# Identify functional patterns in high throughput binding assays
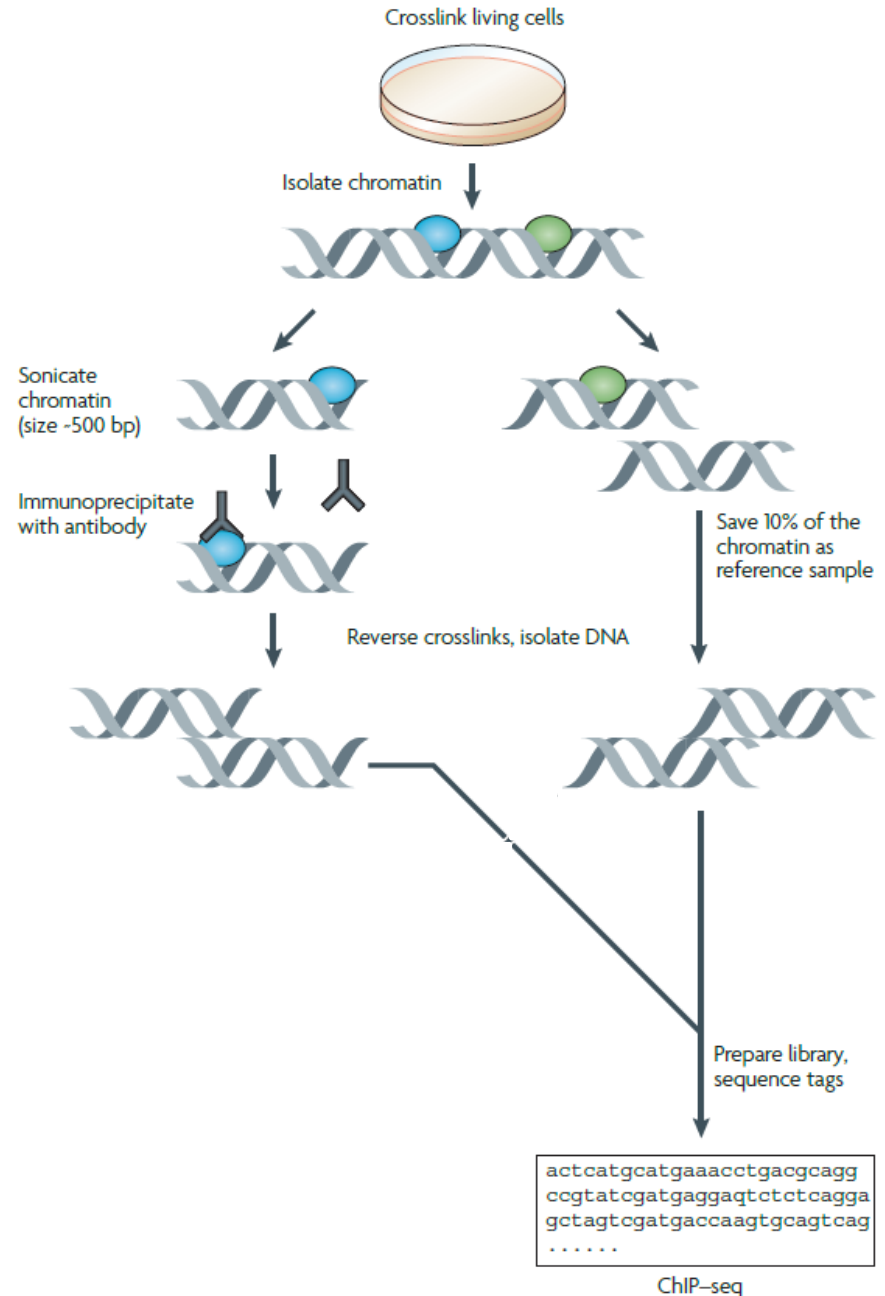
Alex Essebier

# Message

- By clustering ChIP-seq peaks we can identify different patterns in transcription factor binding
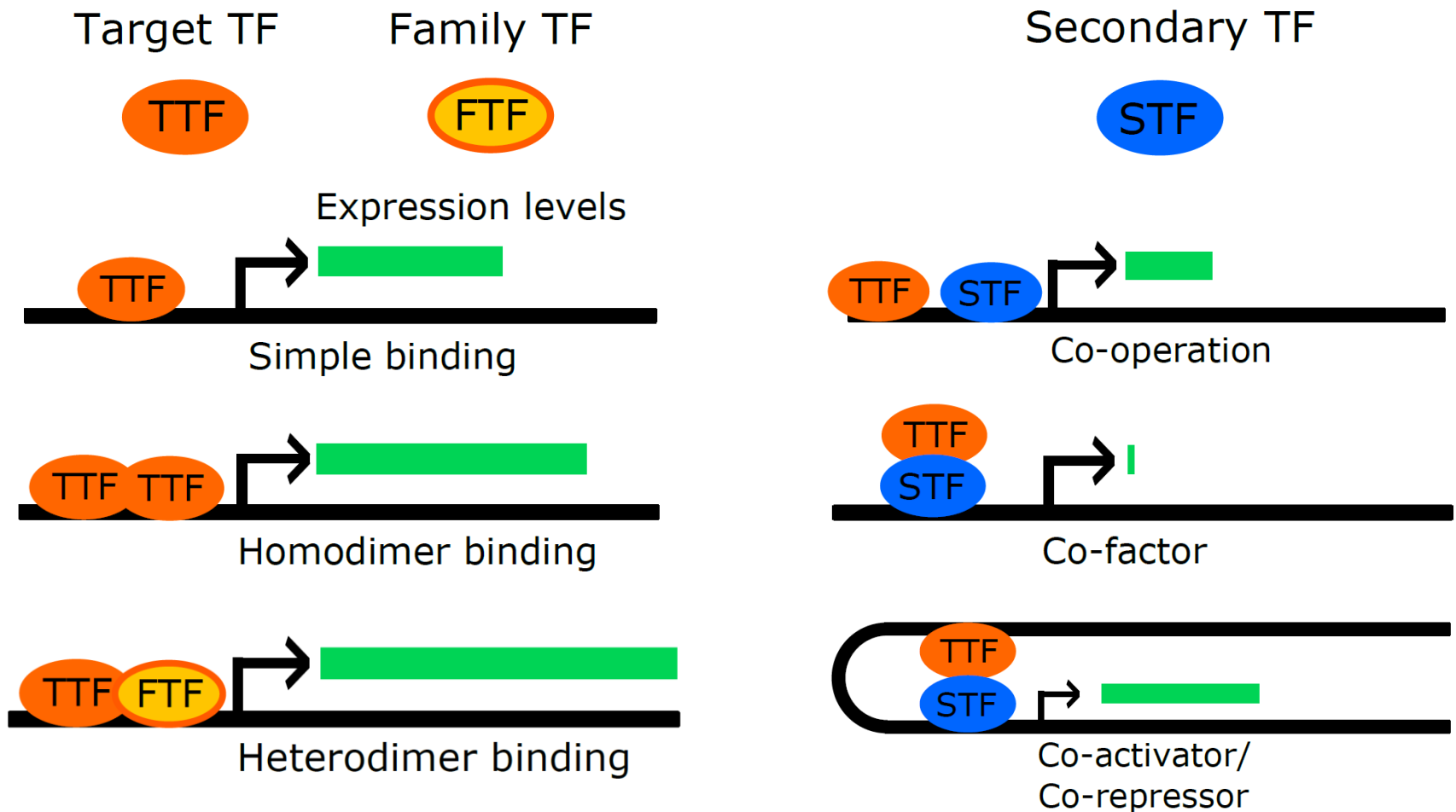
# ChIP-seq experiment

- Chromatin immunoprecipitation followed by sequencing
- To determine where a protein binds the genome
- e.g. for a **single TF** or histone modification



Crosslink living cells

Isolate chromatin

Sonicate chromatin (size ~500 bp)

Immunoprecipitate with antibody

Reverse crosslinks, isolate DNA

Save 10% of the chromatin as reference sample

Prepare library, sequence tags

```
actcatgcatgaaacctgacgcagg
ccgtatcgatgaggaqtctctcagga
gctagtcgatgaccaagtgcagtcag
......
```
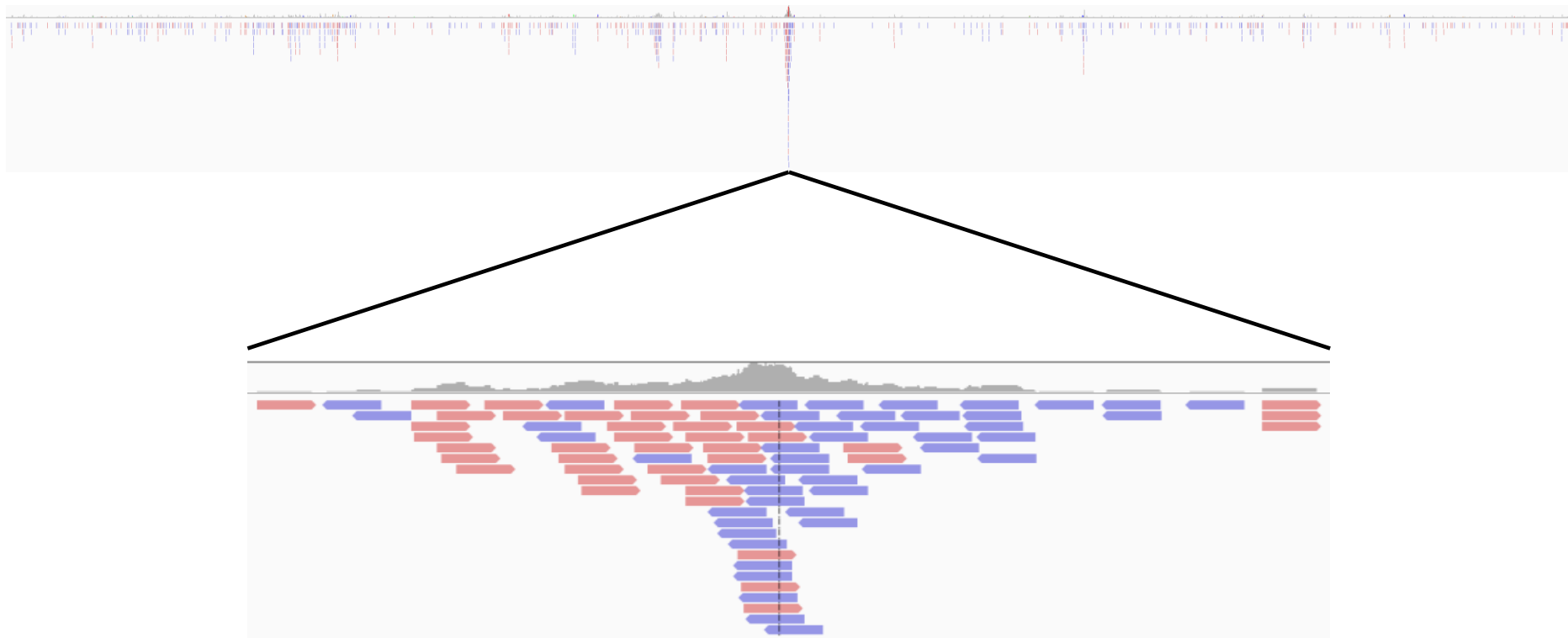
ChIP-seq

# Why TFs?

- Important role in gene expression, cell differentiation and homeostasis
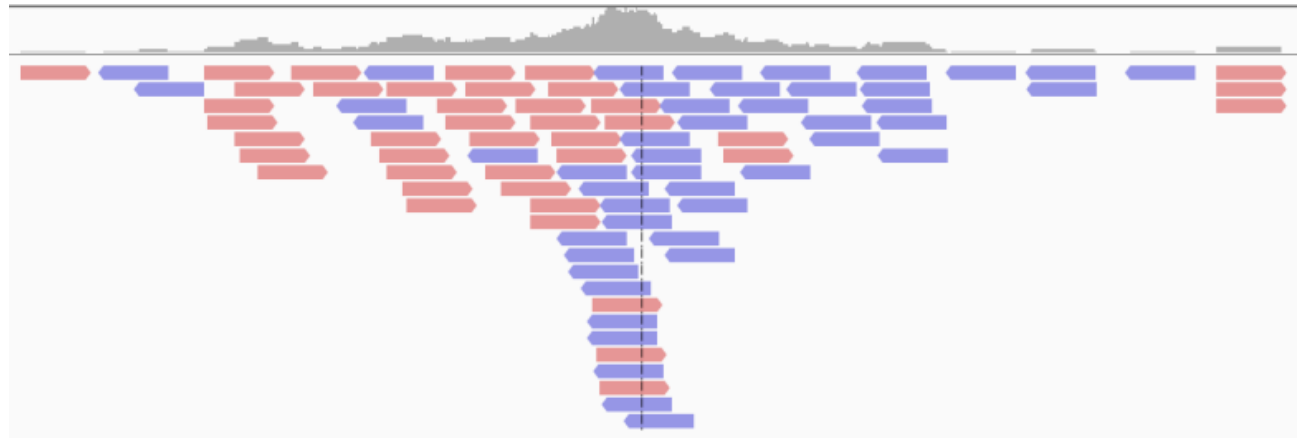
# Peak calling

- Raw sequencing data
  - Single end reads
  - Red mapped 'forward', blue mapped 'reverse'
  - Distribution across genome

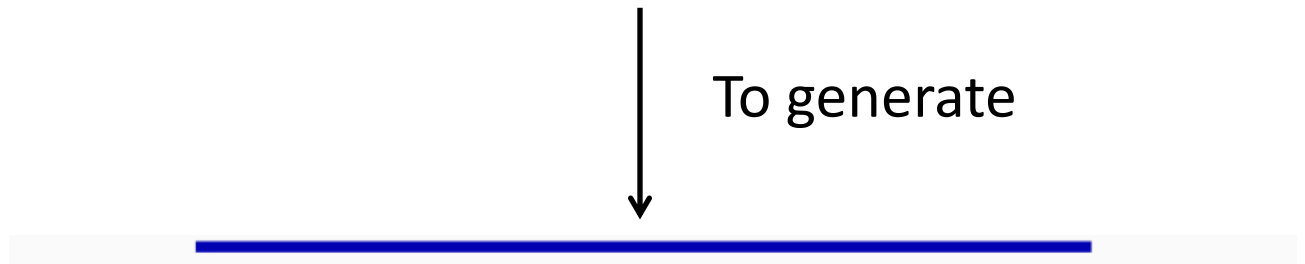# Peak calling

**Sample** – exposed to antibody



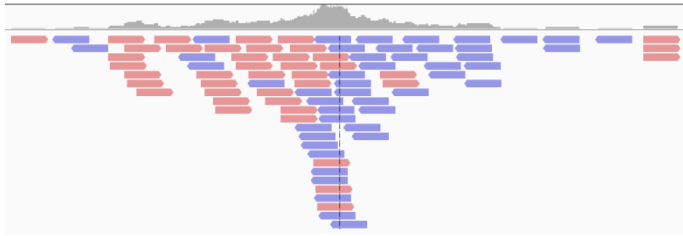Compared to

**Input** – no antibody exposure

To generate

**Peak** – with statistical significance
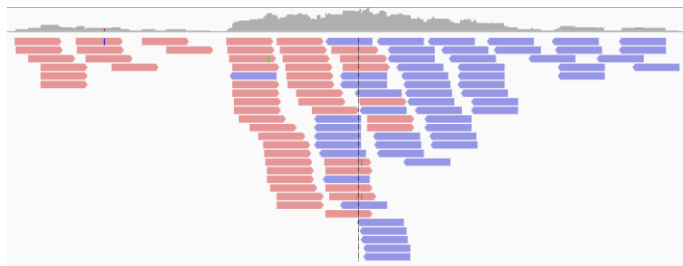
# Peak calling

A



ATTGCC

B



ATTTCC

C



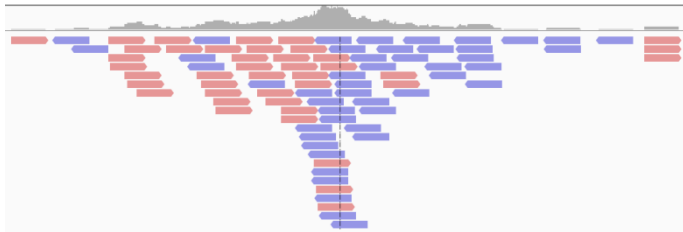ATACCC

- Peaks have different features within a ChIP-seq experiment

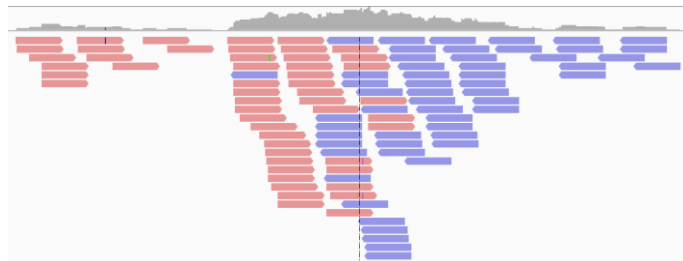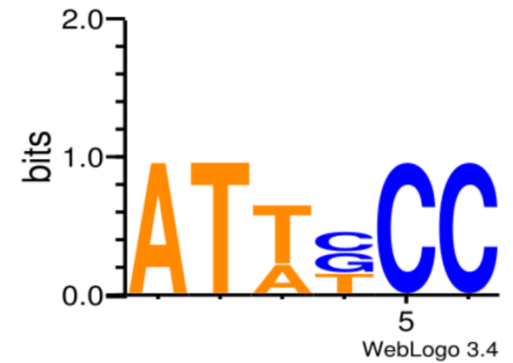| Peak | Width | Enrich-ment | Location | Epigenetic Environment |
|------|-------|-------------|----------|------------------------|
| A | 205 | 298 | Distal Intergenic | Insulator |
| B | 162 | 218 | Promoter | Active Promoter |
| C | 194 | 361 | Promoter | Weak Promoter |

# Peak calling



ATTGCC

ATTTCC

ATACCC

Analyse

Confirm *in vitro* results

Identify consensus motif

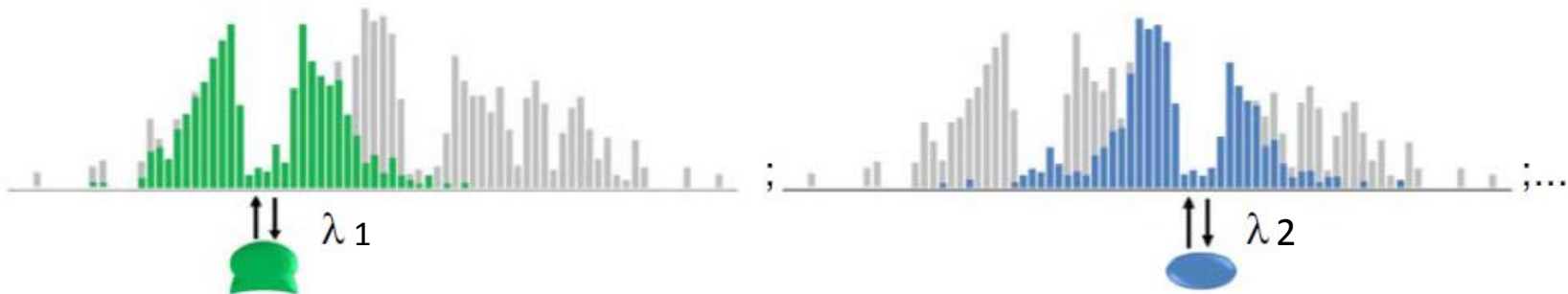Identify target genes of TF

# Hypothesis

- We propose that ChIP-seq peaks from a TF experiment can be clustered based on their read density or 'shape' leading to identification of different binding modes and functional patterns of a TF

# Previous use of peak shape

- Differential binding
  - Compare two conditions
  - Compare two TFs
  - Based on read depth

- TF binding from DNase I hypersensitivity



TF binding estimation from modelled DNase I hypersensitivity profiles

Adapted from: Sherwood, R.I., et al. *Nature Biotechnology* **32(2),**171-178 (2014)
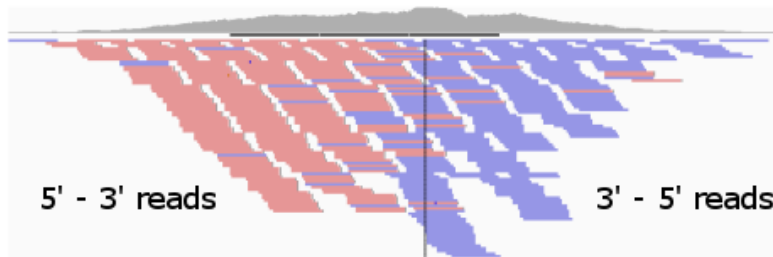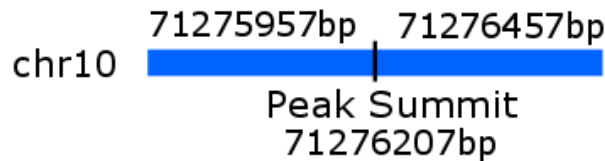
# Aims

- Develop a modelling technique to identify functionally relevant clusters, based on ChIP-seq read density, defining TF binding events

- Identify functional patterns associated with clusters and provide more information about TF binding from ChIP-seq data
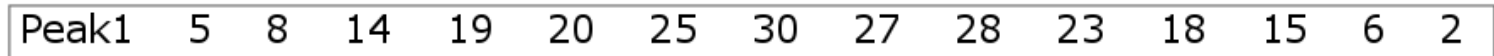
# Processing peak data

## Binding Site aka Peak



5' - 3' reads          3' - 5' reads

Read pile up showing distribution of reads in peak

71275957bp          71276457bp

chr10

Peak Summit
71276207bp

500bp window around summit

chr10

Split window into even segments
Count read depth in each segment

| Peak1 | 5 | 8 | 14 | 19 | 20 | 25 | 30 | 27 | 28 | 23 | 18 | 15 | 6 | 2 |

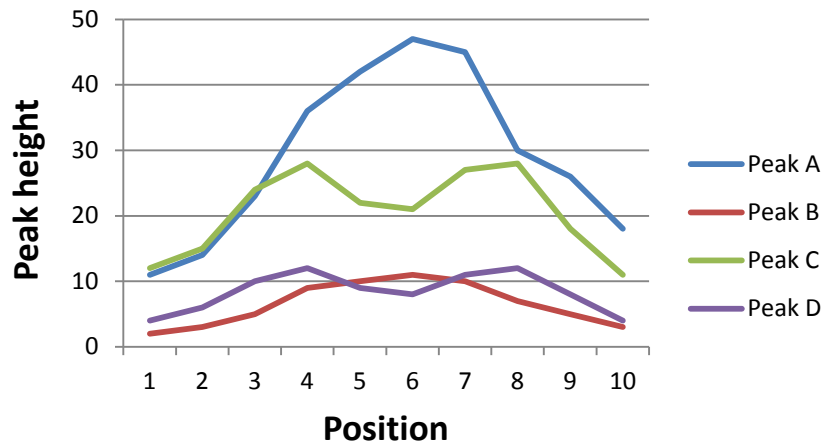This creates a count vector for each peak with equal numbers of columns

Combining all count vectors creates a **Dirichlet distribution** that can be clustered
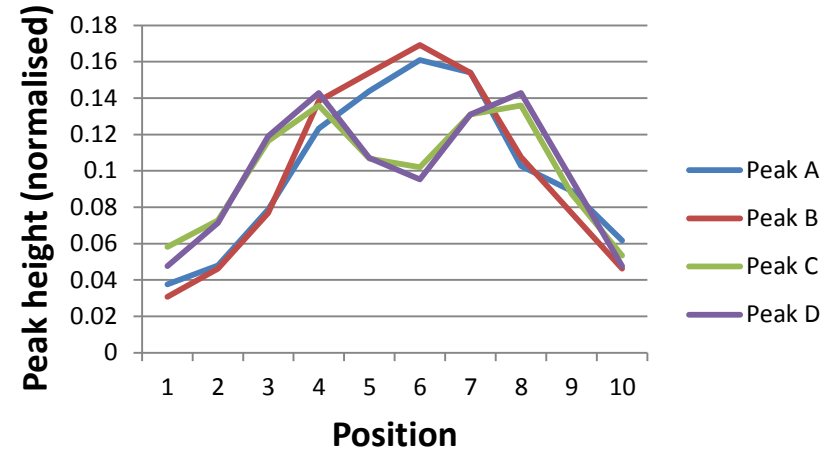
# Dirichlet clustering

- Dirichlet distribution – distribution of distributions
- The model is a Dirichlet mixture
- Unsupervised clustering of peaks
- Evidence based clustering using raw counts
- **No normalisation** of data

# Evidence based clustering

### Shape of read counts



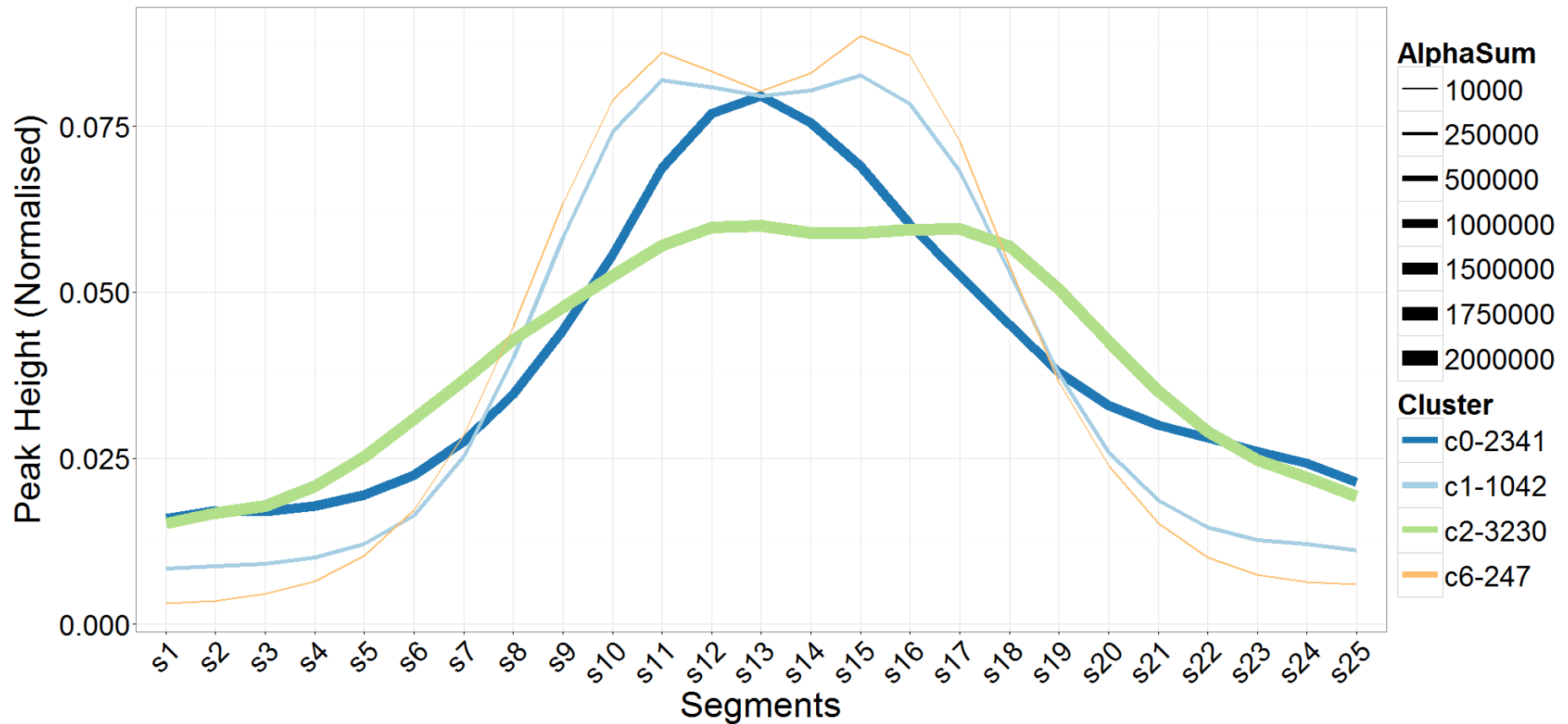### Shape of normalised read counts



| K-means | |
|---|---|
| Peak | Cluster |
| A | 1 |
| B | 1 |
| C | 2 |
| D | 2 |

| Dirichlet | |
|---|---|
| Peak | Cluster |
| A | 1 |
| B | 2 |
| C | 3 |
| D | 4 |

- Read depth is key and can be masked by normalisation
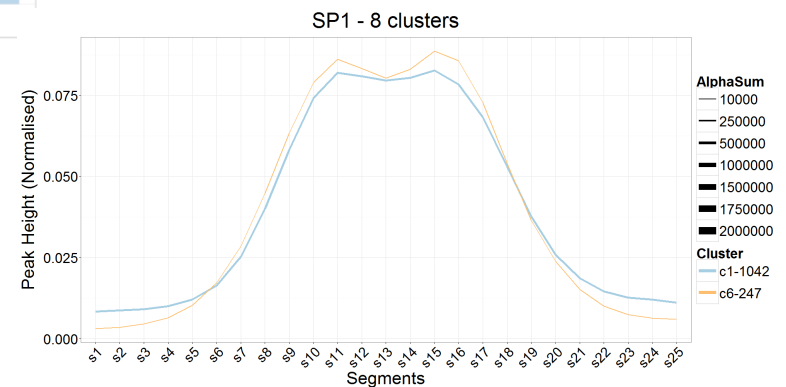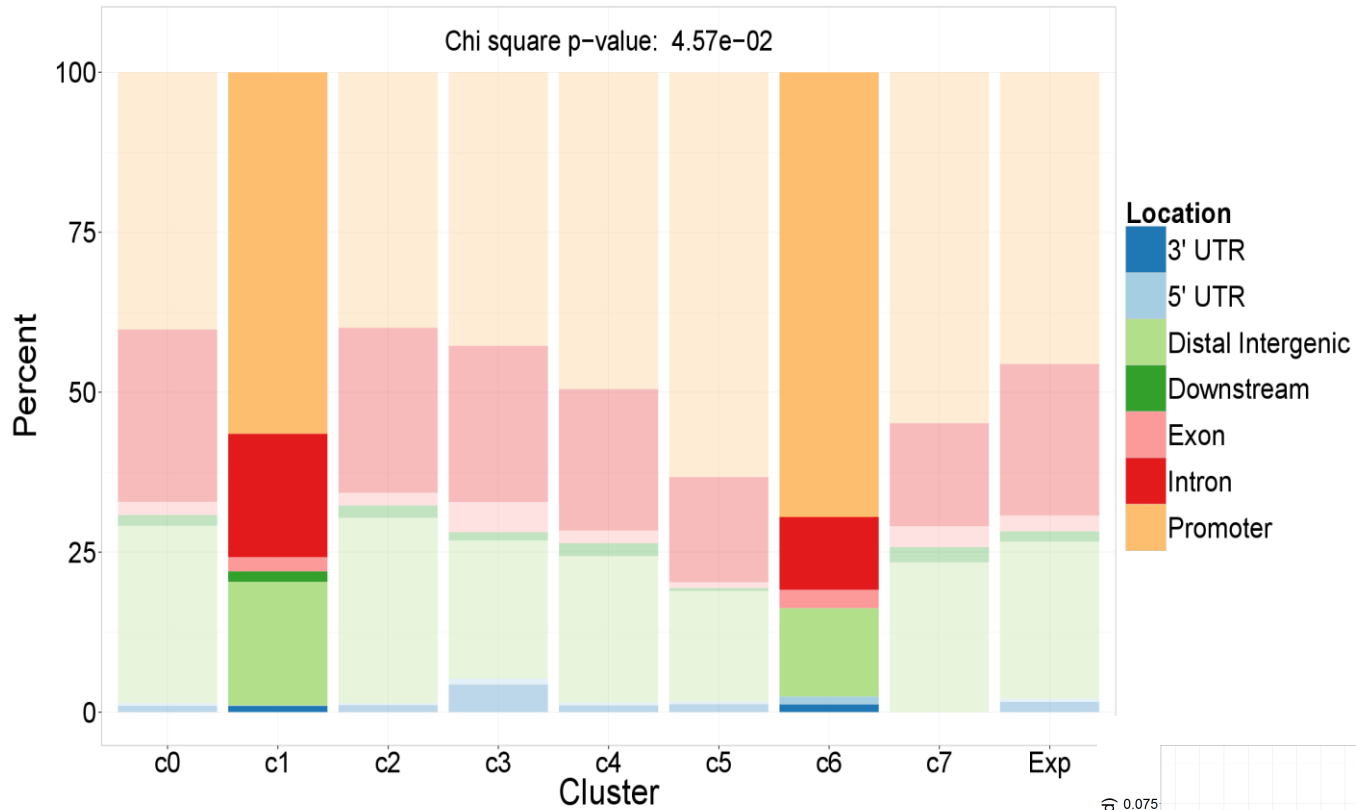- Dirichlet approach does not require normalisation

# Clustering example – SP1



SP1 - 8 clusters

# Genomic location



Genomic Locations of Binding Sites
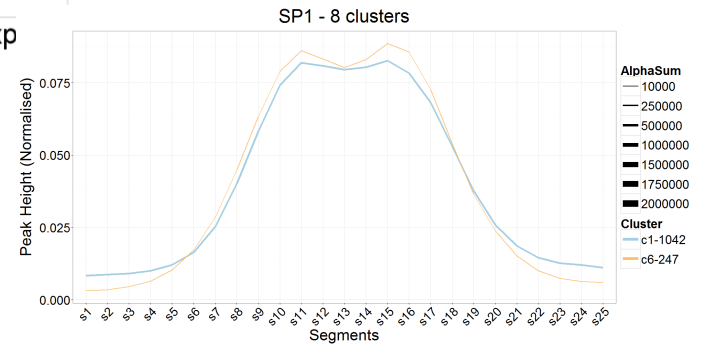
Chi square p-value: 4.57e-02

SP1 - 8 clusters

# Epigenetic environment



Chromatin State Across Binding Sites
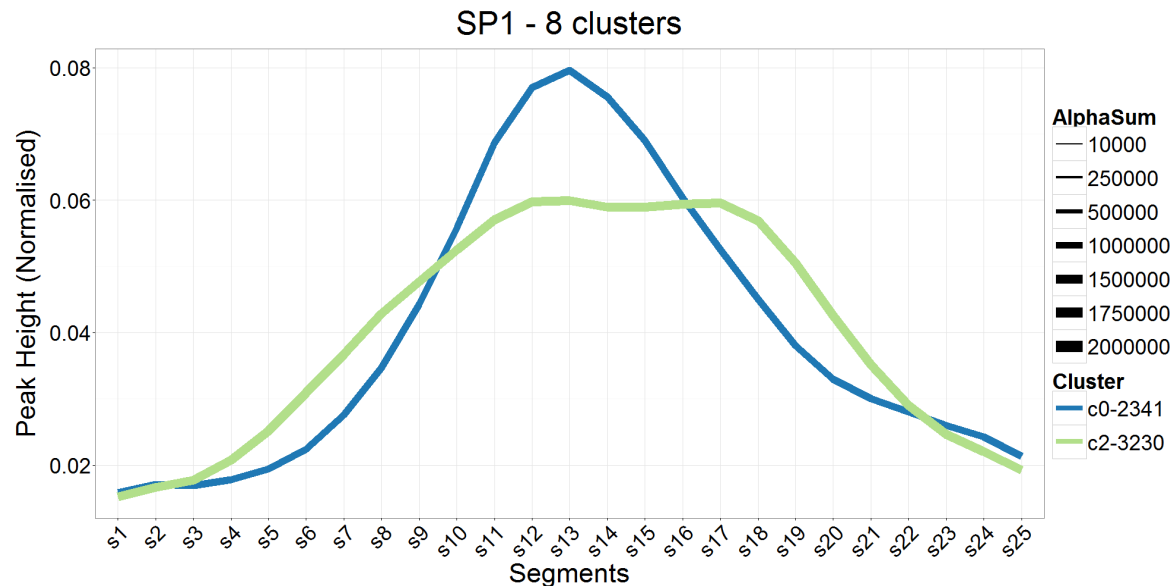
# Consensus motifs

Cluster    Motif

0

2

- SP1 motif

- Differentiating feature in c0 and c2 is binding affinity or read depth
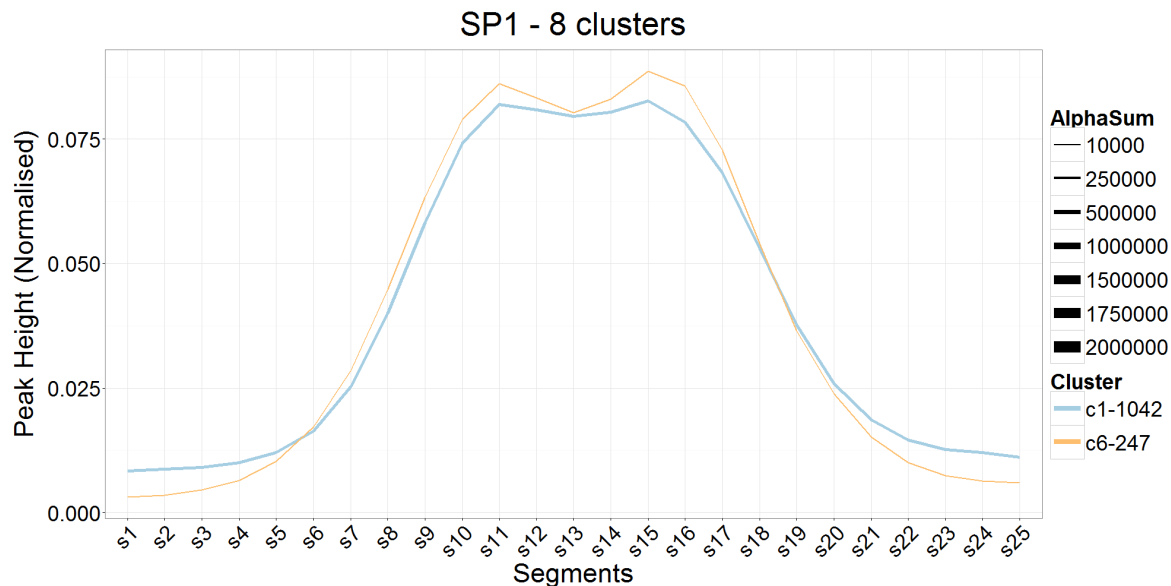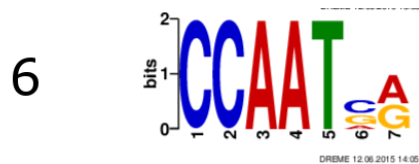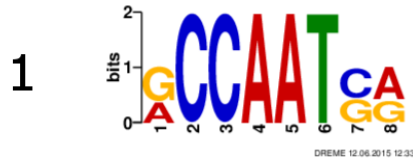
# Consensus motifs

Cluster     Motif

1

6

- NFY motif

- Known interaction between two TFs

- A bimodal peak shape indicates increased NFY binding

# Applications

- Explore TF families by comparing clustering outcomes
- Explore TF dimers using clustering in combination with *in vitro* sequence data
- Explore cooperative interactions

# Summary

- We successfully clustered ChIP-seq peaks based on their shape, density and magnitude then demonstrated how each cluster contains unique, biologically relevant, features

# Thanks

**Supervisor**

Mikael Bodén

**Bodén Group**

Ralph Patrick

Tim O'Connor

Julian Zaugg

Gabe Foley

**Piper Group**

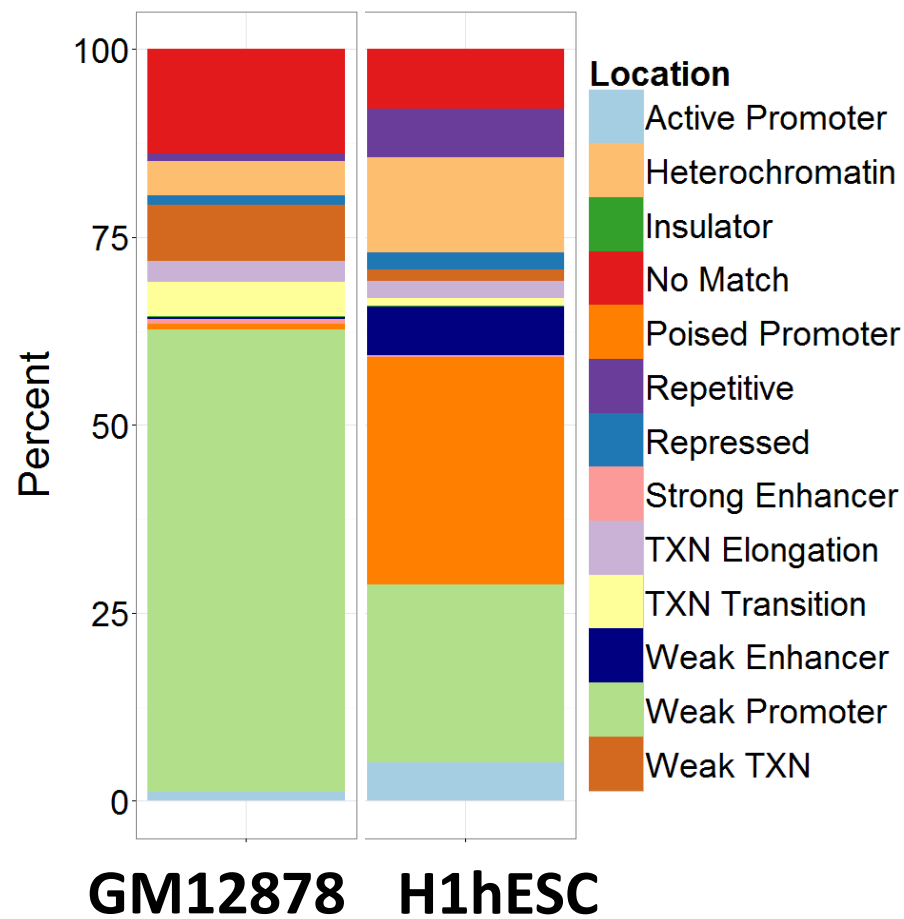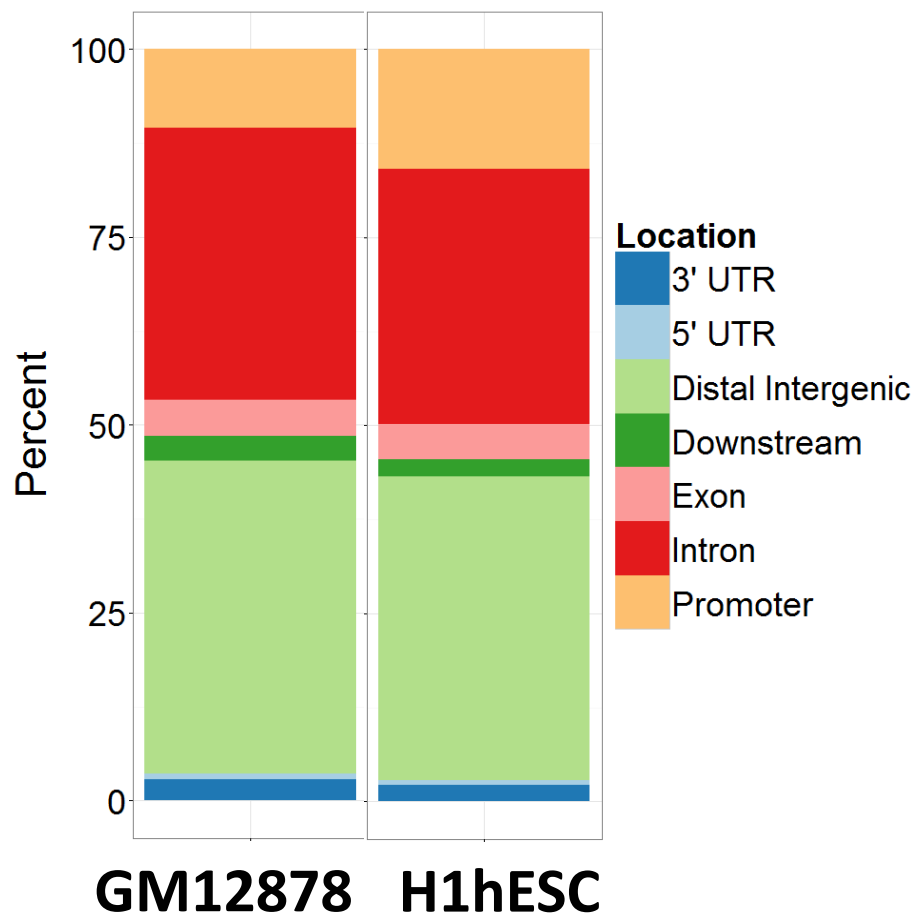Michael Piper

**Rostlab (TUM)**

Burkhard Rost
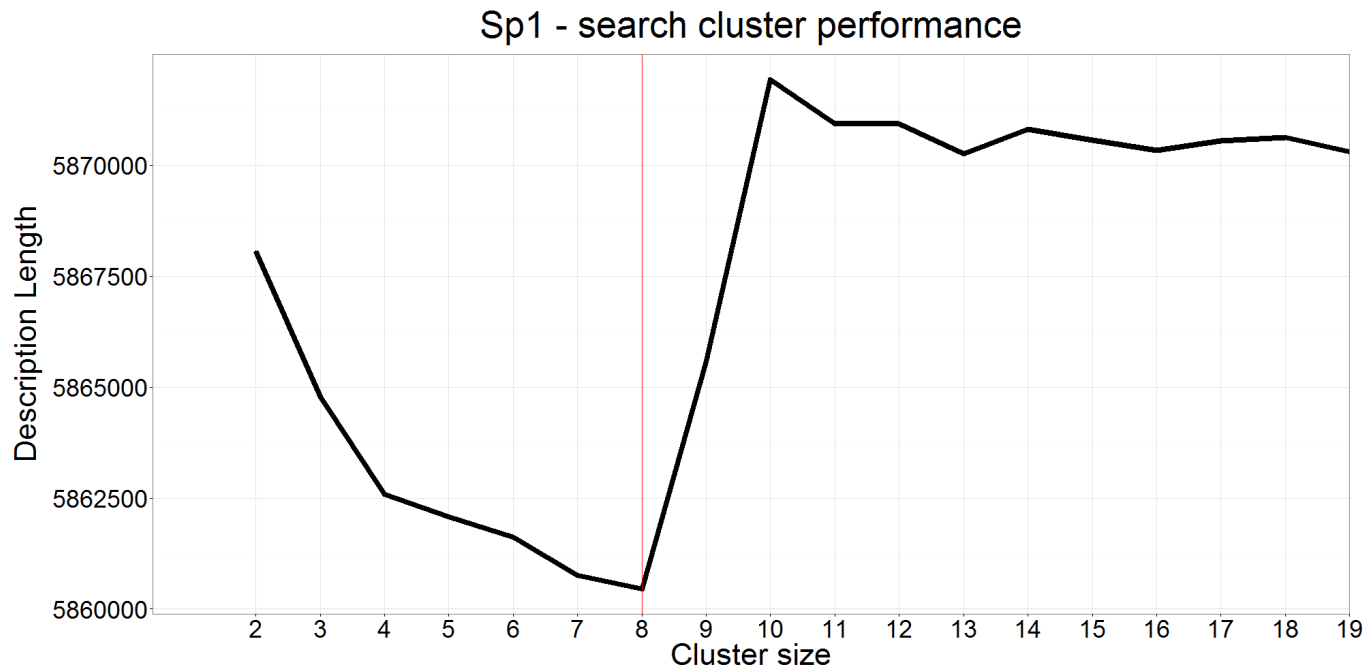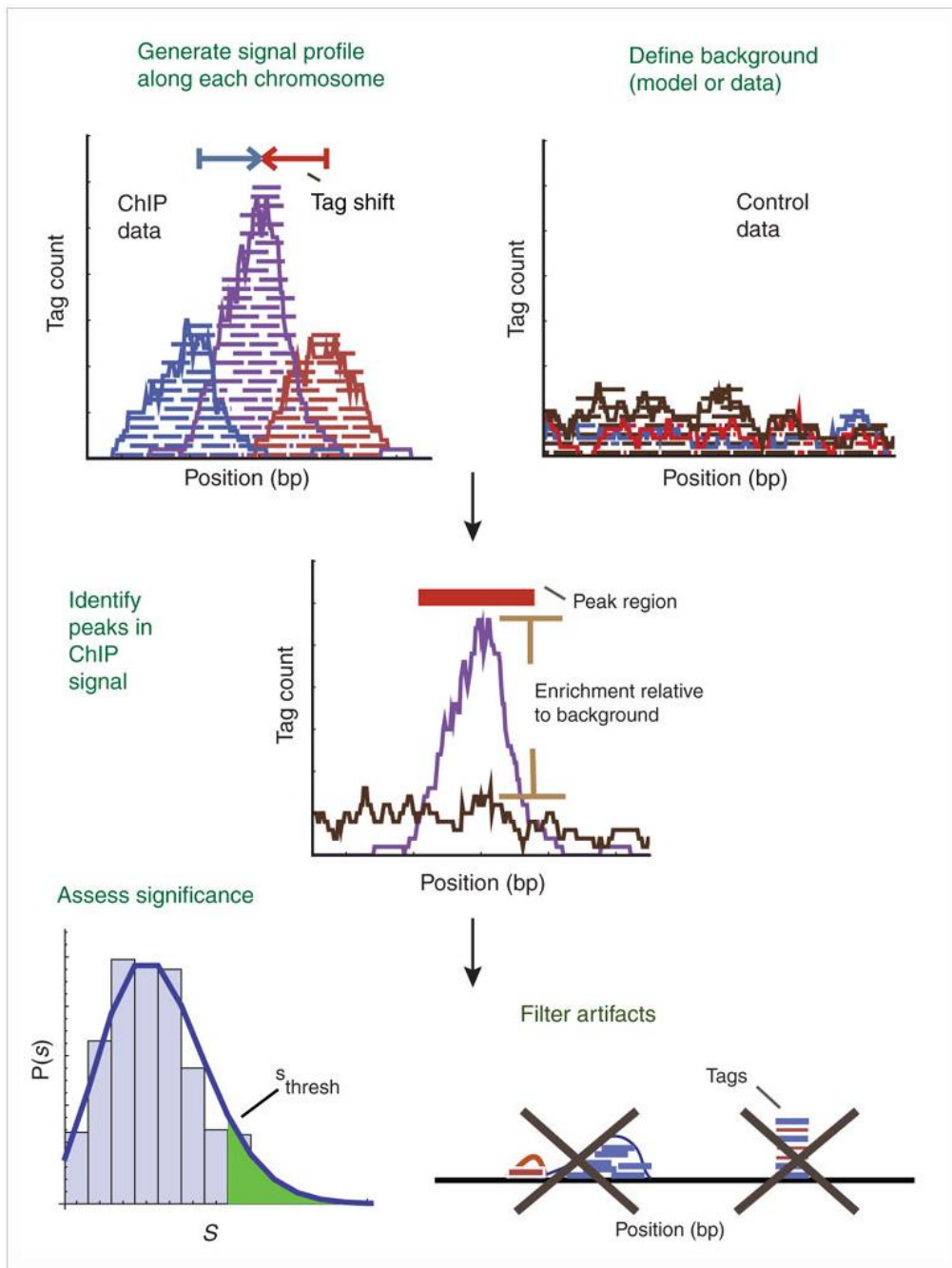
Tatyana Goldberg

# Cell type

# Cluster size optimisation

- Minimum description length (MDL)

- Description length (DL)

  - A measure of information content and model complexity

- Larger models will always be more complex



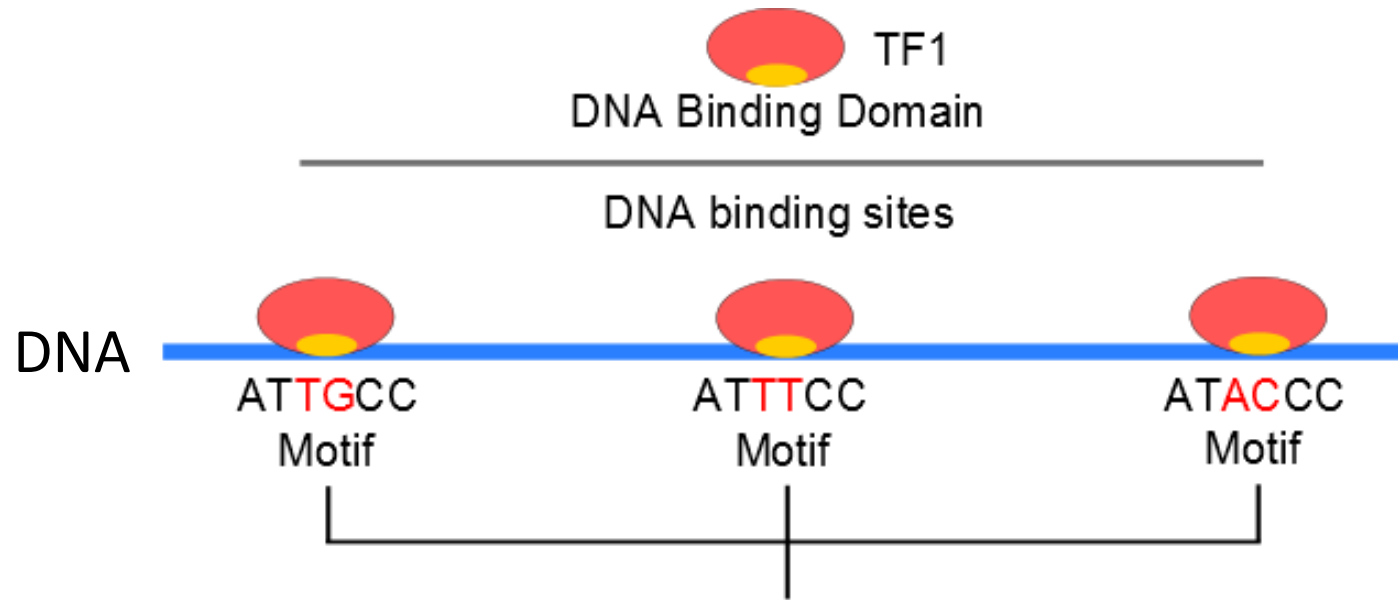Sp1 - search cluster performance

# ChIP-seq peak calling

- Shift reads on both strands to find peak
- Compare to control reads
- Identify significant hits according to a threshold
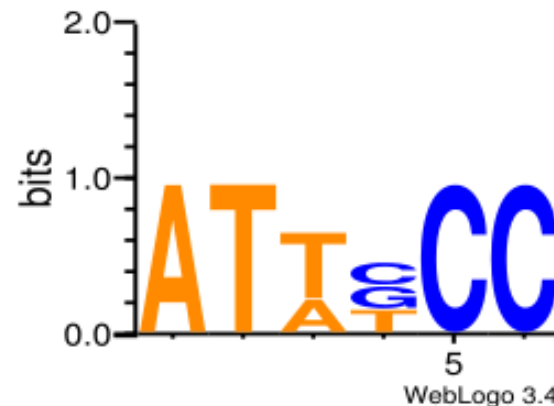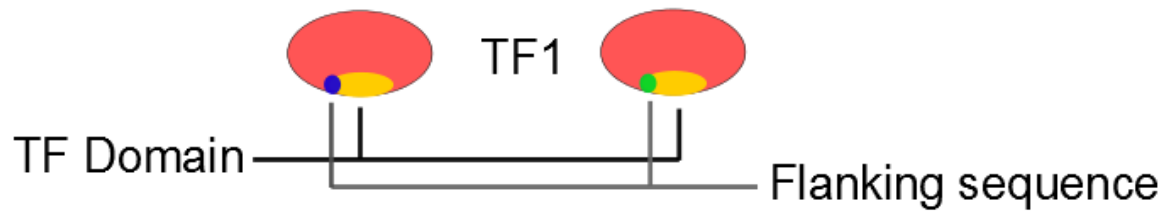- Remove potential artefacts

# DNA binding sites

Variation in flanking sequence of TF domains

TF1

TF Domain

Flanking sequence

Sequence polymorphism in target DNA

ATTGCC          ATTTCC          ATACCC

Site accessibility within chromatin landscape

Euchromatin
('open')

Heterochromatin
('closed')

ATTTCC          ATTTCC

# Minimum description length (MDL) Principle

- Calculate model complexity

- Calculate smallest data description length (DDL)

- Total DL = sum of complexity and DDL

- Plot total DL as number of clusters increases and search for global minima

- Global minima = optimal number of clusters