

A APPENDIX

Algorithm 1 realSEUDO Algorithm

```

1: Initialize:  $\mathbf{X}_{temp} \leftarrow []$ ;  $\mathbf{X}_{stab} \leftarrow []$ 
2: for each frame  $\mathbf{y}_t$  do
3:   Denoise  $\mathbf{y}_t$ 
4:   Identify the profiles in the current frame
5:   for all new profiles do
6:     if current profile overlaps any of  $\mathbf{X}_{temp}$  then
7:       Merge new  $\mathbf{X}_{temp}$  profiles into the current  $\mathbf{X}_{temp}$ 
8:     else
9:       Add the current profile to  $\mathbf{X}_{temp}$ 
10:    end if
11:  end for
12:  for all profiles in  $\mathbf{X}_{temp}$  that have not been updated in the last few frames do
13:    Move them from  $\mathbf{X}_{temp}$  to  $\mathbf{X}_{stab}$ 
14:  Merge the moved profiles with existing  $\mathbf{X}_{stab}$  profiles
15:  end for
16:   $\phi_t, \mathbf{r}_t \leftarrow \text{SEUDO}(\mathbf{y}_t, \mathbf{X}_{stab})$ 
17:  Report  $\phi_t$  as detections
18:   $\phi'_t \leftarrow \text{SEUDO}(\mathbf{r}_t, \mathbf{X}_{temp})$ 
19:  for each profile  $k$  in  $\mathbf{X}_{temp}$  do
20:    if the  $\phi'_{kt} > \gamma$  or this profile was previously active then
21:      Report this activation as early detection
22:    end if
23:  end for
24: end for

```

A.1 APPLICATION OF MODIFIED FISTA TO NEURAL NETWORK OPTIMIZATION

We further tested the modified FISTA momentum descent algorithm to problems outside of neuroscience—the training of neural networks—to evaluate the scope of applicability of our improvements. We used the problem of recognition of handwritten digits on a data set from AT&T Research available at https://hastie.su.domains/StatLearnSparsity_files/DATA/zipcode.html, reduced to 8x8 pixels, with a training set of 7291 images. The neural network (NN) model used the Leaky ReLU activation, with layer sizes 64, 64, 32, 10.

The common approach to training NNs uses stochastic gradient descent, including stochastic momentum methods. Thus to compare to a non-stochastic momentum method we established a non-stochastic baseline.

The dynamic estimation of the Lipschitz constant L from TFOCS cannot be applied to the neural network optimization because the optimization cost is highly non-linear. The multi-dimensional estimation of L is also computed only for the specific SEUDO function. The common practice is to use a fixed descent rate, which serves as an analog of $\frac{1}{2L}$. Estimating the highest descent rate is still not a fully solved problem. Algorithms for dynamic evaluation of the descent rate do exist (e.g., the ADAM algorithm in Kingma & Ba, 2015), however they rely on a constant to be picked for a particular problem.

The advantage of stochastic methods is that they can use a higher descent rate without diverging, as seen by observing the dependency of logarithm of mean square error from the number of training passes for various descent rates (Fig 6). In these results we ensure that we start from the same fixed randomized initial state since different initial states can produce wildly different results.

We observe that for the same descent rate (0.05 per pass), the stochastic and non-stochastic methods produce very similar error values, however the graph for the stochastic method is more smooth. The roughness of the graph represents the small divergences that manage to converge again over time, and shows that the descent rate is close to the maximum. However the stochastic method can accommodate a 100 times higher descent rate without diverging, and even a 1000 higher descent

Algorithm 2 Modified FISTA algorithm

```

1: Initialize  $t = 1$ ;  $x[] = (\text{initial values})$ ;  $\text{diff}[] = [0]$ ;  $\text{gradient\_last}[] = [0]$ 
2: for  $\text{step} = 1$  to  $\text{maxstep}$  do
3:    $t_{\text{next}} = \frac{1 + \sqrt{1 + t^2 * 4}}{2}$ 
4:    $\eta = \frac{t-1}{t_{\text{next}}}$ 
5:   if  $\text{step} \neq 1$  then
6:      $t = t_{\text{next}}$ 
7:   end if
8:    $x[] = x[] + \eta * \text{diff}[]$ 
9:   for each  $i$  in dimensions of  $x$  do
10:    if  $x[i] < 0$  then
11:       $x[i] = 0$ ;  $\text{diff}[i] = 0$ ;
12:    end if
13:  end for
14:   $\text{gradient}[] = \text{compute\_gradient\_f}(x[])$ 
15:   $x[] = x[] - \frac{\text{gradient}[]}{L}$ 
16:  for each  $i$  in dimensions of  $x$  do
17:    if  $x[i] < 0$  then
18:       $x[i] = 0$ ;  $\text{diff}[i] = 0$ ;  $\text{gradient}[i] = 0$ 
19:    else
20:      if  $\text{gradient}[i] * \text{gradient\_last}[i] < 0$  then
21:         $\text{diff}[i] = 0$ 
22:      else
23:         $\text{diff}[i] = \text{diff}[i] - \frac{\text{gradient}[i]}{L}$ 
24:      end if
25:    end if
26:  end for
27:   $\text{gradient\_last}[] = \text{gradient}[]$ 
28: end for

```

Method	error	log(error)
stochastic	0.0956	-2.3476
baseline non-stochastic	0.0952	-2.3515
FISTA	0.1662	-1.7946
momentum+stop on gradient sign change	0.0699	-2.6607
momentum+stop on gradient sign change + $\eta = 1$	0.0701	-2.6578
momentum+stop on gradient sign change + $\eta = 1$ at 4x rate	0.0751	-2.5889
auto-adjusted rate for momentum+stop on gradient sign change	0.0631	-2.7630
auto-adjusted rate for momentum+stop on gradient sign change + $\eta = 1$	0.0654	-2.7272

Table 1: Performance of algorithms on handwritten dataset

rate becomes rough but still converges. The non-stochastic method is able to make the passes faster, because it performs the same accumulation of partial gradients, but saves the overhead of updating the weights after each training case (or batch). However even adjusted for time, the stochastic method performs faster. It is possible to compute the non-stochastic gradient in parallel by multiple threads but we have not implemented this. We used this example of the non-stochastic descent as a baseline for the FISTA-based momentum methods.

The summary of training errors in the momentum methods can be found in Figure 7A, and the mean square errors after 10,000 training passes are listed in the Table 1.

The unmodified FISTA algorithm with $\lambda = 0$ performed on this task out of its domain worse than the non-momentum baseline. Adding the momentum stop in the dimensions with gradient sign change produced a substantial improvement over the baseline. Fixing the parameter $\eta = 1$ produced a close result to not fixing η , but with a less rough curve. We tested if the smoothness indicated that setting

$\eta = 1$ could accommodate a substantial increase in descent rate by re-running the algorithm at a 4x rate (0.2 instead of 0.05), and while this run did not diverge, we did observe a higher error rate.

Finally, we attempted to devise an algorithm that acts similar to the TFOCS dynamic evaluation of L but using the ratios of mean square values of gradient dimensions that change or not change sign as an indication of roughness. The rapid growth of gradient dimensions after sign change is seen as a beginning of a divergence, that causes the reduction of descent rate. This algorithm allowed training at a substantially higher rate in the first few thousands of passes but then flattened out. The automatically determined rate is close to the empirically found 0.05, and is higher in the initial passes where it reaches higher values, but then drops to the lower values (Fig. 7B). It is possible that the chosen criteria were not aggressive enough, and can be improved.

While even with the momentum descent the stochastic methods specialized for NN training can still achieve faster speeds (Fig. 7C-D), we have demonstrated that our more general optimization still represent a major improvement over both simple gradient descent and plain FISTA in a different domain.

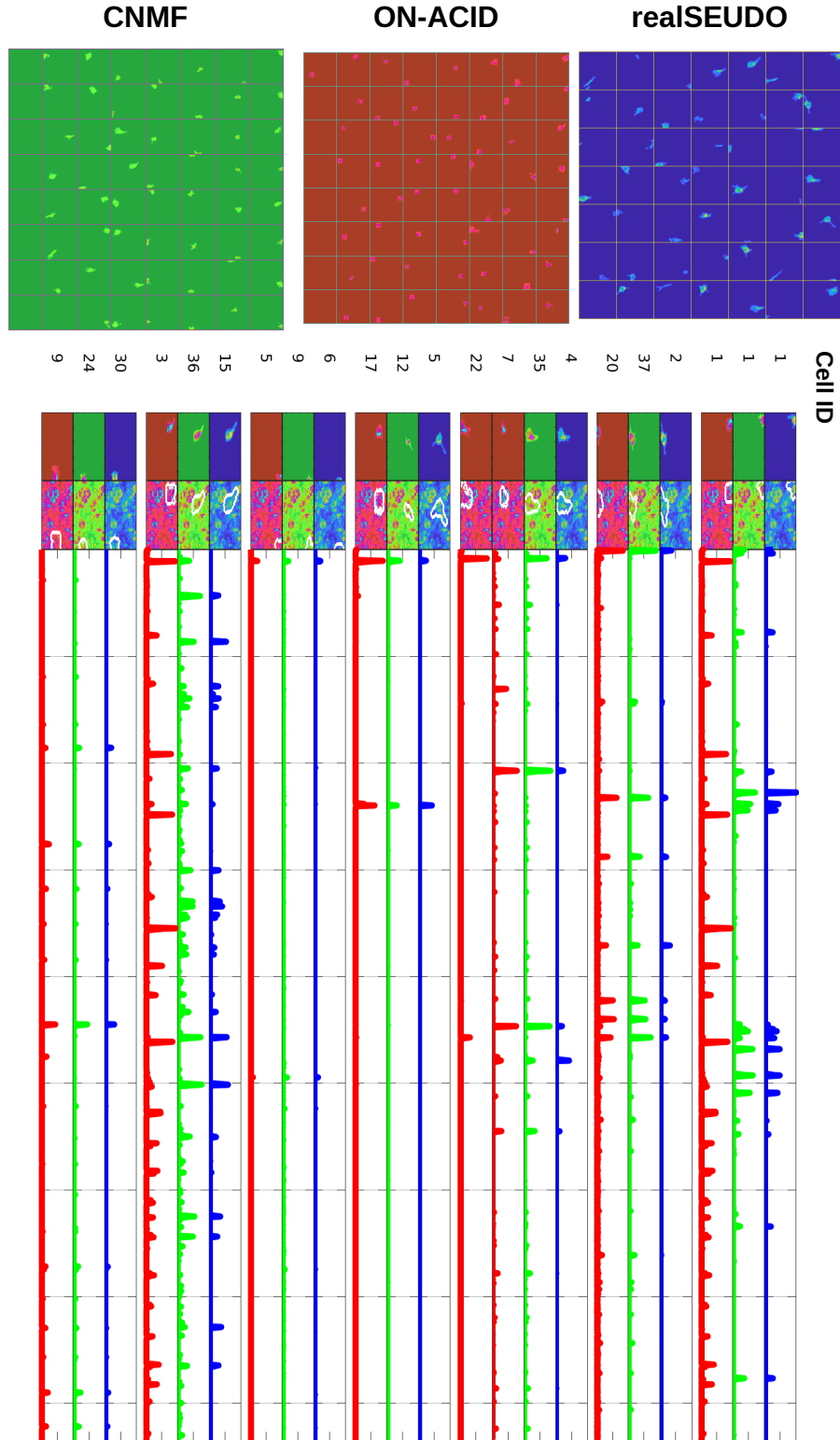


Figure 4: The whole set of cells detected by realSEUDO (blue), OnACID (red), and CNMF (green) in one movie, with selected time traces of matching cells.

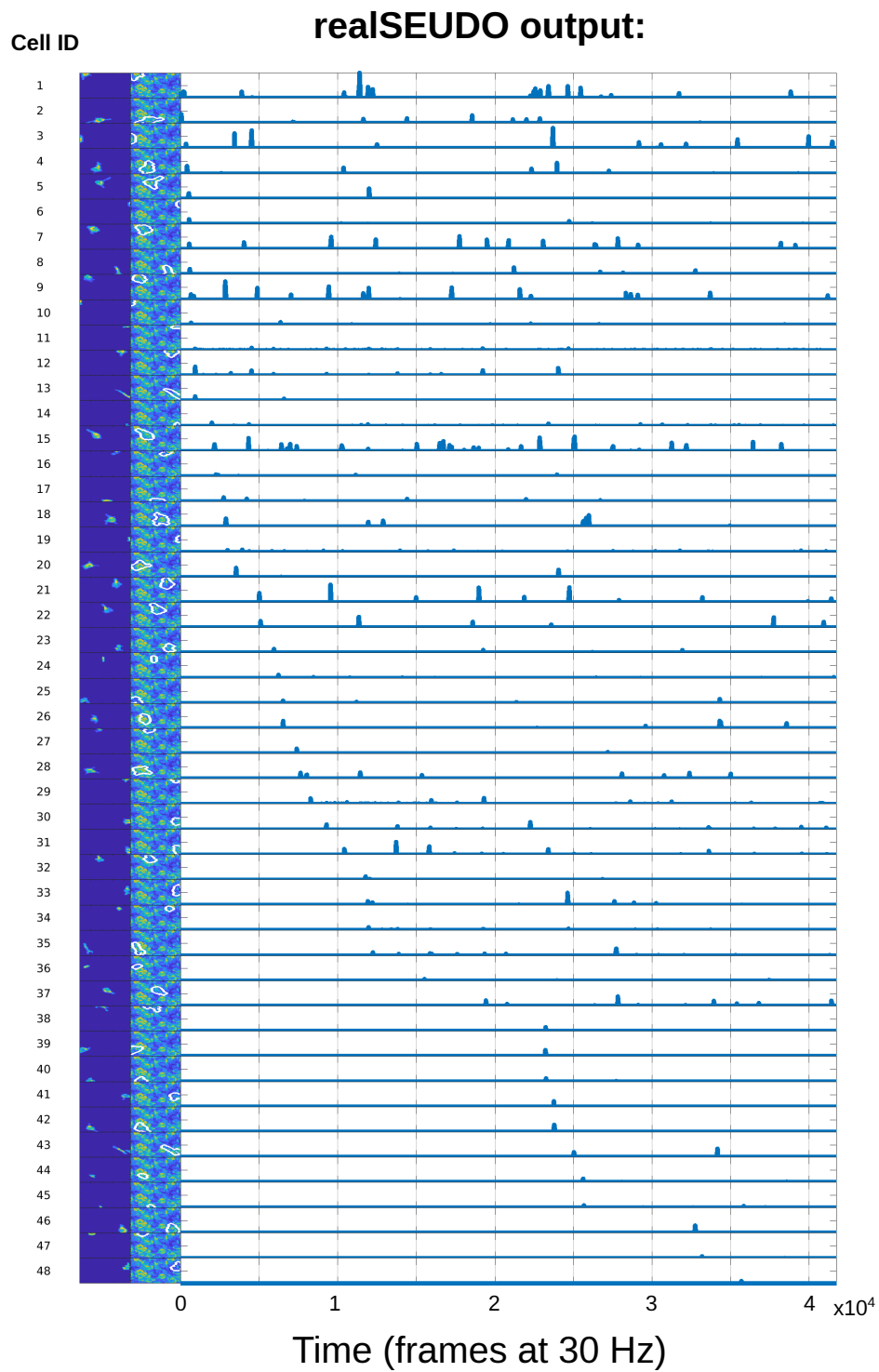


Figure 5: Example realSEUDO cells and traces from a single patch, ordered by the discovery time.

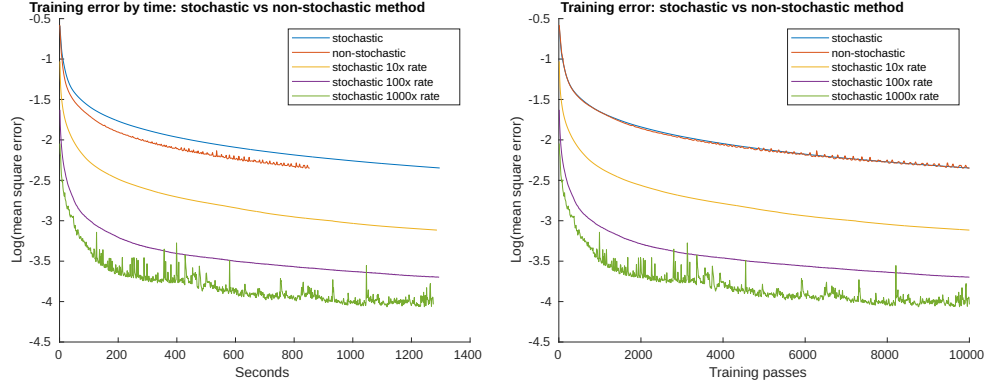


Figure 6: Comparison of different algorithm's learning curves on handwritten datasets. Left: mean-squared error (MSE) as a function of optimization time. Right: MSE as a function of training passes.

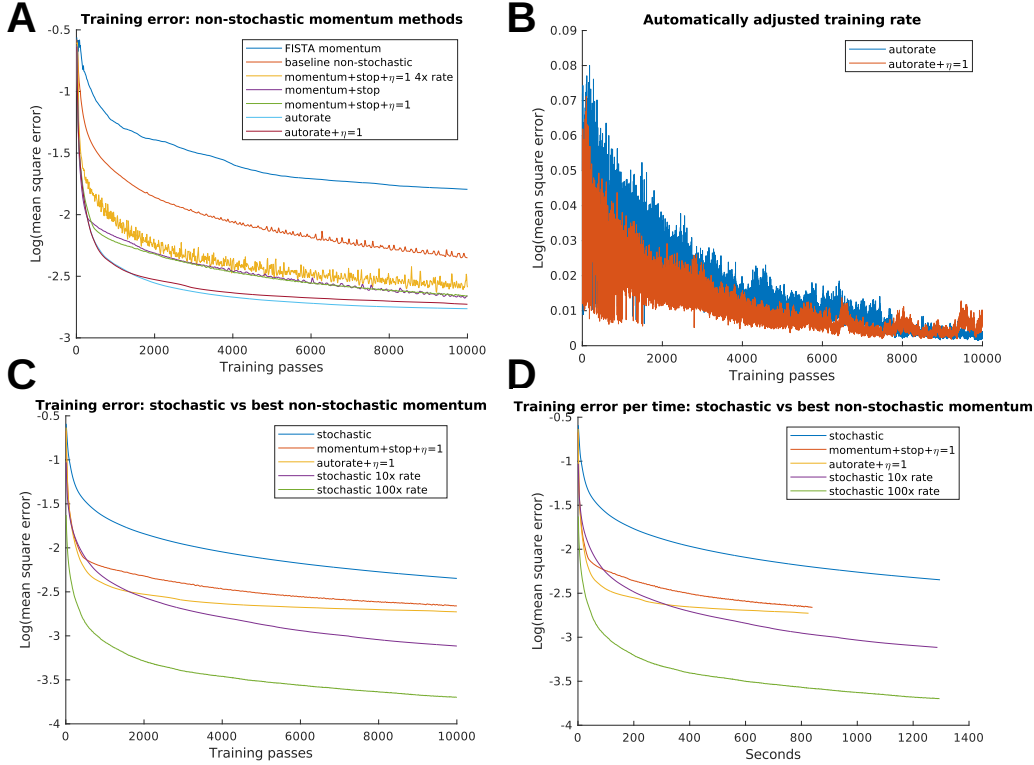


Figure 7: Comparison of training curves for different algorithms. A: Training MSE as a function of training passes for different variants of the improved FISTA algorithm. B: Training MSE improvement when setting $\eta = 1$. C: Training MSE as a function of training passes for the best tested non-stochastic methods vs. momentum-improved FISTA. D: Training MSE as a function of optimization time for the best tested non-stochastic methods vs. momentum-improved FISTA.

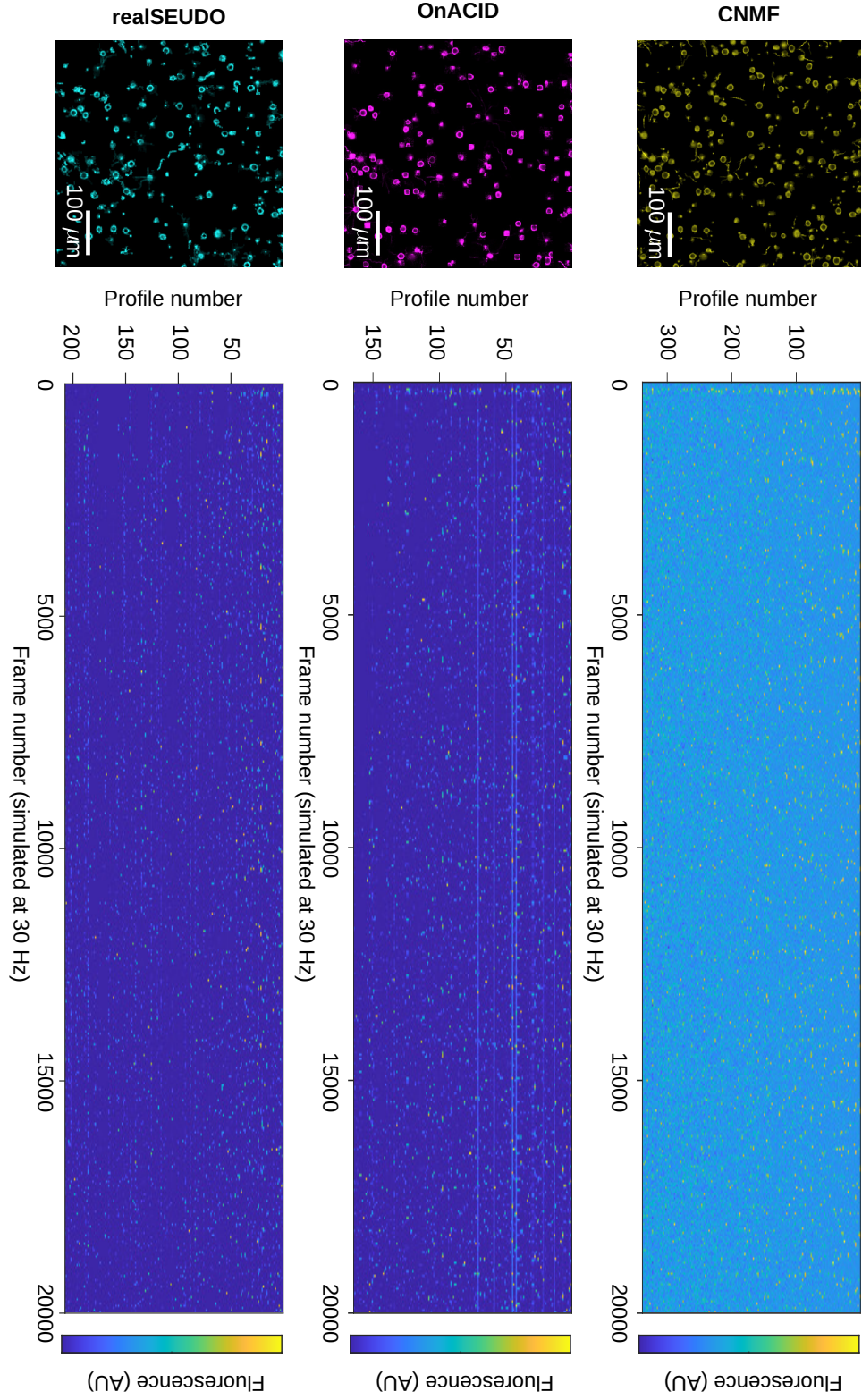


Figure 8: Full sets of strongly-paired traces from CNMF, OnACID and realSEUDO.