# DATA SCIENCE IN PHARMA RESEARCH

ALVARO SEBASTIAN

WOMEN IN MACHINE LEARNING & DATA SCIENCE

21 FEBRUARY 2020 - POZNAN

**SIXTH RESEARCHER**
www.sixthresearcher.com

# THE ORIGINS…

# THE HUMAN GENOME



Craig Venter    Bill Clinton    Francis S. Co

On June 26, 2000, a 'rough draft' of the genome was ann
jointly by U.S. President Bill Clinton (photo) and the British
Minister Tony Blair (via satellite).

https://en.wikipedia.org/wiki/Human_Genome_Project

# THE HUMAN GENOME

# THE HUMAN GENOME

**The total length of the haploid human genome is 3.3 billion base pairs (3.3E9).**

*Don Quixote*, the Spanish novel by Miguel de Cervantes contains around 2 million of letters, so **the human genome has as many letters as 1500 copies of *Don Quixote*.**

There are a lot of letters in that stack and a lot of information that we are trying to understand. For example, **a genetic disease is like having a typo in one of those copies of *Don Quixote*.**

# THE HUMAN GENOME

**The human genome occupies around 750 Megabytes that is about 1 CD of space.**

3×10^9 base pairs/haploid genome x 2 bits/base pairs x 1 byte/8 bits = 0.75E9 bytes

**That is nothing! The Canopy Plant Genome is 50 times bigger!**

| Species | T2 phage | Escherichia coli | Drosophila melanogaster | Homo sapiens | Paris japonica |
|---|---|---|---|---|---|
| Genome Size | 170,000 bp | 4.6 million bp | 130 million bp | 3.2 billion bp | 150 billion bp |
| Common Name | Virus | Bacteria | Fruit fly | Human | Canopy Plant |

# THE HUMAN GENOME

# THE METHODS...

# GENOME SEQUENCING



DNA isolation

DNA fragmentation

Library preparation

Sequencing

Bioinformatics analysis

DNA Reads

Reconstructed Genome

SIXTH RESEARCHER
www.sixthresearcher.com

# GENOME SEQUENCING

# THE PRICE…

# HUMAN GENOME COST

https://www.genome.gov/about-genomics/fact-sheets/Sequencing-Human-Genome-cost

# HUMAN GENOME COST

# MORE GENOMES…

# THE 1000 GENOMES PROJECT

The 1000 Genomes Project, launched in January 2008, consisted in sequencing the genomes of at least one thousand anonymous participants from a number of different ethnic groups.

In 2012, the sequencing of 1092 genomes was announced in a Nature publication



https://en.wikipedia.org/wiki/1000_Genomes_Project

# THE 100,000 GENOMES PROJECT

The project was established by the UK government to sequence 100,000 genomes from NHS patients affected by a rare disease, or cancer.

Recruitment of participants to the 100,000 Genomes Project was completed in 2018, with the 100,000th sequence achieved in December 2018.



The 100,000 Genomes Project

Genomics England & Partners

https://en.wikipedia.org/wiki/100,000_Genomes_Project

# EUROPEAN 1+ MILLION GENOMES

https://ec.europa.eu/digital-single-market/en/european-1-million-genomes-initiative

# THE PURPOSE…

# PERSONALIZED MEDICINE

Genome sequencing can reveal **alterations in DNA that influence diseases** ranging from cystic fibrosis to cancer.

**Personalized medicine** takes advantage of the results from these techniques to design the most appropriate therapy for each patient.

# PERSONALIZED MEDICINE

Percentage of the patient population for which a particular drug in a class is ineffective, on average:



| ANTI-DEPRESSANTS SSRIs | 38% |
| ASTHMA DRUGS | 40% |
| DIABETES DRUGS | 43% |
| ARTHRITIS DRUGS | 50% |
| ALZHEIMER'S DRUGS | 70% |
| CANCER DRUGS | 75% |

# CANCER…

# TYPE OF MUTATIONS

A **germline mutation** is a constitutional mutation that and is transmitted to offspring via the germ cells. **is inherited, present in all the body cells.**

A **somatic mutation** is not inherited from a parent, **is spontaneously generated during life,** and also not passed to offspring.

# MUTATIONS AND CANCER

Genomic Landscape of 5000 Human Cancers:



MacConaill, L. E., Garcia, E., Shivdasani, P., Ducar, M., Adusumilli, R., Breneiser, M., … Lindeman, N. I. (2014). Prospective Enterprise-Level Molecular Genotyping of a Cohort of Cancer Patients. *The Journal of Molecular Diagnostics*, 16(6), 660–672.

# MUTATIONS AND CANCER

Still we do not know many of the cancer driven genomic alterations:



| | | |
|---|---|---|
| Unknown | 34% |
| KRAS Mt | 32% |
| EGFR Mt | 23% |
| ALK Fusion | 3% |
| HER2 Mt | 3% |
| Double Mt | 1% |
| BRAF Mt | 1% |
| ROS1 Fusion | 1% |
| RET Fusion | 1% |
| PIK3CA Mt | <1% |
| MET Amp | <1% |
| MEK1 Mt | <1% |
| NRAS Mt | <1% |

Representative pie charts from molecular diagnostic testing of NSCLC using a combination of assays at Memorial Sloan Kettering Cancer Center (MSKCC). Sanger sequencing, IHC, FISH, multiplex hotspot mutational testing, and multiplex sizing assays were used as part of a diagnostic algorithm for lung adenocarcinomas.

Naidoo, J., & Drilon, A. (2014). Molecular Diagnostic Testing in Non-Small Cell Lung Cancer. The American Journal of Hematology/Oncology, 10(4)(september), 4–11.

# THE RESULTS…

# CANCER THERAPIES

Ongoing improvements in cancer treatments, survivorship up, mortality down:



Sources: US Mortality Files, National Center for Health Statistics, CDC. DeSantis C, Chunchieh L, Mariotto AB, et al. (2014). Cancer Treatment and Survivorship Statistics, 2014. CA: A Cancer Journal for Clinicians.

# CANCER THERAPIES

**Imatinib opened the new era of Cancer Targeted Therapies.** A simple pill putting an end to treatments with serious side effects that had limited success in prolonging life beyond the first year of diagnosis.

A 2011 study concluded that Chronic Myeloid Leukemia (CML) patients whose disease is in remission after 2 years of imatinib treatment have the same life expectancy as those who never had this disease.

Gambacorti-Passerini, C., Antolini, L., Mahon, F.-X., Guilhot, F., Deininger, M., Fava, C., … Kim, D.-W. (2011). Multicenter independent assessment of outcomes in chronic myeloid leukemia patients treated with imatinib. Journal of the National Cancer Institute, 103(7), 553–61.

# CANCER THERAPIES



Chu, H., Zhong, C., Xue, G., Liang, X., Wang, J., Liu, Y., … Bi, J. (2013). Direct sequencing and amplification refractory mutation system for epidermal growth factor receptor mutations in patients with non-small cell lung cancer. Oncology Reports, 30(5), 2311–2315.

# THE CODE…

# CODE EXAMPLES

**Google Colaboratory Notebooks:**

- Exploring human mutations related with cancer:

  https://colab.research.google.com/drive/1xOkGnrLVPiqwj1BfcMgfKdVilUOES5gd

- Looking for EGFR gene mutations at NGS data from lung cancer patients:

  https://colab.research.google.com/drive/1jffxhQoswPEW5-JMMk_y6HFbBL0dLjqD

# THE FUTURE...

# GENE EDITING

# GENE EDITING



Novartis wins approval for world's most expensive drug

May 24 2019

US FDA gives green light for $2.1m treatment of spinal muscular atrophy

FINANCIAL TIMES

# ARTIFICIAL INTELLIGENCE…

# AI IN CANCER DETECTION



McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* **577**, 89–94 (2020)

# AI IN CANCER DETECTION

An AI system that is capable of surpassing human experts in breast cancer prediction. It provides a reduction of 5.7% and 1.2% (USA and UK) in false positives and 9.4% and 2.7% in false negatives.



McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* 577, 89–94 (2020)

# FIRST AI DESIGNED DRUG

In 2016, the pharmaceutical firm Sunovion gave a group of seasoned employees an unusual assignment. At the firm's headquarters in Marlborough, Massachusetts, the chemists were all asked to play a game to see who could discover the best leads for new drugs.

Of the 11 players, 10 struggled through the task for several hours. But one breezed through in milliseconds… because it was an algorithm.



The Drug Makers Guide to the Galaxy: How machine learning and big data are helping chemists search the vast chemical universe for better medicines. Nature 26 SEP, 2017

# FIRST AI DESIGNED DRUG

A drug molecule invented entirely by artificial intelligence is set to enter human clinical trials for the first time, marking a critical milestone for the role of machine learning in medicine.

Four times faster than a typical Drug Discovery process, this AI designed drug to treat patients with obsessive-compulsive disorder will enter its clinical trials 12 months after reseach started.

The new compound was developed by Oxford-based AI start-up Exscientia in collaboration with the Japanese pharmaceutical firm Sumitomo Dainippon Pharma.

Technology

**Artificial intelligence-created medicine to be used on humans for first time**

By Jane Wakefield
Technology reporter

30 January 2020

The drug was much quicker to market than ones developed in more traditional ways

GETTY IMAGES

https://www.bbc.com/news/technology-51315462

# WILL BE AI THE NEW THERANOS?

# WILL BE AI THE NEW THERANOS?



TCR-ANTIGEN MAP

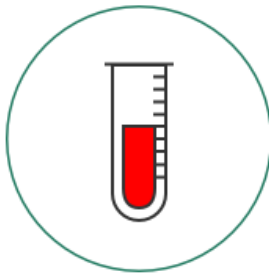Early detection of multiple diseases from a single blood test

Adaptive biotechnologies® + Microsoft

# WILL BE AI THE NEW THERANOS?

Microsoft Healthcare NExT initiative has partnered with Adaptive to map and decode the human immune system, nature's most finely tuned diagnostic. Together we are using immunosequencing, proprietary computational modeling, and machine learning to map T-cell receptor (TCR) sequences to the antigens they bind. Using this data, we aim to translate the natural diagnostic capability of the immune system into the clinic.

## Learning to decode the immune system to diagnose disease

**Blood sample**

The immune system is nature's most finely-tuned diagnostic, providing a fingerprint of a person's health in their blood

**Immunosequencing**

We read immune signatures that store the diagnostic information

**Machine learning**

We generate a map of the immune system by matching trillions of T cells to the diseases they recognize

**Empowering care**

This map of the immune system may be used by doctors and researchers to improve disease diagnosis

# THANK YOU DZIEKUJE GRACIAS

SIXTHRESEARCHER@GMAIL.COM