

# FunSeq3: design, plan and results

Lou Shaoke

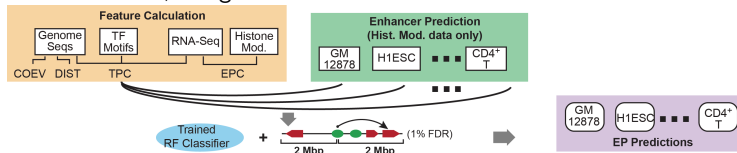
Department of Molecular Biophysics and Biochemistry

*[loushaoke@gmail.com](mailto:loushaoke@gmail.com)*

July 29, 2015

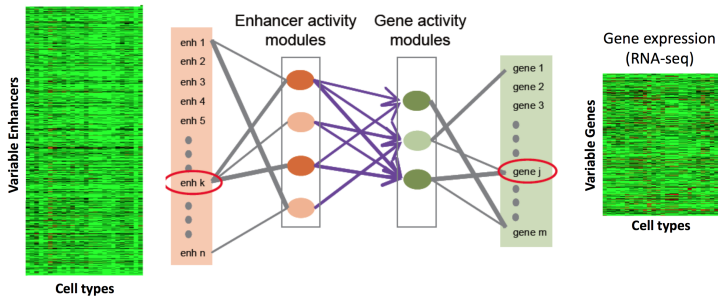
Yale

- Linear correlation between enhancer marks and gene expression
- Random Forest, using ChIA-PET data



- mixed-membership problem

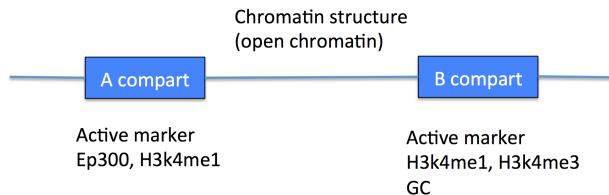
Enhancer activity  
(H3K4me1/H3K27me3)



- why it forms 3D struct?
- All distal interaction need a 3D structure
- how to form 3D structure? Two factors:
  - capability of chromatin structure, not too loose and too tight. or special chromatin structure?
  - mediator and cofactors.

- ① enhancer-gene linkage doesnot randomly happened. The set of linkages is a collection of linkages from all cell lines/tissues.
- ② tissue-specific effect follow the same mechanism to form 3D struct: some tissue don't have the same linkage because the changes of associated factors.
- ③ the linkage can have positive or negative effect (activation or repression)

How to find the associated features?



A compartment: The leftmost interaction region

B compartment: The rightmost interaction region

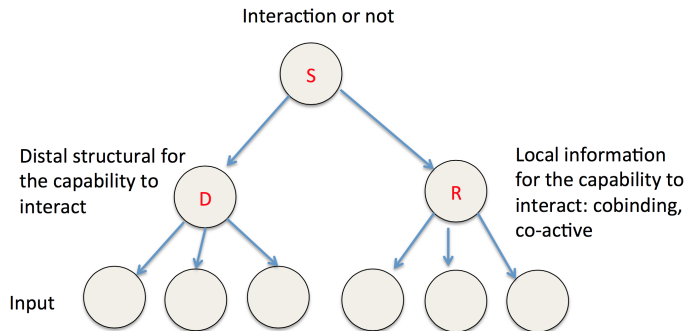
Distal(determine whether can fold): structural information (DNase signal)

Local (determine whether can co-bind by mediators/factors): H3k4me1, GC contents,

Co-evolution(sequence and synteny), TF motif, footprint data?

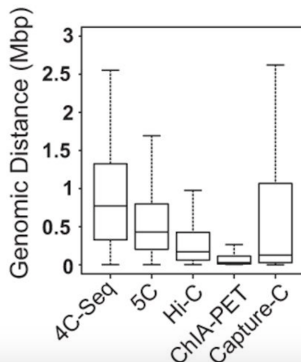
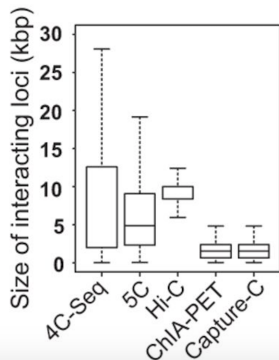
Gene expression

(Not a real model)



Interaction map:

4D map: <http://4dgenome.int-med.uiowa.edu/Download.html>



The number of interactions for HiC, ChIA-PET and 5C

GM12878	H1ESC	HELA	IMR90	K562	MCF7
1177	1221	1849	1114640	66516	64487

Use all predicted enhancer-gene pairs from Wang et al.

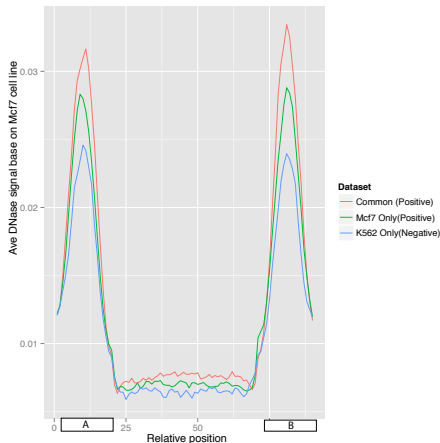
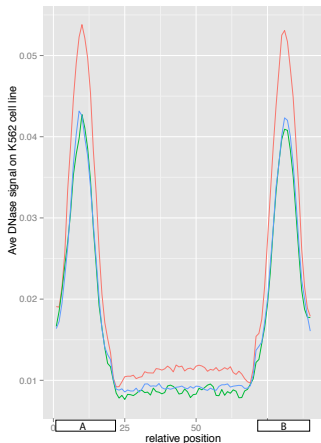
ChIA-PET data for: MCF7 K562

## Define the positive and negative set

- 1 Take the intersection of predicted enhancer-gene linkage and ChIA-PET interactions to define positive set
- 2 The linkage prediction algorithm should not consider 3d interaction: linear correlation or mix-membership relationship
- 3 Use ChIA-PET for Mcf7 and K562
- 4 Define datasets:  
Mcf.K562.com(Positive set) : Interaction overlaps between K562 and Mcf7  
Mcf or K562 only: Only have interaction in one cell cell, as Positive set for itself, Negative set for the other.  
For example: Mcf7-only interaction is positive for Mcf7 cell but to be treated as negative set for K562.



A,B are group into 20 equal-size bins; The intra-region between AB are split into 50 equal-size bins, and then compare the average signal for Positive , cell-specific interaction(Positive) and negative set. (K562 specific interactions are thought as Mcf7-Negative, vice versa)



The negative set might contain positive set because of the ChIA-PET sensitivity.

- extract other related features, especially for enhancer-gene local co-binding features
- build models: bayesian network model, and also try deep learning and other data mining algorithms

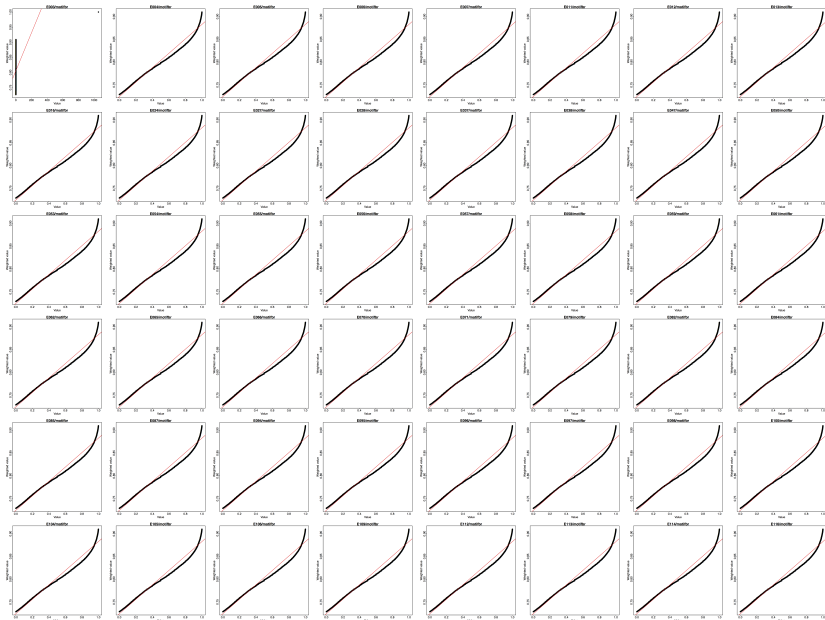
Any question or advices?

Tissue-specific enhancer-gene linkage:

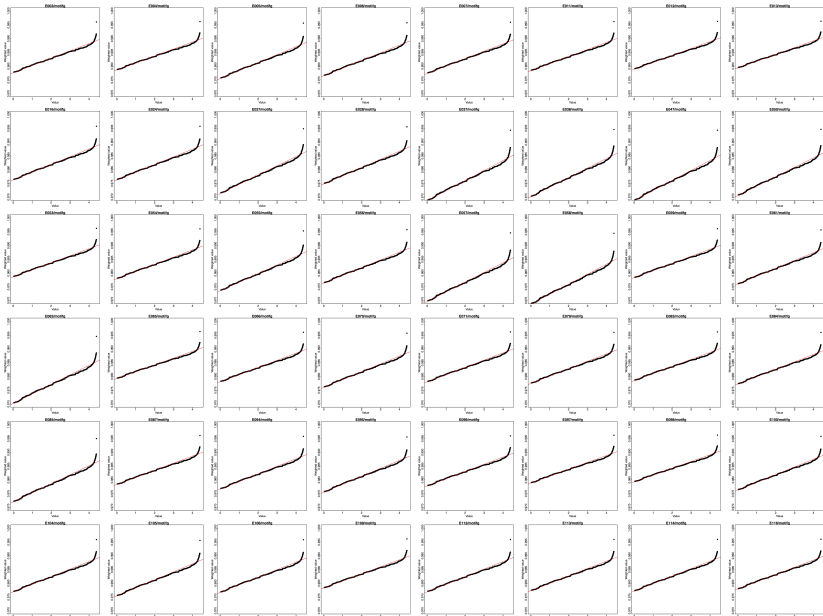
- Data from "Learning three-dimensional regulation of gene expression" (Jianrong Wang, Manolis et al), contains enhancer-gene linkage from 56 tissues.
- Variants Sampling from 1KG
- FunSeq 2.1.2



# Motif breaking score fitting



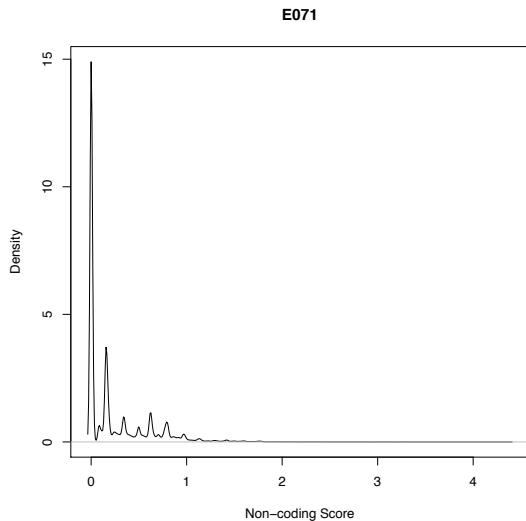
# Motif gain score fitting



## Case study: Glioblastoma

1. exome somatic mutation (public data)
2. test on three set of gene-enhancer linkages:  
Original enhancer-gene linkage;  
E070 Brain\_Germinal\_Matrix  
E071 Brain\_Hippocampus\_Middle  
Results from different linkages set are the same

# Distribution of Noncoding Score





# Enrichment analysis

Sublist	Category	Term	RT	Genes	Count	%	P-Value	Benjamini
<input type="checkbox"/>	UP_TISSUE	Epithelium	RT		43	35.5	2.9E-9	3.4E-7
<input checked="" type="checkbox"/>	GNF_U133A_QUARTILE	Cardiac Myocytes_3rd	RT		86	71.1	5.0E-8	3.9E-6
<input type="checkbox"/>	KEGG_PATHWAY	<a href="#">Pathways in cancer</a>	RT		18	14.9	6.7E-8	5.5E-6
<input checked="" type="checkbox"/>	GNF_U133A_QUARTILE	PB-CD14+Monocytes_3rd	RT		97	80.2	2.1E-7	8.3E-6
<input type="checkbox"/>	KEGG_PATHWAY	<a href="#">Small cell lung cancer</a>	RT		10	8.3	3.3E-7	1.3E-5
<input checked="" type="checkbox"/>	SMART	<a href="#">HLH</a>	RT		10	8.3	1.2E-6	8.3E-5
<input type="checkbox"/>	KEGG_PATHWAY	<a href="#">Chronic myeloid leukemia</a>	RT		9	7.4	1.6E-6	4.4E-5
<input checked="" type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">vascular_high-grade comedo DCIS endothelium_3rd</a>	RT		23	19.0	1.9E-6	6.7E-4
<input type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">stem cell null_3rd</a>	RT		22	18.2	5.1E-6	8.8E-4
<input checked="" type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">stem cell null_3rd</a>	RT		21	17.4	1.5E-5	1.7E-3
<input type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">brain_3rd</a>	RT		23	19.0	3.2E-5	2.8E-3
<input checked="" type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">kidney_normal epithelium_3rd</a>	RT		25	20.7	3.6E-5	2.5E-3
<input type="checkbox"/>	GNF_U133A_QUARTILE	leukemiapromyelocytic(hl60)_3rd	RT		43	35.5	3.9E-5	1.0E-3
<input checked="" type="checkbox"/>	GNF_U133A_QUARTILE	Whole Brain_3rd	RT		51	42.1	4.9E-5	9.5E-4
<input type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">stem cell null_3rd</a>	RT		22	18.2	5.0E-5	2.9E-3
<input checked="" type="checkbox"/>	KEGG_PATHWAY	<a href="#">Prostate cancer</a>	RT		8	6.6	5.8E-5	1.2E-3
<input type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">stem cell null_3rd</a>	RT		18	14.9	6.8E-5	3.4E-3
<input checked="" type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">stem cell null_3rd</a>	RT		20	16.5	6.9E-5	3.0E-3
<input type="checkbox"/>	CGAP_SAGE_QUARTILE	<a href="#">vascular_normal liver_3rd</a>	RT		21	17.4	7.0E-5	2.7E-3

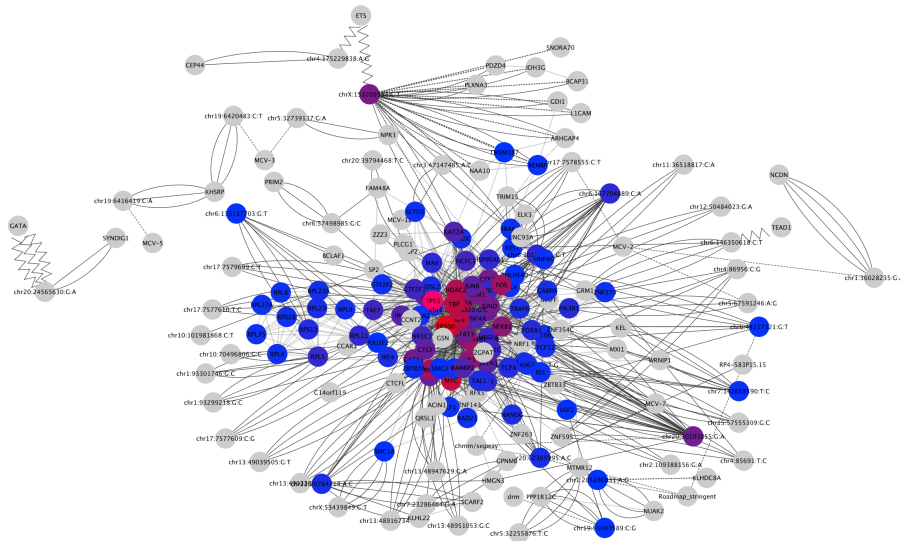
Variant	TF	diff-score
chr20:24565630:G:A	GATA	3.819
chr4:175229838:A:G	ETS	3.271
chr6:146350618:C:T	TEAD1	4.394
chr6:30131441:C:T	ZNF354C	0.811
chrX:153208532:C:T	ETS	4.395

ETS high expression lead to glioma proliferation, and we found the motif gain of ETS followed by highly expression of ETS.

TEAD1, MAPK14 and SERP1 promote glioblastoma progression.

GATA, ZNF354C

# Network visualization



Features are independent use hierarchical rule-based scoring to decrease the dependency of different features, such as: motif gain/break versus gene-link(promoter, enhancer) etc however, the annotation and linkage is quite important, annotation currently cannot  
Todo:

- 1 redefine enhancer-gene linkage
- 2 code rewrite
- 3 flexible on-the-fly weight updating, weight schema and relative importance.
- 4 Add more features, as indicated in nvar grants
- 5 Functional annotation by enhanced gene-based network/visualization

Extend the current output to a network-like view:  
construct gene-based annotation and network, integrate with network analysis for the functional annotation of non-coding variants.  
Motif break/gain; GENE-link,promoter, distal; hot/encode region;

