

# **FastGSK: Fast and Efficient Sequence Analysis using Gapped String Kernels**

---

A Thesis

Presented to

the faculty of the School of Engineering and Applied Science

University of Virginia

---

in partial fulfillment  
of the requirements for the degree

Master of Science

by

Derrick Blakely

December 2019

# APPROVAL SHEET

This Thesis  
is submitted in partial fulfillment of the requirements  
for the degree of  
Master of Science

Author Signature: 

This Thesis has been read and approved by the examining committee:

Advisor: Yanjun Qi

Committee Member: Yangfeng Ji

Committee Member: Samira Khan

Committee Member: \_\_\_\_\_

Committee Member: \_\_\_\_\_

Committee Member: \_\_\_\_\_

Accepted for the School of Engineering and Applied Science:



Craig H. Benson, School of Engineering and Applied Science

December 2019

---

# FastGSK: Fast and Efficient Sequence Analysis using Gapped String Kernels

---

**Derrick Blakely**

Department of Computer Science  
University of Virginia  
Charlottesville, VA 22903  
dcb7xz@virginia.edu

## Abstract

String kernel methods achieve strong classification performance on DNA, protein, and natural language data using modestly-sized training sets. However, existing kernel function algorithms suffer from slow kernel computation times, as they depend exponentially on the sub-sequence feature length and the task’s alphabet size. In this work, we introduce a new string kernel algorithm using gapped  $k$ -mers called FastGSK. Compared to previous string kernels, it uses a simplified kernel algorithm that decomposes into a set of independent counting operations over the possible mismatch positions. This algorithm is naturally parallelizable, is performant on any sequence analysis task, and allows us to devise a fast Monte Carlo approximation method to scale to much greater feature lengths. We evaluate FastGSK on 10 DNA transcription factor binding site (TFBS) prediction tasks, 10 protein remote homology detection tasks, and 7 English-language medical named entity recognition tasks. FastGSK consistently matches or outperforms the state-of-the-art string kernel algorithms in AUC, while achieving average speedups in kernel computation of  $\sim 100\times$  and speedups of  $\sim 800\times$  for large feature lengths. We further show that FastGSK consistently outperforms character-level recurrent neural networks across these sequence analysis tasks. Our algorithm is available as a Python package and as C++ source code<sup>1</sup>.

## 1 Introduction

String kernels in conjunction with Support Vector Machines (SK-SVM) achieve strong prediction performance across a variety of sequence analysis tasks, with widespread use in bioinformatics and natural language processing (NLP). SK-SVMs are a dominant technique for DNA regulatory element identification [10, 23, 7, 18, 25, 8], protein function prediction [25, 6], and medical named entity recognition [21]. SK-SVMs are also popular baselines for evaluating the quality of deep learning models [2, 10, 23] and because they extract interpretable sequence features, they are popular for analyzing variant impacts and explaining learning methods [24, 18].

The key to the success of string kernel methods is their use of simple, yet expressive, substring features to compute a similarity function between sequences. In turn, the similarity function defines an inner product space, where an SVM classifier can be trained. The approach easily enables comparison of arbitrary length sequences, obviates sequence alignment issues, captures task-relevant pattern information, and is simpler than other pattern detection tools, such as position-weight matrices [26, 1]. Viewed as a type of "feature engineering" and model bias, string kernels yield simpler and lower variance models than deep learning. One consequence is that they show strong performance without consuming vast amounts of training data (for example, see figure 1).

In greater detail, string kernels use substring features to map sequences to fixed-dimension feature vectors. Such a mapping is referred to as a *feature map*. For example, one of the most popular string kernels is the  $k$ -spectrum kernel [20], which maps sequences to vectors using contiguous length- $k$  substrings, or  $k$ -mers. In this vector space, the  $i$ -th component of a sequence’s feature vector is simply the number of times the  $i$ -th possible  $k$ -mer occurs in that sequence. The *kernel function* then takes a pair of sequences and returns their similarity score—the inner product of their feature

---

<sup>1</sup>Available for download at <https://github.com/Qdata/FastSK/> or with the command `pip install fastsk`

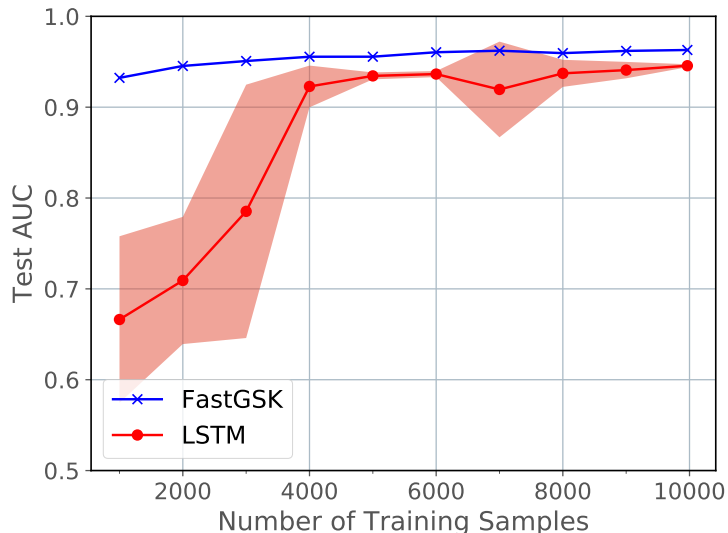


Figure 1: String kernel methods obtain excellent results, greatly outperforming recurrent neural networks given modest training sets. We vary the proportion of training samples used for a DNA sequence classification task (TFBS prediction on the ZZZ3 ENCODE ChIP-Seq dataset) used for model training. Training sizes vary from 1,000 to 9,996 samples (the maximum). Each LSTM point is the average of 5 runs, with the shaded region showing standard error. This error exemplifies the high model variance endemic to deep learning when training data is limited.

vectors in the vector space induced by the feature map. Importantly, a kernel function serves a "kernel trick" in the Support Vector Machine framework. In turn, this provides several benefits of its own: SVMs learn an optimal separating hyperplane to classify sequences and show excellent generalization [27]. Moreover, they obtain excellent classification performance without needing as many training samples as state-of-the-art deep learning models (see figure 1).

The spectrum string kernel and derivatives such as the  $(k, m)$ -mismatch kernel are responsible for many of the successes of string kernels. However, a major problem for spectrum kernels is that the dimensionality of the feature space is  $|\Sigma|^k$ , where  $|\Sigma|$  is the alphabet of characters appearing in the sequences. This presents three major problems. First, kernel computation becomes infeasible for even modest values of  $|\Sigma|$  or  $k$ . For example, spectrum kernels are infeasible for roughly  $k > 10$ , yet many transcription factor binding sites are up to 20 basepairs long. Furthermore, spectrum methods scale poorly to large alphabets, such as protein or natural language [4, 19, 28]. The second challenge is that as  $k$  increases, the odds of observing any particular  $k$ -mer within a sequence rapidly goes to zero. Therefore, the feature vectors are both extremely large and extremely sparse, which makes models trained on these feature vectors highly prone to overfitting and poor generalization. Third, existing algorithms that overcome these challenges leave much to be desired; for example, popular trie-based approaches still exhibit exponential dependence on  $|\Sigma|$  [19, 25],  $k$ , and  $m$ . On the other hand, counting-based methods rely on complex "mismatch statistics" to indirectly obtain feature counts [16, 6, 25] and still fail to scale to greater feature lengths.

Together, these issues present major limitations to the practical utility of  $k$ -mer string kernel methods. To solve these problems, we introduce a novel string kernel algorithm called FastGSK. It makes four high-level contributions:

1. We take inspiration from [7] and use *gapped  $k$ -mer* features, or *gkmers* for short, which are a more compact feature set that greatly reduces the size of the feature space and risk of overfitting. These gapped  $k$ -mer features consist of an overall length of  $g = k + m$ , length- $k$  non-contiguous substrings within the features, and  $m$  mismatch positions.
2. We decompose the kernel function into a set of  $\binom{g}{m}$  independent counting operations to count the gapped  $k$ -mer features shared between sequences. This formulation is simpler and faster than the state-of-the-art string kernel algorithms. Furthermore, the independence of the counting operations mean that the algorithm is naturally parallelizable. We exploit this advantage to create a fast multithreaded implementation.
3. We take advantage of the independence of each mismatch position to create a scalable Monte Carlo approximation algorithm. This approach randomly samples the possible mismatch positions until the variance of the kernel matrix converges. We show the approximation algorithm, called FastGSK-Approx, converges rapidly irrespective of the parameters  $g$  and  $m$ . Therefore, FastGSK-Approx allows scaling to greater feature lengths.
4. Empirically, we show that FastGSK matches or outperforms state-of-the-art string kernel methods and LSTM models in test performance across 10 DNA TFBS datasets, 10 protein remote homology detection datasets, and 7 medical named entity recognition tasks.

## 2 Background

### 2.1 Support Vector Machines

Support Vector Machines learn a linear predictive model  $f(x) = \hat{y} = \mathbf{x} \cdot \mathbf{w} + b$ . In the case of linearly-separable data, SVMs optimize the parameters  $\mathbf{w}$  by learning a pair of max-margin hyperplanes given by:

$$\mathbf{x} \cdot \mathbf{w} + b = 1 \quad (1)$$

and

$$\mathbf{x} \cdot \mathbf{w} + b = -1 \quad (2)$$

This is achieved by minimizing  $\|\mathbf{w}\|_2$ ; because we want to maximize the distance between the planes, which is  $2/\|\mathbf{w}\|_2$ , we minimize  $\|\mathbf{w}\|_2$ . To keep training points from being inside the margin, we also impose the constraint:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i) \geq 1 \text{ for } 1 \leq i \leq n \quad (3)$$

where  $n$  is the number of training samples.

A non-linear, or kernelized, SVM follows roughly this same structure, but uses a kernel function  $K(\cdot, \cdot)$  to compute the pair-wise similarities between samples. As such, a string kernel function  $K$  is easily applied in the Support Vector Machine (SVM) framework [27]. In this case, the predicted class of a sample  $x$  is given by

$$f(x) = \sum_{i=1}^n \alpha_i y_i K(x_i, x) + b \quad (4)$$

where  $x_i$  and  $y_i$  are the  $i$ th training sample and its label, respectively. Each  $\alpha_i$  is a weight, where if  $\alpha_i \neq 0$ ,  $\alpha_i$  corresponds to a *support vector* and  $b$  is a learned additive bias.

### 2.2 String Kernels

String kernel methods compare arbitrary-length sequences by mapping them to a fixed inner product space. The key component is the feature map  $\phi : \mathcal{S} \rightarrow \mathbb{R}^p$ , where  $\mathcal{S}$  is the set of all strings composed from the alphabet  $\Sigma$ . The dimensionality  $p$  of the vector space depends on the particular string kernel's feature map. The canonical example is the spectrum kernel by [20], which uses simple length- $k$  substrings, or  $k$ -mers. Given a string  $x = (s_1, s_2, \dots, s_{|x|})$  with each  $s_i \in \Sigma$  and a substring size  $k$ , the spectrum kernel  $\phi_S$  maps  $x$  to a vector indexed by all possible length- $k$  substrings from  $\Sigma^k$ . Under this feature map, the  $i$ th dimension of the vector  $\phi_S(x)$  holds the number of times the  $i$ th possible  $k$ -mer  $\in \Sigma^k$  occurs in  $x$ . The spectrum kernel function  $K_S$  then provides a similarity score of two sequences  $x$  and  $y$  as the inner product of their spectrum feature vectors:

$$K_S(x, y) = \langle \phi_S(x), \phi_S(y) \rangle \quad (5)$$

A string kernel function  $K$  is easily applied in the Support Vector Machine (SVM) framework [27]. Importantly, the inner product of a string kernel can be defined *implicitly* as part of the kernel function. That is, without invoking the explicit mapping of strings to their full feature vectors in the vector space. For example, the spectrum kernel is also given by

$$K_S(x, y) = \sum_{\alpha \in \Sigma^k} c_x(\alpha) c_y(\alpha) = \langle \phi_S(x), \phi_S(y) \rangle \quad (6)$$

where  $c_x(\alpha)$  and  $c_y(\alpha)$  return the counts of  $k$ -mer  $\alpha$  in sequences  $x$  and  $y$ , respectively. This view provides an important intuition: the feature function  $\phi(\cdot)$  can be evaluated *implicitly*; that is, without fully mapping the sequences to their feature vectors. Another important intuition is that a  $k$ -mer  $\alpha$  only contributes to  $K_S(x, y)$  if it is present in *both* sequences; that is, both  $c_x(\alpha)$  and  $c_y(\alpha)$  have to be non-zero. Therefore, evaluation of  $K_S(x, y)$  reduces to counting the co-occurrences of each possible  $\alpha$ . In fact, this view shows we must compare the  $k$ -mers between samples in order to determine if some  $\alpha$  is present in both. Finally, to allow imprecise matching (e.g., permitting robustness to noise or single nucleotide polymorphisms), we can permit some number of *mismatches* when comparing the  $k$ -mers between sequences. This approach, called the  $(k, m)$ -mismatch kernel, defines the kernel function as

$$K_{(k,m)}(x, y) = \sum_{\alpha \in \Sigma^k} c_x(\alpha; m) c_y(\alpha; m) \quad (7)$$

where  $c_x(\alpha; m)$  is the number of times the  $k$ -mer  $\alpha$  appears in the sequence  $x$  with up to  $m$  mismatches. In the next section, we further elaborate upon this formulation.

## 2.3 (k,m)-mismatch Kernel

Since the introduction of the spectrum kernel, many string kernel variants and generalizations have been proposed, usually involving *mismatches* to incorporate noise into string comparisons. For example, the  $(k, m)$ -mismatch kernel from [19] retains the  $k$  parameter for substring lengths, while adding an  $m$  parameter to denote a number of *mismatches* permitted when comparing the  $k$ -mers of a pair of sequences.

A simple generalization of the  $k$ -spectrum kernel is the  $(k, m)$ -mismatch kernel [19], which permits up to  $m$  *mismatches* when determining if a pair of  $k$ -mers should contribute to the similarity of their respective sequences. Under the  $(k, m)$ -mismatch feature map, a string  $x$  is mapped to a  $|\Sigma|^k$ -dimensional space by

$$\phi_{(k,m)}(x) = \left( \sum_{\alpha \in x} I_m(\alpha, \gamma) \right)_{\gamma \in \Sigma^k} \quad (8)$$

where  $I_m(\alpha, \gamma) = 1$  if the  $k$ -mer  $\gamma$  is in the "mismatch neighborhood" of  $\alpha$ , denoted by  $N_{k,m}(\alpha)$ . The mismatch neighborhood  $N_{k,m}(\alpha)$  is simply the set of all  $k$ -mers that differ from  $\alpha$  by at most  $m$  characters. The  $i$ th index of  $\phi_{(k,m)}(x)$  is simply a count of how many times the  $i$ th possible  $k$ -mer occurs in  $x$  if we allow up to  $m$  mismatches.

An influential approach by [16] uses this formulation to compute the kernel function *indirectly*. Intuitively, the similarity of sequences  $x$  and  $y$  is given by how many "neighboring"  $k$ -mers they share. Following this intuition, the trick is to compute the kernel function by counting how many  $k$ -mers from sequence  $x$  are contained in the mismatch neighborhoods of sequence  $y$ 's  $k$ -mers. As shown in [16], this reduces to counting the number of  $k$ -mers shared between  $x$  and  $y$  at each Hamming distance  $d \in \{0, 1, \dots, \min(2m, k)\}$  and then multiplying the counts (which the authors call "statistics") by an appropriate combinatorial coefficient. We dub this approach and its derivatives "mismatch statistic string kernels." The upside of the mismatch statistic approach is that it works well in the case of the  $(k, m)$ -mismatch kernel and is not computationally dependent on the alphabet size  $|\Sigma|$  (the feature space, however, is). The downside is that it has been applied to situations where it is not actually beneficial. We show this in section 3.4.

As for the  $(k, m)$ -mismatch kernel generally, there are significant shortcomings. First, because the feature space is of size  $|\Sigma|^k$ , operating in this space becomes deleterious for even moderately sized  $\Sigma$  or  $k$ . Second, this is an extremely sparse feature space, as the probability of any particular  $k$ -mer appearing in a sequence quickly approaches 0 as  $k$  grows. As such,  $(k, m)$ -mismatch SK-SVMs are highly prone to overfitting. Third, most implementations use trie-based data structures, which also grow exponentially with  $\Sigma$  and  $k$  [4, 19, 28]. Implementation that do not use tries often use a complex set of *mismatch statistics* to *indirectly* compute the co-occurrence counts [15, 16, 6].

Ultimately, none of these approaches are effective at meeting all of the following three criteria: (1) Feature set that is not exponential in the alphabet size  $|\Sigma|$ . (2) Fast kernel computation algorithm that is scalable in  $\Sigma$ , and conceptually simple. (3) It should scale to greater feature lengths and numbers of mismatches. For example, efficiently handling transcription factor binding sites with 20 basepairs.

## 3 Proposed Method

### 3.1 Gapped k-mer Kernel

Like [7], we use a gapped  $k$ -mer, or gkmer, feature set. These features have three parameters: an overall feature length  $g$ ,  $m$  mismatch (or gap<sup>2</sup>) positions, and  $k$  informative, non-mismatch positions (note that  $k + m = g$ ). For example, as shown in figure 2, the string  $S = ACACA$ , contains the  $g$ -mer  $g_1 = ACA$ , where  $g = 3$ . Now, the parameter  $k$  specifies a number of informative, non-mismatch positions within the  $g$ -mer. For example, if  $k = 2$ , then  $\{AC_, A_A_, _CA\}$  are the set of possible *gapped k-mers* within  $g_1$ . We use  $m$  (in this case  $m = 1$ ) to denote the number of mismatch positions inside the  $k$ -mers. As with the  $k$ -spectrum and  $(k, m)$ -mismatch kernel, our kernel function is determined by the co-occurrences of  $k$ -mers, except in our case the  $k$ -mers are not contiguous features. The gapped  $k$ -mer string kernel function is given by

$$K_{GSK}(x, y) = \sum_{\gamma \in \Theta_{g,m}} c_x(\gamma) c_y(\gamma) \quad (9)$$

where  $\Theta_{g,m}$  is the set of gapped  $k$ -mers with  $m$  mismatch positions appearing in the dataset; in contrast, to the spectrum methods, we do not need to consider the entire feature space. The function  $c_x(\gamma)$  gives the count of the gapped  $k$ -mer  $\gamma$  in  $x$ .

<sup>2</sup>We hereafter use "gap" and "mismatch" interchangeably.

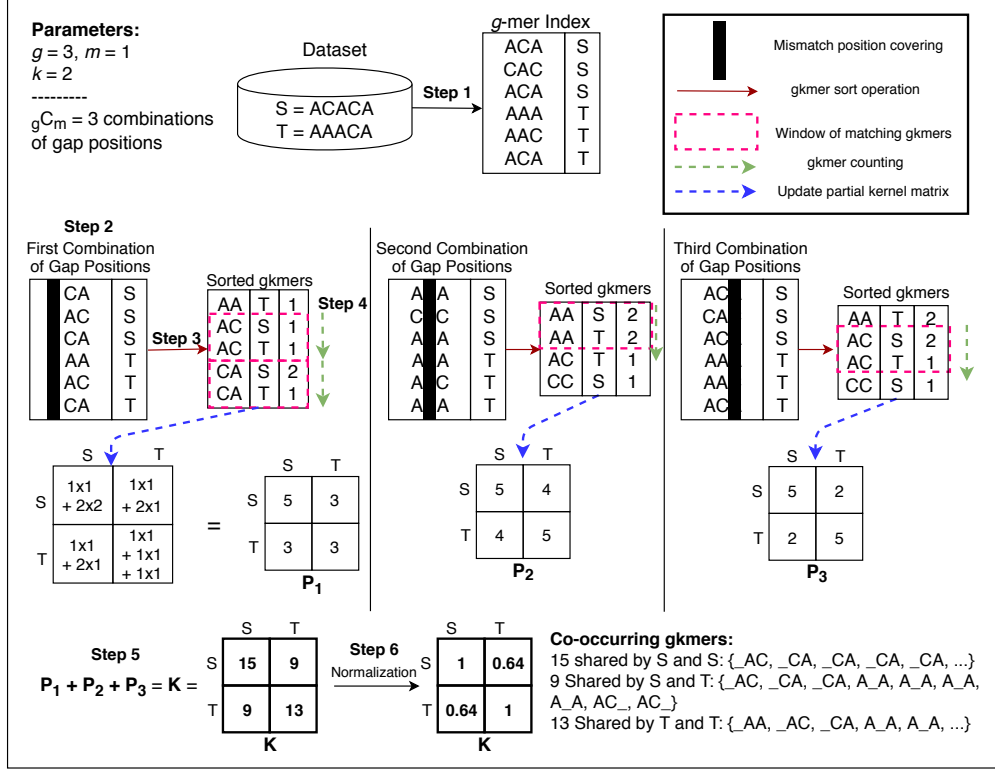


Figure 2: A demonstration of the FastGSK-Exact algorithm for  $g = 3, m = 1$ , and  $k = 2$  using a small dataset with only two sequences (S and T). (1) All  $g$ -mers are extracted from the dataset and stored in a table, along with the IDs of the sequences from which they were extracted. (2) One of the possible  $\binom{g}{m}$  combinations of mismatch positions is removed to produce the gkmers. (3) The gkmers are sorted lexicographically. (4) Co-occurrences of the gkmers are counted and stored in the corresponding *partial kernel matrix*  $P_i$ . Steps (2-4) are repeated until each  $P_i$  is computed. (5) Once all partial matrices are computed, they are summed to produce the unnormalized kernel matrix  $K$ . (6) Kernel values are normalized.

### 3.2 Sort-and-Count Kernel Algorithm

We decompose equation 9 into a summation of multiple independent counting operations, where each operation handles a combination of mismatch positions. Our kernel function is given by:

$$K_{GSK}(x, y) = \sum_{i=1}^{gC_m} \sum_{\gamma \in \Theta_i} c_x(\gamma) c_y(\gamma) \quad (10)$$

Here  $gC_m = \binom{g}{m}$ ,  $\Theta_i$  denotes the set of gapped  $k$ -mers induced by the  $i$ th combination of  $m$  mismatch positions. Note that the sets  $\Theta_i$  compose the full feature set  $\Theta_{g,m}$ . We give a demonstration of our algorithm for computing equation 10 in figure 2 and briefly summarize the algorithm here. First, we extract all  $g$ -mers from the dataset and store how many times each  $g$ -mer occurs in each sequence. Then for each of the  $\binom{g}{m}$  combinations of mismatch positions, we remove  $m$  mismatch positions from the  $g$ -mers to obtain a set of gkmers. Next, we sort the gkmers and count when they are shared in common between pairs of sequences. When we find that a gkmer  $\gamma$  occurs in both sequences  $x$  and  $y$ , we store the product  $c_x(\gamma) c_y(\gamma)$  in a *partial kernel matrix*  $P_i$ . Denoting  $P_i$  as a function, the partial similarity score of  $x$  and  $y$  is given by:

$$P_i(x, y) = \sum_{\gamma \in \Theta_i} c_x(\gamma) c_y(\gamma) \quad (11)$$

Importantly, we compute each  $P_i$  independently. This way the algorithm is both easy to parallelize and easy to approximate using random sampling, as we show in section 3.3. Ultimately, once each  $P_i$  for  $i \in \{1, 2, \dots, \binom{g}{m}\}$  is

---

**Algorithm 1** FastGSK

---

**Require:**  $L, g, k$  ( $L$ =matrix of all  $g$ -mers from the dataset)

```
1: procedure CALCULATEKERNEL( $L, g, k$ )
2:    $M \leftarrow g - k$ 
3:    $N \leftarrow$ MISMATCHPROFILE( $L, g, M$ )
4:    $K \leftarrow 0$ 
5: procedure MISMATCHPROFILE( $L, g, M$ )
6:    $n_{pos} \leftarrow \binom{g}{m}$  ▷ Number of positions
7:   for  $i : 0 \rightarrow n_{pos}$  do
8:      $P_i \leftarrow 0$ 
9:      $L^i \leftarrow$ removePosition( $L, i$ )
10:     $L^i \leftarrow$ sort( $L^i$ )
11:     $P_i \leftarrow$ countAndUpdate( $L^i$ )
12:  for  $i : 0 \rightarrow n_{pos}$  do
13:     $K \leftarrow K + P_i$ 
return  $K$ 
```

**Ensure:**  $K$  ▷ Kernel Matrix

---

computed, the full kernel matrix is given by

$$K_{GSK} = \sum_{i=1}^{gCm} P_i \quad (12)$$

where  $gCm = \binom{g}{m}$ . We show these steps in algorithm 1. Finally, we normalize the kernel matrix using

$$K_{GSK}(x, y) \leftarrow \frac{K_{GSK}(x, y)}{\sqrt{K_{GSK}(x, x)K_{GSK}(y, y)}} \quad (13)$$

for pair of sequences  $(x, y)$ . We refer to this algorithm as FastGSK-Exact.

### 3.3 Approximation Algorithm

Because FastGSK-Exact runs with a coefficient of  $\binom{g}{m}$ , it is exponential in  $g$  and  $m$ . As such, it is unable to handle features roughly of size  $g > 15$ . However, many TF binding sites are up to 20 basepairs. If we let  $k = 6$ , we would have to compute more than 38,000 partial mismatch kernels. Moreover, even if  $g = 20$  is not optimal, a thorough grid-search must include large values of  $g$  in the search space in order to rule them out. Therefore, there is a strong need to create gapped  $k$ -mer algorithms that can scale to greater feature lengths. To solve this problem, we introduce a Monte Carlo approximation algorithm called FastGSK-Approx. FastGSK-Approx is extremely fast even for large values of  $g$ , as it requires only a small random subset of the  $\binom{g}{m}$  partial kernels  $P_i$ . Empirically, we show that FastGSK-Approx is roughly  $\mathcal{O}(1)$  with respect to  $g$ .

To compute FastGSK-Approx, we sample possible mismatch combinations for up to  $I_{max} \leq \binom{g}{m}$  iterations. That is, at iteration  $1 \leq t \leq I_{max}$ , we randomly sample (without replacement) a mismatch combination  $i \leftarrow \binom{g}{m}$  and compute the corresponding partial kernel matrix  $P_i$ . We then compute the online mean kernel matrix  $\bar{K}^{(t)}$  using  $P_i$  and  $\bar{K}^{(t-1)}$ . Furthermore, we compute a matrix of online standard deviations corresponding to the entries of  $\bar{K}^{(t)}$  and use the average of these values, which we denote as  $\sigma^{(t)}$ , to satisfy a convergence condition. Convergence is achieved when there is an approximately 95% probability that the online sample mean kernel  $\bar{K}^{(t)}$  is within  $\delta$  units of the true mean kernel matrix  $\mu_K$ . Here,  $\delta$  is a user-determined parameter. In practice, we use  $\delta = 0.025$ .

The idea rests on the Central Limit Theorem: we assume that for sufficiently large  $t$ , the sample mean kernel is normally distributed. Standardizing the variable  $\bar{K}^{(t)}$ , we have  $Pr[|\frac{\bar{K}^{(t)} - \mu_K}{\sigma^{(t)}}| > 1.96] \approx Pr[|z| > 1.96] = 0.05$ , where 1.96 is the z-score for a 95% confidence interval. Therefore, the convergence condition is satisfied when

$$1.96\sigma^{(t)} < \delta \quad (14)$$

We show that FastGSK-Approx converges rapidly, even for large values of  $g$  or  $m$ ; it typically converges when  $t \approx 50$ , which roughly corresponds to the number of samples needed to invoke the Central Limit Theorem. Furthermore, this means that FastGSK-Approx is roughly  $\mathcal{O}(1)$  with respect to  $g$ . Therefore, FastGSK-Approx solves the common string kernel problem of poor scalability with respect to the feature length.



### 3.4 Connecting to Related Work

**Mismatch Statistic-Based String Kernels** While FastGSK *directly* counts the gapped  $k$ -mers shared between sequences, previous works (e.g. [7, 8, 17, 18, 25]) *indirectly* compute the kernel function by inferring the counts from a set of *mismatch statistics*. These methods take inspiration from [16], which uses the notion of a *mismatch neighborhood* to efficiently compute the  $(k, m)$ -mismatch kernel. A mismatch neighborhood  $N_d(x, y)$  is simply the number of pairs of  $g$ -mers from sequences  $x$  and  $y$  that have a Hamming distance  $\leq d$ . Using this "statistic," the kernel function is inferred by multiplying  $N_d(x, y)$  with an appropriate coefficient for each value of  $d$ . Given a value of  $d$ , the required coefficient is some kernel-dependent combinatorial value. Though the mismatch statistic idea was created to improve the efficiency of the  $(k, m)$ -mismatch kernel, we argue that it actually harms efficiency in the case of the gapped  $k$ -mer kernel.

To compute the gapped string kernel, gkmSVM [7] applies the mismatch neighborhood idea to gapped  $k$ -mers. The key observation is that a pair of  $g$ -mers with a Hamming distance of  $d$  share  $\binom{g-d}{k}$  gkmers. Using this observation, the gkmSVM kernel is given by

$$K_{gkm}(x, y) = \sum_{d=0}^g N_d(x, y) h_{gk}(d) \quad (15)$$

where

$$h_{gk}(d) = \begin{cases} \binom{g-d}{k} & g-d \geq k \\ 0 & \text{otherwise} \end{cases} \quad (16)$$

The key point is that this kernel function does not count the features of interest, which are gapped  $k$ -mers, or gkmers. Rather, it counts  $g$ -mers at various Hamming distances  $d$  and uses a coefficient  $h_{gk}(d)$  to infer the gkmer counts.

**Theorem 1.** *The FastGSK kernel function is equivalent to the gkmSVM kernel function [7, 8].*

*Proof.* Because  $h_{gk}(d) = 0$  for  $d > m = g - k$ , equation 15 is equal to

$$K_{gkm}(x, y) = \sum_{d=0}^{m=g-k} N_d(x, y) h_{gk}(d) = \sum_{d=0}^{m=g-k} \binom{g-d}{k} N_d(x, y) \quad (17)$$

To *directly* count the gkmers, we replace  $\binom{g-d}{k} N_d(x, y)$ ,  $d \in 0, 1, 2, \dots, m$  with  $\binom{g}{m}$ , the number of possible mismatch positions. Therefore, at the  $i$ th mismatch position, we simply count the gapped  $k$ -mers in the set  $\Theta_i$ . Therefore, we have

$$K_{gkm}(x, y) = \sum_{d=0}^{m=g-k} \binom{g-d}{k} N_d(x, y) = \sum_{i=1}^{gCm} \sum_{\gamma \in \Theta_i} c_x(\gamma) c_y(\gamma) \quad (18)$$

as desired. Our method is both faster and simpler because we cut out the middleman: we directly count the gkmers. It also allows us to create a better approximation algorithm than what is possible using equation 15, because we can sample mismatch positions to estimate the kernel function.  $\square$

**Trie-based Implementations** Several implementations, including gkmSVM and gkmSVM-2.0 [7, 8], use a  $k$ -mer tree (or trie) to compute the mismatch neighborhood  $N_d(x, y)$  for each  $d \in \{0, 1, 2, \dots, g\}$  and pair of sequences  $(x, y)$ . These approaches have a branching factor equal to the alphabet size  $|\Sigma|$  (e.g., 4 for DNA) and depth of  $g$ . As such, they are impractical for tasks with larger alphabets, such as protein ( $\Sigma = 20$ ) or English text ( $\Sigma \geq 26$ ) classification. Furthermore, they scale poorly with the parameter  $g$ .

**Counting Implementations** GaKCo [25] is similar to FastGSK in that it uses a sort-and-count algorithm. However, it differs from FastGSK in that it follows the mismatch statistic formulation from [7]. The result is that GaKCo's time complexity has a  $\sum_{d=0}^{m=g-k} \binom{g}{d}$  coefficient. In contrast, we simply have a coefficient of  $\binom{g}{m}$ .

### 3.5 Implementation

**Lower Triangular Kernel** Because a kernel matrix is symmetric by definition, we only store the lower triangle to save memory. In practice, the lower triangular kernel matrix is actually treated as a one-dimensional contiguous array in memory. This simplifies the code and improves cache utilization through better access locality.

**Multithreading** As shown in equation 12, we decompose the kernel function into a summation of  $\binom{g}{m}$  partial kernel matrices  $P_i$ , with each  $P_i$  holding the gkmer counts for a single combination of mismatch positions. Observing that each  $P_i$  is a completely independent subproblem, we exploit the decomposition to create a multithreaded implementation.

To map threads to subproblems, we create a "work queue"—a length  $\binom{g}{m}$  array, where each element indicates a mismatch combination subproblem to solve. If there are  $t$  threads (a user-specified argument), we divide the work queue into  $t$  equally sized portions. In practice, each thread simply moves forward  $t$  positions in the work queue when it completes a subproblem. An illustration is shown in figure 3. A thread's local results are aggregated into its own temporary matrix; since all  $P_i$  are added together in the end, each thread simply aggregates its partial results as it finishes them. Then once each thread finishes, it adds its partial results into the full kernel matrix  $K$ .

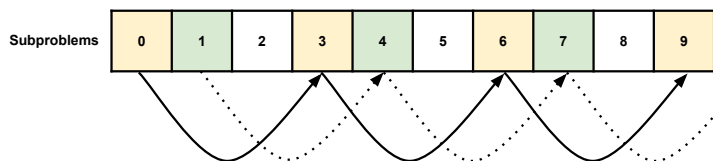


Figure 3: An illustration of how subproblems (partial kernel matrices  $P_i$ ) are divided between threads. If there are 9 mismatch position combinations and  $t = 3$  threads, then thread 0 (solid arrow and yellow cells) handles the combinations/subproblems 0, 3, 6, and 9. After finishing each subproblem, it aggregates the results into a local sub-sum matrix and skips forward  $t$  positions. The dashed arrow (green cells) shows the procedure for thread 1. Once each thread has run out of subproblems to handle, it proceeds to aggregate its results into the full kernel matrix  $K$ .

**Synchronization** Once each thread finishes computing its subset of the  $\binom{g}{m}$  partial kernel matrices, these are aggregated into the kernel matrix  $K$  as per equation 12. But because each thread must access the same  $K$  in memory, we must obviate potential race conditions. To ensure synchronization and maximize the number of threads able to access  $K$  simultaneously, we use a set of mutex locks. Intuitively, the idea is to create a set of evenly sized memory regions such that only one thread is allowed in each region at a time; a thread enters a region, locks it, and performs its aggregations. When it finishes it unlocks the previous regions, permitting a new thread to enter, and then locks the next region. An illustration using a small kernel matrix is shown in figure 4.

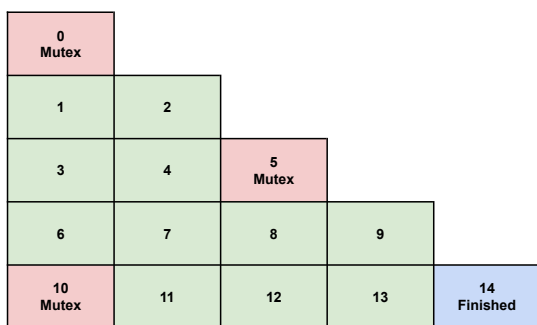


Figure 4: An example of our matrix aggregation and synchronization scheme using the lower triangle of a kernel matrix. The indexes show the order in which cells are accessed by the threads. Each time a thread reaches one of the red cells (labeled "Mutex"), it locks the indices up until the next mutex location. When it reaches the next mutex location, it unlocks the previous memory region to permit a new thread to enter. It also locks the next memory region. Though shown here as a triangular matrix, the kernel matrix is actually a single contiguous array in practice. This maximizes cache performance.

**Software Details** Our algorithms are easily available as a PyPi package<sup>3</sup> for Python. The package, called "fastsk", consists of a well-optimized C++ implementation with a Python interface. We bind the C++ backend and Python interface using the Pybind11 library<sup>4</sup> [13]. Therefore, we combine the simplicity and convenience of Python with the speed and low-level control of C++. To our knowledge, this is the only string kernel package that does so. We believe FastGSK is the best string kernel toolbox for use within modern data science/machine learning toolboxes.

**Usage** Unlike virtually all extant string kernel toolboxes, FastGSK is designed to be easy to use in modern data science/machine learning workflows. For one, it is installed using a simple command (`pip install fastsk`). Two, kernel

<sup>3</sup>Available at <https://github.com/QData/FastSK> or with the command "pip install fastsk"

<sup>4</sup><https://github.com/pybind/pybind11>

---

```

from fastsk import Kernel
from sklearn.svm import LinearSVC

# Train and test samples can be Python lists or numpy arrays
x_train, y_train, x_test, y_test = read_samples("train.fasta", "test.fasta")

# Set the parameters for computing the gapped string kernel
g, m, threads = 8, 4, 20
kernel = Kernel(g=g, m=m, t=threads, approx=False)

# Compute the train and test kernels
kernel.compute(x_train, x_test)
train_kernel, test_kernel = kernel.train_kernel(), kernel.test_kernel()

# Train and evaluate linear SVM
svm = LinearSVC(C=1)
svm.fit(train_kernel, y_train)
accuracy = svm.score(test_kernel, y_test)

```

---

Figure 5: FastGSK is easy to use in Python in conjunction with the Scikit-Learn library. Here, we show an example of computing a gapped string kernel and then using the train and test kernels to train and evaluate a linear SVM. Training many other models (kernel SVM, logistic regression, etc.) is straightforward as well.

matrices can be computed in just a few lines of Python. Three, it is easy to use FastGSK in conjunction with the standard classifiers in the Scikit-Learn library, including linear SVM (LinearSVC), kernel SVM (SVC), stochastic gradient descent (SGDClassifier), and logistic regression (LogisticRegression). Finally, evaluation and analysis are also quite simple, as metrics such as AUC and F1 score are trivial to compute using the Scikit-Learn metrics library. We illustrate a typical use-case in figure 5.

## 4 Experiments

### 4.1 Experimental Setup

**Datasets** We evaluate FastGSK using 10 transcription factor binding site (TFBS) DNA datasets from the ENCODE project, 10 protein remote homology datasets from the SCOP project, and 7 English-language medical named entity recognition datasets. Summary statistics and speedup results for our DNA, protein, and NLP datasets are shown in tables 1, 2, and 3, respectively.

**Baselines** We compare the prediction accuracy and efficiency of FastGSK with 3 state-of-the-art string kernel baselines. We use gkmSVM-2.0 [8] for DNA experiments because it uses gkmer features and was designed for use on DNA tasks [7, 8]. As a string kernel baseline for protein datasets, we use GaKCo [25], as it was designed to extend the gkmSVM-2.0 kernel function to datasets with larger alphabets. As an NLP string kernel baseline, we use the Blended Spectrum Kernel [12, 11], as has it recently achieved strong results. For FastGSK, gkmSVM-2.0, and GaKCo, we perform a grid search over the hyper-parameters  $g \in \{5, 6, \dots, 15\}$ ,  $m \in \{0, 1, \dots, g - 1\}$ , and the SVM margin parameter  $C \in \{0.001, 0.01, 0.1, 1, 10, 100\}$ . We use 5-fold cross-validation on each training set. For the Blended Spectrum Kernel, we use the authors' string kernel package<sup>5</sup> and use the parameters  $k_{min} = 3$  and  $k_{max} = 5$ , as recommended by the authors.

We also compare FastGSK with bidirectional Long Short-Term Memory (LSTM) memory networks. To train a model on each dataset, we perform a grid-search over the embedding and hidden dimensionality  $\in \{32, 64, 128, 256\}$ , and number of layers  $\in \{1, 2, 3, 4\}$ . We use the Adam optimizer [14] with a learning rate of 0.001. We perform 5-fold cross-validation on each training set. We use the probability scores of the final layer to compute AUC.

---

<sup>5</sup>Available for download at: <http://string-kernels.herokuapp.com/>

**SVM Training** To train SVM models using each string kernel method, we use Liblinear [5] via the Scikit-Learn LinearSVC implementation<sup>6</sup> with L2 regularization. We also used a non-linear SVM, namely LIBSVM<sup>7</sup> [3], but found it was both impractically slow and showed inferior performance. In order to use a kernel method with a linear SVM, we use the empirical kernel map (EKM) strategy. That is,  $K \leftarrow K_{GSK}^\top K_{GSK}$ . To evaluate the SVM models, we use area under the ROC curve (AUC) with Platt Scaling [22] to obtain probabilities. Unless otherwise noted, we run FastGSK-Approx with  $I_{max} = 50$  and  $\delta = 0.025$ . All timing experiments are performed on a server with 12 Intel i7-6850K 3.60GHz CPUs.

## 4.2 Dataset Statistics

Table 1: DNA Datasets

| Dataset     | Train | Test | Total | Avg Speedup $\times$ | Max Speedup $\times$ |
|-------------|-------|------|-------|----------------------|----------------------|
| CTCF        | 2000  | 2000 | 4000  | 151                  | 723                  |
| EP300       | 2000  | 2000 | 4000  | 16                   | 92                   |
| JUND        | 2000  | 2000 | 4000  | 151                  | 757                  |
| RAD21       | 2000  | 2000 | 4000  | 161                  | 808                  |
| SIN3A       | 2000  | 2000 | 4000  | 16                   | 88                   |
| Pbde        | 4500  | 5500 | 10000 | 18                   | 107                  |
| EP300_47848 | 6506  | 724  | 7230  | 162                  | 833                  |
| KAT2B       | 6318  | 702  | 7020  | 161                  | 809                  |
| TP53        | 4432  | 494  | 4926  | 23                   | 81                   |
| ZZZ3        | 9966  | 1108 | 11074 | 18                   | 111                  |

We evaluate FastGSK using 10 DNA transcription factor binding site datasets from the ENCODE project. We measure the factor of speedup in kernel computation time relative to gkmSVM-2.0, showing that FastGSK-Approx is typically hundreds of times faster than gkmSVM-2.0. For each dataset, we vary  $g$  from 6 to 20, while fixing  $k = 6$ . We run FastGSK-Approx with  $I_{max} = 50$  and  $\delta = 0.025$ , the parameters used to obtain the results shown in figure 15. The speedups are computed from the averages of the kernel computation as  $g$  increases. An example of the kernel computation times as we increase  $g$  is shown in figure 9. Timing figures for all 10 DNA datasets are shown in figure 12. Note: these factors of speedup are actually biased against FastGSK-Approx, because as we vary  $g$ , we halt gkmSVM-2.0 once its kernel computation time exceed 3600s (30 minutes). This is also why gkmSVM-2.0’s curve stops early in all of the plots shown in figure 12.

Table 2: Protein Datasets

| Name | Train | Test | Total | Speedup $\times$ |
|------|-------|------|-------|------------------|
| 1.1  | 2339  | 1235 | 3574  | 3.4              |
| 1.34 | 2075  | 1237 | 3312  | 3.1              |
| 2.19 | 1345  | 1215 | 2560  | 3.4              |
| 2.31 | 2298  | 1202 | 3500  | 426.6            |
| 2.34 | 1501  | 1237 | 2738  | 0.9              |
| 2.41 | 1427  | 1219 | 2646  | 10.4             |
| 2.8  | 1241  | 1239 | 2480  | 24.5             |
| 3.19 | 2103  | 1238 | 3341  | 6.7              |
| 3.25 | 2395  | 1242 | 3637  | 464.4            |
| 3.33 | 1680  | 1238 | 2918  | 3.6              |

Table 3: NLP Datasets

| Name        | Train | Test | Total | Speedup $\times$ |
|-------------|-------|------|-------|------------------|
| AIMed       | 1500  | 1500 | 3000  | 6.5              |
| BioInfer    | 2534  | 2534 | 5068  | 43.6             |
| CC1-LLL     | 3785  | 330  | 4115  | 13.1             |
| CC2-IEPA    | 3298  | 817  | 4115  | 8.6              |
| CCC3-HPRD50 | 3682  | 433  | 4115  | 4.5              |
| DrugBank    | 2472  | 2472 | 4944  | 7.7              |
| MedLine     | 635   | 635  | 1270  | 20.7             |

We evaluate FastGSK using 10 SCOP project protein remote homology detection datasets and 7 medical named entity recognition datasets. We show the kernel computation time factor of speedup achieved by FastGSK over GaKCo and Blended Spectrum on each dataset. The average factors of speedup are  $94.7\times$  and  $15.0\times$ , respectively.

<sup>6</sup><https://scikit-learn.org/stable/modules/generated/sklearn.svm.LinearSVC.html>

<sup>7</sup>Available at <https://www.csie.ntu.edu.tw/~cjlin/libsvm/> or with the Scikit-Learn Python interface at <https://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>

### 4.3 Prediction Performance

In this section we show the FastGSK-Approx algorithm rivals or outperforms previous string kernel algorithms—as well as long short-term memory (LSTM) networks—in test AUC for DNA, protein, and NLP tasks.

**FastGSK matches or outperforms gkmSVM-2.0 and LSTMs at TFBS prediction** For DNA, we baseline against gkmSVM-2.0 [8], a popular gapped  $k$ -mer string kernel toolbox for DNA sequence analysis. We show that FastGSK-Approx achieves nearly identical results as gkmSVM-2.0 across TFBS tasks. This is despite the fact that we run FastGSK-Approx with the maximum number of mismatch positions  $I_{max} = 50$ , which is often a tiny fraction of the exact computation performed by gkmSVM-2.0 or FastGSK-Exact. As shown in figure 6, FastGSK-Approx nearly matches gkmSVM-2.0 and outperforms LSTMs across all 10 DNA datasets.

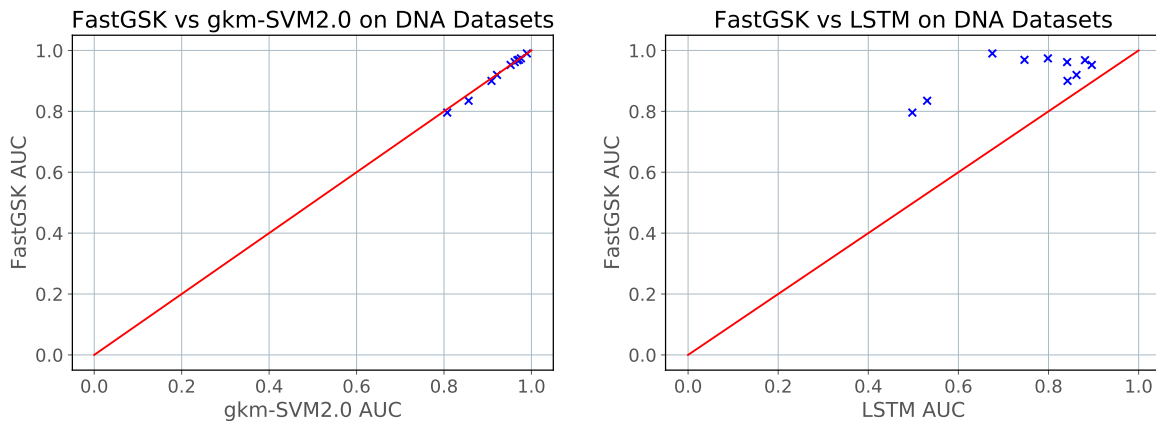


Figure 6: Left: across all 10 DNA datasets, FastGSK matches gkmSVM-2.0 in AUC. Right: FastGSK substantially outperforms our deep learning baseline.

**FastGSK outperforms GaKCo and LSTMs at protein function identification** On protein datasets, we baseline against GaKCo [25]. GaKCo uses the same kernel function as gkmSVM-2.0 (equation 15), but was designed to extend it to larger alphabets, such as amino acids for protein data. We show that FastGSK-Approx outperforms GaKCo at protein remote homology detection. As shown in figure 7, FastGSK-Approx typically outperforms both GaKCo and LSTMs at protein homology identification.

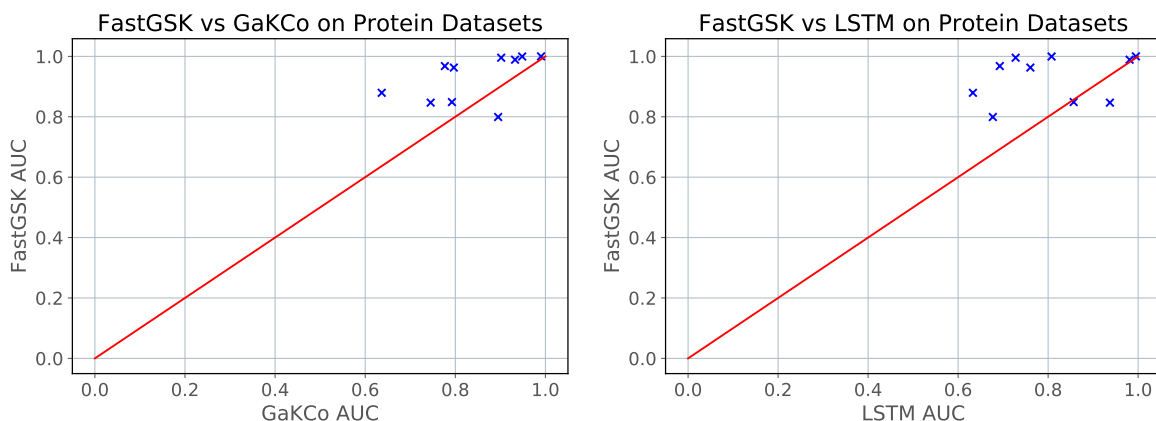


Figure 7: Left: FastGSK obtains AUC scores competitive with GaKCo on all 10 protein datasets. Right: FastGSK outperforms our deep learning baseline.

**FastGSK outperforms Blended Spectrum Kernel and LSTMs at medical NLP** Finally, for medical named entity recognition, we baseline against the Blended Spectrum string kernel [12, 11], as it has recently shown strong performance

in NLP. We show that FastGSK-Approx outperforms the Blended Spectrum kernel. As shown in figure 8, FastGSK-Approx consistently outperforms the Blended Spectrum Kernel and LSTM baselines in test AUC across 7 medical named entity recognition tasks.

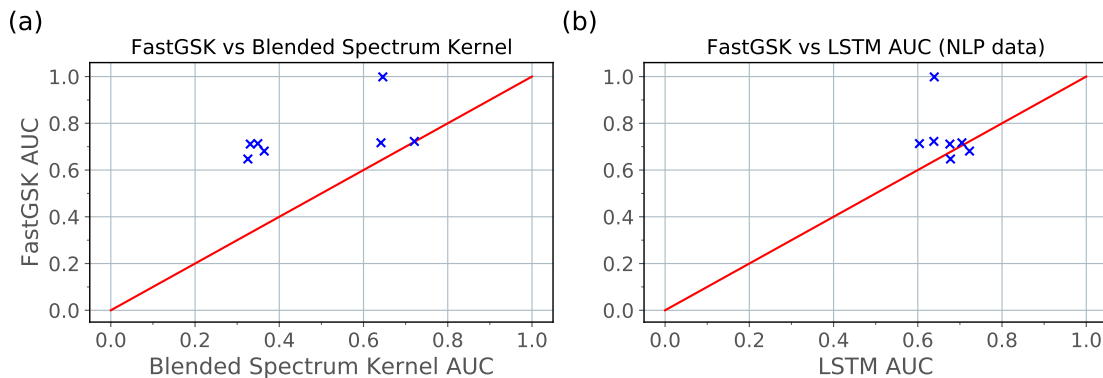


Figure 8: We compare the AUC obtained with the FastGSK kernel across 7 medical named entity recognition tasks. (a) FastGSK outperforms the Blended Spectrum Kernel baseline [12, 11] by 0.259 AUC on average. To compute the kernels, we use the best  $g$  and  $m$  for FastGSK and  $k_{min} = 3$ ,  $k_{max} = 5$  for Blended Spectrum Kernel, as recommended by the authors. (b) FastGSK outperforms the LSTM baseline by 0.075 AUC on average.

Taken together, these results demonstrate FastGSK (1) is performant across a range of sequence analysis tasks and (2) consistently rivals or outperforms the state-of-the-art string kernels and LSTM models.

#### 4.4 Scalability

In this section we showcase a key advantage of FastGSK-Approx: its ability to scale to greater feature lengths  $g$ . In contrast, the state-of-the-art gapped  $k$ -mer algorithms, such as gkmSVM-2.0 [8] and GaKCo [25], both fail to scale to greater feature lengths. This is because their time complexities depend exponentially on the feature length  $g$  and number of mismatch positions  $m$ . In contrast, FastGSK-Approx is approximately  $\mathcal{O}(1)$  with respect to these parameters. As such, it is efficient irrespective of the feature length  $g$ . Moreover, we demonstrate that FastGSK-Approx often shows high test AUCs for feature lengths of up to  $g = 20$ .

**FastGSK efficiently handles greater feature lengths** As we show in figure 9, FastGSK scales to much greater values of  $g$  than gkmSVM-2.0. There are two reasons for this. First, FastGSK uses *direct feature counting*, which is more efficient than the mismatch statistics formulation from equation 15. Ironically, the mismatch statistic formulation was intended to be *faster* than direct feature counting, yet this is indeed not the case. As shown in figure 12, FastGSK-Exact’s *direct* counting method is faster than gkmSVM-2.0 across all values of  $g$ .

Second, FastGSK-Approx converges quickly, using just a small fraction of the possible mismatch combinations. Indeed, the FastGSK-Approx algorithm effectively runs in time  $\mathcal{O}(1)$  with respect to the feature length  $g$ , as shown in figure 9. This is because the algorithm converges after using just a small number of the possible mismatch positions—irrespective of  $g$ .

**FastGSK achieves strong prediction performance for large values of  $g$**  As we show in figures 9, 10, and 11, strong test AUC values are often achieved when the feature length  $g$  is large. Although we observe that test AUC is rarely *optimal* for very large feature lengths, this does not diminish the importance of creating a scalable algorithm. There are two reasons for this. First, we cannot rule out the existence of datasets or tasks where  $g \geq 20$  is optimal. Second, a thorough grid search to find the optimal  $g$  for a given task must check the performance when  $g$  is large. For this reason, it hardly matters that the optimal feature length for TFBS tasks is usually  $g \approx 14$ ; it is still extremely valuable for the algorithm to scale to  $g \geq 20$ .

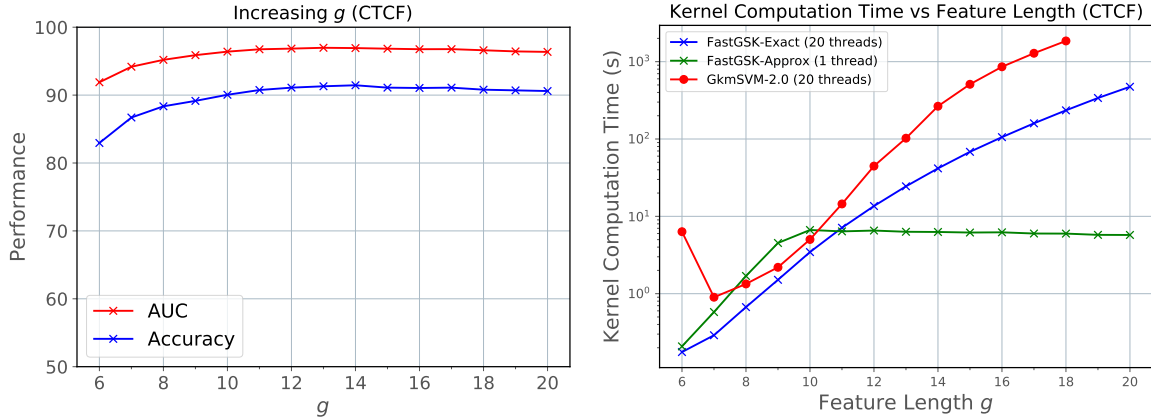


Figure 9: On the CTCF-bound transcription factor dataset, we vary  $g$  from 6 to 20, keeping the number of informative positions fixed at  $k = g - m = 6$ . Left: we show the AUC obtained by an SVM using FastGSK-Approx with  $I = \binom{g}{m}$  (i.e. no upper bound) and  $\delta = 0.025$ . Right: we time the kernel matrix computation for gkmSVM-2.0, FastGSK-Exact, and FastGSK-Approx. We stop each algorithm early once it reaches 1800s (30 minutes). We show that as  $g$  increases, FastGSK-Exact is faster than gkmSVM-2.0 by an order of magnitude, while FastGSK-Approx is faster by multiple orders of magnitude. Because FastGSK-Approx’s time plateaus at  $g = 10$  with  $k = 6$ , we show that the algorithm convergence after using around  $\binom{10}{4} = 210$  mismatch positions. As such, FastGSK-Approx is approximately  $\mathcal{O}(1)$  in terms of  $g$ .

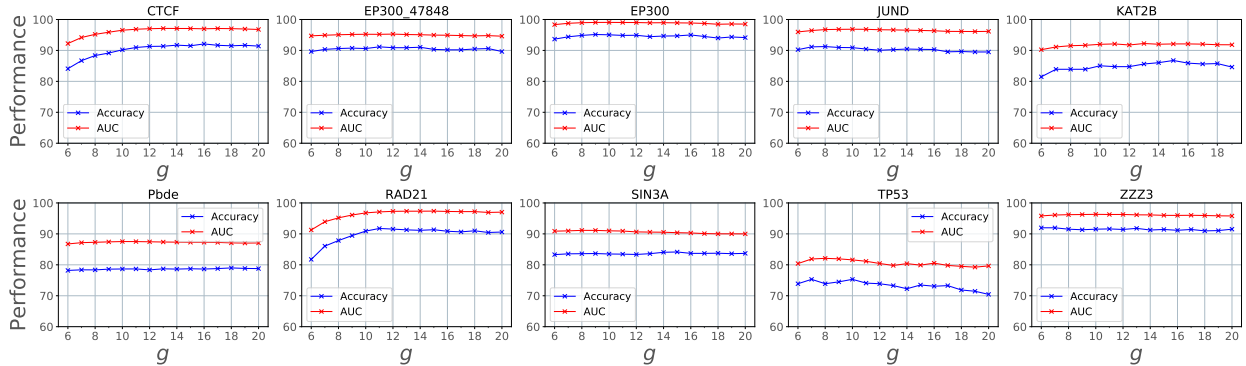


Figure 10: We vary the feature length  $g$  while keeping the number of informative feature positions fixed at  $k = 6$  (found empirically to work well) for the 10 TFBS datasets. For each dataset, we compute the kernel function using FastGSK-Approx with  $I_{max} = 50$ . AUC performance is evaluated using Liblinear [5] via the Scikit-Learn LinearSVC class with an empirical kernel map.

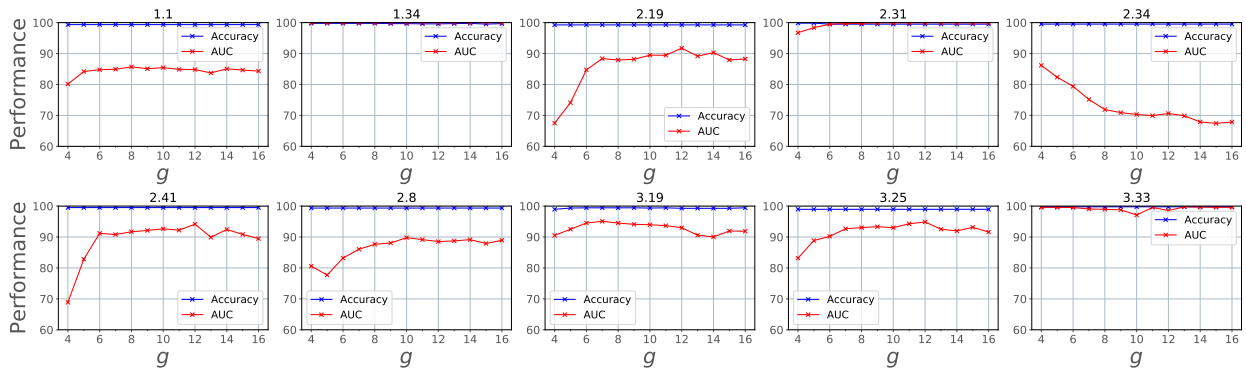


Figure 11: We vary the feature length  $g$  while keeping the number of informative feature positions fixed at  $k = 4$  (found empirically to work well) for the 10 protein SCOP datasets. We ran FastGSK-Approx with  $I_{max} = 50$ . AUC performance is evaluated using Liblinear [5] via the Scikit-Learn LinearSVC class with an empirical kernel map.

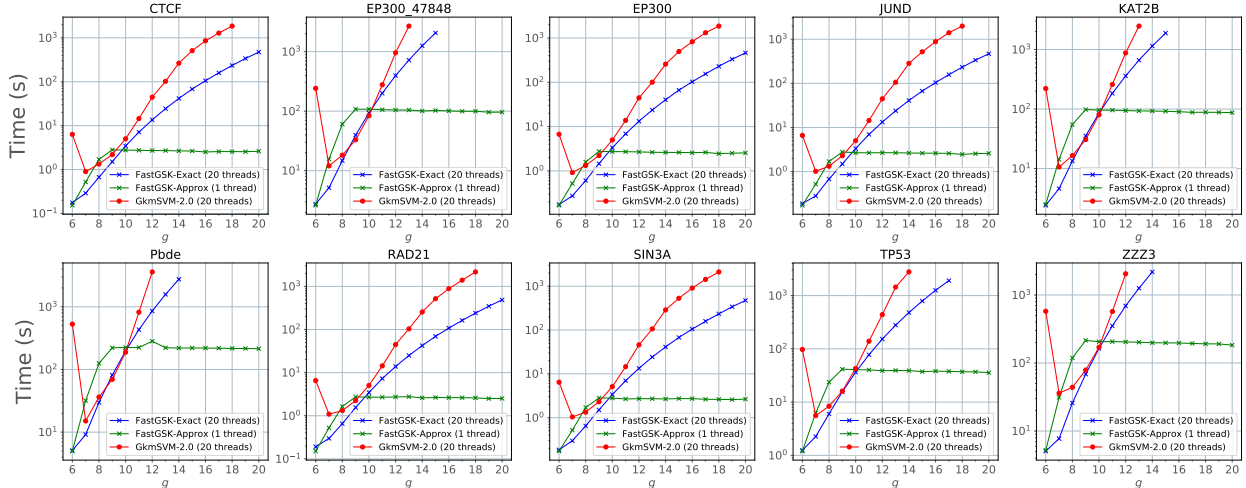


Figure 12: We vary the feature length  $g$  while keeping the number of informative feature positions fixed at  $k = 6$  for each of the 10 TFBS datasets. We run both FastGSK-Exact and gkmSVM-2.0 using 20 threads and stop each algorithm early once the kernel computation time exceeds 1800s (30 minutes). We run FastGSK-Approx using just 1 thread and the number of iterations set to  $I_{max} = 50$ .

#### 4.5 Multithreading

In this section we examine how kernel computation time is a function of the number of threads for FastGSK-Exact, FastGSK-Approx, and gkmSVM-2.0. To start, we point out that FastGSK-Approx generally does not benefit from using multiple threads. There are two reasons. First, our primary implementation of the multithreaded FastGSK-Approx algorithm improves approximation quality, not speed. This is because each thread simply iterates independently until it converges. After all threads converge, their mean kernels are aggregated and normalized. Nothing about this improves speed per se. Second, we did create an option to *accelerate* FastGSK-Approx, but the speedups were marginal. This option is identical to FastGSK-Exact, except with the  $I_{max}$  parameter set to a user-specified value. That is, it simply skips the variance and mean kernel matrix computations, using  $t$  threads to compute  $I_{max}$  partial kernels  $P_i$  in total. In other words, each thread computes  $\approx I_{max}/t$  partial kernels. However, we have shown  $I_{max} \approx 50$  is usually sufficient. This is such a small number of iterations that it does not stand to benefit much from multithreading. In fact, the synchronization overhead in use-cases like this typically negate the benefits of multithreading.

That being said, FastGSK-Approx with just *one* thread is still vastly faster than either FastGSK-Exact or gkmSVM-2.0 with *twenty* threads. Unlike FastGSK-Approx, these algorithms benefit greatly from more threads, as we show. But nevertheless, FastGSK-Approx is much faster. We claim this is one of the key merits of this work: fast kernel computation without multithreading at all. Yet even in the case of exact kernel computation, we claim another key merit: at each number of threads FastGSK-Exact is significantly faster than gkmSVM-2.0, which we show in figures 13 and 14.

**FastGSK outperforms gkmSVM-2.0 for any number of threads** As shown in figure 13, FastGSK-Exact is substantially faster than gkm-SVM2.0 for each number of threads. FastGSK-Approx is vastly faster still, even when using just 1 thread. We show similar behavior across all 10 DNA datasets in the appendix in figure 14.

**FastGSK is better for multithreading than GaKCo or Blended Spectrum** Neither GaKCo nor the Blended Spectrum kernel make adequate use of multithreading. GaKCo limits the number of possible threads to  $m$ , while Blended Spectrum only allows one thread.



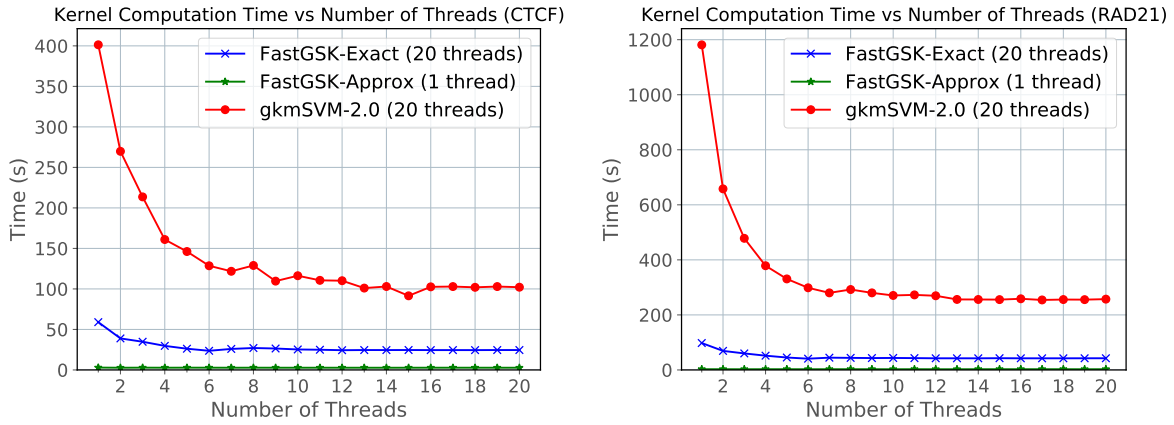


Figure 13: While FastGSK-Exact is faster than gkmSVM-2.0 for each number of threads, we further demonstrate that FastGSK-Approx using just 1 thread is much faster still.

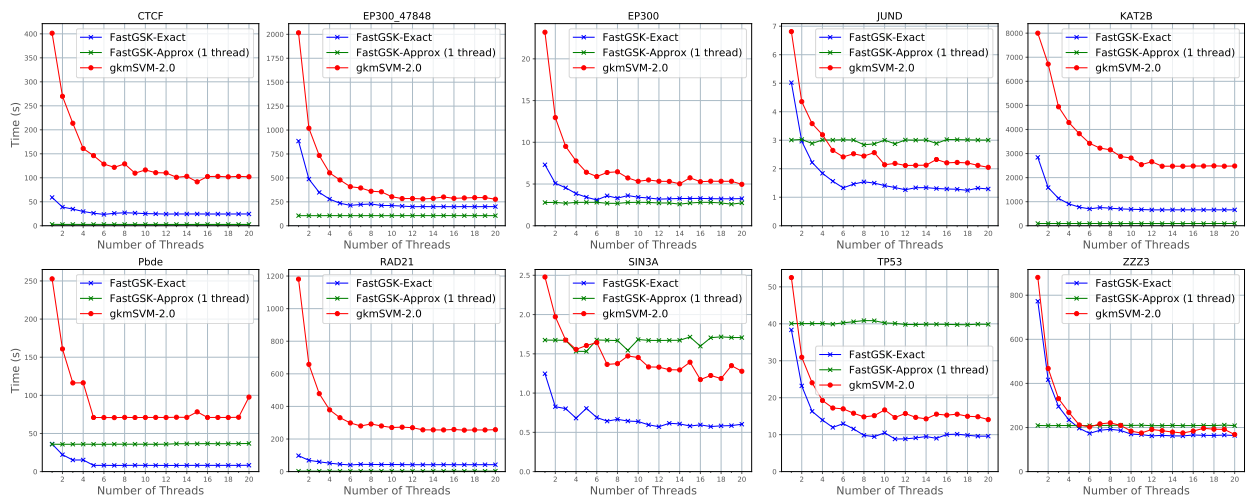


Figure 14: We vary the number of threads used for kernel computation across all 10 DNA datasets. Each kernel is computed using the optimal parameters.

## 4.6 Approximation Quality

In this section we show that FastGSK-Approx is nearly identical to an exact gapped  $k$ -mer string kernel SVM. We demonstrate this by comparing FastGSK-Approx with FastGSK-Exact, showing that FastGSK-Approx yields test AUCs just 0.003 points lower than FastGSK-Exact. Moreover, we validate our explanation as to why FastGSK-Approx works so well: only a very small number of mismatch combinations (iterations of FastGSK-Approx) is needed to achieve excellent test performance. Taken together, these findings prove that FastGSK-Approx is an excellent approximation algorithm and explain why this is so.

**FastGSK-Approx achieves excellent approximations** In figure 15 (right) we demonstrate that FastGSK-Approx yields test AUCs nearly identical to those obtained via FastGSK-Exact. We compare the exact and approximation algorithms across all 27 datasets (DNA, protein, and NLP), showing the AUC penalty incurred is just 0.003 points on average. This illustrates a key merit of this work: FastGSK is orders of magnitude faster than gkmSVM-2.0, GaKCo, and the Blended Spectrum kernel and this speedup comes at an extremely minimal cost.

**Mismatch combinations have diminishing returns** Figures 15 and 16 both show this: the computational cost from handling all  $\binom{g}{m}$  mismatch combinations is unnecessary. Figure 15 (left) shows that the test AUC quickly increases with the number of mismatch combinations (iterations) used, while the standard deviation of the sample mean kernel matrix quickly decays to 0. This simultaneously illustrates the point that few mismatch combinations are needed and explains why FastGSK-Approx converges rapidly. Moreover, we claim that these points hold true across many datasets—which we justify with figure 16.

We point out two interesting findings. First, the AUCs and standard deviations are roughly in accordance with the Central Limit Theorem, as we observe several dozen samples are typically sufficient. Second, figure 16 points to a direction for future research. *Extremely* few mismatch combinations are sufficient in some cases—sometimes just one mismatch combination is good enough. This is a very surprising finding that warrants further investigation. We suspect it could point to even better gapped  $k$ -mer algorithms and refine our understanding of SK-SVM models.

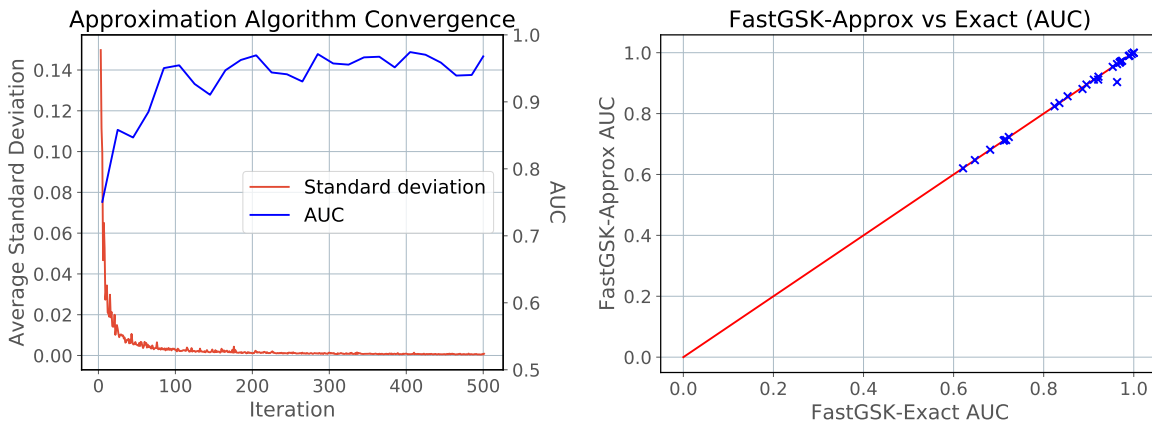


Figure 15: Left: We use the 2.31 protein dataset with  $g = 15$  and  $m = 10$  to show that test performance quickly converges with the number of mismatch positions sampled. Right: FastGSK-Approx with  $I_{max} = 50$  and  $\delta = 0.025$  achieves almost identical test AUC as FastGSK-Exact. This is because very few partial kernels  $P_i$  are needed to achieve a good approximation of the true kernel. Furthermore, we note that  $I_{max} = 50$  is roughly the number of samples needed to invoke the Central Limit Theorem. We found setting  $\delta = 0.025$  works well empirically. On average, the test AUC of FastGSK-Approx is 0.003 points lower than that of FastGSK-Exact.

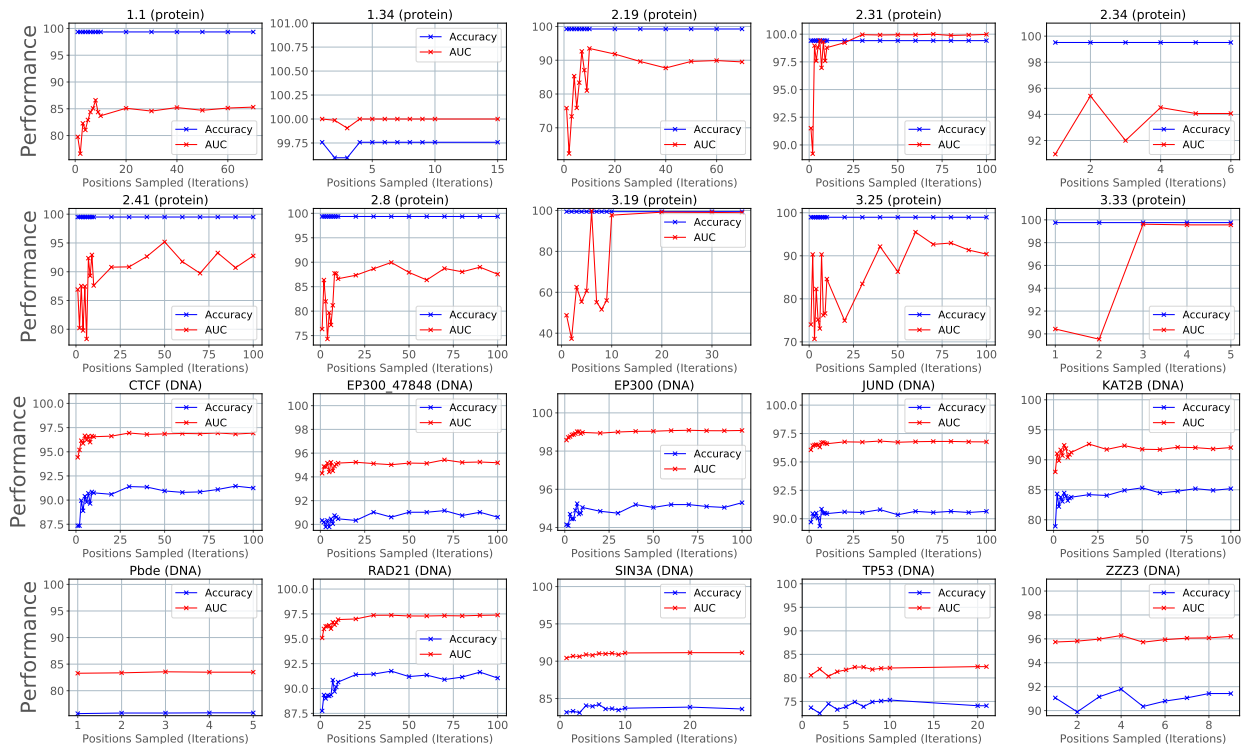


Figure 16: We vary the maximum number of mismatch positions sampled (determined by the  $I_{max}$  parameter) used for kernel computation. We then train and test an SVM and show the resultant test accuracy and AUC. These results show that extremely few mismatch positions for kernel computation usually results in poor test performance, but that test performance increases with the mismatch positions sampled up until several dozen. This finding provides insight into why FastGSK-Approx works so well. However, these results also showcase an unexpected property of gapped  $k$ -mer SVMs: in some cases extremely few mismatch combinations are actually necessary for excellent performance. We propose exploring this phenomenon as a direction for future research.

## 5 Conclusion and Future Work

In this work we introduced a fast and scalable string kernel SVM algorithm called FastGSK. FastGSK rivals state-of-the-art string kernel SVMs in test performance, while running 2-3 orders of magnitude faster on average. It also outperforms deep learning. FastGSK has four high-level contributions:

1. Like [7], it uses a compact feature space using gapped  $k$ -mer features that are less prone to overfitting than traditional  $k$ -mer features.
2. A fast and efficient sort-and-count algorithm that *directly* computes the kernel function. Not only is this approach simpler than previous works, but it is also naturally parallelizable, a property we exploit to create a fast multithreaded implementation. It is also independent of the alphabet size  $\Sigma$ .
3. A fast kernel approximation algorithm that scales to greater feature lengths than the state-of-the-art algorithms. We experimentally showed that FastGSK matches or outperforms existing string kernel SVMs in AUC on DNA, protein, and NLP data while running  $\sim 100\times$  faster on average and up to  $800\times$  faster for large feature lengths.
4. FastGSK outperforms LSTMs on datasets with under 10,000 samples. We showed this result held across DNA, protein, or NLP tasks.

### 5.1 Future Directions

We identify several promising areas where future work could be directed.

1. Finding even simpler gapped  $k$ -mer formulations. The success of FastGSK-Approx and the curiously small number of mismatch combinations needed in many cases suggest simpler exact algorithms could exist. The

goal here would be to shed the  $\binom{g}{m}$  coefficient from the exact algorithm. An intuition is that it seems only a few mismatch positions are necessary to obtain most or all of the *unique* gapped  $k$ -mers and that the number of unique gapped  $k$ -mers is what ultimately matters.

2. Reducing memory usage. We currently use  $\mathcal{O}(n^2)$  memory to store the kernel matrix. This is intractable for  $n$  greater than several tens of thousands. Future work should implement techniques to avoid storing the full kernel matrix in memory at once. For example, [17] only computes small batches of the kernel matrix at a time, greatly reducing the memory footprint.
3. Capitalizing on interpretability. String kernel methods have interpretable features and many works have identified the important features from string kernel methods [24, 18]. Future work can use methods such as Data Shapley [9] to analyze the most salient features. The fact that so few features appear to be necessary in our results suggests a small number of features are critical, while most are disposable. Future interpretability work could study this hypothesis.
4. Low rank approximations of the gram matrix  $K$ . The efficiency of kernelized SVMs can be improved via the Nyström method [29], which creates a low-rank approximation of the kernel matrix. This method still requires  $\mathcal{O}(n^2)$  space for the kernel matrix  $K$ , but it ultimately provides an approximation  $\tilde{K}$  to be used by the SVM optimizer. For inspiration, [30] shows a linear SVM method inspired by Nyström’s method.

## References

- [1] Michael A Beer and Saeed Tavazoie. Predicting gene expression from sequence. *Cell*, 117(2):185–198, 2004.
- [2] Zhen Cao and Shihua Zhang. Probe efficient feature representation of gapped k-mer frequency vectors from sequences using deep neural networks. *IEEE/ACM transactions on computational biology and bioinformatics*, 2018.
- [3] Chih-Chung Chang and Chih-Jen Lin. Libsvm : a library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2(27), 4 2011.
- [4] Eleazar Eskin, Jason Weston, William S Noble, and Christina S Leslie. Mismatch string kernels for svm protein classification. In *Advances in neural information processing systems*, pages 1441–1448, 2003.
- [5] Rong-En Fan, Kai-Wei Chang, Cho-Jui Hsieh, Xiang-Rui Wang, and Chih-Jen Lin. Liblinear: A library for large linear classification. *Journal of machine learning research*, 9(Aug):1871–1874, 2008.
- [6] Muhammad Farhan, Juvaria Tariq, Arif Zaman, Mudassir Shabbir, and Imdad Ullah Khan. Efficient approximation algorithms for strings kernel based sequence classification. In *Advances in Neural Information Processing Systems*, pages 6935–6945, 2017.
- [7] 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.
- [8] Mahmoud Ghandi, Morteza Mohammad-Noori, Narges Ghareghani, Dongwon Lee, Levi Garraway, and Michael A Beer. gkmsvm: an r package for gapped-kmer svm. *Bioinformatics*, 32(14):2205–2207, 2016.
- [9] Amirata Ghorbani and James Zou. Data shapley: Equitable valuation of data for machine learning. *arXiv preprint arXiv:1904.02868*, 2019.
- [10] Fumitaka Inoue, Martin Kircher, Beth Martin, Gregory M Cooper, Daniela M Witten, Michael T McManus, Nadav Ahituv, and Jay Shendure. A systematic comparison reveals substantial differences in chromosomal versus episomal encoding of enhancer activity. *Genome research*, 27(1):38–52, 2017.
- [11] Radu Tudor Ionescu and Andrei Butnaru. Improving the results of string kernels in sentiment analysis and Arabic dialect identification by adapting them to your test set. In *Proceedings of EMNLP*, 2018.
- [12] Radu Tudor Ionescu, Marius Popescu, and Aoife Cahill. Can characters reveal your native language? A language-independent approach to native language identification. In *Proceedings of EMNLP*, pages 1363–1373. Association for Computational Linguistics, October 2014.
- [13] Wenzel Jakob, Jason Rhinelander, and Dean Moldovan. pybind11—seamless operability between c++ 11 and python, 2017.
- [14] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [15] Pavel P. Kuksa, Pai-Hsi Huang, and Vladimir Pavlovic. Efficient use of unlabeled data for protein sequence classification: a comparative study. *BMC Bioinformatics*, 10(S-4), 2009.
- [16] Pavel P Kuksa, Pai-Hsi Huang, and Vladimir Pavlovic. Scalable algorithms for string kernels with inexact matching. In *Advances in neural information processing systems*, pages 881–888, 2009.
- [17] D. Lee. Ls-gkm: a new gkm-svm for large-scale datasets. *Bioinformatics*, 2016.
- [18] Dongwon Lee, David U Gorkin, Maggie Baker, Benjamin J Strober, Alessandro L Asoni, Andrew S McCallion, and Michael A Beer. A method to predict the impact of regulatory variants from dna sequence. *Nature genetics*, 47(8):955, 2015.
- [19] Christina Leslie and Rui Kuang. Fast string kernels using inexact matching for protein sequences. *The Journal of Machine Learning Research*, 5:1435–1455, 2004.
- [20] Christina S. Leslie, Eleazar Eskin, and William Stafford Noble. The spectrum kernel: A string kernel for svm protein classification. In *Pacific Symposium on Biocomputing*, pages 566–575, 2002.

- [21] Raymond J Mooney and Razvan C Bunescu. Subsequence kernels for relation extraction. In *Advances in neural information processing systems*, pages 171–178, 2006.
- [22] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- [23] Qian Qin and Jianxing Feng. Imputation for transcription factor binding predictions based on deep learning. *PLoS computational biology*, 13(2):e1005403, 2017.
- [24] Avanti Shrikumar, Eva Prakash, and Anshul Kundaje. Gkmexplain: fast and accurate interpretation of nonlinear gapped k-mer svms. *Bioinformatics*, 35(14):i173–i182, 2019.
- [25] Ritambhara Singh, Arshdeep Sekhon, Kamran Kowsari, Jack Lanchantin, Beilun Wang, and Yanjun Qi. GaKCo: A fast gapped k-mer string kernel using counting. In Michelangelo Ceci, Jaakko Hollmén, Ljupčo Todorovski, Celine Vens, and Sašo Džeroski, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 356–373. Springer International Publishing, 2017.
- [26] Gary D Stormo. Dna binding sites: representation and discovery. *Bioinformatics*, 16(1):16–23, 2000.
- [27] Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley-Interscience, 1998.
- [28] SVN Vishwanathan, Alexander Johannes Smola, et al. Fast kernels for string and tree matching. *Kernel methods in computational biology*, pages 113–130, 2004.
- [29] Christopher KI Williams and Matthias Seeger. Using the nystrom method to speed up kernel machines. In *Advances in neural information processing systems*, pages 682–688, 2001.
- [30] Kai Zhang, Liang Lan, Zhuang Wang, and Fabian Moerchen. Scaling up kernel svm on limited resources: A low-rank linearization approach. In *Artificial intelligence and statistics*, pages 1425–1434, 2012.