



# Annotating and Cataloging CNVs in VarSeq

Dr. Nathan Fortier – Director of Research

**CIOReview**

20 most promising  
Biotech Technology  
Providers

**pharma**  
TECH OUTLOOK

Top 10 Analytics  
Solution Providers

**Gartner.**

Hype Cycle for  
Life sciences



**Please enter your questions into your GoToWebinar panel**



# 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
  - Award Number 2R44 GM125432-01
  - Award Number 2R44 GM125432-02
- 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.

# Golden Helix – Who We Are



Golden Helix is a global bioinformatics company founded in 1998.



- Variant Calling
- Filtering and Annotation
- Variant Interpretation
- Clinical Reports
- CNV Analysis
- Pipeline: Run Workflows

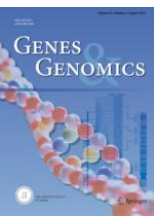
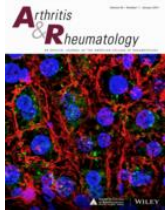
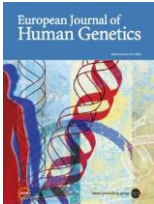
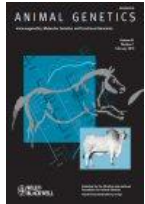


- Variant Warehouse
- Centralized Annotations
- Hosted Reports
- Sharing and Integration



- CNV Analysis
- GWAS
- Genomic Prediction
- Large-N-Population Studies
- RNA-Seq
- Large-N CNV Analysis

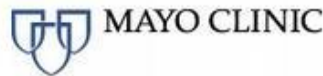
# Cited in over 1300 peer-reviewed publications



# Over 400 customers globally



Stanford University





## When you choose a Golden Helix solution, you get more than just software

- REPUTATION
- TRUST
- EXPERIENCE



- INDUSTRY FOCUS
- THOUGHT LEADERSHIP
- COMMUNITY

- TRAINING
- SUPPORT
- RESPONSIVENESS







- INNOVATION and SPEED

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 concordance with ACMG Guidelines</li> </ul>
 VSWarehouse	Data Warehousing Web-Enabled Interface + Powerful API: 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>





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



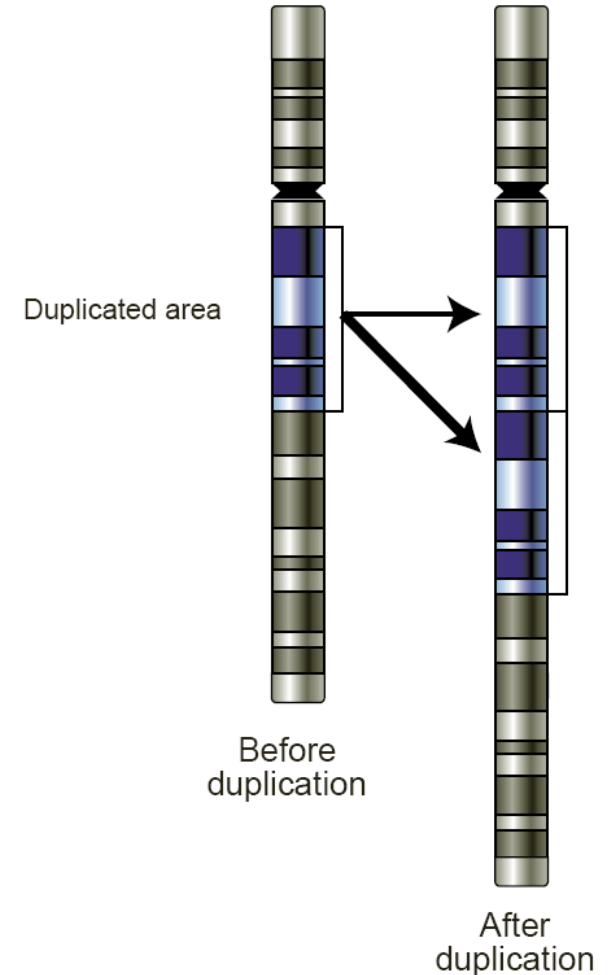


- **Chromosomal microarray**

- Current best practice
- Slow
- Additional expense
- Only detects large events

- **CNV calling from NGS data**

- Calls from existing coverage data
- Detects small single-exon events
- Provides faster results, simplified clinical workflow

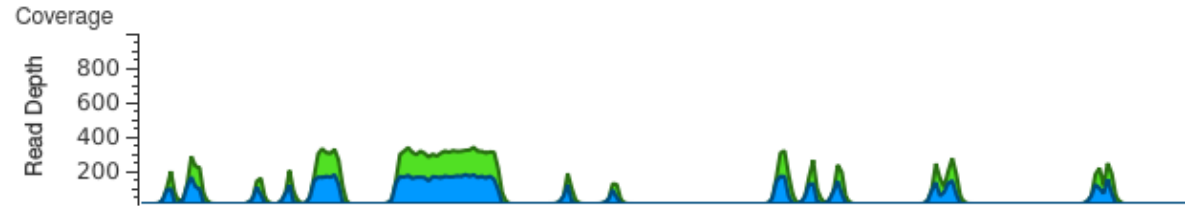


# CNV Detection via NGS

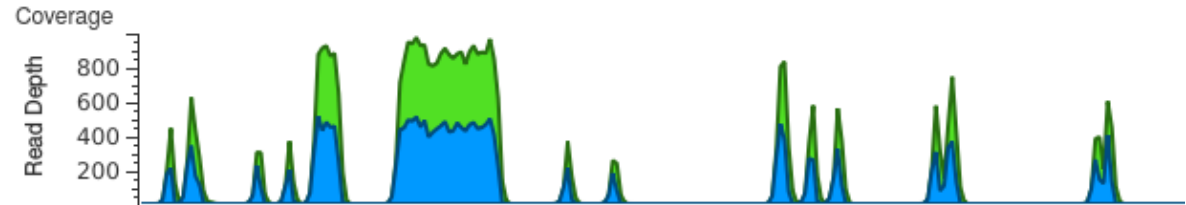


- **CNVs are called from coverage data**
- **Challenges**
  - Coverage varies between samples
  - Coverage fluctuates between targets
  - Systematic biases impact coverage
- **Solutions**
  - Data Normalization
  - Reference Sample Comparison

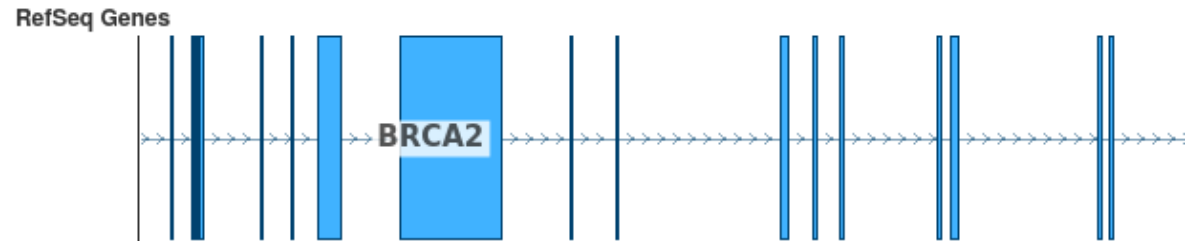
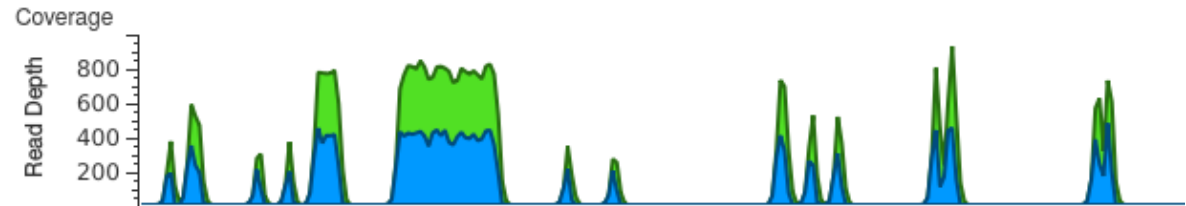
Current Sample: RD-NGSPROGENITYCANCER-SAMPLE11



Current Sample: RD-NGSPROGENITYCANCER-SAMPLE12



Current Sample: RD-NGSPROGENITYCANCER-SAMPLE13



# CNV calling in VarSeq



- **Reference samples used for normalization**

- **Metrics**

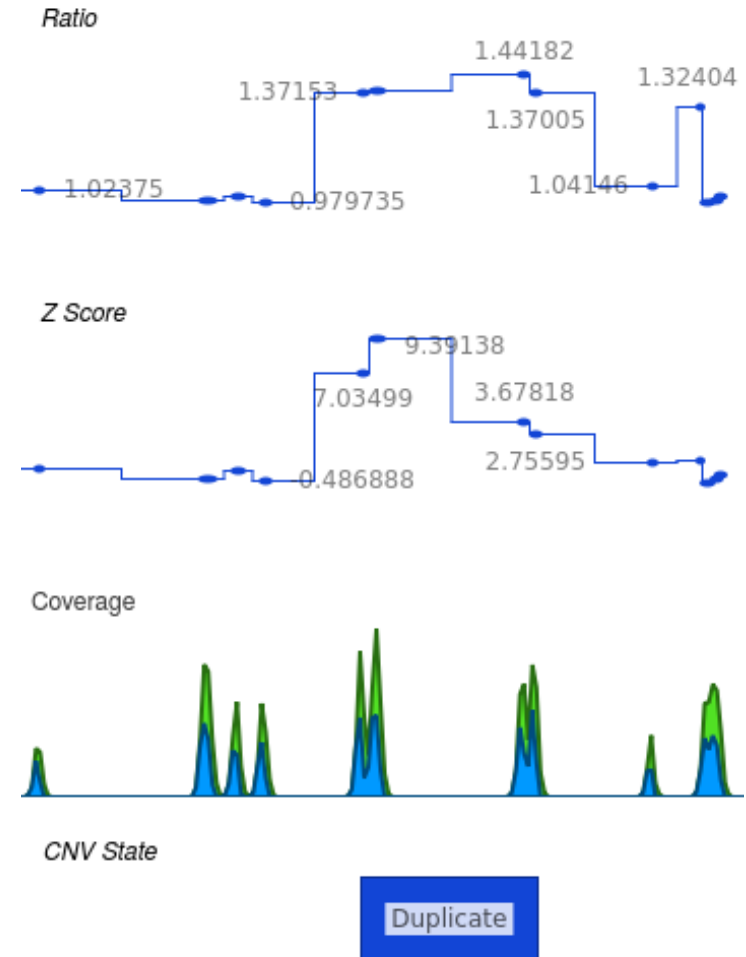
- Z-score: number of standard deviations from reference sample mean
- Ratio: sample coverage divided by reference sample mean
- VAF: Variant Allele Frequency

- **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





- **Low quality events can be flagged if**
  - Event targets have low coverage
  - There is high variation between samples at event targets
  - Event cannot be differentiated from noise at a region
  
- **Samples can be flagged if**
  - The sample does not match the references
  - The sample has extremely low coverage
  - There is high variance across the target regions
  
- **Filtering flagged events improves precision**

# Reference Samples



- **Match references are chosen for each sample**
- **Samples with lowest percent difference chosen**
- **Performance affected if controls don't have matching coverage profile**
- **Samples are flagged if the average percent difference is above 20%**



- **100x Coverage**
- **Reference samples**
  - Recommend at least 30 references
  - From same platform and library preparation
  - Automatic gender matched references for non-autosomal calls

# Sources for Annotating CNVs



- **CNV calls in Populations:**
  - 1000 Genomes Phase3 Large Variants
  - ExAC per-sample CNV calls
  - DGV large-cohort studies
- **Clinical Interpretations:**
  - ClinVar Large Variants
  - ClinGen (Previously ISCA)
- **Genes**
  - Gene track, which transcripts/exons
  - Special considerations considering large sizes
- **Regions**
  - Genomic Superdups (Large Scale)
  - Low Complexity Regions (Smaller Scale)

Select Data Source

Select tracks to use as annotation sources against the imported variant set.

Locations: Local

Filter: \* (Any typ) Homo sapiens (Human), GRCh37 g1k (Fe)  Current

Name	Type
<input type="checkbox"/> 1kG Phase3 - CNVs and Large Variants 5b, GHI	In
<input type="checkbox"/> Cancer Hotspot Panel v2 - Hotspots	In
<input type="checkbox"/> Cancer Hotspot v2 Panel Design	In
<input type="checkbox"/> CIViC - Region Clinical Evidence Summaries 2017-06-01, WUSTL	In
<input type="checkbox"/> ClinGen (ISCA) 2017-09-10, USCS	In
<input type="checkbox"/> ClinVar CNVs and Large Variants, NCBI	In
<input type="checkbox"/> CNV Catalog	In
<input type="checkbox"/> COSMIC Cancer Gene Census 71, GHI	In
<input type="checkbox"/> CpG Islands	In
<input type="checkbox"/> DAC Blacklisted Regions, ENCODE	In
<input type="checkbox"/> Danger Track Regions	In
<input type="checkbox"/> dbNSFP Gene Annotation with Entrez Gene Coordinates and MedGen 2.9, GHI	In
<input type="checkbox"/> DGV SupportingVariants 2016-05-15, DGV	In
<input type="checkbox"/> DGV Variants 2016-05-15, DGV	In
<input type="checkbox"/> DNase Hypersensitivity Sites	In
<input type="checkbox"/> Ensembl Genes 75v2, Ensembl	G
<input type="checkbox"/> ExAC XHMM CNV Calls 0.3.1, BROAD	In
<input type="checkbox"/> GENCODE Genes 19, GENCODE	G
<input type="checkbox"/> Gene Ontology 2017-05-09	Ta
<input type="checkbox"/> Genomic Super Dups 2011-10-25, UCSC	In

Information showing (38/152), 0 selected (0 bytes) Clear

Convert... Download Select Cancel Help



# Annotation Algorithms: Overlapping Regions

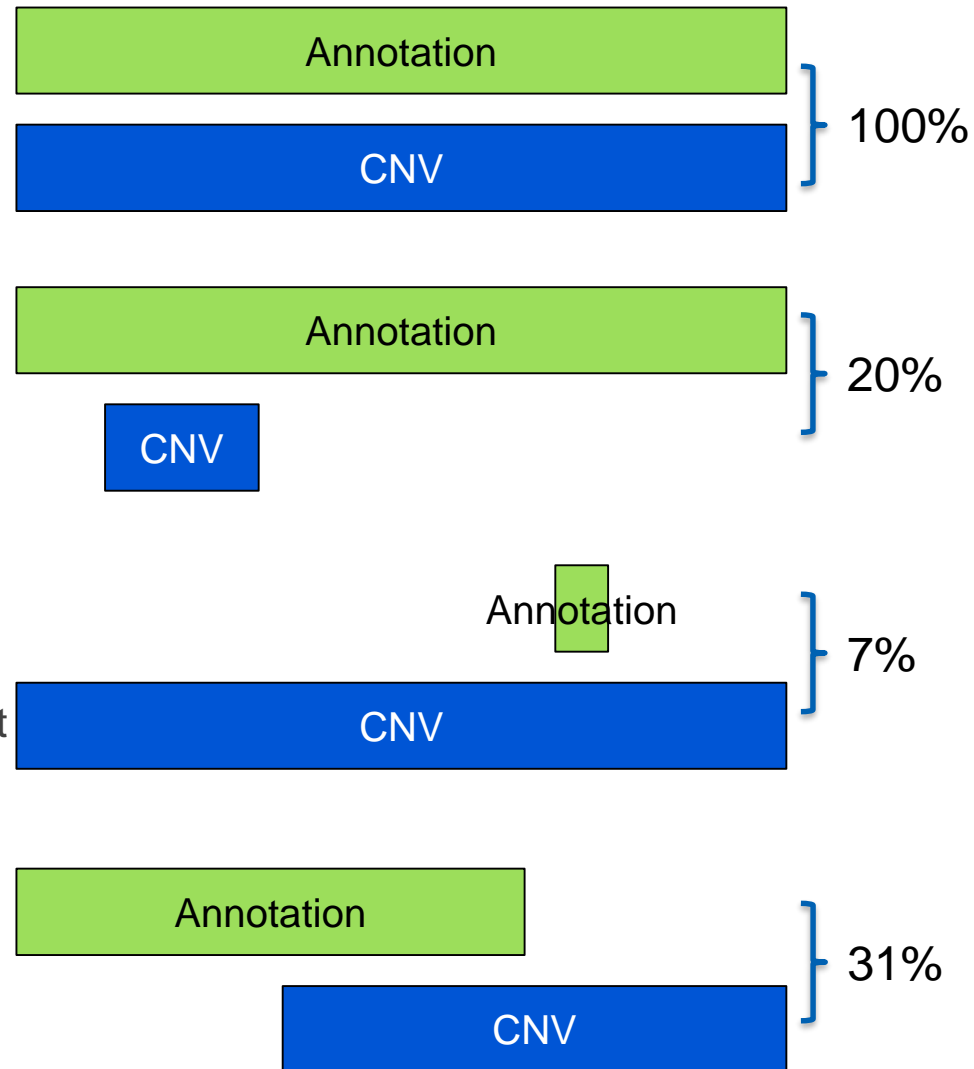


- Not expect exact matches
- Percent overlap not correct metric
- Need metric of “sameness”
- Jaccard index:

- “similarity coefficient”

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

- For fully overlapped regions, the percent overlap of the smaller to the larger
- Default value of 20% for annotations
- If set to 0%, then any overlap matches
- If set to 100%, then exact matches





# 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
  - Award Number 2R44 GM125432-01
  - Award Number 2R44 GM125432-02
- 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.





**AMP**2018 ANNUAL MEETING & EXPO

Precision Medicine Starts Here

**NOVEMBER 1-3, 2018** Henry B. Gonzalez Convention Center  
San Antonio, TX, USA

**See us at AMP 2018 in Booth #1801**

**November 1-3, 2018 | San Antonio, TX**

**Booth 1801**

**5-minute demos at the Golden Helix booth**

**Free t-shirts & chance to win an iPad Pro**

