



のペアを、最も近縁であるとして結び、親ノードを作ることの繰り返しで系統樹を推定する。

次の近縁を求める際は、結ばれた二つのノードの配列の代わりに、親ノードの配列を用いて総当たりを行う。

## 2.2 親ノードのアミノ酸配列の仮決定

アラインメントされた配列 a1, a2 を用いて、親ノードにおけるアミノ酸配列を推定する。a1 と a2 のサイト n を比較し、同じであれば親ノードのサイト n も同じ文字であるとし、異なれば決めることができないということで図3のように候補として両方の文字を保持する。

a1, a2 の少なくとも一方がサイト n で候補を持っている場合、a1 と a2 のサイト n における候補の中に重複があれば、それを親ノードのサイト n として採用する。重複がなければ全てを候補として保持する。

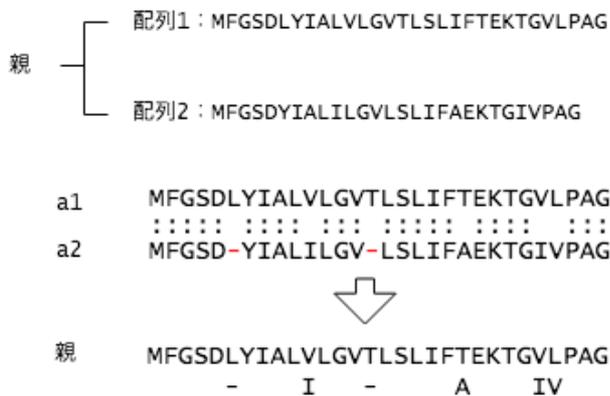


図3: 配列1と配列2から予測したノード「親」におけるアミノ酸配列。アラインメントの結果アミノ酸が一致しているサイトを重点で示した。重点がないサイトは親において候補を持っている。

## 2.3 各ノードにおけるアミノ酸配列の尤度の決定

各ノードにおいて、候補を持つサイトは子孫ノードで候補のアミノ酸を持っている。つまり、分岐していく過程で候補のアミノ酸の全てに進化していくことになる。このことから、アミノ酸置換行列である PAM1 を用いて、進化の起こりやすさを算出した[5]。

## 3 精度検証実験と結果

進化の起点とするアミノ酸配列と任意に形を決めた系統樹を用意し、PAM1 を用いて変化させることで子ノードのアミノ酸配列を決めることを繰り返し、擬似的な進化系統のデータを作成した。

このデータの変異後のアミノ酸配列を用い、本手法を適用することで、各ノードのアミノ酸配列を正しく予測できたか確認し、ANCESCON の適用結果と比較した。

この予測結果について、各ノードを、(1)葉を結んだもの、(2)中央付近で葉に近いもの、(3)根に近いもの、(4)根にあたるものの4種類に分類し、精度を検証した。

結果、4種類のノードすべてで本手法が高い精度で予測することができていた。

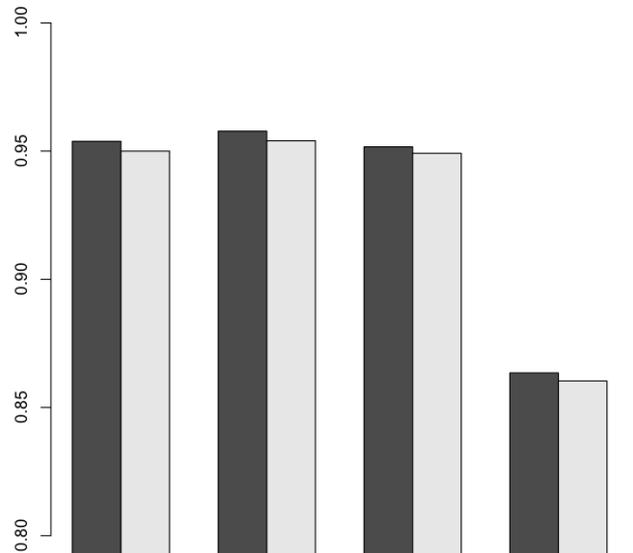


図4: 本手法と ANCESCON の精度比較。黒いバーが本手法、グレーのバーが ANCESSCON

## 4 今後の課題

結果で挙げたように(4)の根にあたるノードでの精度が著しく落ちていたため、原因を検証し、精度向上を計りたい。

また、本研究の手法では文字列としてのアミノ酸配列の類似度のみを考えているが、本来タンパク質は立体構造を取ることで機能を持つものである。そのため、アミノ酸配列は直前、直後のアミノ酸との関係だけではなく、立体構造を取った時に近くに位置するアミノ酸の影響も受けるため、文字列としての小さな変化が機能的に大きな変化につながることもある。しかし、現在立体構造の判明していないタンパク質も未だ多く存在する。このため、離れた位置にあるアミノ酸と変異の相関関係をアミノ酸配列から構造の違いも予測することができれば、より本来の進化系統に近い、機能的に似ているものを近縁として予測することができるのではないかと考える。

## 参考文献

- [1] Joseph W Thornton, Resurrecting ancient genes: experimental analysis of extinct molecules, *Nature Reviews Genetics*(2004), 5:366:375
- [2] D.Sadava(Eds.), 大学生物学の教科書 第4巻 進化生物学 石崎泰樹・斎藤成也(訳), ブルーバックス(2014), pp.112-116
- [3] Wei Cai, Jimin Pei, and Nick V. Grishin, Reconstruction of ancestral protein sequences and its applications, *BMC Evol Biol.*(2004), 17:4:33
- [4] Needleman, Saul B. & Wunsch, Christian D. (1970)"A general method applicable to the search for similarities in the amino acid sequence of two proteins"*JMB* 48 (3): 443-53.
- [5] M. O. Dayhoff, R. M. Schwartz, and B. C. Orcutt (1997) *Atlas of Protein Sequence and Structure* 第5巻, *National Biomedical Research Foundation*, pp.348