

## BIOINF101 – Aufgabe 5:

Nichanok Auevechanichkul (4781404)

### Aufgabe 2:

Anhand der Webseite (<https://www.ncbi.nlm.nih.gov/nuccore/LC378575.1>) lässt sich eine Genomesequenz (**Human T-cell leukemia virus type I OATL9B proviral DNA, complete genome**) darunter darstellen.

```
1  tgacaatgac catgagcccc aaatatcccc cgggggctta gagcctccca gtgaaaaaca
61  tttccgcgaa acagaagtct gaaaagggtca gggcccagac taaggctctg acgtctcccc
121 ccggagggac agctcagcac cggctcaggc taggccctga cgtgtccccc tgaagacaaa
181 tcataagctc agacctccgg gaagccaccg gaaccaccca tttcctcccc atgtttgtca
241 agccgccttc aggcgttgac gacaaccctt cacctcaaaa aacttttcat ggcacgcata
301 tggctgaata aactaacagg agtctataaa agcgtggaga cagttcagga gggggctcgc
361 atctctcctt cagcgccccg ccgccttacc tgaggccgcc atccacgccc gttgagtcgc
421 gttctgcgcg ctcccgcctg tgggtgcttc tgaactgcgt ccgccgtcta ggtaagttaa
481 gagctcaggt cgagaccggg cctttgtccg gcgtccctt ggagcctacc tagactcagc
541 cggttctcca cgtttgcctt gaccctgctt gctcaactct gcgtctttgt ttcgttttct
601 gttctgcgcc gctacagatc gaaagttccg cccctttccc tttcattcac gactgactgc
661 cggcttggcc cacggccaag taccggcgac tccgttggtt cggagccagc gacagcccat
721 tctatagcac tctccaggag agaaacttag tacacagttg ggggctcgtc cgggatacga
781 gcgccccctt attccctagg caatgggcca aatcttttcc cgtagcgcta gccctattcc
841 gcggccgccc cgggggctgg ccgctcatca ctggcttaac ttccctccaag cggcatatcg
901 cctagaaccc ggtccctcca gttacgattt ccaccagttg aaaaaatttc ttaaaatagc
961 tttagaaaca ccggtctgga tctgtcccat taactactcc ctctagcca gcctactccc
1021 aaaaggatac cccggccggg tgaatgaaat tttaacata ctcatccaaa cccaagccca
1081 gatcccgctc cgtcccgcgc caccgcccgc gtcaccccc acccacgacc ccccgatttc
1141 tgatccacaa atccccctc cctatgttag gcctacggcc cccaagtcc tacagtcct
1201 gcaccacatc ggtgcccctc ccaaccatcg cccatggcaa atgaaggacc tacaggccat
1261 taagcaagaa gtctcccaag cagcccctgg gagccccag tttatgcaga ccatccggct
1321 tgcggtgcag cagtttgacc ccactgccaa agacctcaa gacctcctgc agtacctttg
1381 ctccctccctc gtggcttccc tccatcacca gcagctagat agccttatac agaggccga
1441 aacccgaggt attacaggtt ataaccctt agccgggtcc ctccgtgtcc aagccaacaa
1501 tccacaacaa caaggattaa ggcgagaata ccagcaactc tggctcgccg ccttcgccgc
1561 cctgccaggg agtgccaaag acccttctct ggctctatc ctccaaggcc tggaggagcc
1621 ttaccacgcc ttcgtagaac gcctcaacat agctcttgac aatgggctgc cagaaggcac
1681 gcccaaagac cccatcttac gttccttagc ctactccaat gcaaacaaag aatgccaaaa
1741 attactacag gcccgaggac acactaatag ccctctagga gatatgttgc gggcttgtca
1801 gacctggacc cccaaagaca aaaccaaagt gttagtgtgc cagcctaaaa aacccccccc
1861 aatcagccg tgcttcgggt gcgggaaagc aggccactgg agtcgggact gcactcagcc
1921 tcgtccccc ccggggccat gccccctatg tcaagacca actcactgga agcgagactg
1981 ccccgcccta agcccacta tcccagaacc agagccagag gaagatgccc tcctattaga
2041 cctccccgct gacatccac acccaaaaaa ctccataggg ggggaggttt aacctcccc
2101 cccacattac agcaagtcct tcttaaccaa gaccagcat ctattctgcc agttataccg
2161 ttagatcccc cccgtcggcc cgtaattaaa gccaggttg acaccagac cagccacca
2221 aagactatcg aagctttact agatacagga gcagacatga cagtccttcc catagccttg
2281 ttctcaagta atactccctt caaaaatata tccgtattag gggcaggggg ccaaacccaa
2341 gatcacttta agctcacctc ccttctctgt ctaatacgcc tccctttccg gacaacgcct
2401 attgttttaa catcttgctt agttgatata aaaaacaact gggccatcat aggtcgcgat
2461 gccttacaac aatgccaggg cgtcctgtac ctccctgagg caaaaaggcc gcctgtaatc
2521 ttgccaatac aggcgccagc cgtccttggg ctagaacacc tcccaaggcc ccccgaaatc
2581 agccagttcc ctttaaacca gaacgcctcc aggccttgca acacttggtc cggaaggccc
2641 tggaggcagg ccatatcgaa ccctacaccg gaccaggaaa taaccagta ttcccagtta
```

2701	aaaaggccaa	tggaacctgg	cgattcatcc	acgacctgcg	ggccactaac	tctctaacca
2761	tagatctctc	atcatcttcc	cccggggccc	ctgacttgtc	cagcctgcca	actacactag
2821	cccacttgca	aactatagac	cttaaagacg	cctttttcca	aatcccccta	cctaaacagt
2881	tccagcccta	ctttgctttc	actgtcccac	agcagtgtaa	ctacggcccc	ggcactagat
2941	acgcctggaa	ggtactaccc	caagggttta	aaaatagtcc	caccctgttc	gaaatgcagc
3001	tggcccatat	cctgcagccc	attcggcaag	ctttccccc	atgcactatt	cttcagtaca
3061	tggatgacat	tctcctggca	agccccctcc	atgaggacct	actactactc	tcagaggcca
3121	caatggcttc	cctaattctc	catgggttgc	ctgtgtccga	aaacaaaacc	cagcaaacc
3181	ctggaacaat	taagttccta	gggcaaataa	tttcacccaa	ccacctcact	tatgatgcag
3241	tccccacggg	acctatacgg	tcccgcctgg	cgctacctga	acttcaagcc	ctacttggcg
3301	agattcagtg	ggtctccaag	ggaactccta	ccttacgcca	gccccctcac	agtctctact
3361	gtgccttaca	aaggcatact	gatccccgag	accaaataa	tttaaattcct	tctcaagtct
3421	aatcattagt	gcagctgcgg	caggccctgt	cacagaactg	ccgcagtaga	ctagtccaaa
3481	ccctgcccc	cctaggggct	attatgctga	ccctcactgg	caccactact	gtagtgttcc
3541	agtccaagca	gcagtggcca	cttgtctggc	tacatgcccc	cctaccccc	actagccagt
3601	gccccctggg	gcagctactt	gcctcagctg	tgttattact	cgacaaatac	accttgcaat
3661	cctatgggct	actctgcaa	accatacatc	ataacatctc	cacccaaacc	ttcaaccaat
3721	tcattcaaac	atctgaccac	cccagtgttc	ctatcttact	ccaccacagt	caccgattca
3781	aaaatttagg	tgcccaaact	ggagaacttt	ggaacacttt	tcttaaaaca	gctgccccat
3841	tggctcctgt	aaaagccctc	atgccagtgt	ttactctttc	cccgggtgatc	ataaacaccg
3901	ccccctgcct	gttttcagac	ggatctacct	cccgggcagc	ctatattctc	tgggacaagc
3961	atacattgtc	acaaagatca	ttcccccttc	cgccaccgca	caagtcggcc	caacggggccg
4021	aacttctcgg	acttttgc	ggcctttcca	gcgcccgttc	gtggcgctgt	ctcaacatat
4081	ttctagactc	caagtatctt	tatcattacc	ttcggacctt	tgccctgggc	accttccaag
4141	gcaggctctc	tcaggcccc	tttcaggccc	tctgccccg	cttactatcg	cgtaaggctcg
4201	tctatttgca	ccacgttcgc	agccatacca	atctacctga	tcccatctcc	aggctcaacg
4261	ctctcacaga	tgccctacta	atcaccctcg	tctgcaagc	ctctcctgca	gaactacaca
4321	gtttcaccca	ttgcggacag	acggccctca	cattgcaagg	ggcaaccaca	actgaggctt
4381	ccaatatcct	gcgctcttgc	cacgcctgcc	gcaaaaataa	cccacaacat	cagatgcctc
4441	ggggacacat	ccgcccgtgg	ctacttcccta	accacatctg	gcaaggcgac	attaccatt
4501	tcaaataata	aaatacgtcg	taccgccttc	atgtatgggt	agacaccttt	tcaggagcca
4561	tctcagctac	ccaaaagaga	aaagaaacaa	gctcagaagc	tatttcctct	ttgcttcagg
4621	ccattgccta	tctaggcaag	cctagctaca	taaacacaga	caacggccct	gcctatattt
4681	ccaagactt	cctcaatatg	tgtacctccc	ttgctattcg	acatactacc	catgtccctt
4741	acaatccaac	cagctcagga	cttgtagaac	gctctaattg	cattcttaaa	accctattat
4801	ataagtactt	tactgacaaa	cccagacctac	ccatggataa	tgctctatcc	atagccctat
4861	ggacaatcaa	ccacctgaat	gtgttaacca	actgccacaa	aaccgatgg	cagcttcacc
4921	actccccccg	actccagccg	atcccagaga	cacattccct	cagcaataaa	caaaccatt
4981	ggtattattt	caagcttcct	ggtcttaata	gccgccagtg	gaaaggacca	caggaggctc
5041	tccaagaagc	tgccggcgct	gctctcatcc	cggtaaagcg	tagttctgcc	cagtggatcc
5101	cgtggagact	cctcaagcga	gctgcatgcc	caagaccctg	cggaggcccc	gccgatccca
5161	aagaaaaaga	ccaccaacac	catgggtaag	tttctcgcca	ctttgatttt	attcttccag
5221	ttctgcccc	tcatectcgg	tgattacagc	cccagctgct	gtactctcac	aattggagtc
5281	tcctcatacc	actctaaacc	ctgcaatcct	gcccagccag	tttgttcgtg	gaccctcgac
5341	ctgccggccc	tttcagcaga	tcaggcccta	cagccccctt	gccctaattc	agtaagttac
5401	tccagctacc	atgccaccta	ttccctatat	ctattccctc	attggattaa	aaagccaaac
5461	cgaaatggcg	gaggctatta	ttcagcctct	tattcagacc	cttgttccct	aaagtgccca
5521	tacctggggg	gccaatcatg	gacctgcccc	tatacaggag	ccgtctccag	cccctactgg
5581	aagtttcagc	aagatgtcaa	ttttactcaa	gaagtttcac	gcctcaatat	taatctccat
5641	ttttcgaaat	gcggttttcc	cttctccctt	ctagtcgacg	ctccaggata	tgaccccatc
5701	tggttcctta	ataccgaacc	cagccaactg	cctcccaccg	ccccctctct	actccccac
5761	tctaacctag	accacatcct	cgagccctct	ataccatgga	aatcaaaact	cctgaccctt
5821	gtccagttaa	ccctacaaag	cactaattat	acttgcatg	tctgtatcga	tcgtgccagc
5881	ctatccactt	ggcacgtcct	atactctccc	aacgtctctg	ttccatcctc	ttcttctacc
5941	ccccctcttt	acccatcggt	agcgcttcca	gccccccacc	tgacgttacc	atttaactgg
6001	acccactgct	ttgaccccc	gattcaagct	atagtctcct	ccccctgtca	taactccctc
6061	atcctgcccc	ccttttcctt	gtcacctgtt	cccaccctag	gatcccgcctc	ccgccgagcg
6121	gtaccgggtg	cggtctggct	tgtctccgcc	ctggccatgg	gagccggggg	ggctggcggg

6181	attaccggct	ccatgtccct	cgccctcagga	aagagcctct	tacatgaggt	ggacaaagat
6241	atttcccaat	taactcaagc	aatagtcaaa	aaccacaaaa	atctactcaa	aattgcgcag
6301	tatgctgccc	agaacagacg	aggccttgat	ctcctgttct	gggagcaagg	aggattatgc
6361	aaagcattac	aagaacagtg	ctgttttctg	aatattacta	attcccatgt	ctcaatacta
6421	caagagagac	ccccctgga	gaatcgagtc	ctgactggct	ggggccttaa	ctgggacctt
6481	ggcctctcac	agtgggctcg	agaggcctta	caaactggaa	tcacccttgt	cgcgctactc
6541	cttcttgтта	tccttgacgg	accatgcac	ctccgtcagc	tacgacacct	cccctcgcgc
6601	gtcagatacc	cccattactc	tcttataaac	cctgagtcac	ccctgtaaac	caagcacata
6661	attattgcaa	ccacatcgcc	tccagcctcc	cctgccaata	attaacctct	cccattaaat
6721	cctccttctc	ctgcagcaac	ttcctccgtt	cagcctccaa	ggactccacc	tcgccttcca
6781	actgtctagt	atagccatca	atccccaaat	cctgcatttt	ttcttttcta	gcactatgct
6841	gttttcgcctt	ctcagcccct	tgtctccact	tgcgctcacg	gcgctcctgc	tcttcctgct
6901	ttctccgggc	gacgtcagcg	gccttcttct	ccgcccgcct	cctgcgcctg	gccttctcct
6961	cttccttctc	tttcaaatac	tcagcaatct	gcttttctct	ctcttttctc	cgtctttttt
7021	ttcgcttctc	cttctcctca	gcccgtcgct	gccgatcacg	atgcgtttcc	ccgcgagggtg
7081	gcgctttctc	ccctggaggg	ccccgtcgca	gccggccgcg	gctttctctc	tctaaggata
7141	gcaaaccgtc	aagcacagct	tcctcctcct	ccttgctcct	taactcttcc	tccaaggata
7201	atagcccgtc	caccaattcc	tccaccagca	ggctcctccg	gcatggcaca	ggcaagcatc
7261	gaaacagccc	tacagataca	aagttaacca	tgcttattat	cagcccactt	cccagggttt
7321	ggacagagcc	ttcttttcgg	ataccagctc	tacgtgtttg	gagactgtgt	acaaggcgac
7381	tgggtgcccc	tctctggggg	actatgttcg	gcccgcctac	atcgtcacgc	cctactggcc
7441	acctgtccag	agcatcagat	cacctgggac	cccatcgatg	gacgcgttat	cggctcagct
7501	ctacagtctc	ttatccctcg	actcccctcc	ttccccaccc	agagaacctc	taagaccctc
7561	aaggtcctta	ccccgccaat	cactcataca	accccccaaca	ttccaccctc	cttccctccag
7621	gccatgcgca	aatactcccc	cttcgcgaat	ggatacatgg	aaccaccctc	tgggcagcac
7681	ctcccaaccc	tgtcttttcc	agaccccgga	ctccggcccc	aaaacctgta	cacctctggt
7741	ggaggctccg	ttgtctgcat	gtacctctac	cagctttccc	cccccatcac	ctggccccctc
7801	ctgcccccag	tgattttttg	ccaccccggc	cagctcgggg	ccttcctcac	caatgttccc
7861	tacaagcgaa	tagaagaact	cctctataaa	atttccctta	ccacaggggc	cctaataatt
7921	ctaccggaag	actgtttgac	caccaccctt	ttccagcctg	ttagggcacc	cgtcacgcta
7981	acagcctggc	aaaacggcct	ccttccgttc	caactcaaccc	tcaccactcc	aggccttatt
8041	tggacattta	ccgatggcac	gcctatgatt	tcggggccct	gccctaaaga	tggccagcca
8101	tcttttagtac	tacagtccctc	ctcctttata	tttcacaaat	ttcaaaccac	ggcctaccac
8161	ccctcattcc	tactctcaca	cggcctcata	cagtactctt	cctttcataa	tttacatctc
8221	ctgtttgaag	aatacaccaa	catccccatt	tctctacttt	ttaacgaaaa	agaggcagat
8281	gacaatgacc	atgagcccca	aatatcccc	gggggcttag	agcctcccag	tgaaaaacat
8341	ttccgcgaaa	cagaagtctg	aaaaggctcag	ggcccagact	aaggctctga	cgtctcccc
8401	cggaggggaca	gctcagcacc	ggctcaggct	aggccctgac	gtgtccccct	gaagacaaat
8461	cataagctca	gacctccggg	aagccaccgg	aaccacccat	ttcctcccca	tgtttgtcaa
8521	gccgcccctca	ggcgttgacg	acaacccctc	acctcaaaaa	acttttcatg	gcacgcatat
8581	ggctgaataa	actaacagga	gtctataaaa	gcgtggagac	agttcaggag	ggggctcgca
8641	tctctccttc	acgcgcccgc	cgccctacct	gaggccgcca	tccacgcggg	ttgagtgcgc
8701	ttctgcccgc	tcccgctgtg	ggtgctcct	gaactgcgtc	cgccgtctag	gtaagtttag
8761	agctcaggtc	gagaccgggc	ctttgtccgg	cgtcccttg	gagcctacct	agactcagcc
8821	ggttctccac	gctttgctg	accctgcttg	ctcaactctg	cgtctttgtt	tcgttttctg
8881	ttctgcgccc	ctacagatcg	aaagttccgc	ccctttccct	ttcattcacg	actgactgcc
8941	ggcttgggcc	acggccaagt	accggcgact	ccgttggctc		

### Aufgabe 3:

Durch das „Such-Tool: (<https://web.expasy.org/cgi-bin/translate/dna2aa.cgi>)“ lässt sich die codierten Aminosäuresequenzen aus der ersten 1000 Nukleinsäure alle 6-möglichen Übersetzungs-Frames durchsuchen. Die ersten und zweiten des 5'3' Frames sind darunter zu sehen.

#### 5'3' Frame 1

Met A R I W L N K L T G V Y K S V E T V Q E G A R I S P S R A R R  
P T **Stop**

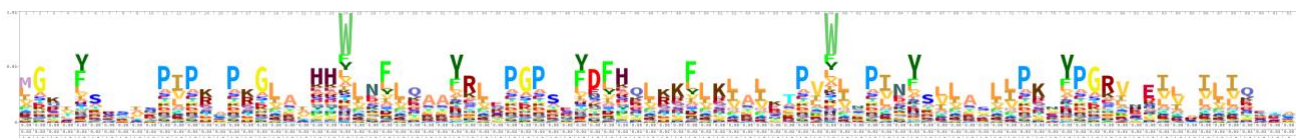
#### 5'3' Frame 2

Met G Q I F S R S A S P I P R P P R G L A A H H W L N F L  
Q A A Y R L E P G P S S Y D F H Q L K K F L K I A L E T P  
V W I C P I N Y S L L A

- Nur anhand Genomsequenzen weiß man nicht direkt welche Protein synthetisiert wird. Man muss die exons und introns Bereiche von dem genome wissen und wie die gene (alternativ) gespleißt wird.
- Außerdem, a Genomen pro 3 nucleotide (Codon) zu einem aminosäure entcodiert werden. ist es sinnvoll alle 6 möglichen Übersetzungs-Frames zu durchsuchen.
- Durch Aminosäuresequenzen weiß man direkt, welche Protein daraus synthetisiert wird. Daher wird Aminosäuresequenzen bevorzugt.

### Aufgabe 4:

Die Aminosäuresequenzen von dem ersten und zweiten des 5'3' Frames werden in dem Such-Tool nach HMM-Profilen aus den Pfam Datenbank eingegeben. Leider gibt es bei der Sequenz des ersten 5'3' Frames kein Match (hits = 34 score). Das HMM Logo zu dem Profil des zweiten 5'3' Frames ist darunter gezeigt (Ref: <http://pfam.xfam.org/family/PF02228.15>).



Ref: <http://pfam.xfam.org/family/PF02228.15#tabview=tab4>

Im Vergleich zu der Suchsequenz in der Aufgabe 3, ist zu erkennen, dass die meiste Sequenz identisch zu der höchst vorkommenden Aminosäure (größte Schrift) von der jeweiligen Position ist.

#### Aufgabe 5:

Gene: **Infectious spleen and kidney necrosis virus strain MV257 major protein capsid (MCP) gene, partial cds**

```
1  tcacacaagg tgaatctgcc attgatggcc accaatcccc tgtccgaggt gtcactcatt
61  tacgagaaca cccctcggct ccaccagatg ggagtagact acttcacatc tgtcgacccc
121 tactactttg cgcccagcat gcctgagatg gatgggtgta tgacctactg ctatacgttg
181 gacatgggca atatcaaccc catgggttca accaactacg gccgcctgtc caacgtcacc
241 ctgtcatgta aggtgtcggg caatgcaaag accaccgcgg cgggcggtgg cggcaacggc
301 tccggctaca cggtgggcca aaagtttgaa ctggtcggt
```

Ref: <https://www.ncbi.nlm.nih.gov/nuccore/KY354078.1>

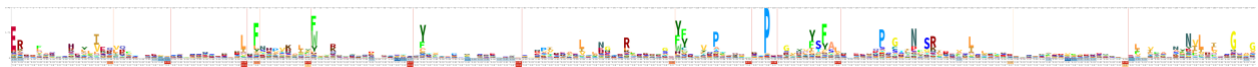
Aminosäuresequenz:

5'3' Frame 1

M A T N P L S E V S L I Y E N T P R L H Q M G V D Y F T S  
V D P Y Y F A P S M P E M D G V M T Y C Y T L D M G N I  
N P M G S T N Y G R L S N V T L S C K V S D N A K T T A A  
G G G G N G S G Y T V A Q K F E L V V

Ref: <https://web.expasy.org/cgi-bin/translate/dna2aa.cgi>

HMM Profile von der Pfam Datenbank:



Ref: <http://pfam.xfam.org/family/PF04451.11#tabview=tab4>

Im Vergleich zu der Suchsequenz aus der ExPASy sind die Sequenz aus beiden Tools ziemlich verschieden. N zum Teil der höchst conserved Aminosäure von beiden Tools sind identisch.