

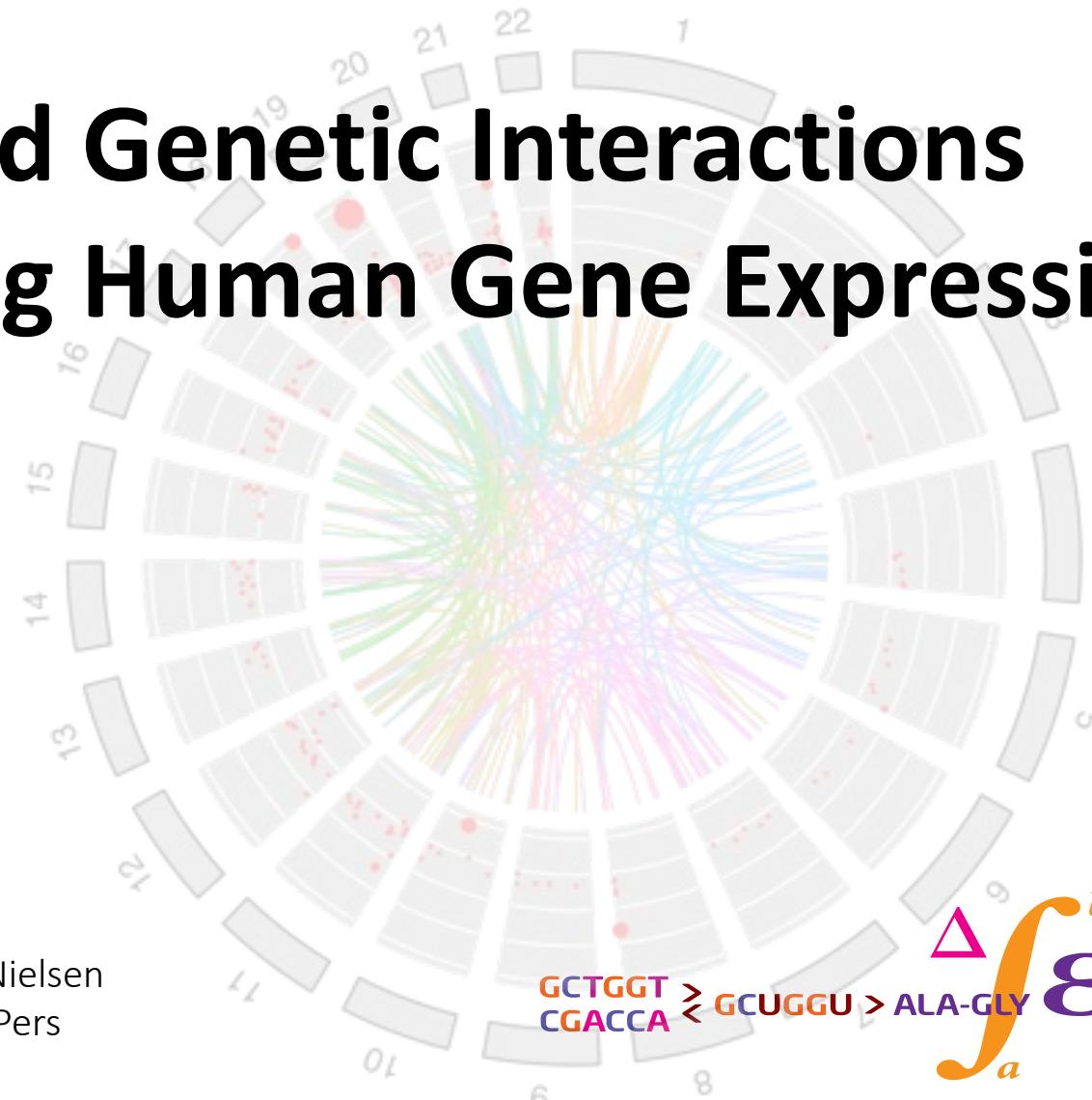
Spatial and Genetic Interactions Influencing Human Gene Expression

Pascal Timshel

26th November 2015

Master Thesis Defense
35 ECTS points

Supervisors: Henrik Bjørn Nielsen
Tune Hannes Pers



GCTGGT
CGACCA

GCUGGU > ALA-GLY

$\sum!$

$\int_a^b \Theta^{\sqrt{17}} + \Omega \int \delta e^{in} =$

$\infty = \{2.7182818284$

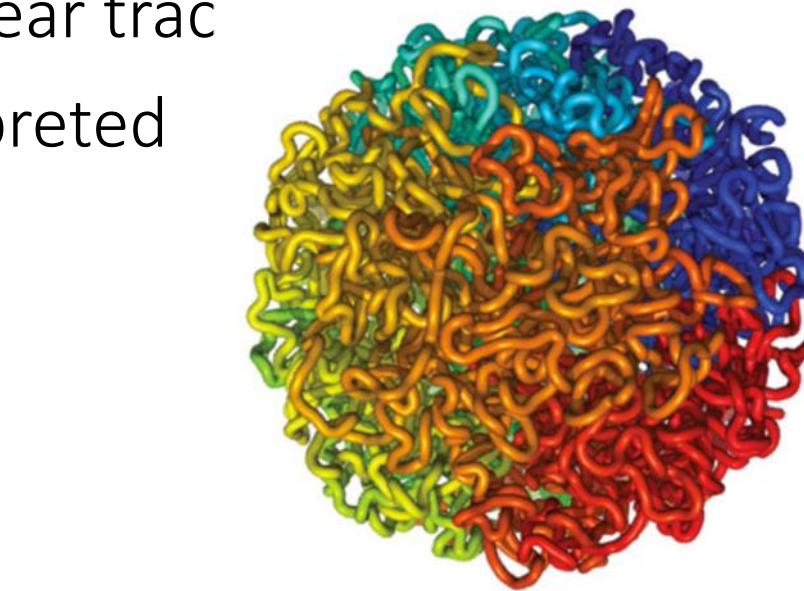
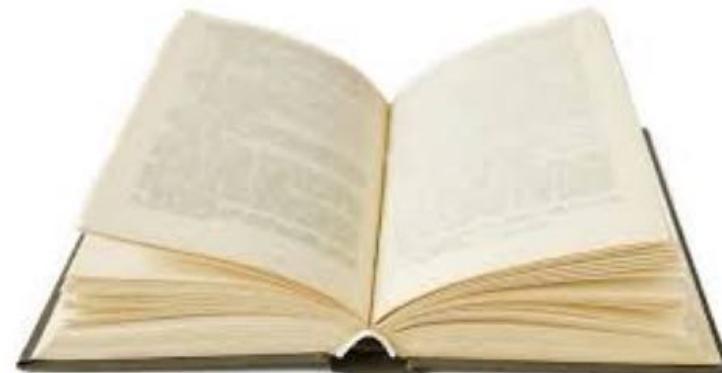
$\Sigma \gg 0, 0,$

Outline

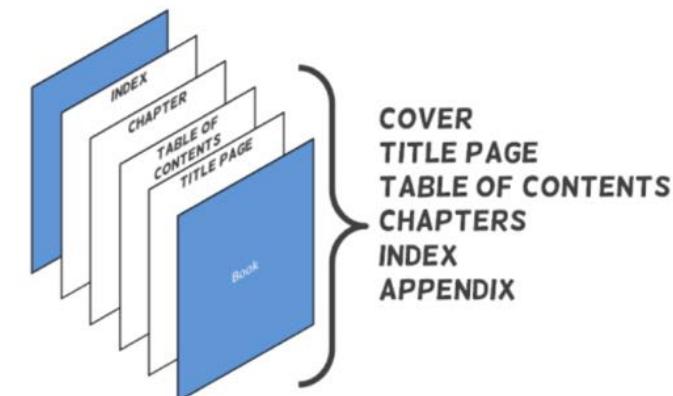
- Motivation
- Spatial Interactions (Genome Organization)
- Genetic Interactions (Epistasis)
- Spatial Epistasis Hypothesis
- Methods and Computational Framework
- Epistatic Signals
- “Meta-analysis”

The Genome has a 3rd Dimension

- The genome is traditionally “read” 1D linear tracks
- The genomic information might be interpreted differently in a 3D context

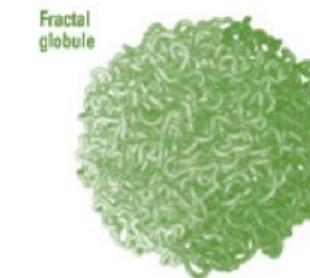
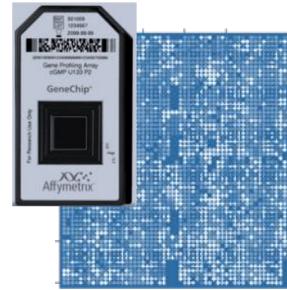


COMMON PARTS OF A BOOK



Motivation and overview

Connecting gene expression and genetic interactions
and chromosomal organization



How do spatial and genetic interactions influence gene expression?

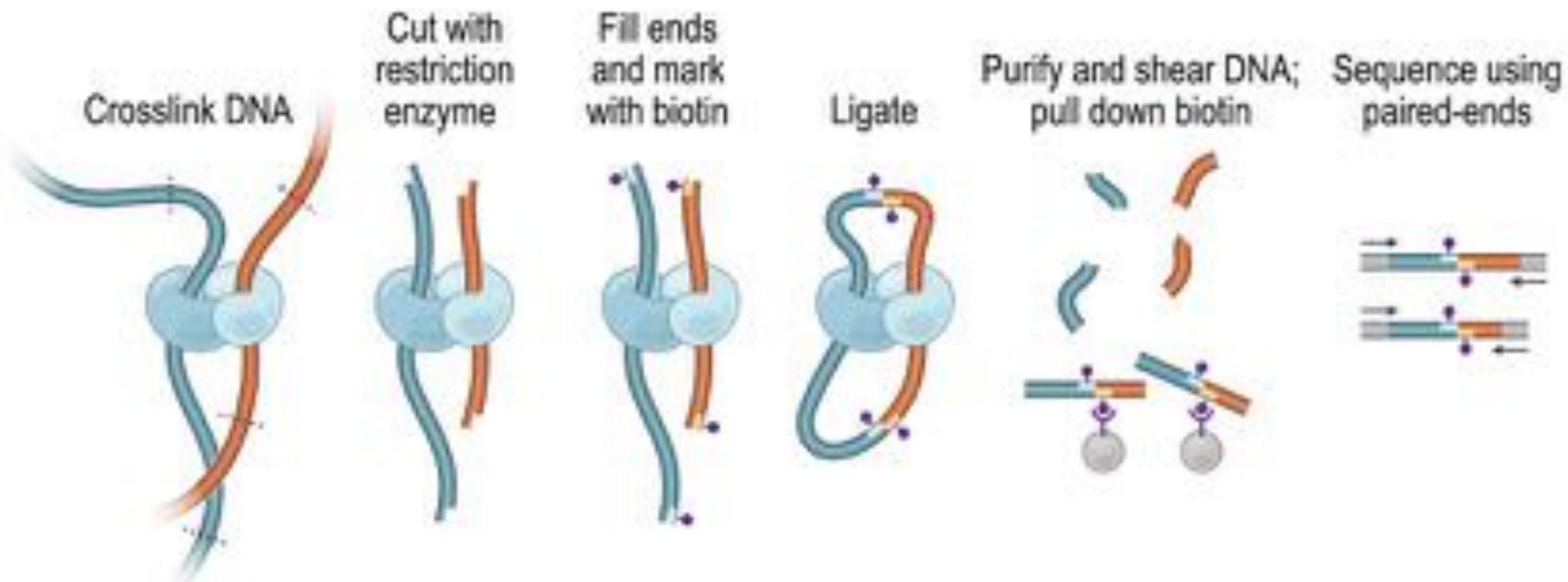
Research questions

1. What are the regulatory molecular mechanisms that cause epistasis?
2. Can regulatory molecular mechanisms facilitate identification of epistatic effects?
3. To which extent does epistasis influences human gene expression?

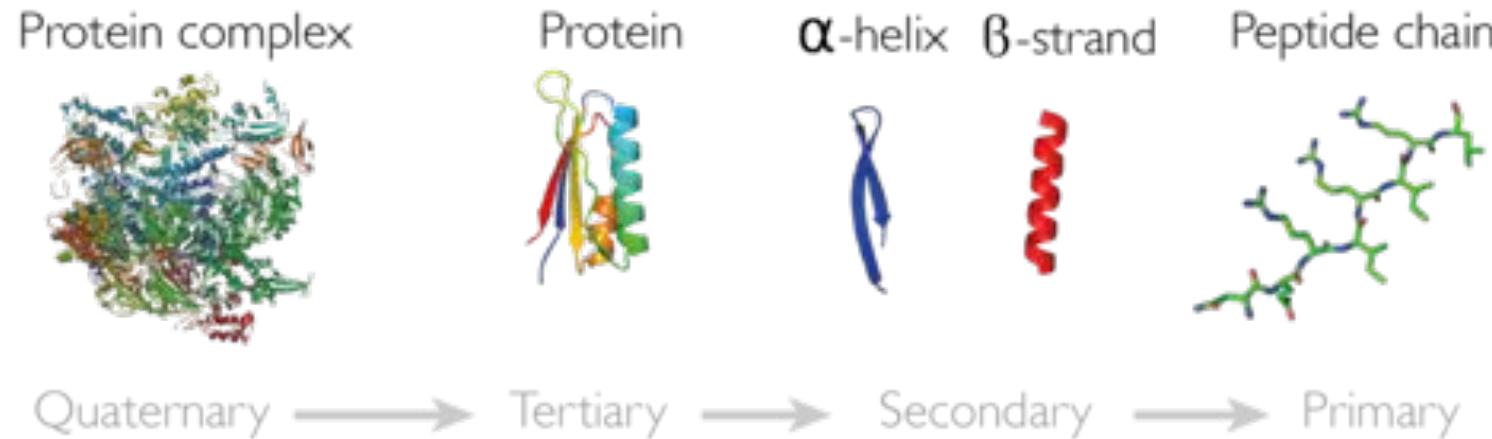
Outline

- Motivation
- Spatial Interactions (Genome Organization)
- Genetic Interactions (Epistasis)
- Spatial Epistasis Hypothesis
- Methods and Computational Framework
- Epistatic Signals
- “Meta-analysis”

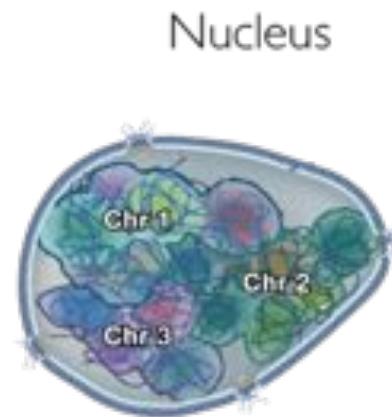
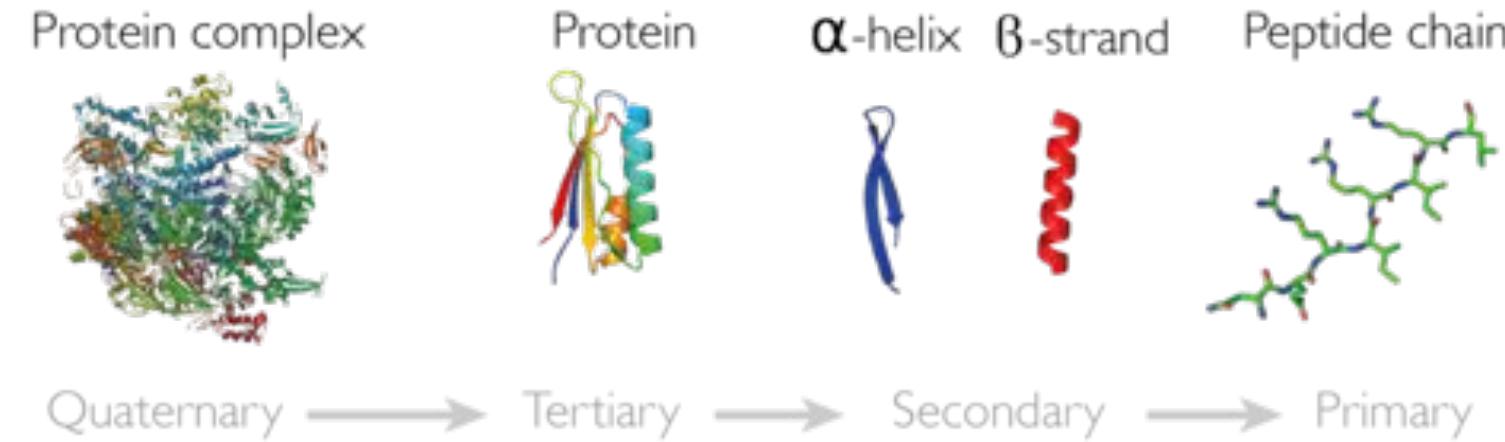
Hi-C protocol (Rao et al., 2014, *Cell*)



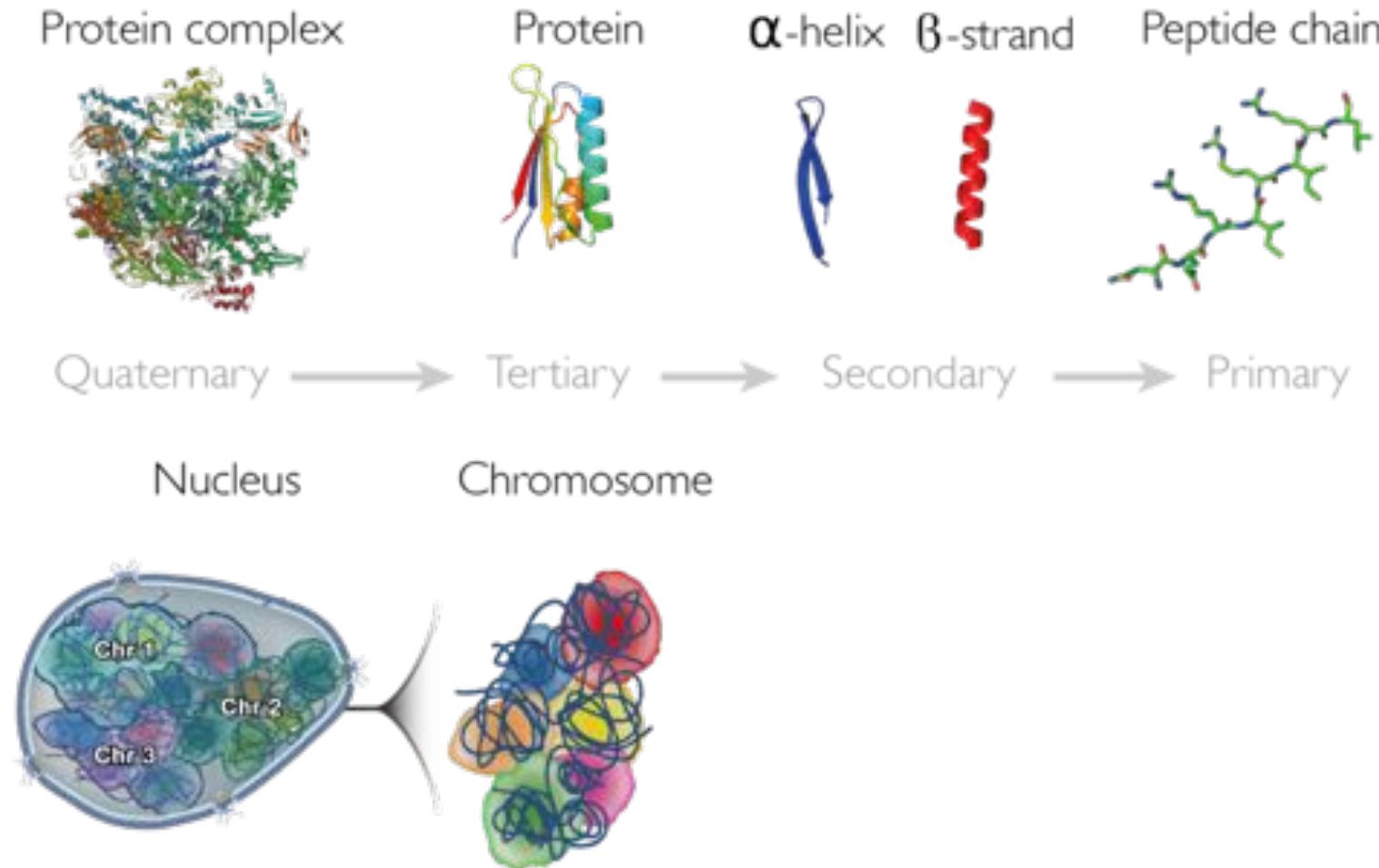
Analogous hierarchical organization of protein and genome structure



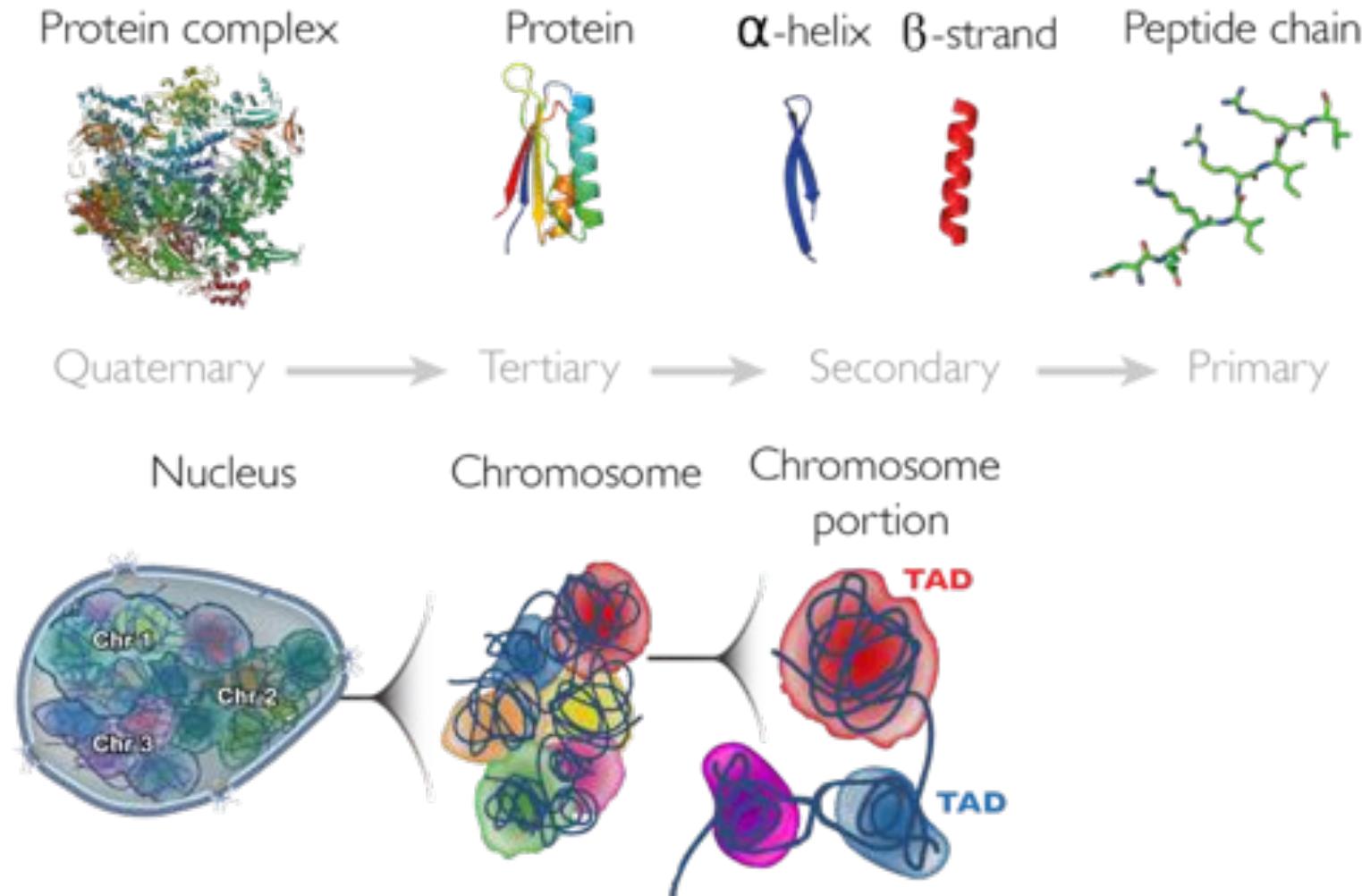
Analogous hierarchical organization of protein and genome structure



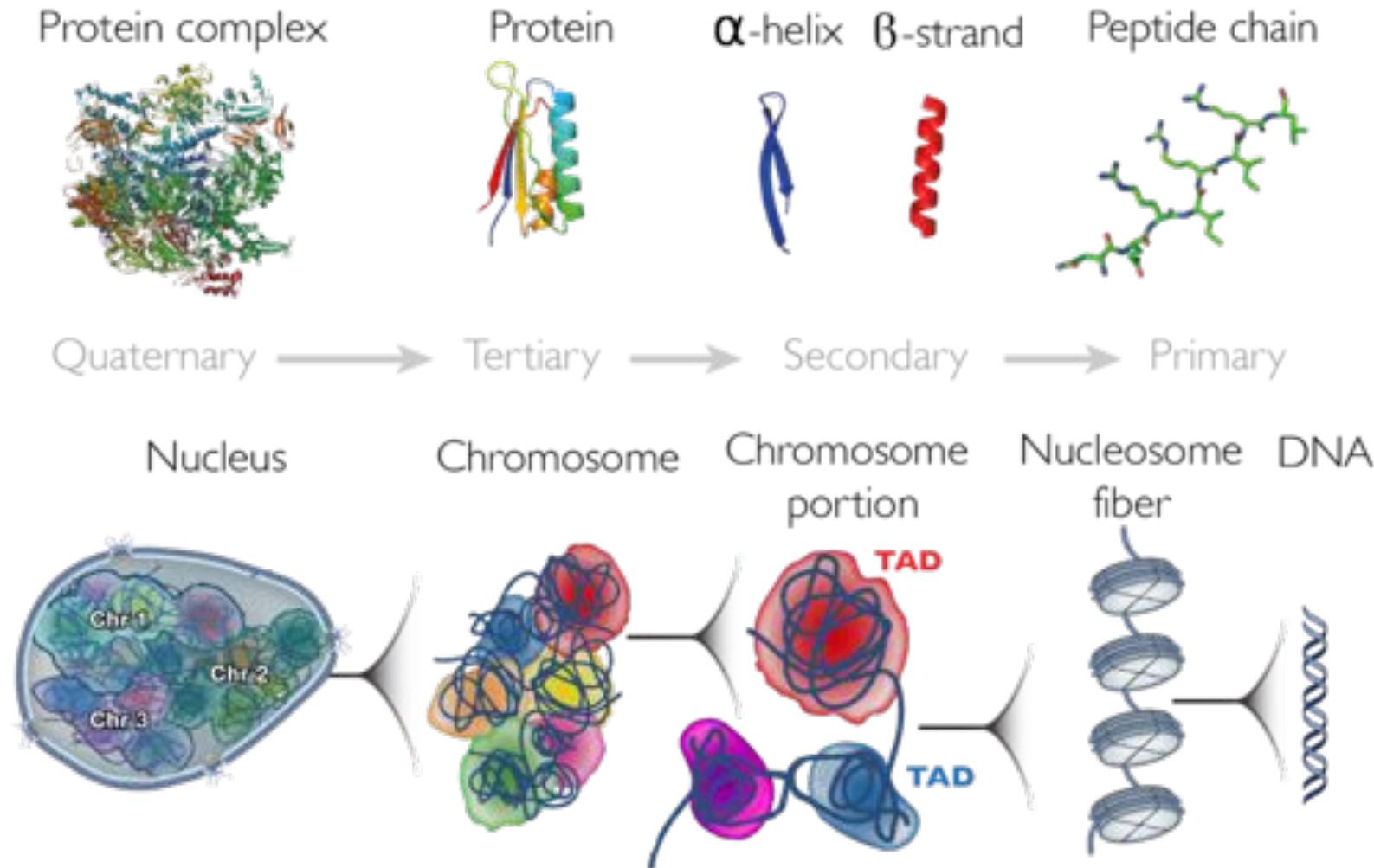
Analogous hierarchical organization of protein and genome structure



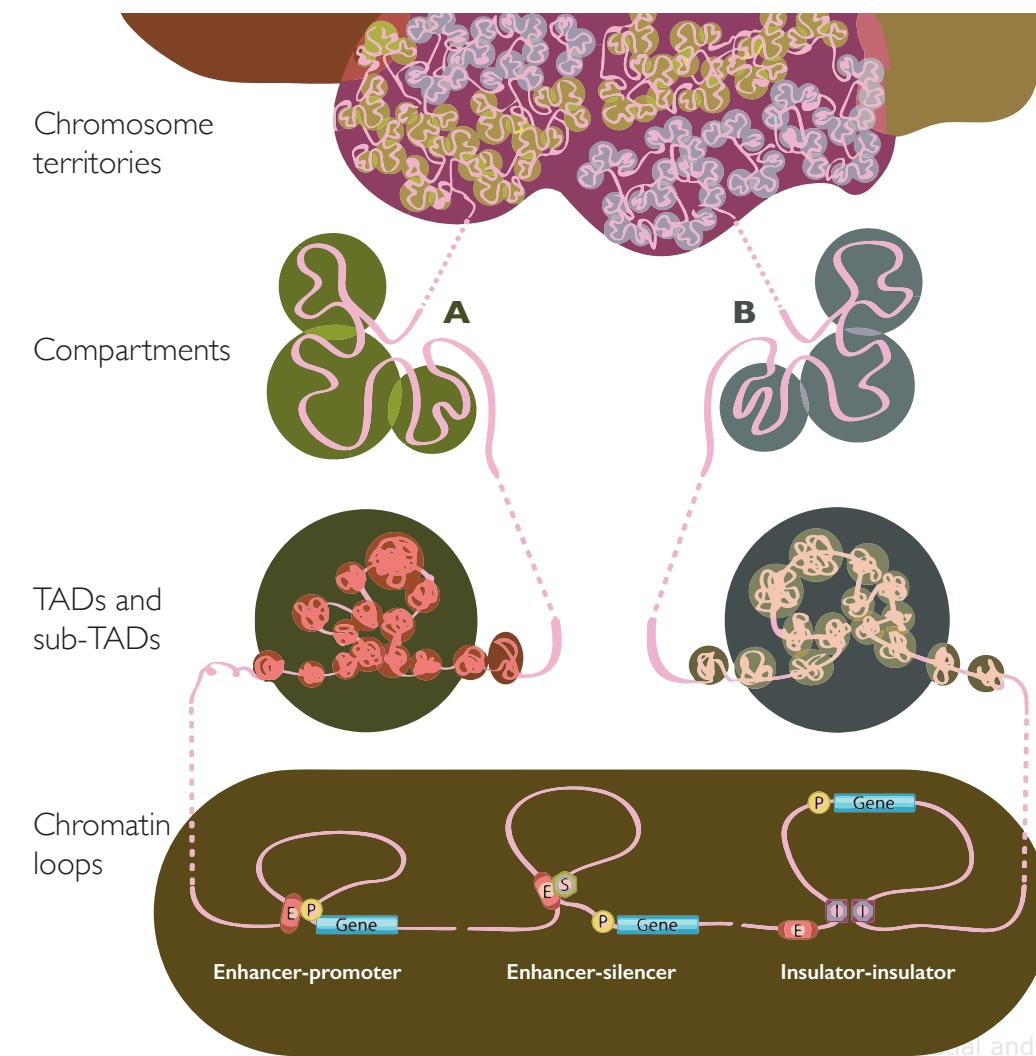
Analogous hierarchical organization of protein and genome structure



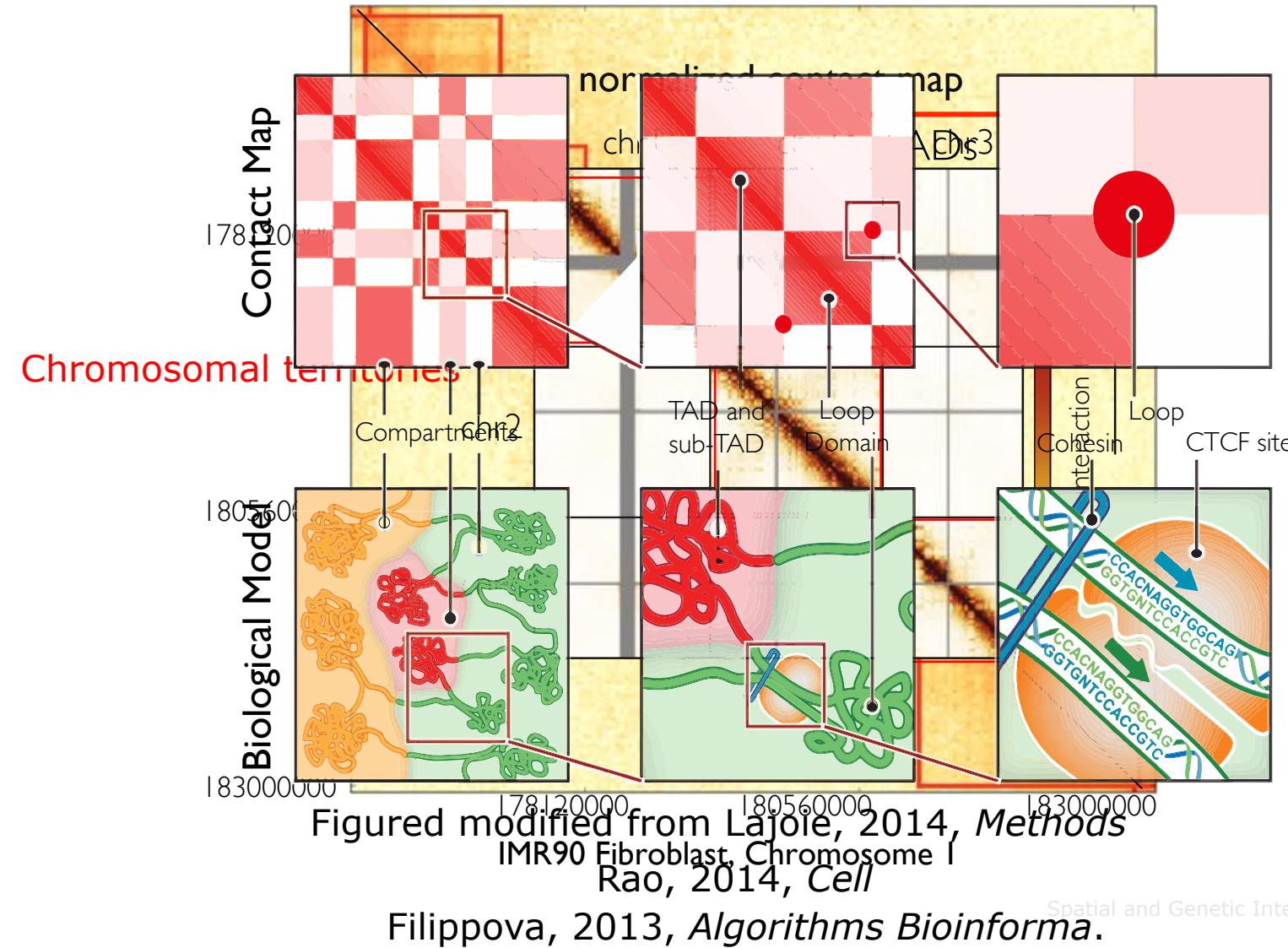
Analogous hierarchical organization of protein and genome structure



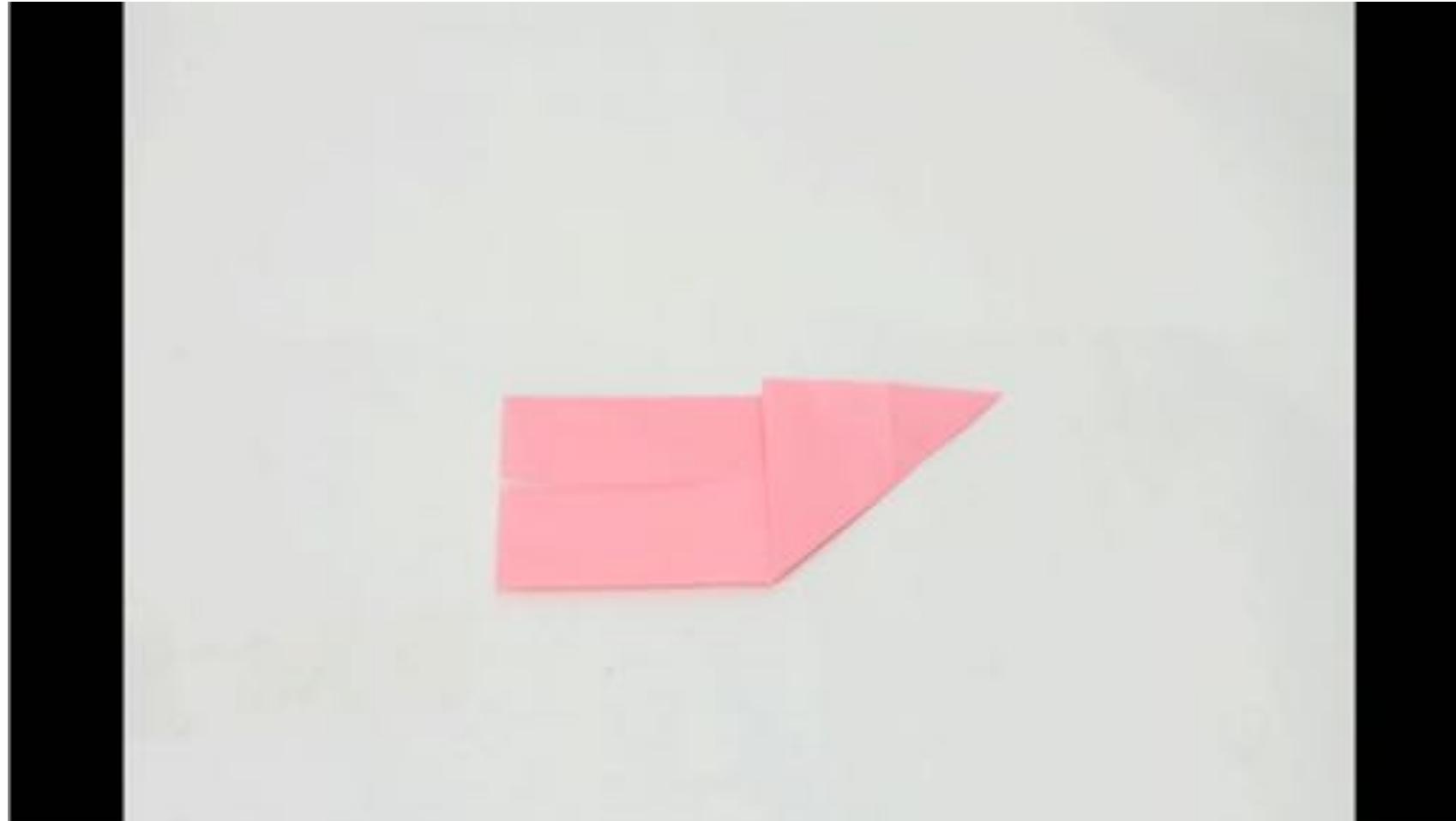
Hierarchical Organization of the Human Genome



Contact Maps Reveals Genome Structures



“Epigenetics is Genome Origami”

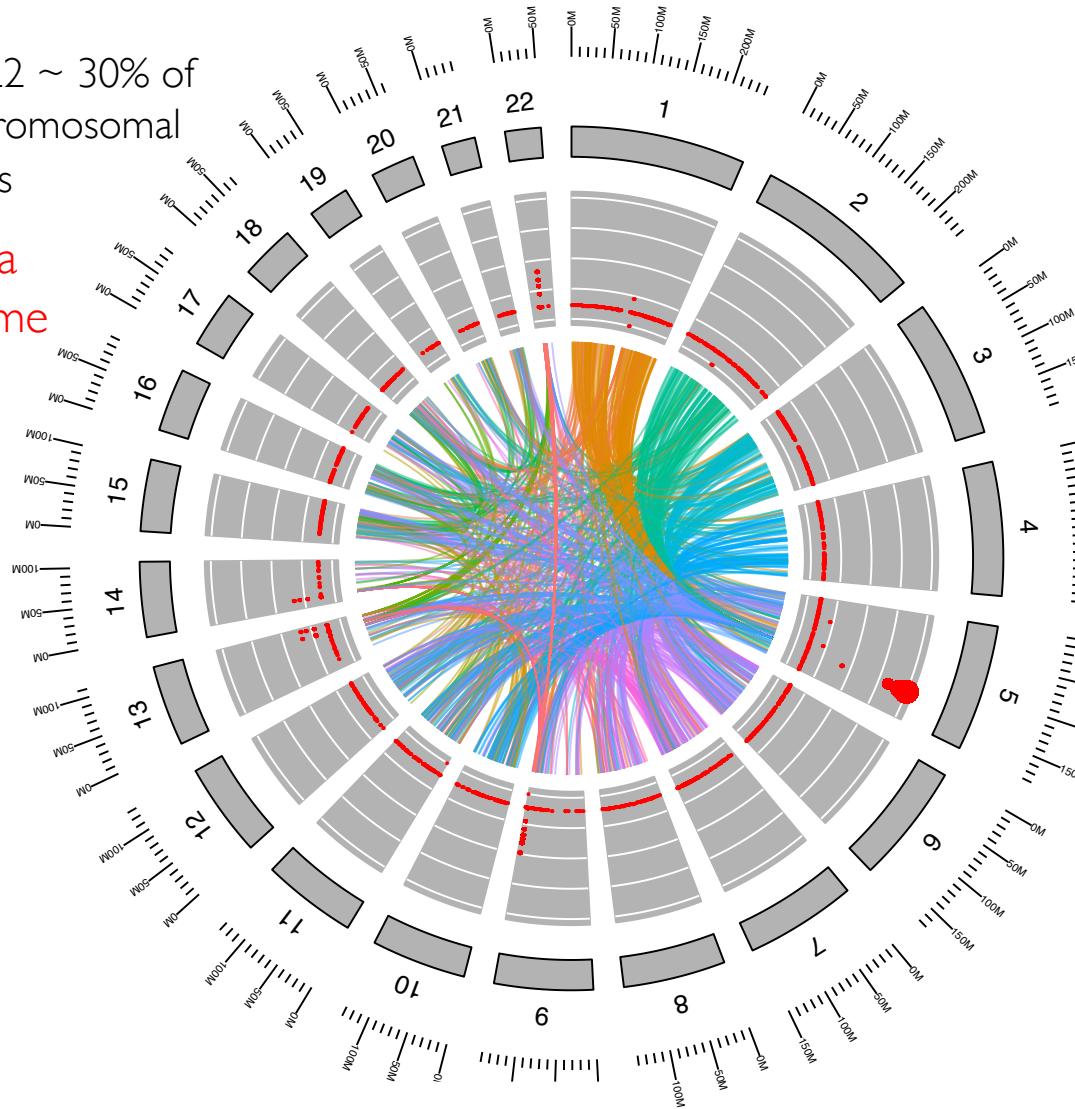


Visualization - circos plot

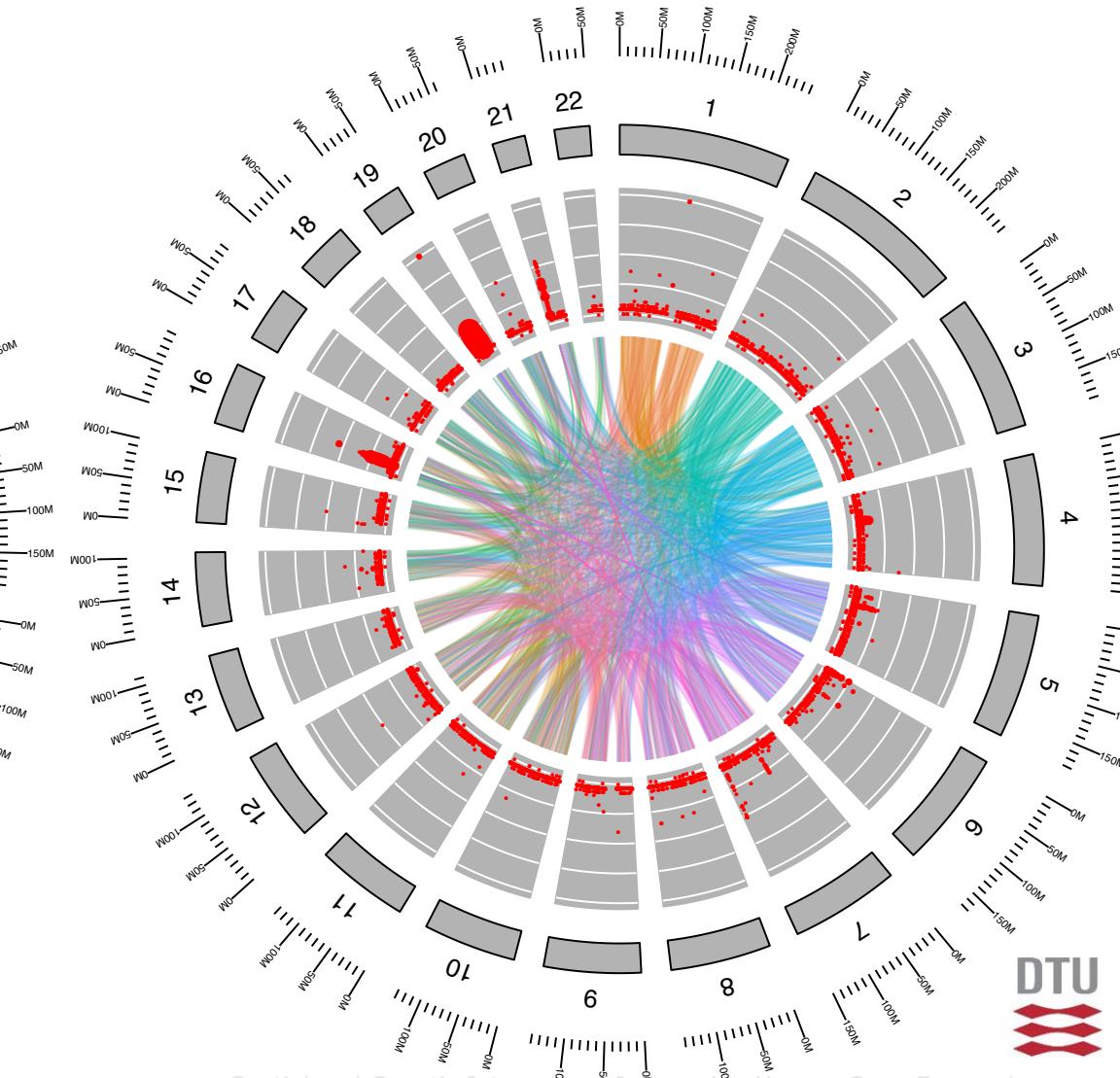
K652 cells Leukemia cancer cell line

Chr9-Chr22 ~ 30% of all inter-chromosomal interactions

Philadelphia chromosome



hESC cells

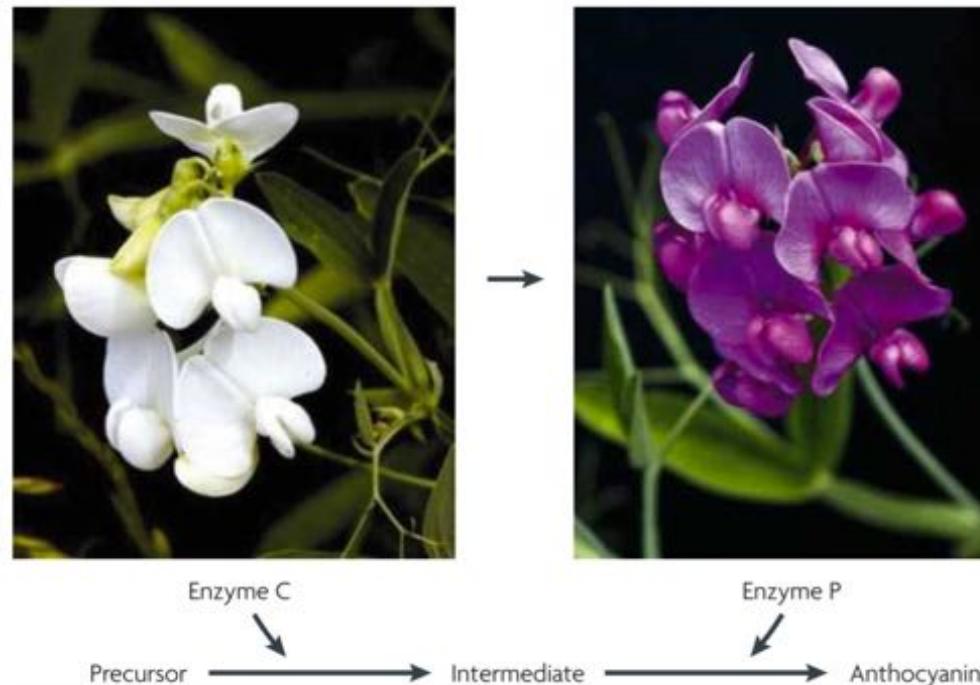


Outline

- Motivation
- Spatial Interactions (Genome Organization)
- **Genetic Interactions (Epistasis)**
- Spatial Epistasis Hypothesis
- Methods and Computational Framework
- Epistatic Signals
- “Meta-analysis”

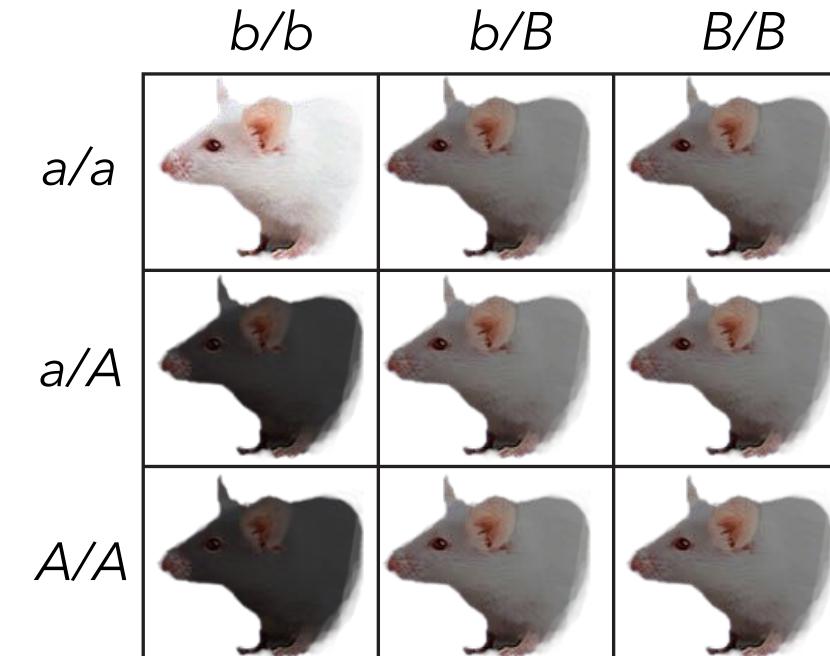
What is Epistasis?

- Biological definition (Bateson, 1909, *Cambridge University Press*)
 - “Standing upon”. Allele at one locus masks the effects of alleles at one or more other loci.



21

Bateson and Punnett classical sweet pea experiment
Adopted from Phillips, 2008, *Nat. Rev. Genet*



CONTENTS.

1. The superposition of factors distributed independently	405
2. Phase frequency in each array	405
3. Partial regression	405
4. Dominance	405
5. Correlation for parent; genetic correlations	405
6. Fraternal correlation	405
7. Correlations for other relatives	405
8. Epistasis	409
9. Asymptotic mating	410
10. Frequency of phases	410
11. Asymptotic factors	411
12. Conditions of inheritance	412
13. Nature of association	412
14. Multiple dielomorphism	415
15. Homozygosity and multiple allelic asciophores	416
16. Coupling	418
17. Theories of marital correlation; ancestral correlations	419
18. Ancestral correlations (second and third theories)	421
19. Numerical values of association	421
20. Fraternal correlations	422
21. Numerical values for environment and dominance ratios; analysis of variance	423
22. Other relatives	424
23. Numerical values (third theory)	425
24. Numerical values of results	427
25. Interpretation of dominance ratio (diagrams)	428
26. Summary	428

Several attempts have already been made to interpret the well-established results of biometry in accordance with the Mendelian scheme of inheritance. It is here attempted to ascertain the biometrical properties of a population of a more general type than has hitherto been examined, inheritance in which follows this

What is Epistasis?

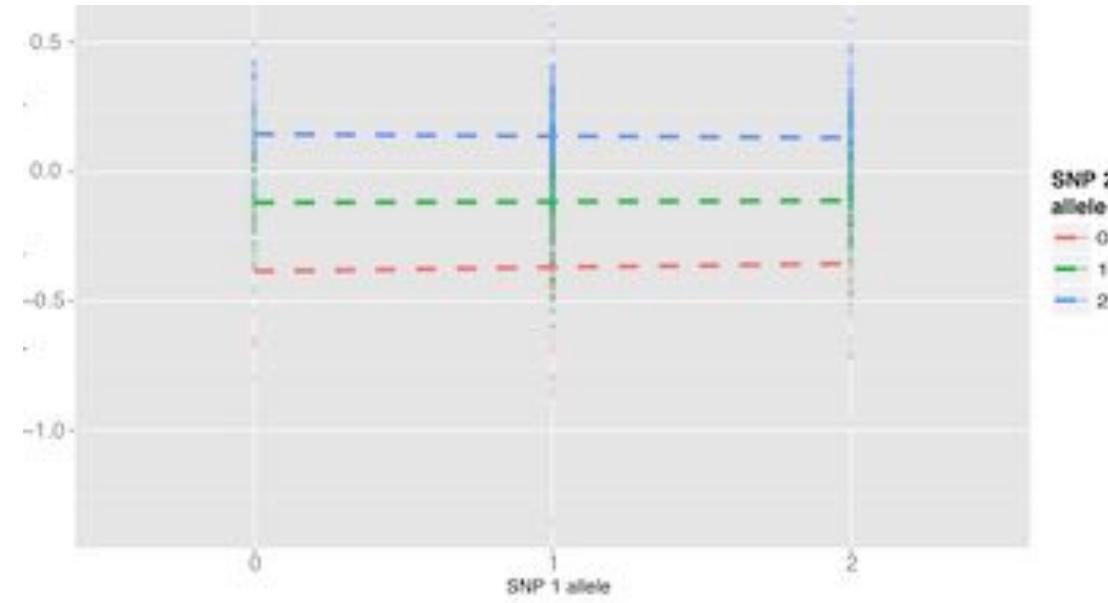
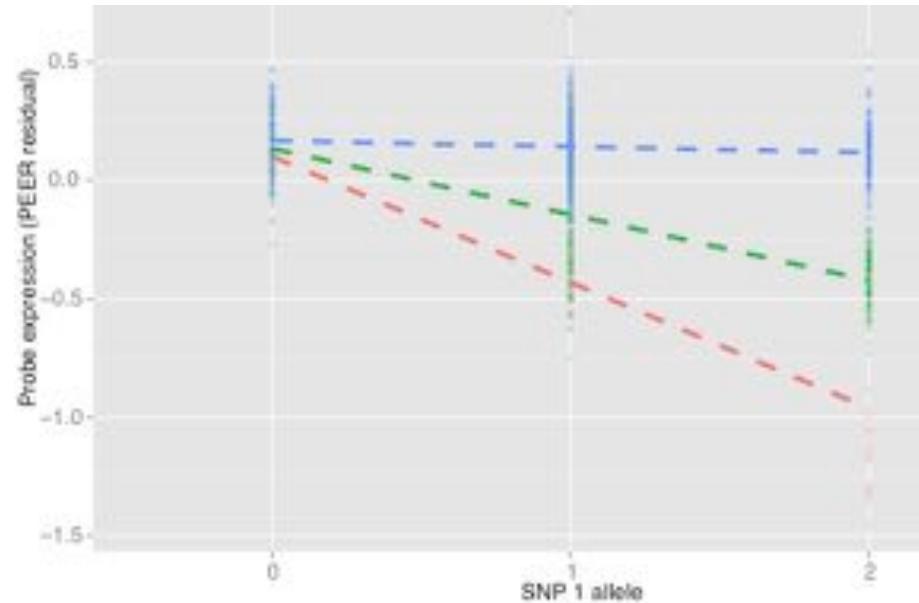
- Statistical definition (Fisher, 1919, *Trans. R. Soc.*)

—deviations from additivity in a linear statistical model

$$\text{Expression(probe)} = \alpha + \beta_1 \cdot \text{SNP}_1 + \beta_2 \cdot \text{SNP}_2 + \beta_{12} \cdot \text{SNP}_1 \cdot \text{SNP}_2$$

$$P < 7 \times 10^{-45}$$

$$P = 0.5$$



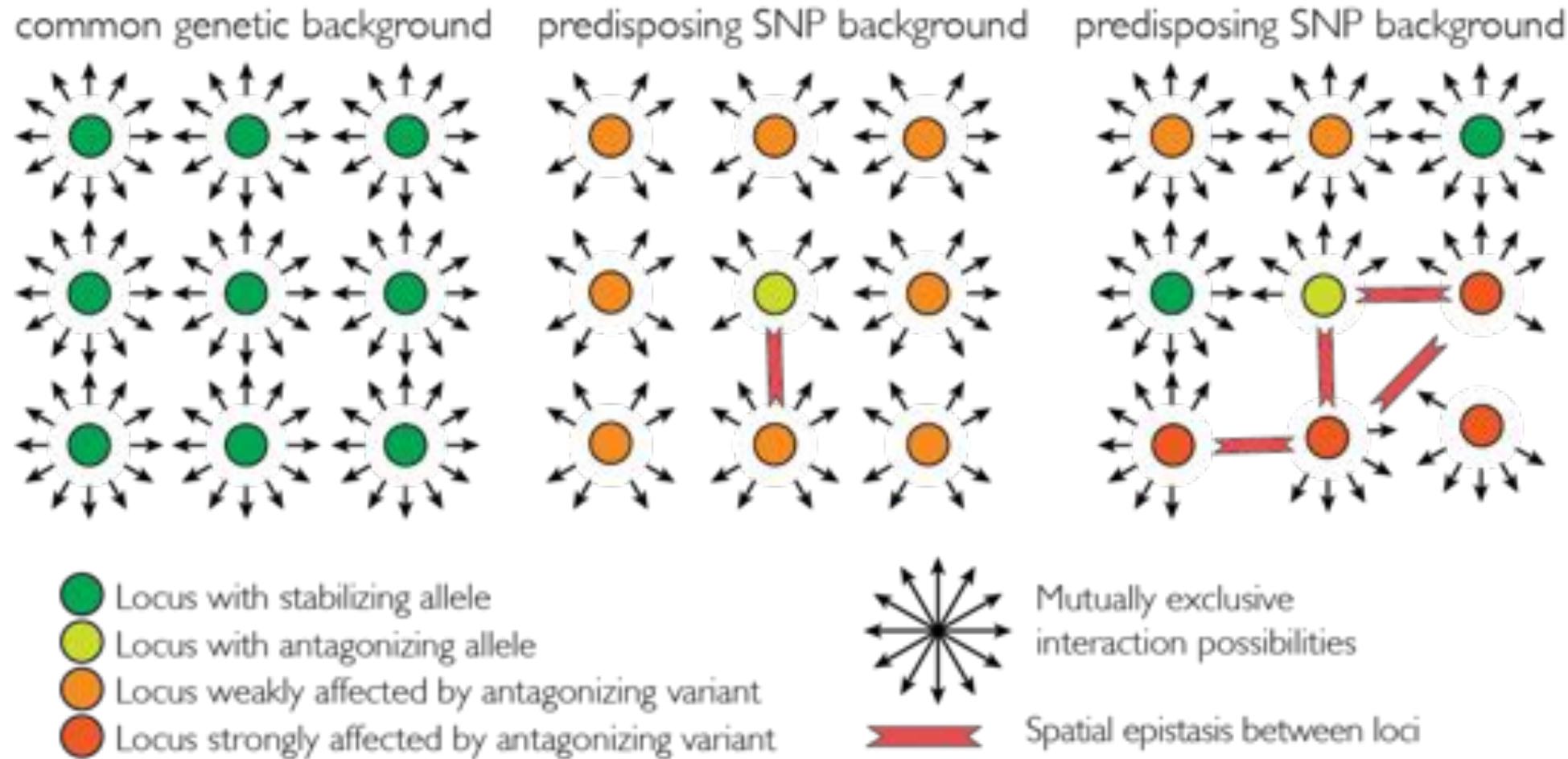
Outline

- Motivation
- Spatial Interactions (Genome Organization)
- Genetic Interactions (Epistasis)
- **Spatial Epistasis Hypothesis**
- Methods and Computational Framework
- Epistatic Signals
- “Meta-analysis”

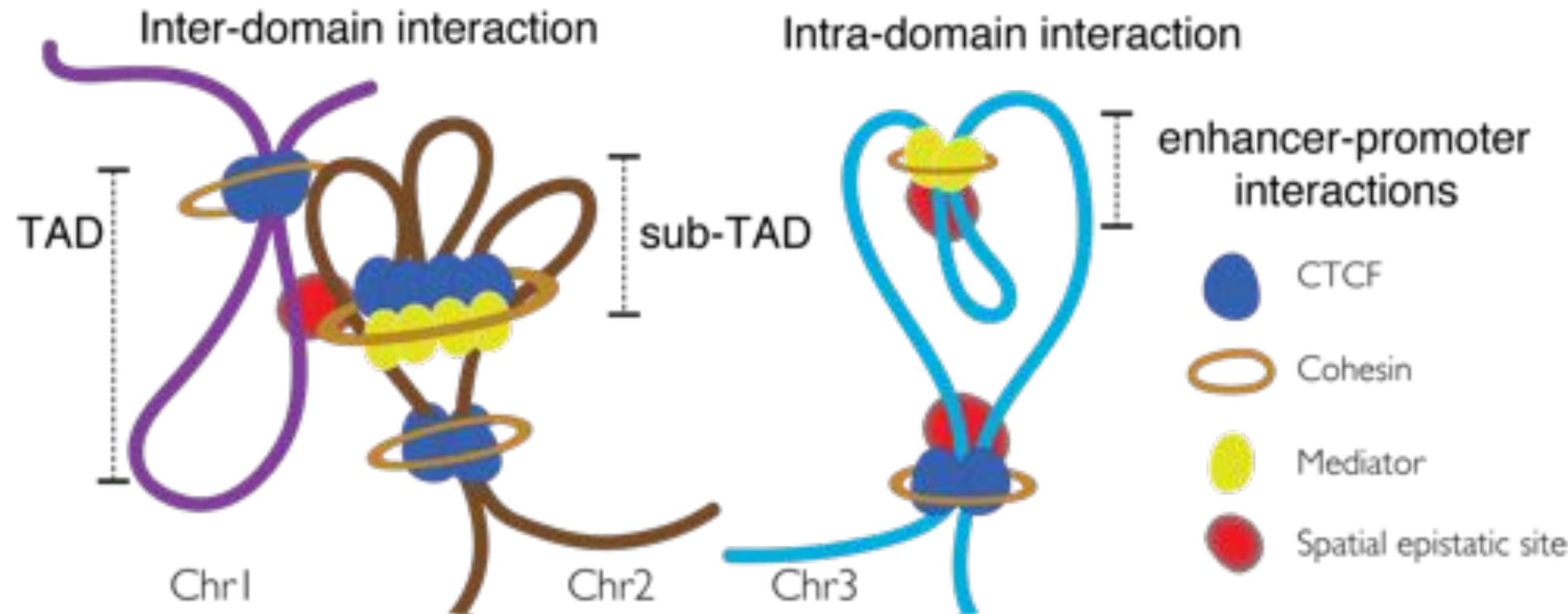
Spatial Epistasis Hypothesis

Spatial genomic interactions are enriched for genetic interactions

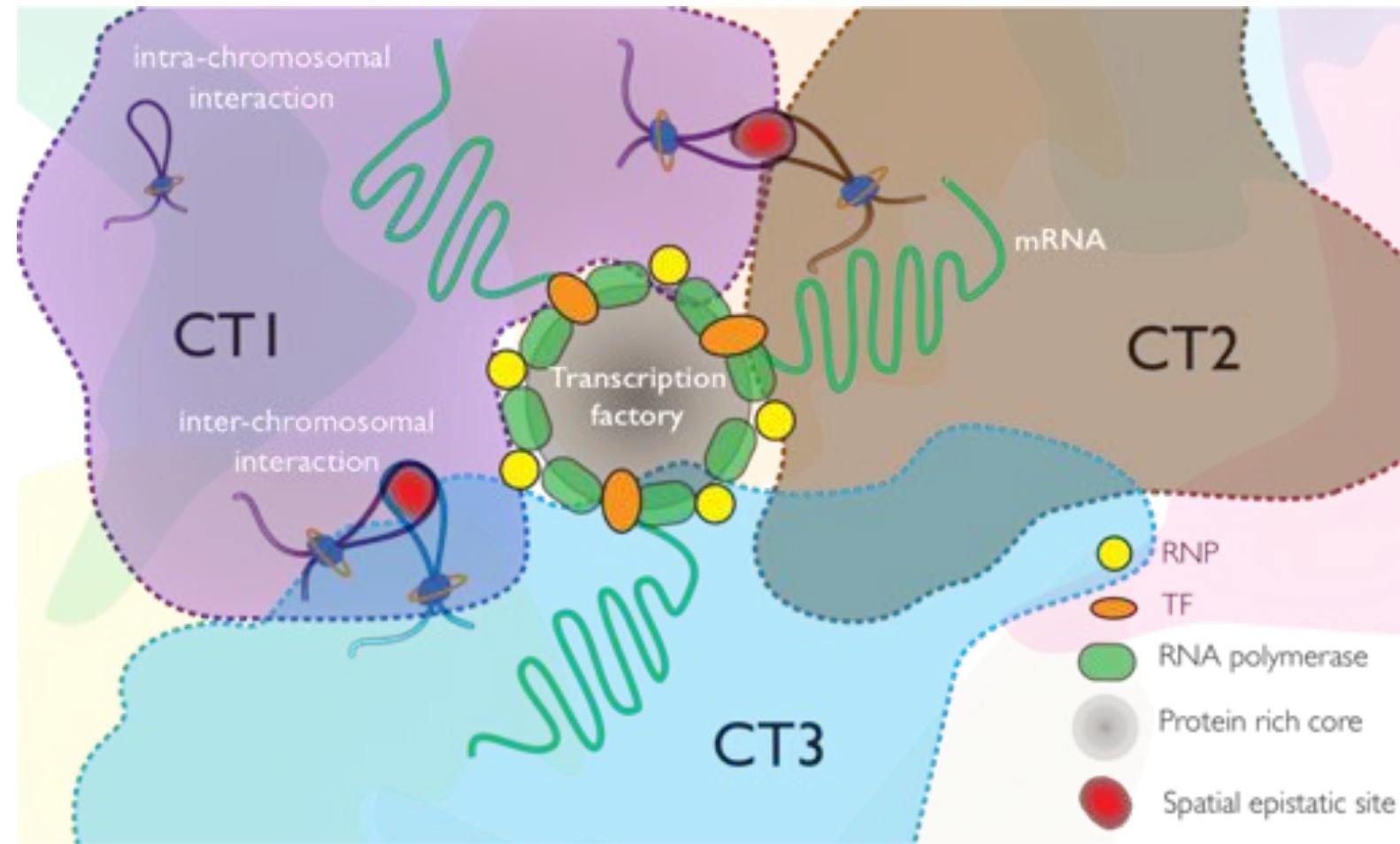
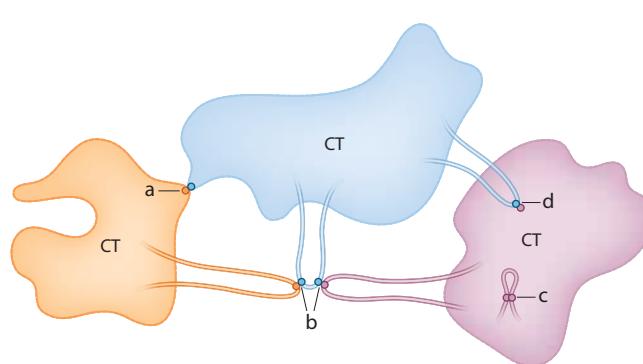
Spatial interactions via chromatin crosstalk



Biological Plausibility



Biological Plausibility



Outline

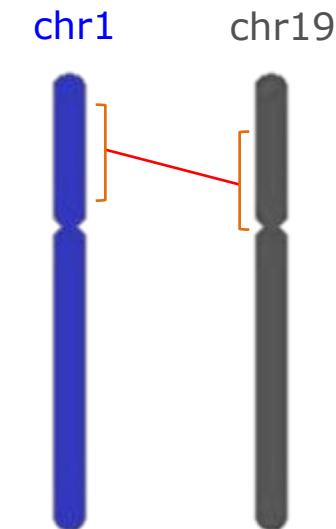
- Motivation
- Spatial Interactions (Genome Organization)
- Genetic Interactions (Epistasis)
- Spatial Epistasis Hypothesis
- **Methods and Computational Framework**
- Epistatic Signals
- “Meta-analysis”

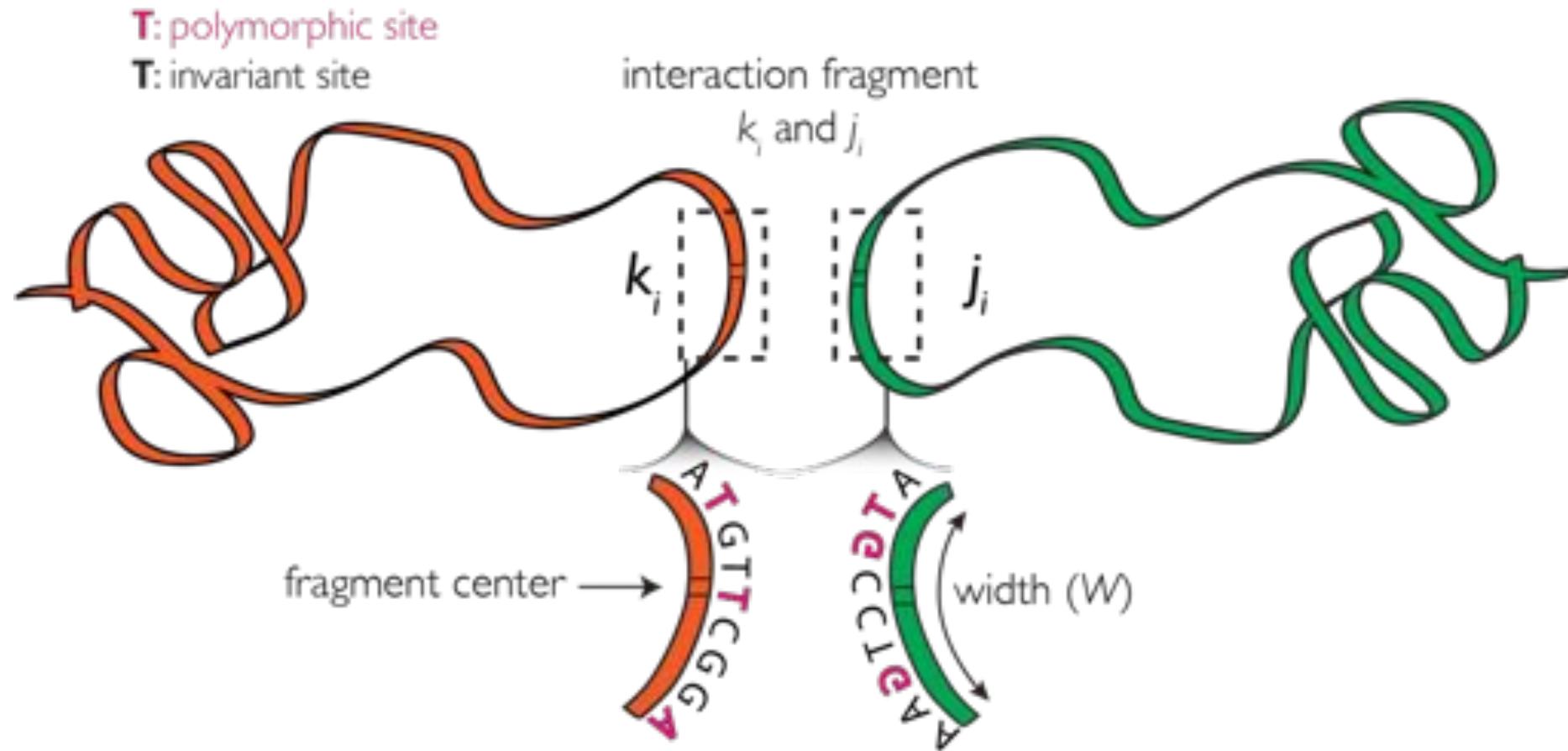
Methods: Hi-C data

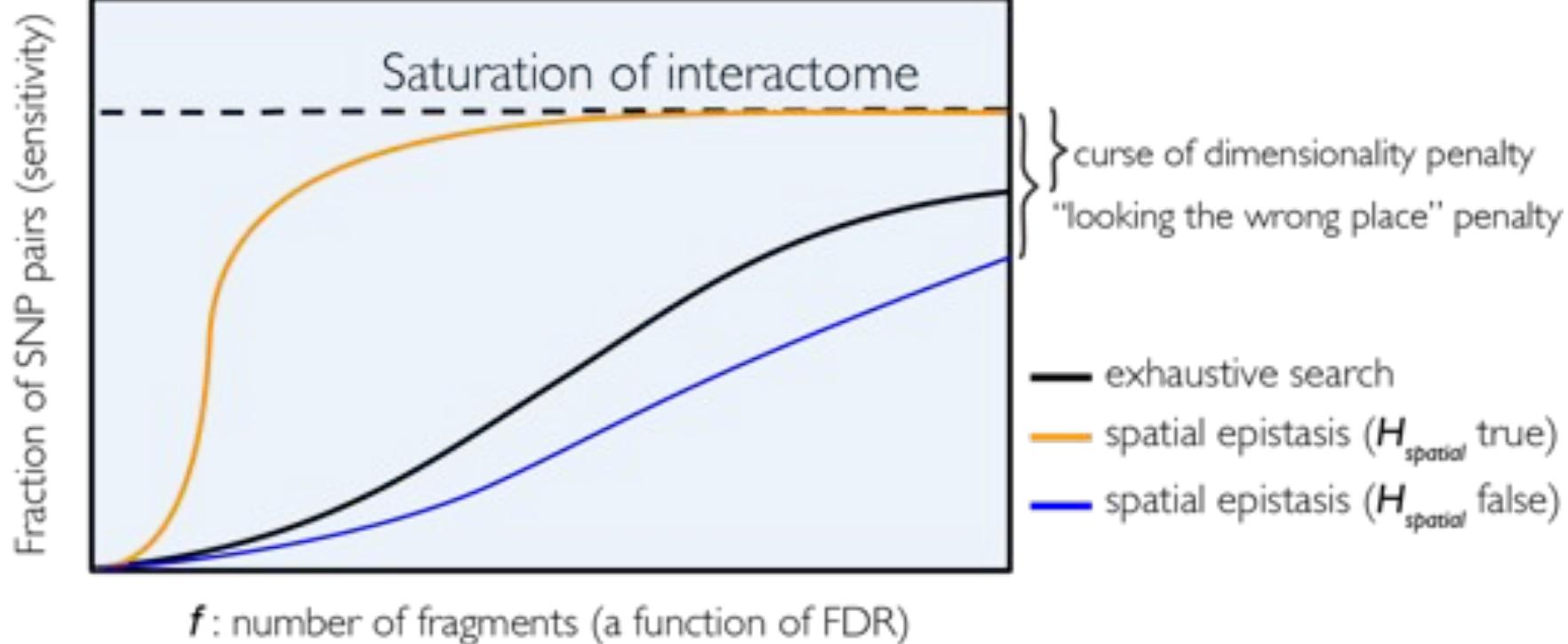
- Subset on inter-chromosomal interactions
 - Avoid problems with LD
 - Reduce the hypothesis space. Make calculations feasible.
- Select significance threshold (q-value) for interactions
- Select spatial interaction width (interaction width)

q-value cut-off →

partnerA	partnerB	contactCount	p.value	q.value (FDR)
3:185305000	10:107435000	21	5.82184E-87	1.57872E-77
10:100105000	12:55885000	8	4.37324E-29	6.50331E-21
4:187965000	10:33385000	8	8.96398E-29	1.21182E-20
3:190745000	10:107585000	8	1.62105E-28	2.01971E-20
2:25205000	10:118945000	8	3.25578E-28	3.68768E-20
10:126375000	11:77885000	7	6.64959E-28	6.92406E-20
10:11495000	16:4115000	7	1.13398E-27	1.10286E-19
10:27665000	11:6745000	8	1.18998E-27	1.14763E-19

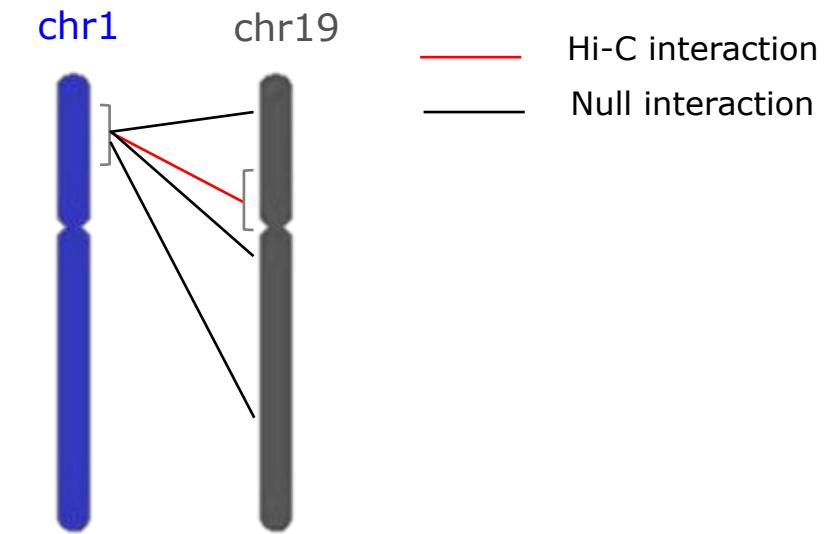






Null distribution

- Criteria for defining null distribution
 - Only inter-chromosomal
 - Distance criteria
- N=1000 permutations



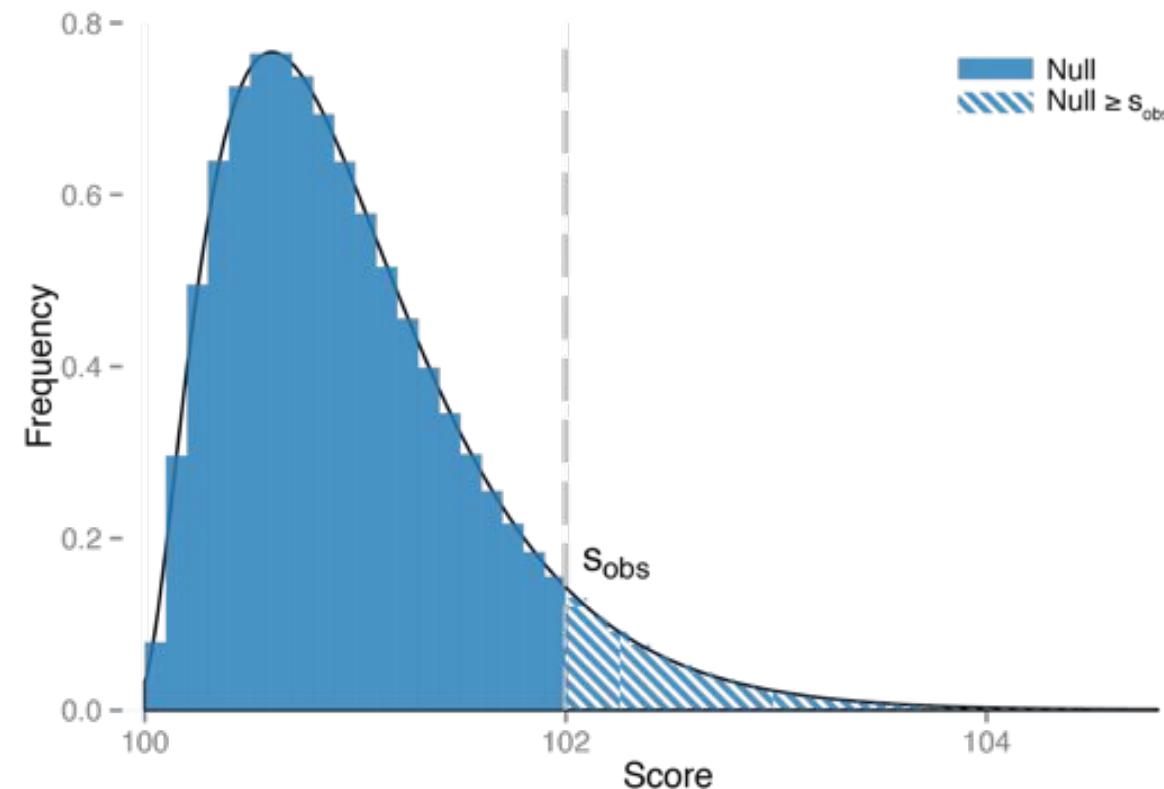
Empirical P -value

$$P_{empirical} = \frac{1 + \sum_{i=1}^N 1 \cdot [S_i \geq S_{obs}]}{N + 1}$$

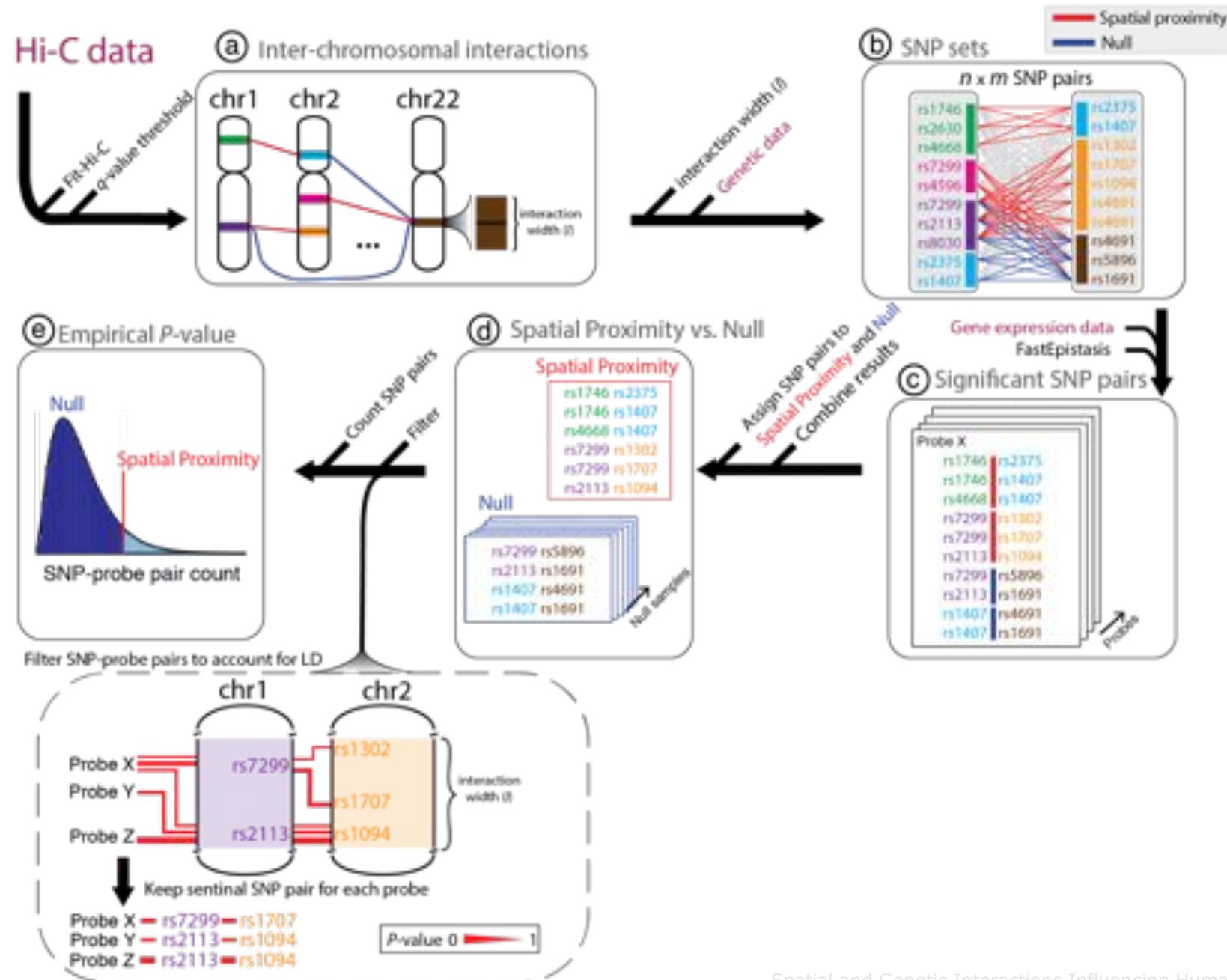
N: Null samples (permutations)

S_{obs} : Observed score

S_i : i^{th} null sample score



Pipeline



Key Data Sets

Genotype and Gene Expression Cohort

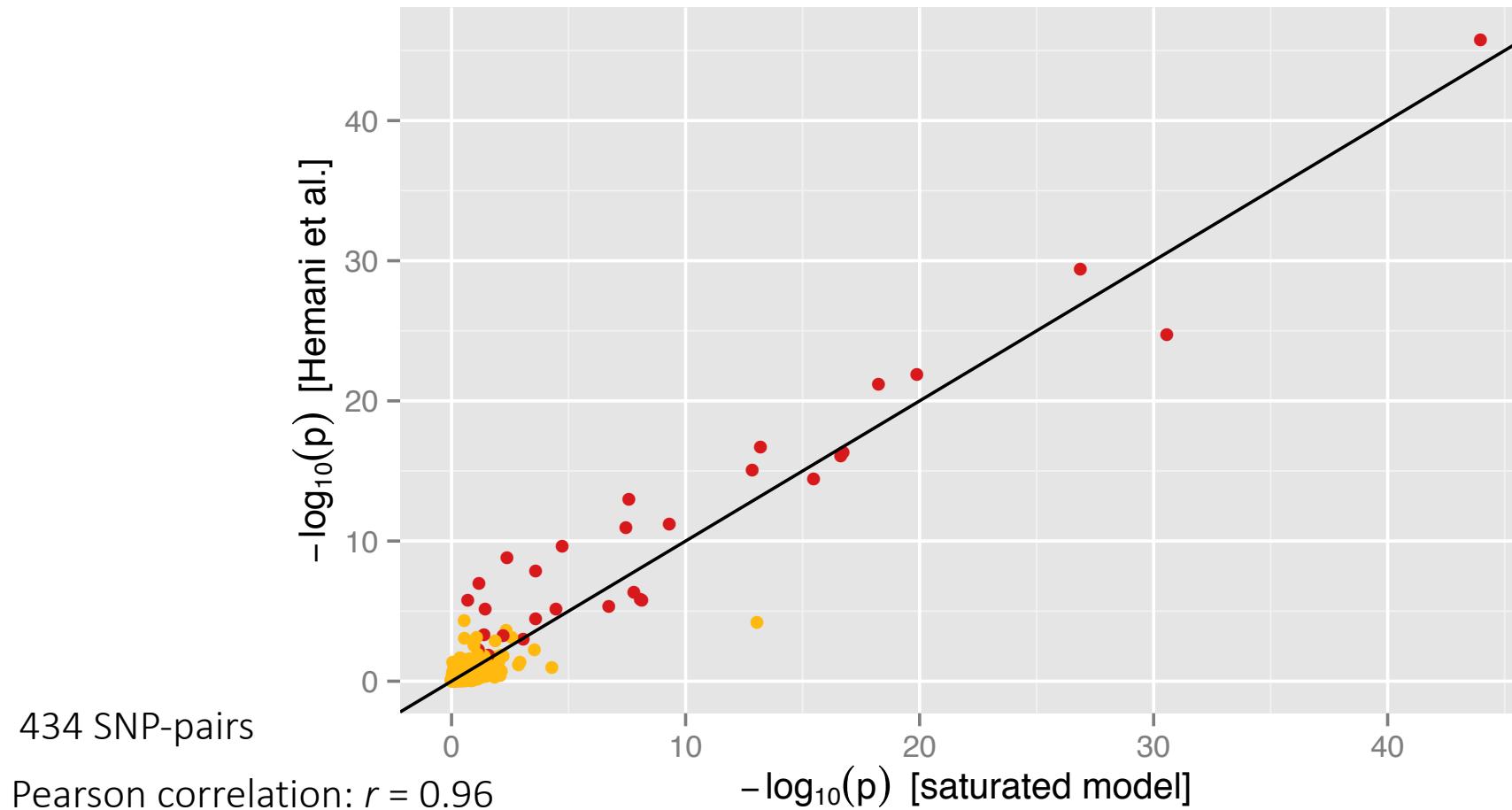
- Estonian Biobank (EGCUT); n=832
- Genotype data
 - Illumina HumanHap370CNV; Imputed to HapMap2
 - 2.5 million SNP
 - Standard QC thresholds (Anderson et al., 2010, *Nat. Protoc.*)
 - MAF \geq 5 %
- Expression data
 - Illumina HT-12 V3.0 platform
 - \sim 48,000 probes
 - Removed 50 (hidden) confounding factors using PEER (Stegle et al., 2012 *Nat. Protoc.*)
 - Probe selection strategy to select “best” \sim 10,000 probes

Hi-C Data

- Lieberman-Aiden, 2009, *Science*
 - K652 (bone marrow; erythroleukemia cell line) [$\sim 95 \times 10^3$ interactions]
Poisson regression model. Available from Lan et al., 2012, *Nucleic Acids Res*
- Dixon, 2012, *Nature*
 - hESC (embryonic stem cells) [90×10^6 interactions]
 - IMR90 (lung fibroblast) [160×10^6 interactions]
 - Fit-Hi-C method (Ay et al., 2014, *Genome Research*).

Epistatic Signals

Replication of Hemani (2014, *Nature*)



Spatial Epistasis Discovery

- Epistasis tests per probe: 1.33×10^{14}
- Single threaded CPU time: >12 years
- Only 0.35% of the possible two-locus SNP-pairs investigated
- Search space covered 6.62% of the genome

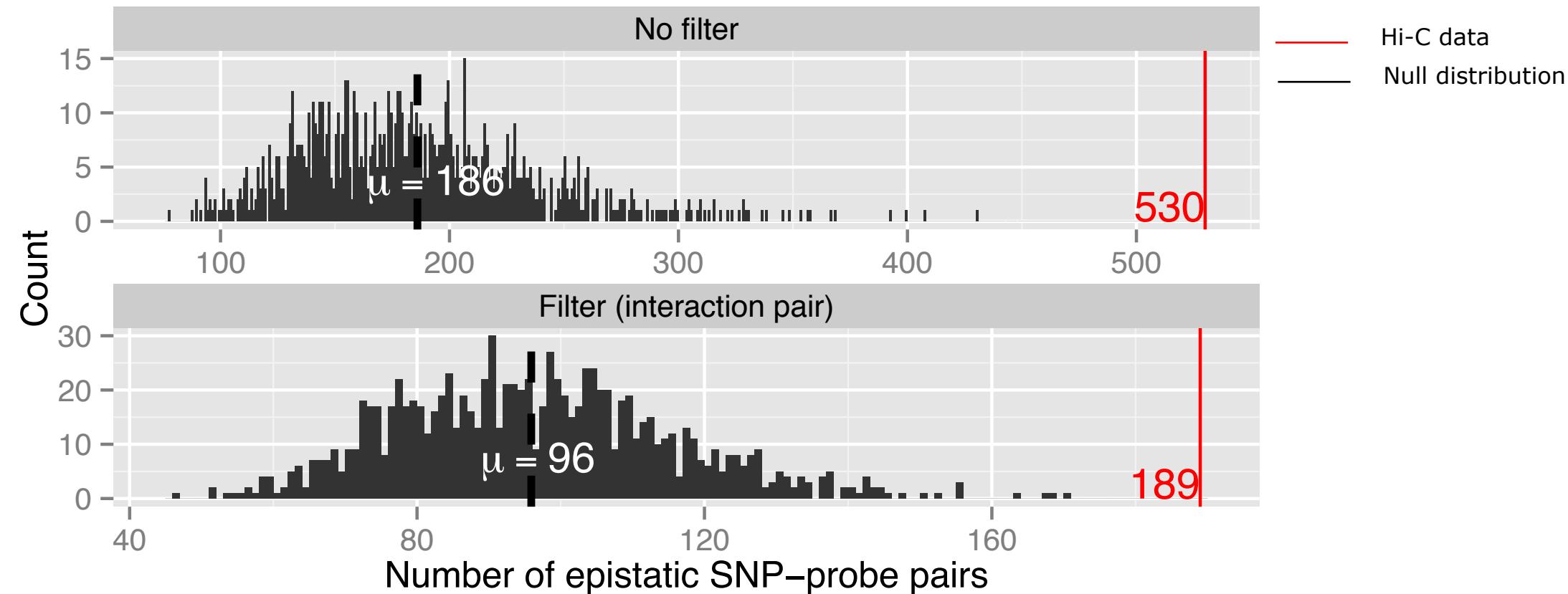
Spatial Epistasis Enrichment

Cell line	Interaction width, l	q -value ($n_{interactions}$)	No. epistasis tests <i>per probe</i> [†]	P -value [‡]
hIMR90	500 bp	10^{-6} (26, 325)	260,273,820	0.078
	1,000 bp	10^{-6} (26, 325)	1,042,826,407	0.001 *
	2,500 bp	10^{-7} (8, 114)	733,442,624	0.020 *
	50,000 bp	10^{-9} (1, 021)	4,284,422,018	0.927
hESC	500 bp	10^{-16} (2, 665)	3,222,025	0.257
	500 bp	10^{-14} (11, 919)	57,787,904	0.620
	1,000 bp	10^{-12} (39, 966)	2,236,957,368	0.643
	2,500 bp	10^{-13} (24, 084)	5,560,830,560	0.143
K562	1,000 bp	NA (1, 263)	1,485,960	0.540
	5,000 bp	NA (1, 263)	28,045,983	0.780
hIMR90 _{control}	1,000 bp	NA (5, 000)	45,668,788	0.169
hESC _{control}	1,000 bp	NA (5, 000)	45,714,894	0.317

[†]total number of statistical tests between interacting and non-interacting loci *per probe*.

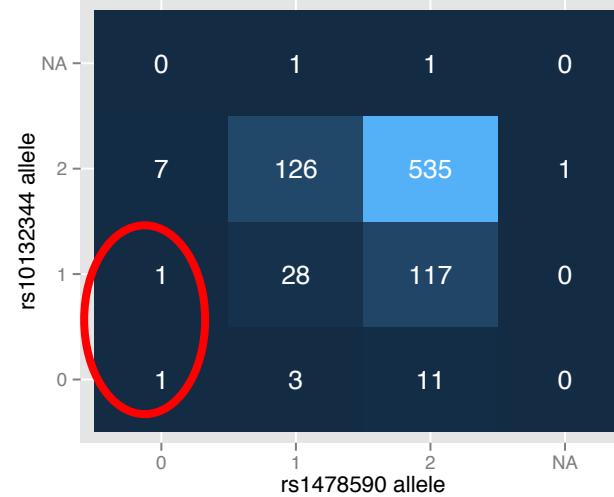
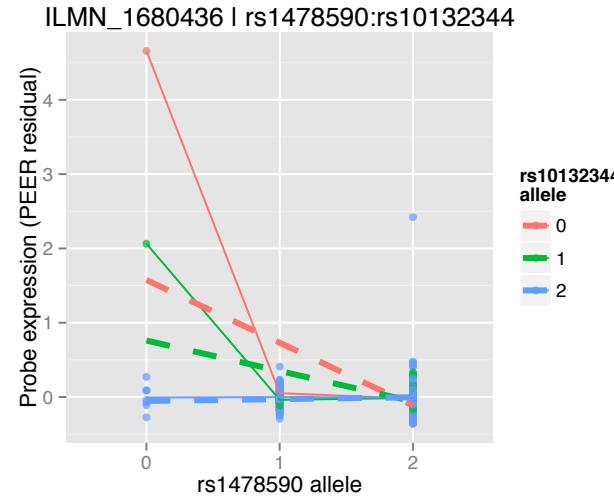
[‡] P -value derived from filtered SNP-probe pairs.

Strong Enrichment Signal for hIMR90 cells

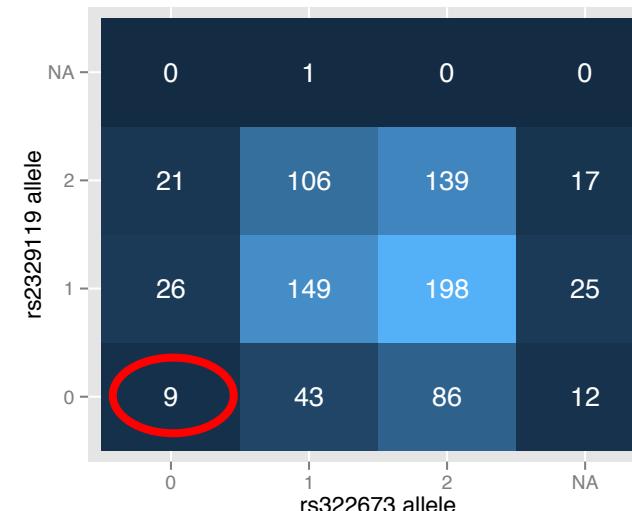
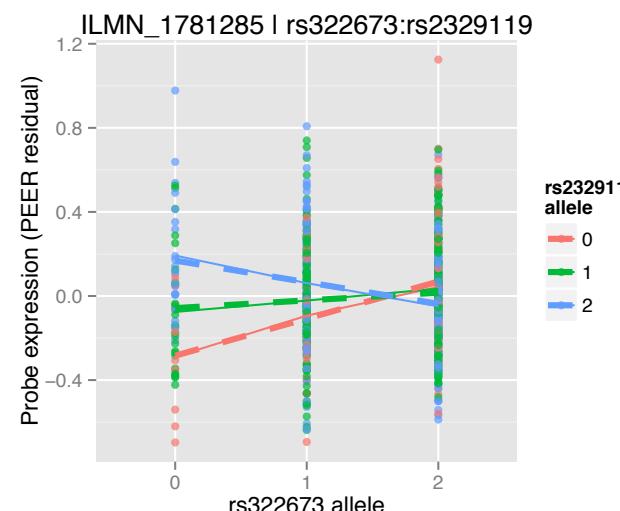


hIMR90[$l = 1,000\text{bp}$; $q = 10^{-6}$]

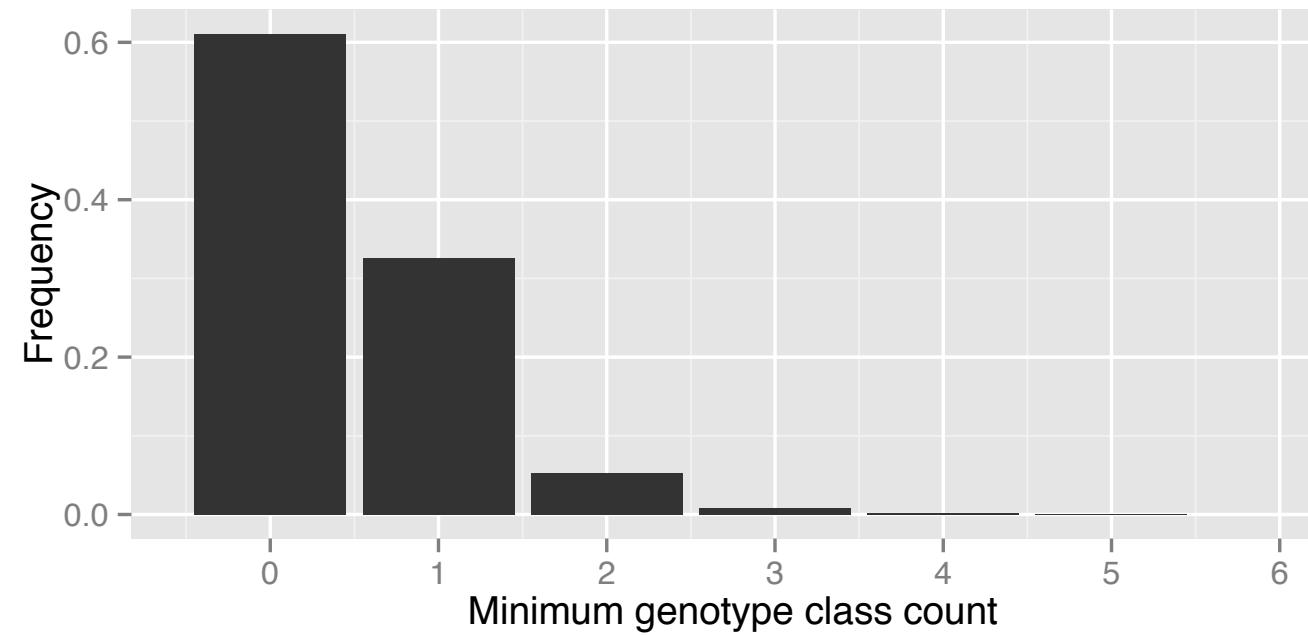
Spurious Epistatic Effects for Rare Variants


 $\beta_{\text{interaction}} = -0.44$
 $P\text{-value} = 1.29 \cdot 10^{-39}$

minimal Genotype Class Count = 1

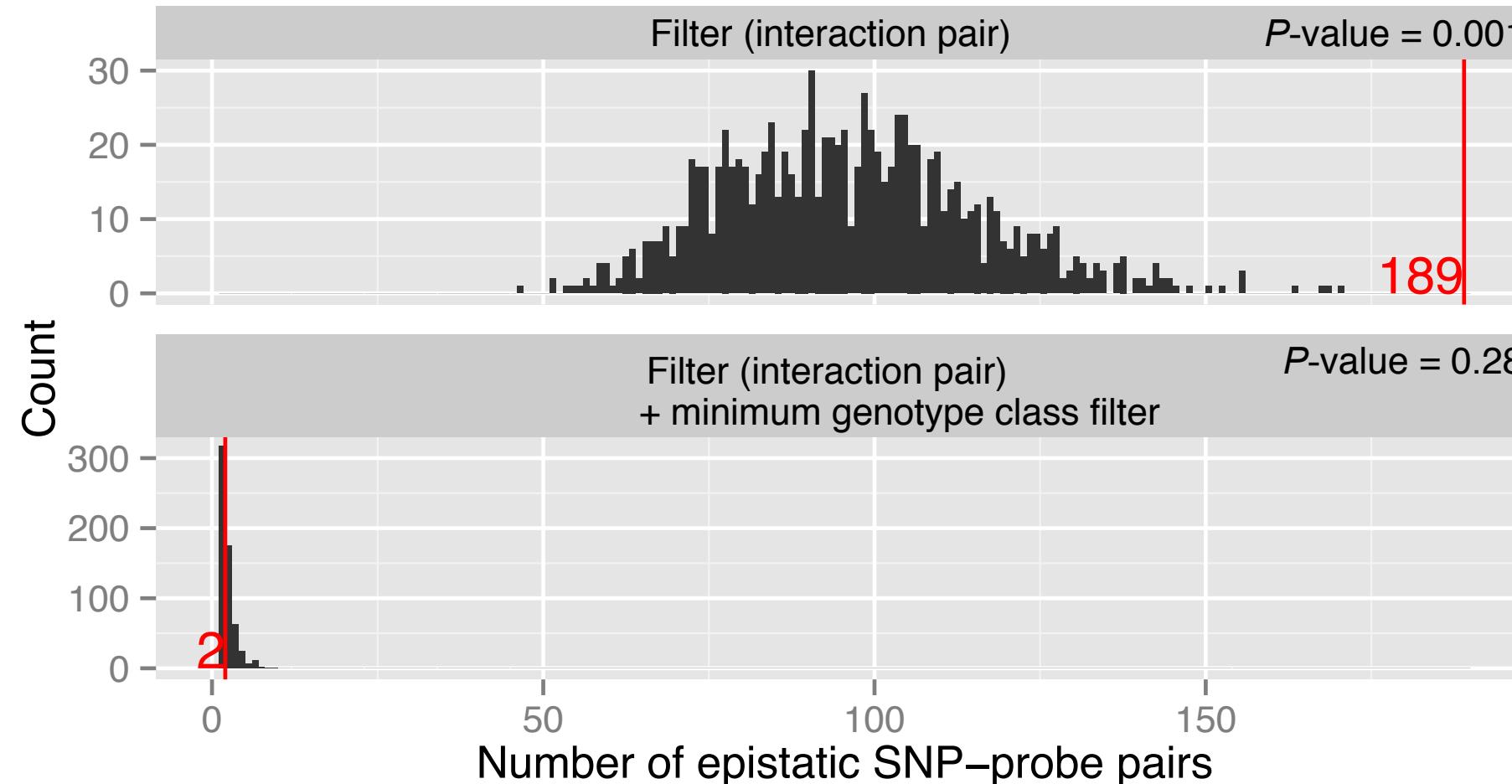

 $\beta_{\text{interaction}} = -0.14$
 $P\text{-value} = 4.07 \cdot 10^{-11}$
 minimal Genotype Class Count = 9

Distribution of Minimum Genotype Class Count



Most apparent epistatic SNP-pairs have minimum genotype class count < 3

Deeper Insights into the hIMR90 Spatial Epistasis Enrichment



Conclusion

Conclusion

- Developed potential powerful computational framework for testing biological hypothesis related to mechanisms of epistasis.
- Insufficient power to detect epistasis
- Cannot reject or accept the spatial epistasis hypothesis

Future Work

- Test for spatial epistasis enrichment
 - TAD boundaries
 - CTCF motifs

Meta-analysis

Meta-analysis: datasets



Connected Apps



Exist.io: Understand your life.

JAWBONE UP



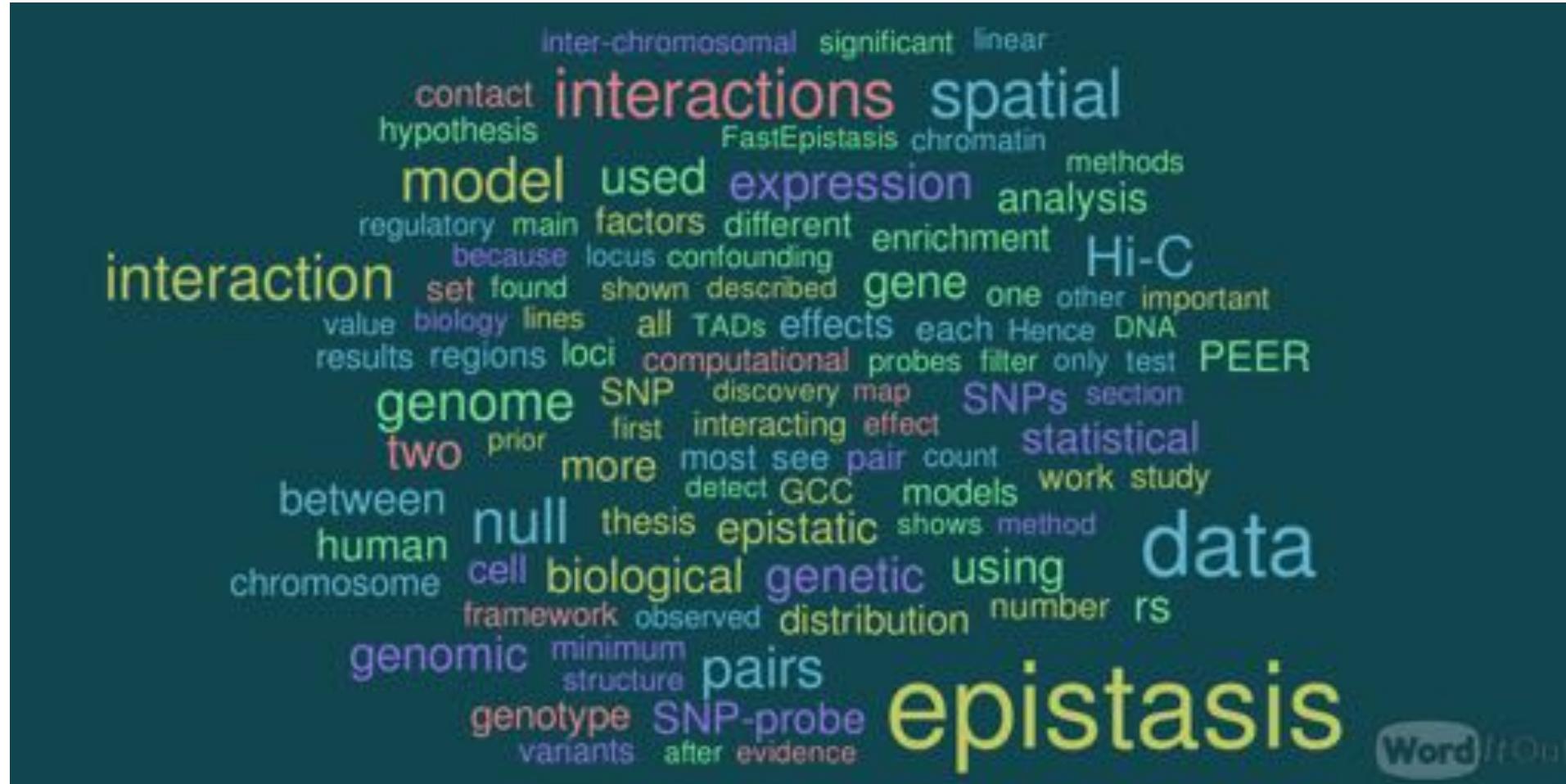
GARMIN



Google 31

Thesis and Thoughts

Two sides of the same coin?



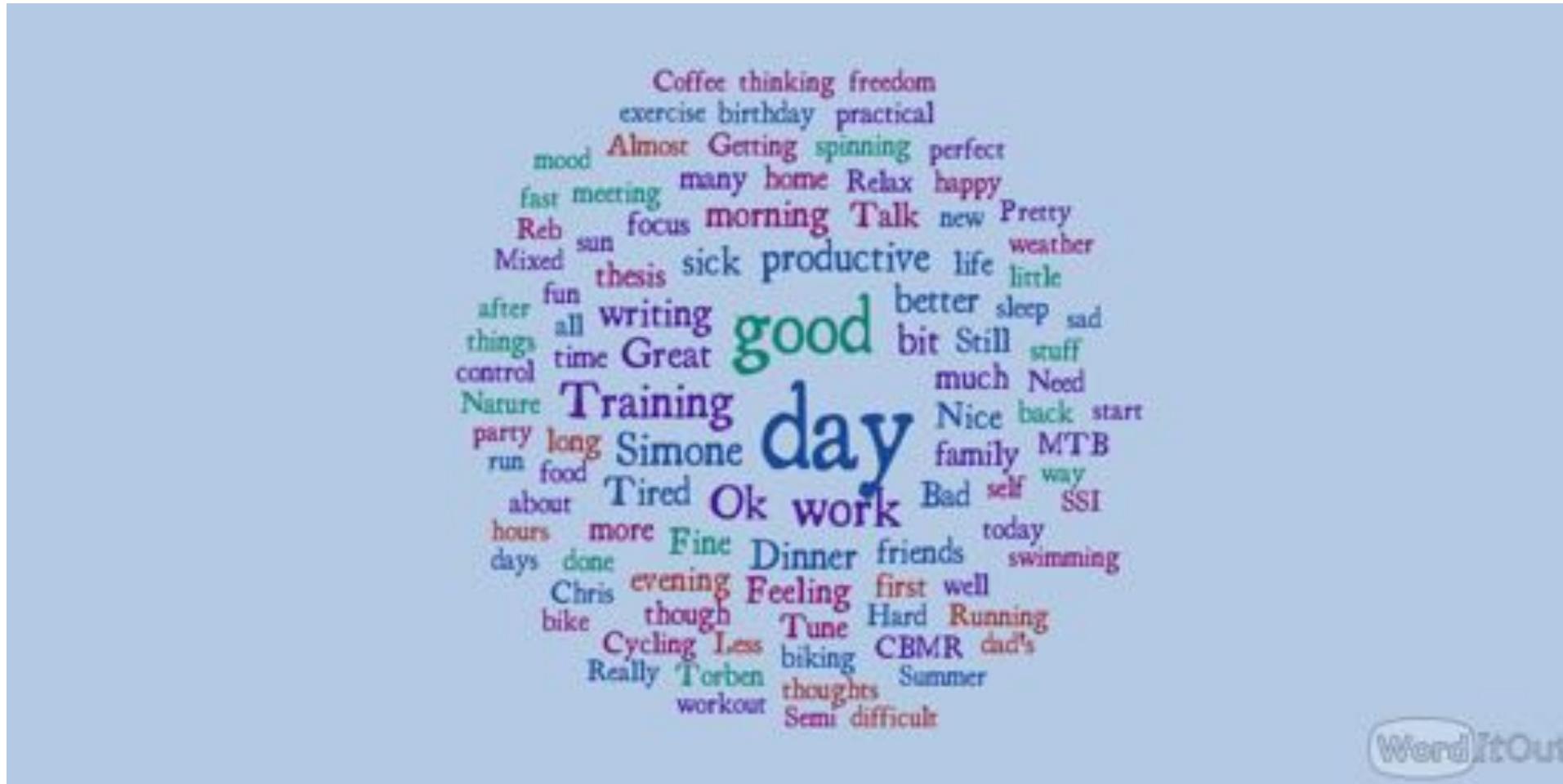
Thesis and Thoughts

Two sides of the same coin?



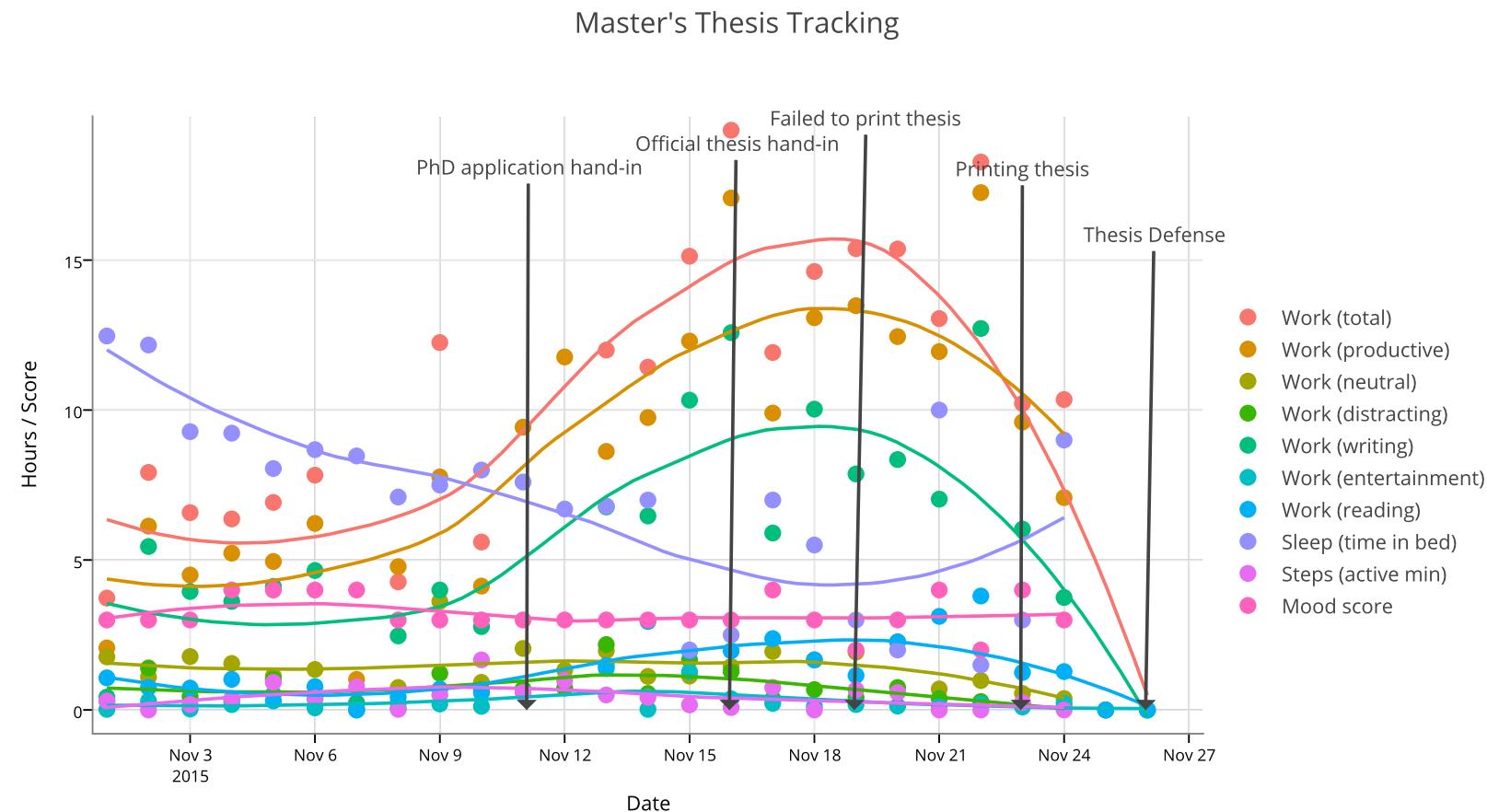
Thesis and Thoughts

Two sides of the same coin?

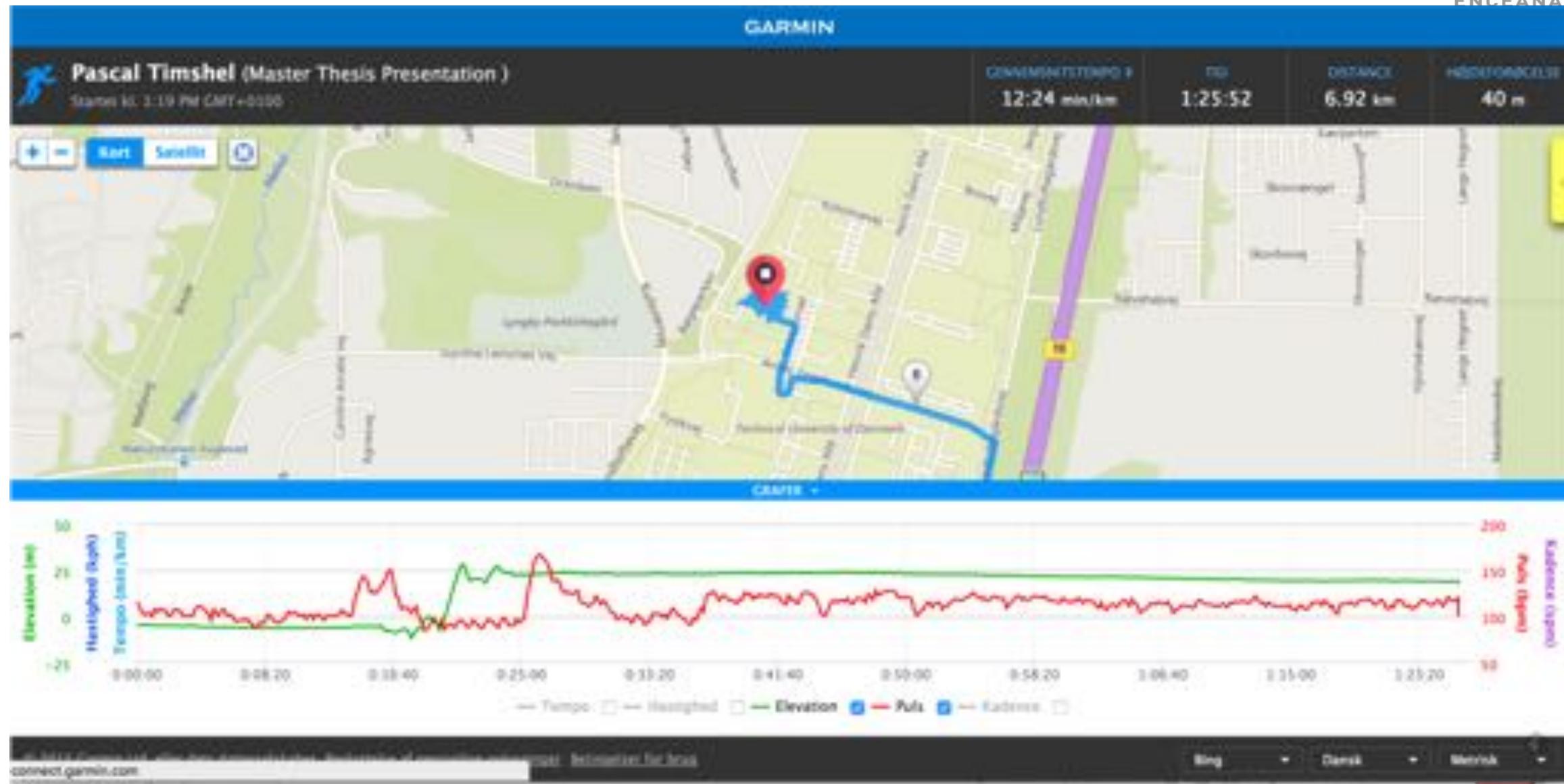


Thesis writing in numbers – a timeline

Your life before a deadline – or lack thereof



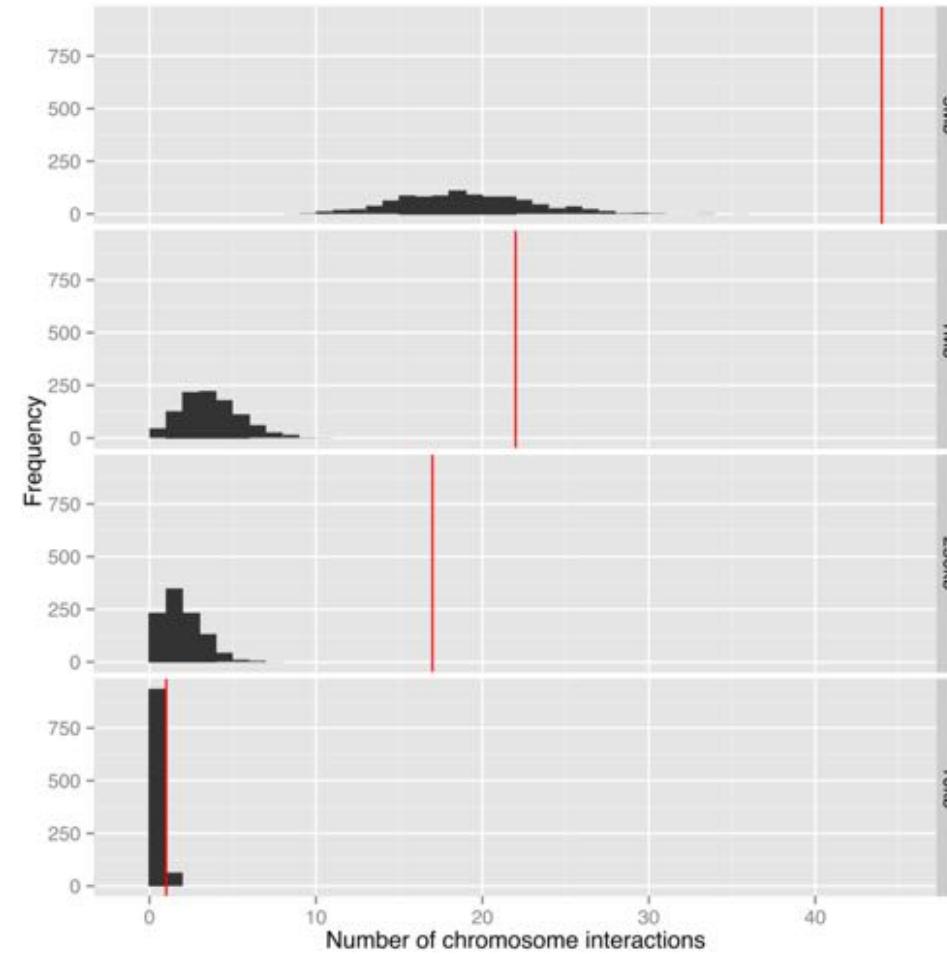




Appendix

Background: Hemani et al., 2014, *Nature*

- Models for biological interaction were examined
- Enrichment of epistatic SNP-pairs within 5 Mb of spatial interacting chromosomal regions.
- 44 out of 501 SNP-pairs co-locate ($P_{\text{empirical}} < 1.8 \times 10^{-10}$)



Epistasis table

SNP1	rsID	MAF1	rsID	MAF2	SNP2	P-value	Beta	minGCC	Gene	Probe	EIID	EID
8:18885955	rs10104492	0.28	rs12451549	0.10	17:46875212	1.72E-11	0.1731	3	CSHL1	17:59341967	hic_1_4837	hic_1
2:309677126	rs1862959	0.20	rs10500890	0.18	11:21075311	2.24E-13	0.18653	3	APOD	3:196777220	hic_1_1090	hic_1

Table B.4: Top confident spatial epistatic SNP-probe pairs identified in this study. That is, these SNP-pairs map to physically interacting genomic loci (in hIMR90 cells). The SNP-probe pairs were selected using the criteria $\text{min GCC} \geq 3$ (these two SNP-probe pairs appear on the right hand side of Table B.1). See Table B.3 for a description of the table headers.

SNP1	rsID	MAF1	SNP2	rsID	MAF2	P-value	Beta	minGCC	Gene	Probe	EIID	EID
4:104066247	rs3074608	0.46	3:163183258	rs4299495	0.24	2.57E-11	0.066	10	LOC727752	1951651	null_86_7175	null_86
3:57250044	rs9840963	0.32	22:47915370	rs80567	0.35	1.13E-11	0.086	12	ATPIF1	1:28436788	null_150_7035	null_150
9:81975039	rs7045654	0.32	16:12633764	rs7194011	0.42	2.64E-11	-0.078	19	SLC46A3	13:28172532	null_198_4075	null_198
1:176775553	rs6700296	0.35	7:825463554	rs6947662	0.4	3.10E-11	0.067	15	PIGR	1:205168794	null_420_1555	null_420
5:60860603	rs10036399	0.29	15:488368316	rs12911143	0.37	1.87E-11	0.064	11	KIAA1147	7:141000222	null_435_7721	null_435
2:135143088	rs4964158	0.39	16:43595417	rs35226	0.31	7.17E-12	-0.050	13	ABCCL1	16:46823387	null_928_6011	null_928

Table B.5: Top confident epistatic SNP-probe pairs identified in this study. The SNP-probe pairs were selected using the criteria $\text{min GCC} \geq 10$. These SNP-probe pairs are the strongest evidence of statistical epistasis found in this thesis. None of these SNP-probe pairs were reported by (Hemani et al., 2014). Notice that only two of the SNP-pairs co-localize with the probe. All the SNP-pairs map to non-interacting genomic loci ("null samples", see Figure 6.2). See Table B.3 for a description of the table headers.

Epistasis table

SNP1	rsID	MAF1	SNP2	rsID	MAF2	P-value	Beta	minGCC	Gene	Probe	EIID	EID
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	2.63E-23	-0.5636	0	ITK	5:156614448	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	1.46E-26	0.5306	0	CHEBIL2	7:137210715	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	8.02E-18	0.5514	0	FAIM3	1:205144550	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	1.07E-31	0.5051	0	FGR2	33:61242090	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	7.56E-25	0.5474	0	LRMP	12:25152225	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	2.00E-110	1.2144	0	KIAA0672	17:12835495	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	8.45E-17	0.3753	0	KIAA1407	3:115166040	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	6.25E-13	0.4408	0	VIL2	6:1559106839	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	1.63E-11	0.6132	0	FCRL2	3:1555983443	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	2.35E-13	0.3674	0	SIASH1	35:60052618	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	4.57E-18	0.4974	0	COLBA2	1:80538993	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	7.00E-36	0.5992	0	PMOD	1:201576664	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	1.77E-30	0.5563	0	GPT2	36:40522097	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	2.80E-18	0.4203	0	STARD7	2:96214694	null	357_4752
5:152894308	rs10041179	0.07	6:8085632	rs9294148	0.07	1.01E-25	0.6053	0	PCRL5	3:1557506161	null	357_4752
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	9.17E-52	0.2780	1	ASS1	9:132310134	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	1.20E-14	0.1545	1	RHOHBTB1	30:62373468	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	8.18E-21	0.1586	1	DGIP3	3:109996100	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	4.13E-15	0.1367	1	FAM80A2P	8:7104545	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	1.73E-21	0.2167	1	GLP1R	6:39163079	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	1.40E-30	0.1971	1	TXNRD2	22:18299943	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	1.29E-39	0.4352	1	CSEHL1	17:59941967	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	8.18E-17	0.1293	1	NCOR2	12:123077955	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	4.90E-15	0.1320	1	NUTT4P1	1:142848976	null	339_6338
2:218646715	rs1478590	0.11	14:69666638	rs10032344	0.11	1.52E-27	0.1737	1	Clef500	3:9037055	null	339_6338
3:25307427	rs923673	0.26	19:79204571	rs2329119	0.42	4.079E-11	-0.1306	9	DUSP1	5:1721284083	null	974_2900

Table B.3: Examples of epistatic effects with different levels of confidence (measured by minimum GCC). The SNP-probe pairs are sorted by their minimum GCC. The table presents data for three exemplary SNP-pairs. SNP-probe pairs marked in bold are plotted in Figure B.6 on page 87. This table serves to illustrate some of the spurious epistatic effects ($\text{minGCC} < 3$) underlying the hIMR90 spatial epistasis enrichment. P-value: 1-df test for the significance of the genetic interaction term; Beta: effect size estimate of the genetic interaction term; Gene: the gene-probe mapping listed in the Illumina manifest file; Probe: chromosomal coordinates of the probe; EID: Experiment Identifier; EIID: Experiment Interaction Identifier.

PEER Factor Correlation

