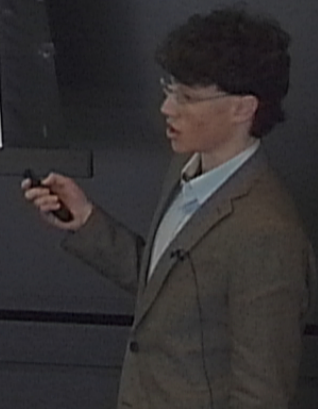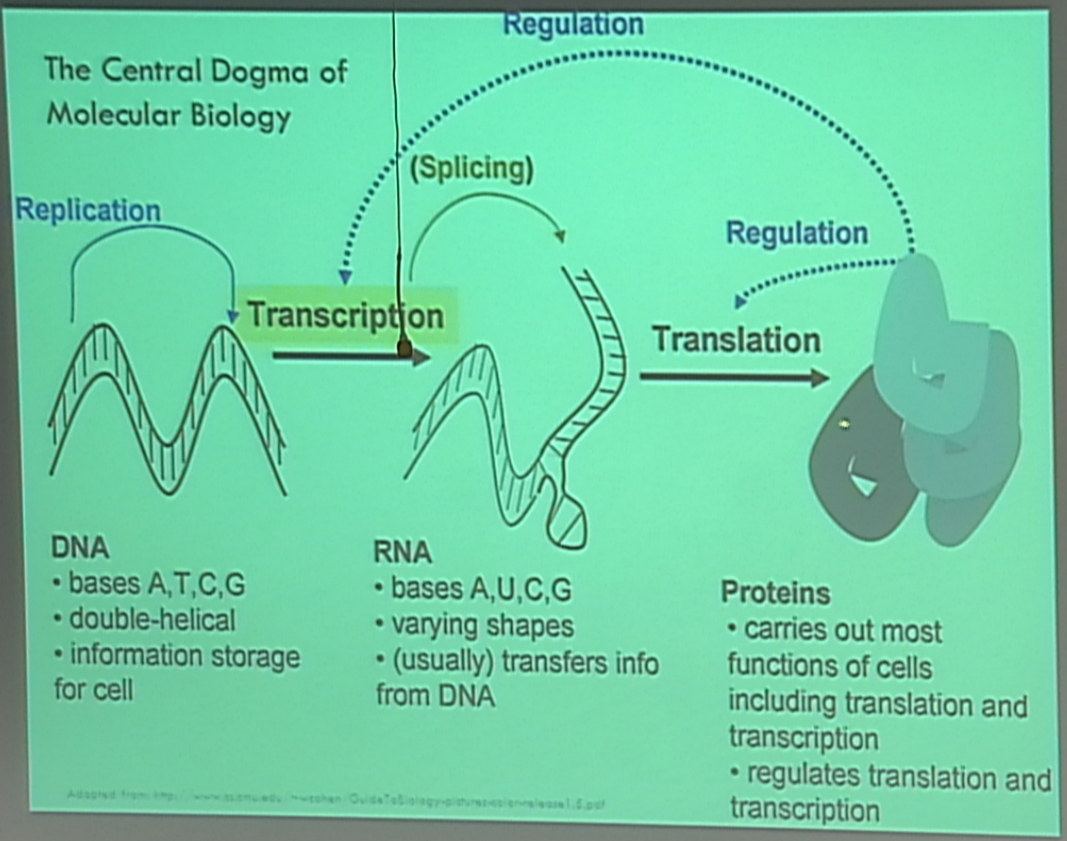Title: Validation of predicted mRNA splicing mutations using high-throughput transcriptome data

Date: May 07, 2014  04:35 PM
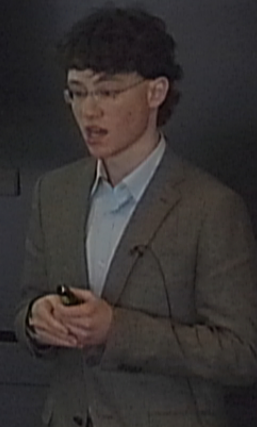
URL: http://pirsa.org/14050059

Abstract: <span>This work has been published:Viner C Dorman SN Shirley BC and Rogan PK (2014)Validation of predicted mRNA splicing mutations using high-throughput transcriptome data [v1; ref status: indexedhttp://f1000r.es/2no]F1000Research20143:8 (doi:10.12688/f1000research.3-8.v1)Additionally this work has been accepted for a highlights presentation at the upcoming Great Lakes Bioinformatics Conference (GLBIO) in Cincinnati Ohio and it was recently presented as a poster at London Health Research Day (LHRD).Abstract:Interpretation of variants present in complete genomes or exomes reveals numerous sequence changes only a fraction of which are likely to be pathogenic. Variants predicted to alter mRNA splicing in particular can be validated by manual inspection of transcriptome sequencing data however this approach is intractable for large datasets. We show that abnormal mRNA splicing patterns are characterized by reads demonstrating either exon skipping cryptic splice site use and high levels of intron inclusion or combinations of these properties. This paper presents Veridical an in silico method for the automatic validation of DNA sequencing variants that alter mRNA splicing. Veridical leverages large numbers of control samples (that lack a putative mutation) applying z-tests to Yeo-Johnson transformed data to normalize read counts of abnormal RNA species in mutant versus non-mutant tissues. With the transformed data the null hypothesis based mainly on either counts of intronic or junctional reads can be rejected for true splicing mutations using conventional parametric statistical methods. </span>
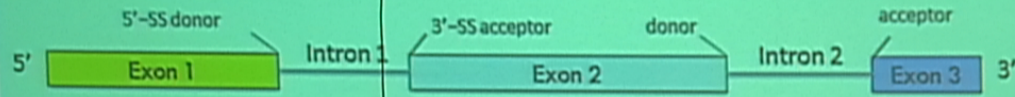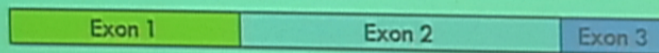
Introduction

- Genome sequencing data from The Cancer Genome Atlas (TCGA)

  - Define mutated and stable genes

  - Enrichment analysis: dysregulated metabolic pathways in solid tumors

- Failure of currently available methods to correctly categorize many gene variants of unknown significance

  - Substantial potential to be pathogenic

- Mutations in coding and non-coding regions (typically near exon/intron boundaries)
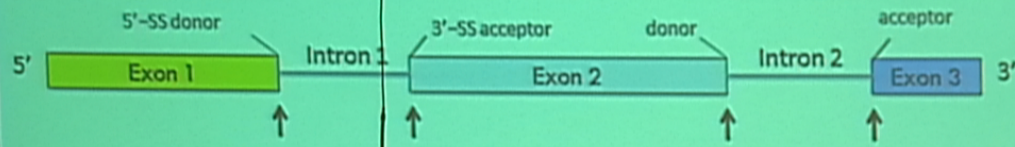
  - Affect mRNA processing can result in aberrant splicing

Splicing Variants

Wild Type:

Introduction

- Shannon Pipeline

  - Implements an algorithm for high-throughput detection and interpretation of these mRNA splicing mutations, using information theory

- Putative variants require empirical confirmation

  - Translate predictions to clinically relevant insights

- Currently, by visual inspection of RNA-Seq for abnormalities

  - Intractable when scaled

- Mutations in DNA corroborated by RNA-Seq from the same patient

Objectives

- To develop a method to automatically validate putative DNA sequencing variants that alter mRNA splicing across multiple patient samples, by using corresponding RNA sequencing data

- To derive novel biological insights from breast carcinoma data via a more in-depth analysis of splicing mutations
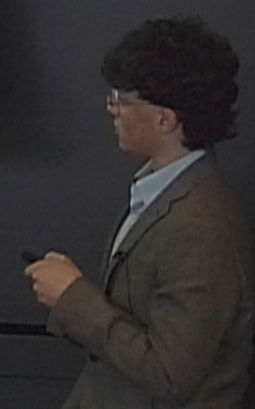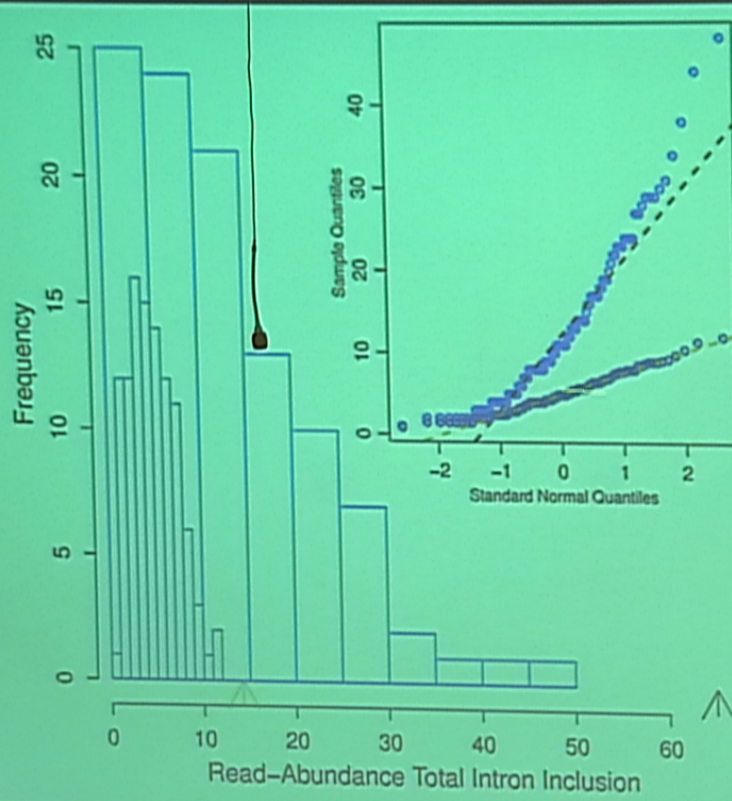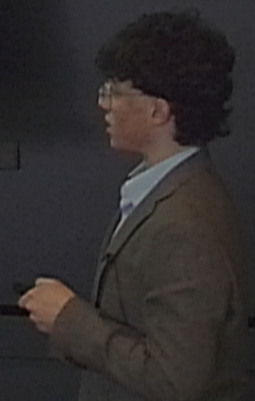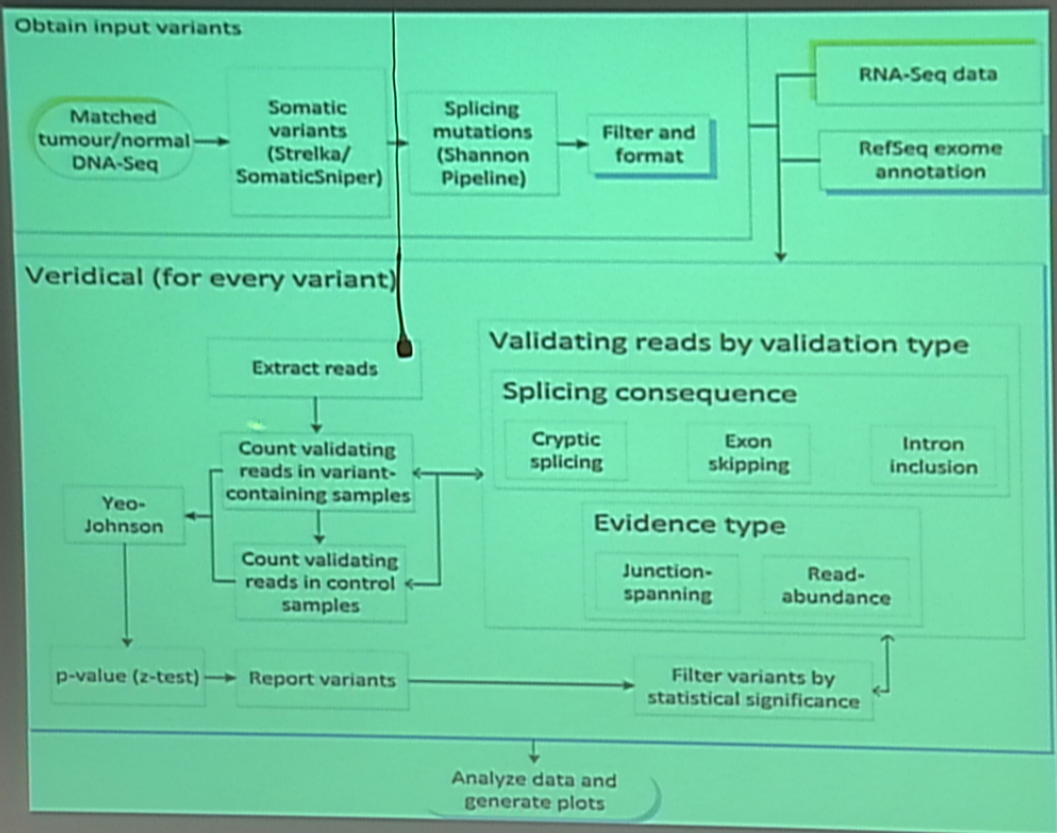
Evidence Types

**Veridical**

- Hypothesis-driven
- Statistically validates mutations throughout entire exome using RNA sequencing data
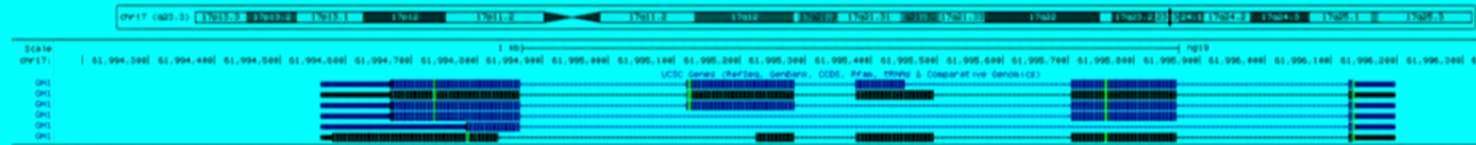- Can be used to validate mutations in any individual/disease

Exon Skipping (Junction-Spanning)

Cryptic Splicing (Junction-Spanning)

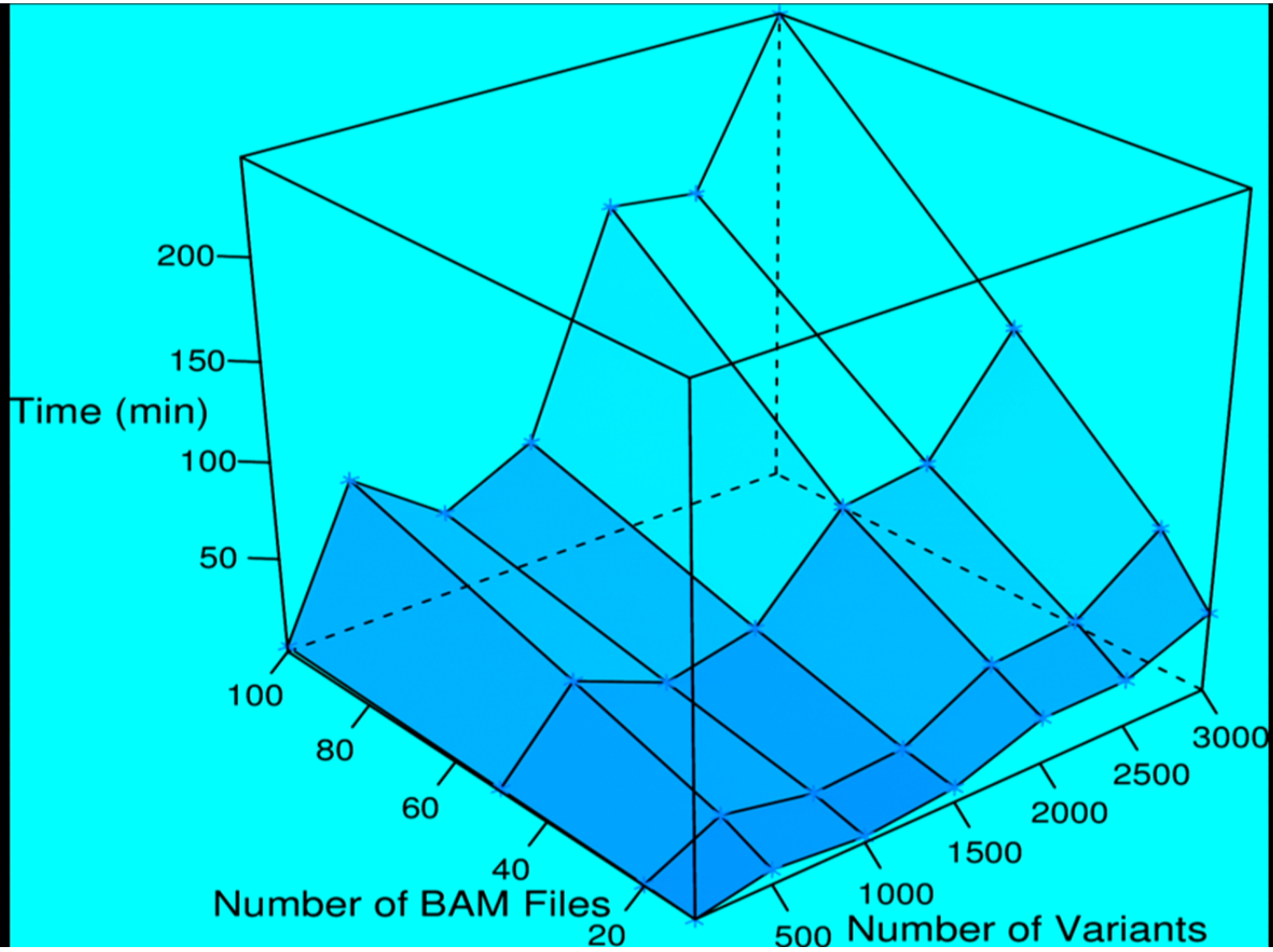**COMPARED SPLICING CONSEQUENCES WITH > 500 TUMOUR/NORMAL CONTROLS
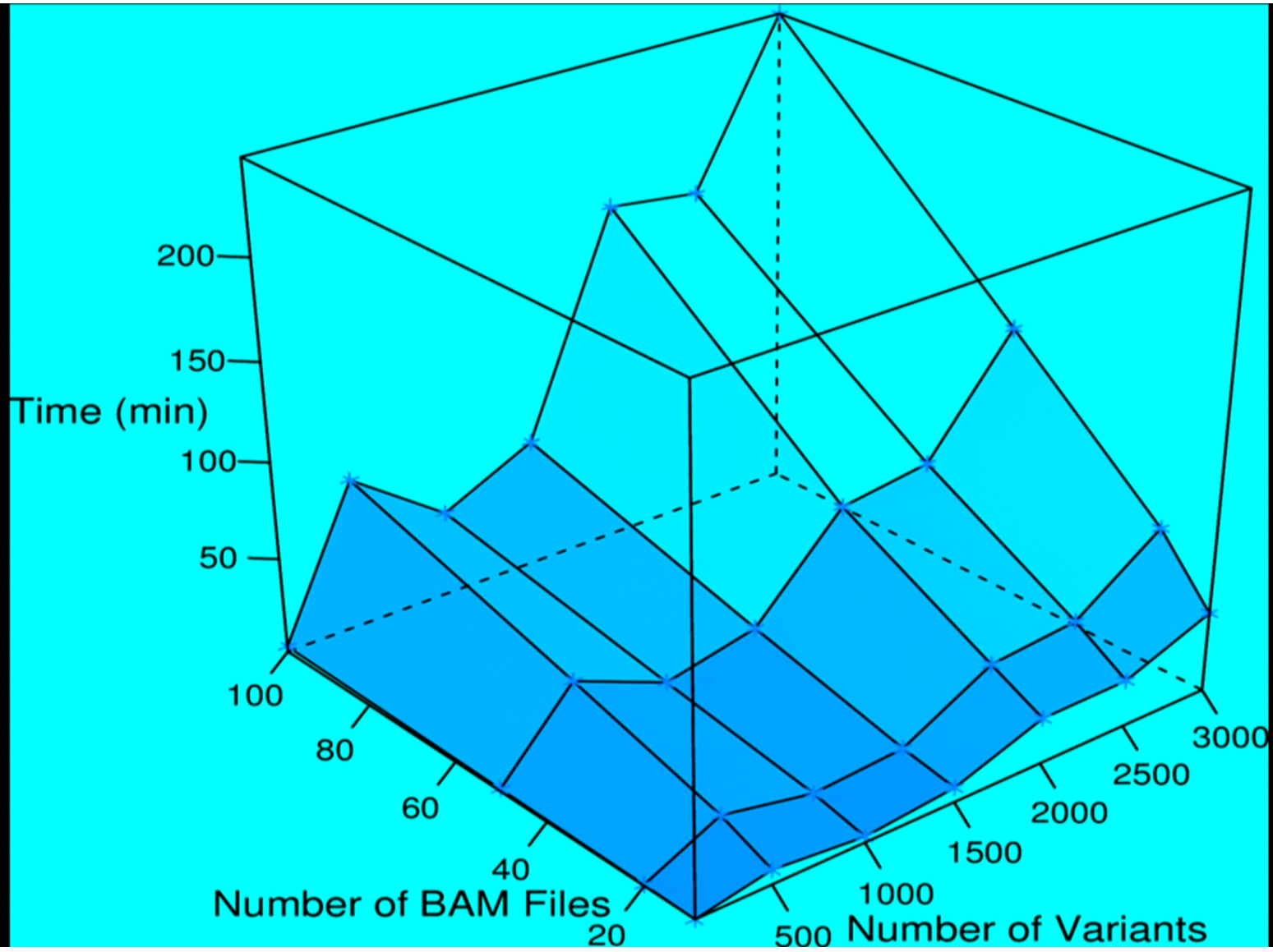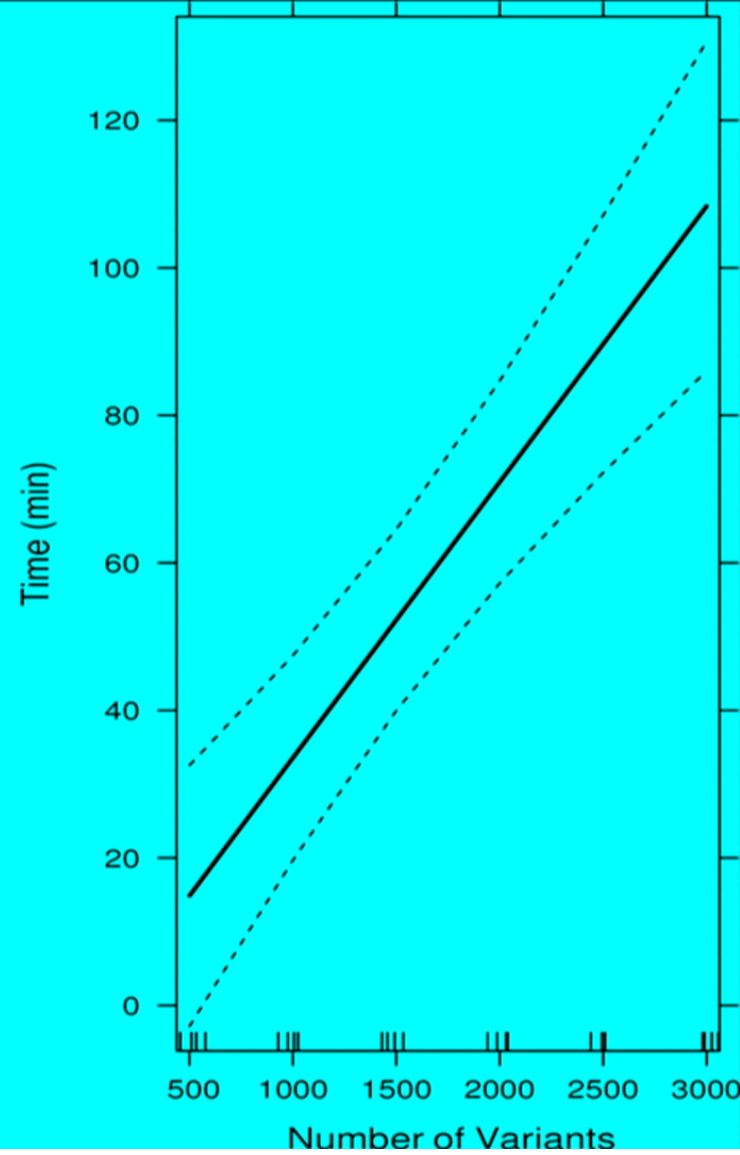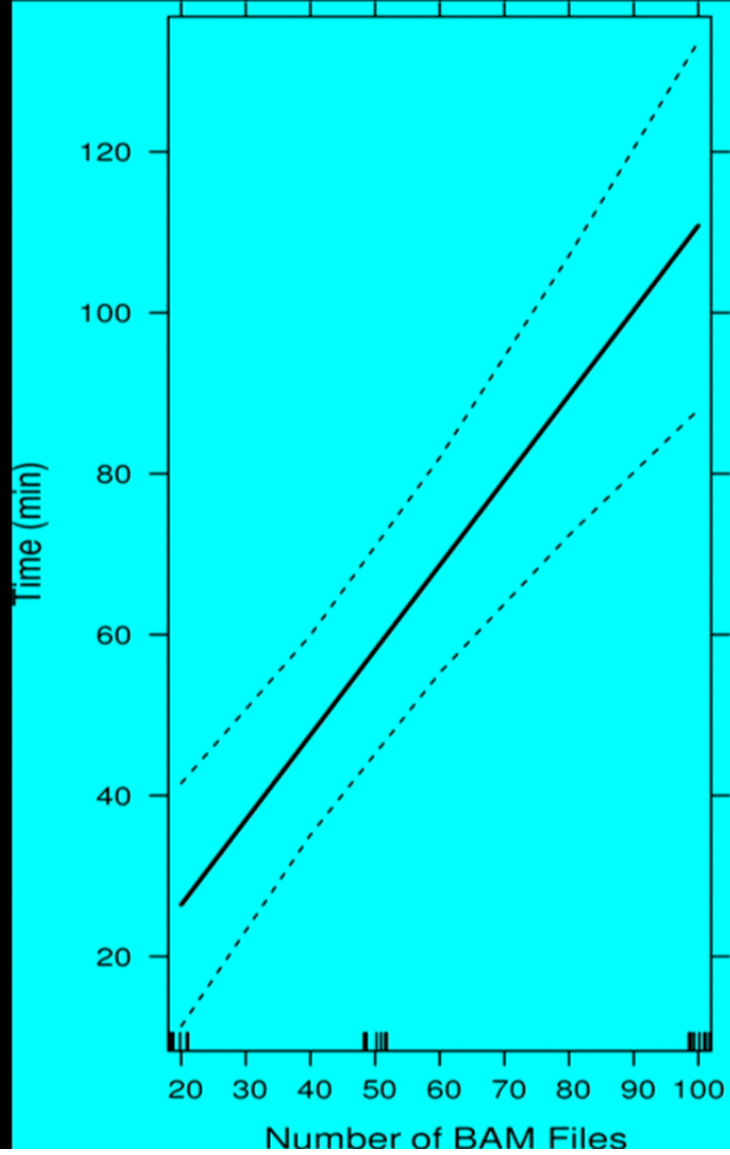
# What we are *not* doing:



- **Alternative splicing is a different problem**

  - Many programs address this

    - Multivariate Analysis of Transcript Splicing (MATS)

      - Detects alternative splicing events via MCMC simulation sampling to compute p-values and FDRs

- **Deriving a set of putative variants is a different problem**

  - Variant Annotation, Analysis and Search Tool (VAAST)

    - Uses a likelihood approach to rank variants by pathogenicity

    - Does not conduct any detailed splicing mutation analyses

# What we *are* doing:

- **Hypothesis testing:** A predicted mutation affects mRNA splicing

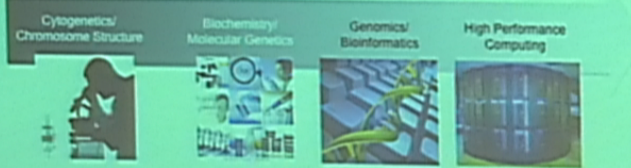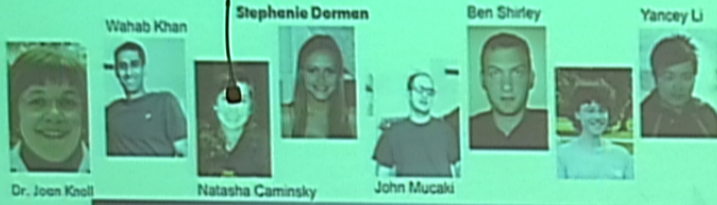  - using variant predictions and an existing exome annotation

Future Work

- Veridical is trivially parallelizable
  - Can use data parallelization at two levels: BAM files and variants
- Improved read-abundance validation for cryptic variants
- Integration of copy number data to inform read count expectations
- Address nonsense-mediated mRNA decay
- Better alignment algorithms may yield better read recognition, particularly with respect to cryptic splice junctions
- Further mining of generated breast carcinoma data