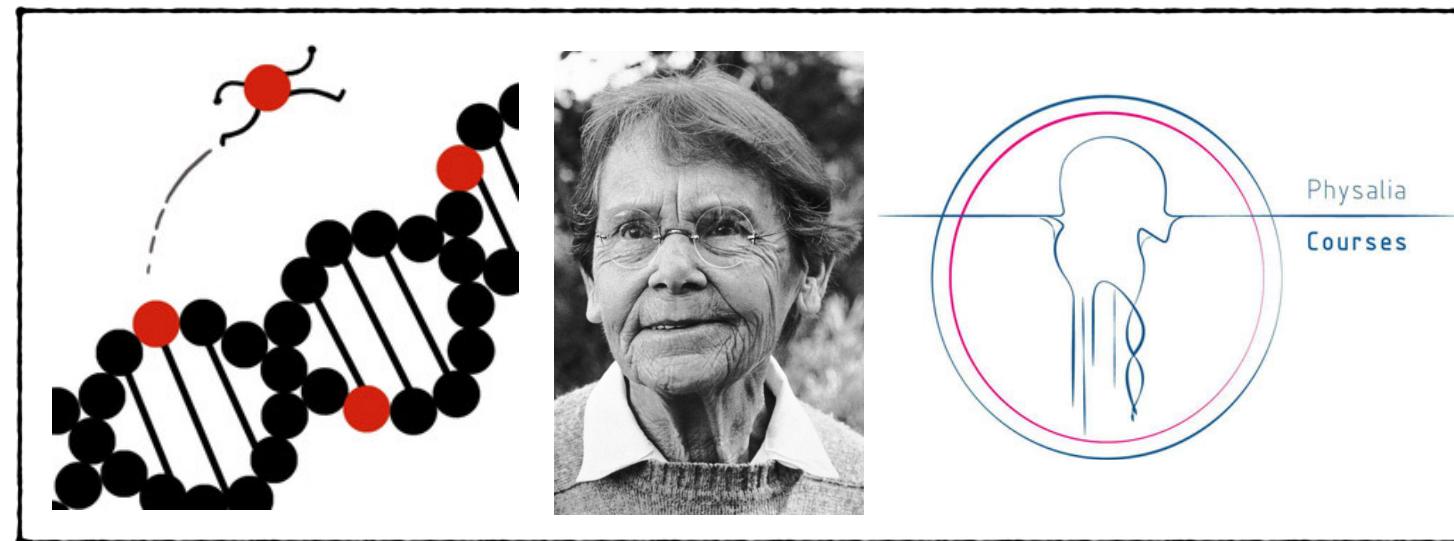


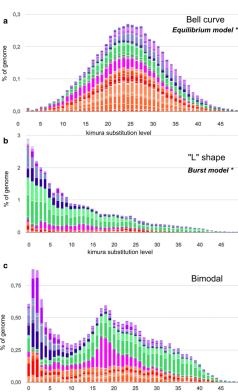
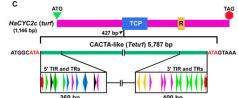
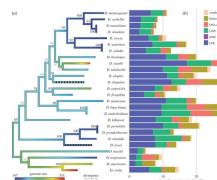
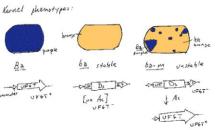
BIOINFORMATIC ANALYSIS OF TRANSPOSSABLE ELEMENTS

3rd-7th November 2025

Lecture 3B TE manual curation



Valentina Peona



What is manual curation?

Library 1.0

```
>cons1  
ATGAT  
>cons2  
GATTACA  
...
```

The automatic generation of a TE library often comes with issues that must be solved:

- Incomplete consensus sequences (e.g., boundaries are wrong)
- Multiple variants were pooled in the same consensus sequence
- Misclassified sequences

What is manual curation?

Library 1.0

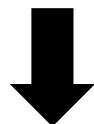
```
>cons1  
ATGAT  
>cons2  
GATTACA  
...
```

The automatic generation of a TE library often comes with issues that must be solved:

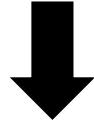
- Incomplete consensus sequences (e.g., boundaries are wrong)
- Multiple variants were pooled in the same consensus sequence
- Misclassified sequences

What can we do?

Collect TE insertions from the genome similar to our consensus sequences



Create a multi-sequence alignment for each consensus sequence and its hits



What is manual curation?

Library 1.0

```
>cons1  
ATGAT  
>cons2  
GATTACA  
...
```

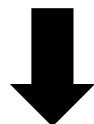
The automatic generation of a TE library often comes with issues that must be solved:

- Incomplete consensus sequences (and its boundaries are wrong)
- Multiple variants
- Misclassified sequences

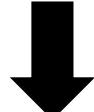
It can be a very long process (and a bit frustrating) but it is worth it!

What can we do?

Collect TE insertions from the genome similar to our consensus sequences



Create a multi-sequence alignment for each consensus sequence and its hits



Why manual curation?

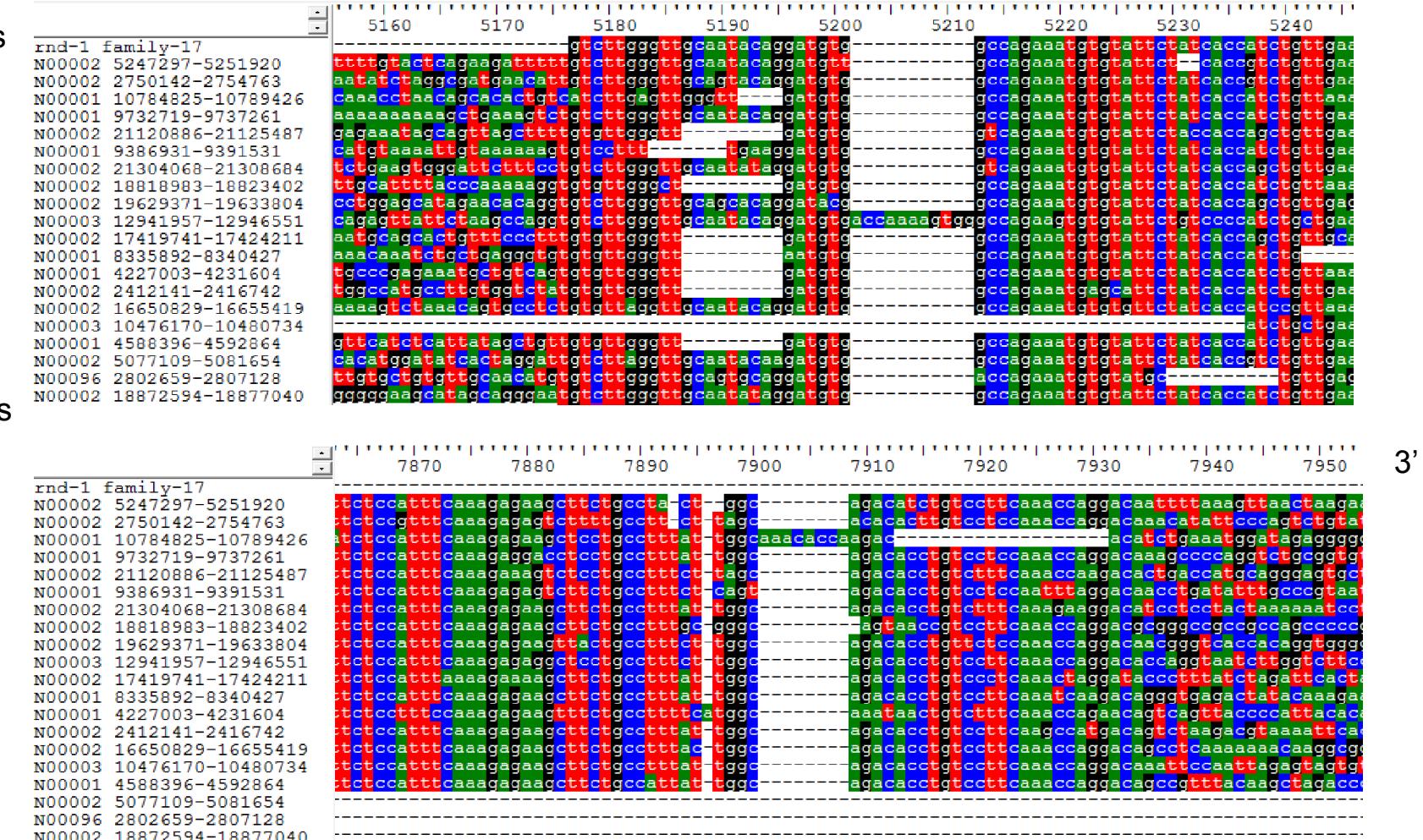
To obtain an as complete and accurate annotation as possible

- 1 Incomplete consensus sequences
 - 2 Misclassified consensus sequences
 - 3 Species-specific libraries are key
 - 4 Complete library equals complete annotation
-  Avoid biological misinterpretations

Incomplete consensus sequences

Old consensus sequence

New consensus sequence



Incomplete consensus sequences

5'

Old consensus sequence

New consensus sequence

rnd-1 family-17

N00002 5247297-5251920
N00002 2750142-2754763
N00001 10784825-10789426
N00001 9732719-9737261
N00002 21120886-21125487
N00001 9386931-9391531
N00002 21304068-21308684
N00002 18818983-18823402
N00002 19629371-19633804
N00003 12941957-12946551
N00002 17419741-17424211
N00001 8335892-8340427
N00001 4227003-4231604
N00002 2412141-2416742
N00002 16650829-16655419
N00003 10476170-10480734
N00001 4588396-4592864
N00002 5077109-5081654
N00096 2802659-2807128
N00002 18872594-1887704n

rnd-1 family-17.LTR 5ntTSD **TCTCCATTCAAAGAGAAGCTTCTGCCTTMT-TGGC** ----- **AGACACCTGTCCCTCAACCCAGGACA**

3'

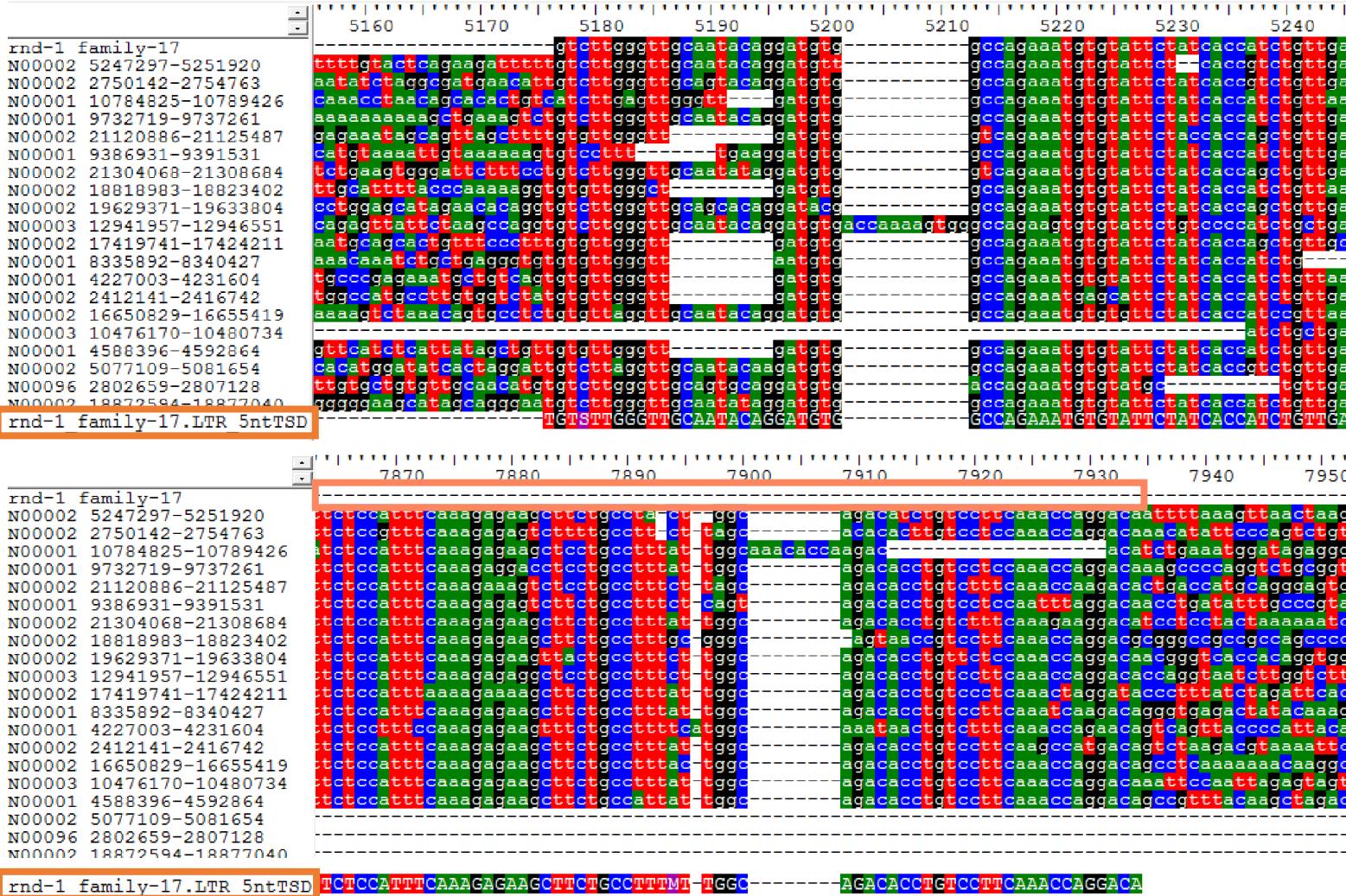
Incomplete consensus sequences

Old consensus sequence

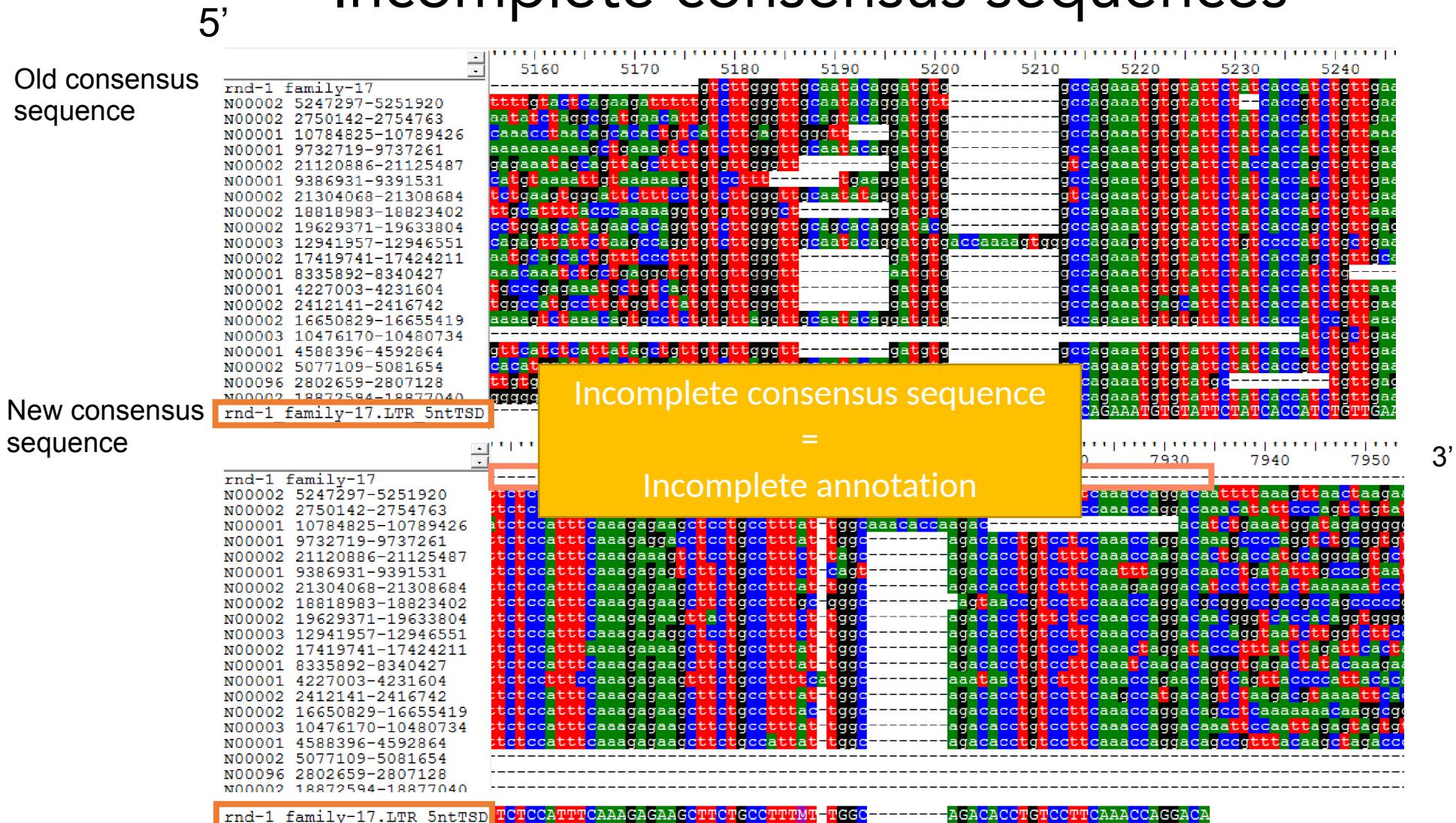
5'

New consensus sequence

3'



Incomplete consensus sequences

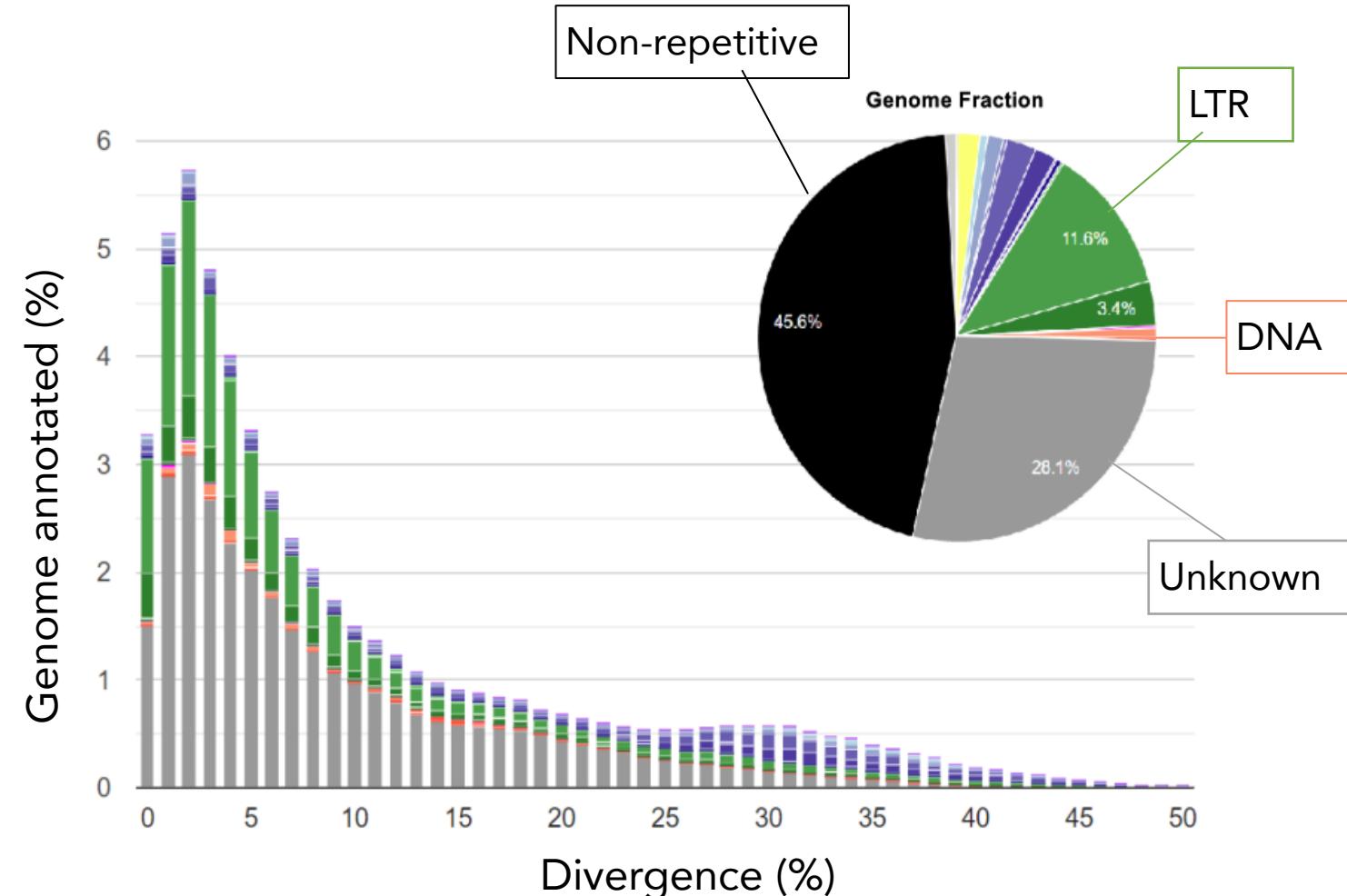


Misclassification

The case of a crustacean

Repeat library of raw consensus sequences from RepeatModeler2

TE category	# consensus sequences
DNA	114
LTR	1909
LINE	378
SINE	3
RC	10
Unknown	3774



Misclassification

The case of a crustacean

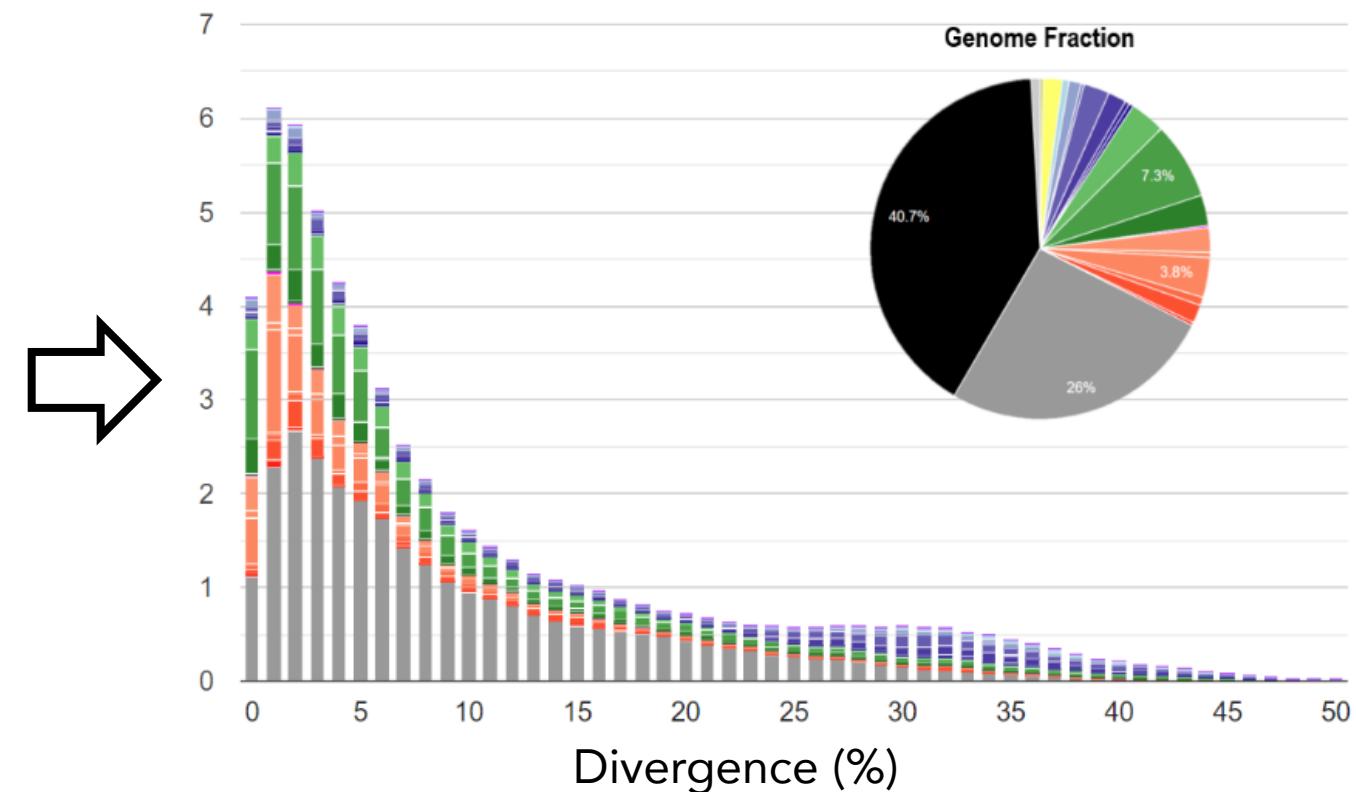
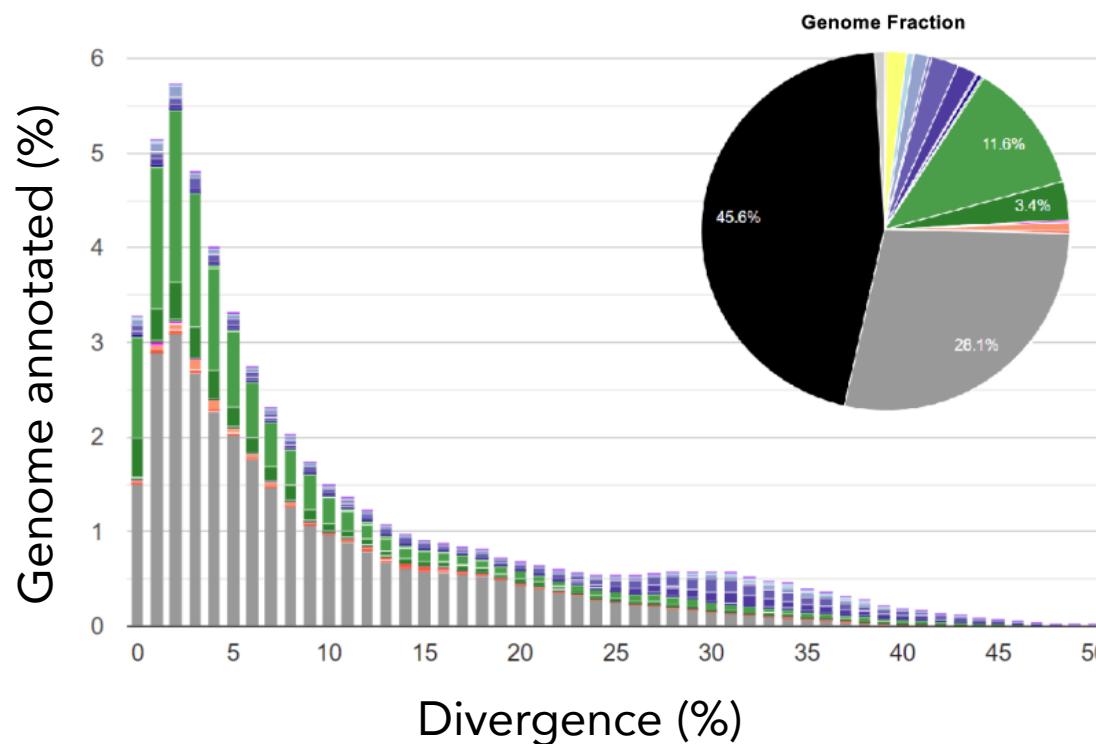
	Curated (RM2)													
Original (RM1)	Unknown	Unknown	LTR	LINE	Mariner	hAT	Merlin	Transib	P	PiggyBac	CACTA	Sola2	Total	Total (%)
	Unknown	14	1	0	14	11	1	1	1	8	4	1	56	62.22%
	LTR	6	1	0	7	1	0	0	0	17	0	0	32	35.56%
	LINE	0	0	2	0	0	0	0	0	0	0	0	2	2.22%
	Mariner	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
	hAT	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
	Merlin	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
	Transib	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
	P	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
	PiggyBac	0	0	0	0	0	0	0	0	0	0	0	0	0.00%
Total		20	2	2	21	12	1	1	1	25	4	1	90	
Total (%)		22.22%	2.22%	2.22%	23.33%	13.33%	1.11%	1.11%	1.11%	27.78%	4.44%	1.11%		
Key: Unknown LTR LINE DNA														

- 32 consensus sequences preliminary classified as LTR by RepeatModeler were curated
 - Only one of them was confirmed to be an LTR
 - Most of them were DNA transposons!
 - Many more remain in the library to be curated...!
 - My solution: label all LTRs as unknown unless there are protein-based evidence that indicate LTR
- Demirci, Peona, Sun et al. unpublished

Annotation improvements

The case of a crustacean

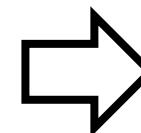
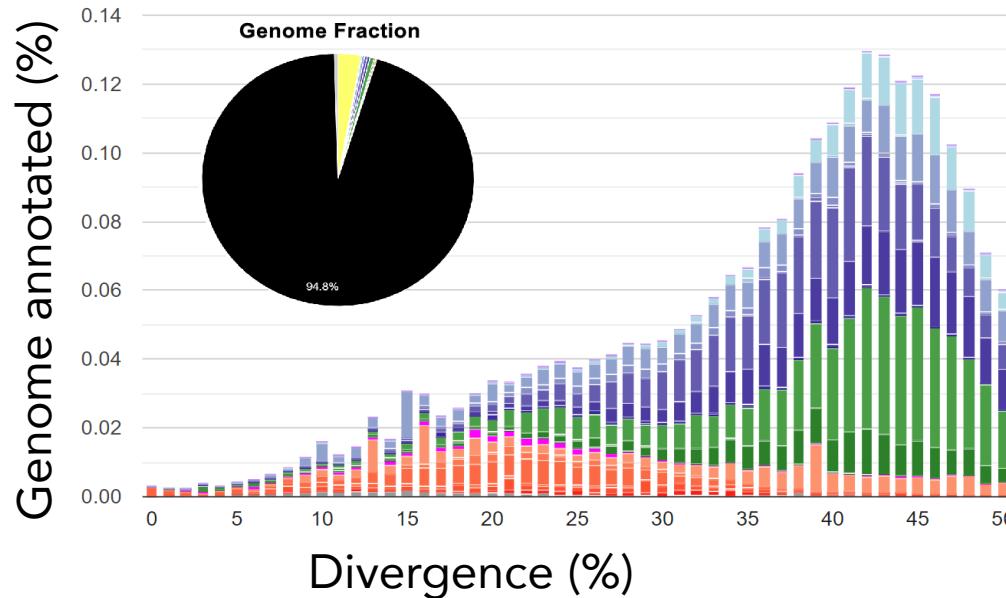
After curating 100 consensus sequences



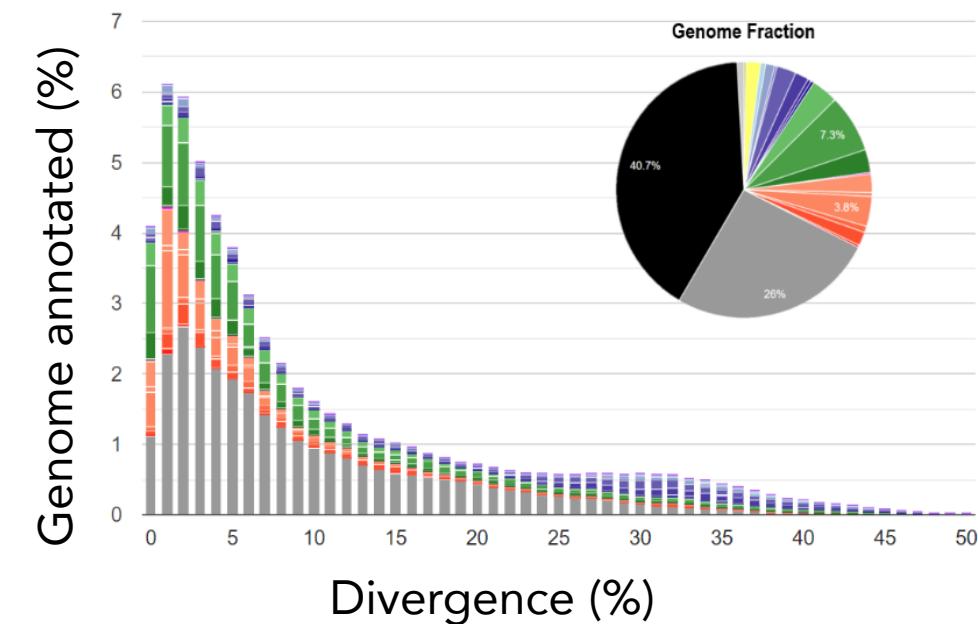
↓ LTR ↑ DNA

Species-specific libraries

Library from database

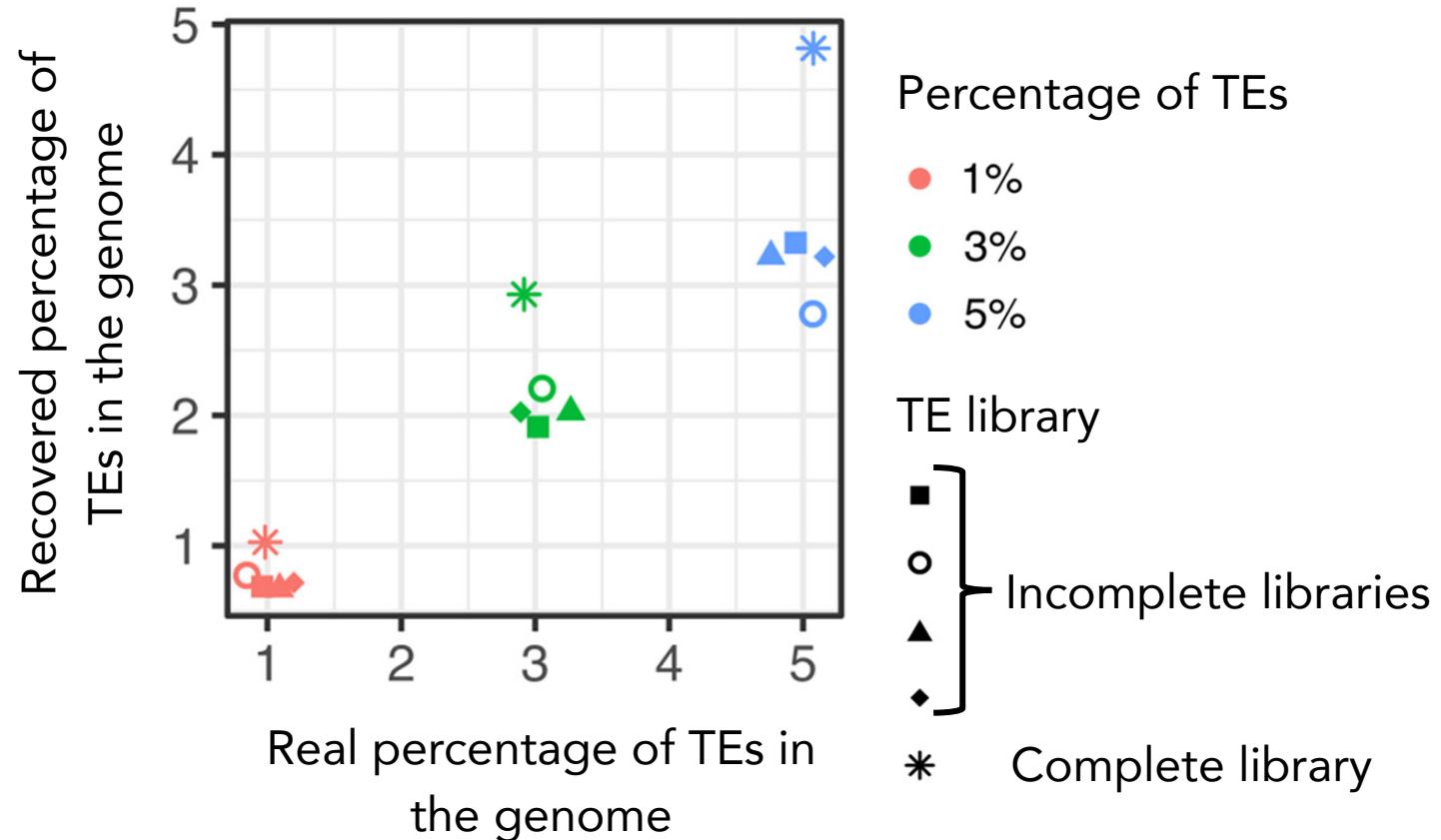


Species-specific library

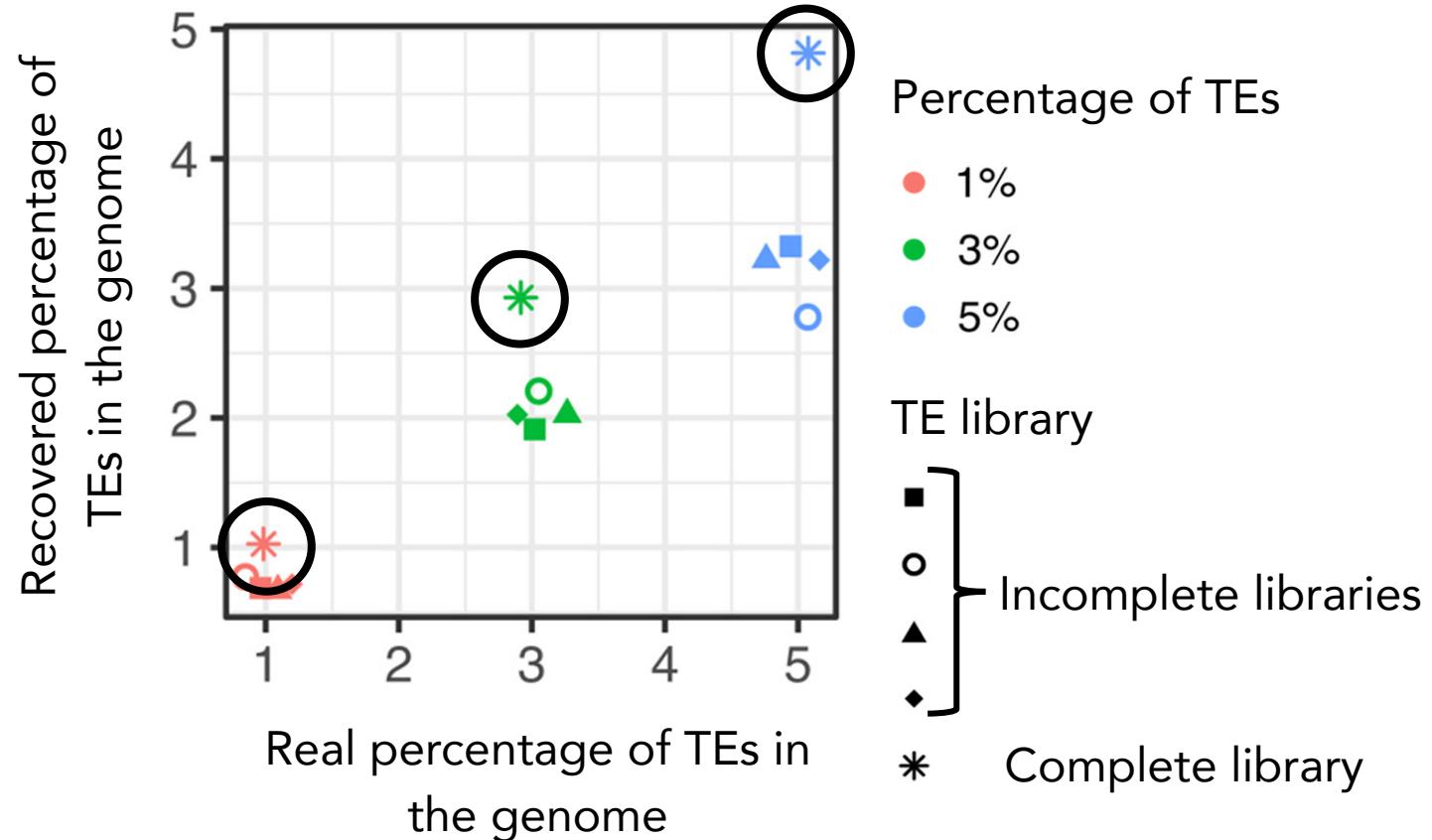


- Overall percentage
- Different accumulation distribution

You have it, you find it!



You have it, you find it!



Some problems of manual curation

- 1 Time consuming
- 2 Difficult to scale up to the number of genomes we have now available
- 3 Always a bit of subjectivity that hinders complete reproducibility

Time consuming: prioritise!

.divsum file can be very useful!

Jukes/Cantor and Kimura substitution levels adjusted for CpG sites

=====
File: haploid_assembly_BUSCOmancur.fasta.align

Weighted average Kimura divergence for each repeat family

Class	Repeat	absLen	wellCharLen	Kimura%
-------	--------	--------	-------------	---------

-----	-----	-----	-----	-----
-------	-------	-------	-------	-------

LTR/Gypsy	ltr-1-1295	3061696	3041108	3.78
-----------	------------	---------	---------	------

LTR/Gypsy	ltr-1-1269	3033238	3016481	4.16
-----------	------------	---------	---------	------

LTR/Gypsy	ltr-1-1589	2885608	2854985	5.77
-----------	------------	---------	---------	------

Unknown	balImpl-86	2847824	2819143	8.79
---------	------------	---------	---------	------

LTR/Gypsy	ltr-1-205	2710211	2684973	3.89
-----------	-----------	---------	---------	------

LTR/Unknown	ltr-1-1379	2496705	2463106	4.88
-------------	------------	---------	---------	------

LTR/Gypsy	ltr-1-1115	2478128	2460958	3.11
-----------	------------	---------	---------	------

Unknown	balImpl-19	2439731	2433087	2.07
---------	------------	---------	---------	------

Unknown	balImpl3-158	2358241	2334267	5.35
---------	--------------	---------	---------	------

LTR/Gypsy	ltr-1-1358	2123961	2083228	6.42
-----------	------------	---------	---------	------

Time consuming: prioritise!

.divsum file can be very useful!

Jukes/Cantor and Kimura substitution levels adjusted for CpG sites
=====

File: haploid_assembly_BUSCOmancur.fasta.align

Weighted aver	Class	Repea	-----	-----	repeat family	%
LTR/Gypsy						
LTR/Gypsy						
LTR/Gypsy						
Unknown balIm						
LTR/Gypsy						
LTR/Unknown	ltr-1-1379		2496705	2463106	3.78	
LTR/Gypsy	ltr-1-1115		2478128	2460958	4.16	
Unknown balImpl-19		2439731	2433087	2.07		
Unknown balImpl3-158		2358241	2334267	5.35		
LTR/Gypsy	ltr-1-1358		2123961	2083228	6.42	

Look for abundant repeats and/or young repeats

Or focus on categories you are already interested in

Scalability

SPECIAL SECTION

ZOONOMIA

RESEARCH ARTICLE

ZOONOMIA

Insights into mammalian TE diversity through the curation of 248 genome assemblies

Austin B. Osmanski¹, Nicole S. Paulat¹, Jenny Korstian¹, Jenna R. Grimshaw¹, Michaela Halsey¹, Kevin A. M. Sullivan¹, Diana D. Moreno-Santillán¹, Claudia Crookshanks¹, Jacquelyn Roberts¹, Carlos Garcia¹, Matthew G. Johnson¹, Llewellyn D. Densmore¹, Richard D. Stevens², Zoonomia Consortium[†], Jeb Rosen³, Jessica M. Storer³, Robert Hubley³, Arian F. A. Smit³, Liliana M. Dávalos^{4,5}, Elinor K. Karlsson^{6,7}, Kerstin Lindblad-Toh^{7,8,9}, David A. Ray^{1*}

Curated two repeat libraries
~40 participants (two editions)

Employed hundreds of students

Peona et al. *Mobile DNA* (2024) 15:10
<https://doi.org/10.1186/s13100-024-00319-8>

Mobile DNA

RESEARCH

Open Access



Teaching transposon classification as a means to crowd source the curation of repeat annotation – a tardigrade perspective

Valentina Peona^{1,2,3*†}, Jacopo Martelossi^{4*†}, Dareen Almojil⁵, Julia Bocharkina⁶, Ioana Brännström^{7,8}, Max Brown⁹, Alice Cang¹⁰, Tomàs Carrasco-Valenzuela^{11,12}, Jon DeVries¹³, Meredith Doellman^{14,15}, Daniel Elsner¹⁶, Pamela Espíndola-Hernández¹⁷, Guillermo Friis Montoya¹⁸, Bence Gaspar¹⁹, Danijela Zagorski²⁰, Paweł Halakuc²¹, Beti Ivanovska²², Christopher Laumer²³, Robert Lehmann²⁴, Ljudevit Luka Boštjančić²⁵, Rahia Mashhoodh²⁶, Sofia Mazzoleni²⁷, Alice Mouton²⁸, Maria Anna Nilsson²⁵, Yifan Pei^{1,29}, Giacomo Potente³⁰, Panagiotis Provatas³¹, José Ramón Pardos-Blas³², Ravindra Raut³³, Tomasa Sbaffi³⁴, Florian Schwarz³⁵, Jessica Stapley³⁶, Lewis Stevens³⁷, Nusrat Sultana³⁸, Radka Symonova³⁹, Mohadeseh S. Tahami⁴⁰, Alice Urzi⁴¹, Heidi Yang⁴², Abdullah Yusuf⁴³, Carlo Pecoraro⁴⁴ and Alexander Suh^{1,45,46*}

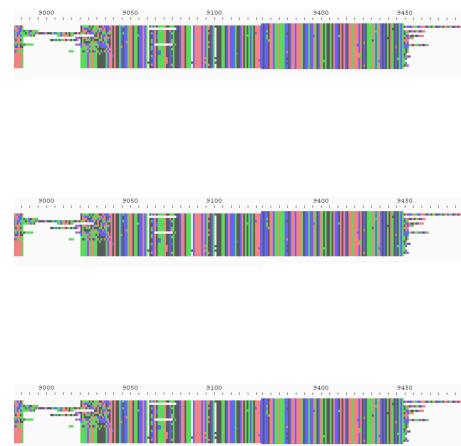
How to curate a repeat library

Overview

Library 1.0

```
>cons1  
ATGAT  
>cons2  
GGATC  
...
```

Align to the reference genome assembly
→
Collect the best hits + 2kb flanks and get a multisequence alignment for each consensus sequence



Library 2.0

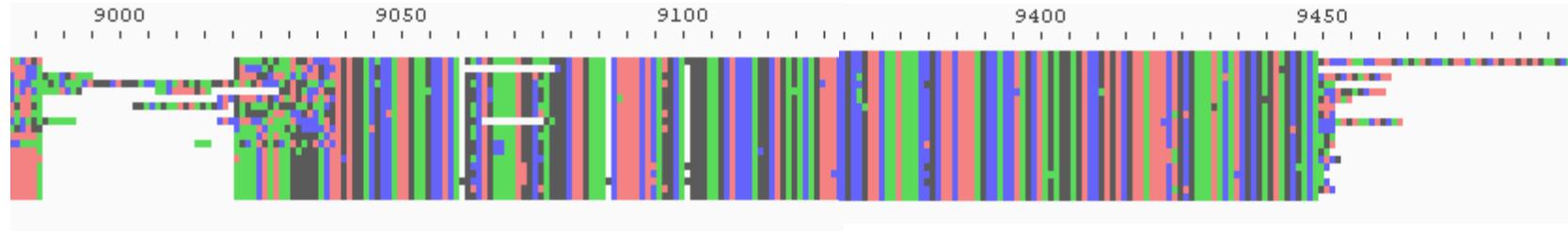
```
>cons1  
ATGAC  
>cons2  
GCATC  
...
```

Get refined consensus sequences
→
DOCUMENT THE PROCESS!!



What to look for and fix during curation

Part 1



- 0 Prepare the alignments and the Excel file to document every step of the curation
- 1 Create a new consensus sequence out of the alignment, e.g., Advanced Consensus Maker, built-in options in your alignment tool, ClAlign
- 2 Establish (temporary) boundaries of the new consensus sequence
- 3 Go through the alignment and fix ambiguous bases and gaps by applying a majority rule
- 4 Pay attention to CpG sites
- 5 Find the correct boundaries + TSDs of the new consensus sequence – remove TSDs from consensus!!

Alignments and

documentation of the library

Align and extend process

RMDL_pipeline.pl + run_RMDL_pipeline.sh

- BLAST to align consensus sequence to reference genome
- Collect up to best 20 hits + flanking regions
- Align the best hits to the consensus sequence with MAFF

The code and citations can be found in the Github repository and on Google Drive (Practical2_part2)

Empty template in the Google Drive folder:

[https://docs.google.com/spreadsheets/d/1FGPWUsdlibTqAK3GVIt3A0ZSXXzJp0XO/edit?
usp=drive_link&ouid=116433763462028904068&rtpof=true&sd=true](https://docs.google.com/spreadsheets/d/1FGPWUsdlibTqAK3GVIt3A0ZSXXzJp0XO/edit?usp=drive_link&ouid=116433763462028904068&rtpof=true&sd=true)

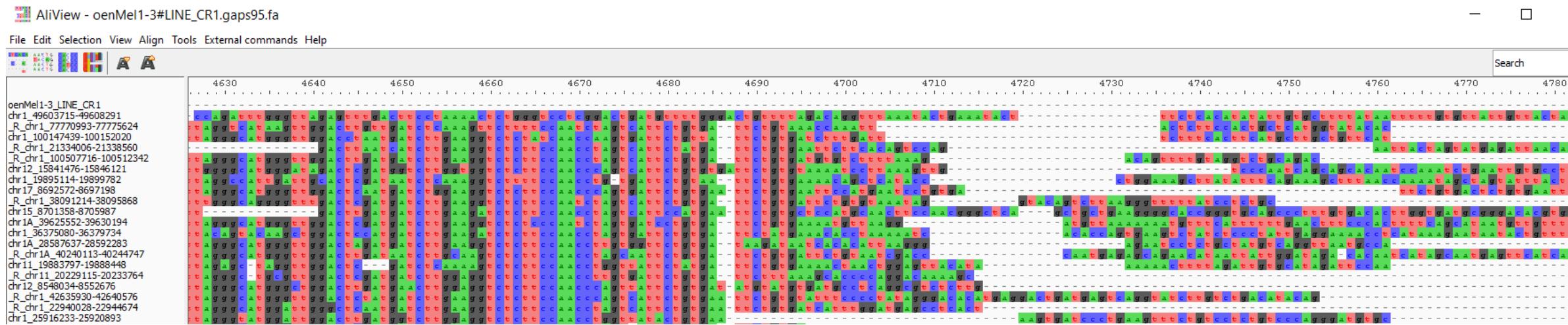
Create a new consensus sequence

Example: Advanced Consensus Maker

The screenshot shows the 'HIV sequence database' homepage with a navigation bar at the top. Below it, the 'Advanced Consensus Maker' tool is displayed. The 'Input' section contains a text area for pasting alignments, a file upload field, and a 'Squeeze gaps in input' option. The 'Consensus-by-Block Options' section includes settings for making consensus for each block (radio button 'No' selected), making a consensus of consensuses (radio button 'No' selected), showing the number of sequences in consensus (radio button 'No' selected), and specifying the minimum number of sequences per block (input field '3'). The 'Consensus Calculation Options' section covers unanimous values (radio button '0.75' selected), majority thresholds (radio button 'Use T if most common character is below Min fraction' selected with 'Min fraction 0.50'), break ties (radio button 'IUPAC characters' selected), character sets for consensus (radio button 'Specify character set' selected with 'AGCTU-' nucleotides and 'ARNDCQGHILKMFPSTWV-' amino acids), and alignment length handling (radio button 'Or check here to use any (all) characters' selected). The 'Output Options' section includes 'Show both consensus + alignment' (radio button 'Yes' selected), 'Output format' (radio button 'Pretty' selected), and a 'Run' button at the bottom.

- Go to the website: <https://www.hiv.lanl.gov/content/sequence/CONSENSUS/AdvCon.html>
- Paste your alignment
- Set parameters (example in figure)
- Get new consensus sequence and copy-paste it back into your alignment

Create a new consensus sequence

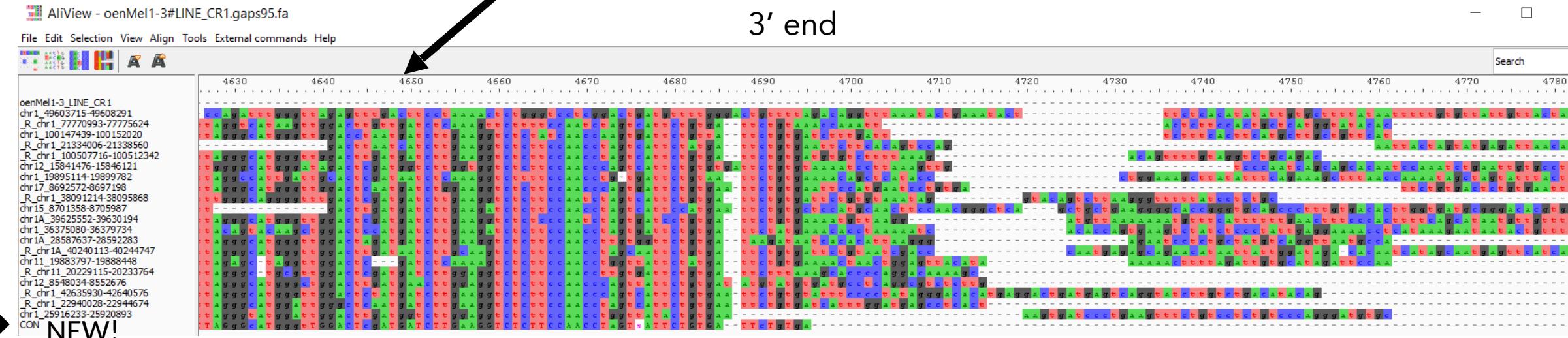


Create a new consensus sequence



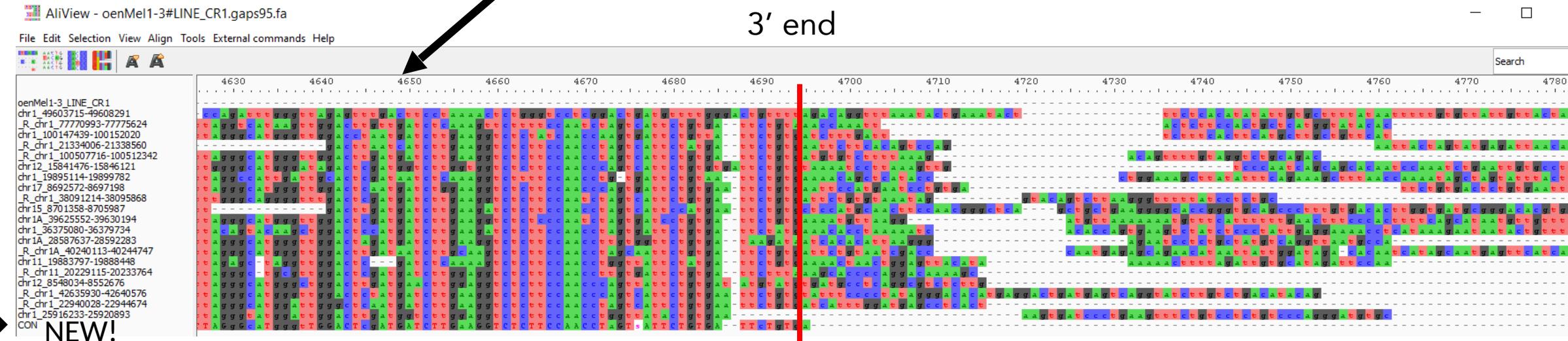
Find the termini of the consensus sequence

Original consensus is incomplete



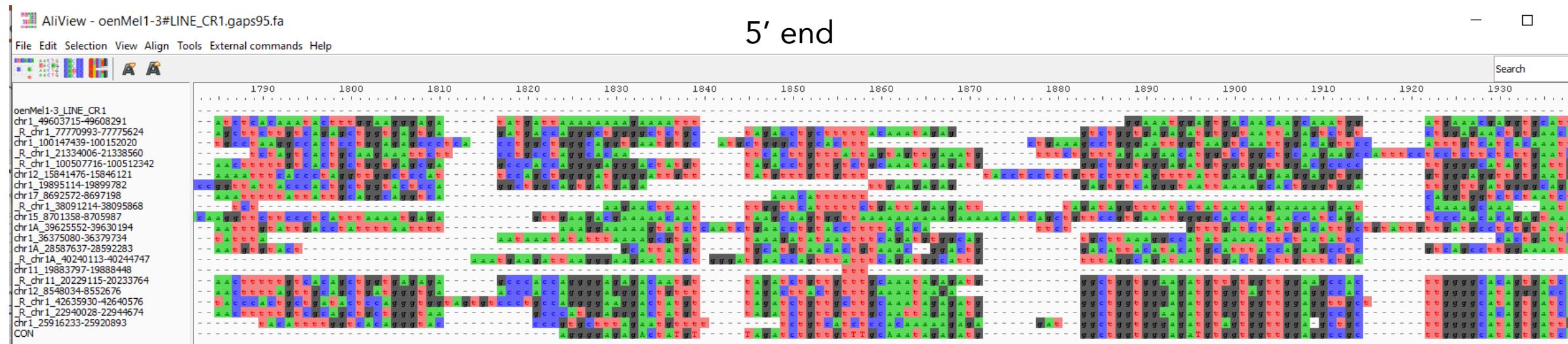
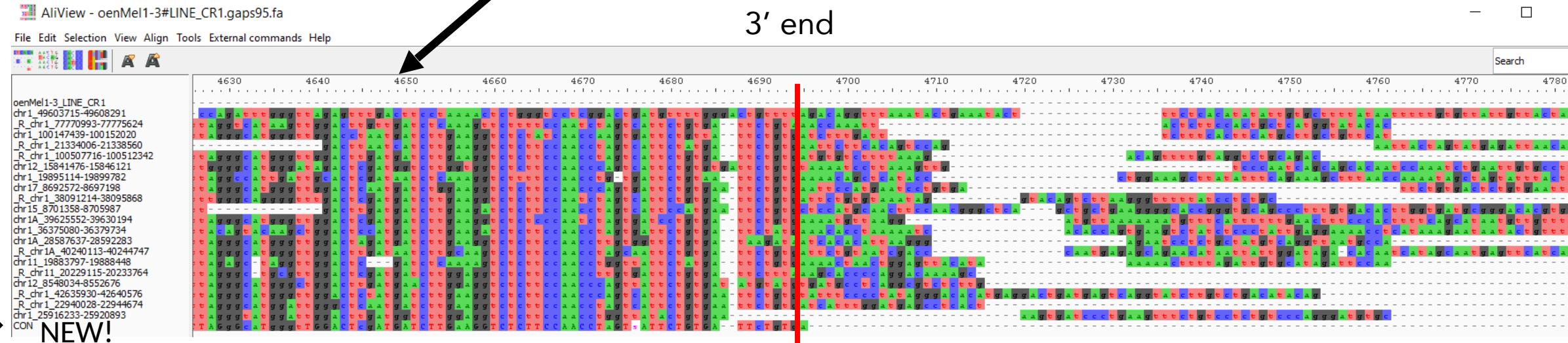
Find the termini of the consensus sequence

Original consensus is incomplete



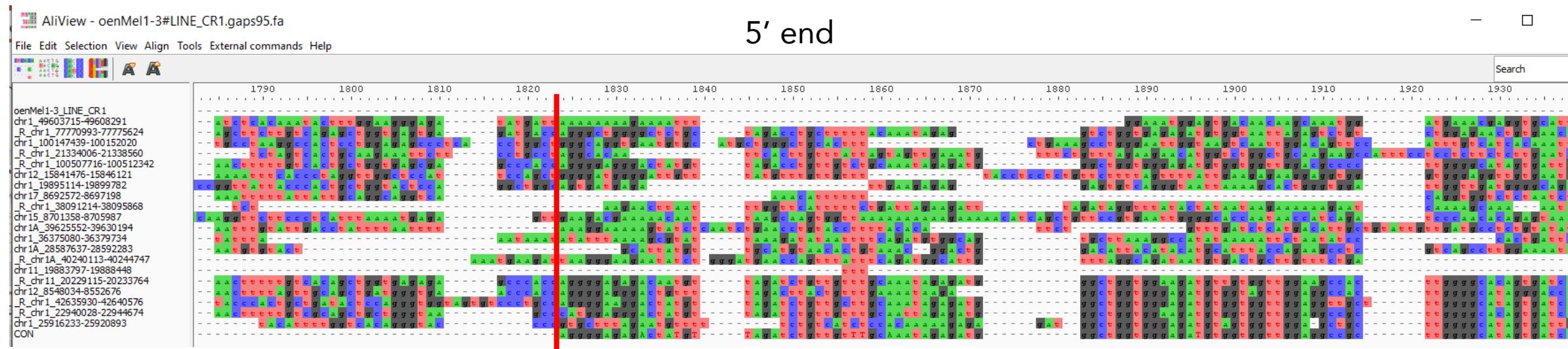
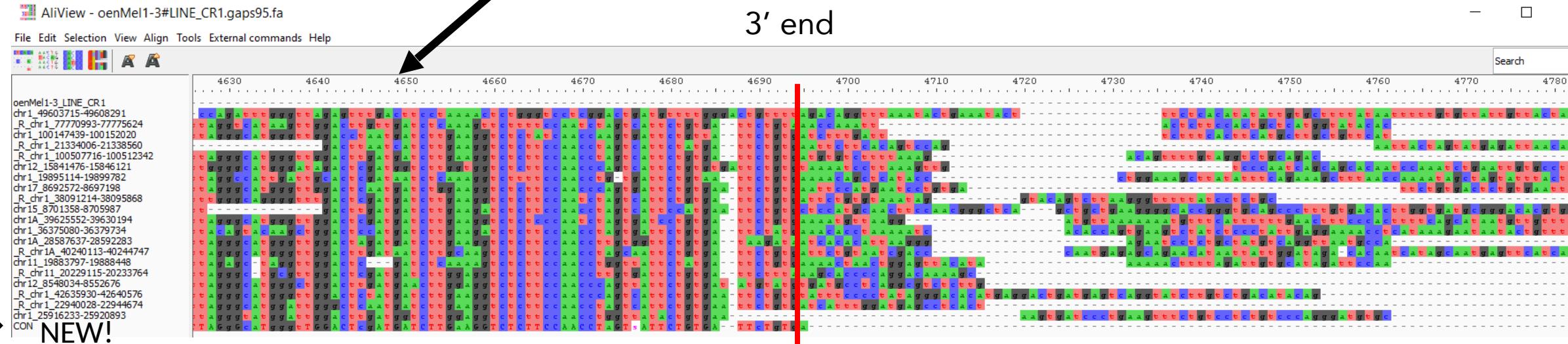
Find the termini of the consensus sequence

Original consensus is incomplete



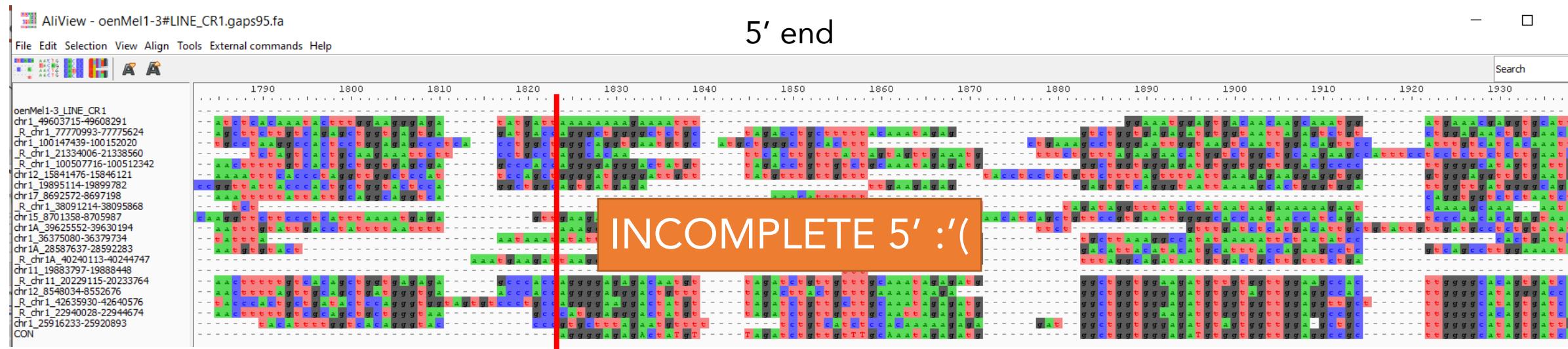
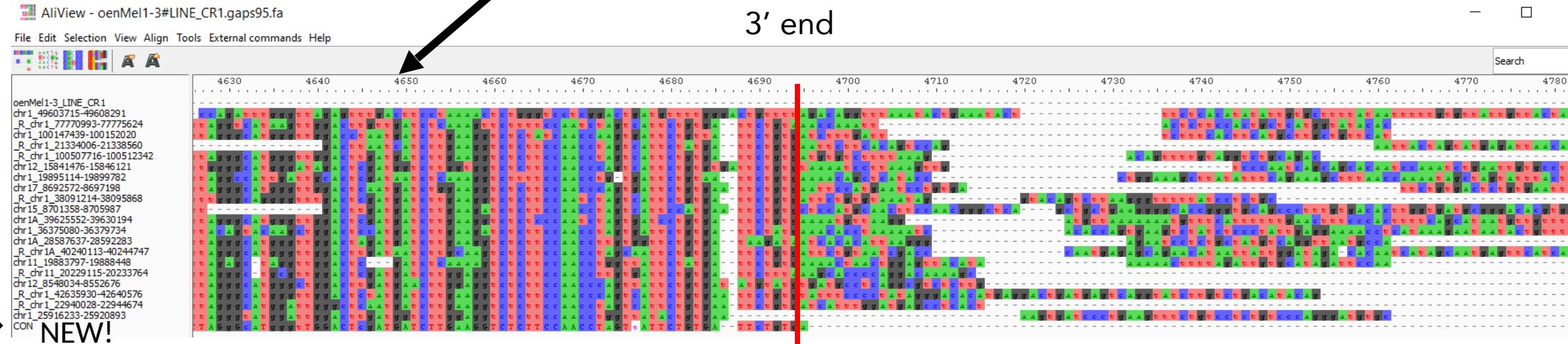
Find the termini of the consensus sequence

Original consensus is incomplete



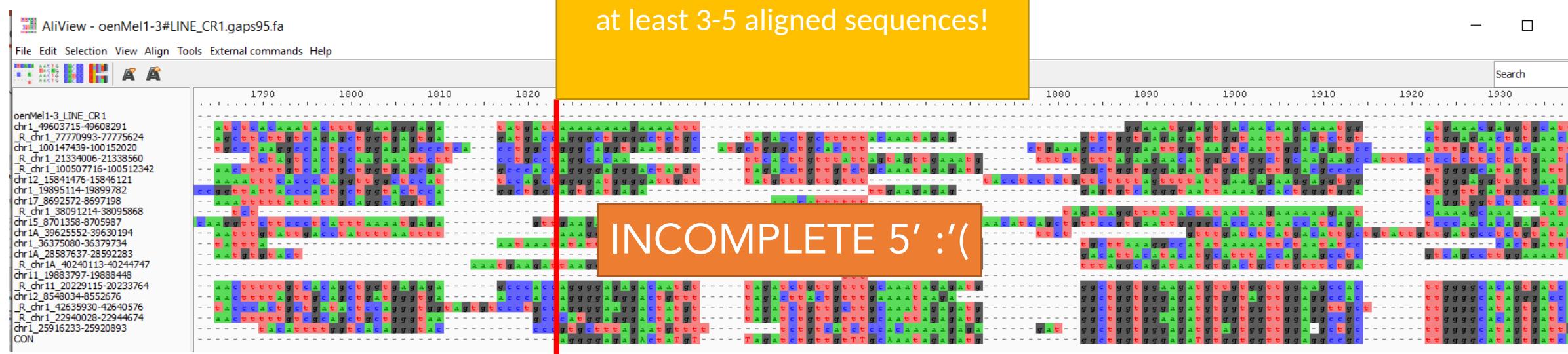
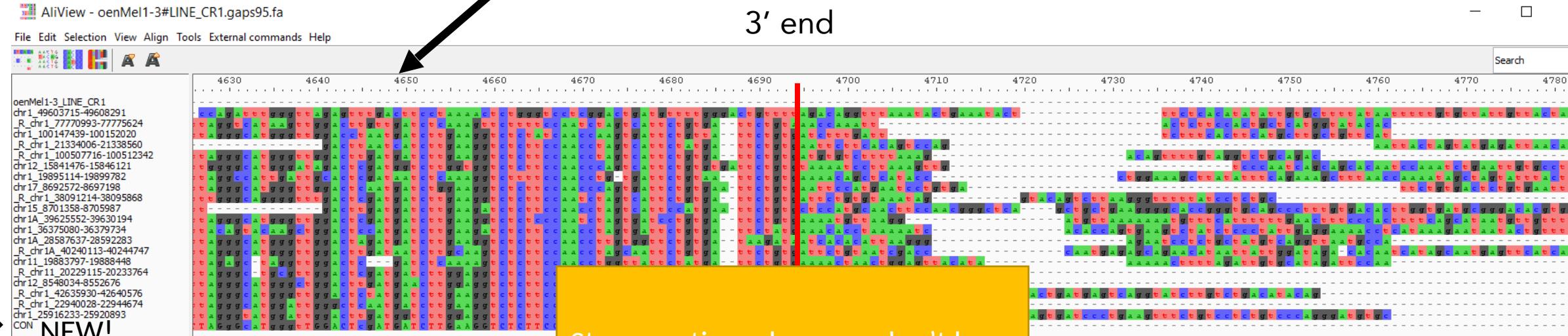
Find the termini of the consensus sequence

Original consensus is incomplete



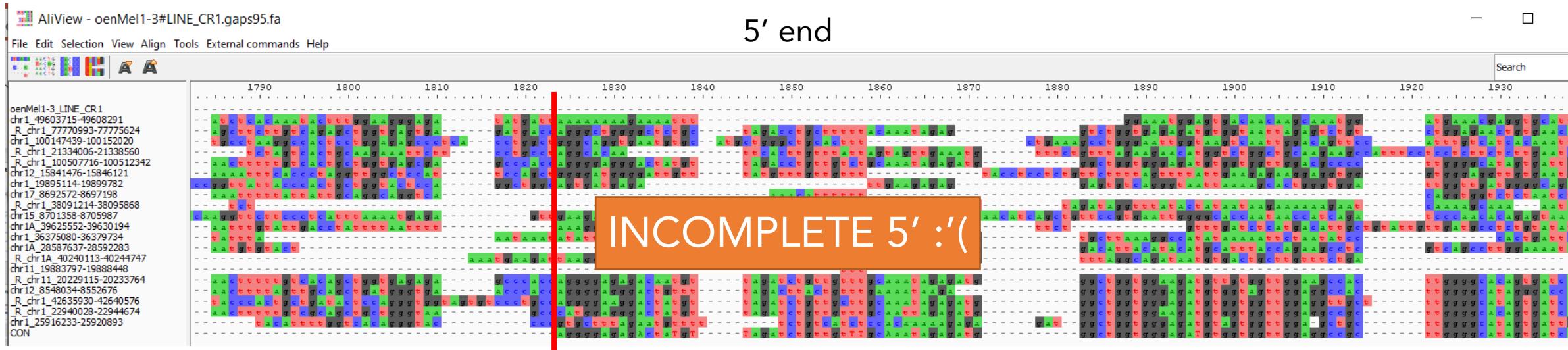
Find the termini of the consensus sequence

Original consensus is incomplete



Find the termini of the consensus sequence

Example: 5' truncation



Run another round of alignments by chopping ~500 bp at the 5' end and extend the consensus

Annotate that the consensus sequence is 5' incomplete (column G) and you can add '.inc' to the name before classification tag

Note: it is very difficult to find full-length LINE elements to build a complete consensus

Majority-rule consensus on SNPs

Raw consensus

1020 11030 11040 11050 11060 11070 11080 11090 111
-----gtcqgcqctgaccqcttggccqcccggngggcctggcaqqcqacttgcgtgtccggtaaggggac-----
caaccacccttgcgtqctgaccqcttggccqcccggagggcctggcaqqcqacttgcgtgtcccaqtcaqaggactgg
-----aqtactgtcacqctgaccqcttggctqcccgggtggcctggcaqqtqactaqtcaqqtgtcccaqtcaqggggac-----
caagacattctgtcacqctgaccqcttggctqcccgggtggcctggcaqqtqactaqtcaqqcgtcccaqtcaqggggac-----
cccac-----cctgtcacqctgaccqcttggccqcccgggtggcctggcaqqtqactaqtcaqgtgtcccaqtcaqggggac-----
cccac-----cctgtcgcqctgaccqcttggccqcccgggtggcctggcaqqtqacttgcgtqgtgtcccaqtcaqggggac-----
-----tgtgtcgcqatgaccacttggccgctgggtggccggcaqqcqacttgcgtgtcccaqtcaqggggac-----
-----tgtgtcgcqctgaccqcttggctqcccgggtggcctggcaqqcqacttgcgtgtcccaqtcaqggggac-----
gtcc-----tgtgttgcgctgatcgcttggccqcccgggtggccggcaqqcqacttgcgtgtcccaqtcaqggggat-----
cccac-----cctgtcgcactgaccqcttggccgggtggcctggcaqqcqacttgcgtgtcccaqtcaqggggac-----
cggacaaaacatgtcacqctgaccqcttggccqcccgggtggcctggcaqqtqacttgcgtgtcccaqtcaqggggac-----
-----ctttgtgttgcgctgaccqcttggccqcccggagggcctggcaqqccacttgcgtgtcccaqtcaqggggactgg-----
-----aaaactgtcacqctgaccqcttggccqcccgggtggcctggcaqqtqactacacttcgtgtcccaqtcaqggggat-----

New consensus from Adv Cons

?ac----?cTGTGCTGGCCGCCCCGGTGGCCTGGCAGGYEACTtGCTctGGTGTCCCAGTCAGGGGAC

Curated consensus sequence

-----TGTCTGGGCTGACCCGCTTGGCCGCCGGTGGCCTGGCAGGGCGACTTGCTCTGGTGTCCCAGTCAGGGGAC-----

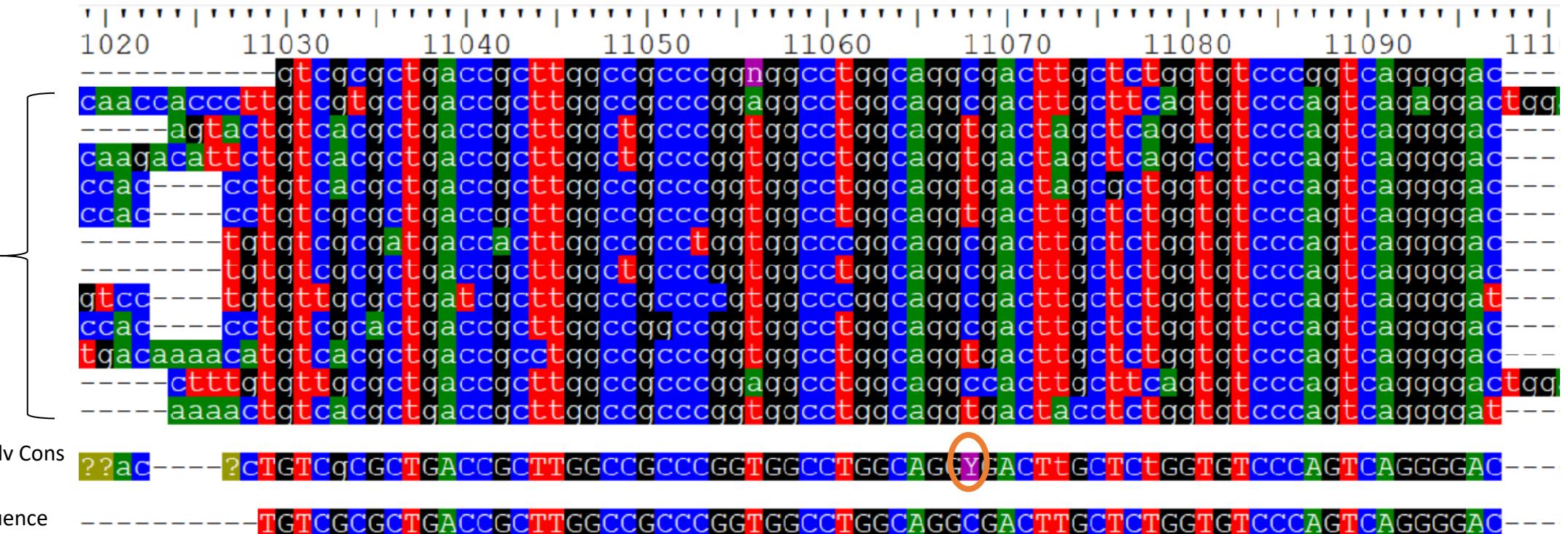
Majority-rule consensus on SNPs

Raw consensus

Copies with 2kb flank

New consensus from Adv Cons

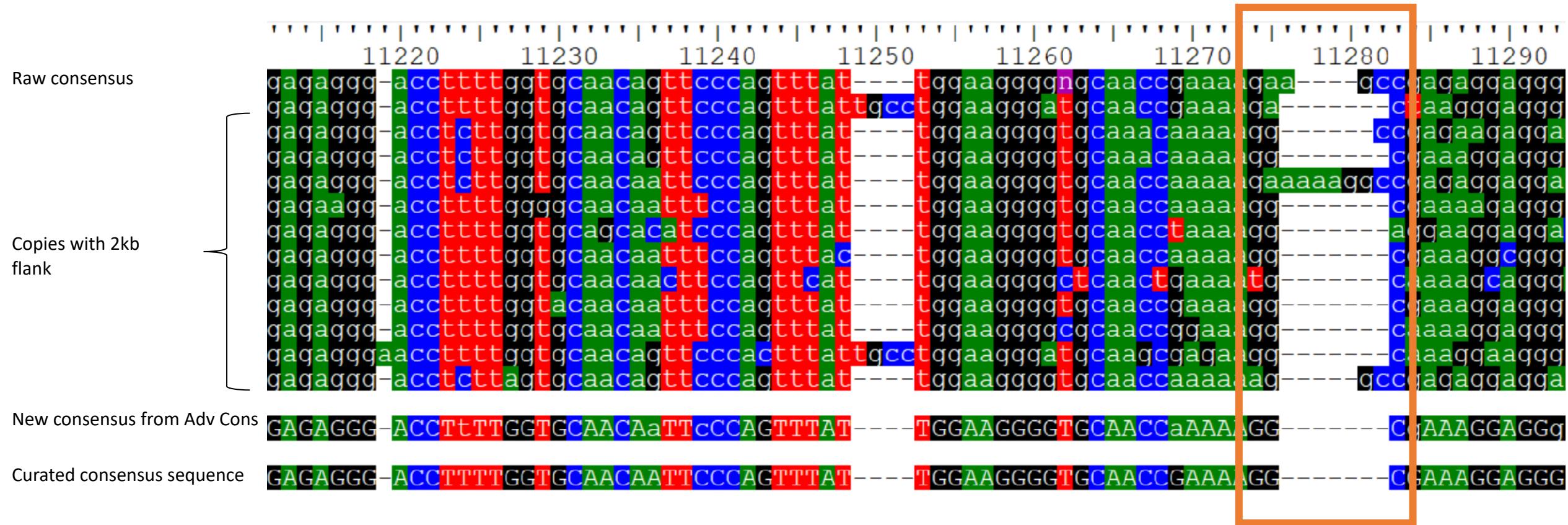
Curated consensus sequence



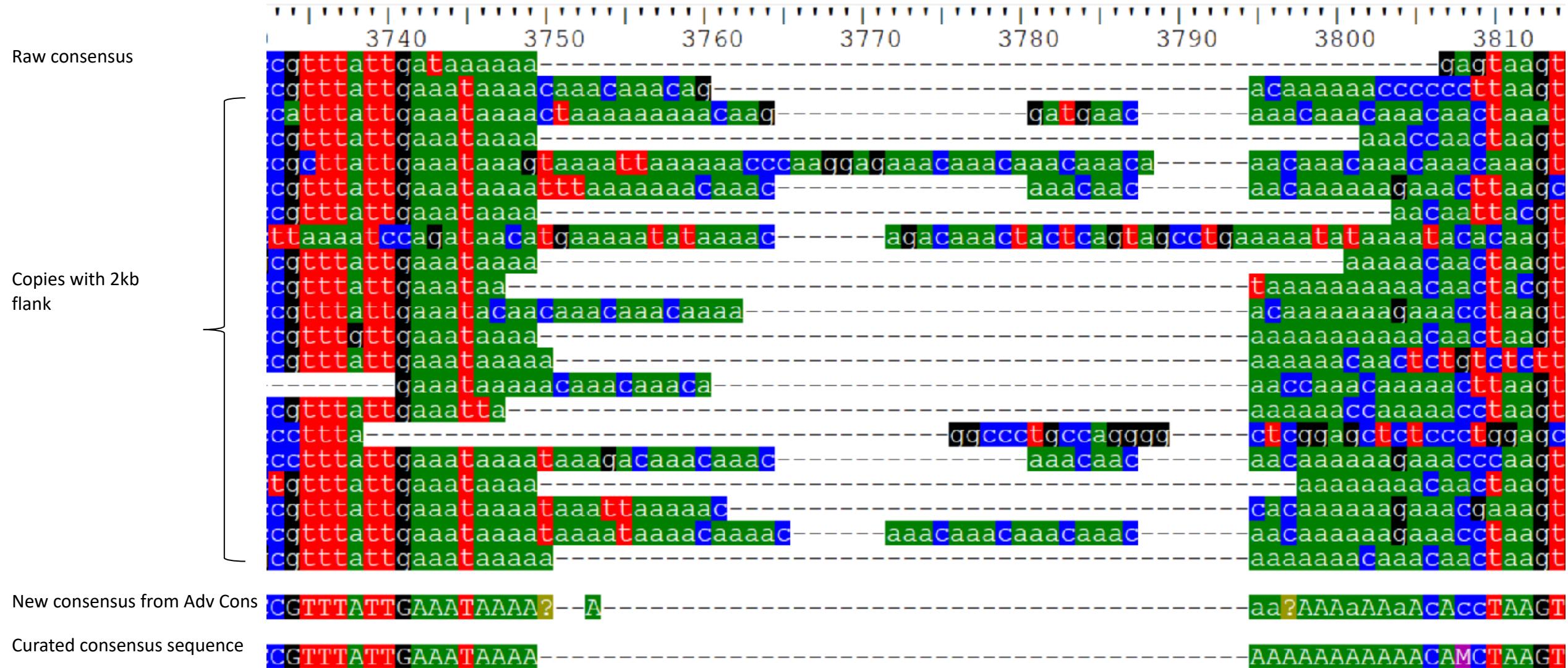
IUPAC Code	Meaning
M	A/C
R	A/G
W	A/T
S	C/G
Y	C/T
K	G/T
V	A/C/G
H	A/C/T
D	A/G/T
B	C/G/T

Get to the most accurate IUPAC character

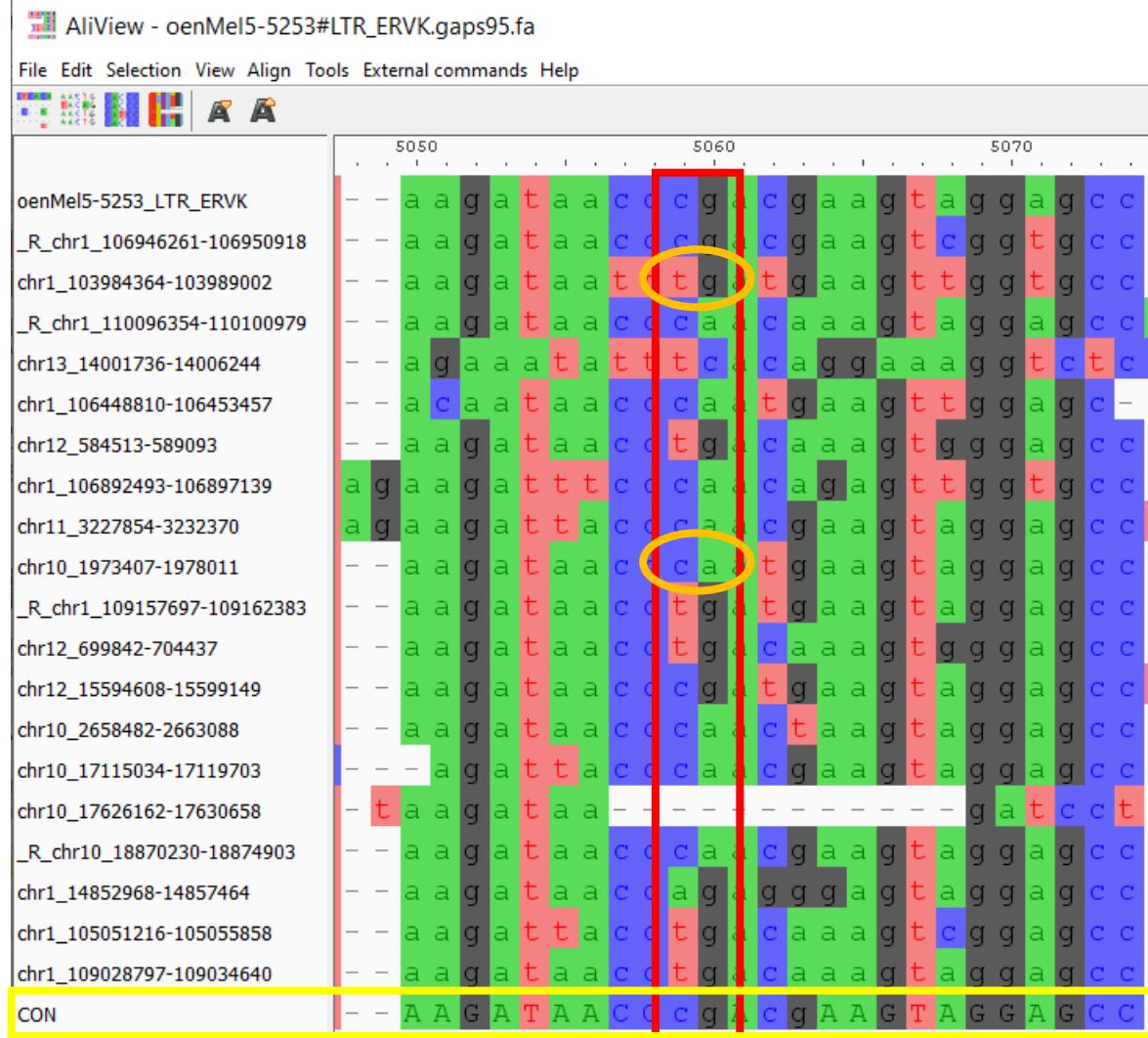
Majority-rule consensus on indels



Tricky regions



Attention to CpG sites



- You can choose to reconstruct the ancestral state CG
- I suggest to do it so to not inflate the divergence from consensus

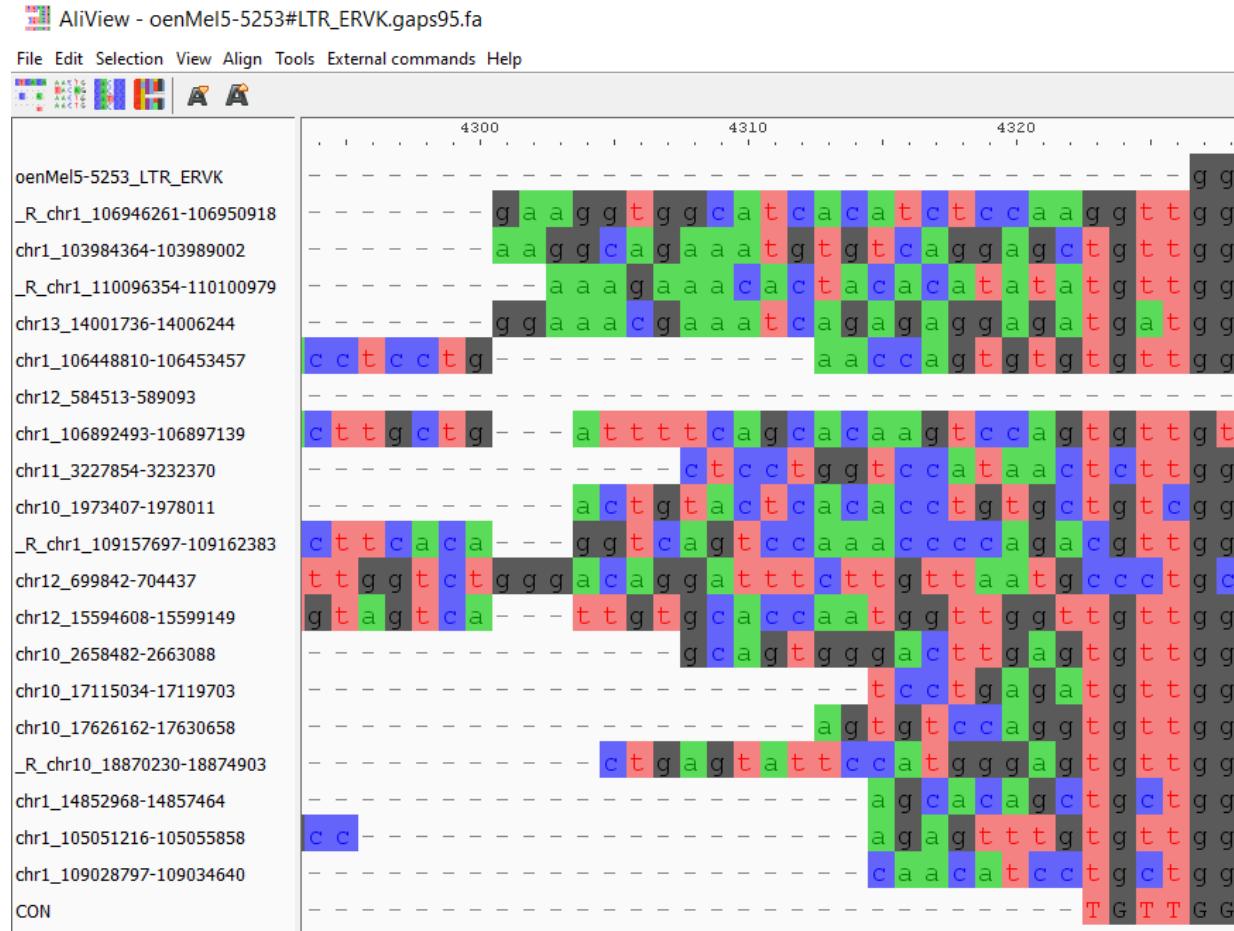
Find the TSDs

5' end

(if any!)

Start looking from the 3' terminus

3' end



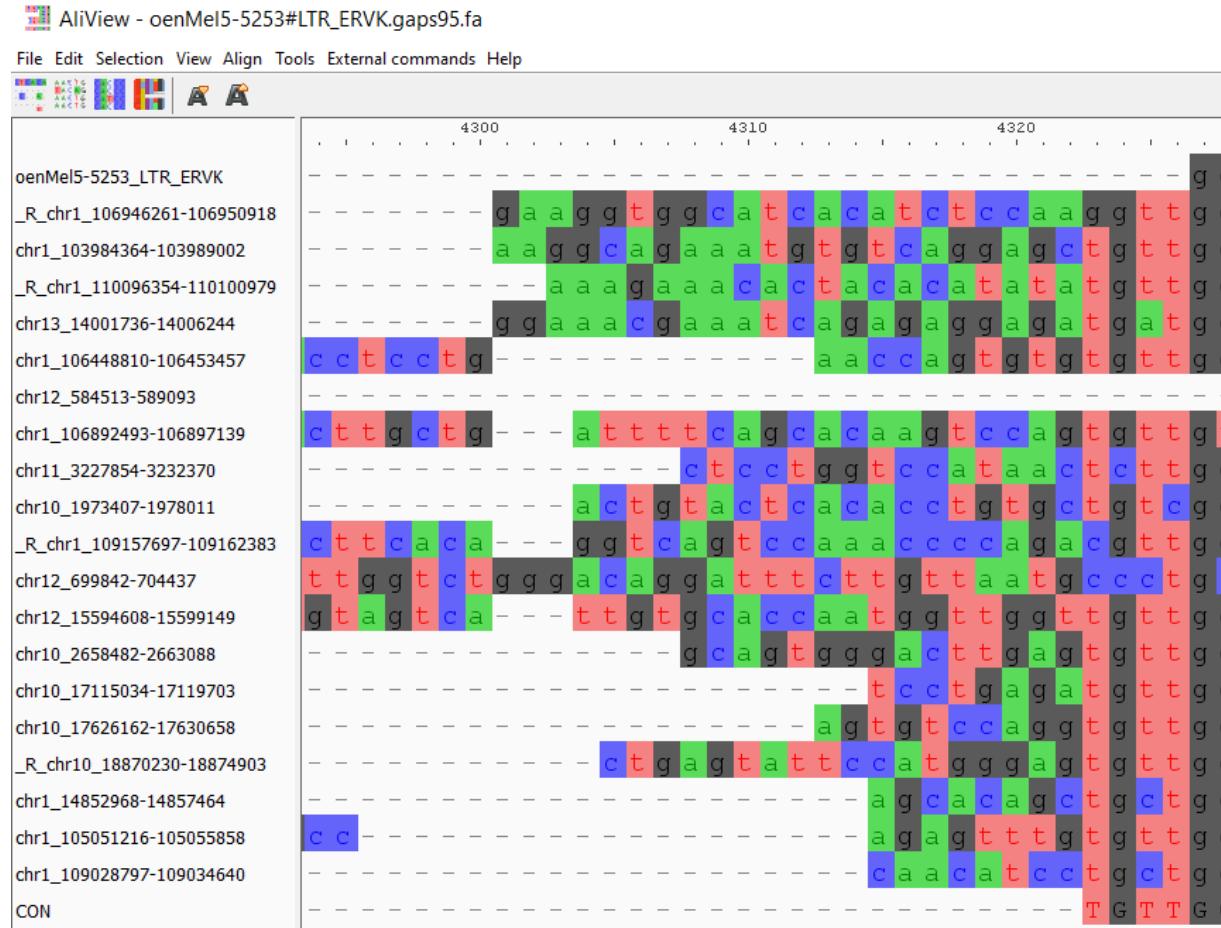
Find the TSDs

5' end

(if any!)

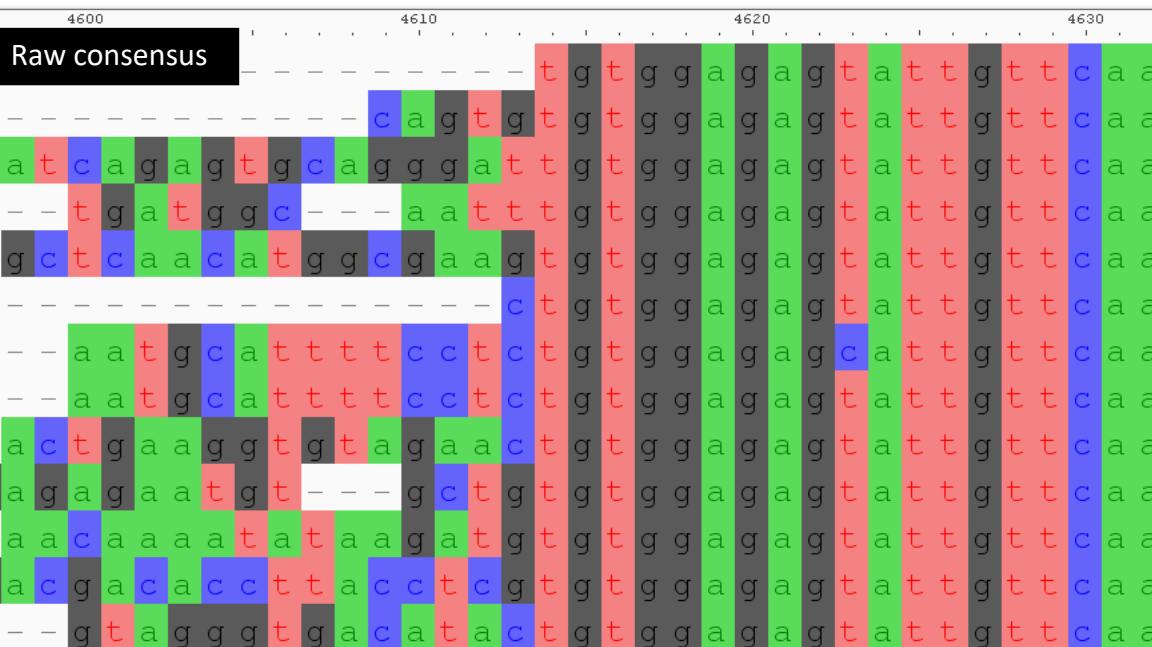
Start looking from the 3' terminus

3' end

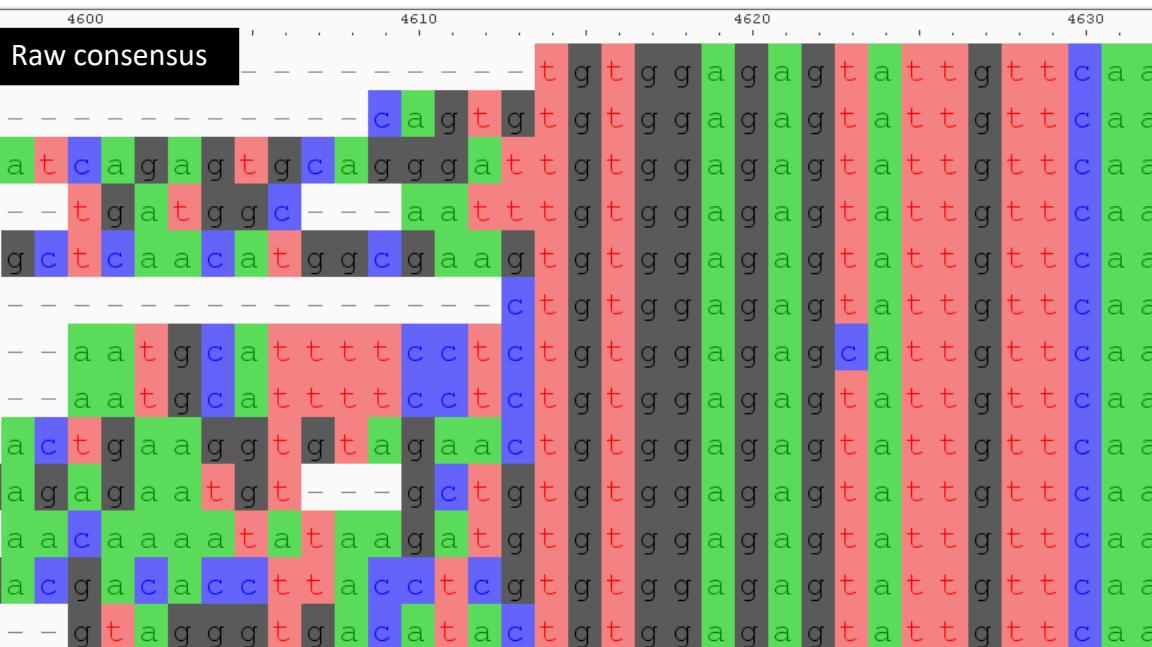


► LTRs and LTR elements usually start with 5'-TG CA-3'

Find the TSDs



Find the TSDs



tccccccctcatcag

tctcgctttcaccgacagcaattga

tctccctttctccgacagcgaagaa

tctccctttctccaaca

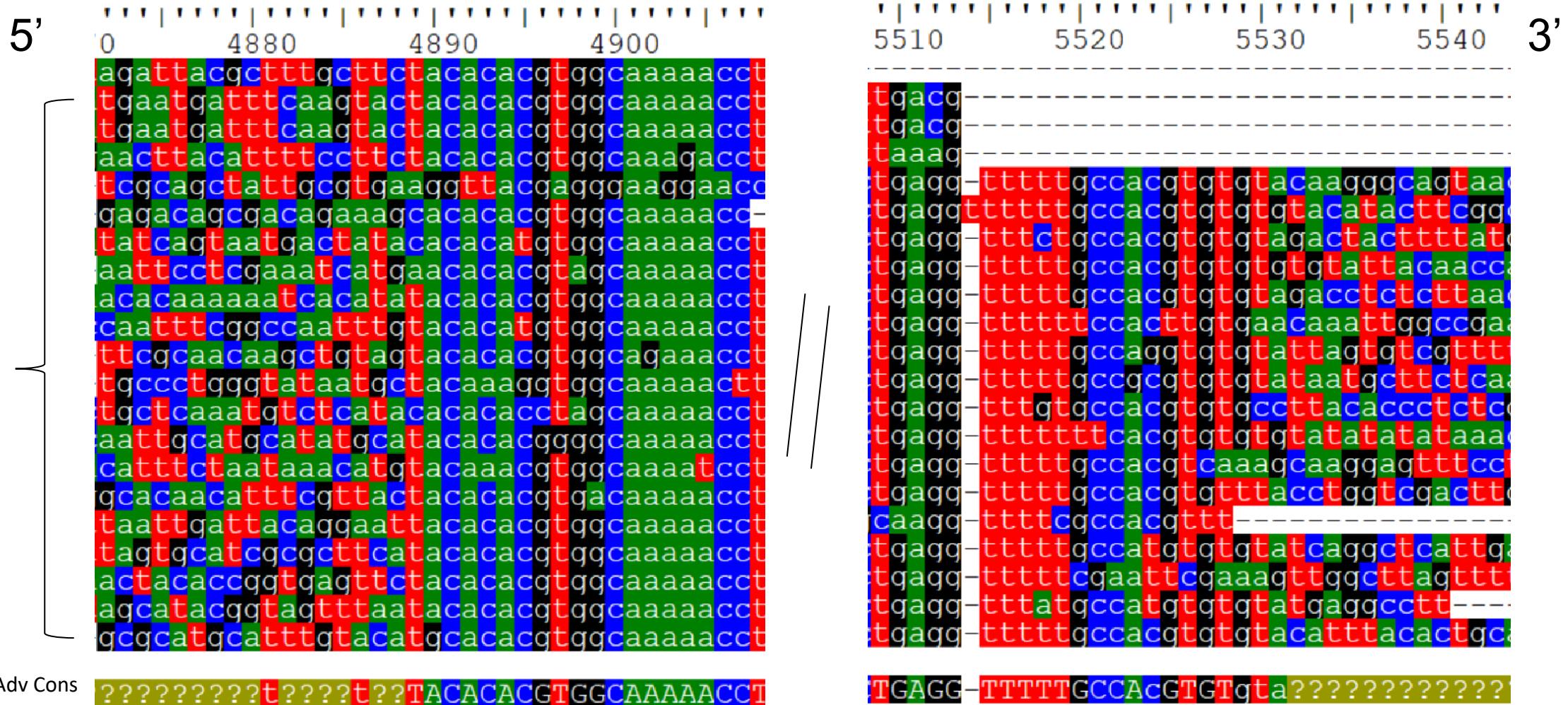
tctcactttcacccgacataagaacaa

tctcactttcacccgacatggaaaaacttagaaca

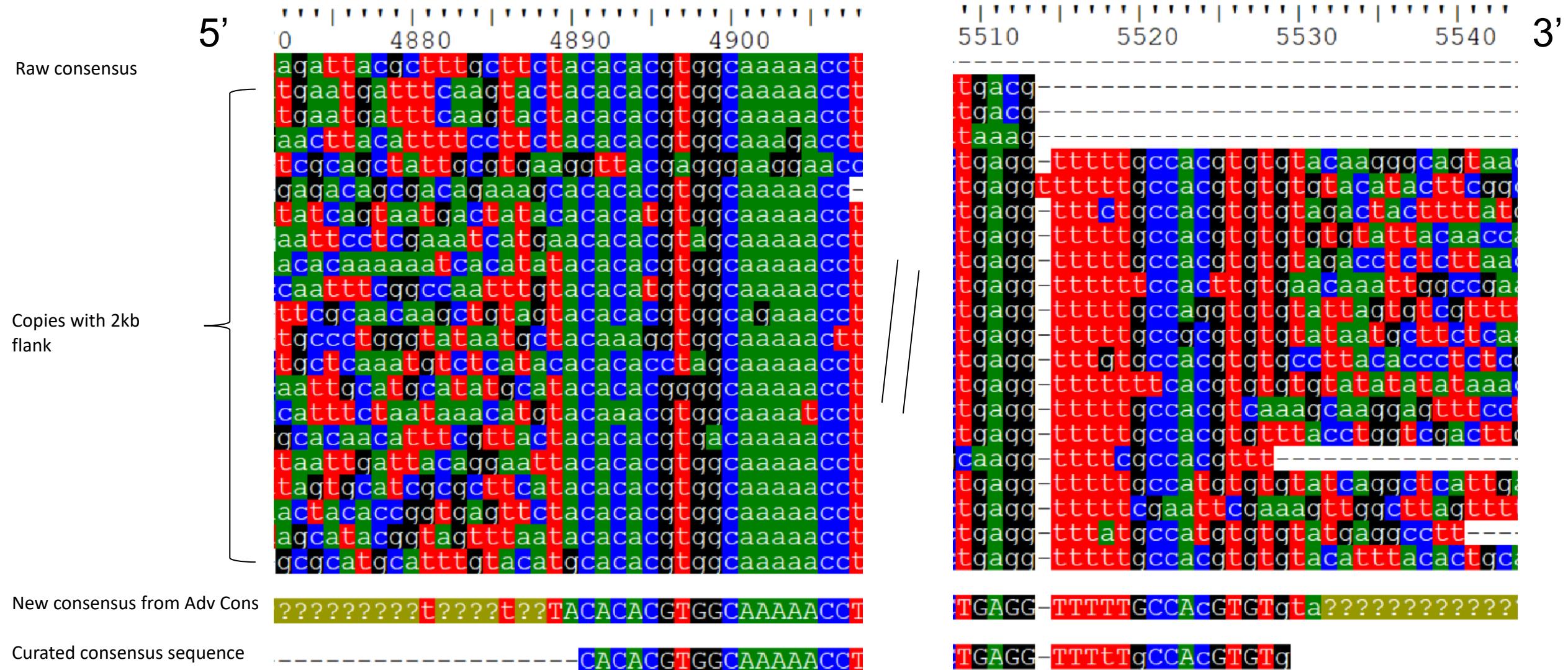
tctcactttcacccgacacaa

tctcactttcacccgacacac

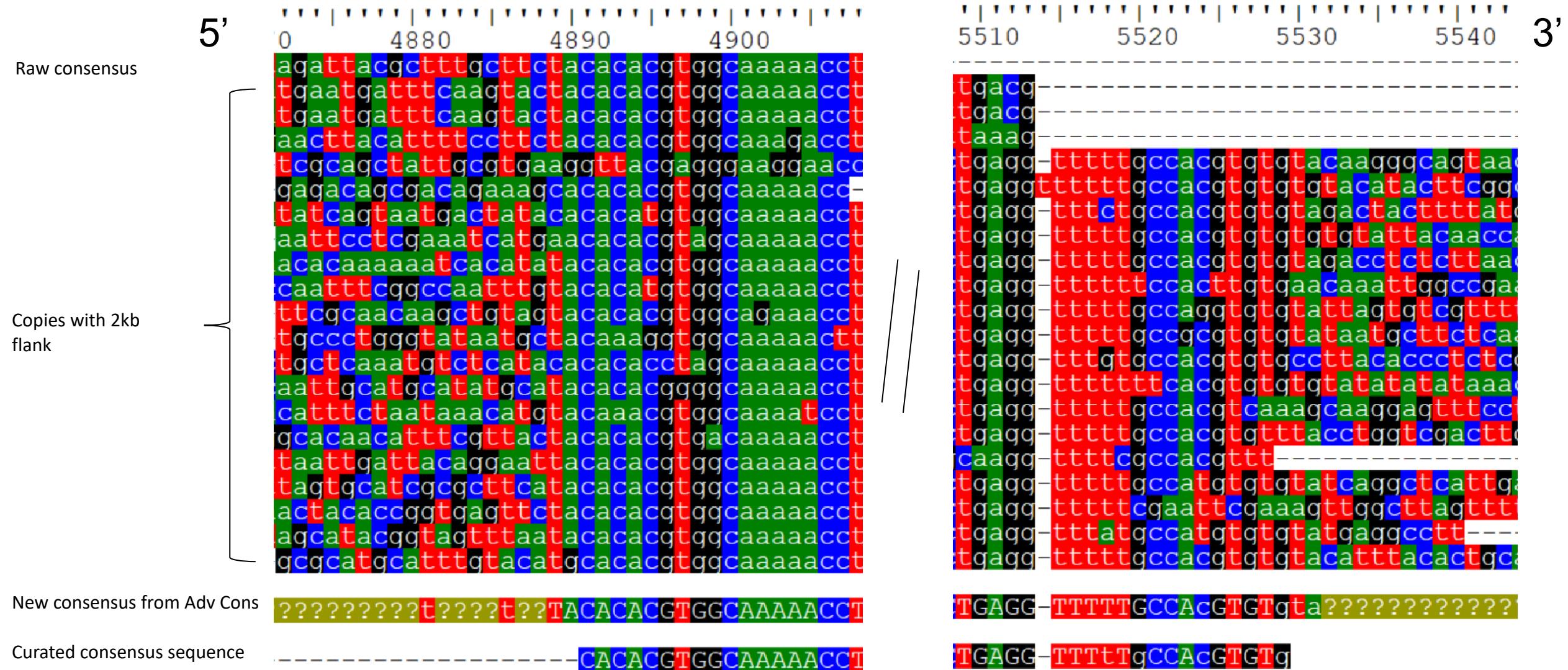
Find the TSDs



Find the TSDs



Find the TSDs



In the presence of TSD motifs, annotate the motif in column K

DNA transposons and TSDs

Table 1 Classification and characteristics of eukaryotic DNA transposons

Superfamily	Related IS	TSD	Length ¹ (kb)	TIRs ¹ (bp)	Terminal motif (5'-3')	TPase ¹ (aa)	Catalytic motif	DNA-binding motif	Additional proteins
<i>Tel/mariner</i>	IS630	TA	1.2–5.0	17–1100	Variable	300–550	DD(30–41)D/E	HTH (cro/paired)	
<i>bAT</i>	nd	8 bp	2.5–5	10–25	YARNG	600–850	D(68)D(324)E ²	ZnF (BED)	
<i>P element</i>	nd	7/8 bp	3–11	13–150	CANRG	800–900	D(83)D(2)E(13)D ³	ZnF (THAP)	
<i>MuDR/ Foldback</i>	IS256	7–10 bp	1.3–7.4	0-sev. Kb	Variable	450–850	DD(~110)E	ZnF (WRKY/GCM1)	
<i>CACTA</i>	nd	2/3 bp	4.5–15	10–54	CMCWR	500–1200	nd	nd	TNPA (DNA-binding protein)
<i>PiggyBac</i>	IS1380	TTAA	2.3–6.3	12–19	CCYT	550–700	DDE?	nd	
<i>PIF/ Harbinger</i>	IS5	TWA	2.3–5.5	15–270	GC-rich	350–550	DD(35–37/47–48)E	HTH	PIF2p (Myb/SANT domain)
<i>Merlin</i>	IS1016	8/9 bp	1.4–3.5	21–462	GGNRM	270–330	DD(36–38)E	nd	
<i>Transib</i>	nd	5 bp	3–4	9–60	CACWATG	650–700	DD(206–214)E	nd	
<i>Banshee</i>	IS481	4/15 bp	3–5	41–950	TGT	300–400 ⁴	DD(34)E	HTH	
<i>Helitron</i>	IS91	none	5.5–17	none	5'-TC...CTAR-3'	1400–3000 ⁵	HHYY ("REP motif")	ZnF-like	RPA (in plants)
<i>Maverick</i>	none	5/6 bp	15–25	150–700	Simple repeat	350–450 ⁴	DD(33–35)E	ZnF (HHCC)	4–10 DNA virus-like proteins

¹Refers to a potentially complete, autonomous element.

²Motif in *Hermes* TPase.

³Motif in *Drosophila P element* TPase.

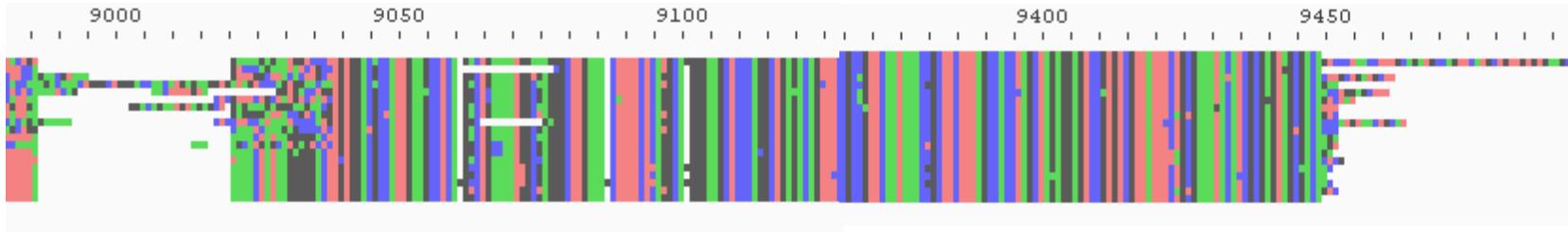
⁴RVE integrase-like.

⁵REP-Helicase.

nd = not determined.

What to look for and fix during curation

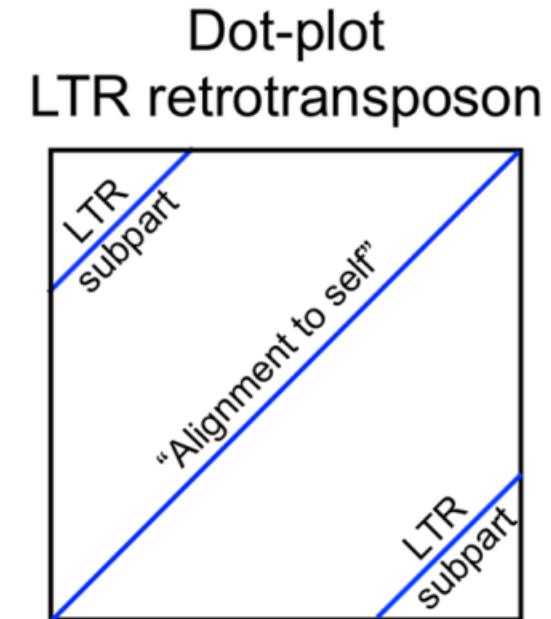
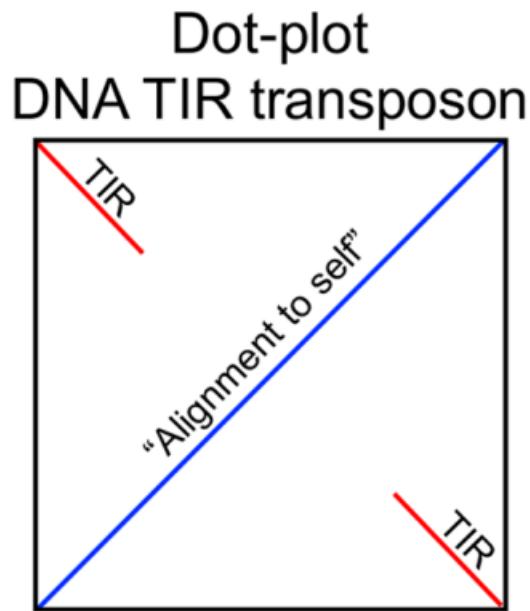
Part 2



- 1 Check for the presence of LTRs or TIRs (e.g., dotplots)
- 2 Check for protein domains
- 3 Check for nucleotide homology
- 4 Re-classify the element

Check for LTRs and TIRs

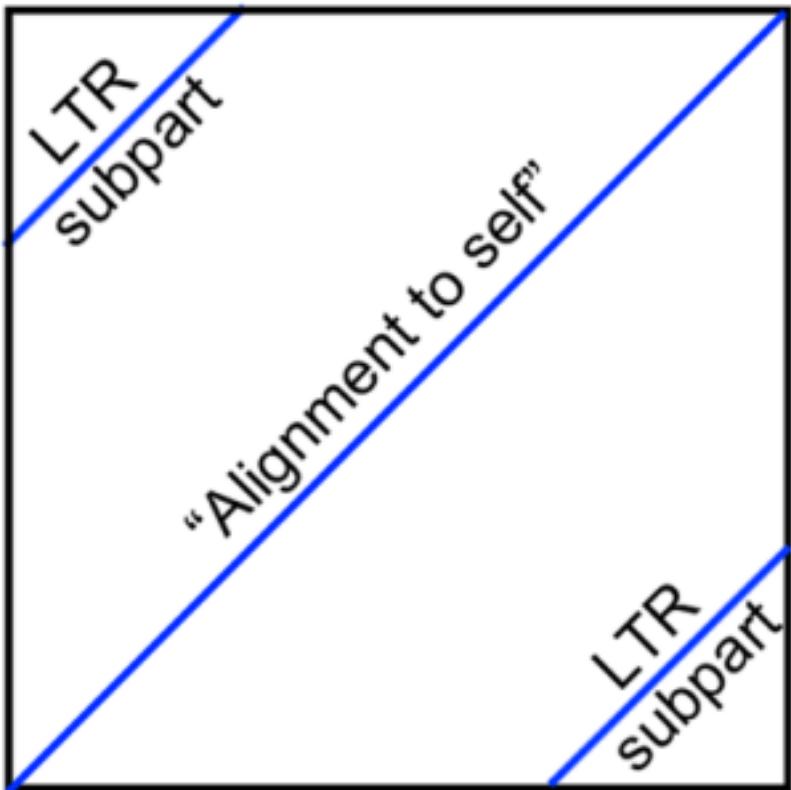
How to read the dotplot of a self-alignment



Use MAFFT for dotplots!
<https://mafft.cbrc.jp/alignment/server/>

Check for LTRs

LTRs = Long Terminal Repeats (direct repeats)



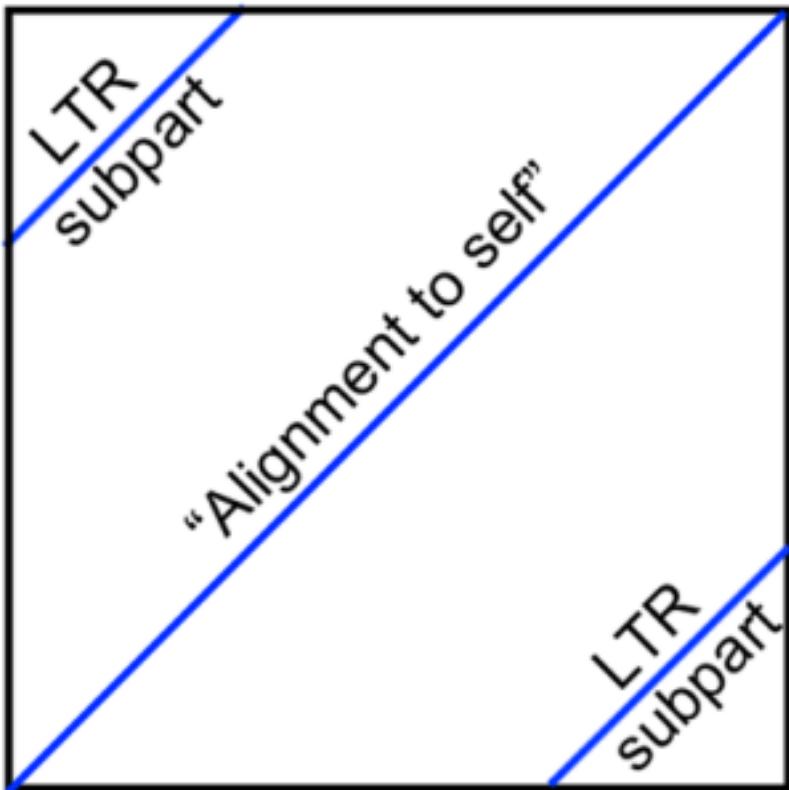
Separate the LTRs from the internal sequence (one of the two LTRs – they should be nearly identical)

You will have a row with the LTR and a second row with the internal region in the Excel table

The name of the consensus sequence remain the same but you add the suffiy '_LTR' and '_int' (or '_I') just before the classification tag

Check for LTRs

LTRs = Long Terminal Repeats (direct repeats)



Separate the LTRs from the internal sequence (one of the two LTRs – they should be nearly identical)

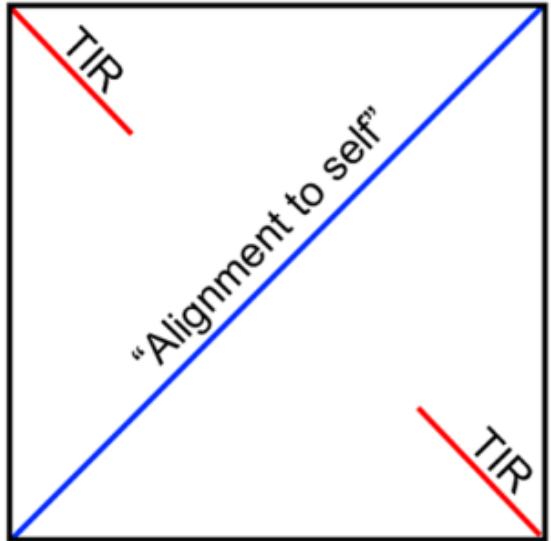
You will have a row with the LTR and a second row with the internal region in the Excel table

The name of the consensus sequence remain the same but you add the suffix '_LTR' and '_int' (or '_I') just before the classification tag

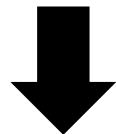
Solo-LTR cannot give this signal! Remember that LTRs usually start with TG and ends with CA

Check for TIRs

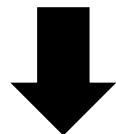
TIRs = Terminal Inverted Repeats



Sometimes TIRs are too small to be seen in a dotplot like this



Copy consensus twice in a new alignment



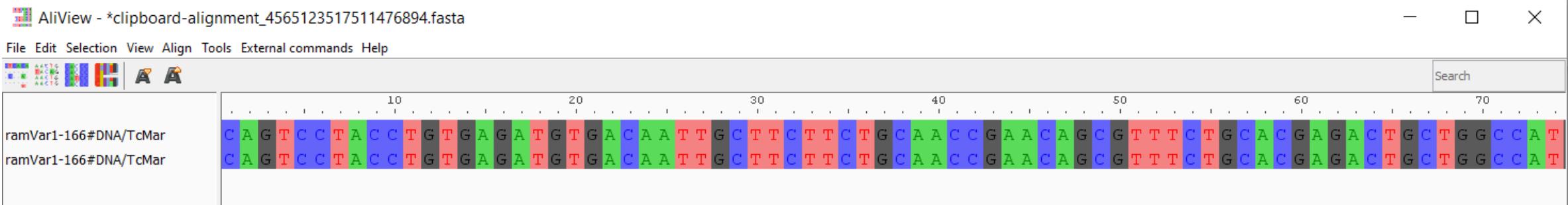
Reverse complement one of them and look for homology

Annotate the length of the TIR in column L

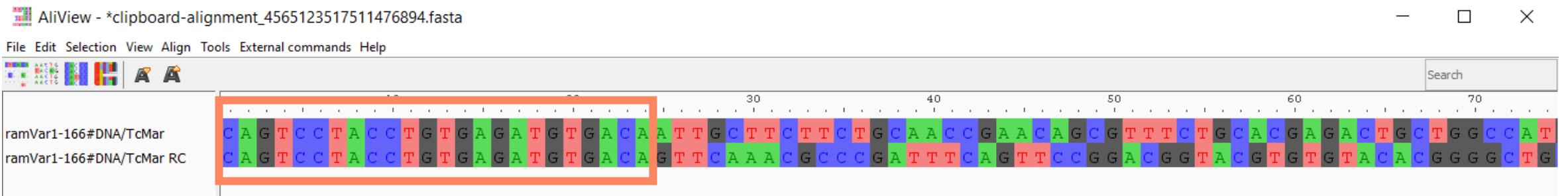
Check for TIRs

TIRs = Terminal Inverted Repeats

Copy consensus twice in a new alignment



Reverse complement one of them and look for homology

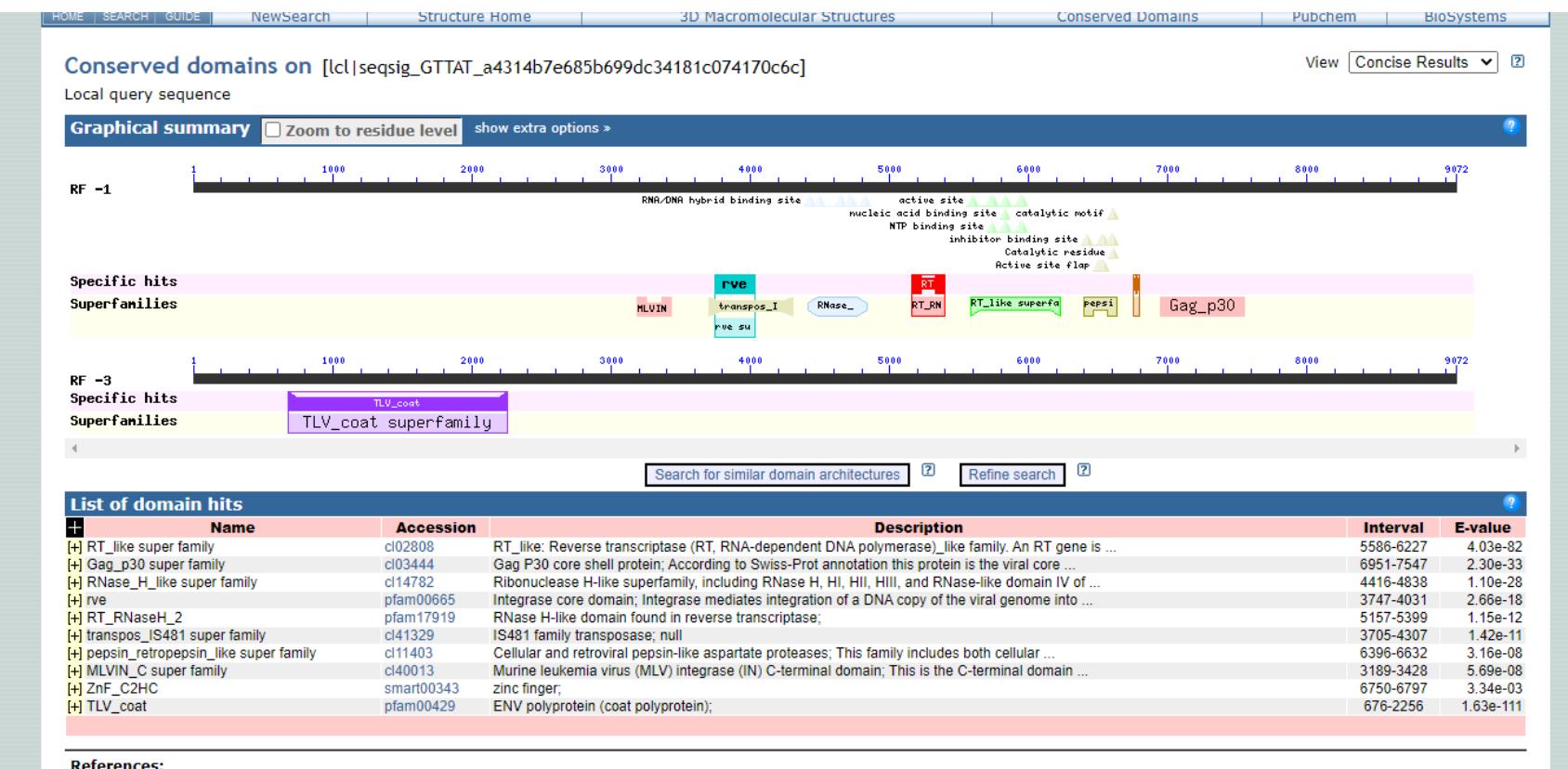


Annotate the length of the TIR in column L

Check for protein domains

Conserved Domains Database

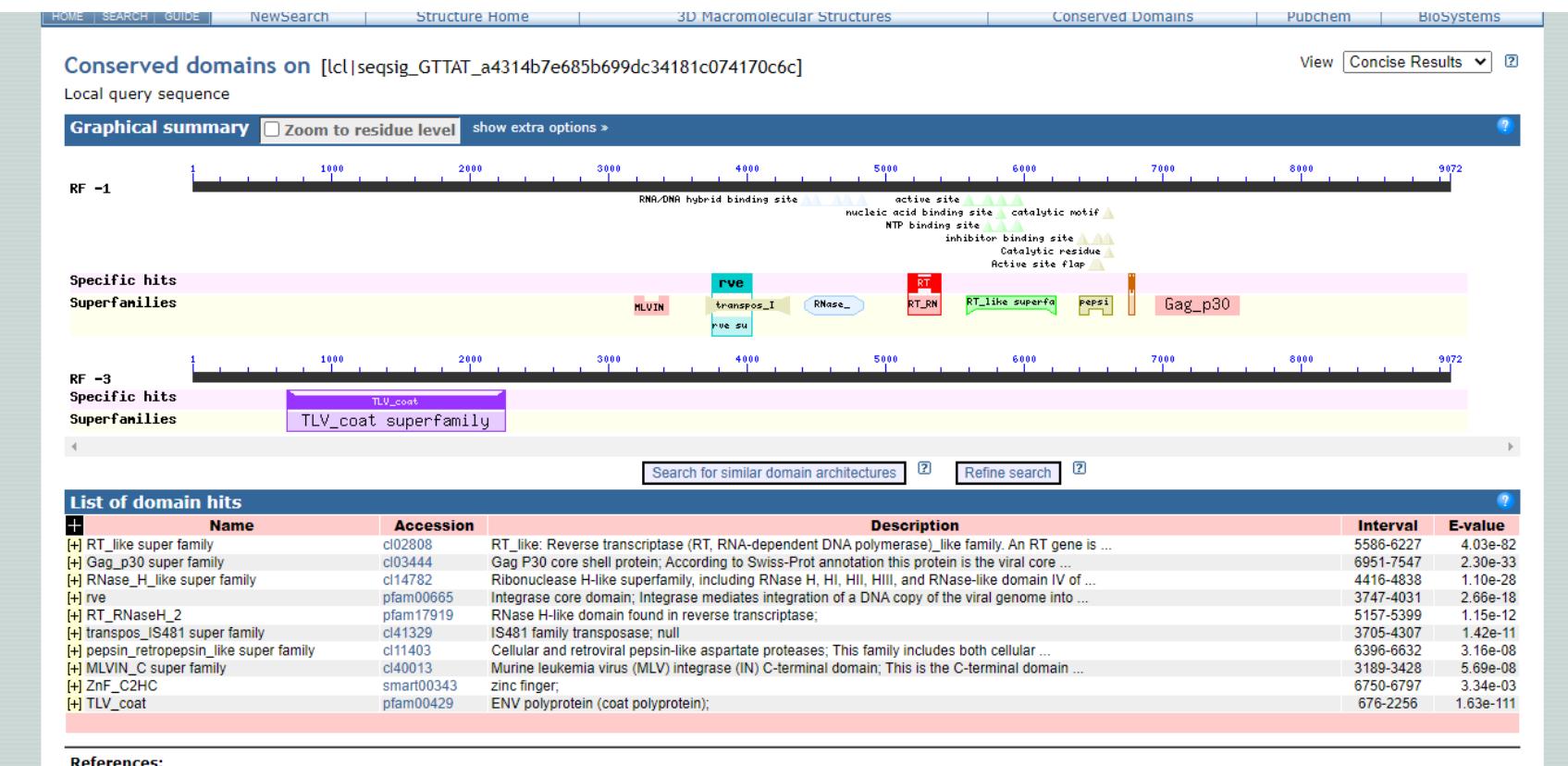
<https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>



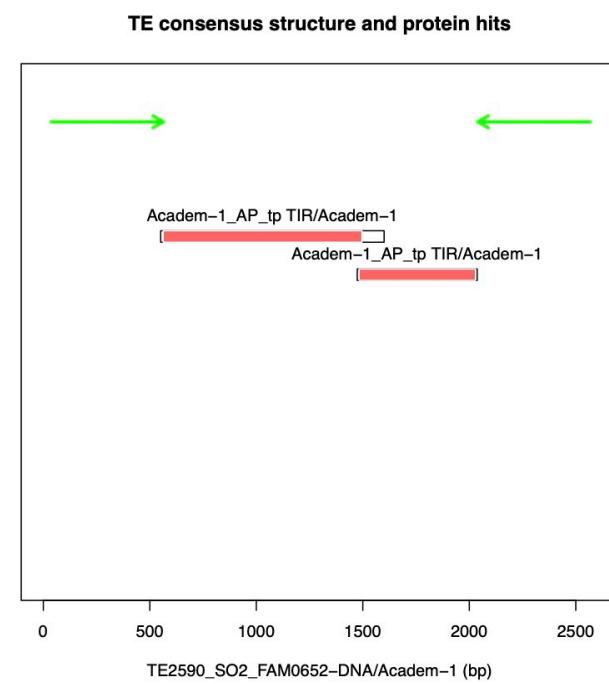
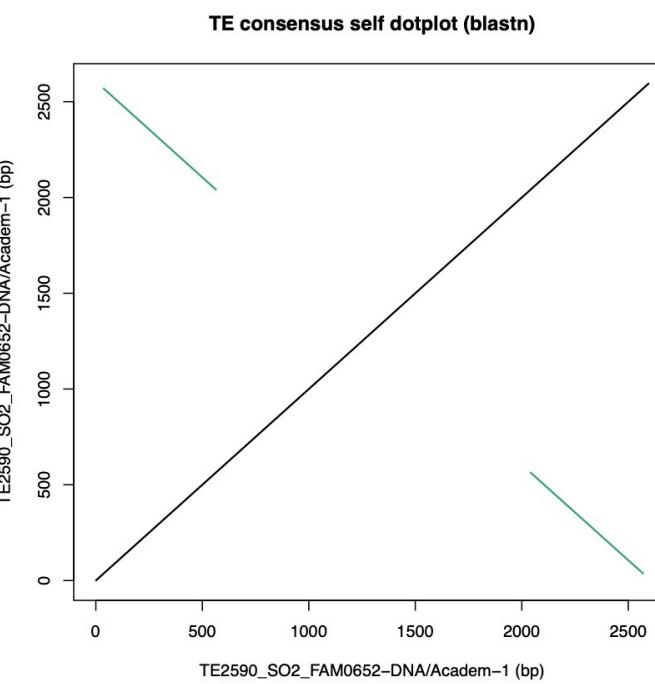
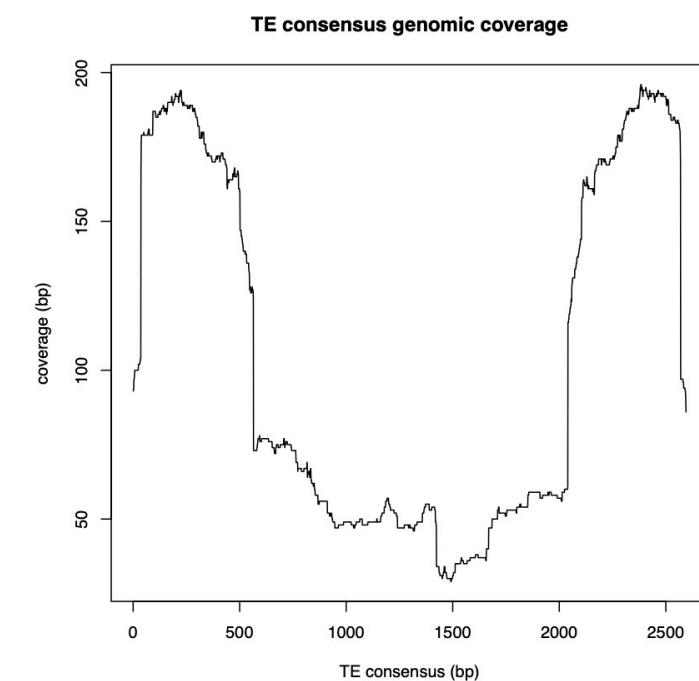
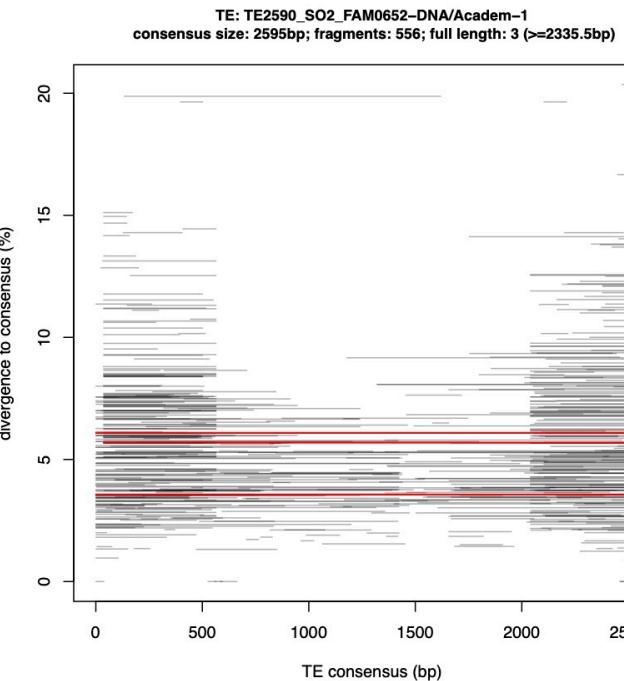
Check for protein domains

Conserved Domains Database

<https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi>



The presence and order of protein domains can help for classification – annotate the presence of important domains in column G



Check for nucleotide homology

Example: Censor and Repbase

<https://www.girinst.org/censor/index.php>

Map of Hits

[SVG viewer](#) is required to view graphical representation of the map as Scalable Vector Graphics (SVG plot).

chr20:12962364-12971435 ([SVG Plot](#); [Alignments](#); [Masked](#))



Name	From	To	Name	From	To	Class	Dir	Sim	Pos/Mm:Ts	Score
chr20:12962364-12971435	1	276	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7451	1.7097	797
chr20:12962364-12971435	288	394	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261
chr20:12962364-12971435	532	8544	IguERV6_I	1	8078	ERV/ERV1	c	0.7176	1.7271	26725
chr20:12962364-12971435	8546	8821	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7490	1.7333	814
chr20:12962364-12971435	8833	8939	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261

Check for nucleotide homology

Example: Censor and Repbase

<https://www.girinst.org/censor/index.php>

Map of Hits

[SVG viewer](#) is required to view graphical representation of the map as Scalable Vector Graphics (SVG plot).

chr20:12962364-12971435 ([SVG Plot](#); [Alignments](#); [Masked](#))



Name	From	To	Name	From	To	Class	Dir	Sim	Pos/Mm:Ts	Score
chr20:12962364-12971435	1	276	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7451	1.7097	797
chr20:12962364-12971435	288	394	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261
chr20:12962364-12971435	532	8544	IguERV6_I	1	8078	ERV/ERV1	c	0.7176	1.7271	26725
chr20:12962364-12971435	8546	8821	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7490	1.7333	814
chr20:12962364-12971435	8833	8939	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261

Check for nucleotide homology

Example: Censor and Repbase

<https://www.girinst.org/censor/index.php>

Map of Hits

[SVG viewer](#) is required to view graphical representation of the map as Scalable Vector Graphics (SVG plot).

chr20:12962364-12971435 ([SVG Plot](#); [Alignments](#); [Masked](#))



Name	From	To	Name	From	To	Class	Dir	Sim	Pos/Mm:Ts	Score
chr20:12962364-12971435	1	276	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7451	1.7097	797
chr20:12962364-12971435	288	394	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261
chr20:12962364-12971435	532	8544	IguERV6_I	1	8078	ERV/ERV1	c	0.7176	1.7271	26725
chr20:12962364-12971435	8546	8821	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7490	1.7333	814
chr20:12962364-12971435	8833	8939	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261

In general, take into consideration hits that are longer than 50 bp and have a score higher than 500

Check for nucleotide homology

Example: Censor and Repbase

<https://www.girinst.org/censor/index.php>

Map of Hits

[SVG viewer](#) is required to view graphical representation of the map as Scalable Vector Graphics (SVG plot).

chr20:12962364-12971435 ([SVG Plot](#); [Alignments](#); [Masked](#))



Name	From	To	Name	From	To	Class	Dir	Sim	Pos/Mm:Ts	Score
chr20:12962364-12971435	1	276	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7451	1.7097	797
chr20:12962364-12971435	288	394	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261
chr20:12962364-12971435	532	8544	IguERV6_I	1	8078	ERV/ERV1	c	0.7176	1.7271	26725
chr20:12962364-12971435	8546	8821	IguERV6_LTR1	407	665	ERV/ERV1	c	0.7490	1.7333	814
chr20:12962364-12971435	8833	8939	hAT6-2_SP	4951	5043	DNA/hAT	d	0.7216	1.5714	261

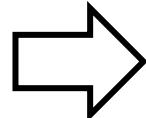
In general, take into consideration hits that are longer than 50 bp and have a score higher than 500

Annotate similarity to known repeats in column F

Re-classification

Collect the evidence

- TSD length and motifs
- TIR length
- Protein domains
- Homology to known transposable elements

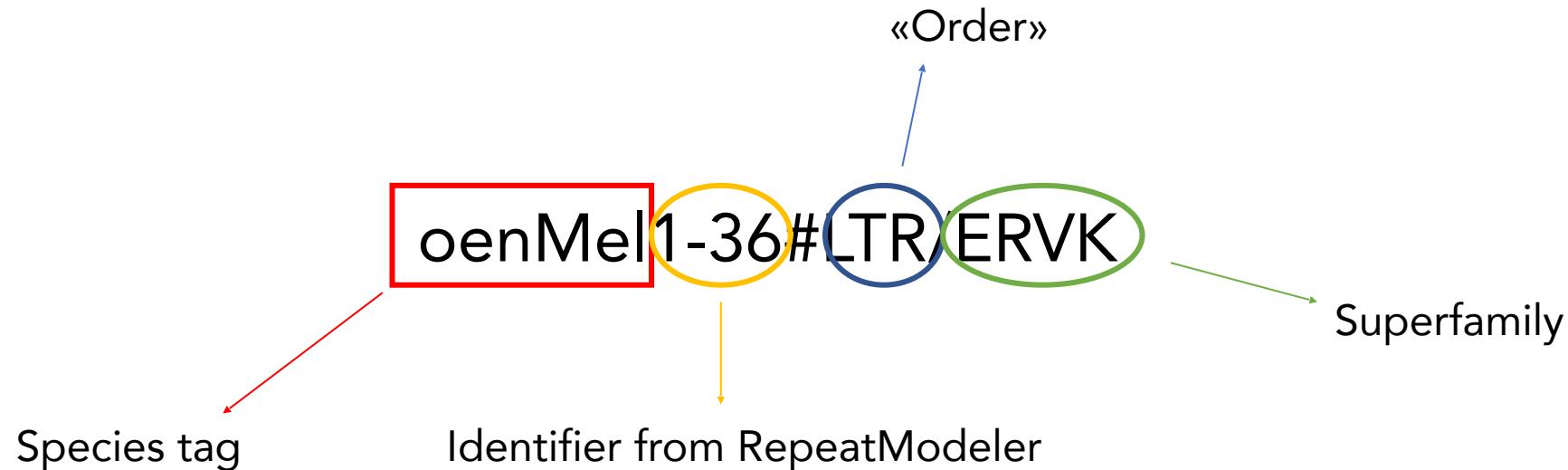


- Check classification tables:
- Wicker et al 2007
 - Feschotte and Pritham 2007
 - Feschotte et al 2020

Change the classification tag according to the collected evidence in column S

It is possible that no better classification can be given to the consensus sequence or that a previously classified repeat becomes an 'Unknown'

Structure of the consensus sequence name



Other examples

oenMel1-36_LTR#LTR/ERVK

oenMel1-36_int#LTR/ERVK

oenMel1-90.inc#LINE/CR1

ramVar1-626#DNA/Piggybac

Practical 2

Part 1 – Explore the diversity of TEs with TEAid plots

Part 2 – Curation and classification of alignments

- Inspect the alignments in the folder “Practical2_part2 > Alignments”
- Create a new consensus sequence with Advanced Consensus Maker (look at the slides for parameters)
- Fix SNPs and indels (majority rule)
- Find termini and TSDs

Classify the element based on TSDs and termini and use dotplots, Repbase, protein hits to help classification!

In the Practical2 folder you can find papers with classification tables.

Curation tips

- Pay attention to **CpG sites**, especially if you know that methylation is "used" by the genome to silence TEs.
- Apply a majority rule to decide how **indels/gaps** in the alignments should be represented in the new consensus sequence. Example: in a position we have 5 gaps ("-", 5 Cs and 5Gs. Break the tie first between insertion and gaps then apply the majority rule on the winner. In this case insertions win over gaps and you write "S" in the consensus (C/G = Strong).
- **Do not leave "?" in the new consensus sequences** because it is not recognised by blast and RepeatMasker as a valid character. If you have a tie between nucleotide frequencies, please use the ambiguity IUPAC characters (e.g., S, W, R, Y, M, K, ...). Also you can look at the nucleotide at that position in the old consensus to decide. In the worst case scenario "N" is a perfectly acceptable character.
- Double check that the new consensus sequences do not contain flanking regions and TSDs (be particularly careful with Mariners).
- **Double check that the features of the curated TEs are recorded** in the Excel file.
- **Adopt a consistent nomenclature within your library and try to be as consistent as possible with repeat libraries from closely related species** (if any) or, in general, to other manually curated repeat libraries.
- Adopt a **nomenclature for non-autonomous elements that is consistent** with other repeat libraries. Always record it in the Excel file.

Curation tips

- When you classify an element **check the naming convention for the class/superfamily on Repbase or Dfam** if possible (it could be inconsistent with RepeatMasker though)
- In case you find **LTR retrotransposons**, please split the consensus in the LTR portion ("_LTR") and internal portion ("_I"); you can check on Repbase examples of LTR naming.
- In case you find **LINE retrotransposons**, don't worry if you cannot find both ends of the element. Because of the 5' truncation phenomenon, it is very difficult to obtain complete consensus sequences of LINEs. Try to extend the consensus as much as you can (do not build a consensus of less than 3 sequences). Try to have at least one end included in the consensus sequence
- In DNA transposons, TIRs must be symmetrical but one or two mismatches are ok
- Don't be afraid to use the "Comment" column to annotate anything you think interesting, you may find unexpected patterns as you keep curating the library.