

WS15: RNA-Seq data analysis with Galaxy for clinical applications

Markus Wolfien and Andrea Bagnacani

Galaxy Training – 4th September 2018 Osnabrück

www.sbi.uni-rostock.de



SYSTEMS BIOLOGY
BIOINFORMATICS
ROSTOCK



GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE



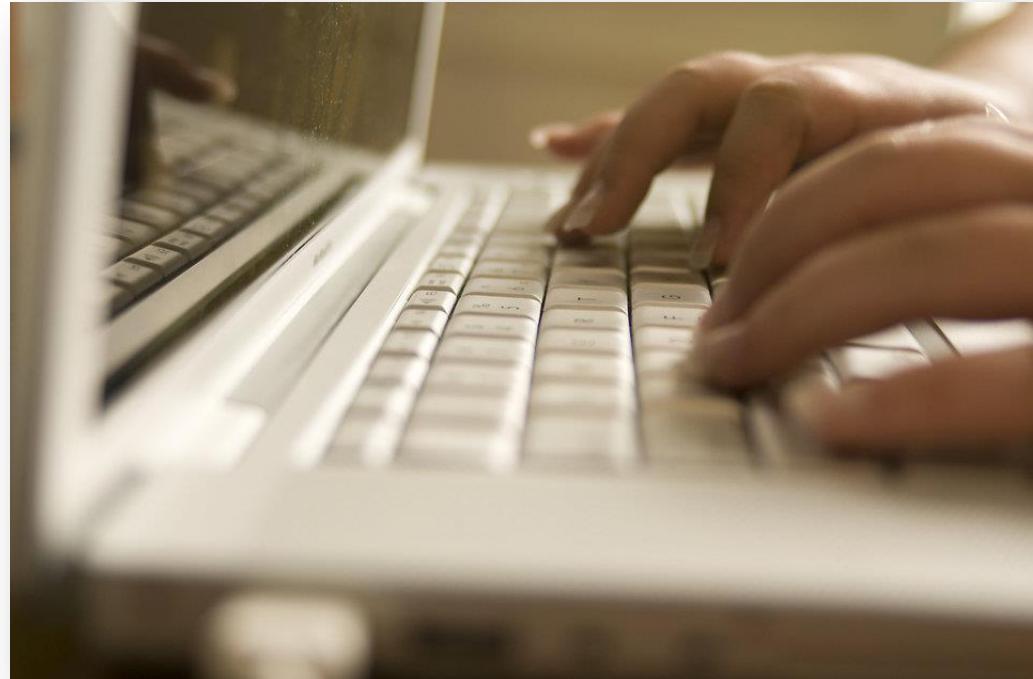
Deutsche Gesellschaft für
Medizinische Informatik,
Biometrie und
Epidemiologie e.V.

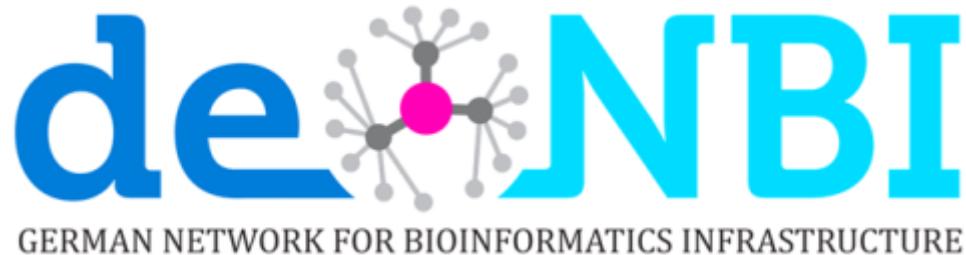
- 8.30 – 10.00
 - Introduction to RNA-Seq and Galaxy
 - Quality control of fastq datasets
- 11.15 – 12.45
 - RNA-Seq mapping algorithms
 - Quantification of alignment files
- 14.30 – 16.00
 - Case study of RNA-Seq patient data
 - Workflow development with Galaxy
- 16.15 -17.30
 - Filtering and visualisation with Galaxy
 - Gene set enrichment analyses

Go and get the slides!

Slides at:

<https://github.com/destairdenbi/trainings>





GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE



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

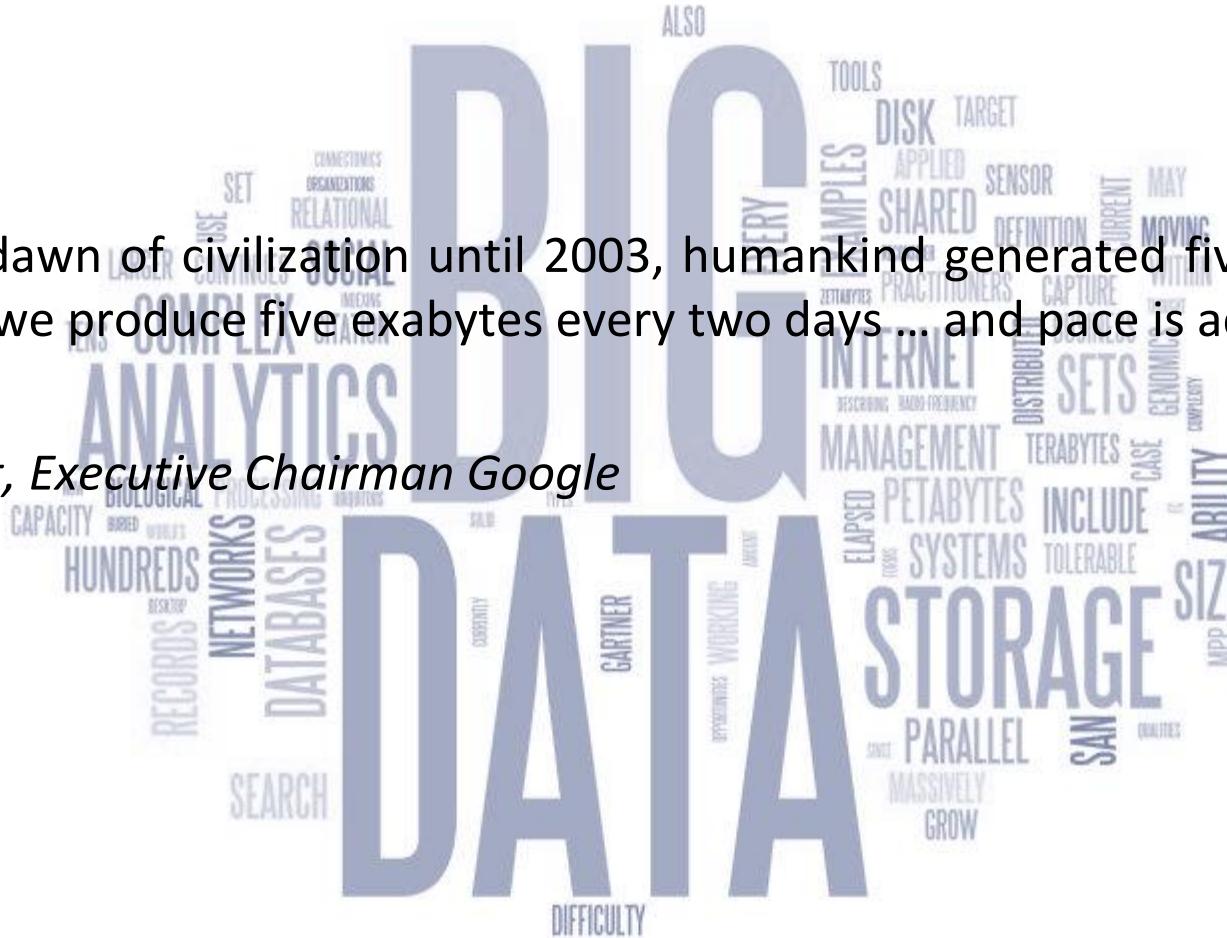
Our aim is to enable a comprehensive **analysis of RNA-Seq experiments as a service**. To enable maximum usefulness, interconnectivity, and accessibility for the developed approaches and services, we will provide dedicated **workshops, training programs and screen casts** for bioinformaticians and other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

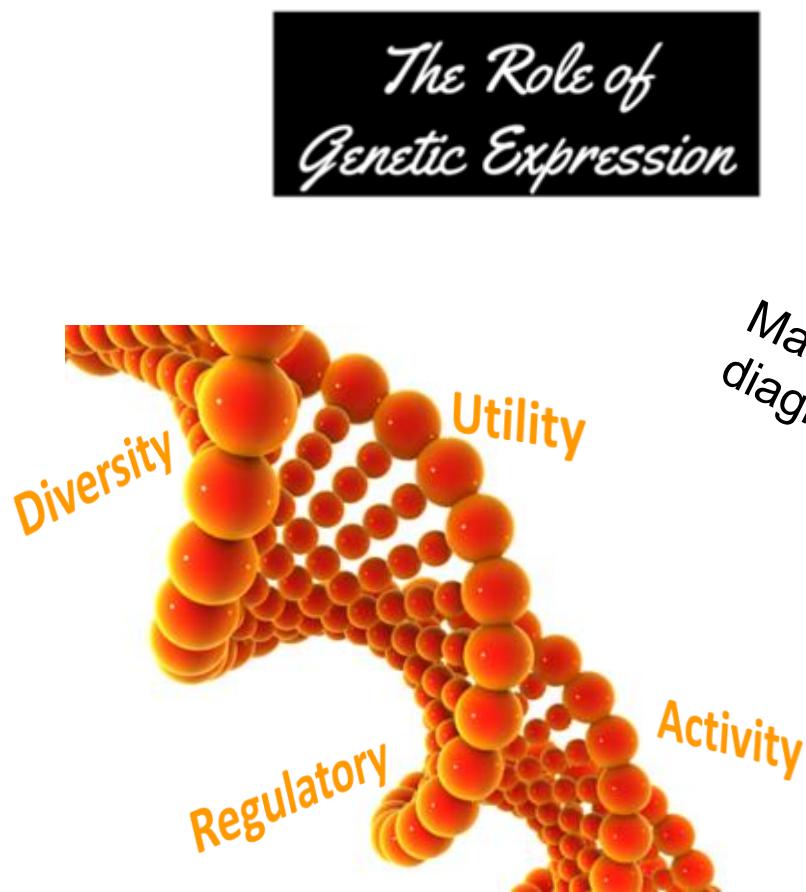
What is big data?



“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



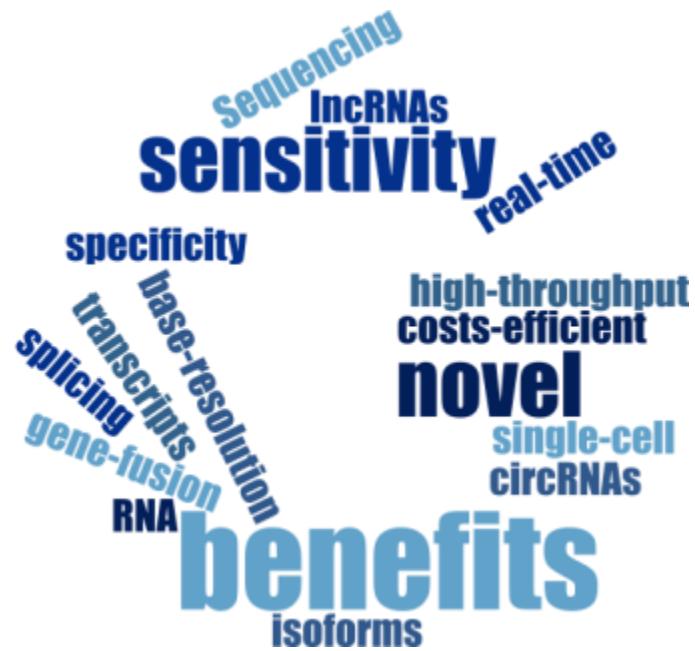


Many different variations and subtypes

Information about regulatory mechanisms

Active and measurable state of the cell ...

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

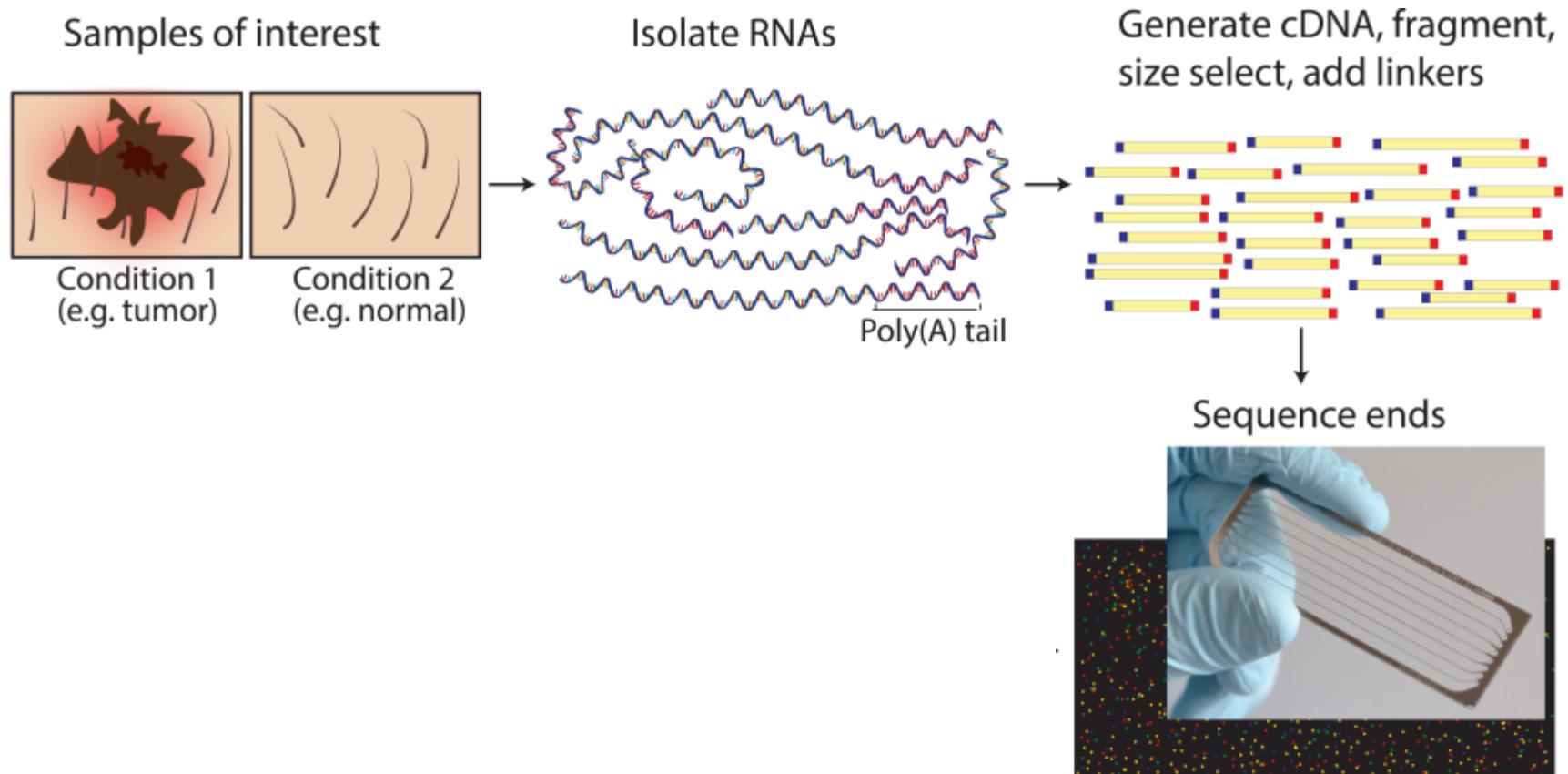


“RNA-Seq is able to identify thousands of differentially expressed genes, tens of thousands of differentially expressed gene isoforms and can detect mutations and germline variations for hundreds to thousands of expressed genetic variants, as well as detecting chimeric gene fusions, transcript isoforms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

Wang, *Nat Rev. Genet.*, 2009



From sample to readout



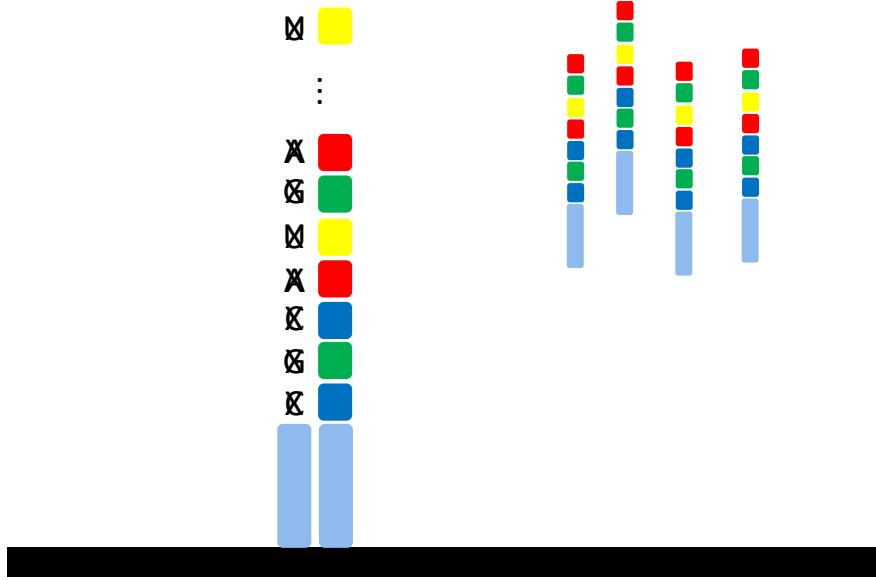
100s of millions of paired reads
10s of billions bases of sequence

Griffith, Plos Comp. Biol., 2015

How NGS works - brief



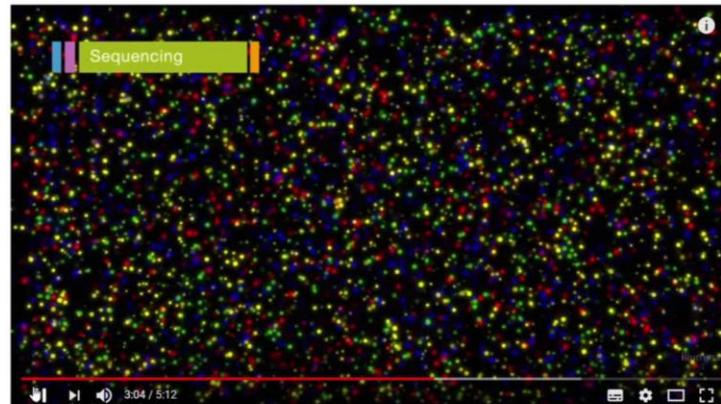
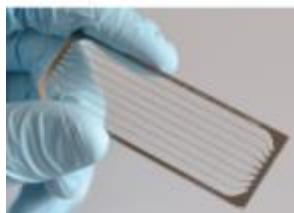
Example: Sequencing



>25 *10⁶ Sequences

RNA-Sequencing
RNA sample of the patient

Adapter
Flow cell

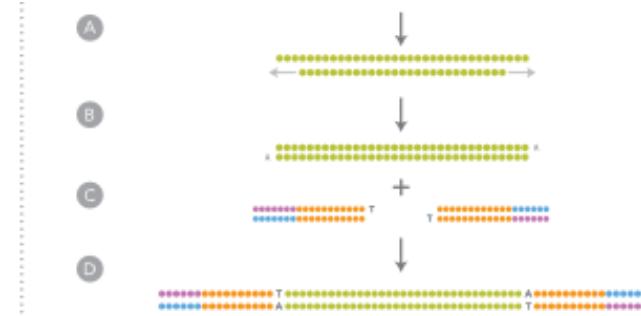


How do I get my NGS data - detailed?



1 Library Preparation

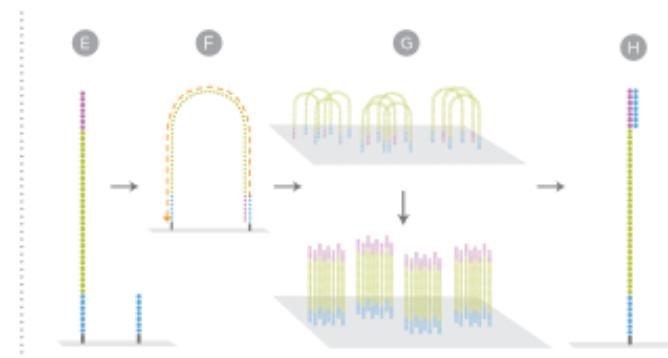
6 hours
3 hours hands-on time



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

2 Cluster Generation

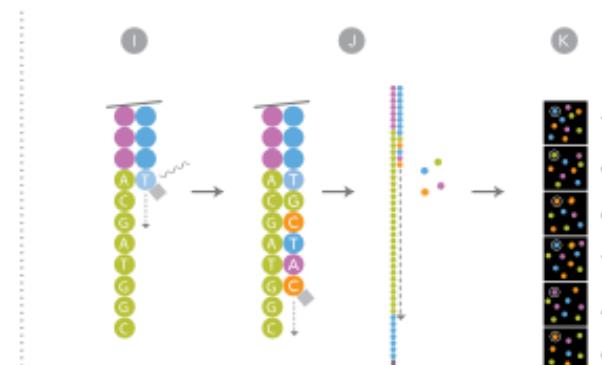
4 hours
< 10 minutes hands-on time
1–96 samples



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

3 Sequencing

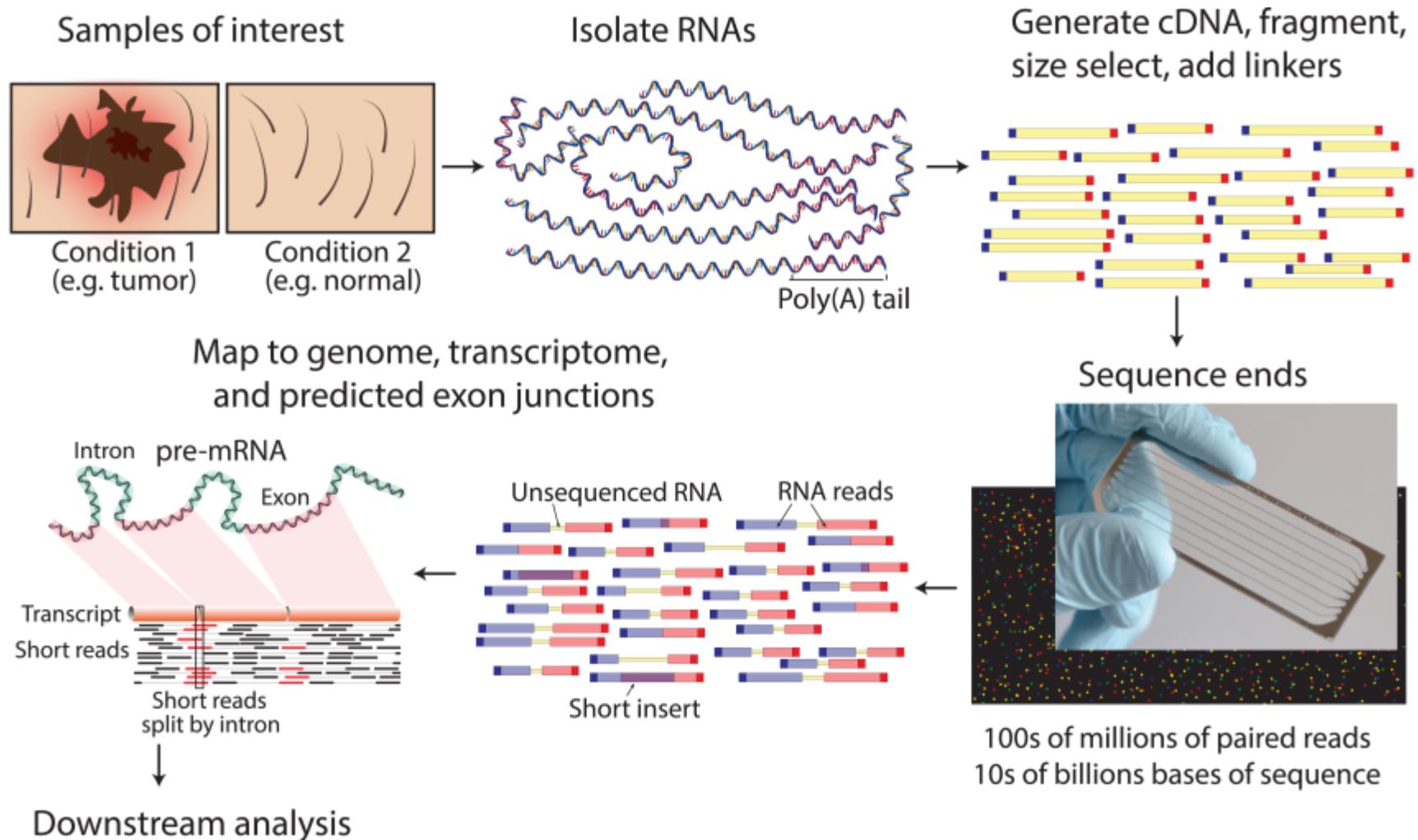
1–3 days single-read run
3–9 days paired-end run
30 minutes hands-on time
8 lanes, up to 96 samples per flow cell (run)



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



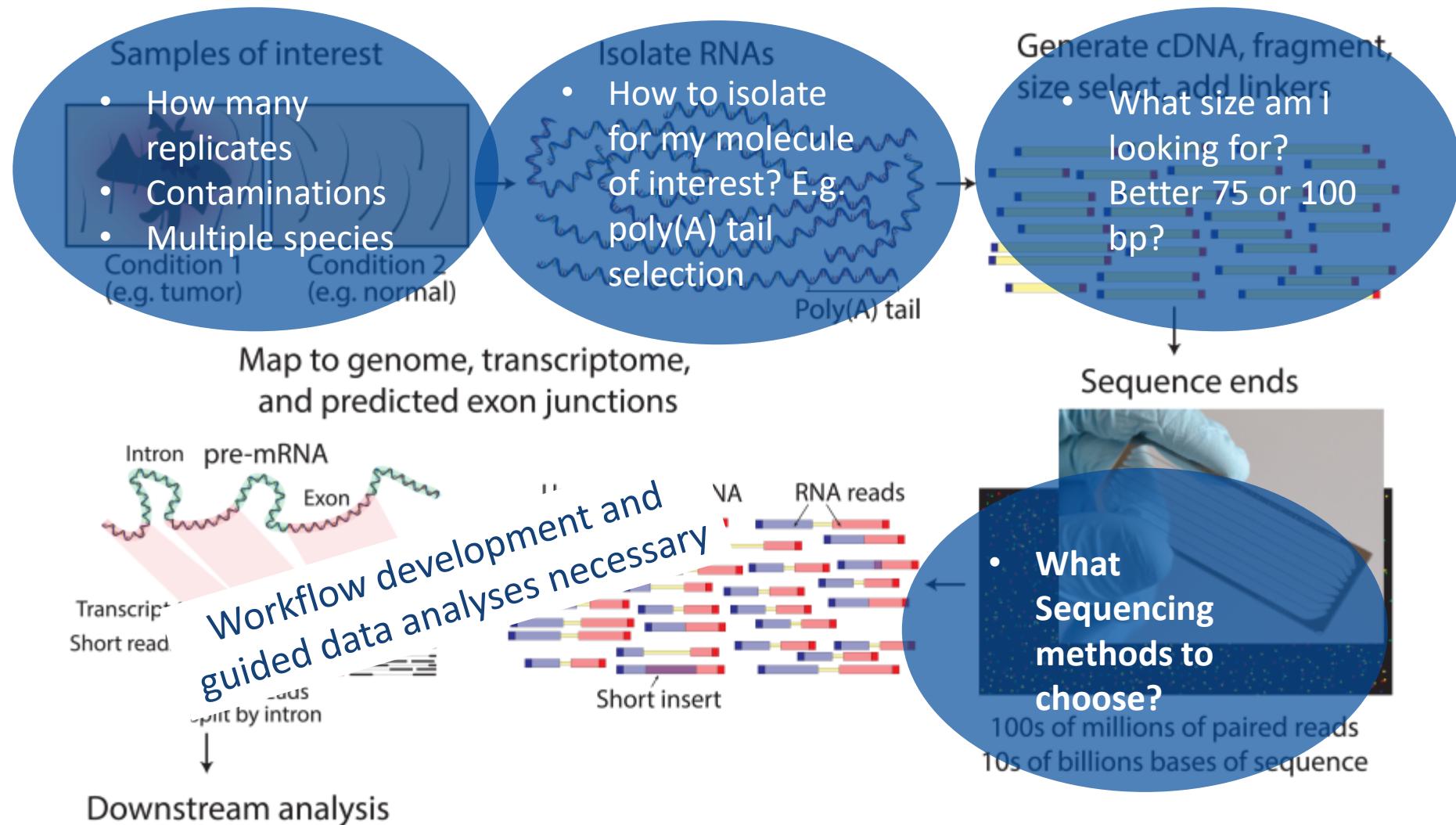
From sample to readout



Griffith, Plos Comp. Biol., 2015



From sample to readout

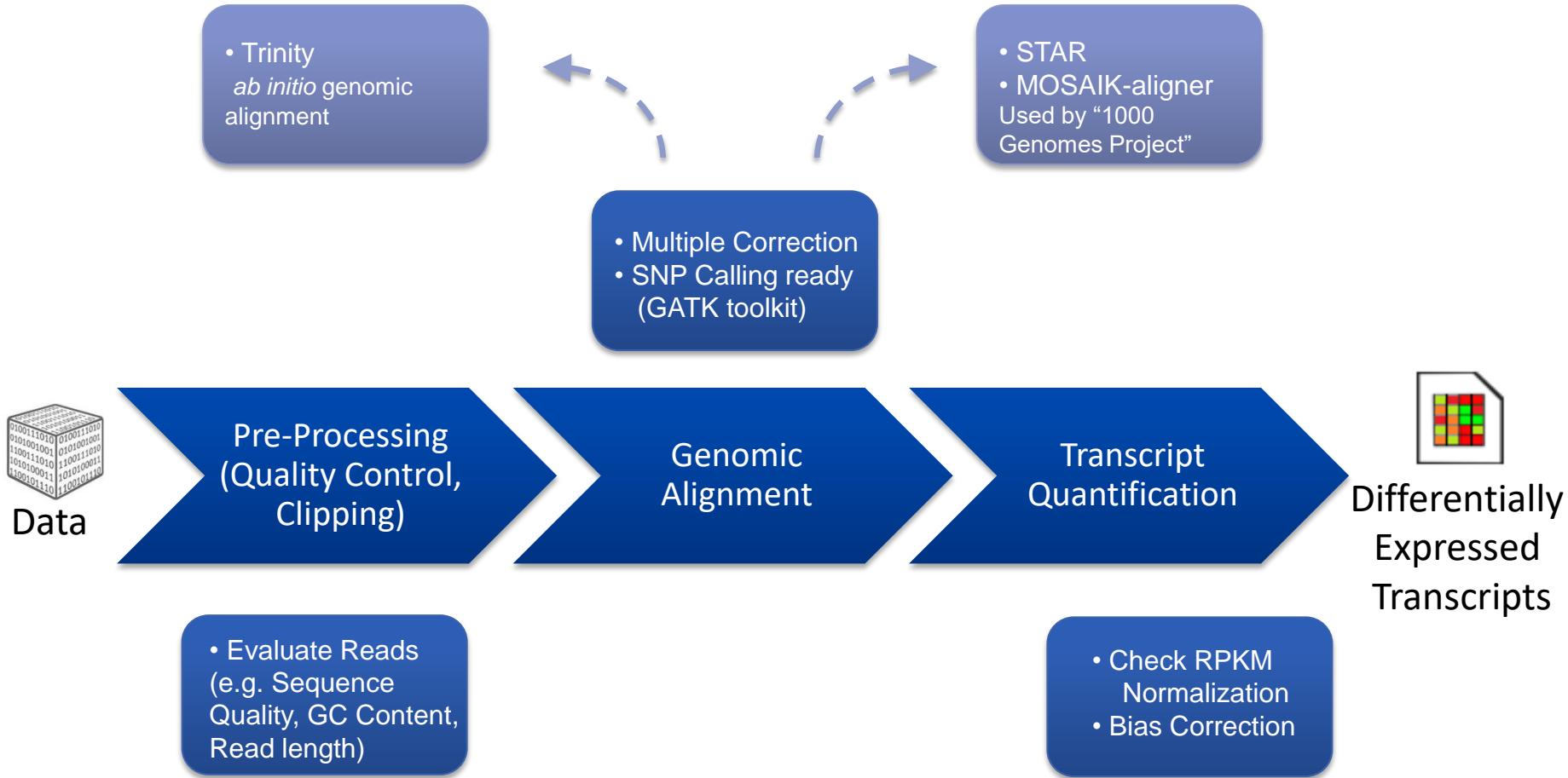


Griffith, Plos Comp. Biol., 2015

Data analysis



Differential expression



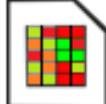
Differential expression



Pre-Processing
(Quality Control,
Clipping)

Genomic
Alignment

Transcript
Quantification



Differentially
Expressed
Transcripts

Condition A
Multiple Copies of a Transcriptome

Condition B

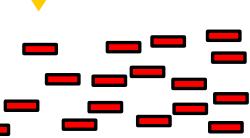
Condition C



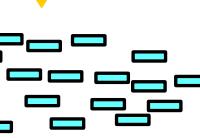
Reads



Reads



Reads



Ref. Sequence



Reads on Ref. Seq.



Differential expression

Big Data and the need for new analyses



Grape

big.crg.cat/services/grape



mapman.gabipd.org



Chipster
Open source platform for data analysis
chipster.csc.fi



r-project.org/



gene-talk.de



illumina.com



bioconductor.org



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

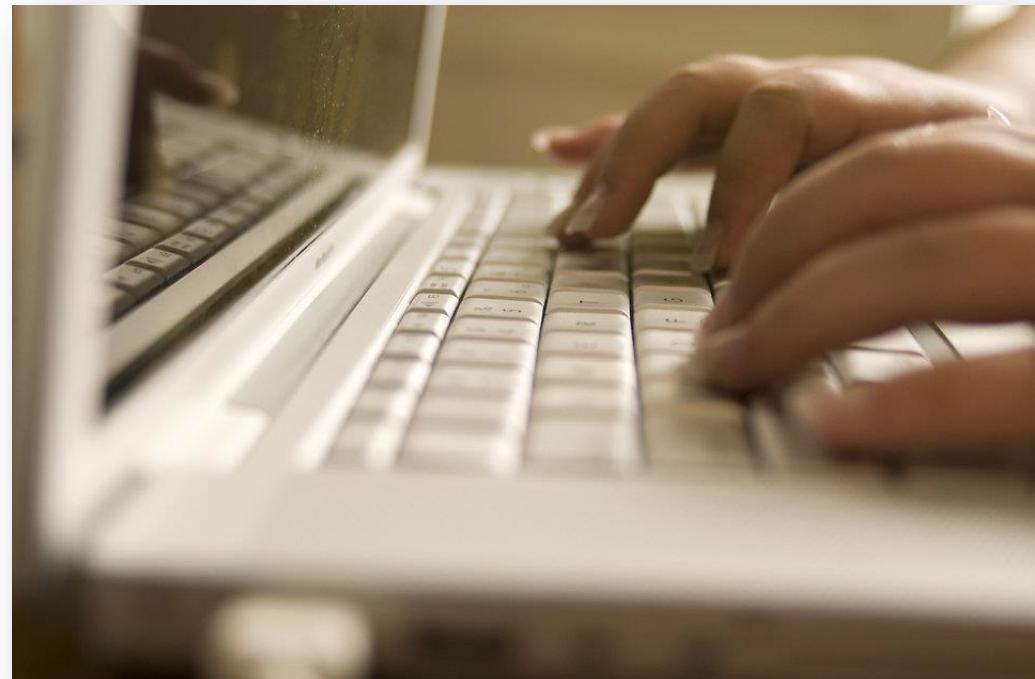
Manuscript: <https://doi.org/10.1016/j.cels.2018.05.012>

Hands on part 1: 9:05 – 10:00

“Introduction to Galaxy & RNA-Seq data preprocessing and quality control”

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

Please visit and explore Galaxy
usegalaxy.org

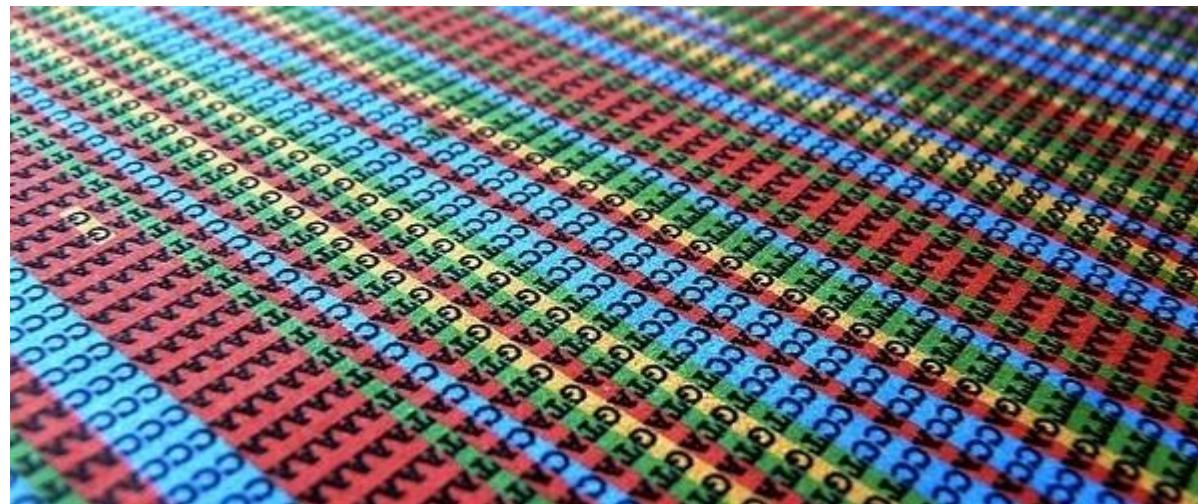


Hands on part 2

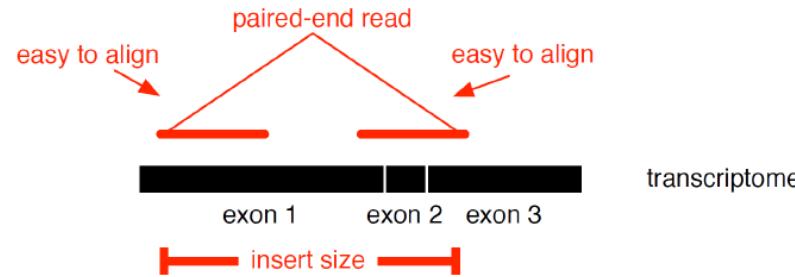
11:15 – 12:45

Introduction

“Application of different read mapping approaches for genomic alignment and subsequent quantification”



Transcriptome alignment

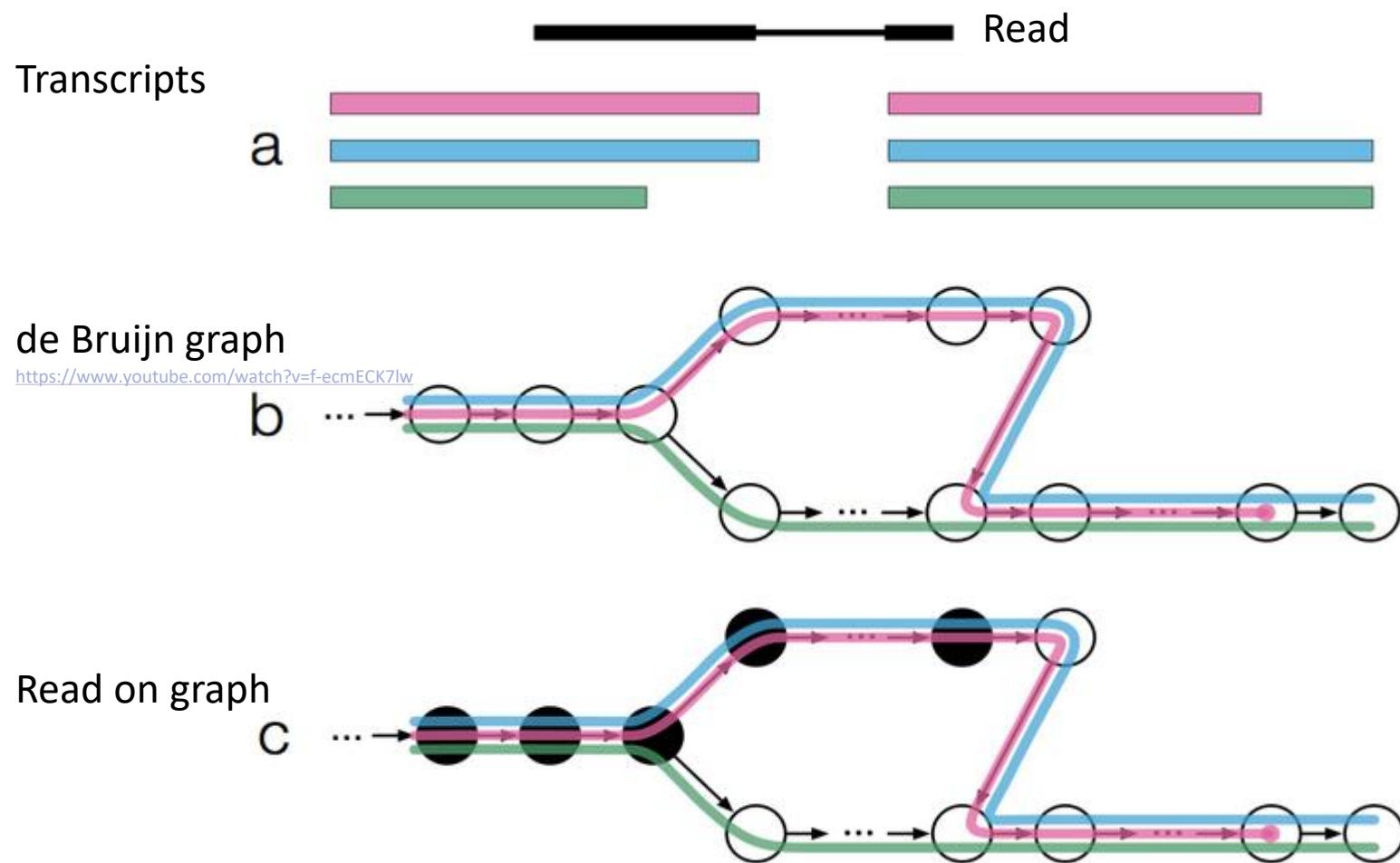


- reliable gene models required
- no detection of novel genes

Turro, EMBO, 2012

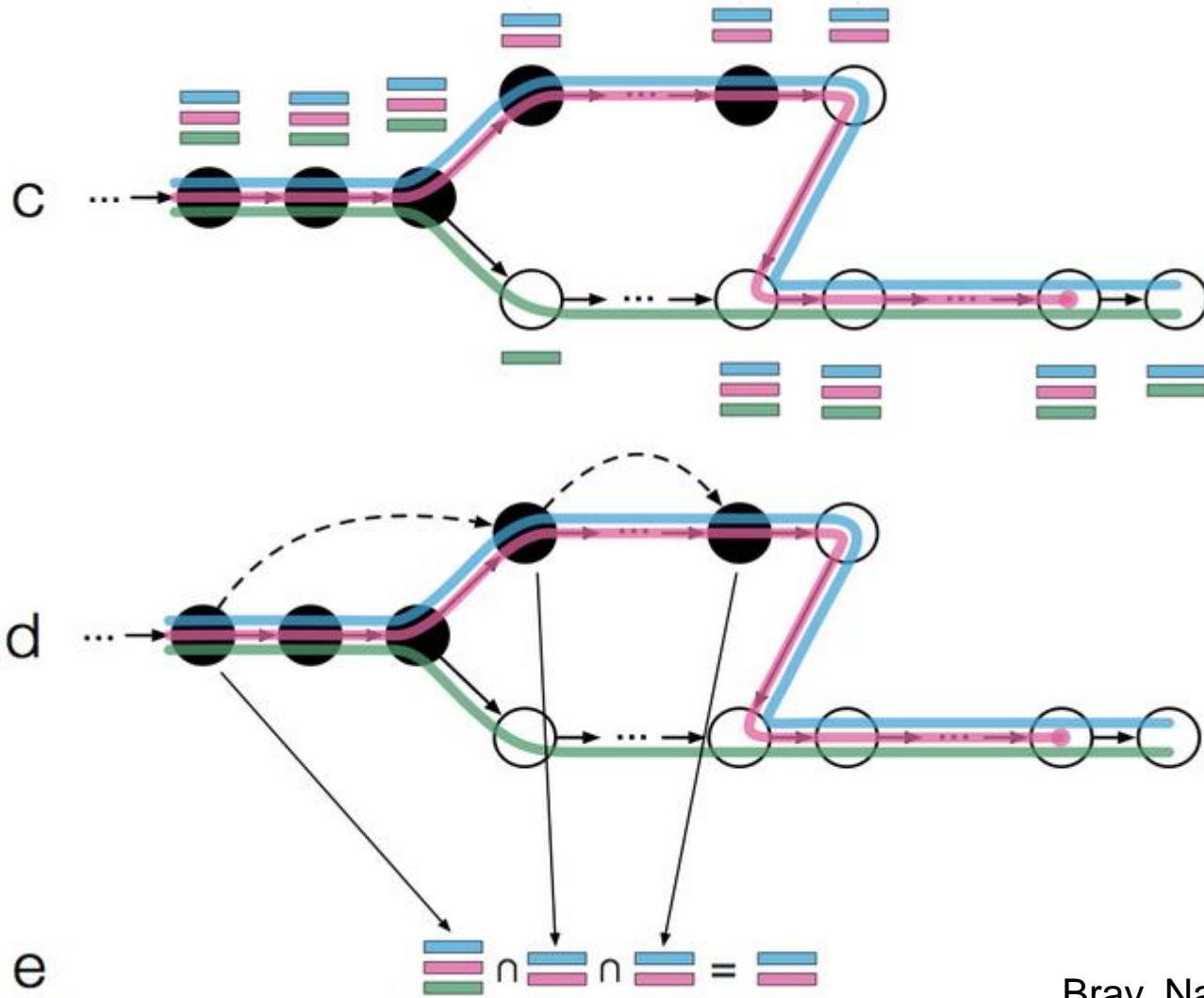
For clinical usage combination of different algorithms possible:

Genomic alignment - pseudoalignment



Bray, Nat Biotech, 2016

Genomic alignment - pseudoalignment



Bray, Nat Biotech, 2016

- 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:
 - TPM/ RPKM / FPKM (Cufflinks /Cuffdiff) (Mortazavi, Nat Meth, 2008)
 - TMM (edgeR) (Robinson & Oshlack, Genome Biol, 2010)
 - DESeq2 (DESeq2) (Love et al., Genome Biol, 2014)

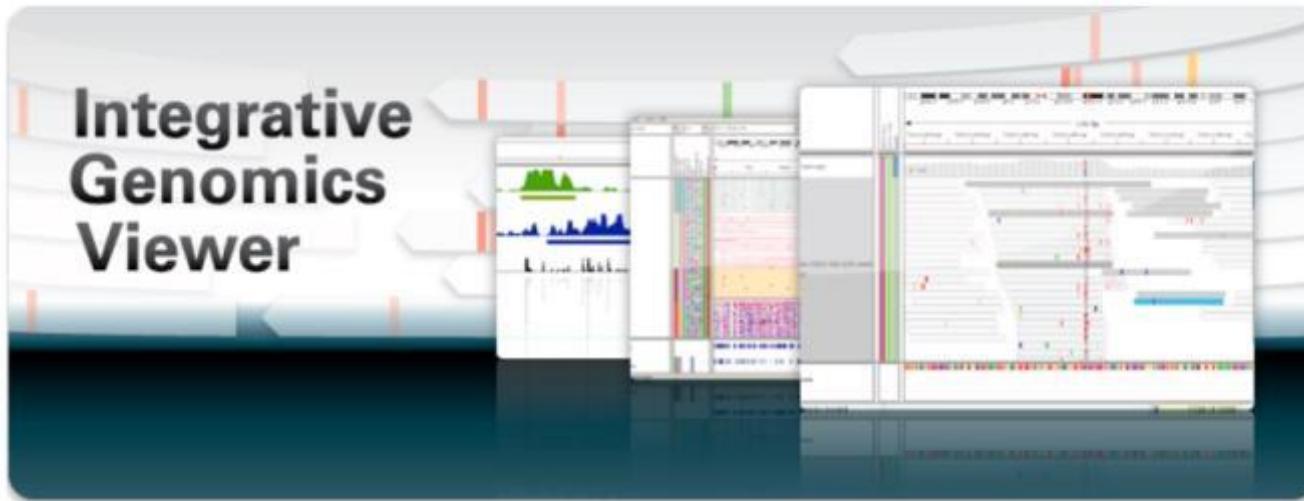


- Reads Per Kilobase per Million

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

- In RNA-seq, the relative expression of a transcript is proportional to the number of cDNA fragments that originate from it, however:
 - 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



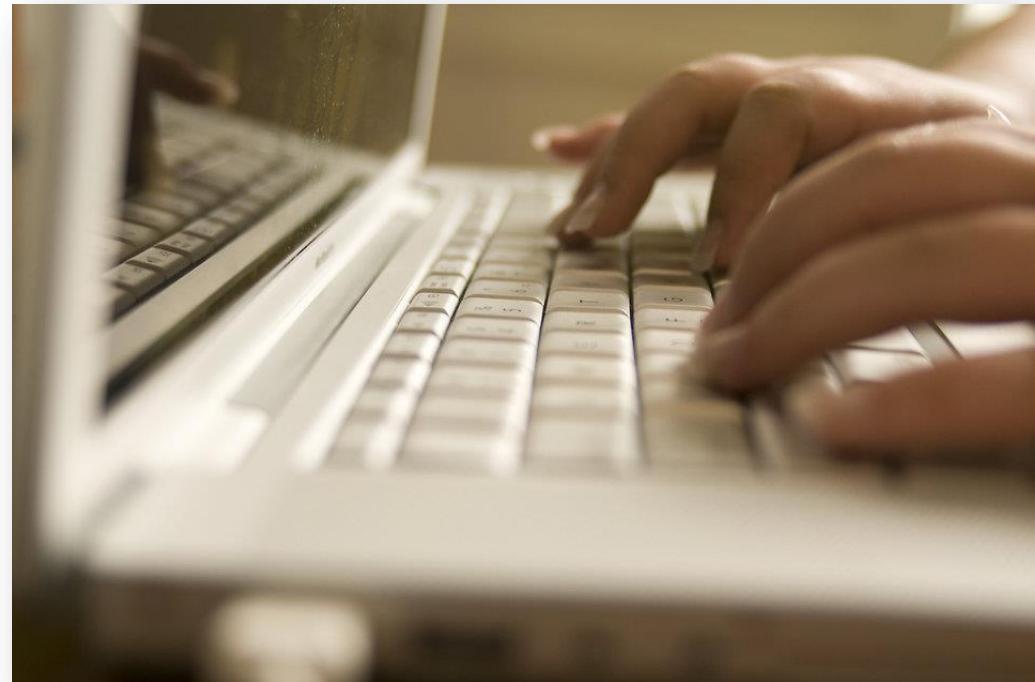
<https://ics.hutton.ac.uk/tablet/>

Hands on part 2

11:15 – 12:45

“Application of different read mapping approaches for genomic alignment”

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





12:45 – 14:30

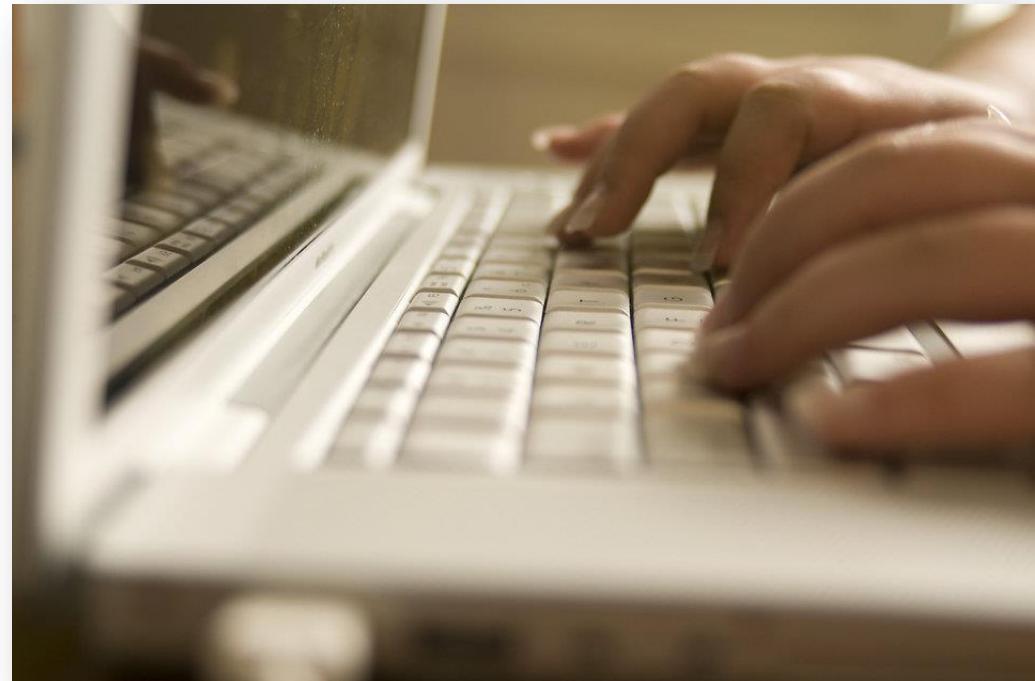


Hands on part 3

14:30 – 16:00

Introduction

“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

Most common cancer in women
worldwide

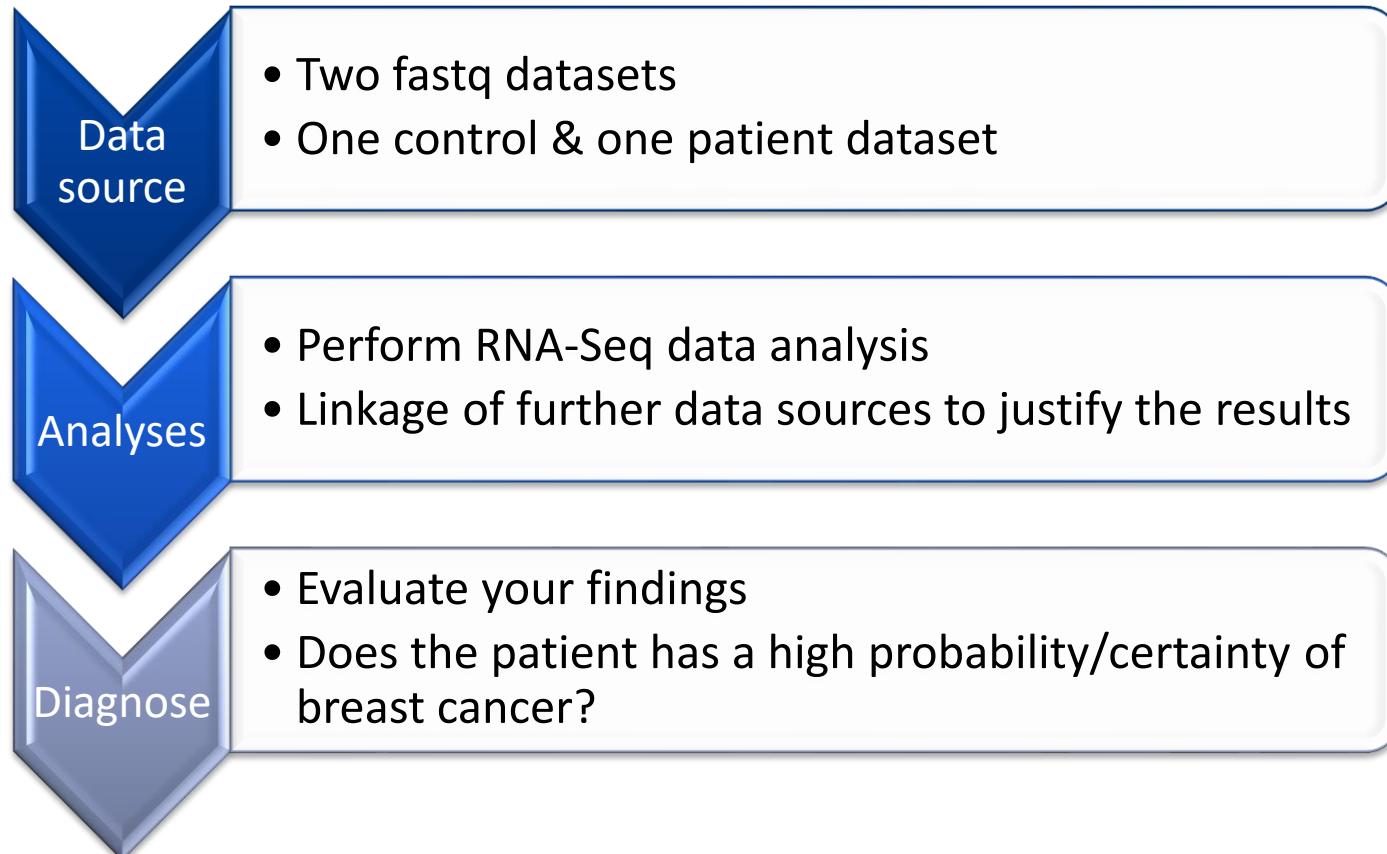
Leading cause of death from cancer in
women worldwide

Predictive factors that identify a benefit

Many different variations and subtypes

*Many different therapeutically
approaches*

Our use case for today – breast cancer screening

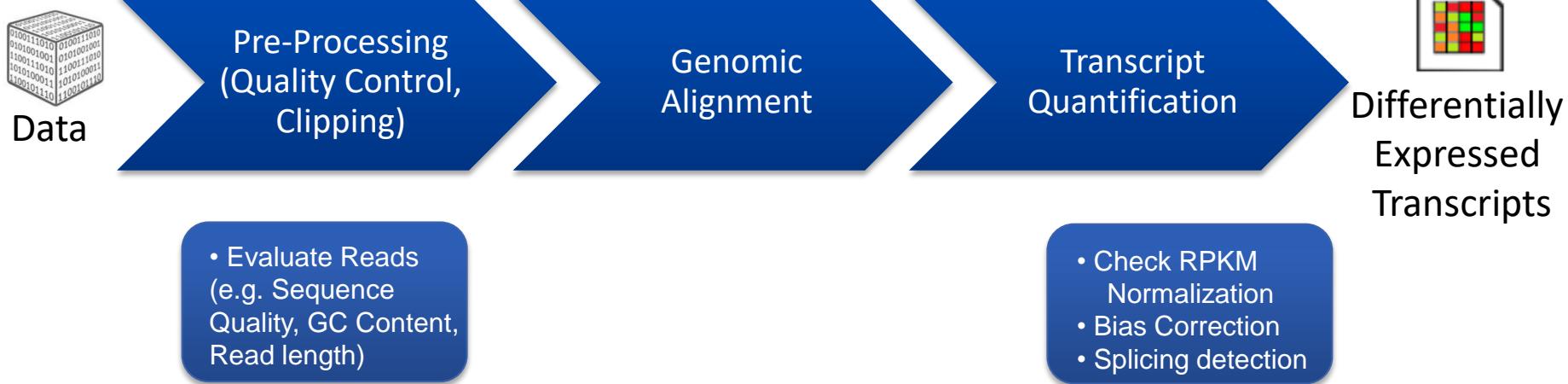




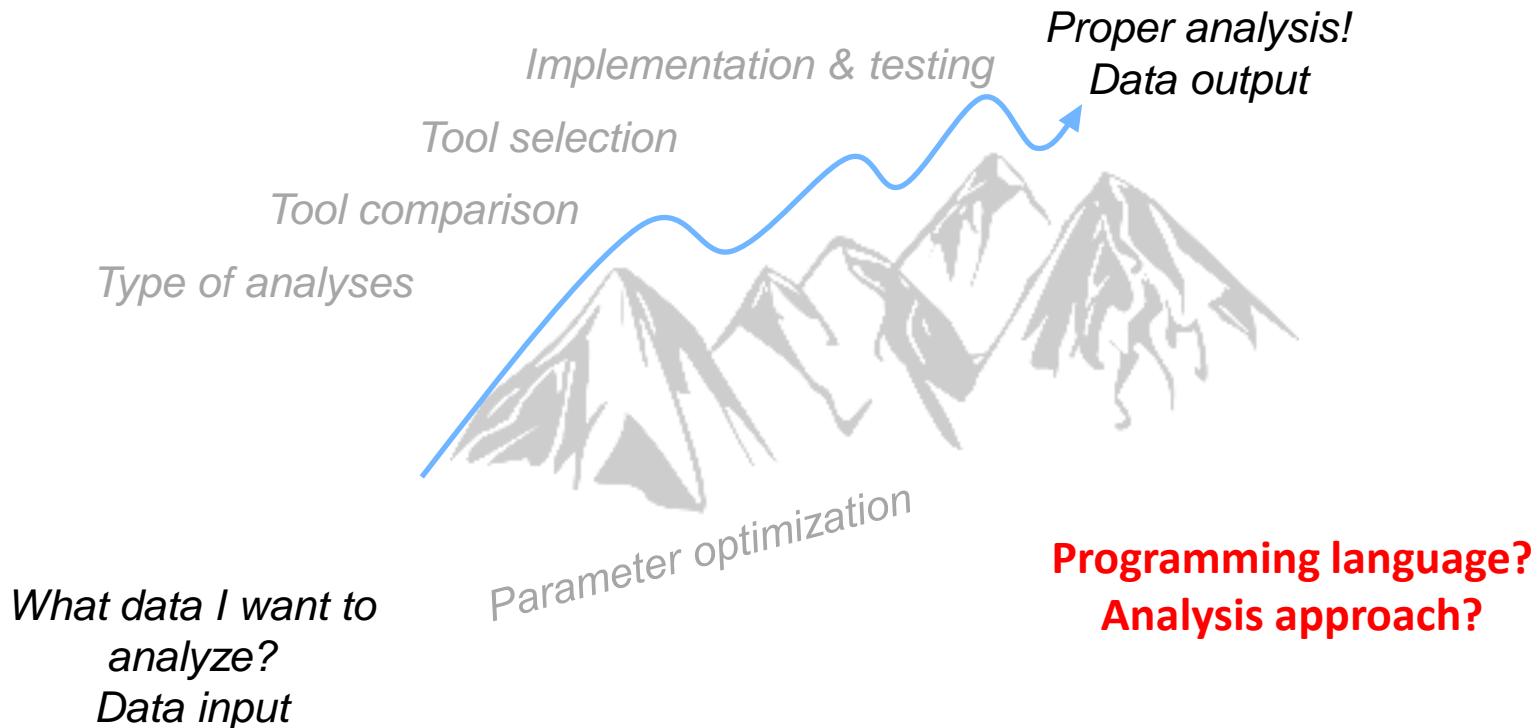
Get your data set!



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



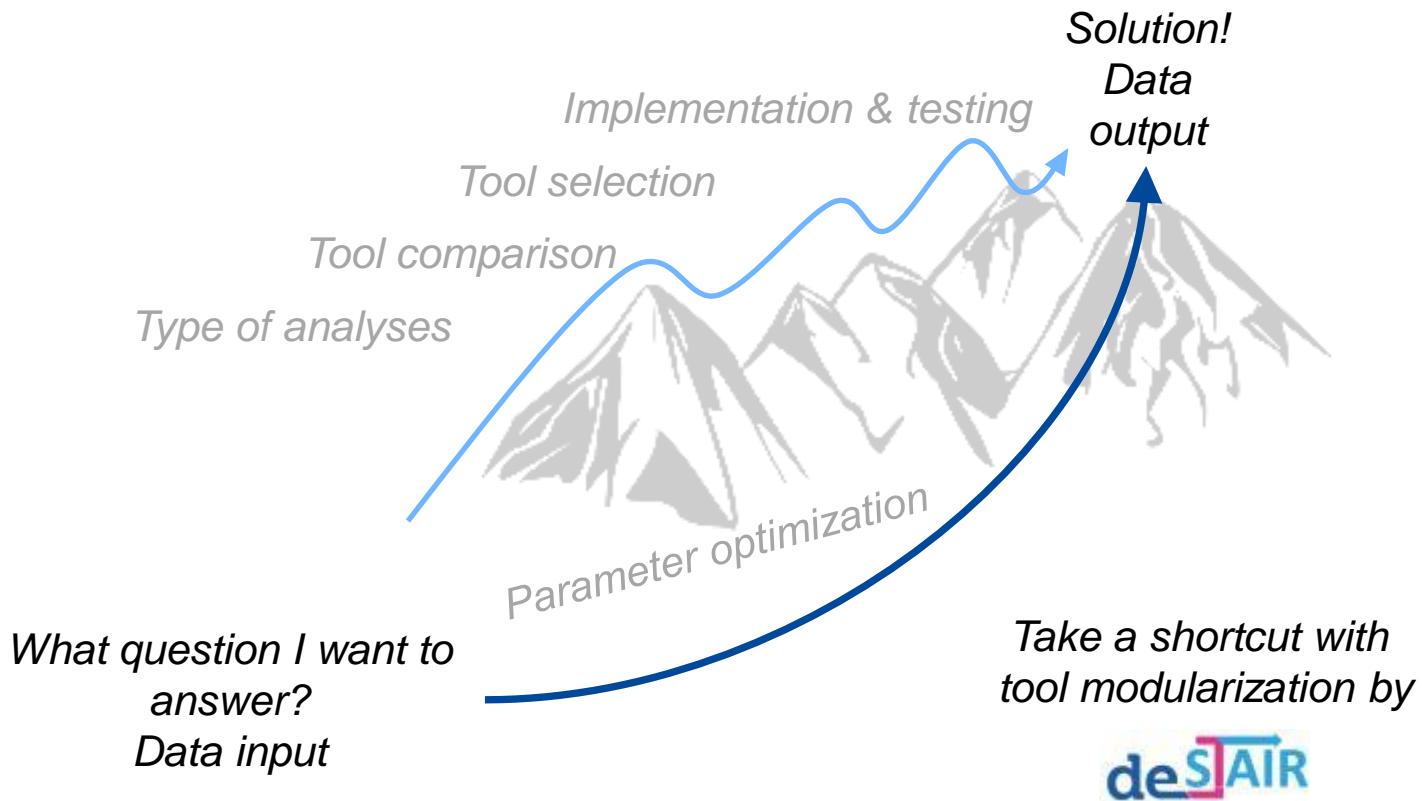
The struggle for the right approaches



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

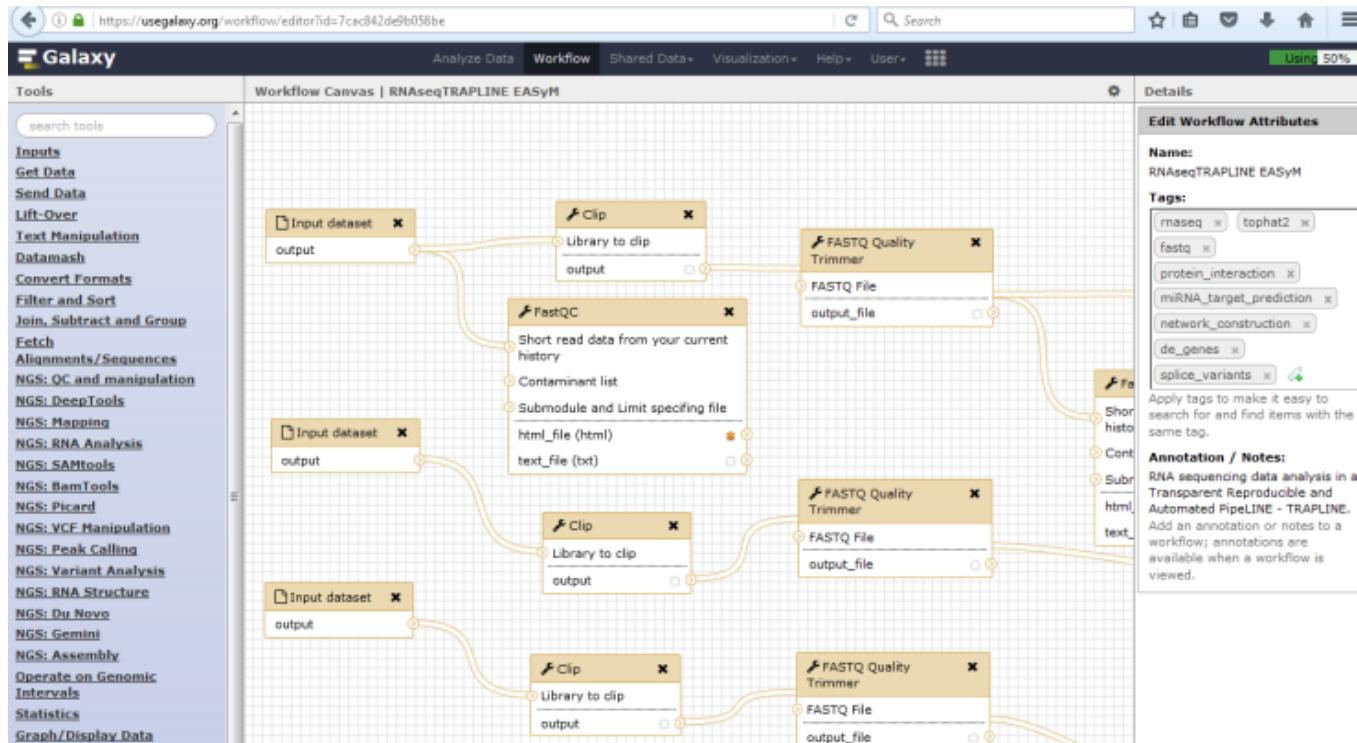


Why using workflows?



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

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



Explore your data!



gene_exp - Microsoft Excel

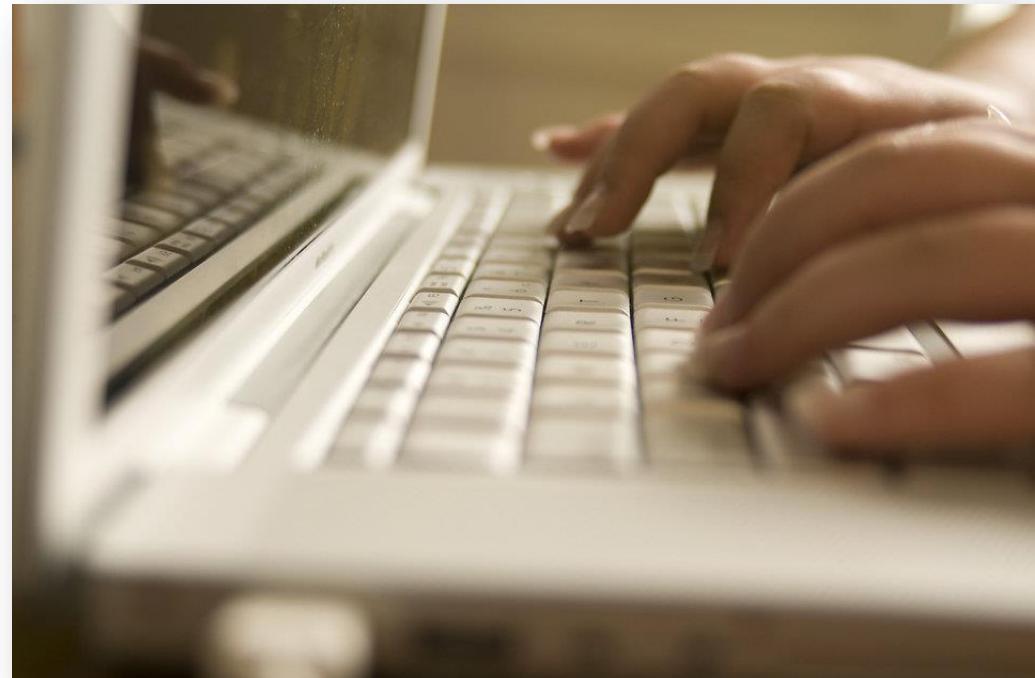
A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
test_id	gene_id	gene	locus	sample_1	sample_2	status	value_1	value_2	log2(fold_ch)	test_stat	p_value	q_value	significant						
1	test_id	gene_id	gene	locus	sample_1	sample_2	status	value_1	value_2	log2(fold_ch)	test_stat	p_value	q_value	significant					
2	ADAMTS4	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	chr5:9522085	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					

Hands on part 4

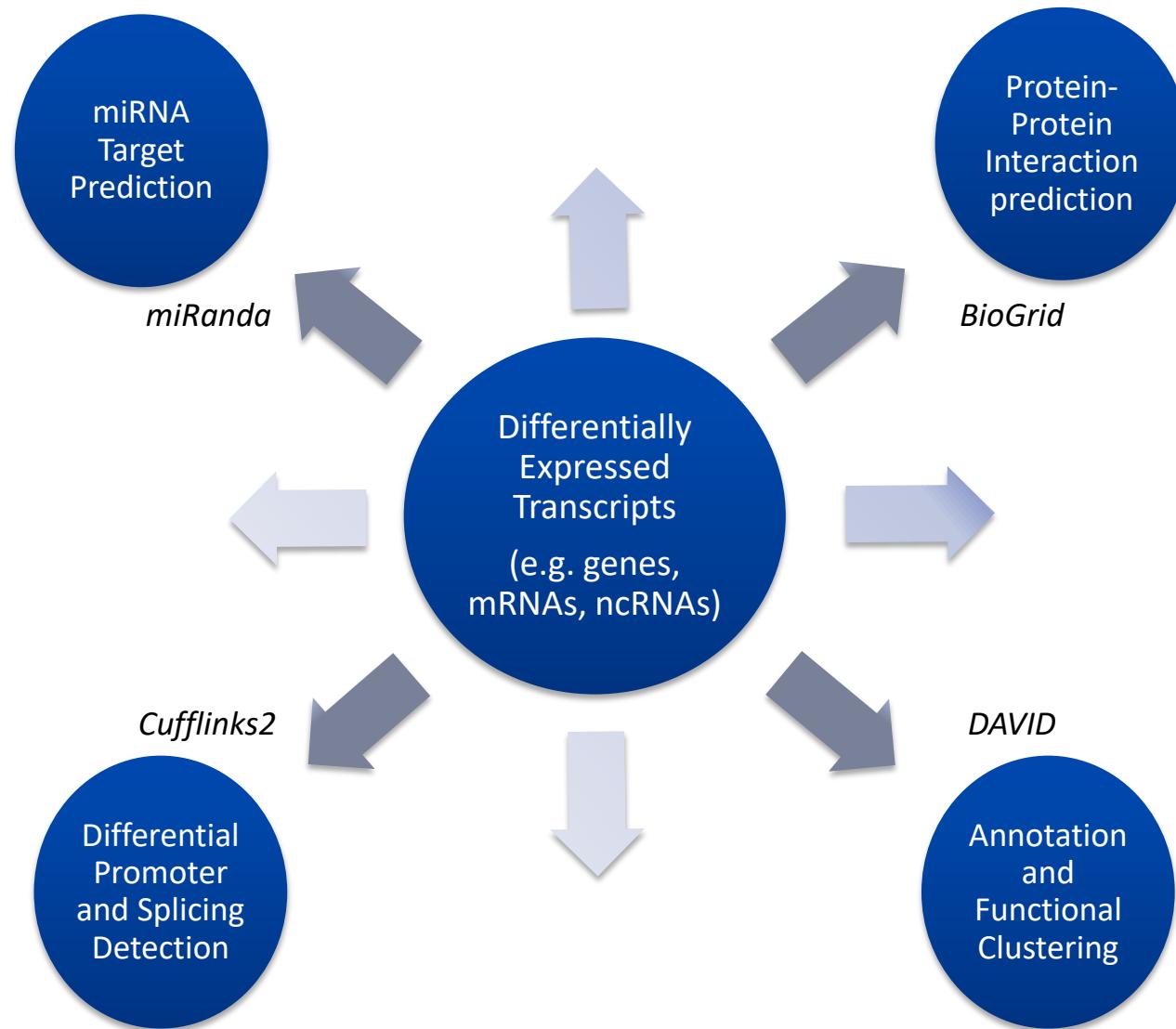
16:15 – 17:30

“Visualizations of RNA-Seq experiments with Galaxy”

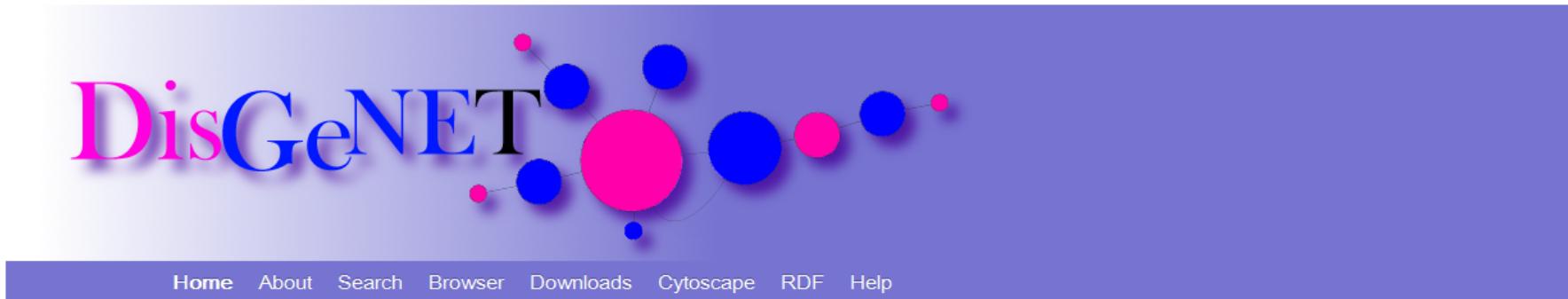
Material: <http://galaxyproject.github.io/training-material/topics/transcriptomics/>



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 computer. 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.

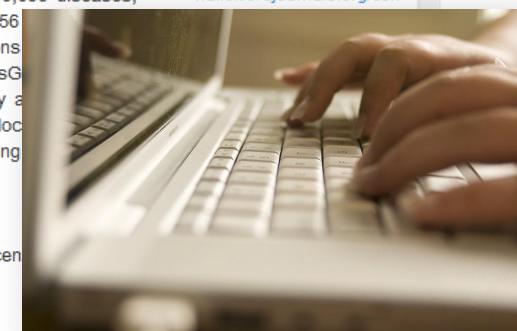
DisGeNET database has been cited by several papers. Some of them can be reviewed [here](#).

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



- Enrichr (<http://amp.pharm.mssm.edu/Enrichr/>)



Analyze What's New? Libraries Find a Gene About Help

Login | Register

12,272,405 lists analyzed

245,575 terms

132 libraries

Input data

Choose an input file to upload. Either in BED format or a list of genes. For a quantitative set, add a comma and the level of membership of that gene. The membership level is a number between 0.0 and 1.0 to represent a weight for each gene, where the weight of 0.0 will completely discard the gene from the enrichment analysis and the weight of 1.0 is the maximum.

Try an example [BED file](#).

Keine ausgewählt

Or paste in a list of gene symbols optionally followed by a comma and levels of membership. Try two examples:
[crisp set example](#), [fuzzy set example](#)

0 gene(s) entered

Enter a brief description for the list in case you want to share it. (Optional)

Submit

Contribute

Please acknowledge Enrichr in your publications by citing the following references:

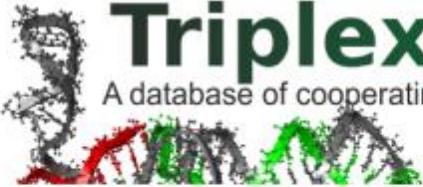
Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, Clark NR, Ma'ayan A. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics*. 2013;128(14).



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/>)



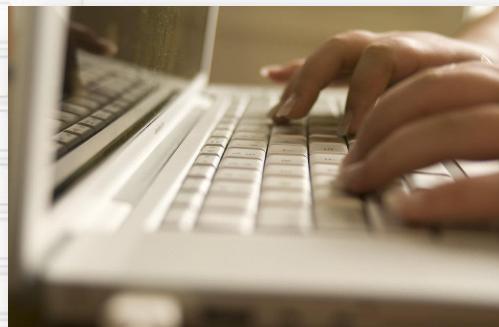
TriplexRNA

A database of cooperating microRNAs and their mutual targets

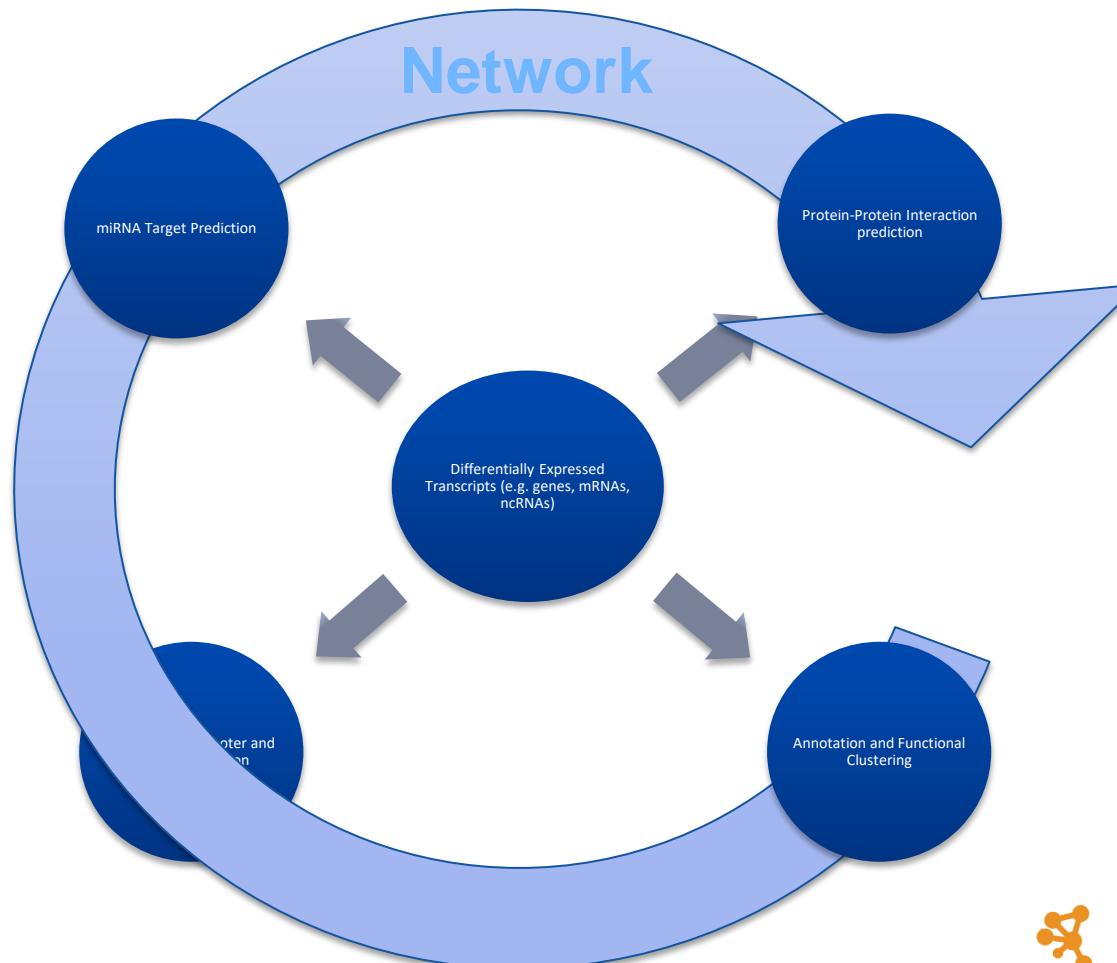
Search targets of synergistic microRNA regulation

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	more >
ATG4B	NM_178326	hsa-miR-140-5p	hsa-miR-346	28	-47.36	-15.58	more >
ZNF705A	NM_001004328	hsa-miR-140-5p	hsa-miR-296-3p	17	-43.76	-14.28	more >
FGR	NM_005248	hsa-miR-140-5p	hsa-miR-326	33	-43.56	-11.58	more >
PTCD1	NM_015545	hsa-miR-140-5p	hsa-miR-339-5p	34	-43.26	-12.98	more >
AARS	NM_001605	hsa-miR-24	hsa-miR-140-5p	32	-43.16	-17.18	more >
WEE1	NM_003390	hsa-miR-15b	hsa-miR-140-5p	16	-42.86	-16.28	more >
WNT1	NM_005430	hsa-miR-31	hsa-miR-140-5p	28	-42.56	-12.78	more >
ZBTB9	NM_152735	hsa-miR-140-5p	hsa-miR-296-3p	29	-41.96	-11.68	more >
ADRA1A	AY491776	hsa-miR-140-5p	hsa-miR-150	21	-41.96	-12.18	more >



There is nothing more practical than a network

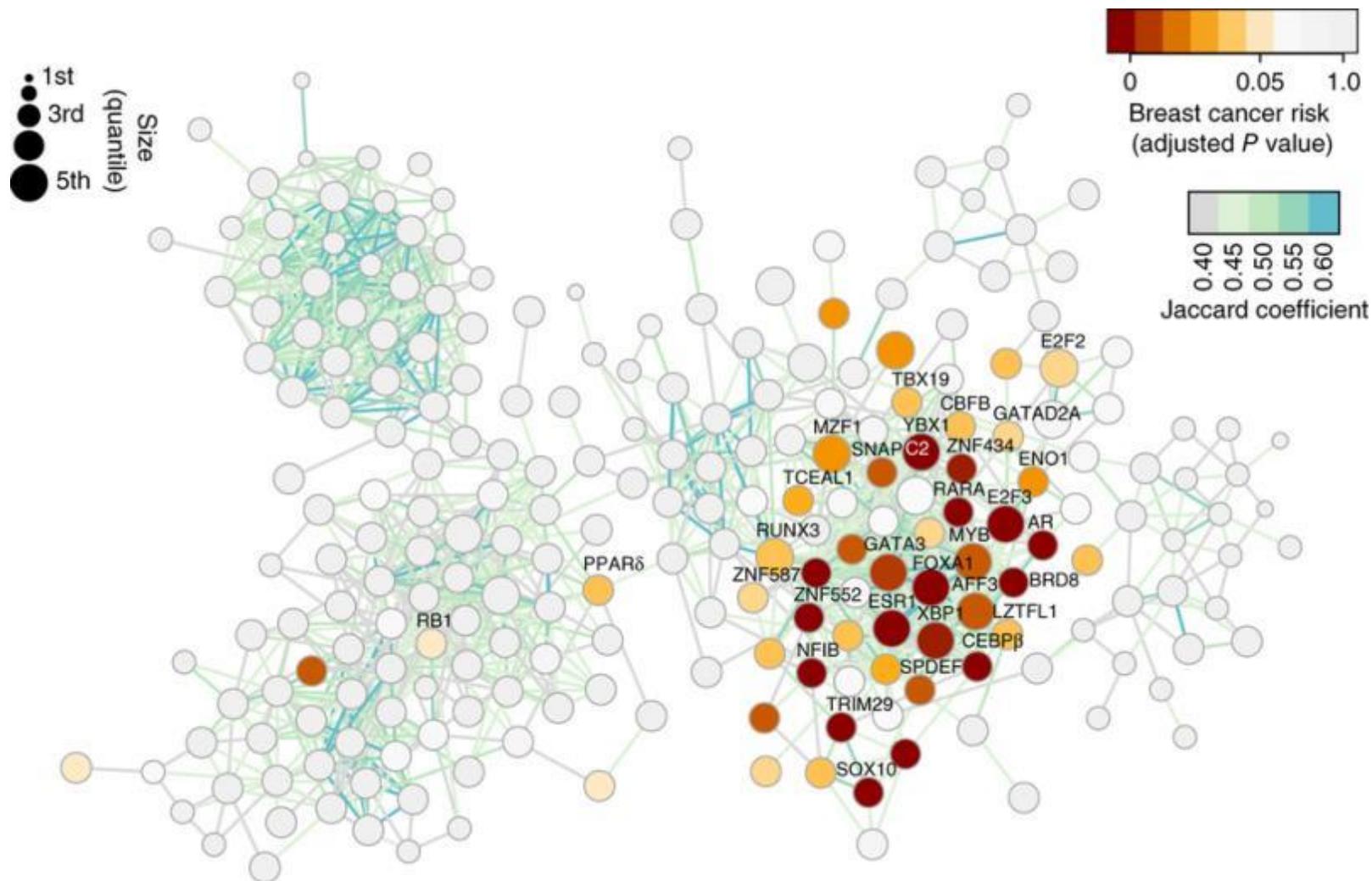


Cytoscape



Vanted / CellDesigner

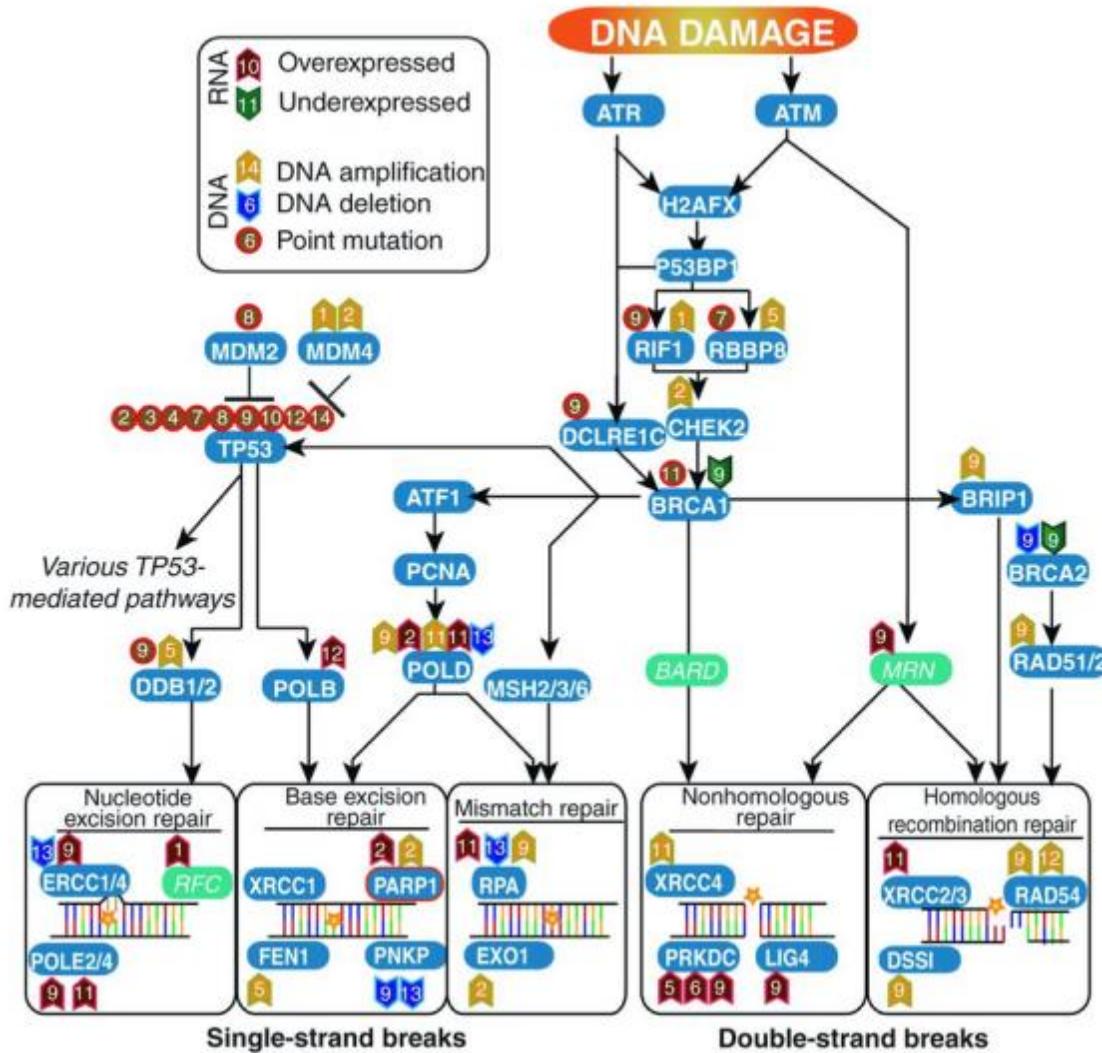
Network comparisons also reveal differences



Castro, Nat. Gen, 2016

Cytoscape

What else is done in this field with NGS?



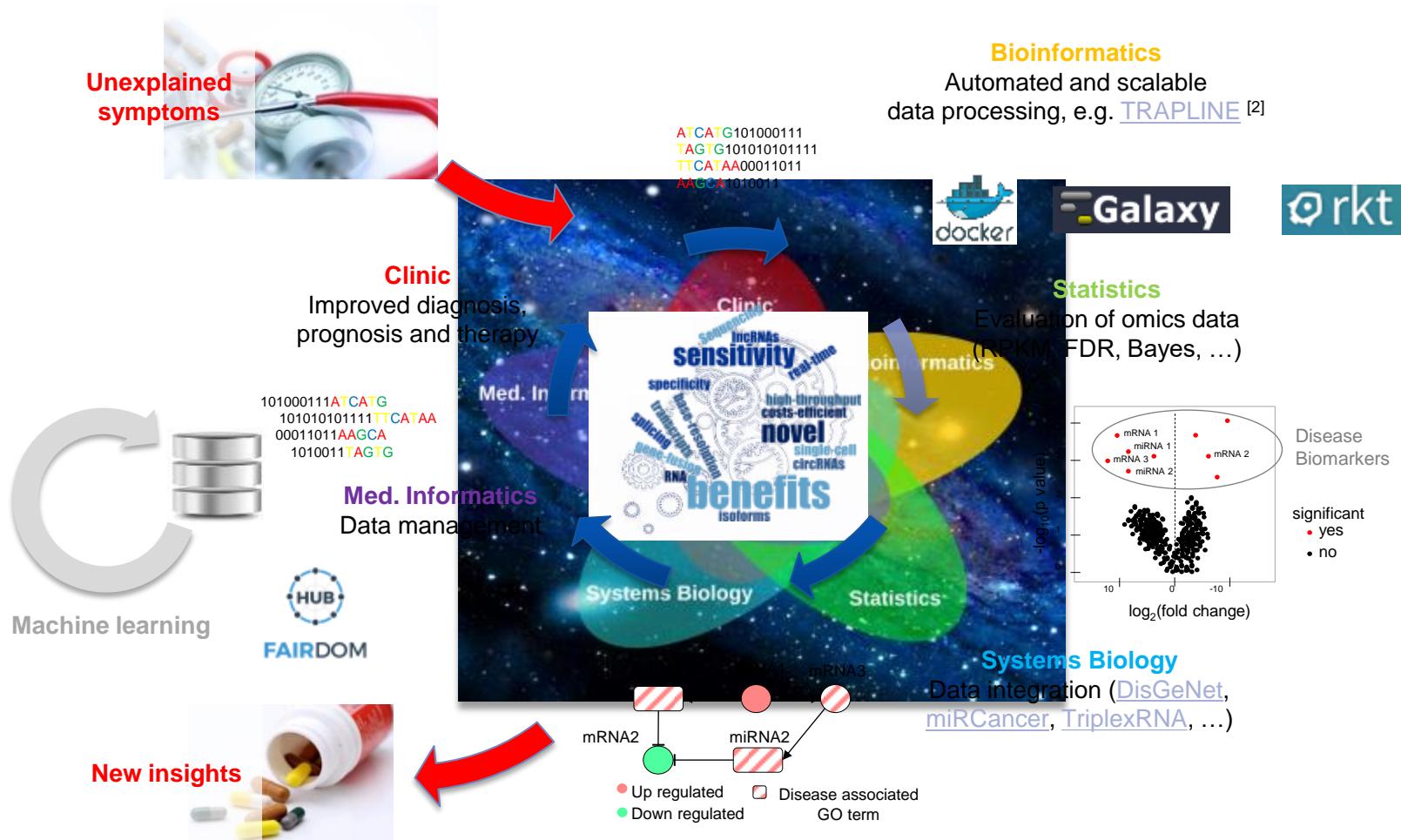
- Craig, D. W. et al. Genome and transcriptome sequencing in prospective metastatic triple-negative breast cancer uncovers therapeutic vulnerabilities. Mol. Cancer Ther. 12, 104–116 (2013). One of the first papers investigating integration of whole-transcriptome sequencing and genome sequencing for targeted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

Does the patient has a high risk of cancer

Patient



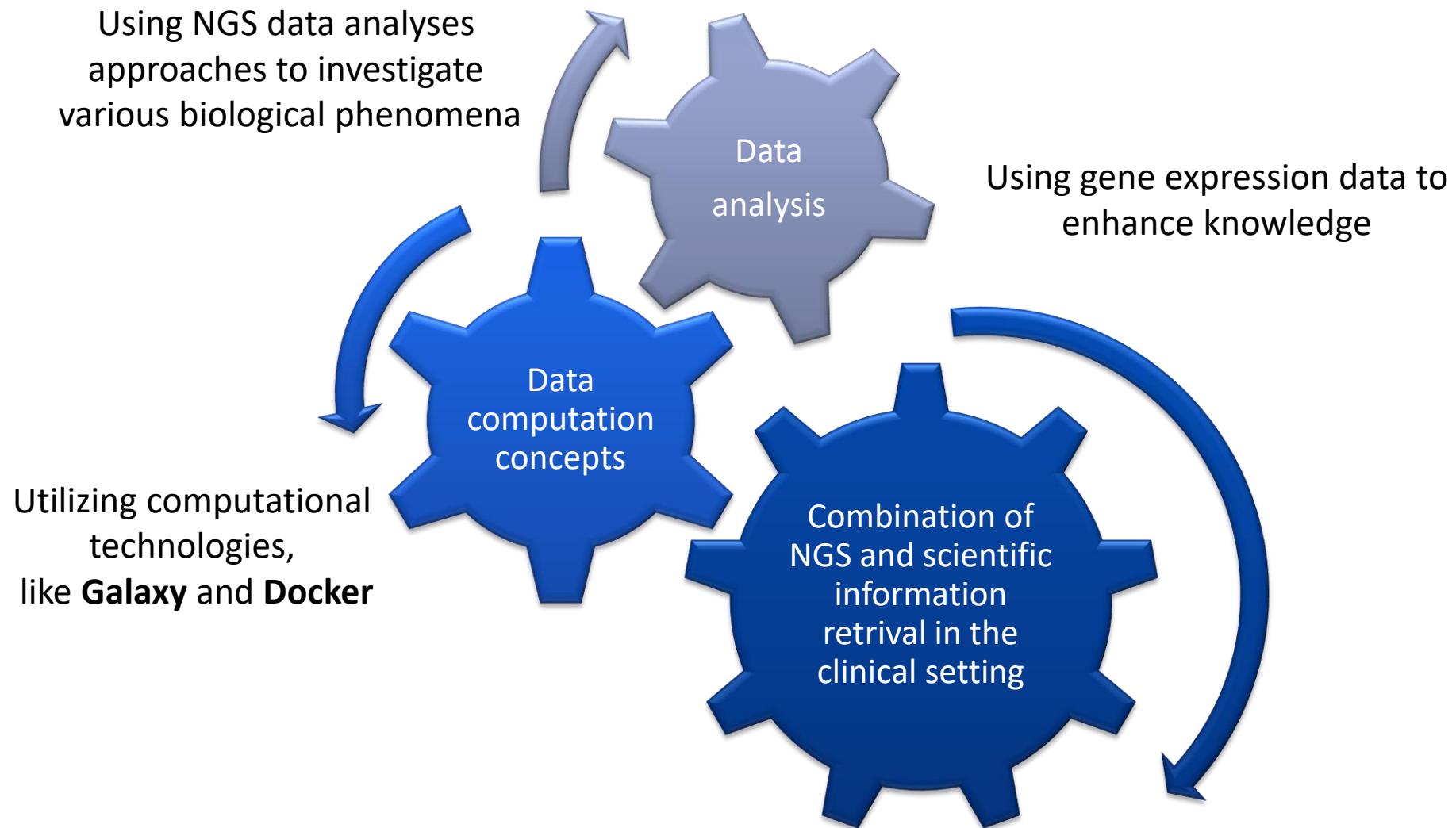
Our implementation strategy





- “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?



www.denbi.de/index.php/training-courses

Training Courses 2018

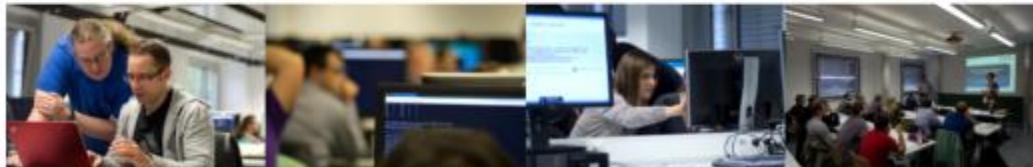
Training Courses 2018

Training Archive sorted
according to date

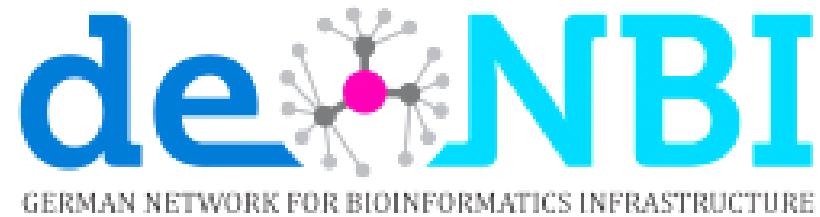
Training Archive sorted
according to de.NBI units

Online training & Media library

de.NBI Youtube channel



2018	Topic	Location
03-07 Sep	de.NBI Summer School 2018 - Riding the Data Life Cycle	Braunschweig
04 Sep	RNA-Seq data analysis with Galaxy for clinical applications - GMDS 2018	Osnabrück
08 Sep	Computational Mass Spectrometry with OpenMS - From Algorithms to Integrated Workflows - ECCB 2018	Athen
09 Sep	From an application to a fully integrated workflow – Comprehensive software engineering with SeqAn - ECCB 2018	Athen
17-21 Sep	6th Galaxy HTS data analysis workshop	Freiburg
19-21 Sep	de.NBI - CeBiTec Nanopore Best Practice Workshop 2018	Bielefeld
19-21 Sep	SeqAn User Meeting 2018	Berlin
19-21 Sep	OpenMS User Meeting 2018	Berlin
19-21 Sep	de.NBI - KNIME User Meeting 2018	Berlin
19-21 Sep	DAIS User Meeting 2018	Berlin
19-21 Sep	MASH User Meeting 2018	Berlin
24-28 Sep	Deep Learning Bootcamp 2018	Dresden



Evaluation at :

<https://de.surveymonkey.com/r/denbi-course?sc=rbc&id=000140>

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)



Deutsche Gesellschaft für
Medizinische Informatik,
Biometrie und
Epidemiologie e.V.



Supported by:



denbi.de



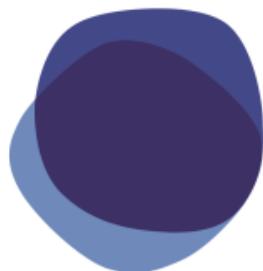
elixir-europe.org



Bundesministerium
für Bildung
und Forschung

bmbf.de

We hope you enjoyed the training!



SYSTEMS BIOLOGY
BIOINFORMATICS
ROSTOCK

deNBI
GERMAN NETWORK FOR BIOINFORMATICS INFRASTRUCTURE

gmds | Deutsche Gesellschaft für
Medizinische Informatik,
Biometrie und
Epidemiologie e.V.

Universität
Rostock



Traditio et Innovatio

