# Using VarSeq to Improve Variant Analysis Research
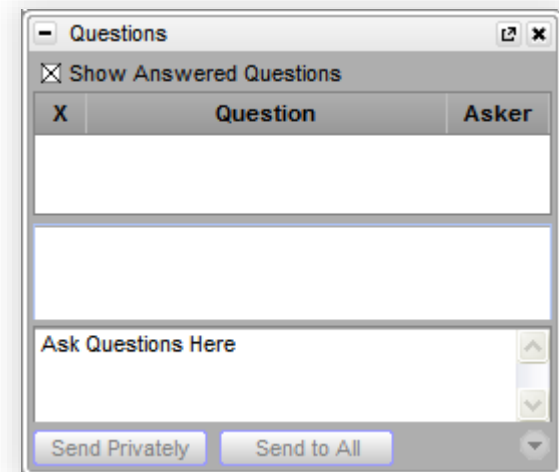
June 10, 2015

G Bryce Christensen
Director of Services

# Questions during the presentation

Use the Questions pane in your GoToWebinar window

# Agenda

**1**    Variant analysis workflows

**2**    What makes a damaging variant?

**3**    QC Considerations

**4**    VarSeq Interactive Demonstration

# What is VarSeq?



Simple · Flexible · Scalable

varSEQ™

- **Variant annotation, filtering and ranking**
- **Repeatable workflows**
- **Rich visualizations with GenomeBrowse integration**
- **Powerful GUI and command-line interfaces**

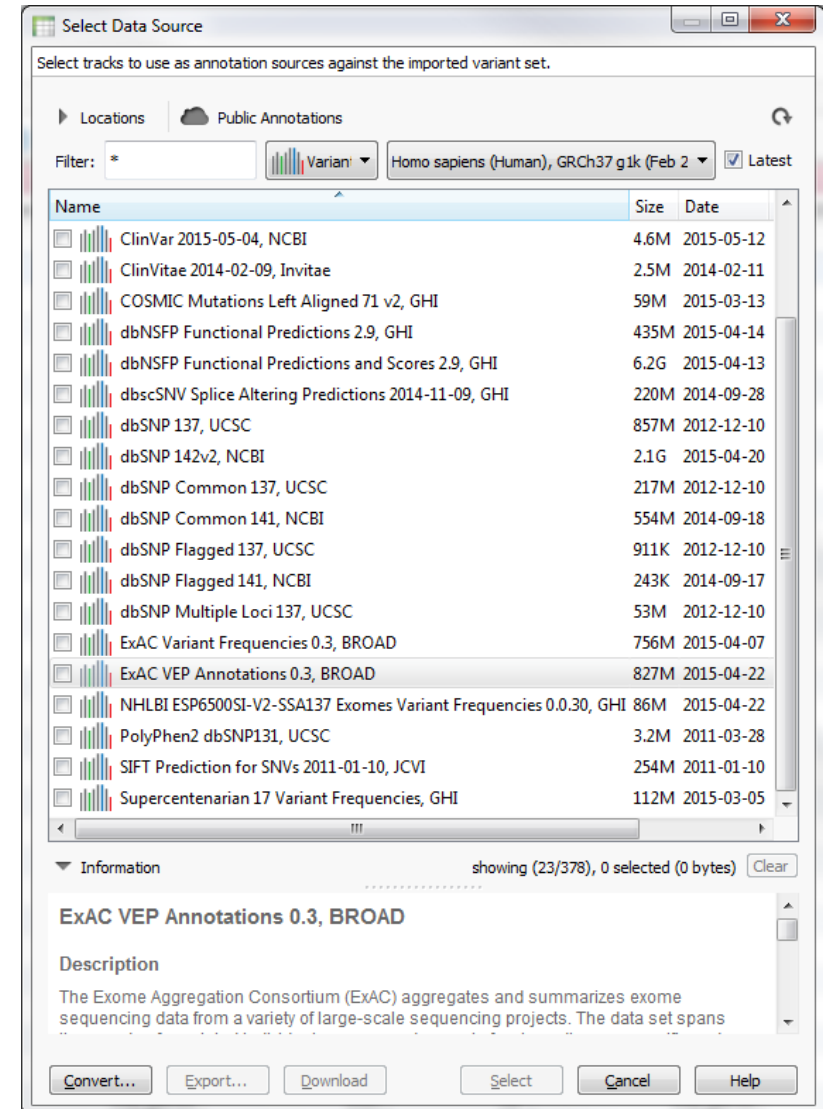# Workflow Development Process in VarSeq

1. Begin from one or many VCF files

2. Annotate variants using public data sources curated by Golden Helix and/or annotate with custom data sources.

3. Run additional computation algorithms
   - Allele counts, genotype zygosity, gene list matching, etc

4. Construct filter chain to identify candidate variants
   - May use combinations of logical operators in filters
   - May have multiple independent filter chains and/or endpoints

5. Process results
   - Gene Ranking with PhoRank
   - Review variant QC
   - Vizualization with GenomeBrowse
   - Commit variants to local database
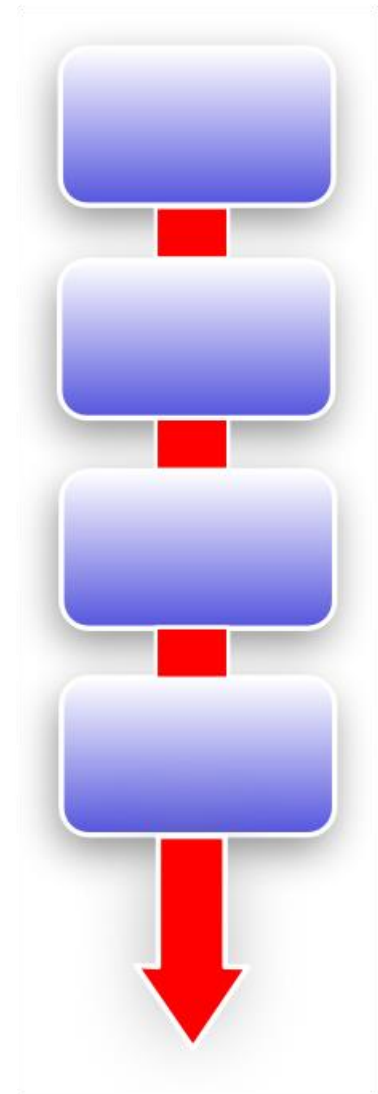   - Etc.

# Annotations are the key

- **Good variant analysis begins with accurate annotations.**

- **Golden Helix invests extensive time and effort in validating and maintaining data sources.**

- **Annotation data sources may be used for either quality control or analytic purposes.**

# Defining Deleteriousness

- **What makes a variant potentially damaging?**

- **Start by defining the search space:**
  - Rare, non-synonymous, homozygous variants?
  - DeNovo mutations in highly conserved genes?
  - Splice-site mutations?
  - Etc.

- **Review annotations for remaining variants to identify causal candidates**

- **Which annotations to use?**
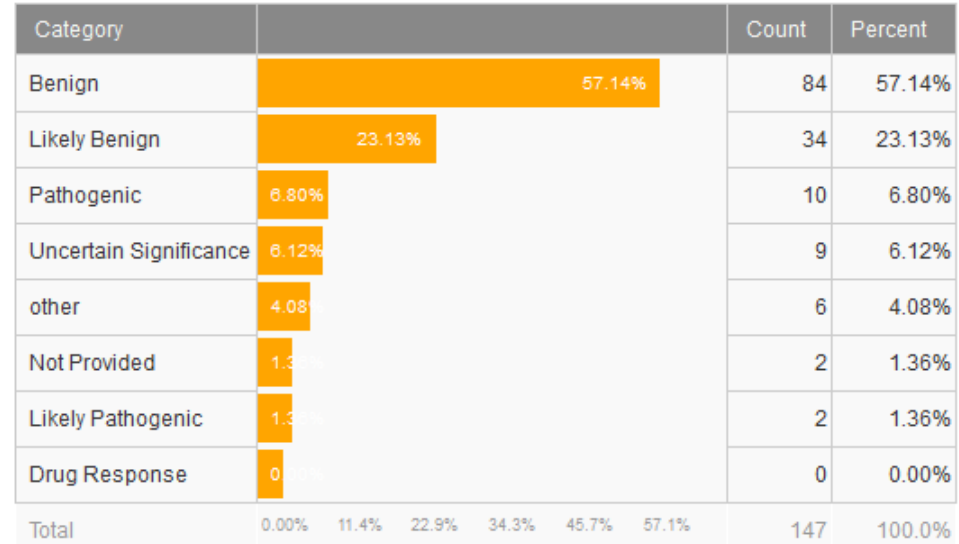
# Variant Classification

| Category | | Count | Percent |
|---|---|---|---|
| 3_prime_UTR_variant | 47.81% | 10282 | 47.81% |
| missense_variant | 19.29% | 4149 | 19.29% |
| intergenic_variant | 11.52% | 2477 | 11.52% |
| synonymous_variant | 11.35% | 2442 | 11.35% |
| 5_prime_UTR_variant | 4.44% | 955 | 4.44% |
| splice_region_variant | 1.78% | 383 | 1.78% |
| intron_variant | 1.56% | 336 | 1.56% |
| frameshift_variant | 0.53% | 115 | 0.53% |
| stop_gained | 0.53% | 113 | 0.53% |
| splice_donor_variant | 0.40% | 87 | 0.40% |
| disruptive_inframe_deletion | 0.22% | 47 | 0.22% |
| splice_acceptor_variant | 0.15% | 33 | 0.15% |
| 5_prime_UTR_premature_start_codon_gain_variant | 0.13% | 28 | 0.13% |
| disruptive_inframe_insertion | 0.11% | 23 | 0.11% |
| inframe_deletion | 0.07% | 14 | 0.07% |
| inframe_insertion | 0.05% | 11 | 0.05% |
| stop_retained_variant | 0.03% | 7 | 0.03% |
| stop_lost | 0.01% | 3 | 0.01% |
| initiator_codon_variant | 0.01% | 2 | 0.01% |
| non_coding_exon_variant | 0.00% | 0 | 0.00% |
| Total | 0.00% 9.56% 19.1% 28.7% 38.2% 47.8% | 21507 | 100.0% |

- **VarSeq classifies variants into 20+ different categories**

- **The categories are further grouped as:**
  - Loss of Function
  - Missense
  - Other

- **Choice of gene transcript reference**
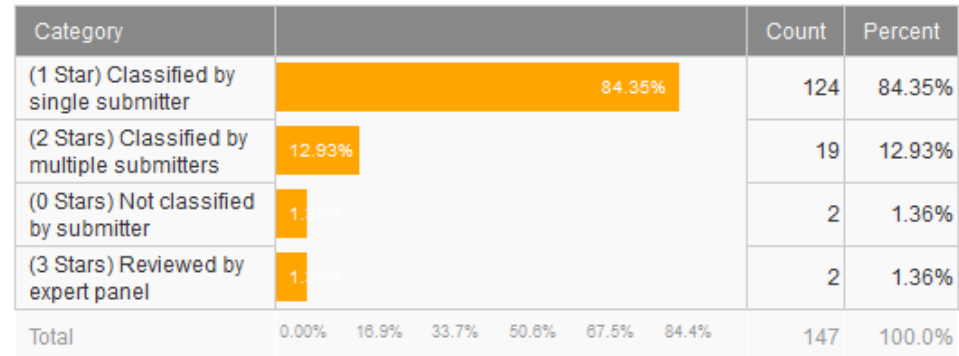  - RefSeq
  - Ensembl
  - Others

# ClinVar

- **ClinVar is a public archive of variants evaluated for potential causal relationships to diseases**

- **Submissions from many sources, including major clinical laboratories**

- **Over 100k records**

- **Updated monthly**

Category Counts (123 Records from a Dataset Total of 98,655)

| Category | | Count | Percent |
|---|---|---|---|
| Benign | 57.14% | 84 | 57.14% |
| Likely Benign | 23.13% | 34 | 23.13% |
| Pathogenic | 6.80% | 10 | 6.80% |
| Uncertain Significance | 6.12% | 9 | 6.12% |
| other | 4.08% | 6 | 4.08% |
| Not Provided | 1.36% | 2 | 1.36% |
| Likely Pathogenic | 1.36% | 2 | 1.36% |
| Drug Response | 0 | 0 | 0.00% |
| Total | 0.00%  11.4%  22.9%  34.3%  45.7%  57.1% | 147 | 100.0% |

Category Counts (123 Records from a Dataset Total of 98,655)

| Category | | Count | Percent |
|---|---|---|---|
| (1 Star) Classified by single submitter | 84.35% | 124 | 84.35% |
| (2 Stars) Classified by multiple submitters | 12.93% | 19 | 12.93% |
| (0 Stars) Not classified by submitter | 1.36% | 2 | 1.36% |
| (3 Stars) Reviewed by expert panel | 1.36% | 2 | 1.36% |
| Total | 0.00%  16.9%  33.7%  50.6%  67.5%  84.4% | 147 | 100.0% |

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Functional Predictions

- **Functional predictions use algorithms to determine the expected consequence of variants (or the resulting amino acid substitutions).**

- **dbNSFP**
  - The Database for NonSynonymous Functional Predictions (dbNSFP) is a free tool developed by Dr. Xiaoming Liu.
  - Catalogs pre-computed conservation and functional prediction scores for all possible missense SNVs in the genome
  - Methods include SIFT, PolyPhen-2, MutationTaster, MutationAssessor, FATHMM, more

- **dbscSNV**
  - Companion to dbNSFP that scores variants in splice consensus regions
  - Variants in these regions may disrupt normal gene expression and/or function

- **dbNSFP and dbscSNV are both accessible in VarSeq**

# Variant/Gene Ranking

- **PhoRank algorithm in VarSeq uses HPO and GO terminology to score relationships between genes and phenotypes**

- **Very useful to prioritize a long list of variants for individual review**

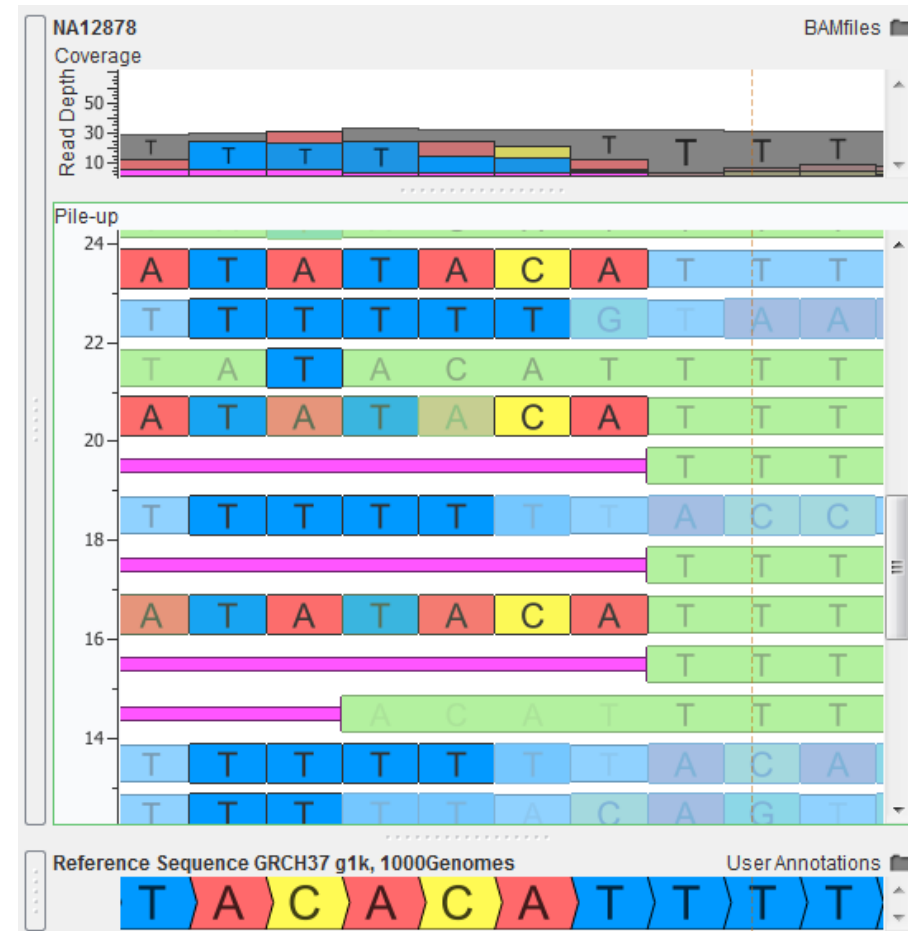- **Based on PHEVOR method.**

# QC Considerations

- **Variant QC**

- **Rare variants deserve special attention**

- **VCF/BAM Data:**
  - Depth - DP
  - Quality - GQ
  - Strand bias
  - Etc.

- **Public Annotations:**
  - "Mappability"
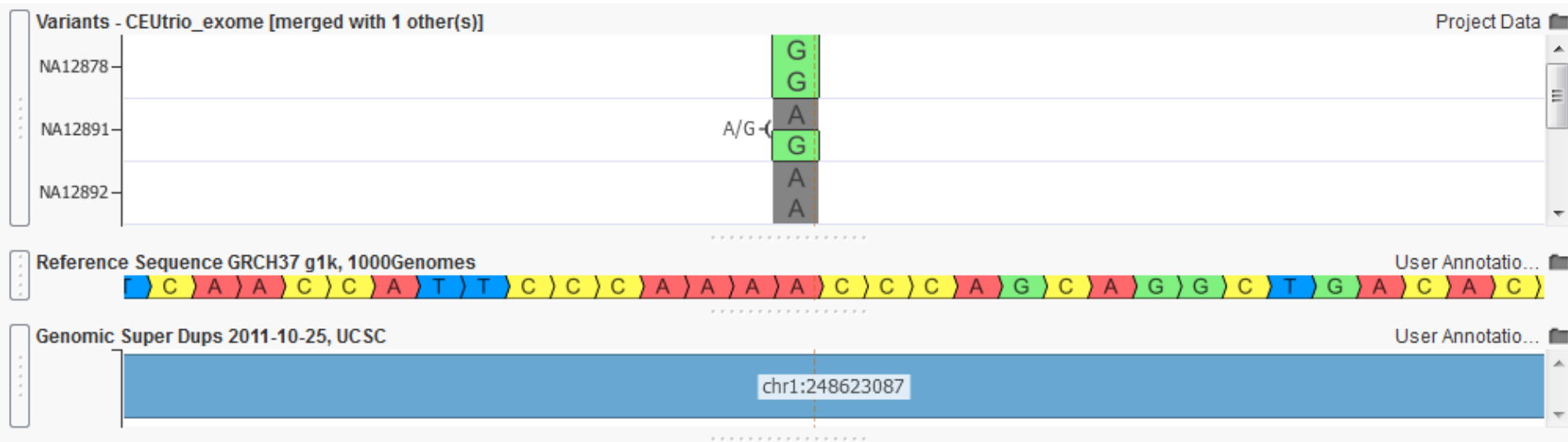
# Mappability Annotations

- **The human reference genome has assembly gaps and other "difficult" regions**

- **NGS technology sequences short DNA fragments which are the aligned to the reference genome**
  - Most sequences are aligned correctly
  - Some sequences can't be aligned uniquely
  - Some sequences may be incorrectly aligned

- **Luckily, we can predict many of the trouble spots**
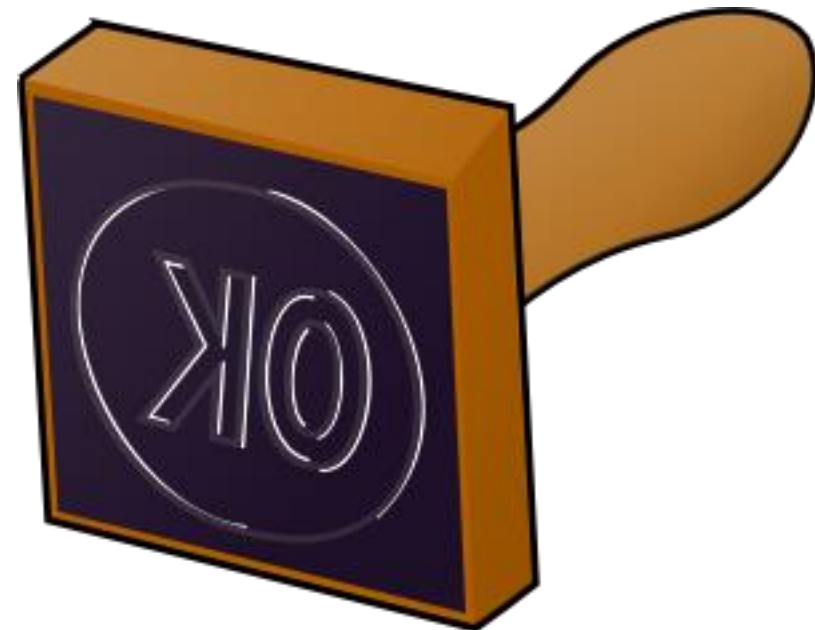
# Segmental Duplications

- **Segmental duplications are a common confounder**

- **UCSC "Genomic Super Dups" annotation available through VarSeq**

- **Recent Example (below):**
  - Apparent UPD feature in family trio was determined to be an artifact of seg. duplication
  - Large chromosome segment duplicated elsewhere with >98% similarity
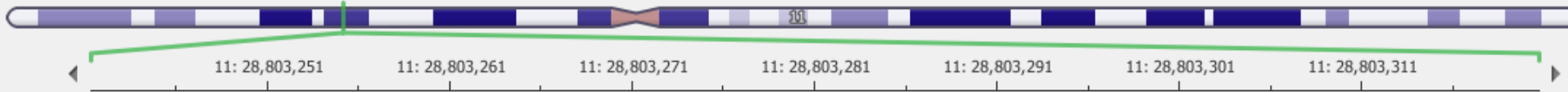
# Emerging Standards

- **Several organizations working on best practices guidelines for genome mappability**
  - 1000 Genomes Project
  - Genome in a Bottle Consortium
  - Global Alliance for Genomics and Health (GA4GH)
  - National Institute of Standards and Technology

- **Downloadable annotations available for many types of features:**
  - Mappability by read length
  - High G-C content regions
  - Low complexity
  - Segmental duplications
  - Etc.

# Example: 1kG Low Complexity Regions

# VarSeq Demonstration Data

- **Exome sequencing of five individuals from family with familial cardiac conduction disease (CCD)**
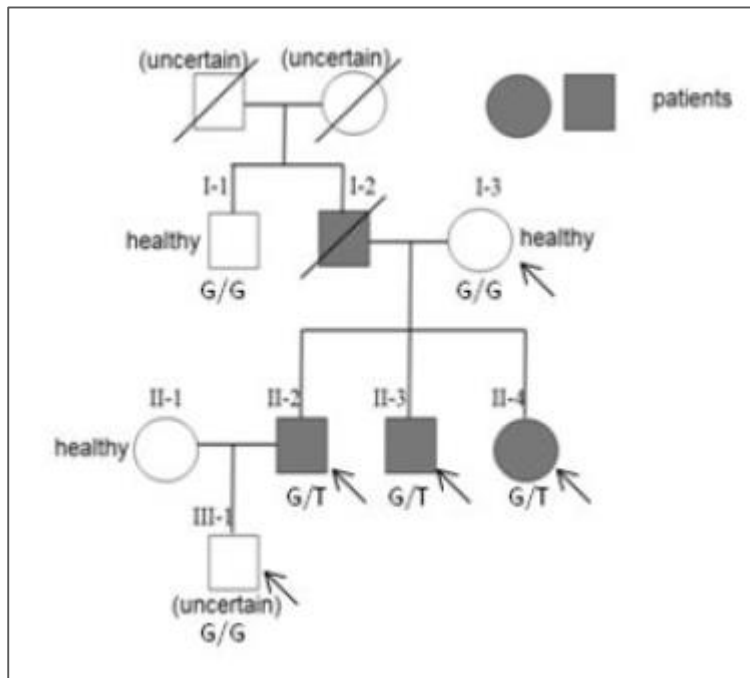
- **Raw sequence data obtained from SRA**

PLOS | ONE

## Whole-Exome Sequencing to Identify a Novel LMNA Gene Mutation Associated with Inherited Cardiac Conduction Disease

Chun-Chi Lai[1,©], Yung-Hsin Yeh[2,©], Wen-Ping Hsieh[3], Chi-Tai Kuo[2], Wen-Ching Wang[4,5], Chia-Han Chu[5], Chiu-Lien Hung[4,5], Chia-Yang Cheng[1,5], Hsin-Yi Tsai[2], Jia-Lin Lee[4], Chuan-Yi Tang[1], Lung-An Hsu[2*]

1 Department of Computer Science, National Tsing Hua University, Hsinchu, Taiwan, 2 First Cardiovascular Division, Chang Gung Memorial Hospital, Chang Gung University College of Medicine, Tao-Yuan, Taiwan, 3 Institute of Statistics, National Tsing Hua University, Hsinchu, Taiwan, 4 Institute of Molecular and Cellular Biology and Department of Life Sciences, National Tsing-Hua University, Hsinchu, Taiwan, 5 Biomedical Science and Engineering Center, National Tsing Hua University, Hsinchu, Taiwan

GOLDEN HELIX
Accelerating the Quest for Significance™

- **Male-to-male transmission makes X-linked model unlikely**

- **May follow dominant or recessive transmission**

- **Inherited forms of CCD are rare**

- **Family has East Asian ancestry**

**[Demonstration]**

# Why VarSeq?

**Flexible**

**Simple**

**Scalable**

varSEQ™

- **Variant annotation, filtering and ranking**

- **Exploratory analysis**

- **Powerful GUI with immediate feedback**

- **Rich visualizations with GenomeBrowse integration**

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Questions or more info:

- Email
  info@goldenhelix.com

- Request an evaluation of the software at
  www.goldenhelix.com

# Questions?

Use the Questions pane in your GoToWebinar window