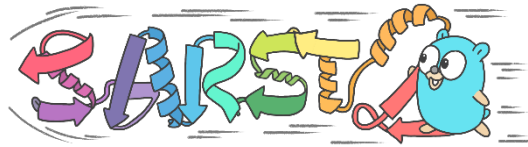


# SARST2 Benutzerhandbuch

## Deutsche Version



## Inhaltsverzeichnis

|   |    |
|---|----|
| Inhaltsverzeichnis .....  | 1  |
| 1. Über diese Software .....                                      | 2  |
| 2. Download, Dekomprimierung und Installation .....               | 2  |
| 3. Inhalt der Software-Verzeichnisse .....                        | 3  |
| 4. Kurzanleitung.....   | 3  |
| 5. Handbuch: sarst2.....  | 4  |
| 5.1 Verwendung .....  | 4  |
| 5.2 Programmooptionen.....  | 4  |
| 5.3 Beispiel: Datenbanksuche nach struktureller Ähnlichkeit ..... | 7  |
| 5.4 Beispielverwendung: One-against-all-Ausrichtungen .....       | 7  |
| 5.5 Beispielverwendung: Paarweises Strukturausrichten .....       | 7  |
| 5.6 Ausgabe von überlagerten Proteinstrukturen.....               | 8  |
| 5.7 Interaktive HTML-Ergebnisseiten erzeugen.....                 | 8  |
| 6. Handbuch: formatdb .....                                       | 10 |
| 6.1 Verwendung .....  | 10 |
| 6.2 Programmooptionen.....  | 10 |
| 6.3 Beispielverwendungen .....                                    | 11 |
| 7. Handbuch: readdb .....   | 12 |
| 7.1 Verwendung .....  | 12 |
| 7.2 Programmooptionen.....  | 12 |
| 7.3 Beispielverwendungen .....                                    | 12 |

## 1. Über diese Software

SARST2 (Structural Similarity Search Aided by Ramachandran Sequential Transformation, Version 2) ist ein Hochleistungsalgorithmus zum Abgleich von Proteinstrukturen. Das Programm unterstützt sowohl Datenbanksuchen nach struktureller Ähnlichkeit anhand eines angegebenen Abfrageproteins (Query) als auch paarweise Strukturvergleiche zwischen zwei Proteinstrukturen.

Diese Software wird zusammen mit der folgenden wissenschaftlichen Publikation veröffentlicht und regelmäßig über die in der Publikation angegebenen URLs aktualisiert:

|          |   |
|----------|---|
| Titel    | SARST2, a high-throughput protein structure alignment algorithm for searching massive databases                               |
| Autoren  | Wei-Cheng Lo*, Arie Warshel, Chia-Hua Lo, Chia Yee Choke, Yan-Jie Li, Shih-Chung Yen, Jyun-Yi Yang and Shih-Wen Weng          |
| Institut | Institute of Bioinformatics and Systems Biology, National Yang Ming Chiao Tung University, Hsinchu, Taiwan, Republic of China |

\*Korrespondenzautor

Download-URLs:

<https://github.com/NYCU-10lab/sarst>

<https://10lab.ceb.nycu.edu.tw/sarst2>

## 2. Download, Dekomprimierung und Installation

Die neueste Version des SARST2-Programms sowie die vorgefertigten Ziel-Datenbanken stehen unter den oben aufgeführten URLs zum Download bereit.

Nach dem Herunterladen eines komprimierten Archivs können Sie die Dateien je nach Archivformat mit den Programmen tar, gzip oder zip extrahieren. Nach dem Entpacken liegen die SARST2-Programmdateien als vorkompilierte, ausführbare Binärdateien vor und erfordern keine separate Installation.

Beispiel: Wenn Sie das Archiv SARST2-v2.0.26-Linux.x86\_64.tar.gz herunterladen, können Sie es unter Linux mit folgendem Befehl entpacken:

```
tar xfp SARST2-v2.0.26-Linux.x86_64.tar.gz
```

Nach dem Entpacken können Sie das SARST2-Programm mit den folgenden Befehlen ausführen:

```
cd SARST2-v2.0.26-Linux.x86_64/bin
chmod +x sarst2
./sarst2 -h
```

### 3. Inhalt der Software-Verzeichnisse

|     |   |
|-----|---|
| bin | 64-Bit-Programme für Linux, Windows und macOS.          |
| dbs | Vorgefertigte Ziel-Datenbanken: PDB-2022 und SCOP-2.07. |
| doc | Benutzerhandbücher in mehreren Sprachen.                |

### 4. Kurzanleitung

Dieses Softwarepaket enthält drei Hauptprogramme, die im Folgenden beschrieben werden:

|          |   |
|----------|---|
| sarst2   | Implementierung des SARST2-Algorithmus zum Abgleich von Proteinstrukturen.  |
| formatdb | Ein Tool zur Datenbankformatierung, mit dem Benutzer eigene Ziel-Datenbanken für Strukturrecherche und -vergleich erstellen können.             |
| readdb   | Ein Tool zum Extrahieren der Aminosäuresequenzen und linear codierten Struktursequenzen von Proteinen aus einer vorformatierten Ziel-Datenbank. |

Das Ausführen eines der oben genannten Programme ohne Parameter (oder mit -h) zeigt eine kurze textbasierte Hilfeseite für das jeweilige Programm an.

#### Linux, macOS

```
./sarst2
./formatdb
./readdb
./sarst2 -h
./formatdb -h
./readdb -h
```

#### Windows (Cmd oder PowerShell)

```
.\sarst2.exe
.\formatdb.exe
.\readdb.exe
.\sarst2.exe -h
.\formatdb.exe -h
.\readdb.exe -h
```

## 5. Handbuch: sarst2

### 5.1 Verwendung

```
./sarst2  Abfrage-Struktur  Zielstruktur(en)  [Optionen]
          -----
          > eine PDB-/CIF-Datei > Zielstruktur(en) kann sein:
                                   1. -db + eine vorformatierte Datenbank
                                   2. eine Liste von PDB-/CIF-Dateien
                                   3. Verzeichnisse mit PDB-/CIF-Dateien
                                   4. eine einzelne PDB-/CIF-Datei (Paarvergleich)
```

### 5.2 Programmoptionen

|          |       |  |
|----------|-------|--|
| -db      | [str] | Ziel-Datenbank der zu durchsuchenden Proteinstrukturen.<br>(Standard: none)  |
| -brief   | [int] | Anzahl der Zielstrukturen, für die eine einzeilige Zusammenfassung angezeigt wird.<br>(Standard: 500)  |
| -detail  | [int] | Anzahl der Zielstrukturen, für die detaillierte Ausrichtungsdaten angezeigt werden.<br>(Standard: 500)   |
| -t       | [int] | Anzahl der Threads.<br>Muss $\geq 0$ sein; bei 0 werden alle Prozessoren verwendet.<br>(Standard: 0, alle Prozessoren)   |
| -w       | [int] | Word-Größe.<br>(Standard: 5)   |
| -orderby | [int] | Sortierung der Trefferliste nach einem der folgenden Faktoren: 1: Conf-Score--, 2: TM-Score--, 3: Sequenzidentität-- oder 4: RMSD++, wobei -/++ für absteigende/aufsteigende Reihenfolge steht<br>(Standard: 1, Conf-score--)<br>(bei paarweisem Abgleich ignoriert) |
| -mode    | [int] | Suchmodus, 1: genau, 2: ausgewogen, 3: schnell oder sonst: auto.<br>(Standard: auto)   |
| -f       | [int] | Minimale Filter aktivieren, 0: aus, 1: ein, oder sonst: auto.<br>(Standard: auto)<br>(bei paarweisem Abgleich immer 0, deaktiviert)  |

|        |         |   |
|--------|---------|---|
| -C     | [float] | Schwellenwert für den Conf-Score (Confidence Score).<br>Muss zwischen 0 und 1 liegen; bei 0 wird kein Schwellenwert angewendet.<br>(Standard: 0.5)<br>(bei paarweisem Abgleich immer deaktiviert)   |
| -pC    | [float] | Schwellenwert für den finalen pC-Wert, d. h. $-\log_2(C)$ .<br>Muss $\geq 0$ sein; bei 0 wird kein Schwellenwert angewendet.<br>(Standard: 1.0, entspricht -C = 0.5)<br>(bei paarweisem Abgleich immer deaktiviert)   |
| -e     | [float] | Schwellenwert für den pC-Wert, angewendet in jedem Filter- und Verfeinerungsschritt. Bei gleichem -e und -pC verwirft -e mehr irrelevante Treffer.<br>Muss $\geq 0$ sein; bei 0 wird kein Schwellenwert angewendet.<br>(Standard: 1.0)<br>(bei paarweisem Abgleich immer deaktiviert) |
| -tmcut | [float] | TM-Score-Schwellenwert.<br>Muss $\geq 0$ sein; bei 0 wird kein Schwellenwert angewendet.<br>Ein TM-Score $\geq 0.7$ gemäß SARST2 deutet möglicherweise auf eine Homologie auf Ebene einer Protein-Familie hin.<br>(Standard: 0.15)<br>(bei paarweisem Abgleich immer deaktiviert)     |
| -mem   | [T/F]   | Puffert alle Daten der Zielproteine im Speicher.<br>(Standard: T)<br>(bei paarweisem Abgleich immer T, aktiviert)   |
| -q     | [T/F]   | Schneller Ausgabestil.<br>Gibt die Ergebnisse in einem vereinfachten und parserfreundlichen Format aus.<br>(Standard: F)  |
| -sa    | [T/F]   | Gibt die strukturbasierte Sequenzausrichtung aus.<br>(Standard: T)  |
| -mat   | [T/F]   | Gibt die Transformationsmatrix für die Überlagerung aus.<br>(Standard: F)   |
| -nmsbj | [T/F]   | Normalisiert den TM-Score basierend auf der Größe der Zielstruktur.<br>(Standard: F)  |
| -nmavg | [T/F]   | Normalisiert den TM-Score basierend auf der durchschnittlichen Größe der Abfrage- und Zielstruktur.<br>(Standard: F)  |
| -nmusr | [float] | Protein-Größe zur Normalisierung des TM-Scores.<br>Sie sollte $\geq$ der minimalen Größe der beiden Strukturen sein; andernfalls kann der TM-Score $> 1$ werden.  |

|           |         |  |
|-----------|---------|--|
| -d        | [float] | d <sub>0</sub> -Wert zur Skalierung des TM-Scores, z. B. 5,0 Ångström (Å).   |
| -ml       | [T/F]   | Maschinelles Lernen anwenden.<br>(Standard: T)<br>(bei paarweisem Abgleich immer F, deaktiviert)   |
| -fdp      | [str]   | Algorithmus für dynamische Programmierung in den Filter-Schritten. Unterstützte Optionen: NW (Needleman-Wunsch), SW (Smith-Waterman)<br>(Standard: NW)   |
| -rdp      | [str]   | Algorithmus für dynamische Programmierung im Verfeinerungsschritt. Unterstützte Optionen: NW (Needleman-Wunsch), SW (Smith-Waterman)<br>(Standard: NW)   |
| -swp      | [str]   | Pfad zur benutzerdefinierten Swap-Datei. Die Verwendung einer Swap-Datei kann den Speicherbedarf reduzieren.<br>(Standard: none)   |
| -Sout     | [str]   | Ordner für die Ausgabe der Überlagerungsdateien (Structure Superimposition).<br>Der Ordner wird erstellt, falls er nicht existiert. Die Anzahl der Überlagerungsdateien wird durch die Option -detail begrenzt<br>(Standard: none)       |
| -html     | [str]   | Erstellung eines HTML-Ausgabeordners.<br>Überlagerungsdateien der Strukturen werden ebenfalls im HTML-Ordner erzeugt.<br>(Standard: none)  |
| -jsmol    | [str]   | Pfad zum JSmol-JavaScript-Paket für die Anzeige von überlagerten Strukturen im HTML-Ausgabedokument.<br>Kann ein lokaler Ordner oder eine HTTP(S)-URL sein.<br>(Standard: none)<br>(Test-URL: "https://10lab.ceb.nycu.edu.tw/ext/jsmol") |
| -pssm_out | [str]   | Datei zum Speichern der PSSMs der Struktur- und Sequenzcodes, die in diesem Algorithmus verwendet werden.<br>(Standard: none)  |
| -pssm_pC  | [float] | Grenzwert des pC-Werts für die PSSM-Erstellung.<br>(Standard: 0.05)  |
| -h        |         | Die Hilfmeldung anzeigen (Kurzanleitung).  |

## 5.3 Beispiel: Datenbanksuche nach struktureller Ähnlichkeit

Die Abfragestruktur gegen eine vorformatierte Ziel-Datenbank durchsuchen

```
./sarst2 Qry.pdb -db my_db/my_proteins.db -brief 10 -w 7 -e 0.1  
./sarst2 Qry.pdb -db my_db/my_proteins.db -brief 10 -d 5.0 -sa F
```

Im obigen Beispiel befindet sich die Ziel-Datenbank im Ordner `my_db`, und `my_proteins.db` ist der Dateistamm der zugehörigen Datenbankdateien. Die tatsächliche Dateiliste der Datenbank finden Sie im Abschnitt **Handbuch: formatdb**.

## 5.4 Beispielverwendung: One-against-all-Ausrichtungen

Die Abfragestruktur gegen aufgelistete Zielstrukturen durchsuchen

```
./sarst2 Qry.pdb Sbj1.pdb Sbj2.cif Sbj3.pdb -mat T
```

Die Abfragestruktur gegen Zielstrukturdateien mit Wildcard-Mustern durchsuchen

```
./sarst2 Qry.pdb "set1/*.pdb" "set2/1a???.cif" -nmavg T
```

In diesem Beispiel sind `"set1/*.pdb"` und `"set2/1a???.cif"` in Anführungszeichen gesetzt und enthalten Wildcard-Zeichen. Das Programm `sarst2` erweitert diese Wildcard-Muster automatisch und ermittelt intern die entsprechenden Dateinamen. Wenn die Muster nicht in Anführungszeichen gesetzt sind, übernimmt das Betriebssystem die Erweiterung der Wildcards. Bei einer großen Anzahl passender Dateien kann die resultierende Befehlszeilenargumentliste die Systemgrenzen überschreiten und dazu führen, dass der Befehl fehlschlägt. Daher empfehlen wir, Wildcard-Muster in Anführungszeichen zu setzen, damit `sarst2` die Dateiliste intern verarbeitet, anstatt sich auf das Standardverhalten des Betriebssystems zu verlassen.

Die Abfragestruktur gegen mehrere Ordner mit Zielstrukturen durchsuchen

```
./sarst2 Qry.pdb set1 set2 -nmavg T
```

In diesem Beispiel sind `set1` und `set2` Ordner, die Proteinstrukturdateien enthalten können. Das Programm `sarst2` ruft automatisch alle Dateien in diesen Ordnern ab (entspricht `set1/*` und `set2/*`). Dateien im PDB- oder CIF-Format werden ausgewählt und gegen die Abfragestruktur ausgerichtet.

## 5.5 Beispielverwendung: Paarweises Strukturausrichten

Die Abfragestruktur mit einer Zielstruktur ausrichten

```
./sarst2 Qry.pdb Sbj.pdb -sa F  
./sarst2 Qry.pdb Sbj.cif -mat T
```

## 5.6 Ausgabe von überlagerten Proteinstrukturen

### Generierung von Überlagerungsdateien (PDB) für Abfrage- und Zielstrukturen

```
./sarst2 Qry.pdb -db prot/myDb -detail 100 -Sout output_folder  
./sarst2 Qry.pdb "set1/*.cif" -detail 100 -Sout output_folder
```

Mit der Option `-Sout output_folder` werden die überlagerten Proteinstrukturen im PDB-Format in den vom Benutzer angegebenen Ordner ausgegeben. Die Anzahl der erzeugten Überlagerungsdateien wird durch die Option `-detail` festgelegt. Jede Ausgabedatei wird als `Qry-SbjSN.pdb` benannt, wobei SN die laufende Nummer der Zielstruktur in der Trefferliste darstellt. In jeder Überlagerungsdatei sind die Ketten-IDs der Abfrage- und der Zielstruktur mit Q bzw. S gekennzeichnet. Die beiden Ketten werden durch einen TER-Datensatz getrennt, wie unten dargestellt:

|      |     |    |     |   |     |        |        |         |   |
|------|-----|----|-----|---|-----|--------|--------|---------|---|
| ATOM | 150 | CA | LEU | Q | 150 | 29.000 | -8.400 | 0.800   | C |
| ATOM | 151 | CA | GLY | Q | 151 | 26.000 | -9.600 | 2.600   | C |
| ATOM | 152 | CA | TYR | Q | 152 | 25.400 | -6.800 | 5.000   | C |
| ATOM | 153 | CA | GLN | Q | 153 | 23.600 | -3.800 | 3.600   | C |
| ATOM | 154 | CA | GLY | Q | 154 | 22.800 | -2.800 | 7.200   | C |
| TER  |     |    |     |   |     |        |        |         |   |
| ATOM | 1   | CA | MET | S | 1   | 24.400 | 9.800  | -10.000 | C |
| ATOM | 2   | CA | VAL | S | 2   | 27.200 | 11.800 | -11.400 | C |
| ATOM | 3   | CA | LEU | S | 3   | 28.800 | 15.200 | -10.400 | C |
| ATOM | 4   | CA | SER | S | 4   | 29.800 | 17.800 | -13.000 | C |
| ATOM | 5   | CA | GLU | S | 5   | 33.400 | 19.200 | -13.000 | C |
| ATOM | 6   | CA | GLY | S | 6   | 32.000 | 22.400 | -11.600 | C |

Wie in der Abbildung gezeigt, enthalten die Überlagerungsdateien nur  $\text{C}\alpha$ -Atome (Alpha-Kohlenstoff,  $\text{C}\alpha$ ). Dies liegt daran, dass SARST2 alle Berechnungen ausschließlich auf den  $\text{C}\alpha$ -Koordinaten durchführt. Die Orientierung der Abfragestruktur bleibt in allen generierten Überlagerungsdateien konsistent erhalten, während jede Zielstruktur gemäß dem Alignment zur Abfragestruktur transformiert (rotiert und translatiert) wird, um die Überlagerung zu erzielen.

Zur Visualisierung der überlagerten Strukturen empfehlen wir die Verwendung von RasMol oder RasWin (<http://www.openrasmol.org/>). Da die Überlagerungsdateien ausschließlich  $\text{C}\alpha$ -Atome enthalten, sollte der Anzeigemodus in RasMol auf "backbone" eingestellt werden.

## 5.7 Interaktive HTML-Ergebnisseiten erzeugen

### Ein HTML-Dokument mit Online-JSmol-Skripten erzeugen

```
./sarst2 Qry.pdb -db my_db/my_proteins.db -html output_folder  
-jsmol "https://101ab.ceb.nycu.edu.tw/ext/jsmol"
```

### Ein HTML-Dokument mit lokalem JSmol-Skriptordner (Windows) erzeugen

```
./sarst2 Qry.pdb "set1/*.cif" -detail 100 -html output_folder  
-jsmol "file:///D:/bioinfo/software/jsmol"
```



Mit der Option "-html output\_folder" werden ein HTML-Ergebnisdokument sowie die zugehörigen überlagerten Proteinstrukturdateien im angegebenen Ordner erzeugt. Die Option "-html" muss zusammen mit "-jsmol" verwendet werden, wobei "-jsmol" die URL oder den lokalen Pfad des JSmol-Pakets (Version 2013) angibt. JSmol ist ein interaktiver 3D-Viewer für molekulare Strukturen, der in Webbrowsern ausgeführt wird und von den meisten modernen Browsern unterstützt wird.

**SARST2: Structural similarity search Aided by Ramachandran Sequential Transformation**  
(v2.0.16, build 20231120)

Please cite our paper(s) should you find this software helpful:  
1. Wei-Cheng Lo, et al. (2007). BMC Bioinformatics, 8:307  
2. Wei-Cheng Lo, et al. (2009). Nucleic Acids Research, 37:W545-551

**## Query: dbs/strucFiles/101m.pdb**  
Protein size = 154 residues

Subjects:  
Dataset size = 713,183 structure(s)  
Database = /SARST/PDB-2022.db

<< Page 1 Go >>

See [Jmol](#) [JSmol](#) if the structure cannot be displayed.  
The structure will be purely in red for perfect superposition.

**Query vs Subject 012**

JSmol

| Ranking | Protein | Size | Ident(%) | ASize | RMSD(Å) | TM-score | Conf   |
|---------|---------|------|----------|-------|---------|----------|--------|
| 1       | 101nA   | 154  | 100.00   | 154   | 0.0000  | 1.0000   | 1.0000 |
| 2       | 1nynA   | 154  | 100.00   | 154   | 0.1504  | 0.9989   | 0.9972 |
| 3       | 1n1jA   | 154  | 100.00   | 154   | 0.1619  | 0.9988   | 0.9968 |
| 4       | 1n1kA   | 154  | 100.00   | 154   | 0.1633  | 0.9988   | 0.9964 |
| 5       | 1n1iA   | 154  | 99.35    | 154   | 0.2007  | 0.9981   | 0.9952 |
| 6       | 2nq1A   | 154  | 99.35    | 154   | 0.2246  | 0.9977   | 0.9948 |
| 7       | 110nA   | 154  | 99.35    | 154   | 0.2669  | 0.9967   | 0.9946 |
| 8       | 1nqnA   | 154  | 98.70    | 154   | 0.1764  | 0.9986   | 0.9946 |
| 9       | 2nbnA   | 154  | 99.35    | 154   | 0.2727  | 0.9966   | 0.9946 |
| 10      | 102nA   | 154  | 98.70    | 154   | 0.1508  | 0.9989   | 0.9944 |
| 11      | 1n1rA   | 154  | 98.70    | 154   | 0.2253  | 0.9976   | 0.9943 |
| 12      | 1j52A   | 154  | 99.35    | 154   | 0.2972  | 0.9960   | 0.9941 |
| 13      | 1co9A   | 154  | 98.70    | 154   | 0.2889  | 0.9962   | 0.9939 |

>> 12  
**## Subject: PDB-2022/1j52A.pdb**

Protein length = 154 residues  
Aligned residues = 154 residues  
RMSD = 0.2972 Å  
Seq identity = 99.35 % (153/154)  
Seq similarity = 99.35 % (153/154)  
Confidence-score = 0.9941  
pC-value = 0.0085

TM-score, nQuery = 0.9960 (norm size: 154 res, d0: 4.62 Å)

Query: 1 NVLSGGEWQLVLRVWAKYEADVAGHGQDILRLFKSHPETLEKFDVRVKKLATEAMKASE 60  
          NVLSGGEWQLVLRVWAKYEADVAGHGQDILRLFKSHPETLEKFDVRVKKLATEAMKASE 60  
Subject: 1 NVLSGGEWQLVLRVWAKYEADVAGHGQDILRLFKSHPETLEKFDVRVKKLATEAMKASE 60  
          NVLSGGEWQLVLRVWAKYEADVAGHGQDILRLFKSHPETLEKFDVRVKKLATEAMKASE 60

Query Subject Backbone Reset HTMLS WEBGL

Die Hauptdatei im HTML-Ausgabeordner ist SarstResults.html, die in einem Webbrowser geöffnet werden sollte. Weitere HTML-Dateien sind als Inline-Frames in die Hauptseite eingebettet. Außerdem wird ein Unterordner mit dem Namen "sup" erstellt. Darin werden die Überlagerungsdateien gespeichert, die die Abfrage- und die jeweilige Zielstruktur aus der Trefferliste enthalten.

Je nach Betriebssystem, Browser oder Antivirensoftware müssen Sie möglicherweise die Sicherheitseinstellungen anpassen, um dem Browser das Ausführen von JavaScript und den Zugriff auf die überlagerten Strukturdateien im Ordner "sup" zu erlauben, sodass der JSmol-3D-Viewer korrekt funktioniert.

## 6. Handbuch: formatdb

### 6.1 Verwendung

```
./formatdb      Zielstruktur(en)    -db Datenbank    [Optionen]
-----
> Zielstruktur(en) kann sein:
1. eine Liste von PDB-/CIF-Dateien
2. Verzeichnisse mit PDB-/CIF-Dateien
3. eine einfache Textdatei mit einer Liste von
   PDB-/CIF-Dateipfaden
```

### 6.2 Programmoptionen

|             |       |   |
|-------------|-------|---|
| -db         | [str] | Ziel-Datenbank der zu erstellenden Proteinstrukturen.<br>(Standard: none)   |
| -flist      | [str] | Einfache Textdatei mit einer Liste von PDB-/CIF-Dateipfaden.<br>Dieses Argument kann zusammen mit den üblichen Argumenten<br>für Zielstrukturen verwendet werden.<br>(Standard: none)   |
| -t          | [int] | Anzahl der Threads.<br>Muss $\geq 0$ sein; bei 0 werden alle Prozessoren verwendet.<br>(Standard: 0, alle Prozessoren)  |
| -split      | [int] | Datenbank in Teilmengen aufteilen, wobei jede Teilmenge die in<br>dieser Option angegebene Anzahl von Zielstrukturen enthält. Das<br>Aufteilen der Datenbank hilft dabei zu verhindern, dass die<br>Datenbankdateien die maximal zulässige Dateigröße des<br>Speichermediums überschreiten.<br>(Standard: none) |
| -save_disk  | [T/F] | Die Atomkoordinaten von drei auf eine Nachkommastelle runden,<br>um Speicherplatz zu sparen.<br>(Standard: F)   |
| -keep_order | [T/F] | Die Reihenfolge der in der Datenbank gespeicherten Zielstrukturen<br>beibehalten, wie sie in den Eingabeargumenten angegeben wurde.<br>Das Setzen von T verlangsamt die Erstellung der Datenbank.<br>(Standard: F)  |
| -h          |       | Die Hilfmeldung anzeigen (Kurzanleitung).   |

## 6.3 Beispielverwendungen

### Eine Ziel-Datenbank aus aufgelisteten Zielstrukturdateien erstellen

```
./formatdb Sbj1.pdb Sbj2.cif Sbj3.pdb -db myDb -keep_order T
```

Es werden mehrere Datenbankdateien mit Dateinamen, die mit myDb beginnen, erstellt. Durch Aktivieren von -keep\_order wird die Reihenfolge der Zielstrukturen in der Ziel-Datenbank entsprechend der Reihenfolge in den Befehlszeilenargumenten beibehalten.

### Eine Ziel-Datenbank aus aufgelisteten Zielstrukturdateien mit Wildcards erstellen

```
./formatdb "set1/*.pdb" "set2/*.cif" Sbj1.pdb Sbj2.cif -db myDb
```

Beim Auflisten von Zielstrukturdateien können Argumente mit und ohne Wildcards gemischt werden. Es wird empfohlen, Wildcard-Argumente in Anführungszeichen zu setzen, damit das Programm die Dateierweiterung korrekt verarbeiten kann.

### Eine Ziel-Datenbank basierend auf einer Dateiliste von Zielstrukturen erstellen

```
./formatdb -flist protlist.txt Sbj1.pdb Sbj2.cif -db myDb
```

Die Datei protlist.txt sollte eine Liste von Dateipfaden enthalten, wobei pro Zeile genau ein Dateipfad angegeben wird.

### Eine Ziel-Datenbank aus Ordnern mit Zielstrukturdateien erstellen

```
./formatdb folder1 folder2 -db myDb -save_disk T -split 50000
```

Durch Aktivieren von -save\_disk werden die C $\alpha$ -Koordinaten auf eine Nachkommastelle gerundet, um Speicherplatz zu sparen. Die Option "-split 50000" führt dazu, dass mehrere Teilmengen-Datenbanken erstellt werden, von denen jede maximal 50000 Strukturen enthält. Diese Aufteilungsoption ist besonders nützlich, wenn die Größe der formatierten Datenbankdateien die maximal zulässige Dateigröße bestimmter Betriebssysteme oder Dateisysteme überschreiten könnte.

## 7. Handbuch: readdb

### 7.1 Verwendung

| <code>./readdb</code> | <u>Zieldatenbank</u> | <u>Ausgabedatei</u> | <u><code>[-seq Sequenztyp]</code></u> |
|-----------------------|----------------------|---------------------|---------------------------------------|
|                       | > Muss sein:         | > Wird im           | > Kann sein:                          |
|                       | vorformatierte       | FASTA-Format        | 1. AA                                 |
|                       | SARST2-Datenbank     | ausgegeben          | 2. AAT                                |
|                       |                      |                     | 3. SARST                              |
|                       |                      |                     | 4. Vier-Symbol-SSE                    |

### 7.2 Programmoptionen

|                   |                    |  |   |
|-------------------|--------------------|--|---|
| <code>-seq</code> | <code>[str]</code> | Der Ausgabesequenztyp.                     |   |
|                   |                    | AA   | Aminosäuresequenz   |
|                   |                    | AAT  | Fünfsymbolige Aminosäuretyp-Sequenz                           |
|                   |                    | SARST                                      | SARST-Ramachandran-Code-Sequenz                               |
|                   |                    | SSE  | Vier-Symbol-Sekundärstrukturelement-Sequenz<br>(Standard: AA) |
| <code>-h</code>   |                    | Die Hilfemeldung anzeigen (Kurzanleitung). |   |

### 7.3 Beispielverwendungen

Zielsequenzen aus einer SARST2-Zieldatenbank extrahieren

```
./readdb my_db/my_proteins.db seqs.fasta
./readdb my_db/my_proteins.db seqs.fasta -seq SARST
./readdb my_db/my_proteins.db seqs.fasta -seq AAT
```

Wenn `-seq` nicht angegeben ist, wird standardmäßig die Aminosäuresequenz als Ausgabetyt verwendet. Die Ausgabedatei wird (`seqs.fasta`) im FASTA-Format erstellt. Falls die Ausgabedatei bereits vor dem Ausführen von `readdb` existiert, wird sie überschrieben.