



Deciphering human non-coding DNA using machine learning approaches

Guillaume Bourque

Department of Human Genetics, McGill University and McGill
University and Genome Quebec Innovation Center

MLPM Summer School
Paris, September 17th 2014

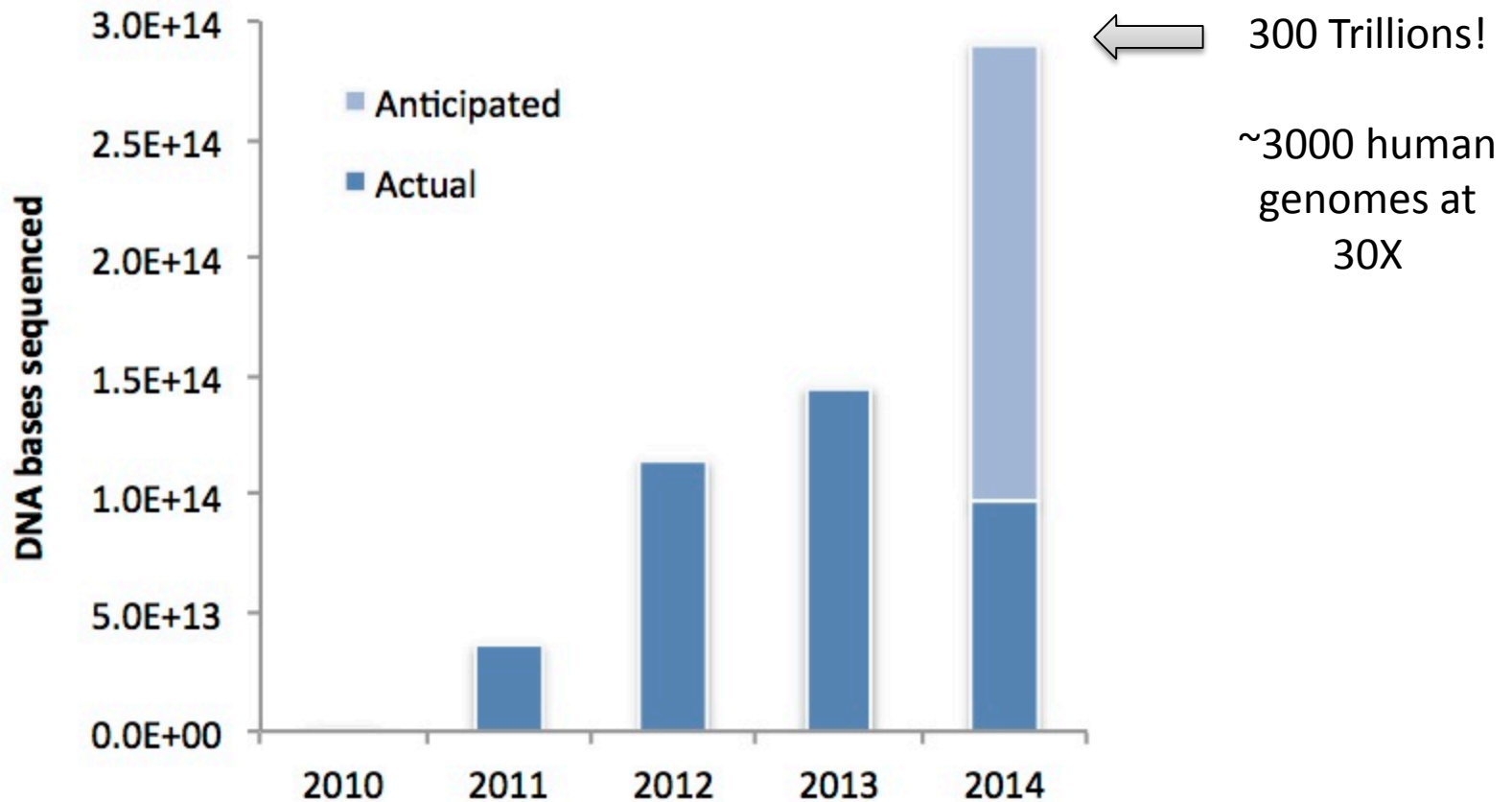
Outline

- Applications of next-generation sequencing
- Functional genomics
- Example of machine learning approaches in functional genomics
- Role of transposable elements in gene regulation

Outline

- Applications of next-generation sequencing
- Functional genomics
- Example of machine learning approaches in functional genomics
- Role of transposable elements in gene regulation

DNA bases sequenced at the Innovation Center



Sequencing human genomes

2001

The
Human
Genome

~ 3 Billion \$

2011

1000
Genomes
Project

~ 10 000 \$

2015 (?)

Your
Genome

< 1000 \$

Big Data

2014



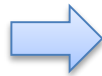
2 X 10 TBytes



1 TBytes

Intensity files

Reads + qualities



300 TBytes



15 TBytes

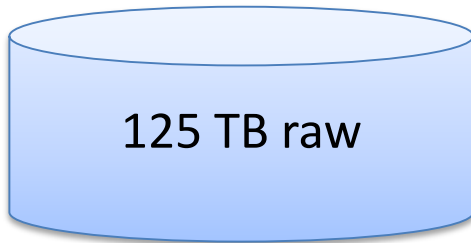
30 TB of raw data / month

360 TB of raw data / year

Large NGS project

Cancer project with whole genome data:

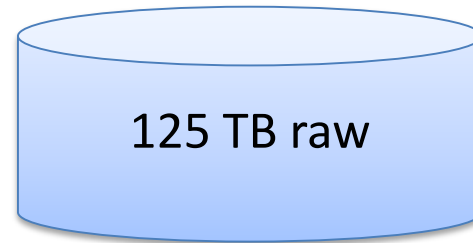
500 tumors



500 X 3 lanes = 500 X 250GB

vs

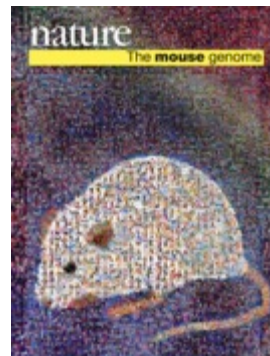
500 matched-normal



500 X 3 lanes = 500 X 250GB

Applications (I)

- *De novo* sequencing
 - From the human genome... To all model organisms... To all relevant organisms (e.g. extreme genomes)... To “all” organisms?



Applications (II)

- Genome re-sequencing
 - Map genomic structural variations across individuals (to understand genetic disorders and also susceptibility factors)
 - Cancer genome sequencing
 - Agricultural crops



1000 Genomes Project



The Cancer Genome Atlas



International
Cancer Genome
Consortium

Complete Resequencing of 40 Genomes Reveals Domestication Events and Genes in Silkworm (*Bombyx*)

Qingyou Xia,^{1,2,*} Yiran Guo,^{3,*} Ze Zhang,^{1,2,*} Dong Li,^{1,2,*} Zhaoling Xuan,^{3,*} Zhuo Li,^{3,*} Fangyin Dai,¹ Yingrui Li,¹ Daojun Cheng,¹ Ruiqiang Li,^{3,4} Tingcai Cheng,^{1,2} Tao Jiang,¹ Celine Becquet,^{1,†} Xan Xu,² Chun Liu,² Xingfu Zha,¹ Wei Fan,² Ying Lin,² Yihong Shen,¹ Lan Jiang,¹ Jeffrey Jensen,⁵ Ines Hellmann,⁵ Si Tang,⁵ Ping Zhao,¹ Hanfu Xu,¹ Chang Yu,¹ Guojie Zhang,² Jun Li,² Jianjun Cao,² Shiping Liu,¹ Ningjia He,² Yan Zhou,² Hui Liu,¹ Jing Zhao,² Chen Ye,² Zhouhe Du,¹ Guoqing Pan,¹ Aichun Zhao,¹ Haojing Shao,^{3,7} Wei Zeng,² Ping Wu,¹ Chunfeng Li,¹ Minhui Pan,¹ Jingjing Li,¹ Xuyang Yin,¹ Duwei Li,² Juan Wang,² Huisong Zheng,¹ Wen Wang,² Xiuqing Zhang,¹ Songgang Li,¹ Huanming Yang,¹ Cheng Lu,¹ Rasmus Nielsen,^{4,5} Zeyang Zhou,^{1,6} Jian Wang,² Zhonghui Xiang,^{1,†} Jun Wang^{1,6,†}

Exome sequencing for Mendelian disease

REVIEWS

 TRANSLATIONAL GENETICS

Exome sequencing as a tool for Mendelian disease gene discovery

Michael J. Bamshad†, Sarah B. Ng†, Abigail W. Bigham*§, Holly K. Tabor*||, Mary J. Emond†, Deborah A. Nickerson† and Jay Shendure†*

Abstract | Exome sequencing — the targeted sequencing of the subset of the human genome that is protein coding — is a powerful and cost-effective new tool for dissecting the genetic basis of diseases and traits that have proved to be intractable to conventional gene-discovery strategies. Over the past 2 years, experimental and analytical approaches relating to exome sequencing have established a rich framework for discovering the genes underlying unsolved Mendelian disorders. Additionally, exome sequencing is being adapted to explore the extent to which rare alleles explain the heritability of complex diseases and health-related traits. These advances also set the stage for applying exome and whole-genome sequencing to facilitate clinical diagnosis and personalized disease-risk profiling.

Exome sequencing for Mendelian disease

REVIEWS

 TRANSLATIONAL GENETICS

Exome sequencing as a tool for Mendelian disease gene discovery

Michael J. Bamshad†, Sarah B. Ng†, Abigail W. Bigham*§, Holly K. Tabor*||, Mary J. Emond†, Deborah A. Nickerson† and Jay Shendure†*

Abstract | Exome sequencing — the targeted sequencing of the subset of the human genome that is protein coding — is a powerful and cost-effective new tool for dissecting the genetic basis of diseases and traits that have proved to be intractable to conventional gene-discovery strategies. Over the past 2 years, experimental and analytical approaches relating to exome sequencing have established a rich framework for discovering the genes underlying unsolved Mendelian disorders. Additionally, exome sequencing is being adapted to explore the extent to which rare alleles explain the heritability of complex diseases and health-related traits. These advances also set the stage for applying exome and whole-genome sequencing to facilitate clinical diagnosis and personalized disease-risk profiling.

“... about one-half to one-third (~3,000) of all known or suspected Mendelian disorders (for example, cystic fibrosis and sickle cell anaemia) have been discovered. However, there is a substantial gap in our knowledge about the genes that cause many rare Mendelian phenotypes.”

“Accordingly, we can realistically look towards a future in which the genetic basis of all Mendelian traits is known, ...”

Cancer genome sequencing

REVIEWS

APPLICATIONS OF NEXT-GENERATION SEQUENCING

Advances in understanding cancer genomes through second-generation sequencing

Matthew Meyerson, Stacey Gabriel and Gad Getz

Abstract | Cancers are caused by the accumulation of genomic alterations. Therefore, analyses of cancer genome sequences and structures provide insights for understanding cancer biology, diagnosis and therapy. The application of second-generation DNA sequencing technologies (also known as next-generation sequencing) — through whole-genome, whole-exome and whole-transcriptome approaches — is allowing substantial advances in cancer genomics. These methods are facilitating an increase in the efficiency and resolution of detection of each of the principal types of somatic cancer genome alterations, including nucleotide substitutions, small insertions and deletions, copy number alterations, chromosomal rearrangements and microbial infections. This Review focuses on the methodological considerations for characterizing somatic genome alterations in cancer and the future prospects for these approaches.

Cancer genome sequencing

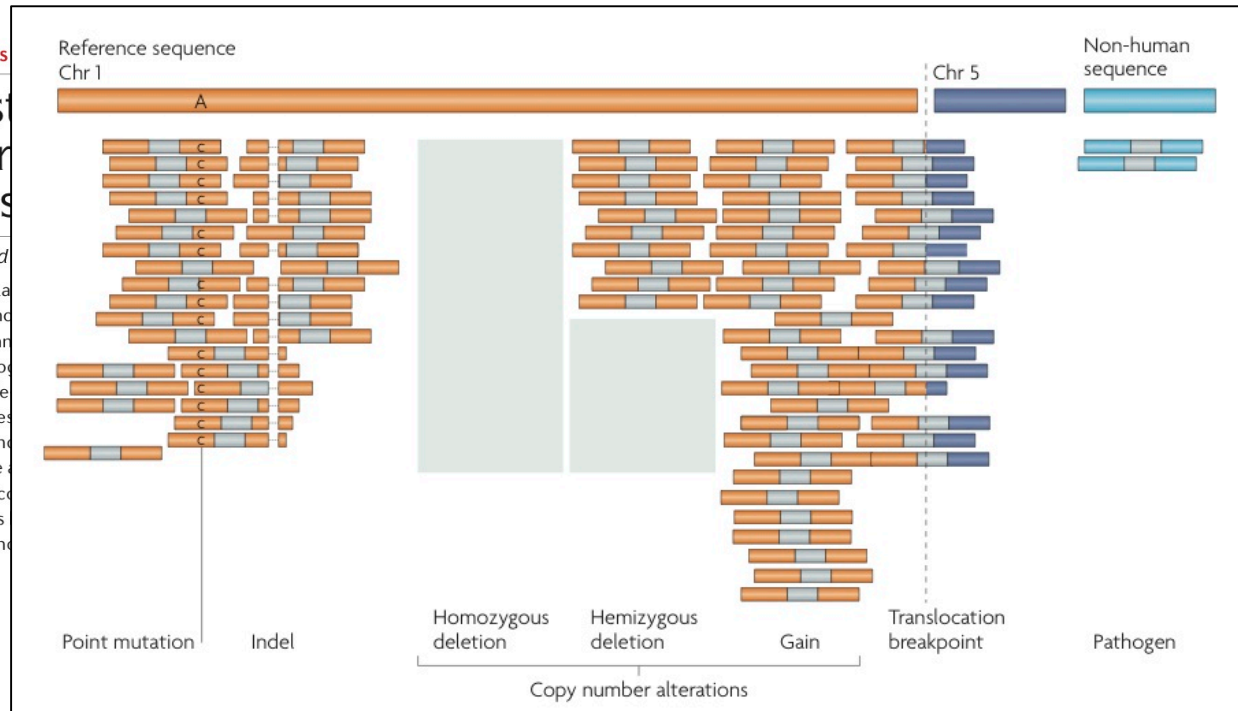
REVIEWS

APPLICATIONS OF NEXT-GENERATION SEQUENCING

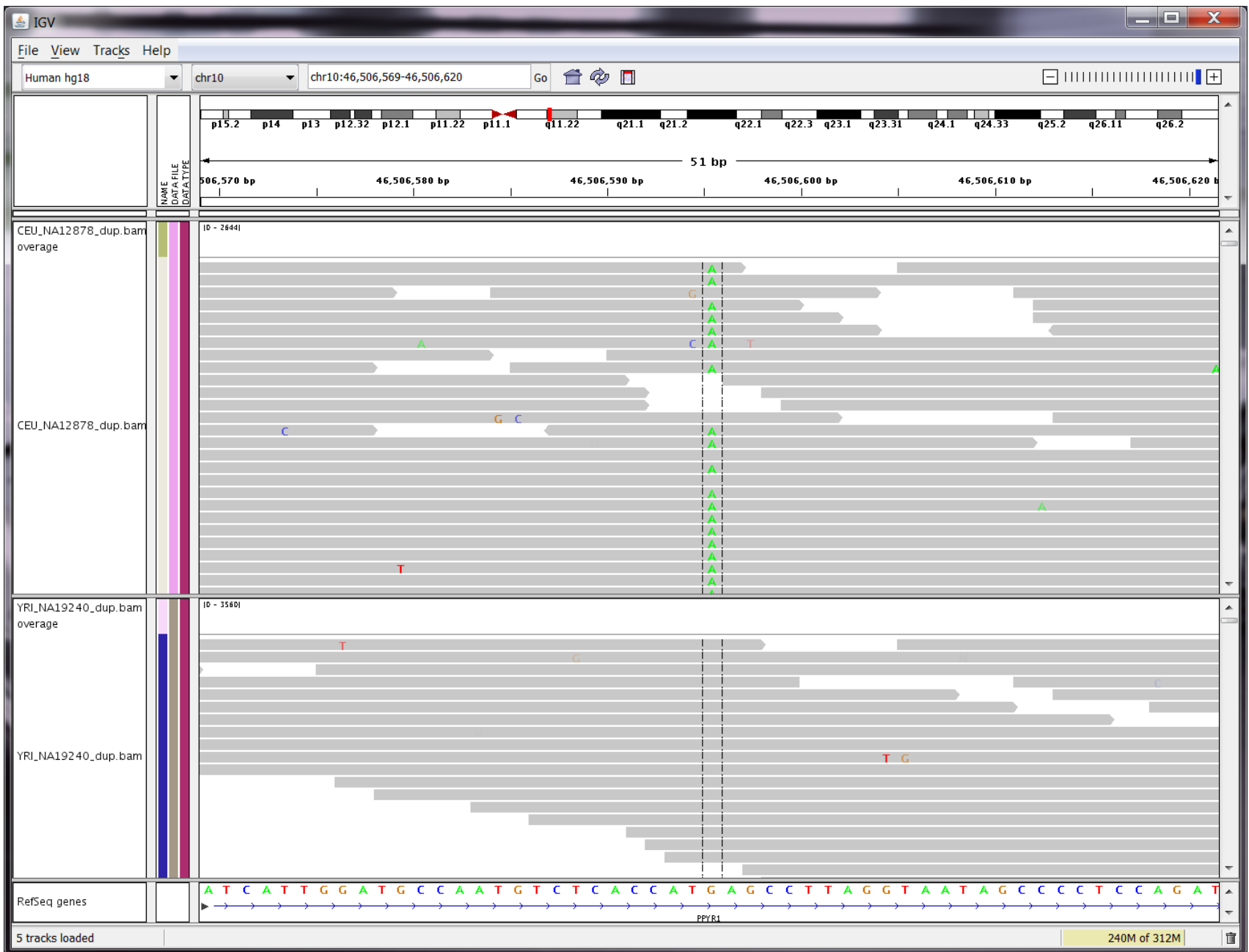
Advances in understanding cancer genomes through second-generation sequencing

Matthew Meyerson, Stacey Gabriel and Gad Getz

Abstract | Cancers are caused by the accumulation of somatic mutations. Therefore, analyses of cancer genome sequences (through whole-genome, whole-exome, or targeted second-generation DNA sequencing technologies) — through whole-genome, whole-exome, or targeted approaches — is allowing substantial advances in our understanding of the principal types of somatic cancer genome alterations, including point mutations, small insertions and deletions, copy number alterations, rearrangements and microbial infections. This review discusses the considerations for characterizing somatic genome alterations and the prospects for these approaches.

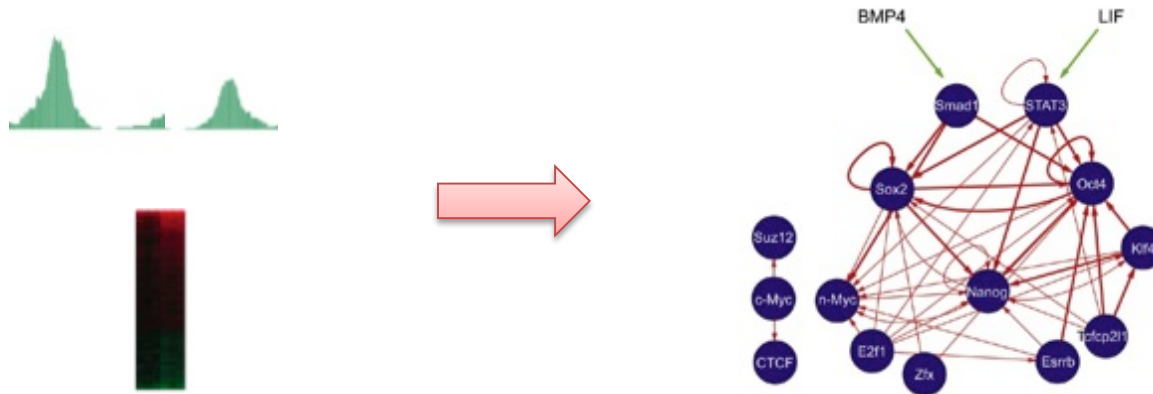


Can obtain a full catalogue of mutations

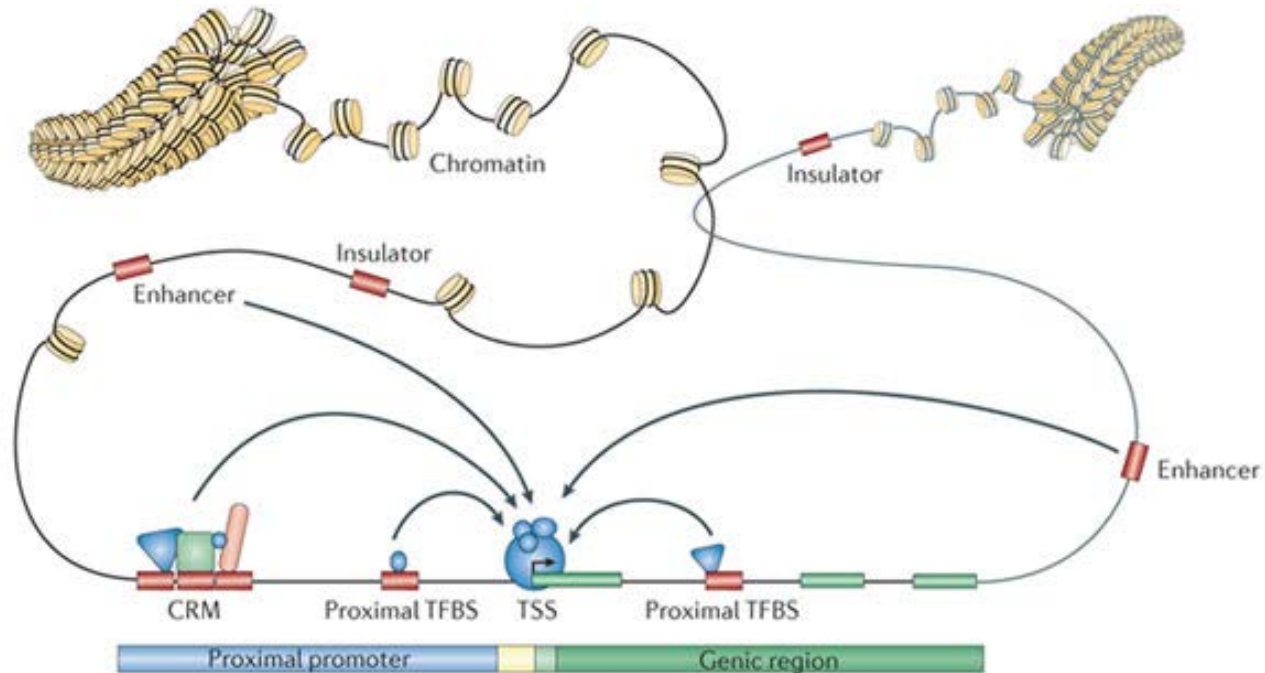


Applications (III)

- Quantitative biology of complex systems
 - New high-throughput technologies in functional genomics: ChIP-Seq, RNA-Seq, ChIA-PET, RIP-Seq, ...
 - From single-gene measurements, to thousands of probes on arrays, to profiles covering all 3B bases of the genome



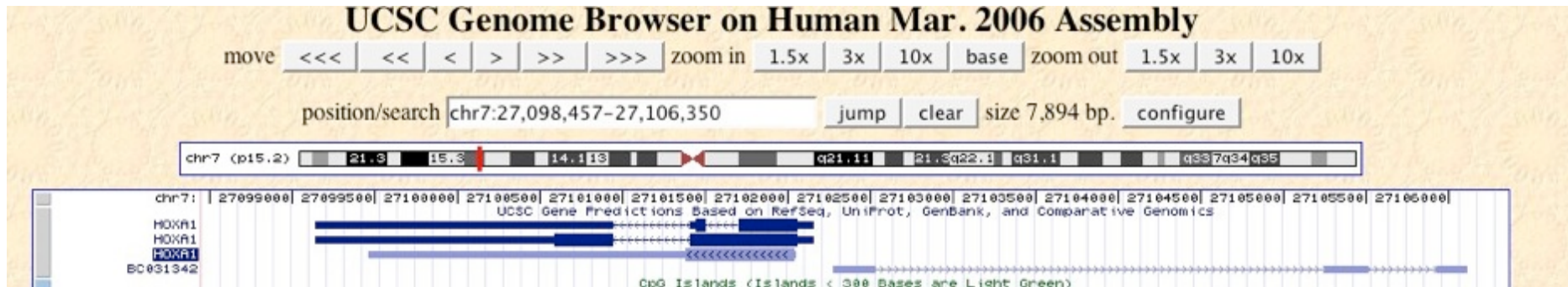
Transcription Regulation



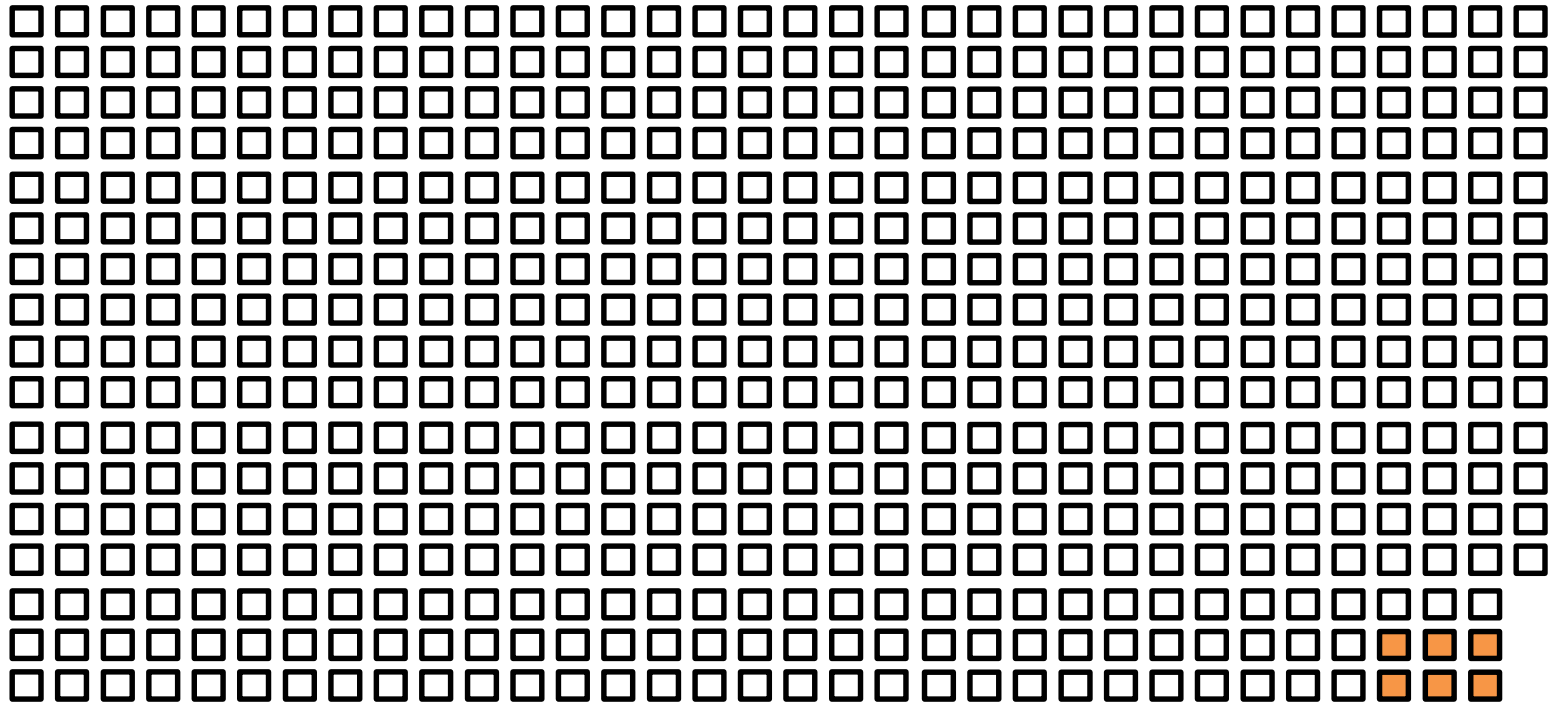
Lenhard, Nat Rev Genet, 2012

Functional genomics

Locate regulatory regions in the human genome in different cell types, describe their functions, and identify how they differ between different groups (i.e. “disease” vs “healthy”)...

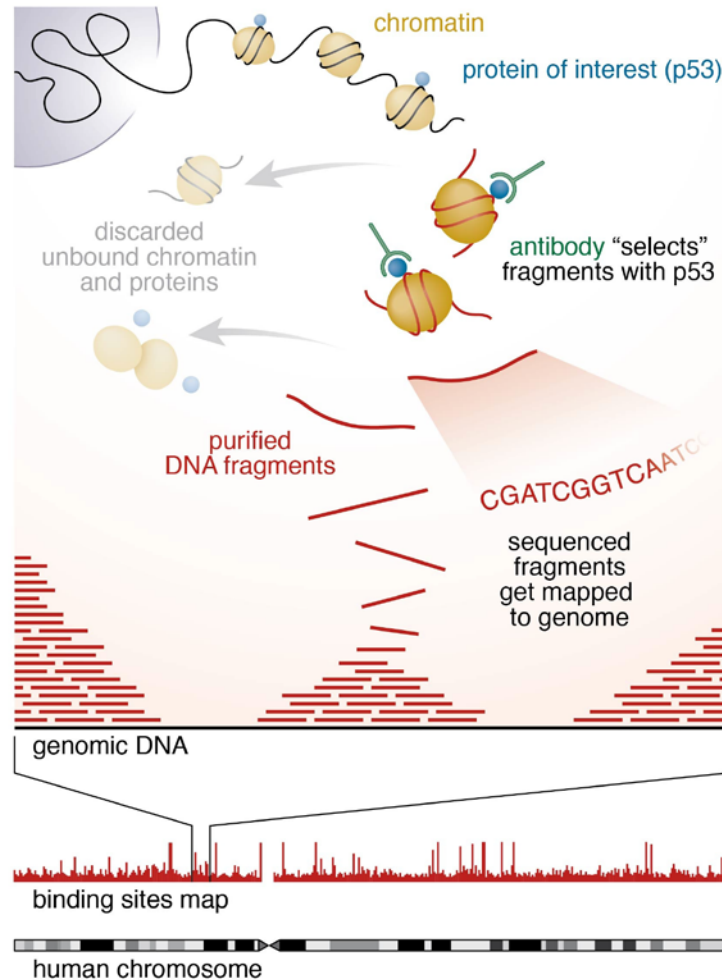


Somatic mutations in 100 kidney tumors

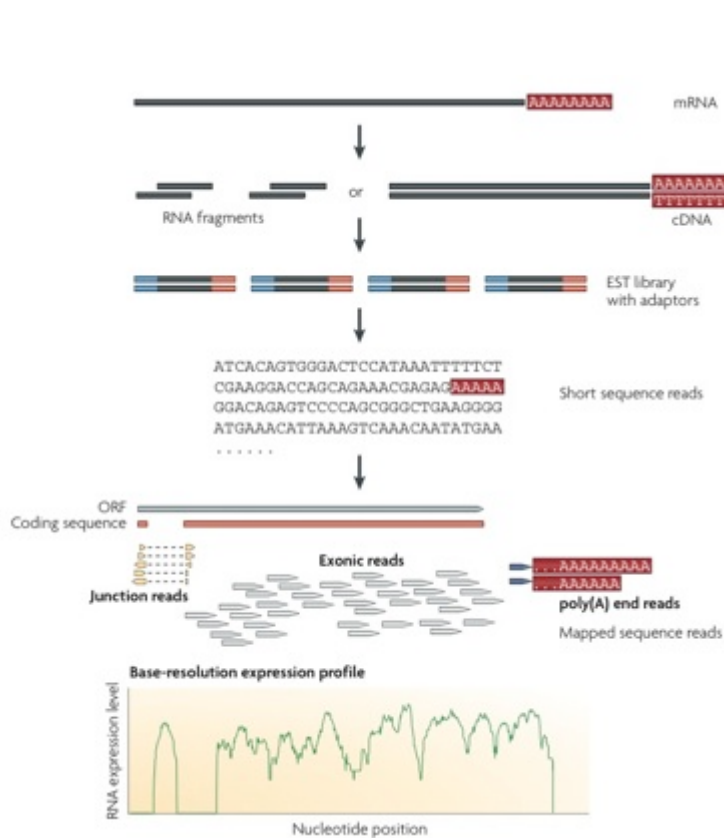


- 1000 mutations (Total 575693)
- 1000 coding mutations (Total 6172)

Chromatin immunoprecipitation – Sequencing (ChIP-Seq)



RNA-Seq: digital expression and much more




Wang et al. *Nat. Rev. Genet.*, 2009

Some of the ENCODE ChIP-Seq

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

HAIB TFBS Track Settings [Preview](#) [Downloads](#) [Subtracks](#) [Description](#)

 **Transcription Factor Binding Sites by ChIP-seq from ENCODE/HAIB** ([▲ENC TF Binding](#))

Maximum display mode: [Reset to defaults](#)

Select views ([help](#)):
 Peaks

Select subtracks by cell line and factor:

*- All	Cell Line	GMI2878 (Tier 1)	HI-hESC (Tier 1)	K562 (Tier 1)	HeLa-S3 (Tier 2)	HepG2 (Tier 2)	A549	ECC-1	GMI2891	GMI2892	HCT-116	PANC-1	PFSK-1	SK-N-MC	SK-N-SH	SK-N-SH RA	T-47D	U87	Cell Line	All *-
Factor		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	Factor	
ATF3		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		<input type="checkbox"/>													ATF3	<input type="checkbox"/>
BATF		<input type="checkbox"/>																	BATF	<input type="checkbox"/>
BCL11A		<input type="checkbox"/>	<input type="checkbox"/>																BCL11A	<input type="checkbox"/>
BCL3		<input type="checkbox"/>		<input type="checkbox"/>															BCL3	<input type="checkbox"/>
BCLAF1 (SC-101388)		<input type="checkbox"/>		<input type="checkbox"/>															BCLAF1 (SC-101388)	<input type="checkbox"/>
BHLHE40		<input type="checkbox"/>				<input type="checkbox"/>													BHLHE40	<input type="checkbox"/>
CTCF		<input type="checkbox"/>														<input type="checkbox"/>			CTCF	<input type="checkbox"/>
CTCF (SC-5916)		<input type="checkbox"/>	<input type="checkbox"/>			<input type="checkbox"/>	<input type="checkbox"/>										<input type="checkbox"/>		CTCF (SC-5916)	<input type="checkbox"/>
CTCF (SC-98982)		<input type="checkbox"/>		<input type="checkbox"/>														<input type="checkbox"/>	CTCF (SC-98982)	<input type="checkbox"/>
E2F6 (SC-22823)		<input type="checkbox"/>		<input type="checkbox"/>															E2F6 (SC-22823)	<input type="checkbox"/>
EBF1 (SC-137065)		<input type="checkbox"/>	<input type="checkbox"/>																EBF1 (SC-137065)	<input type="checkbox"/>
Egr-1		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>															Egr-1	<input type="checkbox"/>
ELF1 (SC-631)		<input type="checkbox"/>		<input type="checkbox"/>		<input type="checkbox"/>													ELF1 (SC-631)	<input type="checkbox"/>
ERalpha a		<input type="checkbox"/>						<input type="checkbox"/>									<input type="checkbox"/>		ERalpha a	<input type="checkbox"/>
ETS1		<input type="checkbox"/>		<input type="checkbox"/>														<input type="checkbox"/>	ETS1	<input type="checkbox"/>
FOSL1 (SC-183)		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>															FOSL1 (SC-183)	<input type="checkbox"/>
FOSL2		<input type="checkbox"/>				<input type="checkbox"/>													FOSL2	<input type="checkbox"/>
FOXA1 (SC-6553)		<input type="checkbox"/>				<input type="checkbox"/>		<input type="checkbox"/>										<input type="checkbox"/>	FOXA1 (SC-6553)	<input type="checkbox"/>
FOXA1 (SC-101058)		<input type="checkbox"/>				<input type="checkbox"/>													FOXA1 (SC-101058)	<input type="checkbox"/>
FOXA2 (SC-6554)		<input type="checkbox"/>				<input type="checkbox"/>													FOXA2 (SC-6554)	<input type="checkbox"/>
FOXP2		<input type="checkbox"/>											<input type="checkbox"/>	<input type="checkbox"/>					FOXP2	<input type="checkbox"/>
GABP		<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>													GABP	<input type="checkbox"/>
GATA2 (SC-267)		<input type="checkbox"/>		<input type="checkbox"/>															GATA2 (SC-267)	<input type="checkbox"/>
GATA3 (SC-268)		<input type="checkbox"/>															<input type="checkbox"/>		GATA3 (SC-268)	<input type="checkbox"/>
GR		<input type="checkbox"/>					<input type="checkbox"/>	<input type="checkbox"/>											GR	<input type="checkbox"/>
HDAC2 (SC-6296)		<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>		<input type="checkbox"/>													HDAC2 (SC-6296)	<input type="checkbox"/>
HEY1		<input type="checkbox"/>		<input type="checkbox"/>		<input type="checkbox"/>													HEY1	<input type="checkbox"/>
HNF4A (SC-8987)		<input type="checkbox"/>				<input type="checkbox"/>													HNF4A (SC-8987)	<input type="checkbox"/>
HNF4G (SC-6558)		<input type="checkbox"/>				<input type="checkbox"/>													HNF4G (SC-6558)	<input type="checkbox"/>
IRF4 (SC-6059)		<input type="checkbox"/>																	IRF4 (SC-6059)	<input type="checkbox"/>



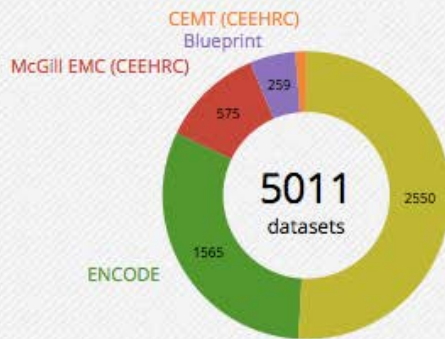
International Human Epigenome Consortium



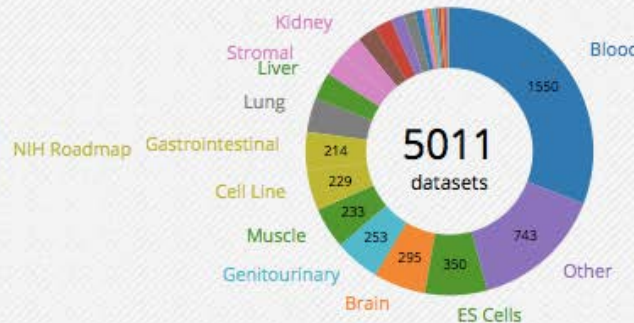
About Overview Data Grid Download Genome Browser IHEC Main Site

Welcome to the IHEC Data Portal

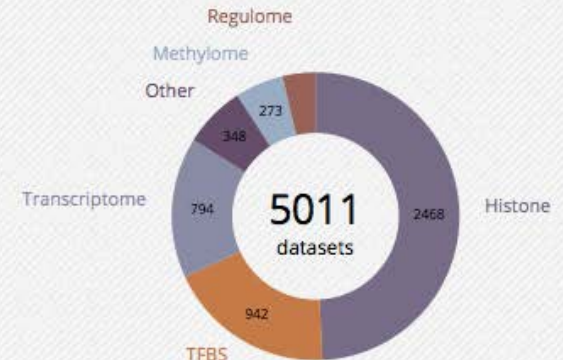
You may select IHEC datasets in these charts to view them in the Data Grid. Alternatively you can download or display them in a Genome Browser.



By Consortium



By Tissue



By Assay Category

View selected

View all

Reset

IHEC Data Grid

Data Grid

Track Hubs

<input type="checkbox"/>	Consortium	Datasets	Policy
<input checked="" type="checkbox"/>	McGill EMC (CEEHRC)	575	
<input type="checkbox"/>	CEMT (CEEHRC)	62	
<input checked="" type="checkbox"/>	Blueprint	259	
<input type="checkbox"/>	ENCODE	1565	
<input type="checkbox"/>	NIH Roadmap	2550	
<input checked="" type="checkbox"/>	Multiple Institutions		

	Histone								Methylome		Transcriptome		Regulome	
	H2A_Zac	H3K27ac	H3K27me3	H3K36me3	H3K4me1	H3K4me3	H3K9_14ac	H3K9me3	Input	WGB-Seq	RNA-Seq	mRNA-Seq	smRNA-Seq	DNase-Seq
Monocytes - Blood		14	12	14	16	14	1	14	20	10	82		2	5
B cell - Blood		1	1	1	2	2			2	1	35			
Brain - Brain											3			
Eosinophils - Blood										3				
Kidney - Kidney		1	1	1	1	1			1	1	1			
Kidney (Renal Carcinoma) - Kidney		1	1	1	1	1		1	1					
Leukemia CD19+CD10+ B Cells - Blood		1			1	1		1	2	2	2			
Naive CD4+ T Cells - Blood											29	39		
Renal Cancer - Kidney									1	1				
Skeletal Muscle (Control) - Muscle		3	2	3	1	3		3	3	7	14			
Skeletal Muscle (Mitochondrial Disease) - Muscle		5	5	5	4	4		3	5	8	14			
T cells - Blood		15	16	16	18	19		14	22	16	66		4	
CD4+_naive - Blood		2	2	2	2	2		2	2					
CD8+_naive - Blood		1	2	2	2	2		2	2	1	1			
CLP - Blood											3			
CMP - Blood											3			

[Click here for instructions.](#)

Visualize in Genome Browser

Get track hub link

Download tracks

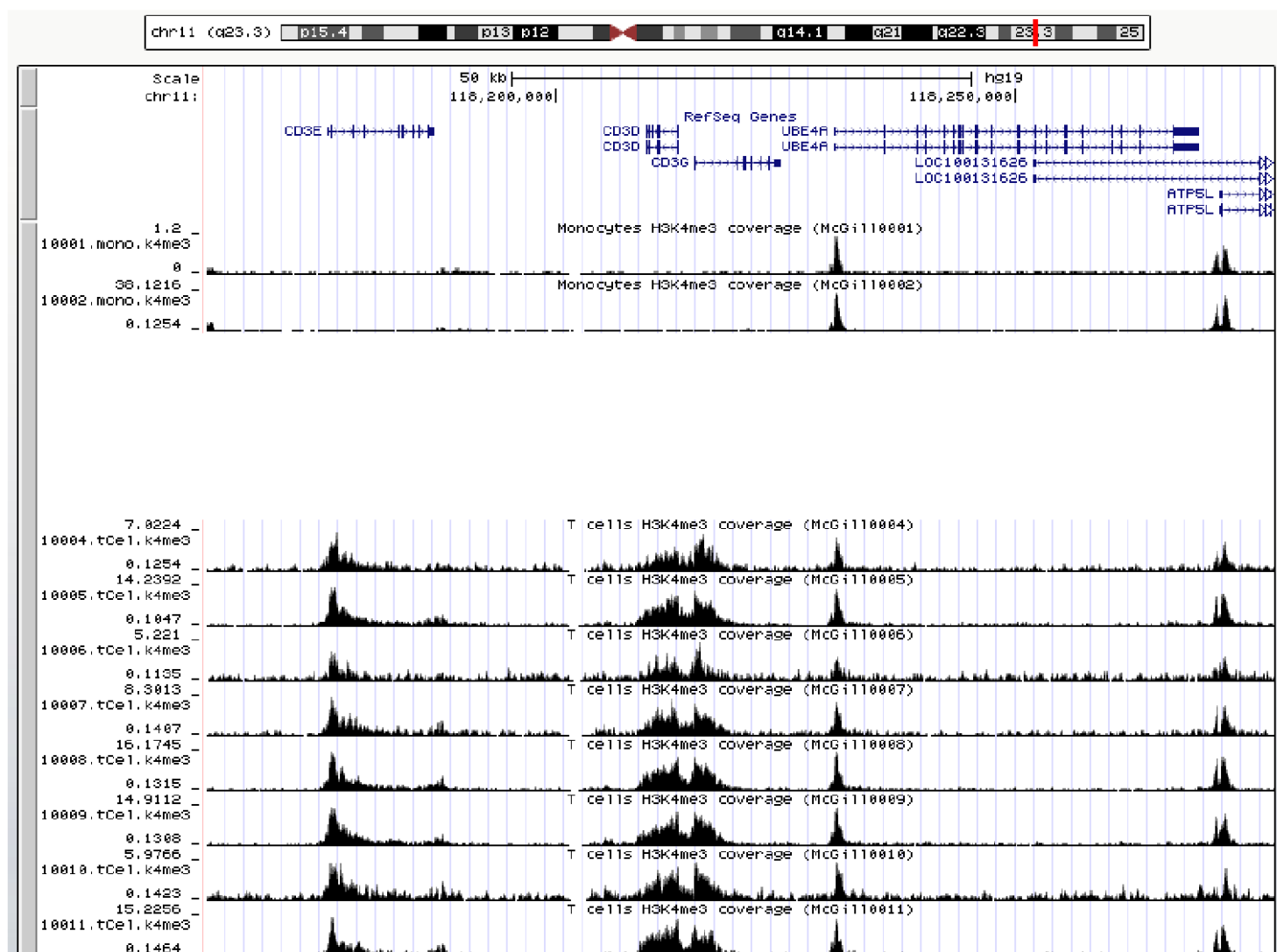
Reset

Order: by Consortium

Outline

- Applications of next-generation sequencing
- Functional genomics
- Example of machine learning approaches in functional genomics (work by Toby Hocking)
- Role of transposable elements in gene regulation

Motivation: visual differences between ChIP-seq profiles



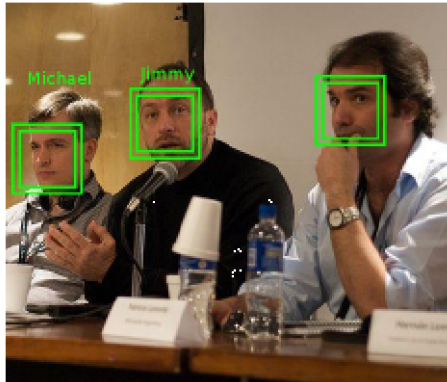
Challenges with peak calling

- Different types of profiles (narrow peaks or broad signal)
- Many peak caller algorithms...
- Lots of parameters...

How to choose the right peak caller and optimal parameters?

Previous work in computer vision: look and add labels to...

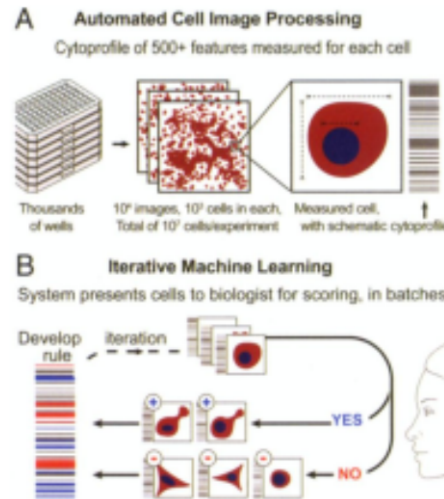
Photos



Labels: names

CVPR 2013
246 papers

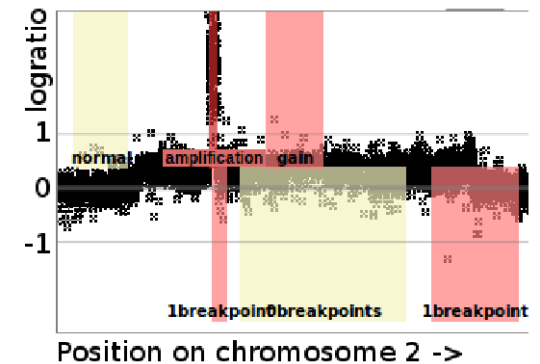
Cell images



phenotypes

CellProfiler
873 citations

Copy number profiles

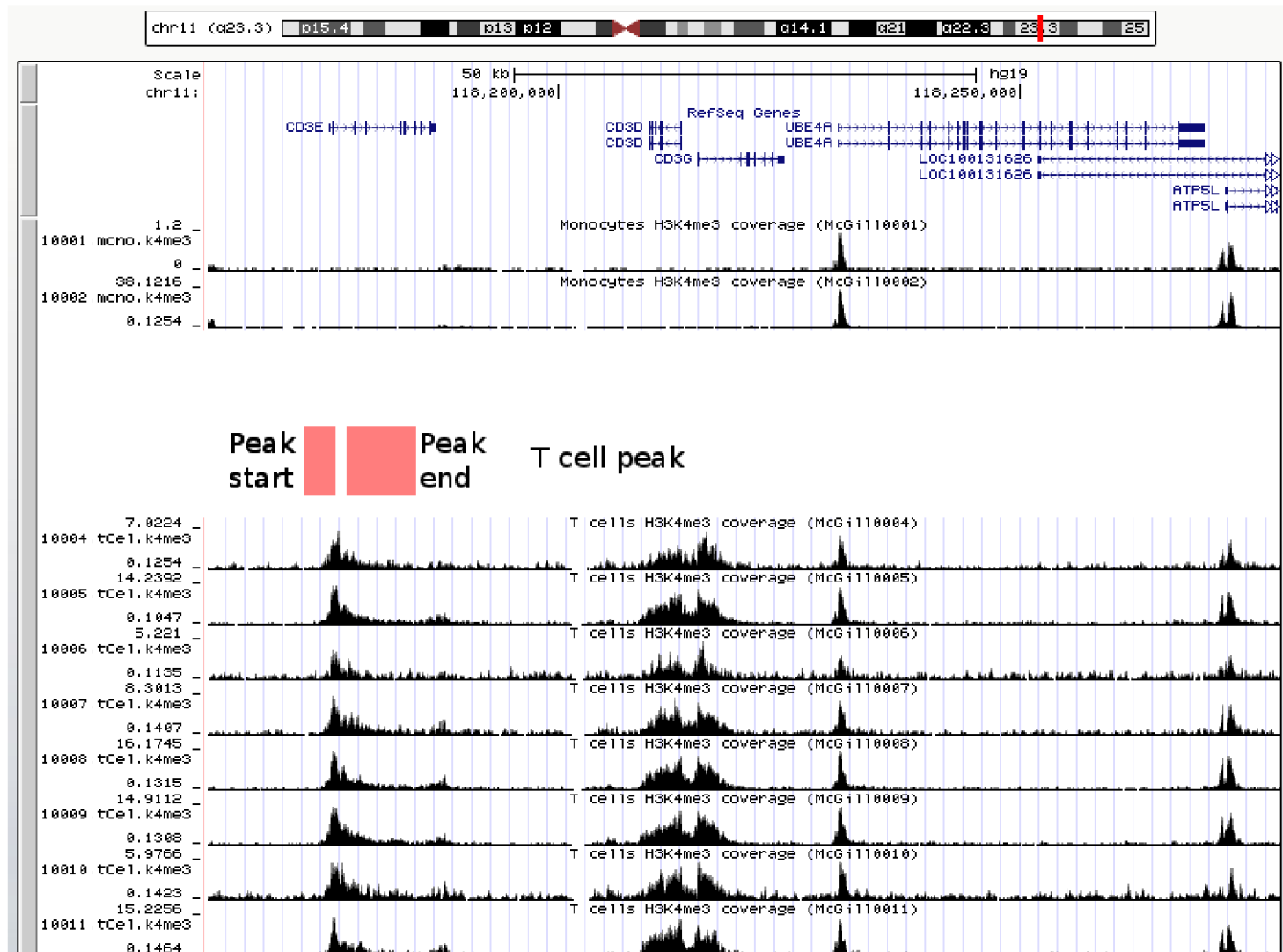


alterations

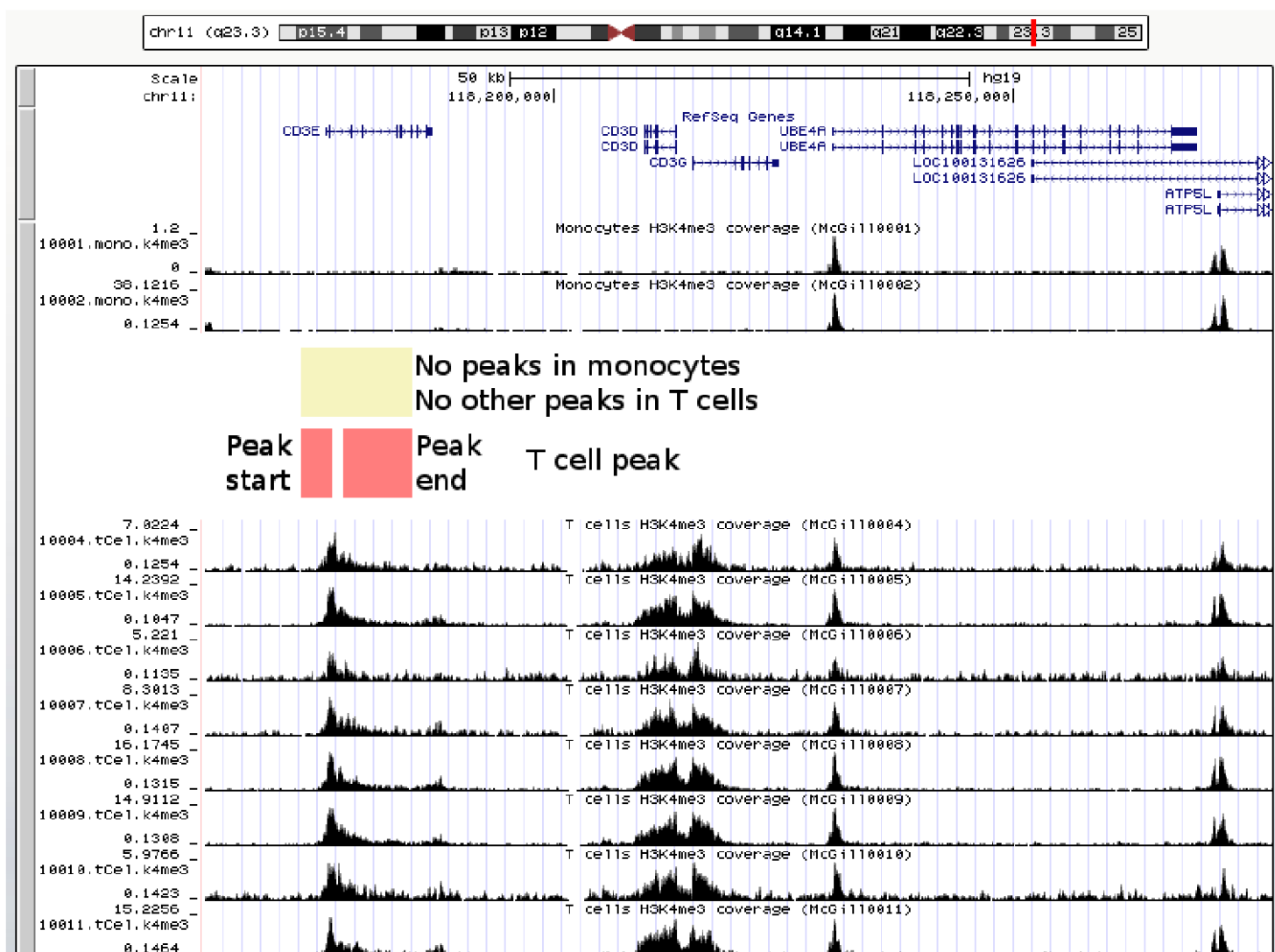
SegAnnDB
H, et. al. 2014.

Sources: http://en.wikipedia.org/wiki/Face_detection
Jones et al PNAS 2009. Scoring diverse cellular morphologies in image-based screens with iterative feedback and machine learning.

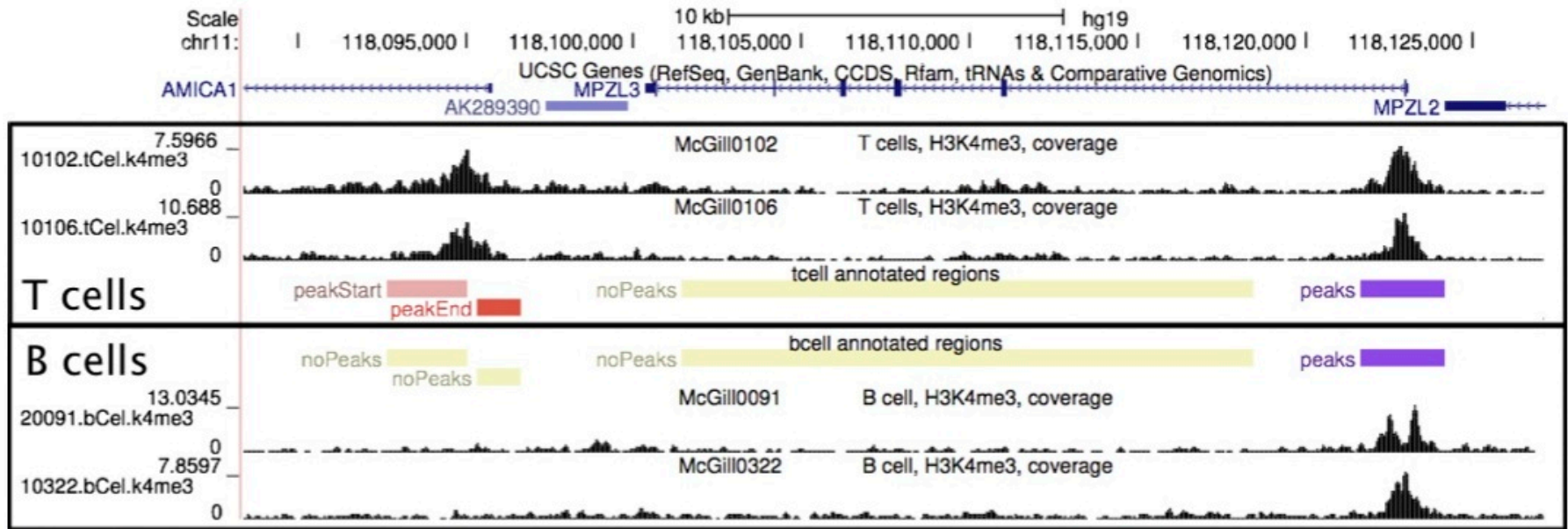
Annotated regions for possible peak starts and ends



Implicit specification of negative regions



Visual peak annotation



Supervised ML problem

A peak detection function or peak caller $c: \mathbb{R}^d \rightarrow \{0,1\}^d$ takes a coverage profile $\mathbf{x} \in \mathbb{R}^d$ as input, and returns a binary peak call prediction $\mathbf{y} = c(\mathbf{x}) \in \{0,1\}^d$ (0 is background noise, 1 is a peak).

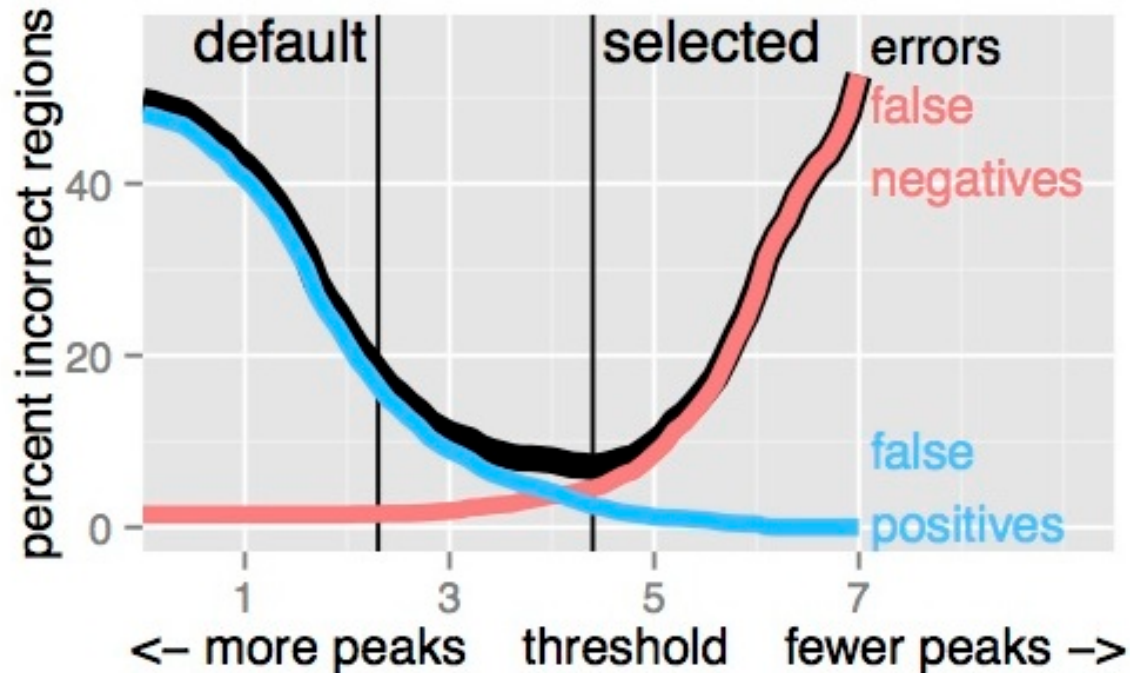
The goal is to learn how to call peaks $c(\mathbf{x}_i)$ which agree with the annotated regions $R_i^+, R_i^-, \underline{R}_i, \overline{R}_i$ for some test samples i . To quantify the error of the peak calls with respect to the annotation data, we define the annotation error as the sum of false positive (FP) and false negative (FN) regions:

$$E(\mathbf{y}, \underline{R}_i, \overline{R}_i, R_i^+, R_i^-) = \text{FP}(\mathbf{y}, \underline{R}_i, \overline{R}_i, R_i^-) + \text{FN}(\mathbf{y}, \underline{R}_i, \overline{R}_i, R_i^+). \quad (1)$$

The supervised machine learning problem can be formalized as the following optimization problem. Find the peak caller c with minimal annotation error on a set of test samples:

$$\underset{c}{\text{minimize}} \sum_{i \in \text{test}} E[c(\mathbf{x}_i), \underline{R}_i, \overline{R}_i, R_i^+, R_i^-]. \quad (2)$$

Choosing the threshold



$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} \sum_{i \in \{1, \dots, n\}} E[c_{\lambda}(\mathbf{x}_i), \underline{R}_i, \bar{R}_i, R_i^+, R_i^-].$$

Creating a benchmark dataset

mark	H3K36me3	H3K4me3	H3K36me3	H3K36me3	H3K4me3	H3K4me3	H3K4me3
annotator	TDH	TDH	AM	TDH	PGP	TDH	XJ
sample.set	other	other	immune	immune	immune	immune	immune
windows	4	29	23	4	30	27	12
noPeaks	72	536	752	230	1653	1656	702
peakStart	68	305	294	200	813	796	216
peakEnd	60	311	294	200	730	933	216
peaks		218	403		638	287	243
tcell			15	15	19	19	19
monocyte			5	5	6	6	6
bcell			1	1	2	2	2
kidney	1	1					
kidneyCancer	1	1					
skeletalMuscleCtrl	3	3					
skeletalMuscleMD	3	4					
leukemiaCD19CD10BCells		1					

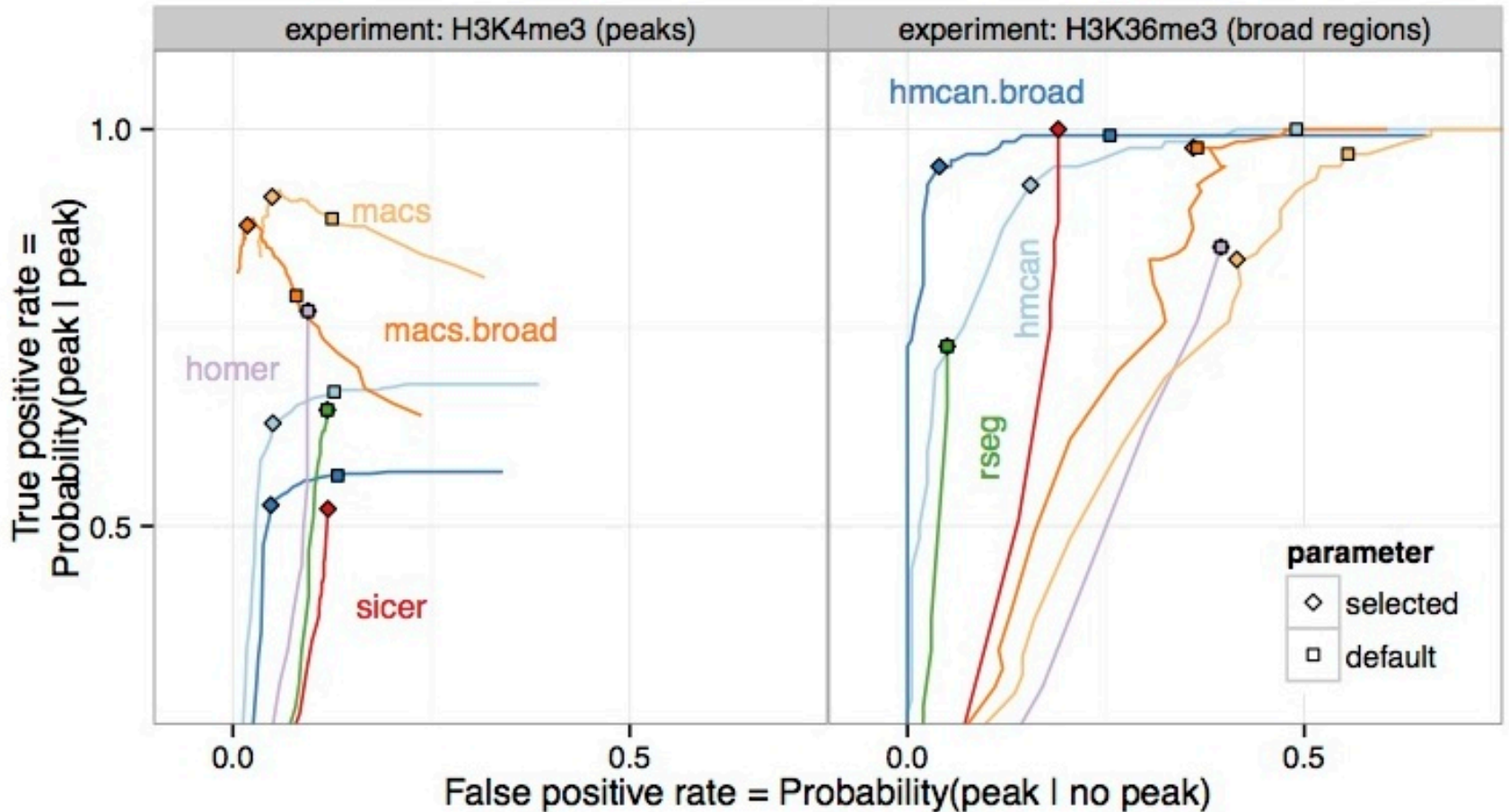
7 genomic regions, 2 histone marks, 4 expert annotators, 8 different cell types

Total of 12,286 annotated regions

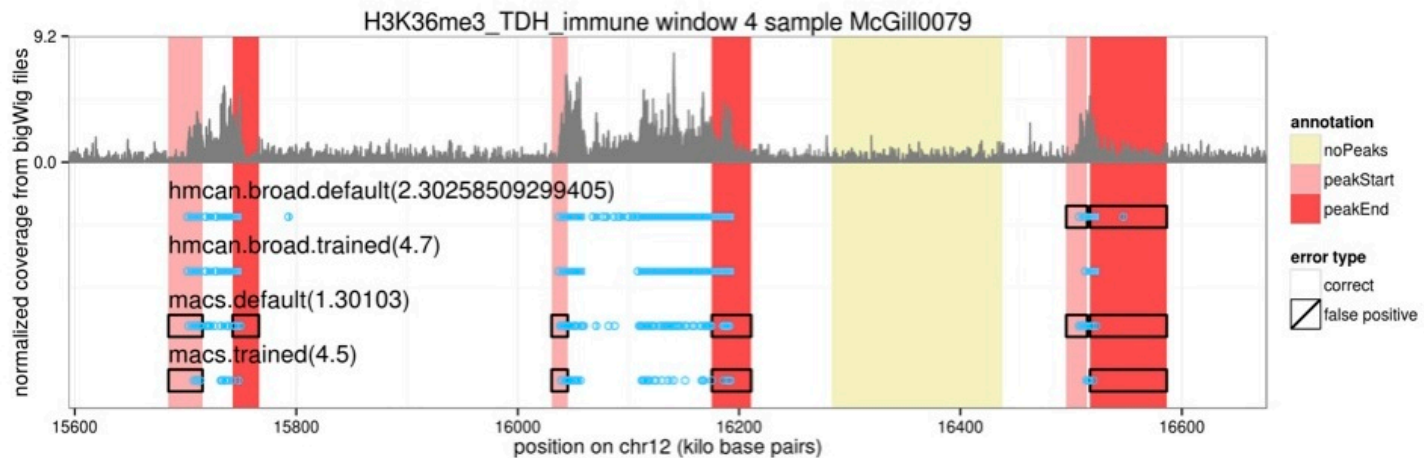
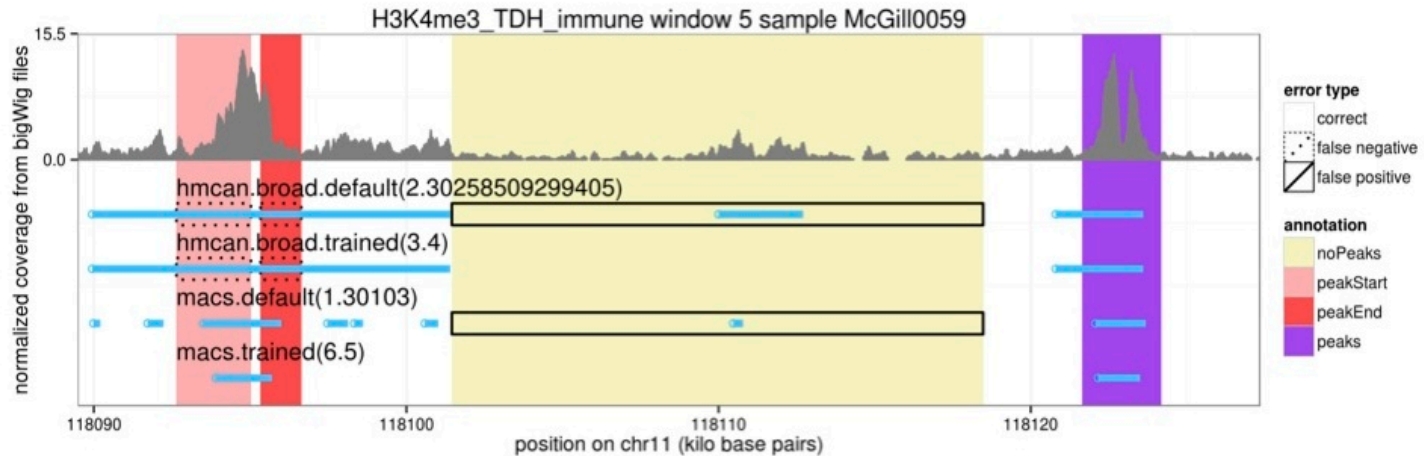
Training and testing

Same	Train	Test
Annotator, cell types	Chr1	Chr2
Annotator, regions	T cells	Kidney cells
Cell types, regions	Annotator 1	Annotator 2

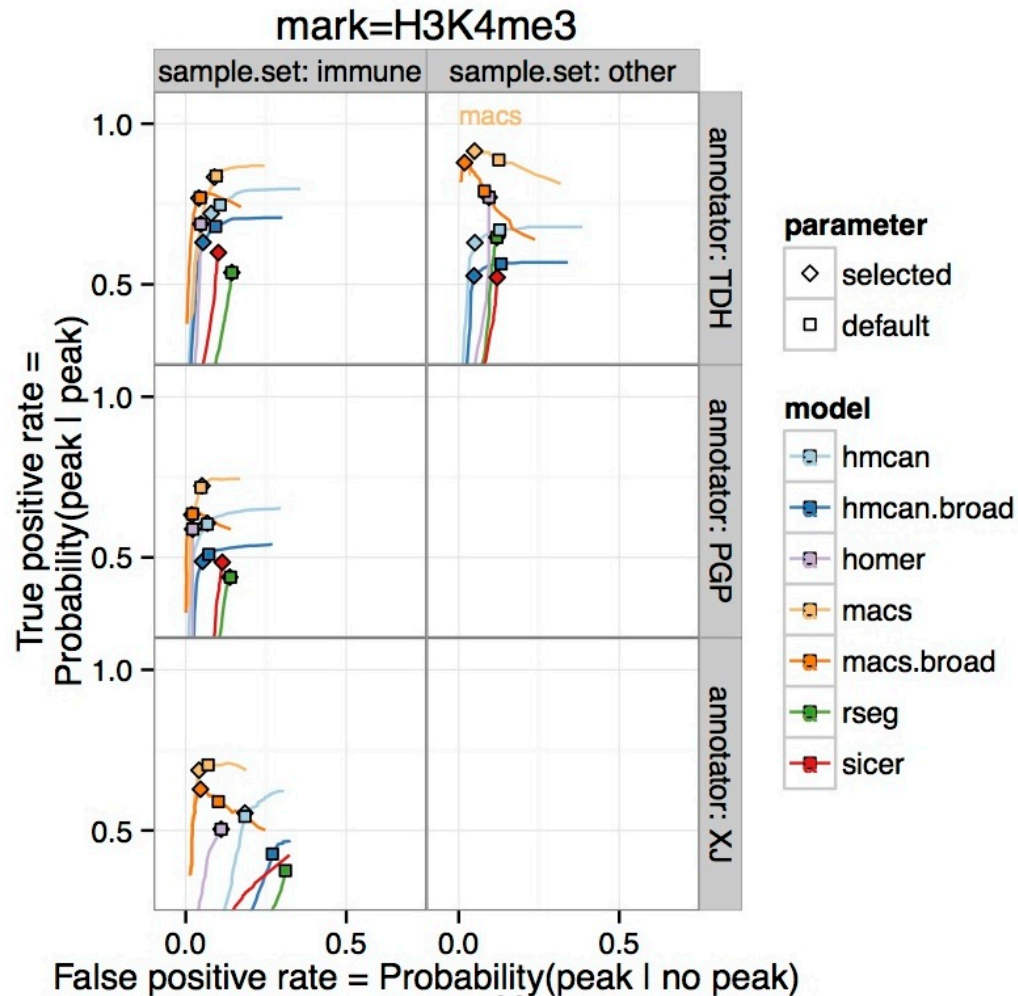
Performance of different peak callers



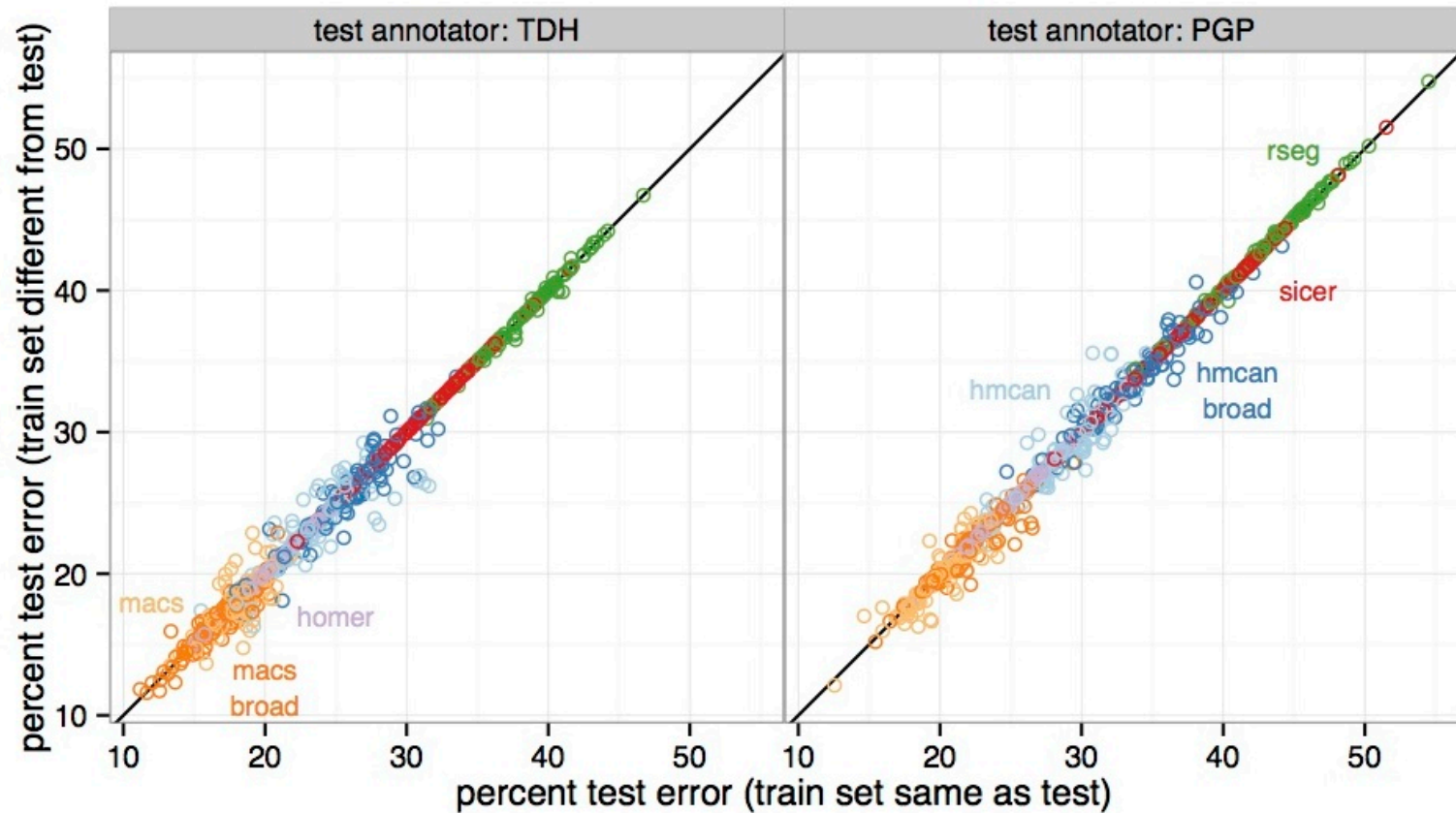
Example region



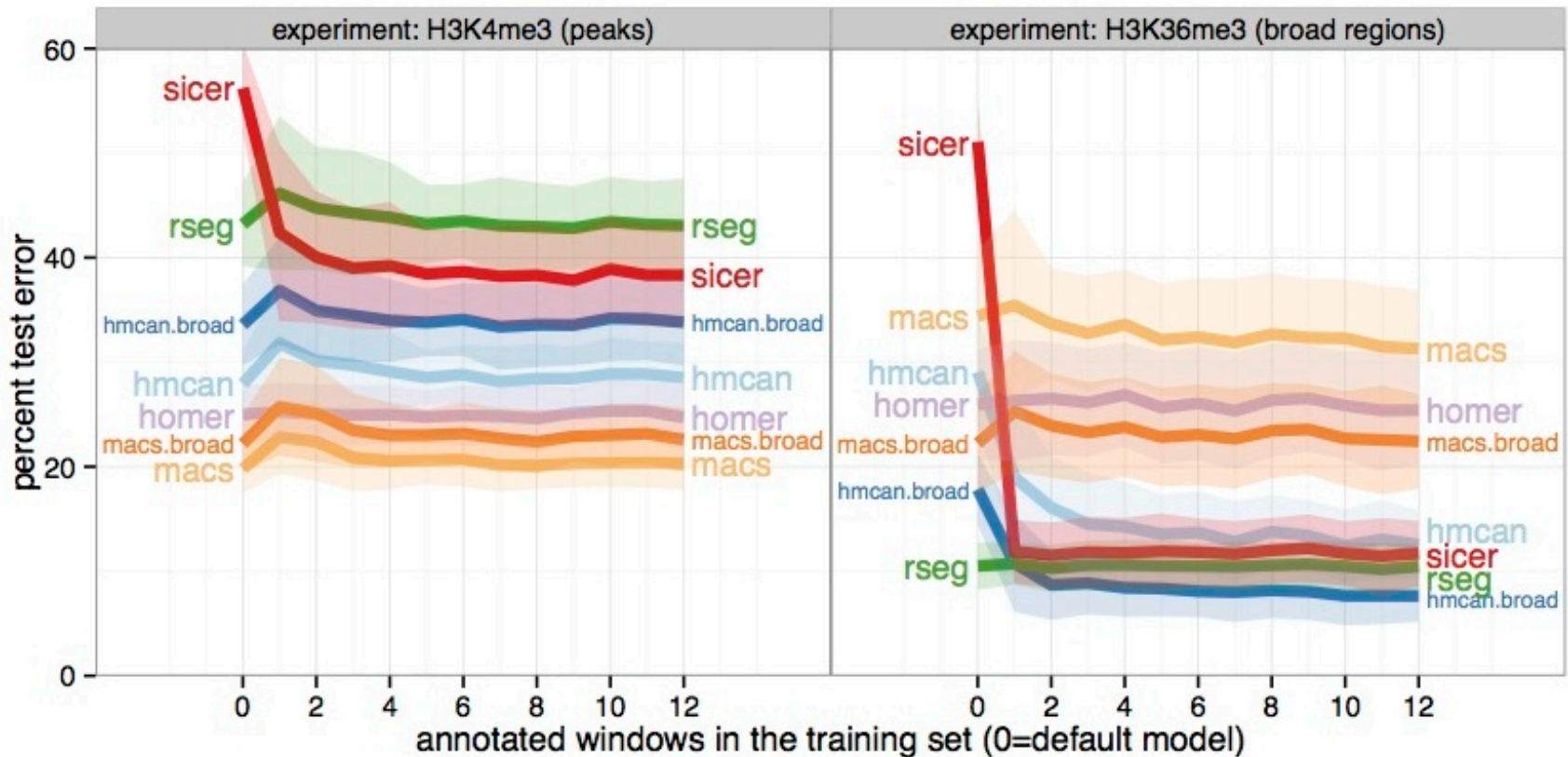
Performance is stable across annotation datasets



Test error comparable for different annotators



Impact of the size of the training set



Conclusions

- The macs algorithm showed the minimum train error of 9–20% across the 4 annotated “peak” data sets.
- The hmcan.broad algorithm showed the minimum train error of 7–20% across the 3 annotated “broad” data sets.
- Error of the different algorithms and model ordering was independent of the annotator
- Each algorithm quickly achieves its model-specific minimum error, after only about 4 annotated windows in the train set
- For some models the trained model parameters were clearly better than the default model parameters

Future directions

- Calling peaks is an “easy” problem
- The real question that we would like to address is:
Can we identify genetic variants that have an impact on the “phenotype”?
- This requires combining genetic and epigenomic data
- Previously, there wasn’t enough data the use of ML approaches to answer this question
- But now, we are generating matched genetic and epigenomic data on 100s of individuals...
- Can we learn from this data and build a predictive model?

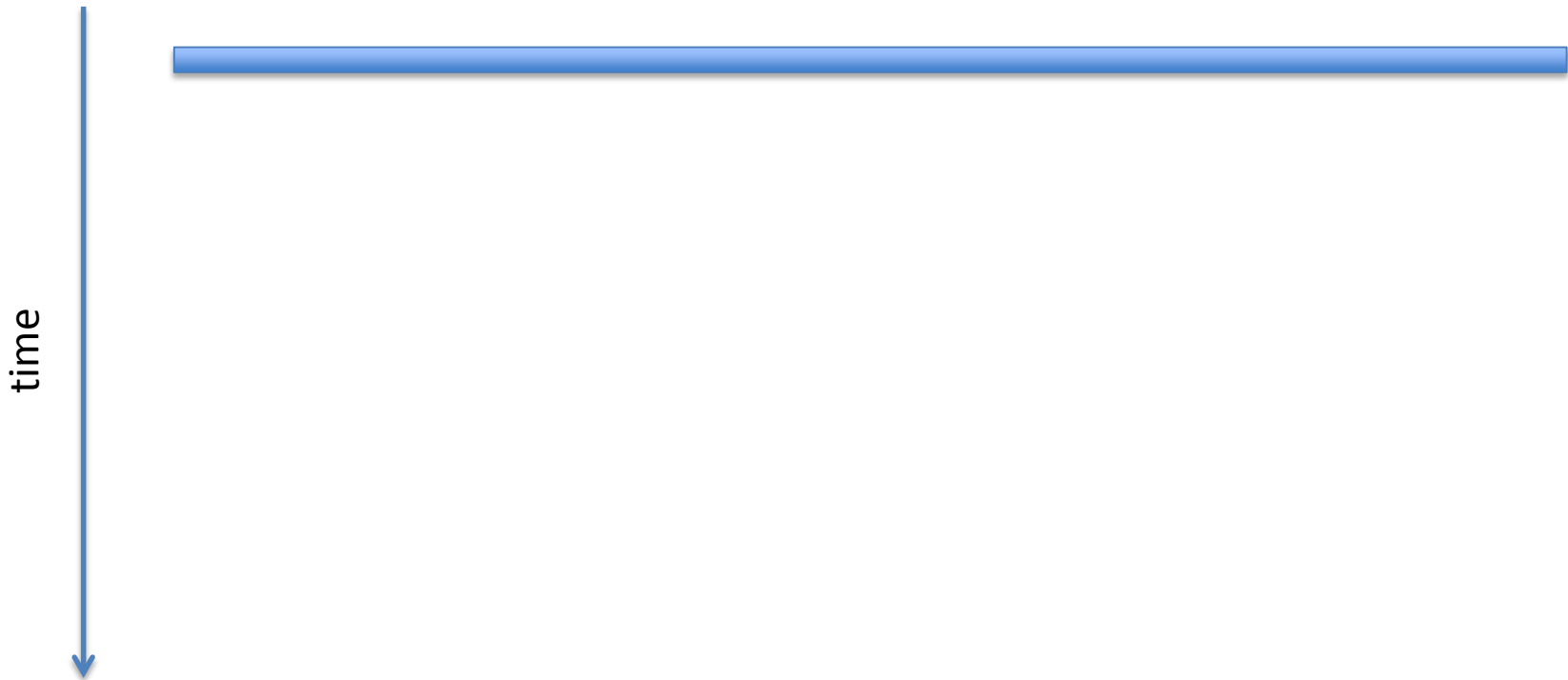
Outline

- Applications of next-generation sequencing
- Functional genomics
- Example of machine learning approaches in functional genomics
- Role of transposable elements in gene regulation

30-50% of every mammalian genome is
'Dark matter'



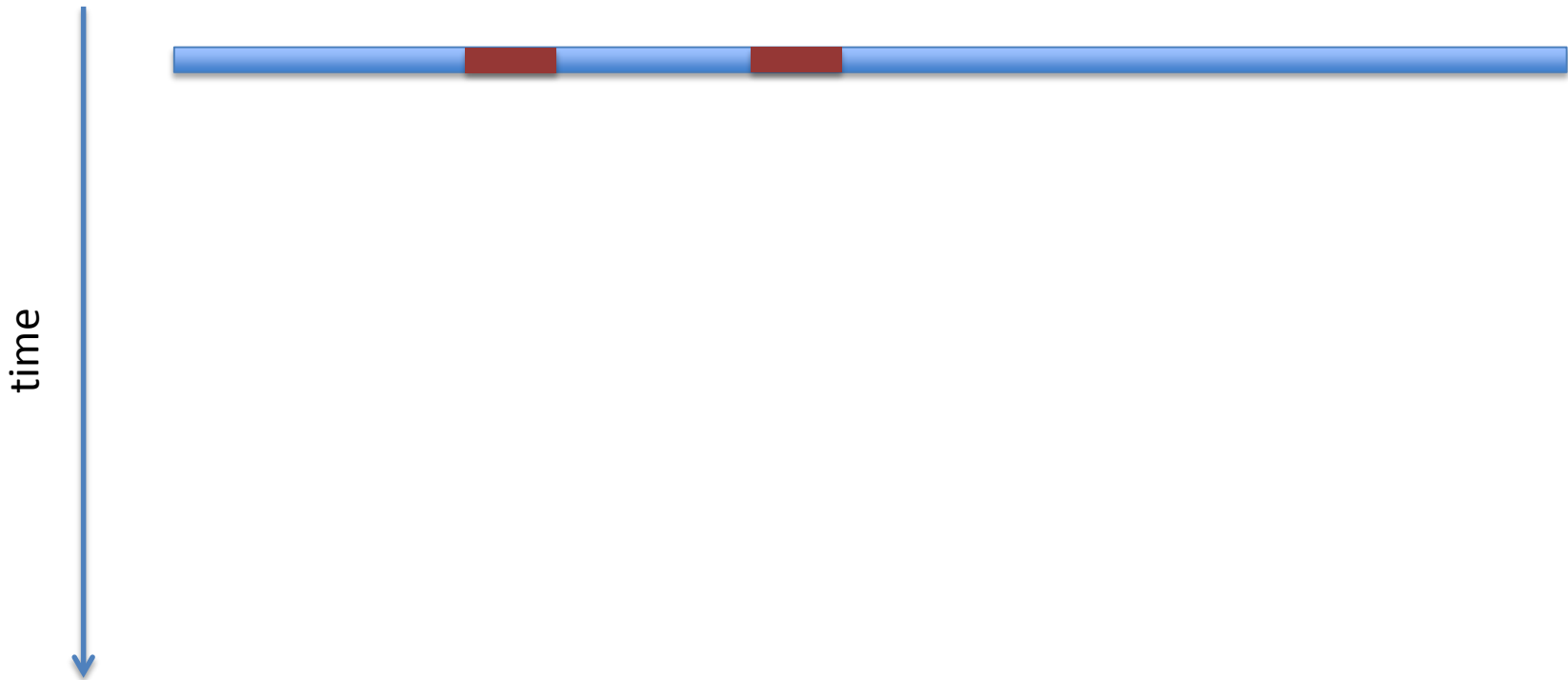
Transposable Elements



Transposable Elements



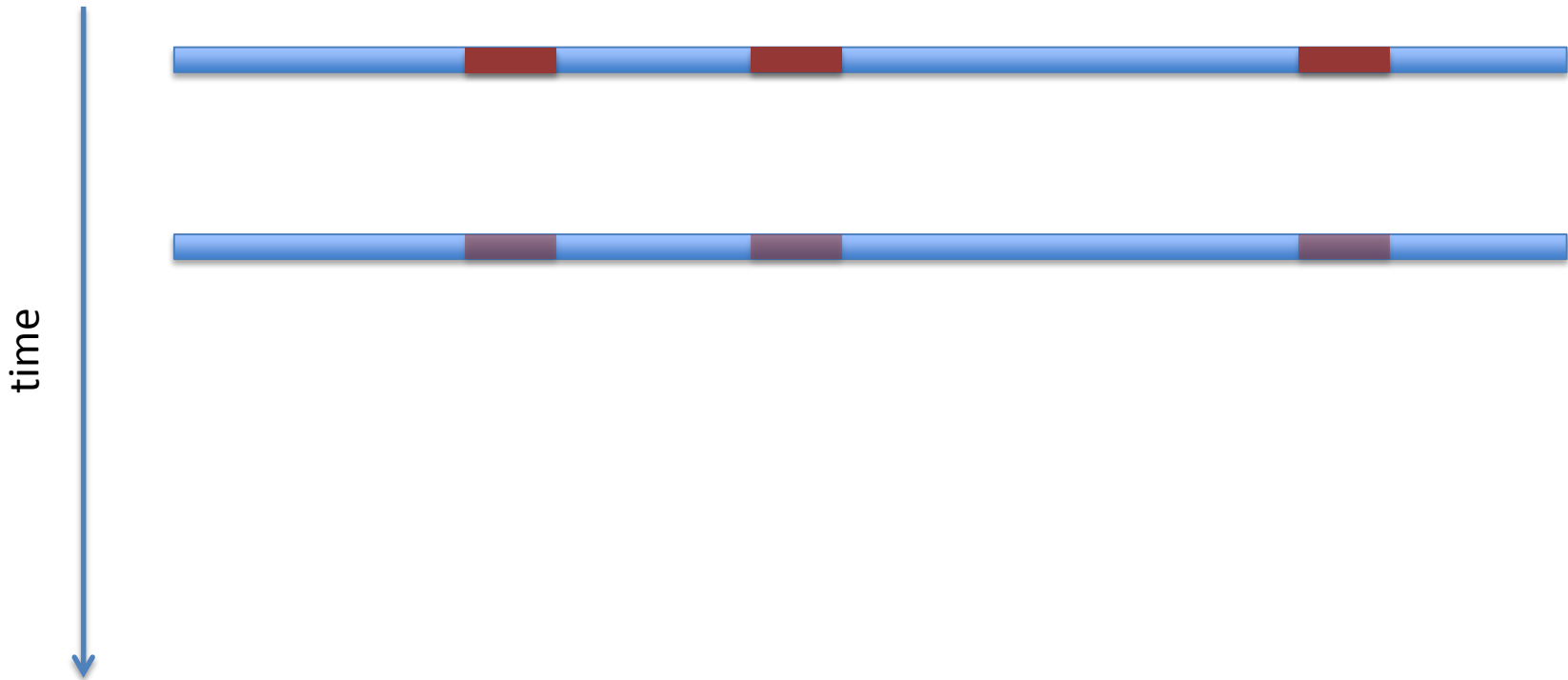
Transposable Elements



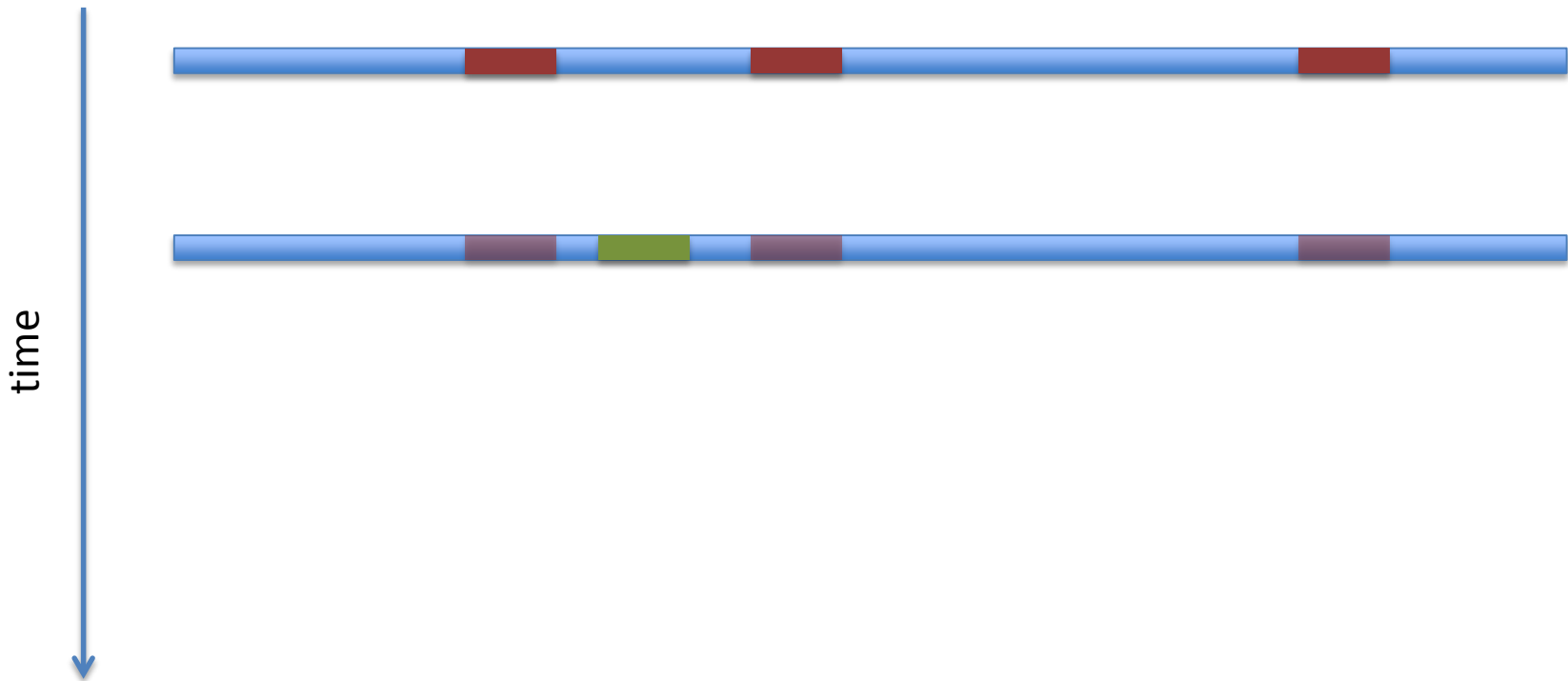
Transposable Elements



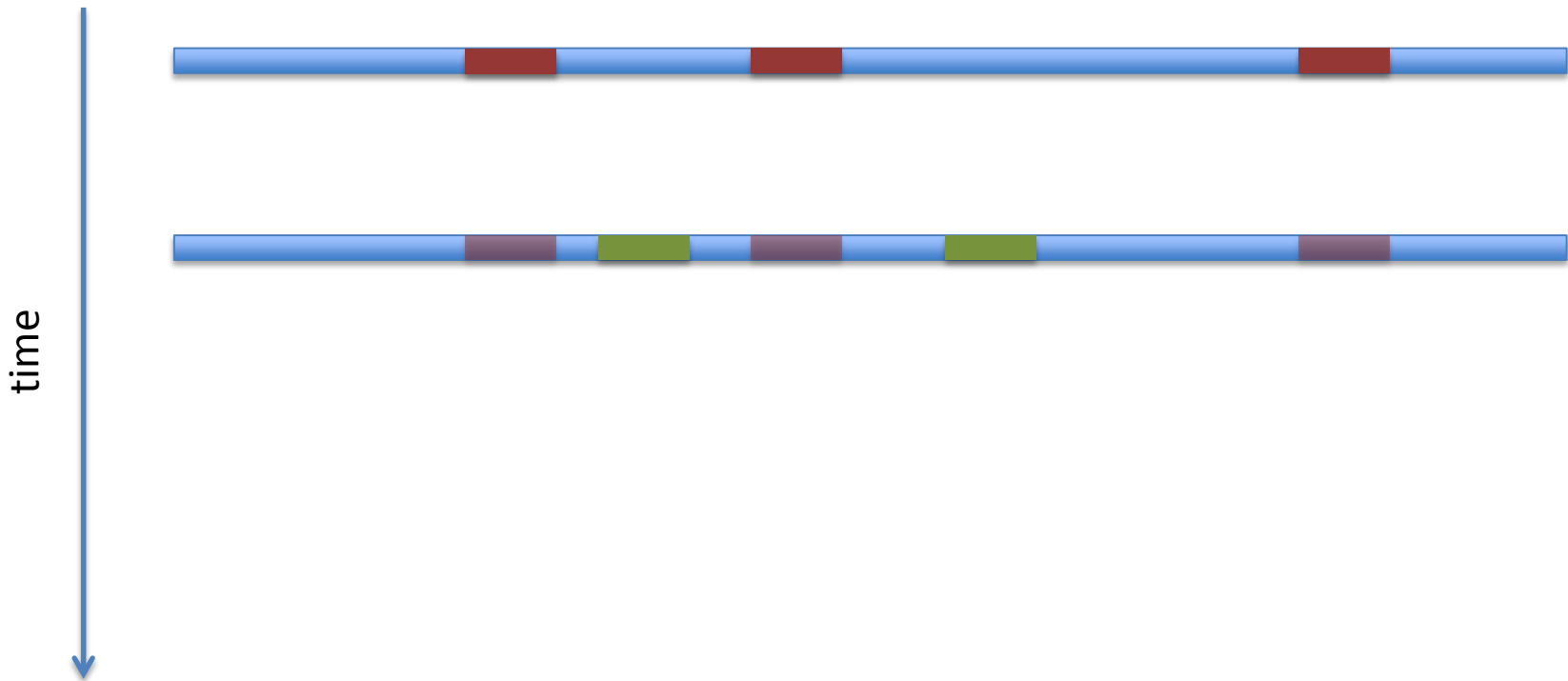
Transposable Elements



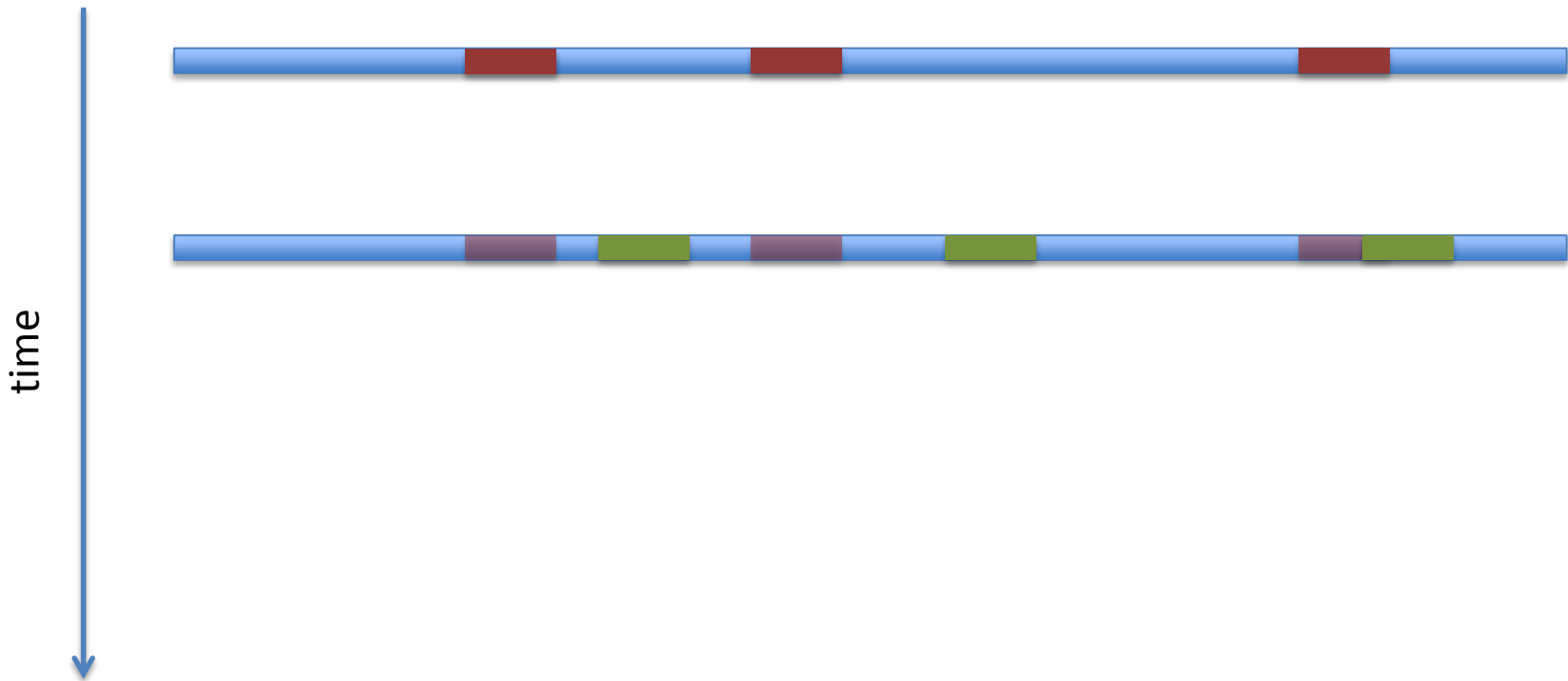
Transposable Elements



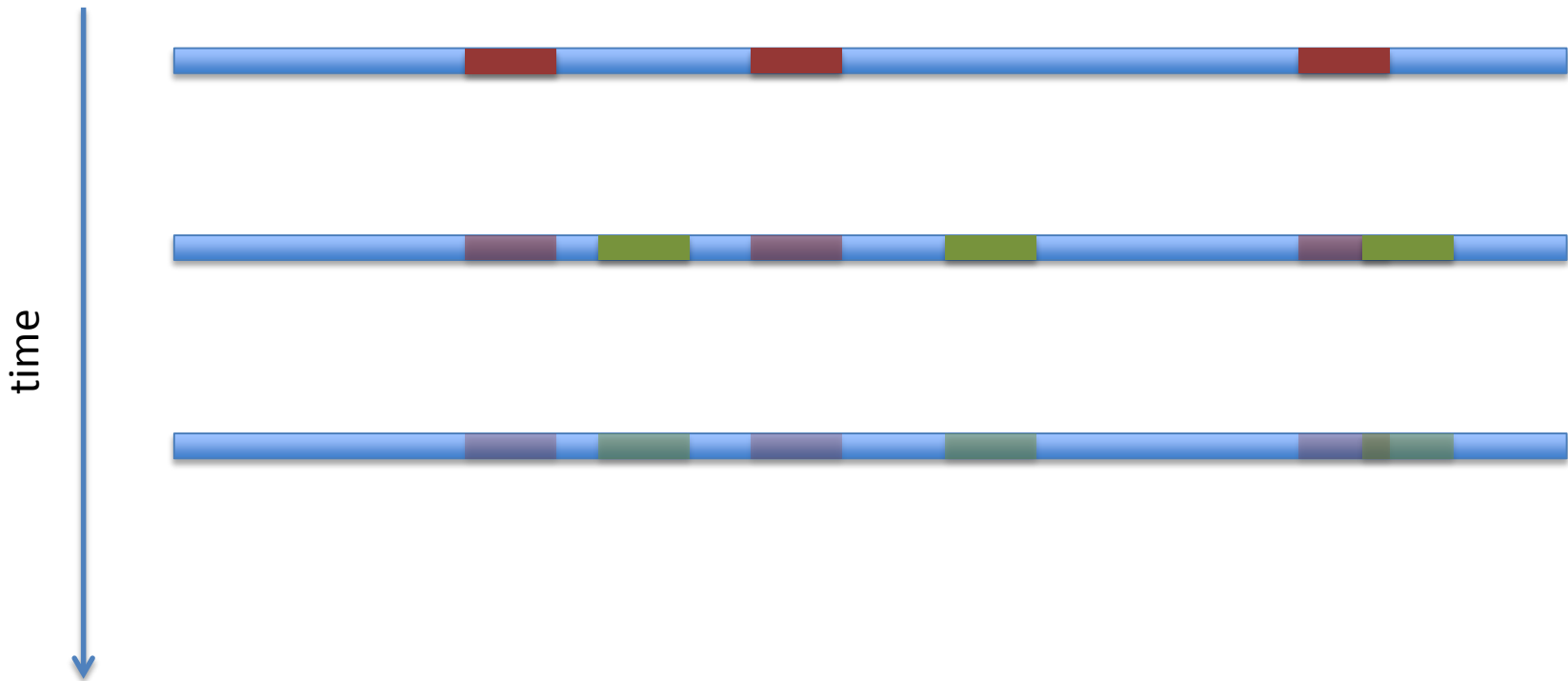
Transposable Elements



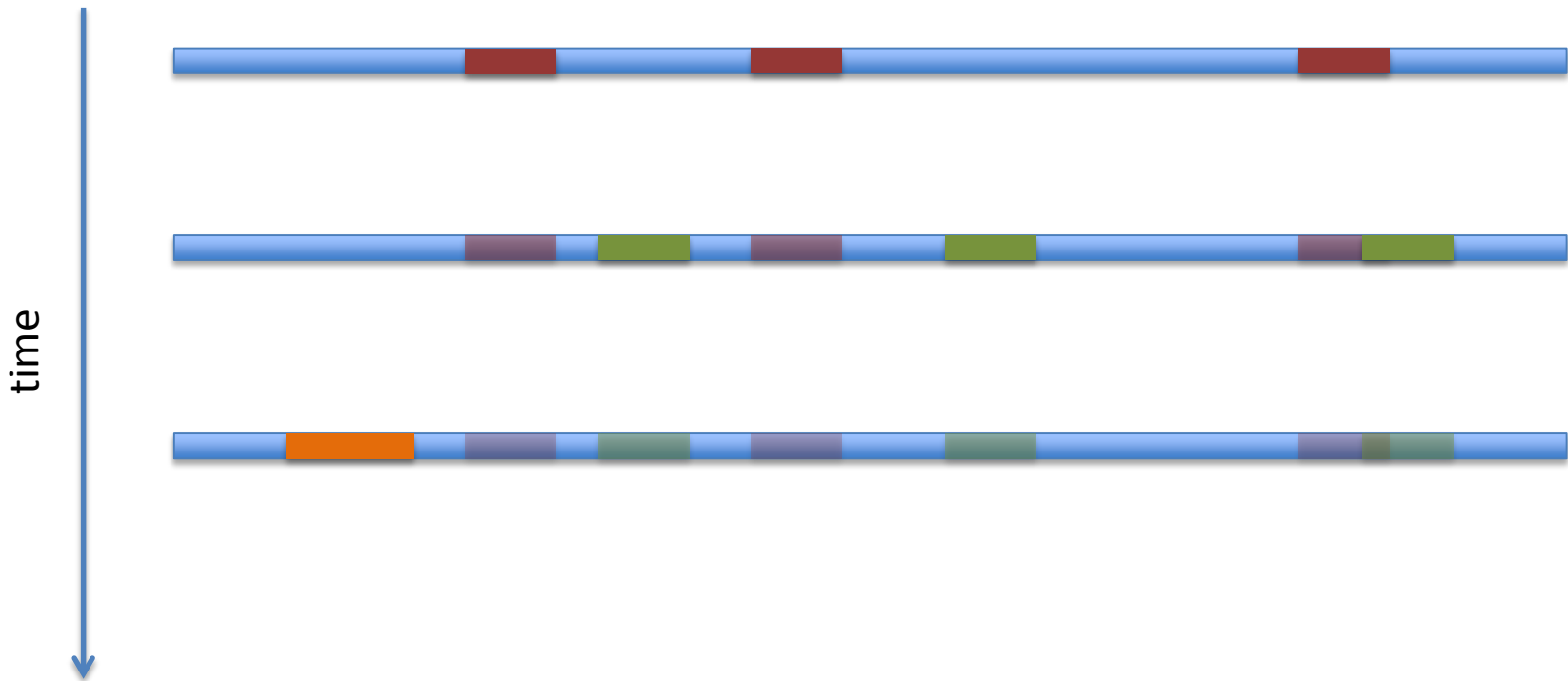
Transposable Elements



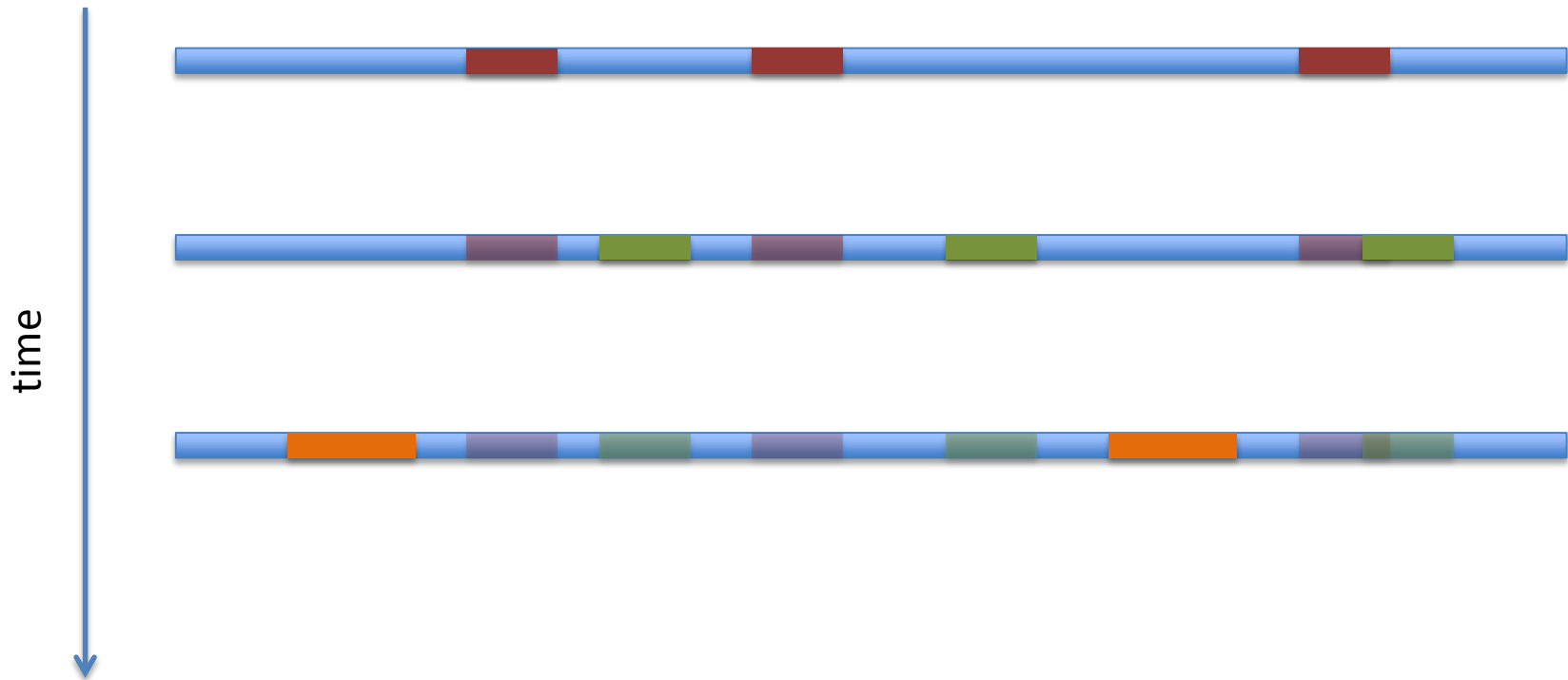
Transposable Elements



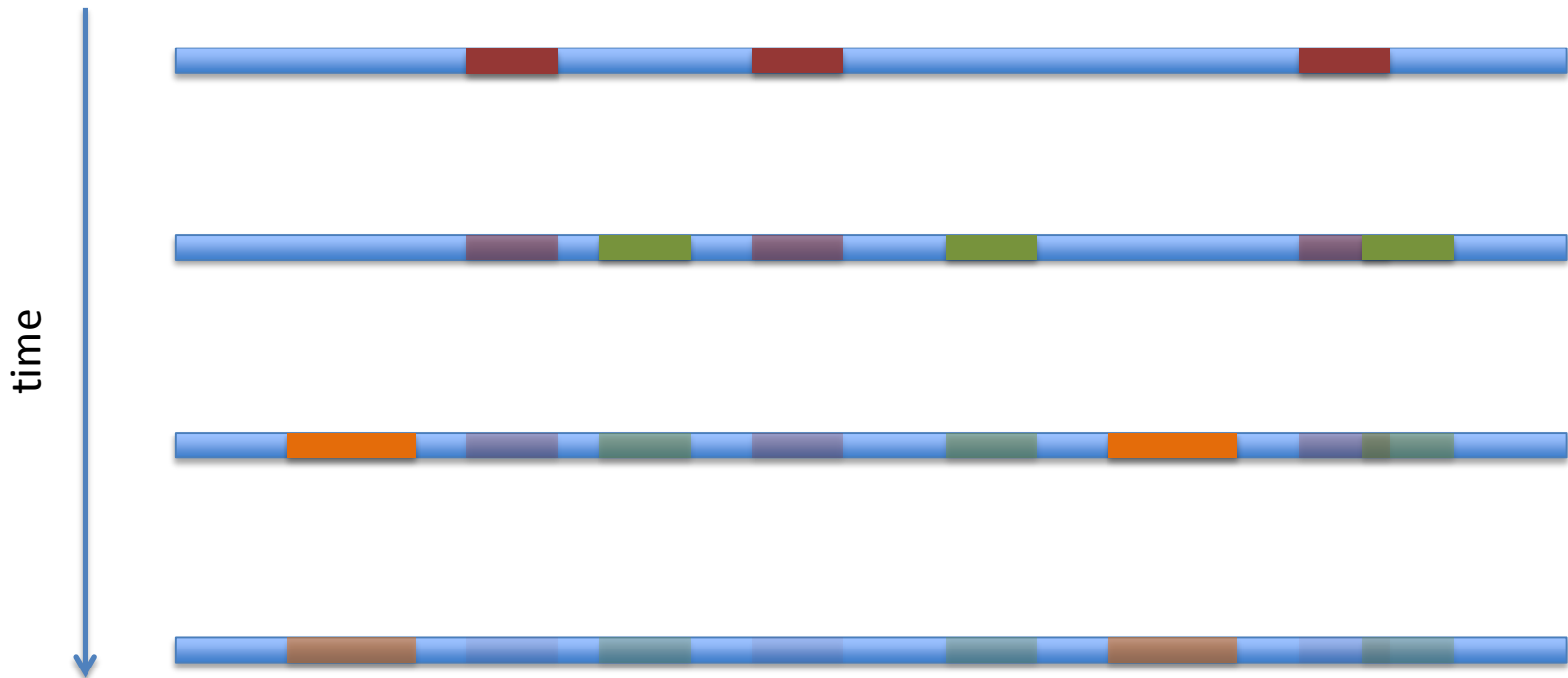
Transposable Elements



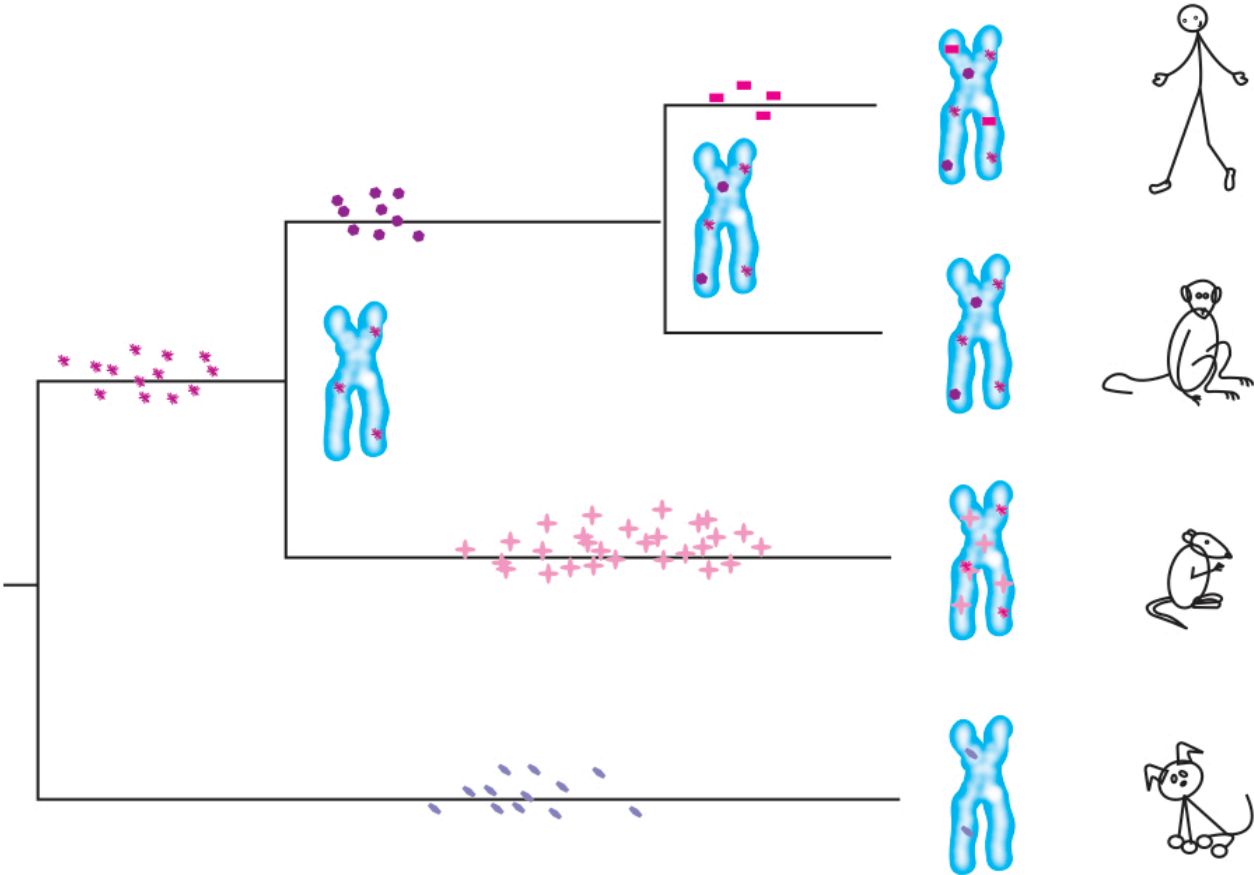
Transposable Elements



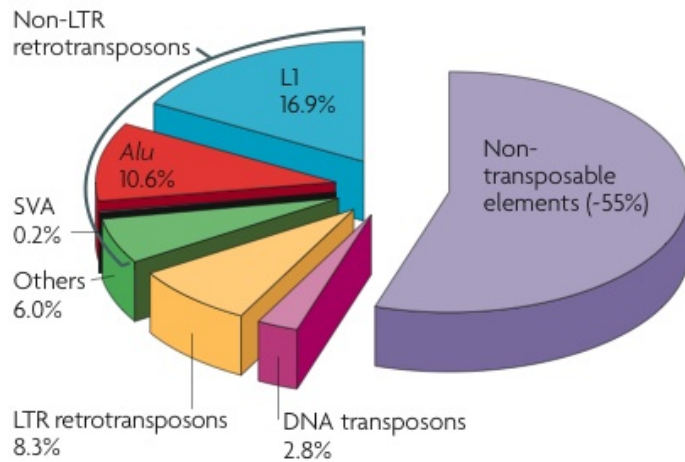
Transposable Elements



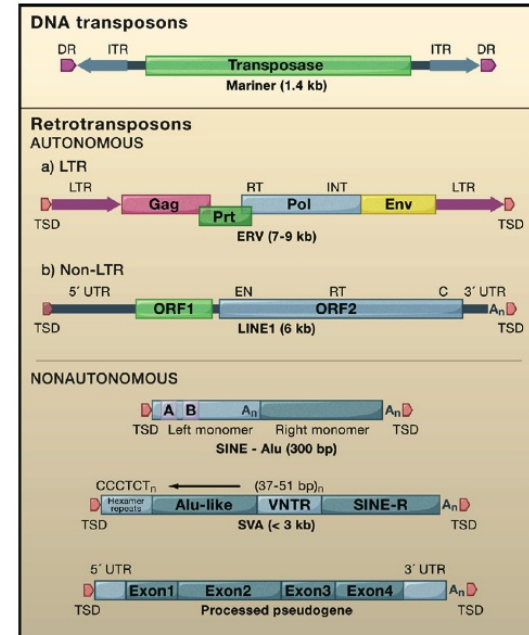
TEs have changed mammalian genomes



Transposable Elements (TEs) in the human genome



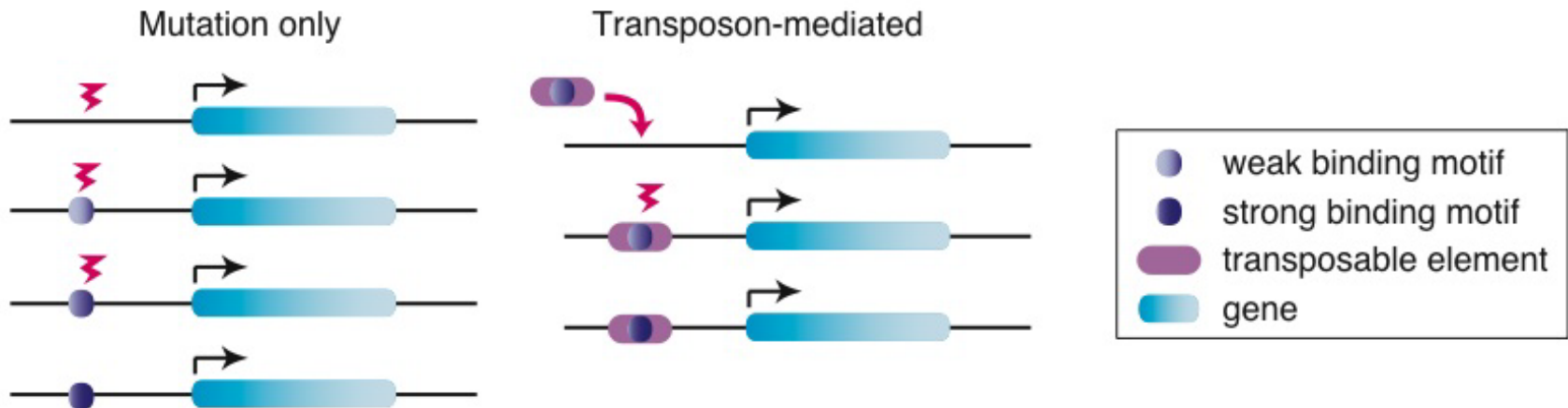
Cordeaux and Batzer,
Nat Rev Genet, 2009



Goodier and Kazazian,
Cell 2008

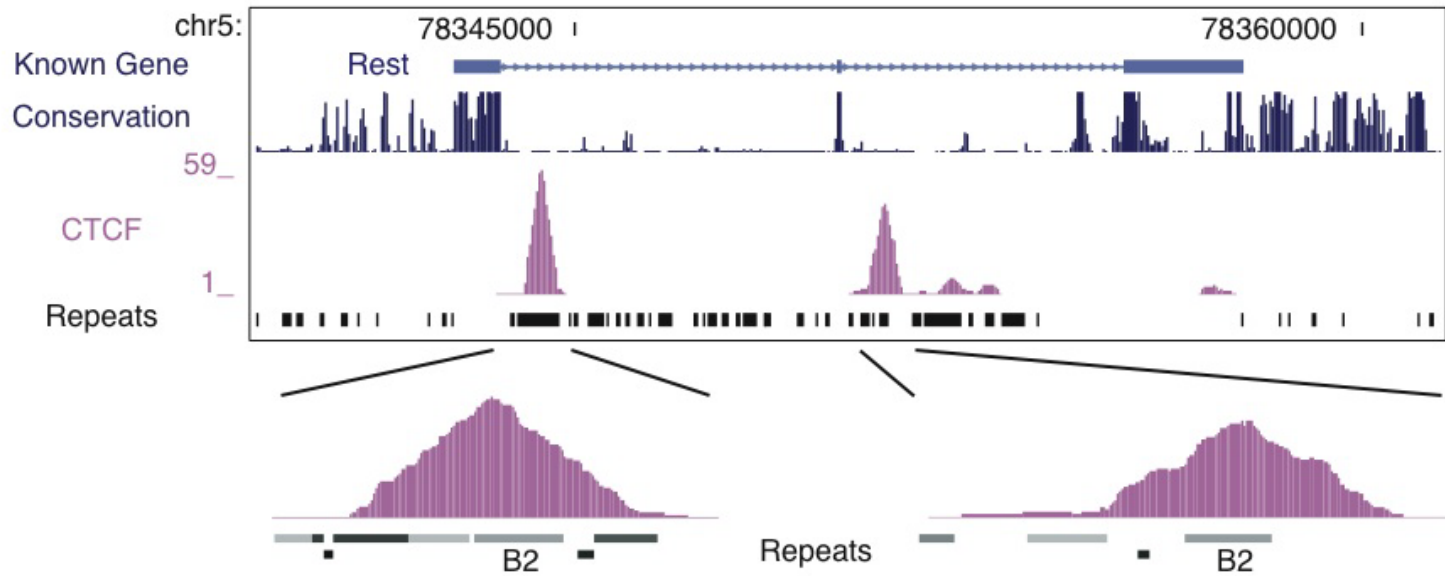
About 25% of the human genome consists of lineage-specific repeats.
For the mouse genome it's 30%.

Two models for regulatory site acquisition



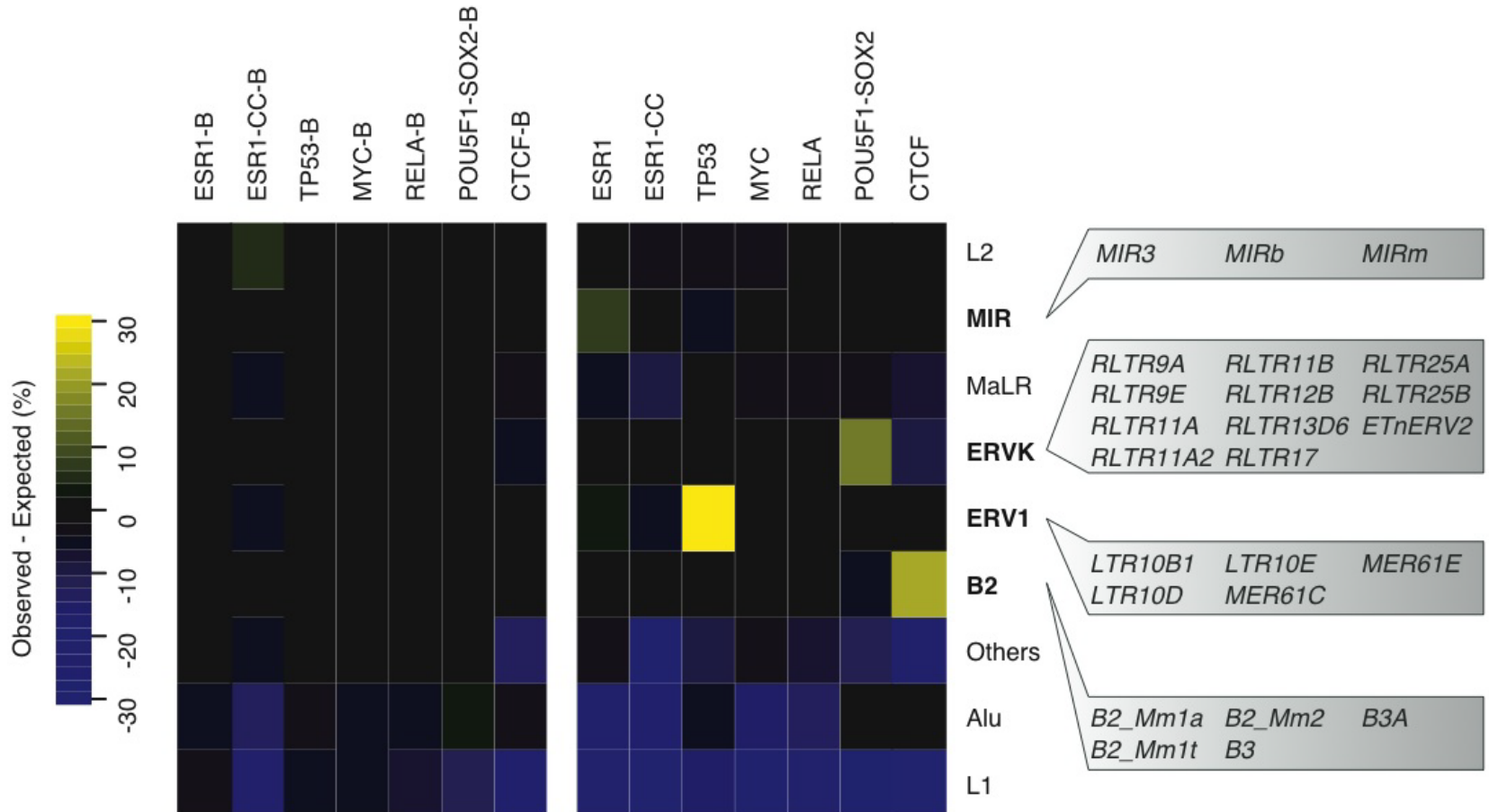
McClintock 1950, Britten & Davidson 1971, Brosius 1991, ...

Genome-wide occupancy maps reveal a strong association to repeats for many TFs



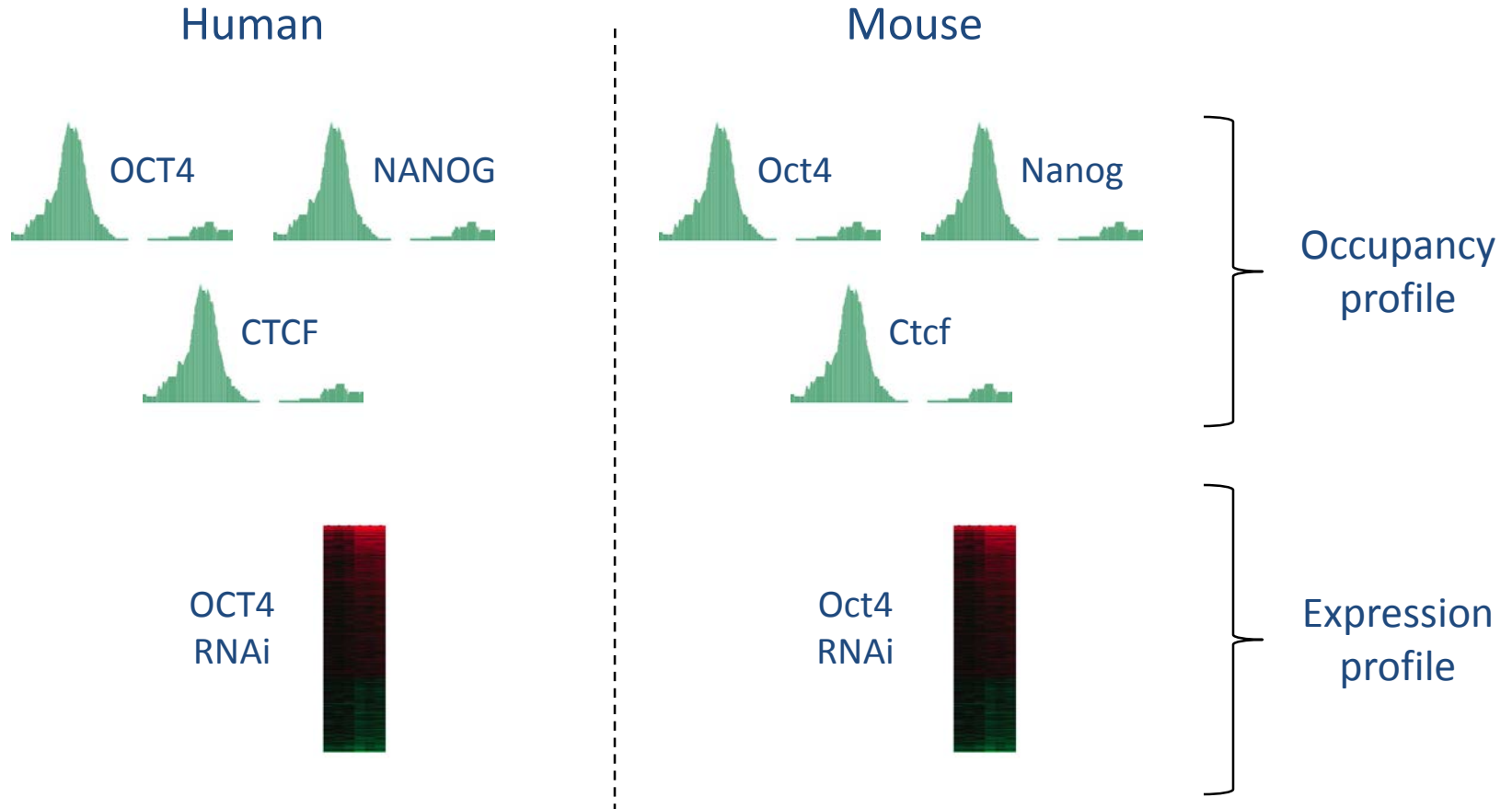
28% of CTCF sites
in B2 repeats

Repeat Associated Binding Sites (RABS)



Bourque et al. *Genome Res*, 2008

Human versus mouse ES cells

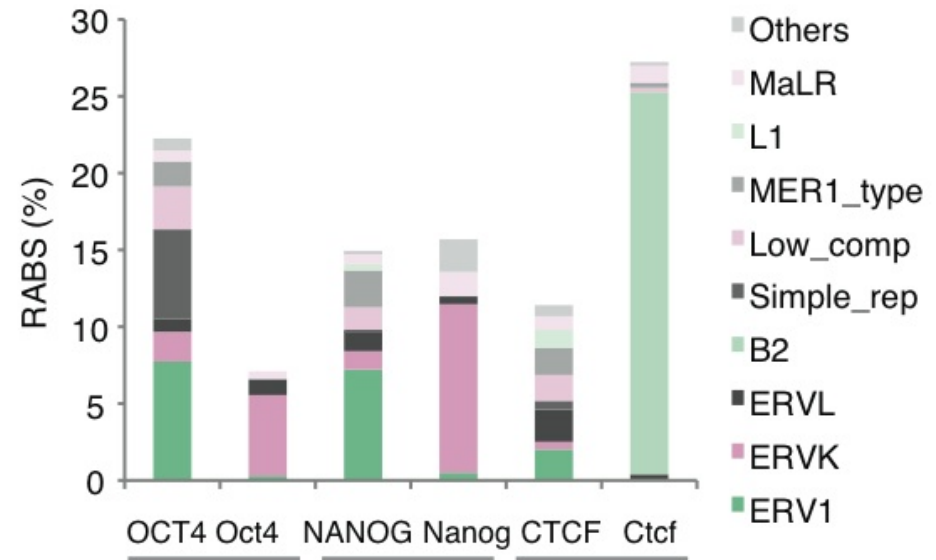
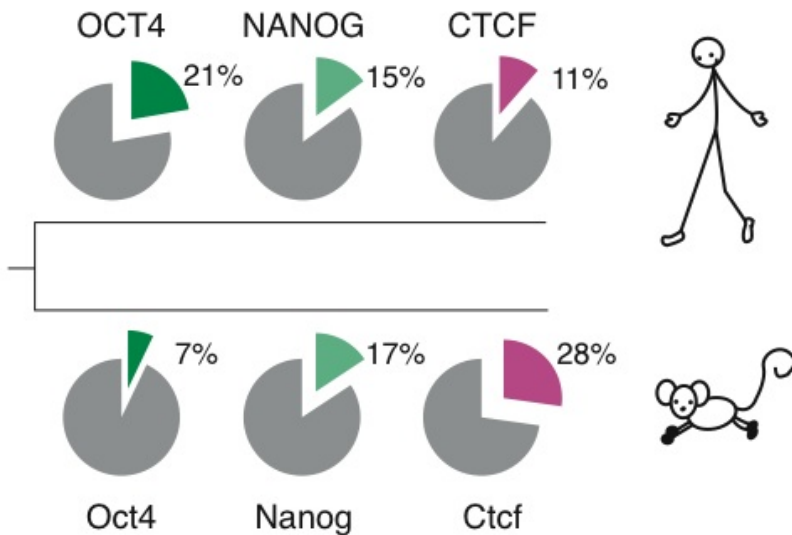


Collaboration with Huck-Hui Ng's lab at GIS

OCT4 and NANOG RABS in human ES cells

	Repeat Class	Repeat Family	Repeats	Bound Repeats	Expected	Ratio
NANOG	Low_complex	GC_rich	13426	1215	624.6	1.9
	ERV1	LTR7	2357	813	32.7	24.9
	MER1_type	MER5A	35661	806	333.1	2.4
	MER1_type	MER5A1	14897	529	136.9	3.9
	L1	L1PA2	4789	361	241.6	1.5
	ERVK	LTR5_Hs	610	340	27.8	12.2
	MER1_type	MER1B	5131	323	80.0	4.0
	MER1_type	MER5B	24784	321	213.8	1.5
	ERV1	MER4A1	2190	319	34.0	9.4
	ERV1	MER21C	4915	235	84.8	2.8
OCT4	Simple_repeat	(TG)n	54553	600	141.6	4.2
	Simple_repeat	(CA)n	54275	548	141.0	3.9
	Low_complex	GC_rich	13426	315	159.4	2.0
	ERV1	LTR9B	757	255	3.1	81.7
	MER1_type	MER5A	35661	249	105.6	2.4
	ERV1	LTR7	2357	208	10.6	19.7
	MER1_type	MER5A1	14897	204	47.4	4.3
	Low_complex	GA-rich	23936	206	80.9	2.5
	ERV1	LTR9	1962	176	14.9	11.8
	Low_complex	CT-rich	23713	181	77.2	2.3

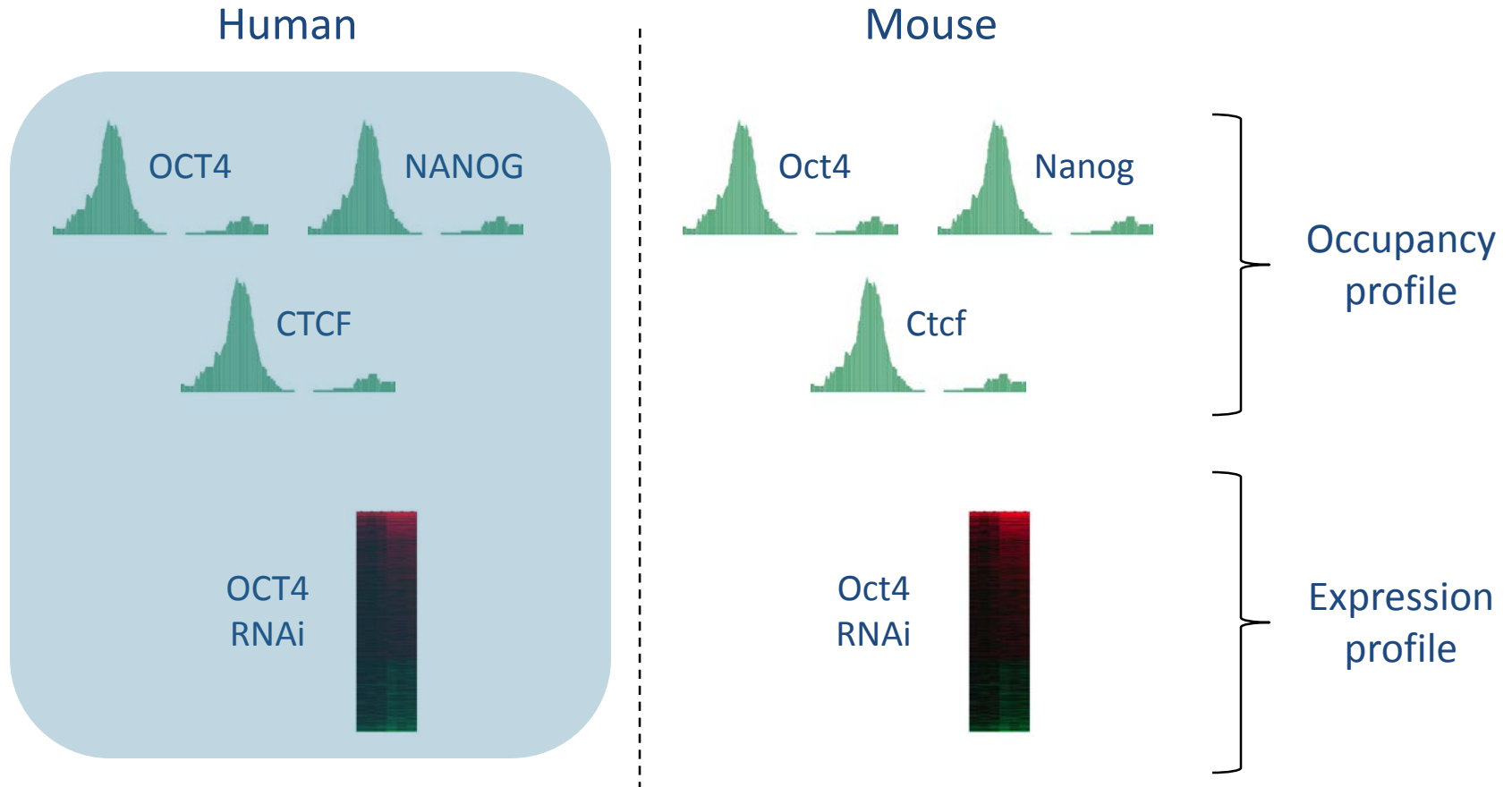
RABS in human and mouse ES cells



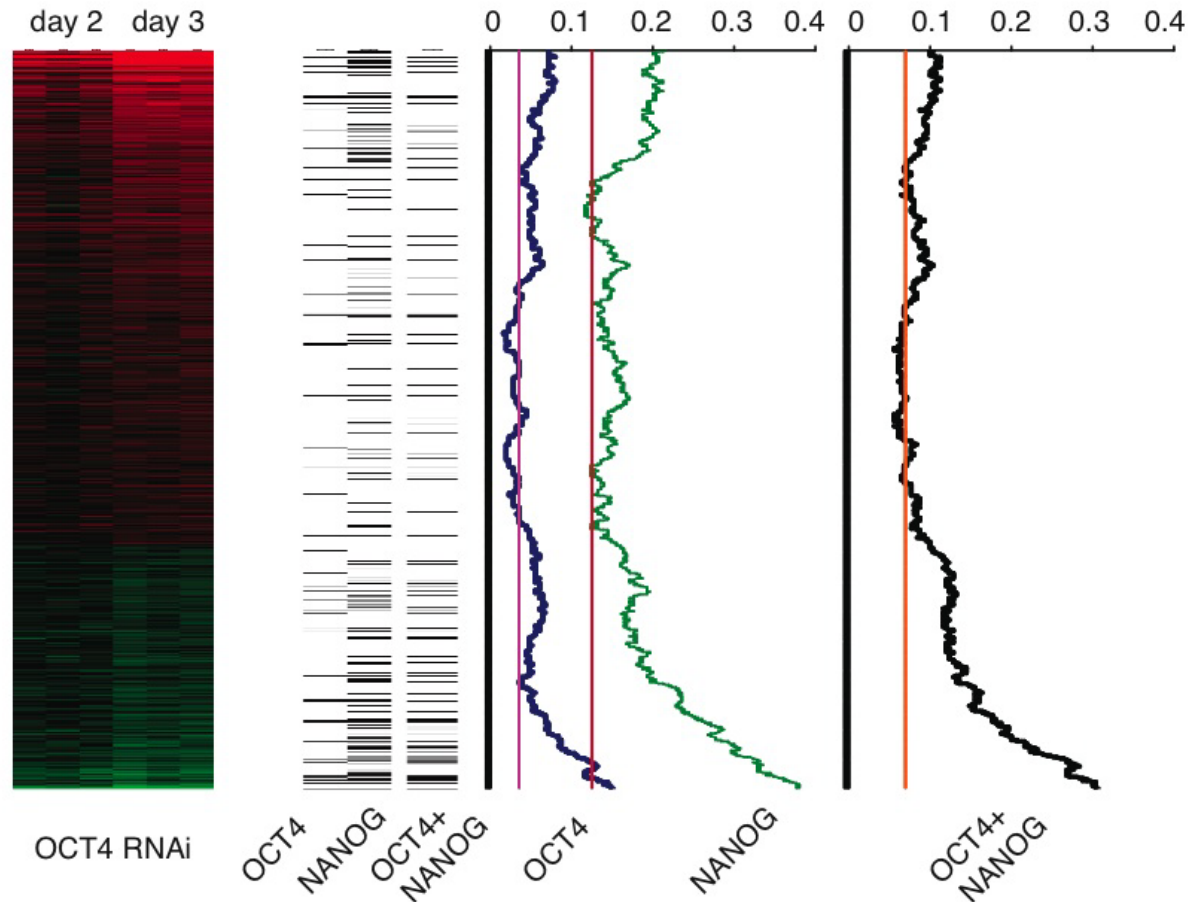
Different families of repeats have contributed a significant fraction of binding sites in both species.

Kunarso et al. *Nat Genet*, 2010

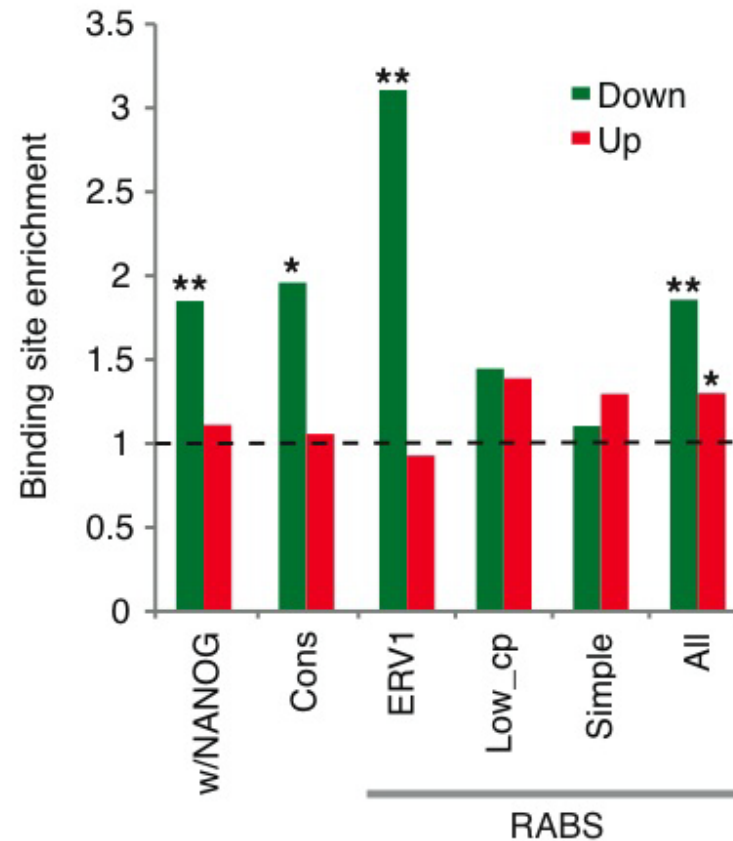
Human versus mouse ES cells



OCT4 sites are associated with genes that are down regulated following RNAi treatment

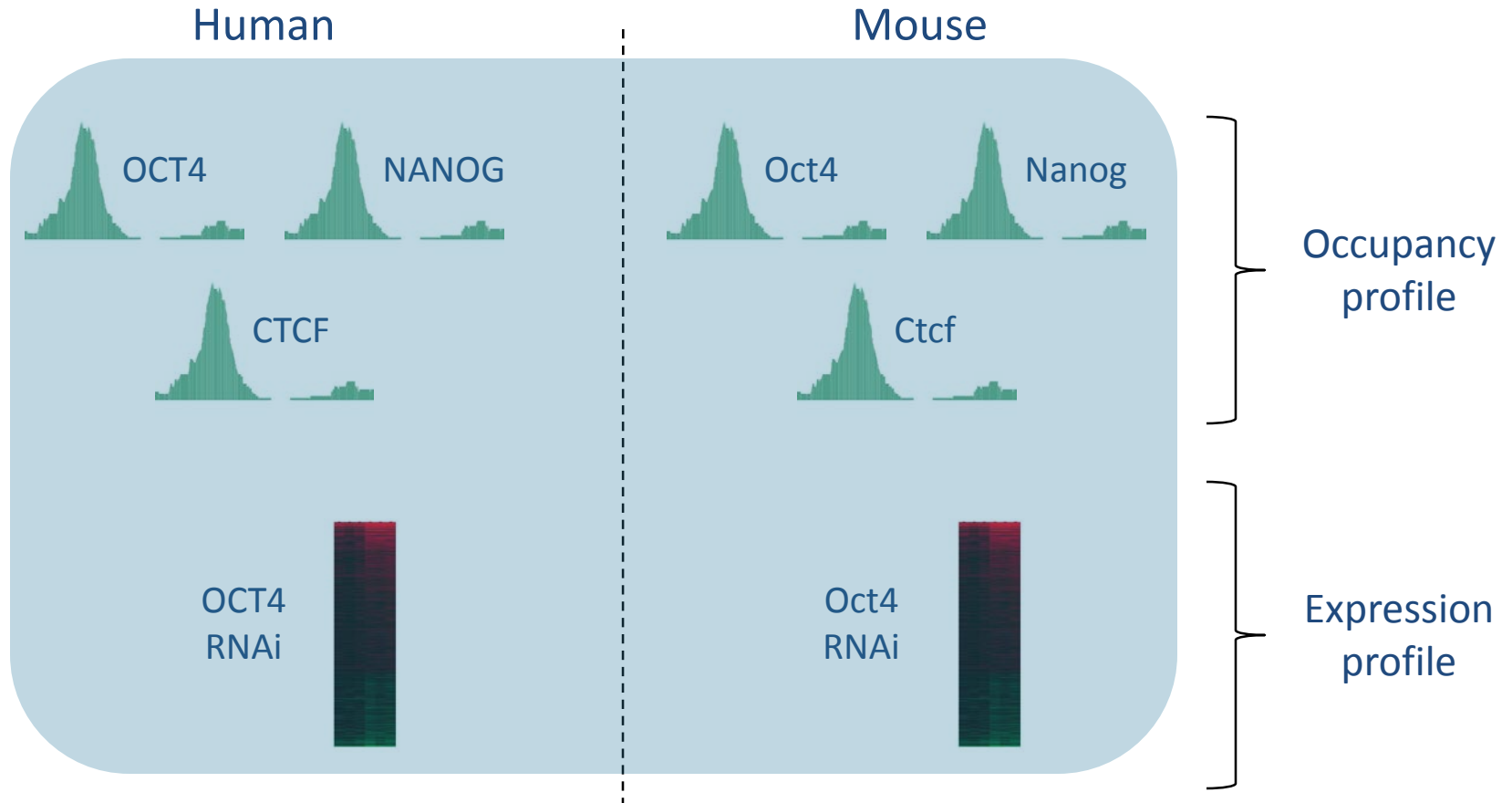


OCT4 RABS are enriched in proximity of down regulated genes



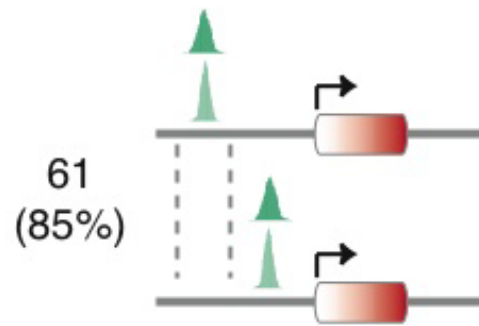
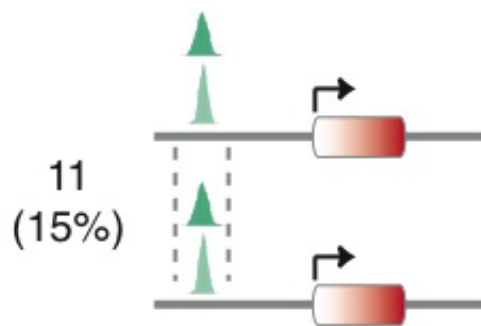
OCT4 binding site categories

Human versus mouse ES cells

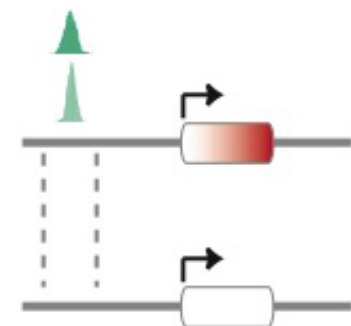


Binding sites around genes with conserved expression profiles

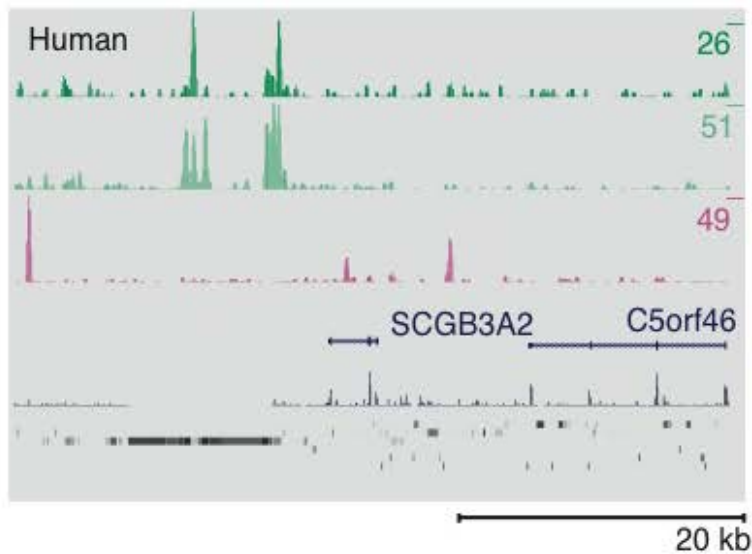
Conserved Targets (72)



Human-Specific Targets (160)

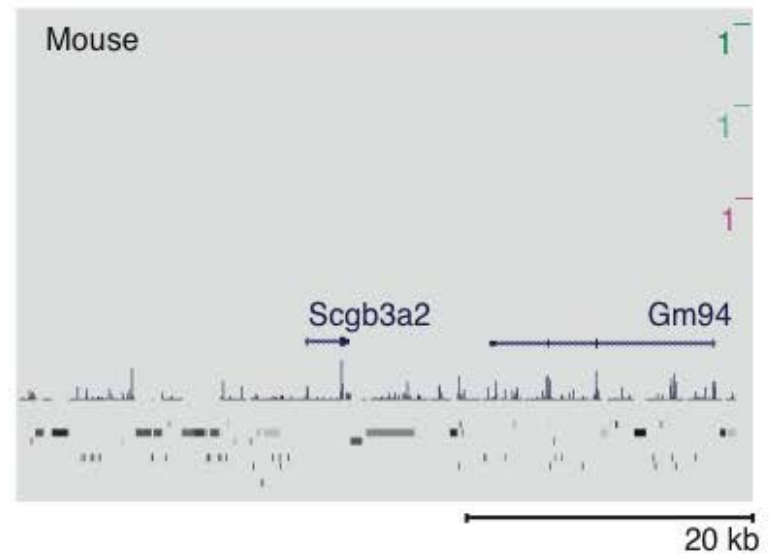


Human-specific target driven by RABS



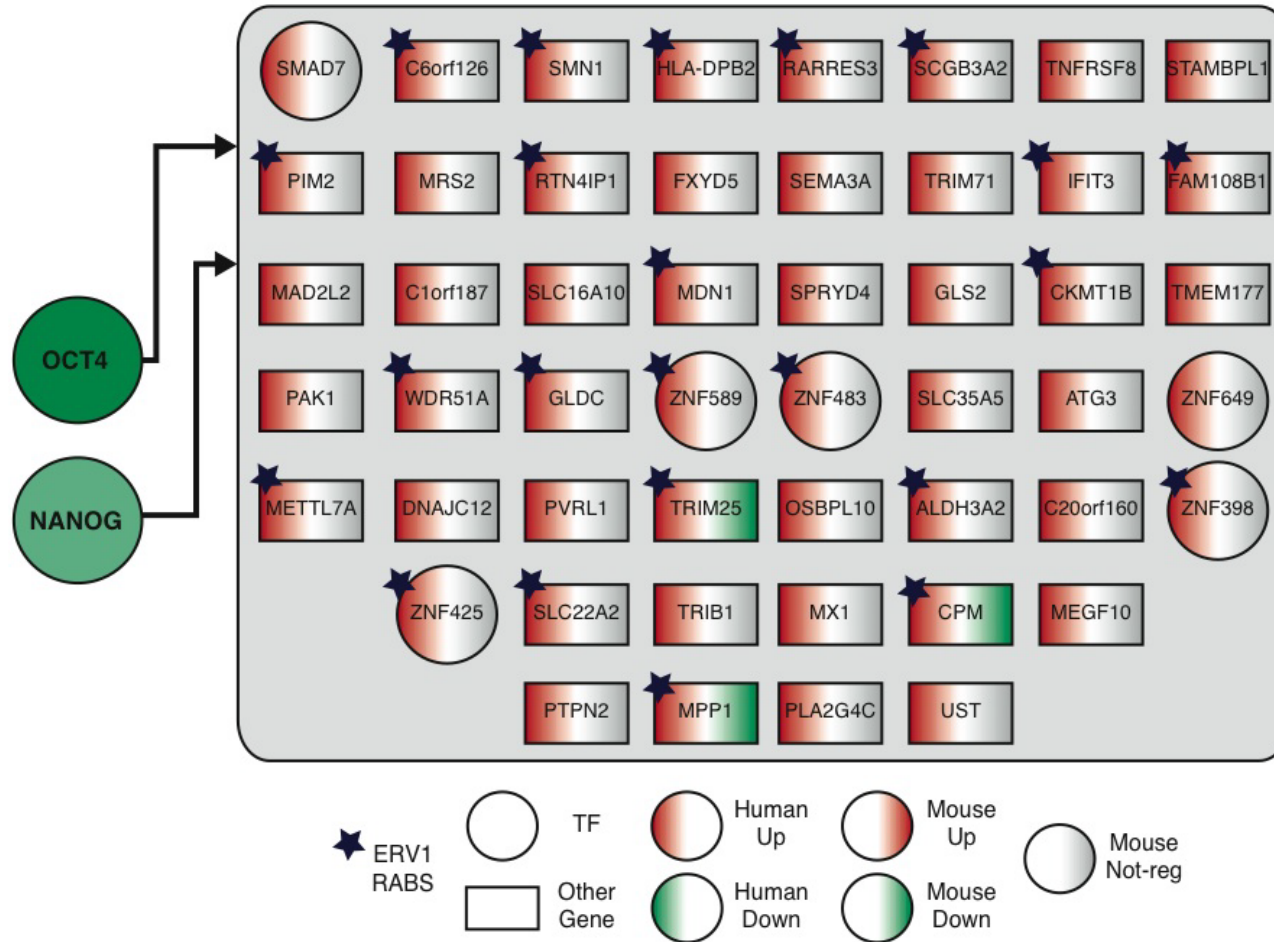
Human locus

OCT4/
Oct4
NANOG/
Nanog
CTCF/
Ctcf
Cons
Repeats

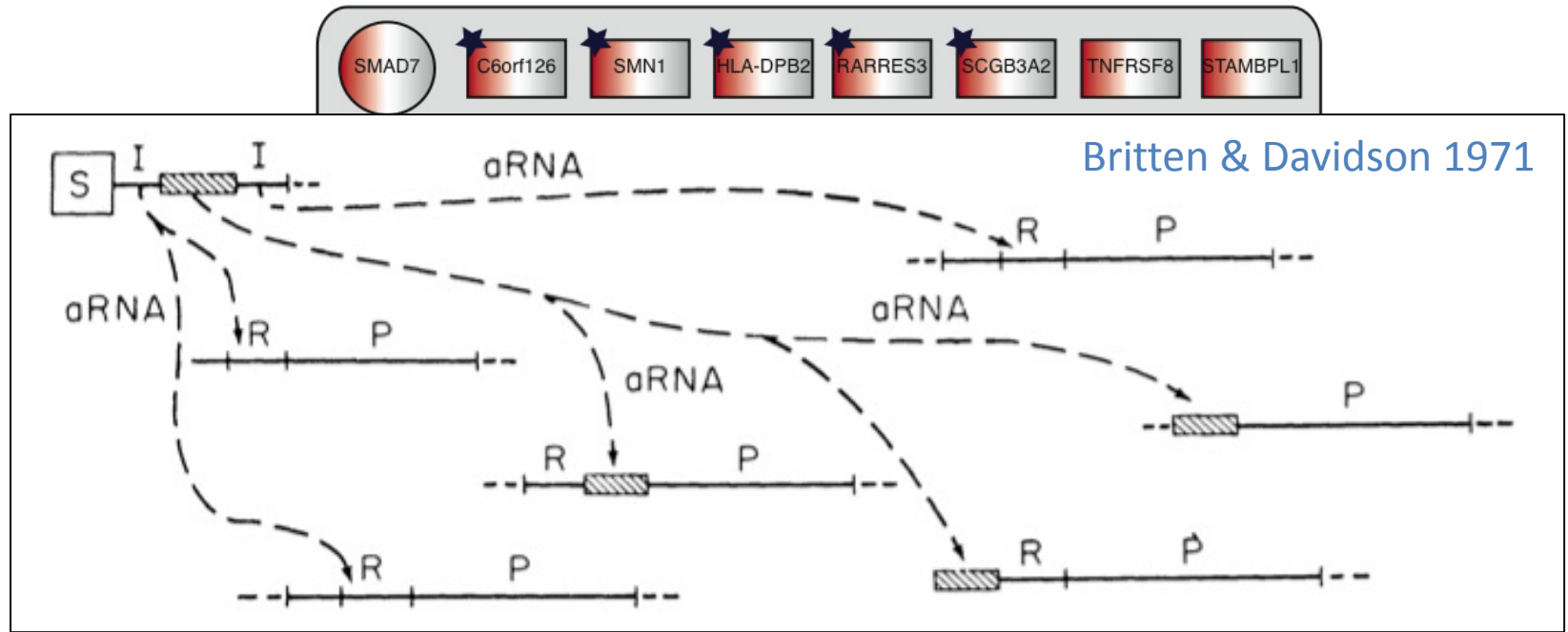


Mouse locus

Transposable Elements have Rewired the Core Regulatory Network of Human ES Cells



Transposable Elements have Rewired the Core Regulatory Network of Human ES Cells



SMAD7
C6orf126
SMN1
HLA-DPB2
RARRES3
SCGB3A2
TNFRSF8
STAMBPL1

PTPN2
MPP1
PLA2G4C
UST

★ ERV1 RABS
 ○ TF
 ○ Other Gene
 ○ Human Up
 ○ Human Down
 ○ Mouse Up
 ○ Mouse Down
 ○ Mouse Not-reg

Conclusions (Are not)

The screenshot shows the top portion of the Institute for Creation Research website. On the left is the logo 'INSTITUTE for CREATION RESEARCH'. To the right is a search bar with 'icr.org' entered and a 'GO' button. Below the search bar is a blue button that says 'FREE SUBSCRIPTION · SIGN UP'. A vertical navigation menu on the left contains the following items: 'About Us', 'Evidence for Creation', 'Resources', 'Search', 'Daily Science Updates', 'Days of Prize', 'Store', 'Give', and 'ICR Home'. The main content area is titled 'Articles' and features a large image of a DNA double helix with colorful base pairs.

Follow Daily Science Updates

Recent Articles

NASA's Ocean Currents Study Confirms Providential Care

The Plague: Birth of a Killer

Design in DNA: Flexibility Is Just Right

Transposon Behavior Negates 'Selfish Gene' Theory

Is Fossil Really a 'Game Changer'

Gene 'Jumps' Serve a Purpose, Study Shows

Send This

by Brian Thomas, M.S. *

In the tiny world of the cell, segments of DNA called transposons copy and reinsert themselves into the DNA. They eventually produce large repetitive sequences that have for many years been considered useless "junk" or remnants of ancient viral infections. But a new study has uncovered an important function for transposons.

Researchers at the Genome Institute of Singapore and other institutions suspected that transposons played an important role in the embryonic development of mammals. They decided to explore how the transposons interact with other key genetic pieces called "transcription factors" during the development of mice and humans.

Conclusions (Are not)

The screenshot shows the website for the Institute for Creation Research. At the top left is the logo "INSTITUTE for CREATION RESEARCH". To the right is a search bar with "icr.org" entered and a "GO" button. Below the search bar is a blue button that says "FREE SUBSCRIPTION · SIGN UP". On the left side, there is a navigation menu with items like "About Us", "Evidence for Creation", "Resources", "Search", "Daily Science Update", "Days of Praise", "Store", "Give", and "ICR Home". Below the menu is a "Recent Articles" section with several article titles. A large blue rounded rectangle is overlaid on the center of the page, containing a quote in white text. To the right of the quote, there is a "Send This" button and a partial article title "reinsert...ces that have...ral infections." Below the quote, there is a paragraph of text starting with "Researchers at the Genome Institute of Singapore..."

INSTITUTE
for CREATION
RESEARCH

Search icr.org GO

FREE SUBSCRIPTION · SIGN UP

Articles

► About Us
► Evidence for Creation
► Resources
► Search
Daily Science Update
Days of Praise
Store
Give
ICR Home

Follow Daily Science Update

Recent Articles

NASA's Ocean Curiosity
Confirms Provident

The Plague: Birth of

Design in DNA: Flexibility Is Just
Right

Transposon Behavior Negates
'Selfish Gene' Theory

Is Fossil Really a 'Game Changer'

reinsert
ces that have
ral infections.

Send This

“However, a biblical perspective predicts high functionality throughout genomes, with traces of degradation having accumulated since the curse that God placed on the earth, ... the high level of functionality of transposons is more consistent with creation.”

Researchers at the Genome Institute of Singapore and other institutions suspected that transposons played an important role in the embryonic development of mammals. They decided to explore how the transposons interact with other key genetic pieces called "transcription factors" during the development of mice and humans.

Conclusions

- Up to 25% of new regulatory sites have been contributed by transposable elements in both human and mouse.
- Only 15% of the genes with conserved expression profiles have a conserved binding site.
- A number of human-specific targets have been wired into the core regulatory network of ES cells by repeats.

Repeats and gene regulation

Nat Genet 2011

Transposon-mediated rewiring of gene regulatory networks contributed to the evolution of pregnancy in mammals

Vincent J Lynch, Robert D Leclerc, Gemma May & Günter P Wagner

Cell 2012

Waves of Retrotransposon Expansion Remodel Genome Organization and CTCF Binding in Multiple Mammalian Lineages

Dominic Schmidt,^{1,2,6} Petra C. Schwalie,^{3,6} Michael D. Wilson,^{1,2} Benoit Ballester,³ Ângela Gonçalves,³ Claudia Kutter,^{1,2} Gordon D. Brown,^{1,2} Aileen Marshall,^{1,5} Paul Flicek,^{3,4,*} and Duncan T. Odom^{1,2,4,*}

OPEN ACCESS Freely available online

PLoS Genet 2013

The Majority of Primate-Specific Regulatory Sequences Are Derived from Transposable Elements

Pierre-Étienne Jacques^{1,2}, Justin Jeyakani¹, Guillaume Bourque^{3,4*}

¹ Computational and Systems Biology, Genome Institute of Singapore, Singapore, Singapore, ² Département de Biologie, Université de Sherbrooke, Sherbrooke, Québec, Canada, ³ Department of Human Genetics, McGill University, Montréal, Québec, Canada, ⁴ McGill University and Génome Québec Innovation Center, Montréal, Québec, Canada

Acknowledgements

Lab

Galih Kunarso

Pierre-Etienne Jacques

LeeAnn Ramsay

Justin Jeyakani

Jean Monlong

Simon Girard

Toby Hocking

Eric Audemard

EDCC team

David Bujold

Bryan Caron

David Morais (UdeS)

Carol Gauthier (UdeS)

Alain Veilleux (UdeS)

ME Rousseau

Analysis team

Louis Letourneau

Mathieu Bourgey

Maxime Caron

Gary Lévesque

Robert Eveleigh

Francois Lefebvre

Johanna Sandoval

Pascale Marquis

Joel Fillon

Julien Tremblay

IT team

Terrance Mcquilkin

Marc-André Labonté

Genevieve Dancausse

Andras Frankel

Alexandru Guja

Development team

Nathalie Émond

Francois Cantin

Catherine Côté

Daniel Guertin

Louis Dumond Joseph

Francois Korbuly

Marc Michaud

Thuong Ngo

Francois Massé

Open positions for
students and postdoc!

guil.bourque@mcgill.ca