# VSClinical

## Using the GRCh38 reference assembly for clinical interpretation in VSClinical

Gabe Rudy | VP of Product & Engineering

# NIH Grant Funding Acknowledgments

# Q & A

**Please enter your questions into your GoToWebinar Panel**

# Golden Helix − Who We Are

**Golden Helix is a global bioinformatics company founded in 1998, celebrating our 20th year!**

## VarSeq

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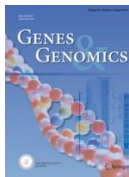**CNV Analysis**
**Pipeline: Run Workflows**

## WARE HOUSE

Variant Warehouse
Centralized Annotations
Hosted Reports
Sharing and Integration

## SNP & VARIATION SUITE

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

GOLDEN HELIX
*Enabling Precision Medicine*

# Cited in over 1,300 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

# Agenda

Genetic Testing with NGS

Variant Representation

Human Reference Genomes

Implications for Variant Interpretation

Demo using VarSeq + VSClinical

Motivation for Using GRCh38

Other Lab Considerations

Thanks! / Q&A

# Representing a Variant

- **Genomic:**
  - chr2: 47,641,560 A/T
  - NC_000002.11:g.47641560A>T
  - chr14: 51,378,590 TT/T
  - NC_000014.8 :g.51378593delT

- **Gene Coding Sequence:**
  - BRAF c.1799T>A
  - NM_058197.4:c.105dupG
  - LRG_218t1(*MSH2*):c.942+3A>T

- **Gene Protein Sequence**
  - DYX1C1 p.E417*
  - NP_000483.3:p.Phe508del

- **Genomic Representation Enables**
  - Precise lookup of annotations
  - Overlap / relation to genomic features
  - Representation of non-genic variants

- **Coding Representation Enables**
  - Genomic reference independent
  - UTR and Intronic variants
  - Informative representation of coding change
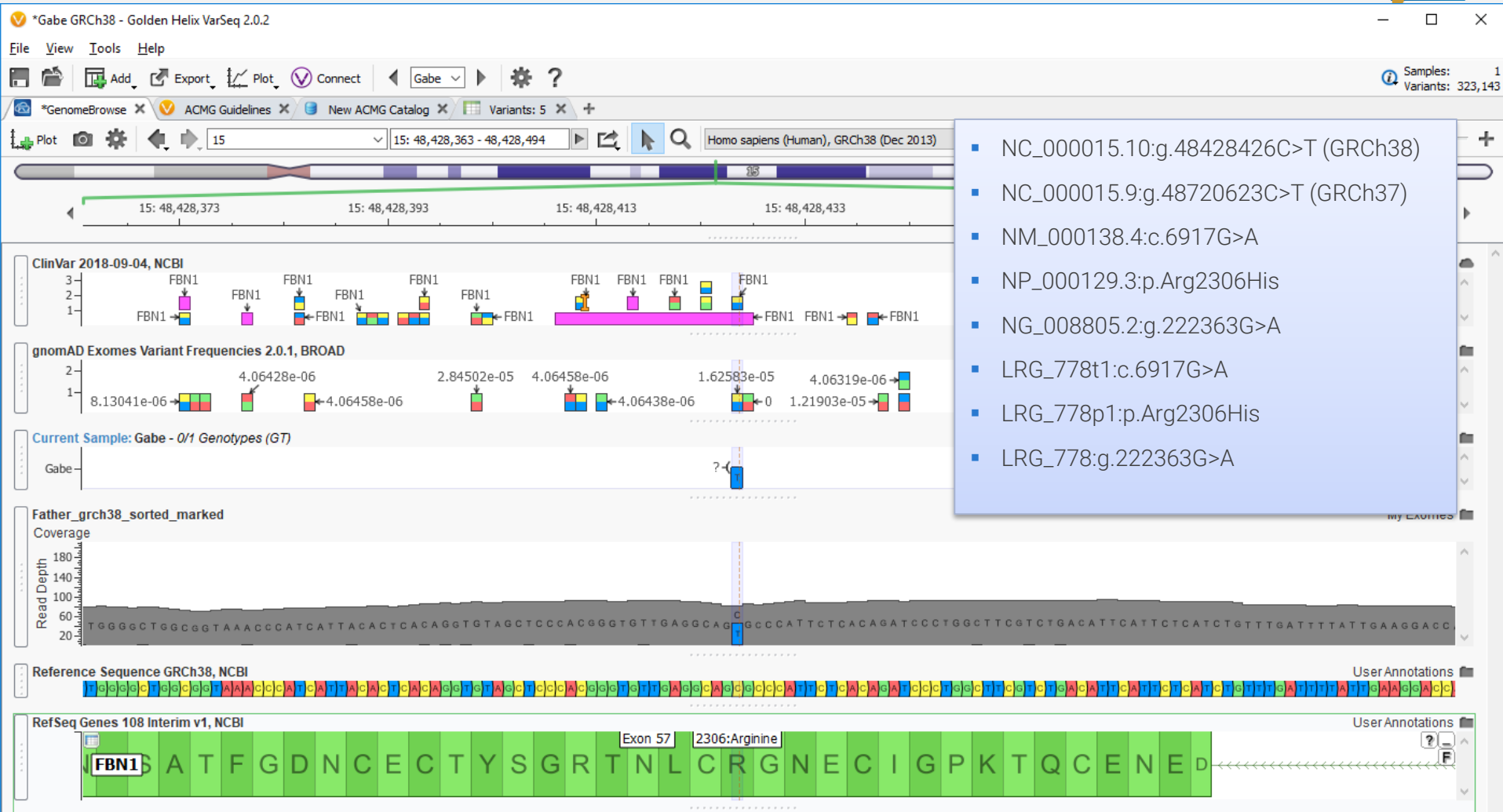
- **Protein Representation Enables**
  - Grouping of variants that result in same protein
  - Descriptive of effect on protein
  - Coordinates match domains and protein DBs

# Genomes Are Just a Means to an End (Genes)

- **RefSeqGenes – mRNA sequence archive, with mappings to genomes**
  - Provided mappings to Locus Reference Gene (LRG) database
  - Use genome mappings by NCBI (through genome annotation builds). NOT UCSC
  - "Clinically Relevant" transcript in VarSeq:
    - Most commonly submitted to ClinVar
    - LRG if available, longest if tied

- **Ensembl – defined directly against the human genome**
  - More inclusive of genes discovered with high-throughput methods
  - Gencode subset – similar to RefSeqGenes in size / definition

- **Each have unique Accessions and Version Numbers**
  - Newer releases are provided only on GRCh38
  - GRCh37 mappings not being updated ("105 Interim" by special request)

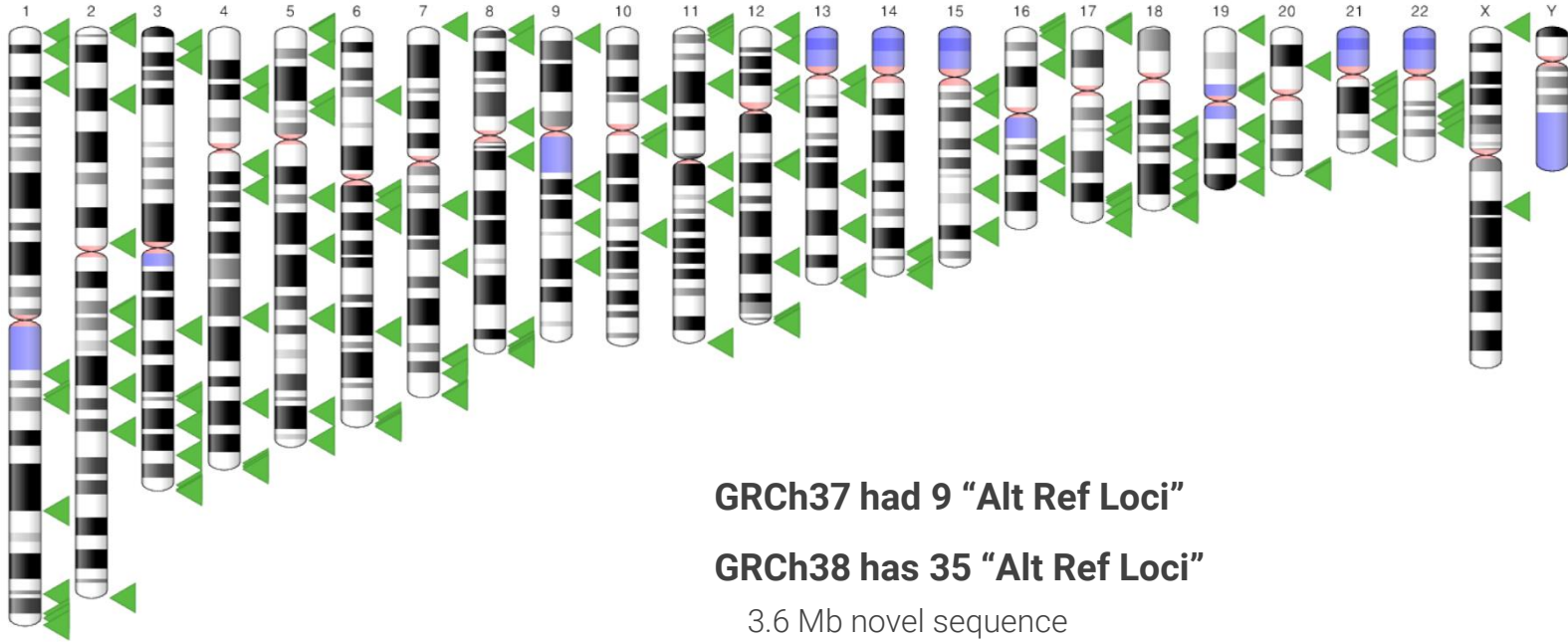# Variant Representation and the Reference Genome

# History of the Human Reference Genome

- **2003: Human Genome Project Declared Done**

- **2006: NCBI36 (hg18)**
  - Produced by the International Human Genome Sequencing Consortium
  - Used by first high throughput sequencers (Illumina GAII), pilot project of 1000 Genomes
  - UCSC uses its own sequential versioning, calling this hg18

- **2009: GRCh37 (hg19)**
  - Handed over to the Genome Reference Consortium (GRCh)
  - Used by the 1000 genome project (Phase I/II/III) in the era of the HiSeq 2000

- **2013: GRCh38 (hg38)**
  - ~100 assembly gaps updated, ~2000 erroneous alleles fixed
  - Included centromere models, mitochondrial reference, alternate sequences

# Alternative Loci / "Haplotypes"



**GRCh37 had 9 "Alt Ref Loci"**

**GRCh38 has 35 "Alt Ref Loci"**

3.6 Mb novel sequence

153 genes

Up to 25% of these genes hare medically interpretable

**Alignment support**

Before using, ensure aligner can support alt loci without flagging "multi-alignment" codes that cause reads to be filtered out / lost. BWA-MEM supports alt loci.

# More than Chromosomes in your FASTA

- **Other bits of the reference:**
  - Un-localized scaffolds assigned to chromosomes
  - Unplaced scaffolds (not assigned to chromosomes
  - Patches Releases (i.e. GRCh37.p13, GRCh38.p12)
    - Types of "alt", "fix" or "novel"
    - Not applied, and do not change the primary sequence
    - You can think of them as "known issues, with proposed fixes for next major release"

- **Other useful things to add for alignment purposes:**
  - A "decoy" reference genome segment as primary reference
    - DNA virus: human herpesvirus 4, type 1, aka Epstein-Barr virus (EBV)
    - Unique sequence found in HuRef (Craig Venter's genome) or de novo assemblies
    - Other novel unaccounted for (or "novel" patch) sequence
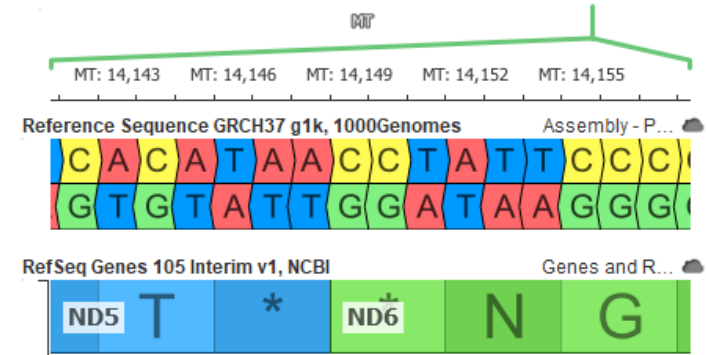  - Full set of HLA "haplotype" sequences, marked as "alternates"

- **Mitochondrial!**

# The Human Mitochondrial

- **Our second genome:**
  - Only 16Kb long
  - Encodes 37 genes (product of energy and its storage in ATP)
  - Slightly different genetic code than nuclear genes:
    - UGA = tryptophan, AUA = methionine, and AGA and AGG = stop



- **Sequence in 1981 as the "Cambridge Reference Sequence" (before HGP)**
  - 2014: "revised Cambridge Reference Sequence" or rCRS
    - 16,569bp long
    - 1000 genome project used with GRCh37 +decoy to create the "g1k" reference
    - This is the default for Golden Helix Sentieon pipeline and VarSeq interpretation

- **NCBI36 (hg18) Included a MT reference NC_001807 in 2006:**
  - Derived from a African (Yoruba) Individual
  - 16,571bp long, differing from the rCRS by 40 variants
  - Removed from GenBank, don't publish with this M!
  - UCSC hg19 includes NC_001807 as "M" and still uses it today!
  - Next VarSeq version drops support for this "hg19" genome

# Variant Interpretation in VSClinical

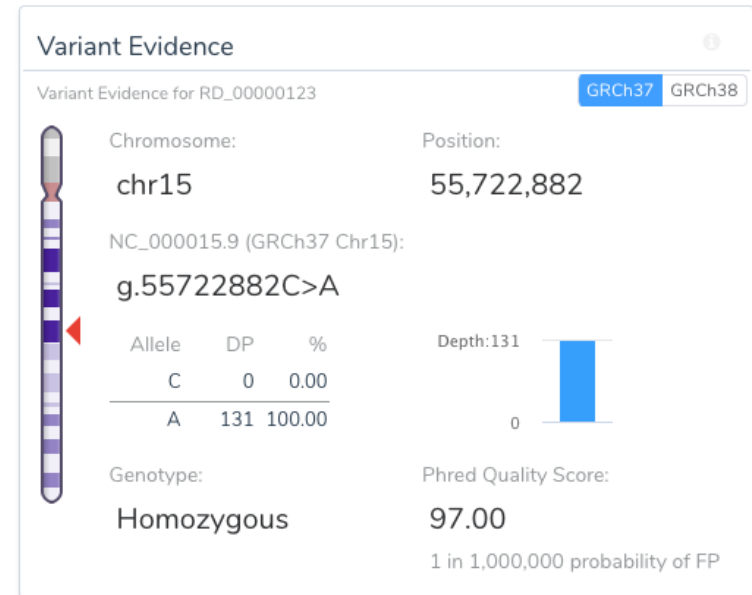- **Evaluate and Classify Variants using ACMG Guidelines:**
  - Focused workflow to evaluate criteria relevant to each variant, resulting in final classification
  - Aggregates annotations from population and clinical resources
  - Customized visualizations and annotation presentations
  - Allows easy look-up and cross reference

- **Save Interpretations into Assessment Catalogs:**
  - New samples have previous classifications brought in
  - See previous interpretations, review and update
  - Can be potted for regional context

- **Use VarSeq's Filter, GenomeBrowse, VSReports:**
  - Customize to lab specific QC, annotation and filtering
  - Genomic context of variant vital to assess
  - VSReports allows custom presentation of VSClinical output



**VSClinical®**

**Variant Evidence**

Variant Evidence for RD_00000123          GRCh37  GRCh38

| Chromosome: | Position: |
|---|---|
| chr15 | 55,722,882 |

NC_000015.9 (GRCh37 Chr15):

g.55722882C>A

| Allele | DP | % |
|---|---|---|
| C | 0 | 0.00 |
| A | 131 | 100.00 |

Depth:131

| Genotype: | Phred Quality Score: |
|---|---|
| Homozygous | 97.00 |

1 in 1,000,000 probability of FP

# GRCh38: Implications for Variant Interpretation

- **Assembly Regions:**
  - Multiple Species Alignment
  - Repeat Regions / Low Complexity Regions
  - Genomic "Super Dups"

- **Genes (and Annotations)**
  - Functional Domains
  - Transcript Counts of Gene Constraint

- **Population Catalogs on GRCh37**
  - dbSNP
  - 1000 Genomes
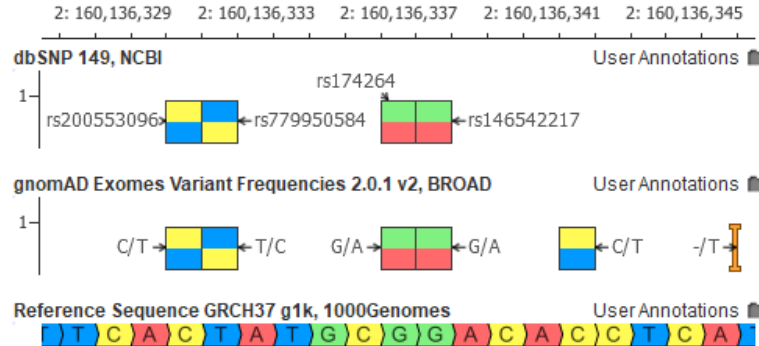  - ExAC / gnomAD Exomes / Genomes

- **Clinical Annotations**
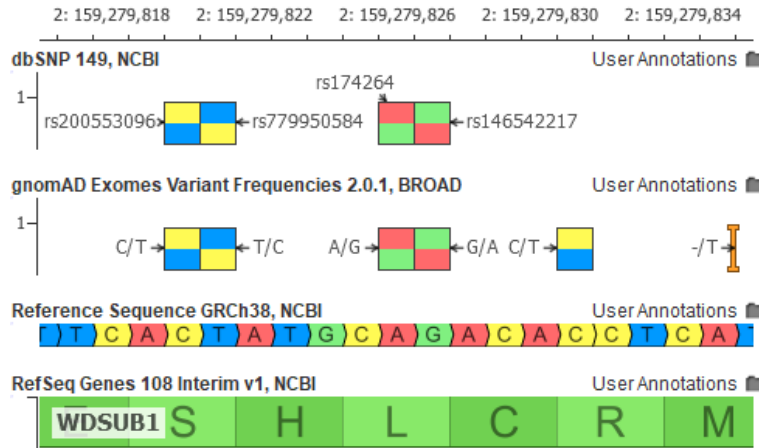  - ClinVar
  - CIVIC
  - OMIM (variants, genes, phenotypes)

- **Functional Annotations / Conservation**
  - CADD
  - SIFT/Polyphen/Missense Badness
  - Conservation scores



Substitution Leu (leucine) → Pro (proline) at 173
Leucine conserved in all vertebrates!

# VarSeq Import LiftOver



Start with GRCh37 VCFs:

LiftOver to GRCh38:

Or the Other Way Around! GRCh38 => GRCh37

[Demo in VarSeq]

# Reasons to Switch to GRCh38

- **Better for alignment**
  - More reads mapped
  - Fewer variants called

- **Better gene representations**
  - Fewer "frame-fixing" introns
  - Some genes fixed/improved

- **Newer annotations are GRCh38**
  - Large consortiums are switching to GRCh38 first:
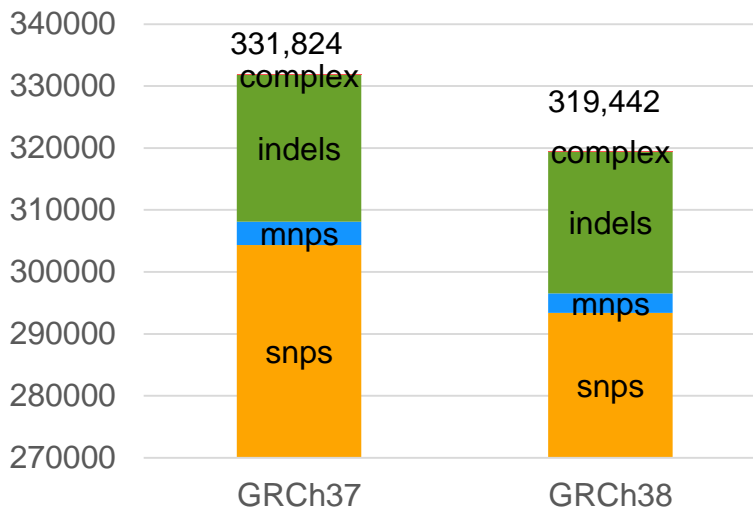    - Cancer: ICGC, COSMIC
    - TopMed (65K WGS)



**Gabe Rudy** @gabeinformatics

Tina Graves: Williams-Beuren Syndrome regions is medically relevant, retiled whole region with valid haplotype. Avail in #GRCh38 #ASHG2013

1:58 PM - 24 Oct 2013



My Exome

# Better Gene Representation

- The human genome does not necessarily contain the mRNA sequence in RefSeq

- "Frame-fixing" intron introduced in alignment of mRNA coding sequence to human reference:

EMG1 on GRCh37:



EMG1 on GRCh38:

# Some Variants are Pure "Reference Artifacts"
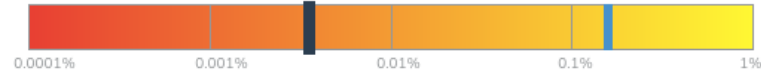
# Some Variants are Pure "Reference Artifacts"

# Considerations for Transitioning your Lab

- **Switching your Secondary Pipeline**

- **Your Genomic Variants Being Saved:**
  - VSClinical Catalog / Assessment Catalogs
  - Catalog of Observed CNVs
  - VSWarehouse Projects (all variants from samples)
  - Target capture annotations
  - Custom in-house annotations

- **Converting Existing Data:**
  - Re-import variants using import Liftover
  - Export/import catalogs using Liftover
  - Convert custom annotations using Liftover

Liftover Using Our Convert Wizard:

# Thank you!

- 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.

# ASHG 2018

- **Booth 408**

- **Live demos and CoLab Sessions**

- **Unveiling our new t-shirt designs**

- **Chance to win some iPads**