

Algorithmische Bioinformatik

Multiple Sequence Alignment

Sum-of-pairs Score

Center-Star Score

Ulf Leser

Wissensmanagement in der
Bioinformatik



Ziel dieser Vorlesung

- Aufgabenstellung „Multiples Sequenzalignment“ verstehen
- Einen elegante (aber nutzlosen) Algorithmus kennenlernen
- Dadurch das Problem der MSA-Zielfunktion verstehen

Inhalt dieser Vorlesung

- Multiples Sequenzalignment
- Sum-Of-Pair Zielfunktion
- Center-Star Zielfunktion

Definition

- Bisher: Alignment zweier Strings
- Jetzt: Alignment von $k > 2$ Strings
- Definition

*Ein **multiple Sequenzalignment (MSA)** von k Strings S_i , $1 \leq i \leq k$, ist eine Tabelle mit k Zeilen und l Spalten, so dass*

- In Zeile i steht String S_i , mit beliebig eingefügten Leerzeichen*
 - Jedes Zeichen jedes S_i steht in exakt einer Spalte*
 - In keiner Spalte stehen nur Leerzeichen*
- Bemerkungen
 - Es folgt, dass $l = |\text{MSA}| \leq \sum(|S_i|)$

Beispiel

S_1	M---	AIDE----	NKQKALAAALGQIEKQFGKGS	SIMRLGEDR-	SMDVETISTGSLSLDI
S_2	MSDN-----	KKQQALELALKQIEKQFGKGS	SIMKLGDG-	ADHSIEAIPSGSIALDI	
S_3	M---	AINDTSGKQKALTMVLNQIERSFGKGA	IMRLGDA-	TRMRVETISTGALTLDL	
S_4	M-----	DRQKALEAAVSQIERAFGKGS	SIMKLGKDKQVVETE	VVSTRILGLDV	
S_5	M-----	DE---NKKRALAAALGQIEKQFGKGA	VMRMGDHE-	RQAIPAISTGSLGLDI	
S_6	MD-----	-----KIEKSF	GKGSIMKM	GEE-VVEQVEVIPTGSIALNA	
S_7	M-----	AL-----	IE--FGKG--	M--G-----	L--

- Uns interessieren „möglichst gute“ MSAs
 - Möglichst **wenig Spalten** – wenig Leerzeichen
 - Möglichst **homogene Spalten** – hohe Übereinstimmung

chite	---	ADKPKRPLSAYMLWLN	SARESIKRENPDFK-	VTEVAKKGGELWRGLKD
wheat	--	DPNKPKRAPSAFFVMG	FEFFKQKNPKNKSVA	AVGKAAGERWKSISE
trybr	KKDSNA	PKRAMTSEMFSSDFRS	----	KHSDLS-IVEMSKAAGAAWKELG
mouse	----	KPKRPRSAJNIYVSE	SFQ----	EAKDDS-AQGKCLKLVNEAWKNLSP
		***. :::	: . . .	: . . . * . * : *

Motivation

- Alignment als Maß für Ähnlichkeitssuche sucht **ähnliche Sequenzen** in einer DB
 - Grund: Ähnliche Sequenz – ähnliche Struktur – ähnliche Funktion
- MSA sucht „**das Ähnliche**“ in vielen Sequenzen
 - Das, was viele Sequenzen ähnlich zueinander macht
 - Man startet mit vielen Sequenzen, bei denen man **ähnliche Funktion** / Struktur vermutet
 - MSA stellt fest, was das Gemeinsame dieser Sequenzen ist – Domänen, Motive, Signaturen, Profile, ...
 - These: Dieses Gemeinsame ist **biologisch relevant**

Konservierte Domänen

- Gedankengang

- Gegeben: Proteine S_1, \dots, S_k mit ähnlicher Funktion
- Annahme: Identischer evolutionärer Ursprung S
- S unterliegt Evolution und generiert S_1, \dots, S_k
- Abschnitte in S_i , die trotz Evolution gleich blieben (konserviert sind), müssen für die gemeinsame Funktion wichtig sein

- Alignment – Multiples Alignment

- Gemeinsamkeiten zwischen zwei Sequenzen oft per Zufall
- Gemeinsamkeiten zwischen vielen Sequenzen eher kein Zufall

```
AAC GTG AT T GAC
AACGAGTGC TTTACACGT
```

```
AAC GTG AT T GAC
AACGAGTGC TTTACACGT
GCCG TGC TA GTTG
TTC AGTGGACGTG GTA
G GTGCA TGACC
```

Blöcke, Domänen, Sites

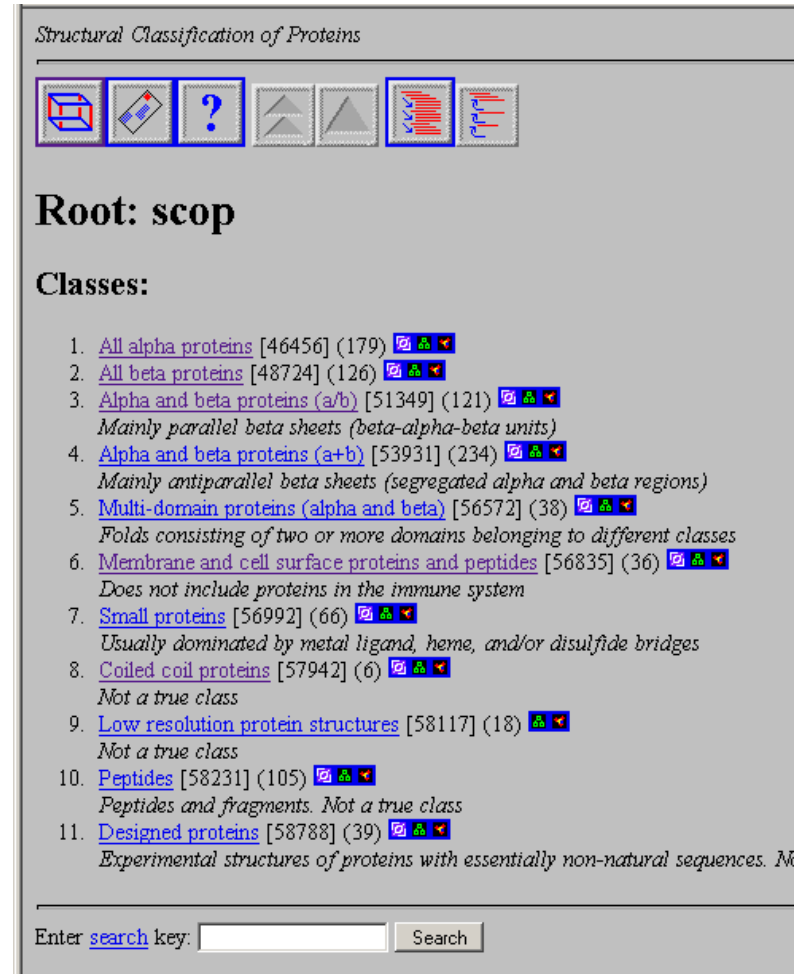
- Proteine
 - Bindungsstellen für andere Proteine / Liganden / Moleküle
 - Bindungsstellen an DNA
 - Signale zur Phosphorylierung / Dephosphorylierung
 - Signale zum Transport des Proteins
 - Signale zum Abbau des Proteins
 - ...
- DNA
 - Bindungsstellen für Transkriptionsfaktoren
 - Start- und Stoppcodons
 - Signal für differenzielles Splicen
 - ...

Proteinfamilien

- Proteine werden in **Familien**, Superfamilien, ... unterteilt
- Diverse Klassifikationen vorhanden (CATH, SCOP, ...)
- Idee: X00.000 Proteine zerfallen in X.000 Klassen ähnlicher Funktion?, Struktur?, Substruktur?
- Finden von Proteinfamilien
 - Starte mit Menge von Proteinen ähnlicher Funktion
 - **Finde das Gemeinsame durch MSA**
 - Suche nur mit konservierten Blöcken nach weiteren Vertretern
 - Modifiziere Familie entsprechend
 - Iteriere, bis Zufriedenheit eintritt

Beispiel: SCOP

- Structural Classification of Proteins
- Hierarchische Anordnung
 - *Fold*: Major structural similarity
 - All Alpha, All Beta, Membrane proteins, ...
 - *Superfamily*: Probable common evolutionary origin
 - Nucleotide-binding domain, Neurotransmitter-gated ion-channel transmembrane pore, ...
 - *Family*: Clear evolutionarily relationship
 - Globins, Death Domain, 4 families of Immunoglobulin ...
 - Protein
 - Spezies



Structural Classification of Proteins

Root: scop

Classes:

1. [All alpha proteins](#) [46456] (179)
2. [All beta proteins](#) [48724] (126)
3. [Alpha and beta proteins \(a/b\)](#) [51349] (121)
Mainly parallel beta sheets (beta-alpha-beta units)
4. [Alpha and beta proteins \(a+b\)](#) [53931] (234)
Mainly antiparallel beta sheets (segregated alpha and beta regions)
5. [Multi-domain proteins \(alpha and beta\)](#) [56572] (38)
Folds consisting of two or more domains belonging to different classes
6. [Membrane and cell surface proteins and peptides](#) [56835] (36)
Does not include proteins in the immune system
7. [Small proteins](#) [56992] (66)
8. [Coiled coil proteins](#) [57942] (6)
Not a true class
9. [Low resolution protein structures](#) [58117] (18)
Not a true class
10. [Peptides](#) [58231] (105)
11. [Designed proteins](#) [58788] (39)
Experimental structures of proteins with essentially non-natural sequences. M

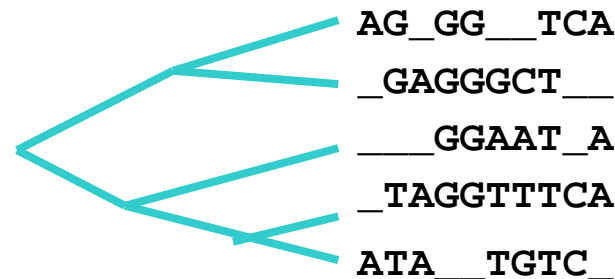
Enter [search](#) key: Search

Inhalt dieser Vorlesung

- Multiples Sequenzalignment
- Sum-Of-Pair Zielfunktion
- Center-Star Zielfunktion

MSA Zielfunktion

- **Zielfunktion** beim einfachen Alignment war klar
 - Möglichst **wenig evolutionäre** Ereignisse – I, R, D
 - Eventuell mit Substitutionsmatrix / Gapmodellen
- Zielfunktion für MSA ist nicht klar
 - Score einer Spalte mit 2 T, zwei G und einem Leerzeichen?
 - Angabe einer Substitutionsmatrix für k Sequenzen über Alphabet Σ würde $O((|\Sigma|+1)^k)$ Werte erfordern
 - Motivation mit evolutionären Ereignissen funktioniert nicht mehr



MSA Überblick

- Es gibt diverse Vorschläge für Zielfunktionen
- Maximiere die Summe aller paarweisen Alignments
- Maximiere die Summe der Alignments jeder Sequenz zu einer Konsensussequenz (Center-Star)
- Maximiere die Summe der Alignments folgend dem phylogenetischen Baum der Sequenzen
 - Allerdings kennen wir die Sequenzen an der Wurzel und an den inneren Knoten nicht

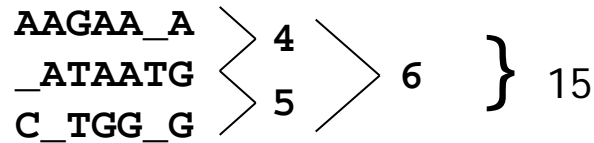
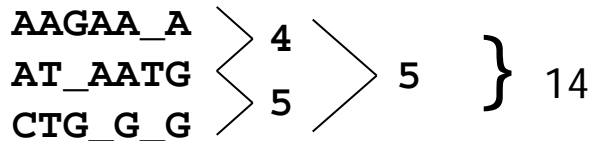
Formal

- Definition

- Gegeben ein MSA M für Sequenzen S_1, \dots, S_k . Das durch M *induziertes Alignment* für zwei Sequenzen S_i und S_j ist durch folgende Vorschrift definiert:
 - Entferne aus M alle Zeilen außer i und j
 - Entferne alle Spalten, die in i und j je ein Leerzeichen enthalten
- Gegeben ein MSA M für Sequenzen S_1, \dots, S_k . Der *Sum-Of-Pairs Score für M (SP-Score)* ist die Summe aller Alignmentsscores der durch M induzierten paarweisen Alignments
- Das *SP-Alignment Problem für Sequenzen S_1, \dots, S_k* sucht das MSA M mit minimalem SP-Score

Beispiel

d/i	=	1
r	=	1
m	=	0



- Die Berechnung des SP-Scores für ein gegebenes MSA über k Sequenzen ist also einfach
 - Komplexität?

Beispiel

d/i	=	1
r	=	1
m	=	0

AAGAA_A } 4 } 5 } 14
AT_AATG } 5 }
CTG_G_G }

AAGAA_A } 4 } 6 } 15
_ATAATG } 5 }
C_TGG_G }

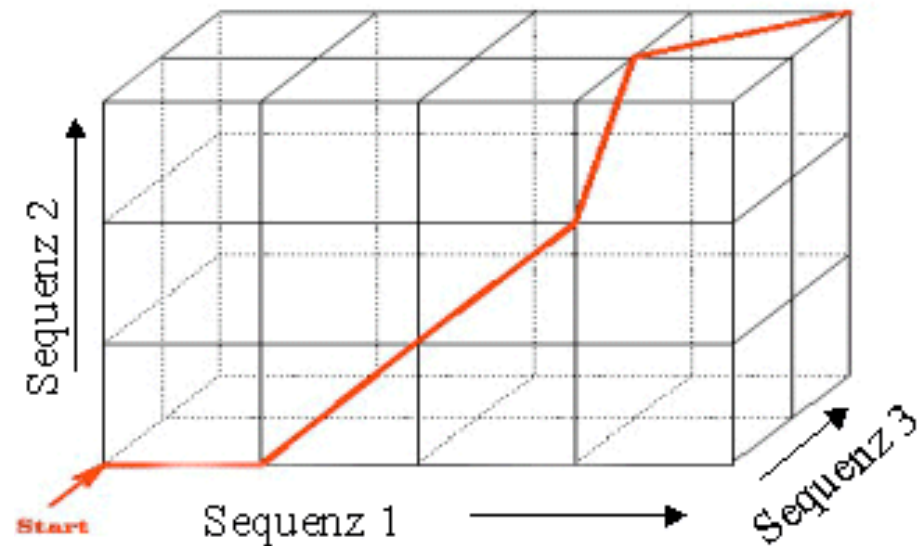
- Die Berechnung des SP-Scores für ein gegebenes MSA über k Sequenzen der Länge l ist also einfach
 - Komplexität $O(k^2 \cdot l)$
- Aber wie findet man das MSA mit **minimalem SP-Score**?

DP in k Dimensionen

k=2: 2-dimensionale Matrix

	0	1	2	3	4	5	6	7
		w	r	i	t	e	r	s
0	0	1	2	3	4	5	6	7
1	v	1	1	2	3	4	5	6
2	i	2	2	2	2	3	4	5
3	n	3	3	3	3	3	4	5
4	t	4	4	4	4	3	4	5
5	n	5	5	5	5	4	4	5
6	e	6	6	6	6	5	4	5
7	r	7	7	6	7	6	5	4

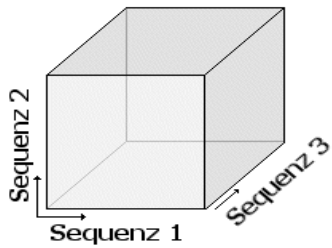
k=3: 3-dimensionale Matrix



- Was kostet ein Schritt?
- Im SoP sind die Kosten einfach zu finden – Paare vergleichen

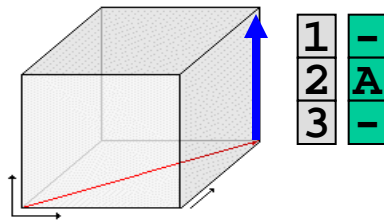
Erinnerung

- Grundidee der DP für zwei Sequenzen S_1, S_2
 - Berechnung des Alignment $d(i,j)$ von $S_1[1..i]$ und $S_2[1..j]$ für steigende Werte (i, j) bis $i=|S_1|$ und $j=|S_2|$
 - Berechnung von $d(i,j)$ aus $d(i-1,j-1), d(i,j-1), d(i-1,j)$
 - Man verlängert $d(i-1,j-1)$ um Match oder Mismatch
 - ... oder man verlängert $d(i,j-1)$ um ein Insert
 - ... oder man verlängert $d(i-1,j)$ um eine Deletion
 - Statische Initialisierung der Werte $d(i,0)$ und $d(0,j)$
- Wir betrachten im Folgenden nur den Fall $k=3$



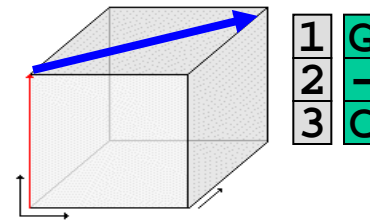
Analogie

$d(i, j-1, k)$

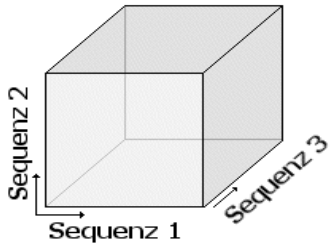


- SP-Alignment von $d(i, j-1, k)$ ist bekannt
- Wir erweitern zu $d(i, j, k)$
- Dazu alignieren wir bei dieser Option $S_2[j]$ zweimal mit Leerzeichen (Inserts)

$d(i-1, j, k-1)$

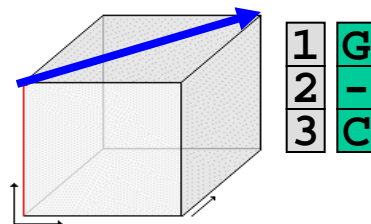
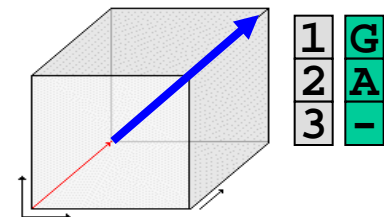
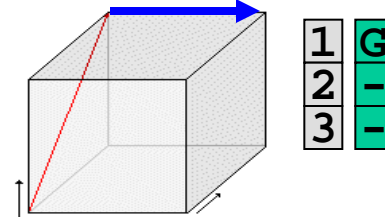
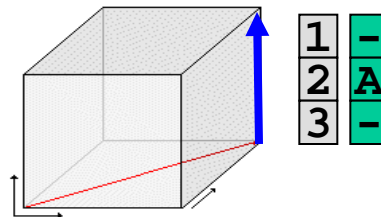
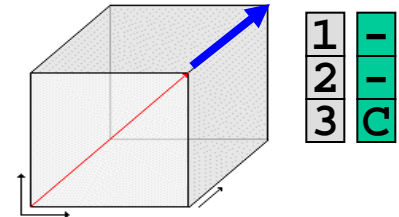
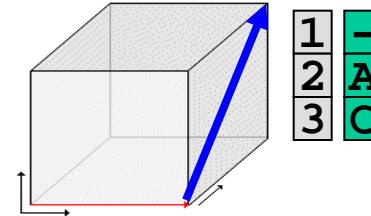
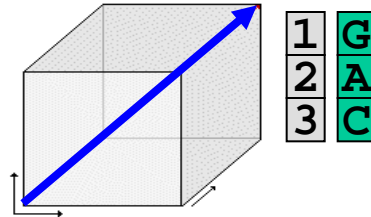


- SP-Alignment von $d(i-1, j, k-1)$ ist bekannt
- Wir erweitern zu $d(i, j, k)$
- Dazu alignieren wir hier ein Leerzeichen mit $S_1[i-1]$ und mit $S_3[k-1]$



Mögliche Schritte zu $d(i,j,k)$

- $d(i-1, j-1, k-1)$
- $d(i, j-1, k-1)$
- $d(i, j, k-1)$
- $d(i, j-1, k)$
- $d(i-1, j, k)$
- $d(i-1, j-1, k)$
- $d(i-1, j, k-1)$



Formal

- Einfaches Kostenmodell (I/D/R=1, M=0)
- Theorem
 - Gegeben Sequenzen S_1, S_2, S_3 .
 - Sei $d(i,j,k)$ der Score des SP-optimalen Alignments der Strings $S_1[1..i], S_2[1..j], S_3[1..k]$
 - Sei $c_{ij} = 0$, wenn $S_1[i] = S_2[j]$, sonst 1
 - Sei $c_{ik} = 0$, wenn $S_1[i] = S_3[k]$, sonst 1
 - Sei $c_{jk} = 0$, wenn $S_2[j] = S_3[k]$, sonst 1
 - Dann berechnet sich $d(i,j,k)$ als:

$$d(i, j, k) = \min \left\{ \begin{array}{lll} d(i-1, j-1, k-1) + c_{ij} & + c_{ik} & + c_{jk} \\ d(i-1, j-1, k) & + c_{ij} & + 2 \\ d(i-1, j, k-1) & + c_{ik} & + 2 \\ d(i, j-1, k-1) & + c_{jk} & + 2 \\ d(i-1, j, k) & & + 2 \\ d(i, j-1, k) & & + 2 \\ d(i, j, k-1) & & + 2 \end{array} \right.$$

Randbedingungen

- Theorem Fortsetzung

- ...

- mit *Initialisierung*

- Sei $D_{a,b}(i,j)$ der optimale Alignmentsscore von $S_a[1..i]$ mit $S_b[1..j]$

- $D(0, 0, 0) = 0$

- $D(i, j, 0) = D_{1,2}(i, j) + (i+j)$

- $D(i, 0, k) = D_{1,3}(i, k) + (i+k)$

- $D(0, j, k) = D_{2,3}(j, k) + (j+k)$

- Bemerkung

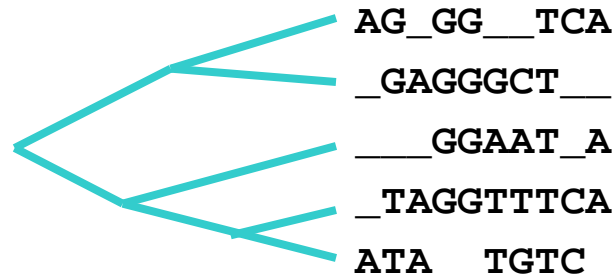
- Alignment eines Leerzeichen mit einem Leerzeichen ist im *induzierten* paarweisen Alignment nicht enthalten

Komplexität

- Für drei Sequenzen der Länge n
 - Würfel hat n^3 Zellen
 - Für jede Zelle sind 7 Berechnungen notwendig
 - Zusammen $O(7 \cdot n^3)$
- Allgemeiner Fall: k Sequenzen der Länge n
 - Hyperwürfel hat n^k Zellen
 - Für jede Zelle sind $2^k - 1$ Vorgänger zu beachten
 - Alle Ecken eines k -dimensionalen Würfels minus eins (Das ist die Ecke die gerade berechnet wird)
 - Zusammen $O(2^k \cdot n^k)$
 - Eigentlich: $O(2^k \cdot n^k \cdot k)$ wg Berechnung des Scores jeder Spalte
 - Häufigkeit aller Zeichen in Spalte in $O(k)$ zählen und Score in $O(|\Sigma|^2)$ ausrechnen
- Tatsächlich: *Das **SP-Alignment Problem** ist NP-vollständig*

MSA unlösbar?

- SP-Score für mehr als eine Handvoll Sequenzen nennenswerter Länge nicht berechenbar
- Aber: SP berechnet gar **nicht die minimale „Menge“ an Evolution**



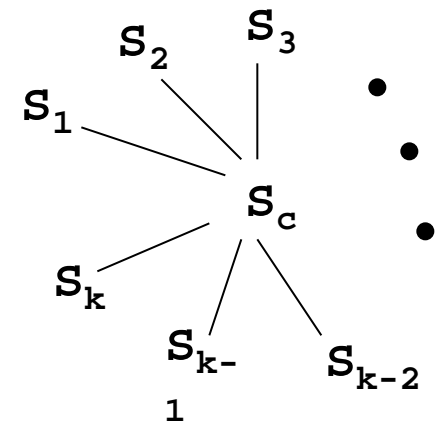
- Andere Zielfunktionen: Center-Star, MSA entlang des phylogenetischen Baums
- Praxis: **Heuristiken** (iterative, lokal-Greedy, ...)

Inhalt dieser Vorlesung

- Multiples Sequenzalignment
- Sum-Of-Pair Zielfunktion
- Center-Star Zielfunktion

MSA mit Konsensussequenz

- Minimiere Summe der **Alignmentscores aller Sequenzen** S_1, \dots, S_k mit einer **Konsensussequenz** S_c
 - S_c kann eine der S_i sein, muss aber nicht
 - Konstruktion von S_c z.B. durch Untereinanderschreiben der S_i ohne Gaps und Wahl des häufigsten Buchstaben (Consensus)
 - Funktioniert nur, wenn alle S_i ungefähr gleich lang – typisch bei MSA
 - MSA wird aus dem „Star“ abgeleitet
- Bei geeigneter Wahl von S_c gilt
 - Der SP Score des berechneten MSA ist **höchstens doppelt so hoch** wie der SP-optimale
 - Center-Star approximiert Sum-of-Pairs also bis auf Faktor 2



Center-Star Verfahren

- Gegeben k Sequenzen der Länge n
- Wähle als Konsensus S_c die Sequenz, die den **kleinsten durchschnittlichen Abstand** zu allen Sequenzen hat
- Beginne: $M = S_c$
- Iteriere
 - Wähle eine noch nicht alignierte Sequenz S beliebig
 - **Aligniere M und S**
 - Genaue Methode führen wir nicht aus (Gusfield, p. 347-)
 - Im Kern aligniert man S mit S_c und fügt dadurch entstehende neue Leerzeichen in S_c auch in M ein
 - Bis alle Sequenzen in M enthalten sind
- **Progressiv**: Sukzessives Hinzufügen von einzelnen Sequenzen zu einem wachsenden MSA

Beispiel

1. ATGGC
2. AGCC
3. TGCGAT
4. GCATG
5. TGCCTA
6. CAACTA

	S1	S2	S3	S4	S5	S6
S1	0	2	4	4	4	5
S2	2	0	4	4	3	4
S3	4	4	0	3	3	5
S4	4	4	3	0	3	4
S5	4	3	3	3	0	3
S6	5	4	5	4	3	0
Durchschnitt	3,8	3,4	3,8	3,6	3,2	4,2

Quelle: Martin Filip, Proseminar, 2005

Beispiel 2

- Kern des MSA: $S_5 = \mathbf{TGCCTA}$

- Wähle Sequenz: $S_3 = \mathbf{TGCGAT}$

- Alignment: $\begin{array}{l} \mathbf{TGCC_TA} \\ \mathbf{TGCGAT_} \end{array}$

- Wähle Sequenz: $S_2 = \mathbf{AGCC}$

- Alignment: $\begin{array}{l} \mathbf{TGCC_TA} \\ \mathbf{TGCGAT_} \\ \mathbf{AGCC____} \end{array}$

- Wähle Sequenz: $S_1 = \mathbf{ATGGC}$

- Alignment: $\begin{array}{l} __\mathbf{TGCC_TA} \\ __\mathbf{TGCGAT_} \\ __\mathbf{AGCC____} \\ \mathbf{ATGGC____} \end{array}$

1. **ATGGC**
2. **AGCC**
3. **TGCGAT**
4. **GCATG**
5. **TGCCTA**
6. **CAACTA**

Beispiel 3

- Wähle Sequenz: S4= **GCATG**

– Alignment:

```
  _TGCC_TA
  _TGCGAT_
  _AGCC___
  ATGGC___
  __GC_ATG
```

- Wähle Sequenz: S6= **CAACTA**

– Alignment:

```
  _TGCC_TA
  _TGCGAT_
  _AGCC___
  ATGGC___
  __GC_ATG
  CAAC__TA
```

Approximationsgüte

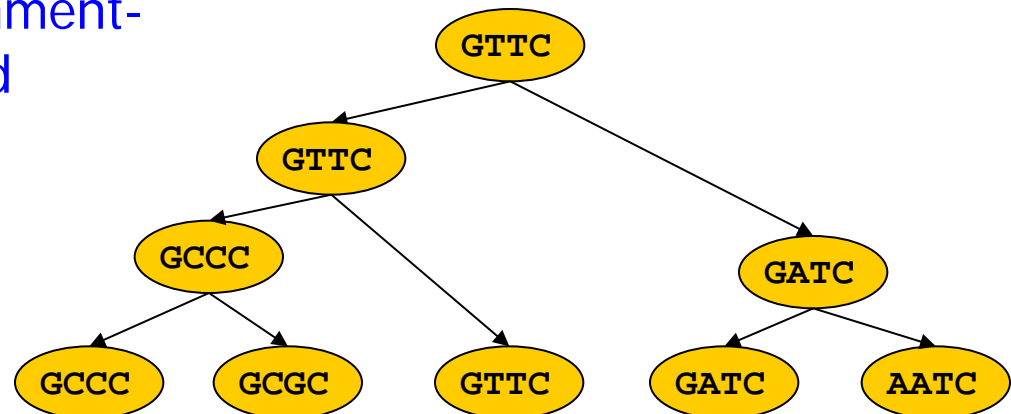
- Wenn für Zeichen i, j, k gilt: $d(i, j) + d(j, k) \geq d(i, k)$
 - Dreiecksungleichung
 - Gilt nicht für viele relevante Substitutionsmatrizen
- Theorem
 - *Gegeben k Sequenzen, $k \geq 3$. Sei d der SP Score des durch den Center-Star Algorithmus berechneten MSA. Sei d^* der optimale Sum-of-Pairs Score. Dann gilt*

$$\frac{d}{d^*} \leq 2 - \frac{1}{k} < 2$$

- Beweis
 - Siehe Gusfield

MSA mit phylogenetischen Bäumen

- These: Sequenzen sind aus einer Ursequenz entstanden
- Natürliche Zielfunktion
 - Suche den Baum T so, dass die **Summe aller Alignmentsscores von benachbarten Sequenzen in T minimiert** wird
 - Aus T kann man ein MSA ableiten (später)
- Leider
 - Wir kennen die inneren Sequenzen nicht
 - Das **phylogenetische Alignmentproblem ist MAX-SNP-hard**
 - Siehe Maximum Parsimony



Suche mit MSA

- Erinnerung: Erzeugung von Proteinfamilien
 - Starte mit Proteinen gleicher/ähnlicher Funktion
 - Finde das Gemeinsame durch MSA
 - Suche mit dem MSA nach weiteren Vertretern
 - Modifiziere Familie entsprechend
 - Iteriere, bis Zufriedenheit eintritt
- Wie sucht man mit einem MSA?
 - Wir müssen entscheiden, wie gut eine (neue) Sequenz S zu einem gegebenen MSA M passt
 - Verschiedene Möglichkeiten: Profile, RegExp, Profile-HMM
 - Nächste Stunde

Selbsttest

- Welche biologische Frage versucht man mit MSA zu beantworten?
- Was ist eine Protein-Domäne? Beispiele? Was wäre eine Entsprechung bei Genen?
- Warum ist SP nicht das sinnvollste MSA-Verfahren?
- Wie könnte ein lokales MSA-Problem definiert sein? Wann könnte so ein Maß sinnvoll sein?
- Erklären Sie den Algorithmus, um das SP-optimale MSA für k Sequenzen zu finden