



A tale of two platforms: An evaluation of the Roche GS Junior and Illumina® MiSeq next-generation sequencing instruments for forensic mitochondrial DNA analysis

Brittania J. Bintz, M.S.; Erin S. Burnside, M.S.; Mark R. Wilson, Ph.D.

Forensic Science Program, Department of Chemistry and Physics – Western Carolina University

ABSTRACT

Next-generation sequencing (NGS) refers to a suite of technologies that enable cost-effective, rapid generation of large amounts of detailed sequence information from clonal populations of individual template molecules. These methods are proving to be particularly well-suited for mitochondrial DNA analysis, and may provide forensic DNA analysts with a powerful tool that enables deconvolution of mtDNA mixtures. Recently, Illumina® has been working with members of the community to establish a human mtDNA forensic genomics consortium (IFGC) for concerted evaluation of NGS methods for potential use in mtDNA casework and databasing. In June, a set of samples was prepared consisting of quantified buccal extracts from two donors, as well as a series of mixtures of the buccal extracts at defined ratios (5, 2, 1 and 0.5%). This sample set has been distributed to participating IFGC laboratories for sequencing on multiple NGS platforms including the Ion PGM™, Roche GS Junior, and Illumina® MiSeq, to enable a cross-laboratory comparison of sequencing methods using identical samples. In our laboratory, the samples were sequenced on both the Roche GS Junior, and Illumina® MiSeq NGS platforms. Libraries from hypervariable regions 1 and 2 (HV1 and HV2) were sequenced on the Roche GS Junior using an amplicon library preparation approach where PCR primers were designed to include required adaptors and multiplexing indexes. For sequencing on the Illumina® MiSeq, libraries were prepared using Nextera® XT in which two large amplicons covering the whole mtGenome as well as HV1 and HV2 amplicons were randomly fragmented, and adaptors and indexes incorporated enzymatically. The resulting data was analyzed using CLC Genomics Workbench software and variant calls were compared. The Illumina® MiSeq resulted in significantly higher coverage across all positions sequenced, giving rise to higher certainty with low-level variant calls. Further, the MiSeq allowed for detection of minor variants in all mixtures where the majority of minor variants were undetected in the 0.5% mixture with the Roche GS Junior. Finally, data from the MiSeq showed lower background noise overall, especially in homopolymeric regions when compared to data from the GS Junior. The Illumina® MiSeq offers a streamlined enzymatic library preparation approach, higher-throughput and more accurate variant detection and bascalling than the Roche GS Junior. As a result, we feel that the MiSeq is better suited for forensic mtDNA analysis in both casework and databasing laboratories.

SAMPLE SET

Buccal swabs (x20) were obtained from two donors (001-CF30 and 003-CM54) whose whole mtGenome had been previously characterized in our laboratory using Sanger methods. The swabs were collected according to approved IRB protocol.

DNA from each set of 20 buccal swabs was extracted independently using the Qiagen DNA Investigator kit surface and buccal swab protocol. A single RB was extracted alongside each set. The resulting extracts from each donor were pooled to create a large volume master extract. Pooled samples were quantified in triplicate using a human mtDNA specific real-time PCR assay developed by Mark Kavliak et al.¹ Quantitative values were averaged after donors were removed. Averages were used to prepare mixtures of donors in defined ratios of 5, 2, 1 and 0.5%. Donor 003-CM54 was used as the minor contributor in all mixtures. Final mixtures and sole source samples were quantified in triplicate using both the Quantifiler® Human kit, and the human mtDNA specific real-time PCR assay mentioned above. All quantified samples were distributed to IFGC laboratories for sequencing as requested.

001-CF30 Sole Source	003-CM54 Sole Source	005-CM54 5:1	010-CM54 2:1	020-CM54 1:1	050-CM54 0.5:1	Region Roche	997A Primer Control

SEQUENCING CHEMISTRIES

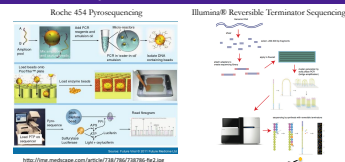
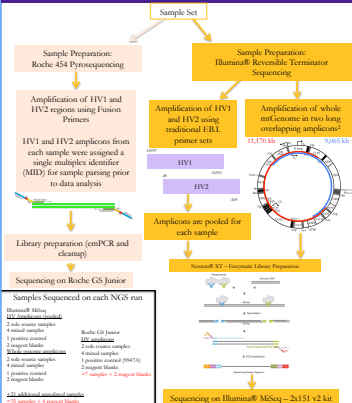


Figure 1: Chemistry and analysis of sequencing NGS instruments.
On the left, Roche 454 Pyrosequencing chemistry is illustrated. Initially, DNA target regions are amplified using PCR with high-fidelity polymerase, and amplified products containing 454 adaptor sequence and multiplexing barcodes. For a template, specific sequence. Multiplexing barcodes enable per-bead sample pooling, while adaptors allow for clonal amplification on the surface of a bead, which then localizes a Zero-Cycle™ plate (ZCP) when sequencing chemistry is initiated. Native dNTPs are flowed across the surface of the ZCP. Release of PPY following incorporation initiates an enzymatic cascade that ends in the ATP catalyzed conversion of luciferin to oxyluciferin, and the production of light. The amount of light produced is directly related to the number of nucleotides in the region being sequenced.
On the right, Illumina® chemistry is shown. DNA target regions are amplified using traditional PCR and a high-fidelity polymerase. The resulting amplicons are fragmented enzymatically using the Illumina® Nextera XT fragmentation kit, which also enables incorporation of sequencing adaptor and multiplexing barcodes during a limited cycle PCR step. Libraries are then added to an optically transparent flowcell, where they bind to methanol oligonucleotide complementary to their incorporated adaptor sequences. Clusters consisting of clonal populations of individual template molecules are generated on the instrument. Sequencing takes place after cluster generation when modified dNTPs with base specific fluorophores and blocking groups on their 3'-OH groups are flowed across the flowcell. After each cycle, LED light excites nucleotide specific fluorophores, and the software images the surface of the flowcell to catalogue the specific base incorporated at each cycle.

NGS LIBRARY PREPARATION



DATA ANALYSIS – CLC GENOMICS WORKBENCH

CLC Genomics Workbench v6.5 was used for all data analysis. Raw data files (SFF files) from the Roche GS Junior were uploaded into CLC Genomics Workbench, and demultiplexed using the software. Fastq files demultiplexed during secondary analysis on the Illumina® MiSeq were also uploaded. All sample files were analyzed using the same pipeline. The data was initially mapped to the revised Cambridge reference sequence (rCRS)² using a local alignment option. Variant calling was performed using the quality-based variant detection method with a 0.1% minor variant detection threshold to enable capture of a majority of sites showing variability. Resulting variant tables were exported as tab-delimited text files, and uploaded into the Galaxy³ open-source cloud computing environment. A custom application within Galaxy was applied to the data set to assist with categorization of unexpected variants. Full analysis parameters, and raw data files are available upon request.

EXPECTED VARIANT FREQUENCIES – SOLE SOURCE SAMPLES

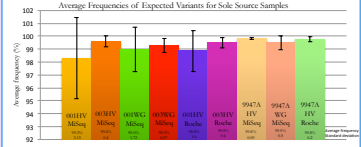


Figure 2: Bar graph showing average frequency for expected variants in HV1 and HV2, divided in sole source samples. The expected variants are those that were previously detected for each donor using Sanger sequencing, and are expected to be detected at or near a frequency of 100% unless otherwise noted. Averages of expected frequency and accompanying standard deviations were calculated for the expected variants of each donor. Error bars represent standard deviations of frequency averages. Donor 001-CF30 is expected to possess a total of 3 variants from the rCRS across the HV regions, with varying levels of heteroplasmy at position 16019 (data not shown). This heteroplasmy is contributing to higher standard deviation values for all donor 001 samples. Donor 003-CM54 is expected to possess a total of 14 variants from the rCRS in the HV region. Differences in frequency were often observed at positions 136C, 132E, 1613G, and 1619T, especially for donor 003-CM54 (figure 6). When this was the case, frequency for these variants were not used and relative frequency averages and standard deviations (overall 150/132/001HV MiSeq, 161/132/001HV MiSeq, 161/132/001HV Roche) were used. The 997A control was expected to possess 5 variants from the rCRS, with no apparent heteroplasmy in the HV region.

EXPECTED VARIANT FREQUENCIES IN HV REGION – MIXED SAMPLES

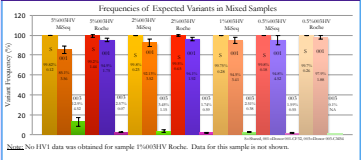


Figure 3: Bar graph showing average frequency for expected variants in HV1 and HV2, divided in mixed samples. The expected variants are those that were previously detected for each donor using Sanger sequencing. All variants shared between donors are expected to be detected at or near a frequency of 100%, barring heteroplasmy sites. Conversely, depending on the mixture, the donor specific variant frequencies change (0.5%: 96.2, 99.1, and 99.6% respectively with donor 003-CM54 representing the minor contributor in all mixtures). In the HV region, donor 001-CF30 and donor 003-CM54 share a total of 3 variants. In addition, donor 001-CF30 possesses 4 unrelated variants from the rCRS, and donor 003-CM54 possesses 9. Similar to the data presented in figure 2, heteroplasmy for donor 001-CF30 at position 16019, is contributing to higher than expected standard deviations in frequency averages. Also, frequencies for variants with observed alignment anomalies were omitted from calculations in this data set. Consistent positions 150/132 from all donor 001 samples and positions 161/132/001HV and 161/001HV data generated with the MiSeq. In all MiSeq HV sample data, frequency of minor contributor variants were higher than expected. This was not observed in data generated using the Roche GS Junior. This was also observed to a much lower extent in whole genome data also generated using the Illumina® MiSeq (see figure 6). This is likely a result of differences in library preparation, since the same samples were used for all sequencing runs. Additionally, in the 0.5% mixture, all minor contributor variants were not seen before the 0.5% minor variant frequency threshold was met. This is a direct result of the higher coverage achieved with the Illumina® MiSeq, and shows that the MiSeq is more sensitive than the Roche GS Junior.

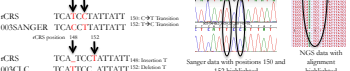
UNEXPECTED VARIANTS IN HV DATA – SUMMARY

	001CF30	003CM54	005CM54	010CM54	020CM54	050CM54	997A
001CF30	112	36	102	82	34	53	32
003CM54	117	36	118	73	99	67	107
005CM54	5	15	3	2	7	0	9
010CM54	60	36	147	135	123	114	67
020CM54	115	34	147	135	28	14	37
050CM54	5	3	1	1	1	1	1
997A	1	1	1	1	1	1	1
Total	204	105	239	238	223	179	227
Control	100	98	100	97	100	98	100

Table 1: Summary of unexpected variants from HV1 and HV2 samples in Roche GS Junior and Illumina® MiSeq data. All unexpected variants were examined using Galaxy open-source cloud computing software. The resulting counts were manually categorized and tabulated. In all cases, when compared to the average coverage calculated for each sample, the number of unexpected variants in Roche GS Junior data sets is higher compared to the total number of unexpected variants in Illumina® MiSeq data sets. This indicates that the Roche GS Junior yields more noise than the Illumina® MiSeq. However, additional research is needed to further characterize the unexpected variants as noise. The 997A control, derived from a human cell line, was not sequenced in the present work, but is included for comparison to the sample obtained from buccal swabs.

ALIGNMENT DIFFERENCES

Figure 4: Alignment differences in Sanger and NGS data at positions 150 and 132 for donor 001-CF30. Differences in alignment at positions 150 and 132 were detected across all samples containing DNA from donor 001-CF30. Alignment discrepancies between bioinformatics algorithms must be considered when implementing NGS into databasing and casework laboratories.



WHOLE GENOME SAMPLES – COVERAGE MAPS

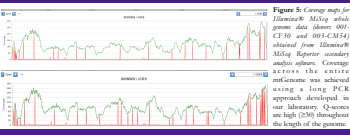


Figure 5: Coverage maps for Illumina® MiSeq whole genome data (donor 001-CF30 and 003-CM54) compared to Roche GS Junior secondary analysis values. Coverage analysis of the HV regions from whole genome data using Illumina® MiSeq was achieved using a long PCR approach developed in our laboratory. Coverage is high (50x) throughout the length of the genome.

UNEXPECTED VARIANTS IN WHOLE GENOME DATA – SUMMARY

	001CF30	003CM54	005CM54	010CM54	020CM54	050CM54	997A
001CF30	396	866	368	360	320	175	912
003CM54	47	124	47	47	47	47	14
005CM54	46	46	46	46	46	46	36
010CM54	201	201	201	201	201	201	118
020CM54	26	26	26	27	26	26	26
997A	1	1	1	1	1	1	1
Total	747	1730	745	747	749	695	1612
Control	100	98	100	97	100	98	100

Table 2: Summary of unexpected variant deconvolution for whole genome samples from Illumina® MiSeq data. Average coverage for whole genome samples is high despite being sequenced in the same run with 7 samples consisting of problem HV amplicons. The number of unexpected variants in Roche GS Junior data sets is higher compared to the total number of unexpected variants in Illumina® MiSeq data sets. This indicates that the Roche GS Junior yields more noise than the Illumina® MiSeq. However, additional research is needed to further characterize the unexpected variants as noise. The 997A control, derived from a human cell line, was not sequenced in the present work, but is included for comparison to the sample obtained from buccal swabs.

EXPECTED VARIANT FREQUENCIES – WHOLE GENOME DATA

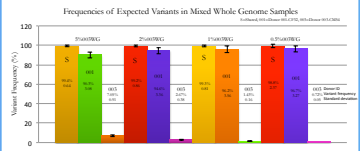


Figure 6: Bar graph showing average frequency for expected variants for whole genome data, divided in mixed samples. All variants shared between donors are expected to be detected at a frequency of 100%. However, depending on the mixture, the donor specific variant frequencies will change (0.5%: 96.2, 99.1, and 99.6% respectively with donor 003-CM54 representing the minor contributor in all mixtures). Frequency averages and standard deviations were calculated for the expected variants of each donor. Error bars represent standard deviations of frequency averages. Across the genome, donor 001-CF30 and donor 003-CM54 share a total of 21 variants. In addition, donor 001-CF30 possesses 7 unrelated variants from the rCRS, and donor 003-CM54 possesses 18. Frequencies for variants with observed alignment anomalies were omitted from calculations in this data set as well (control 150/132 in all samples, 161/132/001HV/001HV). Furthermore, several data points were not detected in 0.5%:001HV. A large portion of data is missing from this set. More research is being conducted to determine the cause. Overall, the frequencies of the expected variants seem to more consistently align when using our whole genome approach, than when using an amplicon approach (figure 3).

CONCLUSIONS

- The Illumina® MiSeq allows for rapid enzymatic library preparation, higher throughput, and longer read-lengths than the Roche GS Junior.
- Much higher coverage was achieved per sample when using the Illumina® MiSeq versus the Roche GS Junior, even with additional whole genome samples included in the MiSeq run. As a result, we feel that the MiSeq can result in higher sensitivity depending on the design of the run, and lead to detection of minor variants at a lower frequency than the GS Junior.
- Additionally, it appears that the Roche GS Junior is more prone to noise than the Illumina® MiSeq, especially in areas containing homopolymers. However, this observation is not new. Additional work is being done to re-analyze the data from the Roche GS Junior with filtering mechanisms designed to remove sequencing artifacts associated with homopolymers. Also, unexpected variants which were labeled as “noise” for this work, are being further studied to identify positions that may be consistent with true biological variations.
- The types of unexpected variants observed in data obtained using both instruments is consistent across samples sequenced.
- Overall, we feel that the Illumina® MiSeq is well-suited for forensic mtDNA analysis through additional validation studies are required to move this technique into the crime laboratory.

REFERENCES

1. M.F. Kavliak, H.S. Lawrence, E.T. Merritt, C. Fisher, A. Isenhardt, J.M. Robertson, B. Badovics. Quantification of human mitochondrial DNA using microfluidic DNA standard. *J Forensic Sci*. 2011;56(6):187-95.
2. B.M. Anderson, J. Kohanski, P.F. Cherny, B.N. Ligonovskaya, D.M. Turchish, N. Howell. Real-time PCR of the Cambridge reference sequence for human mitochondrial DNA. *Nat Commun*. 2012;3:1047.
3. H. Sawada, B.J. Bentz, E.S. Burnside, M.R. Wilson. Paired-whole genome human mitochondrial libraries using Illumina® Nextera XT. *Proceedings of AAFS*, 2013.
4. Public Catalog Intron. <http://ncicb.nci.nih.gov/>

ACKNOWLEDGMENTS

The authors of this work wish to gratefully acknowledge Illumina® for their continuing support. Special thanks to Cynide Holt, Joe Variano, Joe Walsh, Kathryn Stevens, Cary Davis and Tom Rutherford. We also wish to thank Western Carolina University for supporting some of this work.

CONTACT

Brittania J. Bintz, M.S.
Forensic Science Program
Western Carolina University
111 Memorial Drive
NSB #213
Cullowhee, NC 28723
bbintz@wcu.edu
828.227.3600