## 2 Combining transition matrix factoring with HMM pruning

3 By making $r$ suitably small, HMM pruning can exhibit better algorithmic complexity than if we factor the

4 transition matrix. However, we believe it is much better to combine these algorithmic enhancements

5 (Fig 6). To do this, we need to switch into tensor notation, replacing our matrices and vectors with the

6 tensor equivalents we constructed previously. This yields:

$$\boldsymbol{f}_{vk}^{(t+1)} = \widehat{\boldsymbol{O}}_{kj}^{(t+1)} \widehat{\boldsymbol{T}}_{vjui}^{(t)} \boldsymbol{f}_{ui}^{(t)} \tag{32}$$

8 We also want to use the factorization from (26), using timestamp specific sub-tensors of each of the

9 factored pieces. The factorization of (26) becomes:

$$\widehat{\boldsymbol{T}}_{wlui}^{(t)} = \widehat{\boldsymbol{D}}_{wlvk}^{(t)} \widehat{\boldsymbol{\mathcal{E}}}_{vkuj}^{(t)} \widehat{\boldsymbol{B}}_{j_1 i_1}^{(t,1)} \widehat{\boldsymbol{B}}_{j_2 i_2}^{(t,2)} \dots \widehat{\boldsymbol{B}}_{j_C i_C}^{(t,C)} \tag{33}$$

11 Substituting into (32) gives:

$$\boldsymbol{f}_{wm}^{(t+1)} = \widehat{\boldsymbol{O}}_{ml}^{(t+1)} \widehat{\boldsymbol{D}}_{wlvk}^{(t)} \widehat{\boldsymbol{\mathcal{E}}}_{vkuj}^{(t)} \widehat{\boldsymbol{B}}_{j_1 i_1}^{(t,1)} \widehat{\boldsymbol{B}}_{j_2 i_2}^{(t,2)} \dots \widehat{\boldsymbol{B}}_{j_C i_C}^{(t,C)} \boldsymbol{f}_{ui}^{(t)} \tag{34}$$

13 Suppose we were to use standard sparse tensor multiplication techniques and carry this operation out

14 from right to left. Each tensor $\widehat{\boldsymbol{B}}_{j_c i_c}^{(t,c)}$ can introduce any entry of input (index $i_c$) into as many as $\Lambda_c$

15 indices of output (index $j_c$). The resulting computational complexity of the forward algorithm, even with

16 the given sparsity, is then $O\left(r\alpha T \prod_{c=1}^{C} \Lambda_c\right)$.

17 If we preprocess the computation, pruning each operation now from both directions, the algorithmic

18 complexity does not improve the way it does in the matrix case, although likely this would behave faster

19 in practice. The problem is that the pruning operation itself needs to determine which rows to

20 propagate forwards, which requires accessing every non-zero entry reachable in the forward direction.

21 Many of these values are later pruned in the backwards direction, so the computation itself has much

22 better sparsity, but the time to prune then dominates the algorithmic complexity result.

23     To improve this further, we add structure to the pruning of $\widehat{\boldsymbol{O}}^{(t)}$. Instead of keeping the $r$ largest values

24     in $\widehat{\boldsymbol{O}}^{(t)}$, we prune each index of $\widehat{\boldsymbol{O}}^{(t)}$ independently. For additional convenience, we limit each index to a

25     contiguous range of values. Then we let each index for any fluorophore color $c$ have $r_c$ values and allow

26     $\bar{r}$ values to index the number of amino acids. These simplifications may cause the pruning to be non-

27     optimal, but we accept this trade-off.

28     We can then prune our tensors using their known structures (for example, the tensors $\widehat{\boldsymbol{B}}^{(t,c)}$ correspond

29     to an upper triangular matrix). This time when we propagate the pruning results in both directions, the

30     time required is only $O(C^2)$ (the number of minimum and maximum indices to be propagated through

31     each tensor scales with $C$, as does the number of tensors to be pruned). For the runtime of the tensor

32     operations, consider each tensor individually. $\widehat{\boldsymbol{D}}^{(t)}$ and $\widehat{\boldsymbol{\mathcal{E}}}^{(t)}$ are both highly sparse, so they contribute a

33     constant modification to the number of rows or columns when propagating in either direction. $\widehat{\boldsymbol{D}}^{(t)}$

34     requires special handling. We track the detached state separately from the ordinary range, to avoid

35     unnecessarily including a large range of states which don't need to be.

36     Then each $\widehat{\boldsymbol{B}}^{(t,c)}$ operates on an independent index, and therefore can be considered on its own. This

37     tensor after pruning will have dimensions that are $O(r_c^2)$, and should have a constant effect on the

38     number of elements input vs output. Therefore, each of these tensors will require an algorithmic

39     complexity of $O\left(r_c \bar{r} \prod_{\tilde{c}=1}^{C} r_{\tilde{c}}\right)$. Bringing this all together we get an algorithmic complexity of

40     $O\left(C^2 + \bar{r}\left(\sum_{c=1}^{C} r_c\right)\left(\prod_{c=1}^{C} r_c\right)\right)$ for processing one timestep. The full forward algorithm then has a

41     complexity of $O\left(T\left(C^2 + \bar{r}\left(\sum_{c=1}^{C} r_c\right)\left(\prod_{c=1}^{C} r_c\right)\right)\right)$.

42     One remaining clarification is the manner of choosing $r_c$ and $\bar{r}$. In fact, these values should not be kept

43     constant; let us refer to the values for time $t$ as $r_c^{(t)}$ and $\bar{r}^{(t)}$. $\bar{r}^{(0)} = 1$ and $\bar{r}^{(t+1)} = \bar{r}^{(t)} + 1$, due to the

44     possibility of amino acid removal. These on average are proportional to $\alpha$. To get $r_c^{(t)}$, we keep all index

45   values where a fluorophore count of that value has the observed fluorescence intensity for color $c$ at

46   time $t$ within a specified confidence interval – perhaps within $3\sigma$ of the mean, where $\sigma$ is the standard

47   deviation of the distribution. These will necessarily be contiguous. The number of indices kept is then

48   $r_c^{(t)}$.

49   The standard deviation of a normal distribution scales with the square root of the intensity, and the

50   number of possible index values is limited by the total possible number of fluorophores of color $c$. It

51   follows that any removal of index values proportional to the standard deviation will satisfy $r_c^{(t)} < \gamma\sqrt{\Lambda_c}$

52   for some constant $\gamma$ dependent on the cutoff. Then the algorithmic complexity is given by

53   $O\left(T\left(C^2 + \alpha\left(\sum_{c=1}^{C}\sqrt{\Lambda_c}\right)\left(\prod_{c=1}^{C}\sqrt{\Lambda_c}\right)\right)\right).$

54   We chose a specific pruning cut-off by sweeping this parameter and balancing the experimental runtime

55   effects and the precision-recall curves which result from simulated data.