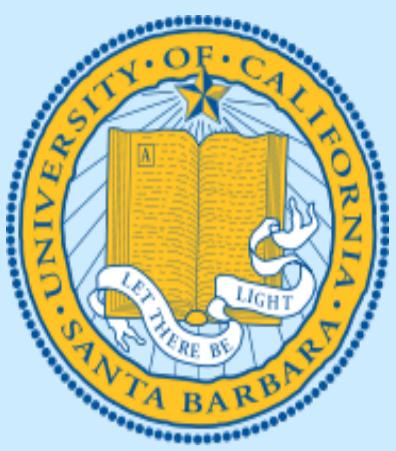


# Fast Motif Discovery in Short Sequences



Honglei Liu<sup>1</sup>, Fangqiu Han<sup>1</sup>, Hongjun Zhou<sup>2</sup>, Xifeng Yan<sup>1</sup>, Kenneth S. Kosik<sup>2</sup>

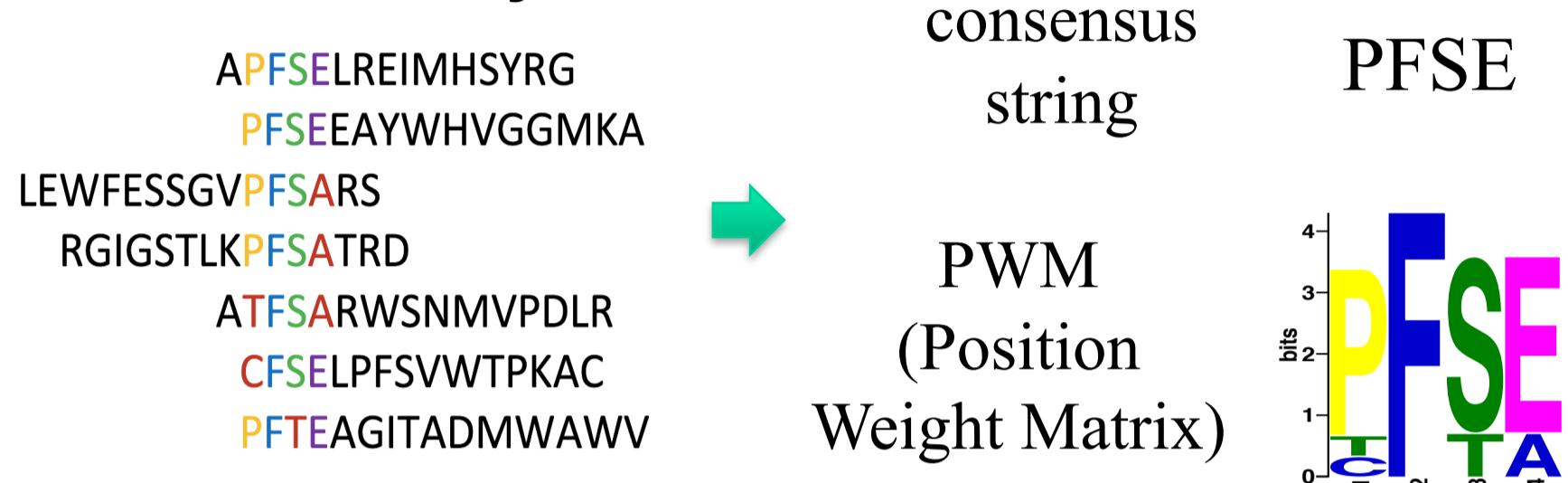
<sup>1</sup>Department of Computer Science, , University of California, Santa Barbara

<sup>2</sup>Neuroscience Research Institute, University of California, Santa Barbara

## Introduction

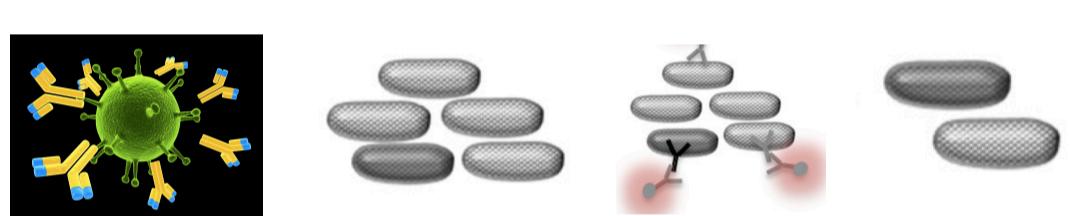
- **Motif:** frequently appearing sequence patterns

- **Motif discovery:**



- **Applications**

- Transcription factor binding sites (TFBSs) discovery
- Antibody biomarkers discovery



ESNTCDL**F**VWQACDGKQ  
AEVACEDN**F**VYQCSDDW  
SSASCD**M**FVYQQCAEFN  
RQGACV**D**DYVYQCGHFE  
GHTACMTD**F**VHQCFPGT  
PCVDA**F**VYQQSGCNIA  
RDGHCADS**F**VNQCVRPL  
GRAACV**D**D**F**VYQCVRQHE

Large scale, Large alphabet set, Short

## Challenges

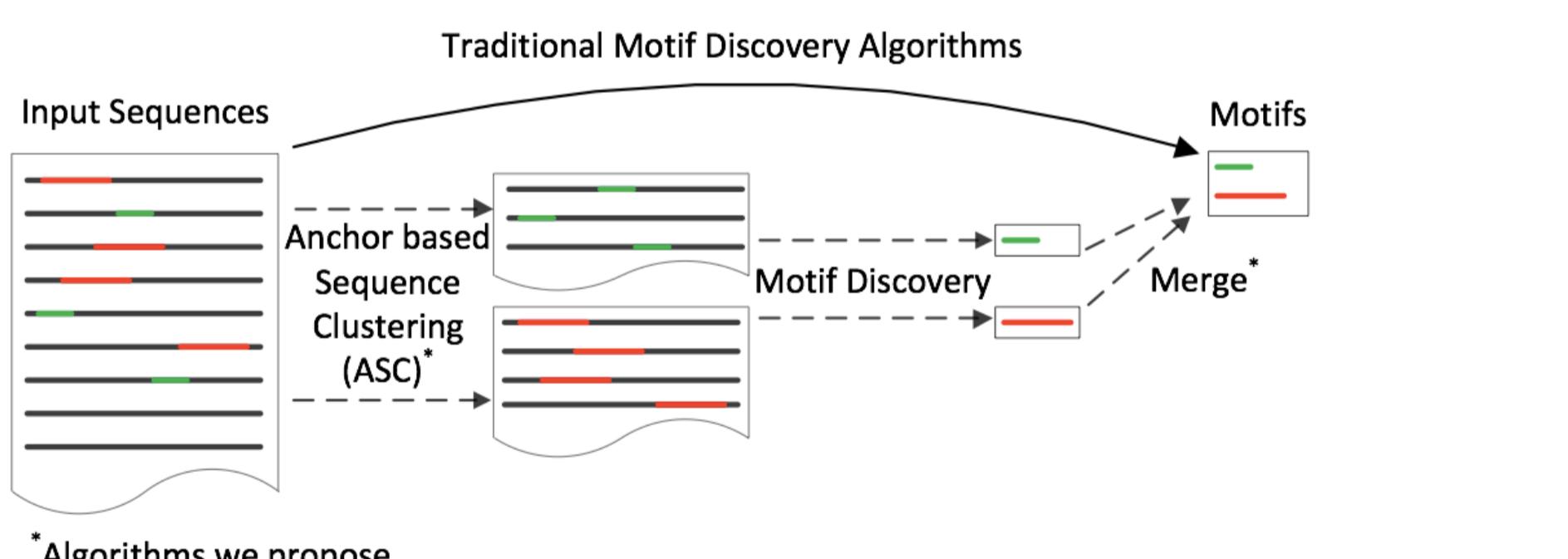
- Unknown: number of motifs, length of motifs, etc
- Before next-generation sequencing era
  - At most several hundred sequences
- After next-generation sequencing era
  - Tens of thousands or even millions of sequences
- Existing methods can not handle the big data challenge very well

	Can handle >10k seq.	Can handle >1M seq.	Can work with protein seq.	Accuracy is as good as MEME
<b>MEME</b> <sup>[Bailey94,06]</sup>			✓	✓
<b>STEME</b> <sup>[Reid11]</sup>	✓	✓		✓
<b>DREME</b> <sup>[Bailey11]</sup>	✓	✓		
<b>GibbsCluster</b> <sup>[Andreatta13]</sup>	✓		✓	
<b>MUSI</b> <sup>[Kim11]</sup>	✓		✓	
<b>Our framework</b>	✓	✓	✓	✓

## Methods

- **Our framework**

- Reuse existing techniques!

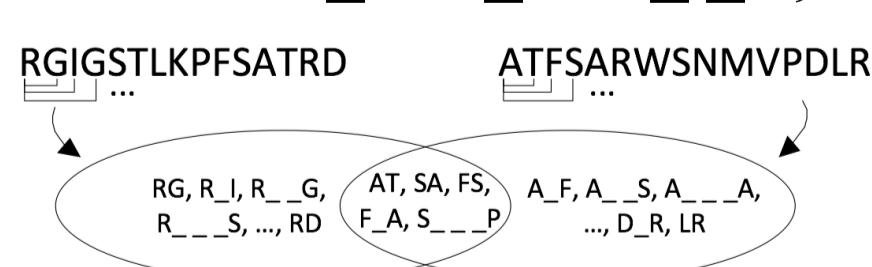


- **Anchor based Sequence Clustering algorithm (ASC)**

- Could capture local similarities
- Avoid pairwise comparisons

- **Anchor based similarity**

- Represent sequences as  $q$ -anchor sets
- e.g. 2-anchors of PFSE are  $\{PF, FS, SE, P\_S, F\_E, P\_E\}$

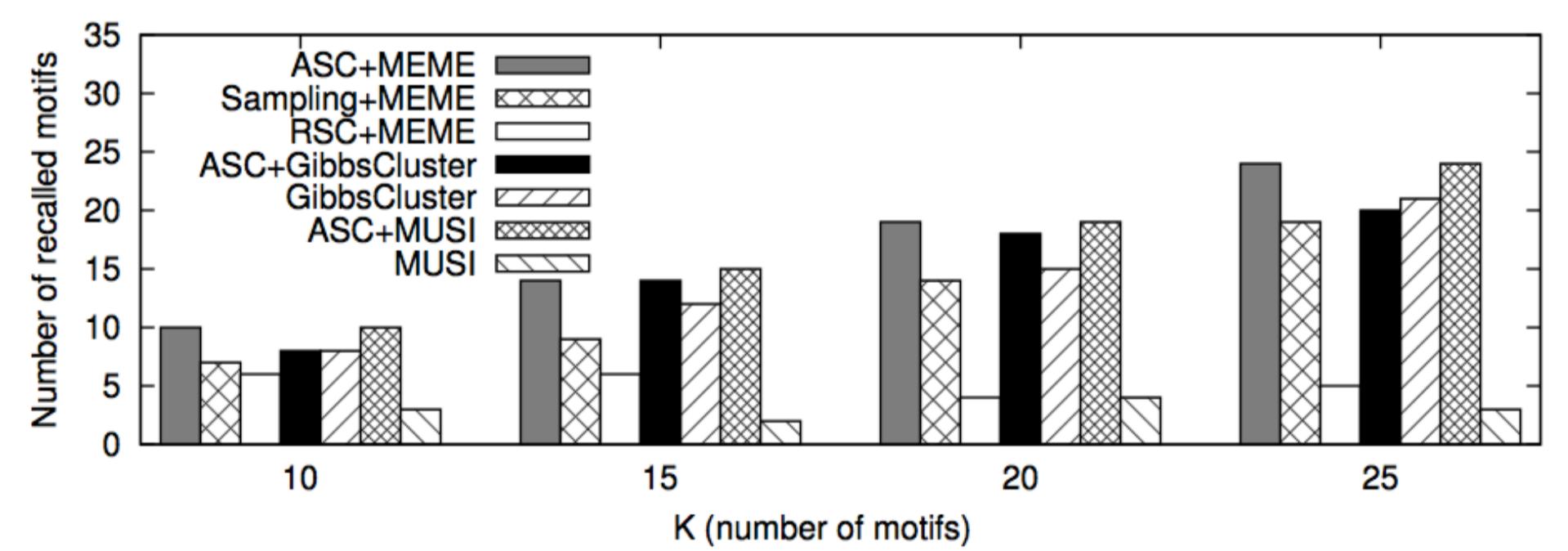


- **Iterative process**

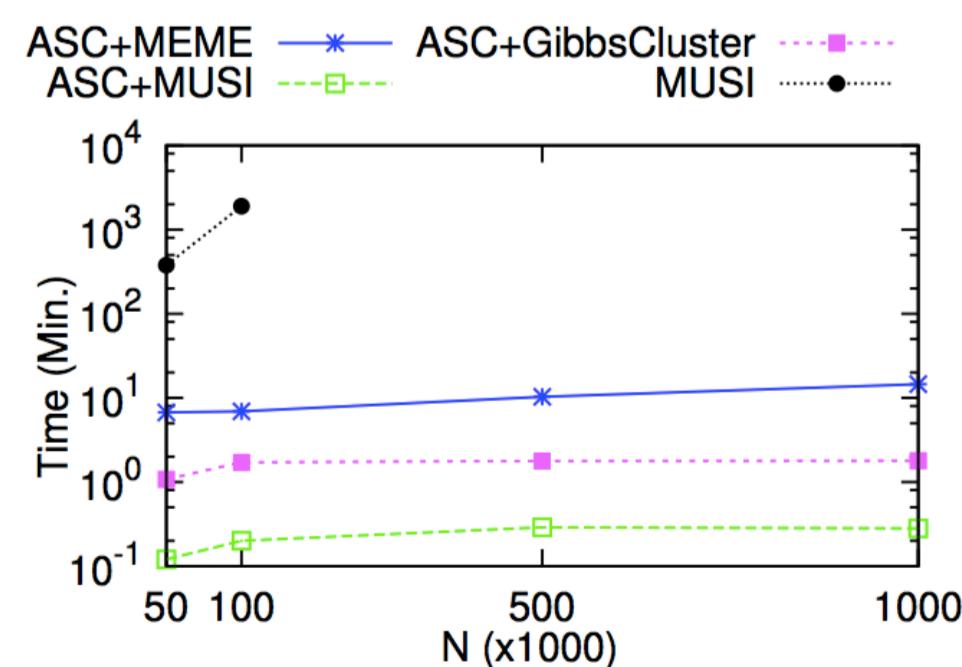
- Choose initial centers using *odd score*
  - Indicates how likely an anchor is from a motif
- Adjust centers using *abundance score*
  - Indicates how unique an anchor is for a motif

## Experiments

- Real data shows that our framework can reduce the runtime of MEME from **weeks** to **minutes** without losing accuracy!
- Apply ASC on top of MEME, MUSI and GibbsCluster
- Number of recalled motifs from different methods using synthetic data (10k seq.)



- Scalability



## Conclusions

- Big data challenge
- Reuse existing techniques
- Huge performance gain without losing accuracy