

Supporting information

Machine learning assisted prediction and generation of antimicrobial peptides

Sukhvir Kaur Bhangu*, Nicholas Welch, Morgan Lewis, Fanyi Li, Brint Gardner, Helmut Thissen, Wioleta Kowalczyk*

S.K. Bhangu, N. Welch, F. Li, H. Thissen, W. Kowalczyk

CSIRO Manufacturing, Research Way, Clayton, Victoria 3168, Australia

E-mail: sukhvir.bhangu@unimelb.edu.au and wioleta.kowalczyk@csiro.au

M. Lewis

CSIRO Information Management & Technology, Kensington, Western Australia 6151, Australia

B. Gardner

CSIRO Information Management & Technology, Research Way, Clayton, Victoria 3168, Australia

Table S1: The physiochemical peptide features selected for model training.

Category	Property
Physiochemical	Overall charge (at pH 7)
	Charge density (at pH 7)
	Molecular weight
	Isoelectric point
	Aromaticity
	Gravy
	Boman index
	Aliphatic index
	Hydrophobicity profile – slope
	Hydrophobicity profile – intercept
Hydrophobicity profile – moment	
Secondary Structures	α -helix (%)
	β -sheets (%)
	Turns (%)

Table S2: The Table below shows the evaluation metrics (accuracy, precision, recall, F1, ROC AUC and confusion matrix's true and false negatives and positives) of four ML models across 5 Cross validation datasets and their respective mean and standard deviation.

<i>XGBOOST</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>mean</i>	<i>Std dev</i>
<i>Accuracy</i>	0.87564	0.87564	0.87564	0.84987	0.89059	0.87347	0.01314
<i>Precision</i>	0.87	0.86634	0.88144	0.87097	0.88833	0.87542	0.00818
<i>Recall</i>	0.88325	0.88833	0.86802	0.82234	0.89286	0.87096	0.02571
<i>F1 score</i>	0.87657	0.87719	0.87468	0.84595	0.89059	0.873	0.01466
<i>ROC AUC</i>	0.94906	0.93257	0.94553	0.92308	0.95944	0.94194	0.01276
<i>True Negatives</i>	171	170	174	172	175	172.4	1.85472
<i>False Positives</i>	26	27	23	24	22	24.4	1.85472
<i>False Negatives</i>	23	22	26	35	21	25.4	5.08331
<i>True Positives</i>	174	175	171	162	175	171.4	4.92341
<i>Random Forest</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>mean</i>	<i>Std dev</i>
<i>Accuracy</i>	0.85279	0.78934	0.82234	0.80407	0.88041	0.82979	0.033
<i>Precision</i>	0.83575	0.77941	0.82902	0.79412	0.87817	0.82329	0.03457
<i>Recall</i>	0.87817	0.80711	0.81218	0.82234	0.88265	0.84049	0.03299
<i>F1 score</i>	0.85644	0.79302	0.82051	0.80798	0.88041	0.83167	0.03215

<i>ROC AUC</i>	0.93561	0.88477	0.92489	0.88861	0.93544	0.91386	0.02256
<i>True Negatives</i>	163	152	164	154	173	161.2	7.57364
<i>False Positives</i>	34	45	33	42	24	35.6	7.39189
<i>False Negatives</i>	24	38	37	35	23	31.4	6.52993
<i>True Positives</i>	173	159	160	162	173	165.4	6.28013
<i>KNN</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>mean</i>	<i>Std dev</i>
<i>Accuracy</i>	0.64975	0.63706	0.63706	0.66667	0.65394	0.64889	0.01116
<i>Precision</i>	0.65775	0.63918	0.63235	0.67742	0.66129	0.6536	0.01615
<i>Recall</i>	0.62437	0.62944	0.65482	0.63959	0.62755	0.63516	0.01108
<i>F1 score</i>	0.64063	0.63427	0.64339	0.65796	0.64398	0.64405	0.00776
<i>ROC AUC</i>	0.6632	0.68379	0.68053	0.70178	0.69517	0.68489	0.01328
<i>True Negatives</i>	133	127	122	136	134	130.4	5.1614
<i>False Positives</i>	64	70	75	60	63	66.4	5.38888
<i>False Negatives</i>	74	73	68	71	73	71.8	2.13542
<i>True Positives</i>	123	124	129	126	123	125	2.28035
<i>SVM</i>	<i>1</i>	<i>2</i>	<i>3</i>	<i>4</i>	<i>5</i>	<i>mean</i>	<i>Std dev</i>
<i>Accuracy</i>	0.533	0.58883	0.59645	0.58015	0.59796	0.57928	0.02399
<i>Precision</i>	0.6383	0.81818	0.85185	0.75807	0.83929	0.78114	0.07835
<i>Recall</i>	0.15228	0.22843	0.2335	0.23858	0.2398	0.21852	0.03336
<i>F1 score</i>	0.2459	0.35714	0.36653	0.36293	0.37302	0.34111	0.04788
<i>ROC AUC</i>	0.54719	0.63929	0.6485	0.57472	0.60212	0.60236	0.03821
<i>True Negatives</i>	180	187	189	181	188	185	3.74166
<i>False Positives</i>	17	10	8	15	9	11.8	3.54401
<i>False Negatives</i>	167	152	151	150	149	153.8	6.67533
<i>True Positives</i>	30	45	46	47	47	43	6.54217

Table S3: Example of engine output for random peptide generation and their relative AMP score as determined using local and online prediction tools.

Sequence	Score
HFEIDAMMVGKMEI	0.140870
ENHTYWDHIVDGDTR	0.889225
GVKQQSWDSQSVE	0.544579
WTLNNSWRQGIKY	0.532396
GWRYQHTFVHQLASM	0.844666
...	...
KPNDALDTWHMRETD	0.203518
PSMDWDEEGHTMFAP	0.051499
PYNLRTTMVVKGRKM	0.419810
IPASLYIHMQQVTFN	0.047382
WATMEPRWRAMEGIH	0.054740

Table S4: Snapshot of top 16 motifs identified using design engine.

'W----T',
 'PY',
 'M--M',
 'L--A',
 'Q-----V',
 'Q-----S',
 'FY',
 'QQ',
 'YM',
 'M-Y',
 'YQ',
 'W-----T',
 'WR',
 'A-M',
 'G---Q',
 'W-----P',

Table S5: Snapshot of final list of peptides generated using top motifs with length 15 and batch size 25 with their respective AMP predicted score.

Sequence	Score
YFDYSNRWRYGILYR	0.082126
HYPDRALYVEFMVMW	0.072012
MPPWMLWIQWVVPYP	0.082951
QLMWQWVMSQKGWSV	0.141375
PLNNVFYVPYDDWMD	0.080597
VMMQYRMIVYGTAWG	0.550192
LGRMFWDDQYLFVHN	0.024660
RGLKTIDPYVYSHRT	0.876791
AWYHHFLWNHQYGPY	0.097113
WEKKRMVNQLTSMPN	0.050804
HAGWNWYVDTPFKS	0.869051
VMWQWRETVYAQSNG	0.030749
GWTQMCPGMTVQSF	0.936140
YQARWVQVELQVFFF	0.072117
LIWAVRLSHQTADIA	0.715828
HLQVSQAQIMWSHYI	0.128904
PMMFDQEQLMVVDVY	0.129583
WKPAKMWQWQNYVQP	0.670909
MLAVMPQSPVDFSW	0.063223
WQDLENDARSNFKYT	0.182409
MTVQWEVQKSSDTMQ	0.026656
AYWRVRLTILMQWPY	0.092487
QKRVVQFQRSSPYPS	0.622008
QFRWWGLVSTPVPW	0.758741

Table S6: List of peptides generated using algorithm and their respective antibacterial activities in terms of exact MIC value in μM .

Peptide name	Sequence	M.W.	Charge	<i>E. Faecalis</i> ATCC19433	<i>S. Aureus</i> ATCC25923	<i>K. Pneumoniae</i> ATCC43816	<i>A. Baumannii</i> AB5075	<i>P. Aeruginosa</i> PAOI	<i>E. Aerogenes</i> ATCC13048
MLWH01	LAKIVPHKIGKQLGT	1601.0111	4						
MLWH02	ALFLGGIWKKLPKAI	1653.0463	4				100	100	
MLWH03	GKIIGIKVATFLKA	1570.0303	4			200			
MLWH04	LLOKIFVVKGKAGLK	1640.0833	5						
MLWH05	GLWKIKLAFGKMFAK	1736.0293	5	25	100	25	12.5	25	25
MLWH06	GKIFGKFLALIKMW	1763.0652	4	12.5	25		25	200	200
MLWH07	GGAFFKLTQKIHKKA	1648.0158	5			200	50	200	
MLWH08	VQKLLKFKLIKITAKA	1727.1515	6	100		50	100	200	50
MLWH09	LLGTAFIKGGQIKTA	1605.9577	3	50	50	100	200	200	100
MLWH10	GIVNVTGAIKLIKLI	1537.9889	3						
MLWH11	AWLKFKKKFGFLFAKI	1823.1305	6	25	100	12.5	12.5	25	12.5
MLWH12	LVLQGVKAVVKAIN	1563.0205	3						
MLWH13	ILKLVANKIFHAVVF	1710.0677	3						
MLWH14	WFKAIAGLGGKFAIK	1661.0151	4	25	25	25	25		25
MLWH15	AVLIAKLFHKLNKIK	1734.1363	5	200		200	100	200	
MLWH16	GGCCWRCNRRVLRFC	1826.8514	5						
MLWH17	IAKVLNKFVFNHGIK	1702.0626	4						
MLWH18	FGNLIAIKQVKALNG	1571.9120	3						
MLWH19	FLPLIILKFIGKLLP	1723.1493	3	1.5	3		200		
MLWH20	LAFVKGILANVGKVK	1554.9944	4						
MLWH21	IKVILKLVTKNGNGI	1608.0416	4	200		200	50		200
MLWH22	TLVKAVAHIAKLLKF	1650.0677	4	100	100	50	12.5	200	200
MLWH23	KFITKFKVGLGLGTL	1620.0096	4			200	50	200	
MLWH24	TKFGQKVIKLIQFL	1732.0731	4	6	12.5	6	3	25	6
MLWH25	ALKAKLKFPGKFLFK	1734.1040	6			100	100	200	200
MLWH26	GLLTTGKKKLLGIIA	1524.0096	4				200		
MLWH27	VHKAFRKGSGAAILW	1638.9444	4			200	200		
MLWH28	FVTTLKKNFAKKKVG	1692.0419	5			200	200		200
MLWH29	GIVKAQLNIATFLTK	1614.9791	3						
MLWH30	KIVFGLAKIGKFIK	1659.0931	5						
MLWH31	FKLIKNVGALAFGTV	1575.9472	3						
MLWH32	KGFAKMFRLKLLAR	1805.1306	7	12.5	6	12.5	3	25	25
MLWH33	NFKIHKVGIKLGKA	1655.0942	6				100		
MLWH34	LGAPLKWIKGKFAIP	1637.0151	4						
MLWH35	KHVFALGGKLVKLLK	1683.0681	5			200	50		
MLWH36	AIVKKGWIKILAWVF	1770.1040	4	100	6	100	12	200	100
MLWH37	LGKVAHWGKLVHPKL	1681.0275	4						
MLWH38	FGLPKIKGIIKAIMI	1640.0543	4						
MLWH39	FVLTKATNGVLKALI	1585.9889	3						
MLWH40	KLKNTAGIFKGGQVK	1658.0325	6						
MLWH41	YFHLKKSFGKLLKRV	1877.1483	7				100		
MLWH42	FKRHGIGKFPKYASK	1762.0126	6						
MLWH43	WFGKFPQIACKRCC	1812.9075	5				25		
MLWH44	GKFFAKALFLAIVKT	1652.0147	4						
MLWH45	KLKGLIFAWFKGIIA	1703.0619	4	25	100	100	25	200	25
MLWH46	IGIKGAMPLKGIHK	1565.0184	5						
MLWH48	GPLMGAFIKLLPKF	1658.0075	4	200		100	25	100	200
MLWH49	IAVNQLTKIAKVAVG	1522.9530	3						

Table S7: List 2 of peptides generated using algorithm and their respective antibacterial activities.

Peptide name	Sequence	M.W.	Charge	<i>E. Faecalis</i> ATCC19433	<i>S. Aureus</i> ATCC25923	<i>K. Pneumoniae</i> ATCC43816	<i>A. Baumannii</i> AB5075	<i>P. Aeruginosa</i> PAO1	<i>E. Aerogenes</i> ATCC13048
MLWH51	LGVKPVFKFFKFFGK	1787.062	5	100	200	100	12.5	50	200
MLWH52	KHFSGFNGFFKFFSK	1803.955	5			200	100		
MLWH53	LVLKIKKKAVKALWE	1765.167	5			100	100	200	
MLWH54	IIKIIIVPNVIWNII	1760.104	2						
MLWH55	FIQMIIIIHYIIFI	1913.072	1						
MLWH56	GSKGSPFFSKKFRNF	1731.918	5						
MLWH57	QGWIKIILPHIRYIK	1876.153	4	25		12.5	12.5		25
MLWH58	AKKKVVAWTKWPPKI	1865.116	6			200	50	200	
MLWH59	KHGIPYIKKITNKI	1764.111	5				200		
MLWH60	KTFIKIKGTIKKWT	1818.157	7			200	6.25	50	
MLWH61	PSKYAKIAIKFHKQI	1770.064	5			50	50	50	50
MLWH62	KGGVWFGRILKHQSK	1739.008	5						
MLWH63	RFVKIGKGGHKMMY	1777.993	6						
MLWH64	FKTLKHKKFKRNPW	1969.186	7			200	50	50	
MLWH65	IQVGDIAYYKIKRIK	1771.116	5						
MLWH66	SKFAKNFKKKYKNKF	1904.112	8				200		
MLWH67	KIAKFVKKIKPSIM	1742.133	6				25	50	
MLWH68	WVPVWPKKLLKLTRG	1834.143	6				50	200	
MLWH69	IVPIKIIKIYKFQK	1857.23	6	200			200	200	
MLWH70	WKGKKKFFDKKKKWF	2027.195	8				50	200	
MLWH71	RQIDWYWKRRGKYAR	2052.125	6			200	100	200	
MLWH72	IMKKILINKIHNMK	1835.133	5				50		
MLWH73	FNKGKKFFKFKWV	1962.136	6	25	25	12.5	12.5	25	25
MLWH74	RMFFKFKKKFKWQ	2032.168	8	200	50	25	25	200	100
MLWH75	GAKFFKDFKRRDKV	1868.087	6						200
MLWH76	LHHIIIGIWKAIYT	1789.074	2						
MLWH77	KKLHHPPFKHFSQKP	1904.066	5						
MLWH78	WKIVKKFKMESHSPYH	1914.042	5			200	200		
MLWH79	KKKFKITGHTKYDYG	1812.002	5			200			
MLWH80	FMKMWKVMGKHKKF	1865.069	7				100	200	
MLWH81	IHAITHIHKFDVIA	1717.037	2						
MLWH82	FSKIIGITLKGTLII	1615.04	3						
MLWH83	GFKKEFDKFWWKKR	2056.149	6			100	200	200	
MLWH84	IHLVIRWIGELATM	1763.025	1						
MLWH85	AKVWVKIHKLKIKP	1819.262	8			100	25	100	
MLWH86	WIDIEITRIIVHI	1845.121	0						
MLWH87	PHFIKIVQIIYIIN	1822.12	2						
MLWH88	KNRIHKHILKIIYI	1915.233	6				50	50	
MLWH89	KFFKKLKDKTWDWI	1951.174	6				200		
MLWH90	KGFGKYKKGKKIYK	1799.127	9						
MLWH91	HPILKIIIMDKHIY	1845.103	2						
MLWH92	FNKAGFKNSKWNI	1808.018	6						
MLWH93	NHTFQIFMRKGNKKY	1910.007	5			100	100	>200	
MLWH94	KQQKFFKFKFFKE	2015.216	8						
MLWH95	DKFVKKPWKKVKKY	1948.21	8						
MLWH96	KAFKGFWKKGKHDK	1831.07	7						
MLWH97	FKLPFNFKKSRKA	1865.148	8				200		
MLWH98	PKWHKNMFKFKRKF	2048.174	8		200	50	100	200	100
MLWH99	VIYSIGLRIFNKIRP	1787.09	4	50		50	50	200	50
MLWH100	MHKVDHIIKKIMKTK	1848.092	5						

Confusion Matrices for Each Cross-Validation Fold

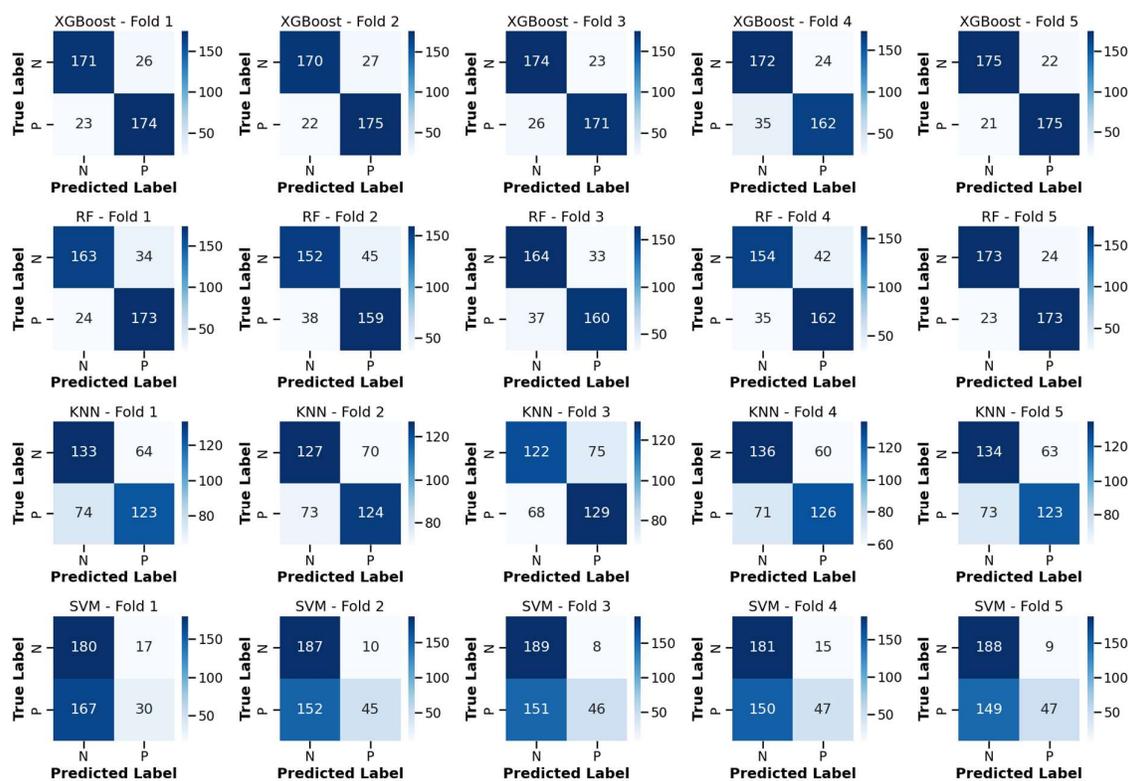


Figure S1: Confusion matrices for each Cross-Validation fold for four Models -XGBoost, RF, KNN, SVM

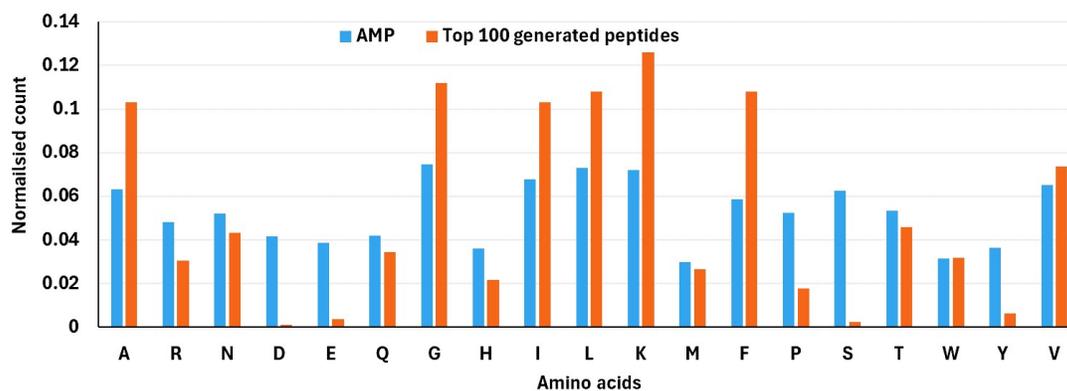


Figure S2: The comparison of amino acid composition of top 100 generated and predicted peptides to the training AMP dataset.

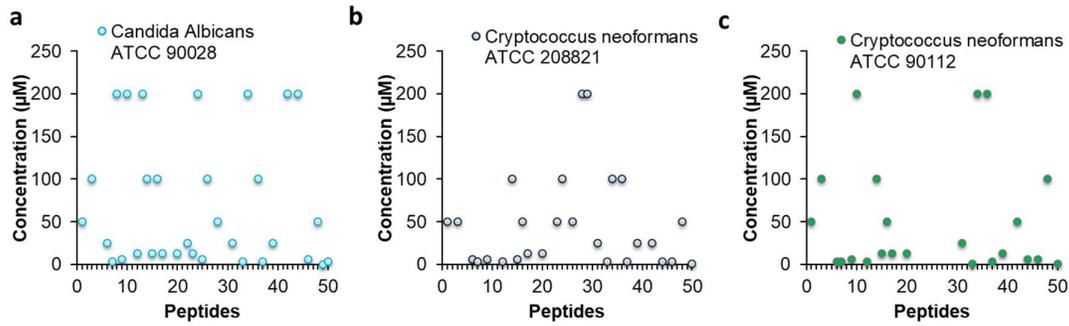


Figure S3: Antifungal activity of the peptide generated using algorithm. The scatter plot showing MIC values of various peptides against 3 different fungal strains, a) *Candida Albicans* ATCC90028, b) *Cryptococcus neoformans* ATCC208821, and c) *Cryptococcus neoformans* ATCC90112. The MIC testing experiments were performed in duplicates.

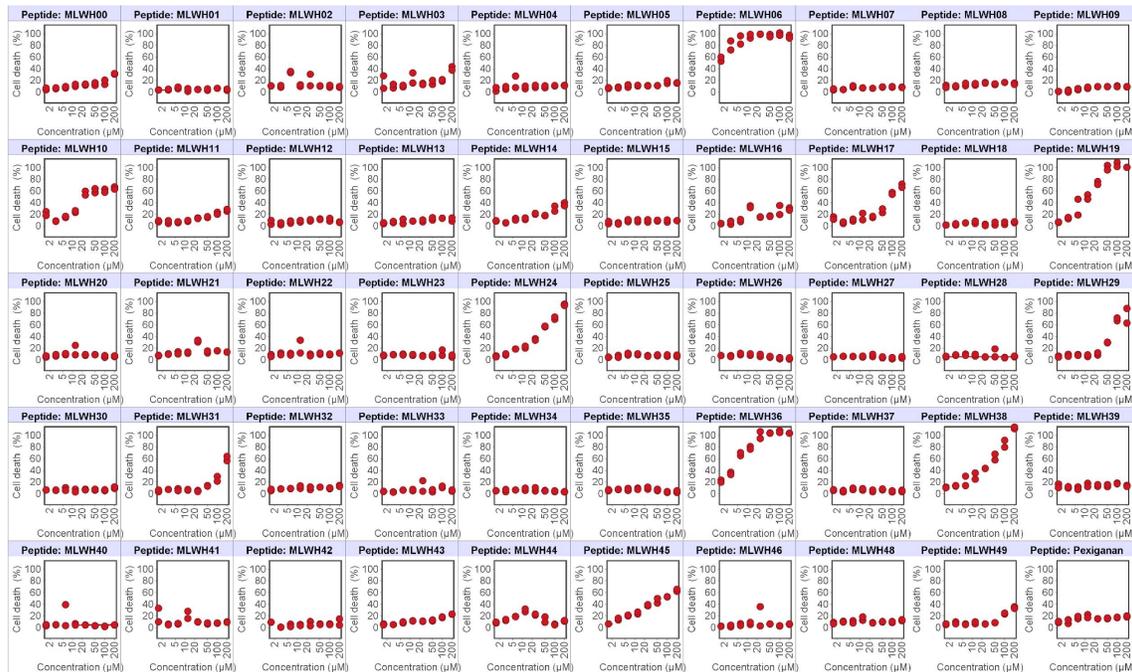


Figure S4: The graphs represent the hemolysis assay data for various peptides after incubation with whole human blood as a function of concentration. The experiment was performed in duplicates

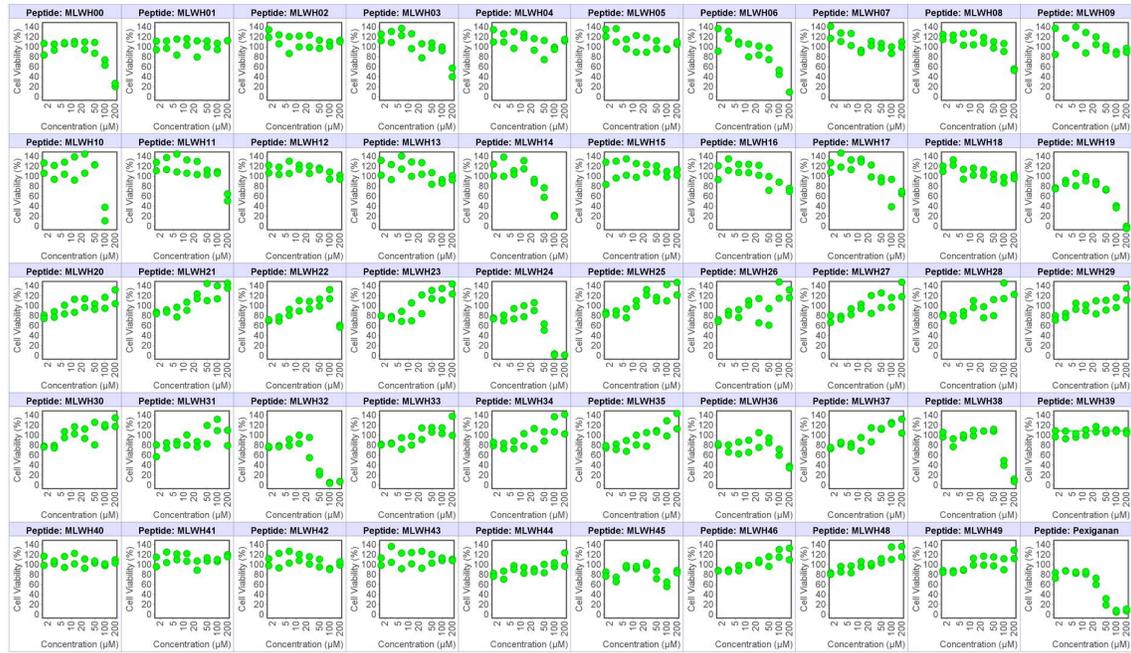


Figure S5: Cell viability of HEK-293 cells after 24 h incubation with various peptides as a function of concentration. The experiment was performed in duplicates