# TopHat-Fusion: An algorithm for Discovery of Novel Fusion Transcripts

### Abstract

TopHat-Fusion is an enhanced version of TopHat [1] designed to handle transcripts resulting from fusion gene products. Fusion genes result from the breakage and re-joining of two different chromosomes, or from rearrangements within a chromosome. Because TopHat-Fusion does not require information about known genes, it can discover novel fusion products from known genes, unknown genes, or unannotated splice variants of known genes. Using RNAseq data collected from breast cancer cells [2], prostate cancer cells [3], and Universal Human Reference (UHR) cells [3], we detected multiple fusion genes including novel and previously-reported ones. TopHat-Fusion is free, opensource software available from tophat.cbcb.umd.edu.



Daehwan Kim and Steven L. Salzberg Center for Bioinformatics and Computational Biology University of Maryland, College Park, MD 20742

#### Methods

TopHat-Fusion implements several major changes to the original TopHat algorithm, all designed to enable discovery of fusion transcripts as shown in the below figure. The first step in analysis of an RNA-seq data set is to align (map) the reads to the genome. For those initially unmapped (IUM) reads, we split each read into multiple of 25-bp pieces, and then use Bowtie to map the 25-bp segments to the genome. TopHat-Fusion looks for cases where the first and last segments are mapped to either (a) two different chromosomes, in the case of inter-chromosomal fusions; or (b) two locations on the same chromosome separated by a user-defined distance. The whole read is then used to identify a fusion point by re-aligning it starting from the mapped first and last segments.

data	sample type		fragment length	read length	number of fragments (or read	
	BT474	paired	100 & 200	50	21,423,69	
Edgren et al.	SKBR3	paired	100 & 200	50	18,140,24	
	KPL4	paired	100	50	6,796,44	
	MCF7	paired	100	50	8,409,78	
	VCaP	paired	300	50	16,894,52	
Maher et al.	UHR	paired	300	50	25,294,16	
	UHR	single		100	56,129,47	

remapped against those fusion points (as well as intron boundaries and indels), and the

TopHat-Fusion's initial output includes many fusion candidates, most of which are induced by repeats or simply due to chance, with few supporting reads and mate pairs. We applied

(1) At least 5 fusion spanning reads or supporting mate pairs for the breast cancer cell lines (BT474, SKBR3, KPL4,

As pointed out in Edgren et al. [2], true fusion transcripts have reads mapping uniformly in a wide window across the fusion point, whereas false positive fusions are narrowly covered. Using this criterion, we examine a 600-bp window around each fusion (300-bp each side), and we reject fusion candidates for which the reads fail to cover this window without any "big" gaps. The final step is to sort candidates based on how well reads are distributed around a fusion point.





#### Four breast cancer cell lines (BT474, SKBR3, KPL4, MCF7) from Edgren [2]



#### VCaP data from Maher et al. [3]

5′ gene	5' chromosome	5' position	3′ gene	3' chromosome	3' position	spanning reads	supporting pairs	four
ZDHHC7	chr16	85023908	ABCB9	chr12	123444867	13	69	yes
TMPRSS2	chr21	42879875	ERG	chr21	39817542	7	285	yes
HJURP	chr2	234749254	EIF4E2	chr2	233421125	3	9	yes
VWA2	chr10	116008521	PRKCH	chr14	61909826	1	10	
RGS3	chr9	116299195	PRKAR1B	chr7	699055	3	11	
SPOCK1	chr5	136397966	TBC1D9B	chr5	179305324	9	31	yes
LRP4	chr11	46911864	FBXL20	chr17	37557613	5	9	
INPP4A	chr2	99193605	HJURP	chr2	234746297	6	12	yes
C16orf70	chr16	67144140	C16orf48	chr16	67700168	2	19	
NDUFV2	chr18	9102729	ENSG00000188699	chr19	53727808	1	35	
NEAT1	chr11	65190281	ENSG00000229344	chr1	568419	1	17	
ENSG0000011405	chr11	17229396	TEAD1	chr11	12883794	7	9	yes
USP10	chr16	84733713	ZDHHC7	chr16	85024243	1	22	yes
LMAN2	chr5	176778452	AP3S1	chr5	115202366	15	2	yes
WDR45L	chr17	80579516	ENSG00000224737	chr17	30439195	1	33	
RC3H2	chr9	125622198	RGS3	chr9	116299072	3	11	yes
CTNNA1	chr5	138145895	ENSG00000249026	chr5	114727795	1	12	
ENSG00000229880	chr21	46097128	IMMT	chr2	86389185	1	50	
ENSG00000214009	chrX	45918367	PCNA	chr20	5098168	1	24	

Maher et al. reported 11 fusion genes in the VCaP sample, 9 of which were identified by TopHat-Fusion. It missed a fusion between two overlapping genes ZNF577 and ZNF649 on chromosome 19. That fusion event appears to be due to readthrough transcription.

#### UHR single-end reads vs. paired-end reads [3]

Using four previously known fusion genes GAS6-RASA3, BCR-ABL1, ARFGEF2-SULF2, and BCAS4-BCAS3, we compared the results of single and paired-end reads from the UHR RNA-seq data. The number of supporting reads per million reads (RPM\*) indicates that single-end reads find the same fusions as paired-end reads. The two figures show read distributions around the BCR-ABL1 fusion for single (upper) and paired-end reads (lower), respectively.

type	5' gene	5' chromosome	5' position	3' gene	3' chromosome	3' position	spanning reads $(RPM^*)$	supporting pairs
single	GAS6	chr13	114529968	RASA3	chr13	114751268	15 (0.267)	
paired	GAS6	chr13	114529968	RASA3	chr13	114751268	$10 \ (0.198)$	43
single	BCR	chr22	23632599	ABL1	chr9	133655755	6 (0.107)	
single	BCR	chr22	23632599	ABL1	chr9	133729450	3 (0.053)	
paired	BCR	chr22	23632599	ABL1	chr9	133655755	2(0.040)	7
paired	BCR	chr22	23632599	ABL1	chr9	133729450	3 (0.059)	10
single	ARFGEF2	chr20	47538548	SULF2	chr20	46365683	17 (0.302)	
paired	ARFGEF2	chr20	47538545	SULF2	chr20	46365686	$10 \ (0.198)$	30
single	BCAS4	chr20	49411707	BCAS3	chr17	59445685	25 (0.445)	
nainad	DCASA	ab m20	40411707	DCAC2	ohn17	50445695	12(0.957)	145

## Acknowledgements

We would like to thank Ryan Kelly for his indel-finding algorithm, which he added to TopHat and which we had used as part of the fusion finding algorithm. This work is supported in part by the National Institutes of Health under grants R01-LM006845 and R01-GM083873 to SLS.

#### References

1. Trapnell, C., L. Pachter, and S.L. Salzberg, TopHat: discovering splice junctions with RNA-Seq. Bioinformatics, 2009. 25(9): p. 1105-11. 2. Edgren, H., et al., Identification of fusion genes in breast cancer by paired-end RNA-sequencing. Genome

biology, 2011. **12**(1): p. R6.

3. Maher, C.A., et al., Chimeric transcript discovery by paired-end transcriptome sequencing. Proceedings of the National Academy of Sciences of the United States of America, 2009. 106(30): p. 12353-8.

TopHat-Fusion found 25 out of 27 fusion genes reported in the four breast cancer cells from Edgren et al [2], where they are indicated by "found\*" in the left table. One of the two missing fusions, DHX35-ITCH, was included in the initial output, but it was filtered out because it was only supported by one read and one



The above figure shows read distribution around the BCAS4-BCAS3 fusion in MCF7 cell

