

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

Amino acid sequence assignment from single molecule peptide sequencing data using a two-stage classifier

Matthew Beauregard Smith¹, Zack Booth Simpson², Edward M. Marcotte³

¹Oden Institute, The University of Texas at Austin, Austin, TX 78712

²Erisyon Inc., Austin, TX

³Department of Molecular Biosciences, The University of Texas at Austin, Austin, TX 78712

Correspondence: mbsmith93@utexas.edu, zack@erisyon.com, marcotte@utexas.edu

21 **Abstract**

22 We present a machine learning-based interpretive framework (*whatprot*) for analyzing single molecule
23 protein sequencing data produced by fluorosequencing, a recently developed proteomics technology
24 that determines sparse amino acid sequences for many individual peptide molecules in a highly
25 parallelized fashion [1] [2]. Whatprot uses Hidden Markov Models (HMMs) to represent the states of
26 each peptide undergoing the various chemical processes during fluorosequencing, and applies these in a
27 Bayesian classifier, in combination with pre-filtering by a k-Nearest Neighbors (kNN) classifier trained on
28 large volumes of simulated fluorosequencing data. We have found that by combining the HMM based
29 Bayesian classifier with the kNN pre-filter, we are able to retain the benefits of both, achieving both
30 tractable runtimes and acceptable precision and recall for identifying peptides and their parent proteins
31 from complex mixtures, outperforming the capabilities of either classifier on its own. Whatprot's hybrid
32 kNN-HMM approach enables the efficient interpretation of fluorosequencing data using a full proteome
33 reference database and should now also enable improved sequencing error rate estimates.

34 Introduction

35 Proteins are key components of living organisms, but their heterogenous chemical natures often
36 complicate their biochemical analyses, and consequently, the state of protein identification and
37 quantification methods (e. g., mass spectrometry, antibodies, affinity assays) has generally tended to lag
38 the remarkable progress exhibited by DNA and RNA sequencing technologies. However, improvements
39 to protein analyses could potentially directly inform better biological understanding and better translate
40 into biomedicine and clinical studies. Thus, the field of single molecule protein sequencing attempts to
41 apply concepts from DNA and RNA sequencing to protein analyses in order to take advantage of the high
42 parallelism, sensitivity, and throughput potentially offered by these approaches [3] [4] [5] [6].

43 Fluorosequencing is one such single-molecule protein sequencing technique inspired by methods used
44 for DNA and RNA [1] [2]. In fluorosequencing, proteins in a biological sample are denatured and cleaved
45 enzymatically into peptides. The researcher then chemically labels specific amino acid types, or
46 alternatively, specific post-translational modifications (PTMs), within each peptide with different
47 fluorescent dyes, then covalently attaches the peptides by their C-termini to the surface of a single-
48 molecule microscope imaging flow-cell (**Figure 1A**). Sequencing proceeds by alternating between
49 acquiring fluorescence microscopy images of the immobilized peptides and performing chemical
50 removal of the N-terminal-most amino acid from each peptide, using the classic Edman degradation
51 chemistry [7] [8] (**Figure 1B**). In this manner, the sequencing cycle (corresponding to amino acid
52 position) at which different fluorescent dyes are removed is measured on a molecule-by-molecule basis,
53 with these data collected in parallel for all the peptide molecules observed in the experiment (**Figure**
54 **1C**).

55 In theory, this process gives a direct readout of each peptide's amino acid sequence, at least for the
56 subset of labeled amino acids (**Figure 1D**), but in practice there are several complications because of the

57 single-molecule nature of this sequencing method. Single molecule fluorescence intensities are
 58 intrinsically noisy, arising from the repeated stochastic transitions of each individual dye molecule
 59 between ground state and excited state, making stoichiometric data inexact, particularly when there are
 60 large fluorophore counts. Typically, no more than 5-6 copies of the same amino acid, hence dye, are
 61 expected for average proteolytic peptide lengths, with the number of distinct colors (*i.e.*, fluorescent
 62 channels) set by the microscopy optics and available dyes, here assumed to be 5 or fewer. However,
 63 inevitably with any chemical process, some fluorophore labeling reactions fail to occur, and

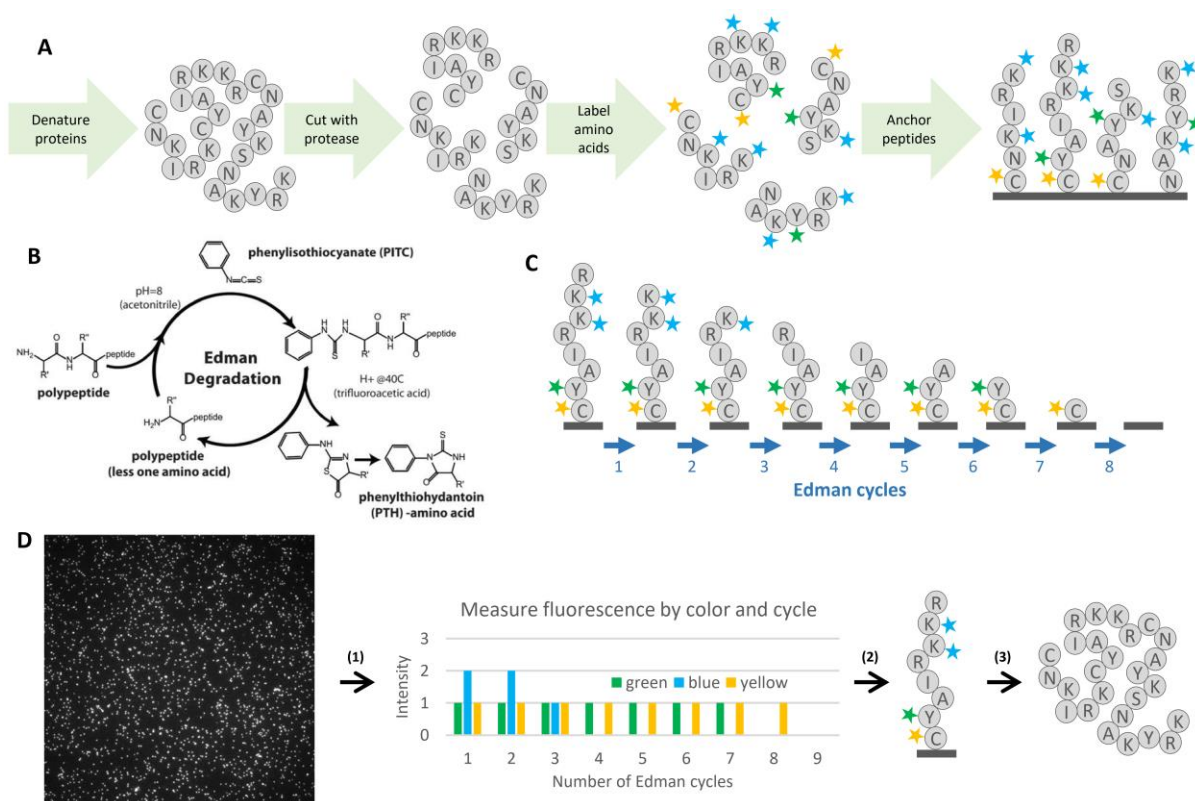


Figure 1. Overview of protein fluorosequencing. (A) illustration of the sample preparation process. Each grey circle represents an amino acid, and the letter in the circle corresponds to the standardized single letter amino acid codes. In the diagram, proteins are denatured, cleaved with protease, labeled with fluorescent dyes, and then labeled peptides are immobilized by their C-termini on the surface of a flow-cell. (B) The Edman degradation chemical reaction cycle, used to predictably remove one amino acid per cycle from each peptide. (C) For a given peptide, the sequencing process removes amino acids one at a time from the N-terminus, taking with them any attached fluorescent dyes. D: Major steps in computational data analysis include: (1) For each field of view, performing image analysis to extract fluorescence intensities for each spot (peptide) in each fluorescent channel across time steps (cycles), collating the fluorescence intensity data per spot across timesteps and colors. A vector of fluorescence intensities is produced, giving a floating-point value for every timestep and fluorophore color combination. (2) These raw sequencing intensity vectors (raw reads) must then be classified as particular peptides from a reference database. This step is the primary concern of this paper. (3)

64 photobleaching or chemical destruction can destroy fluorophores in the middle of a fluorosequencing
65 run. At some low rate, peptides may detach from the flow-cell during sequencing, and Edman
66 degradation can skip a cycle. These error rates, while individually small (approximately 5% each in
67 published analyses [2]), collectively add difficulty to peptide identification, necessitating computational
68 methods to process these data.

69 Currently, there are no published algorithms for mapping fluorosequencing reads to a reference
70 proteome to identify the proteins in a sample. The first analyses of fluorosequencing data used Monte
71 Carlo simulations to generate realistic simulated data as a guide for data interpretation and fitting of
72 experimental error rates [1] [2]. While this strategy did not scale well computationally to full proteomes,
73 it suggested that probabilistic modeling of the fluorosequencing process could provide a powerful
74 strategy for interpreting these data. In this paper, we explore the application of machine learning to
75 develop a classifier that correctly accounts for the characteristic fluorosequencing errors but is
76 computationally efficient enough to scale to the full human proteome.

77 Viewing this as a machine learning problem is challenging due to the large numbers of possible peptides
78 in many biological experiments. For example, in the human proteome, there are about 20,000 proteins,
79 which when processed with an amino-acid specific protease such as trypsin can correspond to hundreds
80 of thousands or even millions of distinct peptides, each of which can potentially vary due to post-
81 translational modifications or experiment-specific processing. This puts fluorosequencing data analysis
82 squarely in the realm of Extreme Classification problems, which are known to be challenging to handle in
83 practice [9].

84 To analyze these data, we took advantage of the ability to generate simulated fluorosequencing data
85 using Monte Carlo simulations [1] [2] to test k-Nearest Neighbors (kNN) classification and found it gave
86 results of poor quality but is able to scale efficiently to the full human proteome while maintaining

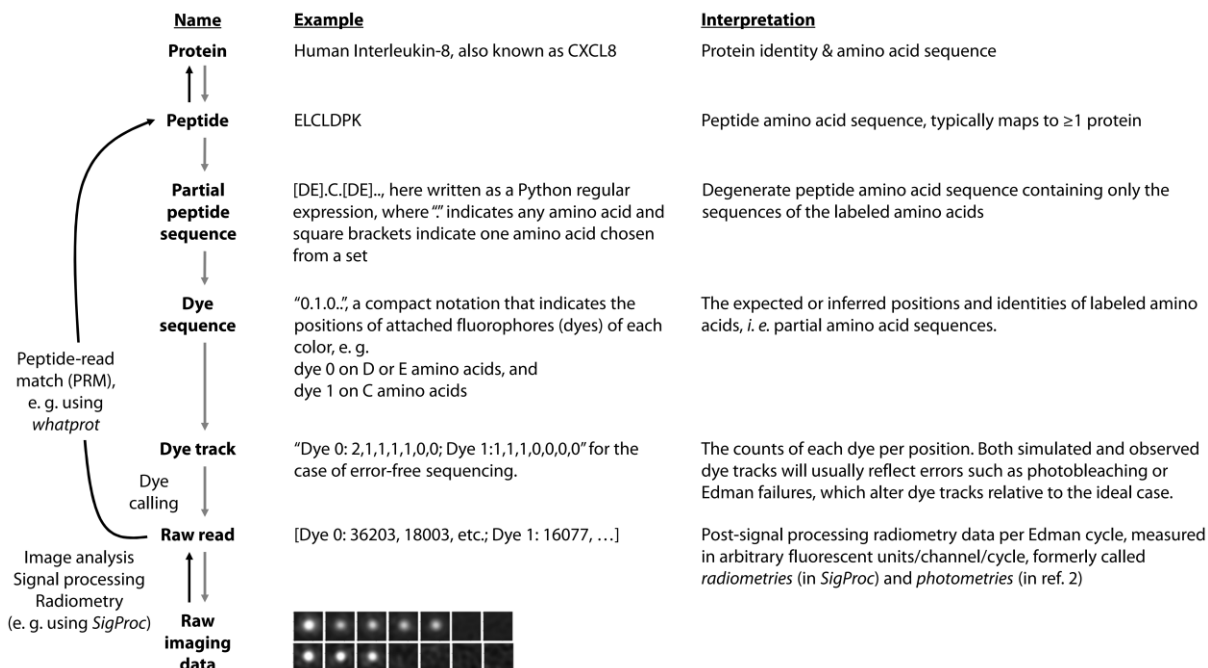


Figure 2. Nomenclature for different stages of fluorosequencing data analysis. The *whatprot* algorithm maps raw single-molecule protein sequencing reads to peptides and their parent proteins in the reference proteome (black arrows) by comparing experimental data (at bottom) to synthetic data generated using a Monte Carlo simulation (gray arrows).

87 reasonable runtimes. These initial explorations motivated the developments presented in this

88 manuscript, which focuses on the specific challenge of matching fluorosequencing reads to peptides

89 from a reference proteome (*peptide-read matching*).

90 Here, we propose a specialized classifier which combines heavily optimized Hidden Markov Models

91 (HMMs) to model the peptide chemical transformations during fluorosequencing, in combination with

92 kNN pre-classification to reduce runtime. We call this tool *whatprot*, compare it with kNN and a

93 classifier which uses HMMs without the kNN based runtime reduction, and demonstrate that the hybrid

94 HMM-kNN approach offers a powerful and scalable approach for interpreting protein fluorosequencing

95 data with the use of a reference proteome.

96 **Methods**

97 **Monte Carlo simulation**

98 To generate training and testing data typical of fluorosequencing experiments, we performed Monte

99 Carlo simulations based on the model and parameters described in [1] [2]. These parameters are the dye
100 loss rates p_c , which differ for each color c , the missing fluorophore rates m_c , the Edman cycle failure
101 rate e , the peptide detachment rate d , the average fluorophore intensity μ_c , and the standard deviation
102 of fluorophore intensity σ_c . We additionally model a background standard deviation σ'_c . Based on prior
103 estimates for the dye Atto647N ([1] [2]), we used the following values unless otherwise noted: $p_c =$
104 $.05$, $m_c = .07$, $e = .06$, $d = .05$, $\mu_c = 1.0$ (arbitrary rescaling of intensity values), $\sigma_c = 0.16$, $\sigma'_c =$
105 $.00667$. Although the code permits different values for different colors c , for our simulations, we
106 modeled each color of fluorophore with identical error values for simplicity.

107 An overview of the process with definitions for key terms is provided in **Figure 2**. We generate simulated
108 data in two formats. The first of these formats we refer to as a *dye track*, and it indicates the number of
109 remaining fluorophores of each color at each time step after considering sequencing errors. Thus, each
110 copy of one particular peptide sequence may give rise to a different specific dye track in a sequencing
111 experiment depending on the details of the labeling schemes and sequencing efficiencies. To simulate a
112 dye track, we randomly alter (with a pseudo random number generator) a representation of a dye
113 sequence in a series of timesteps, writing to memory the count of each color of fluorophore as we
114 progress until we reach a pre-set number of timesteps. In this simulation, we initially remove
115 fluorophores with a probability of m_c before beginning sequencing. We then additionally perform a
116 series of random events after logging fluorophore counts for each timestep: we remove the entire
117 peptide and all fluorophores with a probability of d to simulate peptide detachment from the flow cell,
118 we remove the last amino acid (and any attached fluorophore) with a probability of $(1 - e)$, and we
119 remove each fluorophore with a probability of p_c , where c is the color of the fluorophore, to simulate
120 fluorophore destruction. Each fluorophore count is stored as a two-byte numeric value.

121 The other format of data we consider is a *raw read*, which consists of radiometry data for each
122 fluorescent color and Edman cycle. Raw reads result experimentally from signal processing and

123 radiometry of the microscope imaging data from a fluorosequencing experiment. To simulate a raw
124 read, we first simulate a dye track, and then we convert each fluorophore count into a double-precision
125 floating point value indicating the fluorescent intensity. When we have a dye track entry indicating Λ_c
126 fluorophores for a given fluorophore color c , we sample a normal distribution with a mean of $\Lambda_c \mu_c$ and
127 a variance of $\sigma_c'^2 + \Lambda_c \sigma^2$. We perform this calculation for each channel at each time-step to simulate a
128 raw read.

129 These radiometry *raw reads* simulate the fluorescent intensity data we would expect to collect from
130 processing raw single molecule microscope images, a process currently performed for experimental data
131 using the algorithm *SigProc* (Part of Erisyon's tool *Plaster*, https://github.com/erisyon/plaster_v1), as in
132 [10] [11].

133 **Bayesian classification with HMMs**

134 Whatprot builds an independent HMM for each peptide in a provided reference proteome dataset. Each
135 state in this HMM represents a potential condition of the peptide, including the number of successfully
136 removed amino acids, and the combination of fluorophores which have not yet photobleached or been
137 destroyed by the chemical processing (**Figure 3**). Transition probabilities between these states can be
138 approximated using previously estimated success and failure rates of each step of protein
139 fluorosequencing.

140 We can use the HMM forward algorithm to associate a specific peptide to each *raw read* (a series of
141 observed fluorescence intensities over time and across different fluorescence channels). We obtain the
142 probability of the peptide given the raw read in two steps. First, we compute the HMM forward
143 algorithm using each possible peptide in the dataset to obtain the probability of the raw read given each
144 peptide. This uses the forward algorithm formula

$$145 \quad \mathbf{f}^{(t+1)} = \mathbf{O}^{(t+1)} \mathbf{T} \mathbf{f}^{(t)} \quad (1)$$

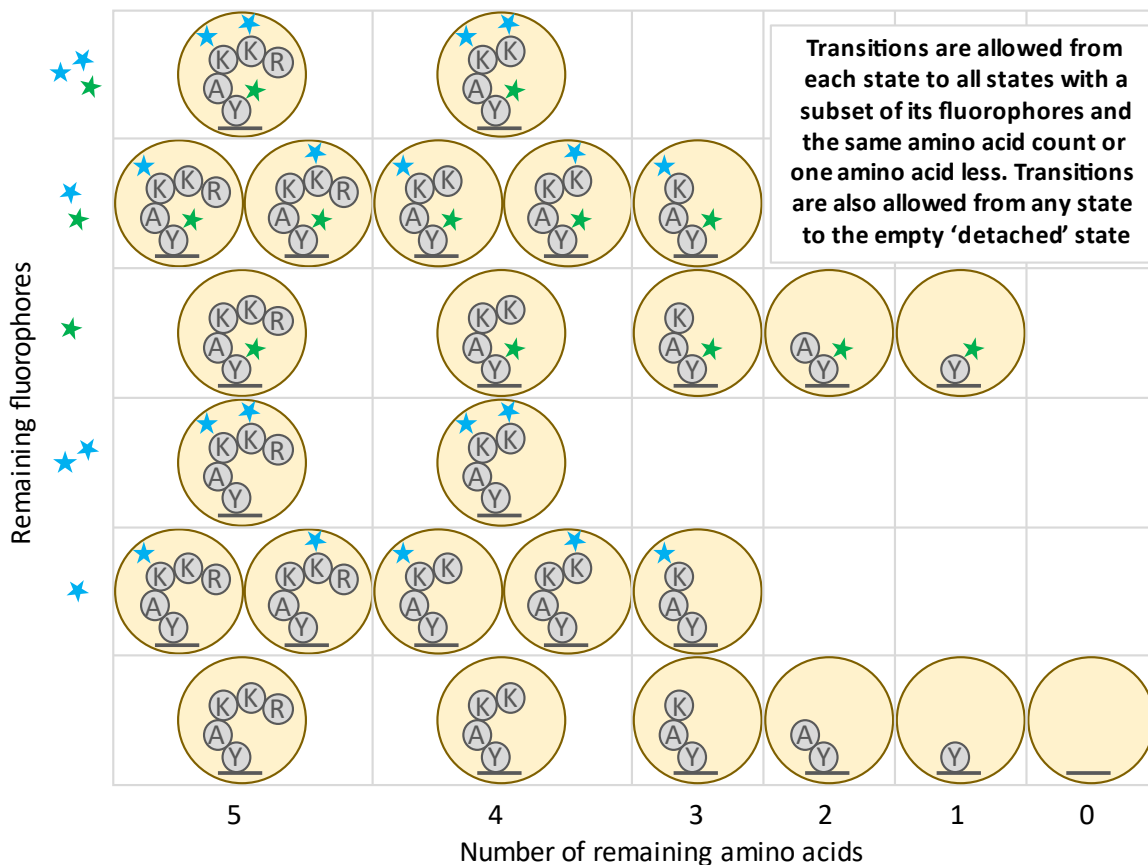


Figure 3. Illustration of the states and transitions of the HMM for an example peptide. For the amino acid sequence RKKAY, we illustrate the case where the lysine (K) residues are labeled with fluorescent dyes of one color (blue stars) and the tyrosine (Y) residue is labeled by a second color (green star).

146 Where $\mathbf{f}^{(t)}$ represents the cumulative probabilities for each state at timestep t , where $0 \leq t \leq T$, $\mathbf{O}^{(t)}$
 147 represents the diagonal emission matrix for the observation seen at timestep t , and \mathbf{T} represents the
 148 transition matrix which is the same at every timestep. The entries in each $\mathbf{f}^{(t)}$, $\mathbf{O}^{(t)}$, and in \mathbf{T} , represent
 149 the following probabilities:

$$150 \quad \mathbf{f}_i^{(t)} = p(Y_{1:t} = y_{1:t}, X_t = i | Z = z) \quad (2)$$

151 Where $Y_{1:T}$ are the random variables for the observations, $y_{1:T}$ are their true values, $X_{1:T}$ are the
 152 random variables for the state in the HMM and Z is the random variable representing the peptide, and z
 153 is a value it can take. We also have diagonal matrices $\mathbf{O}^{(t)}$ defined as:

154
$$\mathbf{O}_{ii}^{(t)} = p(Y_t = y_t | X_t = i, Z = z) \quad (3)$$

155 And :

156
$$\mathbf{T}_{ij} = p(X_{t+1} = i | X_t = j, Z = z) \quad (4)$$

157 We start from an initial state $\mathbf{f}^{(0)}$ which we compute by taking into account the missing fluorophore
158 rate m_c . Applying (1) repeatedly starting with the initial state $\mathbf{f}^{(0)}$ yields a value for $\mathbf{f}^{(T)}$, and we can
159 sum the entries to compute:

160
$$p(Y_{1:T} = y_{1:T} | Z = z) = \sum_i p(Y_{1:T} = y_{1:T}, X_T = i | Z = z) = \sum_i \mathbf{f}_i^{(T)} \quad (5)$$

161 Then, by using Bayesian inversion to normalize the data, we compute the probability of the peptide
162 given the raw read, as given by:

163
$$p(Z = z | Y_{1:T} = y_{1:T}) = \frac{p(Y_{1:T} | Z = z)p(Z = z)}{\sum_{\tilde{z}} p(Y_{1:T} = y_{1:T} | Z = \tilde{z})p(Z = \tilde{z})} \quad (6)$$

164 We implemented several algorithmic optimizations to this approach to reduce runtime. These included
165 reducing the number of states in the HMMs, factoring the HMMs' transition matrices into a product of
166 matrices with higher sparsity, pruning the HMM forward algorithm to consider only reasonably likely
167 states at each timestep, and combining the HMM classifier with a kNN pre-filter that can rapidly select a
168 short-list of candidate peptides for re-scoring by the HMM. We implemented the linear algebra and
169 tensor operations being performed in a manner that makes productive use of spatial and temporal
170 locality of reference. We describe these optimizations in more detail in the following sections and in the
171 supplemental **Appendices**.

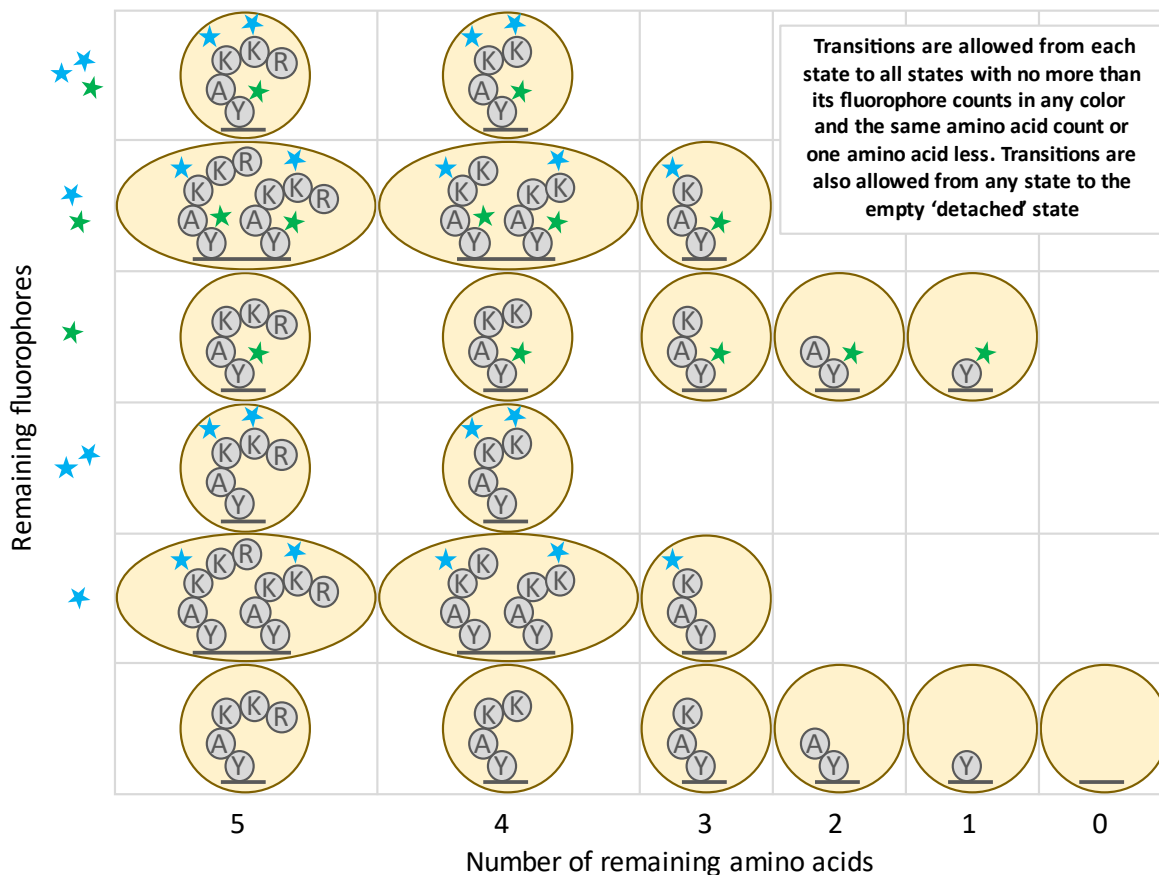


Figure 4. Illustration of HMM state space reduction for the peptide of Figure 3. States are combined that have both the same number of amino acids remaining and the same fluorophore counts for each color of fluorophore.

172 HMM state space reduction

173 We combine states of certain peptide conditions into more inclusive states in our model. For peptide
 174 conditions to be combined into these more inclusive states, they must have experienced the same
 175 number of *successful* Edman degradation events (so that they will have the same number of amino acids
 176 remaining), and they must have the same numbers of fluorophores of all colors. An example of the
 177 resulting HMM for a sample peptide is shown in **Figure 4**.

178 A similar state reduction to ours was previously described by Messina and colleagues in [12]. The
 179 reduction requires fluorophores to behave independently of each other, so that the status of one
 180 fluorophore is uncorrelated with the status of any other. While this is not true in practice due to FRET

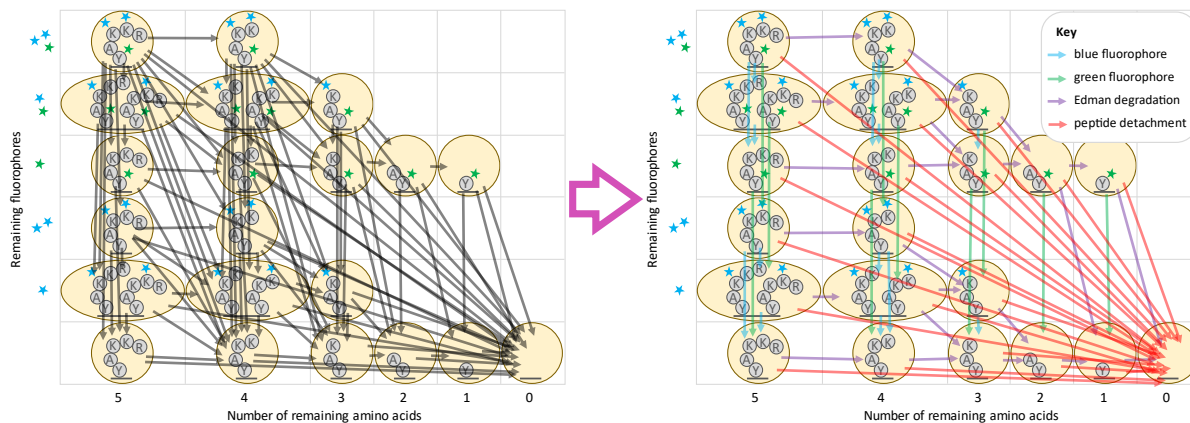


Figure 5. An illustration of the factoring of the transition matrix for the peptide from Figure 4. Note especially the reduction in the total number of transitions (arrows) when the transition matrix is factored. At left, black arrows represent non-zero entries in the unfactored transition matrix. At right, colored arrows (see key) represent non-zero entries in each of the matrices in the factored product. In both diagrams, arrows from a state to itself are omitted for visual clarity.

181 (Förster resonance energy transfer) and other dye-dye interactions, quantification of this effect in the
182 imaging conditions used for fluorosequencing suggests that these effects are negligible enough to ignore
183 [2]. The authors of [12] also require the fluorophores to be indistinguishable to reduce the numbers of
184 states. This is not true in our case because we use Edman degradation and because we use multiple
185 colors of fluorophores.

186 Nonetheless, we demonstrate in **Appendix A1** that despite these complications, this state space
187 reduction incurs no loss in the theoretical accuracy of the model. Further, we demonstrate that this
188 reduces the algorithmic complexity from what would otherwise be exponential with respect to the
189 number of fluorophores, to instead be tied to the product of the counts of fluorophores of each color.

190 **Transition matrix factoring**

191 In the HMM forward algorithm, a vector of probabilities with one value for each state in the HMM's
192 state space is repeatedly multiplied by a square transition matrix. This operation is the dominant
193 contribution to the algorithmic complexity of the HMM forward algorithm. Therefore, by making
194 multiplication by the transition matrix more algorithmically efficient, we can improve the theoretical
195 complexity of our computational pipeline.

196 We factor this transition matrix into a product of highly sparse matrices. This factorization is done by
197 creating a separate matrix for each independent effect under consideration, including loss of each color
198 of dye (where each color is factored separately), Edman degradation, and finally, peptide detachment
199 (**Figure 5**). As with the state space reduction, this optimization incurs no loss in the accuracy of the
200 model, and furthermore, these matrix factors, even in combination, are far sparser than the original
201 transition matrix when computed for larger peptides. This greater sparsity can be leveraged to achieve
202 superior algorithmic complexity results (see **Appendix A2**).

203 HMM pruning

204 Despite significant improvements in the algorithmic complexity of an HMM for one peptide given so far
205 from state space reduction and matrix factorization, performance can be improved if we consider

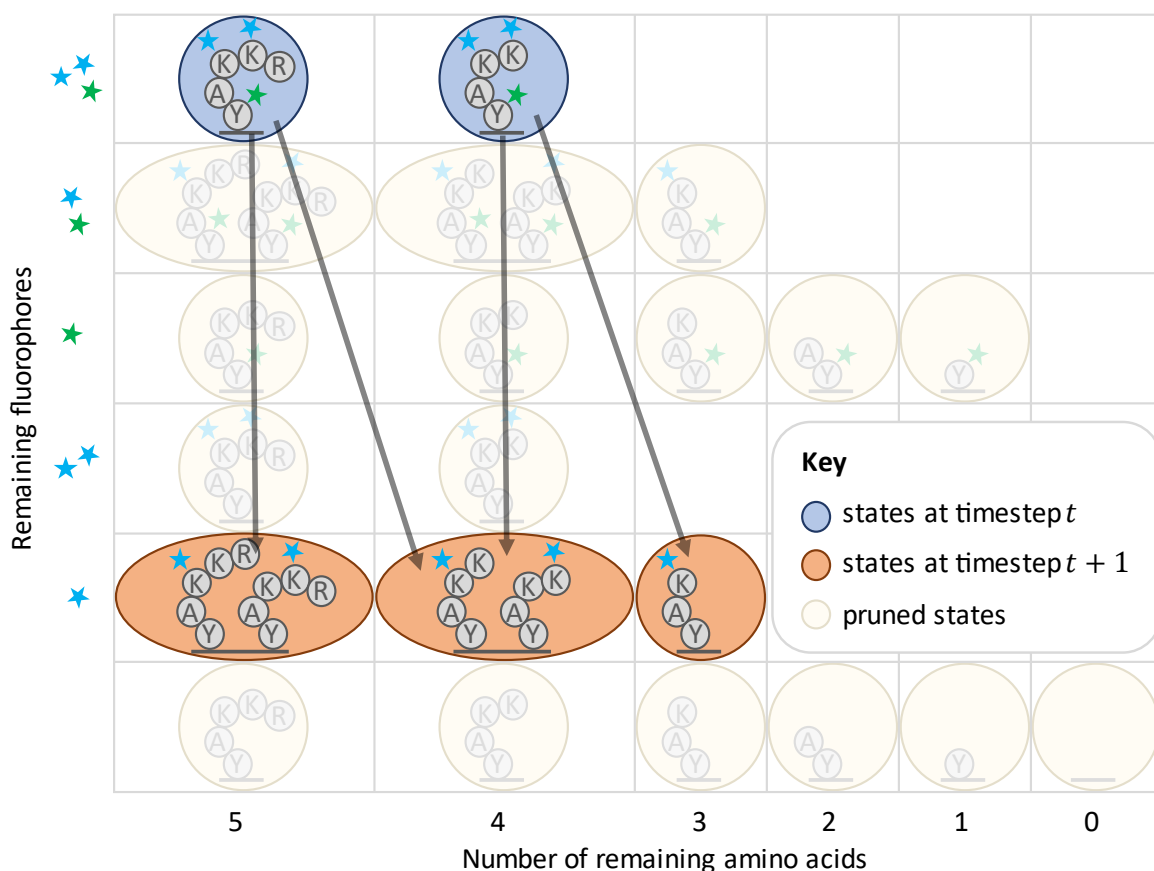


Figure 6. Illustration of the effects of HMM pruning for the peptide of Figure 5.

206 approximations. Intermediate computations contain mostly values close to zero, which will have
207 inconsequential impact on the result of the HMM forward algorithm. The most significant contribution
208 to this occurs for the HMM emission calculations. While there may be many states of a peptide which
209 have a significant probability of producing a particular observation value, in most states (particularly for
210 larger peptides) the observed value is extremely unlikely.

211 Emission computations can be viewed as multiplication by a diagonal matrix, different for each emission
212 in a raw read. The entries represent the probability of the indexed state producing the known emission
213 value for that timestep. We prune this matrix by setting anything below a threshold to zero, which
214 increases the sparsity of the matrix. Although use of a *naïve*, but standard, sparse matrix computational
215 scheme would reduce runtime, we show that better algorithmic complexity can be achieved with a more
216 complicated bi-directional approach in **Figure 6** and **Appendix A3**. While we did not implement this
217 approach precisely, a consideration of this effect served as inspiration for a technique combining
218 pruning with matrix factoring, as described next.

219 **Combining transition matrix factoring with HMM pruning**

220 Both transition matrix factoring and HMM pruning appear, at first glance, to be incompatible
221 improvements. These approaches can be combined, but the bi-directional sparse matrix computational
222 scheme introduces significant additional difficulties.

223 We view the various factored matrices as tensors and propagate contiguous blocks of indices forwards
224 and backwards before running the actual tensor operations to avoid unnecessary calculations (**Figure 7**).
225 Contiguous blocks of indices are needed because propagating lists of indices across the various factors
226 of the matrices has the same computational complexity as multiplying a vector by these matrices. This
227 may make the pruning operation seem less optimal in a sense, as some values that get pruned may be

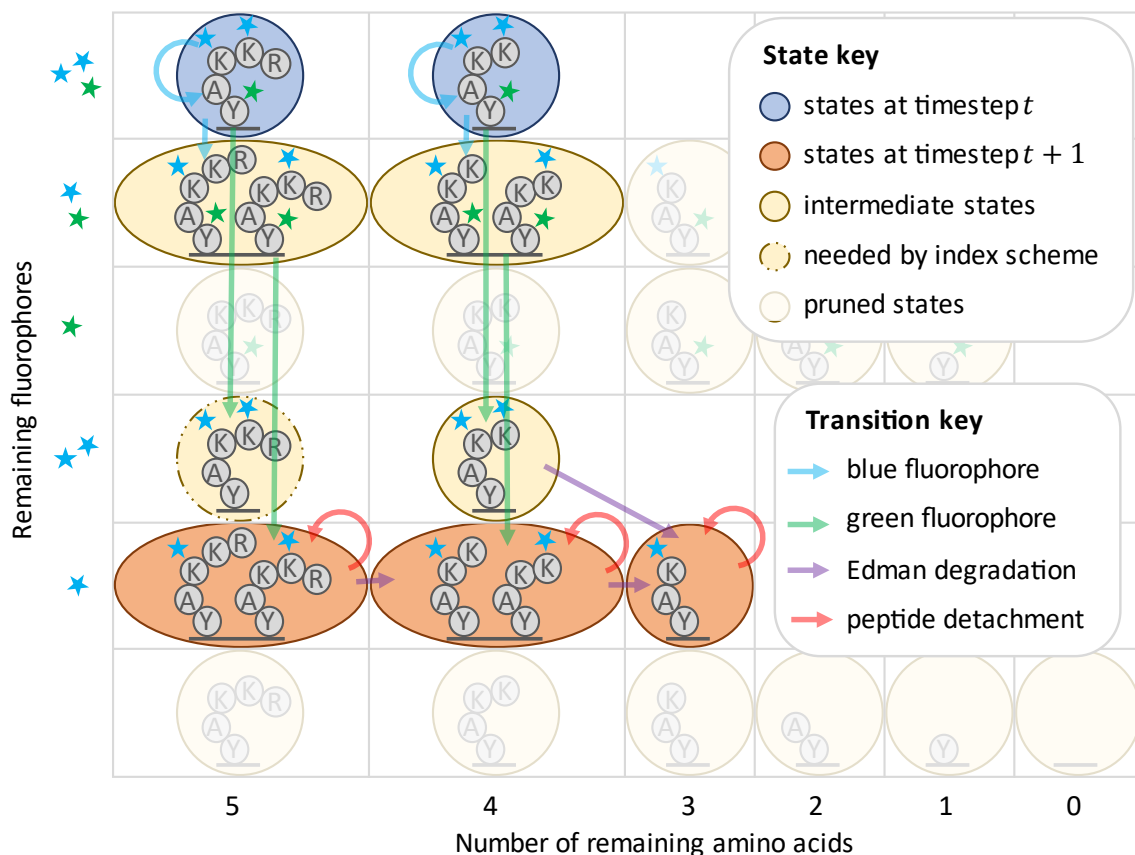


Figure 7. Illustration of HMM pruning combined with transition matrix factoring for the peptide of Figure 5. We emphasize that this is an anecdotal example; while there are more arrows here than in Figure 6, this strategy provides an improvement in asymptotic complexity, as described in Appendix A4 and shown in experiments with simulated data.

228 bigger than some that are kept due to this form of indexing. Nevertheless, we found the tradeoff to be
 229 favorable in practice (Appendix A4).

230 Of note, instead of pruning by the raw values, we prune all states such that the known emission value is
 231 outside of their pre-configured confidence interval. In this way we provide some confidence that the
 232 fraction of true data inadvertently zeroed out is negligible.

233 ***k*-Nearest Neighbors classification**

234 Most traditional machine learning classifiers have an algorithmic complexity which scales proportionally
 235 or worse to the number of classification categories. The Bayesian classifier we have so far described is
 236 no exception; each raw read must be compared against every peptide in the reference dataset to be

237 classified. There are many problems in biology which require large reference datasets, human proteomic
238 analysis being one example. The human proteome has 20,000 proteins, which when trypsinized
239 generate hundreds of thousands of peptides. Classification against these many categories is
240 computationally intractable with a fully Bayesian approach.

241 In contrast, the algorithmic complexity of kNN scales logarithmically with the number of training points
242 used. For this reason, tree-based methods are common in other Extreme Classification applications [9],
243 where similarly massive numbers of categories are under consideration. Unfortunately, the resulting
244 faster runtimes come at a significant cost; kNN often gives far worse results in practice than a more
245 rigorous Bayesian approach.

246 For purely kNN based classification, we simulate 1000 raw reads per peptide in the reference to create a
247 training dataset and put these into a custom KD-Tree implementation for fast and easily parallelizable
248 nearest neighbor lookups. We do not allow edits in our KD-Tree after it is built so as to allow parallelized
249 lookups to occur without any concern for locks or other common issues in parallel data structures. We
250 also reduce the memory footprint of the KD-Tree through an unusual compression scheme. For our
251 training data, we use dye tracks instead of raw read radiometry data; this alone reduces the memory
252 footprint of the KD-Tree by a factor of four (dye tracks have a two-byte numeric value for every
253 timestep/color combination, while radiometry data have an eight-byte double-precision floating point
254 number). But this allows another further compression technique; we find all identical dye tracks and
255 merge them into one entry. With these dye tracks entries in the KD-Tree we store lists of peptides that
256 produced the dye track when we simulated our training data, along with how many times each peptide
257 produced that dye track.

258 To classify an unknown raw read, the k nearest dye track neighbors to a raw read query are retrieved.
259 These neighbors then vote on a classification, with votes weighted using a Gaussian kernel function,

260 $\exp\left(-\frac{\delta^2}{2\sigma_{kNN}^2}\right)$, where δ is the Euclidean distance between the query raw read and the neighbor, and
261 σ_{kNN} is a parameter of the algorithm. A neighbor is also weighted proportionally to the number of times
262 it occurred as a simulation result and will split its voting weight among all of the peptides that produced
263 that dye track proportionally to the numbers of times each peptide produced the dye track during
264 simulation of training data.

265 Once voting is complete, the highest weighted peptide is then selected as the classification, with its
266 classification score given as a fraction of its raw score over the total of all the raw scores. We have
267 explored multiple choices of k and σ values to optimize the performance.

268 **Hybridizing kNN with Bayesian HMM classification**

269 To combine the computational efficiency of kNN with the accuracy of the HMM model, we defined a
270 classifier which hybridizes these two disparate methods. We use a kNN classifier to reduce the reference
271 dataset, for each raw read, down to a smaller shortlist of candidate peptides. These candidates can then
272 be used in the Bayesian classifier by building HMMs to compare them against the specific raw read.
273 While this can result in the true most likely peptide not being in the shortlist and therefore not being
274 selected by this hybrid classifier, with a sufficiently long shortlist this is highly unlikely. A larger problem
275 is in performing Bayes' rule, as in (6). An exact formula for Bayes' rule requires an exhaustive set of
276 probability values for every potential outcome, which are summed in the denominator. Avoiding
277 determining every probability makes this impossible. Instead, we can estimate Bayes' rule as follows:

$$278 \quad p(Z = z | Y_{1:T} = y_{1:T}) = \frac{p(Y_{1:T} | Z = z)p(Z = z)}{\sum_{\tilde{z} \in \zeta_h} p(Y_{1:T} = y_{1:T} | Z = \tilde{z})p(Z = \tilde{z})}$$

279 Where ζ_h is the set of up to h peptides selected by the kNN method; we require $z \in \zeta_h$. Although we
280 lose theoretical guarantees of optimal accuracy given the model, this change provides a considerable
281 improvement to the algorithmic complexity. The algorithmic complexity to classify one raw read using a

282 fully Bayesian approach is $O(RW)$, where R is the number of peptides in the reference dataset, and W
283 is the average amount of work needed to run an HMM for one peptide fluorophore combination. In
284 comparison, with the hybridized classifier, the algorithmic complexity is $O(\log(RQ) + hW)$ where Q is
285 the number of raw reads in the training dataset simulated for each possible peptide.

286 We chose specific values for h , σ , and k , by comparing the runtime and PR curves on simulated datasets.

287 **Maintaining spatial locality of reference**

288 Spatial and temporal locality of reference is the tendency of some computer programs to access nearby
289 data points at similar times. Modern CPUs are designed to make this extremely efficient through multi-
290 level batch caching schemes which cache data from RAM that is nearby a memory address being
291 accessed, so that nearby data can be read more quickly. Programs which exploit this in read-write
292 intensive pieces of code can often achieve significant runtime acceleration compared to programs which
293 do not.

294 We wrote highly optimized kernel functions to perform our structured matrix/tensor operations which
295 exploited the sparse nature of the problem while also iterating over elements in what we believe to be
296 an optimal or near optimal fashion for most computer architectures. We believe this provided
297 considerable improvements in performance, though this has not been rigorously tested.

298 **Results**

299 We simulated the fluorosequencing of peptides to obtain labeled training and testing data (**Figure 2**).
300 We generated several datasets, each with a randomized subset of the proteins in the human proteome.
301 We selected 20 proteins (.1% of the human proteome), 206 proteins (1% of human proteome), 2,065
302 proteins (10%), and 20,659 proteins (the full human proteome). We repeated this randomized selection
303 scheme to examine several protease and labeling schemes. These were (1) trypsin (which cleaves after

304 lysine (K) and arginine (R) amino acids) with fluorescent labels for aspartate (D) and glutamate (E) (these
305 share a fluorophore color due to their equivalent reactivities), cysteine (C), and tyrosine (Y), (2)
306 cyanogen bromide (which cleaves after methionine (M) amino acids) with D/E, C, Y, and K, (3) EndoPRO
307 protease (which cleaves after alanine (A) and proline (P) amino acids) with D/E, C, and Y, (4) EndoPRO
308 with D/E, C, Y, and K, (5) EndoPRO with D/E, C, Y, K, and histidine (H). Thus, in the schemes examined, of
309 the 20 canonical amino acid types found in most proteins, either one or two were recognized by the
310 protease and up to 6 additional amino acids were labeled by fluorescent dyes.

311 These databases of peptides were used to generate databases of idealized *dye sequences*, *dye tracks*,
312 and *raw reads*, used for training and testing purposes for the various models. For our test data for each
313 dataset, we generated 10,000 raw sequencing reads by randomly selecting peptides with replacement
314 from the dataset and simulating sequencing on them using the methods described in the *Monte Carlo*
315 *simulation* section of this paper. For both dye tracks and simulated fluorescent intensity measurements,
316 results where there were zero fluorophores throughout sequencing were discarded, as these would fail
317 to be observed in an actual sequencing experiment.

318 We collected and compared runtime data and precision-recall curves for several different purposes.
319 With the trypsinized 3-color dataset, we performed a parameter sweep of the pruning cut-off for the
320 HMM Bayesian classifier (**Figure B1**). Losses in precision and recall performance were negligible for cut-
321 off values of 5 and greater, though the precision recall curves grew worse at smaller values. Runtimes
322 shrank rapidly as the cut-offs were decreased. The pruning cut-off parameter sweep was also performed
323 on the 20-protein cyanogen bromide dataset (**Figure B2**). We saw that in this second dataset, runtime
324 improvements for lower cut-offs were even more extreme; a speed-up factor of about 1000 could be
325 achieved with minimal effects on the precision recall plots. From these two simulations, we chose a
326 cutoff value of 5 as providing the optimal trade-off between runtime and precision recall performance.

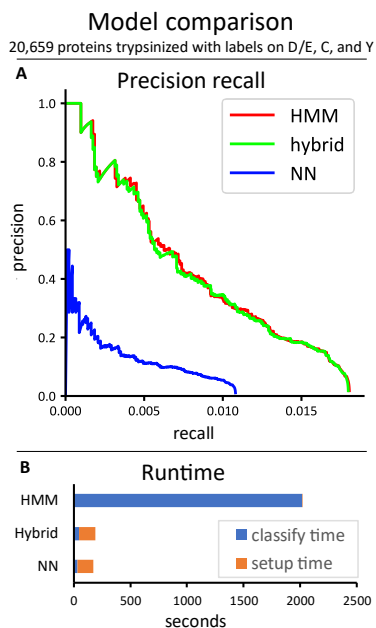


Figure 8. Comparison of the HMM (Bayesian), hybrid, and NN classifiers on a dataset of 10K reads from peptides chosen randomly from all 20,659 human proteins trypsinized and labeled on D/E, C, and Y. (A) The precision recall curves. (B) Runtimes.

On the trypsinized dataset (full human proteome), we also swept the k and σ_{kNN} parameters of the NN classifier (**Figures B3, B4**).

Here large values of k introduce modest reductions in precision recall performance, while the model is extremely sensitive to the selection of σ_{kNN} . Based on this analysis, we suggest that good choices of these parameters are $k = 10$ and $\sigma_{kNN} = 0.5$.

We swept all parameters of the hybrid classifier for the trypsinized dataset as well (hybrid h parameter, k , σ_{kNN} , and cut-off) (**Figures B5-8**). Here, the HMM cut-off parameter had less impact on runtime than for the pure HMM Bayesian classifier, but we still found a cut-off of 5 to be optimal. Higher values of k improved precision recall performance for the hybrid model, contrary to the results

340 of the NN classifier on its own, and we therefore suggest setting k
 341 to 10000. σ_{kNN} had minimal impact of any kind, in contrast to its
 342 significant impact on the precision recall of the NN classifier; we
 343 nevertheless chose to set it to 0.5 in light of the data from
 344 parameter tuning for the NN classifier on its own. We also found
 345 that higher values of h improved performance, though the impact
 346 plateaus after h of about 1000, and we used that value for later
 347 experiments.
 348 After tuning parameters, we compared the performance of the
 349 different classifiers when applied to 10,000 simulated
 350 fluorosequencing reads of peptides drawn randomly from all tryptic

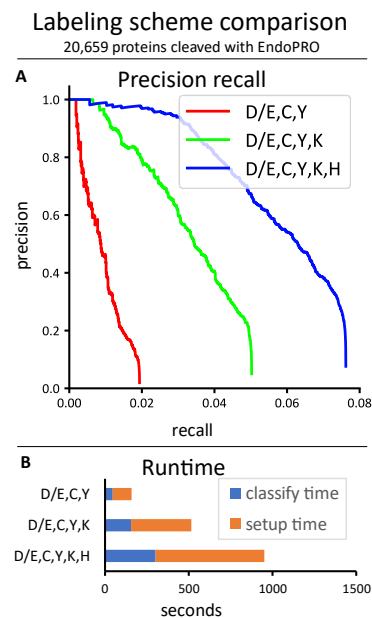


Figure 9. Comparison of the hybrid classifier on a dataset of 10K reads from peptides chosen randomly from all 20,659 human proteins cleaved with EndoPRO and labeled with three different labeling strategies. (A) The precision recall curves. (B) Runtimes.

351 peptides in the human proteome (**Figure 8**). The hybrid classifier
352 achieved similar precision recall curves to the Bayesian HMM
353 classifier, which was much better than the precision recall curve of
354 the NN classifier. The hybrid classifier also achieved runtimes of the
355 same order of magnitude as the NN classifier, which was
356 significantly faster than the runtime of the Bayesian HMM classifier.

357 We also studied how the number of fluorophore colors affected the
358 runtime and precision/recall of the hybrid classifier (**Figure 9**). We
359 found that improvements in precision recall were possible with
360 each additional color of fluorophore, but this did come at the cost
361 of longer runtimes.

362 We also investigated the effect of varying sizes of reference
363 proteomes on the hybrid classifier's performance, using the three

364 color trypsinized dataset (**Figure 10**). We found that significantly better performance was possible when
365 the reference database was smaller.

366 The precision/recall curves plotted above (**Figures 8-10**) show the actual precision/recall scores based
367 on data with known peptide classifications. When working with real data this will typically not be
368 possible, because the real classifications will not be known. It is therefore important that the assignment
369 probabilities produced by the classifier be well-calibrated, so that an estimate of the precision/recall (or
370 as is more often the case in protein mass spectrometry, the false discovery rate (FDR)) can be computed
371 in the absence of known labels. We verified that the probabilities output by the hybrid HMM classifier
372 were indeed well-calibrated relative to the true assignment probabilities (**Figure 11A**). This in turn
373 allowed us to compute a predicted precision/recall curve assuming that each classification is fractionally

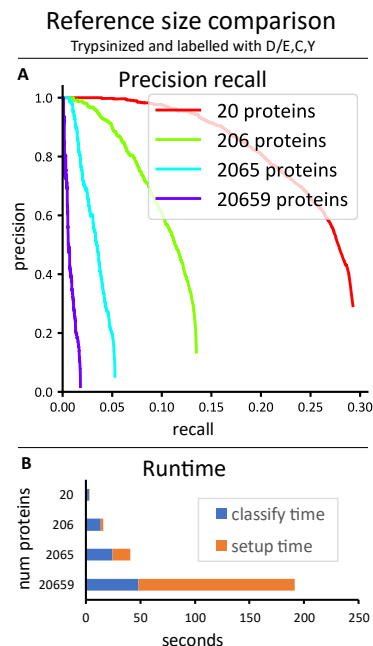


Figure 10. Comparison of the hybrid classifier across datasets of 10K reads each from peptides chosen randomly from different numbers of proteins treated with the same protease and labeling scheme. (A) The precision recall curves. (B) Runtimes.

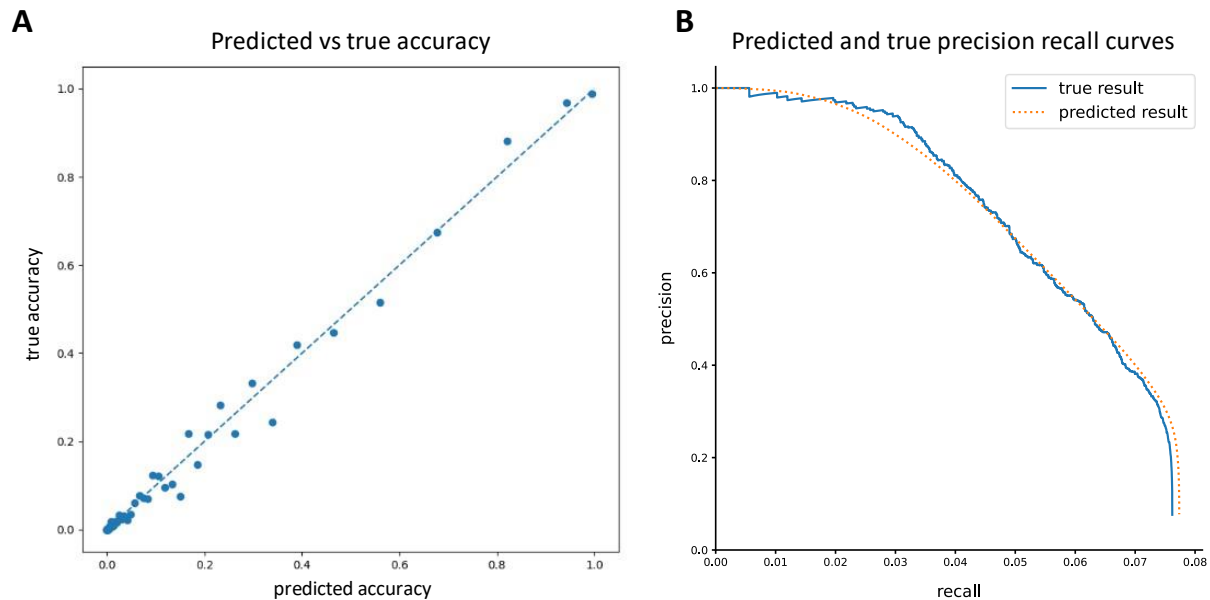


Figure 11. Analysis of the accuracy of probability estimates given as scores by the classifier. Based on 10K reads from peptides chosen randomly from all 20,659 human proteins cleaved with EndoPRO and labeled on D/E,C,Y,K,H. **(A)** Classification results were sorted by their predicted accuracy scores, and then equally distributed between 100 buckets. The average predicted and true accuracy scores were then computed for each bucket and plotted. **(B)** The true result precision/recall curve was computed as normal, while the predicted result precision/recall curve was plotted assuming each classification was fractionally correct according to its predicted accuracy score.

374 correct with a probability given by its classification score. A comparison of this predicted P/R with the
375 actual precision/recall curve for the same set of reads shows excellent agreement (**Figure 11B**).

376 Whatprot specifically attempts to assign each raw fluorosequencing read to one or more peptides from
377 the reference database, *i.e.*, to identify and score *peptide-read matches* (PRMs), a process highly
378 analogous to analytical interpretation of shotgun mass spectrometry (MS) proteomics data in which a
379 key step is comparing experimental peptide mass spectra to a reference proteome (finding *peptide-*
380 *spectral matches*, or PSMs [13]). However, observing multiple reads mapping to the same peptide will
381 tend to increase the confidence that peptide is present in the sample, just as observing multiple
382 peptides from the same protein will similarly increase confidence in that protein being present. Thus, we
383 asked if considering the PRMs collectively led to performance increases for identifying peptides and
384 proteins.

385 As shown in (Figure 12), proteins can
386 be identified correctly at much higher
387 rates than peptides, which are similarly
388 identified at higher rates than
389 individual reads. In fact, provided that
390 a protein possesses some well-
391 identified peptides that are unique, it
392 can typically be identified with very
393 high accuracy. For this test, we used a
394 very simple protein inference scheme.
395 First each peptide was scored to the
396 maximum score of all reads identifying
397 it, while penalizing reads which
398 identified more than one peptide
399 (dividing by n if n peptides were
400 identified). Second each protein was
401 scored as the maximum score of all peptides it contains, penalizing peptides which are associated with
402 more than one protein (again dividing by n if n proteins were associated). However, the problem of
403 integrating peptide level observations to protein observations has been studied extensively for MS [14]
404 [15], and it is likely that these techniques will offer similarly strong interpretive power to the case of
405 single molecule protein sequencing.

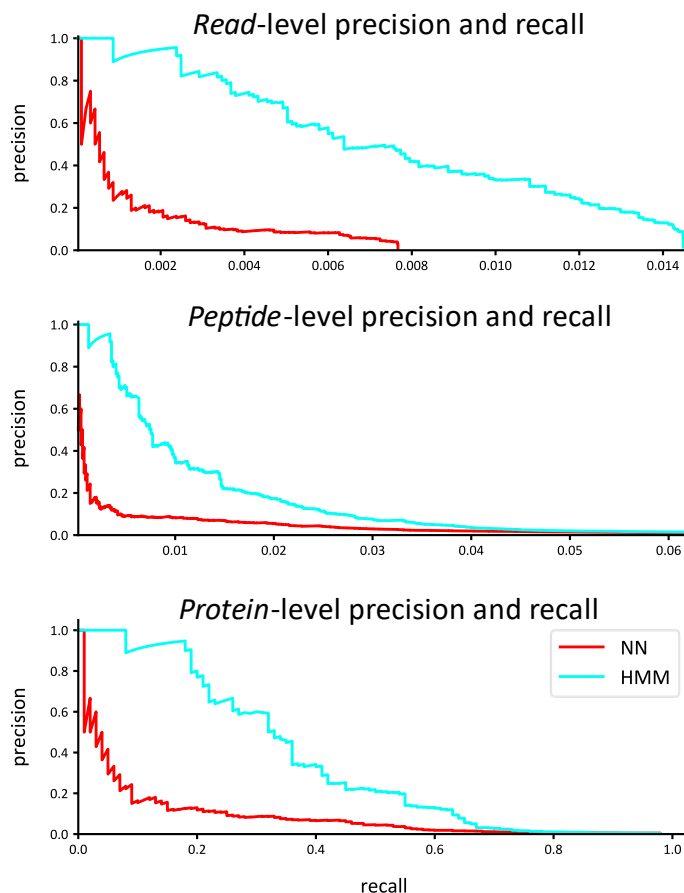


Figure 12. Precision and recall are improved for proteins by integrating identifications across peptides. The example shows 10K reads from peptides derived from 100 proteins randomly selected from the human proteome, considering trypsin digestion and labels on D/E, C, and Y.

406 Discussion

407 We developed an HMM for interpreting single molecule protein fluorosequencing data and showed that
408 a hybrid HMM/kNN model can achieve a high precision and recall comparable the HMM alone while
409 maintaining a runtime comparable to the much faster kNN.

410 It is worth emphasizing that these analyses were performed on datasets of 10,000 raw fluorosequencing
411 reads. In practice, users will likely want to analyze millions to billions of reads, so runs that completed in
412 a seemingly reasonable amount of time might still be intractable in these scenarios with larger datasets,
413 or at a minimum require computing clusters with high parallelization. For analyzing the current datasets
414 in the runtime charts (**Figures 8-10**), note that the *blue* part of the bar graphs indicates the classify time,
415 which will scale with the number of reads being classified (if all else remains equal), and the *orange* part
416 of the graph indicates setup time, which should remain constant regardless of the number of reads
417 (though it changes depending on the model and the size of the reference set).

418 It is also interesting that the HMM pruning operation is more necessary with longer peptides and more
419 colors of fluorophores; with the trypsinized dataset labeling D/E, C, and Y, omitting pruning had little
420 consequence, but in moving to cyanogen bromide with D/E, C, Y, and K, we observed a runtime speedup
421 of about 1000-fold.

422 Finally, our data demonstrates that with a proper selection of parameter values, the hybrid model can
423 achieve precision and recall performance virtually identical to the HMM Bayesian approach alone, while
424 providing those results in a fraction of the time. Similarly, the pruning operation employed in the HMMs
425 has no noticeable positive or negative effect on the precision recall curves while providing a
426 considerable improvement in runtime performance.

427 A number of analytical techniques common in related fields were not explored. In tandem mass
428 spectrometry (MS/MS), peptide spectral mapping is typically done either through database lookups
429 and/or the use of simulated outcomes. Simulated mass spectra, consisting of ion pairs for the C- and N-
430 terminal fragments for each potential breakage point, can be compared to the real data collected from
431 the instrument [13]. Recent advances have been achieved by using deep learning to predict
432 fragmentation behavior with higher quality than is possible with more traditional methods [16][24].
433 While we use the notion of matching fluorosequencing reads to a reference database, the specific
434 algorithms are distinct.

435 Nevertheless, of possible relevance from the field of MS/MS is the analysis of the false discovery rate
436 (FDR) [25][26]. The FDR is affected by two distinct sources: a peptide may be misattributed to the wrong
437 peptide, even when the true peptide is present in the reference dataset, and MS/MS datasets contain
438 significant amounts of modified peptides or contaminants, whose spectra may be mistakenly assigned to
439 peptides in the reference set [27]. FDR is typically evaluated using a decoy database, such as is
440 generated using reversed proteins from the target database. The FDR can then be set by referring to the
441 number of hits in the decoy database given a particular score, as the decoy database is designed such
442 that it should in theory never have hits for the biological sample being analyzed [17]. While an estimate
443 of FDR based in theoretical analysis of the problem could find the misattribution rate of true peptides,
444 even this estimate would be incomplete, because there are errors in mass spectra of peptides that
445 cannot be accounted for by existing theory; furthermore, any effect of modifications or contaminants
446 would likely be omitted.

447 The utility of a similar decoy database strategy for estimating FDR for fluorosequencing is unknown and
448 remains to be established. We note however that, due to the rigorous probabilistic nature of our
449 analysis, a reasonable estimate of FDR can be performed by subtracting the sum of PRM scores from the
450 number of PRMs. This is the same as one minus the precision in a predicted precision/recall curve, and

451 the proximity of our predicted precision/recall curve to the real curve for a known dataset demonstrates
452 the feasibility of this approach (**Figure 11**). This analysis likely fails to account for the contributions of
453 modifications and contaminants. We therefore plan to explore this problem more extensively in future
454 work.

455 We also considered techniques for DNA sequence reconstruction. In general, DNA sequencing provides
456 *de novo* sequence reconstructions and does not use reference database matching, and therefore is not a
457 good model for fluorosequencing. Nevertheless, base calling strategies may have some relevance. For
458 example, methods for base-calling from conventional (e.g. Illumina style) DNA sequencing are
459 straightforward [18] [19], and although errors occur, they are rare [20]. Analysis of errors in DNA
460 sequencing is typically performed using multiple sequence alignment or k-mer based methods [21].
461 Because the error rates are typically much lower in DNA sequencing than in fluorosequencing, we
462 believe existing software is unlikely to be effective in this new domain.

463 Nanopore DNA sequence analysis methods could also be considered. Nanopores, similar to
464 fluorosequencing, deal with single molecule data and the concomitant statistical noise that process
465 involves. However, nanopore data is on a real time continuum, with a DNA fragment which may move
466 through the nanopore at variable rates during sequencing. Base-calling, the assignment of nucleic acid
467 bases to chunks of sequencing information, is again the step most analogous to fluorosequencing. State
468 of the art base-calling methods for nanopore sequencing typically use either HMMs or recurrent neural
469 networks (RNNs). Comparisons of existing approaches suggest that RNNs slightly outperform HMMs in
470 this domain [22]. While this suggests that RNNs are worth exploring for fluorosequencing data, we have
471 avoided this approach for two reasons. First, RNNs are a deep learning technique, which invariably
472 requires access to massive amounts of data; this is not currently feasible with fluorosequencing unless
473 that data is simulated. Second, our approach suggests the possibility of direct estimation of parameters

474 using some variant of the Baum-Welch algorithm adapted to our use case, which we believe would be
475 significantly more difficult in an RNN based approach [23].

476 **Conclusions**

477 We have developed a powerful computational tool for the analysis of protein fluorosequencing data,
478 which significantly increases the complexity of applications available to this new technology. This tool
479 includes critically important optimizations which make our approach feasible in practice.

480 In future work we plan to implement a variation of the Baum-Welch algorithm to fit the parameters to
481 data from a known peptide. We also wish to explore peptide and protein inference methods using
482 peptide data classified using these methods. We may also explore *de novo* recognition of labels without
483 use of a reference database.

484 **Acknowledgements**

485 Code is available at <https://github.com/marcottelab/whatprot>. The authors gratefully acknowledge
486 Jagannath Swaminathan, Angela Bardo, Brendan Floyd, Daniel Weaver, and Eric Anslyn for helpful
487 guidance and discussion throughout the course of this project. M.B.S. acknowledges support from a
488 Computational Sciences, Engineering, and Mathematics graduate program fellowship. E.M.M.
489 acknowledges support from Erisyon, Inc., the National Institute of General Medical Sciences
490 (R35GM122480), the National Institute of Child Health and Human Development (HD085901), and the
491 Welch Foundation (F-1515).

492 **Disclosures**

493 E.M.M. and Z.B.S. are co-founders and shareholders of Erisyon, Inc., and are co-inventors on granted
494 patents or pending patent applications related to single-molecule protein sequencing. E.M.M. serves on
495 the scientific advisory board.

496 References

- 497 **1.** Swaminathan J, Boulgakov AA, Marcotte EM. A theoretical justification for single molecule
498 protein sequencing. *PLoS Computational Biology*. 2015;11(2):1-17.
- 499 **2.** Swaminathan J, Boulgakov A, Hernandez ET, Bardo AM, Bachman JL, Marotta J, *et al.* Highly
500 parallel single-molecule identification of proteins in zeptomole-scale mixtures. *Nature*
501 *Biotechnology*. 2018;36(11):1076-1082.
- 502 **3.** Callahan N, Tullman J, Kelman Z, Marino J. Strategies for development of a next-generation
503 protein sequencing platform. *Trends in Biochemical Sciences*. 2020;45(1):76-89.
- 504 **4.** Floyd BM, Marcotte EM. Protein sequencing, one molecule at a time. *Annual Review of*
505 *Biophysics*. 2022;51:181-200.
- 506 **5.** Restrepo-Pérez L, Joo C, Dekker C. Paving the way to single-molecule protein sequencing. *Nature*
507 *Nanotechnology*. 2018;13:786-796.
- 508 **6.** Vistain LF, Tay S. Single-cell proteomics. *Trends in Biochemical Sciences*. 2021;46(8):661-672.
- 509 **7.** Edman P, Method for determination of the amino acid sequence in peptides. *Acta Chemica*
510 *Scandinavica*. 1950;4:283-293.
- 511 **8.** Edman P, Begg G. A protein sequenator. *European Journal of Biochemistry*. 1967;1(1):80-91.
- 512 **9.** Prabhu Y, Kag A, Harsola S, Agrawal R, Varma M. Parabel: Partitioned label trees for extreme
513 classification with applications to dynamic search advertising. *Proceedings of the international*
514 *world wide web conference*. 2018 Apr 23-27; Lyon, France.
- 515 **10.** Zhang L, Floyd BM, Chilamari M, Mapes J, Swaminathan J, Bloom S, *et al.* Photoredox-catalyzed
516 decarboxylative C-terminal differentiation for bulk- and single- molecule proteomics. *ACS*
517 *Chemical Biology*. 2021;16(11):2595-2603.
- 518 **11.** Hinson CM, Bardo AM, Shannon CE, Rivera S, Swaminathan J, Marcotte EM, *et al.* Studies of
519 surface preparation for the fluorosequencing of peptides. *Langmuir*. 2021;37(51):14856-14865.

- 520 **12.** Messina TC, Kim H, Giurleo JT, Talaga DS. Hidden Markov Model analysis of multichromophore
521 photobleaching. *The Journal of Physical Chemistry B*. 2006;110(33):16366-16376.
- 522 **13.** Eng JK, McCormack AL, Yates JR. An approach to correlate tandem mass spectral data of
523 peptides with amino acid sequences in a protein database. *Journal of the American Society of*
524 *Mass Spectrometry*. 1994;5(11):976-989.
- 525 **14.** Keller A, Nesvizhskii AI, Kolker E, Aebersold R. Empirical statistical model to estimate the
526 accuracy of peptide identifications made by MS/MS and database search. *Analytical Chemistry*.
527 2002;74(20):5383-5392.
- 528 **15.** Käll L, Canterbury JD, Weston J, Noble WS, MacCoss MJ. Semi-supervised learning for peptide
529 identification from shotgun proteomics datasets. *Nature Methods*. 2007;4:923-925.
- 530 **16.** Dincer AB, Lu Y, Schweppe DK, Oh S, Noble WS. Reducing peptide sequence bias in quantitative
531 mass spectrometry data with machine learning. *Journal of Proteome Research*. 2022;21(7):1771-
532 1782.
- 533 **17.** Elias JE, Gygi SP. Target-decoy search strategy for mass spectrometry-based proteomics. *Methods*
534 *in Molecular Biology*. 2010;604:55-71.
- 535 **18.** Saiki RK, Gelfand DH, Stoffel S, Scharf SJ, Higuchi R, Horn GT, *et al.* Primer-directed enzymatic
536 amplification of DNA with a thermostable DNA polymerase. *Science*. 1988;239(4839):487-491.
- 537 **19.** Fedurco M, Romieu A, Williams S, Lawrence I, Turcatti G. BTA, a novel reagent for DNA
538 attachment on glass and efficient generation of solid-phase amplified DNA colonies. *Nucleic*
539 *Acids Research*. 2006;34(3):e22.
- 540 **20.** Stoler N, Nekrutenko A. Sequencing error profiles of Illumina sequencing instruments. *NAR*
541 *Genomics and Bioinformatics*. 2021;3(1):1-9.

- 542 **21.** Mitchell K, Brito JJ, Mandric I, Wu Q, Knyazev S, Chang S, *et al.* Benchmarking of computational
543 error-correction methods for next-generation sequencing data. *Genome Biology*. 2020;21(71):1-
544 13.
- 545 **22.** Cali DS, Kim JS, Ghose S, Alkan C, Mutlu O. Nanopore sequencing technology and tools for
546 genome assembly: computational analysis of the current state, bottlenecks and future
547 directions. *Briefings in Bioinformatics*. 2018;20(4):1542-1559.
- 548 **23.** Alzubaidi L, Zhang J, Humaidi AJ, Al-Dujaili A, Duan Y, Al-Shamma O, *et al.* Review of deep
549 learning: concepts, CNN architectures, challenges, applications, future directions. *Journal of Big*
550 *Data*. 2021;8(53):1-74.
- 551 **24.** Gessulat S, Schmidt T, Zolg DP, Samaras P, Schnatbaum K, Zerweck J, *et al.* Prosit: proteome-
552 wide prediction of peptide tandem mass spectra by deep learning. *Nature Methods*.
553 2019;16(6):509-518.
- 554 **25.** Käll L, Storey JD, MacCoss MJ, Noble WS. Assigning significance to peptides identified by tandem
555 mass spectrometry using decoy databases. *Journal of Proteome Research*. 2008;7(1):29-34
- 556 **26.** Keich U, Kertesz-Farkas A, Noble WS. Improved False Discovery Rate Estimation Procedure for
557 Shotgun Proteomics. *Journal of Proteome Research*. 2015;14(8):3148-61
- 558 **27.** Kong AT, Leprevost FV, Avtonomov DM, Mellacheruvu D, Nesvizhskii AI. MSFragger: ultrafast
559 and comprehensive peptide identification in mass spectrometry-based proteomics. *Nature*
560 *Methods*. 2017;14(5):513-520
- 561

562 **Appendix A – detailed**
 563 **descriptions and proofs of**
 564 **algorithms**
 565 **A0 – crib sheet for variables used**
 566 **throughout appendix A**
 567

Variable crib sheet	
$f^{(t)}$	Cumulative probabilities for each state at timestep t
$\mathbf{f}^{(t)}$	Tensor form of $f^{(t)}$
$\mathbf{O}^{(t)}$	HMM Emission matrix for timestep t
$\mathcal{O}^{(t)}$	Tensor form of $\mathbf{O}^{(t)}$
$\widehat{\mathcal{O}}^{(t)}$	A pruned approximation of $\mathcal{O}^{(t)}$
\mathbf{T}	The HMM transition matrix
\mathcal{T}	Tensor form of \mathbf{T}
$\widehat{\mathcal{T}}^{(t)}$	A pruned approximation of \mathcal{T} at timestep t
T	The number of timesteps
$Y_{1:T}$	Random variables representing the series of observations
$y_{1:T}$	The true values of $Y_{1:T}$
$X_{1:T}$	Random variables representing the state of the HMM across time
Z	Random variable representing the peptide
z	True value of Z
s	The number of states of a fluorophore
Λ	The number of fluorophores
C	The number of colors of fluorophore
Λ_c	The number of fluorophores of color c
ρ	Number of amino acids successfully removed by Edman degradation
\bar{c}	Color of N-terminal amino acid
$\lambda_{\rho,c}$	Number of amino acids which can accept a label of color c when ρ amino acids have been removed from the peptide
ϕ_c	Number of functioning fluorophores of color c remaining for the peptide
α	Number of amino acids in the peptide (before sequencing)
p_c	Fluorophore loss rate for color c
$\mathcal{B}^{(c)}$	A factored component of \mathcal{T} representing loss of fluorophores of color c
$\widehat{\mathcal{B}}^{(t,c)}$	A pruned approximation of $\mathcal{B}^{(c)}$ at timestep t
e	Edman cycle failure rate
\mathcal{E}	A factored component of \mathcal{T} representing Edman degradation success and failure
$\widehat{\mathcal{E}}^{(t)}$	A pruned approximation of \mathcal{E} at timestep t
d	Peptide detachment rate
\mathcal{D}	A factored component of \mathcal{T} representing peptide detachment
$\widehat{\mathcal{D}}^{(t)}$	A pruned approximation of \mathcal{D} at timestep t
r	Number of values of \mathcal{O} kept during pruning
r_c	Number of fluorophore counts of color c kept during pruning
$r_c^{(t)}$	r_c at time t
\bar{r}	Number of amino acid counts kept during pruning
$\bar{r}^{(t)}$	\bar{r} at time t
h	The number of peptides selected by the kNN in the hybrid model
ζ_h	The set of peptides selected by the kNN in the hybrid model

568 **A1 – HMM state space reduction**

569 **The basics**

570 For fluorophores that are distinguishable and dependent on each other, where s is the number of states
571 of one fluorophore, and Λ is the number of fluorophores, the number of states of the whole system is
572 given as in [12] by:

$$573 \quad s^\Lambda \quad (7)$$

574 In contrast, if the fluorophores are indistinguishable and independent of each other, the number of
575 states in the system is instead given by the combinatoric equation from [12]:

$$576 \quad \binom{\Lambda + s - 1}{s - 1} \quad (8)$$

577 Imaging is typically performed with high concentrations of the antioxidant Trolox [2] and for relatively
578 short time intervals (100 msec); reversibly photobleached fluorophores do not occur frequently in our
579 data and we ignore them as a first approximation. Therefore, we take s to be 2, with one state for a
580 functioning fluorophore, and another for a missing, photobleached, or chemically destroyed
581 fluorophore. This reduces (8) to $\Lambda + 1$ states given Λ indistinguishable and independent fluorophores.

582 **More colors**

583 Obviously, a red fluorophore is distinguishable from a blue one. But we would still like to benefit from
584 the indistinguishability of red fluorophores from red fluorophores, and of blue fluorophores from blue
585 fluorophores. This is modeled by first considering the states for each color of fluorophore
586 independently, and then taking the cartesian product of these state spaces. For C colors of fluorophore,
587 where Λ_c is the number of fluorophores of color c , this results in the number of states given by:

$$588 \quad \prod_{c=1}^C (\Lambda_c + 1) \quad (9)$$

589 Each state then represents the number of remaining active fluorophores for each of our C colors of
590 fluorophore.

591 **Edman degradation**

592 Our second challenge is the inclusion of Edman degradation in the HMM. Sequential removal of the N-
593 terminal amino acid from each peptide breaks the assumption of indistinguishable fluorophores, which
594 is the basis for the state reduction performed in [12]. However, through inductive reasoning we show
595 that our model meets a weaker criterion, which can be used to merge states together as desired:

*Any two states with the same numbers of fluorophores of each color and the same
number of amino acids are equally likely.* (10)

596 If we ignore Edman degradation, this follows directly from the assumed indistinguishability property of
597 fluorophores of the same color; if two fluorophores behave identically, they are equally likely to be
598 missing, photobleached, or chemically destroyed, thus it follows by symmetry that any two states with
599 the same numbers of indistinguishable fluorophores of each color are equally likely. If we consider
600 Edman degradation, then (10) is true for all states where no amino acids have yet been successfully
601 removed. Let ρ indicate the number of amino acids removed from the original peptide. We have then
602 shown that (10) is true when $\rho = 0$.

603 If states with identical fluorophore counts are equally probable for all states with ρ amino acids
604 removed, it can be shown that all states with equal fluorophore counts are equally probable for all
605 states with $\rho + 1$ amino acids removed. For removal of an amino acid that can't accept fluorophores
606 under the experimental setup this is trivial, so consider a peptide with ρ removed amino acids, an N-
607 terminal amino acid which accepts fluorophores of color \bar{c} , and $\lambda_{\rho, \bar{c}}$ amino acids total which can accept a

608 label of color \bar{c} . Then let $\phi_{\bar{c}}$ represent the number of remaining functional fluorophores for the peptide,
609 satisfying $0 \leq \phi_{\bar{c}} \leq \lambda_{\rho, \bar{c}}$.

610 There are several conditions of the peptide with $\phi_{\bar{c}}$ functioning fluorophores scattered among the $\lambda_{\rho, \bar{c}}$
611 amino acids that can accept a label. When we remove the N-terminal amino acid, we may or may not
612 remove with it a functioning fluorophore. The states which do have a functioning fluorophore in the N-
613 terminal position (only possible when $\phi_{\bar{c}} > 0$) will have their other $\phi_{\bar{c}} - 1$ fluorophores distributed
614 between the $\lambda_{\rho, \bar{c}} - 1$ remaining amino acids which can be labeled. Furthermore, these states are
615 equally likely, as they are a subset of the equally likely states with $\phi_{\bar{c}}$ fluorophores. Since trivially $\lambda_{\rho, \bar{c}} -$
616 $1 = \lambda_{\rho-1, \bar{c}}$, these states map one-to-one with the states for the peptide with one less amino acid
617 remaining, when it has $\phi_{\bar{c}} - 1$ dyes.

618 When $\phi_{\bar{c}} < \lambda_{\rho, \bar{c}}$, there are states with no fluorophore in the N-terminal position, even though the N-
619 terminal amino acid can accept one. Then the $\phi_{\bar{c}}$ fluorophores will be distributed with equal
620 probabilities among the $\lambda_{\rho, \bar{c}} - 1 = \lambda_{\rho-1, \bar{c}}$ remaining amino acids which can be labeled. Similarly to the
621 other case, these states map one-to-one with the states for the peptide less one amino acid when it has
622 $\phi_{\bar{c}}$ dyes.

623 The equally distributed probabilities and one-to-one correspondence between states across this amino
624 acid removal ensures that these transformations do not break our guarantees of equal probabilities for
625 $\rho + 1$ amino acids removed. Iteratively applying this reasoning, starting with $\rho = 0$, until we prove that
626 states where $\rho = \alpha$ are equally likely if they have the same fluorophore counts, demonstrates that (5) is
627 true under the assumptions we have taken.

628 This proves (10), which allowed us to safely merge states that share both the same fluorophore counts
629 by color and the same numbers of amino acids.

630 **Transition probabilities**

631 We also need to know the transition probabilities for our new reduced state space. To deal with peptide
632 detachment is trivial. Dye-loss, either for dyes missing before sequencing begins, or from chemical
633 destruction during sequencing, can be modeled with a binomial distribution. This follows from the
634 assumption that the fluorophores behave independently of each other.

635 For Edman degradation, there is of course a probability of success or failure of the degradation step,
636 which we model as a Bernoulli random variable. In the case of success, we employ an additional
637 Bernoulli random variable to model the probability of losing or not losing a functioning fluorophore.
638 Because the true states within a merged state are equally likely, we can use combinatorics to count the
639 number of states which will lose a dye, and the number that won't. Together these values can be used
640 to find the probability of losing a fluorophore given a successful Edman degradation, as shown in the
641 following formula, which conveniently reduces to a simple fraction:

$$642 \quad \frac{\binom{\lambda_{\rho, \bar{c}-1}}{\phi_{\bar{c}-1}}}{\binom{\lambda_{\rho, \bar{c}}}{\phi_{\bar{c}}}} = \frac{\phi_{\bar{c}}}{\lambda_{\rho, \bar{c}}} \quad (11)$$

643 **State reduction conclusions**

644 This state reduction provides a considerable algorithmic complexity improvement to the HMM forward
645 algorithm. The complexity of the forward algorithm is $O(S^2T)$, where S is the number of states, and T is
646 the number of timesteps. Then, if implemented with the true state space of a labeled peptide, the
647 number of states S is $O(\alpha^{2\Lambda})$, and we get a complexity of $O(\alpha^2 4^\Lambda T)$ for the HMM forward algorithm,
648 where α is the number of amino acids and Λ is the total number of fluorophores (of any color).

649 However, if we use the reduced state space, then S is $O(\alpha \prod_{c=1}^C \Lambda_c)$, giving an algorithmic complexity of
650 $O(\alpha^2 (\prod_{c=1}^C \Lambda_c^2) T)$ for the forward algorithm, where C is the number of fluorophore colors being used
651 and Λ_c is the number of fluorophores of color c . The scaling in either case is dominated by values of Λ or

652 Λ_c , which ranges from 1 to about 25 for human tryptic peptides, though in rare cases Λ_c can exceed
653 100.

654 **A2 – Transition matrix factoring**

655 **The concept**

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

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

670 However, we did better by factoring this matrix (**Figure 4**). We used the independence of our different
671 forms of error, with one matrix in the factored product for each type of error. To demonstrate this
672 factorization, we reformulated our problem in tensor notation. The vector for the state space of a
673 peptide with C colors not undergoing Edman degradation or peptide detachment can be viewed as a
674 tensor of order C . Each index of the tensor maps to the fluorophore counts of a different color, and the

675 value of an index i_c indicates the number of functioning fluorophores of color c , and satisfies $0 \leq i_c \leq$
 676 Λ_c . We also have indices j_c which are similarly defined. Since the transition matrix is a linear mapping
 677 from and to this tensor of order C , it is necessarily of order $2C$. We use the Einstein summation
 678 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.
 679 the matrix vector multiplication operation for one step of the HMM forward algorithm is then given by:

$$680 \quad \mathbf{f}_k^{(t+1)} = \mathbf{O}_{kj}^{(t+1)} \mathcal{T}_{ji} \mathbf{f}_i^{(t)} \quad (12)$$

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

684 **Considering fluorophore loss only**

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

$$687 \quad \mathcal{T}_{ji} = \begin{cases} \prod_{c=1}^C \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c}, & \text{if } \mathbf{j} \leq \mathbf{i} \\ 0, & \text{otherwise} \end{cases} \quad (13)$$

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

$$691 \quad \mathbf{B}_{ji}^{(c)} = \begin{cases} \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c}, & \text{if } j \leq i \\ 0, & \text{otherwise} \end{cases} \quad (14)$$

692 This produces a factorization of \mathcal{T} :

$$693 \quad \mathcal{T}_{ji} = \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_C i_C}^{(C)} \quad (15)$$

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

$$695 \quad \mathbf{f}_j^{(t+1)} = \mathbf{B}_{j_1 i_1}^{(1)} \mathbf{B}_{j_2 i_2}^{(2)} \dots \mathbf{B}_{j_c i_c}^{(c)} \mathbf{f}_i^{(t)} \quad (16)$$

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

$$697 \quad O\left(\left(\prod_{c=1}^C \Lambda_c\right)\left(\sum_{c=1}^C \Lambda_c\right)\right).$$

698 **Fluorophore loss and Edman degradation**

699 We can expand on this to consider the Edman degradation: In that case we need more indices for the
 700 number of remaining amino acids. We modify (12) with additional indices u and v which satisfy $0 \leq$
 701 $u \leq \alpha$ and $0 \leq v \leq \alpha$, indicating the number of successful amino acid removals, or alternatively the
 702 position of an amino acid in the peptide (i. e., the amino acid at the N-terminus of the peptide when u
 703 amino acids have been removed). This gives:

$$704 \quad \mathbf{f}_{vk}^{(t+1)} = \mathbf{O}_{kj}^{(t+1)} \mathcal{T}_{vjui} \mathbf{f}_{ui}^{(t)} \quad (17)$$

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

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

$$712 \quad \mathcal{T}_{vjui} = \begin{cases} e\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u \\ (1 - e)\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u = 0 \\ (1 - e) \left(\left(1 - \frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right) \beta(\mathbf{i}, \mathbf{j}) + \left(\frac{i_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}}\right) \bar{\beta}(\mathbf{i}, \mathbf{j}, u) \right), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u > 0 \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

713 Where:

$$714 \quad \beta(\mathbf{i}, \mathbf{j}) = \prod_{c=1}^C \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c} \quad (19)$$

715 And:

$$716 \quad \bar{\beta}(\mathbf{i}, \mathbf{j}, u) = \binom{i_{\bar{c}_u} - 1}{j_{\bar{c}_u}} p_{\bar{c}_u}^{i_{\bar{c}_u} - 1 - j_{\bar{c}_u}} (1 - p_{\bar{c}_u})^{j_{\bar{c}_u}} \prod_{\substack{1 \leq c \leq C \\ c \neq \bar{c}_u}} \binom{i_c}{j_c} p_c^{i_c - j_c} (1 - p_c)^{j_c} \quad (20)$$

717 The probability of an Edman degradation failure is essentially the same as in (13), but multiplied by e to
 718 account for the probability of failure. The probability for a transition involving a successful Edman
 719 degradation event which removes an unlabelable amino acid is similarly just like in (13) but multiplied
 720 by $(1 - e)$, the probability of success. If the amino acid in question is labelable by a color \bar{c}_u , then we
 721 may or may not remove a fluorophore of that color in the transition, so we need to take the sum of both
 722 possibilities. β in (19) gives the standard product of binomials formula from (13), but needs to be
 723 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 $\bar{\beta}$
 724 from (20) which gives the product of binomial probabilities starting with one less fluorophore of the
 725 color \bar{c}_u , which in (18) is multiplied with the probability of losing a fluorophore with the Edman
 726 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
 727 degradation success, given by $(1 - e)$.

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

$$730 \quad \mathcal{E}_{vkuj} = \begin{cases} e, & \text{if } v = u \text{ and } \mathbf{k} = \mathbf{j} \\ 1 - e, & \text{if } v = u + 1 \text{ and } \mathbf{k} = \mathbf{j} \text{ and } \bar{c}_u = 0 \\ (1 - e) \left(1 - \frac{j_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}} \right), & \text{if } v = u + 1 \text{ and } \mathbf{k} = \mathbf{j} \text{ and } \bar{c}_u > 0 \\ (1 - e) \left(\frac{j_{\bar{c}_u}}{\lambda_{u, \bar{c}_u}} \right), & \text{if } v = u + 1 \text{ and } k_{\bar{c}_u} = j_{\bar{c}_u} - 1 \text{ and } k_c = j_c \forall c \neq \bar{c}_u \text{ and } \bar{c}_u > 0 \\ 0, & \text{otherwise} \end{cases} \quad (21)$$

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

$$732 \quad \mathcal{T}_{vkui} = \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_c i_c}^{(c)} \quad (22)$$

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

$$734 \quad \mathcal{f}_{vl}^{(t+1)} = \mathcal{O}_{lk}^{(t+1)} \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_c i_c}^{(c)} \mathcal{f}_{ui}^{(t)} \quad (23)$$

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

742 **Everything all together**

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

$$744 \quad \mathcal{T}_{vjui} = \begin{cases} (1 - d)e\beta(\mathbf{i}, \mathbf{j}, p), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u \\ (1 - d)(1 - e)\beta(\mathbf{i}, \mathbf{j}), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u = 0 \\ (1 - d)(1 - e) \left(\left(1 - \frac{i_{\bar{c}_u}}{\lambda_u} \right) \beta(\mathbf{i}, \mathbf{j}) + \left(\frac{i_{\bar{c}_u}}{\lambda_u} \right) \bar{\beta}(\mathbf{i}, \mathbf{j}, u) \right), & \text{if } \mathbf{j} \leq \mathbf{i} \text{ and } v = u + 1 \text{ and } \bar{c}_u > 0 \\ d, & \text{if } j_c = 0 \forall c \text{ and } v = \alpha \\ 0, & \text{otherwise} \end{cases} \quad (24)$$

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

$$749 \quad \mathcal{D}_{whvk} = \begin{cases} 1 - d, & \text{if } \mathbf{h} = \mathbf{k} \text{ and } w = v \leq \alpha \\ d & \text{if } h_c = 0 \forall c \text{ and } w = \alpha + 1 \end{cases} \quad (25)$$

750 Then we find that:

$$751 \quad \mathcal{J}_{wlui} = \mathcal{D}_{wlvk} \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_C i_C}^{(C)} \quad (26)$$

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

$$753 \quad \mathcal{f}_{wm}^{(t+1)} = \mathcal{O}_{ml}^{(t+1)} \mathcal{D}_{wlvk} \mathcal{E}_{vkuj} \mathcal{B}_{j_1 i_1}^{(1)} \mathcal{B}_{j_2 i_2}^{(2)} \dots \mathcal{B}_{j_C i_C}^{(C)} \mathcal{f}_{ui}^{(t)} \quad (27)$$

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

755 Thus, \mathcal{D} has no impact on the algorithmic complexity of this operation. Although \mathcal{D} and \mathcal{E} could be

756 combined to achieve this same algorithmic improvement, we found that this separation made our

757 model easier to reason about and work with.

758 **Transition matrix factoring conclusions**

759 One of the benefits of this approach to algorithmic complexity reduction is that this factorization

760 provides no loss to the theoretical accuracy of the forward algorithm. No theoretical approximations

761 were necessary, aside from the unavoidable differences in floating-point round-off errors. This allows

762 for highly accurate results with much more efficient runtime characteristics than a naïve

763 implementation.

764 **A3 – HMM pruning**

765 Because the emission matrix is diagonal, it is equivalent to the diagonal part of its Singular Value
766 Decomposition (SVD), but with a reordering of its indices. This makes sparsification of this matrix
767 equivalent to the Eckart-Young-Mirsky theorem; we can keep the largest r values for some chosen value
768 of r , and replace the rest of the matrix entries with zeros, having the minimum possible impact on the
769 spectral and Frobenius norms for the chosen value of r .

770 Furthermore, we can propagate this sparsification to the transition matrix. Consider the forward
771 algorithm, with \mathbf{T} representing the transition matrix, and $\mathbf{O}^{(t)}$ representing the diagonal emission matrix
772 for time t . Then if $\mathbf{f}^{(t)}$ represents the vector of intermediate probabilities at time t , we have:

$$773 \quad \mathbf{f}^{(t+1)} = \mathbf{O}^{(t+1)}\mathbf{T}\mathbf{f}^{(t)} \quad (28)$$

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

$$775 \quad \mathbf{f}^{(t+1)} = \widehat{\mathbf{O}}^{(t+1)}\mathbf{T}\mathbf{f}^{(t)} \quad (29)$$

776 Note that we have many copies of \mathbf{T} , which are equal. For our next improvements we need these to be
777 different for each timestep, so we can rewrite (29) with $\mathbf{T}^{(t)}$ for each timestep t , giving

$$778 \quad \mathbf{f}^{(t+1)} = \widehat{\mathbf{O}}^{(t+1)}\mathbf{T}^{(t)}\mathbf{f}^{(t)} \quad (30)$$

779 Here the values of many rows and columns of each $\mathbf{T}^{(t)}$ have been made unnecessary by the
780 sparsification of its neighboring $\widehat{\mathbf{O}}^{(t+1)}$ and $\widehat{\mathbf{O}}^{(t)}$, as any vector product with $\widehat{\mathbf{O}}^{(t)}$ will necessarily have
781 zeros except for the r entries retained, such that we need only keep the corresponding r columns of
782 $\mathbf{T}^{(t)}$. Similarly, any entry in the vector product with $\mathbf{T}^{(t)}$ which is not multiplied by one of the r entries
783 retained in $\widehat{\mathbf{O}}^{(t+1)}$ is multiplied by zeros, and is thus unnecessary, so we need only keep the
784 corresponding r rows of $\mathbf{T}^{(t)}$. Calling these approximations $\widehat{\mathbf{T}}^{(t)}$, we get

785
$$\mathbf{f}^{(t+1)} = \widehat{\mathbf{O}}^{(t+1)} \widehat{\mathbf{T}}^{(t)} \mathbf{f}^{(t)} \quad (31)$$

786 This allows significant sparsity to be used (**Figure 5**). Previously this formula would have been
 787 $O(\alpha^2 T \prod_{c=1}^C \Lambda_c^2)$ to compute, while this reduces the algorithmic complexity to $O(r^2 T)$. This
 788 improvement is beyond what is possible in a more traditional usage of sparse matrix multiplication. For
 789 sparse matrix multiplication, we would need to first multiply $\widehat{\mathbf{T}}^{(t)}$ by $\widehat{\mathbf{O}}^{(t)}$ or multiply $\widehat{\mathbf{O}}^{(t+1)}$ by $\widehat{\mathbf{T}}^{(t)}$. This
 790 will only permit you to sparsify your operations on the rows or the columns of \mathbf{T} but not both, giving a
 791 complexity of $O(r \alpha T \prod_{c=1}^C \Lambda_c)$. While this is better than not using this inherent sparsity at all,
 792 preprocessing the transition matrix in consideration of the emission matrices on either side gives better
 793 results in algorithmic complexity.

794 In practice, we use a more complicated pruning scheme, as detailed next.

795 **A4 – Combining transition matrix factoring with HMM pruning**

796 By making r suitably small, HMM pruning can exhibit better algorithmic complexity than if we factor the
 797 transition matrix. However, we believe it is much better to combine these algorithmic enhancements
 798 (**Figure 6**). To do this, we need to switch into tensor notation, replacing our matrices and vectors with
 799 the tensor equivalents we constructed previously. This yields:

800
$$\mathbf{f}_{vk}^{(t+1)} = \widehat{\mathbf{O}}_{kj}^{(t+1)} \widehat{\mathbf{T}}_{vju}^{(t)} \mathbf{f}_{ui}^{(t)} \quad (32)$$

801 We also want to use the factorization from (26), using timestamp specific sub-tensors of each of the
 802 factored pieces. The factorization of (26) becomes:

803
$$\widehat{\mathbf{T}}_{wlu}^{(t)} = \widehat{\mathbf{D}}_{wlvk}^{(t)} \widehat{\mathbf{E}}_{vku}^{(t)} \widehat{\mathbf{B}}_{j_1 i_1}^{(t,1)} \widehat{\mathbf{B}}_{j_2 i_2}^{(t,2)} \dots \widehat{\mathbf{B}}_{j_C i_C}^{(t,C)} \quad (33)$$

804 Substituting into (32) gives:

805
$$\mathbf{f}_{wm}^{(t+1)} = \widehat{\mathbf{O}}_{ml}^{(t+1)} \widehat{\mathbf{D}}_{wlvk}^{(t)} \widehat{\mathbf{E}}_{vku}^{(t)} \widehat{\mathbf{B}}_{j_1 i_1}^{(t,1)} \widehat{\mathbf{B}}_{j_2 i_2}^{(t,2)} \dots \widehat{\mathbf{B}}_{j_C i_C}^{(t,C)} \mathbf{f}_{ui}^{(t)} \quad (34)$$

806 Suppose we were to use standard sparse tensor multiplication techniques and carry this operation out
807 from right to left. Each tensor $\widehat{\mathcal{B}}_{j_c i_c}^{(t,c)}$ can introduce any entry of input (index i_c) into as many as Λ_c
808 indices of output (index j_c). The resulting computational complexity of the forward algorithm, even with
809 the given sparsity, is then $O(r\alpha T \prod_{c=1}^C \Lambda_c)$.

810 If we preprocess the computation, pruning each operation now from both directions, the algorithmic
811 complexity does not improve the way it does in the matrix case, although likely this would behave faster
812 in practice. The problem is that the pruning operation itself needs to determine which rows to
813 propagate forwards, which requires accessing every non-zero entry reachable in the forward direction.
814 Many of these values are later pruned in the backwards direction, so the computation itself has much
815 better sparsity, but the time to prune then dominates the algorithmic complexity result.

816 To improve this further, we add structure to the pruning of $\widehat{\mathcal{O}}^{(t)}$. Instead of keeping the r largest values
817 in $\widehat{\mathcal{O}}^{(t)}$, we prune each index of $\widehat{\mathcal{O}}^{(t)}$ independently. For additional convenience, we limit each index to a
818 contiguous range of values. Then we let each index for any fluorophore color c have r_c values and allow
819 \bar{r} values to index the number of amino acids. These simplifications may cause the pruning to be non-
820 optimal, but we accept this trade-off.

821 We can then prune our tensors using their known structures (for example, the tensors $\widehat{\mathcal{B}}^{(t,c)}$ correspond
822 to an upper triangular matrix). This time when we propagate the pruning results in both directions, the
823 time required is only $O(C^2)$ (the number of minimum and maximum indices to be propagated through
824 each tensor scales with C , as does the number of tensors to be pruned). For the runtime of the tensor
825 operations, consider each tensor individually. $\widehat{\mathcal{D}}^{(t)}$ and $\widehat{\mathcal{E}}^{(t)}$ are both highly sparse, so they contribute a
826 constant modification to the number of rows or columns when propagating in either direction. $\widehat{\mathcal{D}}^{(t)}$
827 requires special handling. We track the detached state separately from the ordinary range, to avoid
828 unnecessarily including a large range of states which don't need to be.

829 Then each $\widehat{\mathcal{B}}^{(t,c)}$ operates on an independent index, and therefore can be considered on its own. This
830 tensor after pruning will have dimensions that are $O(r_c^2)$, and should have a constant effect on the
831 number of elements input vs output. Therefore, each of these tensors will require an algorithmic
832 complexity of $O(r_c \bar{r} \prod_{\bar{c}=1}^c r_{\bar{c}})$. Bringing this all together we get an algorithmic complexity of
833 $O\left(C^2 + \bar{r}(\sum_{c=1}^C r_c)(\prod_{c=1}^C r_c)\right)$ for processing one timestep. The full forward algorithm then has a
834 complexity of $O\left(T\left(C^2 + \bar{r}(\sum_{c=1}^C r_c)(\prod_{c=1}^C r_c)\right)\right)$.

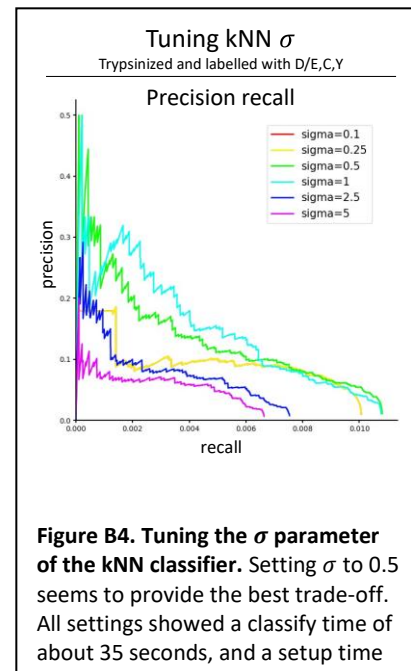
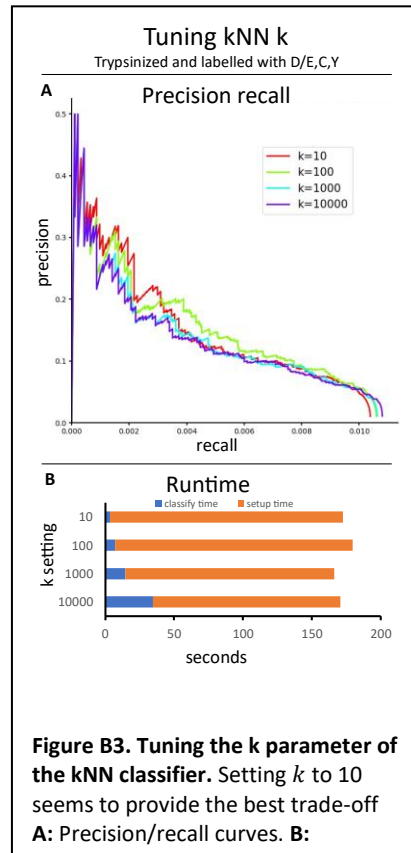
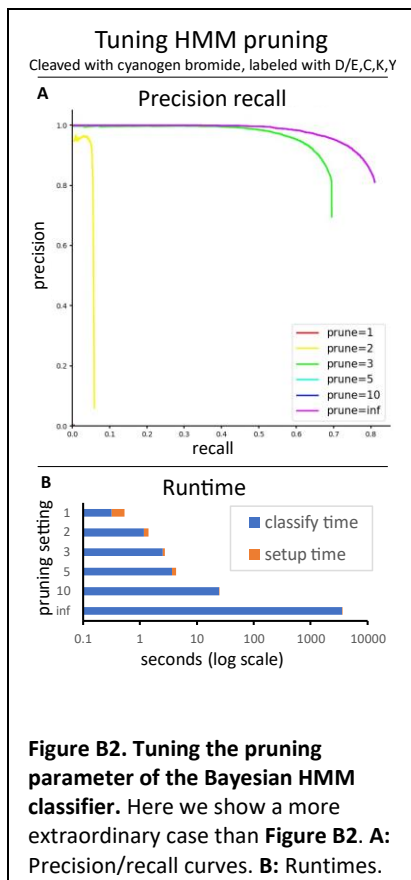
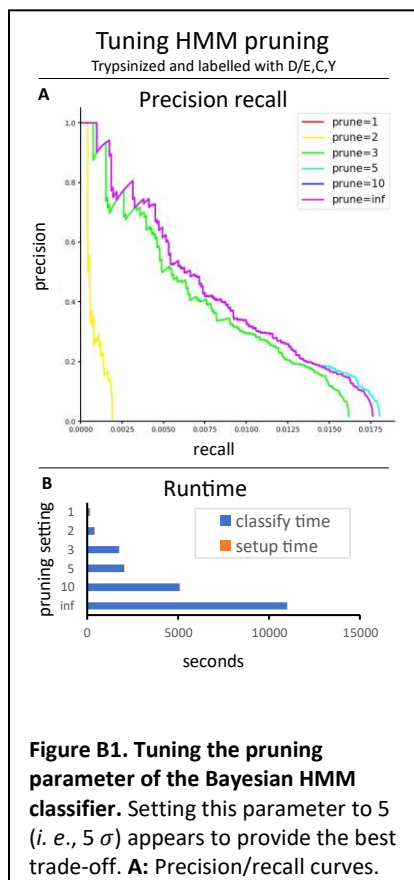
835 One remaining clarification is the manner of choosing r_c and \bar{r} . In fact, these values should not be kept
836 constant; let us refer to the values for time t as $r_c^{(t)}$ and $\bar{r}^{(t)}$. $\bar{r}^{(0)} = 1$ and $\bar{r}^{(t+1)} = \bar{r}^{(t)} + 1$, due to the
837 possibility of amino acid removal. These on average are proportional to α . To get $r_c^{(t)}$, we keep all index
838 values where a fluorophore count of that value has the observed fluorescence intensity for color c at
839 time t within a specified confidence interval – perhaps within 3σ of the mean, where σ is the standard
840 deviation of the distribution. These will necessarily be contiguous. The number of indices kept is then
841 $r_c^{(t)}$.

842 The standard deviation of a normal distribution scales with the square root of the intensity, and the
843 number of possible index values is limited by the total possible number of fluorophores of color c . It
844 follows that any removal of index values proportional to the standard deviation will satisfy $r_c^{(t)} < \gamma\sqrt{\Lambda_c}$
845 for some constant γ dependent on the cutoff. Then the algorithmic complexity is given by
846 $O\left(T\left(C^2 + \alpha(\sum_{c=1}^C \sqrt{\Lambda_c})(\prod_{c=1}^C \sqrt{\Lambda_c})\right)\right)$.

847 We chose a specific pruning cut-off by sweeping this parameter and balancing the experimental runtime
848 effects and the precision-recall curves which result from simulated data.

849

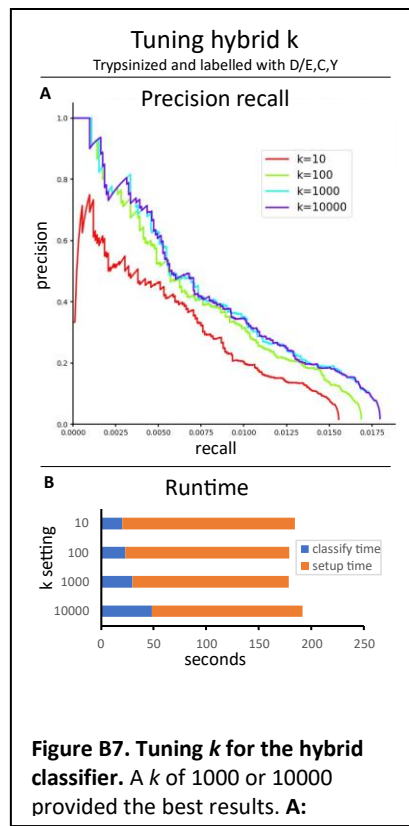
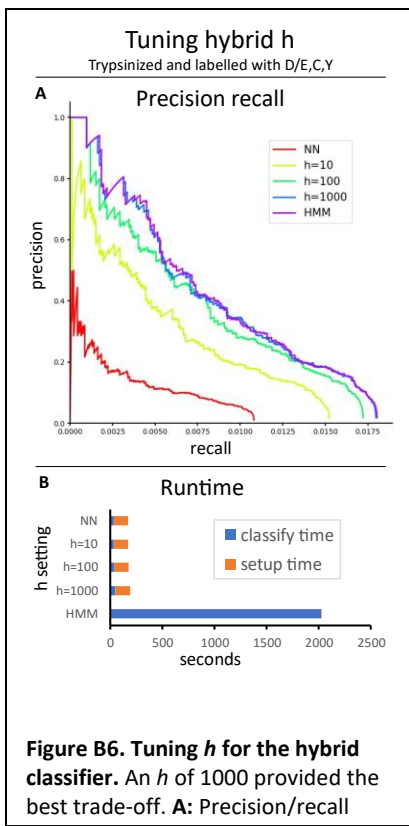
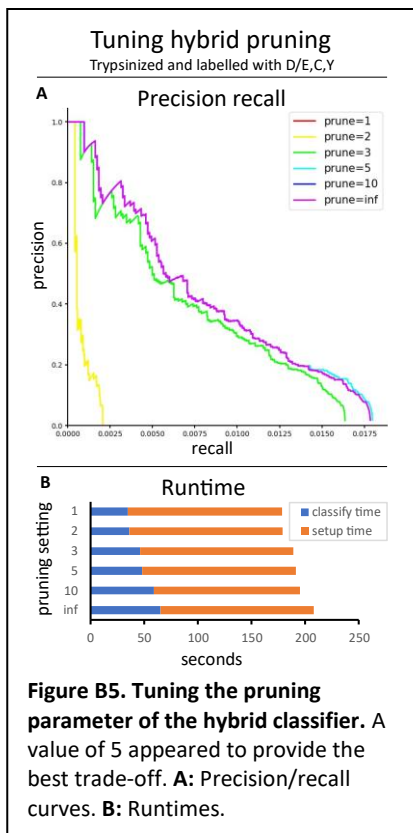
850 **Appendix B – Supporting figures**



851

852

853



854

