

# ЛЕКЦИЯ:

# Медицинская биоинформатика - настоящее и тенденции

Часть 1: Геномные данные - в поисках клинического смысла.

Павлов Александр

Руководитель проекта

Часть 2: Практические аспекты разработки ПО

Брагин Антон

Руководитель группы биоинформатики

Москва 1 августа 2013 г.

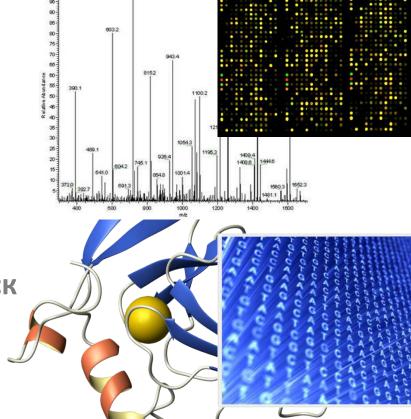
# Биоинформатика в медицине

### Области

- Масс-спектрометрия
- Анализ данных microarray
- Анализ данных секвенирования
- Биомоделирование
- Базы данных и эффективный поиск
- Анализ изображений

## Предпосылки

- Рост объемов генерируемых данных
- Потребность в комплексной аналитике data mining
- Виртуализация wet lab
- Data-driven research







## Роль геномики в медицине

### Области

- Наследственные заболевания
- Онкодиагностика
- Персонализированная медицина
- Пренатальная диагностика
- Medical research

### Задачи

- Поиск генетических вариантов
- Установление ассоциаций (генотип фенотип)
- Пациент, как «омный» профиль

## Препятствия на пути

- Гетерогенность входных данных
- Критерии качества результатов
- Проблема стандартов (форматы, описания фенотипов, гайдлайны)
- Базы данных и их валидация

- Трудность в формулировании задач биологами
- Трансляция результатов
- Уровень достоверности
- Культура работа с данными (scientific data management)

# Что можно сделать?

### Уровень техники

- Введение стандартов для входных данных и метрики их качества
- Обеспечение прослеживаемости путей миграции данных
- Создание набора инструментов (инфраструктура + ПО) для разных задач и пользователей

## Уровень концепции

- Стандарты для трансляционных исследований
- Введение принципов «информационной гигиены» при работе с клиническими данными
- Выработка требований к валидации процессов и результатов

# Эволюция представлений



### октябрь 2007

If you would like to purchase a private copy of your personal genome, please <u>contact us.</u>



### декабрь 2007

Knome Launches First Commercial Whole-Genome Sequencing and Analysis Service for Individuals

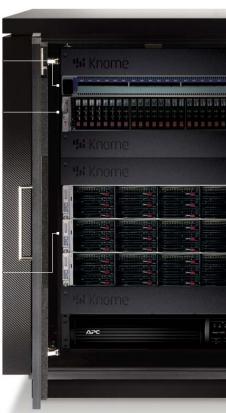


13.07.2011









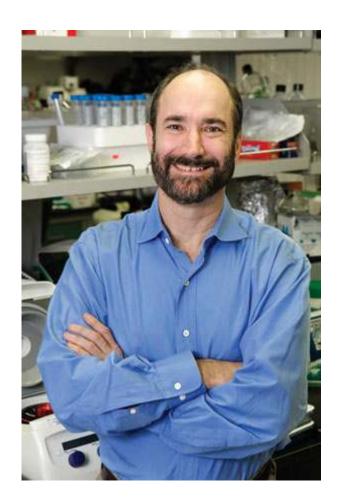
11.07.2013

### Критическое мышление

Новый технологический уровень исследований при старой парадигме интерпретации

### **Boston University:**

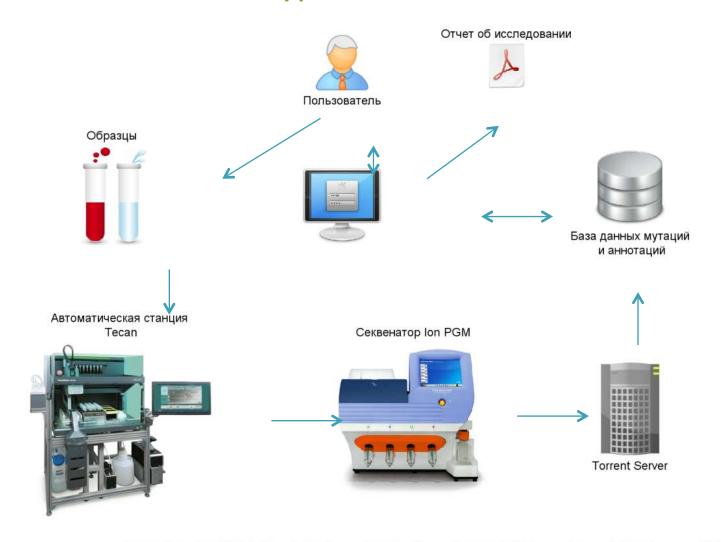
**GMS GE 705: Critical Thinking in Genetics** and Genomics



Mike Snyder

## Наш вклад

### Неонатальная NGS-генодиагностика



# Медицинская биоинформатика - настоящее и тенденции

Практические аспекты разработки ПО

**Антон Брагин Sequoia Genetics abragin@sequoiag.com** 



Mu Bo Cl Pe Mi Oz  MUMANA BOWTIE  OLISTALW  PECAN  MIRA 97 IL OLISOZEP  OS BW  Cf  COMBINE RCOFFEE  FCOFFEE MCOFFEE  FCOFFEE MCOFFEE  FOR CELEBA IL FORGE IL  FORGE IL	ALLPATHS-LG [S]   SOAP-DE NOVO [S]   UNIGENE   GLIMMER   N		Eb Bp Pl Ga
	-     - 9	Au Mp In Gx	Ap Ig Ut Sa
BFAST  MA  MV  Ar  MERACULOUS  Sm  MI  Pb  Ph  Qu  QUAKE [1]  QUAKE [1]  QUAKE [1]	Le Co Pa Oa Ep Os US II O OS BASES OS Ji	Sn Mr Ch Kn SNAP LADO KNIME	APPOLIO   IGB

TOOLS FREE FOR ACADEMICS ONLY	02	xx	01	xx	xx	xx	09	99	11	97	02
	BI	Vm	Sh	Fa	Cm	No	Ab	Fg	Mk	Gn	Uc
COMMERCIAL	BLAT XX		SSAHA2	FASTA XX	CROSSMATCH XX	NOVOALIGN 05	ABYSS [S]	FGENESH	MAKER	GENSCAN	ucsc
10025	Zm	Pc	Ne	Pp	Gs	Cb					
	ZOOM	PCAP	NEWBLER	PIPELINE PILOT	GENESPRING	CLC BIO					

www.eagiegenomics.con

This table is distributed under the Creative Common License, it can be downloaded from our web - Eagle Genomics Ltd 2012, All rights reserved.

# Специалисты из разных областей по-разному воспринимают одни и те же вещи



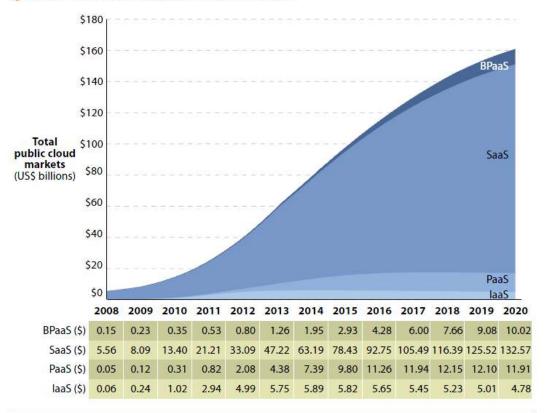
# Текущие представления о работе биомедицинского ПО

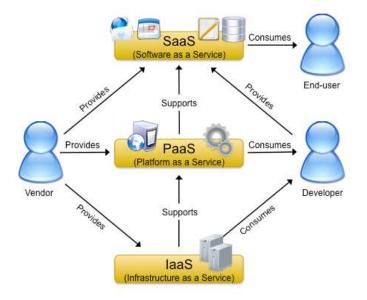


# Общая тенденция – от продукта к сервису

Figure 3 Forecast: Global Public Cloud Market Size, 2011 To 2020

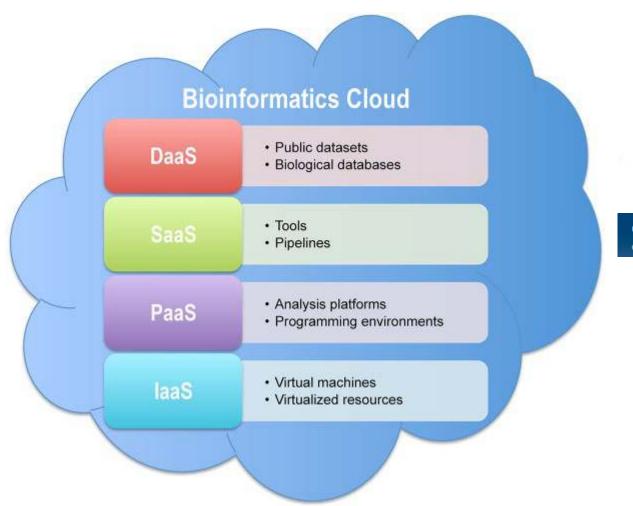
The spreadsheet detailing this forecast is available online.





Source: Forrester Research, Inc.

# Облачные сервисы в биоинформатике



e.g.

dbSNP Short Genetic Variations



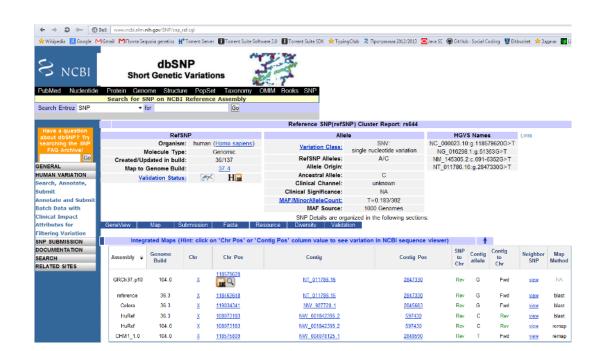
Crossbow
Genotyping from short reads using cloud computing





## DaaS: dbSNP

- •Публично доступные данные (часто с возможностью внесения данных)
- •Наличие API, позволяющего осуществлять программный доступ к данным
- •Возможность развертывания локального экземпляра базы данных
- •Наличие публичных копий БД в облаке (например, AWS)



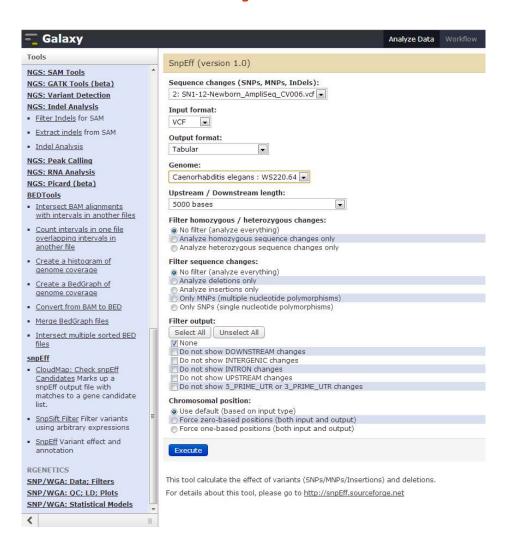
# SaaS: Crossbow

- •Предназначено для конечного пользователя
- •Ограниченный, четко определенный функционал

#### Crossbow 1.2.1

AWS ID *	
AWS Secret Key *	
AWS Keypair Name	gsg-keypair Look it up
,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	Check credentials
Job name	Crossbow
Job type	Crossbow
	Just preprocess reads
Input URL *	s3n://
	Check that Input URL exists
Output URL *	s3n://
	Check that output URL doesn't exist
Input type	Preprocessed reads
	Manifest file
Truncate length	(If blank or 0, truncation is disabled)
	Skip reads shorter than truncate length
Discard fraction	0
Quality encoding	Phred+33 🔻
Genome/Annotation	
GETOTIC/ MITOGRADI	Human (v38, dbSNP 130)
	Specify reference jar URL:
	s3n://
	Check that reference jar URL exists
Bowtie options	-m 1
SOAPsnp options	-2 -u -n -q
Additional SOAPsnp options for haploids	-r 0.0001
Additional SOAPSNP options for diploids	-r 0.00005 -e 0.0001
Chromosome ploidy	All chrosmosomes are diploid
	All are haploid
	All are diploid except:
Options	Keep cluster running after job finishes/aborts
# EC2 instances	1
Instance type	c1.xlarge (recommended)
Made with the help of	
/JotForm	(Submit)
Please cite: Langmead computing. Genome 8/o	B, Schatz MC, Lin J, Pop M, Salzberg SL. <u>Searching for SNPs with cloud</u> logy 10:R134.

# **PaaS: Galaxy**



- •Возможность построения собственных pipeline'ов анализа
- •Расширяющийся функционал
- •Интеграция функциональных модулей из разных источников (samtools, snpEff и др.)
- •Возможность хранения данных

## **PaaS: Cloud BioLinux**

- •Возможность развертывания в пределах произвольной IT- инфраструктуры
- •Расширенные возможности настройки
- •Необходимость построения собственных pipeline'ов анализа
- •Интеграция функциональных модулей из разных источников
- •Возможность хранения данных в облаке



# Сложность процесса не всегда коррелирует с полезностью результатов



#### **Browse your search results** ## Home > IP 2012-11-21 18:57:39 > Full results of workflow IP 2012-11-21 18:57:39 Result View Export Applications Filtering Grouping . Sorting -Results 1-50 of 2,500 Results are Ungrouped 24-Jul-2008 GLAXOSMITHKLINE BIOLOGICALS S Detection method and materials therefor probable disclosure ( Perfect match A, DELFT DIAGNOSTIC LABORATORY B V found by automated GLAXOSMITHKLINE BIOLOGICALS S.A. DETECTION METHOD AND MATERIALS THEREFOR 100 W02006077102 27-Jul-2006 probable disclosure ( Perfect match (BE) : DELFT DIAGNOSTIC found by automated LABORATORY BV (NL) parsing) 100 JF2008526244 24-Jul-2008 GLAXOSMITHKLINE BIOLOGICALS S Detection method and materials therefor probable disclosure (r Perfect match A, DELFT DIAGNOSTIC LABORATORY B V found by automated parsing) 27-Jul-2006 100 WQ2006077102 GLAXOSMITHKLINE BIOLOGICALS S.A. DETECTION METHOD AND MATERIALS THEREFOR probable disclosure () Perfect match (BE) : DELFT DIAGNOSTIC found by automated LABORATORY BY (NL) parsing) 100 JF2011518333 23-Jun-2011 QIAGEN GAITHERSBURG INC COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC TBD (information not PROBES FOR DETERMINING THE PRESENCE OF A TARGET GQ-Pat) 100 EF2262911 22-Dec-2010 QIAGEN GAITHERSBURG, INC. (US) COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC TBD (information not Perfect match PROBES FOR DETERMI NING THE PRESENCE OF A TARGET GQ-Pat) NUCLEIC ACID 100 US20090298187 03-Dec-2009 QIAGEN GAITHERSBURG, INC. COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC claim: 10: 12: 14: 15 GAITHERSBURG, MD PROBES FOR DETERMINING THE PRESENCE OF A TARGET QIAGEN GAITHERSBURG INC COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC TBD (information not PROBES FOR DETERMINING THE PRESENCE OF A TARGET GQ-Pat) NUCLEIC ACID QIAGEN GAITHERSBURG, INC. (US) COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC TBD (information not 100 EP2262011 22-Dec-2010 PROBES FOR DETERMI NING THE PRESENCE OF A TARGET GQ-Pat) NUCLEIC ACID 03-Dec-2009 COMPOSITIONS, METHODS, AND KITS USING SYNTHETIC claim: 10; 12; 14; 15 Perfect match

PROBES FOR DETERMINING THE PRESENCE OF A TARGET

NUCLEIC ACID

# ПО для анализа данных секвенирования следующего поколения (NGS)

Свойства, которыми должно обладать программное обеспечение для анализа клинического генетических данных:

- 1. Глубина анализа и акцент на клинически-значимых аспектах
- 2. Дружественность к пользователю и простота работы
- 3. Стабильность
- 4. Соответствие предоставляемой информации актуальным данным
- 5. Кроссплатформенность
- 6. Локализуемость

Два типа трудностей: технические и концептуальные

# Как обозначить генетический вариант?

ATGCGTGGACTGATGCGAGTCGGATGTAATAGGTGCTGAGAGG ATGCGTGGACTGATGCGAGTCGCATGTAATAGGTGCTGAGAGG

# Идентификация генетических полиморфизмов в общедоступных базах данных

Несмотря на несколько попыток формализации, идентификация мутаций в настоящее время основана на *de facto* стандартах, таких как форматы файлов и структура баз данных



# Mutation Nomenclature Extensions and Suggestions to Describe Complex Mutations: A Discussion

Johan T. den Dunnen<sup>1\*</sup> and Stylianos E. Antonarakis<sup>2\*</sup>

<sup>1</sup>MGC-Department of Human and Clinical Genetics, Leiden University Medical Center, Leiden, The Netherlands <sup>2</sup>Division of Medical Genetics, University of Geneva Medical School, Geneva, Switzerland

Consistent gene mutation nomenclature is essential for efficient and accurate reporting, testing, and curation of the growing number of disease mutations and useful polymorphisms being discovered in the human genome. While a codified mutation nomenclature system for simple DNA lesions has now been adopted broadly by the medical genetics community, it is inherently difficult to represent complex mutations in a unified manner. In this article, suggestions are presented for reporting just such complex mutations. Hum Mutat 15:7–12, 2000. © 2000 Wiley-Liss, Inc.

KEY WORDS: complex mutation; mutation detection; mutation database; nomenclature; MDI

Letters to the Editor

GERALDINE A. McDowell\* AND MIRIAM G. BLITZER<sup>†</sup>

\*Section on Human Biochemical Genetics, Human Genetics Branch, National Institute of Child Health and Human Development, National Institutes of Health, Bethesda; and †Division of Human Genetics, Departments of Pediatrics and Obstetrics and Gynecology, University of Maryland School of Medicine, Baltimore

#### Reference

Deiler JH (1909) The settlement of the German coast of Louisiana and the Creoles of German descent. Americana Germanica, Philadelphia

@ 1993 by The American Society of Human Genetics. All rights reserved. 0002-9297/93/5303-0027\$02.00

Am. J. Hum. Genet. 53:783-785, 1993

#### The Designation of Mutations

#### To the Editor:

Many different conventions are used for the primary designation of mutations. Commonly, the amino acid change that has been deduced from the nucleotide substitution is employed, but often the cDNA number, the genomic number, or even a "nickname" based on a restriction site or a patient's name has been used. These notations are not, of course, mutually exclusive, and several of them are used in first describing a mutation.

### Структура VCF-файлов

##fileformat=VCFv4.1 3020ЛОВОК
##fileDate=20090805
... Bapuahmы
#CHROM POS ID REF ALT QU

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFC
4	14370	rs6054257	G	Α	29	PASS	•••
16	87465	rs6040355	ATGC	Α	30	PASS	•••
21	14489	rs5896654	G	GAT	28	PASS	

Преимуществом VCF является возможность добавления структурированной информации в поля INFO

# dbSNP представляет собой наиболее полный общедоступный источник информации о генетической вариации



Представленная в виде VCF информация о генетических полиморфизмах человека представляет из себя файл размером около 8 Гб, который содержит:

53 417 385 записей, среди которых

SNP/MNP 46 177 146

(86,4%)

Инсерции3 088 167 (5,8%)Делеции4 152 072 (7,8%)

75 126 записей аннотированы как клинически-значимые

Пример 1:

# GATTTCGA... -> GATTTTCGA...

Это изменение последовательности может быть представлено <u>четырьмя</u> различными записями VCF

Pos	Ref	Alt
2	Α	AT
3	Т	TT
4	Т	TT
5	Т	TT

Пример 2:

# TCAGAGCAT... -> TCAGCAT...

Это изменение последовательности может быть представлено <u>тремя</u> различными записями VCF

Pos	Ref	Alt
2	CAG	С
3	AGA	Α
4	GAG	G

Пример 3:

# ACGTCCTCAG... -> ACGTCAG...

Это изменение последовательности может быть представлено <u>тремя</u> различными записями VCF

Pos	Ref	Alt
3	GTCC	G
4	TCCT	Т
5	ССТС	С

VCF допускает двусмысленность в обозначении одних и тех же генетических вариантов, что приводит к драматическим последствиям при его использовании для аннотации:

- Одно и то же фактическое изменение последовательности может иметь несколько аннотаций, противоречащих друг другу.
- При сравнении обнаруженных при секвенировании полиморфизмов с записями базы данных возможна потеря данных аннотации (известные мутации могут быть классифицированы как неизвестные).

Для того, чтобы оценить масштабы вырожденности, была написана программа, анализирующая файлы VCF dbSNP

### Оценка вырожденности записей dbSNP

Среди записей dbSNP:

Инсерции 3 088 167 (5.8%) в том числе 385437 синонимов (0.7% от общего числа мутаций, 12.5% от числа инсерций)

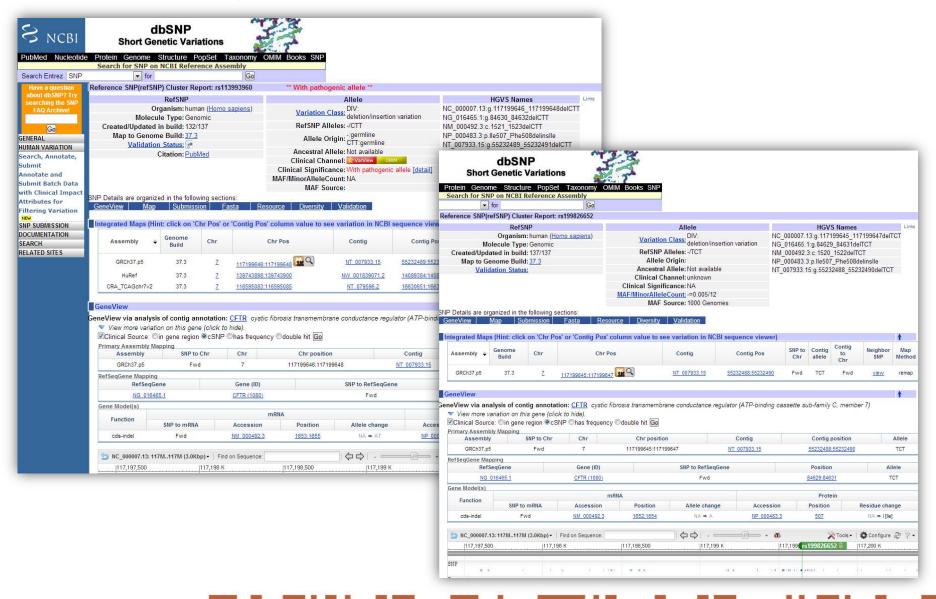
Делеции 4 152 072 (7.8%) в том числе **317437** синонимов (0.6% от общего числа записей, 7.6% от числа делеций) Общее число мутаций 53417385

Среди 75 126 клинически-значимых записей:

Инсерции 9 синонимов среди 3970 вариантов (0.2%)

Делеции 109 синонимов среди 1380 вариантов (7.9%)

### Оценка вырожденности записей dbSNP



# Формат VCF не подходит для хранения полиморфизмов в базе данных

1. Какой формат использовать для однозначной идентификации полиморфизмов?

2. Как обеспечить импорт данных из VCF, который является *de facto* стандартом?

Какой формат использовать для однозначной идентификации полиморфизмов?

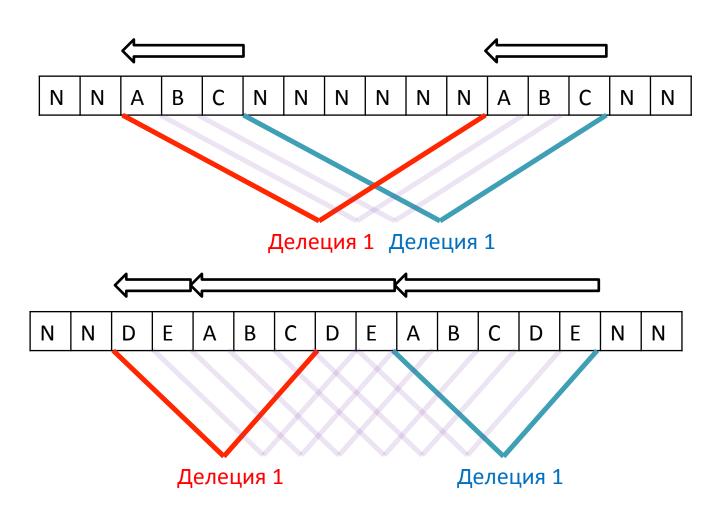
GATTTCGA... -> GATTTTCGA...

TCAGAGCAT... -> TCAGCAT...

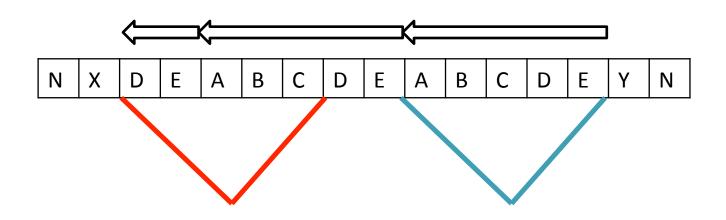
ACGTCCTCAG... -> ACGTCAG...

Вариант однозначно идентифицируется с помощью двух якорных оснований и последовательности, лежащей между ними в измененном варианте.

# Как обеспечить импорт данных из VCF, который является *de facto* стандартом?



# Как обеспечить импорт данных из VCF, который является de facto стандартом?



### Ключ мутации

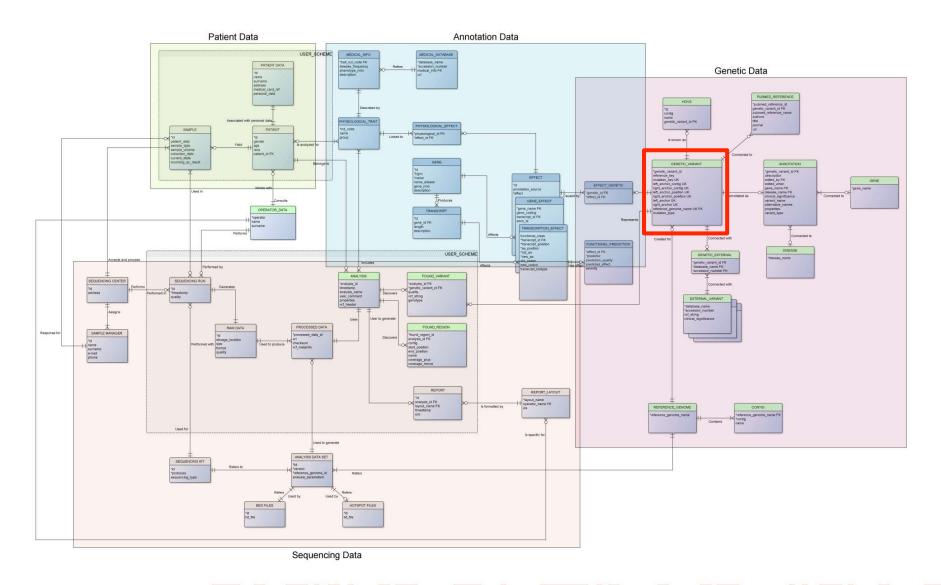


Все варианты (как аннотированные, так и полученные в результате проведения секвенирования) перед помещением в базу данных преобразуются в недвусмысленный формат

# Алгоритм преобразования VCF

Algorithm 1 Creation of AKA identifier Require: ref corresponds to reference in appropriate positions Require: ref and alt do not have similar bases in the same positions 1: procedure CreateIdentifier(ref[], alt[], coord, reference[])▶ Left anchor  $la \leftarrow coord - 1$ 2:  $ra \leftarrow coord + length(ref)$ ▶ Right anchor 3:  $change \leftarrow alt$ ▷ Sequence change 4: if length(alt) = 0 then ▶ This is deletion 5:  $la, ra \leftarrow ExpandBorders(ref, la, ra, reference)$ 6:  $change \leftarrow reference[la+1, ra-length(ref)]$ 7: else if length(ref) = 0 then ▶ This is insertion 8:  $la, ra \leftarrow ExpandBorders(alt, la, ra, reference)$ 9:  $change \leftarrow reference[la+1:coord] + alt + reference[coord:ra]$ 10: end if 11: return la, ra, change 12: 13: end procedure 14: procedure ExpandBorders(seq[], la, ra, reference[])  $l \leftarrow length(seq)$ 15:  $L \leftarrow length(reference)$ 16:  $i \leftarrow la$ 17:  $j \leftarrow ra$ 18: while  $i \geq 0$  and  $seq[(l - (la - i - 1)) \mod l] = reference[i]$  do 19:  $i \leftarrow i - 1$ 20: end while 21: while j < L and  $seq[(j - la - 1) \mod l] = reference[j]$  do 22: 23:  $j \leftarrow j + 1$ end while 24: return i, j25: 26: end procedure

# Секвенирование – это не только GATC



# Спасибо за внимание!

info@sequoiag.com