

次世代シーケンスについて

Sequencerの比較

| Provider | Sequencer | Maximum Read Length (bp) | Maximum sequence yield per run (Gb) | |
|-----------------------|------------------------------|--------------------------|---|------------------------------|
| Short read sequencing | Illumina | HiSeq 4000 | 2 × 150 | 1,500 |
| | | HiSeq 3000 | 2 × 150 | 750 |
| | | HiSeq 2500 | 2 × 125 | 1,000 |
| | | NextSeq 500 | 2 × 150 | 120 |
| | | MiSeq | 2 × 300 | 15 |
| | | MiniSeq | 2 × 150 | 7.5 |
| | ThermoFisher Scientific | Ion Proton | 200 | 10 |
| Ion PGM | | 400 | 2 | |
| Ion Torrent S5 | | 600 | 10-15 (though there is a trade-off between read length and output) | |
| Long read sequencing | Pacific Biosciences (PacBio) | PacBio RSII | Half of data in reads >20,000 Max length >60,000 | 1 |
| | | PacBio Sequel | Half of data in reads >20,000 Max length >60,000 | 7 |
| | Oxford Nanopore (ONT) | MinION | Read length = DNA Fragment length Longest reported approaching 1 x 10 ⁶ | 10–20 |
| | | GridION X5 | Read length = DNA Fragment length Longest reported approaching 1 x 10 ⁶ | 50–100 |
| | | PromethION | Read length = DNA Fragment length Longest reported approaching 1 x 10 ⁶ | 11,000 (theoretical maximum) |

Short read sequencing

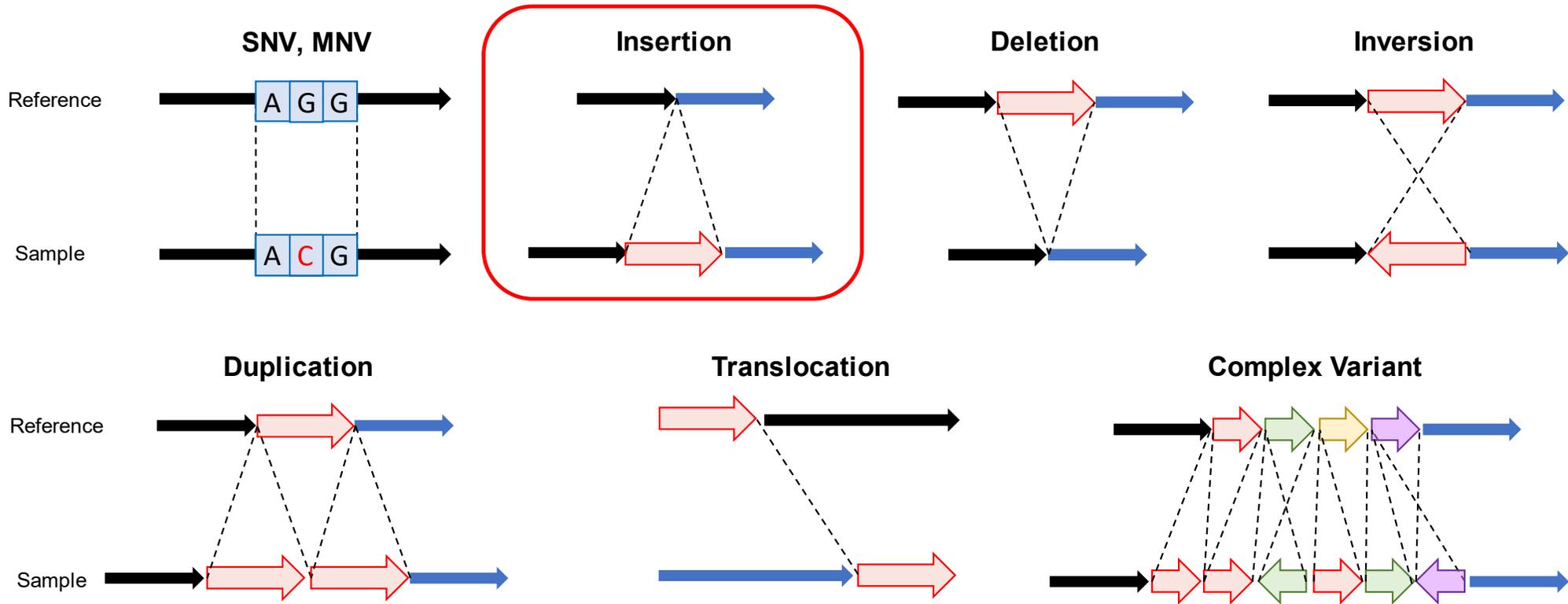
- 正確性が高い
- リピート配列が読めない



Long read sequencing

- 比較的リピート配列が読める
- 正確性がやや低い

起こりうる主な遺伝子変異



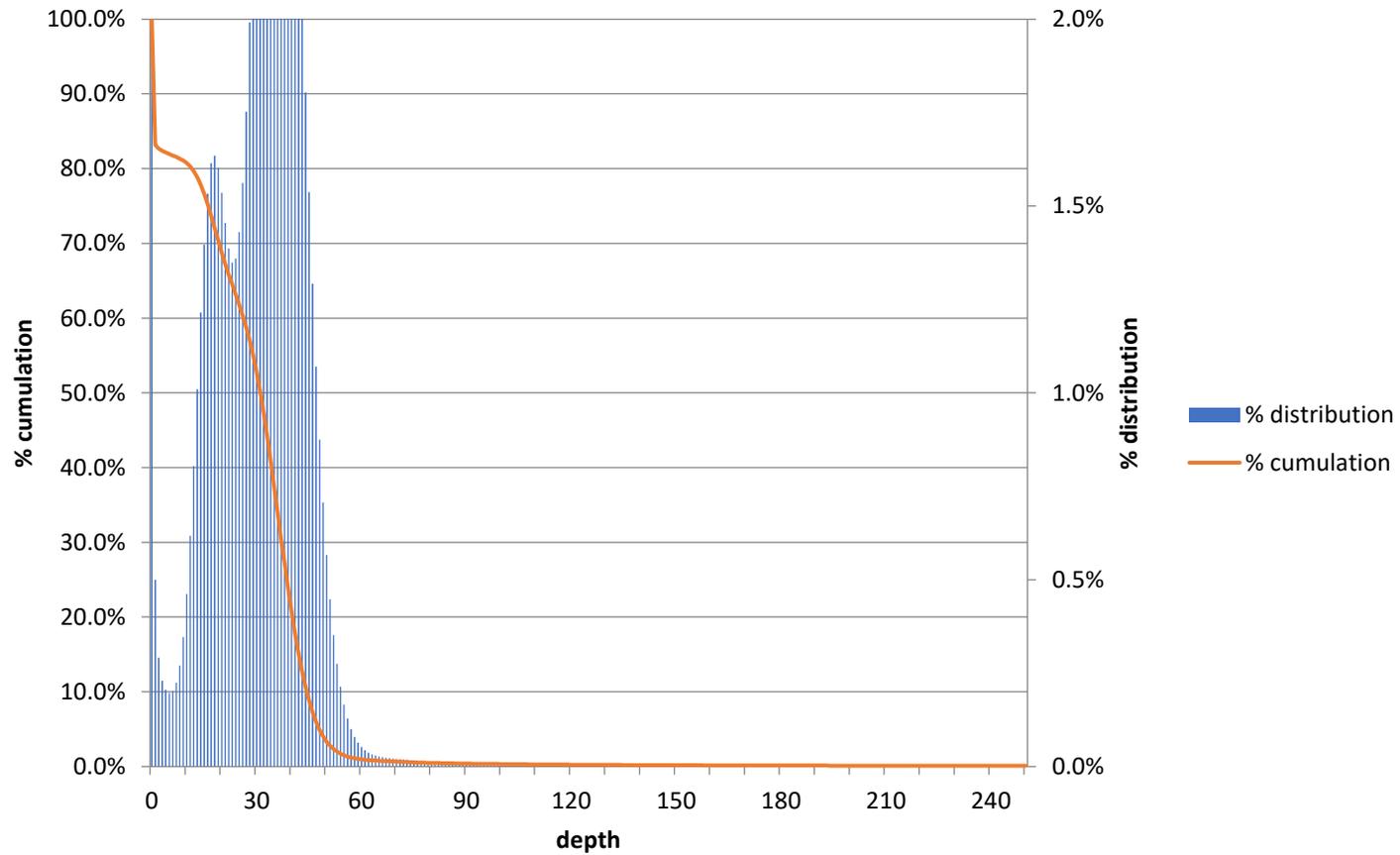
| Map summary | PG3676_01_a | PG3676_03_a |
|---|-------------|-------------|
| Total input reads | 795,743,304 | 714,037,020 |
| Total input reads(%) | 100 | 100 |
| Number of duplicate marked reads | 124,654,147 | 104,371,473 |
| Number of duplicate marked reads(%) | 16 | 15 |
| Number of unique reads (excl. duplicate marked reads) | 671,089,157 | 609,665,547 |
| Number of unique reads (excl. duplicate marked reads)(%) | 84 | 85 |
| Reads with mate sequenced | 795,743,304 | 714,037,020 |
| Reads with mate sequenced(%) | 100 | 100 |
| Reads without mate sequenced | 0 | 0 |
| Reads without mate sequenced(%) | 0 | 0 |
| QC-failed reads | 0 | 0 |
| QC-failed reads(%) | 0 | 0 |
| Mapped reads | 760,352,741 | 683,134,856 |
| Mapped reads(%) | 96 | 96 |
| Number of unique & mapped reads (excl. duplicate marked reads) | 635,698,594 | 578,763,383 |
| Number of unique & mapped reads (excl. duplicate marked reads)(%) | 80 | 81 |
| Unmapped reads | 35,390,563 | 30,902,164 |
| Unmapped reads(%) | 4 | 4 |

ave coverage=40x

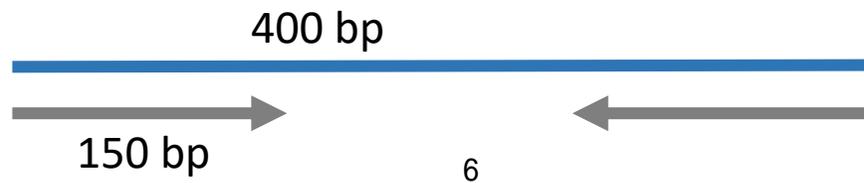
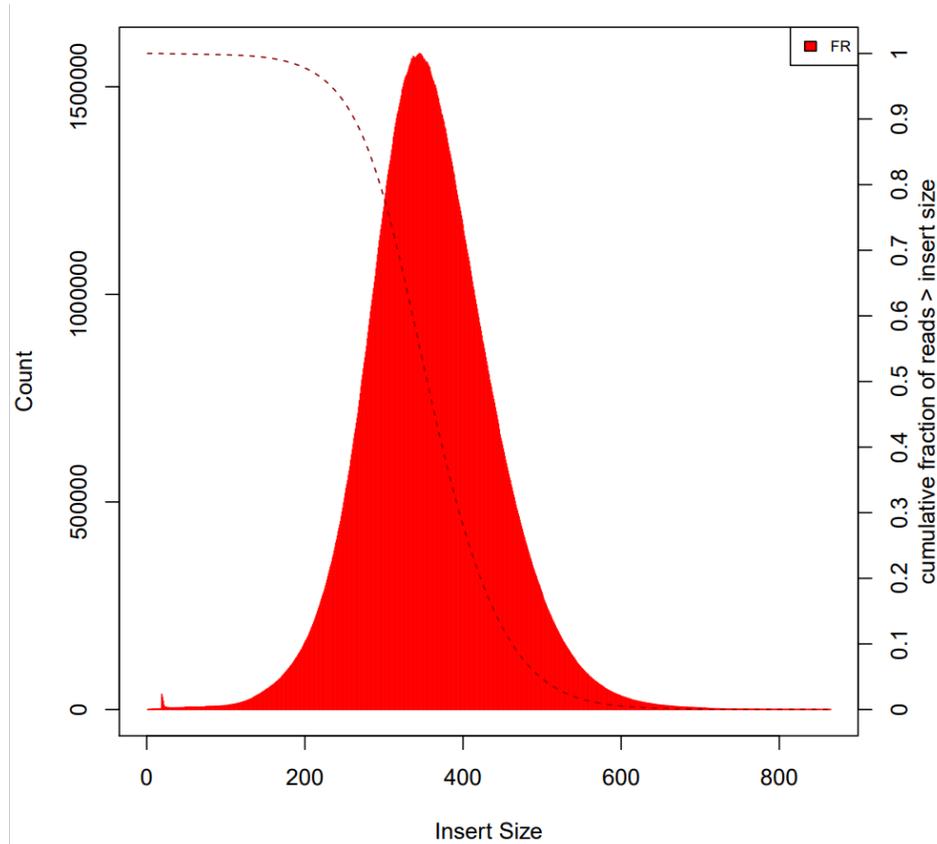
pair-end 150 bp

ave coverage=32x

Coverage分布

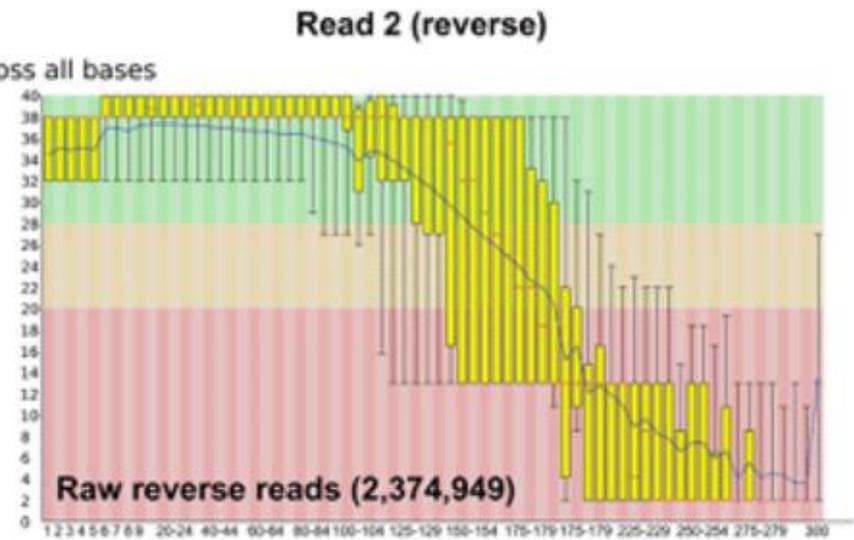
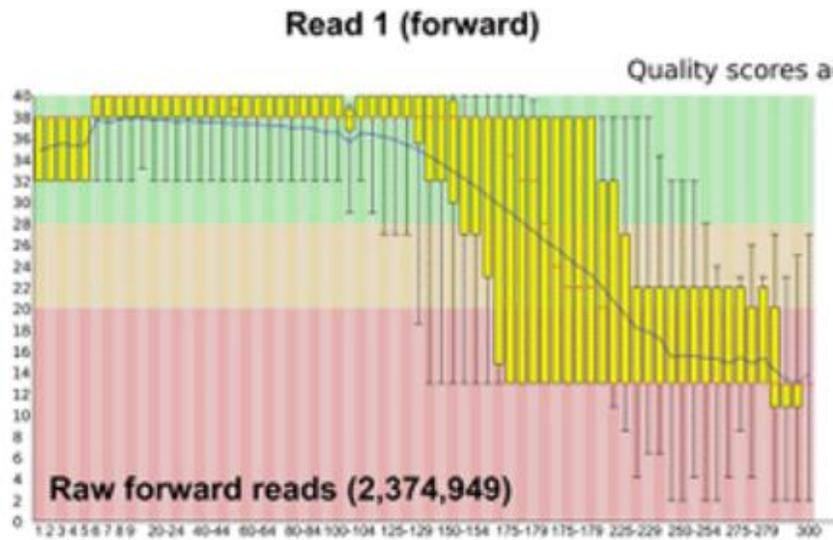


library size分布

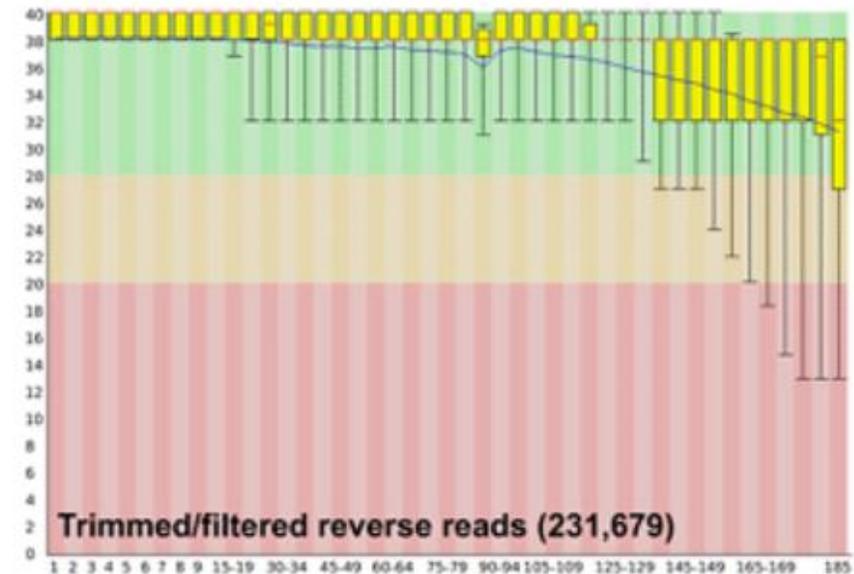
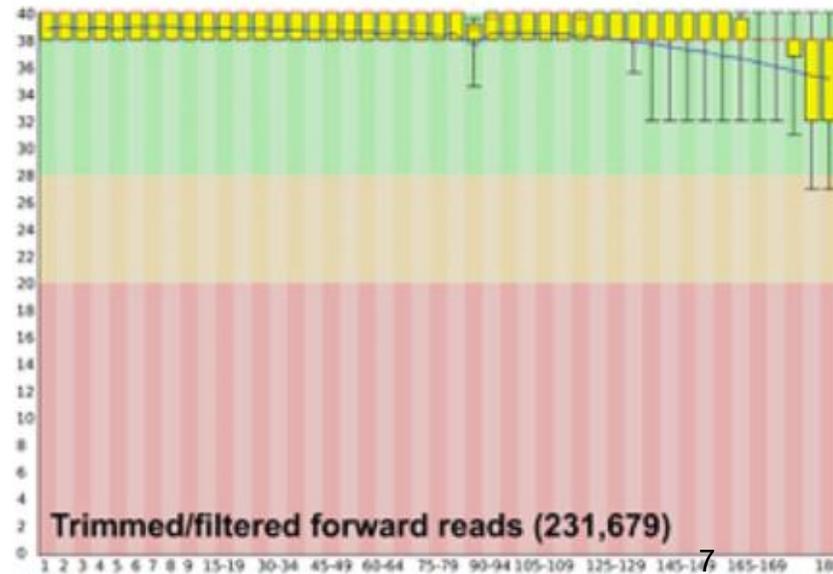


read quality分布

Before trimming



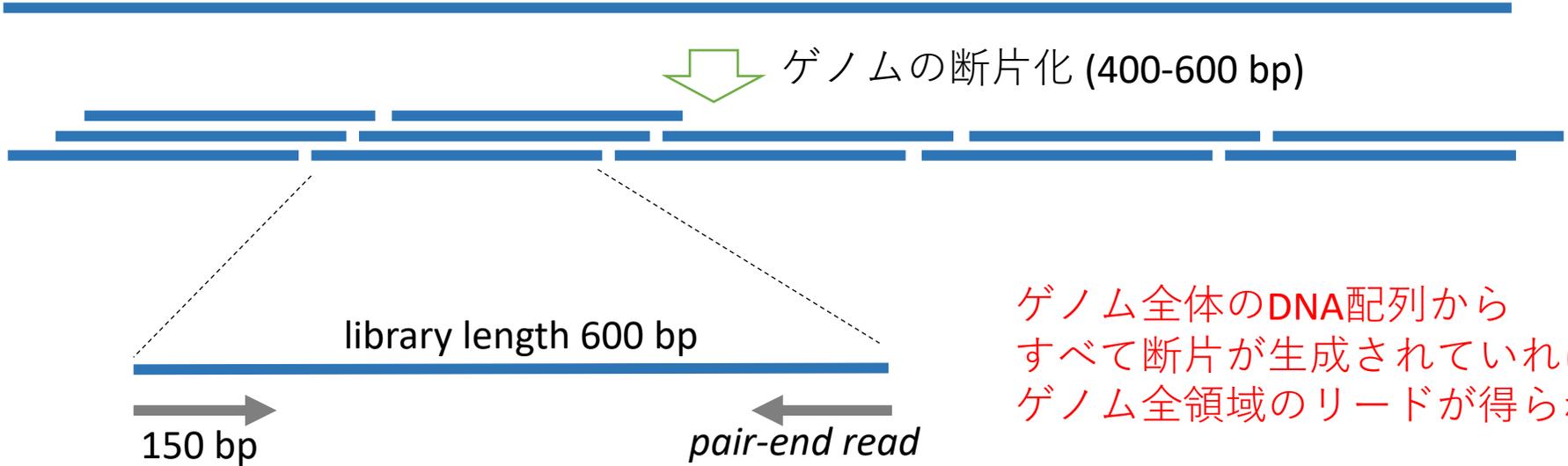
After trimming



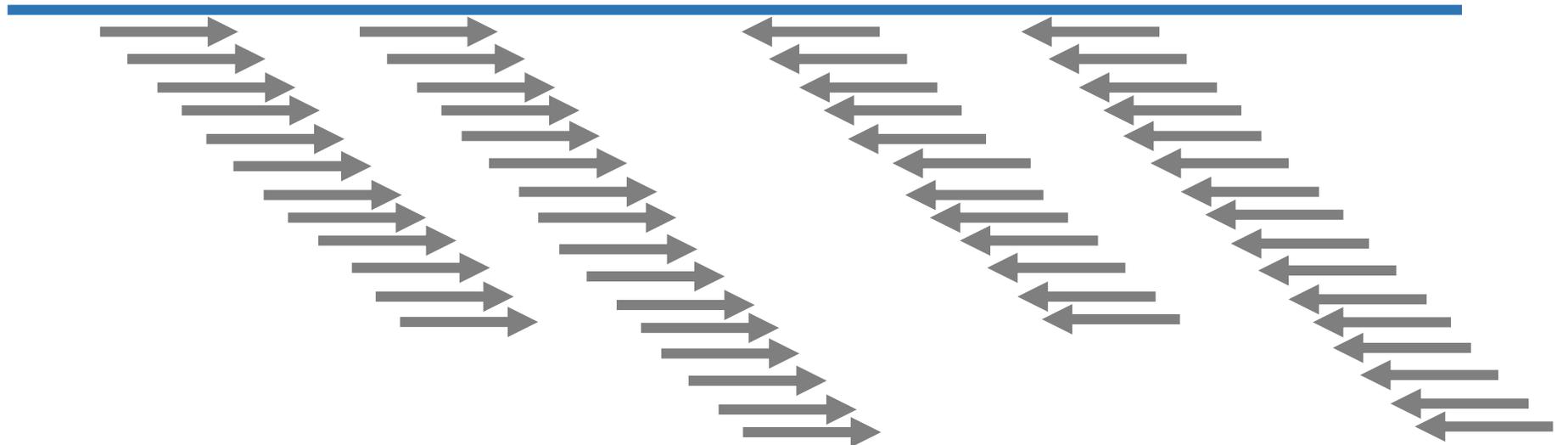
Position in read (bp)

Illumina short readはどのように読まれるか

↓ ゲノムの断片化 (400-600 bp)

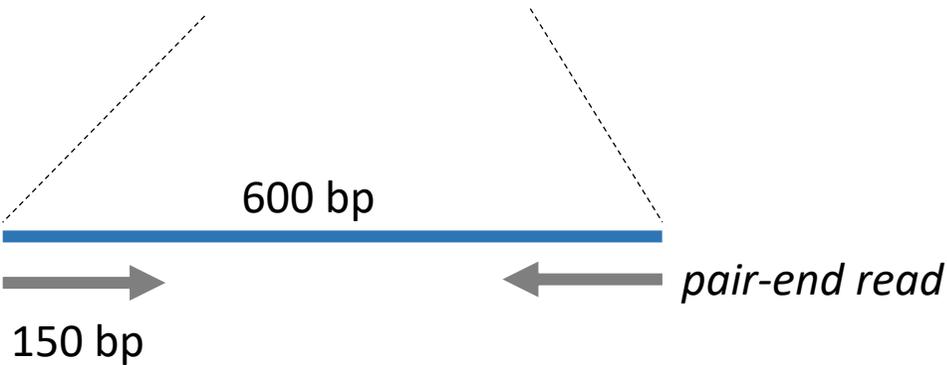


ゲノム全体のDNA配列から
すべて断片が生成されていれば
ゲノム全領域のリードが得られる



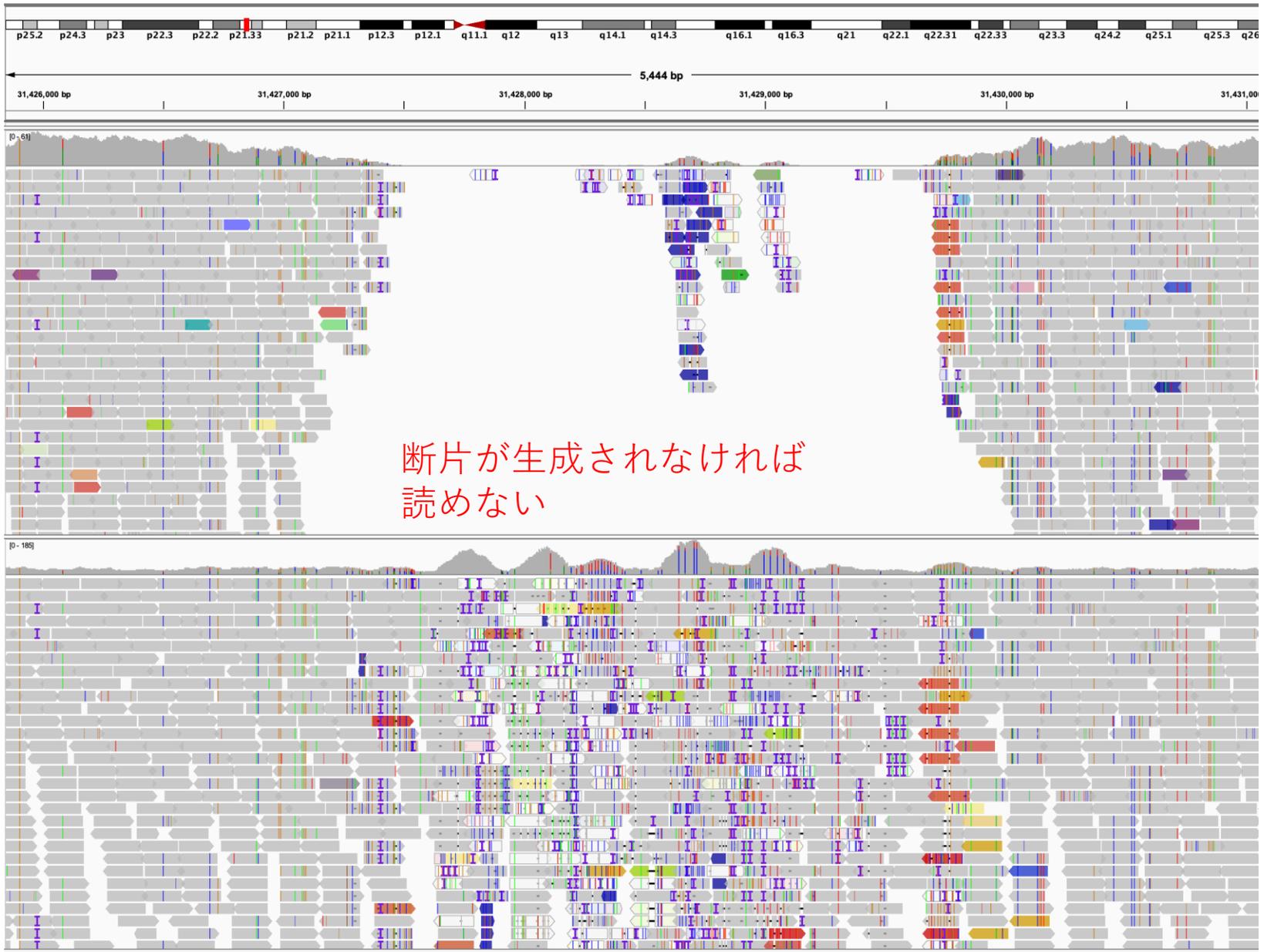
libraryができなかった領域は読めない

断片化 (600 bp)



ゲノム全体のDNA配列から
すべて断片が生成されなければ
読まない領域が存在する

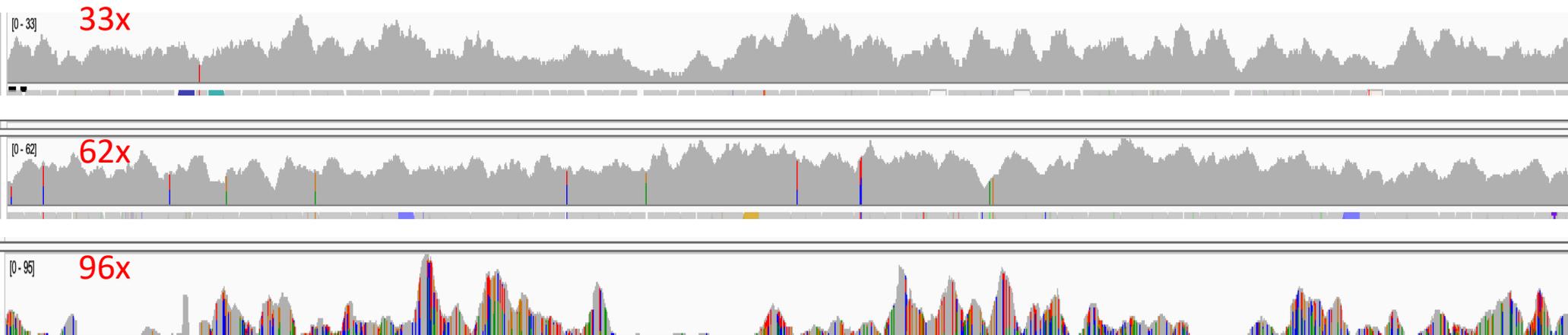




read coverageは全ゲノム領域で均一ではない

average read coverage（平均冗長度）とは、
全リード塩基数をゲノムサイズ（塩基数）で割っただけの値

average read coverage = 50 としても、各領域でcoverageは大きく違う

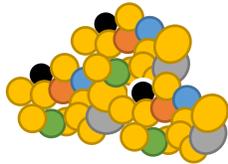


大事なものは解析したい領域で十分なcoverageを確保できているか

Q: coverageはどれくらい必要か？

A: どういうサンプルで何を見たいかで変わる！

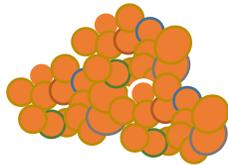
例えば、がんゲノム解析では100 x または300 x で全ゲノム解析
組織から取ったがん細胞は、それぞれ違う変異を持っている



100個細胞があり、その中の1つの細胞がある遺伝子に変異があったときに
その変異の検出確率は100分の1のため最低100リード必要。
通常、その変異をサポートするリードが複数必要なので、3個とすると
300 x で解析していないと検出できない

Q: coverageはどれくらい必要か？

A: どういうサンプルで何を見たいかで変わる！



逆に、homogenousな細胞集団ですべての細胞が同じ変異を持っているのであれば10 x もあれば十分とも言える
ゲノムのある領域ではその半分以下のcoverageであるため、実際にはその数倍のcoverageが必要 (30-40xあれば)

2倍体でヘテロ体であれば、その倍のcoverageで十分

そこに、各ゲノム領域でcoverageは変化するので、+アルファする

ある特定の領域だけを解析する、または
全ゲノムは30xで解析して、特定領域（挿入領域など）をdeep解析
したい場合はAmplicon seqで5000x-30,000xで解析できる

全ゲノム領域を厚く読むのは費用が急に上がる
最初から高深度でNGS解析する例は少ない

***実行可能性**の点から：

現在の多くの受託NGS全ゲノム解析の価格はかなり安くなっている
が、このときのデータ取得量はヒトゲノム換算で35-40xである

仮に75-80xでデータを取得するには数倍の費用と、解析リソースが
必要になる

マッピング結果の表示はIGV (Integrative Genome Viewer) が標準

control



一つのread

CRISPR



CA G C T A A C A G C A T T T C A C C A A A C C C A A A A G G A A A G G C T T G A A A C A A A G T A A A T C A C A A A T G T G G C A A A A C A G A T A G T G T A C T A T C C T T C T C A G G C A A T T A G C C T C A G C A G A G T G A C C C T C A T A G G A C C T T T G C A G G A G A T G G A C T G C T G A A T G T C

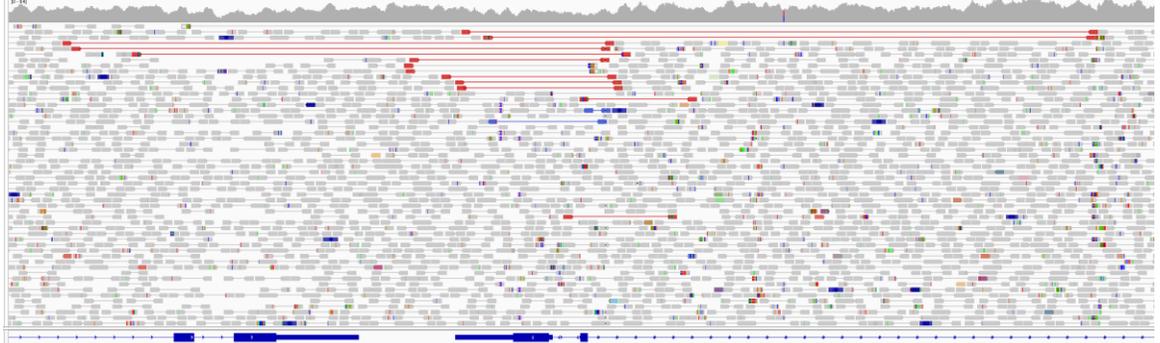
IGV 表示方法の違いで 見えるものが見えな くなる

一つのread

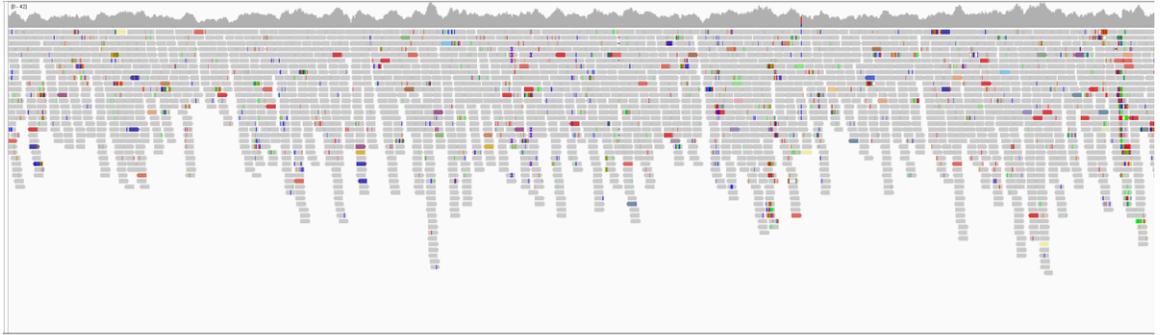
control



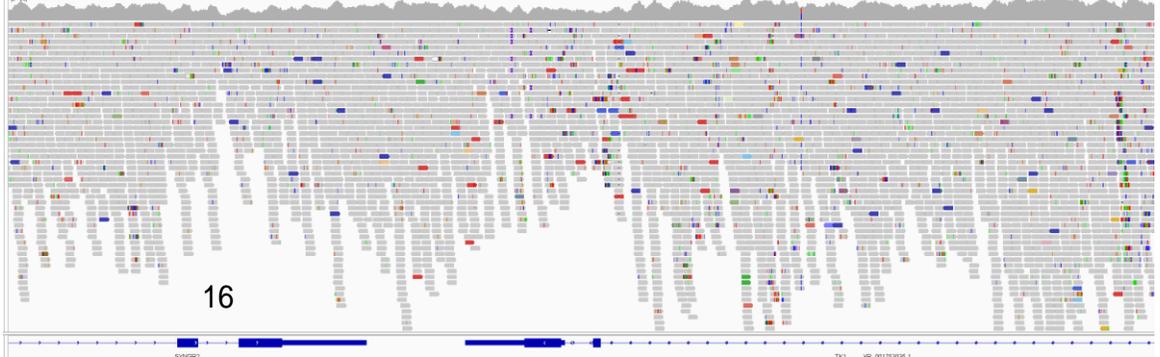
CRISPR



control



CRISPR



- Change Track Color...
- Experiment Type >
- Linked read view (BX)
- Linked read view (MI)
- Link supplementary alignments
- Link by tag...
- Group alignments by >
- Sort alignments by >
- Color alignments by >
- Shade alignments by >
- Re-pack alignments
- Shade base by quality
- Show mismatched bases
- Show all bases
- View as pairs
- Go to mate
- View mate region in split screen
- Set insert size options ...
- Show insertion markers



- normal_markeddup_nonSyncontrol.sorted.bam
- Rename Track...
 - Copy read details to clipboard
 - Change Track Color...
 - Experiment Type >
 - Linked read view (BX)
 - Linked read view (MI)
 - Link supplementary alignments
 - Link by tag...
 - Group alignments by >
 - Sort alignments by >
 - Color alignments by >
 - Shade alignments by >
 - Re-pack alignments
 - Shade base by quality
 - Show mismatched bases
 - Show all bases
 - View as pairs
 - Set insert size options ...
 - Show insertion markers
 - Quick consensus mode
 - Hide small indels
 - Small indel threshold...
 - Collapsed
 - Expanded
 - Squished

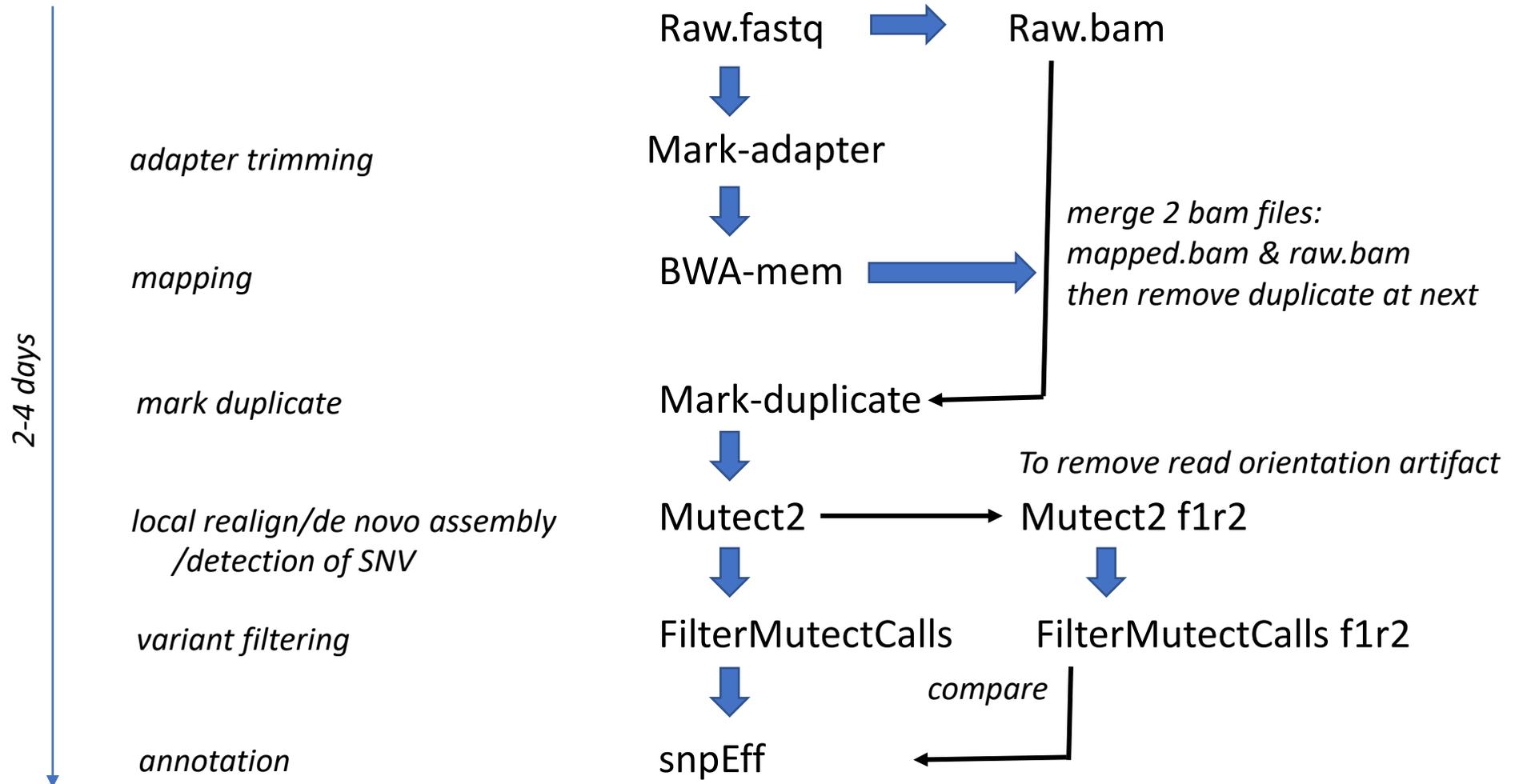


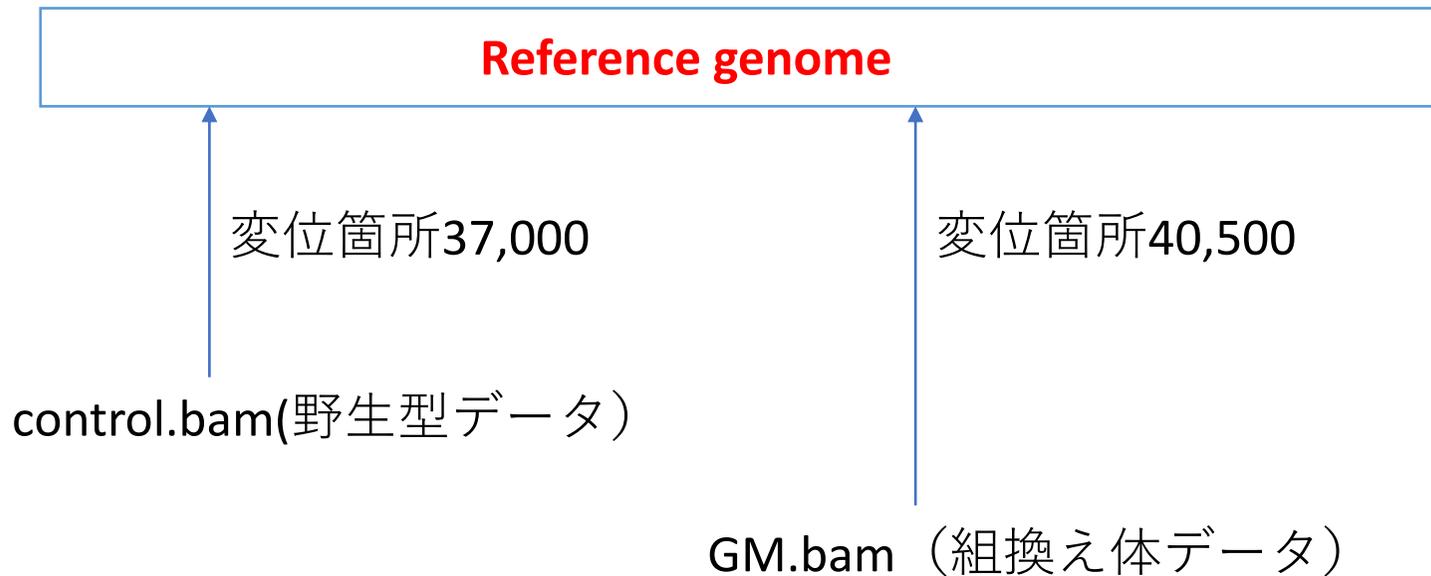
Nanopore sequencing



Long-read
では数10kbを
一気に読む
ことが可能

Workflow_1 (picard/gatk)





controlになかった変異だけを検出すればいい

(精度の高いreference genomeでなくとも問題ない
ある程度完成したdraft genomeがあれば精度よく
GMのみにある変異を高精度に検出できる)

Q: **NGS解析に求める要件は？**

1. read qualityの分布図
2. library調整法とlibrary (insert size)の分布とread length
insert size=600 bp, pair-end 150bp
3. coverageの分布図
4. 統計情報一覧
5. 解析ソフトと使用したparameter
6. マッピングIGV図と表示設定
7. 用いた参照ゲノム (versionなど)

-
- 1 0. 1-6を求めても本当にその通り解析しているかは確認できない
 - 1 1. blindのポジコンデータセット (いろんな変位が入った) を作製して、これを同時に解析してもらうことでデータの解析精度を保証できるが。。。