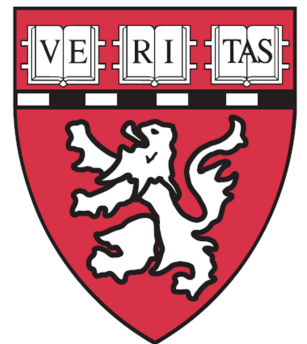


# Conservation of tissue-specific gene expression predicts transcription factor binding sites in human and mouse

Martin Hemberg

Children's Hospital Boston, Harvard Medical School  
Department of Ophthalmology  
Swartz Center for Theoretical Neuroscience

Kreiman Lab

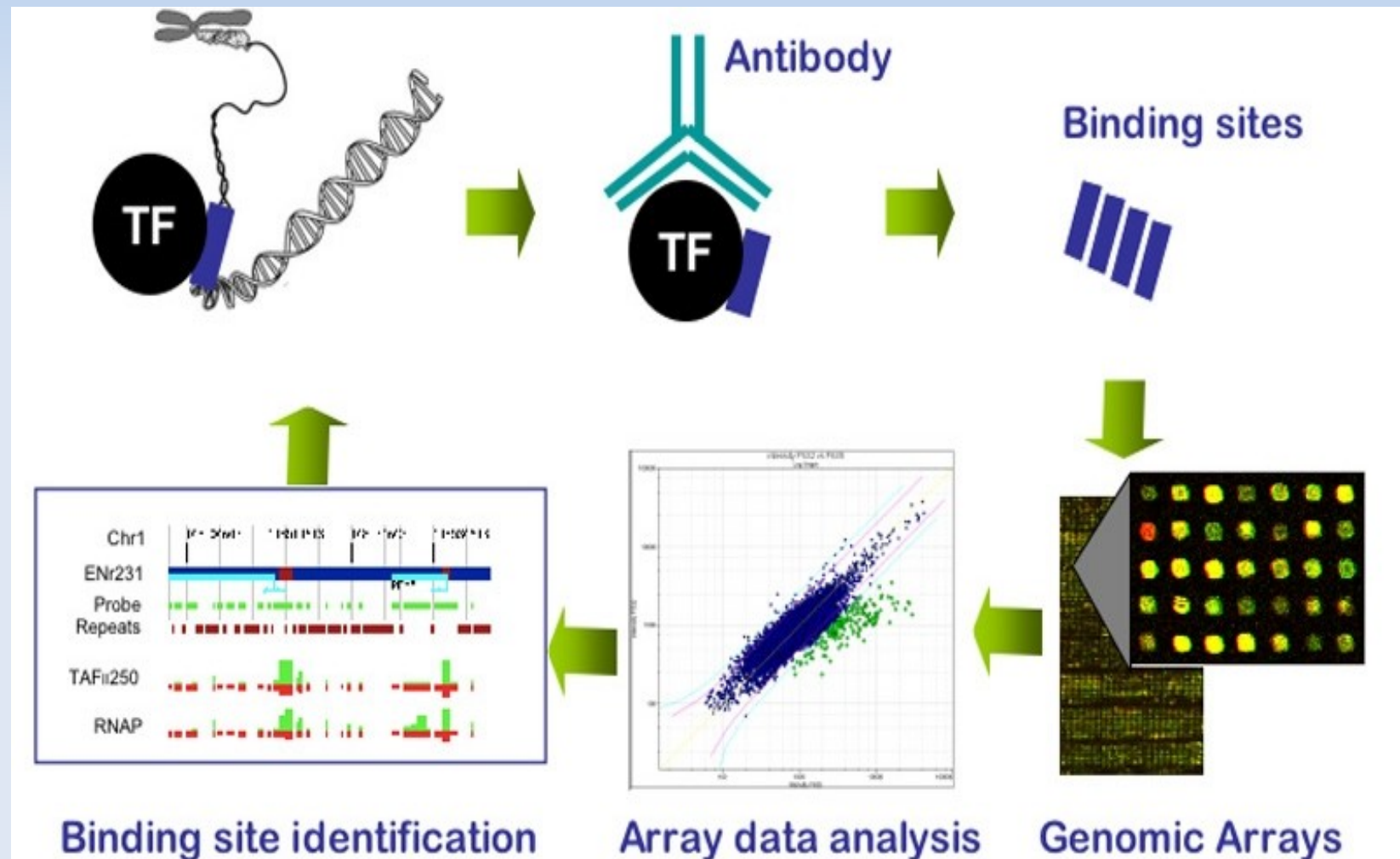


# A talk of two parts

- Using gene expression to predict conservation of transcription factor binding sites (TFBSs)
- Using conservation of TFBSs to predict gene expression

# Genome wide TFBSs

- ChIP-chip, ChIP-Seq experiments



# Genome wide TFBSs

- CHIP-chip, CHIP-Seq experiments
  - Widespread binding of TFs

The Oct4 and Nanog transcription network regulates pluripotency in mouse embryonic stem cells

Yuin-Han Loh<sup>1,2,7</sup>, Qiang Wu<sup>1,7</sup>, Joon-Lin Chew<sup>1,2,7</sup>, Vinsensius B Vega<sup>3</sup>, Weiwei Zhang<sup>1,2</sup>, Xi Chen<sup>1,2</sup>, Guillaume Bourque<sup>3</sup>, Joshy George<sup>3</sup>, Bernard Leong<sup>3</sup>, Jun Liu<sup>4</sup>, Kee-Yew Wong<sup>5</sup>, Ken W Sung<sup>3</sup>, Charlie W H Lee<sup>3</sup>, Xiao-Dong Zhao<sup>4</sup>, Kuo-Ping Chiu<sup>3</sup>, Leonard Lipovich<sup>3</sup>, Vladimir A Kuznetsov<sup>3</sup>, Paul Robson<sup>2,5</sup>, Lawrence W Stanton<sup>5</sup>, Chia-Lin Wei<sup>4</sup>, Yijun Ruan<sup>4</sup>, Bing Lim<sup>5,6</sup> & Huck-Hui Ng<sup>1,2</sup>

Cell, Vol. 122, 947–956, September 23, 2005, Copyright ©2005 by Elsevier Inc. DOI 10.1016/j.cell.2005.08.020

**Core Transcriptional Regulatory Circuitry  
in Human Embryonic Stem Cells**

# Genome wide TFBSs

- CHIP-chip, CHIP-Seq experiments
  - Widespread binding of TFs
    - Functional significance of TFBSs?

The Oct4 and Nanog transcription network regulates pluripotency in mouse embryonic stem cells

Yuin-Han Loh<sup>1,2,7</sup>, Qiang Wu<sup>1,7</sup>, Joon-Lin Chew<sup>1,2,7</sup>, Vinsensius B Vega<sup>3</sup>, Weiwei Zhang<sup>1,2</sup>, Xi Chen<sup>1,2</sup>, Guillaume Bourque<sup>3</sup>, Joshy George<sup>3</sup>, Bernard Leong<sup>3</sup>, Jun Liu<sup>4</sup>, Kee-Yew Wong<sup>5</sup>, Ken W Sung<sup>3</sup>, Charlie W H Lee<sup>3</sup>, Xiao-Dong Zhao<sup>4</sup>, Kuo-Ping Chiu<sup>3</sup>, Leonard Lipovich<sup>3</sup>, Vladimir A Kuznetsov<sup>3</sup>, Paul Robson<sup>2,5</sup>, Lawrence W Stanton<sup>5</sup>, Chia-Lin Wei<sup>4</sup>, Yijun Ruan<sup>4</sup>, Bing Lim<sup>5,6</sup> & Huck-Hui Ng<sup>1,2</sup>

Extensive low-affinity transcriptional interactions in the yeast genome

Cell, Vol. 122, 947–956, September 23, 2005, Copyright ©2005 by Elsevier Inc. DOI 10.1016/j.cell.2005.08.020

Amos Tanay

**Core Transcriptional Regulatory Circuitry in Human Embryonic Stem Cells**

# Genome wide TFBSs

- CHIP-chip, CHIP-Seq experiments
  - Widespread binding of TFs
    - Functional significance of TFBSs?
  - Low degree of conservation between species
    - Regulation and function of a gene decoupled?

## Divergence of Transcription Factor Binding Sites Across Related Yeast Species

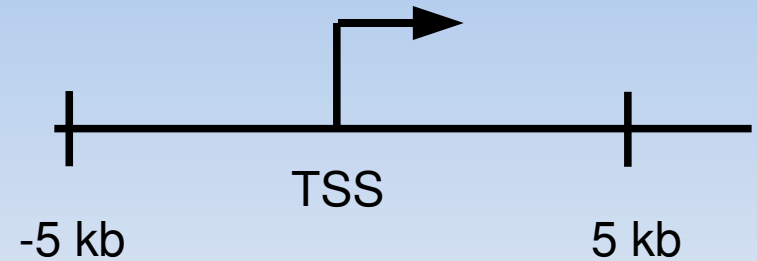
Anthony R. Borneman,<sup>1\*</sup> Tara A. Gianoulis,<sup>2</sup> Zhengdong D. Zhang,<sup>3</sup> Haiyuan Yu,<sup>3</sup> Joel Rozowsky,<sup>3</sup> Michael R. Seringhaus,<sup>3</sup> Lu Yong Wang,<sup>4</sup> Mark Gerstein,<sup>2,3,5</sup> Michael Snyder<sup>1,2,3†</sup>

## The Evolution of Combinatorial Gene Regulation in Fungi

Brian B. Tuch<sup>1,2©</sup>, David J. Galgoczy<sup>1,2©</sup>, Aaron D. Hernday<sup>1,2</sup>, Hao Li<sup>1\*</sup>, Alexander D. Johnson<sup>1,2\*</sup>

# Conservation of TFBSs in liver

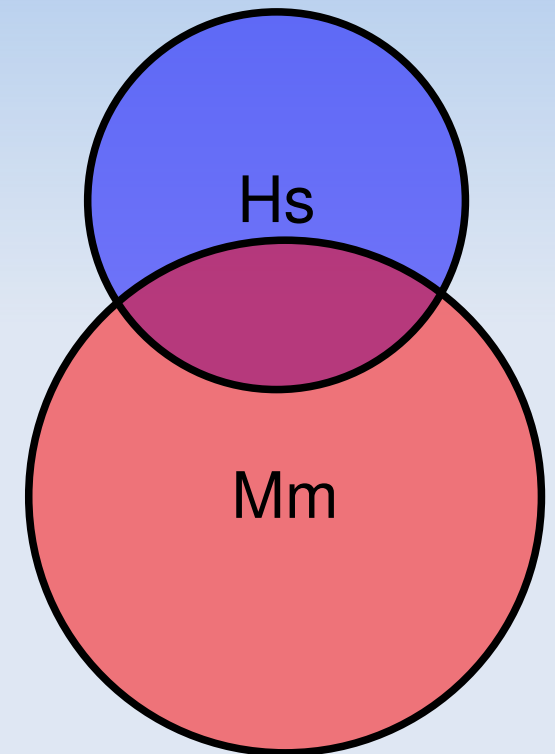
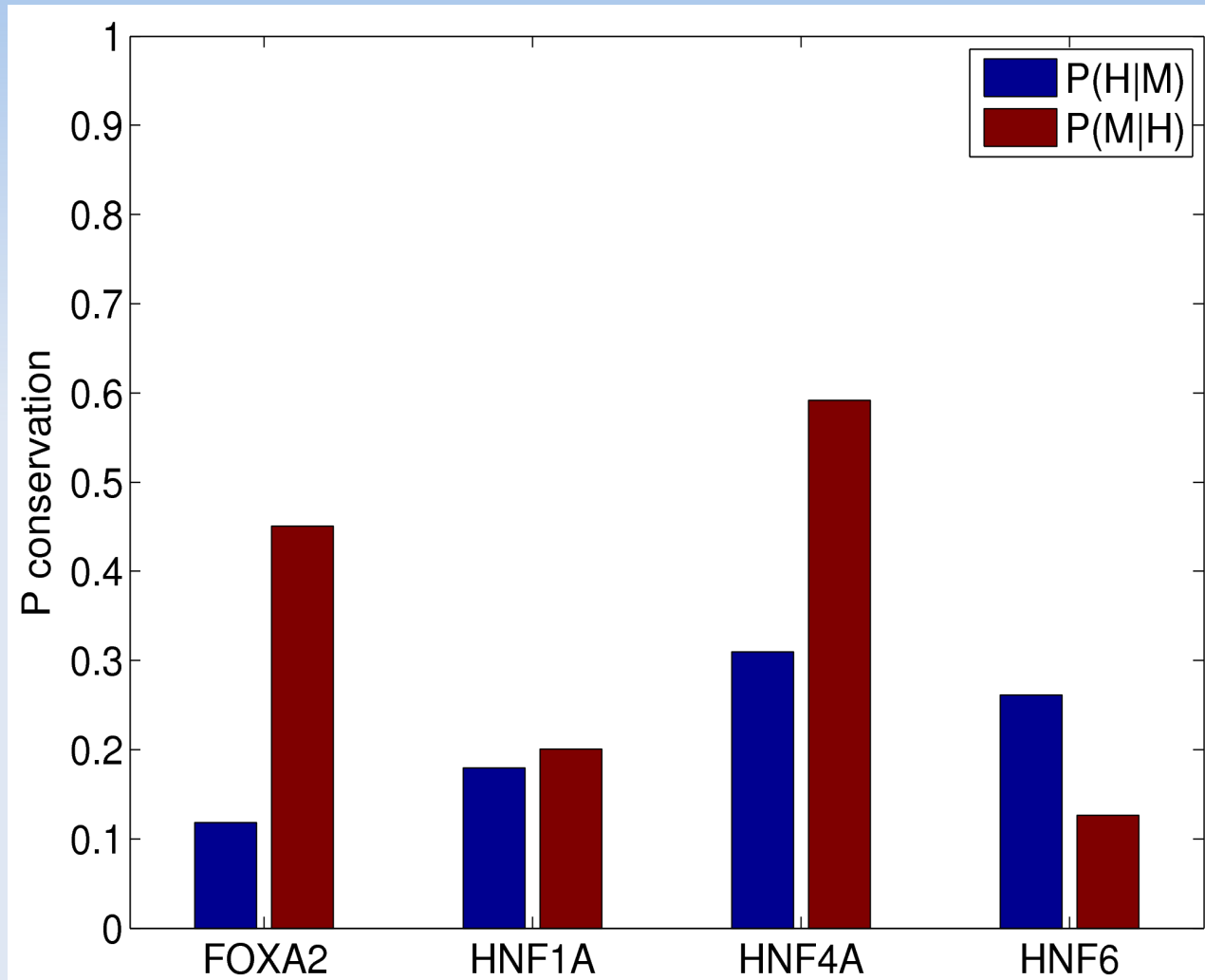
- Liver
  - Expression profiles similar
  - Homogeneous cell types
- Compare TF binding sites
  - ChIP-chip
  - **4022** homologous genes



Tissue-specific transcriptional regulation has diverged significantly between human and mouse

Duncan T Odom<sup>1,5,6</sup>, Robin D Dowell<sup>2,6</sup>, Elizabeth S Jacobsen<sup>1</sup>, William Gordon<sup>3</sup>, Timothy W Danford<sup>2</sup>, Kenzie D MacIsaac<sup>4</sup>, P Alexander Rolfe<sup>2</sup>, Caitlin M Conboy<sup>1,5</sup>, David K Gifford<sup>1,2</sup> & Ernest Fraenkel<sup>2,3</sup>

# Poor Conservation of TFBSs



**$P(H|M)$**  = probability of observing TFBS in Hs if TFBS in Mm  
= #conserved TFBS / #Mm TFBS



# Higher Conservation Expected for Genes Expressed in Both Species

- Functional aspects not considered
- Conservation of TFBS for genes not expressed in tissue?
  - Binding **not** expected for genes **not** expressed
  - These TFs are known to be activators in liver

# Higher Conservation Expected for Genes Expressed in Both Species

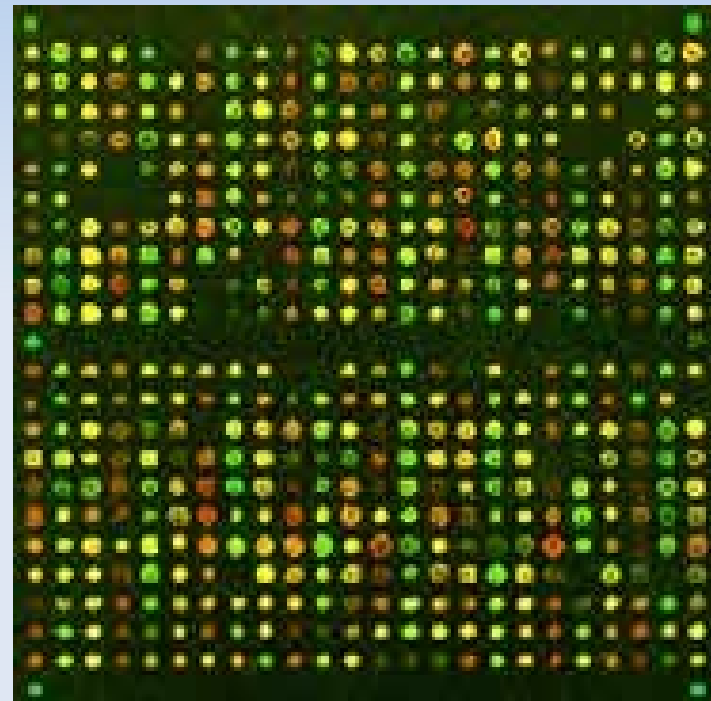
- Functional aspects not considered
- Conservation of TFBS for genes not expressed in tissue?
  - Binding **not** expected for genes **not** expressed
  - These TFs are known to be activators in liver

**Hypothesis:** Greater conservation for subset of genes which are expressed in both species

$$P(T_{Hs} | T_{Mm}, E_{Hs}, E_{Mm}) > P(T_{Hs} | T_{Mm})$$

# Microarray Expression Data

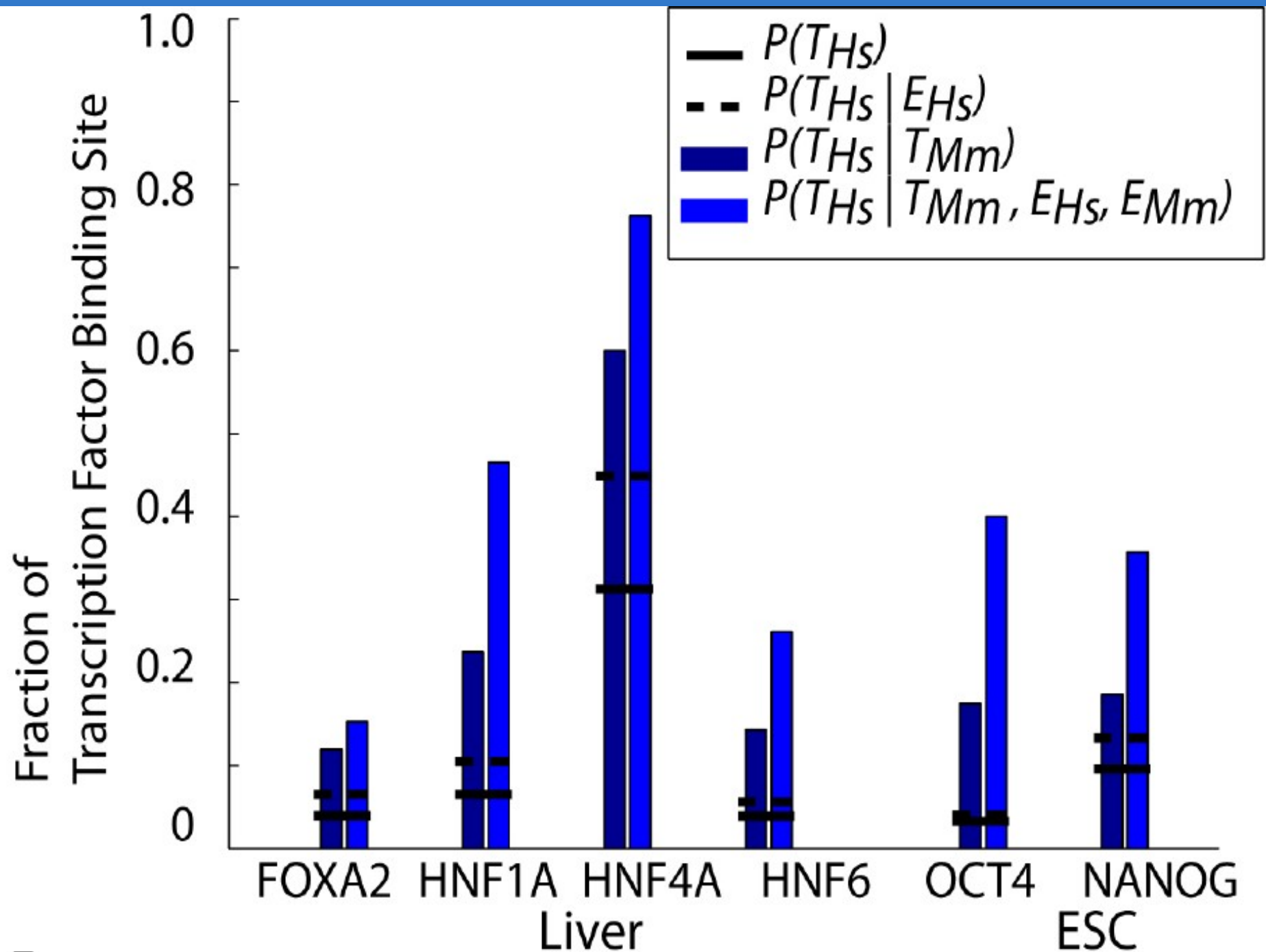
- GNF SymAtlas for liver
  - Microarray data from human and mouse
    - **3051/4022** pairs with expression for both

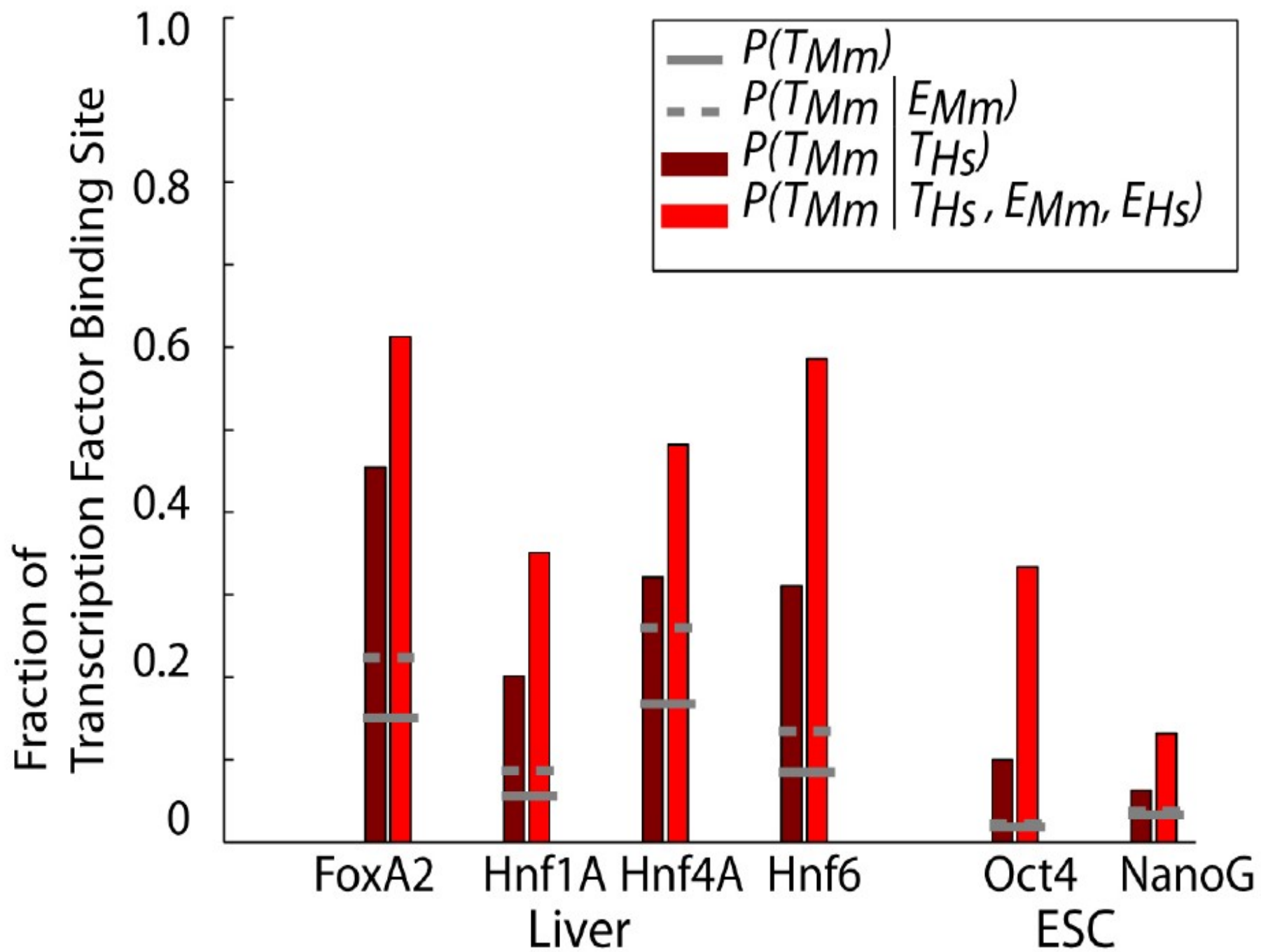


A gene atlas of the mouse and human protein-encoding transcriptomes

Andrew I. Su<sup>\*†</sup>, Tim Wiltshire<sup>\*†</sup>, Serge Batalov<sup>\*†</sup>, Hilmar Lapp<sup>\*</sup>, Keith A. Ching<sup>\*</sup>, David Block<sup>\*</sup>, Jie Zhang<sup>\*</sup>, Richard Soden<sup>\*</sup>, Mimi Hayakawa<sup>\*</sup>, Gabriel Kreiman<sup>\*‡</sup>, Michael P. Cooke<sup>\*</sup>, John R. Walker<sup>\*</sup>, and John B. Hogenesch<sup>\*§¶</sup>







# Predicting Expression from TFBS

- Use linear regression to predict expression
  - Real value for expression  $E$  of gene  $i$
  - Binary variable  $I$  for binding of TF  $j$  at gene  $i$

$$\log E_i = \sum_{j=1}^4 a_j I_{ij}$$

(4 parameters)

$$\log E_i = \sum_{j=1}^4 \sum_{c=0}^1 a_j^c I_{ij}^c$$

(8 parameters)

# Predicting Expression from TFBS

- Use linear regression to predict expression
  - Real value for expression  $E$  of gene  $i$
  - Binary variable  $I$  for binding of TF  $j$  at gene  $i$

$$\log E_i = \sum_{j=1}^4 a_j I_{ij}$$

(4 parameters)

$$\log E_i = \sum_{j=1}^4 \sum_{c=0}^1 a_j^c I_{ij}^c$$

(8 parameters)

$$\log E_i = \sum_{j=1}^4 \left[ a_j I_{ij} + \sum_{k=j+1}^4 a_{jk} I_{ijk} \right]$$

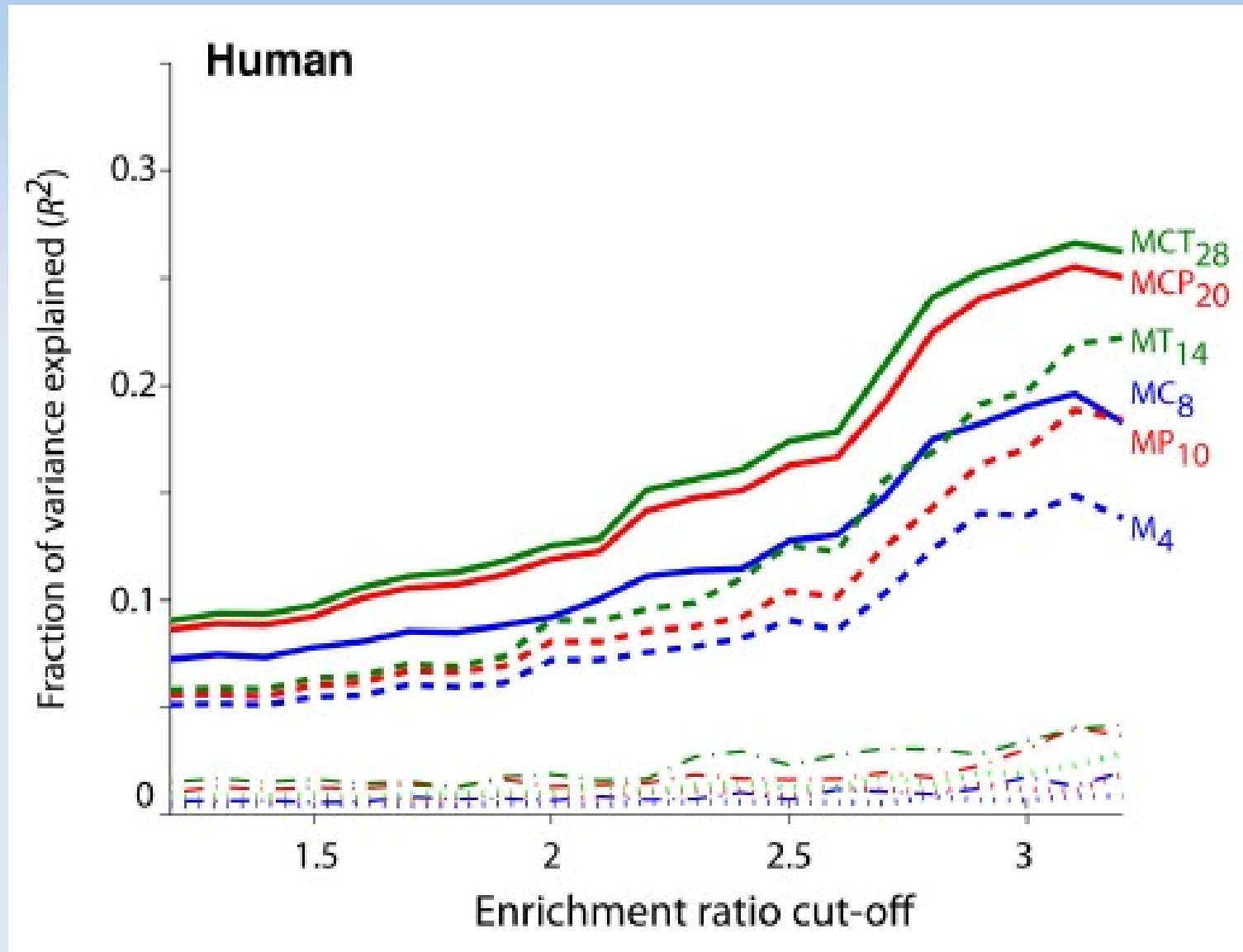
(10 parameters)

$$\log E_i = \sum_{j=1}^4 \sum_{c=0}^1 \left[ a_j^c I_{ij}^c + \sum_{k=j+1}^4 a_{jk}^c I_{ijk}^c \right]$$

(20 parameters)



# Predicting Expression from TFBSs



# Model complexity

- Adding parameters improves prediction
- Use AIC to penalize model complexity

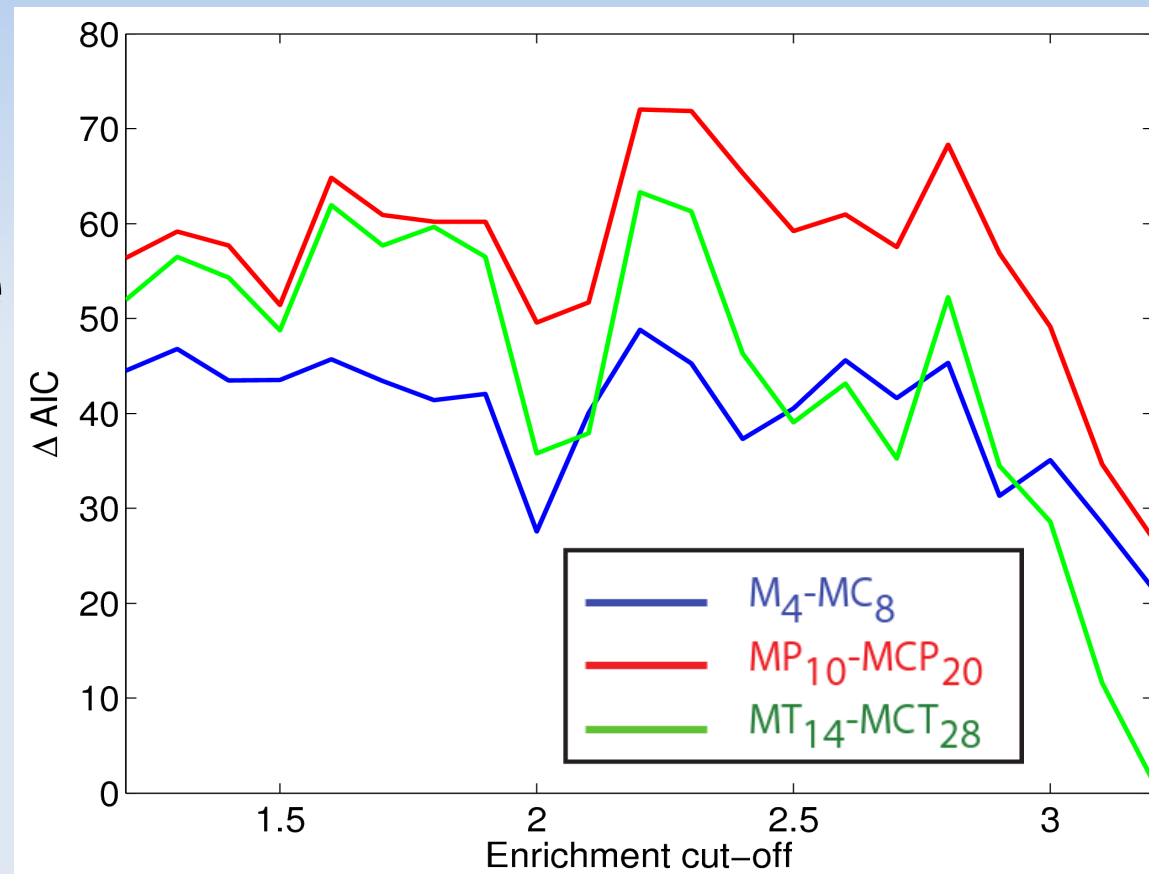
$$AIC = 2k + n \log(RSS/n)$$

$$L \propto \exp(-\Delta AIC/2)$$

$k$  - #parameters

$n$  - #genes

- Coeffs for conserved TFBSs 8 times higher



# Summary

- Combined several studies containing TF binding and gene expression data
- Genes expressed in both human and mouse twice as likely to have conserved TFBSs
  - Expression proxy for function
- Conserved TFBSs have a significantly greater impact on expression
  - Prediction of expression levels increased by ~30%