# Tumor profiling testの実践
# MSK-IMPACT

北海道大学病院　臨床研究開発センター
特任助教　天野虎次

# MSK-IMPACT

**Memorial Sloan Kettering-Integrated Mutation Profiling of Actionable Cancer Targets (MSK-IMPACT)**
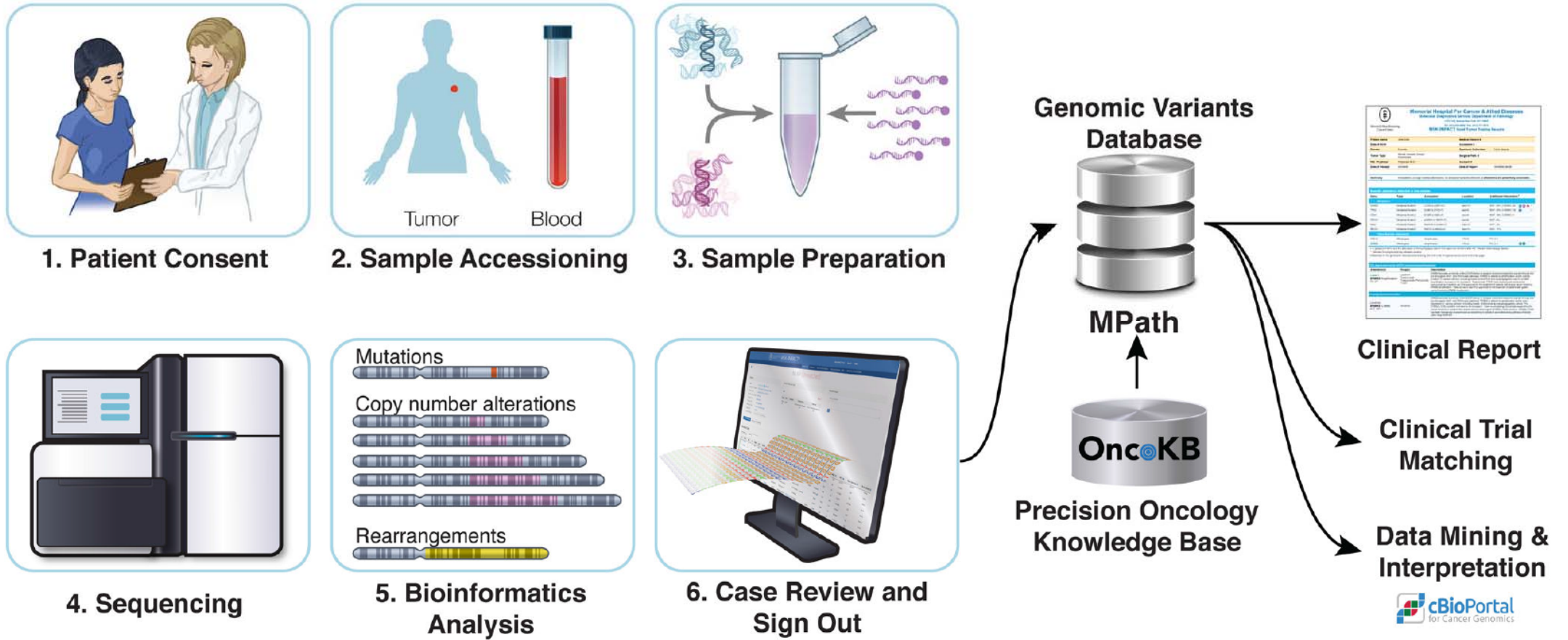
CrossMark

*A Hybridization Capture-Based Next-Generation Sequencing Clinical Assay for Solid Tumor Molecular Oncology*

Donavan T. Cheng,* Talia N. Mitchell,* Ahmet Zehir,* Ronak H. Shah,* Ryma Benayed,* Aijazuddin Syed,* Raghu Chandramohan,* Zhen Yu Liu,* Helen H. Won,* Sasinya N. Scott,* A. Rose Brannon,* Catherine O'Reilly,* Justyna Sadowska,* Jacklyn Casanova,* Angela Yannes,* Jaclyn F. Hechtman,* Jinjuan Yao,* Wei Song,* Dara S. Ross,* Alifya Oultache,* Snjezana Dogan,* Laetitia Borsu,* Meera Hameed,* Khedoudja Nafa,* Maria E. Arcila,* Marc Ladanyi,*[†] and Michael F. Berger*[†]

HOKKAIDO UNIVERSITY

# MSK-IMPACT Clinical Workflow



Zehir et al., Nat Med. 2017 June ; 23(6): 703–713

HOKKAIDO UNIVERSITY

# MSK-IMPACT Bioinformatics Team

＜ **Pathology : Diagnostic Molecular Pathology Service**＞

**Dr. Ladanyi (PI)　Lab**

- 主な**Stuff : 8**人
- **Researcher : 10+**人
- **Technician: 4+**人（**MSK-IMPACT**専任）

＜ **CMO (Center for Molecular Oncology)** ＞

**Dr. Berger (PI)　Lab**

- **Researcher : 6+**人
- **MSK-IMPACT Bioinformatics team : 8**人

**cBioPortal Bioinformatics team : 20**人　**+　SE(10+**人**)**

HOKKAIDO UNIVERSITY

# 実際の運用規模 (2016年)

- IMPACT update（2016/06/13）; 4 batch 12pool （n=12142 ）

| | | Mon | Tue | Wed | Thr | Fri |
|---|---|---|---|---|---|---|
| **6/13** | New batch | 229,230,231 | 232,233,234 | 235,236,237 | 238,239,240 | - |
| | Sequence | 222,223,224, 225,226 | - | 227,228 | - | 229,230,231 |

- **100-140 patients/week ( 1 pool = 10~12 case )**

- **Turn Around Time : 14+日（accession ~ sign out）**

- **QC meeting ： 毎週月曜:60分程度**

- **sign-out     : Pathologistと日程調整して適宜 (20 case / mtg)**

HOKKAIDO UNIVERSITY

# MSK-IMPACT-pipeline

● 使用している解析ツールは一般的なもの。
　　　　GATK, Picard, samtools 等々。

● 使用するpackageの選択は、比較検討を行って確認して決定している。( ABRA、Delly など）

● Filteringおよび一部の解析に自作プログラムを使用

# Analysis (GATK Best Practices)



https://software.broadinstitute.org/gatk/

HOKKAIDO UNIVERSITY

# Example of the Package Comparison

# Assembly-base and Mapping-base

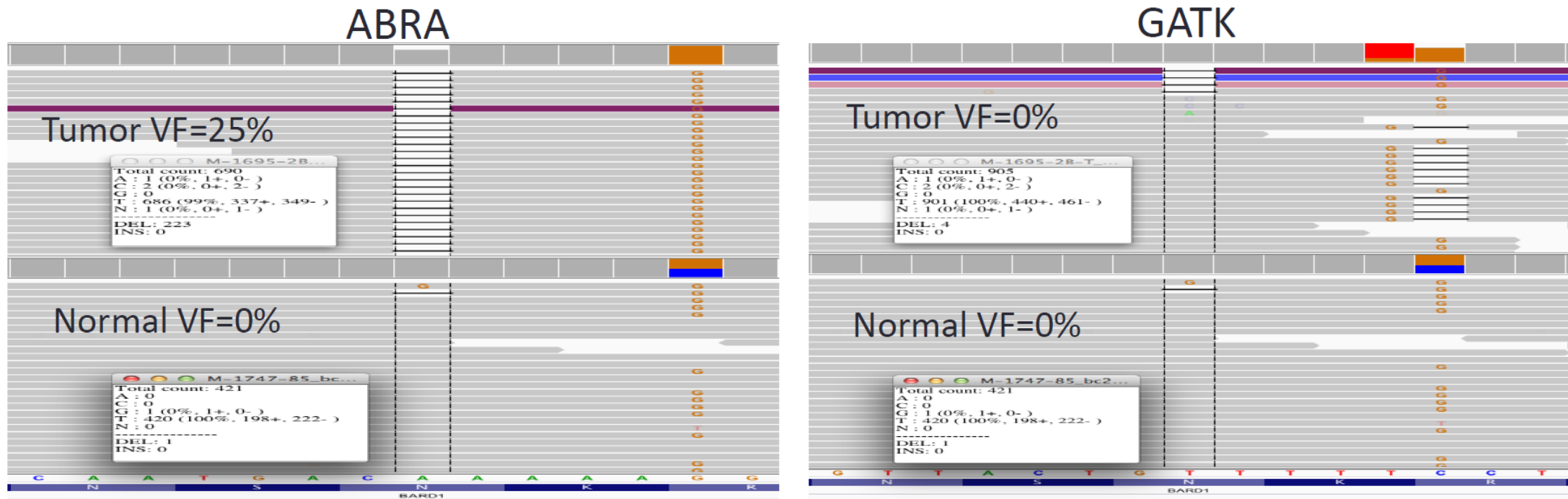**ABRA resolves poorly aligned regions**



Figure 8: A deletion and another mutation event called by GATK in the BARD1 gene was resolved far more cleanly by ABRA into a single deletion event with a separate SNV farther away.

**https://www.slideshare.net/rshah7/comparison-of-lumpy-vs-delly-for-structural-variant-detection**

HOKKAIDO UNIVERSITY

The poster content is shown below as a full-page figure.



# A Comparison of Genomic Structural Variant Detection using LUMPY and DELLY

Lance Tan[1], Ronak H. Shah[2], Michael F. Berger[2]

[1]Newark Academy, Livingston, NJ [2]Department of Pathology, Memorial Sloan Kettering Cancer Center, New York, NY

Memorial Sloan Kettering Cancer Center

**A** — Common and unique SVs with varying breakpoint difference threshold

- SVs found by both
- SVs found by LUMPY only
- SVs found by DELLY only

(Number of SVs vs Breakpoint difference threshold (bases): 10, 20, 30, 40, 50, 60, 70, 80, 90, 100)

**B**

LUMPY calls (780 total) — DELLY calls (1120 total)

685 | 95 | 1025

https://www.slideshare.net/rshah7/comparison-of-lumpy-vs-delly-for-structural-variant-detection

HOKKAIDO UNIVERSITY

**Figure 3. Venn diagrams summarizing called variants by different callers.** The mean percentage with standard deviation of confidence variant calls with equal to or higher than the quality score threshold of 20 are represented for (**A**) Illumina data sets and (**B**) Ion Proton data set.

# Variant Caller Comparison in WES

HOKKAIDO UNIVERSITY

**Fig 6. Variant caller sensitivity.** Variant caller sensitivity for detecting the manually curated mutations for SNVs and indels are shown in left and right panels, respectively. The y-axis depicts the number of variant calls. The dark and light grey bars represent calls in the exome and targeted deep sequencing data, respectively.

AB Krøigård et al., PLoS One. 2016; 11(3)

HOKKAIDO UNIVERSITY

# Error Rate Comparison among Platforms



Error Rates for Substitutions in R1 Reads

Error Rates for Substitutions in R2 Reads

HOKKAIDO UNIVERSITY

# Error Related Motif (3mers preceding errors)



(a) R1 Substitutions

HOKKAIDO UNIVERSITY

# Mutect Defaults Filter Settings

**Filters used in high-confidence mode**

    1. Proximal Gap

    2. Poor Mapping

    3. Strand Bias

    4. Clustered Position

    5. Observed in Control

**Filters applied in all MuTect modes**

    1. Tumor and normal LOD scores
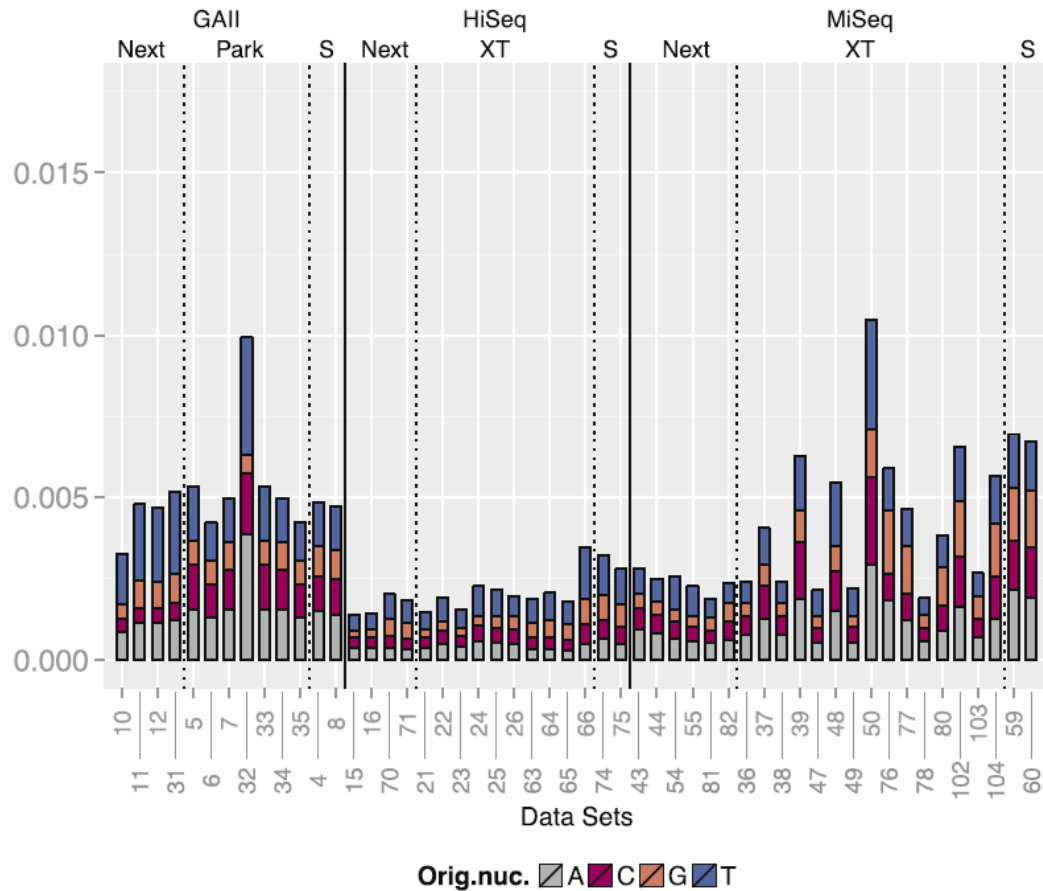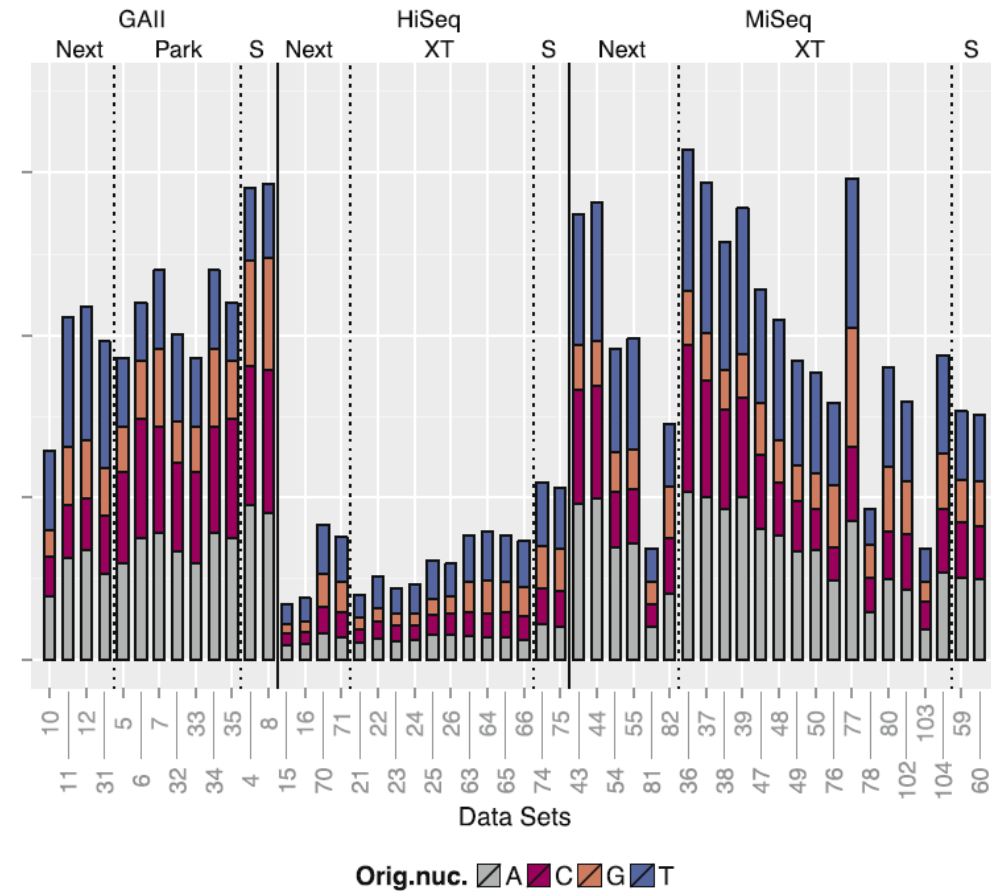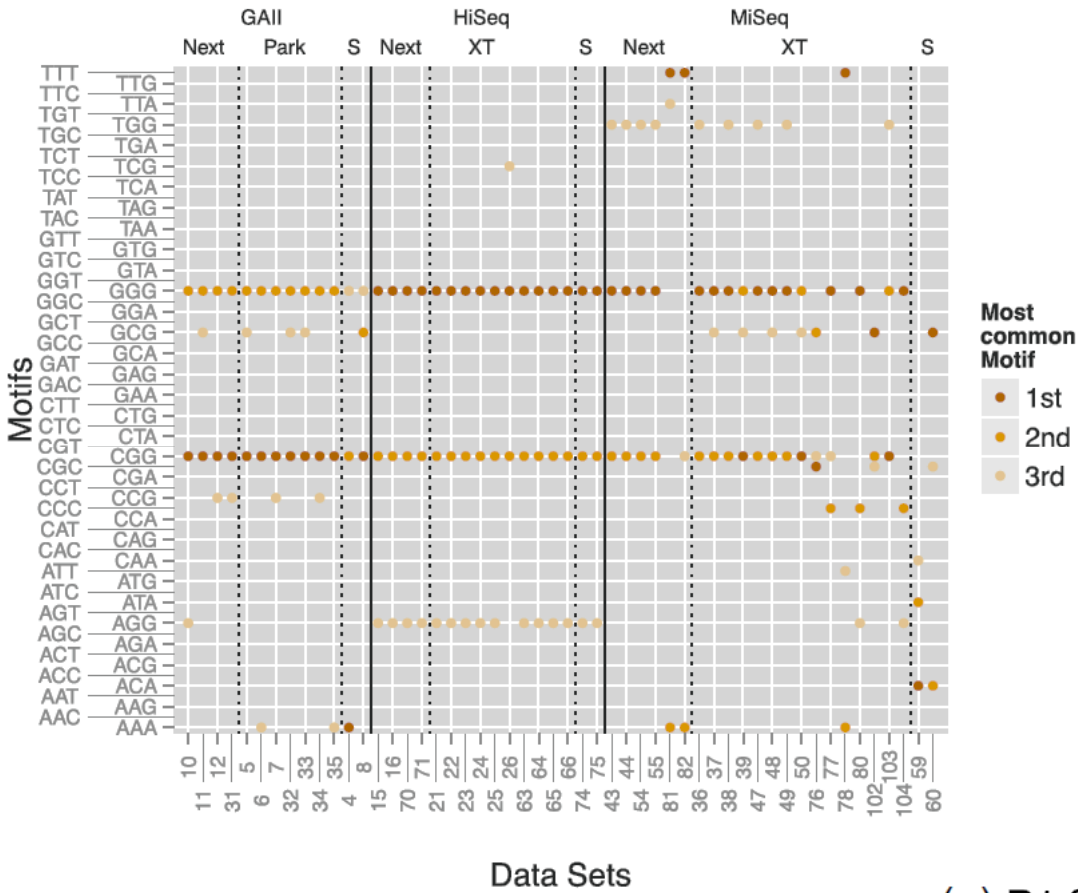
    2. Possible contamination

    3. Normal LOD score and dbsnp status

    4. Triallelic Site Filter

https://gatkforums.broadinstitute.org/gatk/discussion/4464/how-mutect-filters-candidate-mutations

HOKKAIDO UNIVERSITY

**Table 1.** Empirically derived filtering parameters for putative somatic mutations
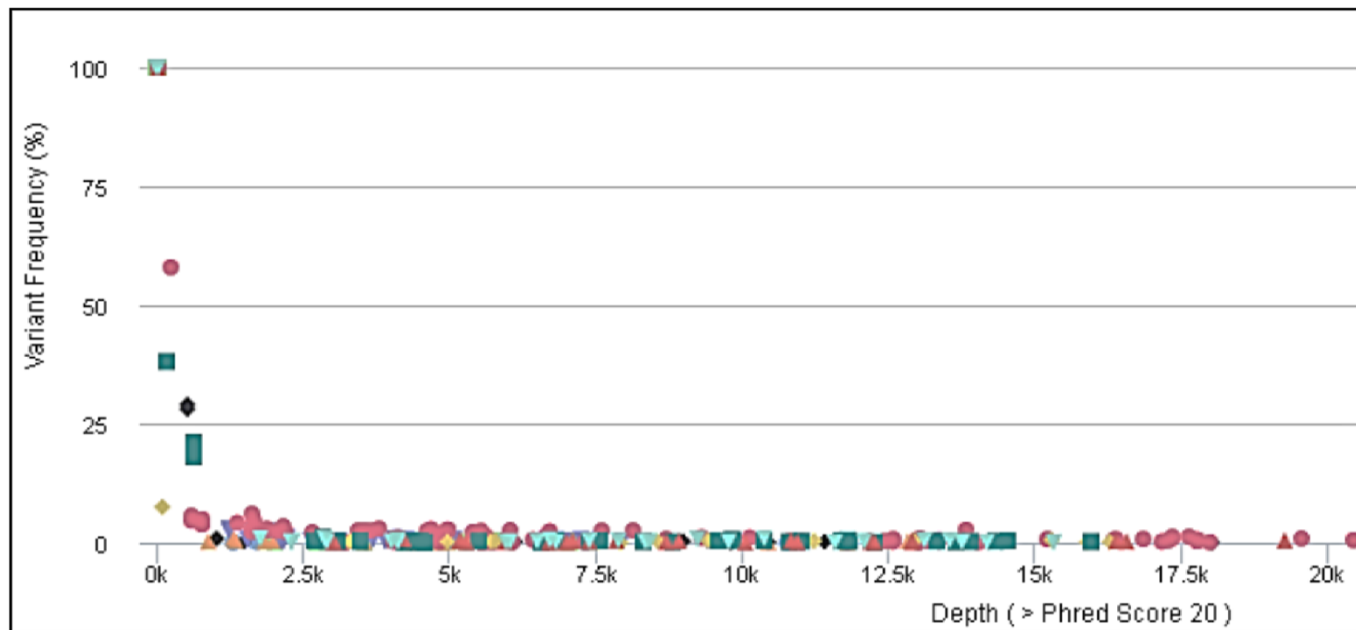
| Parameter | Description | Requirement |
|---|---|---|
| Read position | Average variant position in supporting reads, relative to read length | Between 10 and 90 |
| Strandedness | Fraction of supporting reads from the forward strand | Between 1%–99% |
| Variant reads | Total number of reads supporting the variant | At least four |
| Variant frequency | Variant allele frequency inferred from read counts | At least 5% |
| Distance to 3′ | Average distance to effective 3′ end of variant position in supporting reads | At least 20 |
| Homopolymer | Number of bases in a flanking homopolymer matching one allele | Less than five |
| Map quality difference | Difference in average mapping quality between reference and variant reads | Less than 30 |
| Read length difference | Difference in average trimmed read length between reference and variant reads | Less than 25 |
| MMQS difference | Difference in average mismatch quality sum between variant and reference reads | Less than 100 |

HOKKAIDO UNIVERSITY

# Variant Filtering in Another Cases

# Variant Filtering in Another Cases

http://www.cbioportal.org/
https://civic.genome.wustl.edu/home

HOKKAIDO UNIVERSITY

# Mannual Review

HOKKAIDO UNIVERSITY

https://genome.sph.umich.edu/wiki/Variant_Normalization

HOKKAIDO UNIVERSITY

# Normalization and Complex Variant

単純な例

複雑な例



Zook JM, Nat Biotechnol. 2014 Mar;32(3):246-51

HOKKAIDO UNIVERSITY