

Breaking the NGS Noise Barrier to Accurately Detect Variants below 1% Allele Frequency

Sukhinder Sandhu, Ashley Wood, Bitu Lahann, Jordan RoseFigura, Justin Lenhart, Vanessa Kelchner, Laurie Kurihara, Vladimir Makarov, Timothy Harkins

Swift Biosciences, 58 Parkland Plaza, Suite 100, Ann Arbor, MI 48103, Tel: 734.330.2568



Introduction

The use of circulating, cell-free DNA (cfDNA) for early detection and monitoring of disease is rapidly growing. This necessitates accurate variant detection at and below 1% allele frequencies due to a low population of tumor DNA within cfDNA samples. Reliable, low-frequency variant detection by next-generation sequencing (NGS) is challenging due to non-specific background noise from PCR and sequencing errors. We have employed molecular identifiers (MIDs) to uniquely label individual DNA molecules prior to amplification, facilitating the distinction of true variants from PCR and sequencing errors. MID incorporation also results in increased data retention by removing PCR duplicates while preserving fragmentation and sister strand duplicates.

Here, we have applied MIDs to both our Accel-NGS® 2S whole genome library prep followed by hybridization capture with various cancer panels, and to our Accel-Amplicon™ library prep that uses multiplex PCR ideal for small panels that are amenable to deep sequencing. To validate this technology, we performed low frequency spike-in experiments at 1% and 0.5% with genomic DNA, as well as various cfDNA samples. 2S libraries were prepared and hybridization capture was performed with an 800 kb pan-cancer panel. Libraries were sequenced to greater than 8000x coverage and a consensus sequence was generated with BMFtools. All known variants present at 1% and 0.5% allele frequencies were maintained in the resulting data. Further, true variants were preserved while PCR and sequencing errors were removed, demonstrating improved specificity using MIDs. We also validated variant calling below a 1% allele frequency using the Accel-Amplicon EGFR Pathway Panel with MIDs. After validation of our MID technology, libraries were prepared with cfDNA and tumor samples from individuals with ovarian, liver, stomach, and colon cancers. Libraries were sequenced to a minimum of 13,000x coverage, and we determined data retention after de-duplication with and without the use of MIDs. We observed a significant increase in data retention that led to a 2- to 5-fold increase in coverage using MIDs. Variant calling identified pathogenic mutations in all cfDNA samples, including those present in a corresponding tumor sample when available. This study highlights the ability of MID technology to enable variant detection at and below 1% allele frequencies, critical to track known variants and identify novel pathogenic mutations in cfDNA samples.

Improved Data Analysis with MIDs

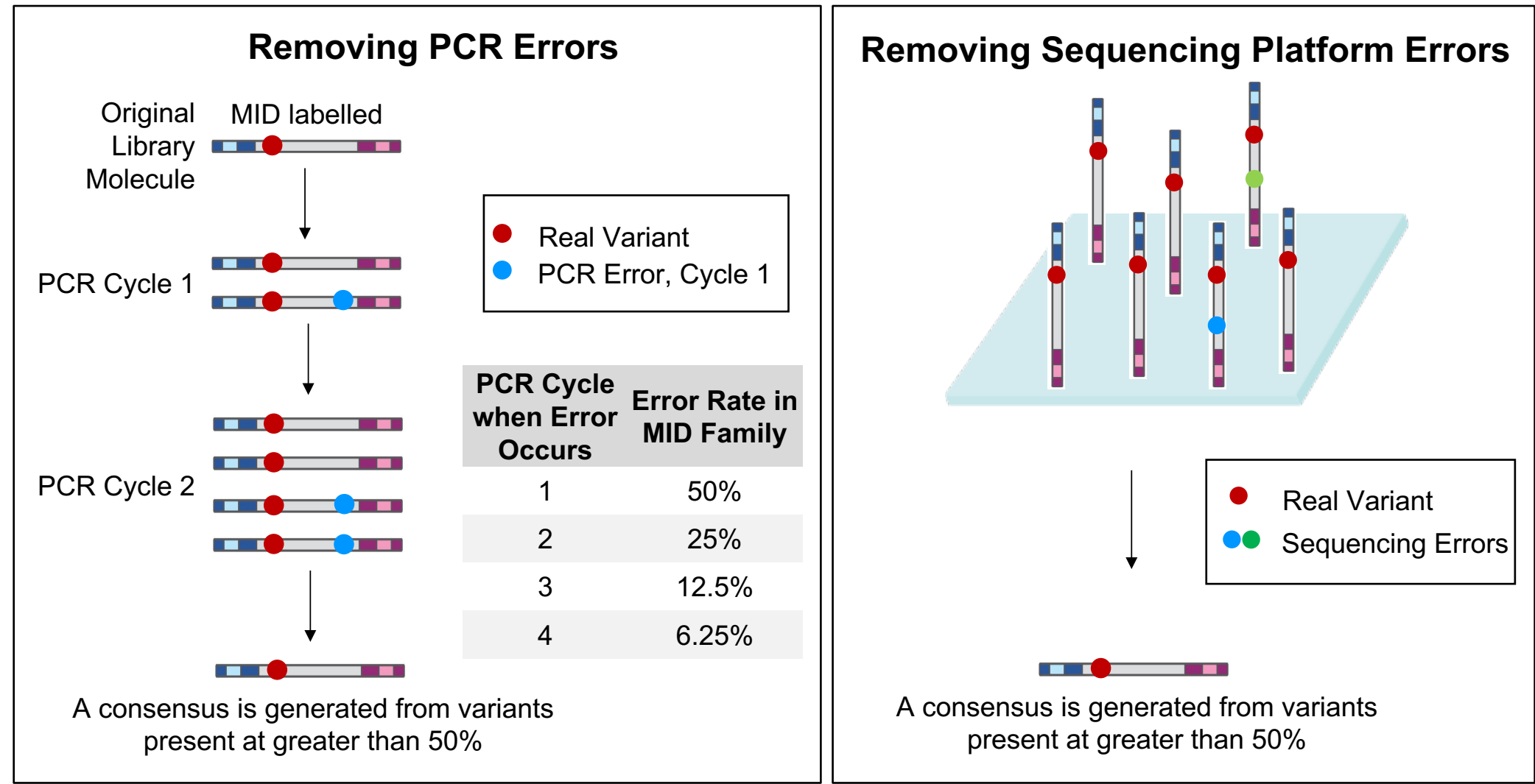


Figure 1. MIDs label individual molecules prior to exponential amplification by PCR facilitating the accurate identification and removal of PCR duplicates. Furthermore, molecules containing the same MID can be used to generate a consensus sequence that retains true variants but removes artificial mutations generated by polymerase errors during PCR amplification and sequencing. Here we depict PCR duplicates from one MID family to demonstrate that PCR and sequencing errors should not exist at greater than 50% and are therefore eliminated in the consensus sequence.

MIDs Label Unique Library Molecules

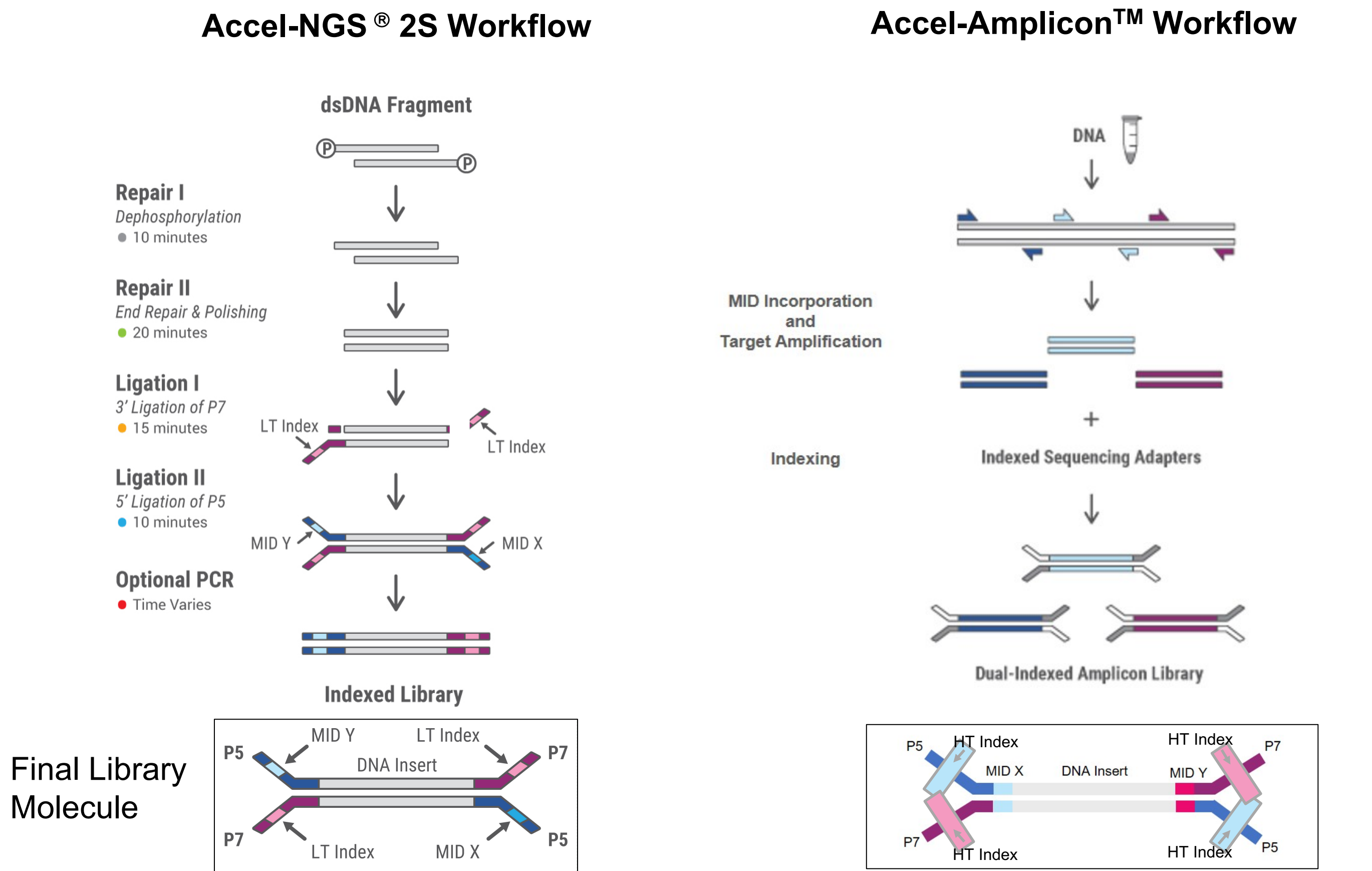


Figure 2. MIDs can be added to both our Accel-NGS 2S whole genome library prep used for hybridization capture and to our Accel-Amplicon library prep that uses multiplex PCR for targeting. Both library preps use the MID to label unique library molecules prior to amplification. Accel-NGS 2S libraries with MIDs are constructed using Illumina®-compatible adapters with a strand-specific 9 base random MID in the i5 index position and a sample index in the i7 position. Each dsDNA substrate receives two independent MID adapters “X” and “Y”. The Accel-Amplicon workflow consists of a 2-cycle MID incorporation step, PCR amplification of targeted amplicon molecules, and an indexing step. The final library molecule consists of two 5-base random N sequences positioned at the start of Read 1 and Read 2. Each original DNA molecule receives a unique MID X and MID Y, yielding a 10 base in-line MID/barcode.

Methods

1. Input Considerations for Low Frequency Allele Detection

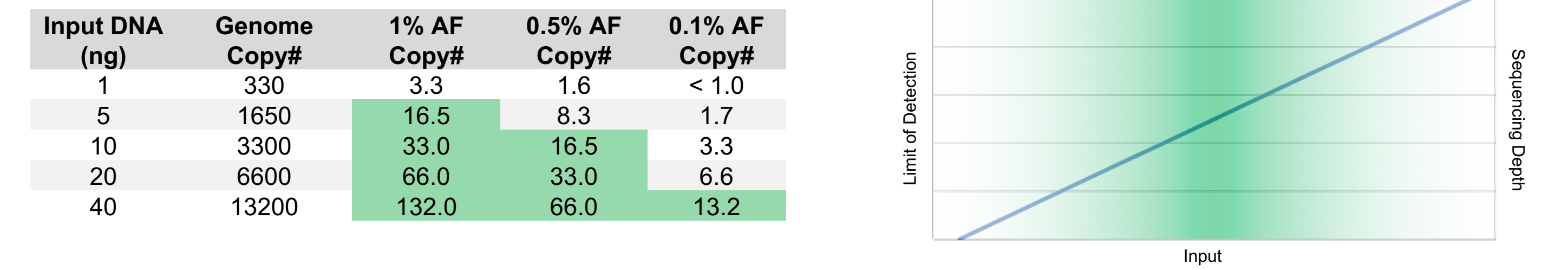


Figure 3. In order to ensure representation of DNA fragments carrying low frequency alleles within a sample, sufficient genome copy number is needed. Sample inputs with a predicted low frequency allele copy number of < 10 may not be present in the DNA sample due to Poisson distribution of DNA fragments in solution. However, higher input genome copy number is proportional to a higher library complexity and impacts the depth of sequencing required. Therefore, a balance between the limit of detection (LOD) and the required sequencing depth must be considered to determine the appropriate input for an experiment. Recommended inputs for specified LOD allele frequency are highlighted in green.

2. Sequencing Depth Considerations for Complexity

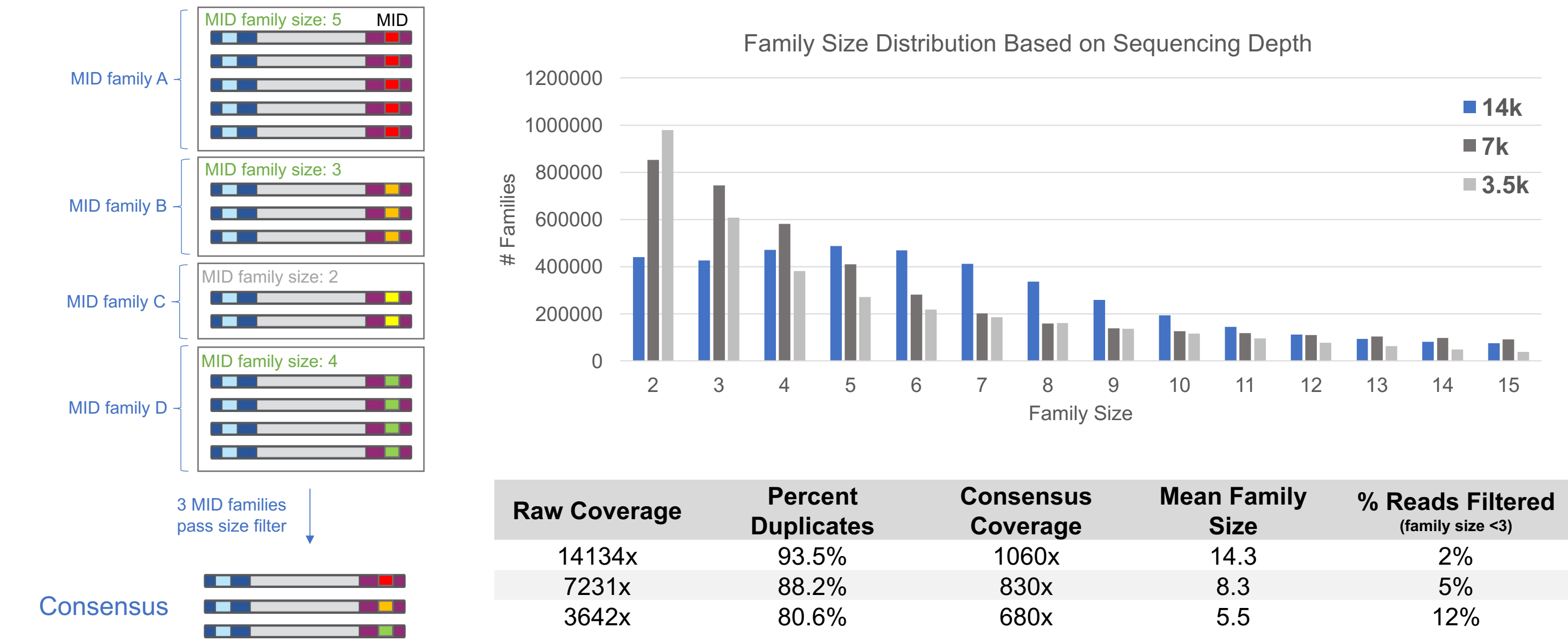
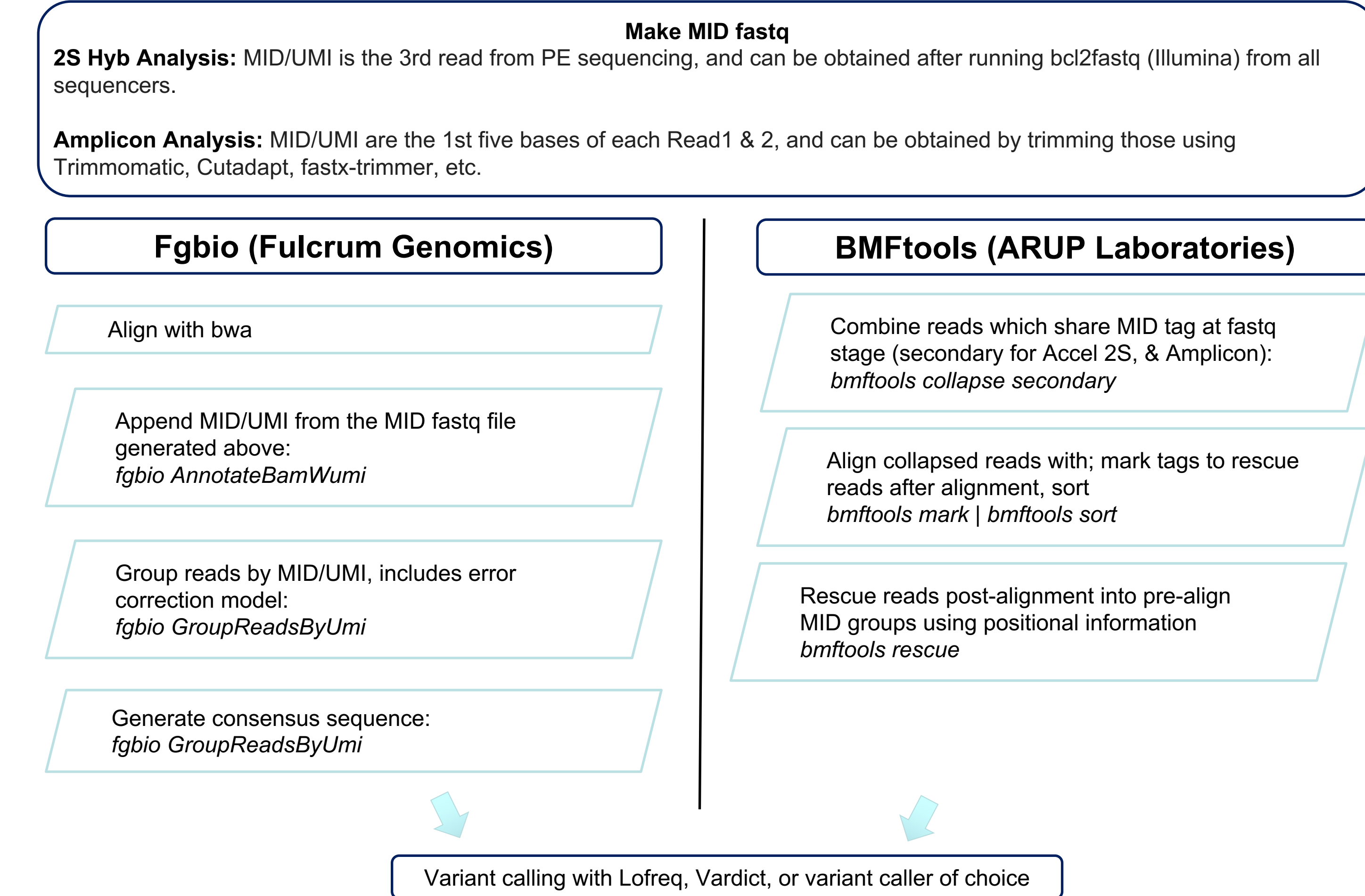


Figure 4. Identification of low frequency variants requires that at least three PCR duplicates (Family Size ≥ 3) are sequenced per MID family to enable consensus sequence analysis since only alleles present in the majority of PCR duplicates are retained in the consensus. To achieve this for the majority of molecules, this requires sequencing DNA libraries to near saturation (i.e., every unique molecule from the library is sequenced more than once). Here we show how sequencing to 14000x, 7000x, and 3500x influences MID family size for a library prepared from 20ng cfDNA and enriched using the 800kb IDT xGen® Pan-Cancer Panel.

3. Bioinformatics Workflow



Results

1. Increased Data Retention with MIDs

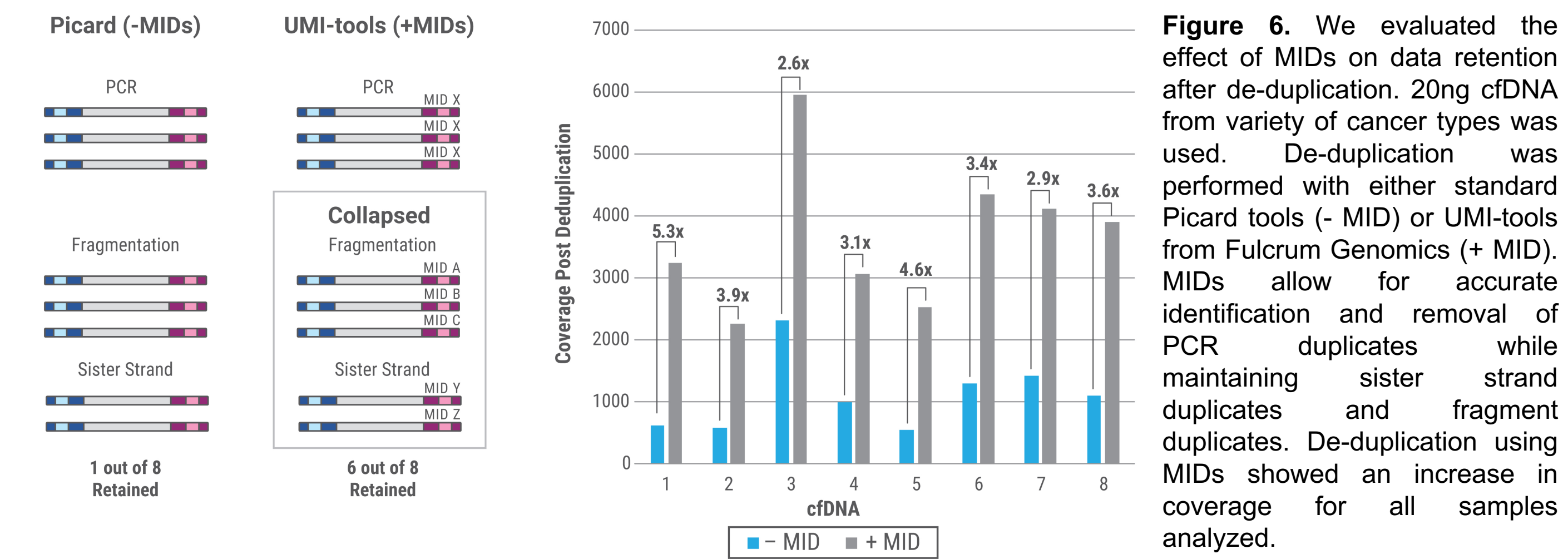


Figure 6. We evaluated the effect of MIDs on data retention after de-duplication. 20ng cfDNA from variety of cancer types was used. De-duplication was performed with either standard Picard tools (- MID) or UMI-tools from Fulcrum Genomics (+ MID). MIDs allow for accurate identification and removal of PCR duplicates while maintaining sister strand fragments and fragment duplicates. De-duplication using MIDs showed an increase in coverage for all samples analyzed.

2. Identification of Variants Down to 0.5%

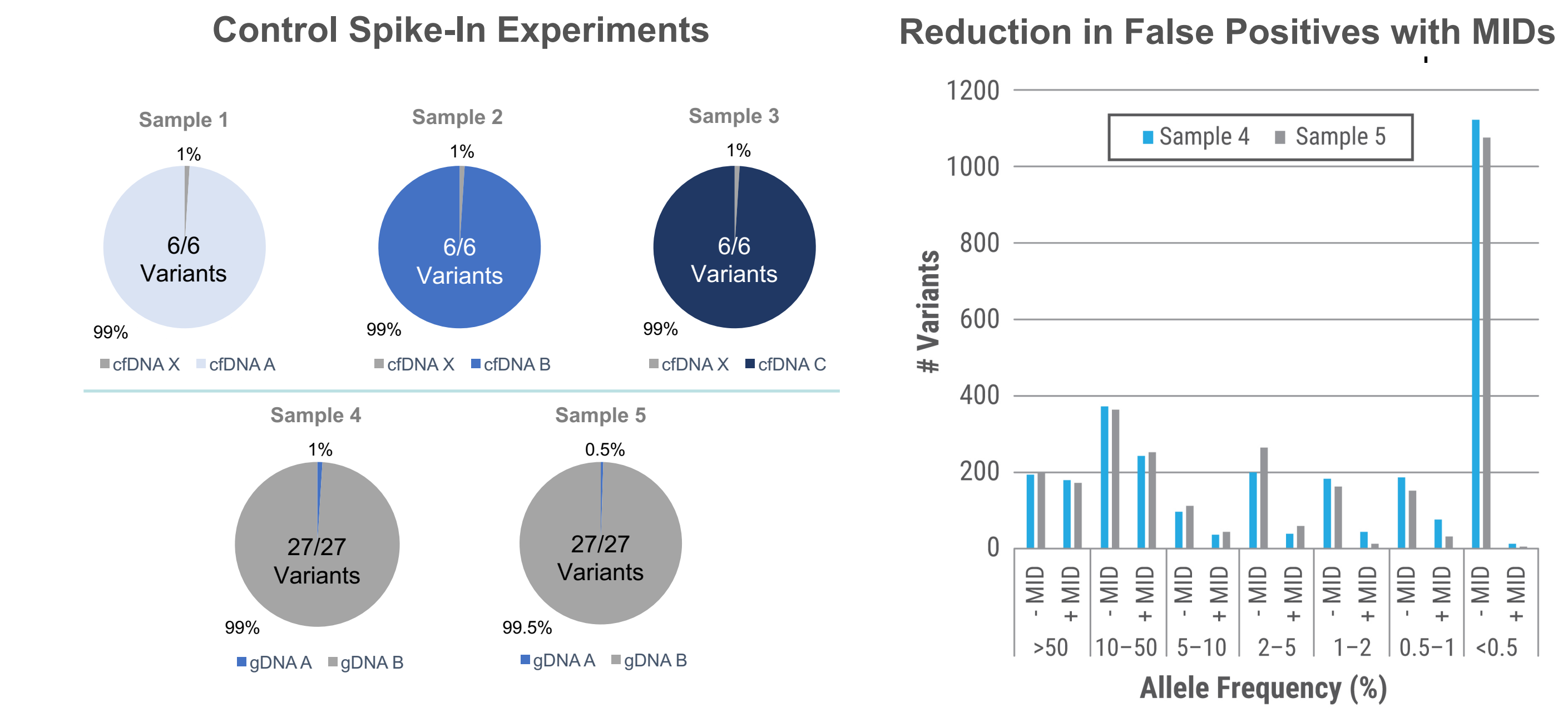


Figure 7. (Left) cfDNA was extracted from blood of four individuals with unique genetic backgrounds and gDNA samples from different genetic backgrounds were obtained (Coriell Institute). Sample spike-ins were performed at 1% or 0.5% frequency into 10 ng cfDNA or 100 ng gDNA. Libraries were prepared with Swift's Accel-NGS 2S Hyb Kit with MIDs, enriched with the IDT xGen Pan-Cancer Panel and sequenced on an Illumina HiSeq® to a minimum of 8000x coverage. Data from 2S Hyb spike-in experiments were analyzed with BMF tools and data were analyzed for homozygous SNPs present in the spike-in sample only. 6/6 known variants were present in all three 1% cfDNA samples and 27/27 known variants were present in both 1% and 0.5% gDNA samples. (Right) Total variants called at various allele frequencies with or without the use of MIDs are depicted from the gDNA spike-in experiments. MIDs only have a subtle effect on the number of variants called at high allele frequencies, but substantially reduce the number of low frequency variants called. This is the result of removing sequencing and PCR errors such that variants called are highly enriched for true variants and the removed variants represent noise.

3. EGFR MID Panel Performance

Variant Calling with the EGFR-MID Amplicon Panel

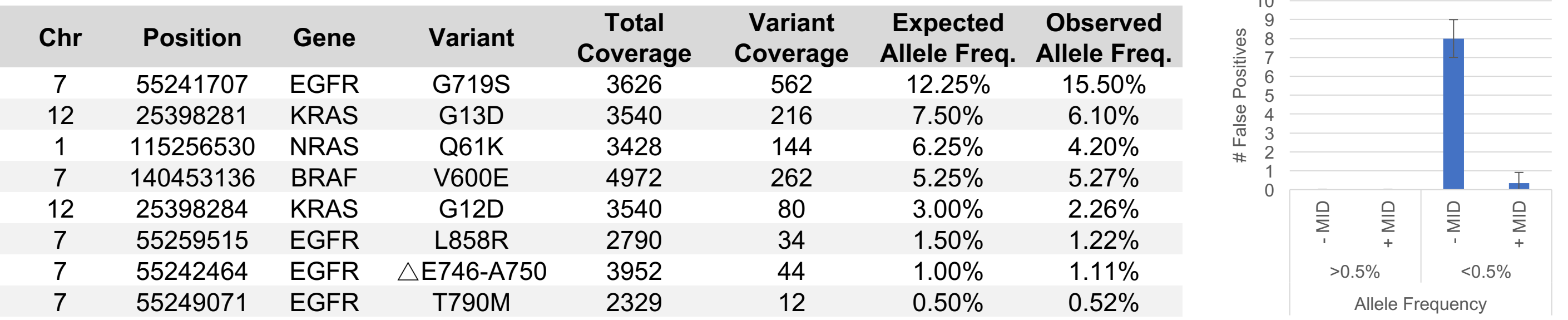


Figure 8. The EGFR MID panel contains 17 amplicons, 107-149bp in size, and shows >95% coverage uniformity (>0.2x mean) and >95% on target reads. Horizon Diagnostics Quantitative Multiplex DNA Standard (HD701) was spiked into Coriell DNA (NA12878) at 50% to obtain expected variants at allele frequencies from 12.25-0.50%. Libraries were prepared using 10 ng of input DNA and the EGFR MID panel. Sequencing was performed on an Illumina MiniSeq® to greater than 100,000x. All expected variants were consistently detected in the consensus sequence and the use of MIDs removed false positives at low allele frequencies. The graph to the right depicts the average number of false positives (n=3) called by LoFreq with and without the use of MIDs.

Conclusion

- We've validated the input DNA quantity and depth of sequencing required to maximize data retention and the LOD when performing MID analysis. We recommend using the minimal input quantity for the desired LOD in order to minimize sequencing cost.
- We have successfully applied Fgbio and BMFtools for MID analysis. Both produce error free consensus bam, compatible with most MID agnostic tools. BMFtools provides more exploratory options and a plethora of summary metrics and filtering options.
- Incorporating MIDs in NGS library preparation increases data retention during de-duplication and improves variant calling by increasing specificity at low allele frequencies. Here we are able to detect known variants at less than 1% allele frequencies using hybridization capture and amplicon approaches for targeted NGS.