

# みんなのポジトリ

国立民族学博物館学術情報リポジトリ National Museum of Ethnology

## Description of Structure of the Folktale : Using a Bioinformatics Multiple Alignment Program

メタデータ	言語: eng 出版者: 公開日: 2009-04-28 キーワード (Ja): キーワード (En): 作成者: 小田, 淳一 メールアドレス: 所属:
URL	<a href="https://doi.org/10.15021/00002828">https://doi.org/10.15021/00002828</a>

## **Description of Structure of the Folktale: Using a Bioinformatics Multiple Alignment Program**

JUN'ICHI ODA

*Institute for the Study of Languages  
and Cultures of Asia and Africa  
Tokyo University of Foreign Studies*

### **INTRODUCTION**

In a recent paper<sup>1)</sup> I produced a rough sketch of an attempt to apply genome informatics (bioinformatics that treats genetic information as a whole) analysis method developed from the field of molecular biology to the study of the folktale. The objective of this attempt was to scientifically reconsider the conventional analogies<sup>2)</sup> of the various concepts used for both analysis of various narrative genres and biology at a concrete analysis level. At that time, I raised the following items as the preliminary stages to be performed before serious application could be carried out.

- (1) A description of folkloristic text based on Fillmore's case grammar and Jackendoff's semantic structure model.
- (2) Setup of a database containing the fundamental materials required for the seme-dictionary (a dictionary with the items that appear on the surface layer of folkloristic text set as headings, that uses articulated semes as reference data for each item).
- (3) Application of an alignment program used for genome informatics to the study of the folktale.

An accurate logic model construction exists for item (1), and currently contributions in related areas are under investigation. With regard to the database construction in item (2), data input has been in progress for several years, and official announcement of an experimental system that will allow searching using the web site of the Institute for the Study of Languages and Cultures of Asia and Africa (Tokyo University of Foreign Studies) will be made in the very near future. This paper is mainly concerned with item (3) - investigation of the effectiveness of the alignment programs used in bioinformatics in analysis of folktales by actual application of the software. The folktale sequence structure used as data for the analysis was provisionally based on Propp's<sup>3)</sup> 'sequence of functions' model (although the actual model used was a modification of this model by Greimas<sup>4)</sup>).

The reason that I presumed to use the Propp model (as of the present a number of problems have been pointed out) is that it is considered that Propp's Functions model produces data suitable for processing.<sup>5)</sup> In addition, it was Propp who first made a 'scientific' proposal to use biological analysis methods in folktale research, and more than anything else, I chose his model to pay him homage for this.

## ALIGNMENT

With regard to the alignment in bioinformatics, each element of the nucleic acids (4 types) that make up DNA, and each element of the amino acids (20 types) that make up proteins is represented as a character, and hyphens (gaps) are inserted at appropriate locations in the sequences in such a way that the common characters in the sequences are aligned in vertical columns. The following example<sup>6)</sup> shows the results of processing on five input protein sequences. An asterisk under a column in the output sequences indicates a maximum similarity value (i.e. all the characters in the column are the same), while a plus sign indicates that the measured similarity exceeds some specified value.

```
.. VKQIESKTAF QEALDAAGDK LVVVD FSATW CGPCKMIKPF FHSLS EKYSN VIFLEV DVDD..
.. SDKIIHLTDD SFDTDVLKAD GAILVDFWAE WCGPCKMIAP ILDEI ADEYQ GKLTVA KLNI..
.. ALEEEEDNVLV LKKS NFEEAL AAHKYLLVEF YAPWCGHCKA LAPEYAKRAA KLKAEGSEIR..
.. DSAVVKLATD SFNEYIQSHD LVLA EFFAPW CGHCKNMAPE YVKA AETLVE KNITLA QIDC..
.. METAMSKGVI TITDAEFESE VLKAEQPVLV YFWASWCGPC QLMSPLINLA ANTYS DRLKV..
```

**Fig. 1** Sequence input

```
-----VKQIE-SKTAFQEALDAAGDKLVVVD FSATWCGPCKMIKPF FHSLS EKYSN VIFLEV-DVDD
-----SDKIIHLTDDSFDTDV-LKADGAILVD FWAEWCGPCKMIAPILDEI ADEYQGK--LTVAKLNI
ALEE--EDNVLVLKKS NFEEAL--AAHKYLLVEF YAPWCGHCKA LAPEYAKRAAKLKAEG-SEI-R---
-----DSAVVKLATDSFNEYI--QSHDLVLA EFAPWCGHCKNMAPEYVKA AETLVEKN-ITLAQIDC
-METAMSKGVI TITDAEFESV-LKAEQPVLVY FWA SWCGPCQLMSPLINLAANTYS DR--LKV-----
          *          ++ * * * * * * *          +
```

**Fig. 2** Alignment output

When a pattern made up of vertically aligned characters in the aligned sequences is determined to bestow a characteristic, this pattern is called a *motif* for the local similarity. When the shaded section in the above sequence is constructed as a pattern, the following description is obtained.

F[SWYF]A[TEPS]WCG[HP]C[KQ]

The characters in the square brackets are the amino acids permitted for that column. When the similarity relationship between sequences for the above description is analyzed, it provides useful data for the motif. In addition, analysis of various relationships such as specific character insertion and deletion provides a lot of information about the evolution process.

The alignment program used for this paper was CLUSTAL W, Version 1.6 (March 1996)<sup>7)</sup>, developed by Julie Thompson, Toby Gibson, and Des Higgins. Another program called TREEVIEW<sup>8)</sup> was used to display and print the phylogenies from the tree file that outputs the distances between the sequences obtained from the alignment process. CLUSTAL W was run on a Hitachi 9000V Series computer (V715S-Tiny) running HP-UX 10.20, and TREEVIEW was run on a Power Macintosh 6200/75. The input and output formats required for the programs are described in the following chapter.

## DATA

The data created for this paper is based on Propp's Function model. Propp proposed the following Function concept based on analysis of Russian fairy-tales (*Narodnye russkie skazki* by A. N. Afanas'ev).

Function is understood as an act of a character, defined from the point of view of its significance for the course of the action.<sup>9)</sup> Propp established the following 31 functions.

**Table 1** 31 Functions by Propp

$\beta$ :	absentation
$\gamma$ :	interdiction
$\delta$ :	violation
$\epsilon$ :	reconnaissance
$\xi$ :	delivery
$\eta$ :	trickery
$\theta$ :	complicity
A:	villainy
a:	lack
B:	mediation, the connective incident
C:	beginning counteraction
$\uparrow$ :	departure
D:	the first function of the donor
E:	the hero's reaction
F:	provision or receipt of magical agent
G:	spatial transference between two kingdoms, guidance
H:	struggle
J:	marking
I:	victory

K : the initial misfortune or lack is liquidated  
 ↓ : return  
 Pr : pursuit, chase  
 Rs : rescue  
 O : unrecognized arrival  
 L : unfounded claims  
 M : the difficult task  
 N : solution  
 Q : recognition  
 Ex : exposure  
 T : transfiguration  
 U : punishment  
 W : wedding

Greimas focused on the Pairs that make up these functions, and based on Propp's catalog, reduced the list to the following 20 functions.

**Table 2** 20 Functions by Greimas

- 1: absention ( $\beta$ )
- 2: interdiction ( $\gamma$ ) vs violation ( $\delta$ )
- 3: reconnaissance ( $\epsilon$ ) vs delivery ( $\xi$ )
- 4: trickery ( $\eta$ ) vs complicity ( $\theta$ )
- 5: villainy (A) vs lack (a)
- 6: mediation, the connective incident (B) vs beginning counteraction (C)
- 7: departure ( $\uparrow$ )
- 8: the first function of the donor (D) vs the hero's reaction (E)
- 9: provision or receipt of magical agent (F)
- 10: spatial transference between two kingdoms, guidance (G)
- 11: struggle (H) vs victory (I)
- 12: marking (J)
- 13: the initial misfortune or lack is liquidated (K)
- 14: return ( $\downarrow$ )
- 15: pursuit, chase (Pr) vs rescue (Rs)
- 16: unrecognized arrival (O)
- 17: the difficult task (M) vs solution (N)
- 18: recognition (Q)
- 19: exposure (Ex) vs transfiguration (T)
- 20: punishment (U) vs wedding (W)

The data presented in this paper was produced by using the Greimas catalog. The materials used for data creation were Propp's actual results of analysis of the fairy tales. These results are from analysis performed on 45 fairy tales<sup>10)</sup> and from tales within the texts. For analysis, Propp does not treat the tale as a whole, but

rather uses the *moves* that comprise the tale as the unit for analysis. In addition, the following type of simplification was performed for presentation of the analysis results.

- 1) The first seven functions ( $\beta, \gamma, \delta, \varepsilon, \zeta, \eta, \theta$ ) which are regarded as the preparatory part of the tale are not shown.
- 2) Trebling (e.g. three tasks, three years of service etc.) is omitted.

For technical reasons relating to data creation, the following simplifications were also used in this paper.

- 1) Although one function (= *genus*) can be sub-divided into multiple species (e.g. A<sup>1</sup> (kidnapping of a person) and A<sup>9</sup> (expulsion) are species of Function A (villainy)), in this paper only functions (= *genera*) are handled (and not species).
- 2) The negative result of a function (e.g. the negative result  $\bar{F}$  of Function F (receipt of magical agent) indicates that nothing is given to the hero) is also treated as the function.
- 3) The considerably complicated 'sequence of functions' that is seen as problematic with respect to both the analysis results and Propp's scheme has been omitted from the data.<sup>11)</sup>

After applying these simplifications, the sequences resulting from the Propp analysis are as follows. The numerals at the start show the enumeration of tales according to the collection of Afanas'ev, and the characters after the hyphen are the moves given by the Propp analysis data.

**Table 3** Schemes by Propp

93-1	AB↑DEF	106-2	AB↑DE↓
93-2	aB↑DEF	108	A↑DE↓P r R s
95-1	AB↑DEF↓	112	A↑DE↓P r R s
95-2	aBC↑DEF↓	113	ABC↑DEFG↓P r DER s
97	AB↑DEF↓	114	A↑DEF↓P r R s OQW
98-1	AB↑DEF↓	115-1	A↑EFK
98-2	aBC↑DEF↓	115-2	ABC↑DEFI W
99-1	AB↑DEFT↓	125-1	DEFA↑O
99-2	AB↑DEF↓	125-2	FABC↑HJ I ↓P r R s LQE x UW
100	ADEFMNW	126-1	ABC↑K↓
101-1	AFMNW	126-2	DEFABC↑FHJ IOQW
101-2	ABC↑KMNE x UW	128-1	aBC↑DEFGJK↓W
102-1	AB↑DEF↓	128-2	AC↑G↓K
102-2	aB↑DEF↓	131	ABC↑HIK↓W
104-1	F aBC↑DEF↓U	132-1	DEABC↑K↓
104-2	aFNMNW	132-2	aCF↑DEFGI ↓
105-1	ABCF↓K↓	132-3	AC↑DEF↓GMNOQE x UW
105-2	aC↑DEF↓P r R s	133-1	ABC↑DEF
106-1	AB↑K↓	133-2	BC↑DEHIK↓

135	AC↑FHIK↓PrRs	155-3	↑SGWF
136-1	aBC↑DEKF↓	155-4	AC↑HIJLQExUW
136-2	ABC↑KF↓	155-5	ABC↑GKU
136-3	ACHIK↓PrRs	156-1	ABC↑DEFGK
137-1	↑DEFACHI↓PrRs	156-2	aBC↑DEFGK↓UW
137-2	aBCF↑FMNK↓UW	156-3	FAC↑GOMNUW
138-1	aBC↑DEFHIK↓PrRs	161-1	ABC↑HIK↓
138-2	AC↑DEFGK↓W	161-2	aC↑FGHIKF
138-3	aBC↑MNK↓	161-3	aBC↑DEFGIK
140	ABC↑DEGFIK↓W	161-4	aC↑GHIW
141-1	AF↑DEFGIK↓	162	↑DEFABCHIK↓W
141-2	AC↑FGK↓UW	163	aFC↑GTW↓
143-1	ACF↑HIK↓	164-1	AC↑K↓
143-2	A↑FHIK↓W	164-2	aFC↑TGIW↓
144	aMC↑DEFFGMNTW	166-1	DEFA↑GW
145	aBC↑GK↓W	166-2	AGTQ
148	ABC↑HIK↓	167-1	aDEFKTUW
149	ABC↑HIK↓	167-2	ABC↑G
150-1	↑DEaBC↑K↓	167-3	aBC↑MNW↓
150-2	aBC↑HIK↓PrRsUW	247	A↑GMNW↓Q
151	aBC↑HIK↓	244	AFExU
152	A↑HIK↓	139	↑DEFGABC↑HIKJW↓
153	ABCFHIUW↓	139-2	ADEFGOLMNQTUQW
154-1	↑aDEFG↓	123-1	DFEAAO
154-2	AC↑FHIUW	123-2	FABC↑FHIKJ↓L
155-1	aBC↑DEF↓	123-3	aBC↑FK↓ExQUW
155-2	aBC↑DEF↓	123-4	ABC↑FK↓ExQUW

As a preliminary step to processing these schemes using the CLUSTAL W software, they were first converted to the Greimas Functions. In this paper, even if one of the functions in the pairs established by Greimas was present, that function was treated as the relevant pair. This is no more than a provisional operation, because the objective of this paper is to investigate the applicability of the alignment program at a technical level. As a result, using the Greimas model, the Propp schemes are represented as 16 functions.<sup>12)</sup>

The following table shows the Propp functions reduced using the Greimas model to 16 functions, and the corresponding amino acid codes (that are recognized by the CLUSTAL W software).<sup>13)</sup>

**Table 4** Conversion table

Functions by Greimas	amino acid codes
A vs a	→ A
B vs C	→ C
↑	→ D
D vs E	→ E
F	→ F
G	→ G
H vs I	→ H
J	→ I
K	→ K
↓	→ L
Pr vs Rs	→ P
O	→ N
M vs N	→ M
Q	→ Q
Ex vs T	→ T
U vs W	→ W

Using this conversion, when the Propp schemes are used as sequence input for the CLUSTAL W program, the following format results.<sup>14)</sup> The “P” before the numerals for each sequence (move of the tale) stands for Propp.

**Table 5** Sequence input used by CLUATAL W

>P093-1	>P138-1
ACDEF	ACDMKL
>P093-2	>P140
ACFDEF	ACDEGFHKLW
>P093-3	>P141-1
ADPHW	AFDEFGHKL
>P095-1	>P141-1
ACDEFL	ACDFGKLW
>P095-2	>P143-1
ACDEFL	ACFDHKL
>P097	>P143-2
ACDEFL	ADFHKLW
>P098-1	>P144
ACDEFL	AMCDEFFGMTW
>P098-2	>P145
ACDEFL	ACDGKLW
>P099-1	>P148
ACDEFTL	ACDHKL



>P099-2	>P149
ACDEFL	ACDHLK
>P100	>P150-1
AEFMW	DEACDKL
>P101-1	>P150-2
AFMW	ACDHKL PW
>P101-2	>P151
ACDKMTW	ACDHLK
>P102-1	>P152
ACDEFL	ADHLK
>P102-2	>P153
ACDEFL	ACFHWL
>P104-1	>P154-1
FACDEFLW	DAEFGL
>P104-2	>P154-2
AFMMW	ACDFHW
>P105-1	>P155-1
ACFDKL	ACDEFL
>P105-2	>P155-2
ACDEFLP	ACDEFL
>P106-1	>P155-3
ACDKL	DGWF
>P106-2	>P155-4
ACDEL	ACDHIQ TW
>P108	>P155-5
ADELP	ACDGKW
>P112	>P156-1
ADELP	ACDEFGK
>P113	>P156-2
ACDEFGLPEFP	ACDEFGLW
>P114	>P156-3
ADEFLPNQW	FACDGNMW
>P115-1	>P161-1
ADEFK	ACDHLK
>P115-2	>P161-2
ACDEFHW	ACDFGHKF
>P125-1	>P161-3
EFADN	ACDEFGHK
>P125-2	>P161-4
FACDHIHLPQ TW	ACDGHW
>P126-1	>P162
ACDKL	DEFACHKLW
>P126-2	>P163

EFACDFHIHNQW	AFCDGTWL
>P128-1	>P164-1
ACDEFGIKLW	ACDKL
>P128-2	>P164-2
ACDGLK	AFCDTGHWL
>P131	>P166-1
ACDHKLW	EFADGW
>P132-1	>P166-2
EACDKL	AGTQ
>P132-2	>P167-1
ACFDEFGHL	AEFKTW
>P132-3	>P167-2
ACDEFGLGMNQTW	ACDG
>P133-1	>P167-3
ACDEF	ACDMWL
>P133-2	>P247
CDEHKL	ADGMWLQ
>P135	>P244
ACDFHKLP	AFTW
>P136-1	>P139
ACDEKFL	DEFGACDHIWL
>P136-2	>P139-2
ACDKFL	AEFGNMQTWQW
>P136-3	>P123-1
ACHKLP	EFEAAN
>P137-1	>P123-2
DEFACHLP	FACDFHKIL
>P137-2	>P123-3
ACFLFMKLW	ACDFKLTQW
>P138-1	>P123-4
ACDEFHKLP	ACDFKLTQW
>P138-2	
ACDEFGKLW	

## DATA PROCESSING AND DISCUSSION

This chapter deals with the processing procedure for sequence input using CLUSTAL W and TREEVIEW, and the analysis of the alignment output. CLUSTAL W outputs two files after processing the sequence input. These are SOMETHING.ALN (the file that contains the alignment) and SOMETHING.DND (for the dendrogram). In general, alignment transformations are performed repeatedly to approach a sensitive result. There was dispersion in the similarity between sequences in the initial results of the alignment performed on the sequence

input given above 'as a whole', and the sensitivity of the result was quite low. In this paper, the procedure given below was used to group suitable numbers of relatively high-similarity sequences from the alignment output as a whole, and perform alignment on these groups again.

① Alignment as a whole

Alignment was performed on the 93 sequences 'as a whole' until the results converged (three iterations), and the results displayed in groups of relatively high similarity.

```

P 1 5 5 - 3      ----- D G W F -----
P 1 6 6 - 1      ----- E F A D G W -----
P 1 2 5 - 1      ----- E F A D N -----
P 1 3 7 - 1      ----- D E F A C H L P -----
P 1 0 5 - 2      ----- A C D E F L P -----
P 0 9 5 - 1      ----- A C D E F L -----
P 0 9 5 - 2      ----- A C D E F L -----
P 0 9 7          ----- A C D E F L -----
P 1 3 2 - 3      ----- A C D E F L G M N Q T W -----
P 1 3 6 - 2      ----- A C D K F L -----
P 0 9 8 - 2      ----- A C D E F L -----
P 1 0 2 - 1      ----- A C D E F L -----
P 1 0 4 - 1      ----- F A C D E F L W -----
P 1 5 5 - 2      ----- A C D E F L -----
P 1 5 5 - 1      ----- A C D E F L -----
P 1 0 2 - 2      ----- A C D E F L -----
P 0 9 8 - 1      ----- A C D E F L -----
P 1 1 5 - 2      ----- A C D E F H W -----
P 1 3 8 - 1      ----- A C D E F H K L P -----
P 1 3 3 - 1      ----- A C D E F -----
P 1 2 8 - 1      ----- A C D E F G I K L W -----
P 1 6 1 - 3      ----- A C D E F G H K -----
P 0 9 3 - 1      ----- A C D E F -----
P 1 3 8 - 2      ----- A C D E F - G K L W -----
P 1 5 6 - 1      ----- A C D E F - G K -----
P 1 4 4          ----- A M C D E F F G M T W -----
P 0 9 9 - 1      ----- A C D E F T L -----
P 1 3 6 - 1      ----- A C D E K F L -----
P 1 0 6 - 2      ----- A C D E L -----
P 1 4 1 - 2      ----- A C D F G -- K L W -----
P 1 4 5          ----- A C D G -- K L W -----
P 1 4 0          ----- A C D E G F H K L W -----
P 1 0 5 - 1      ----- A C F D K L -----

```

P 1 4 3 - 2	-----ADFHKLW-----
P 1 5 3	-----ACFHWL-----
P 1 3 7 - 2	-----ACFLFMKLW-----
P 1 3 2 - 1	-----EACDKL-----
P 1 5 0 - 1	-----DEACDKL-----
P 1 0 6 - 1	-----ACDKL-----
P 1 2 6 - 1	-----ACDKL-----
P 1 6 4 - 1	-----ACDKL-----
P 1 0 1 - 2	-----ACDKMTW-----
P 1 5 5 - 5	-----ACDGKW-----
P 1 6 1 - 4	-----ACDGHW-----
P 1 5 6 - 3	-----FACDGNMW-----
P 1 6 7 - 2	-----ACDG-----
P 1 2 8 - 2	-----ACDGLK-----
P 1 2 5 - 2	-----FACDHIHLPQTW-----
P 1 5 5 - 4	-----ACDHIQTW-----
P 1 5 4 - 2	-----ACDFHW-----
P 1 2 3 - 2	-----FACDFHKIL-----
P 1 2 6 - 2	-----EFACDFHIHNQW-----
P 1 3 3 - 2	-----CDEHKL-----
P 1 3 5	-----ACDFHKL P-----
P 1 3 6 - 3	-----ACHKL P-----
P 1 6 2	-----DEFACHKLW-----
P 1 4 3 - 1	-----ACFDHKL-----
P 1 5 2	-----ADHKL-----
P 1 6 1 - 2	-----ACDFGHKF-----
P 1 4 8	-----ACDHKL-----
P 1 5 0 - 2	-----ACDHKL PW-----
P 1 4 9	-----ACDHKL-----
P 1 5 1	-----ACDHKL-----
P 1 6 1 - 1	-----ACDHKL-----
P 1 3 1	-----ACDHKLW-----
P 1 3 9	--DEFGACDHKIWL-----
P 1 2 3 - 3	-----ACDFKLTQW-----
P 1 2 3 - 4	-----ACDFKLTQW-----
P 1 3 8 - 3	-----ACDMKL-----
P 1 6 7 - 3	-----ACDMKL-----
P 0 9 3 - 2	-----ACFDEF-----
P 1 3 2 - 2	-----ACFDEFGHL-----
P 1 4 1 - 1	-----AFDEFGHKL-----
P 1 0 8	-----ADELP-----
P 1 1 2	-----ADELP-----
P 1 1 4	-----ADEFLPNQW-----

```

P 1 1 5 - 1      -----A D E F K-----
P 0 9 3 - 3      -----A D P - H W-----
P 2 4 7          -----A D G - M W L Q-----
P 1 6 3          -----A F C D G - T W L-----
P 1 6 4 - 2      -----A F C D T G H W L-----
P 1 0 1 - 1      -----A F M W-----
P 2 4 4          -----A F T W-----
P 1 0 0          -----A E F M W-----
P 1 6 7 - 1      -----A E F K T W-----
P 1 0 4 - 2      -----A F M M W-----
P 1 6 6 - 2      -----A G T Q-----
P 1 5 4 - 1      -----D A E F G L-----
P 1 3 9 - 2      -----A E F G N M Q T W Q W-----
P 1 2 3 - 1      -----E F E A A N-----

```

**Fig. 3** Alignment as a whole

## ② Division into two groups

The above alignment output was broadly divided into two groups based on similarity. For the first group (Group A), 28 sequences were extracted from the 5th sequence in the above results (P105-2), and for the second group (Group B), 18 sequences were extracted from the 40th sequence (P132-1). When alignment was performed again on these two groups the following results were obtained (Group A converged after two iterations and Group B after one iteration).

```

P 1 2 8 - 1      -----A C D E F G I K L W-----
P 1 6 1 - 3      -----A C D E F G H K-----
P 0 9 3 - 1      -----A C D E F-----
P 1 5 5 - 2      -----A C D E F L-----
P 1 0 5 - 2      -----A C D E F L P-----
P 1 0 2 - 1      -----A C D E F L-----
P 0 9 7          -----A C D E F L-----
P 1 3 2 - 3      -----A C D E F L G M N Q T W-----
P 1 3 6 - 2      -----A C D K F L-----
P 1 0 2 - 2      -----A C D E F L-----
P 0 9 8 - 1      -----A C D E F L-----
P 0 9 5 - 2      -----A C D E F L-----
P 1 0 4 - 1      -----F A C D E F L W-----
P 0 9 5 - 1      -----A C D E F L-----
P 1 5 5 - 1      -----A C D E F L-----
P 0 9 8 - 2      -----A C D E F L-----
P 0 9 9 - 2      -----A C D E F L-----
P 1 1 5 - 2      -----A C D E F H W-----

```

```

P 1 3 8 - 1      ----- A C D E F H K L P -----
P 1 3 3 - 1      ----- A C D E F -----
P 0 9 9 - 1      ----- A C D E F T L -----
P 1 1 3          ----- A C D E F G L P E F P -----
P 1 3 6 - 1      ----- A C D E K F L -----
P 1 0 6 - 2      ----- A C D E L -----
P 1 3 8 - 2      ----- A C D E F - G K L W -----
P 1 5 6 - 2      ----- A C D E F - G K L W -----
P 1 5 6 - 1      ----- A C D E F - G K -----
P 1 4 4          ----- A M C D E F F G M T W -----

```

\* \*

**Fig. 4** Alignment output of Group A

```

P 1 3 2 - 1      ----- E A C D K L -----
P 1 5 0 - 1      ----- D E A C D K L -----
P 1 0 6 - 1      ----- A C D K L -----
P 1 2 6 - 1      ----- A C D K L -----
P 1 6 4 - 1      ----- A C D K L -----
P 1 0 1 - 2      ----- A C D K M T W -----
P 1 2 5 - 2      ----- F A C D H I H L P Q T W -----
P 1 5 5 - 4      ----- A C D H I Q T W -----
P 1 5 5 - 5      ----- A C D G K W -----
P 1 6 1 - 4      ----- A C D G H W -----
P 1 5 6 - 3      ----- F A C D G N M W -----
P 1 6 7 - 2      ----- A C D G -----
P 1 2 8 - 2      ----- A C D G L K -----
P 1 5 4 - 2      ----- A C D F H W -----
P 1 2 6 - 2      ----- E F A C D F H I H N Q W -----
P 1 3 3 - 2      ----- C D E H K L -----
P 1 3 5          ----- A C D F H K L P -----
P 1 2 3 - 2      ----- F A C D F H K I L -----

```

\* \*

**Fig. 5** Alignment output of Group B

It is apparent that the Function C (mediation, the connective incident vs. beginning counteraction) and Function D (departure) chains are common elements in the alignment for both Group A and Group B. However, the three Group A functions (A, E, F) and one Group B function (A) that exhibit high similarity are not marked as complete common elements. This is due to the fact that their various sequence inputs were grouped roughly.

When more than a certain number of sequences are aligned in this way and there is dispersion in the distances between sequences, dividing the alignment output into groups with fewer sequences and performing alignment again is an

effective means of achieving convergence. In addition to this type of two-stage alignment, for specific functions for which micro-level analysis is desired, manually editing the alignment output based on co-occurrence of functions to give a more sensitive result, then running the alignment again may also be considered. In either case, it is clear that grouping the sequences allows more information to be obtained.

A number of analyses for the alignment output given above are shown below.

With the exception of two examples (P104-1 and P144), all of the 28 Group A sequences have the function chain A-C-D at the start. In addition, with the exception of one example (P136-2), Function D (departure) is followed by Function E (the first function of the donor vs. the hero's reactions). Further, with the exception of two examples (P136-1 and P106-2), Function E is followed by Function F (provision or receipt of magical agent). As a result, the chain A-C-D-E-F appears as a motif at the start of 24 of the 28 Group A sequences.

**Table 6** Motif of Group A

A (villainy VS lack)

C (mediation, the connective incident VS beginning counteraction)

D (departure)

E (the first function of the donor VS the hero's reaction)

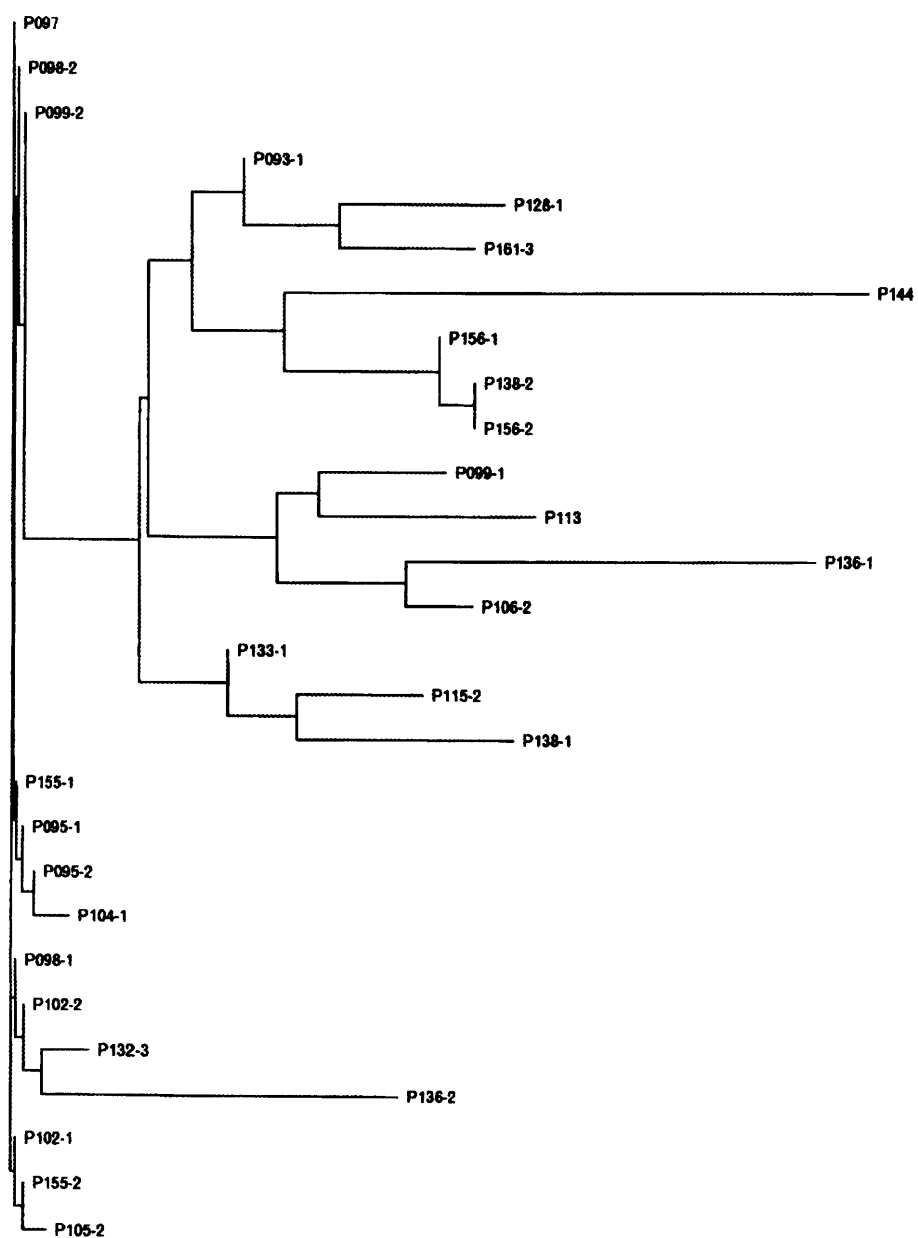
F (provision or receipt of magical agent)

In the case of the 18 Group B examples on the other hand, some have Function E or Function F inserted before Function A at the start, and on the whole, rather than Function D (departure) being followed by Function E (the first function of the donor vs. the hero's reactions) as in the case of Group A (with the exception of P133-2), it is followed by a function such as Function K (the initial misfortune or lack is liquidated), Function G (spatial transference between two kingdoms, guidance), or Function F (provision or receipt of magical agent). From these patterns the following rough description for the motif for the start sections of Group B is obtained.

[EF]ACD[KGF] . . .

The previous discussion was in regard to the parts at the start of the sequences. It is also probably possible to convert the common elements Function C and Function D as fixed chains to specific codes and increase the alignment sensitivity by performing alignment again, and to extract other specific parts and perform alignment to perform local comparisons.

TREEVIEW can be used to output the alignment results in Radical and Phylogram format and facilitate observation. Figures 6 and 7 show the Radicals and Phylograms based on the alignment output for Group A, while Figures 8 and 9 show the Radicals and Phylograms based on the alignment output for Group B.



**Fig. 6** Phylogram of Group A



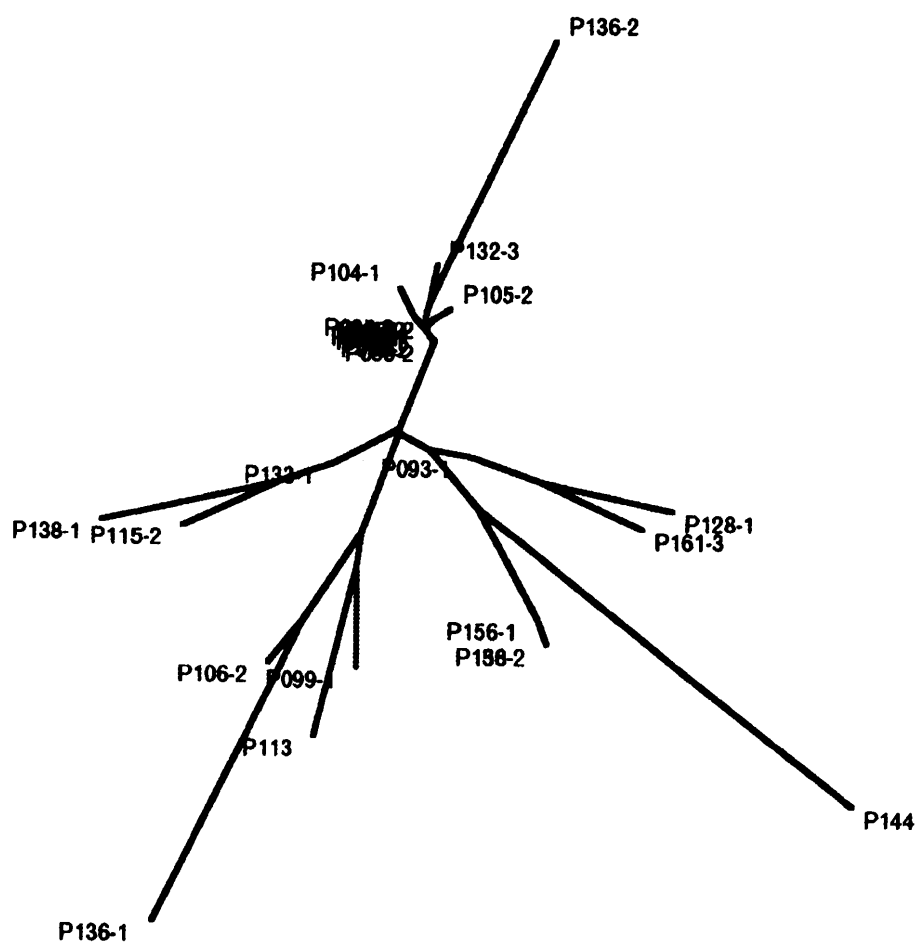
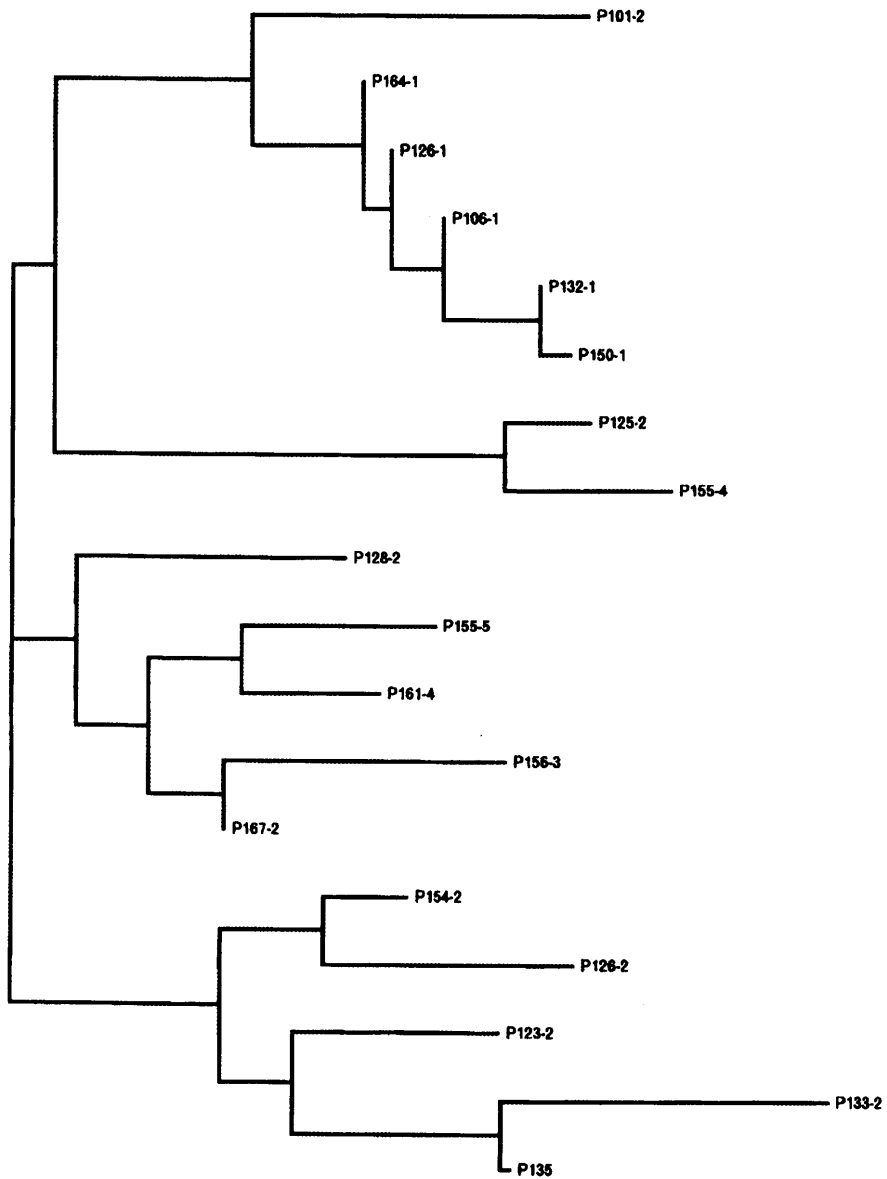
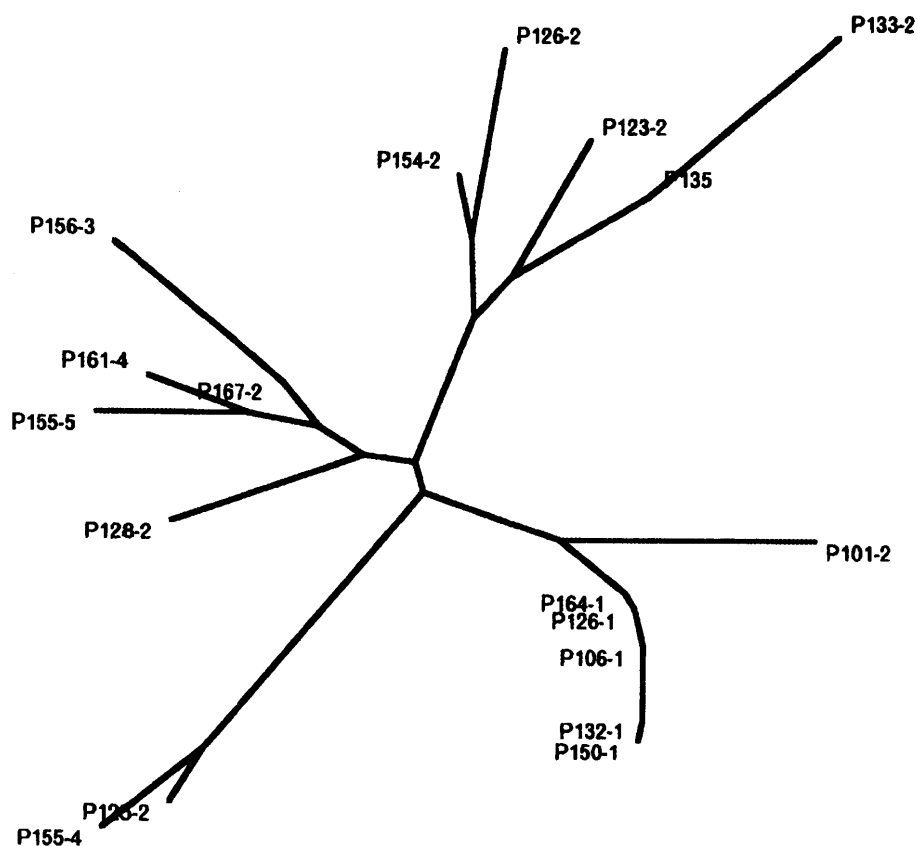


Fig. 7 Radical of Group A



**Fig. 8** Phylogram of Group B



**Fig. 9** Radical of Group B

With the Phylograms, the high-similarity sequences are concentrated on the left side of the tree, and the low-similarity sequences are separated and placed on the right side of the tree. On the other hand, with the Radicals, the distance from the central point exhibits an inversely proportional relationship with the similarity (in other words, when the similarity is high, the mutual distance is short). These figures indicate clearly that all sequences determined to be 'exceptional' in the alignment output are far away. For example, P144, P136-1, P136-2, and P138-1 in Figures 6 and 7, and P101-2, P133-2, and P155-4 in Figures 8 and 9, are separated from the other sequences. It is also possible to observe that the grouping of sequences was performed using function types that followed Function F in the case of the Group A alignment output, and function types that followed Function D in the case of the Group B alignment output. Combining this with the above alignment output seems to provide suggestions for indexing when classifying sequences.

With regard to determining whether the analysis results are significant or not, it is clear that lining up the text analysis of individual tales and treating the sequences of elements that form the high-similarity sequences that require judgment as motifs, then comparing these to low-similarity sequences and identifying the insertions and deletions of specific functions provides a starting point for more detailed analysis.

## CONCLUSION

Through the alignment and its analysis attempted in this paper, I have been able to clearly establish that the analysis methods used for bioinformatics are an effective means for analyzing folktales from several perspectives, and visual presentation of the phylogenetic relationship of sequences enables a lot of information to be garnered. To apply the analysis methods used for bioinformatics more effectively to the analysis of folktales in the future, it will be necessary to investigate both technical problems from an informatics perspective, and narratology problems relating to how to best describe the sequence structure of folktales as data.

An example that relates to the first issue involves the speed and sensitivity limitations of CLUSTAL W when it applied as is to analysis of the structure of the folktale, and this has been clearly established in this paper. As CLUSTAL W was developed for operation on relatively small computer systems, it is effective when used for alignment of sequences that are similar to a certain degree, but to process sequences such as those for folktale structure used in this paper, additional operations are required to improve the sensitivity. Among the alignment programs used for bioinformatics, there are examples that employ different algorithms to CLUSTAL W, and other analysis tools and new software are continually being developed.<sup>15)</sup> I am testing these, as far as possible, while waiting for more sensitive tools to be developed. In addition, I am also looking at tools that do not run on the system currently being used (e.g. tools developed at the ICOT<sup>16)</sup> etc.) and am investigating the possibility of transferring them to high-performance computer systems and mainframes. For this purpose I plan to enlist the cooperation of specialists in the fields of informatics and biology.

With regard to the second problem, there is a need to construct a model that optimally describes the sequences used in the structure of folktales. As I stated at the start of this paper, other papers using description models from the results of case grammar, semantic structure and artificial intelligence are currently being prepared.

Finally, I would like to add a few words about unification of analysis methods for folktale structure. Research into analysis methods for folktales and similar texts can be divided into two types: research that focuses on syntagmatic structure carried out by Propp and others, and research that emphasizes the paradigmatic structure such as the work of Lévi-Strauss<sup>17)</sup> (there are also analysis methods that handle both these structures<sup>18)</sup>). Methodologically speaking, this paper has focused on analysis that uses the syntagmatic structure of the folktale, but by combining this with paradigmatic structure analysis methods, I believe that more fruitful folktale

analysis will be possible. In my past numerical research into paradigmatic structure,<sup>19)</sup> I obtained some results with regard to the conversion procedure stage for converting folkloristic text into data for processing (individuals [= tale] × variables), but the problem of the arbitrariness of the settings of variables arose. In that research, by handling Propp-like functions[constants] and the variables that form the functions in the same dimension, I was able to obtain the trend characteristics of the genre being analyzed from the co-occurrence pattern of the specific constants and specific variables, but this was nothing more than a rough result.

The objectives of the unified description model that I am currently hypothesizing for the structure of the folktale are to reevaluate the syntagmatic structure as a paradigmatic structure with the entire potential insertion/deletion pattern of the manifestable elements included, and conversely, to describe the variables handled as the configuration in the text as the syntagmatic structure in their order of appearance for the text surface structure. By doing this, it will be possible to simultaneously analyze both the syntagmatic structure and paradigmatic structure. In other words, it will realize an analysis method that unifies morphology and transformation. In addition, it will probably become possible for conventional narratology to unify the two different directions of deep structure [content] and surface structure [expression] that have become the subject of analysis.

## NOTES

- 1) Oda 1997.
- 2) Seki 1981 (1955), p.60. Boogaard 1984, p.667.
- 3) Propp 1968 (1928).
- 4) Greimas 1966.
- 5) Beatie 1979, p.194.
- 6) Example from Uchiyama, I. *et al.* (Mitaku, S. and M. Kanehisa (eds.) 1995).
- 7) Thompson, J. D., Higgins, D. G. and T. J. Gibson 1994.  
The program was downloaded from the European Bioinformatics Institute site (<http://www.ebi.ac.uk/>)
- 8) Page 1996.  
The program was downloaded from the Ribosomal Database Project FTP Server (University of Illinois at Urbana-Champaign).
- 9) Propp 1968 (1928), p.21.
- 10) *ibid.*, pp.135-143.
- 11) When alignment was actually performed on these sequences, it was discovered that the distances between sequences were abnormal.
- 12) The 31 functions established by Propp were reduced to 11 function pairs (22 functions) by Greimas. Also, the seven functions ( $\beta$ ,  $\gamma$ ,  $\delta$ ,  $\varepsilon$ ,  $\zeta$ ,  $\eta$ ,  $\theta$ ) for the preparatory part of the tale of Propp's Schemes were omitted. In addition, Greimas eliminated Propp's Function L (unfounded claims), and as a result of this reduction, when Propp's Schemes are converted to the Greimas model, 16 functions result.

- 13) For proteins, CLUSTAL W recognizes the following 20 amino acids. A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, and Y.
- 14) This is called the FASTA format (Pearson and Lipman).
- 15) [http://www.ebi.ac.uk/biocat/Alignment\\_Search\\_software.html](http://www.ebi.ac.uk/biocat/Alignment_Search_software.html)
- 16) Institute for New Generation Computer Technology  
<http://www.icot.or.jp/AITEC/3/IFS/list/catalogue-J.html#experimental>
- 17) Lévi-Strauss 1958.
- 18) Barthes 1977 (1966).
- 19) Oda *et al.* 1984 and Oda 1984.

## REFERENCES

- Barthes, R.  
1977 (1966) Introduction à l'analyse structurale des récits. *Poétique du Récit*, pp.7-57. Paris: Éditions du Seuil.
- Beatie, B.A.  
1979 Measurement and the Study of Literature. *Computers and the Humanities* 13, pp.185-194.
- Boogaard, N. van den  
1984 La définition du fabliau dans les grands recueils. *Epopée, animale, fable, fabliau. Cahiers d'Études Médiévales* 2-3, pp.657-668. Paris: P.U.F.
- Greimas, A.-J.  
1966 *Sémantique Structurale*. Paris: Larousse.
- Lévi-Strauss, C.  
1958 *Anthropologie Structurale*. Paris: Plon.
- Mitaku, S. and M. Kanehisa (eds.)  
1995 *Human Genome Program and Knowledge Information Processing*. Tokyo: Baifukan. [in Japanese]
- Oda, J. *et al.*  
1984 Methods of Bibliographic Information Analysis and Structure of Literary Text. *Proceedings of International Computer Symposium 1984 (ICS '84)*. vol. II, pp.908-915. Tamkang Univ.
- Oda, J.  
1984 Structure du narration du fabliau: Essai d'un reclassement des pièces par Cluster-Analysis. *Mathematical Linguistics* 14-7, pp.281-303.  
1997 Description of Structure of the Folktale I: From Bioinformatics to Text Theory. In J. Oda (ed.), *Genesis of Narrative*, pp.67-92. Tokyo: Institute for the Study of Languages and Cultures of Asia and Africa, Tokyo University of Foreign Studies. [in Japanese]
- Page, R. D. M.  
1996 TREEVIEW: An Application to Display Phylogenetic Trees on Personal Computers. *Computer Applications in the Biosciences* 12, pp.357-358.
- Propp, V.  
1928 *Morfologija skázki*. Leningrad. (English edition translated by Laurence Scott, 1968, *Morphology of the Folktale*. Austin: University of Texas Press)
- Seki, K.  
1981 (1955) *Structure of the Folktale*. Tokyo: Dohosha. [in Japanese]

- Thompson, J.D., Higgins, D.G. and T.J. Gibson  
1994 CLUSTAL W: Improving the Sensitivity of Progressive Multiple Sequence Alignment through Sequence Weighting, Positions-specific Gap Penalties and Weight Matrix Choice. *Nucleic Acids Research* 22, pp.4673-4680.