

# Consequence

Сергей Лебедев

Руководитель: Николай Вяххи

<http://github.com/bioinf/consequence>

4 июня 2012 г.

- aka **S**ingle **N**ucleotide **P**olymorphism,
- точечное отличие в последовательности ДНК представителей одного вида.
- rs334<sup>1</sup>, HBB (haemoglobin beta) gene, position 70614:  
TGACTCCTG**A**GGAGAAGTCT => TGACTCCTG**A**GGAGAAGTCT
- Выявлено значительное количество SNP, влияющих на вероятность развития отдельных заболеваний; примеры: *серповидно-клеточная анемия, глаукома, болезнь Альцгеймера, диабет* – см. [SNPedia](#);
- Крупные проекты по исследованию SNP и др. вариаций: [1000 Genomes](#), [23andMe](#), [deCODE](#).

---

<sup>1</sup>[http://www.ncbi.nlm.nih.gov/SNP/snp\\_ref.cgi?rs=334](http://www.ncbi.nlm.nih.gov/SNP/snp_ref.cgi?rs=334)

- Общая задача: найти в данных ридов все SNP, по сравнению с референсным геномом; когда референс отсутствует – в двух наборах ридов;
- Для человека референс имеется, поэтому большинство существующих решений для поиска SNP *de-novo* предполагают на входе риды, выравненные на референсный геном, а затем максимизируют  $Pr(\textit{genotype}|\textit{reads})$ , примеры: [samtools](#), [GATK](#), [SOAP2](#);
- Часто, требуется проверить наличие уже *известных* SNP, а не искать новые, например для того, чтобы определить предполагаемый фенотип по ридам.

## Общеиспользуемый вариант

reads -> BAM -> sorted BAM -> raw VCF -> SNPs

reads -> BAM -> SNPs

- Выравниваем риды на **частичный** референсный геном, из которого исключены участки, не содержащие ни одного SNP из **индекса**.
- После выравнивания, для каждого рида известно его местоположение в *частичном* референсном геноме, из которого можно получить соотв. позицию в полном геноме;
- а затем поискать в *индексе* все базы, не совпавшие с референсом.
- На выходе получаем список SNP с частотой встречаемости ( $MAPQ \cdot occ$ ), для каждого SNP в индексе хранится соотв. множество геномов →

- Индекс: *Genomic position* → 4 tuple;
- Индекс поддерживается для каждой хромосомы в отдельности;
- Обновление индекса можно производить в фоне, параллельно для всех новых геномов, которые необходимо добавить.

Пример:

```

      A   C       G           T
      |   |       |           |
1049631 -> ([], NULL, ["HG00096"], [])
              ^
              |
reference base
```

## Результаты: Время работы

	<i>E.coli</i>		<i>HG00096</i>	
	samtools	consequence	samtools	consequence
Reference indexing	<1m	<1m	191m	69m
Read mapping	20m	16m	42m	31m
BAM sorting	13m	—	30m	—
SNP calling	15m	3m	103m	61m
BCF → VCF	<1m	<1m	12m	—
	48m	19m	378m	161m

Таблица: Сравнение<sup>2</sup> времени работы на ридах для *E.coli* (?) и человека (1000 Genomes, HG00096)

---

<sup>2</sup>Для выравнивания ридов в обоих случаях использовался [Bowtie](#)

	Общие	samtools	consequence	Чувствительность
<i>E.coli</i>	23	139	19	14.2%
<i>HG00096</i> <sup>3</sup>	10987	860454 <sup>4</sup>	3962	1.2% / 75.8%

Таблица: Количество SNP, найденных каждым из методов

## Sensitivity and specificity

$$\text{sensitivity} = \frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$$

$$\text{specificity} = \frac{\text{number of true negatives}}{\text{number of true negatives} + \text{number of false positives}}$$

---

<sup>3</sup>Индекс содержал данные openSNP для пользователя [Mark Davis](#)

<sup>4</sup>Реальное количество SNP, присутствующих в индексе – 14488

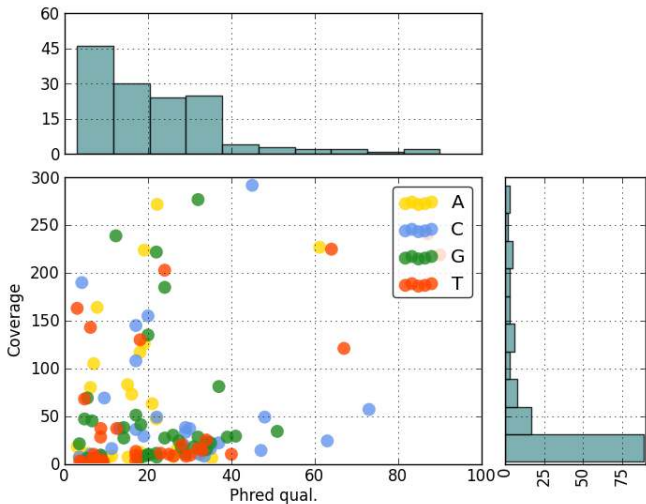


Рис.: Точечный график покрытия **позиций**, в которых не были найдены SNP и **качества** этих SNP, в выводе samtools



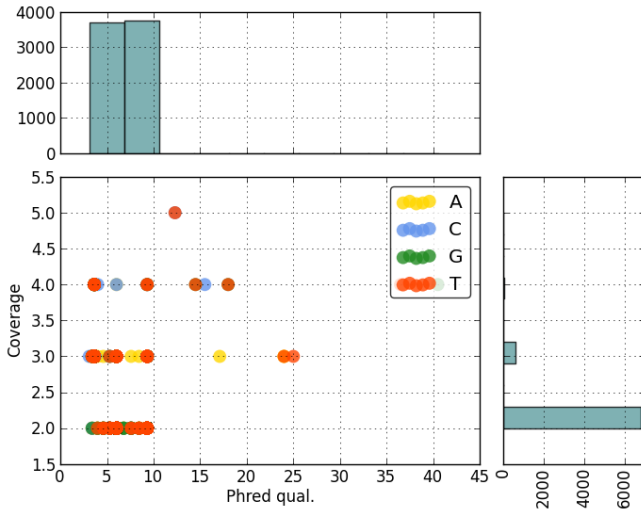


Рис.: Точечный график покрытия **позиций**, в которых не были найдены SNP и **качества** этих SNP, в выводе samtools

- Текущая реализация не обладает достаточной чувствительностью для поиска всех SNP в исходных ридов, но даже такой чувствительности достаточно для поиска “родственных” геномов и определения фенотипа.
- Количество false positives свидетельствует о недостаточной точности фильтрации – хотелось бы получить меньше, но достоверных SNP.

## Планы

- Улучшить точность фильтрации найденных SNP, например добавив к итоговому score для каждого SNP phred. quality из соотв. ридов.
- Реализовать поддержку индексации форматов, доступных на [openSNP](#) – phenotype, finally!

Спасибо за внимание!