

Alignement multiple

Équipe Bonsai

<http://www.lifl.fr/bonsai>

année 2012



Définition de l'alignement multiple

Définition de l'alignement multiple

- entrée : k séquences

```
C A T G C G A G T A G T A G
C A T G G T A G T A G
C C T G G A G T A C G T A G
C A T G A G C G T A G
```

Définition de l'alignement multiple

- entrée : k séquences

```
C A T G C G A G T A G T A G
C A T G G T A G T A G
C C T G G A G T A C G T A G
C A T G A G C G T A G
```

- sortie : un tableau contenant les k séquences, avec des indels

```
C A T G C G A G T A - G T A G
C A T G - - - G T A - G T A G
C C T G - G A G T A C G T A G
C A T G - - A G - - C G T A G
```

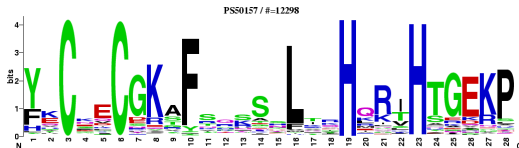
Motif doigt de zinc (*C2H2*-type)

```
TTY1_HUMAN YVCPFDGCKKFAQSTNLKSHILT--H
YKQ8_CAEEL YKCT--VCRKDISSESRLTHMFKQHH
BASO_HUMAN FQCD--ICKKTFKNACSVKIHHKN-MH
ZG2-9_XENL FVCT--VCGKTYKYKHGLNTHLHS--H
P43_XENBO LKCSVPGCKRSFRKKRALRIHVSE--H
IKAR_MOUSE FECN--MCGYHSQDRYEFSSHITRGEH
TRA1_CAEEL YKCEFADCEKAFSNASDRAKHQNR--TH
ZN10_HUMAN YKCN--QCGIIFSQNSPFIHVHQA--H
XFIN_XENLA FRCS--ECSRSFTHNSDLTAHMRK--H
TF3A_BUFAM CKCETENCNLAFTTASNMRHLHFKR-AH
ZG58_XENLA FVCT--ECNLSFAGLANLRSHQHL--H
P43_XENBO YRCSYEDCQTVSPTWTALQTHLKK--H
TSH_DROME FRCV--WCKQSFPITLEALTTHMKDSKH
ZN76_HUMAN FRCGYKGCGRLYTTAHHLKVHERA--H
TF3A_BUFAM YRCPRENCDRITYTKFNLKSHILT-FH
SUHW_DROAN YACK--ICGKDFTRSYPHLKRHHQYSSC
ZN76_HUMAN YTCPEPHCGRGFTSATNYKNHVRI--H
SRYC_DROME FKCN--YCPRDFTNFPNWLKHTRR--RH
EVI1_HUMAN YRCK--YCDRSFSISSNLQRHVRN-IH
```

Motif doigt de zinc (C2H2-type)

```

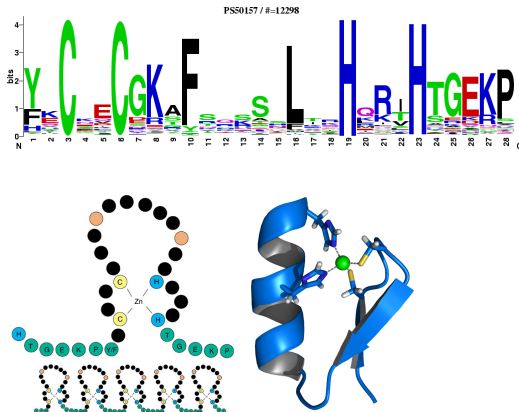
TTY1_HUMAN  YVCPFDGCNKKFAQSTNLKSHILT--H
YKQ8_CAEEL  YKCT--VCRKDISSESLRTHMFKQHH
BASO_HUMAN  FQCD--ICKTTFKNACSVKIHHKN-MH
ZG2-9_XENL  FVCT--VCGKTYKYKHGLNTHLHS--H
P43_XENBO   LKCSVPGCKRSFRKKRALRIHVSE--H
IKAR_MOUSE  FECN--MCGYHSQDRYEFSSHITRGEH
TRA1_CAEEL  YKCEFADCEKAFSNASDRAKHQNR-TH
ZN10_HUMAN  YKCN--QCGIIFSQNSPFIVHQIA--H
XFIN_XENLA  FRCS--ECSRSFTHNSDLTAHMRK--H
TF3A_BUFAM  CKCETENCNLAFTTASNMRLHFKR-AH
ZG58_XENLA  FVCT--ECNLSFAGLANLRSHQHL--H
P43_XENBO   YRCSYEDCQTVSPITWALQTHLKK--H
TSH_DROME   FRCV--WCKQSFPTLEALTTHMKDSKH
ZN76_HUMAN  FRCGYKGCGRLYTTAHHLKVHERA--H
TF3A_BUFAM  YRCPRENCDRYTTKFNLKSHILT-FH
SUHW_DROAN  YACK--ICGKDFTRSYHLKRHQKYSSC
ZN76_HUMAN  YTCPEPHCGRGFTSATNYKNVRI--H
SRYC_DROME  FKCN--YCPRDFTNFPNWLKHTRR-RH
EVI1_HUMAN  YRCK--YCDRSFSISSNLQRHVRN-IH
    
```



Motif doigt de zinc (C2H2-type)

```

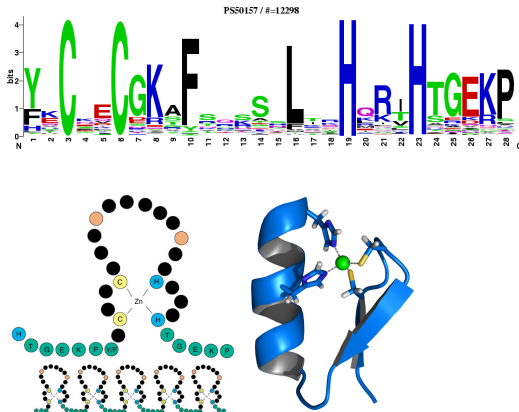
TTY1_HUMAN  YVCPFDGCKKKFAQSTNLKSHILT--H
YKQ8_CAEEL  YKCT--VCRKDISSESRLTHMFKQH
BASO_HUMAN  FQCD--ICKKTFKNACSVKIHHKN-MH
ZG2-9_XENL  FVCT--VCGKTYKYKHGLNTHLHS--H
P43_XENBO   LKCSVPGCKRSFRKKRALRIHVSE--H
IKAR_MOUSE  FECN--MCGYHSQDRYEFSSHITRGEH
TRA1_CAEEL  YKCEFADCEKAFSNASDRAKHQNR-TH
ZN10_HUMAN  YKCN--QCGIIFSQNSPFIHQIA--H
XFIN_XENLA  FRCS--ECSRSTHNSDLTAHMRK--H
TF3A_BUFAM  CKCETENCNLAFTTASNMRLEHFKR-AH
ZG58_XENLA  FVCT--ECNLSFAGLANLRSHQHL--H
P43_XENBO   YRCSYEDCQTVSPITWALQTHLKK--H
TSH_DROME   FRCV--WCKQSFTLEALTTHMKDSKH
ZN76_HUMAN  FRCGYKGCGRLYTTAHLKLVHERA--H
TF3A_BUFAM  YRCPRENCDRTYTTKFNLKSHILT-FH
SUHW_DROAN  YACK--ICGKDFTRSYPHLKRHHQYSSC
ZN76_HUMAN  YTCPEPHCGRGFTSATNYKNHVRI--H
SRYC_DROME  FKCEN--YCPRDFTNFPNWLKHTRR-RH
EVI1_HUMAN  YRCK--YCDRSFSISSNLQRHVRN-IH
    
```



Motif doigt de zinc (C2H2-type)

```

TYY1_HUMAN  YVCPFDGCKKKFAQSTNLKSHILT--H
YKQ8_CAEEL  YKCT--VCRKDISSSESLRTHMFKQHH
BASO_HUMAN  FQCD--ICKKTFKNACSVKIHKHN-MH
ZG2-9_XENL  FVCT--VCGKTYKYKHGLNTHLHS--H
P43_XENBO   LKCSVPGCKRSFRKKRALRIHVSE--H
IKAR_MOUSE  FECN--MCGYHSQDRYEFSSHITRGEH
TRA1_CAEEL  YKCEFADCEKAFSNASDRAKHQNR--TH
ZN10_HUMAN  YKCN--QCGIIFSQNSPFIHVQIA--H
XFIN_XENLA  FRCS--ECSRSTHNSDLTAHMRK--H
TF3A_BUFAM  CKCETENCNLAFTTASNMRLEHFKR-AH
ZG58_XENLA  FVCT--ECNLSFAGLANLRSHQHL--H
P43_XENBO   YRCSYEDCQTVSPITWTALQTHLKK--H
TSH_DROME   FRCV--WCKQSFPTEALTTMHKDSKH
ZN76_HUMAN  FRCGYKGCGRGLYTTAHLKLVHERA--H
TF3A_BUFAM  YRCPRENCDRITYTTKFNLKSHILT-FH
SUHW_DROAN  YACK--ICGKDFTRSYPHLKRHHQYSSC
ZN76_HUMAN  YTCPEPHCGRGFTSATNYKNHVRI--H
SRYC_DROME  FKCEN--YCPRDFTNFPNWLKHTRR-RH
EVI1_HUMAN  YRCK--YCDRSFSISSNLQRHVRN-IH
    
```



modélisation : motif **Prosite**

C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H

Site de fixation de la cellulose

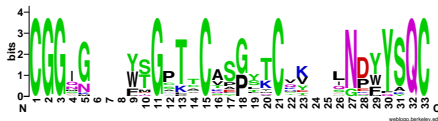
HWGQCGGI---GYSGCKTCTSGTTCCQYSNDYYSQCL
HYGQCGGI---GYSGPTVCASGTTCCVLNPPYYSQCL
QWGQCGGI---GYTGSTTCASPYTCHVLNPPYYSQCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCKVQNWYSQCL
IWGQCGGN---GWTGATTCASGLKCEKINDWYYQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDYYSQCL
DWAQCGGN---GWTGPTTCVSPYTCTKQNDWYSQCL
QWGQCGGQ---NYSGPTTCKSPFTCKKINDYYSQCL
RWQQCGGI---GFTGPTQCEEPYICTKLNDWYSQCL
HWAQCGGI---GFSGPTTCPEPYTCAKDHDYYSQCV
LYEQCGGI---GFDGVTCCEGLMCMKMPYYSQCR
VWAQCGGQ---NWSGTPCCTSGNKCVELNDYYSQCL
PYGQCGGM---NYSGMTMCSGPKCVELNEFFSQCD
AYYQCGGSKSAYPNGNLACATGSKCVKQNEYYSQCV
EYAACGGE---MFMGAKCKFGLVCYETSGKWSQCR

extrait de Prosite, entrée PS00562

Site de fixation de la cellulose

```
HWGQCGGI---GYSGCKTCTSGTTCCQYSNDYYSQCL
HYGQCGGI---GYSGPTVCASGTTCCQLNPYYSQCL
QWGQCGGI---GYTGSTTCASPYTCHVLNPYYSQCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCVKQNQWYSQCL
IWGQCGGN---GWTGATTCASGLCKEINDWYYQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDYYSQCL
DWAQCGGN---GWTGPTTCVSPYTCCKQNDWYSQCL
QWGQCGGQ---NYSGPTTCKSPFTCKKINDYYSQCC
RWQQCGGI---GFTGPTTCEEPYICTKLNDWYSQCL
HWAQCGGI---GFSGPTTCEPEYTC AKDHDYYSQCV
LYEQCGGI---GFDGVTCSEGLMCKMGPYYSQCR
VWAQCGGQ---NWSGTPCCTSGNKC VKLNDYYSQCC
PYGQCGGM---NYSGKTMCS PGFKVELNEFFSQCD
AYYQCGGSKSAYPNGNLACATGSKCVKQNEYYSQCV
EYAA CGGE---MFMGAKCKFGLVCYETSGKWSQCR
```

extrait de Prosite, entrée PS00562

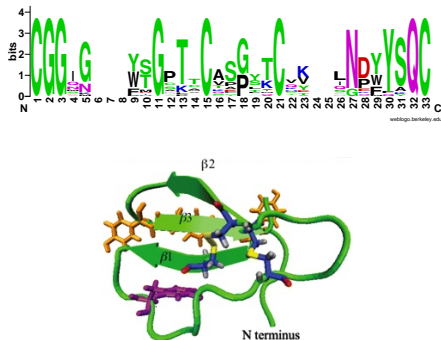


Site de fixation de la cellulose

```

HWGQCGGI---GYSGCKTCTSGTTCQYSNDYYSQCL
HYGQCGGI---GYSGPTVCASGTTCCVLNPYYSQCL
QWGQCGGI---GYTGSTTCASPYTCHVLNPYYSQCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCVKQNQWYSQCL
IWGQCGGN---GWTGATTCASGLCKEINDWYYQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDFYSQCL
DWAQCGGN---GWTGPTTCVSPYTCIKQNDWYSQCL
QWGQCGGQ---NYSGPTTCKSPFTCKKINDFYQQCQ
RWQQCGGI---GFTGPTTCEEPYICTKLNDWYSQCL
HWAQCGGI---GFSGPTTCEPYTCAKDHDYISQCV
LYEQCGGI---GFDGVTCSEGLMCMKMGPPYQQCR
VWAQCGGQ---NWSGTPCTSGNKCCKLNDFYQQCQ
PYGQCGGM---NYSGKTMCSPGFKCVELNEFFSQCD
AAYQCGGSKSAYPNGNLACATGSKCVKQNEYYSQCV
EYAA CGGE---MFMGAKCKFGFLVCYETSGKWSQCR
    
```

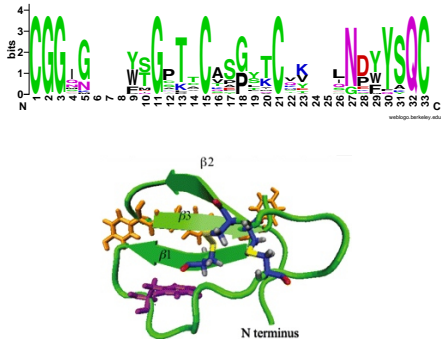
extrait de Prosite, entrée PS00562



Site de fixation de la cellulose

HWGQCGGI--GYSCKTCTSGTTCQYSNDYYSQCL
HYGQCGGI---GYSGPVTCASGTCQVLNPYYYSQCL
QWQGCGGI---GYTGSTCTCASPYTCHVLNPYYYSQCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCVKVQNWYYSQCL
IWGQCGGN---GWTGATTCASGLCKEKNIDWYYQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDFYYSQCL
DWAQCGGN---GWTGPTTCVSPYTC TKQNDWYYSQCL
QWQGCGGQ---NYSGPTTCKSPFTCKKINDFYYSQCL
RWQQCGGI---GFTGPTQC EEPYCTKLNIDWYYSQCL
HWAQCGGI---GFSGPTTCPEPYTCAKDHD IYSQCV
LYEQCGGI---GFDGVTCTSEGLCMCMGPYYYSQCR
VWAQCGGQ---NWSGPTCTSGNCKVLKNDFYYSQCL
PYGQCGGM---NYSKTM CSPGFCV LNEFFYSQCL
AYYQCGGSKSAYPNGLNCAATGSKCVQYNEEYYSQCV
EYAA CGGE---MFMGAKCKFGLVCYTSGKWSQCL

extrait de Prosite, entrée PS00562



C-G-G-x(4,7)-G-x(3)-C-x(5)-C-x(3,5)-[NHG]-x-[FYWM]-x(2)-Q-C

```

      +-----+
      |           +-----|-----+
      |           |           |
xxxxxxCxxxxxxxxxxCxxxxxCxxxxxxxxxCx
*****

```

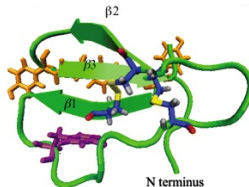
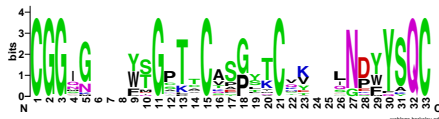
les 4 cystéines sont impliquées
dans des liaisons di-sulfures

Site de fixation de la cellulose

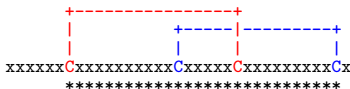
```

HWGQCGGI---GYSGCKTCTSGTTCQYSNDYYSQCL
HYGQCGGI---GYSGPTVCASGTTCCVLPNPPYYSQCL
QWGQCGGI---GYTGSTTCASPYTCHVLPNPPYYSQCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCVKQNDQWYSQCL
IWGQCGGN---GWTGATTCASGLCKEINDWYYSQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDQWYSQCL
DWAQCGGN---GWTGPTTCVSPYTCVKQNDQWYSQCL
QWGQCGGQ---NYSGPTTCCKSPFTCKKINDQWYSQCL
RWGQCGGI---GFTGPTTCCEPYTCTKLNQWYSQCL
HWAQCGGI---GFSGPTTCPEPYTCAKDNDIYSQCV
LYEQCGGI---GFDGVTCSEGLMCMKMGPPYYSQCR
VWAQCGGQ---NWSGTPCTSGNKCCKLNQWYSQCL
PYGQCGGM---NYSGKTMCSGPFKCVELNEFFSQCD
AAYQCGGSKSAYPNGLACATGSKCVKQNEYYYSQCV
EYAA CGGE---MFMGAKCKKFLVCYETSGKWSQCR
    
```

extrait de Prosite, entrée PS00562



C-G-G-x(4,7)-G-x(3)-C-x(5)-C-x(3,5)-[NHG]-x-[FYWM]-x(2)-Q-C



les 4 cystéines sont impliquées dans des liaisons di-sulfures

Structure d'ARN

Structure d'ARN

- on dispose d'une famille d'ARN possédant la même structure

G	A	G	C	C	C	A	G	U	U	C
	A	G	G	A	C	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U

Structure d'ARN

- on dispose d'une famille d'ARN possédant la même structure
- pour un appariement de la structure donné :
 - si une base mute dans la structure d'ARN, la base qui s'y apparie doit muter aussi . . .

G	A	G	C	C	C	A	G	U	U	C
	A	G	G	A	C	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U

Structure d'ARN - méthode comparative

1 construction de l'alignement multiple

G	A	G	C	-	C	C	A	G	U	U	C
-	A	G	G	A	C	-	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U	-

2 détection de positions corrélées

Structure d'ARN - méthode comparative

1 construction de l'alignement multiple

G	A	G	C	-	C	C	A	G	U	U	C
-	A	G	G	A	C	-	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U	-

2 détection de positions corrélées

Structure d'ARN - méthode comparative (ARNt)

```
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAAGGCAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAAGGCAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCCTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAAGCAGGTGT-CGTCGGTTTCGAATCCGTCTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
```

Structure d'ARN - méthode comparative (ARNt)

```
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAAGGCAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAAGGCAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCCTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAAGCAGGTGT-CGTCGGTTTCGAATCCGTCTGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
```

Structure d'ARN - méthode comparative (ARNt)

```
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAAGGAGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAAGGCAAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAAGGCAAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAAGGAGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAAGGAGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAAGGAGTTCAGGAGTTTCGATCCTCCTTGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCCTTGCACGCAAGGAGTCAACGGTTTCGATCCCGTTTTGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAAGGAGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAAGCAAGTGT-CGTCGGTTTCGAATCCGTCTGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAAGCAAGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAAGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAAGCAAGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
```

Structure d'ARN - méthode comparative (ARNt)

GGGG^AAATTAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGGTTCAGCGTTTCGATCCCGCTAT^TCTCCA---
GGGG^{CT}ATAGCTCAGCT-GGGAGAGCGC^{CT}GCTTTGCACGC^{AG}GAGGTCTGCGTTTCGATCCCGCATAG^{CT}CCACCA
GGGG^{CT}ATAGCTCAGCT-GGGAGAGCGC^{TT}GCATGGCATGC^{AA}GAGGTTCAGCGTTTCGATCCCGCTTAG^{CT}CCACCA
GGGG^AAATTAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGGTTCAGCGTTTCGATCCCGCTAT^TCTCCA---
GGGG^CCTTAGCTCAGTC-GGTAGAGCAC^{TG}CCTTTGCAAGG^{CA}GATGTCAGGGGTTTCGATTCCCGCTAG^GCTCCA---
GGGG^GTATAGCTCAGTT-GGTAGAGCGC^{TG}CCTTTGCAAGG^{CA}GAGTTCAGCGTTTCGATTCCCGCTTA^CCCCCA---
GGGG^{CT}ATAGCTCAGCT-GGGAGAGCGC^{CT}GCTTTGCACGC^{AG}GAGGTCTGCGTTTCGATCCCGCATAG^{CT}CCACCA
GGGG^{CT}ATAGCTCAGCT-GGGAGAGCGC^{CT}GCTTTGCACGC^{AG}GAGGTCTGCGTTTCGATCCCGCATAG^{CT}CCACCA
GGGG^GCATAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGT-CGTCGGTTTCGATCCCGTCTG^CCTCCACCA
GGGG^CCATAGCTCAGCT-GGGAGAGCGC^{CT}GCTTTGCACGC^{AG}GAGTTCAGGAGTTCGATCCTCCTTG^GCTCCACCA
GGGG^GCATAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGT-CGTCGGTTTCGATCCCGTCTG^CCTCCACCA
GGGG^CCATAGCTCAGCTGGGGAGAGCGC^{CT}GCCTTGCACGC^{AG}GAGGTCAACGGTTTCGATCCCGTTTG^GCTCCA---
GGGG^GCATAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGT-CGTCGGTTTCGATCCCGTCTG^CCTCCACCA
GGGG^CATTAGCTCAGCT-GGGAGAGCGC^{CT}GCTTTGCACGC^{AG}GAGGTTCAGCGTTTCGATCCCGCTAT^TCTCCACCA
GGGG^CCATAGCTCAGTT-GGTAGAGCGC^{CT}GCTTTGCAAGC^{AG}GTGT-CGTCGGTTTCGAATCCGTCTG^GCTCCACCA
GGGG^CCGTAGCTCAGCTGGG-AGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGGTTCGGAGGTTTCGATCCCGTCCG^GCTCCACCA
GGGG^CCGTAGCTCAGCT-GGGAGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGGTTCGTCGGTTTCGATCCCGTCCG^GCTCCACCA
GGGG^CCGTAGCTCAGCT-GG-AGAGCAC^{CT}GCTTTGCAAGC^{AG}GGGGTTCGTCGGTTTCGATCCCGTCCG^GCTCCACCA

Structure d'ARN - méthode comparative (ARNt)

GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTG CATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCTTTGCAAGGCAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCTTTGCAAGGCAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTTGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAAGCAGGTGT-CGTCGGTTTCGAATCCGTCTGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAAGCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA

Structure d'ARN - méthode comparative (ARNt)

```
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAA GGCAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAA GGCAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTCGATCCTCCTTGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTTGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTCTGCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTTGCAA GCAGGTGT-CGTCGGTTTCGAATCCGTCTGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTTGCAA GCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
```


Structure d'ARN - méthode comparative (ARNt)

```
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAA GGCAGATGTCAGGGGTTTCGATTCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAA GGCAGAAGTCAGCGGTTTCGATTCCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCCATAGCTCAGCTGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCTCCACCA
GGGGCAATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTTGCAA GCAGGTGT-CGTCGGTTTCGAATCCGTCTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTTGCAA GCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
```

Structure d'ARN - méthode comparative (ARNt)

GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGTTTCGATCCCGCTTAGCTCCACCA
GGGGAAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAA GGCAGATGTCAGGGGTTTCGATCCCGCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAA GGCAGAAGTCAGCGTTTCGATCCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCTGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAA GCAGGTGT-CGTCGGTTTCGAATCCGTCTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA

Structure d'ARN - méthode comparative (ARNt)

GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGTTTCGATCCCGCTTAGCTCCACCA
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAA GGCAGATGTCAGGGGTTTCGATTCCCCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAA GGCAGAAGTCAGCGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAA GCAGGTGT-CGTCGGTTTCGAATCCGCTTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA

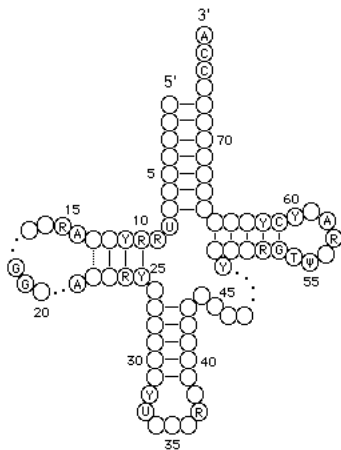
Structure d'ARN - méthode comparative (ARNt)

GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGGTTTCGATCCCGCTTAGCTCCACCA
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCAGCGGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCAC TGCC TTTGCAA GGCAGATGTCAGGGGTTTCGATTCCCCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCC TTTGCAA GGCAGAAGTCAGCGGTTTCGATTCCGCTTACCCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCTTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGT-CGTCGGTTTCGATCCCGCTTGCCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAA GCAGGTGT-CGTCGGTTTCGAATCCGCTTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGGAGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAA GCAGGGGGTCGTCGGTTTCGATCCCGTCCGGCTCCACCA

Structure d'ARN - méthode comparative (ARNt)

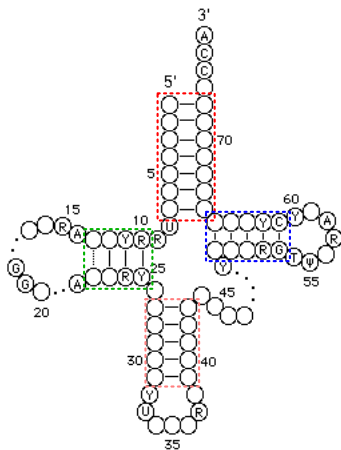
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCTTGCATGGCATGCAAGAGGTTCAGCGTTTCGATCCCGCTTAGCTCCACCA
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTTCAGCGTTTCGATCCCGCTATTCTCCA---
GGGGCCTTAGCTCAGTC-GGTAGAGCACTGCCTTTGCAAGGCAGATGTTCAGGGGTTTCGATTCCCCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCTTTGCAAGGCAGAAGTTCAGCGTTTCGATTCCGCTTACCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTCTGCGGTTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGGAGTTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCCTTGCACGCAGGAGGTCAACGGTTTCGATCCCGTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGT-CGTCGGTTTCGATCCCGTCTGCCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTGCACGCAGGAGGTTCAGCGTTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTGCAAGCAGGTGT-CGTCGGTTTCGAATCCGTCTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTGCAAGCAGGGGGTTCGGAGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTGCAAGCAGGGGGTTCGTCGGTTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTGCAAGCAGGGGGTTCGTCGGTTTCGATCCCGTCCGGCTCCACCA

Structure d'ARN - méthode comparative (ARNt)



Structure secondaire de l'ARNt

Structure d'ARN - méthode comparative (ARNt)



Structure secondaire de l'ARNt

Alignement 2 à 2

Deux séquences quelconques



Détection d'une similarité **syntactique**



Y a-t-il une **fonction** commune ?

Alignement multiple

Famille de séquences avec la même **fonction**



À quelle conservation **syntactique** cela correspond-il ?

Score d'un alignement multiple

- doit rendre compte de la qualité de l'alignement multiple
- habituellement les colonnes sont considérées **indépendantes**

⇒ la somme des scores associés à chaque colonne

somme des paires (Sum of Pairs)

$$SP(m_i) = \sum_{1 \leq j < k \leq n} s(m_i^j, m_i^k)$$

m_i = la i -ème colonne de l'alignement

m_i^j = j -ème aa dans la colonne i

Exemple

jeu de scores :

$$s(x,x)=1, \quad s(x,y)=-1, \quad s(x,-)=s(-,x)=-2, \quad s(-,-)=0$$

	A	A	C	G	T	A	C	G	A	T	A	
	A	-	C	G	T	A	-	A	A	T	G	
	G	T	C	G	T	A	-	-	T	T	A	

(1-2)	1	-2	1	1	1	1	-2	-1	1	1	-1	
(1-3)	-1	-1	1	1	1	1	-2	-1	-1	1	1	
(2-3)	-1	-2	1	1	1	1	0	-2	-1	1	-1	
	=	=	=	=	=	=	=	=	=	=	=	
	-1	-5	3	3	3	3	-4	-4	-1	3	-1	= -1

Définition alternative (équivalente)

- α : alignement multiple pour les séquences s_1, \dots, s_n
- α_{ij} : projection de l'alignement pour s_i et s_j

$$SP(\alpha) = \sum_{1 \leq i < j \leq n} score(\alpha_{ij})$$

Exemple

jeu de scores :

$$s(x,x)=1, \quad s(x,y)=-1, \quad s(x,-)=s(-,x)=-2, \quad s(-,-)=0$$

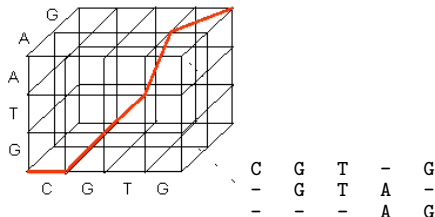
A	A	C	G	T	A	C	G	A	T	A
A	-	C	G	T	A	-	A	A	T	G
G	T	C	G	T	A	-	-	T	T	A

(1-2)	1	-2	1	1	1	1	-2	-1	1	1	-1	=	1
(1-3)	-1	-1	1	1	1	1	-2	-1	-1	1	1	=	0
(2-3)	-1	-2	1	1	1	1	0	-2	-1	1	-1	=	-2
												=	-1

Algorithme exact

⇒ par programmation dynamique

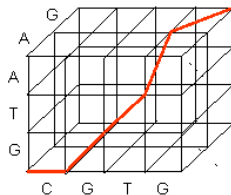
- alignement deux à deux ⇒ chemin dans une matrice de dimension 2
- alignement multiple de n séquences ⇒ chemin dans une matrice de dimension n



Algorithme exact

⇒ par programmation dynamique

- alignement deux à deux ⇒ chemin dans une matrice de dimension 2
- alignement multiple de n séquences ⇒ chemin dans une matrice de dimension n



C	G	T	-	G
-	G	T	A	-
-	-	-	A	G

⇒ impossible de l'utiliser en pratique.

≈ 100 ans pour 8 séquences de 100bp.

Definition (Heuristique)

Algorithme utilisant des règles simples pour diminuer l'espace de recherche des solutions (mais ne donnant pas forcément la meilleure solution)

- Clustal (*le plus populaire : Clustal, ClustalW, Clustal-Omega*)
- Dialign2 (*complémentaire à Clustal*)
- T-coffee, Muscle, Pima, Multalin, MA-FFT...

⇒ autant de programmes qui produisent des alignements différents !

Alignement basé sur un arbre guide

- idée : reconstruire l'alignement multiple à partir d'un **arbre guide** (clusters)
 - feuilles = séquences
 - noeuds = alignements
- partir des feuilles puis remonter dans l'arbre
 - utilisation de la technique de **profile alignment** \Rightarrow produire un seul alignement multiple avec deux (par prog. dyn.)

MultAlin

CLUSTer + **AL**ignement \Rightarrow CLUSTAL

F. Corpet, 1988

principe :

- 1 calcule une matrice de similarité des paires
- 2 construit un arbre de clustering hiérarchique (UPGMA)
- 3 construit l'alignement multiple en suivant l'arbre
- 4 reconstruit un arbre de clustering hiérarchique avec les nouveaux alignements paire à paire issus de l'alignement trouvé
- 5 réitère le processus jusqu'à stabilisation de l'arbre de clustering

- agglomère les 2 séquences de score maximal (\sim distance minimale dans UPGMA)
- calcule les nouveaux scores entre ce cluster et les autres en faisant la moyenne

$$s(C_1, C_2) = \frac{1}{\text{Card}(C_1) \times \text{Card}(C_2)} \sum_{c \in C_1, c' \in C_2} s(c, c')$$

MultAlin - exemple 1

① TACCATGA

② TACCATA

③ GACGACCA

④ GACCATCTCA

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④



MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④



① TACCATGA
② TACCAT-A

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 1

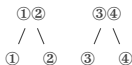
① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

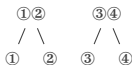
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	2	①②	.	0	2.5
②	.	.	0	3	③	.	.	4
③	.	.	.	4	④	.	.	.
④				



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

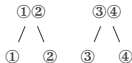
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④		①②	③④
①	.	6	0	2	①②	.	0	2.5	①②	.	1.25
②	.	.	0	3	③	.	.	4	③④	.	.
③	.	.	.	4	④	.	.	.			
④							



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

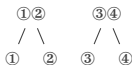
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

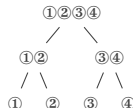
	①	②	③	④		①②	③	④		① ②	③④
①	.	6	0	2	①②	.	0	2.5	①②	.	1.25
②	.	.	0	3	③	.	.	4	③④	.	.
③	.	.	.	4	④	.	.	.			
④							



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA



MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

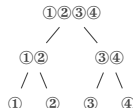
	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA



① TACCAT--GA
② TACCAT--A
③ GACGA-C-CA
④ GACCATCTCA

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

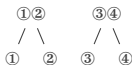
	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

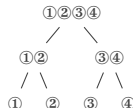
	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA



① TACCAT--GA
② TACCAT---A
③ GACGA-C-CA
④ GACCATCTCA

3 nouvelle matrice des scores et on recommence :

MultAlin - exemple 1

① TACCATGA ② TACCATA ③ GACGACCA ④ GACCATCTCA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

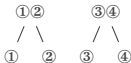
	①	②	③	④
①	.	6	0	2
②	.	.	0	3
③	.	.	.	4
④

	①②	③	④
①②	.	0	2.5
③	.	.	4
④	.	.	.

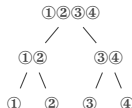
	① ②	③④
①②	.	1.25
③④	.	.



① TACCATGA
② TACCAT-A



③ GACGA-C-CA
④ GACCATCTCA



① TACCAT--GA
② TACCAT---A
③ GACGA-C-CA
④ GACCATCTCA

3 nouvelle matrice des scores et on recommence :

① TACCAT--GA
② TACCAT---A
③ GACGAC--CA
④ GACCATCTCA

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④



MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④



① TACCATGA
② TACCAT-A

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	5	①②	.	0	4
②	.	.	0	3	③	.	.	3
③	.	.	.	3	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	5	①②	.	0	4
②	.	.	0	3	③	.	.	3
③	.	.	.	3	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	5	①②	.	0	4
②	.	.	0	3	③	.	.	3
③	.	.	.	3	④	.	.	.
④				



① TACCATGA
② TACCAT-A

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

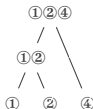
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	5	①②	.	0	4
②	.	.	0	3	③	.	.	3
③	.	.	.	3	④	.	.	.
④				



① TACCATGA
② TACCAT-A



MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

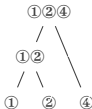
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④
①	.	6	0	5	①②	.	0	4
②	.	.	0	3	③	.	.	3
③	.	.	.	3	④	.	.	.
④				



① TACCATGA
② TACCAT-A



① TACCAT-GA
② TACCAT--A
④ TACGATCGA

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

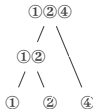
scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④		①②	③	④		①②④	③
①	.	6	0	5	①②	.	0	4	①②④	.	1
②	.	.	0	3	③	.	.	3	③	.	.
③	.	.	.	3	④	.	.	.			
④							



① TACCATGA
② TACCAT-A



① TACCAT-GA
② TACCAT--A
④ TACGATCGA

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

2 construction d'un arbre de clustering (et de l'alignement) :

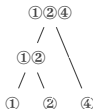
	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④

	①②	③	④
①②	.	0	4
③	.	.	3
④	.	.	.

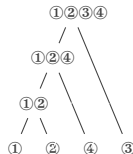
	①②④	③
①②④	.	1
③	.	.



① TACCATGA
② TACCAT-A



① TACCAT-GA
② TACCAT--A
④ TACGATCGA



MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

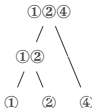
2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④



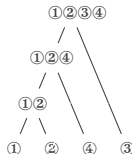
① TACCATGA
② TACCAT-A

	①②	③	④
①②	.	0	4
③	.	.	3
④	.	.	.



① TACCAT-GA
② TACCAT--A
④ TACGATCGA

	①②④	③
①②④	.	1
③	.	.



① TACCAT-GA
② TACCAT--A
④ TACGATCGA
③ GACGACC-A

MultAlin - exemple 2

① TACCATGA ② TACCATA ③ GACGACCA ④ TACGATCGA

1 calcul des meilleurs alignements 2 à 2 :

scores (Match = 1, Mismatch = -1, Indel = -1)

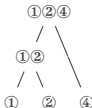
2 construction d'un arbre de clustering (et de l'alignement) :

	①	②	③	④
①	.	6	0	5
②	.	.	0	3
③	.	.	.	3
④



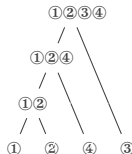
① TACCATGA
② TACCAT-A

	①②	③	④
①②	.	0	4
③	.	.	3
④	.	.	.



① TACCAT-GA
② TACCAT--A
④ TACGATCGA

	①②④	③
①②④	.	1
③	.	.



① TACCAT-GA
② TACCAT--A
④ TACGATCGA
③ GACGACC-A

Thompson et al., 1994

principe :

- 1 calcule une matrice de similarité des paires par prog. dyn.
- 2 convertit les similarités en distances
- 3 construit l'arbre guide (méthode du Neighbor-Joining)
- 4 aligne progressivement les noeuds de l'arbre par ordre décroissant de similarité

■ 4 séquences

s_1	cgatgagtcattgtgactg
s_2	cgagccattgtagctactg
s_3	cgaccattgtagctacctg
s_4	cgatgagtcactgtgactg

■ jeu de scores :

indel : -2, substitution : -1, identité : 1

Etape 1

calcul des scores de similarité de tous les alignements

s1 cgatgagtcattgt-g--actg
||| | ||||| | |||
s2 cga-g--ccattgtagctactg

s1 cgatgagtcattgt-tgactg
||| | | | | |||
s3 cgacca-ttgtagctacctg

s1 cgatgagtcattgtgactg
||||||| |||||
s4 cgatgagtcactgtgactg

s2 cgagccattgtagctac-tg
||| ||||| ||||| ||
s3 cga-ccattgtagctacctg

s2 cga-g--ccattgtagctactg
||| | || ||| | |||
s4 cgatgagtcactgt-g--actg

s3 cgaccattgtagctacctg
||| | | | |||
s4 cgatgagtcactgtgactg

tableau des scores d'alignement :

	s ₁	s ₂	s ₃	s ₄
s ₁		2	0	17
s ₂	2		14	0
s ₃	0	14		-1
s ₄	17	0	-1	

n séquences

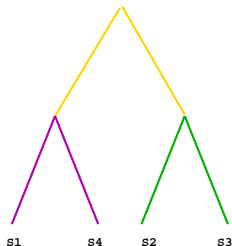
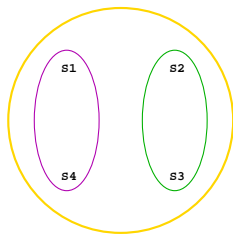
↓
 $\frac{n(n-1)}{2}$ calculs

Etape 2

construction de l'arbre guide

arbre obtenu avec l'algorithme de Neighbor-Joining

dendrogramme



regroupement des séquences suivant leur similarité à partir de la matrice des scores 2 à 2.

Etape 3

construction de l'alignement multiple final

cgatgagtc^{s1}cattgtgactg

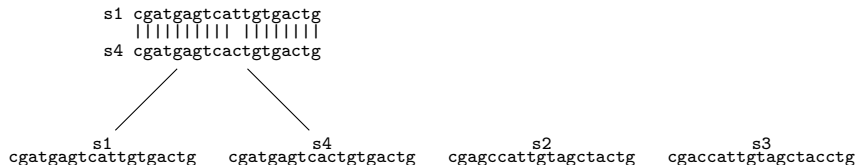
cgatgagtc^{s4}cactgtgactg

cgagccattg^{s2}tagctactg

cgaccattg^{s3}tagctacctg

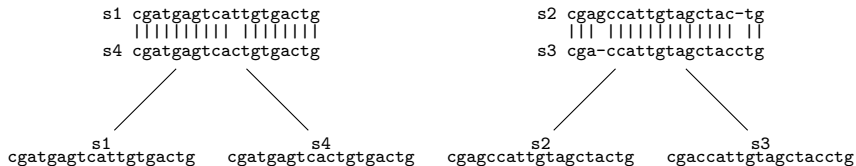
Etape 3

construction de l'alignement multiple final



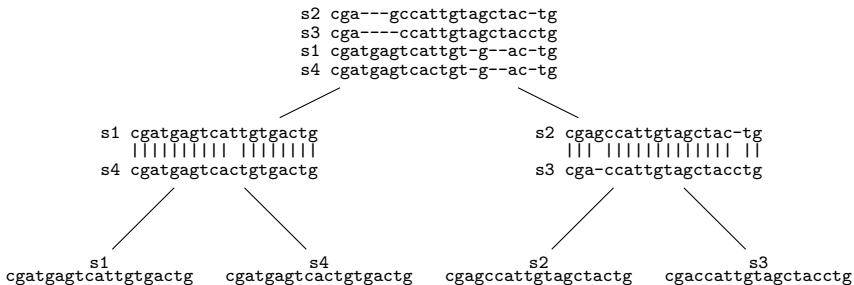
Etape 3

construction de l'alignement multiple final



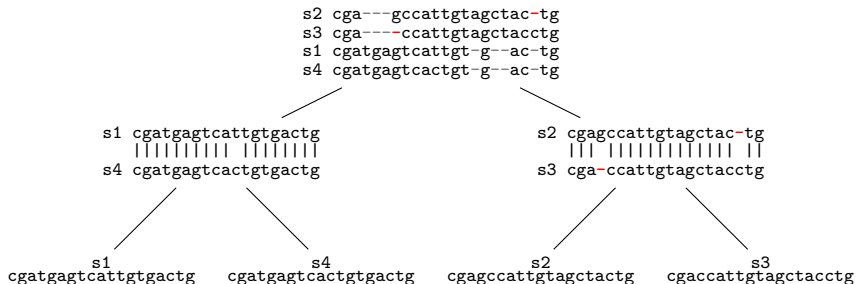
Etape 3

construction de l'alignement multiple final



construction de l'alignement multiple final

"Once a gap, always a gap."

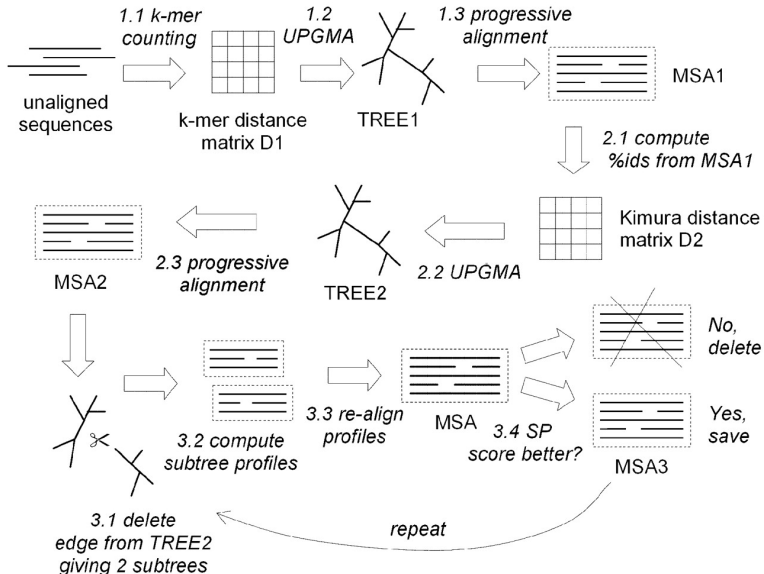


ClustalW est optimisé pour les protéines

- Pondération des séquences en fonction de leur sur- ou sous-représentation
- Adaptation des matrices de similarité au fil de l'algorithme en fonction de la divergence des séquences à aligner
BLOSUM 80 pour aligner les séquences proches,
BLOSUM 50 pour aligner des séquences distantes, par exemple.
- Pénalités de gaps spécifiques à chaque résidu.
Par exemple, les glycines sont davantage susceptibles d'avoisiner un gap que les valines.
- Pénalités de gaps réduites dans les régions hydrophiles
Encourage la formation de gaps dans des boucles plutôt que dans des régions structurées.
- Pénalités de gaps augmentées dans le voisinage d'autres gaps
Evite la formation de petits gaps voisins, au profit de longs gaps

- *Slow/Fast* : qualité des alignements 2 à 2
- Matrice de similarité (PAM, BLOSUM, Gonnet)
- Pénalités de gaps :
 - Ouverture et d'extension de gaps
 - Distance de voisinage entre deux gaps
 - Gaps hydrophiles
 - Ouvertures de gaps spécifiques

Arbre guide : autres méthodes - Muscle



FFT : Fast Fourier Transform

- *Progressive alignment* :

- 1/2 arbre UPGMA à l'aide d'une distance rapide (k -mots).

- 2/2 construction guidée d'un 1er alignement multiple (FFT-NS1).

- *Iterative refinement* :

- 1 réutiliser la matrice de distance de FFT-NS1 pour refaire un alignement multiple (FFT-NS2).

- 2 heuristique de réaligement par groupes (FFT-NSi).

alignement multiple de n séquences

alignement multiple de n séquences

construction de l'arbre guide $\implies \frac{n(n-1)}{2}$ “comparaisons” (alignements, distances par k -mots)

alignement multiple de n séquences

construction de l'arbre guide $\implies \frac{n(n-1)}{2}$ “comparaisons” (alignements, distances par k -mots)

Construction très longue si $n > 1000$ séquences

(l'alignement est souvent plus rapide que la construction de l'arbre guide).

alignement multiple de n séquences

construction de l'arbre guide $\implies \frac{n(n-1)}{2}$ "comparaisons" (alignements, distances par k -mots)

Construction très longue si $n > 1000$ séquences

(l'alignement est souvent plus rapide que la construction de l'arbre guide).

- 1 Éviter ces $\mathcal{O}(n^2)$ comparaisons $\rightarrow \mathcal{O}(n \log(n))$
- 2 Utilisation d'un alignement simple de *profils-HMM* (voir *prochain cours*) plutôt que de *blocks d'alignements*.

À partir des alignements locaux

- idée : repérer des similarités locales fortes entre les séquences
- typiquement : les diagonales du dotplot
- incorporer les diagonales dans l'alignement multiple
- conséquence : les gaps inter-diagonales sont considérés moins importants

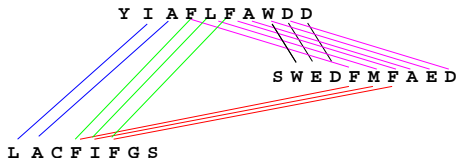
Diagonal + ALIGNment \Rightarrow DIALIGN

Morgenstern et al., 1996

principe :

- 1 alignement des paires avec optimisation des poids des diagonales
- 2 tri des diagonales selon leur poids et leur chevauchement
- 3 reconstruction gloutonne
 - 1 insertion des diagonales par poids décroissants
 - 2 vérification de la consistance avec les diagonales déjà introduites
- 4 recommencer

- Étape 1 : détection des diagonales dans les paires de séquences



- Étape 2 : sélection d'un ensemble cohérent de diagonales pour construire l'alignement
 - pas de croisements
 - pas de chevauchements

■ score maximal

y	I	A	-	F	L	F	A	W	D	d
-	L	A	c	F	I	F	g	s	-	-
s	w	e	d	F	M	F	A	E	D	-

CLUSTAL vs. DIALIGN

Exemple (C. Notre-Dame)

```
GARFIELD THE LAST FAT CAT
GARFIELD THE FAT CAT
GARFIELD THE VERY FAST CAT
THE FAT CAT
```

CLUSTAL vs. DIALIGN

Exemple (C. Notre-Dame)

GARFIELD THE LAST FAT CAT
GARFIELD THE FAT CAT
GARFIELD THE VERY FAST CAT
THE FAT CAT

Alignement fourni par Clustal

seq1	GARFIELDTHELASTFA-TCAT
seq2	----GARFIELDTHEFA-TCAT
seq3	GARFIELDTHEVERYFASTCAT
seq4	-----THEFA-TCAT

CLUSTAL vs. DIALIGN

Exemple (C. Notre-Dame)

```
GARFIELD THE LAST FAT CAT
GARFIELD THE FAT CAT
GARFIELD THE VERY FAST CAT
THE FAT CAT
```

Alignement fourni par Clustal

```
seq1      GARFIELDTHELASTFA-TCAT
seq2      ----GARFIELDTHEFA-TCAT
seq3      GARFIELDTHEVERYFASTCAT
seq4      -----THEFA-TCAT
```

Alignement fourni par Dialign2

```
seq1      GARFIELD THE LAST FA-T CAT
seq2      GARFIELD THE ---- FA-T CAT
seq3      GARFIELD THE VERY FAST CAT
seq4      ----- THE ---- FA-T CAT
```


Quelle méthode utiliser ? (1/2)

⇒ cela dépend du type de séquences à aligner ...

BaliBASE : base de données d'alignements multiples pour *benchmark*

- plus de 150 familles de protéines
- alignements basés sur la structure secondaire
 - Référence 1 séquences équidistantes avec différents niveaux de conservation
 - Référence 2 protéines homologues + 1 séquence orpheline
 - Référence 3 sous-groupes avec moins de 25% d'identité entre les groupes
 - Référence 4 extensions N/C-terminales
 - Référence 5 insertions internes
 - Référence 6 répétitions internes
 - Référence 7 protéines transmembranaires
 - Référence 8 permutations de domaines
- Réf. 1, 2 et 3 : préférer **Clustal** à **Dialign2**
- Réf. 4 et 5 : préférer **Dialign2** à **Clustal**

Quelle méthode utiliser ? (2/2)

- plus les séquences sont divergentes, moins le résultat est fiable
 - quand le taux d'identité est supérieur à 35%, toutes les méthodes sont satisfaisantes
alignements corrects à plus de 90%
 - **twilight zone** : 10-20 % identité
Aucune méthode n'assure un alignement avec plus de 50% de correction
- **Clustal** a tendance à autoriser moins de gaps que **Dialign2**
- similarité locale : **Dialign2**
- similarité globale : **Clustal**

Pas de méthode universelle
Pas de confiance aveugle vis-à-vis
du résultat obtenu

Exemple : domaine SH3

SH3 (Src homology 3) domains are often indicative of a protein involved in signal transduction related to cytoskeletal organization. The SH3 domain has a characteristic fold which consists of five or six beta- strands arranged as two tightly packed anti-parallel beta sheets. The linker regions may contain short helices.

Prosite PS50002

Séquences à aligner		longueur
=====		=====
1aboA	P00520	57
1ycsB	P04637	60
1pht	P27986	80
1ihvA	P00383	49
1vie	P12497	51

- séquences courtes
- similarité faible (< 25%) et diffuse

SH3 - Véritable Alignement

basé sur l'alignement des éléments de structure secondaire

```
1aboA  -NLFVALYDfvasgdntlsitkGEKLRVLgynhn-----
1ycsB  kGVIYALWDyepqnddelpmkeGDCMTIIhrede-----
1pht   gYQYRALYDykkereedidlhlGDILTVNkgs lvalgfsd
1ihvA  -NFRVYYRDSrd-----pvwkGPAKLLWkg-----
1vie   -drvrkksga-----awqQIVGWYctnlt-----

1aboA  -----gEWCEAQt--kngqGWVPSNYITPVN-----
1ycsB  -----deiEWWARl--ndkeGYVPRNLLGLYP-----
1pht   gqearpeeiGWLNGYnettgerGDFPGTYVEYIGrkkip
1ihvA  -----eGAVVIQd--nsdiKVVPRRKAKIIRd-----
1vie   -----peGYAVESeahpgsvQIYPVAALERIN-----
```

SH3 - Alignement fourni par Clustal

```
1aboA  -NLFV-ALYDFVASGDNTLSITKGEKLRV-----LGYNHNG
1ycsB  KGVIIY-ALWDYEPQNDELPMKEGDCMTI-----IHREDED
1pht   -GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQ
1ihvA  -----NFRVYYRDSRD--PVWKGPAKLL-----WKGEG
1vie   -----DRVRRKKSG--AAWQGQIVGW-----YCTNL
```

```
1aboA  -----EWCEA--QTKNGQGWVPSNYITPVN-----
1ycsB  EI-----EWWA--RLNDKEGYVPRNLLGLYP-----
1pht   EARPEEIEWGLNGYNETTGERGDFPGTYVEYIGRKKISP
1ihvA  -----AVVIQ--DNSDIKVVPRRKAKIIRD-----
1vie   TP----EGYAVESEAHPGSVQIYPVAALERIN-----
```

SH3 - Alignement fourni par Dialign2

```
1aboA  n-LFVALYDFVASGDNTLSITKGEKLRVL-----
1ycsB  kgVIYALWDYEPQNDELPMKEGDCMTIIhr---EDEDEI-----
1pht   gyQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDgqearpeei
1ihvA  --NFRV---YYRDSRDPVWKGPAKLLWKGE GAVVIQDNSDI-----
1vie   -----DRVRKKS Gaa-W-----QGQI-----
1aboA  ----GYNhngEWCEAQTKNGQGWV-----PSNYItp-----VN
1ycsB  -----EWWARLNDKEGYV-----PRNLLgLYP-----
1pht   gwlnGYN-----ETTGERGDF-----PGTYV-EYigRKKIsp--
1ihvA  -----Kv-----V-----PRr-----KAKIIRd-
1vie   -----VGWYCTNLTPEGYAveseahPGSVQ-IYPv-AALERIN
```

Exemple : 5 protéines, domaine HLH

domaine *helix - loop - helix*

Séquences à aligner:

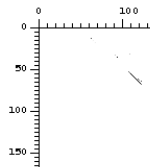
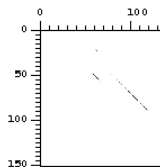
=====

longueur:

=====

1)	HEN1-Human	133
2)	CBF1-Yeast	351
3)	HES5-Mouse	167
4)	INO4-Yeast	151
5)	ESC1-Yeast	413

- longueurs dissemblables
- similarité locale



helix-loop-helix, alignement Clustal

```
-----MMLNSDTMELD-----LPPTHSETESG-----FSDCGGG
--MNSLANNNKLSTEDEEIIHSARKRGYNEEQNYSEARKKQRDQGLLSQESNDGNIDSALLSEGATLKGTQSQYESG-----LTSNKDE
MSSYALPSMQPTTSSIPLRQMSQPTTSAPSNSASSTPYSPQQVPLTHNSYPLSTPSSFQHGQTRLPPINCLAEPFNRPPQWHSNSAAP
-----MAPSTVAVEMLSPEKKNRLRKPVVEKMRDR-----INSSIEQ
-----MTNDIKEIQTIQPGLSEIKEIKGELANVKKR-----

AGPD-----GAGPGG-----
KGSDDDEDASVAEAAVAATVNYTDLIQGE-----DSSDAHTSNQTNANGEHKDSLNGERAITPSNEGVPKNTSLEGMTSSPMEST
ASSSPTSATLSTAHPVHTNAAQVAGSSSSSYVYVPPTNSTTSQASAKHSAVPHRSSQFQSTTLTPSTTSSSDVSSSDSVSTSASSS
LKLL-----

-----PGGGQARGPEPEPGRKD-----LQHLSREERRRRRRAT-----AKYRTA-----
QQSKNDMLIPLAEHGRGPEHQDDDEDNDADID-----LKKDISMQPGRRGKPTTLATDEWKKQKDKS-----
NASNTVSVTSPASSSATPLNPQSQQQLVSKNDAFTTFVHSHVHNTPMQQSMYVPQQQTSHSSGASYQNESANPPVQSPMQYSYSGQP
-----LEQEFARHQPNKLEKAD-----ILEMAVSYLKHSKAFAAAAGPKSLHQDYSEG-----
-----KRRSKKINKLTDGQIR-----INHVSSEKKRRELERAIFDELVAVVPDLQPQ-----

-----HATRERIRVEAFNLFA--ELRKLLPTLPP-----DKKLSKIEILR
-----HKEVERRRRRENINTAIN--VLSDLLPVRESSKAAILARAAEYIQLKETDEANIEKWTQLKLLSEQNASQ
FSYPQHKNQSFSA SPIDPSMSYVYRAPESSFSSINANVPYGRNEYLRVTS LVPNQPEYTG PYTRNPELRTSHKLAERKRKEIKELFDDLKDA
-----YSWCLQEA VQFLT LHAASDTQM KLLYHFQRP-----APAAPAKEPPA
-----ESRSELI IY LKSLSYLSWLYERNEKLR-----KQIIAKHEAKT

LAIC-----YISYLNHVL DV-----
LASANEKLQEELGNAYKEIEYMKRVL RKEGIEYEDMHTHKKQENERKSTRSDNPHEA
LPLDKSTKSSKWGLLTRA IQYIEQLKSEQVALEAYVKSLEENMQSNKEVTKGT-----
PGAAPQPARSSAKAAAAAVTSRQPACGLWRPW-----
GSSSSSDPVQEQNGNIRD LVPKELIWELGDGQSGQ-----
```


helix-loop-helix, alignement Dialign2 (1/2)

```
mml-----  
m-----NSLANNNKLS  
MA-----  
MT-----  
mssyalpsmqptptssiplrqmsqpttsapsnsasstpyspqvplthnsyplstpsffqhggqtrlppinclaePFNRPQWHSNSAAPA SSSPTSATLS  
-----NSDTMELD-----LPPTHSETESGFSDCGGGAGPDgagpggpgggqarg-----  
TEDEEIHSAKRGRGYNEEQNYsearkkqrdqgllsqesndgnidsallsegatLKGTQSQYESGLTSNKDEKGSDdedasvaeaavaatvnytdliqgQED  
-----  
TAAHPVHTNAAQVAGSSSSYVYS-----VPPTNSTTSQAsakhsavphrSSqfqtstltptst-----DSS  
-----PEPGEPRK-----  
SSDAHTSNQTNANGEHKDSLNGERAITPSNEGVP-----NTSLEGMTSSPMESTQQSKNdmliplaehdrg-----  
-----  
STDVSSSDSVSTSASSSNASNTVSVTSPASSSATPLNPQSQqflvskndafttfvhsvhNTPMQQSMYVPQQQTSHSSGasyqnesanppvqspmqys
```

helix-loop-helix, alignement Dialign2 (2/2)

```
-----DLQHL---SREERRRRRRATA-----K
-----PEHQddednddadidlkdkdismqpgrgrkPTTLAtt dew-KKQR-----
-----PSTVAVEMLSPEKN-----
-----NDIKEIQTIQPGLSEIKEIKGELANVKKR---KRRSKKINKLTDG-----Q
ysqgqpfpsyPQHK-----NQSFSASPIDPSMSYVYRAPESFSSINANvpyGRNEYLRrvtslvnpqpeytgpytrnpE

YRTAHATRERIRVEAFNLAFELRKLLPTL----PPDKKLSKIEILRLAICYISYLNHVldv-----
-KDSHKEVERRRRRENINTAINVLSDLLP-V---RESSKAA---ILARAAEYIQKLKETDEanieKWTLQKLLSEQNASQLASANEKLQEELGNaykeie
-RLRKPVVEKMRRDRINSSIEQLKLLLeqefarhQPNSKLEKADILEMAVSYLKHSKAFAA---Aag-----P
IRINHVSSSEKKRRRELAIFDELVAVVPDL---QPQESRSELIIYLSLSYLSWLYERNE---KLRRKQIIAKHEAKTGSSSSSDPVQEQQNGNirdlvP
LRTSHKLAERKKRKEIKELFDDLKDALP-L---DKSTKSSKWGLLTRAIQYIEQLKSEQV---ALEAYVKSLEEnmqsnkevtkgt-----

ymkrv1r-----KEGIEYEDMHThkkqenerkstrsdnphea-----
KSLHQDYSEGYSwclQEAVQFLTLHAasdtqmklllyhfqrppapaapakeppagaapqparssakaaaaavstsrqpacglwrpw
KELIWELGDGQSgq-----
```

Résumé des méthodes

Méthode	Idée	Stratégie
MSA DCA	Extension de l'algorithme de Needleman et Wunsch	Simultanée
Clustal PIMA PILEUP MULTALIGN Dialign	Ajout successif de séquences ou groupes de séquences	Progressive
Saga/Coffee PRRN HMMT MUSCLE MA-FFT	Réalignement lors de l'ajout successif de séquences ou groupes de séquences	Itérative

la comparaison de génomes

- 1 taille de séquences bien plus longue (1000bp \rightarrow $>$ 1000000bp)
- 2 présence de réarrangements/duplications (combinées)

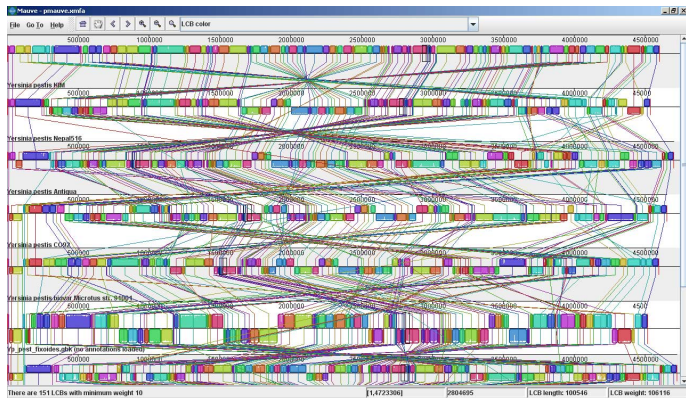
⇒ autant de programmes qui produisent des alignements différents !

Différentes écoles :

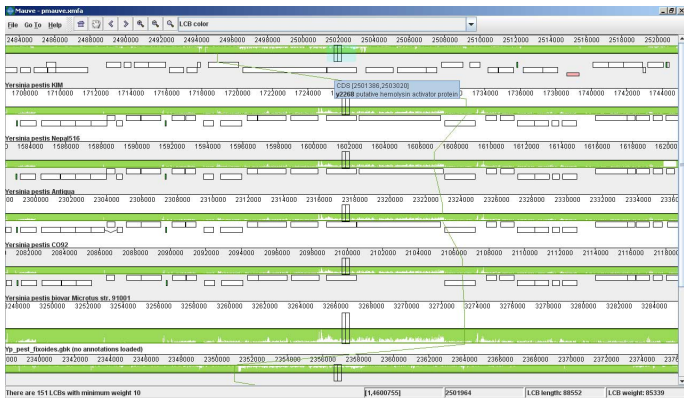
- MGA (Bielefeld)
- MUMmer (Baltimore/Celera genomics)
- Lagan, Multilagan (Stanford)
- MAUVE (Wisconsin-Madison)
- GLASS AVID (Berkeley)
- et bien d'autres ...

La comparaison de génomes - exemple de Mauve

La comparaison de génomes - exemple de Mauve



La comparaison de génomes - exemple de Mauve



La comparaison de génomes - exemple de Mauve

