	V9 – Genomics data	
Progra	am for today:	
(1)	Read mapping	
(2)	SNP calling	
(3)	SNP frequencies in 1000 Genomes data -> consider overlapping genes	
(4)	Isoforms of genes (alternative splicing)	
(5)	Non-canonical translation -> not all translated sequences start with AUG	
(6)	Removing sequence redundancy	
You w lecture	ill hear much more about sequences and how to deal with them in the es by Prof. Sven Rahman and Prof. Fabian Müller.	
V9	Processing of Biological Data WS 2021/22	1

Today, in lecture #9, we will deal with several topics around genomic sequences.

Points (1) and (2) deal with the two most common tasks in sequence analysis: mapping of reads and identification of single nucleotide polymorphisms (SNPs).

Sometimes, we need to take special care when preparing a dataset for a statistical analysis.

E.g. we will mention the issue of overlapping genes in point (3) and the issue of removing sequence redundancy in point (6).

In point (5), we will touch on a point that you may consider for granted: translation starts "always" with a AUG codon that is translated into a methionin (M) amino acid. This is what you read in molecular biology textbooks. It turns out that this is not always the case.

In point (4), we briefly comment on the importance of mRNA and protein isoforms that result from alternative splicing.



Unless we talk about de novo assembly of a genomic sequence, the first task in an NGS project is usually the alignment of sequencing reads to an existing reference genome.

Listed here are some workflows where read mapping is a crucial part.



These are some of the well-known software tools used for mapping of NGS reads.

Read mapping techniques: (1) Hash tables							
For most of the existing tools, the mapping process starts by building an index either for the reference genome or for the reads.							
Then, the index is used to find the corresponding genomic positions for each read.							
There are two main types of techniques for this: Hash tables + BWT							
(1) The hash based methods e	ither hash the	reads or the gen	ome.				
The main idea for both types is reads/genome.	to build a ha s	sh table for subs	equences of the				
The key of each entry is a sub	sequence						
while the value is a list of posit	ions						
where the subsequence can be	e found.						
	Key	Hashed index	Genomic location				
Hatem et al.	"GCTAGC"	Key1	Chr1 123412				
BMC Bioinformatics (2013) 14:184	"TTTAGC"	KeyN	Chr6 988472				
BMC Bioinformatics (2013) 14:184 "TTTAGC" KeyN Chr6 988472 V9 Processing of Biological Data WS 2021/22 4							

In principle, one could simply scan the genomic sequence for each read sequence. However, this would be quite inefficient.

Therefore, the existing tools typically construct an index either for the reference genome or for the reads. This index is then used during the string search.

The first type of indexing techniques use a **hash table**, see the table shown on the bottom right.

In this example, the genomic sequence is indexed. Different 6-letter words are each given a hash index and where they are located in the genome.

Read mapping techniqu	ues: (2) Bu	irrows Wheeler	transform					
The BWT of the string T = "abracadabra\$" is "ard\$rcaaaabb. It is represented by the matrix M where each row is a rotation of the text, and the rows have been sorted lexicographically.								
The transform corresponds to the last column labeled L.								
Modern alignments use an extension of BWT named FM index after Ferragina & Manzina	I 1 2 3 4 5 6 7 8 9 10 11 12	F \$ abracadabr a \$abracadab a bra\$abraca a bracadabra a cadabra\$ab a dabra\$abra b ra\$abracad b racadabra\$ c adabra\$abr d abra\$abrac r a\$abracada r acadabra\$a	L a r d \$ r c a a a b b					
www.wikipedia.org								
V9 Processin	ng of Biological Dat	a WS 2021/22	5					

The second technique does not index the string itself, but uses its so-called **Burrows Wheeler transform (BWT)**.

According to Wikipedia, the Burrows–Wheeler transform is an algorithm to prepare data for use with data compression techniques such as bzip2. It was invented by Michael Burrows and David Wheeler in 1994.

The algorithm can be implemented efficiently using a suffix array thus reaching linear time complexity.

An **FM-index** is a compressed full-text substring index based on the Burrows– Wheeler transform, with some similarities to the suffix array. It was created by Paolo Ferragina and Giovanni Manzini.

Read mapping techniques: (2) Burrows Wheeler transform

C[c] is a table that, for each character c in the alphabet, contains the number of occurrences of lexically smaller characters in the text.

C	\$	а	b	с	d	r
C[c]	0 (no character is smaller than \$)	1 (1 time \$)	6 (5 times a plus 1 time \$)	8	9	10

	а	r	d	\$	r	с	а	а	а	а	b	b
	1	2	3	4	5	6	7	8	9	10	11	12
\$	0	0	0	1	1	1	1	1	1	1	1	1
а	1	1	1	1	1	1	2	3	4	5	5	5
b	0	0	0	0	0	0	0	0	0	0	1	2
с	0	0	0	0	0	1	1	1	1	1	1	1
d	0	0	1	1	1	1	1	1	1	1	1	1
r	0	1	1	1	2	2	2	2	2	2	2	2

These are two auxiliary tables used to construct the FM index.



According to

http://pages.di.unipi.it/ferragina/Libraries/fmindexV2/index.html

The FM-index combines compression and indexing by encapsulating in a single compressed file both the original file plus some *indexing information*. The space occupancy of the FM-index is close to the one required by the best known compressors, like bzip2. But additionally to a compressor, the FM-index is able to efficiently support substring search operations, and the decompression of portions of the original file. Every such operation is executed on the FM-index by looking *only at a small portion of the compressed file*, thus requiring few *milliseconds* on a commodity PC over files of several megabytes.



One complication in read alignment is that we are not only looking at positions that align perfectly or exactly. Often, the two sequence may differ "a little bit" due to either the normal biological variation between an individual and the reference genome or due to technical sequencing errors.

The genetic difference between individual humans today is minuscule – on average about 0.1% of all positions (https://humanorigins.si.edu/evidence/genetics).

Based on this, we instruct the alignment algorithm to search for almost perfectly matching positions of read and reference genome and allow for a small given number of **mismatches**.

From Hatem et al paper: *"Seeding* represents the first few tens of base pairs of a read. The **seed** part of a read is expected to contain less erroneous characters due to the specifics of the NGS technologies. Therefore, the seeding property is mostly used to maximize performance and accuracy."

The sequence in b extracted from. Aft to the original locat	Read alignment: evaluation criteria olue is the original genomic position where the simulated read ter applying sequencing errors, the read does not exactly matches tion (3 mismatches marked in red).	was :h
Reference		
Read	C C G C C G G G A A	
3 possible alignme Reference Alignments	ent locations for the read with their mapping quality score (MQ CCCGCCGGAAATTCCGCCGGGAA IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII). 2=50
Naïve criterion: on Hatem et al.	aly consider the alignment (1) as the correct alignment.	
V0	Processing of Piological Data WS 2021/22	0
¥7	Frocessing of Biological Data w 5 2021/22	9

Link to Hatem et al paper:

https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-14-184

Alignment (3) is a perfect match of the altered string to an alternative position in the genome. In this position, the match achieves the highest mapping score MQ = 50. MQ stands for "mapping quality".

Alignment (2) is a position shifted by one base pair where it has only 1 mismatch G-A and a slightly better score (MQ = 45) than in the original position (MQ = 40).

The naïve criteria would judge the tool as incorrectly mapping the read if the tool returned either alignment (2) or (3) while in fact it picked a more accurate matching.



Link to the Li & Durbin paper: https://pubmed.ncbi.nlm.nih.gov/18714091/

Phred values were introduced by Ewing & Green (https://genome.cshlp.org/content/8/3/186.full) as quality values for base calling:

 $q = -10 \ x \ \log_{10}(p)$

where p is the estimated error probability for that base-call. Thus a base-call having a probability of 1/1000 of being incorrect is assigned a quality value of 30.

Text of this slide was copied from

https://genome.sph.umich.edu/wiki/Mapping_Quality_Scores

For paired end reads, we calculate SUM_BASE_Q as the sum of base quality scores at mismatched bases for both reads.

	Read alignment: evaluation criteria	
Reference	CCCGCCGGAAATTCCGCCGGGAA	
Alignments	(1) \dot{C} \dot{C} \dot{G} \dot{C} \dot{C} \dot{G} \dot{G} \dot{G} \dot{G} \dot{A} \dot{A} MQ=40 (3) \dot{C} \dot{C} \dot{G} \dot{C} \dot{G} \dot{G} \dot{G} \dot{A} \dot{A} MQ=50 (2) $CCGCCGG$ G G AA MQ=45	
Ruffalo et al. cr	riterion: consider also the mapping quality.	
If the used thre incorrectly map	eshold is 30, then (1) is <i>correctly mapped</i> while (2) and (3) are <i>oped-strict</i> .	
If the threshold correct mappin	d is 40, then (3) is considered as <i>incorrectly mapped relaxed</i> (no ng available higher than the threshold).	
Holtgrewe et al	I. criterion: considers all matches with distance k.	
Here, it would o would be consi	detect (1) and (2) and consider them <i>correctly mapped</i> while (3) idered as <i>incorrectly mapped</i> .	
Hatem <i>et al:</i> "W violating the ma	Ve define a read to be correctly mapped if it is mapped while not apping criteria."	
V9 Hatem et al. BMC Bioinform	Processing of Biological Data WS 2021/22 matics (2013) 14:184	11

Ruffalo et al.

(https://academic.oup.com/bioinformatics/article/27/20/2790/201940) classify the accuracy of the mapping(s) of a read as follows.

Correctly mapped read (CM): the read is mapped to the correct location in the genome and its quality score is greater than or equal to the threshold.

Incorrectly mapped read—*strict (IM-S)*: the read is mapped to an incorrect location in the genome and its quality score is greater than or equal to the threshold.

Incorrectly mapped read—relaxed (IM-R): the read is mapped to an incorrect location in the genome, its quality score is greater than or equal to the threshold and there is no correct alignment for that read with quality score higher than the threshold.



Shown are the mapping results using the default options of each tool. The tools try to use the options that yield a good performance while maintaining a good output quality.

For instance, Bowtie achieves a throughput of around $1.6 \cdot 10^5$ bps/s at the expense of mapping only 67.58% of the reads. On the other hand, BWA maps 91% of the reads at the expense of having only a throughput of $0.1 \cdot 10^5$ bps/s. Additionally, SOAP and mrsFAST look like they provide the smallest mapping. However, they are only allowing 2 mismatches while other tools such as mrFAST and GSNAP are allowing more than 5 mismatches. Therefore, using only the default options to build our conclusions would be misleading. Indeed, further experiments show that BWA obtains a high throughput when allowed to use the same options as Bowtie. Moreover, BWA achieves a higher throughput than Bowtie in other experiments. Therefore, it is important to use the same options to truly understand how the tools behave.



For synthetic data generated with the software wgsim, quality thresholds of 60, 80, 100, 120, and 140 should correspond to 3, 4, 5, 6, and 7 mismatches. Here, all tools were allowed a maximum of 7 mismatches while using a quality threshold of 140. The figure shows that the tools map the reads with the same maximum number of mismatches while having similar mapping rates.

The differences in the mapping rates shown in the previous slide are due to the pruning of the search space done by the default options for some of the tools. In addition, other tools incorrectly mapped some of the reads causing an increase in the mapping percentage.

From the throughput point of view, the tools behave differently. For instance, Bowtie, MAQ, RMAP, and mrsFAST are able to maintain almost the same throughput while the throughput increases for SOAP2 and GSNAP and decreases for BWA. The degradation in BWA's performance is due to exceeding the default number of mismatches leading to excessive backtracking to find mismatch locations.



Longer reads tend to have more mismatches beside requiring more time to be fully mapped. In general, for a fixed number of mismatches, increasing the read length decreases the percentage of mapped reads. Therefore, the aim of this experiment is to understand the read length effect.

The figure shows that the mapping percentage decreases with the increase in the read length while the *error* percentage increases.

Bowtie, Bowtie2, and BWA were the only short sequence mapping tools that managed to map long reads. In particular, the max read length was 128 for MAQ, 300 for RMAP, and 200 for GSNAP, 199 for mrsFAST, while SOAP2 took more than 24 hours to map the reads with length 300 and hence not reported.

From the throughput point of view, tools do not maintain the same behavior. For instance, the throughput of Bowtie and SOAP2 decreases for long read lengths. This is due to the backtracking property and the split strategy used by Bowtie and SOAP2, respectively, to find inexact matches.

	Tools	accurately detected SNPs		
	Bowtie	1171		
	Bowtie2	2035		
	BWA	2067		
	SOAP2	1941		
	Novoalign	941		
	GSNAP	2602		
from the Spretu Then Partek wa A quality thresh tools were allow GSNAP detected	is mouse strain. as used to detect hold of 70 was us ved 5 mismatche ed the largest nur	SNPs against mouse genome ed for Bowtie and Novoalign w s. nber of accurate SNPs while N	version mm9. hile the remaining ovoalign detected	
the smallest.				
V9 Hatem et al. BMC Bioinform	Processin natics (2013) 14:184	ng of Biological Data WS 2021/22		15

This example illustrates that using a different mapping tool can greatly affect the number of obtained results.

Read alignment: conclusion	
Mapping of short sequences is still subject of active development.	
Genome indexing tools performed better than read indexing tools.	
In general, there is no <i>best tool</i> among all of the tools; each tool was <i>the-best</i> in certain conditions.	
Hatem et al. BMC Bioinformatics (2013) 14:184	
V9 Processing of Biological Data WS 2021/22	16

(2) Variant calling benchmark
-> Accurately detecting SNPs is critical e.g. for medical diagnostics.
Genome in a Bottle (GIAB) consortium:
public-private-academic consortium to develop the technical infrastructure (reference standards, reference methods, and reference data) to enable translation of whole human genome sequencing to clinical practice.
GIAB generated a set of highly confident variant calls for one individual in the 1000 Genome project:
they integrated 14 variant data sets from 5 NGS technologies, 7 read mappers and 3 variant calling methods, and manually cleaned up discordant data sets.
This highly accurate set of SNP and indel genotype calls can be used as gold standard variant genotype data set for systematic comparisons of variant callers.
Hwang et al., Scientific Reports 5, 17875 (2015)
V9 Processing of Biological Data WS 2021/22 17

Website of the GIAB consortium: https://www.nist.gov/programs-projects/genome-bottle

GIAB publication: https://www.nature.com/articles/sdata201625



Hwang et al paper: https://www.nature.com/articles/srep17875

To compare the overall performance among thirteen pipelines, the authors compared the distributions of APR scores of multiple data sets for each pipeline on SNPs and indels.

The Ion Proton data set has much lower exome coverage ($<10 \times$) than those of Illumina data sets ($43.6 \times -298.5 \times$).

For SNP variant calls, BWA-MEM-Samtools pipeline showed the best performance and Freebayes showed good performance across all aligners for both Illumina platforms.

For Ion Proton data, Samtools outperformed all other callers, including TVC, which is the Ion Proton's own variant calling method. Interestingly, the best variant caller of each data set varies. This observation of variation in best performed pipelines across data sets clearly demonstrates a data-specific effect of benchmarking results. Therefore, benchmarking performance of each variant calling pipeline needs to be based on multiple data sets to avoid misleading conclusions. The tested variant pipelines showed larger performance difference in calling indels. For indel calls, GATK-HC with any aligner outperformed Freebayes and Samtools on both Illumina platforms, while Samtools performed best on Ion Proton data. Although TVC is the official variant caller for Ion Proton data, it performed no better than other

callers on both SNPs and indels.



The authors then assessed the concordance (overlap) among the four variant callers for each NGS platform.

For Illumina data sets, they observed ~92% of concordance among the variant calls by three variant callers (see GATK-HC \cap Samtools \cap Freebayes) based on the average score of data sets. Concordance levels among variant calling pipelines varied across the data sets (82~97% overlap of called variants). These results indicate that not only the variant calling pipelines but also the data sets affect concordance of the identified variants. Therefore, caution is advised in interpreting concordance levels based on a single data set.

For Ion Proton data set, four callers showed 15.5% of overlap for the same quality score threshold (see GATK-HC \cap Samtools \cap Freebayes \cap TVC). This low overlap among called variants is likely to originate from the high false positive rates for calling indel variants by Freebayes and Samtools.

Variant calling: recommendation
The authors recommend the use of BWA-MEM and Samtools pipeline for SNP calls and BWA-MEM and GATK-HC pipeline for indel calls.
Low coverage data is not suitable for reliable SNP calling.
Indels are detected at lower accuracy than SNPs.
Hwang et al. Scientific Reports
5, 17875 (2015)
V9 Processing of Biological Data WS 2021/22 20

Concluding remarks by the authors.



In the third application, we wanted to characterize the locations of SNPs found in genomes. We used the largest public data source available, the 1000 Genomes project which in fact sequenced the genomes of around 2500 individuals from the countries marked on the map.



We focused on the European super-population with ca. 500 individuals.

The reason was that we also analyzed in parallel the 500 parent genomes from the "Genomes of the Netherland" project. The results for both data sets were very similar (data not shown).

We felt that data for the European cohort from the 1000G project would be more compatible with the GoNL data.

Also, we omitted the sex chromosomes X and Y because they appear to behave differently from autosomes.

E.g. the International SNP Map Working Group

(https://www.nature.com/articles/35057149) found that the sex chromosomes have a lower diversity than autosomes. They suggested that the lower rate of polymorphism on the X chromosome may be explained by a lower effective population size, a lower mutation rate or by strong selection acting on the sex chromosomes in males.

Also, we filtered for genes annotated to have more than one allele, excluded SNO and MIR genes, or erroneous genes.



Link to Veeramachaneni et al. paper:

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC327103/

According to Veeramachaneni et al., it is believed that 3.2 billion bp of the human genome harbor ~35,000 protein-coding genes. On average, one could expect one gene per 300,000 nucleotides (nt).

Although the distribution of the genes in the human genome is not random, it is rather surprising that a large number of genes overlap in the mammalian genomes.

Veeramachaneni et al. identified >774 gene pairs sharing a locus in the human genome and 542 in the mouse genome.

	Overlapping genes	
One could speculate tha species than non-overla region would cause cha	at overlapping genes would be more conserved between pping genes because a mutation in the overlapping nges in both genes.	
Then, one would expect stronger.	that evolutionary selection against these mutations is	
However, Veeramachan	eni <i>et al</i> . found that this is not the case.	
Overlapping human and overlapping genes.	I mouse genes were similarly conserved as non-	
Veeramachaneni et al. Genome Res. (2004) 14: 280	D-286	
V9	Processing of Biological Data WS 2021/22	24

The origin of overlapping genes is not clear. Interestingly, the mutation rates in the overlap regions are similar to non-overlap regions.

How to deal with overlapping genes In the case of overlapping genes, it is problematic to define the genomic region because they have a different meaning for the 2 overlapping genes.	S
Therefore, we distinguished 2 cases: (1) Overlaps where one gene is located inside another gene . Such genes inside other genes were excluded from the SNP analysis. (2) staggered overlaps (genes overlap partially).	
We collected all genes with staggered overlap. From each "bundle", only one gene was selected randomly to avoid overlapping genes. In total, about 5% of all genes were removed due to overlaps.	
Neininger K, Marschall T, Helms V (2019). PLoS ONE 14(4): e0214816 V9 Processing of Biological Data WS 2021/22	25

We wanted to analyze the location of SNPs with respect to certain genomic regions, see slide 27.

However, a SNP in an overlapping region may belong to different regions with respect to either of the two genes.

To avoid confusion, we excluded shorter genes that are located inside longer genes and we randomly selected one of the genes showing staggered overlaps.

Since we had "enough" data for our analysis, we rather prefered to analyze a "purified" data set.

Refseq The Reference Sequence (RefSeq) collection at NCBI provides a comprehensive, integrated, non-redundant, well-annotated set of sequences, including genomic DNA, transcripts, and proteins.									
BofSeg transprint and protein records are generated in different waves									
- Computation	Eukarvotic Genome Annotation Pineline								
Computation	Prokarvotic Genome Annotation Pipeline								
- Manual curation									
- Propagation from	m annotated genomes that are submitted to members of the leotide Sequence Database Collaboration (INSDC)								
Research quest	ion:								
Are the Single N genomic regions	ucleotide Polymorphism (SNP) frequencies in different similar to eachother or not?								
https://www.pchi.plm	nih anv/refead/ahout/								
V9	Processing of Biological Data WS 2021/22	26							

The RefSeq annotation from NCBI is a very comprehensive, sophisticated and reliable annotation source for the location of genes and exons.



Based on the Refseq annotations, we analyzed the frequency of transition and transversion SNPs as well as indels in nine types (regions) of coding and non-coding genomic elements in the human genome.



Considering 1000G data, median SNP densities were ~ 8-9 SNPs per kb for each genomic element and all variant types.

<u>Protein-coding regions are conserved with a median SNP density of about 7</u> SNPs/kb for all SNP types. The boxplot for the 5' UTR contains some outliers with a maximum SNP density of up to about 35 SNPs per kb for 1000G data. This effect is due to the short 5' UTR length of 230 bp on average (median 180 bp).

Our findings of smaller SNP densities in genetically important gene regions such as coding exons or 5' UTRs are compatible with purifying selection to preserve their functionality.



Transitions refer to point mutations that change a purine nucleotide to another purine $(A \leftrightarrow G)$, or a pyrimidine nucleotide to another pyrimidine $(C \leftrightarrow T)$. Approximately two out of three single nucleotide polymorphisms (SNPs) are transitions. **Transversions** interchange a purine with a pyrimidine and are less frequent. This was also observed by us.

Indels might have more severe effects on transcription factor binding sites than base exchanges. Hence, the low frequency of indels in CpG islands might be related to a strict conservation of functional sequences within this genomic (regulatory) element especially in CpG islands in the promoter regions of the mammalian genes.

(4) Isoforms of genes	
Gene isoforms are mRNAs that are produced from the same locus but are different in their	
- transcription start sites (TSSs),	
 protein coding DNA sequences (CDSs) and/or 	
- untranslated regions (UTRs),	
All these processes may potentially alter gene function.	
Alternative splicing (AS) of mRNA can generate a wide range of mature RNA transcripts.	
It is estimated that AS of pre-mRNA occurs in 95% of multi-exon human genes.	
There is abundant evidence for the expression of multiple transcripts in cells.	
However, it is less clear whether these transcripts are expressed more or less equally across tissues or whether it would be biologically relevant to designate one transcript per gene as dominant and the rest as alternative .	
V9 Processing of Biological Data WS 2021/22	30

Alternative splicing is a promiment mechanism to enlarge the complexity of gene regulation.

About 95% of all genes with more than one exon are alternatively spliced.

One important question is now whether (1) all these isoforms will be expressed in one tissue, (2) only some of them, or (3) only one of them.

Detect isoforms in proteomic data Ezkurdia <i>et al.</i> re-analyzed 8 HT proteomics MS data sets.
At least 2 peptides were detected for 12 716 (63.9%) of the protein-coding genes but alternative protein isoforms only for 246 genes (1.2%).
\rightarrow the vast majority of genes had peptide evidence for just one protein isoform .
The isoform with the highest number of peptides was the main proteomics isoform.
A unique main proteomics isoform was identified for 5011 genes.
Ezkurdia et al J Proteome Res. (2015) 14: 1880–1887.
V9 Processing of Biological Data WS 2021/22 31

Probably this issue has not been completely settled yet.

For the moment, it is safe to assume that there will exist one major protein isoform in each tissue.

(5) Alternative t human chimpanzee gibbon dag LPGGAPEEE rat RSSDIDAQQ mouse GAPETQAQQ Chinese hamster guinea pig cow rabbit African clawed frog trout red swamp crawfish zebrafish pufferfish	CARACTERISTICS CONTROLOGY CONTROL CONT	GRRQETGPLQGGGGPALGGADVAPRLSPVRVWPRPQAPKEPALHPMGLSLPK GRRQETGPLQGEGGPALGGADVAPRLSPVRVWPRPQAPKEPALHPMGLSLPK GRRQETGPLQGEGRPALGGADVAPRLSPVRVWPRPQAPKEPALHPMGLPLPK DGRQETGPLQGEGRPALGGADVAPRLSPVRVWPRPQAPKEPALHSMGLPLPK DGRQETGPLQGEGRPALGADVAPRLSPVRVWPRPQAPKEPALHSMGLPLK GRTQETGPLQGEGRPALGADVAPRLSPVRVWPRPQAPKEPALHSMGLPLK GRTQETGPLQGEGRPALGADVAPRSPVGVWRVQPQPKEPALBSMGLPLK ARSQETGPLQREGRPALGANVAPRSPVGVWRVQPQPKEPALBSMGLSLSK GRTQETGPLQREGRPALGANVAPRSPVGVWRVQPQPKEPALBSMGLSLSK ARSQETGPLQREGRPALGANVAPRSPVGVWRVQPPKEPALBSMGLSLSK CFLQETGPLQAEDRPAFGAMVAPRFSPVGVWRVQPPKEPALBSKGLSLSK STAHTPFSNRAGGKWPWTL FLKSA*RCMFP*VLTV*EF*RINCTLL*KPFQDSPKEPALAGSMGKWPMTL VHLFSSVLDIFCSPSTSLVWKTIRDSGLLLPFKVESPGVR-MSPSLAR GCPPADKQTCYSSVTKITLGESI*-DFCKSCWSRCPPEI-MPPAISA KDISLVCWIFFSPPLLIVWTEDVQG*WSVTFVV*GVPAPASHSPSLAR					
MUSCLE multiple seq	uence alignment of the	The mammalian sequences upstrear	n of				
translated 5'-UTR of T	RPV6	the first AUG codon are conserved, but th					
Identical aa residues (compared with the	codon. In contrast, sequences from t	stop he				
human sequence) are	shaded;	other organisms contain several stop)				
annotated N termini w	ith the first Met ⁺¹ are in	codons upstream of the annotated A and are not conserved. Sequence ide	UG entity				
rea;		is highest among the 40 amino acids					
* : stop codon in frame	9	upstream of the first Met residue (pos	sition				
– : gap	Fecher-Trost et al. J. Biol. Chem. (2013) 288: 16629	+1). This suggests that translation in mammals may start at a non-AUG					
V9	Processing of Biological	Data WS 2021/22	32				

Now we come to something that you may not have expected.

Sometimes, protein translation may start at a non-AUG codon. This is called "alternative translation".

Shown here is an alignment of the calcium channel protein TrpV6 from different species.

The red colored sequence region on the right is annotated in databases as the protein-coding region.

It is surprising to find that the sequence upstream of the translation start site is highly conserved and extends 40 amino acids upstream.



The group of Prof. Veit Flockerzi from Homburg discovered some years ago that the TrpV6 protein is 40 amino acids longer than they and the rest of the world previously thought.

In principle, this could have drastic consequences. Fortunately, for them, it turned out that the biological properties of the TrpV6 channel that they characterized for decades using a cloned version (that was 40 amino acids too short) were practically the same as those of the full-length protein.



This slide explains the ribosome profiling protocol that was invented in the lab of Jonathan Weissman at Stanford.

In an operational cell, ribosomes will constantly bind to mRNA messenger molecules and translate them into protein sequences.

To monitor the occupancy of ribosomes, one applies small chemical molecules that act as ribosome inhibitors and stall the further processing of mRNAs.

One can imagine that the conformations of the ribosomes get "frozen" in a particular state, such as stopping a video clip.

This situation is shown in the figure below the writing "a Ribosome profiling".

The rest of the protocol is very similar to the ChIP-seq protocol.

PreTIS: predict alternative translation initiation sites 1 CGGUGAGGGU UCUCGGGCGG GGCCUGGGAC AGGCAGCUCC GGGGUCCGCG GUUUCACAUC 61 GGAAACAAAA CAGCGGCUGG UCUGGAAGGA ACCUGAGCUA CGAGCCGCGG CGGCAGCGGG 121 GCGGCGGGGA AGCGUAUACC UAAUCUGGGA GCCUGCAAGU GACAACAGCC UUUGCGGUCC 181 UUAGACAGCU UGGCCUGGAG GAGAACACAU GAAAGAAAGA ACCUCAAGAG GCUUUGUUUU 241 CUGUGAAACA GUAUUUCUAU ACAGUUGCUC CAAUGACAGA GUUACCUGCA CCGUUGUCCU 301 ACUUCCAGAA UGCACAGAUG UCUGAGGACA ACCACCUGAG CAAUACUGUA CGUAGCCAGA 361 AUGACAAUAG AGAACGGCAG GAGCACAACG ACAGACGGAG CCUUGGCCAC CCUGAGCCAU 421 Suppose that a ribosome profiling experiment detected 2 start sites for this mRNA sequence: CUG at position -78 and CUG at position -120 (blue colored codons). These start sites are then considered TP start sites. All near-cognate start sites not listed in the ribosome profiling dataset and upstream of the most downstream reported true start site are then considered TN (red colored codons). Light red colored codons : start sites not considered as false starts in the analyses since they are located downstream of the most downstream reported true start site. Grey colored downstream part : annotated CDS sequence Italic (purple) upstream part : -99 upstream window needed to calculate some features. All marked start sites (TP and TN) exhibit a surrounding window of ±99 nucleotides as well as a downstream in-frame stop codon. In total, this mRNA sequence would provide 2 true start sites and 9 false start sites out of 23 putative starts. Processing of Biological Data WS 2021/22 Reuter et al Plos Comput Biol 35 V9 (2016) 12: e10005170

Which positions are considered as potential alternative translation initiation sites (aTIS)?

In this example, AUG at position 273-275 (colored light grey) would be the annotated translation start site in the database. All subsequent sequence is translated into protein.

But there are many potential alternative start sites upstream of this AUG that differ by one nuclotide. They are termed "near-cognate" sites.

The first one is ACG at position 100-102 (colored red).

Let us assume that ribosome profiling detected two true start sites: CUG at position 153-155 (which means 120 positions upstream of the canonical start site) and CUG at position 195-197 (78 upstream).

These are colored blue. It is actually not easy to detect experimentally if both of them are used or if only the first one is used. We will ignore this complication.

All other alternative sites upstream of the first aTIS and between the two are assumed to be true negatives because they are apparently not used.

For the other aTIS candidates downstream of the second true positive aTIS site, we cannot make a statement whether they are also used or not because they are "overshadowed" by the two aTIS sites upstream of them.

Cell line	Description	Genes	Start codons	TPs	TNs	Used for
HEK293	Human embryonic kidney cells	3,566	AUG and near-cognate	4,482	49,520	Human prediction model
HEK293	Human embryonic kidney cells	391	AUG	Validation set		
Mouse ES	Mouse embryonic stem cells	1,632	AUG and near-cognate	3,009	19,864	Mouse prediction model
We o identi	nly included curated m fier (starting with NM_)	RNA s∉).	equences with ava	ilable	mRNA	RefSeq
We o identi Raw o	nly included curated m fier (starting with NM_)	RNA s∉). d (numl	equences with ava	ilable	mRNA	RefSeq
We o identi Raw o	nly included curated m fier (starting with NM_) data is very unbalance	RNA s∉). d (numl	equences with ava	ailable Is very	mRNA	RefSeq nt)
We o identi Raw o → ne	nly included curated m fier (starting with NM_) data is very unbalance ed to balance data set	RNA se). d (numl ːs (selec	equences with avan ber of TPs and TN ct random TN data	ilable Is very i points	mRNA differei s)	RefSeq nt)
We o identi Raw o → ne	nly included curated m fier (starting with NM_) data is very unbalance ed to balance data set	RNA s∉). d (numl :s (seleo	equences with ava ber of TPs and TN ct random TN data	ailable Is very a points	mRNA differei s)	RefSeq nt)
We o identi Raw o → ne	nly included curated m fier (starting with NM_) data is very unbalance ed to balance data set	RNA s∉). d (numl :s (seled	equences with ava	ailable Is very a points	mRNA differei s)	RefSeq nt)
We o identi Raw ∉ → ne	nly included curated m fier (starting with NM_) data is very unbalance ed to balance data set	RNA s∉). d (numl s (sele¢	equences with ava	ilable Is very a points	mRNA differei s)	RefSeq nt)
We o identi Raw (→ ne	nly included curated m fier (starting with NM_) data is very unbalance ed to balance data set	RNA se). d (numl :s (selec	equences with ava	ls very a points	mRNA differen इ)	RefSeq nt)

These were the available suitable ribosome profiling data sets in 2015 when we conducted this project.

The number of TNs is 7-12 fold larger than TPs. Therefore, we downsampled the TNs by randomly selecting the same number of data points.

If one would not do this, a "successful" classifier could alway predict "negative" and would achieve around 90% accuracy on an imbalanced test set, simply because there are about 10 fold more negatives in the full data set.

If one would balance the test set (50:50), then this classifier would fail completely.



This is the flowchart used to train a classifier that predicts which candidate alternative start sites are used and which ones are not.

As true positives, we used the mRNA sequences detected in ribosome profiling to be bound to ribosomes.

As true negatives, we used all remaining mRNA sequences that were not detected.

Note that both steps include assumptions. There may be different reasons why mRNAs appear to be bound although they are in fact not translated, and the opposite.

Then, we compute a large number of features for the elements of both sets. These will be explained on the next slide in more detail.

We select the 50 (out of 1252) features showing the largest differences between both datasets in order to avoid over-training, and also check for correlation between them.

		Feature	True starts	False starts	P-value
Features used	1.	S' UTR length	414.41±270.48	675.41±545.35	< 10 ⁻⁰¹⁰
i catales asea	2.	5' UTR conservation	0.4±0.16	0.33±0.16	8.2 × 10 ⁻¹⁹⁰
	3.	PWM positive	2.75±1.5	-0.14±2.82	5.5 × 10 ⁻¹⁷³
	4.	K-mer: upstream AUG	0.22±0.57	0.59±0.9	5.1 × 10 ⁻¹⁴⁴
Mean value and standard	5.	5' UTR: percentage A	0.18±0.05	0.2±0.05	9.6 × 10 ⁻¹⁰⁰
	6.	Kozak sequence context	2.67±1.07	2.3±1.11	9.2 × 10 ⁻⁹⁵
doviation of the 44 features	7.	Translational efficiency of flanking sequence	83.75±20.11	77.12±21.4	1.1 × 10 ⁻⁶⁰
Jeviation of the 44 leatures	8.	K-mer: position -12 is C	0.13±0.34	0.3±0.46	2.7 × 10 ⁻⁷⁷
	9.	K-mer: upstream Asparagine	1.25±1.37	1.61±1.61	4.0 × 10 ⁻⁴⁵
hat were used in the best	10.	K-mer: downstream AUG	1.14±1.15	0.92±1.1	9.2 × 10 ⁻⁴¹
	11.	K-mer: upstream A	17.24±7.43	18.81±7.89	4.0 × 10 ⁻⁴⁷
uman madal	12.	K-mer: in-frame upstream Alanine	3.69±2.6	3.16±2.29	4.0 × 10 ⁻³⁷
iuman model.	13.	K-mer: upstream Alanine	10.27±4.5	9.38±4.6	6.2 = 10 **
	14.	S OTH: percentage G	0.3220.06	0.31±0.05	7.1 × 10 **
	15.	Coon conservation	0.2380.42	0.1280.32	3.2 * 10 **
	10.	K-mer, postori 13 to A	0.0120/46	0.220.4	3.4 × 10 ⁻⁰
	10	K-mer, downstream CCA	2.0012/43	2.00mc.31	1.1 × 10 ⁻³²
DMM · probability weight	10.	K-mer oneition -12 is A	0.3+0.45	0.19+0.4	4.0 × 10 ⁻²²
- wiw . probability weight	20	Kimer in frame unstream Methionine	0.07+0.29	0.2+0.48	3.3 × 10 ⁻³¹
	21.	K-mer: upstream Arginine	12.15+4.34	11.3364.64	1.5 × 10 ⁻²⁹
matrix	22.	K-mer: upstream Histicine	1.7±1.52	1.97±1.65	2.2 × 10 ⁻²⁷
	23	K-mer: GOC	6.4+3.87	5.77+3.75	1.1 × 10 ⁻²⁵
$(PFM_{(m,n)})$	24.	K-mer: position 4 is G	0.37±0.48	0.28±0.45	2.3 × 10 ⁻²⁵
$PWM = log \left(\frac{11 m(nt,t)}{11}\right)$	25.	K-mer: upstream Threonine	3.56±2.08	3.91±2.19	4.9×10^{-25}
$r_{w_{1w_{nt,i}}} = log \begin{bmatrix} 1 \\ 1 \end{bmatrix}$	26.	K-mer: upstream CGG	3.14±2.51	2.77±2.41	3.2 × 10 ⁻²⁴
Dg _{nt}	27.	K-mer: upstream C	30.4±8.98	28.96±9.04	1.0 × 10 ⁻²⁰
· · · · · /	28.	K-mer: position -2 is G	0.23±0.42	0.32±0.47	1.2 × 10 ⁻¹²
	29.	K-mer: upstream Stop	2.3±1.71	2.66±2.0	1.4 × 10 ⁻²⁹
	30.	K-mer: UAG	1.34±1.2	1.57±1.35	5.6×10^{-22}
Entries of position-	31.	K-mer: upstream CAU	0.58±0.85	0.73±0.95	3.4 × 10 ⁻²²
	32.	K-mer: upstream Serine	9.44±3.29	8.93±3.14	5.7 × 10 ⁻²²
requency_matrix (PFM) ·	33.	K-mer: downstream Glutamine	3.57±2.01	3.26±1.88	2.4×10^{-21}
requeries—matrix (FFIM).	34.	K-mer: AGG	4.29±2.51	4.7±2.69	2.1 × 10 ⁻³⁰
	35.	K-mer: AGC	4.4±2.43	4.02±2.19	2.1 × 10 ⁻²⁰
sum of occurrences of a	36.	K-mer: downstream ACC	1.45±1.26	1.27±1.17	2.0 × 10 ⁻¹⁹
	37.	K-mer: UAA	1.22±1.42	1.51±1.76	6.2 × 10 ⁻¹⁹
nucleotide at position i	38.	K-mer: downstream Proline	9.3±5.63	8.56±5.47	3.5 × 10 ⁻¹⁸
	39.	K-mer: upstream CAA	0.75±0.92	0.91±1.05	1.3 × 10 ⁻¹⁷
	40.	K-mer: in-frame upstream Histidine	0.54±0.77	0.67±0.87	1.7 × 10 ⁻¹⁷
divided by the total number	41.	K-mer: upstream GAU	0.63±0.85	0.77±0.96	2.1 × 10 ⁻¹⁶
	42.	K-mer: In-frame upstream GCC	1.21±1.4	1.02±1.22	6.7 × 10 ⁻¹⁶
of sequences contained in S	43.	K-mer: in-frame upstream GCG	1.14±1.42	0.97±1.27	6.2 × 10 ⁻¹⁴
n sequences contained in S.	44.	PWM negative	1.94±1.34	1.59±1.09	1.6 × 10 ⁻⁰⁰
Pautar at al Plas Comput Pial (2016)	Mean va bold). All	Aue and standard deviation of the 44 features that were used in 14,482 true and 49,520 talse start sites were considered for the	n the best human model (bio is analysis. All listed feature	slogically-motivated and PW s showed significant different	M features are show noes between true ar
12: e10005170	language	e (scip) version 0.17.0). The PWM-scores are based on the te 71 (score) activity 1005170 4000.	est data (compare to Fig 4).) are represented as 0.0	n pyron programm
D Drosser	00: nl.13	Dialogical Data WE 2021/22	,		
9 Process	ang of	Biological Data WS 2021/22	<u>_</u>		

These are possible features by which true translation start sites and false start sites may potentially differ.

Obvious criteria are the length of the 5'UTR region and its conservation.

If the considered codon is actually a false start and real translation starts in front of it, the annotated UTR may be too long. This matches the observation that the 5'UTR in front of false starts is much longer (675 nt) than in front of true starts (414 nt).

If a UTR regions is highly conserved, this also suggests that it may in fact be translated.

The k-mer counts are raw counts in a 99 nt upstream or downstream window from the central codon.

	Accuracy	Specificity	Sensitivity	Precision	AUC	Threshold				
			HEK293							
Linear SVR	0.80±0.01	0.80±0.01	0.81±0.01	0.80±0.01	0.80±0.01	0.62±0.01				
RBF SVR	0.82±0.01	0.81±0.01	0.83±0.02	0.82±0.01	0.82±0.01	0.55±0.02				
Polynomial SVR	0.80±0.01	0.80±0.01	0.81±0.02	0.80±0.01	0.80±0.01	0.59±0.02				
Linear Regression	0.80±0.01	0.80±0.01	0.81±0.01	0.80±0.01	0.80±0.01	0.55±0.01				
			Mouse ES							
Linear SVR	0.75±0.01	0.75±0.01	0.76±0.01	0.75±0.01	0.76±0.01	0.65±0.03				
RBF SVR	0.76±0.01	0.76±0.01	0.76±0.02	0.76±0.01	0.76±0.01	0.58±0.03				
Polynomial SVR	0.75±0.02	0.75±0.01	0.76±0.02	0.75±0.02	0.75±0.02	0.62±0.03				
inear Regression	0.76±0.01	0.75±0.01	0.76±0.01	0.75±0.01	0.76±0.01	0.55±0.01				
doi:10.1371/journal.pcbi.1005170.002 All human models perform very similarly with accuracies of about 80% while the average performance of the mouse model is lower with average accuracies of about 76%,										
doi:10.1371/journal.pcbi.10 All human m while the ave accuracies c	nodels perforr erage perforr of about 76%,	m very similar nance of the i	ly with accura mouse model	acies of abou	t 80% average					

Support vector regression gave only slightly better results than standard linear regression. Hence, we used the robust linear regression for the Webserver version of PreTIS.

The accuracies for the mouse data set were slightly lower.

				Unbalan	ced datasets			
erformance of the best		Mouse ES		Mouse ES				
uman HEK293 model	Threshold	t=0.54		t = 0.52				
		TP	TN	TP	TN			
pplied to the mouse ES	Predicted positive	2,161	4,569	2,273	5,072			
ataset	Predicted negative	848	15,295	736	14,792			
alasel	Total	3,009	19,864	3,009	19,864			
	Accuracy		0.76		0.75			
	Sensitivity		0.72		0.76			
model is reasonably	Specificity		0.77		0.74			
ranafarabla	Precision		0.32	0.31				
ansierable,				Balanc	ed datasets			
uggests universal		Mouse ES		Mouse ES				
uggests universal	Threshold	t=0.54		t=0.52				
ranslation code		TP	TN	TP	TN			
	Predicted positive	2,161	689	2,273	763			
	Predicted negative	848	2,320	736	2,246			
	Total	3,009	3,009	3,009	3,009			
	Accuracy		0.74		0.75			
	Sensitivity		0.72		0.76			
	Specificity		0.77		0.75			
	Precision		0.76		0.75			
euter et al. Plos Comput Biol 016) 12: e10005170	doi:10.1371/journal.pcbi.10	05170.t004						
		D . N/C 202	1/22					

Interestingly, when applying the trained human model to the mouse embryonic stem cell data set from ribosome profiling, the results were almost as good as with the model trained on mouse data.

On the one hand, this suggests that the mouse data set is maybe not so good.

On the other hand, this suggests that the translation code in human and mouse is quite similar.



As an example for the application of PreTIS, we show here the predictions for the human gene GIMAP5.

The annotated start site is a AUG codon at position 0 (right of the shown codons).

Listed are all alternative start sites upstream of the annotated translation start: AUGs and codons differing by one nt.

For each putative alternative start site, we show the "translation initiation confidence" predicted by PreTIS's linear regression model.

The predictions are color coded according to the confidence score.

AUG at position -203 is assigned the highest score.

Mutation matrix showing the impact	(A)	CL	IG at	nos	ition	/i	rtı	Ja	al -	S	NI	Ρ	aı	na	ıly	/S	is	C	of	ge	en	ie	G	;	M/	41	>5		
sequence context of	(,,)	-1: U	5 -14 C	-13 A	-12 G	-11 U	-10 G	-9 A	-8 C	-7 U	-6 G	-5 C	-4 C	-3 A	-2 C	-1 C	1 C	2 U	3 G	4 G	5 A	6 G	7 G	8 A	9 C	10 A	11 G	12 G	13 G
4 putative start sites	A	0.80	0.80		0.80	0.83	0.82		0.73	0.84	0.82	0.81	0.84		0.85	0.83			17	0.80		0.82	0.86		0.83		0.86	0.89	0.85
of gene GIMAP5 on	ч с	0.81		0.80	0.64	0.83	0.81	0.75		0.80	0.82			0.67					17	0.78	0.81	0.82	0.86	0.79		0.80	0.82	0.81	0.83
the predicted	NS G	0.80	0.79	0.79		0.77		0.78	0.78	0.78		0.76	0.80	0.74	0.80	0.80					0.73			0.77	0.79	0.73			
initiation confidence.	U		0.76	0.78	0.83		0.81	0.82	0.80		0.84	0.83	0.81	0.70	0.83	0.80		Ĺ		0.74	0.77	0.86	0.86	0.81	0.80	0.77	0.85	0.83	0.84
In each case, only	(B)	CU	IG at	pos	ition	-44	40	~		-	~			2	2			~			-	~	-		0	10		10	10
In each case, only		-1: C	0 - 14 C	-13 A	-12 G	-11 A	-10 G	-9	-8 C	-/ U	-6 C	-5 A	-4 G	-3 U	-2 G	-1 A	C 1	2 U	G	4 C	č	ь А	ć	č	č	10 U	11 G	12 G	13 A
one nucleotide is	. A	0.45	0.49		0.57		0.49	0.55	0.49	0.49	0.51		0.46	0.66	0.61		Ž	Ž	Ē	0.54	0.52		0.52	0.50	0.54	0.51	0.54	0.56	
mutated with respec	t ° C			0.50	0.34	0.49	0.47			0.48		0.50	0.52	0.50	0.58	0.48			62	_		0.48				0.48	0.54	0.52	0.47
to the reference	NG G	0.45	0.48	0.49		0.42		0.51	0.46	0.45	0.51	0.45		0.57		0.46			6	0.60	0.44	0.49	0.47	0.45	0.47	0.45			0.48
sequence (top line).	Ű	0.51	0.49	0.48	0.52	0.47	0.48	0.56	0.49		0.49	0.51	0.49		0.55	0.45				0.50	0.46	0.51	0.50	0.50	0.50		0.56	0.53	0.50
Grey : start was																													
predicted as true	(C)	AU	IA at	pos	ition	-237	7																						
translational start		-15	5 -14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	1	2	3	4	5	6	7	8	9	10	11	12	13
(predicted initiation		U	G	G	G	G	G	A	C	A	C	А	C	0	C	C	A	U	A	A	U	C	U	C	U	А	C	U	U
confidence > 0.54).	m c	0.48	3 0.49	0.50	0.56	0.54	0.49		0.48		0.50		0.50	0.63	0.51	0.50		4			0.53	0.48	0.48	0.49	0.50		0.48	0.50	0.48
white : start was	Å C	0.4	3 0.51	0.50	0.33	0.52	0.46	0.45		0.46	0.50	0.50		0.46	0.40	0.68	4	4	K	0.44	0.54	0.40	0.47		0.48	0.45		0.47	0.46
classified as false	00	0.46	2					0.44	0.44	0.44	0.52	0.45	0.45	9.55	0.40	0.46	4		4	0.50	0.47	0.49	0.43	0.44	0.45	0.43	0.42	0.43	0.45
start.	U		0.50	0.48	0.52	0.52	0.48	0.48	0.47	0.49	0.50	0.52	0.45		0.46	0.45			/	0.45		0.51		0.49		0.47	0.49		
Mutations at the	(D)	СЦ	IG at	pos	ition	-16	0																						
start sites itself wore	` `	-15	5 -14	-13	-12	-11	-10	-9	-8	-7	-6	-5	-4	-3	-2	-1	1	2	3	4	5	6	7	8	9	10	11	12	13
start sites itself were		С	С	U	С	С	U	U	А	А	С	U	G	С	G	U	С	U	G	С	U	С	А	А	С	С	U	С	С
not considered. The	A	0.23	0.24	0.25	0.47	0.26	0.26	0.25			0.20	0.25	0.31	0.46	0.33	0.30	\angle		\mathbb{Z}	0.29	0.31	0.24			0.25	0.26	0.26	0.30	0.28
numbers reflect the	S C			0.25			0.24	0.20	0.26	0.24		0.24	0.30		0.31	0.28		\angle	\mathbb{Z}		0.29		0.24	0.24			0.23		
predicted initiation	‰ G	0.23	3 0.23	0.24	0.40	0.21	0.25	0.22	0.21	0.22	0.28	0.20		0.34		0.26	\square	\angle		0.32	0.23	0.25	0.20	0.20	0.22	0.21	0.20	0.24	0.24
confidence values	U	0.25	5 0.25		0.44	0.25			0.25	0.27	0.26		0.27	0.25	0.30		\square		\sim	0.25		0.27	0.24	0.23	0.24	0.25		0.27	0.27
V9 Reuter et al P (2016) 12: e1	los C 0005	om 170	iput)	Bja	əbc	essi	ng	of E	Biol	ogi	cal	Dat	ta V	vs	202	1/2	2											42	

Here, we tested how the PreTIS prediction changes if one nt is exchanged to an alternative nucleotide that could result e.g. from a SNP.

In the upper line, CUG at position -36, has a PreTIS score of 0.81 (see previous slide).

The largest decrease to 0.64 would result from replacing G in position -12 by C.

On the other hand, G at position +7 appears unfavorable. Replacing this by any other nt increases the score to 0.86.

The second line considers CUG at position -44 with PreTIS score of 0.50, which stands for low confidence.

Some mutations, e.g. U->A at position -3 increase the score to a much better value of 0.66.

In the last line, CUG at position -160 has a low score of 0.25. Most mutations show practically no change, except for C->A in position -12 (0.47) and C->A in position -3 (0.46).



Wang et al. identified in 29 human tissues 117 aTIS peptides mapping to 89 genes and 99 alternative translation start sites.

Fifty-five of these aTIS peptides represent 5' N-terminal extensions of the original gene, 32 peptides represent novel (acetylated) N-termini downstream of the canonical start site, 17 represent frame-shifts potentially leading to an entirely new sequence, five peptides likely represent upstream ORFs (uORF) with a stop codon before the canonical start site and 8 peptides with mixed annotation.

One can validate the existence of aTIS peptides by comparing their spectra to synthetic peptides.

Panel E shows an example for a peptide (ac)ATTQISKDELDELKEAFAK derived from the actin-binding protein plastin-3 (PLS3).

Note the lacking y1 peaks (very left) for the experimentally detected (shorter) peptide.

(6) Remov	ving sequence redundancy						
Let's assume we want to know	ow whether the amino acid composition of certain one genomic region from the other regions.						
For example, we want to known membrane proteins are more	ow whether transmembrane (TM) segments of e hydrophobic than the rest of the protein sequence						
To check this, we could simply analyze all protein sequences from NCBI, pre- the TM segments in them and compare the amino acid compositions.							
However, this search would - what proteins have been - by duplicated sequencing	likely be biased by sequenced and which ones not, and experiments.						
→ It is very important to rem This can be done by softwar	ove sequence redundancy before such analyses! e tools such as CDhit or BlastClust						
V9 Proce	essing of Biological Data WS 2021/22	44					

For many bioinformatics analyses, we need to process the considered data set and remove redundant sequences.

Here, we briefly explain for which applications this is important and mention how this can be done with tools such as CDhit or BlastClust.



The Link given at the bottom of the slide links to a page explaining how BlastClust works.

But BlastClust is apparently no longer included in the latest release of the Blast program.



Today, we addressed several points that may be relevant if you analyze genomic data sets.