

Discovering structural motifs using a structural alphabet

Application to magnesium-binding sites

Sebastian Bartschat

21. Juli 2007

Inhalt

- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets
 - unsupervised cluster analyser
 - Strukturvorhersage mittels Bayes-Theorem
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen
 - Vorbereitungen
 - Analyse

Wozu dient Magnesium?

Magnesium ist einer der vielseitigsten metallischen Co-Faktoren:

- an Reaktionen mit Phosphatgruppenübertragung beteiligt (ATP-Magnesium liegt vor)
- nötig für die Nukleinsäurebiosynthese, aber auch für die Stabilität
- Stabilisierung von Proteinstrukturen
- besitzt calciumantagonistische Wirkung

Bindungseigenschaften

- Bindung an Seitenketten von Asp und Glu bzw. Asn und Gln
- bis jetzt nur Studien mit Proteinen mit hoher Sequenzähnlichkeit
2 Sequenzmotive :
 - 1 NADFDGD RNA Pol. und DNA Pol. I
 - 2 YXDD / LXDD Reverse Transkriptase und Telomerase

Fragestellungen

- 1 Besitzen Mg^{2+} -Bindungsstellen strukturelle Neigungen?
- 2 Existieren strukturelle Motive, auch wenn keine Sequenzähnlichkeit vorhanden ist?
- 3 Können strukturelle Motive bestimmten Proteinfunktionen zugeordnet werden?
- 4 Welche Spezifität weisen die gefundenen Motive auf?

- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets**
 - **unsupervised cluster analyser**
 - Strukturvorhersage mittels Bayes-Theorem
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen
 - Vorbereitungen
 - Analyse

- **Ziel** : Generierung eines Sets von Strukturblöcken, um die Struktur eines Proteins anhand der Sequenz bestmöglich zu approximieren

- **Ziel** : Generierung eines Sets von Strukturblöcken, um die Struktur eines Proteins anhand der Sequenz bestmöglich zu approximieren
- **Vorüberlegung** : Zerlegung des Proteins in überlappende Blöcke von $M = 5$ AS (ProteinBlock - PB)

- **Ziel** : Generierung eines Sets von Strukturblöcken, um die Struktur eines Proteins anhand der Sequenz bestmöglich zu approximieren
- **Vorüberlegung** : Zerlegung des Proteins in überlappende Blöcke von $M = 5$ AS (ProteinBlock - PB)
- **ProteinBlock** : Beschreibung durch Vektor aus Diederwinkeln

$$V(\psi_{n-2}, \phi_{n-1}, \psi_{n-1}, \phi_n, \psi_n, \phi_{n+1}, \psi_{n+1}, \phi_{n+2})$$

um das zentrale $C\alpha_n$ des ProteinBlocks

- self-organizing maps (SOM) mit 2 Lernphasen
- Ähnlichkeitsmessung zwischen 2 Vektoren mittels RMSDA (root mean square deviation on angular values)

$$\text{RMSDA}(V_1, V_2) = \sqrt{\frac{\sum_{i=1}^{i=M-1} [\psi_i(V_1) - \psi_i(V_2)]^2 + [\phi_{i+1}(V_1) - \phi_{i+1}(V_2)]^2}{2(M-1)}}$$

- ProteinBlöcke PB_k der Karte sind zu Beginn durch einen Vektor $W(k)$ definiert

first training

- Einlesen von C Proteinen als konsekutive Signale der einzelnen PBs
- Bestimmung des ProteinBlocks W_k mit dem kleinsten Abstand zu dem induzierten Vektor $V(m)$ (mittels RMSDA)
- Neugewichtung des PB_k :

$$W(k) \leftarrow W(k) + (V(m) - W(k)) * \nu(c)$$

wobei $\nu(c)$ eine Lernrate darstellt, die mit der Zeit abnimmt

second training

- Ermittlung der Transitionsmatrix der PBs
- erneutes Einlesen von C Proteinen als Signal konsekutiver Vektoren
- Bestimmung von n ähnlichen Blöcken zu einem induzierten Vektor $V(m)$ (ergibt Menge von Vektoren W_k)
- Neugewichtung desjenigen $W(k)$ mit der höchsten Übergangswahrscheinlichkeit von $W(j)$ ($W(j)$ mit $V(m - 1)$ assoziiert)

shrinking process

Ziel : Reduzierung der Neurone auf eine optimale Anzahl

- 1 untersuche alle PB_i und PB_j der Karte K auf Ähnlichkeit ($i, j \in \{1 \dots |K|\}$)
→ strukturelle Ähnlichkeit
→ Translationsähnlichkeit
- 2 falls PB_i und PB_j ähnlich: lösche den PB der seltener ist und gehe zu Schritt 1
sonst: ENDE

- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets**
 - unsupervised cluster analyser
 - **Strukturvorhersage mittels Bayes-Theorem**
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen
 - Vorbereitungen
 - Analyse

Ziel: Genaue Approximation der 3D-Struktur aus Proteinsequenz

Vorüberlegungen

- jedem ProteinBlock ist eine Menge von Sequenzen zugeordnet
- Erstellen eines Sequenzfenster der Größe $[-\omega, +\omega]$ auf der Proteinsequenz um das $C\alpha$ des jeweiligen PB

Ziel: Genaue Approximation der 3D-Struktur aus Proteinsequenz

Vorüberlegungen

- jedem ProteinBlock ist eine Menge von Sequenzen zugeordnet
- Erstellen eines Sequenzfenster der Größe $[-\omega, +\omega]$ auf der Proteinsequenz um das $C\alpha$ des jeweiligen PB
- Vorkommensmatrix für AS an Pos. j im ProteinBlock k : $n_{i,j}^k$

- Berechnung von $P(a_i \text{ in } j | PB_k) = \frac{n_{i,j}^k}{N_k}$

k ... Nummer des ProteinBlocks

i ... Nummer der Aminosäure

j ... Position im Sequenzfenster

Beschreibung der Daten:

- Kullback-Leibler-Divergenz

$$K_k(p_j, q) = \sum_i p_{ji} * \ln \left(\frac{p_{ji}}{q_i} \right)$$

Beschreibung der Daten:

- Kullback-Leibler-Divergenz

$$K_k(p_j, q) = \sum_i p_{ji} * \ln \left(\frac{p_{ji}}{q_i} \right)$$

Unterschied zwischen der AS-Verteilung an der Stelle j im Block k zur allgemeinen Verteilung

dient der Suche nach Positionen mit großer Spezifität

- Z-Score

$$z = \frac{(n_{ij}^k - n_{ib})}{\sqrt{n_{ib}}}$$

$n_{ib} = N_k * f_i$ beschreibt den Erwartungswert der i-ten AS

dient der Suche nach AS mit großer Spezifität für eine bestimmte Pos. in einem bestimmten ProteinBlock

- gesucht ist ProteinBlock k so dass $P(PB_k | X_s)$ maximal ist

- $$P(PB_k | X_s) = \frac{P(X_s | PB_k) * P(PB_k)}{P(X_s)}$$

$$P(X_s) = \prod_{j=-\omega}^{j=+\omega} P(a_j)$$

Wahrscheinlichkeit die Sequenz X_s zu beobachten; ohne gegebene Strukturinformationen

$$P(X_s) = \prod_{j=-\omega}^{j=+\omega} P(a_j)$$

Wahrscheinlichkeit die Sequenz X_s zu beobachten; ohne gegebene Strukturinformationen

$$P(X_s | PB_k) = \prod_{j=-\omega}^{j=+\omega} P(a_j | PB_k)$$

Wahrscheinlichkeit die Sequenz $X_s(a_{-\omega}, \dots, a_{+\omega})$ zu beobachten, wenn ein bestimmter ProteinBlock k gegeben ist

Berechnung des optimalen ProteinBlockes PB* für eine gegebene Sequenz X_s mittels des Verhältnisse R_k :

$$R_k = \frac{P(X_s | PB_k)}{P(X_s)} = \frac{P(PB_k | X_s)}{P(PB_k)}$$

Berechnung des optimalen ProteinBlockes PB^* für eine gegebene Sequenz X_s mittels des Verhältnisse R_k :

$$R_k = \frac{P(X_s | PB_k)}{P(X_s)} = \frac{P(PB_k | X_s)}{P(PB_k)}$$

- 1 $PB^* = k$ mit $\ln(R_k)$ maximal
- 2 PB^* ist in der Menge der r -besten Blöcke k

- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets**
 - unsupervised cluster analyser
 - Strukturvorhersage mittels Bayes-Theorem
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen
 - Vorbereitungen
 - Analyse

- das Training resultierte in einem Alphabet aus 16 Protein-Blöcken

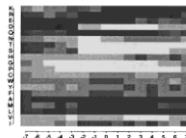
PB label	Frequency (%)	RMSD (Å)	anr	Transitions (%)			Str. II (%)			Coarse char.
				1st	2nd	3rd	α	Coil	β	
<i>a</i>	3.93	0.52	1.01	54.8(c)	16.5(f)	8.0(b)	0.1	76.7	23.3	N-cap β
<i>b</i>	4.58	0.51	1.00	44.4(d)	17.9(c)	13.7(f)	0.2	86.7	13.1	N-cap β
<i>c</i>	8.63	0.51	1.28	62.2(d)	24.4(f)	5.6(e)	0.1	58.2	41.7	N-cap β
<i>d</i>	18.84	0.48	2.74	51.9(f)	25.6(c)	19.2(e)	0.0	28.4	71.6	β
<i>e</i>	2.31	0.54	1.11	80.4(h)	9.1(d)		0.0	49.8	50.2	C-cap β
<i>f</i>	6.72	0.50	1.00	60.7(k)	36.3(b)		0.0	72.5	27.5	C-cap β
<i>g</i>	1.28	0.74	1.05	37.5(h)	28.0(c)	19.1(o)	6.9	83.8	9.3	mainly coil
<i>h</i>	2.35	0.62	1.04	62.4(i)	18.1(j)	10.2(k)	0.0	81.5	18.5	mainly coil
<i>i</i>	1.62	0.56	1.01	87.7(a)			0.0	94.5	5.5	mainly coil
<i>j</i>	0.96	1.03	1.01	17.0(a)	16.6(b)	16.1(l)	3.7	87.9	8.4	mainly coil
<i>k</i>	5.46	0.59	1.00	76.2(l)	13.6(b)		35.1	64.2	0.7	N-cap α
<i>l</i>	5.35	0.63	1.01	68.5(m)	9.2(p)	7.0(c)	44.4	54.9	0.7	N-cap α
<i>m</i>	30.04	0.43	6.74	33.8(n)	18.5(p)	9.7(b)	86.7	13.2	0.1	α
<i>n</i>	1.93	0.61	1.03	90.9(o)			68.4	31.3	0.3	C-cap α
<i>o</i>	2.60	0.60	1.02	74.7(p)	8.3(m)		43.1	56.8	0.1	C-cap α
<i>p</i>	3.41	0.46	1.00	58.1(a)	22.7(c)	11.1(m)	11.2	87.5	1.3	C-cap α to N-cap β

- PBm weist kleinen rmsd-wert auf (0.43 Å)
- überrepräsentiert: aliphatische AS
unterrepräsentiert: α -Helixbrecher
- KLd zeigt breites Spektrum von spezifischen Stellen

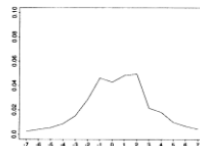
(d) BP m



(h)



(i)





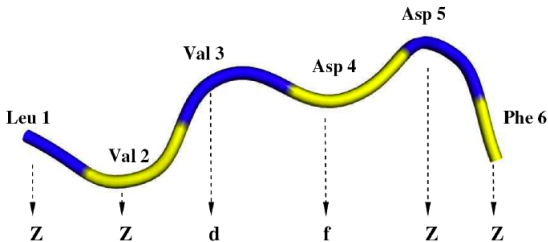
Ergebnisse der Vorhersagestrategie am Beispiel des Proteins 2aak

Left	Sequence		True PB	Neg	Predicted PBs		
	c. window	Right			1st	2nd	3rd
DMSTP	ARKLM	RDFKR	<i>l</i>	2.74	<i>m</i> (11.21)	<i>l̄</i> (0.63)	<i>d</i> (0.36)
MSTPA	RKLMR	DFKRL	<i>m</i>	2.43	<i>m̄</i> (17.51)	<i>d</i> (0.74)	<i>f</i> (0.43)
STPAR	KLMRD	FKRLQ	<i>m</i>	2.06	<i>m̄</i> (33.68)	<i>f</i> (0.38)	<i>d</i> (0.18)
TPARK	LMRDF	KRLQQ	<i>m</i>	2.63	<i>m̄</i> (11.05)	<i>l</i> (0.51)	<i>k</i> (0.36)
PARKL	MRDFK	RLQQD	<i>m</i>	2.48	<i>m̄</i> (22.13)	<i>f</i> (1.25)	<i>b</i> (0.40)
ARKLM	RDFKR	LQQDP	<i>m</i>	3.78	<i>m̄</i> (7.77)	<i>k</i> (1.90)	<i>c</i> (0.54)
RKLMR	DFKRL	QQDPP	<i>m</i>	2.92	<i>m̄</i> (12.48)	<i>b</i> (0.94)	<i>c</i> (0.34)
KLMRD	FKRLQ	QDPPA	<i>m</i>	3.49	<i>m̄</i> (12.98)	<i>n</i> (2.60)	<i>p</i> (0.73)
LMRDF	KRLQQ	DPPAG	<i>m</i>	6.32	<i>m̄</i> (3.51)	<i>n</i> (0.55)	<i>d</i> (0.38)
MRDFK	RLQQD	PPAGI	<i>m</i>	8.61	<i>p</i> (2.02)	<i>b</i> (1.08)	<i>m̄</i> (1.02)
RDFKR	LQQDP	PAGIA	<i>m</i>	4.82	<i>p</i> (3.08)	<i>d</i> (1.12)	<i>c</i> (0.44)
DFKRL	QQDPP	AGIAG	<i>c</i>	2.55	<i>c̄</i> (4.43)	<i>d</i> (0.19)	<i>p</i> (0.12)
FKRLQ	QDPPA	GIAGA	<i>c</i>	3.10	<i>f</i> (13.43)	<i>c̄</i> (2.87)	<i>k</i> (0.23)
KRLQQ	DPPAG	IAGAG	<i>e</i>	5.45	<i>b</i> (7.14)	<i>ḡ</i> (1.94)	<i>g</i> (1.68)
RLQQD	PPAGI	AGAGI	<i>h</i>	5.34	<i>b</i> (12.39)	<i>h̄</i> (6.72)	<i>l</i> (3.16)
LQQDP	PAGIA	GAGIS	<i>i</i>	4.97	<i>ī</i> (11.29)	<i>p</i> (5.29)	<i>c</i> (1.40)
QQDPP	AGIAG	AGISG	<i>a</i>	4.75	<i>g</i> (15.58)	<i>ḡ</i> (6.16)	<i>e</i> (3.24)
QDPPA	GIAGA	GISGA	<i>c</i>	6.75	<i>b</i> (7.15)	<i>h</i> (4.01)	<i>c̄</i> (2.32)

- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets
 - unsupervised cluster analyser
 - Strukturvorhersage mittels Bayes-Theorem
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen**
 - Vorbereitungen**
 - Analyse

- Erstellen des Datensatzes:
 - Sequenzähnlichkeit $< 30\%$
 - Auflösung $< 2,5 \text{ \AA}$
 - Anzahl AS, die an Bindung beteiligt sind ≥ 3 AS
- Datensatz enthält 77 Bindungen in 70 Proteinen

- Erstellen des Datensatzes:
 - Sequenzähnlichkeit $< 30\%$
 - Auflösung $< 2,5 \text{ \AA}$
 - Anzahl AS, die an Bindung beteiligt sind ≥ 3 AS
- Datensatz enthält 77 Bindungen in 70 Proteinen
- Umwandlung mittels PBE-Webinterface in Strukturalphabet



first-shell und second-shell Liganden

1^{st} -shell : Entfernung zwischen Metallion und Donor $\leq 2.5 \text{ \AA}$

2^{nd} -shell : Entfernung $\leq 3.5 \text{ \AA}$

first-shell und second-shell Liganden

1^{st} -shell : Entfernung zwischen Metallion und Donor $\leq 2.5 \text{ \AA}$

2^{nd} -shell : Entfernung $\leq 3.5 \text{ \AA}$

Strukturmotive

- $k \geq 3$; k beschreibt die Anzahl der Wiederholungen
- gleiche Strukturbuchstaben
- ähnlich große Zwischenräume

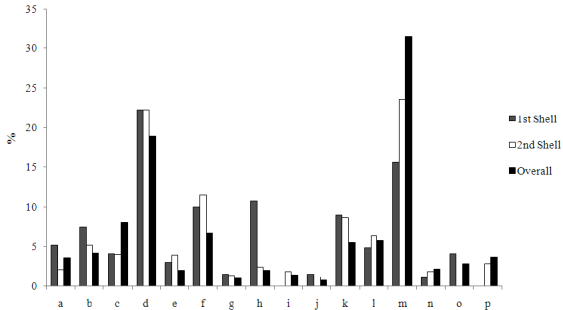
- 1 Einführung
- 2 Fragestellungen
- 3 Generierung des Struktur-Alphabets
 - unsupervised cluster analyser
 - Strukturvorhersage mittels Bayes-Theorem
 - Ergebnisse
- 4 Anwendung auf Mg^{2+} -Bindungsstellen
 - Vorbereitungen
 - Analyse

strukturelle Präferenzen der Mg^{2+} -Bindungststellen

- Analyse der Struktur von 1st- und 2nd shell Bindungsstellen
- Vergleich der Frequenzen der beiden Schalen mit dem allgemeinen Auftreten der Struktur

strukturelle Präferenzen der Mg^{2+} -Bindungststellen

- Analyse der Struktur von 1st- und 2nd shell Bindungsstellen
- Vergleich der Frequenzen der beiden Schalen mit dem allgemeinen Auftreten der Struktur



Mg^{2+} -Bindungsstellen präferieren Loops statt Helizes und Faltblätter

Suche nach strukturellen Motiven

1st-shell : 4 Motive, die 21% aller Bindungsstellen repräsentieren

→ $e(24-47)h(24)k$

→ $f(1)h(109-349)b$

→ $f(2)h(126-158)m$

→ $k(26-29)h(1)a$

Suche nach strukturellen Motiven

1st-shell : 4 Motive, die 21% aller Bindungsstellen repräsentieren

→ $e(24-47)h(24)k$

→ $f(1)h(109-349)b$

→ $f(2)h(126-158)m$

→ $k(26-29)h(1)a$

- gleiche Motive haben meist gleiche CATH-Nummern (hierarchische Klassifizierung von Proteindomänen)
- Möglichkeit für noch nicht klassifizierte Proteine?

Ergebnisse

Motif ^a	PDB code	Mg^{2+} -Ligands	CATH number ^b	Functional Group ^c	EC code ^d
e(24-47)h(24)k	ISJC	D ¹⁸⁹ , E ²¹⁴ , D ²³⁹	3.20.20.120	Lyase ^e , Isomerase ^f	-
	ITKK	D ¹⁹¹ , E ²¹⁹ , D ²⁴⁴	3.20.20.120	Isomerase ^f	-
	2AKZ	D ²⁴⁴ , E ²⁹² , D ³¹⁷	-	Lyase ^e	4.2.1.11
f(1)h(109-349)b	1O08	D ¹⁰⁰⁸ , D ¹⁰¹⁰ , D ¹¹⁷⁰	3.40.50.1000	Isomerase ^f	5.4.2.6
	1U7P	D ¹¹ , D ¹³ , D ¹²³	NYC	Hydrolase ^g	-
	1WPG	D ³⁵¹ , T ³⁵³ , D ⁷⁰³	3.40.50.1000	Hydrolase ^g	3.6.3.8
	2B82	D ⁴⁴ , D ⁴⁶ , D ¹⁶⁷	3.40.50.1000	Hydrolase ^g	3.1.3.2
	2C4N	D ⁹ , D ¹¹ , D ²⁰¹	NYC	Hydrolase ^g	-
f(2)h(126-158)m	1KA1	D ¹⁴² , D ¹⁴⁵ , D ²⁹⁴	3.30.540.10	Hydrolase ^g	3.1.3.7
	1NUY	D ¹¹¹⁸ , D ¹¹²¹ , E ¹²⁸⁰	3.30.540.10+	Hydrolase ^g	3.1.3.11
	2BJJ	E ¹⁰⁹⁰ , D ¹⁰⁹³ , D ¹²²⁰	3.30.540.10+	Hydrolase ^g	3.1.3.25
k(26-29)h(1)a	1ITZ	D ¹⁶⁸ , N ¹⁹⁸ , I ²⁰⁰	3.40.50.970	Transferase ^h	2.2.1.1
	1POX	D ⁴⁴⁷ , N ⁴⁷⁴ , Q ⁴⁷⁶	3.40.50.970+	Oxidoreductase ⁱ	1.2.3.3
			3.40.50.1220		
	1UMD	D ¹⁷⁵ , N ²⁰⁴ , Y ²⁰⁶	3.40.50.970	Oxidoreductase ⁱ	1.2.4.4
	1ZPD	D ⁴⁴⁰ , N ⁴⁶⁷ , G ⁴⁶⁹	3.40.50.970	Lyase ^e	4.1.1.1
2C3M	D ⁹⁶³ , T ⁹⁹¹ , V ⁹⁹³	3.40.50.970	Oxidoreductase ⁱ	1.2.7.1	

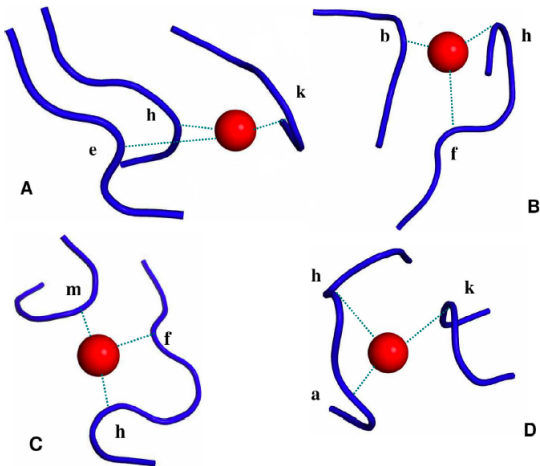
Wie spezifisch sind die Resultate?

- Abgleich gegen Nichtmetallproteine aus de Brevern's DB (Auswahl der Proteine nach vorab definierten Kriterien)
 - matches für $f(1)h(109-349)b$ und $k(26-29)h(1)a$
 - $e(24-47)h(24)k$ & $f(2)h(126-158)m$ sind metallspezifisch

Wie spezifisch sind die Resultate?

- Abgleich gegen Nichtmetallproteine aus de Brevern's DB (Auswahl der Proteine nach vorab definierten Kriterien)
 - matches für $f(1)h(109-349)b$ und $k(26-29)h(1)a$
 - $e(24-47)h(24)k$ & $f(2)h(126-158)m$ sind metallspezifisch
- Abgleich gegen calciumbindende Proteine (Anwendung des gleichen Verfahrens wie für die Mg^{2+} -Bindungsstellen)
 - $f(1)h(109-349)b$ und $k(26-29)h(1)a$ wurden in 1 bzw. 2 Proteinen gefunden
 - $e(24-47)h(24)k$ & $f(2)h(126-158)m$ sind metall- und magnesiumspezifisch

4 extrahierte Bindungsmotive :



Literatur



Minko Dudev, Carmay Lim *Discovering structural motifs using a structural alphabet: Application to magnesium-binding sites*, Bioinformatics 2007 Mar 28; 8(106)



de Brevern AG, Etchebest C, Hazout S *Bayesian Probabilistic Approach for Predicting Backbone Structures in Terms of Protein Blocks* PROTEINS: Structures, Function and Genetics 2000