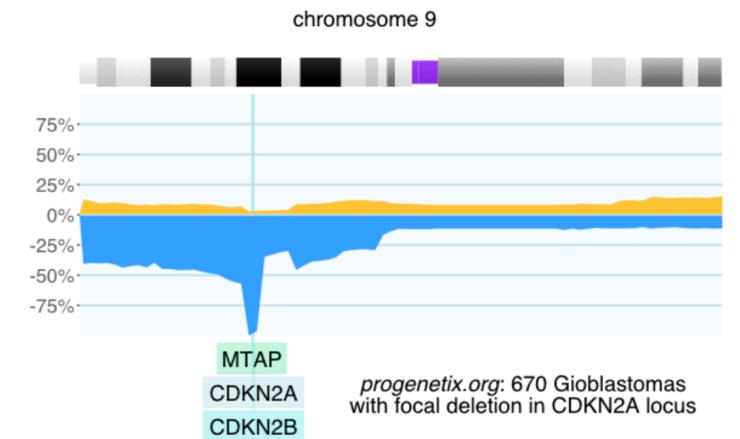


ELIXIR hCNV

First Implementation Study and Ongoing Work

Michael Baudis | ELIXIR Human Data Communities | 2022-03-15



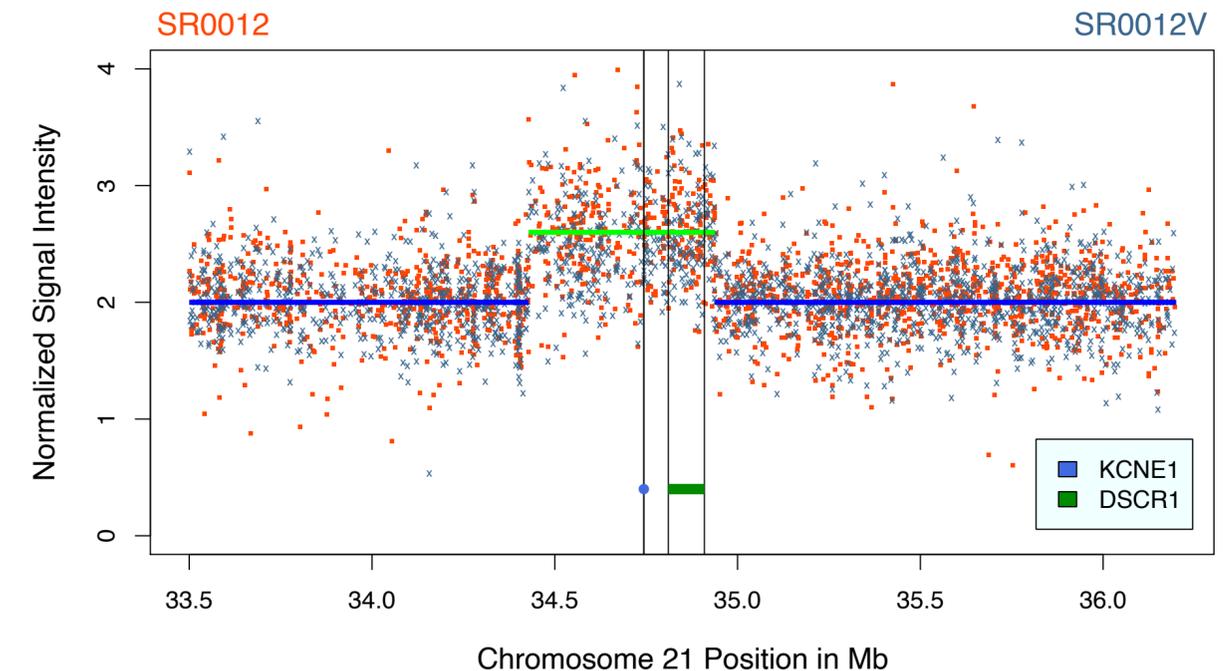
Why hCNV Community?

Structural Genome Variation Data :: Resources and Technologies

- structural genome variations are a major contributor to genetic diseases and cancer
- knowledge about and standards for copy number variations / aberrations (CNV/CNA) has not been in step with NGS & GWAS driven SNV/SNP assessment

Mission statement

Despite the fact that **Copy Number Variations** are the **most prevalent genetic mutation type**, identifying and interpreting them is still a major challenge. The ELIXIR human Copy Number Variation (hCNV) Community aims to implement processes to make the **detection**, **annotation** and **interpretation** of these variations easier



CNV with unknown clinical impact in a case of Silver-Russel Syndrome

Local Affymetrix Genotyping 6 signal distribution pattern and segmentation result in patient SR12 (SR0012, orange data) and his father (SR0012V, steelblue data). In both samples a duplication in the DSCR can be observed, affecting the whole KCNE1 and DSCR1/RCAN coding regions. In contrast, DYRK1A lays ~2.5 Mb distal of the duplication. Only the genes discussed in this article are shown.

RESEARCH ARTICLE

AMERICAN JOURNAL OF
medical genetics PART A

Identification of a 21q22 Duplication in a Silver–Russell Syndrome Patient Further Narrows Down the Down Syndrome Critical Region

Thomas Eggermann,^{1*} Nadine Schönherr,¹ Sabrina Spengler,¹ Susanne Jäger,¹ Bernd Denecke,² Gerhard Binder,³ and Michael Baudis⁴

¹Institute of Human Genetics, RWTH Aachen, Aachen, Germany

²Interdisciplinary Centre for Clinical Research, IZKF "BIOMAT," RWTH Aachen, Aachen, Germany

³Section of Paediatric Endocrinology and Diabetology, University Children's Hospital, Tuebingen, Germany

⁴Institute of Molecular Biology, University of Zürich, Zürich, Switzerland

Received 26 June 2009; Accepted 6 November 2009

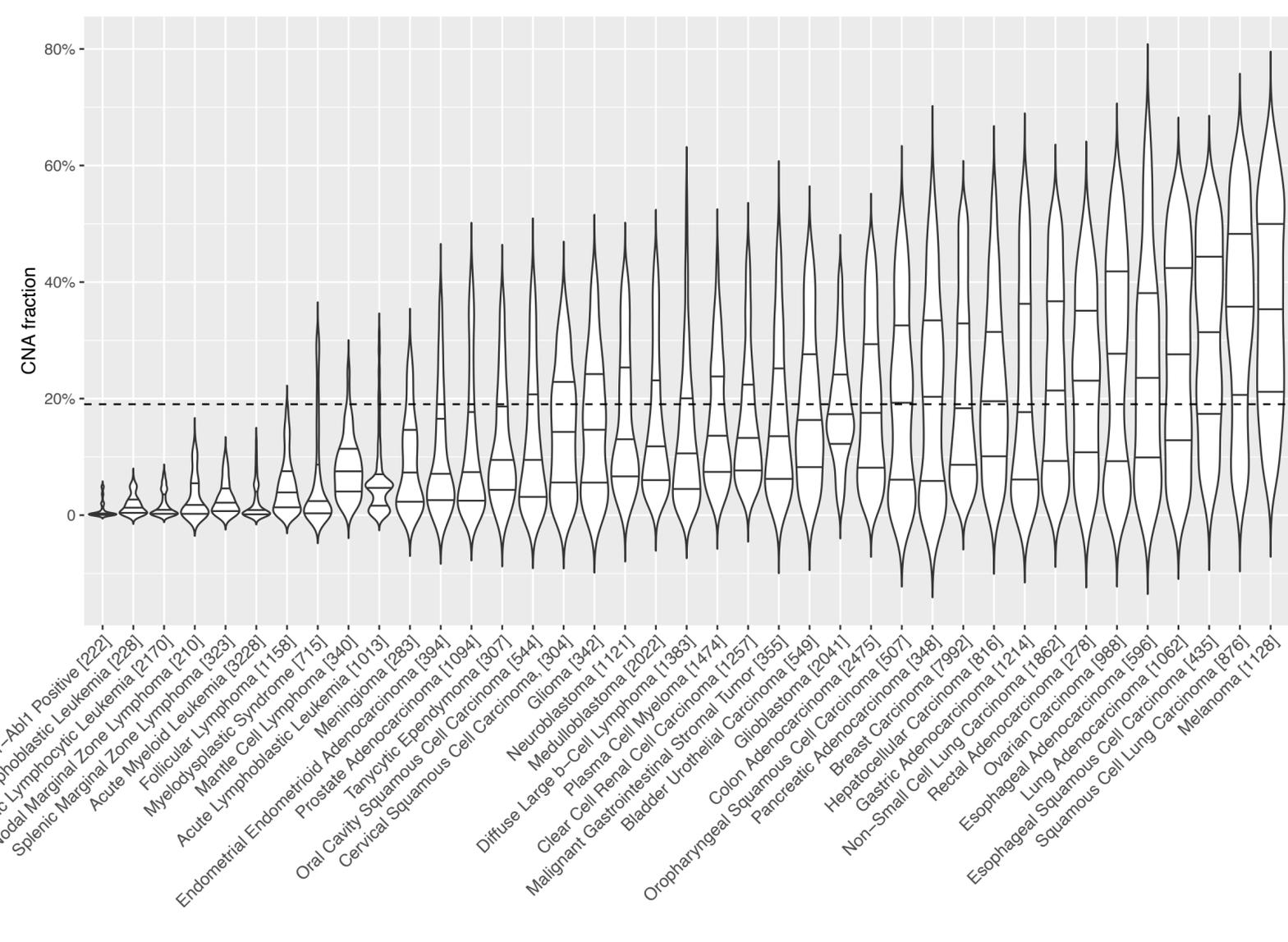
Why hCNV Community?

Structural Genome Variation Data :: Resources and Technologies

- structural genome variations are a major contributor to genetic diseases and cancer
- knowledge about and standards for copy number variations / aberrations (CNV/CNA) has not been in step with NGS & GWAS driven SNV/SNP assessment

Mission statement

Despite the fact that **Copy Number Variations** are the **most prevalent genetic mutation type**, identifying and interpreting them is still a major challenge. The ELIXIR human Copy Number Variation (hCNV) Community aims to implement processes to make the **detection**, **annotation** and **interpretation** of these variations easier



Genome CNV coverage in Cancer Classes

- 43654 out of 93640 CNV profiles; filtered for entities w/ >200 samples
- Single-sample CNV profiles were assessed for the fraction of the genome showing CNVs (relative gains, losses)
- range of medians 0.001 (CML) - 0.358 (malignant melanomas) knowledge about and standards for copy number variations / aberrations (CNV/CNA) has not been in step with NGS & GWAS driven SNV/SNP assessment



ELIXIR hCNV Community

Structural Genome Variation Data :: Resources and Technologies

- First meeting of group in 2018
- ELIXIR Human Copy Number Variation (hCNV) approved in 2019
- initial implementation study (2019-2021) for community set-up, gap analysis and exploration of technical deliverable

Mission statement

Despite the fact that **Copy Number Variations** are the **most prevalent genetic mutation type**, identifying and interpreting them is still a major challenge. The ELIXIR human Copy Number Variation (hCNV) Community aims to implement processes to make the **detection**, **annotation** and **interpretation** of these variations easier

Purpose

The human CNV community (h-CNV) has been officially created in December 2018. It aims to address the major challenge of NGS data interpretation in the era of whole genome sequencing for the most frequent mutation type: Copy Number Variation. Seven topics have been identified during the kick-off meeting and further refined with all h-CNV partners. This ultimately led to the proposal described in this implementation study.

Node	Name of PI
ELIXIR-FR	Christophe Bérout, David Salgado, Marc Hanauer, Victoria Dominguez
ELIXIR-CH	Michael Baudis
ELIXIR-DE	Jan Korbel
EMBL-EBI	Thomas Keane, Fiona Cunningham
ELIXIR-ES	Joaquin Dopazo, Alfonso Valencia, Salvador Capella, Sergi Beltran, Steven Laurie, Gemma Bullich, Laura I. Furlong, Janet Piñero
ELIXIR Hub	John Hancock, Gary Saunders, Kathi Lauer, Leyla Garcia
ELIXIR-NL	Bauke Ylstra, Daoud Sie, Leon Mei, Morris Swertz (UMCG), Lennart Johansson
ELIXIR-NO	Eivind Hovig, Pubudu Samarakoon
ELIXIR-HU	Attila Gyenesi, Katalin Monostory
ELIXIR-SI	Brane Leskošek, Polonca Ferk, Marko Vidak
ELIXIR-UK	Krzysztof Poterlowicz
Delivery	Starting from June 2019 for a period of 24 months.



Christophe Bérout
(ELIXIR France)



David Salgado
(ELIXIR France)



Gary Saunders
(Human Data Coordinator,
ELIXIR Hub)



Michael Baudis
(ELIXIR Switzerland)



hCNV Implementation Study 2019-2021

Setting the Scope | Solidifying the Community | First Deliveries

- challenge participants and define the wider landscape as well as future directions
- set of 7 work packages
 - ➔ landscape analysis
 - ➔ technical products
 - ➔ resource improvement
 - ➔ community building & outreach
- regular meetings, website, hackathons...
- ▶ WP1 - Optimal CNV detection pipelines for research and diagnostics
- ▶ WP2 - Definition of reference datasets
- ▶ WP3 - Improvement of community formats for CNV exchange
- ▶ WP4 - Enabling CNV data discovery in diagnostic and phenotypic context
- ▶ WP5 - Creation of innovative tools
- ▶ WP6 - FAIRification of h-CNV databases and datasets
- ▶ WP7 - Dissemination



hCNV Implementation Study 2019-2021

Setting the Scope | Solidifying the Community | First Deliveries

- highly ambitious goals, beyond available support
 - ➔ especially reference / benchmarking dataset generation and pipeline development
- emerging interactions and collaborations with ELIXIR platforms & communities and beyond
 - ➔ Galaxy
 - ➔ GA4GH / ELIXIR Beacon project
- ▶ WP1 - Optimal CNV detection pipelines for research and diagnostics
- ▶ WP2 - Definition of reference datasets
- ▶ WP3 - Improvement of community formats for CNV exchange
- ▶ WP4 - Enabling CNV data discovery in diagnostic and phenotypic context
- ▶ WP5 - Creation of innovative tools
- ▶ WP6 - FAIRification of h-CNV databases and datasets
- ▶ WP7 - Dissemination



hCNV Implementation Study 2019-2021

Some Achievements and Deliveries

- Benchmarking tools and OpenEbench TransBioNet testing event
- demonstration of CNV detection tools in clinical (cancer) setting
- amending *bio.tools* for extensive list of CNV related analysis tools
 - ➔ <https://bio.tools/t?domain=elixir-hcnv>
- updating / registering shared hCNV resources at fairsharing.org
- consensus collection of perceived requirements for efficient and effective CNV file and data exchange formats



With the ELIXIR Tools Platform: Bio.tools & EDAM ontology



Within the first commissioned service granted to the community (First hCNV Community IS)

The community created a list of 245 CNV detection tools for various detection technologies NGS (WGS, WES, panel), CGHarray, ...

- We wanted to share this list of tools → Bio.tools
- Started to collaborate with ELIXIR tools platform members (Jon Ison / Hervé Menager)
- We created a specific bio.tools subdomain and listed/annotated (with EDAM terms) 109/245 CNV tools <http://elixir-hcnv.bio.tools/>
- We contributed to the EDAM ontology to include about 20 specific terms to describe CNV and Structural variations in (topics/operation branches)

Download	Recently Modified	Software name	Tests dataset	Last Version / Last Update ??	OS	Language	Control set required - note in function	Input format	Output format	EMBL/NCBI ID	Ref	Publication	PMID	Experiment type - Add collection	Executable	Download link - link when repo download page otherwise	Running time	archiver	Sequencer
100		Radial		1.0.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
101		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
102		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
103		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
104		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
105		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
106		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
107		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
108		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
109		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
110		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
111		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
112		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
113		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
114		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
115		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
116		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
117		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
118		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
119		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
120		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
121		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
122		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
123		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
124		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
125		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
126		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
127		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
128		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
129		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
130		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
131		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
132		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
133		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
134		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
135		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
136		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
137		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
138		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
139		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
140		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
141		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
142		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
143		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
144		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
145		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
146		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
147		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
148		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
149		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
150		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
151		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
152		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
153		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
154		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
155		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
156		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
157		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
158		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
159		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
160		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
161		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
162		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
163		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
164		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
165		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
166		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			
167		Radial		0.1.0							Radial et al.	2007				http://www.ncbi.nlm.nih.gov/pubmed/17488888			

hCNV Implementation Study 2019-2021

Some Achievements and Deliveries

- HGVS satellite meeting – Human CNV – June 14th 2019 – Göteborg Sweden
- hCNV community workshop ELIXIR All-Hands Lisbon – June 2019
- survey of data annotation formats, including comments on VCF development
- start FAIRification of CNV national / reference databases (BANCCO, Progenetix)
- Community white paper published
- Biohackathon Paris 2019
- in 2021 start of shared meetings of subgroup with Beacon variants scout team

F1000Research

F1000Research 2020, 9(ELIXIR):1229 Last updated: 01 JUN 2021



OPINION ARTICLE

The ELIXIR Human Copy Number Variations Community: building bioinformatics infrastructure for research [version 1; peer review: 1 approved]

David Salgado ¹, Irina M. Armean², Michael Baudis ³, Sergi Beltran^{4,5}, Salvador Capella-Gutierrez ^{6,7}, Denise Carvalho-Silva ^{2,8}, Victoria Dominguez Del Angel ⁹, Joaquin Dopazo ¹⁰, Laura I. Furlong ¹¹, Bo Gao ³, Leyla Garcia ^{2,12,13}, Dietlind Gerloff¹⁴, Ivo Gut^{4,5}, Attila Gyenesei¹⁵, Nina Habermann¹⁶, John M. Hancock ¹³, Marc Hanauer¹⁷, Eivind Hovig ^{18,19}, Lennart F. Johansson²⁰, Thomas Keane², Jan Korbel¹⁶, Katharina B. Lauer ¹³, Steve Laurie⁴, Brane Leskošek²¹, David Lloyd ¹³, Tomas Marques-Bonet²², Hailiang Mei²³, Katalin Monostory²⁴, Janet Piñero ¹¹, Krzysztof Poterlowicz ²⁵, Ana Rath¹⁷, Pubudu Samarakoon²⁶, Ferran Sanz¹¹, Gary Saunders ¹³, Daoud Sie²⁷, Morris A. Swertz²⁰, Kirill Tsukanov ², Alfonso Valencia^{6,7,28}, Marko Vidak²¹, Cristina Yenyxe González², Bauke Ylstra²⁹, Christophe Bérout^{1,30}

¹Aix Marseille Univ, INSERM, MMG, Marseille, France

²European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, UK

³Department of Molecular Life Sciences and Swiss Institute of Bioinformatics, University of Zurich, Zurich, Switzerland

⁴CNAG-CRG, Centre for Genomic Regulation (CRG), Barcelona Institute of Science and Technology (BIST), Baldri Reixac 4, Barcelona 08028, Spain

⁵Universitat Pompeu Fabra (UPF), Barcelona, Spain

⁶Barcelona Supercomputing Center (BSC), Barcelona, Spain

⁷Spanish National Bioinformatics Institute (INB)/ELIXIR-ES, Barcelona, Spain

⁸Open Targets, Wellcome Genome Campus, Hinxton, Cambridgeshire, CB10 1SD, UK

⁹Institut Français de Bioinformatique, UMS3601-CNRS, CNRS, Paris, France

¹⁰Clinical Bioinformatics Area, Fundación Progreso y Salud, CDCA, Hospital Virgen del Rocío, Sevilla, Spain

¹¹Research Programme on Biomedical Informatics (GRIB), Hospital del Mar Medical Research Institute (IMIM), Department of Experimental and Health Sciences, Pompeu Fabra University (UPF), Barcelona, Spain

¹²ZB MED Information Centre for Life Sciences, Cologne, Germany

¹³ELIXIR Hub, Hinxton, UK

¹⁴Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Belvaux, Luxembourg

¹⁵Szentágothai Research Center, University of Pécs, Pécs, Hungary

¹⁶Genome Biology, European Molecular Biological Laboratory, Heidelberg, Germany

¹⁷Orphanet, INSERM, Paris, France

¹⁸Department of Tumor Biology, Institute for Cancer Research, Oslo University Hospital, Oslo, Norway

¹⁹Centre for bioinformatics, Department of Informatics, University of Oslo, Oslo, Norway

²⁰Department of Genetics, University of Groningen, University Medical Center Groningen, Groningen, The Netherlands

²¹Faculty of Medicine - ELIXIR Slovenia, University of Ljubljana, Ljubljana, Slovenia

²²Institute of Evolutionary Biology (UPF-CSIC), Catalan Institution for Research and Advanced Studies, Barcelona, Spain

²³Sequencing Analysis Support Core, Leiden University Medical Center, Leiden, The Netherlands

²⁴Institute of Enzymology, Research Centre for Natural Sciences, Budapest, Hungary

²⁵Centre for Skin Sciences, University of Bradford, Bradford, UK



hCNV Implementation Study 2019-2021

Some Achievements and Deliveries

- survey about genomic variation file formats and their use, suitability for representing CNV data
- part of the survey was focused specifically on VCF, a key GA4GH standard at the intersection of human and computer readable formats
- Results
 - ➔ BED-like formats are frequently used, but the better defined flavours are not optimal for CNVs and other SVs
 - ➔ JSON w/ schema has potential, but still misses finalized GA4GH schemas (VRS emerging) and suffers "readability" issues for non-bioinformatics customers
 - ➔ VCF was considered as a/the variant standard file format, but not "CNV-friendly" in v4.2 and in the existing tools for I/O handling of CNV data

ELIXIR hCNV 2019-21 Deliverable D3.2

Project Title:	First hCNV Community Implementation Study
Deliverable title:	Create a consensus collection of perceived requirements for efficient and effective CNV file and data exchange formats
WP No.	3
WP Title	Improvement of community formats for CNV exchange
Contractual delivery date:	30.11.2019
Actual delivery date:	12.12.2019
WP leads:	Thomas Keane
Partner(s) contributing to this deliverable:	EMBL-EBI

Report authors: Kirill Tsukanov¹, Sundararaman Venkataraman, Giselle Kerry, Thomas Keane (EMBL-EBI)

2. Results	3
2.1. Feedback overview	3
2.2. Terminology	4
2.3. Existing file formats	4
2.3.1. VCF (Variant Call Format)	4
2.3.2. BED and related tab-separated formats	5
2.3.3. JSON with a schema	5
2.3.4. Other formats	5
2.4. Opinion on CNV representation in VCF	6
2.5. Requirements for CNV formats of the future	6
2.6. Conclusions. Note about use cases	7
3. Impact	8

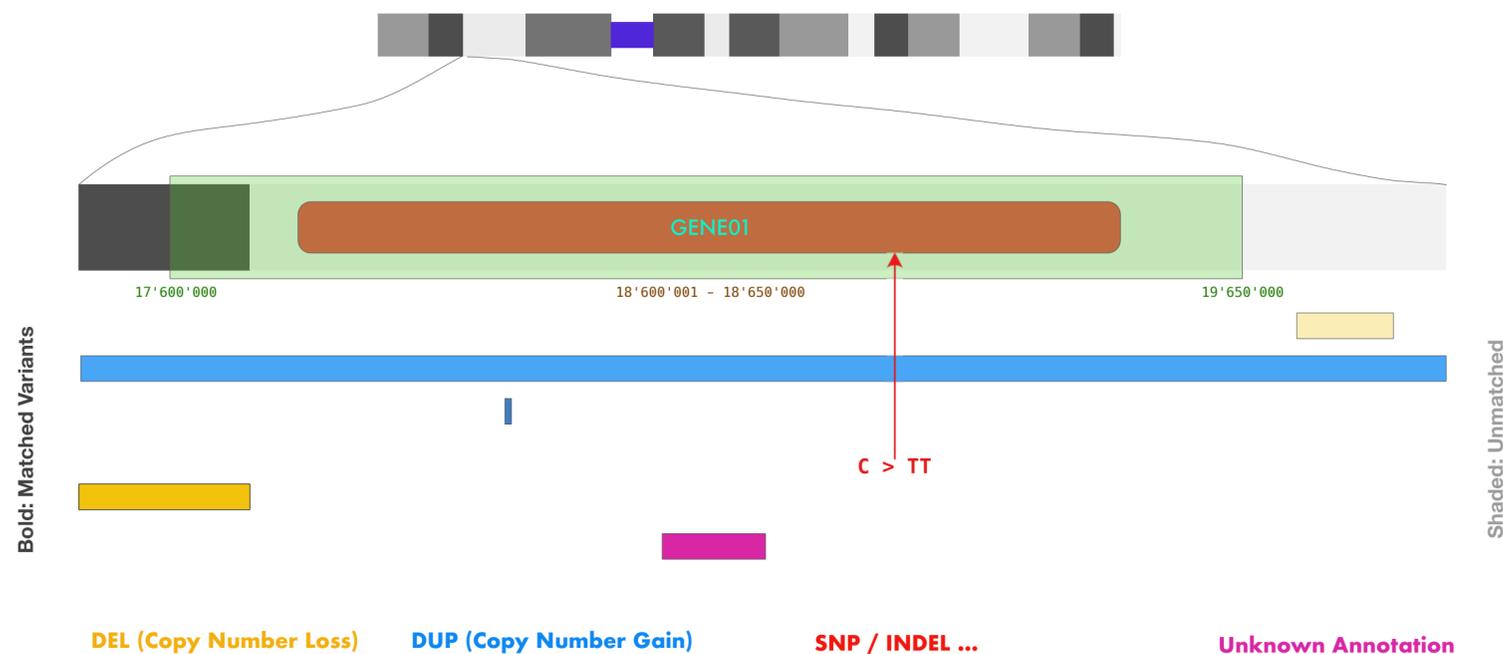


Beacon v2: Extended Variant Queries



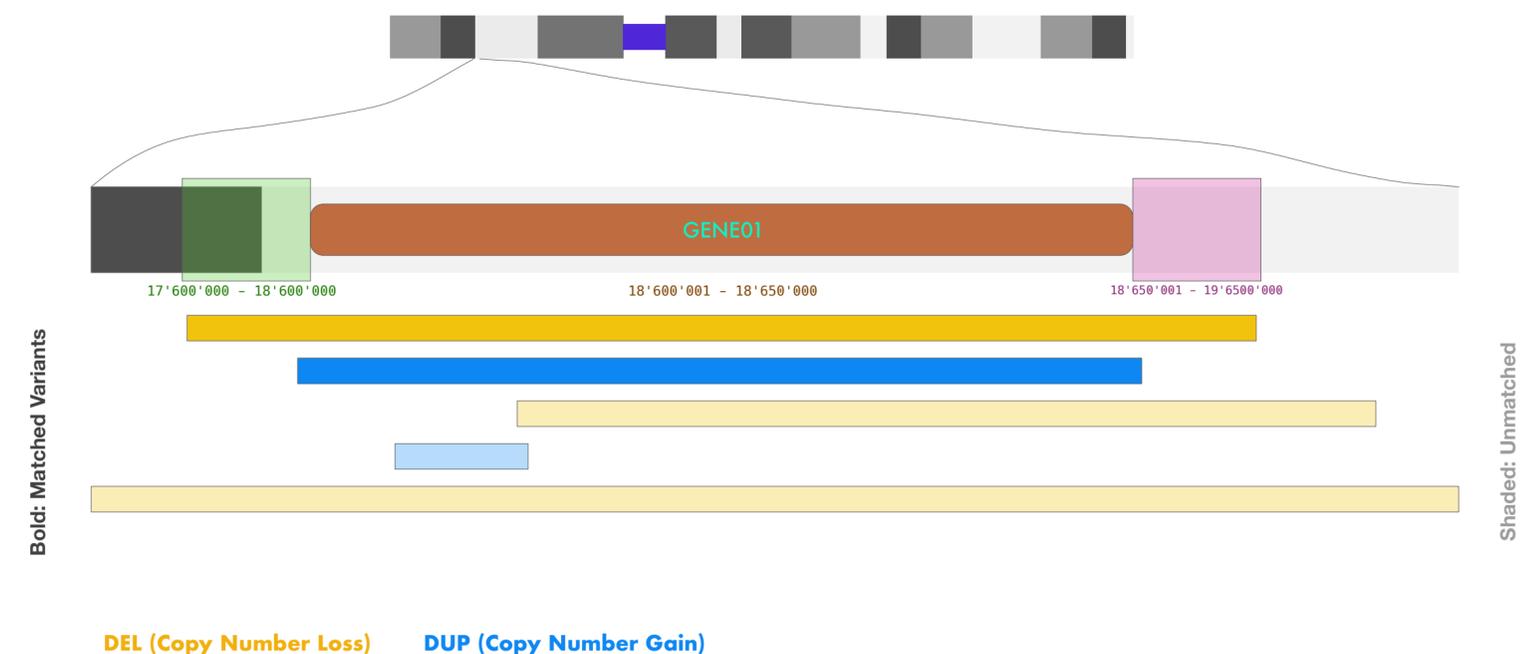
Range and Bracket queries enable positional wildcards and fuzziness

Genome Range Query (matching variants in a region)



- Genome Range Queries provide a way to "fish" for variants overlapping an indicated region, e.g. the CDR of a gene of interest
- Additional parameters (e.g. variant type, reference or alternate bases) limit the scope of the responses
- new Beacon v2 size parameters to limit structural variants (e.g. "focal" CNVs)

Genome Bracket Query (full match)



- Genome Bracket Queries allow to search for structural variants with start and end positions falling into defined sequence ranges
- allows to query any contiguous genomic variant (and in principle also can step in for range queries)
- typical use case is e.g the query for variants such as duplications covering the whole CDR of a gene, while limiting the allowed start or end regions

Links with other projects



Participation to various taskforces

- Variant representation
- Beacon
- Future of VCF
- Adopting new standards (phenopackets, DUO)



Communities & services

Galaxy
Beacon
Human data communities

Platforms

Tools

Data *

Interoperability *

Training *

* Links to be strengthened through future IS

National links reinforcements

e.g.

- BANCCO national CNV db for diagnostic
- Links to national sequencing projects (PFMG, GOLD)
- ...



LINKS to be reinforced

hCNV community



Access to CNV resources (Beacon)
Demonstrator for HOOM ontology



CNV workshop organisation, Göteborg Sweden 2019



2x hCNV Implementation Studies 2021-2023

Reference hCNV datasets, use-case workflows and benchmarking

The ELIXIR human CNV Community (hCNV) was created in December 2018. In two years contributions to the field have been numerous (ELIXIR IS, Rare Diseases, Federated Human Data, Beacons, GA4GH, EJP-RD and Beyond 1 Million Genomes - B1MG). The Community now aims to address the major challenge of NGS data interpretation in the era of whole genome sequencing: Copy Number Variation. During the first commissioned service offered as a starting grant, the Community has identified various gaps to proceed with CNV tools benchmarking and in particular for Exome and targeted sequencing, which are by far the most widely used technologies in diagnostic laboratories and in research. Within this implementation study we want to provide solutions and bioinformatic infrastructure solutions to fill identified gaps, and to make these biomedical reference materials available (i.e. via Open Science) to the various communities and platforms.

Interactions and utility to other projects

ELIXIR platforms:

Data, Tools, Interoperability, Training

ELIXIR Communities:

hCNV, Galaxy, Rare diseases, Federated Human Data

National and International projects:

EJP-RD, B1MG, EOSC-Life, EOSC-Pillar

Beacon and beyond – Implementation-driven standards and protocols for CNV discovery and data exchange

The initial 2019-2021 hCNV community implementation study employed a set of perceived needs to a) deliver first community standards and procedures; b) identify intersections with other ELIXIR communities and stakeholders in ELIXIR connected organizations, such as GA4GH; and c) to streamline priorities for relevant, achievable deliveries of hCNV community projects.

This proposal for an hCNV implementation study focuses on those potential high-value targets for data access and delivery, using reference resources and community stakeholder engagement to directly implement and test hCNV resources aligned with ELIXIR ecosystems.

The main target here will be the empowerment of the Beacon protocol, to act as standard for federated hCNV discovery and data delivery, in conjunction with additional GA4GH derived standards.

Intersecting ELIXIR Platforms, Communities and Projects:

- ELIXIR Galaxy Community
- ELIXIR AAI Infrastructure Service
- ELIXIR Compute Platform
- ELIXIR Training Platform
- ELIXIR FHD Community
- ELIXIR Health Data Focus Group
- ELIXIR Beacon Strategic Implementation Study
- ELIXIR Interoperability Platform

External Projects and Partners:

- EJP-RD
- GA4GH (Discovery, Genomic Knowledge Standards, Phenopackets)



hCNV Implementation Studies 2021-2023 No. 1

Reference hCNV datasets, use-case workflows and benchmarking

- only limited datasets exist to test and benchmark tools for the analysis of CNV and structural variations
 - recent datasets focused on high-quality Whole Genome Sequencing (WGS) analyses but not on the most commonly used Whole Exome Sequencing (WES) or genomic array technologies
 - generation of publicly accessible reference sets (raw and interpreted CNV data) for a variety of technological platforms will allow the hCNV community to generate the mandatory material
 - creation of “control datasets” required by many detection tools
 - complement standardization and benchmarking efforts such as the “Genome in a Bottle” initiative
 - integrate with Galaxy community & platforms
- ▶ WP1 - Dataset selection and generation
 - ▶ WP2 - Analyse and Compare CNV with other Benchmarking initiatives
 - ▶ WP3 - Exploitation of the datasets by the Galaxy Community
 - ▶ WP4 - Training and dissemination



hCNV Implementation Studies 2021-2023 No. 2



Beacon and beyond – Implementation-driven standards and protocols for CNV discovery and data exchange

- reinforce work on priority areas established in the current hCNV Implementation Study
 - extend collaborations with the Rare Diseases and Galaxy Communities, EJP-RD and GA4GH
 - Expected outcomes
 - ➔ shared CNV resources testing advanced versions of the Beacon protocol
 - ➔ integration of GA4GH standards such as Phenopackets in such resources
 - ➔ tools for data ingestion and export for standard formats (e.g. VCF, Phenopackets) and CNV-specific improvements of such standards
 - ➔ ELIXIR AAI demo on clinical and research hCNV resources
 - ➔ demonstration of Galaxy pipeline adoption for real-world hCNV data analysis projects
 - connecting to international partners, e.g. Cancer Genomics Consortium (U.S.)
- ▶ WP1 - hCNV community reference resources
 - ▶ WP2 - hCNV Resources and Beacon
 - ▶ WP3 - Galaxy Community Intersection and Data Exchange
 - ▶ WP4 - Workflows and Tools for hCNV Data Exchange Procedures
 - ▶ WP5 - Training and dissemination





Ongoing... hCNV & Intl. Community

- contributions to ontologies and standard definitions
- close ongoing interactions with GA4GH work streams
- influencing the development of the GA4GH VRS variant standard

hCNV Community

Genomic Copy Number Variations in Humans

News & Events

ELIXIR All Hands 2022 - h-CNV Representation
 CNV Ontology Proposal - Now Live at EFO
 hCNV Site now at cnvar.org
 hCNV Implementation Study 2021/2: Beacon and Beyond
 all ...

Participants

Standards and Guidelines

Studies & Resources

Examples, Guides & FAQ

Contacts

Related Sites

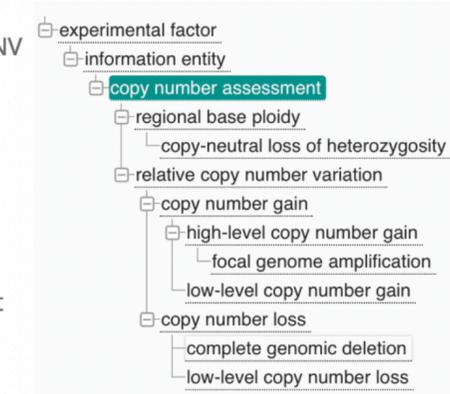
[h-CNV @ ELIXIR](#)
[Beacon Project](#)
[Beacon @ ELIXIR](#)
[SchemaBlocks](#)

Github Projects

[h-CNV](#)

CNV Ontology Proposal - Now Live at EFO

As part of the hCNV-X work - related to "Workflows and Tools for hCNV Data Exchange Procedures" and to the intersection with Beacon and GA4GH VRS - we have now a new proposal for the creation of an ontology for the annotation of (relative) CNV events. The CNV representation ontology is targeted for adoption by Sequence Ontology (SO) and then to be used by an updated version of the VRS standard. Please see the discussions linked from the [proposal page](#). However, we have also contributed the CNV proposal to EFO where it has gotten live on January 21.



Everybody is welcome to contribute to the editing of the proposal at the SO & VRS Github repositories!

2021-01-21: copy number assessment term tree now live on EFO

The [copy number assessment](#) term tree has been accepted into the Experimental Factor Ontology and can be used for referencing CNV types.

More ontologies...

... with h-CNV contributions ca

2022-01-21



larrybabb commented 18 days ago

per a discussion between [@ahwagner](#) and [@larrybabb](#)
 Dreaft Relative Copy Number class proposal

```
-- the target region/gene/feature
subject: region/gene/feature/allele/haplotype

--5 quantifiable values that correspond to the EFO copy number assessi
copy number assessment: (http://www.ebi.ac.uk/efo/EFO_0030063)
  -2 = complete loss (http://www.ebi.ac.uk/efo/EFO_0030069)
  -1 = partial loss (http://www.ebi.ac.uk/efo/EFO_0030068)
  0 = copy-neutral (http://www.ebi.ac.uk/efo/EFO_0030064)
  1 = low-level gain (http://www.ebi.ac.uk/efo/EFO_0030071)
  2 = high-level gain (http://www.ebi.ac.uk/efo/EFO_0030072)
```

RelativeCopyNumber

Relative Copy Number Variation captures a classification of copies of a molecule within a system, relative to a baseline. These types of Variation are common outputs from CNV callers, particularly in the somatic domain where Absolute Copy Counts are difficult to estimate and less useful in practice than relative statements.

Computational Definition

The relative copies of a [Molecular Variation](#), [Feature](#), [Sequence Expression](#), or a [CURIE](#) reference against an unspecified baseline in a system (e.g. genome, cell, etc.).

Information Model

Some RelativeCopyNumber attributes are inherited from [Variation](#).

Field	Type	Limits	Description
_id	CURIE	0..1	Variation Id. MUST be unique within document.
type	string	1..1	MUST be "RelativeCopyNumber"
subject	Molecular Variation Feature Sequence Expression CURIE	1..1	Subject of the Copy Number object
relative_copy_class	string	1..1	MUST be one of "complete loss", "partial loss", "copy neutral", "low-level gain" or "high-level gain".



hCNV Implementation Studies 2021-2023

Focus on Integration with ELIXIR Platforms and Communities - and beyond

- original 2019-2021 implementation study provided visibility and established connections for new studies
- instrumental were Biohackathons, use case & standards surveys and co-participation of group members
- future work plans to leverage the resources of participants through pre-established interactions and synergies
- 2 independent studies provide clearer definitions of deliverables and individual scopes

Michael Baudis	CH
Christophe Béroud	FR
David Salgado	FR
Alexander Kanitz	CH
Anthony Brookes	UK
Babita Singh	ES
Björn Grüning	DE
Jordi Rambla	ES
Kirill Tsukanov	EMBL-EBI
Krzysztof Poterlowicz	UK
Salvador Capella-Gutierrez	ES
Sergi Beltran	ES
Steven Laurie	ES
Tim Beck	UK
Timothee Cezard	EMBL-EBI



{BEH}

BIOHACKATHON EUROPE

7 - 11 November 2022

CALL FOR PROPOSALS

[@ELIXREurope](#)

[#BioHackEU22](#)



