



CEITEC

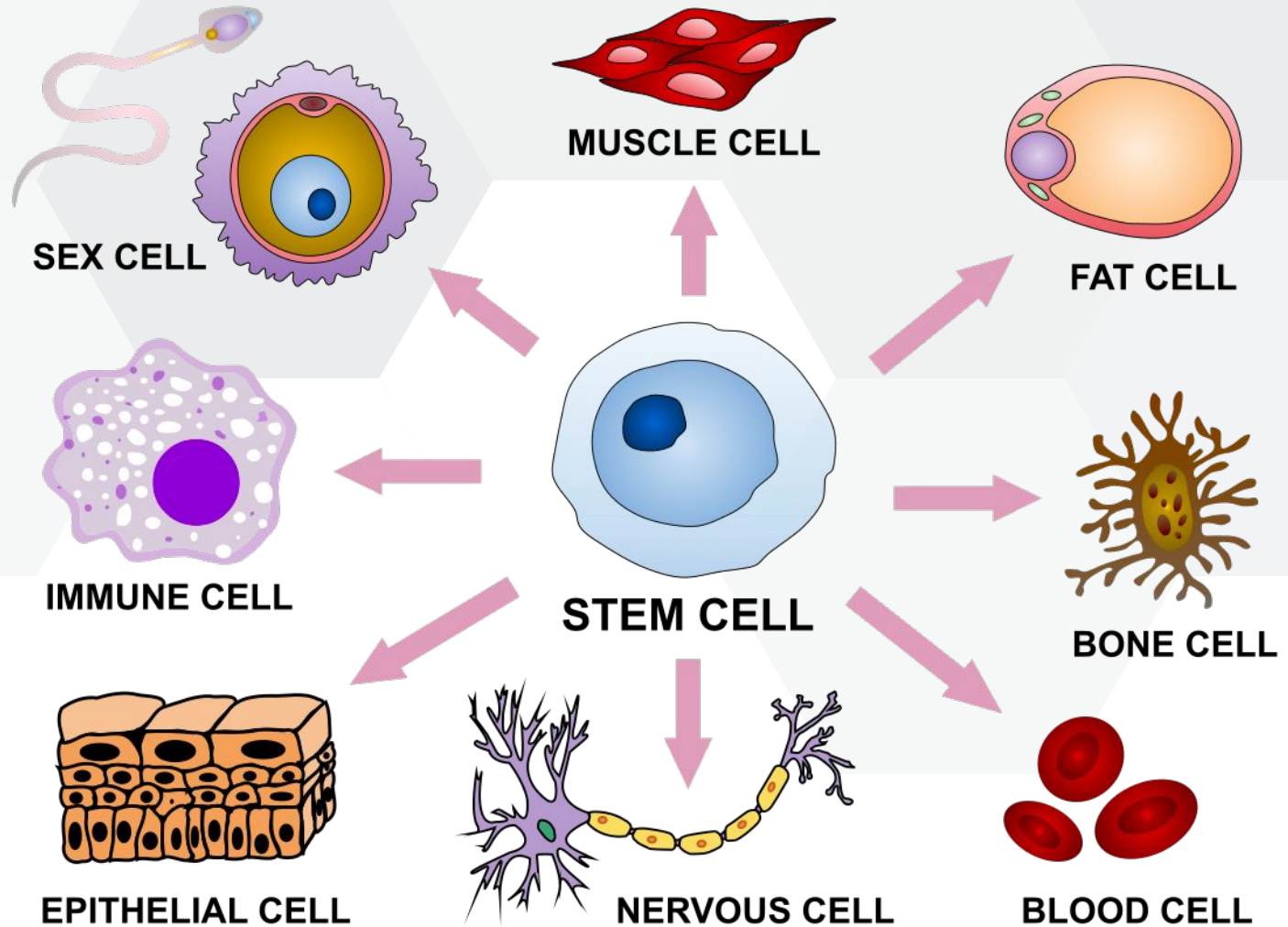
Central European Institute of Technology
BRNO | CZECH REPUBLIC

MUNI

Understanding miRNA binding behaviour through Deep Learning

David Čechák, Katarina Grešová
CEITEC, Masaryk University, Brno, Czech republic

Cells



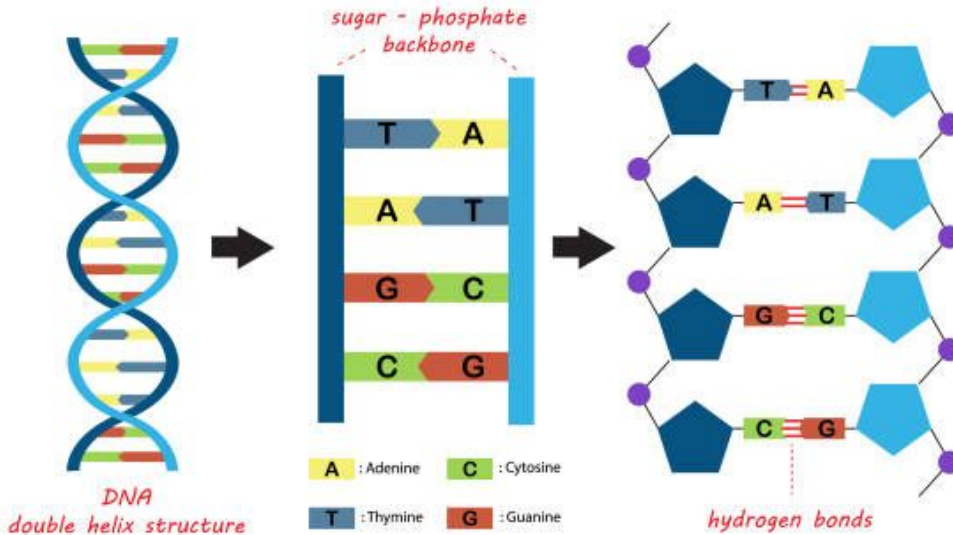




A - T
G - C

BIOLOGY ● ● ●

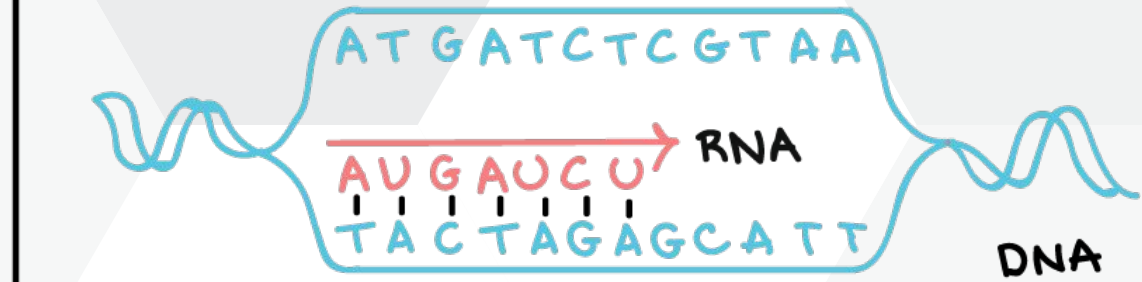
DNA Structure



Transcription



A - T
G - C



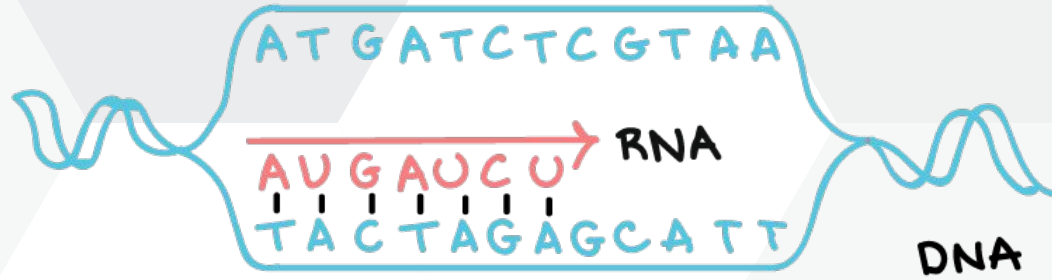
A - U
G - C



Transcription



A - T
G - C



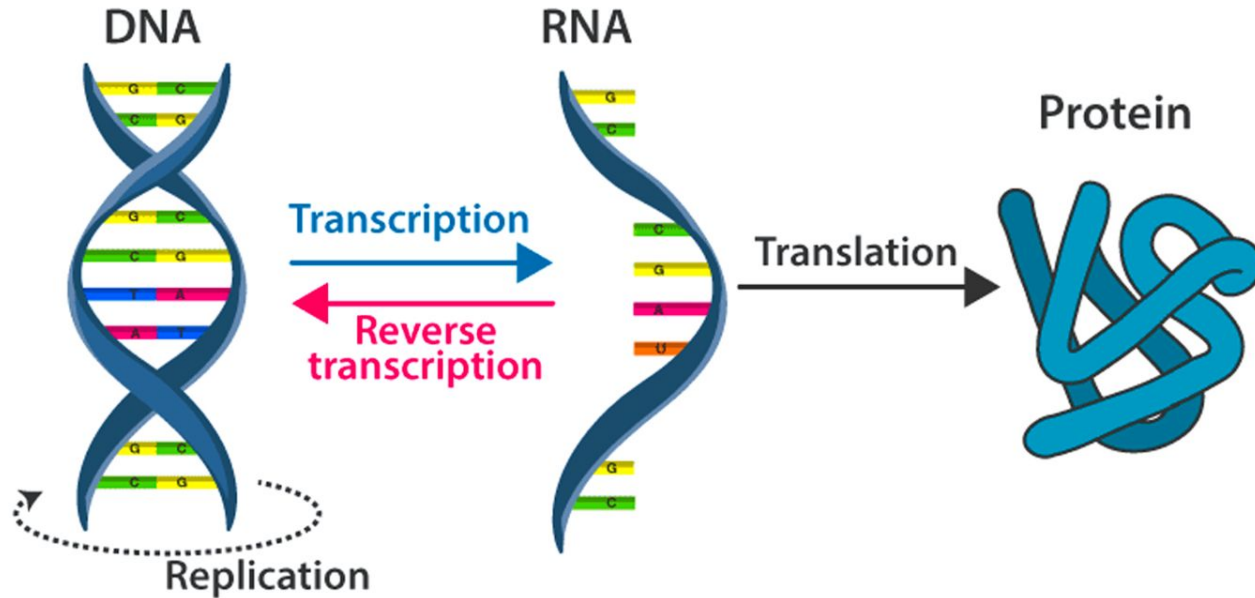
A - U
G - C



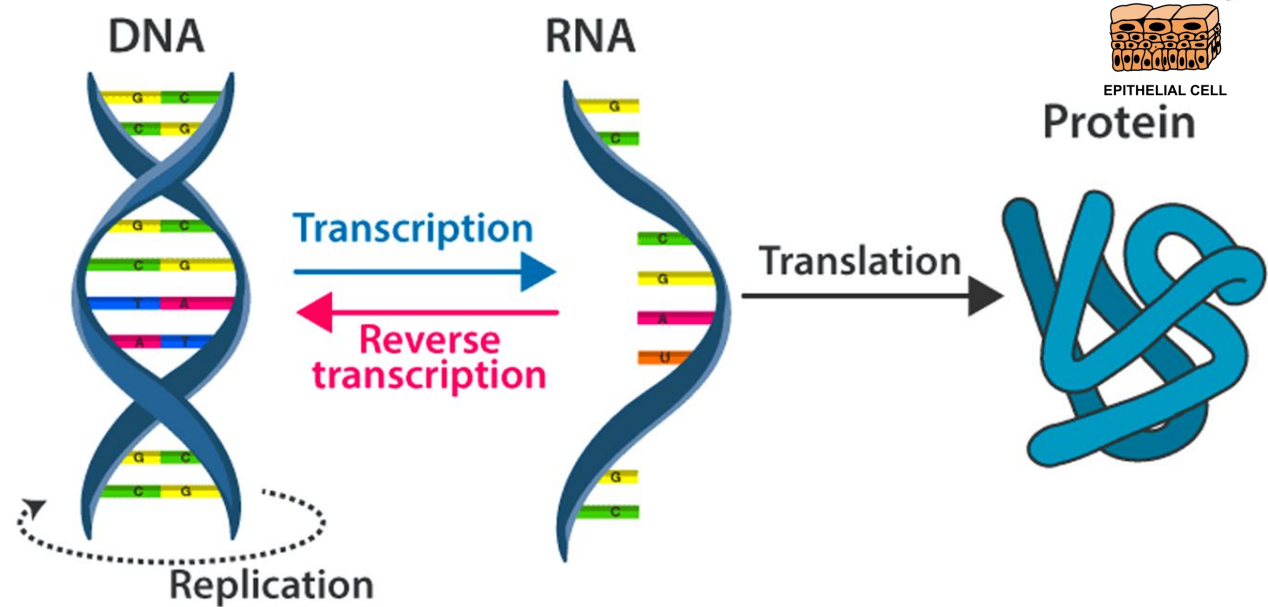
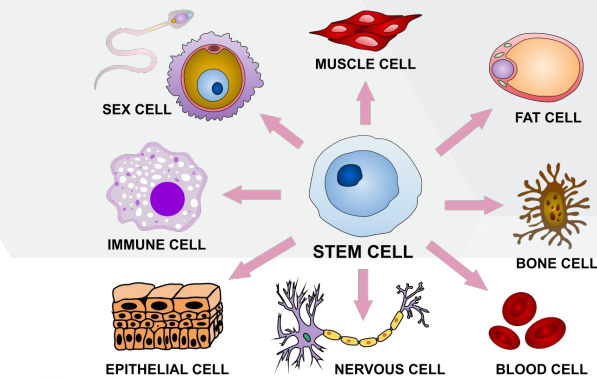
Translation



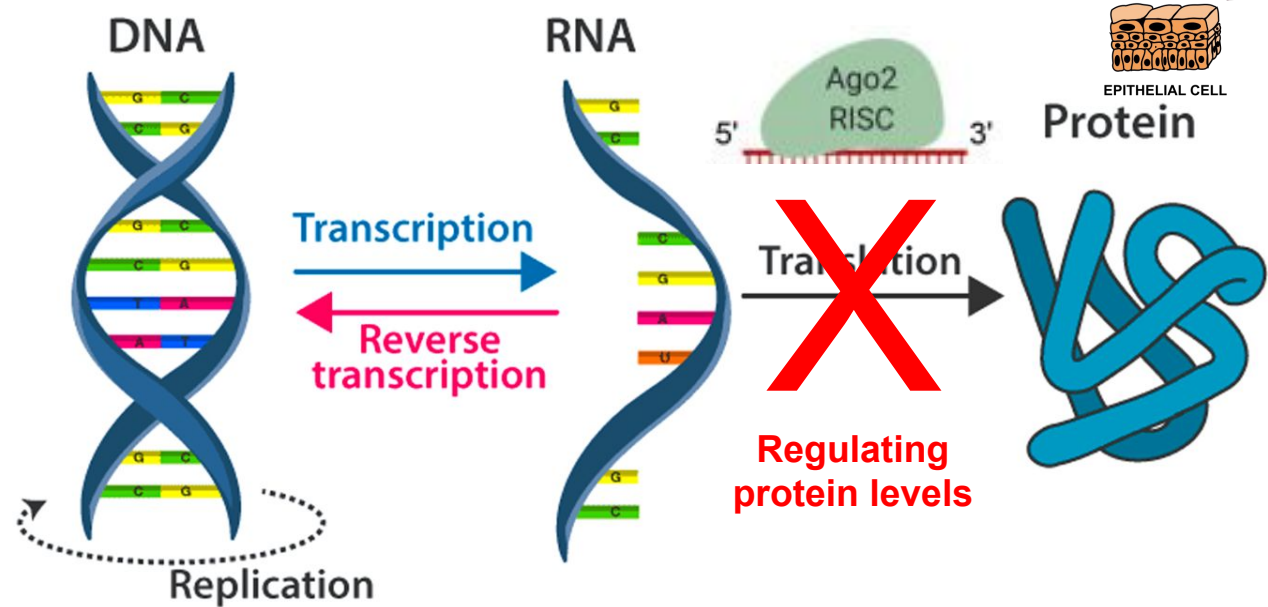
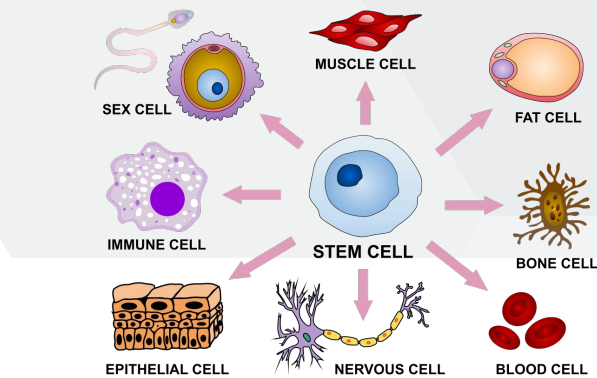
Central Dogma of Molecular Biology



Central Dogma of Molecular Biology



Central Dogma of Molecular Biology



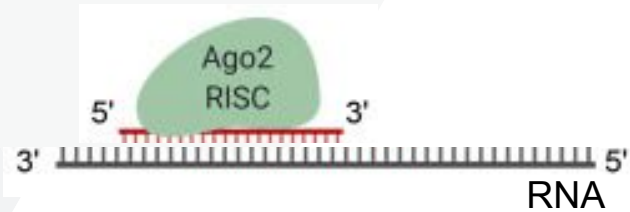
RISC (RNA-induced silencing complex)

CELL



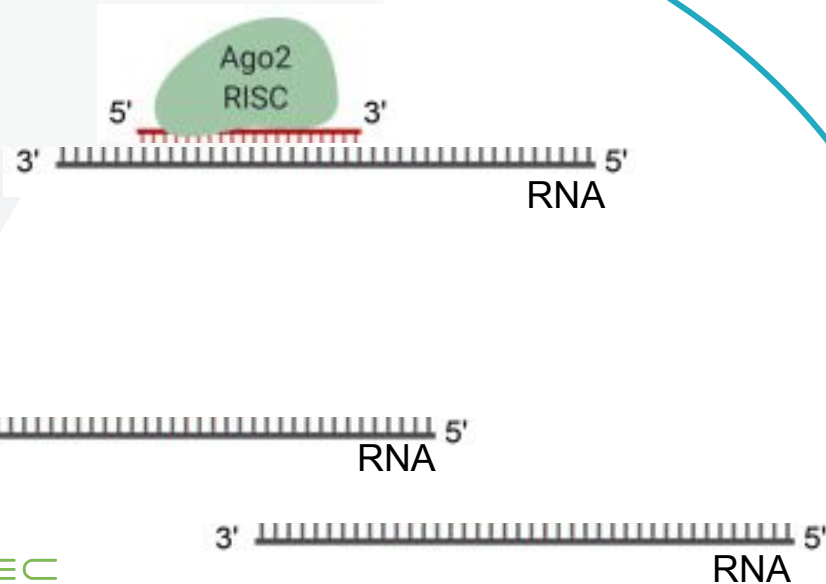
RISC (RNA-induced silencing complex)

CELL



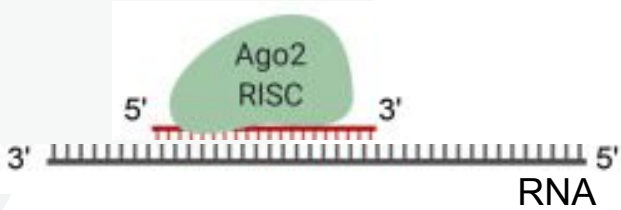
RISC (RNA-induced silencing complex)

CELL



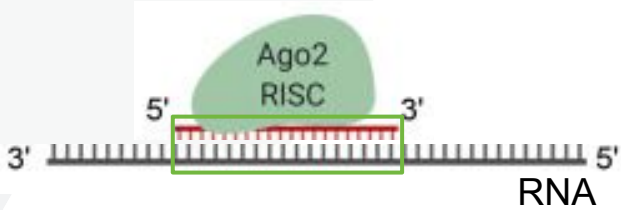
RISC (RNA-induced silencing complex)

CELL



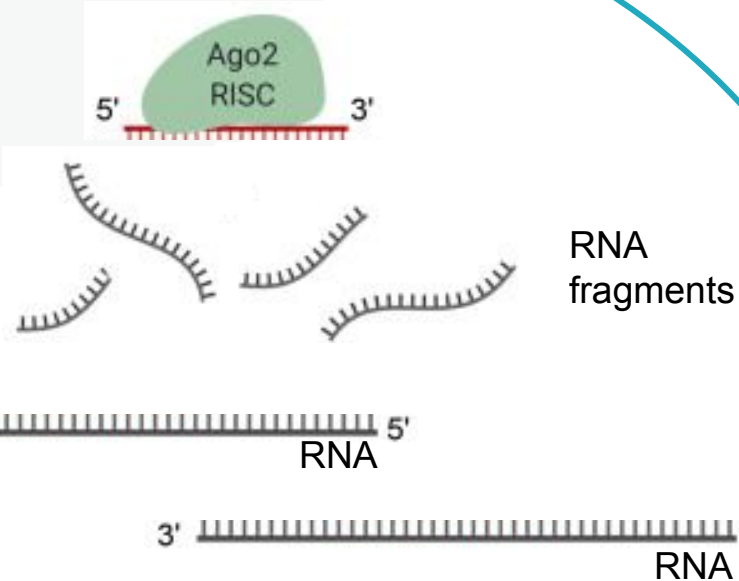
RISC (RNA-induced silencing complex)

CELL



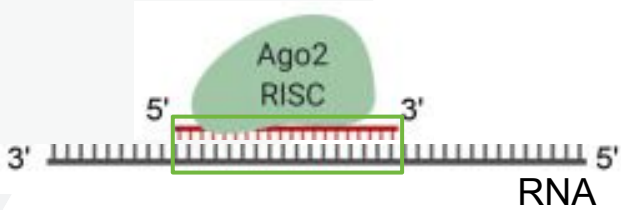
RISC (RNA-induced silencing complex)

CELL

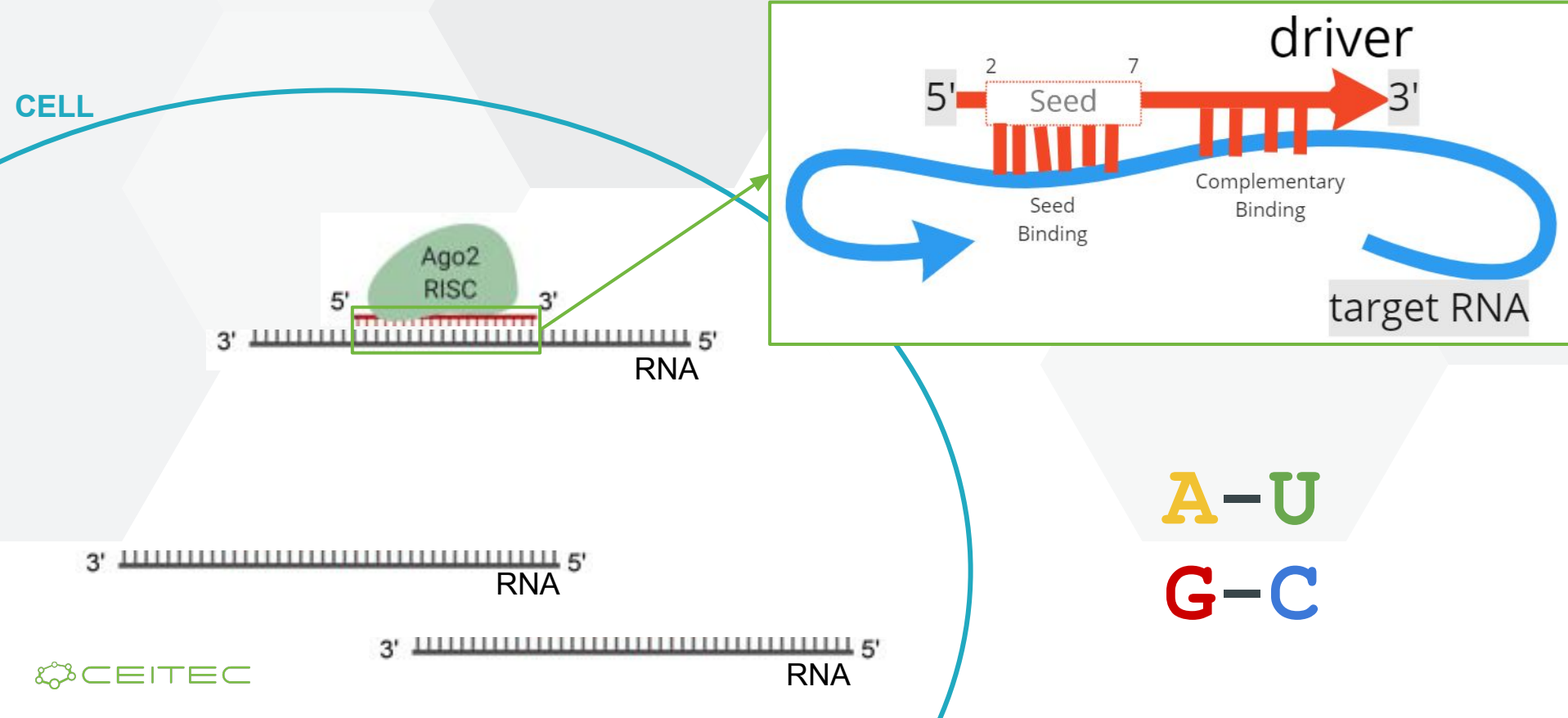


RISC (RNA-induced silencing complex)

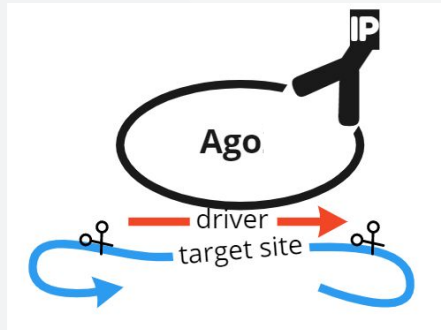
CELL



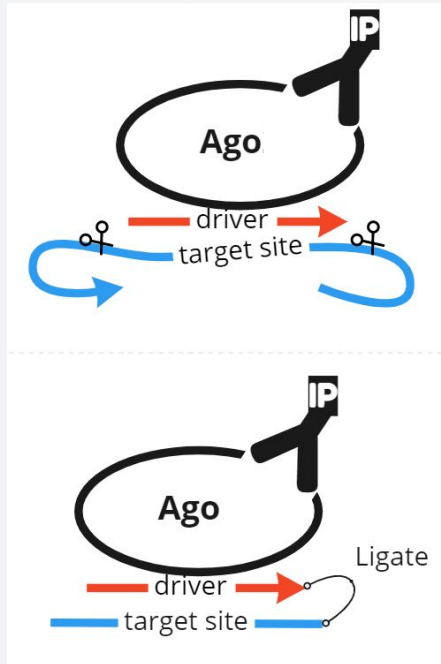
RISC (RNA-induced silencing complex)



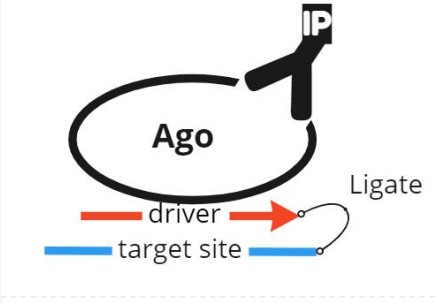
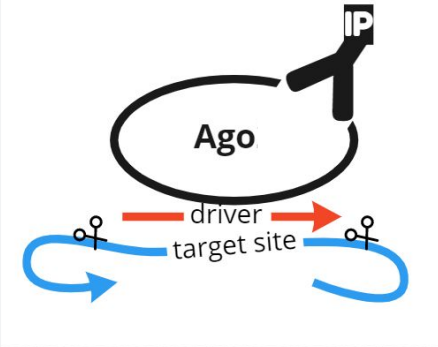
Biological experiment - CLASH



Biological experiment - CLASH



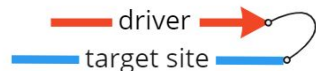
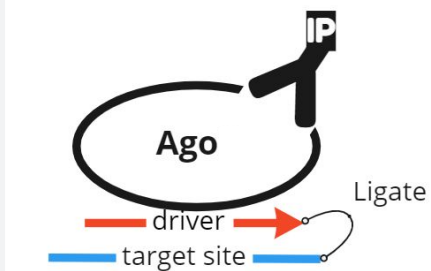
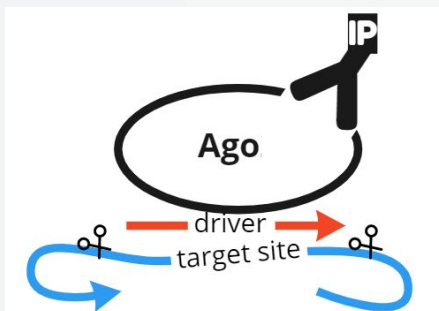
Biological experiment - CLASH



Chimeric
Read



Biological experiment - CLASH - positive and negative samples creation



Chimeric
Read

miRNA	gene	label
AACTGGCCCTCAAAGTCCCG	TGGAGAGCGGGCTTAAGAAGTGGCGGTTTCGGCCGGAGGTTCCATCGTATC	1
ATCAGGGCTTGGAATGGG	CTCGCTGGCGTTCTCCGGGTGGTTGGCATTGTGTCCTGGAAGCGGCCAT	0
TGGGGAGCTGAGGCTCTGGG	CTACACCTCAGCCCGGGGCTGCACTGCCACCCTGGGCAACTTCGCCAAGG	0
GTGAGGGCATGCAGGCCTGG	GTAAGGAGCTGGAGTCGCTGGTAGAGAACGAGGGCAGTGAGGTGCTGGCG	0
ATGCACCTGGGCAAGGATTC	GCATATGGGGCCTTAAGGAATAACAGTGTGCGTGGTGGTGTGCAGGAGA	0
TGCACGGCACTGGGGACACG	TCAGGGTTTCTGGGGGCTTATGAGTCTCACCGGTCAACCCAGGAGGCCT	0
AACTGGCCCTCAAAGTCCCG	ACCTCTTAATGGGCCAGTGAATAACACTCACTGCTGGCATTTAATGTGCA	1
TGGGTTCTGGCATGCTGAT	CACCTGCTGCCCCTTACCCAGCTCCACCACCTGCAGTCCCTAAAGAA	0
TCAGTGCATCACAGAACTTT	ACCCGCACAGCAAGCACCTGTACACGGCCGACATGTTACGCACGGGATC	0
CTGGCCCTCTGCCCCTCC	CTGATTGTGGCAGAGGGGCCACTACCCAAGGTCTAGCTAGGCCCAAGACC	1
TGAGGTAGTAGGTTGTATAG	ATGACCCAACCTACCACCCTGTTTTACATATCCAATTCCAGTAACTCTC	1
TAAAGTGCTTATAGTGCAGG	CAAAGCATACTACCTTCCCCTAGAGGTCTGTAACATTGTGGCTGGGCA	1
TGAGAACTGAATCCATGGG	CCTGGGACCCCCAGGCGTGGAGGACAGTCAAGCCGTGGAGGCCGTGGAGG	0
TGAGGTAGTAGGTTGTATAG	CCCAACCTCAACCTCAACCTCCCAGCACACATCATGCCAGGGGTTGG	1
CTGTACAGGCCACTGCCTTG	GAAGGTAAAGAGGGTCATTGGGGTCGAGCTATGCCAGAGGCTGTGGAGG	0
GTCCCTCTCAAATGTGTCT	GCTGGCCAGCGGACTTCTGGAGTTAGCCTTTGCTTTTGGAGGACTGTGTG	0
TTAGGGCCCTGGCTCCATCT	ACACAGGAAGAGGAGCCAGGCCCTTGTACCTATGGGATTGGACAGGACTG	1
TAGGTAGTTTCATGTTGTTG	TCCGCCCTCTTTTGGCAGCCAGCCCTCCATGCACATTTGGACGCTGTC	0
TAAAGAGCCCTGTGGAGACA	TCCTGAGGCCTGGGGCACCTTTCGTCTGATGAGCCTCTGCATGGAGAGAG	0
GTGGGTACGGCCCAGTGGGG	CATCTTGTCTCACAGCCCAGAGCATGTTCCAGATCCCAGAGTTTGAGCC	0



TACGTCAGTTCATGAAGCT

A

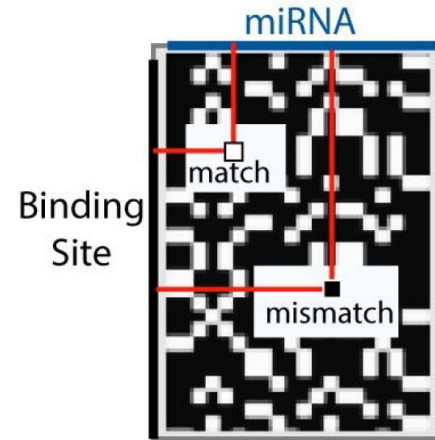
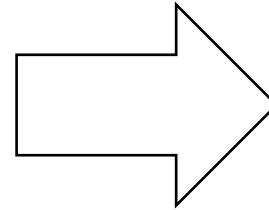
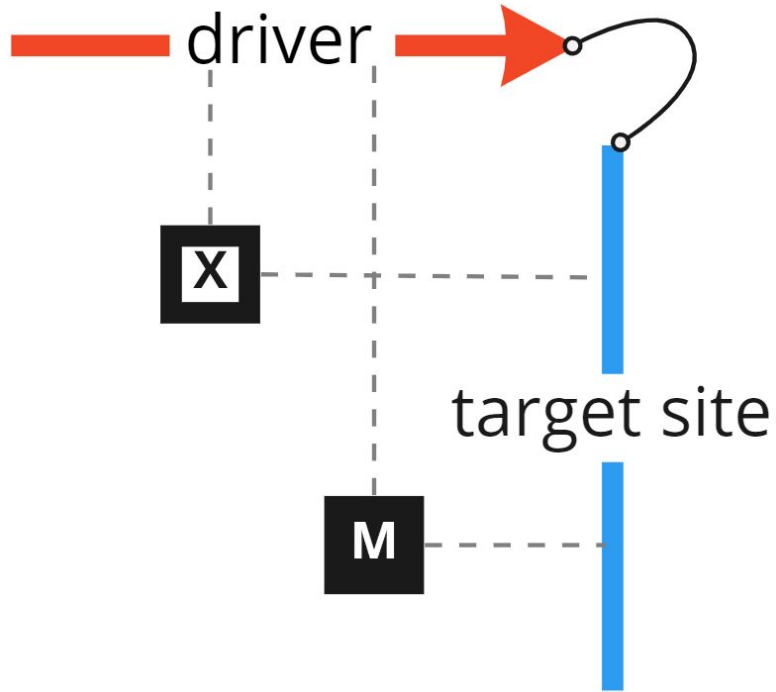
(driver ~20nt)

+

AGTTCTAGTTCGTCCTCGTCAGTGTCAG

TTCATGAGCACCAGTCACGTTTCGTCTA

(target ~50nt)



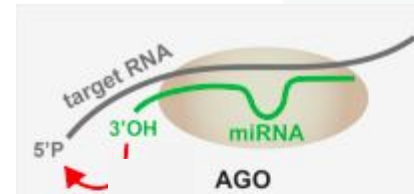
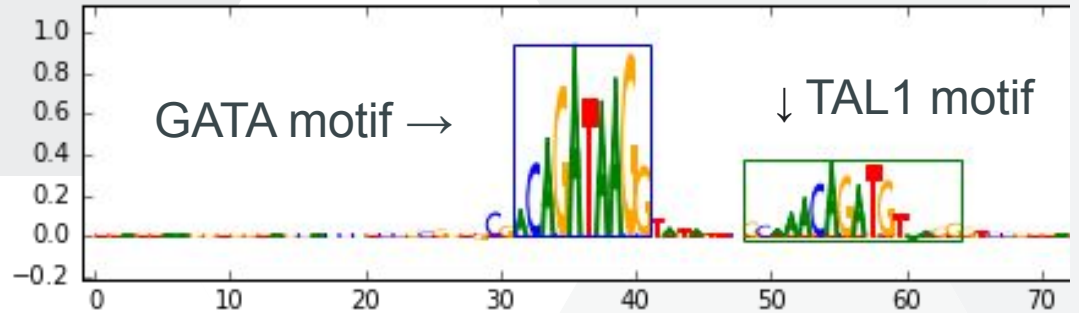
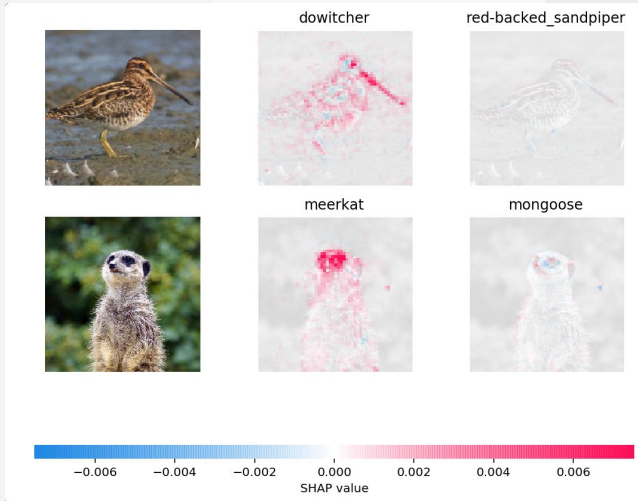
A-T

G-C

Model interpretation - SHAP values



Model interpretation - SHAP values



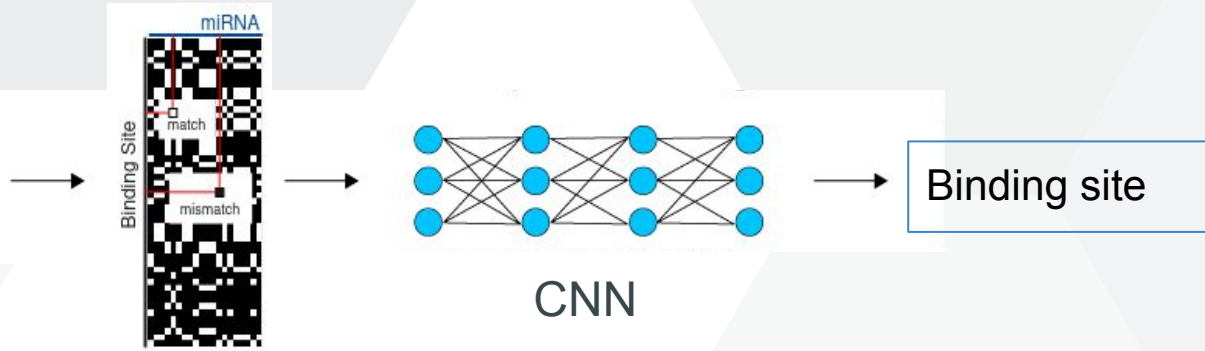
How to visualize interaction between sequences



miRBind model - interpretation

miRNA:
TGAGGTAGTAGGTTGTATAG

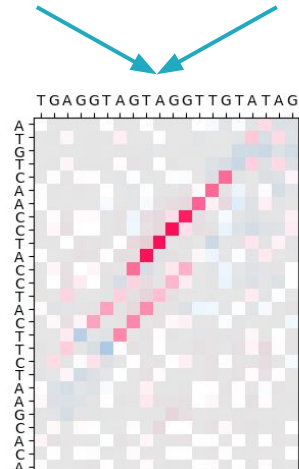
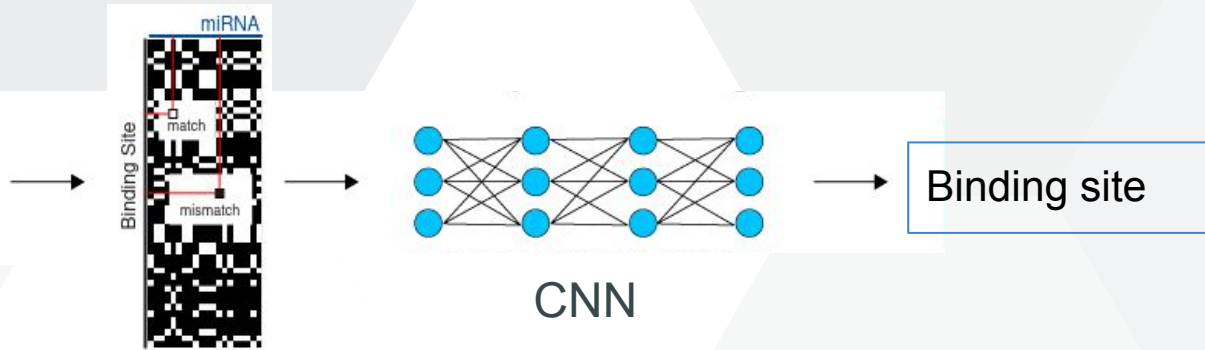
Binding site:
ATGTCAACCTACCTACTTCTAAGCA
CAGGGTATGAAGCTCTCTTTCCACT

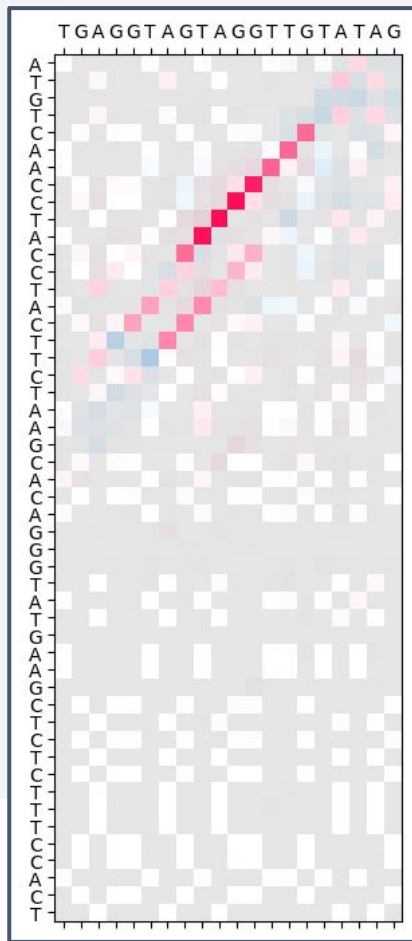


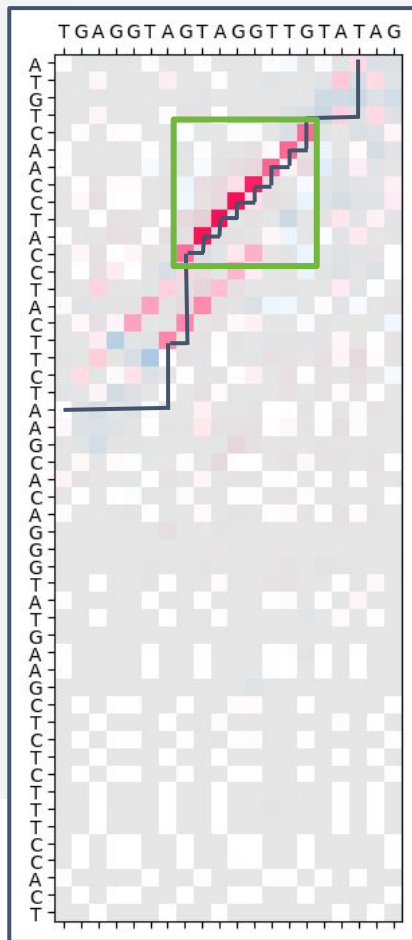
miRBind model - interpretation

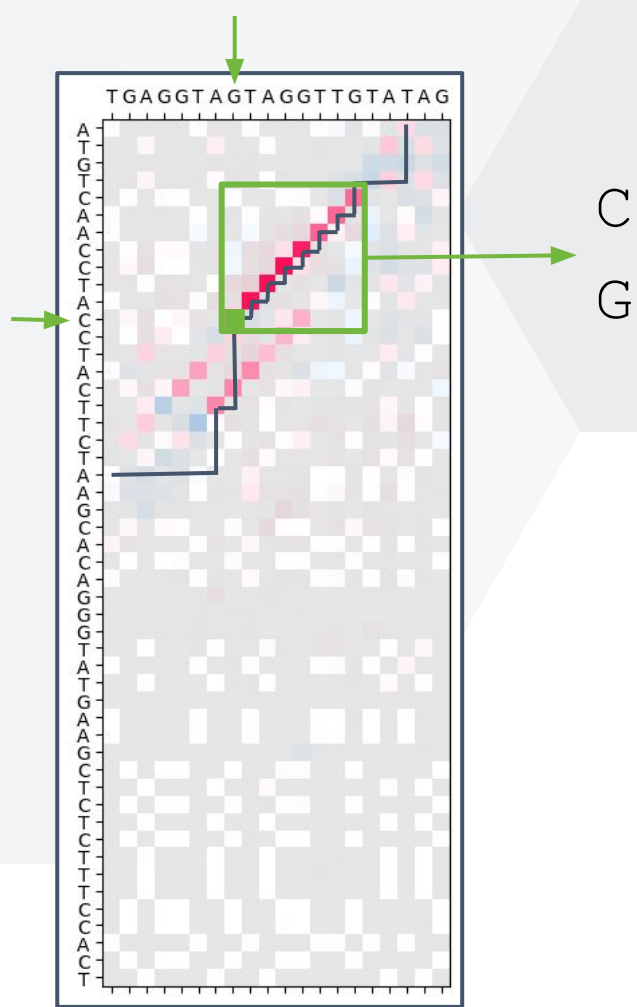
miRNA:
TGAGGTAGTAGGTTGTATAG

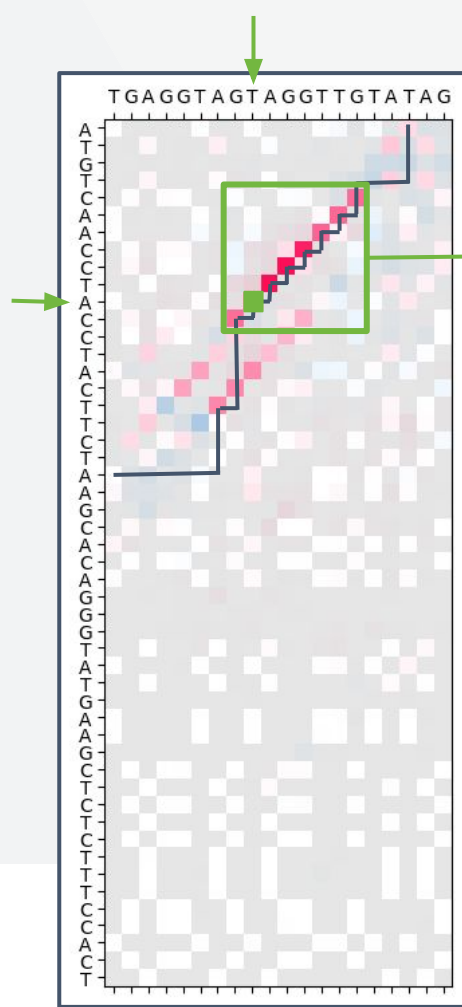
Binding site:
ATGTCAACCTACCTACTTCTAAGCA
CAGGGTATGAAGCTCTCTTTCCACT



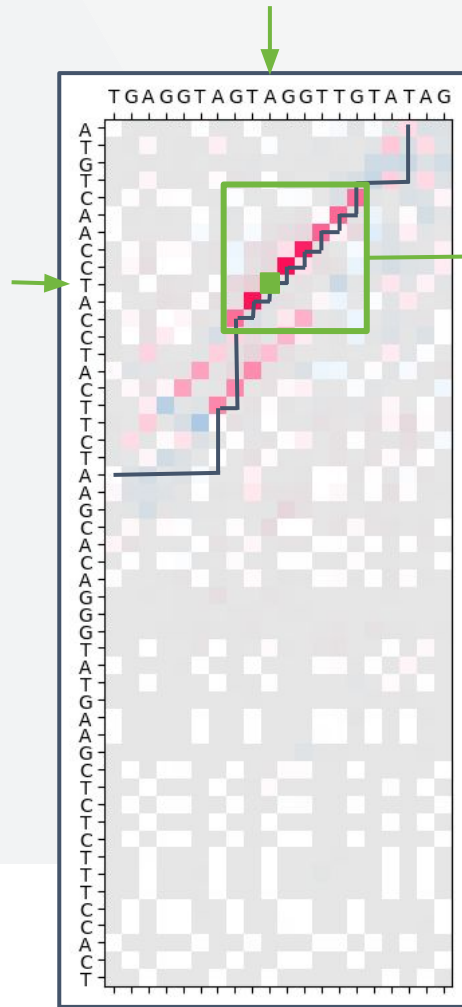


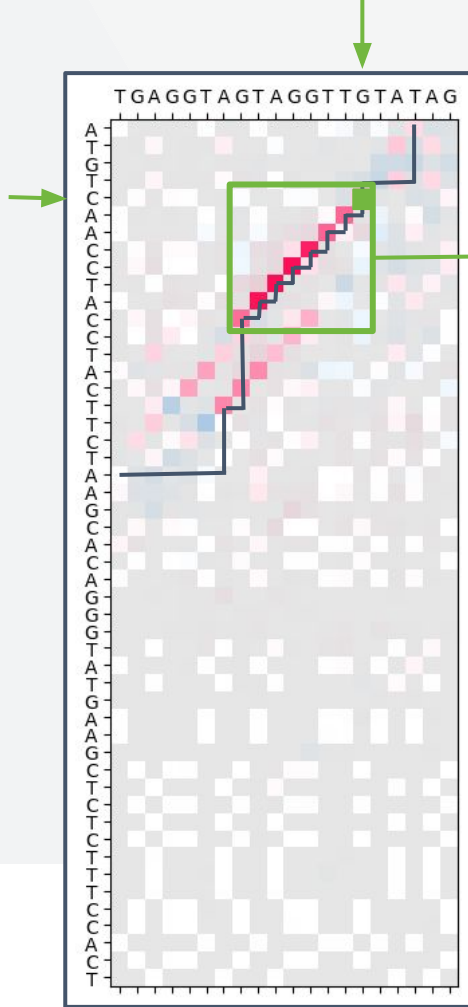




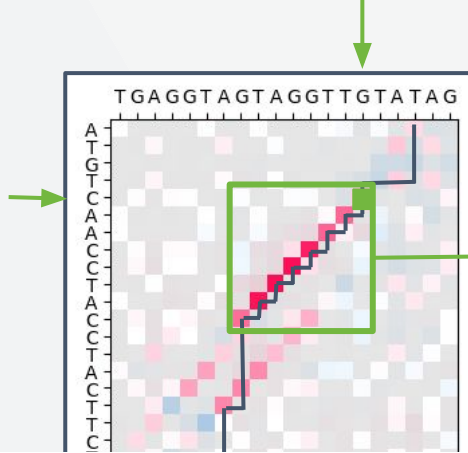


CA
GT



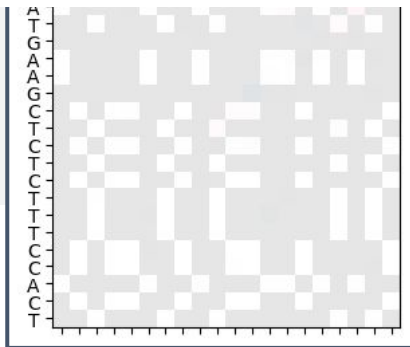


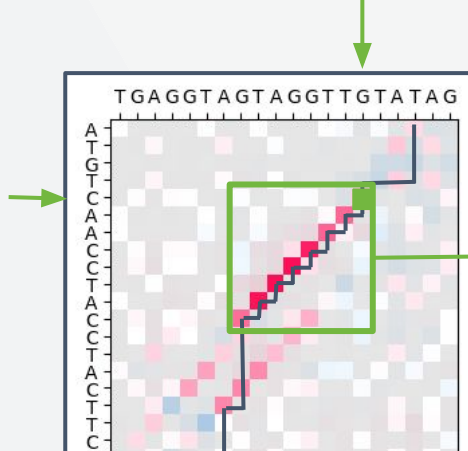
CATCCAAC
 GTAGGTTG



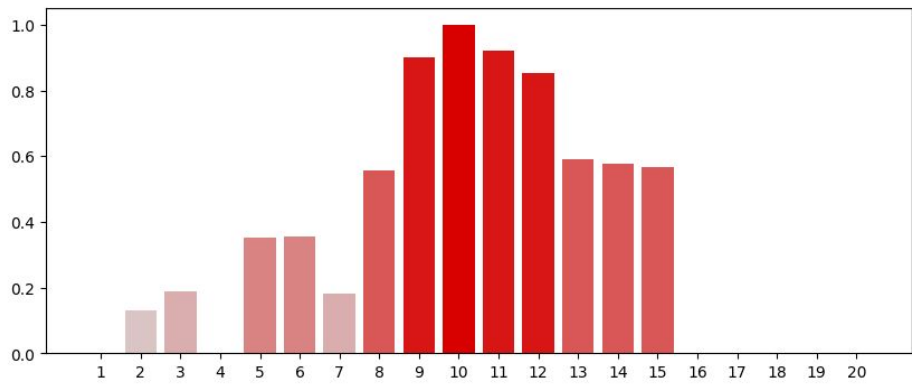
CATCCAAC
GTAGGTTG

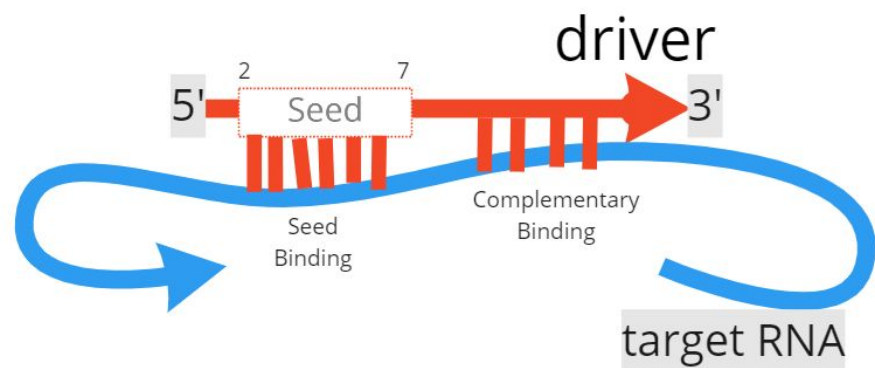
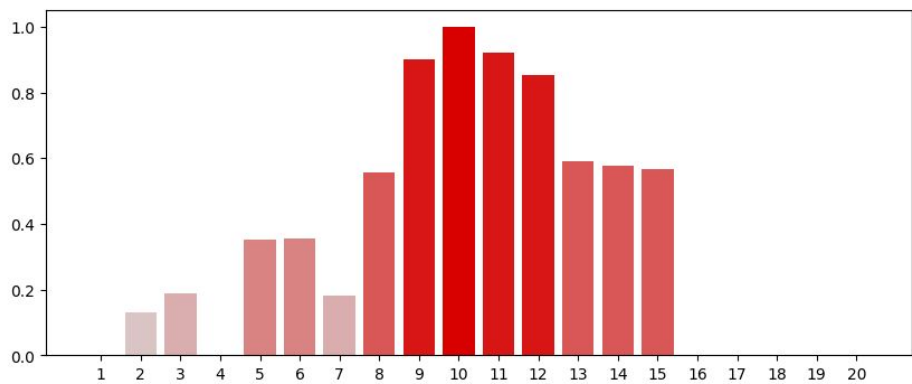
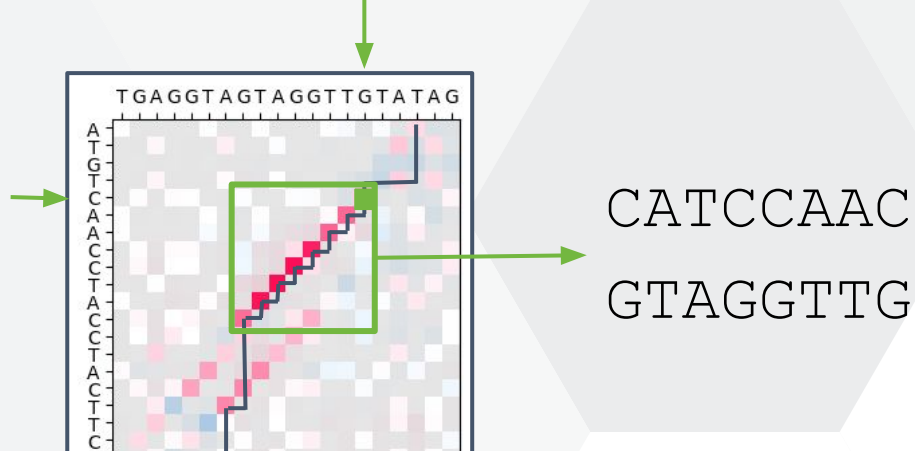
TCACCTTTCTCTCGAAGTATGGGACACGAATCTTCATCCATCCAACGTGA-
 . | | . | | | | | | | | | . . .
 -----TGAGGTA-GTAGGTTGTATAG



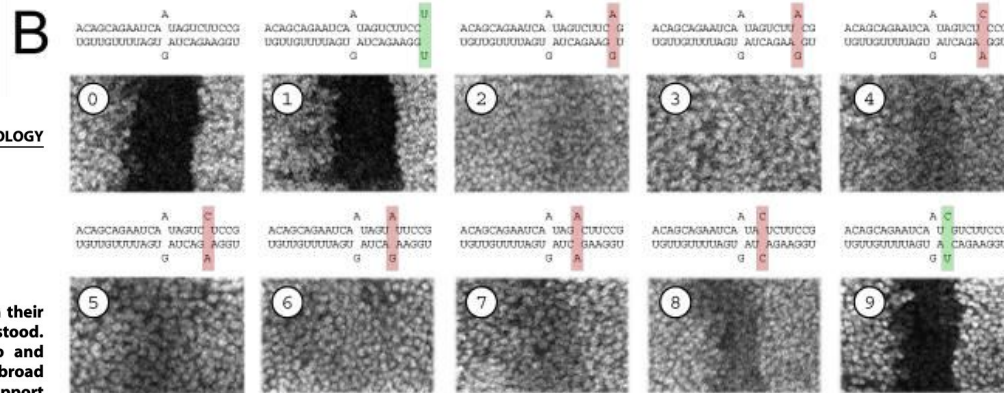
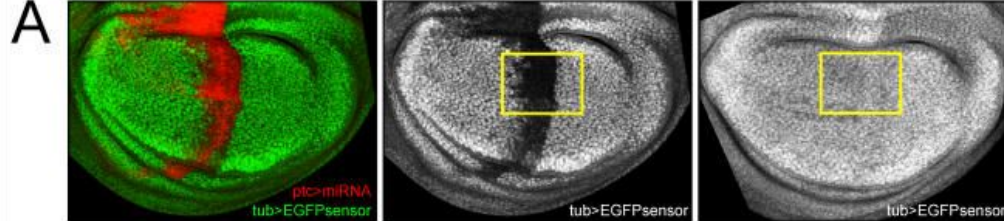


CATCCAAC
GTAGGTTG





Mutagenesis experiment



Open access, freely available online PLOS BIOLOGY

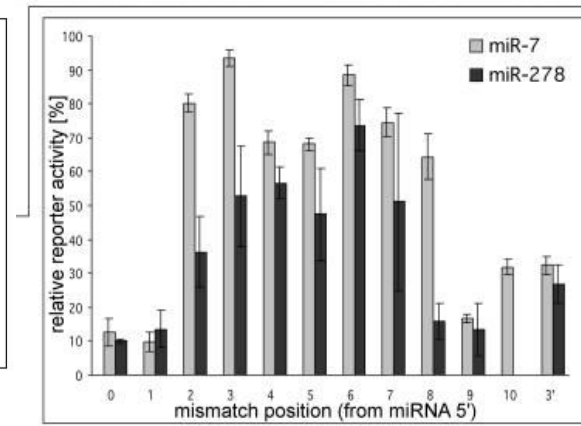
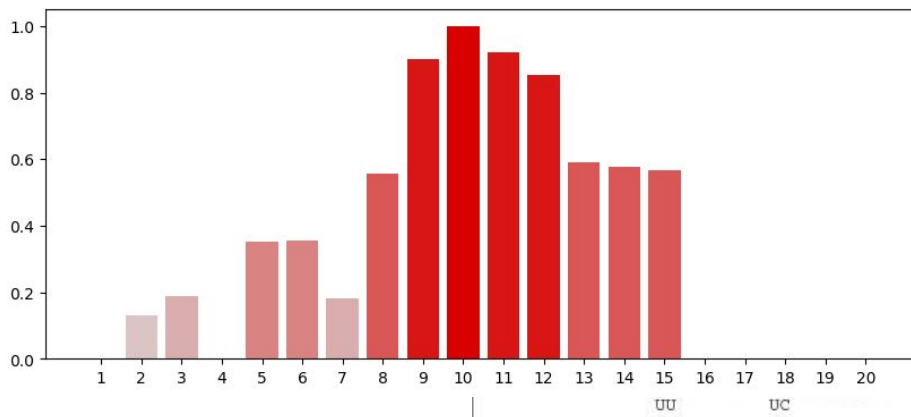
Principles of MicroRNA–Target Recognition

Julius Brennecke¹, Alexander Stark¹, Robert B. Russell, Stephen M. Cohen¹

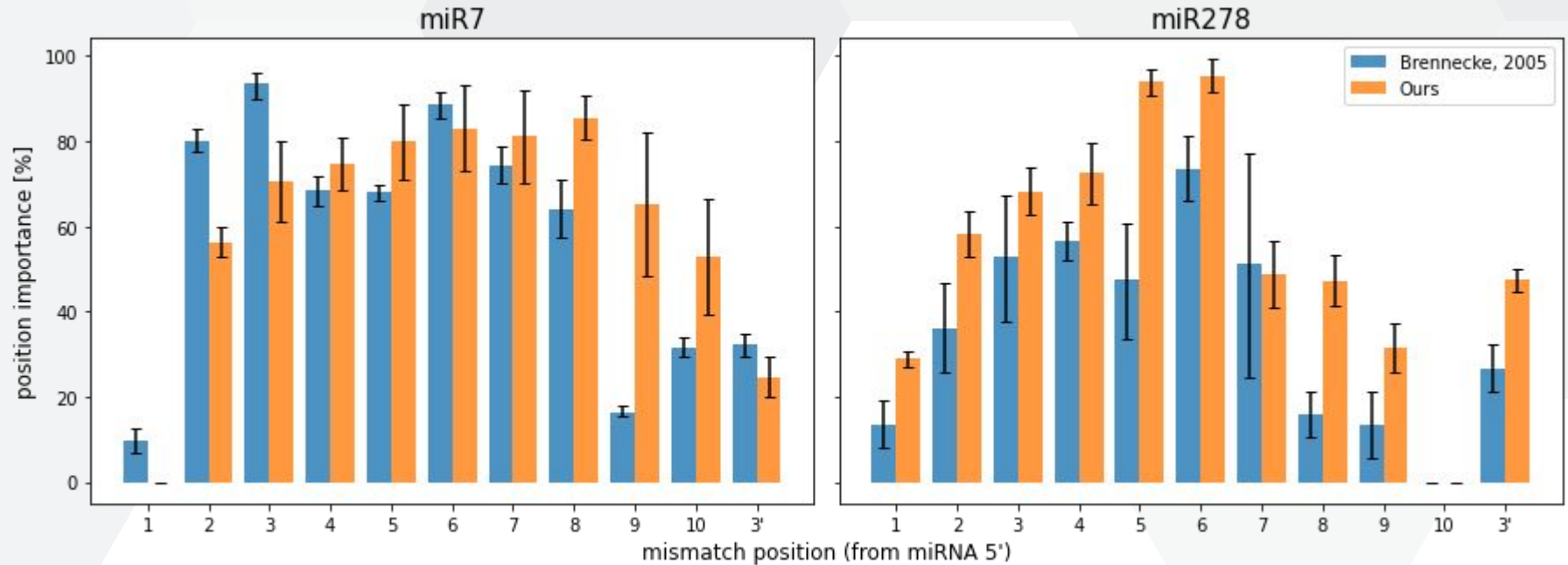
European Molecular Biology Laboratory, Heidelberg, Germany

MicroRNAs (miRNAs) are short non-coding RNAs that regulate gene expression in plants and animals. Although their biological importance has become clear, how they recognize and regulate target genes remains less well understood. Here, we systematically evaluate the minimal requirements for functional miRNA–target duplexes in vivo and distinguish classes of target sites with different functional properties. Target sites can be grouped into two broad categories. 5' dominant sites have sufficient complementarity to the miRNA 5' end to function with little or no support from pairing to the miRNA 3' end. Indeed, sites with strong 5' end. In contrast, 3' compensatory sites have strong 3' end. We present examples and genome-wide statistical analysis of target sites. We provide evidence that an average miRNA regulates a large fraction of protein-coding genes with specificity within miRNA families.

Citation: Brennecke J, Stark A, Russell RB, Cohen SM (2005) Principles of miRNA–Target Recognition. *PLoS Biol* 3(12): e177. doi:10.1371/journal.pbio.0050177



Verification - correlation with in vitro experiment



miRNA	miR-7	miR-278
correlation	0.59	0.85

Functional MicroRNA Targeting

Simple miRNA-mRNA binding model



How will the amount of the products (proteins) of a gene change if a certain miRNA is introduced into the environment in larger quantities?

Task overview

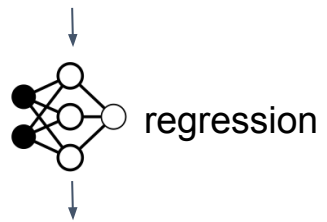


ucagcauagcuacgacguc miRNA, ~20nt long

**Search for binding.
If binds → suppress the mRNA.**

auggacacgcggggcgcgau cgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

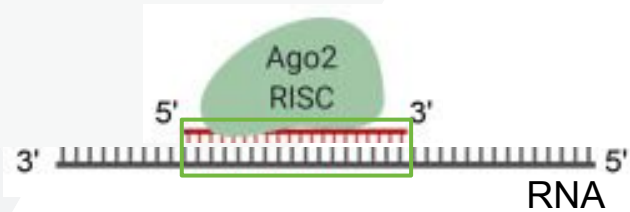
Messenger RNA, 100s – 100,000s nt long



How much less protein products will we get?
(in comparison to a normal cell); Approximate range <0, -2>

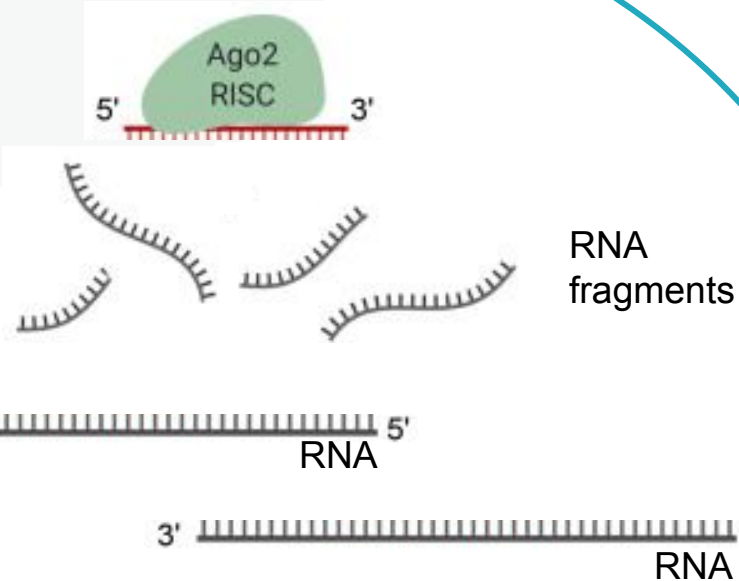
RISC (RNA-induced silencing complex)

CELL

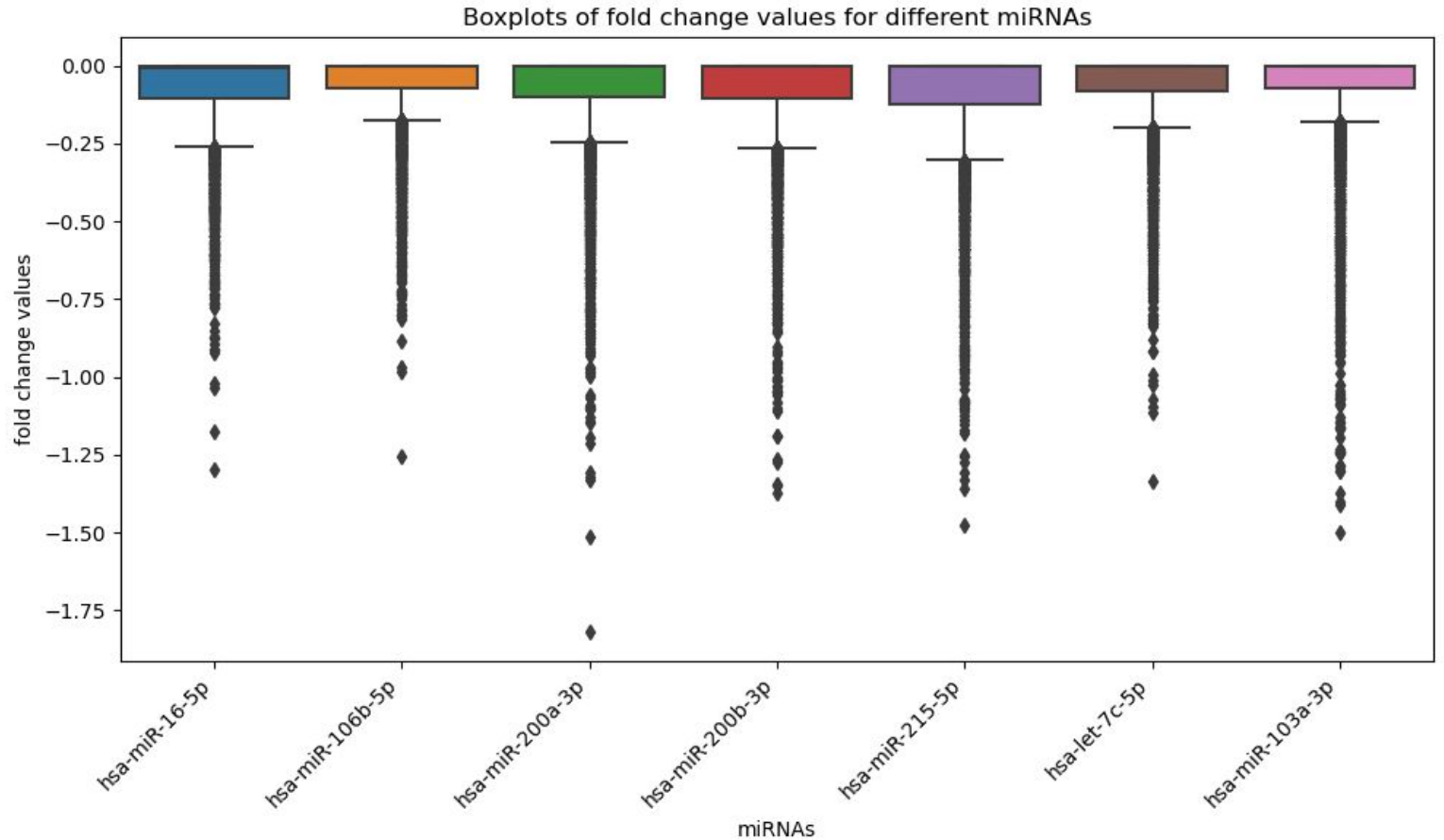


RISC (RNA-induced silencing complex)

CELL



Dataset - labels



Dataset - labels

microRNA	count	mean	std	min	25%	50%	75%	max
hsa-miR-16-5p	7915	-0.072	0.122	-1.297	-0.104	-0.004	0	0
hsa-miR-106b-5p	7902	-0.056	0.109	-1.255	-0.07	0	0	0
hsa-miR-200a-3p	7934	-0.075	0.144	-1.82	-0.098	0	0	0
hsa-miR-200b-3p	7966	-0.077	0.139	-1.372	-0.105	0	0	0
hsa-miR-215-5p	7976	-0.089	0.164	-1.477	-0.122	0	0	0
hsa-let-7c-5p	8002	-0.063	0.119	-1.334	-0.079	0	0	0
hsa-miR-103a-3p	7489	-0.069	0.152	-1.498	-0.072	0	0	0
average	7883.429	-0.07158	0.135495	-1.43614	-0.09286	-0.00057	0	0

f.e. hsa-let-7c-5p

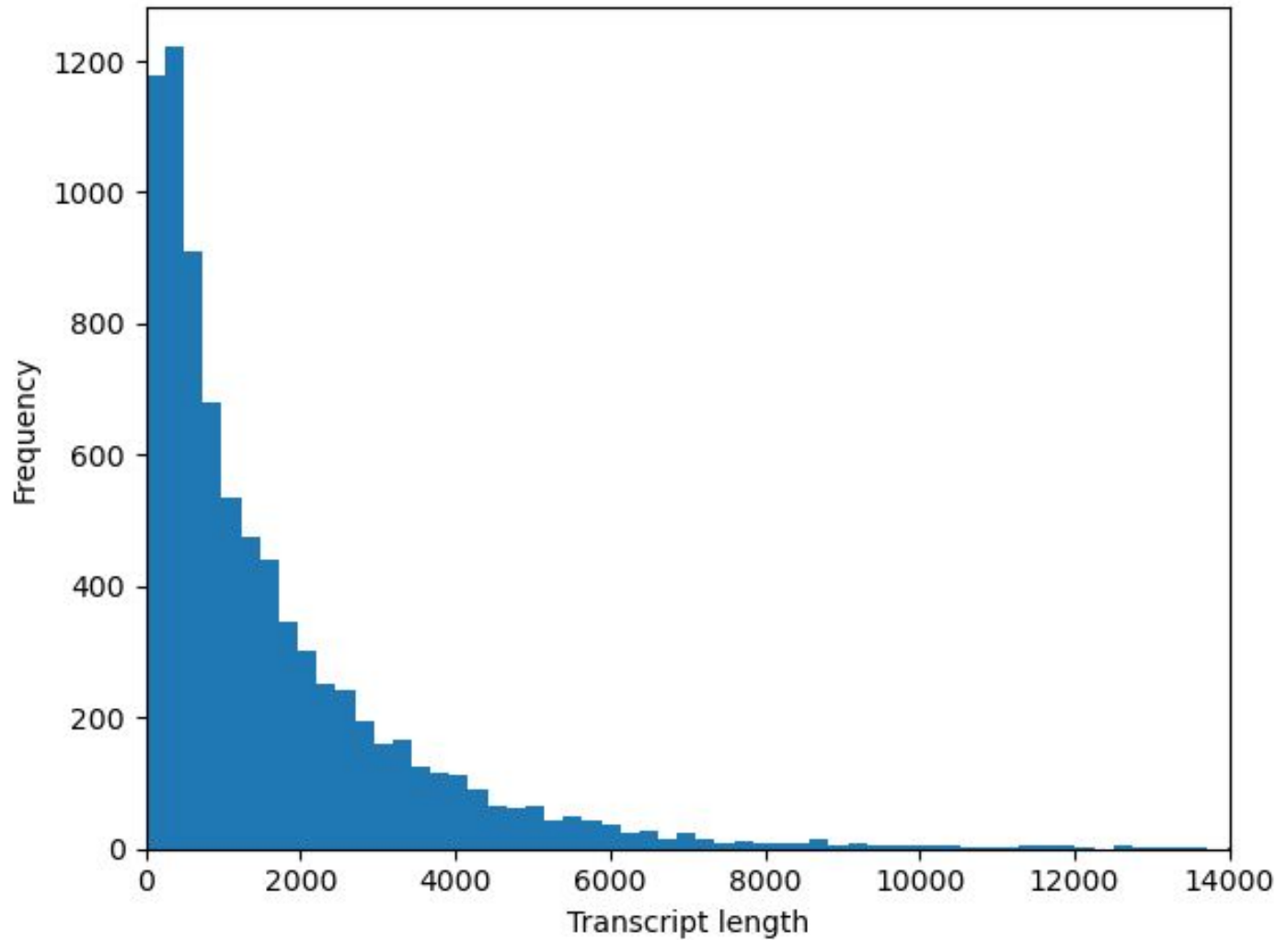
train:4046

test:2050

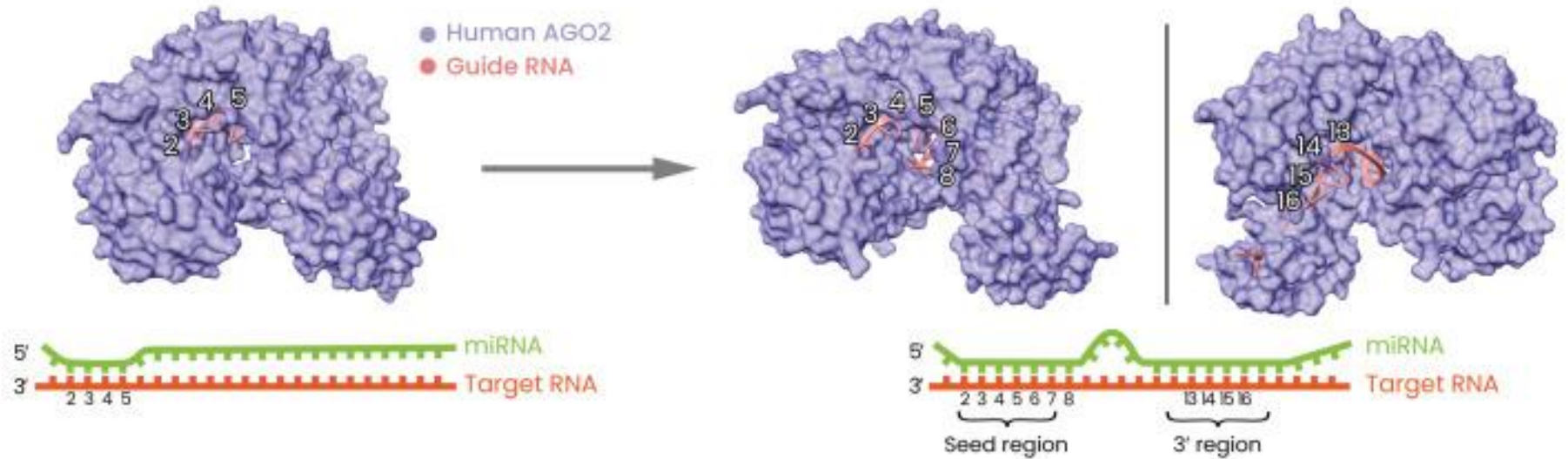
-- because removing transcripts without signal

Dataset - inputs

Lengths of
transcript
sequences



State-of-the-art so far - manual feature extraction



Canonical sites			1	2	3	4	5	6	7	8	9
	8-mer	A	O	O	O	O	O	O	O	O	N
7-mer m8	B	O	O	O	O	O	O	O	O	N	
7-mer A1	A	O	O	O	O	O	O	Ø	N		
6-mer	B	O	O	O	O	O	Ø	N			

Noncanonical sites			1	2	3	4	5	6	7	8	9
	6-mer A1	A	O	O	O	O	Ø	Ø	N		
offset 7-mer	B	Ø	O	O	O	O	O	O	N		
offset 6-mer	B	Ø	O	O	O	O	O	N			
CDNST 1	N	N	Ø	O	O	O	O	B			
CDNST 2	N	N	O	O	Ø	O	O	A			
CDNST 3	O	O	O	Ø	O	Ø	Ø	N			
CDNST 4	N	O	Ø	Ø	Ø	O	O	A			

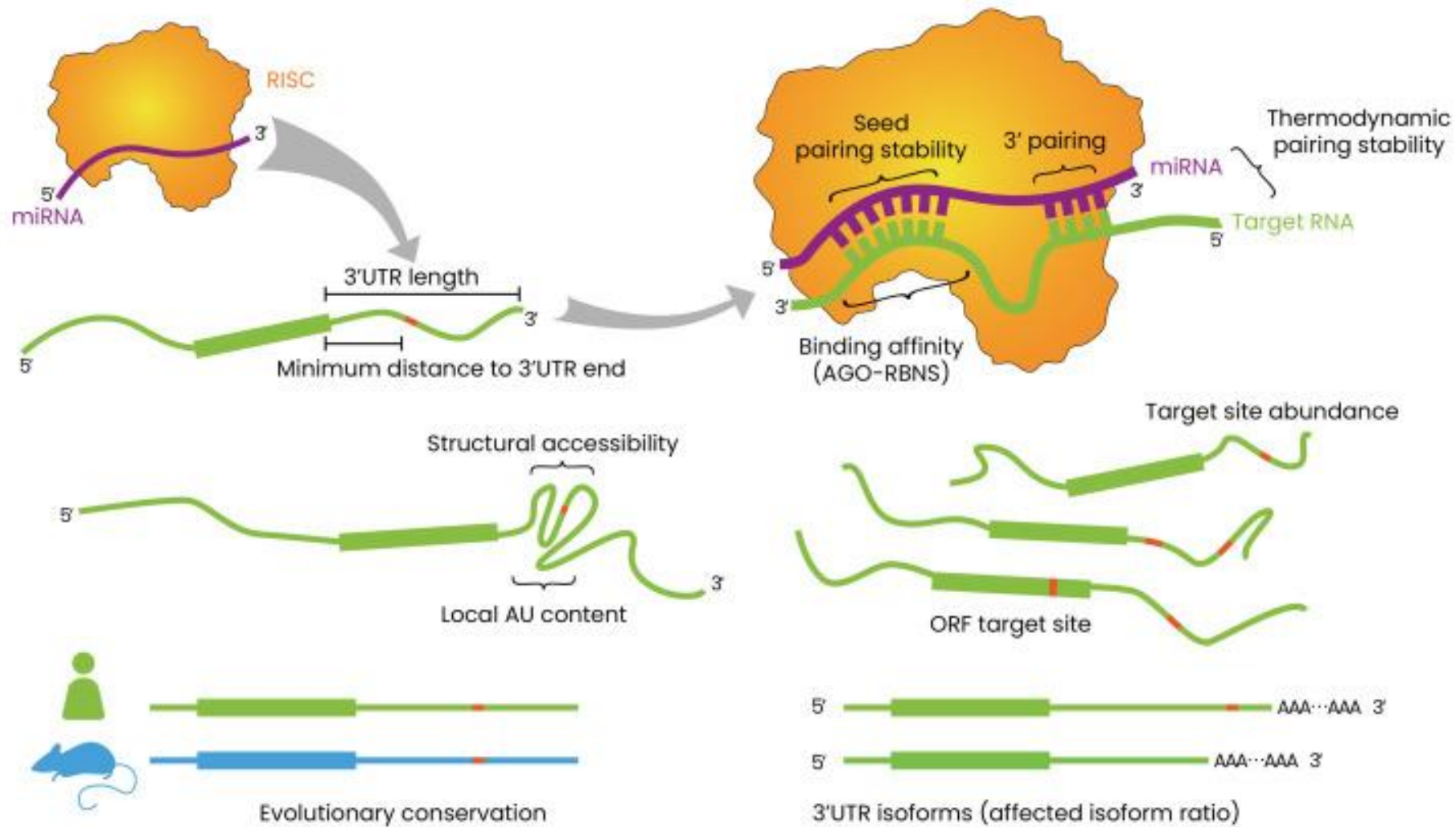


Table 1

A table of representative computational tools for miRNA target prediction and the determinants they use

Model	Seed	TPS	EC	SA	Dist.	AU	Len.	3Sup.	TA	ORFS
TargetScan7	0	SPS	0	0	0	0	0	0	0	8m
miRanda-mirSVR	0	X	0	0	0	0	0	0	X	X
DIANA-microT-CDS	0	0	0	0	0	0	X	X	X	0
MIRZA-G	0	0	0	0	0	X	X	X	X	X
PITA	Opt.	0	X	0	X	X	X	X	X	X
PicTar	0	0	0	X	X	X	X	X	X	X
RNAhybrid	Opt.	0	X	X	X	X	X	X	X	X
MicroTar	0	0	X	X	X	X	X	X	X	X

[Open in a separate window](#)

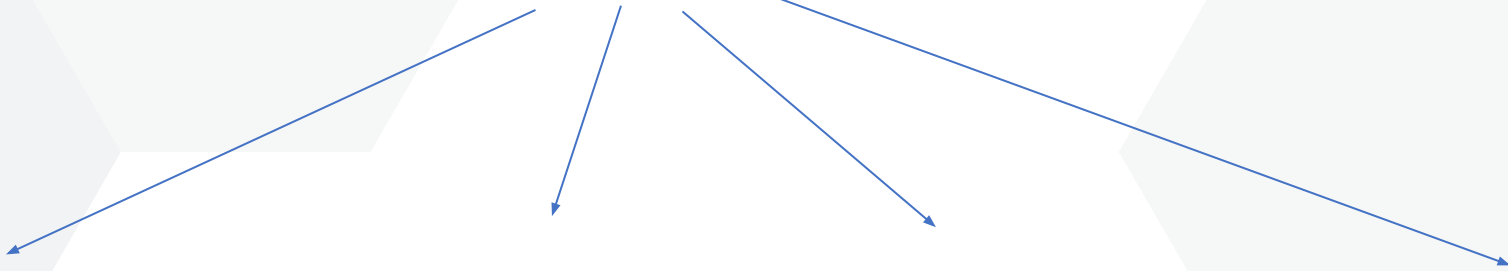
Seed, seed match or site type; TPS, thermodynamic pairing stability; EC, evolutionary conservation; SA, structural accessibility; Dist., distance to 3'UTR ends or relative position of the target sites in the 3'UTR; AU, AU or GC content; Len., length of transcript or UTR; 3Sup., 3' supplementary pairing; TA, target abundance; ORFS, ORF or CDS sites; Opt., optional; SPS, seed pairing stability; 8m, number of 8-mer sites in the ORF.

Scanning - prediction only



Ago

ucagcauagcuacgacguc miRNA, ~20nt long



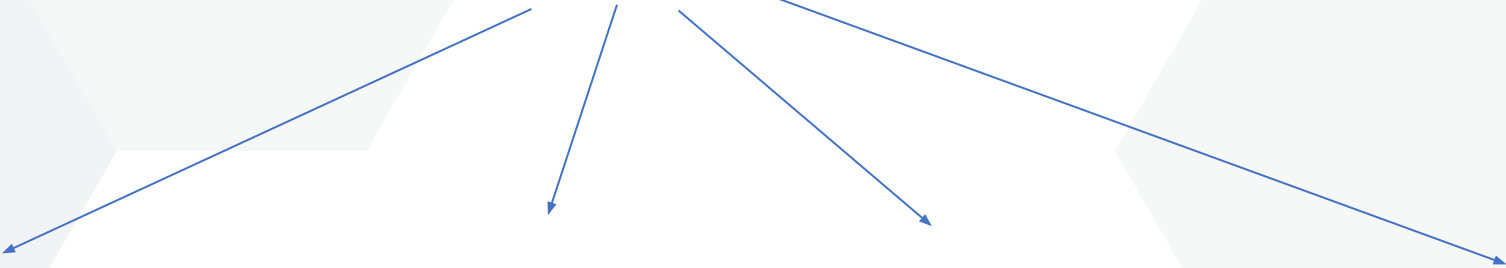
auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

Messenger RNA, 100s – 100,000s nt long

Scanning - prediction only

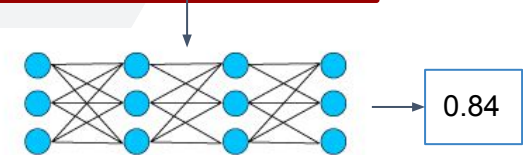


ucagcauagcuacgacguc miRNA, ~20nt long



auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

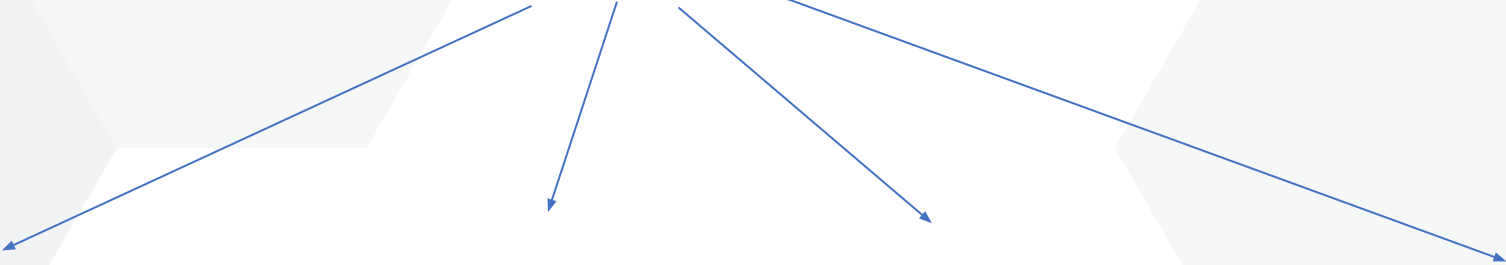
Messenger RNA, 100s – 100,000s nt long



Scanning - prediction only

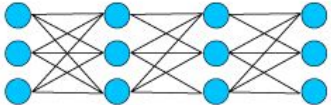


ucagcauagcuacgacguc miRNA, ~20nt long



auggaacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaacacaaauucgac...

Messenger RNA, 100s – 100,000s nt long

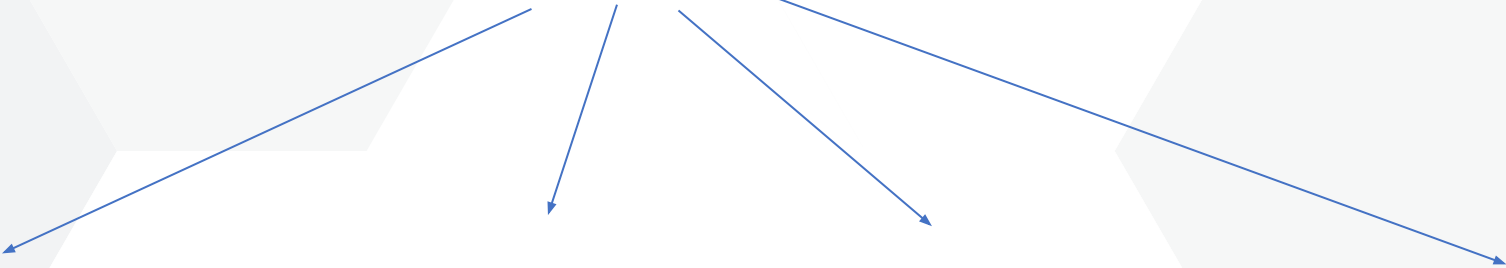


0.52

Scanning - prediction only



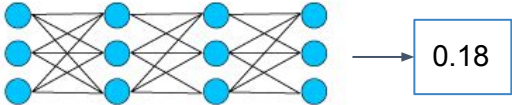
ucagcauagcuacgacguc miRNA, ~20nt long



auggacacgcggggcgcgau cgugucacguagcua tagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...



Messenger RNA, 100s – 100,000s nt long

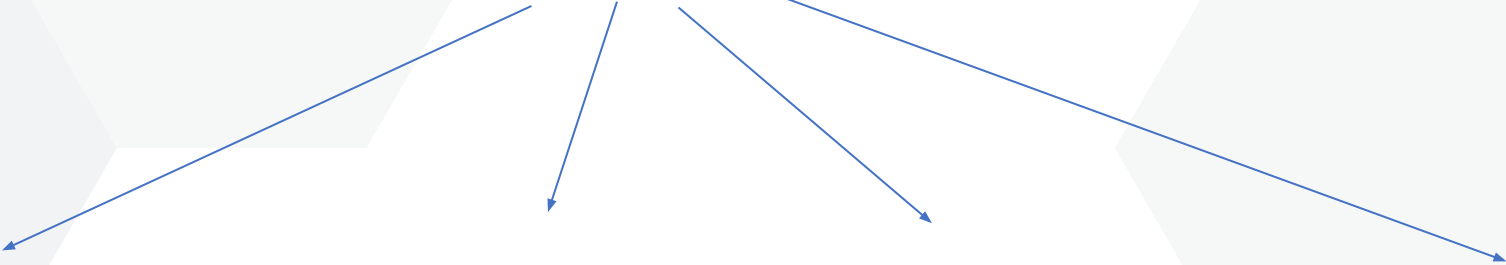


Scanning - prediction only



Ago

ucagcauagcuacgacguc miRNA, ~20nt long



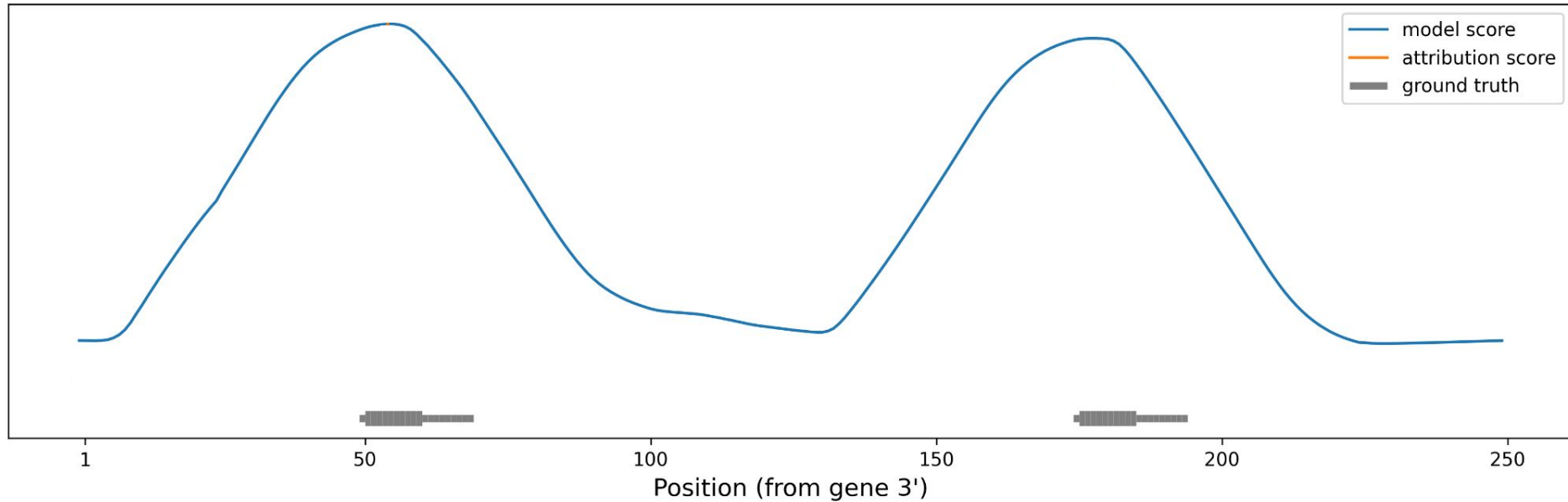
auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

Messenger RNA, 100s – 100,000s nt long



Scanning - prediction only

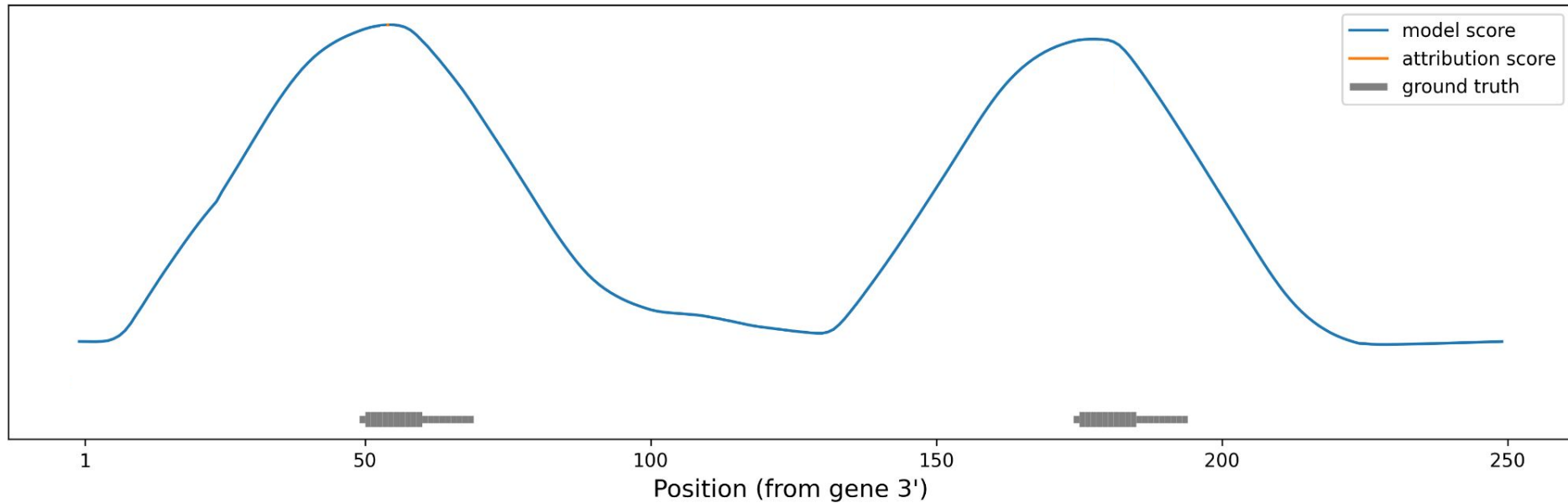
Narrowing the peaks



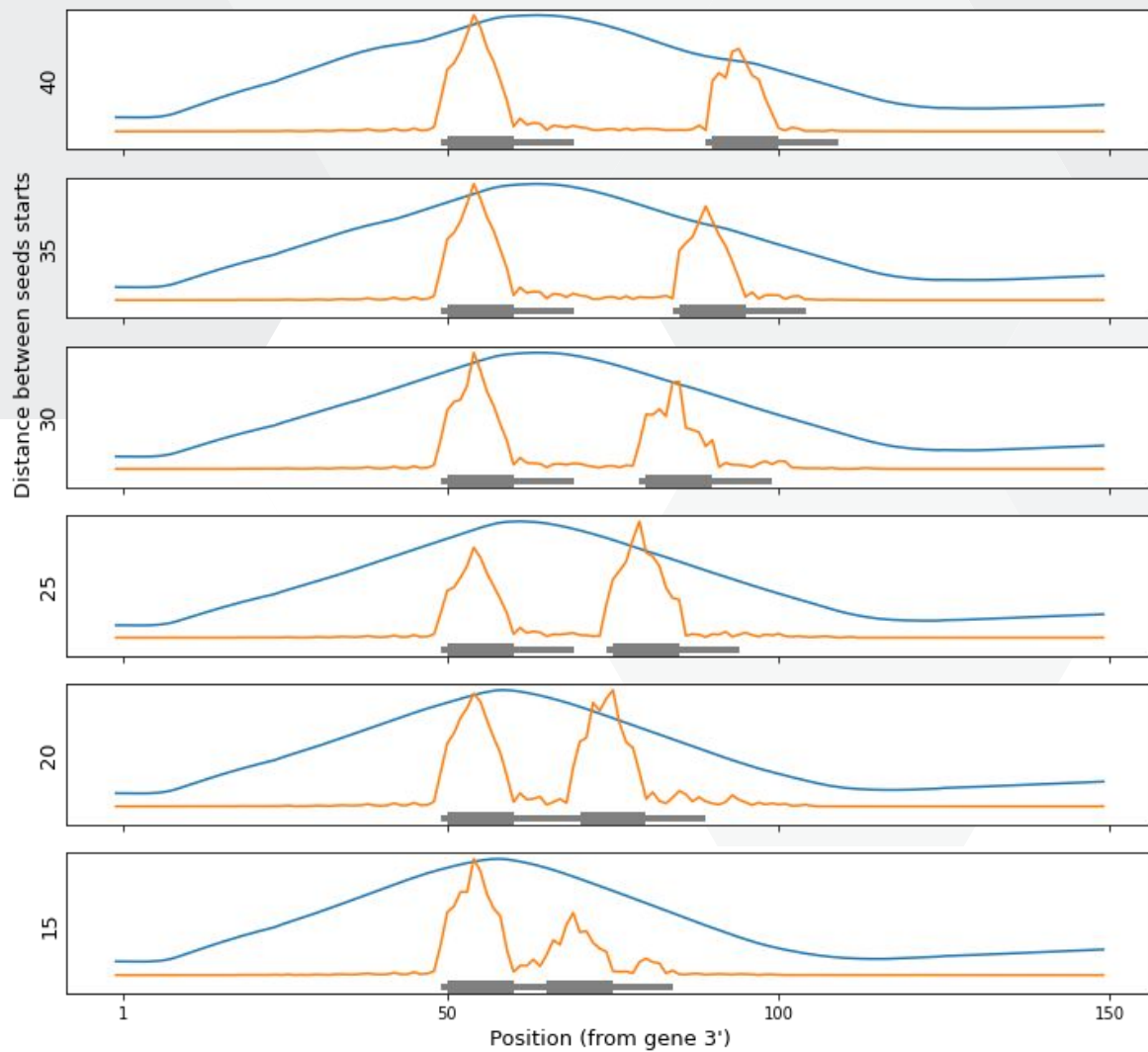
Scanning - including attribution score

Narrowing the peaks

TCACCTTTCTCTCGAAGTATGGGACACGAATCTT**CAT****CAT****CCAAC**TGTÄ→
-----TGAG**GTA**-**G**TAG**GTT**GTATAG
→



Close by peaks

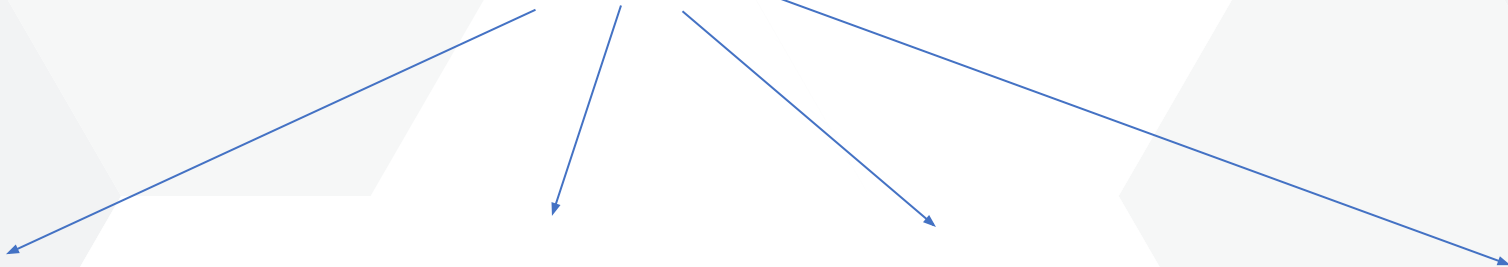


Scanning - using attribution score

Math:
Score at a position =
prediction * attribution_score



ucagcauagcuacgacguc miRNA, ~20nt long



auggacacgcggggcgcgaucguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

Messenger RNA, 100s – 100,000s nt long

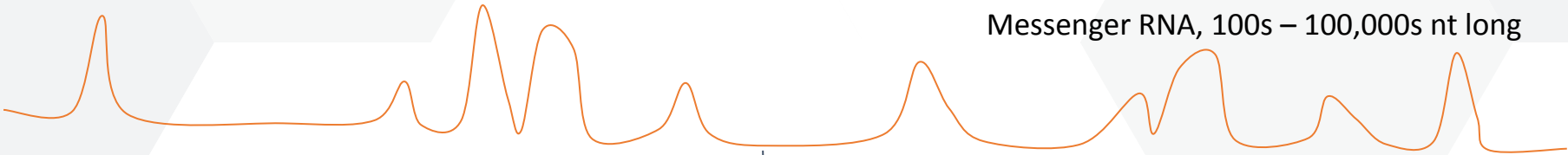


Scanning - using attribution score

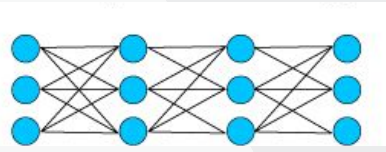
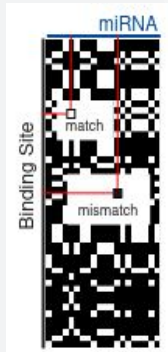


auggacacgcggggcgcgau cgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaaccacaauucgac...

Messenger RNA, 100s – 100,000s nt long



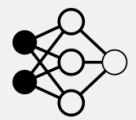
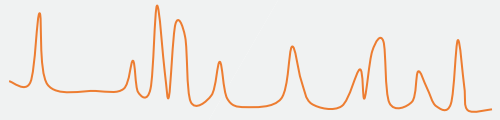
0.82



CNN

ucagcauagcuacgacguc

miRNA, ~20nt long



regression model

prediction

auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaacCacaauucgac...

Messenger RNA, 100s – 100,000s nt long

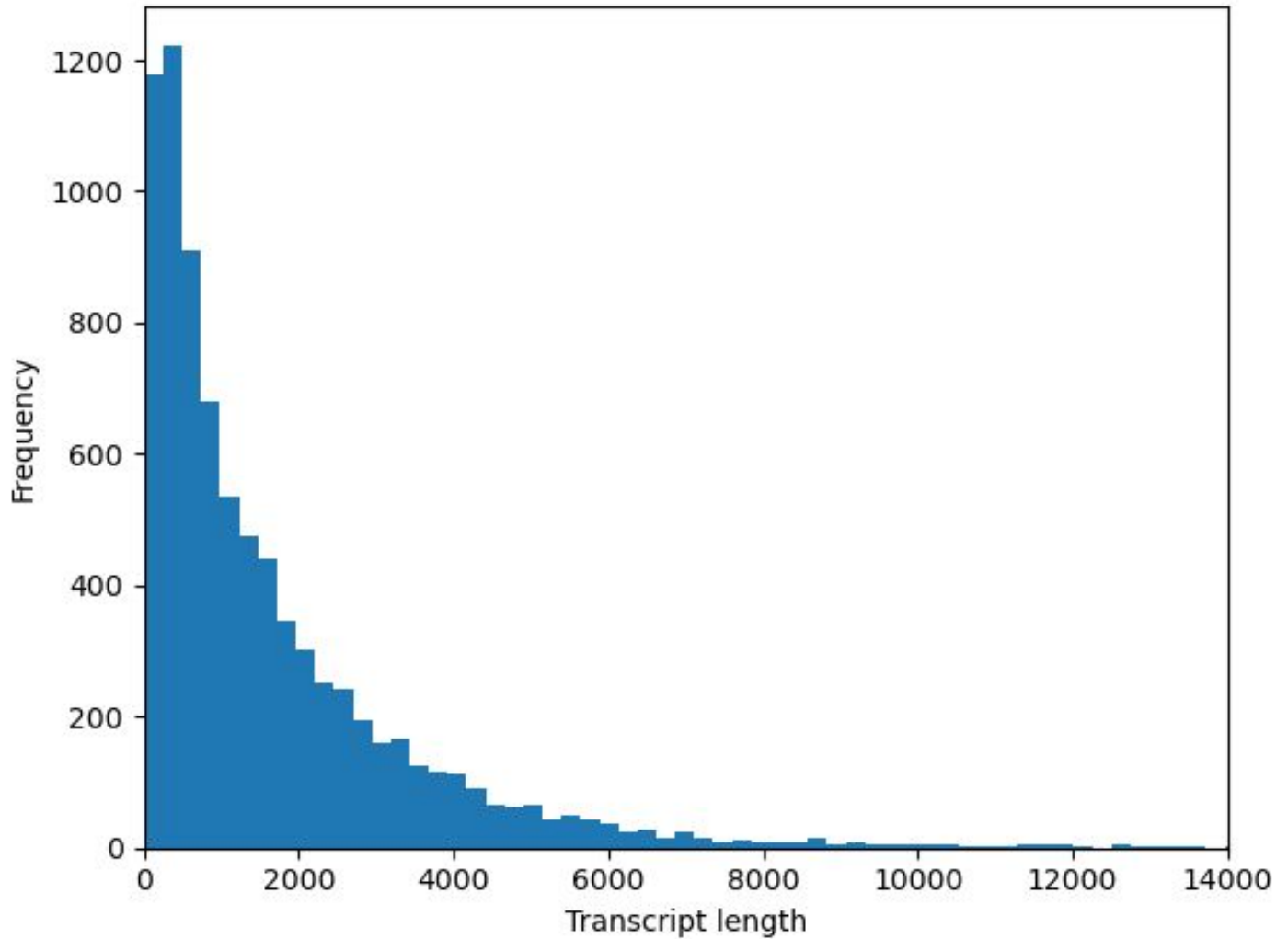


Dataset

- inputs

Lengths of transcript sequences

Sequences too long

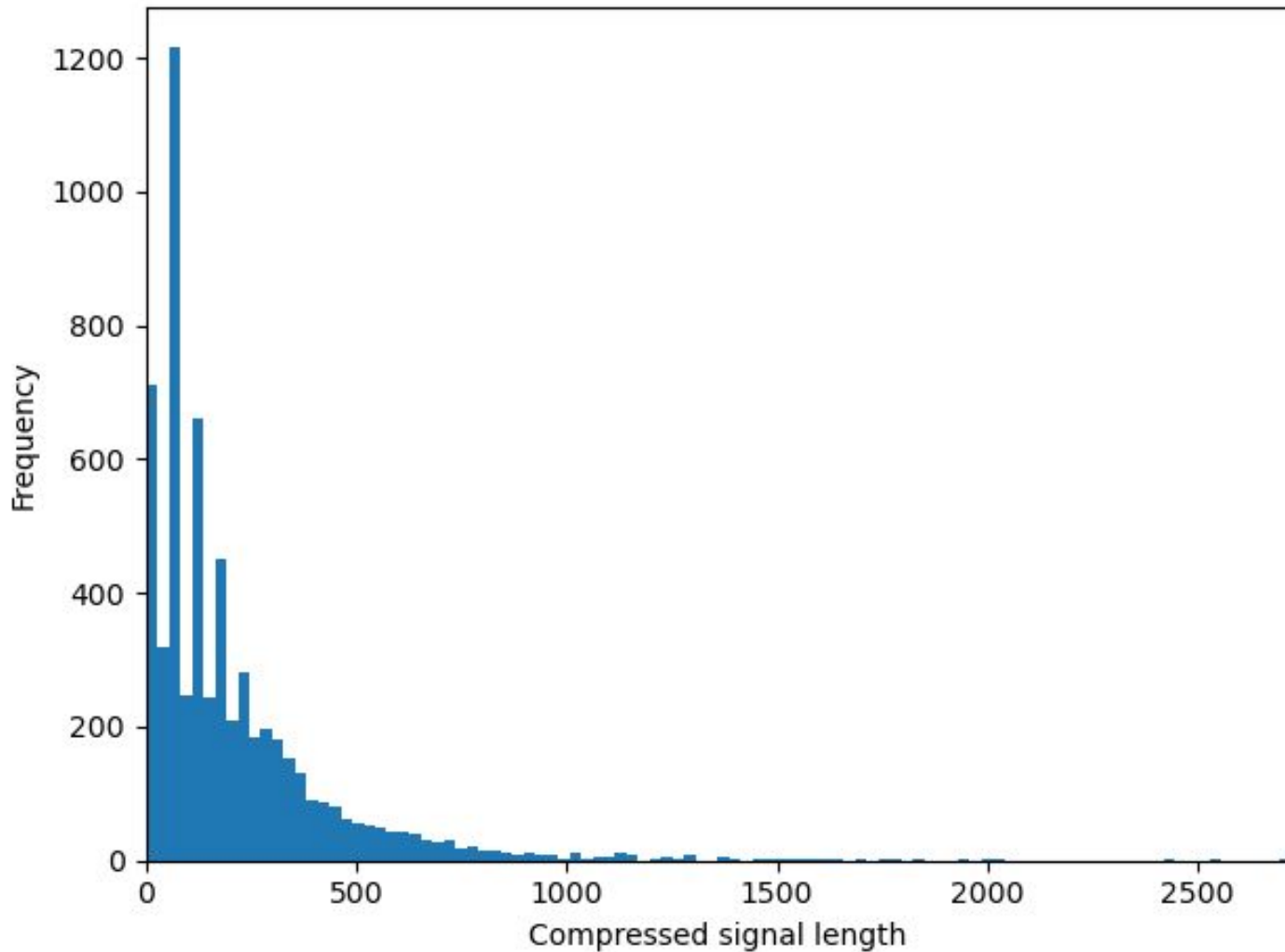


Compressed inputs

Lengths of signals compressed

Per transcript

Longest: 2719



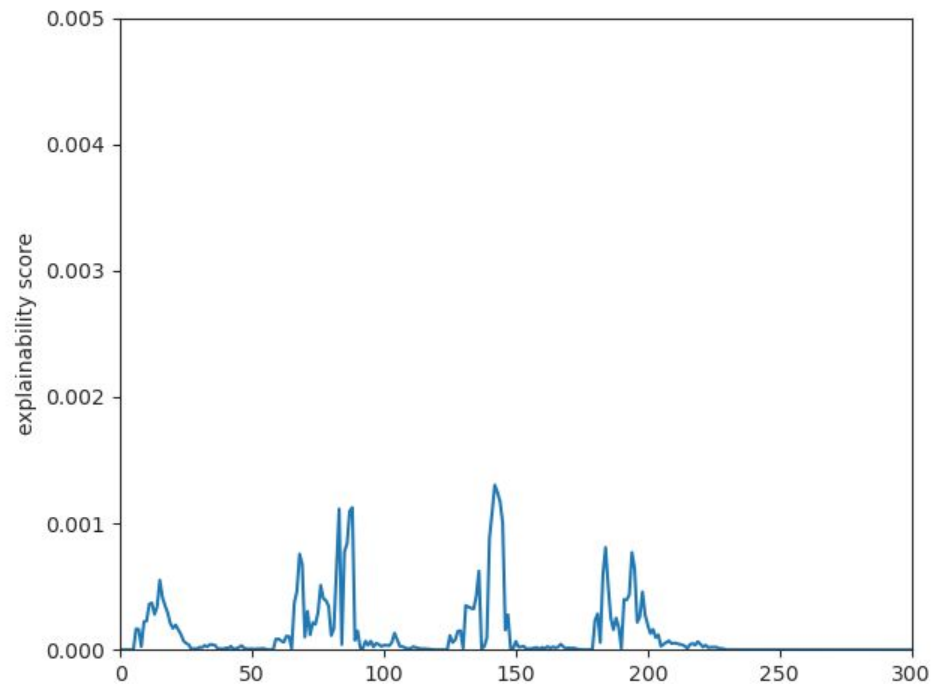
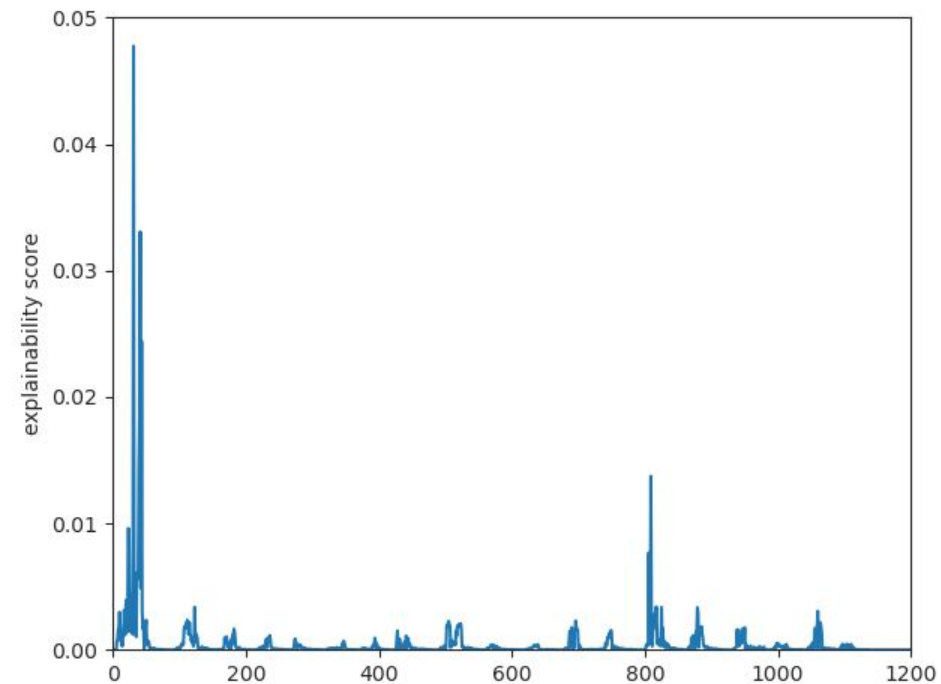
Signals preprocessing

Highly sparse → compression: $(number_of_zeroes \% 100) + 1$

Normalization to $\langle 0.00001, 1 \rangle$

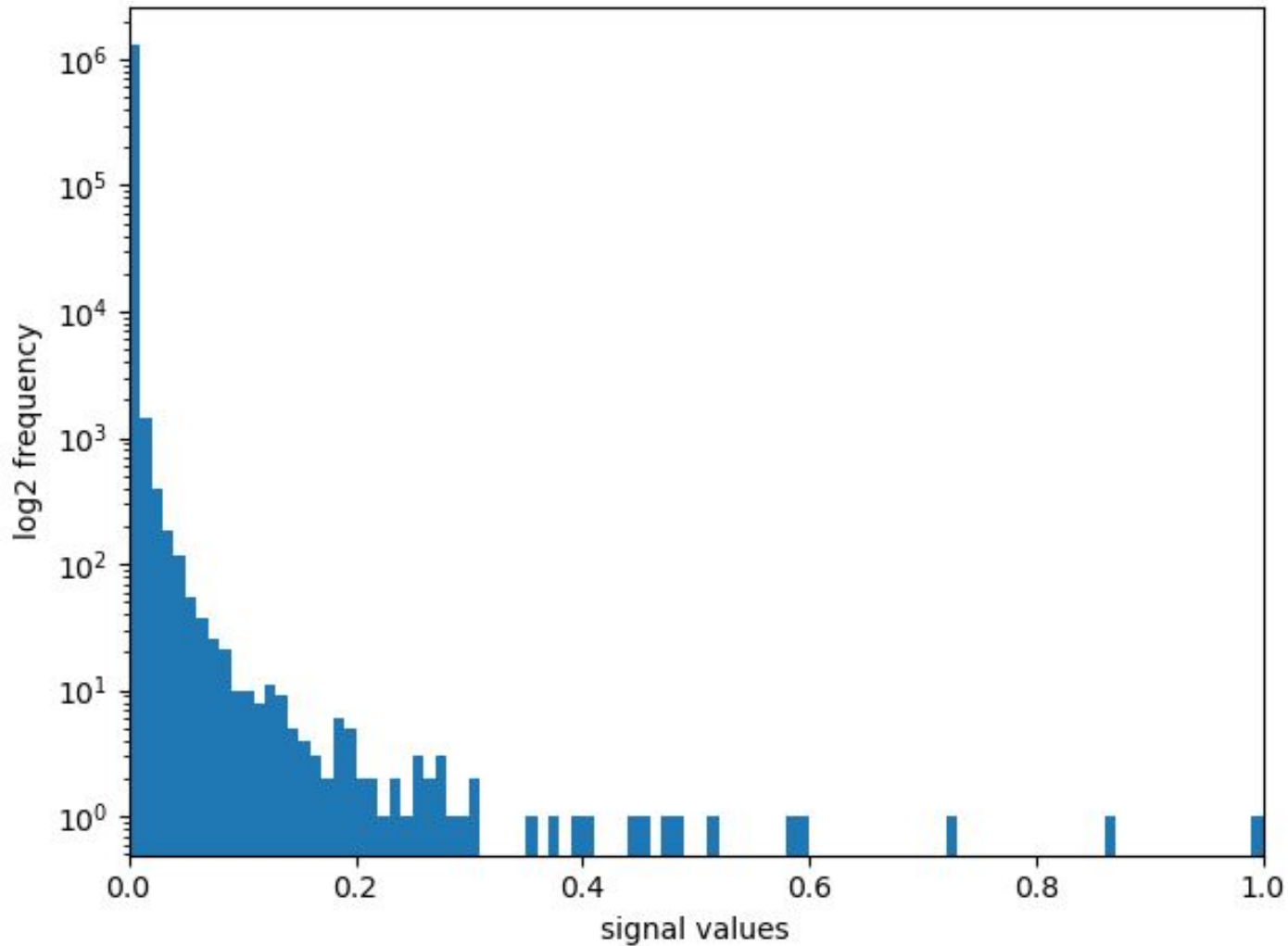
0 used for padding

Signal samples (compressed)



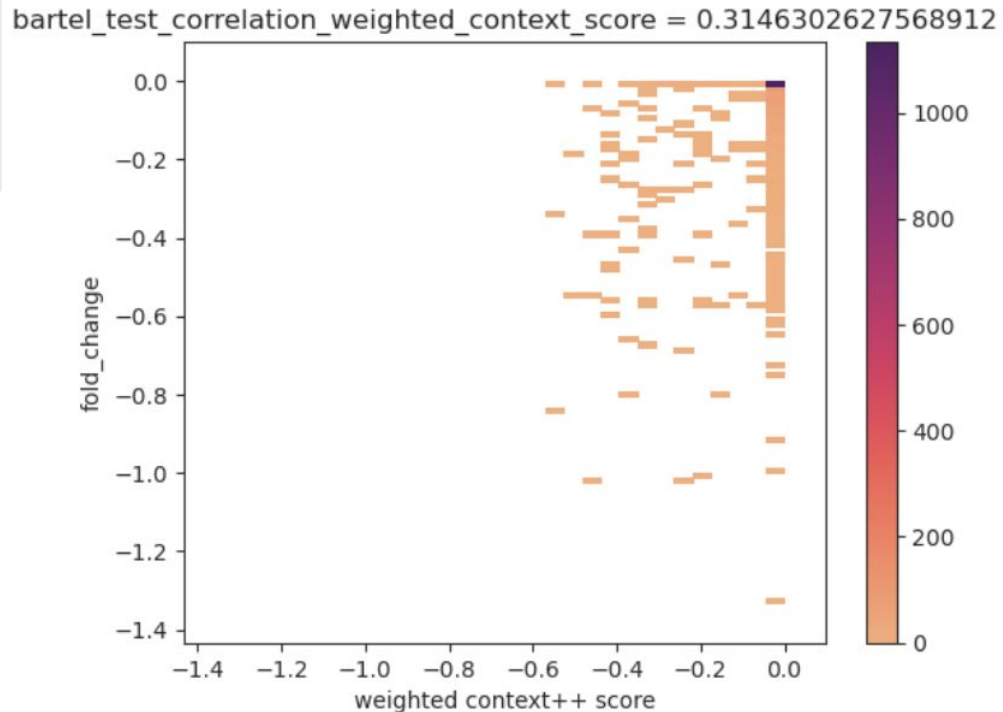
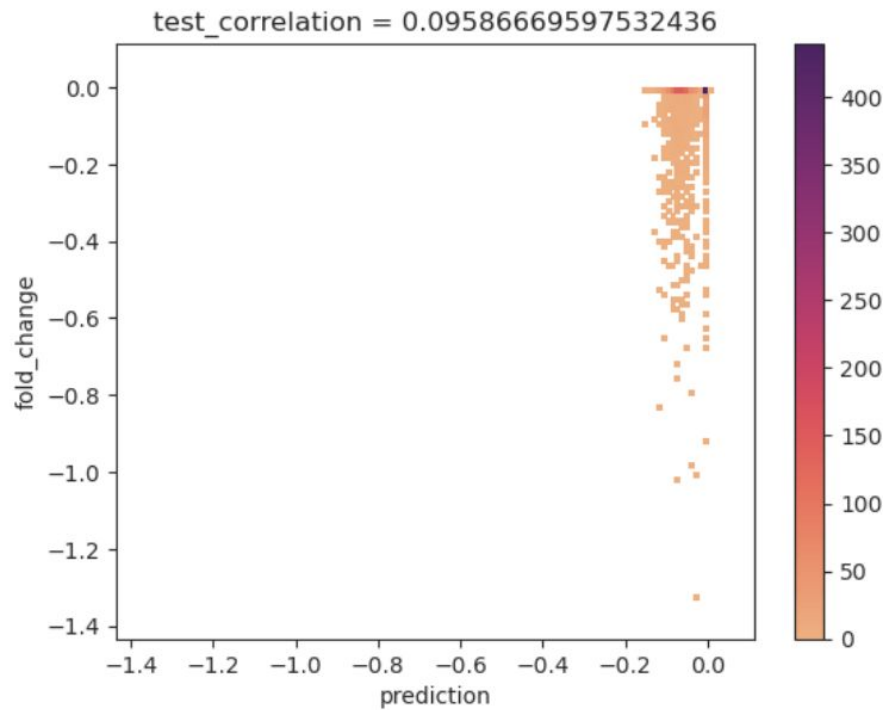
Signal
Values
histogram
over all
samples

(before
padding)

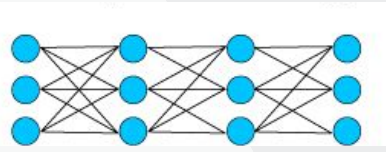
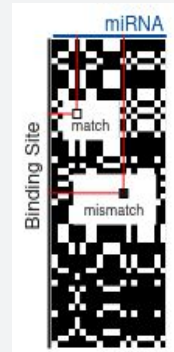


CNN + RNN + pooling

State-of-the-art based on feature selection (TargetScan)



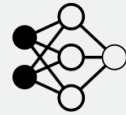
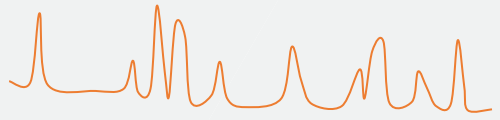
Summary



CNN

ucagcauagcuacgacguc

miRNA, ~20nt long

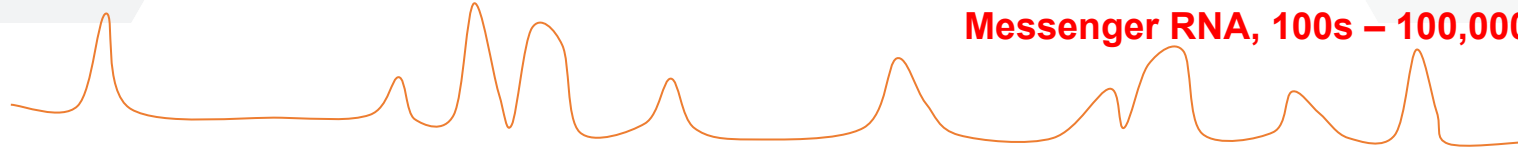


regression model

prediction

auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaacCacaauucgac...

Messenger RNA, 100s – 100,000s nt long

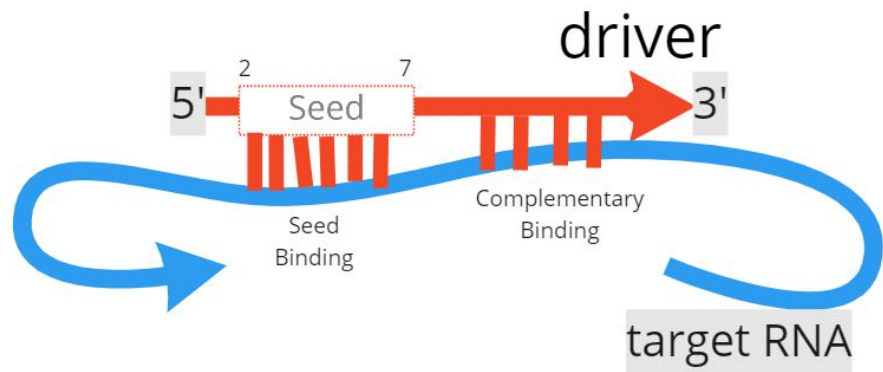


Advantages & Disadvantages of our two-part approach

1. Shields from sequence and overfitting on simple patterns like seed binding
2. Generalizes across miRNAs

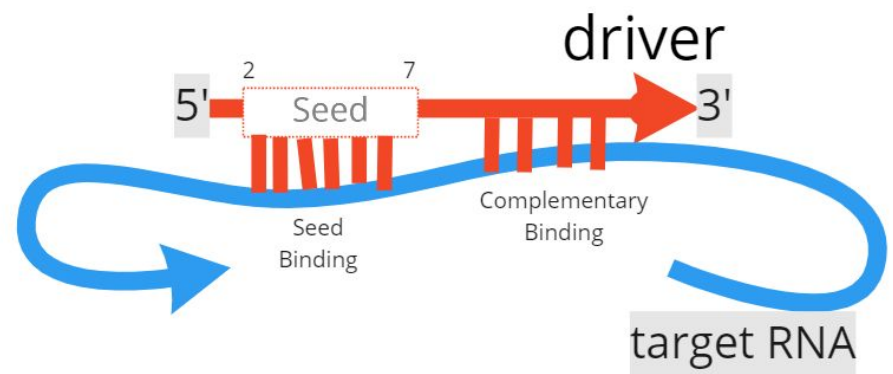
1. First model is not perfect which leads to accumulation of mistakes to the second model
2. Cannot propagate error through second model to the first model

Summary



miRNA:
TGAGGTAGTA
GGTTGTATAG

Binding site:
ATGTCAACCTA
CCTACTTCTAA
GCACAGGGTAT
GAAGCTCTCTT
TCCA

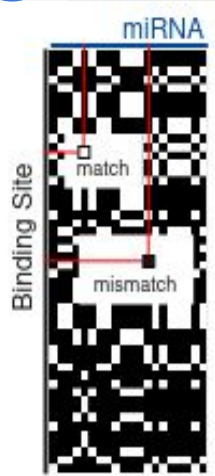


A-T

G-C

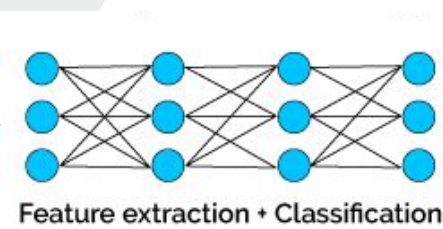
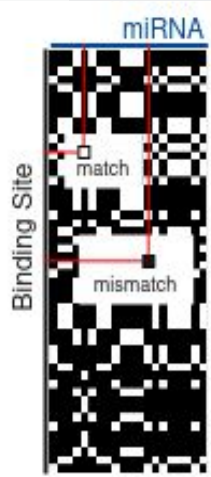
miRNA:
TGAGGTAGTA
GGTTGTATAG

Binding site:
ATGTCAACCTA
CCTACTTCTAA
GCACAGGGTAT
GAAGCTCTCTT
TCCACT

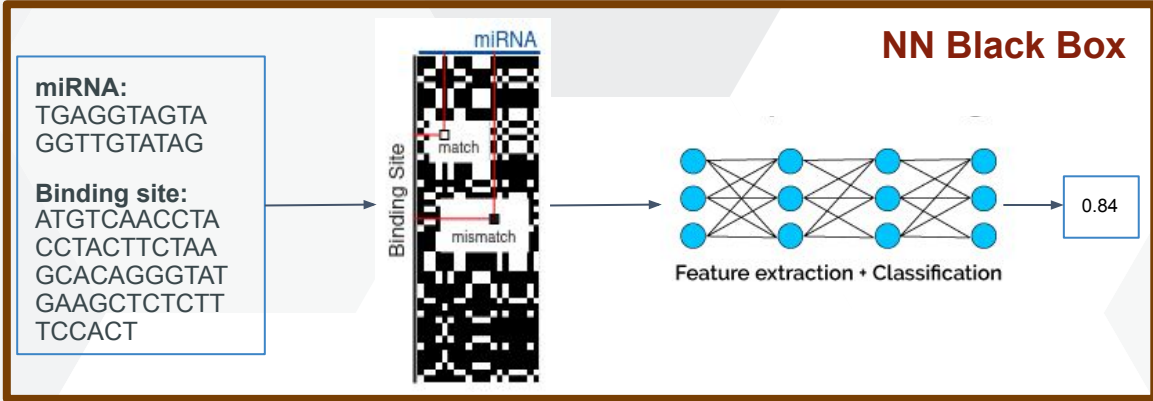


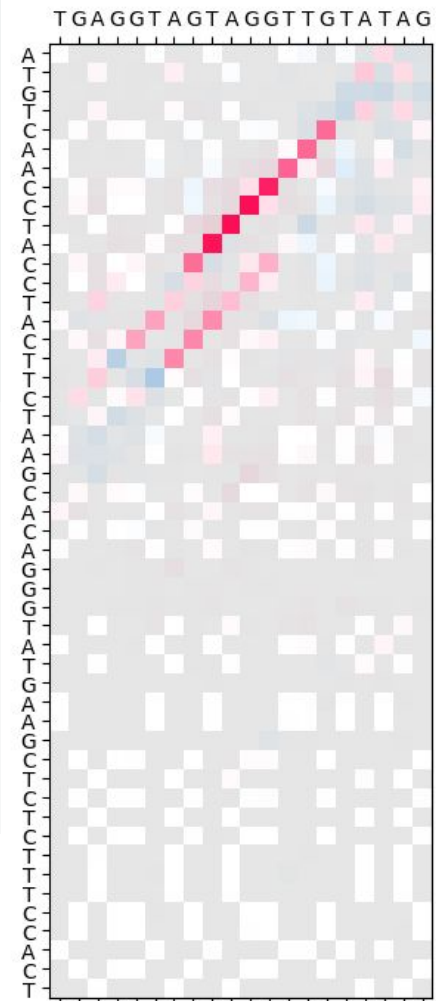
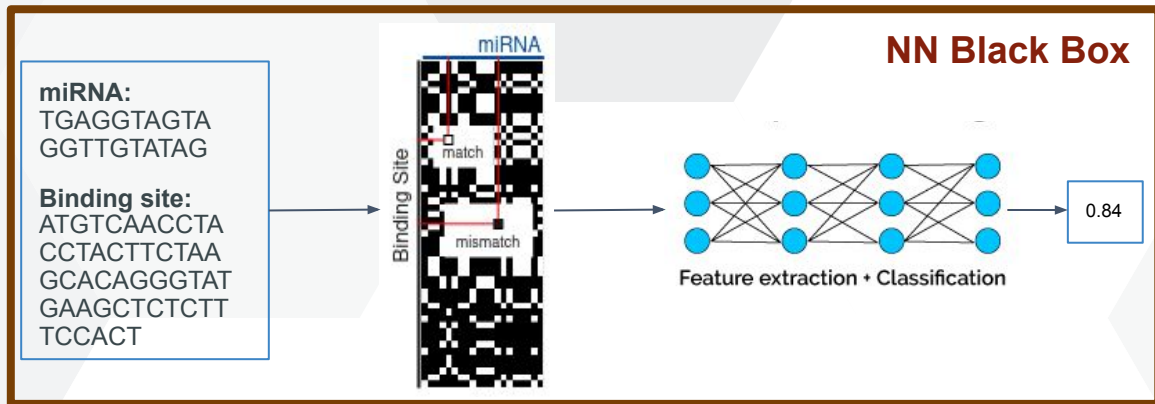
miRNA:
TGAGGTAGTA
GGTTGTATAG

Binding site:
ATGTCAACCTA
CCTACTTCTAA
GCACAGGGTAT
GAAGCTCTCTT
TCCA

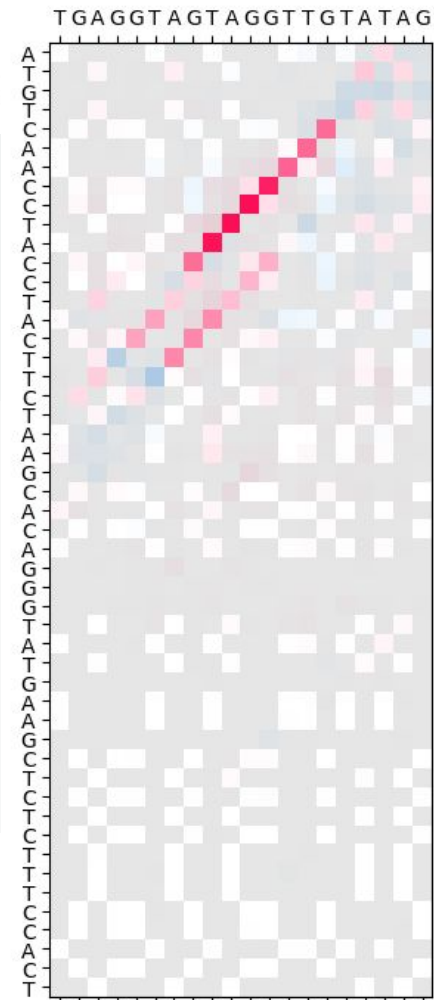
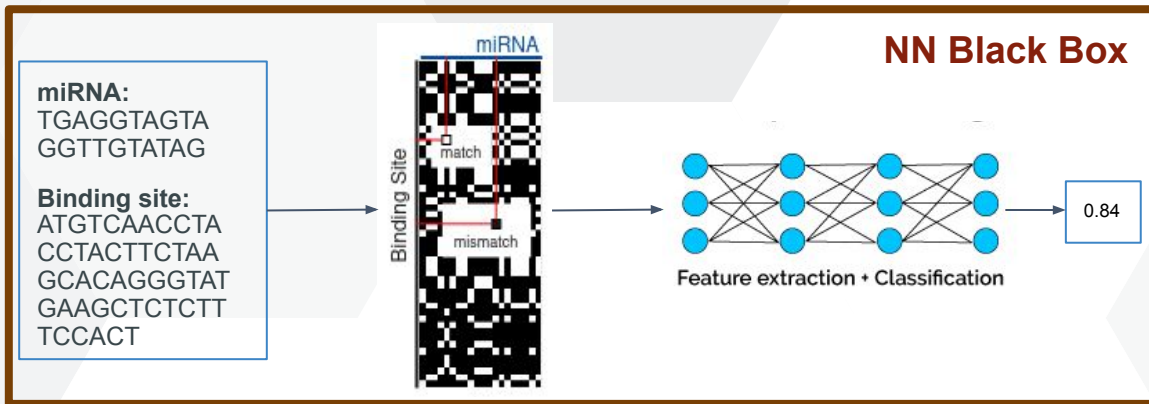


0.84

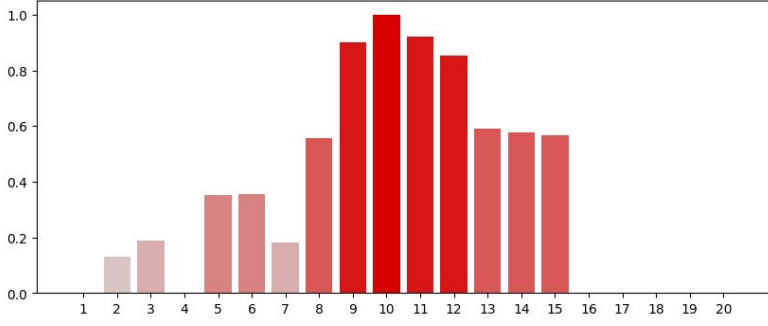
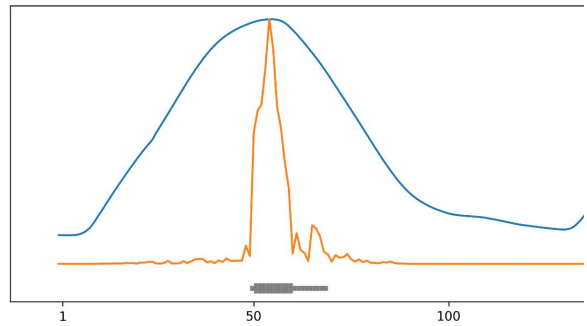
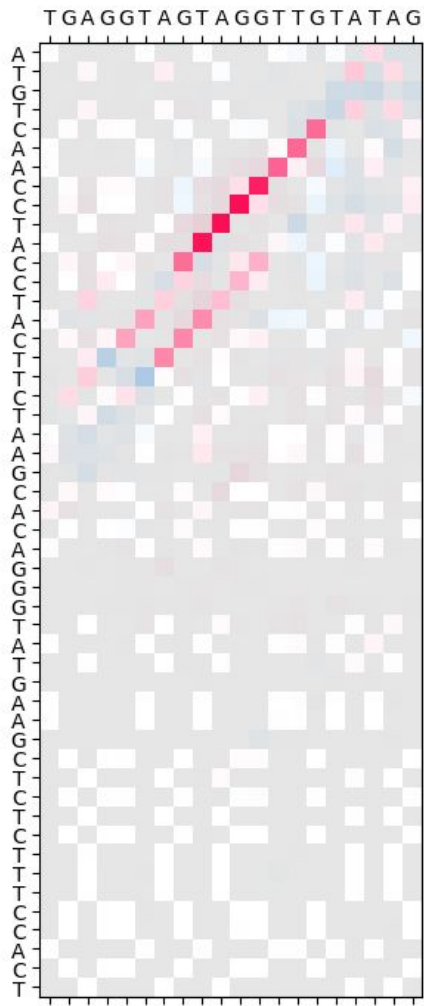
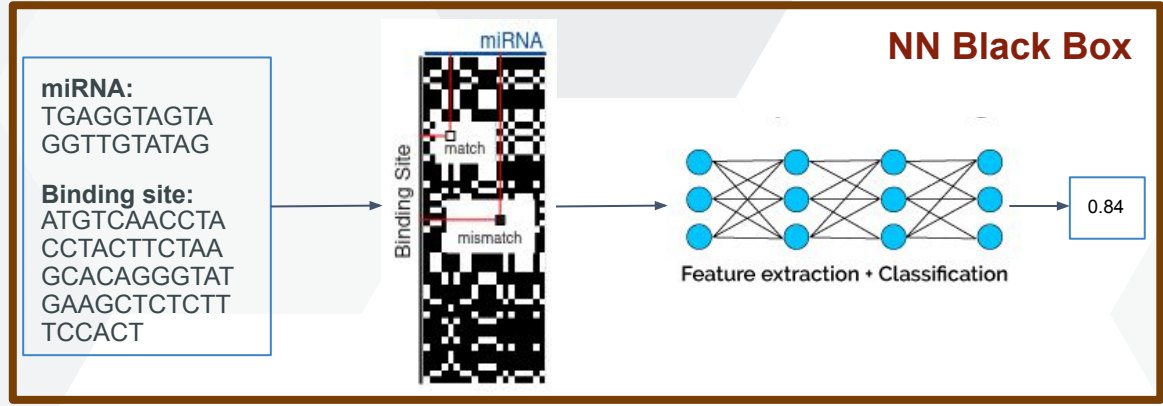
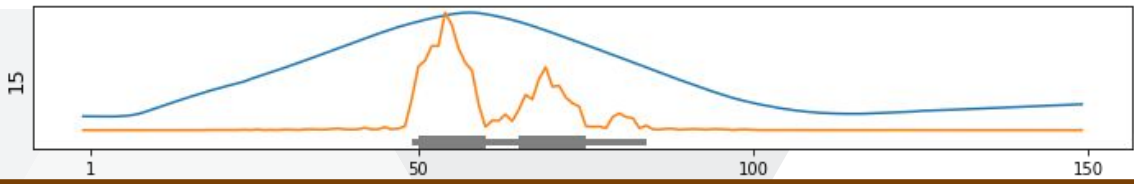




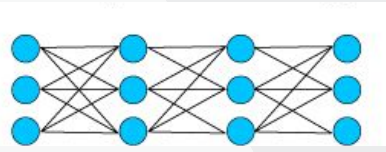
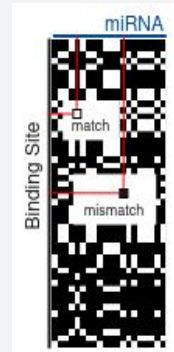
TCACCTTTCTCTCGAAGTATGGGACACGAATCTTCATCCATCCAACGTGTA-
- - - - - TGAGGTA-GTAGGTTGTATAG



TCACCTTTCTCTCGAAGTATGGGACACGAATCTT**CATCCATCCAAC**TGTÄ-
 · | | · | | · | | | | | | | | · · ·
 - - - - - **TGAGGTA-GTAGGTTGTÄTÄG**



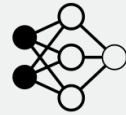
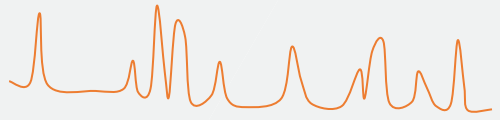
Summary



CNN

ucagcauagcuacgacguc

miRNA, ~20nt long

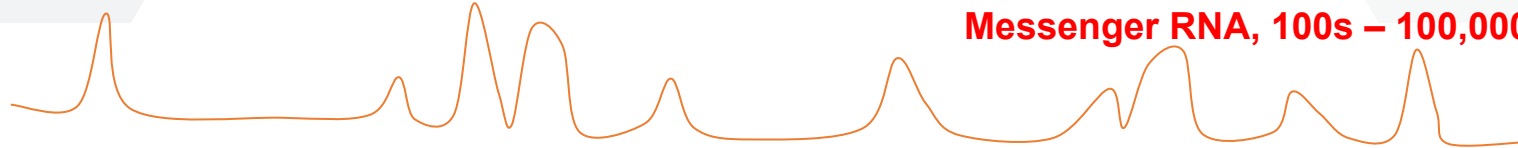


regression model

prediction

auggacacgcggggcgcgaucgugucacguagcuacagucaugcaugucguagcuagcacucgucgucgagcuacgugggagacugcgaaaaaaacCacaauucgac...

Messenger RNA, 100s – 100,000s nt long



Future work

1. Include other features
 - a. Genomic conservation - score / multiple sequence alignment / tree
 - b. RNA Binding Proteins - binding sites
 - c. Sequence?
 - d. ...
 - e. Ablation studies
2. If two-part approach does not work
 - a. Simplify regression to classification task?
 - b. Skip the two-part approach and go with sequence? (in progress)
 - i. HyenaDNA - pretrained single nucleotide resolution transformer for long sequences
 1. Use for embeddings
 2. Use full with regression head

Sources

SHAP <https://github.com/shap/shap>

Determinants of Functional MicroRNA Targeting
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9880601/>

miRBind: A Deep Learning Method for miRNA Binding Classification
<https://www.mdpi.com/2073-4425/13/12/2323>

Using Attribution Sequence Alignment to Interpret Deep Learning Models for miRNA Binding Site Prediction <https://www.mdpi.com/2079-7737/12/3/369>



Thank you for your Attention!



83

