

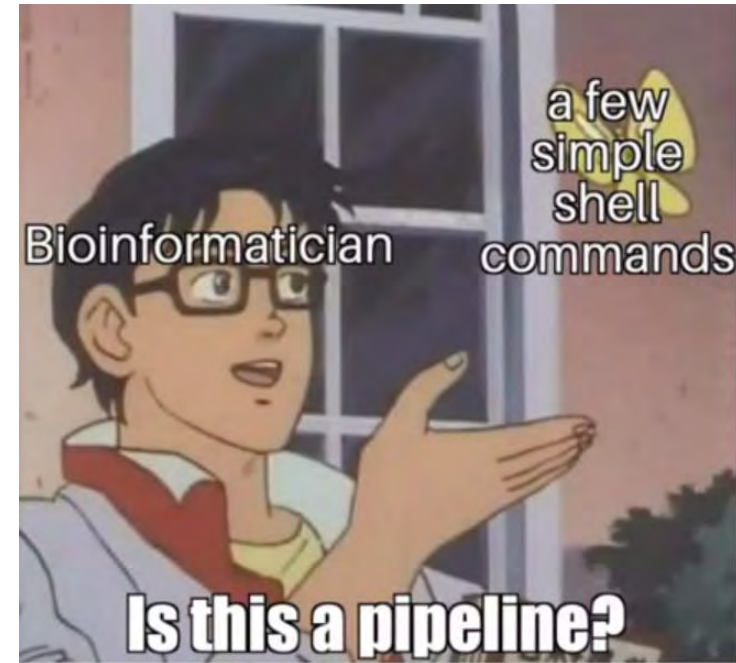
Пайпланы и пайпланы в биоинформатике

Пайплайны

Набор инструментов, объединенных для анализа и обработки данных

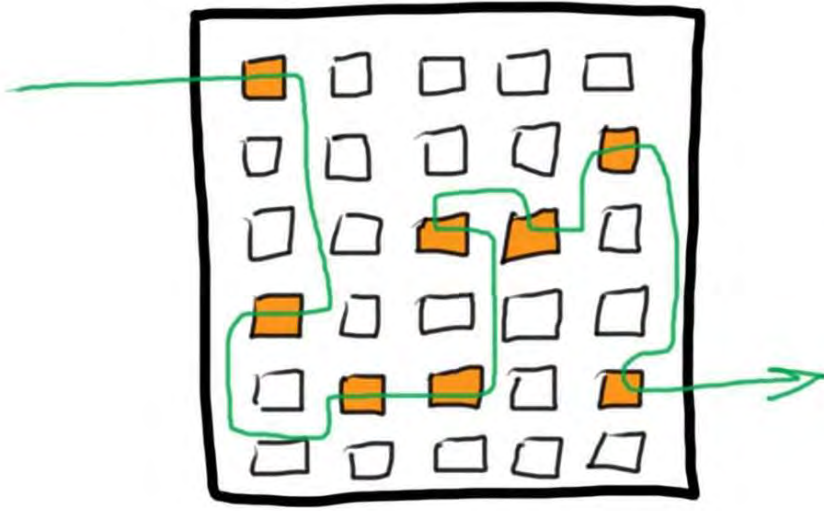
- Детерминированный
- Использует кэши
- Масштабируемый
- Простые и гибкие

Не нужно реализовывать функции снова и снова

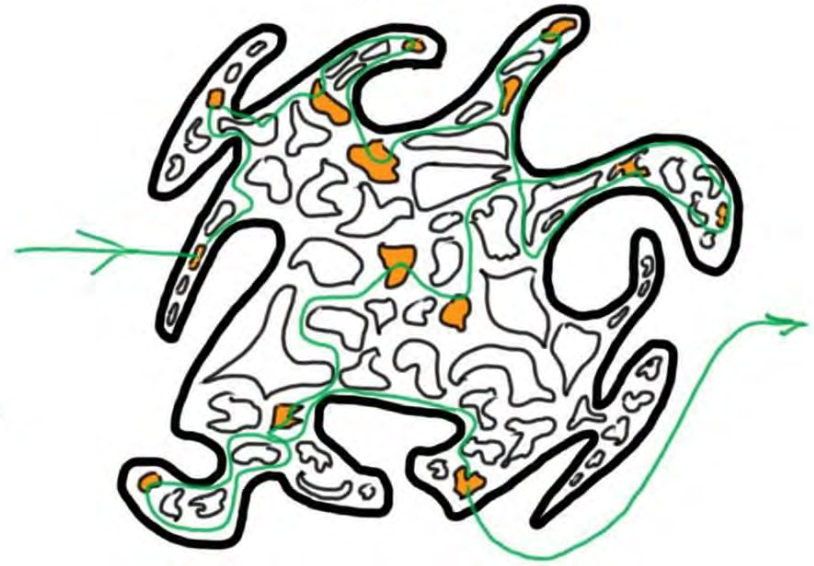


Зачем? У меня уже есть мои скрипты

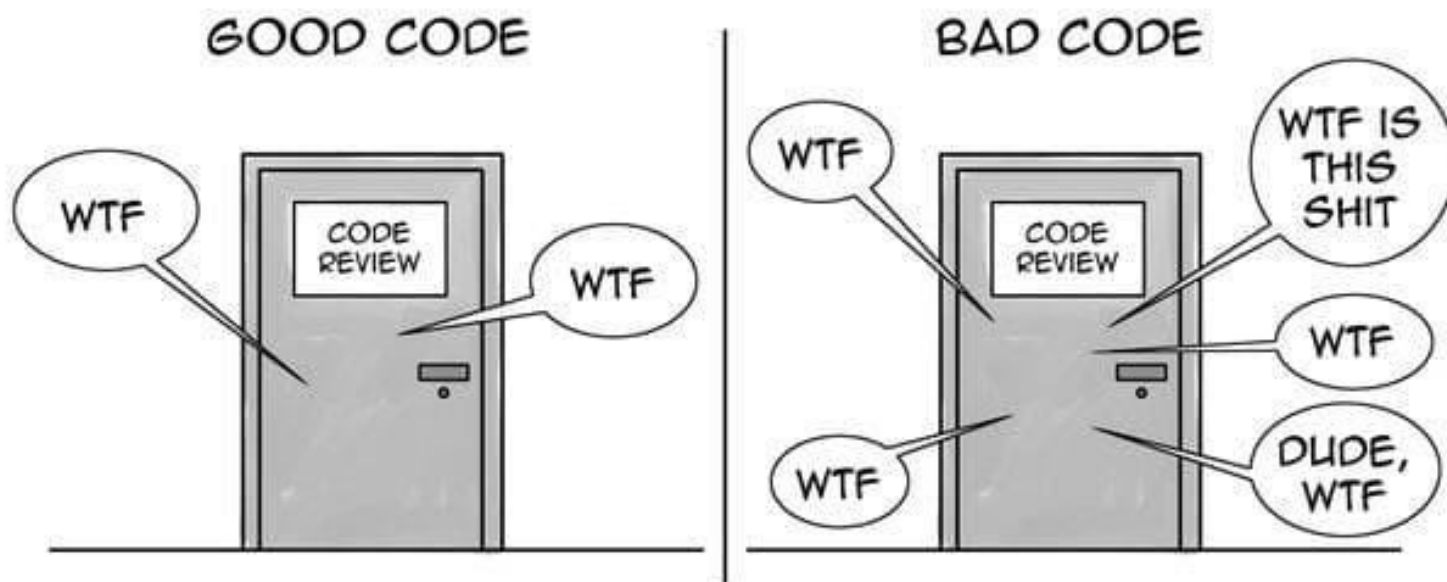
Finding our way
through clean code



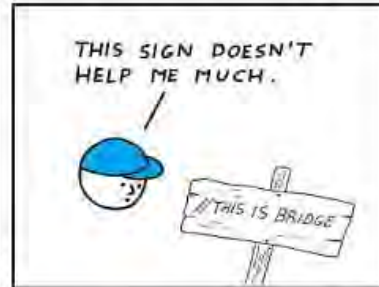
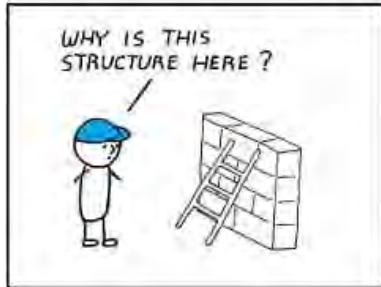
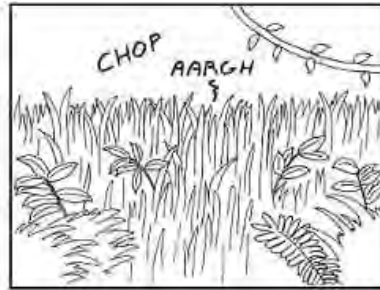
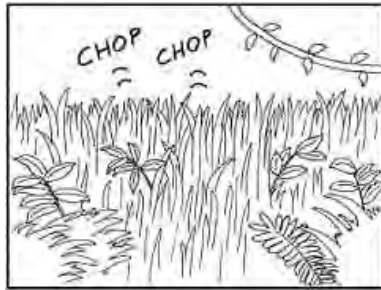
Finding our way
through bad code



Хороший или плохой



THE ONLY VALID MEASUREMENT OF CODE QUALITY: WTFs/MINUTE



I hate reading
other people's code.

Критерий	Пайплайн	Кастомный код
Гибкость	Высокая гибкость с использованием конфигурационных файлов.	Зависит от уровня навыков пользователя в программировании и настройке.
Скорость разработки	Быстрый старт с использованием готовых правил и компонентов.	Требуется времени на создание и отладку собственных скриптов.
Контроль и прозрачность	Четкий контроль зависимостей и прозрачность выполнения задач.	Полный контроль, но требует аккуратности в организации кода.
Сообщество и поддержка	Активное сообщество, множество готовых решений и обновлений.	Ограниченная поддержка, зависит от опыта пользователя.
Работа с уникальными задачами	Может быть адаптирован для уникальных и специфических задач.	Идеально подходит для уникальных и сложных задач, если написан правильно.
Зависимость от ресурсов	Требуется доступ к высокопроизводительным вычислительным ресурсам.	Зависит от конфигурации и требований пользователя.
Обучение и навыки	Требуется изучение синтаксиса, но относительно доступен.	Требуется глубоких навыков в программировании и биоинформатике.

Почему важны пайплайны

- Экономия Времени и Ресурсов
- Стандартизация и Надежность
- Обновления и Поддержка
- Общие Протоколы и Процессы
- Доступность для Неопытных Пользователей
- Расширенные Возможности

**MANUAL
EXECUTION**



**BASH OR
PYTHON SCRIPT**



**EMBARRASSING
PARALLELISM**

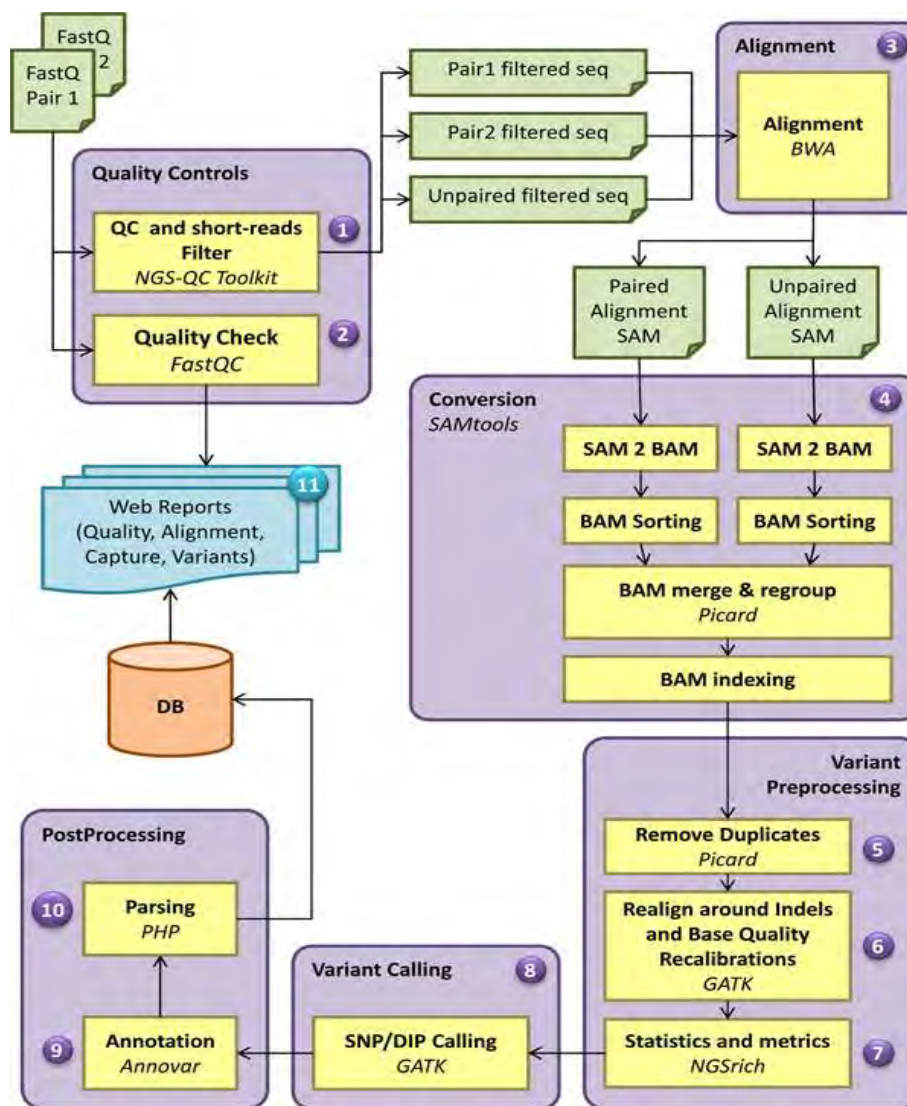


SNAKEMAKE



imgflip.com

Вычислительные графы



Пайпланы в биоинформатике

Snakemake

Snakemake

a

```

1  configfile: "config.yaml"
2
3  rule all:
4      input:
5          expand(
6              "results/plots/{country}.hist.pdf",
7              country=config["countries"]
8          )
9
10 rule download_data:
11     output:
12         "data/worldcitiespop.csv"
13     log:
14         "logs/download.log"
15     conda:
16         "envs/curl.yaml"
17     shell:
18         "curl -L https://burnisushi.net/stuff/worldcitiespop.csv > {output} 2> {log}"
19
20 rule select_by_country:
21     input:
22         "data/worldcitiespop.csv"
23     output:
24         "results/by-country/{country}.csv"
25     log:
26         "logs/select-by-country/{country}.log"
27     conda:
28         "envs/xsv.yaml"
29     shell:
30         "xsv search -s Country '{wildcards.country}' "
31         "{input} > {output} 2> {log}"
32
33 rule plot_histogram:
34     input:
35         "results/by-country/{country}.csv"
36     output:
37         "results/plots/{country}.hist.svg"
38     container:
39         "docker://falzanbashir/python-datascience:3.6"
40     log:
41         "logs/plot-hist/{country}.log"
42     script:
43         "scripts/plot-hist.py"
44
45 rule convert_to_pdf:
46     input:
47         "{prefix}.svg"
48     output:
49         "{prefix}.pdf"
50     log:
51         "logs/convert-to-pdf/{prefix}.log"
52     wrapper:
53         "0.47.0/utis/cairosvg"

```

b

c

```

import sys
sys.stderr = open(snakemake.log[0], "w")

import matplotlib.pyplot as plt
import pandas as pd

cities = pd.read_csv(snakemake.input[0])

plt.hist(cities["Population"], bins=50)

plt.savefig(snakemake.output[0])

```

d

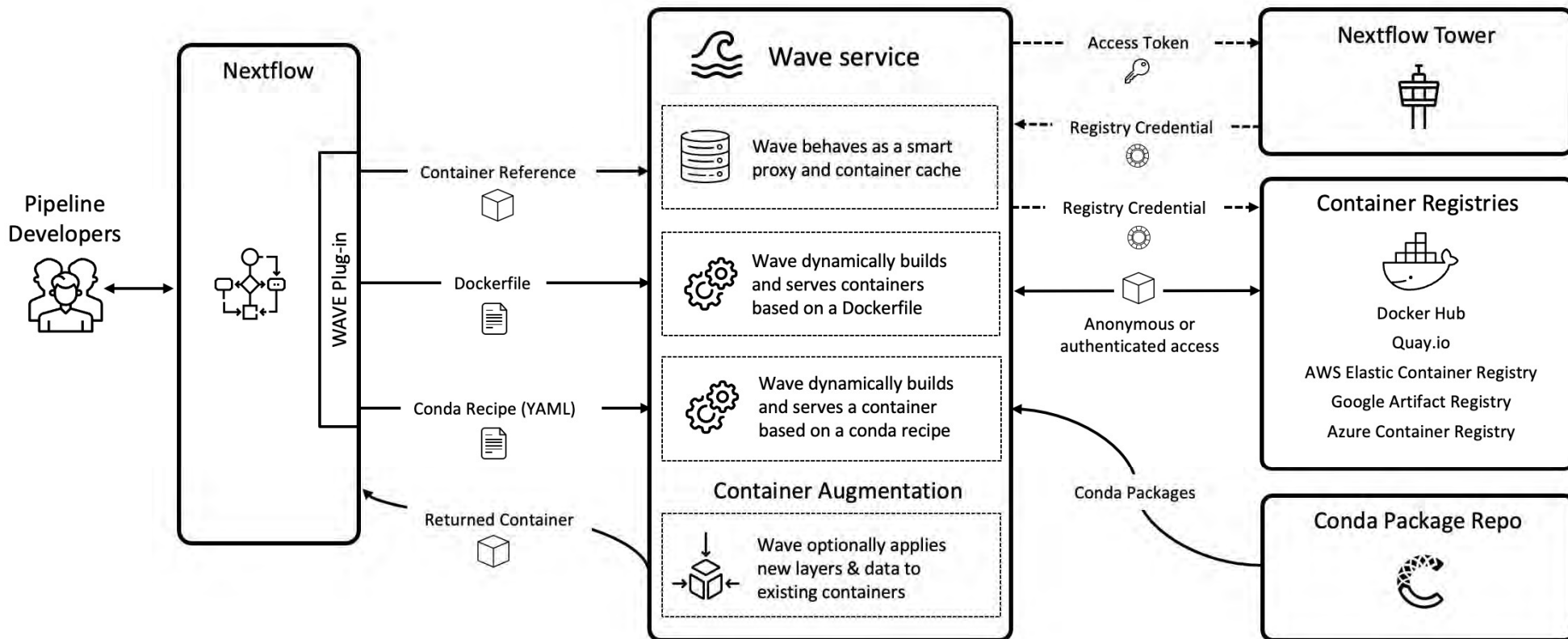
category	trivial	snakemake	technical	domain
1	18	0	0	0
2	4	0	1	0
3	4	0	2	15
4	0	0	2	4
5	0	0	1	1
6	0	0	1	1
7	0	2	0	0

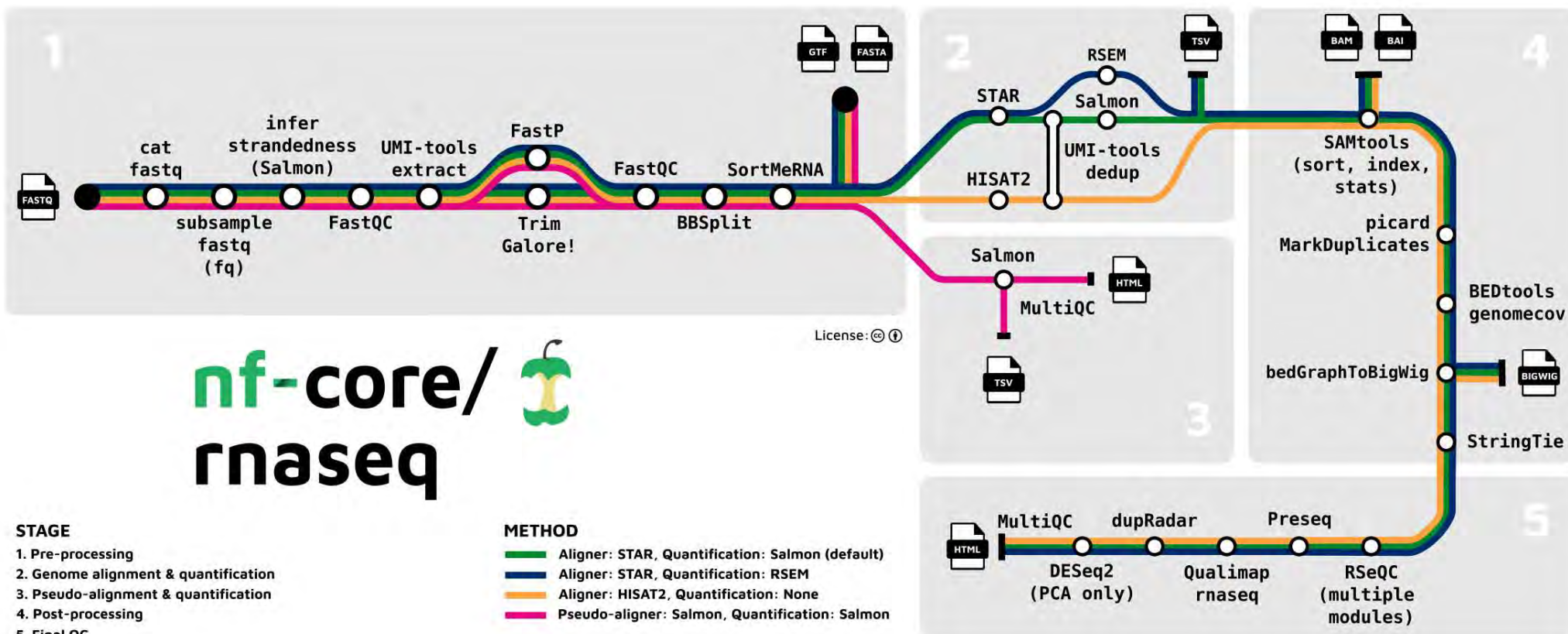
Standardized usage **196**

All workflows **2561**

Workflow	Description	Topics	QC	Stars	Watchers
<div>Usage</div> snakemake-workflows/rna-seq-star-deseq2	RNA-seq workflow using STAR and DESeq2	snakemake , sciworkflows , reproducibility , gene-expression-analysis , deseq2	<div>license MIT</div> <div>last commit november</div> <div>linting passed</div> <div>formatting passed</div>	282	11
<div>Usage</div> snakemake-workflows/dna-seq-gatk-variant-calling	This Snakemake pipeline implements the GATK best-practices workflow	reproducibility , snakemake , sciworkflows , genomic-variant-calling , gatk	<div>license MIT</div> <div>last commit may 2021</div> <div>linting passed</div> <div>formatting failed</div>	207	9

Nextflow





Семинар

