

ortho_seqs

Miles Woollacott



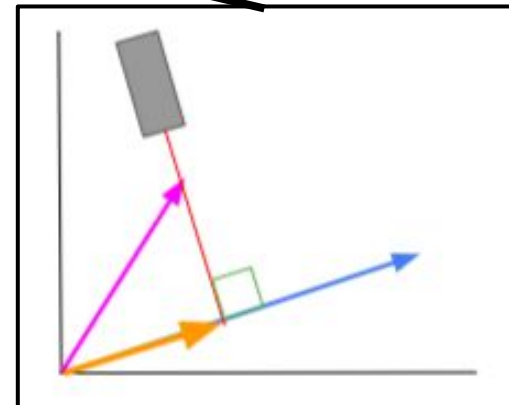
What is ortho_seqs? What does ortho_seqs do?

Results



ic50_seq.txt	ic50_phi.txt
ARDLDVYGLDV	42.0
ARDYGDLYFDY	140.0
ARDFGDFYFDY	25.0
ATESWVYSGSYSSGAFDI	1300.0
ARGPRYSGTIFYDY	480.0
ARGPRYSGTHFDY	160.0
ARDPYGYSSYWDGQGGH	50.0
ARGGYCSGSCYVQDLIYYSGLDV	39.0
ARDRRRRYCTNGVCYRPEEDY	4900.0
ARDPFPGAVAGTGYLQY	590.0
AKSSGSYYYYYGMDV	44.0
ARLHCGGDCYLDY	21.0
AKGSGSGSYPNYYYYGMDV	19.0
AKANKYSSSEFDF	66.0
AVYYYYDSGSGPWFDY	15.0
ARDFRYCSSTRCYFWFDY	12.0
ARWYDSTGSDY	25.0
AKDGSGSYGNYFDY	13.0
ALRNQWDLVY	1500.0
ARDLVVYGMDV	24.0
ARDPIRNGMDVW	18100.0
ARDAMSYGMDV	10.0
ARDAAVYIDV	16.0
ARDLISRGMDV	33.0
ARDRVVYGMDV	12.0
ARDLVSYGMDV	553.0
ARDAQNYGMDV	15.0
ARDRLVSDYW	22.0
ARPQGGSSWYRDYYGMDV	4800.0
ARDLSVRGGMDV	1900.0

ortho_seqs is a command line interface (CLI) written in Python that applies the **tensor-based orthogonal polynomial method*** to sequence & phenotype data.



```
ortho_seq orthogonal-polynomial ../data/cov2/big_mabs_ic50s/ic50_seq.txt
--pheno_file ../data/cov2/big_mabs_ic50s/ic50_phi.txt --molecule protein
--poly_order first --out_dir ../data_sm/home/miles/results/CDRH3/cov2_unique_ic50
```

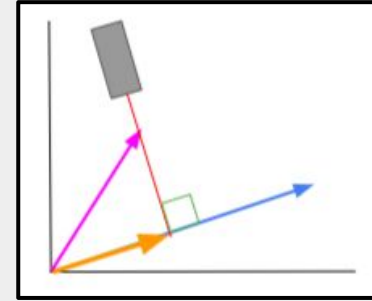
Sequence	Vectors	ϕ
GGATGA...	$\begin{bmatrix} G \\ 0 \\ 1 \\ 0 \end{bmatrix}$	0.2
AAGTGA...	$\begin{bmatrix} A \\ 0 \\ 0 \\ 0 \end{bmatrix}$	
...	$\begin{bmatrix} A \\ 1 \\ 0 \\ 0 \end{bmatrix}$	1.3

...

*See [Nafees et al., 2020](#) for method application and [Rice, 2020](#) for method background

Why is ortho_seqs important?

ortho_seqs converts each site in a sequence into a vector, which allows for mathematical computations and analyses.



Matrix of Covariances

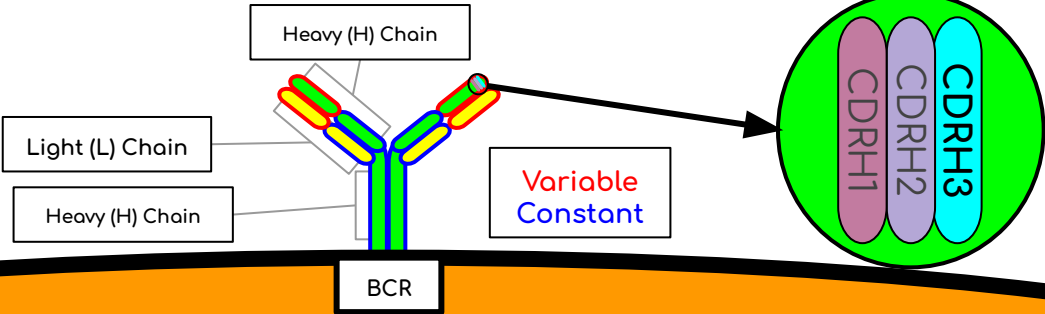
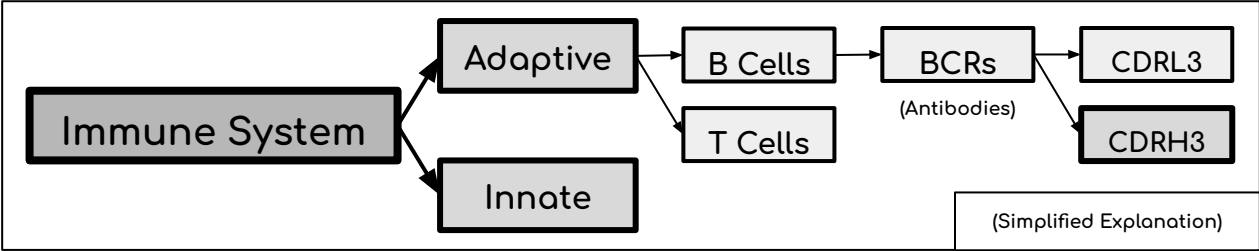
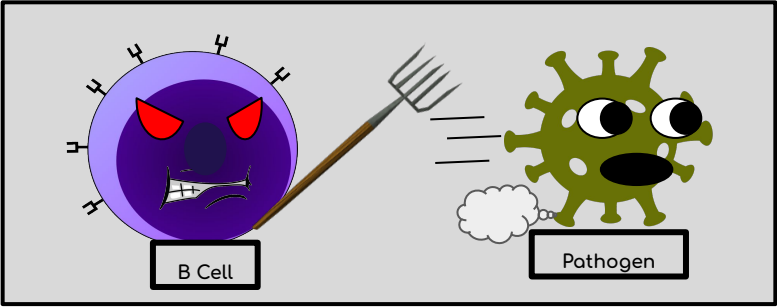
Covariance measures the relationship between two items.

Having a positive covariance for a pair of site \rightarrow values of the first item increase with the second (and vice versa).

Regressions of Phenotype Onto the X Order Conditional Polynomial (rFonXD)

rFon1D: measures the impact of having a given amino acid/nucleotide at that site. This impact can be positive or negative.

Regression onto second order conditional polynomial (rFon2D): measures the impact of having a given nucleotide *pair* at pairs of sites (not supported for proteins yet).



The cov2 dataset is a conglomeration of antibodies that neutralize SARS-CoV-2, along with the sequences, IC50, and EC50 values (where provided).

IC50 = Half Maximal Inhibitory Concentration
EC50 = Half Maximal Effective Concentration

SHARE

RESEARCH ARTICLE

f

t

in

u

Isolation of potent SARS-CoV-2 neutralizing antibodies and protection from disease in a small animal model

Thomas F. Rogers^{1,2,*}

Fangzhu Zhao^{1,3,4,*}

Deli Huang^{1,*}

Nathan Beutler^{1,*}

Alison Burns^{1,3,4}

Wan-ting He^{1,3}...

See all authors and affiliations

Science

21 Aug 2020

Vol. 369, Issue 6506, pp. 956-963

DOI: 10.1126/science.abc7520

mAb ID	VH Gene	% SHM (nt)	CDRH3 (AA)	mAb name	Predicted Unique Epitope	ELISA IC ₅₀ (µg/mL)		
						SARS-CoV-2 S	SARS-CoV-2 RBD	SARS-CoV-1 S
CC12.1	IGHV3-53*01	1.05	CARDLDVYGLDVW	CC12.1	RBO-A	0.017	0.042	>50
CC12.2	IGHV3-53*01	0.70	CARDYGLDYFDYW	CC12.2	RBO-A	0.021	0.14	>50
CC12.3	IGHV3-53*01	1.40	CARDFGDFYFDY	CC12.3	RBO-A	0.018	0.025	>50
CC12.4	IGHV1-2*04	1.74	CATESWVYSGSYSGAFDW	CC12.4	RBO-A	0.062	1.3	>50
CC12.5	IGHV1-2*02	4.86	CARGPRYSGTIFYDYW	CC12.5	RBO-A	0.33	1.4	>50
CC12.6	IGHV1-2*02	4.17	CARGPRYSGTIFYDYW	CC12.6	RBO-A	0.36	3.1	>50
CC12.7	IGHV1-2*02	7.99	CARGPRYSGTIFYDYW	CC12.7	RBO-A	0.091	1.2	>50
CC12.8	IGHV1-2*02	3.12	CARGPRYSGTIFYDYW	CC12.8	RBO-A	0.10	0.48	>50
CC12.9	IGHV1-2*02	3.47	CARGPRYSGTIFYDYW	CC12.9	RBO-A	20	>50	>50
CC12.10	IGHV1-2*02	4.51	CARGPRYSGTIFYDYW	CC12.10	RBO-A	0.013	0.16	>50
CC12.11	IGHV1-2*02	3.47	CARGPRYSGTIFYDYW	CC12.11	RBO-A	0.36	1.2	>50
CC12.12	IGHV1-2*02	3.47	CARGPRYSGTIFYDYW	CC12.12	RBO-A	0.15	0.67	>50
CC12.13	IGHV3-53*01	0.35	CARDPYYSWDDGGGHW	CC12.13	RBO-A	0.044	0.058	>50
CC12.14	IGHV3-21*01	3.47	CARGGYCSGSGCYVQDLIYYSGLDW	CC12.14	RBO-A	0.014	0.039	>50
CC12.15	IGHV3-48*03	1.04	CARDRRRYCTNGVCYRPEEDY	CC12.15	RBO-A	0.87	4.9	>50
CC12.16	IGHV3-33*01	1.74	CARDPFGAVAGTGYLQY	CC12.16	RBO-B	0.073	0.59	>50
CC12.17	IGHV3-30*03	1.39	CAKSGSYYYYYYGMVW	CC12.17	RBO-B	0.020	0.044	>50
CC12.18	IGHV1-46*01	0.00	CARLHCGGDCYLDY	CC12.18	RBO-B	0.017	0.021	0.018
CC12.19	IGHV3-23*04	1.04	CAKSGSGSYPNYYYYYGMVW	CC12.19	RBO-B	0.012	0.019	>50
CC12.20	IGHV3-30*03	1.04	CARDQAYDLTYLWPRYYVYGMVW	CC12.20	SPIKE-A	0.036	>50	>50
CC12.21	IGHV1-24*01	1.39	CATAFSFGVPPVW	CC12.21	SPIKE-A	0.0050	>50	>50
CC12.22	IGHV1-24*01	1.04	CATGFAGNALLTPYW	CC12.22	SPIKE-A	>50	>50	>50
CC12.23	IGHV4-39*01	0.69	CARGGDCSTTSCAYDYW	CC12.23	SPIKE-A	0.014	>50	>50
CC12.24	IGHV3-30*03	1.04	CAKDRTNYYVGMVW	CC12.24	ND	0.013	>50	>50

	A	B	C	D	E	F	G
1	Status	name	author	phenotype_ty	ic50_ng	ec50_ng	cdr-h3
10	Verified	CC12.1	Rogers	IC50		42	ARLDVYGLDV
11	Verified	CC12.2	Rogers	IC50		140	ARDYGLDYFDY
12	Verified	CC12.3	Rogers	IC50		25	ARDFGDFYFDY
13	Verified	CC12.4	Rogers	IC50		1300	ATESWVYSGSYSGAFDI
14	Verified	CC12.5	Rogers	IC50		1400	ARGPRYSGTIFYDY
15	Verified	CC12.6	Rogers	IC50		3100	ARGPRYSGTIFYDY
16	Verified	CC12.7	Rogers	IC50		1200	ARGPRYSGTIFYDY
17	Verified	CC12.8	Rogers	IC50		480	ARGPRYSGTIFYDY
18	Verified	CC12.10	Rogers	IC50		160	ARGPRYSGTIFYDY
19	Verified	CC12.11	Rogers	IC50		1200	ARGPRYSGTIFYDY
20	Verified	CC12.12	Rogers	IC50		670	ARGPRYSGTIFYDY

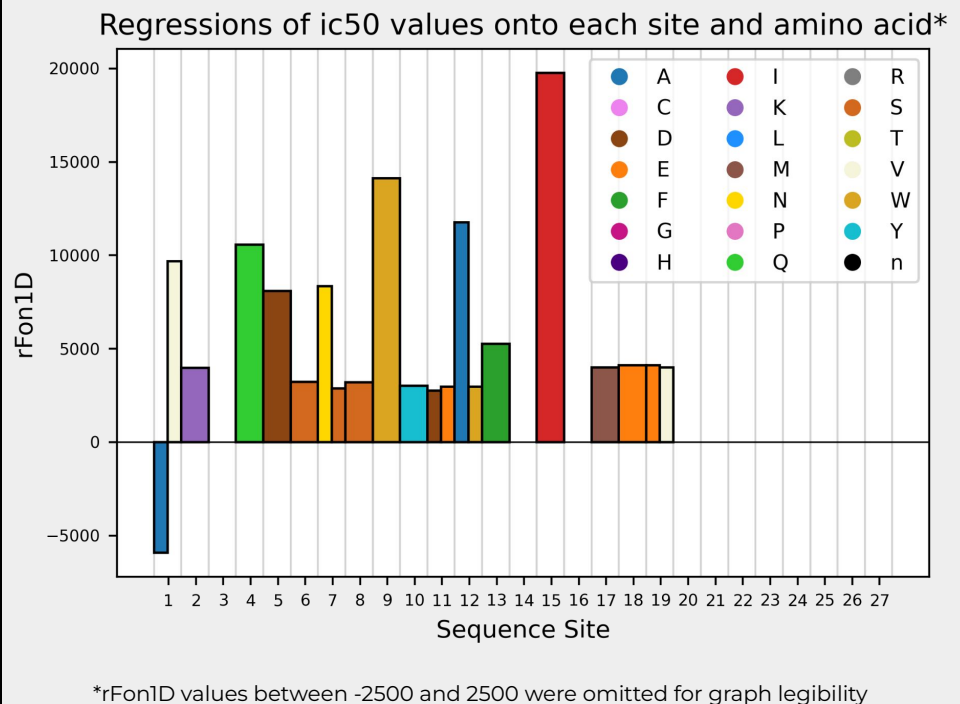
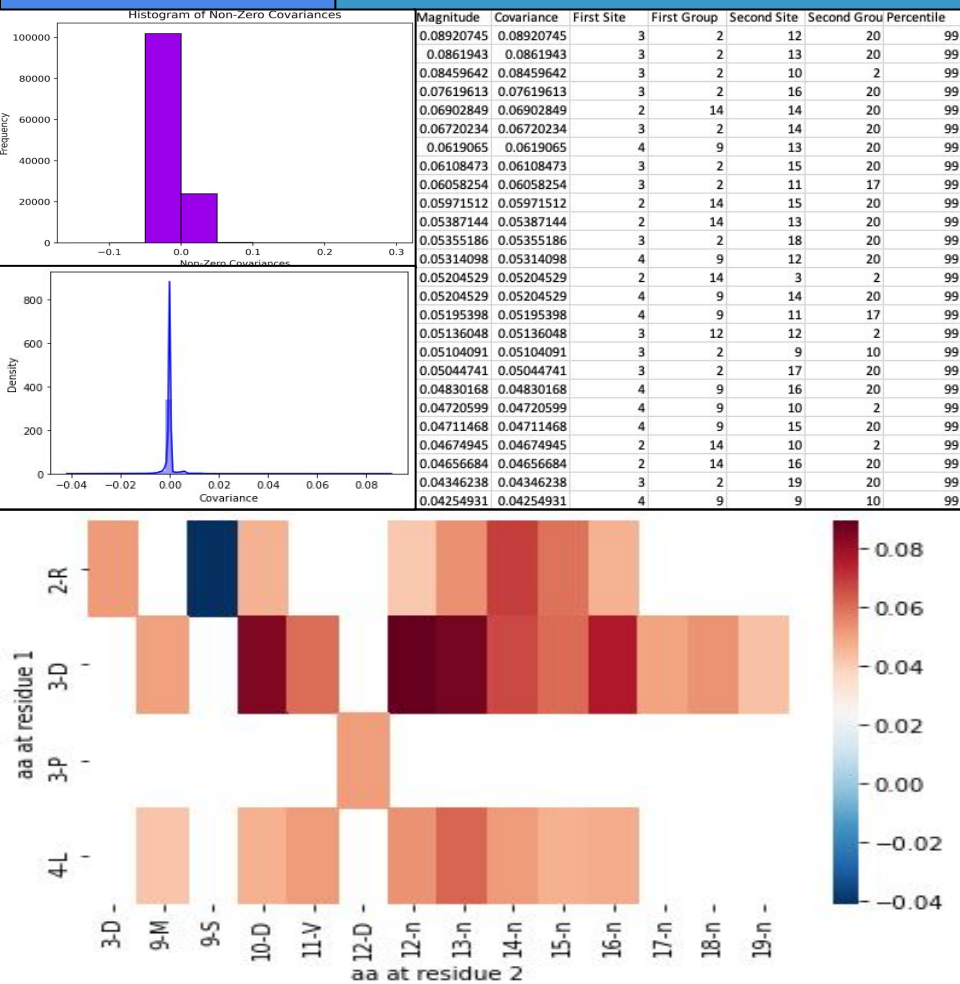
```
ic50_seq.txt ic50_phi.txt
ARLDVYGLDV 42.0
ARDYGLDYFDY 140.0
ARDFGDFYFDY 25.0
ATESWVYSGSYSGAFDI 1300.0
ARGPRYSGTIFYDY 480.0
ARGPRYSGTIFYDY 160.0
ARGPRYSGTIFYDY 310.0
ARDPYGSSIDWQGGH 58.0
ARGGYCSGSGCYVQDLIYYSGLDV 39.0
ARDRRRYCTNGVCYRPEEDY 4900.0
ARDPFGAVAGTGYLQY 590.0
AKSSGSSYYYYYGMV 44.0
ARLHCGGDCYLDY 21.0
AKGSGSGSYPNYYYYYGMV 19.0
AKANKYSSSEFDF 66.0
AVYYDSGSGPWFD 15.0
ARDFRYCSSTRCYFWFD 12.0
ARWYDSTGSDY 25.0
AKDGYSYGNFYDY 13.0
```

CLI Input:

```
ortho_seq orthogonal-polynomial ../data/cov2/big_mabs_ic50s/ic50_seq.txt --pheno_file
../data/cov2/big_mabs_ic50s/ic50_phi.txt --molecule protein --poly_order first --out_dir
/data_sm/home/miles/results/CDRH3/cov2_unique_ic50
```

cov2 Dataset

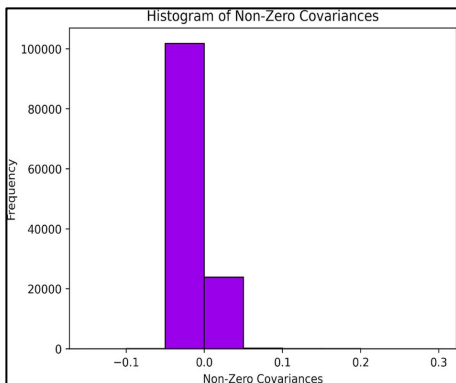
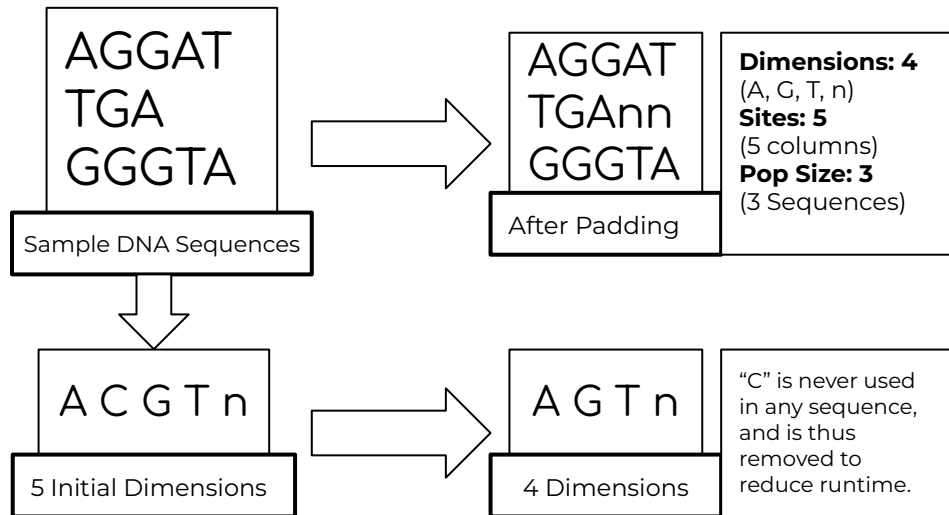
Application Time!



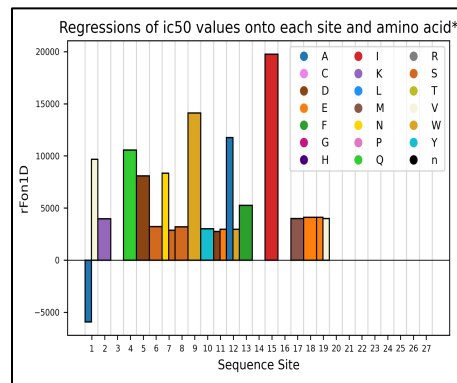
Updates to ortho_seqs

What I've Done:

- User-Friendly Updates and Efficiency
 - Automatic Padding
 - --dm, --sites, --pop_size flag Removal
 - Dimension Reducer
 - No Overwriting Data from --out_dir
- Features
 - Covariance Histogram
 - Covariance .csv File
 - --min_pct flag
 - rFon1D Bar Plot
 - --alphbt_input flag



Magnitude	Covariance	First Site	First Group	Second Site	Second Group	Percentile
0.08920745	0.08920745	3	2	12	20	99
0.0861943	0.0861943	3	2	13	20	99
0.08459642	0.08459642	3	2	10	2	99
0.07619613	0.07619613	3	2	16	20	99
0.06902849	0.06902849	2	14	14	20	99
0.06720234	0.06720234	3	2	14	20	99
0.0619065	0.0619065	4	9	13	20	99
0.06108473	0.06108473	3	2	15	20	99
0.06058254	0.06058254	3	2	11	17	99
0.05971512	0.05971512	2	14	15	20	99
0.05387144	0.05387144	2	14	13	20	99
0.05355186	0.05355186	3	2	18	20	99
0.05314098	0.05314098	4	9	12	20	99
0.05204529	0.05204529	2	14	3	2	99
0.05204529	0.05204529	4	9	14	20	99
0.05195398	0.05195398	4	9	11	17	99
0.05136048	0.05136048	3	12	12	2	99
0.05104091	0.05104091	3	2	9	10	99
0.05044741	0.05044741	3	2	17	20	99
0.04830168	0.04830168	4	9	16	20	99
0.04720599	0.04720599	4	9	10	2	99
0.04711468	0.04711468	4	9	15	20	99
0.04674945	0.04674945	2	14	10	2	99
0.04656684	0.04656684	2	14	16	20	99
0.04346238	0.04346238	3	2	19	20	99
0.04254931	0.04254931	4	9	9	10	99



What Will Be Done:

- User-Friendly Updates and Efficiency
 - GUI
 - Histogram Improvements
 - Runtime Efficiency
 - Only one file for both seq and phi
- Features
 - Third-Order Calculations for DNA
 - Second-Order Calculations for Proteins

...and more!



Special Thanks:

- Saba Nafees: Mentor
- Advisors: Eric Waltari, Joan Wong
- Code review: Pranathi Vemuri
- Server help: Saransh Kaul
- CZ Biohub

Contact Me:

- miles.woollacott@gmail.com

