

RiboAI

-

Ribo-Seq

meets

Machine Learning

by:

Denis Skibinski

supervised by:

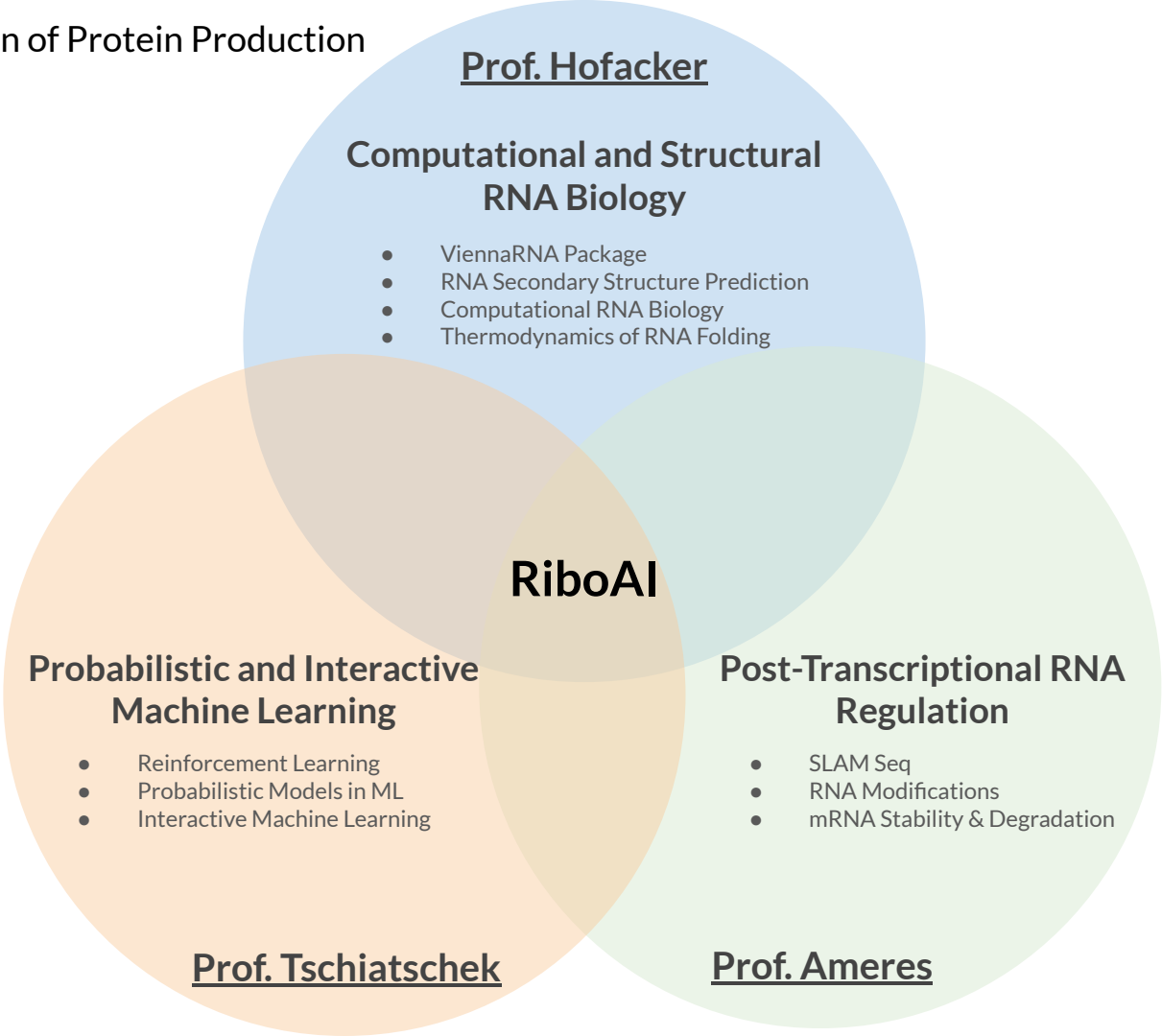
Prof. Hofacker

University Vienna

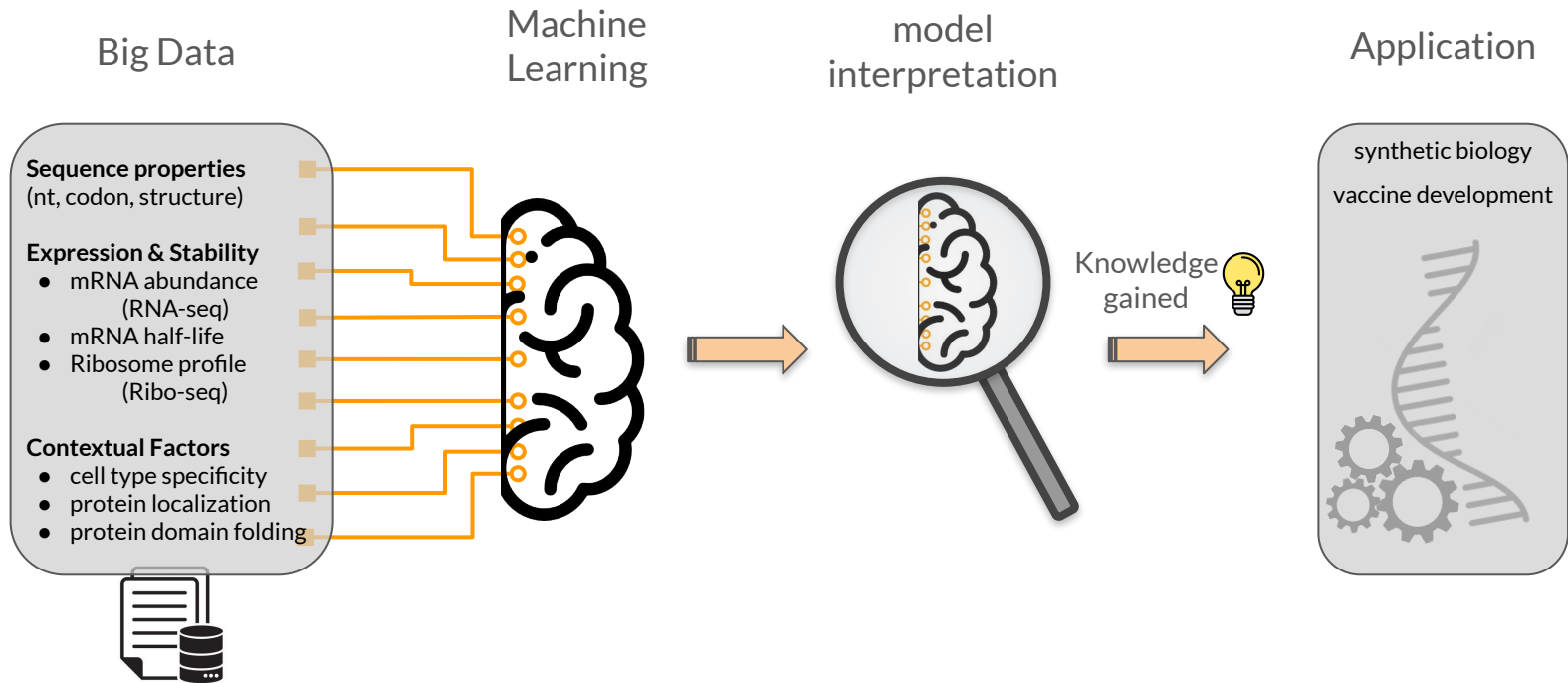
TBI

RiboAI: Data-Driven Exploration of Protein Production

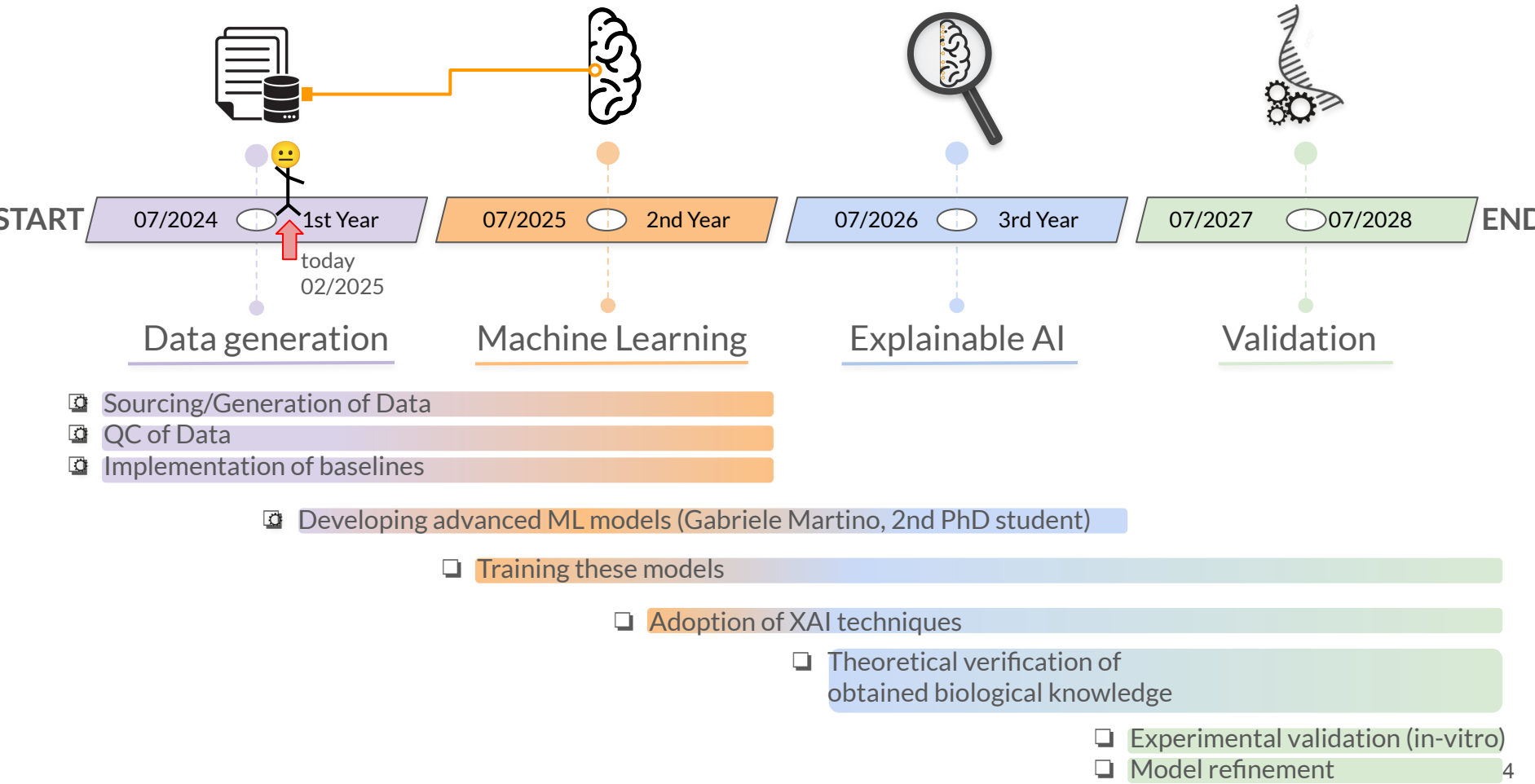
3 Professors
+
1 Senior Scientist
+
2 PhD Student
=
1 RiboAI project



RiboAI: Data-Driven Exploration of Protein Production



RiboAI - PhD Timeline



Data (Hot-)encoded, the machine vision

per transcript:

Nucleotide-Resolution Features

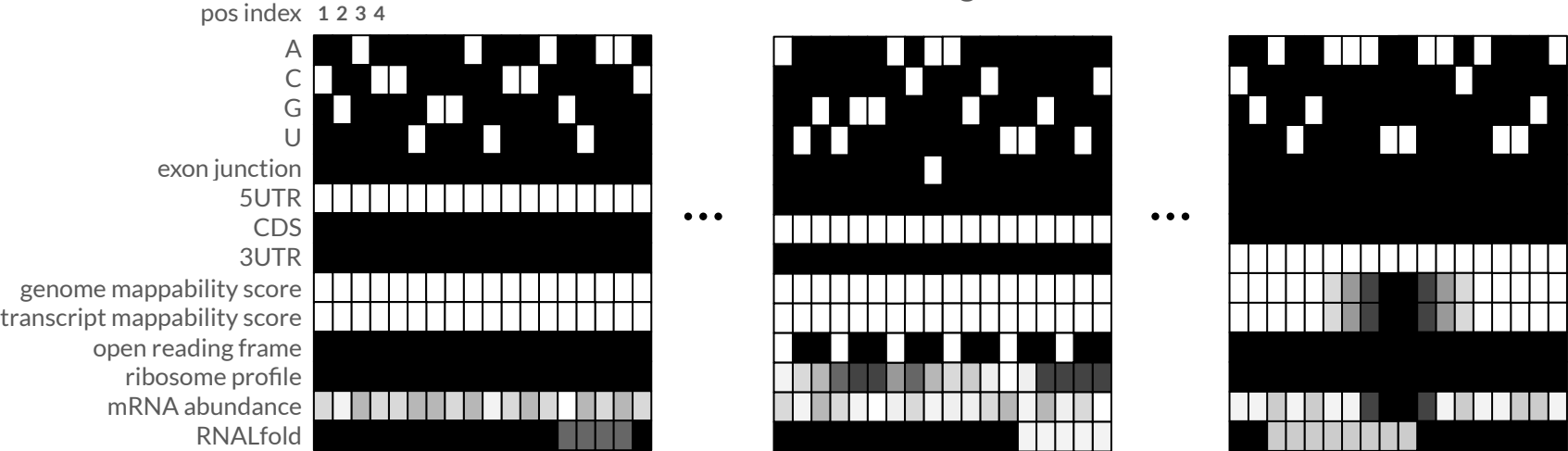
position	A	C	G	U	exon junction	5UTR	RNAfold	mRNA abundance	ribosome profile
1	0	1	0	0	0	1	-5.6	0.3	0
2	0	0	1	0	0	1	-5.6	0.8	0
3	1	0	0	0	0	1	-5.6	0.2	0
4	0	0	1	0	0	1	-5.6	1.8	0

target

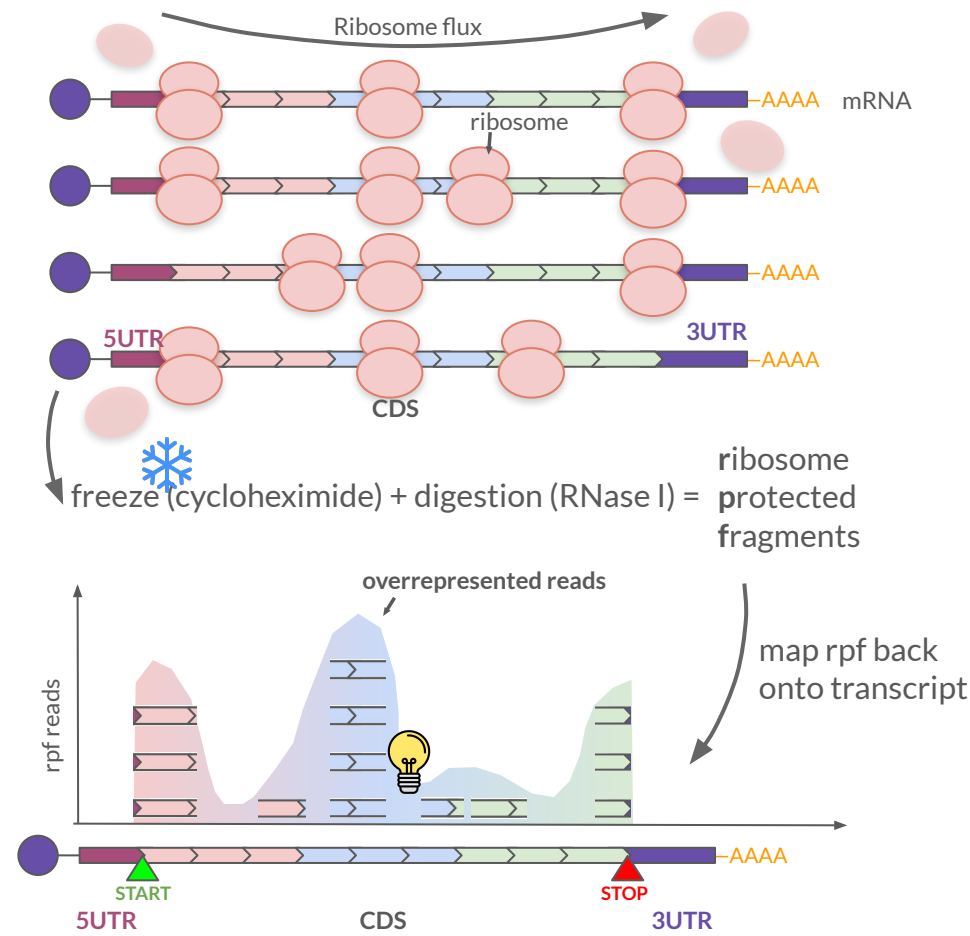
...

⋮

Embedding:



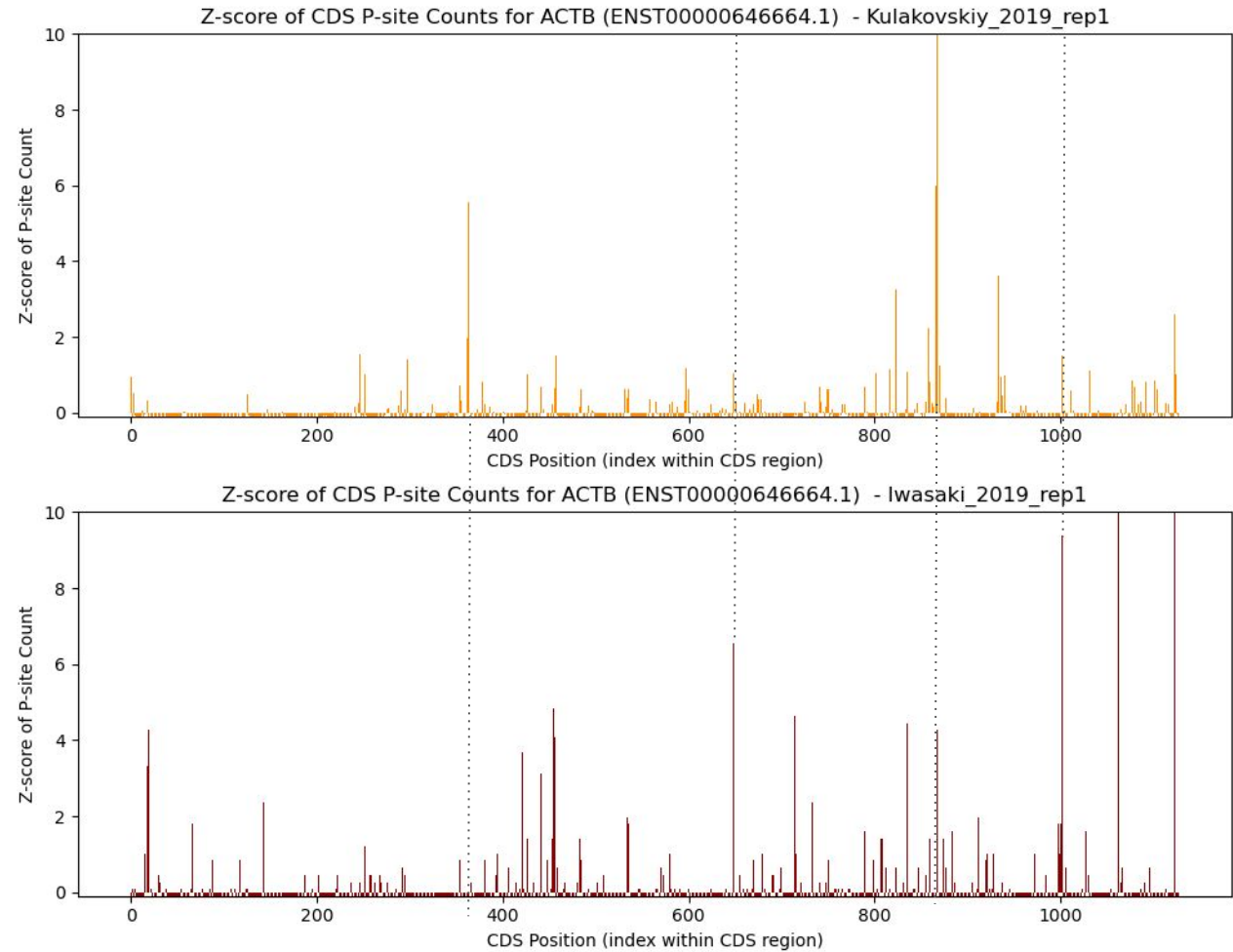
Ribosome Profile - snapshot of translation in action



How-to ribosome profile:

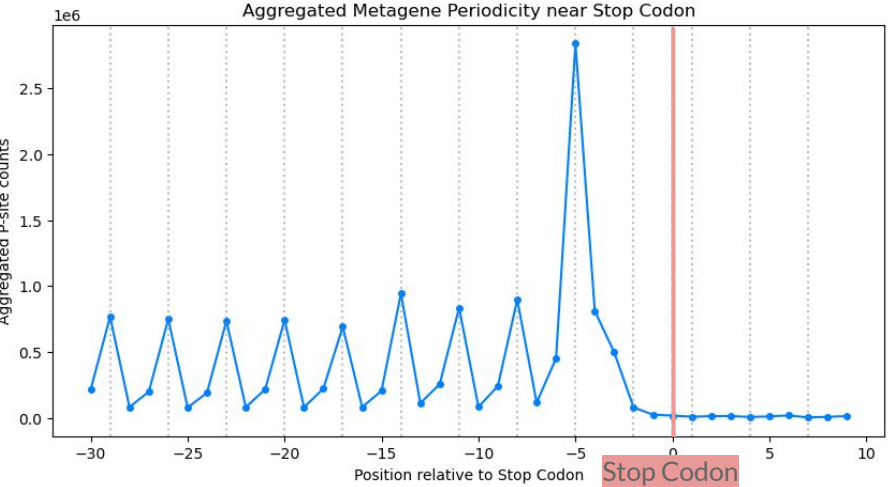
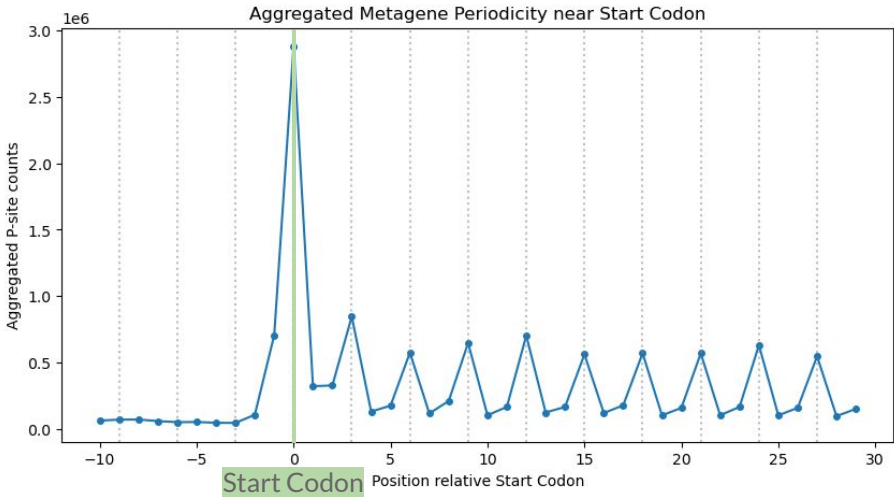
1. created by freezing/stopping ribosomes on mRNA
2. digestion of not protected mRNA around ribosome-> creating ribosome protected fragments (reads)
3. sequence, filter & map reads back onto the transcript

Ribo-seq low reproducibility - between datasets



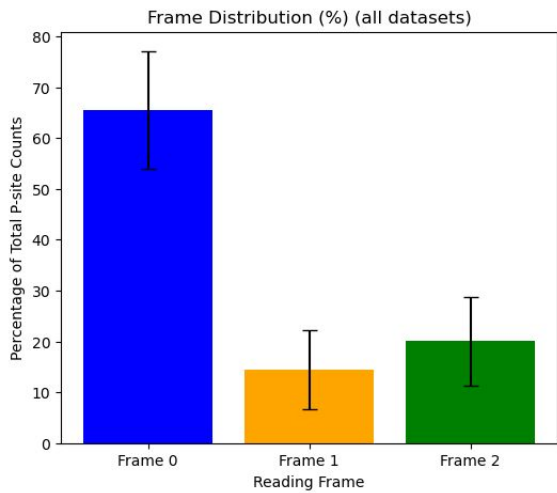
- all HEK datasets
- Low peak reproducibility around ($p=.2$) between datasets
- need to assure high quality ribosome profiles

Quality Control - Ribo-seq

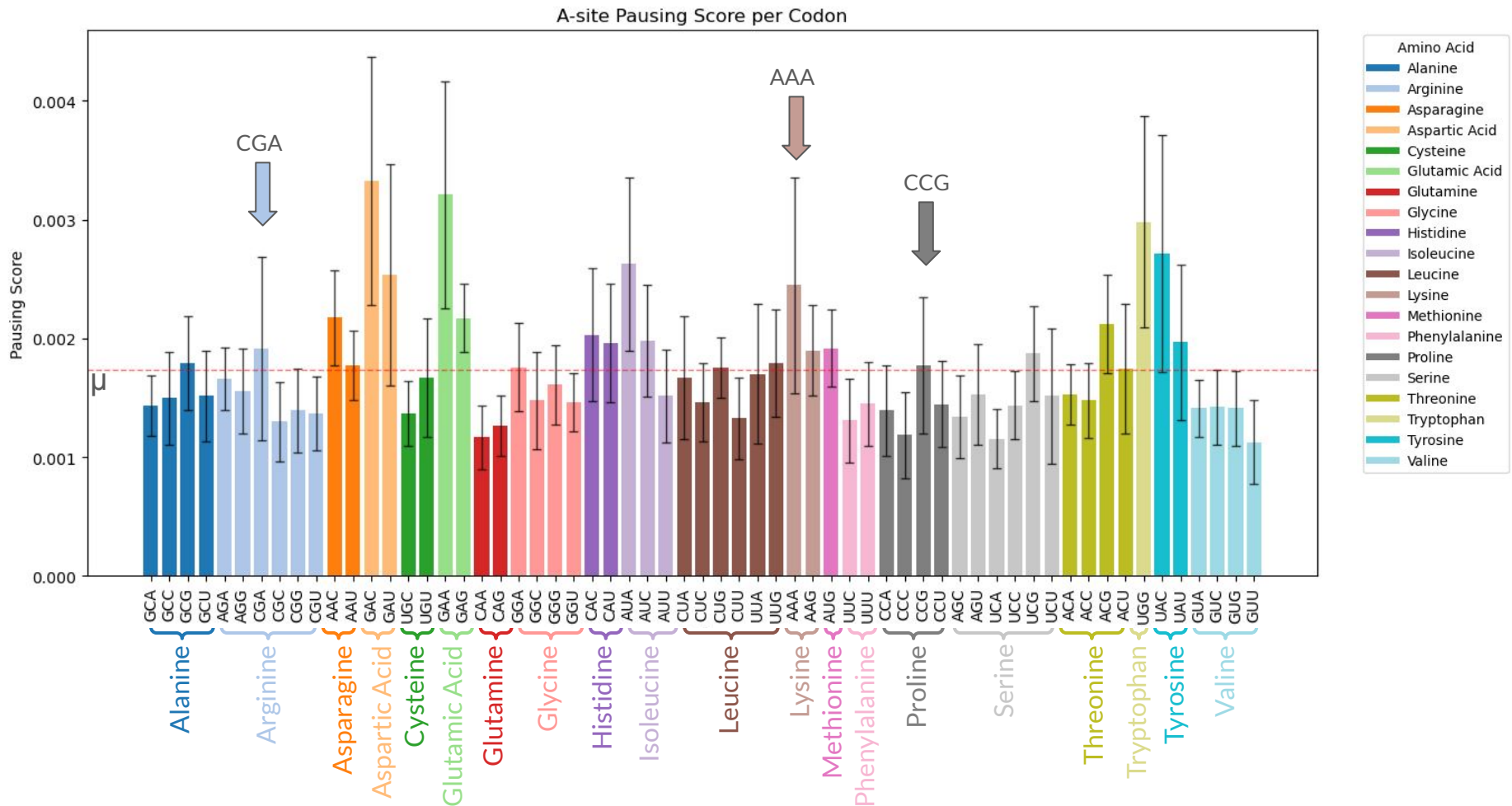


Indicator of high quality Ribosome Profile Data:

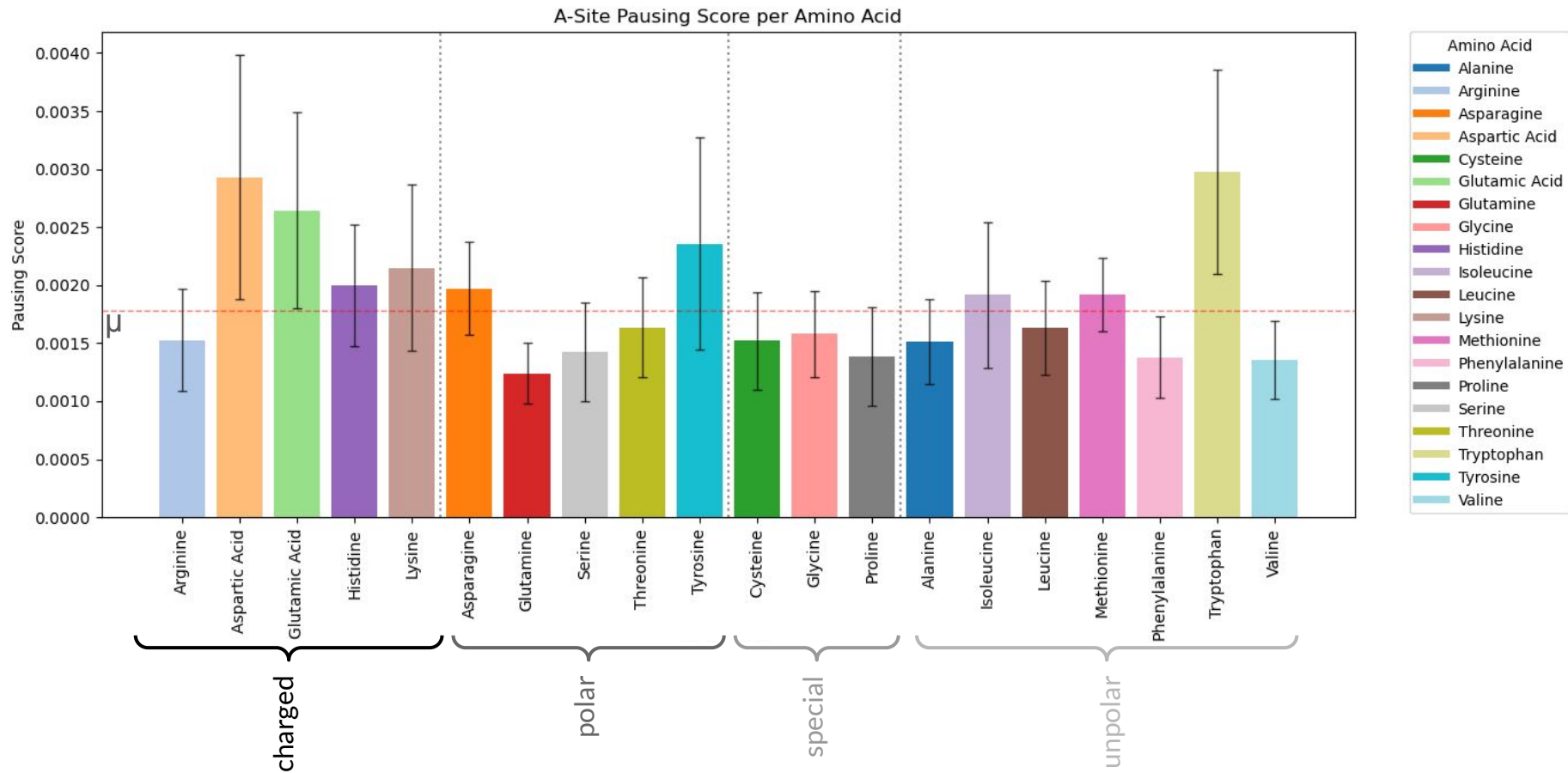
- no mapping outside of CDS
- clean cut at Start Codon/ before Stop Codon
- mean % of reads mapped in OpenReadingFrame (frame 0) > 60%
- high freq every 3nt in CDS



QC - Pausing Score - A-site Codon - sanity check



QC - Pausing Score - Amino Acid - sanity check



Translation Efficiency in Feature space

per transcript:

Nucleotide-Level-Resolution Features

position	A	C	G	U	exon border	5UTR	RNALfold	mRNA abundance	ribosome profile
1	0	1	0	0	0	1	-5.6	0.3	0
2	0	0	1	0	0	1	-5.6	0.8	0
3	1	0	0	0	0	1	-5.6	0.2	0
4	0	0	1	0	0	1	-5.6	1.8	0

calculated
out of
Ribo-seq

targets


Transcript-Level Features

Tx_ID	Ribo_RPKM	RNA_RPKM	GC_Content	A:AU_Ratio	C:CG_Ratio	Exon Junct Density	HalfLife	Translation Efficiency	Protein abundance
ENST00001	16.5	74.3	0.43	0.95	0.87	2	3.4	1.4	1.4
ENST00002	2.8	12.1	0.52	1.11	0.90	4	1.5	0.8	0.8

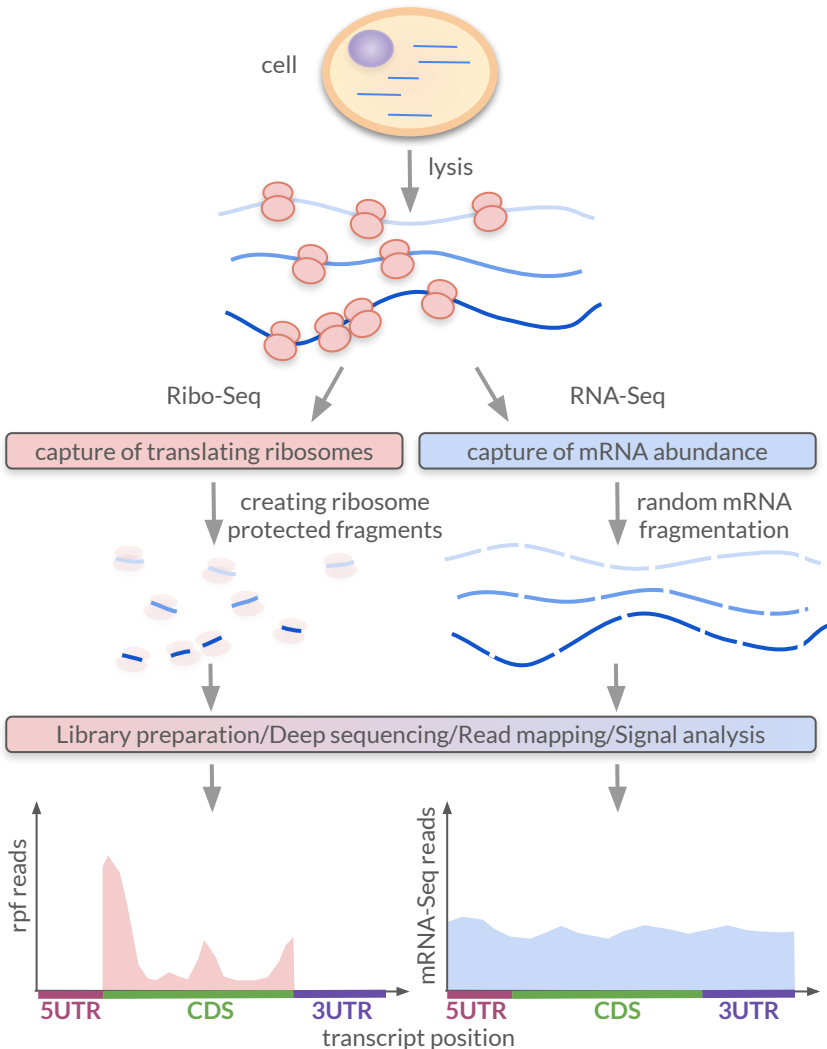
Translation Efficiency

Calculating Translation Efficiency:

$$\log(TE_{tx}) = \log\left(\frac{RiboSeq_{tx}}{RNASeq_{tx}}\right)$$

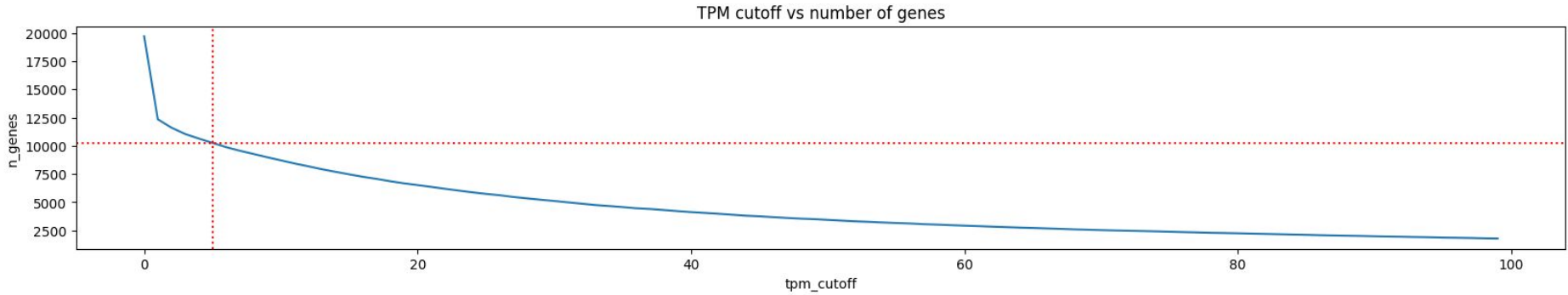
 High TE: Suggests strong protein synthesis

 Low TE: Suggest weak protein synthesis

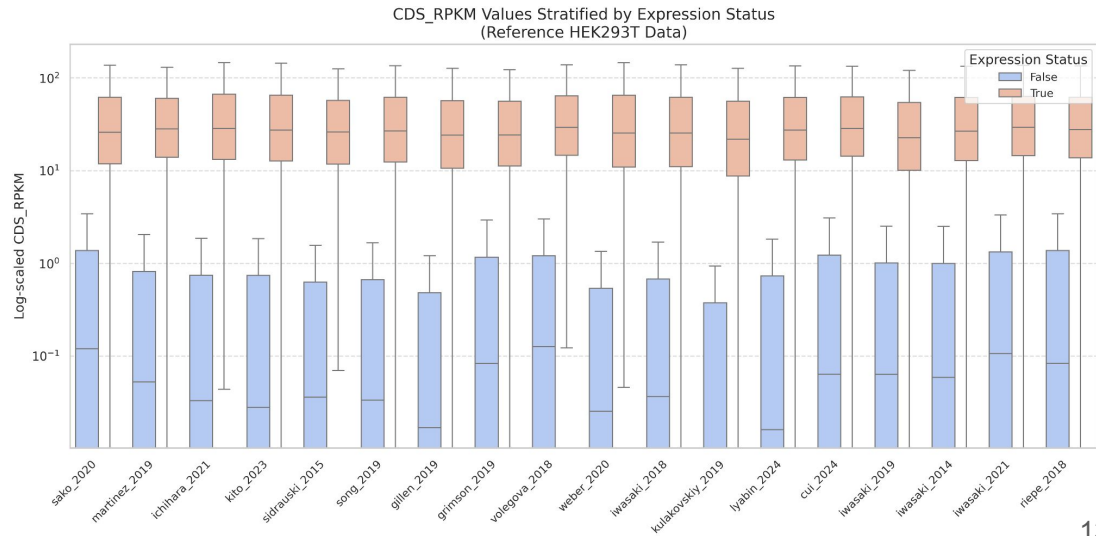


Brar, G., Weissman, J.
Ribosome profiling reveals the what, when, where and how of protein synthesis.
Nat Rev Mol Cell Biol 16, 651–664 (2015)
<https://doi.org/10.1038/nrm4069>

Improve accuracy through highly expressed Genes in HEK cells

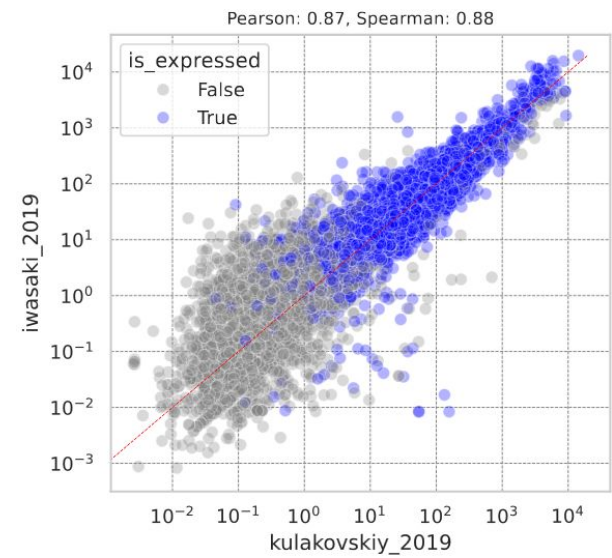


tpm_cutoff = 5
results in 10255 highly expressed genes

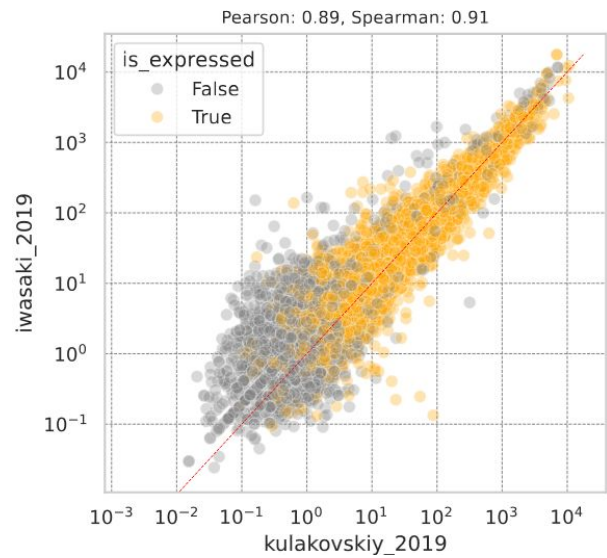


Cut-off leads to higher correlation between datasets (all HEK)

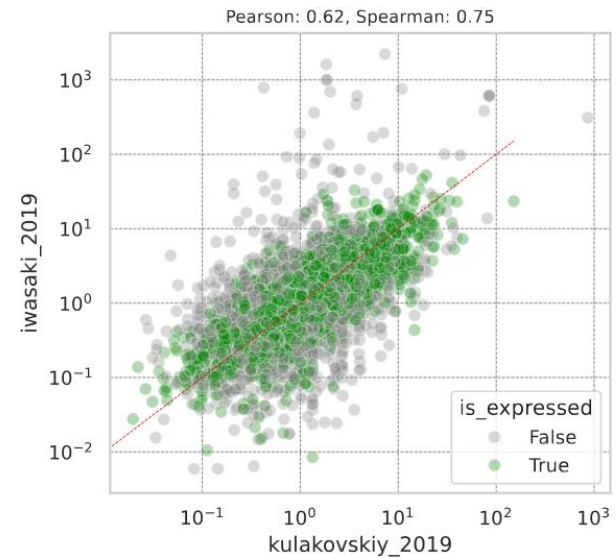
mRNA abundance (rpkm)



ribosome occupancy (rpkm)

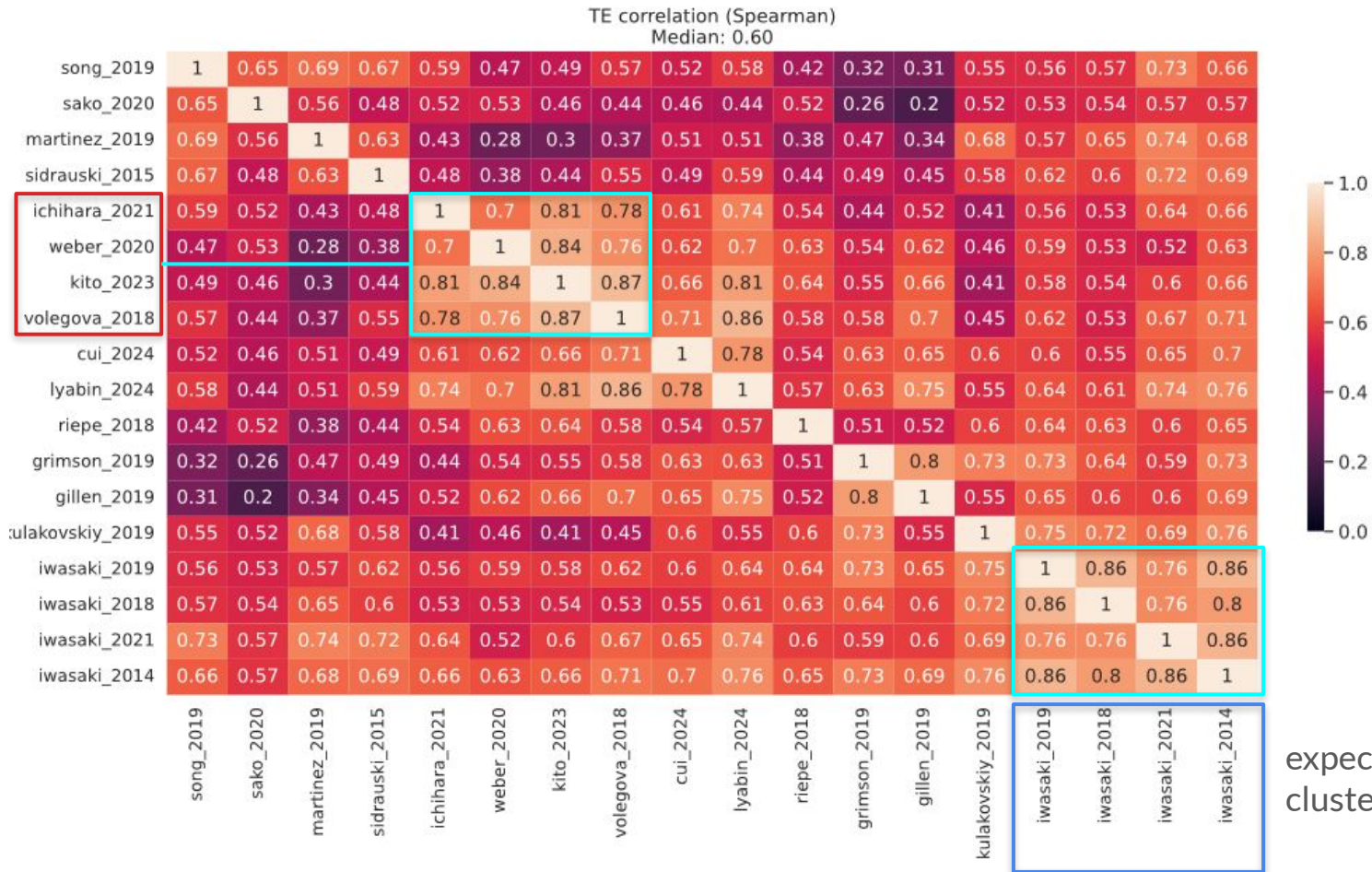


translation efficiency



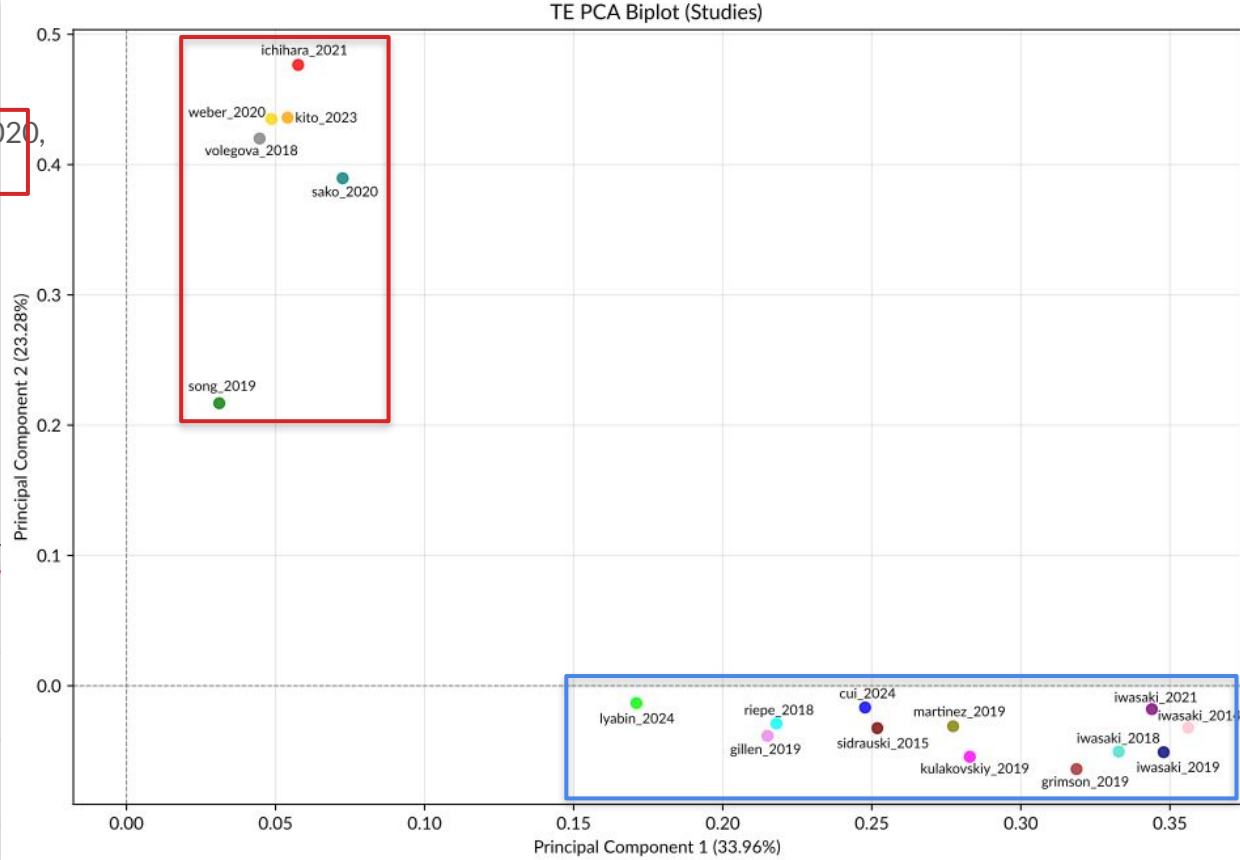
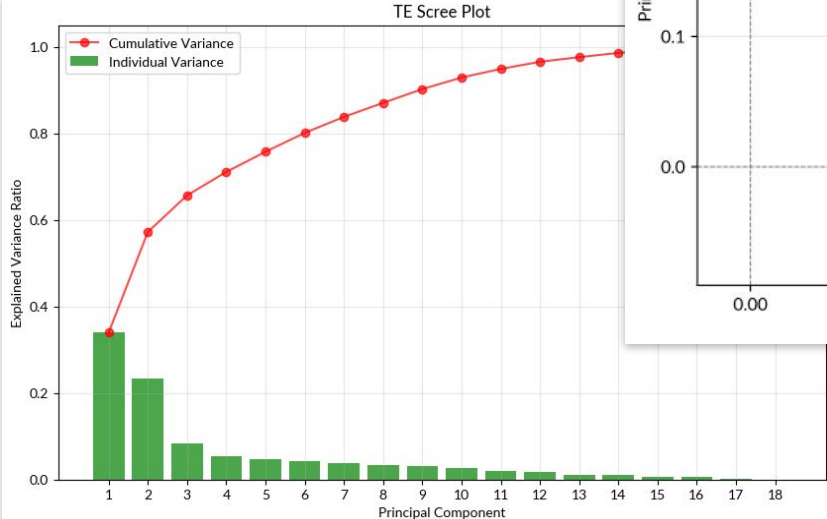
Correlation - Translation Efficiency

unexpected clustering



Looking deeper: PCA - TE

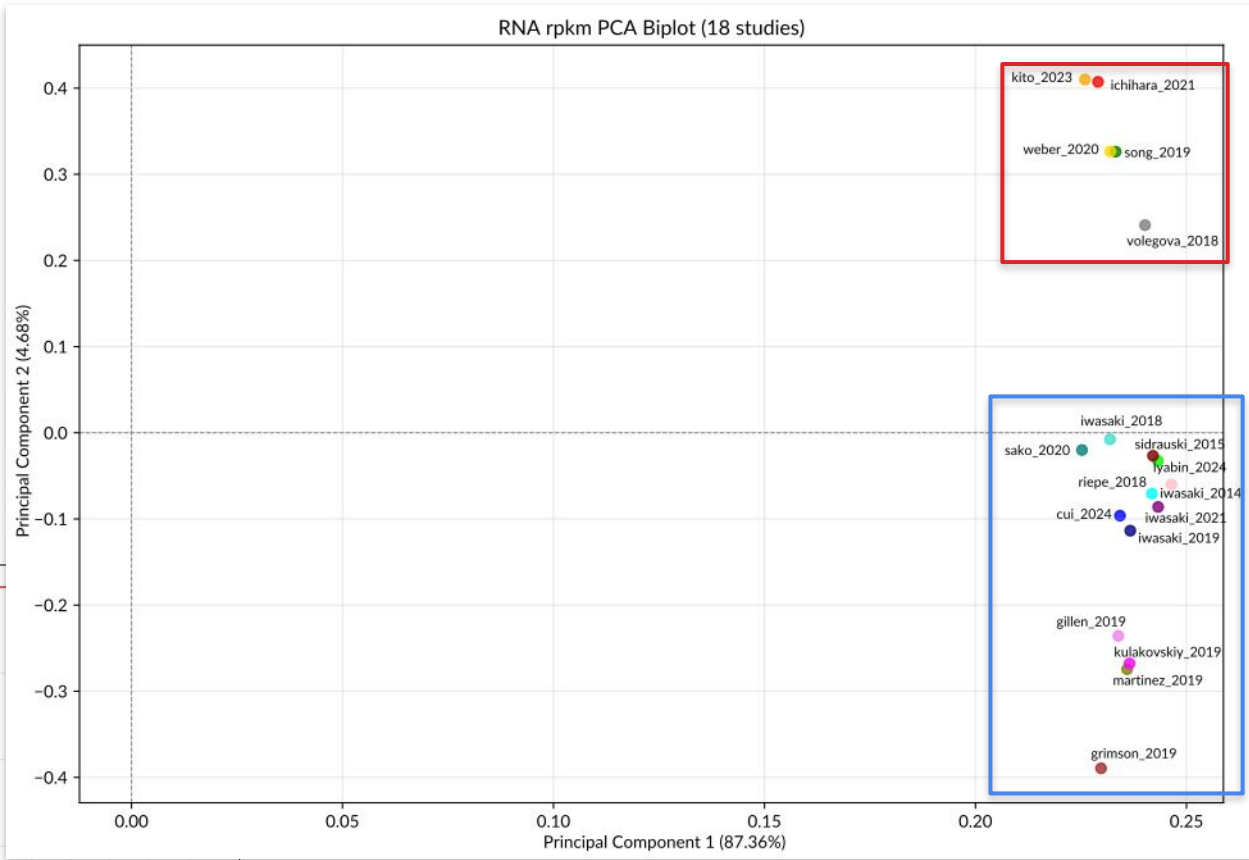
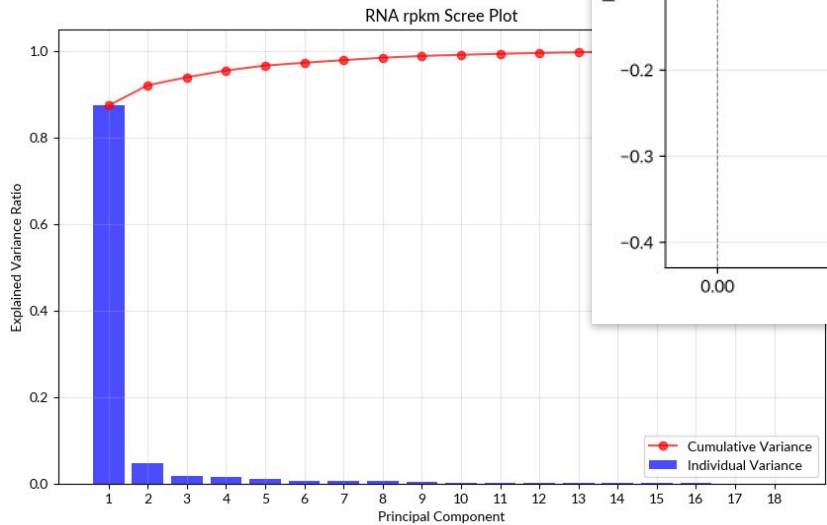
ichihara_2021, kito_2023, vologova_2018, sako_2020,
weber_2020, song_2019 - PC2 cluster



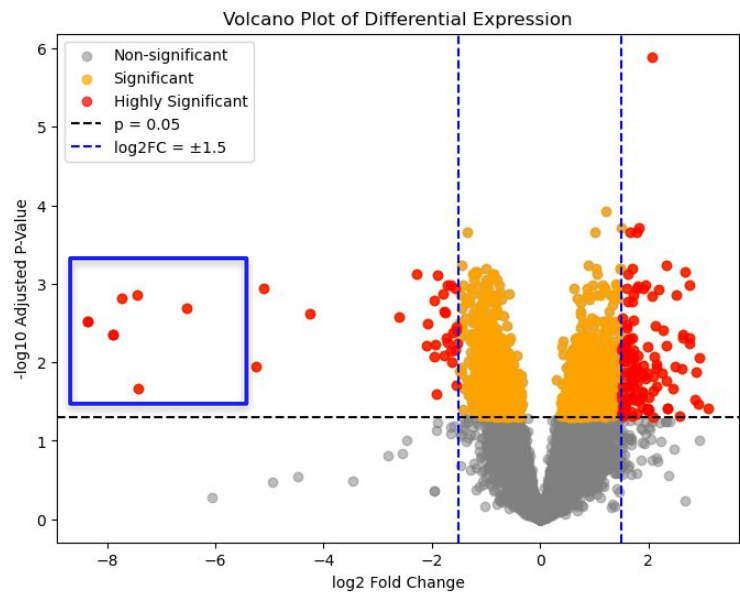
rest - PC1 cluster

PCA - mRNA abundance (rpkm)

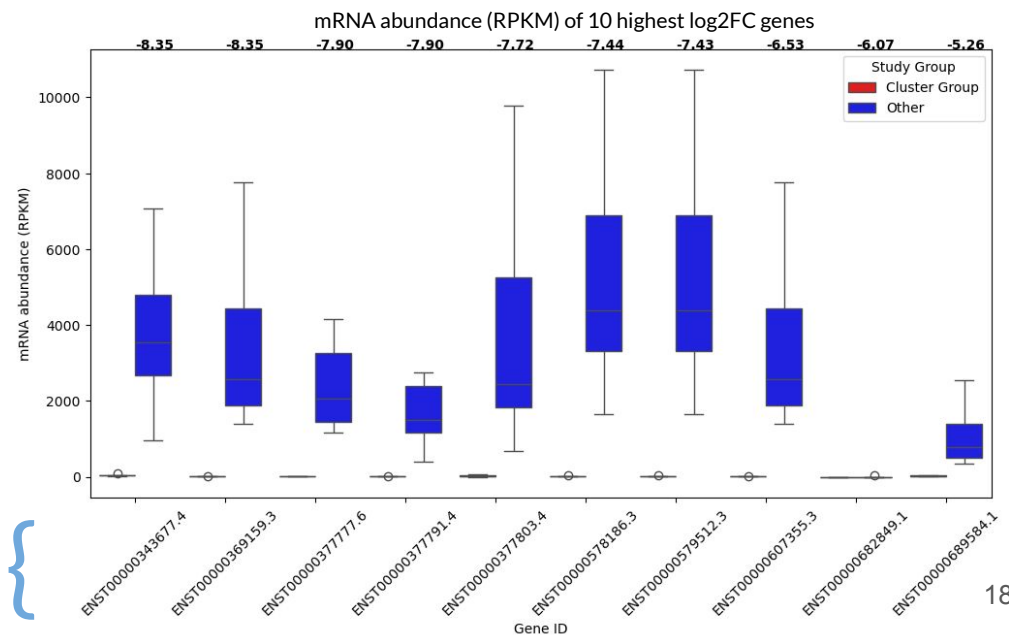
Same clustering as in
Translation Efficiency PCA



PCA - differential expression analysis - do a few genes drive the PC2?

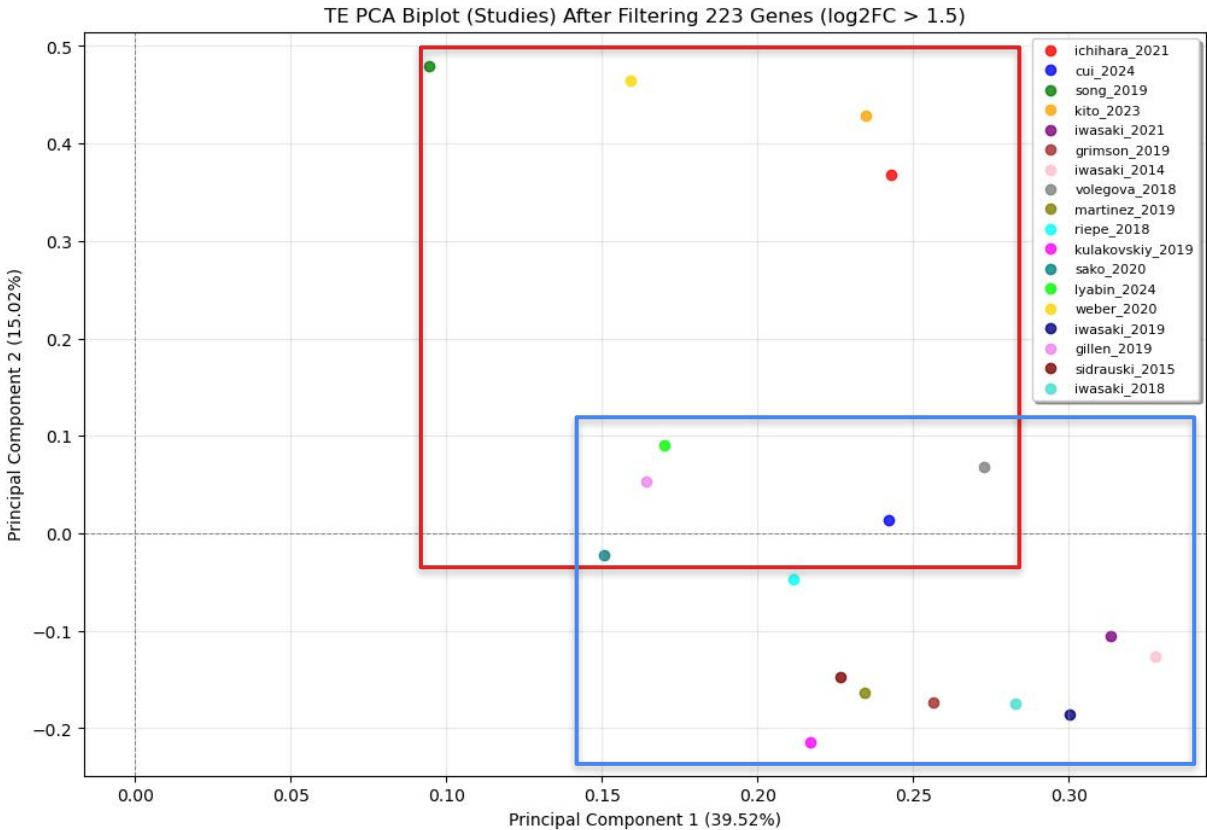


? I could find no explanation in papers ?



histone-encoding genes
(H1, H2, H3, H4 variants)

PCA - removing the culprits



- PC2 cluster dissolves
- PC2 decreased/PC1 increased

Maybe the machine does not care - time will tell!

per transcript:

Nucleotide-Resolution-Level Features

position	A	C	G	U	exon border	5UTR	RNALfold	mRNA abundance	ribosome profile
1	0	1	0	0	0	1	-5.6	0.3	0
2	0	0	1	0	0	1	-5.6	0.8	0
3	1	0	0	0	0	1	-5.6	0.2	0
4	0	0	1	0	0	1	-5.6	1.8	0

targets

Transcript-Level Features

Tx_ID	Ribo_RPKM	RNA_RPKM	GC_Content	A:AU_Ratio	C:CG_Ratio	Exon Junct Density	HalfLife	Translation Efficiency	Protein abundance
ENST00001	16.5	74.3	0.43	0.95	0.87	2	3.4	1.4	1.4
ENST00002	2.8	12.1	0.52	1.11	0.90	4	1.5	0.8	0.8

Going forward...

Ribo-Seq Data QC is challenging

Variability across datasets requires careful Data curation

which features to add/discard for ML

which genes/transcripts to keep for ML

what ML architecture to choose? (nucleotide level vs transcript level)

Thanks to...

Ivo Hofacker (supervisor)

Gabriele Martino (PhD in crime)

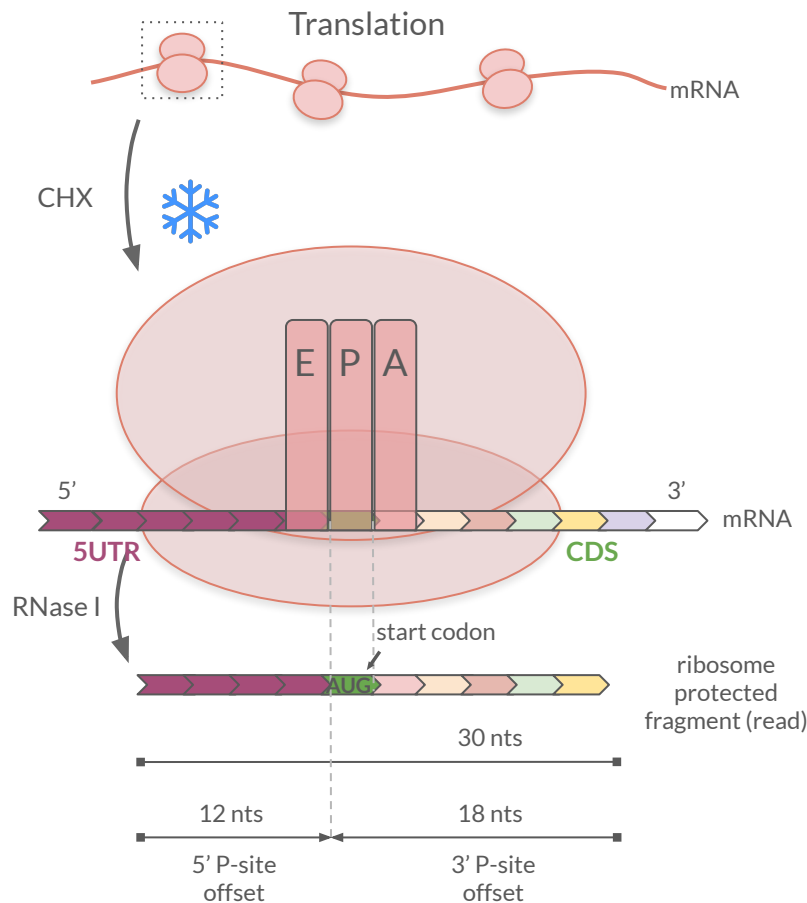
Niko Popitsch (Senior Prophet Scientist at AmeresLab)

TBI (emotional support group)

And you for listening



Problem with variability in Ribo-seq read length



Accurate P-site localization is the foundation of interpreting ribosome profiling data

Variability in rpf lengths! caused by: experimental conditions, ribosome conformations and nuclease biases

→ riboWaltz: most accurate p-site finder!

Aligns reads of the same length and finds with the help of the start codon the ORF and the P-site offsets, which are generalized over the whole rpf length bin.