

Translational bioinformatics and NGS

Center for Health Bioinformatics, Harvard School of Public Health

<intro>



Win Hide



Oliver Hofmann



Shannan Ho Sui



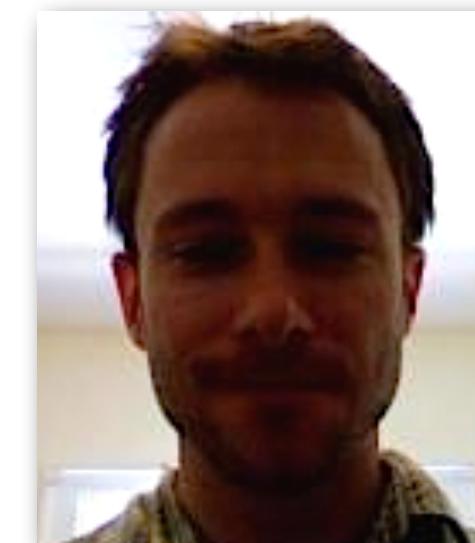
John Hutchinson



Meeta Mistry



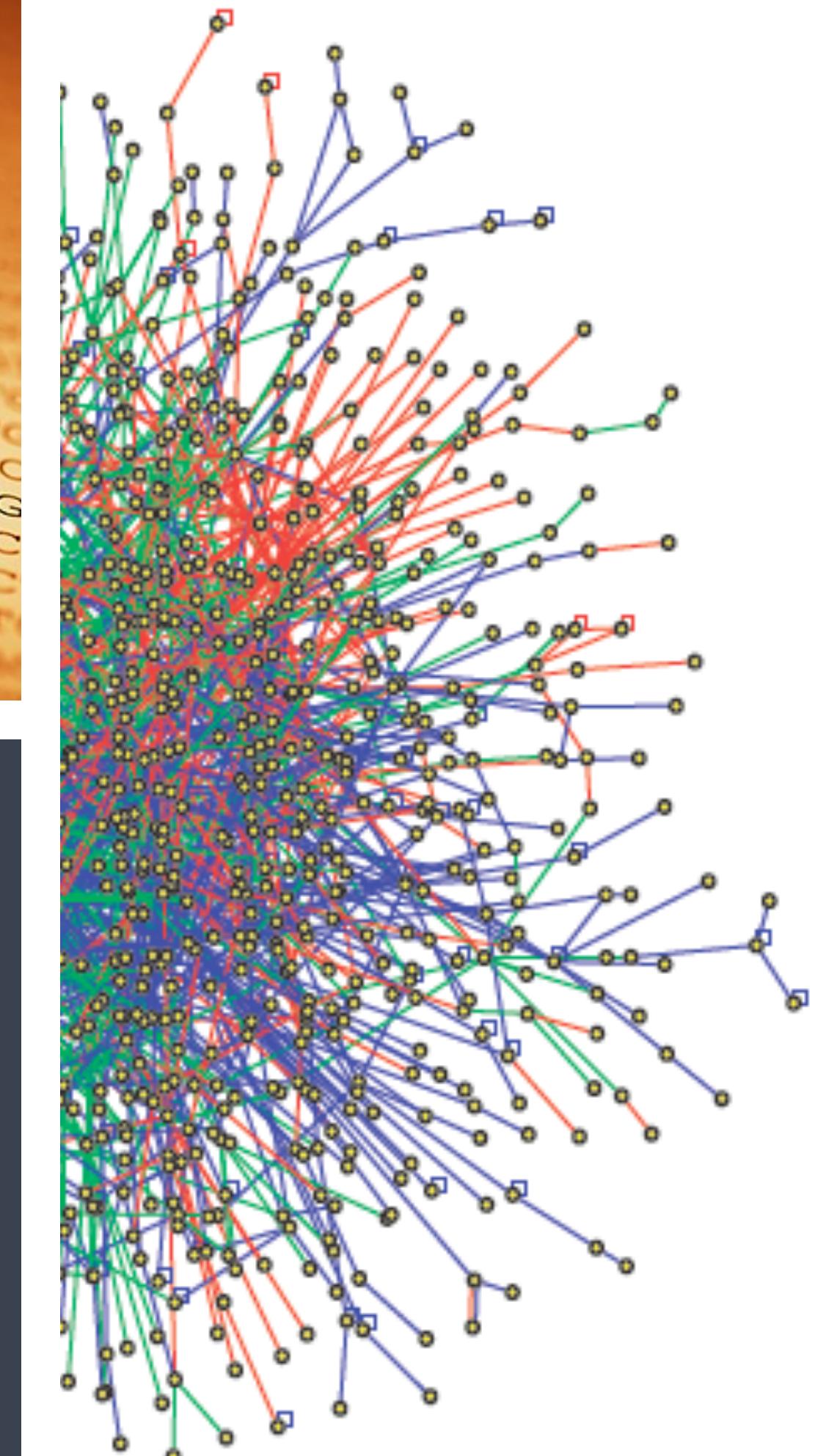
John Morrissey



Rory Kirchner

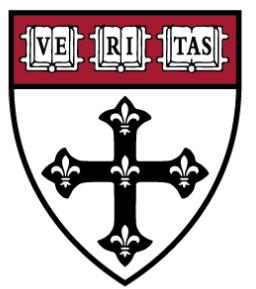


Brad Chapman



Three focus areas

- Research computing
- Next-gen sequencing
- Functional significance



HARVARD
SCHOOL OF PUBLIC HEALTH

HSCI

HARVARD STEM CELL
INSTITUTE



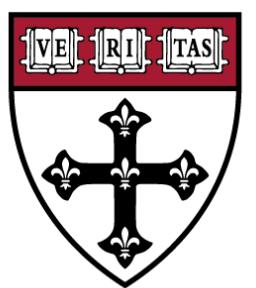
THE HARVARD CLINICAL
AND TRANSLATIONAL
SCIENCE CENTER

NIEHS / CFAR
Bioinformatics
Core

Center for
Stem Cell
Bioinformatics

Harvard
Catalyst
Bioinformatics
Consulting

more than 300
consults from
almost all
Harvard-
affiliated
institutions



HARVARD
SCHOOL OF PUBLIC HEALTH

NIEHS / CFAR
Bioinformatics
Core

HSCI
HARVARD STEM CELL
INSTITUTE

Center for
Stem Cell
Bioinformatics

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

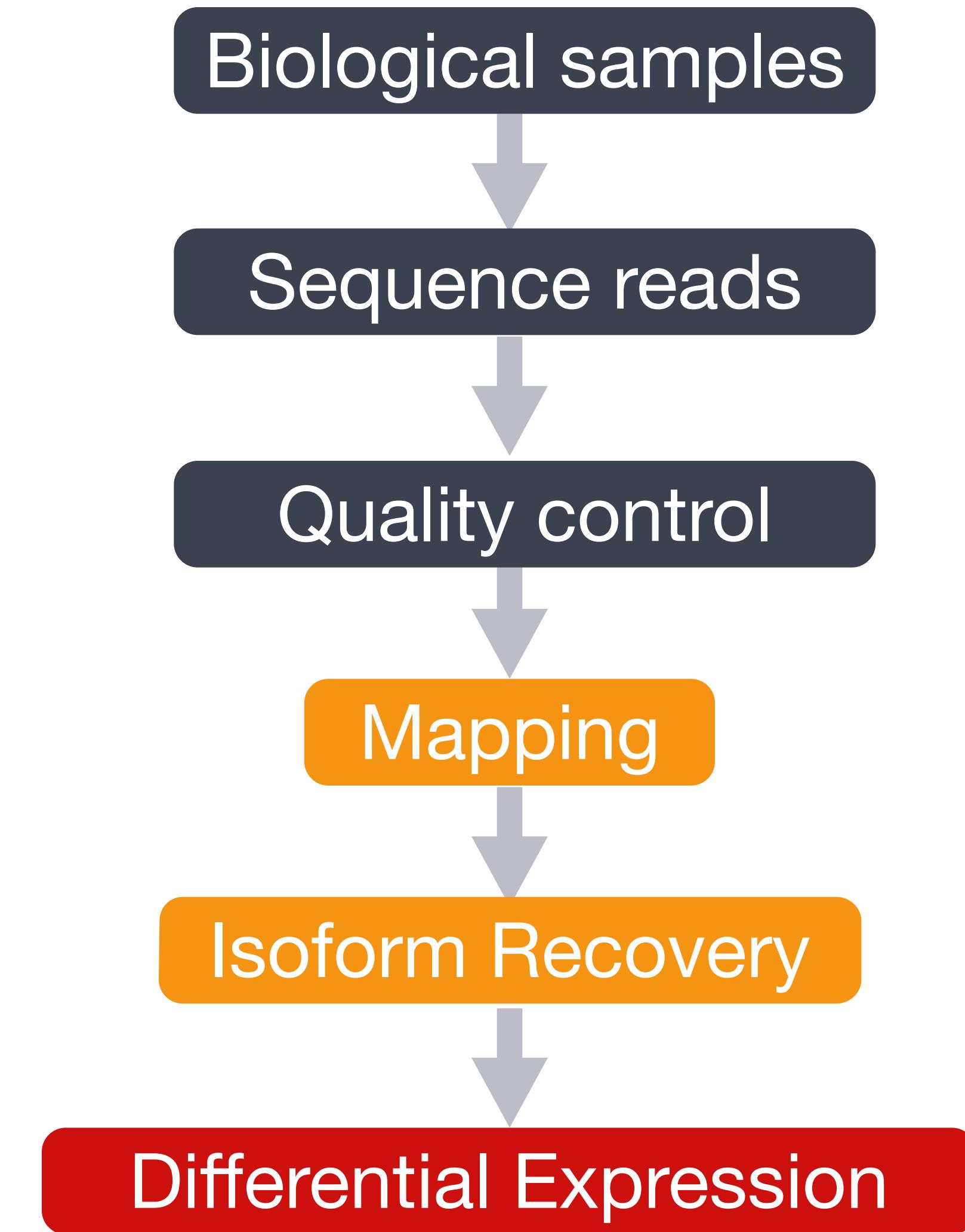
Harvard
Catalyst
Bioinformatics
Consulting


HARVARD MEDICAL SCHOOL
Tools and Technology

more than 300
consults from
almost all
Harvard-
affiliated
institutions

Scope

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

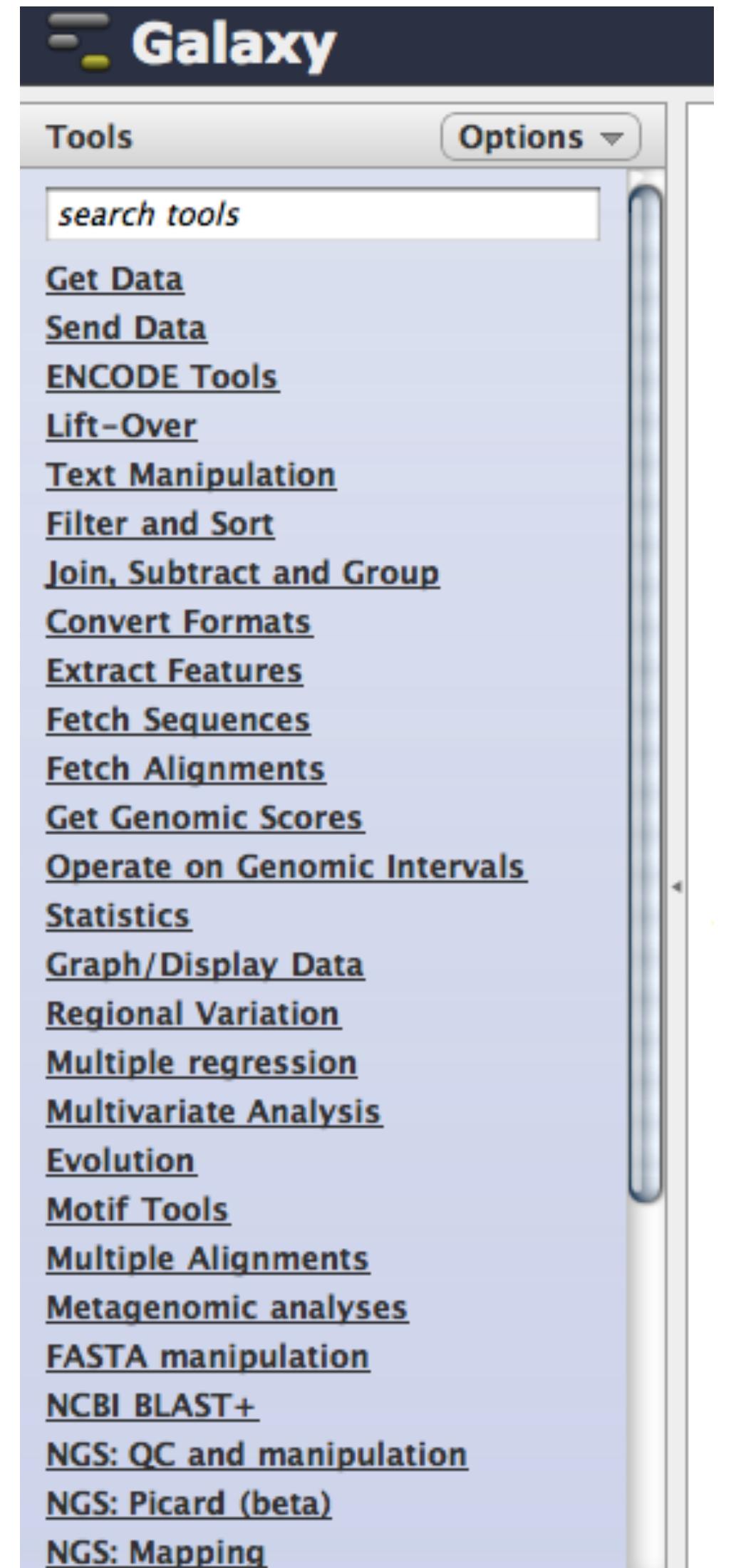


Workflows

Using Galaxy

EC2 lab environment

<http://23.23.134.25/>



Do-it yourself

No black box

More control

Better understanding of the experiment

No command line / UNIX experience required





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>

The screenshot shows a Galaxy web application interface. At the top, there is 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 displays a manuscript page. The title of the page is "Dynamics of mitochondrial heteroplasmy in three families investigated via a repeatable re-sequencing study". Below the title, the authors listed are Hiroki Goto¹, Benjamin Dickins², Enis Afgan^{3,5}, Ian M. Paul⁴, James Taylor^{3,5}, Kateryna D. Makova¹, and Anton Nekrutenko^{2,5}. It is published in Genome Biology on June 23, 2011. Correspondence should be addressed to [KDM](#), [JT](#), or [AN](#). The first section, "1. How to use this document", explains that this is a live copy of supplementary materials for the manuscript and provides access to data and workflows. It lists links for datasets, workflows, and histories. It also mentions two longer screencasts. The second section, "2. Accessing the Data", discusses datasets found in a Galaxy library and an S3 bucket. A note at the bottom explains the naming convention for datasets.

Dynamics of mitochondrial heteroplasmy in three families investigated via a repeatable re-sequencing study

Hiroki Goto¹, Benjamin Dickins², Enis Afgan^{3,5}, Ian M. Paul⁴, James Taylor^{3,5}, Kateryna D. Makova¹, and Anton Nekrutenko^{2,5}

Published in Genome Biology on June 23, 2011

Correspondence should be addressed to [KDM](#), [JT](#), or [AN](#).

1. How to use this document

This document is a live copy of supplementary materials for [the manuscript](#). It provides access to all the data as well as to exact analyses and workflows discussed in the paper, so you can play with them by re-running, changing parameters, or even applying them to your own sequencing data. To import workflows you must [create a Galaxy account](#) (unless you already have one) – a hassle-free procedure where you are only asked for a username and password. To make this even easier, we created several screencasts (very short movies) to help you:

- [access our datasets](#)
- [re-use workflows listed on this page](#)
- [view and import histories listed on this page](#)

In addition, we created two longer screencasts:

- [Watch the analysis of one family \(F7\) from start \(Illumina reads\) to finish \(a list of variable position\);](#)
- [Watch how the complete analysis can be performed on the Amazon Cloud.](#)

If you experience any problems while using this page, please e-mail our [bug report list](#) and we will get back to you.

2. Accessing the Data

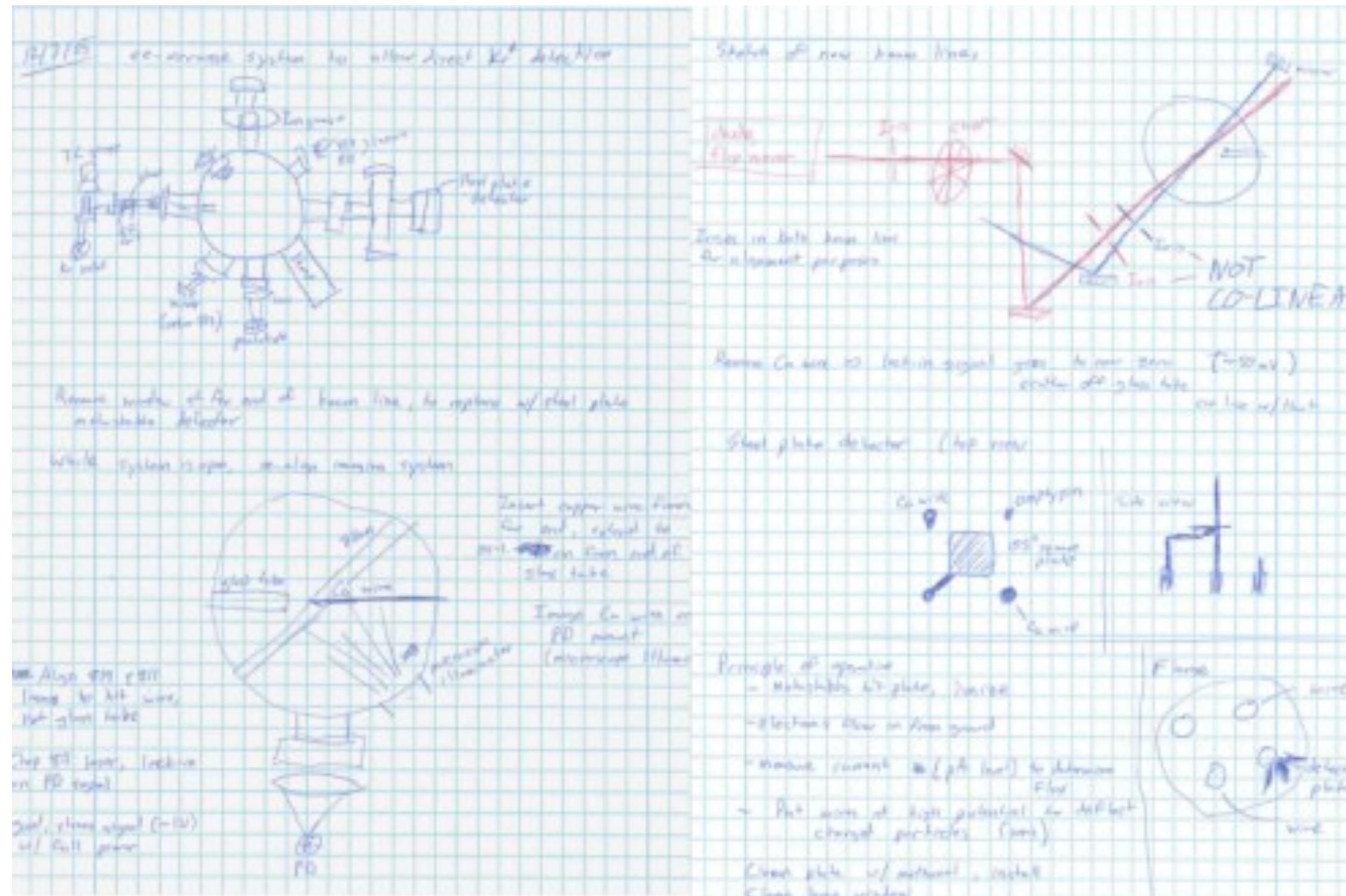
All datasets discussed in the paper can be found in two places:

- [A Galaxy Library called mtProject;](#)
- [An S3 bucket on the Amazon Cloud.](#)

From there these datasets can either be downloaded or re-analyzed with Galaxy as [described here](#). The name of each dataset is formatted as [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. The relationship among

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

Don't be afraid



Test... 2... 3

Course in pilot phase



Collaborate

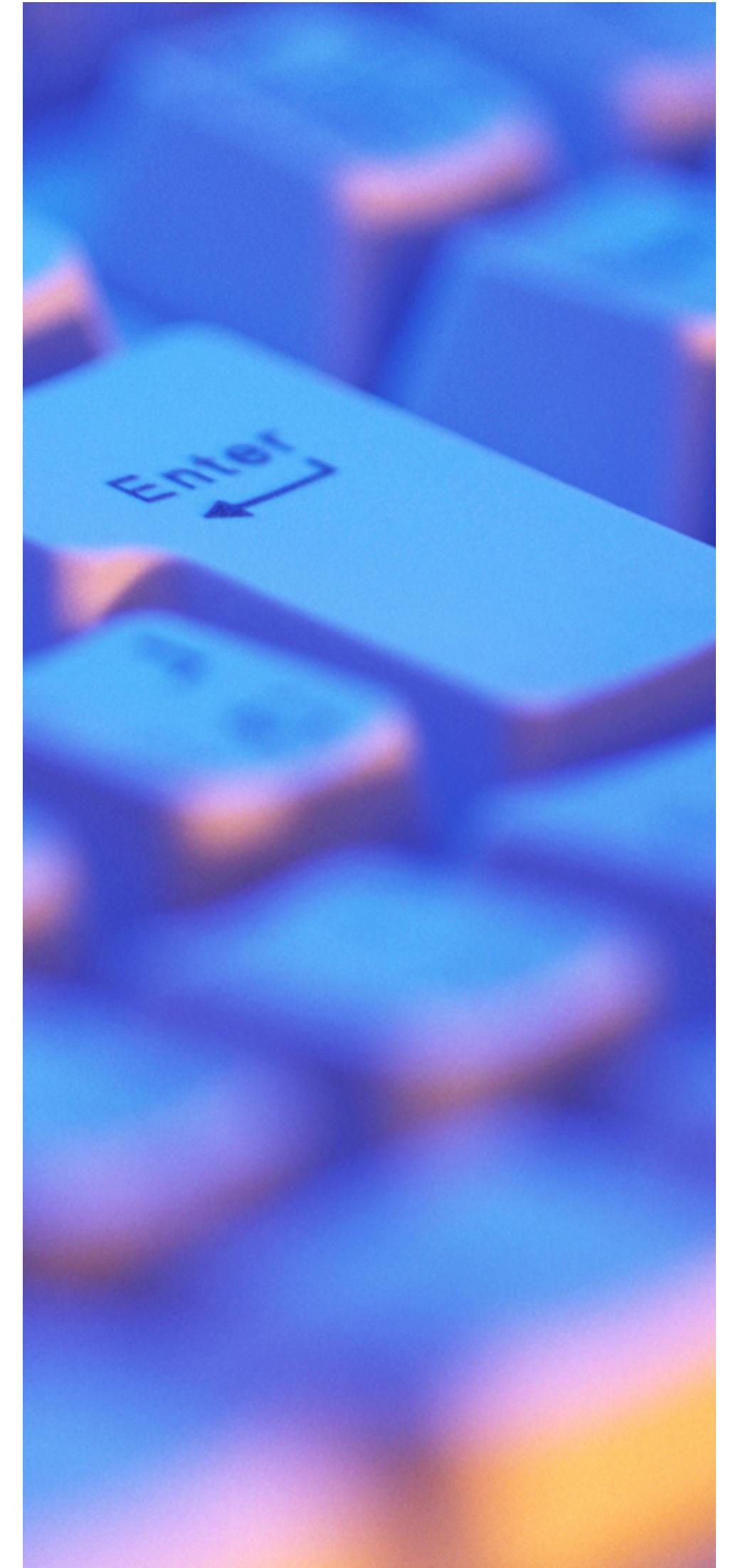
Feedback

- ▶ Survey
- ▶ Minute Cards



More information

See the course website for pointers



Contact

ohofmann@hsph.harvard.edu

@fiamh



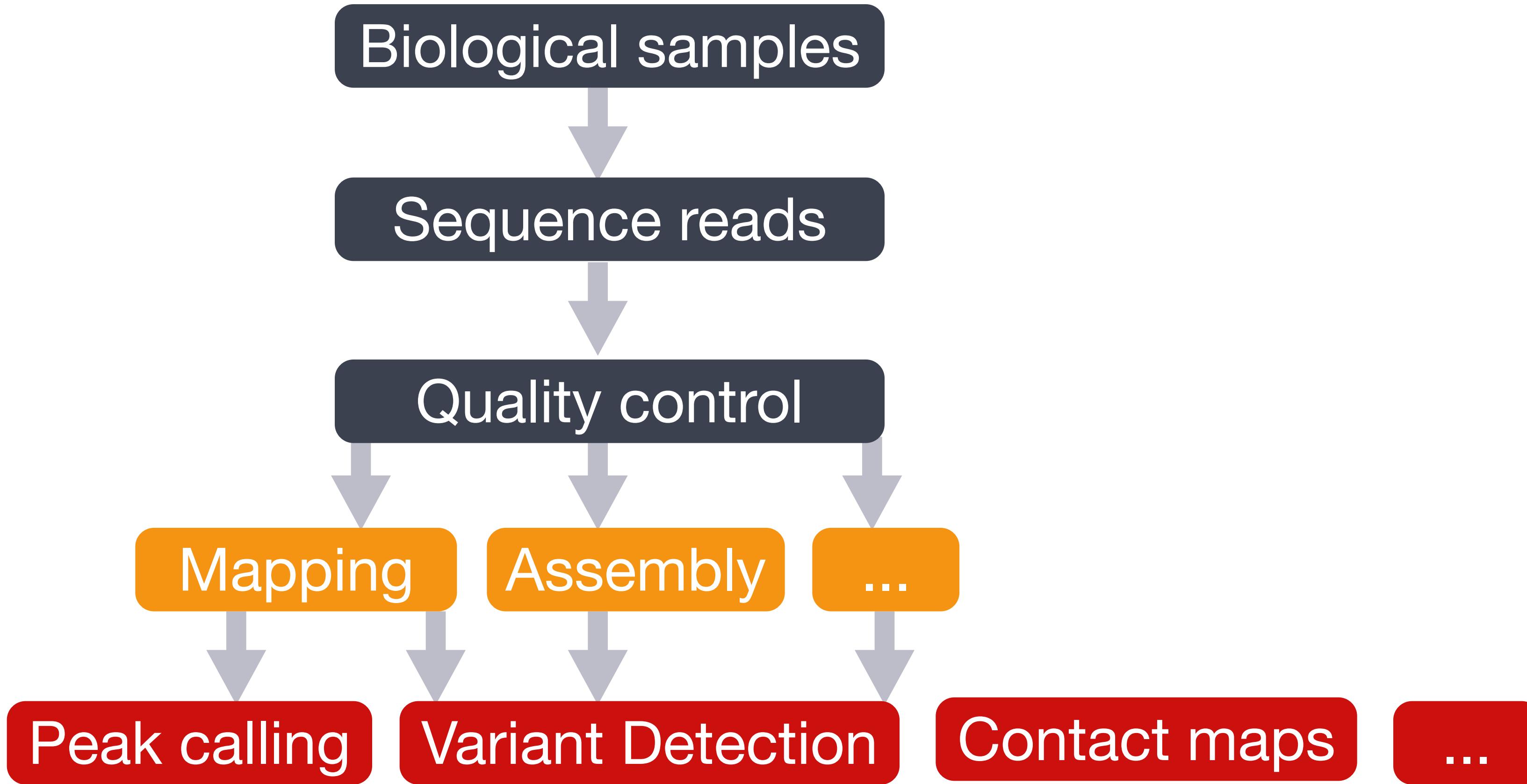
</intro>

Need for standards

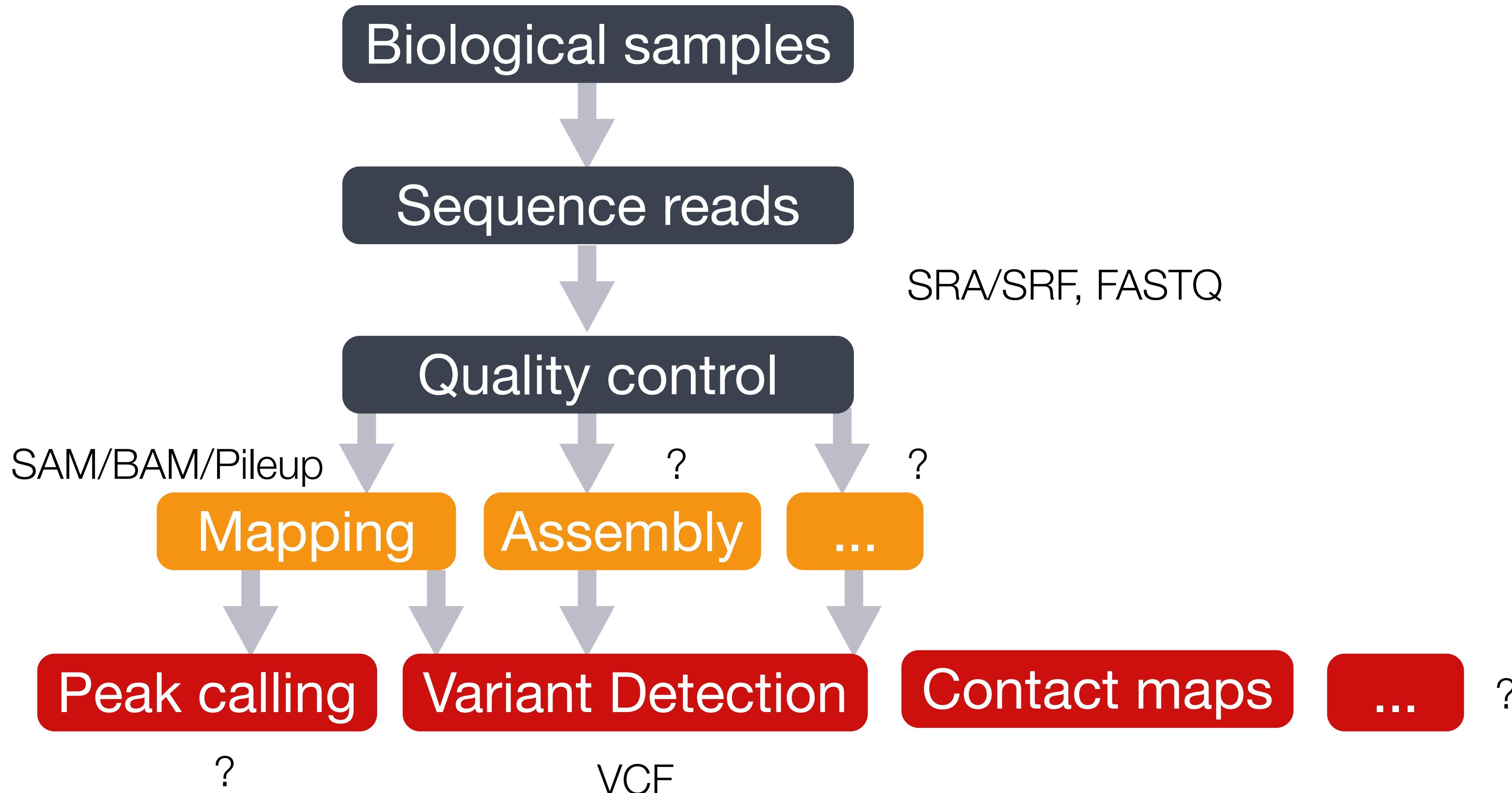
Plug and play: modular approach to tools

Vital factor in application acceptance





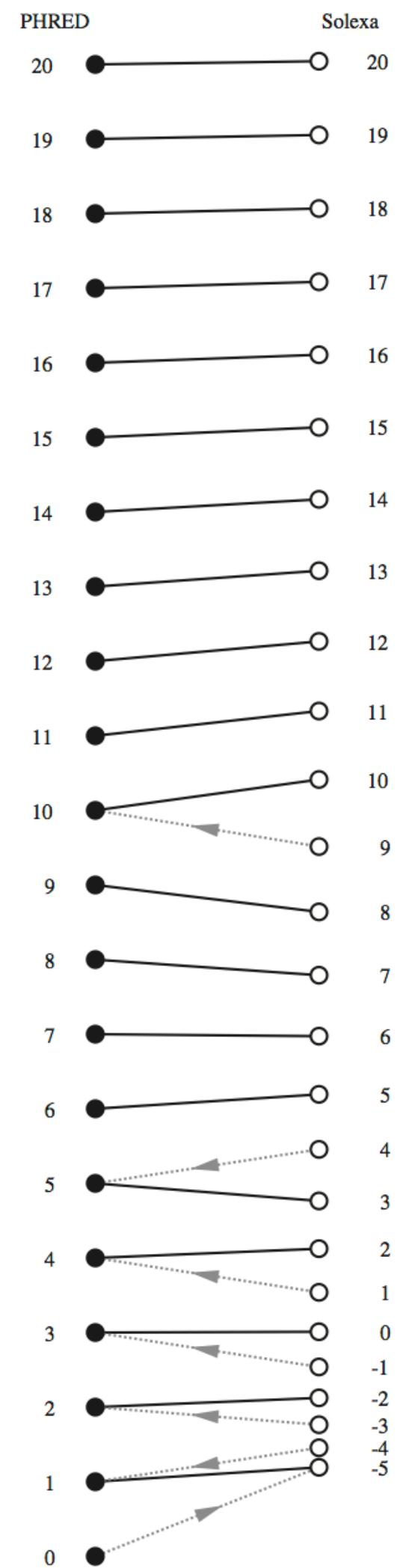
Standards



FASTQ: a “standard”

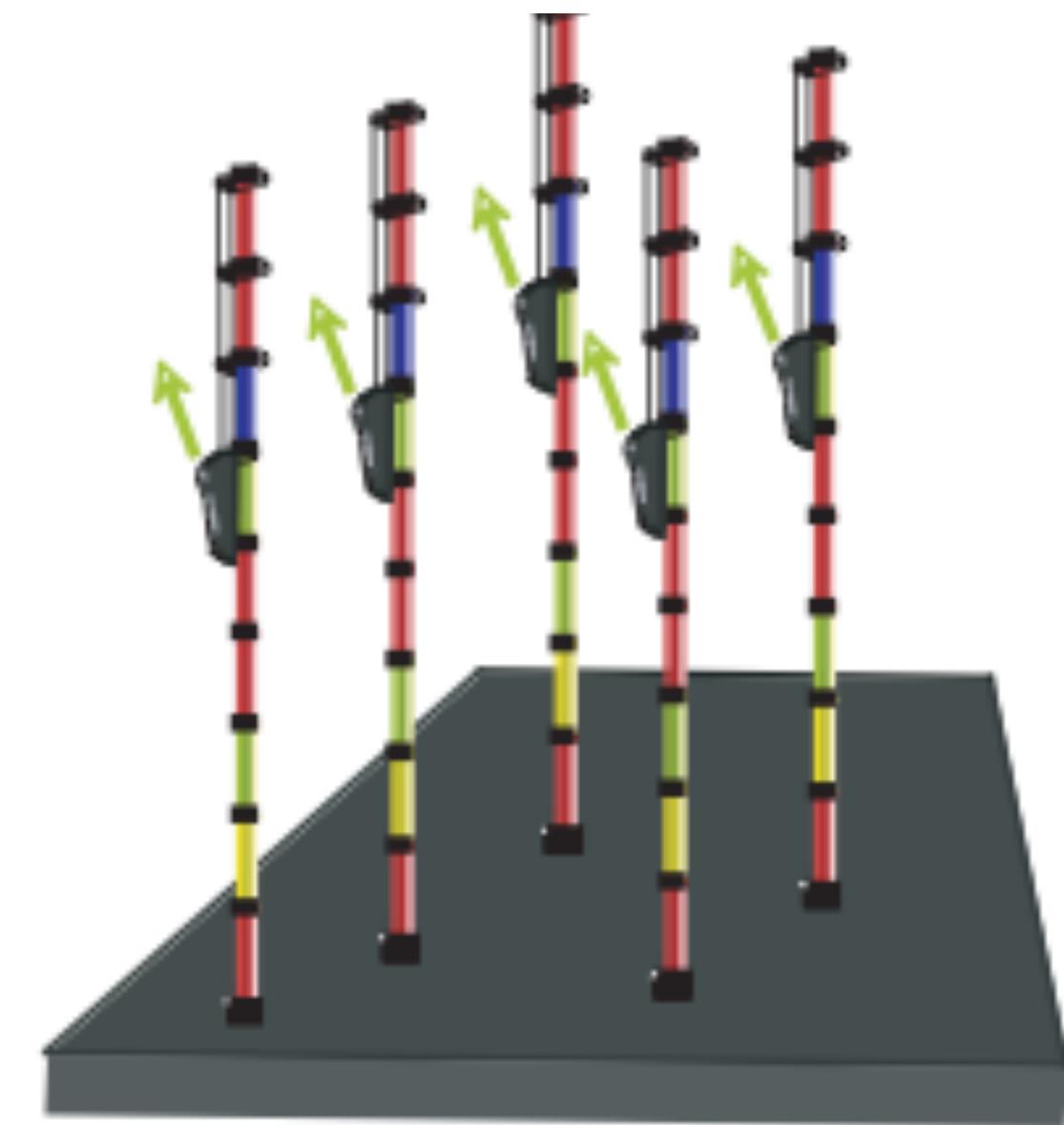
Sanger FASTQ, Solexa FASTQ, ABI Colour Space FASTQ, ...

```
@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==
```

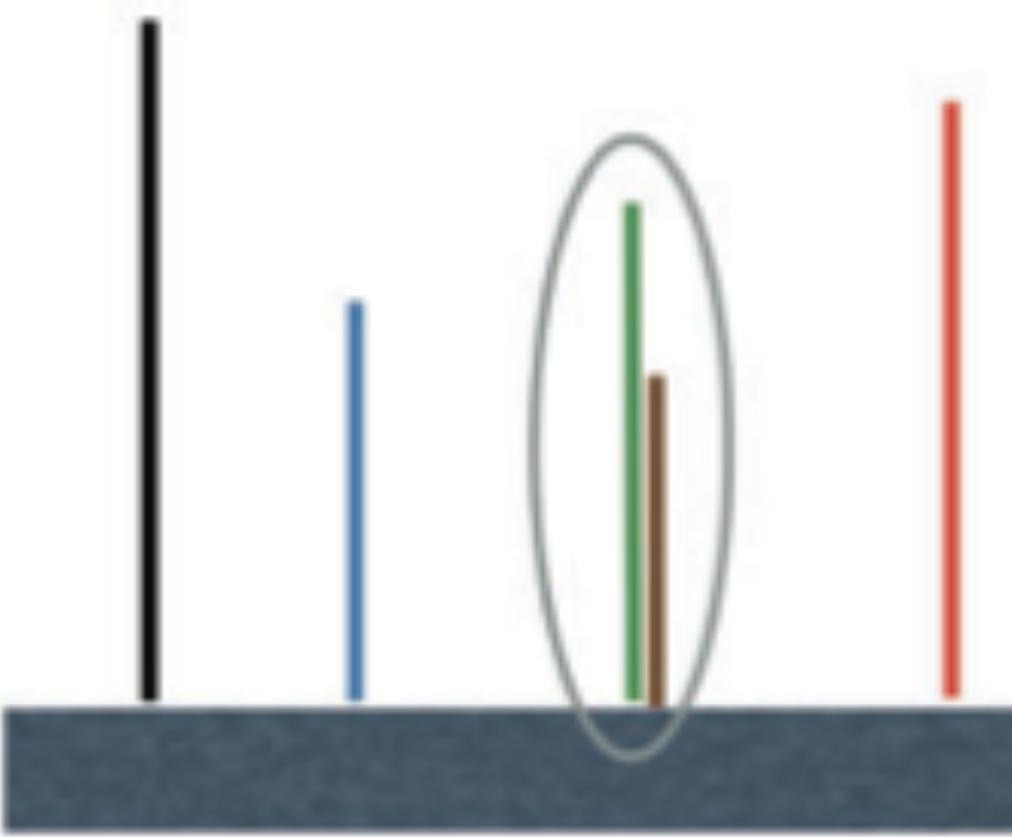


Error profiles

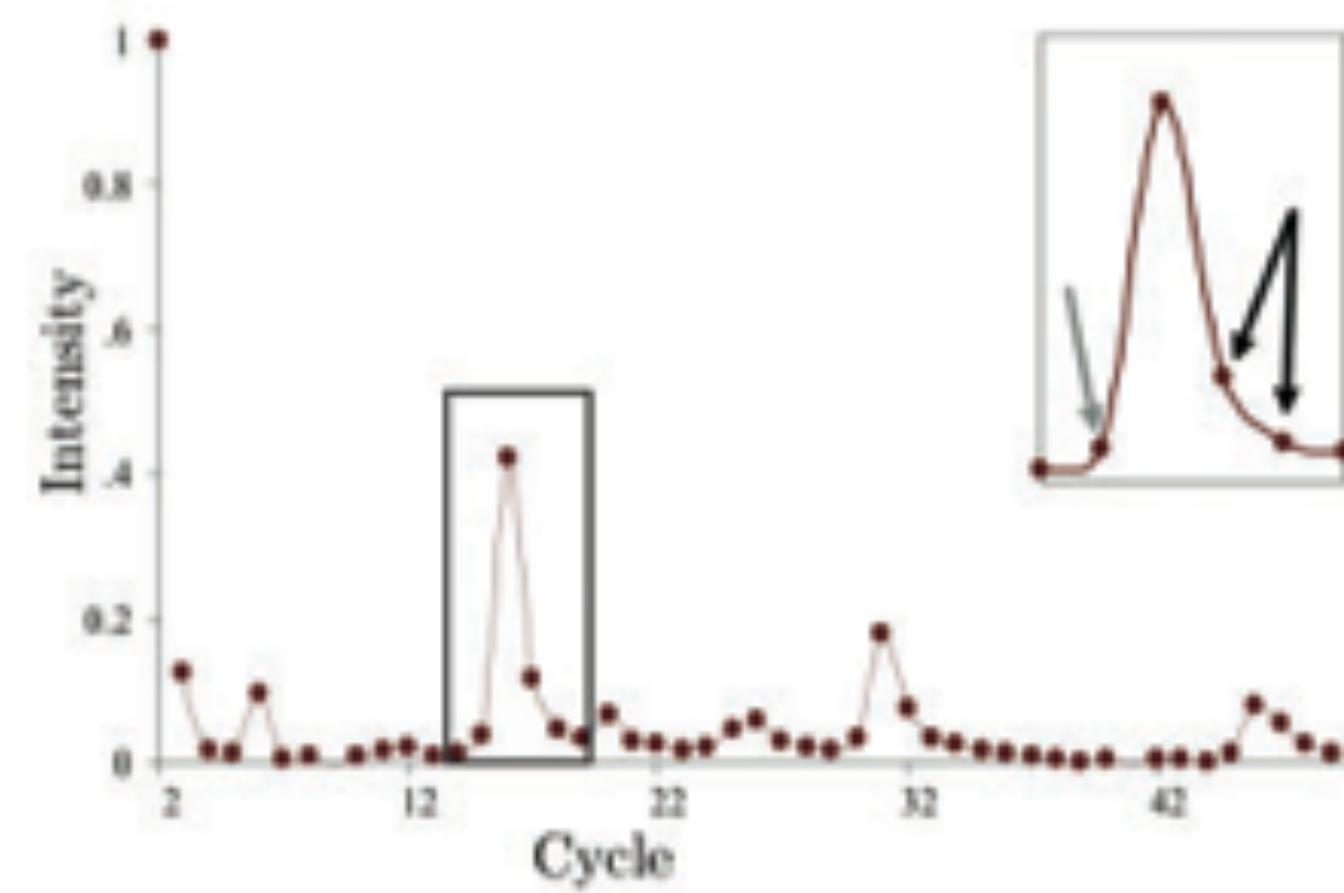
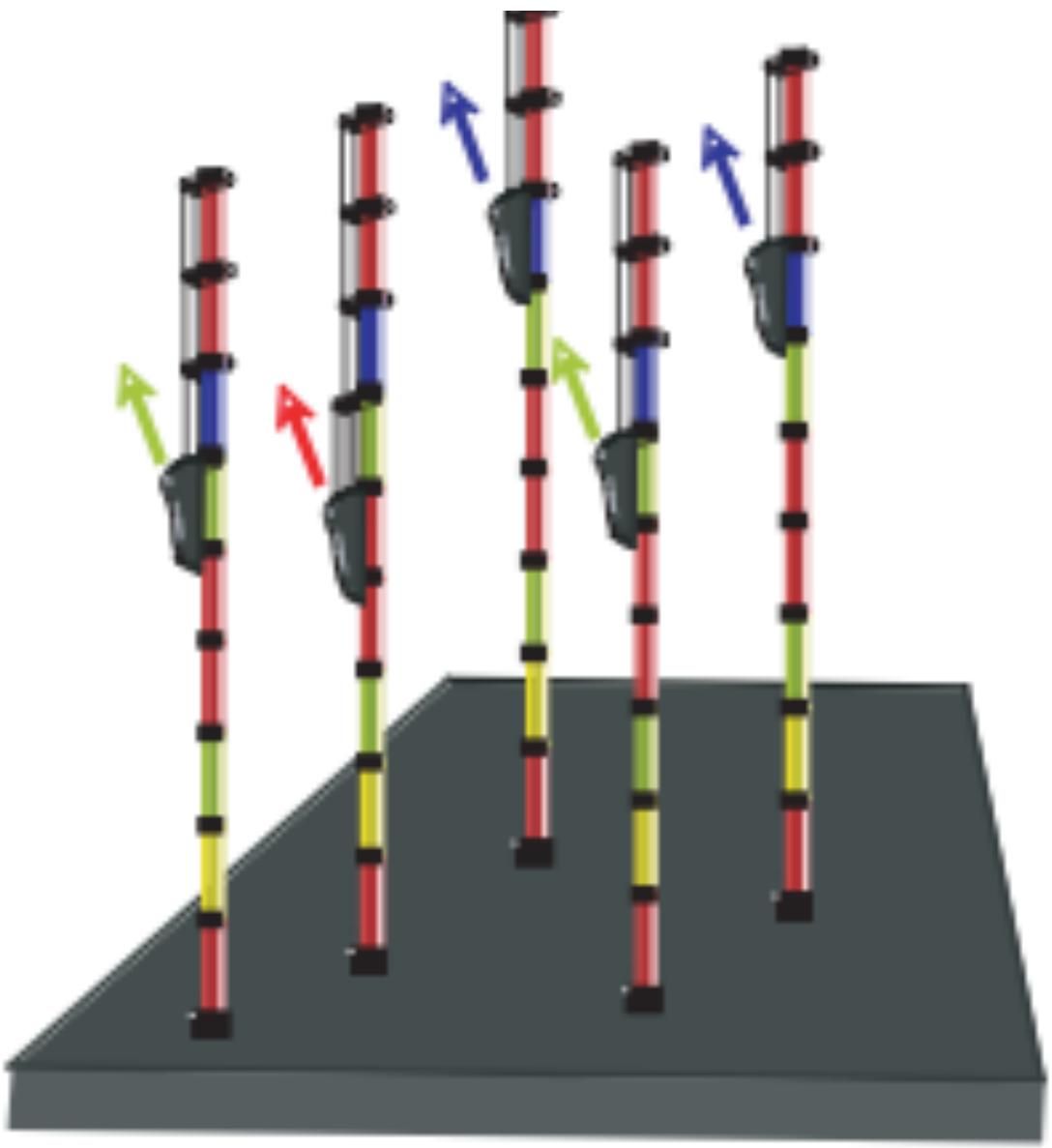
- ▶ PCR artifacts
- ▶ Error dependency on technology



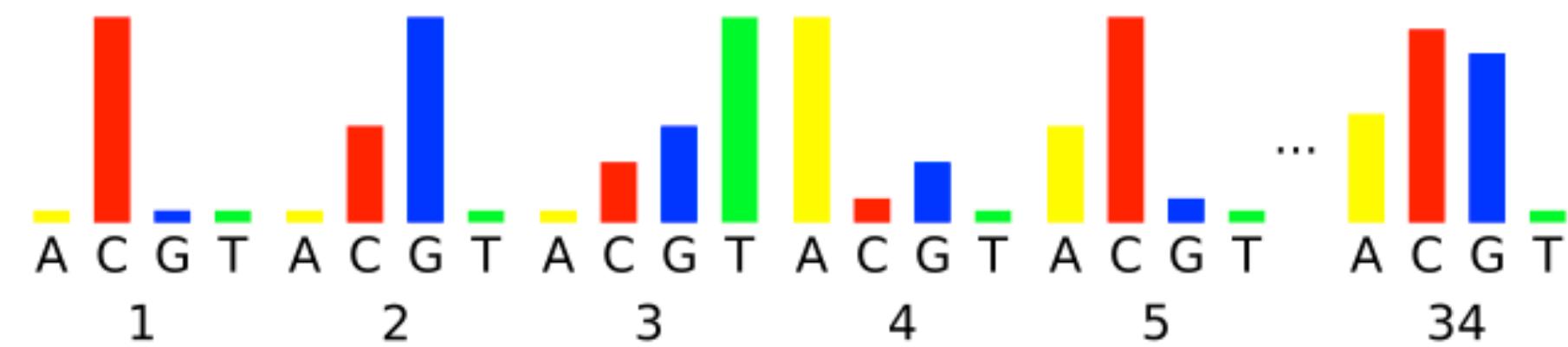
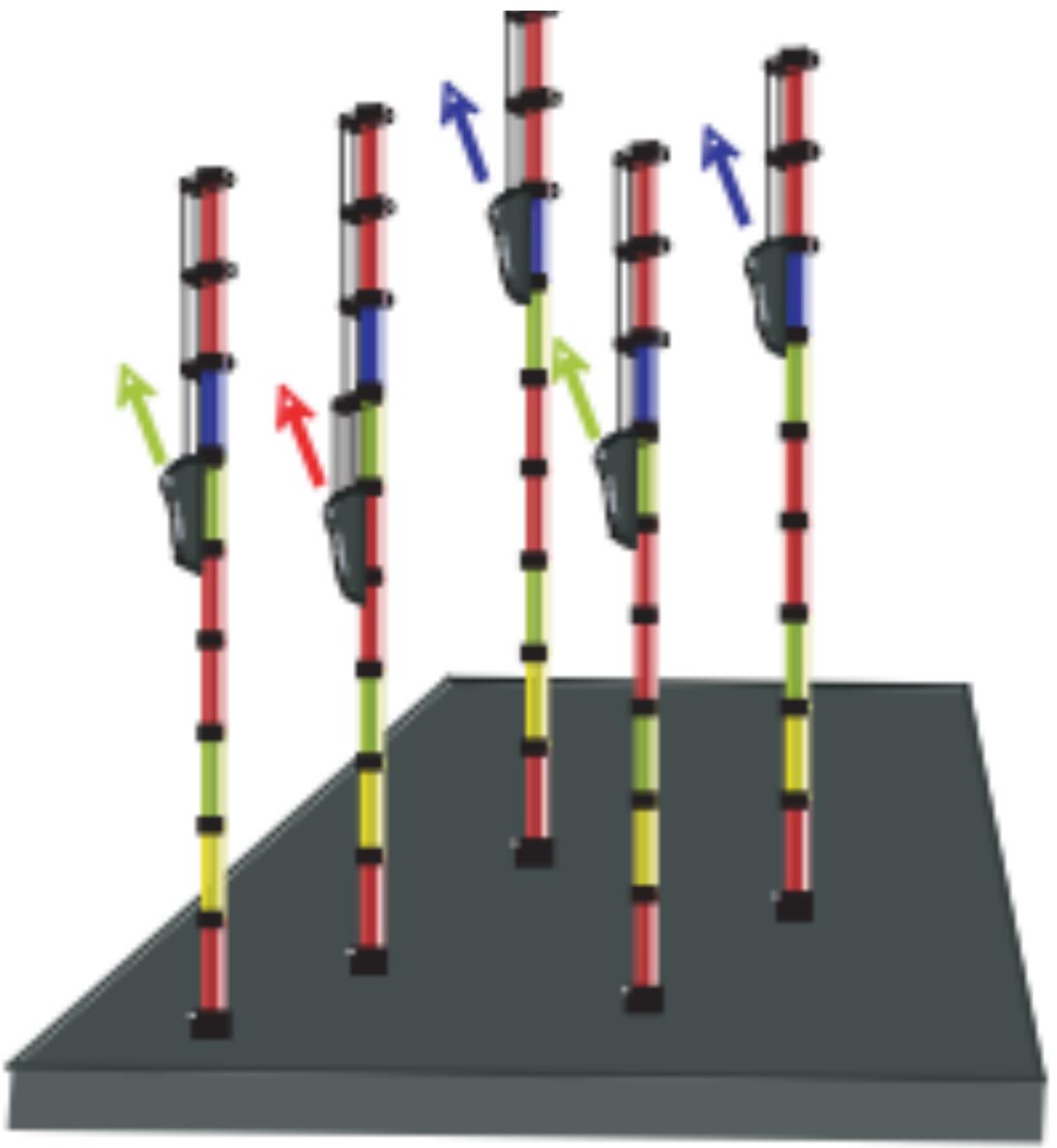
Illumina



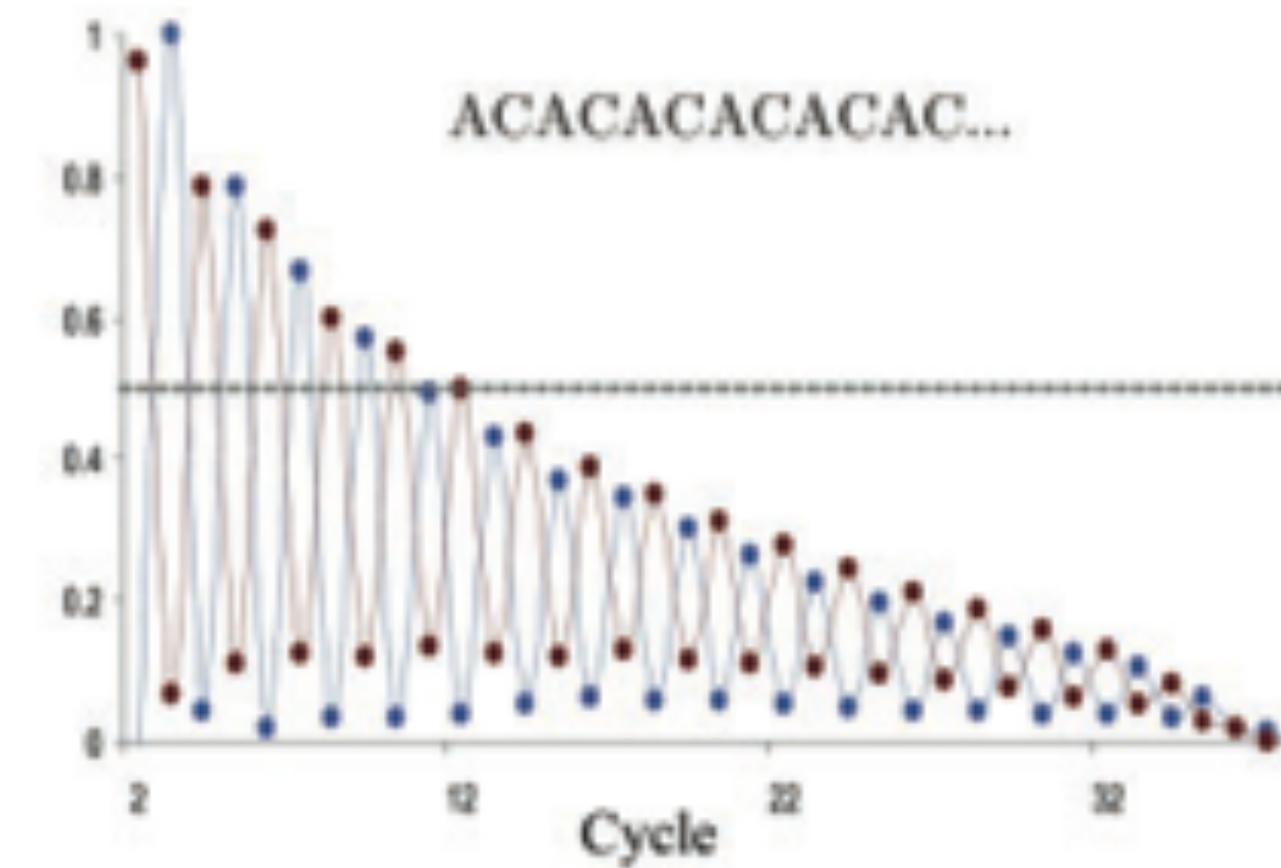
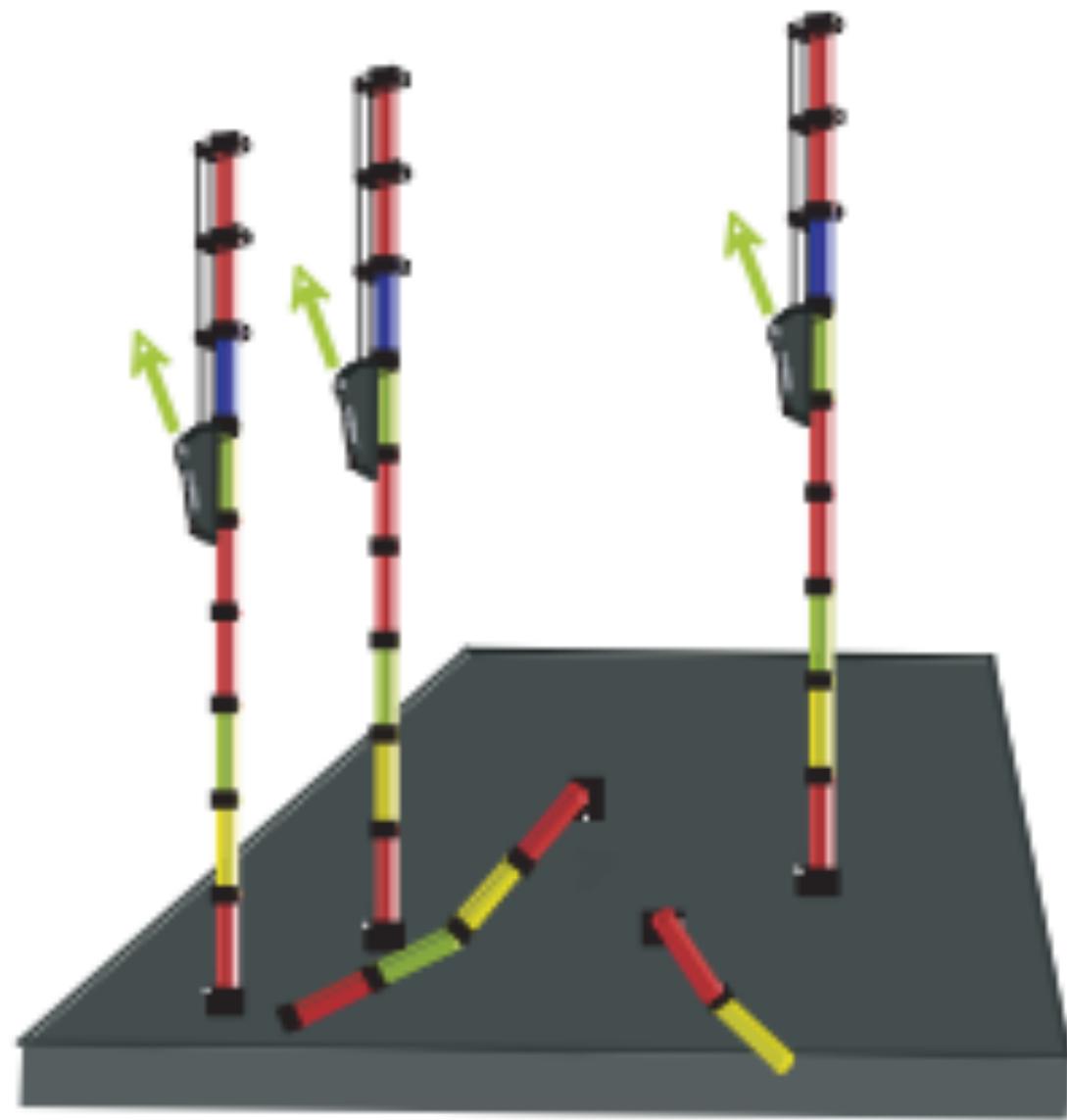
Illumina: mixed clusters



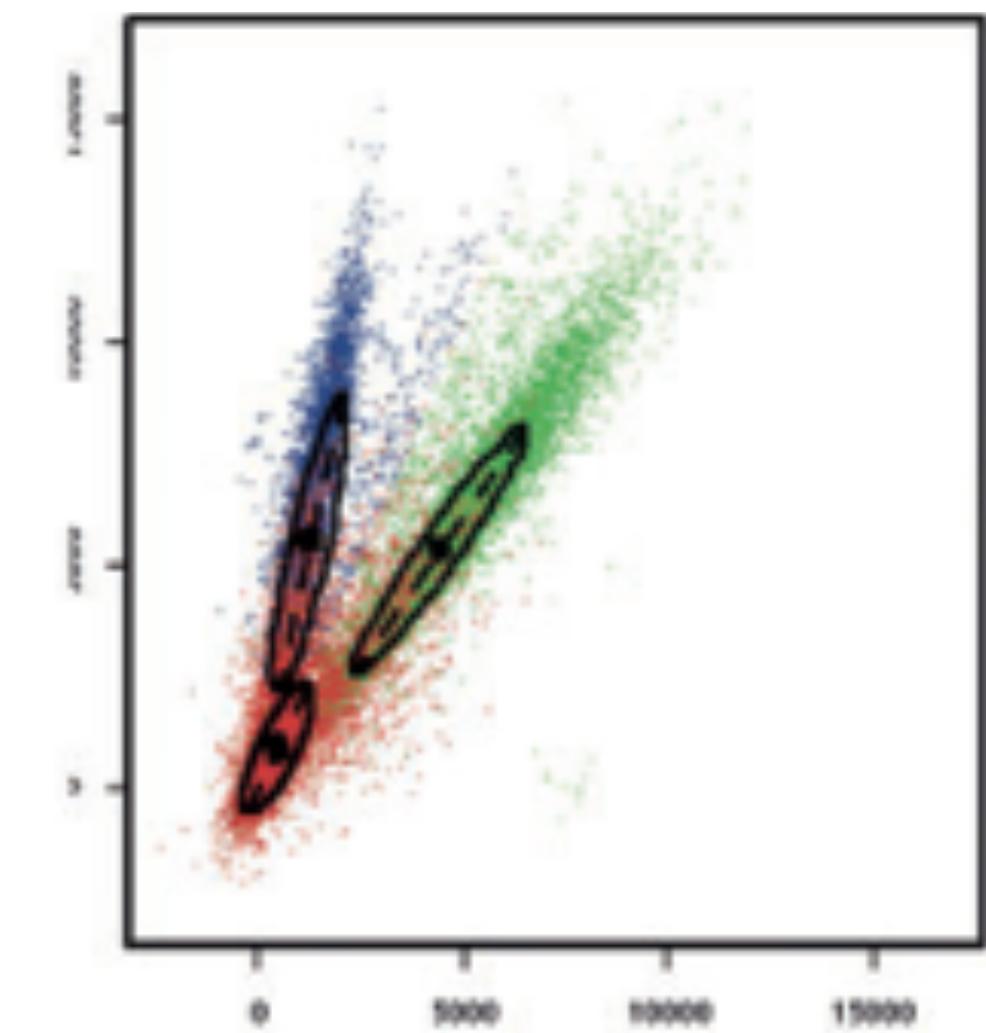
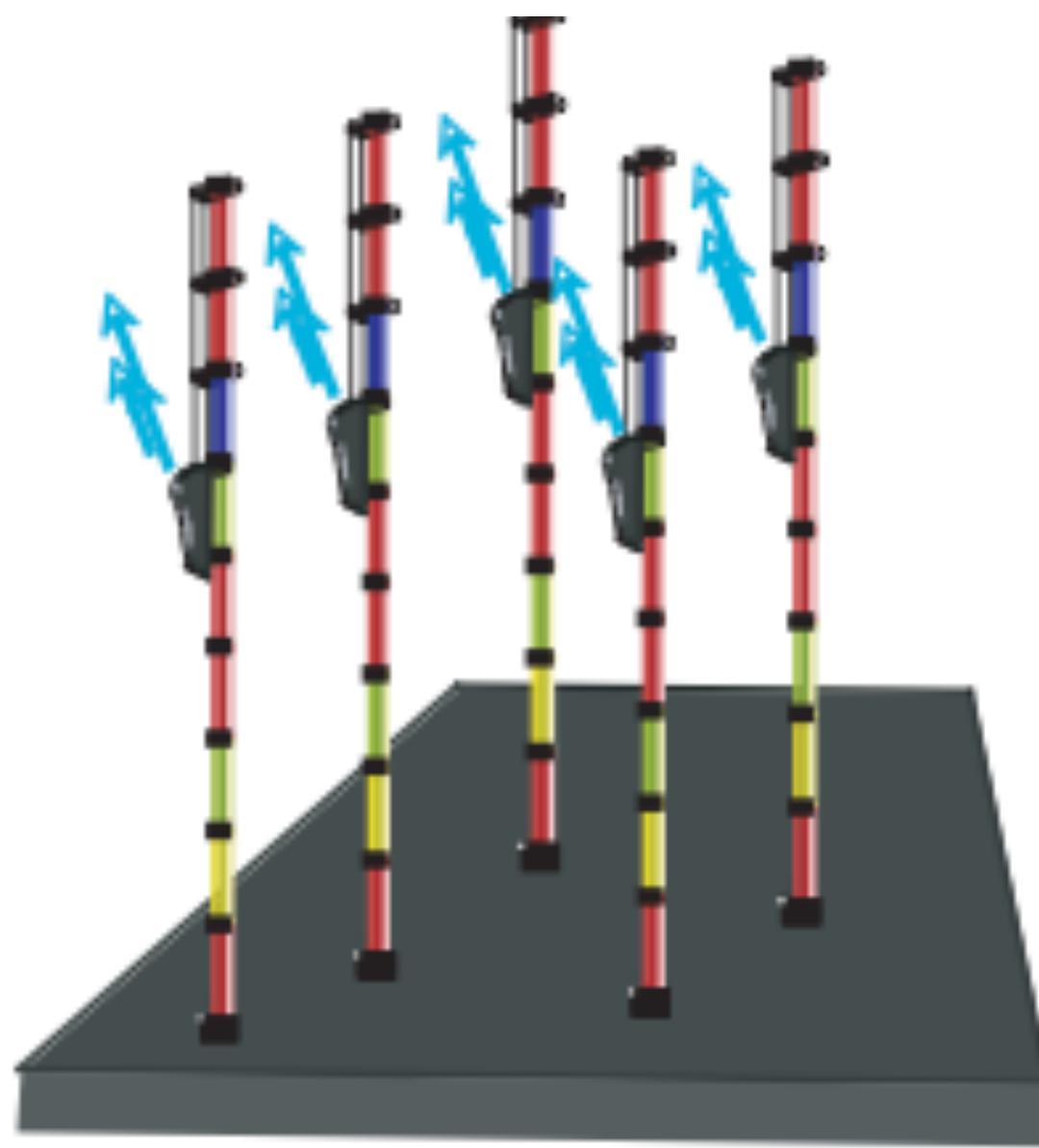
Illumina: phasing



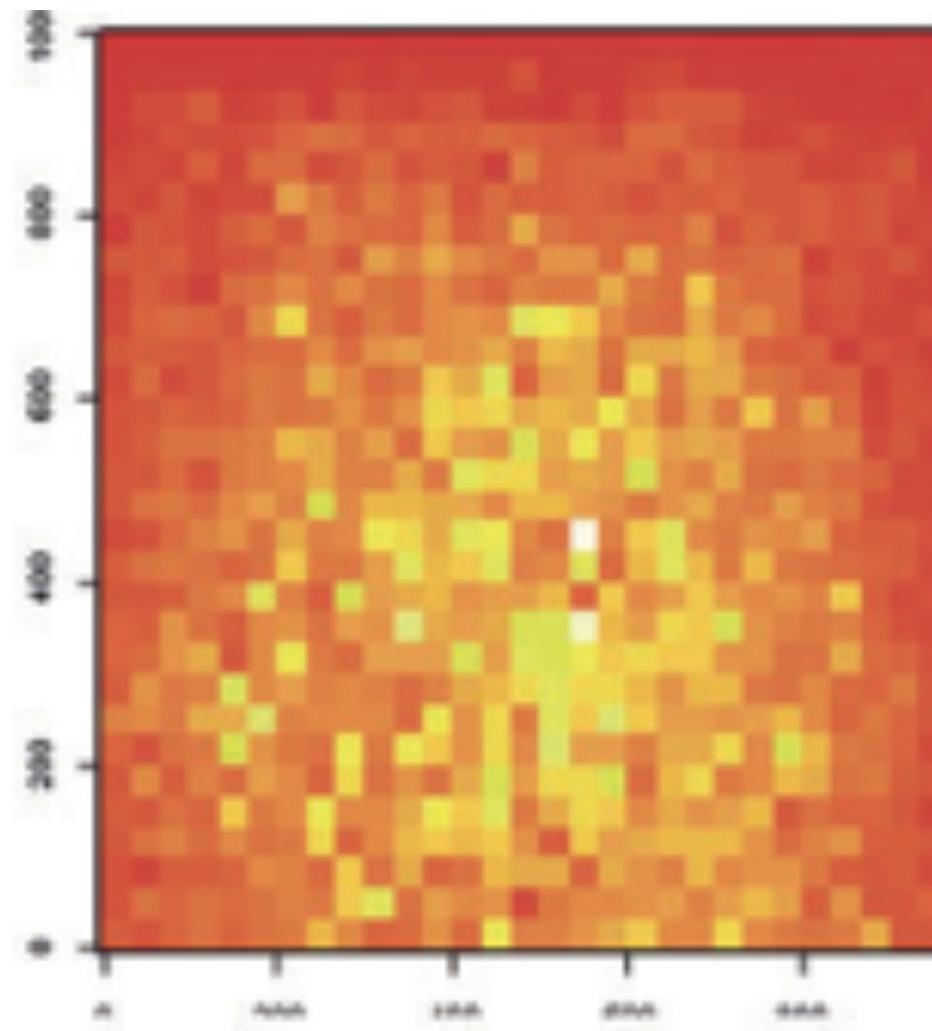
Illumina: phasing



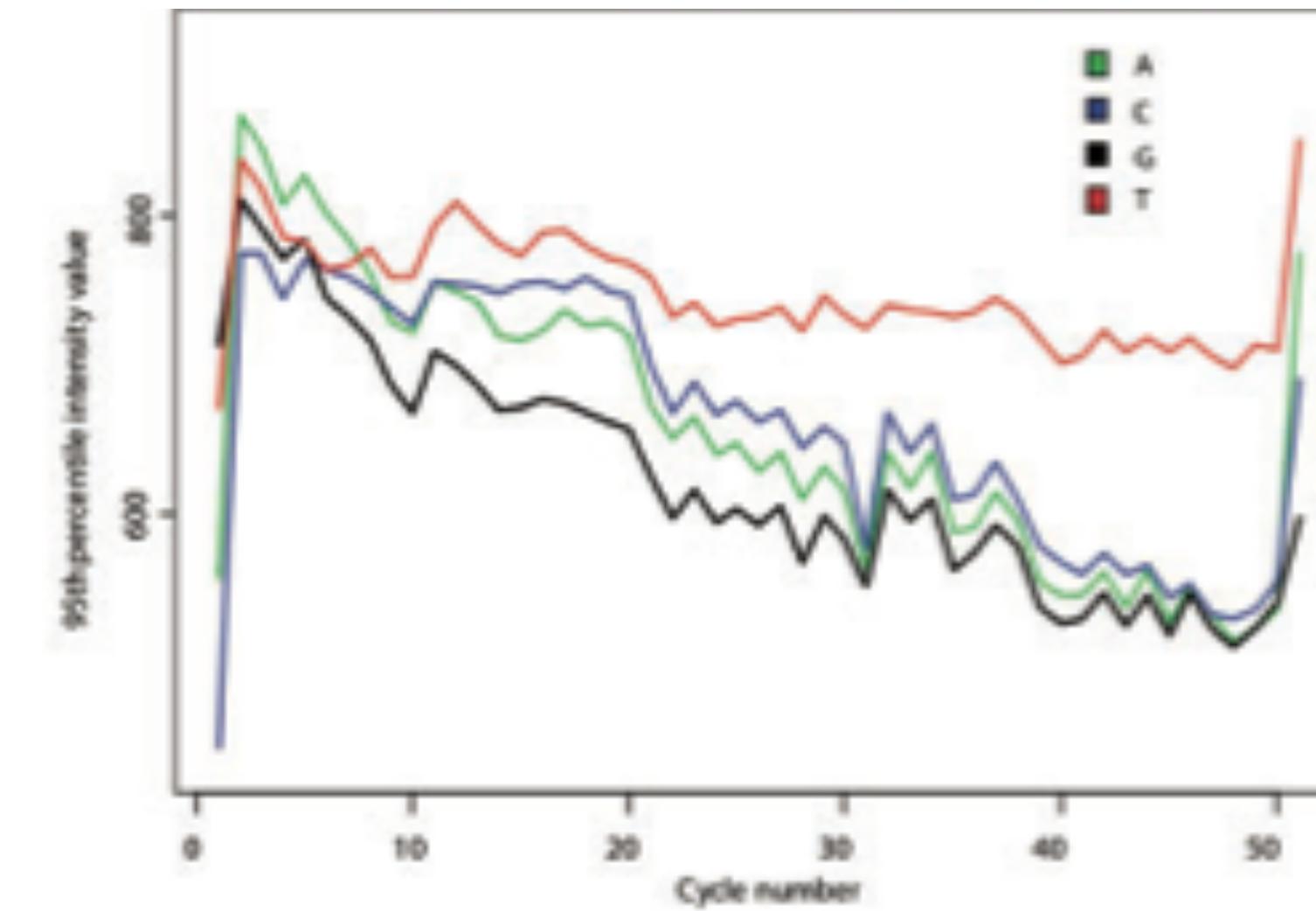
Illumina: signal decay



Illumina: cross-talk

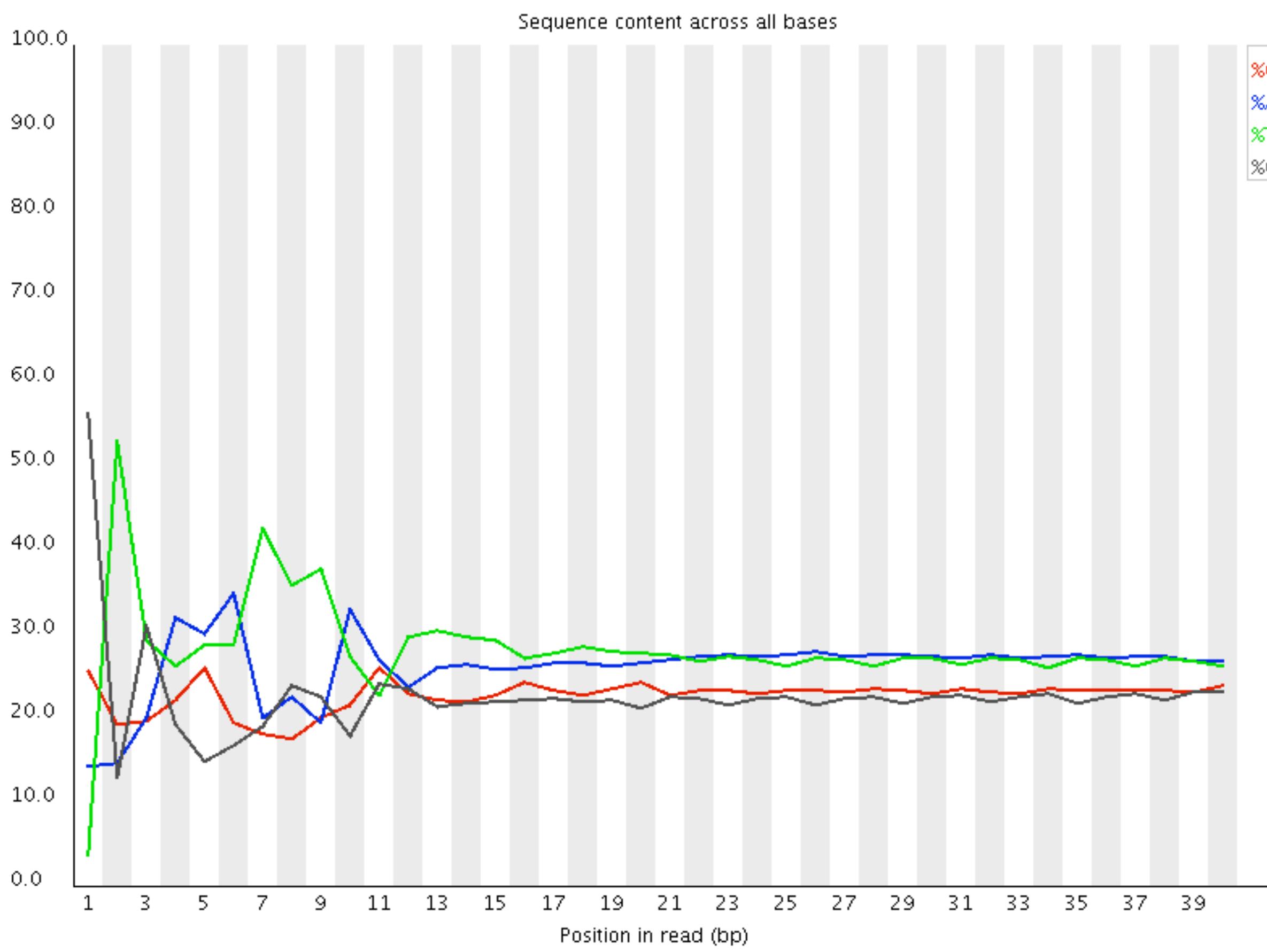


Boundary effects



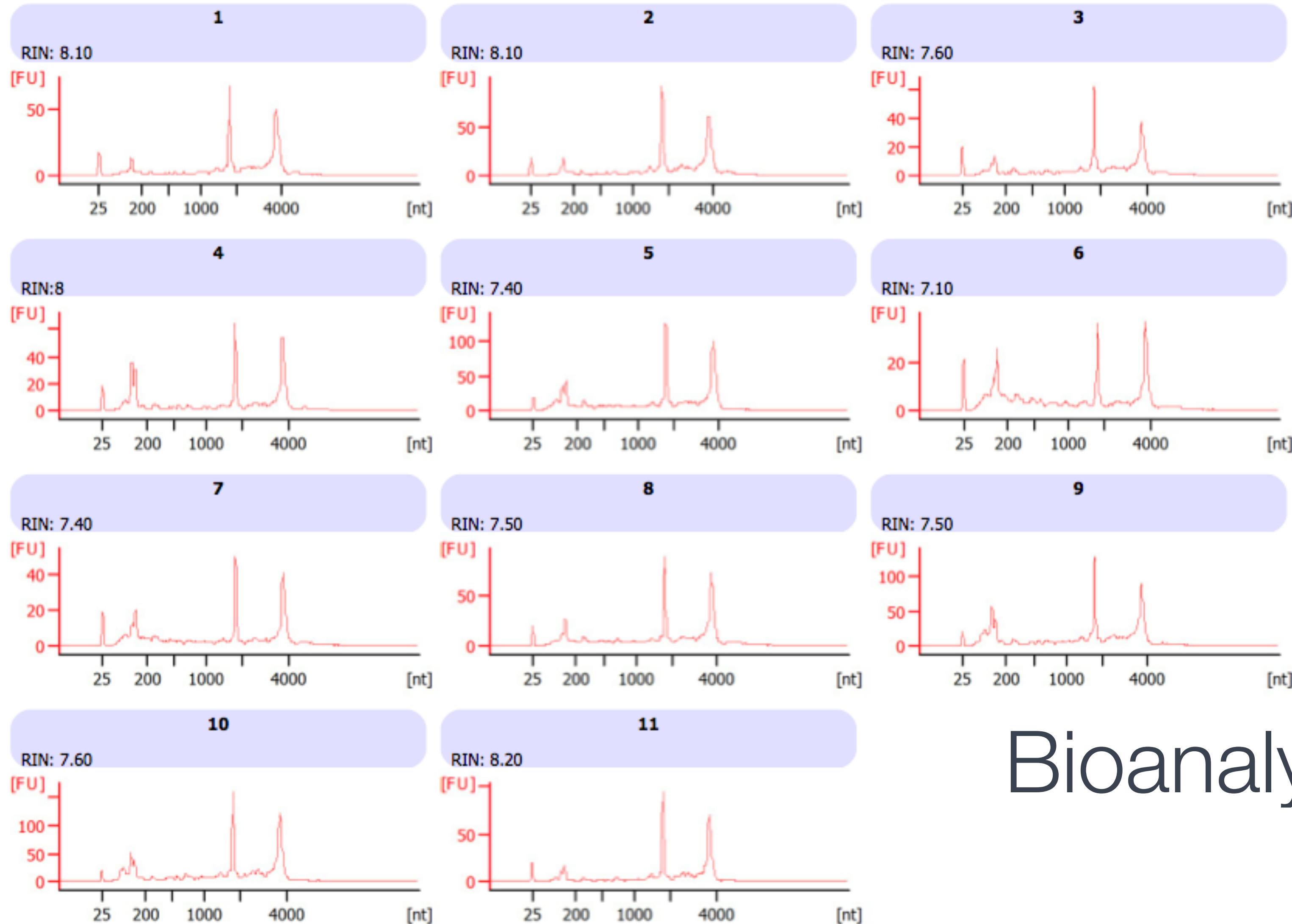
Fluorophore accumulation

Illumina: physical/chemical problems



Q&A

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



Bioanalyzer

	sequence	count
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

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

QA: filtering

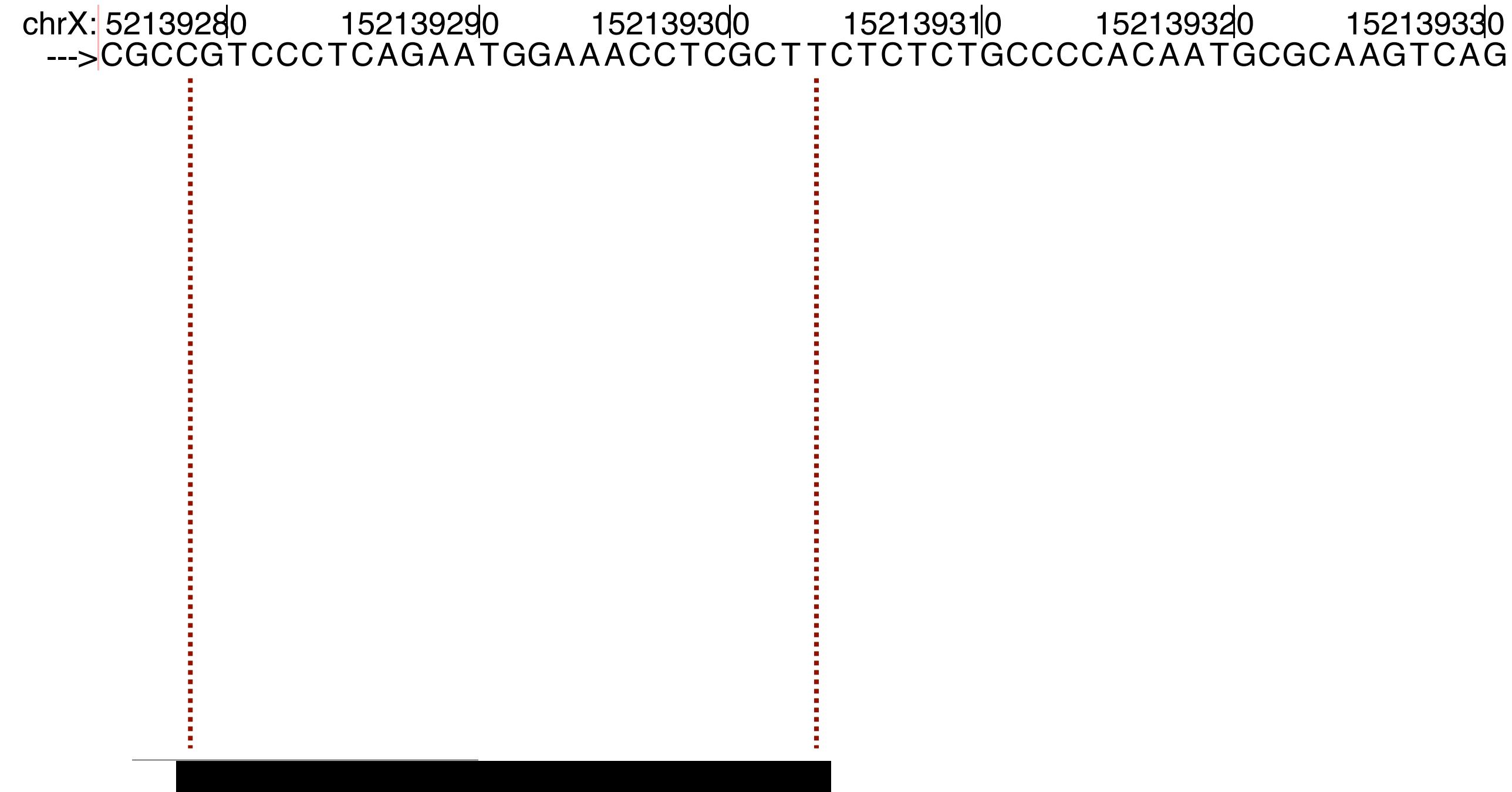
chrX: 52139280 152139290 152139300 152139310 152139320 152139330
---> CGCCGTCCCTCAGAATGGAAACCTCGCTTCTCTGCCCAATGCGCAAGTCAG

Genome

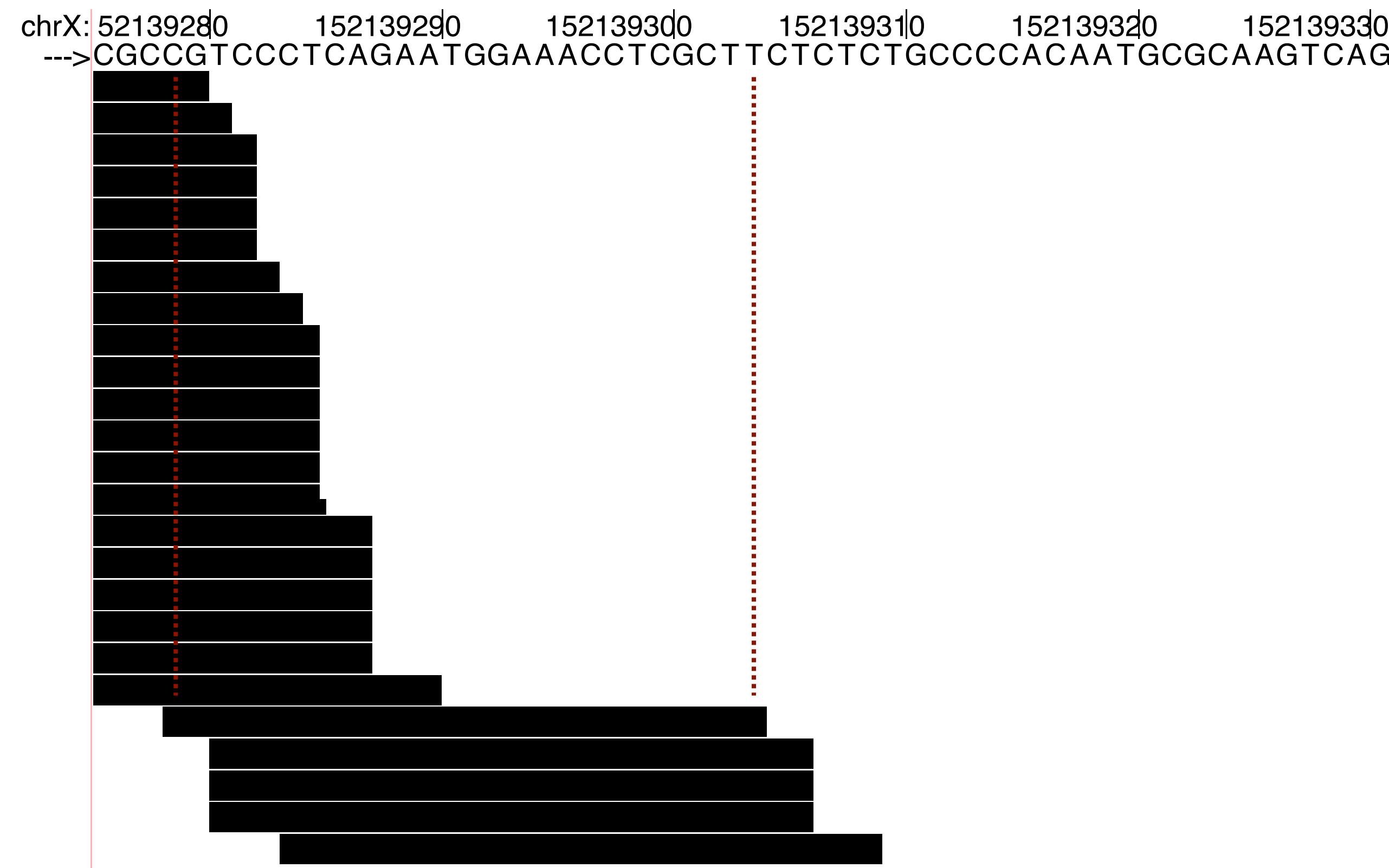
CGTCCCTCAGAATGGAAACCTCGCTT

Sequence tag

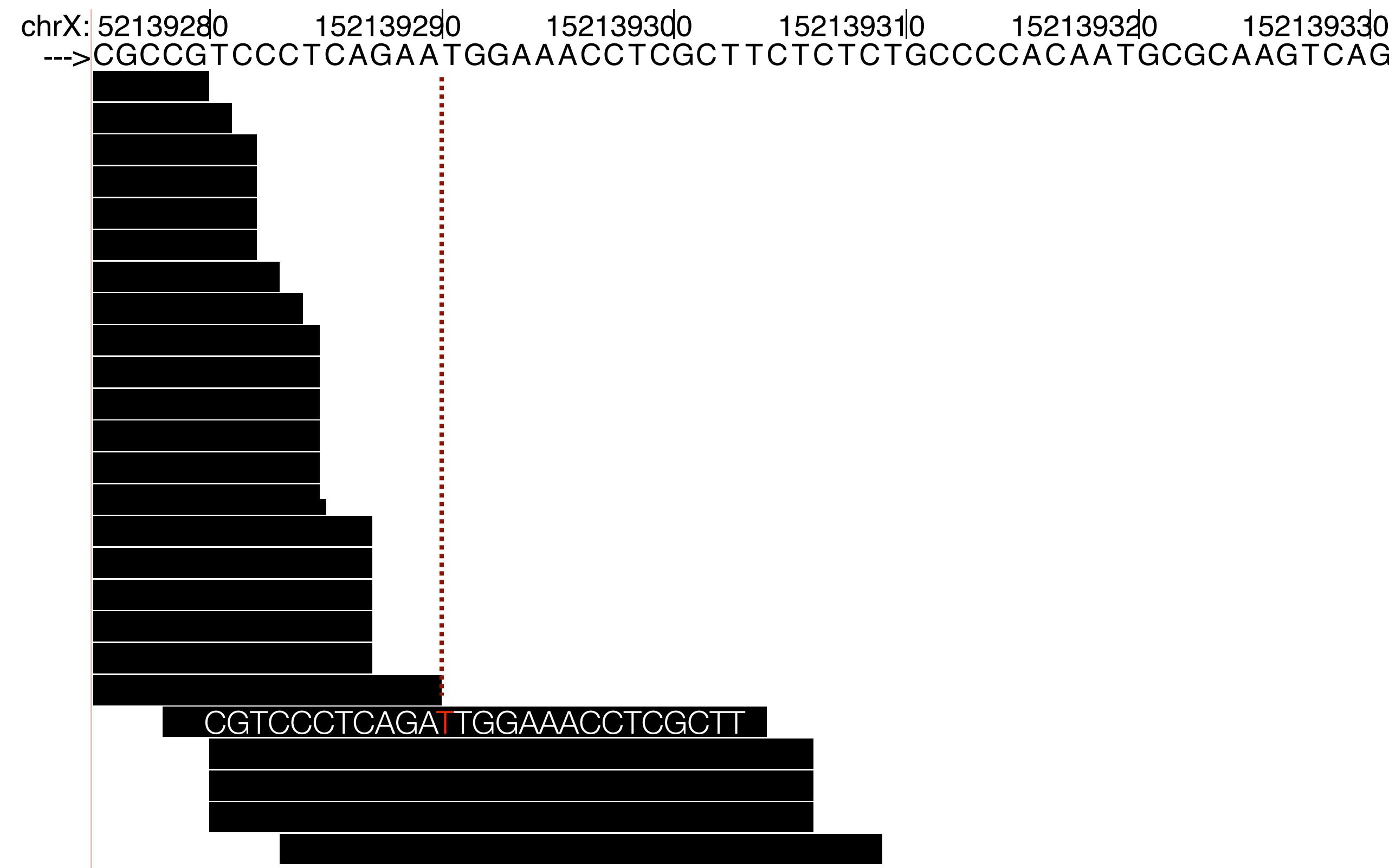
Tool evolution: mapping reads to a genome



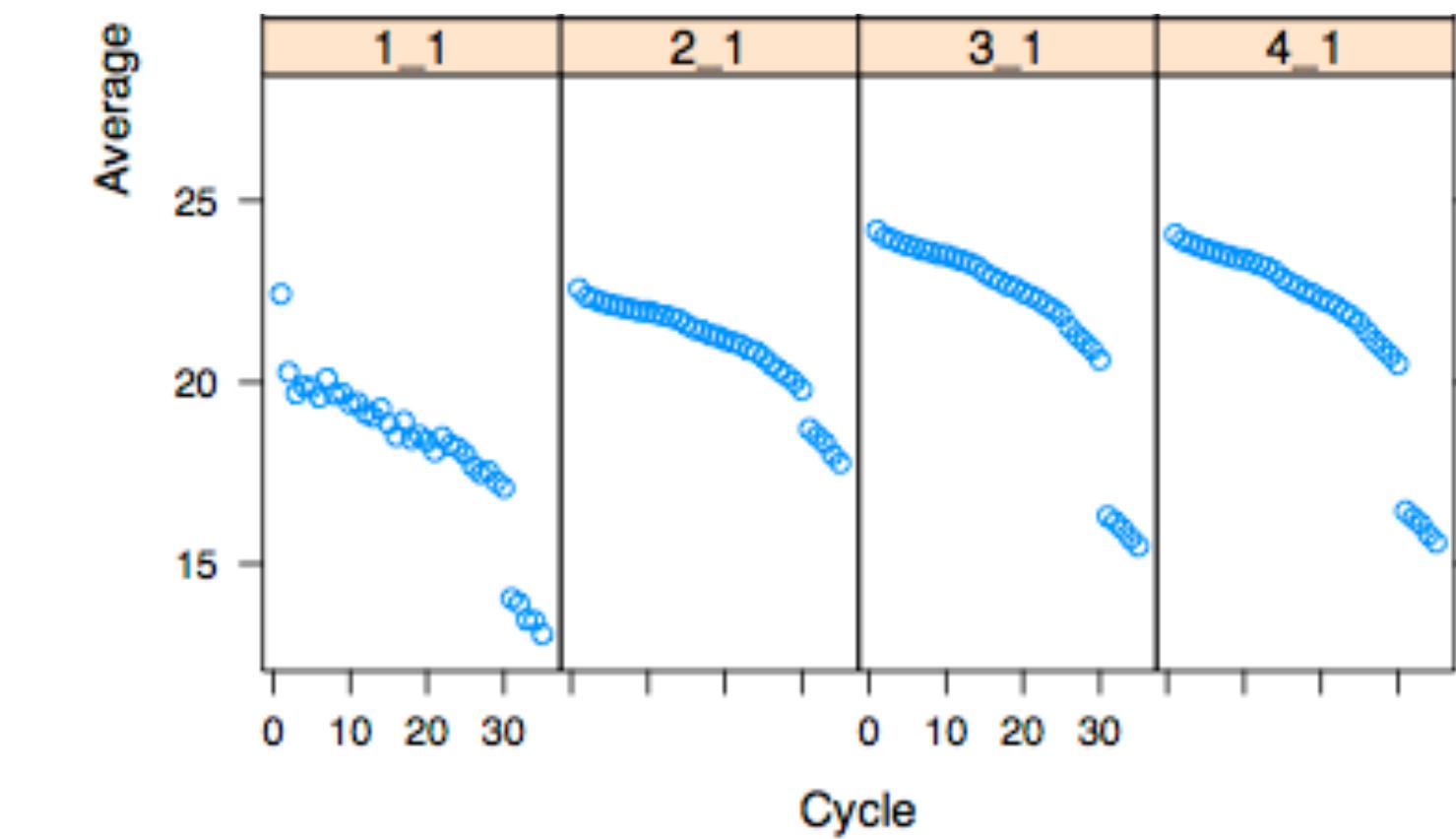
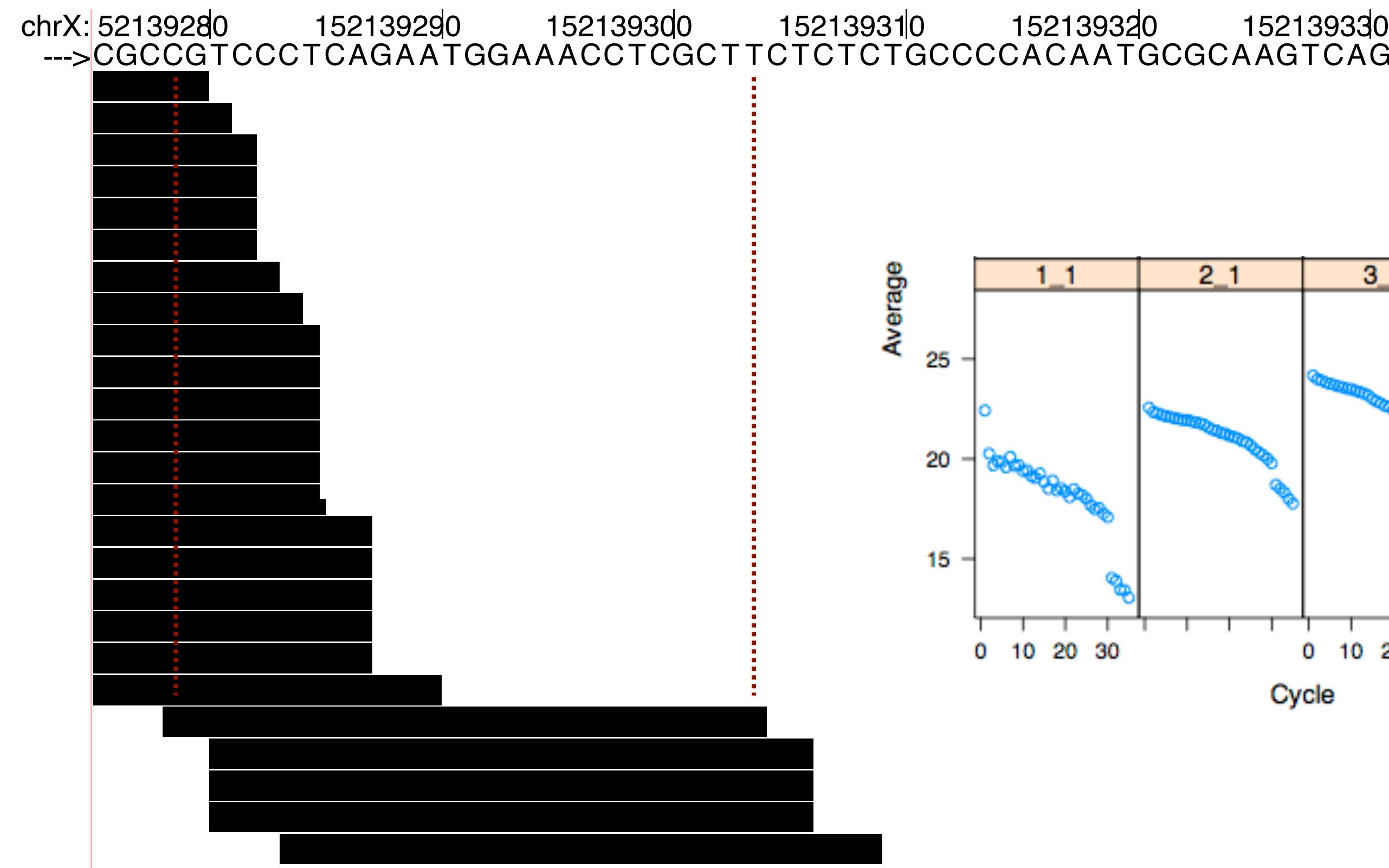
Mapping to a Reference Genome



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.cbcb.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/marthlab/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	3.3.5.2	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		SHRIMP 2	Shrill	TSV	-	-	Y	Y	B	L	N Y	N	Lib and de novo
RazerS	DNA	OS	1.1	Linux,Mac,Windows	58	18,18		Stampy	TSV ELAND	11	*	Score	Score	Y	A,B,S	G	N N	Y	N
mrsFAST	miRNA	OS	2.3.0	Linux	32	18,18		MapSplice	TSV ELAND	25	200	Y	0	N	A	G	N N	Y	N
Stampy	DNA	Bin	1.0.16	Linux,Mac	26	18,18		REAL	TSV	4	4K	0.15	30	N	B,R,S	G	N Y	Y	N
PASS	DNA	Bin	1.6.2	Linux,Mac,Windows	45	18,18		BS Seeker	TSV ELAND	23	1K	Y	Y	Y	A,B	G	SM	Y Y	De novo
SOCS	DNA	OS	2.1.1	Linux,Mac,Windows	49	18,18		Supersplat	TSV	64	Y	0	N	A,B		SM	Y N	N	
GenomeMapper	DNA	OS	0.4.3	Linux,Mac	31	18,18		SpliceMap	TSV	12	2K	10	10	Y	A,B,R	G	SM	N N	N
Slider	DNA	OS	0.6	Linux,Mac,Windows	39	18,18		BRAT	TSV	62	3	0	N	B,S		N Y	Y	N	
BSMAP	Bisulfite	OS	2.43	Linux,Mac	31	18,18		BFAST	TSV	8	144	15	0	N	B,U,S		SM	N Y	N
PerM	DNA	OS	0.4.0	Linux,Unix,Mac,Windows	30	18,18		GNUUMAP	TSV	20	128	9	0	Y	A,U	G	DM	Y Y	N
BWT-SW	DNA	OS	20070916	Linux	15	18,18		GenomeMapper	TSV		1K	Score	Score	Y	A		N N	N N	N
SHRIMP 2	DNA	OS	2.2.2	Linux, Unix, Mac	15	18,18		mrFAST	TSV	30	1K	Y	Score	N	B,U,S	G	SM	Y Y	N
RNA-Mate	RNA	OS	1.1	Linux,Mac	28	18,18		PerM	TSV		10K	Y	Score	N	S	-	DM	Y N	Lib
Supersplat	RNA	OS	1.0	Linux,Mac	28	18,18		X-Mate	TSV			0	0	Y	A,U	G	N N	N N	De novo
PatMaN	miRNA	OS	1.2.2	Linux,Mac	28	18,18		BFAST	TSV	1	*	Y	Y	N	A	G	N N	N N	N
BS Seeker	Bisulfite	OS	1.2	Linux,Mac	28	18,18		RazerS	TSV	-	-	3	0	N	U	-	SM	Y N	N
Slider II	DNA	OS	1.1	Linux,Mac,Windows	28	18,18		SHRIMP	TSV	93	Y		N	B,S		N N	Y	N	
GNUMAP	DNA	OS	3.0.2	Linux,Mac	15	18,18		BWA	TSV	16	1K	Score	Score	Y	B	G	SM/DM	Y N	N
MOM	DNA	Bin	0.6	Linux,Mac,Windows	15	18,18		BWA-SW	TSV			Y	0	N	A	L	SM	N Y	N
Bismark	Bisulfite	OS	0.7.3	Linux,Mac	15	18,18		CloudBurst	TSV	16	10K	Score	Score	N	U	-	SM	Y Y	N
BRAT	Bisulfite	OS	1.2.3	Linux	15	18,18		ProbeMatch	TSV			Y	0	N			N N	Y	N
SOAPSplice	RNA	Bin	1.8	Linux,Mac	15	18,18		TopHat	TSV			Y	0	N			N N	Y	N
WHAM	DNA	OS	0.14	Linux, Unix	15	18,18		Bowtie	TSV	13	3K	5	2	Y	U	-	SM	Y Y	De novo
PASS	DNA	OS	0.1	Linux	15	18,18		MOM	TSV	5	128	5	3	N	A,B,R,U,S	G	N Y	Y	De novo
MicroRazerS	miRNA	OS	1.11	Linux,Mac	15	18,18		PASS	TSV	10	*	Score	0	N	S	G	N N	N N	N
RUM	RNA	OS	1.11	Linux,Mac	15	18,18		Slider II	TSV	-	-	Y	Y	Y	B	-	SM	N Y	De novo
ProbeMatch	DNA	OS	1	Linux,Mac	15	18,18		QPALMA	TSV	36	50	3	Y	N	A,B		N N	N N	N
X-Mate	DNA	OS	1	Linux,Mac	15	18,18		SOCS	TSV	-	-	Y	0	N	S	-	DM	Y N	Lib
SSAHA2	DNA	Bin	2.5.5	Linux,Mac	15	18,18		MAQ	TSV	-	-	Y	0	N	B,S	L	N N		

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

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

SAM/BAM

Courtesy of Nick Croucher, HSPH

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

	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9
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	1 1 1 1 1 1 1 1 1 1 1 1								

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

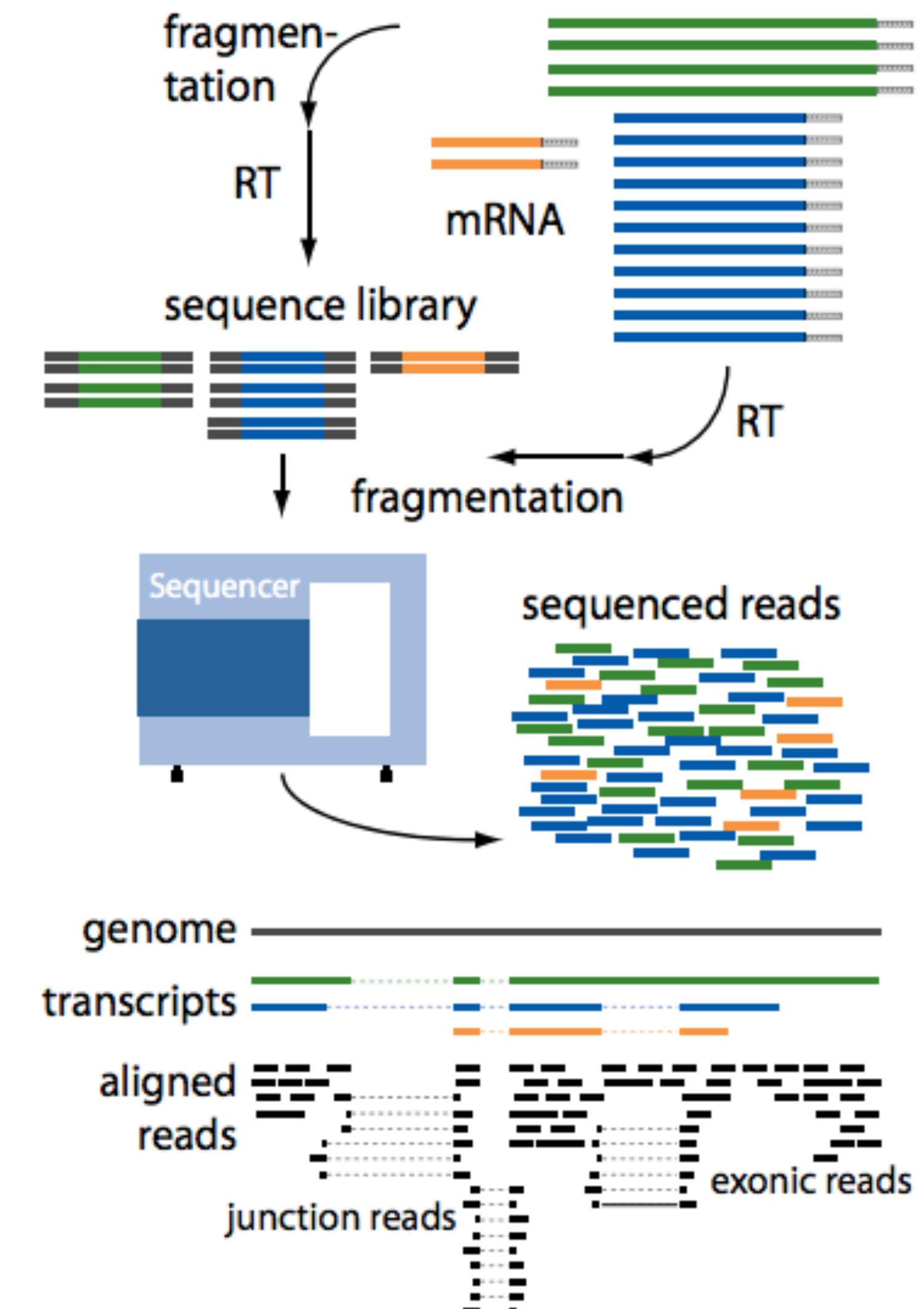
Mate mapped
to same
reference?

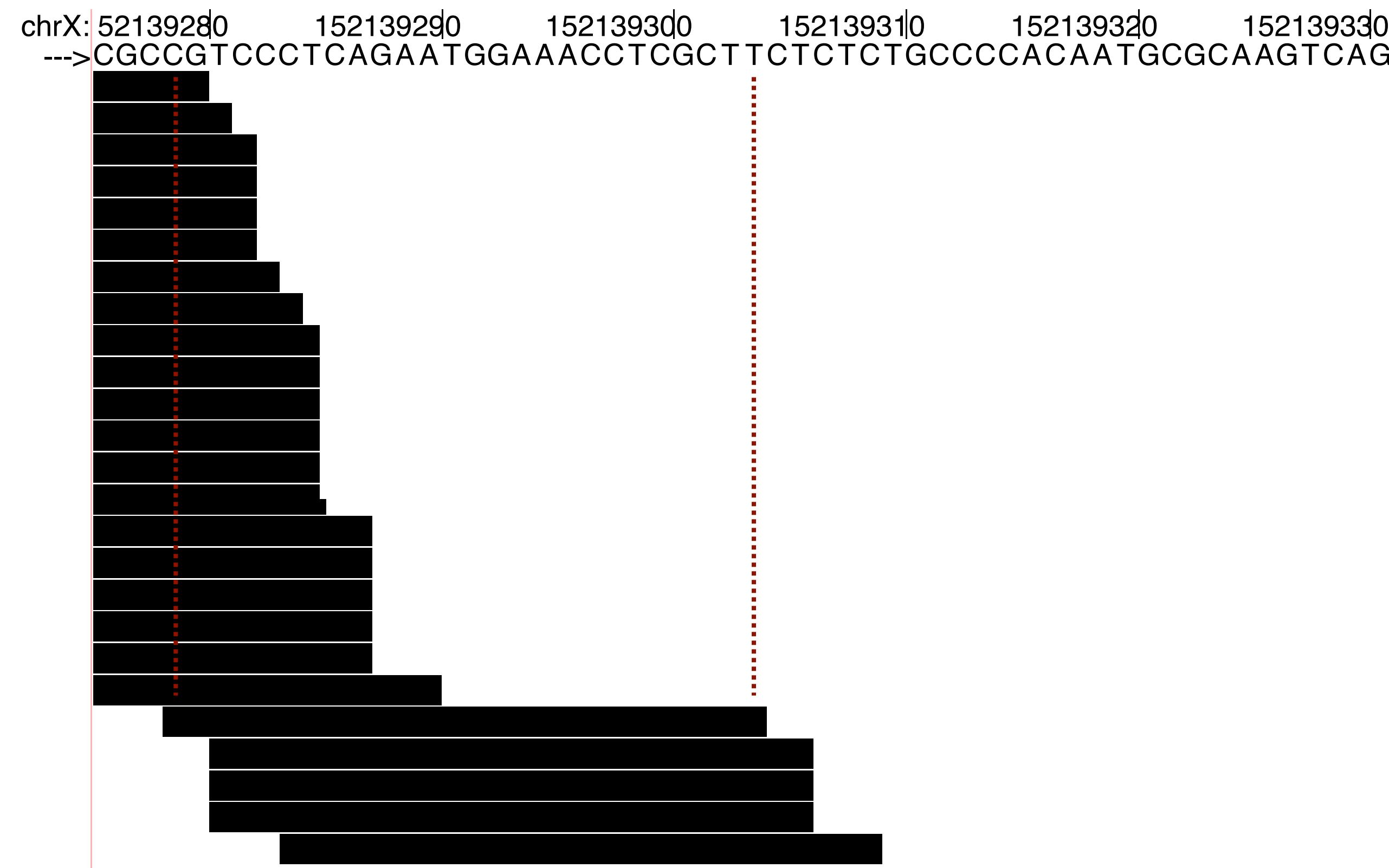
Total length

SAM/BAM

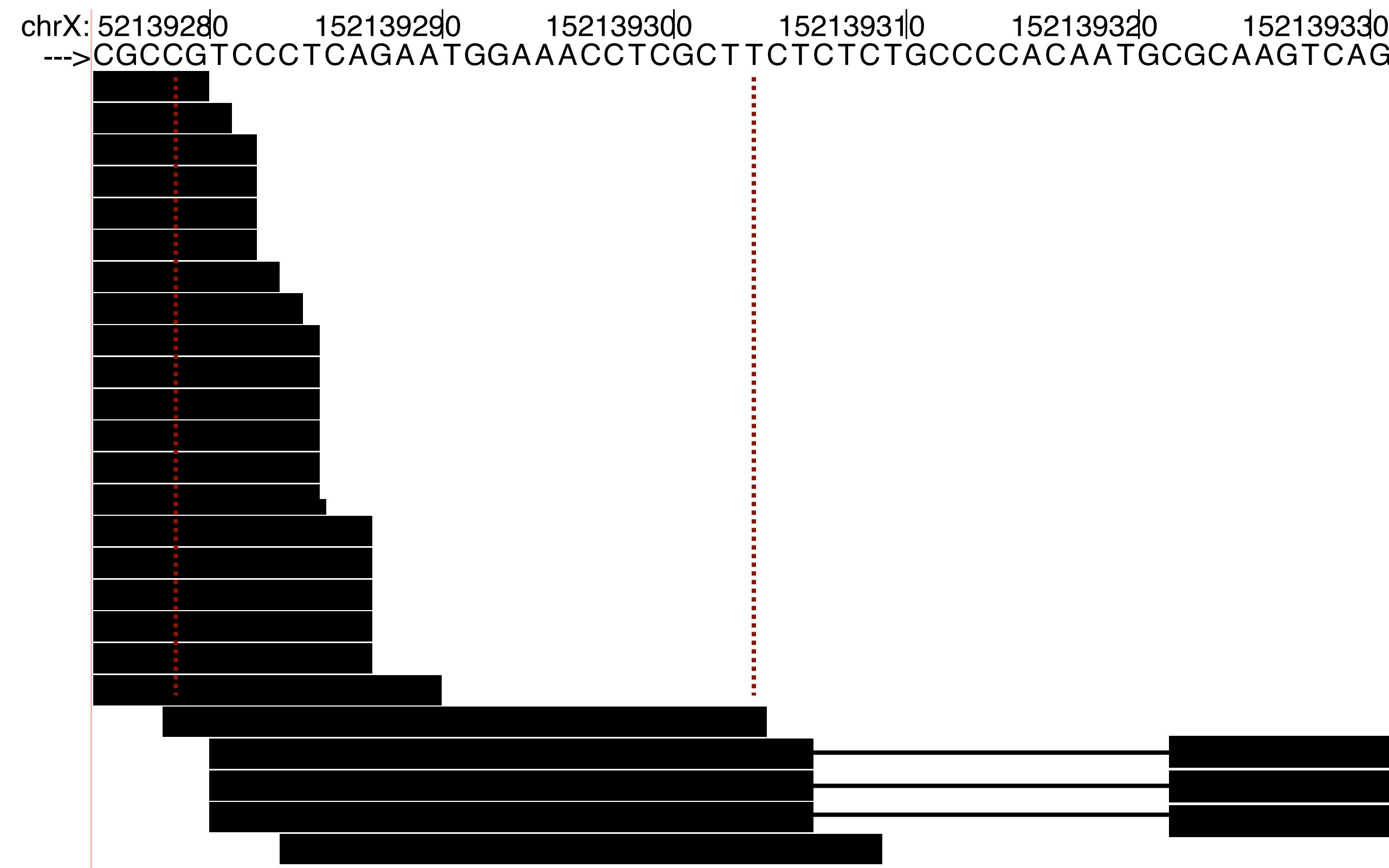
Courtesy of Nick Croucher, HSPH

RNA-seq

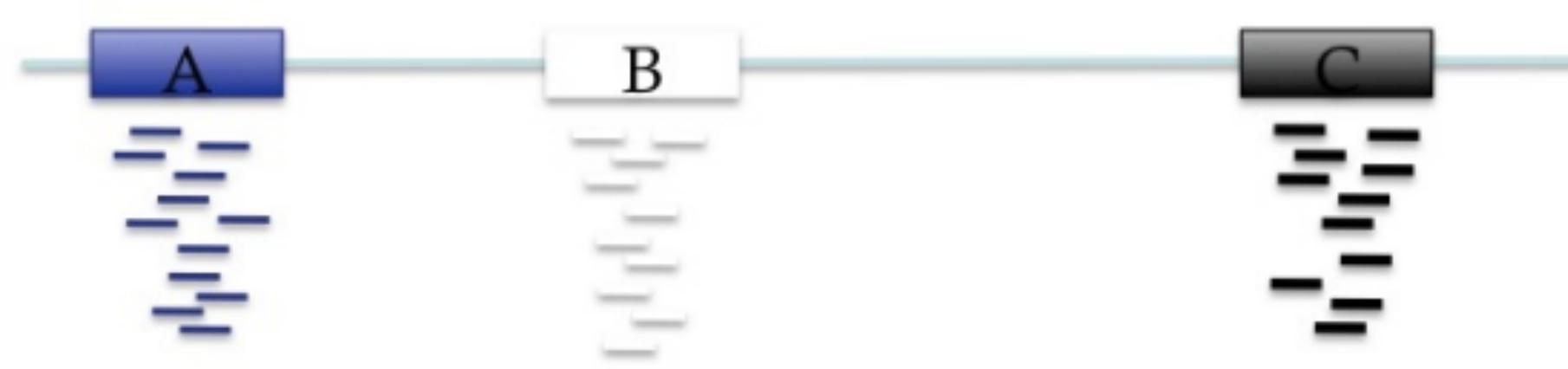




Mapping to a Reference Genome



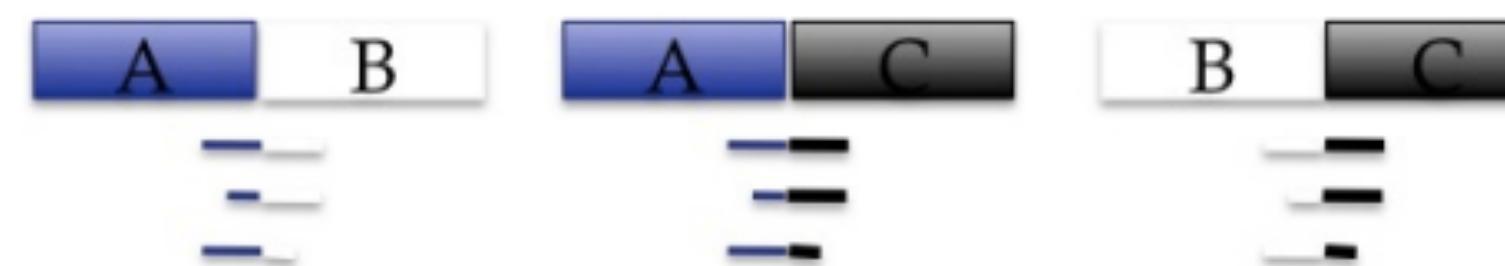
Mapping to a Reference Genome



identify candidate exons
via genomic mapping



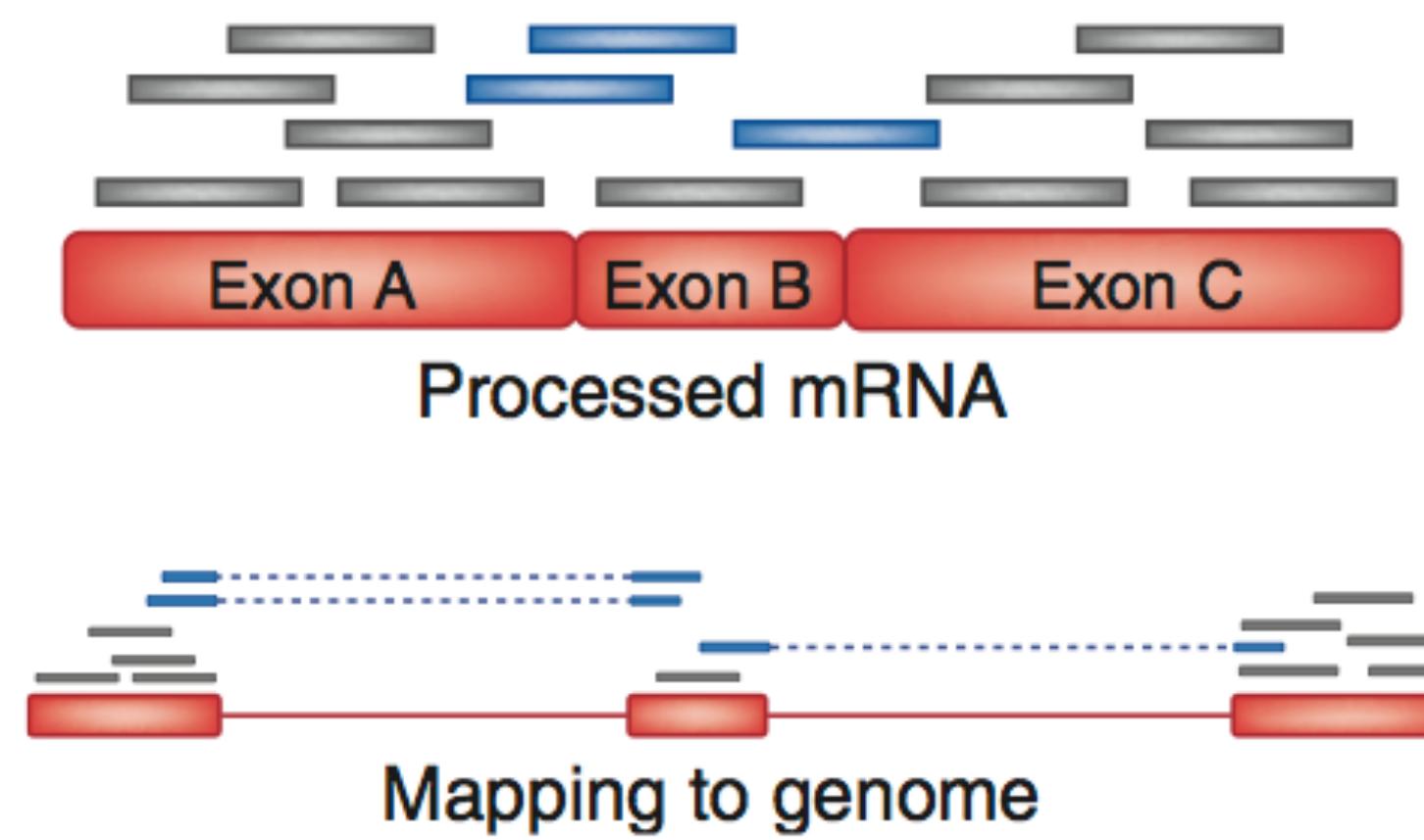
Generate possible pairings
of exons



Align reads to possible
junctions

(Trapnell, 2010)

TopHat alignment concept

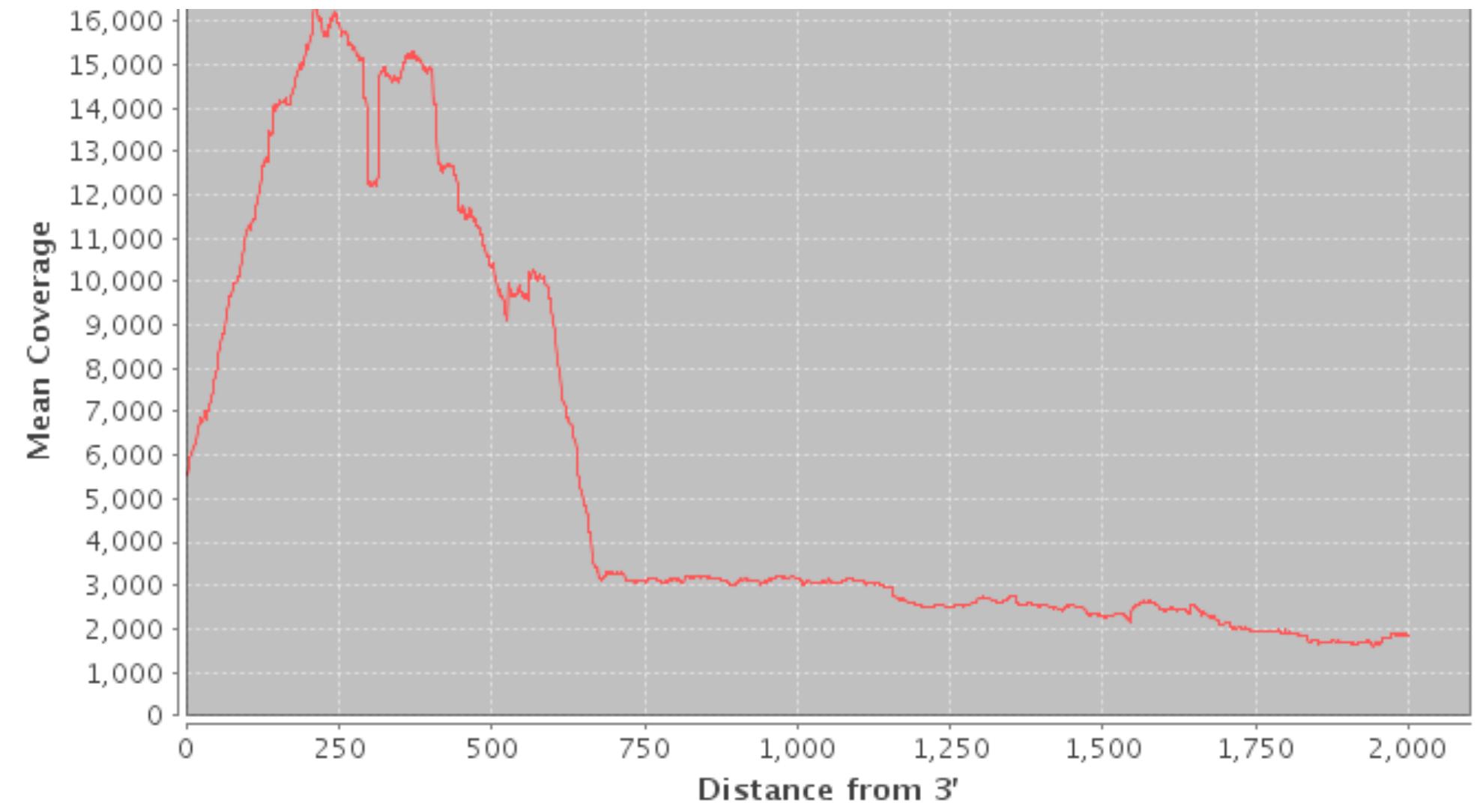
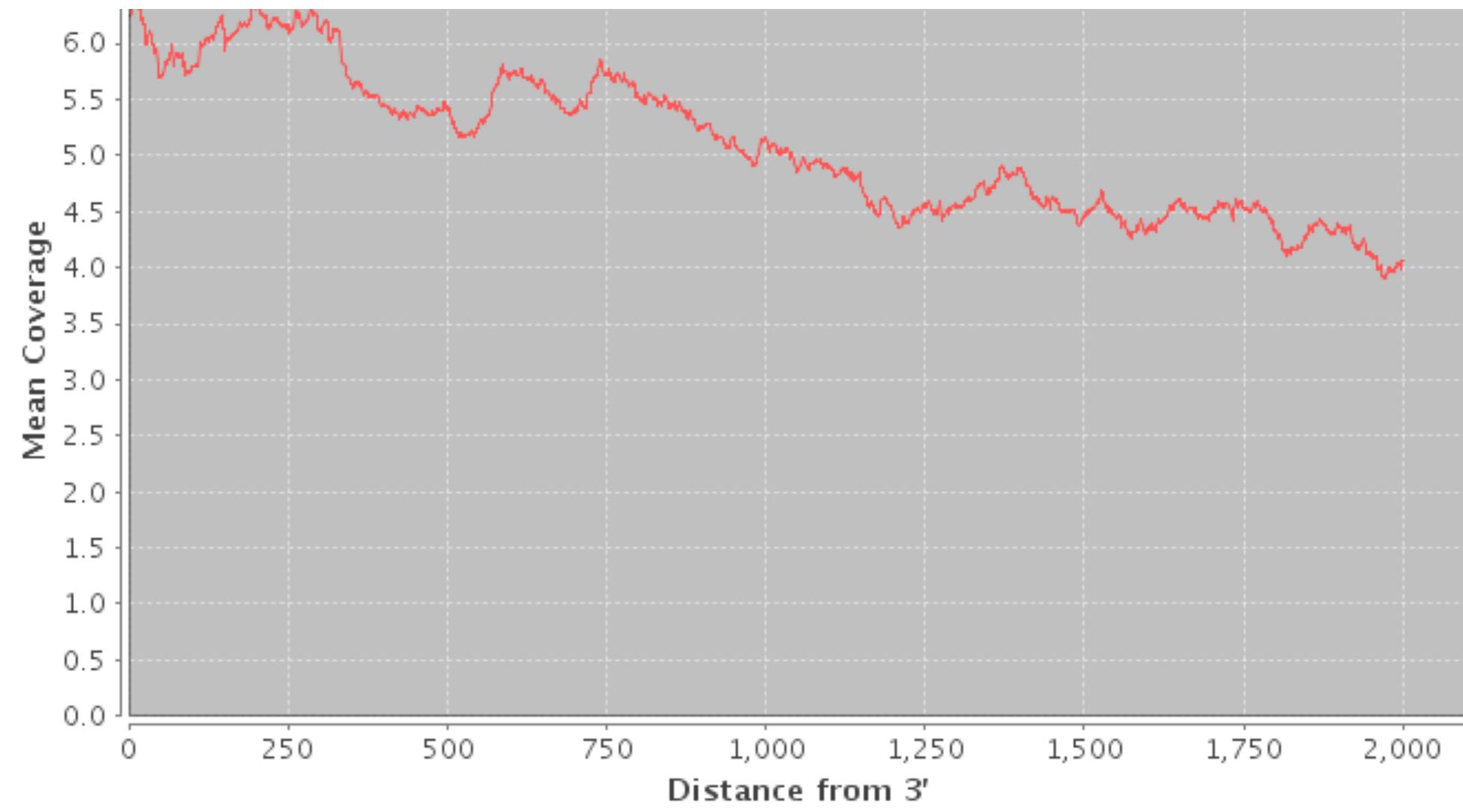


... GCAAACCACTGACCTGACTACTACGTCGTAACGTACACGGTAGCT... CCGTAGAATTGACTGTGTTG...

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

A
AA
AAT
AATT
AATTG
AATTGA

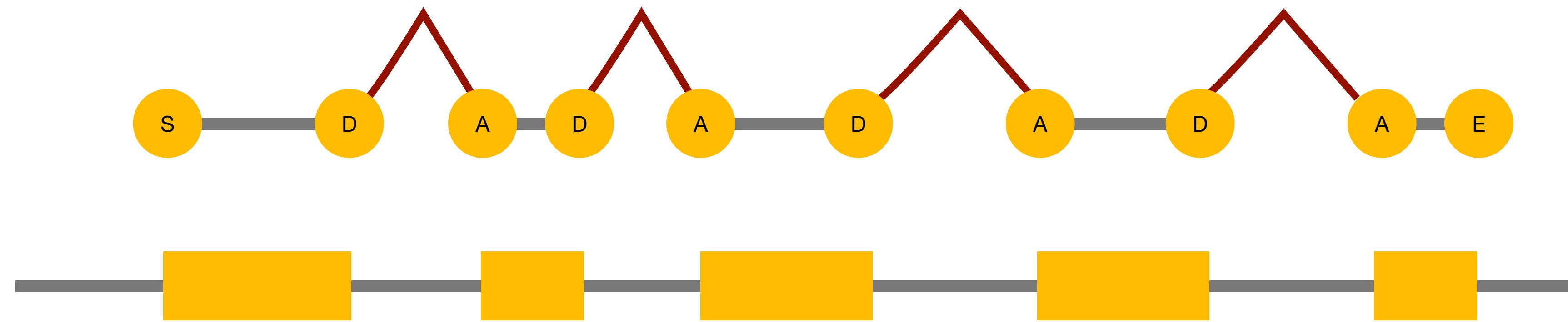
Transcript discovery



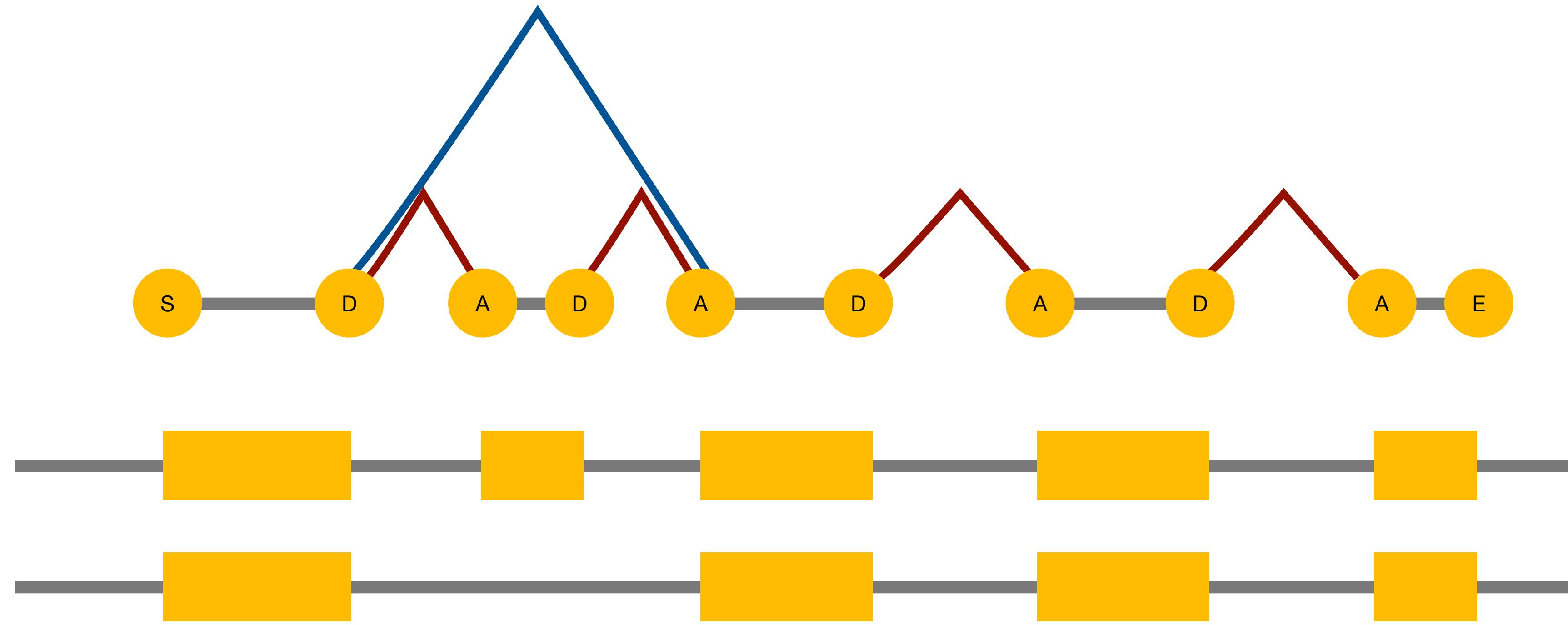
Additional QC



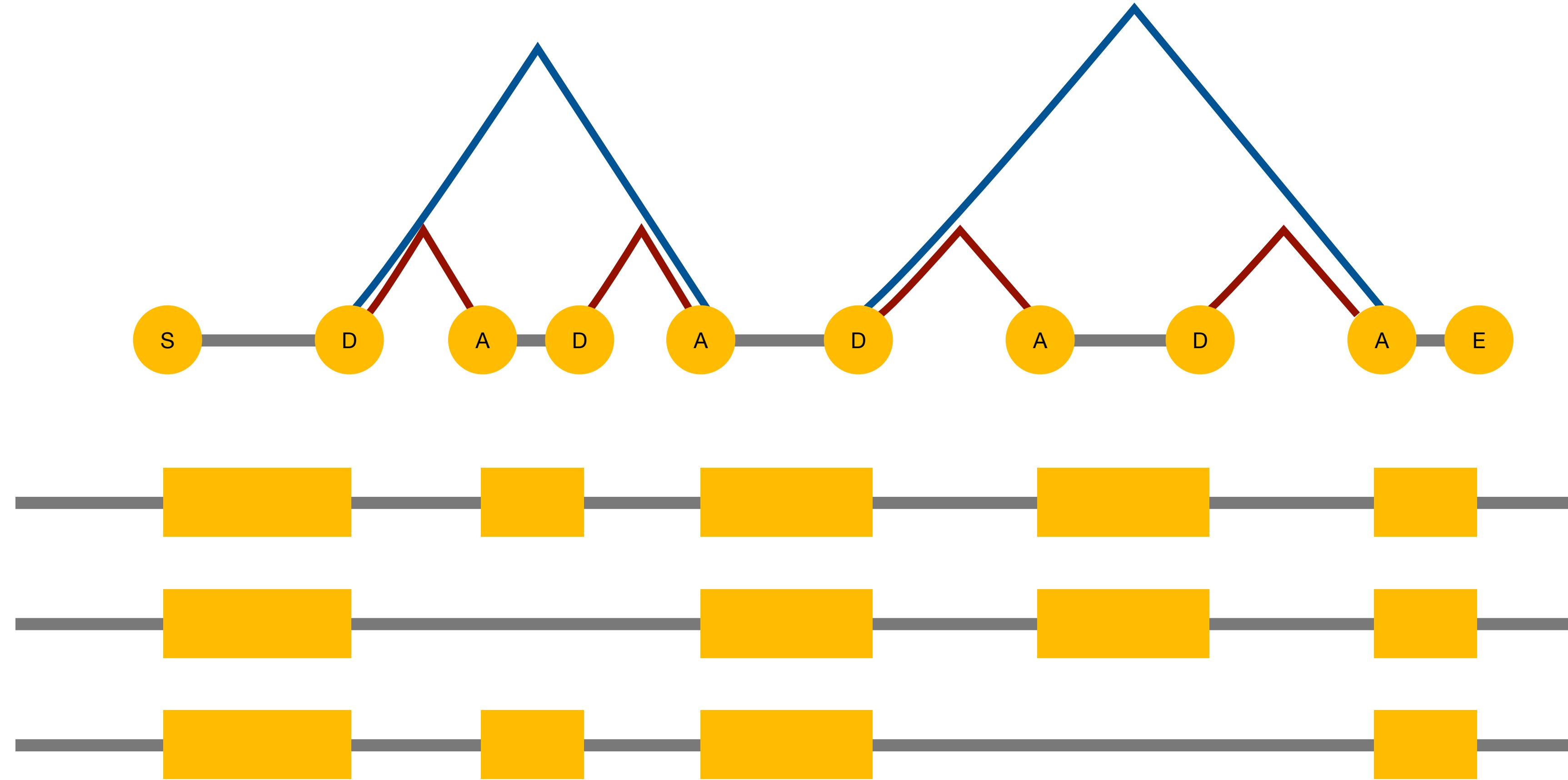
Splice graphs



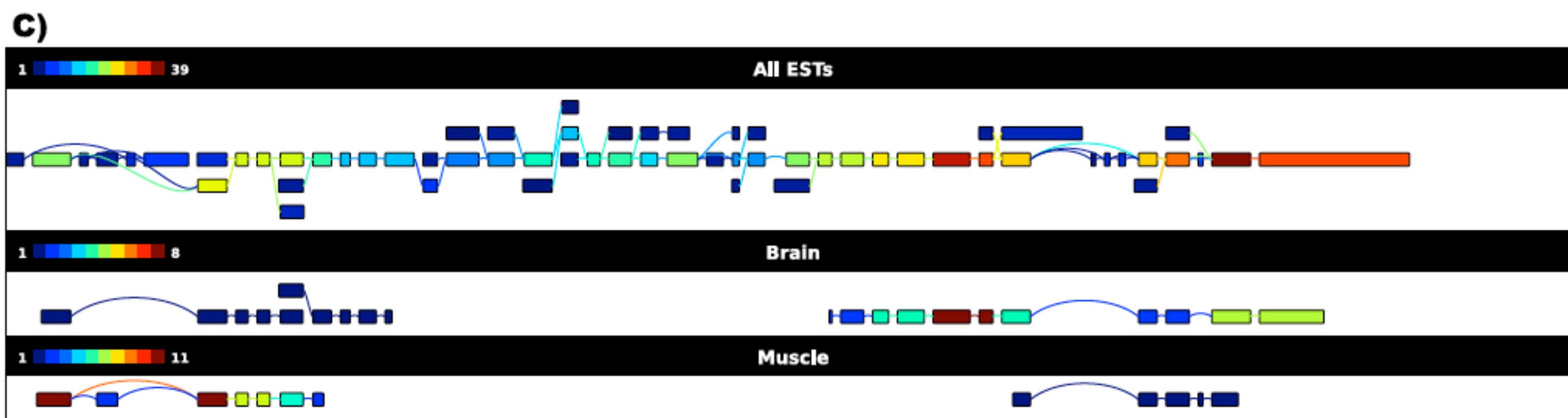
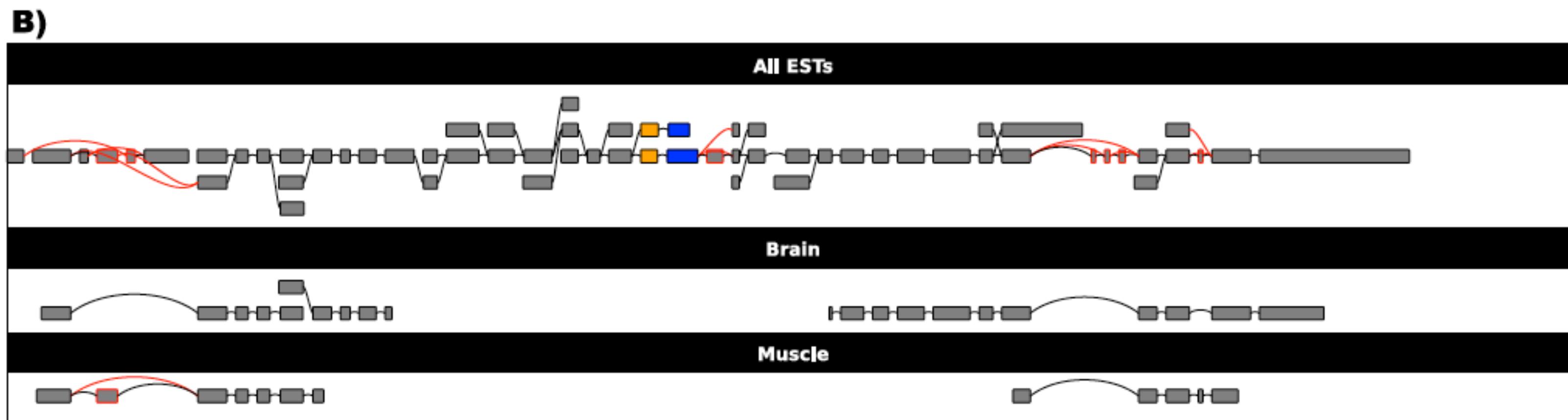
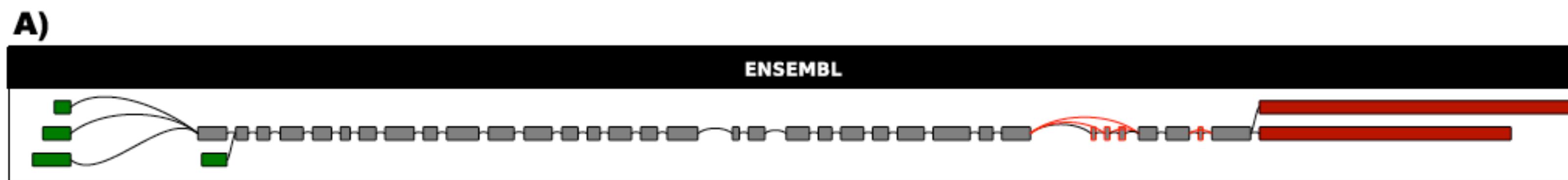
Splice graphs



Splice graphs

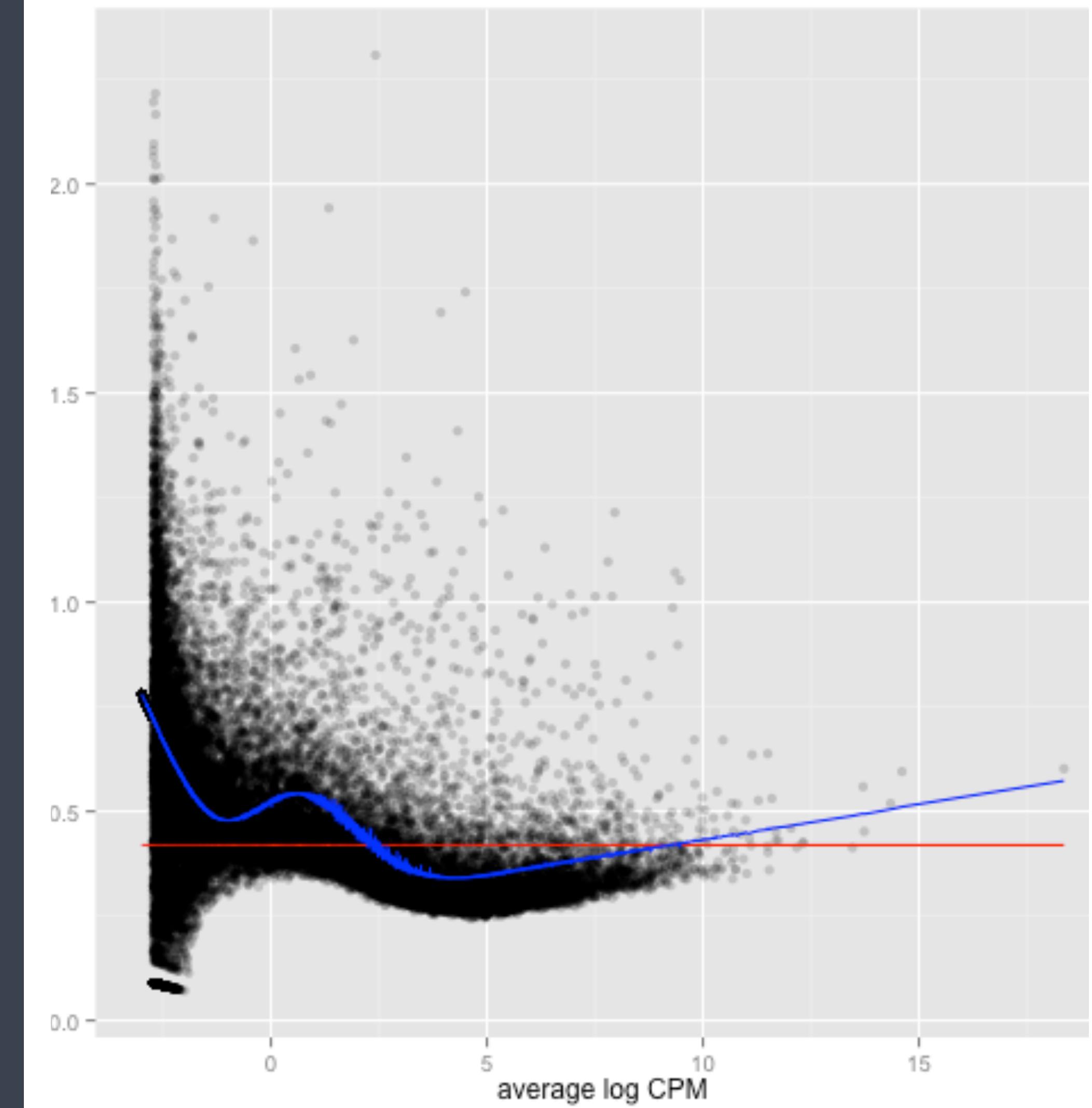


Splice graphs



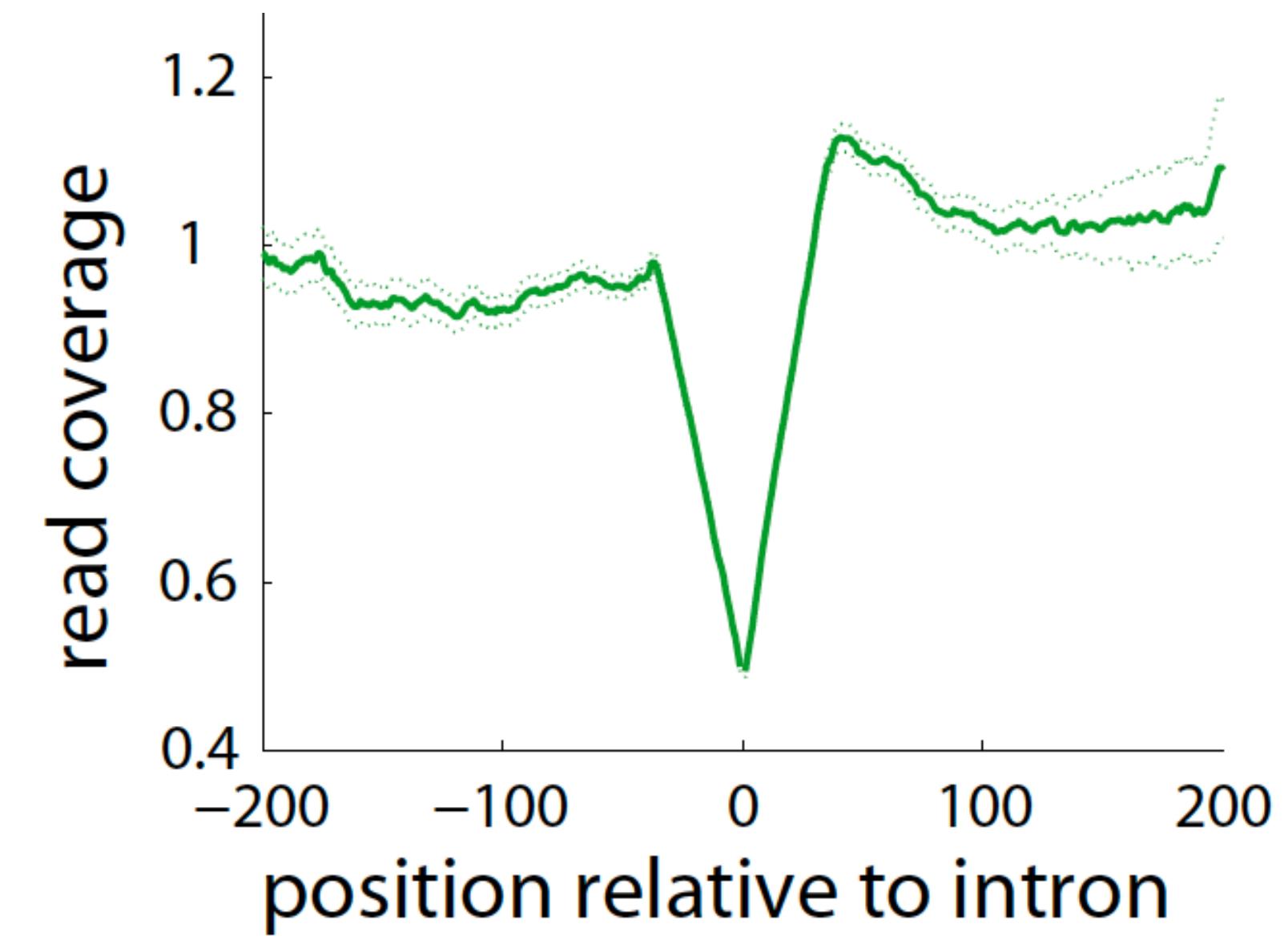
Quantification

Count normalization



Biases

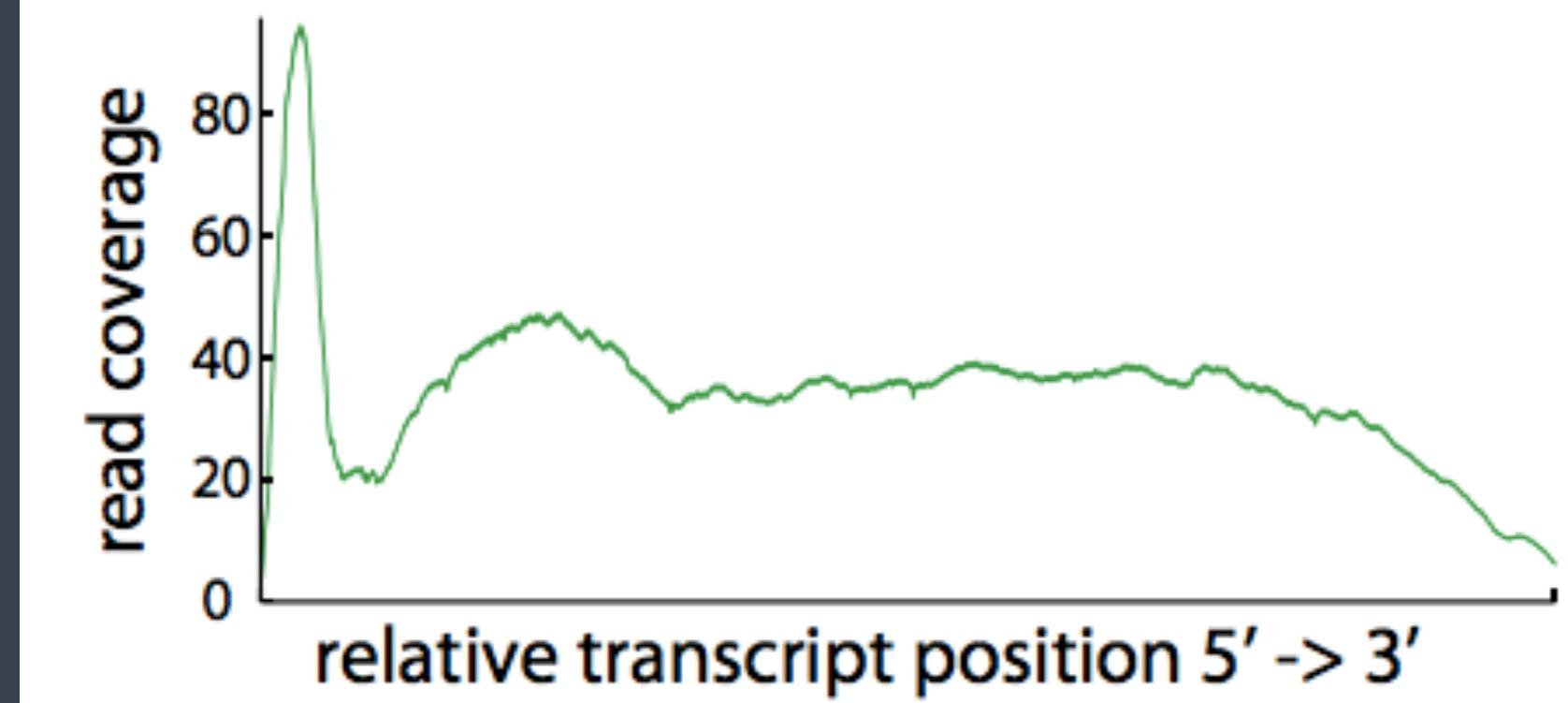
Read mapping

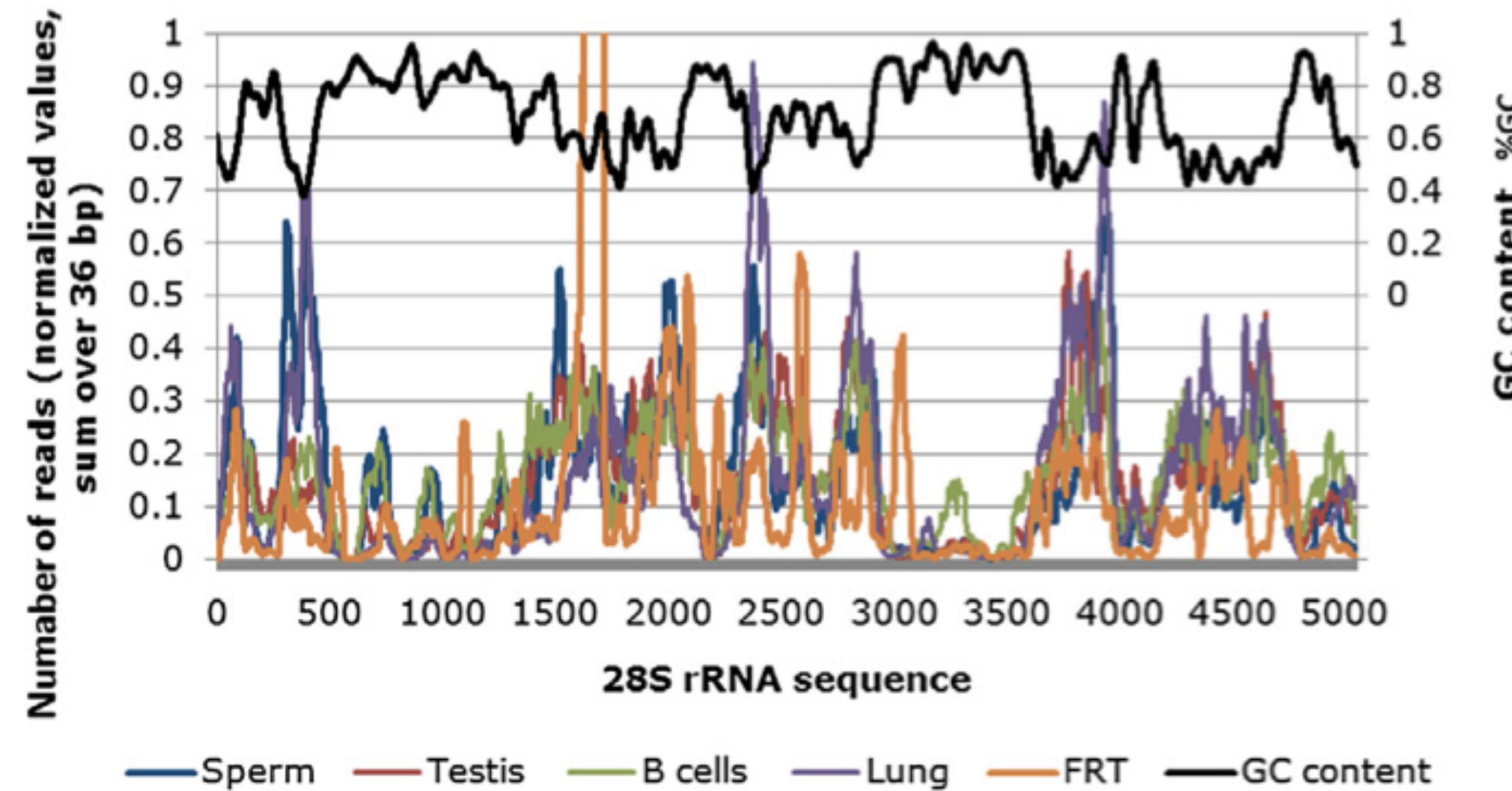


Biases

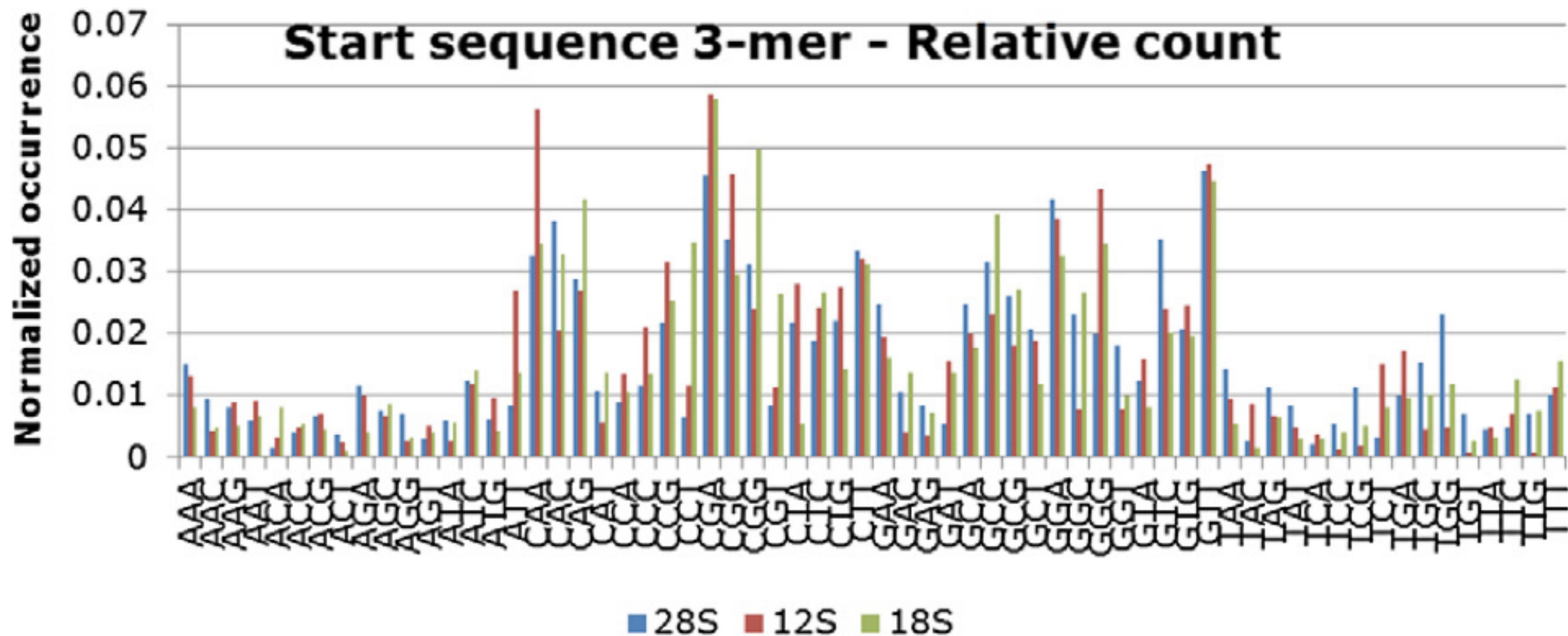
Library construction

Transcript length

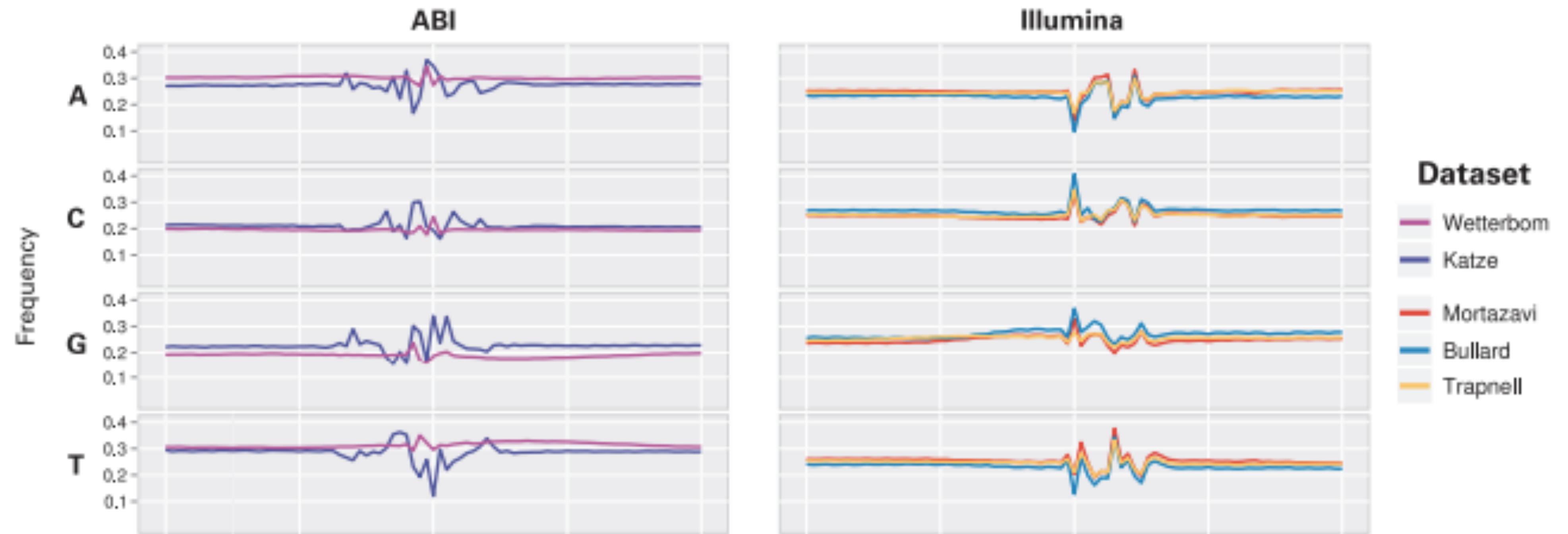




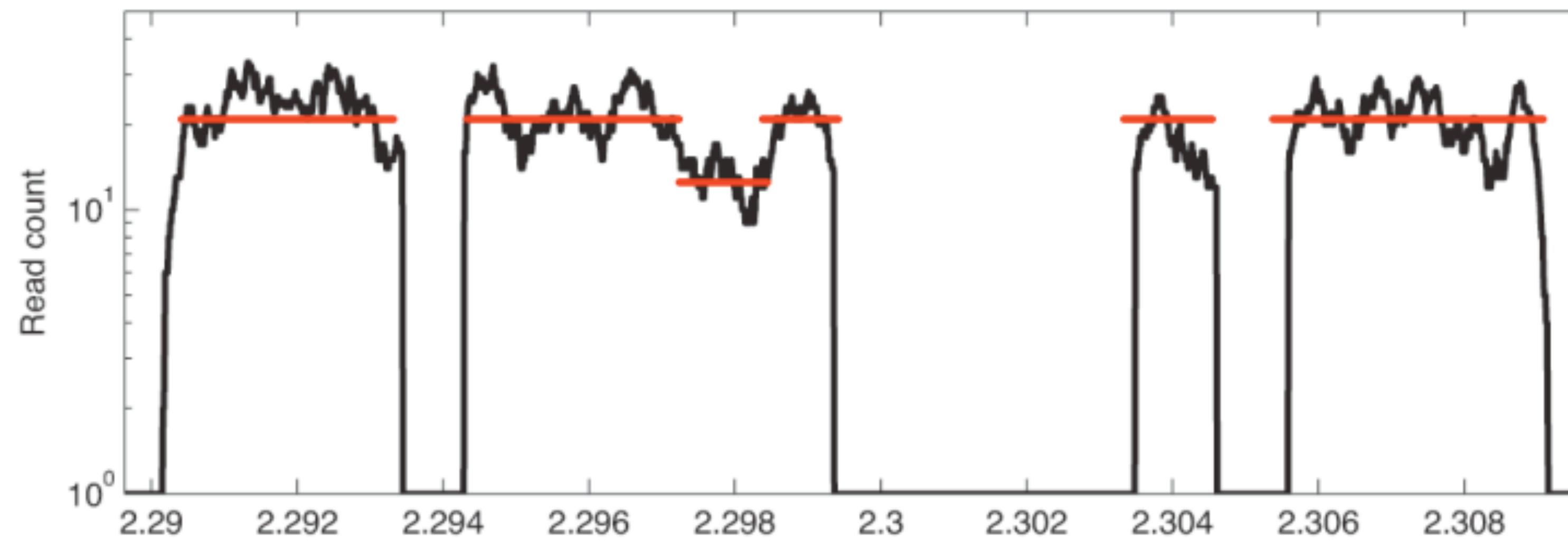
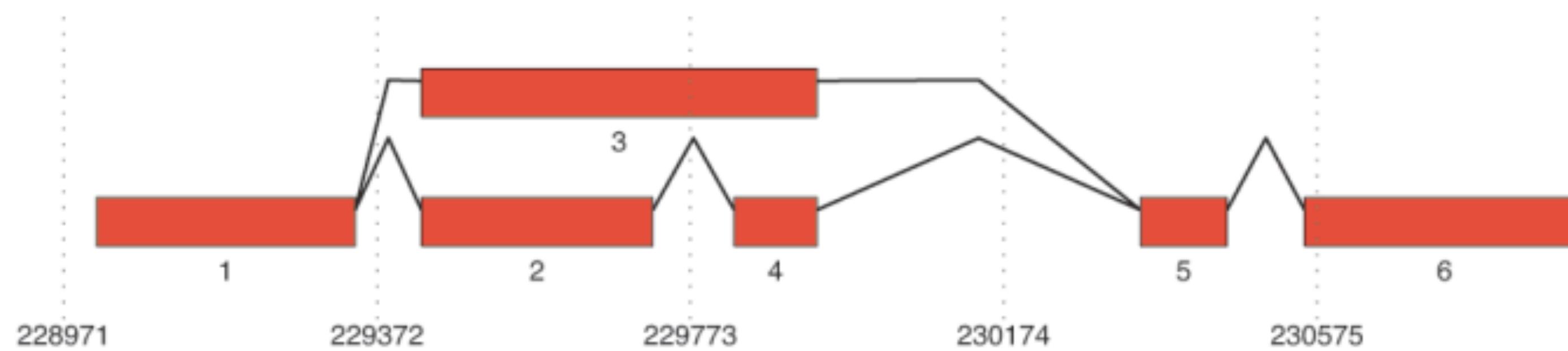
Nucleotide composition



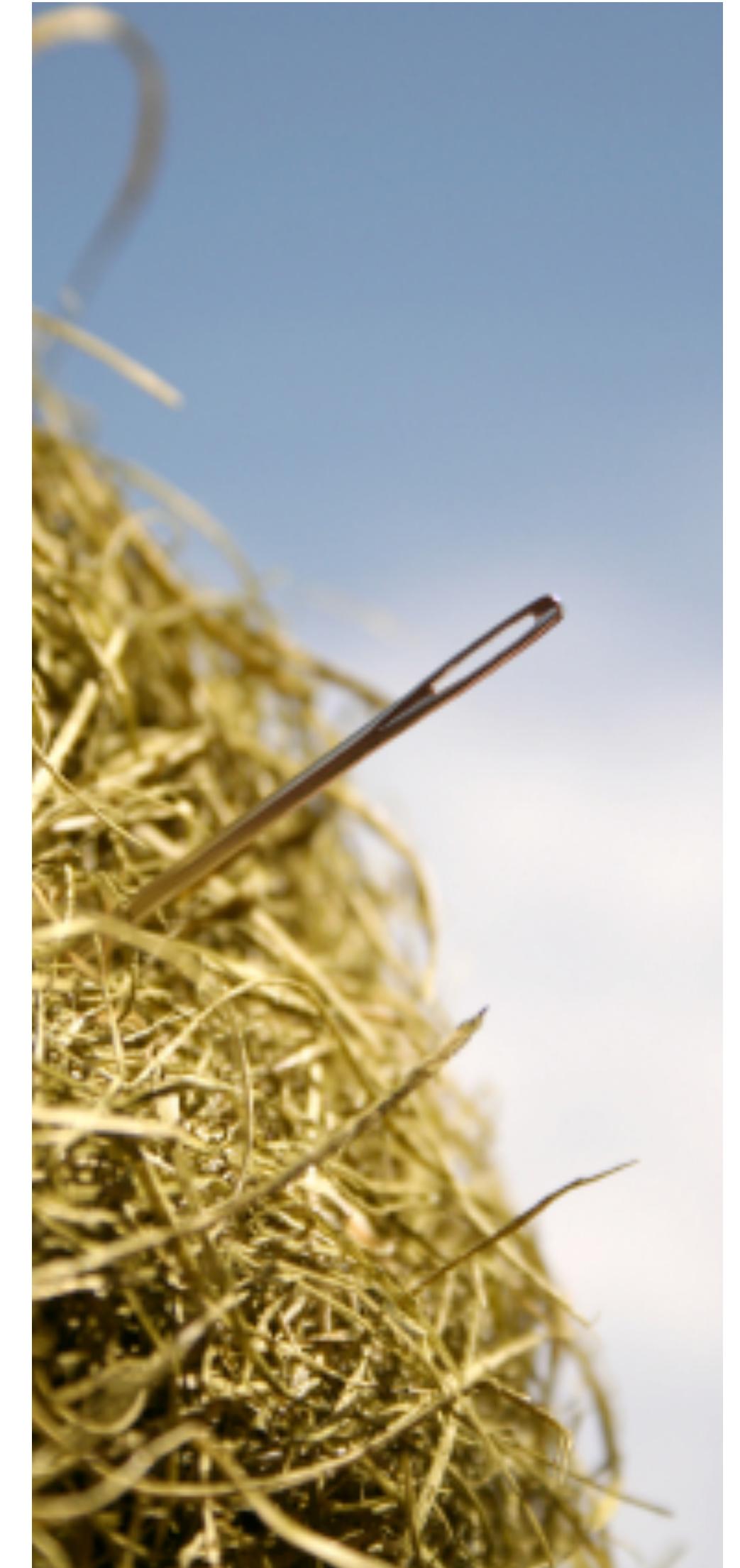
Not-so-random priming



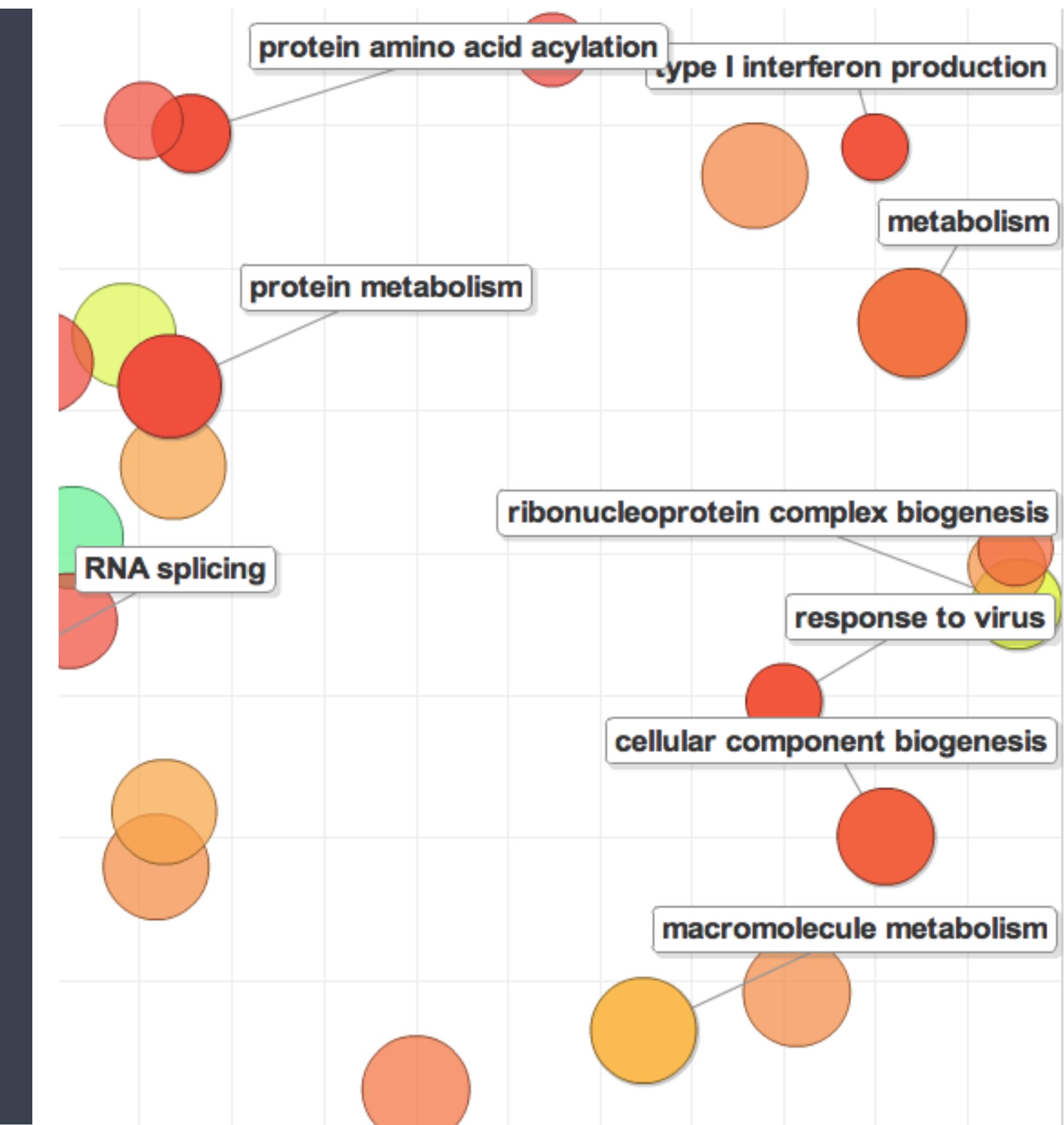
Start/end bias differs by technology

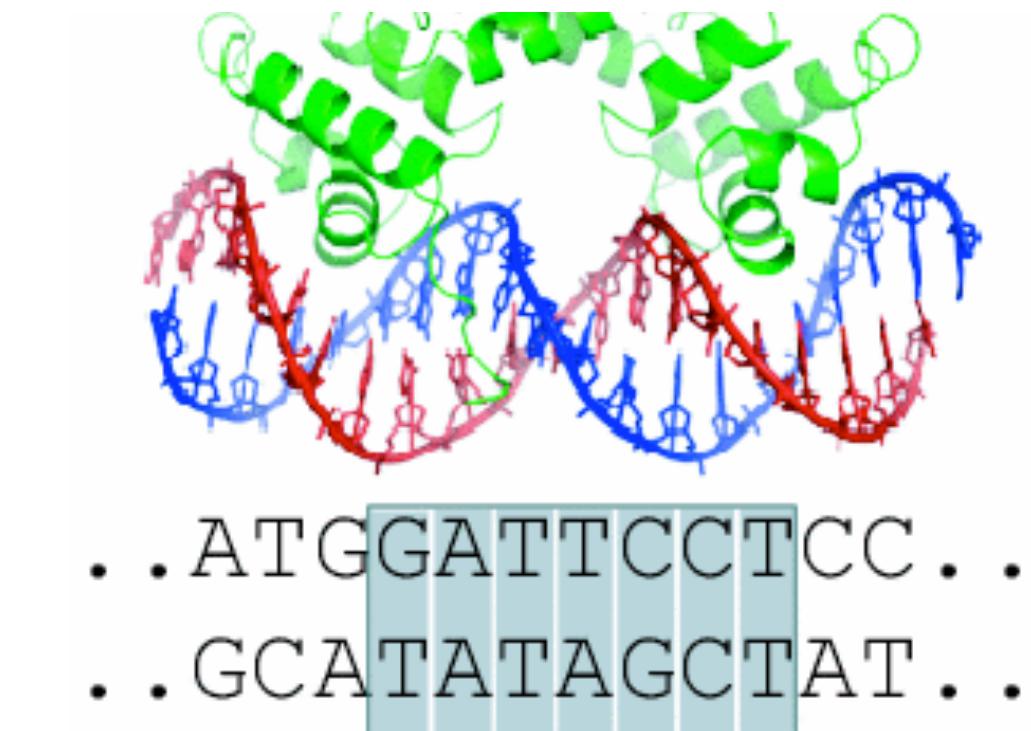
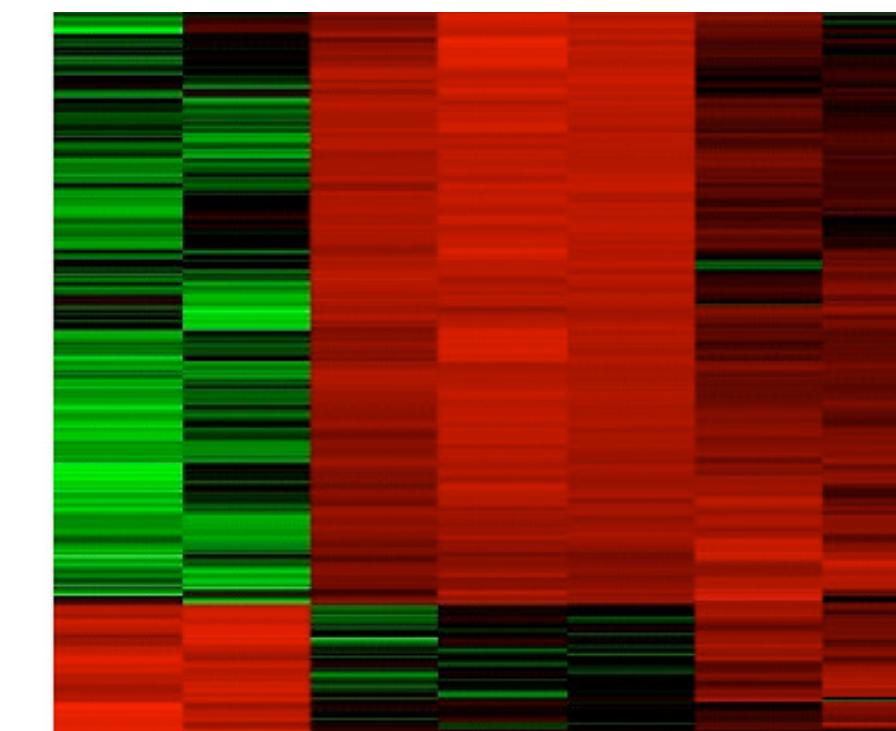
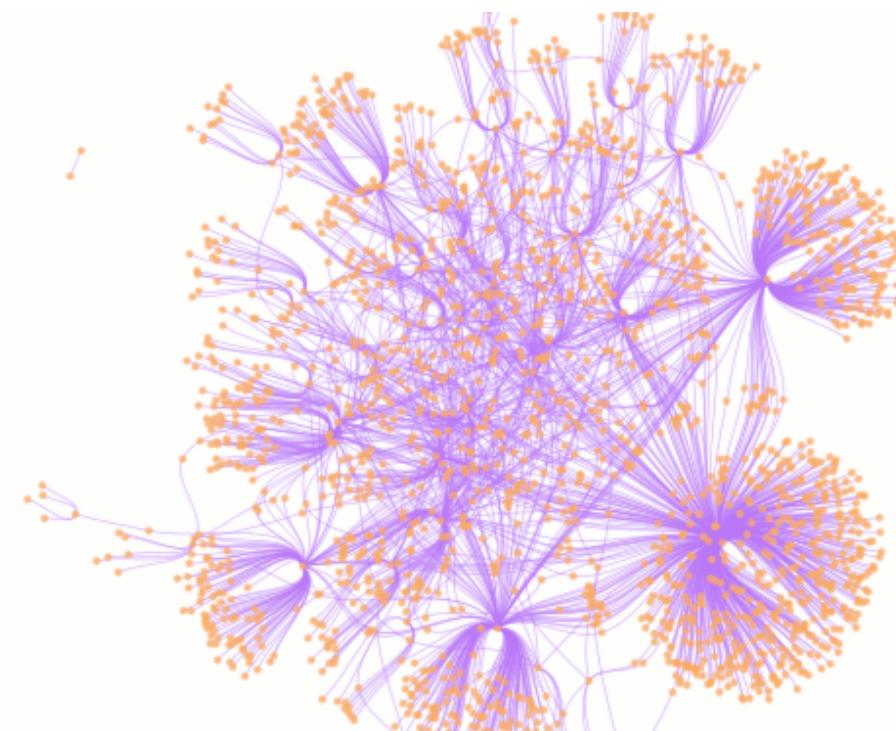
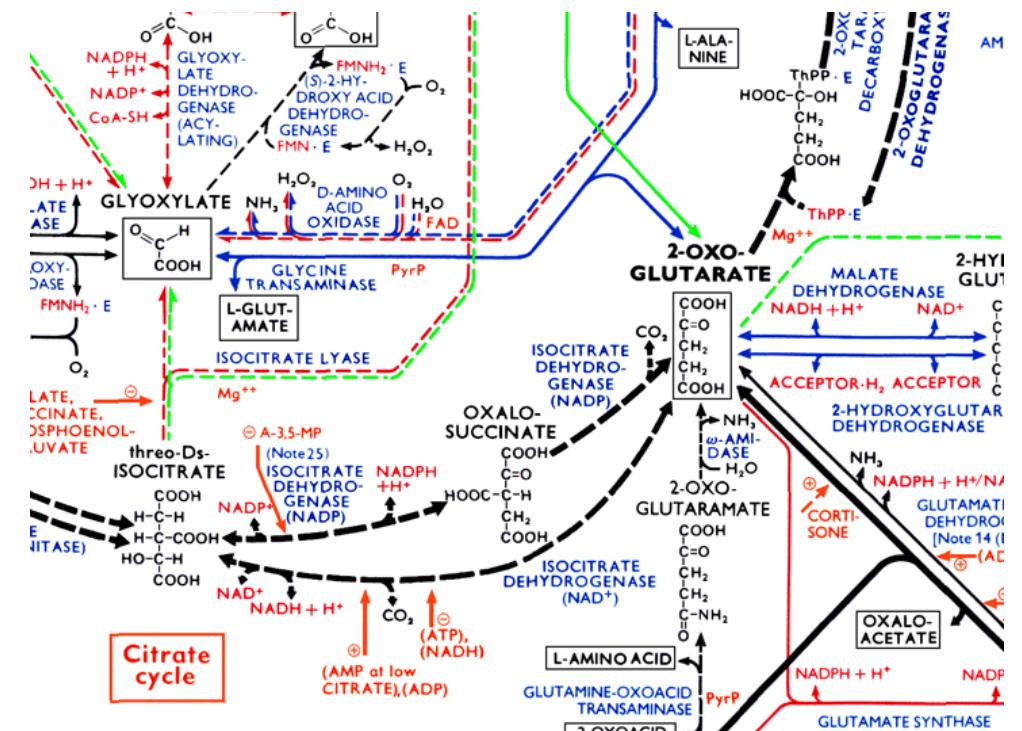


Making sense of your gene list



Functional enrichment



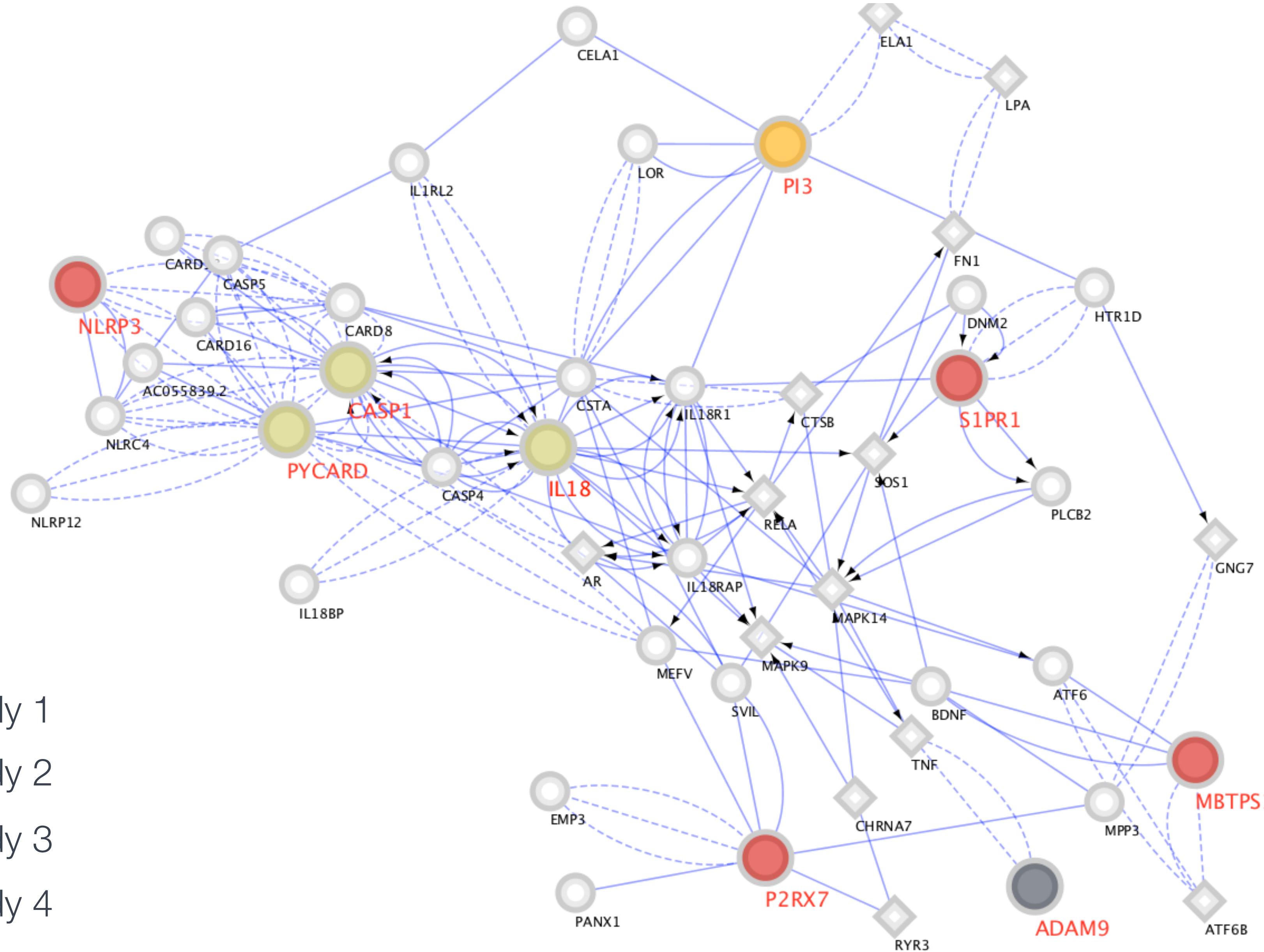


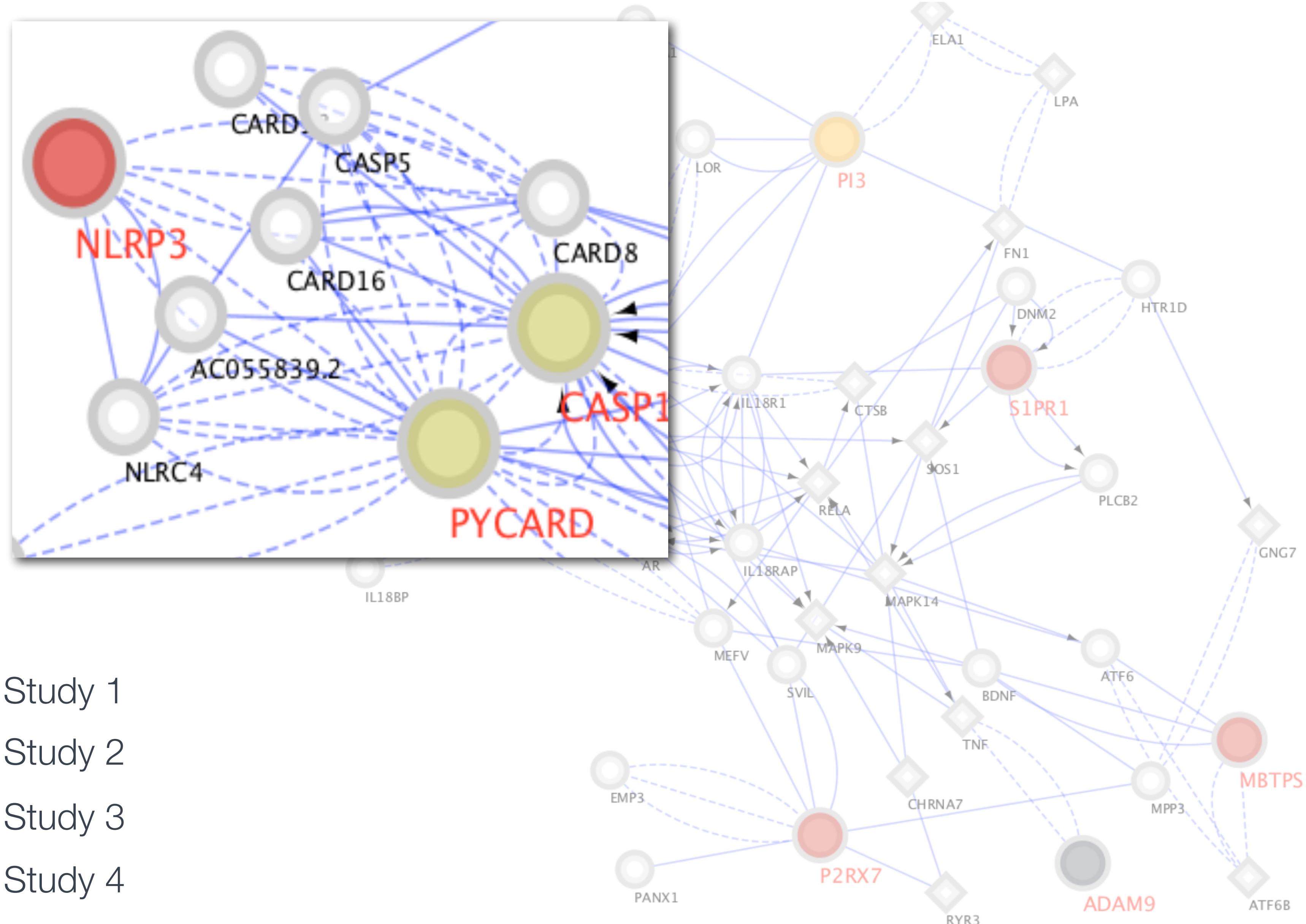
Known interactions

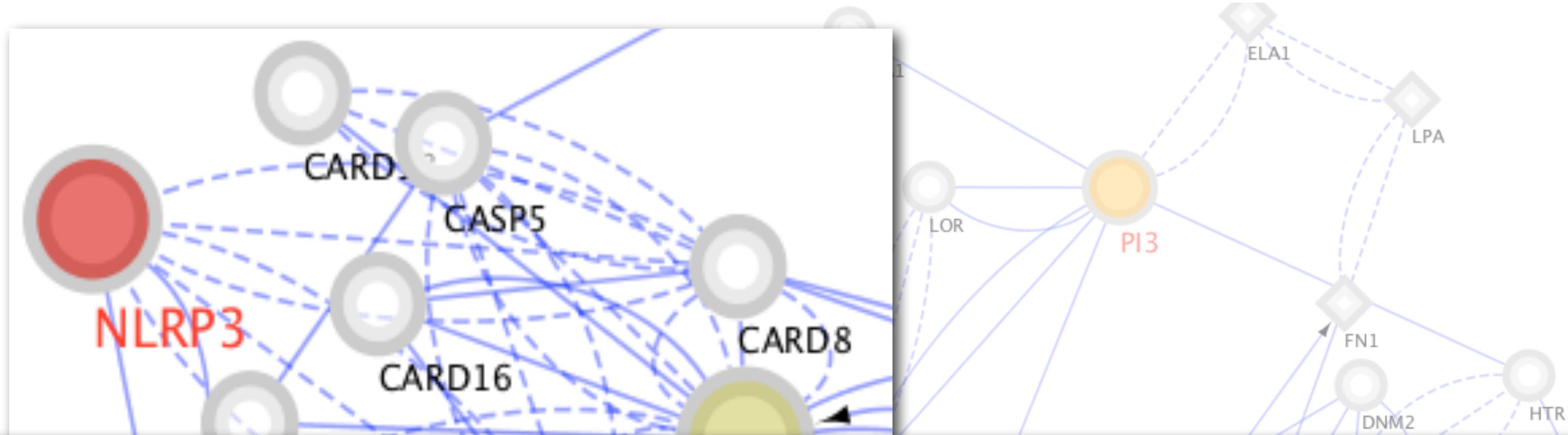
Protein-protein
binding

Consistent co-
expression

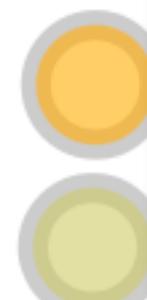
Transcriptional
regulation







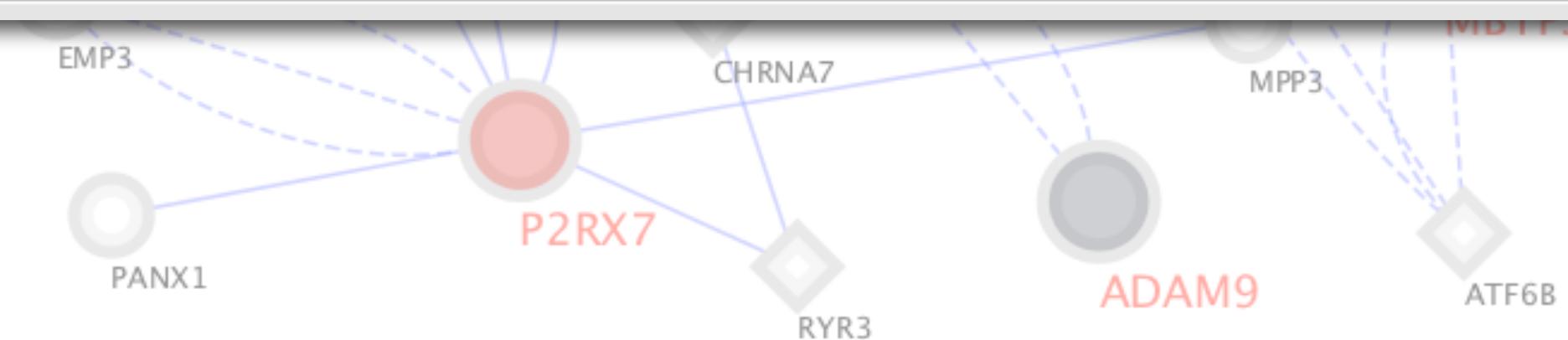
ID	networks	rank	data type	normalized max weight
CASP4 (Co expression H sapiens 2) PYCARD	[Hannenhalli-Cappola-2006]	0	Co-expression	3.365162984771238E...
IL18R1 (Co expression H sapiens 2) CARD8	[Wang-Maris-2006]	0	Co-expression	0.003545847771278...
CASP1 (Pathway H sapiens 2) IL18	[PATHWAYCOMMONS-NCI_NATURE]	2	Pathway	0.03513637205764857
IL18 (Co expression H sapiens 2) PYCARD	[Ross-Perou-2001]	0	Co-expression	0.004547999700193...
CASP5 (Co expression H sapiens 2) IL1RL2	[Wang-Maris-2006]	0	Co-expression	0.001120727970863...
IL18BP (Physical interactions H sapiens 2) IL18	[PATHWAYCOMMONS]	1	Physical interactions	0.05554565914902503
CASP4 (Pathway H sapiens 2) CASP1	[PATHWAYCOMMONS-NCI_NATURE]	2	Pathway	0.17224503150392498
CARD16 (Co expression H sapiens 2) CASP1	[Zangrando-Basso-2009]	0	Co-expression	8.76819811945632E-4
IL1RL2 (Physical interactions H sapiens 2) IL18	[PATHWAYCOMMONS]	1	Physical interactions	0.11139493885906686
CARD16 (Physical interactions H sapiens 2) CARD8	[PATHWAYCOMMONS]	1	Physical interactions	0.06700058824264392
IL18R1 (Pathway H sapiens 2) IL18	[PATHWAYCOMMONS-NCI_NATURE]	2	Pathway	0.020903080601769...
CASP4 (Physical interactions H sapiens 2) IL18	[PATHWAYCOMMONS]	1	Physical interactions	0.014599488029586...
CARD18 (Physical interactions H sapiens 2) CARD8	[PATHWAYCOMMONS]	1	Physical interactions	0.08697183339241087
CASP5 (Physical interactions H sapiens 2) CASP1	[PATHWAYCOMMONS]	1	Physical interactions	0.015267667745317...
CSTA (Co expression H sapiens 2) CASP1	[Ross-Perou-2001, Hannenhalli-Cappola-2006]	0	Co-expression	0.002196498810603...
CARD16 (Co expression H sapiens 2) PYCARD	[Zangrando-Basso-2009]	0	Co-expression	5.060063714509976E...
CASP1 (Co expression H sapiens 2) IL18	[Wang-Maris-2006]	0	Co-expression	0.001987161993368...
CSTA (Co expression H sapiens 2) IL18	[Wang-Maris-2006, Zangrando-Basso-2009, Burcz...]	0	Co-expression	0.003217535435700...

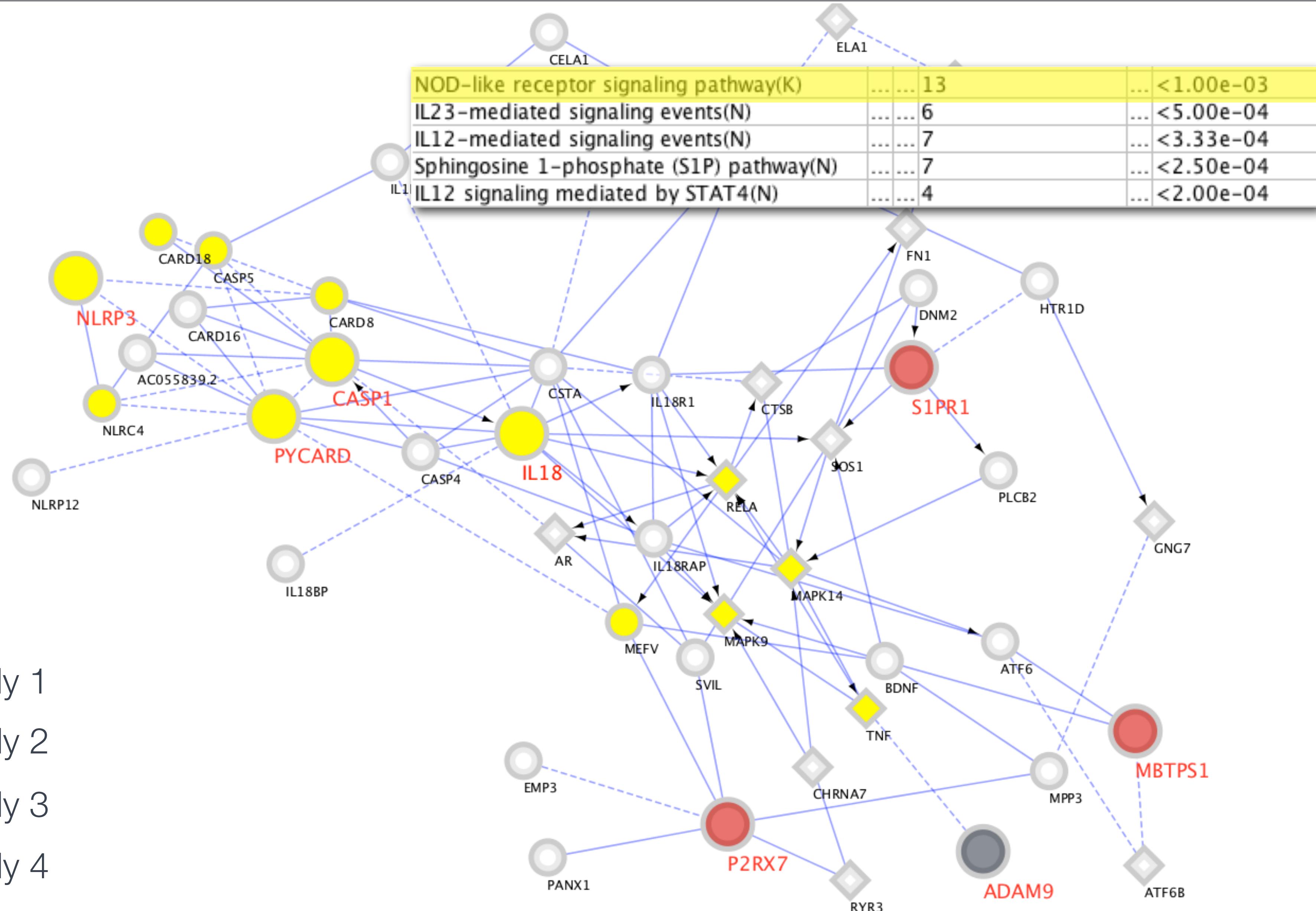


Study 3



Study 4





Next steps

PROCEEDINGS

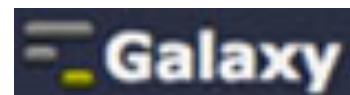
Open Access

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



There is a new version of CloudMan: [What's New](#) | [Update CloudMan](#) Info: [report bugs](#) | [wiki](#) | [screenshots](#)

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](#)

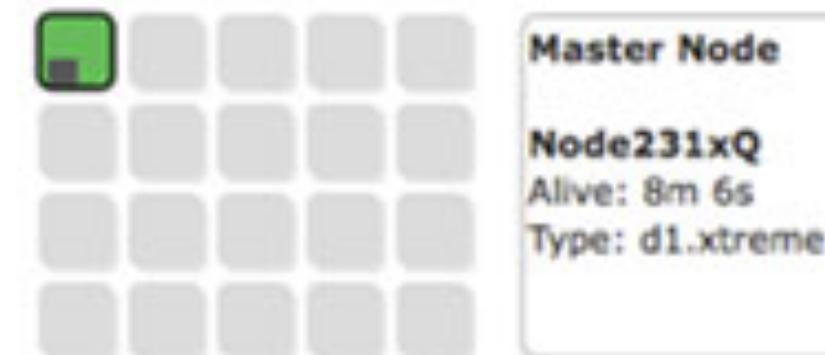
Status

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: 100px;" 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>

crw/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

0hos-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

Not everything as complicated...

- ▶ ... but sequencing is becoming a commodity



Talk to us early

Involvement in study design to optimize experiments



Contact

ohofmann@hsph.harvard.edu

@fiamh

<http://compbio.sph.harvard.edu/chb/>

