

Clinical whole genome sequencing: a validation study

Gerard Irzyk¹, C. Alexander Valencia¹, Zeqiang Ma¹, Edward Szekeres Jr¹, Zdenek Markovic¹, Dhara Shah¹, Yang Wang¹, Alice Tanner¹, Christin Collins¹, Madhuri Hegde^{1,2}
¹PerkinElmer Genomics, ²Emory University

ABSTRACT

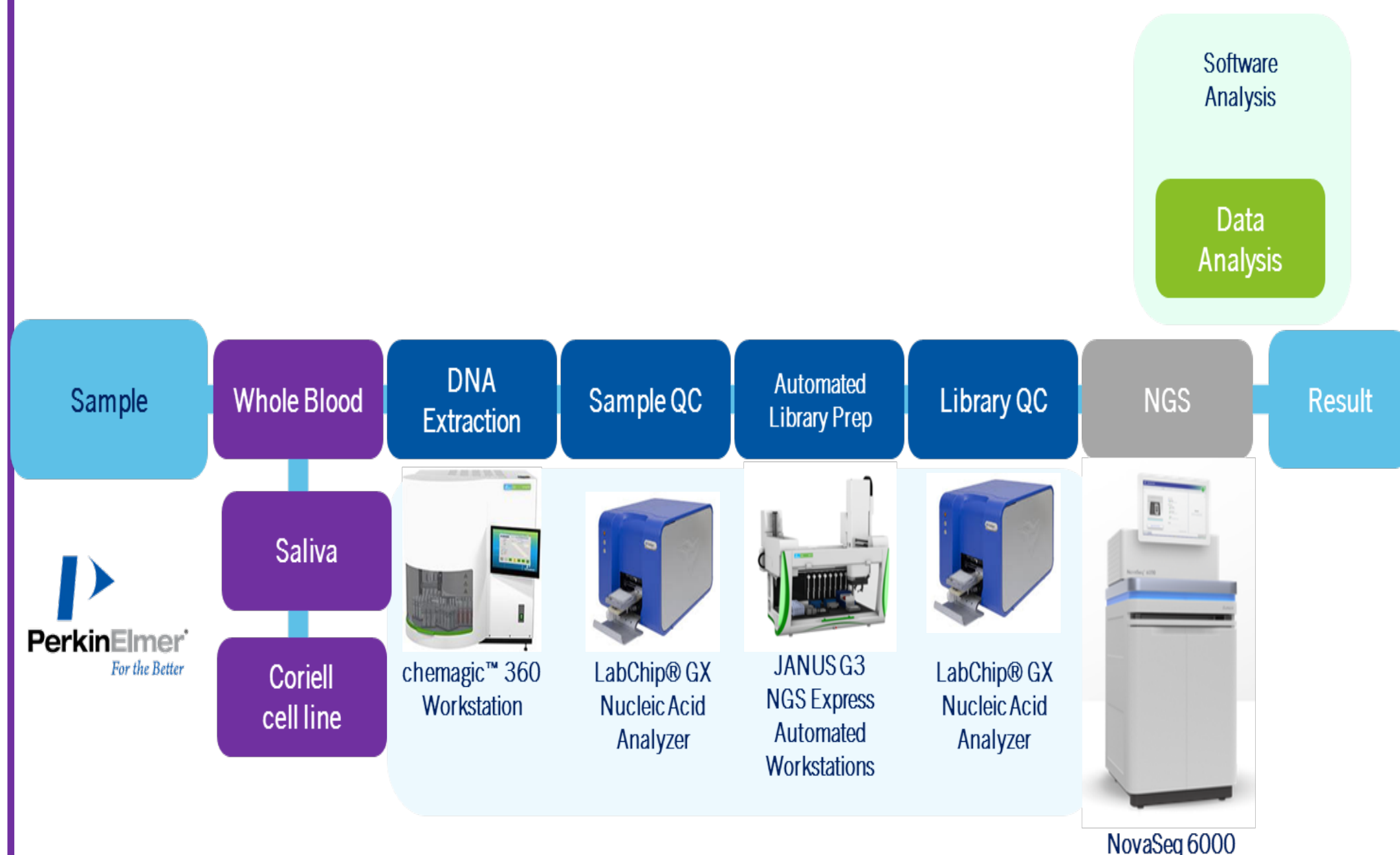
Whole-exome sequencing (WES) is now routinely used in clinical genetic testing and is optimized for the detection of rare and common genetic variants in humans. However, whole genome sequencing (WGS) is becoming increasingly attractive as an alternative, due to its uniform coverage, ability to detect all types of variants and decreasing cost. The objective of this study was to validate WGS in a clinical laboratory using a PCR-free library preparation protocol followed by sequencing on the Illumina NovaSeq™, primary data processing using the Edico Dragen system, and bioinformatic analysis using our in-house proprietary program ODIN. Data was parsed into various categories: Genes causing disease (GCD), Genes of unknown significance (GOUS), ACMG 59, known common and founding pathogenic changes, intragenic and intergenic variants with tagged variants which have been established to be diseases causing. A total of 52 samples, including positive control Coriell samples with known pathogenic variants (26), three control samples from Coriell cell lines NA12878, NA12892 and NA12891, DNA extracted from whole blood samples (23) obtained from positive controls with known pathogenic variants or DNA extracted from unaffected saliva (3). All samples were sequenced at >30X average coverage and analyzed. Subsequently, performance parameters for accuracy, precision, sensitivity, and specificity were calculated. The average global genome metrics for the 52 samples were 38.56X average coverage and 90.60% coverage of the genome at 10X. In addition, the average number of SNPs and indels detected were 4,132,282 (81.52% of calls) and 936,598 (18.48% of calls), respectively. The average number of homozygous and heterozygous calls were 3,243,145 and 1,825,734, respectively. To calculate inter-run accuracy, NA12878 was run twice and compared to the genome in a bottle (GIAB) NA12878 sample, yielding an accuracy value of 99.7% (3,207,518 SNPs in common). Of the unique SNPs (19,951) found in our NA12878 samples, ~30% were determined to be artifacts by bioinformatics analysis. The precision value of 99.2% was calculated by running 11 samples, including Coriell samples NA12878, NA12892 and NA891, two times. Specifically, there was an average number of 3,379,210 SNPs that were concordant in the duplicate samples. The sensitivity, defined as the percentage of low coverage exons, was 99.00% based on 52 samples. From a total of 63,962 exons (4802 genes associated with disease), an average of only 704 exons had coverage less than 8X belonging to 160 genes, which included *SMN1* and *SMN2* due to high homology. The average specificity was 99.30%. Moreover, in this validation study samples with specific variant(s) were included such as missense, splicing and nonsense variants as well as exon-level deletions and duplications to see if they could be identified. All specific variants were identified by WGS. For example, the expected c.1448T>C (p.L483P) and c.115+1G>A were detected in *GBA*, which has homologous sequence with 99% overlap. Moreover, we detected a hemizygous deletion of exons 8-10 in *ABCD1*, X-linked adrenoleukodystrophy. In this WGS validation study it was demonstrated that all of the performance parameters were ≥99%, thus, it is acceptable as a clinical testing. As we learn more about the non-coding regions of the genome we anticipate better clinical interpretation and thus more powerful re-analysis using WGS data in the near future.

INTRODUCTION

- Whole-exome sequencing (WES) is now routinely used in clinical genetic testing and is optimized for the detection of rare and common genetic variants in humans with an approximately 25-30% diagnostic yield
- CHALLENGES:** WES can miss major types and regions of disease-causing genomic variation (INDELs, structural variants, intronic SNVs)
- Whole genome sequencing (WGS) is becoming increasingly attractive as a clinical alternative
- WGS ADVANTAGES:** (i) Higher diagnostic yield than WES, (ii) detection of non-exonic sequence variants, (iii) improved CNV detection and (iv) added lifetime value of re-analysis over time, as clinical features develop in the future of an individual
- APPROACH:** WGS was performed by using the KAPA HyperPlus PCR-free library construction kit on DNAs obtained from a number of samples types: whole blood, saliva, and Coriell cell lines. The libraries were subsequently sequenced Illumina NovaSeq™ 6000 operating in 2 x 150 bp mode

METHODS

WGS sequencing workflow



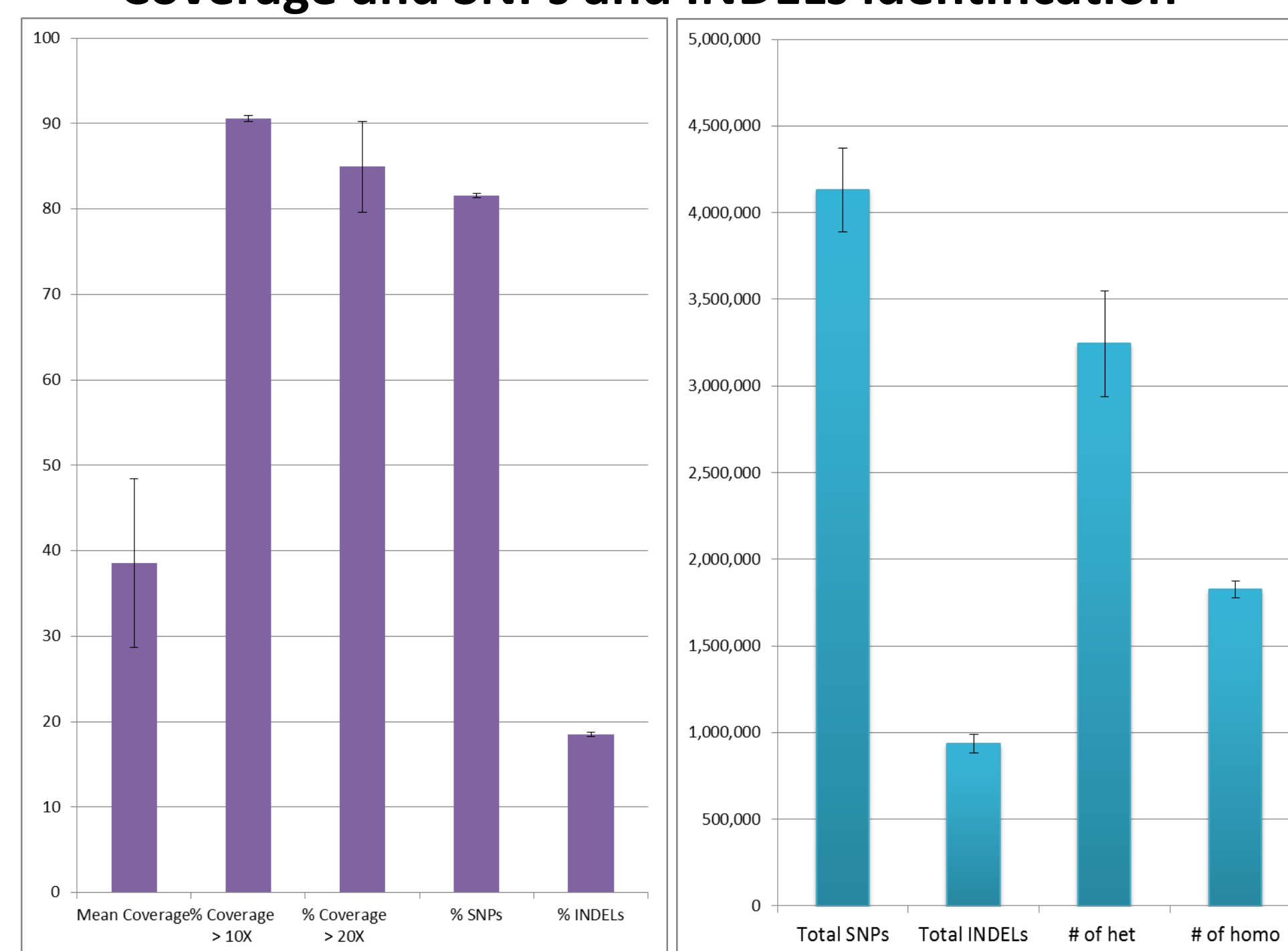
WGS validation plan

| Sample Name | Sample type | Accuracy | Precision | Sensitivity | Specificity |
|--------------|-------------|----------|-----------|-------------|-------------|
| NA12878-GIAB | CL | + | + | + | + |
| NA12878-1 | CL | + | + | + | + |
| NA12878-2 | CL | + | + | + | + |
| NA12892-1 | CL | + | + | + | + |
| NA12892-2 | CL | + | + | + | + |
| NA12893-1 | CL | + | + | + | + |
| NA12893-2 | CL | + | + | + | + |
| NA12894-1 | CL | + | + | + | + |
| NA12894-2 | CL | + | + | + | + |
| NA12895-1 | CL | + | + | + | + |
| NA12895-2 | CL | + | + | + | + |
| NA12896-1 | CL | + | + | + | + |
| NA12896-2 | CL | + | + | + | + |
| NA12897-1 | CL | + | + | + | + |
| NA12897-2 | CL | + | + | + | + |
| 43313-1 | WB | + | + | + | + |
| 43314 | WB | + | + | + | + |
| 5656-1 | WB | + | + | + | + |
| 5656-2 | WB | + | + | + | + |
| 5647-1 | WB | + | + | + | + |
| 5647-2 | WB | + | + | + | + |
| NA20310 | CL | + | + | + | + |
| NA20270 | CL | + | + | + | + |
| NA00694 | CL | + | + | + | + |
| NA03251 | CL | + | + | + | + |
| NA00327 | CL | + | + | + | + |
| NA22496 | CL | + | + | + | + |
| NA13537 | CL | + | + | + | + |
| NA13653 | CL | + | + | + | + |
| NA06314 | CL | + | + | + | + |
| NA03814 | CL | + | + | + | + |
| 7 | WB | + | + | + | + |
| 5 | WB | + | + | + | + |
| 5642 | WB | + | + | + | + |
| 14 | WB | + | + | + | + |
| 5548 | WB | + | + | + | + |
| 5635 | WB | + | + | + | + |
| 3639 | WB | + | + | + | + |
| 2A | WB | + | + | + | + |
| 16 | WB | + | + | + | + |
| 3069 | WB | + | + | + | + |
| 5C | WB | + | + | + | + |
| JF21G | WB | + | + | + | + |
| B-SLV | SLV | + | + | + | + |
| B-WVB | WVB | + | + | + | + |
| E-SLV | SLV | + | + | + | + |
| E-WVB | WVB | + | + | + | + |
| T-SLV | SLV | + | + | + | + |
| T-WVB | WVB | + | + | + | + |
| Total | | 3 | 22 | 52 | 52 |

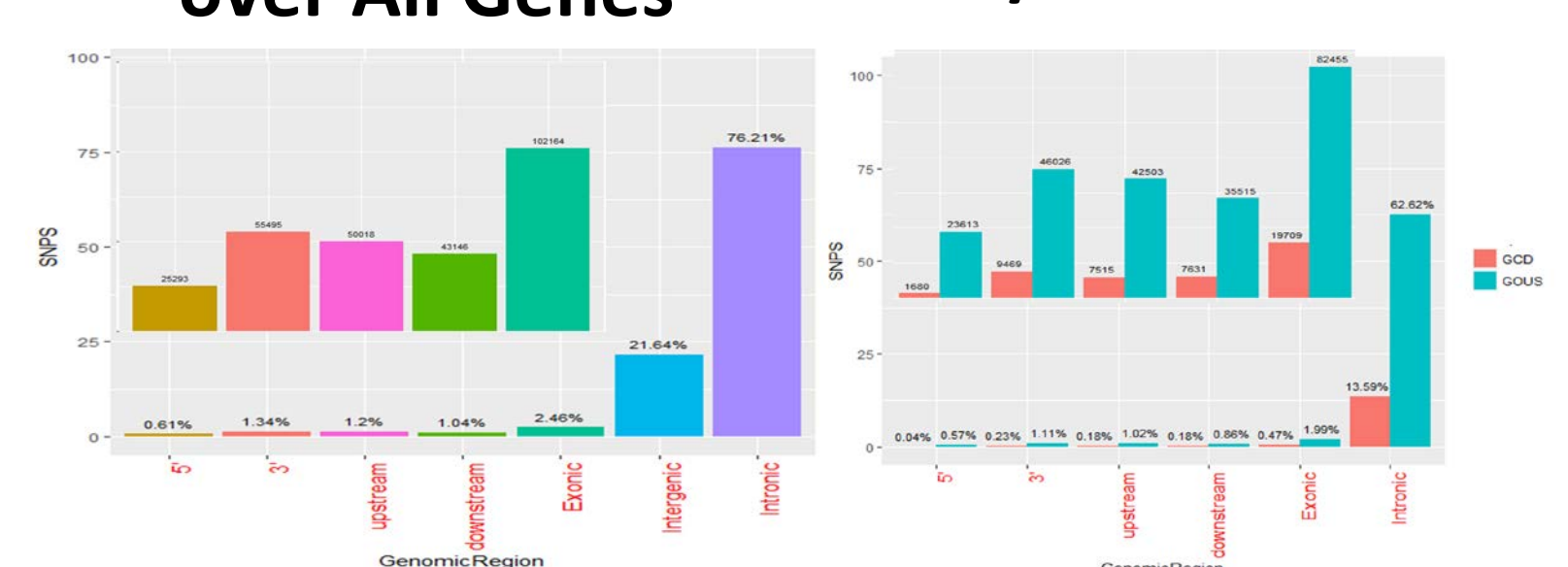
RESULTS

Global Statistics: High quality global key measurements

Coverage and SNPs and INDELs identification



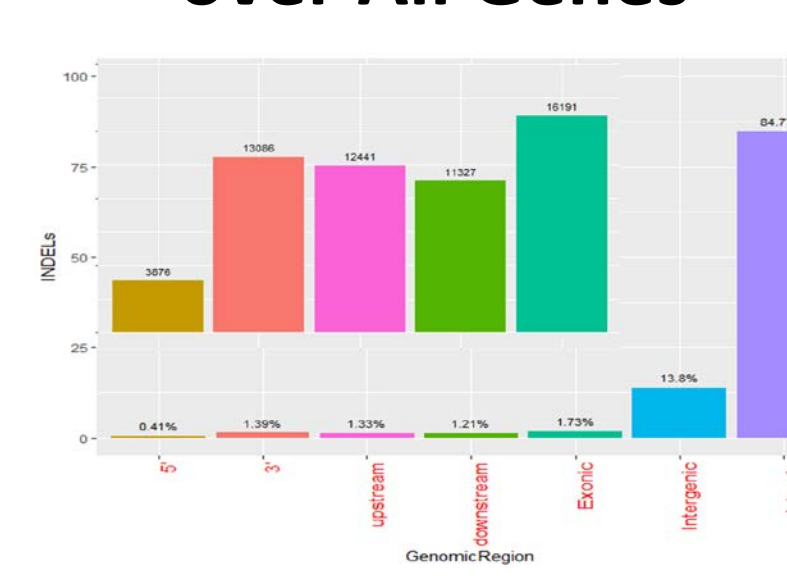
SNPs Distribution over All Genes



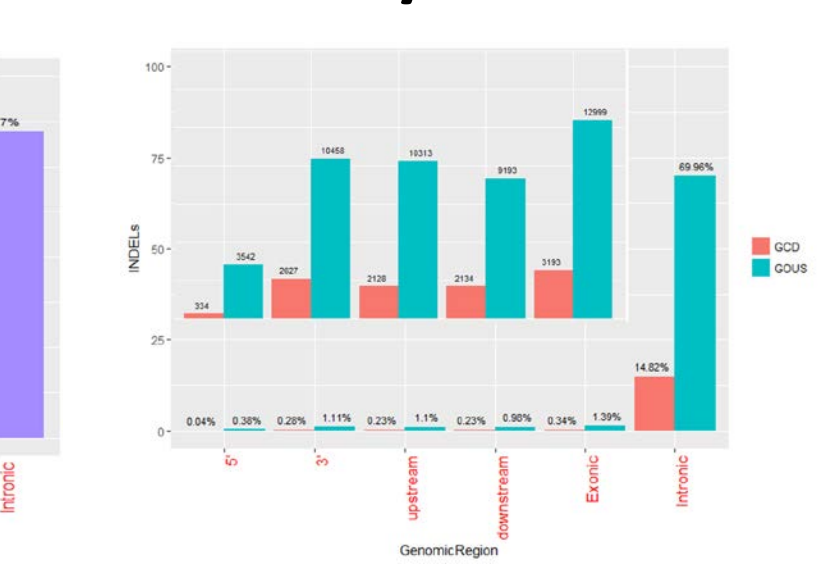
SNPs Distribution over GCD/GOUS Genes



INDELs Distribution over All Genes



INDELs Distribution over GCD/GOUS Genes



- High coverage and sensitivity detection of SNPs and INDELs

- Comprehensive detection of SNPs and INDELs across entire genome and genes causing disease (GCD) genes

Excellent performance and comprehensive detection of variants

Accuracy is 99.7%

| Terms | NA12878-1 | NA12878-2 |
|-------------------------------|--------------|--------------|
| False positive (NA12878 only) | 19222 | 20580 |
| False negative (GIAB only) | 4076 | 1388 |
| True positive (Overlap) | 3280574 | 3280662 |
| True negative (Overlap) | 3800617 | 3801909 |
| Accuracy | 99.7% | 99.7% |
| Average accuracy | 99.7% | |

*GIAB - not confirmed by orthogonal method

Sensitivity is 99.0%

| Sample Name | Variant Type | Count | Percentage |
|-------------|--------------|---------|------------|
| NA12878-1 | SNP | 4132282 | 81.52% |
| NA12878-1 | INDEL | 936598 | 18.48% |
| NA12878-1 | SV | 10 | 0.00% |
| NA12878-1 | CNV | 10 | 0.00% |
| NA12878-1 | STR | 10 | 0.00% |
| NA12878-1 | Other | 10 | 0.00% |

Precision is 99.2%

| Sample Name | Variant Type | Count | Percentage |
|-------------|--------------|---------|------------|
| NA12878-1 | SNP | 3207518 | 99.7% |
| NA12878-1 | INDEL | 936598 | 18.48% |
| NA12878-1 | SV | 10 | 0.00% |
| NA12878-1 | CNV | 10 | 0.00% |
| NA12878-1 | STR | 10 | 0.00% |
| NA12878-1 | Other | 10 | 0.00% |

Specificity is 99.3%

| Sample Name | Variant Type | Count | Percentage |
|-------------|--------------|-------|------------|
| NA12878-1 | SNP | 19951 | 99.2% |
| NA12878-1 | INDEL | 19951 | 99.2% |
| NA12878-1 | SV | 19951 | 99.2% |
| NA12878-1 | CNV | 19951 | 99.2% |
| NA12878-1 | STR | 19951 | 99.2% |
| NA12878-1 | Other | 19951 | 99.2% |

Detection of a spectrum of variants

| Sample Name | Disorder (Reference) | Gene | Allics |
|-------------|--|-------|---|
| NA17819 | Adrenoleukodystrophy (ARL) | ABCD1 | Exon 8-10 deletion (Homo) |
| NA8732 | Micropolyarcharidiosis type VI (Marsheam-Lang) (ARL) | AKR3B | c.14536A>G (p.R4780G) (Homo) |
| NA8161 | Micropolyarcharidiosis IVA (ARL) | GALNS | c.181C>T (p.R61W) (Het) |
| NA20088 | Glycogen storage disease II (ARL) | GAA | c.1213_1218delGTGACC (p.W402_T406del) (Het) |
| NA8472 | Krabbe disease (ARL) | GALC | IVS18 + 2T-A (not detected) 8,264 kb deletion extending from IVS5 to IVS15 (Homo) |
| 8 | Krabbe disease (ARL) | GALC | Exon 11-17 deletion (Homo) common 90 kb deletion (Het) |
| 5656 | Becker muscular dystrophy/Duchenne muscular dystrophy (XLR) | DMD | Exon 45-55 deletion (Homo) |
| 5647 | Becker muscular dystrophy/Duchenne muscular dystrophy (XLR) | DMD | Exon 18 duplication |
| NA3010 | Glycogen storage disease II (ARL) | GAA | 94 kb deletion extending from IVS15 to 44b downstream of exon 20 (Homo) |
| NA12879 | Gusher disease (ARL) | GRA | c.1448T>C (p.L483P) (Het) c.115+1G>A (Het) |
| NA8084 | Acyl-CoA dehydrogenase, medium-chain deficiency of (ARL) | ACADM | c.446G>A (p.C149Y) (Het) c.398A>G (p.R332E) (Het) |
| NA8251 | GMI-gangliosidosis IHL Micropolyarcharidiosis type VII (Morgani) (ARL) | GLB1 | c.1527G>T (p.W50C) (Het) c.1800T>C (Homo) |
| NA8017 | Fabry disease (XLR) | GLA | c.485D>A (p.Trp61I) (Homo) c.600 A>G (Homo) c.1080T>C (Homo) |
| NA12296 | Propionicacidemia (ARL) | PCCB | c.1218delA (Homo) TGACACAGGA (c.1218delA Hom 1) (Homo) |
| NA1307 | Spirochetemia (AD) | SCA1 | 401 CG repeats 32 CAG repeats |
| NA1461 | Glycogen storage disease II (ARL) | GAA | IVS8-11E-G (Het) Exon8 deletion |
| NA8434 | Glycogen storage disease II (ARL) | GAA | IVS13T-G (Het) c.1551+10C>G (Het) |
| NA8814 | Spinal muscular atrophy-I (ARL) | SMN2 | Exon 7-9 deletion (Het) |

DISCUSSION/CONCLUSIONS

- NGS in the past decade has led to an enormous increase in the understanding of the human genome and its relation to disease
- Improved technologies continuously provide faster, cheaper and more accurate results, allowing us to move from gene panels to WGS to routinely sequencing WGS in the clinic
- In this WGS validation study it was demonstrated that all of the performance parameters were **≥99%, thus, it is acceptable as a clinical test**
- We have successfully used the documented clinical cases to demonstrate that WGS performance translates to a useful clinical sensitivity
- We demonstrated the analytical sensitivity for SNVs and small indels and demonstrated that intra- and inter-gene CNVs were accurately detected by improved calling algorithms, thus, increasing the diagnostic yield of genetic disease burden
- As we learn more about the non-coding regions of the genome, large scale efforts, for example gnomAD, we anticipate improved clinical interpretation and thus more powerful re-analysis using WGS data.