# Appendix A

# Derivation of the Proton Linkage Model

The derivation given here is taken from Laskowski & Finkenstadt, 1972 and was slightly modified without changing the basic principles. The original derivation was, however, done earlier (Wyman, 1964). I give it here for the sake of completeness.

Given is the following association reaction:

$$A + B \rightleftharpoons AB \tag{A.1}$$

The equilibrium constant is

$$K_A = \frac{[AB]}{[A][B]} \tag{A.2}$$

Each of the species A, B, and AB can bind protons. The total concentration of A is then given by

$$[A] = \sum_{i=0}^{n_A} [AH_{n_A - i}] \tag{A.3}$$

Analogous expressions are used for $[B]$ and $[AB]$, respectively. Then eq A.2 can be rewritten as

$$K_A = \frac{\sum\limits_{i=0}^{n_{AB}} [ABH_{n_{AB} - i}]}{\sum\limits_{i=0}^{n_A} [AH_{n_A - i}] \sum\limits_{i=0}^{n_B} [BH_{n_B - i}]} \tag{A.4}$$

where $n_A$, $n_B$, and $n_{AB}$ are the numbers of titratable groups of A, B, and AB, respectively. The equilibrium between the fully-protonated species is defined as reference[1].

$$A_{n_A} + B_{n_B} \rightleftharpoons AB_{n_{AB}} \tag{A.5}$$

---

[1]Laskowski & Finkenstadt, 1972 allow that also additional proton binding sites can be created upon association. If however two protein associate, the number of proton binding sites does not change. While Laskowski & Finkenstadt, 1972 more generally assume that the number of proton binding sites in the docked complex is given by $n_{AB} = c + n_A + n_B$, where $c$ represents the number of proton binding sites created upon association, I assume here that $n_{AB} = n_A + n_B$. The extension to the form used in Laskowski & Finkenstadt, 1972 is in principle possible, but makes the derivation more complicated.

The equilibrium constant for the reaction of the fully-protonated species is

$$K_A^o = \frac{[AB_{n_{AB}}]}{[A_{n_A}][B_{n_B}]} \tag{A.6}$$

Now an ionization constant as shown for A in eq A.7 is introduced for each species.

$$[AH_{n_A}] \rightleftharpoons [AH_{n_A-i}] + i[H^+]; \quad L_{0A} = 1; \quad L_{iA} = \frac{[AH_{n_A-i}][H^+]^i}{[AH_{n_A}]} \tag{A.7}$$

Using eq A.6 and eq A.7, eq A.4 becomes

$$K_A = K_A^o \frac{\sum\limits_{i=0}^{n_{AB}}[H^+]^{-i}L_{iAB}}{\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA} \sum\limits_{i=0}^{n_B}[H^+]^{-i}L_{iB}} \tag{A.8}$$

Taking the logarithm of both sides of eq A.8 and differentiating with respect to pH $= -\lg[H^+]$, we get

$$\frac{d\lg K_A}{dpH} = \frac{d\lg(K_A^o)}{dpH} + \frac{d\lg(\sum\limits_{i=0}^{n_{AB}}[H^+]^{-i}L_{iAB})}{dpH} - \frac{d\lg(\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA})}{dpH} - \frac{d\lg(\sum\limits_{i=0}^{n_B}[H^+]^{-i}L_{iB})}{dpH} \tag{A.9}$$

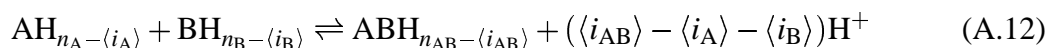The derivation of the last three terms is demonstrated for A

$$
\begin{aligned}
\frac{d\lg(\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA})}{dpH} &= \frac{-[H^+]}{\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA}} \frac{d\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA}}{d[H^+]} \\[2mm]
&= \frac{[H^+]\sum\limits_{i=0}^{n_A}i[H^+]^{-i-1}L_{iA}}{\sum\limits_{i=0}^{n_A}[H^+]^{-i}L_{iA}} \\[2mm]
&= \langle i_A \rangle
\end{aligned} \tag{A.10}
$$

where $\langle i_A \rangle$ represent respectively the average number of protons released from the fully-protonated form at the pH of interest. Therefore, we obtain

$$\frac{d\lg K_A}{dpH} = \overline{q} = \langle i_{AB} \rangle - \langle i_A \rangle - \langle i_B \rangle \tag{A.11}$$

where $\langle i_A \rangle$, $\langle i_B \rangle$, and $\langle i_{AB} \rangle$ represents the average number of protons released from the fully-protonated form of A, B, or AB at the pH of interest. Therefore, $\overline{q}$ is the number of protons released upon association. The reaction equilibrium at a defined pH is then given by

$$AH_{n_A-\langle i_A \rangle} + BH_{n_B-\langle i_B \rangle} \rightleftharpoons ABH_{n_{AB}-\langle i_{AB} \rangle} + (\langle i_{AB} \rangle - \langle i_A \rangle - \langle i_B \rangle)H^+ \tag{A.12}$$

Integrating eq A.11 over pH results in

$$\lg K_A(pH_2) - \lg K_A(pH_1) = \int\limits_{pH_1}^{pH_2} \overline{q}\,dpH \tag{A.13}$$

Therefore, if the association constant at $pH_1$ is known, the association constant at $pH_2$ is supposed to be given by the expression above. An analogous expression for the pH dependence of redox potential can be easily derived.

# Appendix B

# Cluster Algorithm

Cluster analysis is used to relate objects on the basis of an object specific quantity using a defined distance measure. A large variety of cluster algorithms exists (Jain & Dubes, 1988). Some of them have been applied to cluster molecular conformations (Torda & van Gunsteren, 1994; Shenkin & McDonald, 1994; Boutonnet *et al.*, 1995). Here, I describe a new type of algorithm that Ernst-Walter Knapp and I developed and implemented during my PhD work to cluster molecular conformations on the basis of their mean square deviations (Ullmann *et al.*, 1998b).

## B.1   Scoring Function

Given is a set of $M$ molecular structures $\{\vec{x}_i\}$, $i = 1, \ldots, M$. Each structure $\vec{x}_i$ consists of $n$ atoms $\vec{x}_i(j)$; $\vec{x}_i = (\vec{x}_i(1), \ldots, \vec{x}_i(n))$. This set can be subdivided into $N$ clusters. Eq B.1 defines the average structure $\langle \vec{x} \rangle_k$ of cluster $k$ with $m_k$ elements.

$$\langle \vec{x} \rangle_k = \frac{1}{m_k} \sum_{i=1}^{m_k} \vec{x}_i; \qquad \sum_{k=1}^{N} m_k = M \tag{B.1}$$

The average structure of a cluster can also be interpreted as the center of this cluster. The squared distance $d^2(o, p)$ between two structures $o$ and $p$ is given by eq B.2.

$$d^2(o, p) = \frac{1}{n} \sum_{j=1}^{n} (\vec{x}_o(j) - \vec{x}_p(j))^2 \tag{B.2}$$

where $\vec{x}_o(j)$ designates the $j$-th atom of structure $o$, $n$ is the number of atoms in the structure. The variance of cluster $k$ $\langle \Delta x^2 \rangle_k$ is defined as the mean square distance between the elements of the cluster and the cluster center (eq B.3).

$$\begin{aligned} \langle \Delta x^2 \rangle_k &= \frac{1}{m_k} \sum_{p=1}^{m_k} \frac{1}{n} \sum_{i=1}^{n} (\vec{x}_p(i) - \langle \vec{x}(i) \rangle_k)^2 \\ &= \frac{1}{m_k} \sum_{p=1}^{m_k} d^2(p, \langle k \rangle) \end{aligned} \tag{B.3}$$

$$\tag{B.4}$$

where $m_k$ is the number of elements of cluster $k$. The squared distance between two clusters $k$ and $l$ can be defined as the distance of the centers of these two clusters as $d^2(\langle k \rangle, \langle l \rangle)$ given in eq B.2.

Now, we define the scoring function $S$ in eq B.5,

$$S = a \sum_{i=1}^{N} \langle \Delta x^2 \rangle_i + \frac{N-1}{M-1} \langle \Delta x^2 \rangle \tag{B.5}$$

where $a$ is a tuning parameter to vary the average cluster size, $n$ is the number of clusters, $m$ is the total number of structures, $\langle \Delta x^2 \rangle_i$ is the variance of cluster $i$, and $\langle \Delta x^2 \rangle$ is the total variance over all structures. Minimization of this scoring function leads to a small variance within the clusters. The number of clusters becomes, however, not too large due to the second term in eq B.5.

An implementation of the equations outlined above would require an update of the average structure after each step. The variance must therefore be calculated from scratch after each step. The mean square deviation between the center of the cluster and each of its elements must be recalculated, because the center of the cluster changes. This would lead to an enormous computational burden. It can, however, be shown that

$$\langle \Delta x^2 \rangle_k = \frac{1}{m_k} \sum_{p=1}^{m_k} d^2(p, \langle k \rangle) \tag{B.6}$$

$$= \frac{1}{2m_k^2} \sum_{p=1}^{m_k} \sum_{o=1}^{m_k} d^2(p, o)$$

The reformulation of the eq B.3 in eq B.6 allows the use of a precalculated mean square deviation matrix in the cluster algorithm. The variance of each cluster can be updated and must not be calculated from scratch after each step.

## B.2   Data Organization and Implementation

The cluster algorithm is written in the programming language C (Kernighan & Ritchie, 1990). To implement the cluster algorithm efficiently, it is required to add and remove elements from a list. For that purpose, we developed a double-linked list, which makes it possible to remove elements even from the middle of a linked list and repair the remaining list. The clusters are grouped in a double-linked list (Figure B.1a). Each cluster is a structure that contains a pointer `prev` pointing to the previous cluster, a pointer `next` pointing to the next cluster, a pointer `Element` pointing to a double-linked list of elements, an integer `no_elements` that counts the number of elements of the cluster, and a real number `variance` that contains the variance of the cluster.

```
struct Cluster_struct {struct Element_struct    *element;
                        struct Cluster_struct    *prev;
                        struct Cluster_struct    *next;
                        int                      no_elements;
                        double                   variance;
                        };
```

The pointer `prev` of the first cluster and the pointer `next` of the last cluster point to NULL. If a cluster is removed from a double-linked list, the pointer are redirected.
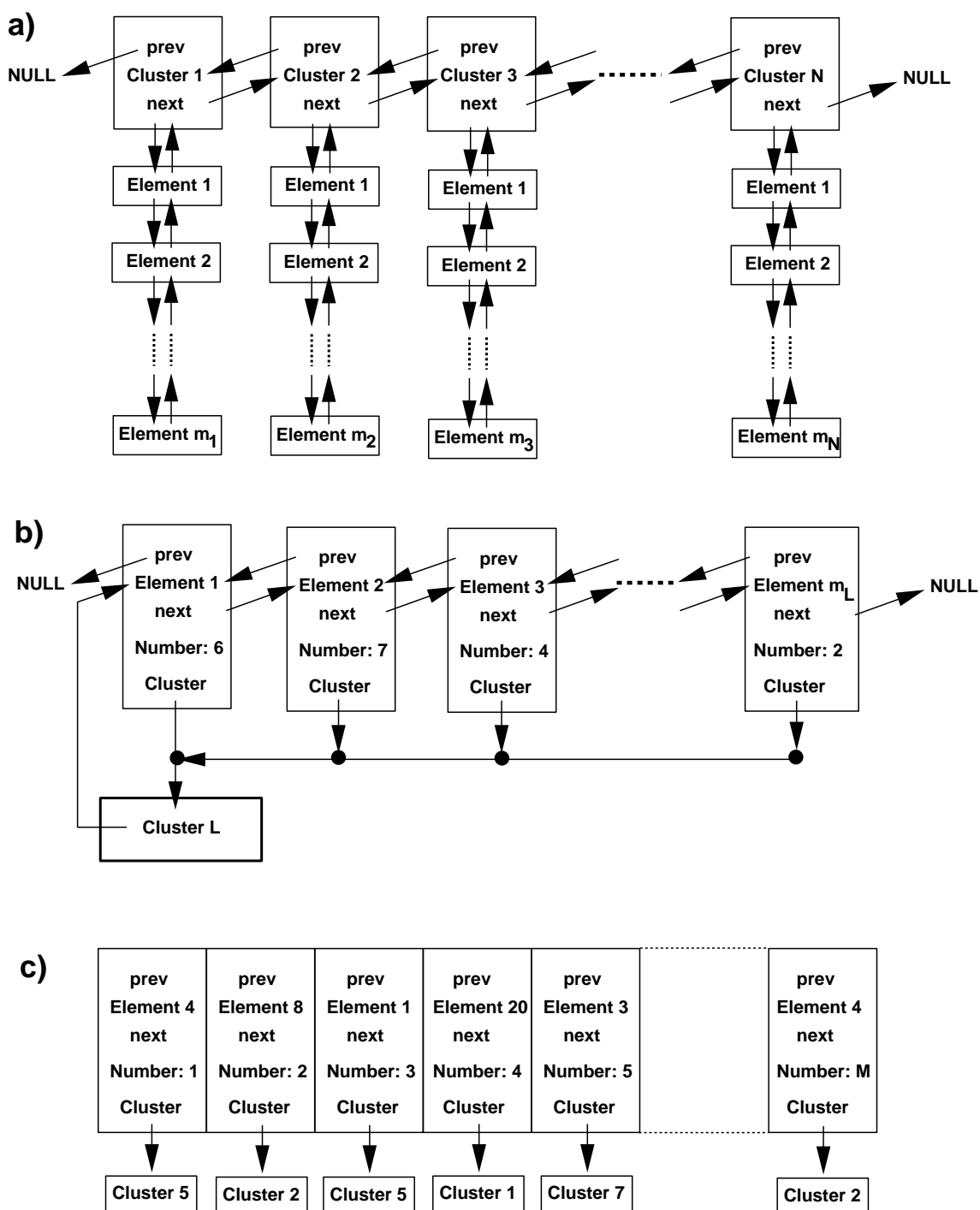
**a)**

| | | | |
|---|---|---|---|
| NULL ← | prev **Cluster 1** next | prev **Cluster 2** next | prev **Cluster 3** next · · · · prev **Cluster N** next | → NULL |

**Element 1**    **Element 1**    **Element 1**    **Element 1**

**Element 2**    **Element 2**    **Element 2**    **Element 2**

**Element $m_1$**    **Element $m_2$**    **Element $m_3$**    **Element $m_N$**

**b)**

NULL ←

| prev **Element 1** next **Number: 6** Cluster | prev **Element 2** next **Number: 7** Cluster | prev **Element 3** next **Number: 4** Cluster | · · · · prev **Element $m_L$** next **Number: 2** Cluster | → NULL |
|---|---|---|---|---|

**Cluster L**

**c)**

| prev **Element 4** next **Number: 1** Cluster | prev **Element 8** next **Number: 2** Cluster | prev **Element 1** next **Number: 3** Cluster | prev **Element 20** next **Number: 4** Cluster | prev **Element 3** next **Number: 5** Cluster | · · · · | prev **Element 4** next **Number: M** Cluster |
|---|---|---|---|---|---|---|
| **Cluster 5** | **Cluster 2** | **Cluster 5** | **Cluster 1** | **Cluster 7** | | **Cluster 2** |

**Figure B.1**: Data organization for the cluster algorithm. a) Double-linked list of clusters. To each cluster a double-linked list of elements is assigned. b) Double-linked list of elements assigned to each cluster. Each element contains a pointer to that cluster to which it is assigned. c) One-dimensional array of elements. The number assigned to each element corresponds to the position of this element in the one-dimensional array.

```
if ((*cluster).next!=NULL) /* if not at the end of the list */
   {help=(*cluster).next;
    (*help).prev=(*cluster).prev;
   }
if ((*cluster).prev!=NULL) /* if not the top of the list */
   {help=(*cluster).prev;
    (*help).next=(*cluster).next;
    /* search for the top of the list */
    while ((*help).prev!=NULL) {help=(*help).prev;}
   }
else {help=(*cluster).next; /* help is now the top of the list */
     }
```

After this procedure, the variable `help` points to the top of the double-linked list of clusters. A new cluster is always added to the top of the double-linked list.

```
(*top).prev=new_cluster; /* Pointer redirection  */
(*new_cluster).next=top;
```

The double-linked list of elements is constructed similarly to the double-linked list of clusters, i. e., each element is a structure that contains two pointers, `prev` and `next` pointing to the previous and the next cluster in the double-linked list respectively. Besides, each element has a pointer `cluster` to the cluster to which the element is assigned and. an integer `number` that assigns this element to a position in an one dimensional array (Figure B.1c).

```
struct Element_struct {int                      number;
                       struct Cluster_struct    *cluster;
                       struct Element_struct    *prev;
                       struct Element_struct    *next;
                      };
```

Addition and elimination of an element to or from the double-linked list is done analogously to the addition and elimination of cluster to or from the double-linked list show above. This integer is required to access the mean square deviation in the precalculated mean square deviation matrix.

The algorithm uses several elementary steps: it moves an element from one cluster to another, it merges two clusters, it divides one cluster into two new clusters, and it sprinkles one cluster, i. e., assigns the elements of one cluster to other clusters. A step is accepted or rejected according to the Metropolis criterion (Metropolis *et al.*, 1953) in order to avoid to get trapped in a local minimum. Figure B.2 shows the flow chart of the program. The following commands can be currently combined in the program.

**distribute_elements.** This command generates the initial distribution of the elements to clusters. In the first step, one cluster is created and the first element is put into this cluster. The other elements are put into the clusters already created in the proceeding steps of this routine if the mean square deviation between this element and all elements in the cluster is smaller than a given threshold, otherwise it is put into a newly-created cluster.

**random_element_move.** This command moves a randomly-selected element to a randomly-selected cluster.

**systematic_element_move1.** This command moves a randomly-selected element to the cluster that has the smallest distance to the selected element but is not the cluster to which the element is assigned.
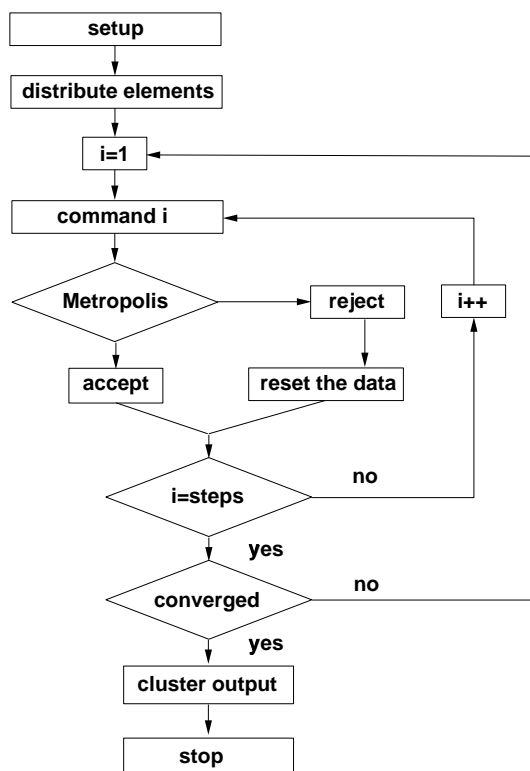
**Figure B.2**: Flow Chart of the Cluster Algorithm Program. A step is accepted or rejected according to the Metropolis criterion. The clustering is iterated until convergence is reached, i. e., until the difference between the old and the new score is smaller than a given threshold.

**random_cluster_divide.** This command divides a randomly-selected cluster. The elements that have the largest distance in the selected cluster are put into two newly-created separate clusters. The other elements are sorted to one of these two clusters according to their mean square deviation from the first element in the newly-created cluster.

**systematic_cluster_divide.** This command divides the cluster with the largest variance. The actual divide occurs in the same way as in random_cluster_divide.

**random_cluster_sprinkle.** This command sprinkles a randomly-selected cluster, i. e., its elements are assigned to the cluster to that the elements have the smallest distance. The selected cluster is eliminated after this move.

**limit_cluster_sprinkle.** This command sprinkles a randomly-selected cluster if the number of elements is smaller than a given number.

**systematic_cluster_sprinkle.** This command sprinkles the cluster with the largest variance.

**random_cluster_merging.** This command merges two randomly-selected clusters.

**systematic_cluster_merging.** This command merges the clusters which have the smallest distance from each other.

**repeat_until_convergence.** The last N step of the script are repeated until the difference between the old score and the new score is smaller than a given threshold, i. e., until the run converged.

108

# Appendix C

# Sequences

All sequences were obtained from SWISSPROT.

## Cytochrome *f*

```
                                   60          70
```

**Higher Plants**
*Brassica rapa*  (turnip)        ...VLANGKKGALN...
*Oenothera hookeri*  (primose)   ...VLANGKKGGLN...
*Nicotiana tabacum*  (tobacco)   ...VLANGKRGGLN...
*Spinacia oleracea*  (spinach)   ...VLANGKKGGLN...
*Glycine max*  (soybean)         ...VLANGKKGALN...
*Pisum sativum*  (graden pea)    ...VLANGKKGALN...
*Vicia faba*  (broad bean)       ...VLANGKKGALN...
*Triticum aestivum*  (wheat)     ...VLANGKKGGLN...
*Zea mays*  (maize)              ...VLANGKKGGLN...
*Oryza sativa*  (rice)           ...VLANGKKGGLN...
*Pinus thunbergii*  (green pine) ...VLANGKKGALN...
*Marchantia polymorpha*          ...VLANGKKGSLN...
**Eukaryotic Algae**
*Chlorophyta (Green Algae)*
*Chlamydomonas reinhardtii*      ...VLANGKKGDLN...
*Bacillariophyta (Diatoms)*
*Odontella sinensis*             ...VGANGKKADLN...
*Glaucophyta*
*Cyanophora paradoxa*            ...IQANGQKGPLN...
*Rhodophyta (Red Algae)*
*Porphyra purpurea*              ...ILGNGSKGGLN...
**Cyanobacteria**
*Synechocystis*  (Strain PCC 6803)  ...VLGDGSKGGLN...
*Synechococcus*  (Strain PCC 7002)  ...VLGDGSKGGLN...
*Nostoc*  (Strain PCC 7906)         ...VGADGSKVGLN...
*Phormidium laminosum*              ...VQADGSKGPLN...

Residues 60 through 70 in all known sequences of cytochrome *f*. The residue Lys65 is missing in the cytochromes *f* from cyanobacteria and some primitive eukaryotic algae.

# Cytochrome $c_6$

|  |  |  |
|---|---|---|
| | 50 | 60 |

**Eukaryotic Algae**

*Chlorophyta (Green Algae)*

*Chlamydomonas reinhardtii*     `...IYQVENGKGAM...`

*Monoraphidium braunii*     `...VYQIENGKGAM...`

*Bryopsis maxima*     `...TSQVRNGKGAM...`

*Euglenophyta*

*Euglena gracilis*     `...EYQVRNGKGPM...`

*Euglena viridis*     `...EYQVRNGKGPM...`

*Phaeophyta (Brown Algae)*

*Petalonia fascia*     `...TYQVTNGKNAM...`

*Alaria esculenta*     `...TYQVTNGKNAM...`

*Rhodophyta (Red Algae)*

*Porphyra tenera*     `...TYQVQNGKNAM...`

*Porphyra purpurea*     `...TYQVTNGKNAM...`

*Xanthophyta (Yellow Algae)*

*Bumilleriopsis filiformis*     `...TYQVTNGKNAM...`

*Chrysophyta (Golden Algae)*

*Monochrysis lutheri*     `...VYQVTNGKNAM...`

**Cyanobacteria**

*Spirulina maxima*     `...AYQVTNGKNAM...`

*Synechococcus lividus*     `...IYQVQHGKNAM...`

*Synechococcus sp.*     `...MYQVQNGKNAM...`

*Synechococcus* (Strain PCC 7942)     `...TTQVTNGKGAM...`

*Synechocystis* (Strain PCC 6803)     `...VAQITNGNGAM...`

*Microcysis aeruginosa*     `...VTQVTKGMGAM...`

*Anacystis nidulans*     `...TTQVTNGKGAM...`

*Aphanizomenon flos−aquae*     `...GAQVTNGKNAM...`

*Plectonema boryanum*     `...IAQVTHGKGAM...`

*Anabaena variabilis*     `...VAQVTNGKGAM...`

*Anabaena* (Strain PCC 7120)     `...IAQVTNGKNAM...`

*Anabaena* (Strain PCC 7937)     `...IAQVTNGKNAM...`

Residues 50 through 60 in all known sequences of cytochrome $c_6$. The residue Tyr51, which is aligned with the dipole moment of the protein, is replaces by non-aromatic amino acids in cytochrome $c_6$ of many species.

# Cytochrome $c_6$

```
                              60          70
```

**Eukaryotic Algae**
*Chlorophyta (Green Algae)*
*Chlamydomonas reinhardtii*     `...MPAWADRLSEE...`
*Monoraphidium braunii*     `...MPAWDGRLDED...`
*Bryopsis maxima*     `...MPAWSDRLDDE...`
*Euglenophyta*
*Euglena gracilis*     `...MPAWEGVLSED...`
*Euglena viridis*     `...MPAWEGVLDES...`
*Phaeophyta (Brown Algae)*
*Petalonia fascia*     `...MPAFGGRLSET...`
*Alaria esculenta*     `...MPAFGSRLAET...`
*Rhodophyta (Red Algae)*
*Porphyra tenera*     `...MPAFGGRLVDE...`
*Porphyra purpurea*     `...MPAFGGRLVDE...`
*Xanthophyta (Yellow Algae)*
*Bumilleriopsis filiformis*     `...MPAFGGRLSDS...`
*Chrysophyta (Golden Algae)*
*Monochrysis lutheri*     `...MPAFGGRLEDD...`
**Cyanobacteria**
*Spirulina maxima*     `...MPGFNGRLSPK...`
*Synechococcus lividus*     `...MPAFAGRLTDE...`
*Synechococcus sp.*     `...MPAFGGRLSEA...`
*Synechococcus* (Strain PCC 7942)   `...MPAFGSKLSAD...`
*Synechocystis* (Strain PCC 6803)   `...MPGFKGRISDS...`
*Microcysis aeruginosa*     `...MPAFGGRLSAE...`
*Anacystis nidulans*     `...MPAFGAKLSAD...`
*Aphanizomenon flos−aquae*     `...MPAFGIRLKAE...`
*Plectonema boryanum*     `...MPAFKGRLSDD...`
*Anabaena variabilis*     `...MPAFKGRLKPE...`
*Anabaena* (Strain PCC 7120)     `...MPAFKGRLKPE...`
*Anabaena* (Strain PCC 7937)     `...MPAFKGRLKPD...`

Residues 60 through 70 in all known sequences of cytochrome $c_6$. Only tryptophane (W) or phenylalanine (F) is found at position 63. The aromatic group may interact with cationic the side chain of arginine (R) or lysine (K) at position 66 to form a cation-π complex. This cationic residue is missing in the two Euglenophyta.

## Plastocyanin

|                                                   | 80          90 |
|---------------------------------------------------|----------------|
| **Higher Plants**                                 |                |
| *Populus nigra*   ( poplar)                       | ...YSF**Y**CSPHQGA... |
| *Phaseolus vulgaris*   (French bean)              | ...YSF**Y**CSPHQGA... |
| *Vicia faba*   ( broad bean)                      | ...YKF**Y**CSPHQGA... |
| *Pisum sativum*   ( graden pea)                   | ...YKF**Y**CSPHQGA... |
| *Nicotiana tabacum*   (tobacco)                   | ...YTF**Y**CAPHQGA... |
| *Solanum tuberosum*   ( potato)                   | ...YTF**Y**CAPHQGA... |
| *Solanum crispum*   ( potato tree)                | ...YSF**Y**CSPHQGA... |
| *Lycopersicon esculentum*   (tomato)              | ...YTF**Y**CAPHQGA... |
| *Lactuca sativa*   (lettuce)                      | ...YSF**Y**CAPHQGA... |
| *Capsella bursa-pastoris*                         | ...YSF**Y**CAPHQGA... |
| *Arabidopsis thaliana*   (cress)                  | ...YGF**Y**CAPHQGA... |
| *Silene pratensis*   ( white campion)             | ...YKF**Y**CAPHAGA... |
| *Spinacia oleracea*   ( spinach)                  | ...YKF**Y**CSPHQGA... |
| *Cucumis sativus*   ( cucumber)                   | ...YSF**Y**CSPHQGA... |
| *Cucurbita pepo*   ( squash)                      | ...YSF**Y**CSPHQGA... |
| *Petroselinum crispum*   ( parsley)               | ...YKF**Y**CEPHAGA... |
| *Daucus carota*   ( carrot)                       | ...YKF**Y**CEPHAGA... |
| *Mercurialis perennis*                            | ...YSF**Y**CSPHQGA... |
| *Sambucus nigra*   (european elder)               | ...YKF**Y**CSPHQGA... |
| *Rumex obtusifolius*   ( bitter dock)             | ...YSF**Y**CSPHQGA... |
| *Hordeum vulgare*   ( barley)                     | ...YGF**Y**CEPHAGA... |
| *Oryza sativa*   ( rice)                          | ...YGF**Y**CEPHAGA... |
| **Eukaryotic Algae**                              |                |
| *Chlorophyta (Green Algae)*                       |                |
| *Chlamydomonas reinhardtii*                       | ...YGY**Y**CEPHQGA... |
| *Chlorella fusca*                                 | ...YGY**F**CEPHQGA... |
| *Scenedesmus obliquus*                            | ...YGY**F**CEPHQGA... |
| *Ulva arasakii*                                   | ...YGV**Y**CEPHAGA... |
| *Enteromorpha prolifera*                          | ...YGV**Y**CDPHSGA... |
| **Cyanobacteria**                                 |                |
| *Prochlorothrix hollandica*                       | ...YSF**Y**CTPH**R**GA... |
| *Synechocystis*   ( Strain PCC 6803)              | ...YTY**Y**CEPH**R**GA... |
| *Anabaena variabilis*                             | ...YTF**Y**CEPH**R**GA... |
| *Anabaena*   ( Strain PCC 7937)                   | ...YSF**Y**CEPH**R**GA... |
| *Anabaena*   ( Strain PCC 7120)                   | ...YTF**Y**CEPH**R**GA... |
| *Phormidium laminosum*                            | ...YTY**Y**CAPH**R**GA... |

Residues 80 through 90 in all known sequences of plastocyanin. The tyrosine (Y) at position 83 is replaced by a phenylalanine (F) in two of the sequences. The aromatic group may interact with the cationic side chain of Arg88, which is conserved in cyanobacterial plastocyanin, to form a cation-π complex.