# Figures and Legends for

## "Blueprint of the "<u>A</u>lignment <u>N</u>eighborhood <u>Ex</u>plorer" (ANEX)
## <span style="color:red">(tentatively named)</span>"
## by  Kiyoshi Ezawa

(Finished on November 24<sup>th</sup>, 2015; added CC4 statement on August 14<sup>th</sup>, 2020)

# Figures and Legends



**Figure 1. Flowchart explaining the overall workflow of the program.**

In the gap configuration in (i), the 'N' and '-' represent the presence of a residue and a gap (*i.e.*, the absence of a residue), respectively.

And the blue and red rectangles enclose the gapless and gapped segments, respectively.

On the right of (i), the alignment of the regions shaded in red and blue indicates a "purge"-type error.

In (ii), the yellow-shaded region likely contains a "complex" error.

**[ I will rewrite the figure later, especially by fleshing out the steps (iii) and (iv). ]**

**Figure 2. Sliding-window analysis to identify regions that likely harbor "purge"-type errors.**

**A**. In this example, we search for possible "purge"s between sub-alignments enclosed by dashed rectangles, which are mutually separated by the red bold branch (on the left).

**B**. The window (shaded gray) contains a specified number of residue pairs, and slides from left to right (indicated by the black solid arrow).

**C**. This dual-colored window shows more predicted substitutions than expected from the branch length.

Thus, the program suspects that a "purge"-type error likely resulted in this region.

**Figure 3. Gap-block, "isolated" gap-block, and sequence-block.**

To focus on the gap-pattern, each residue was represented as 'N' regardless of its identity.

**A.** The yellow-shaded portion of the MSA is a gap-block. It is delimited by two inter-column positions (the vertical red dashed lines) and a branch in the tree (the red branch). (The identifiers of the involved sequences are also shaded in yellow.)

**B, C.** "Isolated" gap-blocks (shaded portions enclosed by dashed boxes). In panel B, the isolation is obvious because the horizontal positions of the gap-blocks do not overlap. In panel C, although the two gap-blocks overlap in their horizontal positions, they are separated with each other by three branches (in red).

**D.** A pair of "non-isolated" gap-blocks, which horizontally overlap and are separated by only two branches (in red).

**E.** A sequence-block (in blue), which is the complement of a gap-block (shaded in yellow).

**(A) "Shift"**

```
1  NNN              1  NNN
2  N--       ➡      2  --N
3  NNN              3  NNN
4  NNN              4  NNN
```

**(B) "Purge"**

```
1  -NN              1  NN
2  NN-       ➡      2  NN
3  -NN              3  NN
4  -NN              4  NN
```

**(C) "Merge" (same)**

```
1  NNNN             1  NNNN
2  -N--      ➡      2  ---N
3  NNNN             3  NNNN
4  NNNN             4  NNNN
```

**(D) "Merge" (complementary)**

```
1  -NNN             1  NNN
2  NN--      ➡      2  -NN
3  -NNN             3  NNN
4  -NNN             4  NNN
```

**(E) "Vertical merge" (of gap-blocks)**

```
1  -NN              1  -NN
2  NN-       ➡      2  -NN
3  NNN              3  NNN
4  NNN              4  NNN
```

**(F) "Vertical merge" (of sequence-blocks)**

```
1  -NN              1  NN
2  NN-       ➡      2  NN
3  -N-              3  N-
4  -N-              4  N-
```

**(G) "Shift + shift" (???)**

???

**Figure 4. Elementary moves that program will attempt.**

The bold "N"s represent the residues that moved.

**A.** A "shift" of a single gap-block (shaded portion enclosed by the dashed box).

**B.** A "purge" of two gap-blocks affecting the complementary sets of sequences. (Its reverse is an "ex-nihilo.")
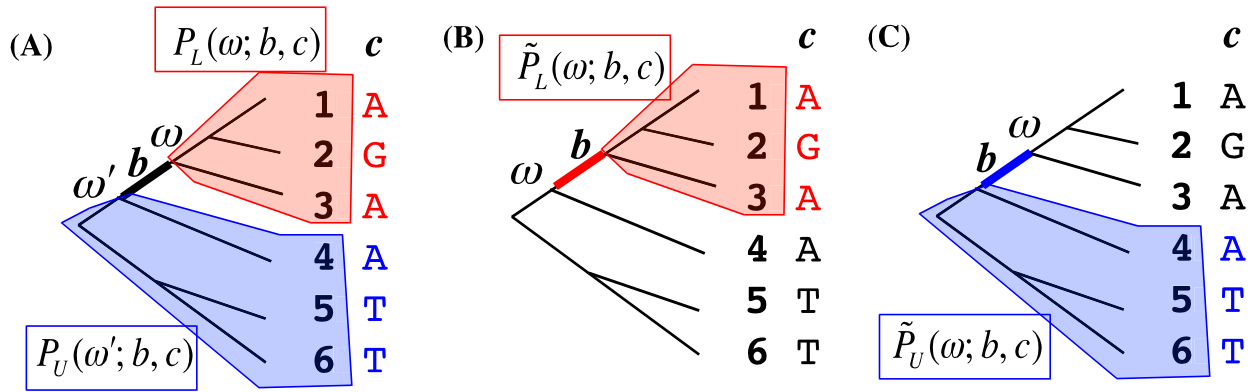
**C.** A "merge" of two gap-blocks affecting the same set of sequences. (Its reverse is a "split" of a gap-block.)

**D.** A "merge" of gap-blocks affecting the complementary sets of sequences. (Its reverse is a "split" of a gap-block.)

**E.** A "vertical-merge" of gap-blocks. (Its reverse is a "vertical-split" of a gap-block.)

**F.** A "vertical-merge" of sequence-blocks. (Its reverse is a "vertical-split" of a sequence-block.)
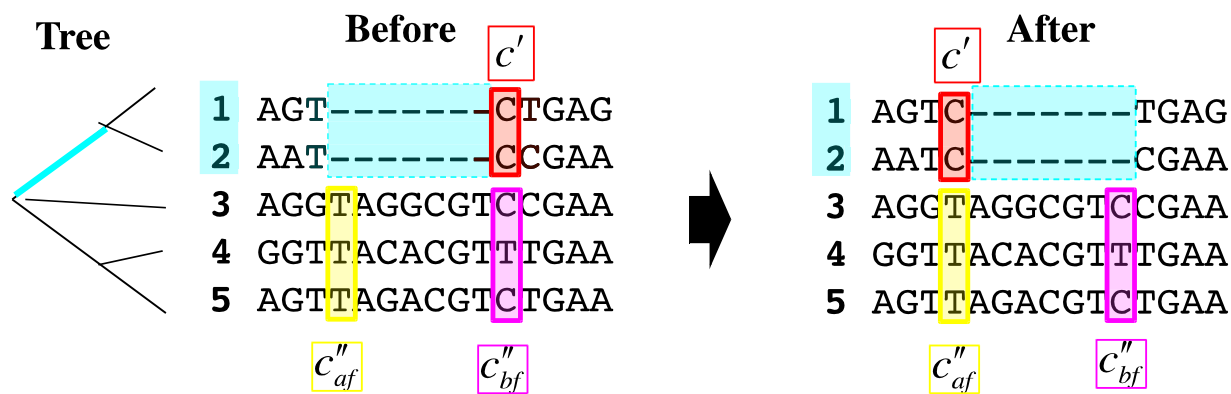
**G.** "Shift + Shift" …???

**Figure 5. Four sets of important probabilities, assigned to each branch-column combination, for fast calculation of MSA residue-pattern probabilities.**

**A.** Probabilities, $P_L(\omega; b, c)$ and $P_U(\omega'; b, c)$, for the residue configurations of two complementary sequence sets in a column ($c$). The branch $b$ (bold) separates the sequence sets. The portions of the tree yielding $P_L(\omega; b, c)$ and $P_U(\omega'; b, c)$ are shaded in red and blue, respectively. And, in column $c$, the color of each residue indicates which of the probabilities it contributes to.

**B.** The extension of $P_L(\omega; b, c)$.

**C.** The extension of $P_U(\omega'; b, c)$.

In each panel, the numbers assigned to the external nodes also specify the sequences in the MSA.
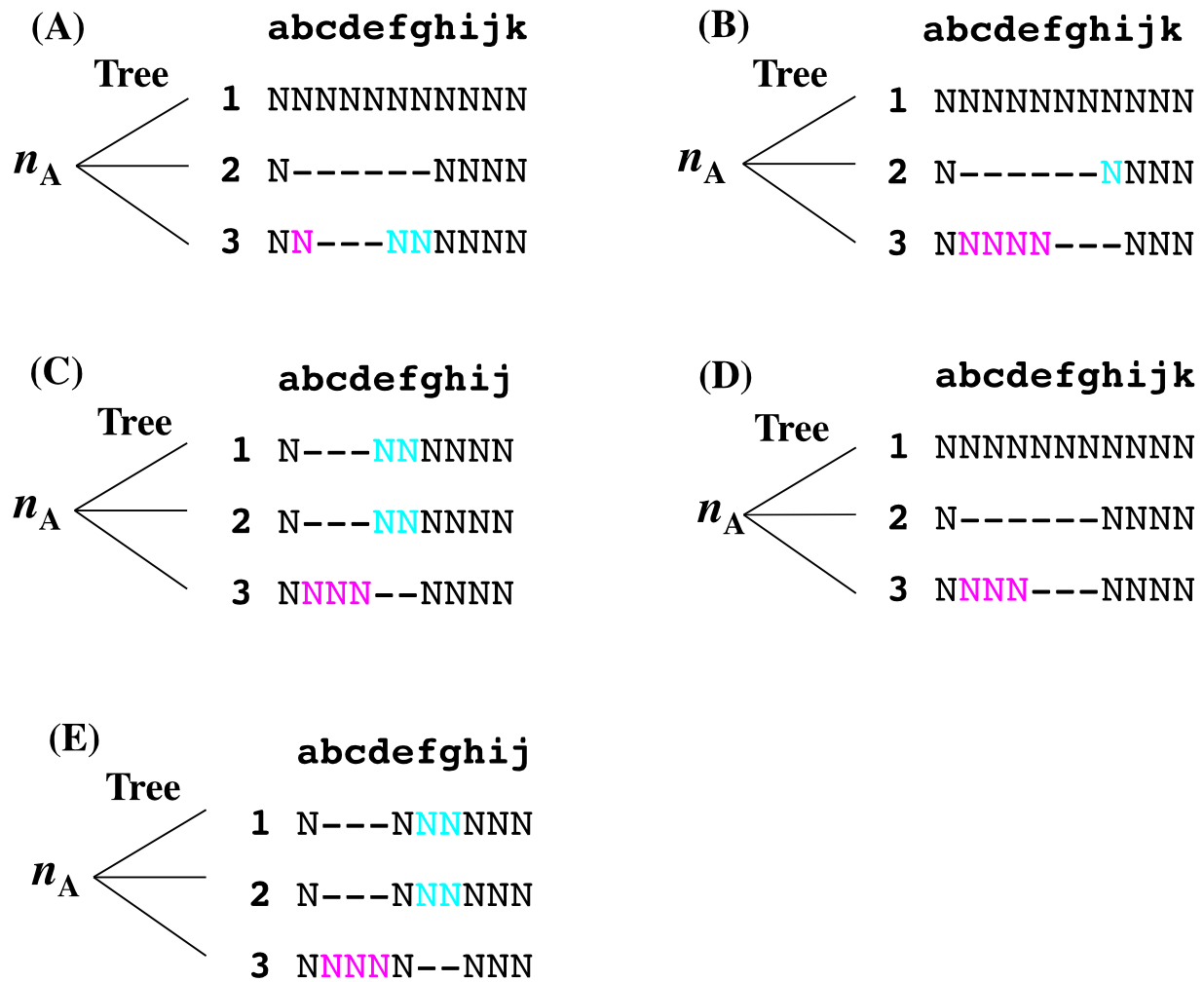
**Figure 6. Change in residue component of MSA probability when gap-block shifts.**

The gap-block in question and the sequences it affects are shaded in cyan.

The red shaded rectangle encloses "semi-column" $c'$.

The magenta and yellow shaded rectangles enclose "semi-columns" $c''_{bf}$ and $c''_{bf}$, respectively.

The cyan branch phylogenetically delimits the gap-block.

**(A)**

```
          abcdefghijk
Tree
      1 NNNNNNNNNNN
nA
      2 N------NNNN
      3 NN---NNNNNN
```

**(B)**

```
          abcdefghijk
Tree
      1 NNNNNNNNNNN
nA
      2 N------NNNN
      3 NNNN---NNN
```

**(C)**

```
          abcdefghij
Tree
      1 N---NNNNNN
nA
      2 N---NNNNNN
      3 NNNN--NNNN
```

**(D)**

```
          abcdefghijk
Tree
      1 NNNNNNNNNNN
nA
      2 N------NNNN
      3 NNNN---NNNN
```

**(E)**

```
          abcdefghij
Tree
      1 N---NNNNNN
nA
      2 N---NNNNNN
      3 NNNNN--NNN
```

**Figure 7. Typical examples of gap-configurations of 3 sequences connected via 3-OTU tree.**

**A.** A short gap in the 3rd OTU (labeled "3") is nested in a long gap in the 2nd OUT ("2").

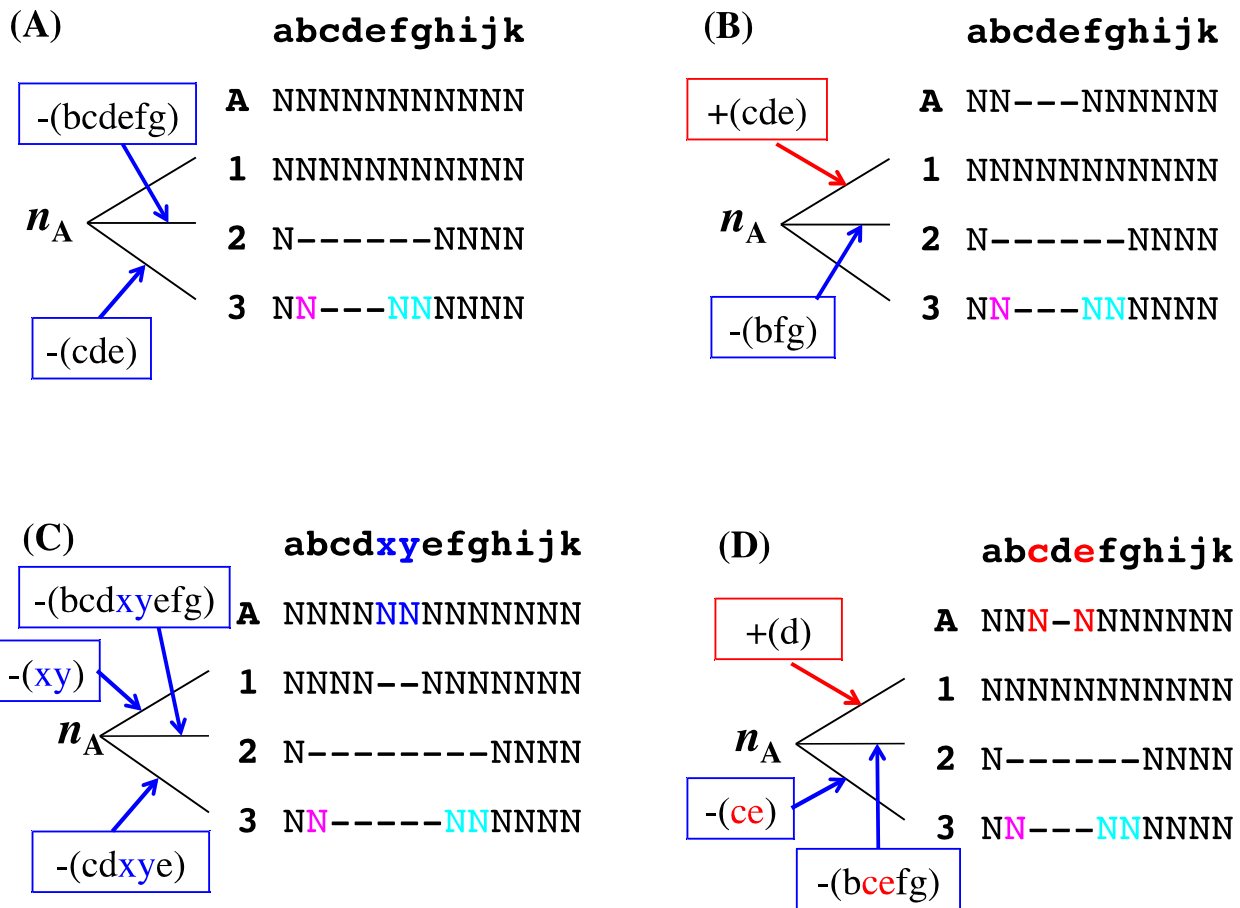**B.** Gaps in two OTUs overlap each other in a non-nested manner.

**C.** Two non-homologous sequence-blocks neighboring each other.

**D.** This pattern is considered as topologically the same as in panel A but not as in panel B, because of the (parsimonious) indel histories that can create the patterns.

**E.** This pattern is considered as topologically different from that in panel C.

In each panel, the " $n_A$ " denotes the "most recent common ancestor (MRCA)" of the three OTUs, and the lower-case letters above the alignment represent the ancestry indices of the sites.

**(A)**

```
                      abcdefghijk

  ┌──────────┐    A  NNNNNNNNNNN
  │ -(bcdefg)│
  └──────────┘    1  NNNNNNNNNNN
          ↘
  n_A  ←──────    2  N------NNNN
          ↗
  ┌──────────┐    3  NN---NNNNNN
  │  -(cde)  │
  └──────────┘
```

**(B)**

```
                      abcdefghijk

  ┌──────────┐    A  NN---NNNNNN
  │  +(cde)  │
  └──────────┘    1  NNNNNNNNNNN
          ↘
  n_A  ←──────    2  N------NNNN
          ↖
  ┌──────────┐    3  NN---NNNNNN
  │  -(bfg)  │
  └──────────┘
```

**(C)**

```
                      abcdxyefghijk

┌────────────┐   A  NNNNNNNNNNNNNN
│ -(bcdxyefg)│
└────────────┘   1  NNNN--NNNNNNN
┌────────┐ ↘
│ -(xy)  │       2  N--------NNNN
└────────┘ ↘
  n_A  ←──────   3  NN-----NNNNNN
          ↗
┌──────────┐
│ -(cdxye) │
└──────────┘
```

**(D)**

```
                      abcdefghijk

  ┌──────────┐    A  NNN-NNNNNNN
  │   +(d)   │
  └──────────┘    1  NNNNNNNNNNN
          ↘
  n_A  ←──────    2  N------NNNN
       ↖  ↑
┌────────┐      3  NN---NNNNNN
│ -(ce)  │
└────────┘
     ┌──────────┐
     │ -(bcefg) │
     └──────────┘
```

**Figure 8. Indel histories that can create pattern A in Figure 7.**

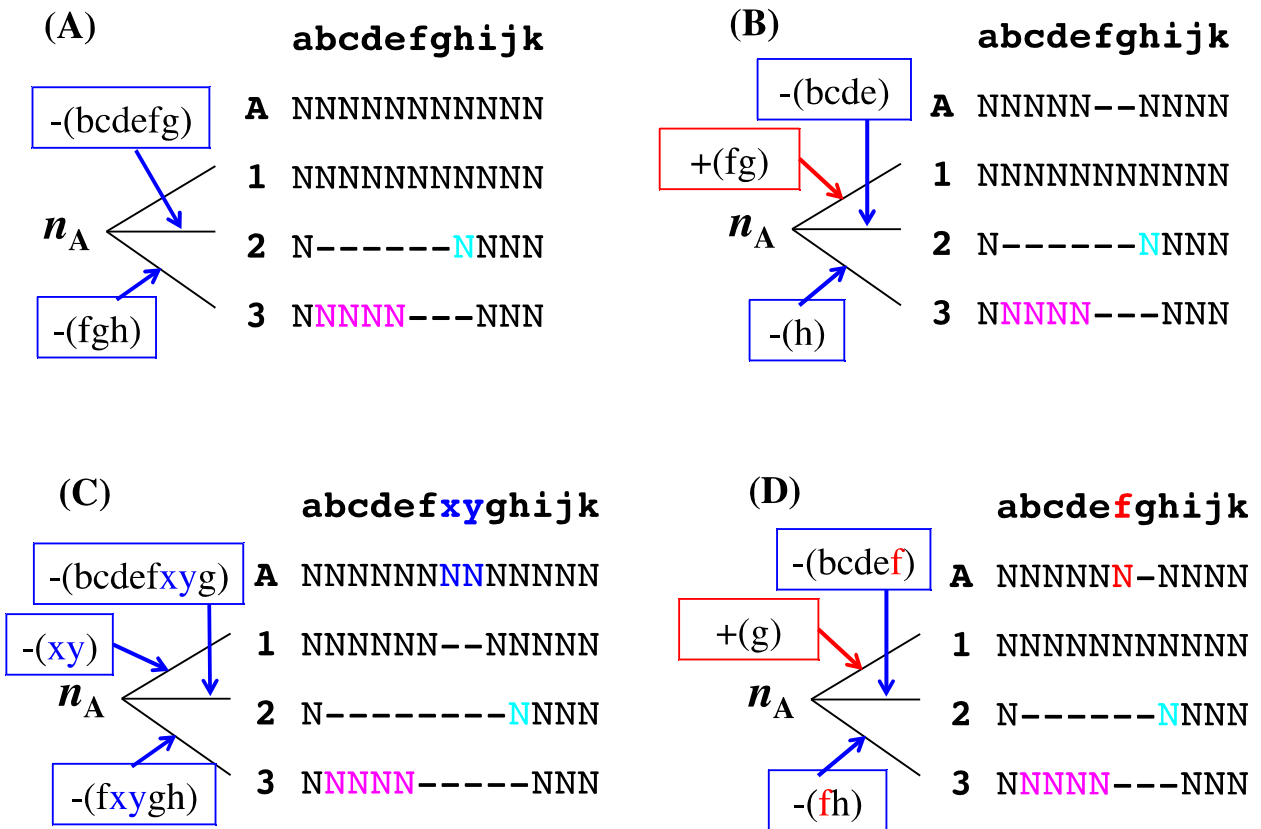Panels **A** and **B** show the parsimonious histories.

Panels **C** and **D** show some next-to-parsimonious histories.

The 'A' on the top-left corner of the alignment indicates the ancestral sequence (at node $n_A$).

"-(abc)" in a blue rectangle shows that the sites a, b and c were deleted along the branch it points, and "+(def)" in a red rectangle shows that the sites d, e and f were inserted.

Some residues and ancestries were colored in order to facilitate the comparisons among the histories.

NOTE: Here we omitted all next-to-parsimonious histories in each of which two or more indels occur along a branch.

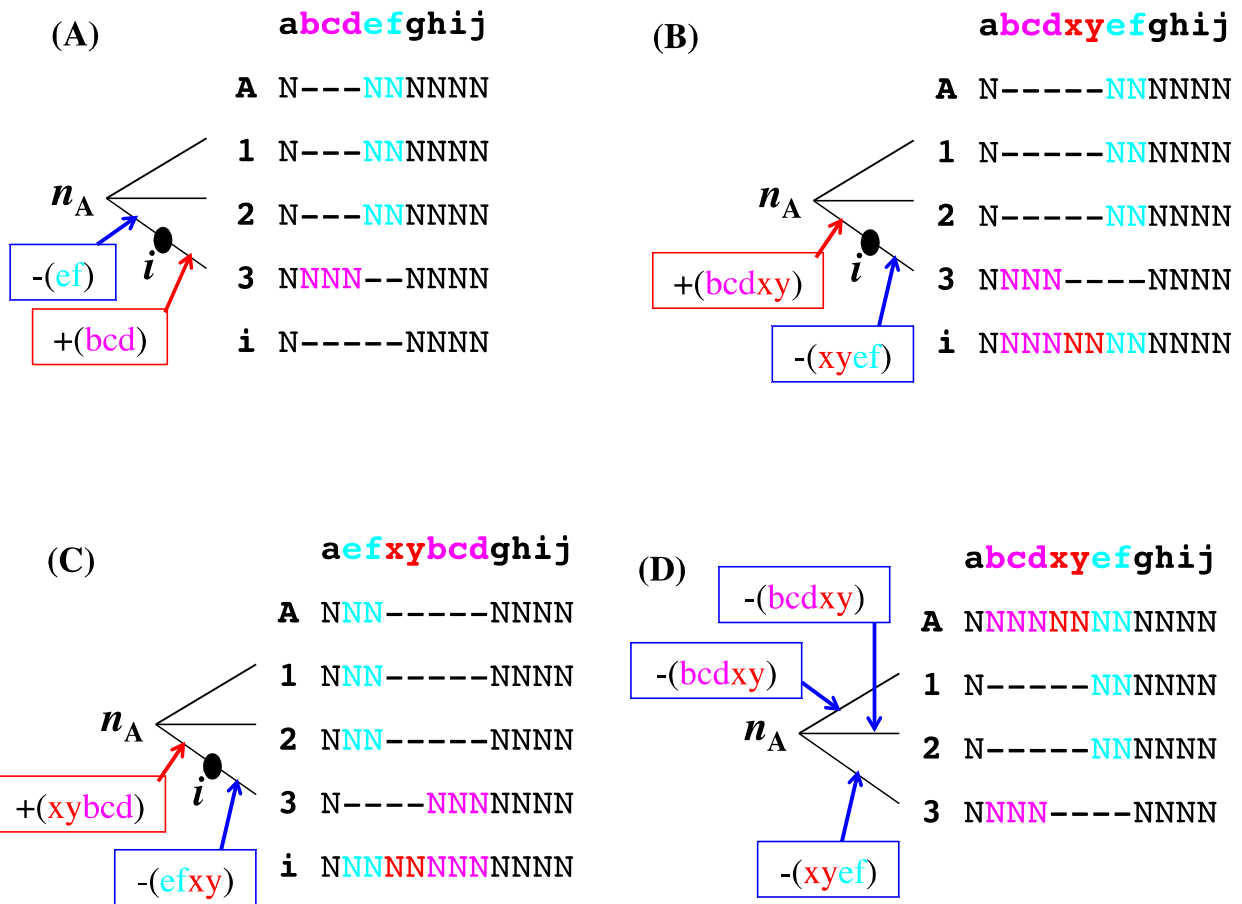**Figure 9. Indel histories that can create patten B in Figure 7.**

Panel **A** shows the parsimonious indel history.

Panels **B, C** and **D** show some next-to-parsimonious histories.

Notations are the same as that for Figure 8.

NOTE: Here, again, we omitted the next-to-parsimonious histories in each of which two indels occur along a branch.

NOTE2: Panel B is actually the history created by concatenating the column-wise Dollo parsimonious indel histories (Farris 1977; Ezawa, Graur and Landan, Part III).

**(A)**

`abcdefghij`

```
A  N---NNNNNN
1  N---NNNNNN
2  N---NNNNNN
3  NNNN--NNNN
i  N-----NNNN
```

$n_A$

-(ef)  $i$  +(bcd)

**(B)**

`abcdxyefghij`

```
A  N-----NNNNNN
1  N-----NNNNNN
2  N-----NNNNNN
3  NNNN----NNNN
i  NNNNNNNNNNNN
```

$n_A$

+(bcdxy)  $i$  -(xyef)

**(C)**

`aefxybcdghij`

```
A  NNN-----NNNN
1  NNN-----NNNN
2  NNN-----NNNN
3  N----NNNNNNN
i  NNNNNNNNNNNN
```

$n_A$

+(xybcd)  $i$  -(efxy)

**(D)**

`abcdxyefghij`

-(bcdxy)

-(bcdxy)

$n_A$

```
A  NNNNNNNNNNNN
1  N-----NNNNNN
2  N-----NNNNNN
3  NNNN----NNNN
```

-(xyef)

**Figure 10. Indel histories that can create pattern C in Figure 7.**

Panels **A**, **B**, and **C** show the three types of parsimonious indel histories.

Panel **D** shows a next-to-parsimonious indel history, which was derived from panel B via a "branching" operation (Ezawa, Graur and Landan, Part III).

(Note that a different history could be derived from panel C in a similar way.)

The notations are basically the same as in Figure 8. In addition, the sequence labelled "*i*" at the bottom of each MSA (in panels A, B and C) is the "intermediate" state at the point marked with a solid circle.

NOTE: The MSA of the extant sequences in panel C is equivalent to that in panel A (and in panel B).

NOTE2: We omitted next-to-parsimonious histories in each of which 3 indels occur along a branch.
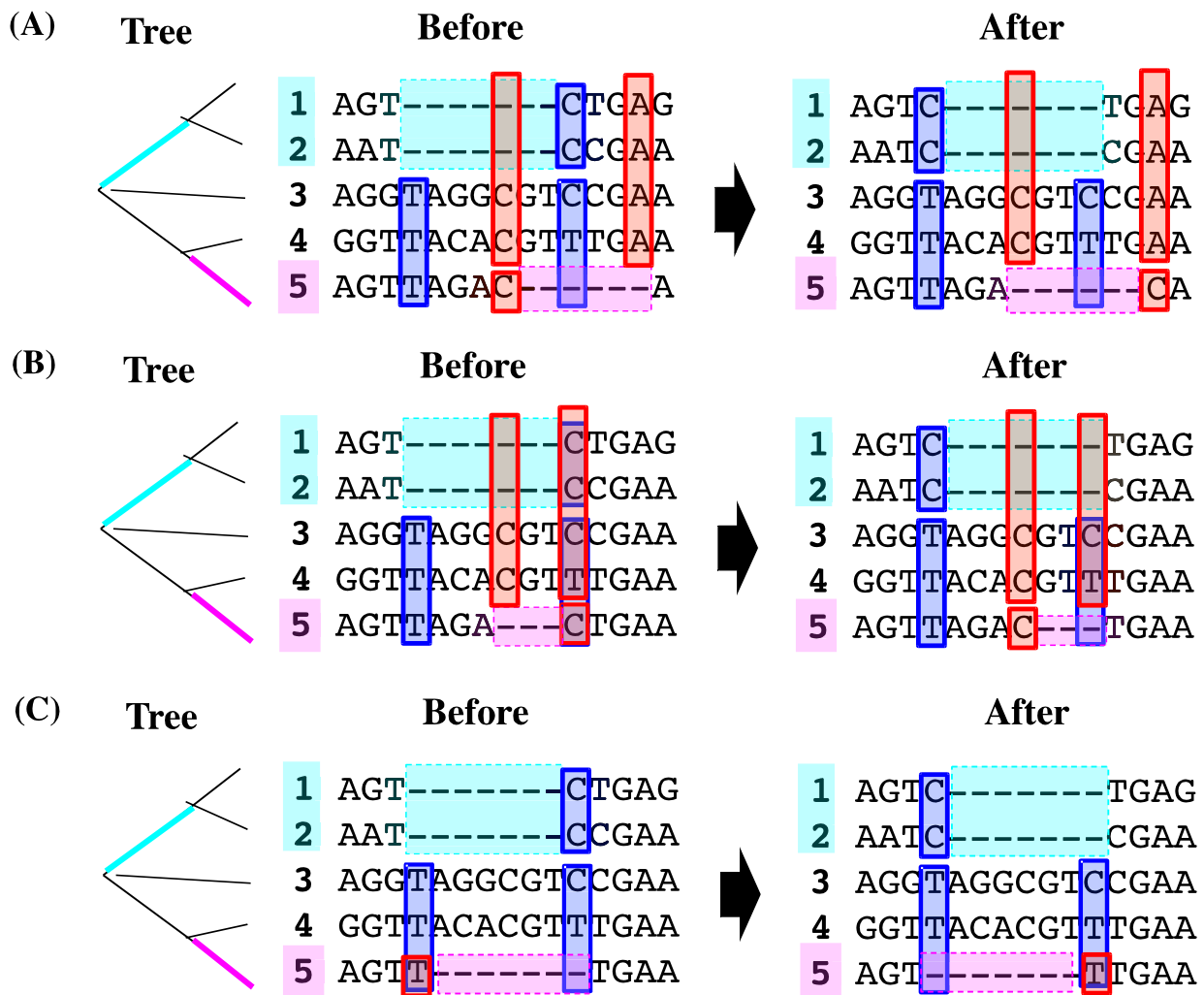
**Figure 11. Indel histories that can create pattern in Figure 7E.**

Panels **A** and **B** show the two parsimonious histories.

Panel **C** shows a next-to-parsimonious history. (It was obtained by applying a "branching" operation to panel B.)

The notations are basically the same as in Figure 8.

NOTE: We omitted those next-to-parsimonious histories each of which requires more than 2 indels in a branch.

**Figure 12. Effects of simultaneous shifts of two isolated gap-blocks on substitution component of MSA probability.**
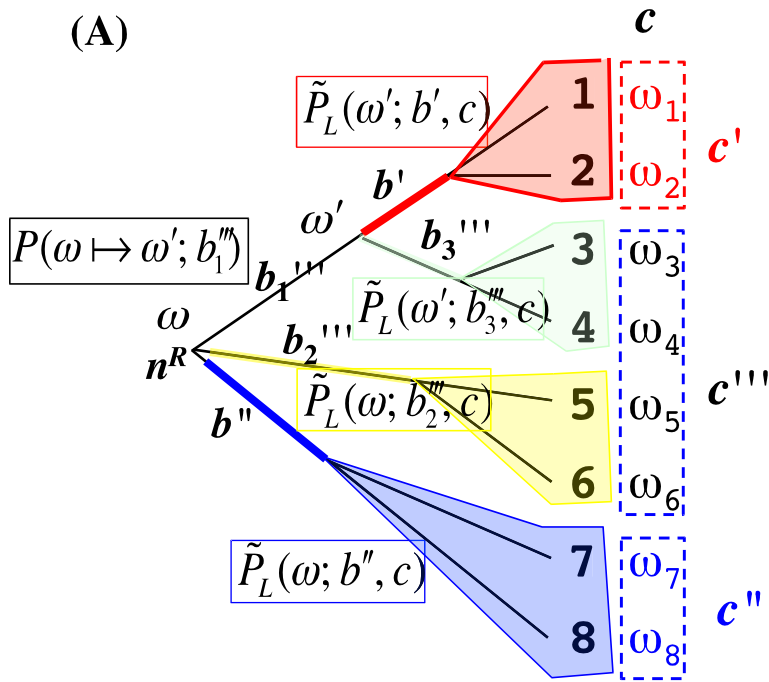
**A.** The shifts affect different columns independently.

**B.** The shifts affect the same column simultaneously.

**C.** The shifts affect the same pair of columns simultaneously.

In each column, shifts of the gap-blocks shaded in cyan and magenta affect the (semi-)columns shaded in blue and red, respectively.

The colored branches in the tree phylogenetically delimits the gap-blocks shaded in the same color.
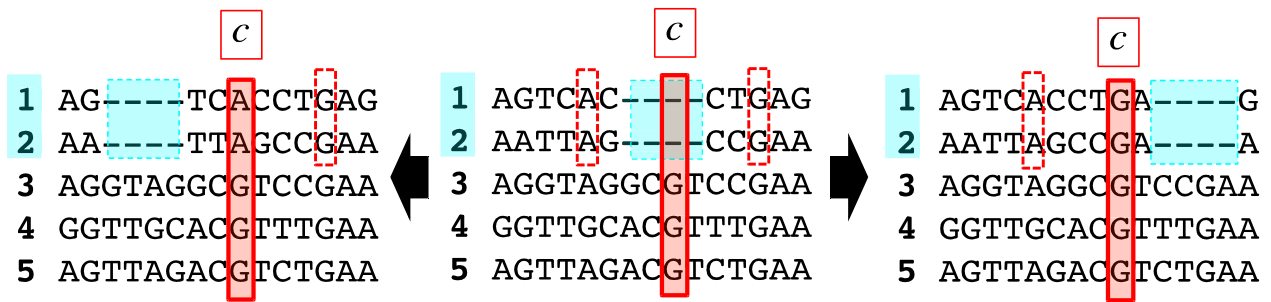
**(A)**

**Figure 13. Model situation for examining whether effects of shifts of two "isolated" gap-blocks are nearly independent or not.**

Each portion of the tree (shaded in respective color) is assigned a building-block probability (enclosed by a rectangle of the same color).

The column ( $c$ ) was divided into three parts: $c'$, $c''$ and $c'''$.

Branches $b'$ and $b''$ delimit $c'$ and $c''$, respectively. And branches $b_i'''$ ( $i=1,2,3$) are important factors determining $c'''$.
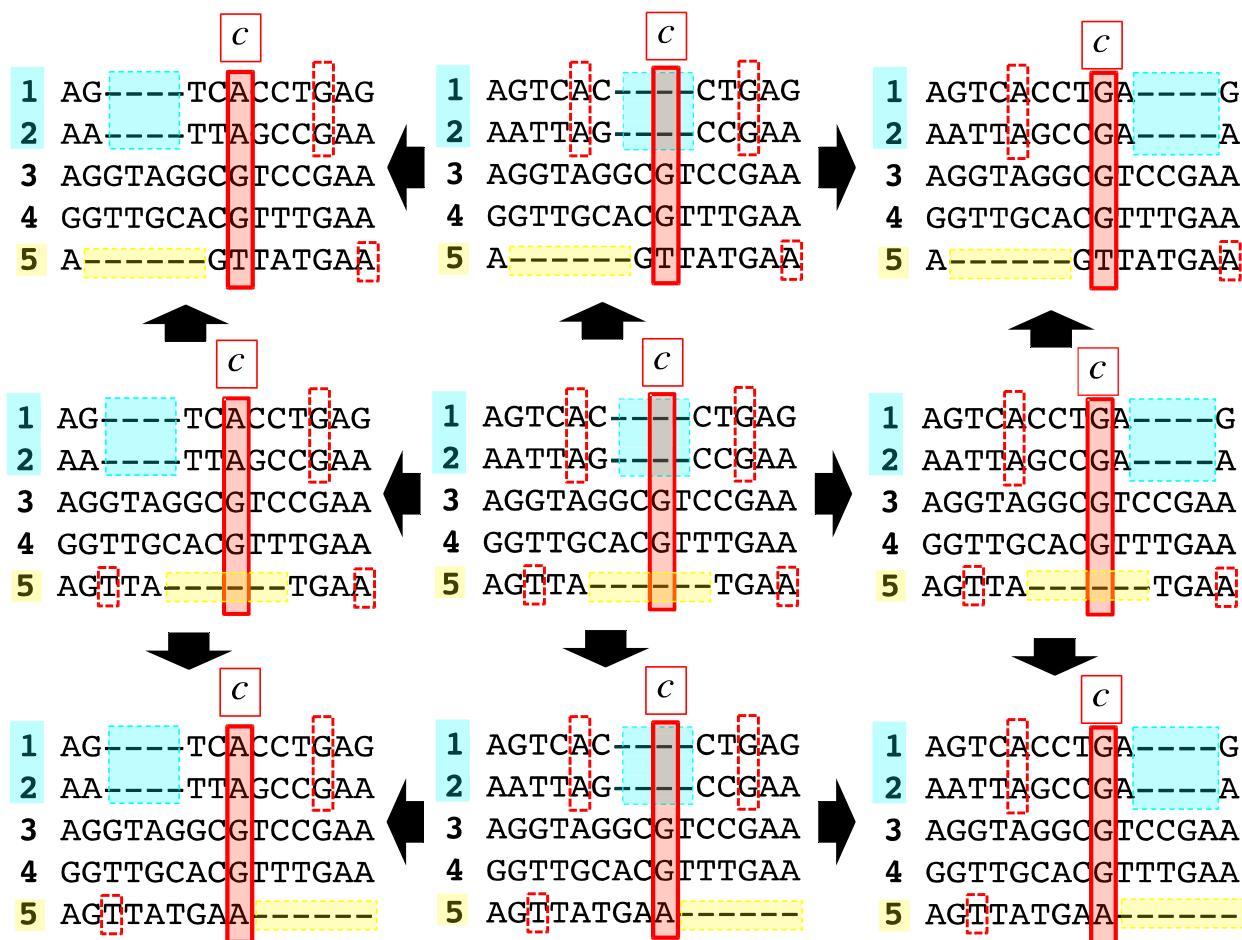
The $\omega$ and $\omega'$ at both ends of branch $b_1'''$ are the residue states over which the probabilities will be summed.

**Figure 14. Possible residue configurations of each column caused by shifts of single gap-block.**
The shifts of a single gap-block (cyan-shaded) give rise to at most 3 possible configurations in each column (red-shaded), because they are "rigid" moves of the gap-block.
The red dotted rectangles enclose the original "ingredients" of the columns that resulted from the shifts.

**Figure 15. Possible residue configurations of each column caused by double-shifts of two gap-blocks.**

The double-shifts of two gap-blocks (cyan- and yellow-shaded) give rise to at most 9 possible configurations in each column (red-shaded),

The red dotted rectangles enclose the original "ingredients" of the columns that resulted from the shifts.