Teaching Computational Genomics: A Tale of Tables and Tests

Texas A&M Research Computing Symposium 2021-05-24

Rodolfo Aramayo - Department of Biology - Texas A&M University

The Genomic Revolution





We Sequenced our own Genome!

Milestone 1 2001

The Human Genome Project



Credit: Rawpixel Ltd / Alamy Stock Photo

Launched in 1990, the Human Genome Project set out to identify the order, that is, sequence, of all DNA bases to obtain the 'genetic blueprint' of humans. In 2001, two pivotal publications reported the first draft of the human genome, obtained by shotgun sequencing, setting the stage for the genomic era. The second phase of the project, which moved from the draft to an essentially finished reference genome, was completed in 2003. Read more.

By Caroline Barranco

We sequenced what was previously impossible!

Milestone 2 2004

Sequencing the unculturable majority

Two key studies unlocked the field of metagenomics – the reconstruction of microbial communities from sequencing data – by providing approaches for unbiased, culture-independent analysis of DNA directly from environmental samples using sequencing technologies. <u>Read more.</u>



Credit: Science Photo Library / Alamy Stock Photo

By lain Dickson

We invented new sequencing technologies!

Milestone 3 2005

Sequencing – the next generation

Two revolutionary studies introduced high-throughput, massively parallel sequencing technologies able to sequence a bacterial genome at a fraction of the cost and time of traditional Sanger sequencing techniques. <u>Read more.</u>

By Joseph Willson



Credit: Zoran Obradovic / Alamy Stock Photo

We understood the need and value of being inclusive!



The dawn of personal genomes



Credit: Redmond Durrell / Alamy Stock Photo

Two studies reported the genomes of an African individual and an Asian individual, respectively, using a new massively parallel sequencing approach based on reversible terminator dyes. Demonstrating the feasibility and resource value of human genome sequences, these studies and the technology they presented paved the way for population-scale genome sequencing. Read more.

By Darren Burgess

Sequencing changed our understanding of disease

Milestone 6 2008

A sequencing revolution in cancer

Ley et al. presented the first wholegenome sequence of a cytogenetically normal acute myeloid leukaemia sample, showing that cancer genome sequencing can identify disease-associated mutations and druggable targets, offering promise for personalized medicine approaches. <u>Read more.</u>

By Safia Danovi



Credit: Sabena Jane Blackbird / Alamy Stock Photo

Sequencing changed our understanding of Gene expression and Gene models



Milestone 7 2008

Transcriptomes – a new layer of complexity

A series of milestone publications reported the development of highthroughput sequencing of whole transcriptomes, known as RNA sequencing (RNA-seq), across different species. <u>Read more.</u>

By Margot Brandt



Credit: Y H Lim / Alamy Stock Photo

We invented even better sequencing technologies!



Long reads become a reality



Credit: Zoonar GmbH / Alamy Stock Photo

Long-read sequencing technologies began to shed light on hidden parts of the human genome by sealing gaps in existing assemblies, allowing modified bases to be detected on native DNA or RNA, and revealing the complexity of the transcriptome. <u>Read more.</u>

By Ivanka Kamenova

We developed single cell Genomics!

Milestone 11 2009

Sequencing one cell at a time

Moving from genomic analysis of tissues or cells in bulk to performing single-cell sequencing provided a whole new perspective on gene regulation, cell-tocell heterogeneity and developmental or disease processes. The difficulty of performing analyses at such resolution required many experimental and computational innovations. <u>Read more.</u>

By Aline Lückgen



Credit: Greenshoots Communications / Alamy Stock Photo

We started sequencing the genomes of our ancestors!



Waking the dead: sequencing archaic hominin genomes



Credit: The Natural History Museum / Alamy Stock Photo

The publication of the first draft genome of a Neanderthal in 2010 marked a turning point for the palaeogenomics field, making it possible to assemble an ancient genome from next-generation sequencing reads by overcoming previous limitations in ancient DNA research such as limited starting material, contamination and degradation. Read more.

By Rebecca Furlong

We re-defined the meaning of reference genomes!



Pan-genomes: moving beyond the reference

Pan-genome studies in a variety of species – from microorganisms to plants to humans – have shown that a large amount of genetic variation can be found in the dispensable genome. This observation has called into question our reliance on single reference genomes for assembling and analysing genomes. <u>Read</u> more.



https://www.nature.com/immersive/d42859-020-0009

html

Credit: Art of Food / Alamy Stock Photo

By Dominique Morneau

We learned how to sequence a chromosome with no gaps!

Milestone 17 2020

Filling in the gaps telomere to telomere

2020 saw the publication of the first gapless, telomere-to-telomere assembly of a human chromosome, the X chromosome. This discovery brought together sequencing technologies and computational tools that had been developed in the preceding decade. <u>Read</u> more.



Credit: PORNCHAI SODA / Alamy Stock Photo

https://www.nature.com/immersive/d42859-020-00099

k.html

By Katharine Wrighton

The Result was...



2009

An explosion of computational tools

As genome sequencing became more affordable and widespread, its applications rapidly expanded, driving the development of new computational tools to accommodate the requirements of transcriptomics, metagenomics or genetic variant discovery. Read mapping tools such as Bowtie and BWA or the splice-aware aligner TopHat were able to align millions of short reads to the reference genome, and downstream analysis software, such as SAMtools and BreakDancer, facilitated the detection of genetic variants.

Related article: <u>Repetitive DNA and next-</u> <u>generation sequencing: computational</u> <u>challenges and solutions</u>



Credit: michalz86 / Alamy Stock Photo

The Price of Genomic Sequencing is Dropping...



Cost per genome data - August 2020



Sequencing cost per megabase - August 2020

Genomic Sequencing Is Generating a Problem The Volume of Incoming Data is Overwhelming



Databases are Growing at an Amazing Rate



Databases are Growing at an Amazing Rate

Homo sapiens

Databasa	Acc	all		
Database	public	controlled	an	
BioSample	<u>3,091,653</u>	<u>651,138</u>	<u>3,742,791</u>	
BioProject	<u>32,639</u>	<u>825</u>	<u>33,464</u>	
dbGaP		<u>18</u>	<u>18</u>	
GEO Datasets	<u>806,746</u>		<u>806,746</u>	

As of: 2021-05-18

https://www.ncbi.nlm.nih.gov/sra/?term=Homo+sapiens

The Reproducibility Versus Replicability Crisis

Reproducibility of Scientific Results

First published Mon Dec 3, 2018

The terms "reproducibility crisis" and "replication crisis" gained currency in conversation and in print over the last decade (e.g., Pashler & Wagenmakers 2012), as disappointing results emerged from large scale reproducibility projects in various medical, life and behavioural sciences (e.g., Open Science Collaboration, OSC 2015). In 2016, a poll conducted by the journal *Nature* reported that more than half (52%) of scientists surveyed believed science was facing a "replication crisis" (Baker 2016). More recently, some authors have moved to more positive terms for describing this episode in science; for example, Vazire (2018) refers instead to a "credibility revolution" highlighting the improved methods and open science practices it has motivated.

The crisis often refers collectively to at least the following things:

- a. the virtual absence of replication studies in the published literature in many scientific fields (e.g., Makel, Plucker, & Hegarty 2012),
- b. widespread failure to reproduce results of published studies in large systematic replication projects (e.g., OSC 2015; Begley & Ellis 2012),
- c. evidence of publication bias (Fanelli 2010a),
- d. a high prevalence of "questionable research practices", which inflate the rate of false positives in the literature (Simmons, Nelson, & Simonsohn 2011; John, Loewenstein, & Prelec 2012; Agnoli et al. 2017; Fraser et al. 2018), and
- e. the documented lack of transparency and completeness in the reporting of methods, data and analysis in scientific publication (Bakker & Wicherts 2011; Nuijten et al. 2016).

https://plato.stanford.edu/entries/scientific-reproducibility/

There is a Divide Between Biology, Computational Sciences and Statistics And

Biology is developing at an amazing rate, becoming increasingly interdisciplinary



How do we do that? Where do we start?

A Solution...

To Gap The Divide Between **Biology, Computer Sciences and Statistics** We should start by Teaching **Hypothesis-driven Computational Genomics Courses**



- Computational Genomics gaps the divide between Biology, Statistics and Computational Sciences
- Computational experiments can and should be designed like any other wet-lab experiment
- The belief that data processing is separated, or different from, wetlab experiments is biased, damaging and wrong
- The rapid evolution and availability of new algorithms, and computational tools make previous data analyses potentially obsolete
- This generates the need to re-test and re-analyze already produced genome data, and provide us with a unique opportunity for training

A Solution...



Graph/Display Data

History of Genomics

Generate tens of millions of sequence reads

Assemble

https://www.nature.com/articles/35084503.pdf

Library Construction

• Sequence reads:

CAGCTGTCCCAGATGAC AACTTCCCTCCCAGCT TCCGCCTTCAGCTCAAGACTTAACTTC CGGCCTTTGGGCTCC ACTTAACTTCCCTCCCAGCTGTCC TCCCAGCTGTCCCAGATGACGCCATC CAGATGACGCC CGGCCTTTGGGCTCCGCCTTCAGCTCAAGA GGGCTCCGCCTTCAGCTC

• Match up overlaps:

CGGCCTTTGGGCTCCGCCTTCAGCTCAAGA AACTTCCCTCCCAGCT CAGATGACGCC TCCGCCTTCAGCTCAAGACTTAACTTC TCCCAGCTGTCCCAGATGACGCCATC GGGCTCCGCCTTCAGCTC ACTTAACTTCCCTCCCAGCTGTCC CGGCCTTTGGGCTCC CAGCTGACC

• Contig:

http://training.ensembl.org/events/2021/2021-05-18-OpenVirtualBrowser_May

http://training.ensembl.org/events/2021/2021-05-18-OpenVirtualBrowser_May

http://training.ensembl.org/events/2021/2021-05-18-OpenVirtualBrowser_May

Final Project Hypothesis: The number of conserved proteins is directly proportional to the evolutionary distance between the proteomes tested

All students get the same (control) proteome and then each student gets two different unique proteomes

Key Concepts Relationships: Tables to Graphs Tables and Text Manipulations The Canonical Gene Table containing pronoter Field 12 Field 01 Record 01 gene Record 01 EDEN 8 9 10 11 12 -Head -2 nRNAs EDEN.1 3 4 EDEN.2 5 Field 01 EDEN.3 (CDS 1) 6 Record 08 7 EDEN.3 (CDS 2) 9 10 11 12 13 Tail -4 ctg123 . gene 1000 9000 . + ID=gene00001;Name=EDEN • Sorting by Colors 14 ctg123 . TF_binding_site 1000 1012 Parent=gene00001 15 ctq123 . mRNA 1050 9000 ID=mRNA00001; Parent=gene00001 On Field 04 + . ID=mRNA00002;Parent=gene00001 .ctg123 .mRNA 1050 9000 . . ctg123 . mRNA 1300 9000 ID=mRNA00003;Parent=gene00001 . . cta123 exon 1300 1500 Parent=mRNA00003 Previous Head 9 10 11 12 1 2 3 New Head -2 1050 1500 Parent=mRNA00001,mRNA00002 .ctg123 . exon . . 3 ctq123 . exon 3000 3902 Parent=mRNA00001,mRNA00003 . . 5 5000 5500 Parent=mRNA00001, mRNA00002, mRNA00003 ctg123 . exon + . 7 ctg123 . exon 7000 9000 Parent=mRNA00001,mRNA00002,mRNA00003 . ctg123 . CDS 1201 1500 . + ID=cds00001:Parent=mRNA00001 14 0 -Previous Tail ID=cds00001;Parent=mRNA00001 ctg123 . CDS 3000 3902 . + 0 15 ctq123 . CDS 5000 5500 . + 0 ID=cds00001;Parent=mRNA00001 Previous Head 2 ctg123 . CDS 7000 7600 . + 0 ID=cds00001:Parent=mRNA00001 6 ctg123 . CDS 1201 1500 . + 0 ID=cds00002;Parent=mRNA00002 9 ctg123 . CDS ctg123 . CDS 5000 5500 . + 0 ID=cds00002;Parent=mRNA00002 10 7000 7600 . + 0 ID=cds00002;Parent=mRNA00002 ctg123 . CDS 3301 3902 . + 0 ID=cds00003;Parent=mRNA00003 ctg123 . CDS 5000 5500 . + 1 ID=cds00003:Parent=mRNA00003 8 ctg123 . CDS 7000 7600 . + 1 ID=cds00003;Parent=mRNA00003 11 ctg123 . CDS 3391 3902 . + 0 ID=cds00004;Parent=mRNA00003 New Tail -4 12 ctg123 . CDS 5000 5500 . + 1 ID=cds00004;Parent=mRNA00003 Previous Tail ctg123 . CDS 7000 7600 . + 1 ID=cds00004;Parent=mRNA00003

12 fields (columns) and 15 records (rows or lines)

History	C + 🗆 🗘
search datasets	00
Exam01_Assemblies	
35 shown, 3 hidden	
83.81 MB	
38:	⊛ & ×
37: Sort on data 34	④ ∦ ×
36: Sort on data 33	● / ×
35: Sort on data 32	● / ×
34: Compute sequence length on data 28	● / ×
33: Compute sequence length on data 20	● / ×
32: Compute sequence length on data 12	● / ×
31:	● / ×
29: VelvetOptimiser on data 5 and data 6: Contig Stats	● / ×
28: VelvetOptimiser on data 5 and data 6: Contigs	● / ×
27: FastQC on data 6: RawData	● / ×
26: FastQC on data 6: Webpage	● / ×

ling							
	Grading Galaxy Work						
	Number of Students	Number of Questions	Number of Processes/ Question	Number of Processes to Grade	Times Six Exams/ Semester		
	1	1	5	5	30		
	50	1	5	250	1500		
	50	10	5	2500	15000		
	65	10	5	3250	19500		

- Grading is Hard because there are different • ways to solve the same problem
- Scripting grading is possible but requires files to have consistent names! (...which requires students to name files correctly, which assumes that all students will follow instructions...)

About Digital Biology (BIOL647) TT :::: Supercomputer :::: :::: **Command-line Driven** Internet **Virtual Machines** 6 13:10:35 on ttys003 84 6 15:11:33 on ttys0 **Students** VPN 346656 May 4 2009 /bin/sh Terminal - tcsh 2 % whereis tcsh GitHub etstream

CYVERSE"

About Digital Biology (BIOL647)

Module 01	Module 02	Module 03
Virtual N (Jetst	Supercomputer (Ada / Grace)	
Master the C	Introduction to Supercomputers	
Basic GNU/Linux/Unix Introduction Basic Scripting		Advanced Scripting

Jetstream work is supported by XSEDE Grant N° BIO210035 Funded by NSF ACI-1445604, Awarded to XSEDE

About Digital Biology (BIOL647)

About Digital Biology (BIOL647)

Genome Files, Genes and Genome Browsers

NCBI Databases and NGS Quality Control

Introduction to and Scripting of Transcriptome Mapping, Assembling and Quantification (Supercomputer)

Who Takes BIOL350, BIOL650, and BIOL647?

Problems Encountered

- Galaxy is intrinsically collaborative, which makes it hard to "isolate" students contributions and/or interactions
- Because we are using the same Galaxy instantiation for Exams and Teaching, we cannot stop students from continuing processing Examrelated data, once a given exam is completed
- Currently, we cannot unmistakably document and isolate a given student's work on a given exam
- Other Supercomputer users sometimes submit unreal amount of jobs to Ada/Terra while an exam is in progress

Potential Solutions

- To have the ability to float one Galaxy instantiation/student/exam
- Dedicated nodes for Galaxy during exams

Other Observations/Questions

• To access the Supercomputer students need to apply for a Supercomputer account.

After the semester is over, the majority of students will not be using their Supercomputer account again

Is this an unnecessary overhead on the Supercomputer facility?

Should we have a temporary Supercomputer accounts for teaching?

We urgently need to implement the local deployment of virtual machines and, ideally, to have those machines outside the firewall

Acknowledgements

This talk is dedicated to the memory of

Dr. Jim Hu

My friend and colleague

With whom I developed and taught Genomics for many, many years

Jim's constant sense of humor, and his fearless attitude towards implementing new ideas was the driving force behind these courses

Acknowledgements

I cannot thank enough the Texas A&M High Performance Research Computing facility for their unwavering, constant, support and continuous help

> Texas A&M High Performance Research Computing is one of the jewels of this institution

> > **Special Thanks go to:**

Michael Dickens Lisa Perez Francis Dang Mark Huang Keith E. Jackson Robin Burns Honggao Liu