Big Data in Genomics and Life Sciences

Alvis Brāzma European Bioinformatics Institute European Molecular Biology Laboratory

Rīga, July 2022



European Molecular Biology Laboratory - EMBL

Europe's centre of excellence in *life sciences* research, services and training

"I believe that international activity is very important in building world peace." Sir John Kendrew, EMBL's 1st DG

Founded in 1974 by 10 states as an intergovernmental organisation to promote the molecular life sciences in Europe and beyond

First Director General Sir John Kendrew Nobel Prize in 1962



EMBL member states in 2022





Six sites with almost 1900 people and >90 nationalities





Wellcome Genome Campus in Hinxton, near Cambridge, UK





What is EMBL-EBI?

- World leading source of public biomolecular data
- Cutting edge research in bioinformatics and computational biology
- Part of the European Molecular Biology Laboratory (EMBL), Europe's flagship laboratory for the life sciences.



Bioinformatics is interdisciplinary field that develops methods and software tools for understanding biological data (Wikipedia 2022)

Data science in biology



Bioinformatics: from fringe to mainstream



Central dogma of molecular biology and OMICs

 $\dots gatgatgatgatgatgatgatgatgatgatgatgaggacaactctctttttccaacaagagagccaagaagccatttttttccatttgatctgtttccaatg\ldots$





Human Genome project was an epic

- 1990 2003 about 85% of the **3 billion letters** of human genome sequenced
 - The first "gapless" human genome sequence was achieved only in 2022
- HGP budget over the 13 years was \$ 3 billion, though the actual sequencing costs were considerably less
- Some lessons from HGP
 - There are 20,000 protein coding genes and some other genes (most shared with other mammals)
 - Less than 2% of the human genome codes for proteins (other parts are interesting too)
 - Over a half of human genome consist of repeated sequences of different length apparently not having any function (junk genome)
- Sequencing one human genome in 2022 is < \$1000
- Genomes of two unrelated individuals differ in 1 "letter" per 3000 on average



Detailed view

www.ensemble.org



Features V Comparative VDAS Sources Repeats Decorations Export Image size Help 🔻 : 52581629 - 52599222 Refresh Refresh Jump to region 4 Band: << 5MB Window 1MB >2MB >5MB >> < 2MB < Window Chr. 4 52.58 Mb 52.58 Mb 52.59 Mb 52.59 Mb 52.59 Mb 52.59 Mb 52.59 Mb 52.60 Mb 52.60 Mb Length Forward strand 17.59 Kb DNA(contigs) < \$GCB_HUMAN Havana Known Protein coding Ensembl trans. < SGCB HUMAN Ensembl Known Protein Coding UniProtKB Eponine FirstEF CpG is lands SNPs Length Revense strand ` 52.58 Mb 52.59 Mb 52.59 Mb 52.59 Mb 52.60 Mb 52.58 Mb 52.59 Mb 52.59 Mb 52.60 Mb Havana Known Protein coding Gene legend Ensembl Known Protein Coding SNP legend 3' UTR Intronic Non-synonymous coding EnsEMBL Homo sapiens version 44.36f (NCBI 36) Chromosome 4 52,581,629 - 52,599,222



Personalized genomics - assembling a new genome vs. mapping sequencing "reads" to a reference genome

- The output of a DNA sequencing machine is a set of many ~100 letter long sequences that jointly cover (most of the genome) 10 to 50 times on average
 - Current sequencing technologies output pairs of such ~100 letter sequences with approximately known distance between them



The "insert" may help to deal with repeated sequences in the genome

EMB

- Once we have one reference genome, to get another we can "map" these "reads" onto the reference to obtain the new one
- Mismatches of letters in the "read" can be genome variation or errors
- If we can "trust" the reference? Many references for different populations

Personal genomics: what do genomes tell us?





Genomics promises a leap forward for rare disease diagnosis

Faster and cheaper DNA sequencing brings new hope to patients



Jessica suffers from a rare condition that was diagnosed through DNA analysis

"Kate Palmer and Simon Wright were in despair. Their four-year-old daughter Jessica was suffering from epilepsy, poorly co-ordinated movement and slow mental development, but doctors had been unable to pinpoint the rare disorder causing these problems."





- Jessica was enrolled in the UK 100,000 genome project
- By analysing Jessica's genome, a mutation in a gene called SLC2A1 was found starving her brain of the sugar
- This condition is extremely rare, but there is a treatment in the form of a diet that enables the brain to maximise glucose production
- After a month on the new diet, Jessica's parents "noticed a big increase in her speech, energy levels and general steadiness"



Sharing personal genome data

- While reference are an average of many genomes, individual genomes are unique and can be linked to person and other information
- Implication to privacy, health insurance, etc
- Pseudonymisation (or even anonymization) is not sufficient to be sure that one's genome identity is not revealed – linking to health records, etc
- Consent and law
- Controlled access by bone-fide researchers
- In many countries personal genome data are not allowed to leave the country by law
- Solutions: federation of data, common standards, cloud computing



Know your genome

• Figuring out which regions are involved in disease – and what they do – is a major challenge.





European Genome-Phenome Archive

- Resource for permanent secure archiving and sharing of all types of potentially identifiable human genetic and phenotypic data
- Distributed data access model
- Access to bona-fide researchers controlled by Data Access Committees

How the EGA is managed

The EGA was launched in 2008 by the EBI



EUROPEAN

ARCHIVE

GENOME-PHENOME

10THANNIVERSARY



Federated European Genome-phenome Archive (FEGA) Central EGA node >20 countries / national Federated EGA node initiatives and counting! Preparing to sign FEGA **Collaboration Agreement** Engaging in work to establish a FEGA node Expressing interest in joining FEGA Network

Beyond 1 Million Genomes Project B1MG

Building towards genomic data infrastructure





Data generated for research project vs. data generated as a part of healthcare

Research data:

Disease Associations Molecular biology resource



Healthcare genomics data and electronic health records

Feedback to patients Actionable variants

Iceland; Denmark; Faroe; Finland; Dundee; UK BioBank; (others)





The Global Alliance for Genomics and Health (GA4GH) is an international, nonprofit alliance formed in 2013 to accelerate the potential of research and medicine to advance human health. **Bringing together** 600+ leading organizations working in healthcare, research, patient advocacy, life science, and **information technology**, the GA4GH community is working together **to create frameworks and standards to enable the responsible, voluntary, and secure sharing of genomic and health-related data.**

https://www.ga4gh.org/



Not only humans



To sequence a 1.5 million known eukaryotic species in 10 years

Earth Biogenome Project (EBP)

The Earth BioGenome Project (EBP), a moonshot for biology, aims to sequence, catalog and characterize the genomes of all of Earth's eukaryotic biodiversity over a period of ten years. More information can be found in the EBP official web portal at https://www.earthbiogenome.org

EBP Contribution to Eukaryotic Genome Sequencing

Progress of Eukaryotic Genome Sequencing by Taxon Rank: all assemblies in INSDC and those submitted under the EBP umbrella (BioProject PRJNA533106)







Contribution of Projects under EBP Umbrella Project Id PRJNA533106, corresponding to 17 out of the 49 affiliated projects in the EBP Network





Cancer Genomics

- Cancer is a genome disease
- Human has ~10¹³ cells, many growing and dividing, there are 2-3 mutations per every cell division
- As important genes, such as DNA-repair genes, get mutated, more mutations accumulate, and genomes get increasingly changed
- Roughly 300 cancer genes are known, mutation in which can lead to cancer
- Some cells escape control and start proliferate tumors
- Some go on to migrate metastasis



PERSPECTIVES

International network of cancer genome projects

The International Cancer Genome Consortium (ICGC) was launched to coordinate large-scale cancer genome studies in tumours from 50 different cancer types and/or subtypes that are of clinical and societal importance across the globe. Systematic studies of more than 25,000 cancer genomes at the genomic, epigenomic and transcriptomic levels will reveal the repertoire of oncogenic mutations, uncover traces of the mutagenic influences, define clinically relevant subtypes for prognosis and therapeutic management, and enable the development of new cancer therapies

A subset of the somatic mutations in cancer cells confers oncogenic properties such as growth advantage, tissue invasion and metastasis, angiogenesis, and evasion of apoptosis². These are termed 'driver' mutations. The identification of driver mutations will provide insights into cancer biology and highlight new drug targets and diagnostic tests. Knowledge of cancer mutations has already led to the development of specific therapies, such as trastuzumab for *HER2* (also known as *NEUL* or *EBBR2*) provide heaven and in stirily arbitrary and in stirily. org/files/ICGC_April_29_2008.pdf). Since then, working groups and initial member projects have further refined the policies and plans for international collaboration.

Bioethical framework

ICGC members agreed to a core set of bioethical elements for consent as a precondition of membership (Box 2). The Ethics and Policy



Pan-Cancer Analysis of Whole Genomes of ICGC completed 2020

The international journal of science / 6 February 2020

nature

CANCER CATALOGUED

Whole-genome sequences for 38 types of tumour

Article

Genomic basis for RNA alterations in cancer

https://doi.org/10.1038/s41586-020-1970-0

Accepted: 11 December 2019

Received: 29 March 2018

Published online: 5 February 2020

Open access

 $\Gamma > G =$

C > T =

PCAWG Transcriptome Core Group^{1,35}, Claudia Calabrese^{2,35}, Natalie R. Davidson^{3,4,5,6,7,35}, Deniz Demircioğlu^{8,9,35}, Nuno A. Fonseca^{2,35}, Yao He^{10,35}, André Kahles^{3,4,6,7,35}, Kjong-Van Lehmann^{3,4,6,7,35}, Fenglin Liu^{10,35}, Yuichi Shiraishi^{11,35}, Cameron M. Soulette^{12,35}, Lara Urban^{2,35}, Liliana Greger², Siliang Li^{13,14}, Dongbing Liu^{13,14}, Marc D. Perry^{15,16}, Qian Xiang¹⁵, Fan Zhang¹⁰, Junjun Zhang¹⁵, Peter Bailey¹⁷, Serap Erkek¹⁸, Katherine A. Hoadley¹⁹, Yong Hou^{13,14}, Matthew R. Huska²⁰, Helena Kilpinen²¹, Jan O. Korbel¹⁸, Maximillian G. Marin¹², Julia Markowski²⁰, Tannistha Nandi⁹, Qiang Pan-Hammarström^{13,22}, Chandra Sekhar Pedamallu^{23,28,29}, Reiner Siebert²⁴, Stefan G. Stark^{3,4,6,7}, Hong Su^{13,14}, Patrick Tan^{9,25}, Sebastian M. Waszak¹⁸, Christina Yung¹⁵, Shida Zhu^{13,14}, Philip Awadalla^{15,26}, Chad J. Creighton²⁷, Matthew Meyerson^{22,28,29}, B. F. Francis Ouellette³⁰, Kui Wu^{13,14}, Huanming Yang¹³, PCAWG Transcriptome Working Group¹, Alvis Brazma^{2,36*}, Angela N. Brooks^{12,23,28,36*}, Jonathan Göke^{9,31,36}, Gunnar Rätsch^{3,4,5,6,7,36*}, Roland F. Schwarz^{2,20,32,33,8}, Oliver Stegle^{218,33,36}, Zemin Zhang^{10,36} & PCAWG Consortium³⁴

Transcript alterations often result from somatic changes in cancer genomes¹. Various forms of RNA alterations have been described in cancer, including overexpression², altered splicing³ and gene fusions⁴; however, it is difficult to attribute these to underlying genomic changes owing to heterogeneity among patients and tumour types, and the relatively small cohorts of patients for whom samples have been analysed by both transcriptome and whole-genome sequencing. Here we present, to our knowledge, the most comprehensive catalogue of cancer-associated gene alterations to date, obtained by characterizing tumour transcriptomes from 1,188 donors of the Pan-Cancer Analysis of Whole Genomes (PCAWG) Consortium of the International Cancer Genome Consortium (ICGC) and The Cancer Genome Atlas (TCGA)⁵. Using matched whole-genome sequencing data, we associated several categories of RNA alterations with germline and somatic DNA alterations, and identified probable genetic mechanisms. Somatic copy-number alterations were the

Cancer Genomics

- Cancer genomics is complicated
 - Because of clonality there is not longer one genome
 - To find what is changed in cancer cells we compare normal genome (e.g., from blood cells) to tumor genomes personalized genome becomes the rererence





Cancer related genome rearrangements cause genes to fuse creating "cancer genes"

• Hybrid genes formed from two previously separate genes as a result from of genomic rearrangements, read-through transcription and trans-splicing.



Fusion genes are tumor-specific and therefore important targets for therapy.



Central dogma of molecular biology - transcriptomics





$\overleftarrow{\leftarrow}$ \rightarrow G													
						🖶 EMBL-EBI	 Services 	🕸 Research	Å Training	i About us Q	EMBL-EBI		
0	Expres Gene express	SION Atla	3S ecies and bio	ological co	onditions				Que To s	Ery single cell	expressior	n	
🔺 Home	Browse experiments	土 Download 🔹	Release notes	£ FAQ Ø	Help 🛛 🗹 Lice	nce 🚯 Abou	t 🔀 Suppo	ort					
Search ad	Search across 65 species, 4,052 studies, 134,900 assays Ensembl 99, Ensembl Genomes 46, WormBase ParaSite 14, EFO 3.10.0												
Search	Gene set enrichment												
Gene / Gene properties Species Biological conditions													
REG1B × Homo sapiens						Enter condition query							
Examples: REG1B, zinc finger, O14777 (UniProt), GO:0010468 (regulation of gene expression) Examples: lung, leaf, valproic acid, cancer													
Searc	h Clear												
Animals	Plants Fungi												
					-		3			J			
	Homo sapiens	Mus muscul	us F	attus norve	gicus	Drosophi	la	Gallus	gallus	Caenorhabdit	s elegans		
	1449 experiments	1153 experime	nts	152 experime	ents	melanogas	ster	36 exper	iments	29 experir	nents		
	Baseline: 59	Baseline: 46	77	Baseline: 3	3	140 experime	ents 1	Baselin	ne: 3	Baseline	9: 1		
					• • •	Differential: 1	136	Dilleren	iai. 00	Differentia	a. 20		



Tissues consists of cells of different types

THE CELL IS THE FUNDAMENTAL UNIT OF LIFE



It is often said that human has ~200 different cell types



Human body contains >10¹³ cells of 200 different types





From looking at gene expression in tissues to individual cells



Single-cell genomics



Single Cell **RNA** sequencing

а

b

Single cells in study

1

2009

2010

2011



2013

Study publication date

2012

2014

2015

2016

2017

EMBI



Single Cell Expression Atlas



EMBL

TSN plots (non-liner PCR)

Human body contains >10¹³ cells of 200 different types





What is a cell type?

- A cell type is a classification traditionally used to distinguish between morphologically or phenotypically distinct cell forms within a species
- Nowadays cell types are typically defined through sets of genes specifically expressed there, called *marker genes*
- Both definitions are only proxies and there is no generally agreed definition of this concept



Marker genes for immune cells





Automated approaches using known marker genes



These methods do not allow for discovering new cell "types" of states



We wanted to automate the cell type identification process in a way that allows for **discovery of new cell "types"** as well as new marker genes automatically

- Problem of finding the cell clusters corresponding to biologically meaningful cell groups (cell "types")
- Building gene expression "models" for each group which are the genes that are differentially expressed for the group?
- Single Cell Clustering Assessment Framework SCCAF



SCCAF in a nutshell:



Miao et al, Nature Methods, 2020

Comparison with expert annotation of the mouse retina dataset from Shekhar et al. 2016



 MG (Mueller Glia) BC5A (Cone Bipolar cell 5A) BC7 (Cone Bipolar cell 7) BC6 BC5C BC1A BC3B BC1B BC2 BC5D BC3A BC5B BC4 BC8/9 (mixture of BC8 and BC9) AC (Amacrine cell) Rod Photoreceptors Cone Photoreceptors

RBC (Rod Bipolar cell)





Adjusted Rand index >0.99.



Annotating human brain data based on model trained on mouse (Allen Brain Atlas)

SCCAF

- GABAergic_IL1RAPL2
- GABAergic_LAMP5 ANO4
- GABAergic_LAMP5 CHST9
- GABAergic_LAMP5 FRAS1
- GABAergic_LHFPL3
- GABAergic_LUZP2
- GABAergic_PVALB SLIT2
- GABAergic_SST GRIK1
- GABAergic_SST NPY
- GABAergic_VIP OLFM3
- GABAergic_ZFPM2
- Non-Neuronal_Astro DPP10
- Non-Neuronal_Astro LGR6 TNC
- Non-Neuronal_Macrophage CD74
- Non-Neuronal_Oligo LHFPL3
- Non-Neuronal_Oligo ST18
- Non-Neuronal_Stellate DCN LAMA2
- Non-Neuronal_mixed DNAH17

- Glutamatergic_KC/MC/PC SAMD5
- Glutamatergic_L2/3 COL5A2 CA10
- Glutamatergic_L2/3/4ab TESPA1 MEIS2
- Glutamatergic_L4/5 TLL1 PDE4B
- Glutamatergic_L4/5 VWC2L SLC35F3
- Glutamatergic_L4abc SPHKAP
- Glutamatergic_L4abc TSHZ2 ZNF804B
- Glutamatergic_L4c/5 EYA4 CDH20
- Glutamatergic_L4c/5 PDE1C FMN1
- Glutamatergic_L5/6 ARHGAP15 PDZRN4
- Glutamatergic_L5/6 GRM8
- Glutamatergic_L5/6 HTR2C
- Glutamatergic_L5/6 RGS12 ITGB8
- Glutamatergic_L5/6ab ANXA1
- Glutamatergic_L6 PARD3B CDH9
- Glutamatergic_L6 PCSK5 ITPR2
- Glutamatergic_L6 SULF1 ADAMTSL1
- Glutamatergic_L6ab/5 SYT6 SERPINE2
- Glutamatergic_L6ab/5 THEMIS
- Glutamatergic_MC/PC ITGA8

SCCAF has been implemented in the Expression Atlas Galaxy pipeline





The Human Cell Atlas (HCA) is a global partnership of scientists who are actively working to create an exhaustive guidebook of the types and properties of all human cells.







Central dogma of molecular biology







What can we say about protein expression from RNA?



Correlation between RNA and protein abundances are not all that high but can we predict protein abundances from RNA abundances? Moreover, can we predict proteins that are difficult to measure diretly?





Research Article 🙃 Open Access 💿 🔅

Using Deep Learning to Extrapolate Protein Expression Measurements

Mitra Parissa Barzine, Karlis Freivalds, James C. Wright, <u>Mārtiņš Opmanis</u>, <u>Darta Rituma</u>, Fatemeh Zamanzad Ghavidel, Andrew F. Jarnuczak, <u>Edgars Celms</u>, Kārlis Čerāns, Inge Jonassen, Lelde Lace, Juan Antonio Vizcaíno, Jyoti Sharma Choudhary **x**, Alvis Brazma **x**, Juris Viksna **x**... **See fewer authors**

First published: 16 September 2020 | https://doi.org/10.1002/pmic.202000009 | Citations: 1

🔧 TOOLS < SHARE PDF

SECTIONS

Abstract

Mass spectrometry (MS)-based quantitative proteomics experiments typically assay a subset of up to 60% of the ≈20 000 human protein coding genes. Computational methods for imputing the missing values using RNA expression data usually allow only for imputations of proteins measured in at least some of the samples. In silico methods for comprehensively estimating abundances across all proteins are still missing.

Here, a novel method is proposed using deep learning to extrapolate the observed protein expression values in label-free MS experiments to all proteins, leveraging gene functional annotations and RNA measurements as key predictive attributes. This method is tested on four datasets, including human cell lines and human and mouse tissues. This method predicts the protein expression values with average R^2 scores between 0.46 and 0.54, which is significantly better than predictions based on correlations using the RNA expression data alone. Moreover, it is demonstrated that the derived models can be "transferred" across experiments and species. For instance, the model derived from human tissues gave a $R^2 = 0.51$ when applied to mouse tissue data. It is concluded that



Volume 20, Issue 21-22 Special Issue: Computational Proteomics: Focus on Deep Learning November 2020 2000009

Advertisement





Details

 $\ensuremath{\mathbb{C}}$ 2020 The Authors. Proteomics published by Wiley-VCH GmbH



This is a second s

But what about predicting protein expression from RNA and other data?



Barzine et al, Proteomics, Nov 2020



Deep learning allows to make good predictions



$$R^2 = 1 - rac{SS_{
m res}}{SS_{
m tot}}$$

$$SS_{
m res} = \sum_i (y_i - f_i)^2$$



Conclusions

- Biology is increasingly becoming a data science (like in most human activities, there is a shift from material processing to data processing)
- Biological and health data are growing faster than Moore's law creating challenges and opportunities for computer and data scientists
- Growing importance of AI and ML to biological data analysis
- Personalized genomics is having an impact on healthcare now

