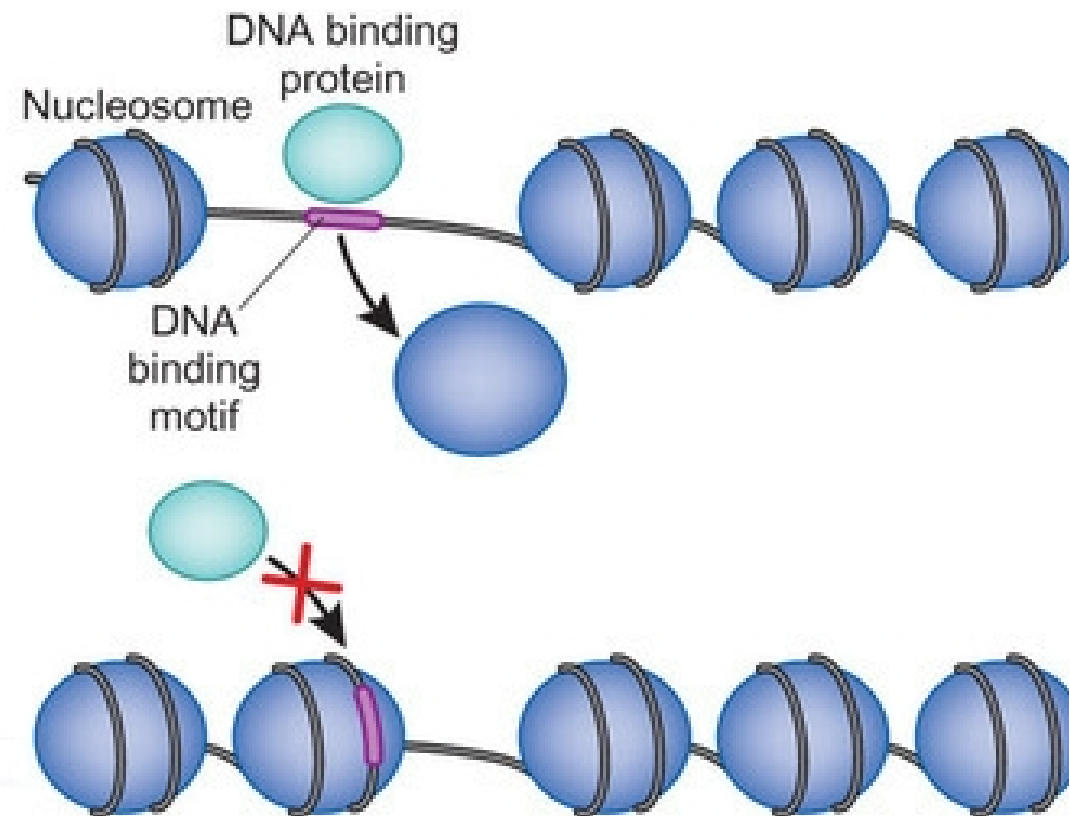


Transcription Factor Binding at Specific and Ubiquitous Open Chromatin Sites Across Multiple Human Cell Types

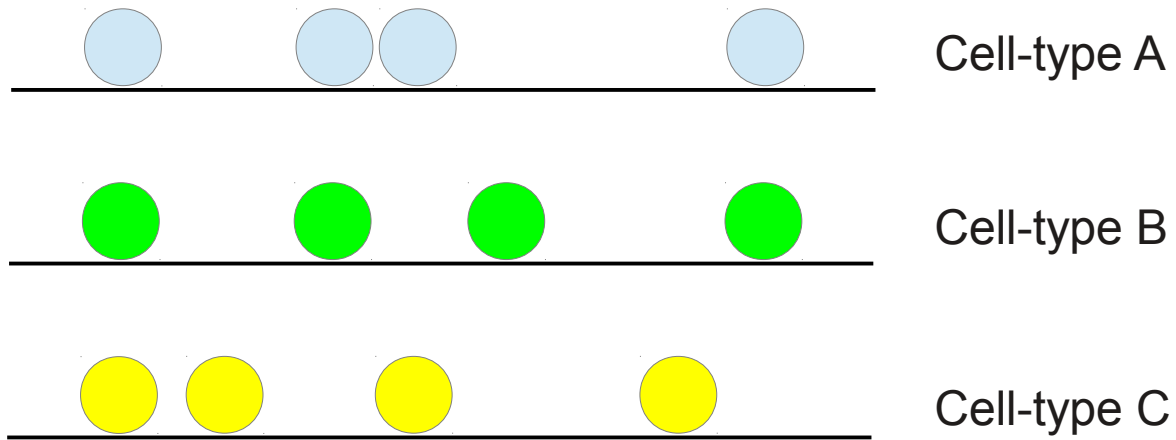
Martin Hemberg
RECOMB Conference
on Systems Biology and Regulatory Genomics
San Francisco
11/13/12

Open chromatin are the accessible locations in the genome

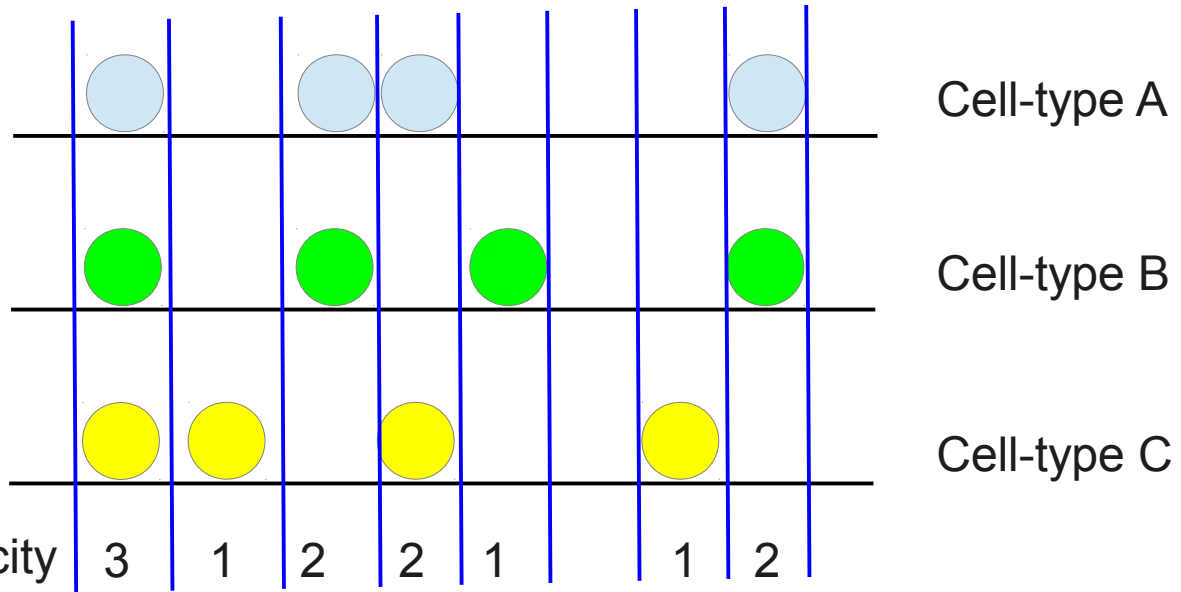
- DNase I Hypersensitive Sites (DHSs)
- Overlap genes and regulatory elements
- ~100,000/cell-type
 - Cell-type specific



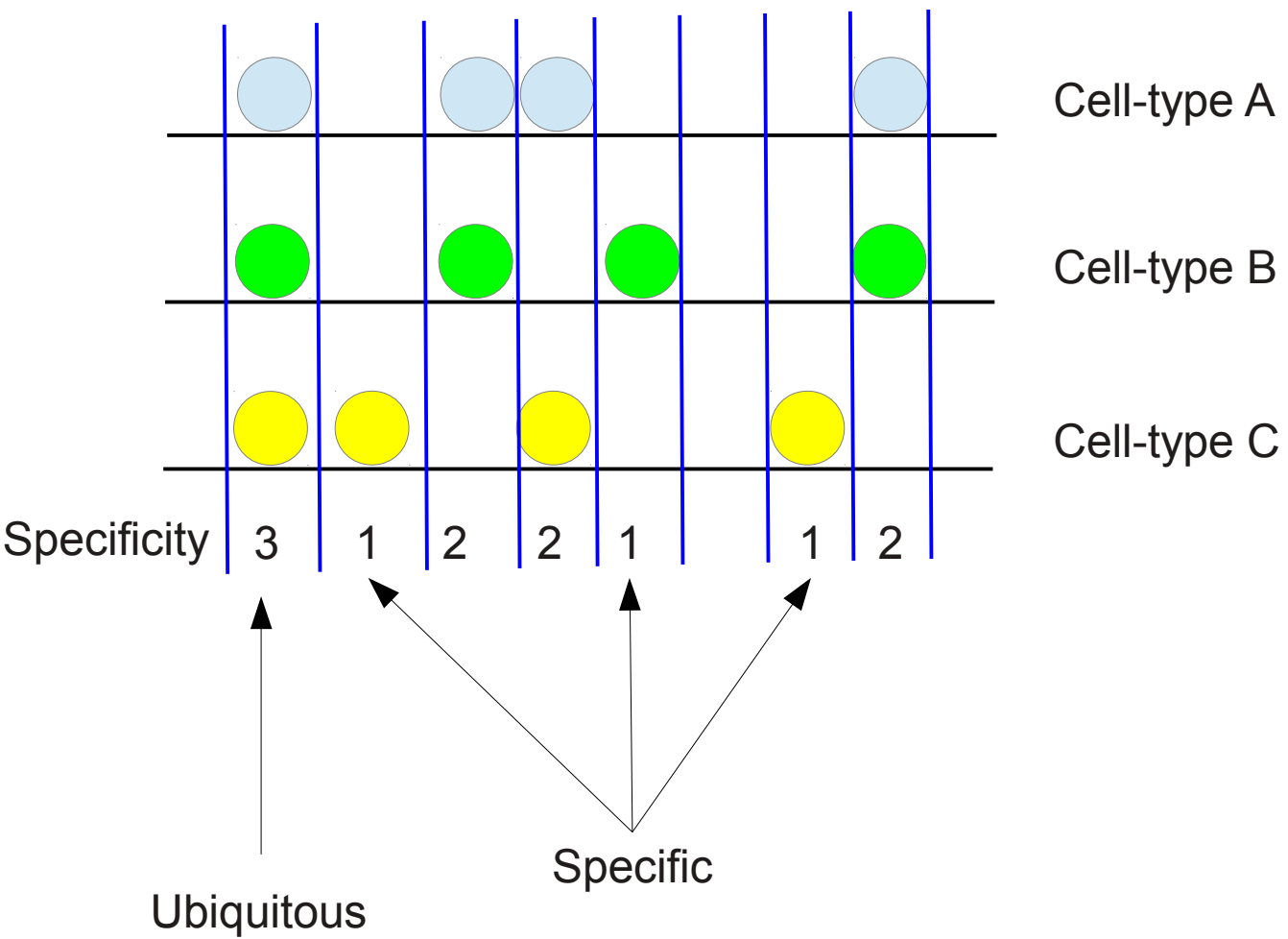
DHS specificity



DHS specificity

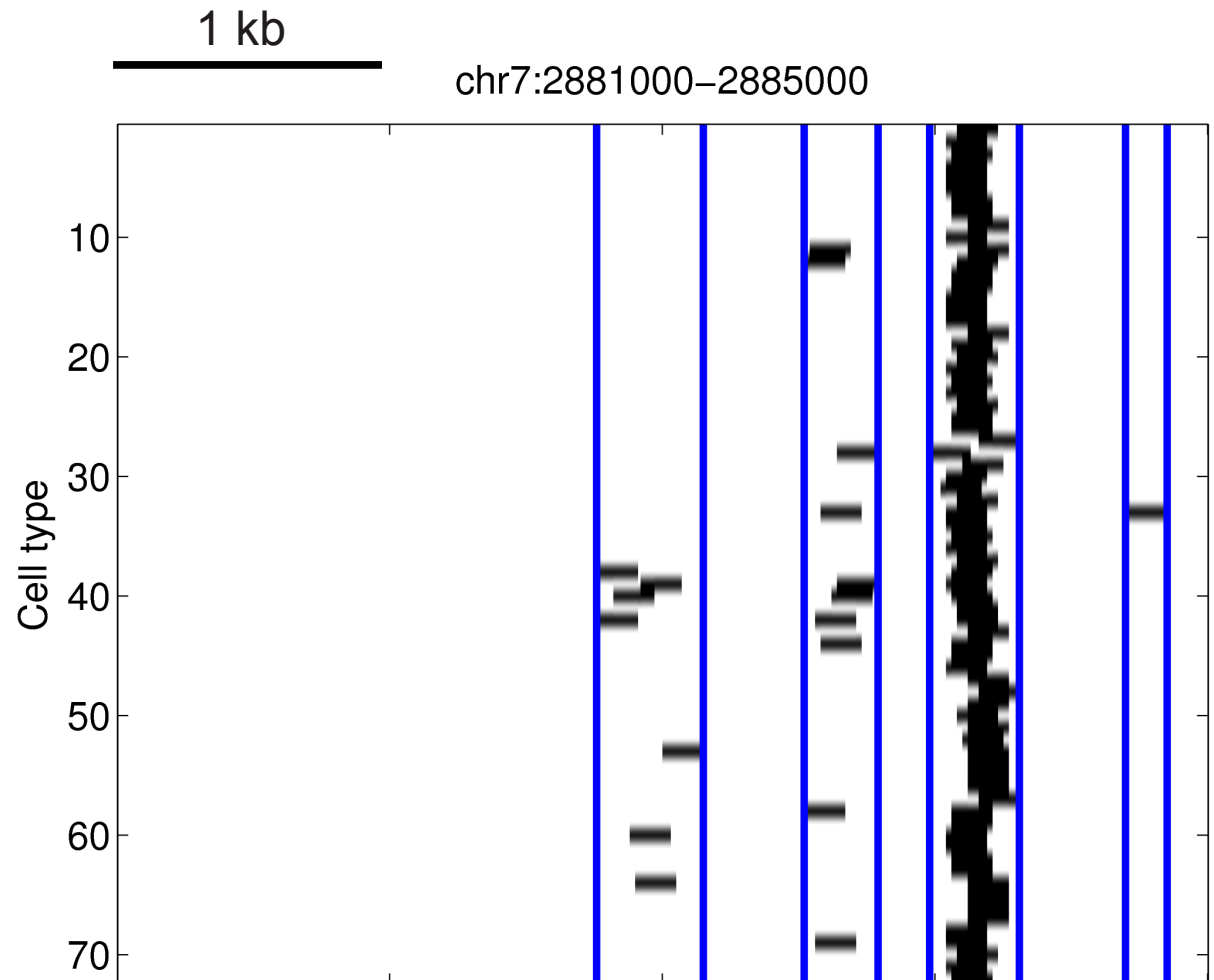


DHS specificity



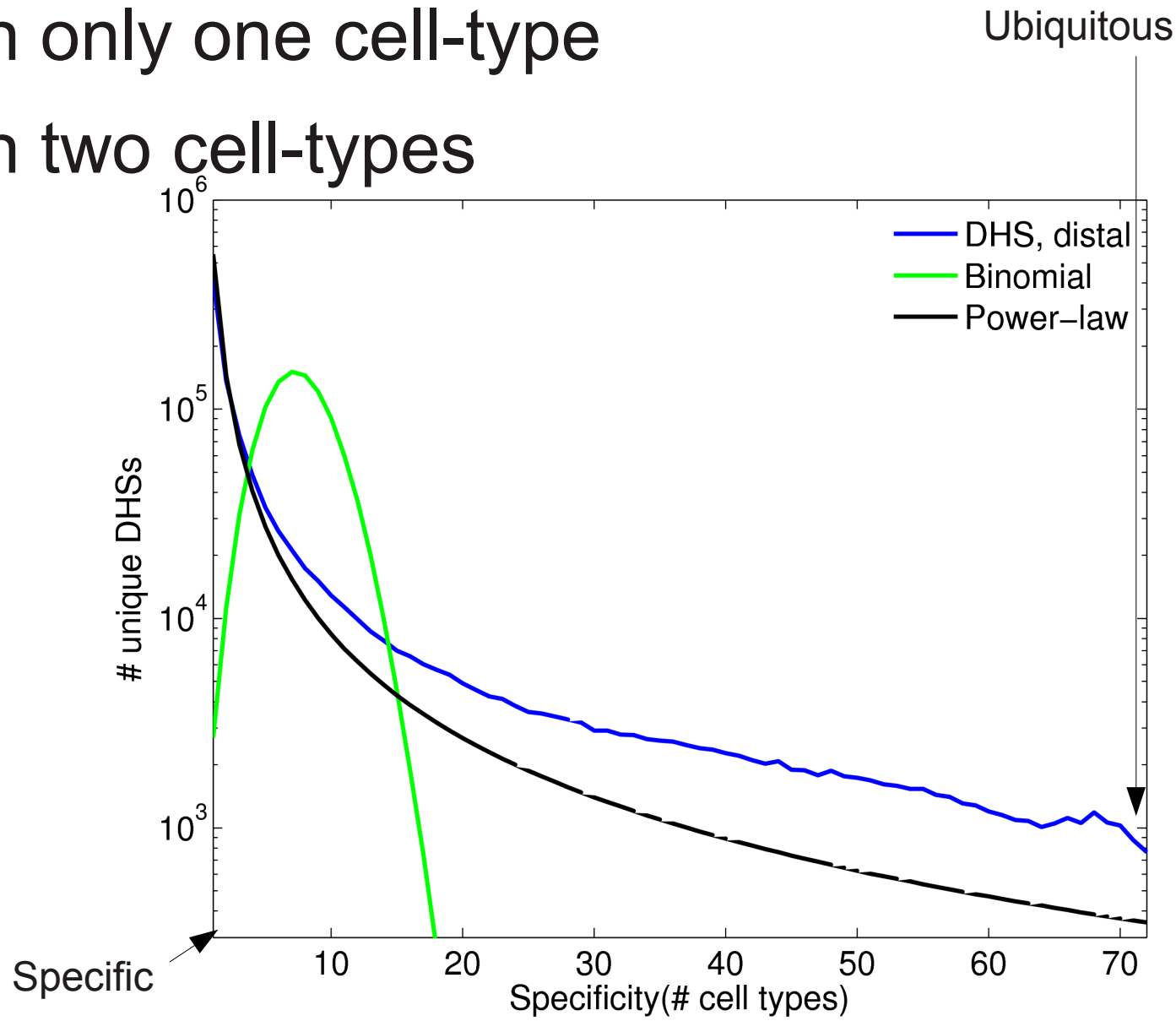
DHS Specificity, 72 cell-types

- ~1 million unique loci, >5 kb from TSS

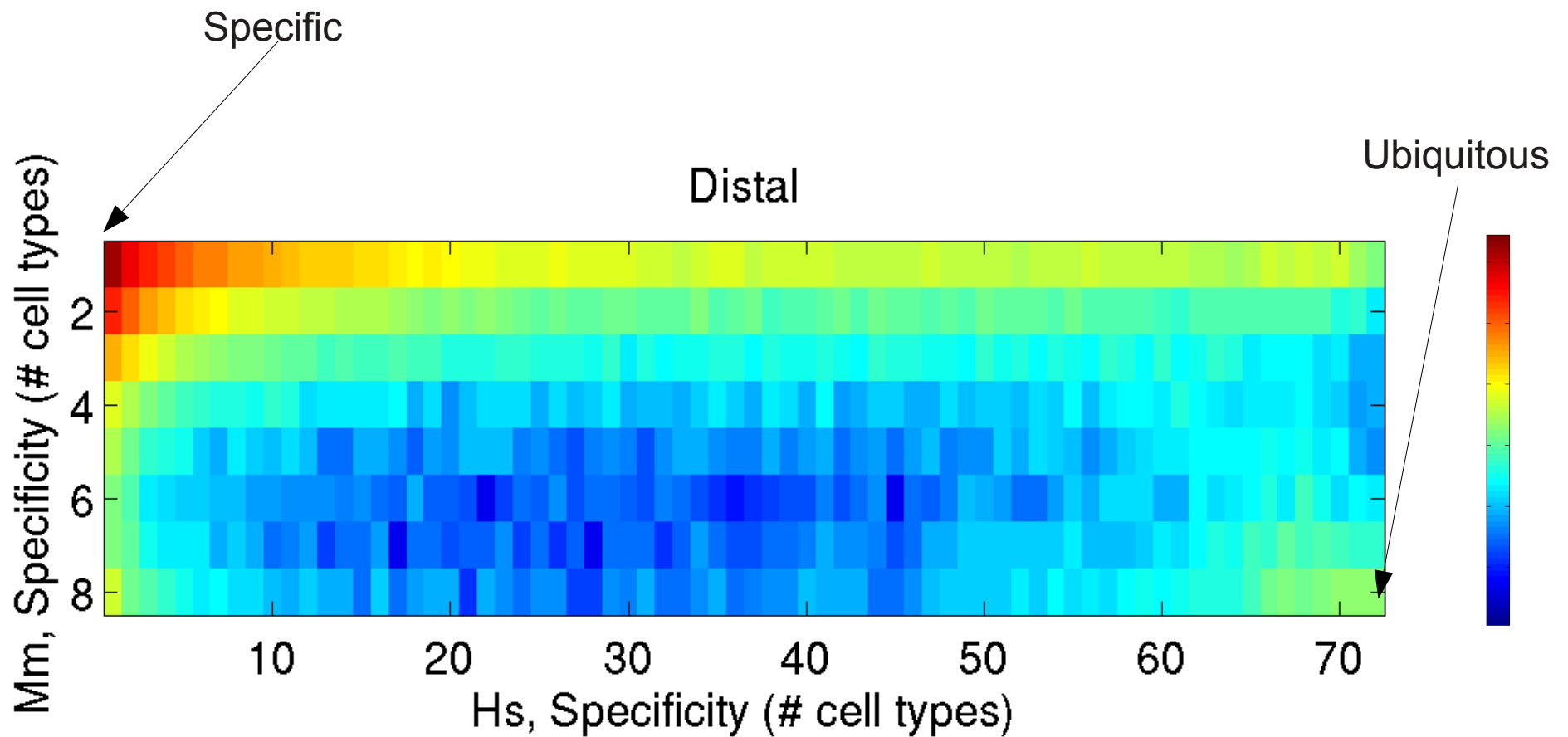


DHS Specificity follows a power-law

- ~42% found in only one cell-type
- ~14% found in two cell-types
- Slope ~ -1.35

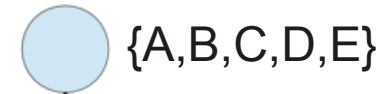


Specificity conserved in mouse



Why do we observe the power-law?

- Start with one cell-type with N DHSs

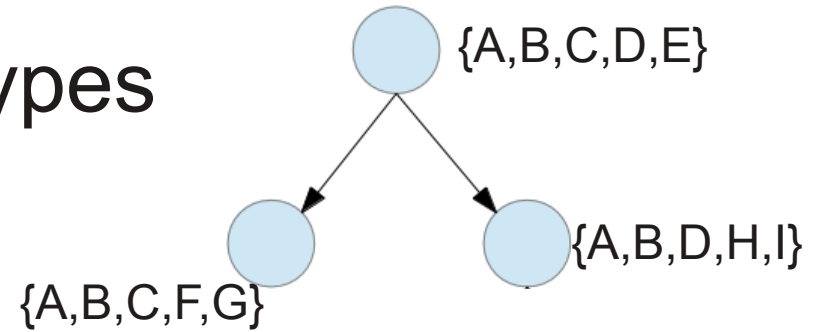


A: 1
B: 1
C: 1
D: 1
E: 1

Why do we observe the power-law?

- Start with one cell-type with N DHSs
- Generate two new cell-types

$$P(\text{keep DHS}) = \frac{k^b}{1 + k^b}$$



- k - # cell-types where DHS present
- b – selective advantage of existing DHSs
- Add new DHSs to get N

A: 3
B: 3
C: 2
D: 2
E: 1
F: 1
...

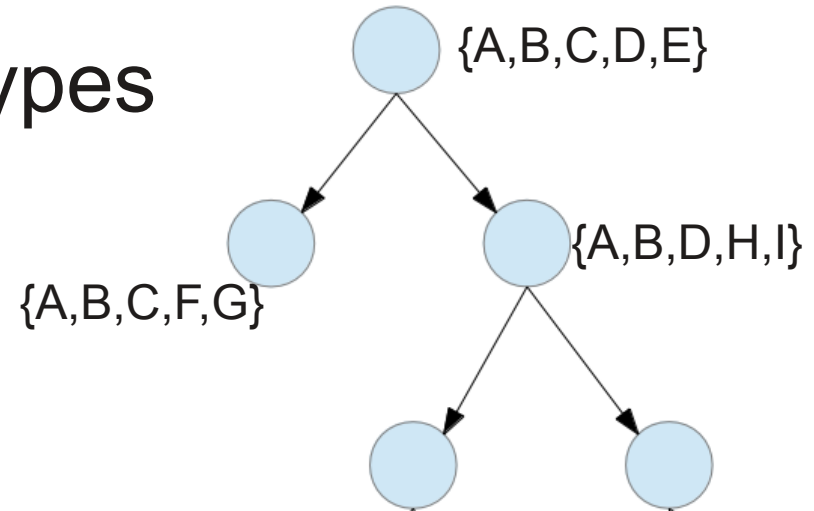
Why do we observe the power-law?

- Start with one cell-type with N DHSs
- Generate two new cell-types

$$P(\text{keep DHS}) = \frac{k^b}{1 + k^b}$$

- k - # cell-types where DHS present
- b – selective advantage of existing DHSs

– Add new DHSs to get N



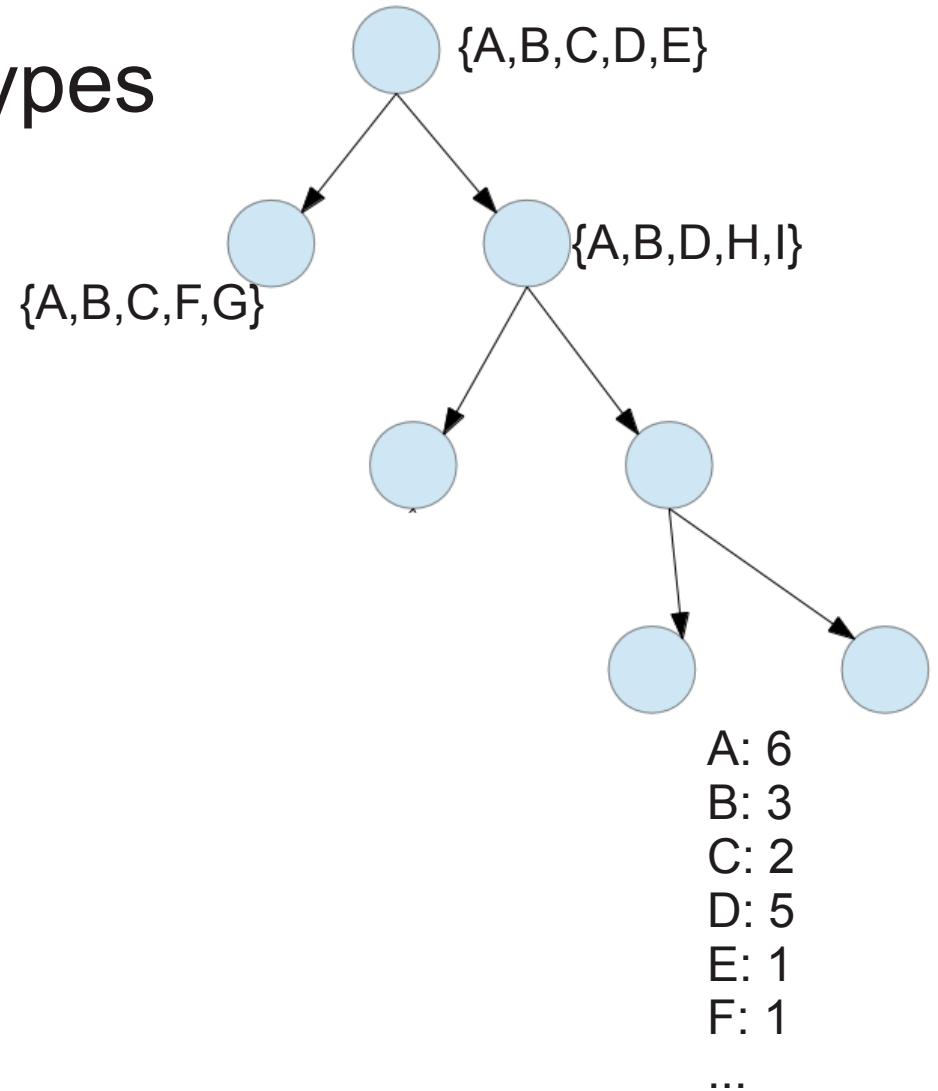
A: 4
B: 3
C: 2
D: 3
E: 1
F: 1
...

Why do we observe the power-law?

- Start with one cell-type with N DHSs
- Generate two new cell-types

$$P(\text{keep DHS}) = \frac{k^b}{1 + k^b}$$

- k - # cell-types where DHS present
- b – selective advantage of existing DHSs
- Add new DHSs to get N
- Expand until C nodes

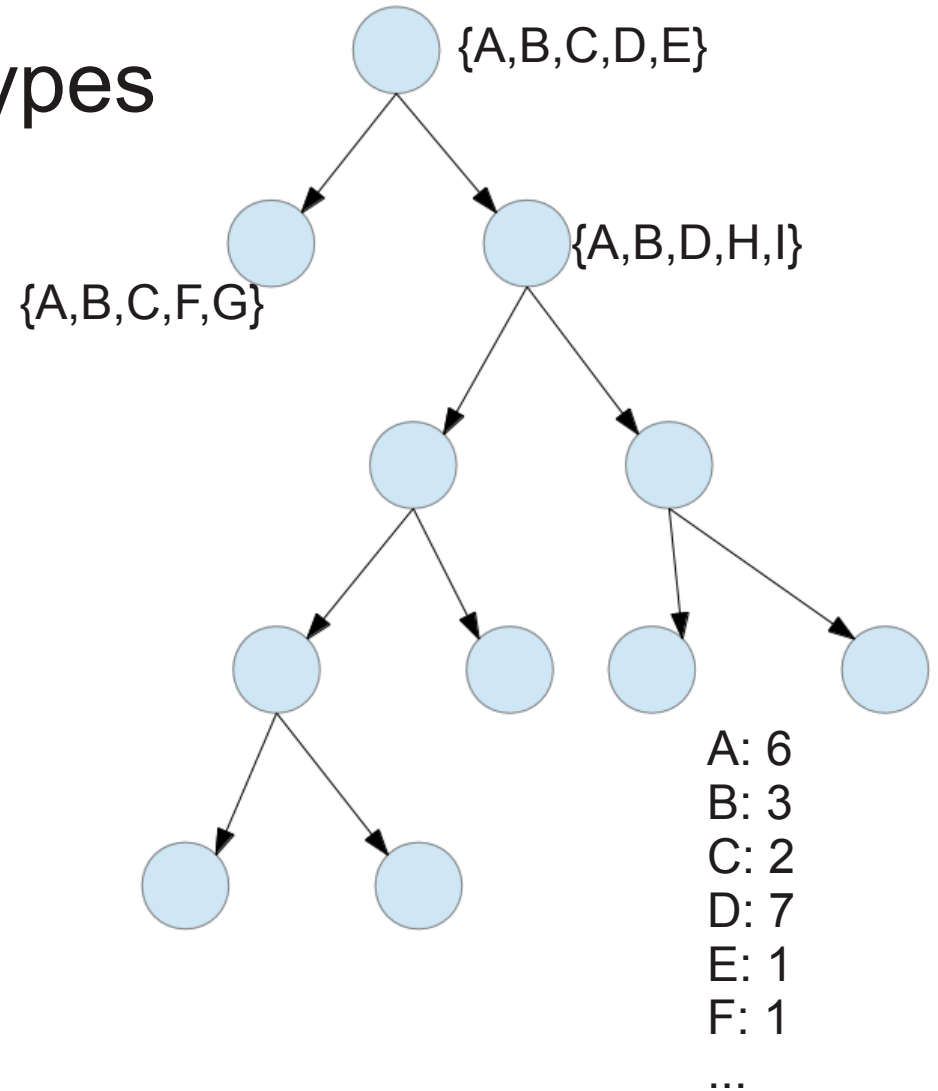


Why do we observe the power-law?

- Start with one cell-type with N DHSs
- Generate two new cell-types

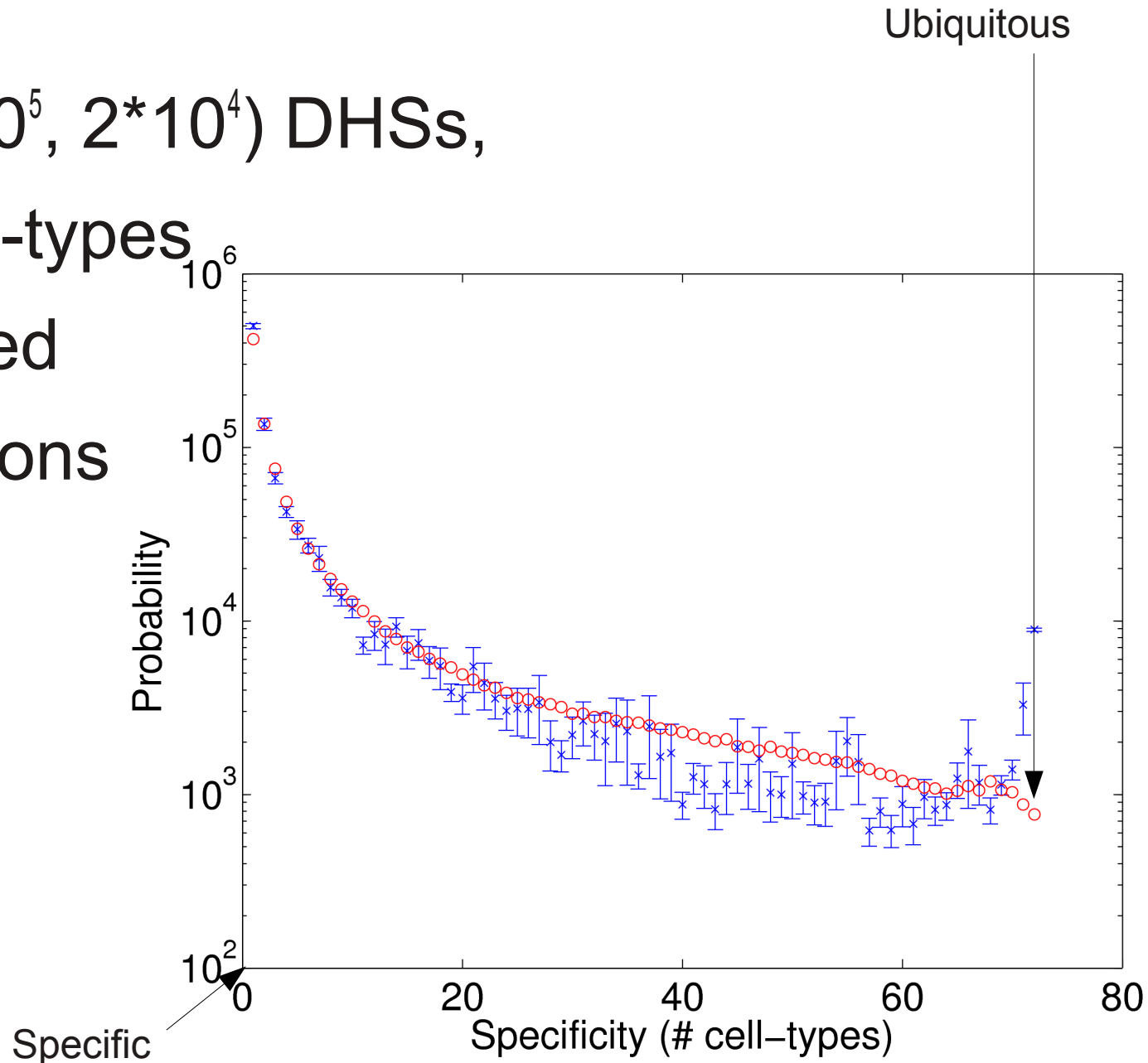
$$P(\text{keep DHS}) = \frac{k^b}{1 + k^b}$$

- k - # cell-types where DHS present
- b – selective advantage of existing DHSs
- Add new DHSs to get N
- Expand until C nodes
- Sample 72 nodes

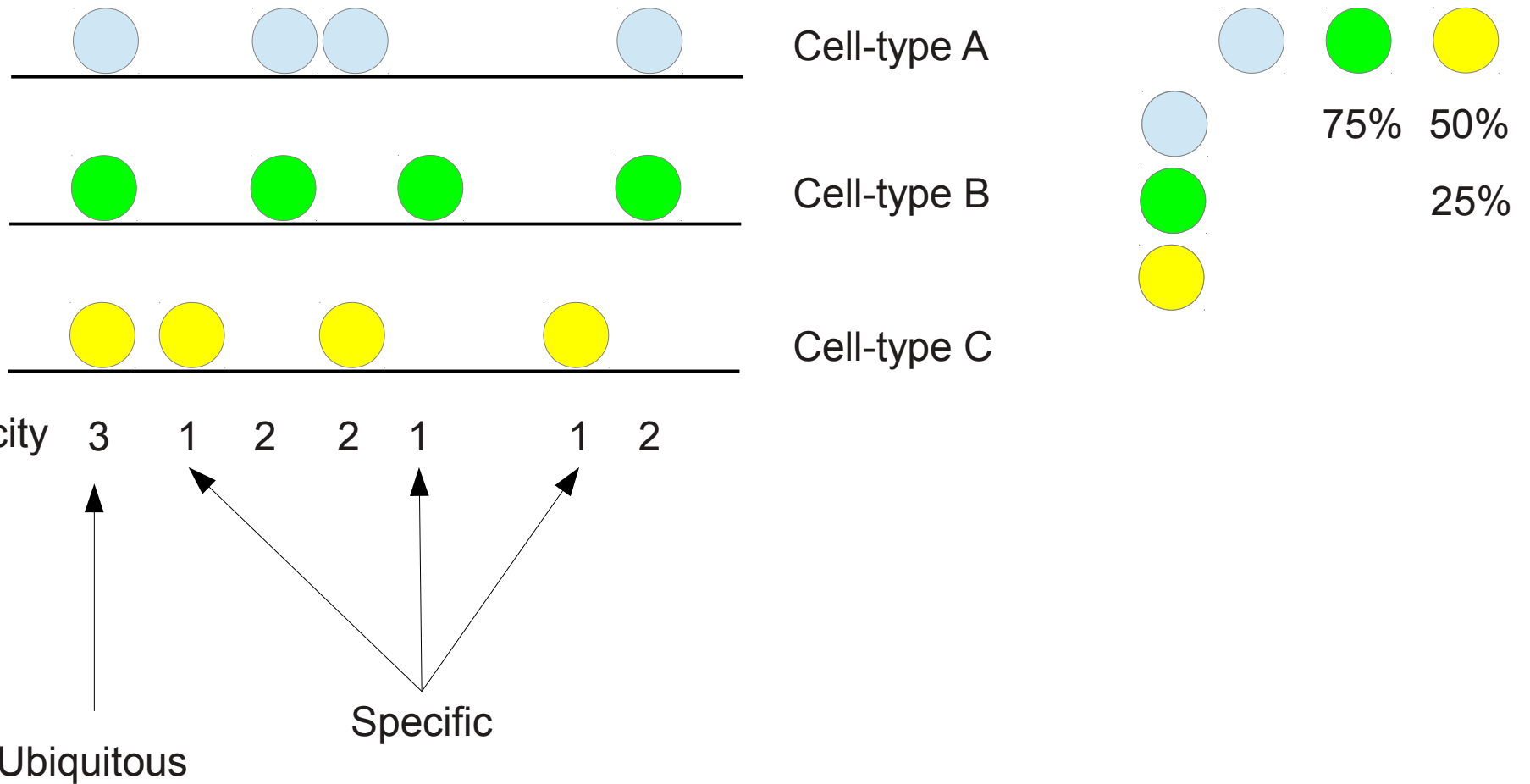


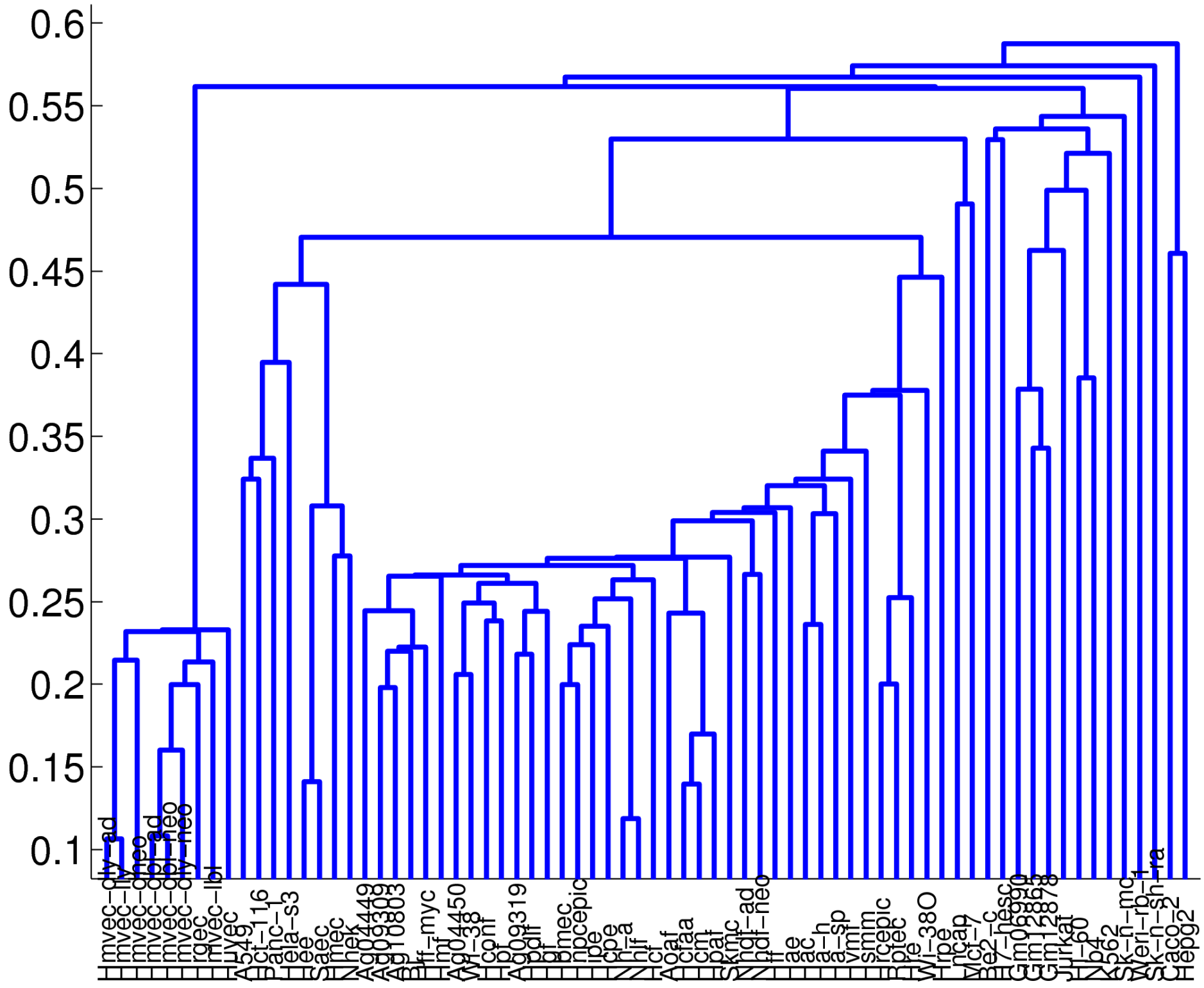
Model yields power-law $\sim k^{-b}$

- $N = \text{Normal}(10^5, 2 \cdot 10^4)$ DHSs,
- $C = 1,000$ cell-types
- $S = 72$ sampled
- $n = 10$ repetitions
 - $b = 1.35$
 - Slope = -1.4

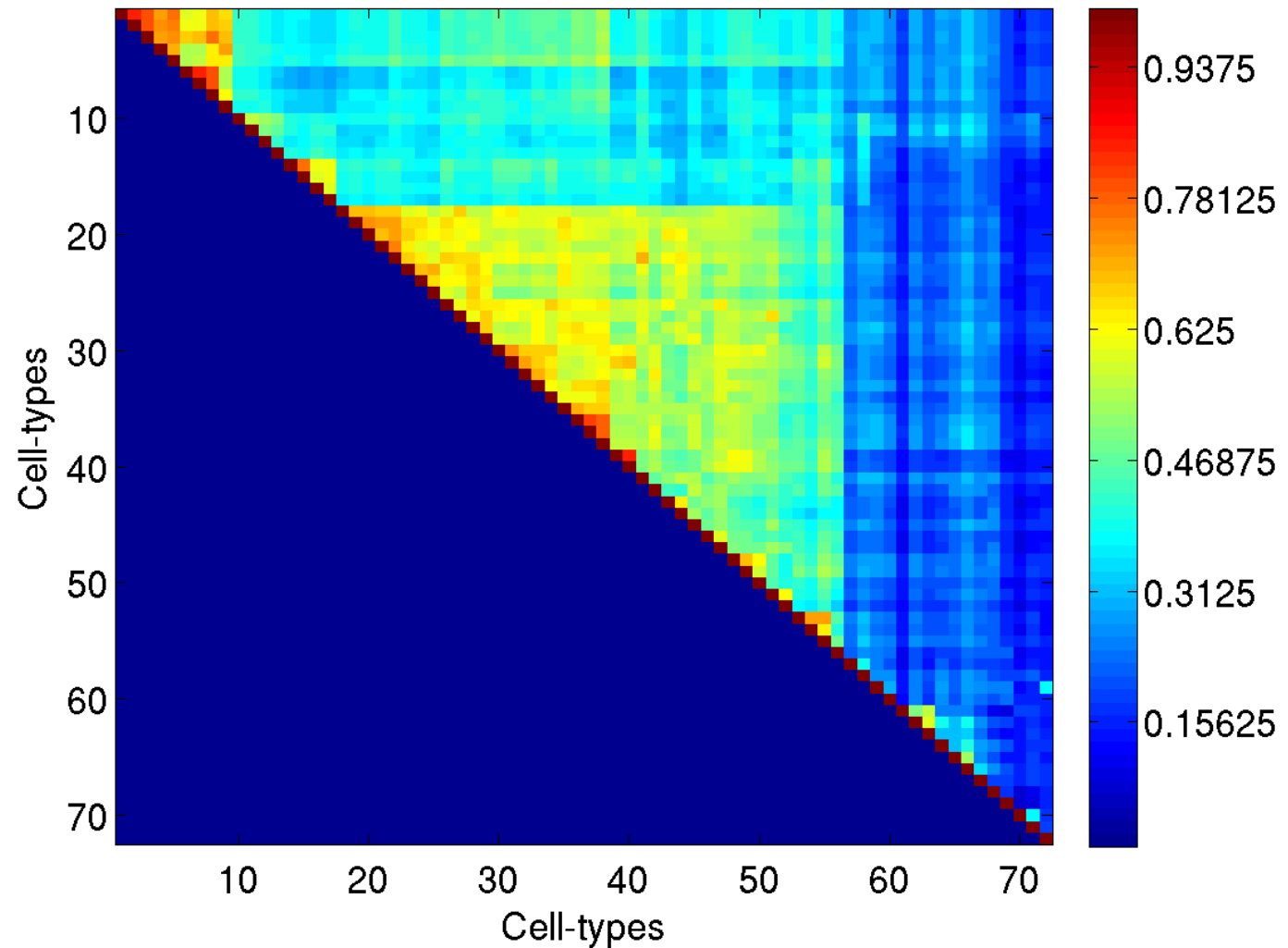


Cell-type distance can be inferred from DHS overlap

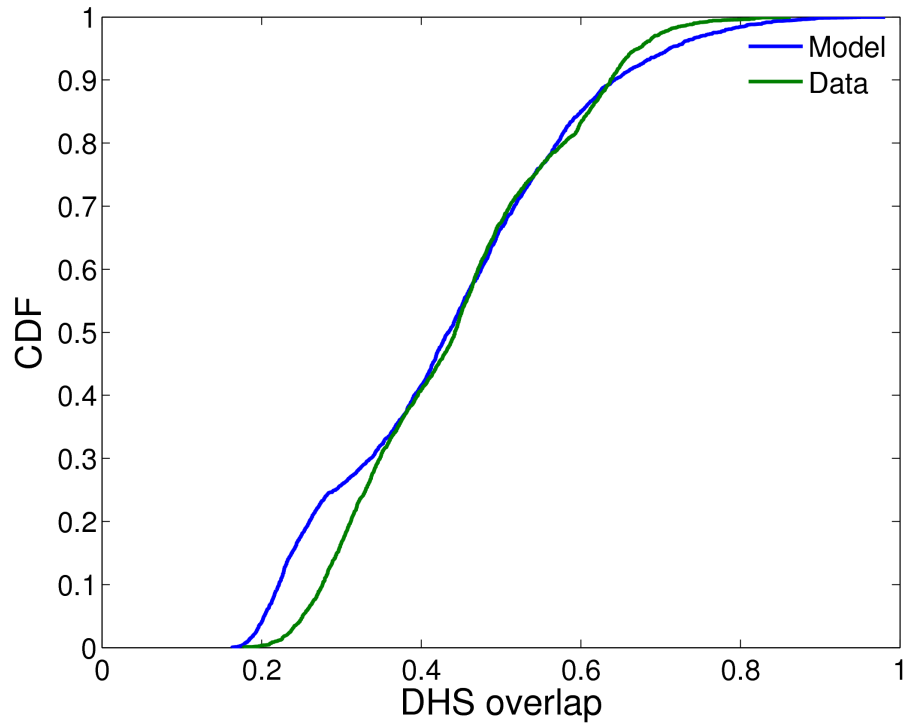




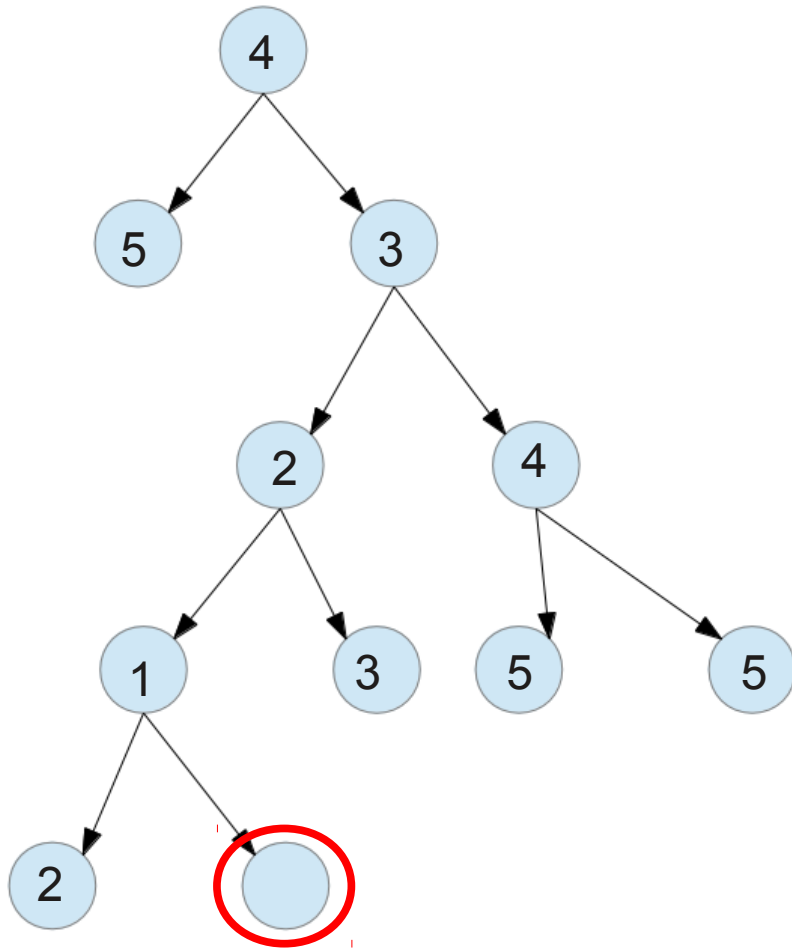
Distribution of pair-wise overlaps



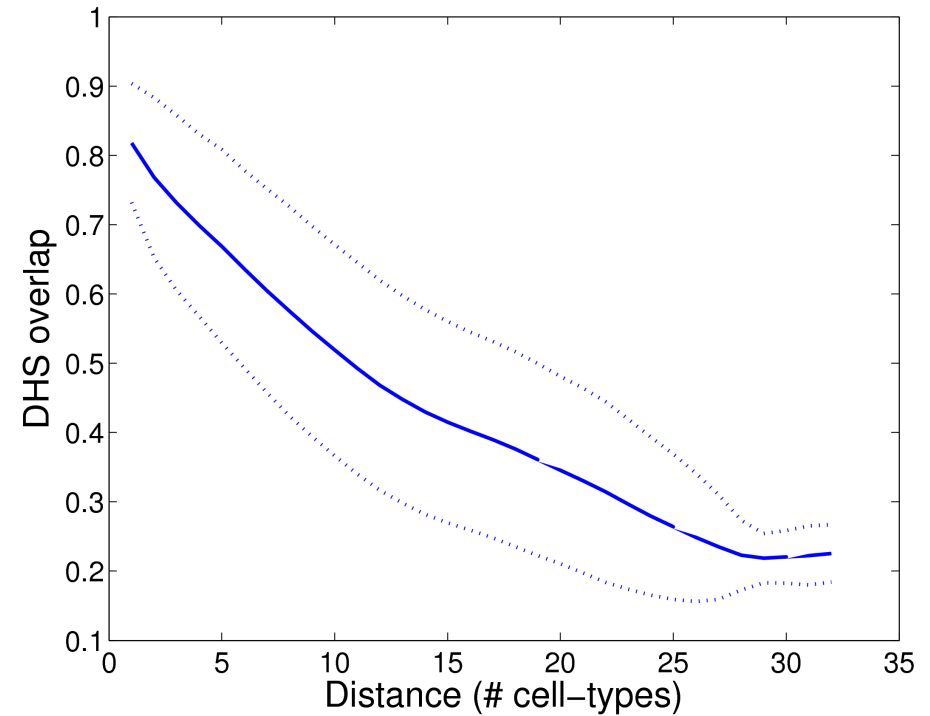
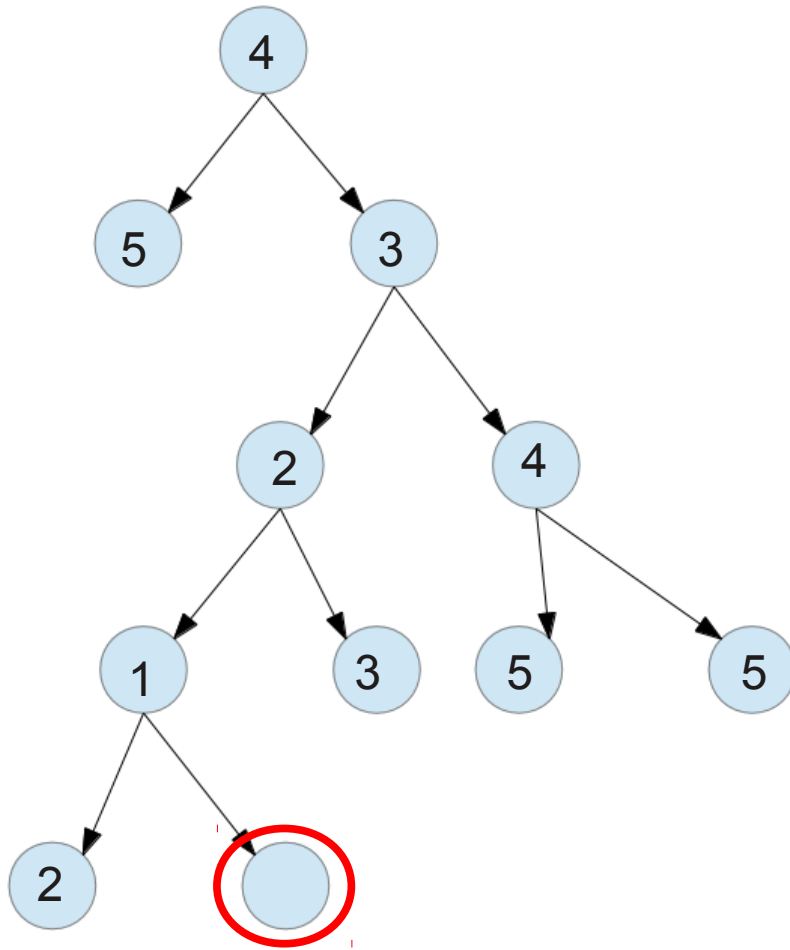
Model predicts distribution of pairwise overlaps



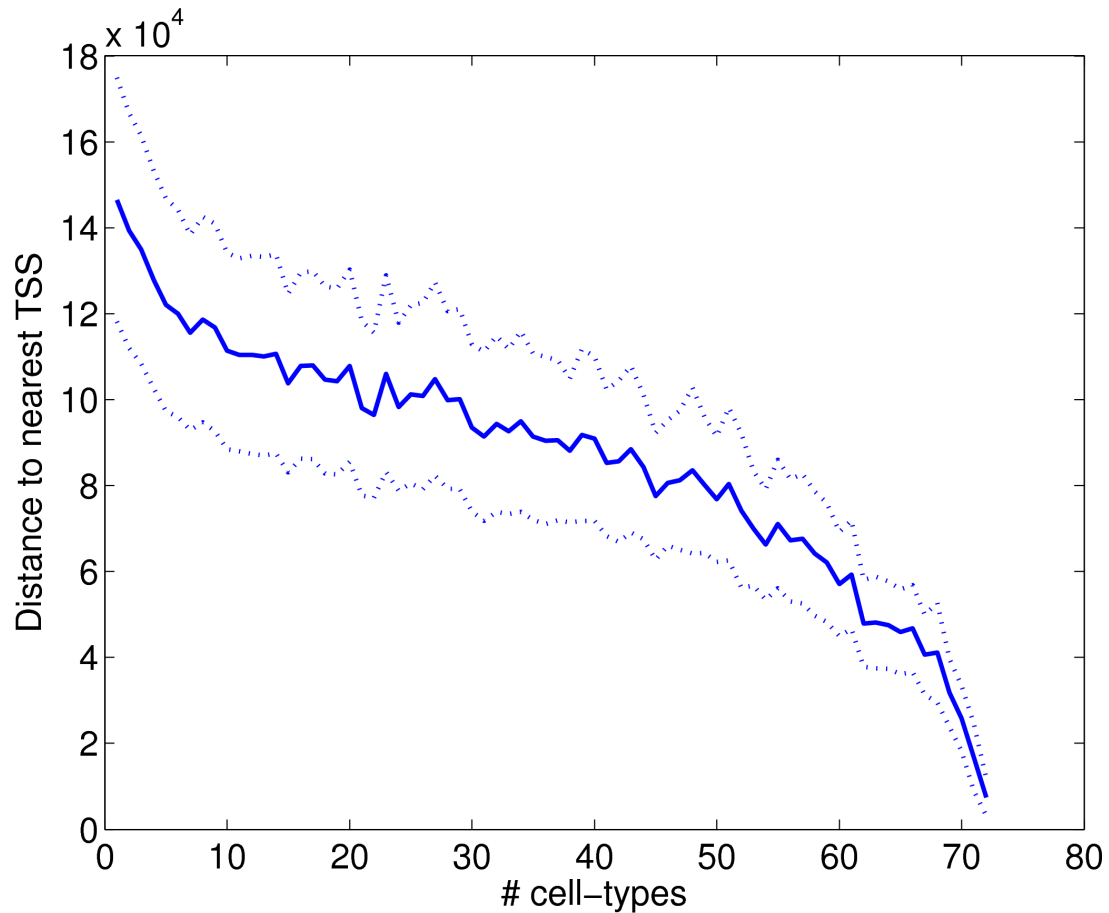
The overlap is inversely proportional to the distance



The overlap is inversely proportional to the distance



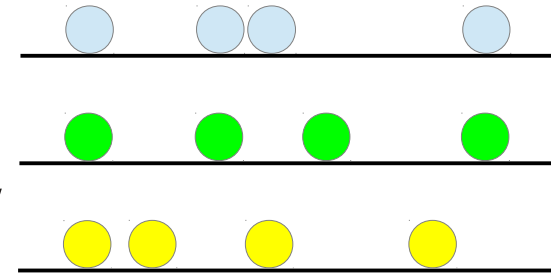
Specific DHSs are located further from known TSSs



WTF! Where are those TFs?



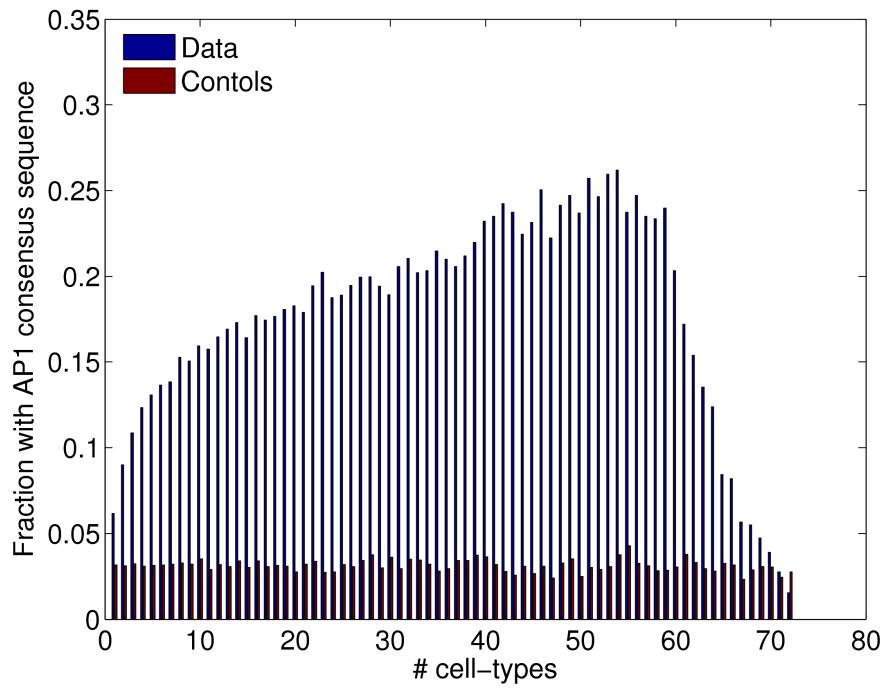
Summary

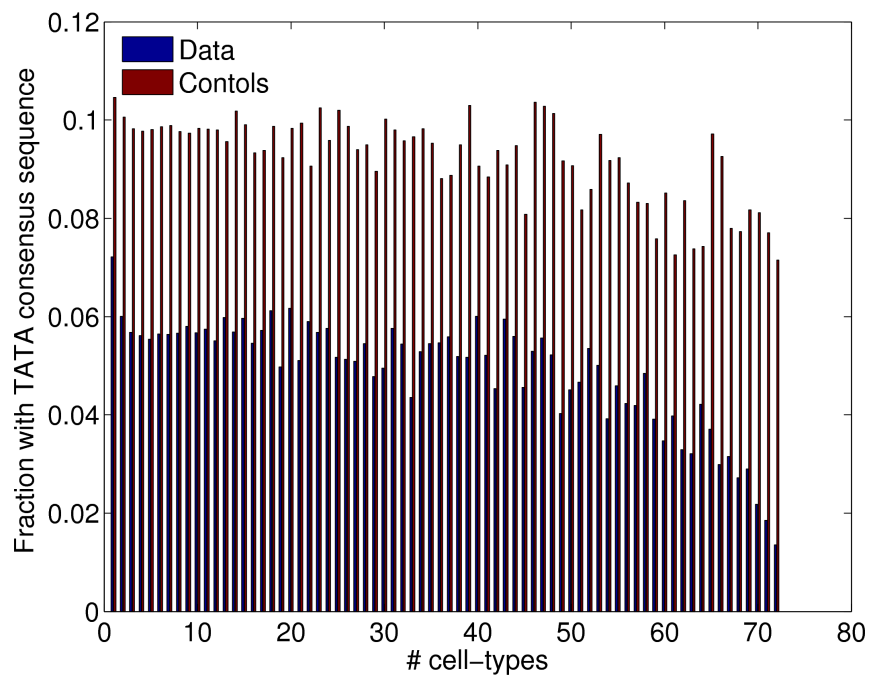
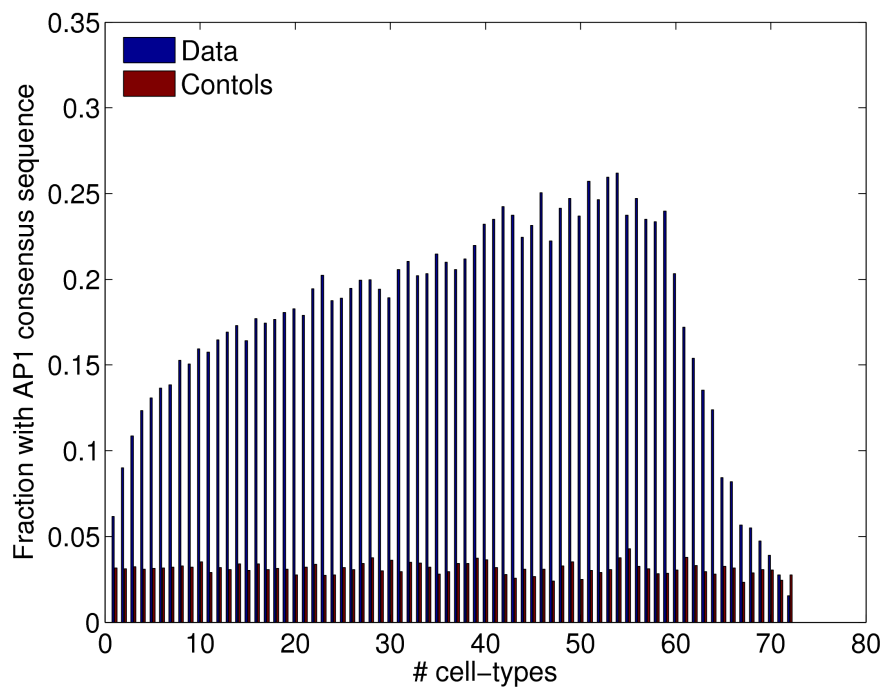


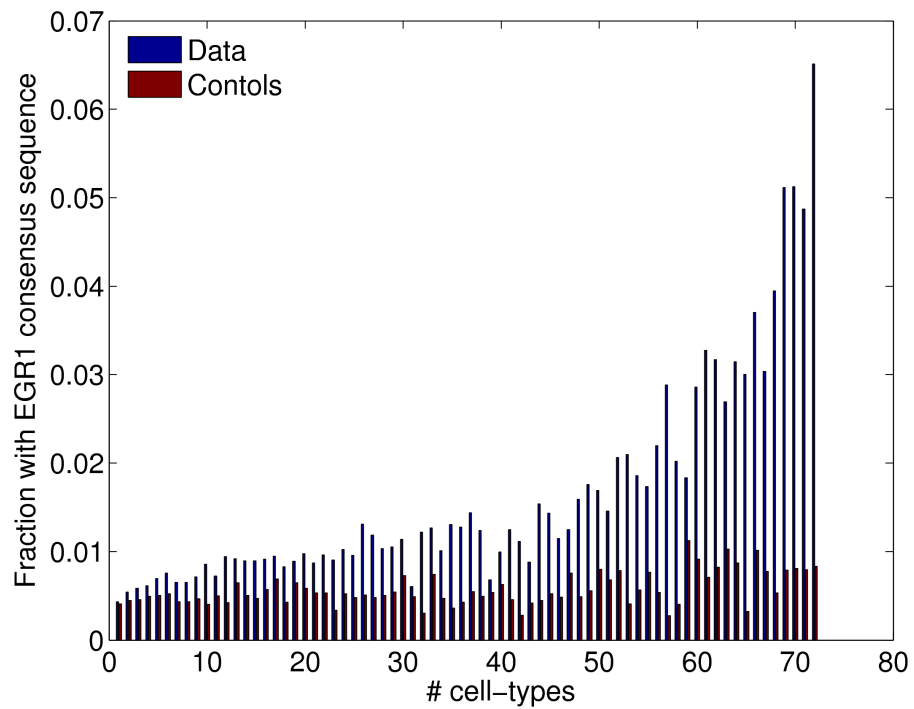
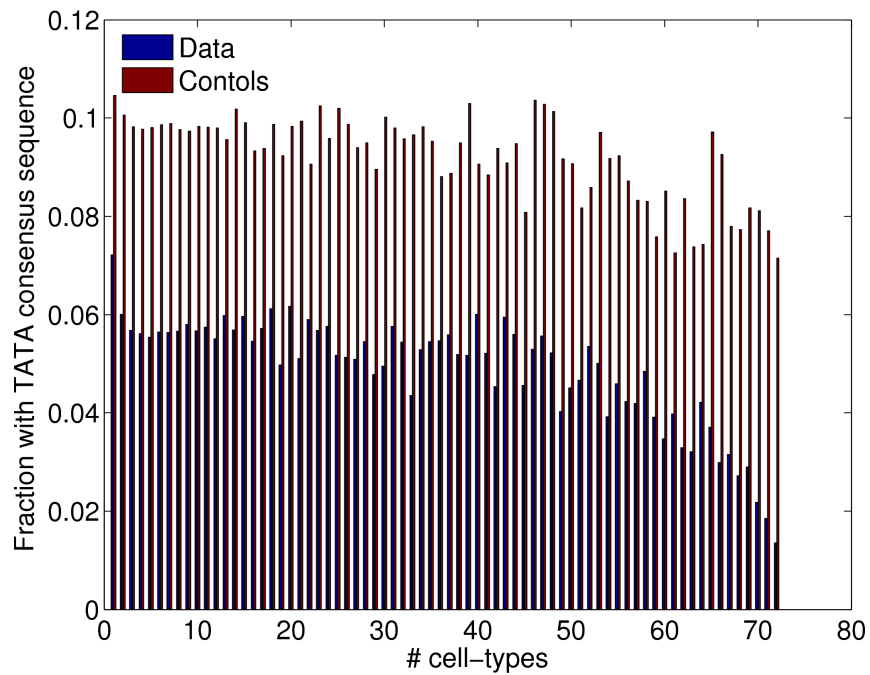
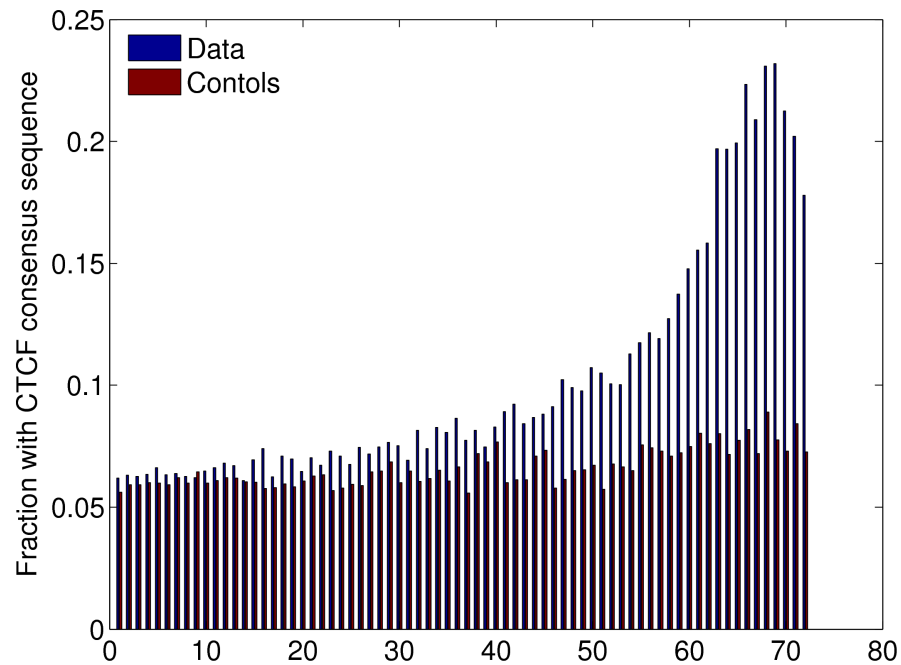
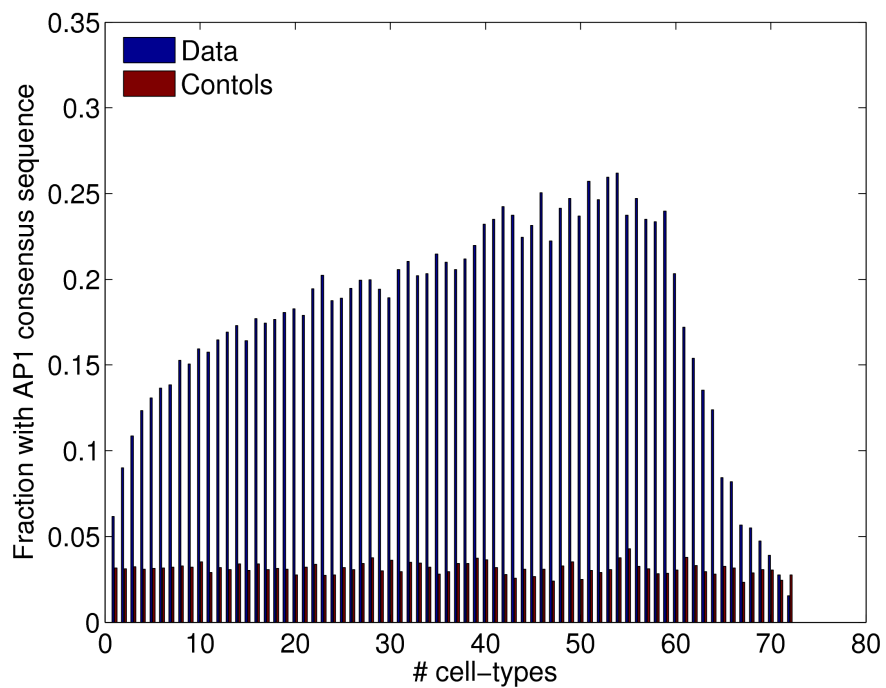
- DHS specificity follows a power-law
- Generative model
 - Parent DHS used as template
 - Selection rule: $P(\text{keep DHS}) = \frac{k^b}{1 + k^b}$
 - Selective advantage b
- Specificity related to function

Acknowledgements

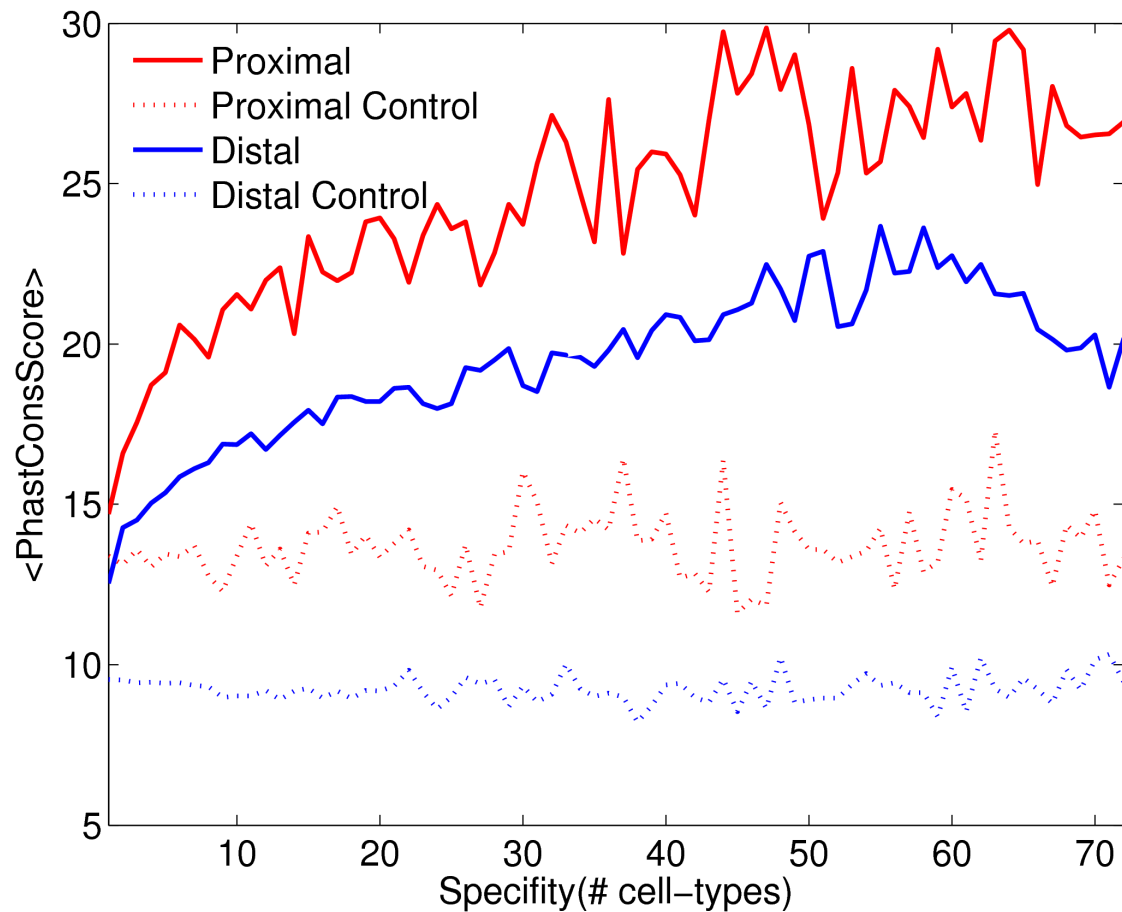
- Gabriel Kreiman
- Wui Ip
- Ben Tsuda
- Enrique Tobis



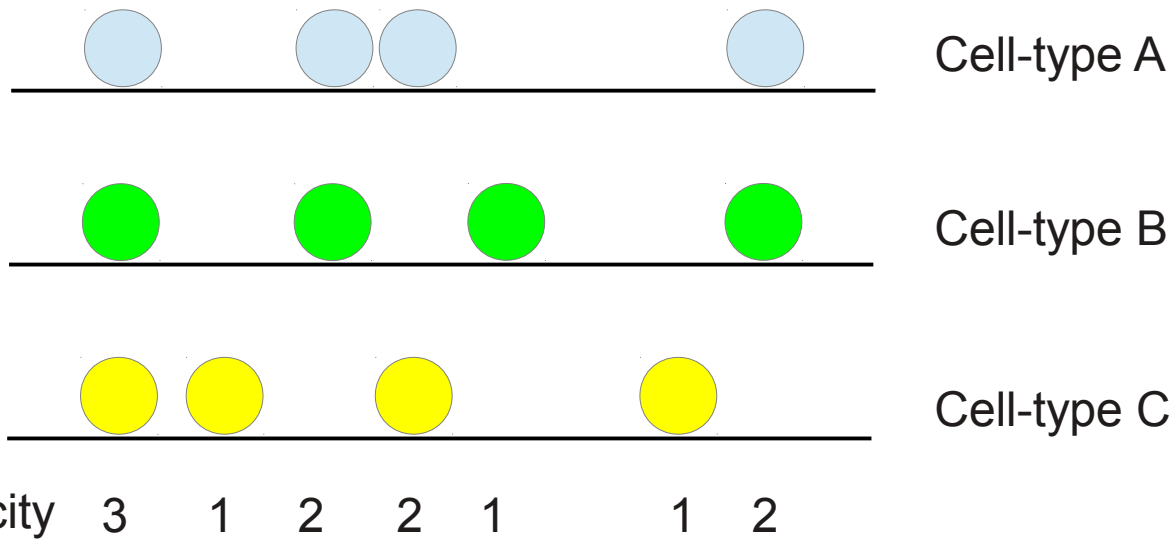




Specific DHSs are less conserved

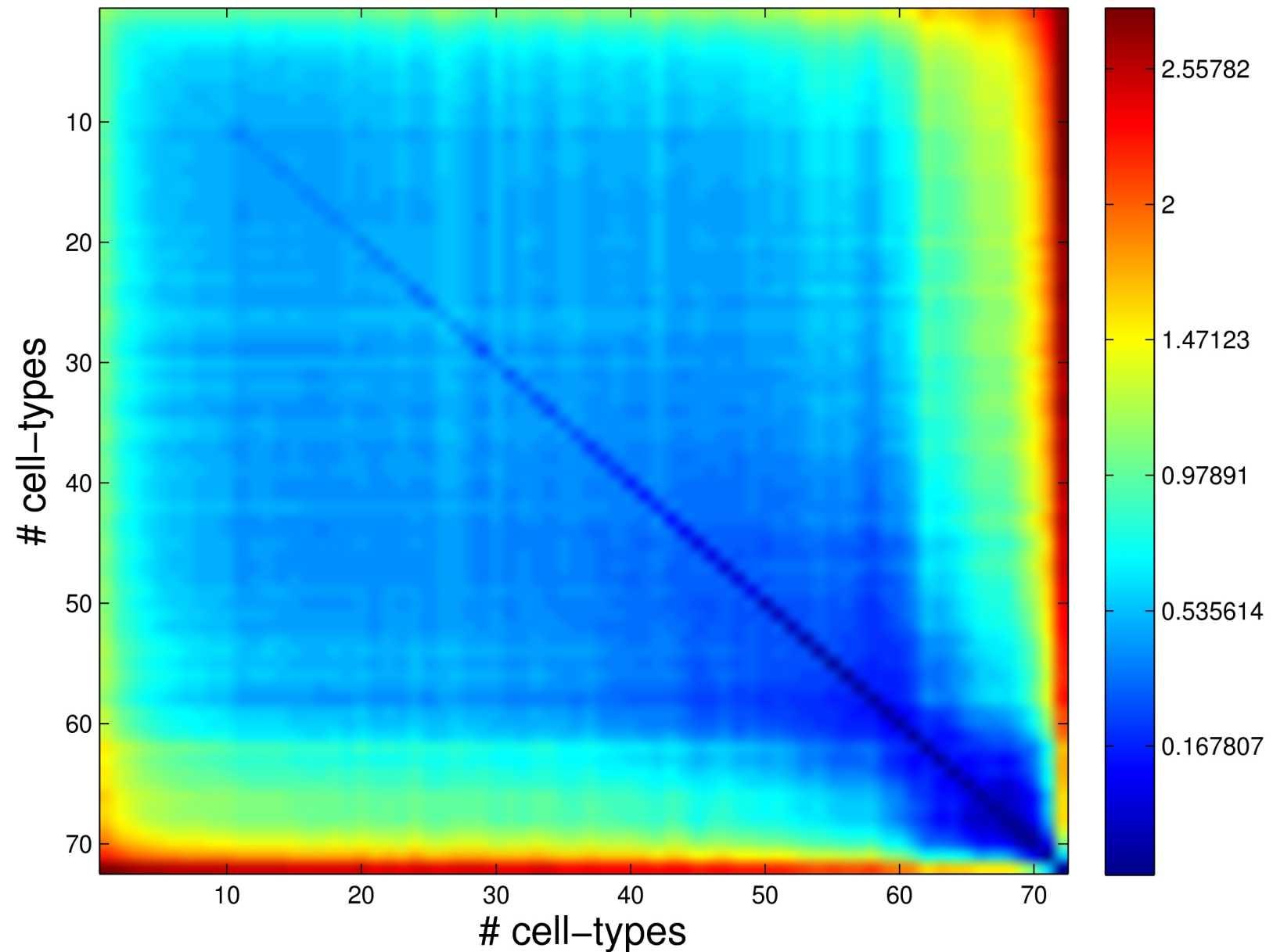


Sequence similarities of DHSs



Word	Specificity		
	1	2	3
AAAAAA	.001	.0002	.0005
AAAAAC	.002	.001	.0043
AAAAAG		
AAAAAT			
AAAACA			
.....			

Difference between sequence distributions

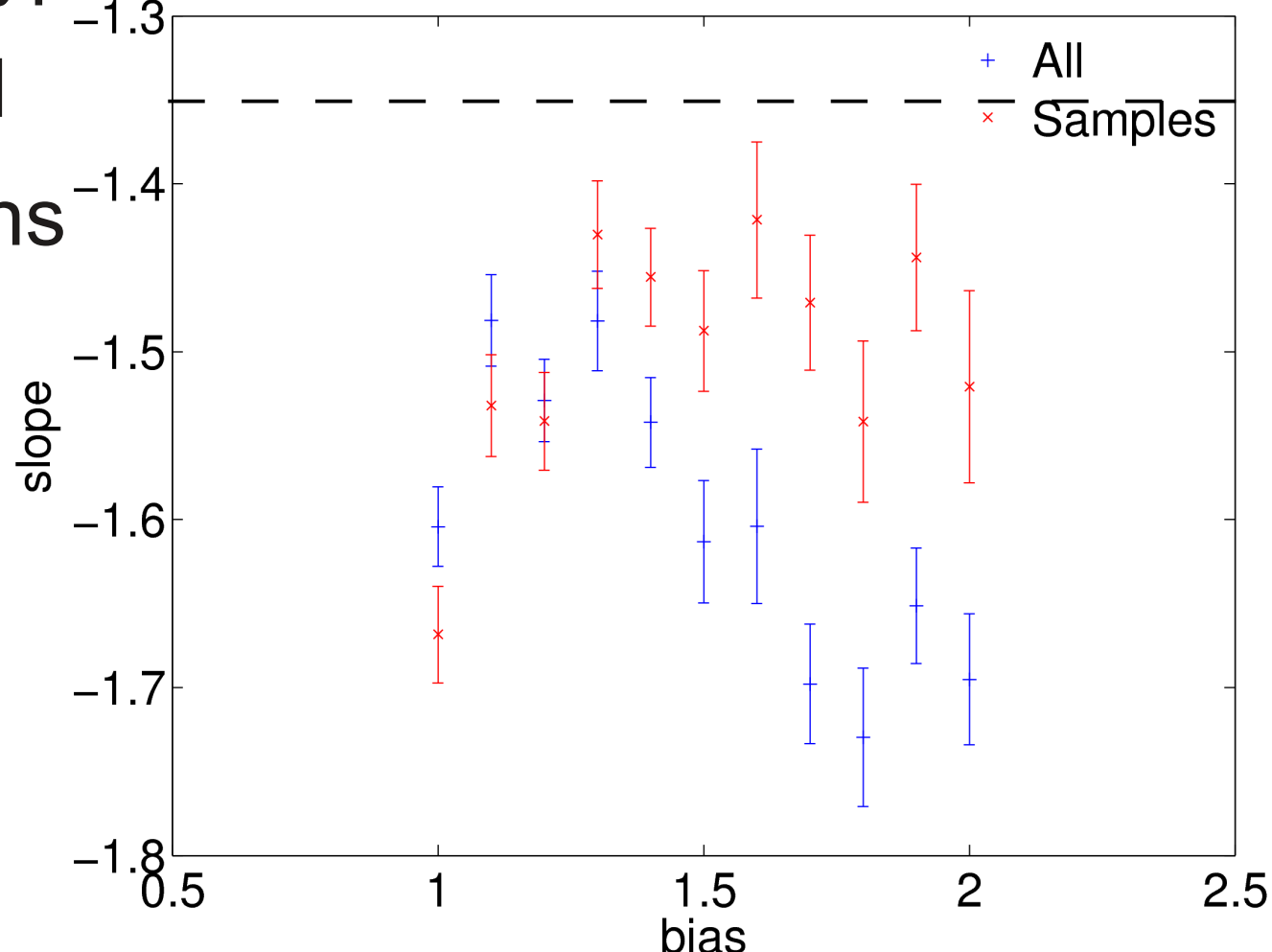


Functional role can be predicted by
histone modifications and TFs

Ubiquitous DHSs have higher GC

Monte Carlo Simulations

- $N = 100,000$ DHSs
- $C = 1,000$ cell-types
- $S = 72$ sampled
- $n = 10$ repetitions



Random branching process

$$P(\textit{Specificity} = k) = \frac{1}{n} \sum_{i=1}^n \underbrace{Q(k, i)}_{\substack{\text{Existing DHS} \\ \text{having} \\ \text{specificity } k}} + \underbrace{R(k, i)}_{\substack{\text{New DHS} \\ \text{having} \\ \text{specificity } k}}$$

Random branching process

$$P(\textit{Specificity} = k) \approx \frac{1}{n} \sum_{i=1}^n S(k) = S(k)$$

- Assume that tree is very large

Random branching process

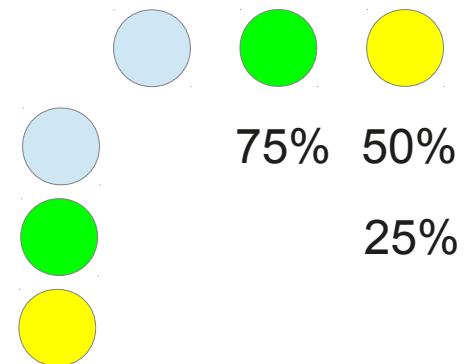
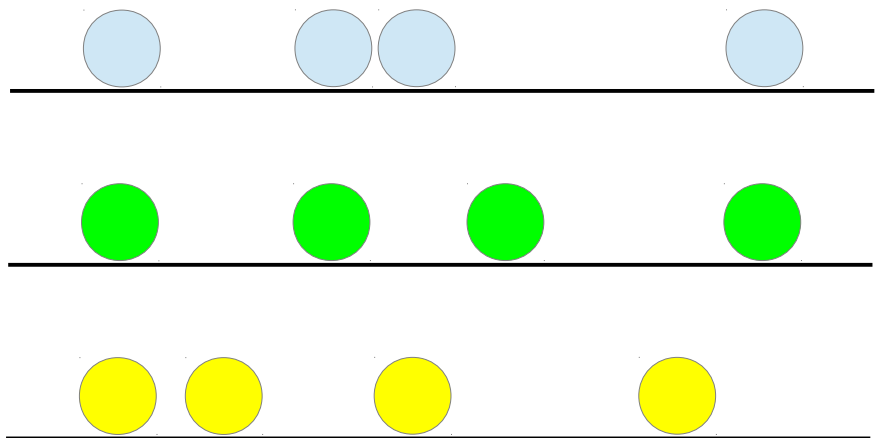
$$P(\textit{Specificity} = k) \approx \frac{1}{n} \sum_{i=1}^n S(k) = S(k)$$

$$S(k) = \frac{1}{1+k^b} \prod_{j=1}^k \frac{j^b}{1+j^b} = \frac{k^b(k-1)^b \dots 2^b}{(1+k^b)^2(1+(k-1)^b) \dots (1+2^b)2} \geq \frac{1}{1+k^b} \geq k^{-b}$$

Signal strength

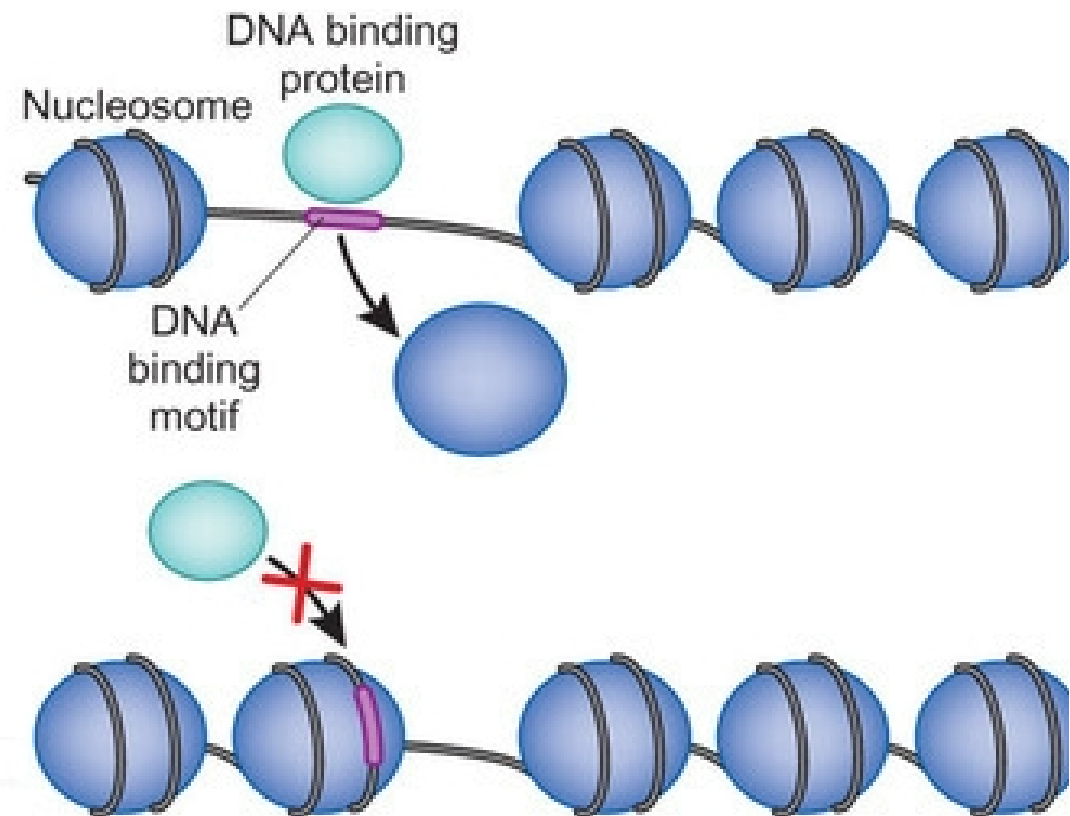
Many near-ubiquitous DHSs are
insulators

DHS overlap



What determines the location of DHSs?

- Identify distribution of DHSs
 - Functional characteristics
 - Generative model



Uniform distribution of specificities

- ~100,000 DHSs per cell-type
 - ~1,000 of each specificity

