

# Leveraging Molecular IDs to Improve Coverage of Low Frequency Variants in Ovarian Cancer Using Next Generation Sequencing

Andrew McUsic<sup>1</sup>, Ashley Wood<sup>1</sup>, Matthew Dashkoff<sup>2</sup>, Timothy Harkins<sup>1</sup>, Sukhinder Sandhu<sup>1</sup>, Olga Camacho-Vanegas<sup>2</sup>, Peter Dottino<sup>2</sup>, Laurie Kurihara<sup>1</sup>, John Martignetti<sup>2</sup>, Vladimir Makarov<sup>1</sup>

<sup>1</sup>Swift Biosciences, 58 Parkland Plaza, Suite 100, Ann Arbor, MI 48103, Tel: 734.330.2568

<sup>2</sup>Icahn School of Medicine at Mount Sinai, Departments of Genetics and Genomic Sciences and Obstetrics/Gynecology & Reproductive Sciences, 1425 Madison Avenue New York, NY 10029

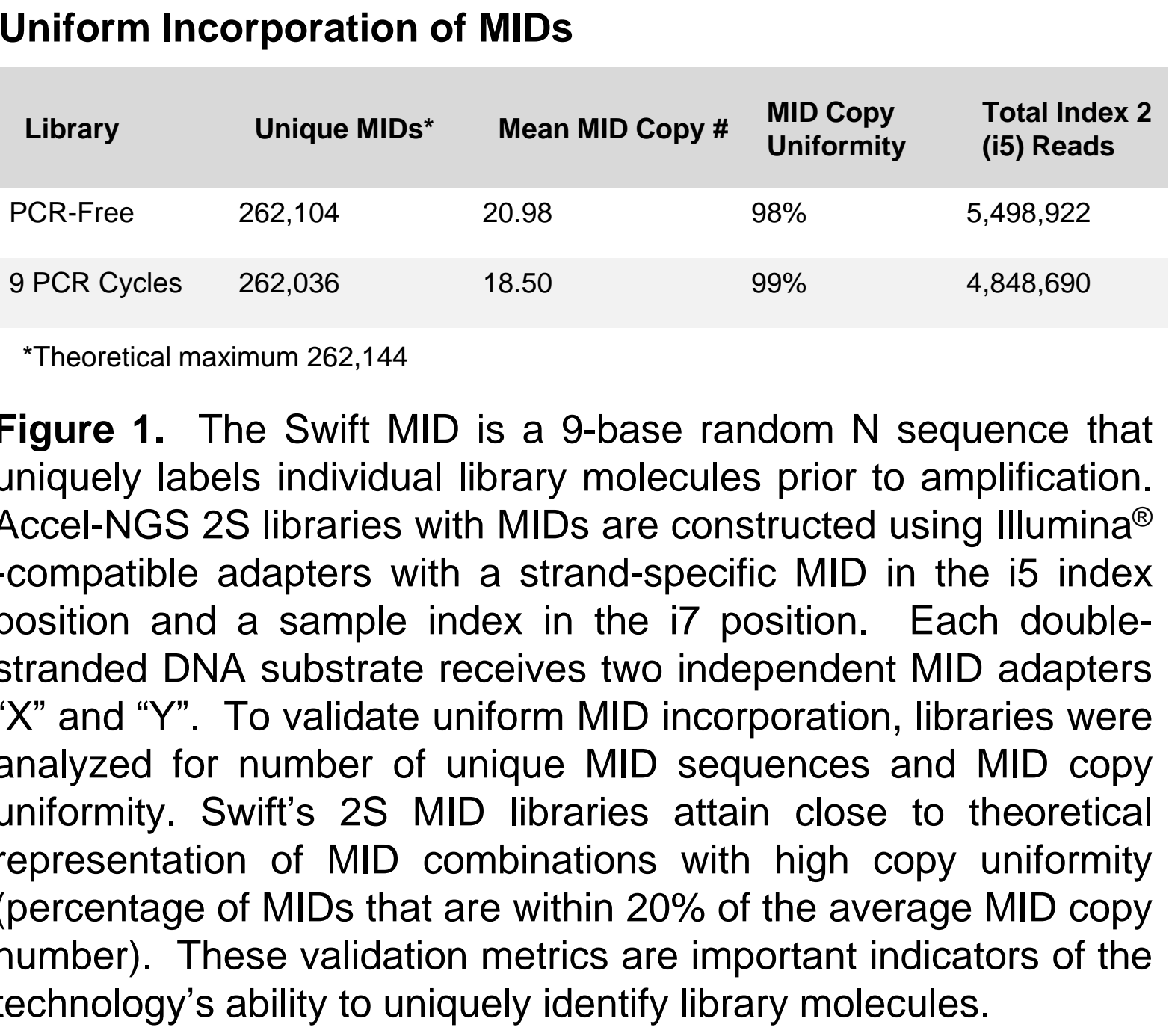
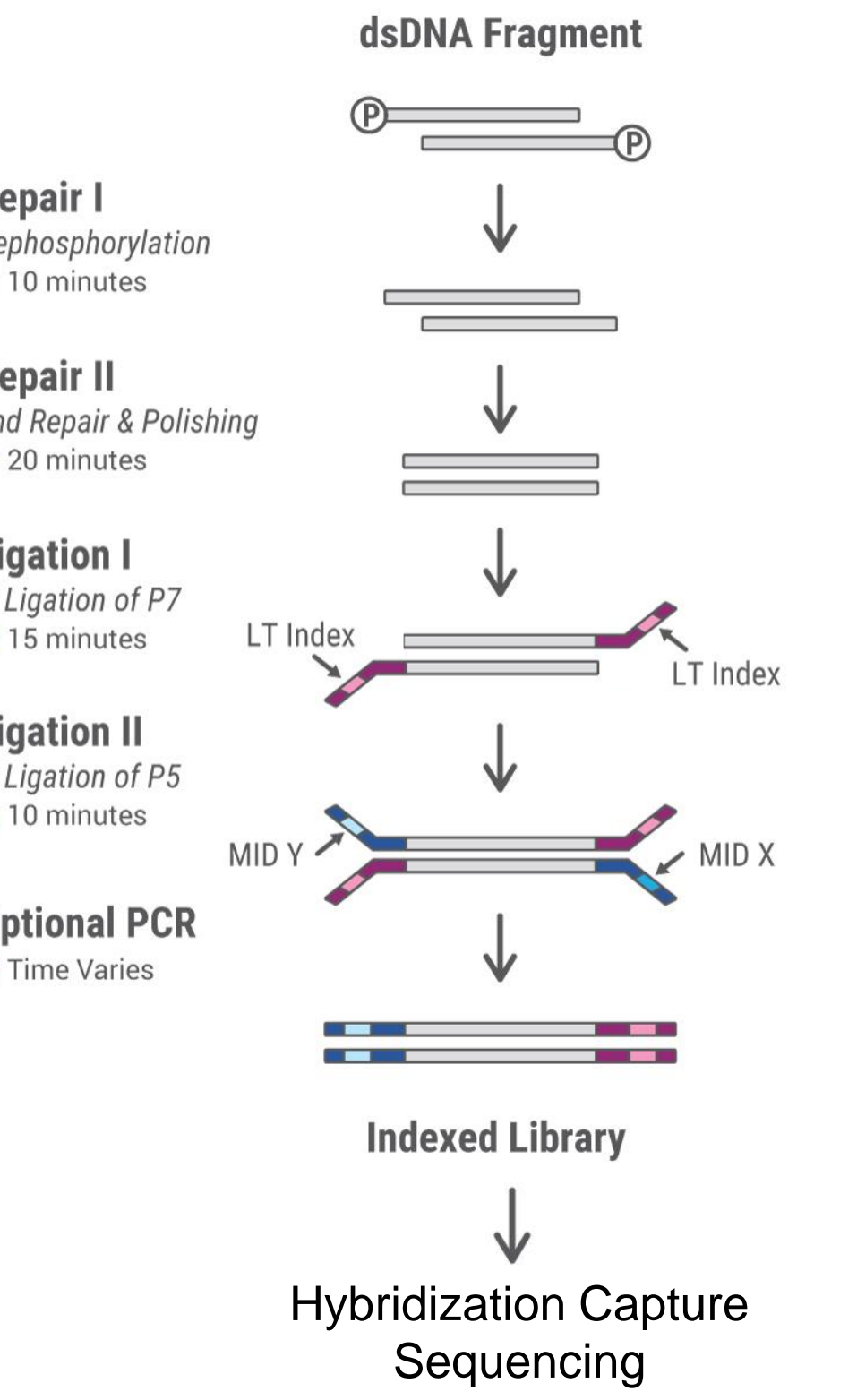
## Abstract

**Introduction:** The detection of low frequency allele variants by next generation sequencing (NGS) holds great promise for elucidating the progression of cancer and other pathologies but also presents significant challenges for sensitivity and specificity. Below a threshold of 5% allelic fraction, abundance levels of mutations approach the assay limit of detection, and true positives are difficult to distinguish from false positives due to errors in both the amplification methods used to produce the NGS library as well as read errors inherent to the sequencing platform. Without a specific technology for uniquely labeling each library molecule, it is difficult to distinguish duplicate molecules produced by PCR from unique templates or otherwise understand the contribution of each clonal family of library molecules, especially for liquid biopsy samples where positional information is not a reliable method for distinguishing fragments due to common sites for enzymatic shearing of circulating, cell-free DNA (cfDNA). The use of Molecular IDs (MIDs) to label unique library molecules provides a method for detecting and removing PCR duplicates from the sequencing data while preserving unique information otherwise conventionally discarded during deduplication due to the presence of duplicates arising from fragmentation and complementary strands. Leveraging these MIDs allows the generation of consensus sequence from PCR duplicates to reduce PCR and sequencing errors and also increases confidence in the detection of low frequency mutations by quantifying the presence of mutations in multiple unique molecules.

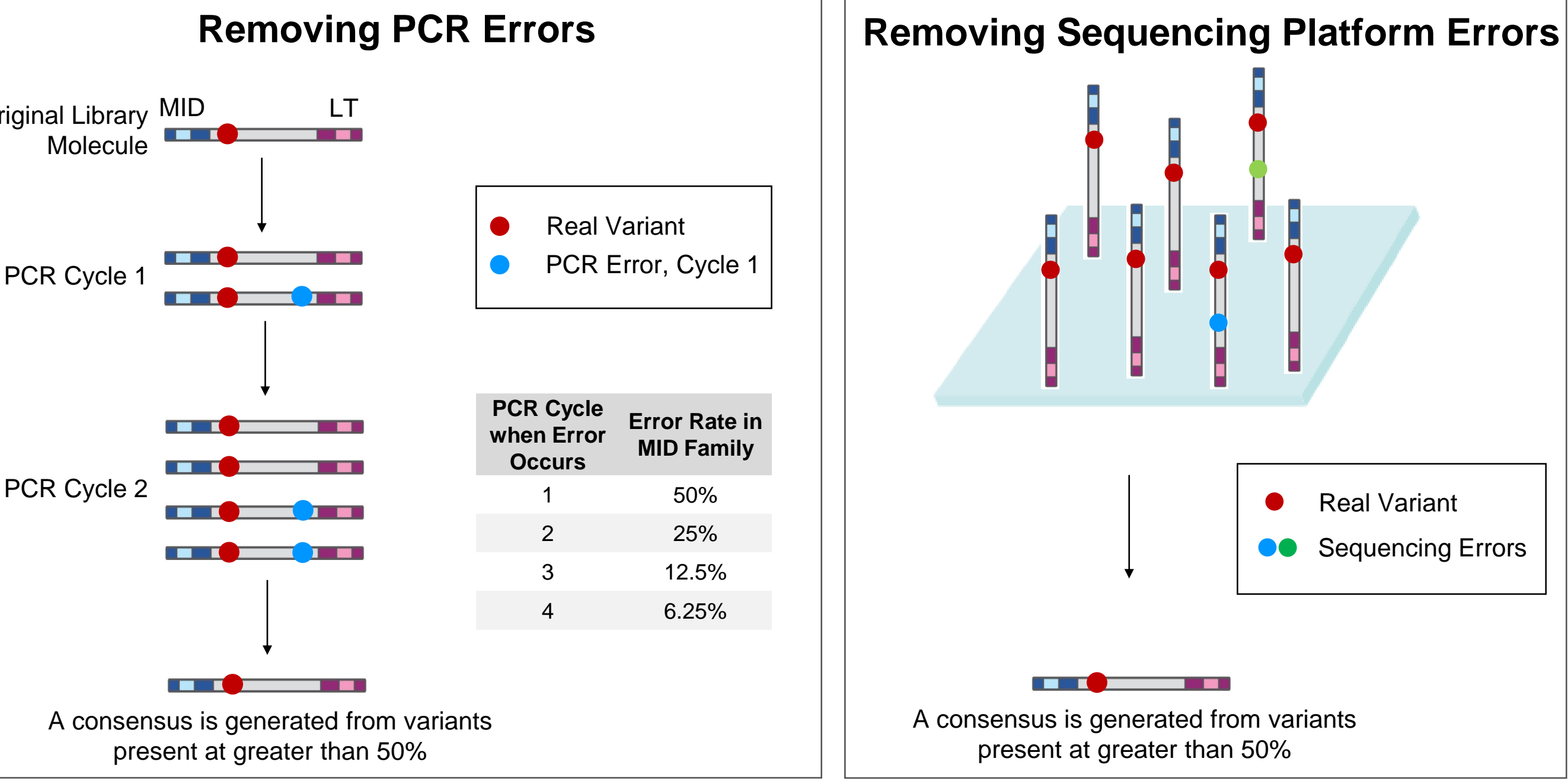
**Methods:** Tumor samples were biopsied from a stage 3B, grade 3 ovarian carcinosarcoma (malignant mullerian mixed tumor from 75-year-old Caucasian female) and cfDNA was obtained from blood plasma collected at the time of surgery. The Accel-NGS® 2S Hyb DNA Library Kit from Swift Biosciences was combined with MID technology and a pan-cancer hybridization capture enrichment panel to prepare DNA isolated from the series of metastatic tumors for sequencing by Illumina® MiSeq® and HiSeq®. Data retention was defined as the increase in useful coverage obtained when considering MIDs versus removing duplicates without consideration of the MIDs. We determined data retention after de-duplication by comparing results from standard Picard tools without the use of MIDs to UMI-tools (Fulcrum Genomics) to utilize MIDs. Variant calling was performed using BMFtools, LoFreq, and VarDict.

**Results:** When sequencing the panel to an average depth of 200-300X by MiSeq, we observed an increase in data retention for both Covaris-sheared FFPE and cfDNA libraries that led to a 1.1-1.3 fold increase in coverage using UMI-tools over Picard de-duplication. This effect was most pronounced at lower DNA input quantities and with higher depth of sequencing. When sequencing the same libraries to an average depth of 1000-2000X by HiSeq, we observed a 1.7-2.9 fold increase in coverage when leveraging MIDs. These libraries were further evaluated for variant calling, and we identified variants with as low as 1% frequency present across all tumor samples and within cfDNA, as well as variants unique to each tumor sample.

## MIDs Label Unique Library Molecules

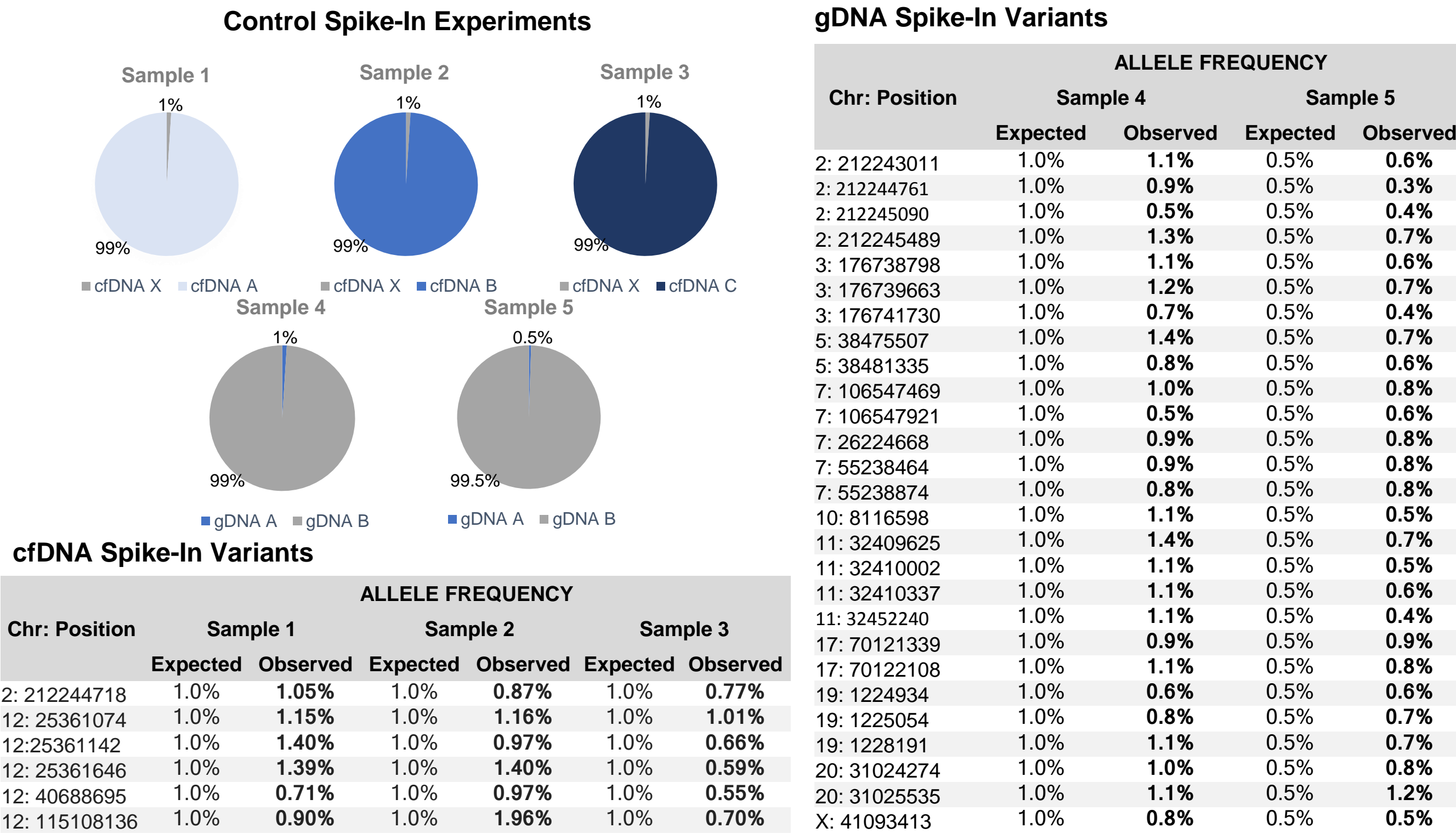


## Improved Data Analysis with MIDs



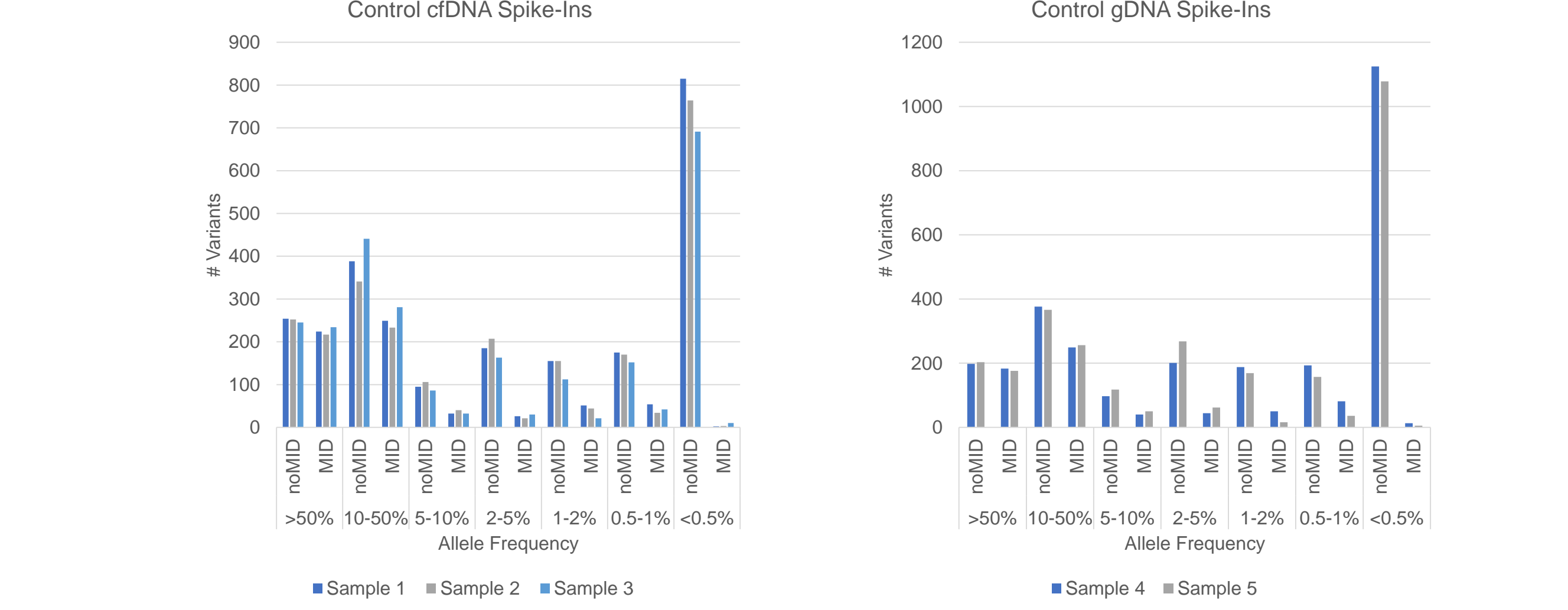
**Figure 2.** MIDs label individual molecules prior to PCR, facilitating the accurate identification and removal of PCR duplicates during data analysis. Furthermore, molecules containing the same MID can be used to generate a consensus sequence that retains true variants but removes artificial mutations 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.

## Identification of Variants Down to 0.5%



**Figure 3.** cfDNA was extracted from blood of four individuals with unique genetic backgrounds and Coriell gDNA samples from different genetic backgrounds were obtained. To determine the effect of MIDs on low frequency variant calling, sample spike-ins were performed at 1% or 0.5% frequency into 10 ng cfDNA or 100 ng gDNA. Libraries were prepared with the Accel-NGS 2S Hyb kit with MIDs, enriched with the IDT xGen® Pan-Cancer Panel that covers an 800 kb target containing 127 genes, and sequenced on an Illumina HiSeq® to a minimum of 8000x coverage. A consensus sequence was generated for each MID family (BMFtools) 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 depicting the power of MIDs for low frequency variant calling.

## Increased Specificity with MIDs



**Figure 4.** Total variants called at various allele frequencies with or without the use of MIDs are depicted from the 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. In this way MIDs lead to increased specificity in low frequency variant calling.

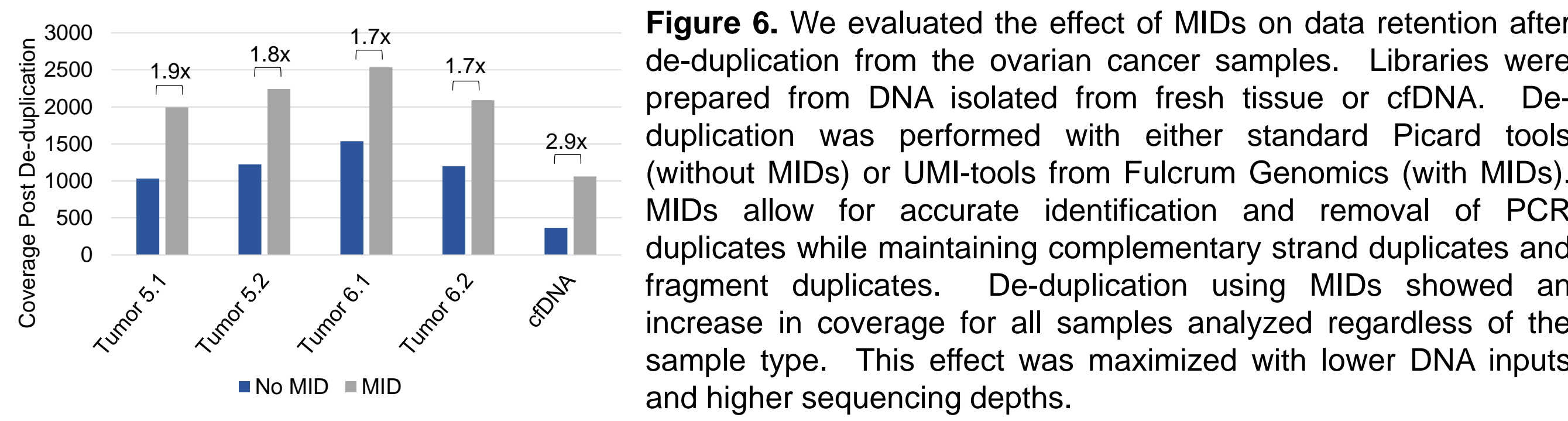
## Variant Analysis from a Single Individual

Sample	Description	Final Pathology	Input	Reads	Average Bait Coverage (Pre-deduplication)	Duplication Rate (Without MID/With MID)	On Target
Tumor 5.1	Small bowel implant (site 1)	malignant	50ng	93,536,722	15,634x	89.7% / 84.1%	74.2%
Tumor 5.2	Small bowel implant (site 2)	malignant	50ng	94,802,633	15,962x	89.1% / 83.4%	74.4%
Tumor 6.1	Left pelvic peritoneum (site 1)	benign	50ng	89,990,124	15,096x	86.9% / 81.0%	74.6%
Tumor 6.2	Left pelvic peritoneum (site 2)	benign	50ng	85,759,502	14,387x	88.3% / 82.7%	74.0%
cfDNA	Plasma cfDNA	N/A	20ng	81,387,833	14,135x	93.5% / 88.6%	74.4%

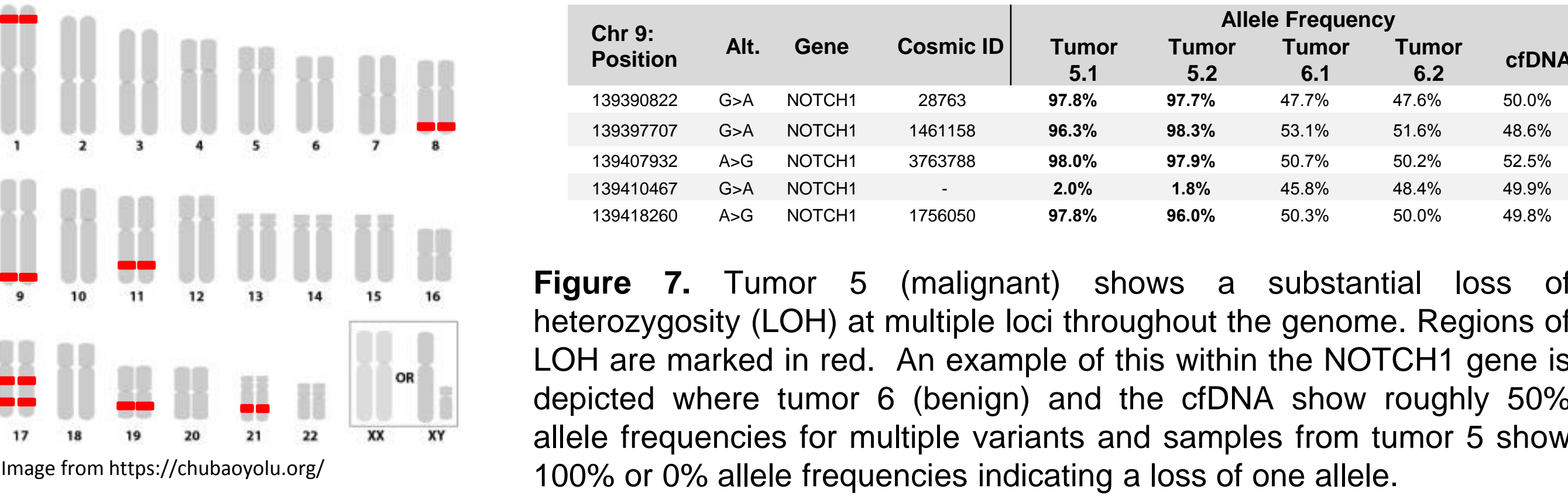
**Figure 5.** Samples were analyzed from a 75-year-old female with stage 3B, grade 3 ovarian carcinosarcoma. Samples from two different sites on two tumors were taken during recurrent surgery performed 9 months after the primary surgery. In addition, plasma cfDNA collected at the time of the recurrent surgery was analyzed. Libraries were prepared using an Accel-NGS 2S Hyb kit with MIDs and enriched for oncology-related genes and hotspots with the IDT xGen Pan-Cancer Panel. Sequencing was performed on an Illumina HiSeq to greater than 14,000x. Deep sequencing maximized the number of PCR duplicates for each molecule used to generate a consensus sequence.



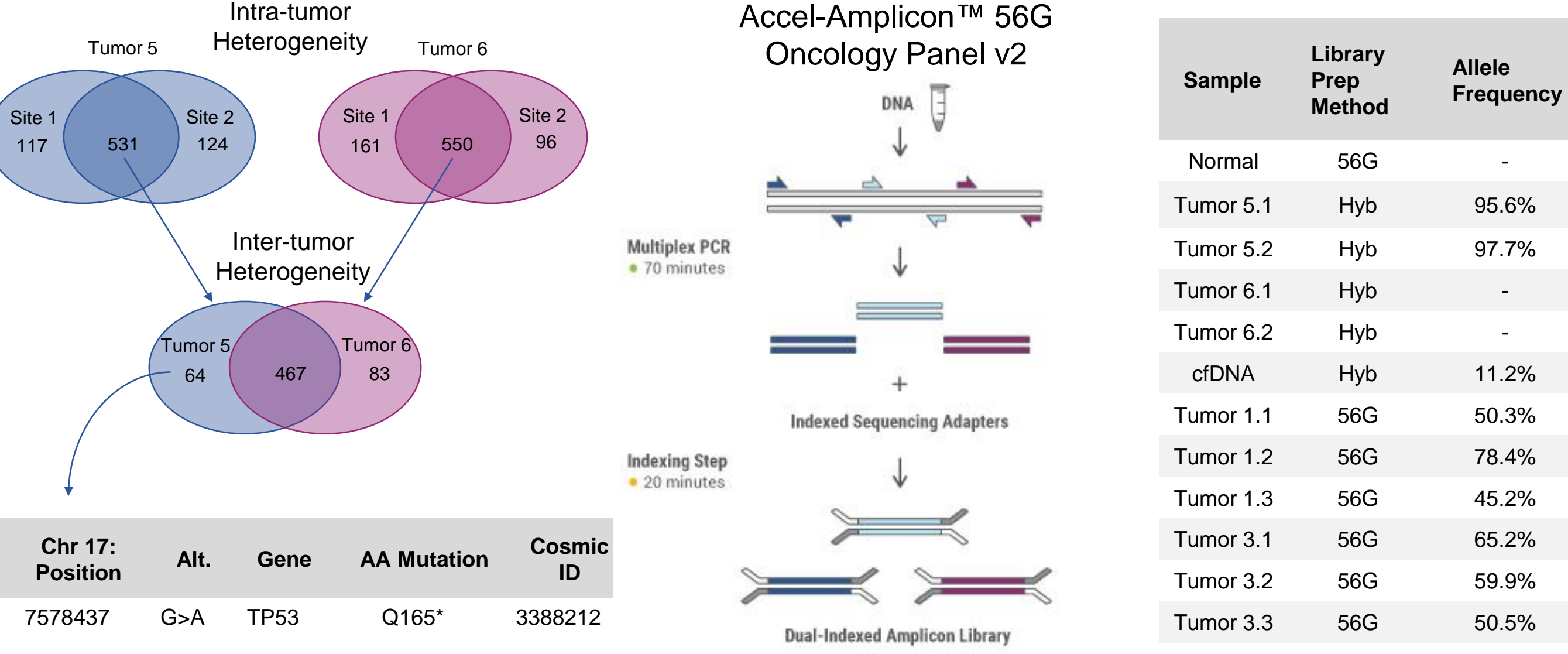
## 2-Fold Increase in Data Retention with MIDs



## Loss of Heterozygosity in Ovarian Cancer Samples



## Tumor Heterogeneity in Ovarian Cancer Samples



## Conclusions

- Labelling unique library molecules with MIDs prior to amplification allows for the removal of sequencing and PCR induced errors during data analysis.
- MIDs improve low frequency variant calling by increasing sensitivity and specificity and are shown here to facilitate detection of known variants present at 1% and 0.5% allele frequencies.
- The use of MIDs for de-duplication results in increased data retention through the accurate distinction of PCR duplicates from fragmentation and complementary strand duplicates.