

*“Programming is about trying to make the future less painful.  
It’s about making things easier for our teammates.”*

The Pragmatic Programmer  
Andy Hunt & Dave Thomas

# CeMM

SCIENCE IS OUR MEDICINE



CENTER FOR MEDICAL DATA SCIENCE  
MEDICAL UNIVERSITY OF VIENNA  
Institute of Artificial Intelligence



# MrBiomics More Time for Science!

---

Stephan Reichl

Predoctoral Fellow @ Bock Lab

November 6th, 2024

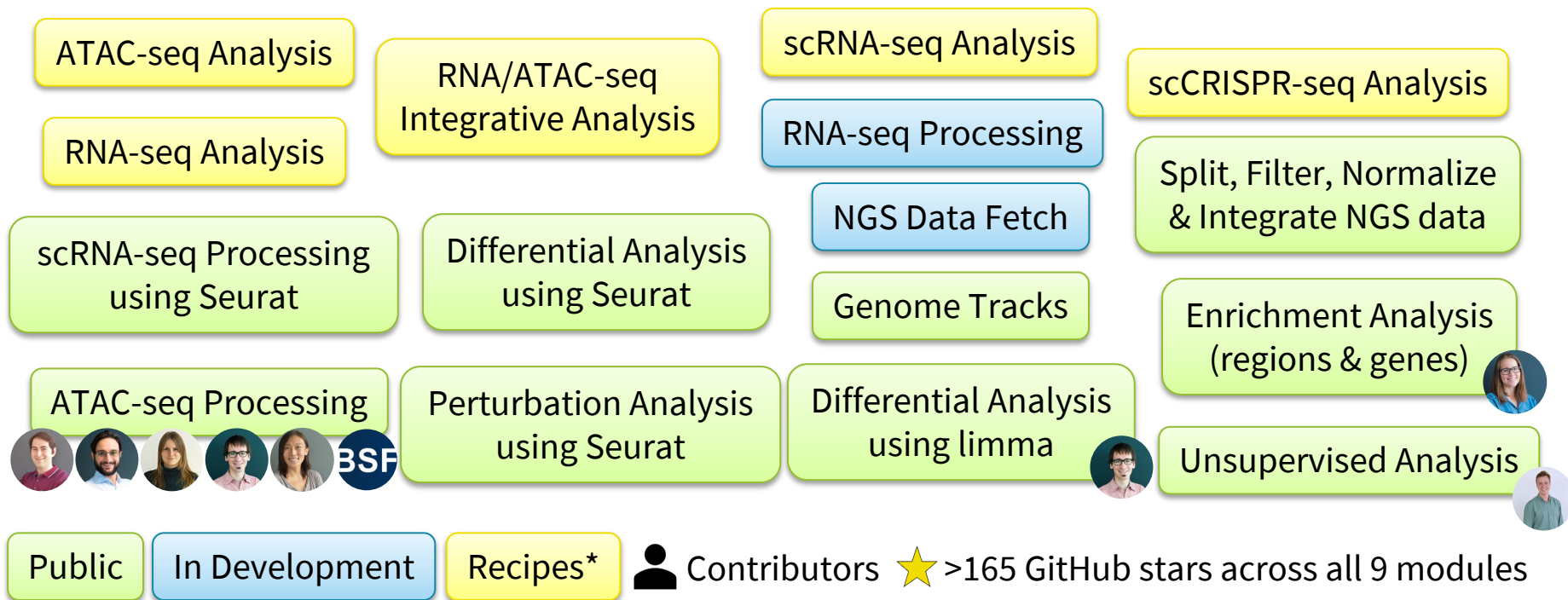
Research Center for Molecular Medicine  
of the Austrian Academy of Sciences

[www.cemm.at](http://www.cemm.at)

# MrBiomics – More Time for Science!

Modules & Recipes augment Bioinformatics for Multi-Omics Analyses at Scale

Achieve 80% of standard biomedical data science analyses semi-automatically with 20% effort by leveraging Snakemake's module functionality to use and combine pre-existing workflows into arbitrarily complex analyses.



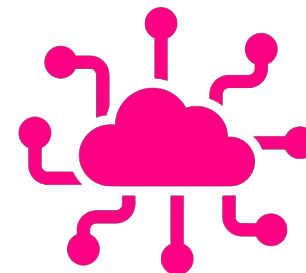
# Motivation – Three Observations at the End of 2021



Increased demand,  
but limited resources.



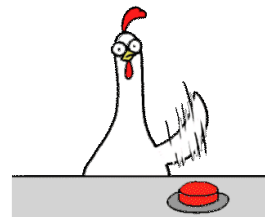
Increased fields of relevance,  
but not more time.



Reproducibility crisis and  
technological developments.

Over time it became clear that I created MrBiomics  
because I need it to handle my many multi-omics projects!

# Scope of MrBiomics



## What it facilitates (in scope):

- Accelerated arbitrarily complex end-to-end best practice analyses.
- Exploration of the computational option space and comparison of different approaches/hypotheses.
- Reproducible, transparent, documented, scalable, portable data analysis.

## What it doesn't facilitate (out of scope):

- Liberate you of thinking, the opposite is the case. It provides you with more time to think.
- Tell you how to use the supported methods e.g., parameters.
- Critically assess and interpret your results.

# Metaprogramming | Separation of Concerns

## Explained Using The World's Simplest Program

*"Out with the details!" Get them out of the code. While we're at it, we can make our code highly configurable and "soft"—that is, easily adaptable to changes.*

Not reusable,  
single-purpose code.

```
print("Hello world!")
```



Hello World!

Bad practice,  
but most first analyses.

Reusable by copy-paste-edit,  
but duplication of code base.


```
message = "Hello world!"  
...  
print(message)
```



Hello World!

Common practice,  
most analyses in papers.

Reusable.

 message: "Hello Mars"  
config.yaml



```
message = config["message"]  
...  
print(message)
```



Hello Mars!

Best practice,  
but rare in papers.

# Metaprogramming | Separation of Concerns

## Advantages

### Reusable.



message: "Hello Mars"



```
message = config["message"]
```

...

```
print(message)
```



Hello Mars!

Best practice,  
but rare in papers.

### Advantages

- “Do not repeat yourself” (DRY-)principle of coding
- Enables continuous improvement and compounding effects (i.e., solve it once for everyone → YAY Science).
- Reusable for different data, do not re-invent the wheel.
- Less error prone (i.e., change parameter, variable, path).
- Consistency, stability, robustness, ...
- Scalable, reproducible, portable, ...



Everyone

But Stephan, this  
looks like SO  
MUCH more work.

Okay, convinced,  
but it looks  
difficult...

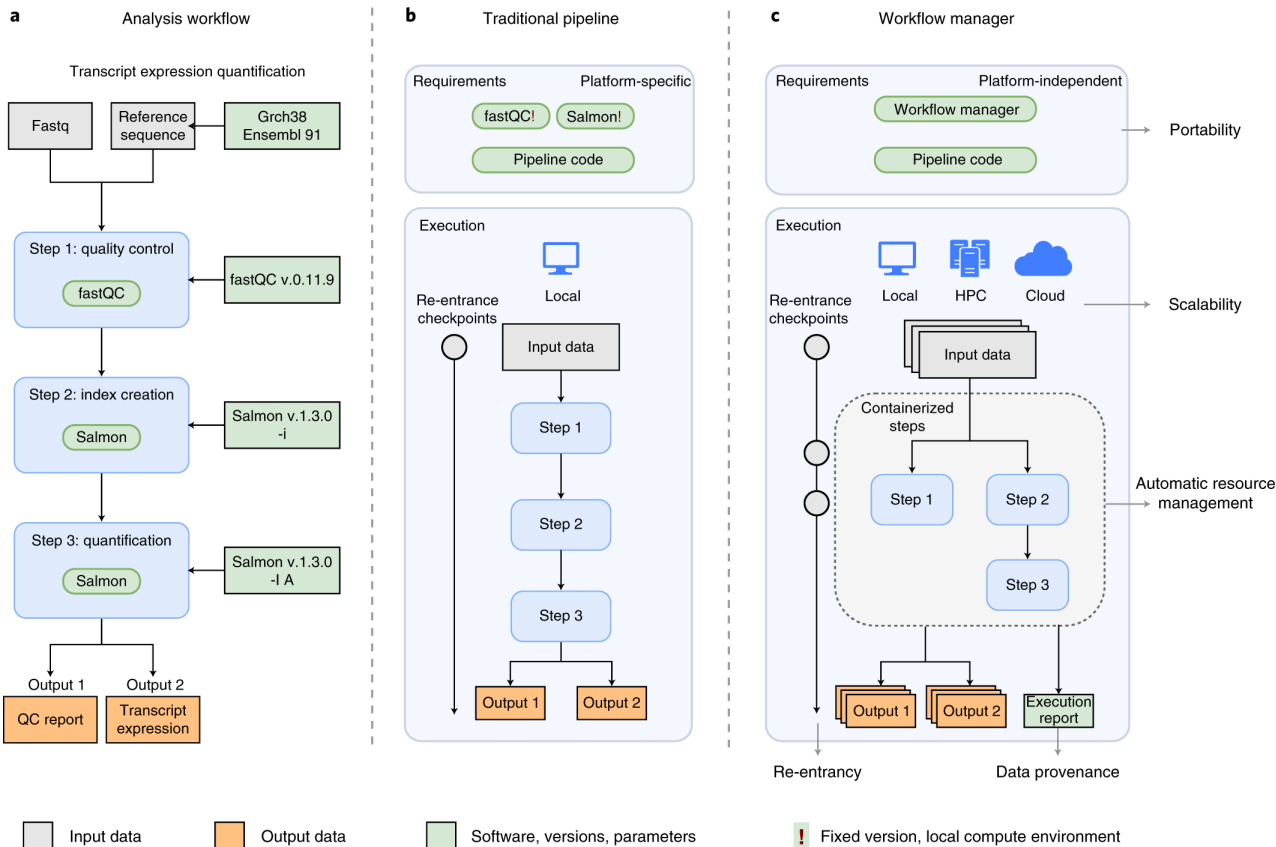
Yes and no. Only  
use it for reusable  
code that works.  
Copy-paste-edit  
doubles code .

True, but luckily  
someone else  
already did the  
heavy-lifting.



CeMM

# Workflow Management Systems & Snakemake



## What is Snakemake?

**A framework for reproducible and scalable data analysis**

- Readability (python based)
- Portability (conda & container)
- Modularization (script, notebooks, wrapper),
- Transparency (reports)
- Scalability (local, cluster, cloud)

## Highly popular

- >11 new citations per week
- >1,000,000 downloads
- Open source (MIT licensed)

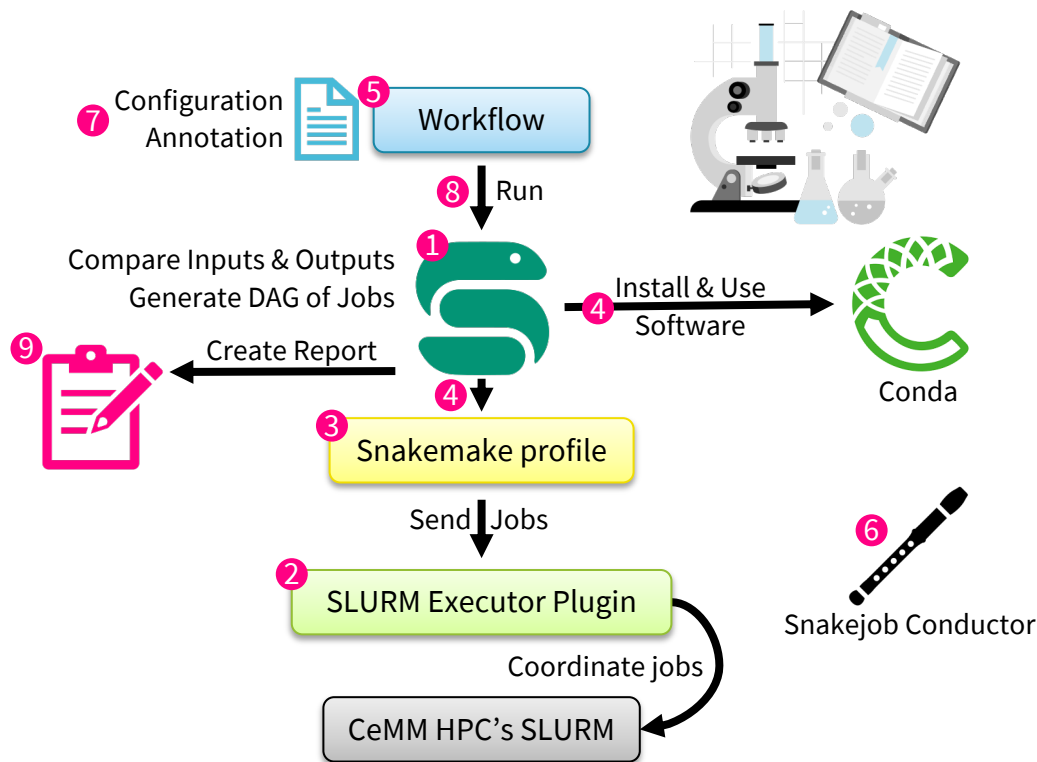


# Workflow Management by Snakemake



## How to run a Snakemake Workflow (at CeMM)

1. Install Snakemake<sup>(1)\*</sup>
2. Install SLURM executor plugin\*.
3. Clone CeMM's global Snakemake profile<sup>(2)\*</sup>
4. Set environment variables<sup>\*/\*\*</sup>
  - a. Conda environments
  - b. Snakemake profile
5. Clone/Deploy workflow
6. Setup Snakejob Conductor<sup>(2)\*\*</sup>
7. Configure workflow for analysis
8. Run workflow within Snakemake
9. Generate Snakemake report
10. Do great science and have fun!

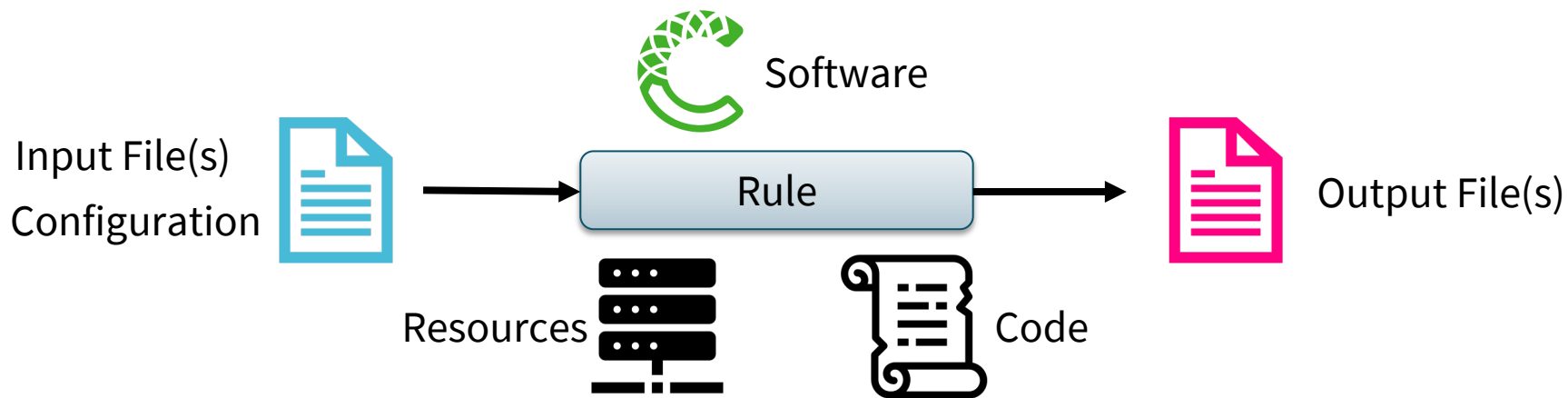


# Rules - Definition

**Rules** are specific computational tasks.

They can be bash commands, scripts, notebooks, or plain (Python) code.

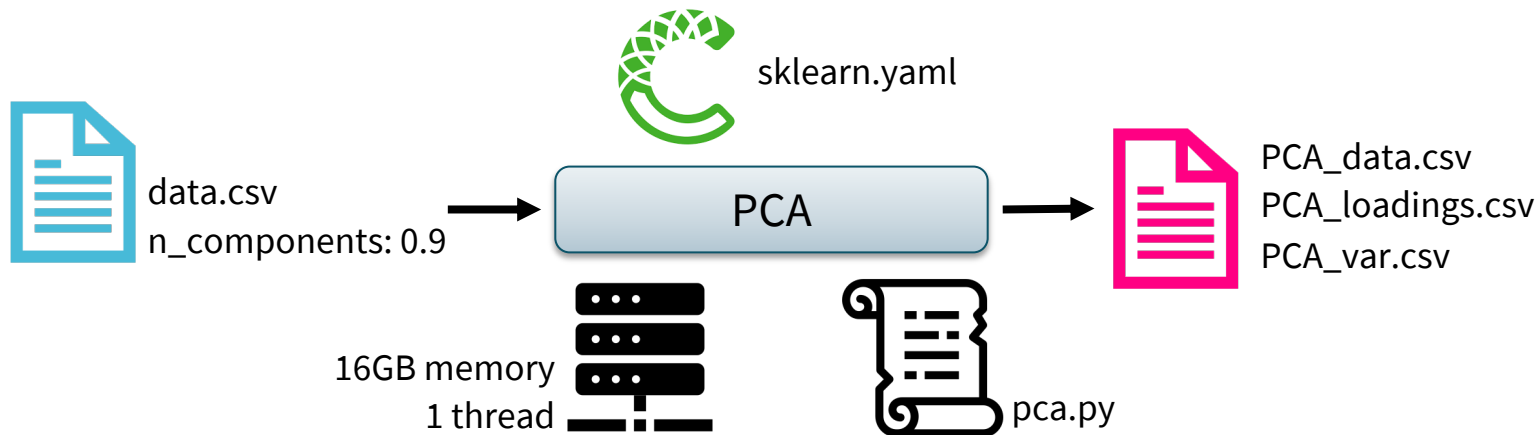
Input, output, software, configuration and computational resources are pre-defined.



# Rules – Example: PCA

Rule

Perform a Principal Component Analysis (PCA), using software provided in `sklearn.yaml` and code in `pca.py`, on `data.csv` with configuration `n_components: 0.9`. The job will get `16GB` of memory and `1` thread.



# Rules – Example: PCA - Code

Rule

```
##### perform Principal Component Analysis (PCA) #####
```

```
rule pca:
```

```
    input:
```

```
        unpack(get_sample_paths),
```

```
    output:
```

```
        ...
```

```
        result_data = os.path.join(result_path, '{sample}', 'PCA', 'PCA_{parameters}_data.csv'),
```

```
        ...
```

```
    resources:
```

```
        mem_mb=config.get("mem", "16000"),
```

```
    threads: config.get("threads", 1)
```

```
    conda:
```

```
        "../envs/sklearn.yaml"
```

```
    log:
```

```
        os.path.join("logs", "rules", "PCA_{sample}_{parameters}.log"),
```

```
    params:
```

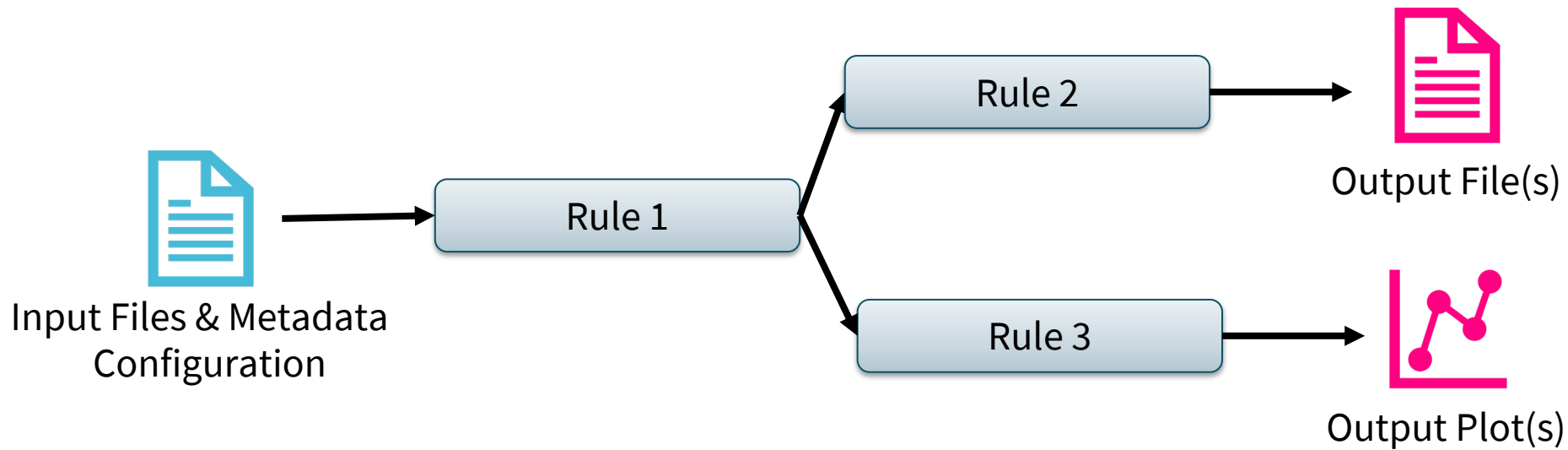
```
        samples_by_features = get_data_orientation,
```

```
    script:
```

```
        "../scripts/pca.py"
```

# Modules - Definition

**Modules** are Snakemake workflows, consisting of **Rules** for multi-step analyses. They can be general-purpose (e.g., Unsupervised Analysis) or modality-specific (e.g., RNA-seq).



Rule

## Visualization



# Sustainable & Reproducible via Documentation & Reports



Rapid changes and updates to projects are great, but there's one aspect of development that has long suffered from quick iteration: **project documentation**. – Techrepublic\*



Authors

Software

**Methods**

Features

Usage

Configuration

Examples

Links

Hook for  
Releases



Zenodo Repository  
for DOI

Automatic  
Curation



Snakemake Workflow Catalog



Snakemake Reports  
Self contained HTML

CeMM

# Projects using (multiple) Modules



You can (re-)use and combine pre-existing workflows within your projects by loading them as **Modules**. i.e., Workflows-of-Workflows.

```
# load local clones of a module
```

```
module MyData_other_workflow:
```

```
    snakefile: "path/to/other_workflow/Snakefile"
```

```
    config: config["MyData_other_workflow"]
```

```
use rule * from MyData_other_workflow as MyData_other_workflow_*
```

```
# load modules directly from GitHub
```

```
module MyData_other_workflow:
```

```
    snakefile: githup("epigen/unsupervised_analysis", path="workflow/Snakefile", tag="v2.0.0")
```

```
    config: config["MyData_other_workflow"]
```

```
use rule * from MyData_other_workflow as MyData_other_workflow_*
```

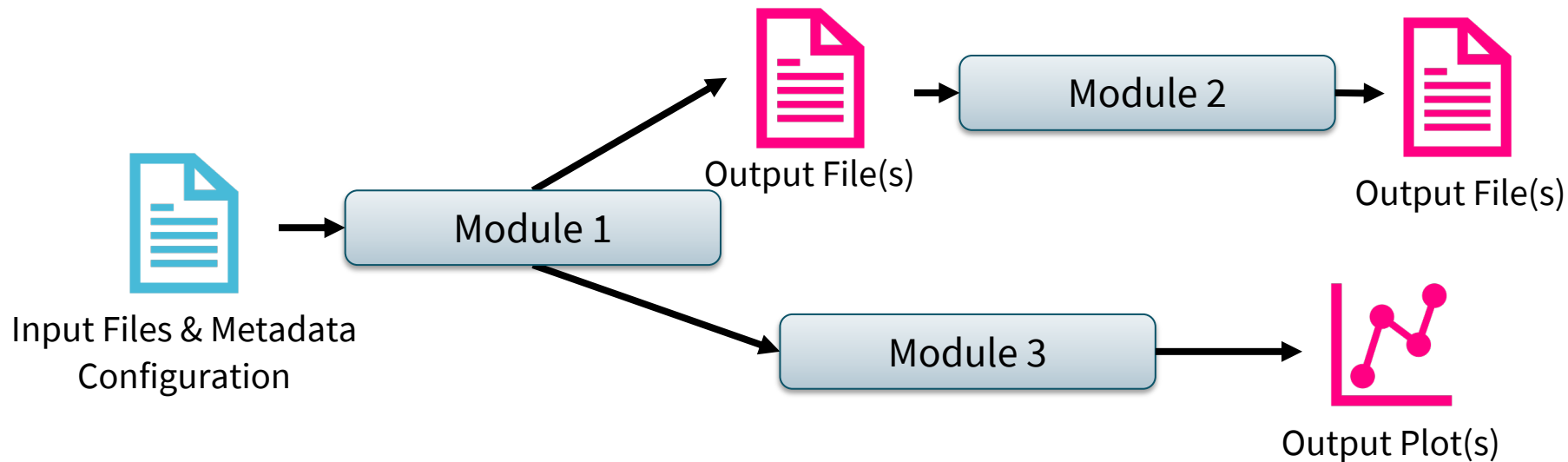
The combination of multiple modules into projects that analyze multiple datasets represents the overarching vision and superpower of MrBiomics.

CeMM

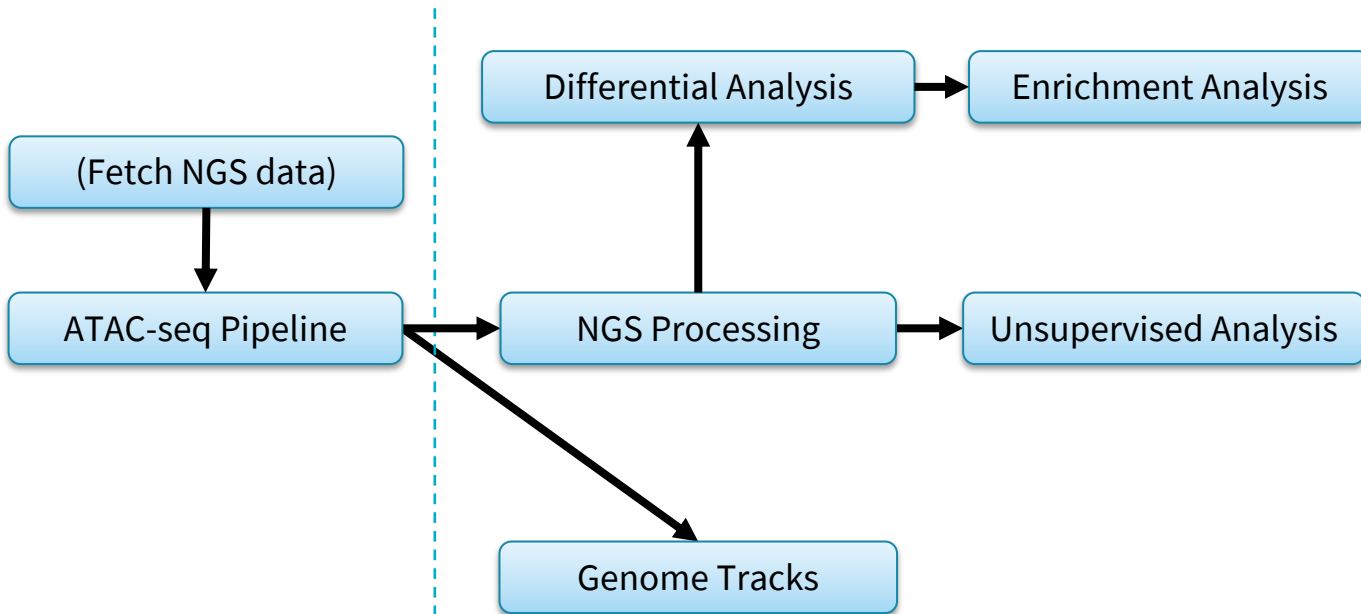


# Recipes - Definition

**Recipes** are combinations of existing modules into end-to-end best practice analyses.

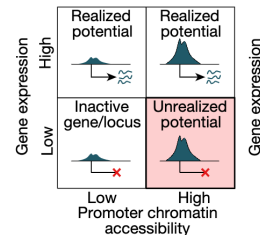


# Recipe for ATAC-seq Analysis

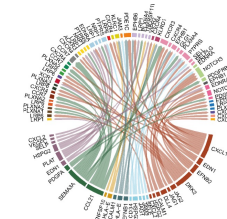


**More time for  
downstream analyses!**

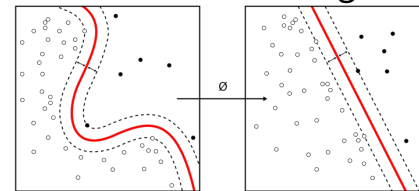
## Epigenetic potential



## Cell-cell communication

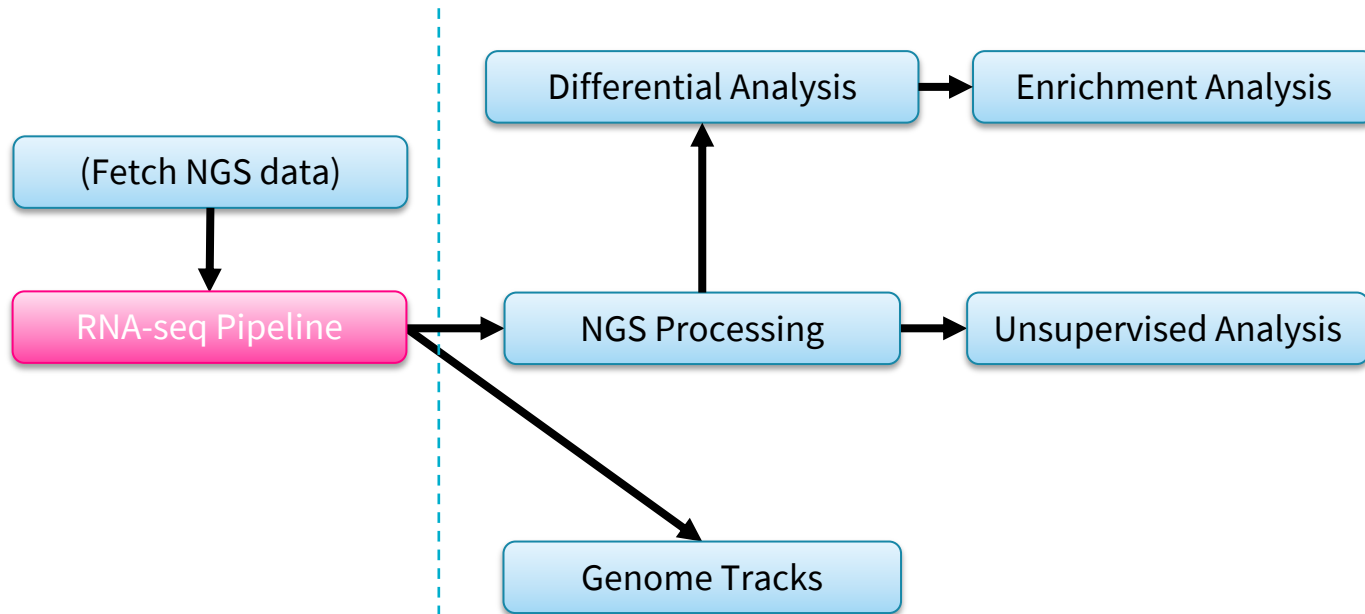


## Machine learning



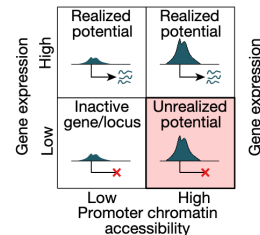
Data → Information → Knowledge

# Recipe for RNA-seq Analysis

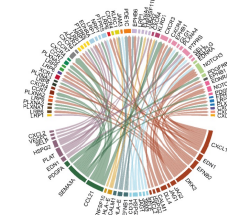


**More time for  
downstream analyses!**

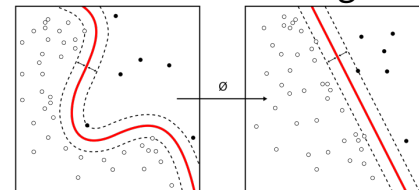
## Epigenetic potential



## Cell-cell communication

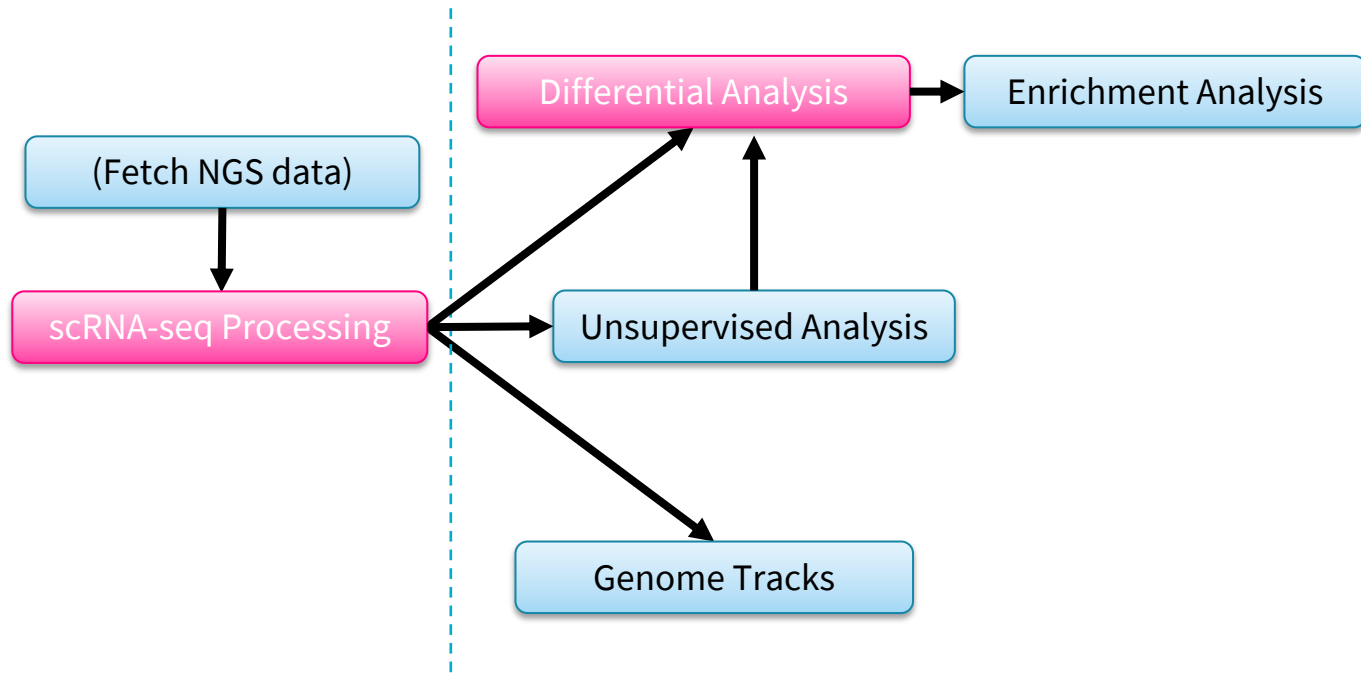


## Machine learning



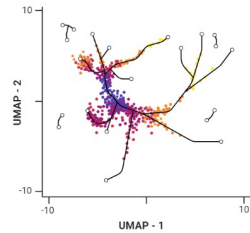
Data → Information → Knowledge

# Recipe for scRNA-seq Analysis

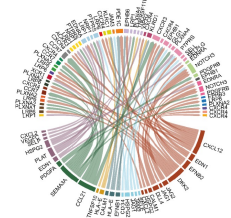


**More time for  
downstream analyses!**

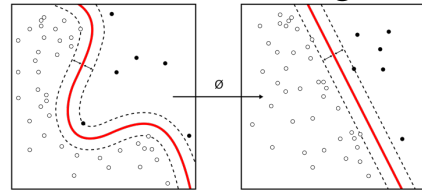
## Trajectory analysis



## Cell-cell communication



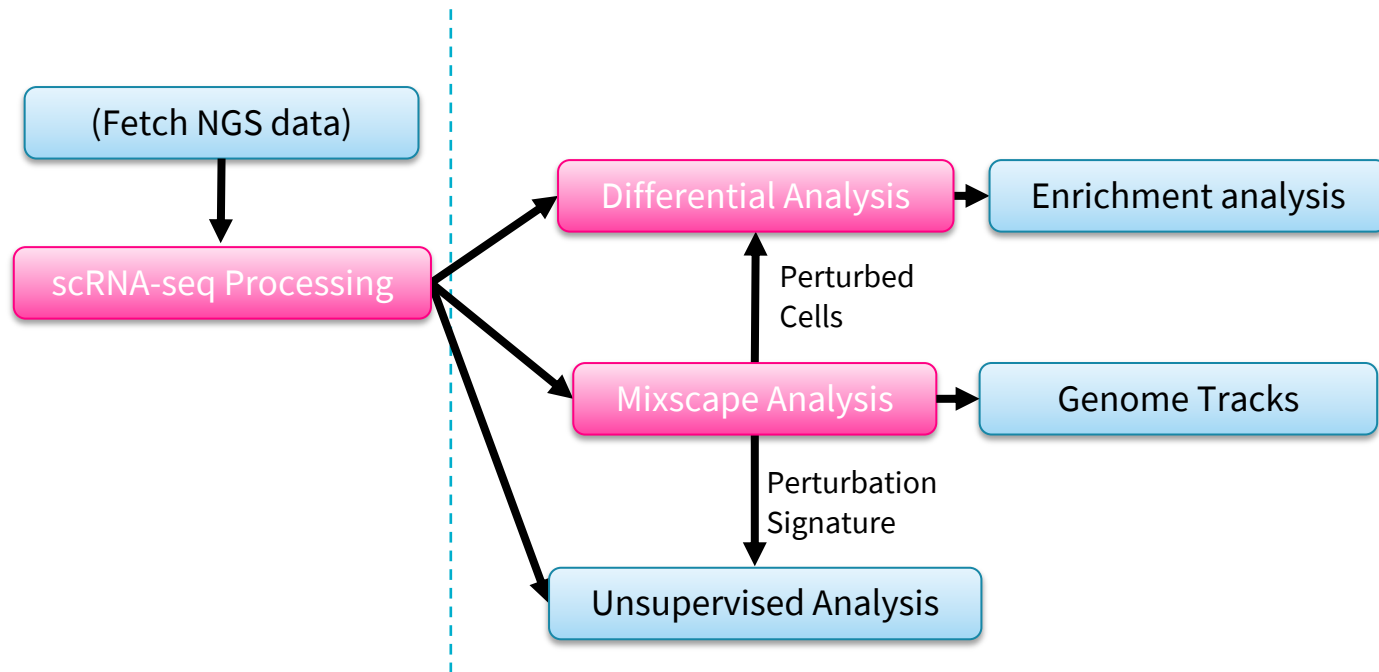
# Machine learning



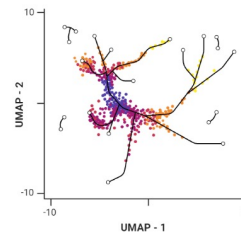
Data → Information → Knowledge

# Recipe for scCRISPR-seq Analysis

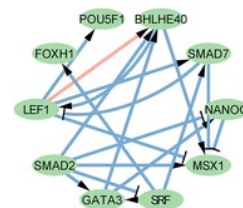
**More time for  
downstream analyses!**



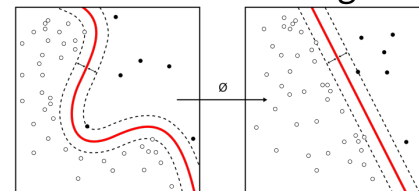
Trajectory analysis



Gene regulatory networks



Machine learning

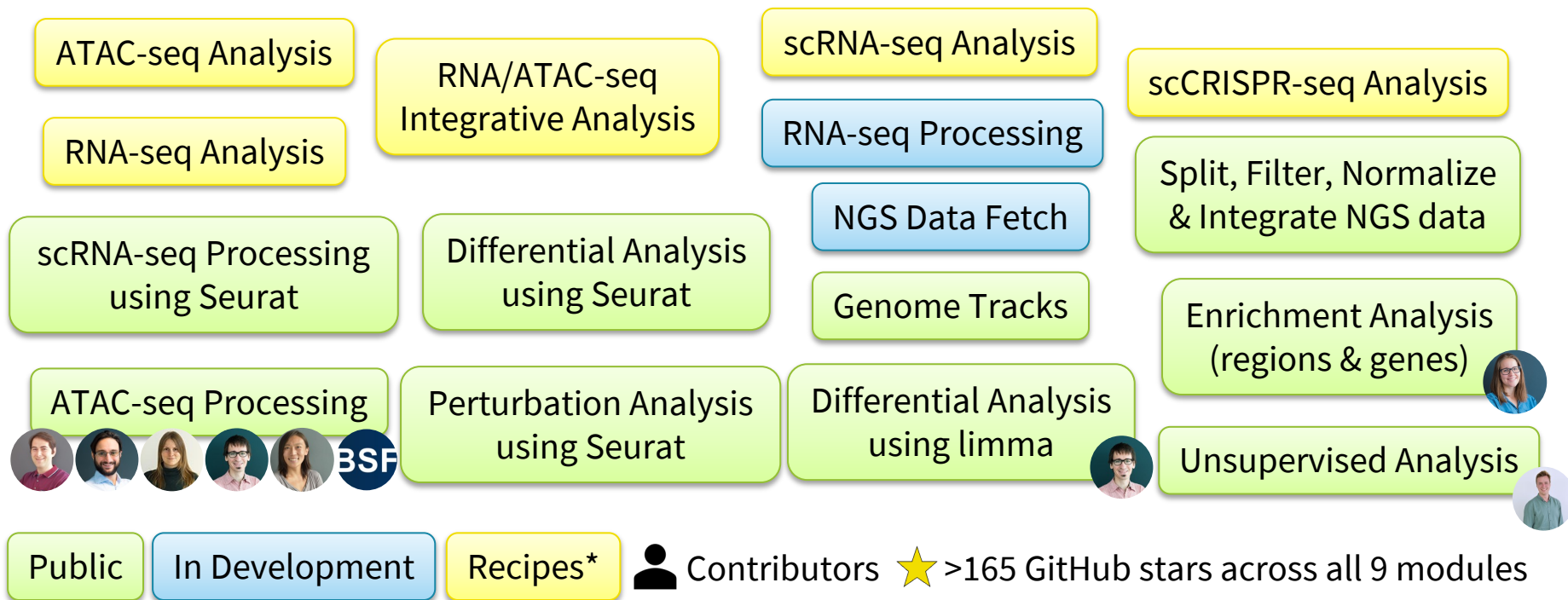


Data → Information → Knowledge

# MrBiomics – More Time for Science!

Modules & Recipes augment Bioinformatics for Multi-Omics Analyses at Scale

Achieve 80% of standard biomedical data science analyses semi-automatically with 20% effort by leveraging Snakemake's module functionality to use and combine pre-existing workflows into arbitrarily complex analyses.



# Resources

---

## MrBiomics

- List of current modules <https://github.com/stars/sreichl/lists/MrBiomics>
- Project repository: <https://github.com/epigen/MrBiomics>

## External Snakemake Workflows:

- snakePipes <https://snakepipes.readthedocs.io/en/latest/>
- seq2science <https://vanheeringen-lab.github.io/seq2science/index.html>
- Snakemake Workflow catalog <https://snakemake.github.io/snakemake-workflow-catalog/>

## Software:

- Snakemake Documentation <https://snakemake.readthedocs.io/en/stable/>
- Conda <https://docs.conda.io/en/latest/>