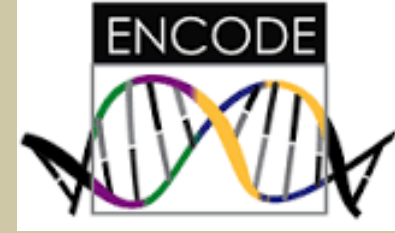




**Indiana State
University**



**Evaluation of U12-type non-canonical
splicing in human ENCODE RNA-Seq
datasets and analysis of biological
functions for spliced sequences by
Read-Split-Fly algorithm**

Yongsheng Bai and Jeff Kinne

Department of Biology

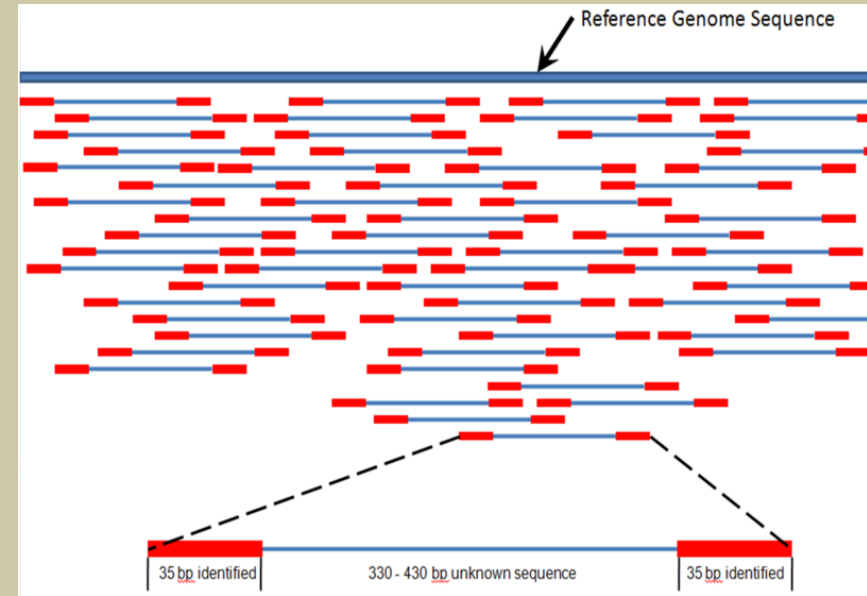
Department of Mathematics & Computer Science



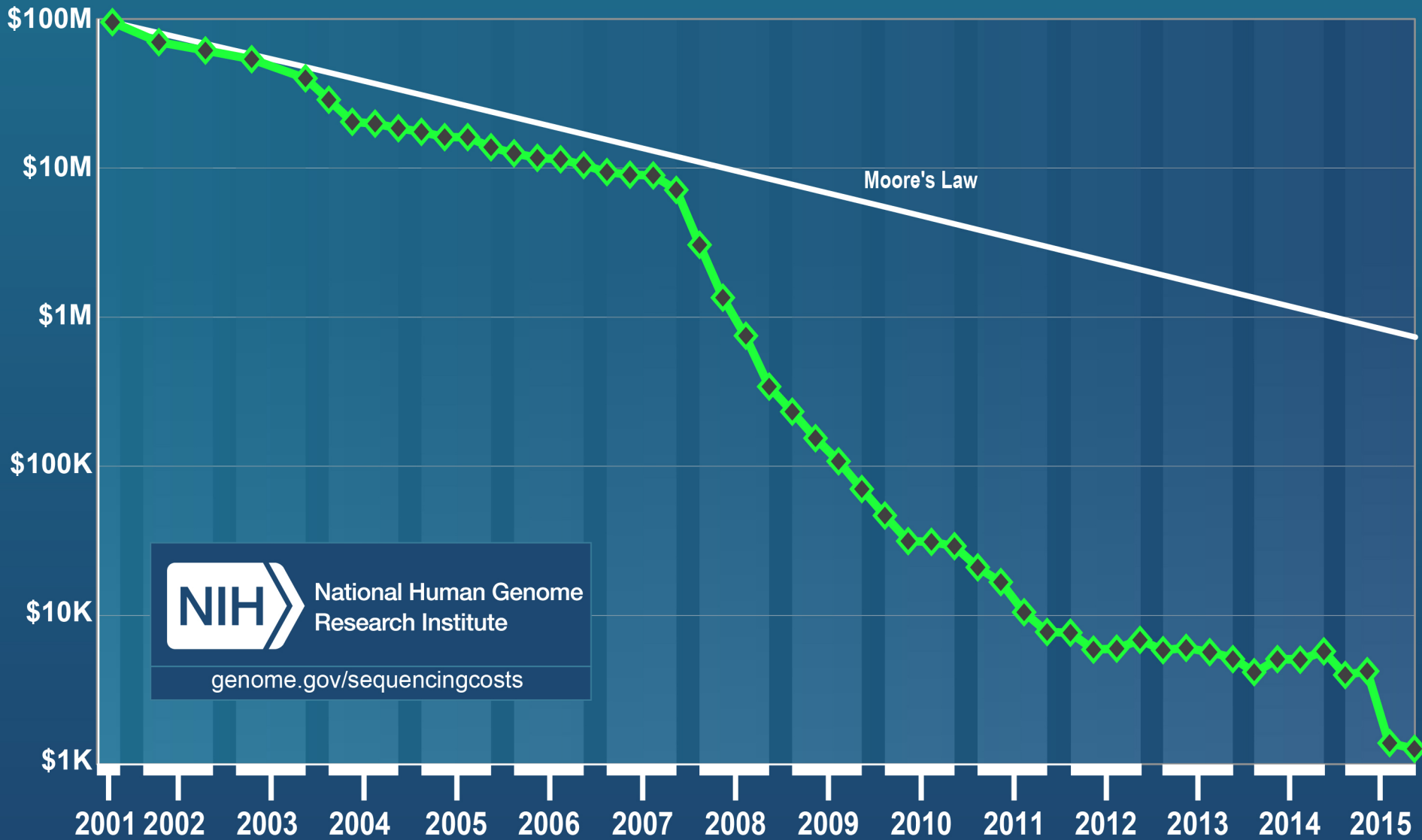
June 8, 2016



Human Genome Project and Next-Generation Sequencing

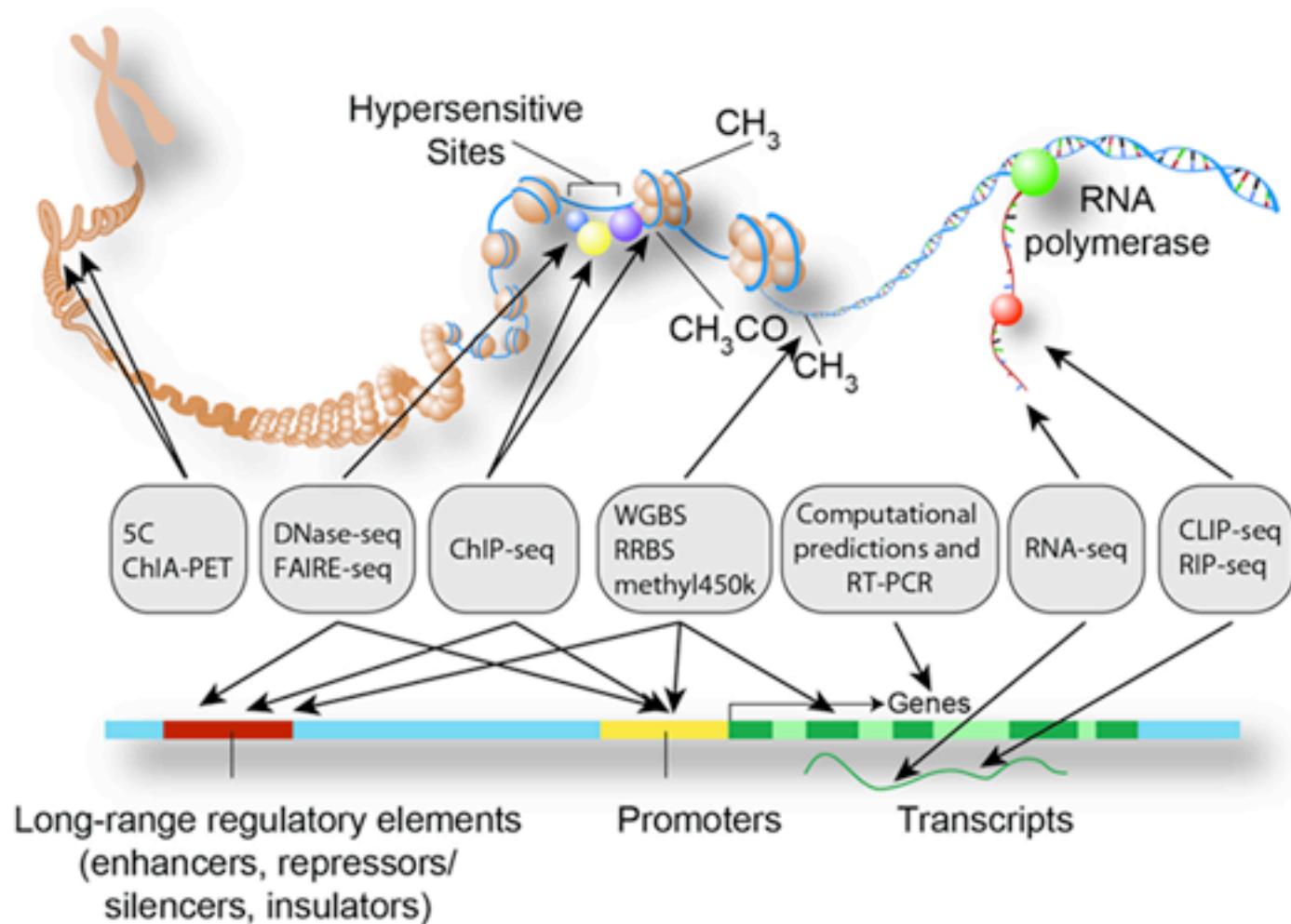


Cost per Genome



ENCODE Project

(<https://www.encodeproject.org/>)

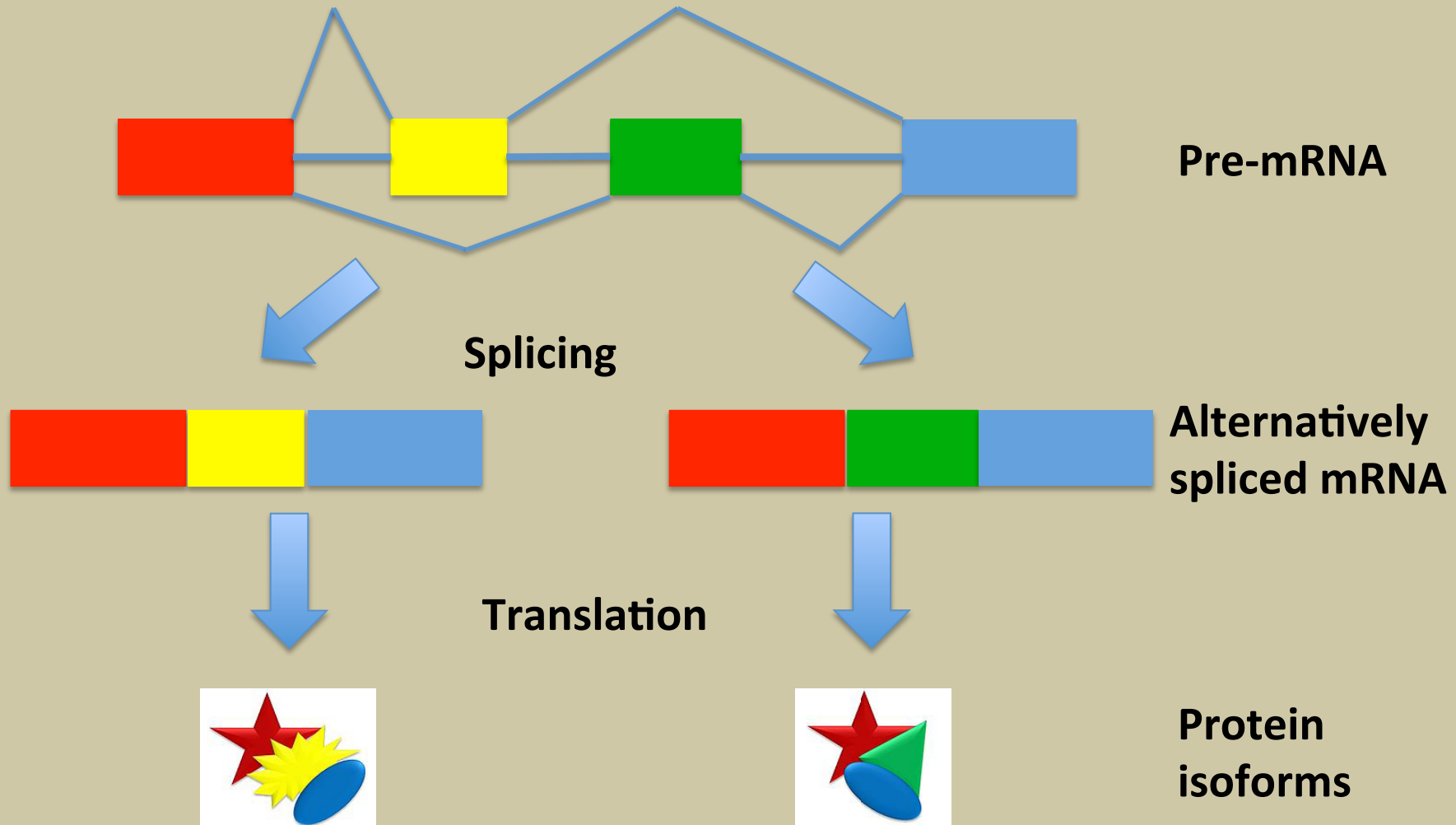


Datasets – ENCODE

(<https://www.encodeproject.org/>)

Dataset Classification	Experiment	File Name	Run type	Forward Or Reverse Mate	Biological Replicate	Technical Replicate
Paired-Ended with Biological Replicate	ENCSR468ION	ENCFF002DJH	paired-ended	1	1	1
	ENCSR468ION	ENCFF002DJU	paired-ended	2	1	1
	ENCSR468ION	ENCFF002DJV	paired-ended	1	2	1
	ENCSR468ION	ENCFF002DJX	paired-ended	2	2	1
Paired-Ended without Biological Replicate	ENCSR001UXR	ENCFF576OBS	paired-ended	1	1	1
	ENCSR001UXR	ENCFF063OPQ	paired-ended	2	1	1
Single-Ended with Biological Replicate	ENCSR000BYG	ENCFF000MVL	single-ended	1	1	1
	ENCSR000BYG	ENCFF000MVN	single-ended	1	2	1

Transcription and Translation



RSR Pipeline



The original code for the RSW algorithm worked by considering all pairs of parts of unmapped reads at once. The software was written in **Perl**.



Results are computed much faster in RSR by considering each unmapped read separately. The algorithm was implemented in **C++**.

Comparison of number of reads for supporting *Xbp1* 26 nt spliced regions reported by RSR and other tools

	500 nM Thapsigargin (Tg)		1 mM Dithiothreitol (Dtt)	
Software	Het (Ire1 α +/-)	KO (Ire1 α -/-)	Het (Ire1 α +/-)	KO (Ire1 α -/-)
Read-Split-Run (RSR)	21	0	173	0
TopHat	23	86	59	289
BWA	0	0	67	0
Bowtie2	0	0	171	0
STAR	0	0	0	0
Alt Event Finder	0	0	0	0

(Bai *et al.*, 2016. To appear in BMC Genomics)

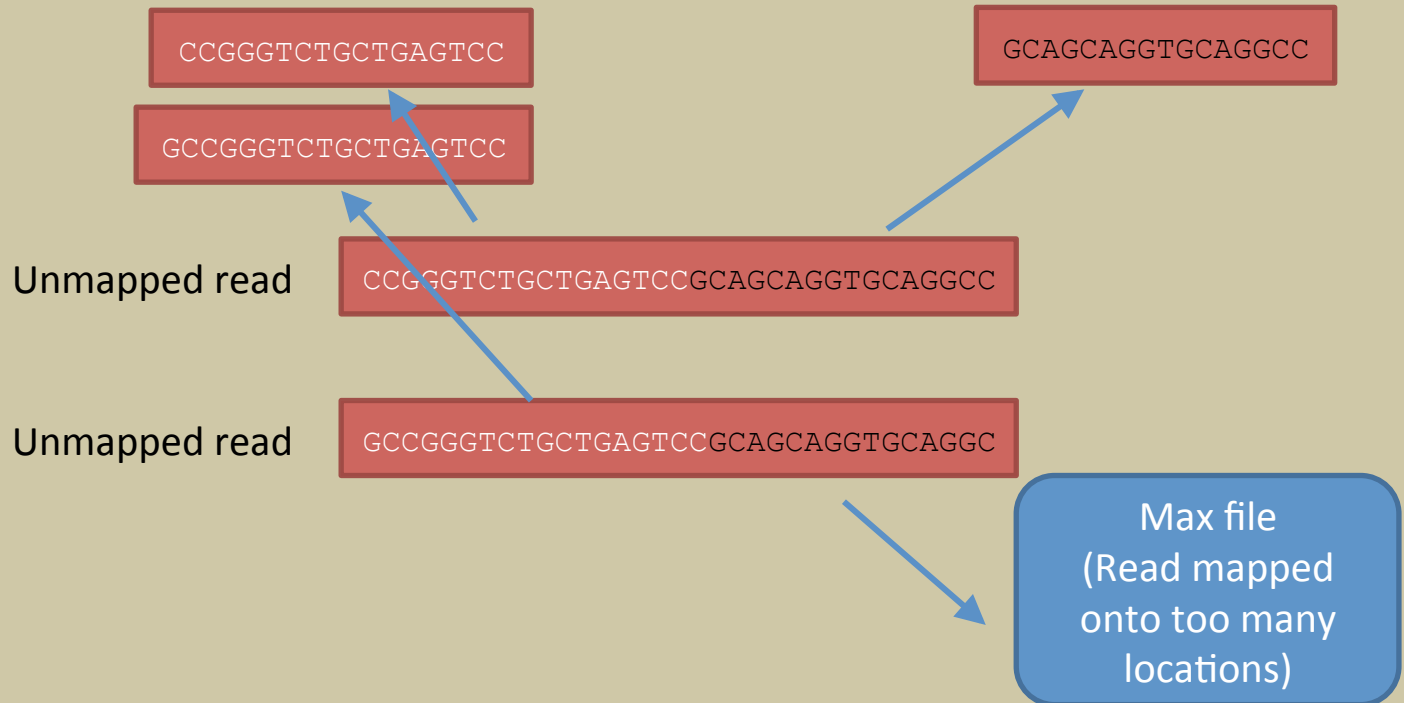
RSF Pipeline – Improved algorithm

- Rescue read halves from bowtie max file

Mouse mm9

Chr 11, position 5,424,242

GGGAGTGGAGTAAGGCTGGTGGCCGGGTCTGCTGAGTCCGCAGCACTCAGACTATGTGCACCTCTGCAGCAGGTGCAGGCCAGT



RSF Pipeline – Improved algorithm

- Increased sensitivity
 - Rescue read halves from bowtie max file

<i>Xbp1</i> splice Supporting reads	Read-Split-Run	Read-Split-Fly
Tg treated sample	21	27
Dtt treated sample	173	209

RSF Pipeline – Improved algorithm

- Increased performance (Memory and Time Usage)

Memory (GB) / Time (Hours) Usage	Read-Split-Run	Read-Split-Fly	Change (%)
ENCSR558LHB	123.2 / 15.8	18.9 / 11.5	-84.7 / -27.2
ENCSR905FLM	110.2 / 13.5	20.7 / 12.6	-81.2 / -6.7
ENCSR000CQY	55.2 / 0.55	46.5 / 0.42	-15.8 / -24

Extracting Spliced Sequences

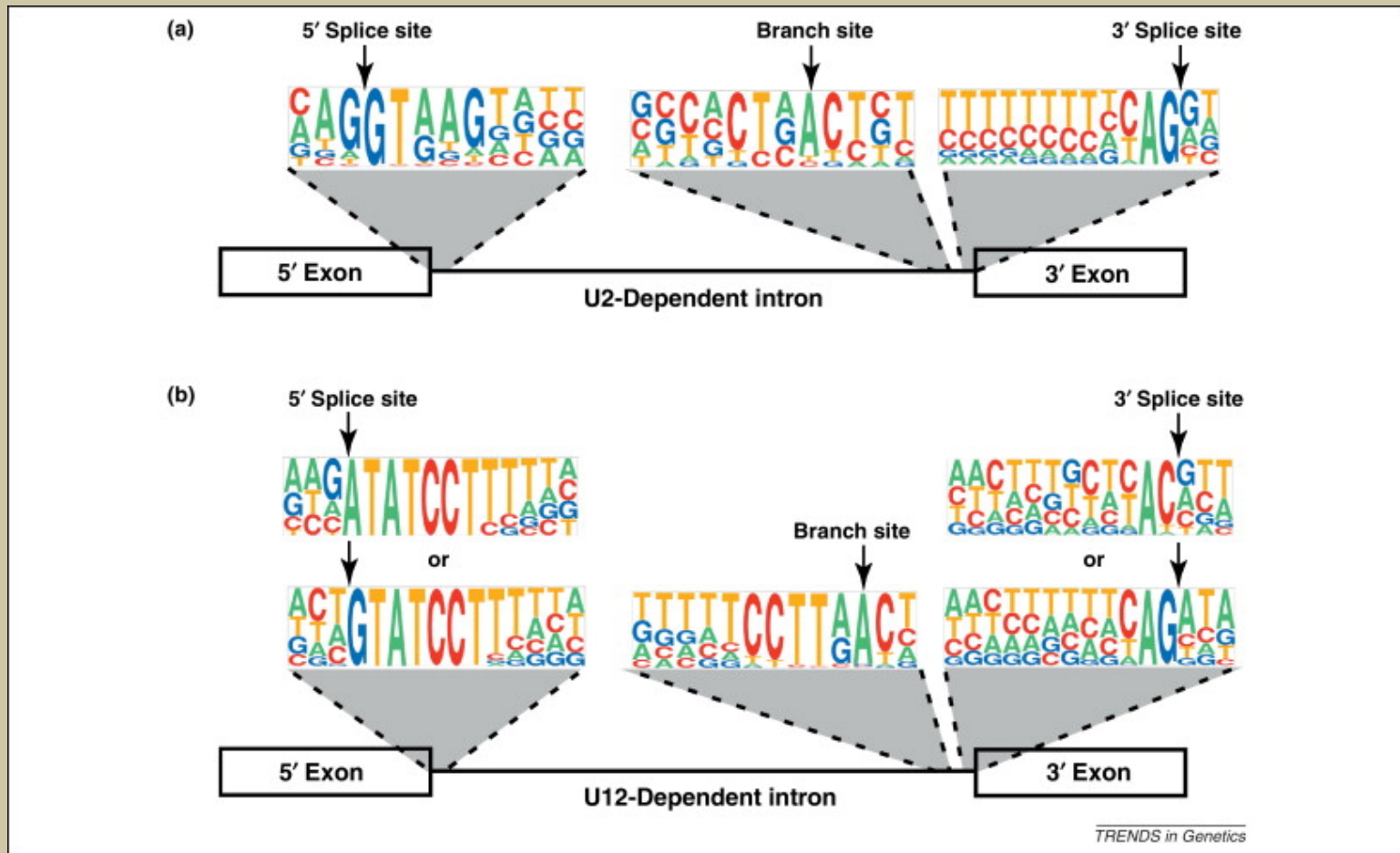
- Sample data before and after bowtie-inspect-RSR file processing:

GeneName	Chromosome	# supporting reads	splice length	range of supporting reads	Novel or not (*)	Bracketed sequence
Xbp1	chr11	22	26	5424280--5424312	Novel	GGGTCTGCTGAGTCCGCA--GCAGGTGCAGGCCAGT
Ddx41	chr13	2	74	55635449--55635524	*	CCGAGCCAGCTCTCGC--GAGGGGCAGATGATGAGCC
Nol7	chr13	4	76	43494024--43494104	*	GCGAGTGCTCGCA--GGGATAAACGCTGTTGAAGGA
Gtpbp4	chr13	2	80	8991097--8991177	*	GCACATAATCTTTAGCAACA--TTGTCCACCAAGTTT
Pdlim7	chr13	5	81	55610171--55610255	*	TGCTCTGAACAGGCTGGGCT--CTGCTAAGACCCAGG

GeneName	Chromosome	# supporting reads	splice length	range of supporting reads	Novel or not (*)	Bracketed sequence	spliced sequence
Xbp1	chr11	22	26	5424279--[5424312	Novel	GGTCTGCTGAGTCC]--[GGTGCAGGCCAGTT	GCAGCACTCAGACTATGTGCACCTCTGCAGCA
Ddx41	chr13	2	74	55635448]--[55635524	*	CCGAGCCAGCTCTCG]--[GAGGGGCAGATGATG	CTACAGGCAGAGACAAGGCTGGTACCGGAAGCCAGCACAGTCTTCTCCCGACCCAGGATTAGCCTATCTTACC
Nol7	chr13	4	76	43494023]--[43494104	*	GCGCGCAGTGCTCG]--[GATAAACGCTGTTG	CAGGTGCAGCGCCGGGAATGTTGGCGCTCGGGACTCCCCAGGCTCGGCCGGCGCTGACCCGACTCCTTCCCTCCAGG
Gtpbp4	chr13	2	80	8991096]--[8991177	*	TAATCTTTAGCAACA]--[TTGTCCACCAAGTTT	CTTAGAAAAAGAGAAAGTAACTAAGGGTTTGTATTTTTGGTAATAACTTATAGATATATAGTAAGAAAAATGCCTTAC
Pdlim7	chr13	5	81	55610170]--[55610255	*	TGAACAGGCTGGGCT]--[CTAAGACCCAGGCTG	CTGGAGAAGAAGAGAGACGACTCACAGCAGGTAAAGAGACCAGAGGCTATCCACACCCCTTCTGCCTGTGAGCACCTGTACCTG

- Consequently, the spliced sequence fits snugly within the “range of supporting reads.”
 - For example, the ***XBP1*** candidate spliced sequence begins at index 5424280 and ends at index 5424311

Consensus sequences of human U12- and U2-type intron



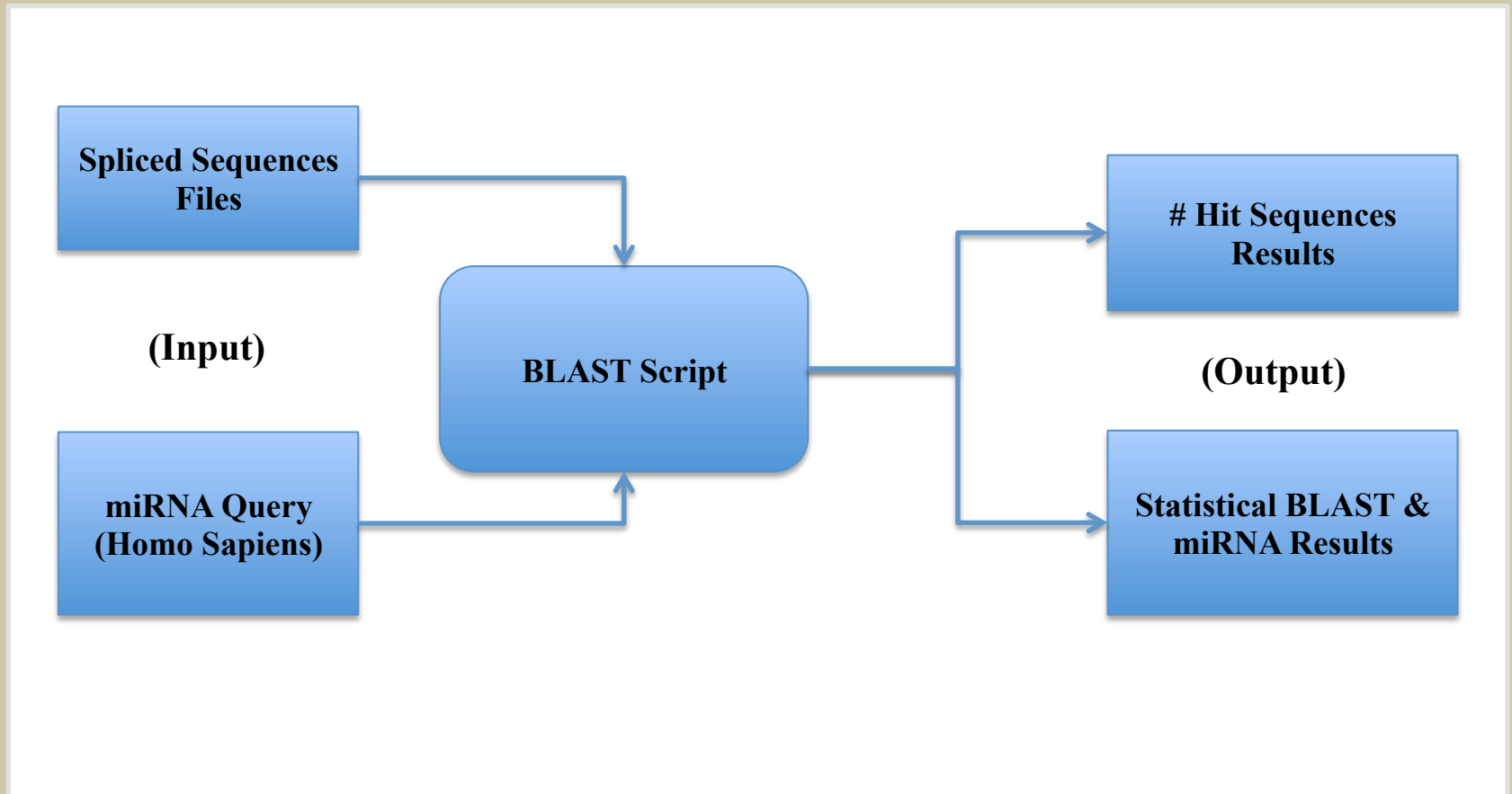
(Padgett, R. 2012)

U12-type non-canonical splicing in 21 human ENCODE RNA-Seq datasets

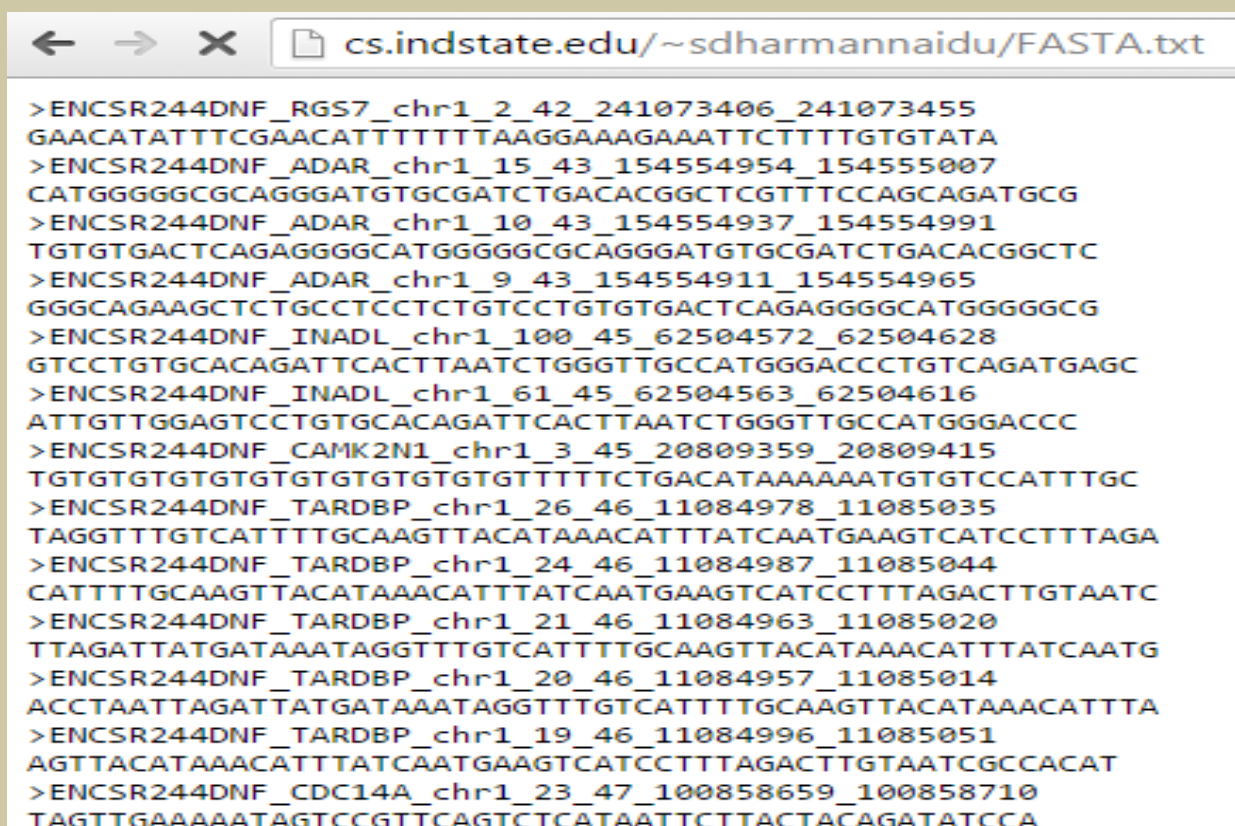
E-Value 1

	u12		u12 Total	u2		u2 Total	Grand Total
Row Labels	known	novel		known	novel		
3pFull	52977	36760	89737	13094	8295	21389	111126
5pFull	8840	11618	20458	3466	2140	5606	26064
branch	222	397	619	NA	NA	NA	619
branch#extend	12968	13342	26310	NA	NA	NA	26310
Grand Total	75007	62117	137124	16560	10435	26995	164119

RSF -BLAST-Workflow



Sample DB File

A screenshot of a web browser window displaying a FASTA file. The browser's address bar shows the URL `cs.indstate.edu/~sdharmannaidu/FASTA.txt`. The main content area displays a list of DNA sequence entries, each starting with a greater-than sign and a header. The entries include identifiers like `ENCSR244DNF_RGS7_chr1_2_42_241073406_241073455` and their corresponding nucleotide sequences. The sequences consist of uppercase letters A, C, G, and T.

```
>ENCSR244DNF_RGS7_chr1_2_42_241073406_241073455
GAACATATTTTTCGAACATTTTTTTAAAGGAAAGAAATTCTTTTGTGTATA
>ENCSR244DNF_ADAR_chr1_15_43_154554954_154555007
CATGGGGGCGCAGGGGATGTGCGATCTGACACGGCTCGTTTCCAGCAGATGCG
>ENCSR244DNF_ADAR_chr1_10_43_154554937_154554991
TGTGTGACTCAGAGGGGCGATGGGGGCGCAGGGATGTGCGATCTGACACGGCTC
>ENCSR244DNF_ADAR_chr1_9_43_154554911_154554965
GGGCAGAAGCTCTGCCTCCTCTGTCCTGTGTGACTCAGAGGGGCGATGGGGGCG
>ENCSR244DNF_INADL_chr1_100_45_62504572_62504628
GTCCTGTGCACAGATTCACTTAATCTGGGTTGCCATGGGACCCTGTCAGATGAGC
>ENCSR244DNF_INADL_chr1_61_45_62504563_62504616
ATTGTTGGAGTCCTGTGCACAGATTCACTTAATCTGGGTTGCCATGGGACCC
>ENCSR244DNF_CAMK2N1_chr1_3_45_20809359_20809415
TGTGTGTGTGTGTGTGTGTGTGTGTGTGTTTTCTGACATAAAAAAATGTGTCCATTTGC
>ENCSR244DNF_TARDBP_chr1_26_46_11084978_11085035
TAGGTTTGTCATTTTGCAAGTTACATAAAACATTTATCAATGAAGTCATCCTTTAGA
>ENCSR244DNF_TARDBP_chr1_24_46_11084987_11085044
CATTTTGCAAGTTACATAAAACATTTATCAATGAAGTCATCCTTTAGACTTGTAATC
>ENCSR244DNF_TARDBP_chr1_21_46_11084963_11085020
TTAGATTATGATAAATAGGTTTGTCATTTTGCAAGTTACATAAAACATTTATCAATG
>ENCSR244DNF_TARDBP_chr1_20_46_11084957_11085014
ACCTAATTAGATTATGATAAATAGGTTTGTCATTTTGCAAGTTACATAAAACATTTA
>ENCSR244DNF_TARDBP_chr1_19_46_11084996_11085051
AGTTACATAAAACATTTATCAATGAAGTCATCCTTTAGACTTGTAATCGCCACAT
>ENCSR244DNF_CDC14A_chr1_23_47_100858659_100858710
TAGTTGAAAAATAGTCCGTTTCAGTCTCATAAATTCTTACTACAGATATCCA
```


Query – Sample Data

>hsa-let-7a-5p MIMAT0000062 Homo sapiens let-7a-5p
TGAGGTAGTAGGTTGTATAGTT

>hsa-let-7a-3p MIMAT0004481 Homo sapiens let-7a-3p
CTATACAATCTACTGTCTTTC

>hsa-let-7a-2-3p MIMAT0010195 Homo sapiens let-7a-2-3p
CTGTACAGCCTCCTAGCTTTCC

Summary

- Preliminary results from 21 samples of ENCODE datasets show that there are several miRNAs are prevalent (> 50%) in studied ENCODE samples. Two of them (*hsa-miR-1273d* and *hsa-miR-548*) are associated with many diseases as suggested in the literature.
- U12-type non-canonical splicing could be noteworthy in ENCODE datasets.
- RSF for identifying U12-type splicing events using ENCODE datasets is applicable to study a range of diseases across biological systems under different experimental conditions.



<http://bioinf1.indstate.edu/RSR>

Read-Split-Run

A pipeline for detecting non-canonical spliced-regions in RNA-Seq data.

Experiment Settings

Mode

Non-comparative

Reads Type

Single

Experiment Replicates

1

Data

Select Genome Reference

mm9

Select RNA-Seq File (FASTQ format)

Tg_Het.txt

Select file...

Quality Encoding

Solexa

Check files

Length of Reads (in bp)

35

Pre-processing

Minimum Split Size (2+)
bp

Maximum Good
Alignments Allowed Per
Read

Candidate Selection

Minimum Distance
Between Candidate
Pairs (2 to 12) bp

Maximum Distance
Between Candidate
Pairs (bp)

Read Mapping Region
Boundary Buffer Size (0
to 11) bp

Minimum Number of
Supporting Reads

Output

Email:





weblogo.berkeley.edu

