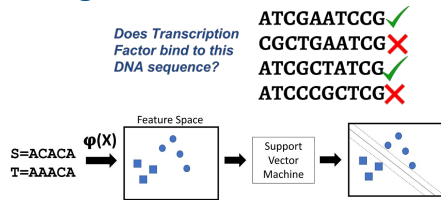


## 1. Overview

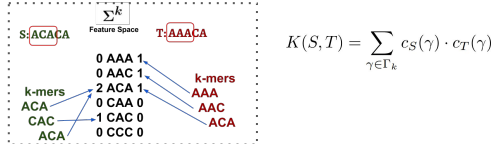
String Kernel techniques, especially those using gapped k-mers as features, have obtained great success in classifying sequences like DNA, protein, and text. However, the state-of-the-art gk-SVM runs extremely slow when we increase the dictionary size ( $\Sigma$ ) or allow more mismatches ( $M$ ). We propose a fast algorithm for calculating Gapped k-mer Kernel using Counting (GaKCo)-

1. *Faster* than state-of-the-art gk-SVM
2. Independent of dictionary size  $\Sigma$  and can scale up to large values of  $M$  and  $\Sigma$ .
3. *Parallelizable*

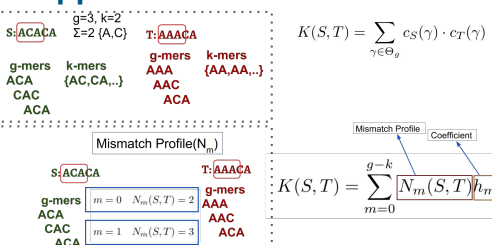
## 2. String Kernel + SVM framework



## 3. Spectrum Kernel



## 4. Gapped k-mer Kernel



## 5. GaKCo

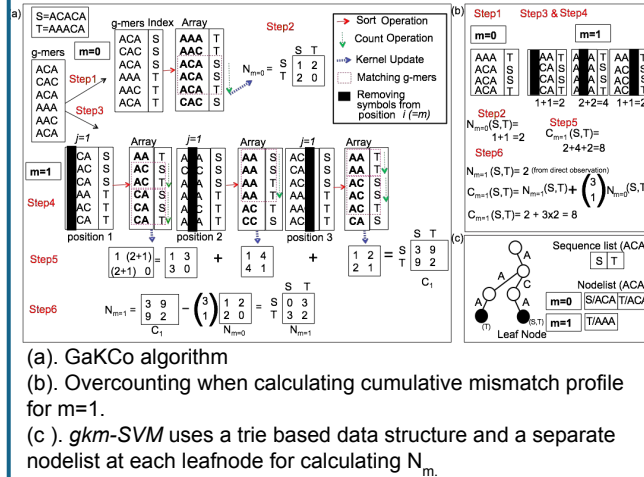
-Algorithm: GaKCo performs g-mer based cumulative counting of cooccurrence to calculate  $N_m$  (independent of dictionary size ( $\Sigma$ )).

-Parallelization: GaKCo groups computations for each value of  $m$  into an independent function, making it naturally parallelizable.

Algorithm 1 GaKCo

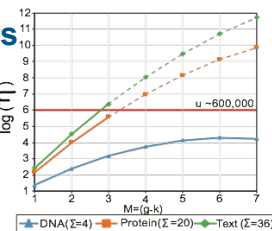
Require:  $L, g, k$   $\Sigma$  L=Array list of g-mers ( $L, g, k$ )

- 1: procedure CalculateKernel
- 2:  $M = g - k$
- 3:  $N = \text{MismatchProfile}(L, g, M)$
- 4:  $K = 0$
- 5: for  $m : 0 \dots M$  do
- 6:  $h_m = \binom{g-m}{k}$
- 7:  $K = K + N_m \cdot h_m$
- 8: procedure MismatchProfile ( $L, g, M$ )
- 9: for  $m = 0 \dots M$  do
- 10:  $C_m = 0$   $\Sigma$  Cumulative Profile
- 11:  $n_{\text{pos}} = \binom{g}{k}$   $\Sigma$  Number of positions
- 12: for  $i : 0 \dots n_{\text{pos}}$  do
- 13:  $C_m = 0$
- 14:  $L' = \text{removePosition}(L, i)$
- 15:  $L' = \text{sort}(L')$
- 16:  $C_m = \text{countAndUpdate}(L')$
- 17:  $C_m = C_m + C_m^L$
- 18: for  $m : 0 \dots M$  do
- 19:  $N_m = C_m$
- 20: for  $j : 0 \dots m-1$  do
- 21:  $N_m = N_m - \binom{g-j}{k-j} N_j$
- 22: return  $N$   $\Sigma N = [N_0, \dots, N_M]$
- 23: Ensure:  $K$   $\Sigma$  Kernel Matrix



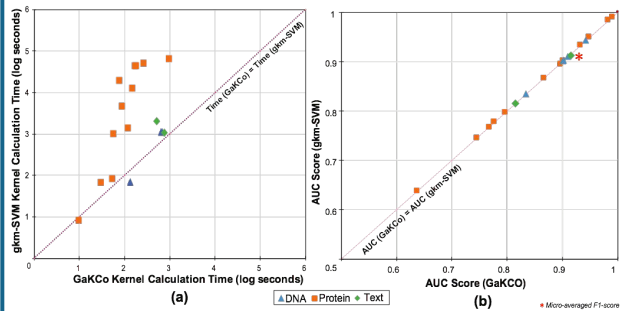
## 6. Theoretical Analysis

	GaKCo	gkm-SVM
Pre-processing	$c_g \cdot gNI$	$gNI + \eta u g$
Kernel updates	$c_g \cdot zN^2$	$\eta u N^2$



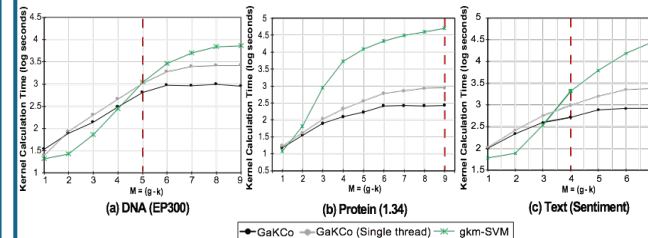
## 7. Experiments and Results

Data: We perform 19 different classification tasks to evaluate the performance of GaKCo. These tasks belong to three categories: (1) Transcription Factor (TF) binding site prediction (DNA dataset), (2) Remote Protein Homology prediction (protein dataset), and (3) Character based English text classification (text dataset).



### Comparisons:

- (a). Kernel Calculation times (log(seconds)) of GaKCo(X-axis) vs gkm-SVM(Y-axis). GaKCo is faster for 16/19 datasets.
- (b). Empirical performance for the same 19 datasets of GaKCo (X-axis) versus gkm-SVM (Y-axis). GaKCo achieves the same AUC-scores as gkm-SVM.



Kernel calculation times with varying mismatches ( $M$ ): scales well with increasing  $\Sigma$  and  $M$ .

### References

1. Mahmoud Ghandi, Dongwon Lee, Morteza Mohammad-Noori, and Michael A Beer. Enhanced regulatory sequence prediction using gapped k-mer features. PLoS Comput Biol, 10(7):e1003711, 2014.
2. Mahmoud Ghandi, Morteza Mohammad-Noori, and Michael A Beer. Robust k-mer frequency estimation using gapped k-mers. Journal of mathematical biology, 69(2):469–500, 2014