

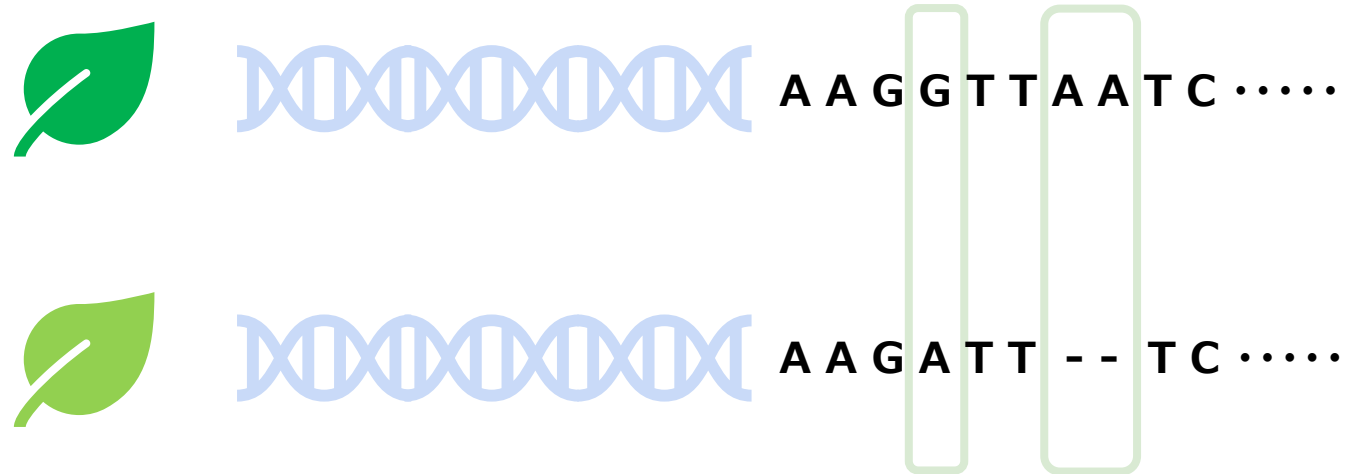


非モデル生物の遺伝子変異解析

フィルジェン株式会社

バイオインフォマティクス部(biosupport@filgen.jp)

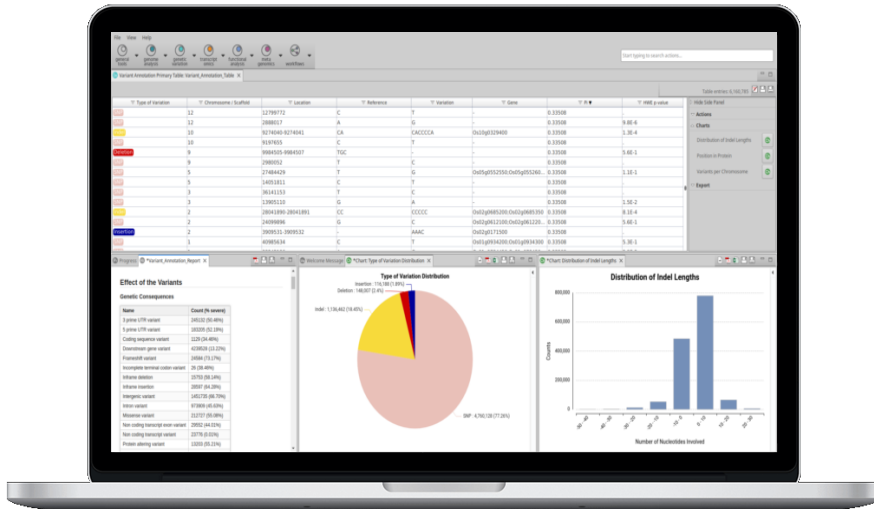
- 変異解析は配列データからバリエーション（SNP、挿入、欠失など）を同定するプロセスである。
- 変異の原因や表現形への影響を理解することができるため、ゲノミクス研究において基本的な解析。



- モデル生物では十分に発展しているが、非モデル生物では、一貫した変異解析パイプラインのための十分なガイダンスが不足している。
- 非モデル生物にはモデル生物に利用可能な遺伝情報が欠けていることが多く、遺伝的変異解析の実行が困難。

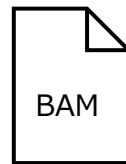


OmicsBoxの遺伝子変異解析機能



- 倍数性に対応したVariant Calling
 - バリアントのフィルタリング
 - バリアント アノテーション
- その他
- ゲノムアライメント
 - 遺伝子の機能予測

必要なファイル



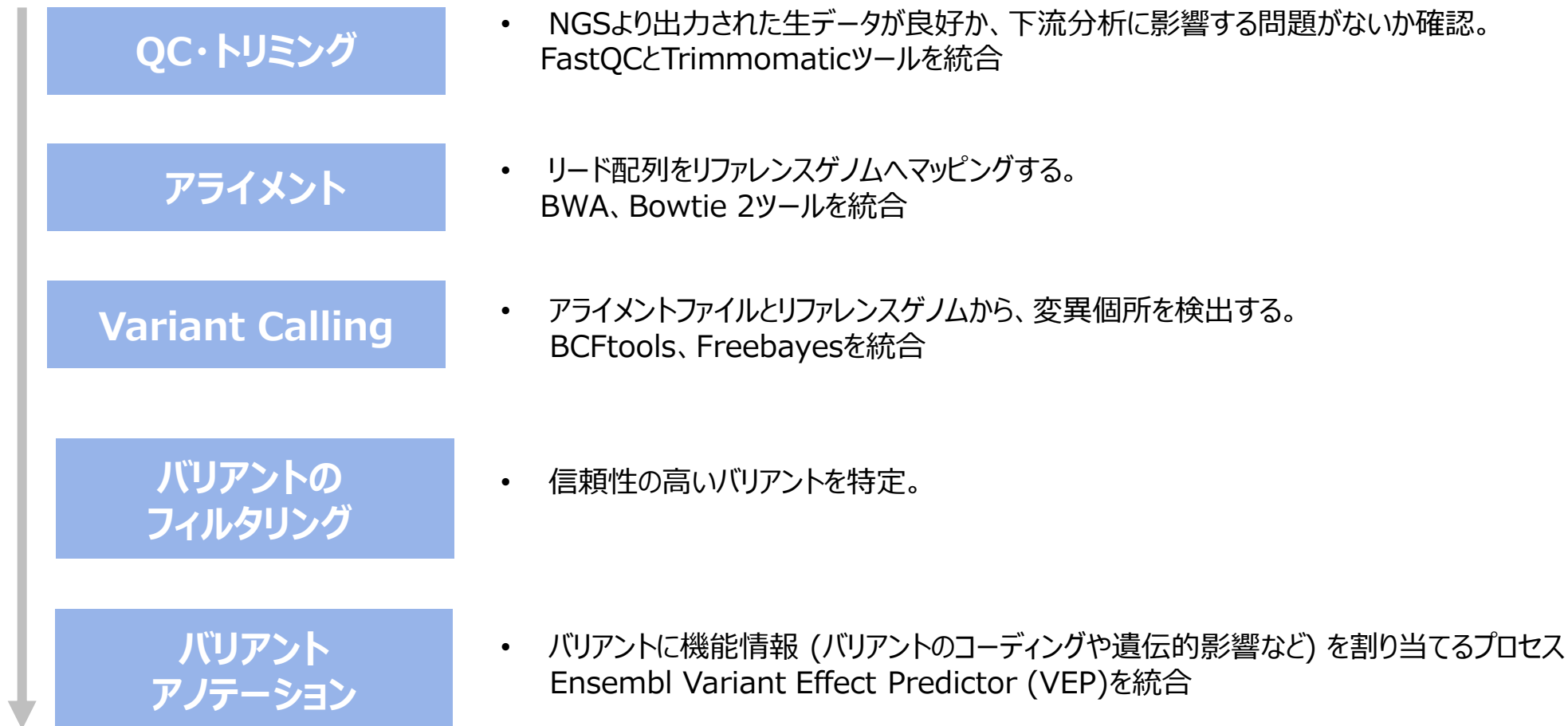
アラインメントされた配列ファイル



リファレンスゲノム



アノテーション ファイル





・データが良好か、下流分析に影響する問題がないか確認

Welcome Message | FASTQ Quality Check (Dataset) | FASTQ Quality Check (ERR1948631_1.fastq) | FASTQ Quality Check (clean_ERR1948631_1.fq) | Chart: Adapter Content

FASTQ Quality Check

Name: Dataset

Overall Results

Name	Per Base Sequence Quality	Per Sequence Quality Scores	Per Base Sequence Content	Per Sequence GC Content	Per Base N Content
ERR1948631_1.fastq	PASS	PASS	FAIL	PASS	PASS
clean_ERR1948631_1.fq	PASS	PASS	FAIL	PASS	PASS

Name	Sequence Length Distribution	Adapter Content	Overrepresented Sequences	Sequence Duplication Levels	Report
ERR1948631_1.fastq	PASS	FAIL	WARNING	FAIL	🔍
clean_ERR1948631_1.fq	WARNING	PASS	WARNING	FAIL	🔍

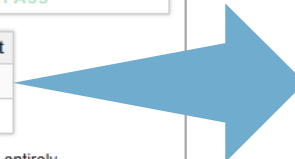
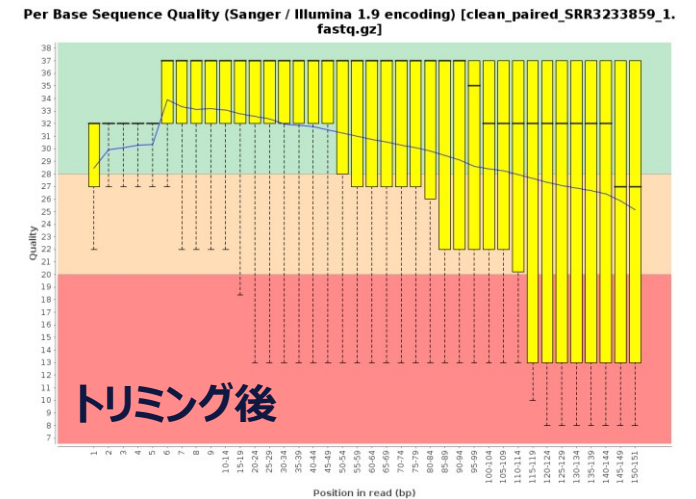
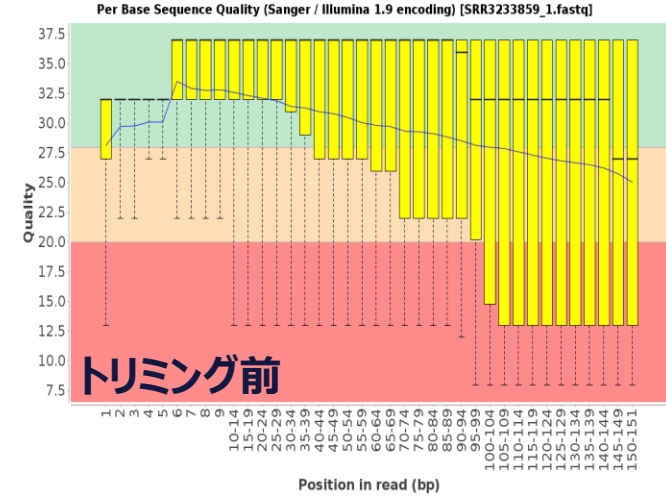
The FASTQ quality check task is performed by nine analysis modules. The table above provides a quick evaluation of whether the results of each module seem entirely normal (pass), slightly abnormal (warning) or very unusual (fail). Note that these evaluations must be taken in the context of what is expected from the library. For example, some experiments may be expected to produce libraries which are biased in particular ways. Therefore, the summary evaluations should be treated as pointers that guide the preprocessing of the libraries.



✓解析が終了するとレポートが作成

正常 (PASS)
わずかに異常 (WARNING)
異常 (FAIL)

シーケンスデータの品質をすばやく評価



レポートのアイコンをクリック→さらに詳細な結果を見ることが可能



・ショート参照ゲノムに効率的にマッピングして、リードの起源となった正しい遺伝子座を特定する
 →Variant Callingに必要なBAMファイルの作成が目的

Read Alignment (BWA) Results

Input 1: Reference Genome Sequences

assembly

Sequences	Minimum Length	Maximum Length	Average Length	Total Length
1	4,045,279	4,045,279	4,045,279	4,045,279

Input 2: FASTQ Files

File Name	Sample Name	Format	Sequencing
ERR2935851_1.fastq.gz, ERR2935851_2.fastq.gz	ERR2935851	FASTQ	Paired End

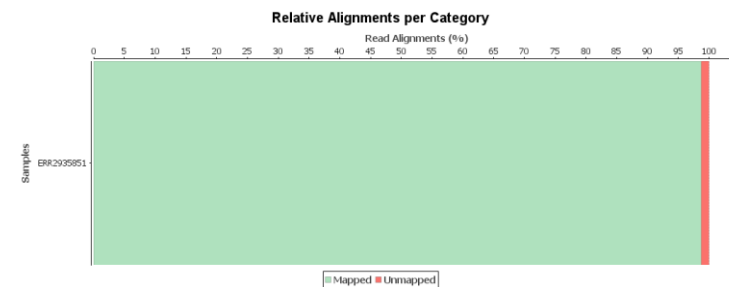
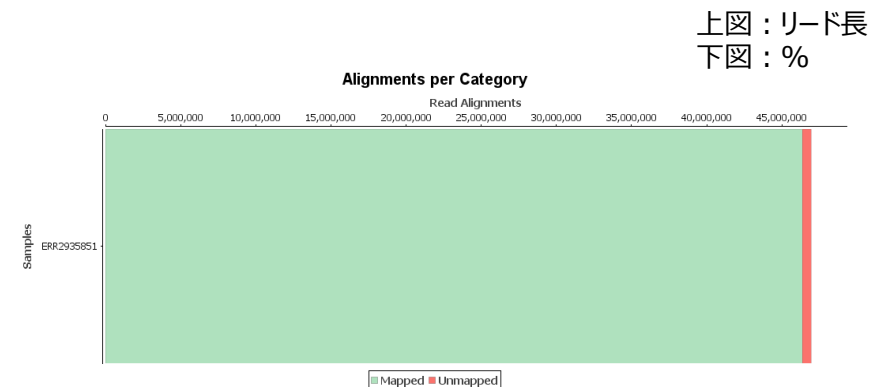
Results Overview

Globals

Sample	Total Alignments	Mapped	Supplementary	Unmapped	Duplicated Reads (estimated)	Duplication Rate
ERR2935851	46,952,430	46,358,042 / 98.734%	623,990 / 1.329%	594,388 / 1.266%	42,318,677 / 90.131%	99.4

どのくらいリードがマッピングされたか？

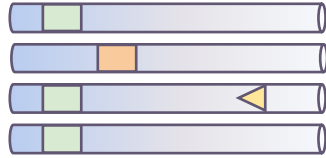
✓解析が終了するとレポートが作成



カテゴリグラフ



アラインメントされた配列ファイル



・前項で作成したアライメントファイル（BAMファイル）とリファレンスゲノムファイルを使用して変異個所を検出する。

Variant Calling Algorithms

BCFtools

Variant calling can be done applying BCFtools in two steps. The first step, BCFtools mpileup, reads the alignments and for each position of the genome constructs a vertical slice across all reads covering the position ('pileup'). Genotype likelihoods are then calculated, representing how consistent are the observed data with the possible diploid genotypes.

The second step, "bcftools call" then evaluates the most likely genotype under the assumption of Hardy-Weinberg equilibrium (in the sample context customizable by the user) using allele frequencies estimated from the data.

FreeBayes

FreeBayes is an haplotype-based variant detector and is a great tool for calling variants from a population. FreeBayes is a Bayesian genetic variant detector designed to find small polymorphisms, specifically SNPs (single-nucleotide polymorphisms), indels (insertions and deletions), MNPs (multi-nucleotide polymorphisms), and complex events (composite insertion and substitution events) smaller than the length of a short-read sequencing alignment.

We recommend using BCFtools for diploid species and FreeBayes for haploid/polyploid species.

Default < Back **Next >** Cancel Run

BCFtools

二倍体の生物向けのVariant Callingパッケージ

Freebayes

一倍体/倍数体の生物向けのVariant Callingパッケージ

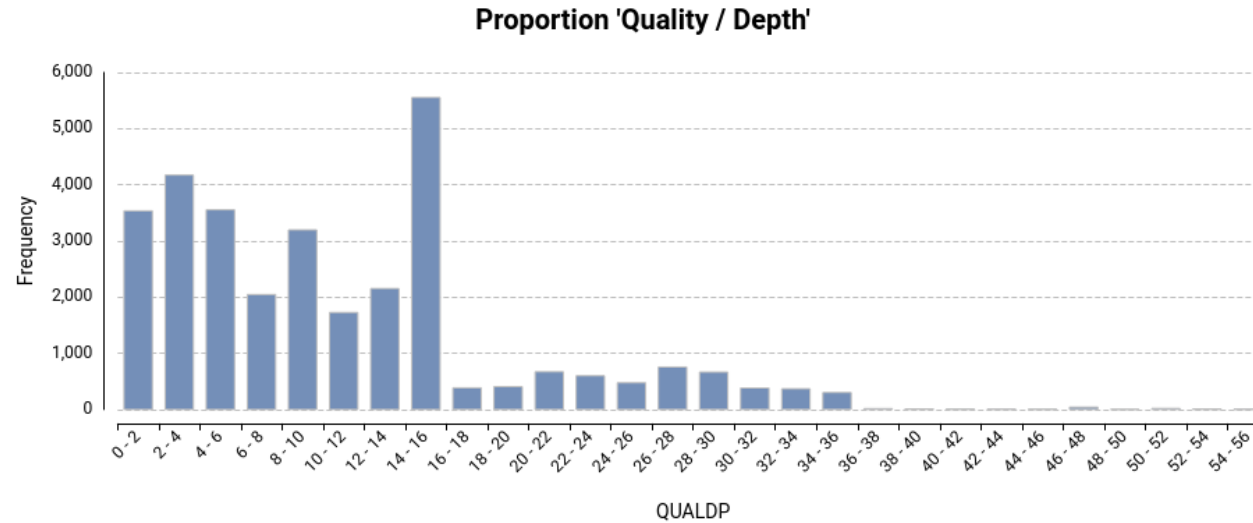
データに適切なパッケージを選択可能

Results

VCF File Saved as: /home/enrique/OmicsBoxWorkspace/freebayes_ew.vcf.gz

Type of variant	Frequency
SNP	3546
MIXED	2
MNP	182
INDEL	269

Number of alleles in a variant	Frequency
2	3989
3	9
4	1



✓解析が終了するとレポートが作成



✓さまざまな品質パラメーターの分布を示すいくつかのグラフ
この情報を参考に次項のフィルタリングを行う。



見つかったすべてのバリエントを含むVCF ファイル



・信頼性の高いバリエントが特定され、誤って検出バリエントが削除する。

Variant Filtering Report

Input Data

freebayes.vcf.gz

Results

Number of variants before filtering	Number of variants after filtering	Percentage filtered
3960	1001	74.7222%

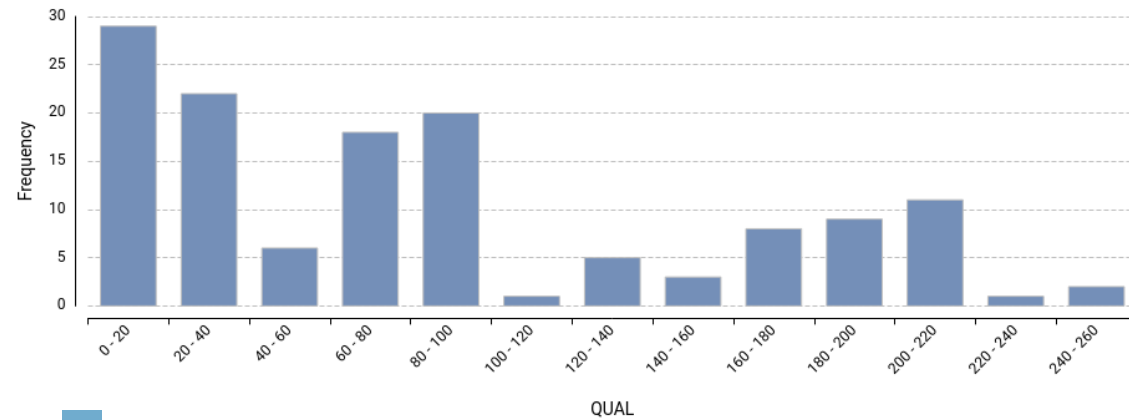
Parameters

Parameter	Value
Proportion 'Quality / Counts'	2.0
Raw Read Depth	5
Check Reads in Both Strands	true
Check if Reads are Balanced	true
Average Mapping Quality	50.0
Phred Quality	1.0
Remove Multiple Alleles	true

References

- Danecek P et al. (2021). Twelve years of SAMtools and BCFtools. *GigaScience*, 10(2).
- OmicsBox - Bioinformatics made easy. BioBam Bioinformatics (Version 3.0.23). March 3, 2019. www.biobam.com/omicsbox.

Phred Quality



✓更新されたグラフが作成される。この情報に基づいて、分析を続行するか、フィルタリング ステップのしきい値を調整するかを決定できる。



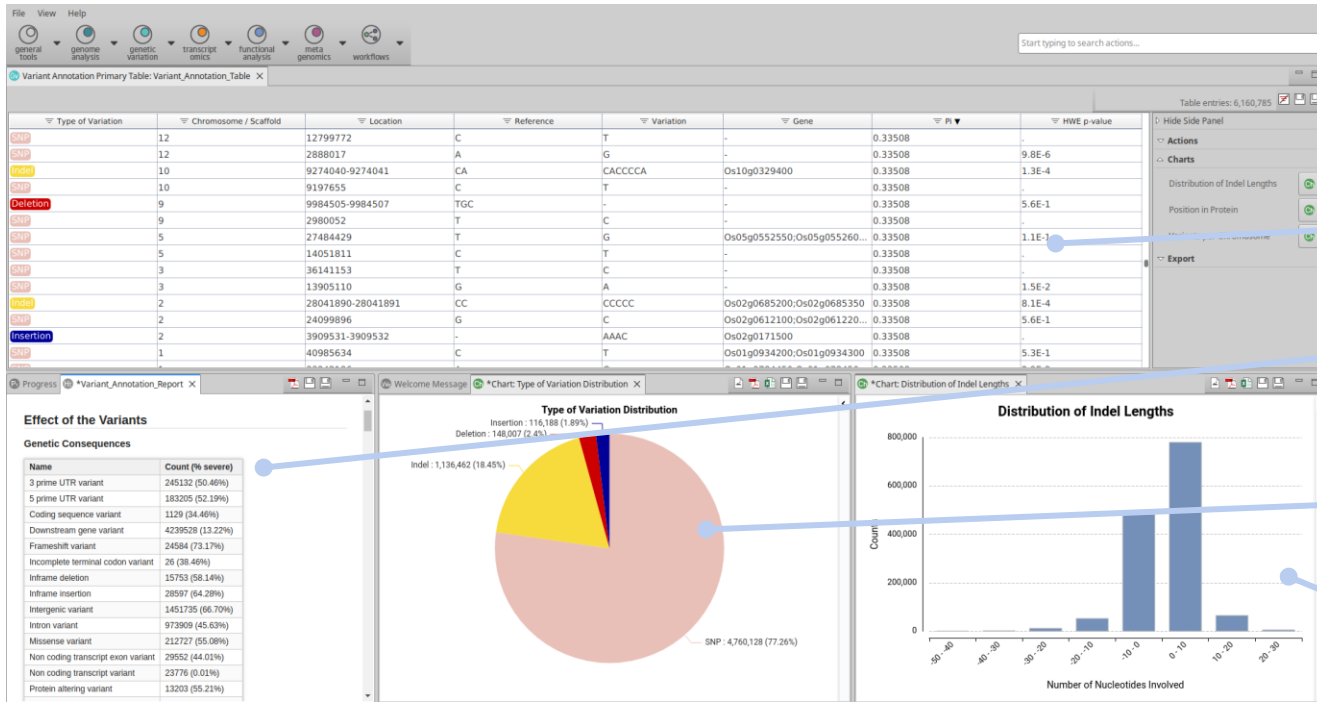
✓保持されているバリエントの数を表すレポートが作成



信頼性の高いバリエントを含むVCF ファイル



• バリエーションに機能情報（バリエーションのコーディングや遺伝的影響など）を割り当てるプロセスであり、ゲノム配列解析において重要なプロセス。別途外部データベースなどからアノテーションファイル（GTFファイル）をダウンロードしておく必要がある。



✓ 次のような包括的な情報が作成

• 各変異の位置、変異によって影響を受ける遺伝子、および2つの集団遺伝学値を示す表。

• 遺伝的およびコーディングの結果とサンプルのヘテロ接合性に関する情報を含むレポート。

• バリエーションの種類を示す円グラフ。

• 品質管理チャート

バリエーションによって影響を受ける遺伝子機能の名前 (または ID)

amino acid change
バリエーションがタンパク質コード配列に影響を与える

Annotation Details

Gene: Os06g0130800

Feature ID	Feature Type	Consequence	cDNA Position	CDS Position	Protein Position	Amino Acids	Codons	Impact	Distance	Strand
Os06t0130800-01	Transcript	upstream_gene_variant	-	-	-	-	-	MODIFIER	2453	-1

Gene: Os06g0130900

Feature ID	Feature Type	Consequence	cDNA Position	CDS Position	Protein Position	Amino Acids	Codons	Impact	Distance	Strand
Os06t0130900-00	Transcript	frameshift_variant	66-71	65-70	22-24	EKA/EKGX	gAGAAAGcc/gAGAAAGGcc	HIGH	-	-1

Gene: Os06g0131001

Feature ID	Feature Type	Consequence	cDNA Position	CDS Position	Protein Position	Amino Acids	Codons	Impact	Distance	Strand
Os06t0131001-00	Transcript	frameshift_variant	567-572	567-572	189-191	GFL/GLSX	ggCTTTCTc/ggCCTTTCTc	HIGH	-	1

Gene: Os06g0131100

Feature ID	Feature Type	Consequence	cDNA Position	CDS Position	Protein Position	Amino Acids	Codons	Impact	Distance	Strand
Os06t0131100-01	Transcript	upstream_gene_variant	-	-	-	-	-	MODIFIER	2338	1
Os06t0131100-02	Transcript	upstream_gene_variant	-	-	-	-	-	MODIFIER	2360	1

✓バリエーションの行を右クリックし、「Show Annotation Details」をクリックすると、そのバリエーションの影響を受ける遺伝子ごとに1つのテーブルが表示されたレポートが開く。

OmicsBox のGenetic Variation

- アラインメントファイル、リファレンスゲノムゲノム、アノテーションファイルの3つのファイルのみで、あらゆる場面で簡単に遺伝子変異解析を行うことができる。
- 非モデル生物のバリエーションのアノテーションに適している。
- Variant Callingとフィルタリングの実行は、コマンド不要のウィザードで簡単に操作できる。さらに、OmicsBoxは、信頼性の低いバリエーションを除外し、最も信頼できるバリエーションのみを選択してさらなる解析を行うことができる品質管理チャートとサマリーレポートを作成できる。



お問い合わせ先：フィルジェン株式会社

TEL 052-624-4388 (9:00～17 : 00)

FAX 052-624-4389

E-mail: biosupport@filgen.jp