# *Customized workflow development and data integration concepts in Systems Medicine*

Markus Wolfien
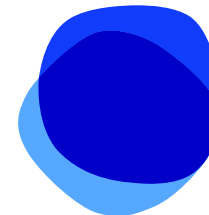
*GMDS Workshop "Methods in Sysems Medicine"*
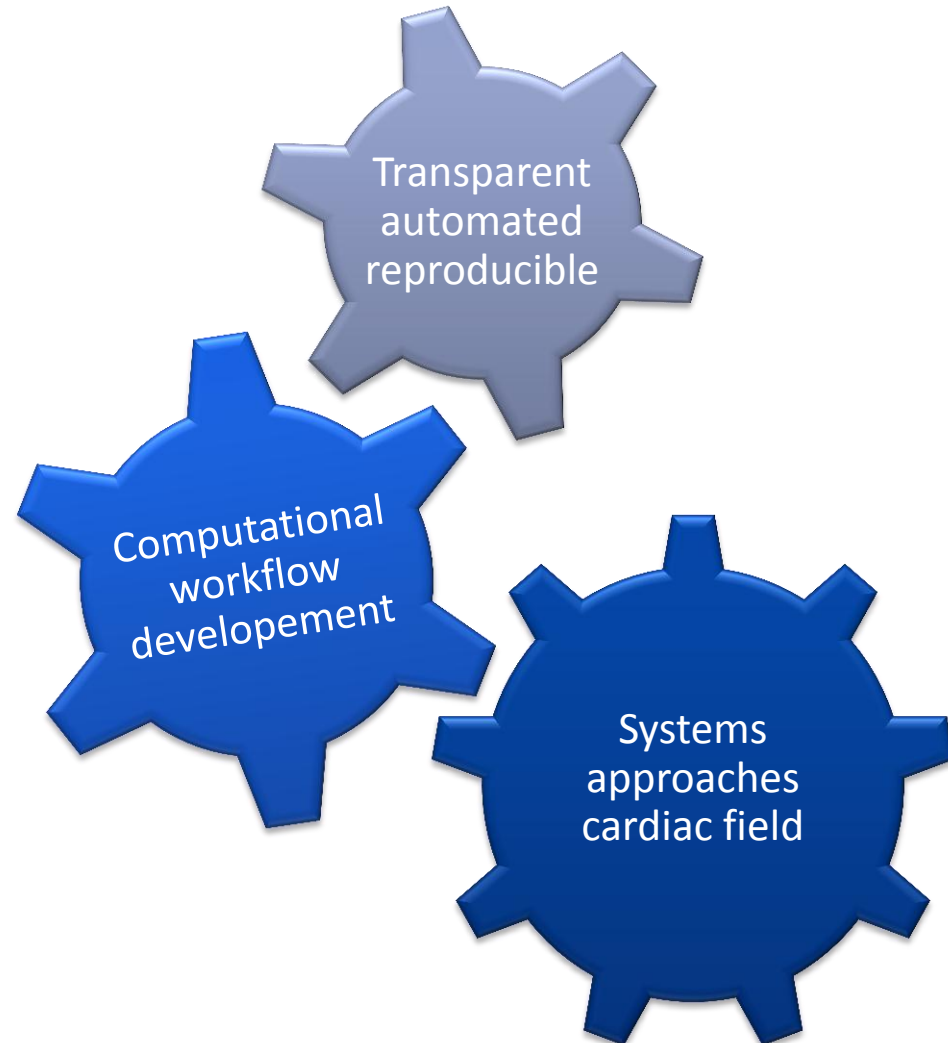www.sbi.uni-rostock.de

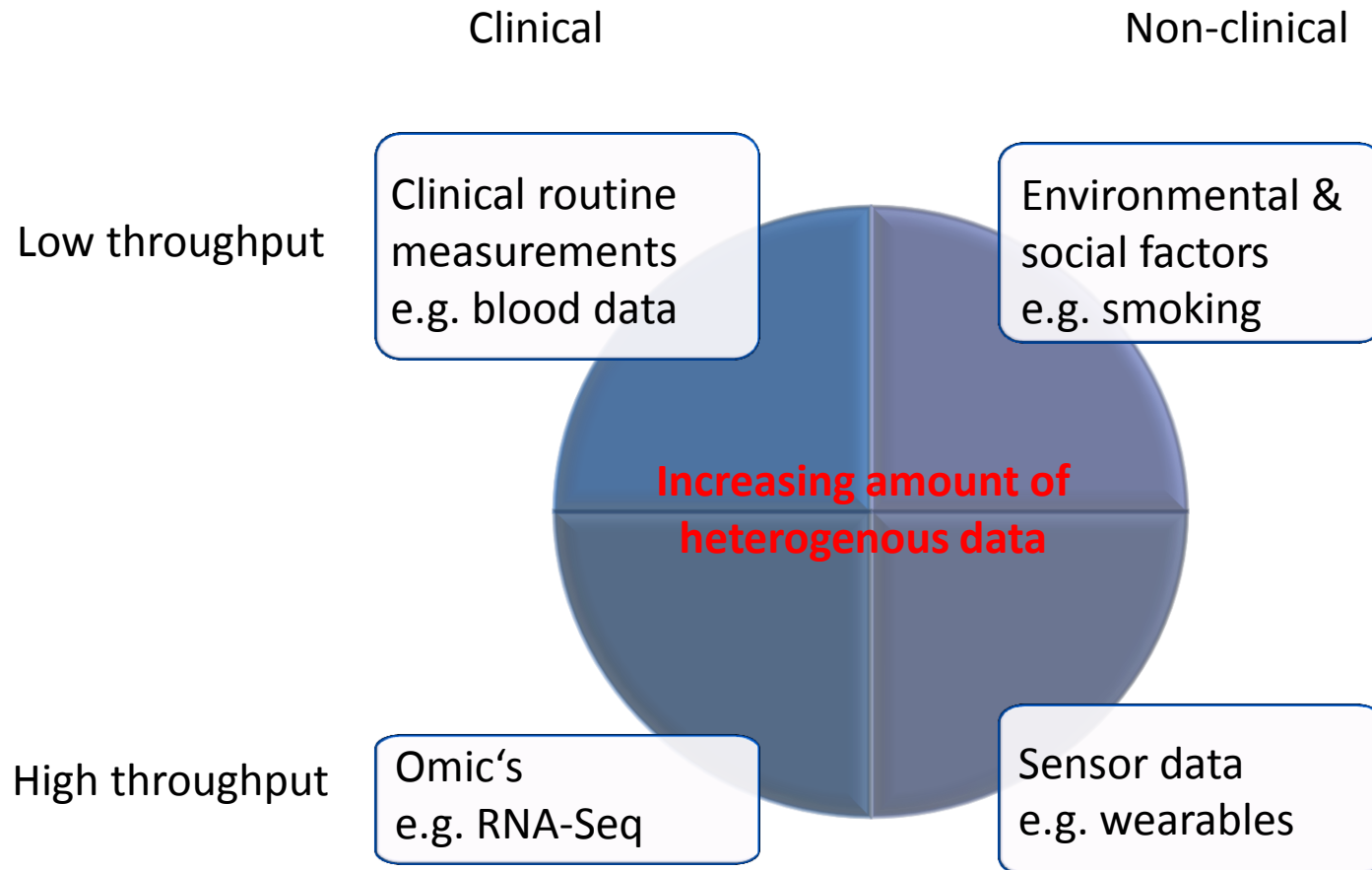# Objectives for my talk

- Why using workflows?

- How are workflows being developed?

- What can be done to analyze and integrate diverse data?

Transparent automated reproducible

Computational workflow developement

Systems approaches cardiac field

# Data being generated is steadily increasing

Clinical             Non-clinical

Low throughput

Clinical routine measurements e.g. blood data

Environmental & social factors e.g. smoking

**Increasing amount of heterogenous data**

High throughput

Omic's e.g. RNA-Seq

Sensor data e.g. wearables

# The struggle for the right approaches

*Implementation & testing*

*Proper analysis!*
*Data output*

*Tool selection*

*Tool comparison*

*Type of analyses*

*Parameter optimization*

*What data I want to analyze?*
*Data input*

**Programming language?**
**Analysis approach?**

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

# Workflows!
Provide an infrastructure to set up, execute and monitor tool environments

# Medical "Big Data" and the need for new analyses



GenePattern
broadinstitute.org

GeneProf
geneprof.org

Grape
big.crg.cat/services/grape

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

python
python.org

R
mapman.gabipd.org

Galaxy
usegalaxy.org

KNIME
Open for Innovation
knime.org

R
r-project.org/

GeneTalk
gene-talk.de

Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS
bioconductor.org

geneXplain
genexplain.com

BaseSpace
illumina.com

# The Galaxy interface and tours: usegalaxy.org

# Basic workflow for data processing

Example: RNA-Seq data analysis



**Data**

**Pre-Processing (Quality Control, Clipping)**

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

**Genomic Alignment**

- Multiple Correction
- Gene fusion
- SNP Calling ready (GATK toolkit)

**Transcript Quantification**

- Check RPKM Normalization
- Bias Correction
- Splicing detection

**Differentially Expressed Transcripts**

# Differential expression is the key of evaluation



**SYSTEMS BIOLOGY BIOINFORMATICS ROSTOCK**

miRNA Target Prediction

*miRanda*

Protein-Protein Interaction prediction

*BioGrid*

Differentially Expressed Transcripts (e.g. genes, mRNAs, ncRNAs)

*Cufflinks2*

Differential Promoter and Splicing Detection

*DAVID*

Annotation and Functional Clustering

# Transparent Reproducible Automated PipeLINE - TRAPLINE

**Galaxy Modules implemented:**

*Data*

*Pre-Processing*

*Alignment*

*Differential Expression (DE)*

*Transcriptomic Analysis*

*Networks & Function*

**Legend:**



## Data Processing

**FASTQ Files**
Upload RNA sequencing data (Illumina, SOLiD, Solexa) & format by FASTQ Groomer

*FASTXclipper*

**Clip Seq. Adapter**
Removal of preliminary adapter sequences
Minimum remaining sequence length ≥ 15 bases

*FASTQ Quality Trimmer*

**Quality Trimming**
Discarding low quality reads
Q ≤ 20 starting at 5' and 3' end; window/step = 1

*TopHat2*

**Genome Mapping**
Single or paired end reads
Default alignment options and built-in genomes

*Picard Toolkit*

**Correcting Reads**
Investigation of duplicated reads for read correction and SNPs
Maximum offset = 100

*Cufflinks2, Cuffdiff2*

**Gene Expression**
Calculating RPKM values and identifying sign. diff. genes with
FDR ≤ 0.05, p ≤ 0.05, FC ≥ 2

## Data Evaluation and Annotation

**Sign. Diff. Expr. Genes**
Lists with up and down-regulated genes

**Promoters & Spliced Genes**
Lists with potential splicing sites & multi-promoter genes

*DAVID*    *miRanda*    *BioGRID*

**Functional Annotation**
Gene annotation with DAVID link (e.g. Go terms, Interpro)
ID = Official Gene Symbol
List Type = Gene List

**miRNA Target Prediction**
Predict binding sites to identify target genes and evaluate their impact on downreg. genes
Conserverd and non-conserved miRNAs + high mirSVR scores (Release 08/10)

**Find Protein Interactions**
Using Protein-Protein interactions from BioGRID
BioGRID Organism Version 3.3.122

*Functional classification*    *Compare Datasets*    *Compare Datasets*

**Enriched Genes**
Genes annotated by GO terms and clustered in functionally enriched groups
Default parameter

**miRNA Target Genes**
List of upregulated miRNAs and their predicted downregulated mRNAs

**Prot.–Prot. Interactions**
List of interacting proteins, between up and down-regulated mRNAs

*Join Datasets*

**Network construction**
Using the information of all data evaluation modules to easily enable the user building networks in Cytoscape or further network analysis tools (e.g. ClueGO)

Wolfien, *BMC Bioinf*., 2016

# 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"
        },
```
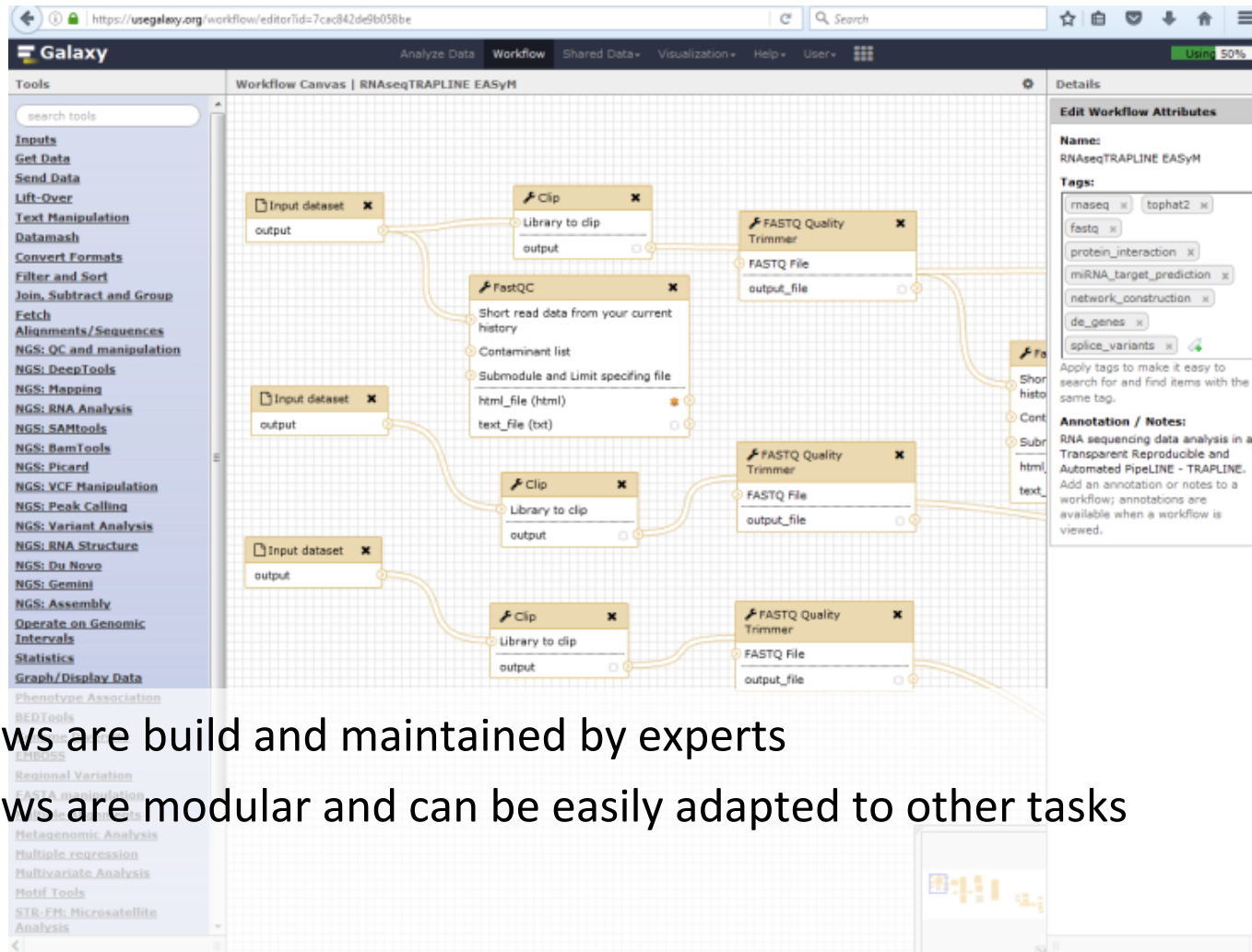
Specific xml file with tools, parameters and meta data!

# Using workflow development



- Workflows are build and maintained by experts
- Workflows are modular and can be easily adapted to other tasks

sbi.uni-rostock.de/destair

*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.

# de.STAIR – RNA-Seq analysis and integration

# Supporting new data analysis approaches

- Key performance of Galaxy
  - Accessibilty
  - Reproducibility
  - Transparency

jupyter.org/

usegalaxy.org

elixir-europe.org

denbi.de

# 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)
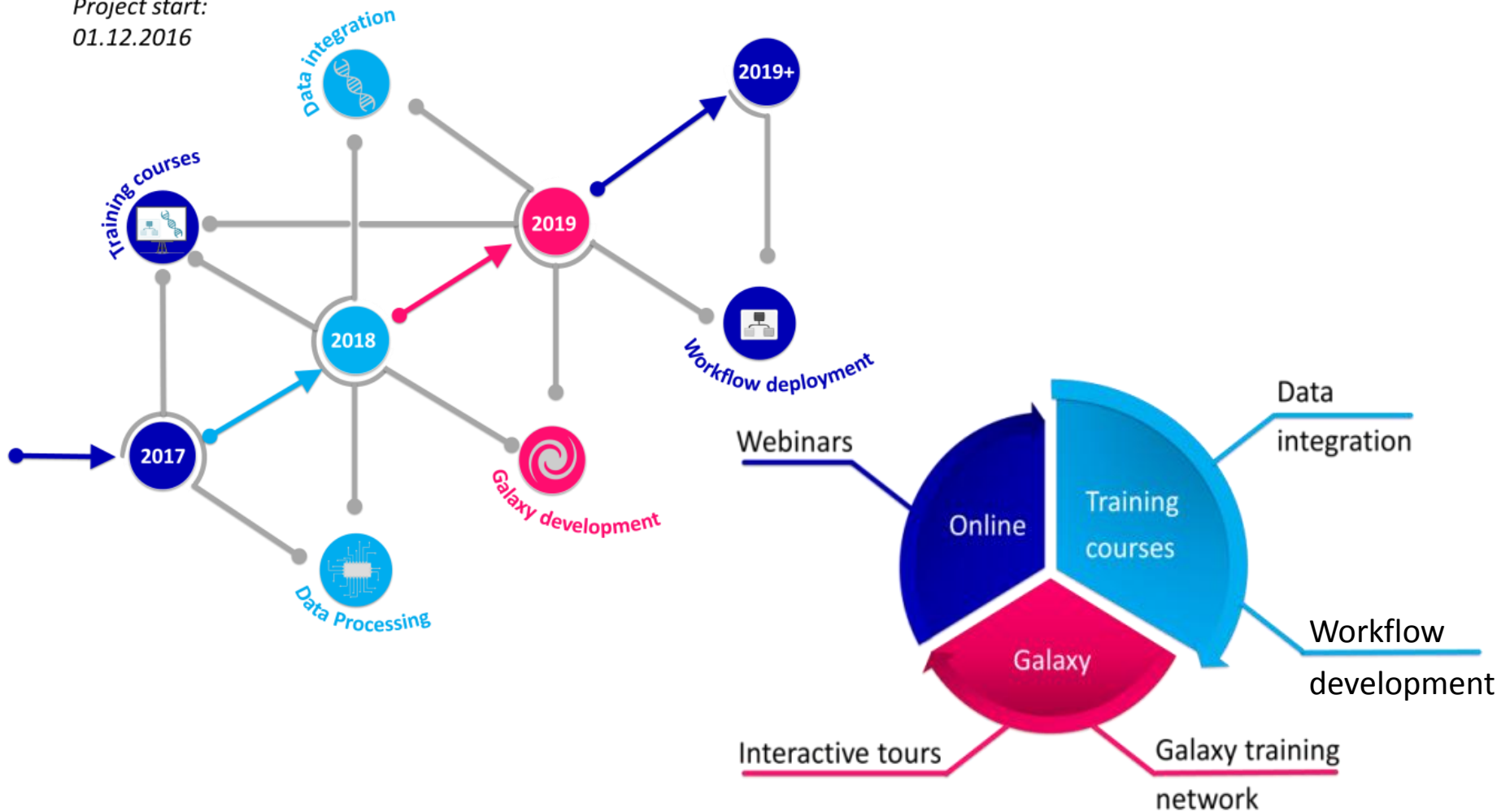- Application of workflows and tools is targeted for non-computational users

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

# Containerization!

- Build for scale
- Extensible and flexible
- Using optimized system libraries
- E.g. used by ebay, GE, illumina, Spotify

biocontainers.pro

# 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"]
```

# One tool to share them all

**Galaxy** + **docker** = **Symbiosis!**

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

- Stand-alone Docker container which "conserves" your tool compilation (for an easy use – one command line or single kitematic.com click!) – Slurm cluster

```
docker run –p 8080:80 mwolfien/trapline
```

- Docker swarm for single docker container & tools for higher flexibility

# Supporting the RNA Galaxy-workbench

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

github.com/bgruening/galaxy-rna-workbench



galaxyproject.github.io/training-material

Gruening *et al., NRA, 2017*

# Why using workflows?



Implementation & testing

Tool selection

Tool comparison

Type of analyses

Parameter optimization

Proper data analysis!
Data output

What data I want to analyze?
Data input

Providing guidence with Workflow development and tool modularization by e.g.

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

# Medical "Big Data" and the need for new …

KNIME
knime.org

Galaxy
usegalaxy.org

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

gene✕plain
genexplain.com

BaseSpace®
illumina.com

+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

- 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

# 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-gdffc67f
```
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

*Workflows can be used for complex data integration and interpretation*

# Our implementation strategy



**Unexplained symptoms**

**Clinic**
Improved diagnosis, prognosis and therapy

**Machine learning**

**Med. Informatics**
Data management

**New insights**

**Bioinformatics**
Automated and scalable data processing, e.g. TRAPLINE [2]

**Statistics**
Evaluation of omics data (RPKM, FDR, Bayes, …)

Disease Biomarkers

$-\log_{10}$(p value)

$\log_2$(fold change)

significant
- yes
- no

**Systems Biology**
Data integration (DisGeNet, miRCancer, TriplexRNA, …)

miRNA1   mRNA1   mRNA3
mRNA2   miRNA2

- Up regulated
- Down regulated
- Disease associated GO term

„Medizin 4.0 - Zur Zukunft der Medizin in der digitalisierten Welt". ISBN: 978-3-9809206-5-0
Poster prize. https://doi.org/10.6084/m9.figshare.4029069.v1

Clinical trial design



**Trial group
n = 77**

Blood analysis



**Trial group
n = 77**

BM aspiration



**Stem cell
purification for
therapy**

CABG & Therapy



**Treatment**
CD133+
Placebo

Laboratory measurements



**HRO efficacy subgroup
n = 31**

**CHRONIC ISCHEMIC HEART FAILURE**

**INDUCED CARDIAC REGENERATION**

*Day -2*          *Day -1*          *Day 0*          *Day 180*

Steinhoff, Nesteruk, Wolfien, et al., *EBioMedicine*, 2017

# Customized workflow strategy for clinical trial



**Unexplained phenotype**

Blood analysis, e.g. thrombocytes, leucocytes
Functional analyses, e.g. LVEF, 5-min-walk test
Environmental factors, e.g. age, sex, other diseases, medications
FACS measurements, e.g. CD133+ stem cells and further subtypes
Protein expressions data, e.g. ELISA for angiogenic factors
Gene expression data, e.g. PCR data, RNA-Seq
Further assay's, e.g. Matrigel Plug, CFU-Hill assay

**Bioinformatics**
Individual data analysis by customized workflows, e.g. TRAPLINE [2]

**Clinic**
Responder vs Non-responder classification

## Response signature

Thrombocytes
$CD^{+133/34+}$ EPC

$CD^{146+}$ CEC
VEGF
EPO
Vitronectin
NT-proBNP
SH2B3

**Statistics**
Evaluation of data
(p-value, RPKM, FDR, …)

Disease Biomarkers

**Med. Informatics**
Data management

**Machine learning**

**Systems Biology**
Data integration and enhancement
(Transfac, GO, …)

**New insights**

Steinhoff, Nesteruk, Wolfien, et al., EBioMedicine, 2017

# Summarizing my talk

- Workflows are a great ressource for data analyses

- Containers are used to simplify data analyses

- Integrative workflows support clinical investigations

Transparent automated reproducible

Computational workflow developement

Systems approaches cardiac field

# Acknowledgements

Olaf Wolkenhauer (University of Rostock)

Andrea Bagnacani (University of Rostock)

Wolfgang Hess (University of Freiburg)

Steve Hoffmann (University of Leipzig)

Rolf Backofen (University of Freiburg)

Björn Grüning (University of Freiburg)

Gustav Steinhoff (University of Rostock)

Robert David (University of Rostock)

Supported by:

denbi.de

elixir-europe.org

cardiac-stemcell-therapy.com

bmbf.de

# Work where others would like to spend their holidays

# Making sense out of data – providing meaning to models



Omic analyses and data integration

Data management & standardisation (Fairdom partner)

eHealth iOS Application

SYSTEMS BIOLOGY BIOINFORMATICS ROSTOCK

www.sbi.uni-rostock.de

Computational modeling & machine learning

Summer schools and workshops

Systems Medicine Pre-clinical trials

# References

- Steinhoff G, Nesteruk J, Wolfien M, Kundt G, The PERFECT Trial Investigators Group. Cardiac Function Improvement and Bone Marrow Response Outcome Analysis of the Randomized Perfect Phase III Clinical Trial of Intramyocardial CD133 + Application After Myocardial Infarction. *EBioMedicine*. 2017. doi.org/10.1016/j.ebiom.2017.07.022

- Lott SC, Wolfien M, Riege K, Bagnacani A, Wolkenhauer O, Hoffmann S, Hess WR. Customized workflow development and data modularization concepts for RNA-Sequencing and metatranscriptome experiments. *Journal of Biotechnology. 2017.* doi.org/10.1016/j.jbiotec.2017.06.1203

- Gruening BA, Fallmann J, Yusuf D, Will S, Erxleben A, Eggenhofer F, Houwaart T, Batut B, Videm P, Bagnacani A, Wolfien M, Lott SC, Hoogstrate Y, Hess WR, Wolkenhauer O, Hoffmann S, Akalin A, Ohler U, Stadler PF, Backofen R. The RNA workbench: best practices for RNA and high-throughput sequencing bioinformatics in Galaxy. *Nucleic Acids Research*. 2017. doi.org/10.1093/nar/gkx409

- Wolfien M, Rimmbach C, Schmitz U, Jung JJ, Krebs S, Steinhoff G, David R, Wolkenhauer O. TRAPLINE: A standardized and automated pipeline for RNA sequencing data analysis, evaluation and annotation. *BMC Bioinformatics. 2016. doi: 10.1186/s12859-015-0873-9*