# Bioinformatics for cancer research

**Bing Zhang, Ph.D.**

*Professor of Molecular and Human Genetics*

*Lester & Sue Smith Breast Center*

*Baylor College of Medicine*

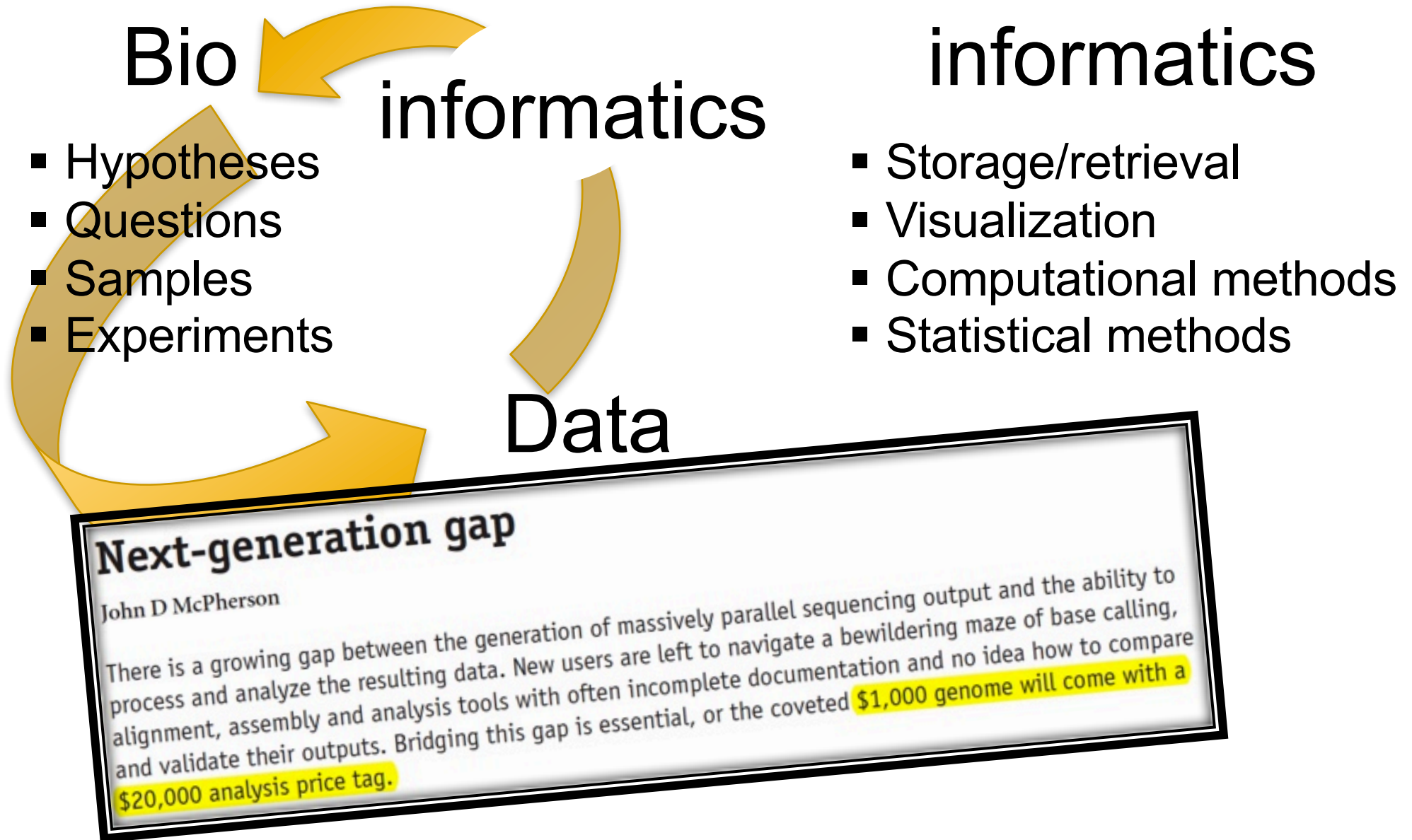*bing.zhang@bcm.edu*

# What is bioinformatics

## Bio    Bioinformatics    informatics

- Hypotheses
- Questions
- Samples
- Experiments

- Storage/retrieval
- Visualization
- Computational methods
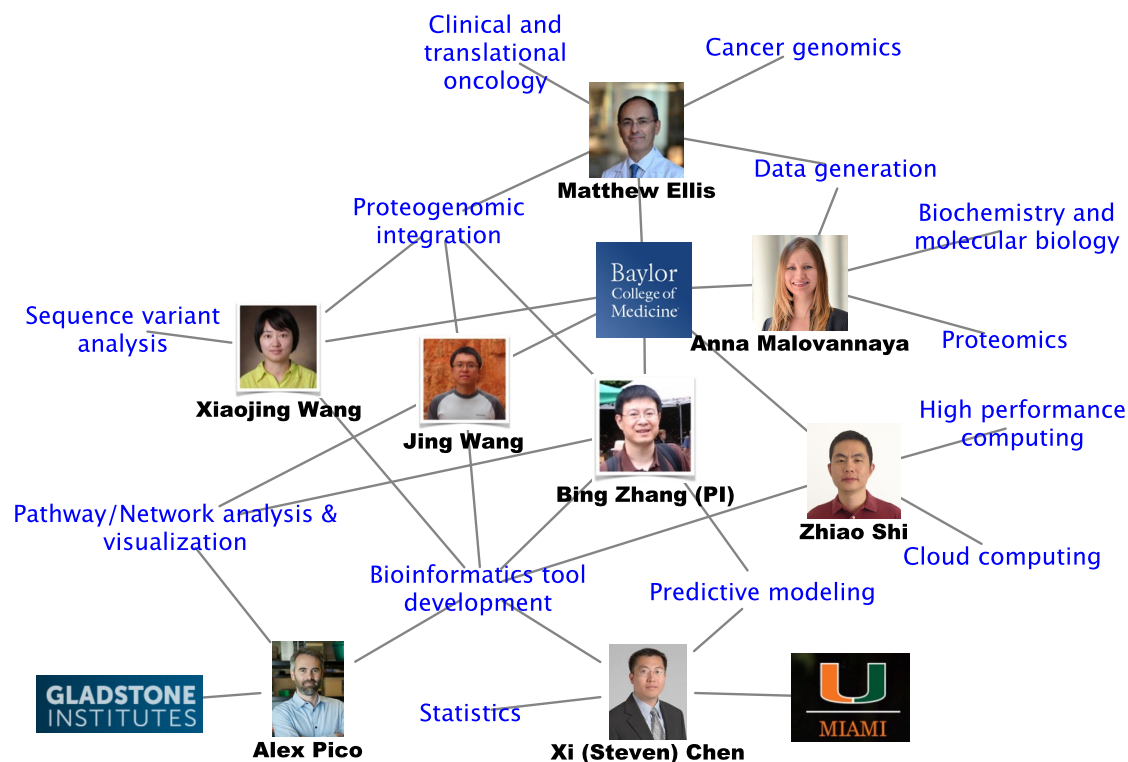- Statistical methods

## Data

- DNA
- RNA
- Protein
- Metabolite
- Phenotype

- Sequence
- Expression
- Structure
- Interaction

Translational Breast Cancer Research, 2016

# Why now?

## Bio informatics

- Hypotheses
- Questions
- Samples
- Experiments

## informatics

- Storage/retrieval
- Visualization
- Computational methods
- Statistical methods

## Data

### Next-generation gap

John D McPherson

There is a growing gap between the generation of massively parallel sequencing output and the ability to process and analyze the resulting data. New users are left to navigate a bewildering maze of base calling, alignment, assembly and analysis tools with often incomplete documentation and no idea how to compare and validate their outputs. Bridging this gap is essential, or the coveted $1,000 genome will come with a $20,000 analysis price tag.

Translational Breast Cancer Research, 2016

# Roles for different investigators in bioinformatics



**iPGDAC team**

- **Algorithm developer**
  - ❑ Statisticians
  - ❑ Mathematicians
  - ❑ Computer scientists
- **Tool developer**
  - ❑ Bioinformaticians
- **Data provider/consumer**
  - ❑ Biologists

Translational Breast Cancer Research, 2016

# Comprehensive list of bioinformatics resources



**Bioinformatics Links Directory**

The Bioinformatics Links Directory features curated links to molecular resources, tools and databases. The links listed in this directory are selected on the basis of recommendations from bioinformatics experts in the field. We also rely on input from our community of bioinformatics users for suggestions. Starting in 2003, we have also started listing all links contained in the NAR Webserver issue.

Hide Resources (176)  Hide Databases (621)  Hide Tools (1548)

**Computer Related** (85)
This category contains links to resources relating to programming languages often used in bioinformatics. Other tools of the trade, such as web development and database resources, are also included here.

**DNA** (604)
This category contains links to useful resources for DNA sequence analyses such as tools for comparative sequence analysis and sequence assembly. Links to programs for sequence manipulation, primer design, and sequence retrieval and submission are also listed here.

**Education** (75)
Links to information about the techniques, materials, people, places, and events of the greater bioinformatics community. Included are current news headlines, literature sources, educational material and links to bioinformatics courses and workshops.

**Expression** (396)
Links to tools for predicting the expression, alternative splicing, and regulation of a gene sequence are found here. This section also contains links to databases, methods, and analysis tools for protein expression, SAGE, EST, and microarray data.

**Human Genome** (240)
This section contains links to draft annotations of the human genome in addition to resources for sequence polymorphisms and genomics. Also included are links related to ethical discussions surrounding the study of the human genome.

**Literature** (87)
Links to resources related to published literature, including tools to search for articles and through literature abstracts. Additional text mining resources, open access resources, and literature goldmines are also listed.

**Model Organisms** (378)
Included in this category are links to resources for various model organisms ranging from mammals to microbes. These include databases and tools for genome scale analyses.

**Other Molecules** (117)
Bioinformatics tools related to molecules other than DNA, RNA, and protein. This category will include resources for the bioinformatics of small molecules as well as for other biopolymers including carbohydrates and metabolites.

**Protein** (1007)
This category contains links to useful resources for protein sequence and structure analyses. Resources for phylogenetic analyses, prediction of protein features, and analyses of interactions are also found here.

**RNA** (203)
Resources include links to sequence retrieval programs, structure prediction and visualization tools, motif search programs, and information on various functional RNAs.

**Sequence Comparison** (271)
Tools and resources for the comparison of sequences (nucleic acid or protein) including sequence similarity searching, alignment tools, classification and general comparative genomics resources.

- **October 2016**

  - 176 Resources

  - 621 Databases

  - 1548 Tools

http://bioinformatics.ca/links_directory/

Translational Breast Cancer Research, 2016

# Sequence and structure databases
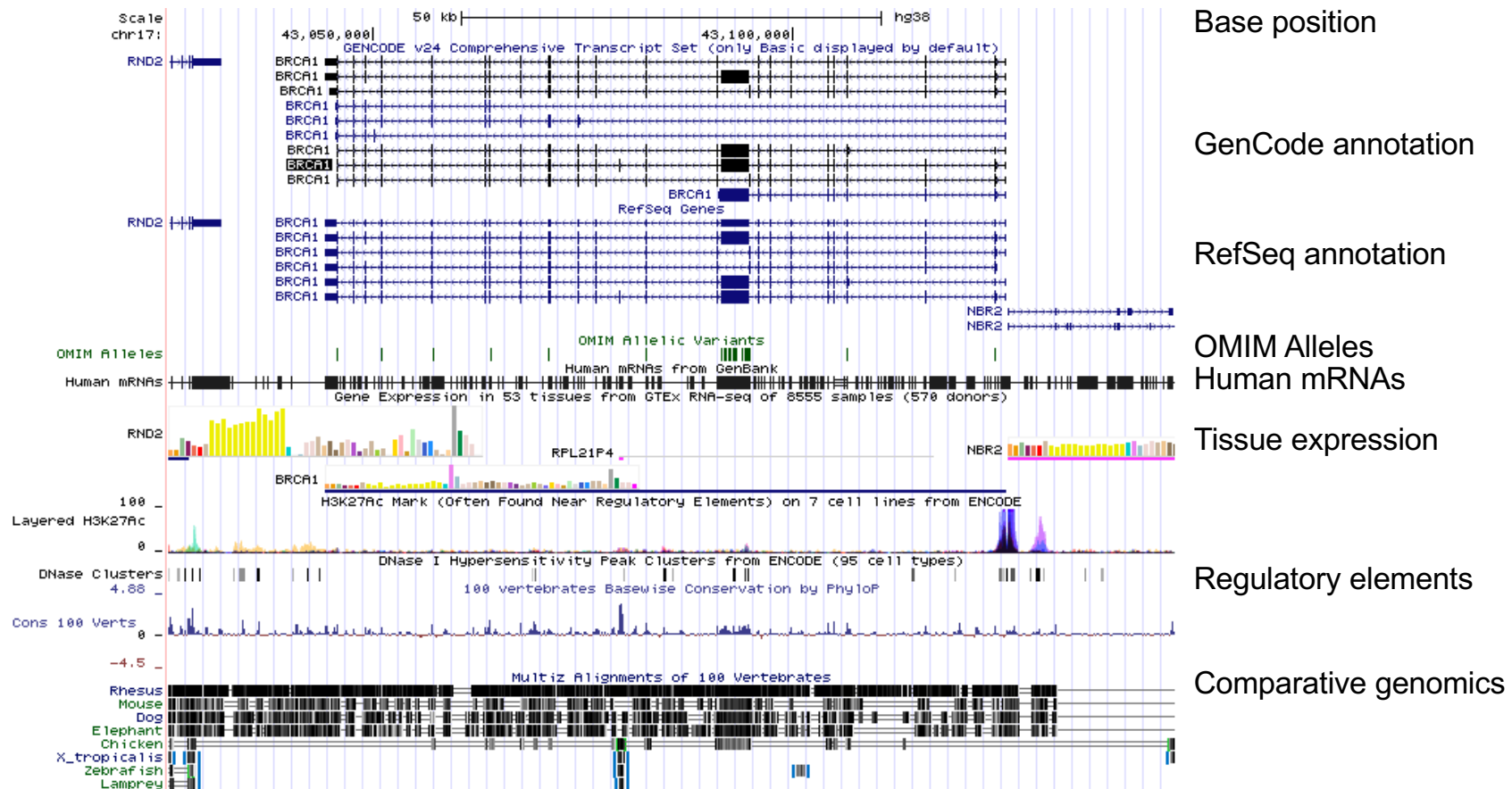
- **Genbank:** http://www.ncbi.nlm.nih.gov/genbank/
  - Annotated collection of all publicly available **DNA sequences**
  - 220,731,315,250 bases in 197,390,691 sequences as of October 2016
  - Whole Genome Sequencing (WGS) data: ftp://ftp.ncbi.nih.gov/ncbi-asn1/wgs ftp://ftp.ncbi.nih.gov/genbank/wgs
  - WGS: 1,676,238,489,250 bases in 363,213,315 sequences as of October 2016

- **UniProt:** http://www.uniprot.org/
  - Comprehensive resource for **protein sequences** and functional information
  - 552,259 reviewed entries as of October 2016

- **PDB:** http://www.rcsb.org/
  - **3D structures** of large biological molecules, including proteins, nucleic acids, and complex assemblies
  - 123,870 structures as of October 2016

- **Pfam:** http://pfam.xfam.org/
  - Collection of **protein families**, each represented by multiple sequence alignments and hidden Markov models (HMMs)
  - 16,306 families as of October 2016

# Genome browsers

Graph interface for browsing and visualizing genome-wide sequence and annotation data.

- UCSC genome browser
    - http://genome.ucsc.edu/cgi-bin/hgGateway
- Ensembl genome browser
    - http://www.ensembl.org/index.html

- Integrative Genomics Viewer (IGV)
    - http://software.broadinstitute.org/software/igv/

## UCSC genome browser screenshot



Translational Breast Cancer Research, 2016

# Genome browsers

IGV: copy number, expression and mutation data grouped by tumor subtype



EGFR

Robinson et al. Nat Biotechnol, 2011

Translational Breast Cancer Research, 2016

# Genome browsers

IGV: view of aligned reads at 20Kb resolution



Robinson et al. Nat Biotechnol, 2011

Translational Breast Cancer Research, 2016

# Gene-centric databases

- **Entrez Gene**
  - http://www.ncbi.nlm.nih.gov/gene
  - NCBI/NIH
  - All completely sequenced genomes
  - **One gene per page**

- **Ensembl BioMart**
  - http://www.ensembl.org/biomart/martview
  - EMBL-EBI and Sanger Institute
  - Vertebrates and other selected eukaryotic species
  - **Batch information retrieval**

# Gene/protein expression data repositories

- Gene Expression Omnibus (GEO)
  - http://www.ncbi.nlm.nih.gov/geo/

- ArrayExpress
  - http://www.ebi.ac.uk/arrayexpress/

- PRIDE
  - https://www.ebi.ac.uk/pride/archive/

# Pathway and network databases

- Gene Ontology (GO): http://www.geneontology.org/

- Pathway databases
  - KEGG: http://www.genome.jp/kegg/pathway.html
  - Reactome: http://www.reactome.org/
  - WikiPathways: http://www.wikipathways.org/

- Protein-protein interaction databases
  - DIP: http://dip.doe-mbi.ucla.edu/
  - MINT: http://mint.bio.uniroma2.it/mint/
  - BioGRID: http://www.thebiogrid.org/
  - HPRD: http://www.hprd.org
  - iRef: http://wodaklab.org/iRefWeb

- Protein-DNA interaction database
  - Transfac: http://www.gene-regulation.com
  - Jaspar: http://jaspar.genereg.net/

Translational Breast Cancer Research, 2016

# Pathway and network analysis: motivation

- **Genomics**
  - Genome Wide Association Study (GWAS)
  - Whole genome or exome sequencing
  - Copy number analysis
- **Epigenomics**
  - DNA methylation
- **Transcriptomics**
  - mRNA profiling
    - Microarray
    - RNA-Seq
  - Protein-DNA interaction
    - Chromatin immunoprecipitation (ChIP)-Seq
- **Proteomics**
  - Protein profiling
    - LC-MS/MS
  - Protein-protein interaction
    - Yeast two hybrid
    - Affinity pull-down/LC-MS/MS

# Pathway and network analysis: tools

- **Pathway analysis**
  - WebGestalt: http://www.webgestalt.org
  - DAVID: https://david.ncifcrf.gov/
  - GSEA: http://software.broadinstitute.org/gsea

- **Network analysis**
  - Cytoscape: http://www.cytoscape.org/
  - NetGestalt: http://www.netgestal.org
  - STRING: http://string-db.org
  - GeneMANIA: http://genemania.org/
  - Gene2Net: http://www.gene2net.org

# WebGestalt: http://www.webgestalt.org

**Gene list**

```
92546_r_at
92545_f_at
96055_at
102105_f_at
102700_at
......
```

~200
ID types

~60K
Functional
categories

Pathways/
functional categories

**8 organisms**
**Human, Mouse, Rat, Dog, Fruitfly, Worm, Zebrafish, Yeast**

**Microarray Probe IDs**
- Affymetrix
- Agilent
- Codelink
- Illumina

**Genetic Variation IDs**
- dbSNP

**Gene IDs**
- Gene Symbol
- GenBank
- Ensembl Gene
- RefSeq Gene
- UniGene
- Entrez Gene
- SGD
- MGI
- Flybase ID
- Wormbase ID
- ZFIN

**Protein IDs**
- UniProt
- IPI
- RefSeq Peptide
- Ensembl Peptide

196 ID types with mapping to Entrez Gene ID

WebGestalt

Analysis/
visualization

59,278 functional categories with genes identified by
Entrez Gene IDs

**Gene Ontology**
- Biological Process
- Molecular Function
- Cellular Component

**Pathway**
- KEGG
- Pathway Commons
- WikiPathways

**Network module**
- Transcription factor targets
- microRNA targets
- Protein interaction modules

**Disease and Drug**
- Disease association genes
- Drug association genes

**Chromosomal location**
- Cytogenetic bands

*Zhang et.al. Nucleic Acids Res. 33:W741, 2005*
*Wang et al. Nucleic Acids Res. 41:W77, 2013*

*http://www.webgestalt.org*

Daily Unique Visitors (Jan 2015 – Dec 2015)

- Users

400

200

April 2015    July 2015    October 2015

Visits

1    13,681

Jan. 1, 2015 – Dec. 31, 2015
63,932 visits from 27,409 visitors
>300 citations

# Cancer-specific resources

- Pavlopoulou et al., **Human cancer databases (Review), Oncology Reports**, 33:3-18, 2015

- Yang et al., **Databases and web tools for cancer genomics study**, Genomics proteomics bioinformatics, 13: 46-50, 2015

- https://www.oxfordjournals.org/our_journals/nar/database/subcat/8/33



Nucleic Acids Research

Oxford Journals · Life Sciences · Nucleic Acids Research · Database Summary Paper

**NAR Database Summary Paper**

Nucleotide Sequence Databases
RNA sequence databases
Protein sequence databases
Structure Databases
Genomics Databases (non-vertebrate)
Metabolic and Signaling Pathways
Human and other Vertebrate Genomes
Human Genes and Diseases
  CancerResource
  DriverDBv2
  Protein Mutant Database
General human genetics databases
General polymorphism databases
Cancer gene databases
  ArrayMap
  Atlas of Genetics and Cytogenetics in Oncology and Haematology
  BCCTBbp
  BreCAN-DB
  Cancer RNA-Seq Nexus
  Cancer3D
  CancerGenes
  CancerPPD
  Candidate Cancer Gene Database
  CanGEM
  CanSAR
  CaSNP
  cBioPortal
  CCDB
  ccmGDB
  CellLineNavigator
  CGED - Cancer Gene Expression Database
  ChimerDB
  CMPD
  Colorectal Cancer Atlas
  COLT-Cancer
  COSMIC
  CTDatabase
  Database of Germline p53 Mutations
  dbDEPC
  DDOC
  DDPC
  DIVAS
  DriverDB
  EHCO
  FusionCancer
  HLungDB
  HPTAA
  Human p53, human hprt, rodent lacI and rodent lacZ databases
  IARC TP53 Database
  IGDB.NSCLC
  intOGen
  ITTACA
  Lnc2Cancer
  MethHC
  MethyCancer
  MoKCa
  Mouse Tumor Biology Database
  MutationAligner
  Network of Cancer Genes
  OncoDB.HCC
  OpenTein
  Pancreas Expression
  PubMeth
  SNP500Cancer
  Stem Cell Discovery Engine
  SynLethDB
  TCGA SpliceSeq
  TSGene
  Tumor Associated Gene database
  Tumor Gene Family Databases (TGDBs)
  UCSC Cancer Genomics Browser
  UMD-BRCA1/BRCA2 databases

› Compilation Paper
› Category List
› Alphabetical List
› Category/Paper List
› Search Summary Papers

Translational Breast Cancer Research, 2016

# Catalogue of somatic mutations in cancer (COSMIC)

- COSMIC is designed to store and display somatic mutation information and related details and contains information relating to human cancers

- Wellcome Trust Sanger Institute

- http://cancer.sanger.ac.uk/cosmic

- Expert curation data and genome-wide screen data

- Search by gene, cancer type, mutation, or sample



Translational Breast Cancer Research, 2016

# COSMIC: breast cancer



- **Understanding mutation frequency**
  - What mutation detection method was employed
  - Was the whole gene screened
  - Has the sample been screened before
  - Are all mutations real?

Translational Breast Cancer Research, 2016

# COSMIC: getting help



http://cancer.sanger.ac.uk/cosmic/help

Translational Breast Cancer Research, 2016

# Cancer Gene Census

- Futreal et al. **A census of human cancer genes**. Nature Reviews Cancer, 4:2004

- The Cancer Gene Census is an ongoing effort to catalogue those genes for which mutations have been causally implicated in cancer.

- 602 genes as of October 2016.

http://cancer.sanger.ac.uk/census



Translational Breast Cancer Research, 2016

# Genomics of Drug Sensitivity in Cancer (GDSC)

- A collaboration between the Cancer Genome Project at the Wellcome Trust Sanger Institute (UK) and the Center for Molecular Therapeutics, Massachusetts General Hospital Cancer Center (USA), funded by the Wellcome Trust.

- Goal: to identify molecular features of cancers that predict response to anti-cancer drugs.



http://www.cancerrxgene.org/

Translational Breast Cancer Research, 2016

# GDSC: drug sensitivity vs Her2 amplification



http://www.cancerrxgene.org/

Translational Breast Cancer Research, 2016

# The Cancer Genome Atlas (TCGA)

- A collaboration between the National Cancer Institute (NCI) and National Human Genome Research Institute (NHGRI)

- To accelerate the understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing.



**NATIONAL CANCER INSTITUTE
THE CANCER GENOME ATLAS**

**TCGA BY THE NUMBERS**

TCGA produced over
**2.5** PETABYTES of data

To put this into perspective, **1 petabyte** of data is equal to
**212,000** DVDs

TCGA data describes
**33** DIFFERENT TUMOR TYPES

...including
**10** RARE CANCERS

...based on paired tumor and normal tissue sets collected from
**11,000** PATIENTS

...using
**7** DIFFERENT DATA TYPES

**TCGA RESULTS & FINDINGS**

MOLECULAR BASIS OF CANCER — Improved our understanding of the genomic underpinnings of cancer

For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the serous subtype of ovarian cancer on a molecular level, suggesting that despite arising from different tissues in the body, these subtypes may share a common path of development and respond to similar therapeutic strategies.

TUMOR SUBTYPES — Revolutionized how cancer is classified

TCGA revolutionized how cancer is classified by identifying tumor subtypes with distinct sets of genomic alterations.*

THERAPEUTIC TARGETS — Identified genomic characteristics of tumors that can be targeted with currently available therapies or used to help with drug development

TCGA's identification of targetable genomic alterations in lung squamous cell carcinoma led to NCI's Lung-MAP Trial, which will treat patients based on the specific genomic changes in their tumor.

**THE TEAM**

**20** COLLABORATING INSTITUTIONS across the United States and Canada

**WHAT'S NEXT?**

The Genomic Data Commons (GDC) houses TCGA and other NCI-generated data sets for scientists to access from anywhere. The GDC also has many expanded capabilities that will allow researchers to answer more clinically relevant questions with increased ease.

*TCGA's analysis of stomach cancer revealed that it is not a single disease, but a disease composed of four subtypes, including a new subtype characterized by infection with Epstein-Barr virus.

www.cancer.gov/ccg

Translational Breast Cancer Research, 2016

# NCI Genomic Data Commons (GDC)

- A product of the NCI Center for Cancer Genomics (CCG)

- Mission: to provide the cancer research community with a unified data repository that enables data sharing across cancer genomic studies in support of precision medicine

- Associated projects: TCGA, Therapeutically Applicable Research to Generate Effective Treatments (TARGET) initiative, and Cancer Genome Characterization Initiative (CGCI)



https://gdc.cancer.gov/

Translational Breast Cancer Research, 2016

# GDC: TCGA breast cancer



TCGA-BRCA

⬇ Download Manifest  ⬇ Download Clinical  ⬇ Download Biospecimen

## Summary

| Project ID | TCGA-BRCA |
|---|---|
| Project Name | Breast Invasive Carcinoma |
| Disease Type | Breast Invasive Carcinoma |
| Primary Site | Breast |
| Program | TCGA |

CASES
1,098

FILES
25,970

ANNOTATIONS
78

### Case and File Counts by Experimental Strategy

| Experimental Strategy | Cases | Files |
|---|---|---|
| ■ Genotyping Array | 1,096 | 4,446 |
| ■ WXS | 1,050 | 10,820 |
| ■ RNA-Seq | 1,092 | 4,888 |
| ■ miRNA-Seq | 1,079 | 3,621 |

### Case and File Counts by Data Category

| Data Category | Cases | Files |
|---|---|---|
| ■ Raw Sequencing Data | 1,098 | 4,604 |
| ■ Transcriptome Profiling | 1,097 | 6,080 |
| ■ Simple Nucleotide Variation | 1,044 | 8,645 |
| ■ Copy Number Variation | 1,096 | 4,446 |
| ■ Clinical | 1,097 | 1,097 |
| ■ Biospecimen | 1,098 | 1,098 |

https://gdc-portal.nci.nih.gov/projects/TCGA-BRCA

Translational Breast Cancer Research, 2016

# NCI Genomic Data Commons



https://gdc.cancer.gov/

Translational Breast Cancer Research, 2016

# cBioPortal for Cancer Genomics

- http://www.cbioportal.org/

- Visualization, analysis and download of large-scale cancer genomics data sets

- 147 cancer genomics studies as of October 2016

- References

    - Gao et al., Sci Signal 2013

    - Cerami et al. Cancer Discov, 2012

# cBioPortal: TCGA breast cancer overview



Translational Breast Cancer Research, 2016

# cBioPortal: query interface



Select cancer study

Select genomic profiles

Select patient/case set

Enter gene set

Translational Breast Cancer Research, 2016

# cBioPortal: oncoprint

- Compact visualization of distinct genomic alterations, including somatic mutations, copy number alterations, and gene expression changes across a set of cases.

ERBB2

p53 signaling



Translational Breast Cancer Research, 2016

# cBioPortal: mutual exclusivity



Mutual exclusivity => functional link
- Alteration to the second gene within the same pathway offers no further selective advantage
- Alteration to the second gene within the same pathway leads to a disadvantage for the cell, *i.e.*, synthetic lethality.

*Ciriello et al., Genome Res, 2012*

Translational Breast Cancer Research, 2016

# cBioPortal: copy number vs mRNA expression



Translational Breast Cancer Research, 2016

# cBioPortal: mutations

# cBioPortal: survival

# cBioPortal: exploring the interactome

# International Cancer Genome Consortium (ICGC)

- To obtain a comprehensive description of genomic, transcriptomic and epigenomic changes in 50 different tumor types and/or subtypes which are of clinical and societal importance across the globe.



Translational Breast Cancer Research, 2016

http://www.icgc.org/

# ICGC Data Portal

# ICGC Data Portal: cancer projects



Translational Breast Cancer Research, 2016

# ICGC Data Portal: breast cancer



Translational Breast Cancer Research, 2016

# ICGC Data Portal: BRAF



Translational Breast Cancer Research, 2016

# ICGC Data Portal: BRAF mutations



Translational Breast Cancer Research, 2016

# ICGC Data Portal: BRAF targeting compounds



| MU672753 | chr7:g.140434586G>A | single base substitution | Intron: BRAF | 6 / 10,638 (0.06%) |
| MU1299736 | chr7:g.140481403C>T | single base substitution | Missense: BRAF G469R, G76R | 5 / 10,638 (0.05%) |
| | | | 3 UTR: BRAF | |
| MU831694 | chr7:g.140453193T>C | single base substitution | Missense: BRAF N581S, N188S | 5 / 10,638 (0.05%) |
| | | | Splice Region: BRAF | |

Showing 10 rows    <<<  <  **1**  2  3  4  5  >  >>>
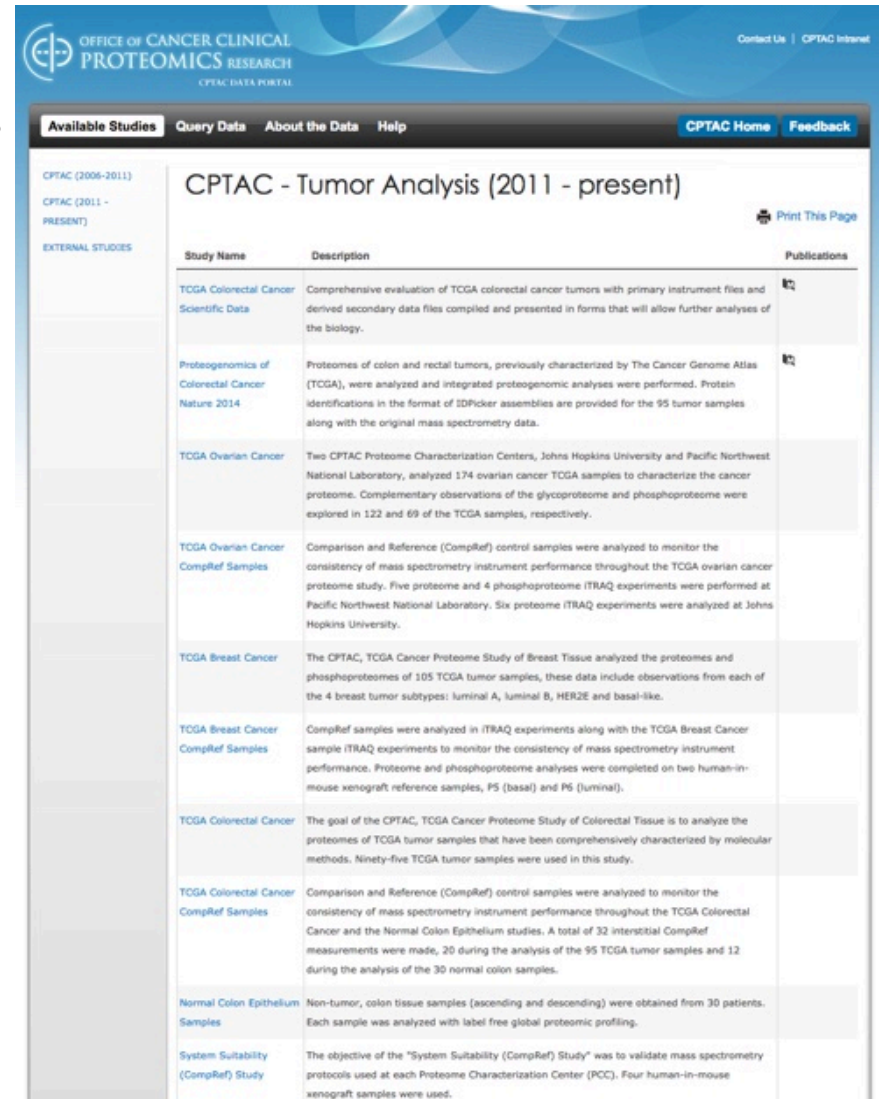
## Targeting Compounds

Showing **9** compounds

| Name ▲ | ATC Level 4 Description | Compound Class | # Clinical Trials |
|---|---|---|---|
| imatinib (ZINC000019632618) | Protein kinase inhibitors | FDA | 200 |
| llagate (ZINC000003872446) | -- | World | 0 |
| nilotinib (ZINC000006716957) | Protein kinase inhibitors | FDA | 67 |
| pazopanib (ZINC000011617039) | Protein kinase inhibitors | FDA | 132 |
| ruxolitinib (ZINC000043207851) | Protein kinase inhibitors | FDA | 26 |
| sorafenib (ZINC000001493878) | Protein kinase inhibitors | FDA | 410 |
| sprycel (ZINC000003986735) | Protein kinase inhibitors | FDA | 155 |
| stivarga (ZINC000006745272) | Protein kinase inhibitors | World | 48 |
| zelboraf (ZINC000052509366) | Protein kinase inhibitors | FDA | 67 |

Translational Breast Cancer Research, 2016

# Clinical Proteome Tumor Analysis Consortium (CPTAC)

- **Goals**
  - Global proteomic characterization of TCGA tumors
  - Proteogenomic data integration

- **Five centers established in 2011**
  - Broad Institute
  - John Hopkins University
  - Pacific Northwest National Laboratory
  - Washington University
  - Vanderbilt University

- **Tumor samples**
  - Breast (Broad and Wash U)
  - Colon and Rectal (Vanderbilt)
  - Ovarian (JHU and PNNL)

- **CPTAC data portal**
  - https://cptac-data-portal.georgetown.edu

# Clinical Proteome Tumor Analysis Consortium (CPTAC)

## Proteogenomic characterization of human colon and rectal cancer

Bing Zhang[1,2], Jing Wang[1], Xiaojing Wang[1], Jing Zhu[1], Qi Liu[1], Zhiao Shi[3,4], Matthew C. Chambers[1], Lisa J. Zimmerman[5,6], Kent F. Shaddox[6], Sangtae Kim[7], Sherri R. Davies[8], Sean Wang[9], Pei Wang[10], Christopher R. Kinsinger[11], Robert C. Rivers[11], Henry Rodriguez[11], R. Reid Townsend[8], Matthew J. C. Ellis[8], Steven A. Carr[12], David L. Tabb[1], Robert J. Coffey[13], Robbert J. C. Slebos[2,6], Daniel C. Liebler[5,6] & the NCI CPTAC*

**Nature, 2014**

## Proteogenomics connects somatic mutations to signalling in breast cancer

Philipp Mertins[1]*, D. R. Mani[1]*, Kelly V. Ruggles[2]*, Michael A. Gillette[1,3]*, Karl R. Clauser[1], Pei Wang[4], Xianlong Wang[5], Jana W. Qiao[1], Song Cao[6], Francesca Petralia[4], Emily Kawaler[2], Filip Mundt[1,7], Karsten Krug[1], Zhidong Tu[4], Jonathan T. Lei[8], Michael L. Gatza[9], Matthew Wilkerson[9], Charles M. Perou[9], Venkata Yellapantula[6], Kuan-lin Huang[6], Chenwei Lin[5], Michael D. McLellan[6], Ping Yan[5], Sherri R. Davies[10], R. Reid Townsend[10], Steven J. Skates[11], Jing Wang[12], Bing Zhang[12], Christopher R. Kinsinger[13], Mehdi Mesri[13], Henry Rodriguez[13], Li Ding[6], Amanda G. Paulovich[5], David Fenyö[2], Matthew J. Ellis[8], Steven A. Carr[1] & the NCI CPTAC†
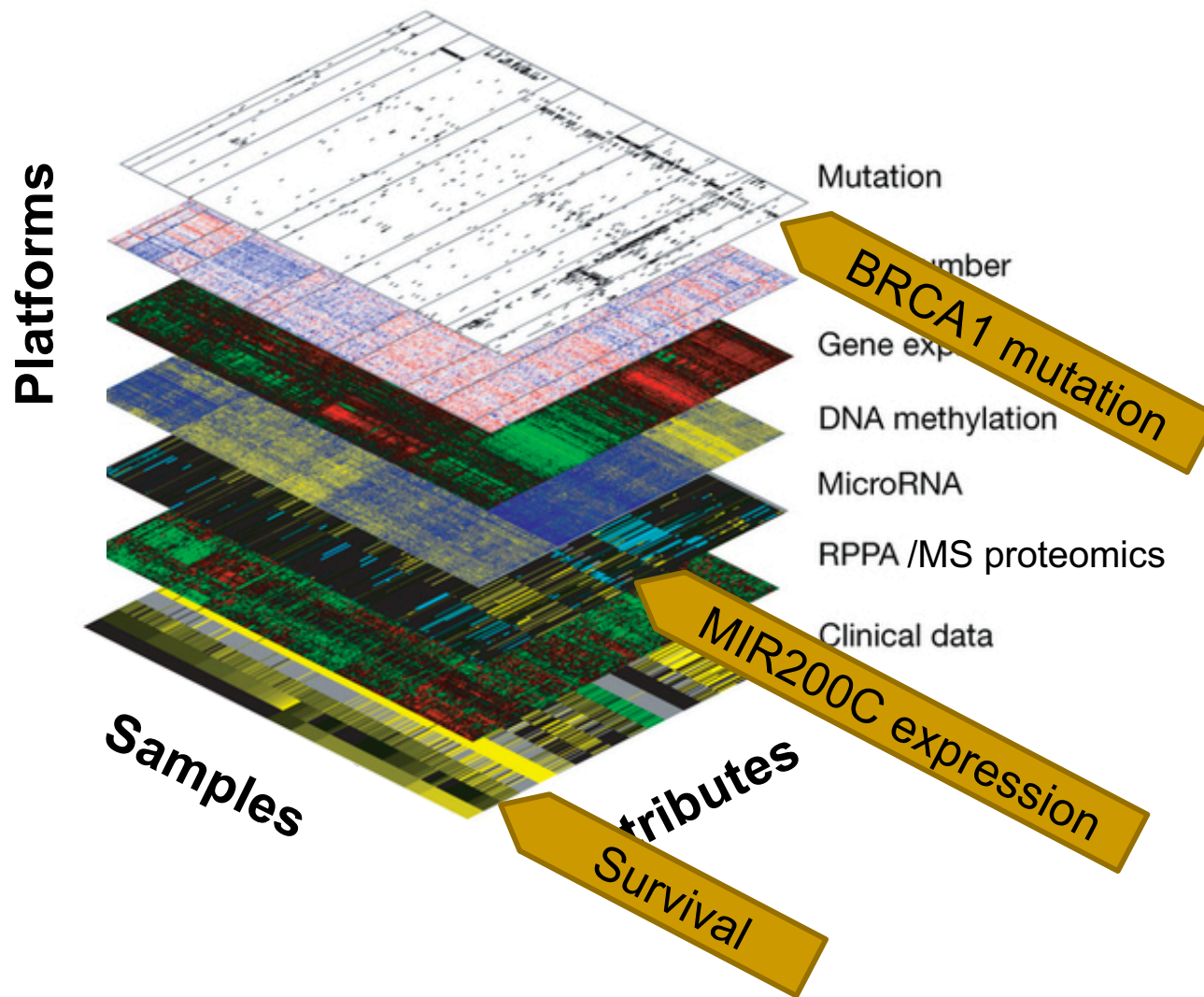
**Nature, 2016**

## Integrated Proteogenomic Characterization of Human High-Grade Serous Ovarian Cancer
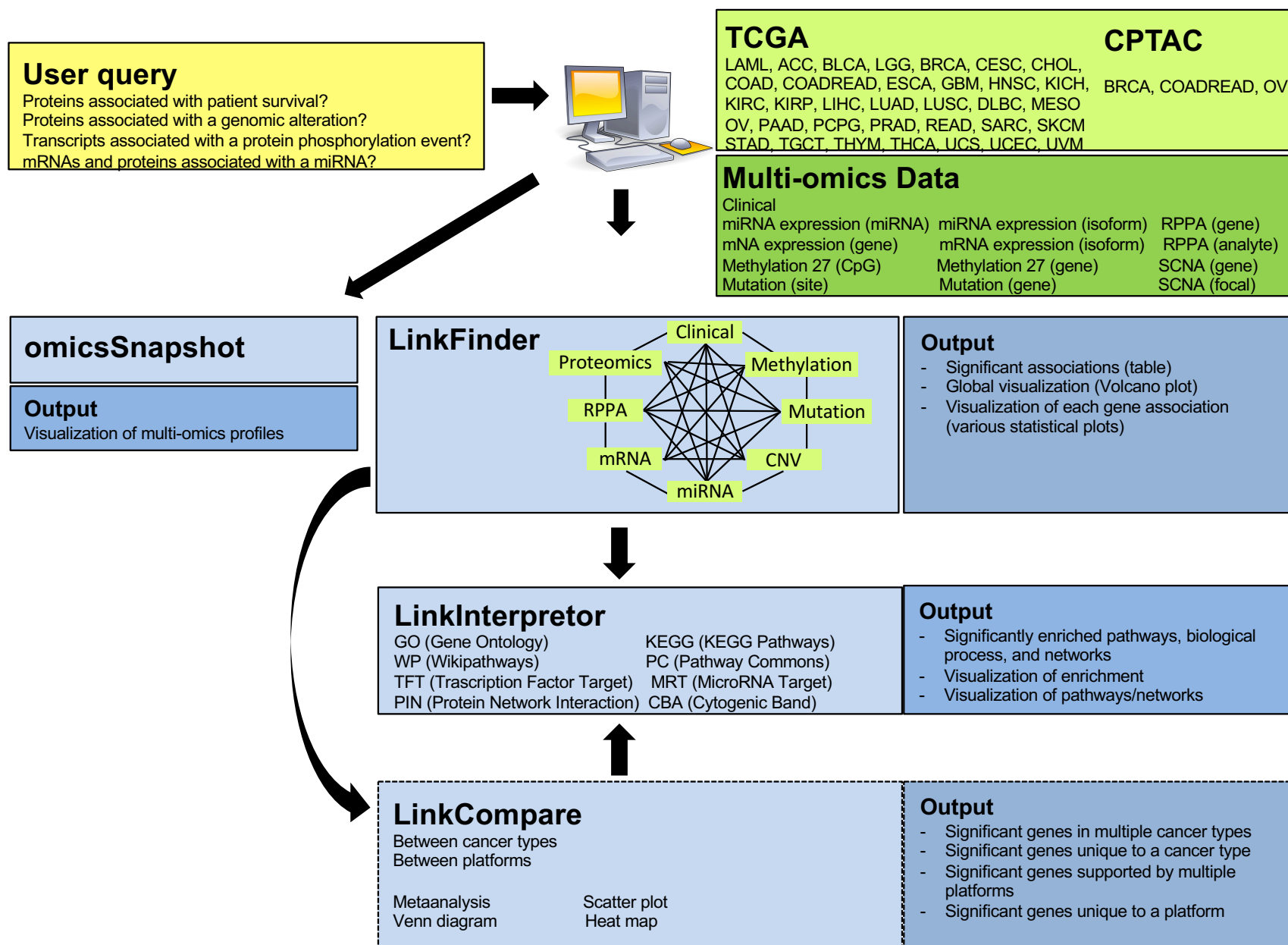
Hui Zhang,[1,15] Tao Liu,[2,15] Zhen Zhang,[1,15] Samuel H. Payne,[2,15] Bai Zhang,[1] Jason E. McDermott,[2] Jian-Ying Zhou,[1] Vladislav A. Petyuk,[2] Li Chen,[1] Debjit Ray,[2] Shisheng Sun,[1] Feng Yang,[2] Lijun Chen,[1] Jing Wang,[3] Punit Shah,[1] Seong Won Cha,[4] Paul Aiyetan,[1] Sunghee Woo,[4] Yuan Tian,[1] Marina A. Gritsenko,[2] Therese R. Clauss,[2] Caitlin Choi,[1] Matthew E. Monroe,[2] Stefani Thomas,[1] Song Nie,[2] Chaochao Wu,[2] Ronald J. Moore,[2] Kun-Hsing Yu,[5,6] David L. Tabb,[3] David Fenyö,[7] Vineet Bafna,[8] Yue Wang,[9] Henry Rodriguez,[10] Emily S. Boja,[10] Tara Hiltke,[10] Robert C. Rivers,[10] Lori Sokoll,[1] Heng Zhu,[1] Ie-Ming Shih,[11] Leslie Cope,[12] Akhilesh Pandey,[13] Bing Zhang,[3] Michael P. Snyder,[6] Douglas A. Levine,[14] Richard D. Smith,[2] Daniel W. Chan,[1,16,*] Karin D. Rodland,[2,16,*] and the CPTAC Investigators

**Cell, 2016**

# LinkedOmics: cross-omics association analysis

# LinkedOmics: cross-omics association analysis

# Become a superuser



*Graph courtesy of http://www.incogen.com/*

- **Algorithm developer**
  - Statisticians
  - Mathematicians
  - Computer scientists
- **Tool developer**
  - Bioinformaticians
- **Data provider/consumer**
  - Biologists

R
Linux
Python

Translational Breast Cancer Research, 2016