**Bioinformatik**

**Stefanie Hartmann**

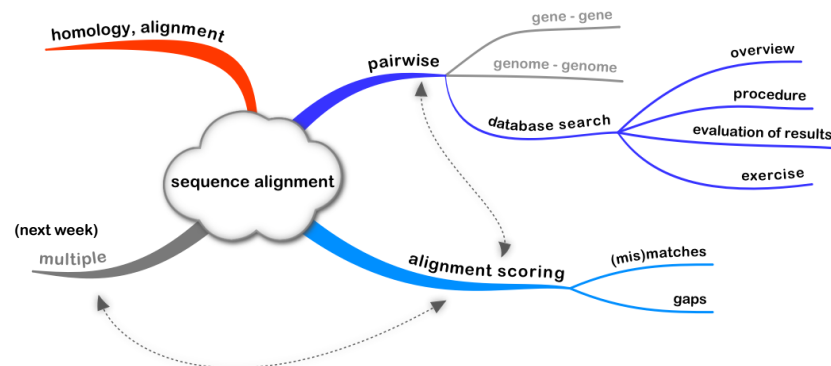Wintersemester 2019 / 2020, Universität Potsdam

# Database searches
# Nov 15, 2019

---

# heute!

10h: 3 Computer Pools (1a, 1b, 2a)

11: 2 Computer Pools (1a, 2a)

---

# Today's topics

---

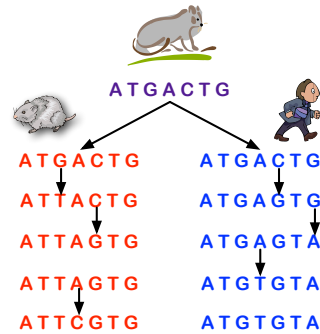# Analysis & comparison of sequences

sequence data: abundant & relevant for many disciplines!

- **database searching (pairwise alignments)**
  - identify / collect **homologous** sequences or domains

- **multiple sequence alignment**
  - study conservation & variability of **homologous** sequences

- **phylogenetic analysis**
  - infer evolutionary history of **homologous** sequences

- **genome alignment and assembly**

- **structure analysis & prediction**

- …

# Homology

**homologs**

- **morphological structures, sequences, domains, ....**
- **derived from a common ancestor (and evolving!)**
- **all-or-nothing condition**

ATGACTG

ATGACTG          ATGACTG
ATTACTG          ATGAGTG
ATTAGTG          ATGAGTA
ATTAGTG          ATGTGTA
ATTCGTG          ATGTGTA

**similarity / identity**

- **a quantitative measure**

ATTCGTG
ATGTGTA

---

# Sequence alignment

**the comparison & relative arrangement of sequences by**

- **searching for similarities between their characters (amino acids or nucleotides) and**
- **possibly inserting gaps in each sequence**

**pairwise sequence alignment**

MYITENGMDEFNNPKVSLERALDDSNR
MYITENGRDEASTGKIDLK----DSER

**multiple sequence alignment**

---

# Pairwise alignment

**two gene sequences**

- **global or local**
- **homologous sequences from different species**
- **homologous sequences from the same species**
- **genomic vs. cDNA sequence**
- **...**

MYITENGMDEFNNPKVSLERALDDSNR
MYITENGRDEASTGKIDLK----DSER

---

# Pairwise alignment

**two entire genomes**

- **global or local**
- **computation, display**

image taken from http://tinyurl.com/y9e6su4x

# Database search

very frequent problem

- **I have a starting (query) sequence**
- **I want to compare it to homologous sequences**
- **…but which other sequences are homologous to it?**

solution: database search!

- **compare query to a collection of sequences**
- **identify / collect homologs for further analysis**

query sequence   database of sequences

online or local

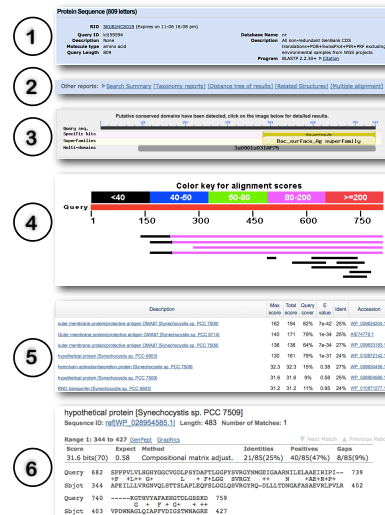4 pairwise alignments

---

# Database search using BLAST

BLAST {
- **compute & score pairwise alignments between a starting (query) sequence and sequences in a database**
- **statistically evaluate the alignments, return the best alignments**
}

YOU →
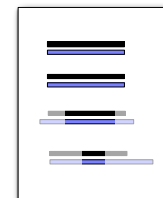- **decide if the similarity is the result of homology (shared ancestry) or of chance**

query sequence   database of sequences

online or local

4 pairwise alignments

---

# BLAST results

1. **Header**
2. **Other reports:**
   Search Summary
   *online only*
   Taxonomy reports
   Distance tree of results
   Multiple alignment
3. **Conserved Domains**
4. **Graphic Summary**
5. **Summary table of alignments**
6. **Alignments**

---



⑤
```
                                                                      Score    E
Sequences producing significant alignments:                          (bits)  Value

TOC75_ORYSJ RecName: Full=Protein TOC75, chloroplastic; AltName:...   1048    0.0
TC754_ARATH RecName: Full=Protein TOC75-4, chloroplastic; AltNam...   493     e-139
MATK_LEPPR RecName: Full=Maturase K; AltName: Full=Intron maturase;   35      0.14
MATK_MACCO RecName: Full=Maturase K; AltName: Full=Intron maturase;   33      0.46
IF3C_EUGGR RecName: Full=Translation initiation factor IF-3, chl...   33      0.54
PDI_MEDSA RecName: Full=Protein disulfide-isomerase; Short=PDI; ...   30      2.8
```

⑥
```
>TOC75_ORYSJ RecName: Full=Protein TOC75, chloroplastic; AltName: Full=75 kDa
           translocon at the outer-envelope-membrane of
           chloroplasts; Flags: Precursor;
           Length = 817

 Score = 1048 bits (2709), Expect = 0.0,   Method: Compositional matrix adjust.
 Identities = 487/695 (70%), Positives = 579/695 (83%), Gaps = 4/695 (0%)

Query: 119 FWSRILSPARAIADEPKSEDWDSHELPADITVLLGRLSGFKKYKISDILFFDRNKKSKVE 178
           FWSRI S  A ADE  S DWD H LPA+I V + +LSG K+YKIS++ FFDR
Sbjct: 123 FWSRIFSGGAAHADEKSSGDWDPHGLPANINVPMTKLSGLKRYKISELKFFDRAAGGGGA 182

Query: 179 ---TQDSFLDMVSLKPGGVYTKAQLQKELESLATCGMFEKVDMEGKTNADGSLGLTISFA 235
              +DSF +MV+L+PGGVYTK+QL KELE+L +CGMFE+VD+EGK   DG+LGLT+SF
Sbjct: 183 FTGPEDSFFEMVTLQPGGVYTKSQLLKELETLVSCGMFERVDLEGKAKPDGTLGLTVSFV 242

Query: 236 ESMWERADRFRCINVGLMGQSKPVEMDPDMSEKEKIEFFRRQEREYKRRISSARPCLLPT 295
           ES+W  A +F+CINVGLM QS  V+ D DM+E+EK+++ R+QER+Y++R+  A+PC+LP
Sbjct: 243 ESVWSAAKQFKCINVGLMSQSGQVDFDQDMTEREKMDYLRKQERDYQQRVRGAKPCILPD 302
Sbjct: 103 -HFVVVKEALLKTIKEVSEDKWSEELNTAWEIAYDGLASAI 142


>PDI_MEDSA RecName: Full=Protein disulfide-isomerase; Short=PDI; EC=5.3.4.1;
           Flags: Precursor;
           Length = 512

 Score = 30.4 bits (67), Expect = 2.8,   Method: Compositional matrix adjust.
 Identities = 16/37 (43%), Positives = 21/37 (56%), Gaps = 1/37 (2%)

Query: 738 PIKGTHVYAFAEHGTDLGSSKDVKGNPTV-VYRRMGQ 773
           P+   V A EH  DL S  DVKG PT+ ++R  G+
Sbjct: 85  PVVLAKVDANEEHNKDLASENDVKGFPTIKIFRNGGK 121
```

## Database search using BLAST

**considerations & challenges:**

- **score similarities?**
- **efficiently search the database?**
  - query: 1 sequence, 500-1000 bp or aa
  - GenBank: 209,656,636 sequences, 279,668,290,132 bases
- **statistically evaluate if the similarity could be due to chance?**

```
MYITENGMDEFNNPKVSLERALDDSNR
|||||||| ||    |  |     || |
MYITENGRDEASTGKIDLK----DSER
```



query sequence    database of sequences

online or local

4 pairwise alignments

## Scoring an alignment

- **matches**
- **mismatches**
- **gaps**

```
MYITENGMDEFNNPKVSLERALDDSNR
|||||||| ||    |  |     || |
MYITENGRDEASTGKIDLK----DSER
```

- **a scoring matrix is used for matches, mismatches,**
- **gap penalties are used for opening / extending gaps**
- **alignment score: sum of scores at each alignment position,** $S = \sum s_{ij}$

- **goal: identify the optimal alignment**

## Scoring an alignment

```
seq 1  MYITENGMDEYNNPKVSLERALDDSNR
       |||||||| ||    |  |     || |
seq 2  MYITENGRDEASTGKIDLK----DSER
```

score: 1+1+1+1+1+1+0+1+1+0+0+0+0+1+0+0+1+0+1+1+0+1 = 14

## Scoring an alignment

```
seq 1  MYITENGMDEYNNPKVSLERALDDSNR
       |||||||| ||    |  |     || |
seq 2  MYITENGRDEASTGKIDLK----DSER
```
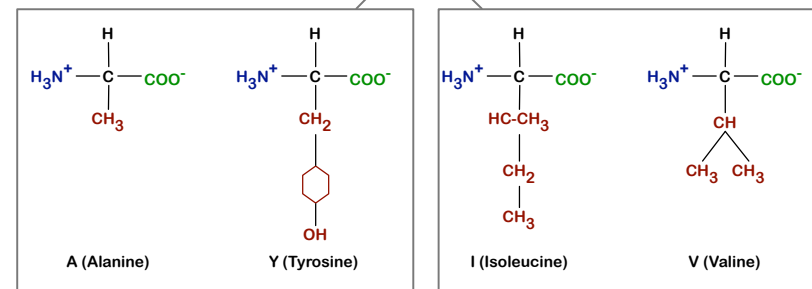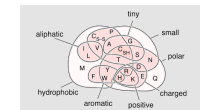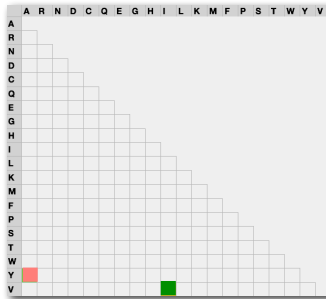




A (Alanine)    Y (Tyrosine)    I (Isoleucine)    V (Valine)

# Scoring an alignment

```
seq 1  MYITENGMDEYNNPKVSLERALDDSNR
       |||||||  ||   | |      || |
seq 2  MYITENGRDEASTGKIDLK----DSER
```

---

# Scoring an alignment

**very similar homologous sequences**

**more divergent homologous sequences**

---

# BLOSUM matrices

- **BLO**cks amino acid **SU**bstitution **M**atrices
- based on local alignments of divergent sequences
- different BLOSUM matrices are based on observed alignments with different degrees of similarity

  - **BLOSUM60 matrix:**

    derived from and best used for sequences that are 60% identical

|  | for proteins of |
| --- | --- |
| **BLOSUM45** | **45% identity** |
| **BLOSUM60** | **60%** |
| **BLOSUM90** | **90%** |

---

# Substitution matrices

➡ **which amino acids occur together in the alignment columns more often than expected by chance?**

**trusted alignment of homologous sequences**



$$s(a,b) = log\left(\frac{p_{ab}}{q_a q_b}\right)$$

**$p_{ab}$ :**     observed frequency of residues a and b aligned

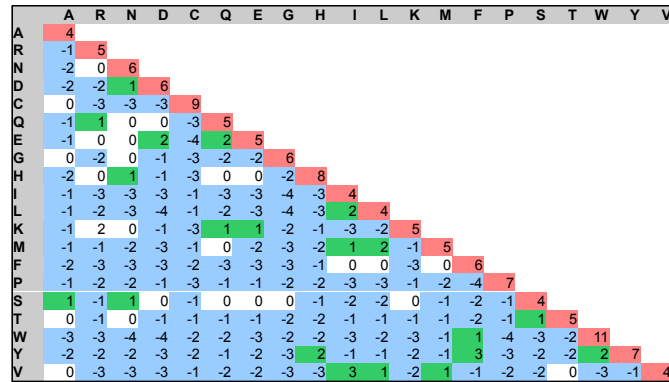**$q_a$, $q_b$ :**   frequencies of residues a and b

```
M:   0.01
L:   0.1
ML:  0.002
```

$$s(M,L) = log(\frac{0.002}{0.01*0.1}) = +1$$

# Slide 1: BLOSUM62

**BLOSUM62**  gaps: -2 each



|   | A | R | N | D | C | Q | E | G | H | I | L | K | M | F | P | S | T | W | Y | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A | 4 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| R | -1 | 5 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| N | -2 | 0 | 6 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| D | -2 | -2 | 1 | 6 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 | -3 | -3 | -3 | 9 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| Q | -1 | 1 | 0 | 0 | -3 | 5 |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| E | -1 | 0 | 0 | 2 | -4 | 2 | 5 |   |   |   |   |   |   |   |   |   |   |   |   |   |
| G | 0 | -2 | 0 | -1 | -3 | -2 | -2 | 6 |   |   |   |   |   |   |   |   |   |   |   |   |
| H | -2 | 0 | 1 | -1 | -3 | 0 | 0 | -2 | 8 |   |   |   |   |   |   |   |   |   |   |   |
| I | -1 | -3 | -3 | -3 | -1 | -3 | -3 | -4 | -3 | 4 |   |   |   |   |   |   |   |   |   |   |
| L | -1 | -2 | -3 | -4 | -1 | -2 | -3 | -4 | -3 | 2 | 4 |   |   |   |   |   |   |   |   |   |
| K | -1 | 2 | 0 | -1 | -3 | 1 | 1 | -2 | -1 | -3 | -2 | 5 |   |   |   |   |   |   |   |   |
| M | -1 | -1 | -2 | -3 | -1 | 0 | -2 | -3 | -2 | 1 | 2 | -1 | 5 |   |   |   |   |   |   |   |
| F | -2 | -3 | -3 | -3 | -2 | -3 | -3 | -3 | -1 | 0 | 0 | -3 | 0 | 6 |   |   |   |   |   |   |
| P | -1 | -2 | -2 | -1 | -3 | -1 | -1 | -2 | -2 | -3 | -3 | -1 | -2 | -4 | 7 |   |   |   |   |   |
| S | 1 | -1 | 1 | 0 | -1 | 0 | 0 | 0 | -1 | -2 | -2 | 0 | -1 | -2 | -1 | 4 |   |   |   |   |
| T | 0 | -1 | 0 | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -2 | -1 | 1 | 5 |   |   |   |
| W | -3 | -3 | -4 | -4 | -2 | -2 | -3 | -2 | -2 | -3 | -2 | -3 | -1 | 1 | -4 | -3 | -2 | 11 |   |   |
| Y | -2 | -2 | -2 | -3 | -2 | -1 | -2 | -3 | 2 | -1 | -1 | -2 | -1 | 3 | -3 | -2 | -2 | 2 | 7 |   |
| V | 0 | -3 | -3 | -3 | -1 | -2 | -2 | -3 | -3 | 3 | 1 | -2 | 1 | -1 | -2 | -2 | 0 | -3 | -1 | 4 |

```
MYITENGMDEFNNPKVSLERALDDSNR          MYITENGMDEFNNPKVSLERALDDSNR
                            →        ||||||||| ||     |  |     ||   →  60
MYITENGRDEASTGKIDLKDSER              MYITENGRDEASTGKIDLK----DSER
```

# Slide 2: (BLOSUM) matrices

**(BLOSUM) matrices**



already done, once

choose matrix

**use to compute & evaluate alignments**

# Slide 3: Scoring matrices

**Scoring matrices**

**identity matrices**
- **easy to understand and implement**
- **biologically not very realistic**

**matrices based on physico-chemical aa properties**
- **better than identity scores**
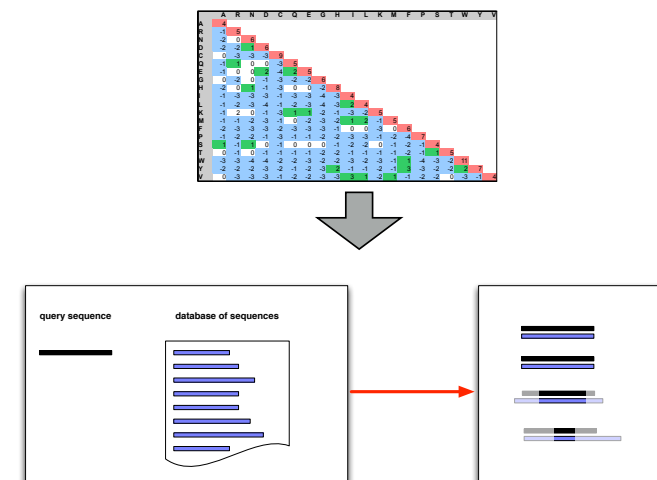- **different divergence times - same scores!**

**BLOSUM matrices**
- **evolutionary perspective:
  sequences accumulate changes over time**

**(many other matrices exist, some are very specific)**

# Slide 4: Finding & evaluating sequence matches

**Finding & evaluating sequence matches**



query sequence    database of sequences

## Pairwise alignment approaches

**exhaustive search**

- **compare every position of the query with every position of every database sequence**

- **not practical**

```
MELPTRD              PPTVKNSN
GCMFA       SRGCMFAAIAAL   GCMFA
 GCMFA         GCMFA         GCMFA
  GCMFA          GCMFA         GCMFA
                  GCMFA          GCMFA
                   GCMFA
                    GCMFA
                     GCMFA
                      GCMFA
                       GCMFA
```

**heuristics**

- **most popular approach: BLAST**
- **exclude unpromising regions from the search**

---

## Main steps of BLAST

- **pre-processing of the database to be searched**

- **seeding**

  - **homologous sequences contain (at least) short stretches of identical or high-scoring amino acids**
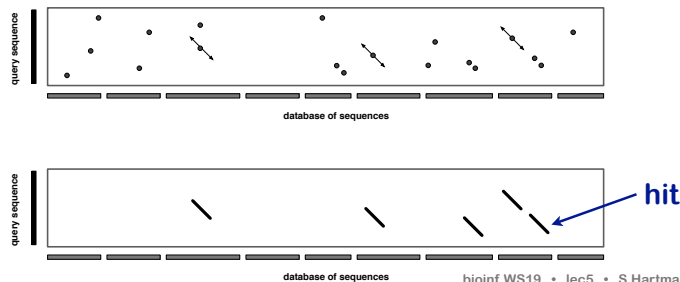
```
KRTDADGAVRG FCFL TASPELQQAMQVQKSASIALERLKEVAYMKQEIRNISQGMIRSREKGLQFVRETT EVKN TLFGDQLRLQQVLADFLTTAVRFTSSS
|||| | + | |||| | + + | | ||| ||++|+++|+|| || +|| || + + ++|+ ++|||++|| ||+||| + || |
KRTDERGNIIG FCFL TMAVDHPQISARDIDDRECLSTLKEFAYIQQQMKNVSQVMIPLKEKNLQLLHDIP QIKS PIYGDQIKLQLVLSDFLLSIVRHAPSP
```

  - **these will be identified and used as seeds**

- **extension to a good longer alignment**

- **evaluation of statistical significance**

- **ranking & presentation of alignments**

---

## Seeding & extension

- **the locations of short, high scoring hits are identified**
- **these are used as alignment seeds**
- **seeds are extended into longer alignments (→"hits")**
- **for each hit, an "E(xpect)-value" is computed**
- **hits are ranked by their E-values and reported**



**hit**

---

## Evaluating alignments

```
DSNRMYIGMDEFNNPKVSLE
DSERMYIGRDEASTGKIDLK
```
→ **query sequence**
→ **database sequence**

**alignment raw score:** $S = \sum s_{ij}$

- **any two sequences (even if they are unrelated!) will have a "best" alignment score**
- **how high of a score can we expect from random (unrelated) sequences? is our current score better?**
- **statistics!**
  - **substitutions, aa distribution, seq lengths, etc**
  - **based on Gumbel extreme value distribution**

## Evaluating alignments

real query sequence →
DSNRMYIGMDEFNNPKVSLE
DSERMYIGRDEASTGKIDLK
← real database sequence

alignment raw score: $S = \sum s_{ij}$

$$E = Kmne^{-\lambda S}$$

$$S' = \frac{(\lambda \times S) - ln(K)}{ln(2)}$$

$$E = mn(2^{-S'})$$

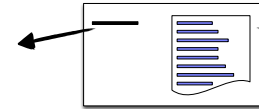**E**: the number of alignments at least as good
that are expected by chance

(under comparable conditions: aa compositon, db & query length, etc )

| 0 | 0.000000003 | 0.002 | 5 | 20 |
|---|---|---|---|---|
| | 3×10⁻⁹ | | | |
| | 3e-09 | | | |

(note: superscript in 3×10⁻⁹)

---

## Evaluating alignments

**query sequence:**
one plant gene
sequence of 809 aa

**database to be searched:**
51,826,119 bacterial
sequences with a total of
18,610,688,330 aa

**db seq: 763 aa; aln: 673** (287/42% | or +)



**E: 7e-42**

**db seq: 409 aa;
aln: 26** (17/65% | or +)

TTSNFLNPQDDLAFKMEYAHPYLDGV
|| + | | + + | | | |++ | |||++
TTYDVLPPSANLAFLMETAVPYVEAV

**E: 9**

---

## BLAST results

- pairwise alignments query–database sequence
  - overview table, ranked by E-value

```
                                                         Score    E
Sequences producing significant alignments:             (bits)  Value
TOC75_ORYSJ    Full=Protein TOC75, chloroplastic; AltName:...  1048   0.0
TC754_ARATH    Full=Protein TOC75-4, chloroplastic; AltNam...   493   e-139
MATK_LEPPR     Full=Maturase K; AltName: Full=Intron maturase;  35    0.14
PDI_MEDSA      Full=Protein disulfide-isomerase; Short=PDI; ...  30    2.8
```

  - alignments, ranked by E-value

```
Score = 1048 bits (2709), Expect = 0.0,  Method: Compositional matrix adjust.
Identities = 487/695 (70%), Positives = 579/695 (83%), Gaps = 4/695 (0%)

Query: 119 FWSRILSPARAIADEPKSEDWDSHELPADITVLLGRLSGFKKYKISDILFFDRNKKSKVE 178
           FWSRI S   A ADE  S DWD H LPA+I V + +LSG K+YKIS++ FFDR
Sbjct: 123 FWSRIFSGGAAHADEKSSGDWDPHGLPANINVPMTKLSGLKRYKISELKFFDRAAGGGGA 182
```

- information about the query and the database
- information about statistical parameters
- different presentation online and in the terminal
- additional information & analysis options online

---

## BLAST search: results

| | online | terminal |
|---|---|---|
| **Header** | ✗ | ✗ |
| Other reports:    Search Summary | ✗* | |
| Taxonomy reports | ✗ | |
| Distance tree of results | ✗ | |
| Multiple alignment | ✗ | |
| **Graphic Summary** | ✗ | |
| Conserved Domains | | |
| Graphical Overview of Hits | | |
| **Descriptions** | ✗ | ✗ |
| **Alignments** | ✗ | ✗ |
| **Footer** | | ✗* |

✗ *     equivalent information

## BLAST programs

| query | database | program |
|---|---|---|
| nucleotide | nucleotide | blastn |
| nucleotide (translated) | nucleotide (translated) | tblastx |
| nucleotide (translated) | peptide | blastx |
| peptide | peptide | blastp |
| peptide | nucleotide (translated) | tblastn |

blastx
translated nucleotide ▶ protein

Nucleotide BLAST
nucleotide ▶ nucleotide

tblastn
protein ▶ translated nucleotide

Protein BLAST
protein ▶ protein

---

## BLAST search in practice

- call the appropriate BLAST program
- select parameters for search
- select database to be searched
- select and submit query sequence

```
blastp
    -db database
    -query querySeq
    -b xxx
    -e yyy
    ...
```

---

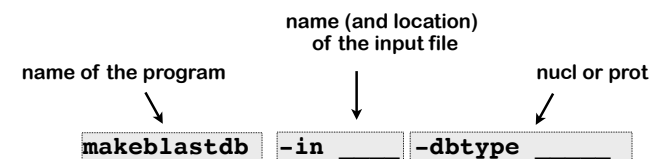## Today's exercise

query:        human beta globin
database:     *Medicago truncatula* genes

1. download the human beta globin in fasta format from NCBI
   - one file each: protein sequence, DNA sequence

2. download medicago DNA and protein sequences from Moodle
   - one file each: protein sequences, DNA sequences

3. convert the medicago sequences into a BLAST database

4. search if the human globin is similar to a medicago sequence
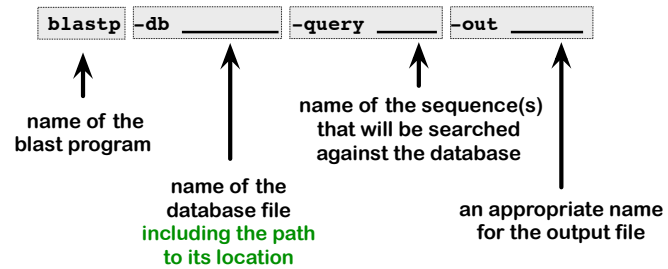   - protein-protein (blastp)
   - DNA-DNA (blastn)

---

## convert a fasta file into a database

name of the program

name (and location)
of the input file

nucl or prot

`makeblastdb` `-in _____` `-dbtype _____`

## command-line BLAST

```
blastp  -db _____  -query _____  -out _____
```

name of the
blast program

name of the
database file
**including the path
to its location**

name of the sequence(s)
that will be searched
against the database

an appropriate name
for the output file

## heute!

**10h: 3 Computer Pools (1a, 1b, 2a)**

**11: 2 Computer Pools (1a, 2a)**

## Key terms and concepts

- homology, similarity
- alignment
  - local, global
- BLOSUM substitution matrices:
  - what are they, how were they derived,
    how are they used in sequence alignment
- heuristics
- database searching / BLAST
  - principle, steps
  - evaluation of results
  - definition and interpretation of "E-value"