## 2   **HMM pruning**

3   Because the emission matrix is diagonal, it is equivalent to the diagonal part of its Singular Value

4   Decomposition (SVD), but with a reordering of its indices. This makes sparsification of this matrix

5   equivalent to the Eckart-Young-Mirsky theorem; we can keep the largest $r$ values for some chosen value

6   of $r$, and replace the rest of the matrix entries with zeros, having the minimum possible impact on the

7   spectral and Frobenius norms for the chosen value of $r$.

8   Furthermore, we can propagate this sparsification to the transition matrix. Consider the forward

9   algorithm, with $\boldsymbol{T}$ representing the transition matrix, and $\boldsymbol{O}^{(t)}$ representing the diagonal emission matrix

10   for time $t$. Then if $\boldsymbol{f}^{(t)}$ represents the vector of intermediate probabilities at time $t$, we have:

$$\boldsymbol{f}^{(t+1)} = \boldsymbol{O}^{(t+1)}\boldsymbol{T}\boldsymbol{f}^{(t)} \tag{28}$$

12   Now we sparsify each $\boldsymbol{O}^{(t)}$ as discussed above, to get a series of $\widehat{\boldsymbol{O}}^{(t)}$. This gives:

$$\boldsymbol{f}^{(t+1)} = \widehat{\boldsymbol{O}}^{(t+1)}\boldsymbol{T}\boldsymbol{f}^{(t)} \tag{29}$$

14   Note that we have many copies of $\boldsymbol{T}$, which are equal. For our next improvements we need these to be

15   different for each timestep, so we can rewrite (29) with $\boldsymbol{T}^{(t)}$ for each timestep $t$, giving

$$\boldsymbol{f}^{(t+1)} = \widehat{\boldsymbol{O}}^{(t+1)}\boldsymbol{T}^{(t)}\boldsymbol{f}^{(t)} \tag{30}$$

17   Here the values of many rows and columns of each $\boldsymbol{T}^{(t)}$ have been made unnecessary by the

18   sparsification of its neighboring $\widehat{\boldsymbol{O}}^{(t+1)}$ and $\widehat{\boldsymbol{O}}^{(t)}$, as any vector product with $\widehat{\boldsymbol{O}}^{(t)}$ will necessarily have

19   zeros except for the $r$ entries retained, such that we need only keep the corresponding $r$ columns of

20   $\boldsymbol{T}^{(t)}$. Similarly, any entry in the vector product with $\boldsymbol{T}^{(t)}$ which is not multiplied by one of the $r$ entries

21   retained in $\widehat{\boldsymbol{O}}^{(t+1)}$ is multiplied by zeros, and is thus unnecessary, so we need only keep the

22   corresponding $r$ rows of $\boldsymbol{T}^{(t)}$. Calling these approximations $\widehat{\boldsymbol{T}}^{(t)}$, we get

23
$$\boldsymbol{f}^{(t+1)} = \widehat{\boldsymbol{O}}^{(t+1)} \widehat{\boldsymbol{T}}^{(t)} \boldsymbol{f}^{(t)} \tag{31}$$

24 This allows significant sparsity to be used (Fig 5). Previously this formula would have been

25 $O(\alpha^2 T \prod_{c=1}^{C} \Lambda_c^2)$ to compute repeatedly across all timesteps, while this reduces the algorithmic

26 complexity to $O\left(\left(r^2 + \alpha\left(\prod_{c=1}^{C} \Lambda_c\right) \log r\right) T\right)$. Here $\alpha\left(\prod_{c=1}^{C} \Lambda_c\right) \log r$ represents the algorithmic cost of

27 determining the $r$ largest elements on the diagonal of the emission matrix using a priority queue. This

28 improvement is beyond what is possible in a more traditional usage of sparse matrix multiplication. For

29 sparse matrix multiplication, we would need to first multiply $\widehat{\boldsymbol{T}}^{(t)}$ by $\widehat{\boldsymbol{O}}^{(t)}$ or multiply $\widehat{\boldsymbol{O}}^{(t+1)}$ by $\widehat{\boldsymbol{T}}^{(t)}$. This

30 will only permit you to sparsify your operations on the rows or the columns of $\boldsymbol{T}$ but not both, giving a

31 complexity of $O\left(r\alpha T \prod_{c=1}^{C} \Lambda_c\right)$ (here sorting considerations are dominated by the rest of the formula

32 and can be omitted). While this is better than not using this inherent sparsity at all, preprocessing the

33 transition matrix in consideration of the emission matrices on either side gives better results in

34 algorithmic complexity.

35 In practice, we use a more complicated pruning scheme, as detailed in S4 Appendix.