

## BIOINF101 - Aufgabe 5

---

### *Human T-cell leukemia virus type I (NC\_001436)*

- Genom - die 100 ersten Nukleotide:

```
GGCTCGCATCTCTCCTTCACGCGCCCGCCGCCTTACCTGAGGCCGCCATCCACGCCGGTTGAGTCGCGTTCTGCC  
GCCTCCCGCCTGTGGTGCCTCCTGA
```

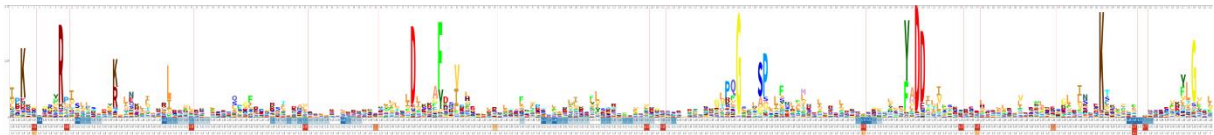
- Peptidsequenz - die ersten 30 Aminosäuren des 1. 5'3' Frames

```
GSHLSFTRPPPYLRPPSTPVESRSAASRLW
```

- Hidden Markov Model Logo

*Reverse transcriptase (RNA-dependent DNA polymerase)* (754 - 925)

<http://pfam.xfam.org/family/PF00078.26#tabview=tab4>



**Einschätzung:** die hochkonservierten Bereiche, an denen bestimmte Aminosäuren mit sehr hoher Wahrscheinlichkeit auftreten, sind in der Sequenz vorhanden und große Bereiche des Profils stimmen weitestgehend mit der Sequenz überein, da an vielen Stellen eine Vielzahl verschiedener Aminosäuren vom Profil zugelassen wird.

### *Human immunodeficiency virus 1 (NC\_001802.1)*

- Genom - die ersten 100 Nukleotide:

```
GGTCTCTCTGGTTAGACCAGATCTGAGCCTGGGAGCTCTCTGGCTAACTAGGGAACCCACTGCTTAAGCCTCAAT  
AAAGCTTGCCTTGAGTGCTTCAAGT
```

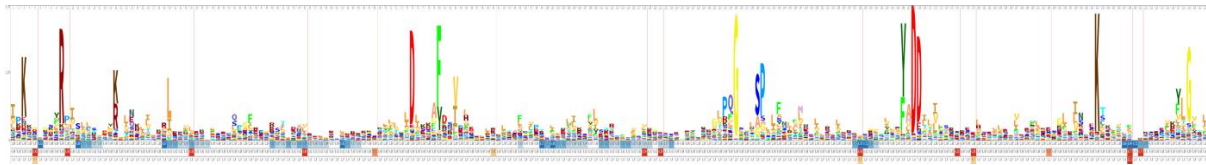
- Peptidsequenz - die ersten 30 Aminosäuren des 2. 5'3' Frames

```
VSLVRPDLGLALWLTREPTA-ASIKLALS
```

- Hidden Markov Model Logo

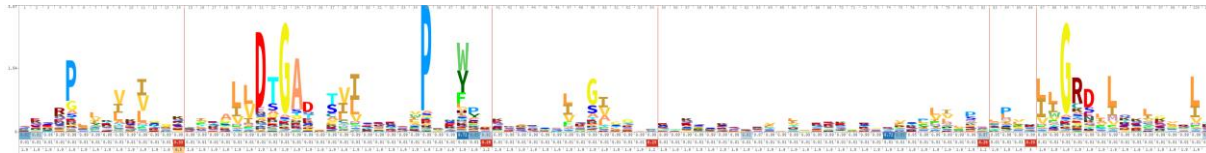
*Reverse transcriptase (RNA-dependent DNA polymerase)*

<http://pfam.xfam.org/family/PF00078.26#tabview=tab4>



## Retroviral aspartyl protease

<http://pfam.xfam.org/family/PF00077.19#tabview=tab4>



**Einschätzung:** Aminosäuresequenz des *immunodeficiency virus* stimmt mit der Sequenz des *T-cell leukemia virus type I* nur in wenigen, hochkonservierten Sequenzen überein (darunter VLPQG im *RVT\_1 profile*), die auch im HMM als konserviert zu erkennen sind. Das HMM lässt an vielen Stellen Variabilität in der Aminosäuresequenz zu. Deshalb passt es zu beiden Sequenzen sehr gut, auch wenn sie sich stark unterscheiden. (siehe *BLAST Results*, oben *T-cell virus* unten *immunodeficiency virus*). Mit der *Retroviral aspartyl protease* verhält es sich ähnlich.

Range 1: 685 to 942 <a href="#">Graphics</a>						▼ Next Match ▲ Previous Match
Score	Expect	Method	Identities	Positives	Gaps	
139 bits(351)	1e-37	Compositional matrix adjust.	83/265(31%)	135/265(50%)	15/265(5%)	
Query 711	APRNQVPFVKPERLQALQHLVRKALEAGHIEPY--TGPGNNPVFPVKKANGT-WRFIHDL	767				
Sbjct 685	GPVKVQWPLTEEEKIKALVEICTEMEKEGKISKIGPENPYNTPVFAIKKKDSTKWRKLVDF	744				
Query 768	RATNSLTIDLSSSSPGPPDLSSLPPTLAHLQITIDLKDAFFQIPLPKQFQPYFAFTVPQQC	827				
Sbjct 745	RELNKRRTQDFWEVQLGIPHPAGLKKKKKS-VTVLDVGDAYFSVPLDEDFRKYTAFTIPSIN	803				
Query 828	NYGPGTRYAWRVLPQGFKNSTPLFEMQLAHILQPIRQAFPOCTILQYMDILLASPSHAD	887				
Sbjct 804	NETPGIRYQYNVLPQGNKGSPIAFQSSMTKILEPFRKQNPDIIVIQYMDLYVGS---D	859				
Query 888	LQLLSEATMAS-----LISHGLPVSENKTOOTPGTIKFLGQIISPHLYDAVPKVPIRS	942				
Sbjct 860	LEIGQHRTKIEELRQLHLLRWGLTTPDKKHQKEPPFL-WMGVELHPDKWTVQPIV-LPEKD	917				
Query 943	RWALPELQALLGEIQWVSKGTPTLR	967				
Sbjct 918	SMTVNDIQKLVGKLNMAISQIYPGIK	942				

Range 2: 267 to 334 <a href="#">Graphics</a>						▼ Next Match ▲ Previous Match ▲ First Match
Score	Expect	Method	Identities	Positives	Gaps	
20.0 bits(40)	0.75	Compositional matrix adjust.	17/68(25%)	29/68(42%)	9/68(13%)	
Query 656	DAYCFNILP-SYQKQLGHH--RSLCTTMRPRVPVPGKKAACNLANTGASC-----PWAR	706				
Sbjct 267	DTHVFSIIRRHPTFRKHHAKHSGGTSSSHANVKRDHQGSCRMGSASSACRAYCTRDER	326				
Query 707	TPPKAPRN	714				
Sbjct 327	TKGKHSRN	334				

Range 3: 936 to 961 <a href="#">Graphics</a>						▼ Next Match ▲ Previous Match ▲ First Match
Score	Expect	Method	Identities	Positives	Gaps	
16.9 bits(32)	6.3	Compositional matrix adjust.	9/26(35%)	11/26(42%)	0/26(0%)	
Query 816	QPYFAFTVPQQCNYPGTRYAWRVLP	841				
Sbjct 936	QIYPGIRVRLCKLLRGTKALTEVIP	961				

Range 4: 799 to 817 <a href="#">Graphics</a>						▼ Next Match ▲ Previous Match ▲ First Match
Score	Expect	Method	Identities	Positives	Gaps	
16.5 bits(31)	8.6	Compositional matrix adjust.	8/21(38%)	11/21(52%)	2/21(9%)	
Query 644	LPSCANTPPFPDDAYCFNILP	664				
Sbjct 799	IPSIINNETP--GIRYQYNVLP	817				

3 a)	3 b)
<ul style="list-style-type: none"> <li>• Genomsequenzen sind größere Datensätze als Aminosäuresequenzen</li> <li>• es würde viele falschpositive Ergebnisse geben, weil die Suche mit 4 Basen unspezifischer ist als mit 20 Aminosäuren</li> <li>• außerdem würden viele unzutreffende Ergebnisse geliefert werden, da der <i>frame</i> nicht bekannt ist (kontextabhängig)</li> <li>• Genomsequenzen weisen nochmal mehr Variabilität auf als Aminosäuresequenzen</li> </ul>	<ul style="list-style-type: none"> <li>• jeder Frame ergibt eine völlig andere Aminosäuresequenz</li> <li>• jeder Frame wird andere HMM-Profile liefern</li> </ul>