

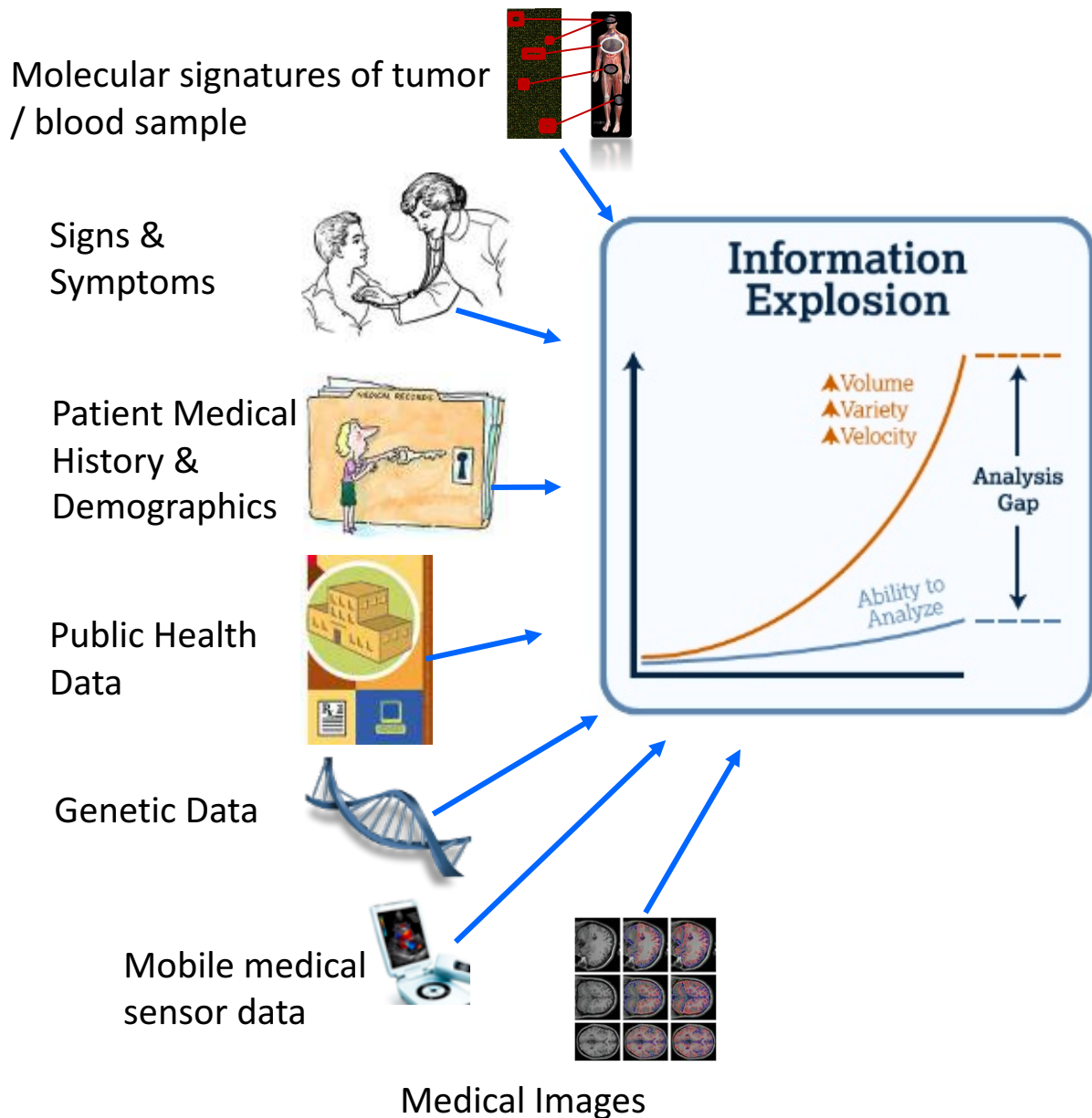


# GaKCo

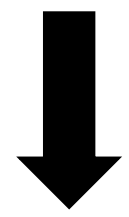
## A Fast Gapped k-mer String Kernel using Counting

Ritambhara Singh, Arshdeep Sekhon, Kamran Kowsari, Jack Lanchantin,  
Beilun Wang, and Yanjun Qi

# Challenge of data explosion



Traditional Approaches



Data-Driven Approaches

**Machine Learning**

# String Classification

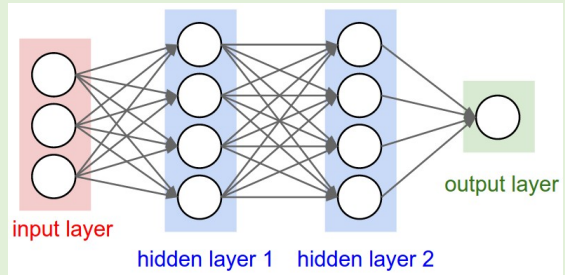
X

This Food is not good.

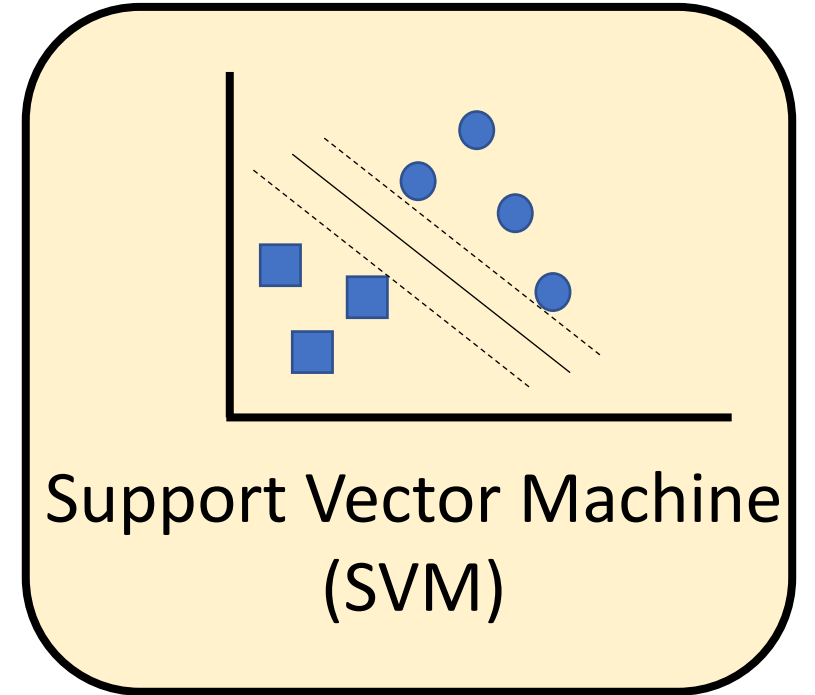


NO  
(-1)

# State-of-the-art Machine Learning Models

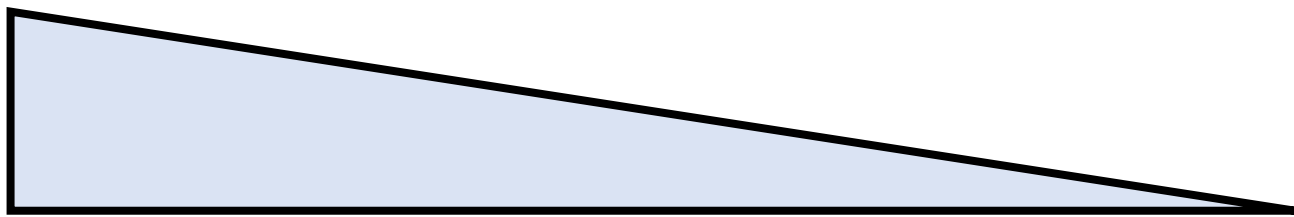
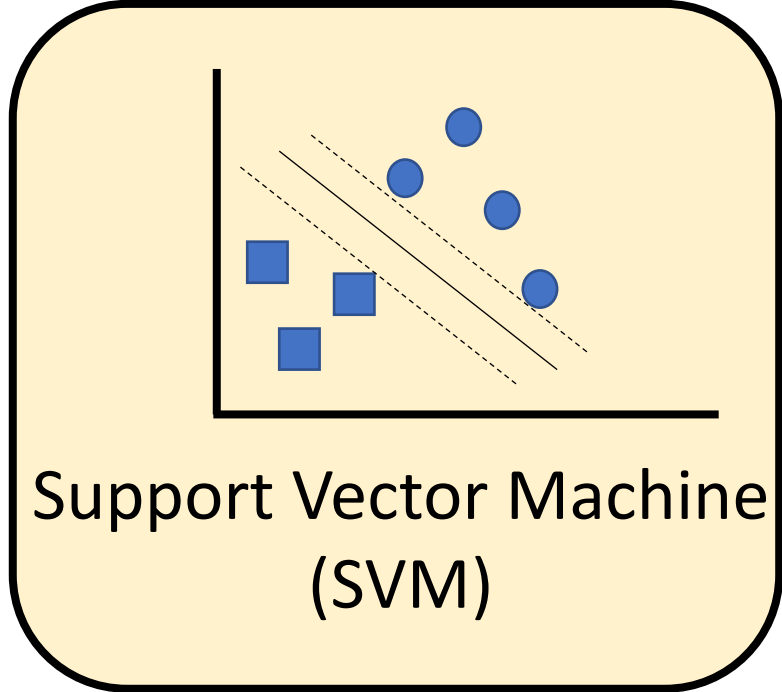
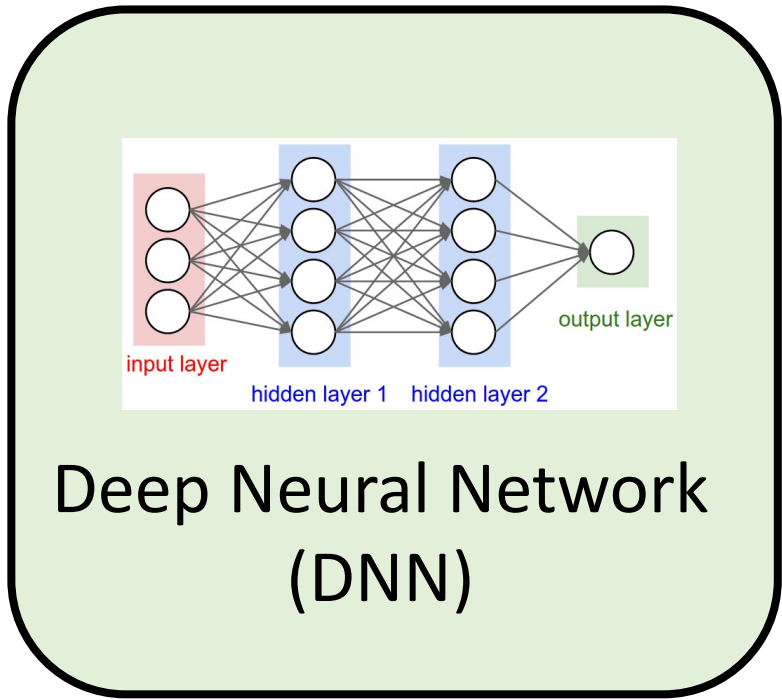


Deep Neural Network  
(DNN)



Support Vector Machine  
(SVM)

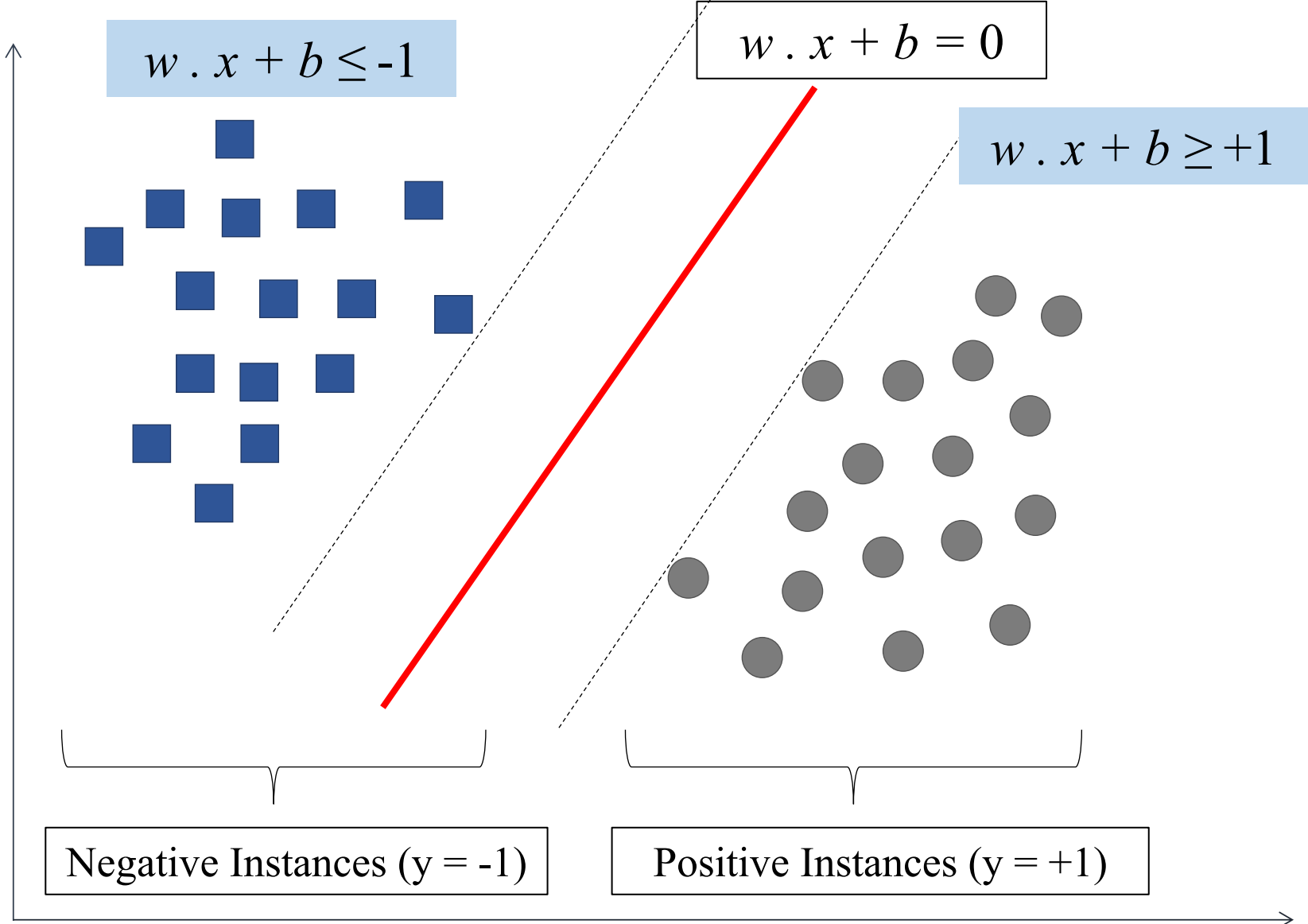
# State-of-the-art Machine Learning Models



Number of Training Samples

# Support Vector Machine (SVM)

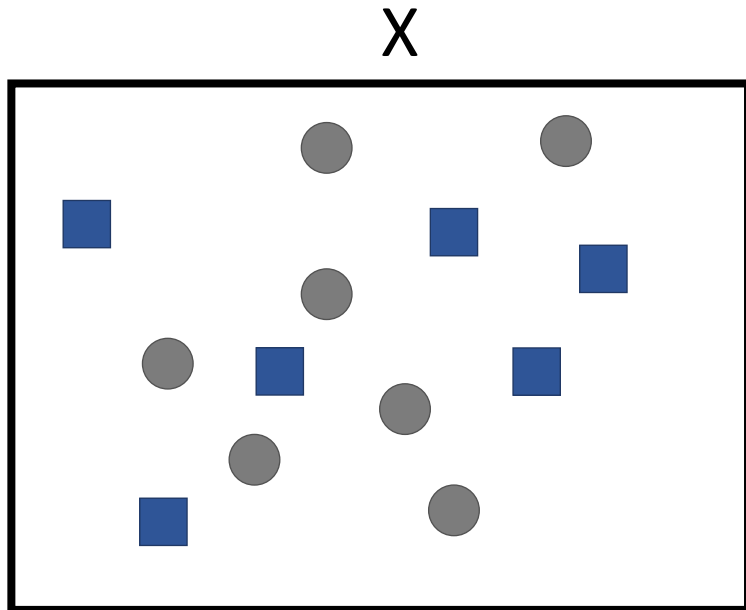
$$f(\mathbf{x}) = \mathbf{w} \cdot \mathbf{x} + b$$



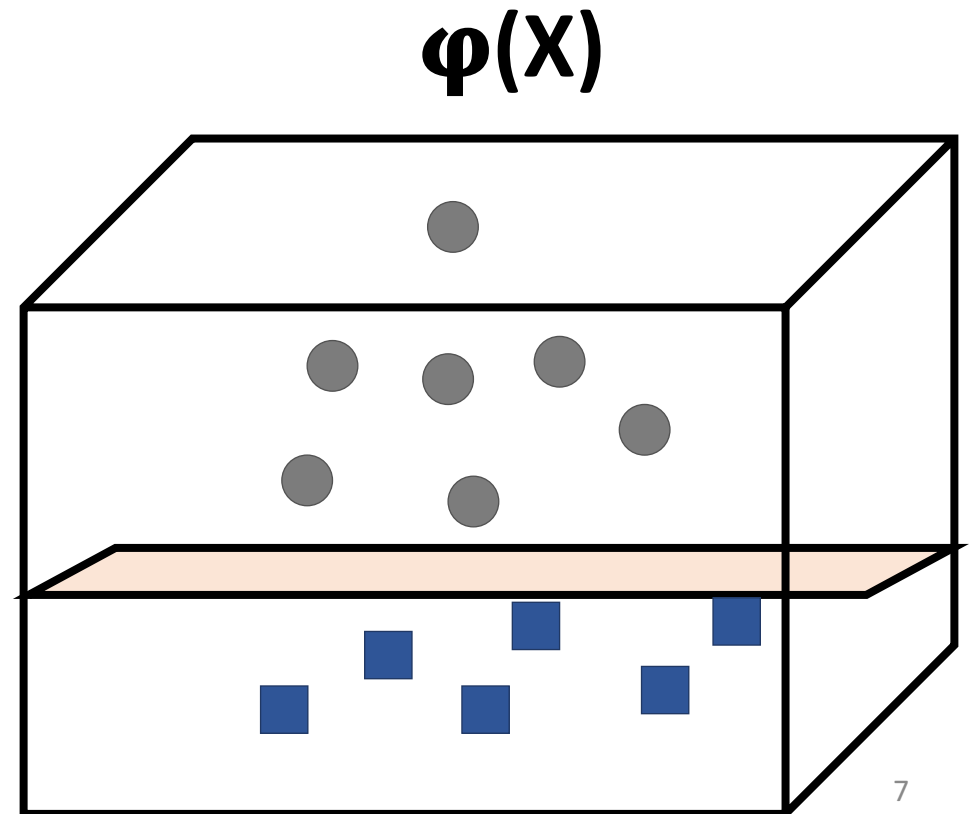
# Kernel

$X$

$$f(x) = w \cdot \Phi(x) + b.$$

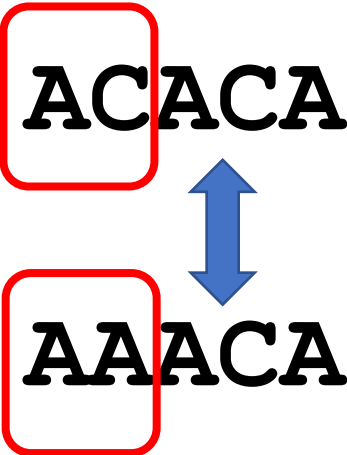
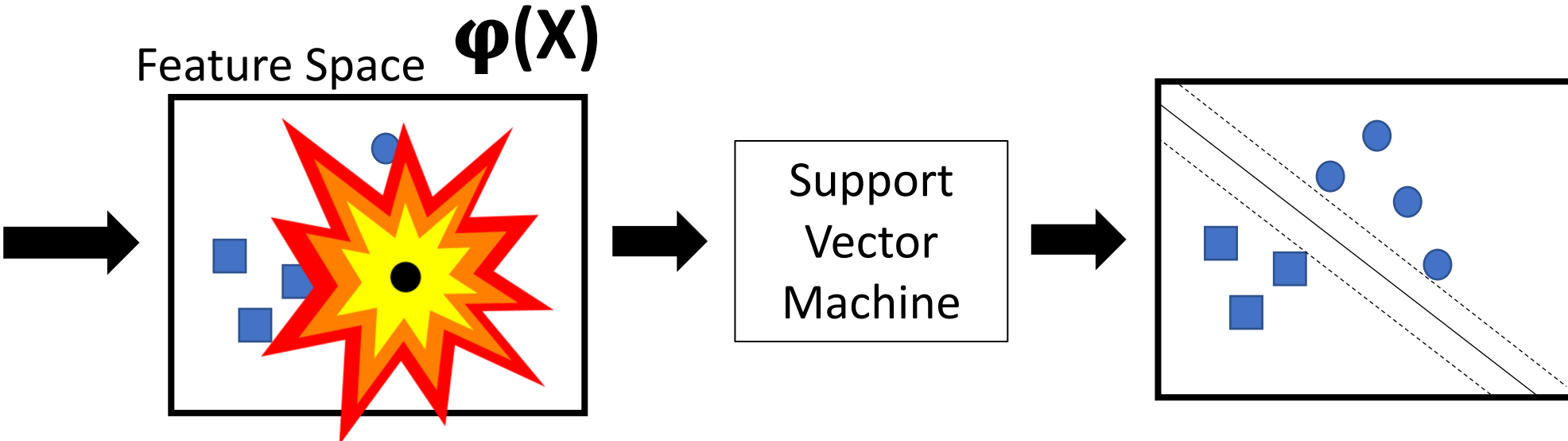


$K(\cdot)$



# Challenge : SK-SVM methods are slow

**S=ACACA**  
**T=AAACA**

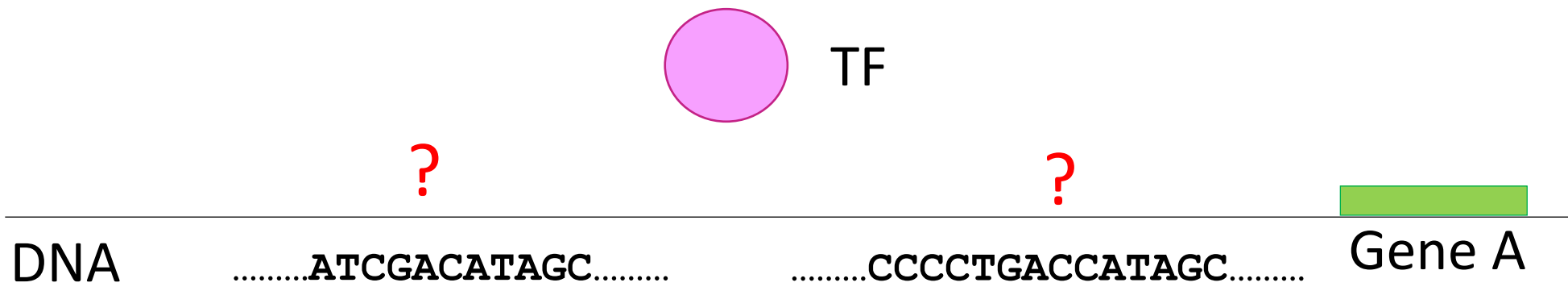




# Solution: GaKCo

- Faster than state-of-the-art string kernel
- Independent of dictionary size
- Parallelizable

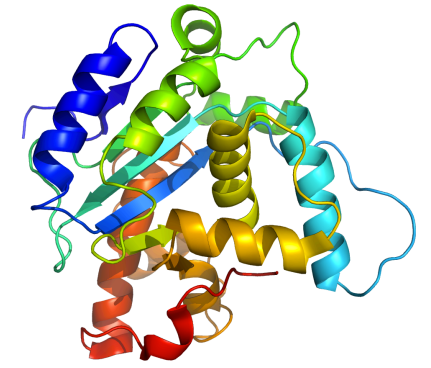
# Does TF bind to this sequence?



- ATCGAATCCG ✓
- CGCTGAATCG ✗
- ATCGCTATCG ✓
- ATCCCGCTCG ✗

$\Sigma=4$   
Minimum Number  
of Training Samples  
**1000**

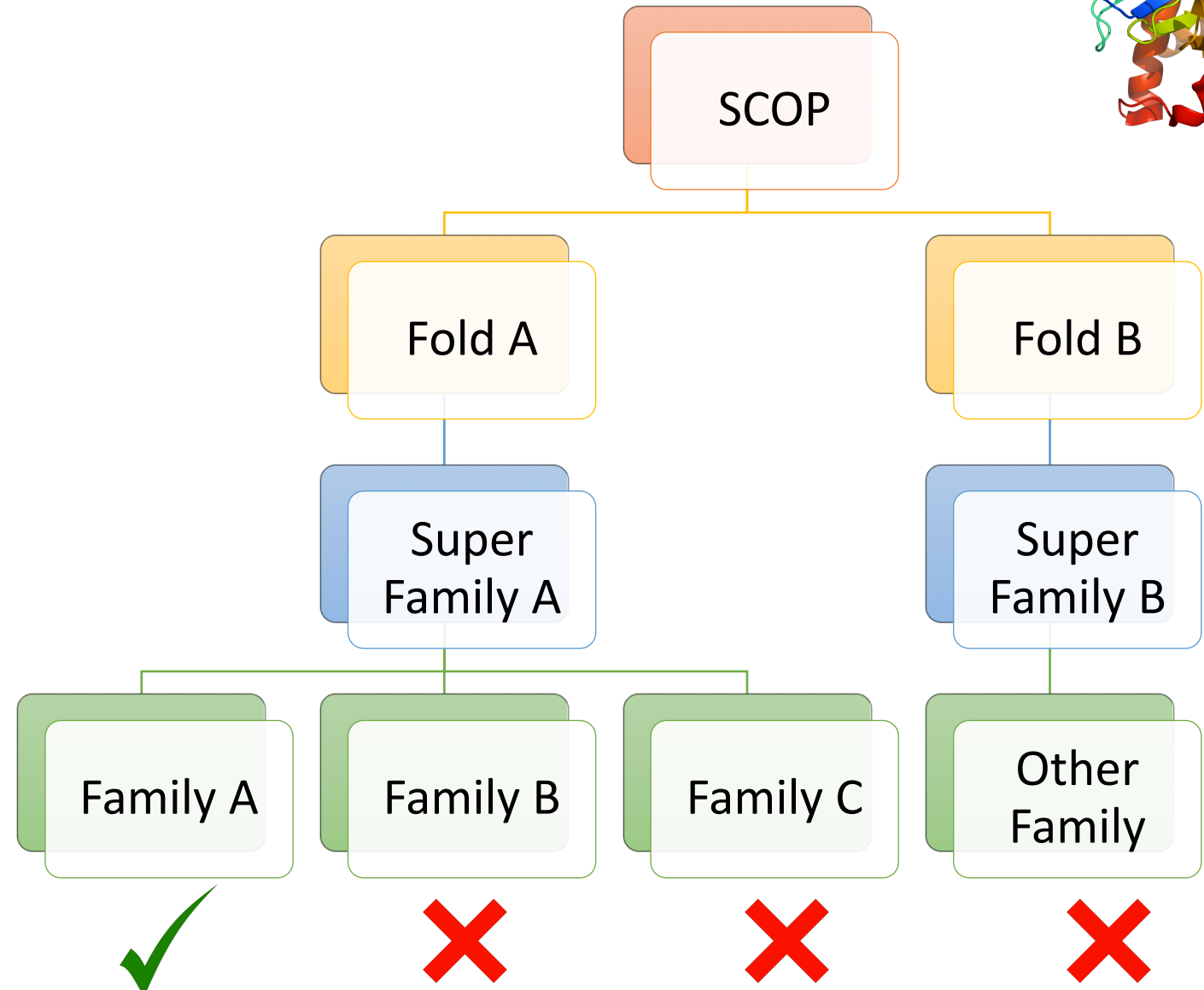
# What family does the protein belong to?



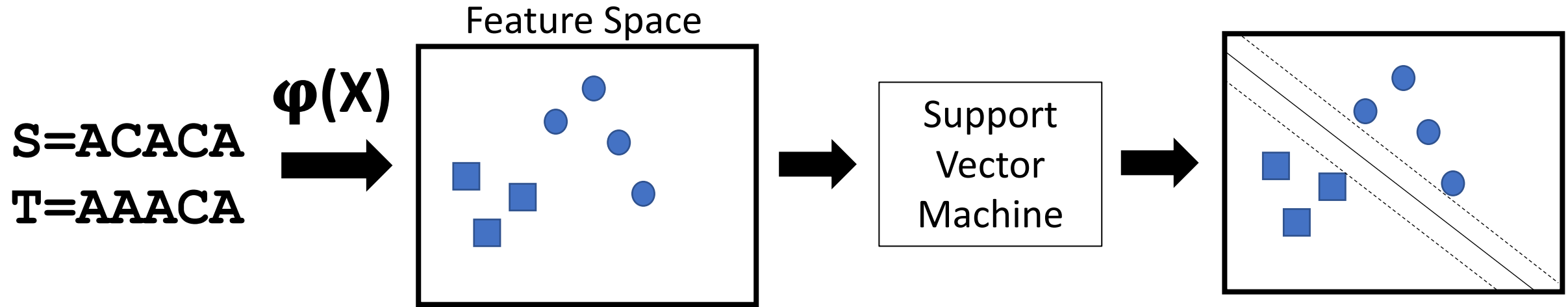
VQGGHACCAKKQQQ

$$\Sigma=20$$

Minimum Number  
of Training Samples  
**500**



# String Kernel + SVM Framework



# String Kernel

## Mismatch Kernel

$\Sigma=2 \{A,C\}$

$k=3$

$m=1$

S: **ACA**ACA

T: **AAA**ACA

k-mers

ACA

CAC

ACA

Feature  
Space

0	AAA	1
0	AAC	1
2	ACA	1
0	CAA	0
1	CAC	0
0	CCC	0

k-mers

AAA

AAC

ACA

**Mismatch  
Neighborhood**  
{AAC, ACA, CAA..}

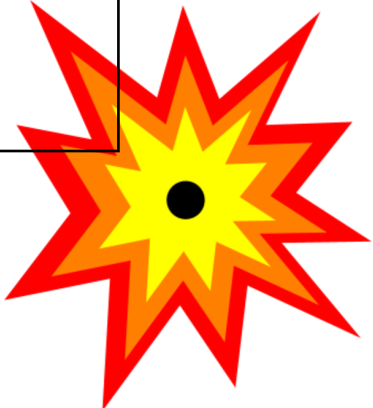
Bag of  
k-mers

# Computational Challenge

$$\Sigma=20$$

$$k=10$$

**Feature Space**  $\Sigma^k$



$$(20)^{10} = 10^{13}$$

	AAA	1
	AAC	1
<b>Sparse!</b>	ACA	0
	CAA	0
	CAC	0
	CCC	0

# Related Work: Gapped k-mer Kernel

$\Sigma=2 \{A,C\}$

$g=3$

$k=2$

S: **ACA**CA

T: **AAA**CA

g-mers    k-mers  
ACA    {AC,CA,...}  
CAC  
ACA

g-mers    k-mers  
AAA    {AA,AA,...}  
AAC  
ACA

Gaps= $g-k=1$

**Feature Space**     $g(|S| + |T|) \ll \Sigma^k$

# Related Work: Gapped k-mer Kernel

$$(1) \quad K(x, x') = \sum_{\gamma \in \Theta_g} c_x(\gamma) \cdot c_{x'}(\gamma)$$

$$(2) \quad K(x, x') = \sum_{i=0}^{l_1} \sum_{j=0}^{l_2} h_{gk}(g_i^x, g_j^{x'})$$

$$(3) \quad K(x, x') = \sum_{m=0}^{g-k} N_m(x, x') h_m$$



# Related Work: Gapped k-mer Kernel

$\Sigma=2 \{A,C\}$

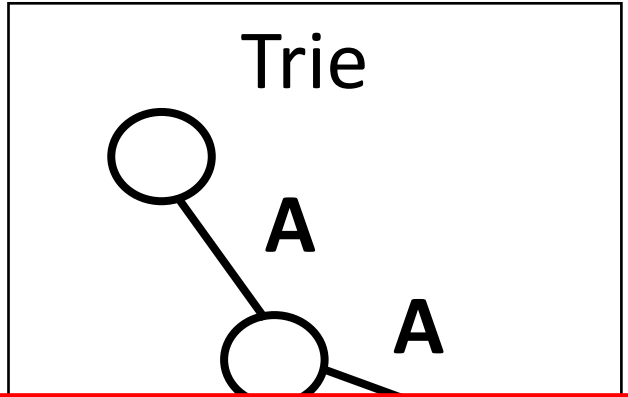
**Dependence on term  $(\Sigma)^m$**

$g=3$

S:ACACA

T:AAACA

gkm-SVM



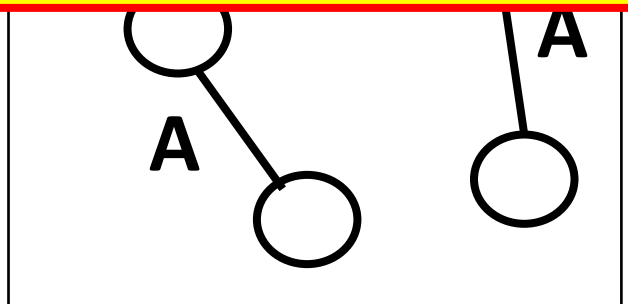
**Time taken to compute Kernel for Protein dataset > 5 HOURS**

ACA  
CAC  
ACA

Mismatch Profile

AAA  
AAC  
ACA

Coefficient



Gaps=g-k=1

$$K(S, T) = \sum_{m=0}^{g-k} h_m N_m$$

m=0

$N_m = 2$

m=1

$N_m = 3$

# GaKCo : Gapped k-mer Kernel using Counting

S:ACACA

T:AAACA

$g=3$

$m=1$

CA  
AC  
CA  
A  
A  
C

AA  
AC  
AC

AA  
AA  
AA

AA  
AA  
AC  
AC  
AC  
CA

**Independent of term  $(\Sigma)^m$**   
**Parallelizable**

$C_m(S,T)=1+2=3$

$C_m(S,T)=4$

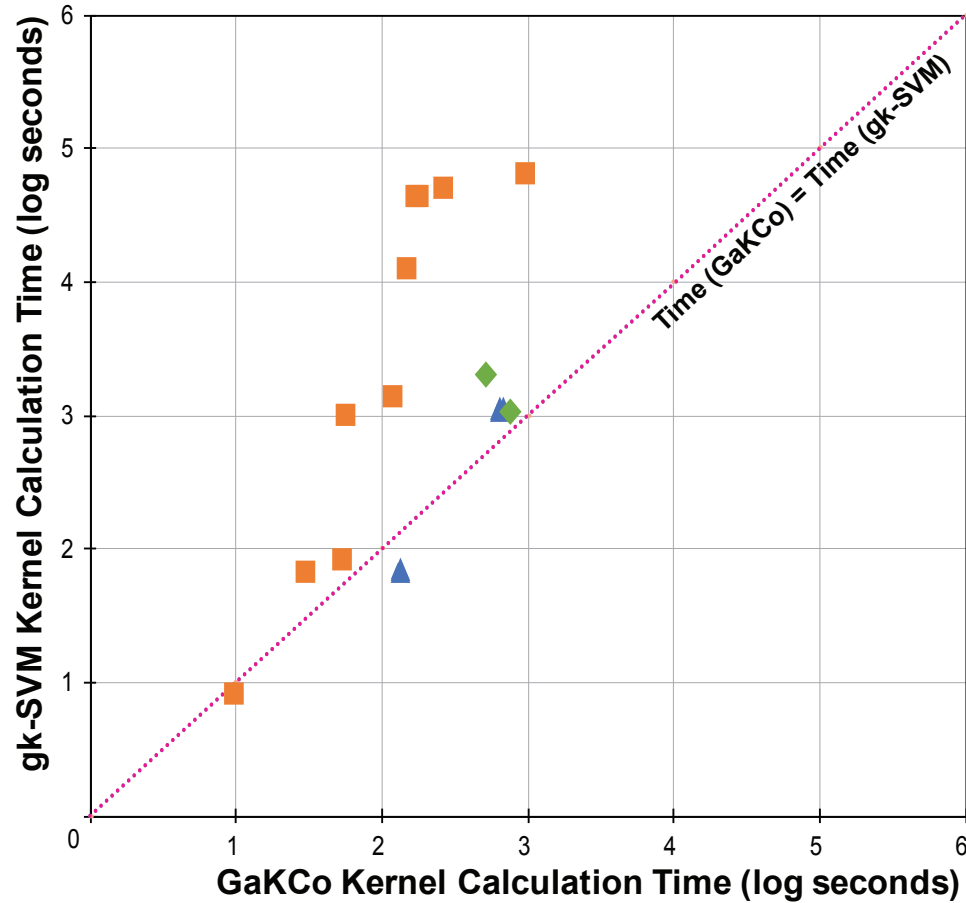
$C_m(S,T)=2$

Total  $C_m(S,T)=3+4+2=9$

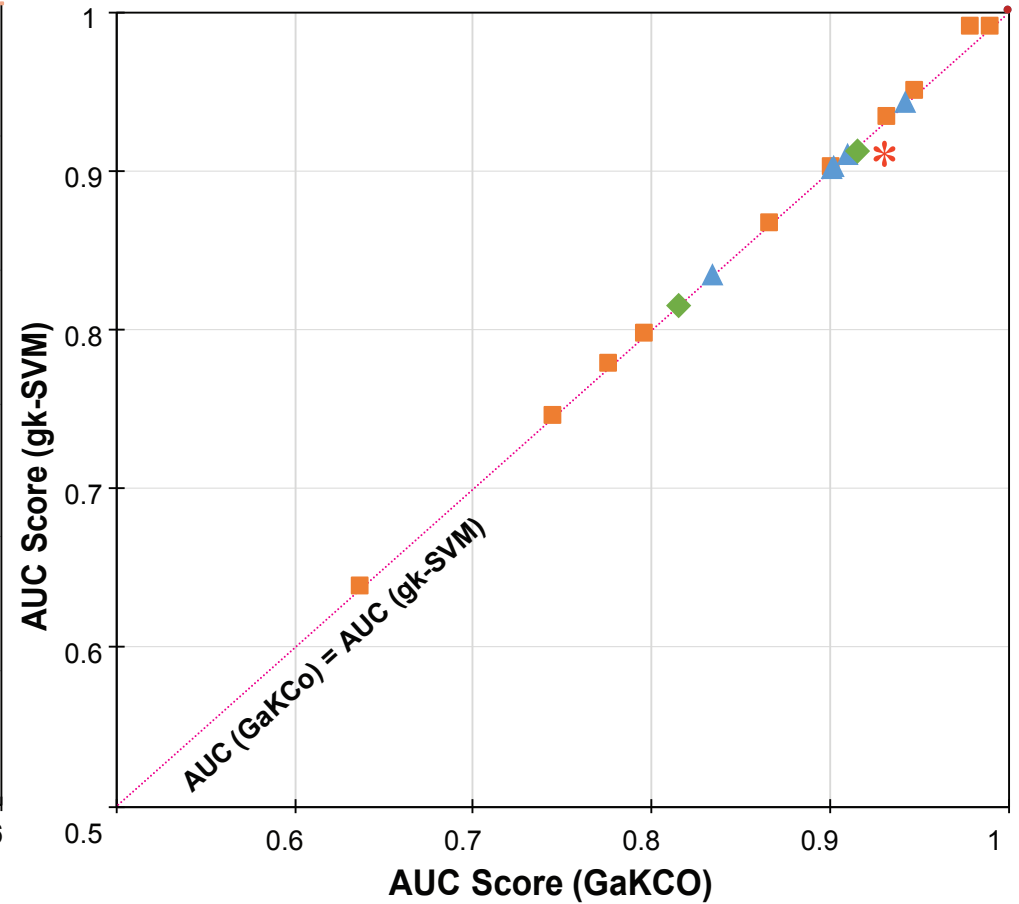
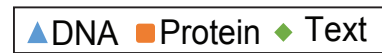
Over Counting =  $\binom{3}{1} N_{m=0} = 3 \times 2 = 6$

$N_{m=1}(S,T) = 9 - 6 = 3$

# GakCo is Faster!



(a)

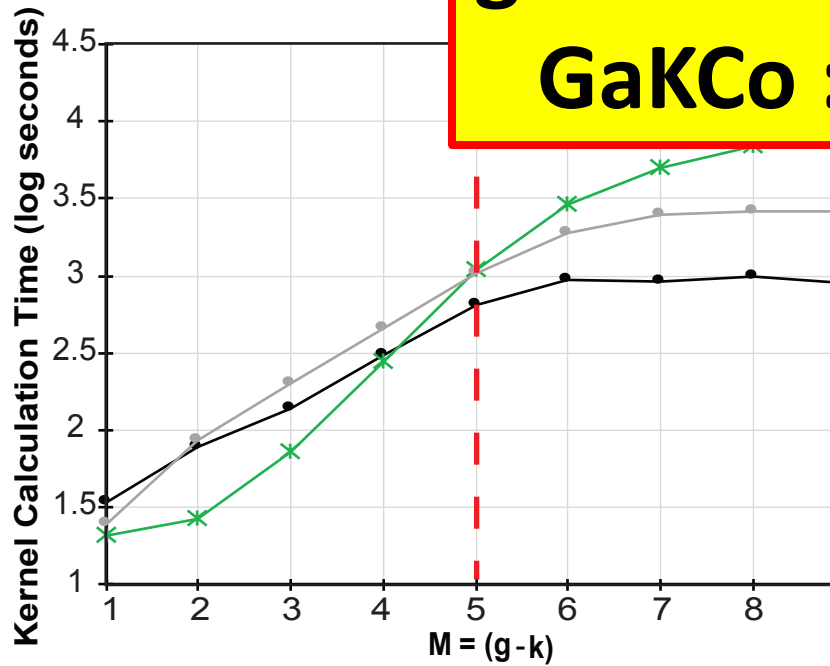


(b)

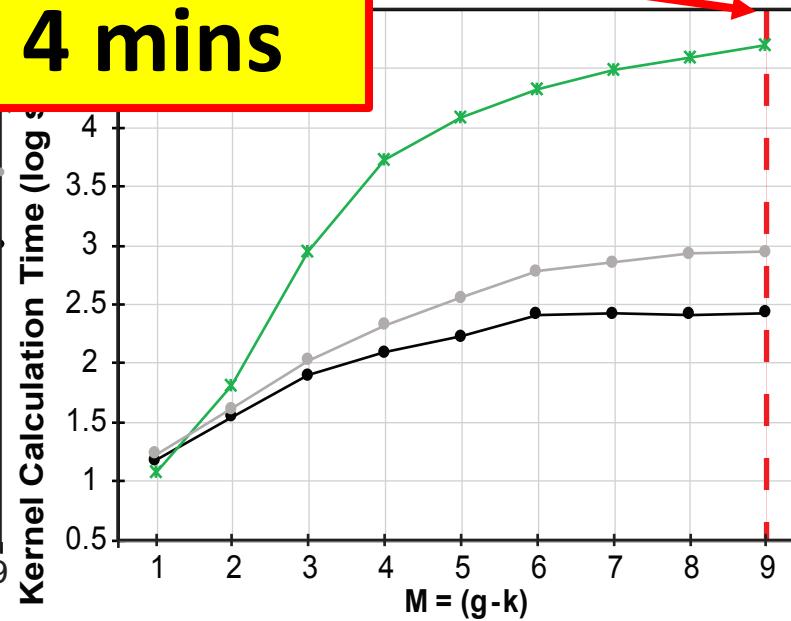
\* Micro-averaged F1-score

# Scales well with increasing $\Sigma$ and $m$

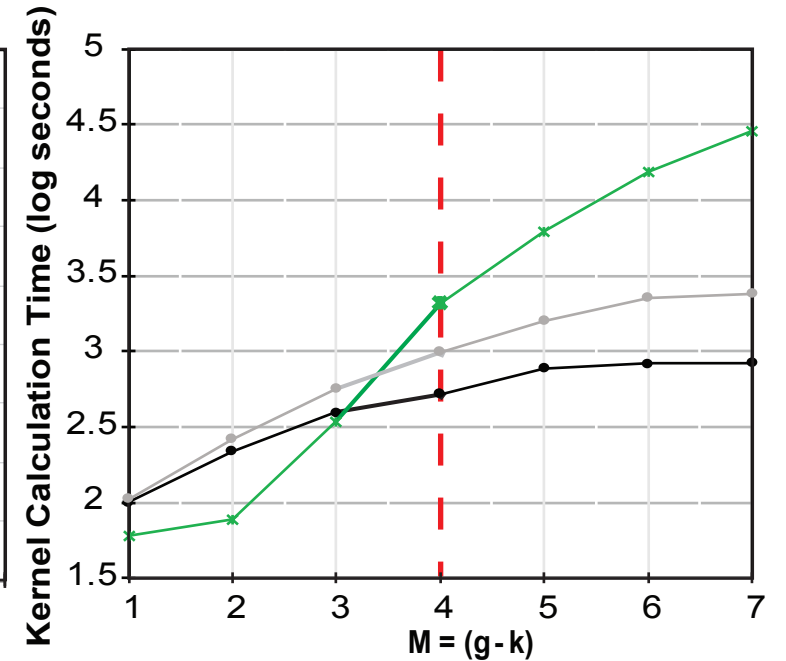
**gkm-SVM : > 5 hrs**  
**GaKCo : 4 mins**



(a) DNA (EP300)



(b) Protein (1.34)

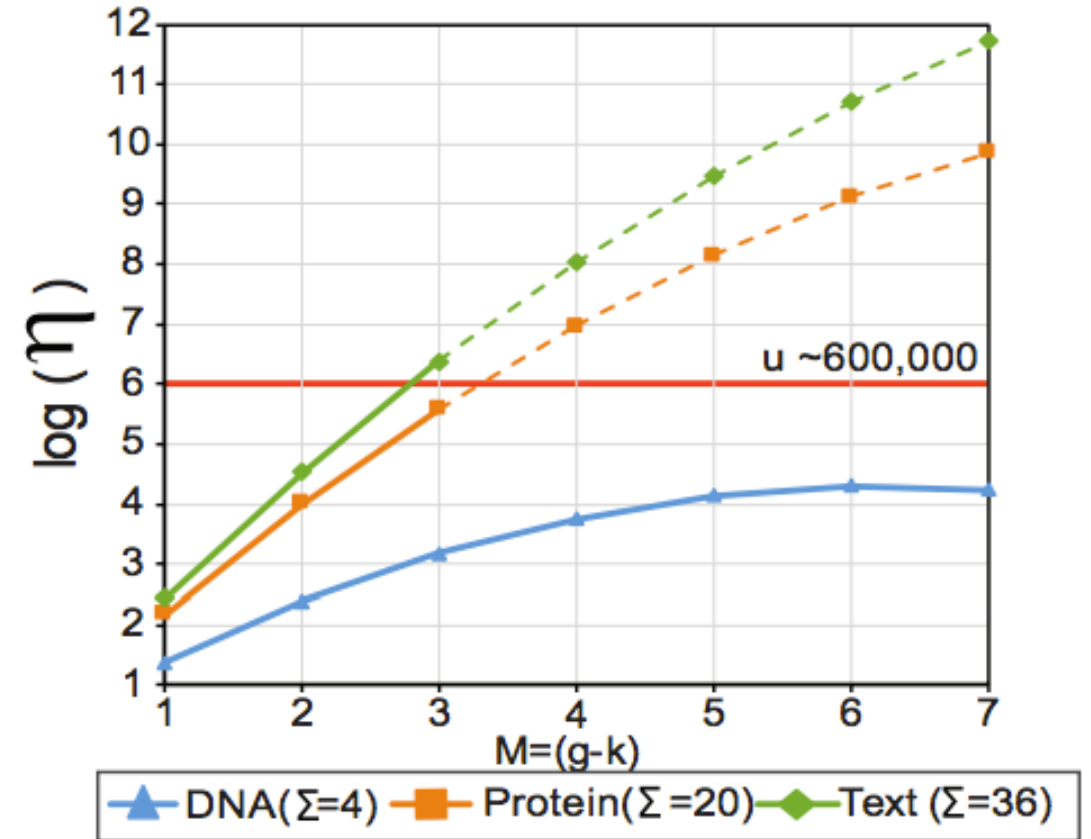


(c) Text (Sentiment)

● GaKCo    ● GaKCo (Single thread)    \* gkm-SVM

# Detailed Theoretical Analysis (in paper)

	GaKCo	gkm-SVM
Pre-processing	$c_{gk}gNl$	$gNl + \eta u g$
Kernel updates	$c_{gk}zN^2$	$\eta u N^2$



# Summary: GaKCo

- Fast – (1) reduced gapped k-mer feature space and (2) counting based method
- Scalable – allows larger dictionary size and gap size
- Parallelizable – naturally parallelizable implementation
- Code and datasets available at: <https://github.com/QData/GaKCo-SVM>
- ArXiv: <https://arxiv.org/abs/1704.07468>

Thank you

