

Fine-tuning CNV Analysis for the Clinical Analysis of NGS Samples

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



Golden Helix – Who We Are



Golden Helix is a global bioinformatics company founded in 1998.



Variant Calling

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

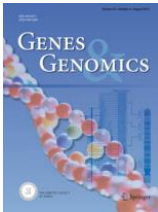
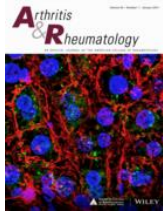
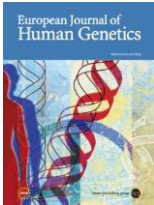
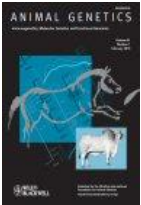


- Variant Warehouse
- Centralized Annotations
- Hosted Reports
- Sharing and Integration

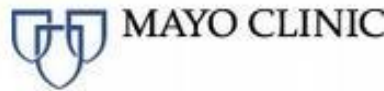


- GWAS
- Genomic Prediction
- Large-N-Population Studies
- RNA-Seq
- Large-N CNV-Analysis

Cited in over 1100 peer-reviewed publications



Over 350 customers globally

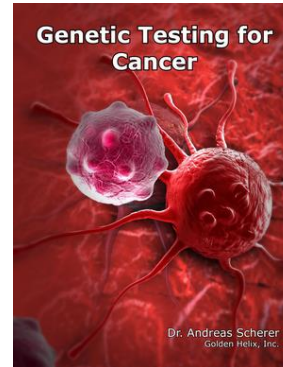


Golden Helix – Who We Are



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
- CUSTOMIZATIONS

SEQUENCER

PRODUCTS

BIOINFORMATICS PIPELINE

FUNCTION



VS-CNV



SENTIEON DNASEQ



SENTIEON TNSEQ

OMIM SIFT & POLYPHEN CLINVAR ENSEMBL GENES
CADD EXAC & GNOMAD EXOMES DBSNP REFSEQ GENES
ONCO MD CONSERVATION SCORES COSMIC

FASTQ

SINGLE NUCLEOTIDE VARIATION

BAM

COPY NUMBER VARIATION & LOSS OF HETEROZYGOSITY

VCF

CHROMOSOMAL ABERRATION

ANNOTATE

PUBLIC & COMMERCIAL ANNOTATIONS
TO ENRICH GENOMIC DATA SETS



VARSEQ

VSREPORTS

VSPipeline

CLINICAL REPORT

ANNOTATE & FILTER
VISUALLY INSPECT ALIGNMENTS
VARIANT PRIORITIZATION
CLINICAL ASSESSMENT



WAREHOUSE

DATA WAREHOUSING

CLINICAL ASSESSMENT CATALOG
ADVANCED DATA QUERYING
VERSIONING

WEB-ENABLED INTERFACE
+ POWERFUL API: JSON, XML
TSV, CSV, SQL, FHIR

INTEROPERABILITY
COMPLIANCE WITH HIPAA, CLIA, & CAP
DATA DISCOVERY



- **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)

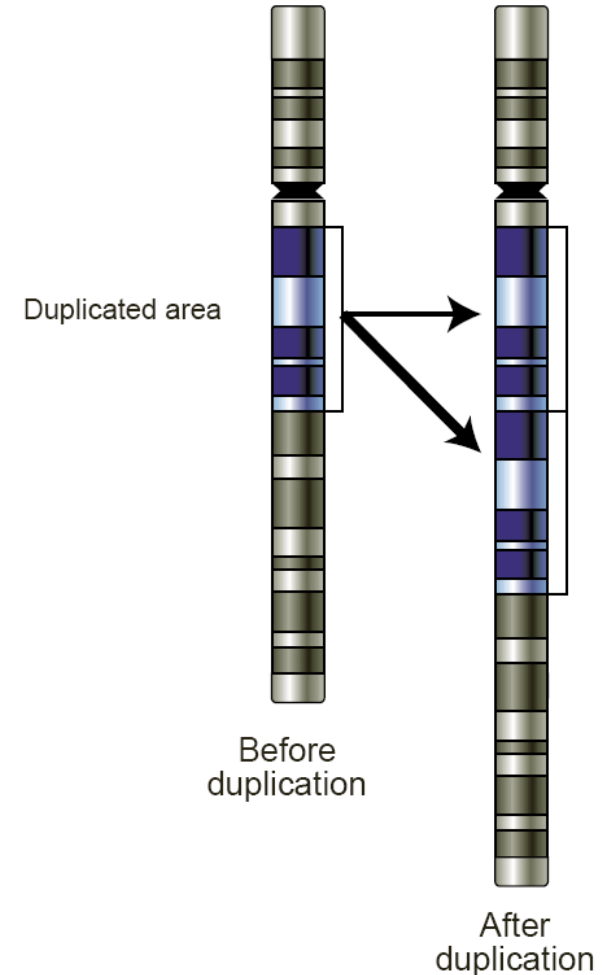


- **Chromosomal microarray**

- 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

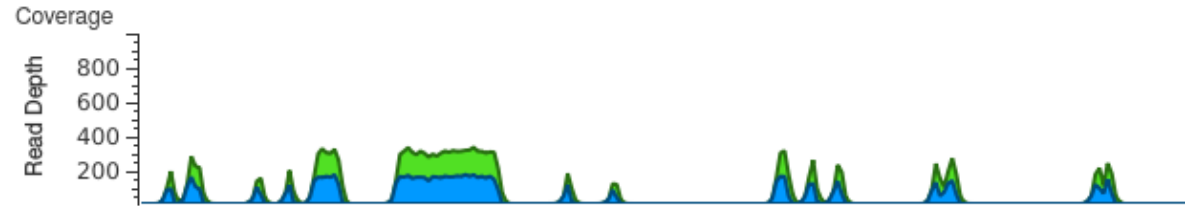


CNV Detection via NGS

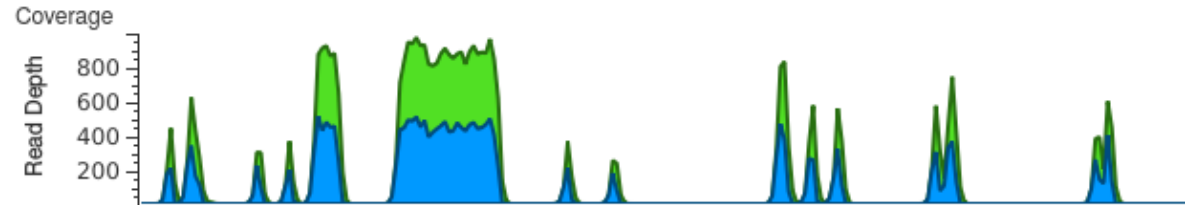


- **CNVs are called from coverage data**
- **Challenges**
 - Coverage varies between samples
 - Coverage fluctuates between targets
 - Systematic biases impact coverage
- **Solutions**
 - Data Normalization
 - Reference Sample Comparison

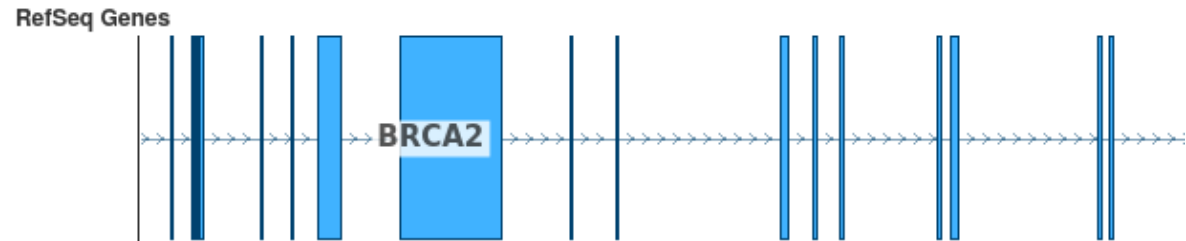
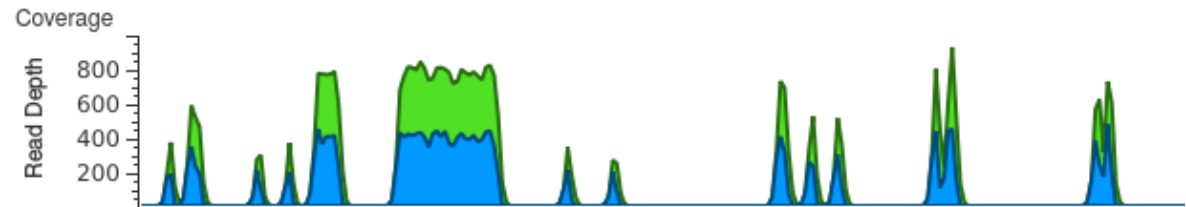
Current Sample: RD-NGSPROGENITYCANCER-SAMPLE11



Current Sample: RD-NGSPROGENITYCANCER-SAMPLE12



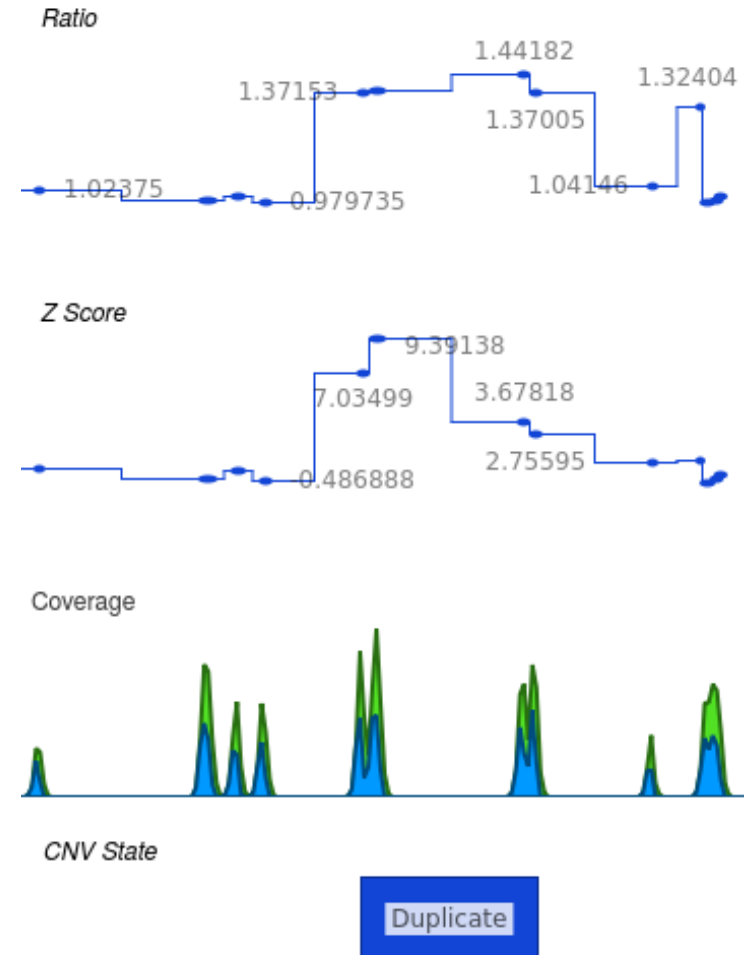
Current Sample: RD-NGSPROGENITYCANCER-SAMPLE13



CNV calling in VarSeq



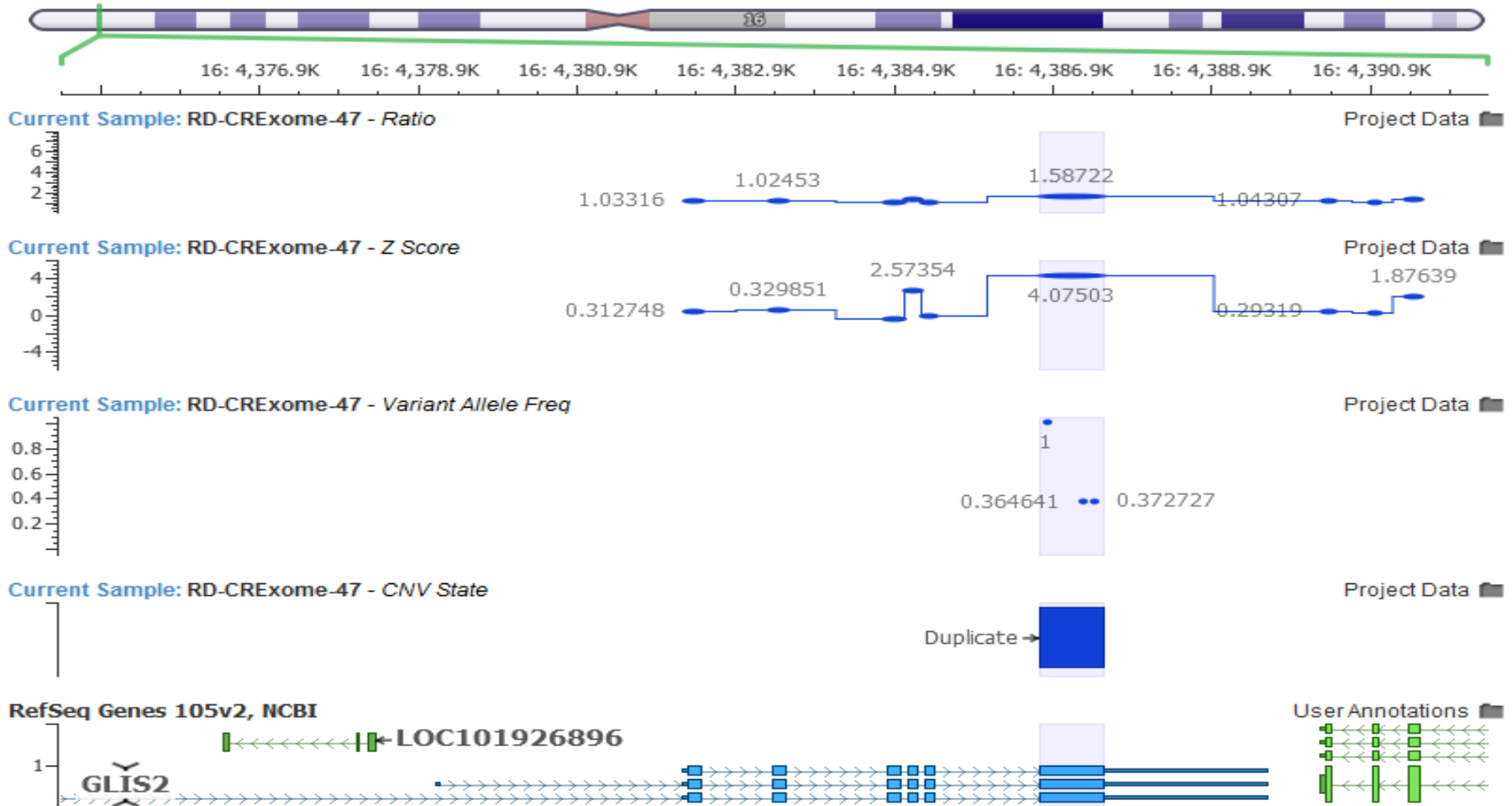
- **Reference samples used for normalization**
- **Metrics**
 - Z-score: number of standard deviations from reference sample mean
 - Ratio: sample coverage divided by reference sample mean
 - VAF: Variant Allele Frequency
- **For Gene Panels and Exomes**
 - Probabilistic model used to call CNVs
 - Segmentation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





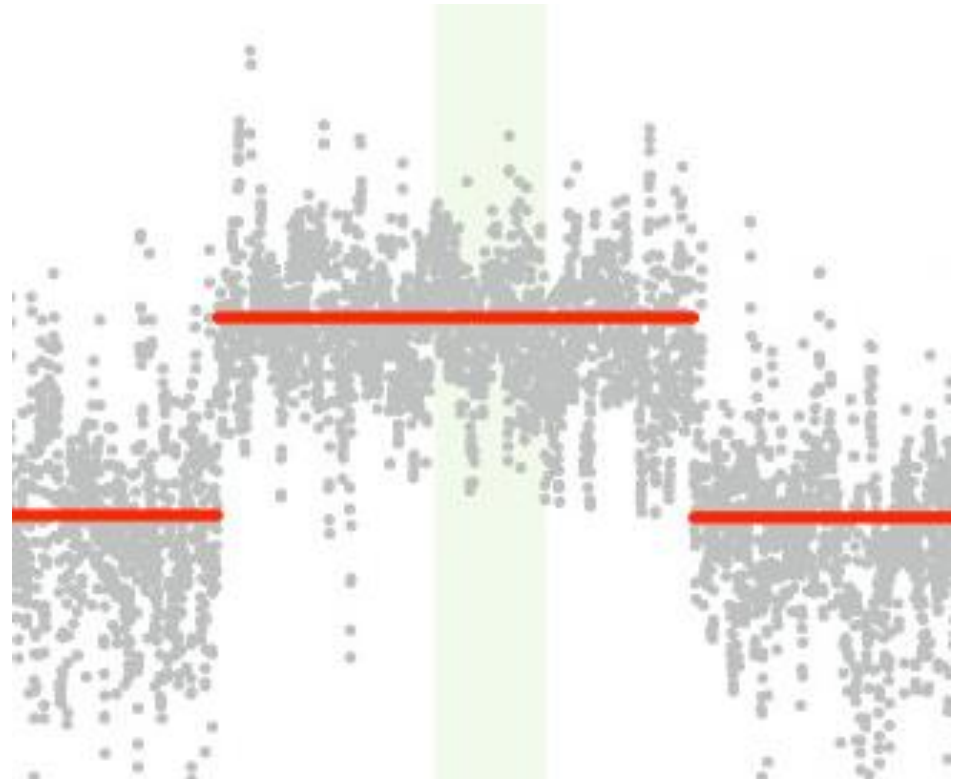
VAF provides supporting evidence

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



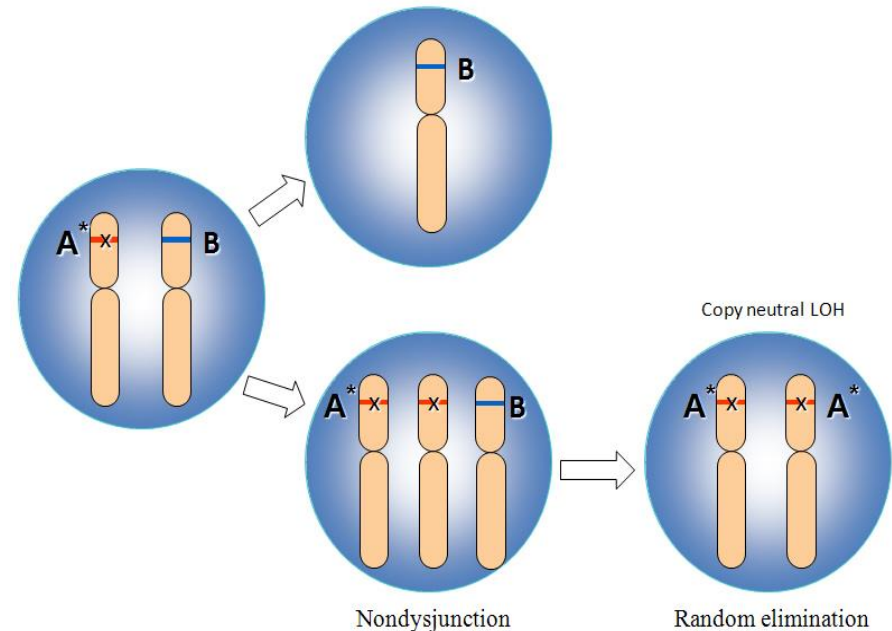


- **Metrics are noisy over large regions**
- **Outliers cause large events to be called as many small events**
- **Addressed using segmentation:**
 - CNAM Optimal Segmentation
 - Regions containing many events are segmented
 - Small events sharing a segmented region are merged

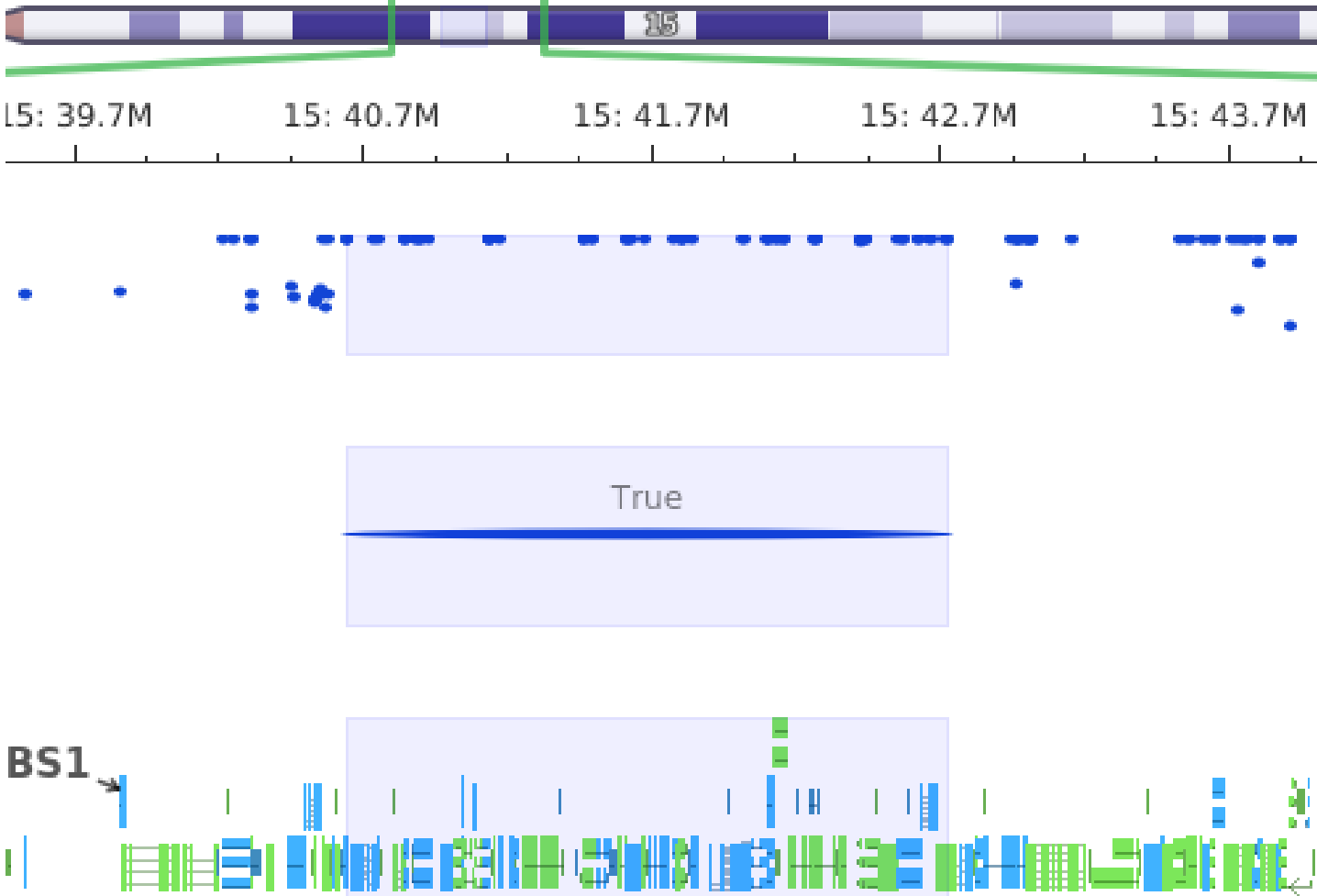




- Large LoH events need to be interpreted in any gene test that covers large CNVs
- New Loss of Heterozygosity(LOH) detection based on H3M2 (Magi *et al.*)
- Calls LoH events using Hidden Markov Model (HMM)
 - Observations are variant allele frequencies
 - States are either Homozygous or Non-Homozygous



LoH Calling



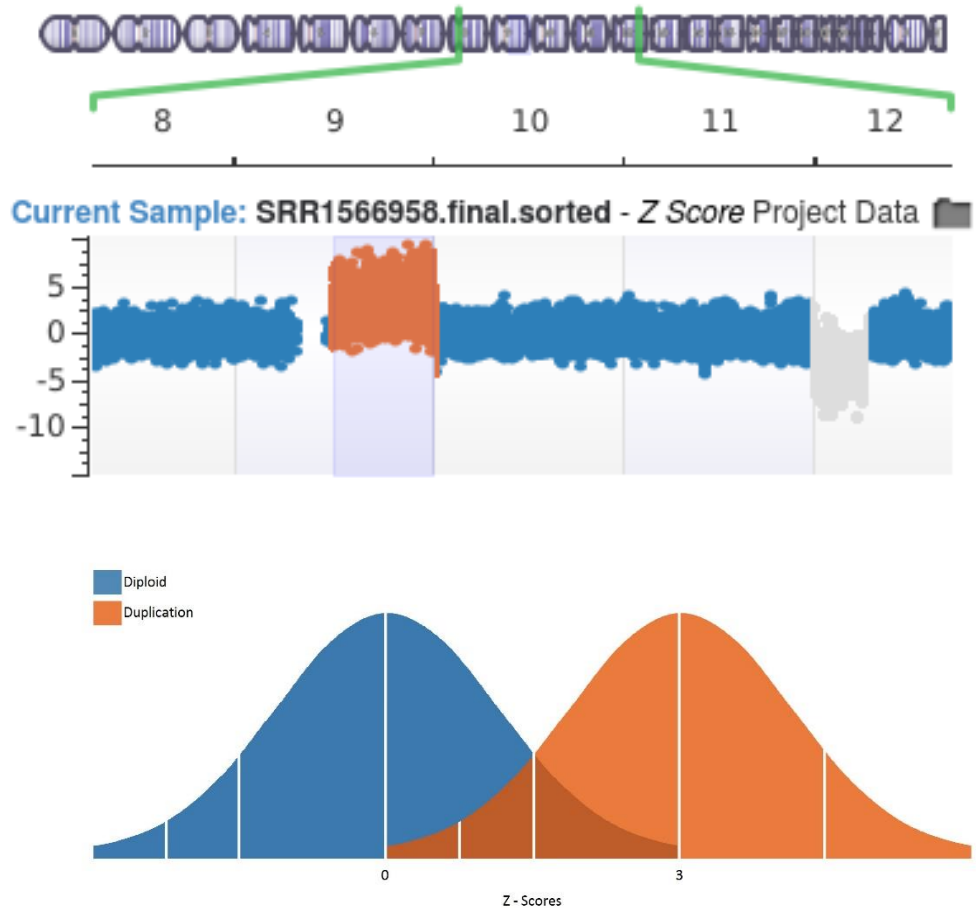


■ P-Values

- Probability of z-scores at least as extreme assuming the event targets are diploid
- Computed using Student's t-test
- Distribution of event z-scores compared to distribution of diploid targets

■ Quantifies CNV Call Confidence

- Values below 0.01 indicate high confidence calls
- Values above 0.01 indicate lower confidence calls

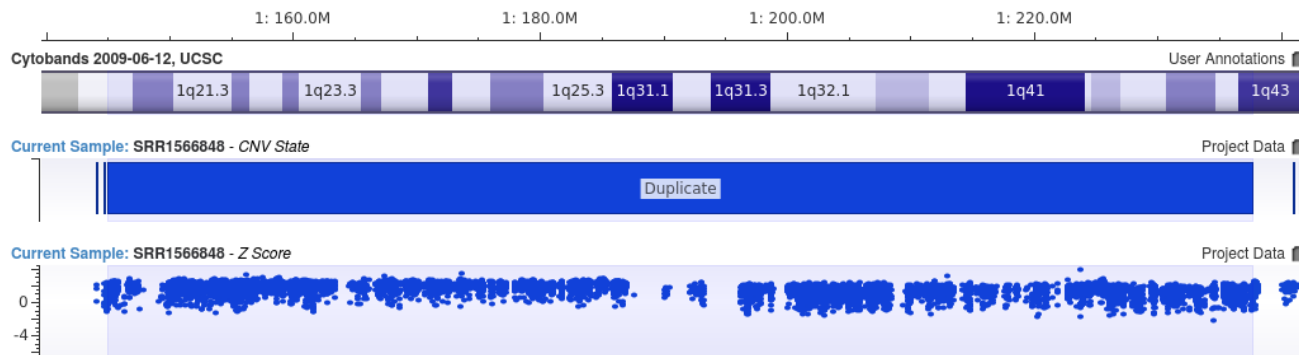


$$p = 1.4 \cdot 10^{-32}$$

Karyotype Notation



- Karyotype notation provided for large cytogenetic events
- Karyotypes provided at both event and sample level
- Uses common notation
- Specifies chromosome, arm, and band for each mutation



46,XY,dup(1)(q21.1q43)



- **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
 - Minimum of 10
 - From same platform and library preparation
 - Gender matched references required for non-autosomal calls



- **Sex is inferred from coverage data**
 - Sample is inferred female if
 - Y chromosome coverage is low
 - X chromosome coverage matches the autosome
 - Otherwise the sample is inferred to be male
- **Samples are matched on inferred sex**
- **Same-sex samples are used for normalization of non-autosomal chromosomes**

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

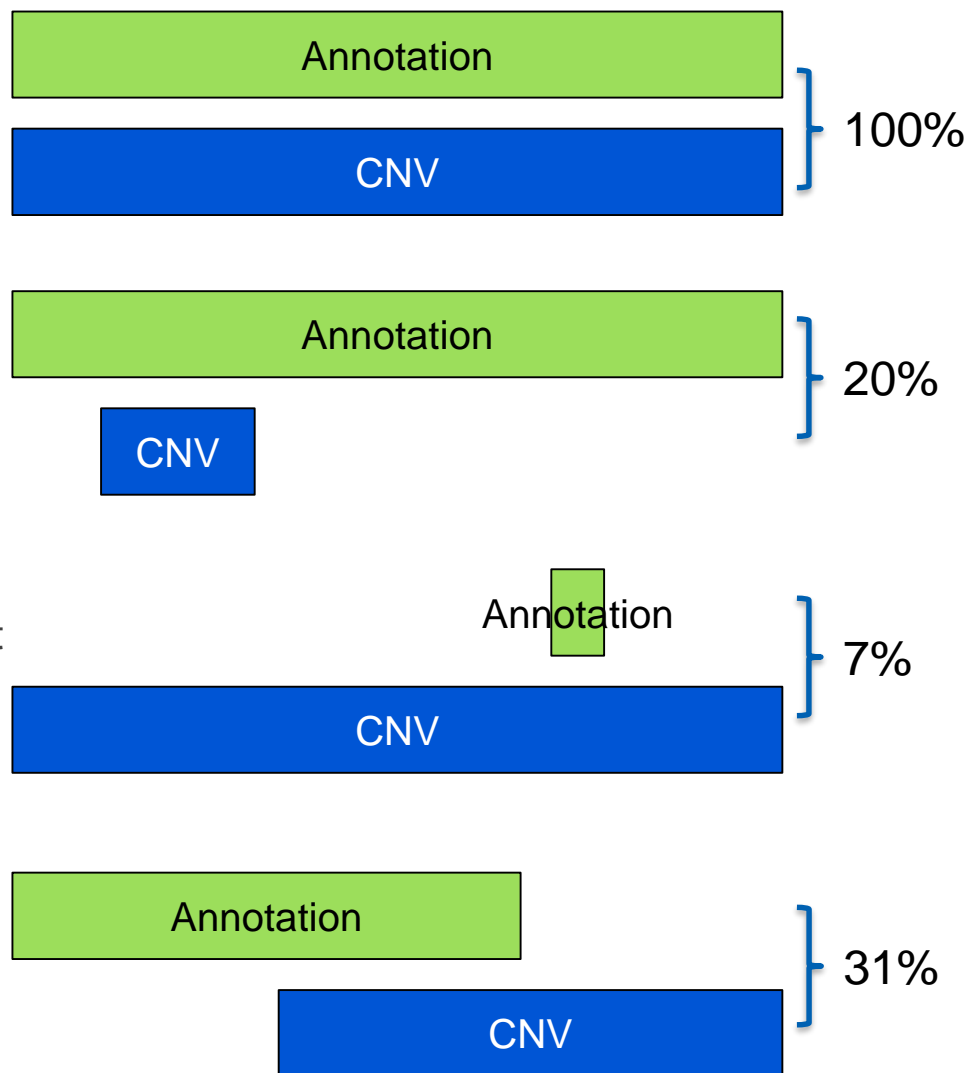
Annotation Algorithms: Overlapping Regions



- Not expect exact matches
- 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





Please enter your questions into your GoToWebinar Panel





2018

ACMG Annual Clinical Genetics Meeting

APRIL 10-14 | EXHIBIT DATES: APRIL 11-13
CHARLOTTE CONVENTION CENTER | CHARLOTTE, NC

We're Exhibiting at ACMG 2018

Find us (and the t-shirts) in booth 1306!



Please enter your questions into your GoToWebinar Panel

