# Bioinformatics workflow management tools

Vojta Bystry
vojtech.bystry@ceitec.muni.cz

# Bioinformatics workflow (pipeline)

variant calling

report

CNV

alignment

trim primers

CEITEC

# Bioinformatics workflow (pipeline)

# Bioinformatic workflow management



**a** Analysis workflow

Transcript expression quantification

- Fastq
- Reference sequence
- Grch38 Ensembl 91

Step 1: quality control — fastQC ← fastQC v.0.11.9

Step 2: index creation — Salmon ← Salmon v.1.3.0 -i

Step 3: quantification — Salmon ← Salmon v.1.3.0 -I A

Output 1 — QC report
Output 2 — Transcript expression

**b** Traditional pipeline

Requirements — Platform-specific
- fastQC!
- Salmon!
- Pipeline code

Execution — Local
Re-entrance checkpoints

Input data → Step 1 → Step 2 → Step 3 → Output 1 / Output 2

**c** Workflow manager

Requirements — Platform-independent
- Workflow manager
- Pipeline code → Portability

Execution — Local / HPC / Cloud → Scalability
Re-entrance checkpoints

Input data

Containerized steps
Step 1 → Step 2 → Step 3 → Automatic resource management

Output 1 / Output 2 / Execution report

Re-entrancy — Data provenance

Legend:
- Input data
- Output data
- Software, versions, parameters
- ! Fixed version, local compute environment
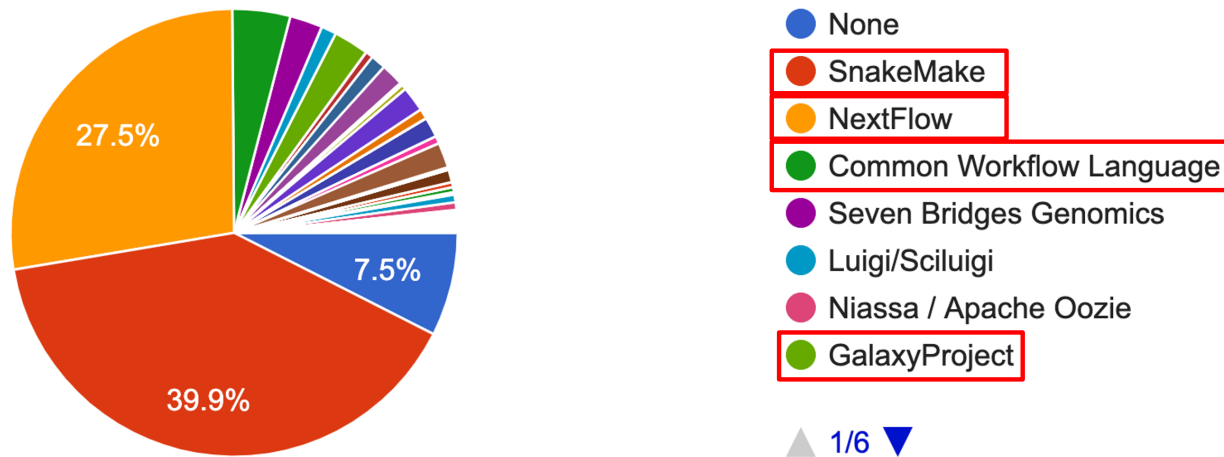
CEITEC

7

# Bioinformatic workflow management

- Reusability and Reproducibility

- Parallelization and Scale

- Error solving / debuging

# Bioinformatic workflow managers

Which Bioinformatics Workflow Manager / Tool / Platform / Standard do you use or prefer?

bit.ly/biowl

549 responses



- None
- SnakeMake
- NextFlow
- Common Workflow Language
- Seven Bridges Genomics
- Luigi/Sciluigi
- Niassa / Apache Oozie
- GalaxyProject

△ 1/6 ▽

# Common Workflow Language (CWL)

- Pushed by EU projects
- Not big grassroots community
- Scripts in .yaml format

```
hello_world.cwl

cwlVersion: v1.2

# What type of CWL process we have in this document.
class: CommandLineTool
# This CommandLineTool executes the linux "echo" command-line tool.
baseCommand: echo

# The inputs for this process.
inputs:
  message:
    type: string
    # A default value that can be overridden, e.g. --message "Hola mundo"
    default: "Hello World"
    # Bind this message value as an argument to "echo".
    inputBinding:
      position: 1
outputs: []
```

# Galaxy project

- Workflow manager with GUI
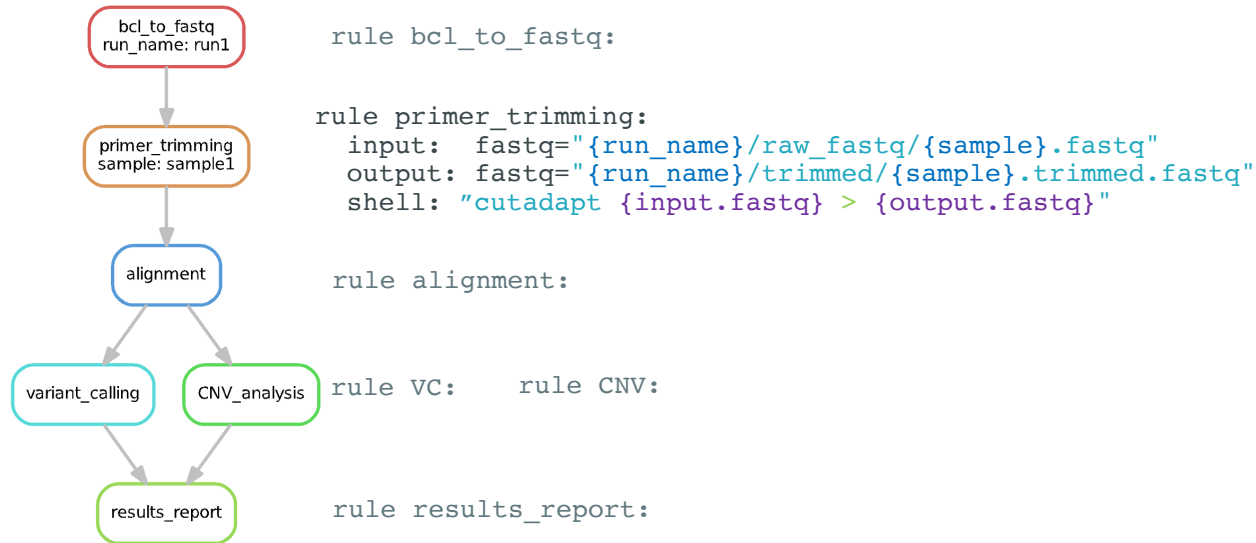- Biologists can do their own analysis ???
- It can work - EMBL

# Nextflow



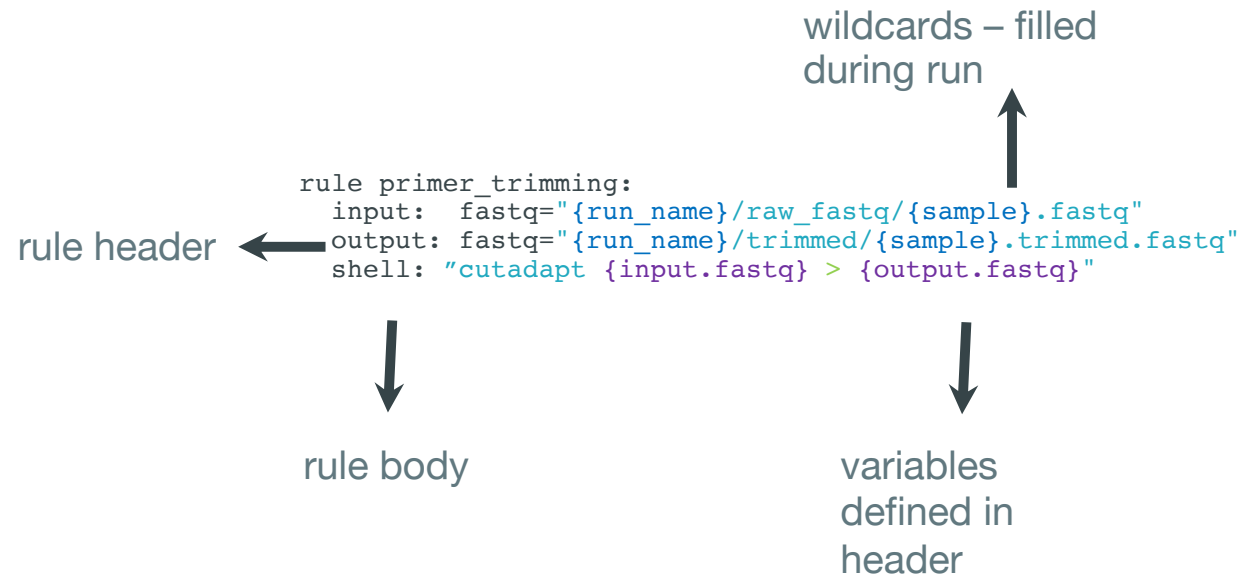- Great deployability
- Great existing workflow repository



```
1
2   #!/usr/bin/env nextflow
3
4   params.in = "$baseDir/data/sample.fa"
5
6   /*
7    * Split a fasta file into multiple files
8    */
9   process splitSequences {
10
11      input:
12      path 'input.fa'
13
14      output:
15      path 'seq_*'
16
17      """
18      awk '/^>/{f="seq_"++d} {print > f}' < input.fa
19      """
20  }
21
22  /*
23   * Reverse the sequences
24   */
25  process reverse {
26
27      input:
28      path x
29
30      output:
31      stdout
32
33      """
34      cat $x | rev
35      """
36  }
```

# Snakemake



```
rule bcl_to_fastq:


rule primer_trimming:
    input:  fastq="{run_name}/raw_fastq/{sample}.fastq"
    output: fastq="{run_name}/trimmed/{sample}.trimmed.fastq"
    shell: "cutadapt {input.fastq} > {output.fastq}"


rule alignment:



rule VC:      rule CNV:



rule results_report:
```

# Snakemake



wildcards – filled
during run

```
rule primer_trimming:
    input:  fastq="{run_name}/raw_fastq/{sample}.fastq"
    output: fastq="{run_name}/trimmed/{sample}.trimmed.fastq"
    shell: "cutadapt {input.fastq} > {output.fastq}"
```

rule header

rule body

variables
defined in
header

# Snakemake



```
rule bcl_to_fastq:

rule primer_trimming:
  input:  fastq="{run_name}/raw_fastq/{sample}.fastq"
  output: fastq="{run_name}/trimmed/{sample}.trimmed.fastq"
  shell: "cutadapt {input.fastq} > {output.fastq}"

rule alignment:

rule VC:        rule CNV:
                  output: "{run_name}/CNV/{sample}.CNV.tsv"

rule results_report:
  input: vcf = "{run_name}/VC/{sample}.vcf"
         cnv = "{run_name}/CNV/{sample}.CNV.tsv"
```

config_file.json

# Snakemake



► Simple shell script

```
shell: "mv —R {input} {output}"
```

► Combine languages

```
run:
if {params.cluster} is TRUE:
  R("cutree(hclust({input}),h = 7)")
else:
  shell("mv —R {input} {output}")
```
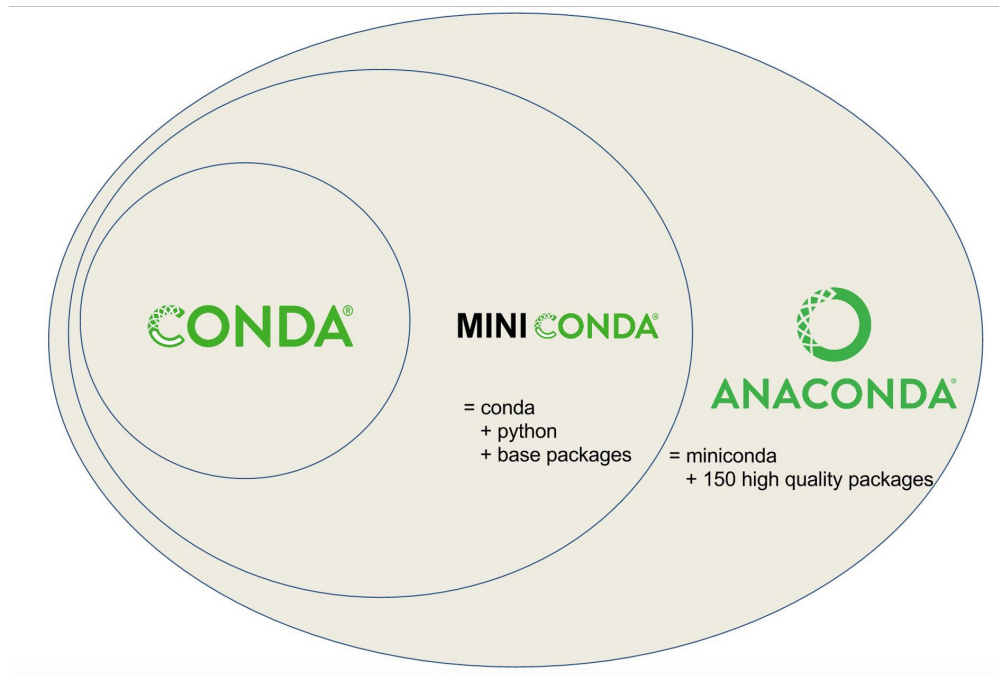
**python**

**R**

**shell**

► Wrap it in separate script

```
script: "my_script.py"
```

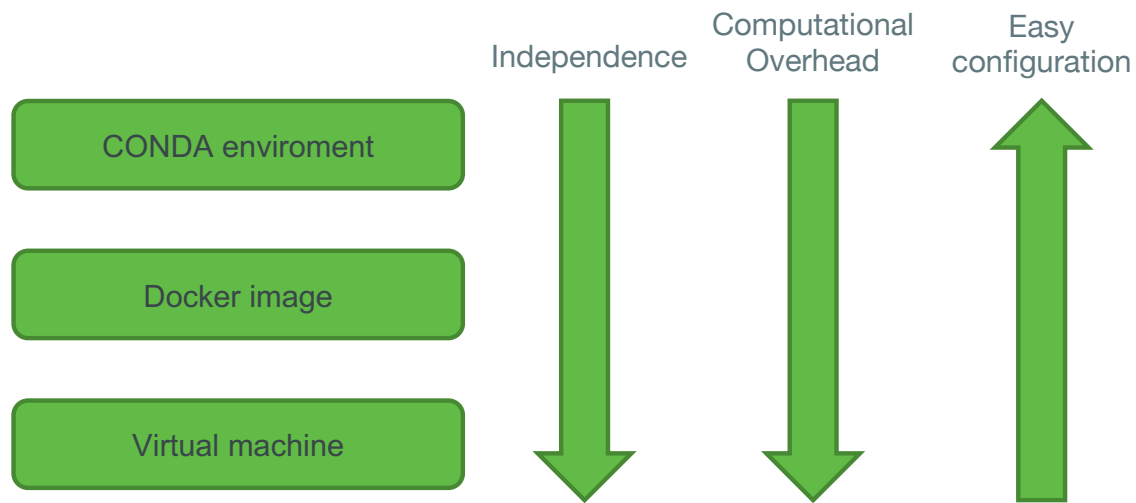► Separation of logic and functionality

  ► Organization

  ► Re-usability

# Conda / Anaconda / Bioconda



Bioconda is a distribution of bioinformatics software realized as a channel for the versatile Conda package manager.

# Conda

# Conda

CONDA

- Easy installation and management

- Instalation recepies:

```
conda install vardict
conda update vardict
conda remove vardict
conda env create -f myenv.yaml -n myenv
```

- Isolated environments:

```
channels:
 - conda-forge
 - defaults
dependencies:
 - pandas ==0.20.3
 - statsmodels ==0.8.0
 - r-dplyr ==0.7.0
 - r-base ==3.4.1
```

CEITEC

# Conda

CONDA

- Cheat sheet
  - https://docs.conda.io/projects/conda/en/4.6.0/_downloads/52a95608c49671267e40c689e0bc00ca/conda-cheatsheet.pdf

- Google it
  - conda [bioinformatics tool name]

# Computational resources and execution

- Snakemake is quite flexible in cluster execution
  - https://snakemake.readthedocs.io/en/stable/executing/cloud.html
  - ! Nothing works as advertise ☺