Paired End Sequencing on the SOLiD[™] Platform

Eileen T. Dimalanta¹, Lei Zhang¹, Lele Sun¹, Kara S. Eusko¹, Dalia M. Dhingra¹, Manning¹, Heather E. Peckham¹, Eric Tsung¹, Steve M. Menchen², Alan P. Blanchard¹, and Kevin J. McKernan¹ [']Life Technologies, 500 Cummings Center, Suite 2400, Beverly, MA 01915 [']Life Technologies, 850 Lincoln Centre Drive, Foster City, CA 94404

ABTSRACT

The SOLiD[™] platform is a revolutionary sequencing system that utilizes sequential ligation of fluorescently labeled oligonucleotide probes, enabling high fidelity and ultra-high throughput sequencing. Previous sequencing protocols for the SOLiD[™] system have only been available in the forward direction (3' to 5'). Sequencing in the reverse direction (5' to 3') is ideal as it enables fragment library paired end sequencing. To this end, novel ligation chemistries were developed to support 5' to 3' read lengths of up to 35 bases. This paired end sequencing technology will be incorporated into the SOLiD V4 platform, increasing effective read length, maximizing throughput per run, and meeting special research interests, such as whole transcriptome studies.

SOLiD™ Overview

SOLiD[™] Sequencing involves the serial ligation of probes in which a dye reports the subset of four possible dibase pairs at the 1st and 2nd positions from the ligation junction. SOLiD™ Sequencing in the forward direction involves: A.) 5'-phosphorylated primer is hybridized to the adapter region of the templates to be sequenced; B.) Fluorophore labeled 8-mer complementary probes, containing 3 universal bases to decrease complexity, are ligated to the primer; a second round of ligation is performed with unlabeled probe to increase the amount of primer extended per bead; C.) Any remaining unextended primer is capped by dephosphorylation to prevent dephasing; beads are then imaged to record fluorophore reporter; D.) A phosphorothiolate bond in the ligated probes is cleaved with AgNO3, reducing the probe to 5 nucleotides and generating a free phosphate for the next round of ligation; E.) Additional cycles of ligation, capping, imaging, and cleavage are performed until the desired read length is obtained.

SOLiD[™] Sequencing in the reverse direction involves: F.) 3'-hydroxylated primer is hybridized to the adapter region of the templates to be sequenced; G.) Fluorophore labeled 8-mer complementary probes (5' phosphorylated), containing 3 universal bases to decrease complexity, are ligated to the primer; a second round of ligation is performed with unlabeled probe to increase the amount of primer extended per bead; H.) Any remaining unextended primer is capped by polymerase incorporation of a ddNTP to prevent dephasing; beads are then imaged to record fluorophore reporter; I.) A phosphorothiolate bond in the ligated probes is cleaved with AgNO3, reducing the probe to 5 nucleotides and generating a free 3' phosphate; J.) The 3' phosphate is removed for the next round of ligation; K.) Additional cycles of ligation, capping, imaging, cleavage, and dephosphorylation are performed until the desired read length is obtained.







F3 Adapter

DNA Fragment ~ 200 Bases

F5 Adapter

Figure 1. Paired end sequencing on the SOLiD[™] platform. Priming from the 5' end of the template for forward reads (F3 Adapter) and from the 3' end for the reverse reads (F5 adapter) of a standard fragment library. Novel reverse ligation chemistry enables paired end sequencing without requiring synthesis of the complementary strand of the DNA template.

> Figure 2. "Satay" Plots for an *E. coli* DH10b 35 base reverse sequencing run. These plots show the spectral quality and intensity of the sample. The axes correspond to the 4 different fluorochromes used in SOLiD[™] sequencing: FAM, CY3, TXR, CY5. Each dot on the plot represents the fluorescent wavelength and intensity of multiple copies of the bead bound DNA template. Beads that fall on or near an axis are monoclonal (i.e. they contain multiple copies of a single DNA template), and beads that are far from the origin are high intensity beads.



Figure 3. Matching Statistics for paired end sequencing of an *E. coli* DH10b genome. For reverse sequencing (F5), 65% of the reads matched the genome for 35 bases allowing up to 4 mismatches; 77% of the reads matched the genome for 25 bases allowing up to 3 mismatches. For forward sequencing (F3) 76% of the reads matched the genome for 50 bases allowing up to 6 mismatches.



- TXR 100,000 beads sampled from 10 panels



Figure 4. Estimated throughput of paired end sequencing based on SOLiD[™] 4 bead densities. A full sequencing run can generate 28 Gigabases of data for 25mers (reverse reads), 35 Gigabases of data for 35mers (reverse reads), and 57 Gigabases of data for 50mers (forward reads). A 25mer (reverse) x 50mer (forward) can generate up to 86 Gigabases of data, and up to 120 Gigabases for 35mer (reverse) x 50mer (forward).

E coli DH10b Pairing Rates

Read Length	Pairing Rate
25 x 50	98.94%
35 x 50	99.05%

 Table 1. Pairing rates for 25mer (reverse)
read) x 50mer (forward read) and 35mer (reverse read) x 50mer (forward read).

Figure 5. Size distribution of DNA template determined from paired reads.

Paired End Sequencing of the Human Genome

Reverse Sequencing Human "Satay" Plots

- DNA extracted from buccal swab
- Constructed paired end library using 1 µg of DNA 1 day library protocol
- 2 SOLiD Sequencing Runs 50 bp Forward Reads 25 bp Reverse Reads
- Generated over 88 Gigabases of aligned reads

Figure 6. "Satay" Plots for a 25 base reverse sequencing run of a human genome.

Figure 7. Distribution of coverage of uniquely placed paired reads across the human genome. The average coverage was 14.5x.

Figure 8. Average coverage for each chromosome. Coverage ranged from 14.7x to 17.4x for the autosomal chromosomes.

Summary of SNPs/Indels

	Count	% in dbSNP	Table 2. Summa
Total SNPS	2972853	82.16%	SNPs and small
Heterozygous SNPs	1963475	73.70%	concordance wit
Homozygous SNPs	1213219	94.42%	
Small Indels *	103027	73.80%	* stringent indel call

Distribution of Small Deletions

Figure 10. Distribution of small deletions (1-11 bp); 71.2% were in dbSNP.

ACKNOWLEDGEMENTS

Christopher Clouser Brittney Coleman Gina Costa Tenzin Dawoe Dave Dupont

were in dbSNP.

- Nick Fantin Tamara Gilbert Sherry Hansen Stephen Hendricks Jeffrey Ichikawa
- Rachel Kasinskas Mike Laptewicz Clarence Lee Liz Levandowsky Stephen McLaughli
- Bashar Mullah Dick Noble Ken Otteson Vaishnavi Panchapakesa Allen Wong Andrew Sheridan
- Jessica Spangler Dean Tsou Tristen Weaver Joon Yang

TRADEMARKS/LICENSING

© 2010 Life Technologies Corporation. All rights reserved.

The trademarks mentioned herein are the property of Life Technologies Corporation or their respective owners.

Cy3 and Cy5 are registered trademarks of GE Healthcare.