Домашнее задание по лекции №9

Тема: «Анализ и интерпретация NGS данных (Dna-seq) в диагностике наследственных заболеваний»

Описание: домашнее задание выполняла на сервере лаборатории МГНЦ (eod-wgs), поэтому я буду прикреплять скриншоты из команд терминала с промежуточным результатом в виде картинок. Открытые варианты можно посмотреть в <u>Google Colab.</u>

ТЕОРЕТИЧЕСКАЯ ЧАСТЬ

1. Опишите основные этапы проведения секвенирования в "мокрой" части лаборатории.

Одни из основных этапов проведения секвенирования в мокрой части включают: выделение ДНК (РНК), подготовка библиотеки и секвенирование (запуск на платформе).

Фрагментация ДНК может быть воспроизведена различными способами, глобально подразделяют но на ферментативную (использование каких-то эндонуклеаз рестрикций) или ультразвуковую После (соникация, примеру). К вносятся адаптерные последовательности, содержащие чаще всего баркоды (для того, чтобы отличать образцы друг от друга) и после подготовленные библиотеки измеряются по концентрации и отправляются на секвенирование.

Кажется очень абстрактно, но я могу описать подготовку библиотек **полногеномного бисульфитного секвенирования (WGBS)**, которыми я занималась:

Выделение геномной ДНК: фенол-хлороформная экстракция из любой доступной ткани, например периферическая кровь. На выходе измерить концентрацию ДНК (нг/мкл) и чистоту (соотношение длин волн A260/280 и A260/230) ¹

Растворить необходимое количество ДНК в MQ, чтобы получился требуемый под протокол **одинаковый объем ДНК** для библиотек на вход

Гидролиз XmaI: один из примеров ферментативной фрагментации ДНК. Получаем на выходе пул фрагментов ДНК

Затупление концов частичное фрагментом Klenow exo-, где добавление

¹ Это важно!! Потому что даже на самом первом этапе гидролиза фермент может не до конца порезать геномную ДНК, и дальнейшие этапы с дорогими реагентами просто в помойку можно (было...)

5-methyl-dCTP нужно для частичного заполнения концов ДНК-фрагментов, полученных после гидролиза XmaI. Это предотвращает самолигирование ДНК-фрагментов и адаптеров и создает липкие не палиндромные концы, которые предотвращают самозамыкание фрагментов

Лигирование фрагментов с адаптерами: аннелированный синтетический адаптер с ДНК-фрагментами

Ник-трансляция. создание одноцепочечных разрывов DNAse I, удаление нуклеотидов с 5'-разрыва и добавляются новые к 3' разрыву, используя Сте. То есть олигонуклеотиды адаптеров заменялись на последовательности с метилированными цитозинами для предотвращения деаминации цитозинов в последующем преобразовании бисульфитом.²

Очистка для инактивации полимеразы протоколом Agencourt AMPure XP

Селекция по длине (Pippin Prep), "дорожка в кассете". Здесь уже оценивается качество полученной библиотеки и отбираются оптимальной длины фрагменты. Например, Ion Torrent требует до 200 п.н. максимум длину, но и короткие (до 100 п.н.) тоже плохо.

Промывка от бромистого этидия AMPureXP

Бисульфитная конверсия "Отмывка от бисульфита с carrier RNA" Отработка циклов BS qPCR

Амплификация библиотеки, циклов

Финальная отмывка и измерение конечной концентрации библиотеки **HS** (нг/мл)

2. Какие платформы для секвенирования вы знаете? Кратко опишите не менее трех платформ, указав их основные характеристики и особенности.

Техноло гия	Чтения	Длина чтений	Скорость	Точность	Описание метода
Sanger	по одному чтению	400-900 bp	несколько часов на чтение	>99.99%	метод, основанный на дидактических цепях, с использованием флуоресцентных меченых нуклеотидов для

² Очень дорого заказывать адаптеры сразу с метилированным цитозином, поэтому дополнительный этап

					определения последовательности ДНК
Illumina/ Solexa	массово паралле льное	50-300 bp	несколько дней на геном	98-99.9%	использует флуоресцентные метки для определения нуклеотидов во время синтеза, чтения фиксируются камерой
Roche 454	массово паралле льное	400-600 bp	несколько дней на геном	99.9%	Пиросеквенирование: определение последовательности на основе высвобождения пирофосфата при инкорпорации нуклеотидов
SOLiD	массово паралле льное	50-75 bp	несколько дней на геном	99.94%	использует лиганды для чтения двух оснований одновременно, улучшая точность
PacBio SMRT	одномол екулярн ое	10,000- 15,000 bp	несколько часов на геном	87-99.99 %	использует одномолекулярное времяпролетное секвенирование (SMRT) для прямого чтения длинных цепей ДНК
Oxford Nanopor e MinION	одномол екулярн ое	До 1,000,00 0 bp	несколько часов на геном	90-98%	использует поры белка для чтения нуклеотидов по мере прохождения через мембрану.
Illumina NovaSeq	массово паралле льное	100-300 bp	Нескольк о дней на геном	99.9%	усиленная версия Illumina/Solexa с улучшенной пропускной способностью и скоростью секвенирования

Примечание: табличку составляла сама, поэтому ссылку не привела

3. Опишите стандартный биоинформатический пайплайн для анализа данных NGS (DNA-seq). Кратко укажите цель каждого из них.

- **1. Контроль качества исходных данных (FastQ).** Задача оценить качество сырых данных (FastQ файлы), включая качество чтений, уровень ошибок, наличие адаптеров и других артефактов. *Инструменты: FastQC, MultiQC*
- **2.** Удаление адаптеров и обрезка низкокачественных областей. Удалить адаптеры и обрезать низкокачественные концы чтений, чтобы повысить точность последующих этапов анализа.

Инструменты: Cutadapt, Fastp

3. Выравнивание (Alignment). Выравнивание чтений на референсный геном для идентификации местоположений, где были сделаны изменения (варианты)

Инструменты: BWA, Bowtie2, STAR (для RNA-seq), Minimap2

- **4. Обработка выровненных данных.** Например, удаление дублирующих чтений и улучшение качества выравнивания. *Инструменты: MarkDuplicates, Samtools*
- **5. Вызов вариантов (Variant Calling).** Идентификация однонуклеотидных полиморфизмов (SNPs) и инделов из выровненных данных.

Инструменты: GATK HaplotypeCaller, Samtools mpileup, FreeBayes

6. Аннотация вариантов. Обогащение информации о вариантах, добавление функциональной информации (например, какие гены или экзоны затронуты) и сопоставление с известными базами данных, такими как ClinVar, dbSNP.

Инструменты: SnpEff, ANNOVAR

7. Фильтрация вариантов. Отбор надежных и потенциально значимых вариантов, исключая низкое качество, низкую глубину покрытия

Инструменты: GATK VariantFiltration, bcftools, vcfutils.

8. Качественная оценка и визуализация данных. Оценка качества вызванных вариантов и визуализация для дальнейшей интерпретации.

Инструменты: IGV (Integrative Genomics Viewer), UCSC Genome Browser. bedtools.

- 9. Статистический и функциональный анализ. Инструменты: PLINK, SNPTEST, GWAS tools, R/Bioconductor
- 4. Кратко опишите основные этапы пайплайна GATK для анализа наследуемых коротких инделов (indels) и однонуклеотидных полиморфизмов (SNPs)

GATK (Genome Analysis Toolkit) используется для анализа наследуемых коротких инделов (indels) и однонуклеотидных полиморфизмов (SNPs) и включает в себя следующие шаги:

- **1. Качество и фильтрация исходных данных (FastQ):** выполняется предварительная фильтрация сырых данных (например, с FastQC или тримминг адаптеров). Если много образцов еще multiQC.
- **2. Выравнивание с использованием BWA.** FastQ файлы выравниваются на референсный геном с помощью программы BWA или Bowtie2. Результат: BAM файл с выровненными прочтениями
- **3. Обработка выравнивания:** Mark Duplicates удаляются дублирующиеся чтения, Indel Realigner, чтобы исправить возможные ошибки выравнивания, и можно выполнить перекалибровку качества баз с использованием Base Recalibrator.
- **4. Вызов вариантов (SNPs и инделов)** с использованием HaplotypeCaller или DeepVar из выровненных ВАМ файлов после удаления дубликатов. Выход vcf файл.
- **5. Варианты аннотируются** с использованием инструментов, таких как SnpEff или ANNOVAR.
- **6. Фильтрация вариантов по качеству** (Variant Filtration, например, качество, глубина покрытия, флаг PASS)
- 7. Анализ и интерпретация данных. Визуализация результатов с помощью таких инструментов, как IGV (Integrative Genomics Viewer)

ПРАКТИЧЕСКАЯ ЧАСТЬ

Прежде чем приступить к работе с ВАМ файлом, мне требовалось посмотреть на сам файл: отображается ли основная информация и целый файл, так как предыдущие этапы подготовки файла (контроль качества, тримминг, выравнивание) были выполнены не мной. Это можно сделать при помощи команды samtools view с выводом первых 4 строк:

samtools view sample_15.bam | head -n 4

Рисунок 1. Содержимое ВАМ-файла

Из Рисунка 1 можно увидеть основные аспекты выравнивания, которые пишутся в бамках, где каждая строка представляет прочтение (read) с его атрибутами. Основные поля включают:

QNAME: идентификатор прочтения (например, SRR22359471.36921505)

FLAG: флаги выравнивания, дают информацию о выравнивании и парности прочтений (99, 147, 83, 163 и тд)

RNAME: референсное имя хромосомы (например, chr1)

POS: позиция выравнивания на референсе

MAPQ: Качество выравнивания (например, 41). Более высокие значения указывают на лучшее качество выравнивания

CIGAR: Строка CIGAR, описывающая операции выравнивания, такие как совпадение (М), вставка (I), удаление (D), мягкое отсечение (S) и т.д.

RNEXT: Референсное имя для следующего прочтения в паре

PNEXT: Позиция следующего прочтения в паре

TLEN: Шаблонная длина

SEQ: Последовательность нуклеотидов прочтения

QUAL: Качество нуклеотидов прочтения

После такой инфы можно оценить качество выравниваний, например по MAPQ и CIGAR. После маркирования дубликатов с помощью MarkDuplicates, можно анализировать метрики дубликатов, искать аномальные выравнивания, такие как инверсии, транслокации и тд.

ЧАСТЬ ПЕРВАЯ. Маркирование дубликатов в bam-файле с помощью MarkDuplicates из GATK

- 1) Прежде чем маркировать дубликаты, я установила GATK при помощи команды wget
 - wget https://github.com/broadinstitute/gatk/releases/download/4.6.1.0/gatk-4.6.1.0.zip
- 2) Команда распаковки архива: unzip gatk-4.6.1.0.zip
- 3) Для маркирования дубликатов я использовала команду с использованием функции **JavaScript**:

```
java -jar picard.jar MarkDuplicates \
    I=sample_15.bam \
    O=sample_15_marked.bam \
    M=sample_15_metrics.txt \
    CREATE INDEX=true
```

```
ast 1,000,000: 12s. Last read position: chr4:1,974,707
INFO 2024-12-02 11:35:12 MarkDuplicates Read 19,000,000 records. Elapsed time: 00:04:09s. Time for ast 1,000,000: 12s. Last read position: chr4:57,031,303
INFO 2024-12-02 11:35:37 MarkDuplicates Read 20,000,000 records. Elapsed time: 00:04:22s. Time for ast 1,000,000: 12s. Last read position: chr4:122,271,013
INFO 2024-12-02 11:35:37 MarkDuplicates Read 20,000,000 records. Elapsed time: 00:04:22s. Time for ast 1,000,000: 12s. Last read position: chr4:122,271,013
INFO 2024-12-02 11:35:50 MarkDuplicates Read 21,000,000 records. Elapsed time: 00:04:35s. Time for ast 1,000,000: 12s. Last read position: chr5:9,252,380
INFO 2024-12-02 11:36:02 MarkDuplicates Read 22,000,000 records. Elapsed time: 00:04:48s. Time for ast 1,000,000: 12s. Last read position: chr5:179,799,526
INFO 2024-12-02 11:36:02 MarkDuplicates Read 22,000,000 records. Elapsed time: 00:04:48s. Time for ast 1,000,000: 12s. Last read position: chr5:179,799,526
INFO 2024-12-02 11:36:15 MarkDuplicates Read 22,000,000 records. Elapsed time: 00:05:00s. Time for ast 1,000,000: 12s. Last read position: chr5:177,388,198
INFO 2024-12-02 11:36:28 MarkDuplicates Read 24,000,000 records. Elapsed time: 00:05:00s. Time for ast 1,000,000: 12s. Last read position: chr5:177,388,198
INFO 2024-12-02 11:36:41 MarkDuplicates Read 25,000,000 records. Elapsed time: 00:05:26s. Time for ast 1,000,000: 12s. Last read position: chr5:177,388,198
INFO 2024-12-02 11:36:41 MarkDuplicates Read 25,000,000 records. Elapsed time: 00:05:26s. Time for ast 1,000,000: 12s. Last read position: chr6:35,462,605
INFO 2024-12-02 11:36:44 MarkDuplicates Read 26,000,000 records. Elapsed time: 00:05:39s. Time for ast 1,000,000: 13s. Last read position: chr6:35,462,605
INFO 2024-12-02 11:36:54 MarkDuplicates Read 26,000,000 records. Elapsed time: 00:05:25s. Time for ast 1,000,000: 13s. Last read position: chr6:79,960,687
INFO 2024-12-02 11:36:54 MarkDuplicates Read 27,000,000 records. Elapsed time: 00:05:52s. Time for 2024-12-02 11:36:54 MarkDu
```

Рисунок 2. Выполнение команды MarkDuplicates

Ha выходе получила новый файл sample_MarkDuplicate.bam, который я также просмотрела:

```
ls view sample_MarkDuplicates.bam | head -n 4

SRR22359471.36921505 99 chr1 10006 41 92M1D45M13S = 10146 171 c

TAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACGCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAACCTAAC
```

Рисунок 3. Результат маркирования дупликатов, видно новый определитель *PG:MarkDuplicates*

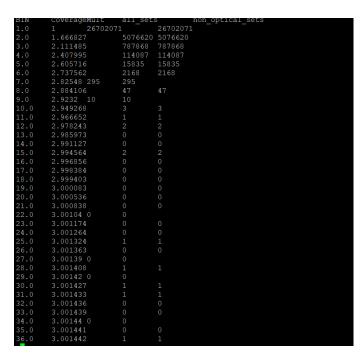
После я просмотрела отдельно созданный файл с метриками sample_MD_dedup.metrics.txt, который содержит обычно информацию о параметрах, с которыми была выполнена команда MarkDuplicates из GATK и метрики, собранные в ходе процесса маркирования дубликатов. Вот что получила на выходе:

Таблица метрик

```
## METRICS CLASS picard.sam.DuplicationMetrics
LIBRARY UNPAIRED READS EXAMINED RAD PAIRS EXAMINED
PERCENT_DUPLICATION ESTIMATED LIBRARY SIZE
UNKNOWN Library 51902 39770473 74177 59938 24251 7071456 0 0.177995 98144242
```

- unpaired_reads_examined: одиночных прочтений: 51,902.
- read_pairs_examined: пар прочтений: 39,770,473.
- **secondary_or_supplementary_rds**: вторичные или дополнительные прочтения: 74,177.
- unmapped_reads: невыравненные прочтения: 59,938.

- unpaired_read_duplicates: дубликаты среди одиночных прочтений: 24,251.
- read_pair_duplicates: дубликаты среди пар прочтений: 7,071,456. Это может быть следствием технических артефактов или особенностей подготовки библиотеки.
- read_pair_optical_duplicates: оптические дубликаты среди пар прочтений: 0. Отсутствие оптических дубликатов указывает на то, что проблемы с дублированием связаны не с техническими ошибками считывания, а, возможно, с этапом подготовки библиотеки или характеристиками исходного материала.
- percent_duplication: процент дублирования: 0.177995 или 17.8%. Процент дублирования 17.8%, значит почти пятая часть данных состоит из дубликатов.
- estimated_library_size: оценочный размер библиотеки: 98,144,242, вообще достаточный объем данных для проведения дальнейших анализов.



Гистограмма покрытий (Coverage Histogram)

ЧАСТЬ ВТОРАЯ. Коллинг вариантов с помощью DeepVariant в режиме WES

```
sudo apt -y update
sudo apt-get -y install docker.io
sudo docker pull google/deepvariant:"1.6.1"
```

Для проведения коллинга воспользовалась инструкцией, приведенной по <u>ссылке</u>. Для этого требовалось предварительно установить <u>Docker</u>³.

Для его установки воспользовалась следующей инструкцией:

1) Установление зависимостей:

```
sudo apt update \
sudo apt install apt-transport-https ca-certificates curl
software-properties-common
```

2) Добавила GPG ключ для официального Docker репозитория:

```
curl -fsSL https://download.docker.com/linux/ubuntu/gpg |sudo apt-key add -
```

3) Docker репозиторий в APT источники:

```
sudo add-apt-repository "deb [arch=amd64]
https://download.docker.com/linux/ubuntu $(lsb_release -cs) stable"
```

- 4) Установка Docker: sudo apt install docker-ce
- 5) Добавила в своего пользователя в группу Docker, чтобы запускать команды Docker без sudo: sudo usermod -aG docker aimanalieva

Теперь, когда Docker установлен, я провела **DeepVariant** для анализа экзомки (WES). Для этого скрипт описан ниже:

```
mkdir -p output
mkdir -p output/intermediate_results_dir
BIN_VERSION="1.6.1"
sudo docker run \
 -v "${PWD}/input":"/media/HEAP-EPI/Diplome_gi_Imanalieyva_Amina/NGS/input" \ \
 -v "${PWD}/output":"/media/HEAP-EPI/Diplome_gi_Imanalieyva_Amina/NGS/deepvariant/output" \
 -v "${PWD}/reference":"/media/HEAP-EPI/Diplome_gi_Imanalieyva_Amina/NGS/input/reference" \
  google/deepvariant:"${BIN_VERSION}" \
  /opt/deepvariant/bin/run_deepvariant \
  --model_type WES \
  --ref /reference/GRCh38_no_alt_analysis_set.fasta \
  --reads /input/sample_15_marked.bam \
  --regions /input/idt_capture_novogene.grch38.bed \
  --output_vcf /output/sample_15_marked.vcf.gz \
  --output_gvcf /output/sample_15_marked.g.vcf.gz \
  --num_shards 4
```

Примечание: выражаю большую благодарность моему коллеге Тимуру Кулагину за помощь и консультацию, так как на сервере лабы не получалось настроить Tensor Flow для запуска Docker.

³ Docker — платформа для контейнеризации, позволяющая запускать приложения в изолированных, переносимых средах.

ЧАСТЬ ТРЕТЬЯ. Фильтрация VCF файла.

Флаг 'PASS' означает, что данный вариант успешно прошел все фильтры и считается надежным для дальнейшего анализа. Например, это помогает отобрать варианты, прошедшие дополнительные или специфические фильтрации (например, генные варианты с низким качеством или низким покрытием).

Input file: sample_15_marked
Output file:output.vcf

Команда для фильтрации:

```
awk -F '\t' '\{if(\$0 \sim /\#/) \text{ print}; \text{ else } if(\$7 == "PASS") \text{ print}' \text{ output/sample}\_15\_marked.vcf > output.vcf}
```

(base) etcetera_mi@AMINA-16100	2:~\$ tail	output.	vcf		
chrX 154652556 .		Α	69	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:67:204:1,2	03:0.9950	98:69,72	, 0		
chrX 154653251 .		G .	69	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:65:89:0,89		0			
chrX 154653499 .		T	32.8	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:16:6:0,6:1					
chrX 154680220 .			56	PASS	GT:G
Q:DP:AD:VAF:PL 0/1:46:154:65,					
chrX 154766321 .			70.2	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:69:252:0,2					
chrX 154791839 .			50	PASS	GT:G
Q:DP:AD:VAF:PL 0/1:50:179:102	•		•	5466	CT C
chrX 154929926 .		G 1 F0 0 F1	50.1	PASS	GT:G
Q:DP:AD:VAF:PL 0/1:49:76:43,3				DAGG	CT C
chrY 14840455 .		С	28.2	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:24:9:0,9:1		6	2/1 0	DACC	CT . C
chrY 14840467 .			24.8	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:23:10:0,10 chrY 14840785 .		T	22.9	PASS	CT+C
		6	22.9	PASS	GT:G
Q:DP:AD:VAF:PL 1/1:16:29:0,29	.1.22,16,	U			

Рисунок 4. Отфильтрованный файл по 'PASS'

ЧАСТЬ ЧЕТВЕРТАЯ. Аннотация с помощью SnpEff.

Аннотация с помощью SnpEff помогает интерпретировать результаты, полученные из данных о генетических вариантах. Он предоставляет тип варианта (например, замена нуклеотида, вставка, делеция), frame-shift, сдвиг стоп-кодона или какие-то миссенс-варианты. Часто использует аннотирование вариантов с использованием с использованием базы данных по типу dbSNP, ClinVar, Ensembl. Он также может давать информацию о возможных эффектах вариантов на регуляцию генов. Это может включать: изменения в промоторах, энхансерах или других регуляторных элементах, влияние на сплайсинг РНК.

Скрип приведен ниже:

4.1 Добавить аннотацию SnpEff

1) Скачать SnpEff

```
wget https://snpeff.blob.core.windows.net/versions/snpEff latest_core.zip
unzip snpEff_latest_core.zip
```

2) Скачать базу данных hg38 для SnpEff

```
java -jar snpEff/snpEff.jar download -v hg38
```

3) Аннотация файла с флагом -canon на hg38

java -jar snpEff/snpEff.jar -canon hg38 output.vcf > annotated_output.vcf

```
tail output.vcf
                                                         PASS
                                                                         GT:G
chrX
        154652556
Q:DP:AD:VAF:PL 1/1:67:204:1,203:0.995098:69,72,0
                                                         PASS
                                                                         GT:G
chrX
        154653251
Q:DP:AD:VAF:PL 1/1:65:89:0,89:1:68,67,0
        154653499
                                                                         GT:G
chrX
                                                 32.8
                                                         PASS
Q:DP:AD:VAF:PL 1/1:16:6:0,6:1:32,15,0
                                                 56
                                                                         GT:G
                                                         PASS
chrX
        154680220
Q:DP:AD:VAF:PL 0/1:46:154:65,87:0.564935:56,0,46
                                                 70.2
                                                         PASS
                                                                         GT:G
        154766321
Q:DP:AD:VAF:PL 1/1:69:252:0,252:1:70,76,0
                                                                         GT:G
        154791839
                                                 50
                                                         PASS
chrX
Q:DP:AD:VAF:PL 0/1:50:179:102,77:0.430168:50,0,63
        154929926
                                                 50.1
                                                                         GT:G
chrX
                                                         PASS
Q:DP:AD:VAF:PL 0/1:49:76:43,33:0.434211:50,0,57
chrY
        14840455
                                                 28.2
                                                         PASS
                                                                         GT:G
Q:DP:AD:VAF:PL 1/1:24:9:0,9:1:28,26,0
                                                 24.8
chrY
        14840467
                                                         PASS
                                                                         GT:G
Q:DP:AD:VAF:PL 1/1:23:10:0,10:1:24,28,0
        14840785
                                                 22.9
                                                         PASS
                                                                         GT:G
chrY
Q:DP:AD:VAF:PL 1/1:16:29:0,29:1:22,16,0
(base) etcetera_mi@AMINA-161002:~$ tail annotated_output.vcf | head -n 3
                                                                 ANN=A|synony
        154652556
                                                 69.0
                                                         PASS
mous_variant|LOW|CTAG2|CTAG2|transcript|NM_020994.5|protein_coding|2/2|c.345
G>T|p.Pro115Pro|409/988|345/633|115/210||
                                                 GT:GQ:DP:AD:VAF:PL
67:204:1,203:0.995098:69,72,0
                                                                 ANN=G|missen
        154653251
                                         G
                                                 69.0
                                                         PASS
se_variant|MODERATE|CTAG2|CTAG2|transcript|NM_020994.5|protein_coding|1/2|c.
265G>C|p.Glu89Gln|329/988|265/633|89/210||
                                                 GT:GQ:DP:AD:VAF:PL
65:89:0,89:1:68,67,0
        154653499
                                                 32.8
                                                         PASS
                                                                 ANN=T|missen
se_variant|MODERATE|CTAG2|CTAG2|transcript|NM_020994.5|protein_coding|1/2|c.
17G>A|p.Arg6Gln|81/988|17/633|6/210||
                                        GT:GQ:DP:AD:VAF:PL
                                                                 1/1:16:6:0,6
:1:32,15,0
```

Вот для сравнения я привела на картинке, как отличается вывод информации до добавления аннотации (output.vcf) и после добавления аннотации (annotated_output.vcf).

4.2 Добавить информацию из ClinVar

Скачать базу данных ClinVar с индексированным файлом:

wget https://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf_GRCh38/clinvar.vcf.gz
wget https://ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf_GRCh38/clinvar.vcf.gz.tbi

Скрипт для аннотации на клинвар:

java -Xmx1g -jar /snpEff/SnpSift.jar annotate -v clinvar.vcf annotated_output.vcf > clinvar_annotation.vcf

ЧАСТЬ ПЯТАЯ. Фильтрация патогенных и вероятно патогенных вариантов Отфильтруйте Pathogenic и Likely pathogenic варианты.

- **Патогенные варианты:** варианты, которые могут быть связаны с заболеваниями
- **Рискованные варианты:** варианты, которые могут увеличивать вероятность возникновения заболевания, но еще не полностью подтверждены

После аннотации на патогенные и вероятно-патогенные варианты полученного **vcf.файла** получила **17 вариантов**:



Патогенных вариантов получилось всего 9, при этом, очень интересный вариант с делецией (всего один)

ЧАСТЬ ШЕСТАЯ. Генетические варианты, ассоциированные с фенотипом

Из проведенного анализа и полученных 17 вариантах я решила остановиться на некоторых патогенных вариантах (из девяти я выбрала после анализа 6, другие не оч были показаны в ИБС) и одном потенциально-патогенным, делеция которого приводит к изменению рамки считывания.

Likely Pathogenic

ITPKB (chr1, позиция 226737174)

- disruptive inframe deletion c.276_284delCAGCGGCAG
- **белковая мутация:** p.Ser93_Ser95del потеря аминокислот серина в позиции 93 и 95 в белке

Ген *ITPKB* (инозитол-трифосфаткиназа В) участвует в **клеточном** метаболизме и может быть связан с регуляцией клеточного сигнала.

Тип вариации: Microsatellite, что указывает на возможную микросателлитную нестабильность, которая может быть связана с различными заболеваниями, включая рак и сердечно-сосудистые заболевания.

CLNDISDB: Связь с несколькими медицинскими идентификаторами и заболеваниями, такими как миелопролиферативная неоплазия, что может указывать на редкие заболевания с возможной сердечно-сосудистой компонентой. *Частота встречаемости варианта VAF (Variant Allele Frequency)* 0.6304.

Показано на NCBI, что вариант в гене может повлиять на метаболизм клеток и процессы передачи сигналов, что в свою очередь может повлиять на развитие болезней, таких как ишемическая болезнь сердца, особенно если этот вариант влияет на регуляцию клеточного ответа в сосудистой ткани или сердечной мышце.

Pathogenic

PROC (chr 2, позиция 127426090)

- $missense_variant (T > G)$
- c.541T>G изменение аминокислоты в белке: p.Phe181Val (замена фенилаланина на валин в позиции 181)

Ген **PROC** (протеин С) участвует в антикоагулянтной системе организма и играет важную роль в регуляции свертывания крови. **CLNDISDB**: вариант ассоциирован с различными заболеваниями: thrombophilia due to protein C deficiency, autosomal dominant тромбофилия, вызванная дефицитом протеина С, в доминантной форме, thrombophilia due to protein C deficiency, autosomal recessive тромбофилия, вызванная дефицитом протеина С, в рецессивной форме, OMIM:176860 и Orphanet:745 — заболевания, связанные с нарушениями свертывания крови и тромбообразованием. В аннотации также указано upstream gene variant MIR4783, наличие ДЛЯ гена что свидетельствовать о влиянии на регуляцию гена, расположенного выше по течению, но это не оказывает значительного воздействия на основной вариант в PROC. Недавние исследования показывают, что дефицит протеина С может увеличивать склонность к образованию тромбов, что, в свою очередь, повышает риск ишемической болезни сердца, инсульта и других сосудистых заболеваний. **Миссенс-мутация**, как в этом варианте, может привести к аномальному функционированию протеина С, что нарушает его антикоагулянтные свойства, увеличивая вероятность тромбообразования.

ALDH5A1 (chr 6, позиция 24503362)

- missense_variant: C > T
- мутация с.538C>Т приводит к замене аминокислоты гистидина (His) на тирозин (Tyr) в позиции 180 (р.His180Tyr)

ALDH5A1 ЭТОТ ген кодирует белок, называемый 5A1. альдегиддегидрогеназа который участвует метаболизме В гамма-аминомасляной кислоты (ГАМК) и играет важную роль в детоксикации токсичных метаболитов. CLNDISDB: вариант ассоциирован заболеваниями, включая дефицит несколькими метаболизме сукцинат-семальдегидрогеназы, нарушениями В расстройствами. нейрологическими Хотя прямое связанное c ишемической болезнью сердца (ИБС) не установлено, нарушения в метаболизме, связанные с *ALDH5A1*, могут оказывать влияние на систему, участвующую в воспалении и регуляции сосудистого тонуса, что в свою очередь может быть связано с развитием заболеваний, таких как ИБС. Мутации в генах, таких как *ALDH5A1*, могут нарушать сосудистую функцию, увеличивая риск сердечно-сосудистых заболеваний, особенно если мутация влияет на клеточные процессы, связанные с воспалением или метаболизмом липидов.

CYP21A2 (chr 6, позиция 32038610)

- missense_variant: A > T
- замена аминокислоты гистидина (His) на лейцин (Leu) в позиции 63 (p.His63Leu)

СУР21А2 — этот ген кодирует 21-гидроксилазу, фермент, который участвует в синтезе стероидных гормонов в надпочечниках, таких как кортизол и альдостерон. Недостаток этого фермента приводит к нарушению гормонального баланса и может вызвать конгенитальную гиперплазию надпочечников. Прямое отношение между СУР21А2 и ишемической болезнью сердца (ИБС) не установлено, так как этот ген в основном ассоциирован с эндокринными нарушениями и гиперплазией надпочечников. Однако, изменения в метаболизме стероидов и воспаления,

связанные с дефицитом 21-гидроксилазы, могут косвенно влиять на сердечно-сосудистую систему. Например, стресс и хроническое воспаление, характерные для дефицита гормонов, могут повышать риски развития сердечно-сосудистых заболеваний, включая ИБС.

NQO1 (chr 16, позиция 69711242)

- missense_variant: G > A
- замена аминокислоты **пролина (Pro)** на **серин (Ser)** в позиции 187 (**p.Pro187Ser**)

кодирует фермент NAD(P)H:quinone oxidoreductase 1, который участвует в детоксикации организма, защищая клетки от повреждений, вызванных окислительным стрессом. Он также играет роль в метаболизме различных химических веществ, включая токсины и Conflicting канцерогены. Вариант имеет статус classifications pathogenicity (противоречивые классификации патогенности), указывает на различные мнения о его патогенности. Вариант описан, как связанный с дефицитом или дисфункцией NQO1, увеличением риска развития рака легких. повышенной восприимчивостью к токсичности бензола, ухудшением прогноза для выживаемости при раке молочной железы после химиотерапии. В случае NQO1 может нарушиться его способность детоксифицировать клетки. Прямое отношение между NQO1 и ишемической болезнью сердца (ИБС) не установлено, но изменение активности NQO1 может повлиять на уровень окислительного стресса в клетках и тканях. Окислительный стресс связан с воспалением и повреждением сосудов, что может быть фактором риска для ИБС. Также, нарушение метаболизма токсинов и канцерогенов может повысить риск сердечно-сосудистых заболеваний через воздействия на организм в целом.

RNF213 (chr 17, позиция 80385145)

- missense variant: G > A
- intron_variant: для гена RNF213-AS1
- замена аминокислоты **аргинина (Arg) на лизин (Lys)** в позиции 4810 (p.Arg4810Lys)

Связан с **болезнью Мойя-Мойя** (Moyamoya disease), кодирует белок, который участвует в различных клеточных процессах, включая регулирование клеточного цикла и поддержание сосудистого здоровья. Белок *RNF213* играет важную роль в поддержании сосудистой

целостности и функциональности, и его нарушение может быть связано с сосудистыми расстройствами, которые могут также повлиять на сердце.

FECH (chr 18, позиция 57571588)

- intron variant

FECH кодирует ключевой фермент в пути биосинтеза гемоглобина, который катализирует последнюю стадию синтеза гема. с.333-48Т>С находится в интронной области, что может повлиять на правильное сплайсирование РНК или нарушить регуляцию экспрессии гена, что может привести к дефициту фермента и нарушениям в метаболизме порфиринов. Вариант ассоциирован с эритропоэтической протопорфирией: накоплением протопорфирина в эритроцитах и других метаболизме порфиринов, нарушения В протопорфирии, эритропоэтической ΜΟΓΥΤ косвенно повлиять сосудистую систему, поскольку накопление порфиринов может привести к повреждению сосудов. В редких случаях такие нарушения могут быть заболеваниями, влияющими на кровоснабжение, связаны теоретически может повысить риск ишемической болезни сердца.

APOE (chr 19, позиция 44908684)

- missense variant: T > C
- замена цистеин (Суѕ) на аргинин (Агд) на позиции 156

АРОЕ участвует в метаболизме липидов и холестерина, а также в функционировании мозга. Он играет ключевую роль в транспортировке холестерина В нейронах. Вариант ассоциирован заболеваниями: болезнь Альцгеймера), ответ на антикоагулянтное лечение (варфарин), нарушения обмена липидов и холестерина. Прямое отношение между АРОЕ и ишемической болезнью сердца (ИБС) в контексте данного варианта не установлено. Однако известно, что АРОЕ играет ключевую роль в метаболизме липидов, и определенные аллели АРОЕ могут увеличивать риск атеросклероза и, как следствие, ИБС, особенно АРОЕ вариант p.Cys156Arg может потенциально функциональные характеристики белка АРОЕ, что влияет на липидный обмен и может быть связано с риском атеросклероза и кардиоваскулярных заболеваний.