# STANDING ON THE SHOULDERS OF GIANTS

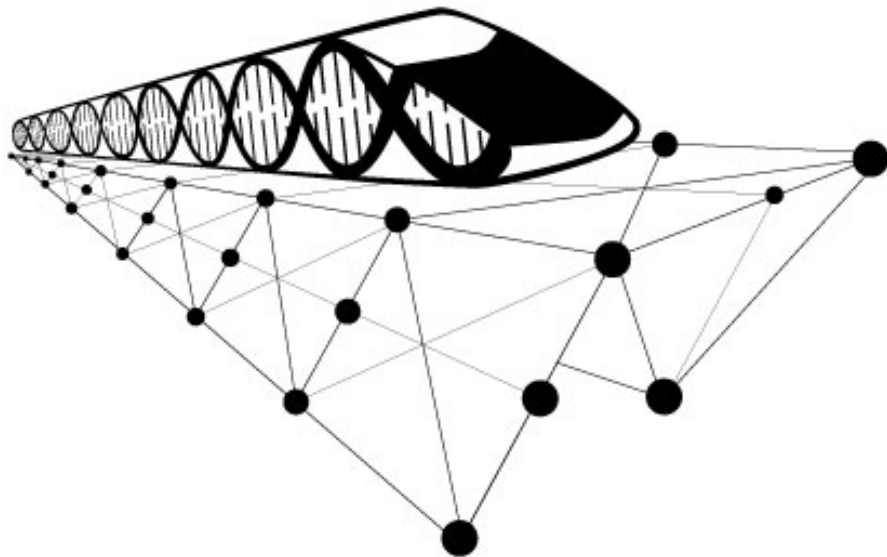# HOW MASSIVE KNOWLEDGE-BASES ARE TRANSFORMING DATA ANALYTICS IN BIOLOGY

SAURABH SINHA

PROFESSOR OF COMPUTER SCIENCE

AND THE CARL R. WOESE INSTITUTE FOR GENOMIC BIOLOGY

CO-DIRECTOR & RESEACH PI, NIH BD2K CENTER OF EXCELLENCE, UIUC & MAYO CLINIC.

IT'S, UH, "GENOMICAL"

# IN THE BEGINNING THERE WAS THE GENOME

# "MACHINE" CODE IS NOT VERY USEFUL

# SO PEOPLE STARTED PROFILING THE CODE

# ALL ROADS LEAD TO A SPREADSHEET

Conditions

| | Tissue 1 | Tissue 2 | ... | Tissue 400 |
|---|---|---|---|---|
| Gene 1 | 20 | 5 | 23 | 37 |
| Gene 2 | 10 | 17 | 201 | 29 |
| ... | 100 | 102 | 99 | 84 |
| Gene 20000 | 20 | 45 | 74 | 62 |

Genes

Conditions

Genes

| | Patient 1 | Patient 2 | ... | Patient 400 |
|---|---|---|---|---|
| Gene 1 | 20 | 5 | 23 | 37 |
| Gene 2 | 10 | 17 | 201 | 29 |
| ... | 100 | 102 | 99 | 84 |
| Gene 20000 | 20 | 45 | 74 | 62 |

# SPREADSHEET ANALYTICS (A.K.A. BIOINFORMATICS)

**Regression**

**Classification**

|  | Patient 1 | Patient 2 | ... | Patient 400 |
|---|---|---|---|---|
| Gene 1 | 20 | 5 | 23 | 37 |
| Gene 2 | 10 | 17 | 201 | 29 |
| ... | 100 | 102 | 99 | 84 |
| Gene 20000 | 20 | 45 | 74 | 62 |

**Clustering**

"A GOOD DECISION IS BASED ON KNOWLEDGE AND NOT ON NUMBERS"
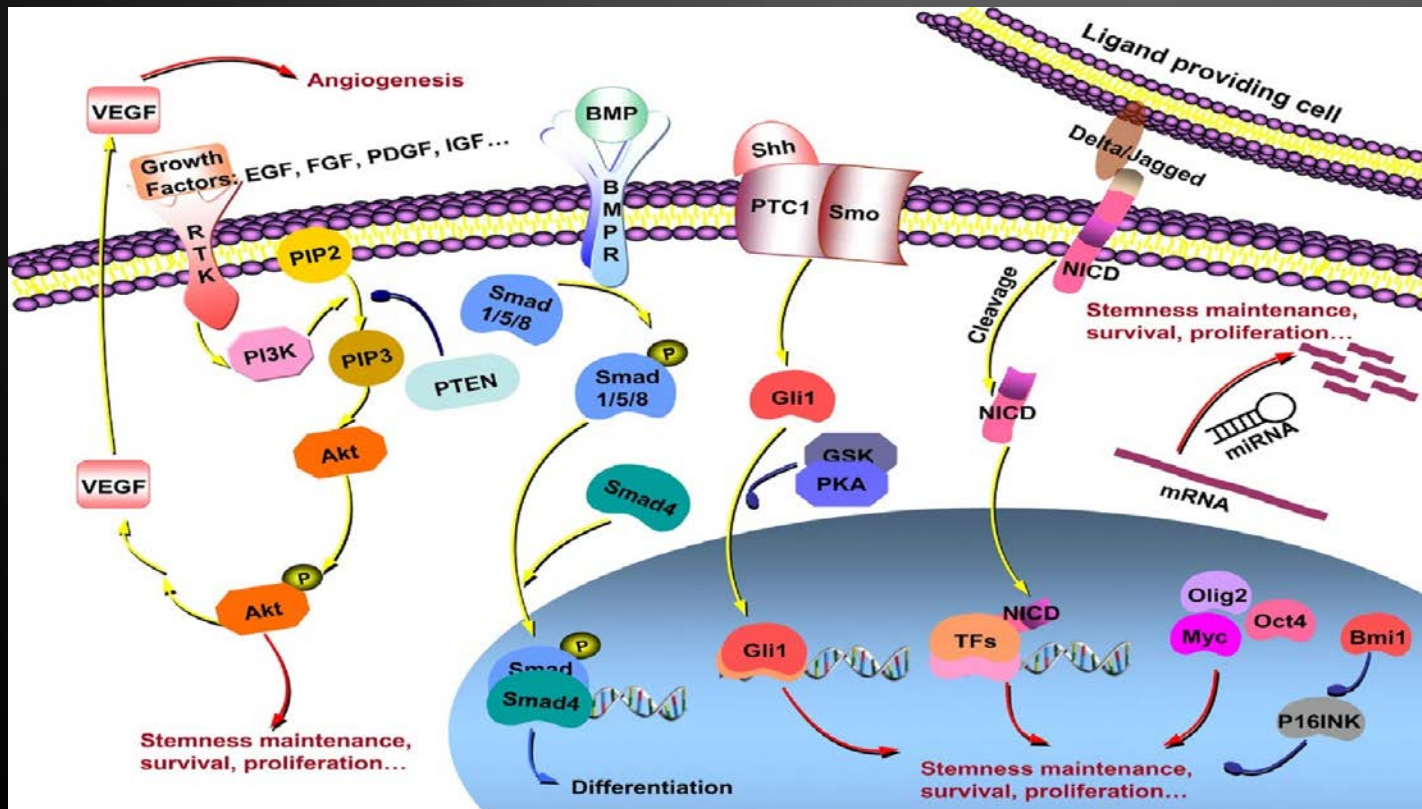- PLATO

# THE GWAS STORY

Disease?

I1:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
I2:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
I3:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
I4:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
I5:    AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT    +
I6:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
I7:    AACGAGCTAGCGATCGATCGACTACGACTACGAGGT    –
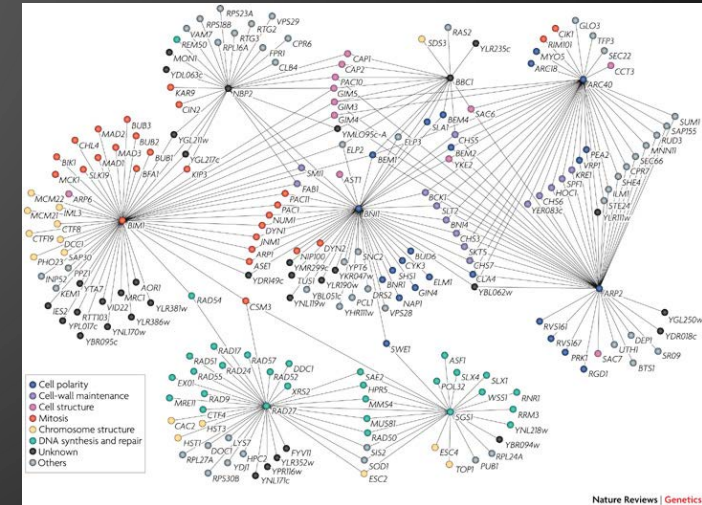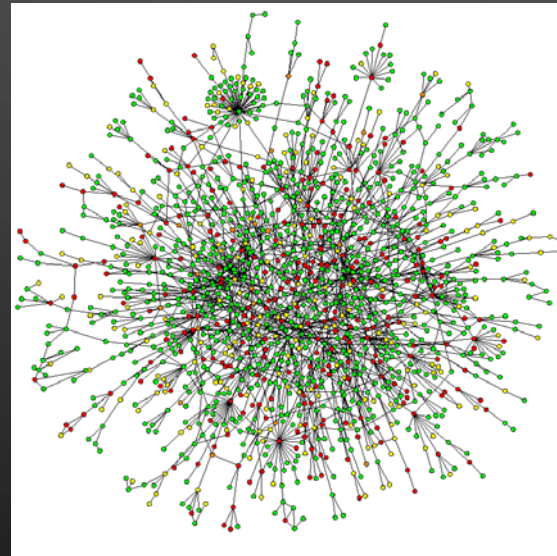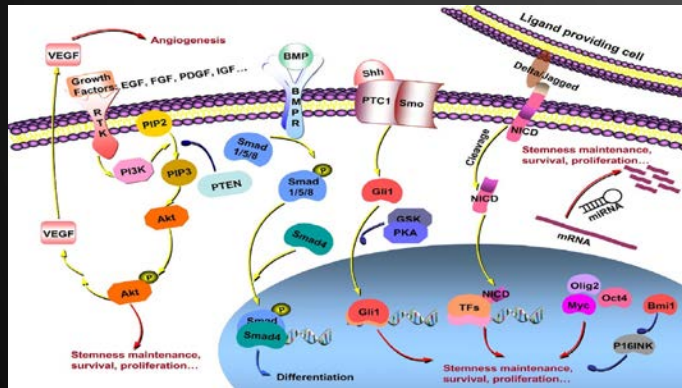I8:    AACGAGCTAGCGATCGATCGAC**A**ACGACTACGAGGT    +

Doesn't work as well
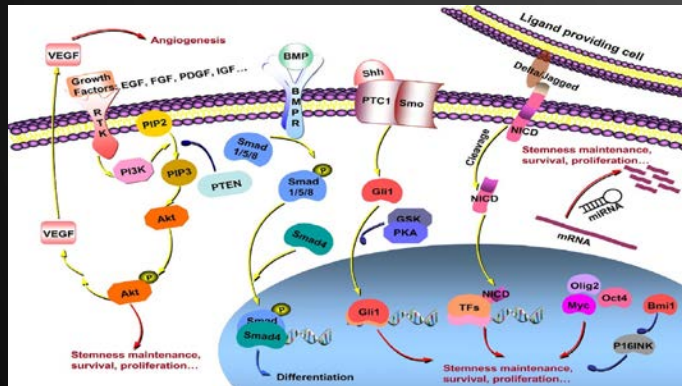as we'd like it to

# SEEK MODULES, NOT INDIVIDUAL GENES
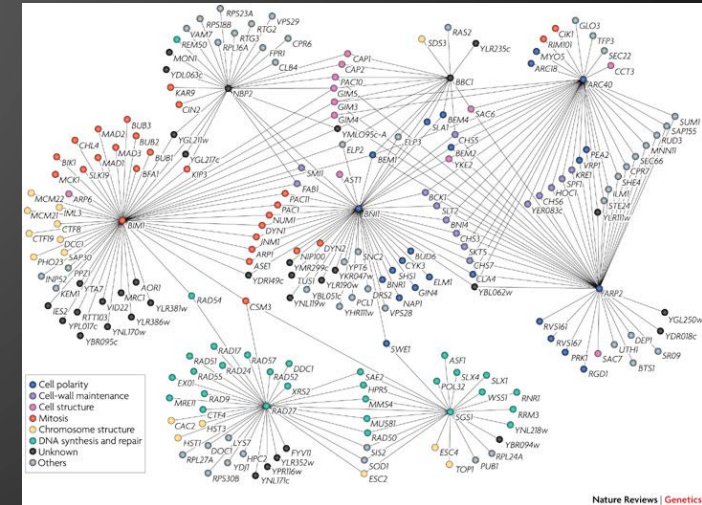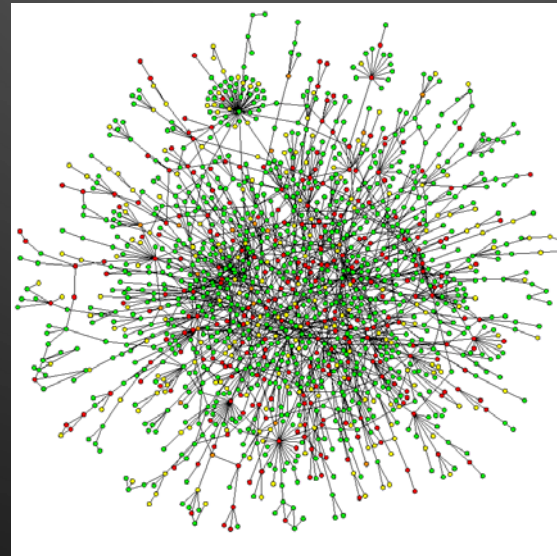


Many ways to 'break' the code
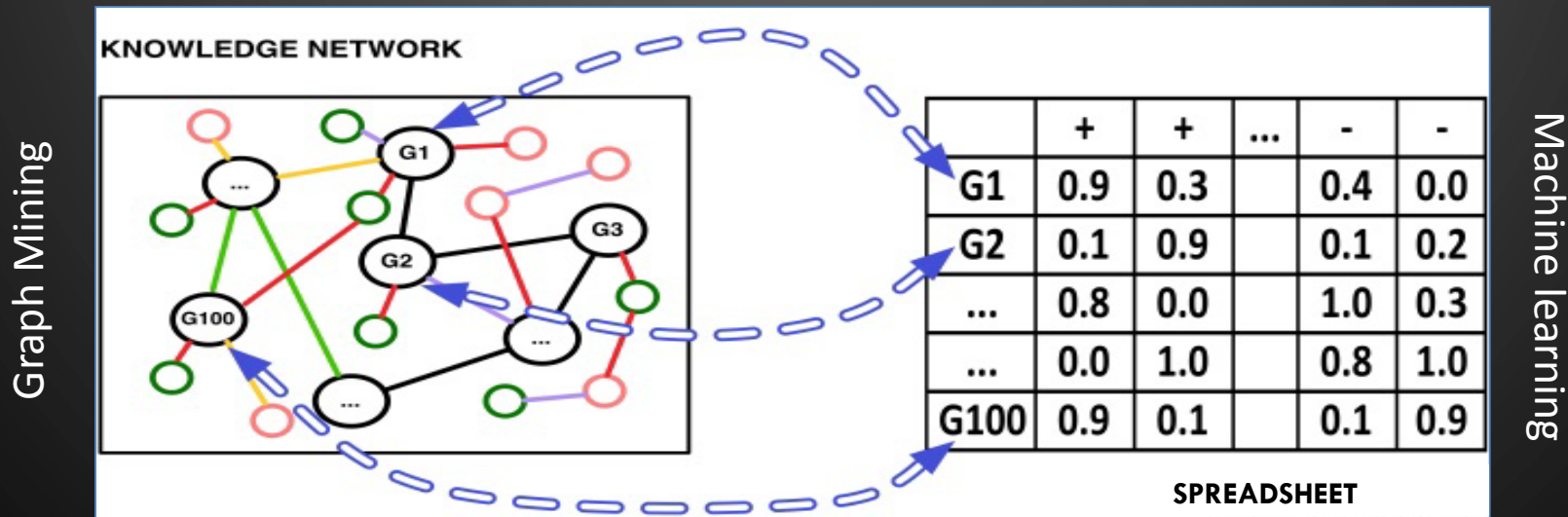
# BIOLOGICAL NETWORKS GALORE

# BIOLOGICAL NETWORKS GALORE







'Knowledge Network':
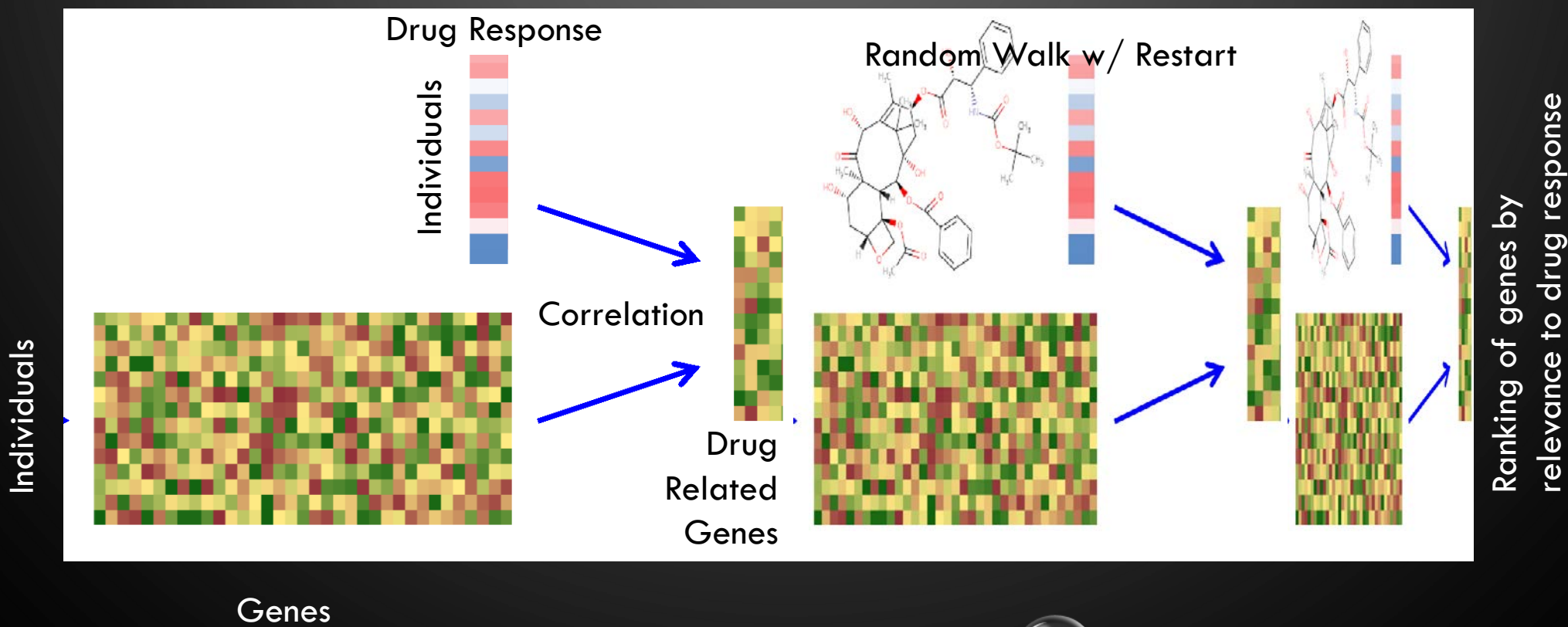3M nodes
80M edges
82 edge types

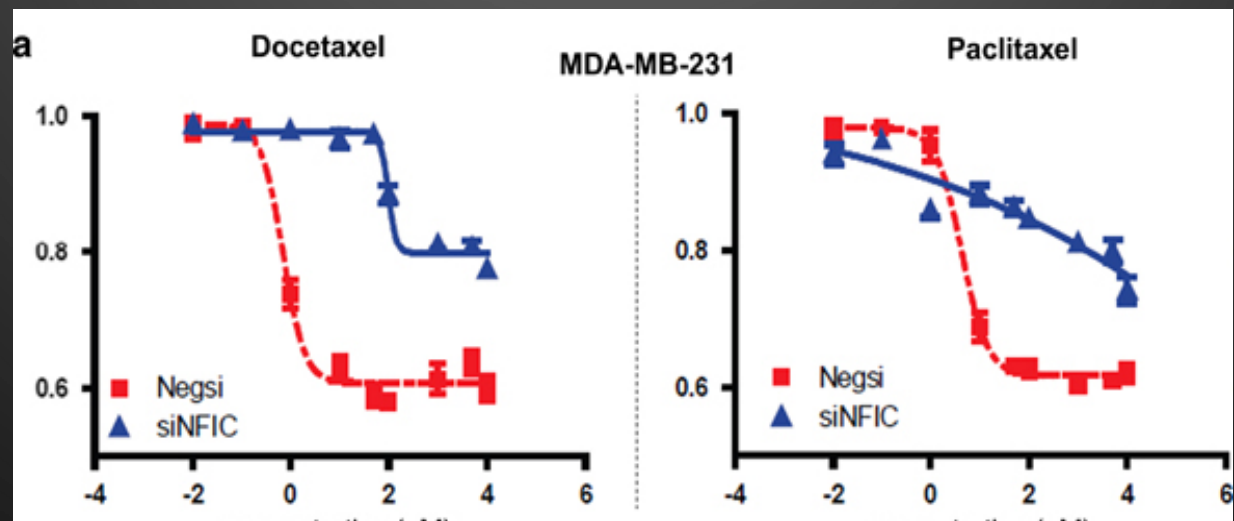# GIANT NETWORK GUIDES BIOLOGICAL ANALYSIS



Knowledge network + user spreadsheet

# EXAMPLE: FINDING GENES INFLUENCING DRUG RESPONSE

# EXAMPLE: FINDING GENES INFLUENCING DRUG RESPONSE

Validated 17 genes for several cancer drugs

# CYBERINFRASTRUCTURE

# HOW BIOLOGISTS DO BIOINFORMATICS TODAY

- HIRE BIOINFORMATICIAN OR SEEK BIOINFORMATICS COLLABORATOR.

- DELEGATE THE FOLLOWING:

  - DOWNLOAD AND INSTALL CODE.

  - BUY COMPUTE CLUSTERS

  - RUN CODE ON CLUSTER

IN SHORT, PAINFUL.

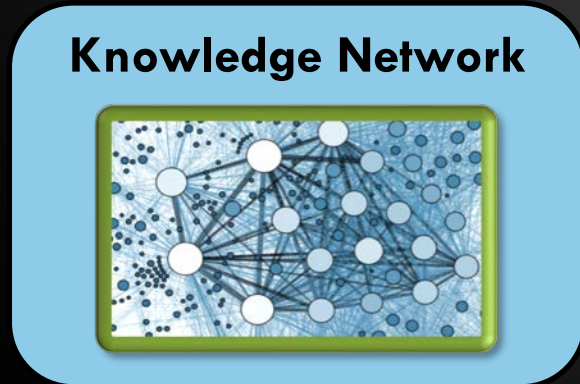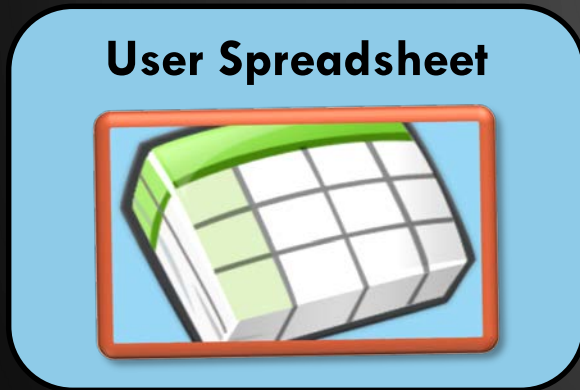# CLOUD-BASED KNOWLEDGE ENGINE FOR GENOMICS

# COMPLEX WORKFLOWS ON THE CLOUD

Docker Containers

Spreadsheet

Knowledge Network

**Network Smoothing**

Random Walk
with Restart

**Clustering Algorithm**

Network NMF

**Aggregate Subtypes**

Hierarchical
Clustering

**Easy, parallel
exploration of
workflow variants**

Network-based stratification

Patients

Patients

Network-based stratification of tumor mutations
Hofree et al. *Nature Methods* 2013

# SOFTWARE IS ONLY AS GOOD AS ITS FRONT END

KNOWENG ET AL.