# Computational Approaches for Big Data Analytics & Integration

March 16, 2017

Jianguo (Jeff) Xia, PhD

McGill University, QC Canada

# Outline

- Overview of big omics data
- Approaches in big data analytics
- Approaches in big data integration
- Web-based tools for common omics data analysis & integration

# The Promises of Big Omics Data

Comprehensive molecular profiles

– Global systems overview

– Less biased

– Robust to noise?

What we can obtain from this type of data?

– Patterns & trends

- Group behaviors
- Collective functions
- Hypothesis generation

– Mechanisms & knowledge

- Mechanistic understanding
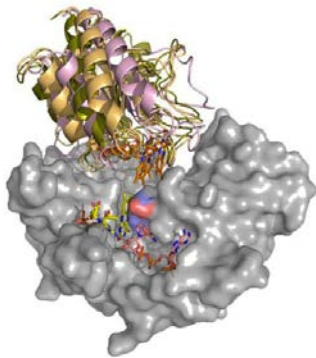- Pathways
- Networks

# Reality: the "Happy" Middle

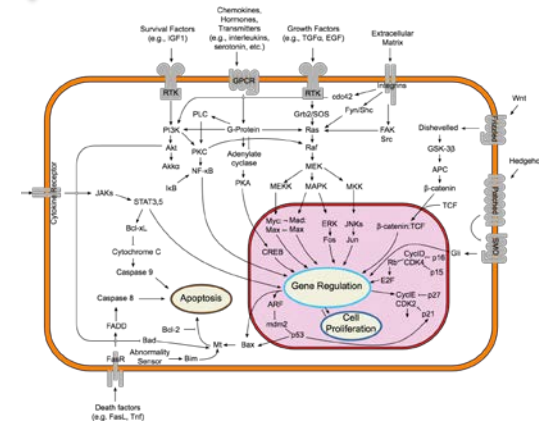Reductionism                                          Knowledge of everything

**Individual Molecule**      ⟵ **Big Data Analytics** ⟶      **Systems Biology**
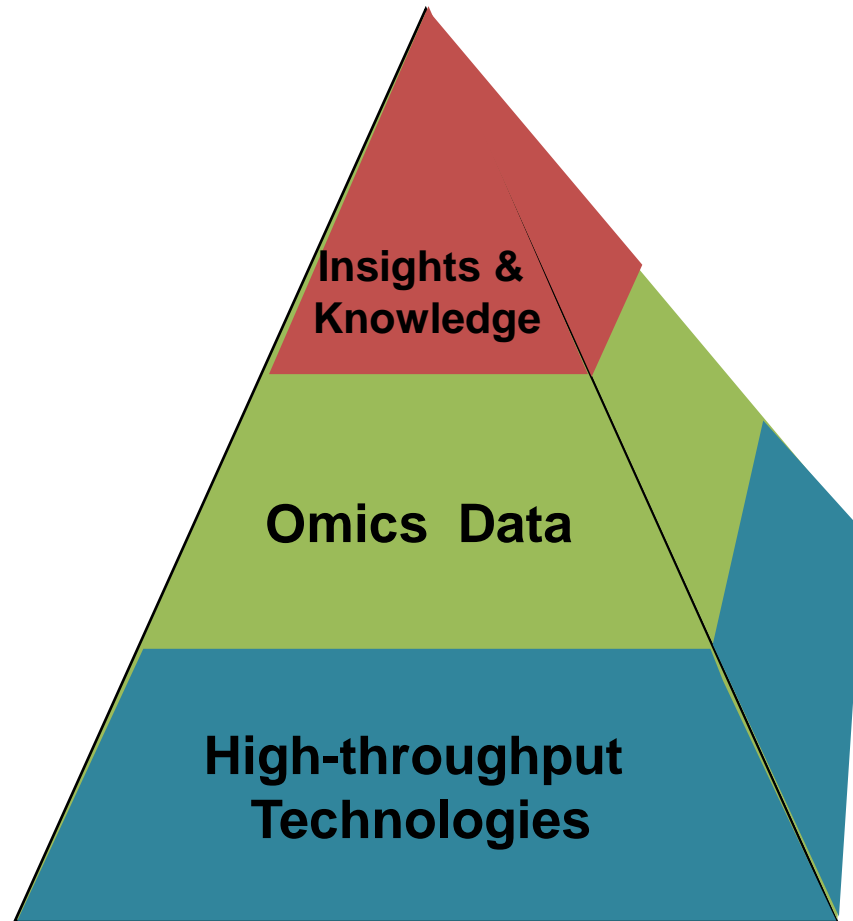
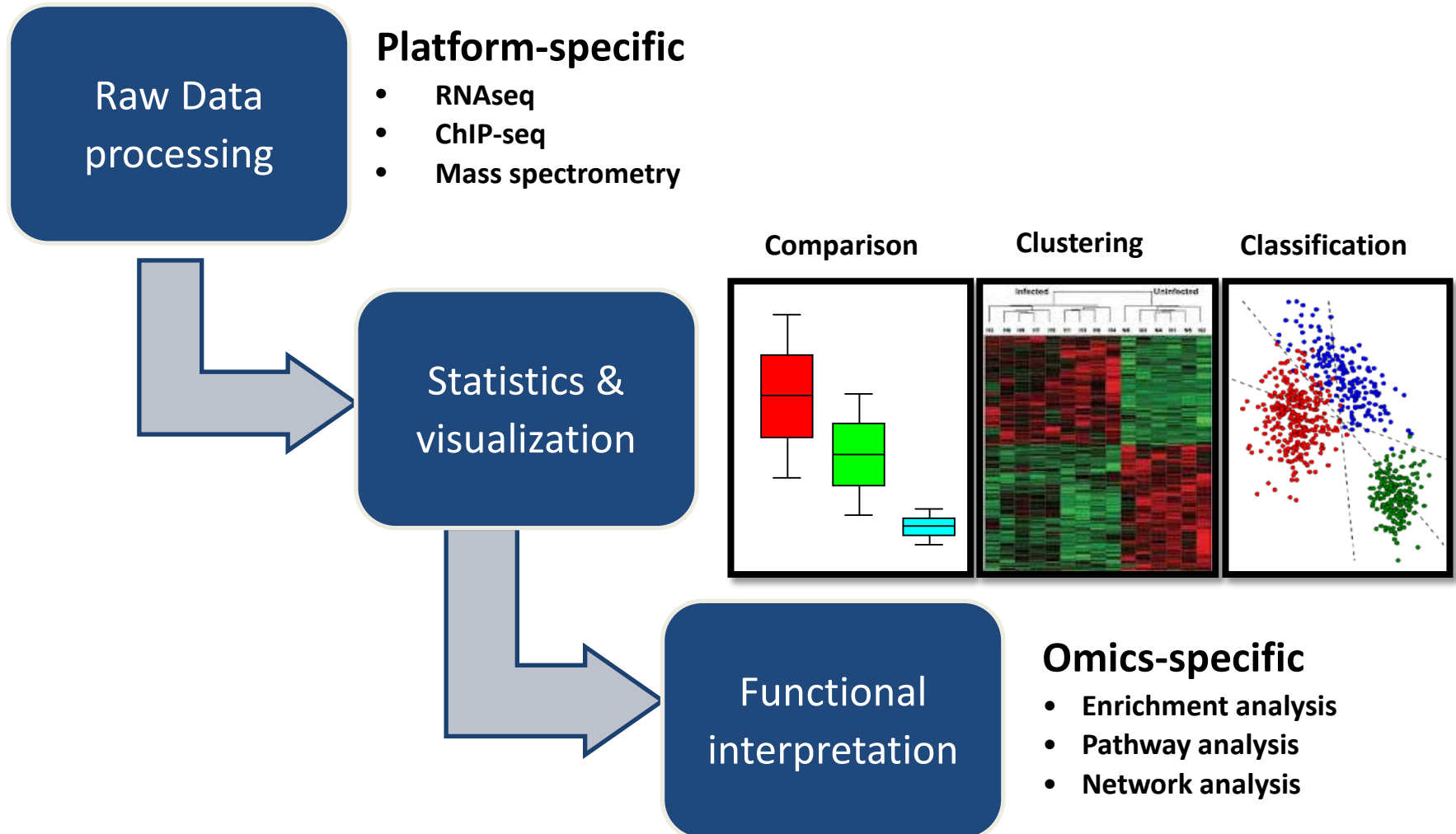- Optimal simplicity
- Empowering
- Achievable

Mechanism      ⟵ Patterns & Trends ⟶      Knowledge

# The Big Omics Data

# Workflow of Omics Data Analytics

**Raw Data processing**

**Platform-specific**

- **RNAseq**
- **ChIP-seq**
- **Mass spectrometry**

**Statistics & visualization**

**Comparison**  **Clustering**  **Classification**



**Functional interpretation**

**Omics-specific**

- **Enrichment analysis**
- **Pathway analysis**
- **Network analysis**

# Two distinct big data challenges

- **Size challenge** (raw data)
  - Raw reads, spectra, images
  - Large (100s MB ~ GB)
  - Large storage and computing resources

- **Complexity challenge** (feature table)
  - Feature table (abundance, intensities)
  - Small (100s KB ~ MB)
  - High-dimensional, missing values
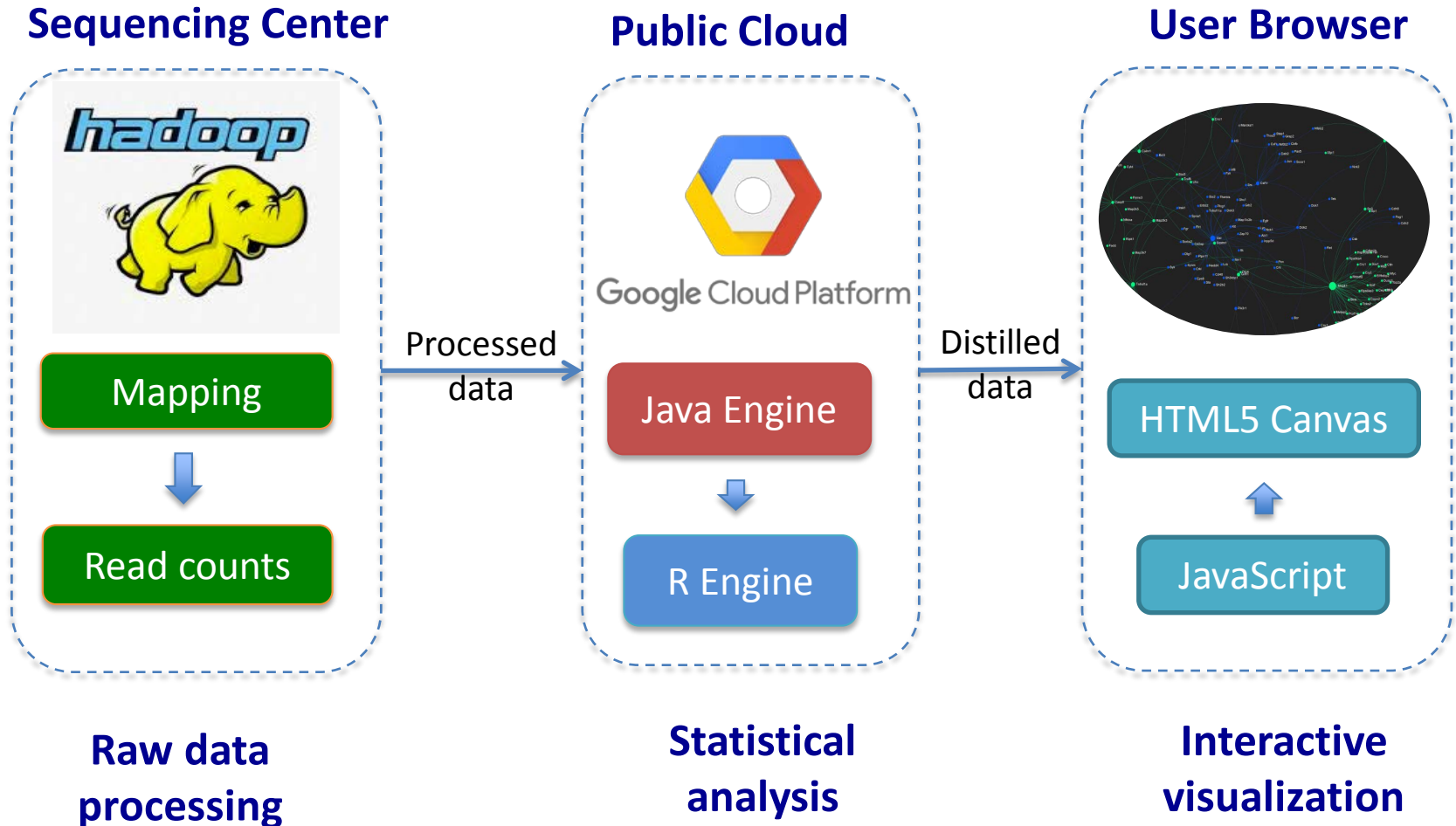  - Data integration usually starts here

➔ **Divide and conquer**

# Dealing with big omics data

- **Bringing analysis to (raw) data**
  - Too big to transfer around
  - Done **locally** at the same place generating the data (i.e. omics centers)
  - Following standardized operation protocols (SOP)
- **Bringing (processed) data to analysis**
  - Usually done by individual researchers
  - Can be easily uploaded to tools deployed on **cloud**
  - Highly domain-specific and individualized

# Local + Cloud + Browser

**Sequencing Center**



Mapping

Read counts

**Raw data processing**

Processed data →

**Public Cloud**


Google Cloud Platform

Java Engine

R Engine

**Statistical analysis**

Distilled data →

**User Browser**



HTML5 Canvas

JavaScript

**Interactive visualization**

# Big Data Analytics

# Big Data Analytics

Three Keys Factors

1. Data distillation

   - Not all data are useful

2. Coupling statistics with visualization

   - Support for user exploration

3. Connecting with prior knowledge

   - Validation previous findings & identification of novelties

# Data distillation

**Raw data ➜ processed data ➜ informative data**

- Only a small portion of molecules will respond to perturbations.

- The majority will still maintain their "normal" states (homeostasis)

# Statistics + Visualization

## Bringing data closer to domain experts

- Big omics data is complex and heterogeneous

- Statistics often fail to capture the data characteristics

- Disruptive discovery are often driven by the <u>hunches</u> and <u>leaps of faith</u> of the researcher by interacting with the data

# Statistics + Domain Knowledge

**Connecting new data with prior knowledge**

- Encode knowledge into computable forms
  - Gene sets, pathways, networks
- Evaluating new data within the context of our knowledge
  - Enrichment analysis
  - Pathway analysis
  - Network analysis

# Big data analytics in a nutshell



The Big Data

Statistics

Data Summary

Prior Knowledge

Visualization

What we know

# Omics Data Integration

# This is our goal ....



Trends in Biotechnology, April 2016, Vol. 34, No. 4

Mechanism ← Patterns & Trends → Knowledge

**Big Data Analytics**

# Current Status

- Done for simple organisms under well-controlled experiments
  - E. coli, yeast …
  - Time series to resolve false patterns
- For human studies, the limiting factor is more technological than computational
  - Longitudinal omics data collections for a large cohorts
  - Data sharing with public for algorithm benchmarking

# The 'ecosystem' of big omics data



Knowledge:
pathways, networks

Published
datasets

Meta-data

Omics data

# Data integration - three common scenarios

**One disease;**

**One omics level;**

**Multiple datasets;**

**One omics data;**

**Multiple clinical parameters:**

- Diagnosis
- Age, Gender, Ethnicity, Smoking ….

**One disease;**

**Multiple omics level;**

**Multiple datasets;**

# Computational Approaches

**Omics Data**



**Statistical integration**

**Visual integration**

**Network integration**

# Web-based Tools for Big Omics Data Analysis & Integration

- MetaboAnalyst (since 2009)
  - [www.metaboanalyst.ca](www.metaboanalyst.ca)
  - Metabolomics data analysis & integration
- NetworkAnalyst (since 2012)
  - [www.networkanalyst.ca](www.networkanalyst.ca)
  - Transcriptomics data analysis & integration
- miRNet (since 2015)
  - [www.mirnet.ca](www.mirnet.ca)
  - miRNA data analysis & integration
- MicrobiomeAnalyst  (since 2017)
  - [www.microbiomeanalyst.ca](www.microbiomeanalyst.ca)
  - Microbiome data analysis & integration

# 2009 - CURRENT

**Metabolomics & MetaboAnalyst**
**http://www.metaboanalyst.ca**

- Metabolomics (and integration with genomics)
- Real-time interactive data analysis
- 100,000 users, > 6,000 jobs submitted per day

# A Roadmap of MetaboAnalyst



MetaboAnalyst
(2009, 2012, 2015)

MetaboMiner
(2008)

MSEA (2010)    MetPA (2010)    MetATT (2011)    ROCCET (2011)

# The power of cloud



A

| 1. | Beijing |
| 2. | Boston |
| 3. | Seoul |
| 4. | Shanghai |
| 5. | Potsdam |
| 6. | New York |
| 7. | Seattle |
| 8. | Munich |
| 9. | Gainesville |

**3,000 jobs/month (2014)**
**150,000 jobs/month (2017)**
**50 times increase in traffic**

**Moved to cloud (Oct. 2014)**

200,000

100,000

January 2015 · July 2015 · January 2016 · July 2016 · January 2017

# MetaboAnalyst 3.0
## – a comprehensive tool suite for metabolomic data analysis

**Home**

**Overview**

**Data Formats**

**FAQs**

**Tutorials**

**Troubleshooting**

**Resources**

**Update History**

**User Stats**

**About**

## Welcome >> click here to start <<

### News & Updates

- Added support for peak filtering based on QC samples for untargeted metabolomics (*03/10/2017*); NEW
- Added support for "flipping" PCA for cross-study comparison (*02/09/2017*); NEW
- Added support for **network summary** of enrichment analysis result (*02/06/2017*); NEW
- Fixed the bug in feature table display in Biomarker Tester module (*01/05/2017*); NEW
- Updated the pathway result table to show all/matched compounds (*11/25/2016*);
- Enhanced Normalization and Data Editor for better user experience (*11/15/2016*);
- Added support for **sparse PLS-DA** (sPLS-DA) analysis (*10/28/2016*);
- Added support for **quantile normalization** (*08/29/2016*);
- Improved name mapping functions for common metabolite names (*08/18/2016*);
- More than **1 million jobs** have been processed since 06/2015 (*06/21/2016*);

Read more ......

**Please Cite:**

Xia, J. and Wishart, D.S. (2016) Using MetaboAnalyst 3.0 for Comprehensive Metabolomics Data Analysis Current Protocols in Bioinformatics, 55:14.10.1-14.10.91.

Xia, J., Sinelnikov, I., Han, B., and Wishart, D.S. (2015) MetaboAnalyst 3.0 - making metabolomics more meaningful . Nucl. Acids Res. 43, W251-257.

Xia, J., Mandal, R., Sinelnikov, I., Broadhurst, D., and Wishart, D.S. (2012) MetaboAnalyst 2.0 - a comprehensive server for metabolomic data analysis . Nucl. Acids Res. 40, W127-133.

# Comprehensive Options for Data Analysis

**Sample normalization**

- ● None
- ○ Sample-specific normalization (i.e. weight, volume) Click here to specify
- ○ Normalization by sum
- ○ Normalization by median
- ○ Normalization by a specific reference sample — PIF_178
- ○ Normalization by a pooled sample from group — cachexic
- ○ Normalization by reference feature — 1,6-Anhydro-beta-D-glucose
- ○ Quantile normalization

**Data transformation**

- ● None
- ○ Log transformation (generalized logarithm transformation or glog)
- ○ Cube root transformation (take cube root of data values)

**Data scaling**

- ● None
- ○ Mean centering (mean-centered only)
- ○ Auto scaling (mean-centered and divided by the standard deviation of each variable)
- ○ Pareto scaling (mean-centered and divided by the square root of standard deviation of each variable)
- ○ Range scaling (mean-centered and divided by the range of each variable)

**Univariate Analysis**

Fold Change Analysis   T-tests   Volcano plot

One-way Analysis of Variance (ANOVA)

Correlation Analysis   Pattern Searching

**Chemometrics Analysis**

Principal Component Analysis (PCA)

Partial Least Squares - Discriminant Analysis (PLS-DA)

Sparse Partial Least Squares - Discriminant Analysis (sPLS-DA)

Orthogonal Partial Least Squares - Discriminant Analysis (orthoPLS-DA)

**Feature Identification**

Significance Analysis of Microarray (and Metabolites) (SAM)

Empirical Bayesian Analysis of Microarray (and Metabolites) (EBAM)

**Cluster Analysis**

Hierarchical Clustering:   Dendrogram   Heatmaps

Partitional Clustering:   K-means   Self Organizing Map (SOM)

**Classification & Feature Selection**

Random Forest

Support Vector Machine (SVM)

# Clustering

# Multivariate Statistics

# Pathway Analysis

# Biomarker Analysis

- Highly relevant for translational studies
- Performance evaluation
  - Receiver operator characteristics (ROC) curve
  - Modern machine learning approaches
    - Cross validation
    - Permutation
    - Predicting new samples



AUC = 87%



AUC = 97%

| Var. | AUC | CI |
|---|---|---|
| 2 | 0.978 | 0.975-0.981 |
| 3 | 0.978 | 0.974-0.981 |
| 5 | 0.975 | 0.967-0.983 |
| 10 | 0.948 | 0.936-0.96 |
| 20 | 0.916 | 0.902-0.931 |
| 39 | 0.908 | 0.894-0.923 |

# Gene-metabolite joint pathway analysis

# Comprehensive report generation

# 2012 - CURRENT

Transcriptomics & NetworkAnalyst
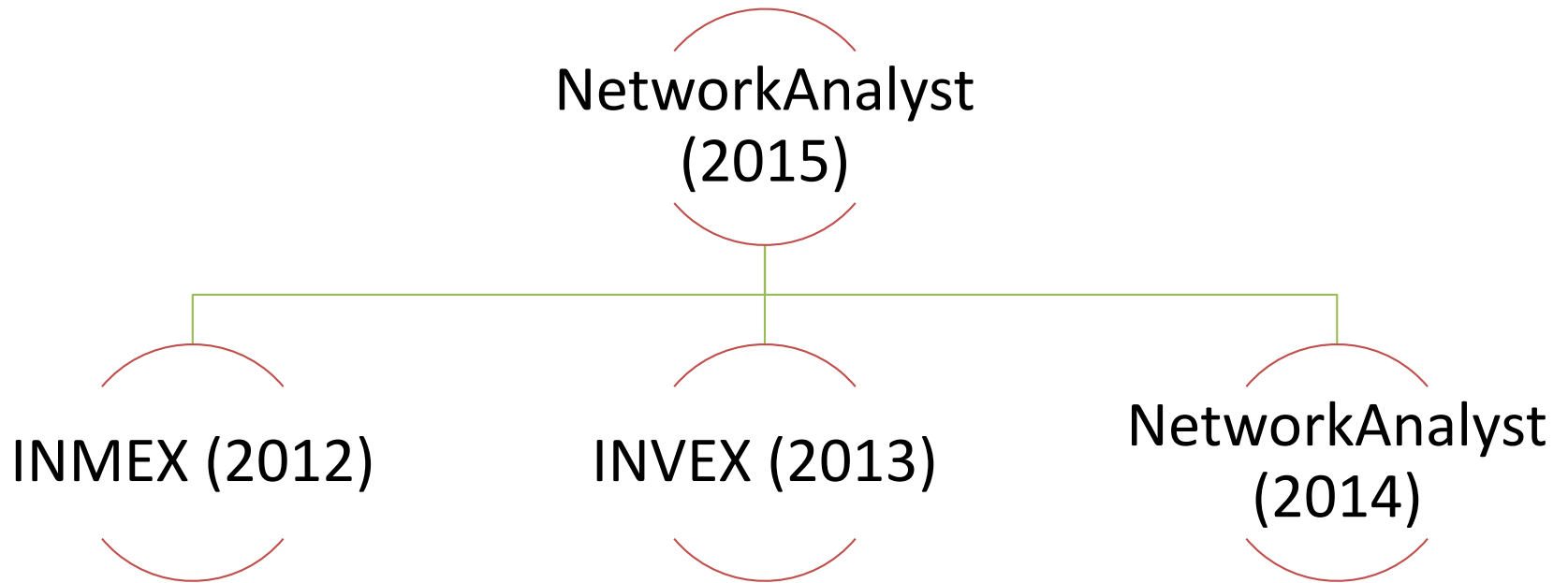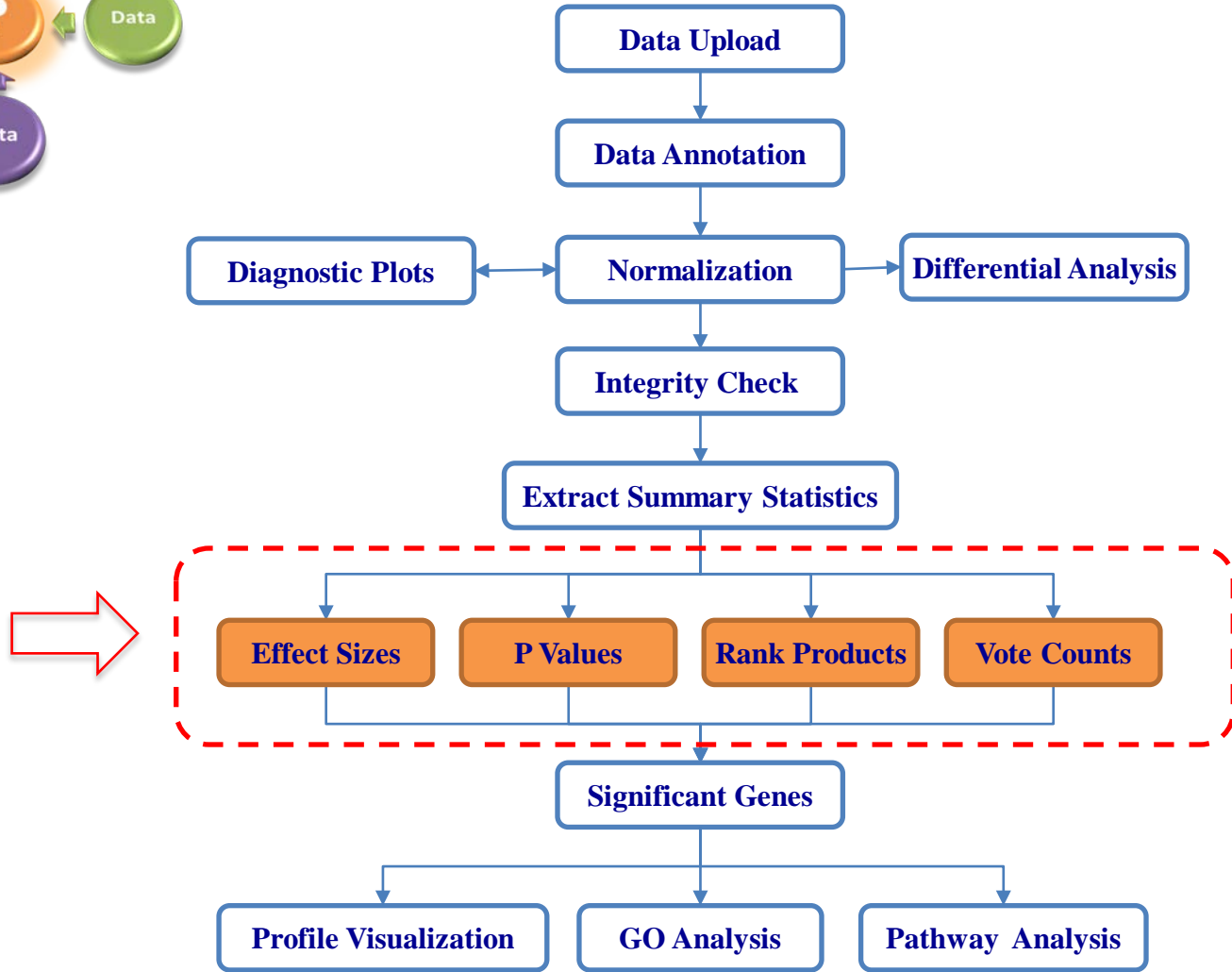http://www.networkanalyst.ca

- Gene expression analysis (microarray & RNAseq)
- Statistical data integration (meta-analysis)
- Visual integration (heatmaps, Venn diagrams, 3D PCA/tSNE)
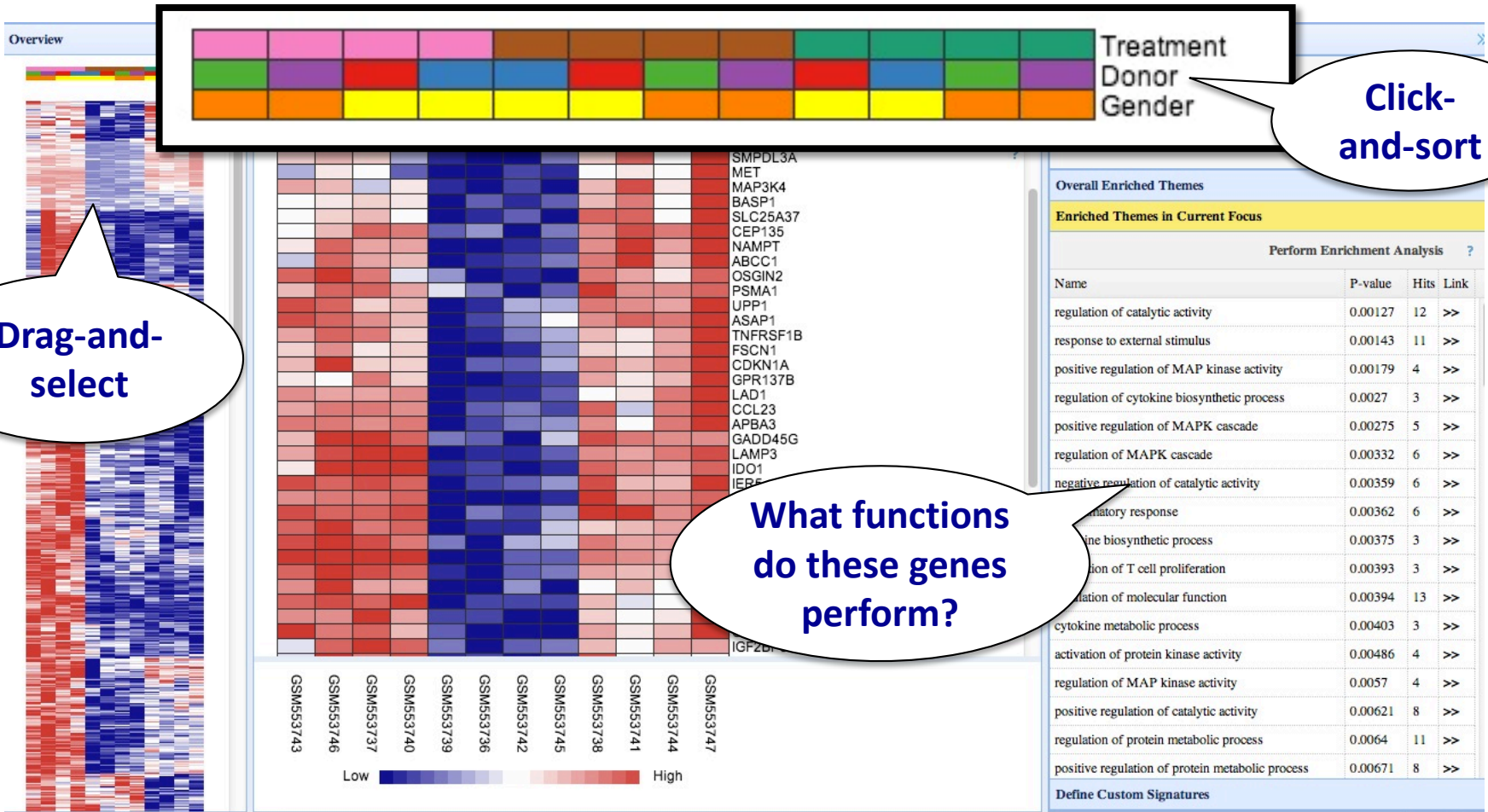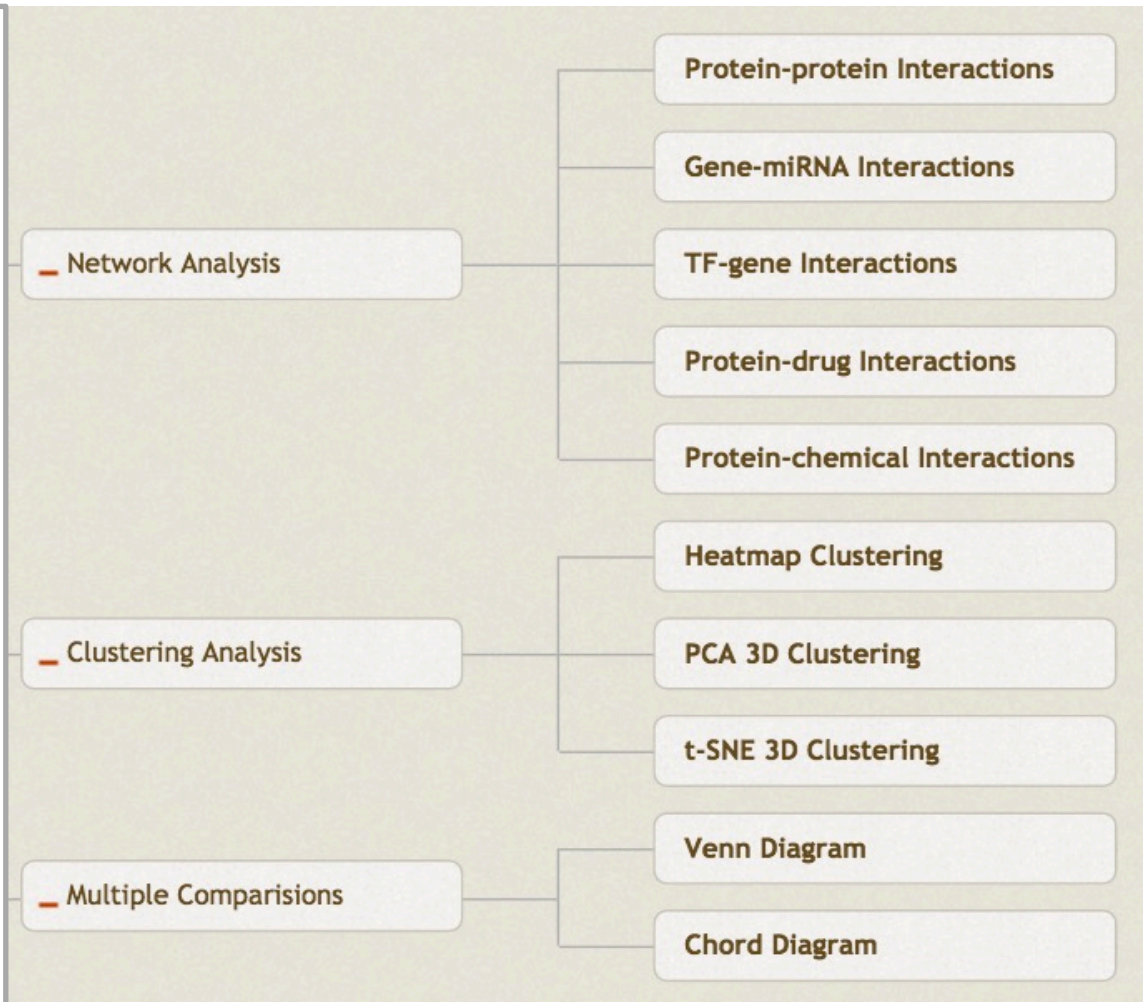- Network integration (PPI, KEGG, miRanda)

# Implementation Roadmap

NetworkAnalyst
(2015)

INMEX (2012)        INVEX (2013)        NetworkAnalyst
(2014)

# Statistical Integration



```
                              Data Upload
                                   │
                                   ▼
                             Data Annotation
                                   │
                                   ▼
   Diagnostic Plots  ◄────►    Normalization    ────►  Differential Analysis
                                   │
                                   ▼
                             Integrity Check
                                   │
                                   ▼
                        Extract Summary Statistics
                                   │
   ┌───────────────────────────────────────────────────────────────┐
   │   Effect Sizes      P Values      Rank Products    Vote Counts │
   └───────────────────────────────────────────────────────────────┘
                                   │
                                   ▼
                            Significant Genes
                                   │
          ┌────────────────────────┼────────────────────────┐
          ▼                        ▼                        ▼
  Profile Visualization        GO Analysis          Pathway Analysis
```

# Visual Integration



*J. Xia. et al. (2013) Bioinformatics 29 (24), 3232-3234*

# Integrating with prior knowledge

- **Protein-protein interactions**
  - ➢ STRING
  - ➢ InnateDB
- **Metabolic pathways**
  - ➢ KEGG
  - ➢ Reactome
- **Chemicals**
  - ➢ DrugBank
  - ➢ CTD
- **Gene regulations**
  - ➢ Trans Factor
  - ➢ miRNAs

**Network Analysis**
- Protein-protein Interactions
- Gene-miRNA Interactions
- TF-gene Interactions
- Protein-drug Interactions
- Protein-chemical Interactions

**Clustering Analysis**
- Heatmap Clustering
- PCA 3D Clustering
- t-SNE 3D Clustering

**Multiple Comparisions**
- Venn Diagram
- Chord Diagram

# Network Visualization & Exploration

# MicroRNA Data (www.mirnet.ca)



Fan, Y et al. (2016) (doi: 10.1093/nar/gkw288)

# Microbiome Data (www.microbiomeanalyst.ca)

Dhariwal, A et al. (2017) (under review)

# Microbiome profiling & integration

# In the near future ......

# Acknowledgements

- Xia Lab @ McGill University
- Wishart Lab @ Univ. Alberta
- Hancock Lab @ Univ. of British Columbia