

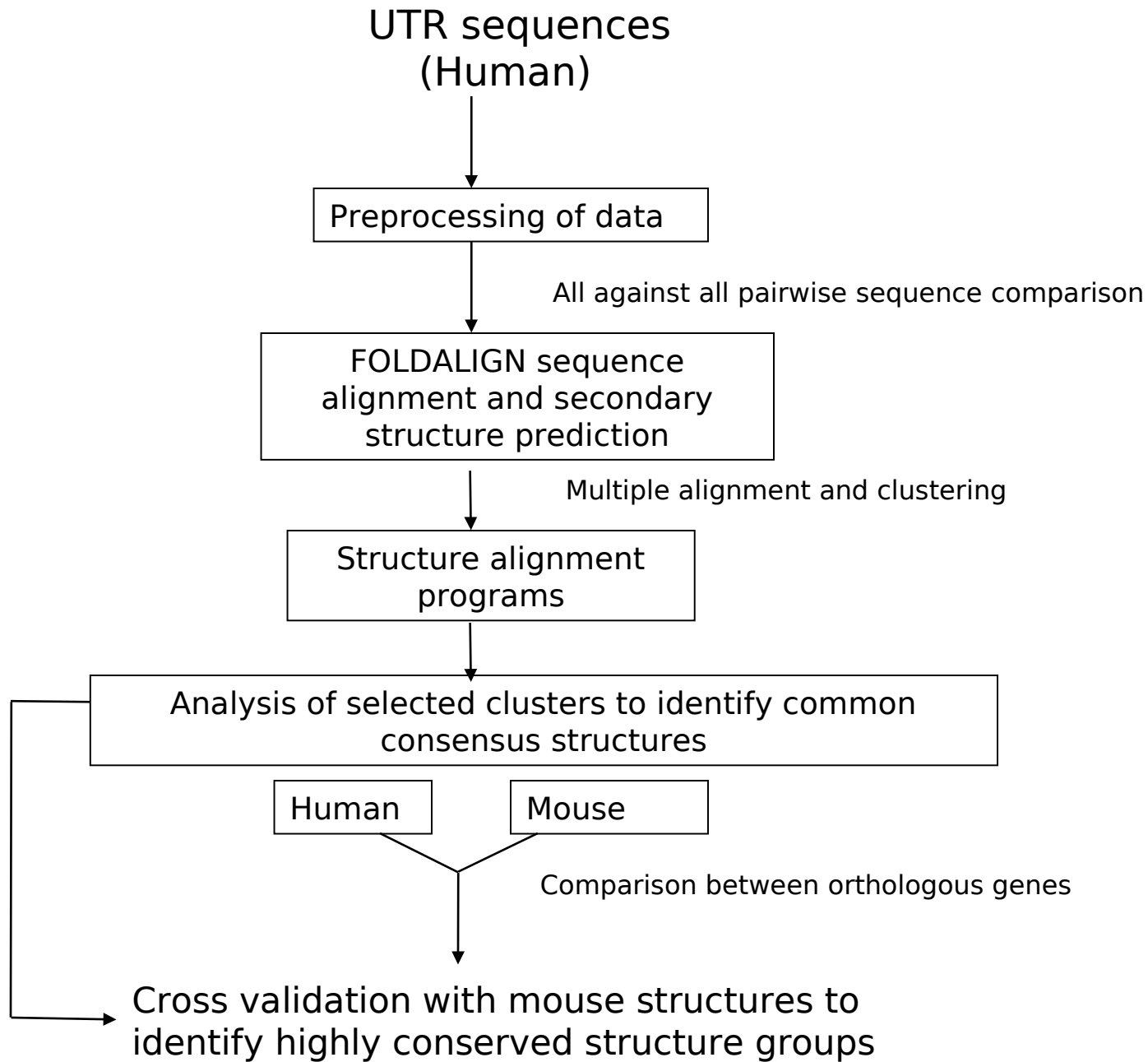
Sequence analysis of Human UTRs

Winterseminar 2012, Bled

Simranjeet Kaur
Centre for non-coding RNA in technology and
Health

Outline

- ▶ An approach to uncover new conserved RNA structures in human UTRs
- ▶ Systematic study of combined sequence and structure similarity among human UTRs.



UTR sequences for human (assembly GRCh37) retrieved from Ensembl release 64 using BioMart

	5' UTRs	3' UTRs	Total
UTRs	78,662	71,427	
Genes	20,399	20,390	20,764
Average transcript per gene	3.85	3.5	
Identical UTRs	11,429	12,343	
UTR candidates	67,233	59,084	126,317

Genes having both 5' and 3' UTRs - 20,025

More insights into the human UTR data-set

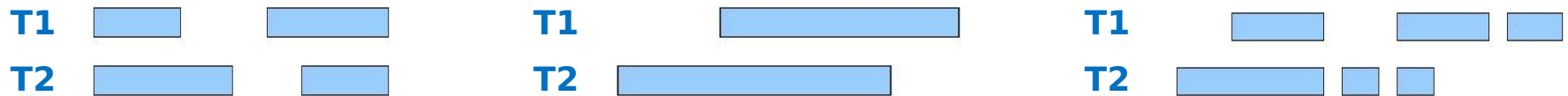
	5' UTRs	3' UTRs	Total
UTRs after removing duplicates	66,665	58,869	125,534
Genes	19,451	19,573	
Average transcript per gene	3.43	3.01	
Genes with single transcripts	6,002	6,283	
Genes having both 5' and 3' UTRs			18,969
UTRs with single exons	35,932	30,733	
Max. length of UTRs	17,184 (average- 256)	17,159 (average- 950)	
Max length for UTRs with single exons	17,184 (average- 179)	14,960(average- 930)	

Removing cis redundancy from the dataset:-

-Merged transcripts that were completely within other transcripts of a gene.



-Merged overlapping coordinates.



-New dataset (non-redundant UTRs for each gene) after merging overlapping and common exons:-

5'UTRs- 38,494

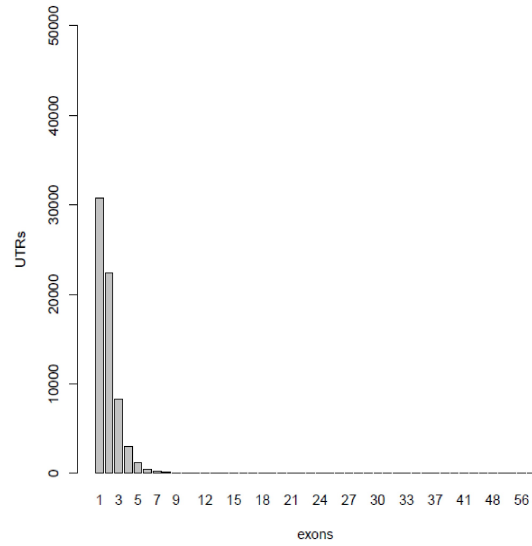
3'UTRs- 31,557

T1- Transcript 1

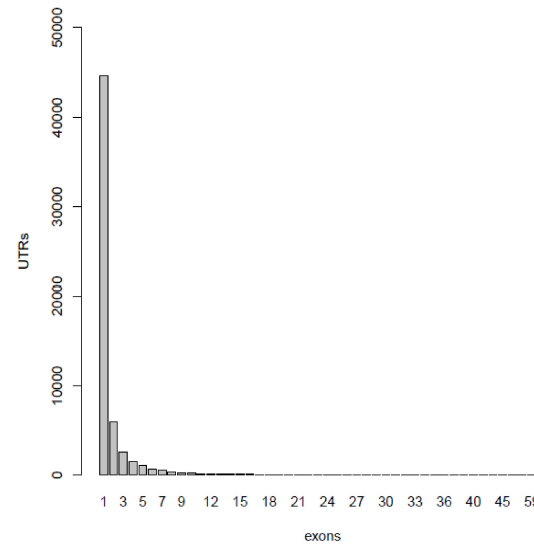
T2- Transcript 2

Distribution plots after merging
overlapping coordinates (cis redundancy
reduction)

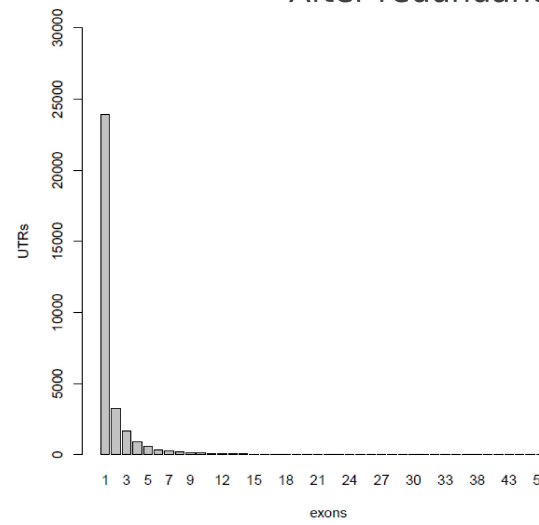
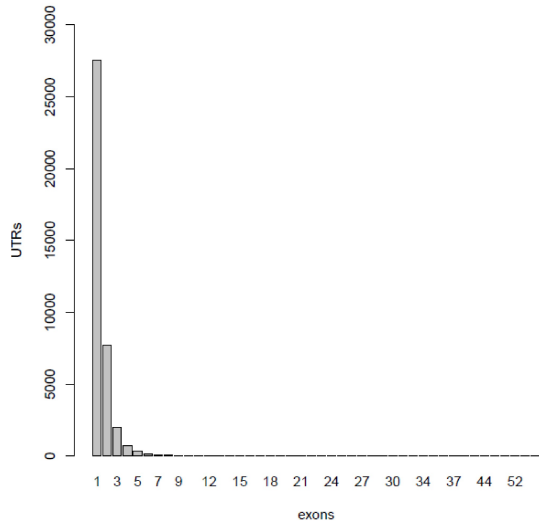
Exons distribution



Before redundancy reduction



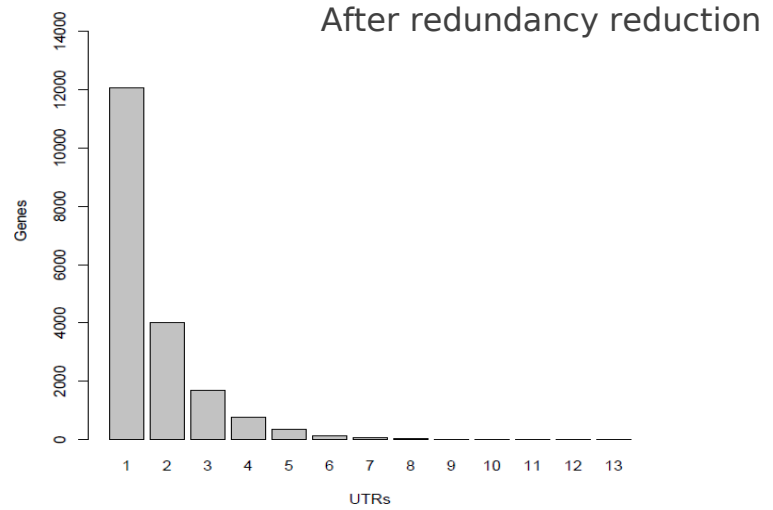
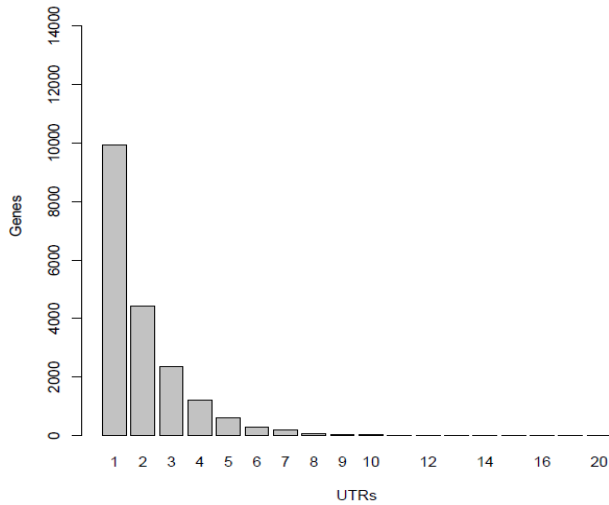
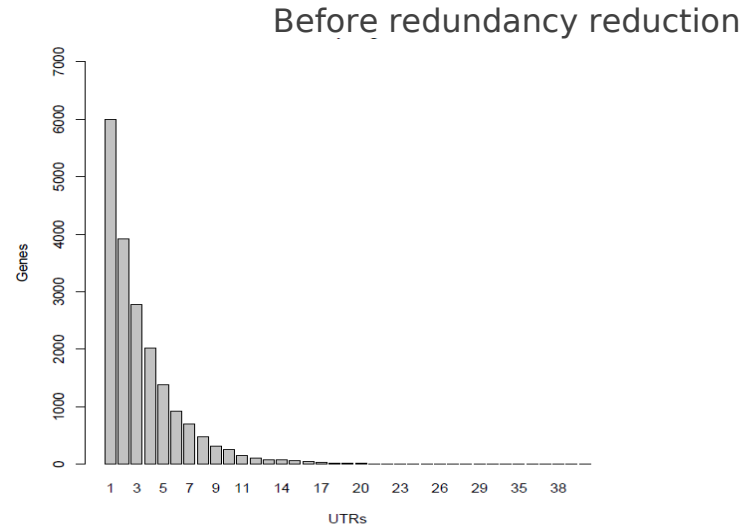
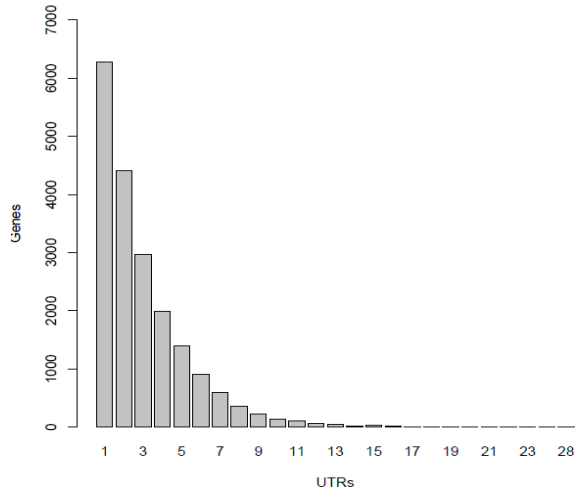
After redundancy reduction



5'UTRs

3'UTRs

UTRs distribution



Identifying repeat elements in UTRs

Repeat elements in UTRs identified and masked (excluding small RNAs) using RepeatMasker Version: 4/26/2011.

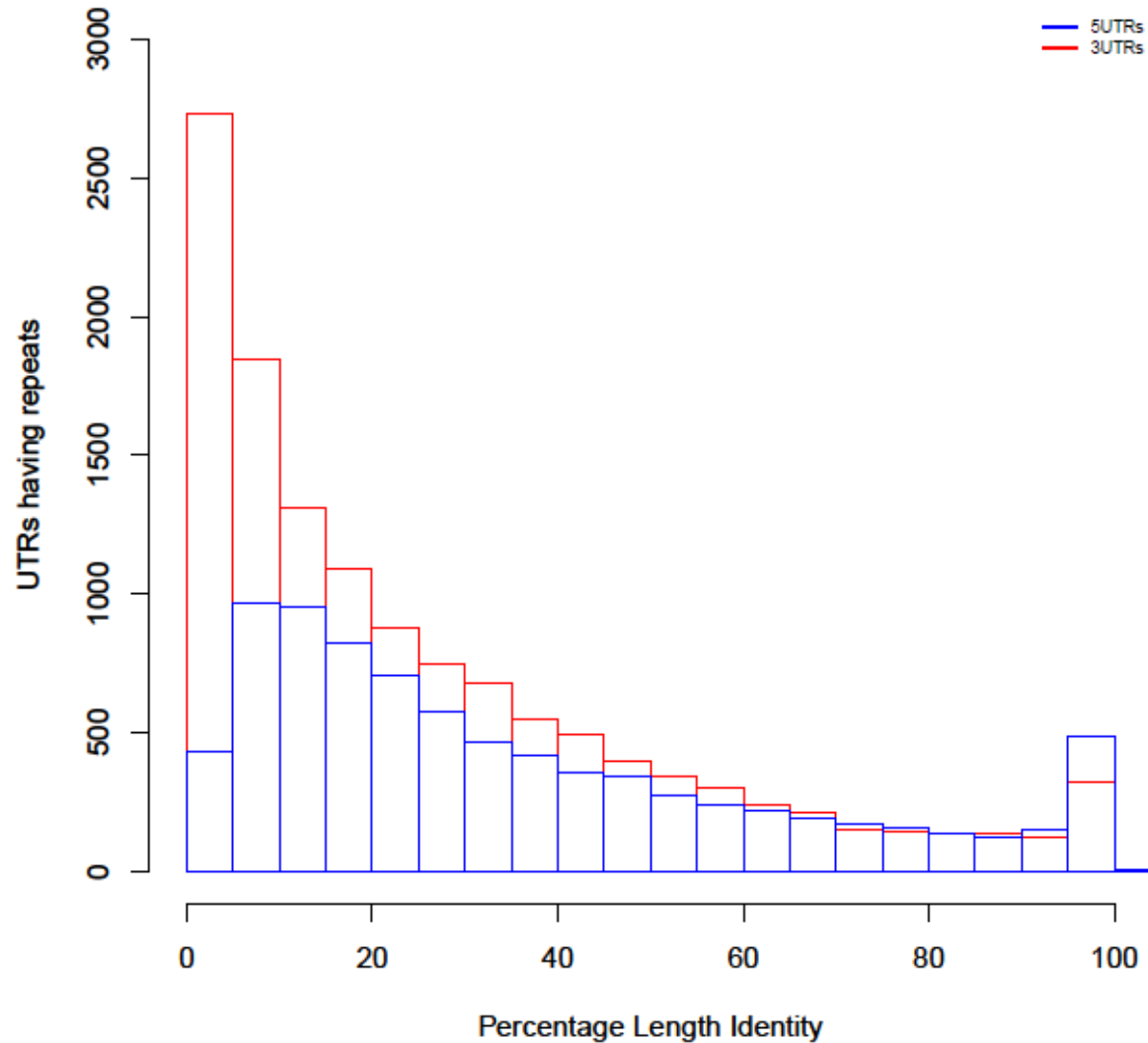
	5'UTRs	3'UTRs
Total UTR Sequences	38,494	31,557
Number of UTRs having repeats	8,191	12,824
Total Number of Repeats in UTRs	11,367	29,401
Average length of alignment with repeats in UTRs	25.00%	11.28%
Bases masked	1,106,390 bp (10.16%)	4,557,296 bp (13.04%)
UTRs having more than 80% length covered by repeats	903	721

Types of Repeats found in human UTRs

	5'UTRs			3'UTRs		
	Number of elements	Length occupied (bps)	Percentage of sequence	Number of elements	Length occupied (bps)	Percentage of sequence
SINEs	2,813	403,185	3.70%	10,456	2,126,539	6.08%
ALUs	1,694	276,006	2.53%	7,047	1,685,870	4.82%
MIRs	1,109	125,883	1.16%	3,307	429,529	1.23%
LINEs	1,258	202,448	1.86%	4,345	1,034,935	2.96%
LINE1	496	98,278	0.90%	2,161	686,766	1.96%
LINE2	663	92,827	0.85%	1,764	285,734	0.82%
L3/CR1	80	9,306	0.09%	320	48,813	0.14%
LTR elements	129,709	129,709	1.19%	1,799	557,845	1.60%
ERVL	25,925	25,925	0.24%	390	120,350	0.34%
ERVL-MaLRs	38,230	38,230	0.35%	888	260,735	0.75%
ERV_classI	59,013	59,013	0.54%	445	148,838	0.43%
ERV_classII	5,009	5,009	0.05%	43	21,954	0.06%
DNA elements	419	59,802	0.55%	2,511	436,966	1.25%
hAT-Charlie	257	33,403	0.31%	1,445	229,544	0.66%
TcMar-Tigger	95	17,669	0.16%	544	126,319	0.36%
Unclassified	17	5,076	0.05%	53	17,862	0.05%
Total Interspersed repeats		800,220	7.35%		4,174,147	11.94%
Small RNAs*	18	1,547	0.01%	108	10,799	0.03%
Satellites	28	5,472	0.05%	19	5,517	0.02%
Simple repeats	2,144	109,978	1.01%	3,485	161,512	0.46%
Low complexity	3,749	189,173	1.74%	5,449	207,974	0.59%

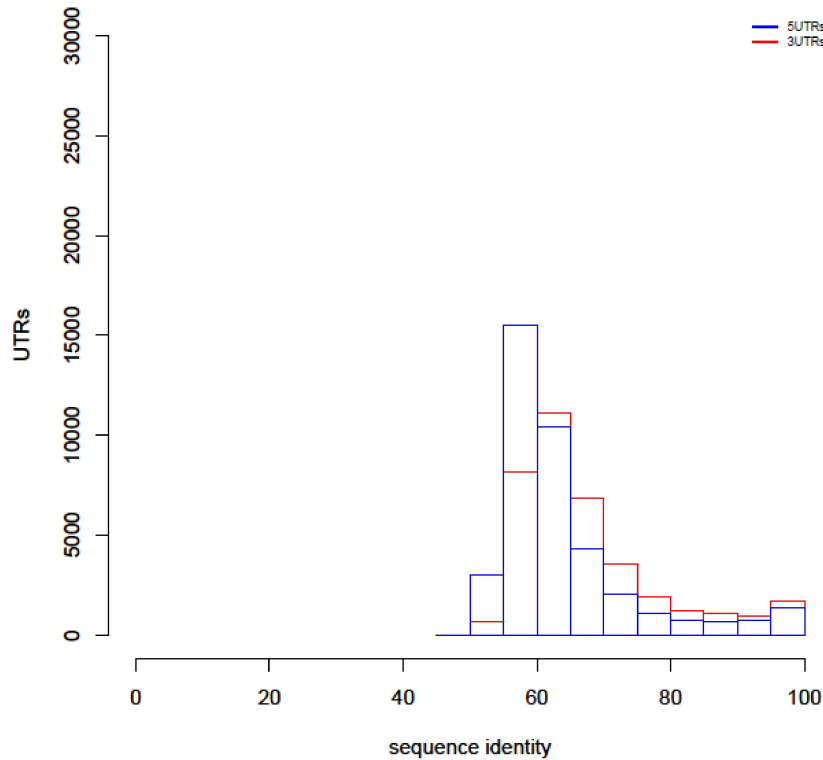
*Small RNAs includes snRNA, srpRNA, scRNA, rRNA, and tRNAs

Percentage length covered by repeats in UTRs

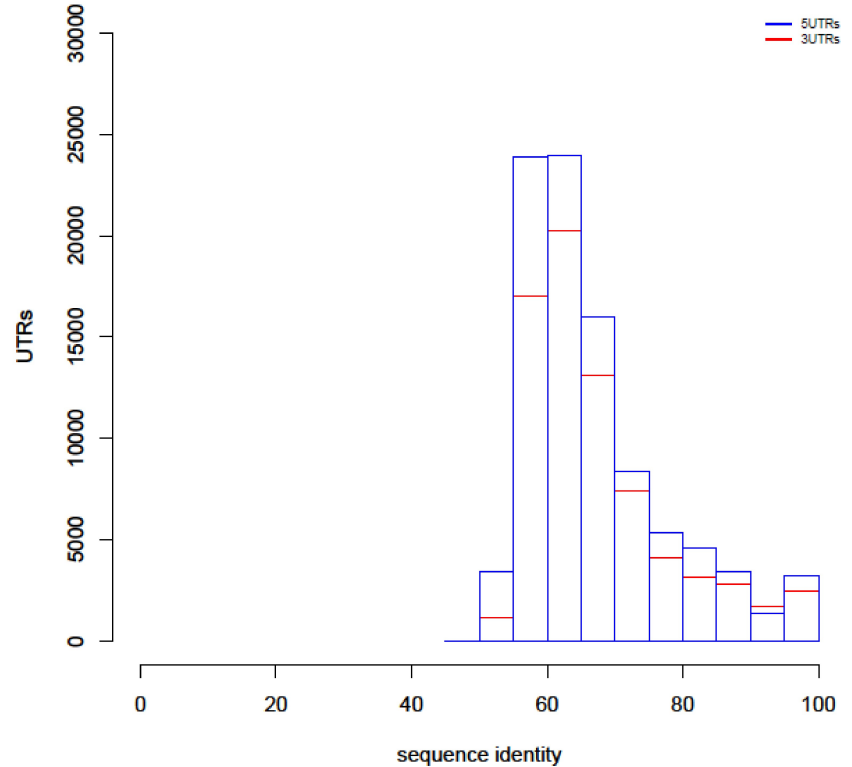


Computing pairwise alignments of UTRs (using ssearch33)

Repeat masked UTRs



Unmasked UTRs



Mouse UTRs (Ensembl release 65, assembly NCBIM37)

	5'UTRs	3'UTRs
Total UTRs	45,601	43,441
Genes	19,921	20,147
Average Length	231	1080
Maximum length	13,770	23,858
Average transcript per gene	2.28	2.15
UTRs after removing redundancy	29,536	25,204

Next steps:-

Removing redundancy from the UTR sequences

Mapping GO terms with UTRs and clustering of genes with overlapping GO terms.

Mapping all known regulatory elements with human and mouse UTRs.

Running structural alignments using Foldalign

Clustering of results based on Foldalign scores

Comparison with orthologous genes (Mouse UTRs)

Acknowledgements

My sincere thanks to

Jakob Hull Havgaard

Jan Gorodkin

THANK YOU