# BIOINFORMATICS

Introduction

David Fazekas

fazekas@netbiol.elte.hu

Department of Genetics (ELTE, HU)

Earlham Institute (UK)

# COURSE INFORMATION

# LECTURERS

- Eszter Ari - chief trainer and administration
- David Fazekas
- Zsuzsa Dosztányi

# TRAINERS

- Amanda Demeter
- Dániel Gerber
- Gábor Erdős

# COURSE MATERIAL

https://genetics.elte.hu

username: genetika2016
password: genetika2016

# SYLLABUS

1. Introduction
2. Data sources
3. Sequence alignment
4. Sequence databases and searching
5. Molecular phylogenetics I.
6. Molecular phylogenetics II.
7. Genomics and transcriptomics I.
   <span style="color:blue">autumn break</span>
8. Genomics and transcriptomics II.
9. Network and systems biology I. - <span style="color:red">exam I. (lecture 1-6)</span>
10. Network and systems biology II.
11. Network and systems biology III.
12. Protein structure bioinformatics
13. <span style="color:red">exam I. (lecture 7-12)</span>

# EXAM

### Lecture

- 2 written exam during the semester
- Average of those
- If either is 1, oral exam is required

### Practice

- Maximum 3 absenteeism
- Submit a project
- Group of 3 student

# CHOOSE ONE

- RASK (KRAS)
- ERK1 (MAPK3)
- JAK1
- IGF1R
- GSK3B
- AXIN1
- SMAD2
- NOTCH1

# WHAT IS BIOINFORMATICS?

# Definition: Bioinformatics

*"Research, or application of computational tools and approaches for expanding the use of biological, medical, behavioral or health data, including those to acquire, store, organize, archive, analyze, or visualize such data." "Bioinformatics applies principles of information sciences and technologies to make the vast, diverse, and complex life sciences data more understandable and useful."*

# Definition: Computational Biology

*"The development and application of data-analytical and theoretical methods, mathematical modeling and computational simulation techniques to the study of biological, behavioral, and social systems."*

*"Computational biology uses mathematical and computational approaches to address theoretical and experimental questions in biology."*

# WHAT IS BIOINFORMATICS?
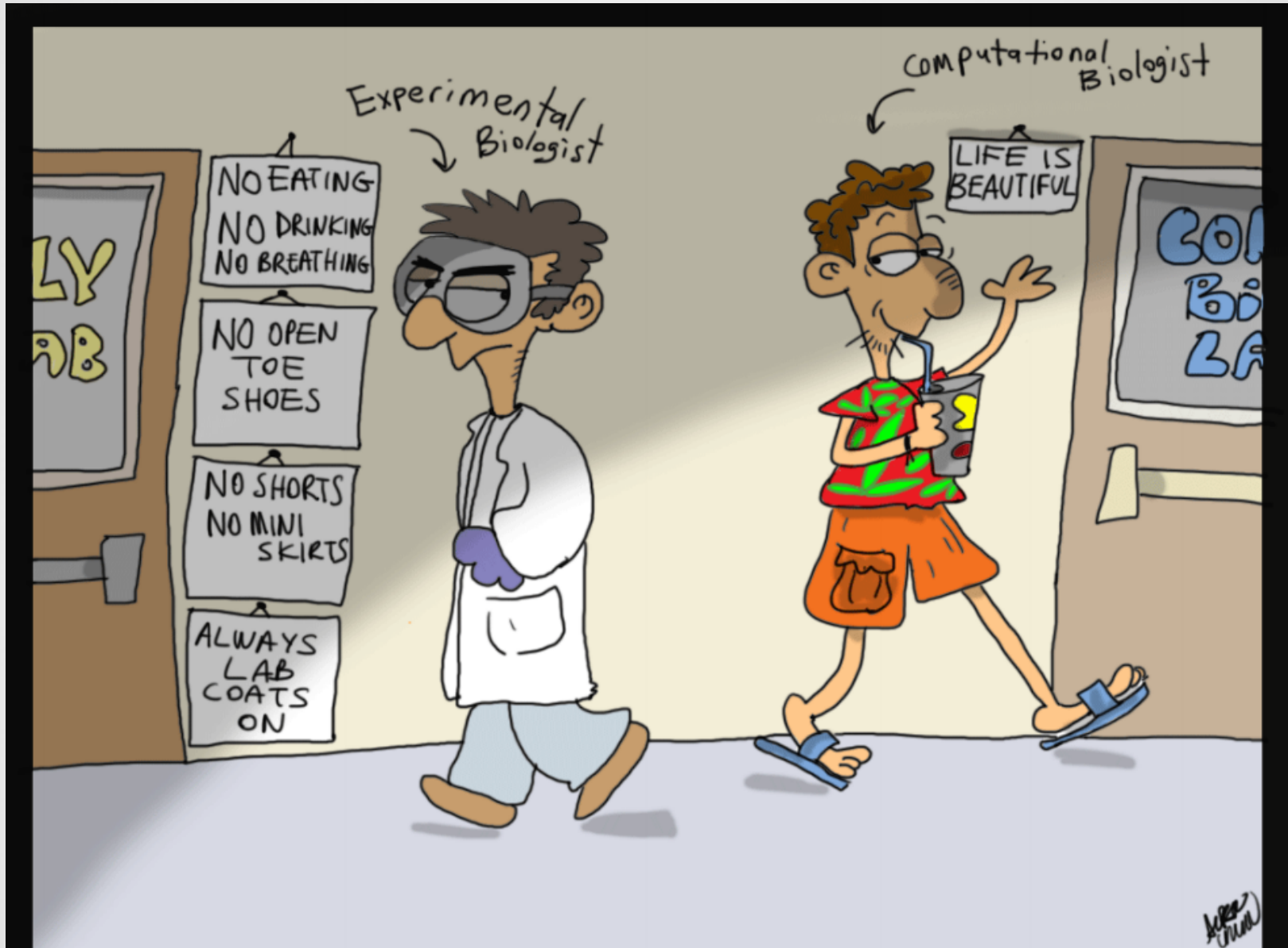
## IN A NARROWER SENSE

- Working with data in life sciences

## IN THE BROADER SENSE

- Molecular bioinformatics
- Sequence and structure of macro molecules
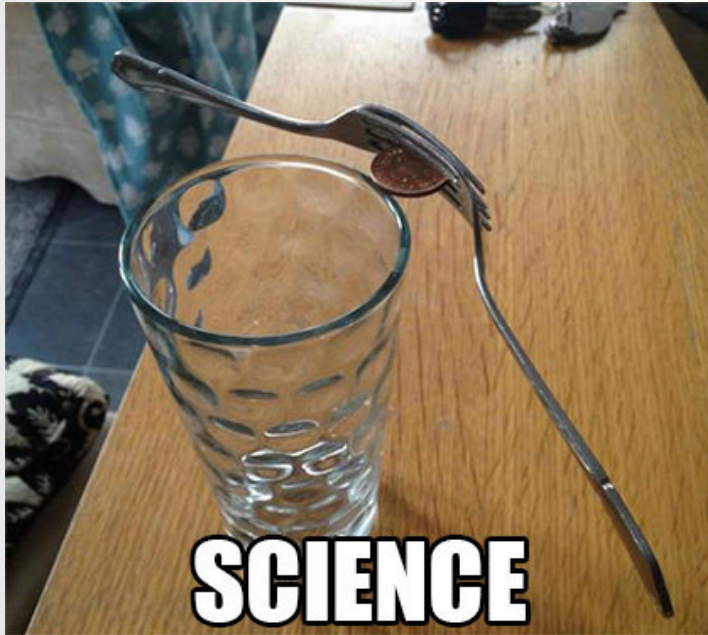- Annotations
- Network biology

# MOLECULAR BIOINFORMATICS



Central dogma of molecular biology

DNA → RNA → protein → cellular phenotype

Central dogma of genomics

genome → transcriptome → proteome → cellular phenotype

DNA → RNA → protein

# WET LAB - DRY LAB

# DATA SCIENCE?

# SCIENCE VS ENGINEERING

Science

Engineer

Bioinformatics

Computational
biology

Data
sience

Business

# BIG DATA

"*Big data is like teenage sex;*
*everyone talks about it,*
*nobody really knows how to do it,*
*everyone thinks everyone else is doing it,*
*so everyone claims they are doing it*".

Dan Ariely, Duke University

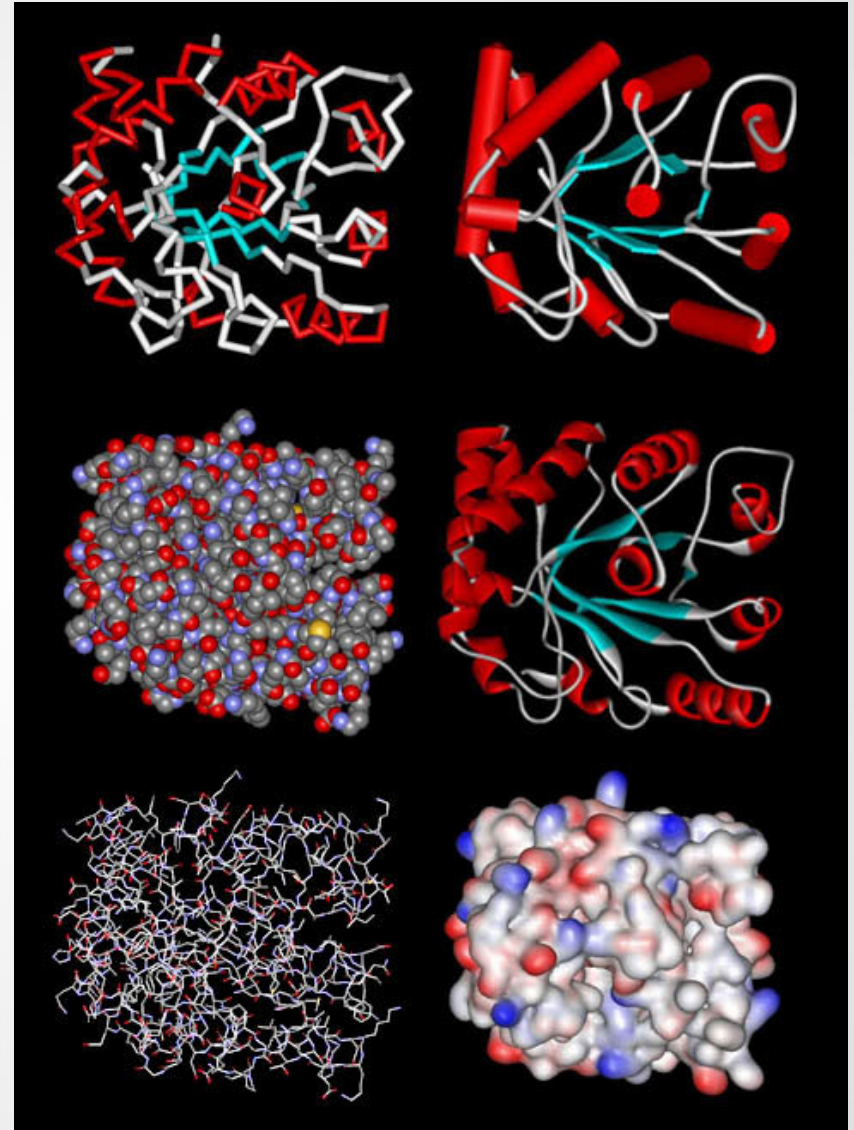# THE SUBJECT OF BIOINFORMATICS

Classic bioinformatics

# CLASSIC BIOINFORMATICS

- Sequence alignment
- Statistical analysis (e.g. CG ratio, gene length)
- Genome annotation:
  - ORF, gene prediction
  - promoter analysis
- Sequence database
- Sequence searching

# CLASSIC BIOINFORMATICS

- 3D structure of macro molecules
- Protein docking

# CLASSIC BIOINFORMATICS

- Molecular phylogenetics

# "OLD" AND "NEW" BIOLOGY

In the (near) past:

For researchers the greatest the challenge was to produce good quality data.

Today:

The biggest challenge for researchers to interpret a massive set of data is because biological data collection is done in bulk, has become industry-standard.

# MODERN BIOINFORMATICS

- Post-genome era
- Next gen sequencing
- Comparative genomics
- Transcriptomics
- Proteomics
- Systems biology

# THE SUBJECT OF BIOINFORMATICS

Data

# NCBI GENEBANK: $2.6*10^{12}$ BASE



Bases

# SEQUENCED GENOMES: 72000 ORGANISM

**Organisms**



Legend:
- Archaea
- Bacteria
- Complete
- Fungi
- Invertebrate
- Mitochondrion
- Plant
- Plasmid
- Plastid
- Protozoa
- Vertebrate Mammalian
- Vertebrate Other
- Viral

Y-axis: 100,000 / 10,000 / 1,000 / 100 / 10

X-axis: 2005 / 2010 / 2015

"And that's why we need a computer."

GOFF

# TOOLS WE'LL USE DURING THIS COURSE

- Database querying, searching
- Data managing in text files and table
- Web services
- Graphical programs
- 3D structure modelling
- Network analysis and visualization tools
- ...
- LINUX

OMG!

# Linux  Windows  Mac

as seen by...

Mac
Fanboys

Windows
Fanboys

Linux
Fanboys

# LINUX COMMAND LINE

# LINUX

- Distributions
- Package managers
- Kernel
- Shell
- File system
- Graphical user interface
- Remote access

# PCLAB

# DATA STORING

# SEQUENCE - FLAT FILE

# SEQUENCE - FASTA FILE

## Homo sapiens chromosome 16, GRCh38 Primary Assembly

NCBI Reference Sequence: NC_000016.10

GenBank    Graphics

```
>gi|568815582:c355241-287440 Homo sapiens chromosome 16, GRCh38 Primary Assembly
CTCCAGGCTTTCTGACCCCCTTCCTGTCCCTGCAGGGCAGGGGCCCCGAGGAGCTGGAGGCCCAAGGGCC
TTGTTCTGGTCCCAGGGCCTGGGGACACCTGCCACCGAGGGCCAAGAGGAGGATGGACGTGGACACAGCC
CCGAGAGCCTGGCCCGGACACAGCAGCAGCTCTGCAGAACCAAGGCAAGCATTGGGGACCTTGTTGGGAG
TCGGGGGGCAGCCCAGGGGCCAGTGTCTGAGGTCCTGCTGTCTGTCTGGCCACCCAGGACTCCCTCATCC
CTGGAAACTGTGCTTTACCATGGAGGCCACCCACTCTGTCTCCTCTAAGGTTCTGAGGCTGAATGGGCTA
GGGGGCTTGCGGGGAGGCCCCAGTGTCCAGCACTGTGGGACCTGGCAGGGTGCCTGCGGCCAGGACCCAG
CGGGGCCAGGTGTTGGTCTAACAGTGCAGCTTCGTTCATATCCCCAGCCCCTGCCCACCTGCTCTGAGCA
CAGTGATGGCCCTGGGAGGTGGGCCTGGGCCCTTGGCAGGCTGGGGACAGCCTAGTGGCCCTTGTCCCAT
GCTACCCCCTTTCCCACACAGCGATGCTGGCATCAGACACCATGCTGAGTGCTGGCAGGGGCGAGGGCTG
GGAGGCTTCCACACATGGTTCCCCATGCAGTCCCACCTGTGGGCATCTGGTTGGGGGTAGGCTGGAAGCT
CGGGGAGCCTGGAGCTGGGACTTCTGTGCTTGCCTGGGAGCTCTGAAGGGTGAGGCTGGGCATCCAGGGT
GACACAGCCCAGGGAAAGACATGGGGGTGACGTGAGAGGTGCCTGGAGGGAGCTGGCAGGTATGGACATG
ATGGACACGGAAGCACGGAGGCGGGCAAGTGGCCAGACGCATCTAGGGGAAGGTGTGGGGGAGGCGCCCT
TAGGAAGGGCCCTGAAGGGCTGTGGGCAGCAGGGAGCCCTGGGAGGCCTAAAGCAGAGGGCAGATAGATG
AGGGCTGTTGTCTCAGGTGTGGGGCGGACGAGGTGGGGAGCCCTCACCCAACAGGAGGCGAGCTGGTCCT
GTGTGGCCTGAACTGCAGCTGTCTCCTCTGTGAAACGGGGGTATAGCTGACCCCAGGGGGCTGCCTGGAG
CATCCCGGGAGGTGCCAGGCCCAATAGTGCTCTGGGAAGGACAGGGCCCTGGGCTGTTGTGGGAGGCGGC
AGATCCTGGTACTCACATCCTCCTCCTTGGGGAGGGCCTGATGGTTGGCTGAGGCCTGGGTGGAGAGCAG
AGGGTTGGTTCTGACAGGGTTGGGCTGGCCAGAGCTGGTGCTGGGGCTGCTGCTGGGGGCCCGTGCCTCT
CTGCCGTGGGGTGCCTGGGGCTGTGACCTCATGCTCTGTGGCCTGCAGGGCAAGTGACACGGATCTGGGC
AGCCAGGGTGGCAGGATCGGACTGGACCCCTTGGCAGGGCCGCTGTGGAGACAGCCCAGGGGAAGGGGTG
```
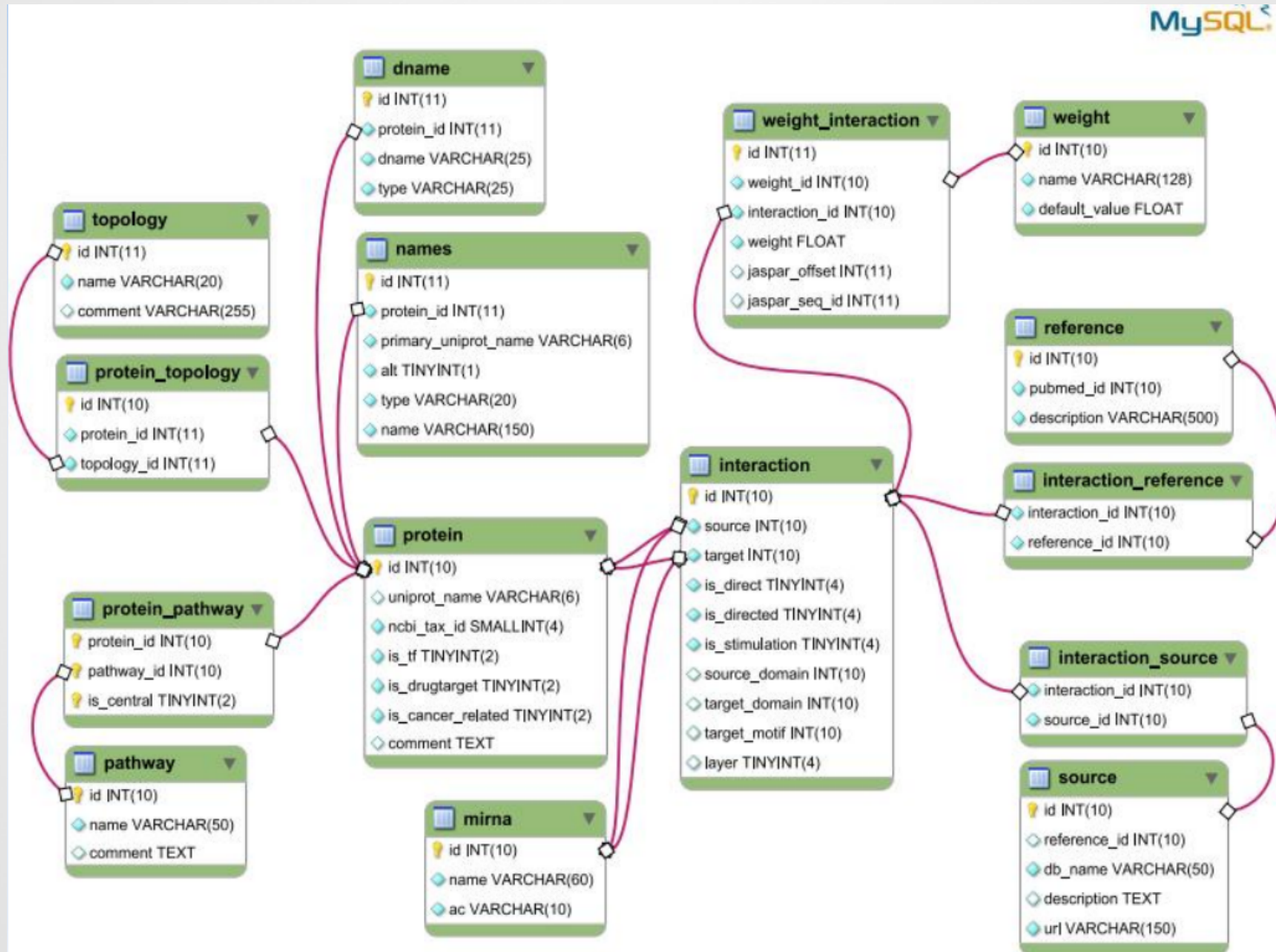
# SEQUENCE - WITH ANNOTATIONS

# NETWORK - PLAIN TEXT

# NETWORK - TABLE

| source_na | source_ur | source_sp | source_sp | source_to | source_pa | target_na | target_un | target_sp | target_sp | target_to | target_pa | layer | interactio |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| JAK2 | O60674 | ENSG0000 | H. sapiens | Mediator | JAK/STAT( | PTPN11 | Q06124 | ENSG0000 | H. sapiens | Co-factor, | RTK(non-c | Interactio | PPI direct( |
| PTPN11 | Q06124 | ENSG0000 | H. sapiens | Co-factor, | RTK(non-c | JAK2 | O60674 | ENSG0000 | H. sapiens | Mediator | JAK/STAT( | Interactio | PPI direct( |
| IRS1 | P35568 | ENSG0000 | H. sapiens | Mediator, | RTK(core) | JAK1 | P23458 | ENSG0000 | H. sapiens | Mediator | RTK(core) | Interactio | PPI direct( |
| JAK1 | P23458 | ENSG0000 | H. sapiens | Mediator | RTK(core) | IRS1 | P35568 | ENSG0000 | H. sapiens | Mediator, | RTK(core) | Interactio | PPI direct( |
| GSK3B | P49841 | ENSG0000 | H. sapiens | Mediator, | RTK(non-c | AXIN1 | O15169 | ENSG0000 | H. sapiens | Mediator, | RTK(non-c | Interactio | PPI direct( |
| AXIN1 | O15169 | ENSG0000 | H. sapiens | Mediator, | RTK(non-c | GSK3B | P49841 | ENSG0000 | H. sapiens | Mediator, | RTK(non-c | Interactio | PPI direct( |
| MAP2K1 | Q02750 | ENSG0000 | H. sapiens | | RTK(core) | MAPK3 | P27361 | ENSG0000 | H. sapiens | Mediator | RTK(core) | Interactio | PPI direct( |
| MAPK3 | P27361 | ENSG0000 | H. sapiens | Mediator | RTK(core) | MAP2K1 | Q02750 | ENSG0000 | H. sapiens | | RTK(core) | Interactio | PPI direct( |
| SMAD3 | P84022 | ENSG0000 | H. sapiens | Mediator, | RTK(core) | ESR1 | P03372 | ENSG0000 | H. sapiens | Receptor, | NHR(core) | Interactio | PPI direct( |
| ESR1 | P03372 | ENSG0000 | H. sapiens | Receptor, | NHR(core) | SMAD3 | P84022 | ENSG0000 | H. sapiens | Mediator, | RTK(core) | Interactio | PPI direct( |
| PEA15 | Q15121 | ENSG0000 | H. sapiens | Co-factor | RTK(non-c | MAPK3 | P27361 | ENSG0000 | H. sapiens | Mediator | RTK(core) | Interactio | PPI direct( |

# NETWORK - RELATIONAL DATABASE

# NETWORK - VISUALIZATION