EHzürich



Data Sharing Strategies & Benefits in ICGC, TCGA and BRCA Challenge

Gunnar Rätsch

Biomedical Informatics Group

@gxr #PrecisionMedicine #Cancer #Genomics #ClinicalData #SPHN





UniversityHospital Zurich



Memorial Sloan-Kettering



Gunnar Rätsch | 30.11.2018 | 1

EHzürich

Biomedical Informatics Lab: Data Science for Biomedical Applications



2

ETH zürich

Towards Comprehensive Patient Models for Precision Medicine



Data Science Research Challenges

Challenge 1: Develop novel data science approaches for medical data

Challenge 2: Provide analysis tools for the community

Challenge 3: Solve important biomedical problems through collaborations

Challenge 4: Create an environment which allows us to work on the above



Source: Center: Google icons search



ETHzürich

a) Large-scale Initiatives



Source: Courtesy of Torsten Schwede



Swiss Personalized Health Initiative Aim: Data interoperability between hospitals and researchers

b) Data Sharing Standards



Global Alliance

GA4GH 6th Plenary Meeting

Source: GA4GH

3 – 5 October 2018, Basel, Switzerland

http://ga4gh.org

Large-scale Cancer Genome Projects (TCGA, ICGC, ICGC-ARGO, ...)

- Two major cancer genome projects started >10 years ago
 - Aim: collect and profile tumor & normal tissue samples
 - Whole Genome Sequencing (WGS), Exome (WXS), RNA-seq, Micro RNAs, some Mass Spec, pathology slides
 - All publicly/controlled accessible
 - <u>TCGA</u>: organized by US National Cancer Institute, total ≈12'000 donors, ≈1PB
 - <u>ICGC</u>: international effort, total ≈25'000 donors, ≈2'500 samples with WGS, ≈1PB
 - ICGC-ARGO: 10x ICGC + detailed medical records
- Hundreds of groups generating and analysing the data, thousands of papers written about it
- Anything left to be done?





Hzürich

Example: International Cancer Genome Consortium



Collaborative Project Principles/Goals

"It is amazing what you can accomplish if you do not care who gets the credit."

— Harry S. Truman

- Be part of something great
- Support junior faculty and students
- Change clinical and research practices
- Have fun working together

Important: Clear rules of collaboration/publication in consortium



Project 1: Integrate Diverse Transcriptomic Alterations to Identify Cancer-relevant Genes and Signatures

PanCancer Analysis of Whole Genomes and Transcriptomes Working Group (PCAWG-3), PCAWG Consortium, Kjong-Van Lehmann^{1,2}, André Kahles^{1,2}, <u>Alvis</u> <u>Brazma³, Angela N. Brooks</u>⁴, Claudia Calabrese³, Nuno A. Fonseca³, Jonathan Göke⁵, Roland F Schwarz^{3,6}, <u>Gunnar Rätsch^{1,2}</u>, Zemin Zhang^{7,8}

¹ETH Zürich, Computer Science Dept, Universitätsstrasse 6, 8092 Zürich, Switzerland;
²Memorial Sloan Kettering Cancer Center, 1275 York Avenue, New York 10065, USA;
³European Molecular Biology Laboratory - European Bioinformatics Institute, Wellcome Trust Genome Campus, Cambridge CB10 1SD, UK;
⁴University of California, Santa Cruz, CA 95060; ⁴BaylorCollege of Medicine, Houston, TX, USA; ⁵Duke-NUS Graduate
Medical School, 8 College Road, Singapore 169857, Singapore; ⁶Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular Medicine, Berlin, Germany; ⁷Peking-Tsinghua Center for Life Sciences, Peking
University, Beijing, 100871, China; ⁸Genome Institute of Singapore, 60 Biopolis Street, Genome #02-01, Singapore 138672, Singapore;

PCAWG-3 Marker Manuscript https://doi.org/10.1101/183889



Gunnar Rätsch | 30.11.2018 | 9

Collaborative Analysis of Multiple Alteration Types



Source: PCAWG-3 Working group

















Courtesy of Natalie Davidson





Gunnar Rätsch | 30.11.2018 | 15

0

1



Binary value for each triplet (sample, gene, alteration).















Gunnar Rätsch | 30.11.2018 | 19

ETHzürich

Compare and contrast samples

➡ Cancer Type





ETH zürich

Identify known and novel recurrently altered genes



Source: PCAWG-3 Working group

<u>Project 2:</u> Cancer-Specific Splicing & Implications in 8,705 Samples

Goals:

- Identify cancer-specific splicing patterns
- Identify variants regulating splicing in the same gene (cis)
- Identify variants regulating splicing in other genes (trans)
- Is splicing relevant for cancer treatments?

The Cancer Genome Atlas provides RNA-seq and matching exome data

- RNA-seq => Find & quantify splicing events
- Exome => Identify variants in exons and flanking intronic regions

Mode of interaction: One group of experts jointly analyzes data from multiple source.

Would not be possible without data sharing.

H-198		 	-		 - A.
				I	
CAL	Life al				



EHzürich

Analysis of Aberrant Splicing in Cancer in 8,512 tumors





Kahles, Lehmann, et al., Rätsch, (Cancer Cell, 2018)

Gunnar Rätsch | 30.11.2018 | 23

ETHzürich

Analysis of Aberrant Splicing in Cancer in 8,512 tumors



Question: Can aberrant splicing be exploited in immunotherapies?



Kahles, Lehmann, et al., Rätsch, (Cancer Cell, 2018)

Project 3: BRCA Exchange



Sharing Global Knowledge about BRCA1/2

Current State, Challenges & Opportunities

Gunnar Rätsch

ETH Zürich (MSKCC New York)



@gxr #GA4GH #BRCAExchange

Motivation for the BRCA Challenge

BRCA variation is relatively common with well known medical implications Problem 1: Many variants lack clear interpretation

Problem 2: Variation databases are disjoint

<u>Problem 3:</u> Too little available data for effective curation





Wouldn't it be nice, if ...

HGVS Variant Lookup





Federated Database







ETHzürich

Goal: One-Stop Shop for BRCA1/2 Variant Data



APIs, data format standards, Security, ELSI mechanisms



Want everything! But only for two well-studied genes. GA4GH driver project!

Types of Data

Variant-level (data annotated to a variant)

- Genotype, classification, allele frequencies, etc.; the majority of our data.
- Well-structured, most problems from inaccurate/ambiguous variant specs.
- Easy to share.

Case-level (data annotated to a case/patient)

- Currently a small percentage of our data. Hoping to grow!
- Detailed data of cancer history, molecular features, family history, pedigrees.
- Heterogeneous clinical data.
- Privacy concerns, may need controlled-access mechanisms or other privacy enhancing mechanisms.



BRCA Exchange: Variant Exchange Platform

Highlights:

- Federated network for variant data exchange: ClinVar, LOVD, BIC, ExAC, ...
- Uniform variant processing and identification.
- Open source, cloud based, fully automatic (<u>https://github.com/BRCAChallenge/brca-exchange</u>).
- Public access, monthly releases and versioning support.
- Programmatic access via GA4GH interfaces.
- Largest public (federated) repository of BRCA1/2 variants.



ETHzürich

Each repository contributes distinct information on BRCA1/2 variation



Combined, BRCA Exchange has **21,691** individual deduplicated variants (11/2018, monthly release)

Largest BRCA1/2 public variant database worldwide.



Cline et al., PLoS Genetics, December 2018, in press.

Gunnar Rätsch | 30.11.2018 | 31

Source: BRCA Challenge working group

Data Flow and Project Aims



ETHzürich

Aim 1: Enable Finding Variant Classifications

- One place for all known BRCA1/2 variants
- Highlight expert-panel reviewed variant interpretations for clinical use
- Simple to use user interface



Variant Lookup via BRCA Exchange App

BRCA Exchange Expert Reviewed



+

following variant

Carrier ᅙ	10:24 AM		
	Home		
search for	"c.1105G>A" or "brca1"		

C

The BRCA Exchange aims to advance our understanding of the genetic basis of breast cancer, ovarian cancer and other diseases by pooling data on BRCA1/2 genetic variants and corresponding clinical data from around the world. Search for BRCA1 or BRCA2 variants above.

This project is supported by the BRCA Exchange of the Global Alliance for Genomics and Health.

> BRCA EXCHANGE

arrier 穼	11:00 AM	
\equiv	Search	
c.956		
26 variants	legend	
Gene	HGVS Nucleotide	
BRCA2	c.9561T>A	
BRCA1	c.4265G>A	
BRCA2	c.9567C>G	
BRCA2	c.729_732delTGAT	
BRCA1	c.956A>G	
BRCA2	c.9564T>A	
BRCA2	c.9560A>C	
BRCA2	c.9565C>G	
BRCA1	c.844_850dupTCATTAC	
BRCA2	c.9334G>A	
BRCA2	c.956dupA	



ETHzürich

Aim 2: Research Data & Curation Environment

- Information necessary for variant classification (allele frequencies, priors, ...)
- Data from many different, possibly disagreeing sources
- Curation tools & partially automatic variant classification
- All public data!





Aim 3: Case Level Data Exchange

- Provide infrastructure to collect & store case-level data
- Genotypes, clinical data, family history, etc.
- Analysis tools: based on family history; multi-factorial
- Controlled access mechanisms





Health Trains!?

WORK IN

BRCA Exchange -- Global Data Sharing demonstrated

Aims:

- Share variant information to clinicians/physicians
- Provide platform to facilitate research & variant curation
- Collect data from case-level data repositories help curation of VUS

We need your help!

- Help connect us to large case-level repositories. National initiatives/consortia.
- Come talk to me or write email raetsch@ethz.ch

Technical, legal, organizational challenges are similar for other diseases:

- Relatively easy to replicate BRCA Exchange for other diseases/genes
 - MMR/InSIGHT variant database
 - Lynch Syndrome
 - Other hereditary cancers

Mode of interaction: Groups of experts solve a global challenge to make clinical variant interpretation more effective.



Summary

Data sharing is key to progress in biomedicine.

- Project 1: Effective collaboration at the highest level
- Project 2: Advanced data aggregation across different technologies and cohorts
- Project 3: Leverage global expert knowledge to make medical genetics more efficient

Data sharing has to be thought about globally, but implemented locally.





Global Alliance for Genomics & Health

Acknowledgements to Collaborators

Biomedical Informatics Cristóbal Esteban <u>Andre Kahles</u> <u>Kjong Lehmann</u> Stephanie Hyland <u>Natalie Davidson</u> Gideon Dresdner Stefan Stark Xinrui Liu Matthias Hüser Vipin Sreedharan David Kuo Francesco Locatello

Alumni Julia Vogt Yi Zhong Linda Sundermann Melanie Fernandez Theofanis Karaletsos Katherine Redfield-Chan MSKCC Cancer Biology Guido Wendel Kamini Singh MSKCC Molecular Oncology Center Niki Schultz David Solit David Hyman

MSKCC IT Services Chris Crosbie Stuart Gardos Juan Perin

ETH IT Services Bernd Rinn Olivier Byrde Stefan Walter

Global Alliance for Genomics and Health

David Haussler/UCSC Benedict Paten/UCSC Melissa Cline/UCSC Stephen Chanock/NCI John Burn/University of Newcastle

International Cancer Genome Consortium Angela Brooks/UCSC Alvis Brazma/EBI Oliver Stegle/EBI

NEXUS@ETH Nora Toussaint Daniel Stekhoven ETH BSSE Dean Bodenham Karsten Borgwardt

University of Tübingen Oliver Kohlbacher

Funding: ETH Zürich, Sloan Kettering Institute, Memorial Hospital, National Institute of Health, National Cancer Institute, Swiss National Science Foundation, Max Planck Society, German Research Foundation, European Union, Geoffrey Beene Foundation, Lucille Castori Center

Thank You!

Courtesy of Mark Rubin

Questions?