

The Top-Scoring ‘N’ Algorithm: A Generalized Relative Expression Classification Method from Small Numbers of Biomolecules

Andrew T. Magis^{1,2} and Nathan D. Price^{1,2}

¹Institute for Systems Biology, Seattle, WA

²Center for Biophysics and Computational Biology, University of Illinois, Urbana, IL

SUPPLEMENTAL SECTION

Decimal Counting System (Fixed radix)						
Place	10^5	10^4	10^3	10^2	10^1	10^0
Max Digit	9	9	9	9	9	9
Digit	0	0	0	0	4	2
						$4 \times 10^1 + 2 \times 10^0 = 42$
Binary Counting System (Fixed radix)						
Place	2^5	2^4	2^3	2^2	2^1	2^0
Max Digit	1	1	1	1	1	1
Digit	1	0	1	0	1	0
						$1 \times 2^5 + 0 \times 2^4 + 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0 = 42$
Factorial Counting System (Mixed radix)						
Place	$5!$	$4!$	$3!$	$2!$	$1!$	$0!$
Max Digit	5	4	3	2	1	0
Digit	0	1	3	0	0	0
						$1 \times 4! + 3 \times 3! + 0 \times 2! + 0 \times 1! + 0 \times 0! = 42$

Figure S1: Counting systems

Representation of a number (42) in three different counting systems: decimal (42), binary (101010), and factorial (13000). Note that decimal and binary are both fixed radix systems, in which the multiplicative distance between each digit place is the same, and the set of digits used for each place is the same. The factorial counting system (factoradics) uses a mixed-radix system in which the multiplicative distance between each digit place is not the same, and the set of digits used for each place is also not the same.

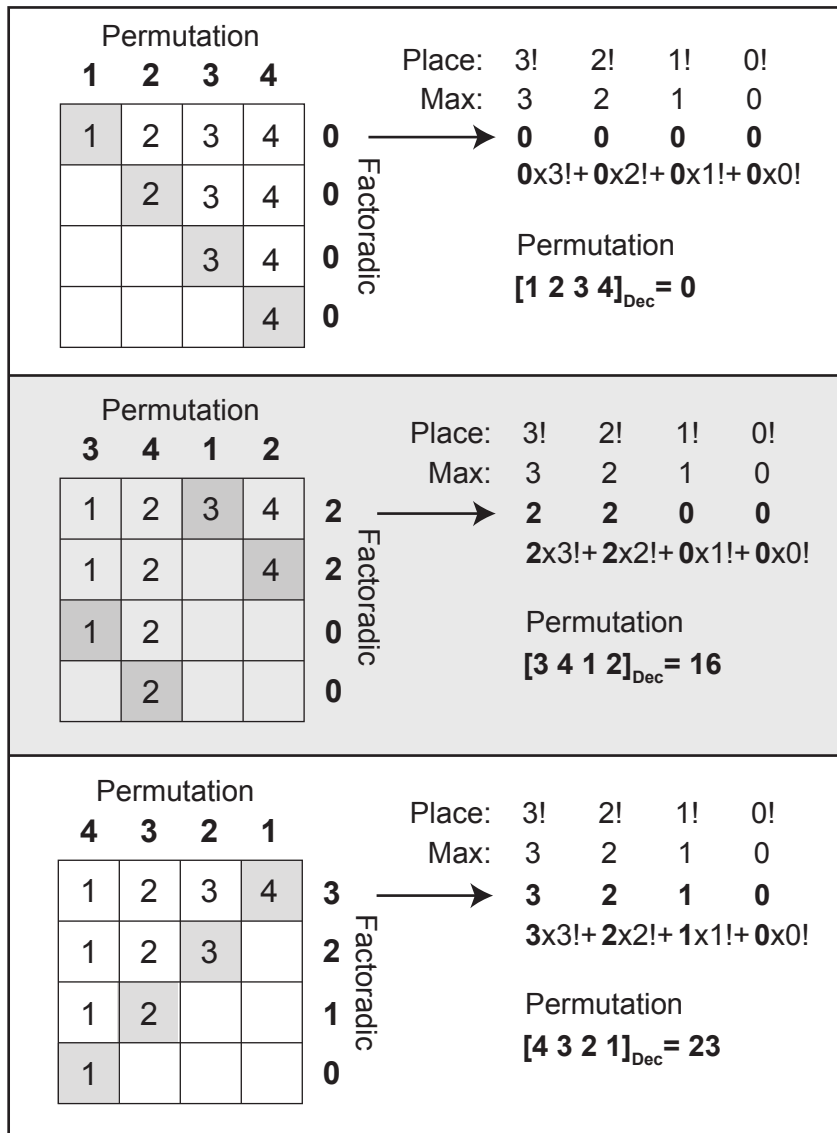


Figure S2: Three complete conversions from permutation to decimal

Three complete translations from permutation to decimal, by way of the factoradic, for a set of size 4. The permutation to be translated is shown above the grid box, and the factoradic representation of that permutation is listed to the right of the grid box. Each row of the grid box determines a digit of the factoradic. The first row of the grid box is the complete sorted list. To perform the translation, consider the elements of the permutation progressively from 1 to 4. The digit of the factoradic is equal to the number of elements of the sorted list to the left of this position. After each digit of the permutation is considered, the corresponding digit of the sorted list is removed. Note that the sorted permutation (**Top**) translates into decimal representation 0. The reverse sorted permutation (**Bottom**) translates into decimal representation $4! - 1 = 23$, the maximum number that can be represented by a factoradic of size 4.

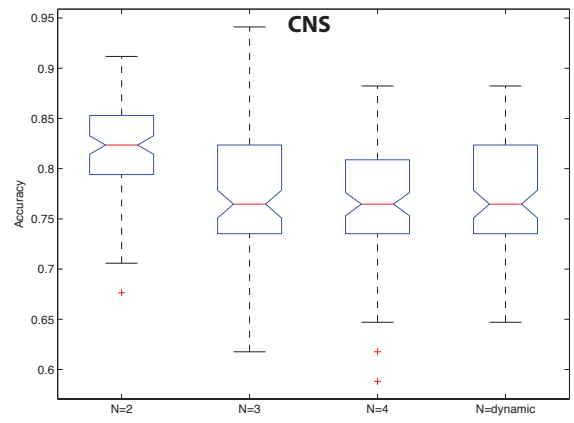
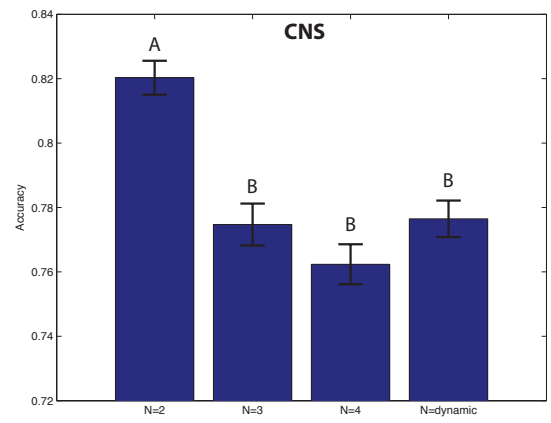
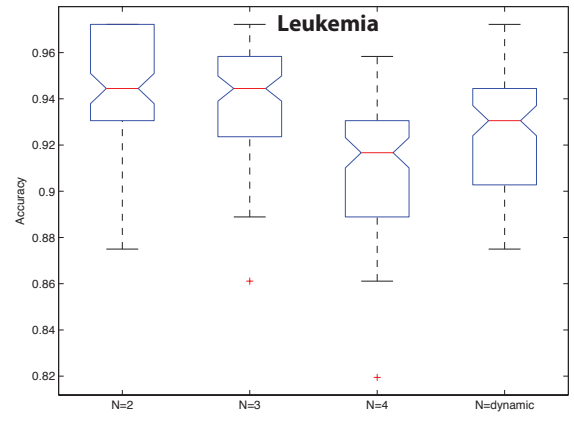
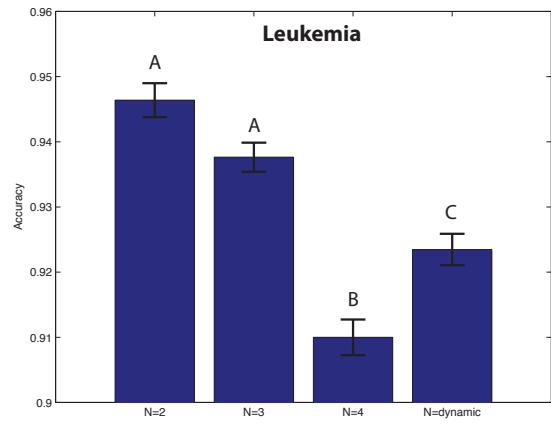
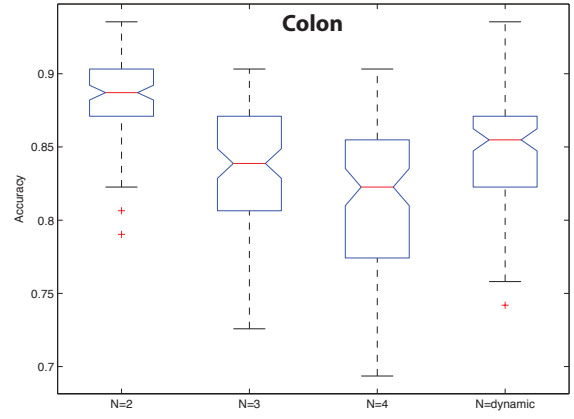
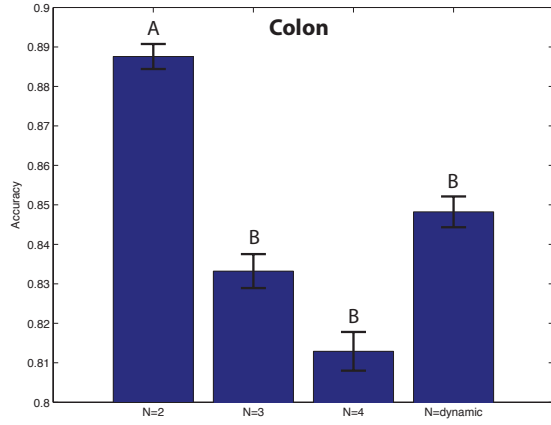
```

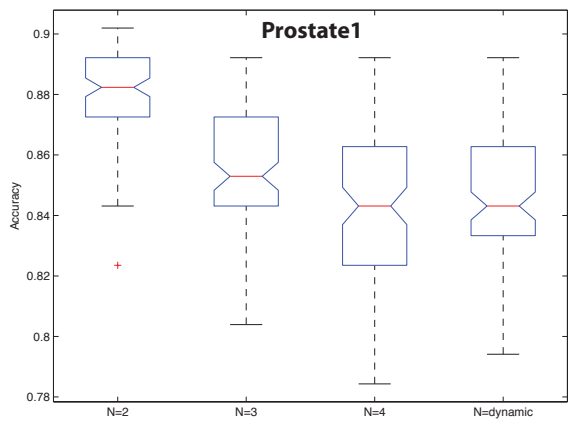
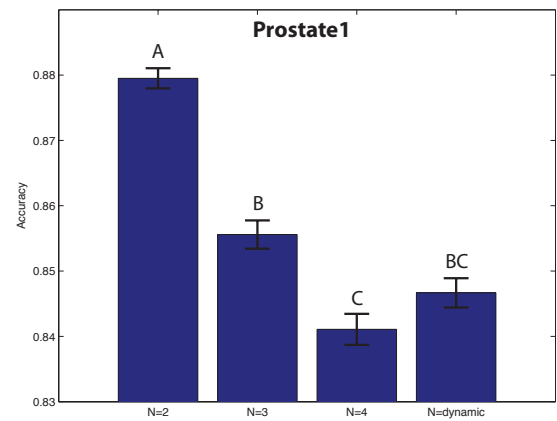
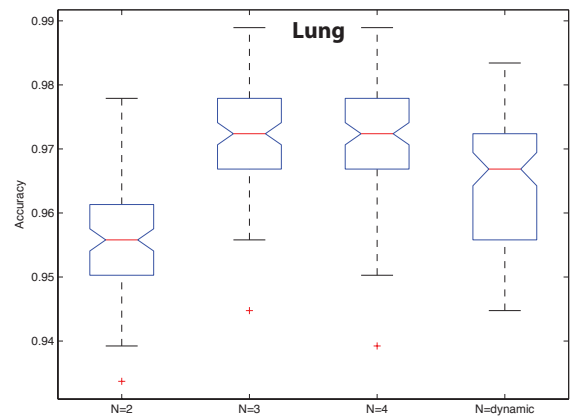
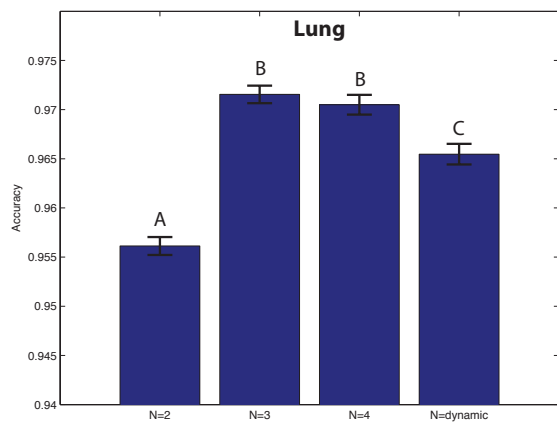
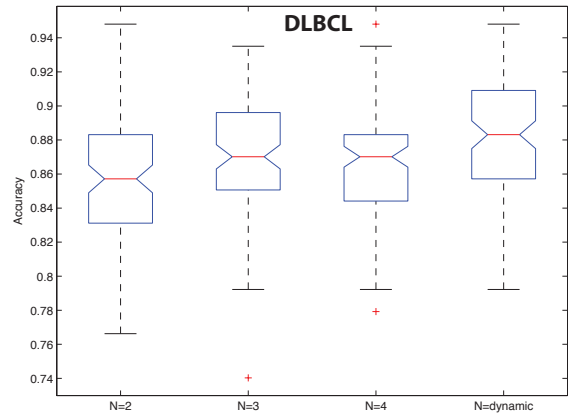
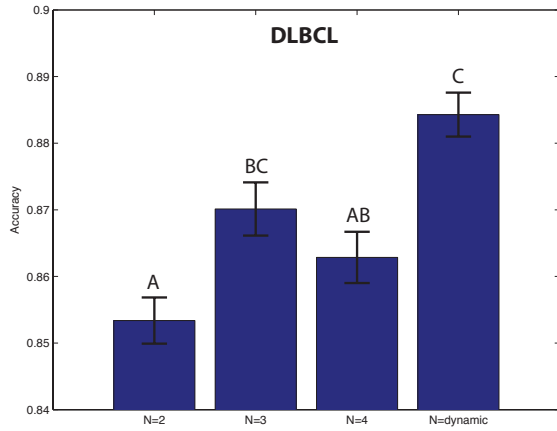
1. INPUTS: C1,C2 = Ranked feature expression data
2.     M = num features of C1,C2
3.     N = classifier size
4. copy C1,C2 to GPU
5. hist1, hist2 = allocate zeroed-out  $N!$  size histograms on GPU
6. num_combinations = (M choose N)
7. num_iterations = num_combinations / GPU memory
8. foreach iteration in num_iterations (CPU)
9.   combs = combinations for this iteration
10.  copy combs to GPU
11.  parallel foreach combination C in combs (GPU)
12.    foreach sample S in C1
13.      fact = factoradic for permutation of C in S
14.      hist1[decimal(fact)] = hist1[decimal(fact)] + 1
15.    foreach sample S in C2
16.      fact = factoradic for permutation of C in S
17.      hist2[decimal(fact)] = hist2[decimal(fact)] + 1
18.    hist1 = normalize(hist1); hist2 = normalize(hist2)
19.    scores = sum(hist1 - hist2)
20.  copy scores, hist1, hist2 to CPU
21.  sort scores and save (CPU)
22. classifiers = top combs over all iterations
23. OUTPUTS: top classifiers, scores, hist1, hist2

```

Figure S3: Pseudocode for the core operation of the TSN algorithm on the GPU.

This pseudocode does not include code to choose the value of N using apparent accuracy or determine classification accuracy using cross validation. Lines 11 through 19 are computed in parallel on the GPU. The GPU is called multiple times depending on the number of combinations requested and the size of the GPU device memory.





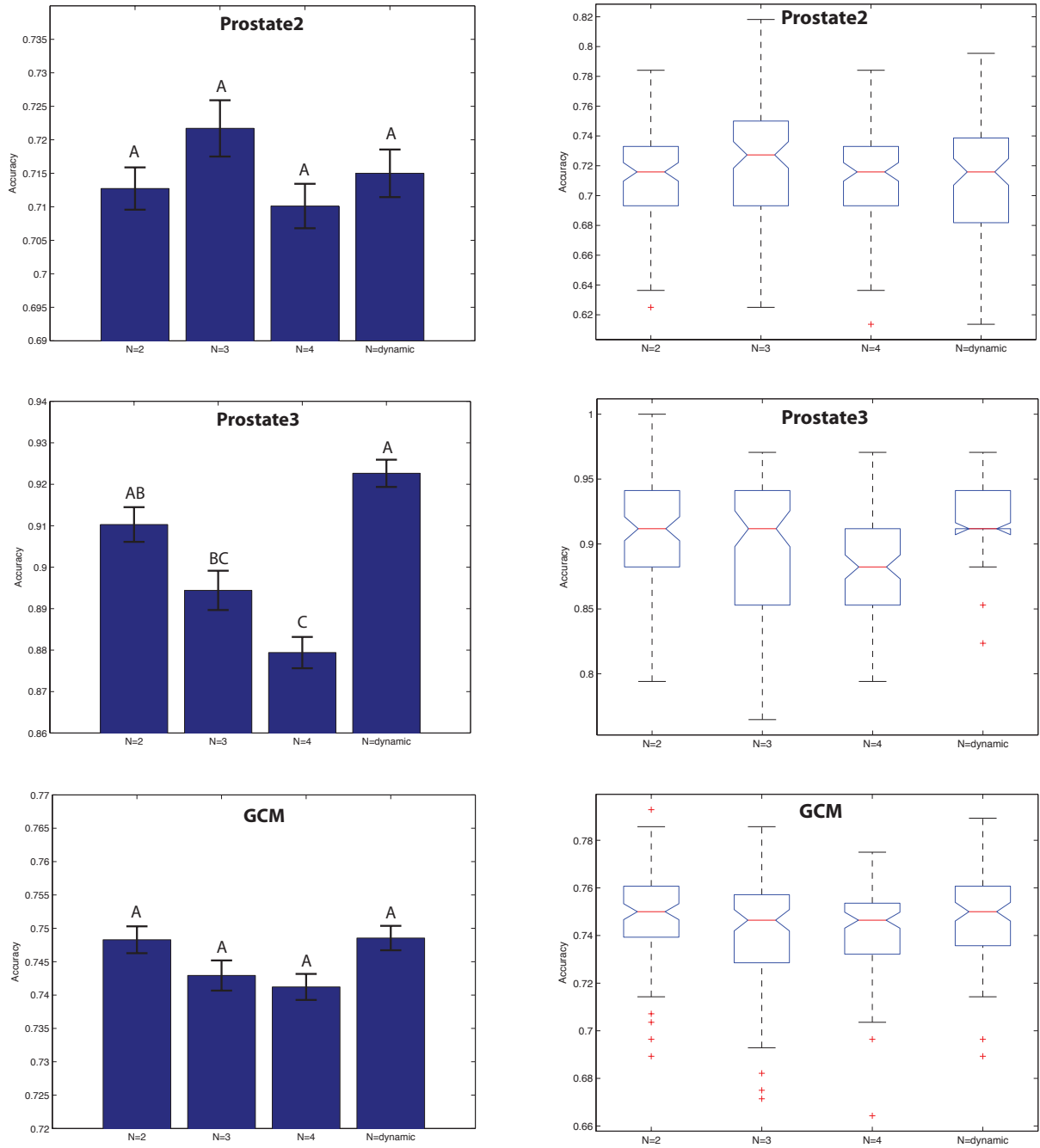
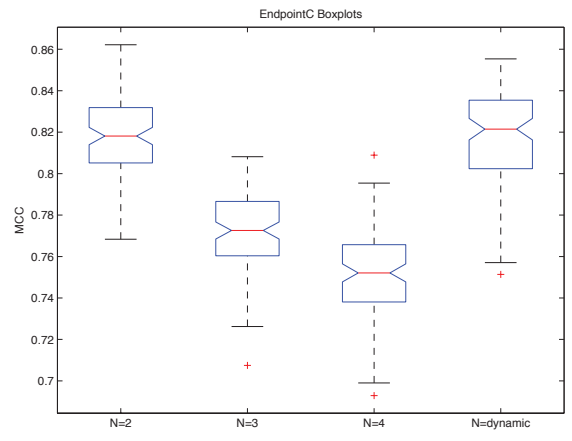
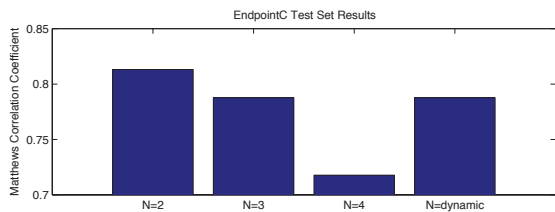
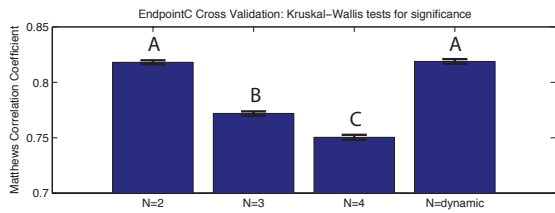
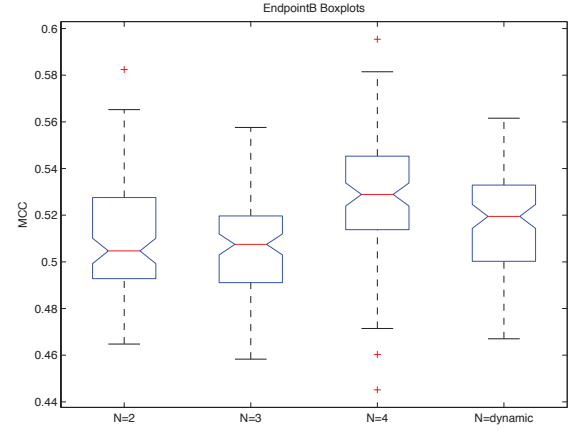
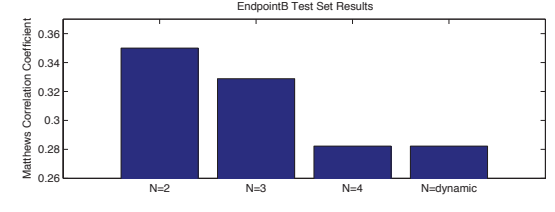
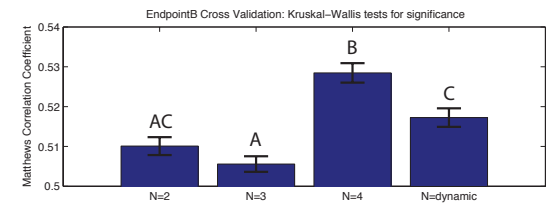
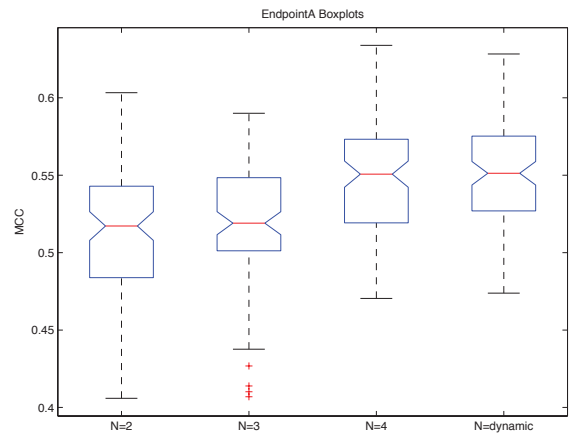
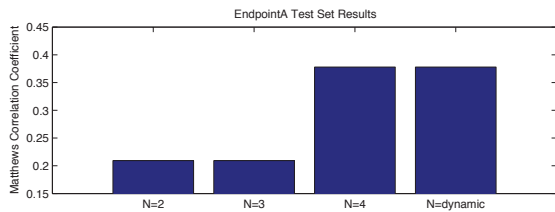
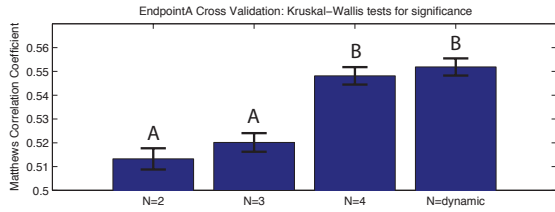
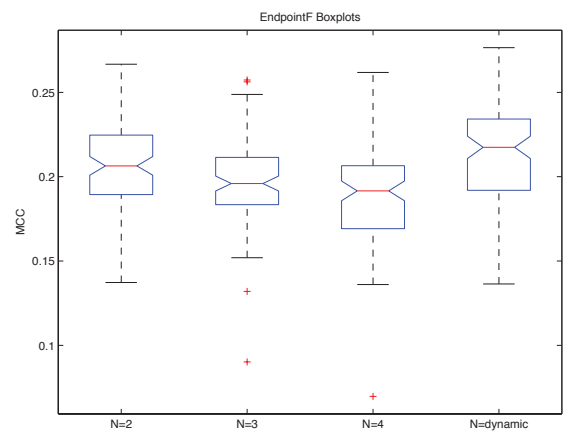
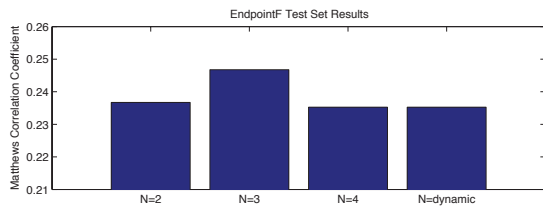
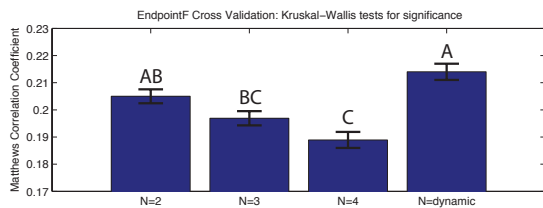
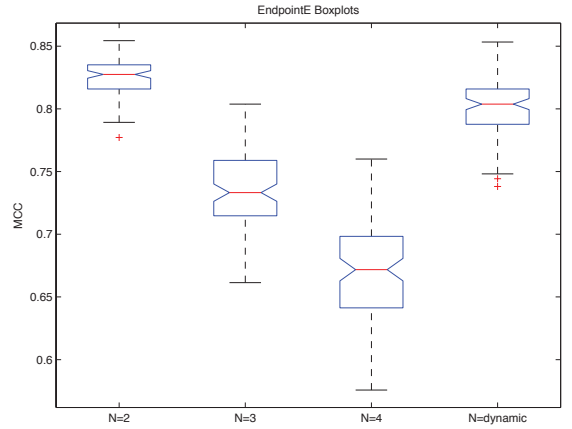
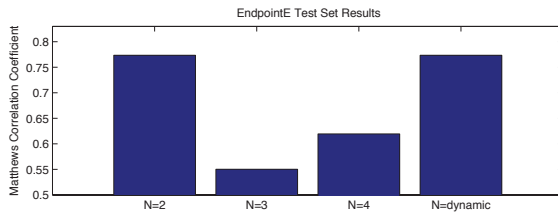
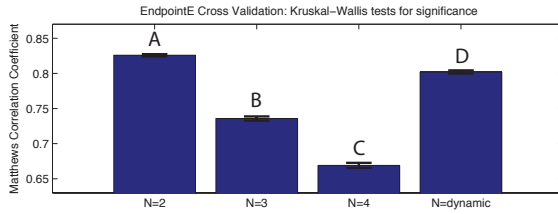
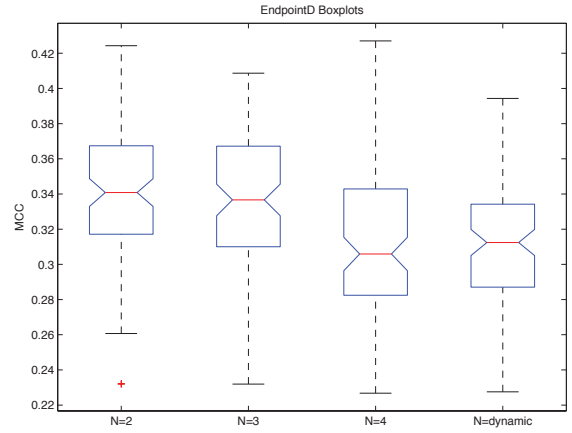
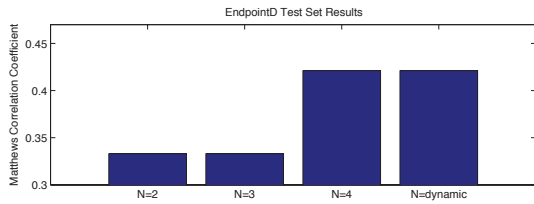
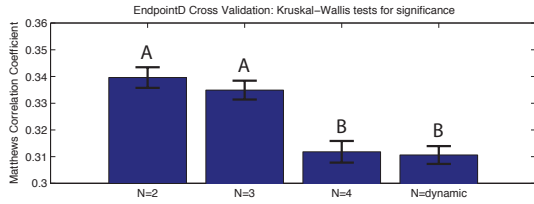


Figure S4: Cancer dataset statistical tests for differences between values of N .

100 iterations of 5-fold cross validation were run using TSN on all nine cancer datasets cited in the paper. Four different types of TSN were run: $N=2$, $N=3$, $N=4$, and *dynamic* N , where the algorithm was allowed to choose the value of N at each iteration of the cross validation loop. These four groups were tested using the nonparametric Kruskal-Wallis test to determine significant differences between the accuracies. A p -value < 0.05 was considered significant. Error bars in the bar plot indicate standard error, not confidence

intervals. Letters above each bar indicate membership in a significance group. If two bars share the same letter, they are not significantly different. If two bars do not share the same letter, they are significantly different. All raw data is included in Additional file 2.





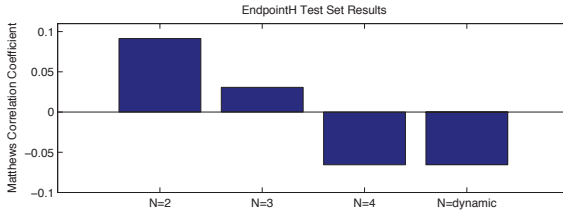
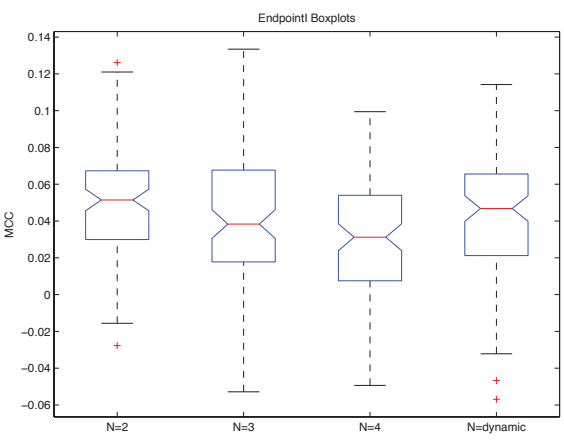
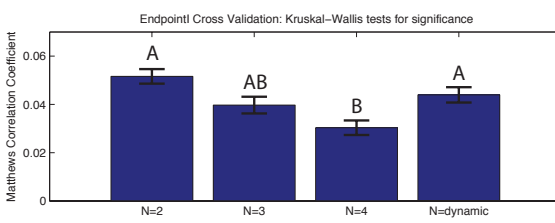
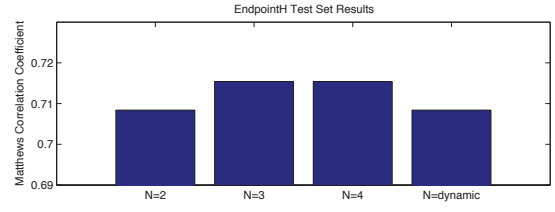
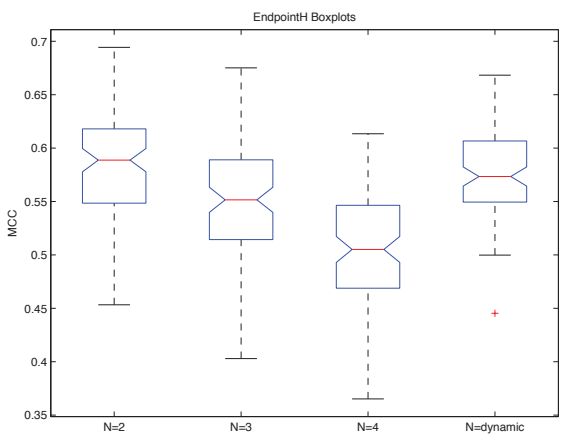
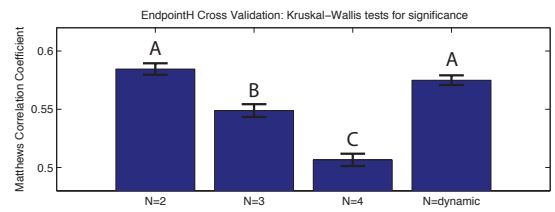
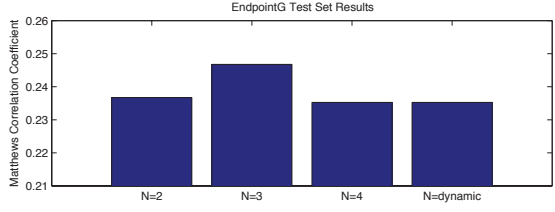
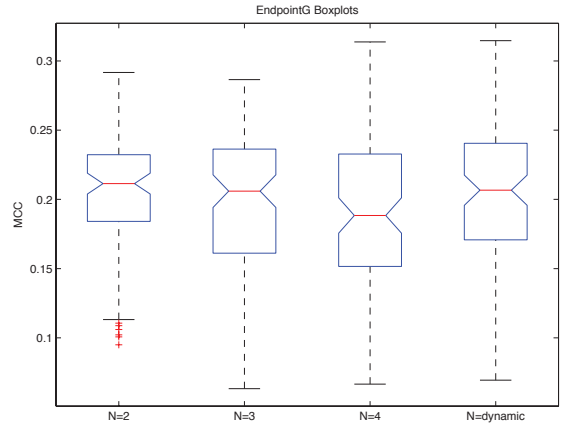
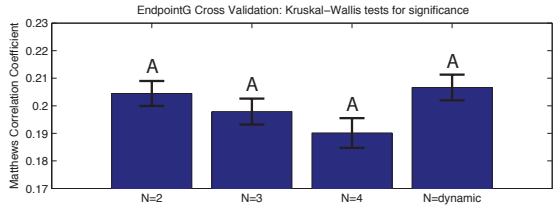


Figure S5: MAQC-II Statistical tests for differences between values of N .

100 iterations of 5-fold cross validation were run using TSN on all nine MAQC-II datasets cited in the paper. Four different types of TSN were run: $N=2$, $N=3$, $N=4$, and *dynamic N*, where the algorithm was allowed to choose the value of N at each iteration of the cross validation loop. These four groups were tested using the nonparametric Kruskal-Wallis test to determine significant differences between the accuracies. A p-value < 0.05 was considered significant. Error bars in the bar plot indicate standard error, not confidence intervals. Letters above each bar indicate membership in a significance group. If two bars share the same letter, they are not significantly different. If two bars do not share the same letter, they are significantly different. All raw data is included in Additional file 4.

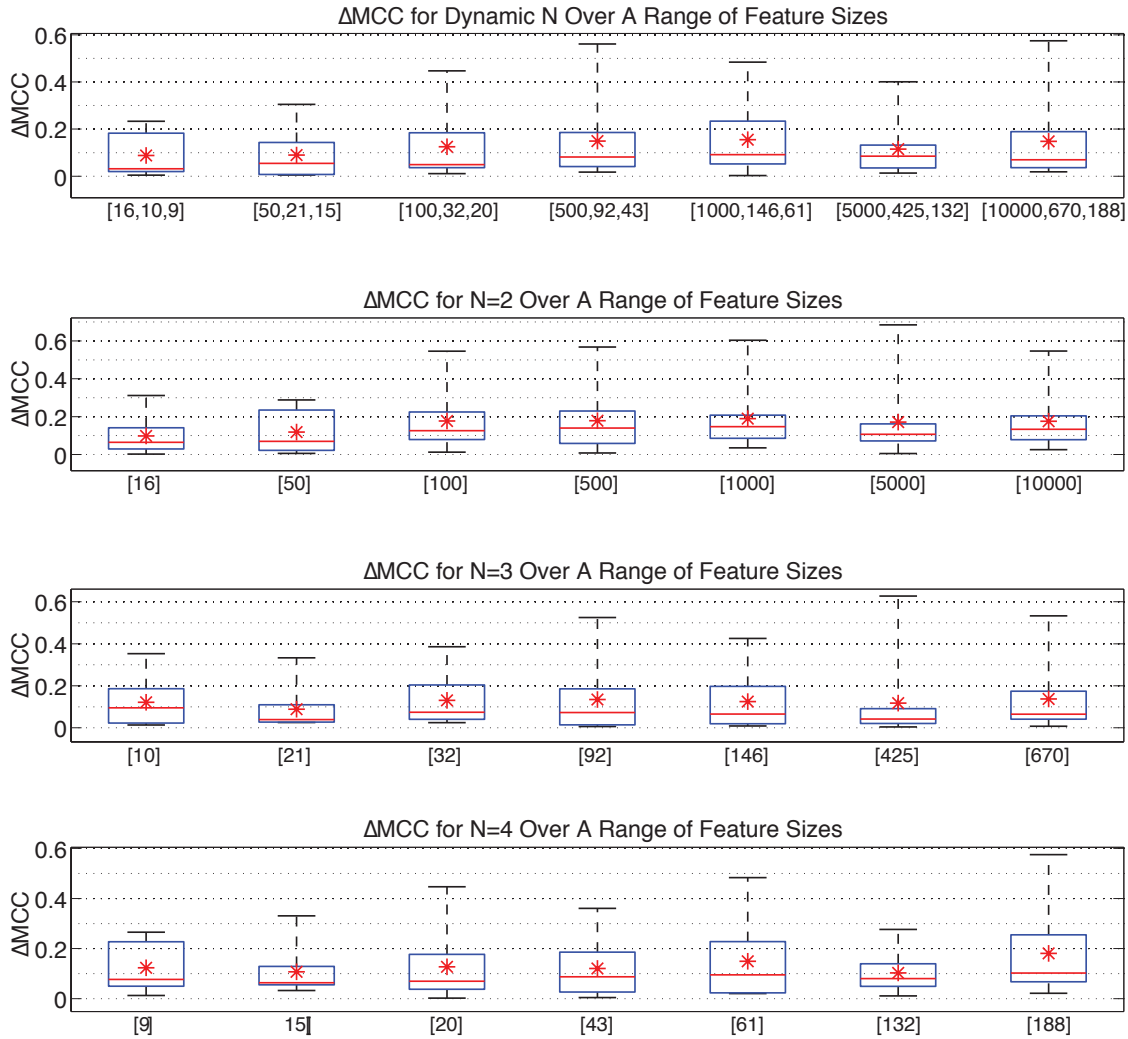


Figure S6: ΔMCC boxplots for different input sizes on the MAQC-II dataset

The TSN algorithm is tested in the paper with very small input sizes (up to 16 features). When the cross validation MCC is compared with the test set MCC, the average difference is very small (Figure 6), indicating that the level of overfitting with the TSN algorithm is low. To determine if the low overfitting of the TSN algorithm is due to the small number of input features we tested all the MAQC-II datasets again using a large range of input sizes. This process was performed for fixed values of $N=2$, $N=3$, and $N=4$, as well as *dynamic* N . The number of features was chosen to yield approximately the same number of combinations for each classifier size. The feature sizes were chosen to span a range of combinations from 120 to 50,000,000. The top panel shows the distribution of ΔMCC scores for *dynamic* N , and the 2nd, 3rd, and 4th panels show the same calculations for fixed $N=2$, $N=3$, and $N=4$, respectively. Despite the wide range of input feature sizes, the mean (asterisk) and median (middle line) ΔMCC scores stay low, among the lowest of any of the MAQC-II participants. The distance between the lower and upper quartiles of the data also remains fairly constant, as indicated by the top and

bottom of the box. The whiskers indicate the extreme values. All raw data is included in Additional file 5.