# KNOWENG

### **BIG DATA TO KNOWLEDGE** CENTER OF EXCELLENCE

### Analysis of Genomic Data

Researcher has experimental genomic profiles of multiple samples stored in a **spreadsheet**. Examples:

- genes x tumor samples somatic mutation matrix
- genes x cell lines basal transcription from RNA-seq

### **Input Data**







The user wants to apply machine learning and graph mining analysis to their spreadsheet in order to gain understanding of their experimental results.



### **Gene Prioritization**



**Goal**: Identify genes whose mRNA expression levels explain the variation of drug sensitivity in different cell lines/individuals

**Knowledge Network Hypothesis**: Interesting genes for follow up analysis may show correlations of their neighbors' expression values with the drug sensitivity measurements

**Approach**: Robust network ranking of network smoothed expression profiles. Based on:

**Knowledge-Guided Prioritization of Genes Determinant** of Drug Resistance. Emad A, Cairns J, Kalari K, Wang L, Sinha S.

# **KnowEnG Cloud-based Scalable Analytics Suite**

C. Blatti<sup>1</sup>, M. Berry<sup>2</sup>, L. Gatzke<sup>2</sup>, A. Emad<sup>1</sup>, N. Sobh<sup>1</sup>, C. Bushell<sup>1,2</sup>, S. Sinha<sup>1,3</sup> <sup>1</sup>KnowEnG Center, <sup>2</sup>Applied Research Institute, <sup>3</sup>Department of Computer Science, University of Illinois, Urbana-Champaign, IL



### Analysis Workflows in KnowEnG

### **Gene Set Characterization**

![](_page_0_Picture_23.jpeg)

**Goal**: Find associations between the researchers novel gene sets and previously known biological annotations in order to provide understanding and hypothesis

Knowledge Network Hypothesis: Extends associations to poorly/incompletely annotated domains by integrating multiple relationship types

**Approach**: Network ranking of public gene sets for specificity to the query. Based on:

Characterizing gene sets using discriminative random walks with restart on heterogeneous biological networks. <u>Blatti C<sup>1</sup>, Sinha S<sup>2</sup>.</u>

Bioinformatics

![](_page_0_Picture_29.jpeg)

**Goal**: Identify cancer subtypes from somatic mutation genomic data that are predictive of clinical outcomes

Knowledge Network Hypothesis: Within a subtype, two tumors may not share the same somatic mutation, but may have somatic mutations affecting the same pathways and protein interaction networks

**Approach**: Robust clustering of network smoothed profiles. Based on:

Network-based stratification of tumor mutations free<sup>1</sup>, John P Shen<sup>2</sup>, Hannah Carter<sup>2</sup>, Andrew Gross<sup>3</sup> & Trey Ideker<sup>1-3</sup> nature

Annotations - Annotations Characteristics t Achilles 🥔 Enrich Outcomes

ALLEN BRAIN ATLAS

### On Scalable Cloud Infrastructure

KnowEnG uses distributed systems, algorithms, and workflows that can scale to handle large-scale bioinformatics analysis on the increasing size and diversity of user data and community knowledge.

We deploy our analysis pipelines: As a series of **Docker containers** for each analysis task Whose execution is orchestrated by the **Chronos** job scheduling

- framework
- On a compute cluster managed by **Apache Mesos**
- That syncs user and community data through AWS S3

## **Subtype Stratification**

### **New Workflows**

- *Text Mining* Find genes most specifically related to different disease terminology
- **Phenotype Prediction** Create model that predicts phenotypic outcomes from genomic data
- **Gene Regulatory Networks** Model interactions between transcripts and transcription factors

### **Integration with Other Clouds**

- Import user spreadsheets directly from other cloudbased datasets like TCGA, LINCS Package analysis workflows for other cloud computation engines like Seven Bridges Cancer
- **Genomics** Cloud

### Interface Improvements

- Automatic mapping of gene/protein identifiers Interactive network-based visualizations to view user and public data
- Save and export analysis results (including support for Bagit data format and publication in BDDS Data Repository)

## **Acknowledgement:**

"This research was supported by grant 1U54GM114838 awarded by NIGMS through funds provided by the trans-NIH Big Data to Knowledge (BD2K) initiative (www.bd2k.nih.gov)."

![](_page_0_Picture_62.jpeg)

![](_page_0_Picture_65.jpeg)

### Upcoming Features

![](_page_0_Picture_69.jpeg)