



ИНСТИТУТ
БИОИНФОРМАТИКИ

Анализ данных геномных секвенирований

Научный руководитель:
Геннадий Захаров

EPAM Systems, Lifesciences department

Студенты:

Скитченко Ростислав

Черенкова Ксения

Федотова Евгения

Санкт-Петербург
2017

Цель проекта:

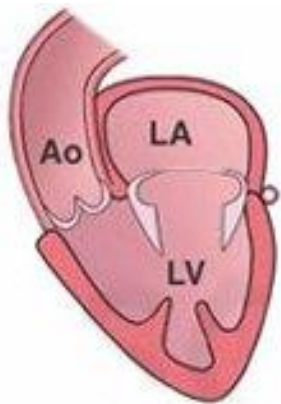
Сконструировать пайплайн для анализа данных, который помог бы медикам находить клинически значимые вариации.

Задачи проекта:

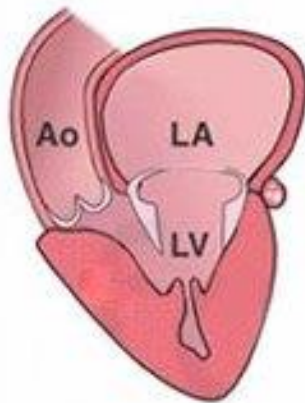
- Обработка ДНК-сиквенсов, нахождение вариаций
- Фильтрация и ранжирование вариаций по значимости
- Разработка переносимой системы, которая делает полную обработку данных автоматически от начала и до конца

Исходные данные

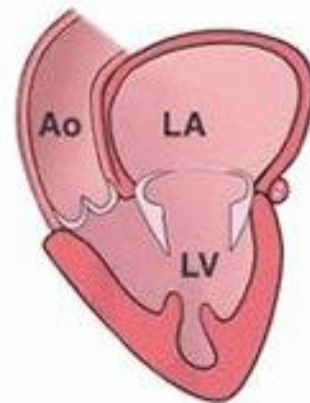
- Результаты направленного сиквенса генов членов трёх семей – кардиоассоциированные гены. Наборы хранятся в формате FASTQ.



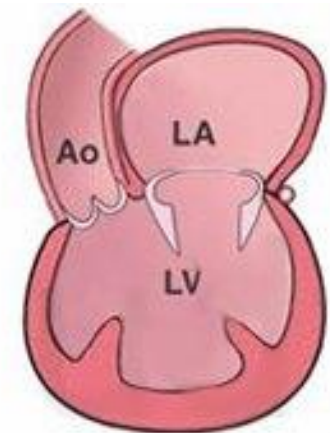
здоровое
сердце



гипертрофическая
КМП



рестриктивная
КМП

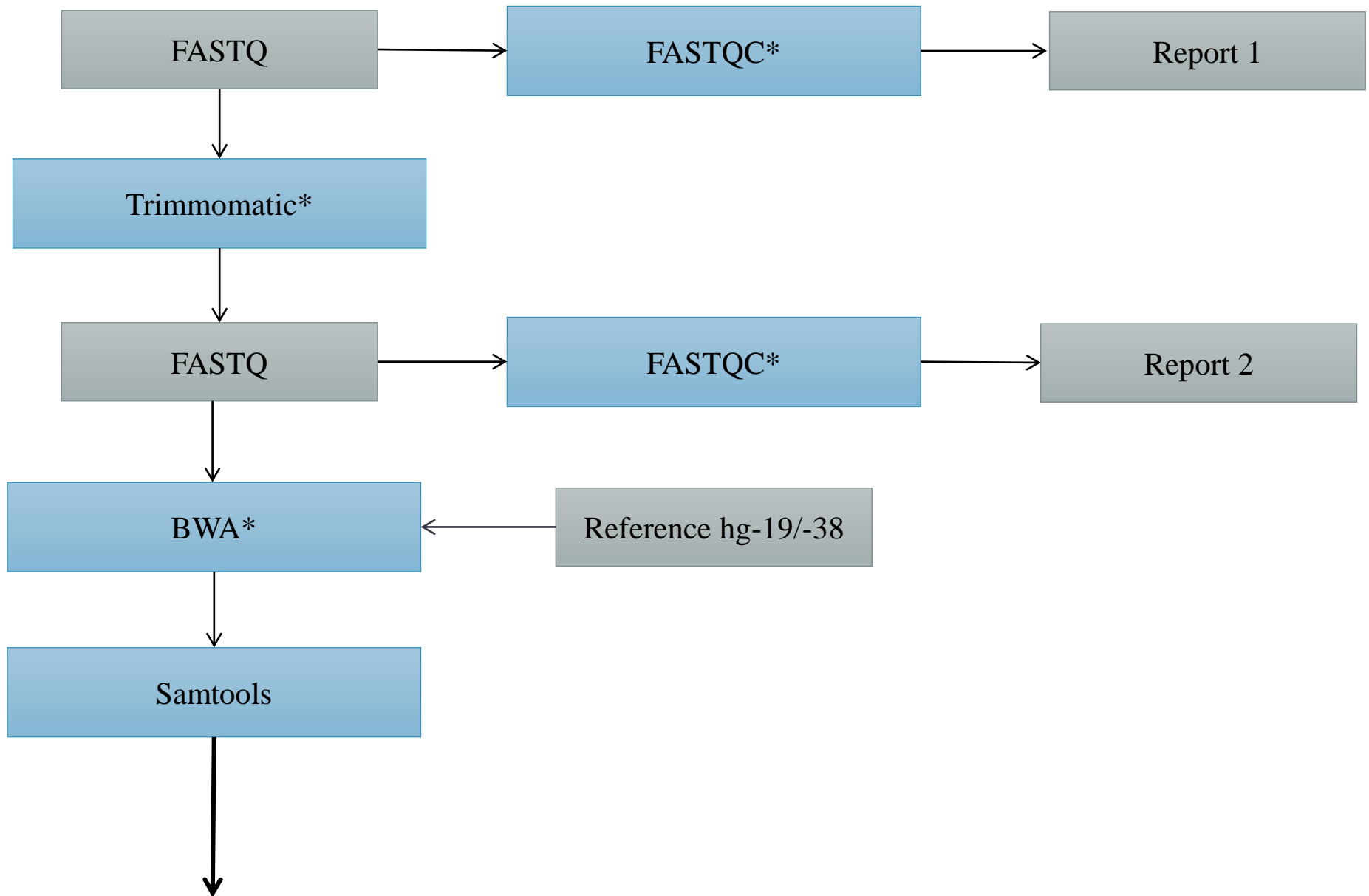


дилатационная
КМП

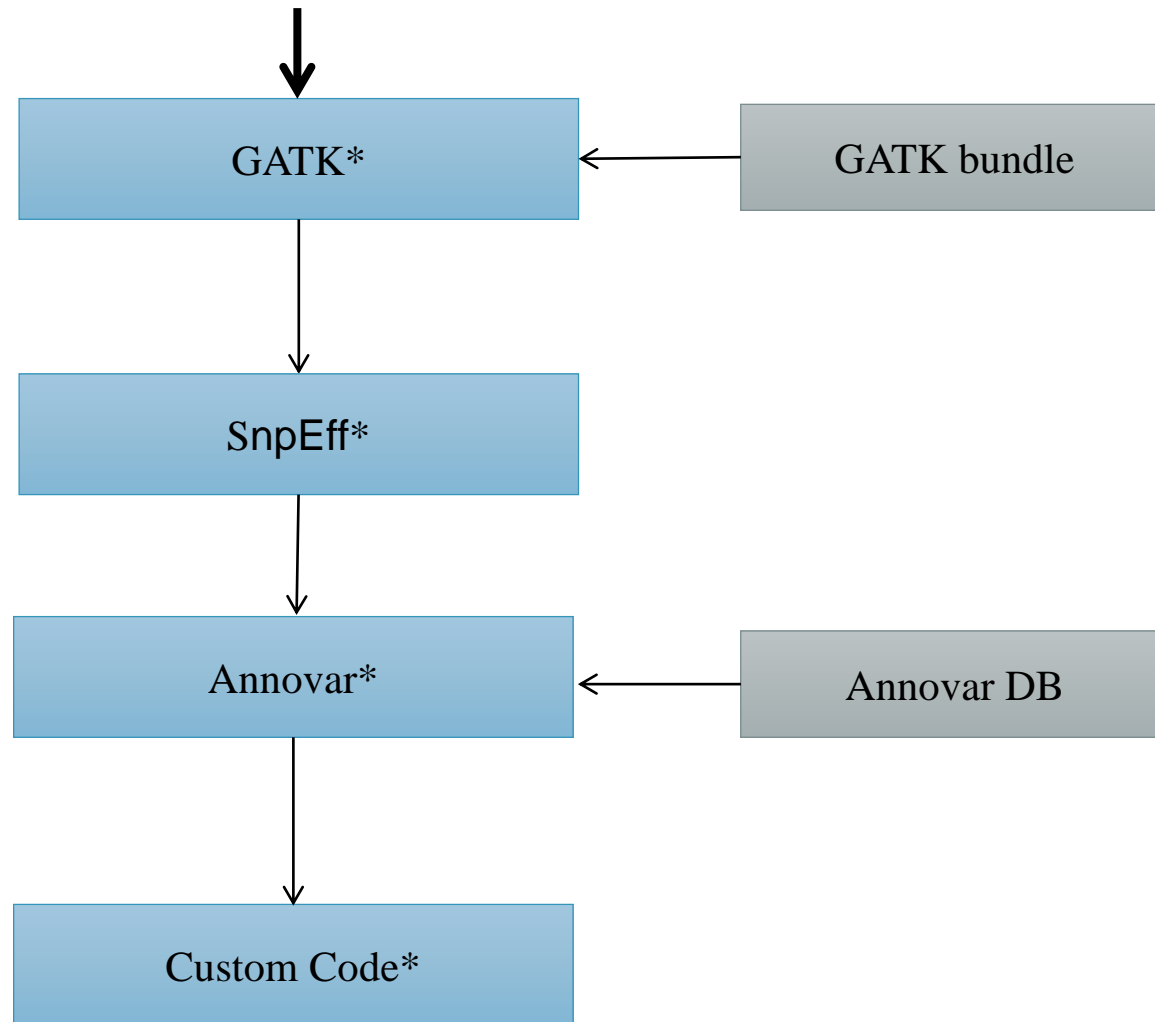
Используемые программы и ход работы

1. Fastqc (fastq)
 2. Trimmomatic (fastq)
 3. BWA (sam)
 4. Samtools (bam)
 5. Picard-tools mark duplicates (bam)
 6. GATK (VCF)
 7. SnpEff (VCF)
 8. Annovar (VCF)
- Custom code (txt)*

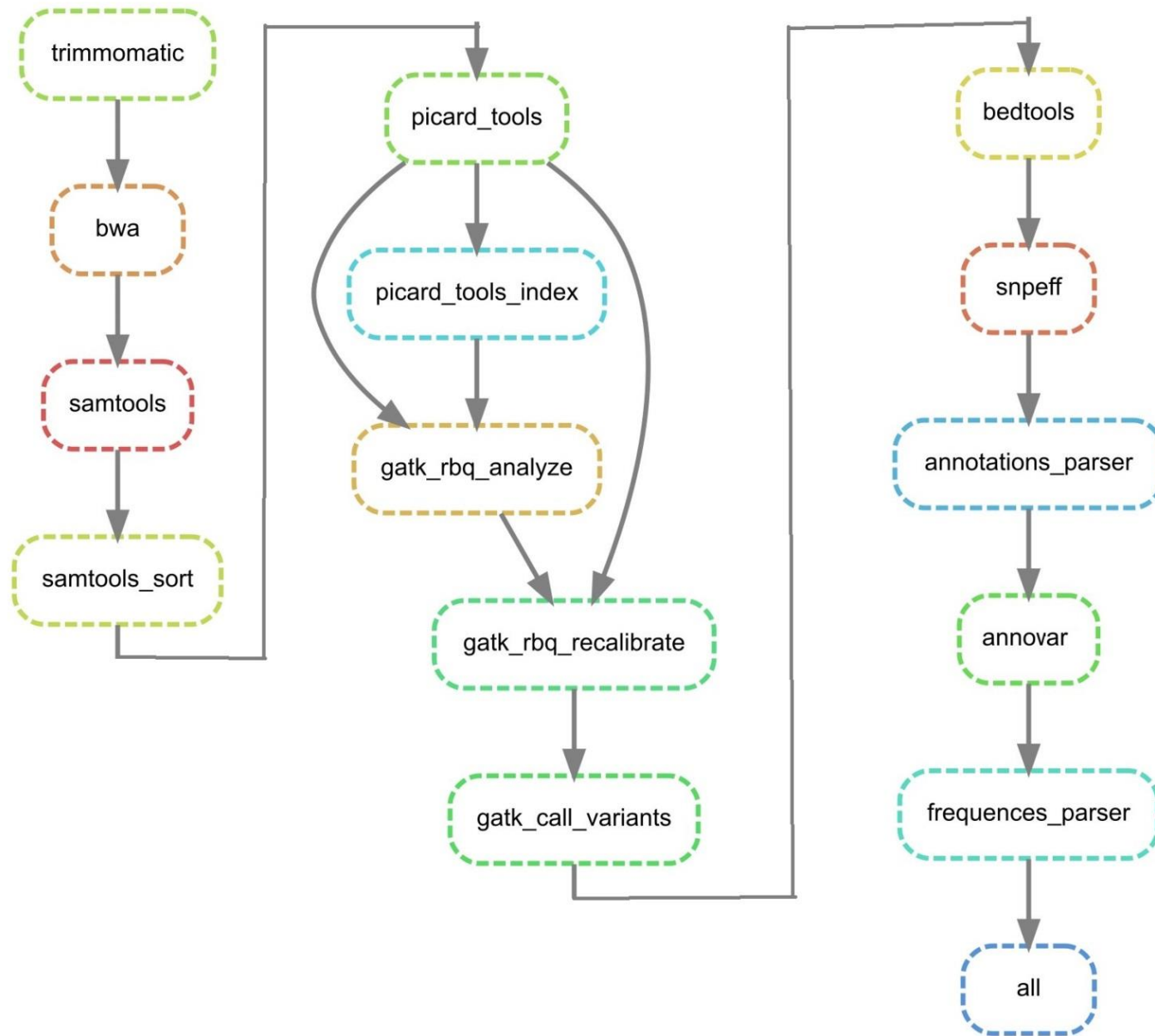
Общая схема анализа (Часть 1):



Общая схема анализа (Часть 2):



Инструменты для пайплайна: язык Snakemake



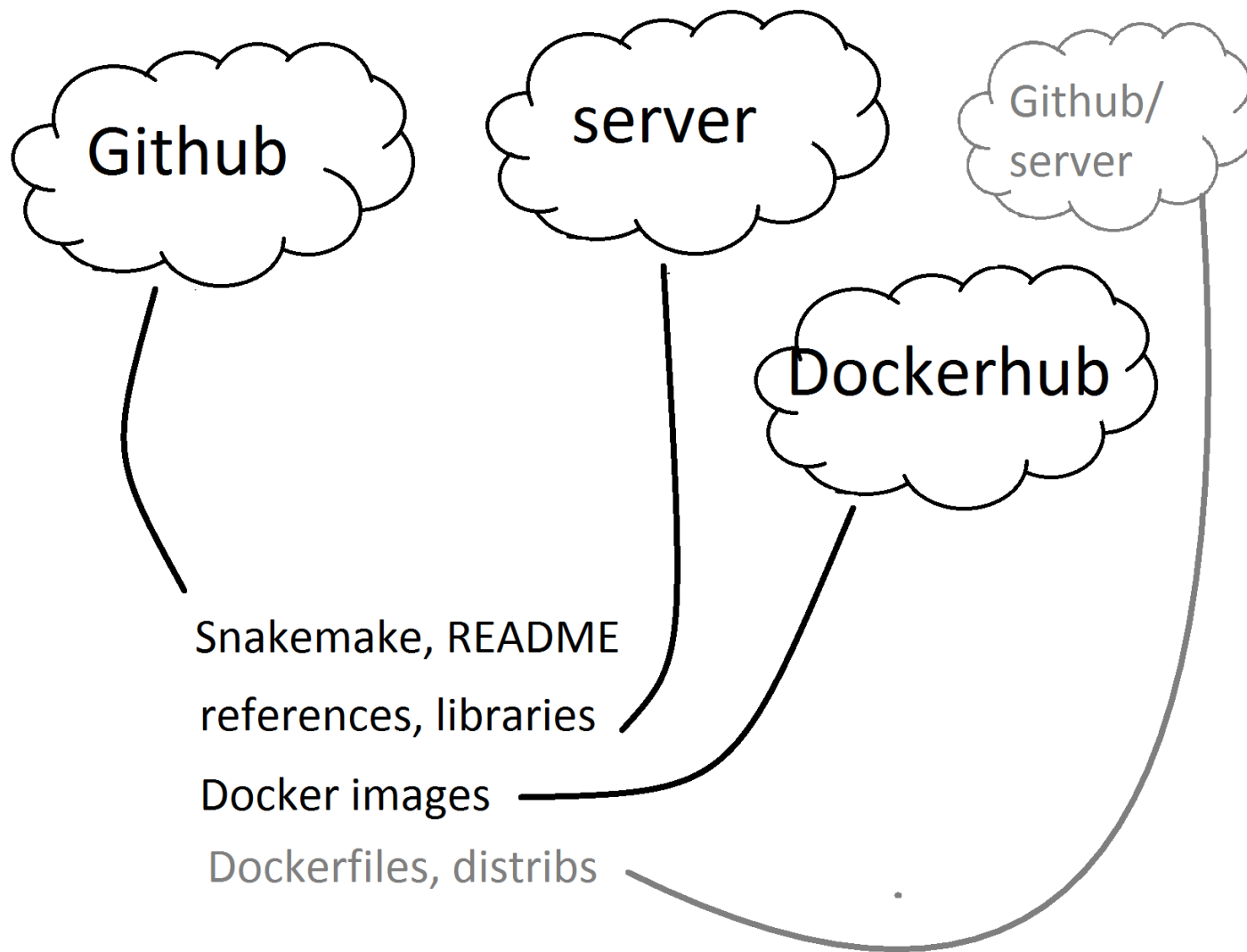
Инструменты для пайплайна: ПО Docker



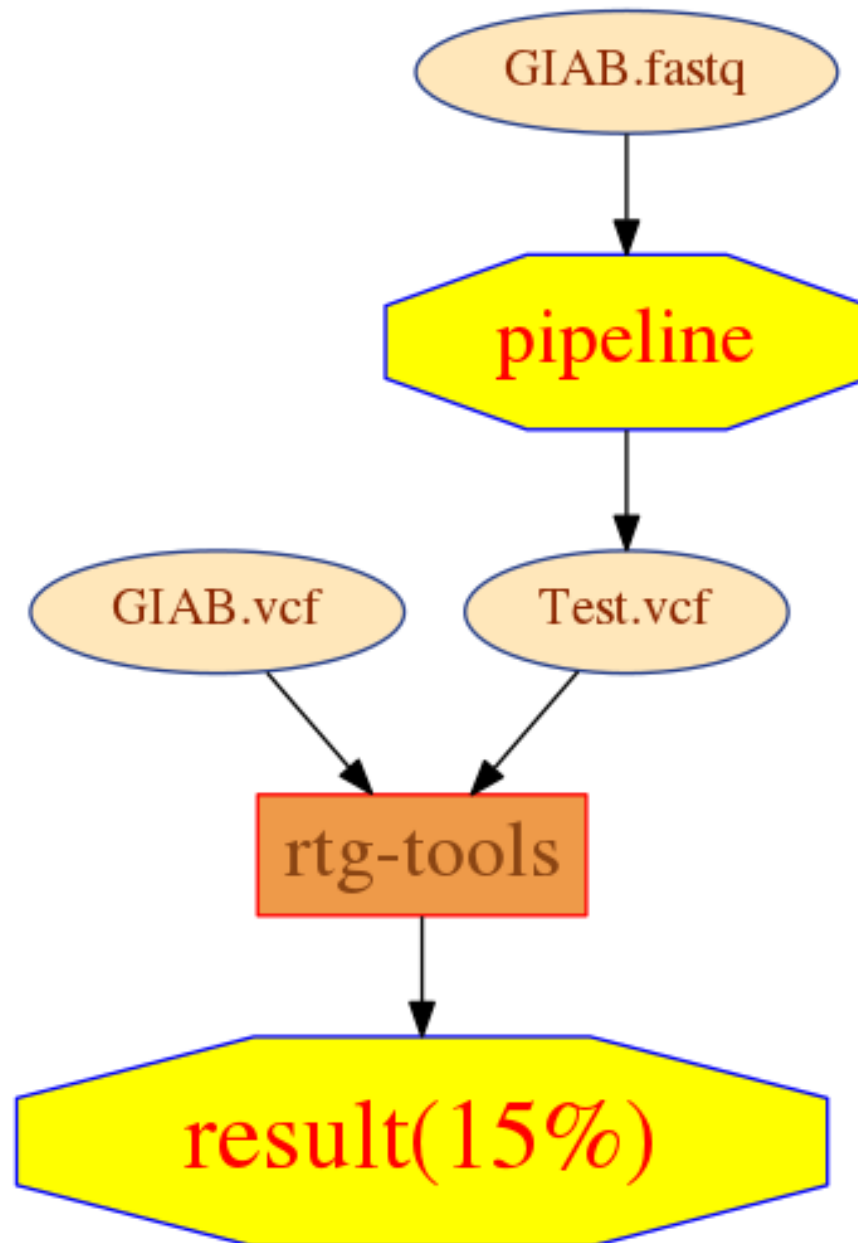
Сложности при работе с пайплайном:

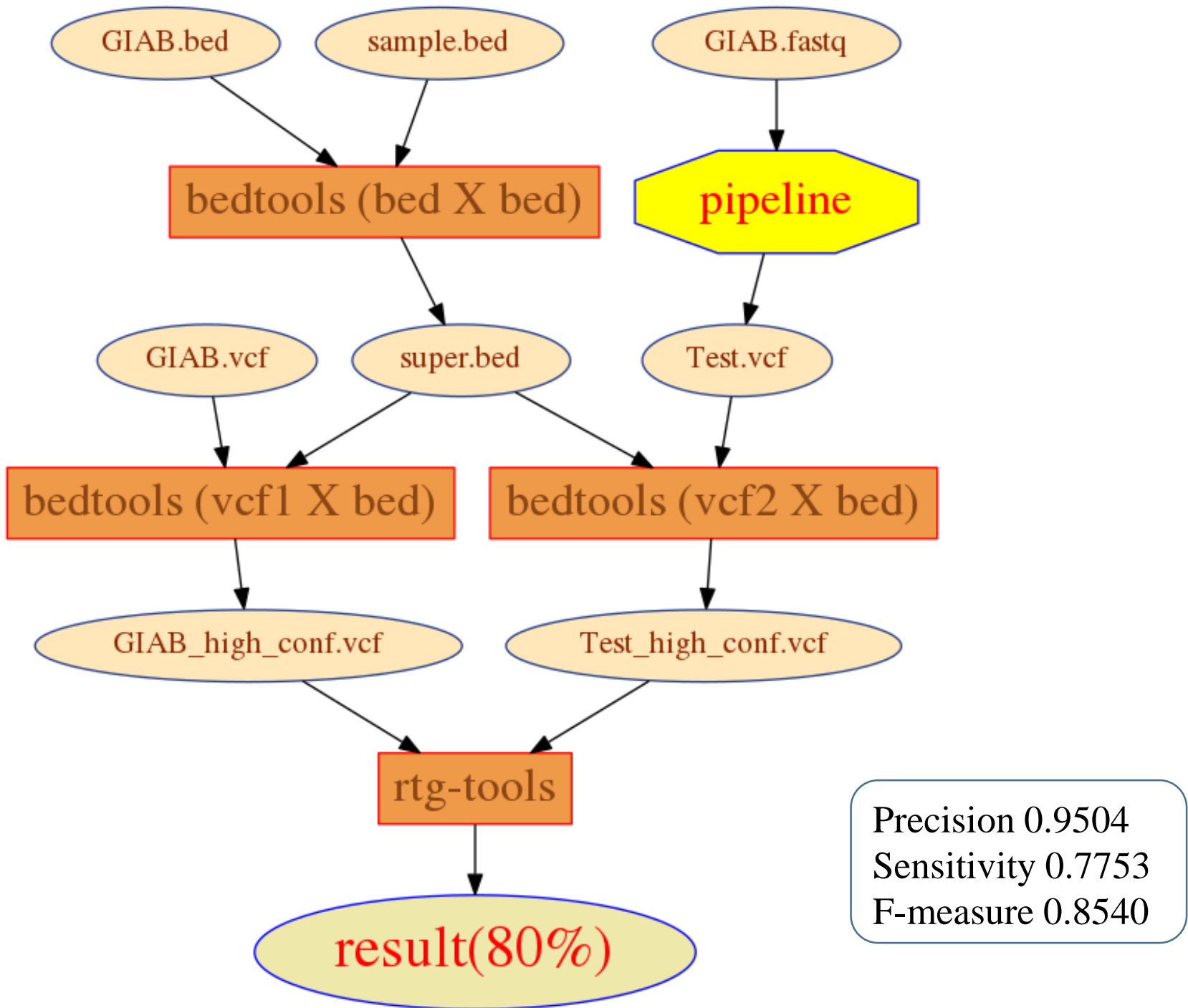
- Надо учитывать версии используемых программ;
- Надо учитывать окружение и зависимости используемых программ;
- Сложность перенесения и развертывания пайплайна на других компьютерах.

Развертывание пайплайна



Алгоритм проверки пайплайна





rs55972010
rs55847238



rs55972010

Clinical significance: [Benign/Likely benign](#)
Last evaluated: Jan 21, 2017
Number of submission(s): 12
Condition(s):

- Dilated cardiomyopathy 1G [\[MedGen - OMIM\]](#)

rs55972010
rs55847238
...

rs55847238

Clinical significance: [Conflicting interpretations of pathogenicity](#)
Benign(3);Likely benign(3);Uncertain significance(7)
Last evaluated: Jun 30, 2017
Number of submission(s): 13
Condition(s):

- Dilated cardiomyopathy 1G [\[MedGen - OMIM\]](#)



rs55972010
...



<u>rs55972010</u>	Dilated Cardiomyopathy, Dominant	2q31
<u>rs55847238</u>	Dilated Cardiomyopathy, Dominant	2q31
<u>rs13087941</u>	Hypertrophic cardiomyopathy	3p25
<u>rs505058</u>	Dilated Cardiomyopathy, Dominant	1q22
<u>rs538089</u>	Dilated Cardiomyopathy, Dominant	1q22

Низкая частота встречаемости в популяции

Высокая частота встречаемости в популяции