

Global Multiple Sequence Alignment

```

HUMAN  MKWVTFISLL FLFSSAYSRG V..FRRDA.H KSEVAHRFKD LGEENFKALV
RABIT  MKWVTFISLL FLFSSAYSRG V..FRREA.H KSEIAHRFND VGEEHFGLV
PIG    ~~WVTFISLL FLFSSAYSRG V..FRRDT.Y KSEIAHRFKD LGEQYFKGLV
CHICK  MKWVTLISFI FLFSSATSRN LQRFARDAEH KSEIAHRYND LKEETFKA
  
```

Align k sequences, so that residues in each column share a property of interest:

- a common ancestor
- a structural or functional role

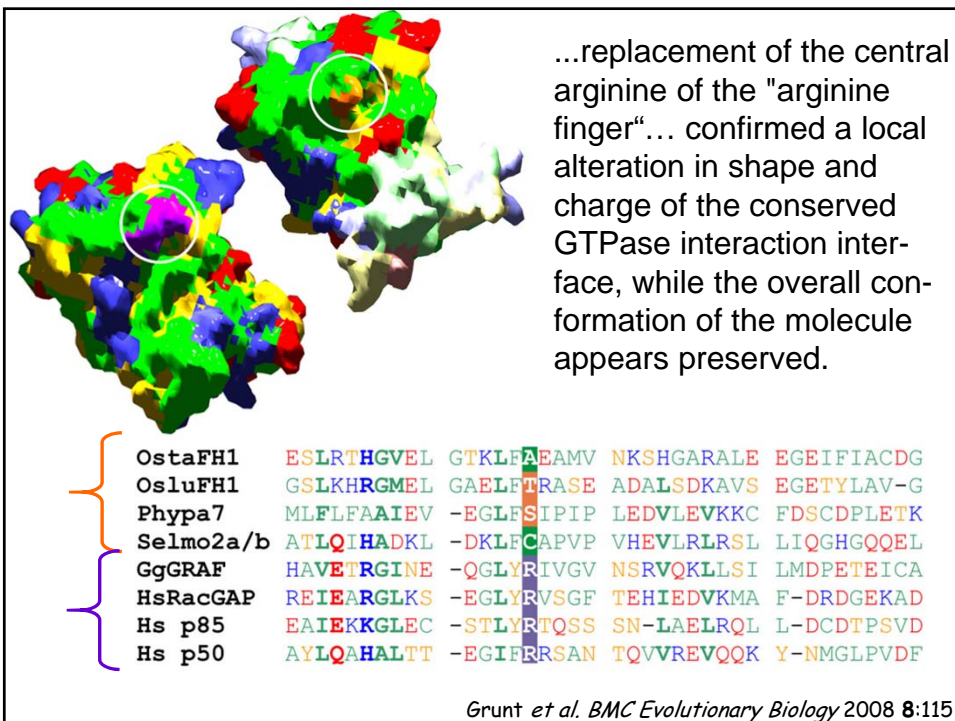
Applications

Global Multiple Alignment

```

HUMAN  MKWVTFISLL FLFSSAYSRG V..FRRDA.H KSEVAHRFKD LGEENFKALV
RABIT  MKWVTFISLL FLFSSAYSRG V..FRREA.H KSEIAHRFND VGEEHFIGLV
PIG    ~-WVTFISLL FLFSSAYSRG V..FRRDT.Y KSEIAHRFKD LGEQYFKGLV
CHICK  MKWVTLISFI FLFSSATSRN LQRFARDAEH KSEIAHRYND LKEETFKA
  
```

- Protein structure and function
- RNA structure
- Evolutionary tree reconstruction



A MULTIPLE SEQUENCE ALIGNMENT EXAMPLE

Ribosome: an RNA/protein complex

rpS14: a ribosomal protein in yeast

Goal: Determine residues responsible for binding rpS14 to ribosomal RNA

Known:

- Sequence of rpS14
- Structure of homolog in bacteria
- Sequences in many species

Pam Bush, PhD CMU, 02

A MULTIPLE SEQUENCE ALIGNMENT EXAMPLE

Strategy: alanine scan

- Find “likely” candidate amino acids
- Replace candidates with alanine
- Hope that
 - Alanine preserves structure
 - Alanine will destroy binding

Seventeen residues are conserved in rpS14

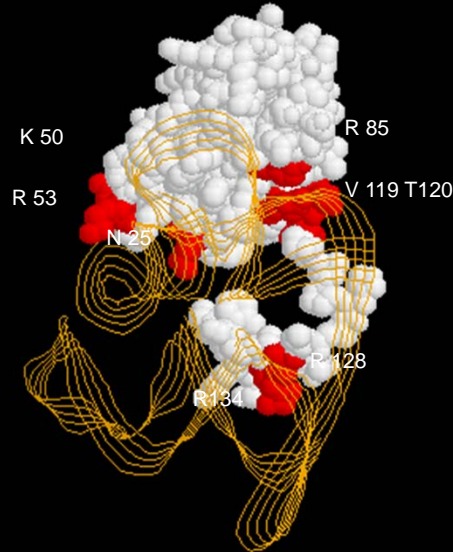
| | | | | | |
|------------------------|------------|------------------|------------------|-----------------|------------------|
| <i>T. thermophilus</i> | -----MAK | KPSKKKVKRQVASGR | AYIHASYNNITIVTIT | DPDGNPITWSSGGVI | GYKGSR-KGTFYAAQ |
| <i>A. aeolicus</i> | -----M | AKKKKKQKRQVTKAI | VHIHTTFNNTIVNVT | DTQGNITAWASGGTV | GFKGTR-KSTPYAAQ |
| <i>P. aeruginosa</i> | -----MAKPA | ARPRKKVKKTIVVDGI | AHIHASFNNITIVTIT | DRQGNALSWATSGGS | GFRGSR-KSTPFYAAQ |
| <i>E. coli</i> | -----MAKAP | IRARKRVKQVSDGV | AHIHASFNNITIVTIT | DRQGNALGWATAGGS | GFRGSR-KSTPFYAAQ |
| <i>H. sapiens</i> | | MAPRKGKEKKEEQVI | SLGPQVAEGENVFVG | CHIFASFNDTFVHVT | DLSGKETICRVTTGM |
| <i>D. melanogaster</i> | | MAPRKAKVQKEEVQV | QLGPQVRDGEIVFVG | AHIYASFNDTFVHVT | DLSGRETIVRTVGM |
| <i>S. pombe</i> | -----MAT | NVGPQIRSGELVFGV | AHIYASFNDTFVHIT | DLTGKETIVRTVGM | KVKIDRDESSPYAAM |
| <i>S. cerevisiae</i> | -----MA | NDLVQARDNSQVFGV | ARIYASFNDTFVHVT | DLSGKETIARVTGM | KVKADREDESSPYAAM |
| <i>S. solfataricus</i> | ----- | MSSRRREIRWGI | AHIYASQNNITLTIS | DLTGAEIISRASGM | VVKADREKSSPYAAM |
| <i>M. jannaschii</i> | ----- | MAEQKKEKWGI | VHIYSSYNNITIIHAT | DITGAETIARVSGGR | VTRNQREDESSPYAAM |

| | | | | | |
|------------------------|-----------------|-----------------|-----------------|-----------------|------------------|
| <i>T. thermophilus</i> | LAALDAAKKAMAYGM | QSVDVIVRG----- | --TGAGREQAIRALQ | ASGLQVKSIVDDTPV | PHNGCRPKKKFFKAS- |
| <i>A. aeolicus</i> | LAAQKAMKEAKEHGV | QEVEIWWKG----- | --PGAGRESAVRAVF | ASGVKVTAIRDVTPI | PHNGCRPPARRRV--- |
| <i>P. aeruginosa</i> | VAAERAGQAALFYGL | KNLDVNVKG----- | --PGPGRESAVRALN | ACGYKIASITDVTPI | PHNGCRPPKKRRV--- |
| <i>E. coli</i> | VAAERCADAVKEYGI | KNLEVMVKG----- | --PGPGRESTIRALN | AAGFRITNITDVTPI | PHNGCRPPKKRRV--- |
| <i>H. sapiens</i> | LAAQDVQRCKELGI | TALHIKLRATGGNRT | KTPGPGAQSALRALA | RSGMKIGRIEDVTPI | PSDSTRKKGRRGRRL |
| <i>D. melanogaster</i> | LAAQDVAEKCKTLGI | TALHIKLRATGGNKT | KTPGPGAQSALRALA | RSSMKIGRIEDVTPI | PSDSTRKKGRRGRRL |
| <i>S. pombe</i> | LAAQDAAAKCKEVGI | TALHIKIRATGGTAT | KTPGPGAQAALRALA | RAGMRIGRIEDVTPI | PTDSTRKKGRRGRRL |
| <i>S. cerevisiae</i> | LAAQDVAACKCEVGI | TAVHVKIRATGGTAT | KTPGPGGQAALRALA | RSGLRIGRIEDVTPI | PSDSTRKKGRRGRRL |
| <i>S. solfataricus</i> | LAANKAASDALEKGI | MALHIKVRAPGGYGS | KTPGPGAQPAIRALA | RAGFLIGRIEDVTPI | PHDTIRRPGGRRGRRV |
| <i>M. jannaschii</i> | QAAPKLAEVLKERGI | ENIHIKVRAPGGSGQ | KNPGPGAQAALRALA | RAGLRIGRIEDVTPI | PHDGTTPKKRFFK--- |

Criteria for selecting conserved amino acids:

- Conserved among all three phylogenetic groups
- Conserved in at least 90% sequences analyzed, allowing for conservative substitution K↔R
- No residue could be an alanine, proline or glycine

Eight conserved residues are on the surface of the bacterial protein interacting with the rRNA



...and these are distributed throughout the protein sequence

| | | | | | |
|------------------------|-------------------|-----------------|-----------------|------------------|------------------|
| <i>T. thermophilus</i> | -----MAK | KPSKKVKRQVASGR | AYIHASYNNTIVTIT | DPDGNPITWSSGGVI | GYKGSR-KGTPYAAQ |
| <i>A. aeolicus</i> | -----M | AKKKKKQKRQVTKAI | VHIHTTFNNTIVNVT | DTQGNNTIAWASGGTV | GFKGTR-KSTPYAAQ |
| <i>P. aeruginosa</i> | -----MAKPA | ARPRKKVKKTVVDGI | AHIHASFNNTIVTIT | DRQGNALSWATSGGS | GFRGSR-KSTPFAAQ |
| <i>E. coli</i> | -----MAKAP | IRARKVRKQVSDGV | AHIHASFNNTIVTIT | DRQGNALGWATAGGS | GFRGSR-KSTPFAAQ |
| <i>H. sapiens</i> | MAPRKGEKKKEEQVI | SLGPQVAEGENVFVG | CHIPASFNDTFVHVT | DLSGKETICRVVTGGM | KVKADRDDESSPYAAM |
| <i>D. melanogaster</i> | MAPRKAKVQKEEVQV | QLGPQVRDGEIVFVG | AHIYASFNDTFVHVT | DLSGKETIARVTGGM | KVKADRDDESSPYAAM |
| <i>S. pombe</i> | -----MAT | NVGQIRSGELVFGV | AHIFASFNDTFVHIT | DLTGKETIVRVVTGGM | KVKIDRDDESSPYAAM |
| <i>S. cerevisiae</i> | -----MA | NDLVQARDNSQVFGV | ARIYASFNDTFVHVT | DLSGKETIARVTGGM | KVKADRDDESSPYAAM |
| <i>S. solfataricus</i> | -----MSSRREIRWGI | AHIYASQNTLLTIS | DLTGAEIISRASGGM | VVKADREKSSPYAAM | |
| <i>M. jannaschii</i> | -----MAEQKKEKWKGI | VHIYSSYNNTIIHAT | DITGAETIARVSGGR | VTRNQDDEGSPYAAM | |

| | | | | | |
|------------------------|-----------------|-----------------|-----------------|-----------------|------------------|
| <i>T. thermophilus</i> | LAALDAAKKAMAYGM | QSVDIVRG----- | --TGAGREQAIRALQ | ASGLQVKSIVDDTPV | PHNGCRPKKKFKAS- |
| <i>A. aeolicus</i> | LAAQKAMKEAKEHGV | QEVEIWNKG----- | --PGAGRESAVRAVF | ASGVKVTAIRDVTPV | PHNGCRPPARRRV--- |
| <i>P. aeruginosa</i> | VAAERAGQAALAYGL | KNLDVNVKG----- | --PGPGRESAVRALN | ACGYKIASITDVTPV | PHNGCRPPKKRRV--- |
| <i>E. coli</i> | VAAERCADAVKEYGI | KNLEVMVKG----- | --PGPGRESTIRALN | AAGFRITNITDVTPV | PHNGCRPPKKRRV--- |
| <i>H. sapiens</i> | LAAQDVQQRCKELGI | TALHIKLRATGGNRT | KTPGPGAQSALRALA | RSGMKIGRIEDVTPV | PSDSTRRKGGRRGRRL |
| <i>D. melanogaster</i> | LAAQDVAEKCKTLGI | TALHIKLRATGGNKT | KTPGPGAQSALRALA | RSGMKIGRIEDVTPV | PSDSTRRKGGRRGRRL |
| <i>S. pombe</i> | LAAQDAAKCKEVGI | TALHIKIRATGGTAT | KTPGPGAQAALRALA | RAGMRIGRIEDVTPV | PTDSTRRKGGRRGRRL |
| <i>S. cerevisiae</i> | LAAQDVAACKKEVGI | TAVHVKIRATGGTRT | KTPGPGGAALRALA | RSGLRIGRIEDVTPV | PSDSTRRKGGRRGRRL |
| <i>S. solfataricus</i> | LAANKAASDALEKGI | MALHIKVRAPGGYGS | KTPGPGAQPAIRALA | RAGFLIGRIEDVTPV | PHDTIRRPGGRRGRRV |
| <i>M. jannaschii</i> | QAAPKLAEVVKERGI | ENIHIKVRAPGGSGQ | KNPGPGAQAALRALA | RAGLRIGRIEDVTPV | PHDGTTPRKKRKK--- |

About the changes to alanine of the conserved residues:

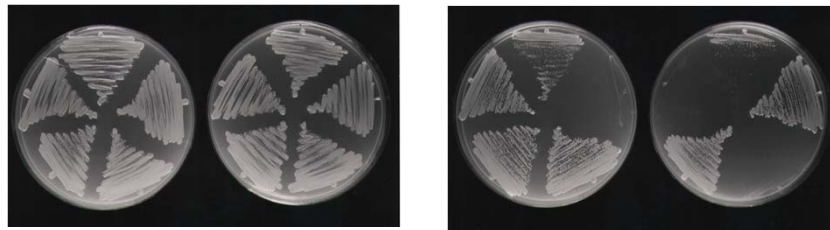
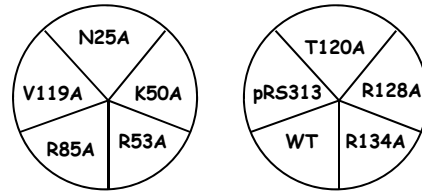
N25, K50, R53, R85, V119, T120, R128, R134

- Like most ribosomal proteins, high percentage of basic residues
- Change to alanine is a big charge difference
- Alanines on surface usually do not have a big structural change

The predicted secondary structure of
rpS14 does not change significantly when the
conserved residues are changed to alanine



Functional assays suggest that K50A and R134 play a role in rpS14/RNA binding



Alignment of the PI3 and cAMP dependent kinases

Local:

PI3-kinase DRHNSNIMVKDDGGLFH:DTG
cAMP PK DLKPEMLLDQGGYIQV:DTG

Global:

PI3-kinase 120 130 140 150 160
QFNSHT-LHRQLKDKKKG ELYDA--IDLFTSCAGYCVATFILGIDRHNSNIMVKD-D
cAMP PK 110 120 130 140 150 160
SFKDNSNLYVMEYVPGGERFSLRRIIGFSEPHARYAAQIVLTFYLSLDLIYHDLK

PI3-kinase 170 180 190 200 210 220
GGLFHIDTGHFLDHKKKFGYKRERVP----FVLTQDFL--IVISKGAQECTKTREFE
cAMP PK 170 180 190 200 210 220
PENLLIDQGGYI--QVDTGFAK-RVKGRIVXLCGTPEYLAPEIILSKGYNKAQDQWALG

Multiple alignment with other kinases:

p110β SYVLGIG-----DRHSDNINVKKTGGLFHIDFGHILGNFKSKFGIKRERVPFILT 136
p110δ TYVLGIG-----DRHSDNIMIRESGGLFHIDFGHFLGNFKTKFGINRERVPFILT 136
p110α IFILGIG-----DRHSDNIMVKDDGGLFHIDFGHFLDHKKKFGYKRERVPFVLT 135
p110γ TFVLGIG-----DRHSDNIMITETGNLFHIDFGHILGNYSKFLGINKERVPFVLT 135
p110_dicti TYVLGIG-----DRHSDNLMVTKGGRLFHIDFGHFLGNYSKFGFKRERAPFVFT 135
cAMP-kinase QIVLTFEYLSLDLIYRDLKPEMLLDQGGYIQVDTGFAKRVKGRIVXLCG--TPEYLA 177

From Zvelebil and Baum, "Understanding Bioinformatics"

Global Multiple Sequence Alignment

Given sequences $s_1 \dots s_k$ of lengths $n_1 \dots n_k$


seek $s'_1 \dots s'_k$ of length $l \geq \max\{n_i\}$ such that

- Obtain s_i from s'_i by removing gaps
- No column contains all gaps
- The score of the alignment is optimal

Scoring function: Sum-of-Pairs

$$\text{Score} = \sum_{a=1}^k \sum_{b=1}^k \sum_{b > a} p(s'_a[i], s'_b[i])$$

(1) **A** T **T**
 (2) **A** T **_**
 (3) **A** **C** **A** **T**





$p[_, _] = 0$

$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= 0 + g + g = 2g \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

Scoring function: Sum-of-Pairs

$$\text{Score} = \sum_{a=1}^k \sum_{a=1}^k \sum_{b>a} p(s'_a[i], s'_b[i])$$




(1) **A_TT**   $p[_,_]=0$

$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= M + m + m = 2m + M \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

Scoring function: Sum-of-Pairs

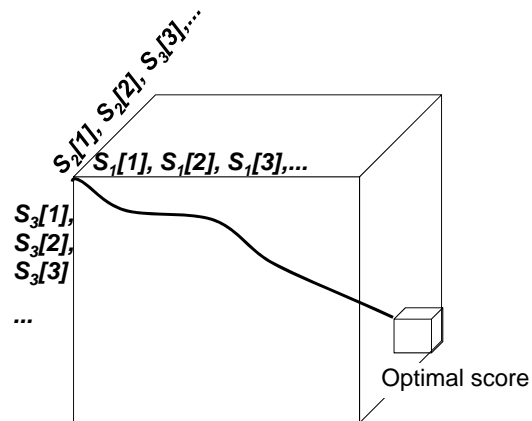
$$\text{Score} = \sum_{a=1}^k \sum_{a=1}^k \sum_{b>a} p(s'_a[i], s'_b[i])$$

(1) **A_TT** 
 (2) **A_T_** 
 (3) **ACA**T****  $p[_,_] = 0$

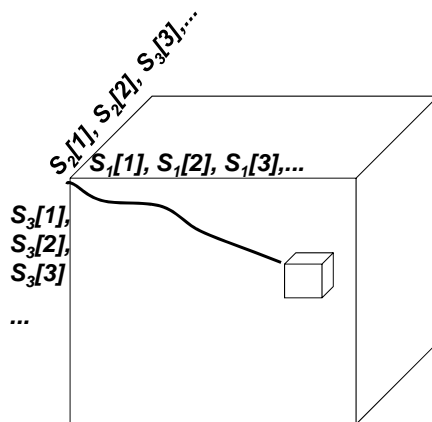
$$\begin{aligned} \text{Score} &= p[s_1, s_2] + p[s_1, s_3] + p[s_2, s_3] \\ &= g + M + g = 2g + M \end{aligned}$$

Note: this example uses a similarity function. We can also use Sum-of-Pairs with distance scoring.

Dynamic Programming for Multiple Alignment



Dynamic Programming for Multiple Alignment



Limits:

- ~ $k = 8 - 10$ sequences
- ~ $n = 500$ residues

Each cell has $O(2^k)$ neighboring cells

Calculating the sum-of-pairs score for each neighbor is $O(k^2)$

Number of cells in matrix: $O(n^k)$

Total computational complexity:
 $O(n^k 2^k k^2)$

MSA is NP-complete for Sum-of-Pairs scoring

Observations

1. A multiple alignment induces pairwise alignments
2. A column in the induced pairwise alignment may contain all gaps, even though no column in the MSA contains all gaps.

(1) **AG** CT
 (2) **AG** CT
 (3) **ACT** T

3. The pairwise alignments induced by the *optimal multiple alignment* are *not* the same as the *optimal pairwise alignments*.

Optimal Pairwise Alignments

(1) **ACT**
 (2) **AGT**

Optimal Multiple Alignment

(1) **AC** T
 (2) **A** GT
 (3) **ACGT**

1 substitution

2 indels

Although this costs more, it may be a biologically more realistic alignment

Since exact methods for MSA have exponential time complexity, heuristic approaches are used.
Progressive alignment is the most commonly used.

Basic progressive alignment strategy:

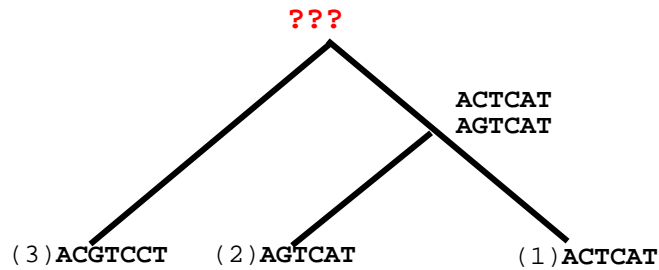
- Compute D , a matrix of distances between all pairs of sequences
- From D , construct a “guide tree” T
- Construct MSA by pairwise alignment of partial alignments (“profiles”) guided by T
- Improve alignment by postprocessing steps.

Optimal Pairwise Alignments

| | | | |
|--------------------|--|--------------------|---|
| | | (1) ACTCAT | 3 |
| | | (2) AGTCAT | |
| (1) ACTCAT | | | |
| (2) AGTCAT | | (2) A_GTCAT | 5 |
| (3) ACGTCCT | | (3) ACGTCCT | |
| | | | |
| | | (1) AC_TCAT | 5 |
| | | (3) ACGTCCT | |

$d(x, y) = 3$
 $d(x, \text{"_"}) = 2$

Progressive Alignment



- Use *profile alignment* to merge sequences according to a guide tree.
- Typically, most closely related sequences are merged first.

Merging strategy:

Align the profile (1,2) with sequence (3)

(1) ACTCAT
(2) AGTCAT
(3) ACGTCCT

| | | |
|-----|--------|---|
| (1) | ACTCAT | 3 |
| (2) | AGTCAT | |

(2) A_GTCAT
(3) ACGTCCT

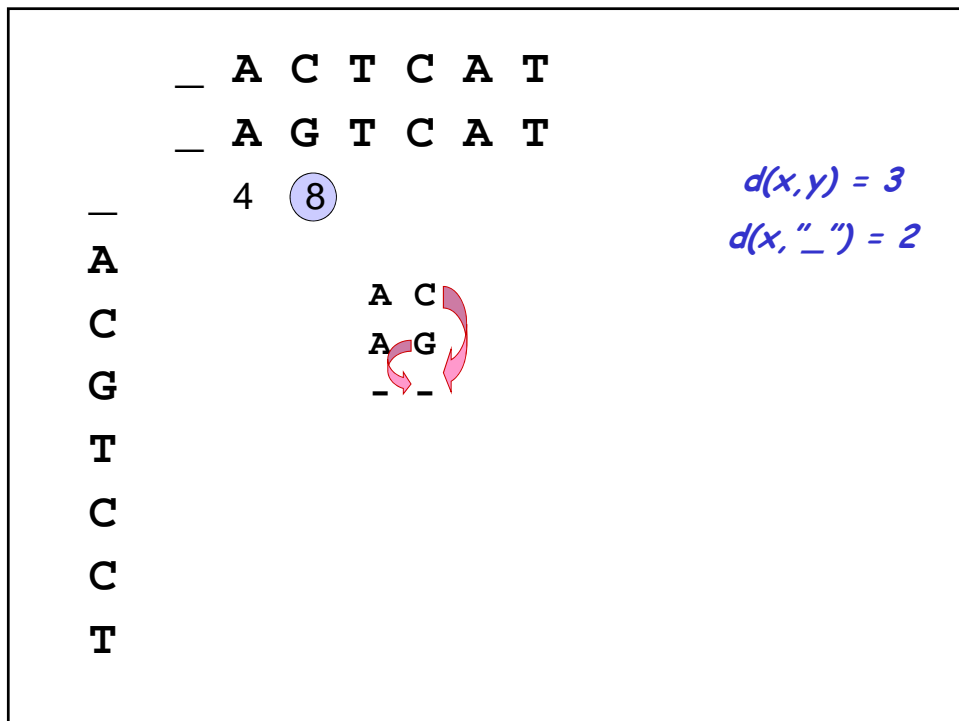
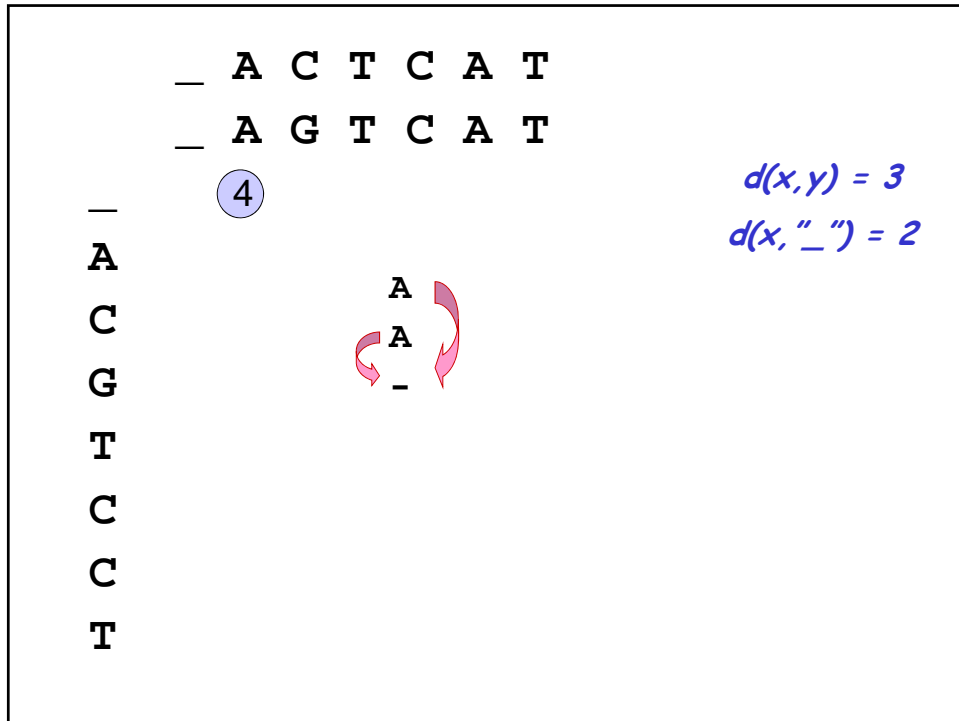
5

$$d(x, y) = 3$$

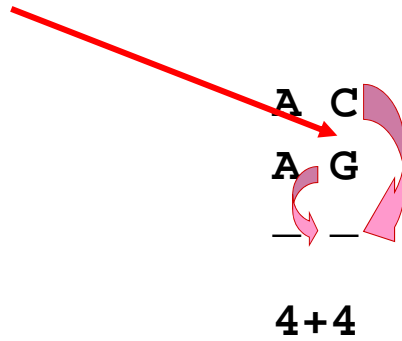
$$d(x, \text{"_"}) = 2$$

(1) AC_TCAT
(3) ACGTCCT

5



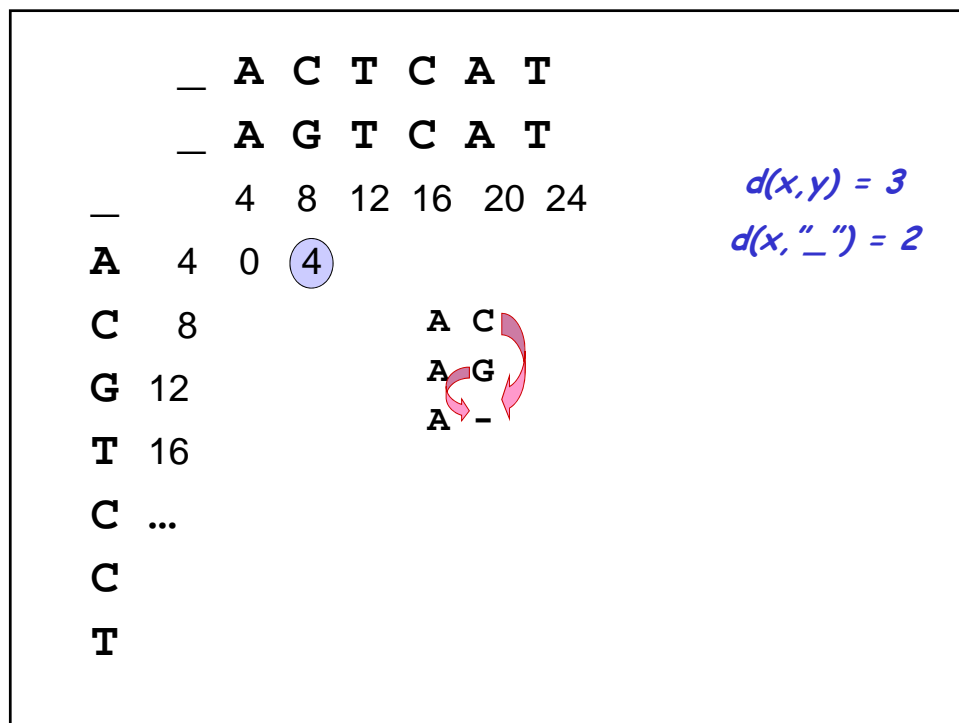
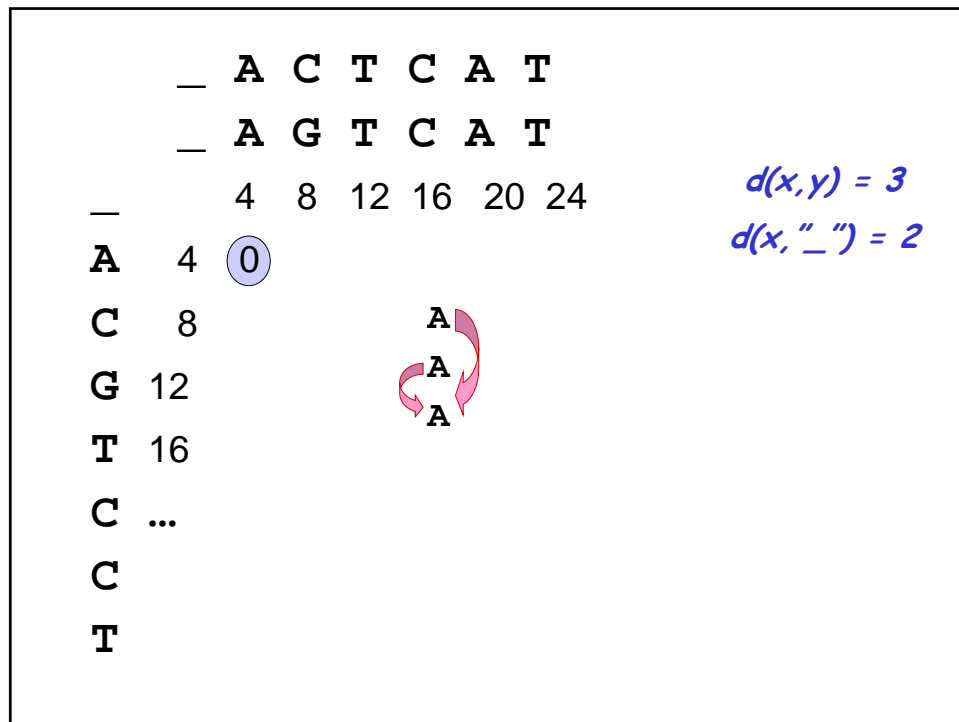
Note: no penalty for mutations in the profile.
We paid for those in a previous step



| | | | | | | | |
|----------|-----|----------|----------|----------|----------|----------|----------|
| | — | A | C | T | C | A | T |
| | — | A | G | T | C | A | T |
| — | | 4 | 8 | 12 | 16 | 20 | 24 |
| A | 4 | | | | | | |
| C | 8 | | | | | | |
| G | 12 | | | | | | |
| T | 16 | | | | | | |
| C | ... | | | | | | |
| C | | | | | | | |
| T | | | | | | | |

$$d(x, y) = 3$$

$$d(x, \text{"_"}) = 2$$



| | | | | | | | | | |
|---|-----|---|---|---|----|----|----|----|------------------------|
| | | — | A | C | T | C | A | T | |
| | | — | A | G | T | C | A | T | |
| | — | | 4 | 8 | 12 | 16 | 20 | 24 | $d(x, y) = 3$ |
| A | 4 | 0 | 4 | 8 | | | | | $d(x, \text{"_"}) = 2$ |
| C | 8 | | | | | | | | |
| G | 12 | | | | | | | | |
| T | 16 | | | | | | | | |
| C | ... | | | | | | | | |
| C | | | | | | | | | |
| T | | | | | | | | | |

| | | | | | | | | | |
|---|-----|---|---|---|----|----|----|----|------------------------|
| | | — | A | C | T | C | A | T | |
| | | — | A | G | T | C | A | T | |
| | — | | 4 | 8 | 12 | 16 | 20 | 24 | $d(x, y) = 3$ |
| A | 4 | 0 | 4 | 8 | 12 | | | | $d(x, \text{"_"}) = 2$ |
| C | 8 | | | | | | | | |
| G | 12 | | | | | | | | |
| T | 16 | | | | | | | | |
| C | ... | | | | | | | | |
| C | | | | | | | | | |
| T | | | | | | | | | |

| | | | | | | | | | |
|---|-----|---|---|---|----|----|----|----|--|
| | | — | A | C | T | C | A | T | |
| | | — | A | G | T | C | A | T | |
| | — | | 4 | 8 | 12 | 16 | 20 | 24 | |
| A | 4 | 0 | 4 | 8 | 12 | 16 | | | |
| C | 8 | | | | | | | | |
| G | 12 | | | | | | | | |
| T | 16 | | | | | | | | |
| C | ... | | | | | | | | |
| C | | | | | | | | | |
| T | | | | | | | | | |

$d(x, y) = 3$
 $d(x, \text{" "}) = 2$

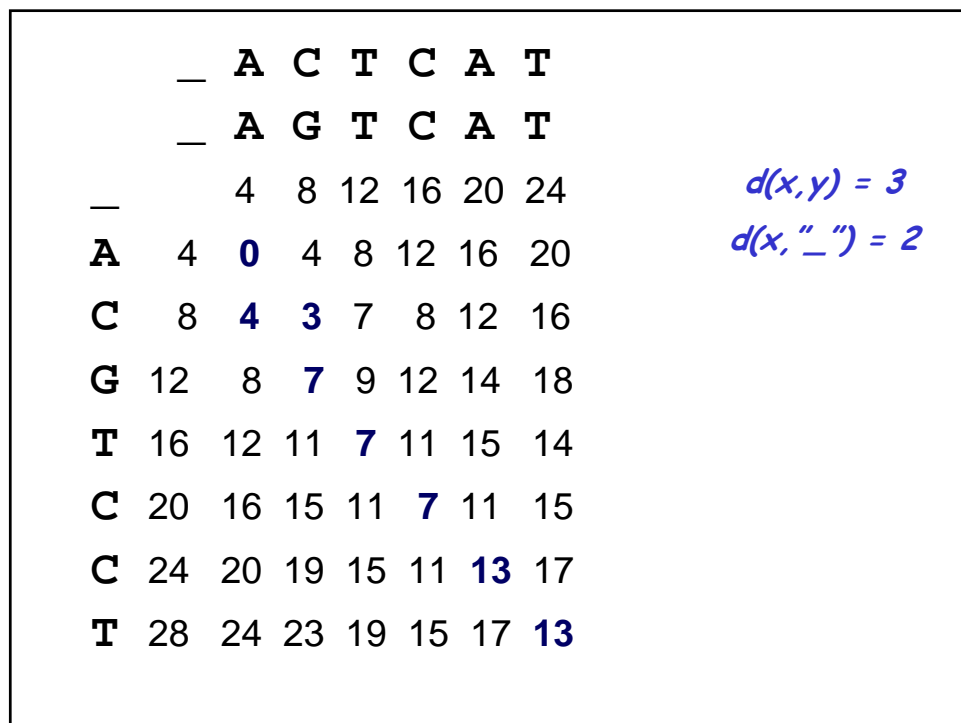
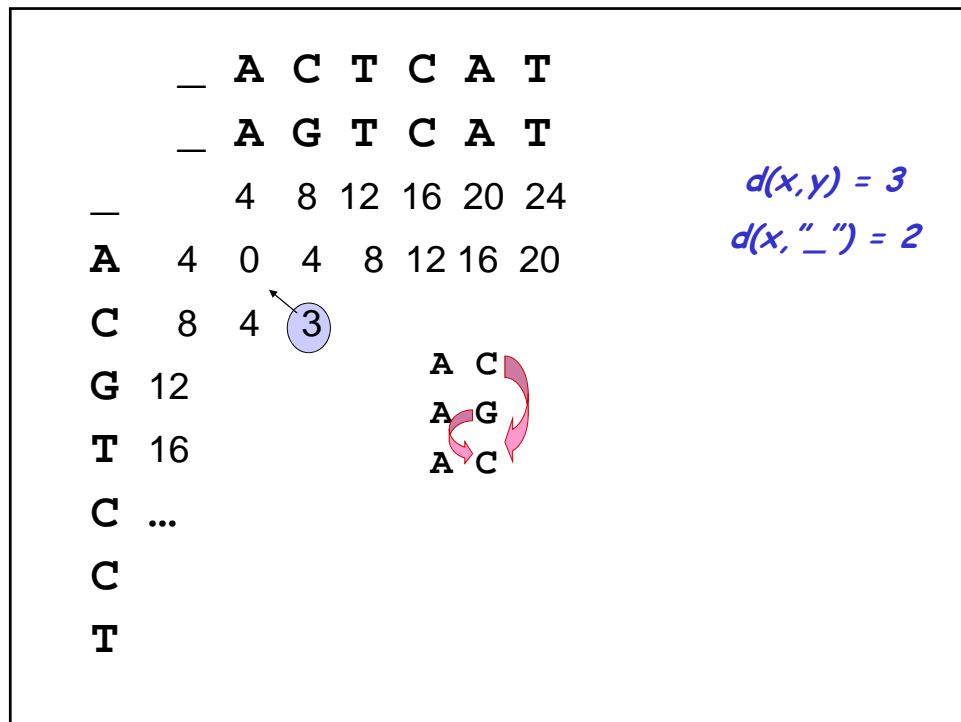
| | | | | | | | | | |
|---|-----|---|---|---|----|----|----|----|--|
| | | — | A | C | T | C | A | T | |
| | | — | A | G | T | C | A | T | |
| | — | | 4 | 8 | 12 | 16 | 20 | 24 | |
| A | 4 | 0 | 4 | 8 | 12 | 16 | 20 | | |
| C | 8 | 4 | | | | | | | |
| G | 12 | | | | | | | | |
| T | 16 | | | | | | | | |
| C | ... | | | | | | | | |
| C | | | | | | | | | |
| T | | | | | | | | | |

$d(x, y) = 3$
 $d(x, \text{" "}) = 2$

A —

A —

A C



| Optimal Pairwise Alignments | | Progressive alignment | |
|-----------------------------|---------|------------------------|-----------------|
| | | (1,2) + (3) | |
| (1) | ACTCAT | (3) | ACGTCCT |
| (2) | AGTCAT | (1) | AC_TCAT $4m+2g$ |
| | | (2) | AG_TCAT |
| (2) | A_GTCAT | An alternate alignment | |
| (3) | ACGTCCT | | |
| (1) | AC_TCAT | (1) | AC_TCAT |
| (3) | ACGTCCT | (2) | A_GTCAT $2m+4g$ |
| | | (3) | ACGTCCT |

| Optimal Pairwise Alignments | | Progressive alignment | |
|-----------------------------|---------|------------------------|--------------|
| | | (1,2) + (3) | |
| (1) | ACTCAT | (3) | ACGTCCT |
| (2) | AGTCAT | (1) | AC_TCAT 16 |
| | | (2) | AG_TCAT |
| (2) | A_GTCAT | An alternate alignment | |
| (3) | ACGTCCT | | |
| (1) | AC_TCAT | (1) | AC_TCAT |
| (3) | ACGTCCT | (2) | A_GTCAT 14 |
| | | (3) | ACGTCCT |

Progressive alignment

- “Once a gap, always a gap”
 - You can’t go back and correct a bad decision at an earlier step.
- Progressive alignment is not guaranteed to give the optimal alignment.
- But it does have better complexity...

Complexity of progressive alignment

- Distance matrix
 - Each pairwise alignment $O(n^2)$
 - Number of pairwise alignments $O(k^2)$
 - Iterative construction of MSA
 - Number of merge steps $O(k)$
 - Each pairwise alignment $O(k^2n^2)$
- Entire method $O(k^2n^2)$

Summary: Progressive alignment heuristics

- Not guaranteed to give the optimal MSA
- Bad choice of gaps propagates
- Complexity
 - Progressive: $O(k^2 n^2)$
 - versus DP: $O(n^k 2^k k^2)$
- Typically, merge the most closely related sequences first.