

## 2 Transition matrix factoring

### 3 The concept

4 Multiplication by sparse matrices is far more efficient than with dense matrices. Matrix vector  
5 multiplication with a dense matrix is  $O(S^2)$  where  $S$  is the size of the vector; for this application vectors  
6 with thousands of entries are not uncommon, and even larger vectors are possible, although this  
7 depends on the protease and labeling scheme used. For a sparse matrix, matrix vector multiplication can  
8 be made to be  $O(V)$ , where  $V$  is the number of non-zero entries in the matrix. For highly sparse  
9 matrices this can be a significant improvement.

10 Since peptides cannot gain amino acids or functioning fluorophores during sequencing, a basic transition  
11 matrix for fluorosequencing has zeros except for entries for transitions in which the numbers of  
12 fluorophores of each color is decreasing or staying the same. While this does reduce the number of  
13 necessary operations, it only does this by a constant factor, with no effect on the asymptotic behavior in  
14 the limit. Additionally, the number of amino acids either stays the same, decreases by one (from a  
15 successful Edman degradation), or decreases to zero (from a peptide detachment event). This did  
16 improve the asymptotic behavior in the number of non-zero entries of the transition matrix, reducing  
17 this from  $O(\alpha^2 \prod_{c=1}^C \Lambda_c^2)$  to  $O(\alpha \prod_{c=1}^C \Lambda_c^2)$ .

18 However, we did better by factoring this matrix (Fig 4). We used the independence of our different  
19 forms of error, with one matrix in the factored product for each type of error. To demonstrate this  
20 factorization, we reformulated our problem in tensor notation. The vector for the state space of a  
21 peptide with  $C$  colors not undergoing Edman degradation or peptide detachment can be viewed as a  
22 tensor of order  $C$ . Each index of the tensor maps to the fluorophore counts of a different color, and the  
23 value of an index  $i_c$  indicates the number of functioning fluorophores of color  $c$ , and satisfies  $0 \leq i_c \leq$   
24  $\Lambda_c$ . We also have indices  $j_c$  which are similarly defined. Since the transition matrix is a linear mapping

25 from and to this tensor of order  $C$ , it is necessarily of order  $2C$ . We use the Einstein summation  
 26 convention, and three multi-indices  $\mathbf{i} = i_1 i_2 \dots i_C$  and  $\mathbf{j} = j_1 j_2 \dots j_C$  and  $\mathbf{k} = k_1 k_2 \dots k_C$  for convenience.  
 27 The matrix-vector multiplication operation for one step of the HMM forward algorithm is then given by:

$$28 \quad \mathbf{f}_k^{(t+1)} = \mathbf{O}_{kj}^{(t+1)} \mathcal{T}_{ji} \mathbf{f}_i^{(t)} \quad (12)$$

29 Where  $(t)$  and  $(t + 1)$  indicate the timestamp of the values in the order  $C$  tensor  $\mathbf{f}^{(t)}$ , which is indexed  
 30 by the numbers of working fluorophores for each color and is the tensor form of  $\mathbf{f}$  from (1),  $\mathcal{T}$  is the  
 31 transition matrix  $T$  converted into tensor form,  $\mathbf{O}$  is the emission matrix  $O$  converted into tensor form.

### 32 **Considering fluorophore loss only**

33 Assuming no interactions between different fluorophores and ignoring Edman degradation and peptide  
 34 detachment,  $\mathcal{T}$  satisfies the following equation:

$$35 \quad \mathcal{T}_{ji} = \begin{cases} \prod_{c=1}^C \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c}, & \text{if } \mathbf{j} \leq \mathbf{i} \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

36 Where  $p_c$  is the per cycle dye loss rate of the fluorophores for color  $c$ . This is simply the product of the  
 37 binomial distributions for each indexed color of fluorophore. To improve the sparsity of this  
 38 representation, we can factor  $\mathcal{T}$  into second order tensors  $\mathbf{B}^{(1)} \mathbf{B}^{(2)} \dots \mathbf{B}^{(C)}$  such that:

$$39 \quad \mathbf{B}_{ji}^{(c)} = \begin{cases} \binom{i}{j} p_c^{i-j} (1 - p_c)^j, & \text{if } j \leq i \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

40 This produces a factorization of  $\mathcal{T}$ :

$$41 \quad \mathcal{T}_{ji} = \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_C i_C}^{(C)} \quad (15)$$

42 We can plug this into (12) and find:

43 
$$\mathbf{f}_j^{(t+1)} = \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_C i_C}^{(C)} \mathbf{f}_i^{(t)} \quad (16)$$

44 This reduces the algorithmic complexity in this simple case from  $O(\prod_{c=1}^C \Lambda_c^2)$  to

45  $O\left(\left(\prod_{c=1}^C \Lambda_c\right)\left(\sum_{c=1}^C \Lambda_c\right)\right).$

46 **Fluorophore loss and Edman degradation**

47 We can expand on this to consider the Edman degradation: In that case we need more indices for the  
 48 number of remaining amino acids. We modify (12) with additional indices  $u$  and  $v$  which satisfy  $0 \leq$   
 49  $u \leq \alpha$  and  $0 \leq v \leq \alpha$ , indicating the number of successful amino acid removals, or alternatively the  
 50 position of an amino acid in the peptide (i. e., the amino acid at the N-terminus of the peptide when  $u$   
 51 amino acids have been removed). This gives:

52 
$$\mathbf{f}_{vk}^{(t+1)} = \mathbf{O}_{kj}^{(t+1)} \mathcal{T}_{vjui} \mathbf{f}_{ui}^{(t)} \quad (17)$$

53 Note that the emission tensor  $\mathbf{O}$  is unaffected by the amino acid count, and depends only on the  
 54 fluorophore counts, so it does not need to be modified.

55 We modify  $\mathcal{T}$  from (13) to model Edman degradation, and the exact form of  $\mathcal{T}$  will depend on the  
 56 peptide under consideration. Let  $\bar{c}_u$  be a number indicating the color of the fluorophore at position  $u$  in  
 57 the peptide, with a value of 0 indicating no fluorophore, and let  $\lambda_{u, \bar{c}_u}$  indicate the number of  
 58 fluorophores of color  $\bar{c}_u$  remaining when  $u - 1$  amino acids have been removed from the peptide. Then  
 59  $\mathcal{T}$  is defined by:

60 
$$\mathcal{T}_{vjui} = \begin{cases} e\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u \\ (1 - e)\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u = 0 \\ (1 - e) \left( \left(1 - \frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right) \beta(\mathbf{i}, \mathbf{j}) + \left(\frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right) \bar{\beta}(\mathbf{i}, \mathbf{j}, u) \right), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u > 0 \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

61 Where:

62 
$$\beta(\mathbf{i}, \mathbf{j}, u) = \prod_{c=1}^C \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c} \quad (19)$$

63 And:

64 
$$\bar{\beta}(\mathbf{i}, \mathbf{j}, u) = \binom{i_{\bar{c}_u} - 1}{j_{\bar{c}_u}} p_{\bar{c}_u}^{i_{\bar{c}_u} - 1 - j_{\bar{c}_u}} (1 - p_{\bar{c}_u})^{j_{\bar{c}_u}} \prod_{\substack{1 \leq c \leq C \\ c \neq \bar{c}_u}} \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c} \quad (20)$$

65 The probability of an Edman degradation failure is essentially the same as in (13), but multiplied by  $e$  to  
 66 account for the probability of failure. The probability for a transition involving a successful Edman  
 67 degradation event which removes an unlabelable amino acid is similarly just like in (13) but multiplied  
 68 by  $(1 - e)$ , the probability of success. If the amino acid in question is labelable by a color  $\bar{c}_u$ , then we  
 69 may or may not remove a fluorophore of that color in the transition, so we need to take the sum of both  
 70 possibilities.  $\beta$  in (19) gives the standard product of binomials formula from (13), but needs to be  
 71 multiplied by the probability of no dye loss, which in (18) is  $\left(1 - \frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right)$ . This is then summed with  $\bar{\beta}$   
 72 from (20) which gives the product of binomial probabilities starting with one less fluorophore of the  
 73 color  $\bar{c}_u$ , which in (18) is multiplied with the probability of losing a fluorophore with the Edman  
 74 degradation,  $\frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}$ . The sum of these two possibilities is then multiplied by the probability of an Edman  
 75 degradation success, given by  $(1 - e)$ .

76 To make this more efficient, we introduce a new tensor  $\mathcal{E}$  which represents a transformation for Edman  
 77 degradation. We define tensor  $\mathcal{E}$  as:

78 
$$\mathcal{E}_{v\mathbf{k}uj} = \begin{cases} e, & \text{if } v = u \text{ and } \mathbf{k} = \mathbf{j} \\ 1 - e, & \text{if } v = u + 1 \text{ and } \mathbf{k} = \mathbf{j} \text{ and } \bar{c}_u = 0 \\ (1 - e) \left(1 - \frac{j_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right), & \text{if } v = u + 1 \text{ and } \mathbf{k} = \mathbf{j} \text{ and } \bar{c}_u > 0 \\ (1 - e) \left(\frac{j_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right), & \text{if } v = u + 1 \text{ and } k_{\bar{c}_u} = j_{\bar{c}_u} - 1 \text{ and } k_c = j_c \forall c \neq \bar{c}_u \text{ and } \bar{c}_u > 0 \\ 0, & \text{otherwise} \end{cases} \quad (21)$$

79 This provides the following factorization of  $\mathcal{T}$ :

$$80 \quad \mathcal{T}_{vkui} = \mathcal{E}_{vkuj} \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_c i_c}^{(c)} \quad (22)$$

81 By substituting into (17) and adding an additional multi-index  $\mathbf{l} = l_1 l_2 \dots l_c$  we get:

$$82 \quad \mathbf{f}_{vl}^{(t+1)} = \mathbf{O}_{lk}^{(t+1)} \mathcal{E}_{vkuj} \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_c i_c}^{(c)} \mathbf{f}_{ui}^{(t)} \quad (23)$$

83 Despite its high dimensionality,  $\mathcal{E}$  is highly sparse, with no more than three non-zero entries per column  
 84 (here, meaning column in the original non-tensor form matrix). This reduces the algorithmic complexity  
 85 from  $O(\alpha \prod_{c=1}^C \Lambda_c^2)$  to  $O\left(\alpha \left(\prod_{c=1}^C \Lambda_c\right) \left(\sum_{c=1}^C \Lambda_c\right)\right)$ . We note that while the extraction of the Edman  
 86 degradation tensor appears to have little direct effect on the algorithmic complexity reduction, which is  
 87 because it has a sparsity effect on the original transition tensor, properly handling Edman degradation is  
 88 critical to this decomposition. We feel this is the easiest way to do this while also factoring the  
 89 fluorophore loss effects into separate tensors.

## 90 **Everything all together**

91 Handling peptide detachment is simpler. We modify  $\mathcal{T}$  to be:

$$92 \quad \mathcal{T}_{vjui} = \begin{cases} (1-d)e\beta(\mathbf{i}, \mathbf{j}, p), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u \\ (1-d)(1-e)\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u = 0 \\ (1-d)(1-e) \left( \left(1 - \frac{i_{\bar{c}_u}}{\lambda_u}\right) \beta(\mathbf{i}, \mathbf{j}) + \left(\frac{i_{\bar{c}_u}}{\lambda_u}\right) \bar{\beta}(\mathbf{i}, \mathbf{j}, u) \right), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u > 0 \\ d, & \text{if } j_c = 0 \forall c \text{ and } v = \alpha \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

93 This creates a new “empty” state which can always be transitioned to with probability  $d$  of detachment.  
 94 The probability of avoiding this state is  $(1-d)$ . The functions  $\beta$  and  $\bar{\beta}$  are the same as before in (19)  
 95 and (20). The matrix vector multiplication step of the HMM forward algorithm has not changed from  
 96 (17). We can then construct a new tensor  $\mathcal{D}$  for peptide detachment which satisfies:

97 
$$\mathcal{D}_{whvk} = \begin{cases} 1 - d, & \text{if } \mathbf{h} = \mathbf{k} \text{ and } w = v \leq \alpha \\ d & \text{if } h_c = 0 \forall c \text{ and } w = \alpha + 1 \end{cases} \quad (25)$$

98 Then we find that:

99 
$$\mathcal{J}_{wlui} = \mathcal{D}_{wlvk} \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_C i_C}^{(C)} \quad (26)$$

100 Substituting into (17) with another multi-index  $\mathbf{m} = m_1 m_2 \dots m_C$  provides:

101 
$$\mathcal{f}_{wm}^{(t+1)} = \mathcal{O}_{ml}^{(t+1)} \mathcal{D}_{wlvk} \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_C i_C}^{(C)} \mathcal{f}_{ui}^{(t)} \quad (27)$$

102  $\mathcal{D}$  is clearly highly sparse, with two entries in each column of the original matrix in non-tensor form.

103 Thus,  $\mathcal{D}$  has no impact on the algorithmic complexity of this operation. Although  $\mathcal{D}$  and  $\mathcal{E}$  could be  
 104 combined to achieve this same algorithmic improvement, we found that this separation made our  
 105 model easier to reason about and work with.

## 106 **Transition matrix factoring conclusions**

107 One of the benefits of this approach to algorithmic complexity reduction is that this factorization  
 108 provides no loss to the theoretical accuracy of the forward algorithm. No theoretical approximations  
 109 were necessary, aside from the unavoidable differences in floating-point round-off errors. This allows  
 110 for highly accurate results with much more efficient runtime characteristics than a naïve  
 111 implementation.