

Introduction to RNA-Seq in Galaxy

HSPH Bioinformatics Core

<intro>



Oliver Hofmann



Shannan Ho Sui



John Hutchinson



Lorena Pantano



Meeta Mistry



John Morrissey



Rory Kirchner



Brad Chapman



Radhika Khetani



Mary Piper



Andreas Sjödin



Peter Kraft



Oliver Hofmann



Shannan Ho Sui



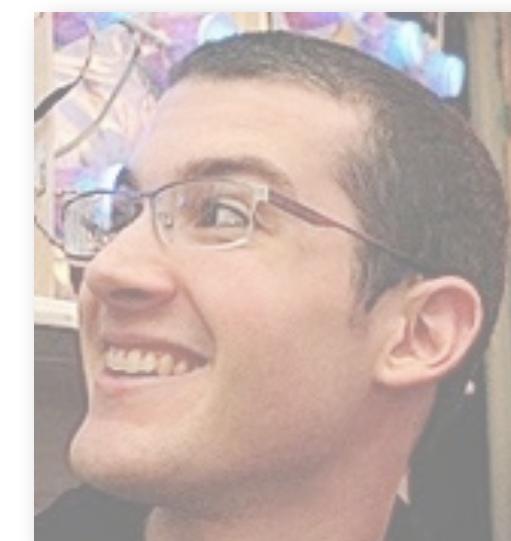
John Hutchinson



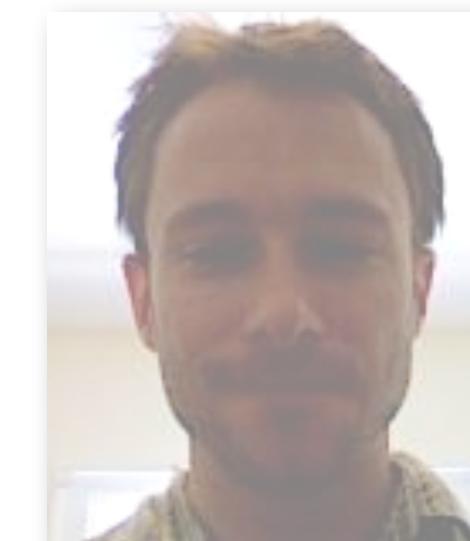
Lorena Pantano



Meeta Mistry



John Morrissey



Rory Kirchner



Brad Chapman



Radhika Khetani



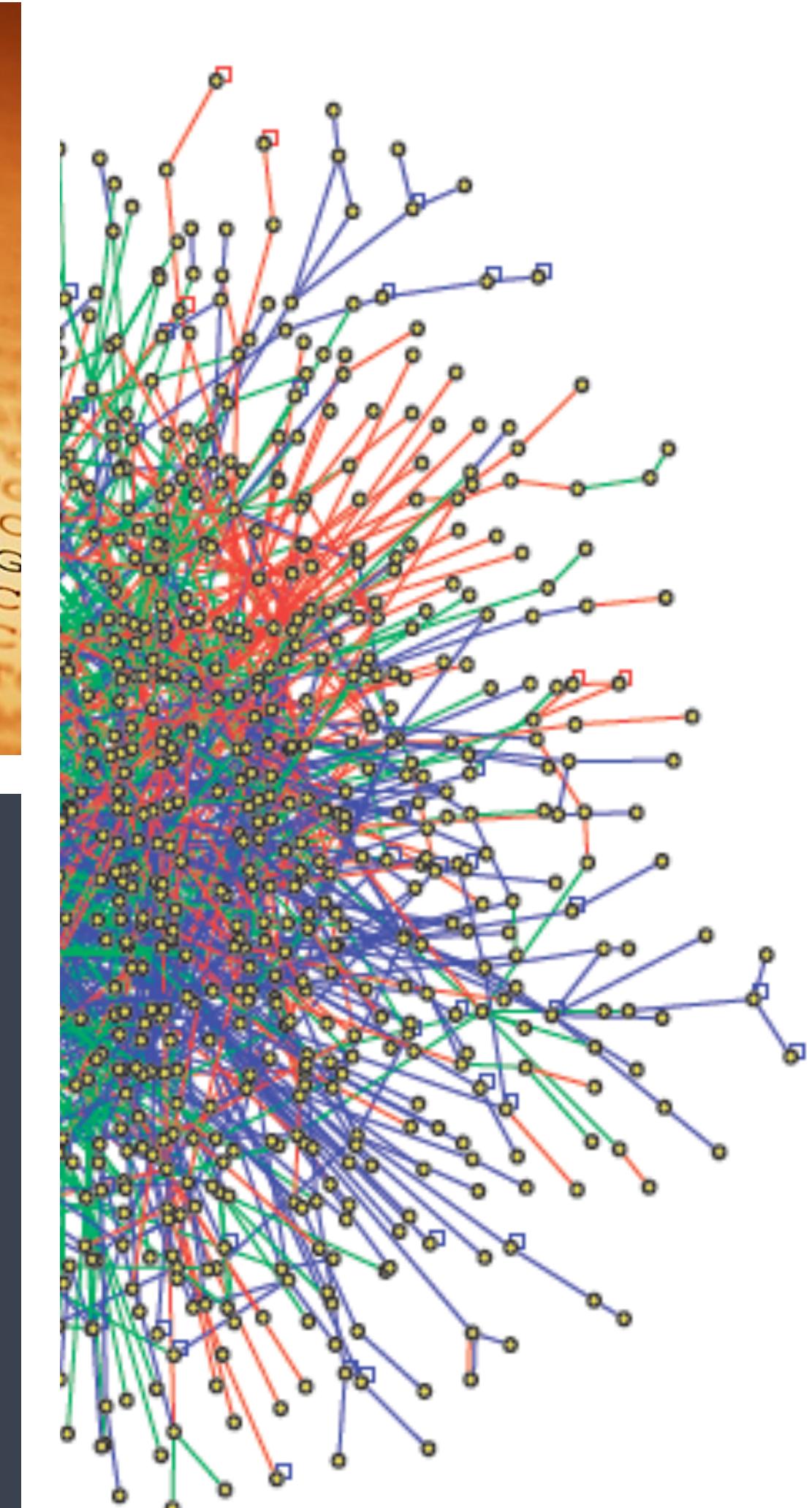
Mary Piper



Andreas Sjödin

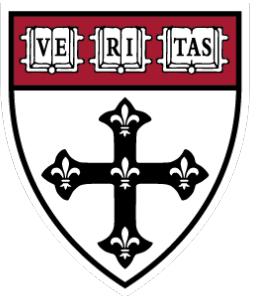


Peter Kraft



Three focus areas

- Research computing
- Next-gen sequencing
- Functional significance



HARVARD
SCHOOL OF PUBLIC HEALTH

HSCI
HARVARD STEM CELL
INSTITUTE

 **HARVARD CATALYST**
THE HARVARD CLINICAL
AND TRANSLATIONAL
SCIENCE CENTER

 **HARVARD**
MEDICAL SCHOOL

NIEHS / CFAR
Bioinformatics
Core

Center for
Stem Cell
Bioinformatics

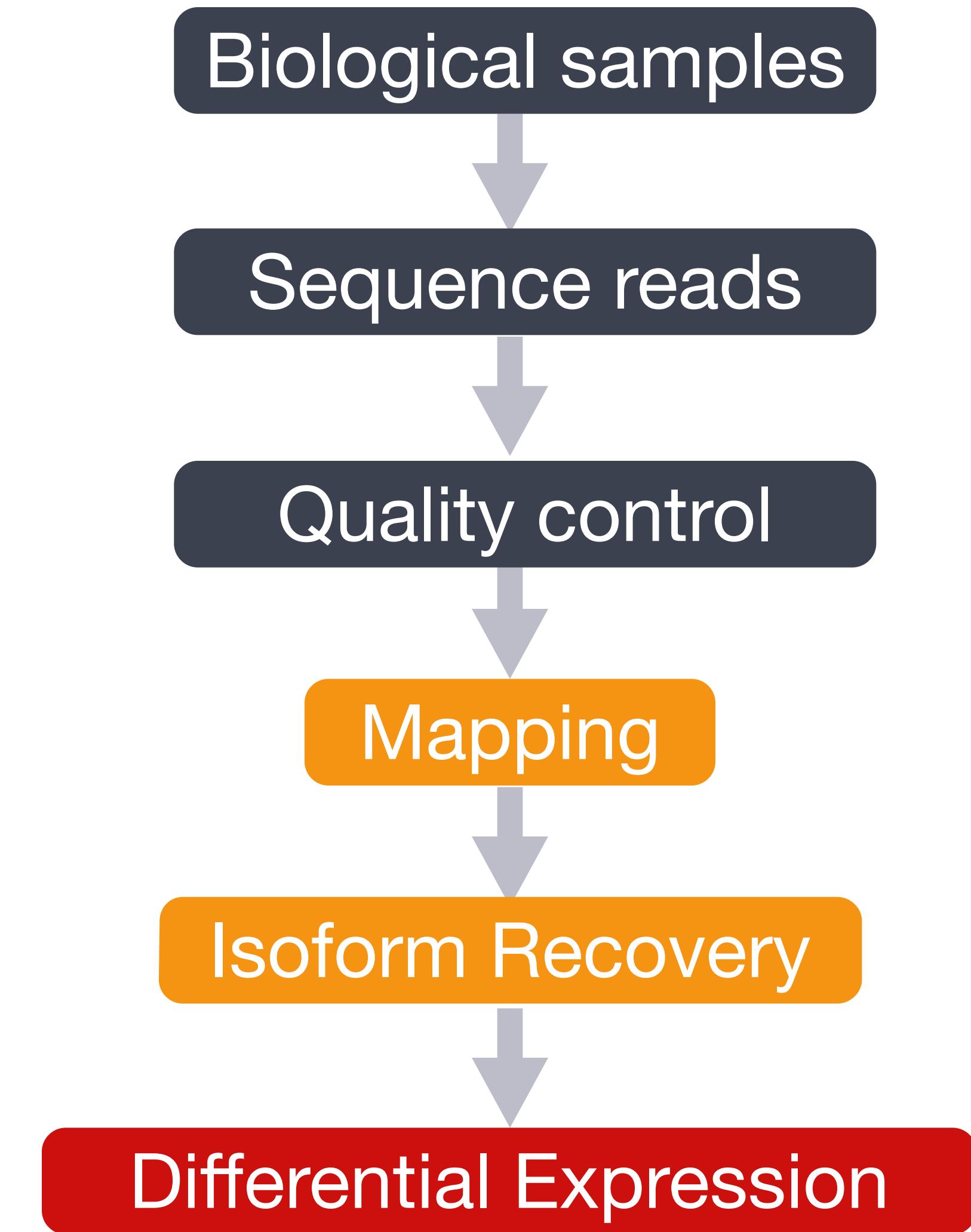
Harvard
Catalyst
Bioinformatics
Consulting

HMS
Tools &
Technology

Harvard
NeuroDiscovery
Center

Scope

$$\frac{dP}{dt} = \frac{1}{C} \frac{1}{P} \frac{dP}{dt}$$
$$\frac{1}{2} \frac{P_0 - P}{P} \sim \frac{1}{P}$$
$$\frac{P_0 - P}{P_0} \sim \frac{1}{t}$$
$$10^{-53}$$
$$10^{-26}$$
$$10^8 \text{ L.J}$$
$$10^{10} (10^{11}) \text{ J}$$



Workflows

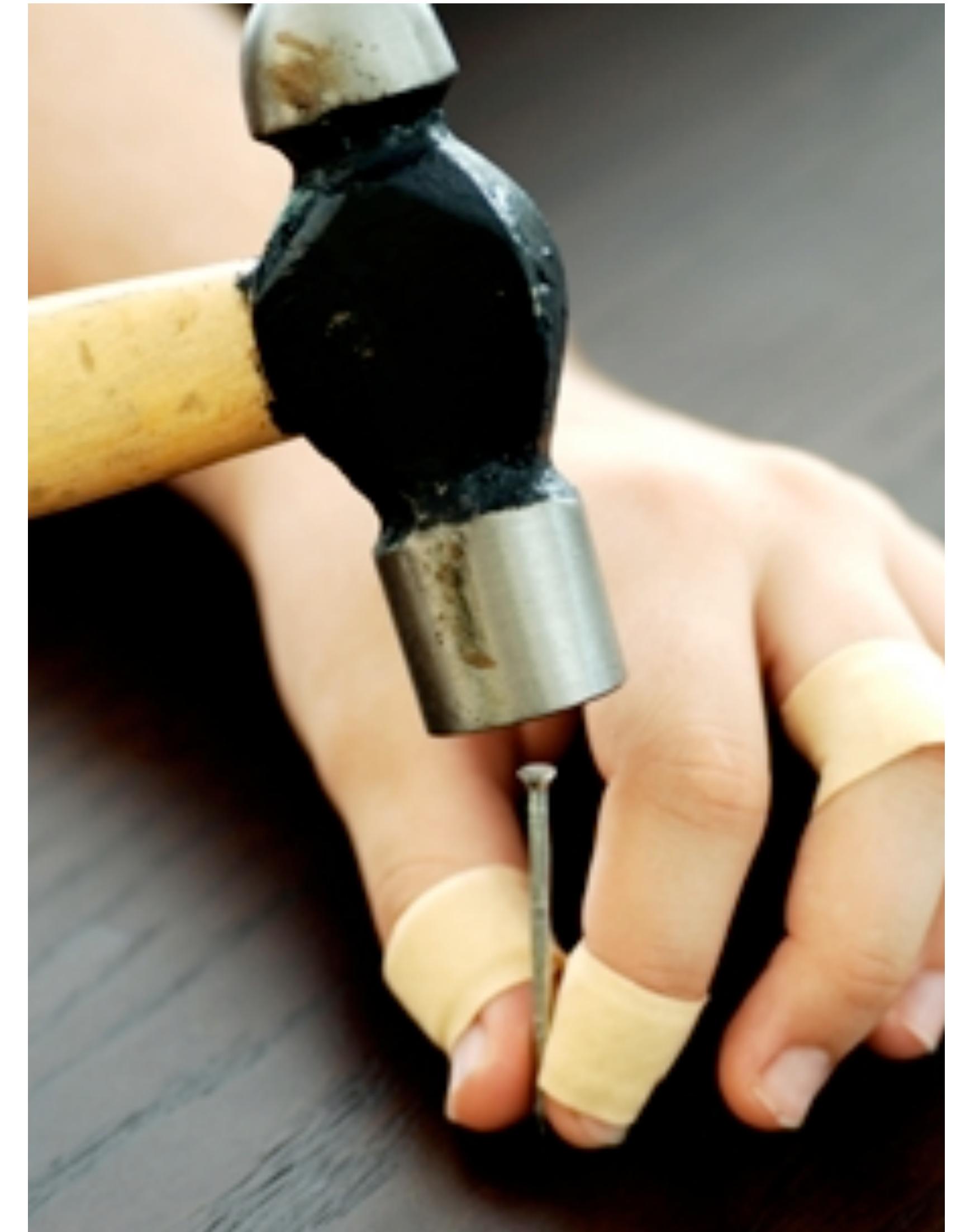
Do-it yourself

No black box

More control

Better understanding of the experiment

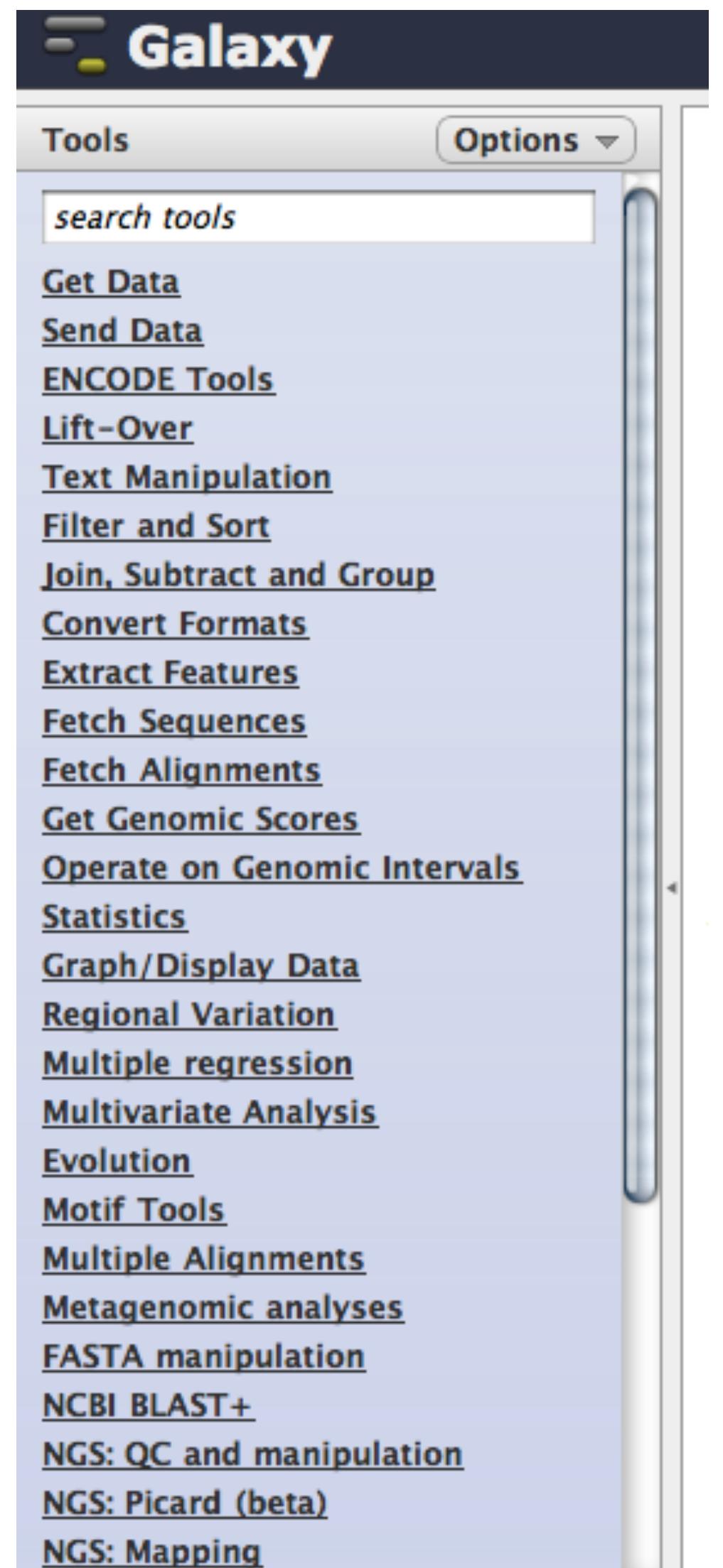
No command line / UNIX experience required

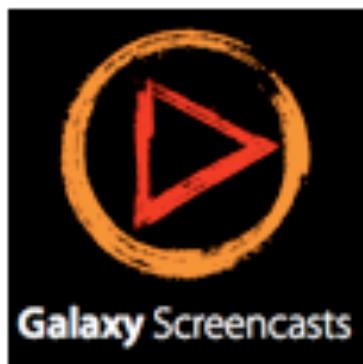


Using Galaxy

EC2 lab environment

<http://23.23.134.25/>





Galaxy Screencasts: The best way to understand how...

ChIP-seq 101
a simple example

Exons and SNPs
a bit more complexity

Saving & Sharing
preserving your data

Workflows...
if you don't like to repeat
yourself

... from scratch
building workflows

DNA
fetching sequences and
alignments

Online tutorials and screencasts

<http://usegalaxy.org>

<http://galaxyproject.org>

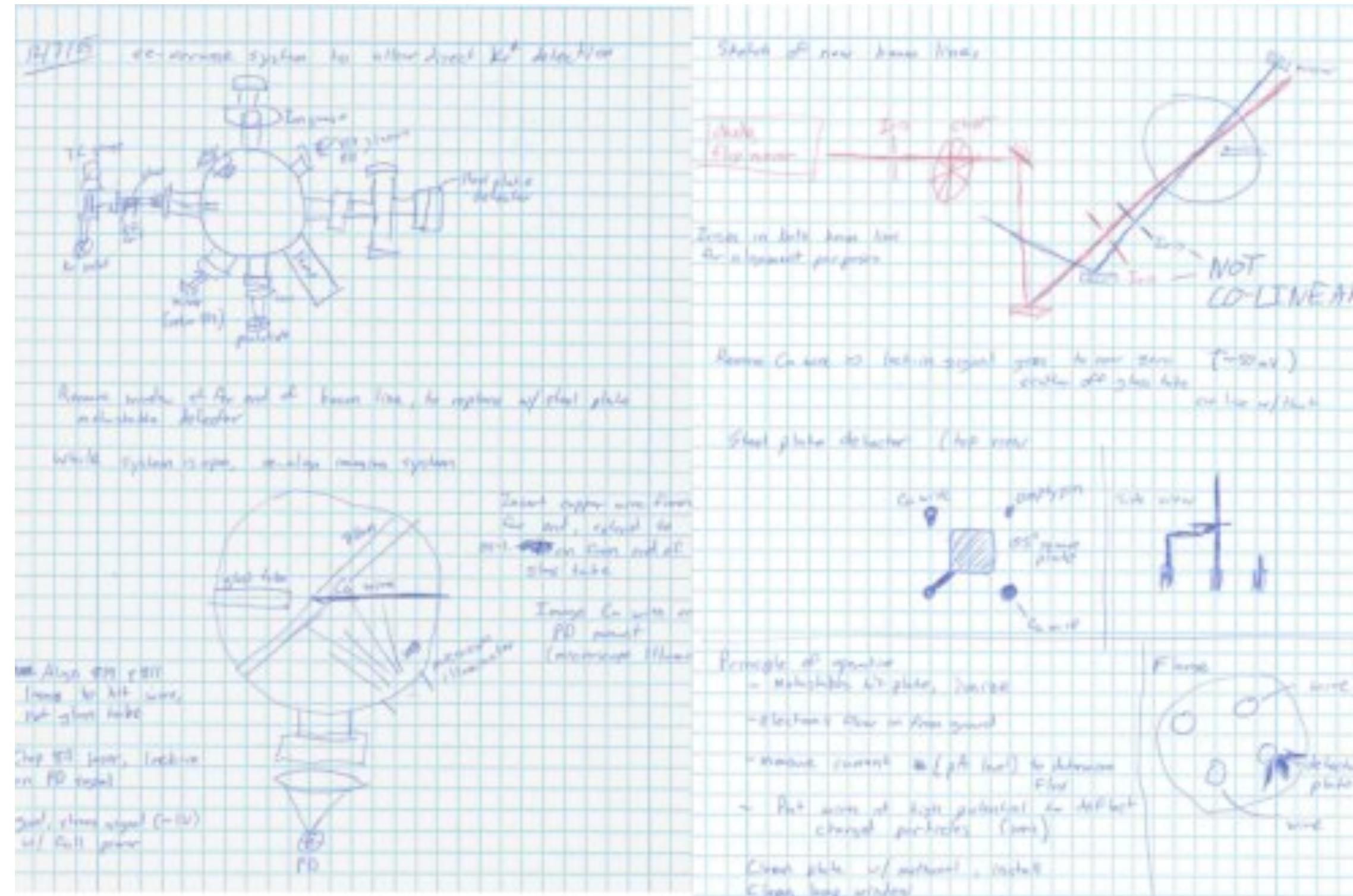
Don't be afraid



The screenshot shows a Galaxy web interface. At the top, there's a dark header bar with the Galaxy logo on the left and navigation links: Analyze Data, Workflow, Shared Data (which is highlighted in yellow), Visualization, Help, and User. Below the header, a light gray navigation bar contains the text "Published Pages | aun1 | heteroplasmy". The main content area has a white background. At the top of the content, the title of the manuscript is displayed in bold black font: "Dynamics of mitochondrial heteroplasmy in three families investigated via a repeatable re-sequencing study". Below the title, the authors' names are listed: Hiroki Goto¹, Benjamin Dickins², Enis Afgan^{3,5}, Ian M. Paul⁴, James Taylor^{3,5}, Kateryna D. Makova¹, and Anton Nekrutenko^{2,5}. A note below the authors states: "Published in Genome Biology on June 23, 2011". Another note says: "Correspondence should be addressed to [KDM](#), [JT](#), or [AN](#)". The first section, titled "1. How to use this document", contains text explaining that it's a live copy of supplementary materials for the manuscript. It provides access to all data and workflows. It mentions that users can play with them by re-running, changing parameters, or applying them to their own sequencing data. It notes that importing workflows requires a Galaxy account. It also mentions several screencasts to help users. A bulleted list follows: "access our datasets", "re-use workflows listed on this page", and "view and import histories listed on this page". Below this, another bulleted list includes: "Watch the analysis of one family (F7) from start (Illumina reads) to finish (a list of variable position)" and "Watch how the complete analysis can be performed on the Amazon Cloud". A note at the bottom of this section says: "If you experience any problems while using this page, please e-mail our [bug report list](#) and we will get back to you." The second section, titled "2. Accessing the Data", contains text about datasets. It says that datasets discussed in the paper can be found in two places: "A Galaxy Library called mtProject" and "An S3 bucket on the Amazon Cloud". A note below this explains the naming convention for datasets: "[family]-[tissue][individual]-[PCR replicate]" where family is "F4", "F7", or "F11", tissue is either "c" (cheek swab of buccal tissue) or "b" (blood), individual is an individual id, and PCR replicate is either 1 or 2. For example, F4-bM4C2-1 means PCR replicate 1 from blood of individual M4C2 from family 4.

Reproducible research

<http://main.g2.bx.psu.edu/u/aun1/p/heteroplasmy>



Reproducible research

Keep a lab book

Repeatability of published microarray gene expression analyses

John P A Ioannidis^{1–3}, David B Allison⁴, Catherine A Ball⁵, Issa Coulibaly⁴, Xiangqin Cui⁴, Aedín C Culhane^{6,7}, Mario Falchi^{8,9}, Cesare Furlanello¹⁰, Laurence Game¹¹, Giuseppe Jurman¹⁰, Jon Mangion¹¹, Tapan Mehta⁴, Michael Nitzberg⁵, Grier P Page^{4,12}, Enrico Petretto^{11,13} & Vera van Noort¹⁴

Reproducible research

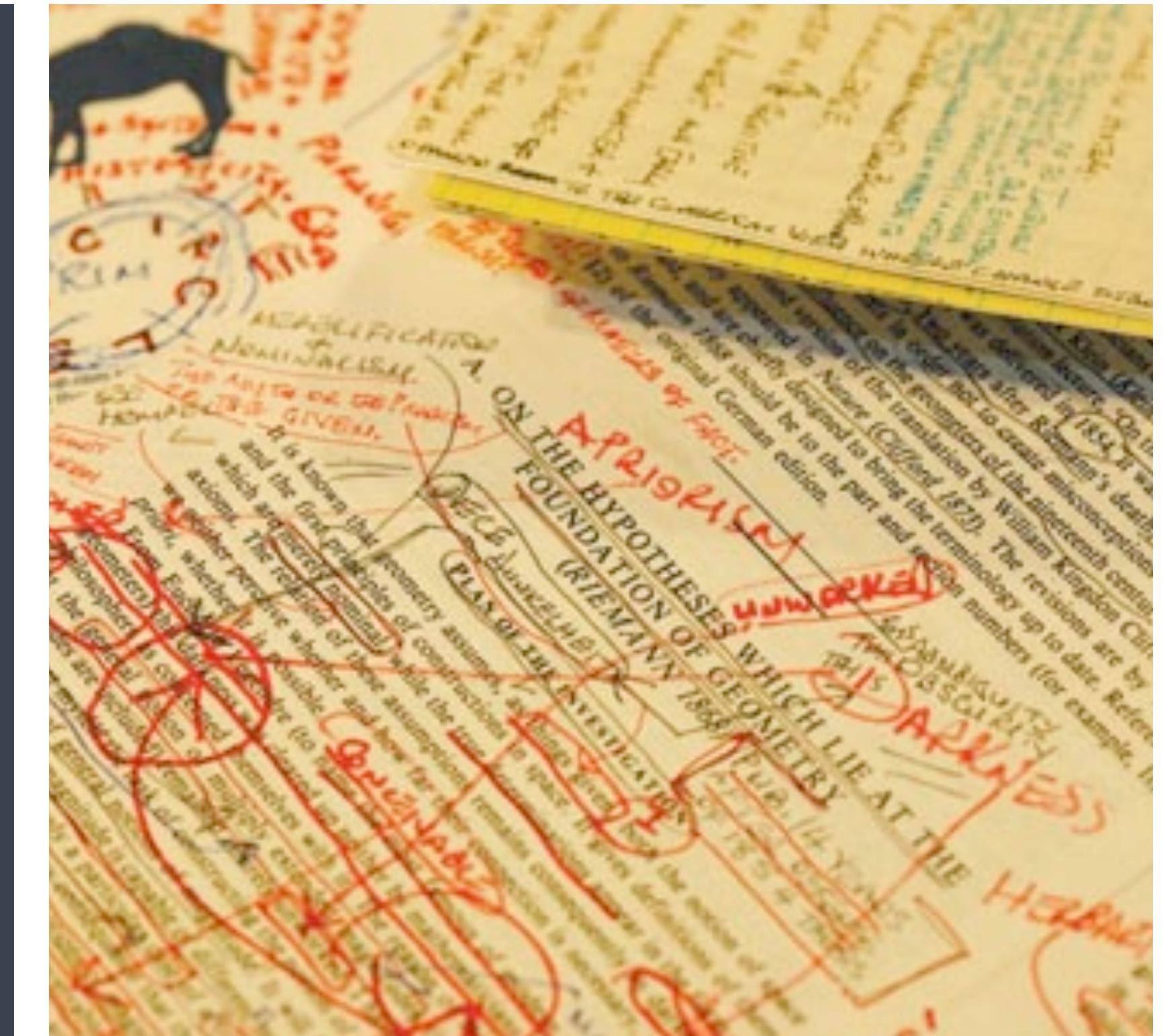
Test... 2... 3

Course in pilot phase



Feedback

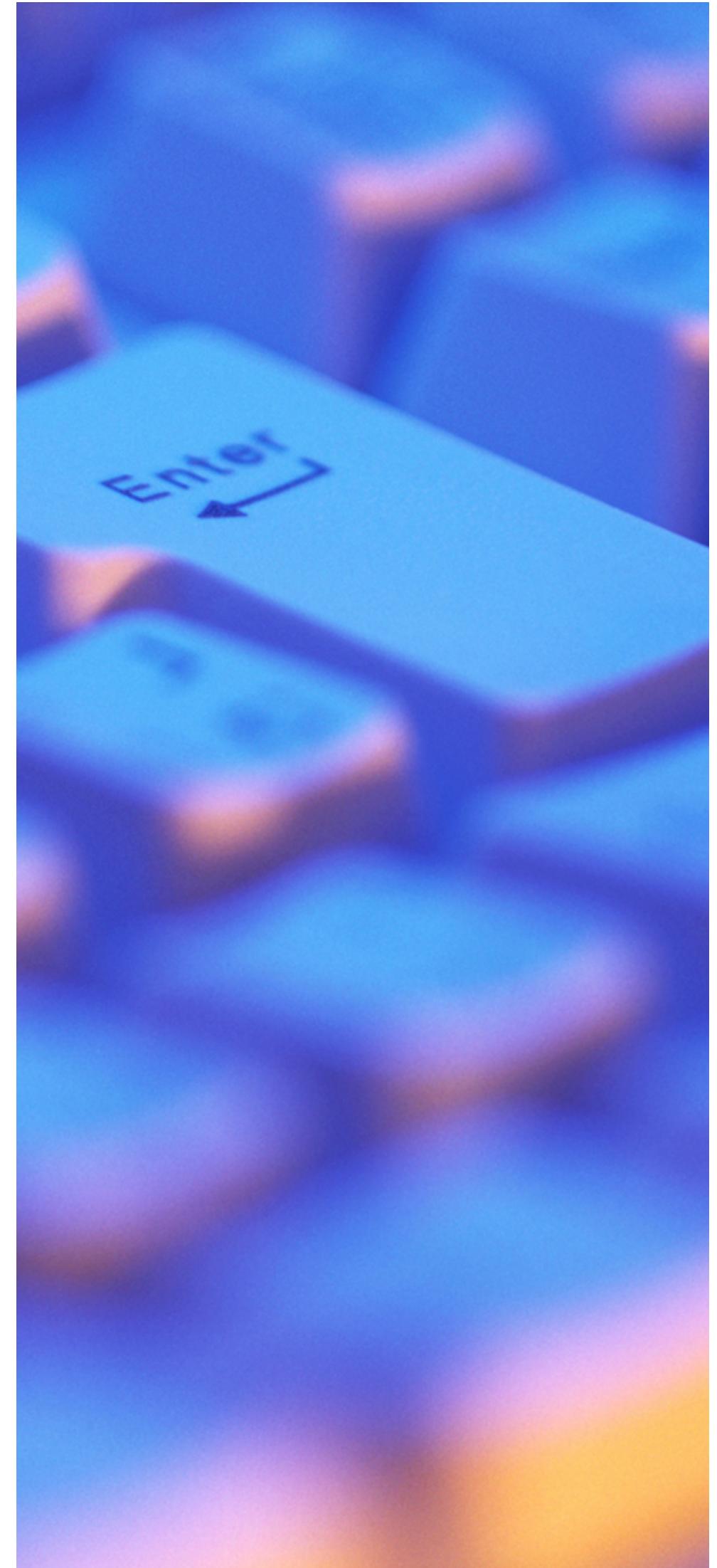
- ▶ Survey coming up



Collaborate

More information

See <http://hbc.github.io/ngs-workshops/> for additional pointers



Contact

bioinformatics@hsph.harvard.edu

<http://bioinformatics.hms.harvard.edu/>



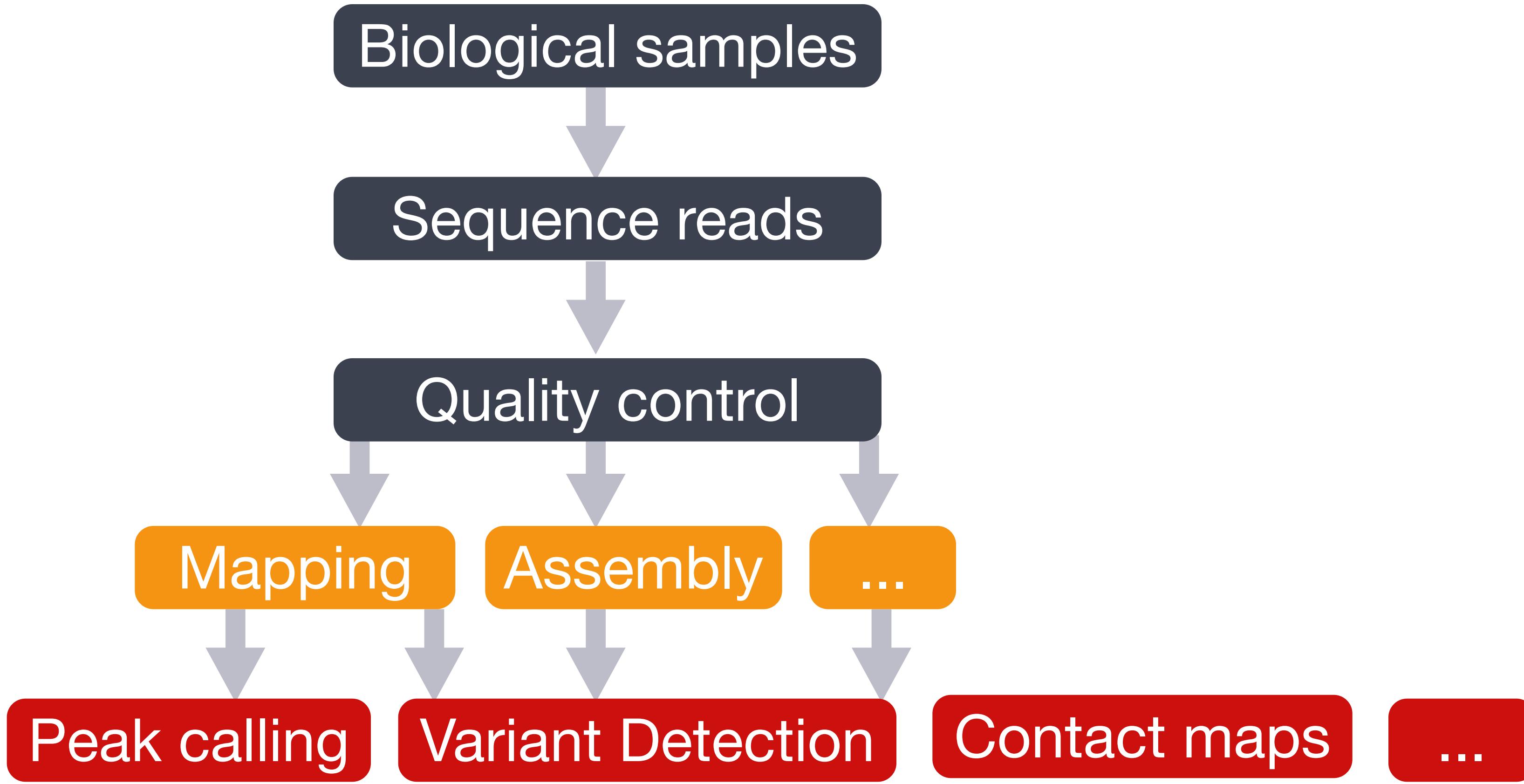
</intro>

Need for standards

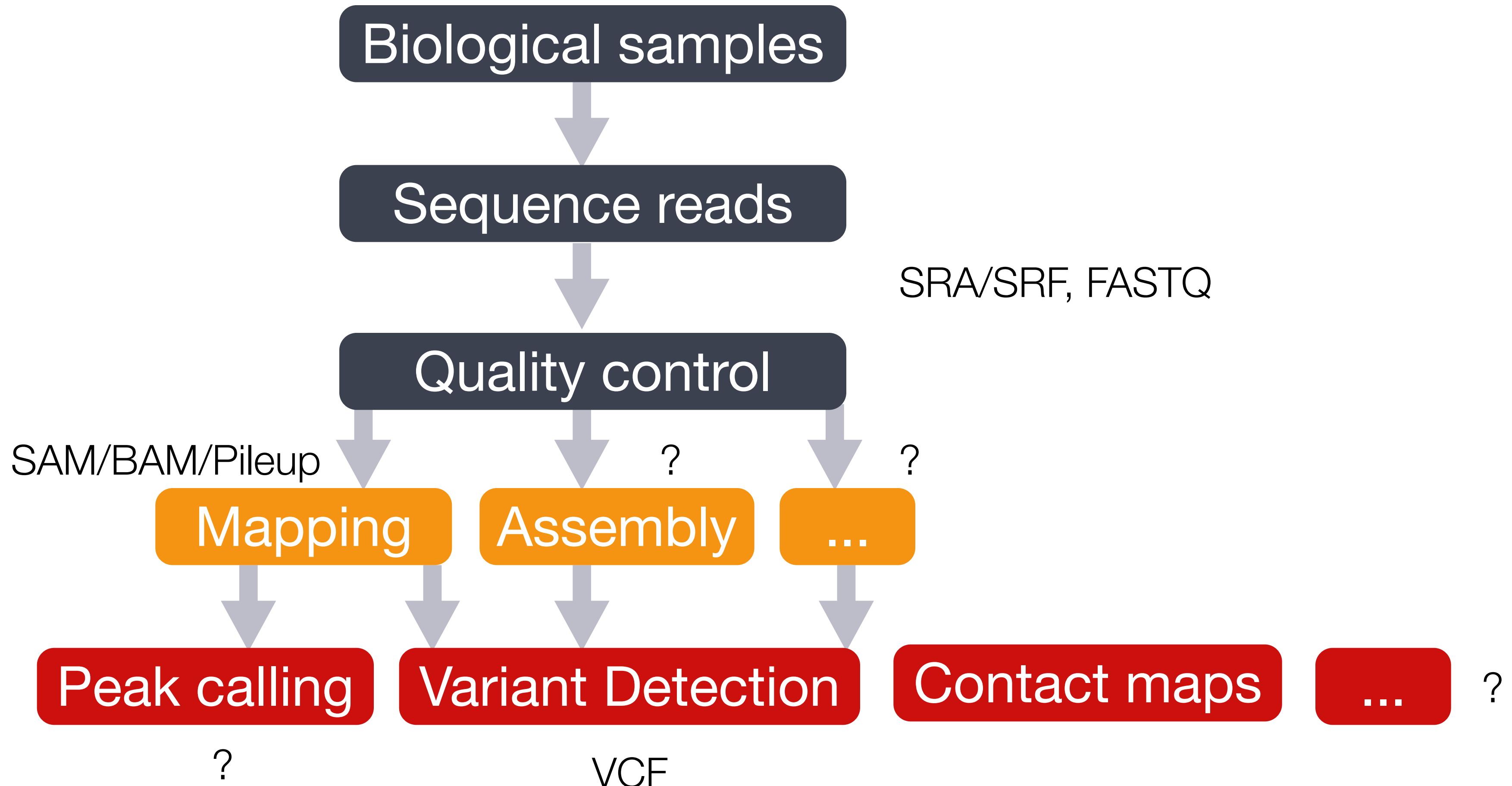
Plug and play: modular approach to tools

Vital factor in application acceptance





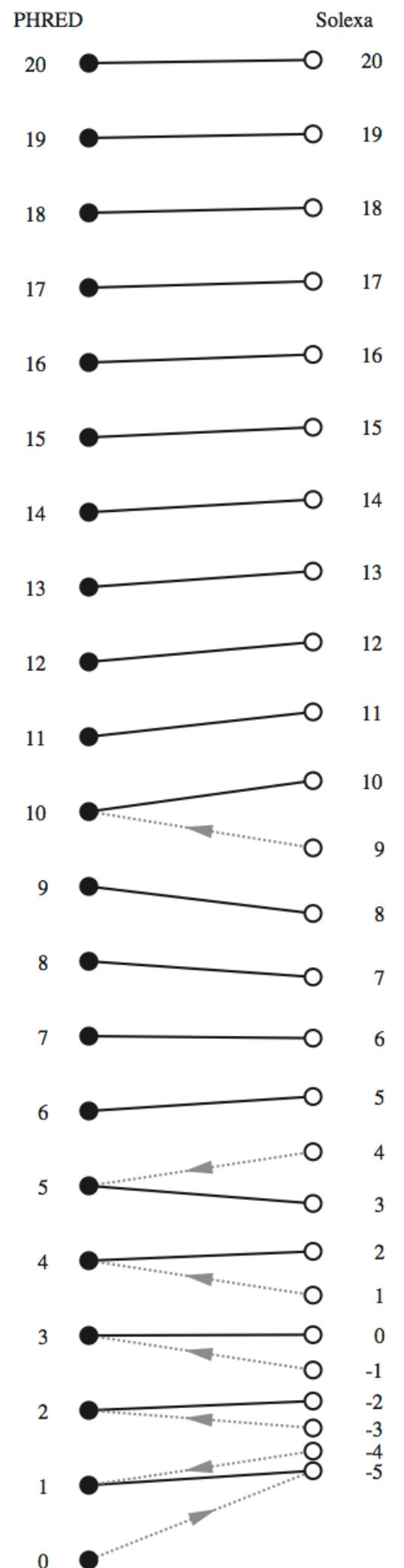
Standards



FASTQ: a “standard”

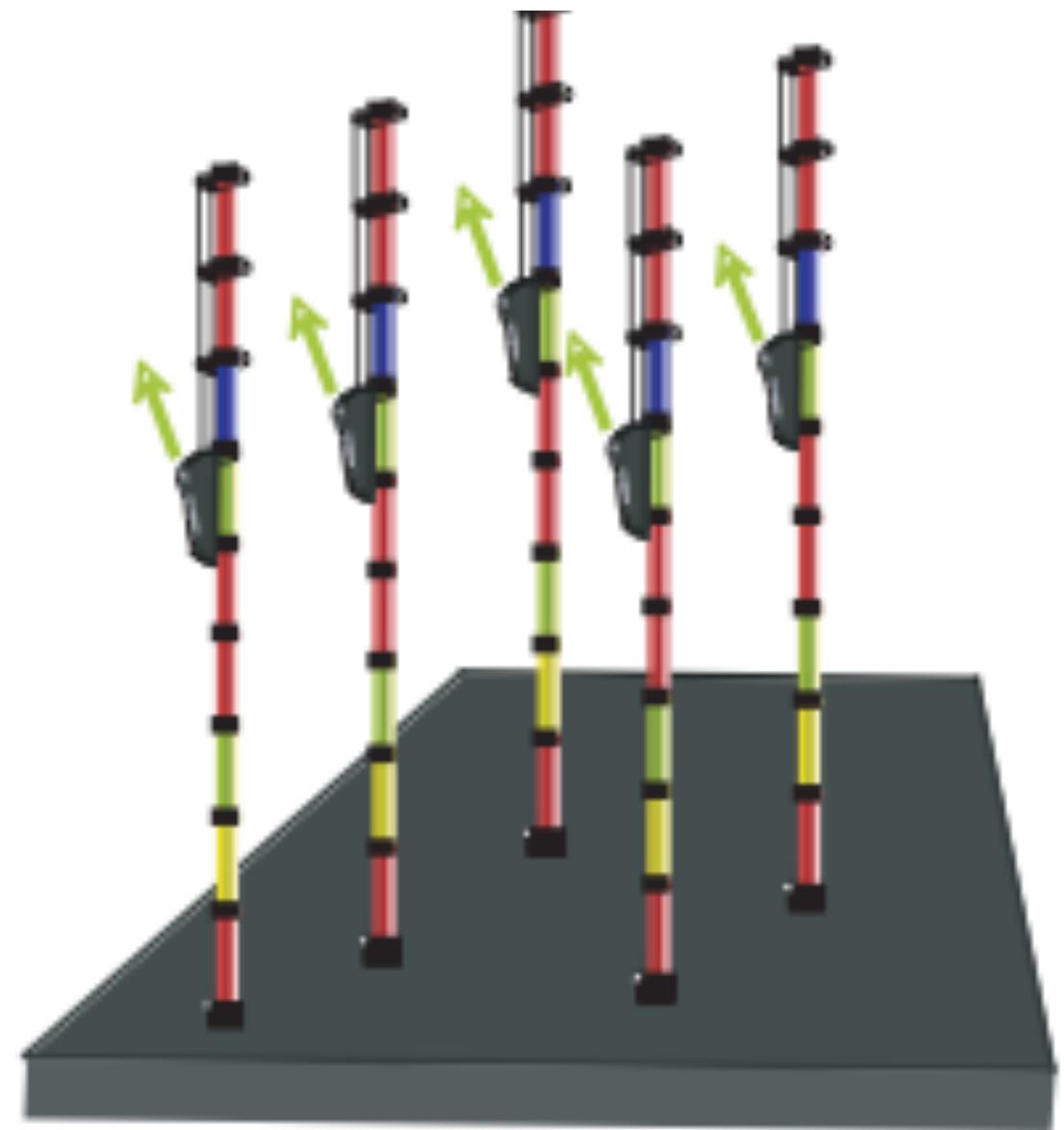
Sanger FASTQ, Solexa FASTQ, ABI Colour Space FASTQ, ...

```
@SRR014849.1 EIXKN4201CFU84 length=93
GGGGGGGGGGGGGGGGCTTTTGTGGAACCGAAAGG
GTTTGAAATTCAAACCCCTTCGGTTCCAACCTCCAA AGCAATGCCAATA
+SRR014849.1 EIXKN4201CFU84 length=93 3+&$#""""""""""7F@71,"";C?,B;?6B;:EA1EA
1EA5'9B;?:#9EA0D@2EA5':>5?:%A;A8A;?9B;D@ /=<?7=9<2A8==
```

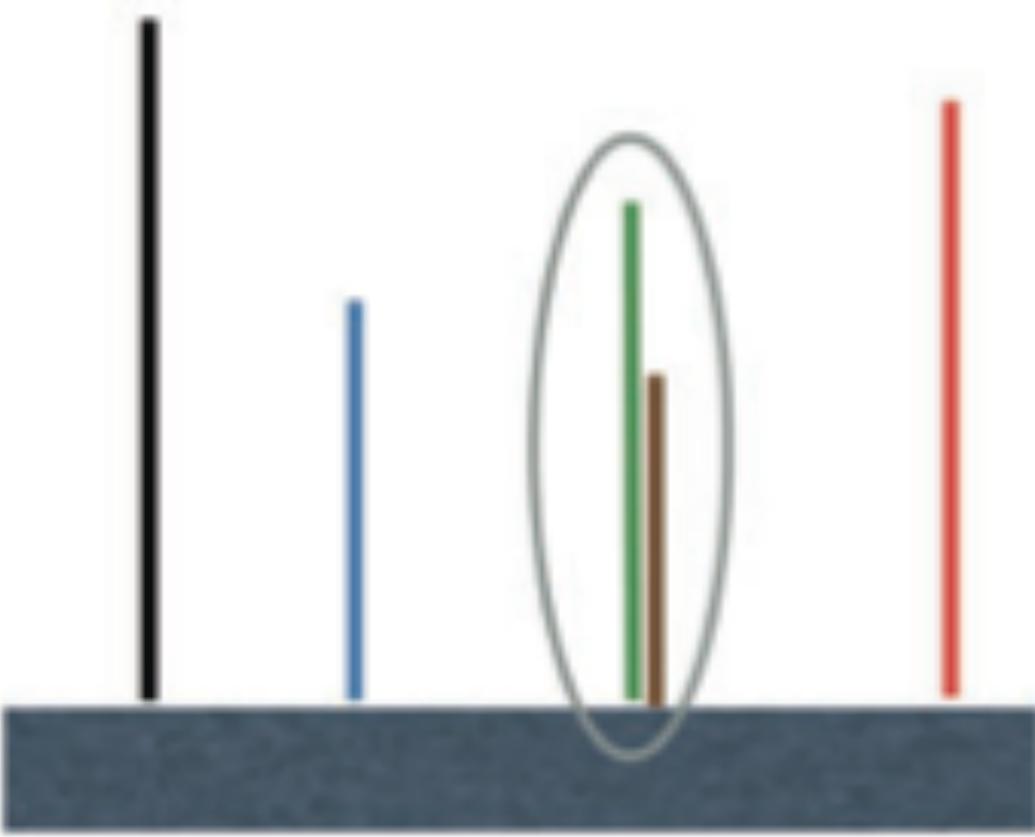


Error profiles

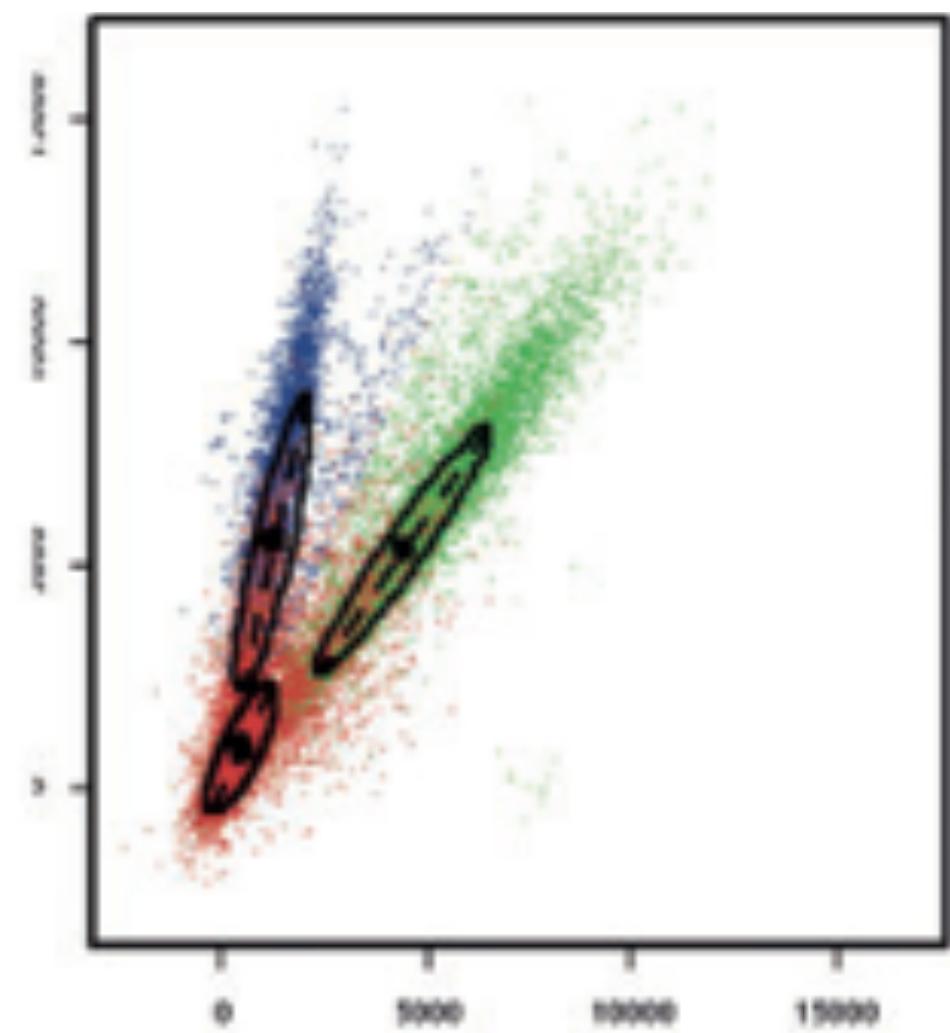
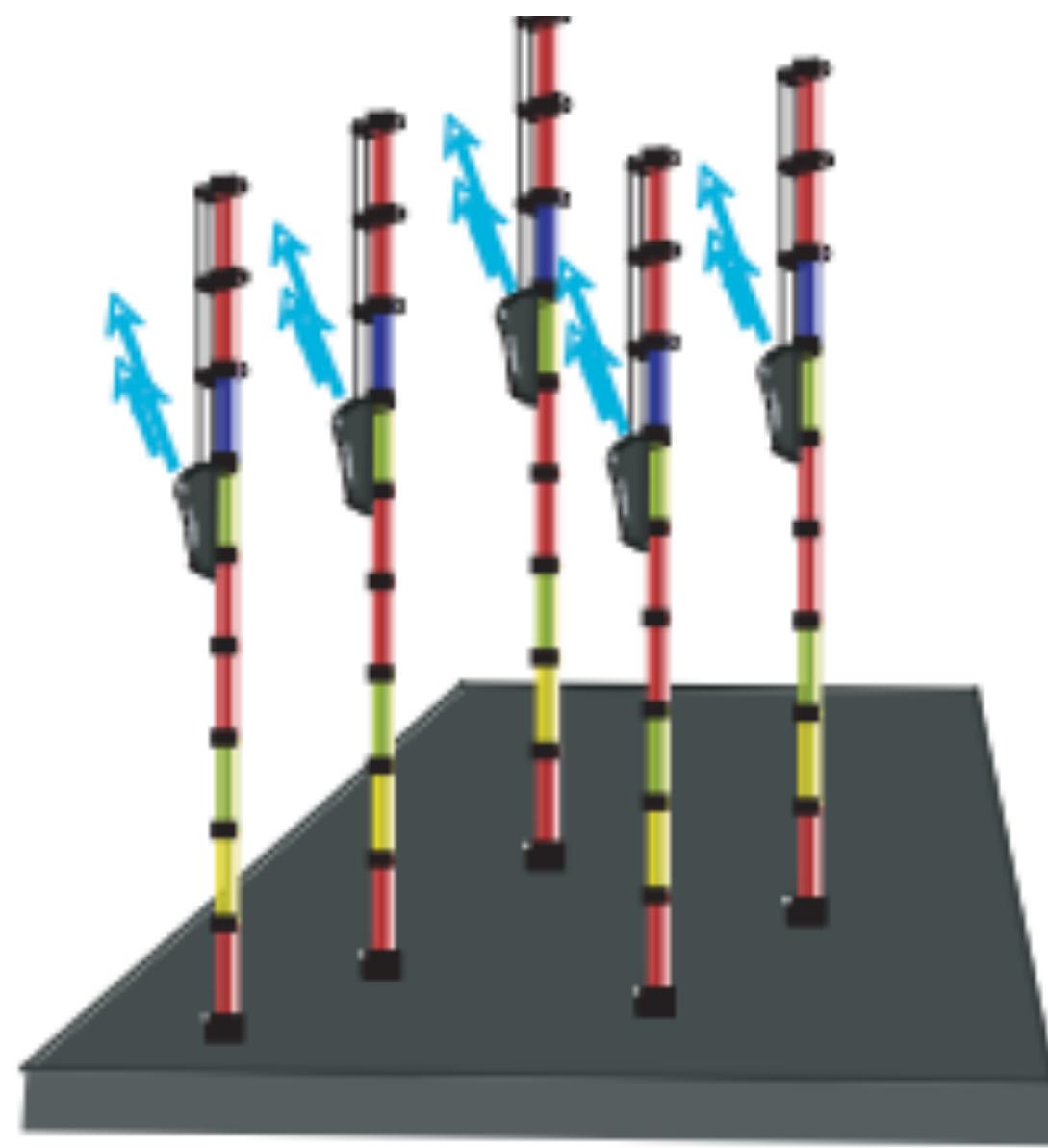
- ▶ PCR artifacts
- ▶ Error dependency on technology



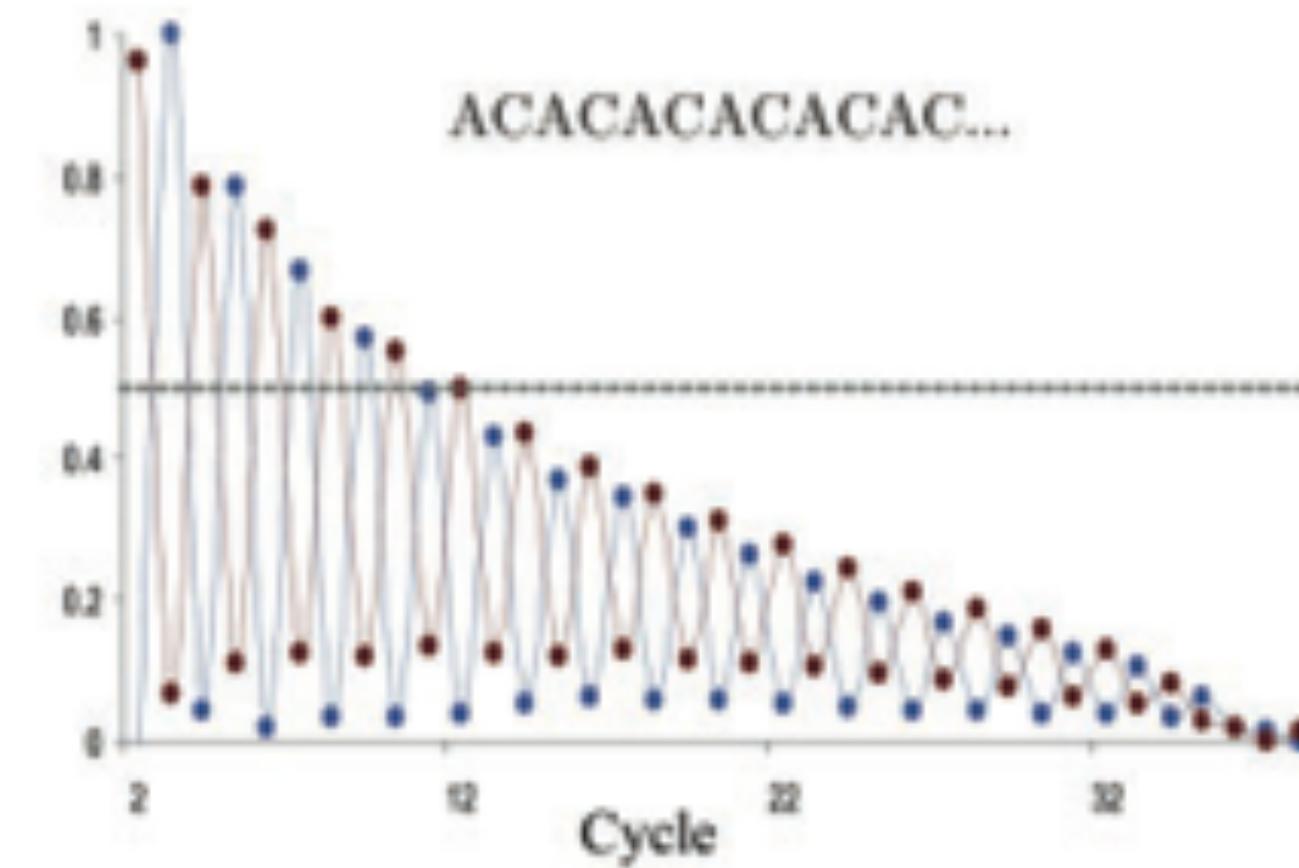
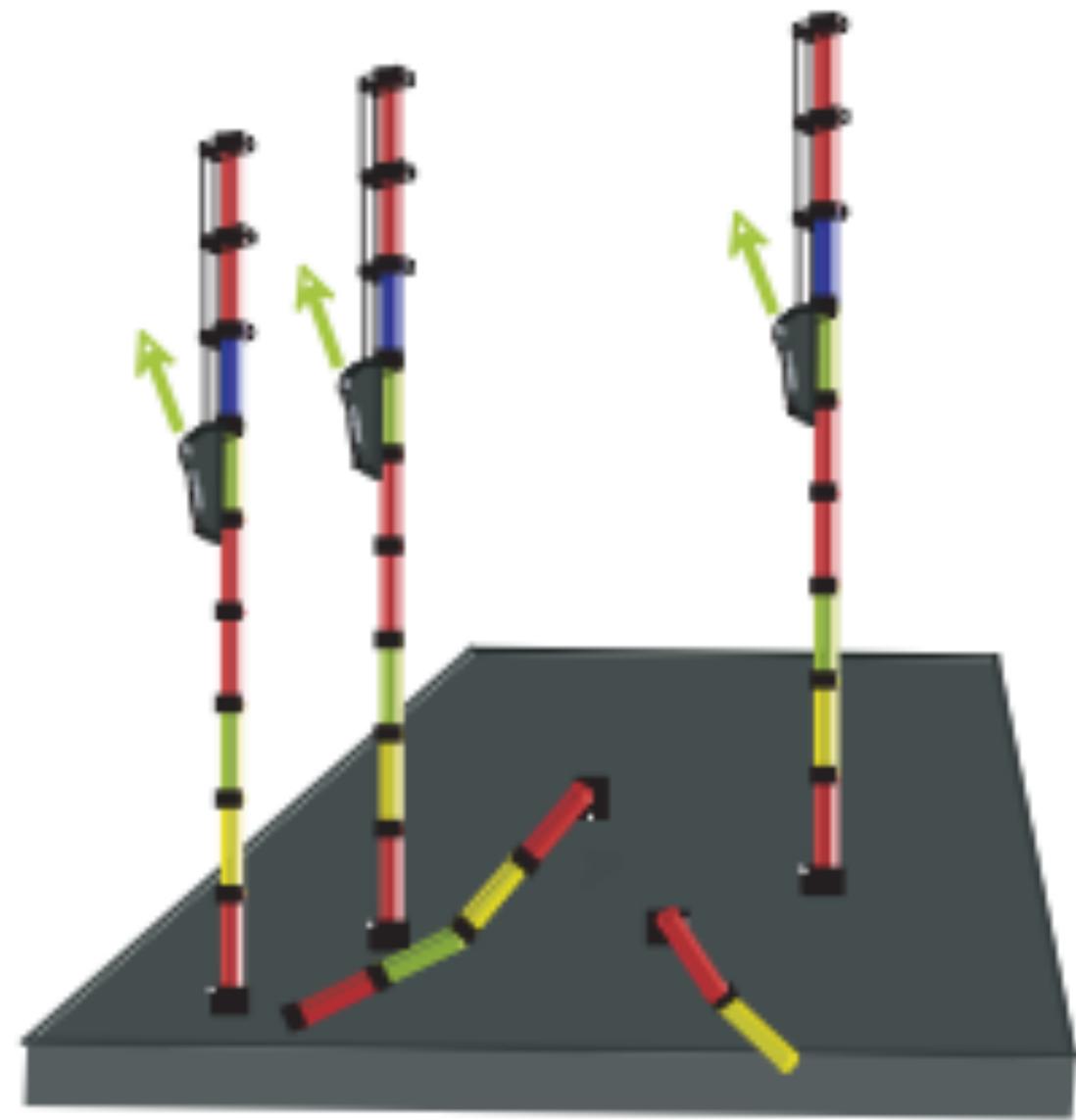
Illumina



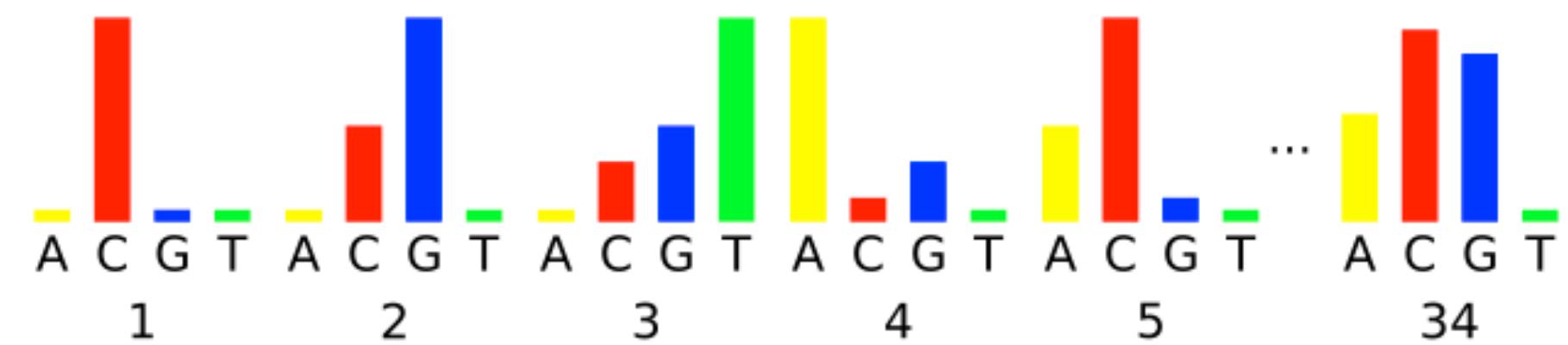
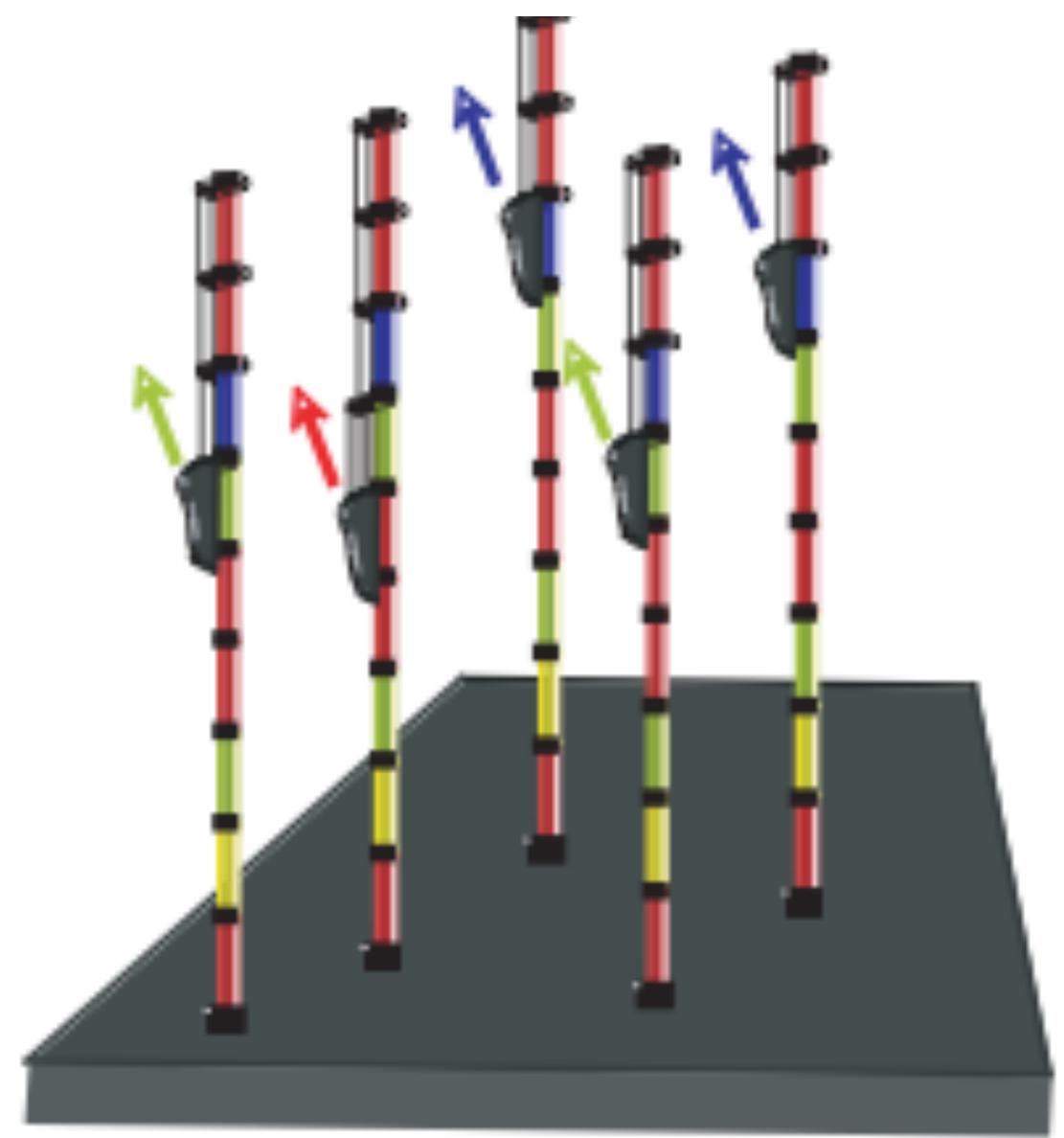
Illumina: mixed clusters



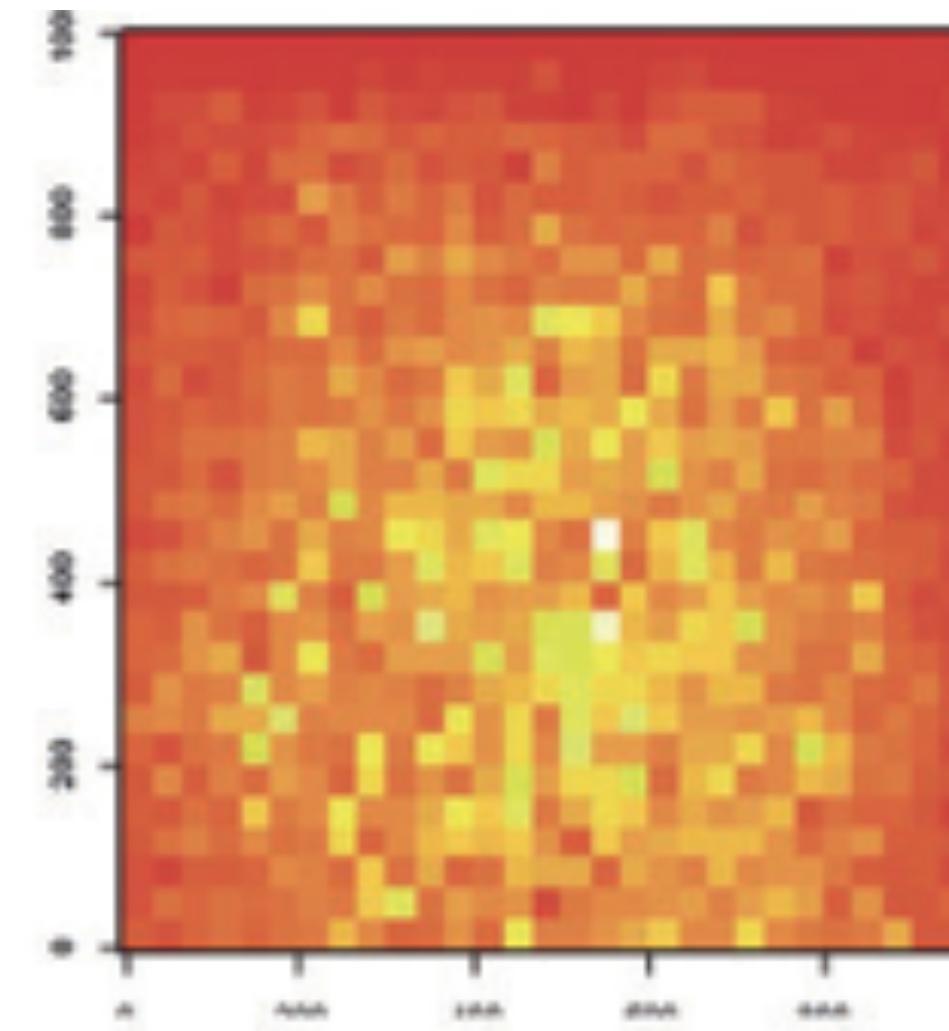
Illumina: cross-talk



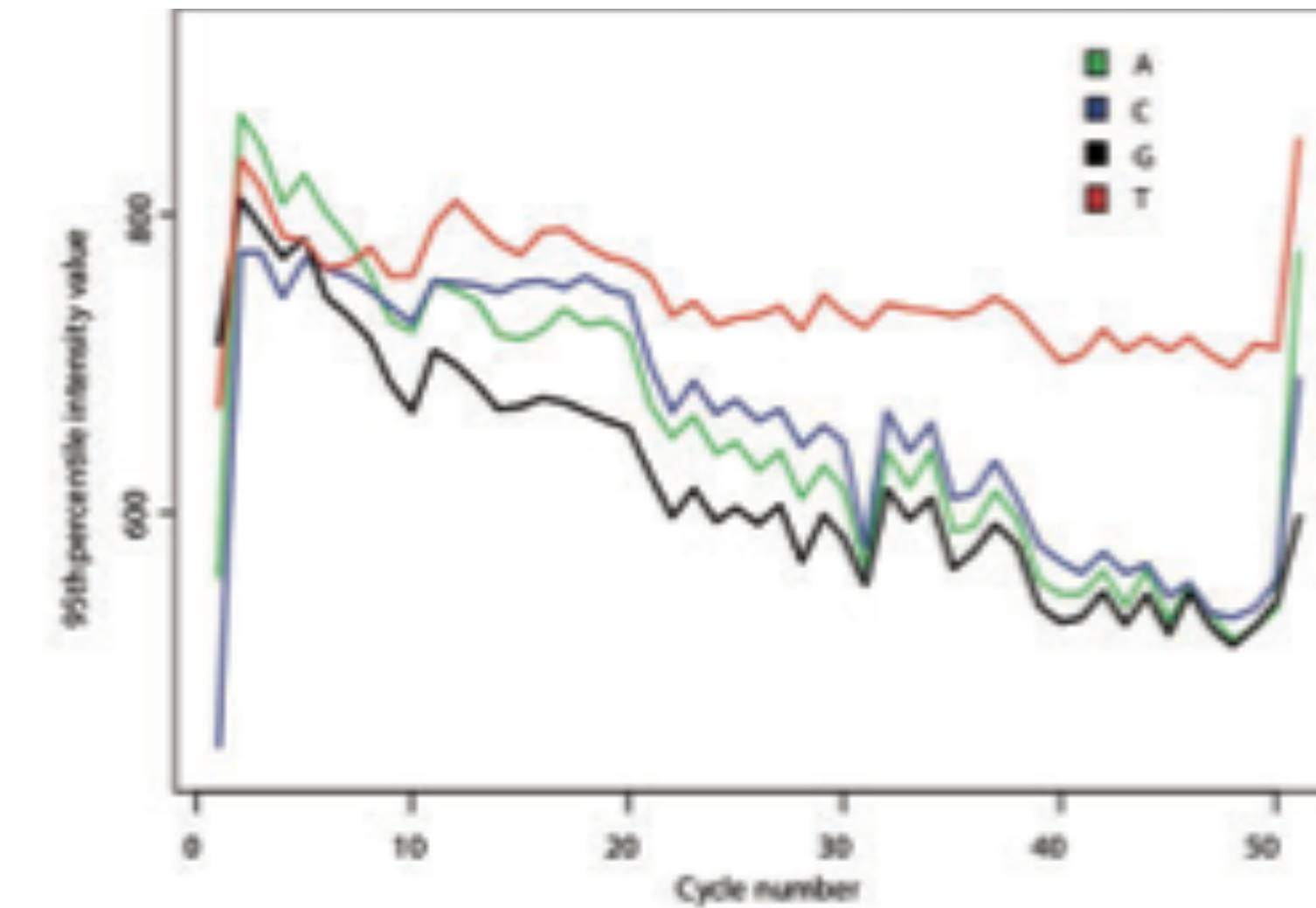
Illumina: signal decay



Illumina: phasing

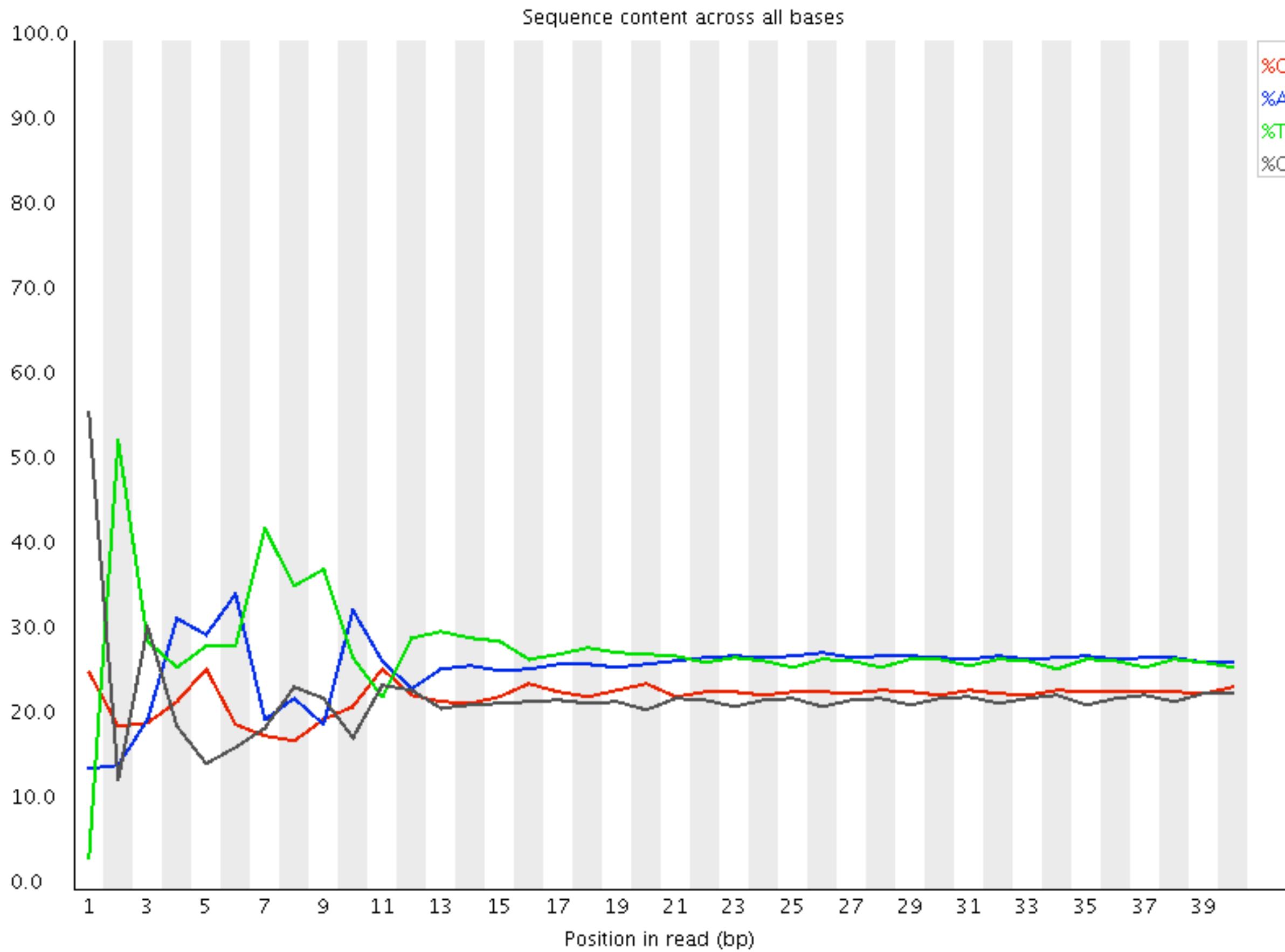


Boundary effects



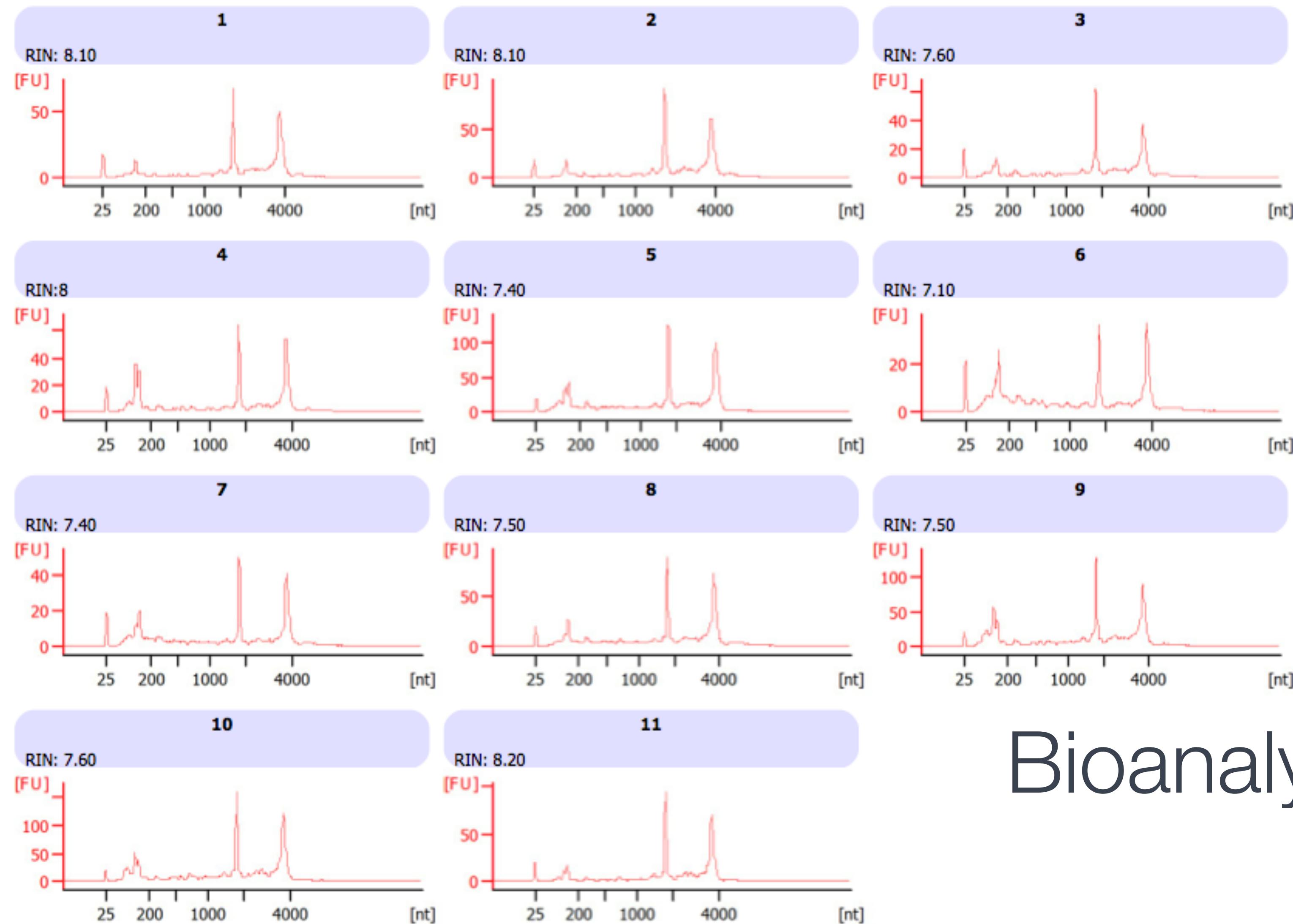
Fluorophore accumulation

Illumina: physical/chemical problems



Q&A

See http://bioinfo-core.org/index.php/9th_Discussion-28_October_2010 for more examples



Bioanalyzer

	sequence	count	lane
1051	AAAAAAAAAAAAAAAAAAAAA	70947	s_5_1_export.txt
451	AAAAAAAAAAAAAAAAAAAAA	69116	s_4_1_export.txt
601	AAAAAAAAAAAAAAAAAAAAA	66776	s_6_1_export.txt
301	AAAAAAAAAAAAAAAAAAAAA	63998	s_3_1_export.txt
751	AAAAAAAAAAAAAAAAAAAAA	55729	s_7_1_export.txt
151	AAAAAAAAAAAAAAAAAAAAA	54828	s_2_1_export.txt
901	AAAAAAAAAAAAAAAAAAAAA	40359	s_8_1_export.txt
1	NNNNNNNNNNNNNNNNNNN	30880	s_1_1_export.txt
152	NNNNNNNNNNNNNNNNNNN	30485	s_2_1_export.txt
153	CNNNNNNNNNNNNNNNNNNN	26476	s_2_1_export.txt
2	TNNNNNNNNNNNNNNNNNNN	25600	s_1_1_export.txt
154	GNNNNNNNNNNNNNNNNNNN	25594	s_2_1_export.txt
3	CNNNNNNNNNNNNNNNNNNN	25063	s_1_1_export.txt
155	TNNNNNNNNNNNNNNNNNNN	24965	s_2_1_export.txt
4	GNNNNNNNNNNNNNNNNNNN	24164	s_1_1_export.txt
302	NNNNNNNNNNNNNNNNNNN	22501	s_3_1_export.txt
5	AAAAAAAAAAAAAAAAAAAAA	20996	s_1_1_export.txt
452	TNNNNNNNNNNNNNNNNNNN	20842	s_4_1_export.txt

QA: filtering

	sequence	count
1	ATTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	482185
151	ATTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	271724
2	TAATACGACTCACTATAGGGCGAATTGAATTAGCGGCCGCGAATTGCC	159936
152	TAATACGACTCACTATAGGGCGAATTGAATTAGCGGCCGCGAATTGCC	105273
153	CTTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	46872
3	CTTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	43212
4	NN	13142

Read Frequency Distribution

QA: filtering

> gnl|uv|NGB00105.1:1-219 pCR4-TOPO multiple cloning site
Length=219

Score = 100 bits (50), Expect = 9e-19
Identities = 50/50 (100%), Gaps = 0/50 (0%)
Strand=Plus/Plus

Query	1	ATTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	50
Sbjct	43	ATTAACCCTCACTAAAGGGACTAGTCCTGCAGGTTAACGAATTGCC	92

VecBase Screen

QA: filtering

chrX:52139280 152139290 152139300 152139310 152139320 152139330
--->CGCCGTCCCTCAGAATGGAAACCTCGCTTCTCTGCCCAATGCGCAAGTCAG

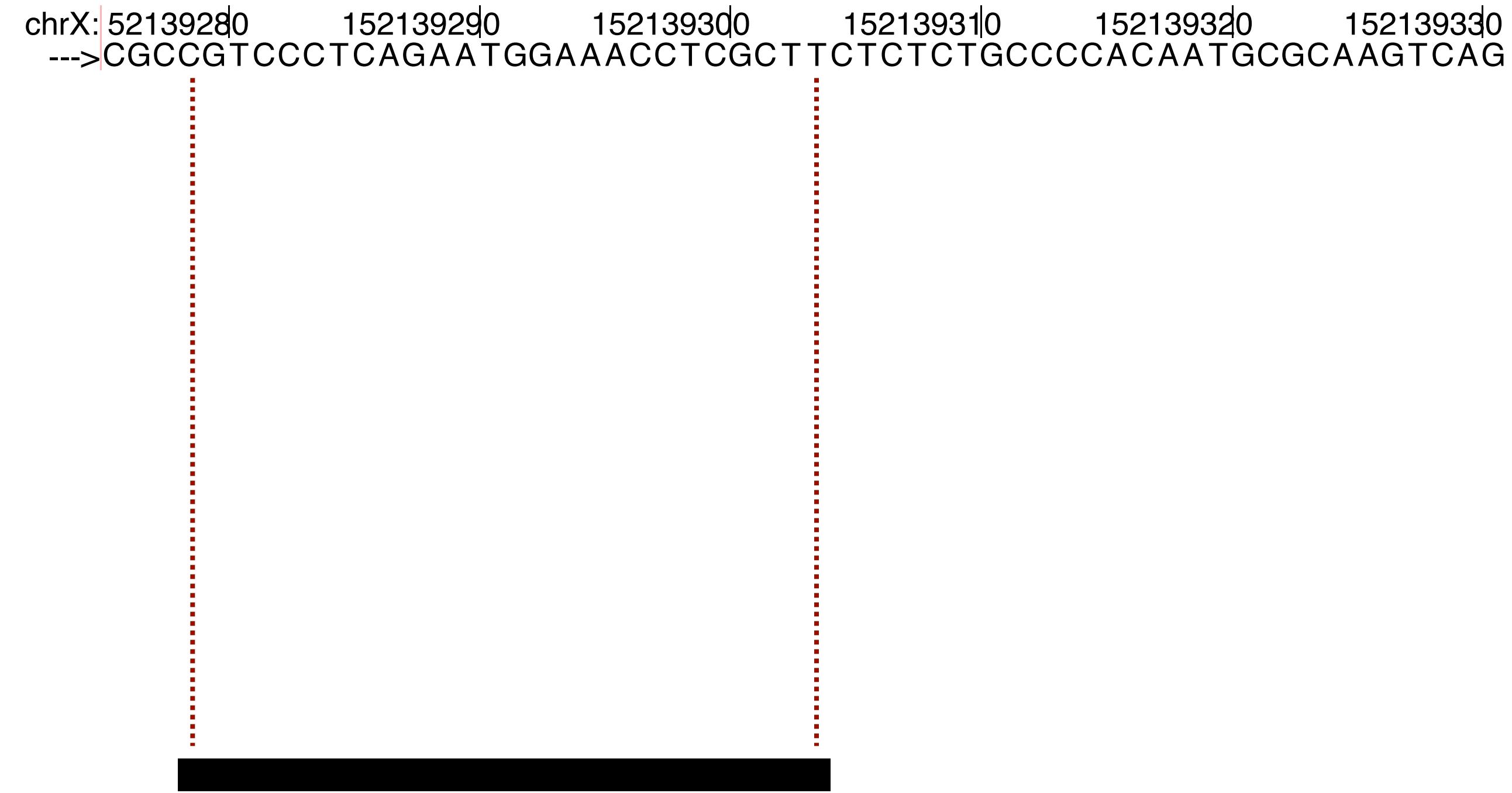
Genome

A horizontal sequence of DNA bases: CGTCCCTCAGAATGGAAACCTCGCTT. The last four bases, TCGCTT, are highlighted with a red oval.

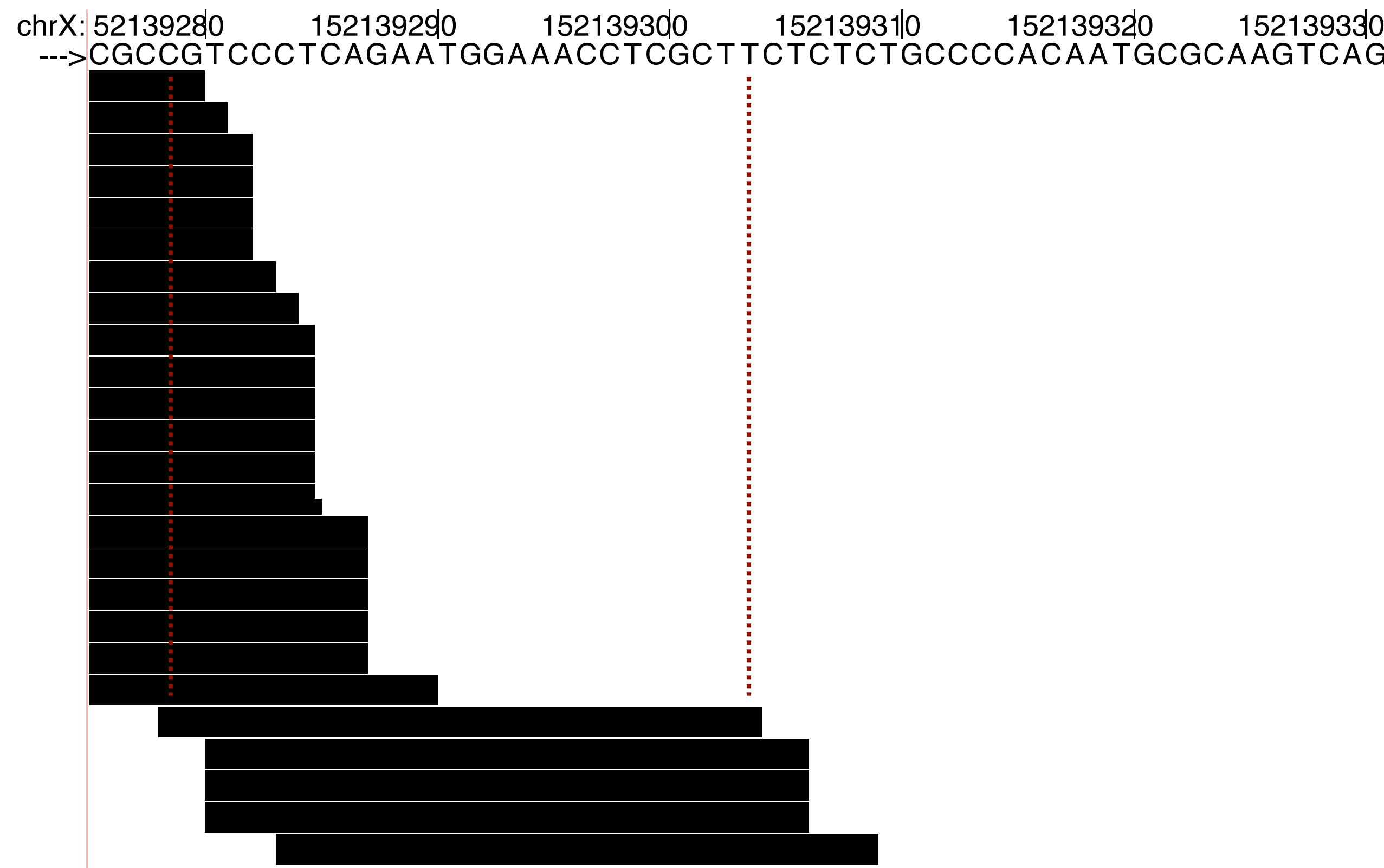
CGTCCCTCAGAATGGAAACCTCGCTT

Sequence tag

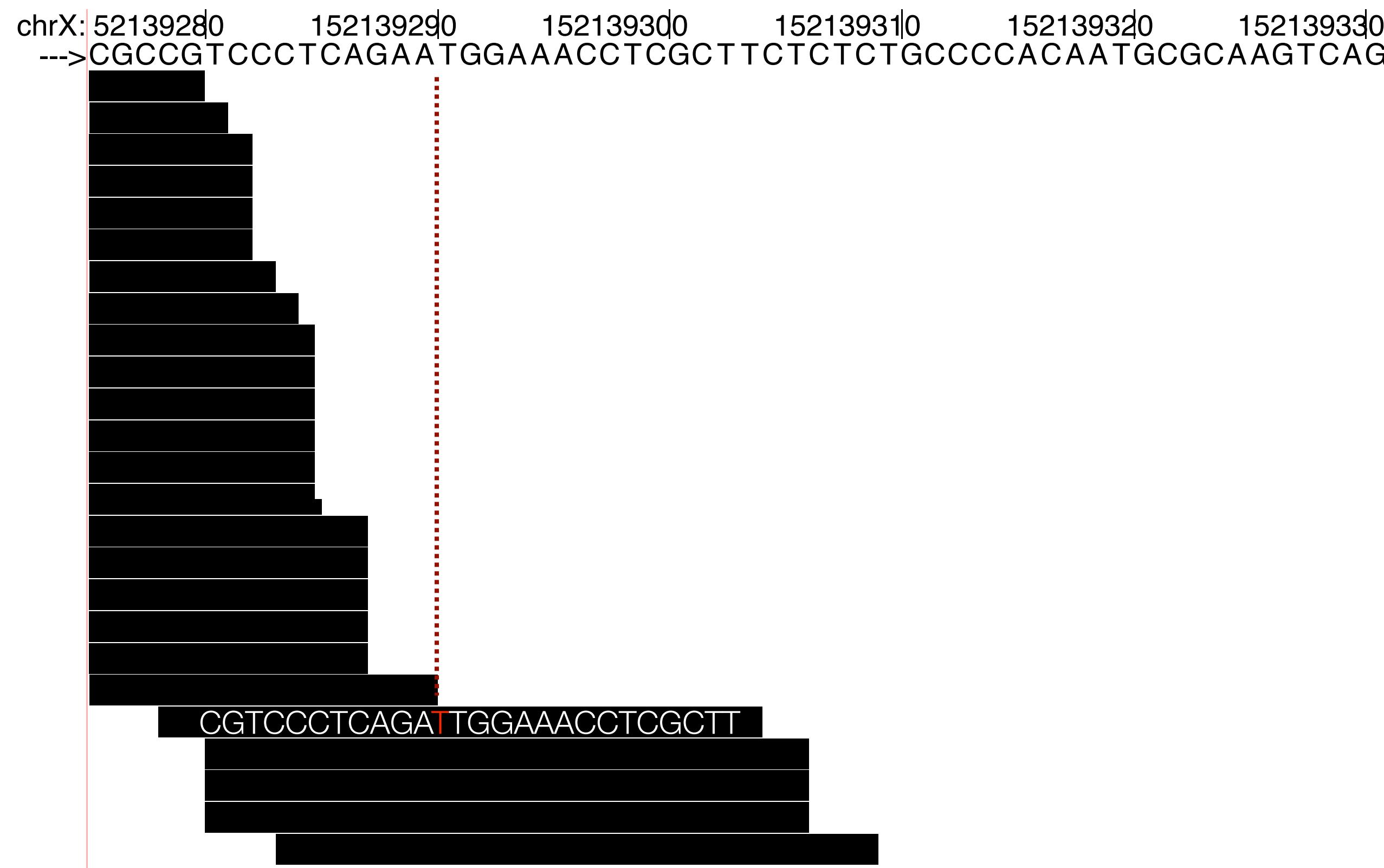
Tool evolution: mapping reads to a genome



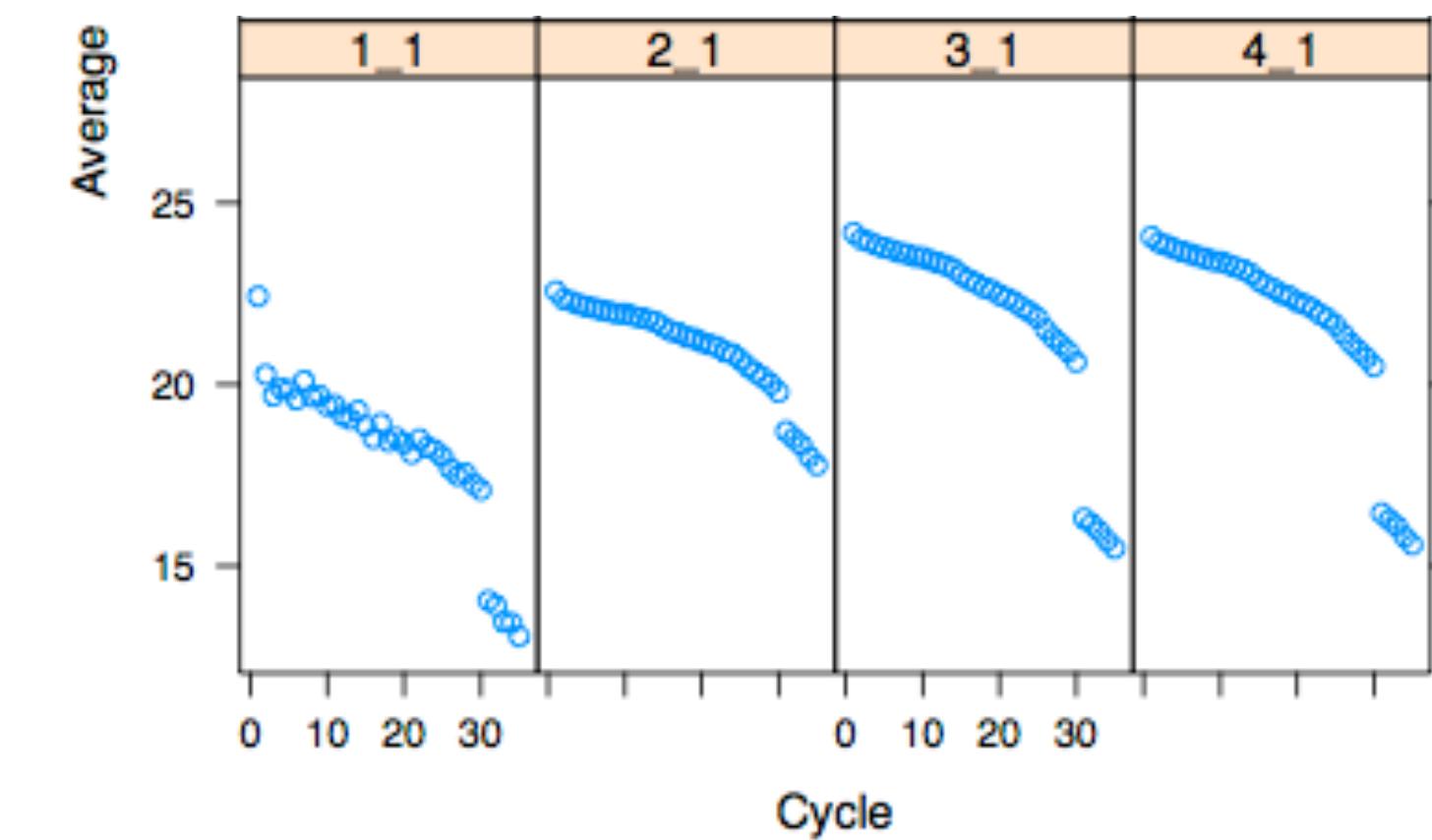
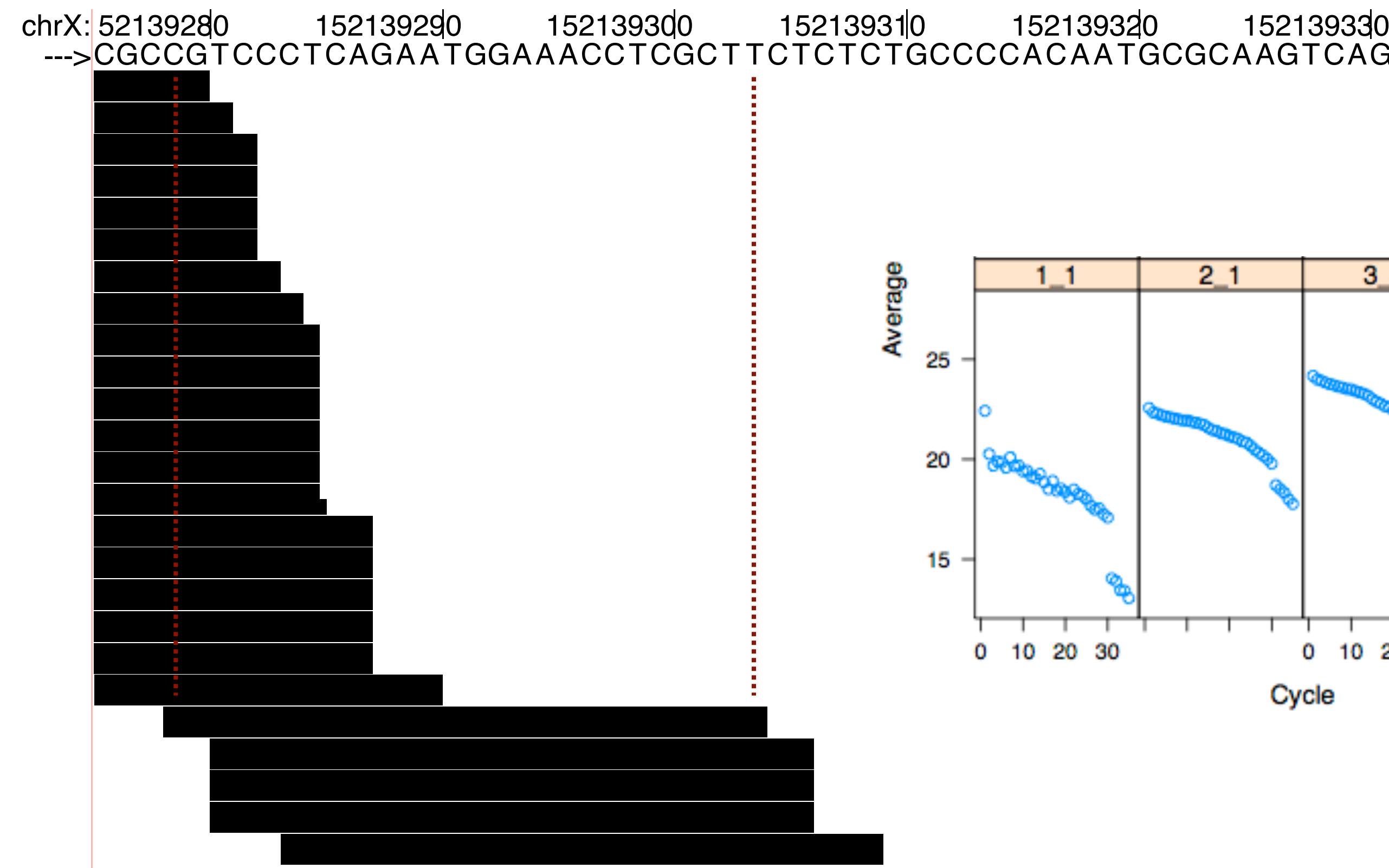
Mapping to a Reference Genome



Mapping to a Reference Genome



Mapping to a Reference Genome



Mapping to a Reference Genome

Tool evolution: mapping approaches

- ▶ Variation in algorithm
- ▶ Alignment speed
- ▶ Memory requirements
- ▶ Error tolerance
- ▶ ...

Table 1 A selection of short-read analysis software

Program	Website	Open source?	Handles ABI color space?	Maximum read length
Bowtie	http://bowtie.cbcn.umd.edu	Yes	No	None
BWA	http://maq.sourceforge.net/bwa-man.shtml	Yes	Yes	None
Maq	http://maq.sourceforge.net	Yes	Yes	127
Mosaik	http://bioinformatics.bc.edu/marbl/Mosaik	No	Yes	None
Novoalign	http://www.novocraft.com	No	No	None
SOAP2	http://soap.genomics.org.cn	No	No	60
ZOOM	http://www.bioinfor.com	No	Yes	240

The standard tools

Mapper	Data	Availability	Version	O.S.	Number Citations	Citations/Years	Seq.Plat.	Input	Output	Min. RL	Max. RL	Mismatches	Indels	Gaps	Align. Reported	Alignment	Parallel	QA PE	Splicing	
Bowtie	DNA	OS	0.12.7	Linux,Mac,Windows	1168	335,04	I,So,4,Sa,P	(C)FAST(AQ)	SAM TSV	4	1K	Score	Score	N	A,B,R,S	G L	SM	Y Y	N	
Blat	DNA	OS	34	Linux,Mac	2844	268,37	N	FASTA	TSV BLAST	11	5000K	Score	Score	Y	B	L	N	N N	De novo	
MAQ	DNA	OS	0.7.1	Linux,Mac	957	237,27	I,So	(C)FAST(AQ)	TSV	8	63	Y	Y	N			N	Y Y	N	
BWA	DNA	OS	0.6.2	Linux,Mac,Windows	738	225,15	I,So,4,Sa,P	FASTA/Q	SAM	4	200	Y	8	Y	R,S	G	SM	Y Y	N	
TopHat	RNA	OS	1.4.1	Linux,Mac	389	112,66	I	FASTA/Q, GFF	BAM	-	-	2	0	N	B,S	-	SM	Y Y	De novo	
SOAP	DNA	OS	1.11	Linux,Mac	451	98,04	I	FASTA/Q	TSV	7	60	5	3	N	B,R,S		SM	N Y	N	
SOAP2	DNA	OS	2.21	Linux	294	90,93	I	FASTA/Q	SAM TSV	27	1K	2	0	Y	A,B,R	L	SM	N Y	N	
Mummer 3	DNA	OS	3.23	Linux,Mac	683	78,93	N	FASTA	TSV	10	*	Y	Y	Y	A,B	G	N	N N	N	
BWA-SW	DNA	OS	0.6.2	Linux,Mac,Windows	160	61,41	I,So,4,Sa,P	FASTA/Q	SAM	4	1000K	0.1	0.1	Y	R,S	L	SM	Y N	N	
mrFAST	miRNA	OS	2.1.0.4	Linux	158	52,86	I	FASTA/Q	SAM	25	300	Score	6	N	A,B	G	N	N Y	N	
SHRIMP	DNA	OS	1.3.2	Linux,Mac	155	47,45	I,So,4,Hel	(C)FAST(AQ)	TSV	14	1K	Score	Score	Y	B,S	G	SM	N Y	N	
SSAHA	DNA	OS	3.1	Linux,Mac	483	43,94	N	FASTA/Q	TSV	15	*	Y	Y	Y	B,S	G L	N	N N	N	
CloudBurst	DNA	OS	1.1	Linux,Mac,Windows	146	43,08	N	FASTA	TSV		1K	Y	Y	Y	A,B	G	Cloud	N N	N	
RMAP	DNA	OS	2.05	Linux,Mac	162	35,89	I,So,4	(C)FAST(AQ)	BED	11	10K	Y	0	N	B,S		N	Y Y	N	
SeqMap	DNA	OS	1.0131	Linux,Mac	142	35,04	I	FASTA	ELAND	15	500	5	3	N	A		SM	N N	N	
BFAST	DNA	OS	0.7.0	Linux,Mac	94	33,74	I,So,4, Hel	(C)FAST(AQ)	SAM TSV		*	Y	Y	Y	B,R,U	G	SM	N Y	N	
Exonerate	DNA	OS	2.2	Linux,Mac	255	33,59	N	FASTA	TSV	20	*	Score	Score	Y	B,S	G L	N	N N	De novo	
GMAP	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	217	28,68	I,4,Sa,Hel,Ion P	FASTA/Q	SAM, GFF	8	*	Y	Y	Y	B	G L	SM	N N	De novo	
GSNAP	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	72	28,42	I,4,Sa,Hel,Ion P	FASTA/Q	SAM	8	250	Y	Y	Y	A,B,U,S	G L	SM	N Y	Lib and de novo	
ZOOM	DNA	Com	1.5	Linux,Mac,Windows	109	26,78		(Q)FAST(AQ)	SAM BED GFF	12	240	Y	Y	N	B,U,S	G	SM/DM	Y Y	N	
SpliceMap	RNA	OS	3.3.52	Linux,Mac	63	26,43		FASTA/Q	SAM BED	-	-	0.1		Y	A	-	SM	N Y	Lib and/or de novo	
MapSplice	RNA	OS	1.15.2	Linux	50	25,25		FASTA/Q	SAM BED	-	-	3		Y	B	-	SM	N Y	De novo	
QPALMA	RNA	OS	0.9.2	Linux,Mac	75	19,19		Bismark	FASTA	-	-	Y	Y	Y	B	L	N	Y N	Lib and de novo	
RazerS	DNA	OS	1.1	Linux,Mac,Windows	58	18,18		Stampy	TSV	-	-	Y	Y	Y	B		N	Y N	Lib and de novo	
mrsFAST	miRNA	OS	2.3.0	Linux	32	18,18		MapSplice	ELAND	11	*	Score	Score	Y	A,B,S	G	N	N Y	N	
Stampy	DNA	Bin	1.0.16	Linux,Mac	26	14,14		REAL	FASTA/Q	SAM TSV	4	4K	0.15	30	N	B,R,S	G	N	Y Y	N
PASS	DNA	Bin	1.62	Linux,Mac,Windows	45	14,14		BS Seeker	(Q)FAST(AQ)	SAM GFF3 BLAST	23	1K	Y	Y	Y	A,B	G	SM	Y Y	De novo
SOCS	DNA	OS	2.1.1	Linux,Mac,Windows	49	14,14		Supersplat	(Q)FAST(AQ)	TSV	64	Y	0	N	A,B		SM	Y N	N	
GenomeMapper	DNA	OS	0.4.3	Linux,Mac	31	14,14		SpliceMap	FASTA/Q	BED TSV	12	2K	10	10	Y	A,B,R	G	SM	N N	N
Slider	DNA	OS	0.6	Linux,Mac,Windows	39	14,14		BRAT	FASTA/Q	TSV	62	3	0	N	B,S		N	Y Y	N	
BSMAP	Bisulfite	OS	2.43	Linux,Mac	31	14,14		BFAST	FASTA/Q	SAM TSV	8	144	15	0	N	B,U,S		SM	N Y	N
PerM	DNA	OS	0.4.0	Linux,Unix,Mac,Windows	30	9,67		GNUUMAP	(Q)FAST(AQ)	SAM TSV	20	128	9	0	Y	A,U	G	DM	Y Y	N
BWT-SW	DNA	OS	20070916	Linux	15	8,64		GenomeMapper	FASTA	TSV		1K	Score	Score	Y	A		N N N	N	
SHRIMP 2	DNA	OS	2.2.2	Linux, Unix, Mac	15	8,64		mrFAST	FASTA/Q	SAM	30	1K	Y	Score	N	B,U,S	G	SM	Y Y	N
RNA-Mate	RNA	OS	1.1	Linux,Mac	28	8,41		PerM	FASTA	BED Counts	-	-	Y	0	N	S	-	DM	Y N	Lib
Supersplat	RNA	OS	1.0	Linux,Mac	28	8,41		X-Mate	FASTA	TSV		0	0	Y	A,U	G	N	N N	De novo	
PatMaN	miRNA	OS	1.2.2	Linux,Mac	31	8,41		BMAP	FASTA	TSV	1	*	Y	Y	N	A	G	N	N N	N
BS Seeker	Bisulfite	OS	1.2	Linux,Mac	31	8,41		RazerS	FASTA	TSV	-	-	3	0	N	U	-	SM	Y N	N
Slider II	DNA	OS	1.1	Linux,Mac,Windows	31	8,41		SHRIMP	FASTA	TSV	93	Y		N	B,S		N N Y	N		
GNUMAP	DNA	OS	3.0.2	Linux,Mac	31	8,41		BWA	FASTA	SAM TSV	16	1K	Score	Score	Y	B	G	SM/DM	Y N	N
MOM	DNA	Bin	0.6	Linux,Mac,Windows	30	8,41		BWA-SW	FASTA	SAM TSV		Y	0	N	A	L	SM	N Y	N	
Bismark	Bisulfite	OS	0.7.3	Linux,Mac	30	8,41		CloudBurst	FASTA	TSV	16	10K	Score	Score	N	U	-	SM	Y Y	N
BRAT	Bisulfite	OS	1.2.3	Linux	30	8,41		ProbeMatch	FASTA	TSV		Y	0	N			N N Y	N		
SOAPSplice	RNA	Bin	1.8	Linux,Mac	30	8,41		TopHat	FASTA	TSV	13	3K	5	2	Y	U	-	SM	Y Y	De novo
WHAM	DNA	OS	0.14	Linux, Unix	30	8,41		Bowtie	FASTA	SAM TSV	5	128	5	3	N	A,B,R,U,S	G	N	Y Y	De novo
MicroRazerS	miRNA	OS	0.1	Linux	30	8,41		MOM	FASTA	SAM TSV	10	*	Score	0	N	S			N N N	
RUM	RNA	OS	1.11	Linux,Mac	30	8,41		PASS	FASTA	SAM TSV	-	-	Y	Y	Y	B	-	SM	N Y	De novo
ProbeMatch	DNA	OS	1	Linux,Mac	30	8,41		Slider II	FASTA	SAM TSV BED	-	-	Y	Y	Y	B	-	SM	N N	De novo
X-Mate	DNA	OS	1	Linux,Mac	30	8,41		QPALMA	FASTA	ELAND	36	50	3	Y	N	A,B		N N N	Lib	
SSAHA2	DNA	Bin	2.5.5	Linux,Mac	30	8,41		SOCS	FASTA	SAM BED Counts	-	-	Y	0	N	S	-	DM	Y N	Lib
Novoalign	DNA	Bin	V2.08.01	Linux	30	8,41		MAQ	FASTA	SAM	15	48K	Score	Score	N	B,S	L	N N Y	N	
VMATCH	DNA	Bin		Linux,Mac	30	8,41</														

M00628:11:00000000-A1P5L:1:1112:26953:13136
163 CP000921 20 60 149M
= 108 239
CCACTATGTTTCGATAAAAAGCTTAATAAAT
?????BBBBBDBDB=?FFECFACCFH>09C

SAM/BAM

Sequence/Alignment Format

Read name

M00628:11:00000000-A1P5L:1:1112:26953:13136

163 CP000921 20 60 149M

= 108 239

CCACTATTTTCGATAAAAAGCTTAATAAT

Read sequence

? ? ? ? ? BBBBBDDB=?FFECFACCCFFHHH>09C

Read quality

SAM/BAM

Courtesy of Nick Croucher, HSPH

Bitwise flag

M00628:11:00000000-A1P5L:1:1112:26953:13136

163 CP000921 20 60 149M

= 108 239

CCACTATGTTTCGATAAAAAGCTTAATAAAT

?????BBBBBDBDB=?FFECFACCFH>09C

SAM/BAM

Courtesy of Nick Croucher, HSPH

	Mapping position		Mapping quality	
M00628:11:00000000-A1P5L:1:1112:26953:13136				
163	CP000921	20	60	149M
=	108	239		
CCACTATGTTTCGATAAAAAGCTTAATAAAT				
?????BBBBBDBDB=?FFECFACCFH>09C				

SAM/BAM

Courtesy of Nick Croucher, HSPH

M00628:11:00000000-A1P5L:1:1112:26953:13136

163 CP000921 20 60 149M

= 108 239

Alignment

CCACTATGTTTCGATAAAAAGCTTAATAAAT

? ? ? ? ? BBBBBDDBDB=?FFECFACCFH>09C

SAM/BAM

Courtesy of Nick Croucher, HSPH

M: match
I: insertion relative to reference
D: deletion relative to reference

RefPos:	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
Ref:	C	C	A	T	A	C	T	G	A	A	C	T	G	A	C	T	A	A	C
Read:				A	C	T	A	G	A	A		T	G	G	C	T			

CIGAR strings

Courtesy of Nick Croucher, HSPH

M: match
I: insertion relative to reference
D: deletion relative to reference

RefPos:	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
Ref:	C	C	A	T	A	C	T	G	A	A	C	T	G	A	C	T	A	A	C
Read:				A	C	T	A	G	A	A		T	G	G	C	T			

CIGAR string: 3M1I3M1D5M

CIGAR strings

Courtesy of Nick Croucher, HSPH

Distance to
mate pair

M00628:11:00000000-A1P5L:1:1112:26953:13136

163	CP000921	20	60	149M
=	108	239		

CCACTATGTTTCGATAAAAAGCTTAATAAAT
?????BBBBBDBDB=?FFECFACCFH>09C

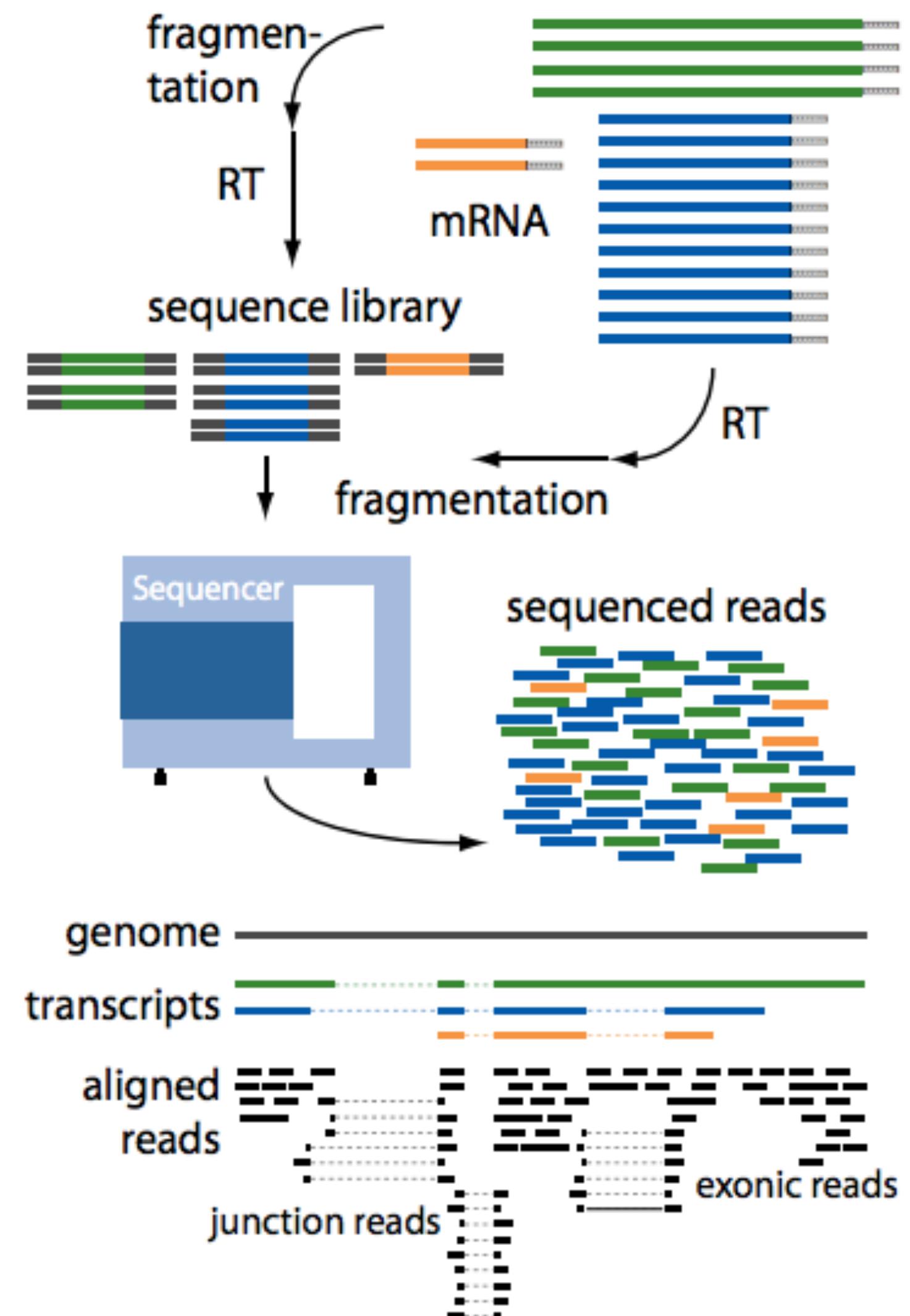
Mate mapped
to same
reference?

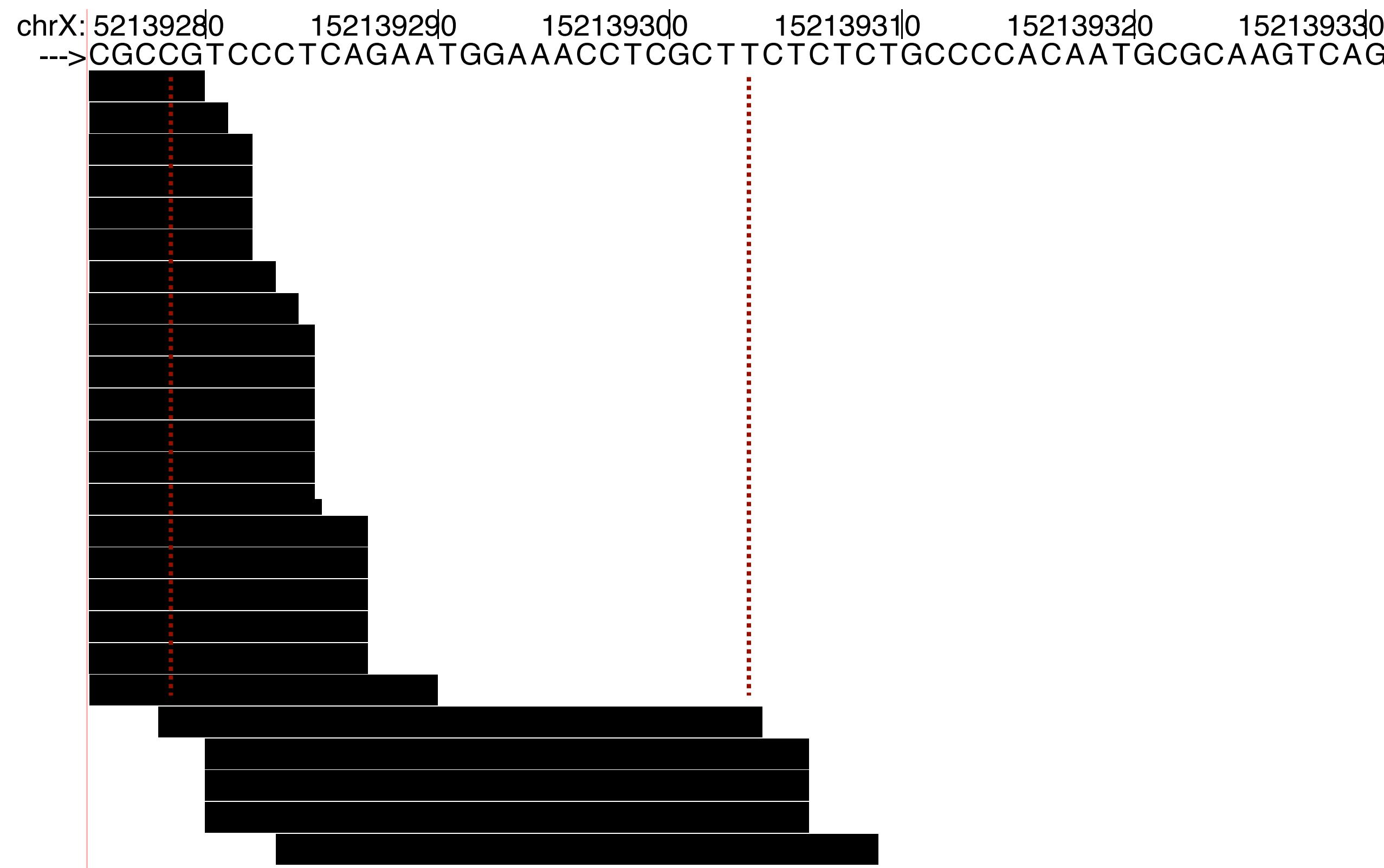
Total length

SAM/BAM

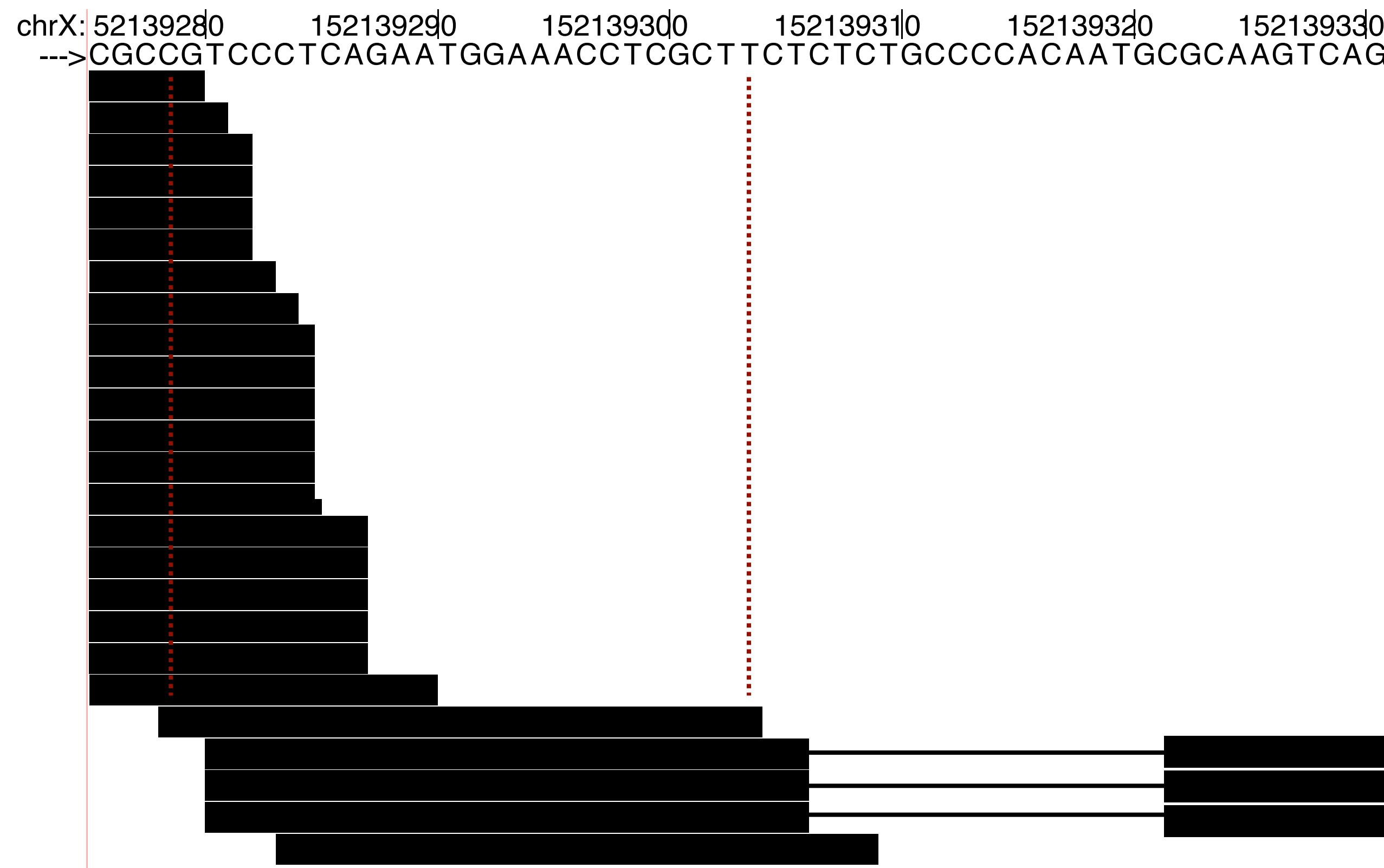
Courtesy of Nick Croucher, HSPH

RNA-seq

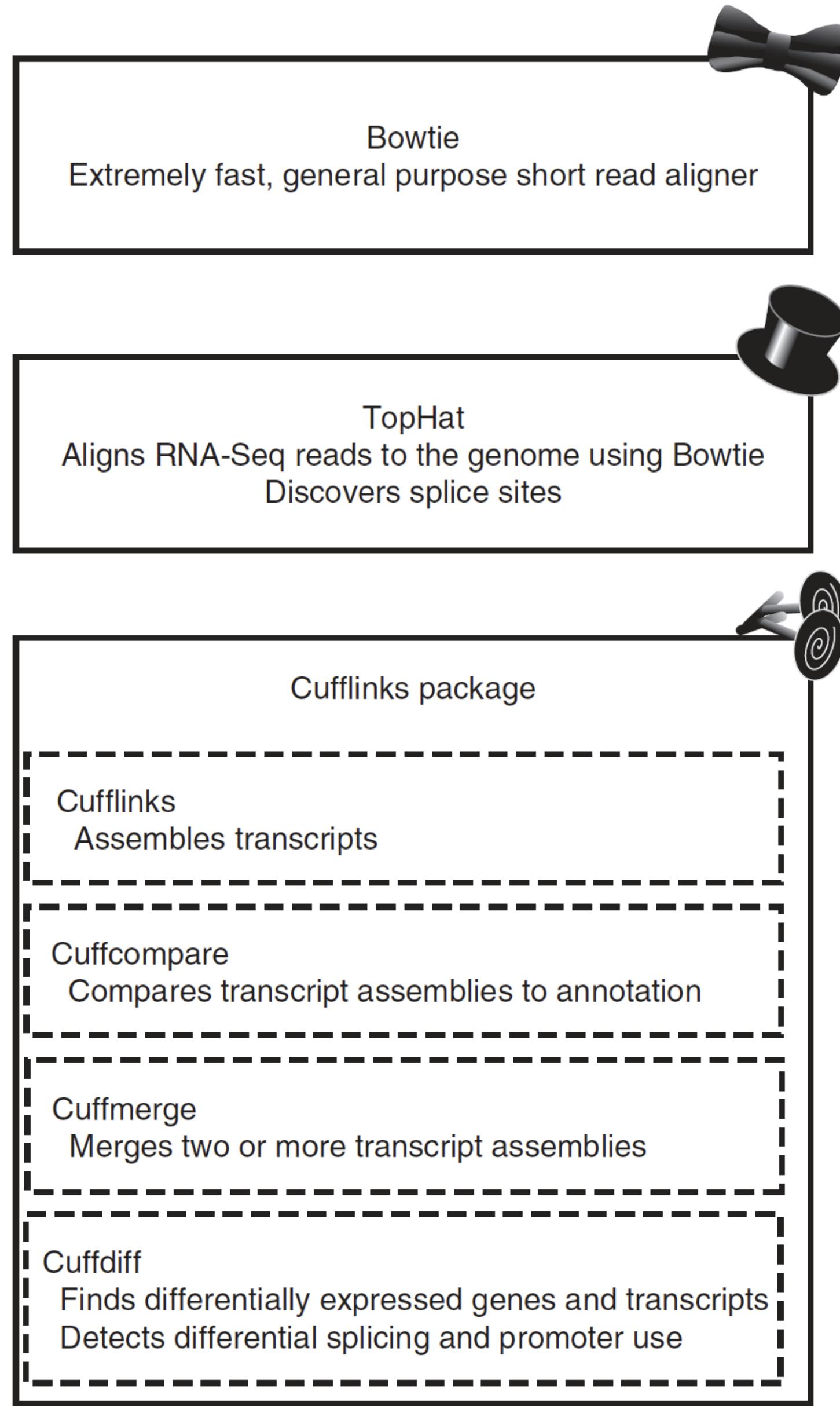




Mapping to a Reference Genome



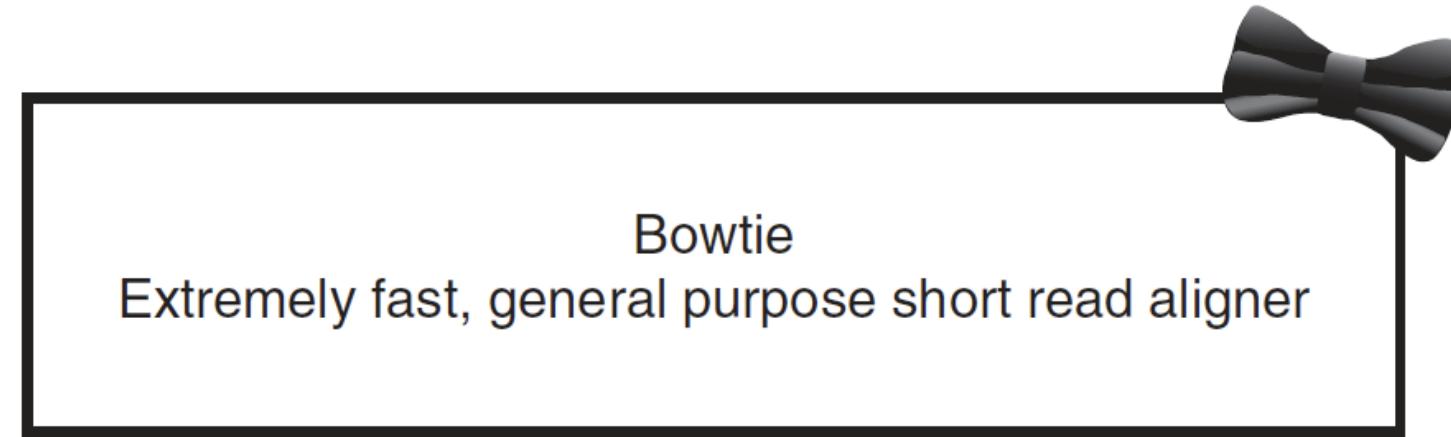
Mapping to a Reference Genome



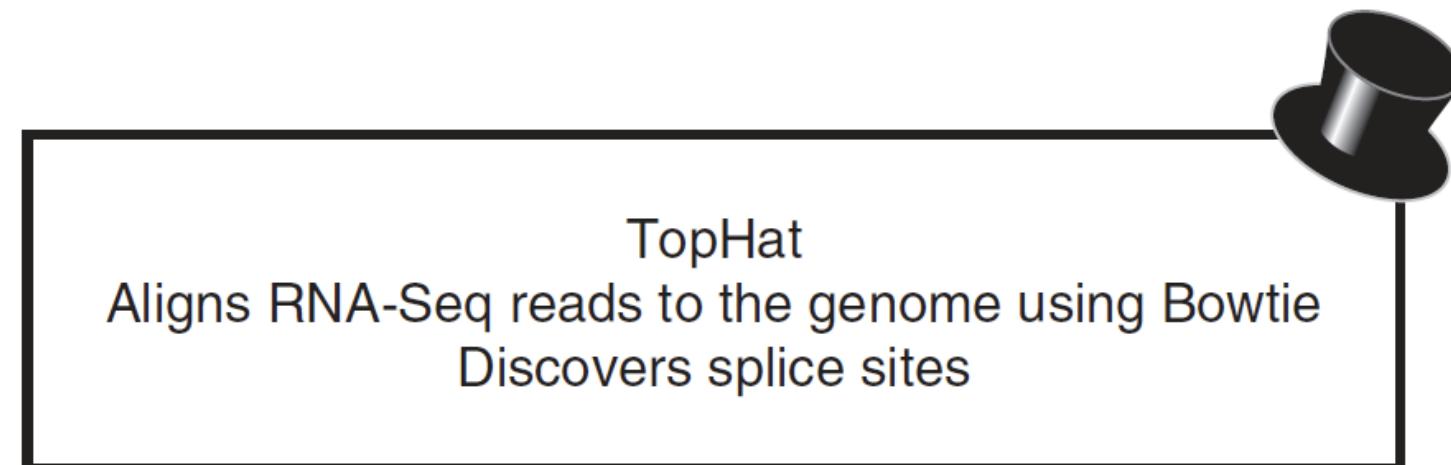
Uses Burrows-Wheeler indexing for aligning reads. With **Bowtie2** there is no upper limit on the read length.

The Tuxedo suite

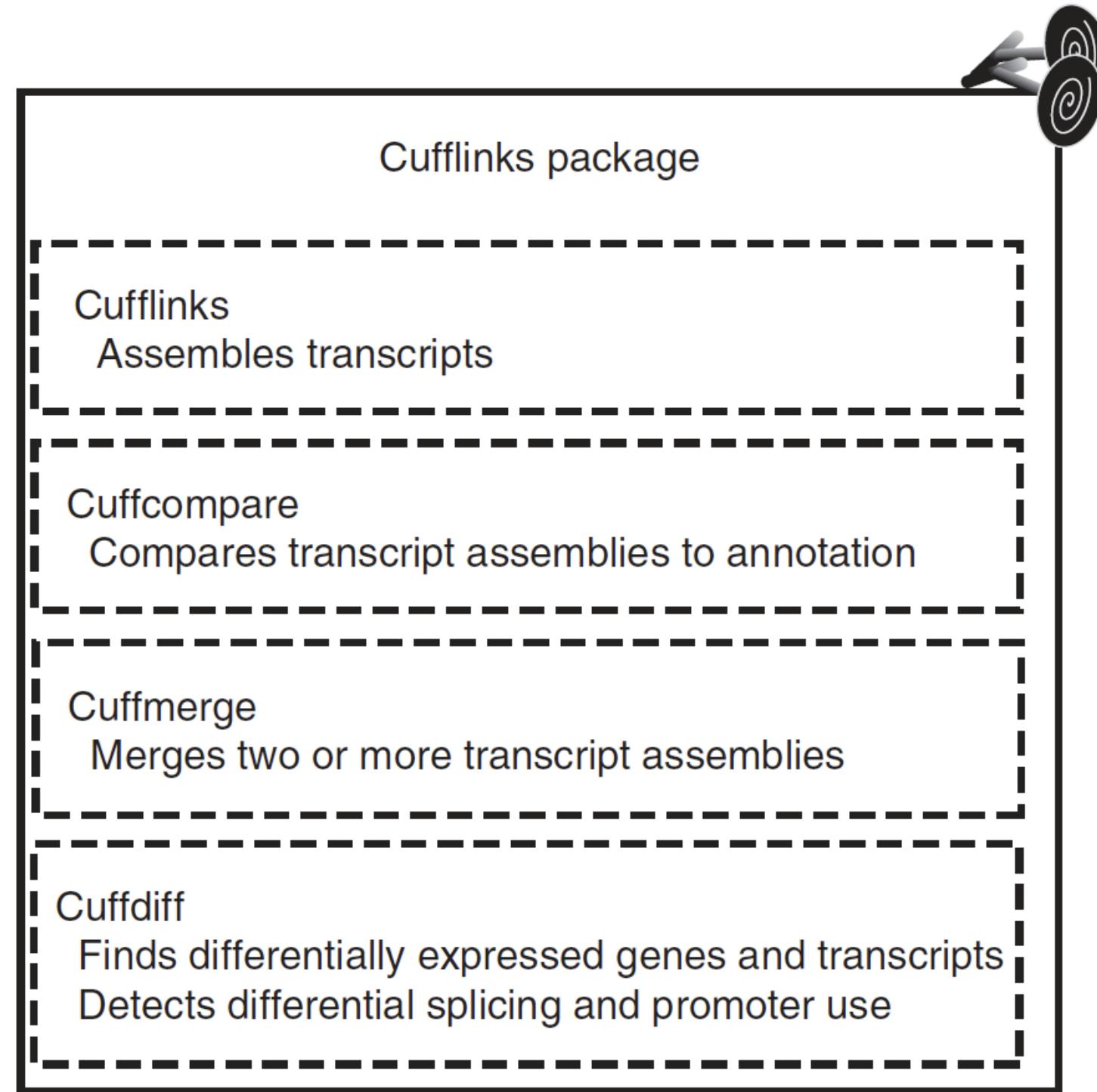
Trapnell et al., Nature Protocols, March 2012



Uses Burrows-Wheeler indexing for aligning reads. With **Bowtie2** there is no upper limit on the read length.

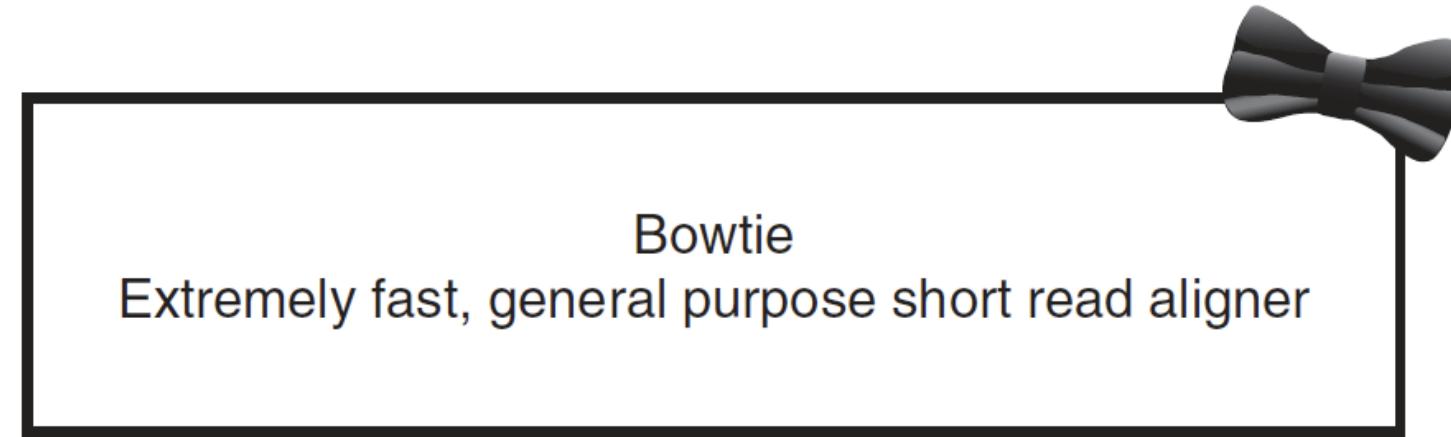


Tophat uses either Bowtie or Bowtie2 to align reads in a splice-aware manner and aids the discovery of new splice junctions

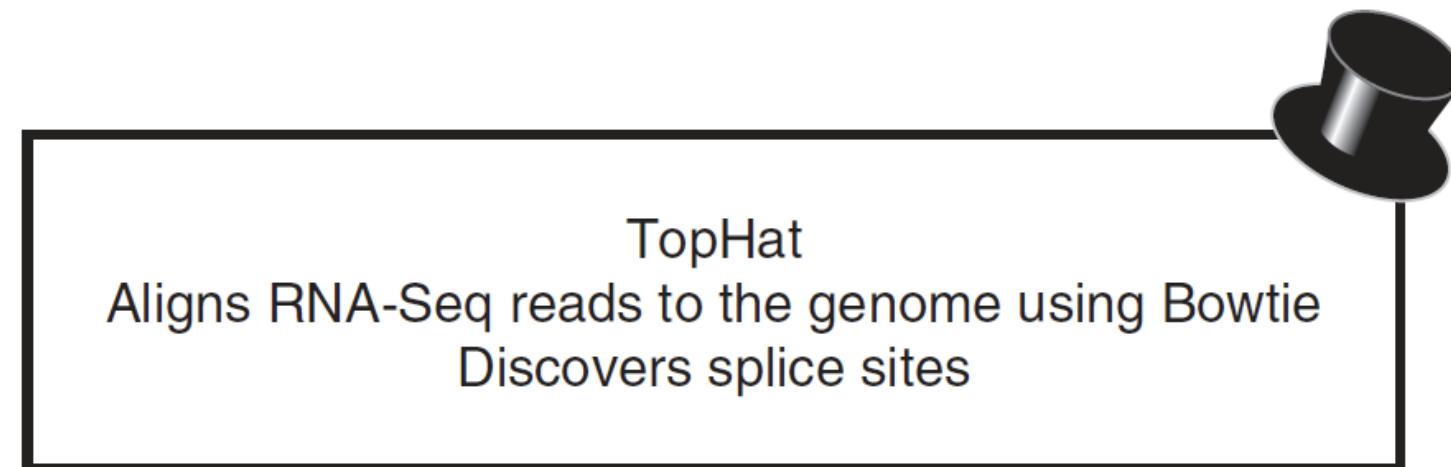


The Tuxedo suite

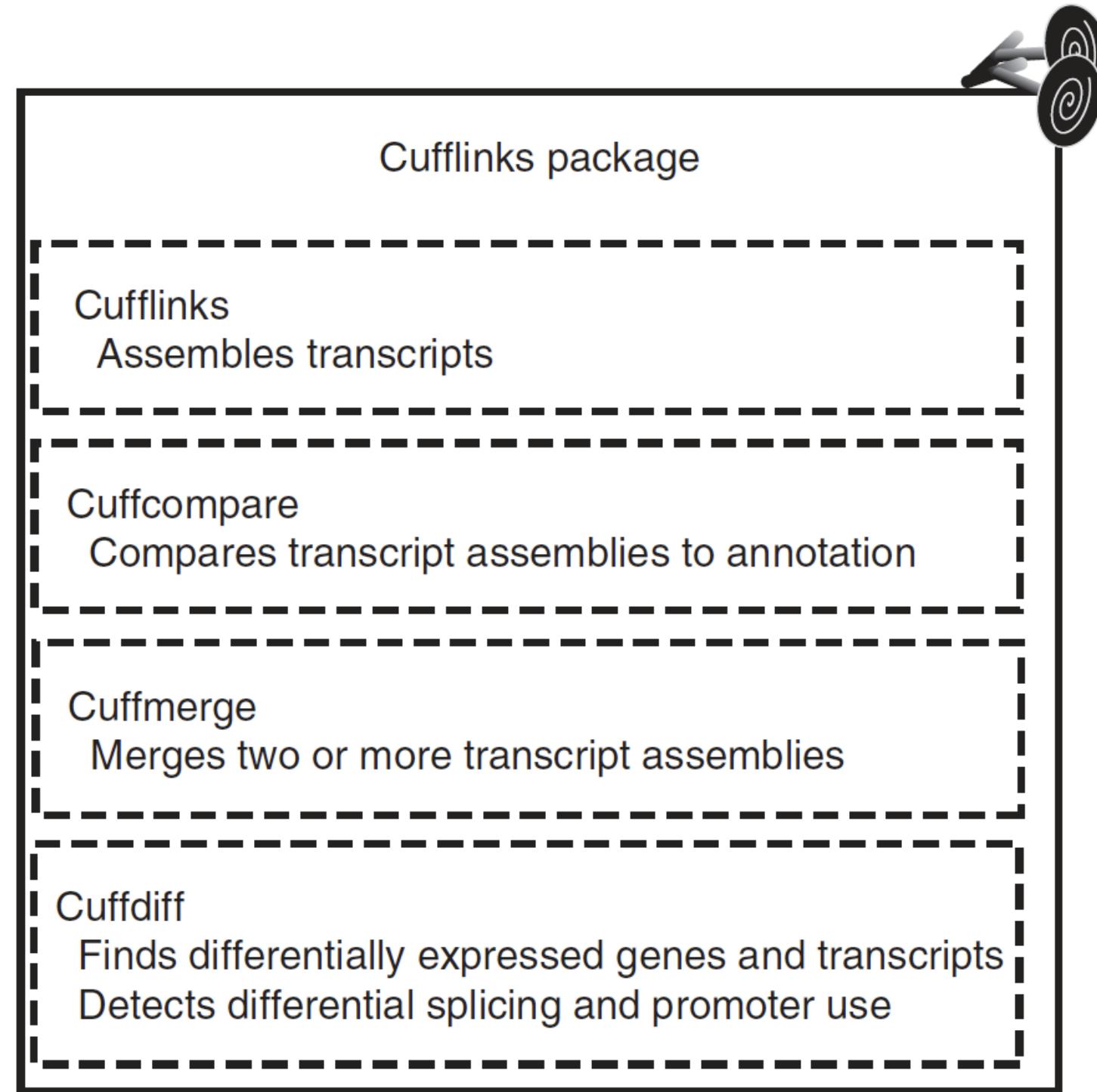
Trapnell et al., Nature Protocols, March 2012



Uses Burrows-Wheeler indexing for aligning reads. With **Bowtie2** there is no upper limit on the read length.

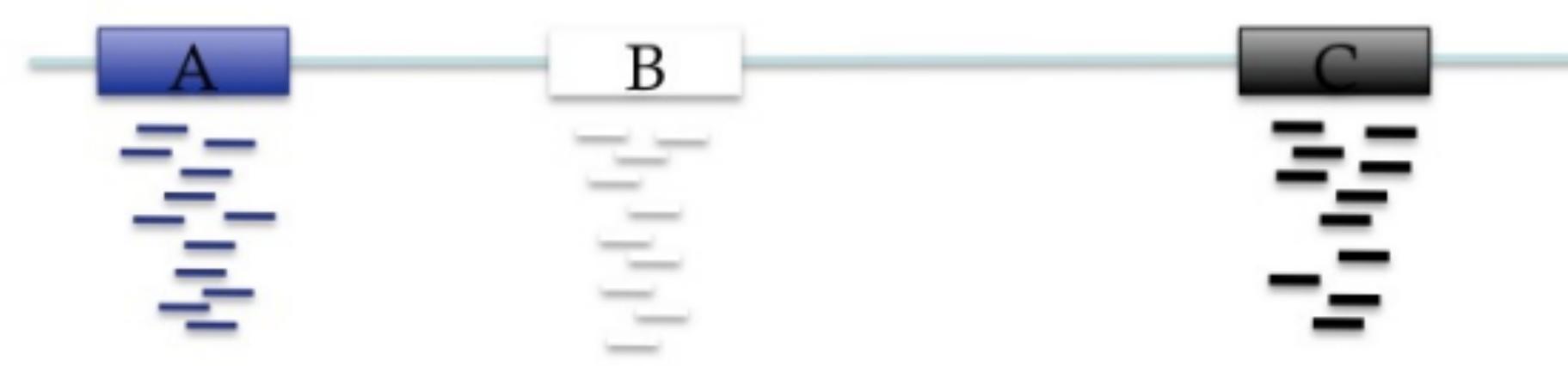


Tophat uses either Bowtie or Bowtie2 to align reads in a splice-aware manner and aids the discovery of new splice junctions



Cufflinks does *de novo* (reference-based) transcriptome assembly

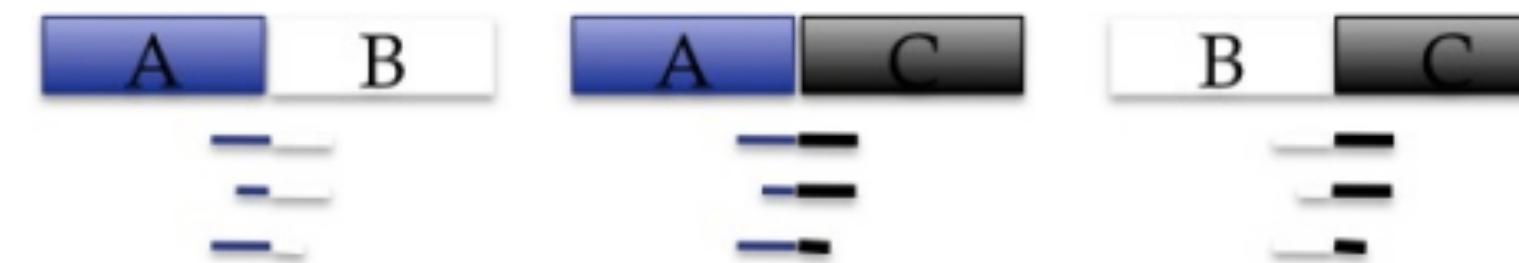
Cuffdiff does statistical analysis and identifies differentially expressed transcripts



identify candidate exons
via genomic mapping



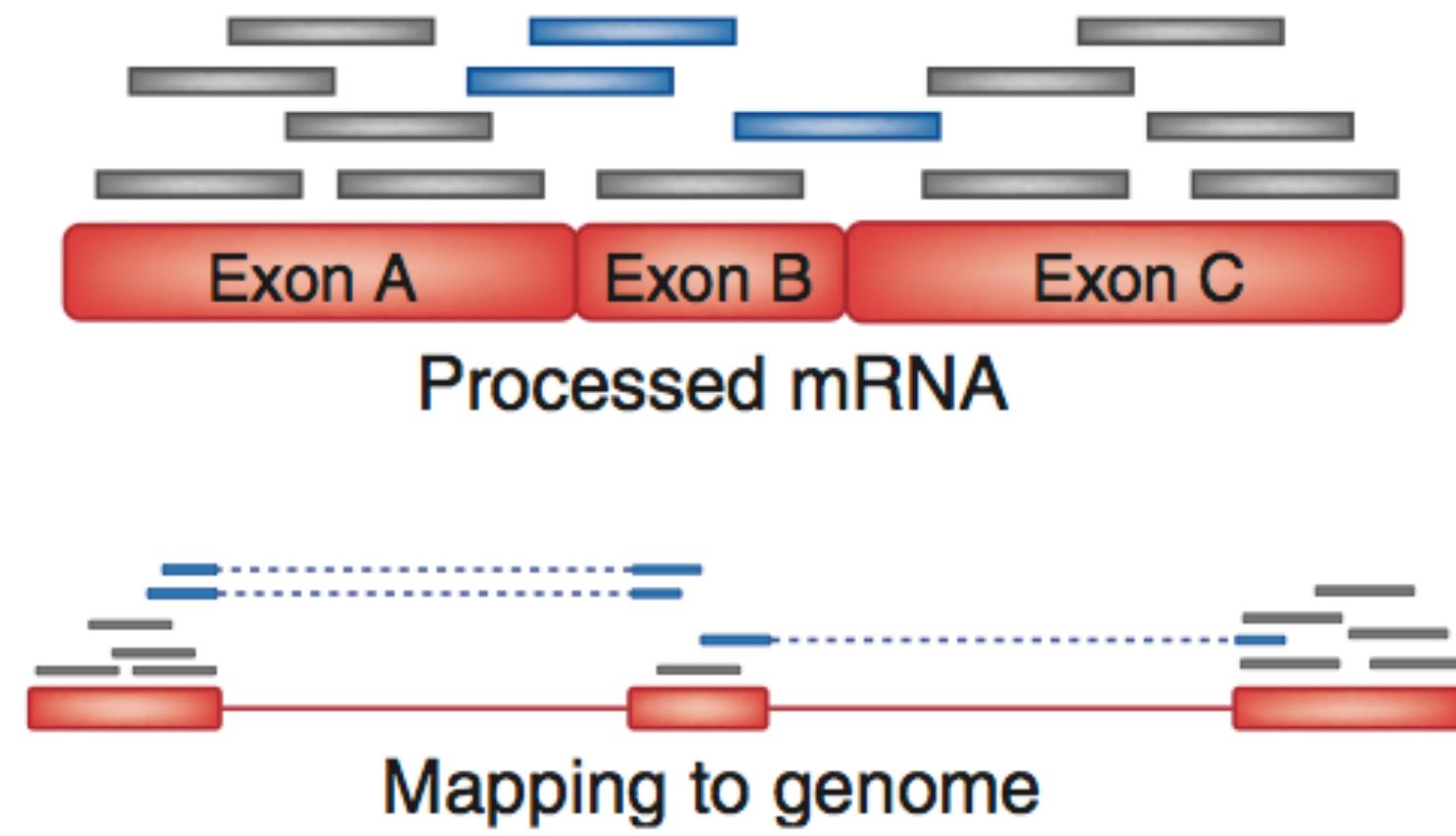
Generate possible pairings
of exons



Align reads to possible
junctions

(Trapnell, 2010)

TopHat alignment concept



... GCAAACCACTGACCTGACTACTACGTCGTAACGTACACGGTAGCT... CCGTAGAATTGACTGTGTTG...

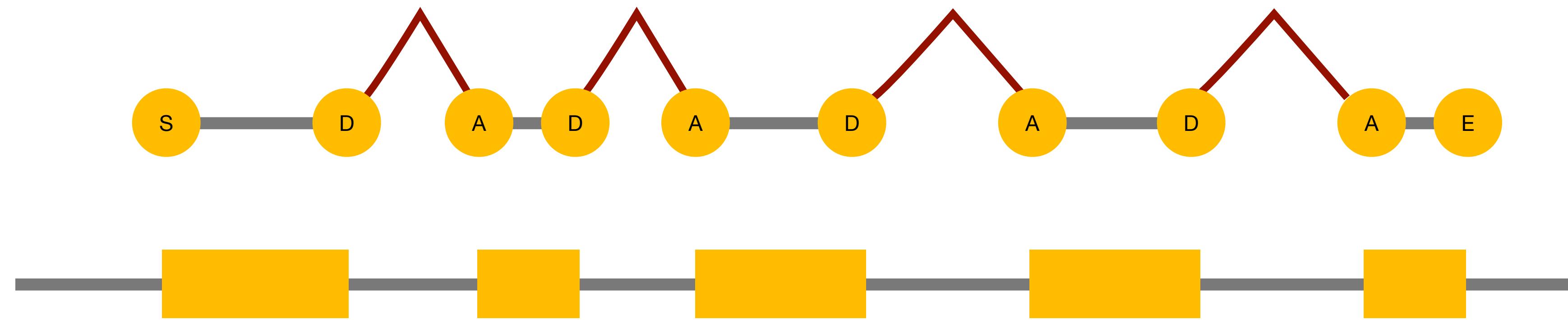
... GCAAACCACTGACCTGACTACTACGTCGTAACGTAC
CAAACCACTGACCTGACTACTACGTCGTAACGTACA
AAACCACTGACCTGACTACTACGTCGTAACGTACAC
AACCACTGACCTGACTACTACGTCGTAACGTACACG
ACCACTGACCTGACTACTACGTCGTAACGTACACG
CCAGTGACCTGACTACTACGTCGTAACGTACACG
CAGTGACCTGACTACTACGTCGTAACGTACACG
AGTGACCTGACTACTACGTCGTAACGTACACG
GTGACCTGACTACTACGTCGTAACGTACACG
TGACCTGACTACTACGTCGTAACGTACACG

A
AA
AAT
AATT
AATTG
AATTGA

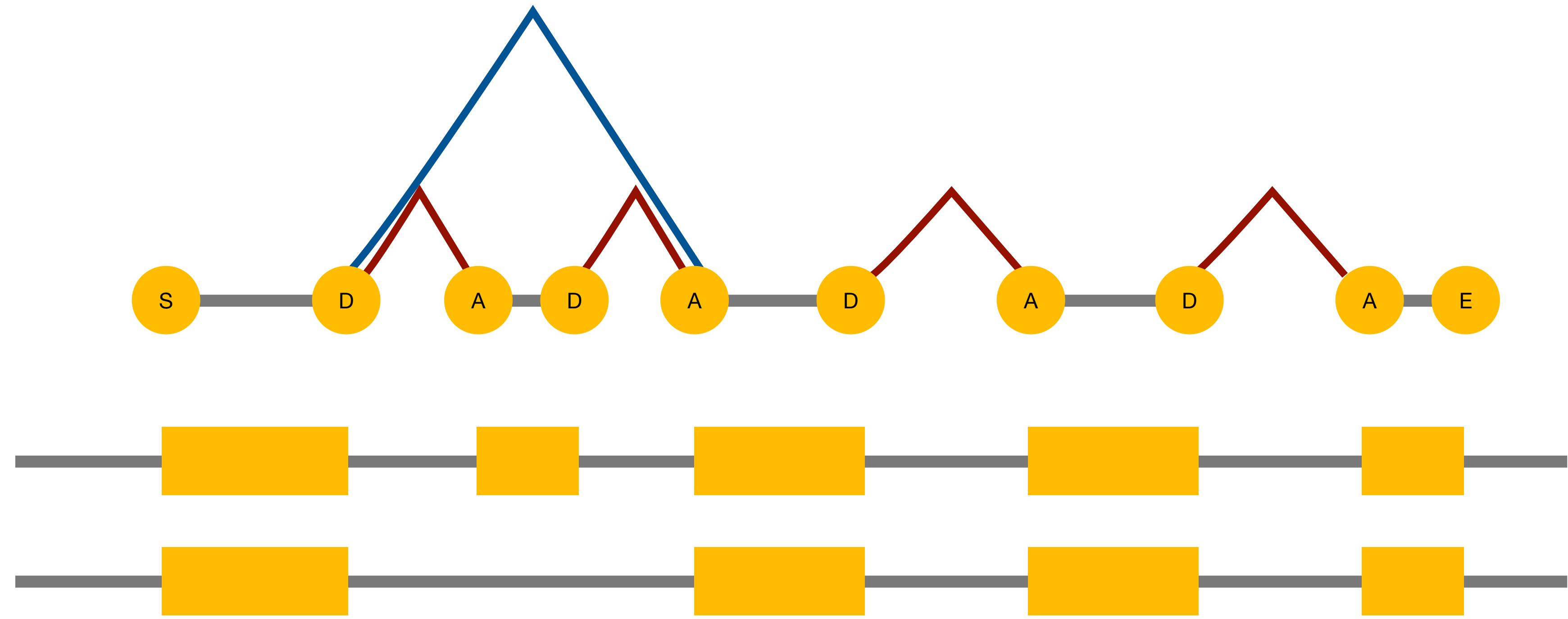
Transcript discovery



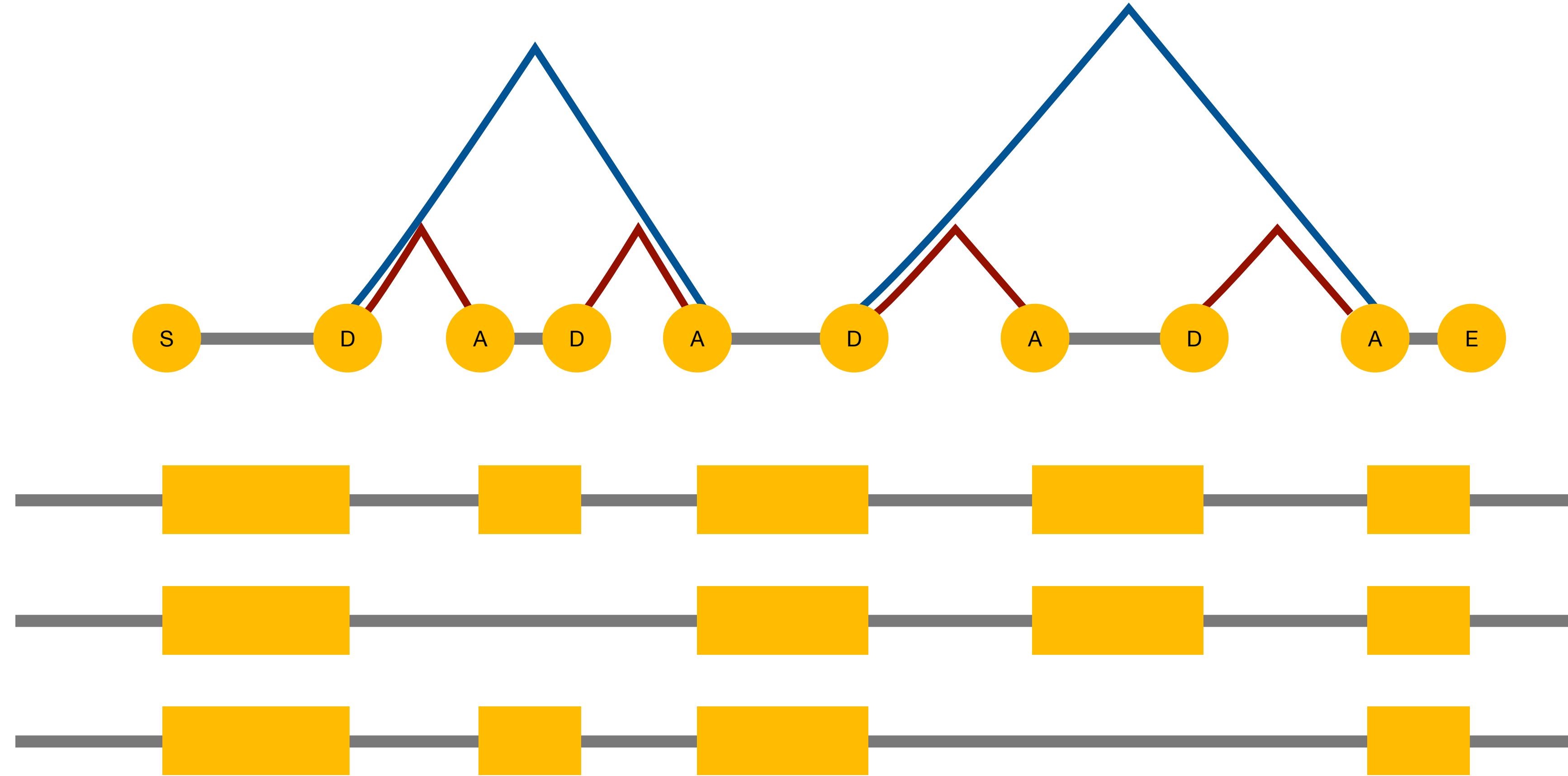
Splice graphs



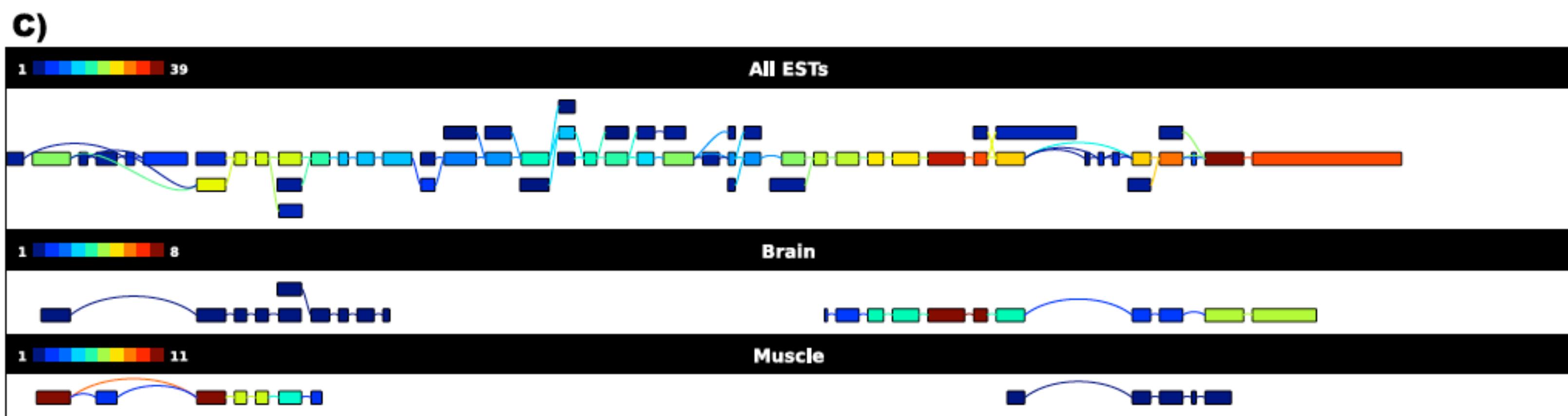
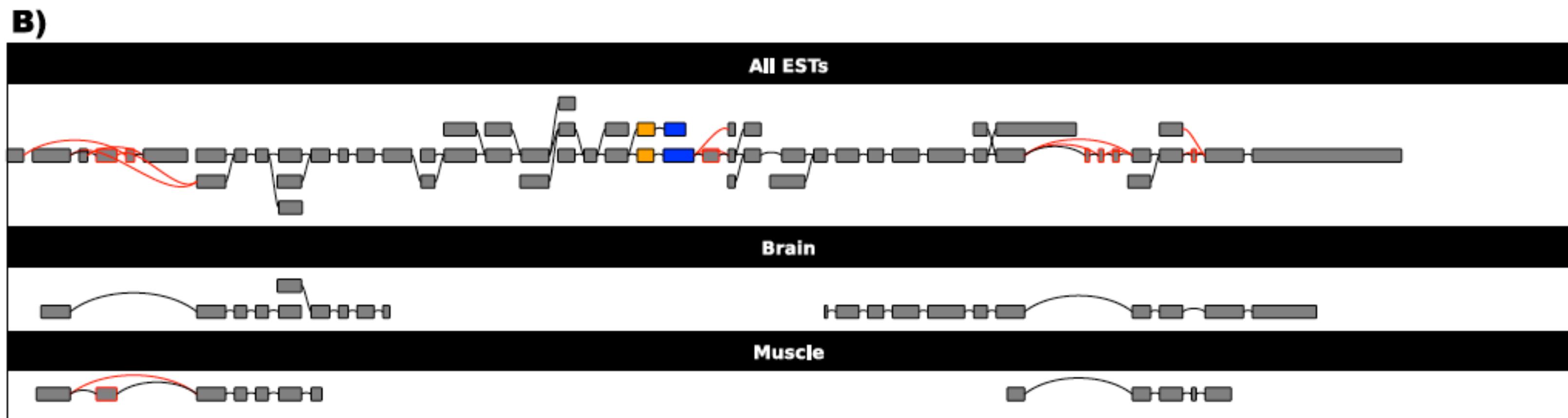
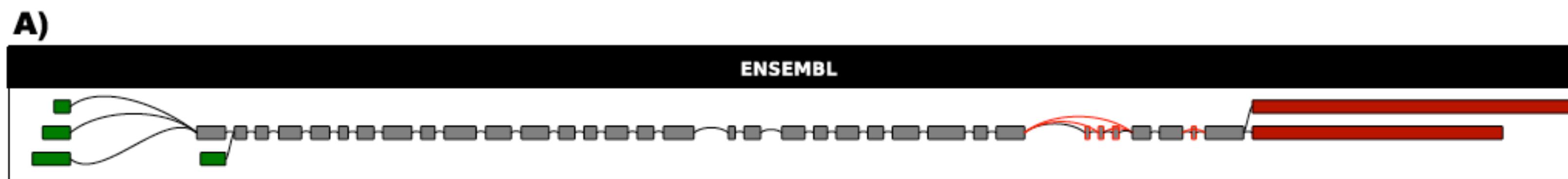
Splice graphs

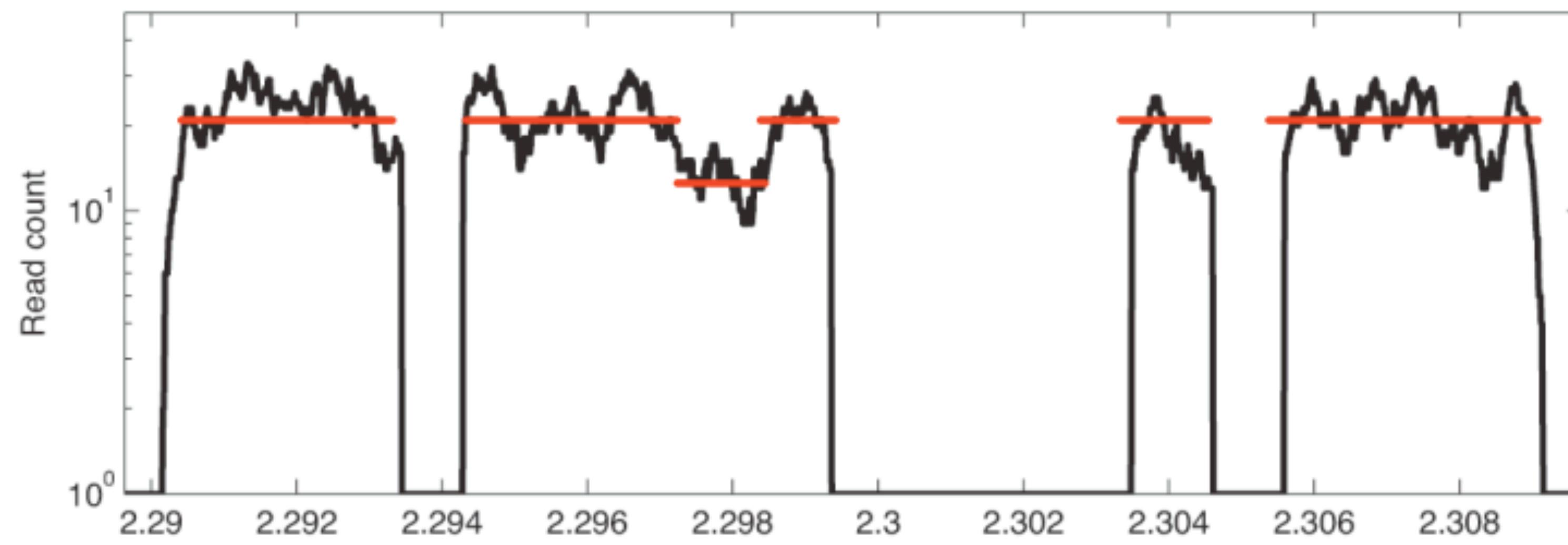
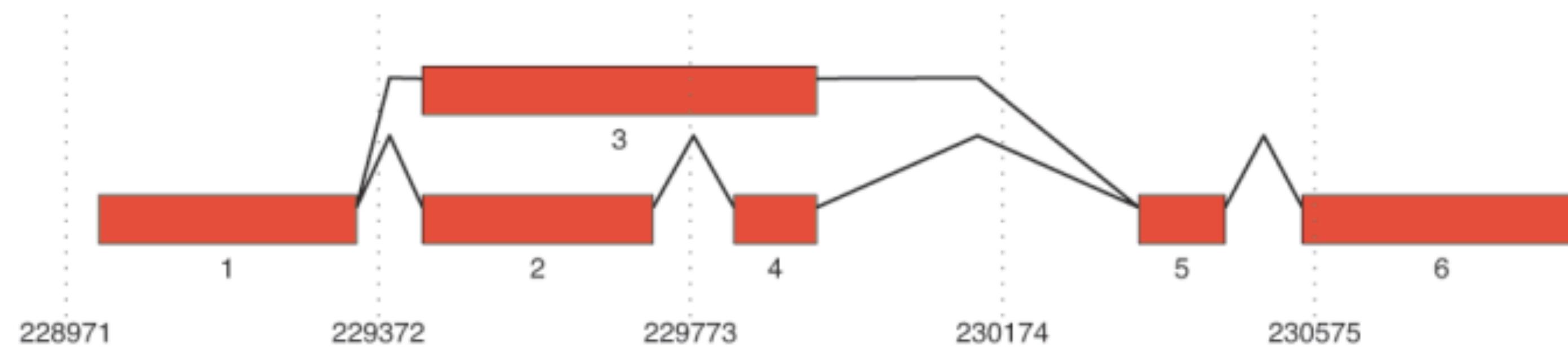


Splice graphs



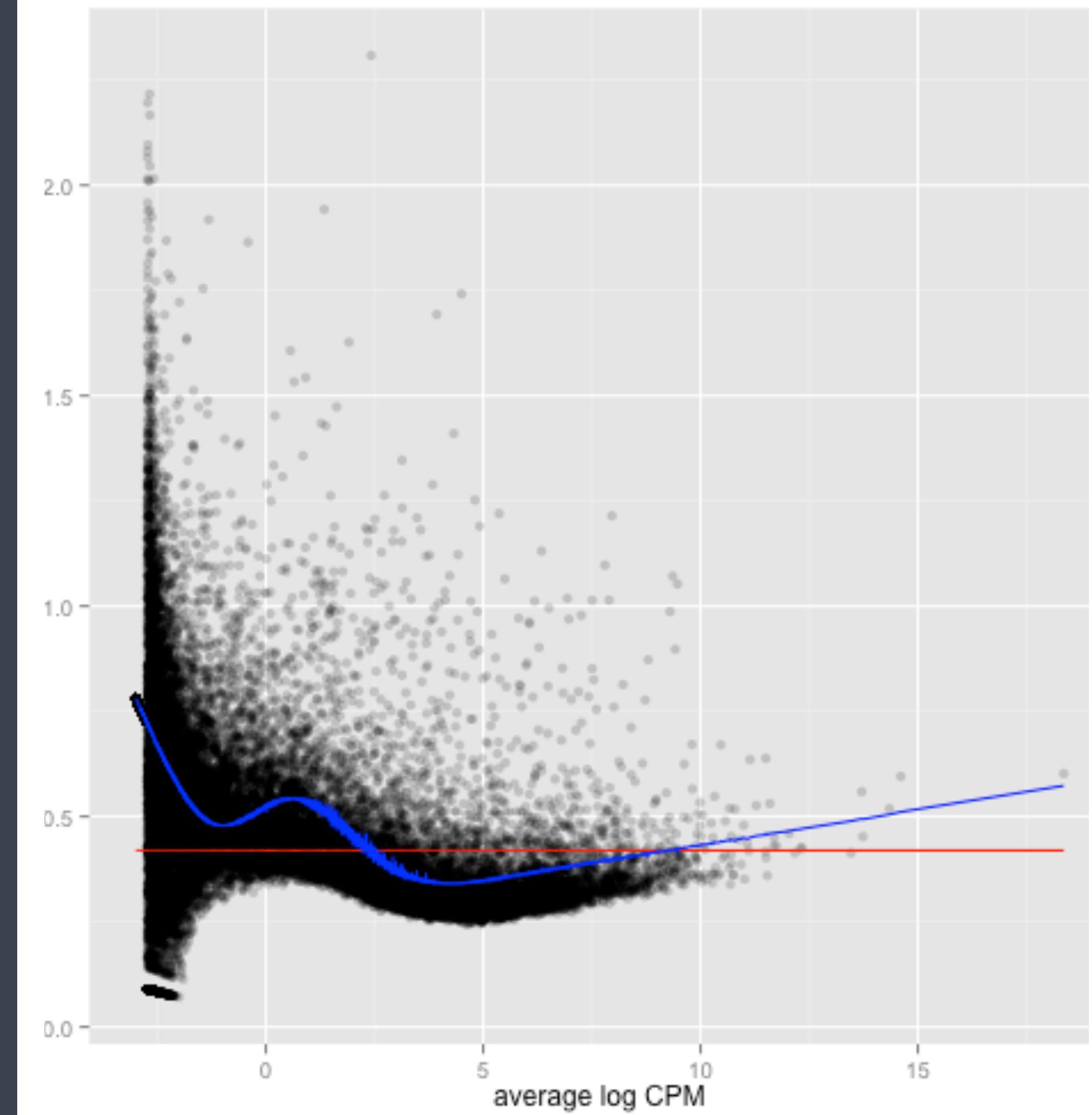
Splice graphs





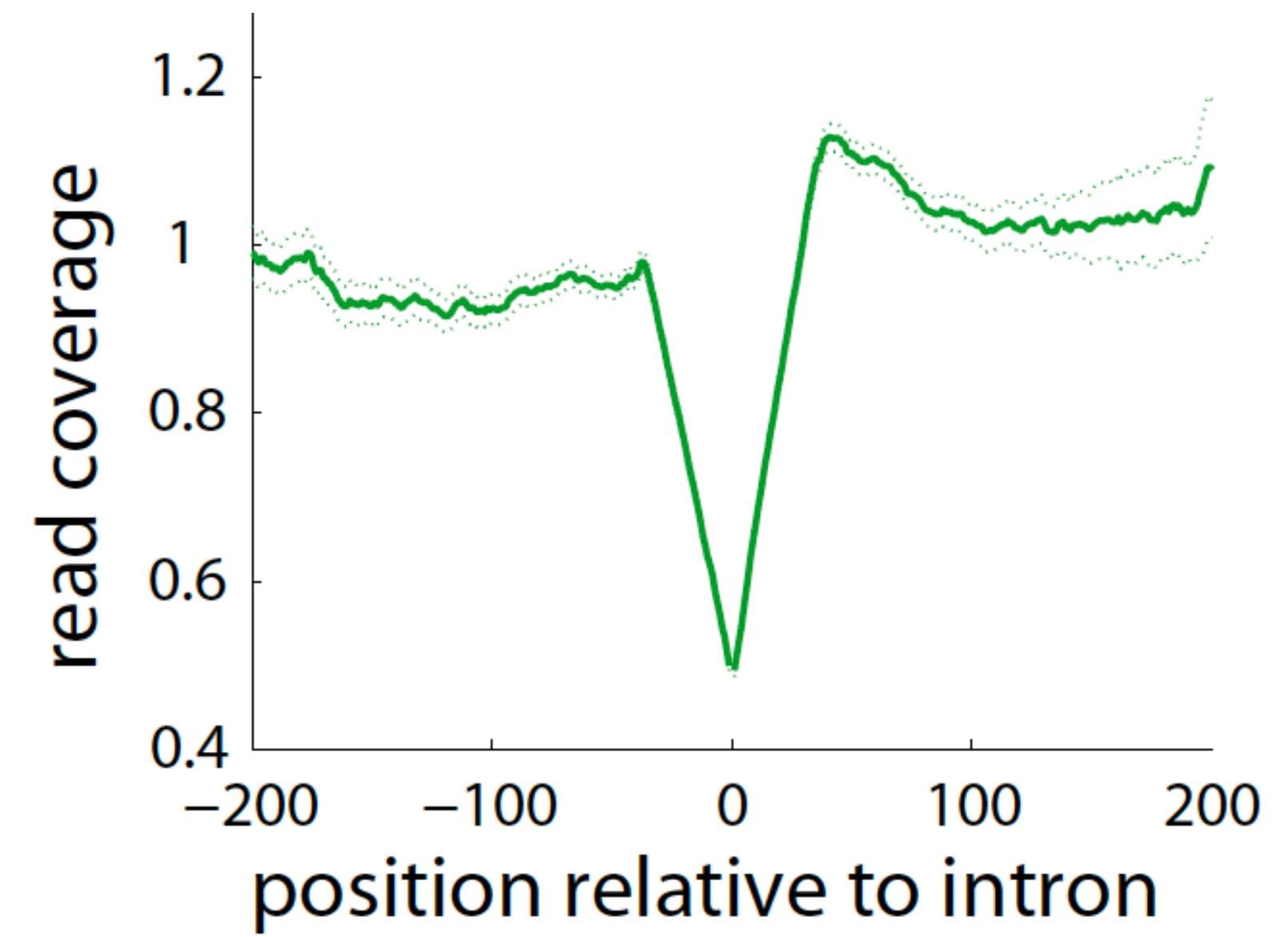
Quantification

Count normalization



Biases

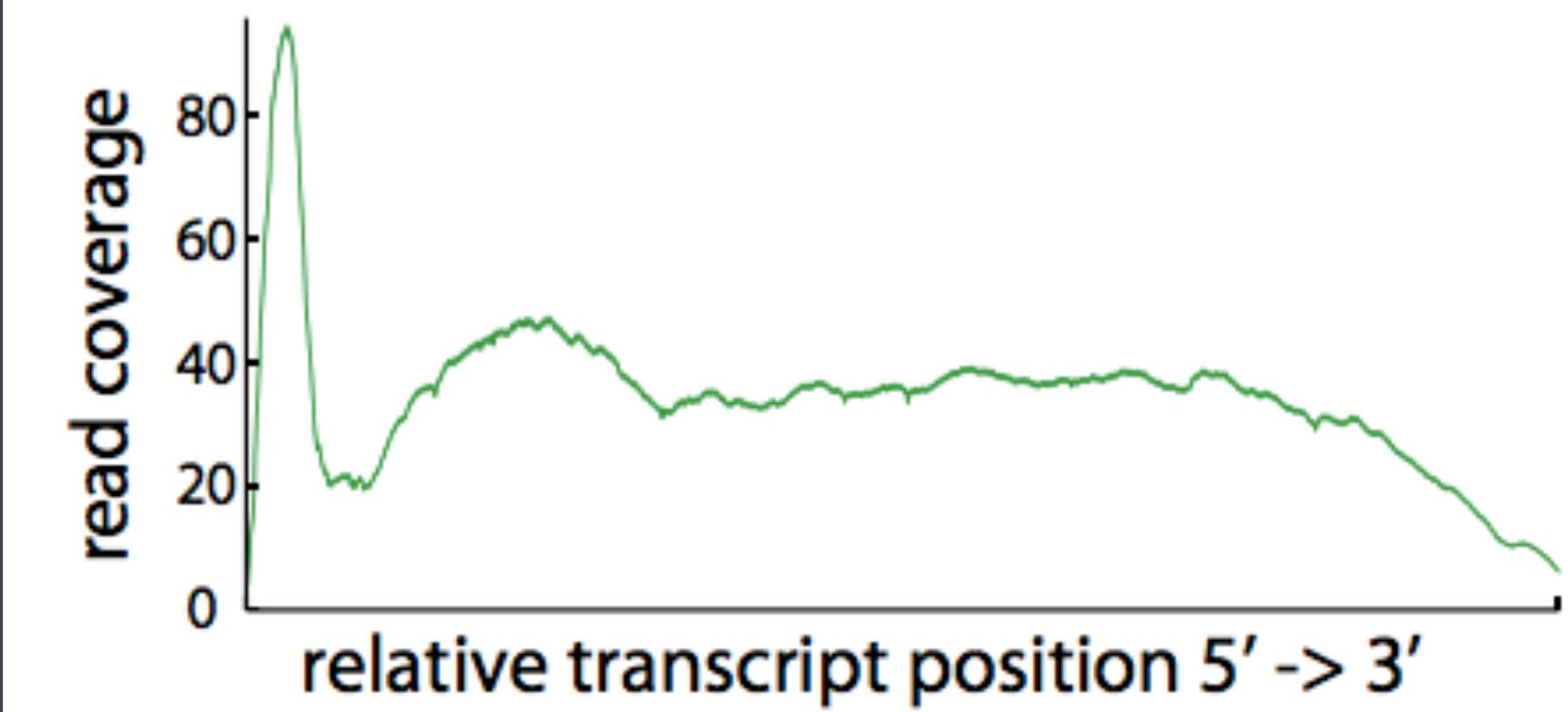
Read mapping

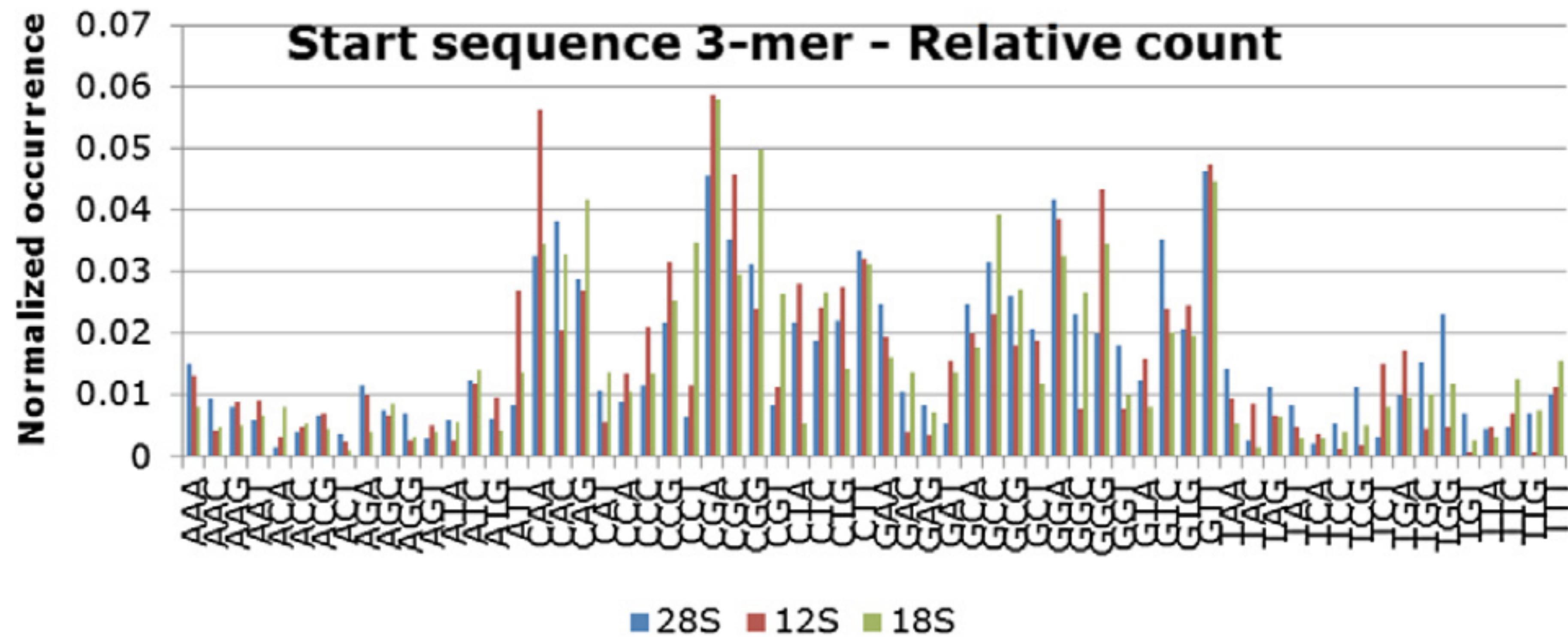


Biases

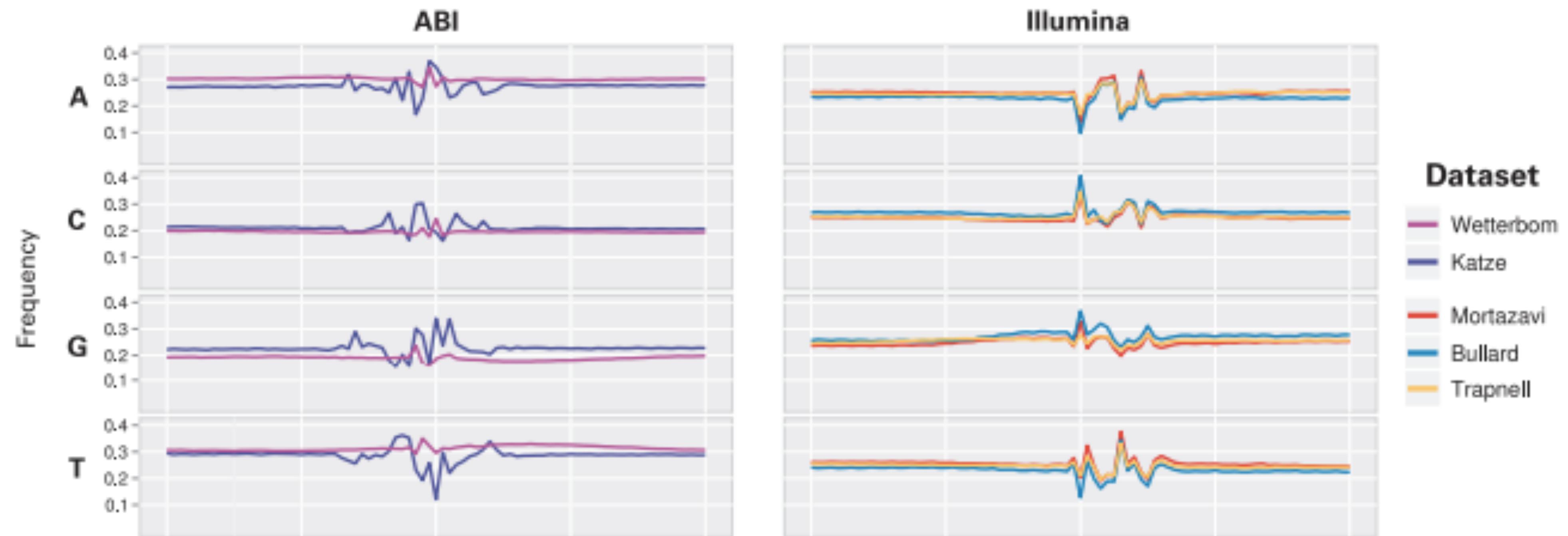
Library construction

Transcript length



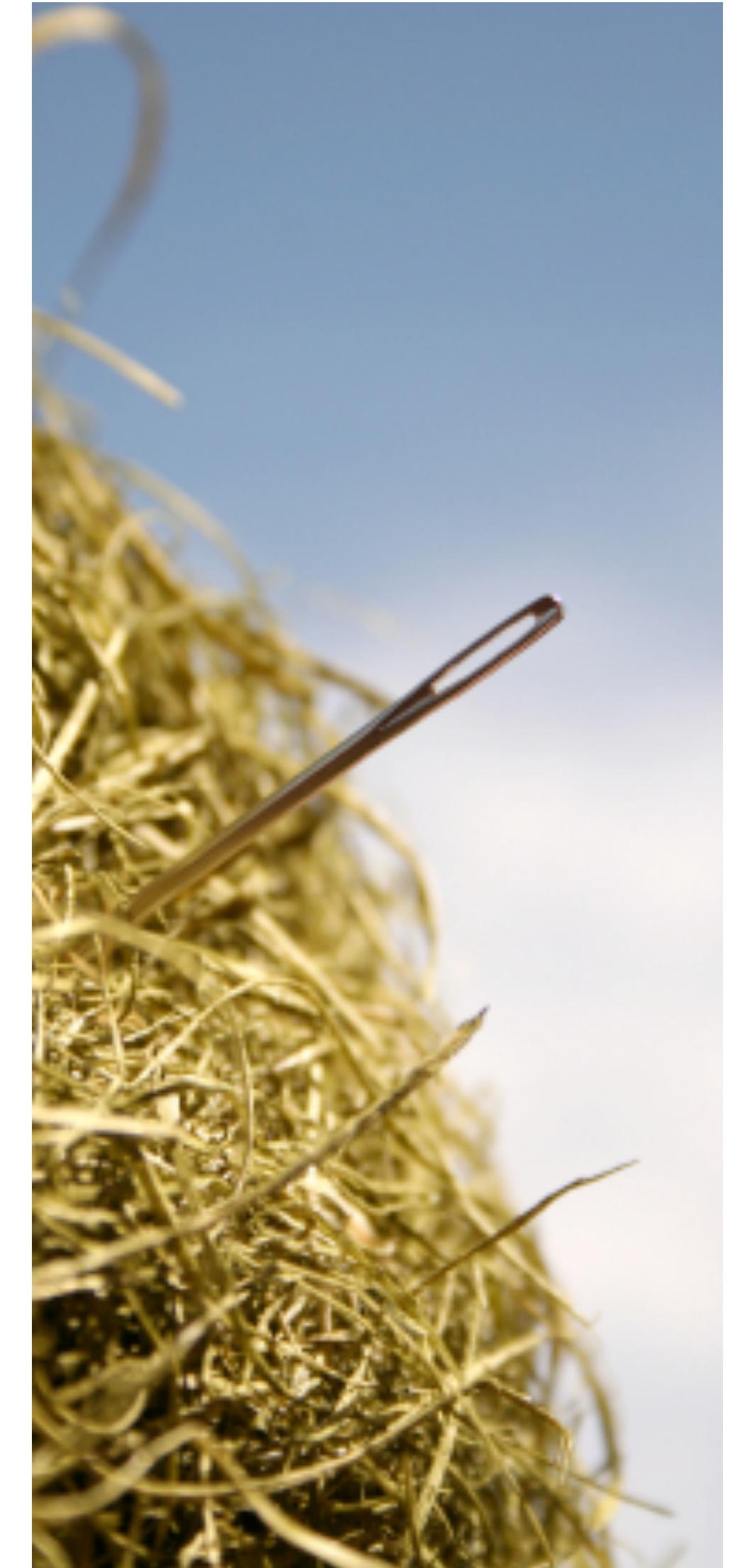


Not-so-random priming



Start/end bias differs by technology

Making sense of your gene list



goo.gl/ENP6w4

Next steps

PROCEEDINGS

Open Access

Galaxy CloudMan: delivering cloud compute clusters

Enis Afgan¹, Dannon Baker¹, Nate Coraor², Brad Chapman³, Anton Nekrutenko², James Taylor^{1*}

From The 11th Annual Bioinformatics Open Source Conference (BOSC) 2010
Boston, MA, USA. 9-10 July 2010

Run your own Galaxy instance

Galaxy CloudMan Console

Welcome to the Galaxy Cloud Manager. This application will allow you to manage this cloud and the services provided within. If this is your first time running this cluster, you will need to select an initial data volume size. Once the data store is configured, default services will start and you will be able to add and remove additional services as well as 'worker' nodes on which jobs are run.

[Terminate cluster](#)[Add instances ▾](#)[Remove Instances](#)[Access Galaxy](#)

Status

Cluster name: local test**Disk status:** 0 / 0 (0%) **Worker status:** Idle: 0 Available: 0 Requested: 0**Service status:** Applications  Data [Cluster status log](#) 

Instant Galaxy cluster

Cluster name	<input type="text"/>	Name of your cluster used for identification. This can be any name you choose.
Password	<input type="password"/>	Your choice of password, for the CloudMan web interface and accessing the Amazon instance via ssh or FreeNX.
Access key	<input type="text"/>	Your Amazon Access Key ID. Available from the security credentials page .
Secret key	<input type="text"/>	Your Amazon Secret Access Key. Also available from the security credentials page .
Instance type	<input style="width: 150px;" type="text" value="Large"/> ▼	Amazon instance type to start.

BioCloudCentral

<https://biocloudcentral.herokuapp.com/launch>

Need for Unix

<http://ritg.med.harvard.edu/classes.html>

2477777 -rw-r--r--	1 root	wheel	500	23 Jun	2009	networks
2476106 -r--r--r--	1 root	wheel	1.6K	18 May	2009	newsyslog.conf
2377284 drwxr-xr-x	4 root	wheel	136B	3 Sep	2009	newsyslog.d
2387394 -rw-r--r--	1 root	wheel	132B	10 Jul	2009	notify.conf
2483297 -rw-r--r--	1 root	wheel	366B	18 May	2009	ntp-restrict.conf
2733572 -rw-r--r--@	1 root	wheel	27B	20 Oct	2009	ntp.conf
2471151 drwxr-xr-x	7 root	wheel	238B	17 Jul	2009	openldap
4929956 drwxr-xr-x	3 root	wheel	102B	30 Nov	2007	opt
2403948 drwxr-xr-x	17 root	wheel	578B	23 Dec	2010	pam.d
2386884 -rw-r--r--	1 root	wheel	3.6K	23 Jun	2009	passwd
2479751 -rw-r--r--	1 root	wheel	45B	23 Jun	2009	paths
2475307 drwxr-xr-x	3 root	wheel	102B	11 Jul	2009	paths.d
2475966 -rw-r--r--	1 root	wheel	1.2K	19 Jul	2009	pear.conf
2466447 drwxr-xr-x	5 root	wheel	170B	18 May	2009	periodic
13701882 -r--r--r--	1 root	wheel	67K	9 Mar	00:49	php.ini.default
2733641 -r--r--r--	1 root	wheel	44K	6 Feb	2009	php.ini.default-f
2485974 drwxr-xr-x	23 root	wheel	782B	27 Jun	2010	postfix
2486265 drwxr-xr-x	2 root	wheel	68B	1 Aug	2009	ppp
2476001 -r--r--r--	1 root	wheel	189B	4 May	2009	profile
2386885 -rw-r--r--	1 root	wheel	5.6K	23 Jun	2009	protocols
2444032 drwxr-xr-x	4 root	wheel	136B	3 Sep	2009	racoon
2387062 -rw-r--r--	1 root	wheel	1.6K	25 Jul	2009	rc.common
2444075 -rw-r--r--	1 root	wheel	5.0K	25 Jul	2009	rc.netboot
2442640 lrwxr-xr-x	1 root	wheel	20B	3 Sep	2009	resolv.conf
						-> /
		lv.conf				
2479752 -rw-r--r--	1 root	wheel	0B	23 Jun	2009	rmtab
2386886 -rw-r--r--	1 root	wheel	971B	23 Jun	2009	rpc
2483260 -rw-r--r--	1 root	wheel	983B	15 Jul	2009	rtadvd.conf
2403837 drwxr-xr-x	7 root	wheel	238B	16 Jun	2009	security
2386887 -rw-r--r--	1 root	wheel	662K	23 Jun	2009	services
2386888 -rw-r--r--	1 root	wheel	179B	23 Jun	2009	shells
4613036 -rw-r--r--	1 root	wheel	2.9K	27 Jun	2010	smb.conf
2497382 -rw-r--r--	1 root	wheel	2.9K	22 May	2009	smb.conf.old
4599486 -rw-r--r--	1 root	wheel	2.9K	6 May	2010	smb.conf.template
2482910 drwxr-xr-x	4 root	wheel	136B	28 Jul	2009	snmp
2403959 -rw-r--r--	1 root	wheel	1.5K	11 Jul	2009	ssh_config
2403960 -rw-r--r--	1 root	wheel	3.6K	11 Jul	2009	sshd_config
2497815 -r--r----	1 root	wheel	1.2K	23 Jun	2009	sudoers
2386889 -rw-r--r--	1 root	wheel	772B	23 Jun	2009	syslog.conf
2386890 -rw-r--r--	1 root	wheel	1.4K	23 Jun	2009	ttys
2475311 drwxr-xr-x	4 root	wheel	136B	28 Jul	2009	xgrid
2479753 -rw-r--r--	1 root	wheel	0B	23 Jun	2009	xtab
2503639 -r--r--r--	1 root	wheel	126B	11 May	2009	zshenv

Ohos-MacBook-Air:etc oho\$ █

TEACHING LAB SKILLS FOR SCIENTIFIC COMPUTING



Who We Are

Our volunteers teach basic software skills to researchers in science, engineering, and medicine. Founded in 1998, we are now part of the Mozilla Science Lab.

What We Do

We run bootcamps all over the world, and provide open access material for self-paced instruction. We also run a training program for people who'd like to help us teach.

How To Help

Like all volunteer organizations, we depend on you to help us help others. You can host a bootcamp, help create new teaching materials, or improve the tools we use.

Software Carpentry



bcbio-nextgen.readthedocs.org

Talk to us early

Involvement in study design to optimize experiments



Contact

bioinformatics@hsph.harvard.edu

<http://bioinformatics.hms.harvard.edu/>

