

# **Detection of Conserved RNA Secondary Structures: Hepatitis B as an Example**

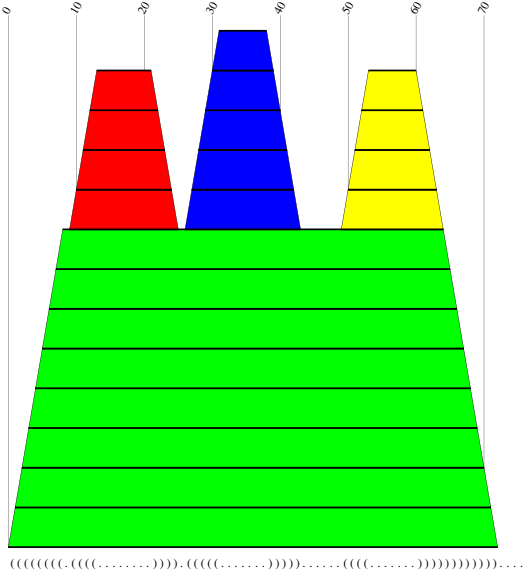
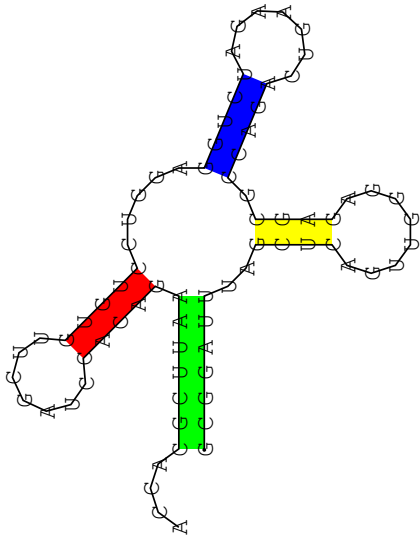
Roman STOCSITS

Theoretical Biochemistry Group

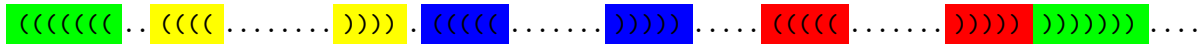
University of Vienna

*Bled, Jan 2002*

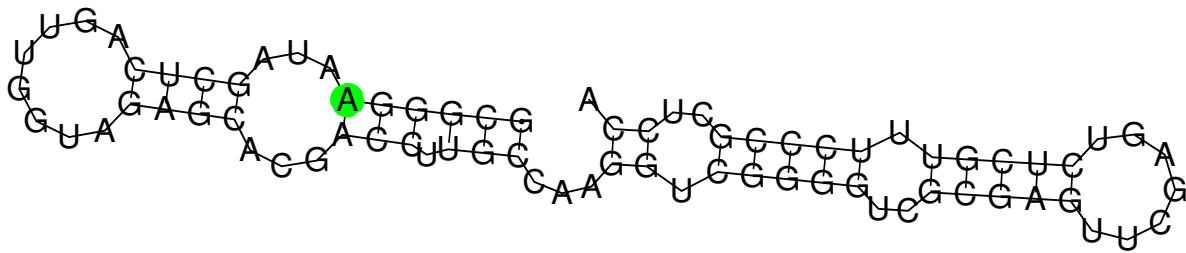
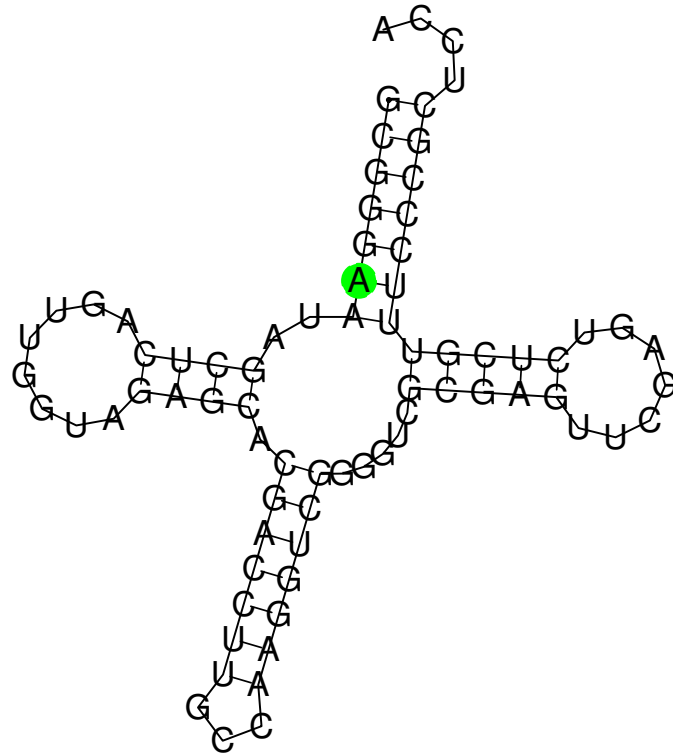
# Representation of Secondary Structures



(((((((.....)))))).....((((.....)))).....((((.....)))))).....



# The Effect of a Single Mutation



SARGLSSTVSLGQFEHWSPR

+AR+LS+TVSL+QF+H SPR

NARNLSDTVSLSQFDHPSPR

AGTGCAAGAGGATTAAGTAGTACAGTAAGTTTAGGACAATTTGAACATTGGAGTCCAAGA

GC G G T AC G T CA TT GA CA CC G

GACGCCGCGACCTCTCCGACACCGCTTCCCTCTCCCAGTTCGACCACCCCTCCCCCGC

Example for the problem of higher sequence heterogeneity on the level of nucleic acids. Amino acid alignment on top with high degree of similarity.

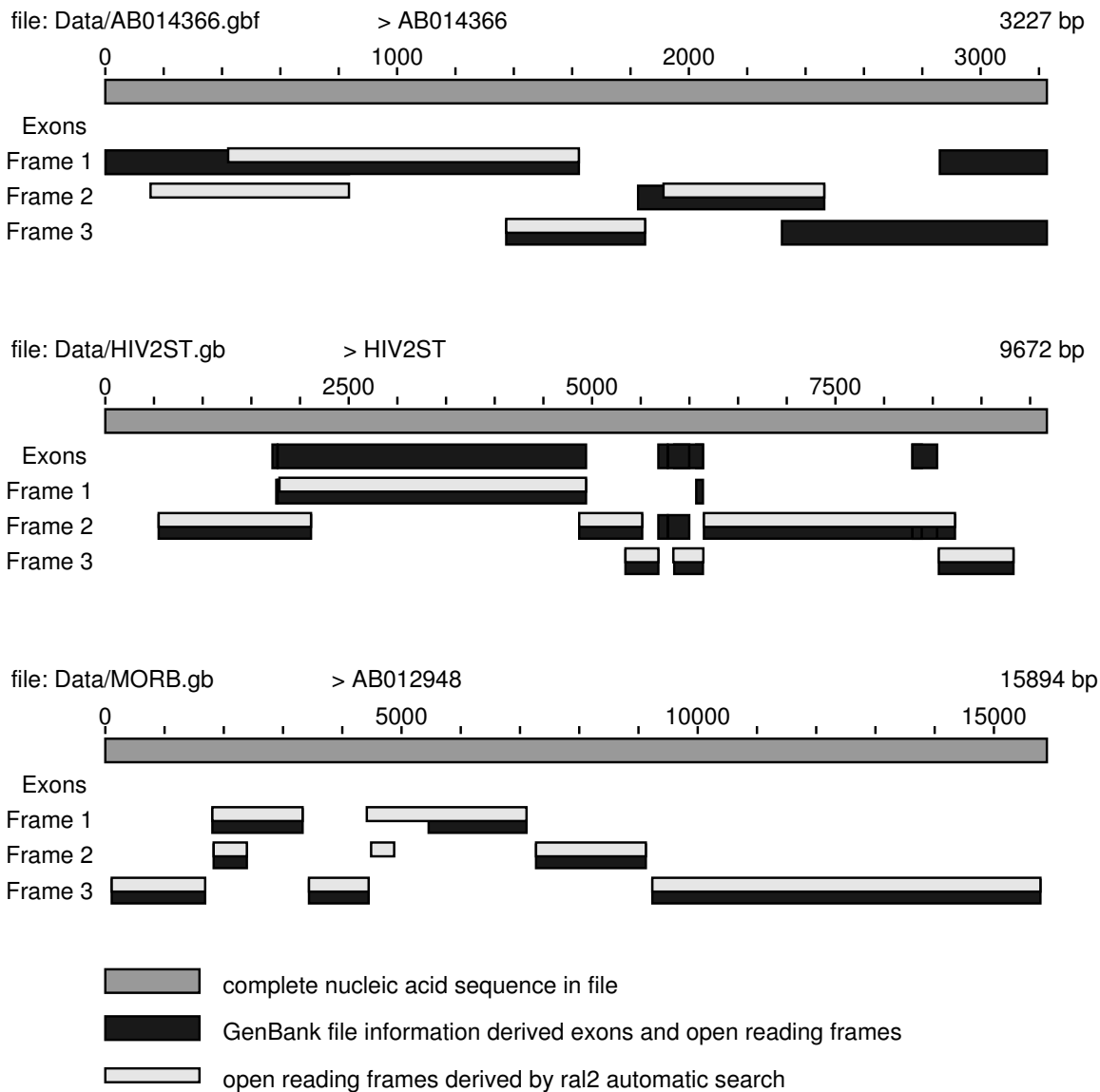
Below the same sequences on level of nucleic acids: much more heterogenous.

Pairwise identity is only 33%.

- Introduction to the problem
- The RALIGN project: A short review
- The RAL2 project: The status quo of development

The following codon tables are available:

univ: universal genetic code (default)  
acet: Acetabularia  
ccyl: Candida cylindrica  
tepa: Tetrahymena, Paramecium,  
Oxytrichia, Stylonychia, Glaucoma  
eupl: Euplotes  
mlut: Micrococcus luteus  
mysp: Mycoplasma, Spiroplasma  
mitocan: canonical mitochondrial code  
mitovrt: Vertebrates - mitochondrial code  
mitoart: Arthropods - mitochondrial code  
mitoech: Echinoderms - mitochondrial code  
mitomol: Molluscs - mitochondrial code  
mitoasc: Ascidians - mitochondrial code  
mitonem: Nematodes - mitochondrial code  
mitopla: Plathelminths - mitochondrial code  
mitoyea: Yeasts - mitochondrial code  
mitoeua: Euascomycetes - mitochondrial code  
mitopro: Protozoans - mitochondrial code

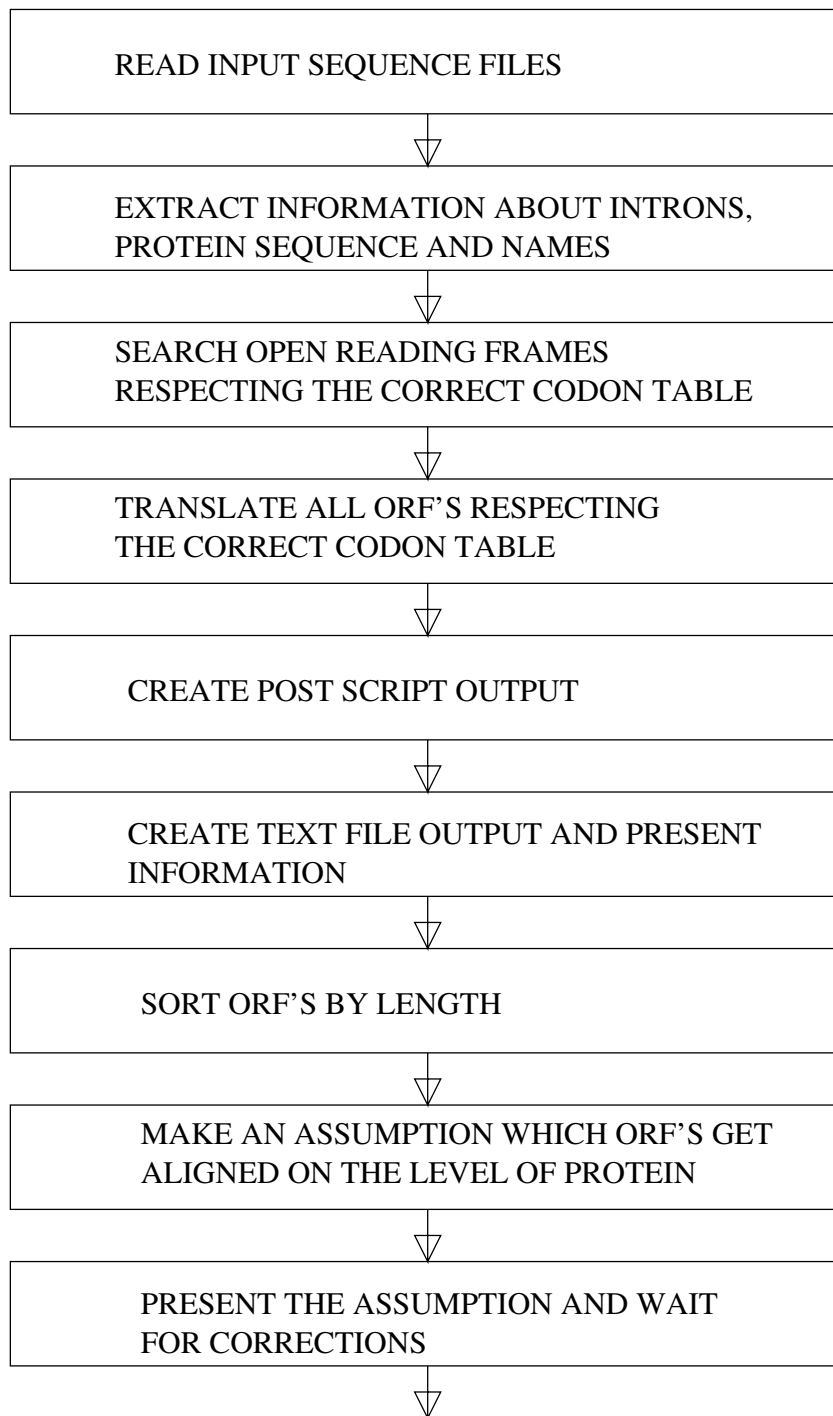


An example for the PostScript output of RAL2. The figure shows a graphical representation of the found open reading frames of three unrelated sequences. The genomes of hepatitis B virus, HIV1 and the measles virus.

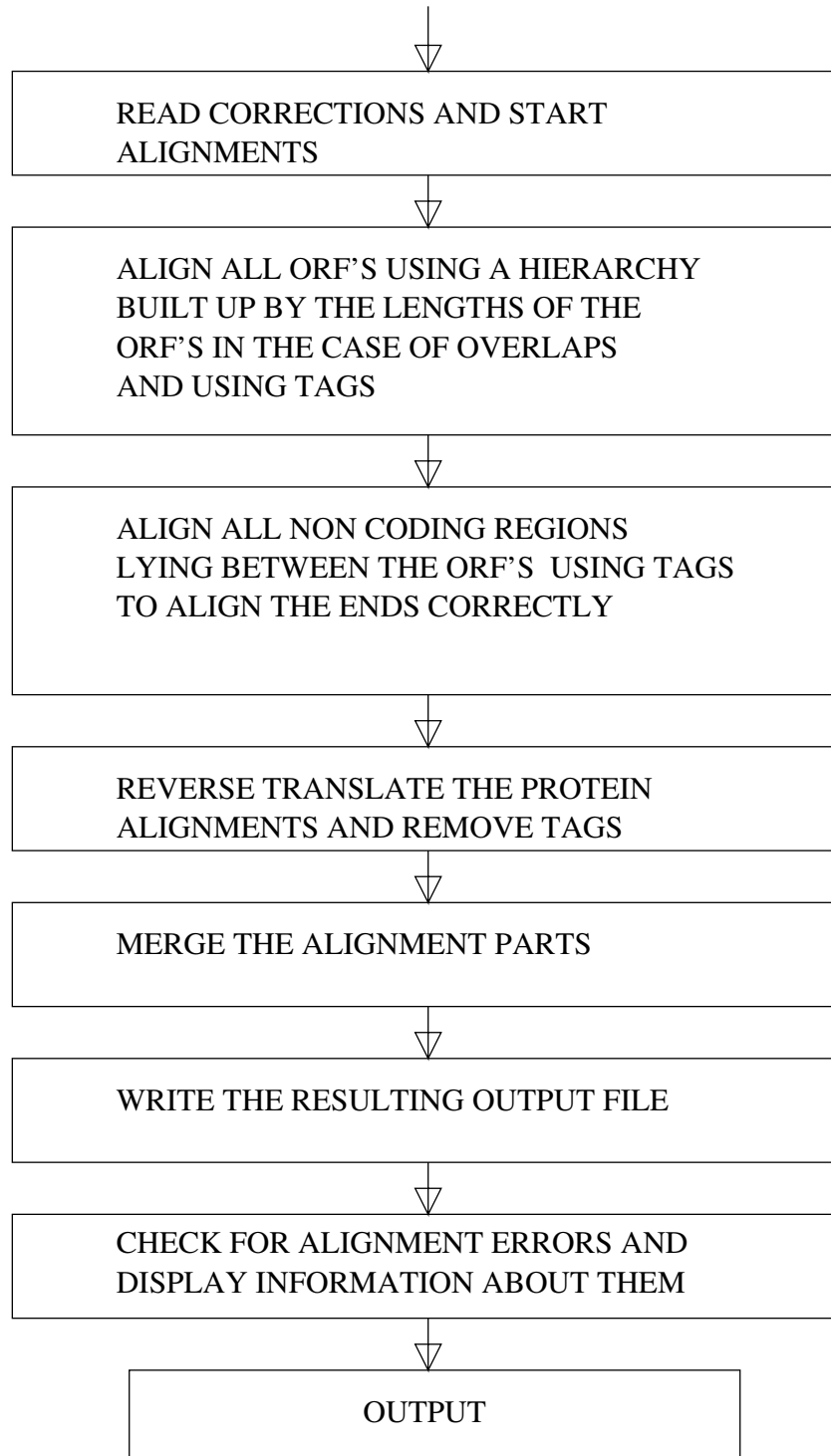
```
File name: Data/AB014366.gbf
Codon Table: univ
Sequence name: > AB014366 [3227bp]
GenBank file information derived exons, if available:
no exons
GenBank file derived CDS marked open reading frames:
START - STOP:
1374 - 1850
1826 - 2464
2319 - 3227
1 - 1623
2860 - 3227
1 - 835
Open reading frames derived by ral2 automatic search:
START - STOP:
155 - 835
421 - 1623
1374 - 1850
1913 - 2464
```

An example for the text file output of RAL2. Various data about the input sequences are shown: the names of files and sequences, the codon table used, the start and stop codons of found and Genbank derived open reading frames, and, if available, the exon data.





This flow chart shows the first main steps of RALIGN.



This flow chart shows the final main steps of RALIGN.

## OUTPUT:

ral.aln: the resulting  
alignment file

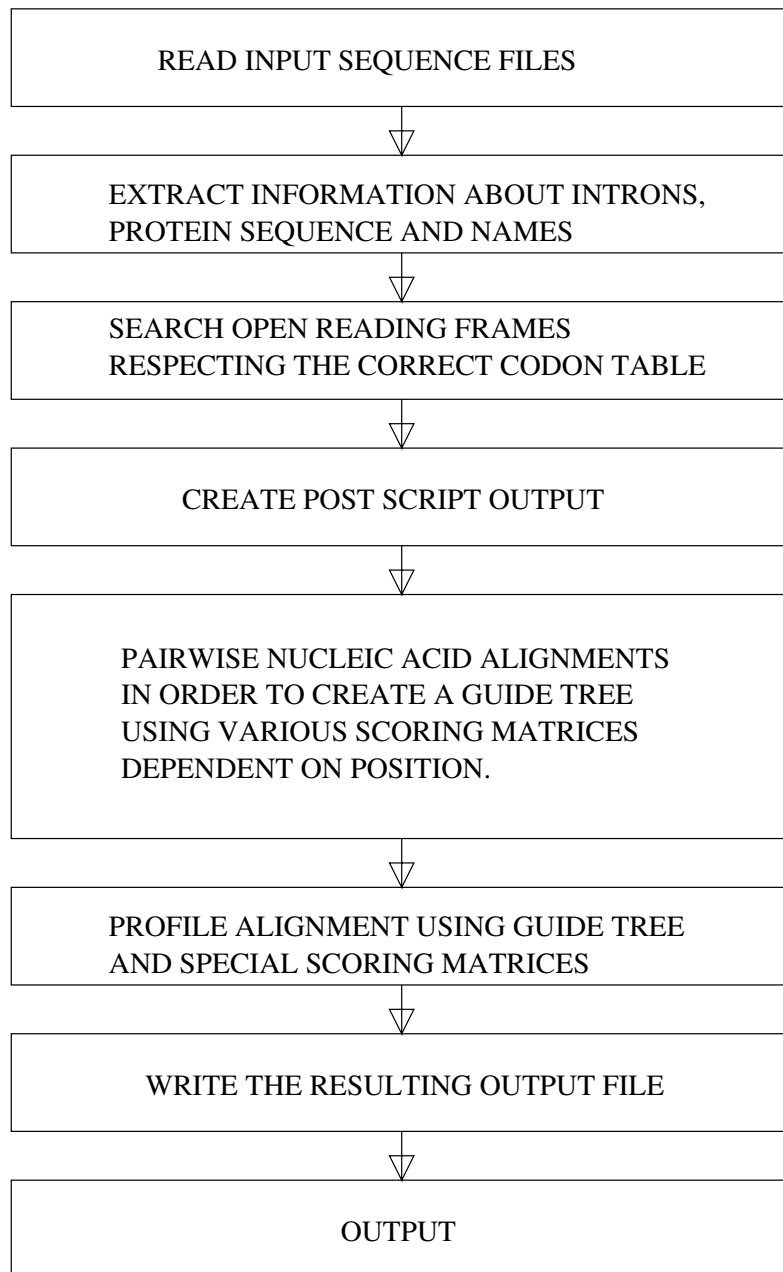
ORF.ps: the Post Script  
output

ORF.ral: the text output

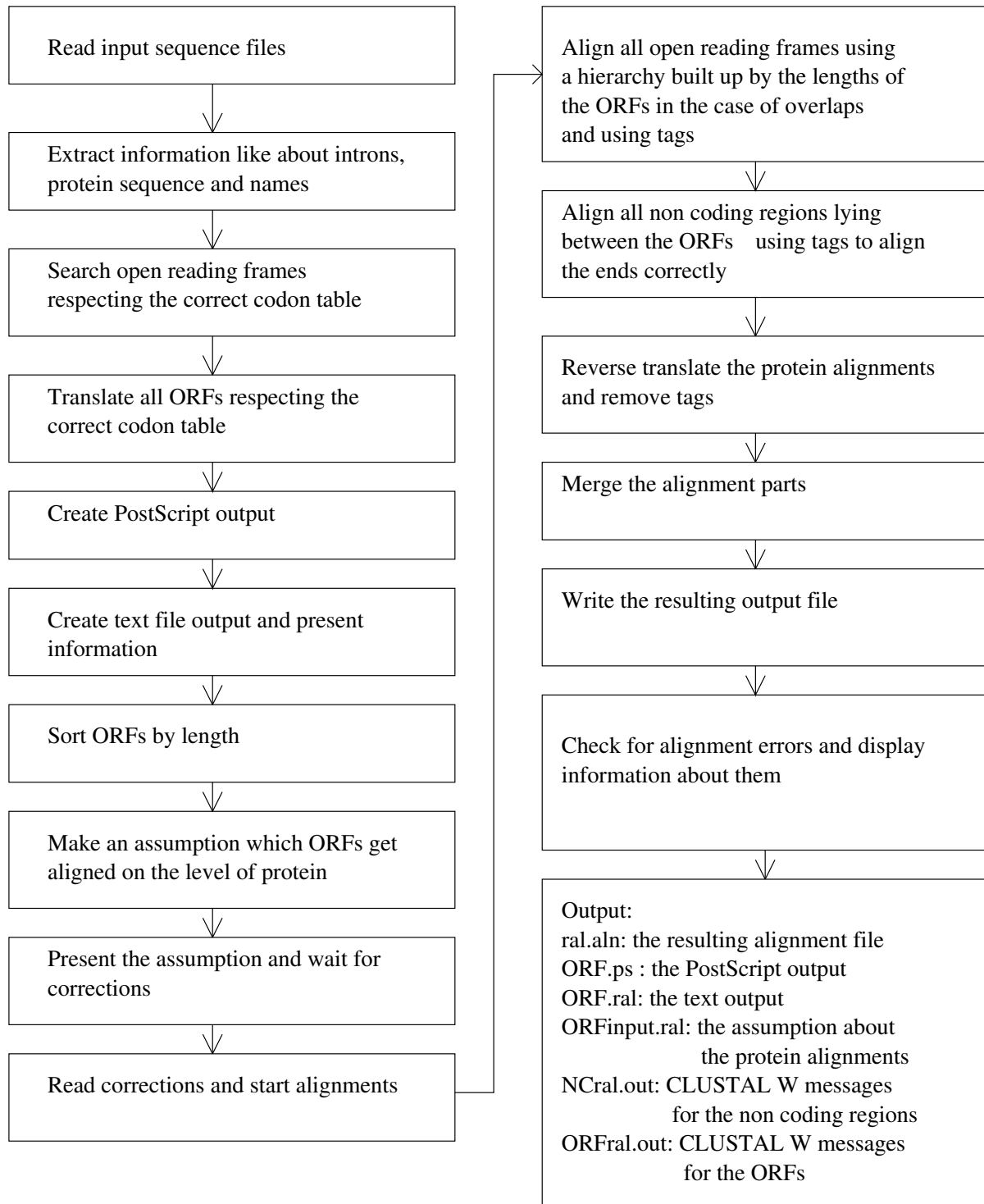
ORFinput.ral: the  
assumption about the  
protein alignments

NCral.out: CLUSTALW  
messages for the non  
coding regions

ORFral.out: CLUSTALW  
messages for the ORF's



This flow chart gives a graphical representation of the main steps of ra12.



This flow chart shows the main steps of RALIGN.

# RALIGN: improved alignments

```
ADI-MAL      UAAACUGCACUAAUGUGAAUGGGACUGCUGU---GAAUGGGACUAAUGCUGGGAGUAAUA
AE-90CF402  UACAUUGUACCAAG--CU-----AGUUUUACUAAUGCCA-----
AE-CM240     UAAAUUGUACCAAG--CU-----AAUUUGACCAAUGGCAGUAGCAAAA
B-896       UAAAUUGCACUAAUUUGAAUAUCUA-----AGAAUACUACUAAUCCACUAGUAGCA
B-ACH320A   UAAAUUGCACUGAU---UUUG-----GGAAUGCUACUAAUACCACUAGUAGUA
B-BCSG3     UAAAUUGCACUGAUGAGUUGA-----AGAAUGCUACUAAUACCACUAGUACUA
B-CAM1      UAAAUUGCACUAAUG-UA-----AAUAAUACUAGGACCAAUAGUAGUG
B-D31       UAAAUUGCACUGAUC-UGAAG-----AAUGCUACUAAUACCAAUAAUAGUA
B-HIV1AD8   UAAAUUGCACUGAUU-UGAGG-----AAUGUUACUAAUAUCAAAUAAUAGUA
B-HXB2      UAAAGUGCACUGAU---UUGA-----AGAAUGAUACUAAUACCAAUAGUAGUA
B-JRCSF     UAAAUUGCACAAAGAUGUGAA-----UGCUACUAAUACCACUAGUAGUA
B-LAI       UAAAGUGCACUGAU---UUGG-----GGAAUGCUACUAAUACCAAUAGUAGUA
B-MANC      UAGAUGUGCACUGAUUAUGUAG-----GGAAUGCUACUAAUACCACUAGCACUA
B-OYI       UAGAUGUGCACUGAUGUUAUA-----CCACUAGUAGUAGUUUGAGGAAUGCUA
B-SF2       UAAAUUGCACUGAUU-UGGGG-----AAGGCUACUAAUACCAAUAGUAGUA
B-WEAU      UAAAUUGCACUAAUGUGAAUGUGACUAAUUUGAAGAAUGAGACUAAUACCAAUAGUAGUA
B-YU2       UAAAUUGCACUGAU---UUAA-----GGAAUGCUACUAAUACCACUAGUAGUA
B-pNL43     UAAAGUGCACUGAU---UUGA-----AGAAUGAUACUAAUACCAAUAGUAGUA
D-ELI       UAAACUGUAGUGAU-----GAAU-UGAGGAACA--AUGGCACUA
D-NDK       UAAACUGCACUGAU-----GAAU-UGAGGAACAGCAAGGGCA--A
O-ANT70     UGGAGUGU-----ACAAACAUAGCUGGAACAA
O-MVP5180   UGAACUGUGUAGAUC-----U-----GCAAAACAAAUAACAGGCCUUAU
SIVCPZGAB   UGCAGUGCAGUAAGG-----CUAACUUUAGCCAGGCAAAAAACCUAA
```

```
ADI-MAL      UAAACUGCACUAAUGUGAAUGGGACUGCUGUGAAUGGGACUAAUGCUGGGAGUAAUAGGA
AE-90CF402  UACAUUGUACCAAG-----GCUAGUUUUACUAAUGCCA
AE-CM240     UAAAUUGUACCAAU-----GCUAAUUUGACCAAUGGCA
B-896       UAAAUUGCACUAAUUUGAAUAUC-----ACUAAAGAAUACUACUAAUCCCA
B-ACH320A   UAAAUUGCACUGAUUUUGGG-----AAUGCUACUAAUACCA
B-BCSG3     UAAAUUGCACUGAU-----GAGUUGAAGAAUGCUACUAAUACCA
B-CAM1      UAAAUUGCACUAAUGUAAAU-----AAUACUAGGACCAAUAGUAGUAGUU
B-D31       UAAAUUGCACUGAUCUGAAG-----AAUGCUACUAAUACCA
B-HIV1AD8   UAAAUUGCACU-----GAUUUGAGGAAUGUUACUAAUAUCA
B-HXB2      UAAAGUGCACUGAUUUUGAAG-----AAUGAUACUAAUACCA
B-JRCSF     UAAAUUGCACAAAGAUGUG-----AAUGCUACUAAUACCA
B-LAI       UAAAGUGCACUGAUUUUGGG-----AAUGCUACUAAUACCA
B-MANC      UAGAUGUGCACUGAUUAUGUAGG-----AAUGCUACUAAUACCA
B-OYI       UAGAUGUGCACUGAUGUUAUACCACUAGUAGU--AGUUUGAGGAAUGCUACUAAUACCA
B-SF2       UAAAUUGCACUGAU-----UUGGGGAAGGCUACUAAUACCA
B-WEAU      UAAAUUGCACUAAUGUGAAUGUGACU-----AAUUUGAAGAAUGAGACUAAUACCA
B-YU2       UAAAUUGCACUGAU-----UUAAGGAAUGCUACUAAUACCA
B-pNL43     UAAAGUGCACUGAUUUUGAAG-----AAUGAUACUAAUACCA
D-ELI       UAAACUGUAGUGAUGAAUUGAGGAACAUGGC-----ACUAGGGGAACAUGUCA
D-NDK       UAAACUGCACUGAUGAAUUGAGGAAC-----AGCAAGGGCAUUGGGAAGG
O-ANT70     UGGAGUGU-----ACAAACAUAGCUGGAACAA
O-MVP5180   UGAACUGUGUAGAUCUGCAA-----ACAAAUAACAGGCCUUAU
SIVCPZGAB   UGCAGUGCAGUAAGGCUAACUUUAGC-----CAGGCAAAAAACCUAACAAACCAGA
```

44% less gaps

# RALIGN: improved alignments

```
ADI-MAL      GGGAGUAAUAGGACUA-----AUGCAG--AAUUGAAA--AUGGAAUUGGAGAAGUGAAAAAC
AE-90CF402   A-----CCAGUG-----ACAGAAUA---AAAAUG---GAAGAUGCAGUAAGAAAC
AE-CM240     AGUAGCAAAAACCA AUGU-----CUCUAACAUAUUAGGAAUAUA---ACAGAUGAAGUAAGAAAC
B-896       ACUAGUAGCA-----GCUGGGGAAUGA-----UGGAGAAAGGAGAAAUAAAAAU
B-ACH320A   ACUAGUAGUA-----GCGGGUUUAUUAUA---G-AGAAAGGAGAAAUAAAAAAC
B-BCSG3     ACUAGUACUAAUACCCCUAGUGGUAGCUGGAAAAAGAU-----GGAAAGAGGAGAAAUAAAAAAC
B-CAM1      AAUAGUAGUGAUU-----GGGACAGGAGGGAAGGAGAAAGAUG---AAAGGAGAAAUAAAAAAC
B-D31       AAUAAUAGUAGU-----UGGACGAUGACAGGAGAAUUG---AAAGGAGAAAUAAAAAAC
B-HIV1AD8   AAUAAUAGUAGU-----GAGGGAA-----UG---AGAGGAGAAAUAAAAAAC
B-HXB2      AAUAGUAGUA-----GCGGGGAAUGAUA---AUGGAGAAAGGAGAGAUAAAAAAC
B-JRCSF     ACUAGUAGUA-----GUGAGGGAAUGA-----UGGAGAGAGGAGAAAUAAAAAAC
B-LAI       AAUAGUAGUAAUACCAAUAGUAGUAGCGGGGAAUUGAUG---AUGGAGAAAGGAGAGAUAAAAAAC
B-MANC      ACUAGCACUAAUAAUACCCGUAGUGGAAGUUGG---GGAGCGAUG---AGAGGGGAAAUAAAAAAC
B-OYI       AGGAAUGCUCUAAUACCAAGUAGUAGUUGG---GAAACGAUGGAGAAAGGAGAAUUAAAAAAC
B-SF2       AAUAGUAGUAAU-----GGAAA--GAAGAAA-----UA---AAAGGAGAAAUAAAAAAC
B-WEAU      AAUAGUAGUAGU-----GGAGGGGAAAGA-----UGGAGGAGGGGAGAAUAAAAAAC
B-YU2       ACUAGUAGUA-----GCUGGGAAACGAU-----GGAGAAAGGAGAAAUAAAAAAC
B-pNL43     AAUAGUAGUA-----GCGGGGAAUGAUA---AUGGAGAAAGGAGAGAUAAAAAAC
D-ELI      -AUGGCACUA-U-----GGGGAACAAUGUCACUACAGAGGAGAAAGGAUAAAAAAC
D-NDK       AAGGGCA--A-U-----GGG-----AAGGUA---GAAGAGAGGAGAAAUAAAAAAC
O-ANT70     GCUGGAACAA-----CAAUGA-----AAACCUUUAUGAAGAAG
O-MVP5180   ACAGGCCUAU-----UAAAUGAGA-----CAAUAAAUGAGAUGAGAAAU
SIVCPZGAB   AAAAAACCUAA-----CAAACCAGACAUCUUCUC-CGCCUCUGGAAUAAAAAAC
```

```
ADI-MAL      ACUAAUGCAGAAUUGAAAAUGGAA-----AUUGGAGAGUGAAAAAC
AE-90CF402   ACCAGUGACAGA-----AUAAAAUUGGAAGUAGCAGUAAGAAAC
AE-CM240     AGUAGCAAAAACCA AUGUCUCU-----AACAUAAUAGGAAUAUAACAGAUAGUAAGAAAC
B-896       ACUAGUAGCAGC-----UGGGGA-----AUGAUGGAGAAAGGAGAAAUAAAAAU
B-ACH320A   ACUAGUAGU-----AGCGGG---GUUAUAAUAGAGAAAGGAGAAAUAAAAAAC
B-BCSG3     ACUAGUACUAAU--ACCCUAGUGGUAGCUGGAAAAAGAUUGGAAAGAGGAGAAUAAAGAAC
B-CAM1      UGGGACAGGAGGGAAGGAGAA-----AAG---AUGAAAGGAGAAAUAAAAAAC
B-D31       AAUAAUAGUAGUUGGACGAUGACAGGA-----GAAUUGAAAGGAGAAAUAAAAAAC
B-HIV1AD8   AAUAAUAGUAGUGAGGGA-----AUGAGAGGAGAAAUAAAAAAC
B-HXB2      AAUAGUAGU-----AGCGGGGAAUGAUAUUGGAGAAAGGAGAGAUAAAAAAC
B-JRCSF     ACUAGUAGUAGU-----GAGGGA-----AUGAUGGAGAGAGGAGAAAUAAAAAAC
B-LAI       AAUAGUAGUAAUACCAAUAGUAGUAGCGGGGAAUUGAUGAUGGAGAAAGGAGAGAUAAAAAAC
B-MANC      ACUAGCACUAAUAAUACCCGUAGUGGAAGUUGG---GGAGCGAUGAGAGGGGAAAUAAAAAAC
B-OYI       ACAAGUAGU-----AGUUGGGAAACGAUGGAGAAAGGAGAAUUAAAAAAC
B-SF2       AAUAGUAGU-----AAUUGGAAAGAAAGAAUAAAGGAGAAAUAAAAAAC
B-WEAU      AAUAGUAGUAGUGGAGGGGAA-----AAGAUGGAGGAGGAGAAUAAAAAAC
B-YU2       ACUAGUAGU-----AGCUGGGAAACGAUGGAGAAAGGAGAAAUAAAAAAC
B-pNL43     AAUAGUAGU-----AGCGGGGAAUGAUAUUGGAGAAAGGAGAGAUAAAAAAC
D-ELI      ACUACAGAGGAGAAAGGAAUG-----AAAAAC
D-NDK       GUAGAAGAGGAGGAAAAAAGG-----AAAAAC
O-ANT70     ACAAUUGAAAACCUU-----AUGAGAAG
O-MVP5180   UUAAAUGAGACAAUAAU-----GAGAUGAGAAAU
SIVCPZGAB   ACAUCUUCUCCG-----CCUCUCGAAAUAAAAAAC
```

50% less gaps

# RALIGN: improved alignments

```
ADI-MAL          CAAU---AG-----AUGAUAGUGAUAAUAG-----AUAG--UACUAAUUUAGGUUAAUAAUUG
AE-90CF402      CAAUUUAGUAGUGUACAAAUAUAUAACAGUAAUACUAGUGGACAAAUAUAGUCAUAAAGUUUAGAUAAUACAUGG
AE-CM240        CAAUU--GAAG-----AUAAGAGACUAGUAGUG-----AGUUAAGGUUAAUAAUUG
B-896           CCAAU---AG-----AAAUAUCUAAUAAUA-----CUAAGUUAAGGUUAAUAAUUG
B-ACH320A       CCAAU---AG-----AUAUAUAUAUACUAAUA-----CCAGCUAUACCAGCUUAGGUUGAUAAAGUUG
B-BCSG3         CCAAU---AG-----AUAAGUUAAGAAUAG-----UACCAAUAUAGGUUGAUAAAGUUG
B-CAM1          CCAAU---AG-----AUAAGGCUAUACAAGU-----UAUACAUUAGUUAACAUGG
B-D31           CCAAU---AG-----AUAAGACAAUACUAG-----CUUAAGGUUGAUAAAGUUG
B-HIV1AD8       CCAAU---AG-----AUAAGUAAUACUAG-----CUUAAGGUUGAUAAUUG
B-HXB2          CCAAU---AG-----AUAAGUAC-----UACCAGCUUAAAGGUUGACAAGUUG
B-JRCSF         CCAAU---AG-----AUAUAUAGAAUAAUA-----CCAAUAUAGGUUAAUAAUUG
B-LAI           CCAAU---AG-----AUAAGUAC-----UACCAGCUUACGUUGACAAGUUG
B-MANC         CCAAU---AG-----AAAAGAAUACUAG-----CUUUAGAUUAGUAAAGUUG
B-OYI           CCAAU---AG-----AUAAGAAUAGUACUAA-----AUUUAGGUUAAUACAUGG
B-SF2           CCAAU---AG-----AUAAGCUAGUACUACUA-----CCAACUUAACCAACUUAAGGUUGAUACAUGG
B-WEAU         CCAAU---AG-----AUCUAGUAAUACAAG-----CUUAAGGUUGAUAAUUG
B-YU2           CCAAU---AG-----AUAAG-----CUAGCUUAGGUUGAUAAAGUUG
B-pNL43         CCAAU---AG-----AUAU-----ACCAGCUUAGGUUGAUAAAGUUG
D-ELI           CCAAU---AG-----ACAAGUAGUAGUACCA-----AUAG--UACCAAUAUAGGUUAAUAAUUG
D-NDK           CCAAU---AG-----ACAAUAUAUAGGACCA-----AUAG--UACUAAUUUAGGUUAAUAAUUG
O-ANT70         GAACUG--AAUGAG--ACAAGCAGCACAAUAAGACAAACAGC-----AAAAGUUAACAUUACAUAUUG
O-MVP5180       AAGGUU--AAUGA--CUCAAUGCAGUAAUUGGA-----ACAACUUAUAGGUUAAUAAUUG
SIVCPZGAB      AACCU---AG-----GGAUAGAGAACAACAC-----AUUAAGGAUAAUAAUUG
```

```
ADI-MAL          CAAUAGAUGAU-----AGUGAUAAUAGUAGUUAUAGGCUAAUAAUUG
AE-90CF402      CAAUUUAGUAGUGUACAAAUAUAUAACAGUAAUACUAGUGGACAAAUAUAGUCAUAAAGUUUAGAUAAUACAUGG
AE-CM240        CAAUUUGAA-----GAU-----AAGAAGACUAGUAGUAGUUAAGGUUAAUAAUUG
B-896           CCAUAGAAAAUACUAAU-----AAUACUAGUUAAGGUUAAUAAUUG
B-ACH320A       CCAUAGAUAAUAAUAAUACUAAU-----ACCAGCUUACCAGCUUAGGUUGAUAAAGUUG
B-BCSG3         CCAUAGAUAAUAGUAAAG-----AAUAGUACCAAUAUAGGUUGAUAAAGUUG
B-CAM1          CCAUAGAUAAAGGCU-----AAUACAAGUUUAACUUGAUACAUGG
B-D31           CCAUAGAUAAAGAC-----AAUACUAGCUUAGGUUGAUAAAGUUG
B-HIV1AD8       CCAUAGAUAAUGAU-----AAUACUAGCUUAGGUUGAUAAUUG
B-HXB2          CCAUAGAUAAUGAU-----ACUACCAGCUUAAAGGUUGACAAGUUG
B-JRCSF         CCAUAGAUAAUAGAAU-----AAUACCAAUAUAGGUUAAUAAUUG
B-LAI           CCAUAGAUAAUGAU-----ACUACCAGCUUACGUUGACAAGUUG
B-MANC         CCAUAGAAAAAGAAAG-----AAUACUAGCUUAGUUAAGUUAAGUUG
B-OYI           CCAUAGAUAAAGAAU-----GAUACUAAUUUAGGUUAAUACAUGG
B-SF2           CCAUAGAUAAUGCUAGUACUACU-----ACCAACUUAACCAACUUAAGGUUGAUACAUGG
B-WEAU         CCAUAGAUCAUGAU-----AAUACAAGCUUACGUUGAUAAUUG
B-YU2           CCAUAGAUAAU-----GCUAGCUUAGGUUGAUAAAGUUG
B-pNL43         CCAUAGAUAAU-----ACCAGCUUAGGUUGAUAAAGUUG
D-ELI           CCAUAGACAAUGAUAGUAGU-----ACCAAUAGUACCAAUAUAGGUUAAUAAUUG
D-NDK           CCAUAGACAAUAAUAAUAGG-----ACCAAUAGUACUAAUUUAGGUUAAUAAUUG
O-ANT70         GAACUGAAUGAGACAGCAGCAGCACAAUAAGACAAACAGC-----AAAAGUUAACAUUACAUAUUG
O-MVP5180       AAGGUUAAUGACUCAAUGCAGUAAUUGGAACAACA-----UAUAGUUAACUAAUUG
SIVCPZGAB      AACCUAGGGAAUGAG-----AACCAACAUUAUAGGAUAAUAAUUG
```

66% less gaps



# RALIGN: improved alignments

```
ADI-MAL      CUCUAUACAACAGGGAUAGU-----A----GGAG--AUAU-AAGAAGAGCAUUAUGUACU
AE-90CF402  UUCCAUACAACAGGAAACAU-----AAAU--GGUG--AUAU-AAGAAAAGCAUUAUGUGAA
AE-CM240     UUCUAUAGAACAGGAGAUUU-----AAUA--GGAA--AUAU-AAGAAAAGCAUUAUGUGAG
B-896        UUUUAUGCAAGAAGAAACAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-ACH320A   UUUUAUGCAACAGGACAAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-BCSG3      UAUUUAACAACAGGAGAAAU-----AGUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-CAM1       GUUUUUGCAACAGACAGAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-D31        UUUUAUACAAAAGGAAAAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-HIV1AD8    UUUUAUACAACAGGAGACAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGCAAC
B-HXB2       UUUGUUACAACUAGGAAAAAU-----AG----GAA--AUAU-GAGACAAGCACAUUGUAAC
B-JRCSF      UUUUAUACAACAGGAGAAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-LAI        UUUGUUACAACUAGGAAAAAU-----AG----GAA--AUAU-GAGACAAGCACAUUGUAAC
B-MANC       UUUCAUGUAACAAGAGCCGU-----ACA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-OYI        UUUCAUACAACAAAACAAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-SF2        UUUCAUACAACAGGAAGAAU-----AAUA--GGAG--AUAU-AAGAAAAGCACAUUGUAAC
B-WEAU       CUUUUAACAACAGGAGAAAU-----AAUA--GGAG--AUAU-AAGACGAGCACAUUGUAAC
B-YU2        UUGUAUACAACAGGAGAAAU-----AAUA--GGAG--AUAU-AAGACAAGCACAUUGUAAC
B-pNL43      UUUGUUACAACUAGGAAAAAU-----AG----GAA--AUAU-GAGACAAGCACAUUGUAAC
D-ELI        CUCUAUACUACAAGAUCAA----GAUCAUA-----AU-AGGACAAGCACAUUGUAAU
D-NDK        CUCUAUACAACUACAGGAAAAAAGAAAGAAAACAGGAU--ACA-AGGACAAGCACAUUGUAAA
O-ANT70      UACAGCAUGGGAAU---AGGGGAAACAGCAGGAAAC--AGCUCAGGGCAGCUUAUUGCAAG
O-MVP5180    CGCAGUAUGACACUUAAAAGAAGUAACAAUACAUCACCAAGAUCAGGGUAGCUUAUUGUACA
SIVCPZGAB    UUUUAUAAUUAAGAAAAUGU-----AGUA--GGAG--AUAC-CAGAUUCGCCUACUGUAAG
```

```
ADI-MAL      CUCUAUACAACAGGGAUAGUAGGA-----GAUUAUAGAAGAGCAUUAUGUACU
AE-90CF402  UUCCAUACAACAGGAAACAUUAUUGGU-----GAUUAUAGAAAAGCAUUAUGUGAA
AE-CM240     UUCUAUAGAACAGGAGAUUAUUAUAGGA-----AAUUAUAGAAAAGCAUUAUGUGAG
B-896        UUUUAUGCAAGAAGAAACAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-ACH320A   UUUUAUGCAACAGGACAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-BCSG3      UAUUUAACAACAGGAGAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-CAM1       GUUUUUGCAACAGACAGAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-D31        UUUUAUACAAAAGGAAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-HIV1AD8    UUUUAUACAACAGGAGACAUUAUAGGA-----GAUUAUAGACAAGCACAUUGCAAC
B-HXB2       UUUGUUACAACUAGGAAAA--AUAGGA-----AAUUAUGAGACAAGCACAUUGUAAC
B-JRCSF      UUUUAUACAACAGGAGAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-LAI        UUUGUUACAACUAGGAAAA--AUAGGA-----AAUUAUGAGACAAGCACAUUGUAAC
B-MANC       UUUCAUGUAACAAGAGCCGUAAACAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-OYI        UUUCAUACAACAAAACAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-SF2        UUUCAUACAACAGGAAGAAUUAUAGGA-----GAUUAUAGAAAAGCACAUUGUAAC
B-WEAU       CUUUUAACAACAGGAGAAAUUAUAGGA-----GAUUAUAGACGAGCACAUUGUAAC
B-YU2        UUGUAUACAACAGGAGAAAUUAUAGGA-----GAUUAUAGACAAGCACAUUGUAAC
B-pNL43      UUUGUUACAACUAGGAAAA--AUAGGA-----AAUUAUGAGACAAGCACAUUGUAAC
D-ELI        CUCUAUACUACAAGAUCAGAUC-----AUAAUAGGACAAGCACAUUGUAAU
D-NDK        CUCUAUACAACUACAGGAAAAAAGAAAGAAAACAGGA---UACAUAGGACAAGCACAUUGUAAA
O-ANT70      UACAGCAUGGGAAUAGGGGAAACAGCAGGAAACAGC-----UCAAGGGCAGCUUAUUGCAAG
O-MVP5180    CGCAGUAUGACACUUAAAAGAAGUAACAAUACAUCACCAAGAUCAGGGUAGCUUAUUGUACA
SIVCPZGAB    UUUUAUAAUUAAGAAAAUGUAGUAGGA-----GAUACCAGAUUCGCCUACUGUAAG
```

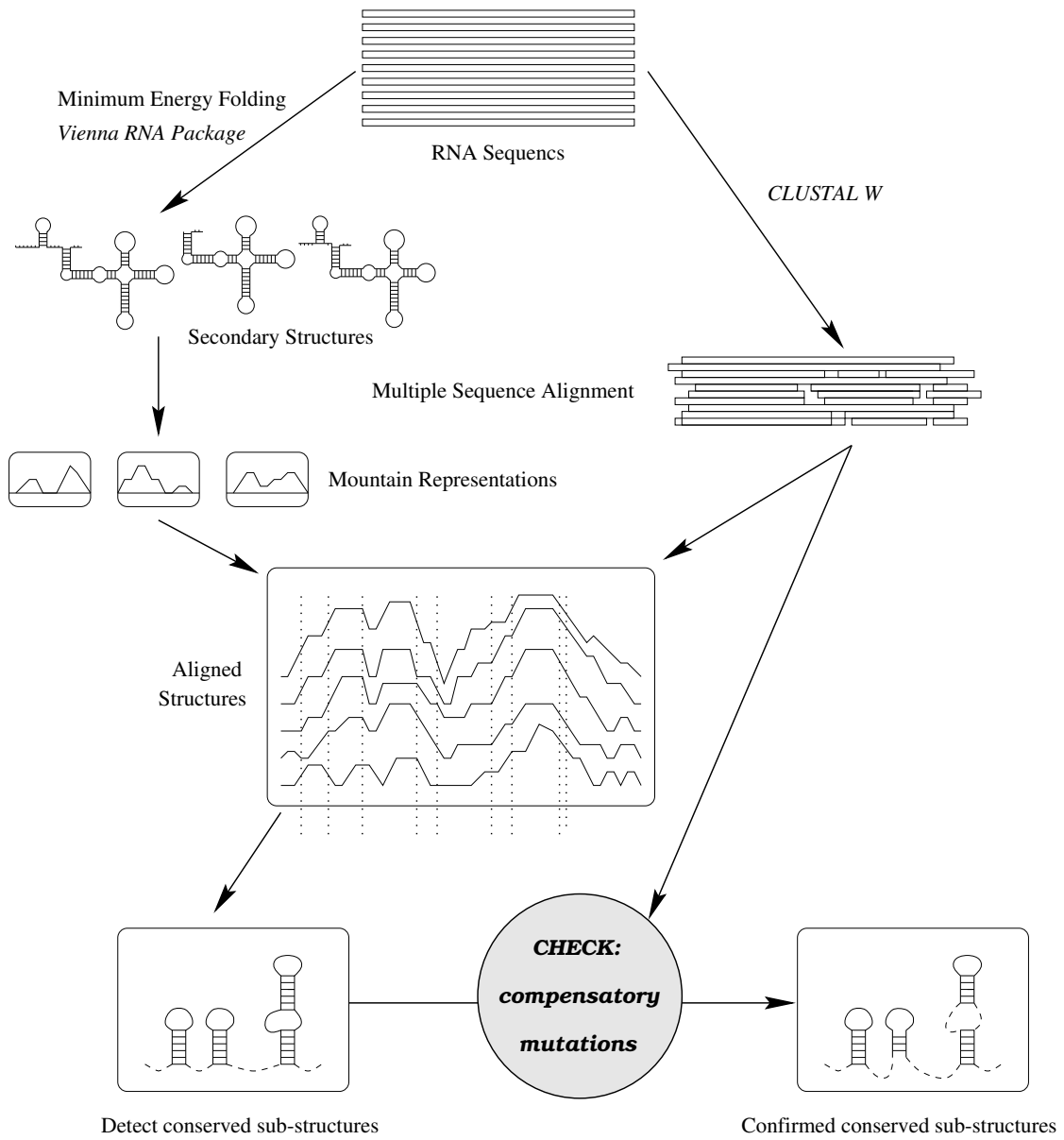
70% less gaps

# RALIGN: improved alignments

```
ADI-MAL      AU-AGUACAUGGCAGAAUAAUGGGUC---AAGA-----CU-AA--GU---AAUAGCACAGAGUC---AACUGGUAGUAUCACACUCCCAUG
AE-90CF402  AU-AGUACUUGGAUA-----AAUGGAACCAUGCAGGAGGUU--AAUGGCACAAACUC---A---GGCAAUAUCACACUCCCAUG
AE-CM240     AU-AAUACUUGCCUAG-----GAAAUGAAACCAUGCCGGGGUGU--AAUGACAC-----UAUCACACUCCCAUG
B-896       AU-AGUACUUGGAUU-----G-----U-UA-----CUGGAGGGACA--AAUGGCACUGAAGG---AAAUGACAUAUACACUCCCAUG
B-ACH320A   AU-AGUACUUGG-----AAUGAUACUGGGAUUGUUA--CUGAAAGGUCA--AAUAACAUGA-----AAU-----AUCACACUCCCAUG
B-BCSG3     AU-AGUACUUGGCUGGGAUAAUACUUGGAUAGUAGUGCUGAAAGGUCA--GAUGACACUGGAGG---AAU-----AUCACACUCCCAUG
B-CAM1      AU-ACUACUUGGCUGUUUAAUGGUACUUGGAUAGUA---CUGAAGGGUUA--AAUAACACUGAAG---AAU-----AUUACACUCCCAUG
B-D31       AU-AGUACUUGGAU-----GAUA---CUAAAGAGUCA--AAUAACACAAAU-----GGAACUAUCACACUCCCAUG
B-HIV1AD8   AU-AGUACUUGGAUUUUAAUGGUACUUGGAUUUUA---CACAUCG-----AAUGGUACUGAAGG---AAAUGACACUAUCACACUCCCAUG
B-HXB2      AU-AGUACUUGGUUU---AAUAGUACUUGGAGUA-----CUGAAGGGUCA--AAUAACACUGAAGG---AAGUGACACAAUCACCCUCCCAUG
B-JRCSF     AU-AGUACUUGGAU-----G-----A-UA-----CUGAAAAGUCA--AGUGGCACUGAAGG---AAUUGACACCAUCAUACUCCCAUG
B-LAI       AU-AGUACUUGGUUU---AAUAGUACUUGGAGUA-----CUGAAGGGUCA--AAUAACACUGAAGG---AAGUGACACAAUCACACUCCCAUG
B-MANC      AU-AGUACUUGGAUACUGGG-----AAUGAUA---CUAGAGAGUCA--AAUGACACAAUUA---UACUGGAAUUAUCACACUCCCAUG
B-OYI       AU-AGUACUUGGAU-----GAUA---CUACAAGGGCA--AAUAGCACUGAA-----GUAACUAUCACACUCCCAUG
B-SF2       AU-AAUACAUGGAGGUUAAU-----CACACU-G-----AA-GGAACUAAGG---AAAUGACACAAUCAUACUCCCAUG
B-WEAU      AU-AGUACUUGGCAUGCAAUGGUACUUGGAAGAUA---CUGAAGGGGCA--GAUAACAAU-----AUCACACUCCCAUG
B-YU2       -----CUUGG-----AAUGAUACUAGAAA-----GUUA--AAUAACACUGGAAG---AAU-----AUCACACUCCCAUG
B-pNL43     AU-AGUACUUGGUUU---AAUAGUACUUGGAGUA---CUGAAGGGUCA--AAUAACACUGAAGG---AAGUGACACAAUCACACUCCCAUG
D-ELI      AU-AGUACAUGGAUUAU---UAGUGCAUGGAUUAUUA---UACAGAGUCA--AAUAAUAGCACAAA---CAC---AAACAUCACACUCCCAUG
D-NDK      AU-AGUACAUGGAU---CA--GACUAAUAG---UACAGGGUUC--AAUAAUGGCACAG-----UCACACUCCCAUG
O-ANT70     AUUA-UACUUUUUCA-----UGUAACCGGAACCCUGUAGUGUUAGUAAUGUUAUGUCAAGG-----UAACAAUGGCACUCUACCUUG
O-MVP5180  ACUA-UACUUUUUCA-----CUGUACAAGUCCGGAUGCCAGGAGAUCAAAGGGAGCAUAGAGACCAUAAAAAUGGUACUAUACCUUG
SIVCPZGAB  CUGACAACAUA-----CUGUACAAGUCCGGAUGCCAGGAGAUCAAAGGGAGCAUAGAGACCAUAAAAAUGGUACUAUACCUUG
CUGACAACAUA-----CA--AAUGGCAU-----AUAUACUGGCAUG
```

```
ADI-MAL      AUAGUACAUGGCAGAAUAAUGGU---GCAAGACUAAAGUAAUAGCACAGAGUCAACUGGU-----AGUAUCACACUCCCAUG
AE-90CF402  AUAGUAUUGGAUA-----AAUGGAACCAUGCAGGAGGUUAAUGGCACAAACUCA-----GGCAAUAUCACACUCCCAUG
AE-CM240     AUAAUACUUGCCUAGGA-----AAUGAAACCAUGCCGGGGUGUAAUGACACU-----AUCACACUCCCAUG
B-896       AUAGUAUUGGAUU-----GUUACUGGAGGGACAUAUGGCACUGAAGGA-----AAUGACAUAUACACUCCCAUG
B-ACH320A   AUAGUAUUGGAUAGUACUGGG-----AAUGUUAUCUGAAGGGUCAAUAACAUGAA-----AAUACACACUCCCAUG
B-BCSG3     AUAGUAUUGGGCUGGG-----AAUAAUACUUGGAUAGUAGUGCUGAAGGUCAGAUGACACUGGAGGAAUUAUCACACUCCCAUG
B-CAM1      AUACUACUUGGCUG-----UUUAAUGGUACUUGGAUAGUACUGAAGGG---UUAUUAAACACUGAAGAAUUAUUAUCACUCCCAUG
B-D31       AUAGUAUUGGAU-----GAUACUAAAGAGUCAAAUAACACAAAU-----GGAACUAUCACACUCCCAUG
B-HIV1AD8   AUAGUAUUGGAU-----UUUAAUGGUACUUGGAUUUAACACAAUUCGAAUGGUACUGAAGGAAUAGACACUAUCACACUCCCAUG
B-HXB2      AUAGUAUUGGUUUAAUAGUACU-----UGGAGUACUGAAGGGUCAAAUAACACUGAAGGA-----AGUGACACAAUCACCCUCCCAUG
B-JRCSF     AUAGUAUUGGAU-----GAUACUGAAAAGUCAAGGGCACUGAAGGA-----AAUGACACCAUCAUACUCCCAUG
B-LAI       AUAGUAUUGGUUUAAUAGUACU-----UGGAGUACUGAAGGGUCAAAUAACACUGAAGGA-----AGUGACACAAUCACACUCCCAUG
B-MANC      AUAGUAUUGGAUUAU---GGG-----AAUGAUAUCAGAGAGUCAAAUAGACACAAUUAU-----ACUGGAAUUAUCACACUCCCAUG
B-OYI       AUAGUAUUGGAU-----GAUACUACAAGGGCAAUAGCACUGAAGUA-----ACUAUCACACUCCCAUG
B-SF2       AUAAUACUUGGAGG-----UUAAAUACACUGAAGGAACUAAAGGAAU-----GAC-----ACAAUCAUACUCCCAUG
B-WEAU      AUAGUAUUGGCAU-----GCUAAUGGUACUUGGAAGAAUACUGAAGGG-----GCAGUAACAAUUAUCACACUCCCAUG
B-YU2       -----ACUUGGAU-----GAUACUAGAAAGUUAAUUAACACUGGAAGA-----AAUACACACUCCCAUG
B-pNL43     AUAGUAUUGGUUUAAUAGUACU-----UGGAGUACUGAAGGGUCAAAUAACACUGAAGGA-----AGUGACACAAUCACACUCCCAUG
D-ELI      AUAGUAUUGGAUUAUAGUGCAUGGAAUUAUUAUACAGAGUCAAAUUAUAGCACAAACACA-----AACAUCACACUCCCAUG
D-NDK      AUAGUAUUGGAUUCAGACU-----AAUAGUACAGGGUUAUUAUAGGCACA-----GUACACUCCCAUG
O-ANT70     AUUAUACCUUUUCA-----UGUAACCGGAACCCUGUAGUGUUAGUAAUGUUAUGUCA-----GGUAACAAUGGCACUCUACCUUG
O-MVP5180  ACUAUACUUUUUCA-----UGUAACAAGUCCGGAUGCCAGGAGAUCAAAGGGAGCAUAGAG-----ACCAAUAAAAAUGGUACUAUACCUUG
SIVCPZGAB  ACUAUACUUUUUCA-----ACUGACAACAUAUCAAUGGC-----AUAUAAUUAUCUGGCAUG
```

60% less gaps



Flow diagram of alidot

## Sorting Base Pairs By “Credibility”

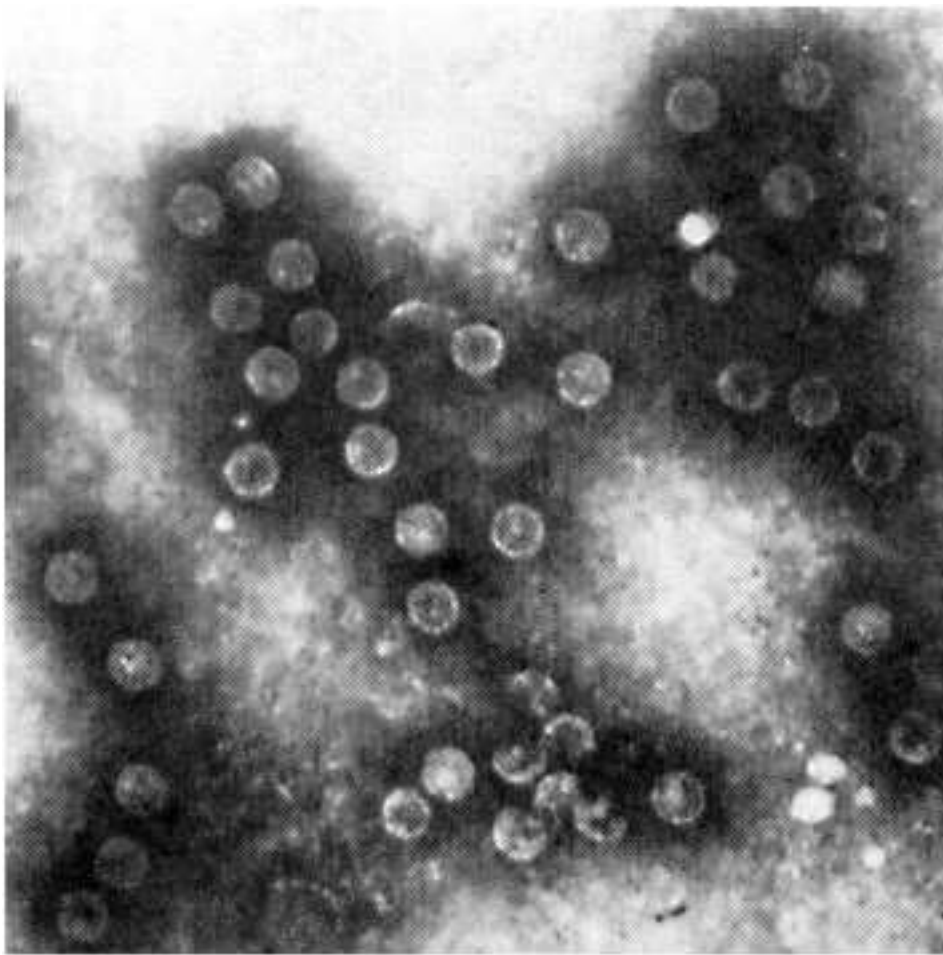
1. The more sequences are non-compatible with  $(i.j)$ , the less credible is the base pair.
2. If the number of non-compatible sequences is the same, then the pairs are ranked by the product  $\bar{p}_{i.j} \times c_{i.j}$  of the mean probability and the number of different pairing combinations.

Then we go through the sorted list and remove all base pairs that conflict with a higher ranked pair by violating conditions (i) or (ii).

## Filtering Steps

1. Remove all pairs with more than two non-compatible sequences, as well as pairs with two non-compatible sequences adjacent to a pair that also has non-compatible sequences.
2. Omit all isolated base pairs.
3. Collect the remaining pairs into helices and retain only those that satisfy the following conditions:
  - The highest ranking base pair must not have non-compatible sequences.
  - For the highest ranking base pair the product  $\bar{p}_{i,j} \times c_{i,j}$  must be greater than 0.3.  
 $\bar{p}_{i,j}$  . . . average pairing probability  
 $\bar{c}_{i,j}$  . . . number of different pairing combinations.
  - If the helix has length 2, it must not have more non-compatible sequences than consistent mutations.

The remaining list of base pairs is the conserved structure predicted by the `pfra1i` program.



Electron micrograph of hepatitis B virus particles, including virions, 20-nm spheres, and filaments.

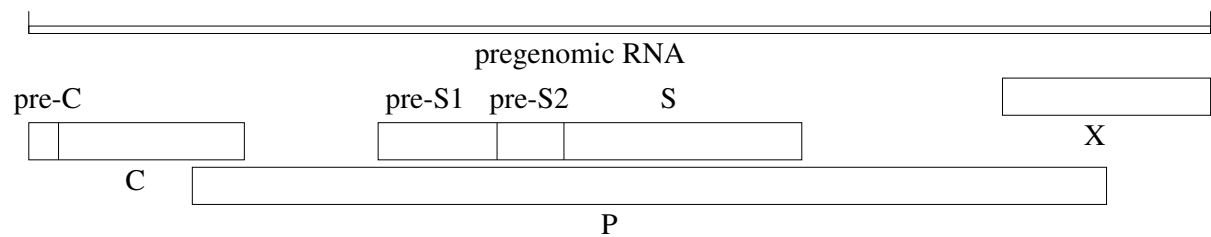


Diagram of the genome organization of hepatitis B virus with the four open reading frames (C, P, S, X).

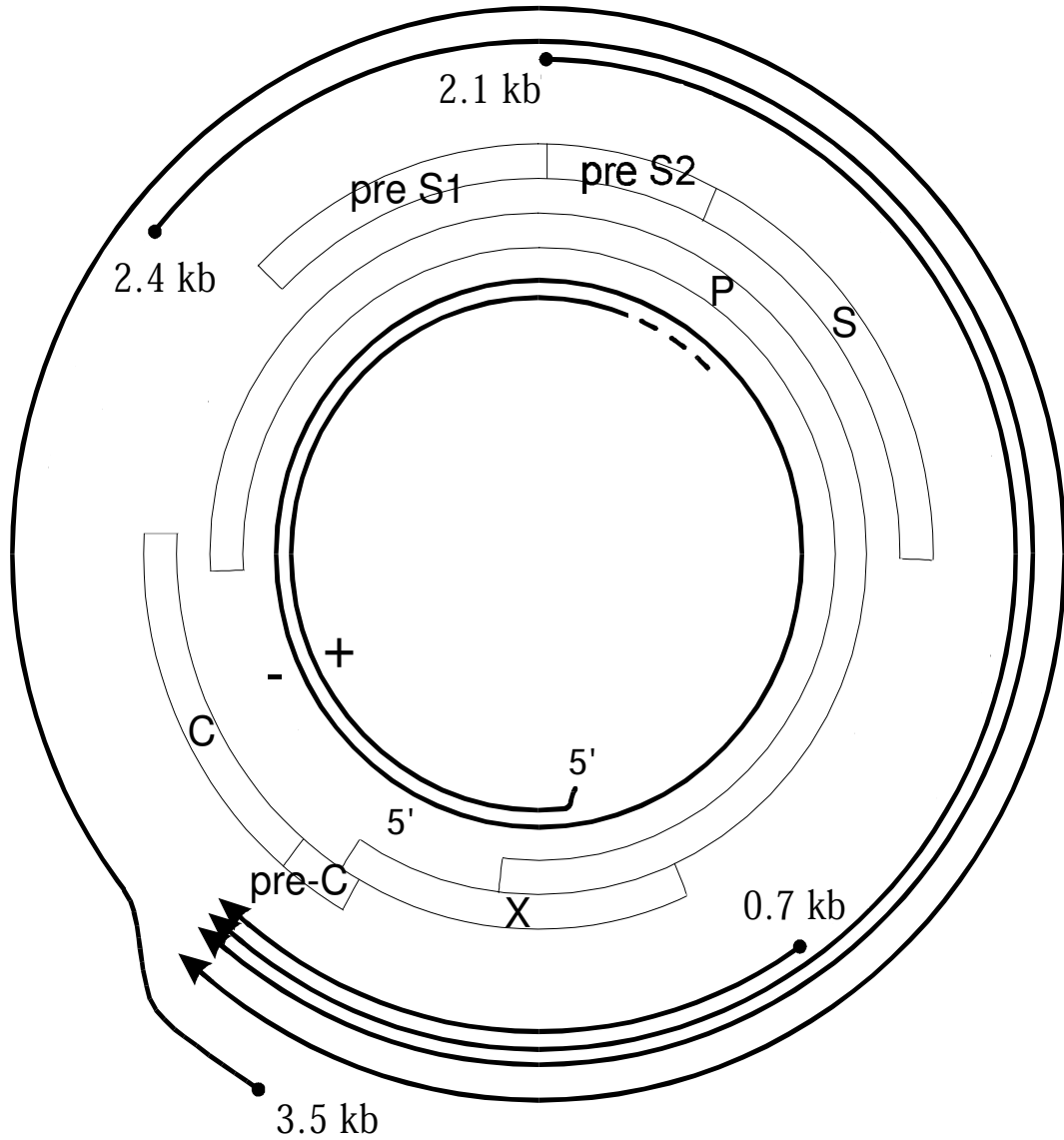
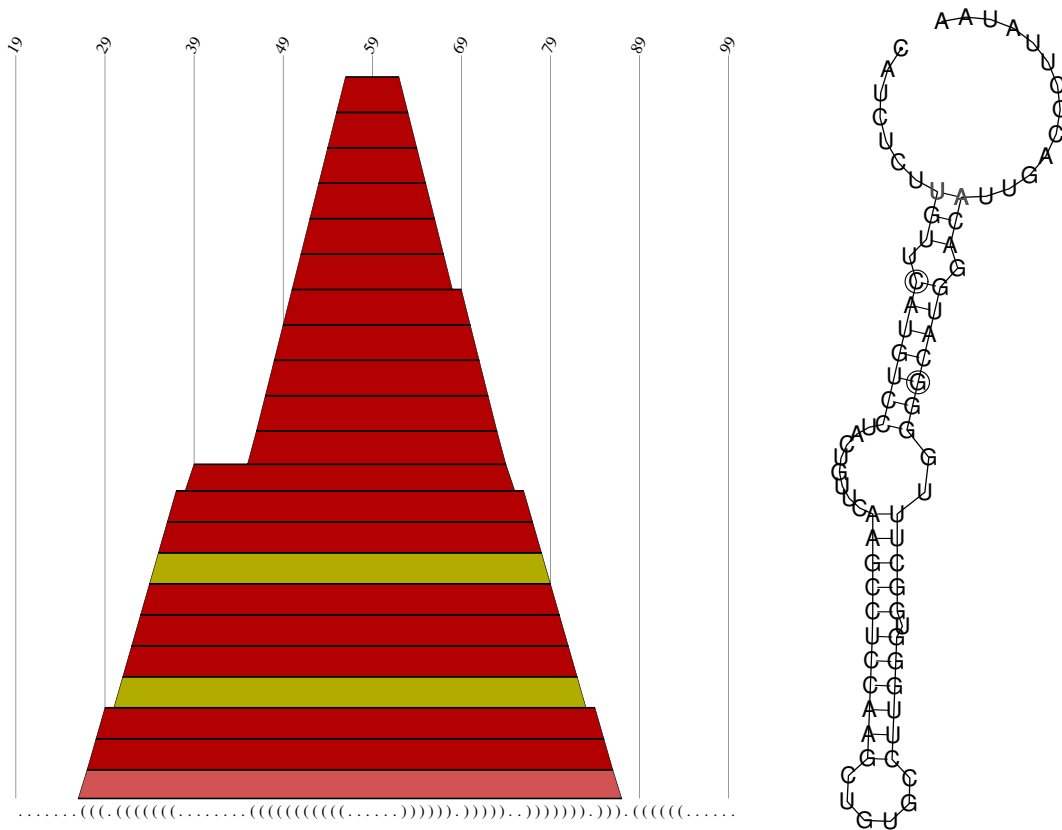


Diagram of the genome organization of hepatitis B virus indicating the DNA arrangements, the positions of the four open reading frames (C, P, S, X) and the mRNA transcripts.

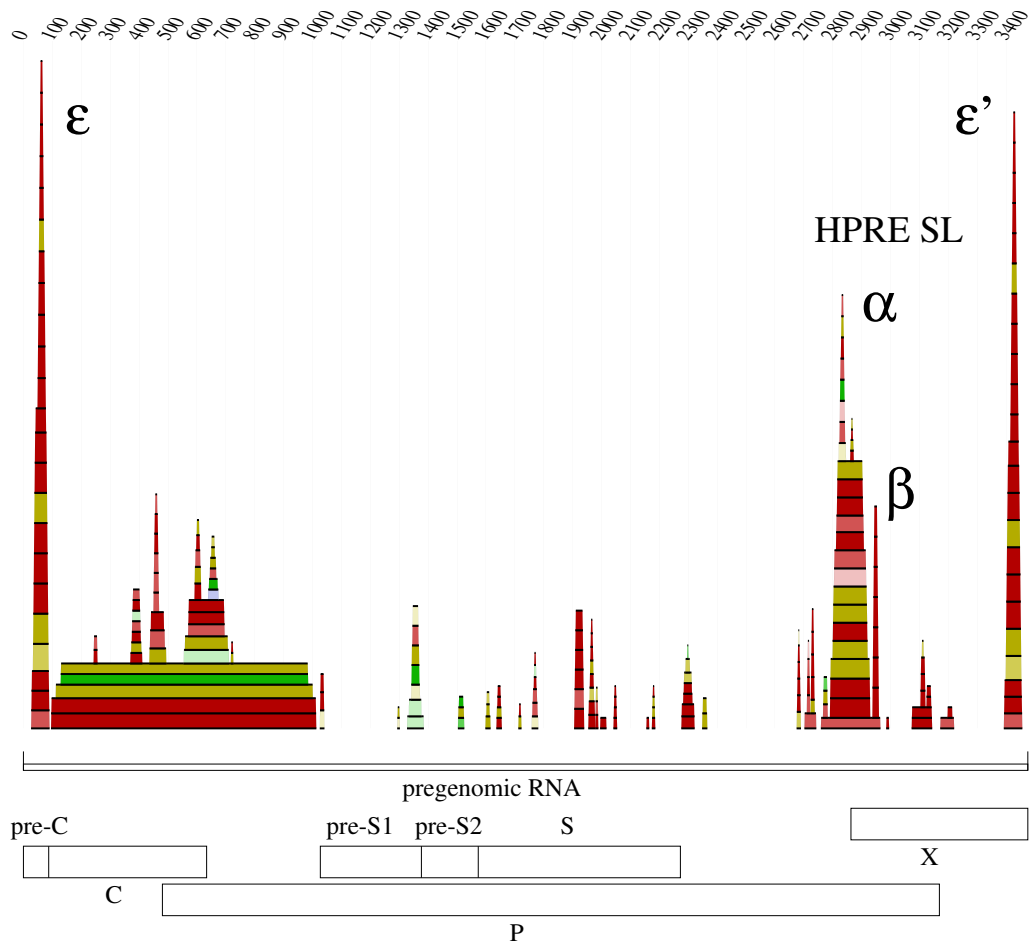


# Hepatitis B Virus RNA Pregenome

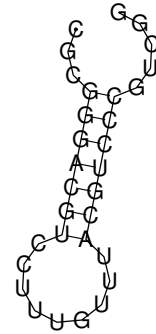
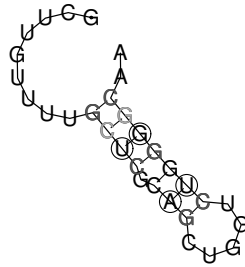
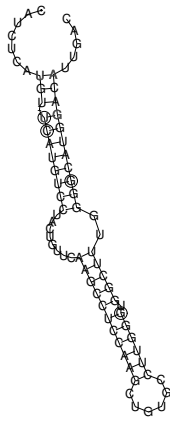


The  $\epsilon$ -element in Human Hepatitis B Virus RNA Pregenomes.

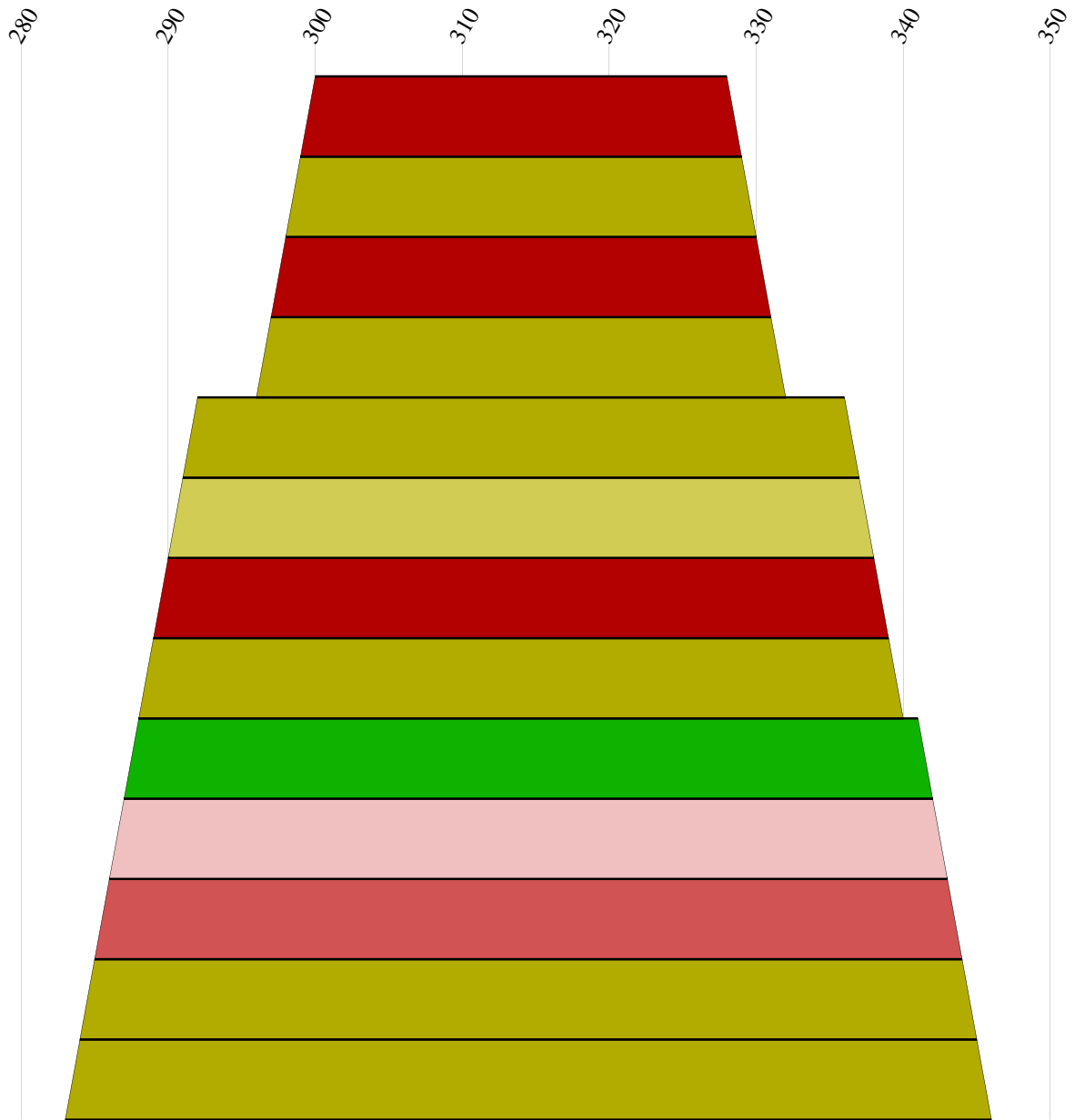
This structure is conserved among all Mammalian Hepatitis B Viruses.



The four most prominent conserved elements are the  $\epsilon$ -element and its copy at the 5' and 3' end of the pregenome and the two stem loop structures  $\alpha$  and  $\beta$  in the HPRE region. Colors indicate compensatory mutations from 0 to 3 different types of base pairs.

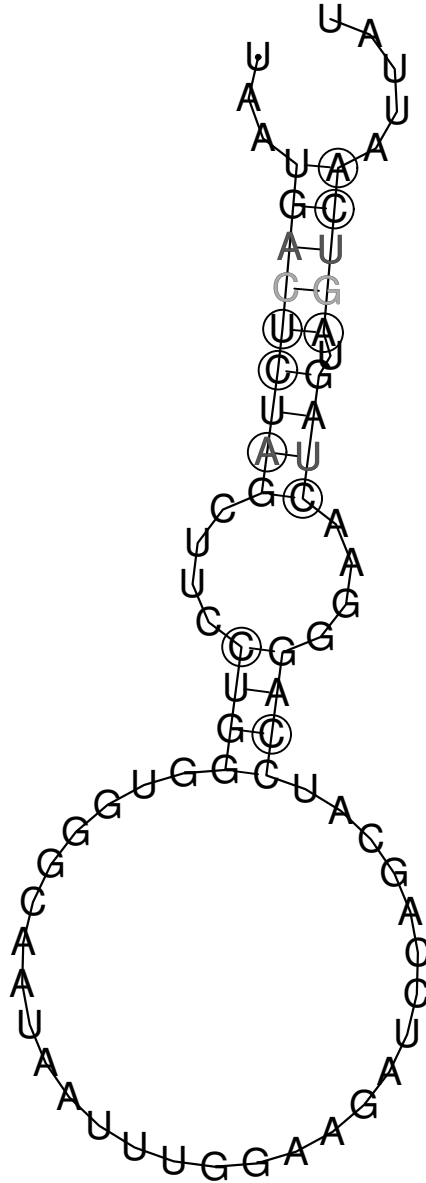


Predicted secondary structure of  $\epsilon$ ,  $\alpha$ , and  $\beta_1$  element. Circles indicate compensatory mutations.



...((((((((.....((((.....)))))).....)))))).....

A secondary structure element with unknown function that appears to be conserved in the C - mRNA but does not appear as conserved structure in complete pregenomic RNA.



A secondary structure element with unknown function that appears to be conserved in the C - mRNA but does not appear as conserved structure in complete pregenomic RNA.