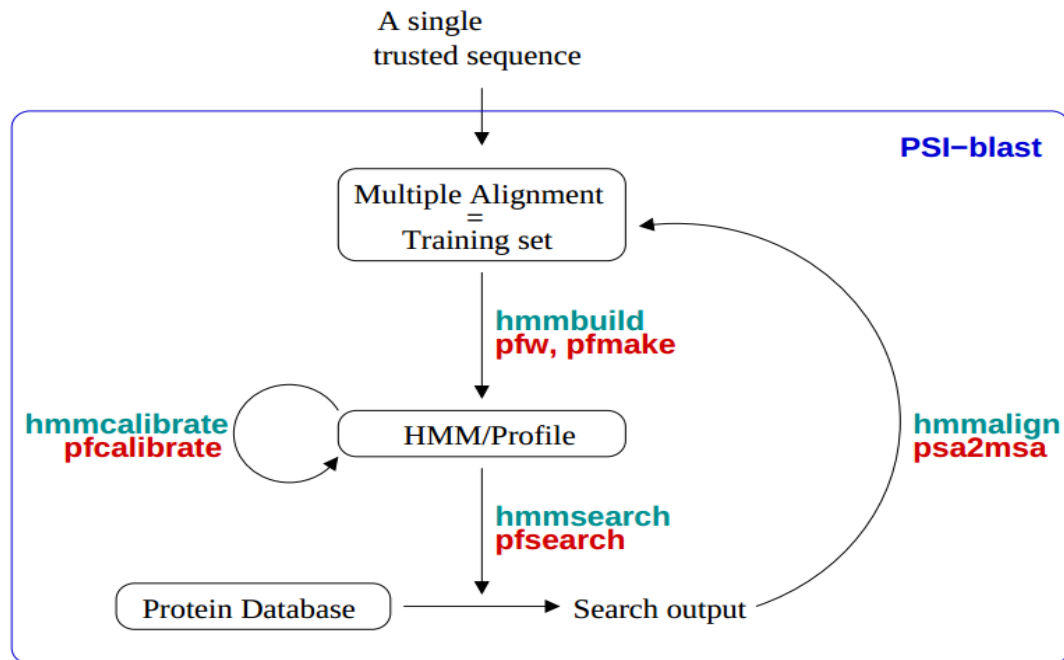


Lab 1

1. Скачать данные - [Plasmodium falciparum chromosome 14](#) или скачиваем chromosome 14 и ее аннотацию с [ncbi](#).
2. Прочитать позиции начал кодирующих регионов из gbk файла (CDS FEATURES). Позиции, в которых встречаются “<” или “>”, игнорировать (обозначают, что точные позиции начала или конца региона не известны).
3. Вывести в файл последовательности нуклеотидов для начал сайтов трансляции (+- 10 нуклеотидов от первой позиции трансляции, т.е. всего 21 нуклеотид). Неизвестные нуклеотиды ‘N’ не надо учитывать. Если кодирующий регион идет “не последовательно”, то надо считать только кодирующую область. (Например, если кодирующая область join(1000..1008,1200..1500), то начальная позиция +10 будет соответствовать координате 1201).
4. Делим хромосому на 2 выборки (сайты трансляции, которые начинаются в первой половине хромосомы - обучающий набор, остальные - тестовый). Далее работаем с обучающей выборкой. Посчитайте частоту встречаемости нуклеотидов для каждой позиции сайта трансляции (PFM).
5. Используя найденные последовательности в пункте 2, постройте лого с помощью [weblogo](#).
6. Посчитать матрицу весов (position specific scoring matrix), используя результаты, полученные в пункте 4.
7. Используя матрицу весов, постройте гистограмму scores для
 - a. известных начал CDS в обучающем наборе (из пунктов 1-2)
 - b. всех подпоследовательностей во всей хромосоме. Приблизить это распределение нормальным.
8. По данным, полученным в предыдущем пункте, постройте график FP, FN в зависимости от отсечки для score, по которой определяется, что последовательность в хромосоме - это CDS.
9. С использованием значения отсечки, которое кажется вам разумным, проведите поиск CDS во второй половине хромосомы и проанализируйте результат (FP, FN, sensitivity, specificity, and PPV).

Lab 2

PSI-BLAST, Generalized profiles, and HMMs



1. Цитохромы P450 – семейство наиболее мощных детоксицирующих ферментов организма. Последовательности для нескольких представителей представлены после условий заданий.
2. Выравниваем последовательности белков и берем в нем достаточно консервативный участок длиной 15-20 аминокислот. Для выравнивания можно использовать любую программу, например, <http://tcoffee.vital-it.ch/apps/tcoffee/do:regular>
3. По выравниванию(для выбранных консервативных участков) строим профиль, например, с помощью [pfmake](#). (Объяснить, какую [матрицу](#) вы использовали и почему)
4. По базе [UniProtKB/Swiss-Prot](#) ищем белки, которые подходят к нашему профилю. (Для поиска по базе используется [pfsearch](#)). Какой cut off вы выбрали?
5. Найдите все цитохромы P450 в этой базе и проанализируйте насколько хорош найденный профиль (вычислить sensitivity, specificity and PPV). Прodelайте еще одну итерацию(начиная с пункта 3), используя найденные последовательности в качестве входных данных. Проанализируйте новые результаты. Стали ли они лучше/хуже и как вы думаете, почему? Что нашлось кроме цитохромов (посмотрите несколько последовательностей)?
6. Используя [PRATT](#), найдите для изначальных цитохромов паттерн и прогоните его по базе UniProtKB/Swiss-Prot с помощью [ScanProsite](#). Проанализируйте результаты.
7. *(необязательное) Прогоните psi-blast с одной из входных последовательностей. Какие последовательности нашел psi-blast?

>1cpt00

MDARATIEPHIARTVILPQGYADDEVIYPAFKWLRDEQPLAMAHIEGYDPMWIATKHADVMQIG
KQPGFLSNAEGSEILYDQNNEAFMRSISGGCPHVIDSLTSM DPPTHTAYRGLTLNWFQPASIRK
LEENIRRIAQASVQRLLD F D G E C D F M T D C A L Y Y P L H V V M T A L G V P E D D E P L M L K L T Q D F F G V E A
ARRFHETIATFYDYFNGFTVDRRSCP KDDVMSLLANSKLDGNYIDDKYINAYYVAIATAGHDTT
SSSSGGAIIGLSRNPEQLALAKSDPALIPRLVDEAVRWTAPVK S F M R T A L A D T E V R G Q N I K R G D
RIMLSYPSANRDEEVFSNPDEFDITRFPNRHLGFGWG AHMCLGQHLAKLEMKIFFEELLPKLKS
VELSGPPRLVATNFVGGPKNVPIRFTKA

>1bvyA0

LNTDKPVQALMKIADELGEIFKFEAPGRVTRYLSSQRLIKEACDES R F D K N L S Q A L K F V R D F A G
DGLFTSWTHEKNWKKAHNILLPSFSQQAMKGYHAMMVDIAVQLVQKWERLNADEHIEVPEDMTR
LTLDTIGLCGFNYRFNSFYRDQPHPFITSMVRALDEAMNKLQRANPDDPAYDENKRQFQEDIKV
MNDLVDKIADRKASGEQSDDLTHMLNGKDPETGEPLDDENIRYQIITFLIAGHETTSGLLSF
ALYFLVKNPVHLQKAAEEAARVLVDPVPSYKQVKQLKYVGMVLNEALRLWPTAPAFSLYAKEDT
VLGGEYPLEKGDELMLIPQLHRDKTIWGDDVEEFRPERFENPSAIPQHAFKPFNGNGQRACIGQ
QFALHEATLVLGMMMLKHDFEDHTNYELDIKETLTLKPEGFVVKAKSKKIPLGGI

>1eupA0

ATVPDLESDFS F H V D W Y S T Y A E L R E T A P V T P V R F L G Q D A W L V T G Y D E A K A A L S D L R L S S D P K K K Y
PGVEVEFPAYLGFPE D V R N Y F A T N M G T S D P P T H T R L R K L V S Q E F T V R R V E A M R P R V E Q I T A E L L
DEVGDSGVVDIVDRFAHPLPIKVICELGVDEAARGAFGRWSSEILVMDPERAEQRGQAAREVV
NFILDLVERRRTEPGDDL S A L I S V Q D D D D G R L S A D E L T S I A L V L L L A G F E A S V S L I G I G T Y L L
LTHPDQLALVRADPSALPNAVEEILRYIAPPETTT R F A A E E V E I G G V A I P Q Y S T V L V A N G A A N R
DPSQFPDPHRF D V T R D T R G H L S F G Q G I H F C M G R P L A K L E G E V A L R A L F G R F P A L S L G I D A D D V V
WRRSLLLRGIDHLPVRLDG

>1akd00

NLAPLP PHVPEHLVDFDFMYNPSNLSAGVQEAWAVLQESNV PDLVWTRCNGGHWIATRGQLIRE
AYEDYRHFSS E C P F I P R E A G E A Y D F I P T S M D P P E Q R Q F R A L A N Q V V G M P V D K L E N R I Q E L A C S
LIESLRPQGQC N F T E D Y A E P F P I R I F M L L A G L P E E D I P H L K Y L T D Q M T R P D G S M T F A E A K E A L Y
DYLIPIEQRRQKPGTDAISIVANGQVNGRPITSDEAKRMCGLLLVGGLDTVVNFLSFSMEFLA
KSPEHRQELIQRPERIPAA CE L L R R F S L V A D G R I L T S D Y E F H G V Q L K K G D Q I L L P Q M L S G L D E
RENACPMHVDFSRQKVSH T T F G H G S H L C L G Q H L A R R E I I V T L K E W L T R I P D F S I A P G A Q I Q H K S
GIVSGVQALPLVWDPATTKAV

>1cmnA0

APSF PFSRASGPEPPAEFAKL RATNPVSQVKLFDGSLAWLVTKHKDVC F V A T S E K L S K V R T R Q G
FPELSASGKQAAKAKPTFVDM DPPEHMHQRSMVEPTFTPEAVKNLQPYIQRTVDDLLEQMKQKG
CANGPVDLVKEFALPVPSYIIYTLLGV PFNDLEYLTQQNAIR TNGSSTAREASAANQELLDYLA
ILVEQRLVEPKDDIISKLC TE Q V K P G N I D K S D A V Q I A F L L L V A G N A T M V N M I A L G V A T L A Q H P D
QLAQLKANPSLAPQFVEELCRYHTAVALAIKRTAKEDVMIGDKLVRANEGH I A S N Q S A N R D E E V
FENPDEFNMNRKWPPQDPLGFGFGDHRCIAEHLAKAELTTVFSTLYQKFPDLKVAVPLGKINYT
PLNRDVGIVDLPVIF