

Targeting Clinically Significant Dark Regions of the Human Genome with High-Accuracy, Long-Read Sequencing

Abstract #: eP273

Ian McLaughlin, John Harting, Lori Aro, Cheryl Heiner
Pacific Biosciences, 1305 O'Brien Drive, Menlo Park, CA 94025

Introduction

There are many clinically important genes in "dark" regions of the human genome characterized by a paucity of NGS coverage as a result of short-read sequencing or mapping difficulties. Low NGS sequencing yield can arise in these regions due to the presence of various repeat elements or biased base composition while inaccurate mapping is attributable to segmental duplications. Long-read sequencing coupled with an optimized, robust enrichment method has the potential to illuminate these dark regions.

Methods and Materials

High-long PCR targeted enrichment was combined with PacBio long-read, high-accuracy (HiFi) sequencing and a novel amplicon analysis tool (pbaa) to develop individual prototype screening assays for *CYP21A2* and *GBA*; genes associated with congenital adrenal hyperplasia and Gaucher disease, respectively. Each gene is difficult to accurately type with short-read sequencing due to proximity to and interaction with a highly homologous pseudogene.

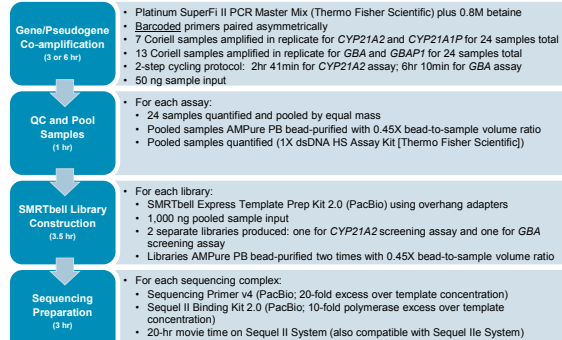


Figure 1. Workflow for *CYP21A2* and *GBA* screening assays. Two days are required starting from purified gDNA that has sequencing complexes ready to sequence. An additional two days are required to sequence the libraries and produce variant-type results.

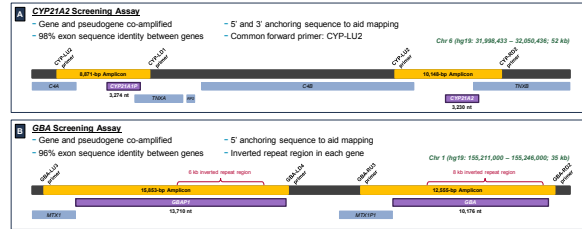


Figure 2. PCR target maps for *CYP21A2* (A) and *GBA* (B) screening assays. Primer sequences were developed in-house. Primer reaction concentrations were adjusted to achieve an approximately balanced yield of co-amplified amplicons in each assay system. Both systems benefited from supplemental betaine at a final concentration of 0.8M, especially the *GBA* system with a large tract of inverted repeats in each gene likely leading to stable secondary structures in the absence of supplemental betaine.

Results

Assay	PCR Input	PCR Yield (mean)	Library Input	Library Yield	HiFi Reads (S020)	%Barcoded HiFi Reads	%Full-length Aligned Reads
<i>CYP21A2</i>	50 ng	431 ± 131 ng	1 µg	671 ng	1,520,782	84% (1,279,425)	84% (1,276,879)
<i>GBA</i>	50 ng	1,086 ± 304 ng	1 µg	673 ng	1,045,191	77% (801,843)	81% (847, 538)

Table 2. Top-line performance metrics for *CYP21A2* and *GBA* screening assays. A high percentage of HiFi reads are barcoded and produce full-length alignments to the target genes.

Gene	Aligned Reads (mean)	CV	Min	Max
<i>CYP21A1P</i>	17,223	71%	526	41,817
<i>CYP21A2</i>	27,471	41%	985	45,696
<i>GBA1P</i>	8,065	51%	3,866	21,355
<i>GBA</i>	22,165	23%	10,341	31,652

Table 3. Mean number of aligned reads per sample for each target gene. The difference in yield between target genes in each assay reflects known PCR yield differences. The high CV level reflects the presence of samples with alleles containing deletions ultimately reducing the number of reads for those samples.

Sample	Expected Gene Amplification: <i>CYP21A1P</i>				Expected Gene Amplification: <i>GBA1P</i>				Outer Primers LU3-RD2			
	<i>CYP21A1P</i>	<i>CYP21A2</i>	<i>CYP21A1P</i>	<i>CYP21A2</i>	<i>GBA1P</i>	<i>GBA</i>	<i>GBA</i>	<i>GBA</i>	<i>GBA1P</i>	<i>GBA</i>	<i>GBA1P</i>	<i>GBA</i>
02241	15,053	326	161	38,273	0,982	9,360	1	6	26,984	147	175	175
02242	17,856	378	116	37,845	0,877	6,145	9	16	23,410	474	453	453
11781	33,736	168	313	30,757	0,078	4,939	18	54	22,317	908	991	991
12217	27,662	73	16,414	17,085	0,131	5,24	20	33	25,102	706	699	699
14732	895	411	20,927	20,728	0,1260	6,938	18	18	20,358	443	421	421
14733	14,989	353	22,491	22,168	0,1607	5,805	10	30	24,768	815	921	921
14734	870	1	46,529	27	0,2627	6,569	5	13	27,590	475	318	318
08752	5,639	8	50	19,438	0,2972	5,639	8	50	19,438	20,250	20,250	20,250
08753	5,279	32	142	12,894	1,117	11,784	11,717	11,784	11,717	11,784	11,784	11,784
10873	7,119	13	31	21,249	865	961	961	961	865	961	961	961
10874	8,247	2	9	26,519	297	356	356	356	297	356	356	356
20270	10,090	4	5	24,745	0	0	0	0	24,745	0	0	0
20273	4,667	2	25	26,531	724	809	809	809	724	809	809	809

Table 4. Mean number of aligned reads per sample replicate for *CYP21A2* (A) and *GBA* (B) screening assays. Highlighted values in purple indicate numbers significantly different than others within their respective columns. *CYP21A2* samples 12217, 14732, 14733, and 14734 contain *CYP21A2* deletion alleles which alter the number of expected reads (e.g., sample 14734 contains a full deletion and partial deletion of *CYP21A2*, forming fusion alleles between *CYP21A1P* and partial *CYP21A2* or sequence downstream of *CYP21A2*). *GBA* samples 08752 and 08753 contain *GBA* deletion alleles but due to the priming strategy the resulting fusion product requires detection with the outer-most primers, LU3 and RD2. For each assay, the presence of a low number of mismatched aligned reads to expected amplicons is due to the high level of homology between gene and pseudogene target regions.

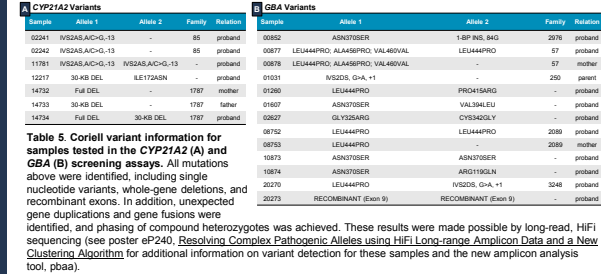


Figure 3. Detection of deletion alleles, gene fusions, and an SNV in *GBA* compound heterozygote samples 08752 and 08753. The deletions and gene fusions were not described in the Coriell database, only the L444P pathogenic variant. Primers are identified at the ends of each aligned read, and the L444P mutation is indicated with a red arrow. The fusion allele is created by a large deletion encompassing *GBA* and *GBA1P* genes.



Figure 3. Detection of deletion alleles, gene fusions, and an SNV in *GBA* compound heterozygote samples 08752 and 08753. The deletions and gene fusions were not described in the Coriell database, only the L444P pathogenic variant. Primers are identified at the ends of each aligned read, and the L444P mutation is indicated with a red arrow. The fusion allele is created by a large deletion encompassing *GBA* and *GBA1P* genes.

Conclusions

- We demonstrate that HiFi sequencing provides new opportunities for sequencing clinically relevant but previously dark regions of the human genome that are underrepresented in short-read sequencing.
- The *CYP21A2* and *GBA* screening assays achieved accurate typing results for all samples and provided additional variant information due to the use of optimized and robust long-range PCR enrichment, HiFi sequencing, and newly developed PacBio amplicon analysis tool, pbaa.
- Accurate long reads provide important phasing information, identify structural variations, and avoid potential confusion with pseudogenes.
- The PCR, library preparation, and sequencing preparation steps of these assays are amenable to automation and a cost-effective workflow enabling high-throughput screening.
- HiFi sequencing of these regions enables a better understanding of the relationship between genetic factors and personal health and has the potential to ultimately help guide health-related decisions.