

# *Chromatin Conformation Prediction from ChIPseq*

Ricky Lim<sup>1</sup>, Samuel Collombet<sup>2</sup>, Agus Salim<sup>3</sup>, Touati Benoukraf<sup>1</sup>

<sup>1</sup>CSI-NUS <sup>2</sup>Ecole Normale Supérieure <sup>3</sup>La Trobe University

CSI-Meeting  
<mailto:rlim.email@gmail.com>

# Contents

## 1 Introduction

### Goal

What is ChIPseq?

What is Mixture Model (MM)?

## 2 Preliminary results

### Pipeline

Input Data

Peak Calls: MACS2 vs jaHMM

### Summary: Peak Calls

Component Calls

### Summary: Component Calls

Motif Calls

### Summary: Motif Calls

Model Update: Biclustering

## 3 Summary and Future

# Chromatin Conformation Prediction

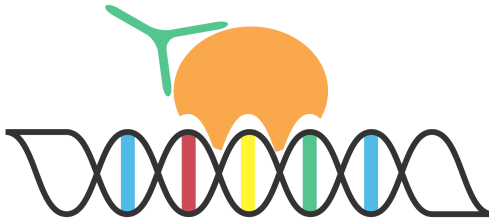
- **Main Question:** Can we use transcription factor (TF)-ChIPseq to predict protein complexes (direct and indirect bindings) on chromatin?

# Chromatin Conformation Prediction

- **Main Question:** Can we use transcription factor (TF)-ChIPseq to predict protein complexes (direct and indirect bindings) on chromatin?
- **Strategy:** Model ChIPseq signal using Mixture Models to cluster the direct and indirect bindings.

What is ChIPseq?

# ChIP-Seq

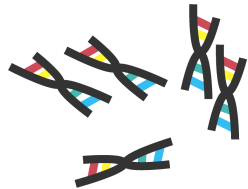


Chromatin  
ImmunoPrecipitation



Sequencing

ATCGTTTAACGCATTAGCAGT...



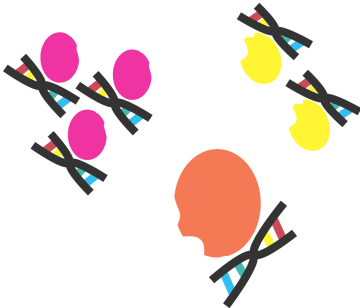
# Chromatin Conformation



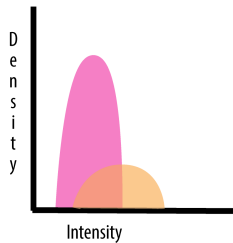
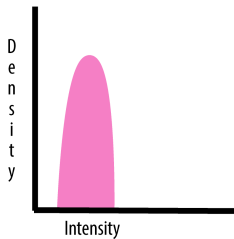
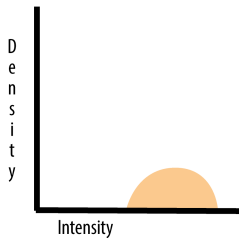
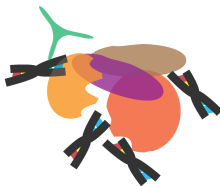
Direct binding sites



Indirect binding sites



# Mixture of Chromatin Conformations





What is Mixture Model (MM)?

# Mixture Model (GMM): Revisited

Types of clustering methods:

- Hard clustering: non-overlapping clusters
- Soft clustering: overlapping clusters

# Mixture Model (GMM): Revisited

Types of clustering methods:

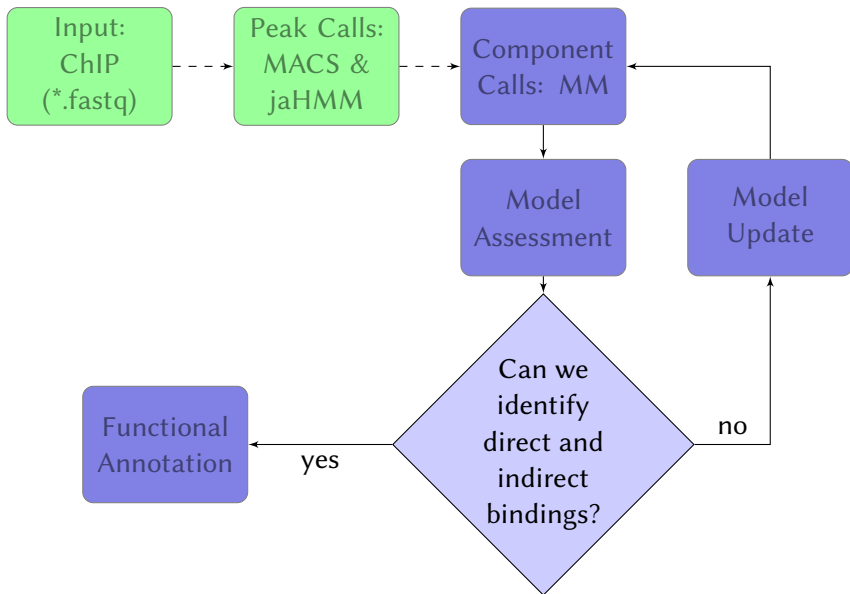
- Hard clustering: non-overlapping clusters
- Soft clustering: overlapping clusters

MM is a probabilistic way of soft clustering. Each cluster is a generative mixture model (pdf) with its parameters.

## Mixture Gaussian pdf:

Key Assumption:

- ChIP-seq peaks are drawn from a finite set of gaussian distributions.
- ChIPseq peaks are fit with gaussian mixture models, with mixing  $\lambda$  parameter.
- Each gaussian corresponds to a cluster of peaks with  $\mu$  and  $\sigma$  parameters.



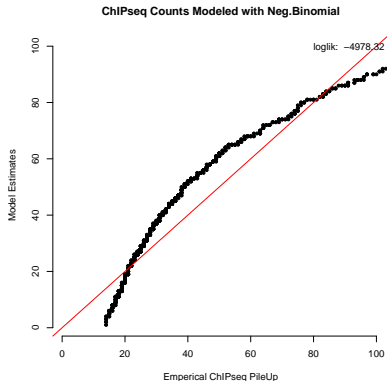
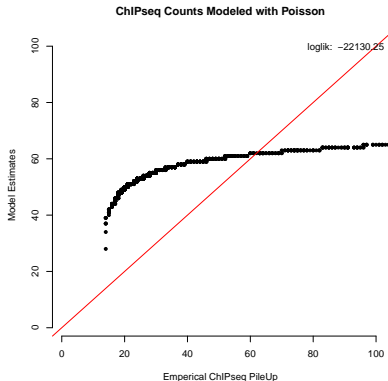
# Input: ChIP-seq of Cebp $\epsilon$ from Koeffler-BM

```
##FastQC 0.10.1
>>Basic Statistics pass
#Measure Value
Encoding Illumina 1.5
Total Sequences 41586141
Sequence length 40
#Summary
PASS Basic Statistics
PASS Per base sequence quality
PASS Per sequence quality scores
WARN Per base sequence content
PASS Per base GC content
PASS Per sequence GC content
PASS Per base N content
PASS Sequence Length Distribution
PASS Sequence Duplication Levels
PASS Overrepresented sequences
WARN Kmer Content
```

# Principles

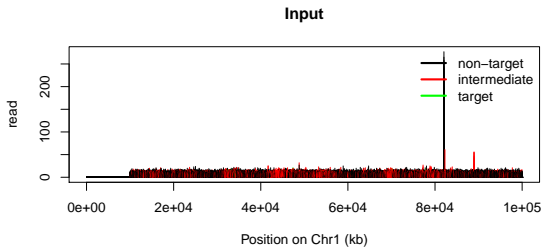
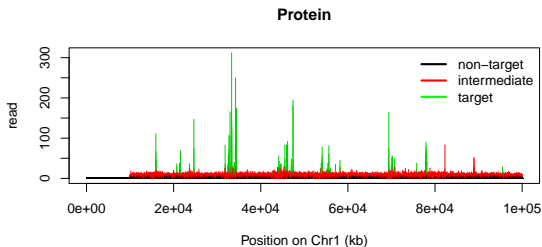
- **MACS2**: *poisson* model-based analysis of Peak calls MACS reference
- **jaHMM**: *negative binomial* model-based analysis of Peak calls jaHMM reference

# Why jaHMM is better?

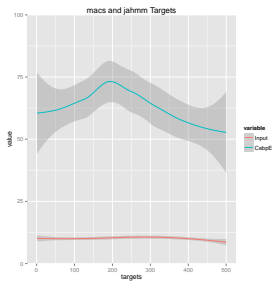
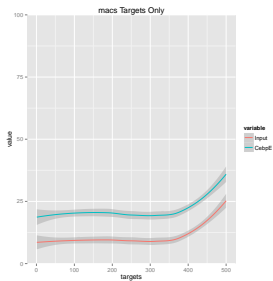
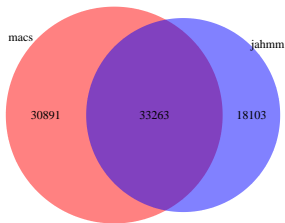




# Targets Identified by jahmm



# Targets Identified by MACS2 vs jahmm



# Why jaHMM is better than MACS2?

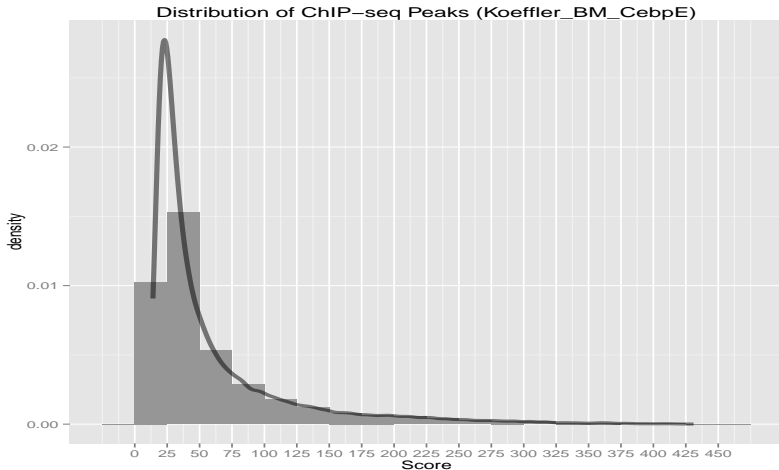
- Given our dataset, negative binomial model assumed by jaHMM fits better than poisson model assumed by MACS2

# Why jaHMM is better than MACS2?

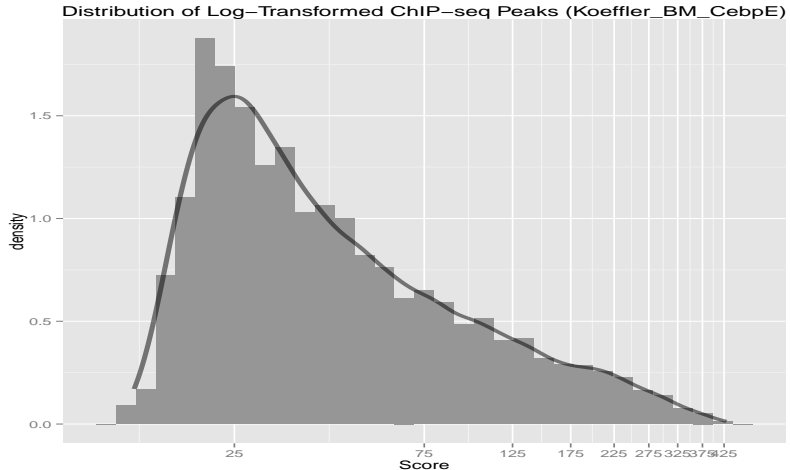
- Given our dataset, negative binomial model assumed by jaHMM fits better than poisson model assumed by MACS2
- Peaks identified solely by jaHMM have scores higher with respect to their input than solely by MACS2

Can we model ChIPseq Peaks using components of MMs?

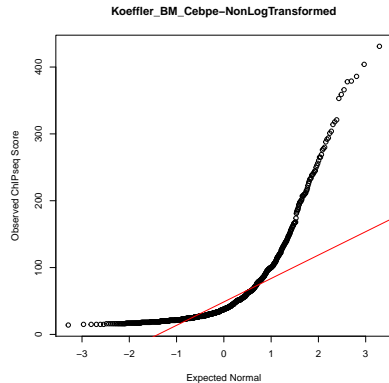
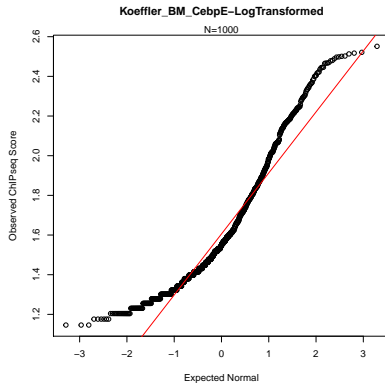
# Input: ChIP-seq of Cebp $\epsilon$ from Koeffler-BM



# Log Transformation of ChIP-seq Input

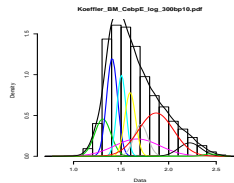
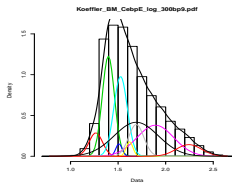
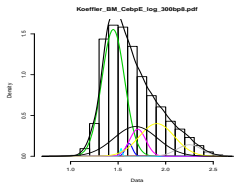
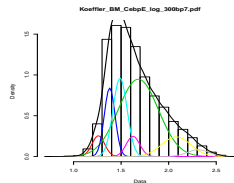
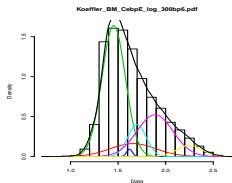
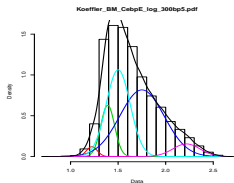
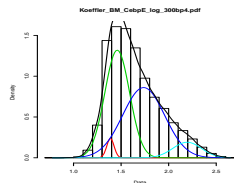
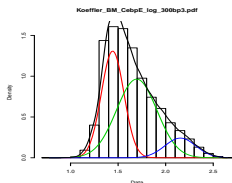
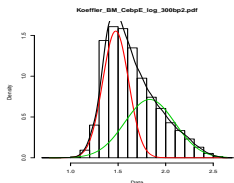


# Check the Normality

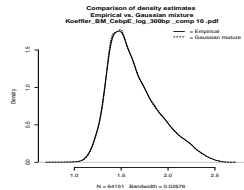
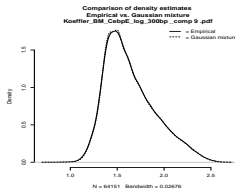
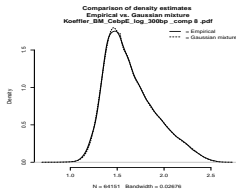
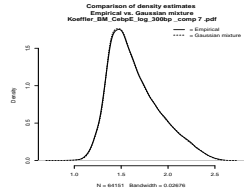
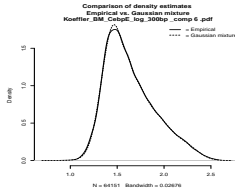
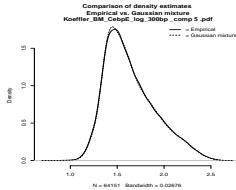
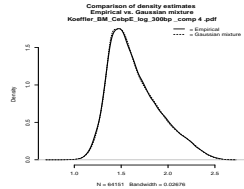
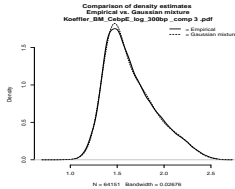
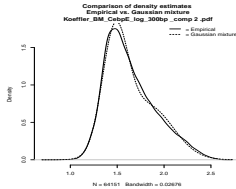




# ComponentCalls: Fit ChIPseq Peaks with GMMs



# GMM-ModelAssessment: Overfit<sup>1</sup>



## Model Assessment: BIC-AIC

AIC and BIC is based on Occam's razor principle, i.e, the simplest the better.

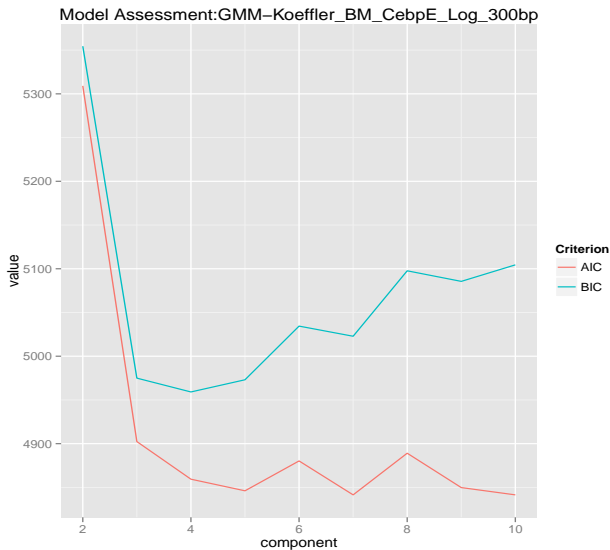
$$\text{AIC} = -2 \times \log L + 2 * P$$

$$\text{BIC} = -2 \times \log L + \log(n) * P$$

$L$  is likelihood

$P$  is the number of parameters

# Model Assessment: BIC-AIC



# Summary

- **Can we model ChIPseq using several components of MMs?**

Yes, our ChIPseq Peaks identified by jaHMM can be fit with GMMs.

# Summary

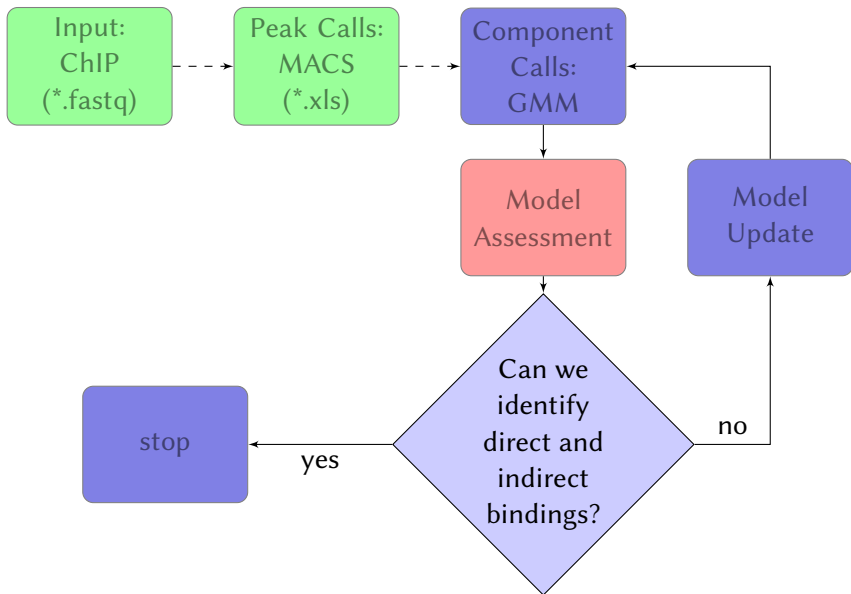
- **Can we model ChIPseq using several components of MMs?**

Yes, our ChIPseq Peaks identified by jaHMM can be fit with GMMs.

- **How many components are required?**

From AIC-BIC and cross-validation, with 3 components are sufficient to fit the ChIPseq.

Note: the lower the AIC and BIC values, the better the fitting.



## Motif Calls using Centdist



# Group1: low peak score (29559 peaks)

2/9/2015

CENTDIST:Koeffler\_BM\_CelpE\_GMM\_ModelAssignment\_log\_300\_group1\_compSorted3.bed

Results for Koeffler\_BM\_CelpE\_GMM\_ModelAssignment\_log\_300\_group1\_compSorted3.bed  
VERSION: 2011.07.08

[Try our De Novo Motif Finding Tool for ChIP-seq \(SEME\)](#)

746 TFs

Show top 50

Rank	Name	Family	Logo	Score	Distribution	%Sequence with motif optimal setting	%Sequence with motif 1e-4 fold within +/- 200bp	Binding Range	P/W/M Score Cutoff	Z0Score	Z1Score
1	V\$janpar_MZF1_1_4	<a href="#">janpar_BetaBeta/Alpha_zinc_finger</a>		12.2743				440	2.7671	6.19578	6.07853
2	V\$janpar_SP1	<a href="#">janpar_BetaBeta/Alpha_zinc_finger</a>		11.5458				480	3.0083	8.28603	3.25976
3	V\$SP1_01	<a href="#">SP1</a>		11.3454				480	2.7192	8.56304	2.78238
4	V\$SP1_Q2_01	<a href="#">SP1</a>		9.69061				480	3.2844	7.55059	2.14002
5	V\$MAZR_01	<a href="#">SP1</a>		9.67953				440	2.9471	5.14373	4.5356
6	V\$MUSCLE_INL_B	<a href="#">MINI</a>		9.64468				480	2.8998	7.04083	2.60384

[http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif\\_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...](http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...) 1/7

# Group2: intermediate peak score (28851 peaks)

2/9/2015

CENTDIST:Koeffler\_BM\_CebpE\_GMM\_ModelAssignment\_log\_300\_group2\_compSorted3.bed

Results for Koeffler\_BM\_CebpE\_GMM\_ModelAssignment\_log\_300\_group2\_compSorted3.bed  
VERSION: 2011.07.08

[Try our De Novo Motif Finding Tool for ChIP-seq \(SEMF\)](#)

746 TFs

Show top 50

Factors

Go

Download As Text

Rank <a href="#">[1]</a>	Name <a href="#">[1]</a>	Family <a href="#">[1]</a>	Logo <a href="#">[1]</a>	Score <a href="#">[1]</a>	Distribution <a href="#">[1]</a>		%Sequence with motif optimal setting	%Sequence with motif 1e-4 fold within +/- 200bp	Binding Range <a href="#">[1]</a>	PWM Score Cutoff <a href="#">[1]</a>	Z0Score <a href="#">[1]</a>	Z1Score <a href="#">[1]</a>	P-value <a href="#">[1]</a>
1	V\$CEBPB_Q2	<a href="#">CEBP</a>		36.4541			0.1615542	0.1398565	480	2.9101	33.7596	2.69456	0
2	V\$CEBP_Q2_Q1	<a href="#">CEBP</a>		30.0046			0.1442931	0.1265121	480	3.1246	27.234	2.7706	0
3	V\$jaque_CEBPA	<a href="#">jaque:Leucine Zipper</a>		29.099			0.09510935	0.0828342	480	2.9262	26.4199	2.67911	0
4	V\$CEBP_Q2	<a href="#">CEBP</a>		27.1049			0.1106028	0.09573325	480	2.9207	23.2332	3.87169	0
5	V\$CEBPA_Q1	<a href="#">CEBP</a>		27.0551			0.1544141	0.135108	480	2.8306	24.272	2.78314	0
6	V\$ETS_Q4	<a href="#">ETS</a>		26.0123			0.2105993	0.1824547	480	3.3301	22.8659	3.14638	0

[http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif\\_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...](http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...) 1/7

# Group3: high peak score (5741 peaks)

2/9/2015

CENTDIST:Koeffler\_BM\_CebpE\_GMM\_ModelAssignment\_log\_300\_group3\_compSorted3.bed

Results for Koeffler\_BM\_CebpE\_GMM\_ModelAssignment\_log\_300\_group3\_compSorted3.bed  
VERSION: 2011.07.08

[Try our De Novo Motif Finding Tool for ChIP-seq \(SEME\)](#)

746 TFs

Show top 50

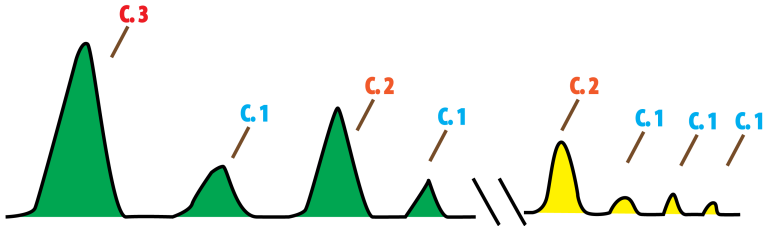
Rank	Name	Family	Logo	Score	Distribution		%Sequence with motif optimal setting	%Sequence with motif 1e-4 fold within +/- 200bp	Binding Range	PW/M Score Cutoff	Z0Score	Z1Score	P- value
1	VSCEBPB_02	<a href="#">CEBP</a>		32.9624					320	2.9101	29.8398	3.12262	0
2	VSCEBP_Q2_01	<a href="#">CEBP</a>		28.6415					360	3.1246	24.4458	4.19579	0
3	VSPEA3_Q6	<a href="#">ETS</a>		28.3666					440	2.8742	21.9021	6.46444	0
4	VSjanus_CEBPA	<a href="#">janus_Louise_Zipper</a>		27.798					360	2.9262	24.5191	3.2789	0
5	VSCEBPB_01	<a href="#">CEBP</a>		27.6113					360	3.1659	23.722	3.88938	0
6	VSCEBPA_01	<a href="#">CEBP</a>		26.0984					360	2.8312	21.6892	4.40913	0

[http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif\\_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...](http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...) 1/7

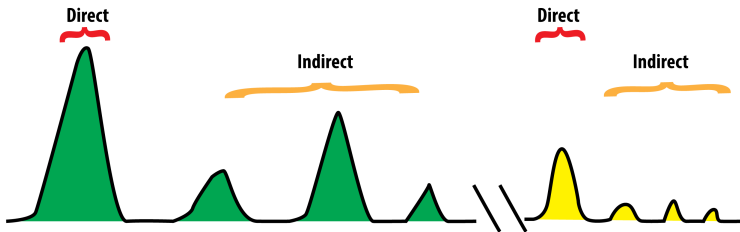
- Cebp motif is found in group 2 and 3 in 3-component GMMS using centdist

- Cebp motif is found in group 2 and 3 in 3-component GMMS using centdist
- Next, can we further segregate these groups into direct and indirect bindings?

## 3 Component-Mixture Model



## Local Clustering



# Direct: 24948 peaks

2/9/2015

CENTDIST:Koeffler\_BM\_CebpE\_GMM\_BiclusterAssignment\_SinglePeakFilteredOut\_log\_300\_compSorted3\_dist3kb\_direct.bed

Results for Koeffler\_BM\_CebpE\_GMM\_BiclusterAssignment\_SinglePeakFilteredOut\_log\_300\_compSorted3\_dist3kb\_direct.bed  
VERSION: 2011.07.08

[Try our De Novo Motif Finding Tool for ChIP-seq \(SEME\)](#)

746 TFs

Show top 50

Factors

Go

Download As Text

Rank <a href="#">[1]</a>	Name <a href="#">[2]</a>	Family <a href="#">[3]</a>	Logo <a href="#">[4]</a>	Score <a href="#">[5]</a>	Distribution <a href="#">[6]</a>		%Sequence with motif optimal setting <a href="#">[7]</a>	%Sequence with motif 1e-4 fdr within +/- 200bp <a href="#">[8]</a>	Binding Range <a href="#">[9]</a>	PWM Score Cutoff <a href="#">[10]</a>	Z0Score <a href="#">[11]</a>	Z1Score <a href="#">[12]</a>	P-value <a href="#">[13]</a>
1	VSCEBPB_02	<a href="#">CEBP</a>		58.3326			 0.1833414	 0.1734007	440	2.9101	47.6979	10.6347	0
2	V\$jaspar_CEBPA	<a href="#">jaspar: Leucine Zipper</a>		46.312			 0.1111111	 0.1045775	440	2.9262	39.0362	7.27579	0
3	VSPEA3_Q6	<a href="#">ETS</a>		41.8737			 0.2650313	 0.2477152	440	2.8742	33.5307	8.34297	0
4	VSCEBP_Q2_01	<a href="#">CEBP</a>		41.8022			 0.1406526	 0.1499519	360	3.1246	39.0291	2.77308	0
5	V\$jaspar_SPI1	<a href="#">jaspar: Fos</a>		40.3973			 0.2049864	 0.1925204	440	3.5842	32.6871	7.71024	0
6	VSCEBPB_01	<a href="#">CEBP</a>		40.0647			 0.1501924	 0.1604938	360	3.1658	36.5447	3.52005	0

[http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif\\_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handie=guest.172.16.227...](http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handie=guest.172.16.227...) 1/7

# Indirect: 26547 peaks

2/9/2015

CENTDIST:Koeffler\_BM\_CelpE\_GMM\_BiclusterAssignment\_SinglePeakFilteredOut\_log\_300\_compSorted3\_dist3kb\_indirect.bed

Results for Koeffler\_BM\_CelpE\_GMM\_BiclusterAssignment\_SinglePeakFilteredOut\_log\_300\_compSorted3\_dist3kb\_indirect.bed  
VERSION: 2011.07.08

[Try our De Novo Motif Finding Tool for ChIP-seq \(SEME\)](#)

746 TFs

Show top 50

Factors

Go

Download As Text

Rank	Name	Family	Logo	Score	Distribution	%Sequence with motif optimal setting	%Sequence with motif 1e-4 fold within +/- 200bp	Binding Range	PWM Score Cutoff	ZScore	Z1Score	P-v
1	V\$jaspar_NFATC2	<a href="#">jaspar_Rel</a>		9.86315				320	3.4936	3.21144	6.65171	0.00
2	V\$FOXO3_01	<a href="#">FOX</a>		9.20345				120	3.1121	2.05819	7.14525	0.00
3	V\$HNF1_Q6	<a href="#">HNF1</a>		8.25904				200	3.1667	2.46373	5.7953	0.00
4	V\$SRV_01	<a href="#">FOX</a>		8.24092				160	2.7795	0.858234	7.38260	0.00
5	V\$PAX4_04	<a href="#">PAX</a>		7.70944				280	3.0252	2.41782	5.29162	0.00
6	V\$FOXPI_01	<a href="#">FOX</a>		7.70615				160	2.2301	1.69717	6.00898	0.00

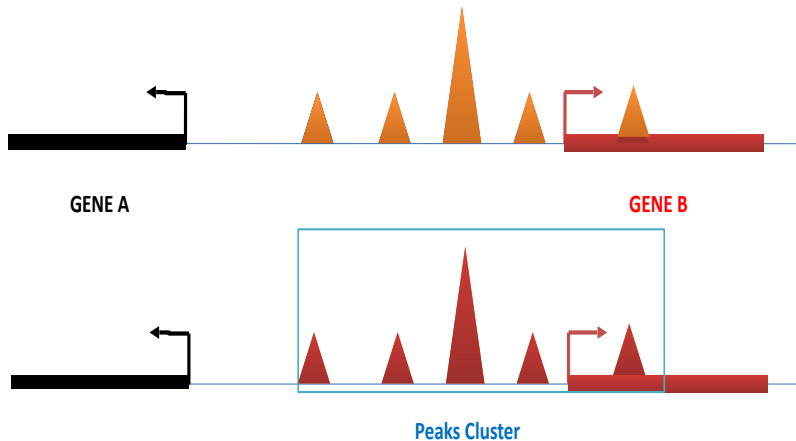
[http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif\\_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...](http://biogpu.ddns.comp.nus.edu.sg/~chipseq/webseqtools2/TASKS/Motif_Enrichment/view.php?top=50&show=factor&submit=Go&email=guest.172.16.227.227&handle=guest.172.16.227...) 1/7



- Our current method could separate direct and indirect bindings

- Our current method could separate direct and indirect bindings
- Next, can we further using peak clusters increase functional annotation?

# Find the targeted genes



*What problems the invention solves and advantages over existing methods?*

*An Example:*

