

Identifying Novel Expressed Gene Fusions in MCF-7 Cell Line Using Next Generation Sequencing



Onur Sakarya¹, Fiona C Hyland¹, Yulei Wang¹, Heinz Breu¹, Paolo Vatta¹, Liviu Popescu¹, Chieh-Yuan Li¹, Matthew Muller¹, Alec Wong¹, Nriti Garg¹, Sowmi Utiremur¹, Zheng Zhang¹, John P Bodeau¹, Robert C Nutter¹, Mark Mooney¹, Milan Radovich², Penn Whitley³, Kelli Bramlett³, Francisco de la Vega¹, and Asim Siddiqui¹

¹ Life Technologies, Foster City, CA, USA, 94404 and ³ Ambion, Austin, TX, USA 78744
² Indiana University School of Medicine, Indianapolis, IN, USA 46202

ABSTRACT

High throughput RNA sequencing (RNA-Seq) enables detection and quantification of novel transcripts including gene fusions. Gene fusions are potential chemotherapeutic sites such as in the case of BCR/ABL and imatinib. We sequenced the breast cancer cell line MCF-7 using paired-end RNA-Seq protocol with the SOLiD™ 4 System. By using a new gene fusion detection algorithm called SASR (Suffix Array Single Read splice detection) implemented in the BioScope™ software, we called forty gene fusions and validated twenty five of them to be expressed specifically in MCF-7 including four novel inter-chromosomal events. We report most MCF-7 fusion breakpoints on the 5' gene had likely occurred at the early introns (median 23% of gene size) while no bias was observed for the 3' fusion genes' breakpoints. Additionally, we ran TaqMan® assays for select gene fusions on a set of cancer cell lines, and on forty-eight clinical normal and breast tumor samples.

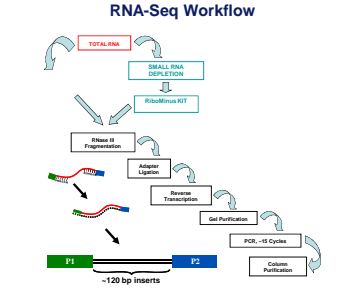
INTRODUCTION

Chromosome aberrations, especially gene fusions, are implicated in the initiation of tumorigenesis. Various gene fusions are important diagnostic and prognostic indicators in leukemia, sarcomas, and other solid tumors. The high throughput of massively parallel sequencers (up to 1 billion mapped reads on a single run of the SOLiD System) enables genome-wide hypothesis-free detection of gene fusions. The availability of DNA barcoded paired-end reads facilitates concurrent sequencing of many samples, reducing per-sample costs.

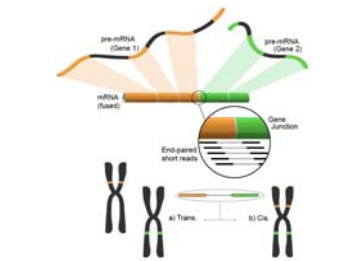
MATERIALS AND METHODS

We sequenced one slide of MCF-7 and two slides of UHR (Universal Human Reference) and HBR (Human Brain) with various insert sizes. We used DNA barcoded samples with paired-end SOLiD System sequencing. We predicted 40 fusions in MCF-7 of which 36 were validated by TaqMan® assays and 25 were specific to MCF-7. Gene fusions were called with the BioScope software.

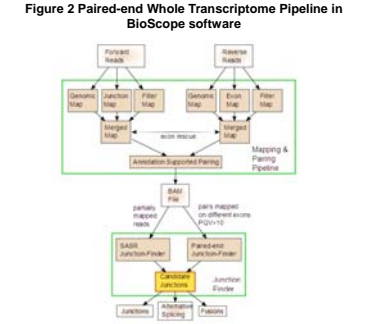
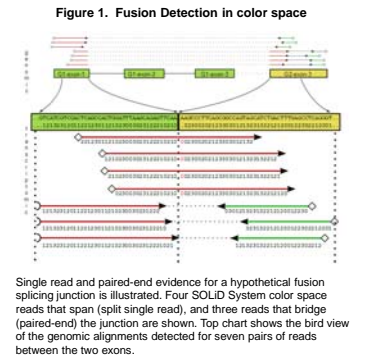
We describe a new suffix array single read (SASR) based intron splicing detection algorithm which allows detection of fusion breakpoints to single base and is not prone to homology based miscalls. In conjunction with the paired-end approach, we devise an integrated, evidence based system to construct splicing graphs and detect fusion transcripts.



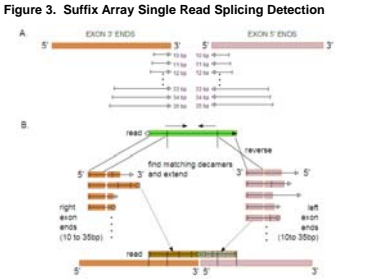
Gene Fusions



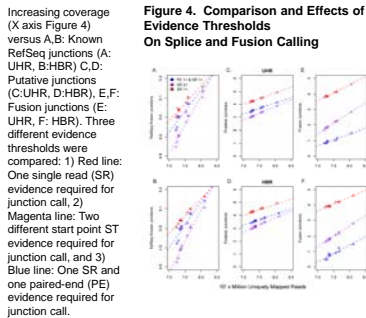
ALGORITHM OVERVIEW



Forward and reverse reads were aligned separately, and paired using an algorithm that finds the best pair based on a probabilistic scheme. This algorithm considered mismatches of individual tags and proximity of the pair in comparison to the expected insert size. After the BAM file is generated, another program is used to detect splicing junctions and gene fusions (Figure 2).



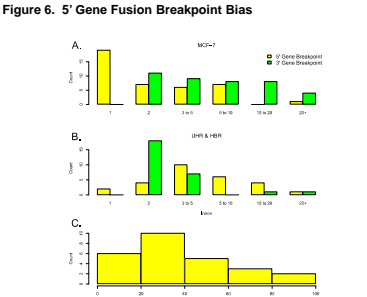
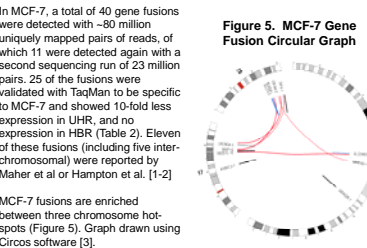
In order to discover single reads that span hypothetical fusion junctions, we constructed a data structure of 12,692,600 left- and right-end suffixes (from hg18 and UCSC gene models) that were sorted lexicographically in arrays allowing logarithmic time string comparison with mismatches (Figure 3)



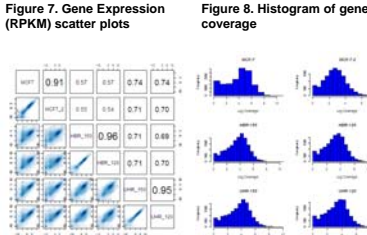
RESULTS

Table 2. List of Validated MCF-7 gene fusions

5' gene-exon	Chr	3' gene-exon	Chr	Distance	MCF7 UHR HBR PC3
ARFGEF2-1	20	SULF2-3	20	Inverted	20.6 24.2 40.0 39.7
SLC25A4-4	1	MBP5-18	1	Inverted	23.9 27.2 40.0 40.0
USP91-1	16	CRYL1-4	13	Inter-chr	27.5 31.5 40.0 40.0
TBL1XR1-1	3	RGS7-2	6	Inter-chr	26.1 30.6 40.0 40.0
TAFA1	20	BRP1-5	17	Inter-chr	26.6 29.2 40.0 40.0
RPS8KB1-6	17	DAPH3-30	13	Inter-chr	22.6 26.1 40.0 36.7
SOX4-1	20	SOX5-24	17	Inter-chr	21.3 25.0 40.0 40.0
ANCYL1-1	1	RAD51C-19	17	Inter-chr	31.0 34.8 40.0 40.0
ANKK4-4	17	PDP4R1L-4	20	Inter-chr	28.1 29.9 40.0 40.0
C10orf45-1	10	ABCC1-15	16	Inter-chr	25.3 29.0 40.0 40.0
C10orf62-8	10	SOX3-10	16	Inter-chr	26.7 30.0 40.0 40.0
C10orf15-1	1	STAF1-2	5	Inter-chr	25.1 32.1 40.0 40.0
MYO6-1	6	SENPP1-15	6	Intra-chr	28.4 31.9 40.0 40.0
RPS6B1-3	17	TMEM31-1	17	Intra-chr	24.4 26.4 40.0 38.9
SMARCA4-7	19	CARM1-2	19	Intra-chr	29.3 33.1 40.0 40.0
POPF2-2	8	MATN2-15	8	Intra-chr	28.5 31.8 40.0 40.0
GATA2B1-1	1	NUP105-28	1	Intra-chr	107.251 29.3 32.4 40.0 40.0
ESR1-2	6	C6orf97-7	6	Intra-chr	116.116 30.3 35.0 40.0 40.0
DEPDC1B-7	5	SLC17A4-6	5	Intra-chr	119.995 26.6 29.0 39.8 40.0
ESR1-2	6	C6orf97-6	6	Intra-chr	128.831 25.2 29.1 40.0 40.0
CDNLL3-2	12	RSN1-14	12	Intra-chr	152.218 25.3 28.2 40.0 28.8
ATXN7L3-1	17	FAM172A-4	17	Intra-chr	159.668 24.8 28.1 40.0 40.0
SYTL2-1	11	PICALM-20	11	Intra-chr	217.187 26.7 30.7 40.0 40.0
ADAMTS19-1	5	SLC27A6-10	5	Intra-chr	434.137 26.5 31.2 40.0 40.0
ADAMTS19-2	5	SLC27A6-10	5	Intra-chr	433.412 26.8 30.5 40.0 40.0



MCF-7 fusion breakpoints were enriched to early 5' introns of the gene (Figure 6 A), but the fusions of UHR and HBR were not (Figure 6 B). On average, first introns of the human genome constitute 22% of the genes and are larger than other introns. We asked whether the large intron size alone would explain the breakpoint bias to 5' introns. The breakpoints fell on average to the first 23% to 36% of the genes (Figure 6 C). In 16 out of 23 fusions, the introns that contained the putative breakpoint were larger than 10Kb, which is considerably larger than the human median intron size of 1334 bp



Lower panels show the scatter plot of logarithm (base 2) of gene RPKMs between the six sequencing experiments. Data points are smoothed color density representation of the scatterplot, obtained through a kernel density estimate. Upper panels show the square of the correlation (R2) between two distributions. Technical replicates show high correlation (Figure 7)

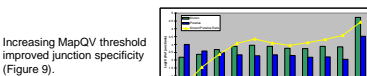
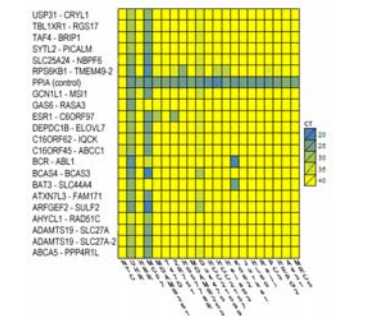


Figure 10. Fusion Assays Tested in 20 Cancer Cell Lines



In UHR, we identified the previously described gene fusions including BCR-ABL1, GAS6-RASA3, ARFGEF2-SULF2, NUP214-XKR3 and BAT3-SLC44A4. We prepared TaqMan assays for these three fusions as well as selected MCF-7 gene fusions and tested them in twenty cancer cell lines (Figure 10)

Two 'adjacent gene' fusions were expressed in multiple samples: ESR1-C6ORF97 and RPS6KB1-TMEM49. The fusion between the estrogen receptor alpha gene ESR1 and its neighboring gene C6ORF97 on Chr.6 was expressed in two other ER+ breast cancer cell lines in addition to MCF-7 and Du4475. This fusion may have occurred due to a rearrangement or trans-splicing. We further tested these fusions in 48 clinical normal and tumor breast cancer samples, and ESR1 fusion was found expressed in only one sample.

CONCLUSIONS

Whole transcriptome paired-end sequencing with the SOLiD system v4.0 analyzed with BioScope software allows easy, low-cost, genome-wide sensitive and specific detection of gene fusions, including novel gene fusions. This allows interrogation of large numbers of tumor samples and detection and discovery of biologically important gene fusions. Comprehensive exon junction detection within genes suggests splice variants and facilitates prediction of alternative splicing. RNA-seq also measures gene expression, and allows quantitation of changes in gene expression between samples with large dynamic range. Strand-specific sequencing disentangles expression of proximate exons on opposite strands.

REFERENCES

1. Maher et al (2009) Transcriptome sequencing to detect gene fusions in cancer. Nature 458(7234):97-101
2. Hampton et al (2009) A sequence-level map of chromosomal breakpoints in the MCF-7 breast cancer cell line yields insights into the evolution of a cancer genome. Genome Research 19:167-177
3. Krzywinski et al (2009) Circo: an information aesthetic for comparative genomics. Genome Research 19(9): 1639-1645

ACKNOWLEDGEMENTS

We thank Benjamin Kong and Caifu Chen (Life Technologies) for designing and running the TaqMan validation of fusion transcripts. We thank Yundan Lou, Goke Ojewole, Brijesh Krishnaswami and the software team for software integration and verification.

TRADEMARKS/LICENSING

© 2010 Life Technologies Corporation. All rights reserved. The trademarks mentioned herein are the property of Life Technologies Corporation or their respective owners. TaqMan is a registered trademark of Roche Molecular Systems, Inc.

For Research Use Only. Not intended for any animal or human therapeutic or diagnostic use.