



# CNV Annotations: a crucial step in your variant analysis

Darby Kammeraad  
Field Application Scientist Manager



20 Most Promising Biotech  
Technology Providers



Hype Cycle for Life sciences



Top 10 Analytics  
Solution Providers

# NIH Grant Funding Acknowledgments

- Research reported in this publication was supported by the National Institute Of General Medical Sciences of the National Institutes of Health under:
  - Award Number R43GM128485-01
  - Award Number R43GM128485-02
  - Award Number 2R44 GM125432-01
  - Award Number 2R44 GM125432-02
- Montana SMIR/STTR Matching Funds Program Grant Agreement Number 19-51-RCSBIR-005
- PI is Dr. Andreas Scherer, CEO Golden Helix.
- The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

# Who Are We?

Golden Helix is a global bioinformatics company founded in 1998



Filtering and Annotation

ACMG Guidelines

Clinical Reports

CNV Analysis

Pipeline: Run Workflows



Variant Warehouse

Centralized Annotations

Hosted Reports

Sharing and Integration



CNV Analysis

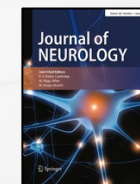
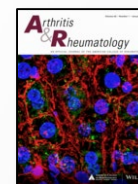
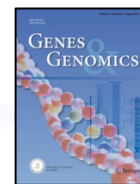
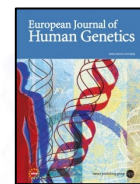
GWAS | Genomic Prediction

Large-N Population Studies

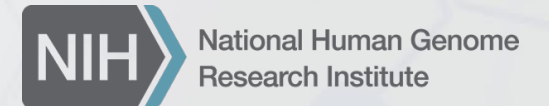
RNA-Seq

Large-N CNV-Analysis

# Cited in 1,000s of Peer-Reviewed Publications



# Over 400 Customers Globally



# When you choose Golden Helix, you receive more than just the software



## SOFTWARE IS VETTED

- 20,000+ users at 400+ organizations
- Quality & feedback



## DEEPLY ENGRAINED IN SCIENTIFIC COMMUNITY

- Give back to the community
- Contribute content and support



## SIMPLE, SUBSCRIPTION-BASED BUSINESS MODEL

- Yearly fee
- Unlimited training & support



## INNOVATIVE SOFTWARE SOLUTIONS









- Cited in 1,000s of publications

GENE PANEL

EXOME

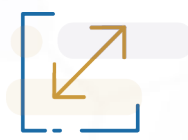
GENOME

SEQUENCER

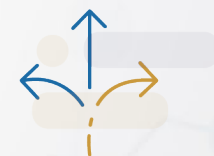
PRODUCTS	BIOINFORMATICS PIPELINE	FUNCTION
 DNaseq (Sentieon)  TNseq (Sentieon)  VS-CNV	FASTQ BAM VCF	<ul style="list-style-type: none"> <li>▶ Single nucleotide variation</li> <li>▶ Copy number variation &amp; loss of heterozygosity</li> <li>▶ Chromosomal aberration</li> </ul>
Annotations	Annotated VCF	<ul style="list-style-type: none"> <li>▶ Public &amp; commercial annotations to enrich genomic data sets</li> </ul>
 VarSeq  VSReports  VSPipeline	Clinical Report	<ul style="list-style-type: none"> <li>▶ Annotate &amp; filter</li> <li>▶ Visually inspect alignments</li> <li>▶ Variant prioritization</li> <li>▶ Clinical assessment</li> </ul>
 VSclinical	Automated ACMG Guidelines	<ul style="list-style-type: none"> <li>▶ Clinical variant interpretation in coordination with ACMG Guidelines &amp; AMP Guidelines</li> </ul>
 VSWarehouse	Data Warehousing  Web-Enabled Interface + Powerful APQ: JSON, XML, TSV, CSV, SQL, FHIR	<ul style="list-style-type: none"> <li>▶ Clinical assessment catalog</li> <li>▶ Advanced data querying</li> <li>▶ Versioning</li> <li>▶ Interoperability</li> <li>▶ Compliance with HIPPA, CLIA &amp; CAP data discovery</li> </ul>



Simple



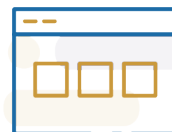
Flexible



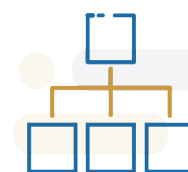
Scalable



Variant annotation  
filtering, and interpretation



Powerful GUI with  
rich visualizations

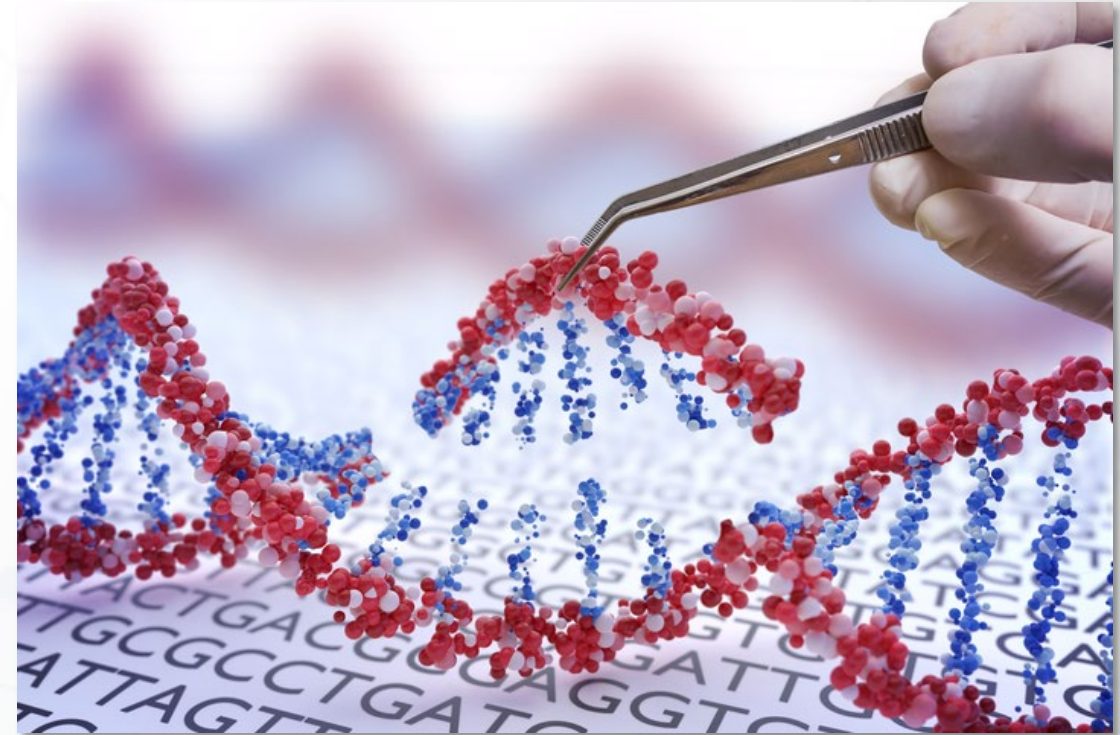


Repeatable  
workflows



# CNVs in Clinical Testing

- Critical evidence needed for many genetic tests
- Common driver specific cancers, causal hereditary variation
  - EGFR Exon 19 deletion common in lung cancer
  - PIK3CA Amplification in breast cancer
- Large events used heavily in diagnostics
  - Chromosome 13 deletion common in melanoma
  - Autism Spectrum Disorder (ASD)
  - Developmental Delay (DD)
  - Intellectual Delay (ID)



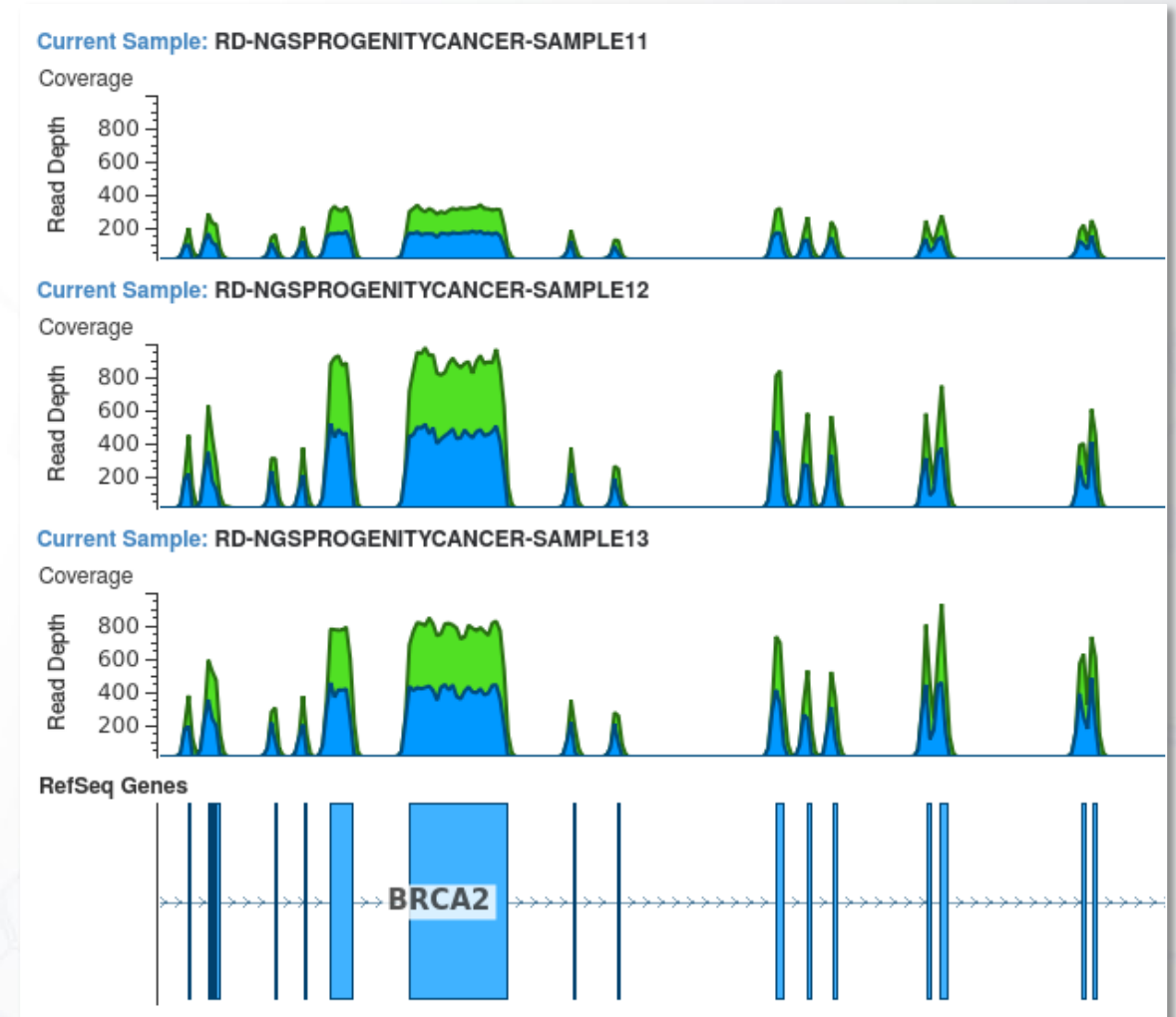
# Power of NGS CNV Detection

	Detectable events			Supported Data types		
	Small: 150b+	Medium: 1 – 10Kb	Large: 10Kb+	Gene panel	Whole exome	Whole genome
MLPA	✓			✓		
CMA			✓			✓
VS-CNV	✓	✓	✓	✓	✓	✓

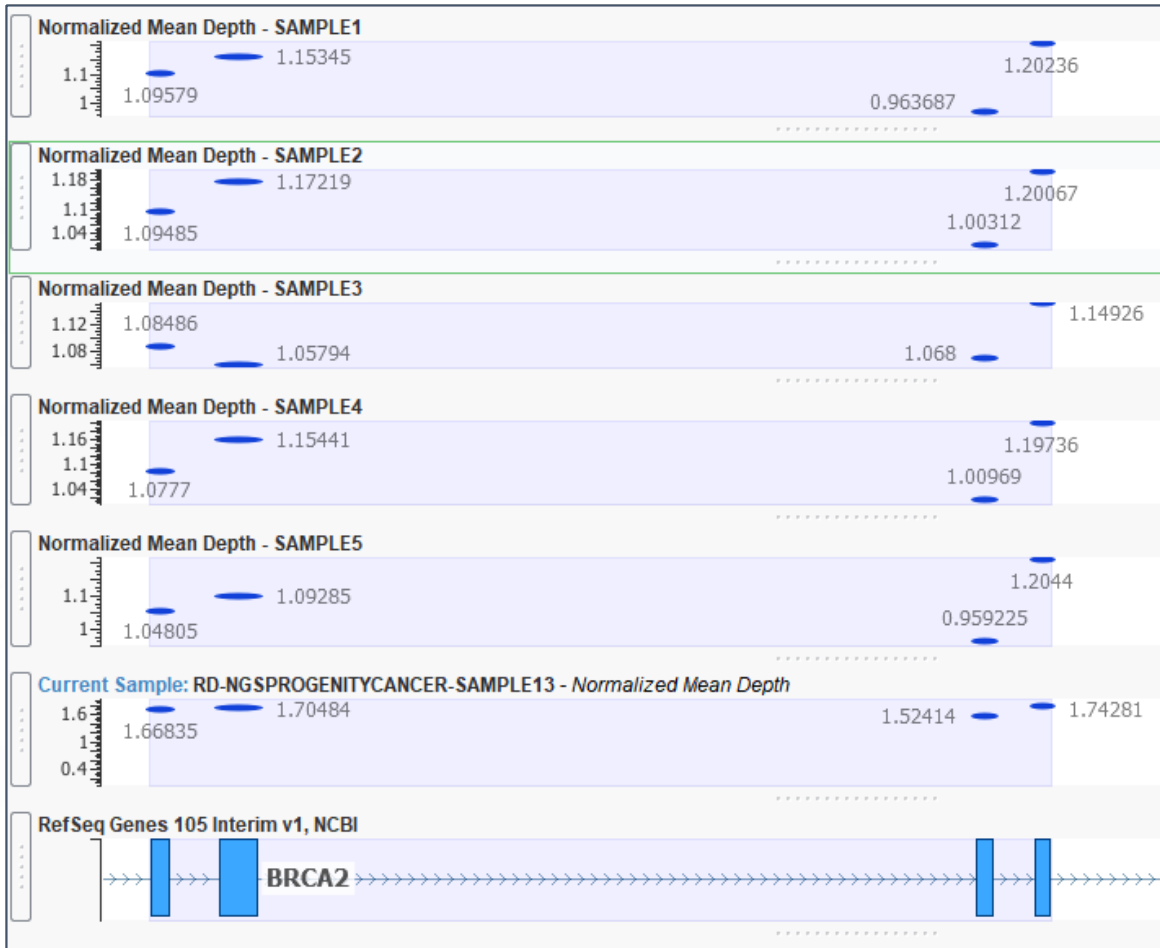
- ✓ One single testing paradigm
- ✓ True simplification of clinical workflow
- ✓ Saves time and money – all on site

# Addressing Issues - CNV Detection via NGS

- CNVs detected from coverage data in BAM
- Challenges
  - Coverage varies between samples
  - Coverage fluctuates between targets
  - \*Systematic biases impact coverage
- Solutions
  - Data Normalization
  - Reference Sample Comparison
  - Algorithm works without case/control data
- Requirements
  - $\geq 30$  ref samples
  - From same library prep method
  - Ideally  $\geq 100X$  coverage

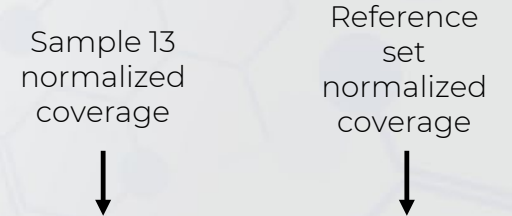


# Principle Approach to CNV Calling



- Reference samples:
  - Coverage normalization and averaging represents diploid/normal regions for comparison
  - Reference set is unique for each sample
  - Reference set is selected based on similarity to the sample
  - Non-autosomal regions matched for gender automatically

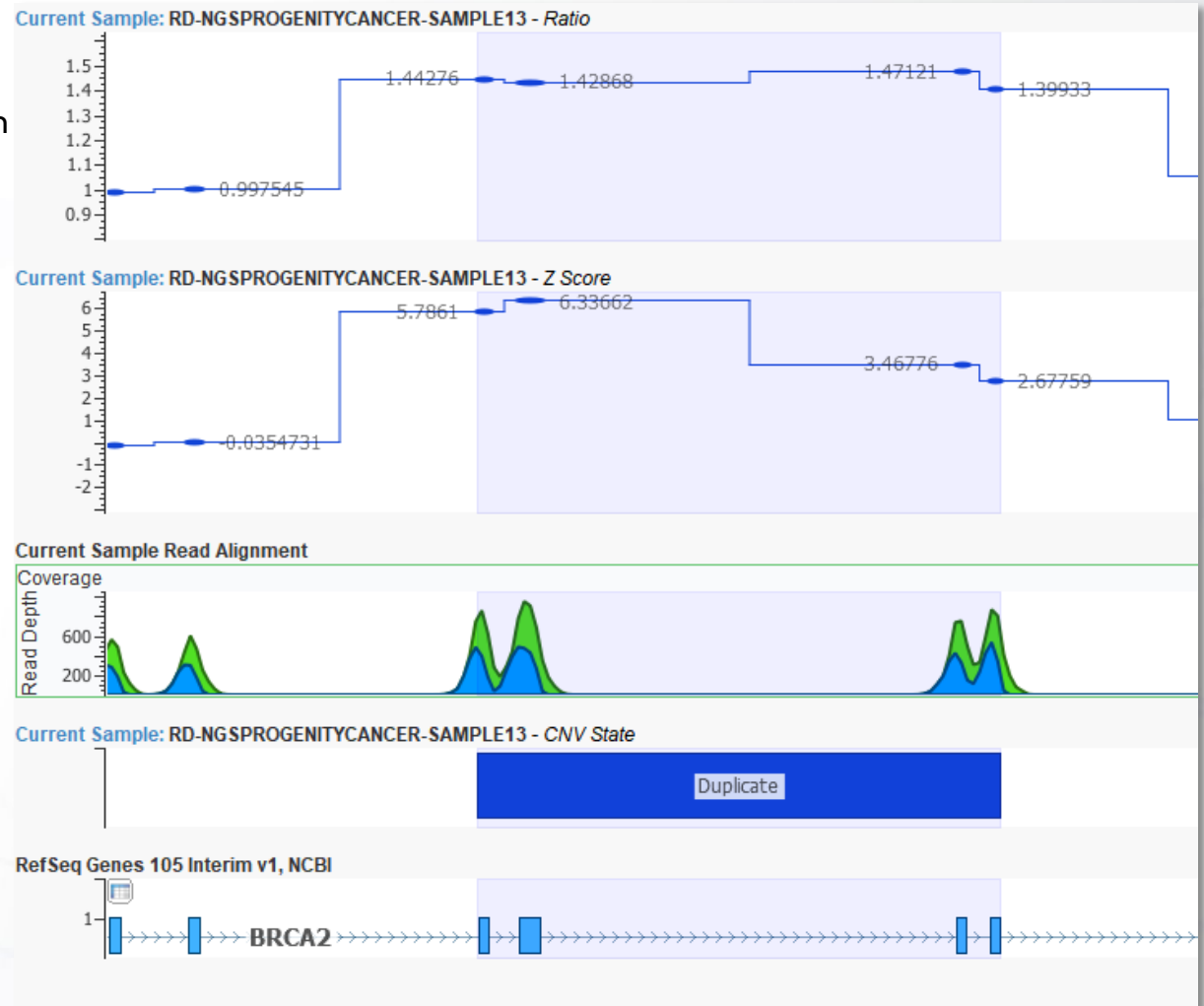
Sample of Interest



Coverage Region I...	Overlapp...	RD-NGSP...	Target Copy Number State for RD-NGSPROGENITYCANCER-SAMPLE13		
Region	Gene Names	Mean Depth	CNV State	Normalized Mean Depth	Avg. Normalized Control Depth
13:32936641-32936850	BRCA2	810.281	Duplicate	1.66835	1.17383
13:32937297-32937690	BRCA2	828.008	Duplicate	1.70484	1.20154
13:32944520-32944714	BRCA2	740.241	Duplicate	1.52414	1.02256
13:32945074-32945257	BRCA2	846.446	Duplicate	1.74281	1.22956

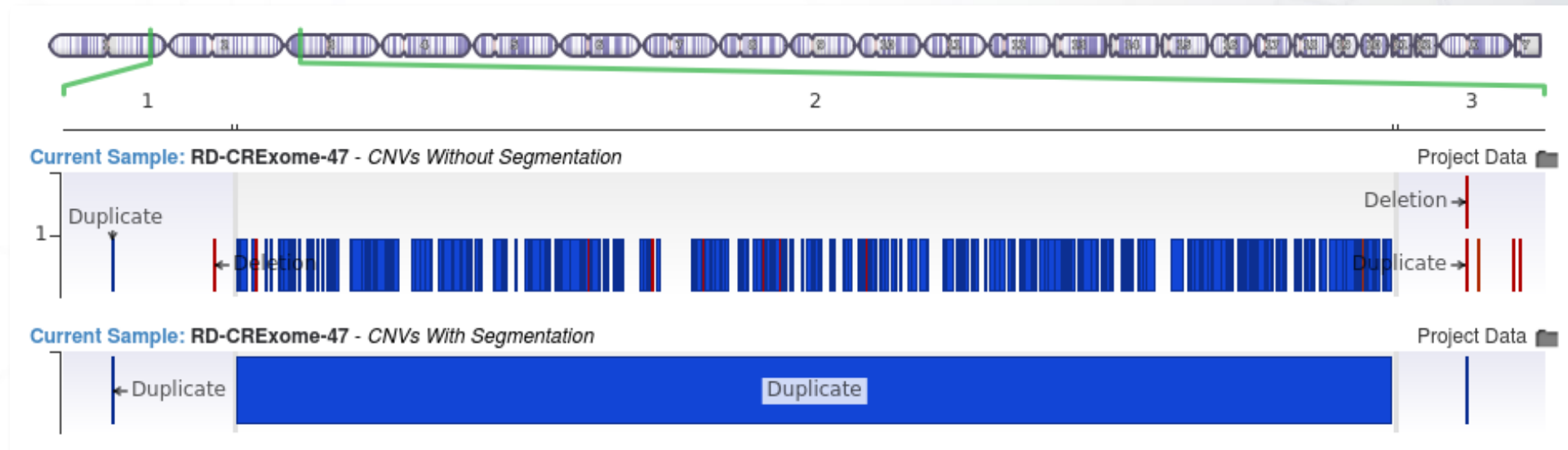
# CNV Detection: Ratio, Z-score, and VAF

- Metrics
  - Ratio: sample coverage divided by reference sample mean
  - Z-score: standard deviations from reference sample mean
  - VAF: Variant Allele Frequency
- For Gene Panels and Exomes
  - Probabilistic model used to call CNVs
  - Segmentation identifies large cytogenetic events
- For Whole Genome Data
  - Targets segmented using Z-scores
  - Events called based on Z-score and Ratio thresholds



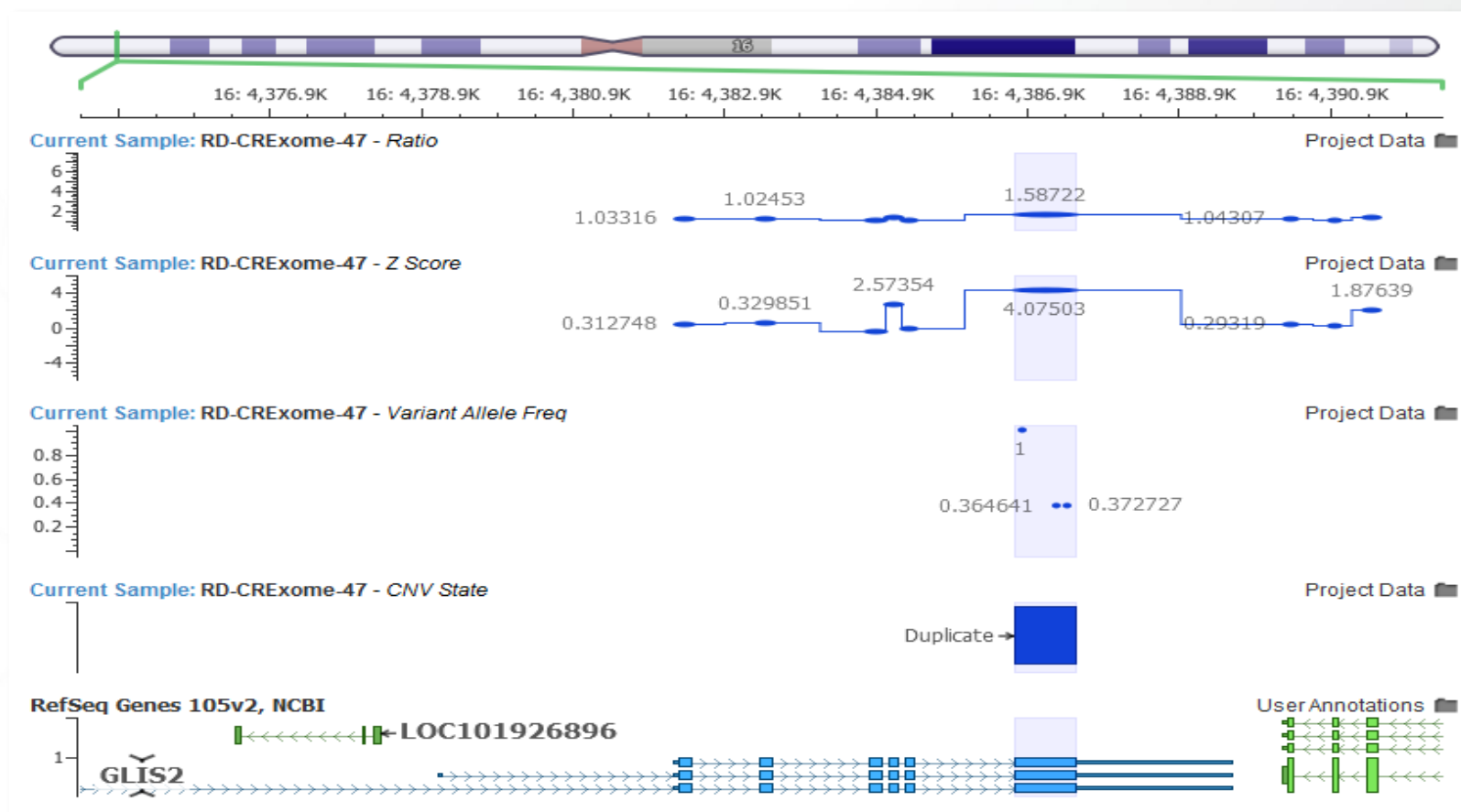
# Optimizing CNV Detection - Segmentation

- Metrics are noisy over large regions
- Outliers cause large events to be called as many small events
- Solved using segmentation:
  - Regions containing many events are segmented
  - Small events sharing a segmented region are merged



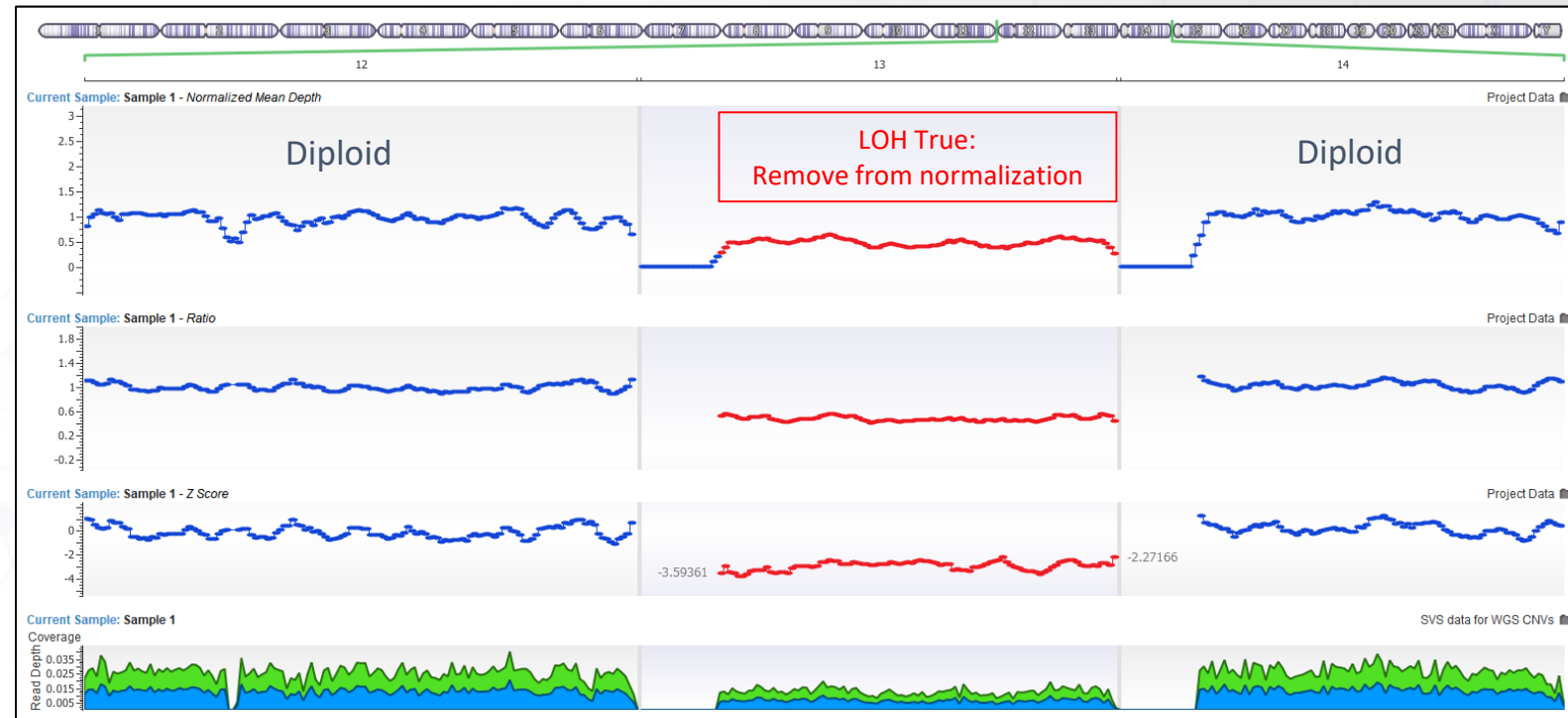
# VAF Provides Supporting Evidence

- Values other than 0 or 1 are evidence against het deletions
- Values of 2/3 and 1/3 are evidence for duplications



# Advanced Optimization for CNV Detection - LOH

- **Issue** - Large chromosomal deletions and duplications can skew the mean coverage of a sample
  - The previous approaches don't account for this
- **Solution** - Detection of LOH areas to optimize definition of normal regions
  - Prior to running CNV caller
  - Probabilistic model based on VAF (Whole Exome/Genome)
  - Identifies and excludes non-diploid regions from normalization
  - Quality control step to improve CNV caller





# CNV Confidence: P-Values

- P-Values

- Introspective ability to define confidence in the CNV event
- Confidence threshold definable in your workflow/protocol
- You can modify the threshold at any point
- Lower the p-value/probability the higher the confidence the event is real

- Typical p-values

- <0.05\*
- <0.01\*\*
- <0.001\*\*\*



The image shows a software interface with a filter menu on the left and a table of CNV data on the right. The filter menu is titled 'p-value (Current) < 0.01' and has a search box containing '0.01'. Below the search box are several filter options with counts: 'Less than 0.01' (2), 'Equal to 0.01' (0), 'Greater than 0.01' (0), and 'Missing' (0). The table on the right is titled 'CNV Info' and has columns for Region, Type, # Targets, # Samples, Span, CNV State, and p-value. The table contains 12 rows of data, with the 7th row highlighted in red.

CNV Info						
Region	Type	# Targets	# Samples	Span	CNV State	p-value
11:108160310...	Loss	21	1	39676	Het Deletion	4.68818626586653e-07
11:108128189...	Gain	9	1	22167	Duplicate	8.93228134373203e-05
17:29683459-...	Gain	6	1	4283	Duplicate	0.00131093873642385
8:90965453-9...	Gain	5	1	17353	Duplicate	0.0050437287427485
3:37053483-3...	Gain	5	1	14036	Duplicate	0.00561221409589052
17:56769986-...	Gain	4	1	10725	Duplicate	0.00574351288378239
11:64573088-...	Gain	4	1	1624	Duplicate	0.0100980764254928
17:41219606-...	Gain	4	1	9046	Duplicate	0.0126859014853835
17:41201106-...	Gain	4	1	14305	Duplicate	0.0156013108789921
11:108151711...	Loss	4	1	6752	Het Deletion	0.0184763930737972
11:108121409...	Gain	3	1	2251	Duplicate	0.0187129583209753
3:37045873-3...	Loss	4	1	7501	Het Deletion	0.0200324188917875

# CNV Workflow and Annotations

- Workflow Issue: How to screen through many CNV events?
  - Especially relevant with WES
  - What steps are necessary to find events relevant to the patient?
- First step
  - Prioritize CNVs specific to sample
  - High quality
  - High confidence
- Second step
  - DGV CNVs - Exclude CNVs in known healthy individuals
  - Genomic Sup Dups - Exclude CNVs in known duplication regions
  - ExAC + 1kG Phase3 - Eliminate common CNVs
  - ClinGen + ClinVar - Eliminate known benign CNVs
- Third step
  - Prioritize individual sample phenotype + gene list

The screenshot shows the 'Data Source Library' window with the 'CNV and Large Variants' category selected. The table below lists the available data sources:

Name	Type	Size	Date	Url
<input type="checkbox"/> 1kG Phase3 CNVs and Large Variants 5b V2, GHI	Interval	2.4M	2017-10-09	rfs://data.golde...
<input type="checkbox"/> ClinGen (ISCA) CNVs 2017-09-10, USCS	Interval	1.4M	2017-09-25	rfs://data.golde...
<input type="checkbox"/> ClinGen Gene Dosage Sensitivity 2017-09-27, NCBI	Interval	172K	2017-09-28	rfs://data.golde...
<input type="checkbox"/> ClinGen Region Dosage Sensitivity 2017-09-27, NCBI	Interval	56K	2019-09-29	rfs://data.golde...
<input type="checkbox"/> ClinVar CNVs and Large Variants 2019-05-01, NCBI	Interval	3.6M	2019-05-03	rfs://data.golde...
<input type="checkbox"/> DGV CNVs - Supporting Variants 2016-05-15, DGV	Interval	64M	2017-08-24	rfs://data.golde...
<input type="checkbox"/> DGV CNVs - Variants 2016-05-15, DGV	Interval	35M	2017-08-30	rfs://data.golde...
<input type="checkbox"/> ExAC XHMM CNV Calls 0.3.1, BROAD	Interval	1.7M	2016-12-07	rfs://data.golde...

The 'Information' panel for the selected track '1kG Phase3 CNVs and Large Variants 5b V2, GHI' is shown below:

**Description**  
This track provides the catalog of CNV and large variant "sites" called by the 1000 Genomes project for 2504 individuals from the 2013-05-02 sequence and alignment release. Besides the description of the variant itself, there is an Alternate Allele Frequency, Allele Counts and the number of Het and Homozygous genotypes for all individuals as well as one for each of specific populations are provided.

**Properties**

Type	Interval
Species	Homo sapiens



Project Demonstration

# NIH Grant Funding Acknowledgments

- Research reported in this publication was supported by the National Institute Of General Medical Sciences of the National Institutes of Health under:
  - Award Number R43GM128485-01
  - Award Number R43GM128485-02
  - Award Number 2R44 GM125432-01
  - Award Number 2R44 GM125432-02
- Montana SMIR/STTR Matching Funds Program Grant Agreement Number 19-51-RCSBIR-005
- PI is Dr. Andreas Scherer, CEO Golden Helix.
- The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.



We're headed to ESHG 2019 in June!