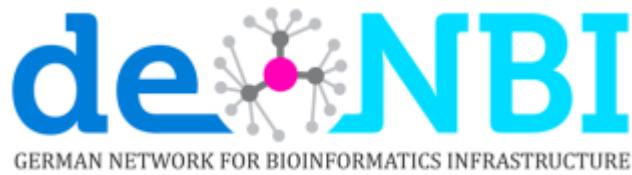


# Introduction to RNA-Seq data analysis with Galaxy

Markus Wolfien and Andrea Bagnacani

*de.NBI Training – 7<sup>th</sup> March 2018 Kiel*

[www.sbi.uni-rostock.de](http://www.sbi.uni-rostock.de)



# Making sense out of data – providing meaning to models



SYSTEMS BIOLOGY  
BIOINFORMATICS  
ROSTOCK



*Omic analyses  
and data  
integration*

*Data management  
& standardisation  
(Fairdom partner)*



*eHealth  
iOS Application*

SYSTEMS BIOLOGY  
BIOINFORMATICS  
ROSTOCK

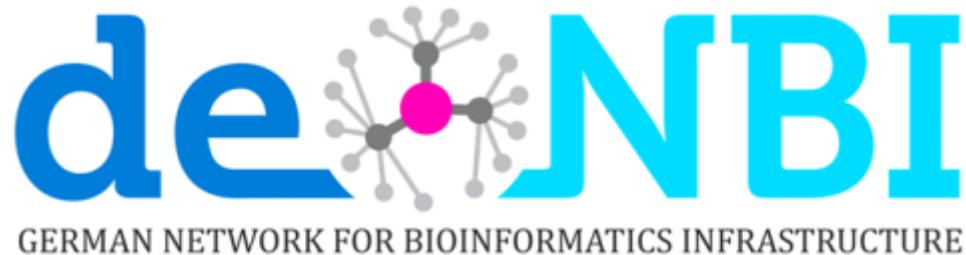
*Computational  
modeling &  
machine learning*



*Summer schools  
and workshops*

*Systems Medicine  
Pre-clinical trials*



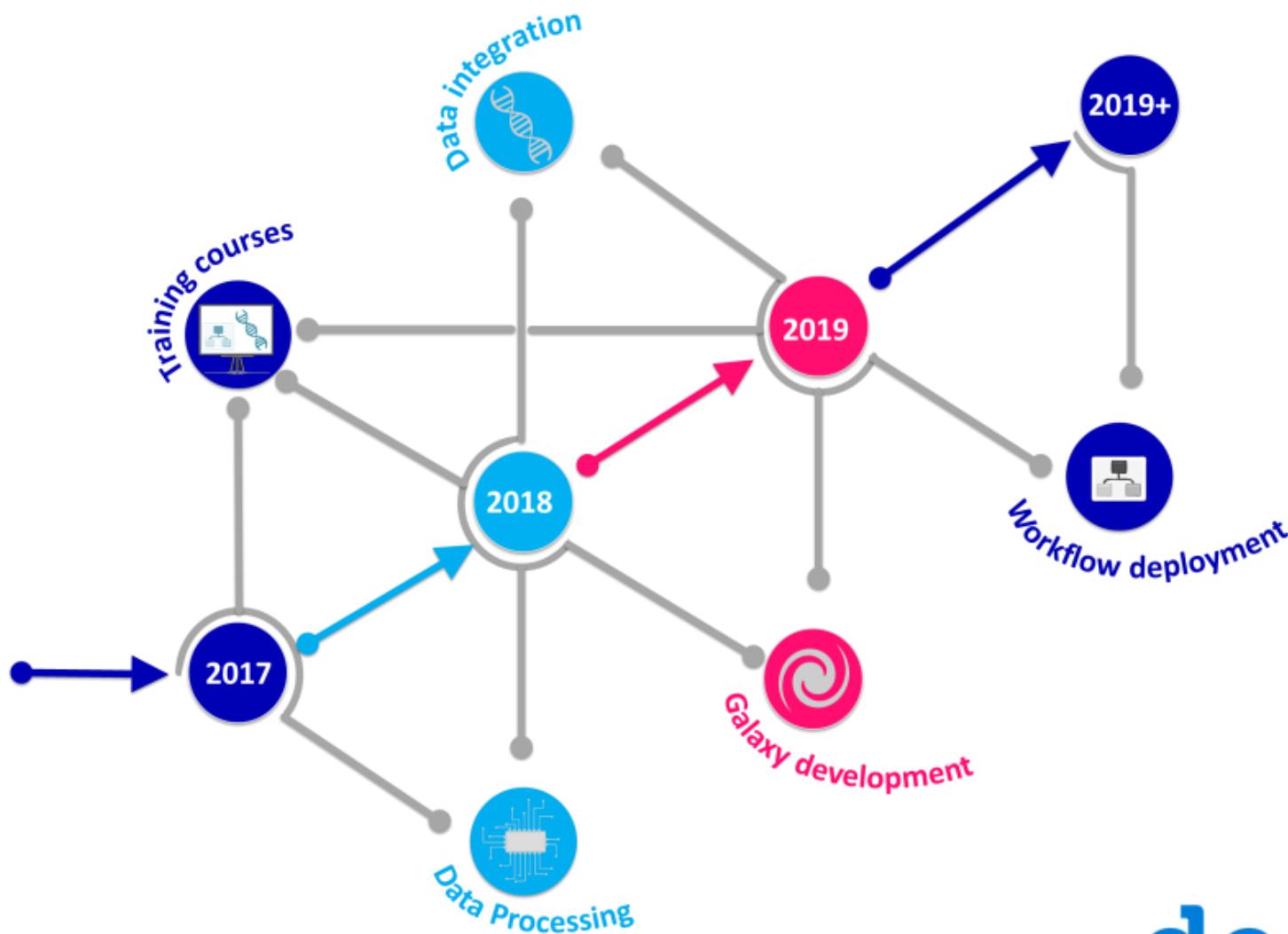


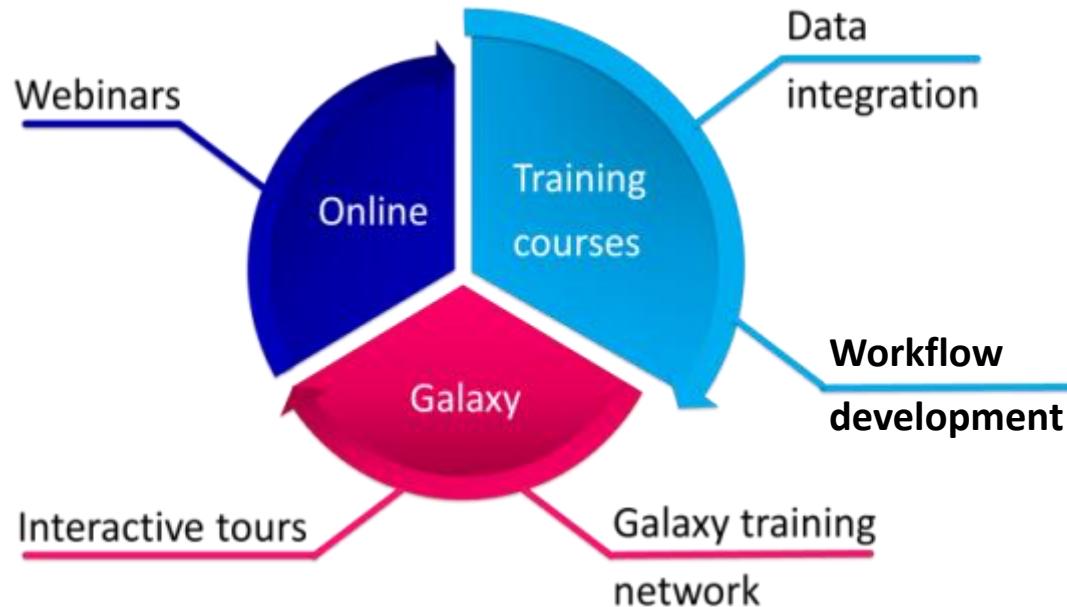
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE



***Structured Analysis and Integration of  
RNA-Seq experiments (de.STAIR)***

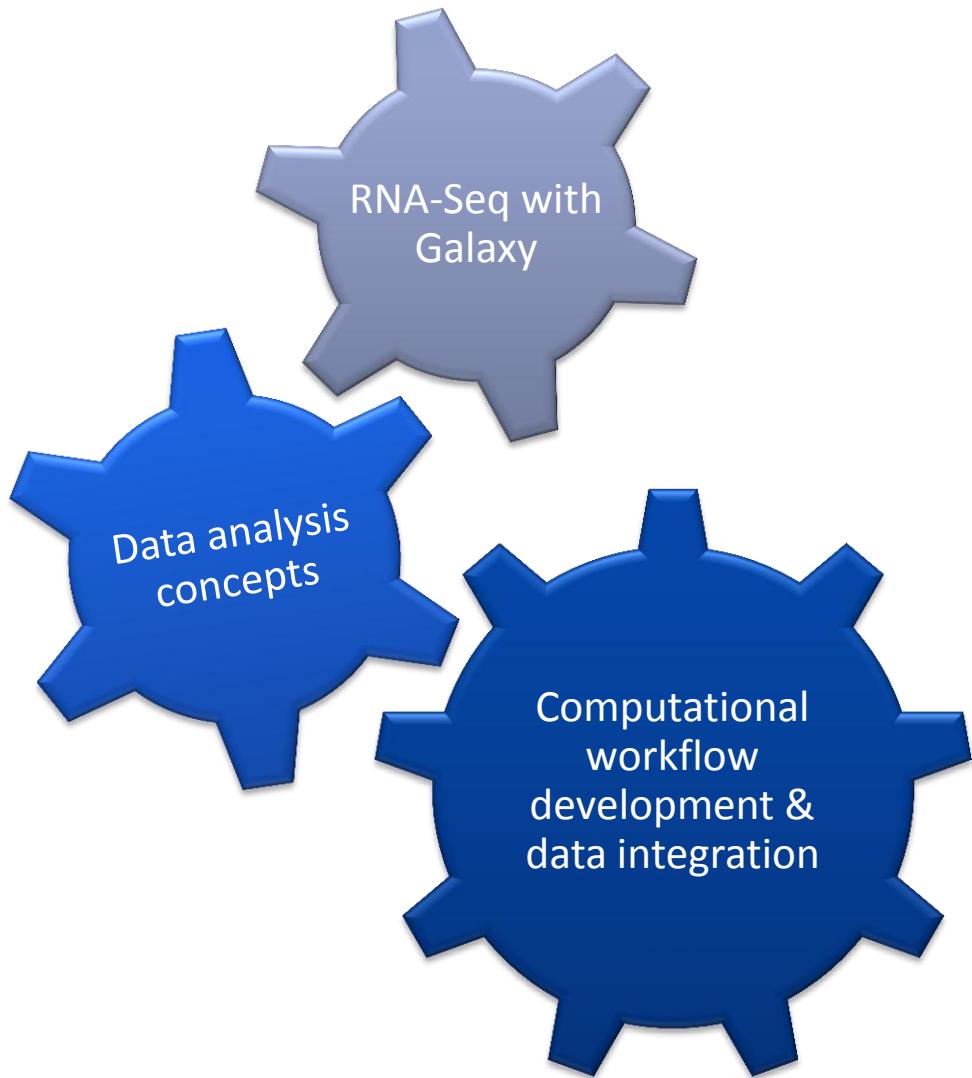
Our aim is to enable a comprehensive **analysis of RNA-Seq experiments as a service**. To enable maximum usefulness, interconnectivity, and accessibility for the developed approaches and services, we will provide dedicated **workshops, training programs and screen casts** for bioinformaticians and other life scientists.





# Objectives for the training today

- What is medical Big data?
- Why using NGS?
- The Galaxy around me
  
- How do we analyze NGS data?
- Are there best practices?
- Are there automation strategies?
  
- Explore a use case in the medical field



# What is big data?

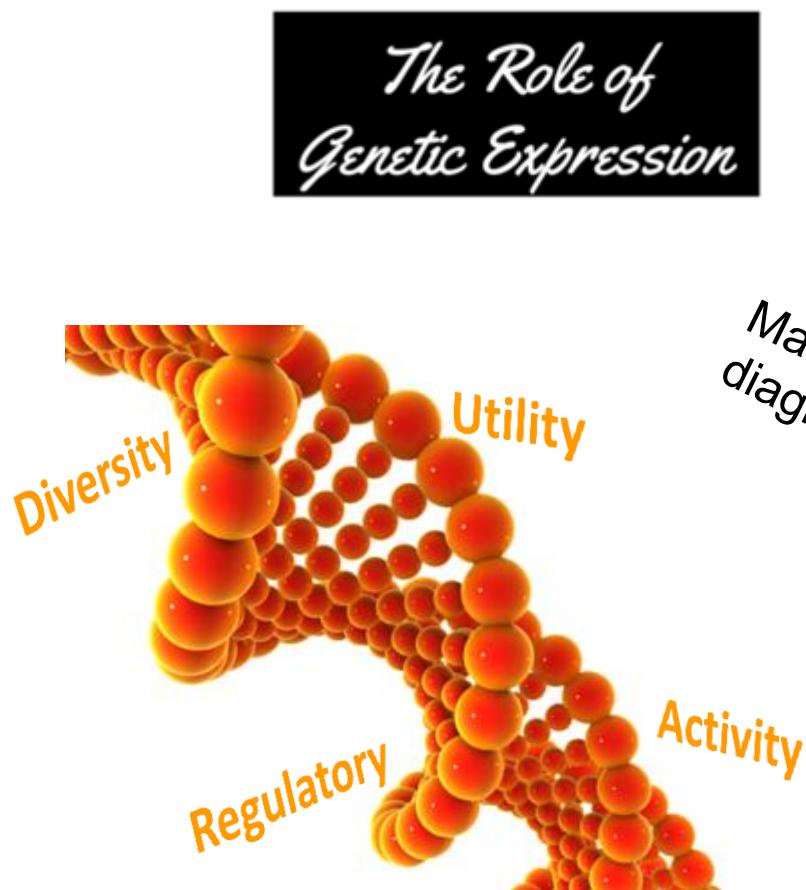
“The basic idea behind the phrase “Big Data” is that everything we do is increasingly leaving a digital trace (or data), which we (and others) can use and analyze. Big data therefore refers to our ability to make use of the ever-increasing volumes of data.”

*Bernhard Marr, Big Data*



“From the dawn of civilization until 2003, humankind generated five exabytes of data. NOW we produce five exabytes every two days ... and pace is accelerating.”

*Eric Schmidt, Executive Chairman Google*



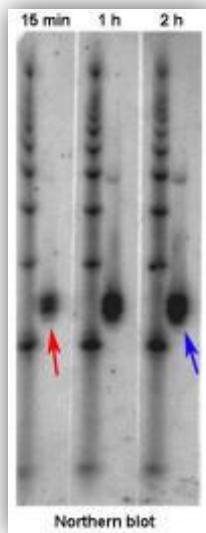
Many different variations and subtypes

Information about regulatory mechanisms

Active and measurable state of the cell ...

... ,but only a snapshot  
Many different therapeutical and  
diagnostical approaches

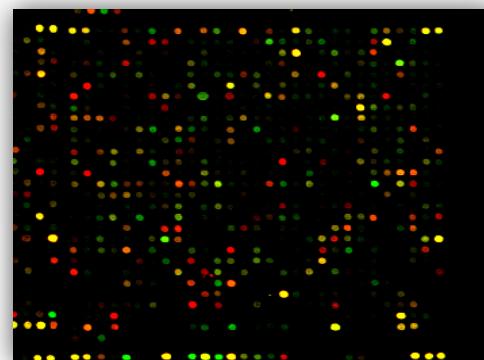
# Measuring gene expression



Northern Blot



Reverse Transcription PCR

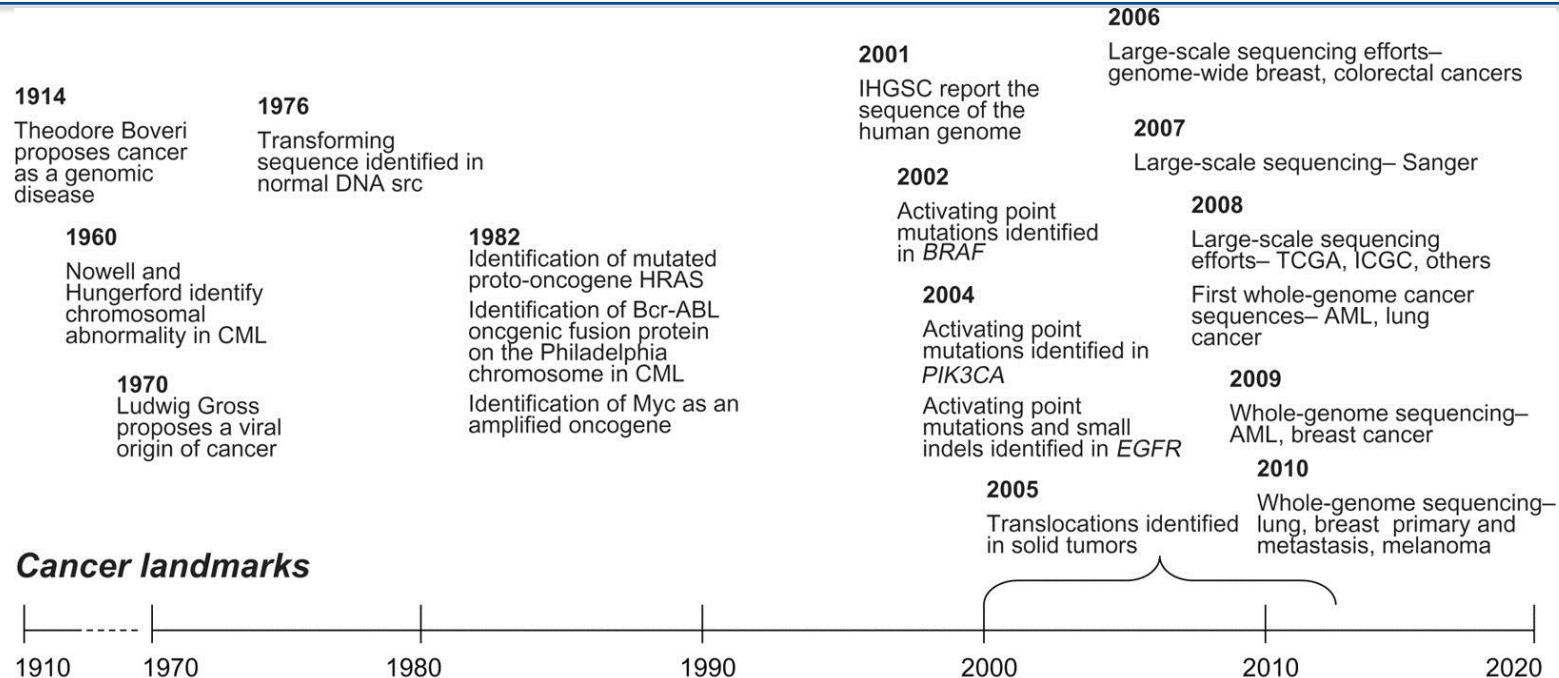


Microarrays



NGS

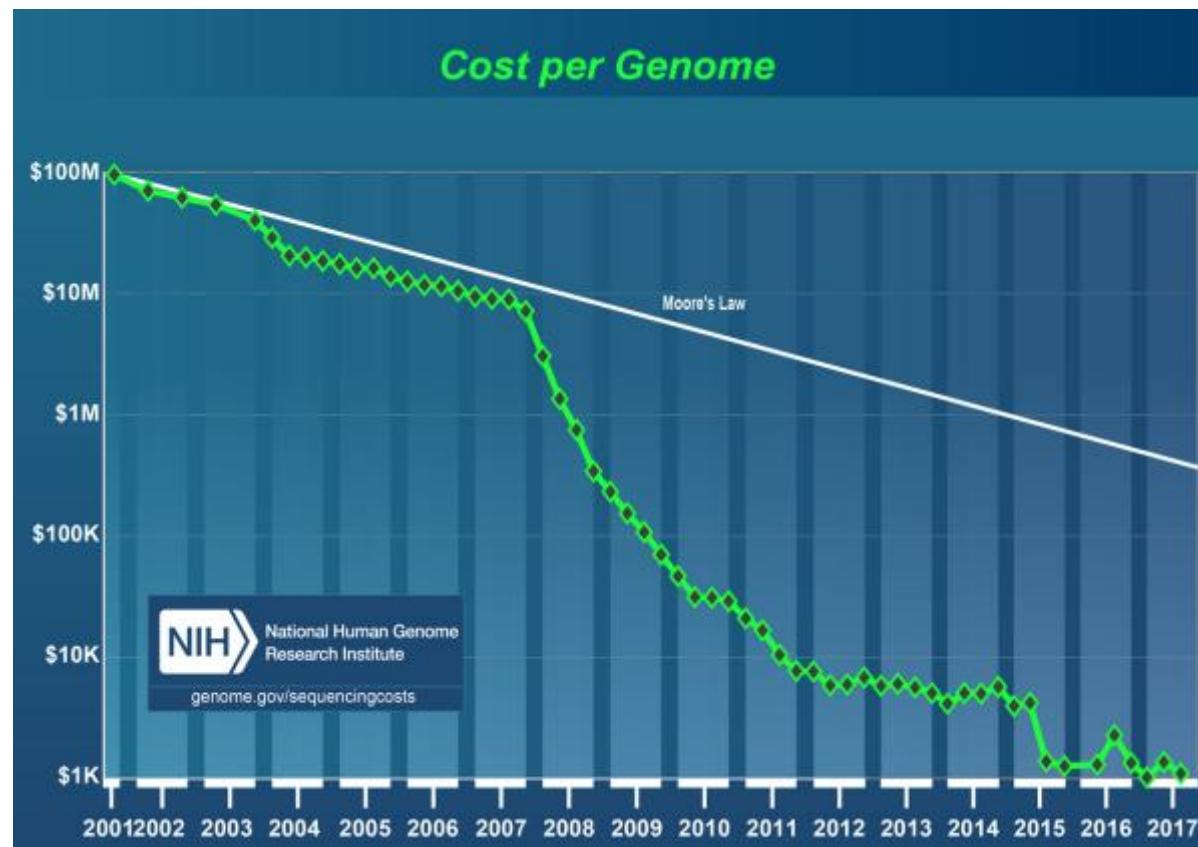
# Evolution of Sequencing technologies

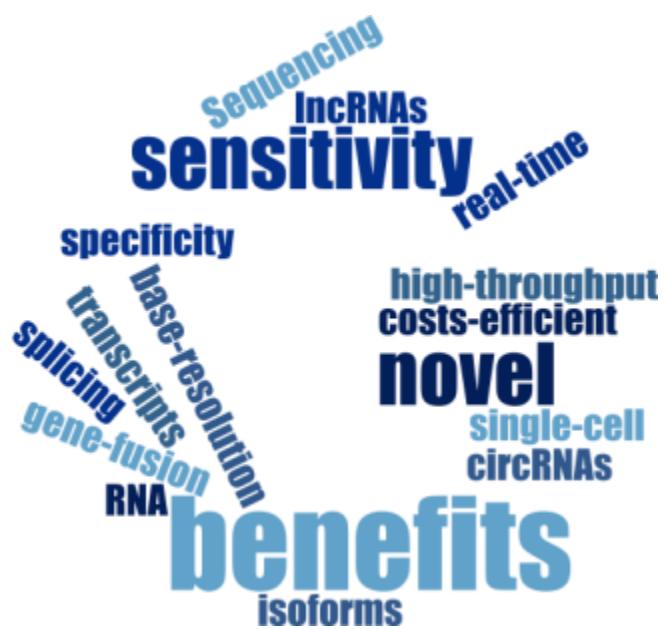


## Technological advances

<b>1972</b> First reported recombinant DNA technologies	<b>1982</b> Archetypes of cancer alterations defined	<b>1994</b> Microarrays for gene expression and sequence analysis	<b>2005</b> Next-generation sequencing: massively parallel sequencing-by-synthesis multiplex polony sequencing four-color DNA sequencing-by-synthesis
<b>1975</b> Sanger reports DNA sequencing method	<b>1986-7</b> CalTech reports first semiautomated DNA sequencing machine	<b>1995</b> Mathies et al. reports high-throughput dye-based DNA sequencing	<b>2007</b> Integrative analytic approaches for multiple types of large datasets
<b>1977</b> Maxam and Gilbert report DNA sequencing method		<b>1998</b> RNAi screening to specify gene function Mass-spectrometric genotyping of SNPs	<b>2008</b> Single-molecule DNA sequencing
	Bachetti and Graham report method for DNA transfer		<b>2010</b> Single-molecule real-time DNA sequencing

# Technical advances lead the way

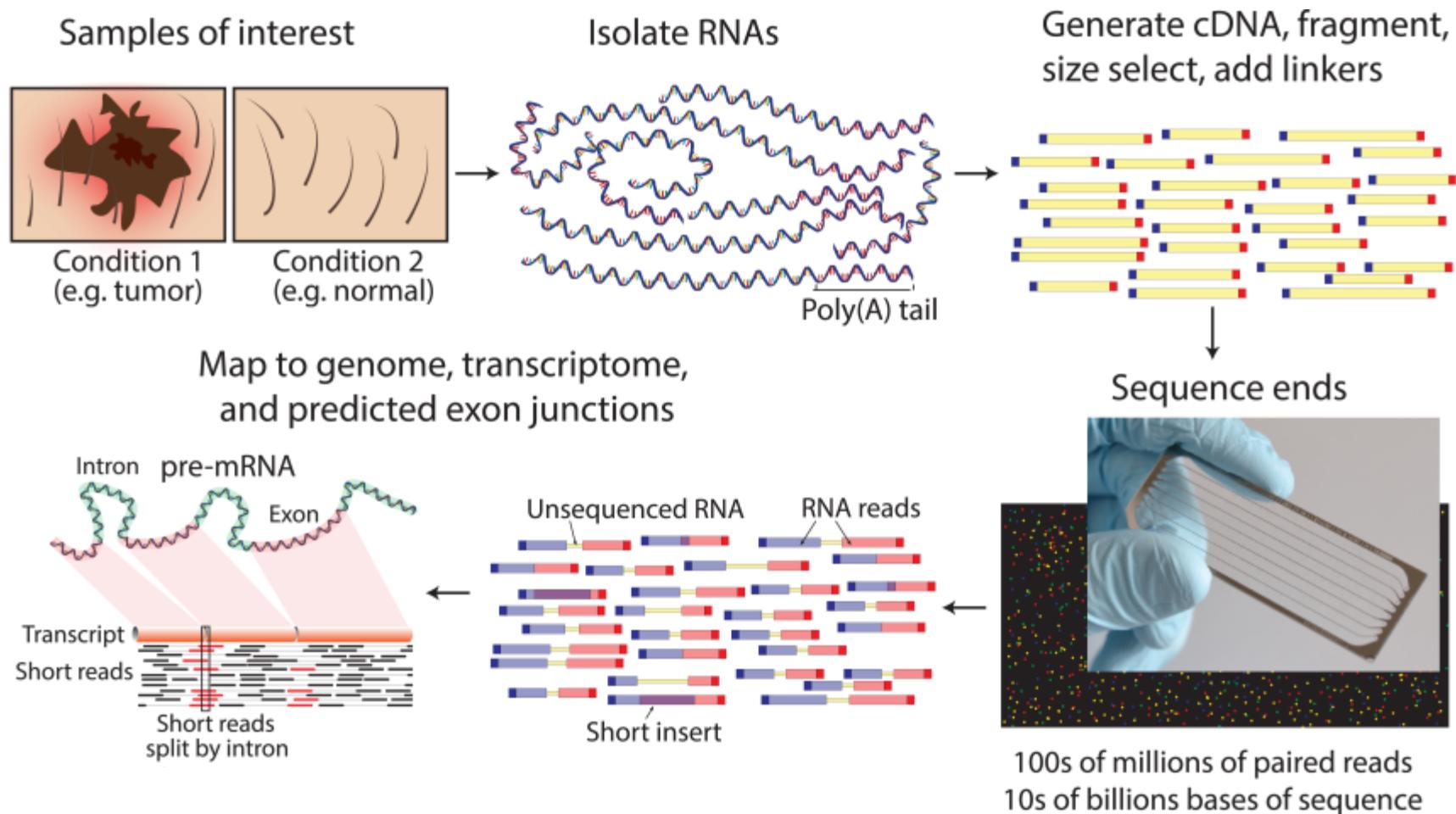




“RNA-Seq is able to identify thousands of differentially expressed genes, tens of thousands of differentially expressed gene isoforms and can detect mutations and germline variations for hundreds to thousands of expressed genetic variants, as well as detecting chimeric gene fusions, transcript isoforms and splice variants.”

Wang, *Nat Rev. Genet.*, 2009

# From sample to readout

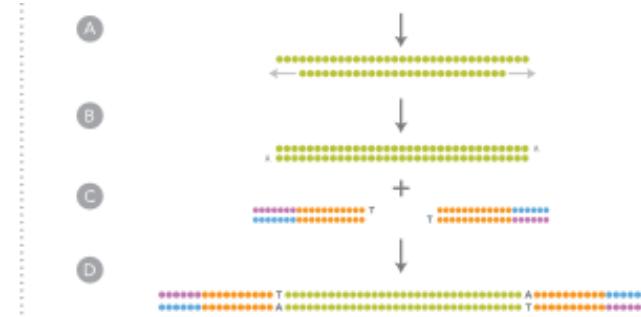


Griffith, Plos Comp. Biol., 2015

# How do I get my NGS data?

## 1 Library Preparation

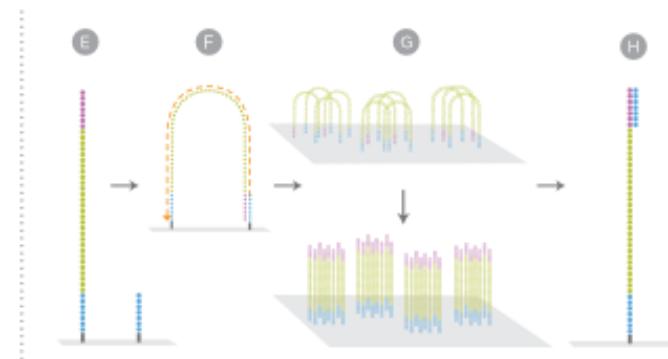
6 hours  
3 hours hands-on time



- A** Fragment DNA
- B** Repair ends  
Add A overhang
- C** Ligate adapters
- D** Select ligated DNA

## 2 Cluster Generation

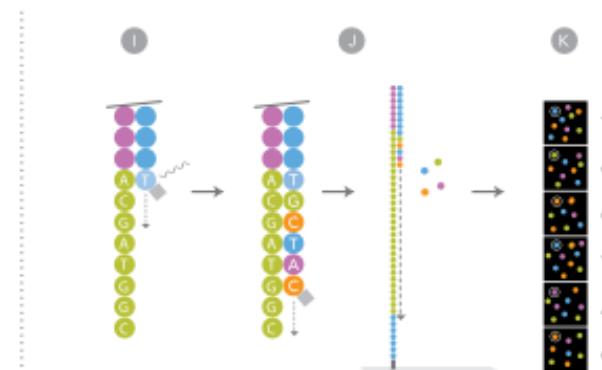
4 hours  
< 10 minutes hands-on time  
1–96 samples



- E** Attach DNA to flow cell
- F** Perform bridge amplification
- G** Generate clusters
- H** Anneal sequencing primer

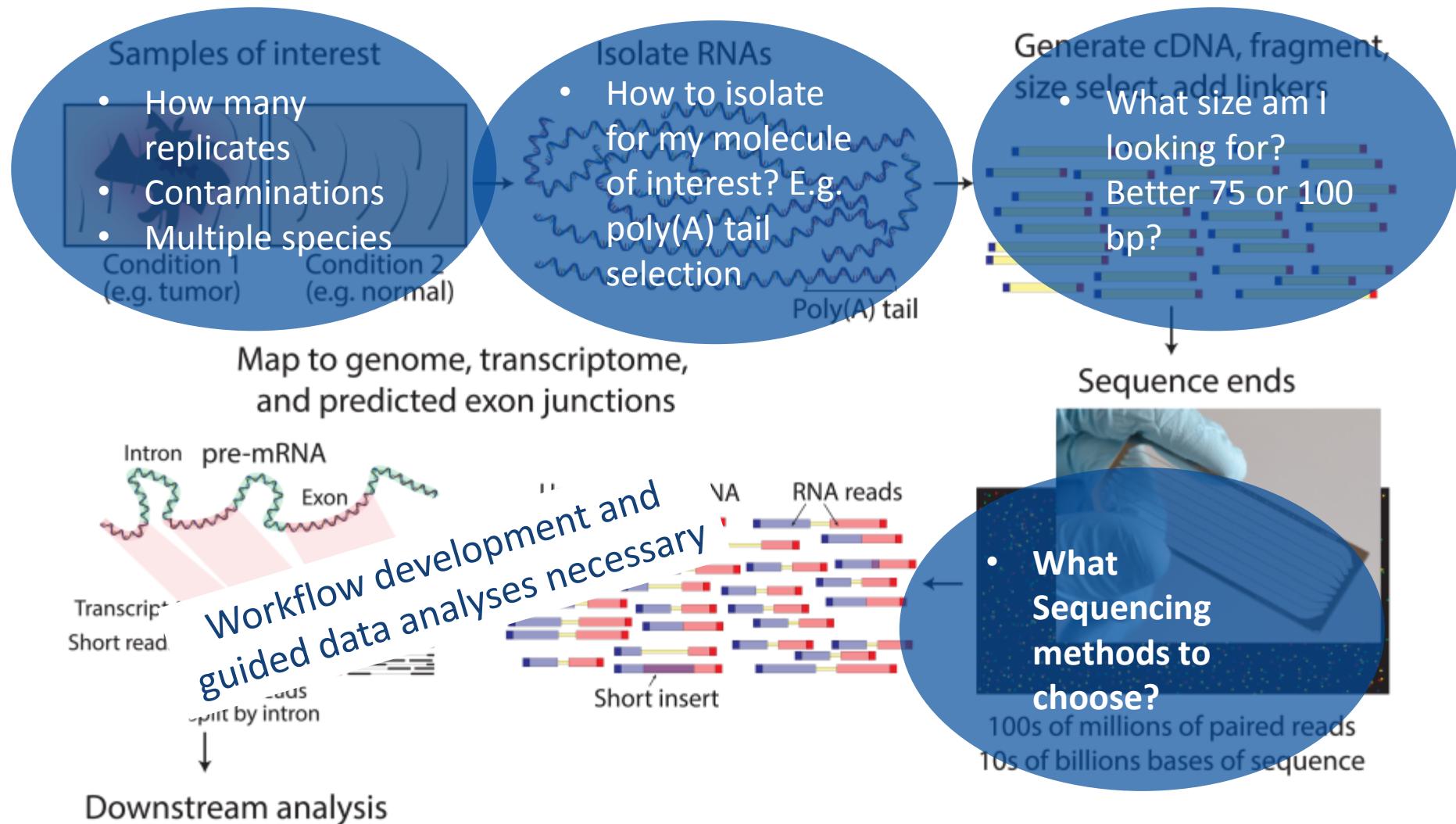
## 3 Sequencing

1–3 days single-read run  
3–9 days paired-end run  
30 minutes hands-on time  
8 lanes, up to 96 samples per flow cell (run)



- I** Extend first base, read, and deblock
- J** Repeat step above to extend strand
- K** Generate base calls

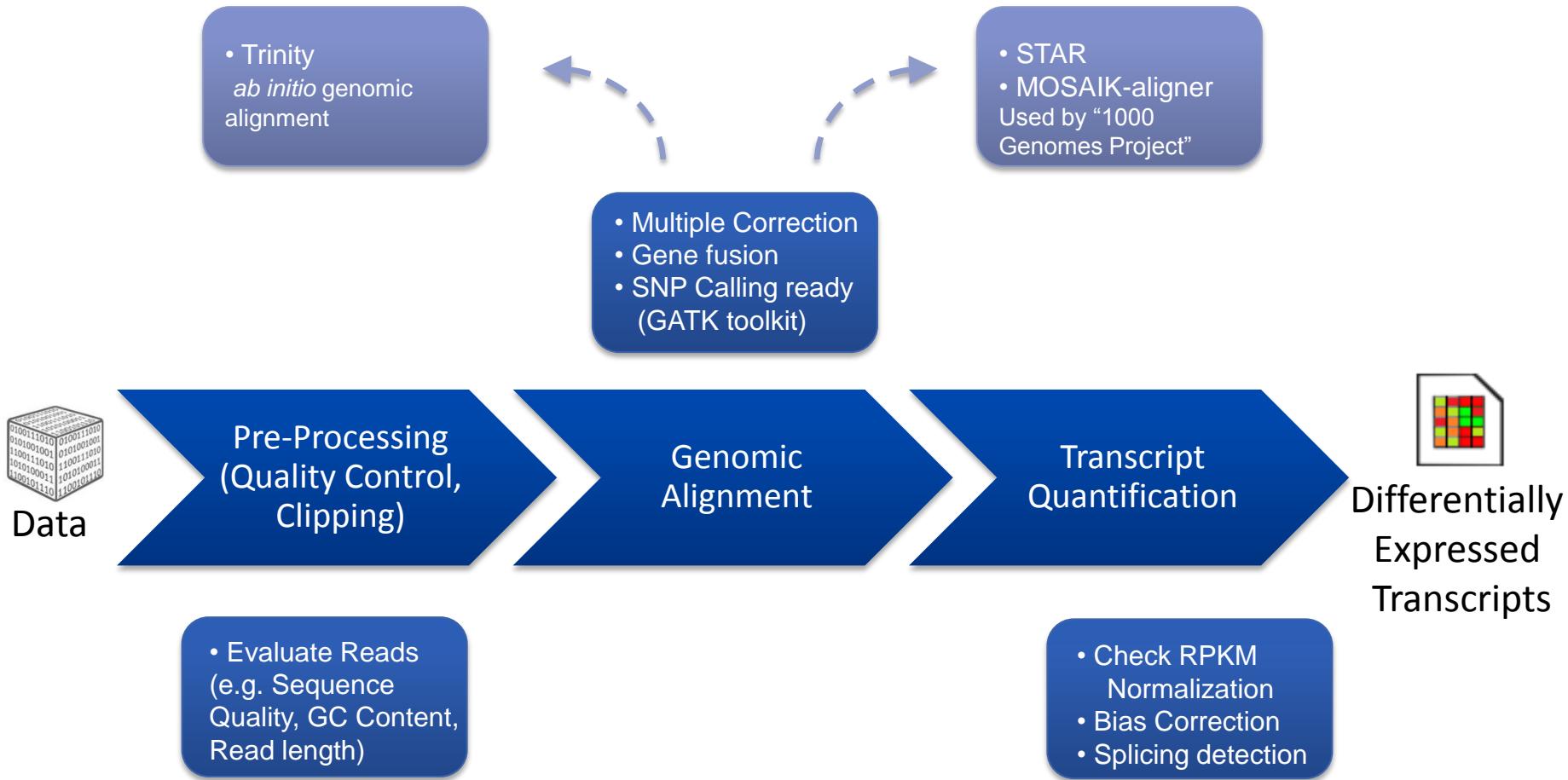
# From sample to readout



Griffith, Plos Comp. Biol., 2015



# Basic workflow for data processing



# Differential expression

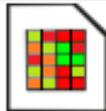


Data

Pre-Processing  
(Quality Control,  
Clipping)

Genomic  
Alignment

Transcript  
Quantification



Differentially  
Expressed  
Transcripts

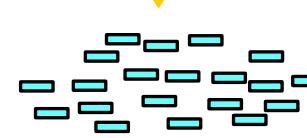
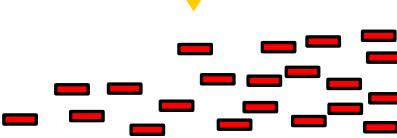
Condition A  
Multiple Copies of a Transcriptome

Condition B

Condition C



Reads



Ref. Sequence



Reads on Ref. Seq.



*Differential expression*

# Big Data and the need for new analyses



Grape

[big.crg.cat/services/grape](http://big.crg.cat/services/grape)



[mapman.gabipd.org](http://mapman.gabipd.org)



Chipster  
Open source platform for data analysis  
[chipster.csc.fi](http://chipster.csc.fi)



[r-project.org/](http://r-project.org/)



[gene-talk.de](http://gene-talk.de)



[illumina.com](http://illumina.com)

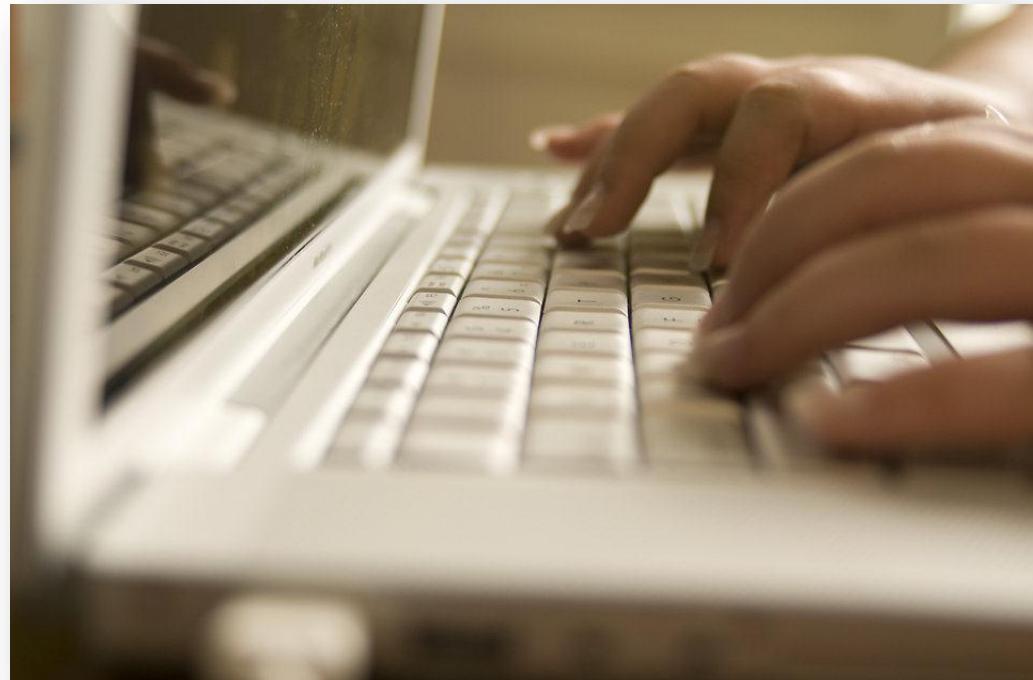


[bioconductor.org](http://bioconductor.org)

## Hands on part 1:

“Register and login to Galaxy. Explore and get familiar with your Galaxy.”

Please visit and explore Galaxy  
[usegalaxy.org](http://usegalaxy.org)





Training material: <http://galaxyproject.github.io/training-material/>

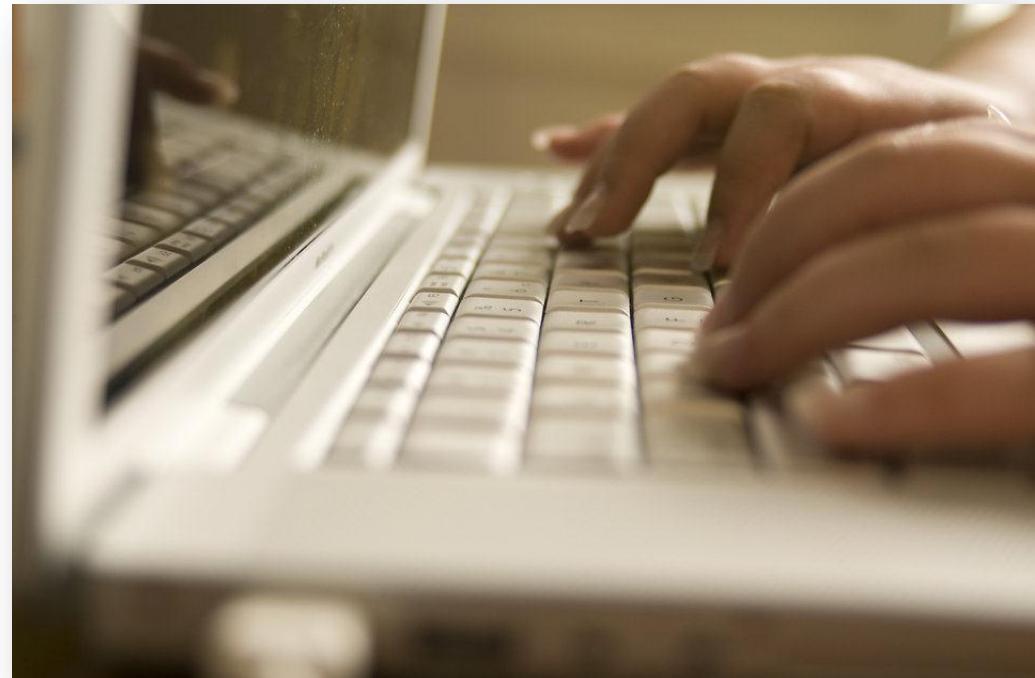
Manuscript: <https://www.biorxiv.org/content/early/2017/11/29/225680>

## Hands on part 2

9:00 – 10:00

“RNA-Seq data preprocessing and quality control”

Material: <http://galaxyproject.github.io/training-material/topics/sequence-analysis/>

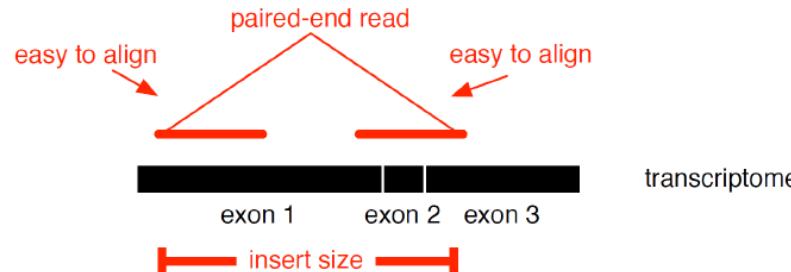


## Hands on part 3

10:15 – 11:45

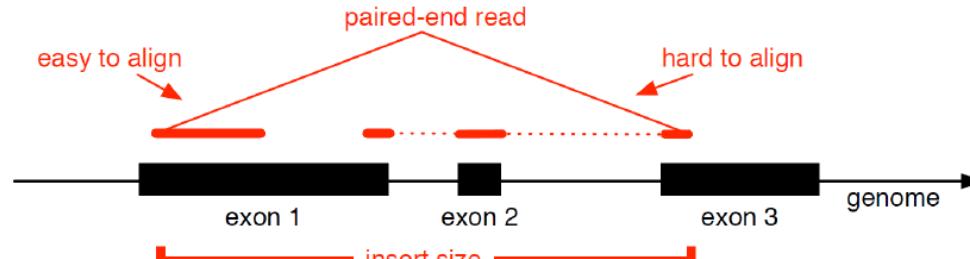
“Application of different read mapping approaches for genomic alignment”

## Transcriptome alignment



- reliable gene models required
- no detection of novel genes

## Genome alignment (splice-aware read alignment)

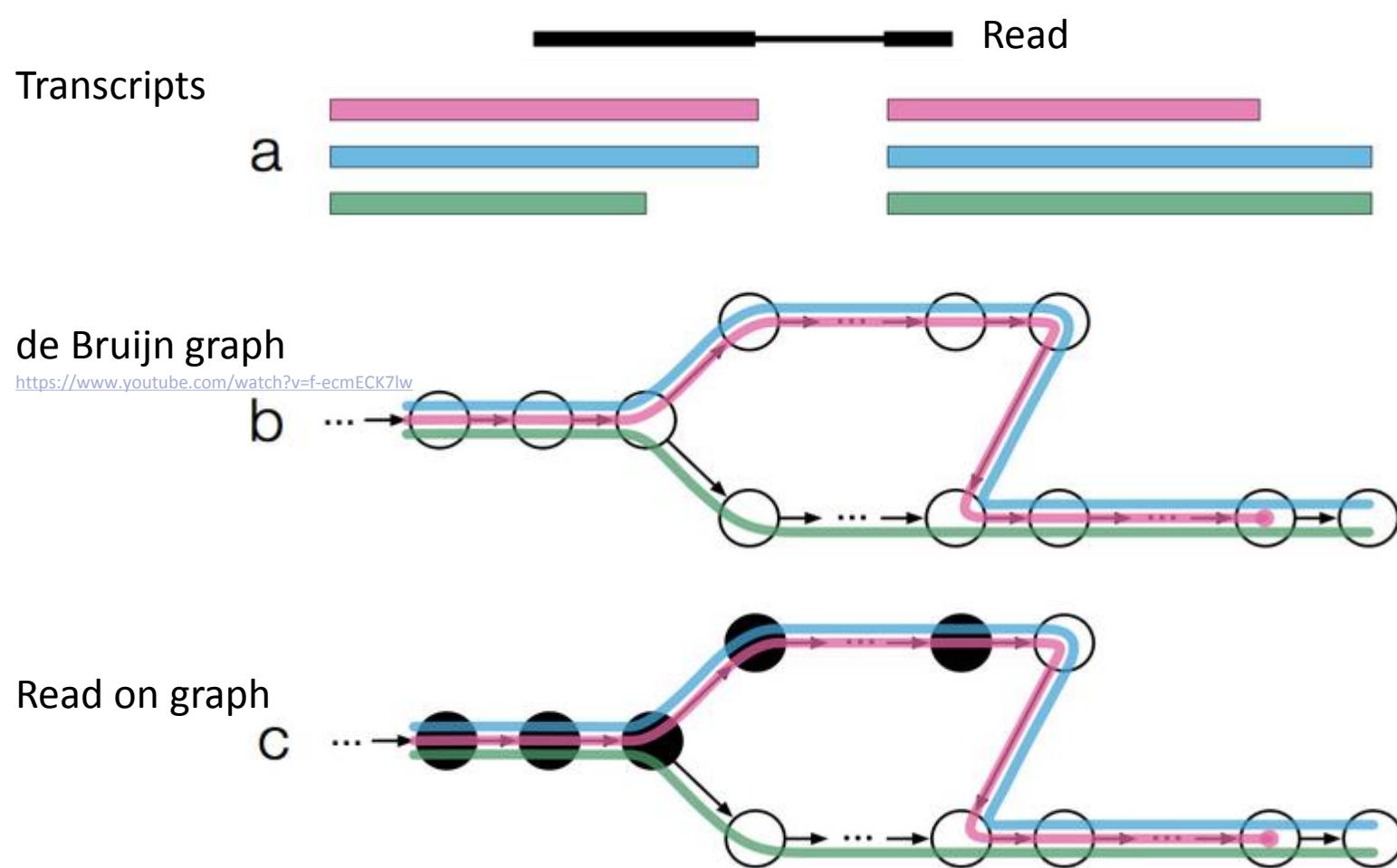


- + detection of novel genes and isoforms

Turro, EMBO, 2012

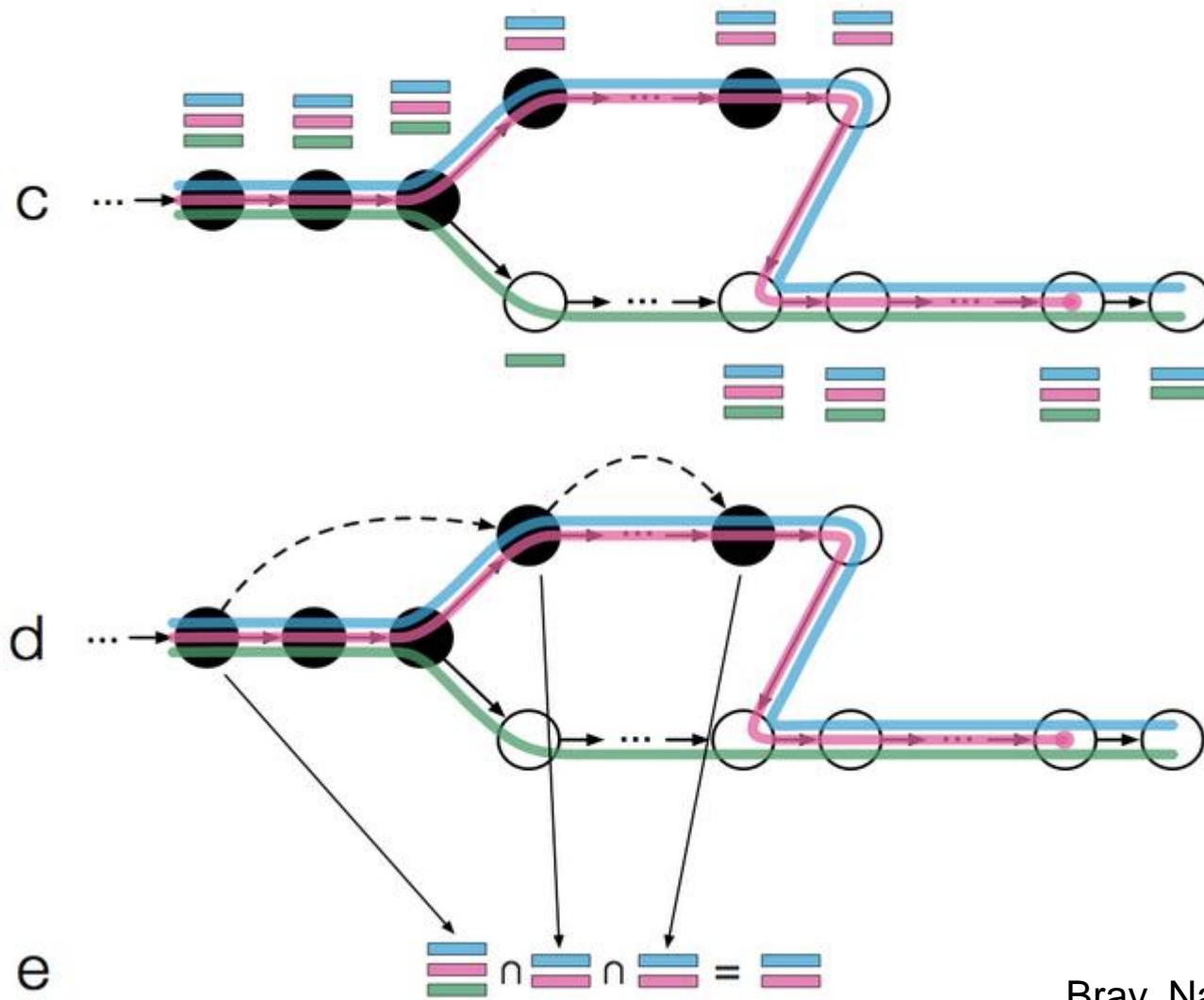
*For clinical usage combination of different algorithms possible:*

# Genomic alignment - pseudoalignment



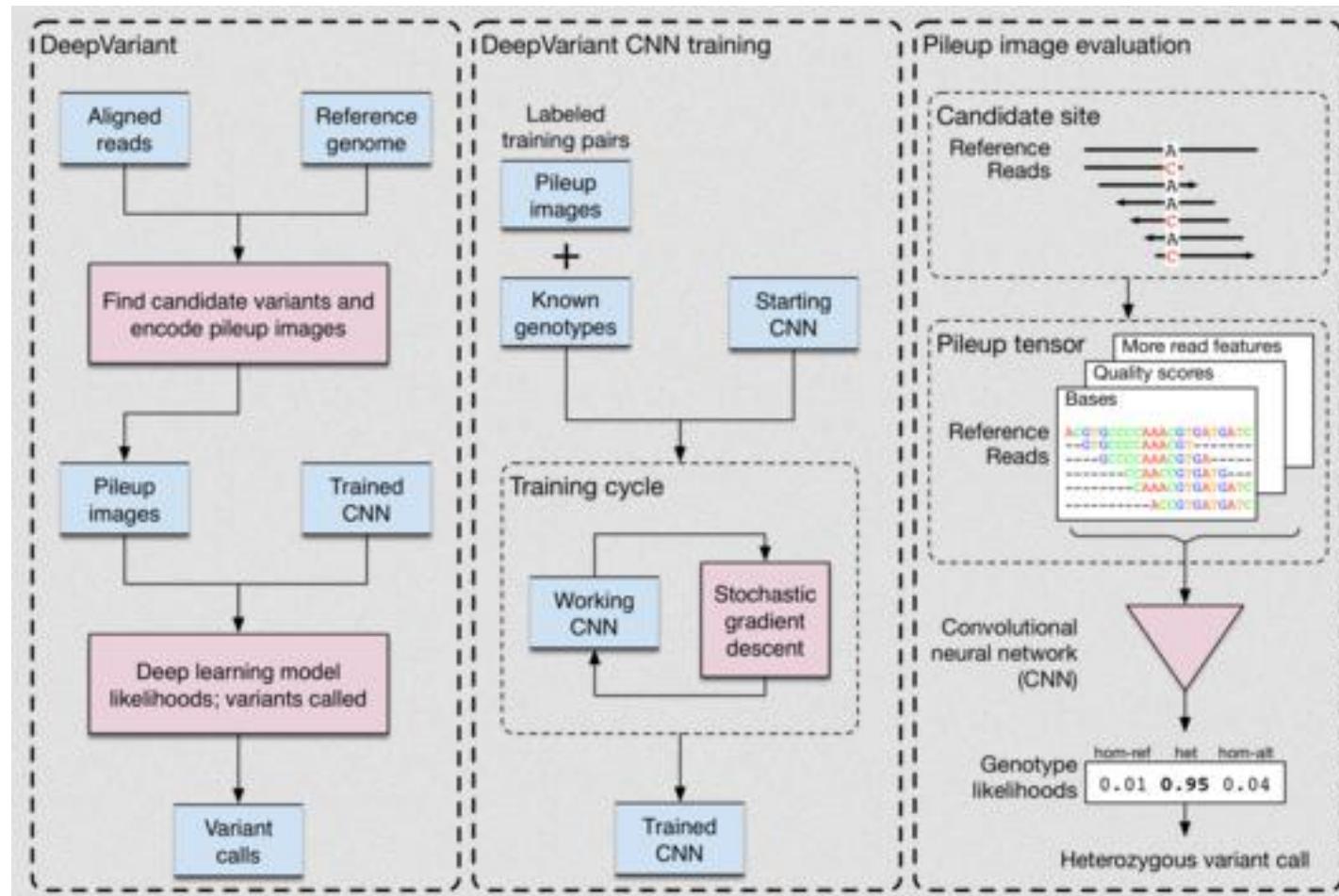
Bray, Nat Biotech, 2016

# Genomic alignment - pseudoalignment



Bray, Nat Biotech, 2016

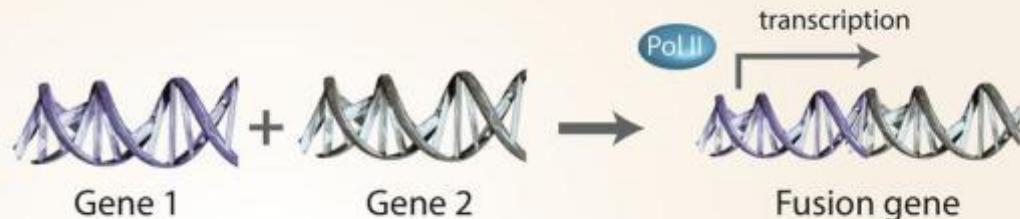
# Genomic alignment – deep neural network



<https://github.com/google/deepvariant>

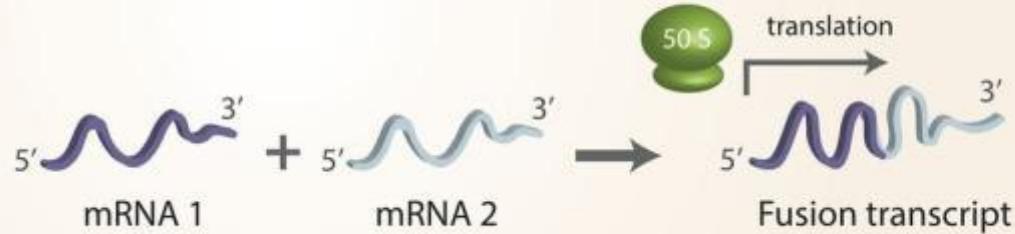
## A Fusion by structural rearrangements

Translocations, inversions, deletions and insertions



## B Fusion by transcription or splicing

Transcription read-through, mRNA *trans*-splicing or *cis*-splicing



Natasha, *Nucl. Acids Res.*, 2016

“ Gene fusions are associated with oncogenic properties, and often act as driver mutations in a wide array of cancer types.”

- Deregulating one of the involved genes
- Forming a fusion protein with oncogenic functionality
- Inducing a loss of function

Yoshihara, *Oncogene*, 2015

- Read counts
  - Count the reads per feature
    - relatively easy: count the number of reads per gene, exon, ...
  - How to handle multi-mapping reads (i.e. reads with multiple alignments)?
- Normalization - aims to make expression levels comparable across:
  - Features (genes, isoforms, ...)
  - RNA libraries (samples)
- Normalization methods:
  - RPKM / FPKM (Cufflinks /Cuffdiff) (Mortazavi, Nat Meth, 2008)
  - TMM (edgeR) (Robinson & Oshlack, Genome Biol, 2010)
  - DESeq2 (DESeq2) (Love et al., Genome Biol, 2014)



- **Reads Per Kilobase per Million**

- $$RPKM = \frac{\text{Raw number of reads}}{\text{Exon length}} * \frac{1.000.000}{\text{Number of reads mapped in the sample}}$$

- In RNA-seq, the relative expression of a transcript is proportional to the number of cDNA fragments that originate from itHowever:
  - The total number of fragments is biased towards larger genes
  - Total number of fragments is related to total library depth
- Differences with and without normalization and differences among them stated at <https://www.youtube.com/watch?v=TTUrtCY2k-w>

# Visualization of .bam files



<http://software.broadinstitute.org/software/igv/>



Tablet



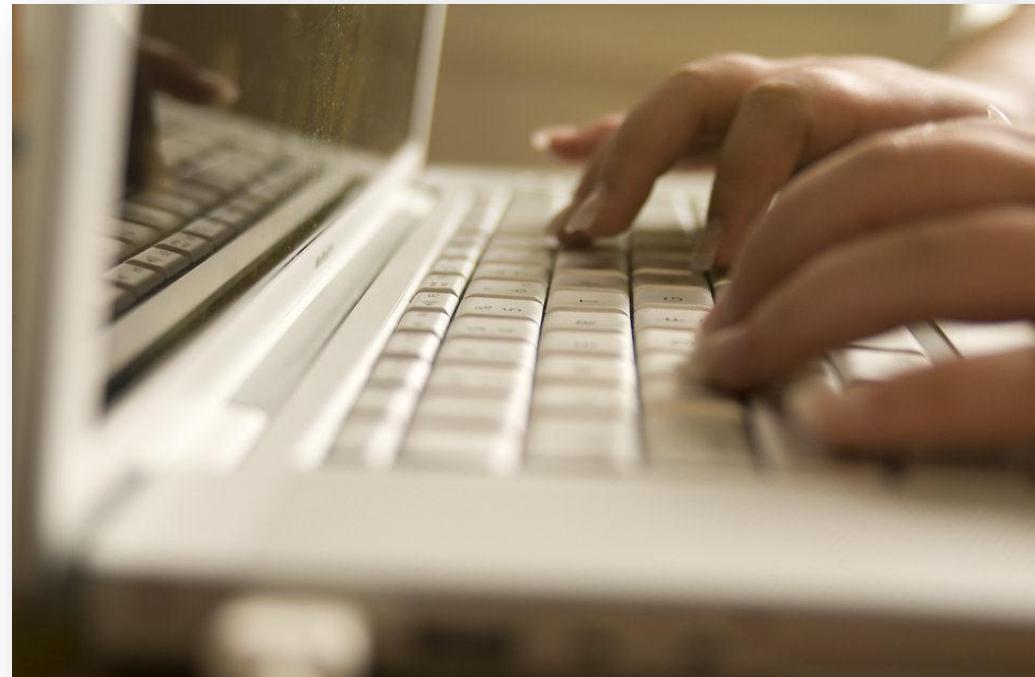
<https://ics.hutton.ac.uk/tablet/>

## Hands on part 3

10:15 – 11:45

“Application of different read mapping approaches for genomic alignment”

Material: <http://galaxyproject.github.io/training-material/topics/sequence-analysis/>



11:45 – 12:30

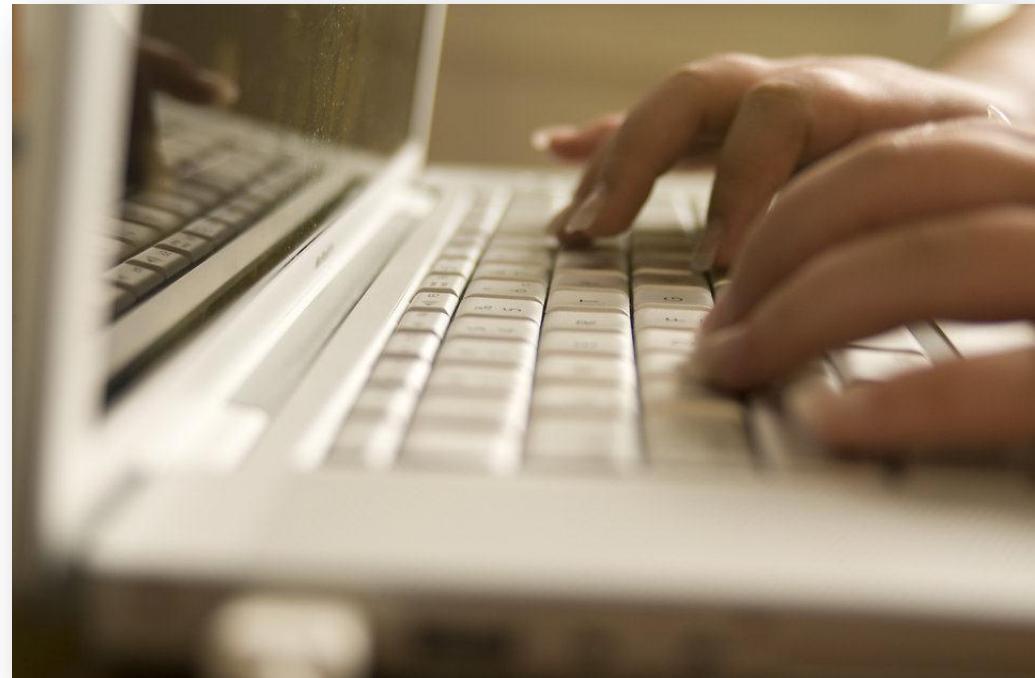


## Hands on part 4

12:30 – 13:30

“Visualizations of RNA-Seq experiments with Galaxy”

Material: <http://galaxyproject.github.io/training-material/topics/transcriptomics/>



## Part 5

13:30 – 14:00

“Introduction into the underlying technology of Galaxy”



# Supporting new data analysis approaches



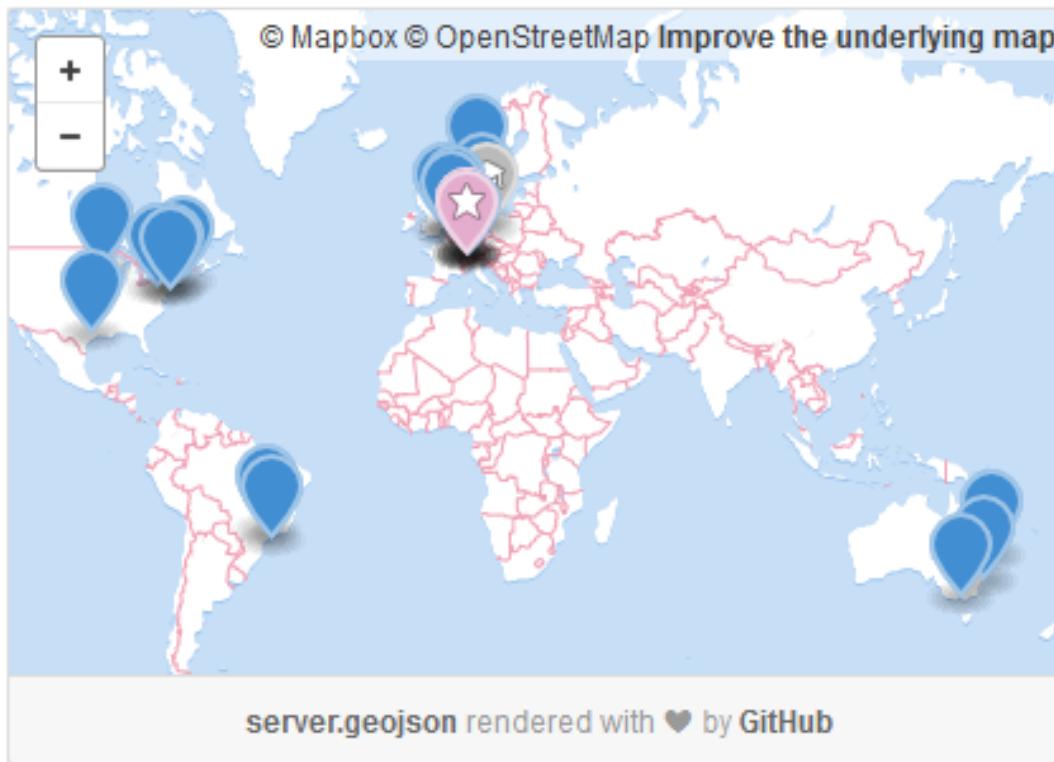
- Key performance of Galaxy
  - Accessibility
  - Reproducibility
  - Transparency



[jupyter.org/](http://jupyter.org/)



[usegalaxy.org](http://usegalaxy.org)



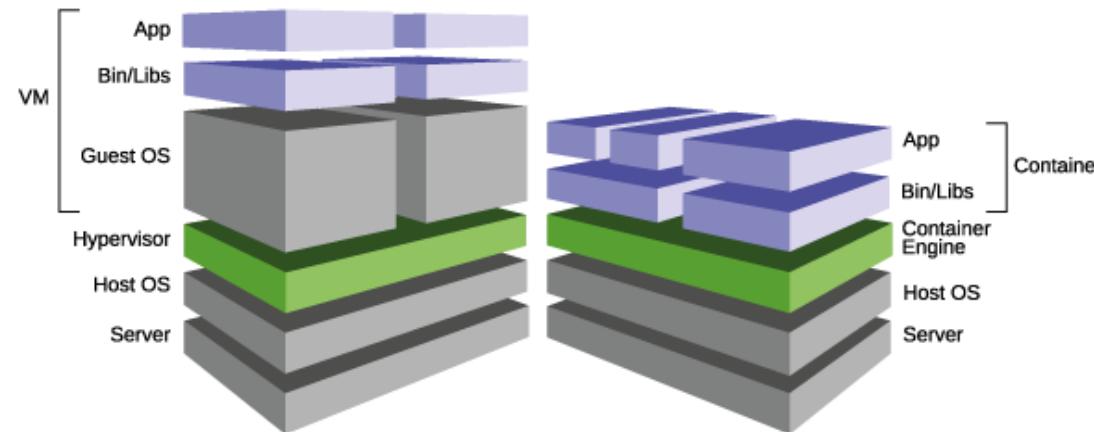
[elixir-europe.org](http://elixir-europe.org)



GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

[denbi.de](http://denbi.de)

- Secure by default
- Build for scale
- Extensible and flexible
- E.g. used by ebay, GE, illumina, Spotify





+



= Symbiosis!



- Tailor-made, user specific and integration into a general framework to develop workflows addressing the users need and facilitating a reuse

- Stand-alone Docker container which “conserves” your whole tool compilation (for an easy use – one command line or [kitematic.com](https://kitematic.com) click!)

```
docker run -p 8080:80 bgruening/galaxy-rna-workbench
```

- Get a single minimized Docker container for every tool and obtain a maximum flexibility (more advanced)

# Example Dockerfile: TRAPLINE + Docker

```
FROM bgruening/galaxy-stable
```

Source  
Container

```
MAINTAINER Markus Wolfien, markus.wolfien@gmail.com
```

```
ENV GALAXY_CONFIG_BRAND "TRAPLINE_160801"
```

```
WORKDIR /galaxy-central
```

```
RUN install-repository \
```

```
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name fastq_groomer" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name fastq_trimmer_by_quality" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name fastx_clipper" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name tophat_fusion_post" \
--url https://toolshed.g2.bx.psu.edu/ -o scottx611x --name tophat2_with_gene_annotations" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name cufflinks" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name cuffmerge" \
--url https://toolshed.g2.bx.psu.edu/ -o devteam --name cuffcompare" \
```

Tools to be  
added to  
the new  
Container

```
VOLUME ["/export/", "/data/", "/var/lib/docker"]
```

```
EXPOSE :80
```

```
EXPOSE :21
```

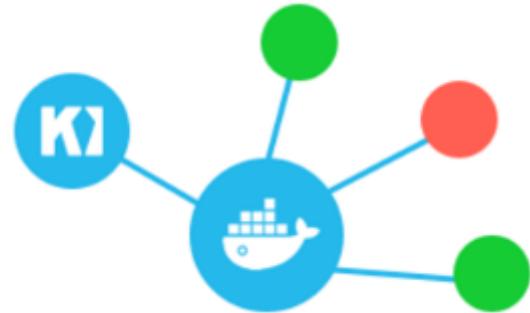
```
EXPOSE :8080
```

```
CMD ["/usr/bin/startup"]
```



Short *Kitematic* and *Docker-Galaxy* showcase

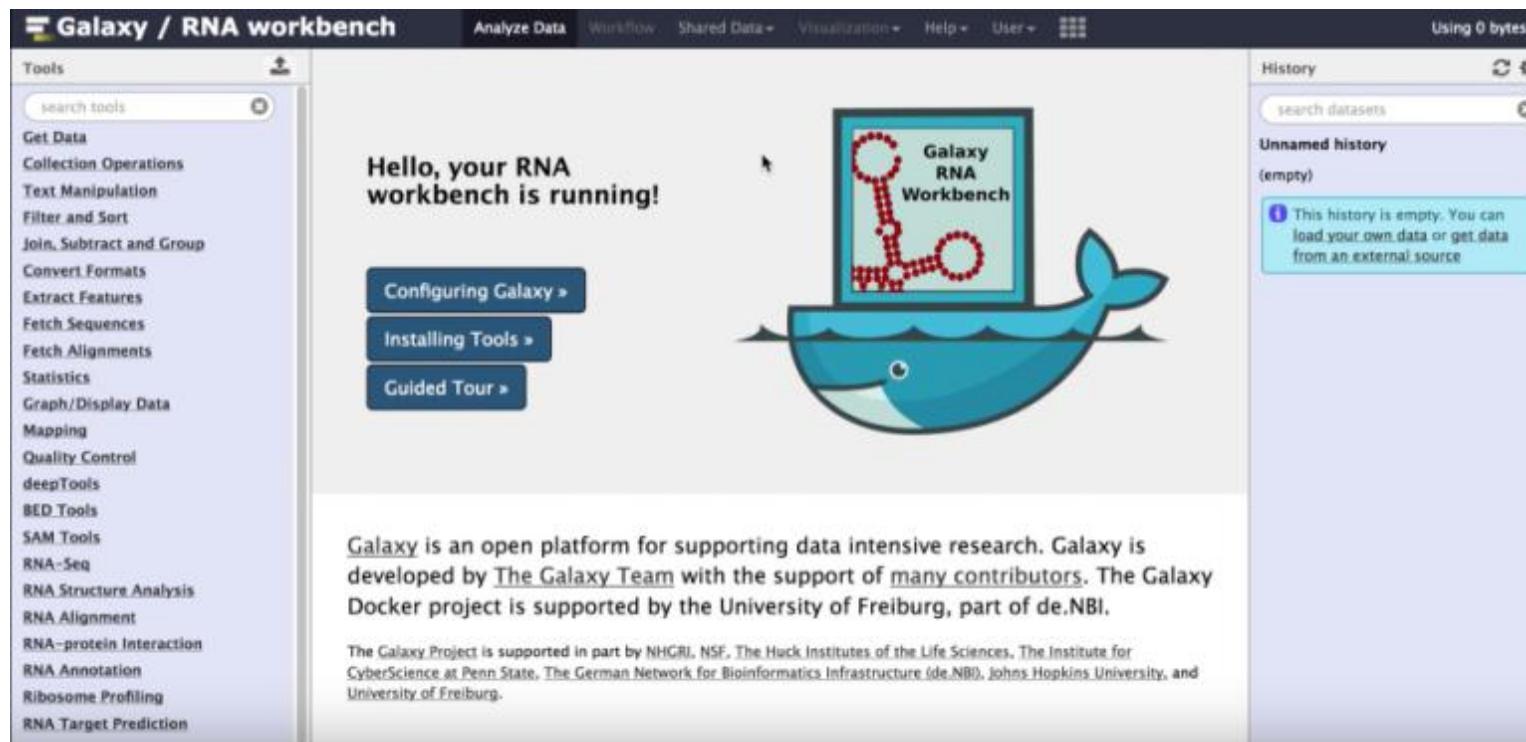
<https://kitematic.com/>



*Containerization!*  
*New workflow technologies eliminate*  
*“works on my machine” problems*

- Specialized Galaxy instance for RNA analyses provided by the RBC
- Contains +50 tools for structure analyses, annotation, alignment and many more

[github.com/bgruening/galaxy-rna-workbench](https://github.com/bgruening/galaxy-rna-workbench)



Galaxy / RNA workbench

Analyze Data Workflow Shared Data Visualization Help User Using 0 bytes

Tools

Get Data

Collection Operations

Text Manipulation

Filter and Sort

Join, Subtract and Group

Convert Formats

Extract Features

Fetch Sequences

Fetch Alignments

Statistics

Graph/Display Data

Mapping

Quality Control

deepTools

BED Tools

SAM Tools

RNA-Seq

RNA Structure Analysis

RNA Alignment

RNA-protein Interaction

RNA Annotation

Ribosome Profiling

RNA Target Prediction

Hello, your RNA workbench is running!

Configuring Galaxy »

Installing Tools »

Guided Tour »

Galaxy RNA Workbench

History

search datasets

Unnamed history (empty)

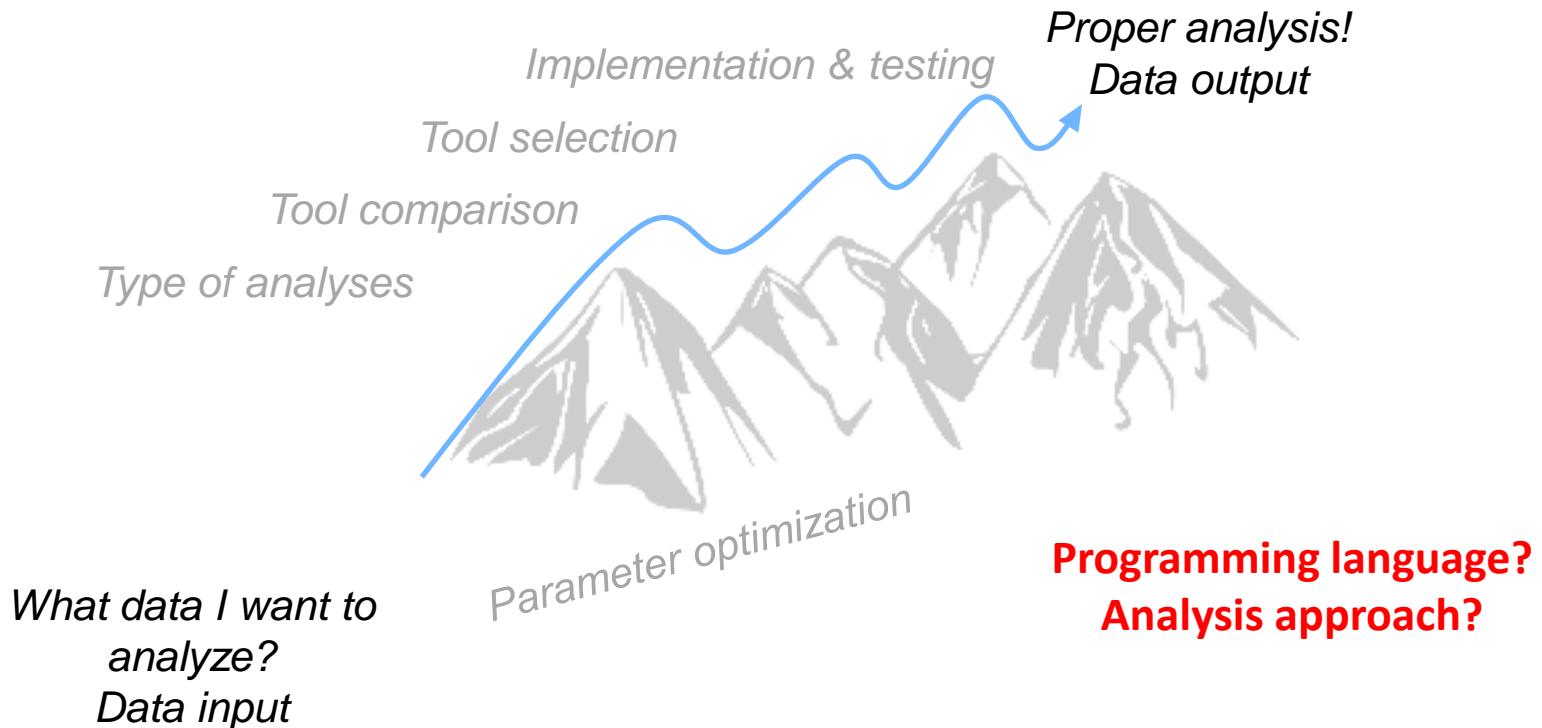
This history is empty. You can load your own data or get.data from an external source

Galaxy is an open platform for supporting data intensive research. Galaxy is developed by [The Galaxy Team](#) with the support of many contributors. The Galaxy Docker project is supported by the University of Freiburg, part of de.NBI.

The Galaxy Project is supported in part by NHGRI, NSF, The Huck Institutes of the Life Sciences, The Institute for CyberScience at Penn State, The German Network for Bioinformatics Infrastructure (de.NBI), Johns Hopkins University, and University of Freiburg.

Gruening et al., NRA, 2017

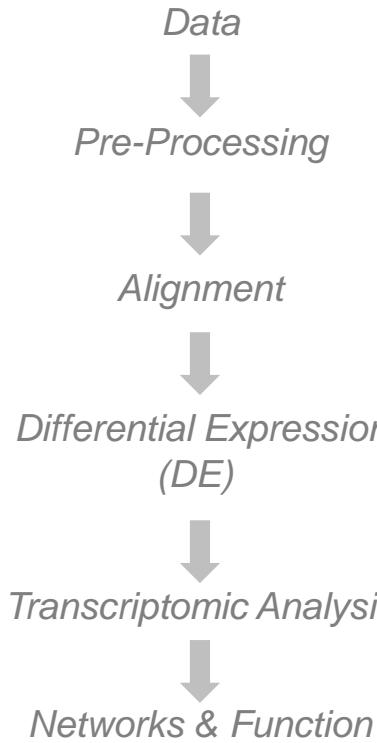
# The struggle for the right approaches



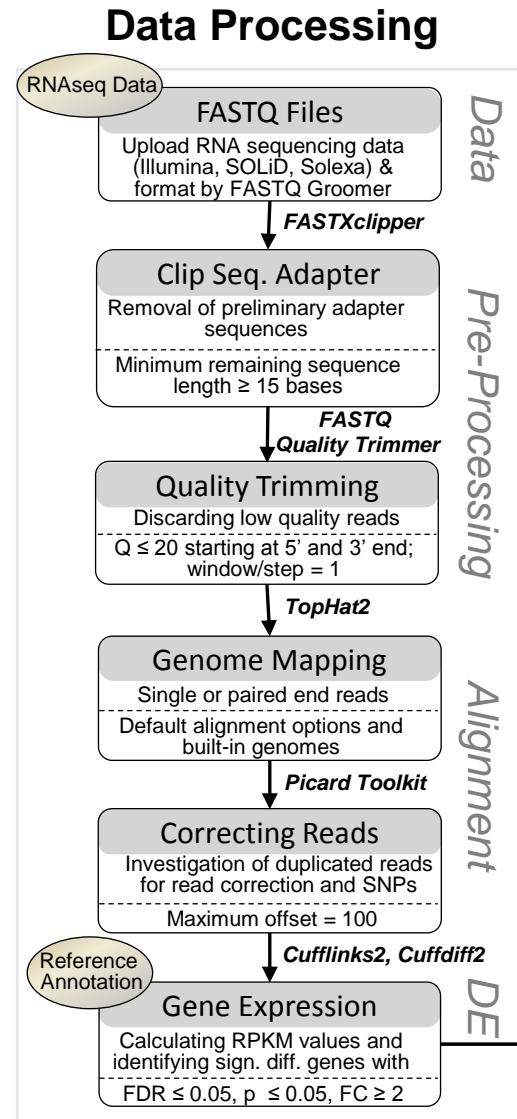
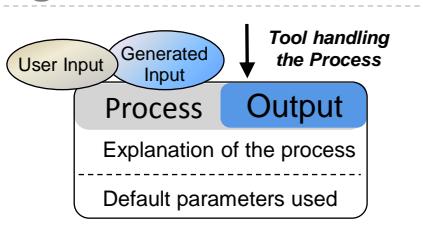
Lott, Wolfien, Riege, Bagnacani, et al., *J.Biotech*, 2017

# Transparent Reproducible Automated PipeLINE - TRAPLINE

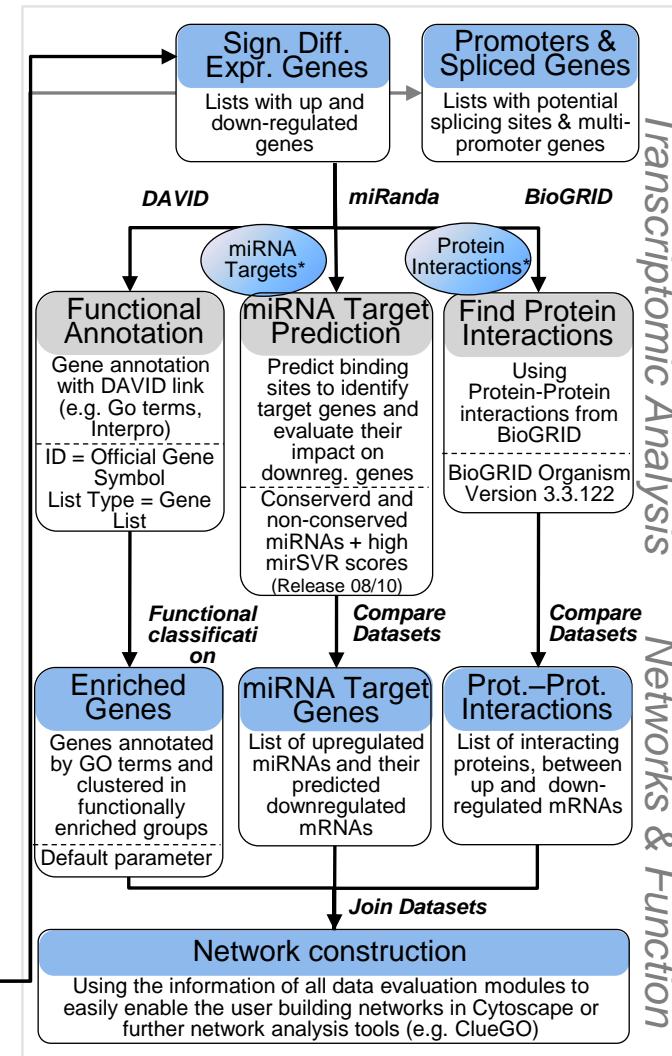
 Galaxy Modules implemented:



## Legend:



## Data Evaluation and Annotation

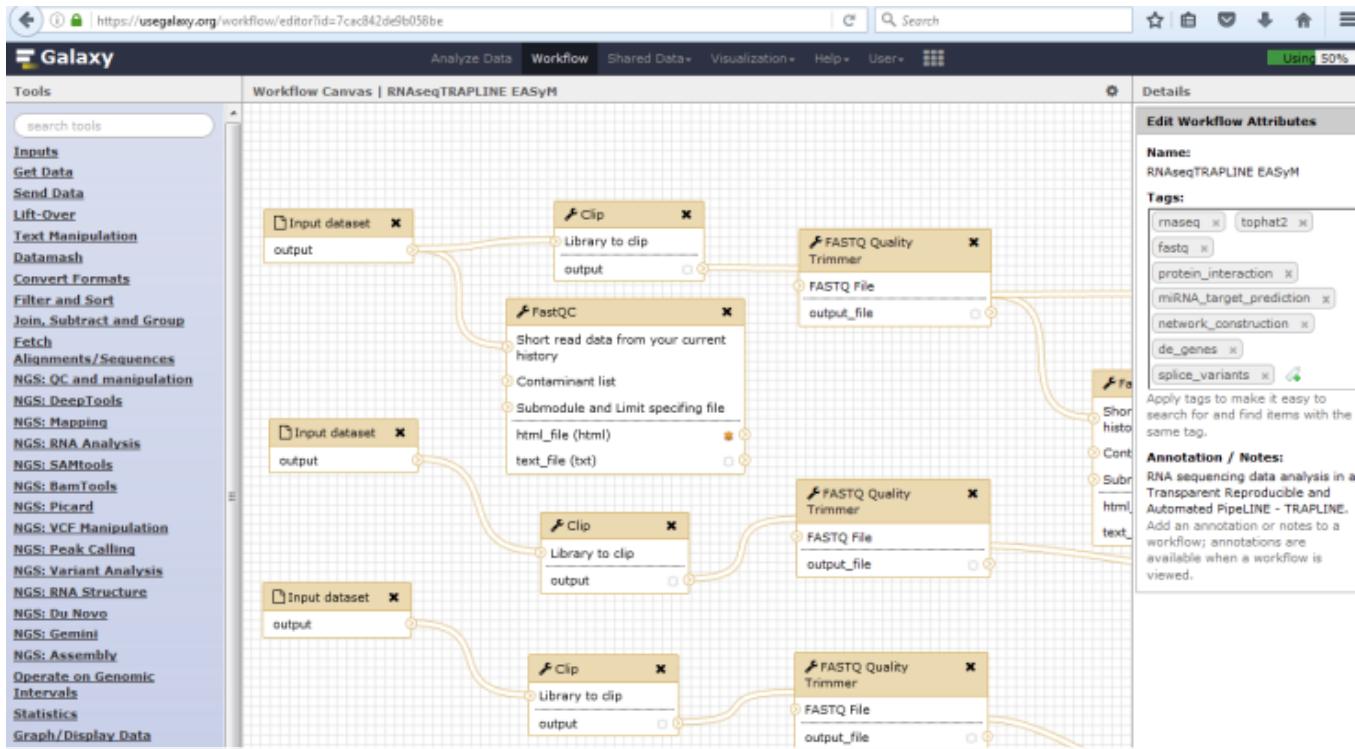


# Example Galaxy workflow: TRAPLINE.ga

```
"12": {
    "annotation": "",
    "content_id": "toolshed.g2.bx.psu.edu/repos/devteam/tophat2/tophat2/2.1.0",
    "id": 12,
    "input_connections": {
        "refGenomeSource|ownFile": {
            "id": 3,
            "output_name": "output"
        },
        "singlePaired|input": {
            "id": 8,
            "output_name": "trimmed_reads_paired_collection"
        }
    },
    "inputs": [
        {
            "description": "runtime parameter for tool TopHat",
            "name": "refGenomeSource"
        },
        {
            "description": "runtime parameter for tool TopHat",
            "name": "singlePaired"
        }
    ],
    "label": null,
    "name": "TopHat",
    "outputs": [
        {
            "name": "align_summary",
            "type": "txt"
        },
        {
            "name": "fusions",
            "type": "tabular"
        }
    ]
},
```

Specific xml file  
with tools,  
parameters and  
meta data!

# Using workflow development



- Workflows are build and maintained by experts
- Workflows are modular and can be easily adapted to other tasks
- Implementation of other tools can be done (quickly)
- Application of workflows and tools is targeted for non-computational users



+133 different workflow management systems!

Almost no interoperability!

Need for a common line!



<https://github.com/common-workflow-language/common-workflow-language/wiki/Existing-Workflow-systems>

# CWL - Open standard designed to express workflows and their tooling groups in YAML structured text files



- Common format for bioinformatics tool execution
  - Inputs & outputs are fully specified
- Community based standards effort, not a specific software package
- Designed for shared-nothing cluster & cloud environments
  - Tool executions are isolated from one another & from parent process
- Designed for containers (e.g. Docker, BioContainer)
- Well defined execution process:
  - 1. Collect & validate inputs
  - 2. Map input file paths to locations inside container
  - 3. Build tool command line
  - 4. Build Docker invocation
  - 5. Execute
  - 6. Collect & validate outputs



COMMON  
WORKFLOW  
LANGUAGE

[commonwl.org](http://commonwl.org)

## Example.yaml: samtools [sort]

```
class: CommandLineTool  
cwlVersion: draft-3  
description: Sort by chromosomal coordinates
```

File type and meta data

```
requirements:  
  - class: DockerRequirement  
    dockerPull: scidap/samtools:v1.2-216-gdfffc67f
```

Runtime environment

```
inputs:  
  - id: input  
    type: File  
    inputBinding:  
      position: 1  
  - id: output_name  
    type: string  
    inputBinding:  
      position: 2
```

Input parameters

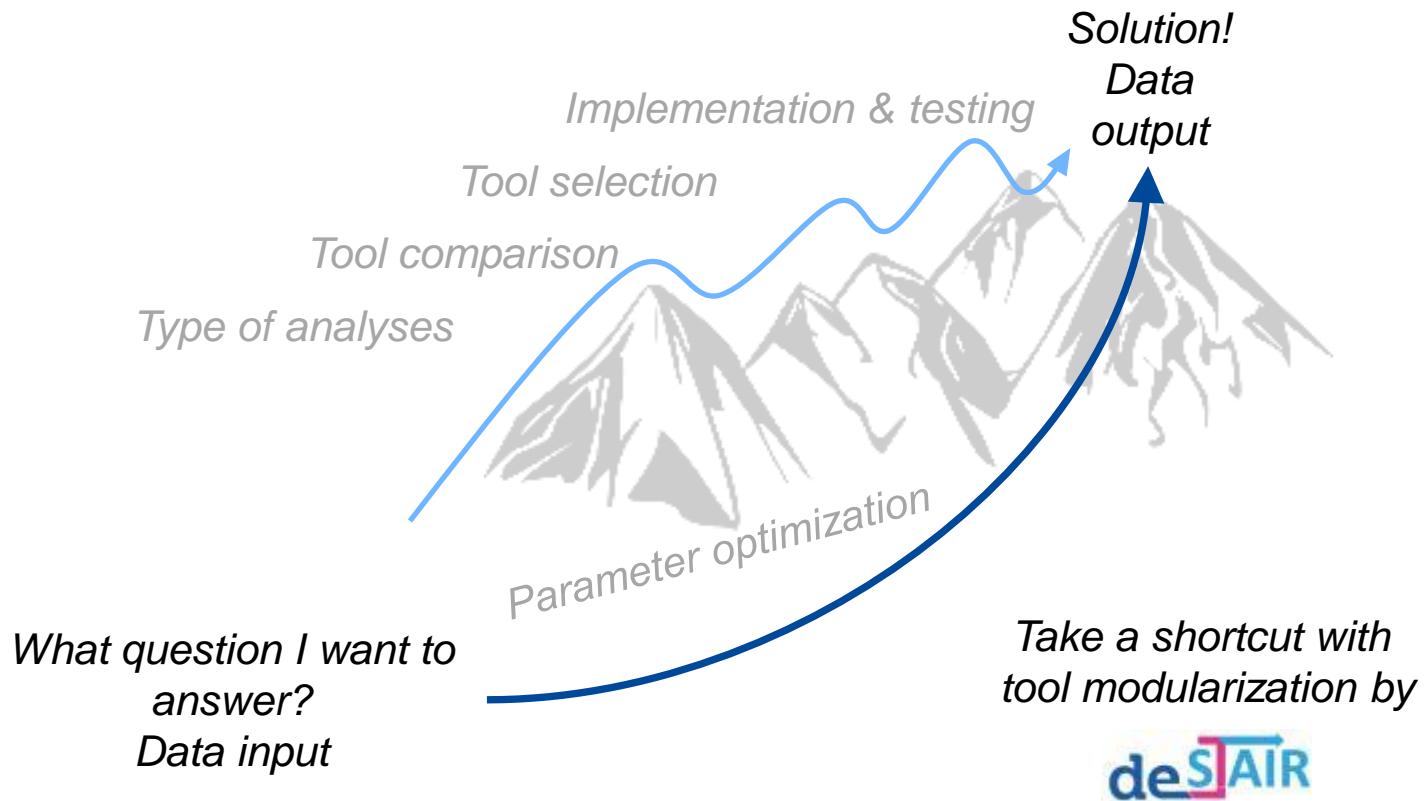
```
outputs:  
  - id: output  
    type: File  
    outputBinding:  
      glob: $(inputs.output_name)
```

Output parameters

```
baseCommand: [samtools, sort]
```

Executable

# Why using workflows?

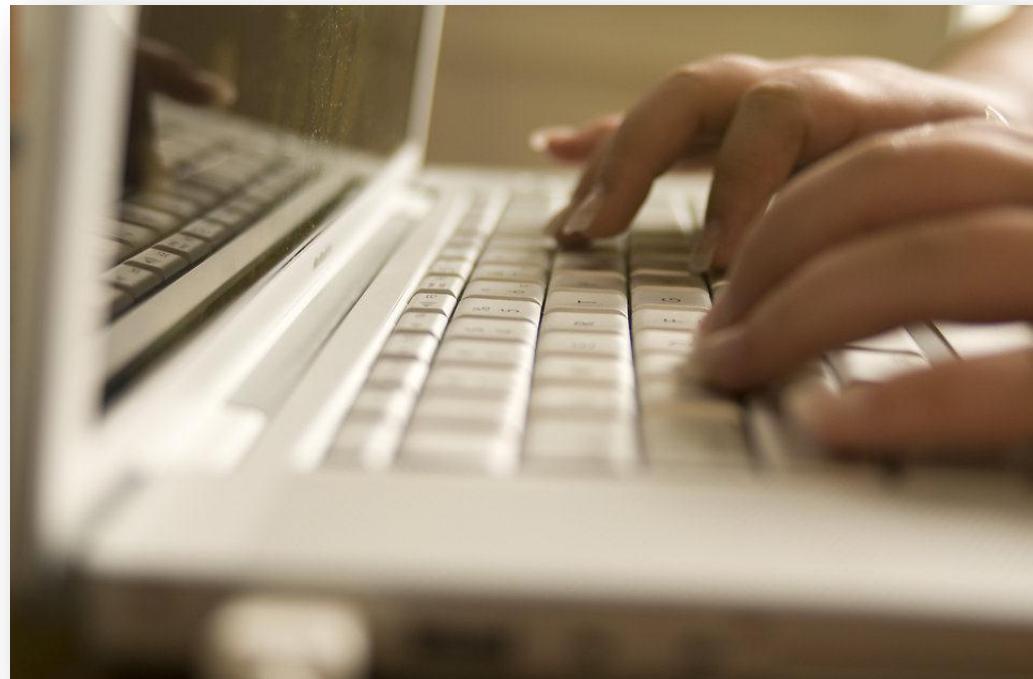


Lott, Wolfien, Riege, Bagnacani, et al., *J.Biotech*, 2017

## Hands on part 6

14:15 – 15:45

“Clinical use case for RNA-Seq, combining all previous processing steps and linking results to further resources”



# Our use case for today



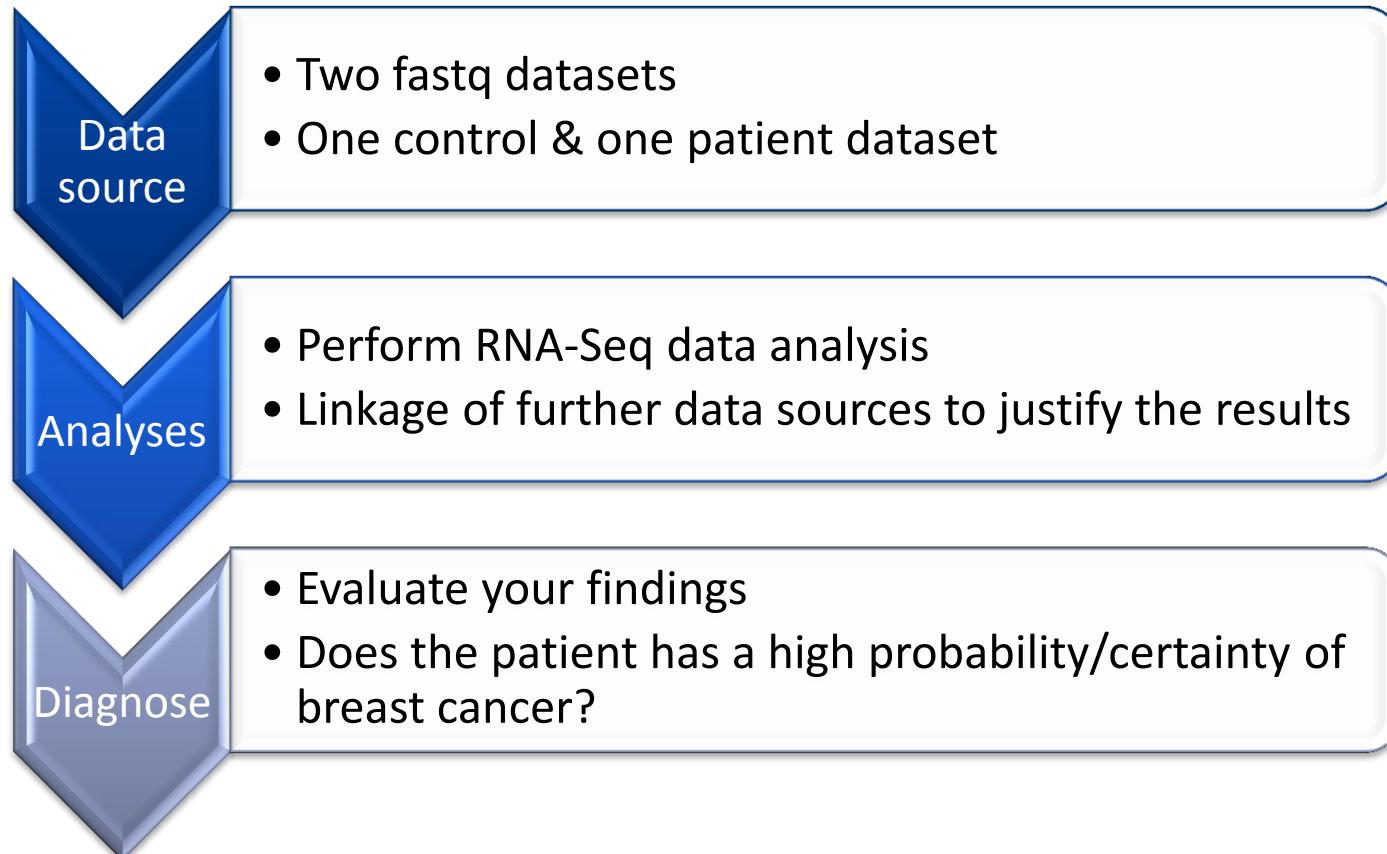
[www.pinkribbon-deutschland.de](http://www.pinkribbon-deutschland.de)

Most common cancer in women  
worldwide

Leading cause of death from cancer in  
women worldwide

Predictive factors that identify a benefit  
*Many different variations and subtypes*  
*Many different therapeutically approaches*

# Our use case for today – breast cancer screening

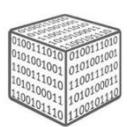


Get your data set!



<https://usegalaxy.org/u/mwolfien/h/rna-seq-workshop-kiel>

# Basic workflow for data processing



Data



- Evaluate Reads  
(e.g. Sequence Quality, GC Content, Read length)

Genomic  
Alignment

Transcript  
Quantification



Differentially  
Expressed  
Transcripts

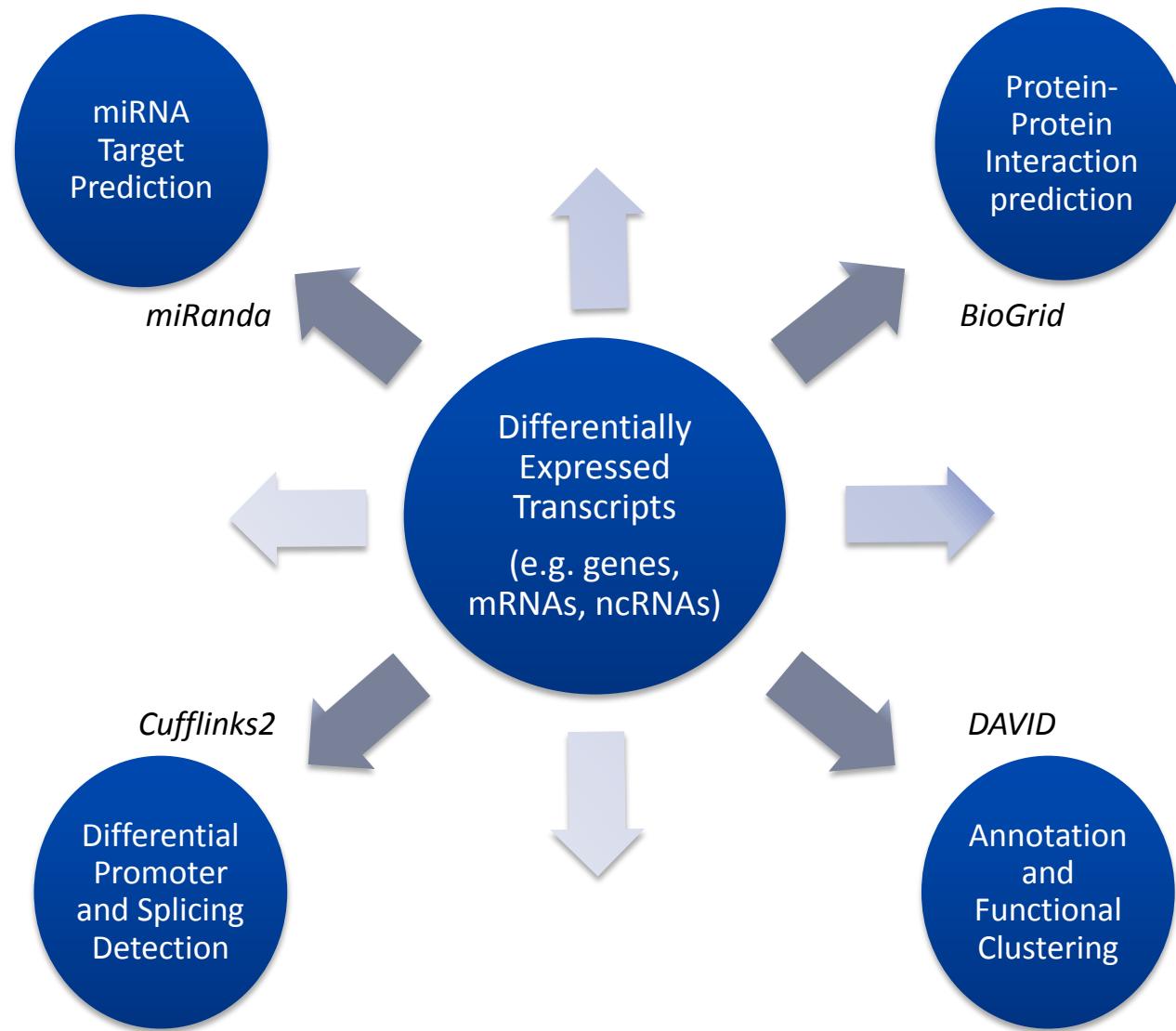
- Check RPKM Normalization
- Bias Correction
- Splicing detection

# Explore your data!

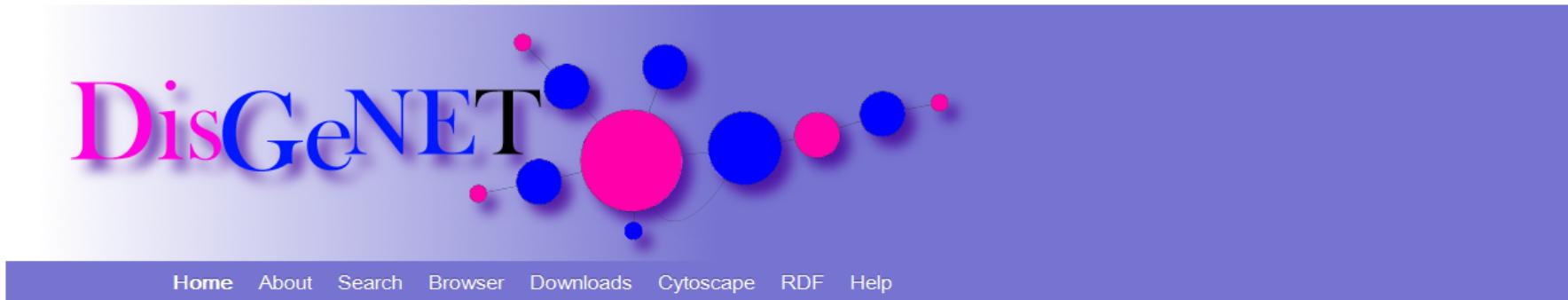
gene\_exp - Microsoft Excel

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
test_id	gene_id	gene	locus	sample_1	sample_2	status	value_1	value_2	log2(fold_ch)	test_stat	p_value	q_value	significant						
1	test_id																		
2	ADAMTS4	ADAMTS4	ADAMTS4	chr1:1611595	patient	control	OK	118.001	48.454	-128.411	-43.619	5,00E-05	0.0140016	yes					
3	AOC3	AOC3	AOC3	chr17:410032	patient	control	OK	11.335	457.446	-130.911	-43.654	5,00E-05	0.0140016	yes					
4	APOD	APOD	APOD	chr3:1952955	patient	control	OK	207.113	540.507	13.839	471.016	5,00E-05	0.0140016	yes					
5	ARHGAP40	ARHGAP40	ARHGAP40	chr20:372305	patient	control	OK	643.756	188.355	154.887	530.044	5,00E-05	0.0140016	yes					
6	ARHGEF19	ARHGEF19	ARHGEF19	chr1:1652459	patient	control	OK	357.156	112.397	-166.795	-60.984	5,00E-05	0.0140016	yes					
7	ARRDC1	ARRDC1	ARRDC1	chr9:1405000	patient	control	OK	699.613	139.438	0.994996	340.275	0.0003	0.0493798	yes					
8	ATL1	ATL1	ATL1	chr14:509997	patient	control	OK	76.407	249.703	-161.349	-441.121	5,00E-05	0.0140016	yes					
9	ATP6V0D2	ATP6V0D2	ATP6V0D2	chr8:8711113	patient	control	OK	141.496	602.768	-123.109	-398.314	0.0002	0.0383333	yes					
10	BCAT1	BCAT1	BCAT1	chr12:249625	patient	control	OK	191.027	586.341	-170.397	-622.628	5,00E-05	0.0140016	yes					
11	BMP2	BMP2	BMP2	chr20:674874	patient	control	OK	961.681	336.814	-151.361	-466.973	5,00E-05	0.0140016	yes					
12	BPIFB1	BPIFB1	BPIFB1	chr20:318705	patient	control	OK	36.651	761.925	10.558	395.806	0.00015	0.03125	yes					
13	C9orf152	C9orf152	C9orf152	chr9:1129618	patient	control	OK	126.281	293.493	121.669	444.641	5,00E-05	0.0140016	yes					
14	CCL5	CCL5	CCL5	chr17:341984	patient	control	OK	263.041	571.128	111.852	400.647	5,00E-05	0.0140016	yes					
15	CD109	CD109	CD109	chr6:7440362	patient	control	OK	245.725	920.931	-141.588	-522.319	5,00E-05	0.0140016	yes					
16	CEMIP	CEMIP	CEMIP	chr15:810717	patient	control	OK	838.152	219.043	-1.936	-651.489	5,00E-05	0.0140016	yes					
17	CHI3L1	CHI3L1	CHI3L1	chr1:203148C	patient	control	OK	521.762	786.714	-272.948	-907.519	5,00E-05	0.0140016	yes					
18	CITED4	CITED4	CITED4	chr1:4132672	patient	control	OK	550.335	145.476	14.024	520.669	5,00E-05	0.0140016	yes					
19	CNN1	CNN1	CNN1	chr19:116495	patient	control	OK	237.007	108.075	-11.329	-374.136	0.0003	0.0493798	yes					
20	COL10A1	COL10A1	COL10A1	chr6:1164215	patient	control	OK	168.116	608.885	-146.522	-479.166	5,00E-05	0.0140016	yes					
21	CRYAB	CRYAB	CRYAB	chr11:117775	patient	control	OK	741.219	291.357	-134.711	-423.311	5,00E-05	0.0140016	yes					
22	CYP4B1	CYP4B1	CYP4B1	chr1:4726466	patient	control	OK	382.688	134.253	181.071	526.254	5,00E-05	0.0140016	yes					
23	CYP4X1	CYP4X1	CYP4X1	chr1:4748922	patient	control	OK	430.477	106.728	130.993	392.501	0.00015	0.03125	yes					
24	DEGS2	DEGS2	DEGS2	chr14:100612	patient	control	OK	123.647	369.609	157.977	528.636	5,00E-05	0.0140016	yes					
25	DKK3	DKK3	DKK3	chr11:119845	patient	control	OK	48.43	250.027	-0.953814	-359.497	0.00025	0.043125	yes					
26	EDIL3	EDIL3	EDIL3	chr5:8323641	patient	control	OK	138.012	653.211	-107.917	-392.637	5,00E-05	0.0140016	yes					
27	EDNRA	EDNRA	EDNRA	chr4:148402C	patient	control	OK	525.935	269.993	-0.961961	-353.657	0.00015	0.03125	yes					
28	ELL2	ELL2	ELL2	chr5:952208	patient	control	OK	312.428	149.467	-10.637	-395.955	5,00E-05	0.0140016	yes					
29	ERBB2	ERBB2	ERBB2	chr17:378443	patient	control	OK	379.353	15.602	-20.451	195.473	0.0004655	0.00632554	yes					
30	ERMN	ERMN	ERMN	chr2:1581751	patient	control	OK	605.981	119.265	-23.451	-60.247	5,00E-05	0.0140016	yes					
31	FBXL16	FBXL16	FBXL16	chr16:742495	patient	control	OK	160.236	341.278	109.075	408.836	0.0001	0.0250727	yes					
32	FGR2	FGR2	FGR2	chr10:123237	patient	control	OK	116.387	364.171	164.569	40.751	5,00E-05	0.0140016	yes					
33	FXYD6	FXYD6	FXYD6	chr11:11769C	patient	control	OK	40.807	190.584	-109.839	-398.591	0.0001	0.0250727	yes					
34	GALNT15	GALNT15	GALNT15	chr3:1621618	patient	control	OK	233.808	821.554	-15.089	-556.106	5,00E-05	0.0140016	yes					
35	GALNT5	GALNT5	GALNT5	chr2:1581143	patient	control	OK	586.511	189.127	-163.281	-438.311	0.00015	0.03125	yes					
36	GFPT2	GFPT2	GFPT2	chr5:179727C	patient	control	OK	344.814	16.149	-109.437	-408.532	0.0001	0.0250727	yes					
37	GJB2	GJB2	GJB2	chr13:207616	patient	control	OK	309.823	147.739	-10.684	-390.409	0.0002	0.0383333	yes					
38	GOLM4	GOLM4	GOLM4	chr3:167727C	patient	control	OK	463.025	227.216	-102.703	-384.509	0.00015	0.03125	yes					
39	GPR68	GPR68	GPR68	chr14:916988	patient	control	OK	227.216	950.869	-125.674	-449.161	5,00E-05	0.0140016	yes					
40	GRAMP2	GRAMP2	GRAMP2	chr15:732521	patient	control	OK	515.005	142.751	-140.002	-500.002	5,00E-05	0.0140016	yes					

# Interconnection of RNA-Seq data



- DisGeNET (<http://www.disgenet.org/>)



Home About Search Browser Downloads Cytoscape RDF Help

One of the most challenging problems in biomedical research is to understand the underlying mechanisms of complex diseases. Great effort has been spent on finding the genes associated to diseases (Botstein and Risch, 2003; Kann, 2009). However, more and more evidences indicate that most human diseases cannot be attributed to a single gene but arise due to complex interactions among multiple genetic variants and environmental risk factors (Hirschhorn and Daly, 2005). Several databases have been developed storing associations between genes and diseases such as CTD™ (Davis, et al., 2014), OMIM® (Hamosh et al., 2005) and the NHGRI-EBI GWAS catalog (Welter et al., 2014). Each of these databases focuses on different aspects of the phenotype-genotype relationship, and due to the nature of the database curation process, they are not complete. Hence, integration of different databases with information extracted from the literature is needed to allow a comprehensive view of the state of the art knowledge within this research field. With this need in mind, we have created DisGeNET.

DisGeNET is a discovery platform integrating information on gene-disease associations (GDAs) from several public data sources and the literature (Píñero et al., 2015). The current version contains (DisGeNET v4.0) contains 429,036 associations, between 17,381 genes and 15,093 diseases, disorders and clinical or abnormal human phenotypes, and 72,870 variant-disease associations (VDAs), between 46,589 SNPs and 6,356 phenotypes. Given the large number of GDAs compiled in DisGeNET, we have also developed a score in order to rank the associations supporting evidence. Importantly, useful tools have also been created to explore and analyze the data contained in DisGeNET. DisGeNET can be queried through Search and Browse functionalities available from this web interface, or by a plugin created for Cytoscape to query a network representation of the data. Moreover, DisGeNET data can be queried by downloading the SQLite database to your local machine. Furthermore, an RDF (Resource Description Framework) representation of DisGeNET database is also available. It can be queried using an endpoint and a Faceted Browser. Follow the link for more information.

DisGeNET database has been cited by several papers. Some of them can be reviewed [here](#).

The DisGeNET database is made available under the [Open Database License](#). Any rights in individual contents of the database are licensed under the [Database Contents License](#).

Tweets by @DisGeNET

 DisGeNET  
@DisGeNET

Check out the new publication describing the DisGeNET platform in NAR database issue [nar.oxfordjournals.org/con](http://nar.oxfordjournals.org/con)



# Linking and integrating data

- David (<https://david.ncifcrf.gov/list.jsp>)

**Gene Name Batch Viewer**  
DAVID Bioinformatics Resources 6.8, NIAID/NIH

Home Start Analysis Shortcut to DAVID Tools Technical Center Downloads & APIs Term of Service Why DAVID? About Us

\*\*\* Welcome to DAVID 6.8 with updated Knowledgebase ([more info](#)). \*\*\*  
\*\*\* If you are looking for [DAVID 6.7](#), please visit our [development site](#). \*\*\*

**Upload List Background**

**Upload Gene List**

[Demolist 1](#) [Demolist 2](#)  
[Upload Help](#)

**Step 1: Enter Gene List**

A: Paste a list

```
ERBB3  
ERBB4  
ERC1  
ERC2  
ERC2-IT1
```

Or

B: Choose From a File

No file selected.  
 Multi-List File ?

**Step 2: Select Identifier**

OFFICIAL\_GENE\_SYMBOL

**Step 3: List Type**

Gene List   
Background

**Step 4: Submit List**

**Gene Name Batch Viewer**

Submit your gene list to start !

Tell us how you like the tool  
Read technical notes of the tool  
Contact us for questions

What does this tool do?

- Quickly translate given gene IDs to corresponding gene names in a batch way
- Provide links for each genes to DAVID Gene Report for in-depth information
- Search functionally related genes within user's input gene list or genome

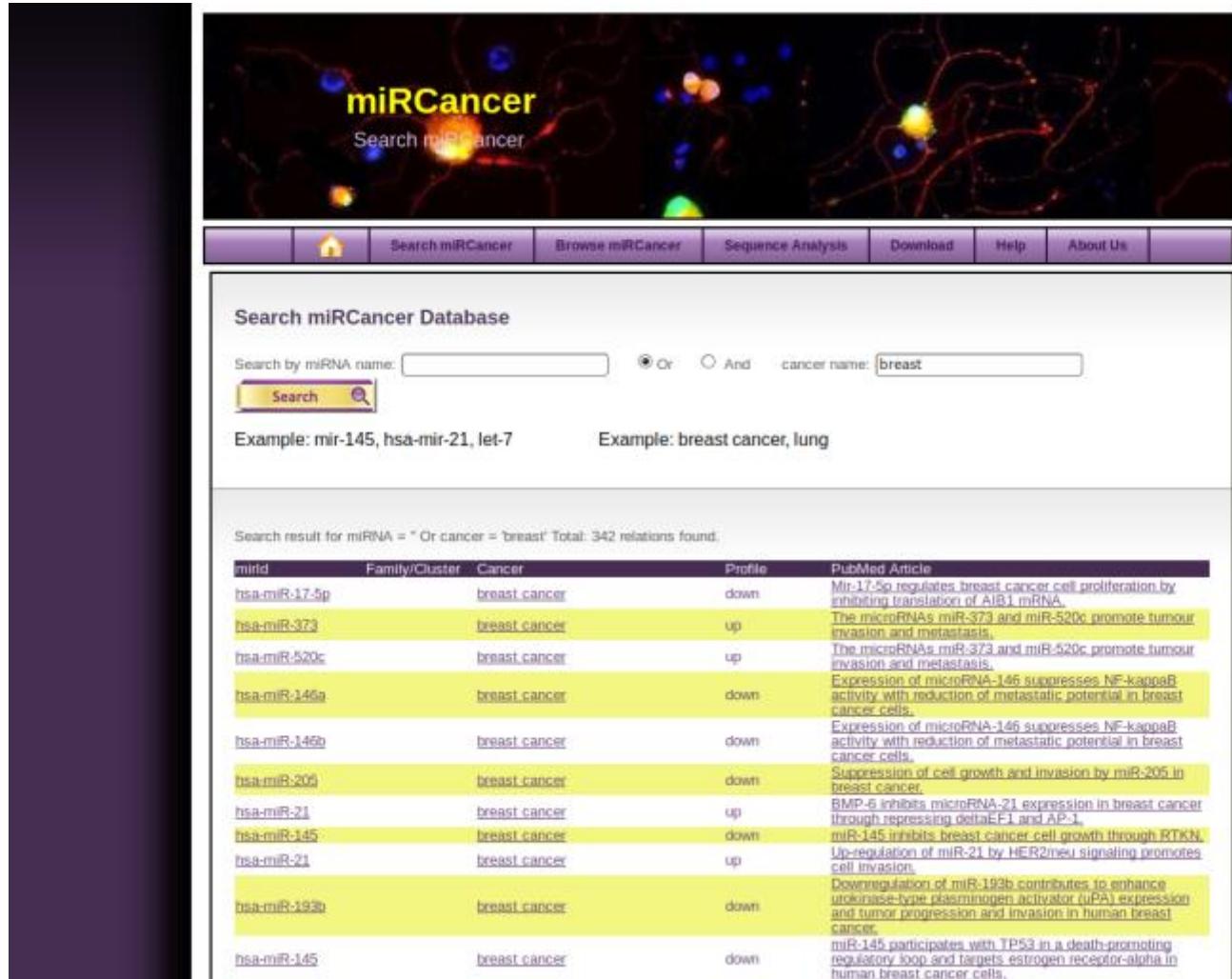
Key Concepts of "Search Related Genes"

Any given gene is associating with a set of annotation terms. If genes share similar set of those terms (annotation profile), they are most likely involved in similar biological mechanisms. The algorithm adopts kappa statistics to quantitatively measure the degree of the agreement how genes share the ~75,000 annotation terms collected by DAVID knowledgebase. For any given gene(s), the tool instantly searches and lists the related genes passed kappa similarity measurement threshold. The searching scope could be within user's input gene list, selected genome or all genomes (~1.2 million genes) as user's choice.

Find Related Genes Tool is very different and complementary to the common gene clustering methods, such as homologous genes based on sequence similarity; protein families based on one common biological activity. The approach provides researchers a new way to group those functional related genes by measuring the similarity of their global annotation profile, which facilitates new understanding of the biological network. [More](#)



- miRCancer db - Find up regulated miRNAs (<http://mircancer.ecu.edu/index.jsp>)



Search miRCancer Database

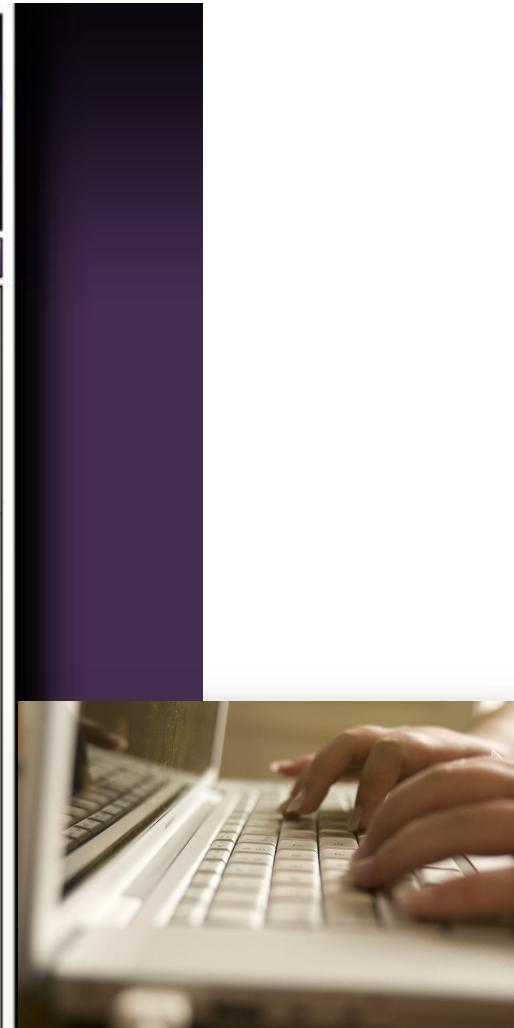
Search by miRNA name:   Or  And cancer name:  breast

Search 

Example: mir-145, hsa-mir-21, let-7      Example: breast cancer, lung

Search result for miRNA = " Or cancer = 'breast' Total: 342 relations found:

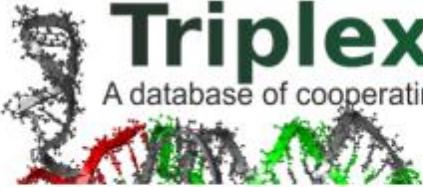
mirid	Family/Cluster	Cancer	Profile	PubMed Article
hsa-miR-175p		breast cancer	down	Mir-175p regulates breast cancer cell proliferation by inhibiting translation of AIB1 mRNA.
hsa-miR-373		breast cancer	up	The microRNAs miR-373 and miR-520c promote tumour invasion and metastasis.
hsa-miR-520c		breast cancer	up	The microRNAs miR-373 and miR-520c promote tumour invasion and metastasis.
hsa-miR-146a		breast cancer	down	Expression of microRNA-146 suppresses NF-kappaB activity with reduction of metastatic potential in breast cancer cells.
hsa-miR-146b		breast cancer	down	Expression of microRNA-146 suppresses NF-kappaB activity with reduction of metastatic potential in breast cancer cells.
hsa-miR-205		breast cancer	down	Suppression of cell growth and invasion by miR-205 in breast cancer.
hsa-miR-21		breast cancer	up	BMP-6 inhibits microRNA-21 expression in breast cancer through repressing deltaEF1 and AP-1.
hsa-miR-145		breast cancer	down	miR-145 inhibits breast cancer cell growth through RDNK.
hsa-miR-21		breast cancer	up	Up-regulation of miR-21 by HER2/neu signaling promotes cell invasion.
hsa-miR-193b		breast cancer	down	Downregulation of miR-193b contributes to enhance urokinase-type plasminogen activator (uPA) expression and tumor progression and invasion in human breast cancer.
hsa-miR-145		breast cancer	down	miR-145 participates with TP53 in a death-promoting regulatory loop and targets estrogen receptor-alpha in human breast cancer cells.



miRNA	Regulation	Target
140-5b	up	
148b	Down	
150	Up	
106b	Up	
143	Down	
19b	Up	
21	up	
...	...	...

# Explore miRNA cooperativity to justify diagnosis

- TriplexRNA database (<https://www.sbi.uni-rostock.de/triplexrna/>)



## TriplexRNA

A database of cooperating microRNAs and their mutual targets

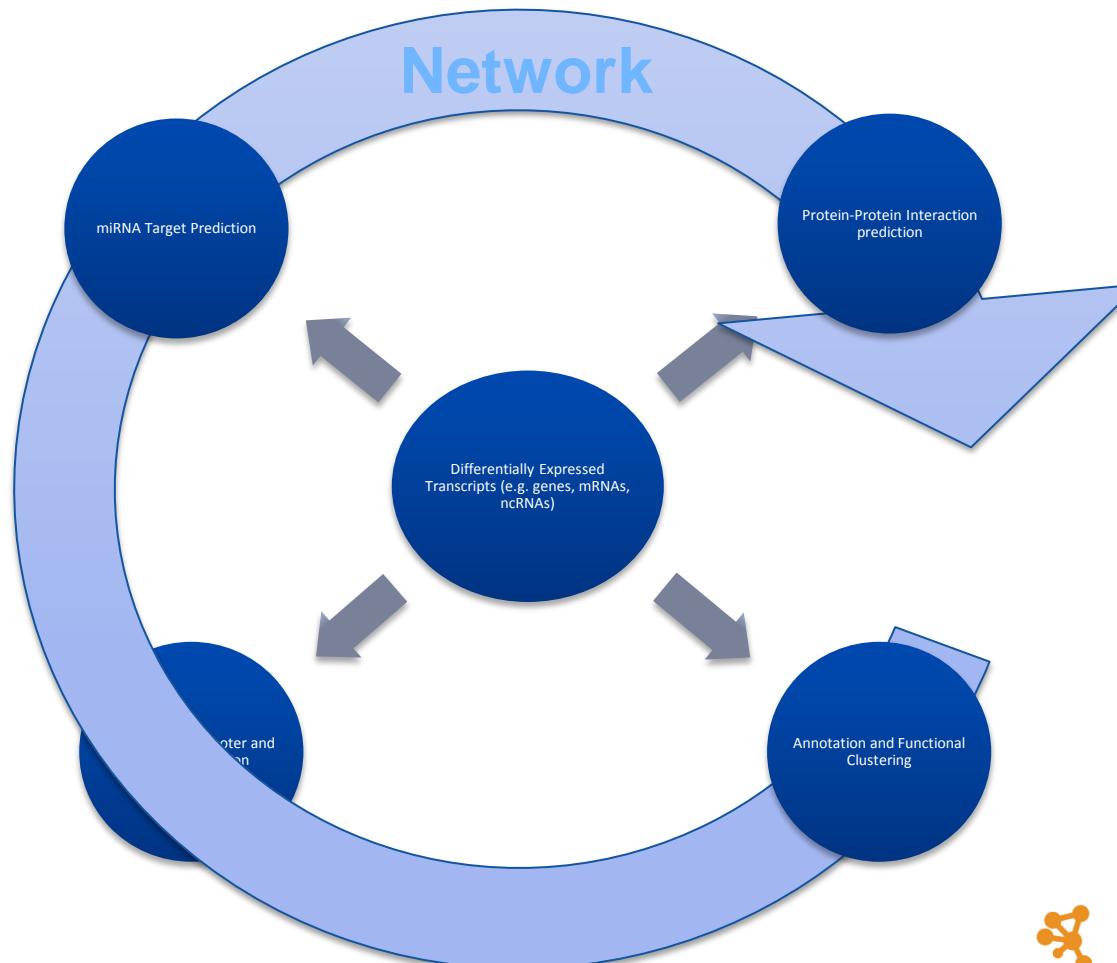
Search targets of synergistic microRNA regulation

Search in: Human for miRNA ID: hsa-miR-140-5p

Gene ID	RefSeq ID	miRNA1 ID	miRNA2 ID	Seed distance (nt)	Free energy (Kcal/mol)	Energy gain (Kcal/mol)	Triplex details
ADCY6	NM_015270	hsa-miR-197	hsa-miR-140-5p	23	-48.66	-14.38	<a href="#">more &gt;</a>
ATG4B	NM_178326	hsa-miR-140-5p	hsa-miR-346	28	-47.36	-15.58	<a href="#">more &gt;</a>
ZNF705A	NM_001004328	hsa-miR-140-5p	hsa-miR-296-3p	17	-43.76	-14.28	<a href="#">more &gt;</a>
FGR	NM_005248	hsa-miR-140-5p	hsa-miR-326	33	-43.56	-11.58	<a href="#">more &gt;</a>
PTCD1	NM_015545	hsa-miR-140-5p	hsa-miR-339-5p	34	-43.26	-12.98	<a href="#">more &gt;</a>
AARS	NM_001605	hsa-miR-24	hsa-miR-140-5p	32	-43.16	-17.18	<a href="#">more &gt;</a>
WEE1	NM_003300	hsa-miR-15b	hsa-miR-140-5p	16	-42.86	-16.28	<a href="#">more &gt;</a>
WNT1	NM_005430	hsa-miR-31	hsa-miR-140-5p	28	-42.56	-12.78	<a href="#">more &gt;</a>
ZBTB9	NM_152735	hsa-miR-140-5p	hsa-miR-296-3p	29	-41.96	-11.68	<a href="#">more &gt;</a>
ADRA1A	AY491776	hsa-miR-140-5p	hsa-miR-150	21	-41.96	-12.18	<a href="#">more &gt;</a>



# There is nothing more practical than a network

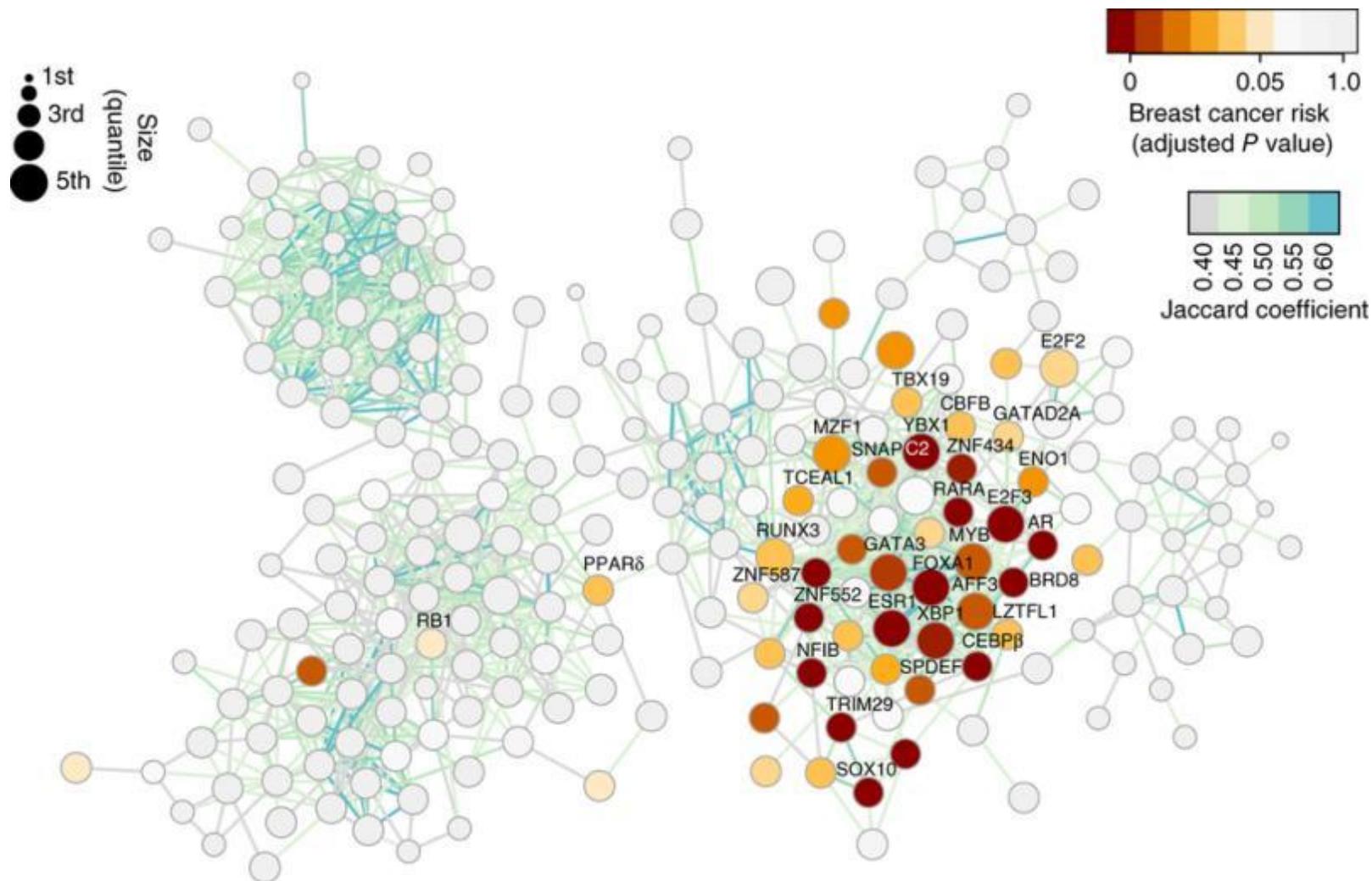


Cytoscape



Vanted / CellDesigner

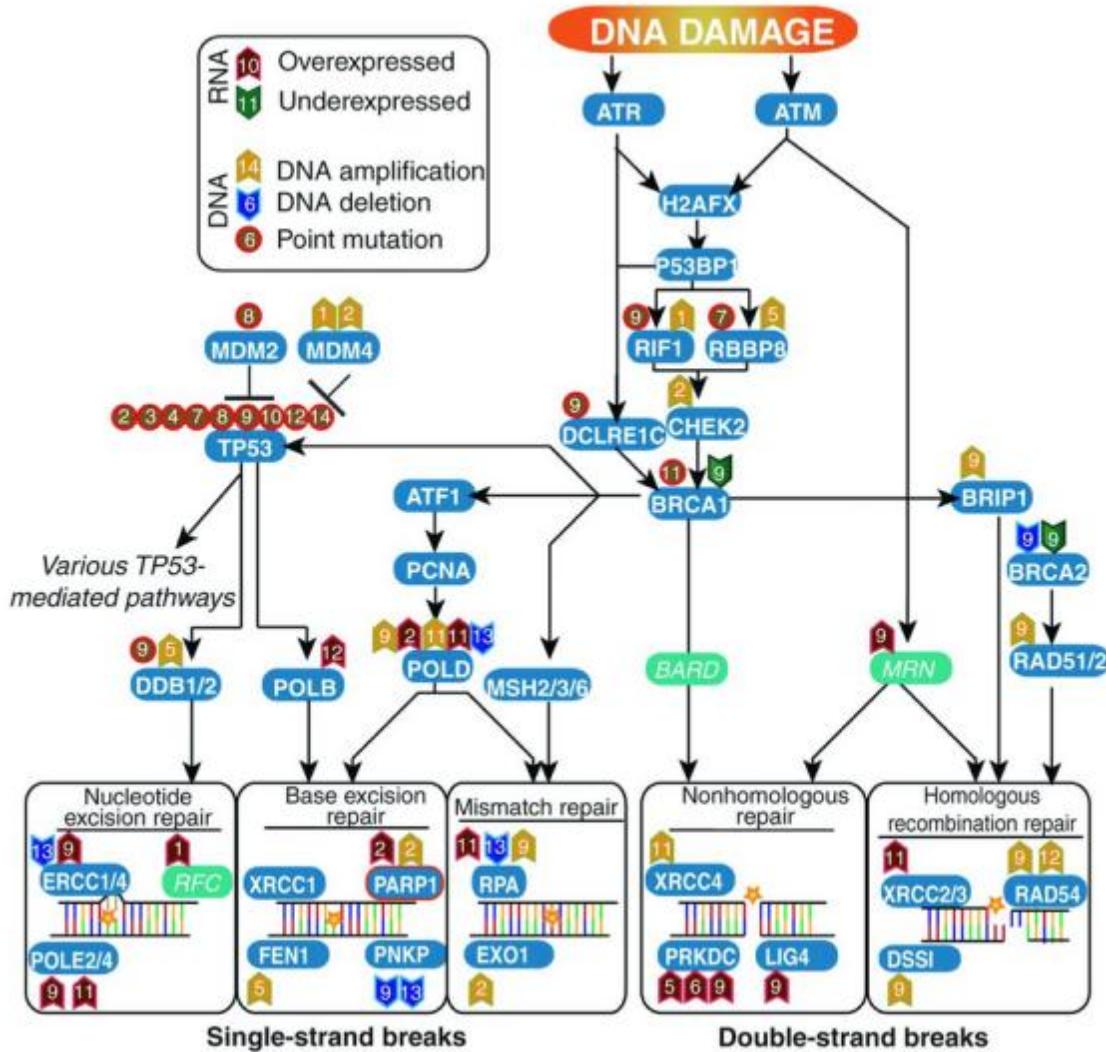
# Network comparisons also reveal differences



Castro, Nat. Gen, 2016

Cytoscape

# What else is done in this field with NGS?



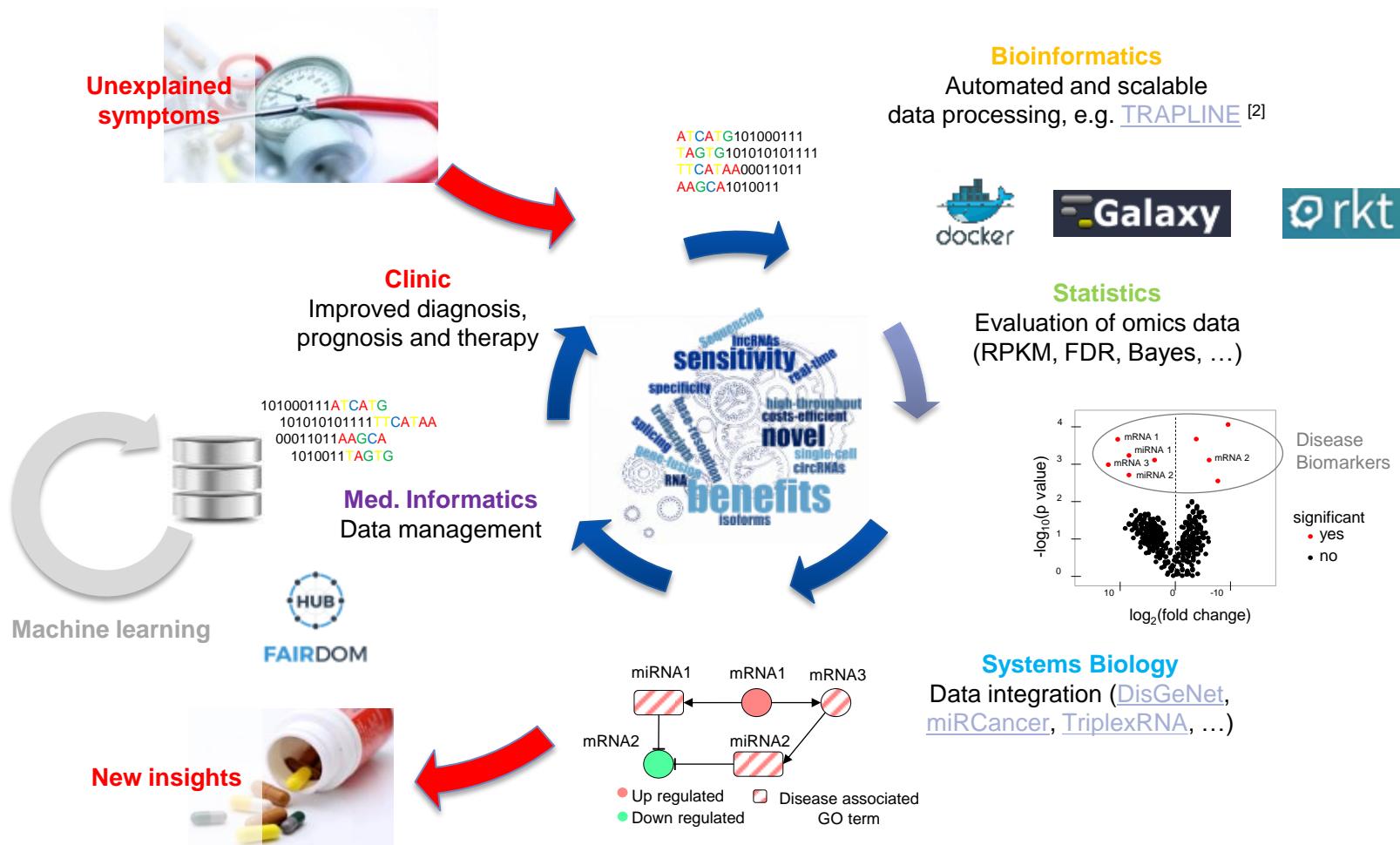
- Craig, D. W. et al. Genome and transcriptome sequencing in prospective metastatic triple-negative breast cancer uncovers therapeutic vulnerabilities. *Mol. Cancer Ther.* 12, 104–116 (2013). One of the first papers investigating integration of whole-transcriptome sequencing and genome sequencing for targeted therapy selection in advanced metastatic triple-negative breast cancer

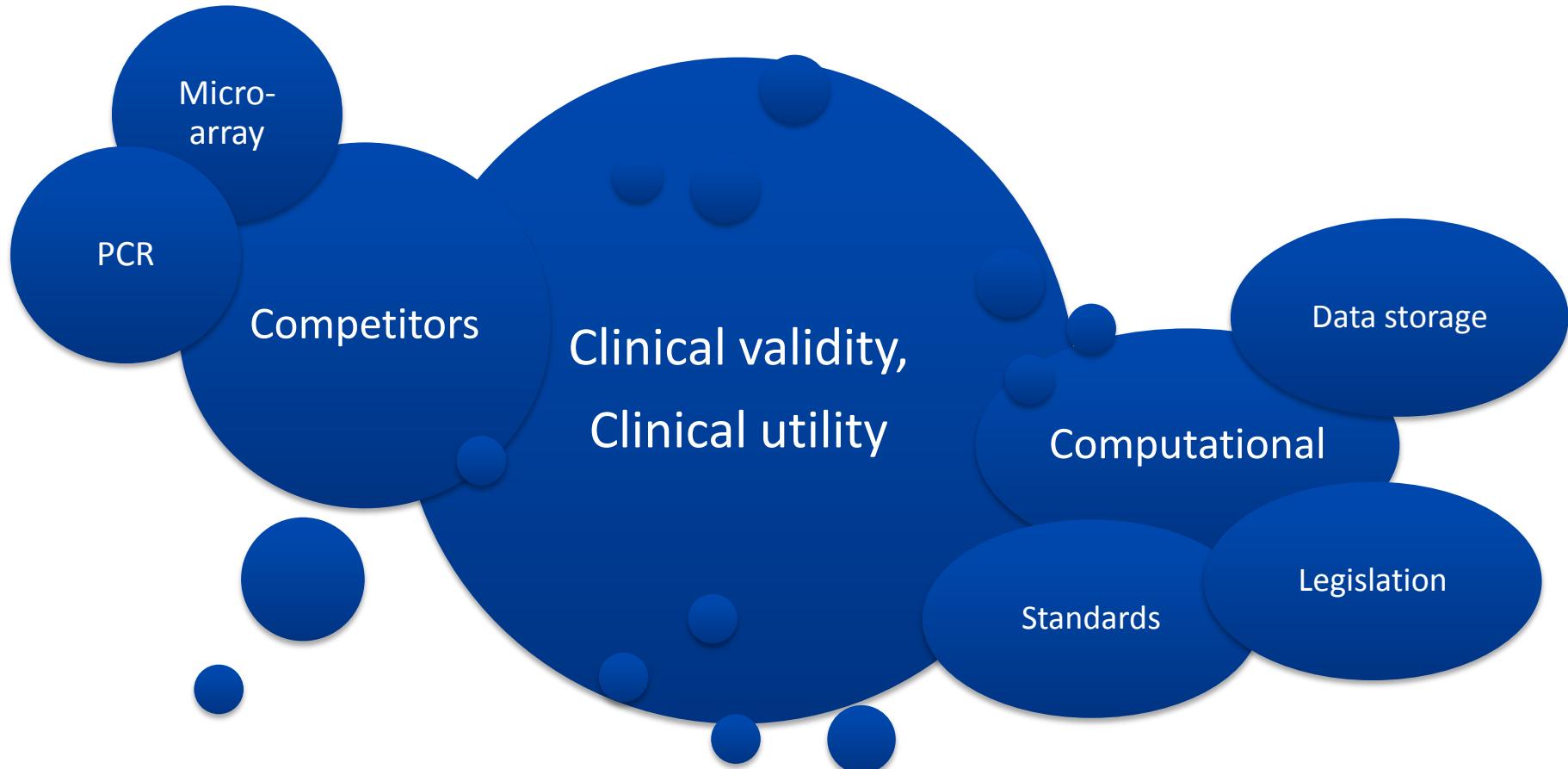
# Does the patient has a high risk of cancer

Patient



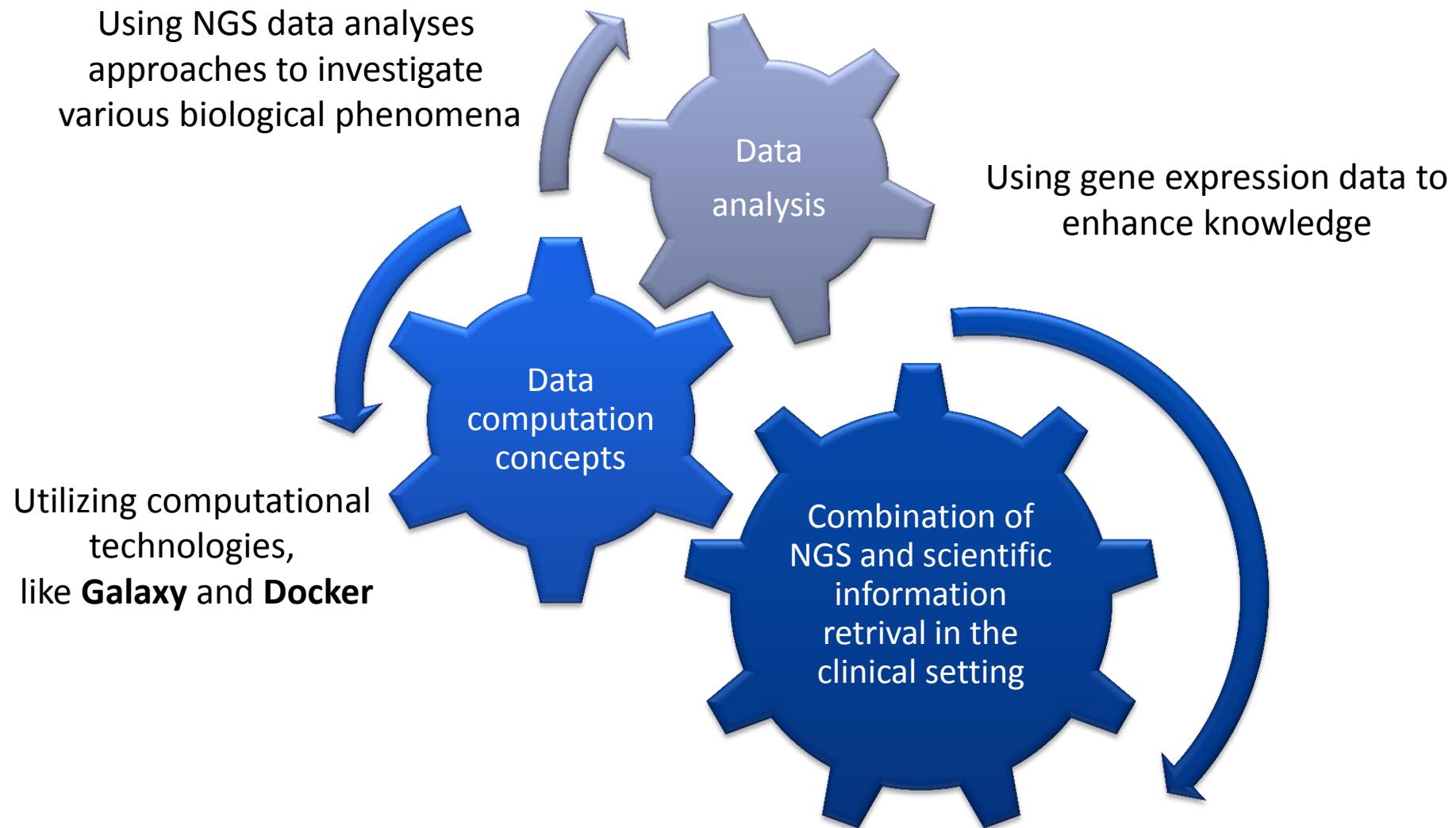
# Our implementation strategy





- “With its unprecedented ability to simultaneously detect global gene transcript levels and diverse RNA species, RNA-seq has the potential to revolutionize clinical testing for a wide range of diseases.” Byron *et al.*, *Nat. Rev. Genet.*, 2016
- “Once the discovery phase is complete, many diagnostic tests will become targeted assays, sensitive enough to detect small numbers of rare transcripts.”  
Andersson *et al.*, *Nat. Genet.*, 2015
- “Feed in latest scientific findings and analyze the same dataset over and over again [...]”. Comment on crowdsourced research in Medicine (*Nature*)
  - “Value of incorporating RNA sequencing (RNA-seq) with DNA sequencing to evaluate the expression of mutant alleles, to detect both known and novel gene fusions, and to detect splice variants.” Robinson *et al.*, *Cell*, 2015

# What did we learn so far?



### Part 7

16:00 – 16:30

“Open discussions about RNA-Seq projects from the participants”



- What data do you have?
- What do you want to analyze?
- Do you feel better informed?



[www.denbi.de/index.php/training-courses](http://www.denbi.de/index.php/training-courses)



The screenshot shows the de.NBI website's main navigation bar at the top with links for Home, Mission, Organization, Network, Services, Training, Cloud, Events, News, Jobs, and Help. Below the navigation is a large blue header banner with the text "Training Courses 2018" in white. To the right of the banner is a graphic of three stylized nodes connected by lines. The main content area below the banner lists various training events for 2018.

## Training Courses 2018

Training Archive sorted according to date

Training Archive sorted according to de.NBI units

Online training & Media library

de.NBI Youtube channel



2018	Topic	Location
05-08 Mar	<a href="#">de.NBI Winter School on Computational Metabolomics</a>	Lutherstadt Wittenberg
06 Mar	<a href="#">Introduction to genome-wide association studies (GWAS)</a>	Kiel
07 Mar	<a href="#">Introduction to RNA-Seq data analysis with Galaxy</a>	Kiel
07-13 Mar	<a href="#">Applied Metaproteomics Workshop 2018</a>	Magdeburg
08 Mar	<a href="#">Getting Started with the de.NBI-Cloud</a>	Bielefeld
09 Mar	<a href="#">Analysis of Mass Spectrometry and Sequence Data with KNIME - KNIME Spring Summit</a>	Berlin
11 Mar	<a href="#">Computational Proteomics Workshop - DGMS2018</a>	Saarbrücken
11-18 Mar	<a href="#">SeqAn Developer Retreat</a>	Mallorca
27-29 Mar	<a href="#">3rd de.NBI Training Course on Metagenome Analysis</a>	Bielefeld
02-06 Apr	<a href="#">OpenMS Developer Meeting 2018</a>	Troia
04-05 Apr	<a href="#">Introduction to Python Programming</a>	Heidelberg
11-12 Apr	<a href="#">Introduction to BRENDA &amp; ProteinsPlus 2018</a>	Braunschweig
23-25 Apr	<a href="#">4th de.NBI Genomics training course</a>	Gießen
24-26 Apr	<a href="#">Tools for Systems biology modeling and data exchange: COPASI, CellNetAnalyzer, SABIO-RK, SEEK</a>	Magdeburg

## Welcome to Galaxy Training!

<http://galaxyproject.github.io/training-material/>

Collection of tutorials developed and maintained by the worldwide Galaxy community

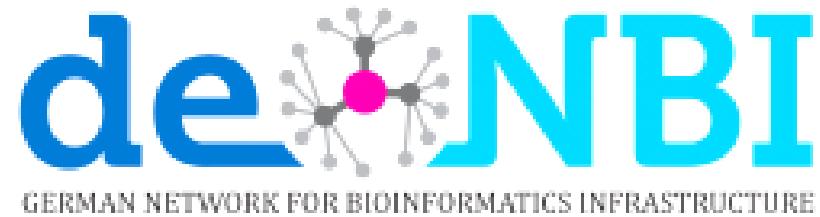
### Galaxy for Scientists

Topic	Tutorials
Introduction to Galaxy	13
Assembly	3
ChIP-Seq data analysis	2
Epigenetics	2
Metagenomics	2
Proteomics	8
Sequence analysis	6
Transcriptomics	5
Variant Analysis	6

### Galaxy for Developers and Admins

Topic	Tutorials
Galaxy Server administration	8
Development in Galaxy	14
Train the trainers	6





Evaluation at :

<https://de.surveymonkey.com/r/denbi-course?sc=rbc&id=000111>

# Acknowledgements



Olaf Wolkenhauer (University of Rostock)

Wolfgang Hess (University of Freiburg)

Steve Hoffmann (University of Leipzig)

Rolf Backofen (University of Freiburg)

Björn Grüning (University of Freiburg)



Supported by:



[elixir-europe.org](http://elixir-europe.org)



Bundesministerium  
für Bildung  
und Forschung

[bmbf.de](http://bmbf.de)

We hope you enjoyed the training!

**Universität  
Rostock**



Traditio et Innovatio



**SYSTEMS BIOLOGY  
BIOINFORMATICS  
ROSTOCK**

**de**NBI  
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

