

Sequencing all life captivates biologists

Project would read genomes of more than a million eukaryotes—just for starters

By Elizabeth Pennisi, in Washington, D.C.

Earth BioGenome Project (EBP)



A head start

The Earth BioGenome Project could coordinate the efforts below and others that are already sequencing broad swaths of the planet's life.

PROJECT	YEAR STARTED	SEQUENCING GOAL	NUMBER SEQUENCED
G10K	2009	9478 vertebrate genera	100
i5K	2011	5000 arthropods	30
GIGA	2013	7000 marine invertebrates	60
GAGA	2016	All 300 ant genera	25
B10K	2016	All 10,500 bird species	300
AOCC	2013	101 African food crops	22

Pennisi (Science, March 3, 2017)

Timeline for the development of sequencing technologies



From: Slideshare presentation of Cosentino Cristian and Jonathan Eisen

Timeline and Comparison of output for sequencing platforms



Numbers inside circles = read-length in bp

Reuter et al. (Mol Cell 2015)

Work by Maxam & Gilbert and Sanger started it all



The Nobel Prize in Chemistry 1980 Paul Berg, Walter Gilbert, Frederick Sanger

The Nobel Prize in Chemistry 1980





Paul Berg

Walter Gilbert

Frederick Sanger

The Nobel Prize in Chemistry 1980 was divided, one half awarded to Paul Berg *"for his fundamental studies of the biochemistry of nucleic acids, with particular regard to recombinant-DNA"*, the other half jointly to Walter Gilbert and Frederick Sanger *"for their contributions concerning the determination of base sequences in nucleic acids"*.

http://www.nobelprize.org/nobel_prizes/chemistry/laureates/1980/

Slide from Jonathan Eisen (UCD Bodega Bay Workshop in Applied Phylogenetics 2015)

Maxam-Gilbert Sequencing



A new method for sequencing DNA

(DNA chemistry/dimethyl sulfate cleavage/hydrazine/piperidine)

ALLAN M. MAXAM AND WALTER GILBERT

Department of Biochemistry and Molecular Biology, Harvard University, Cambridge, Massachusetts 02138 Contributed by Walter Gilbert, December 9, 1976



Maxam and Gilbert (PNAS 1977)

Capacity = For 40 cm gel, 100 bp

Mathews and van Holde (Biochemistry 4th ed.)

Sanger Sequencing

Proc. Natl. Acad. Sci. USA Vol. 74, No. 12, pp. 5463–5467, December 1977 Biochemistry

DNA sequencing with chain-terminating inhibitors

(DNA polymerase/nucleotide sequences/bacteriophage \$\phi X174\$)

F. SANGER, S. NICKLEN, AND A. R. COULSON Medical Research Council Laboratory of Molecular Biology, Cambridge CB2 2QH, England Contributed by F. Sanger, October 3, 1977

- Developed in 1977
- Chain termination using ddNTPs
- First sequenced genome: Bacteriophage \$\op\$X174 (5,368bp)



Sanger Sequencing





Sanger Sequencing – Automation I (e.g. ABI sequencers)



https://www.abmgood.com/marketing/knowledge_base/next_generation_sequencing_introduction.php

Sanger Sequencing Trace/ Chromatograms



Up to 800 – 1000 bases per read

High Throughput Sanger Shotgun Sequencing

Not really high throughput if it's manual



Nature Biotechnology 26, 1135 - 1145 (2008)

Sanger Sequencing – Automation II



Slide from Jonathan Eisen (UCD Bodega Bay Workshop in Applied Phylogenetics 2015)

Automated Sanger Sequencing Highlights

- 1991: ESTs by C. Venter
- 1995: H. influenzae shotgun genome
- 1996: Yeast, archaeal genomes
- 1998: 1st animal genome *C. elegans*
- 1999: Drosophila melanogaster shotgun genome
- 2000: Arabidopsis thaliana genome
- 2000: Human genome (Lander et al. 2001) (\$2.7 billion in FY 1991 – genome.gov)

Next-generation sequencing (NGS)



Sequence by Synthesis (SBS)

Pyrosequencing

Sequence by Ligation

Slide from Jonathan Eisen (UCD Bodega Bay Workshop in Applied Phylogenetics 2015)

Next-generation sequencing (NGS) – paradigm shift but still relies on PCR



https://www.abmgood.com/marketing/knowledge_base/next_generation_sequencing_introduction.php

Clonal amplification step









Cyclic Array Sequencing

Each method uses different chemistry

AA

sequencing introduction.php

Illumina/Solexa



Illumina HiSeq

8-lane flow cell

Source: www.illumina.com

Prepare Genomic DNA Sample – Library prep



Prepare genomic DNA sample

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

- Randomly fragment genomic DNA and ligate adapters to both end of fragments
- Size selection

Load samples into Flow Cell



Bridge amplification to generate clusters (amplify signal)



Cluster Generation





Cluster Generation



Genome Analyzer Sequencing reaction





Sequencing over Multiple Chemistry Cycles

a Illumina/Solexa — Reversible terminators



****** Generate lots of image files

Metzker, Nat. Rev. Genet, 11:31-46

Multiplexing is easy – barcode and pool your libraries

- Allow sequencing of multiple samples per lane
- Add 4-6 base barcode/index for sample identification
- Single index vs Dual-index (pooling more samples per lane)
- Reads are "de-multiplexed" before assembly

Single-end vs. Paired-end sequencing



http://nextgen.mgh.harvard.edu/IlluminaChemistry.html

Illumina – Model Specifications

	MiSeq	NextSeq 500		HiSeq 2500		HiSeq 3000	HiSeq 4000	HiSeq X	NovaSeq
Run Mode	N/A	Mid-Output	High-Output	Rapid Run	High-Output	N/A	N/A		
Flow Cells Per Rnu	1	1	1	1 or 2	1 or 2	1	1 or 2		
Output Range	0.3-15 Gb	20-39 Gb	30-120 Gb	10-300 Gb	50-1000 Gb	125-750 Gb	125-1500 Gb	1800 Gb	6000Gb
Run Time	5-55 hrs	15-26 hrs	12-30 hrs	7-60 hrs	<1-6 days	<1-3.5 days	<1-3.5 days	< 3 days	19-40 hrs
Reads per Flow Cell	25million	130 million	400 million	300 million	2 billion	2.5 billion	2.5 billion	3 billion	10 billion
Maximum Read Length	2 x 300bp	2 x 150bp	2 x 150bp	2 x 250bp	2 x 125bp	2 x 150bp	2 x 150bp	2X150bp	2X150bp

14



New ways to improve assembly even with Illumina short reads



10X genomics LIT00003 Rev B Chromium Genome Solution Application Note

GemCode[™] Technology for Partitioning High Molecular Weight DNA

Develop high-quality libraries from as little as 1 ng of genomic DNA, ready for standard short-read sequencing.



GemCode[™] Technology for Single Cell Partitioning

Utilize an efficient droplet-based system to encapsulate up to 100-80,000+ cells in a single 10-minute run.



LIT00001 Rev C 10X Genomics Chromium System Brochure

454 Pyrosequencing Chemistry and Base Calling



4-mer

3-mer

2-mer

1-mer

Ion Torrent: Ion Chip Non-Optical Sequencing

- Thermo Fisher
- Natural nucleotides
- Leverage semiconductor manufacturing techniques
- Non-optical but still requires amplification





Ion Chef (library prep)





Ion Proton

lon S5

Sequencing by detecting pH change (release of proton during extension)

sequentially flow unlabeled dNTPs





Reuter et al. (Mol Cell 2015)

Sensor, Well and Chip Architecture



3rd Generation Sequencing

- Longer read lengths, easier assembly
- Single molecule sequencing amplification not necessary, can handle regions with high levels of repeat sequences



Pacific Biosciences (PacBio RS II)



Pacific Biosciences

(Sequel)



Oxford Nanopore (MinION) - \$1000

Oxford Nanopore (PromethION)



Oxford Nanopore (SmidgION) – use with smartphone

PacBIO -SMRT[™] Sequencing



Regions with lots of repeats = easier assembly

From: PacBio_RSII_Brochure

Key Innovations for PacBIO -SMRT[™] Sequencing

Phospholinked nucleotides



ZMWs = Zero-mode Waveguides

DNA Polymerase as a Sequencing Engine





ZMW with DNA polymerase DNA polymerase is immobilized

ZMW with DNA polymerase and phospholinked nucleotides

DNA Polymerase Processive Synthesis with Phospholinked Nucleotides



Base-specific fluorescence and DNA sequence can be detected in real-time

Reuter et al. (Mol Cell 2015)

Nanopore Sequencing

- Oxford Nanopore
- First commercial product = MinION





http://www.medgadget.com/

Different Nanopore Base-Readers



Schneider & Dekker (Nat Methods 2012)

Nanopore Sequencing



Library prep (~1.5 hours):

- DNA fragmentation
- Ligation of 2 adaptors for attachment to motor enzyme and HP motor
- HP adaptor (red) enables the sequencing of both DNA strands (blue and yellow)

How long are we talking about?



http://omicsomics.blogspot.com/2017/03/minion-leviathan-reads-update.html

Start using MinION for \$1000



Min**ION**

- Pocket-sized, portable device for biological analysis
- Up to 512 nanopore channels
- Simple 10-minute sample prep available
- Real-time analysis for rapid, efficient workflows
- Adaptable to direct DNA or RNA sequencing

Start using MinION

https://nanoporetech.com/products

What to consider when using sequencing platforms?



Others??

Van Dijk et al (Trends Genet 2014)

Types of Data for Molecular Phylogenetics

A few mitochondrial or nuclear genes



Maybe somewhere in the middle? Reduced Representation sequencing

Whole Genomes

Reduced Representation Sequencing for Molecular Phylogenetics

- Transcriptome Sequencing (RNA)
- Restriction Site Associated DNA (RAD) Sequencing
- Ultra Conserved Element (UCE) Sequencing
- Anchored Hybrid Enrichment (AHE) Sequencing



Restriction Site Associated DNA (RAD) Sequencing



 Ligation of P1 adapters (one barcoded adapter/individual)

3. Pooling of individual, shearing (300-800 bp) and ligation of P2 adapters



5(a) Single end assemblies 108 bp contigs





5(b) Paired end assemblies 108 + 400 bp contigs

2.2.2.2.2

Shotgun sequencing

For review, also see: Andrews et al. (Nat Rev Genet, 17: 81-92, 2016)

Rowe et al. (Mol Ecol, 2011)

RAD-seq for phylogenomics



Fig. 3 The phylogeny produced based on the largest supermatrix analysed, which contains a minimum of 15 individuals out of the total 156 with sequence data per locus ('min individuals 15'; Table 1g).

Wagner et al. (Mol Ecol, 22:787-798, 2013)

Probe capture and library enrichment (UCE and AHE)

- DNA sequencing library is heat-denatured in the presence of adapter-specific blocking oligonucleotides
- Library and blockers are dropped to the hybridization temperature, allowing blockers to hybridize to the library adapters

 Biotinylated RNA baits are introduced and allowed to hybridize to targets for several hours



http://www.mycroarray.com/mybaits/mybaits-technology.html

Probe capture and library enrichment

 Bait-target hybrids are pulled out of the solution with streptavidin-coated magnetic beads

 Beads are stringently washed several times to remove non-hybridized and nonspecifically-hybridized molecules

- Captured DNA library is released from the beads and amplified
- ****** Design Custom Probe Sets

http://www.mycroarray.com/mybaits/mybaits-technology.html

streptavidincoated magnetic bead

Ultra Conserved Element (UCE) for Phylogenetics



- 481 regions perfectly conserved over 200bp or more between human, mouse, and rat
- 20-fold fewer SNPs than the human average
- Most are in non-coding regions, but some in exonic regions
- May be enhancer sequences, control gene expression
- The exonic UCEs tend to overlap with alternative spliced exons

Human Genome Ultraconserved Elements Are Ultraselected

Sol Katzman,¹* Andrew D. Kern,²* Gill Bejerano,²† Ginger Fewell,³ Lucinda Fulton,³ Richard K. Wilson,³ Sofie R. Salama,^{2,4} David Haussler^{1,2,4}‡



Data for the human genome

Data for primates (Faircloth et al. Syst Biol 61(5): 717-726, 2012)



Workflow for UCE phylogenomics



Probe sets to resolve different levels of phylogenetic relationships need to be adjusted

Faircloth et al. Syst Biol 61(5): 717-726, 2012



n 🤉

UCEs were used to generate phylogeny for Hymenoptera

- Designed 1510 UCEs
- 2749 RNA probes
- Avg. 721 UCEs from 30 taxa
- 400 ng DNA input (avg)
- 2 runs of PE250 on MiSeq
- ~4X coverage

Faircloth et al. (Mol Ecol Res, 15:489-501, 2015)

Anchored Hybrid Enrichment (AHE) for Phylogenetics



CENTER FOR ANCHORED PHYLOGENOMICS

ACCELERATING THE RESOLUTION OF LIFE™



Agilent SureSelect Target Enrichment

Nuclear loci mostly designed from a few reference taxa

As low as 50ng DNA But best 2ug DNA

http://anchoredphylogeny.com/



AHE target enrichment using Agilent SureSelect

Originally used for exome capture

Bird data set (Prum et al. 2015):

- 394 vertebrate loci (use 259)
- ~1350 bp per locus
- HiSeq 2000 PE150, 4 lanes



Which method to use? (Caution)

A few mitochondrial or nuclear genes

Reduced Representation sequencing UCE, AHE, transcriptome, RADseq

** Method used could result in different Tree toplogy

Whole Genomes

STRUCTURE DETAILS



http://nextgen.mgh.harvard.edu/IlluminaChemistry.html