

# Towards physicochemical bioinformatics

Bojan Žagrović  
Max Perutz Labs & University of Vienna

**Computational Chemistry Days**  
IRB Zagreb, May 9<sup>th</sup> 2026

# Biomolecular sequences analysis: a foundation of modern molecular biology

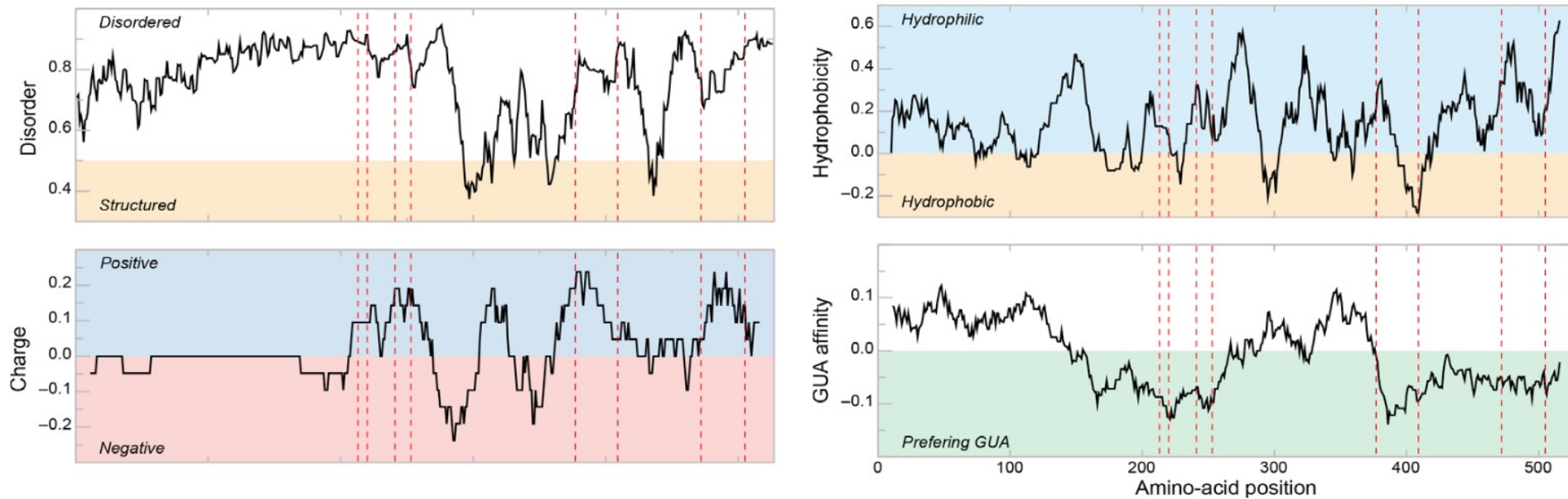
FUS (human)

MASNDYTQQATQSYGAYPTQPGQGYSQQSSQPYGQQSYSGYSQSTDTSGYGQSSYSSYGQSQNTGYGTQSTPQGYGS  
TGGYGSSQSSQSSYGQQSSYPGYGQQPAPSSTSGSYGSSSQSSSYGQPQSGSYSQQPSYGGQQQSYGQQQSYNPPQGY  
GQQNQYNSSSGGGGGGGGGNYGQDQSSMSSGGGSGGGYGNQDQSGGGGSGGYGQQDRGGRGRGGSGGGGG  
GGGGYNRSSGGYEPRGRGGGRGGRGGMGGSDRGGFNKFGGPRDQGSRDSEQDNSDNNTIFVQGLGENVTIESV  
ADYFKQIGIIKTNKKTGQPMINLYTDRETGKLGKGEATVSFDDPPSAKAAIDWFDGKEFSGNPIKVSFATRRADFNRRGGGNG  
RGGRRGGPMGRGGYGGGGSGGGGRGGFPSGGGGGGGQQRAGDWKCPNPTCENMNFSWRNECNQCKAPKPDG  
PGGGPGGSHMGGNYGDDRRGGRGGYDRGGYRGRGGDRGGFRGGRRGGDRGGFGPGKMDSRGEHRQDRRERPY

# Biomolecular sequences analysis: a foundation of modern molecular biology

FUS (human)

```
MASNDYTQQATQSYGAYPTQPGQGYSQQSSQPYGQQSYSGYSQSTDTSGYGQSSYSSYGQSQNTGYGTQSTPQGYGS  
TGGYGSSQSSQSSYGQQSSYPGYGQQPAPSSTSGSYGSSSQSSSYGQPQSGSYSQQPSYGGQQQSYGQQQSYNPPQGY  
GQQNQYNSSSGGGGGGGGGGNYGQDQSSMSSGGGSGGGYGNQDQSGGGGSGGYGQQDRGGRGRGGSGGGGG  
GGGGYNRSSGGYEPRGRGGGRGGRMGGSDRGGFNKFGGPRDQSRHDSEQDNSDNNTIFVQGLGENVTIESV  
ADYFKQIGIIKTNKKTGQPMINLYTDRETGKLGKGEATVSFDDPPSAKAAIDWFDGKEFSGNPIKVSFATRRADFNRGGGNG  
RGGRRGGGPMGRGGYGGGGSGGGGRGGFSGGGGGGGQQRAGDWKCPNPTCENMNFWRNECNQCKAPKPDG  
PGGGPGGSHMGGNYGDDRRGGRRGGYDRGGYRGRGGDRGGFRGGRRGGDRGGFGPGKMDSRGEHRQDRRERPY
```



understanding biomolecular sequences as physicochemical objects

## Example 1

### Frameshifts: major impact on reading of genetic message

#### mRNA

```
atggggaaat tgatcaggat ggggccgcaa gagaggtggt tactccggac  
aaagcggctt cattggagtc gcctcctctt cttactggga atgttgatca  
tcggttctac ttatcagcac ctaggagac cccggggcct ttcctcattg
```

#### protein

```
wt  1MGKLIRMGPQERWLLRTKRLHWSRLLFLLGMLIIGSTYQHLRRPRGLSSL50  
+1  1WGN-SGWGRKRGGYSGQSGFIGVASSSYWEC-SSVLLISTLGDPGAFPHC50  
-1  1GEIDQDGAAREVVTPDKAASLESPELLTGNVDHRFYLSAP-ETPGPFLIV50
```

## Example 1

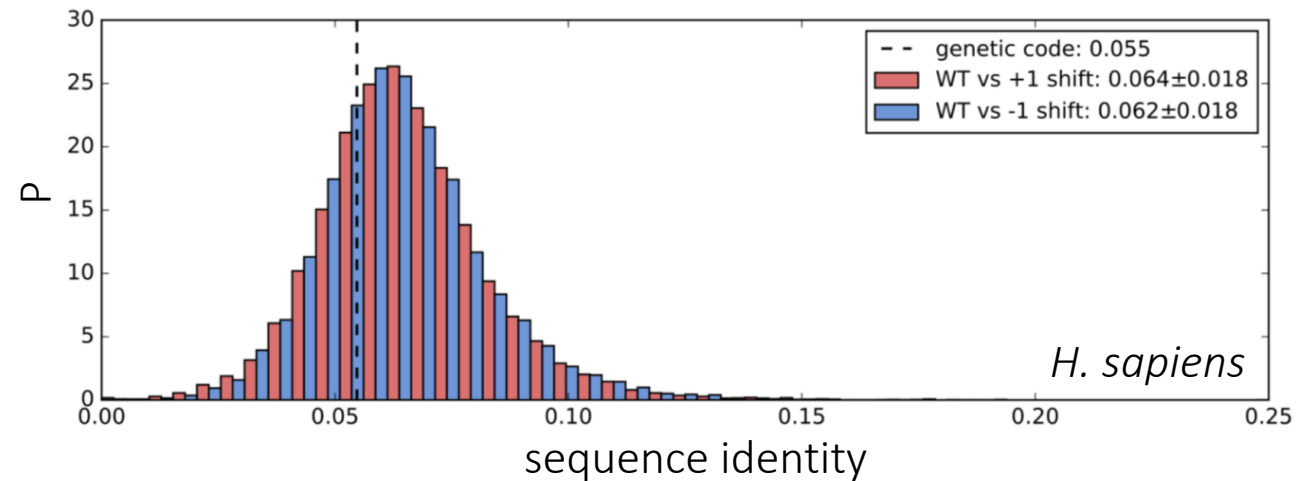
### Frameshifts: major impact on reading of genetic message

#### mRNA

```
atggggaaat tgatcaggat ggggccgcaa gagaggtggt tactccggac  
aaagcggctt cattggagtc gcctcctctt cttactggga atgttgatca  
tcggttctac ttatcagcac cttaggagac cccggggcct ttctcattg
```

#### protein

```
wt  1MGKLIRMGPQERWLLRTKRLHWSRLLFLLGMLIIGSTYQHLRRPRGLSSL50  
+1  1WGN-SGWGRKRGGYSGQSGFIGVASSSYWEC-SSVLLISTLGDPGAFFHC50  
-1  1GEIDQDGAAREVVTPDKAASLESPELLTGNVDHRFYLSAP-ETPGPFLIV50
```

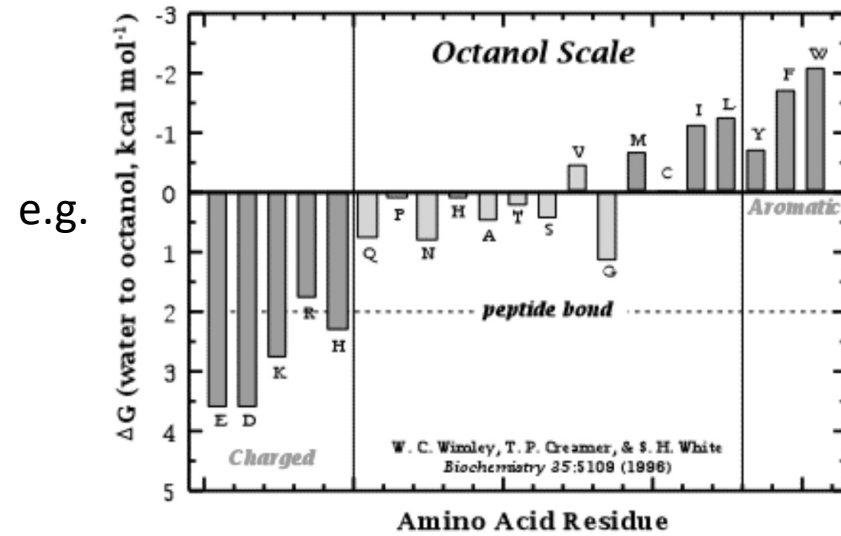


## How related are amino acids encoded by frameshifted codons?

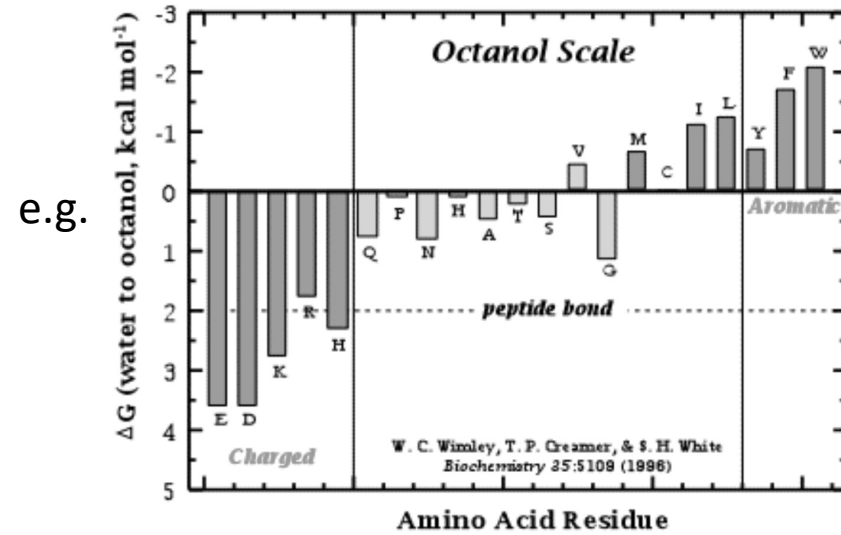
		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA Stop UAG Stop	UGU } Cys UGC } UGA Stop UGG Trp	U C A G
	C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }	U C A G
	A	AUU } AUC } Ile AUA } AUG Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }	U C A G
	G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }	U C A G

e.g. Phe → Phe, Leu, Ile, Val, Ser

>600 different amino-acid properties (Kawashima et al., **NAR**, D202, 2008)



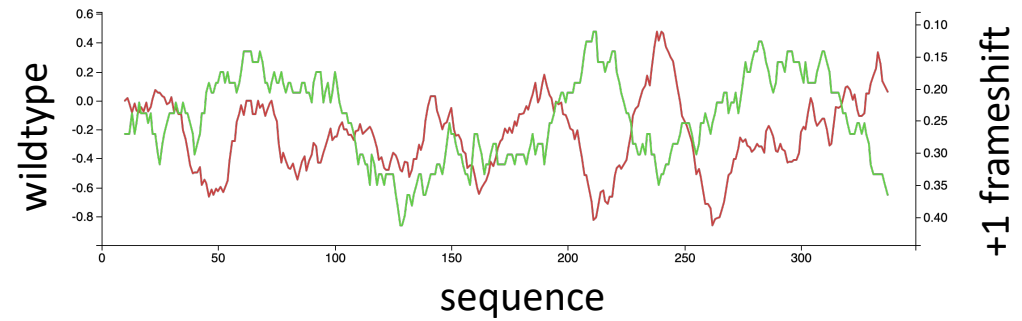
>600 different amino-acid properties (Kawashima et al., **NAR**, D202, 2008)



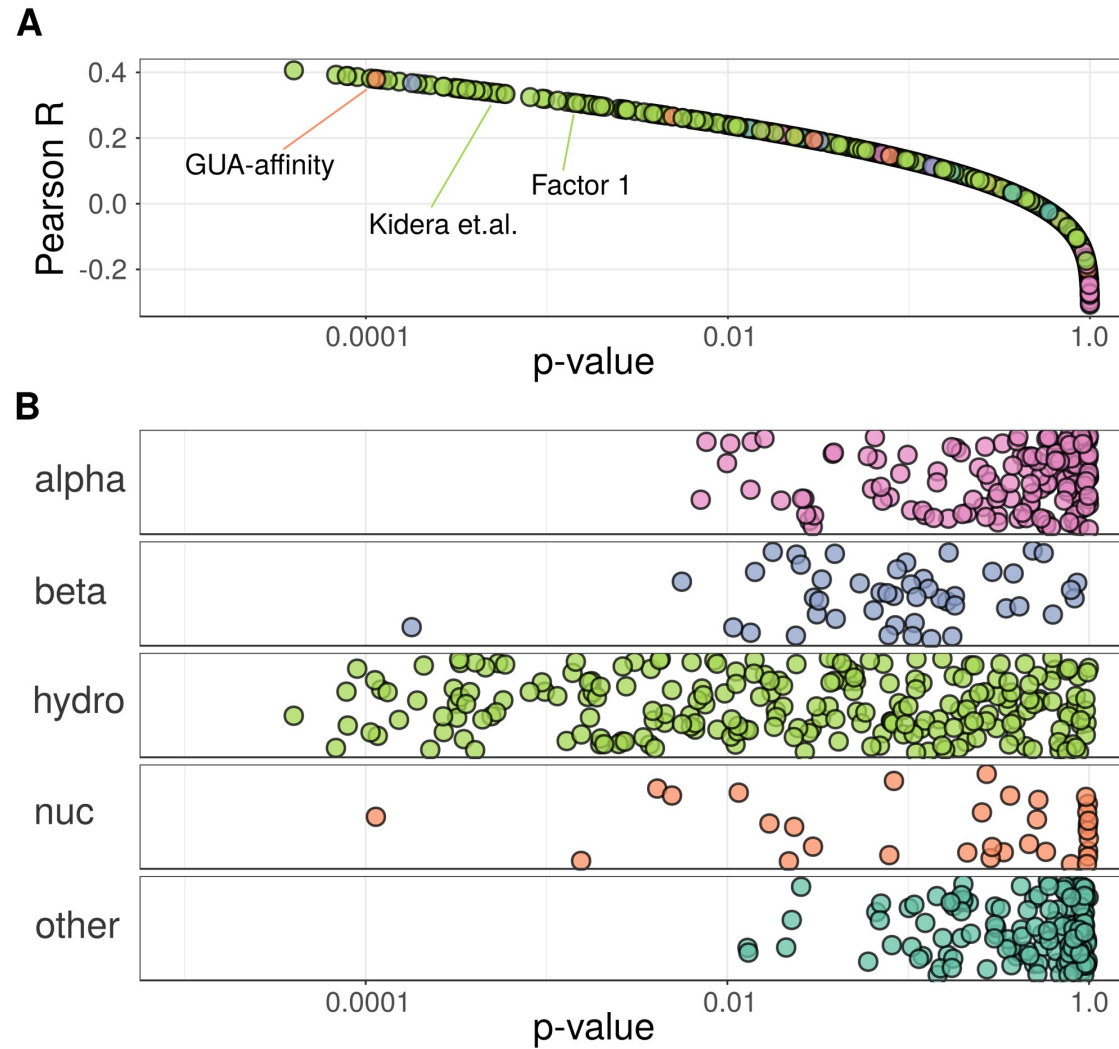
### 1. genetic code



### 2. sequence profiles

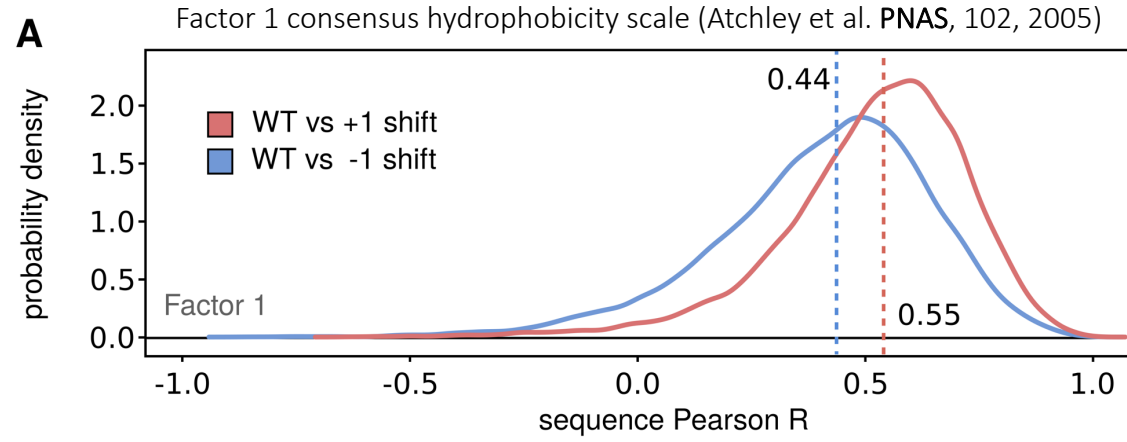


# Impact of frameshifting at the level of genetic code

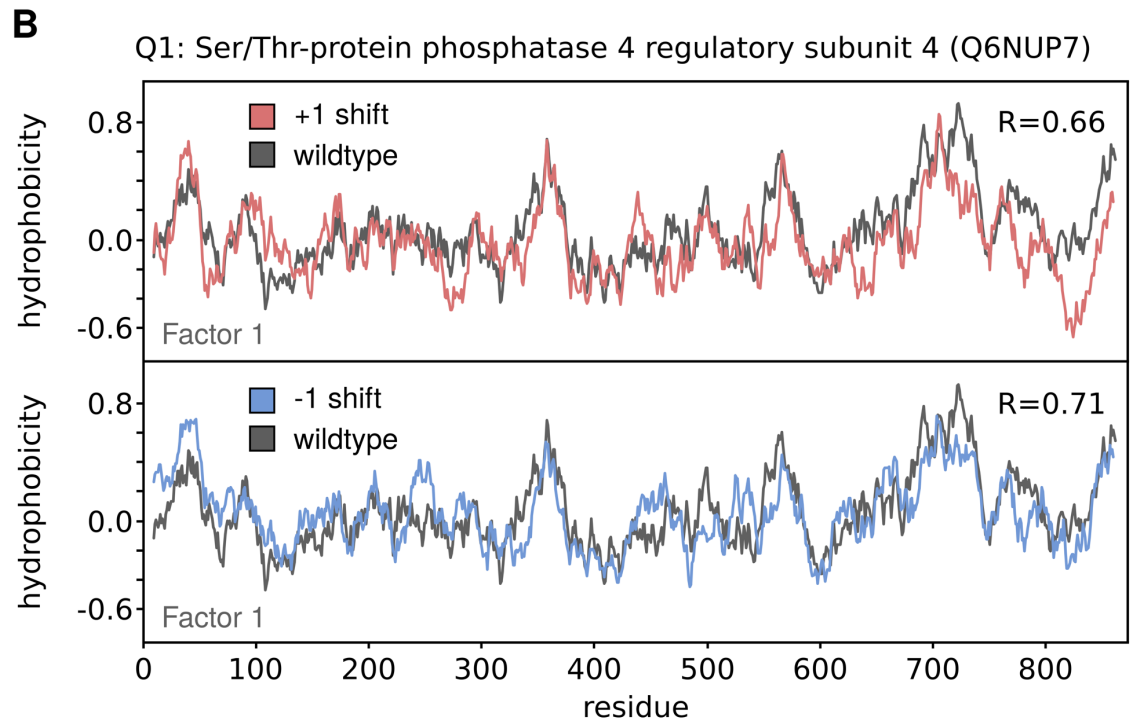
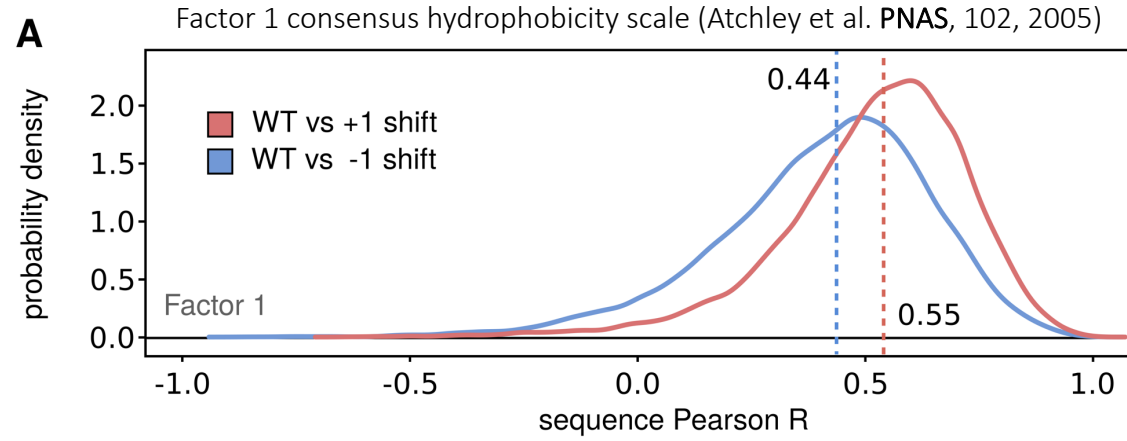


Hydrophobicity is significantly preserved upon frameshifting at the genetic code level

# Sequence hydrophobicity profiles are robust against frameshifting

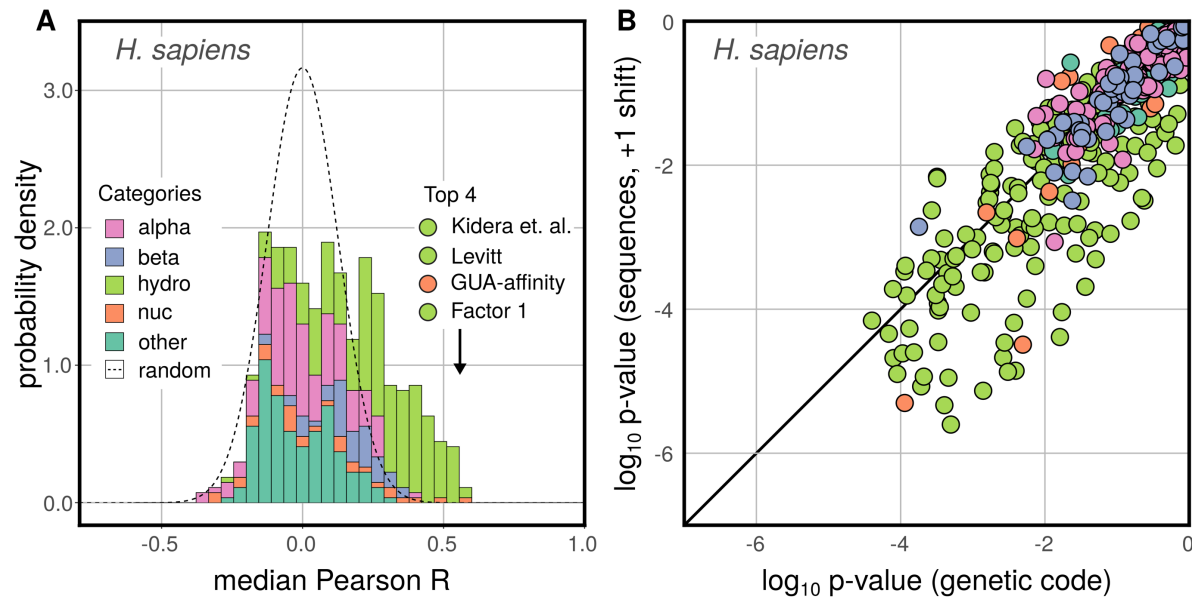


# Sequence hydrophobicity profiles are robust against frameshifting

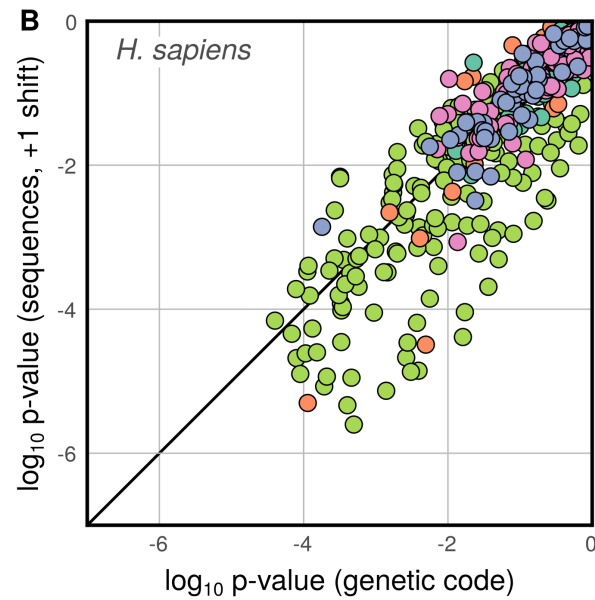
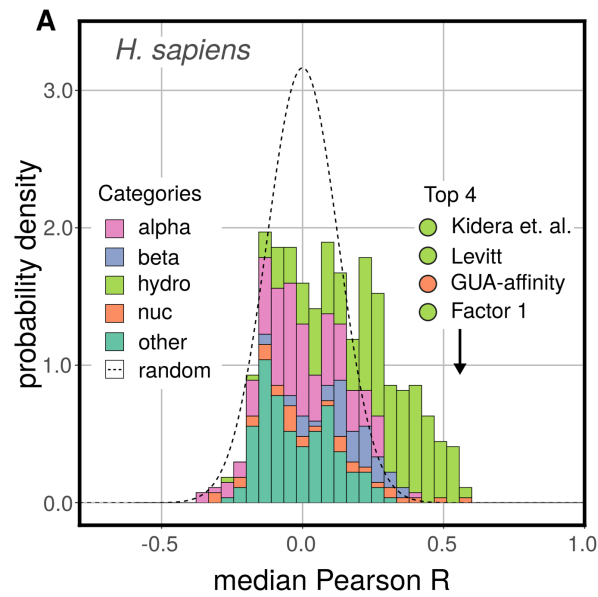


1<sup>st</sup> quartile

## Sequence hydrophobicity profiles are robust against frameshifting



# Sequence hydrophobicity profiles are robust against frameshifting



*top 25%*

**C** GO enrichment for Factor 1

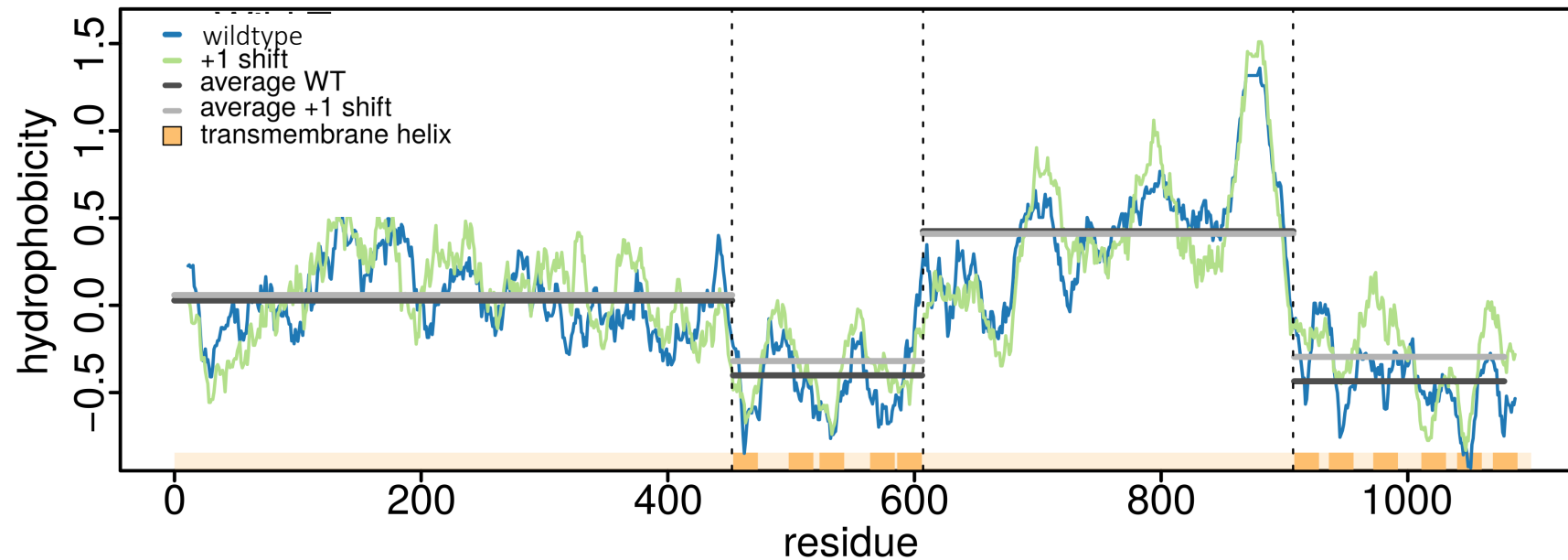
	p-value
<b>+1 shift</b>	
integral component of membrane	$1.1 \times 10^{-26}$
keratin filament	$2.2 \times 10^{-4}$
endoplasmatic reticulum	$2.8 \times 10^{-3}$
nucleolus	$3.8 \times 10^{-3}$
<b>-1 shift</b>	
nucleolus	$2.0 \times 10^{-9}$
nucleoplasm	$1.7 \times 10^{-7}$
ribonucleoprotein complex	$5.8 \times 10^{-5}$
chromosome	$1.4 \times 10^{-4}$
intracellular non-membrane organelle	$2.2 \times 10^{-3}$

## An example of a membrane protein

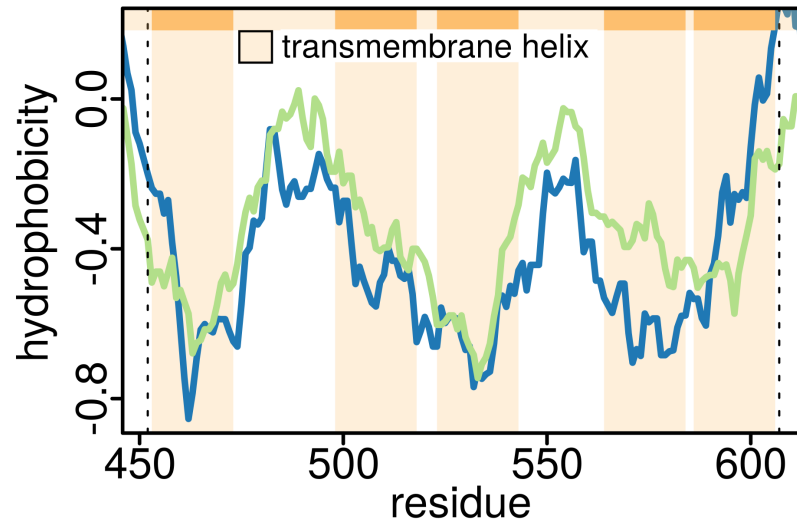
sodium/potassium/calcium exchanger 1 (O60721)

**wt** 1 MGKLIRMGPQERWLLRTKRLHWSRLLFLLGMLIIGSTYQHLRRPRGLSSL 50  
**+1** 1 WGN-SGWGRKRGGYSGQSGFIGVASSSYWEC-SSVLLISTLGDPGAFPHC 50

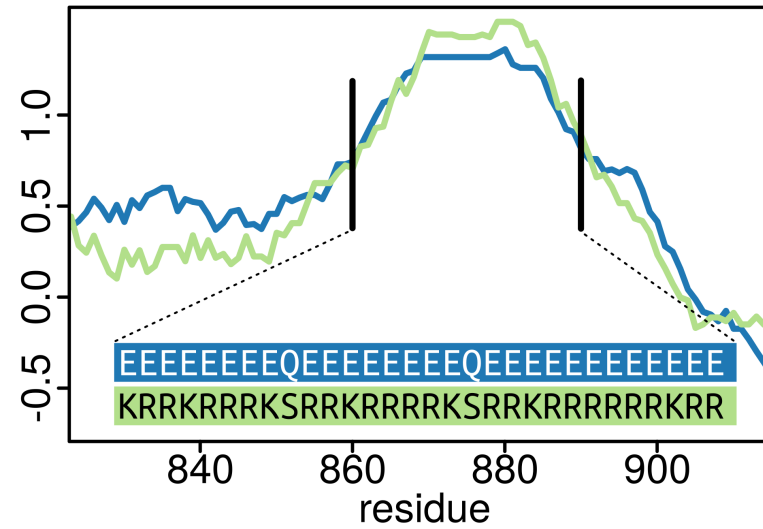
.....



sodium/potassium/calcium exchanger 1 (O60721)

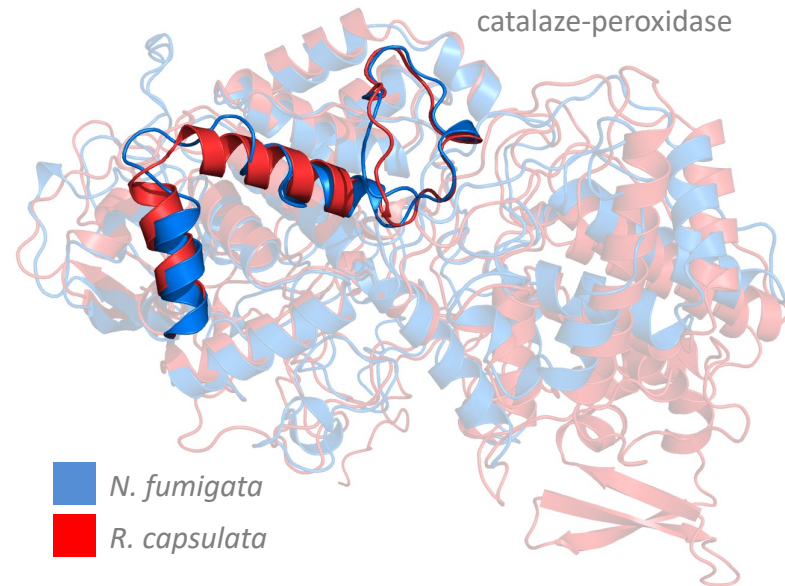


retained molecular topology



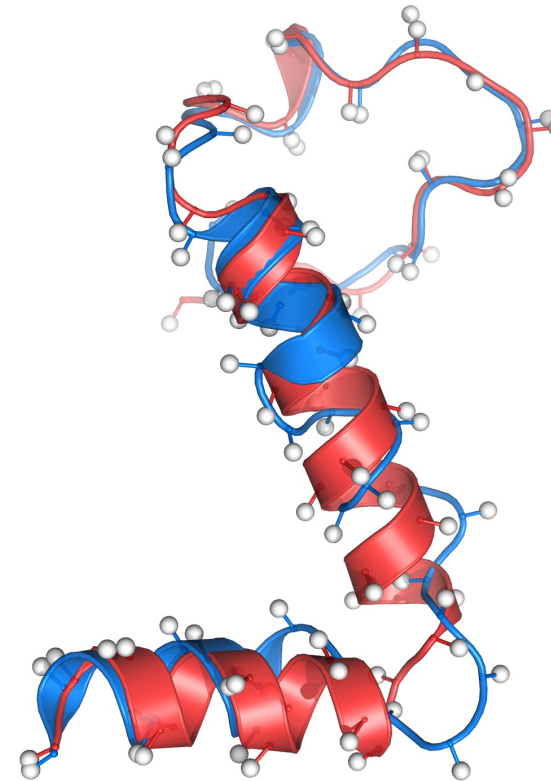
retained hydrophobicity, but  
charge inversion

## Frameshifted regions could be structurally related



```
PAYEKIARRFLEHPDQFADAFARAWFKLTHRDMGPRARYLGPEVPSEVLI  
HHHHHHHHHHH-HHHHHHHHHHHHHHHHTTS-SGGG-BSTT--SS--G  
HHHHHHHHHH--SS--TTHHHSHHHHHHHHT---GGG--STT--S-SSS  
PSTMRSVRSSWPPIRPPSTTLRAPGSSCCCIATWGRRRATSAPMCPPRIWS
```

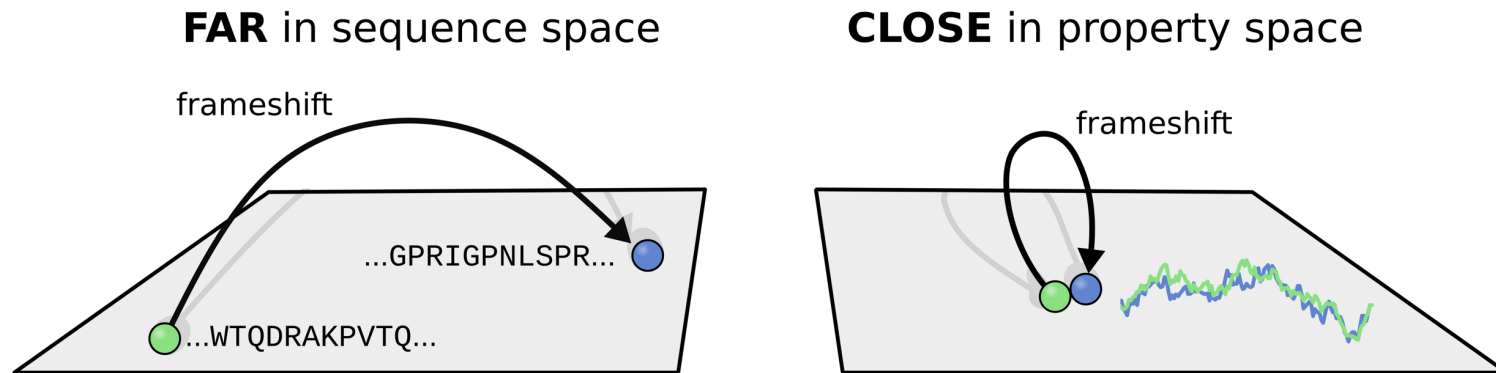
sequence identity = 12%



backbone RMSD = 2.8 Å

## Conclusion 1

- different protein sequence properties e.g. hydrophobicity, structural disorder and nucleobase affinity are significantly robust against frameshifting
- implications: overprinted genes; natural frameshifts; stop-codon reassignment; gain-of-function
- a potentially powerful way for evolution of novel sequences from optimized starting points

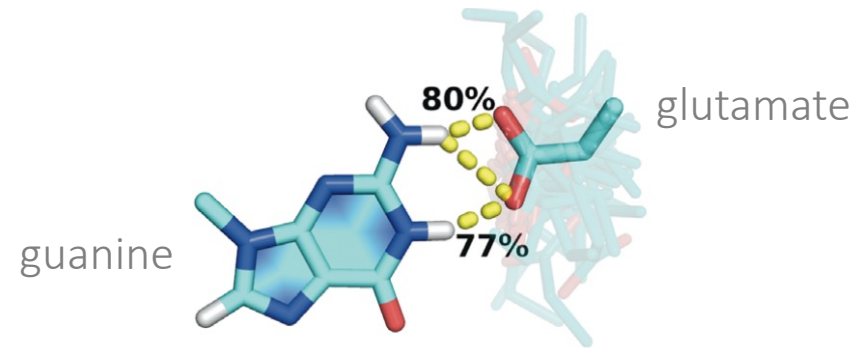


## Stereochemical hypothesis of the origin of the genetic code

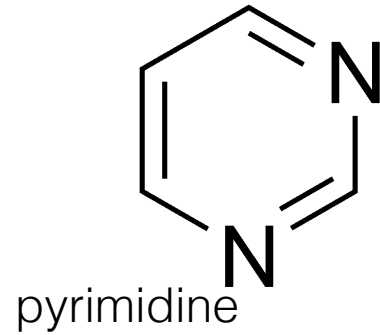
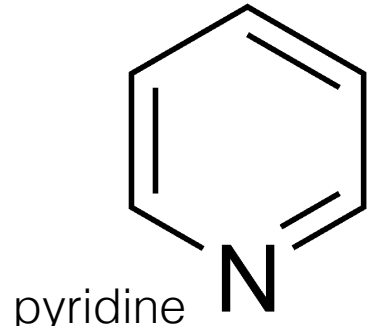
Gamow, Woese, Yarus and others ...

		Second letter				
		U	C	A	G	
First letter	U	UUU } Phe	UCU } Ser	UAU } Tyr	UGU } Cys	Third letter
		UUC } Leu	UCC } Ser	UAC } Tyr	UGC } Cys	
		UUA } Leu	UCA } Ser	UAA Stop	UGA Stop	
		UUG } Leu	UCG } Ser	UAG Stop	UGG Trp	
	C	CUU } Leu	CCU } Pro	CAU } His	CGU } Arg	
		CUC } Leu	CCC } Pro	CAC } His	CGC } Arg	
		CUA } Leu	CCA } Pro	CAA } Gln	CGA } Arg	
		CUG } Leu	CCG } Pro	CAG } Gln	CGG } Arg	
	A	AUU } Ile	ACU } Thr	AAU } Asn	AGU } Ser	
		AUC } Ile	ACC } Thr	AAC } Asn	AGC } Ser	
		AUA } Met	ACA } Thr	AAA } Lys	AGA } Arg	
		AUG } Met	ACG } Thr	AAG } Lys	AGG } Arg	
	G	GUU } Val	GCU } Ala	GAU } Asp	GGU } Gly	
		GUC } Val	GCC } Ala	GAC } Asp	GGC } Gly	
		GUA } Val	GCA } Ala	GAA } Glu	GGA } Gly	
		GUG } Val	GCG } Ala	GAG } Glu	GGG } Gly	

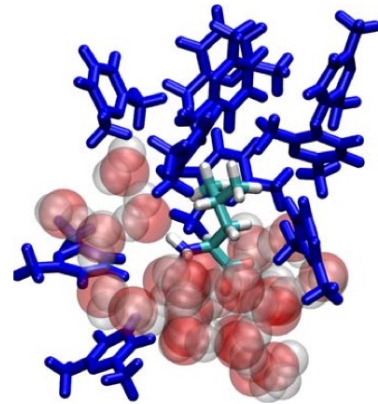
**universal genetic code**  
↔  
**nucleotide/amino-acid affinities**



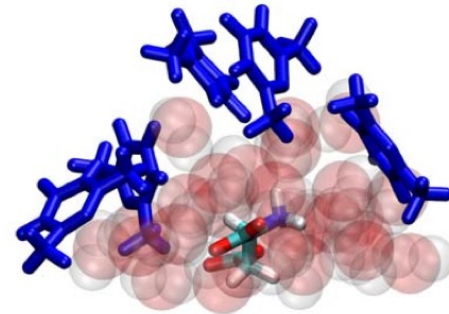
# What is the affinity of amino acids for pyrimidines?



$$\Delta\Delta G \sim -\log(C_{\text{DMP}}/C_{\text{water}})$$

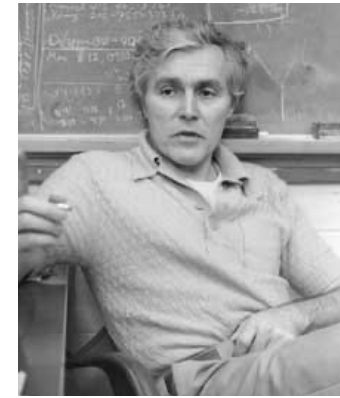


LEU in DMP:H<sub>2</sub>O



ASP in DMP:H<sub>2</sub>O

\*DMP: 2,6-dimethylpyridine

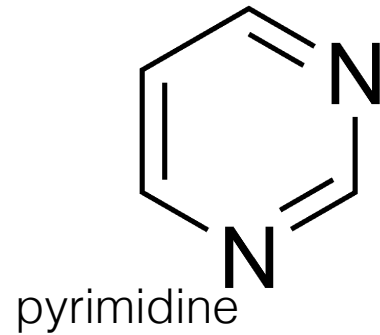
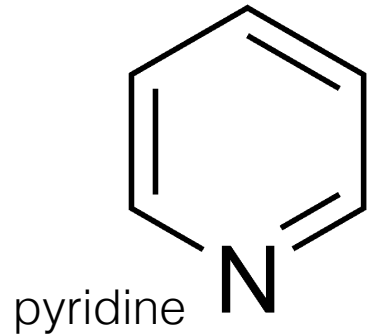


Carl Woese

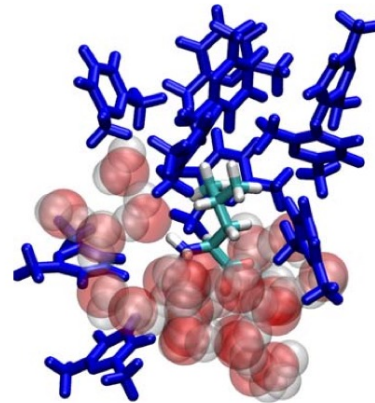
# What is the affinity of amino acids for pyrimidines?



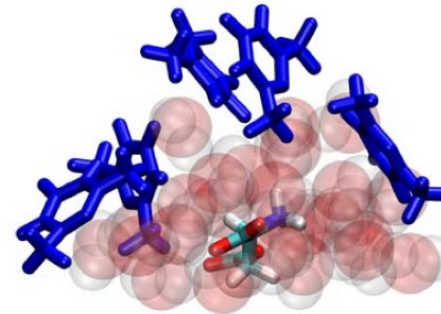
Carl Woese



$$\Delta\Delta G \sim -\log(C_{\text{DMP}}/C_{\text{water}})$$



LEU in DMP:H<sub>2</sub>O



ASP in DMP:H<sub>2</sub>O

\*DMP: 2,6-dimethylpyridine

<b>C</b>	<b>L</b>	<b>F</b>	<b>W</b>	<b>I</b>	<b>M</b>	<b>P</b>	<b>V</b>	<b>T</b>	<b>A</b>	<b>S</b>	<b>Y</b>	<b>H</b>	<b>R</b>	<b>Q</b>	<b>G</b>	<b>N</b>	<b>K</b>	<b>D</b>	<b>E</b>
4.3	4.4	4.5	4.9	5.0	5.0	6.1	6.2	6.2	6.5	7.5	7.7	7.9	8.6	8.9	9.0	9.6	10.2	12.2	13.6

a. u.

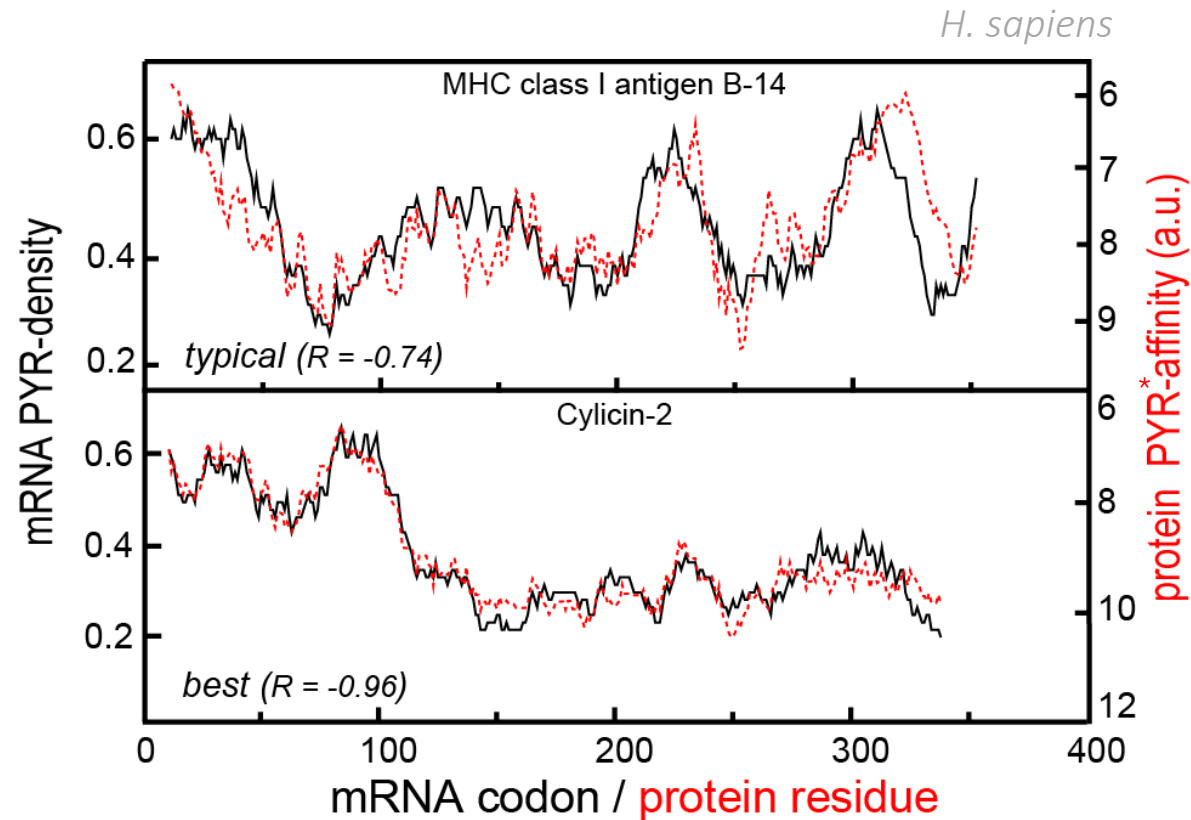
PYR-density = 0.7

mRNA ...UGGGCC **CUGGGCUUCUACCCU** GCGGAG...

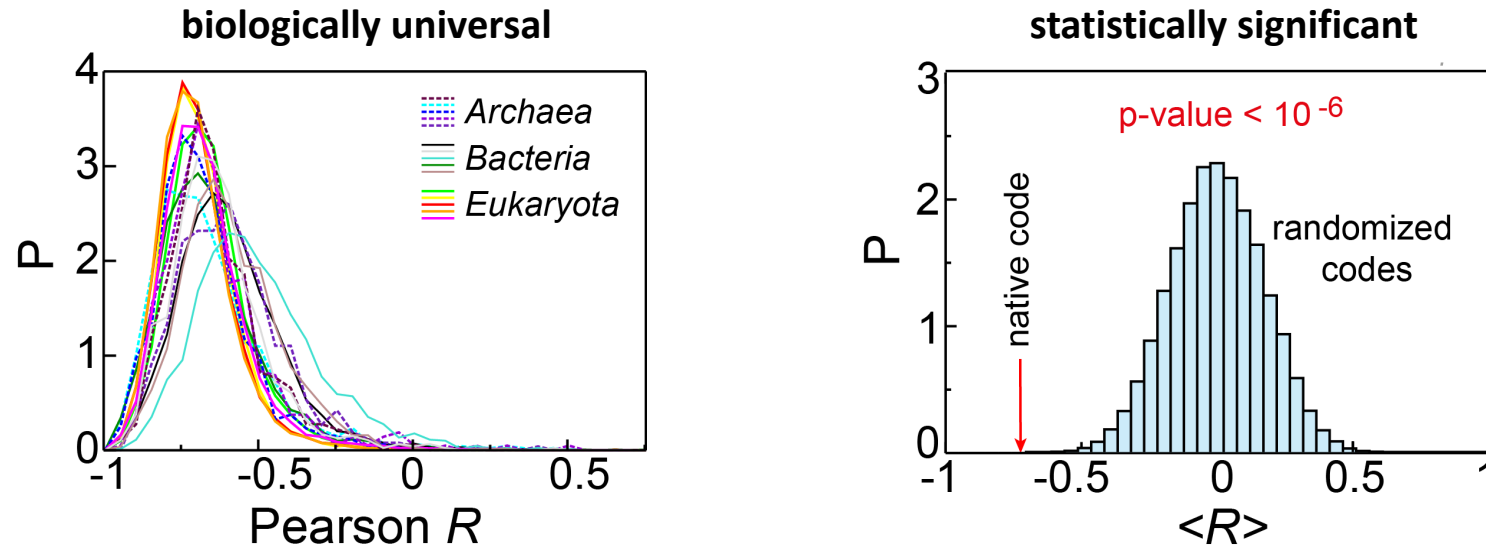
protein ...Trp - Ala - **Leu - Gly - Phe - Tyr - Pro** - Ala - Glu...

PYR\*affinity = 6.3

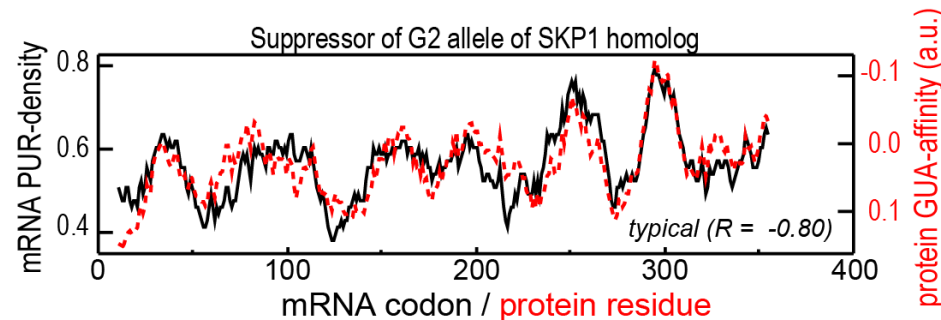
mRNA ...UGGGCC **CUGGGCUUCUACCCU** GCGGAG...  
 protein ...Trp - Ala - **Leu - Gly - Phe - Tyr - Pro** - Ala - Glu...  
 PYR-density = 0.7  
 PYR\*-affinity = 6.3



# A robust and reproducible result of general validity



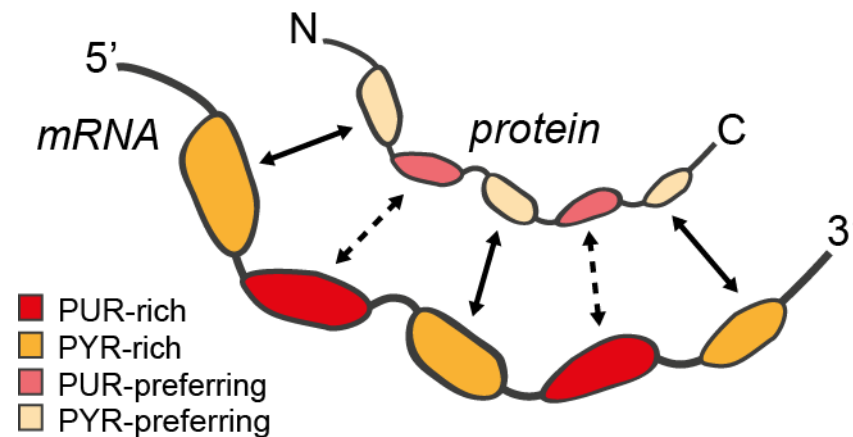
## confirmed with multiple affinity scales



Hlevnjak, Polyansky & Zagrovic, **NAR**, 40, 2012; Polyansky & Zagrovic, **NAR**, 41, 2013; de Ruiter & Zagrovic, **NAR**, 43, 2015; Hlevnjak & Zagrovic, **NAR**, 43, 2015; Bartonek & Zagrovic, **PLOS CB**, 13, 2017; Zagrovic, Bartonek & Polyansky, **FEBS Letters**, 592, 2018; Bartonek, Braun & Zagrovic, **PNAS**, 117, 2020; Kapral, Farnhammer, Zhao, Lu & Zagrovic, **NAR**, 50, 2022; Zagrovic, Adlhart & Kapral, **Ann Rev Biophys**, 52, 2023; Adlhart, Hoffmann, Polyansky & Zagrovic, **PNAS**, 122, 2025

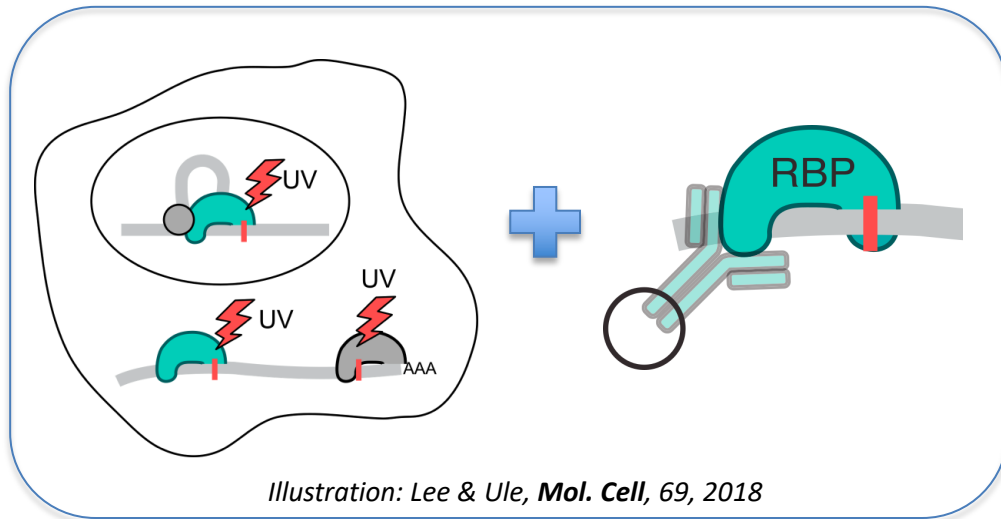
## CENTRAL HYPOTHESIS

mRNA CDS and their autogenous proteins are **complementary to each other and bind in a co-aligned fashion**, especially if unstructured, reflecting the driving forces behind the origin of the genetic code. Complementarity is negatively regulated by the mRNA adenine content.



# Comparison with experiment

- crosslinking and immunoprecipitation (CLIP-seq) informs on RNA binding targets of individual proteins
- 341 different RNA-binding proteins (RBPs) tested so far



POSTAR3 Modules ▾ Submit Help About Contact

### RBP binding sites of TAF15

Export data to CSV file Search:

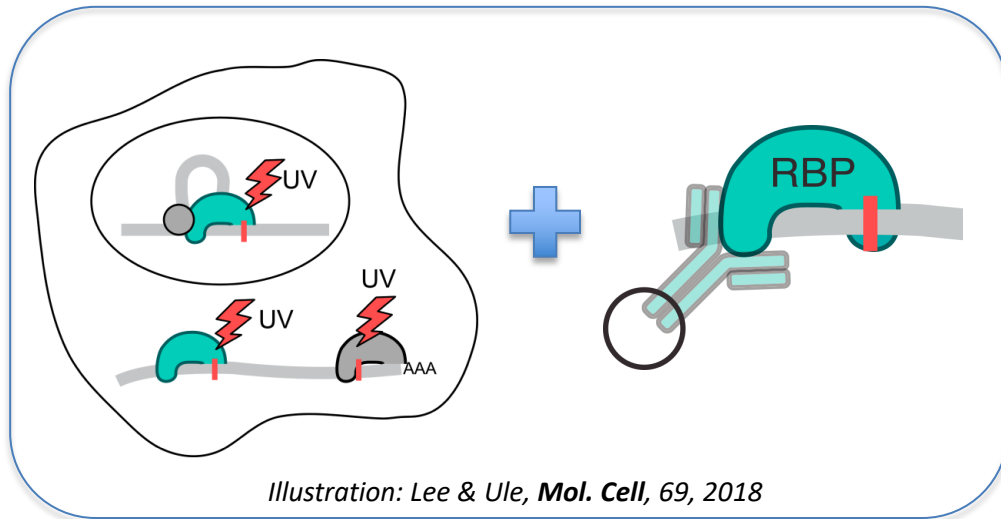
Target gene symbol	Target gene ID	Target gene type	Protocol	Target gene exp. level	Binding site records
DLEU1	ENSG00000176124	processed_transcript	eCLIP	<a href="#">display</a>	154
CASC19	ENSG00000254166	processed_transcript	eCLIP	<a href="#">display</a>	107
SNHG1	ENSG00000255717	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	98
GAS5	ENSG00000234741	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	58
SNHG14	ENSG00000224078	processed_transcript	HITS-CLIP,Piranha_0.01	<a href="#">display</a>	53
PROX1-AS1	ENSG00000230461	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	38
GAS5	ENSG00000234741	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	36
SNHG1	ENSG00000255717	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	35
LRRC75A-AS1	ENSG00000175061	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	34
DLEU1	ENSG00000176124	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	25

Showing 1 to 10 of 143 entries (filtered from 13,176 total entries) Previous  2 3 4 5 ... 15 Next

Jump to  page

# Comparison with experiment

- crosslinking and immunoprecipitation (CLIP-seq) informs on RNA binding targets of individual proteins
- 341 different RNA-binding proteins (RBPs) tested so far



## Two key predictions:

1. autogenous mRNA/protein interactions are enriched over background
2. autogenous interactions occur preferentially in the coding sequence (CDS)

POSTAR3 Modules ▾ Submit Help About Contact

### RBP binding sites of TAF15

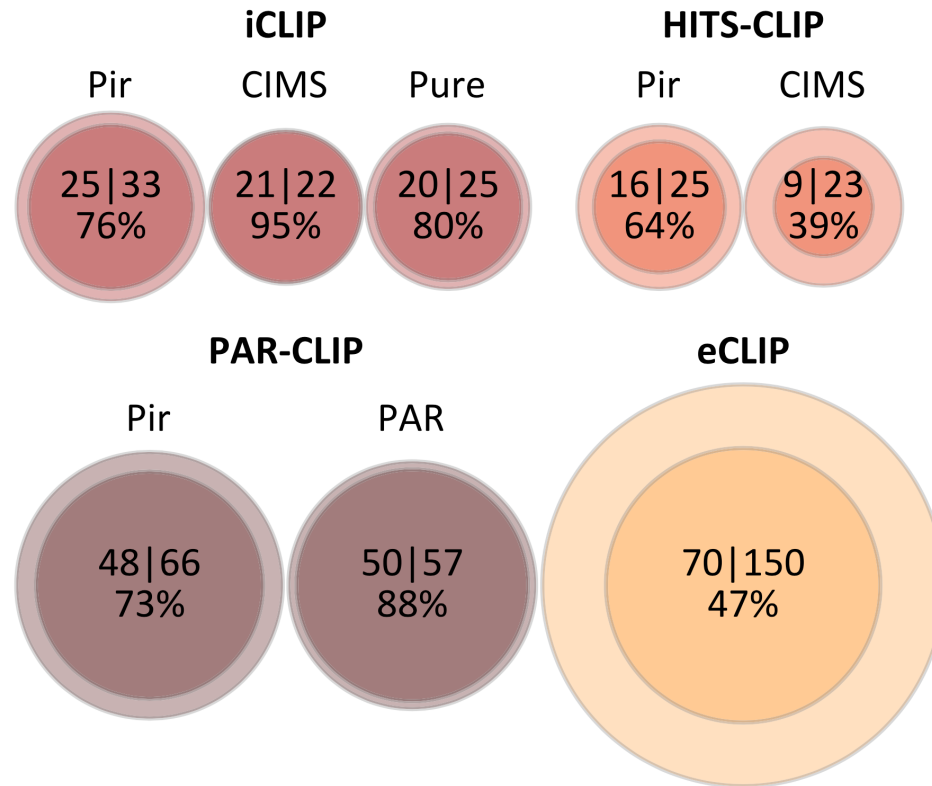
Export data to CSV file Search:

Target gene symbol	Target gene ID	Target gene type	Protocol	Target gene exp. level	Binding site records
DLEU1	ENSG00000176124	processed_transcript	eCLIP	<a href="#">display</a>	154
CASC19	ENSG00000254166	processed_transcript	eCLIP	<a href="#">display</a>	107
SNHG1	ENSG00000255717	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	98
GAS5	ENSG00000234741	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	58
SNHG14	ENSG00000224078	processed_transcript	HITS-CLIP,Piranha_0.01	<a href="#">display</a>	53
PROX1-AS1	ENSG00000230461	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	38
GAS5	ENSG00000234741	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	36
SNHG1	ENSG00000255717	processed_transcript	PAR-CLIP,Piranha_0.01	<a href="#">display</a>	35
LRRC75A-AS1	ENSG00000175061	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	34
DLEU1	ENSG00000176124	processed_transcript	PAR-CLIP,PARalyzer	<a href="#">display</a>	25

Showing 1 to 10 of 143 entries (filtered from 13,176 total entries) Previous  2 3 4 5 ... 15 Next

Jump to  page

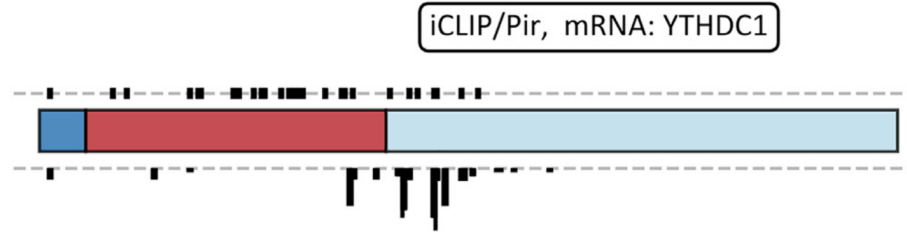
# Autogenous mRNA-protein interactions are seen frequently in CLIP-seq experiments



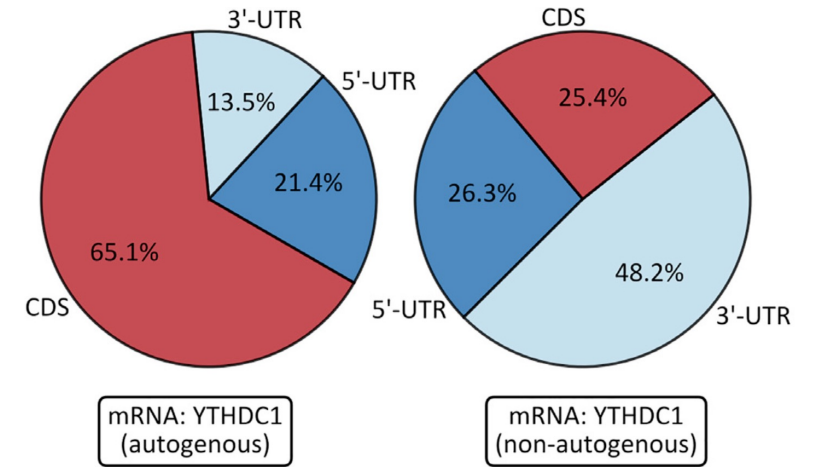
**auto binding:**  
**230/341 (67%)**

peak callers: Pir- Piranha, PAR – Paralyzer, CIMS, CLIPPER

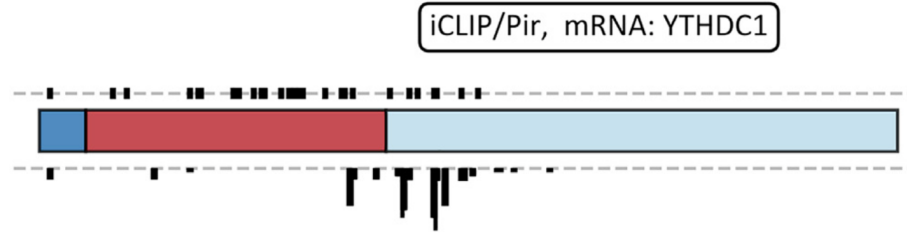
## Autogenous interactions occur preferentially in the CDS



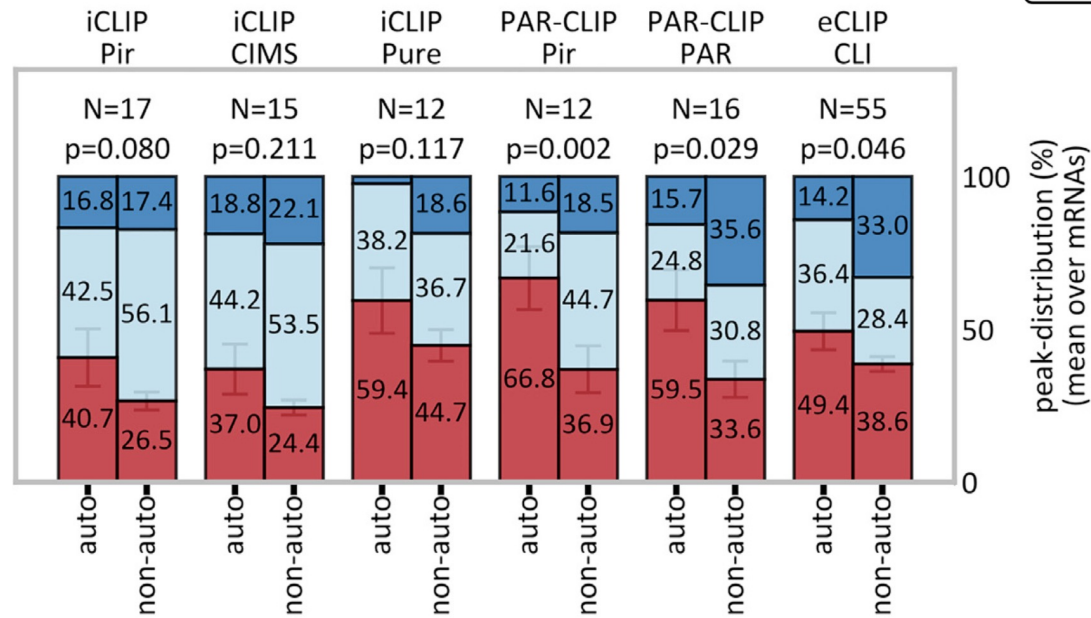
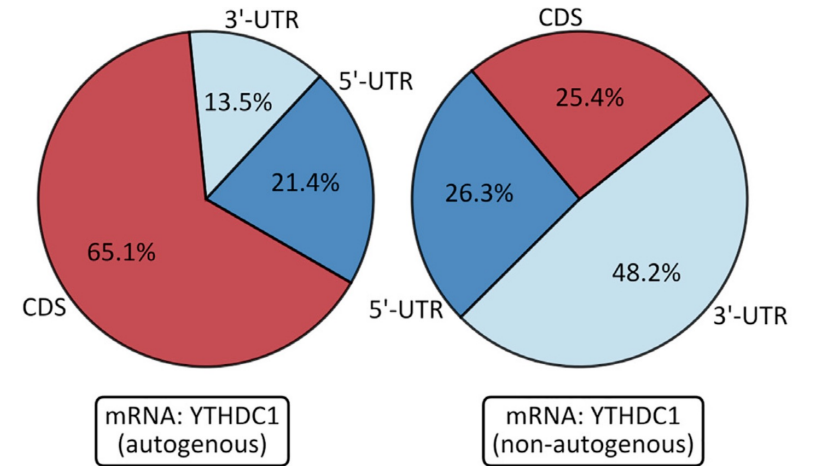
a .. autogenous peaks  
n .. non-autogenous peaks



# Autogenous interactions occur preferentially in the CDS



a .. autogenous peaks  
n .. non-autogenous peaks

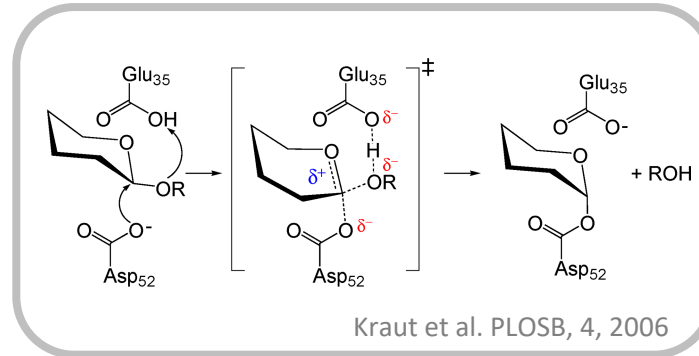


# Physicochemical complementarity

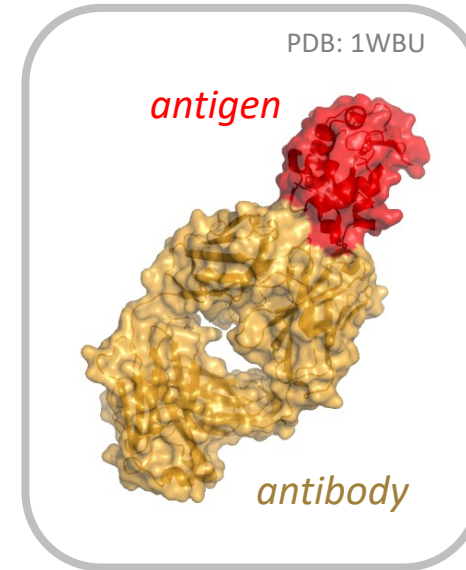
one of the most powerful paradigms in biology



DNA replication



enzymatic catalysis



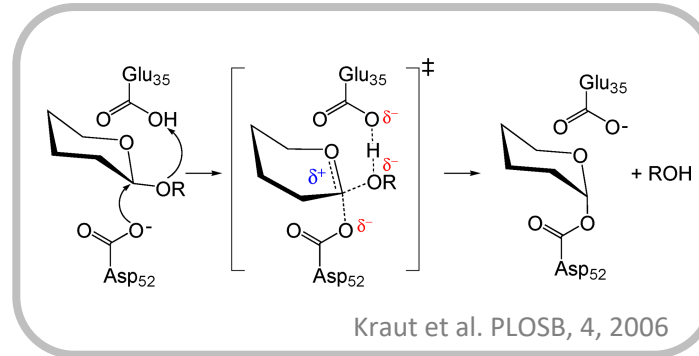
immune response

# Physicochemical complementarity

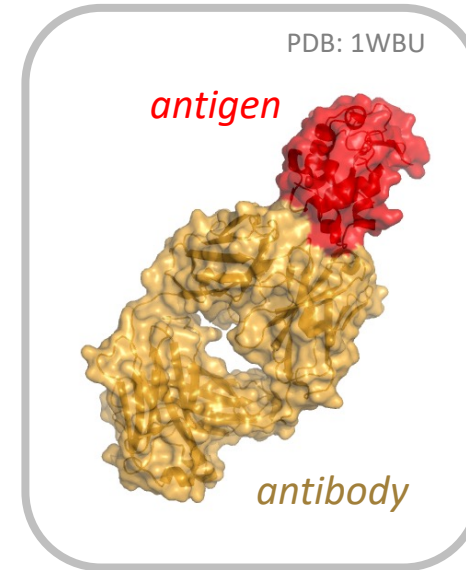
one of the most powerful paradigms in biology



DNA replication

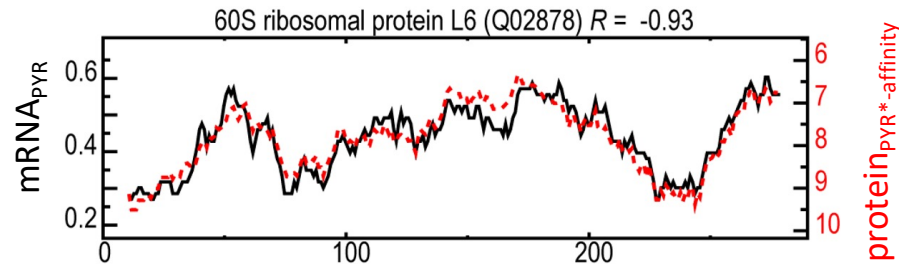


enzymatic catalysis



immune response

a novel complementarity?



## Conclusion 2: coding and binding as two faces of the same coin?

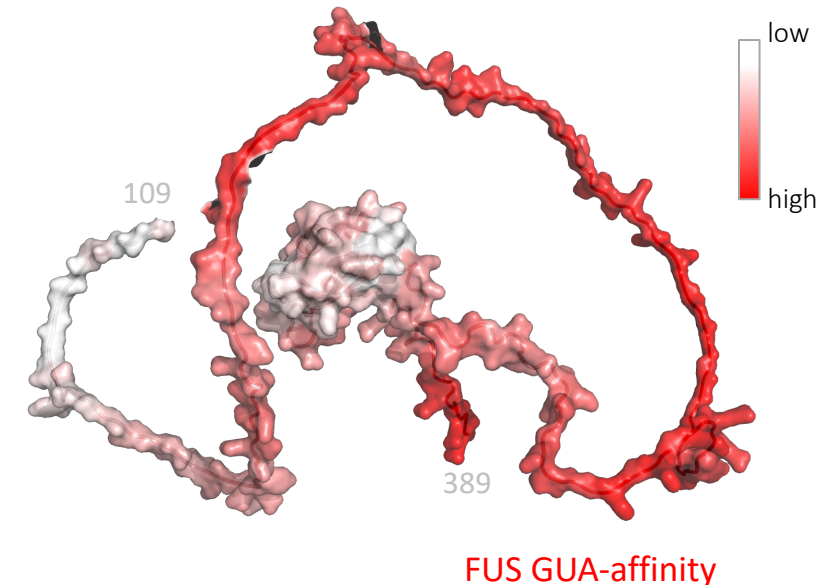
- evidence of complementary interactions between mRNAs and own proteins in unstructured state
- proteins interact with RNAs that are compositionally similar to their own mRNA and *vice versa*
- genetic code as a key for understanding large-structure of the cellular RNA-protein interactome

### OPEN CHALLENGES

geometry;  $\Delta G$ s; functional significance; experimental testing

## Towards physicochemical bioinformatics

- **Analyzing biomolecular sequences from a physicochemical perspective:** a rich source of unexpected patterns
- **Particularly relevant in an unstructured context:** sequences as physicochemical objects
- **Advances required when it comes to:**
  - robust algorithms for sequence profile alignment
  - new measures for profile comparison
  - development of fast methods for detection of historical frameshifts
  - experimental testing of complementarity hypothesis



# Thanks

## Present members

- **Thomas Kapral**
- **Anton A. Polyansky**
- Fran Miočić Stošić
- Felix Fischer
- Hannah Dydziul
- Christian Stelmach

## Former members

- **Marlene Adlhart**
- **Mario Hlevnjak**
- **Theres Friesacher**
- **Daniel Hoffmann**
- **Anita de Ruiter**
- **Florian Poetsch**
- **Matea Hajnic**
- **Megan Hoogmoed**
- **Fiona Farnhammer**

& other members of the lab

## Collaborators

- Renee Schroeder, Max Perutz Labs & University of Vienna
- Dea Slade, Max Perutz Labs & University of Vienna
- Thomas Leonard, Max Perutz Labs & University of Vienna
- John D. Sutherland, LMB Cambridge
- Zoya Ignatova, University of Hamburg
- Markus Zweckstetter, Max Planck Göttingen
- John Zhi Lu, Tsinghua University, Beijing
- Vanja Nagy, Medical University of Vienna



## Funding



**FWF**

