

I. RICIN ENHANCER SCREEN

OBJECTIVES:

- All candidate enhancer/promoter elements around the set of 884 CRISPRk ricin genes
- Ricin provides a strong phenotype that should result in real hits (also, off-target effects should be mitigated somewhat)
- Format: CRISPRi primarily, but would not hurt to also run CRISPRk/a screens
- single-sgRNA, multiple sgRNAs per enhancer element

Summary of possible libraries

Number sgRNAs*

cRE definition	radius around ricin gene:	30kb		50kb		100kb	
	screen type:	10 sgRNAs	fine mapping	10 sgRNAs	fine mapping	10 sgRNAs	fine mapping
	DNase	168,354	2,091,317	181,367	2,263,783	209,612	2,646,083
H3K27ac	137,814	1,532,074	149,030	1,661,616	173,104	1,941,066	
DNase+H3K27ac	74,631	1,314,458	80,762	1,422,407	94,035	1,665,511	

Number cCREs*

cRE definition	radius around ricin gene:	30kb	50kb	100kb
	DNase	19,342	20,733	23,800
	H3K27ac	14,594	15,741	18,231
DNase+H3K27ac	7,933	8,565	9,941	

Final library:

Two-stage strategy:

1. CRISPRi/a with 5 guides per element
2. Follow up fine-mapping analysis on hits

Final library:

1. CRISPRi/a with 5 guides per element
2. Positive controls: 5 guides per TSS of ricin gene
3. 1000 safes + other positives

Total: ~55K sgRNAs

II. ENHANCER COMBINATORICS SCREENS

OVERALL DESIGN AND PREMISE:

- Picked a dozen surface marker genes + *MYC* and *GATA1*
- Targeted all DNase HS sites in the neighborhood with the 5 best guides available
- All pairwise guide combinations
- For each site, all pairwise guide combinations with safes (assaying the single-sgRNA effects)
- CRISPRi, CRISPRa screens

III. Enhancer combinatorics screens

Summary for surface markers:

gene	#guides
CD114/CSFR3	4096
CD117/c-Kit	13696
CD146/MCAM	1126
CD235a+b	11828
CD24	10586
CD30/TNFRSF8	4751
CD31	4851
CD33	1651
CD41/ITGA2B	2556
EEA1	9455
Ki-67	4371

In addition:

GATA1 – 12,374 guides

MYC -- 34,151 guides

These two have been ordered and exist on a chip

III. CRISPR-X ENHANCER MUTAGENESIS

OVERALL DESIGN AND PREMISE:

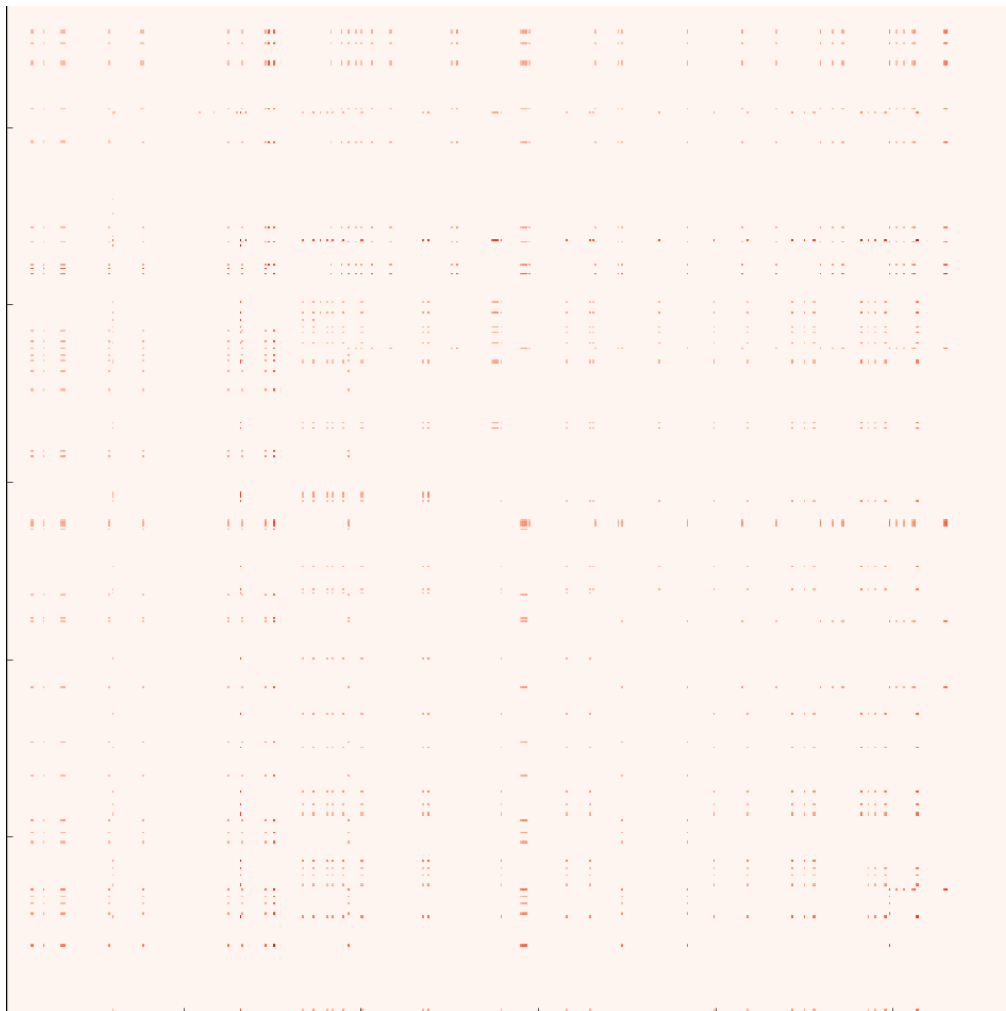
- Generate all pairwise pgRNA pairs within an enhancers that is known to have a growth effect
- Use CRIPSR-X to carry out saturation mutagenesis in a pgRNA format
- Run a growth screen
- Sequence the DNA in the region, see which mutations are enriched/depleted

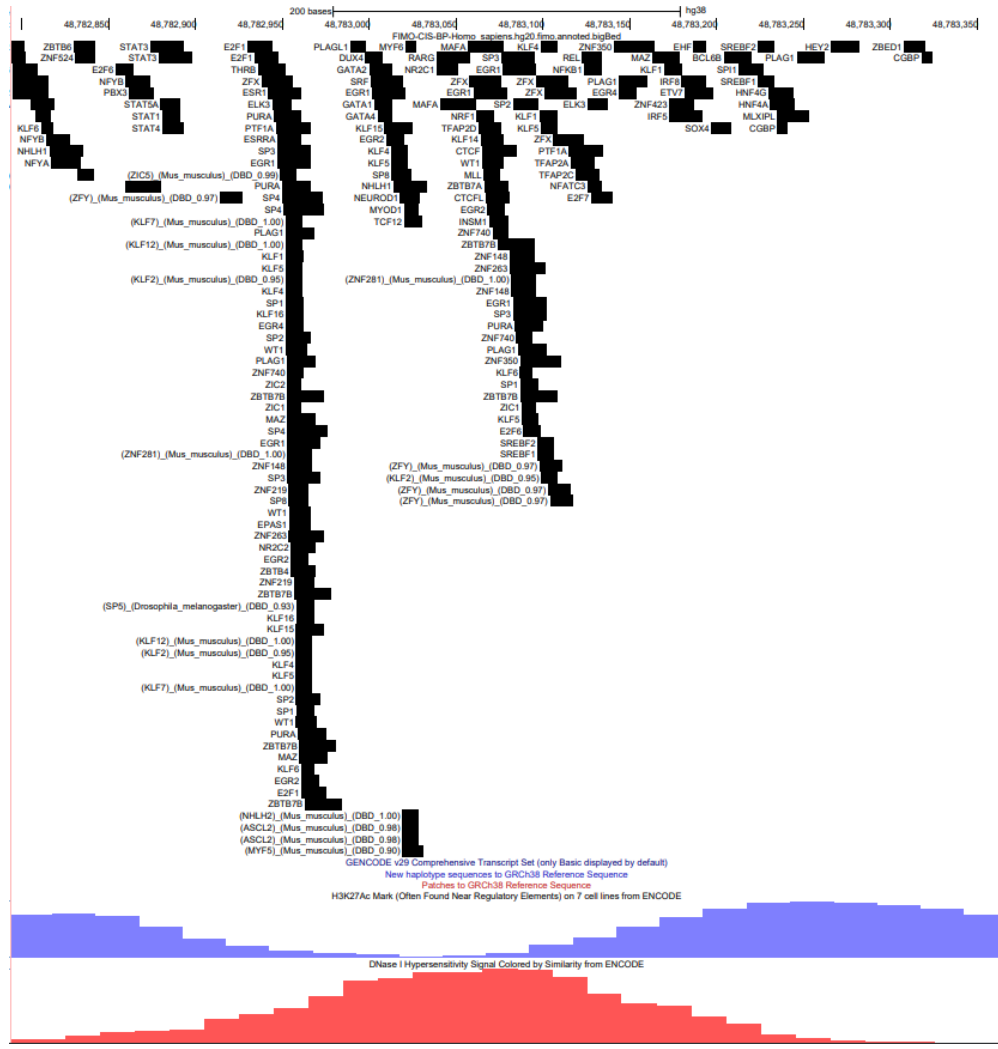
Two libraries were designed and ordered:

eGATA1 – 1,108 guides
eHDAC6 – 1,876 guides

In addition, 4 libraries were designed and ordered for each of the b-globin LCR enhancers (with the idea of the doing flowFISH)

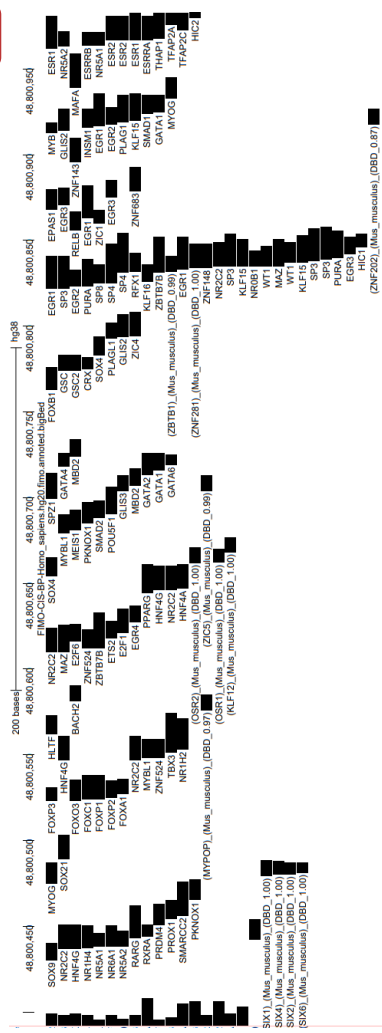
1,108 pgRNAs





1,876 pgRNAs

eHDAC6



H9K27Ac Mark (Often Found Near Regulatory Elements) on 7 cell lines from ENCODE

Disease I Hypersensitivity Signal Colored by Similarity from ENCODE

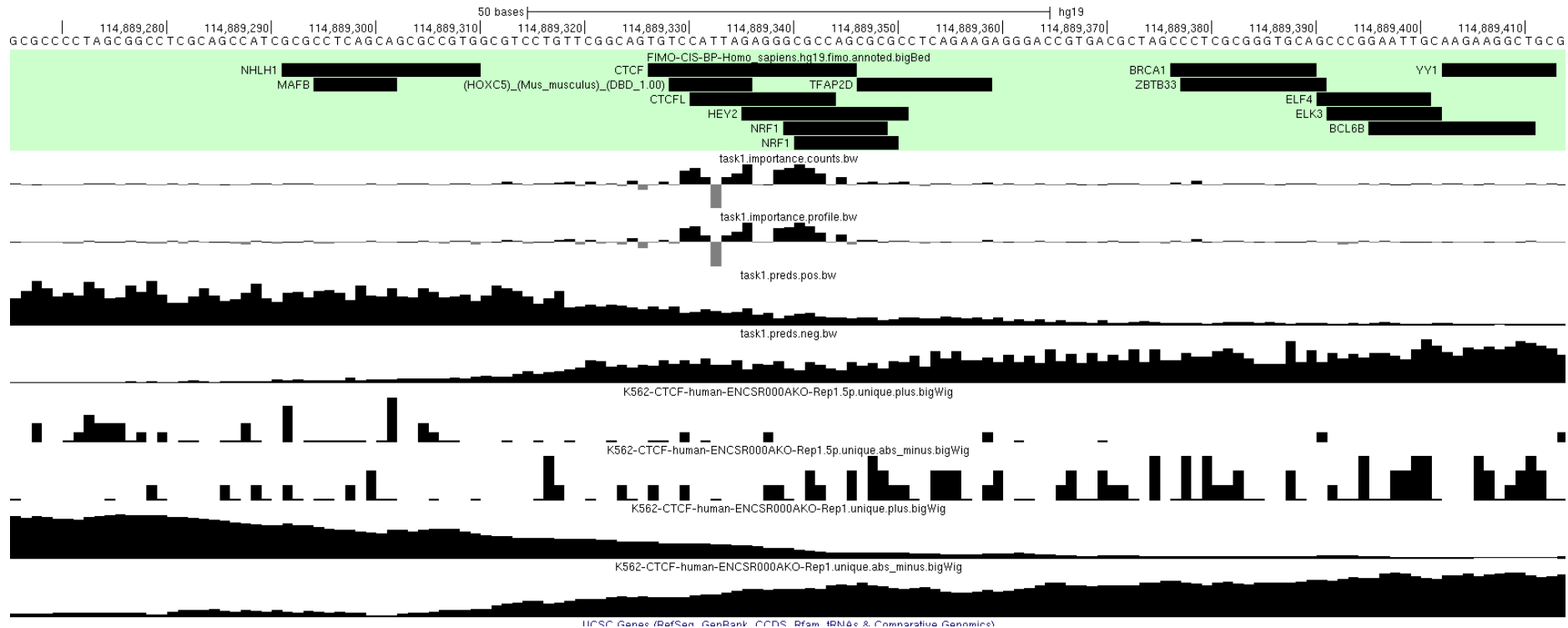


IV. TRANSCRIPTION FACTOR BINDING SATURATION MUTAGENESIS

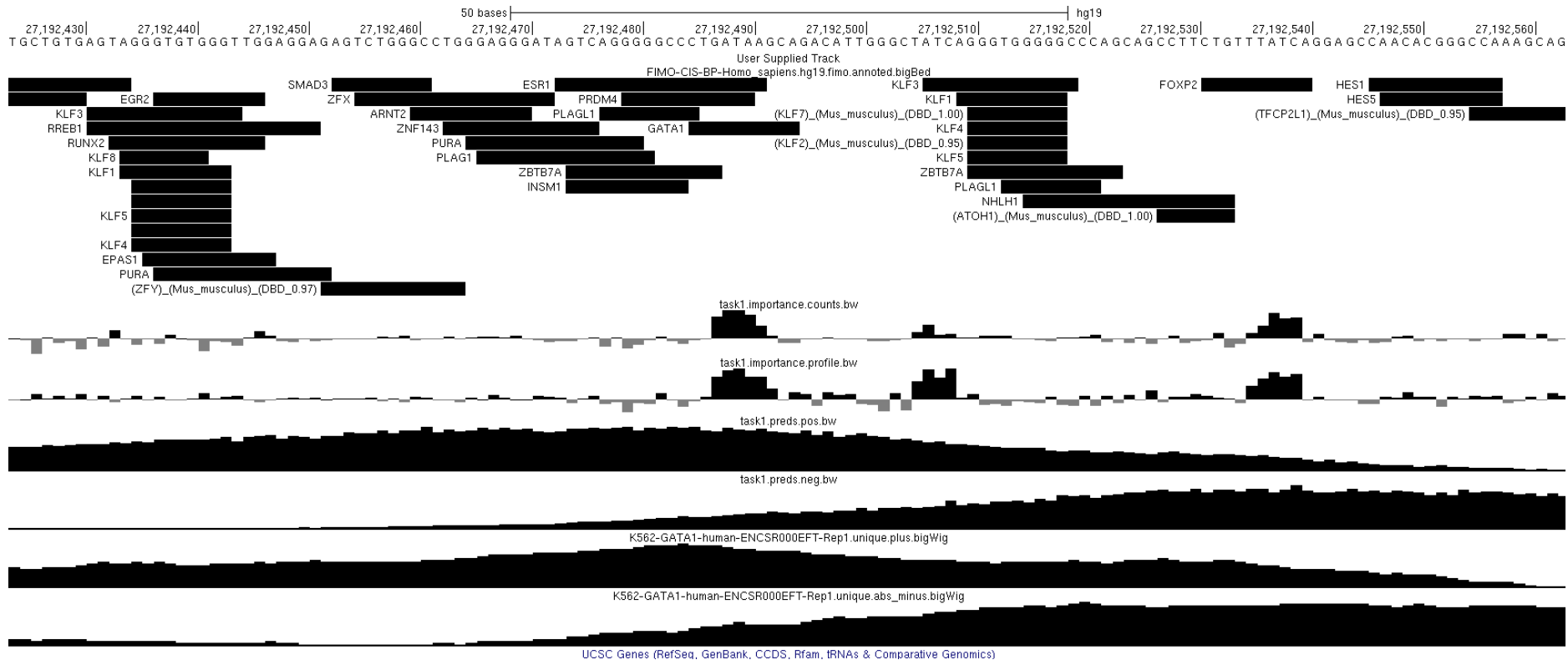
OVERALL DESIGN AND PREMISE:

- Tile a transcription factor binding site and its neighborhood with sgRNAs
- Use CRISPR-X to generate point mutations
- Carry out ChIP experiments against the population of cells (using restriction digestion for fragmentation)
- Do amplicon sequencing on the region(s) of interest in the ChIP and in an input sample
- Look at enriched/depleted variants in the ChIP relative to the input
- The goals are to validate deep learning models and to find interesting cases of interaction between neighboring motifs where binding of one TF is dependent on adjacent binding of other TFs
- Two sites were picked for an initial test – a CTCF site and a GATA1 site (this for technical validation)

CTCF site



GATA1 site



V. PROMOTER-AI VALIDATION WITH ILLUMINA

OVERALL DESIGN AND PREMISE:

- Building on the SpliceAI approach, the Illumina group had generated predictions for sequence variants in the promoter of genes that will disrupt gene expression (“PromoterAI” models)
- The agreement was that we would install those variants using base editing (or more direct approaches) and directly test gene expression effects
- We were given a list of some 400 such sites
- After manual curation and accounting for limitations imposed by editing windows and sequence composition, we were left with 29 sites that could be installed with base editing
- Of course, since then prime editing became available, so if that turns out to be straightforward to adopt, now we can go back to the original list and test many more. But we have not begun to do any of the original 29 either.

VI. PROTEIN CODING GENES SATURATION MUTAGENESIS

OVERALL DESIGN AND PREMISE:

- When those designs were being finalized, we also discussed with Illumina doing saturation mutagenesis on disease-relevant genes
- This would involve tilling a gene with all valid sgRNAs and doing CRISPR-X
- Capture probes would be designed by Illumina, and the sequencing (which would be massive) would also be done by them
- We generated a design for 10 growth genes and also for a set of ricin genes (summarized below)

DESIGN DETAILS

- Used GENCODE V29 (hg38) as annotation
- Split genes by biotypes:
 - protein_coding
 - lincRNAs
 - miRNA
 - Ribozyme
 - snoRNA
 - snRNA
 - scaRNA
 - misc_RNA
- Extended exons by 20bp on each side
- Picked all guides within the extended exons for each gene (so that the site 12-16bp from the PAM is within the extended exon)

Top 10 growth genes

#genesymbol	totalSNPs	nClinvarBenignSNPs	nClinvarPathogenicSNPs	nClinvarTotal
POLG	8098	14	44	58
ACTA1	2499	0	39	39
BRCA2	22765	211	37	248
DKC1	3363	1	37	38
SOD1	1008	0	34	34
KIF1A	11741	3	30	33
SOS1	8792	11	29	40
DYNC1H1	30703	7	23	30
SDHB	1835	5	22	27
GNE	4966	1	21	22

#GeneID	Symbol	GeneInfo	Localization	Process	Function	# elements 1	casTLE Effect 1	casTLE Score 1	# elements 2	casTLE Effect 2	casTLE Score 2	Combo casTLE Effect	Combo casTLE Score	casTLE p-value	Minimum Effect Estimate	Maximum Effect Estimate
ACTA1	ACTA1	N/A	None	None	None	10	-6.1	19.6	10	-5.2	20	-6.1	38.8	0.0273	-11.5	-5.7
SOS1	SOS1	N/A	None	None	None	10	-0.8	11.6	10	-1.1	35.8	-1.1	46.4	0.0119	-1.4	-0.8
KIF1A	KIF1A	N/A	None	None	None	10	-3.4	20.8	10	-2.4	27.1	-3.3	45.7	0.013	-5.7	-2.8
POLG	POLG	N/A	None	None	None	10	-1.7	21.7	10	-1.3	33.9	-1.6	52.9	0.00552	-2.1	-1.2
DKC1	DKC1	N/A	None	None	None	10	-4.2	175	10	-4.8	124	-4.8	296	1.00E-05	-5.6	-4
SDHB	SDHB	N/A	None	None	None	10	-2	16.9	9	-2.5	35.3	-2.5	50.2	0.00758	-3.7	-1.8
BRCA2	BRCA2	N/A	None	None	None	10	-1.1	54.8	10	-1	55.4	-1.1	110	2.00E-05	-1.3	-0.9
DYNC1H1	DYNC1H1	N/A	None	None	None	10	-4.1	67.4	10	-3.3	90.1	-4	154	1.00E-05	-4.9	-3.2
GNE	GNE	N/A	None	None	None	10	-4.3	113	10	-3.2	135	-4.2	243	1.00E-05	-4.9	-3.4
SOD1	SOD1	N/A	None	None	None	10	-2.2	61	10	-2.3	104	-2.2	164	1.00E-05	-2.7	-1.9

All guides, i.e. CFD > 0.0

#gene	exons	total_length	total_guides_passing_CFD	fraction_of_bases_covered
ACTA1	11	1768	236	0.29638009
BRCA2	43	13404	778	0.184049537
DKC1	41	4788	522	0.300960735
DYNC1H1	167	39010	4122	0.29671879
GNE	26	6467	553	0.244008041
KIF1A	107	18957	3761	0.47618294
POLG	81	6912	1311	0.457899306
SDHB	19	2891	323	0.312348668
SOD1	10	2253	315	0.373280071
SOS1	38	10277	781	0.231390484

Total guides: 12702

All guides, i.e. CFD > 0.2

<u>#gene</u>	<u>exons</u>	<u>total_length</u>	<u>total_guides_passing_CFD</u>	<u>fraction_of_bases_covered</u>
ACTA1	11	1768	173	0.238122172
BRCA2	43	13404	437	0.1087735
DKC1	41	4788	357	0.219298246
DYNC1H1	167	39010	2733	0.211407331
GNE	26	6467	378	0.175506417
KIF1A	107	18957	2179	0.311019676
POLG	81	6912	871	0.335069444
SDHB	19	2891	220	0.227948807
SOD1	10	2253	234	0.286284953
SOS1	38	10277	477	0.151211443

Total guides: 8059

Top 10 ricin genes

#genesymbol	totalSNPs	nClinvarBenignSNPs	nClinvarPathogenicSNPs	nClinvarTotal
NF1	18784	12	68	80
PIK3CA	7111	2	35	37
CDKL5	6792	8	34	42
SMAD4	3631	0	31	31
ACTB	2483	0	26	26
PMM2	1636	2	25	27
SMCHD1	13282	6	22	28
HMBS	2354	0	22	22
TBC1D24	3678	7	20	27
MEN1	3952	5	20	25

					Combo casTLE Effect		casTLE p-value
chr16	8882680	8882681	+	PMM2	8.2	1.00E-06	9.33E-05
chr16	2525147	2525148	+	TBC1D24	4.8	0.000453	0.016388
chr17	29421945	29421946	+	NF1	4.7	1.00E-06	9.33E-05
chr11	64578766	64578767	-	MEN1	4.1	1.00E-06	9.33E-05
chrX	18443703	18443704	+	CDKL5	3.7	0.00362	0.086138
chr7	5603415	5603416	-	ACTB	3.6	7.60E-05	0.004327
chr18	2655737	2655738	+	SMCHD1	2.9	1.00E-06	9.33E-05
chr3	1.79E+08	1.79E+08	+	PIK3CA	2.6	0.000241	0.010178
chr18	48494410	48494411	+	SMAD4	2.5	4.00E-06	0.000339
chr11	1.19E+08	1.19E+08	+	HMBS	2.1	0.00125	0.037083

All guides, i.e. CFD > 0.0

#gene	exons	total_length	total_guides_passing_CFD	fraction_of_bases_covered
ACTB	44	3831	417	0.266249021
CDKL5	43	18167	1590	0.251665107
HMBS	77	5317	817	0.397404551
MEN1	31	4597	942	0.486839243
NF1	102	29548	1830	0.19074049
PIK3CA	30	10430	524	0.14928092
PMM2	41	8152	789	0.263493621
SMAD4	38	13715	1131	0.242726941
SMCHD1	79	14541	967	0.202118149
TBC1D24	29	12642	2018	0.405711122

Total guides: 11025

All guides, i.e. CFD > 0.2

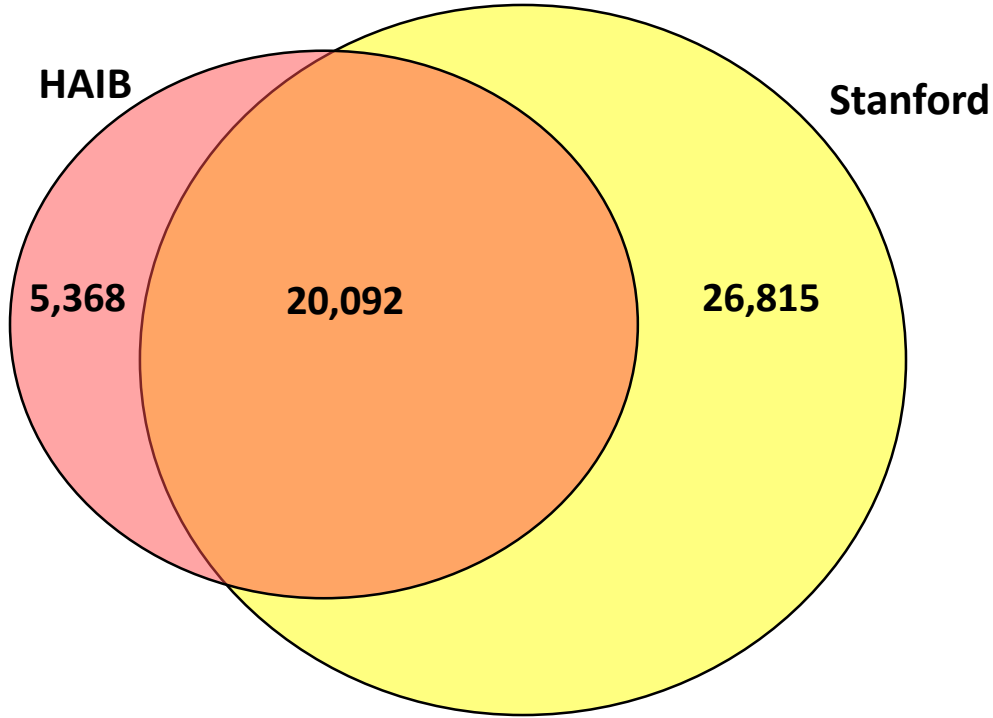
<u>#gene</u>	<u>exons</u>	<u>total_length</u>	<u>total_guides_passing_CFD</u>	<u>fraction_of_bases_covered</u>
ACTB	44	3831	267	0.184547116
CDKL5	43	18167	1024	0.173721583
HMBS	77	5317	505	0.277223998
MEN1	31	4597	596	0.350663476
NF1	102	29548	986	0.109313659
PIK3CA	30	10430	290	0.085522531
PMM2	41	8152	521	0.18326791
SMAD4	38	13715	649	0.149325556
SMCHD1	79	14541	557	0.124269307
TBC1D24	29	12642	1046	0.23991457

Total guides: 6441

VII. MOTIF SCREENS

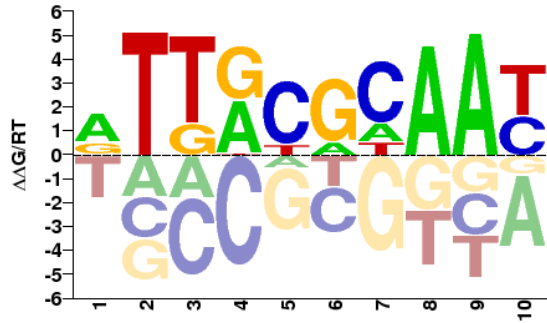
OVERALL DESIGN AND PREMISE:

- We did the CTCF screen, but now we know that CTCF is not an ordinary TF.
- The goal of this phase of ENCODE is *functional characterization*. It would still be a valuable contribution to properly characterize essentiality for a set of regular TFs.
- These would be CRISPRk screens (with an option to carry out CRISPRi in parallel)
- The screens would be done assaying for both growth and ricin sensitivity in parallel



Used the union

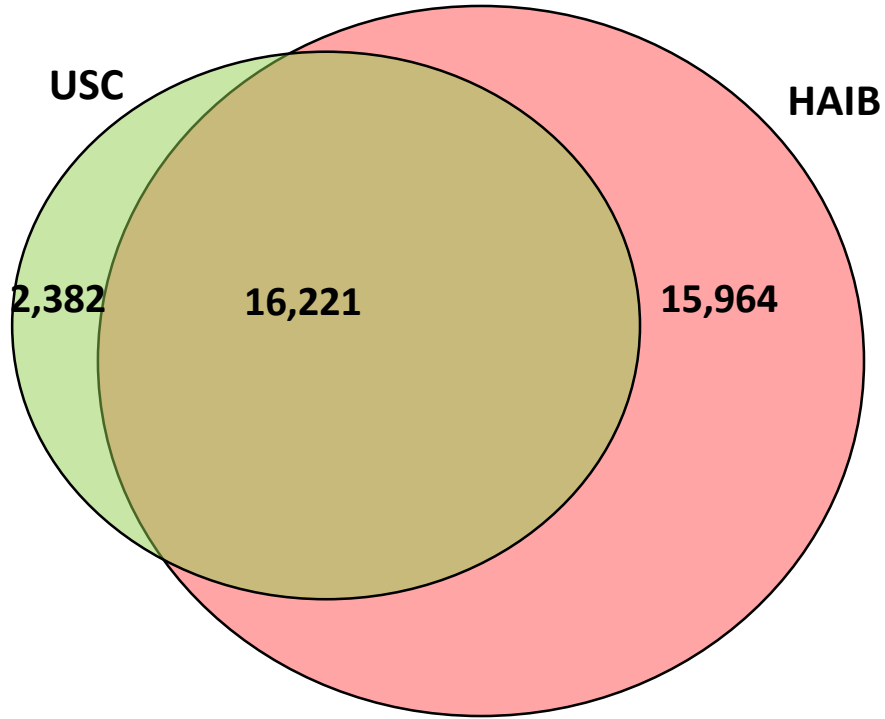
CEBP motif (CIS-BP)



- present in 29,185 out of 52,275 ChIP peaks

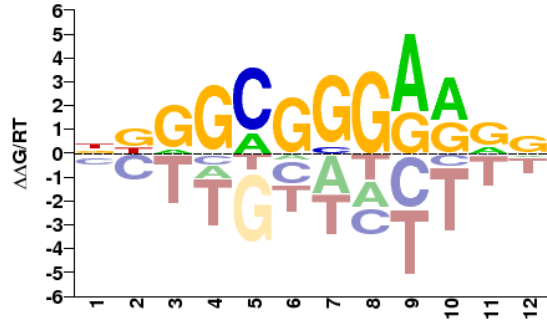
Library summary (without safes)

- 45,837 sgRNAs
- 17,802 ChIP-seq peaks
- 19,035 motifs



Used the union

E2F6 motif (CIS-BP)

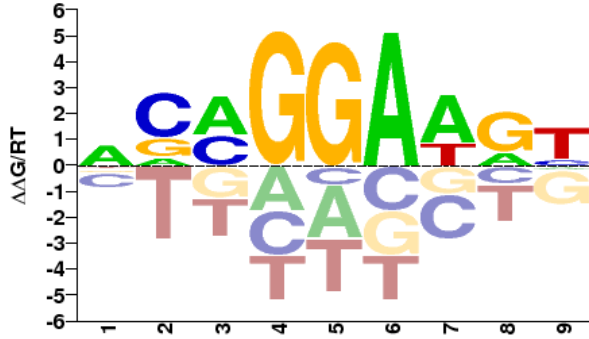


- present in 17,018 out of 34,567 ChIP peaks

Library summary (without safes)

- 161,050 sgRNAs
- 15,938 ChIP-seq peaks
- 23,822 motifs

ETS1 motif (CIS-BP)

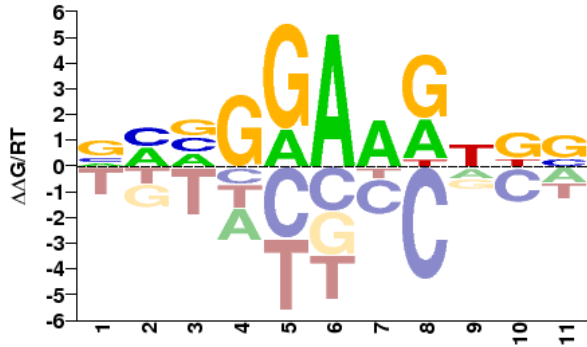


- present in 2,998 out of 11,327 ChIP peaks

Library summary (without safes)

- 30,260 sgRNAs
- 2,878 ChIP-seq peaks
- 3,566 motifs

ETS2 motif (CIS-BP)

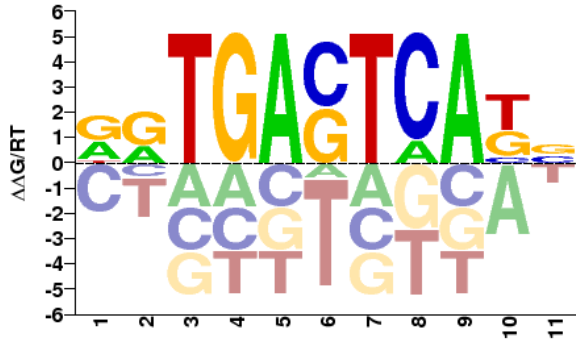


- present in 454 out of 1,705 ChIP peaks

Library summary (without safes)

- 3,070 sgRNAs
- 363 ChIP-seq peaks
- 458 motifs

FOSL1 motif (CIS-BP)

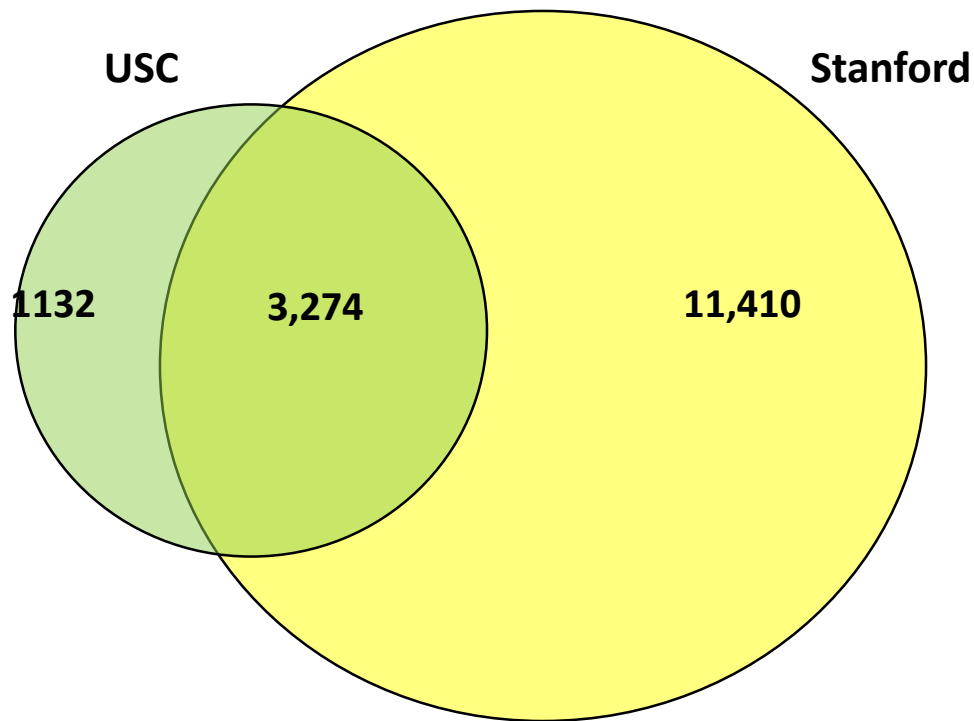


- present in 6,568 out of 8,194 ChIP peaks

Library summary (without safes)

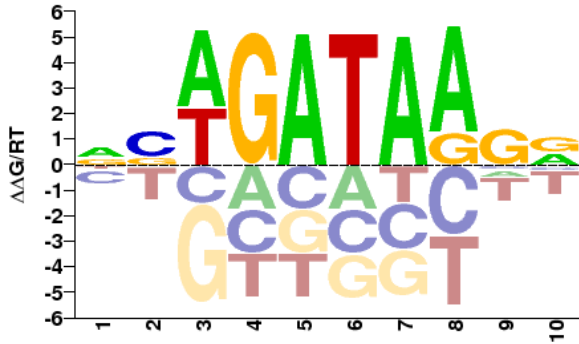
- 24,198 sgRNAs
- 5,649 ChIP-seq peaks
- 10,164 motifs

K562 GATA1



Used the union

GATA1 motif (CIS-BP)

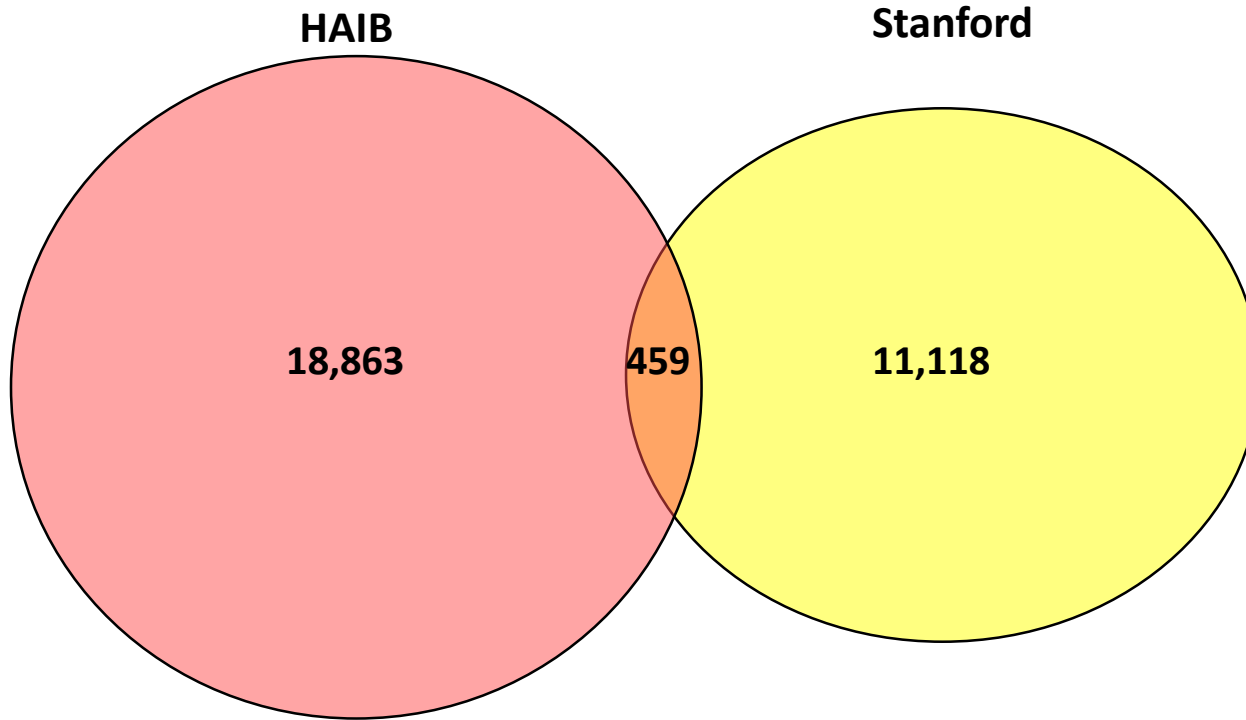


- present in 6,749 out of 15,816 ChIP peaks

Library summary (without safes)

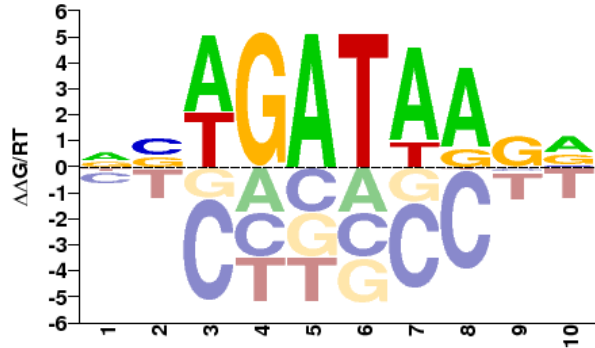
- 22,620 sgRNAs
- 5,874 ChIP-seq peaks
- 6,599 motifs

K562 GATA2



Used the union, results to be parsed after screen (will be interesting)

GATA2 motif (CIS-BP)



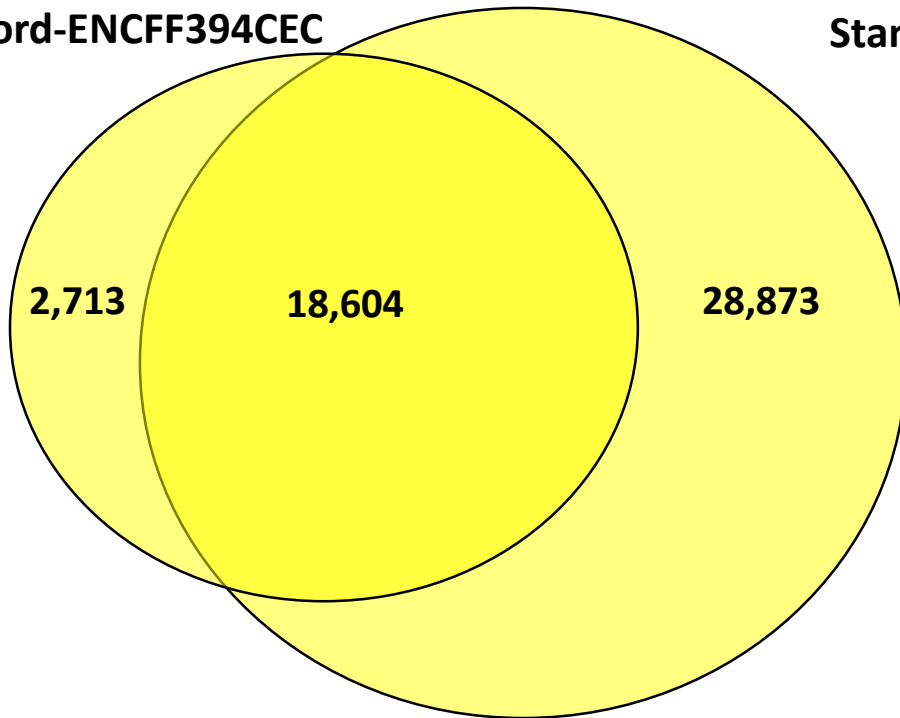
- present in 5,547 out of 30,440 ChIP peaks

Library summary (without safes)

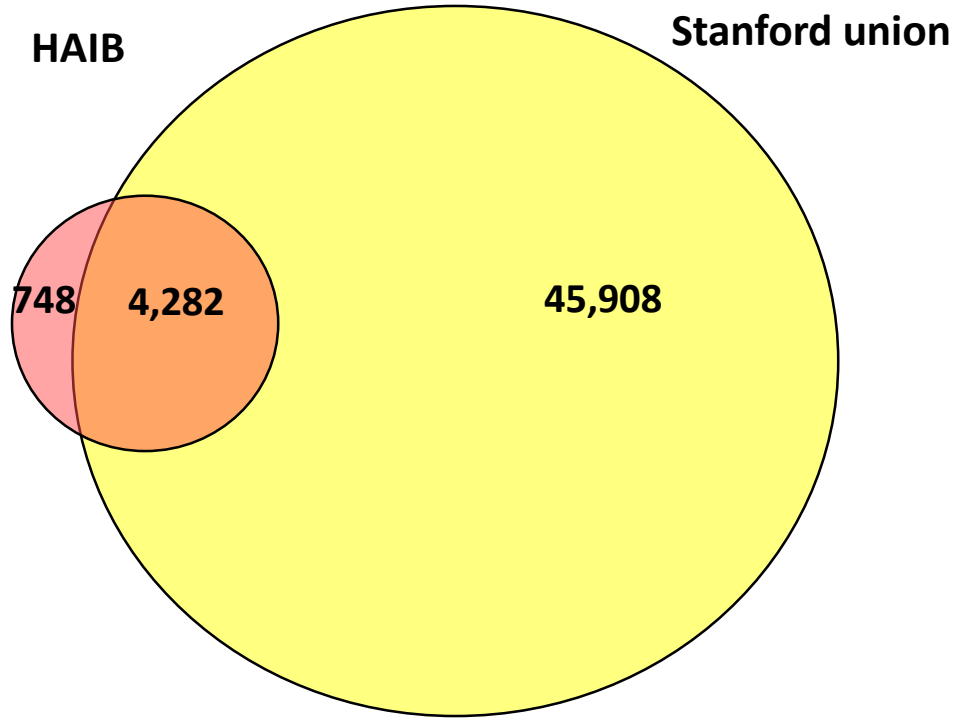
- 17,898 sgRNAs
- 4,601 ChIP-seq peaks
- 4,988 motifs

Stanford-ENCFF394CEC

Stanford-ENCFF213EYD

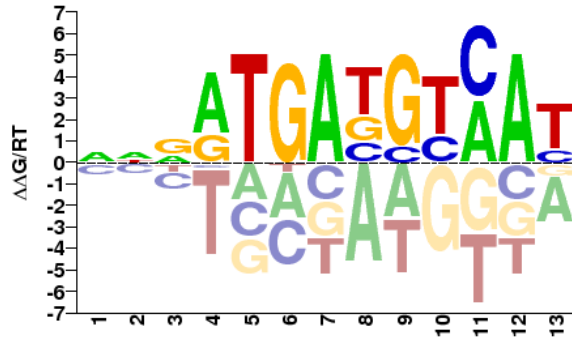


Used the union



Used the union

JUN motif (CIS-BP)

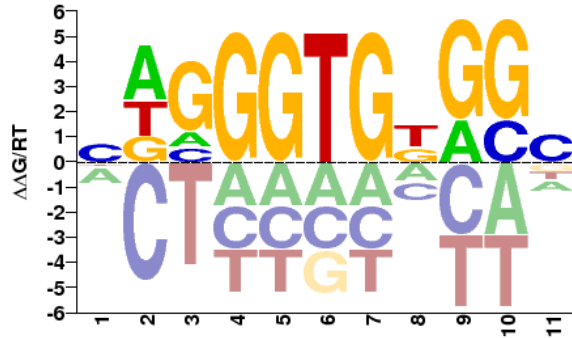


- present in 7,535 out of 50,938 ChIP peaks

Library summary (without safes)

- 23,103 sgRNAs
- 6,190 ChIP-seq peaks
- 7,571 motifs

KLF1 motif (CIS-BP)

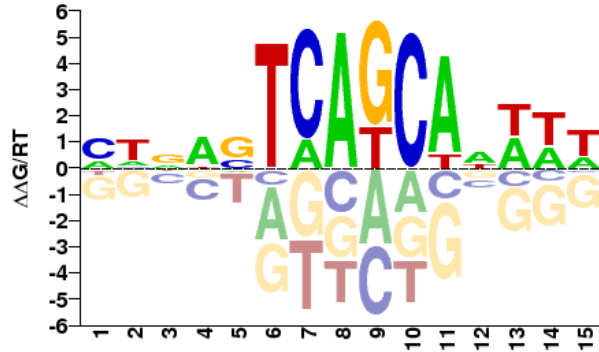


- present in 5,561 out of 11,755 ChIP peaks

Library summary (without safes)

- 52,008 sgRNAs
- 5,038 ChIP-seq peaks
- 6,624 motifs

MAFK motif (CIS-BP)

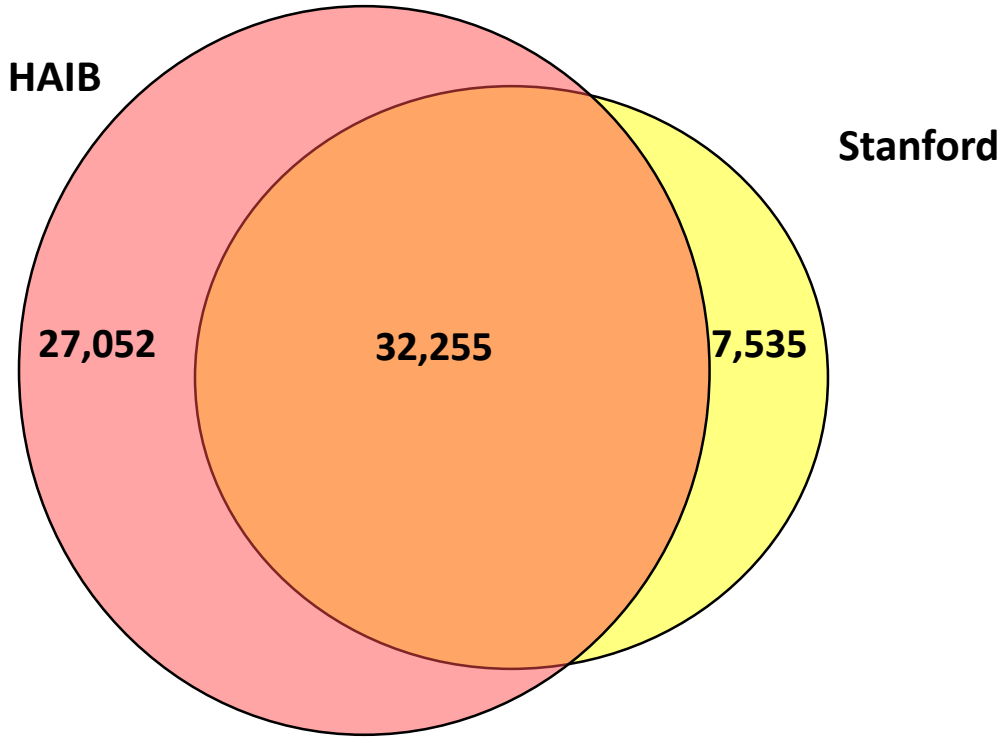


- present in 15,772 out of 26,862 ChIP peaks

Library summary (without safes)

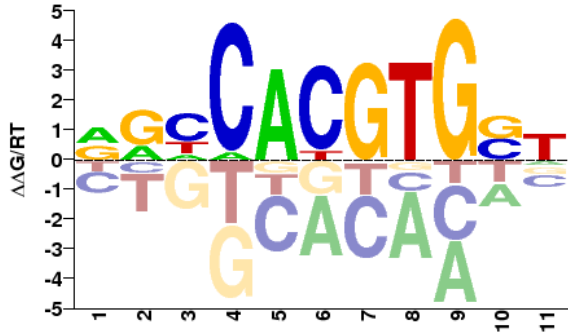
- 32,103 sgRNAs
- 11,583 ChIP-seq peaks
- 13,246 motifs

K562 MAX



Used the union

MAX motif (CIS-BP)

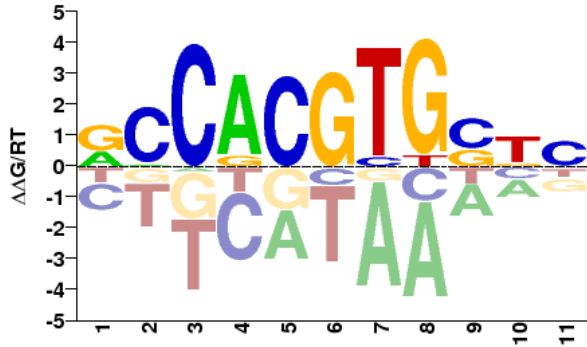


- present in 17,144 out of 66,842 ChIP peaks

Library summary (without safes)

- 71,600 sgRNAs
- 15,336 ChIP-seq peaks
- 17,391 motifs

MYC motif (CIS-BP)

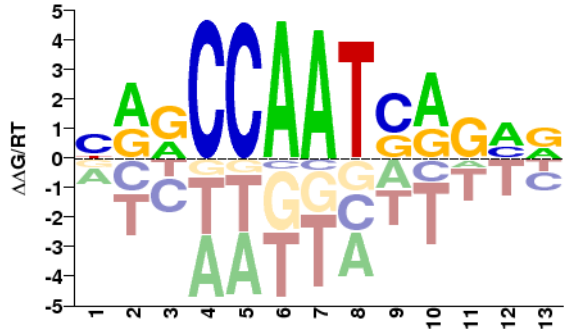


- present in 7,854 out of 31,139 ChIP peaks

Library summary (without safes)

- 44,714 sgRNAs
- 7,305 ChIP-seq peaks
- 10,060 motifs

NFYB motif (CIS-BP)

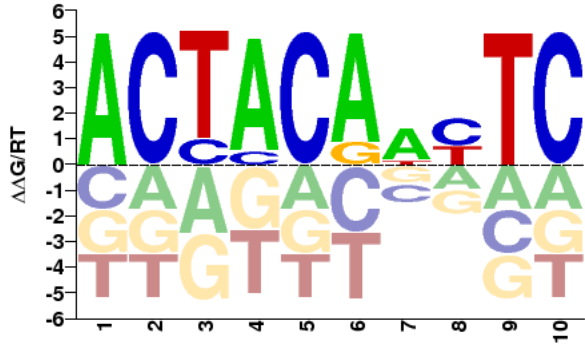


- present in 6,937 out of 9,159 ChIP peaks

Library summary (without safes)

- 45,783 sgRNAs
- 5,696 ChIP-seq peaks
- 8,285 motifs

SIX5 motif (CIS-BP)

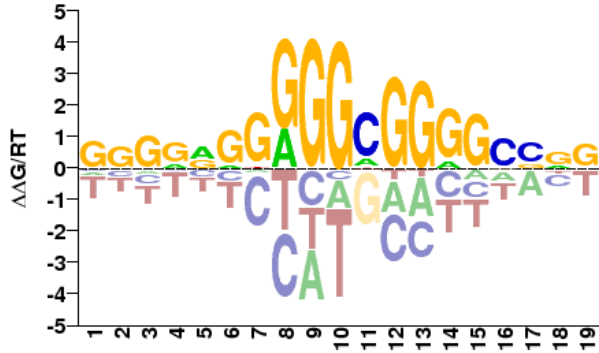


- present in 1,420 out of 3,590 ChIP peaks

Library summary (without safes)

- 4,645 sgRNAs
- 1,142 ChIP-seq peaks
- 1,353 motifs

SP1 motif (CIS-BP)



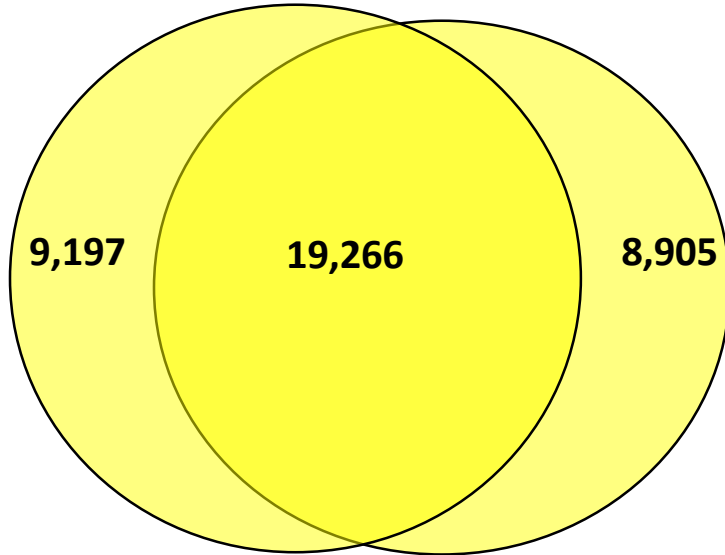
- present in 11,871 out of 14,782 ChIP peaks

Library summary (without safes)

- 218,798 sgRNAs
- 10,884 ChIP-seq peaks
- 27,977 motifs

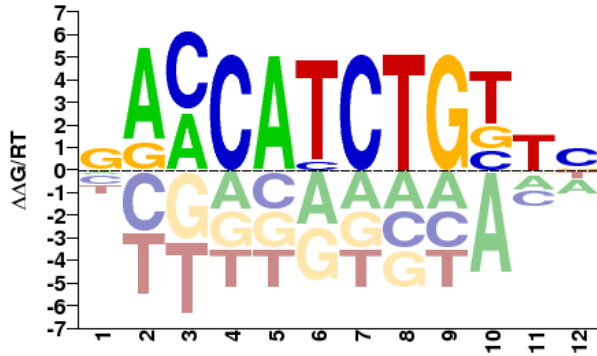
Stanford-ENCFF078OUD

Stanford-ENCFF475LFH



Used the union

TAL1 motif (CIS-BP)



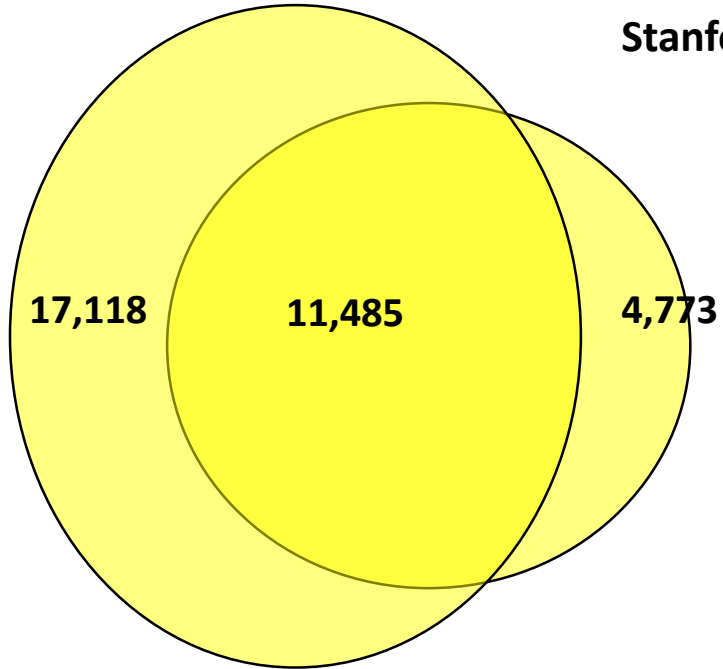
- present in 4,294 out of 37,368 ChIP peaks

Library summary (without safes)

- 12,537 sgRNAs
- 3,401 ChIP-seq peaks
- 3,615 motifs

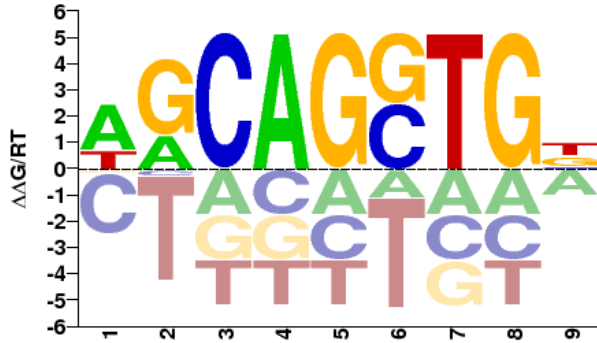
Stanford-ENCFF912LXU

Stanford-ENCFF952JIK



Used the union

TCF12 motif (CIS-BP)

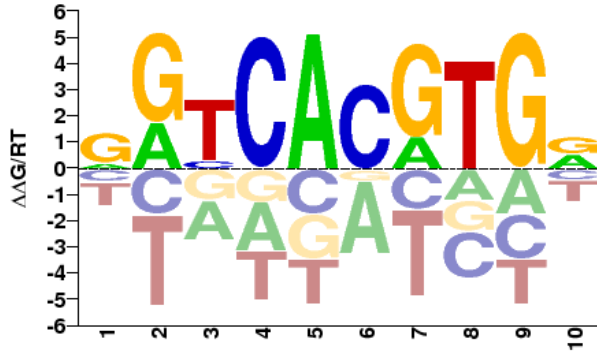


- present in 7,086 out of 33,376 ChIP peaks

Library summary (without safes)

- 30,362 sgRNAs
- 5,821 ChIP-seq peaks
- 8,340 motifs

USF1 motif (CIS-BP)

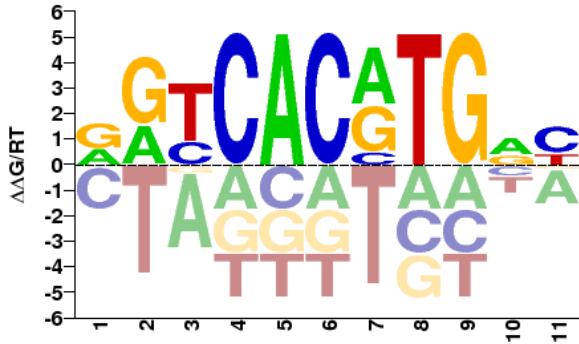


- present in 14,014 out of 21,382 ChIP peaks

Library summary (without safes)

- 61,147 sgRNAs
- 12,715 ChIP-seq peaks
- 18,503 YY1 motifs

USF2 motif (CIS-BP)

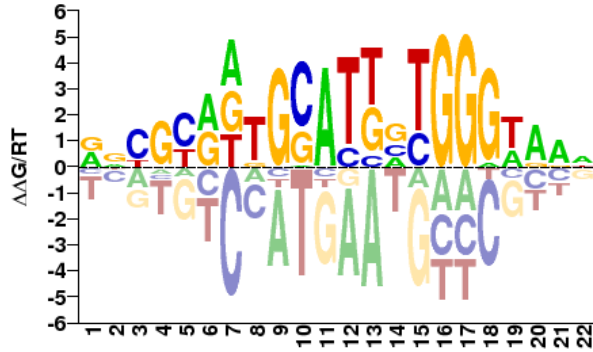


- present in 2,674 out of 3,542 ChIP peaks

Library summary (without safes)

- 17,912 sgRNAs
- 2,526 ChIP-seq peaks
- 4,606 motifs

ZNF143 motif (CIS-BP)



- present in 4,765 out of 29,840 ChIP peaks

Library summary (without safes)

- 28,957 sgRNAs
- 4,267 ChIP-seq peaks
- 5,244 motifs