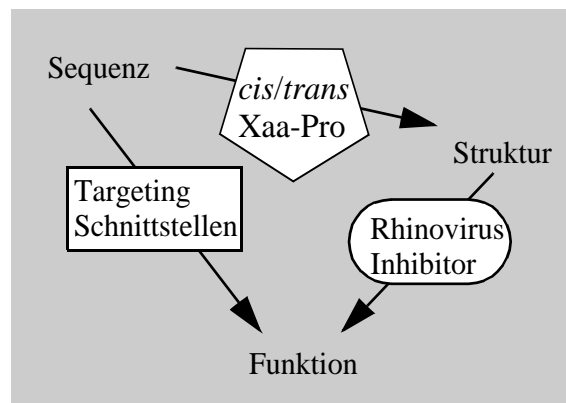


Bernd Alfred Jagla

# Funktionen und Strukturen abgeleitet aus Aminosäuresequenzen mit bioinformatischen Methoden



# Funktionen und Strukturen abgeleitet aus Aminosäuresequenzen mit bioinformatischen Methoden

Dissertation zur Erlangung des Doktorgrades des  
Fachbereichs Biologie, Chemie, Pharmazie  
der Freien Universität Berlin

vorgelegt von  
Bernd Alfred Jagla  
aus Berlin  
1999

Diese Arbeit wurde in der Zeit von Januar 1997 bis April 1999 am Institut für Chemie der Freien Universität Berlin und am Institut für Medizinisch Technische Physik und Lasermedizin des Fachbereichs Humanmedizin der Freien Universität Berlin angefertigt. Der Verfasser versichert, die Arbeit eigenständig durchgeführt und alle Hilfsmittel angegeben zu haben.

Gutachter: Prof. Dr. Paul Wrede  
Prof. Dr. Hanspeter Herzel

Tag der Disputation: 24.9.1999

## Danksagungen

Mein besonderer Dank gilt Herrn Prof. Dr. Paul Wrede für seine wohlwollende Unterstützung und Betreuung bei dieser Arbeit. Herrn Prof. Dr. Ing. Gerhard Müller, Prof. h.c. Dr. h.c., der die Arbeit möglich machte, möchte ich ebenso danken.

Herrn Dr. Johannes Schuchhardt möchte ich herzlich danken für die Betreuung und die stete geduldige Hilfsbereitschaft. Herrn Dr. Gisbert Schneider danke ich für die Unterstützung während der ersten Phase der Doktorarbeit.

Herrn Prof. Dr. Volker Erdmann danke ich für seine idelle und materielle Unterstützung beim Zustandekommen und der Durchführung dieser Arbeit.

Herrn Stefan Müller und Herrn Dr. Arif Malik danke ich für die Zusammenarbeit in der AG Bioinformatik.

Ich danke des weiteren den Herren Prof. Dr. Heinz Zeichhardt und Herrn Dr. Hans-Peter Grunert für das Testen der nona-Peptide sowie für die Unterstützung meiner Bemühungen.

Herrn Prof. Dr. Reinhart Heinrich danke ich stellvertretend für das gesamte Graduiertenkolleg "Dynamik und Evolution Makromolekularer und zellulärer Prozesse", an dem ich teilhaben durfte.

Dem ZIB, insbesondere Herrn Dr. Manfred Stolle danke ich für die Betreuung meiner Projekte an der CRAY.

Herrn Dr. Dieter Riedel danke ich für Unterstützung bei der Benutzung der Rechner der ZEDAT.

Darüber hinaus bedanke ich mich bei allen, die mich bei der Durchführung dieser Arbeit tatkräftig unterstützt haben.

# Abkürzungsverzeichnis

$\text{Å}$	Angström ( $=10^{-10}$ m)	$N_c$	Anzahl der Kodierungsvektoren
3D	dreidimensional	$N_m$	Anzahl der Muster
ACN	Adaptiv kodierendes KNN	NMR	Kernmagnetresonanz
C	Kohlenstoffatom	NN	Neuronales Netz
CATH	class(C), architecture(A), topology(T) and homologous superfamily (H)	$N_s$	Anzahl der Aminosäuren in einer Sequenz
CsA	Cyclosporin A	O	Sauerstoffatom
Cyp	Cyclophilin	PDB	Protein Datenbank
D <sub>2</sub> O	Schweres Wasser	PPIase	Peptidylprolylisomerasen
DSSP	Verzeichnis von Protein Sekundärstrukturen	$\theta$	Schwellenwert (threshold)
ER	Endoplasmatisches Retikulum	$\sigma$	Lernschrittweite
FK506	Immunsuppressivum	s	Standardabweichung
FKBP	FK506 bindendes Protein	$\Sigma$	Summenfunktion
H	Wasserstoffatom	$s^2$	Varianz
H <sub>2</sub> O	Wassermolekül	SAR	Sequenz-Aktivitäts-Relation
HRV	Humanes Rhinovirus	SME	Simulierte Molekulare Evolution
i	Index i	SOK	Selbst Organisierende Karte
$K, \vec{k}$	Eingabezahl, -vektor	SPase	Signalpeptidase
kD	Kilo Dalton	SPC18, SPC21, SPC22/23, GP23, P19, SPC12, Sec11p	Signalpeptidaseuntereinheiten
KNN	Künstliches Neuronales Netz	SRP	Signalerkennungspartikel
$\lambda$	Lernrate in einem evolutionären Algorithmus	w, $\vec{w}$	Gewichte
$\mu$	Index der Sequenzen	WWW	World Wide Web
MHC	Major Histocompatibility Complex	X, Xaa	eine der 20 natürlichen Aminosäuren
n	Index der Kodierungsvektoren	x, y	beliebige reelle Zahl
N	Stickstoffatom		

Tabelle der hier verwendeten Aminosäuren im Ein- und Drei-Buchstabenkode

A	Ala	Alanin
C	Cys	Cystein
D	Asp	Asparaginsäure
E	Glu	Glutaminsäure
F	Phe	Phenylalanin
G	Gly	Glycin
H	His	Histidin
I	Ile	Isoleucin
K	Lys	Lysin
L	Leu	Leucin
M	Met	Methionin
N	Asn	Asparagin
P	Pro	Prolin
Q	Gln	Glutamin
R	Arg	Arginin
S	Ser	Serin
T	Thr	Threonin
V	Val	Valin
W	Trp	Tryptophan
Y	Tyr	Tyrosin

Weiterhin wurden allgemein übliche Abkürzungen verwendet.

# Inhaltsverzeichnis

<b>KAPITEL 1</b>	<b>Einleitung</b>	5
	Schnittstellen humaner Signalpeptide	6
	Konformation der Peptidylprolylbindung	7
	Peptid Docking	8
<b>KAPITEL 2</b>	<b>Daten</b>	11
	Sequenzen der Schnittstellenregion humaner sekretorischer Proteine	12
	Peptidylprolylsequenzen	14
	Liganden-Docking	15
<b>KAPITEL 3</b>	<b>Methoden</b>	17
	Gütekriterien	17
	Kodierung der Daten	19
	Vergleich von Kodierungen	22
	Informationsanalyse	24
	Lineare Trennverfahren	25
	Schwerpunktanalyse (Zentroid-Verfahren)	25
	Bayes-Prediktor	26
	Mahalanobis-Distanzanalyse	26
	Künstliche Neuronale Netze	28
	Perzeptron	28
	Mehrlagige Neuronale Netze	31
	Neuronale Netze mit adaptiver Kodierung	32
	Lernstrategien	34
	Identifikation von Schnittstellen in Sequenzen unter Verwendung von Künstlichen Neuronalen Netzen	36

	Selbstorganisierende Karten (Kohonenkarten, SOK)	36
	Peptid Docking	37
<b>KAPITEL 4</b>	<b>Ergebnisse</b>	<b>39</b>
	Schnittstellen humaner Signalpeptide	39
	Informationsanalyse	40
	Positionsabhängige Häufigkeitsverteilung von Aminosäuren in Schnittstellenpeptiden	41
	Schwerpunktanalyse	44
	Bayes-Prediktor	45
	Hauptkomponentenanalyse (PCA)	45
	Mahalanobis-Distanzanalyse	48
	Künstliche Neuronale Netze	48
	Lernstrategien	49
	Vergleich von Kodierungsmethoden	51
	Kohonen Netze	61
	Vorhersage der <i>cis/trans</i> - Konformation von Peptidylprolylbindungen	64
	Informationsanalyse	67
	Schwerpunktanalyse	72
	Hauptkomponentenanalyse	73
	Mahalanobis-Distanzanalyse	76
	Adaptive Kodierung	76
	Kohonennetze	77
	Klassifikation und Information der Oberflächeneigenschaften	78
	3D-Umgebung	79
	Sekundärstruktur am Prolin	80
	Design von Proteinliganden	83
	Simulierte Molekulare Evolution	83
	Peptid-Docking	85
<b>KAPITEL 5</b>	<b>Diskussion</b>	<b>93</b>
	Schnittstellen humaner Signalpeptide	93
	Konformationsvorhersage der Peptidylprolylbindungen	104
	Ligandendesign	114
	Simulierte Molekulare Evolution	115
	Peptid-Docking	115



<b>KAPITEL 6</b>	<b>Literaturverzeichnis</b>	117
<b>KAPITEL 7</b>	<b>Lebenslauf</b>	125
<b>KAPITEL 8</b>	<b>Anhang I Humane Schnittstellendaten</b>	127
<b>KAPITEL 9</b>	<b>Anhang II Physikochemische Eigenschaften von Aminosäuren</b>	139
<b>KAPITEL 10</b>	<b>Anhang III: Gruppierung der PDB Datenbank nach Wechselwirkungsklassen</b>	151



---

## KAPITEL 6 Zusammenfassung

---

Verschiedene Methoden zur Funktions- und Strukturvorhersage ausgehend von der Primärstruktur wurden vorgestellt.

Ein adaptiv kodierendes Neuronales Netz (ACN) ist entwickelt worden, mit dem es möglich ist, die charakteristischen physikochemischen Eigenschaften in Proteinsequenzen zu bestimmen. Hiermit ist es möglich, 96% der experimentell verifizierten Schnittstellen von humanen sekretorischen Proteinen vorherzusagen. Bisher waren nur Vorhersagen mit einer Genauigkeit von bis zu 68% möglich. Außerdem können mit den ACN Mutationen in den Schnittstellenregionen humaner sekretorischer Proteine korrekt vorhergesagt werden. Es wurde somit erfolgreich gezeigt, wie aus der Sequenzinformation auf eine Funktion geschlossen werden kann.

Die Vorhersage der Konformation von Peptidylprolylbindungen ist mit dem vorliegenden Datensatz nicht zufriedenstellend lösbar. Für eine erfolgversprechende Analyse müßte die 10- bis 20-fache Menge an Daten vorliegen. Trotzdem konnte gezeigt werden, daß die MHC-I Proteine eine charakteristische *cis*-Prolylgruppe besitzen, die eine Unterscheidung von anderen Proteinen erlaubt.

Die Simulierte Molekulare Evolution wurde erfolgreich zum Entwurf zellprotektiver Peptide angewendet. Ein Algorithmus zum Design inhibitorischer Proteine, basierend auf der Analyse von Protein-Protein-Wechselwirkungen, wurde entwickelt und erfolgreich auf zellprotektive Experimente getestet.

---

---

## KAPITEL 7    Summary

---

Several methods were presented for the prediction of function and structure based on the primary structure of proteins.

An adaptive coding artificial neural network (ACN) was developed that characterizes the physico chemical features in protein sequences. It is now possible to identify 96% of experimentally verified cleavage sites in human secretory proteins. Hitherto only 68% of an independent test data were correctly identified. Furthermore mutations within the cleavage site region of human secretory proteins are correctly predicted. It is shown how to use ACN to build up a sequence activity relation.

It was not possible to predict the conformation of the peptidyl prolyl bond using the available data. However there should be 10 - 20 times more non-homologues sequence data to make reliable predictions. Nevertheless, it was possible to show that MHC-I proteins have a characteristic *cis*-proline group that can be used to distinguish those proteins from others.

The simulated molecular evolution was successfully applied to design cell protective peptides against the Human Rhinovirus. An algorithm for the design of inhibitory peptides using the structural information of the receptor protein has been developed. Using this method other peptides were designed. All peptides were experimentally tested, showing that theoretical and practical results are in agreement.

---



---

## KAPITEL 7    Lebenslauf

---

Name:                    Bernd Alfred Jagla  
Geburtstag:            12. April 1970  
Geburtsort:            Berlin  
Eltern:                   Eheleute Ludwig und Dr. med. Maria Jagla

### Schule

1976 - 1977            Matthias-Claudius-Schule in Düsseldorf  
1977 - 1978            Städt. Kath. Grundschule Düsseldorf  
1978 - 1980            Mühlenau-Grundschule Berlin, Zehlendorf  
1980 - 1982            Erich-Kästner-Grundschule Berlin, Zehlendorf  
1982 - 1986            Arndt-Oberschule Berlin, Zehlendorf  
1986 - 1989            Walther-Rathenau-Oberschule Berlin, Wilmersdorf  
23.5.1989              Abitur

### Studium

WS 1989                Aufnahme des Studiums der Chemie (Diplom) an der Freien  
                              Universität Berlin  
WS 1990-                Studium der Informatik für Lehramt als 2. Fach  
SS 1991  
25.9.1991                Vordiplom in Chemie  
WS 1991                Fortführung des Chemiestudiums an der Universität Erlangen-  
                              Nürnberg  
SS 1992                Fortsetzung des Studiums an der Technischen Universität München  
24.10.1995              Diplom in Chemie an der TU-München (Prof. Ugi)  
                              Zusammenarbeit mit der FU-Berlin (Prof. Wrede)  
                              Titel: Sequenzbasierte Vorhersage der Konformation von  
                              prolinhaltigen Peptidgruppen in Proteinen

---

