

Probabilistic Models of Structured Data

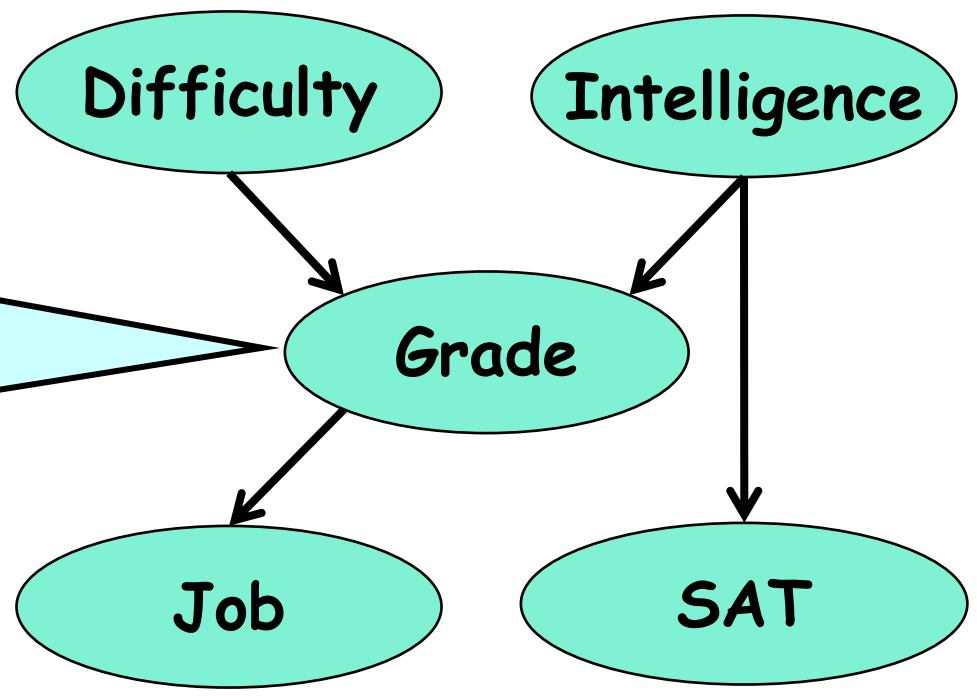
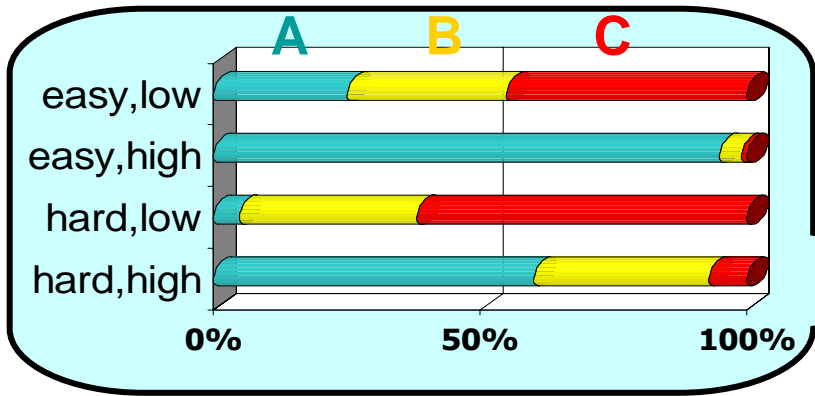
Daphne Koller

Stanford University



Bayesian Networks

CPD $P(G|D,I)$

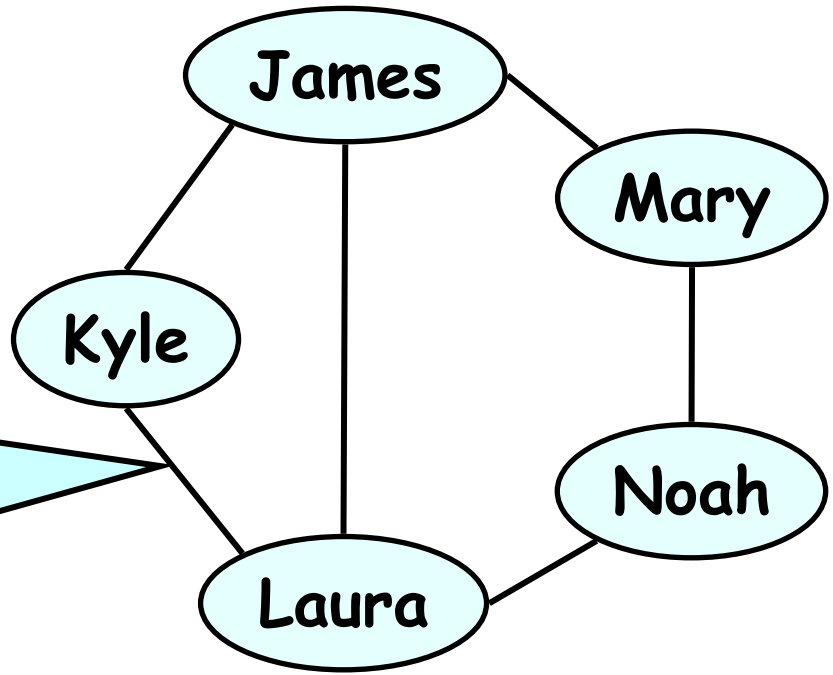
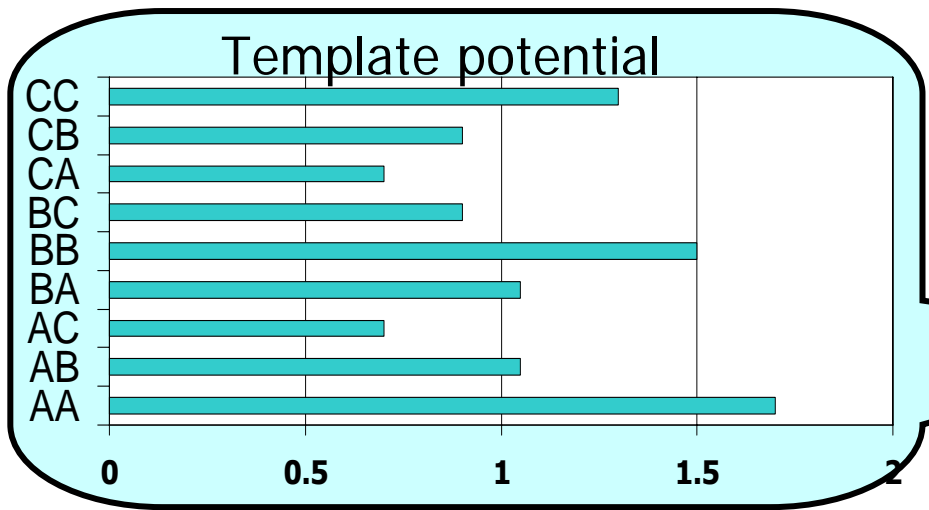


nodes = variables
edges = direct influence

Graph structure encodes independence assumptions:
Job conditionally independent of *Intelligence* given *Grade*



Markov Networks



$$P(J, K, L, M, N) = \frac{1}{Z}$$

$$\phi(J, K)\phi(J, L)\phi(K, L)\phi(J, M)\phi(M, N)\phi(L, N)$$



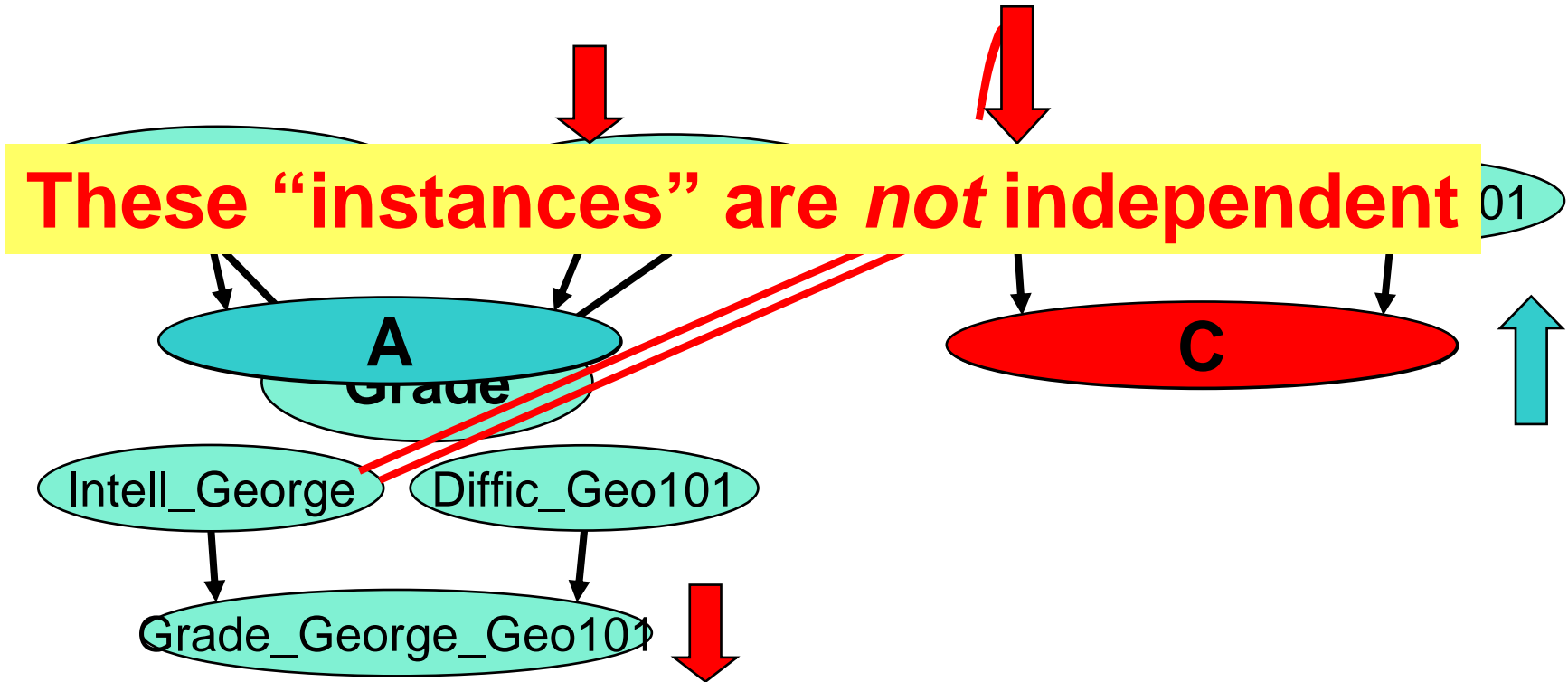
The World is Richly Structured

- The web
 - Webpages (& the entities they represent), hyperlinks
- Biological data
 - Genes, proteins, interactions, regulation
- Physical environments
 - People, rooms, objects
- Natural language



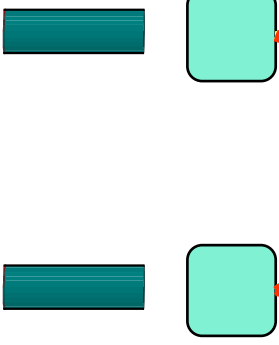
Problem

- Bayesian/Markov nets use attribute representation
- Real world has objects, related to each other



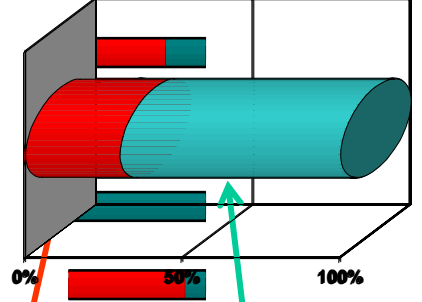
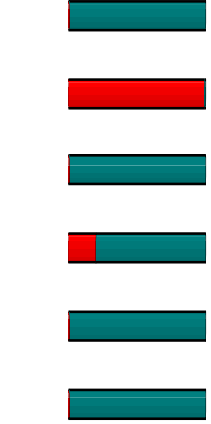


Holistic Reasoning



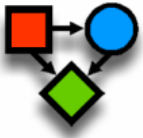
C

A



easy / hard

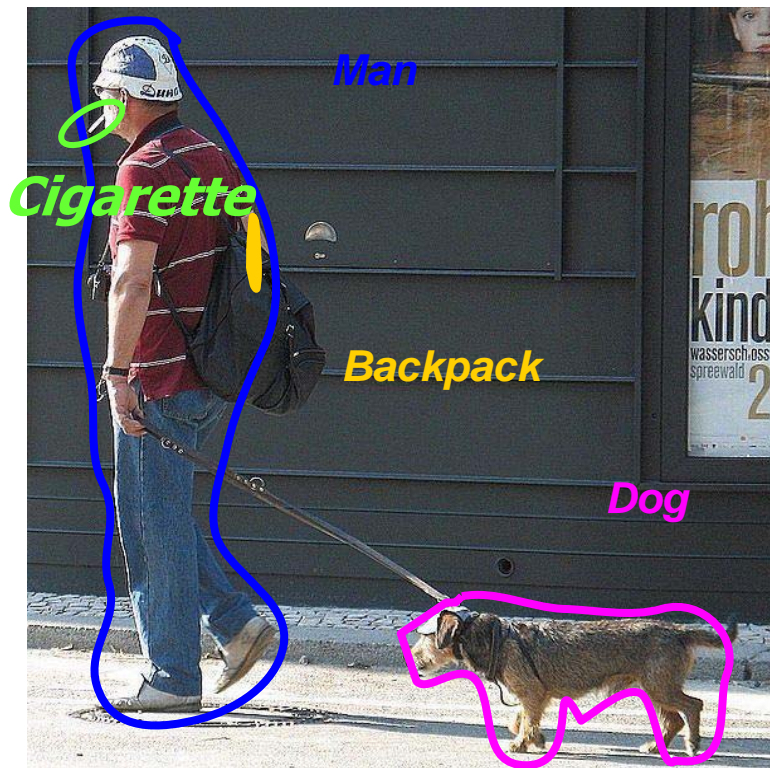
low / high



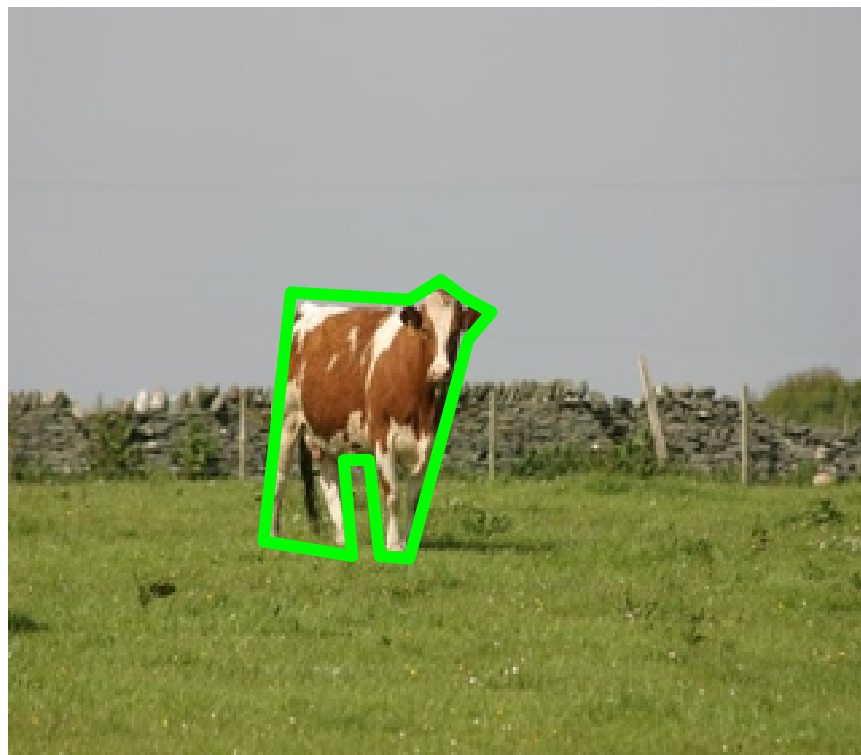
Main Application Domains

- Symbolic understanding of the physical world
- Understanding and reconstructing cellular processes from genomic data

Long-Term Goal: Scene Understanding



***"man wearing a
backpack,
smoking a cigarette,
walking a dog"***



***"A cow walking
through the grass
on a pasture by the sea"***

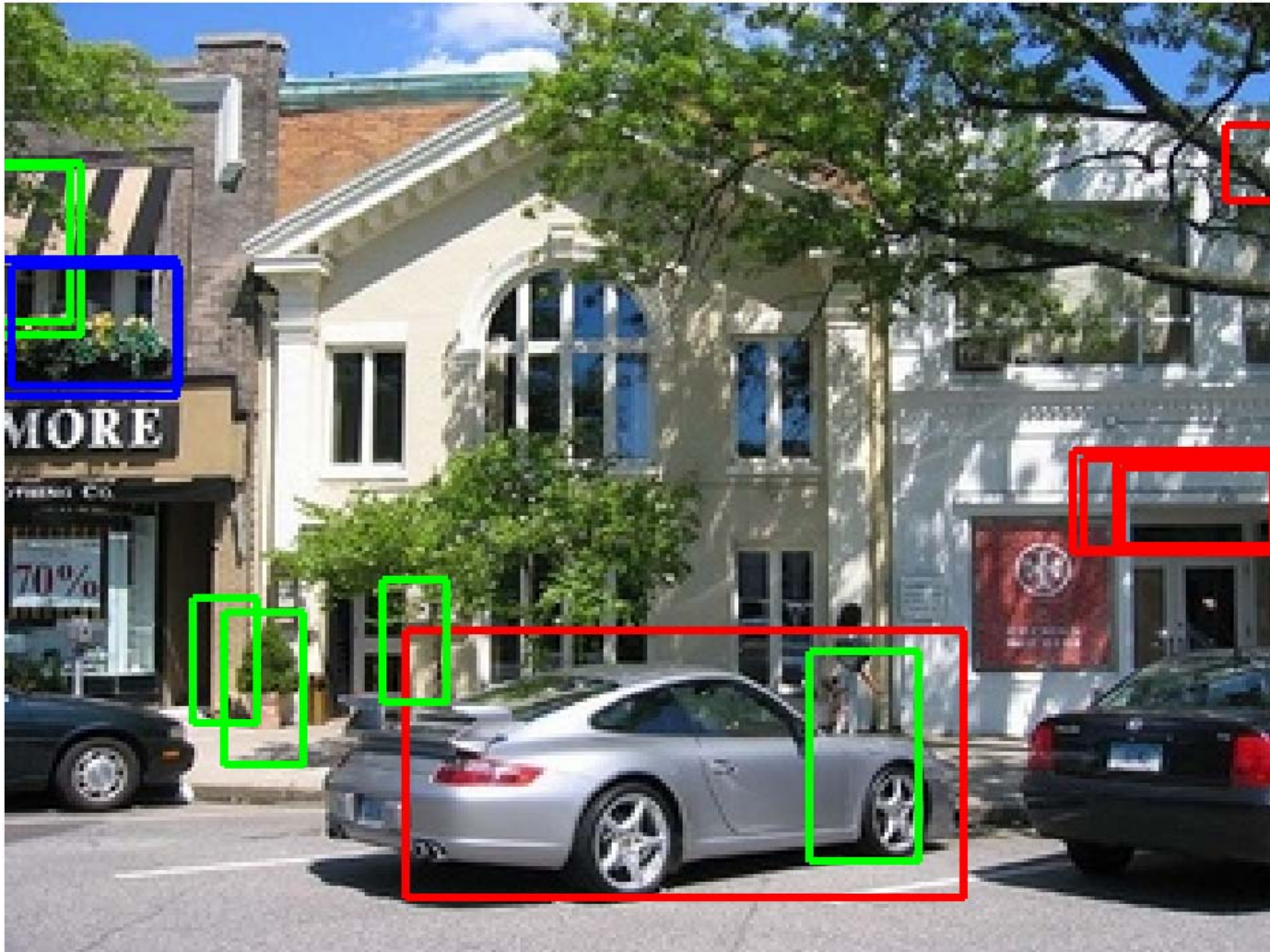


Object Detection





Basic Object Detection



car

person

*motorcycl
e*



Scene Segmentation



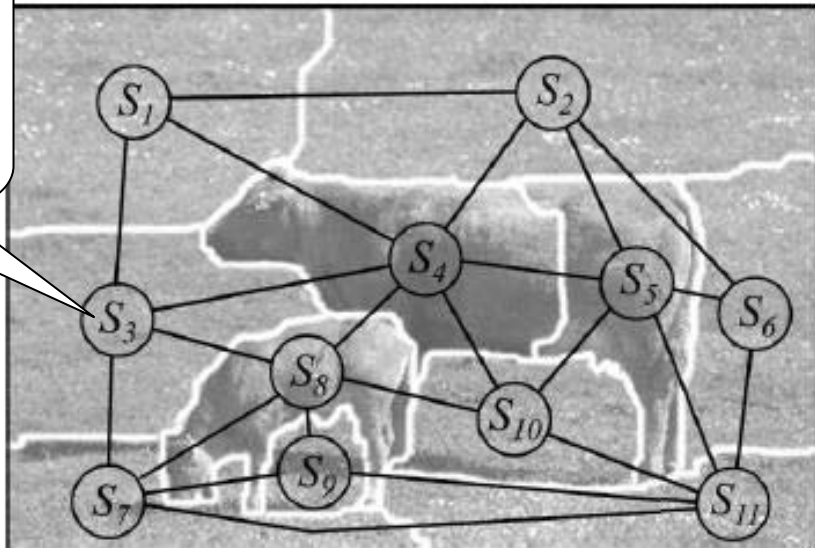
'grass', 'road', 'tree', 'sky', 'water', 'building', 'foreground'



Segmentation CRF



*grass,
road, tree,
water,
building,
sky,
foreground*



(a) Over-segmentation

(b) CRF Construction

$$P(S_1, \dots, S_N) \propto \prod_i \underbrace{\phi_i(S_i : I)}_{\text{Singleton energy}} \prod_{i,j} \underbrace{\phi_{i,j}(S_i, S_j : I)}_{\text{Pairwise energy}}$$

Singleton energy:

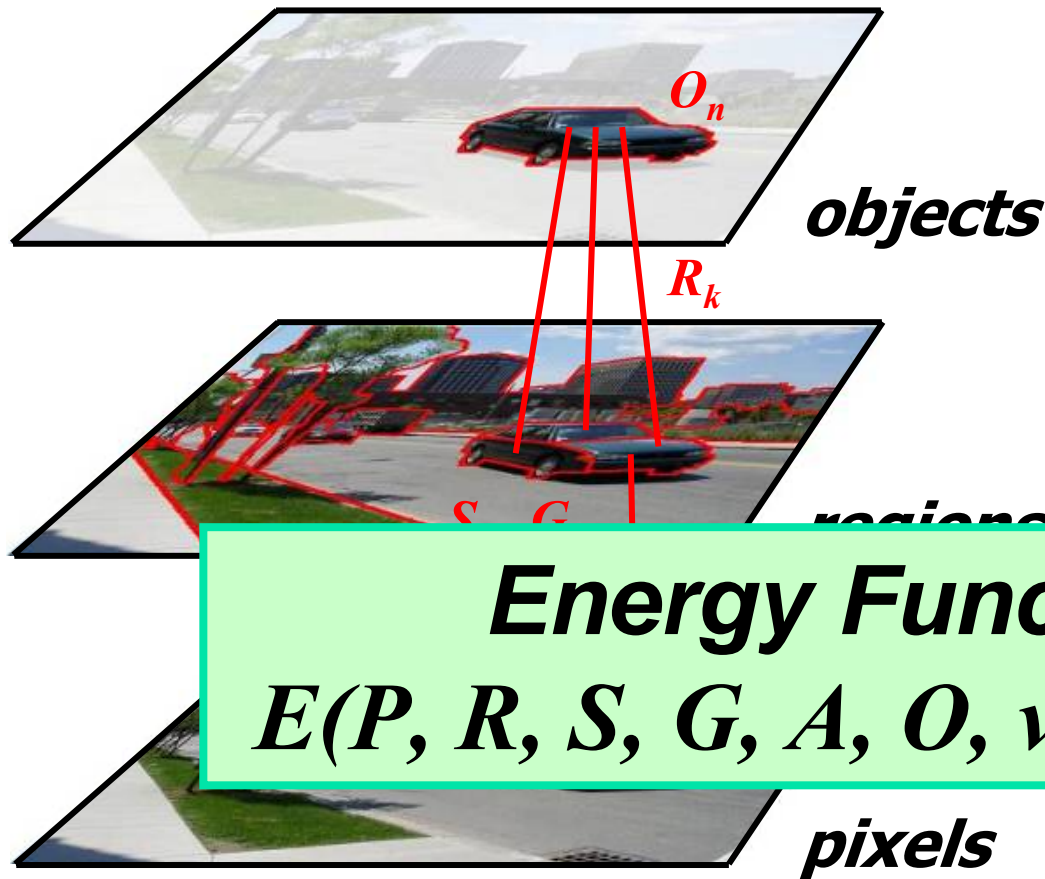
- Mean R,G,B
- Mean H,U,V
- Texture Responses
- ...

Pairwise energy:

- Delta R,G,B
- Offset Vector
- ...



Hierarchical Scene Model



Variables

α_i : pixel appearance
 P_i : pixel-to-region correspondence
 A_k : region appearance
 S_k : region semantic class
 G_k : region geometry
 R_k : region-to-object correspondence
 O_n : object class
 v^{hz} : location of horizon

Energy Function

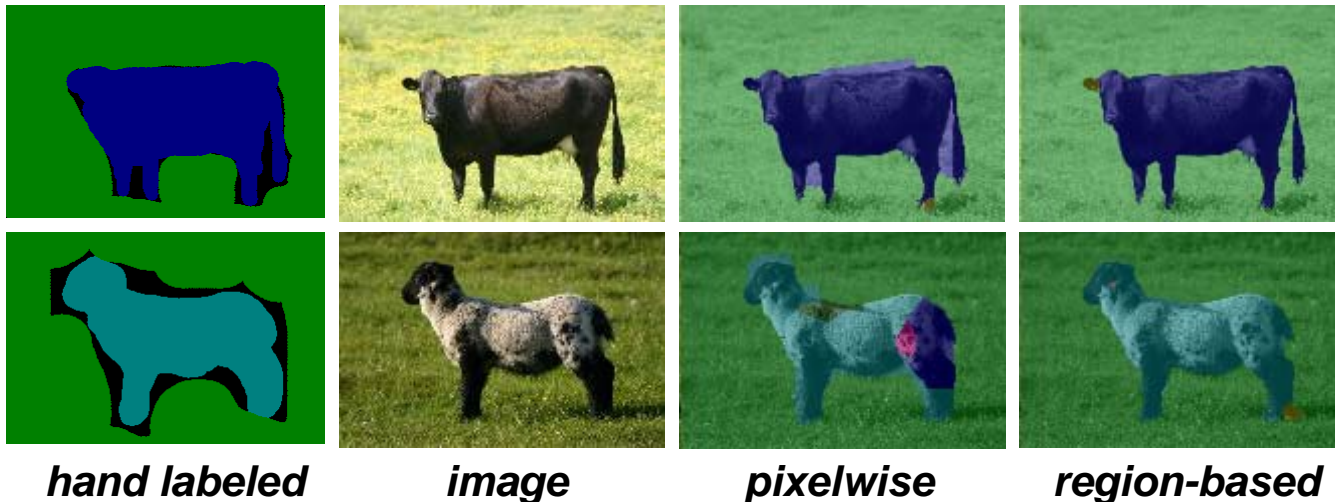
$$E(P, R, S, G, A, O, v^{hz}, K | I, \theta)$$



Results: 21-class MSRC

- Validate against state-of-the-art approaches
- Region/pixel class only
- Ground truth labels are approximate
- No geometry information

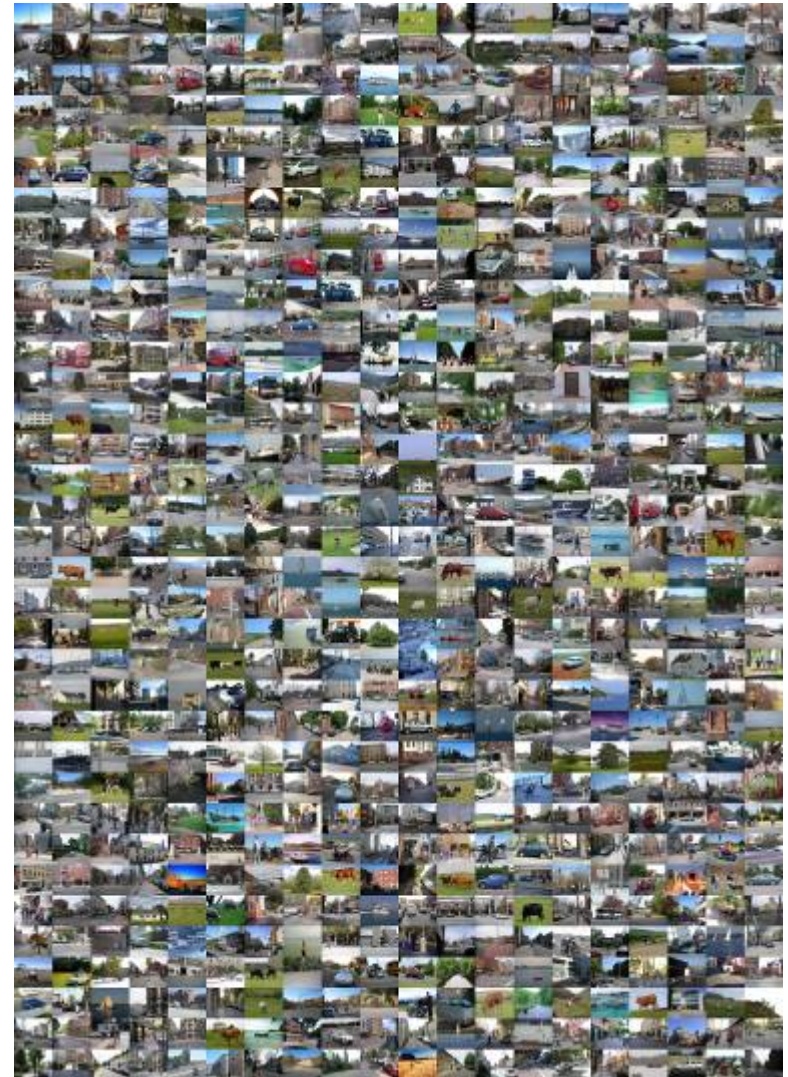
21 CLASS	Mean
<i>Shotton et al.</i>	72.2
<i>Gould et al.</i>	76.5
Pixelwise	75.3
Region-based	75.4





High Quality Dataset

- MSRC dataset is limited
 - poorly labeled boundaries
 - many missing pixels (void)
 - no geometry information
- Collected images from MSRC, Hoiem et al., Pascal VOC
- 715 outdoor scenes with high-quality labels
 - region boundaries
 - region class and geometry
 - horizon



[Gould, Fulton, Koller, ICCV 2009]



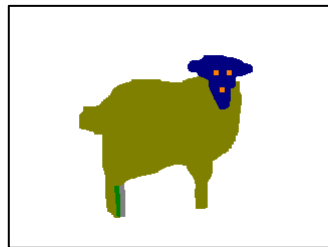
Amazon Mechanical Turk (AMT)

- \$0.10 per task (regions, classes, surface types)
- 5-10 minutes per task
- 24-48 hour turn-around time (for 715 images)
- Less than 10% of tasks needed rework
- **Total cost for labels:** under \$250 (includes \$40 textbook on Adobe Flash)
- **Saving Steve from having to label images:** priceless.





AMT: Label Quality



Typical quality (hand labeled)



You don't always get what you want

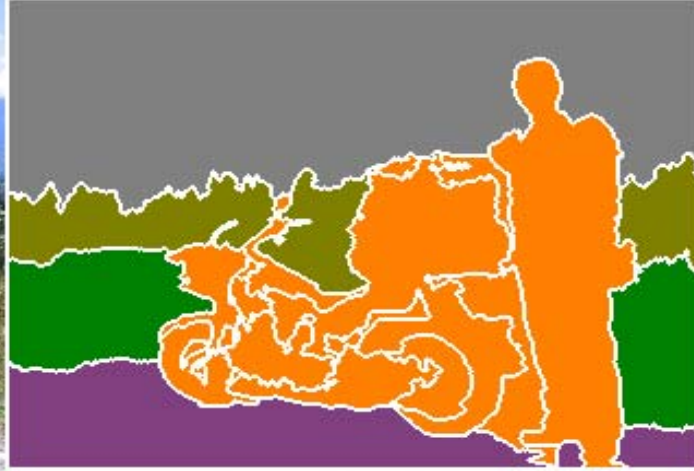
Comparison with MSRC labels



[Gould, Fulton, Koller, ICCV 2009]



Example Results

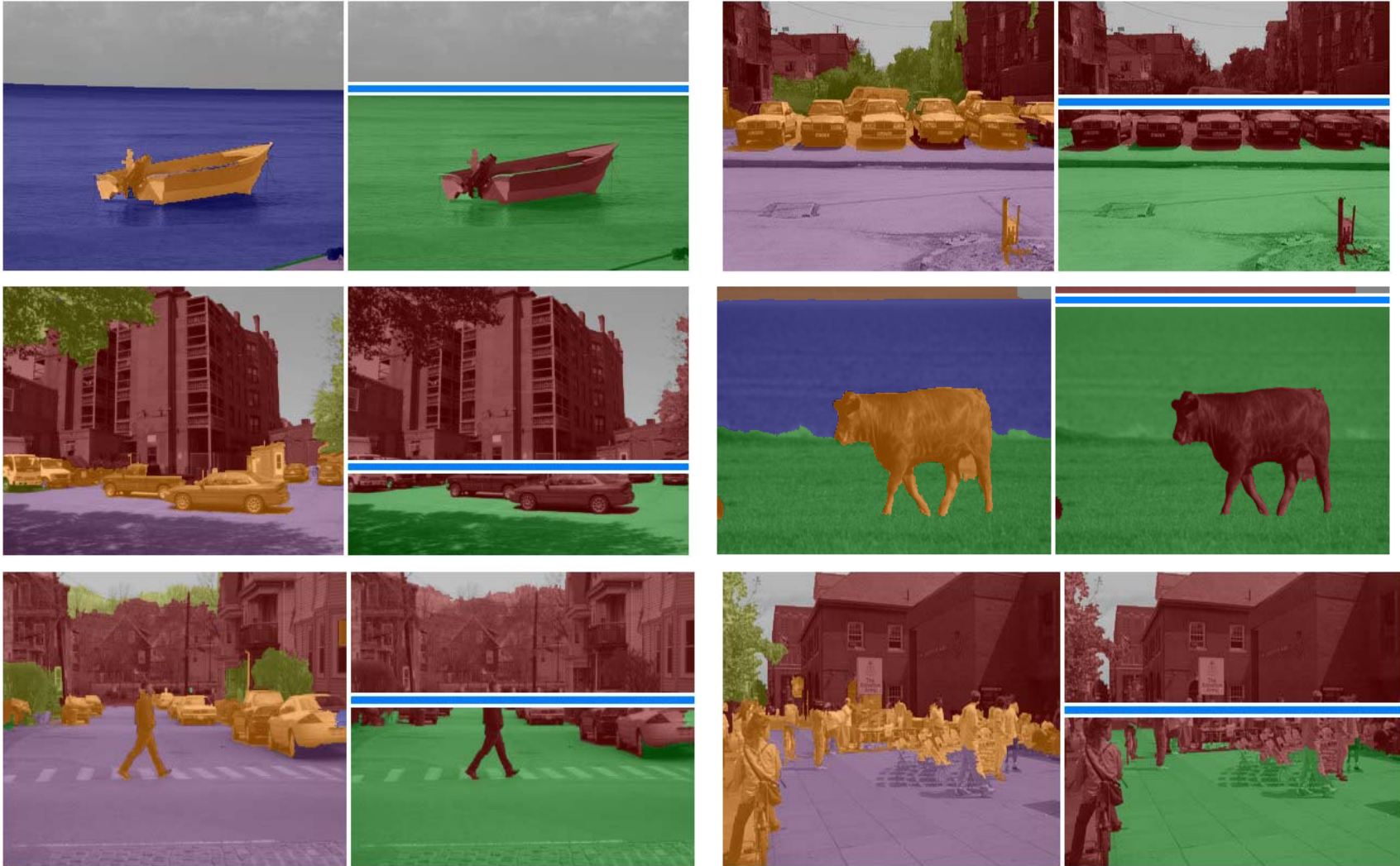


■ sky ■ tree ■ road ■ grass ■ water ■ bldg ■ mntn ■ fg obj. ■ sky ■ horz. ■ vert.

[Gould, Fulton, Koller, ICCV 2009]



More Example Results

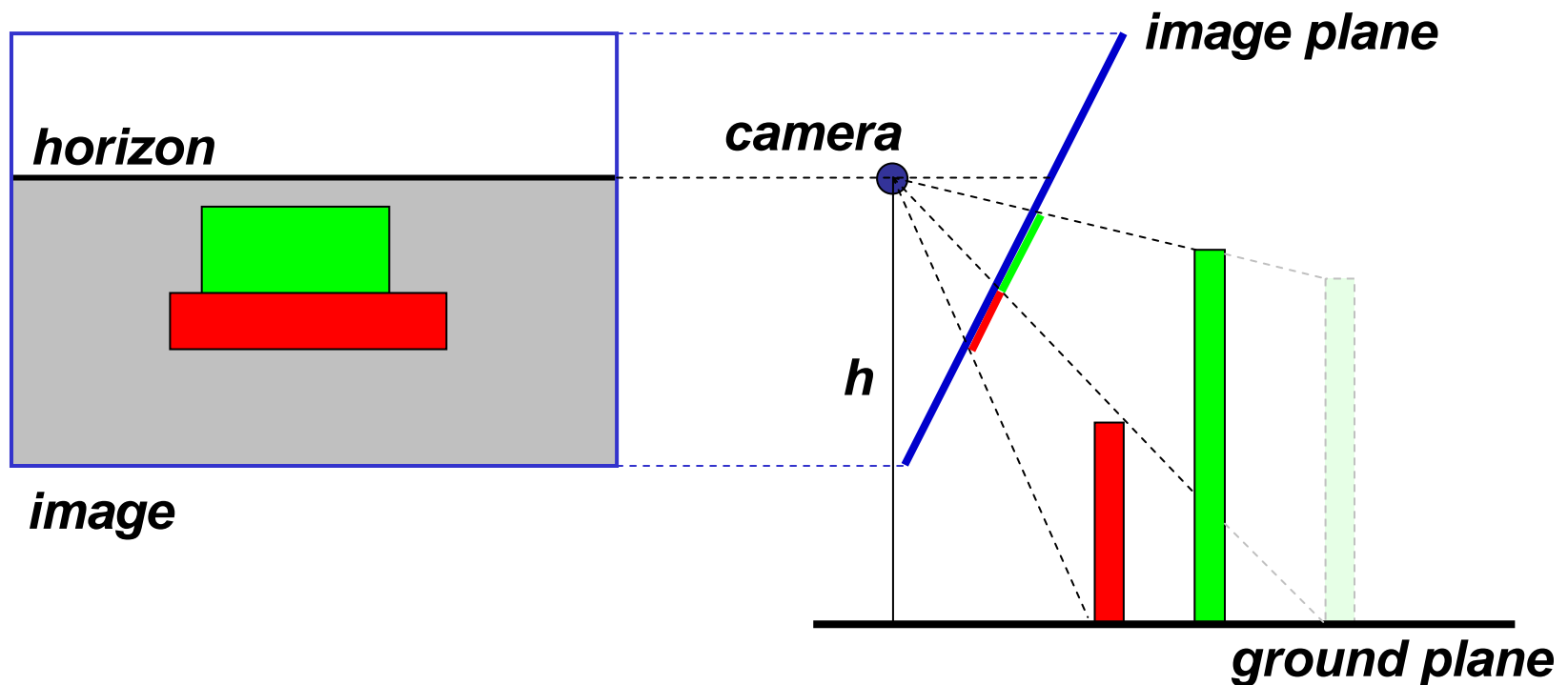


[Gould, Fulton, Koller, ICCV 2009]



Application: 3d Reconstruction

- Estimate camera tilt from location of horizon
- Predict region 3D position using ray projected through camera plane





Example 3D Reconstructions





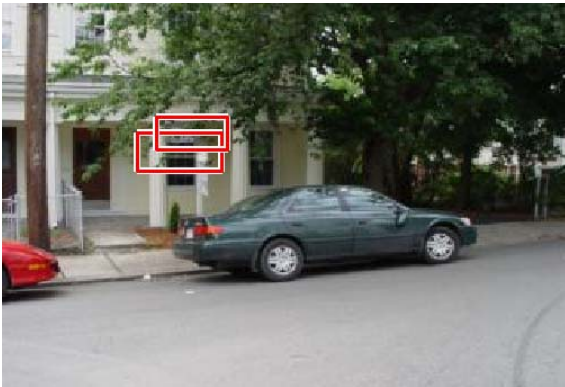
Object Detection Examples



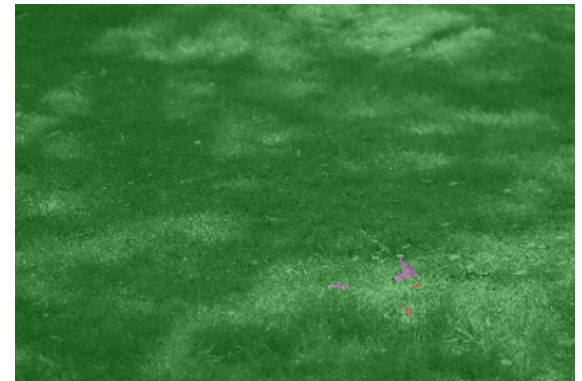


Sliding-Window Failures

Sliding-window detector top results

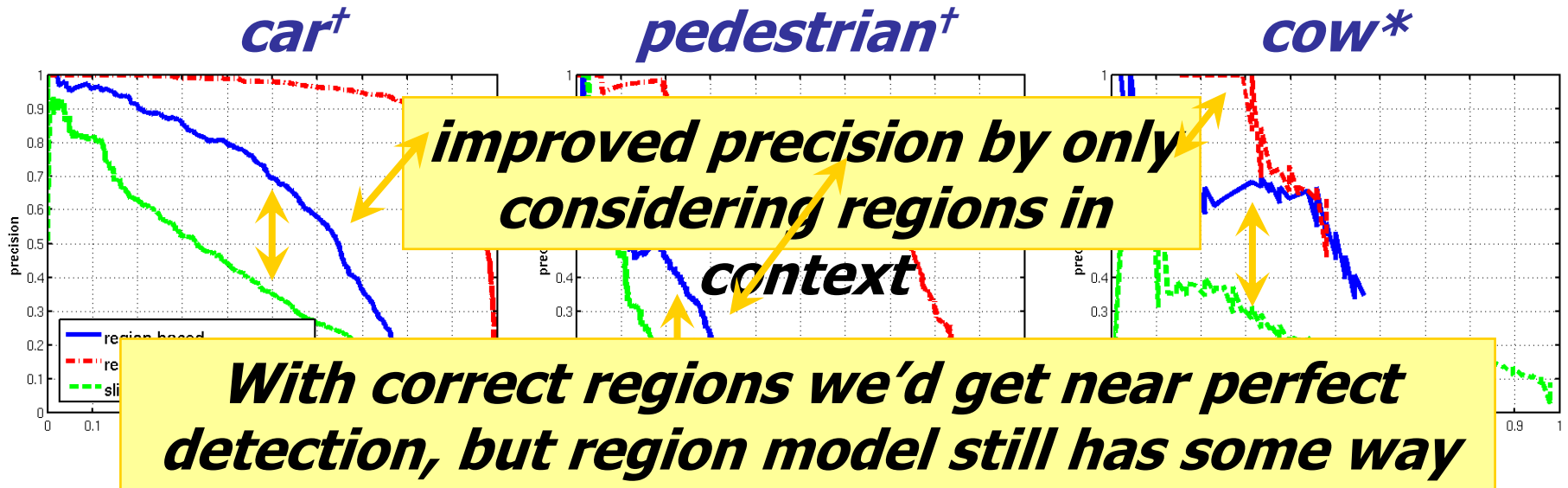


Our region-based object detector results





Object Detection



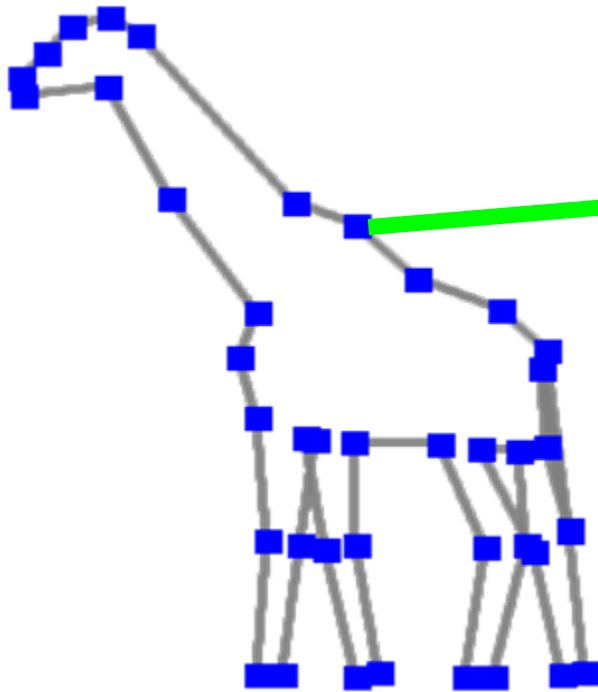
[†] run on Street Scene dataset

^{*} run on subset of 21-class MSRC data

[Gould, Gao, Koller submitted]



Shape-Based Segmentation

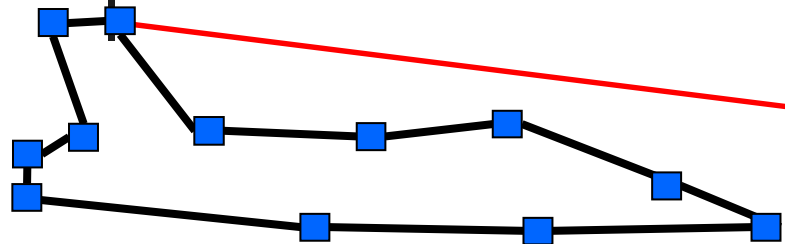


- Set of shape landmarks
- Shape defined by connecting piecewise-linear contour
- Semantic outlining = assignment L of landmarks to pixels

[Heitz, Elidan, Packer, Koller, NIPS-08b]



The LOOPS Model



L – assignment of model landmarks to pixels



$$P(L \mid O = 1, I, w, \mu, \Sigma) =$$

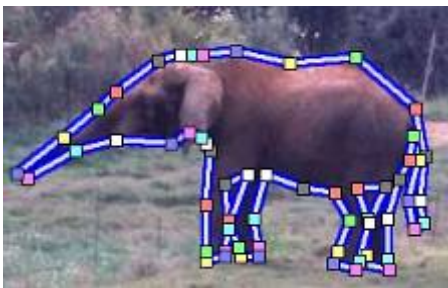
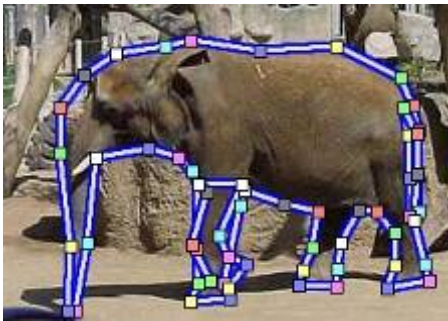
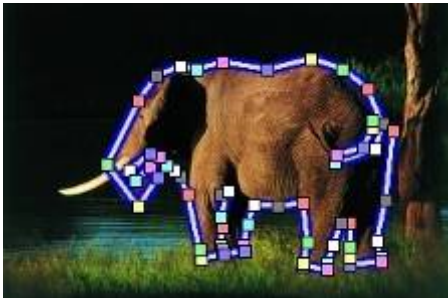
pairwise image features

$$\frac{1}{Z} \underbrace{P_{\text{SHAPE}}(L; \mu, \Sigma)}_{\text{prior}} \prod_{i,j} \exp \left\{ w_{i,j} \underbrace{F_{i,j}(l_i, l_j; I)}_{\text{pairwise image features}} \right\}$$

Outlining = MAP Inference over *L*

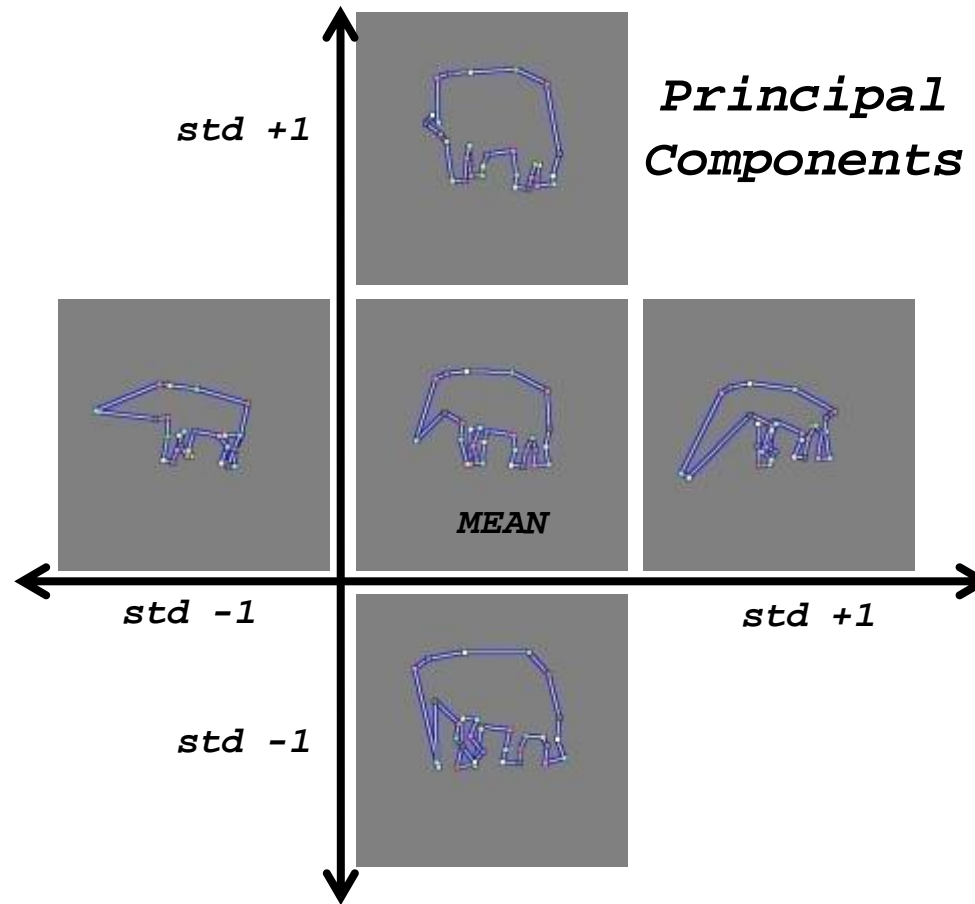
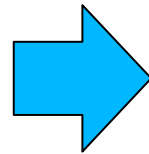
Learning the Shape Model

Training Set:



Problem:

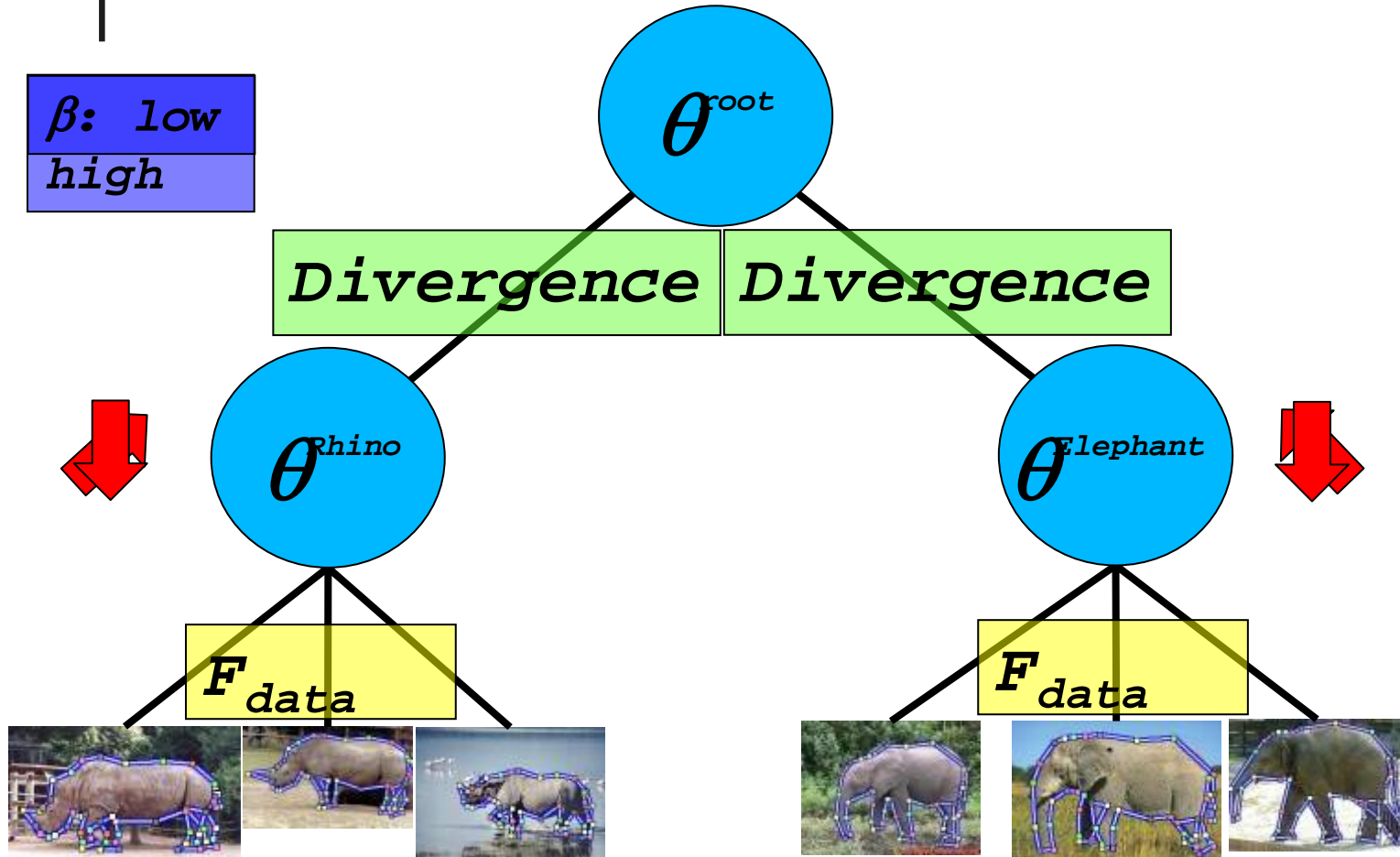
With few instances, learned models aren't robust





Undirected Probabilistic Model

β : low
high



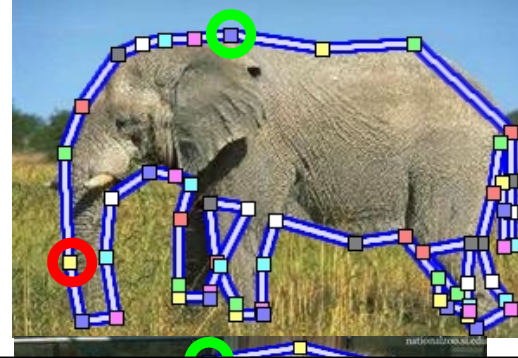
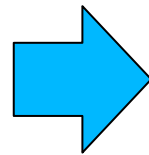
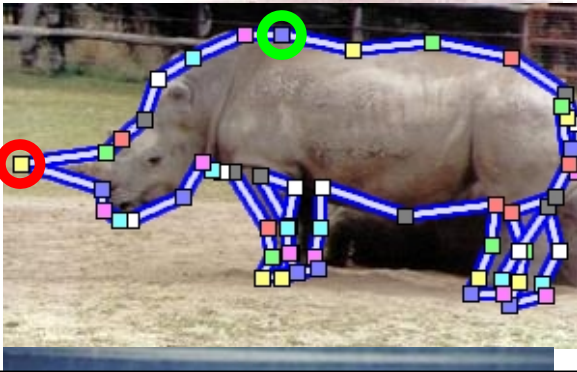
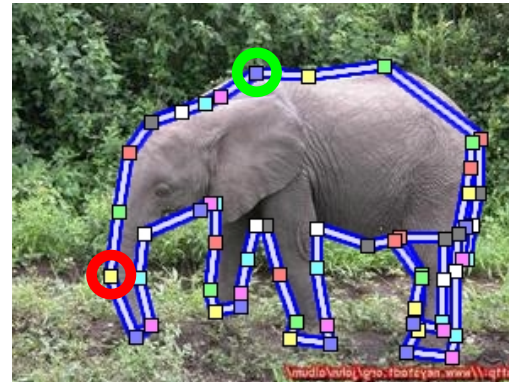
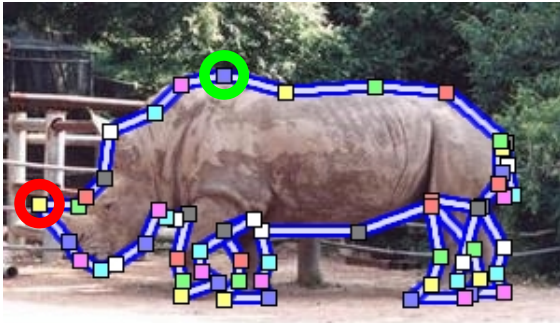
F_{data} :

Encourage parameters to explain data

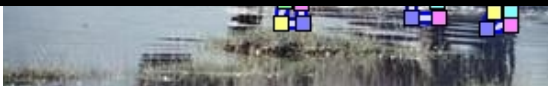
Divergence:

Encourage parameters to be similar to parents

Degrees of Transfer



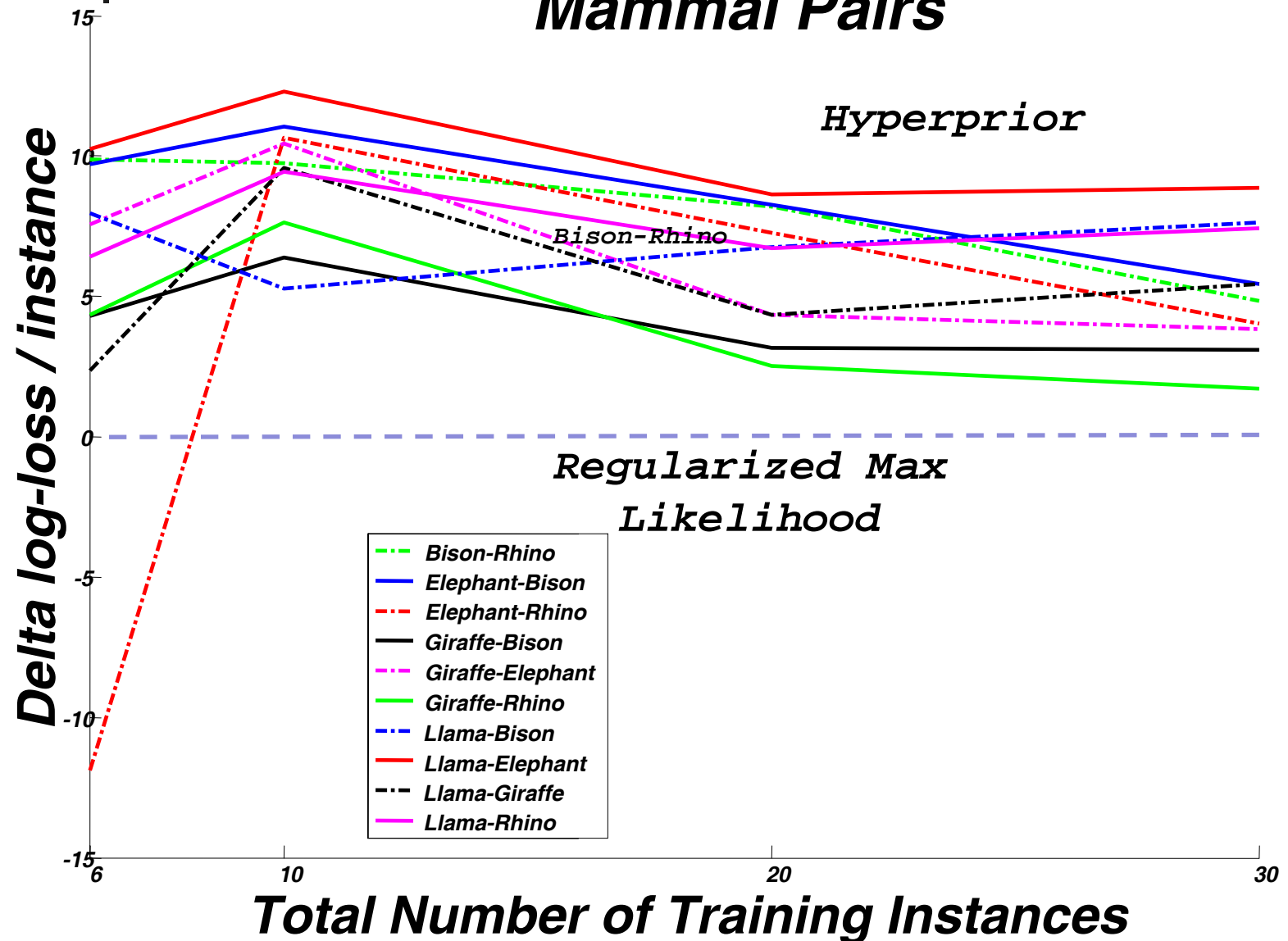
Not all parameters deserve equal sharing



Do Degrees of Transfer Help?



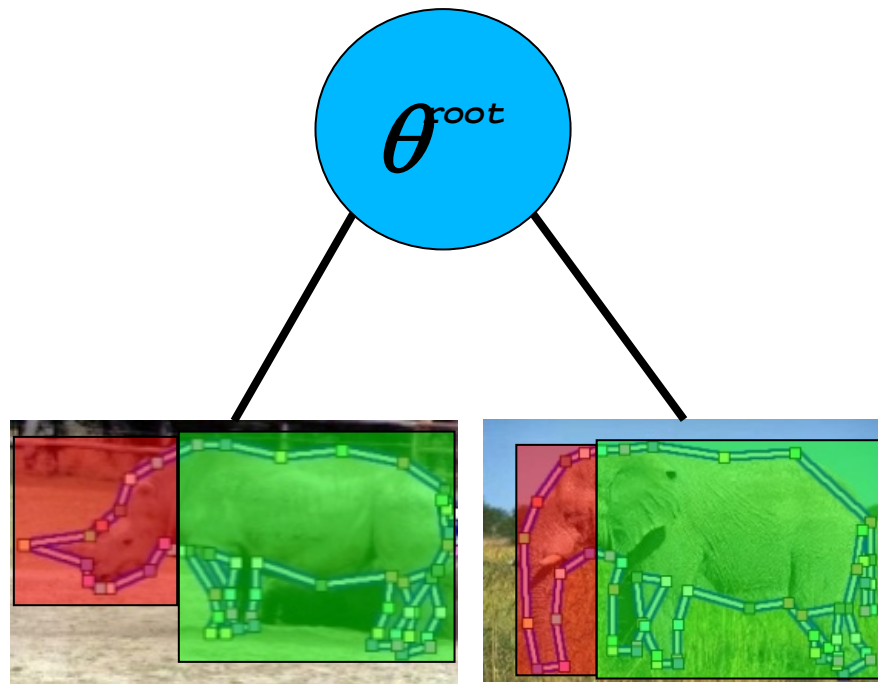
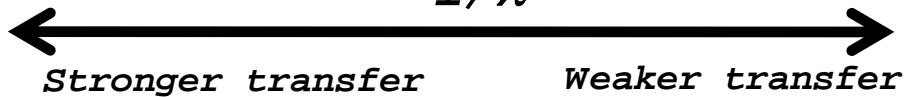
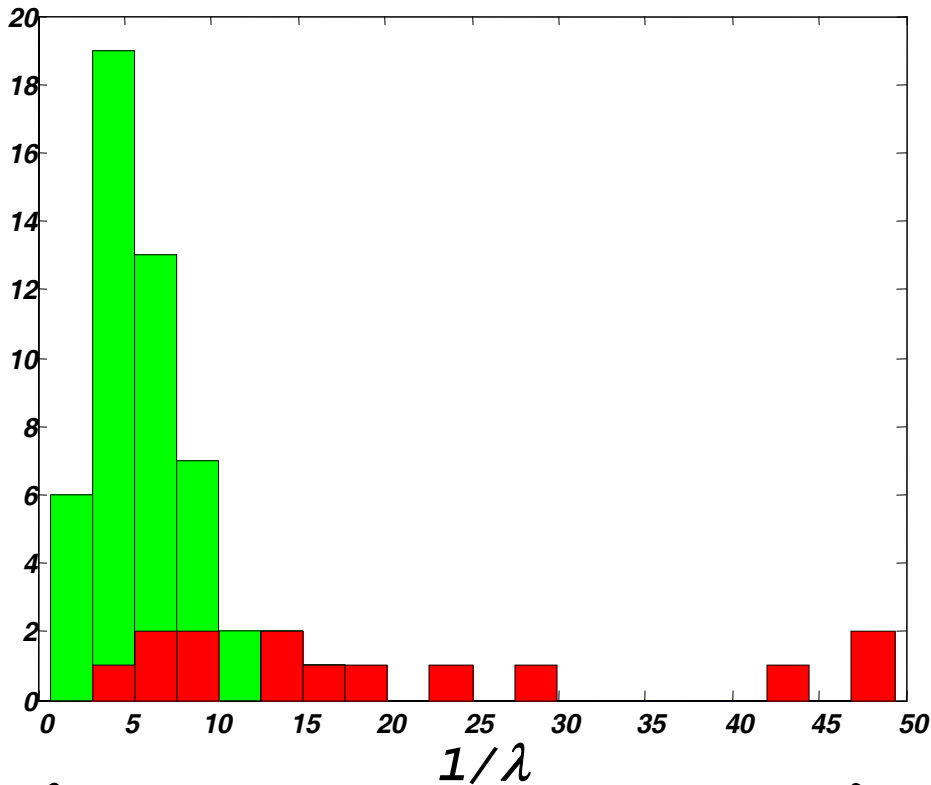
Mammal Pairs





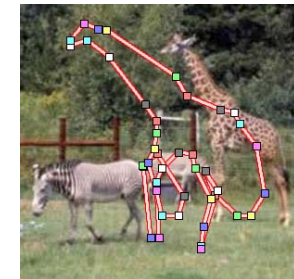
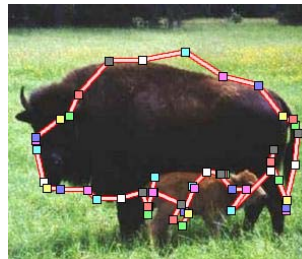
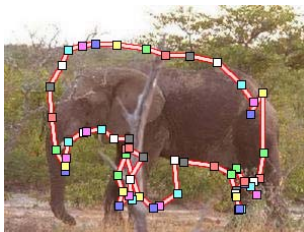
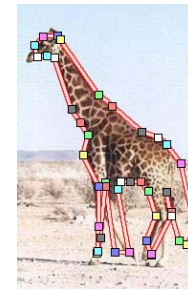
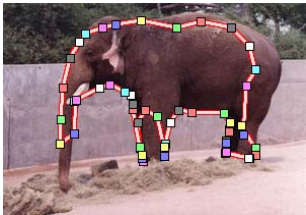
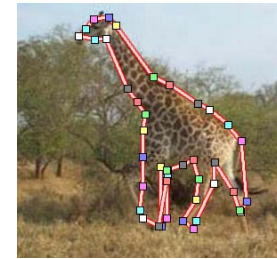
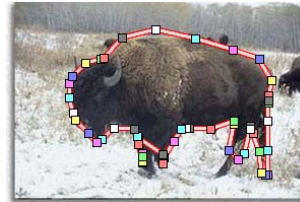
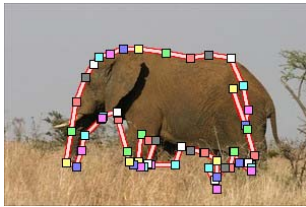
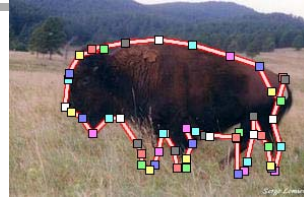
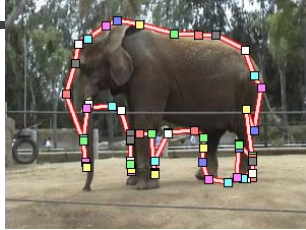
Degrees of Transfer

Distribution of DOT coefficients using Hyperprior





Outlining Results: Mammals

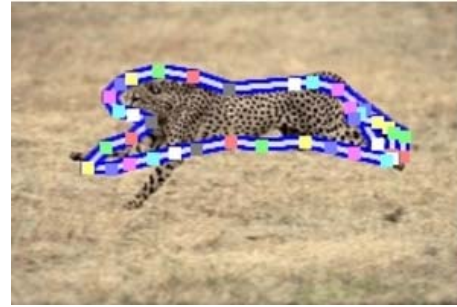


[Heitz, Elidan, Packer, Koller, NIPS-08b]

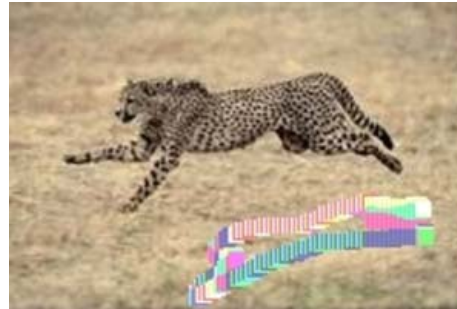
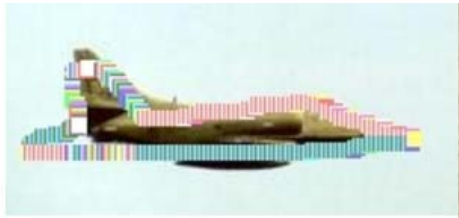
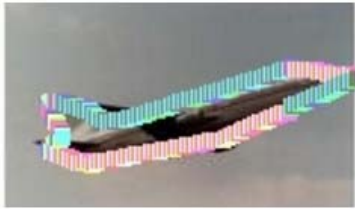


Outlining: Comparison

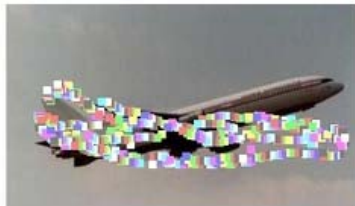
LOOPS



OBJ CUT [Kumar et al., CVPR 05]



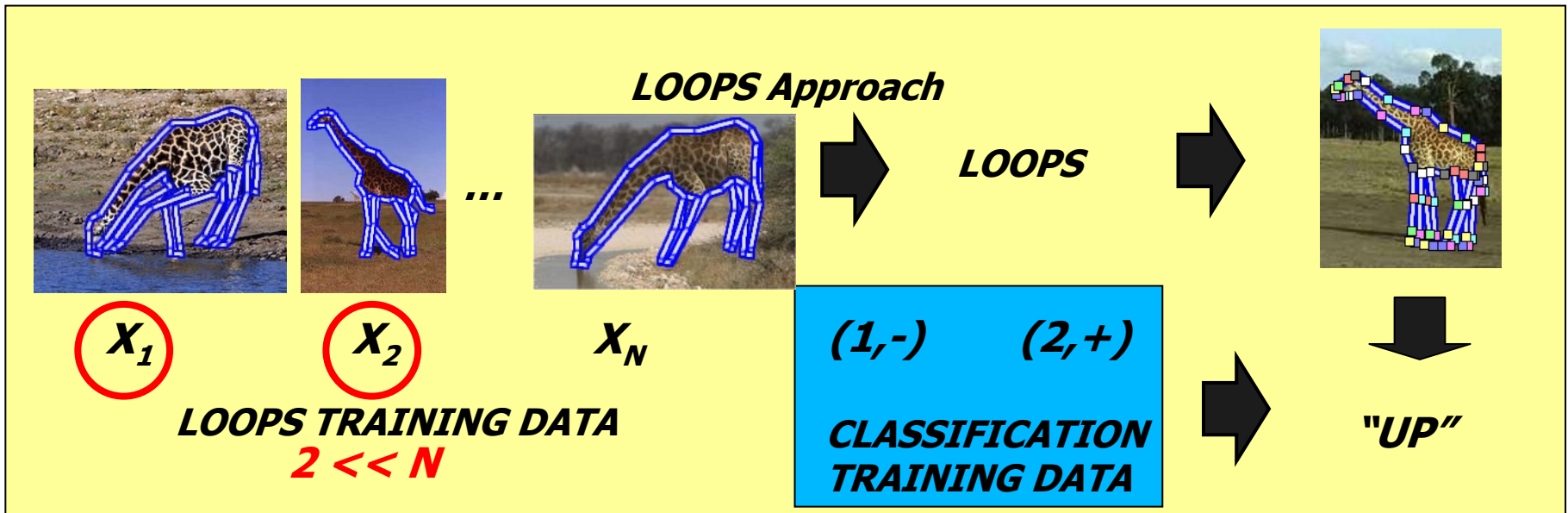
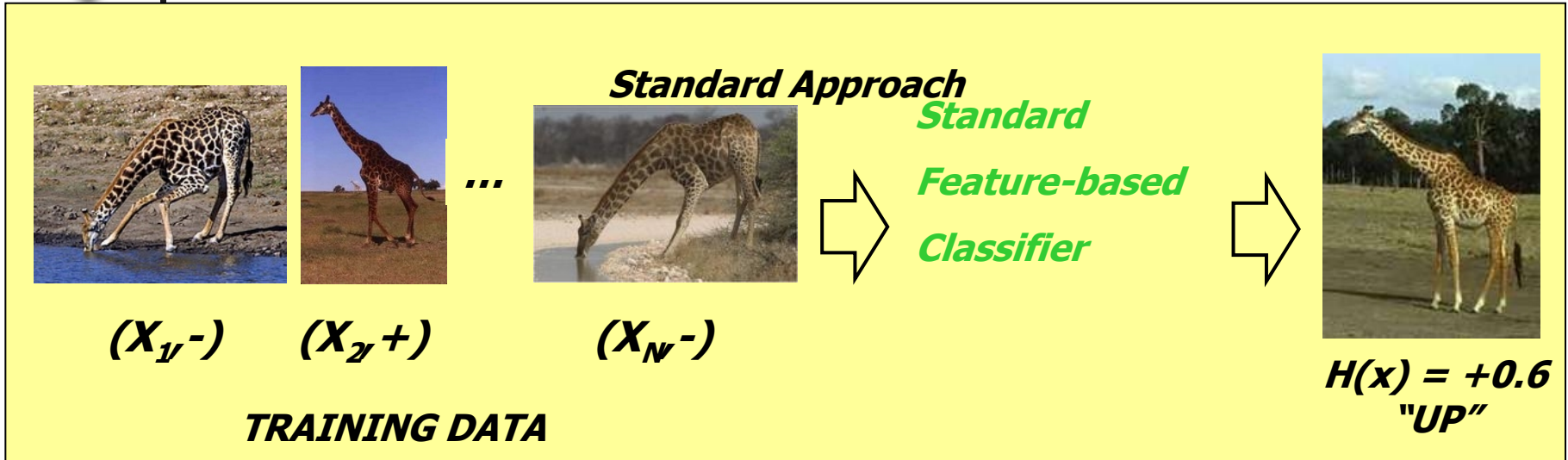
kAS [Ferrari et al., CVPR 07]



[Heitz, Elidan, Packer, Koller, NIPS-08b]

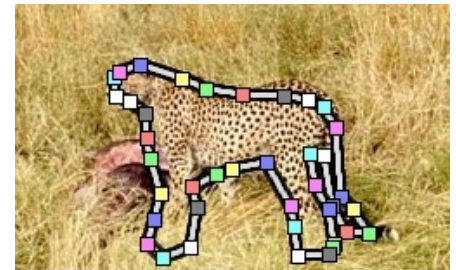
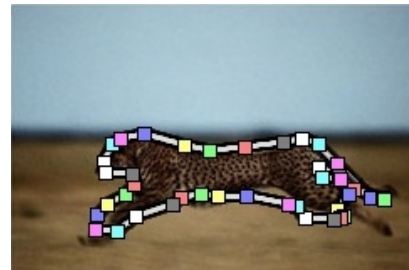
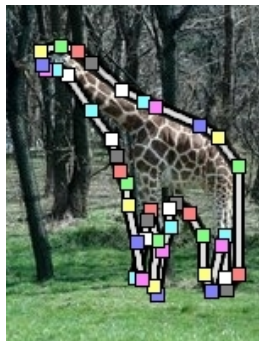
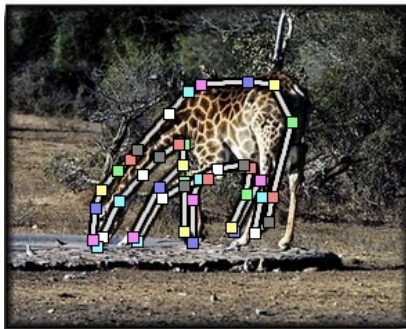
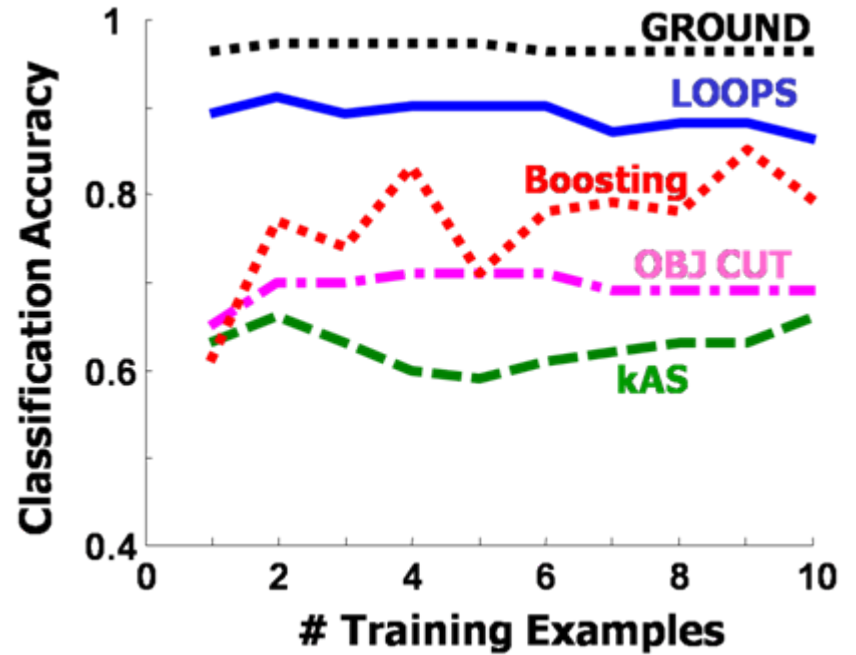
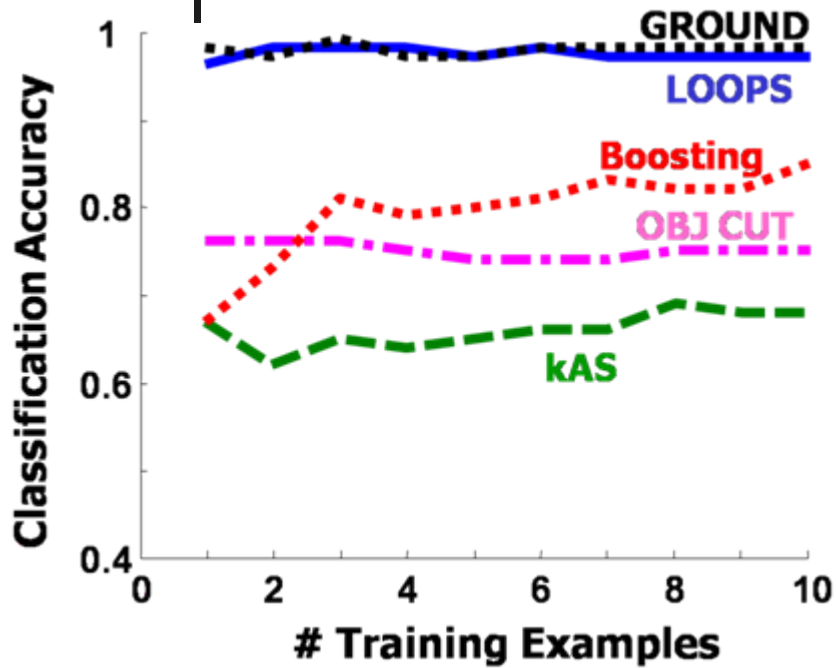


Descriptive Classification



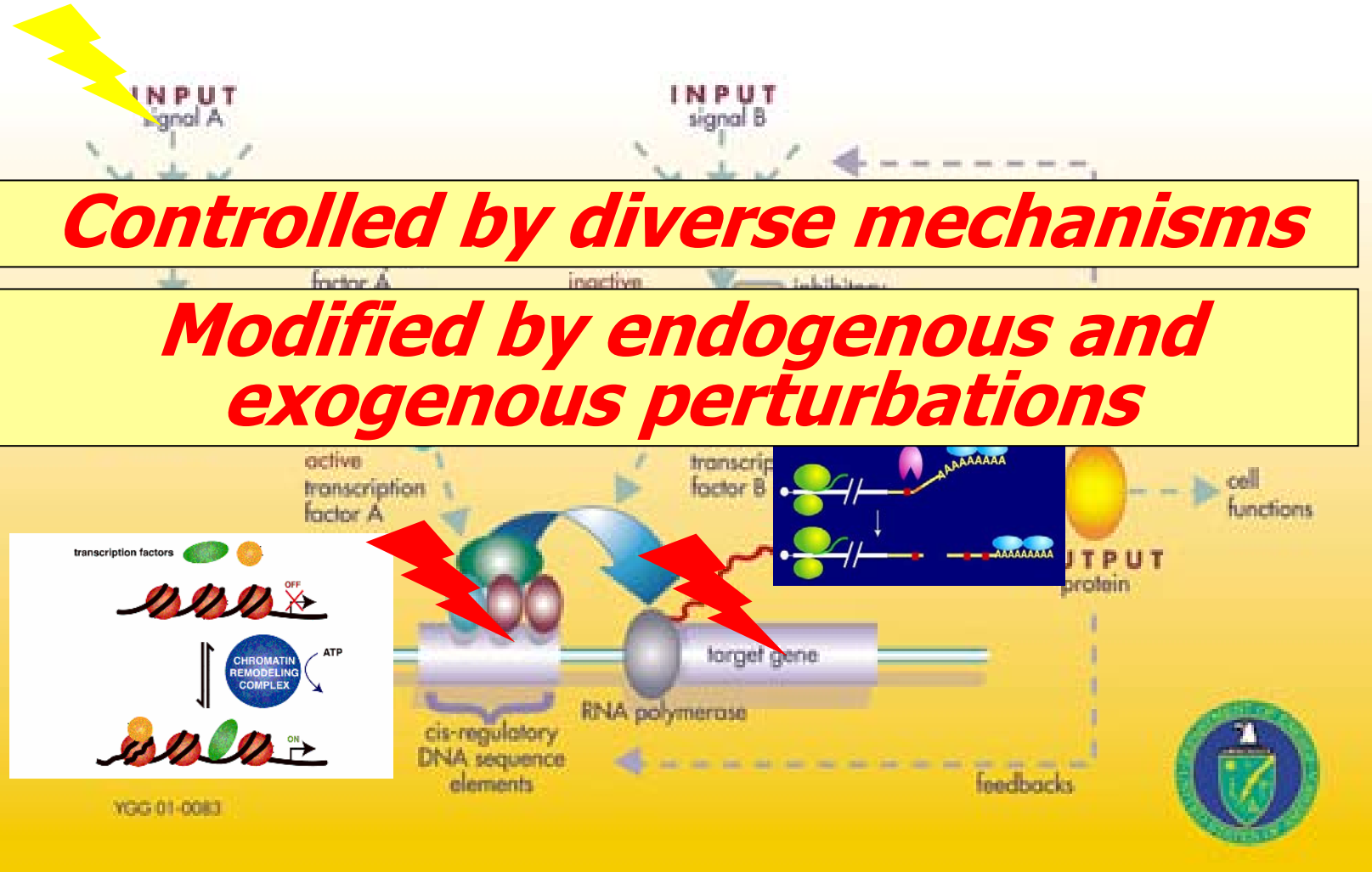


Mammals





Gene Regulatory Networks





Goals

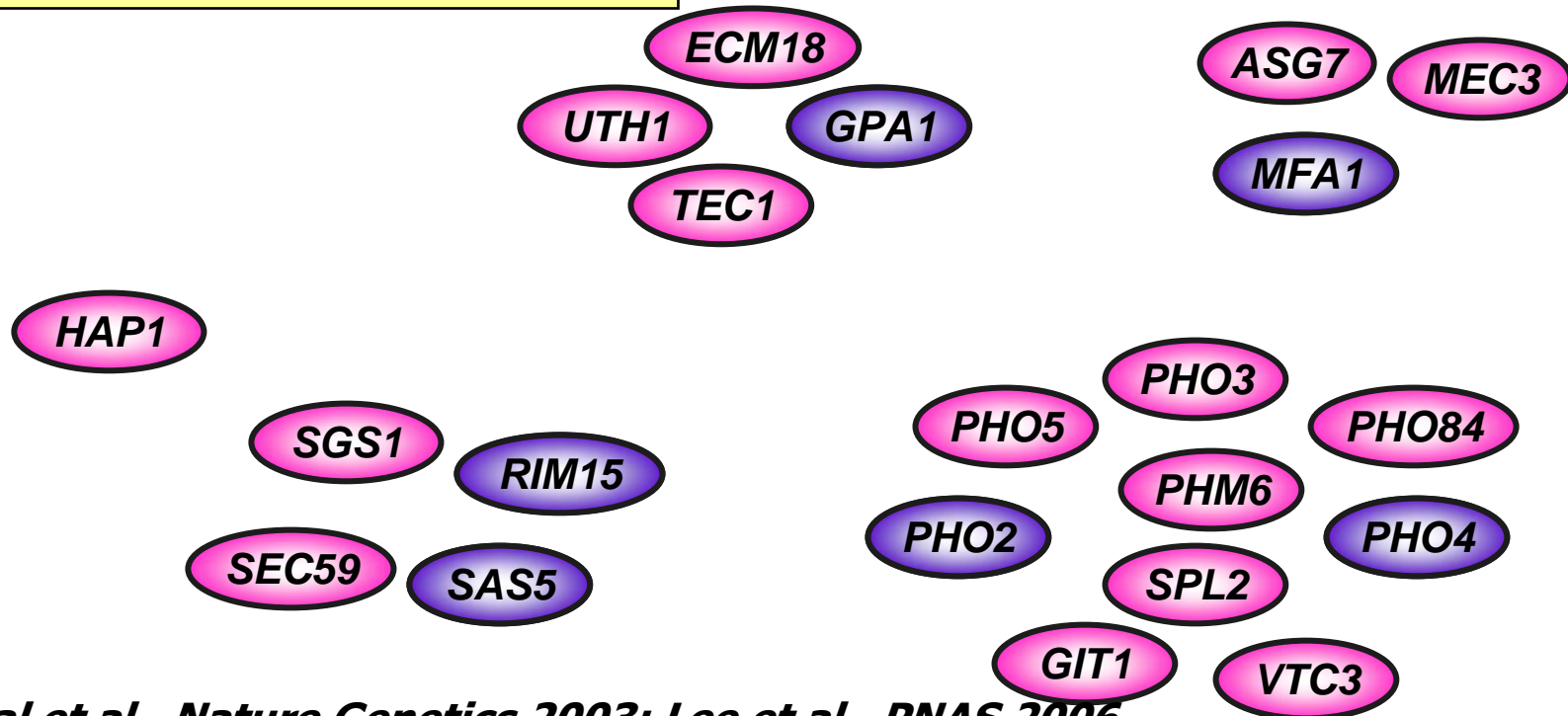
- Infer **regulatory network and mechanisms** that control gene expression
- Identify **effect of perturbations** on network
- Understand effect of gene regulation on **phenotype**



Regulatory Network I

- mRNA level of regulator can indicate its activity level
- Target expression is predicted by expression of its regulators
- Use expression of regulatory genes as regulators

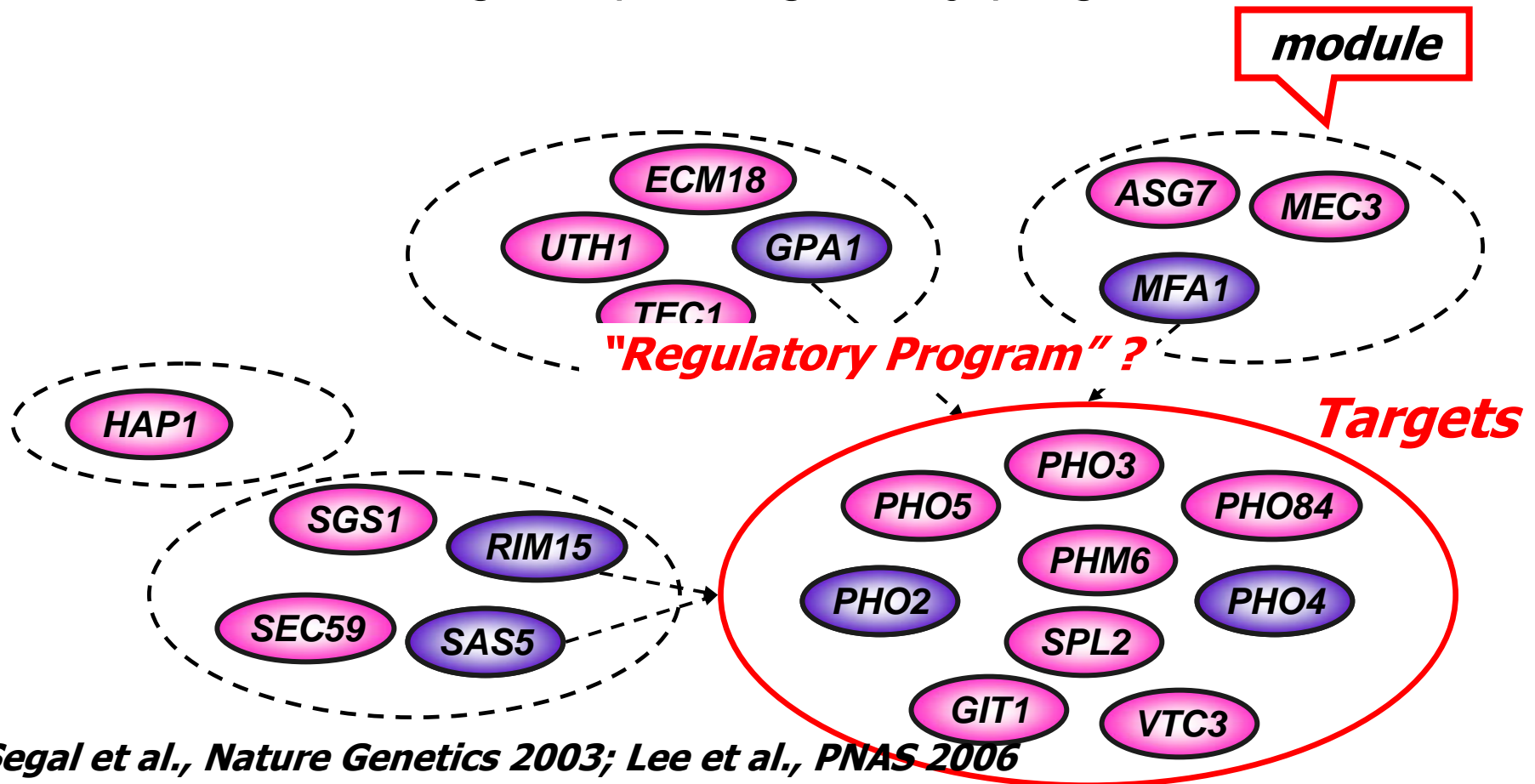
Transcription factors, signal transduction proteins, mRNA processing factors, ...





Regulatory Network II

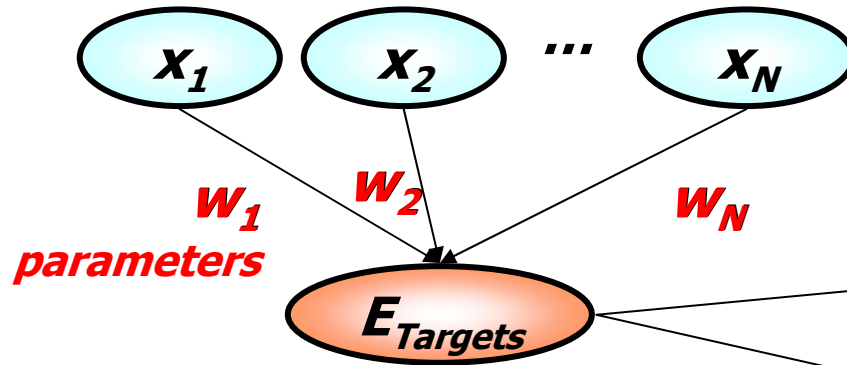
- Co-regulated genes have similar regulation program
- Exploit modularity and predict expression of entire module
- Allows uncovering complex regulatory programs



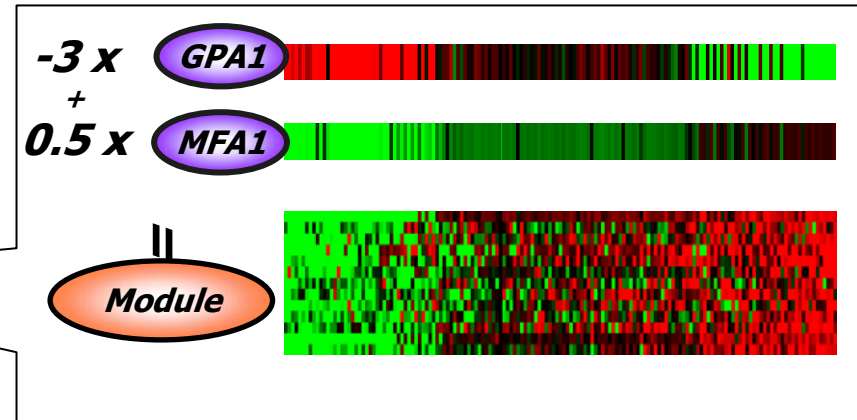


Regulation as Linear Regression

$$\text{minimize}_w (\sum w_i x_i - E_{\text{Targets}})^2$$



$$E_{\text{Targets}} = w_1 x_1 + \dots + w_N x_N + \epsilon$$



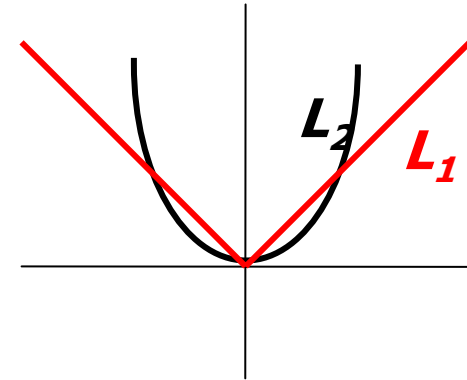
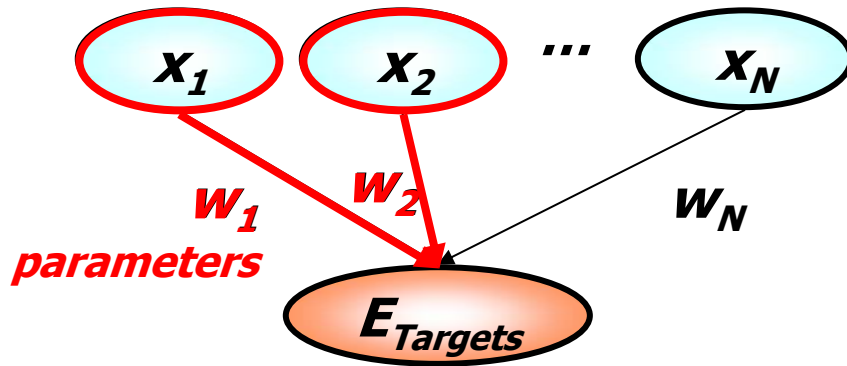
- But we often have hundreds or thousands of regulators
- ... and linear regression gives them all nonzero weight!

Problem: This objective learns too many regulators



Lasso* (L_1) Regression

$$\text{minimize}_{\mathbf{w}} (w_1x_1 + \dots + w_Nx_N - E_{\text{Targets}})^2 + \sum \mathbf{C} |w_i|$$

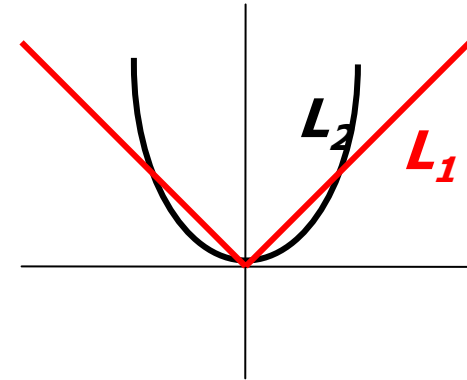
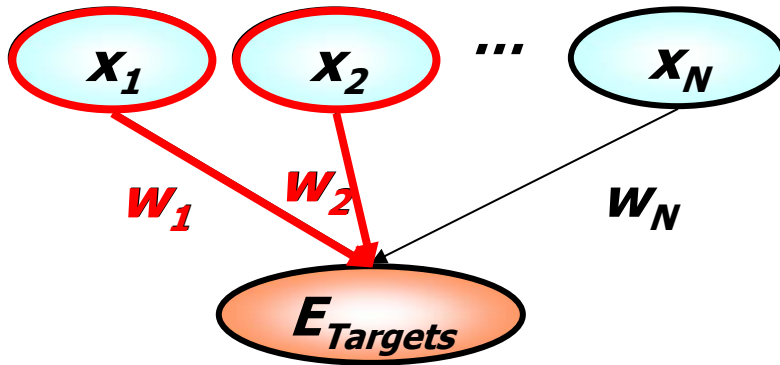


- Induces sparsity in the solution \mathbf{w} (many w_i 's set to zero)
 - Provably selects "right" features when many features are irrelevant
- Convex optimization problem
 - Unique global optimum
 - Efficient optimization
- But, arbitrary choice among correlated regulators



Elastic Net* Regression

$$\text{minimize}_{\mathbf{w}} (w_1x_1 + \dots + w_Nx_N - E_{\text{Targets}})^2 + \sum \mathbf{C} |w_i| + \sum \mathbf{D} w_i^2$$

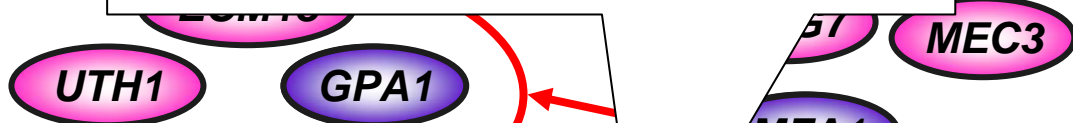
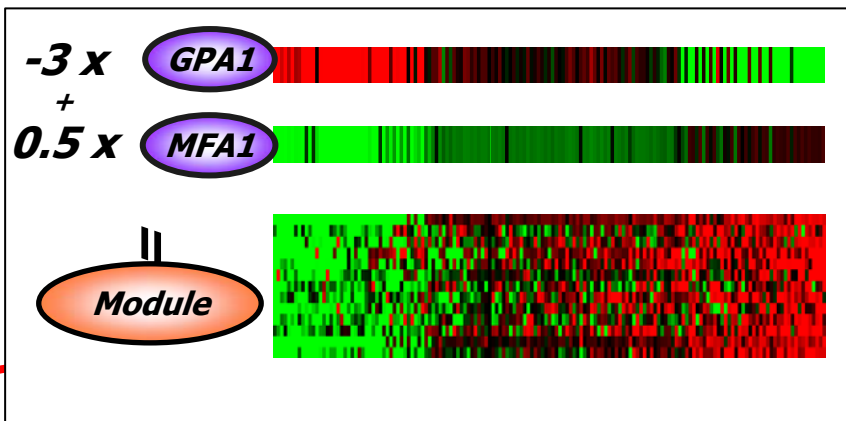


- Induces sparsity
- But avoids arbitrary choices among relevant features
- Convex optimization problem
 - Unique global optimum
 - Efficient optimization algorithms



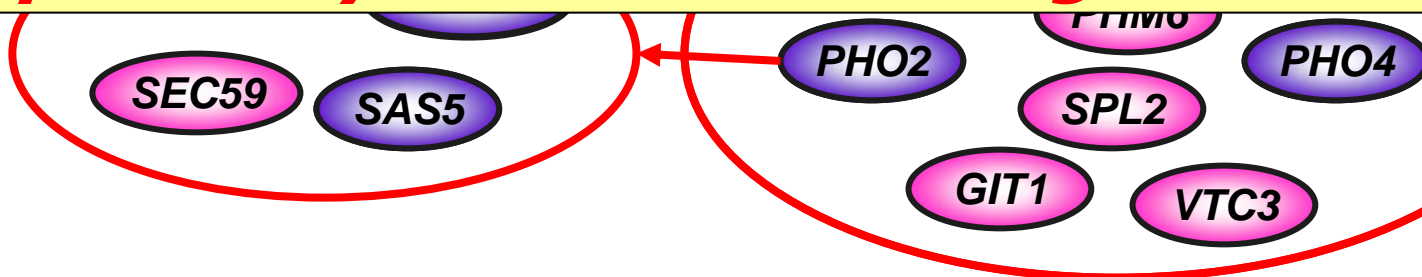
Learning Regulatory Network

- Cluster genes into
- Learn a regulatory



This is a Bayesian network

- But multiple genes share same program***
- Dependency model is linear regression***





Genotype → phenotype

***Different
sequences***

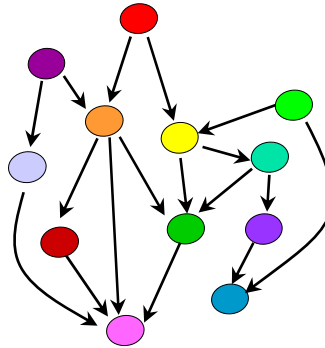
...ACTCGGTTGGCCTAAATTCGGCCCGG...

...ACCCGGTAGGCCTTAATTCGGCCCGG.

:

...ACTCGGTAGGCCTATATTCGGCCGGG...

***Perturbations to
regulatory network***



***Different
phenotypes***

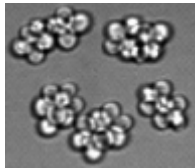
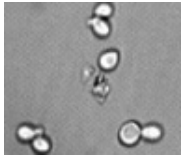




eQTL Data [Brem et al. (2002) Science]

BY

RM



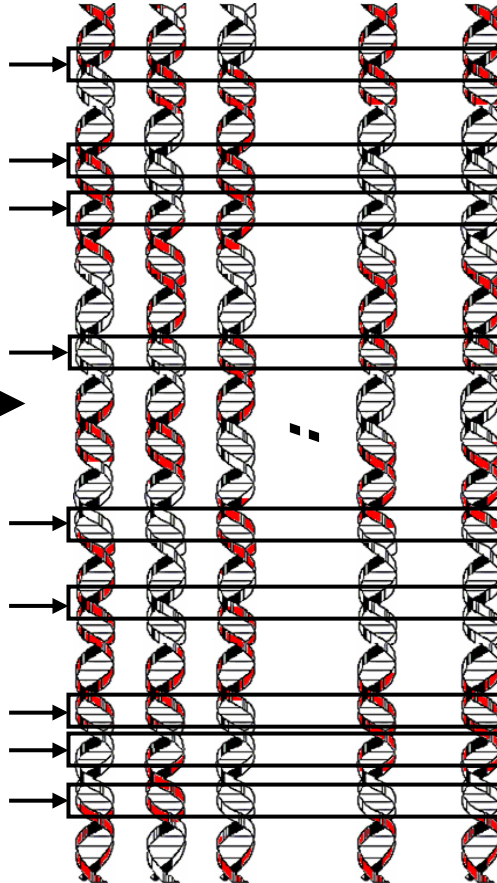
**112
progeny**



X



Markers



Genotype data

112 individuals

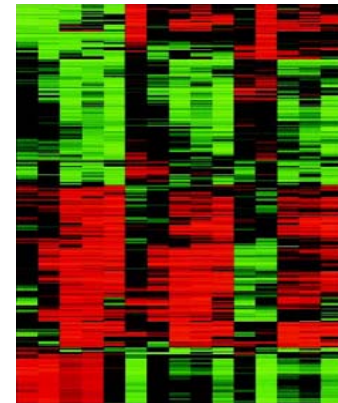
3000 markers

```
0101100100...011
1011110100...001
0010110000...010
:
0000010100...101
0010000000...100
```

Expression data

112 individuals

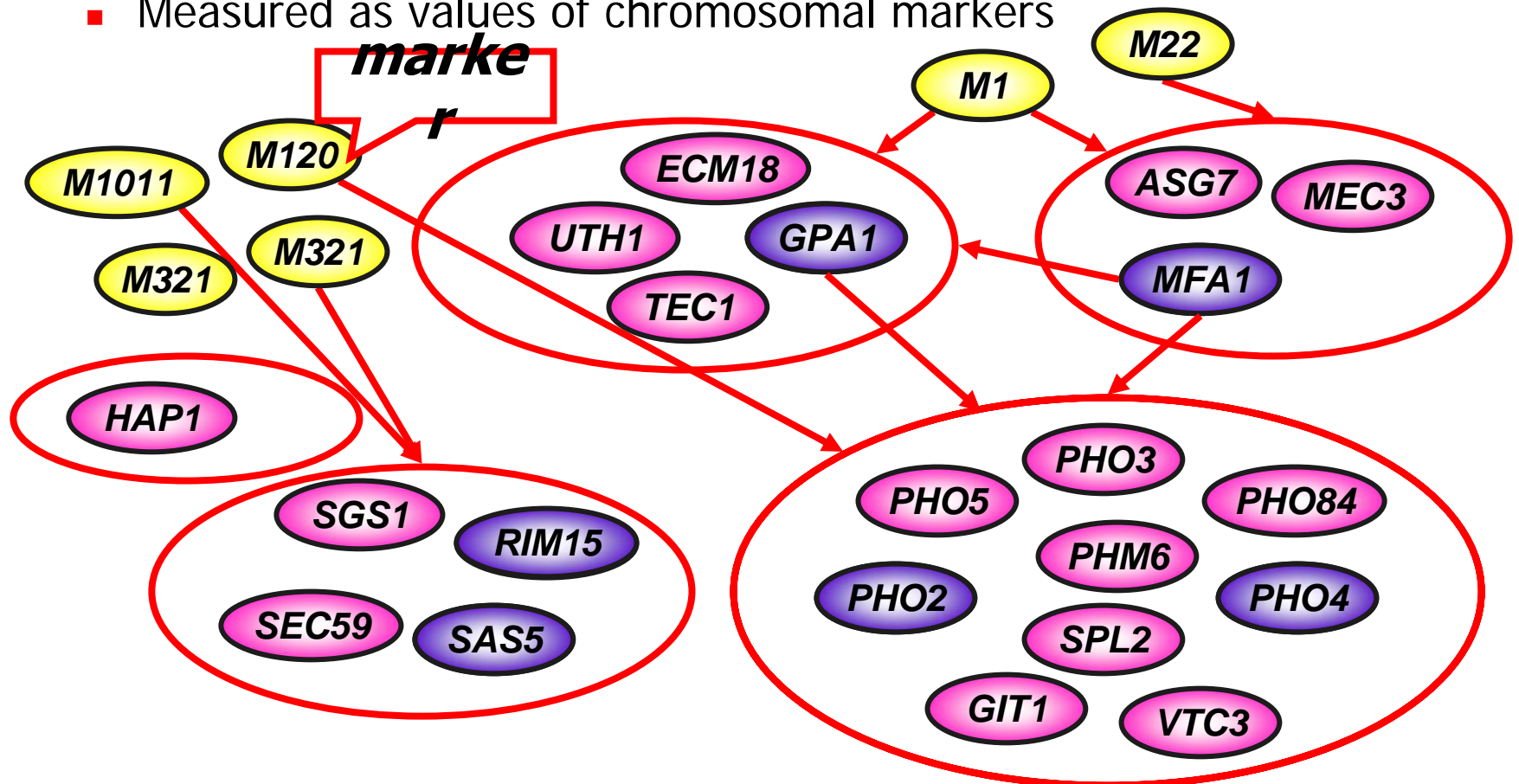
6000 genes





LirNet Regulatory network

- E-regulators: Activity (expression) of regulatory genes
- G-regulators: Genotype of genes
 - Measured as values of chromosomal markers

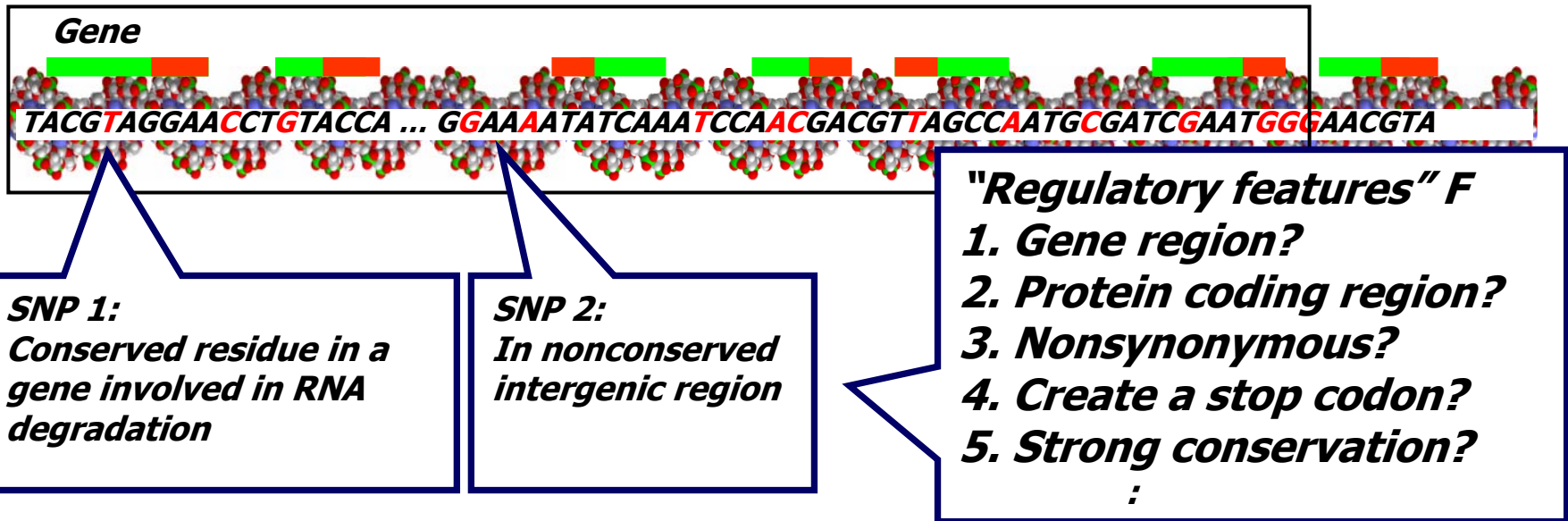




Motivation

- Not all SNPs are equally likely to be causal.

ChrXIV: 449,639-502,316

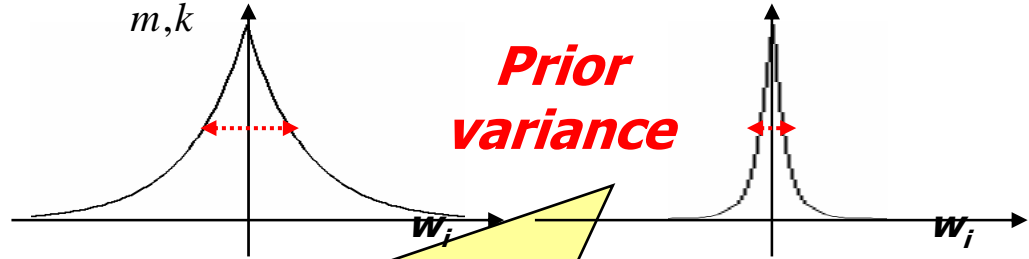


- ***Idea: Prioritize SNPs that have "good" regulatory features***
- ***But how do we weight different features?***



Bayesian L₁-Regularization

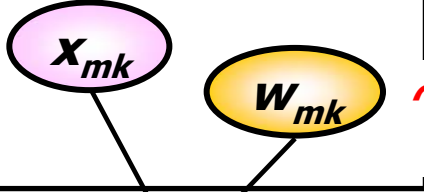
$$\sum_{m=1}^M \log P(y_m | \mathbf{x}_m, \mathbf{w}_m) - \sum_{m,k} C |w_{m,k}|$$



higher prior variance
 \Rightarrow **weight can more easily deviate from 0**
 \Rightarrow **regulator more likely to be selected**

Module m

Regulator k

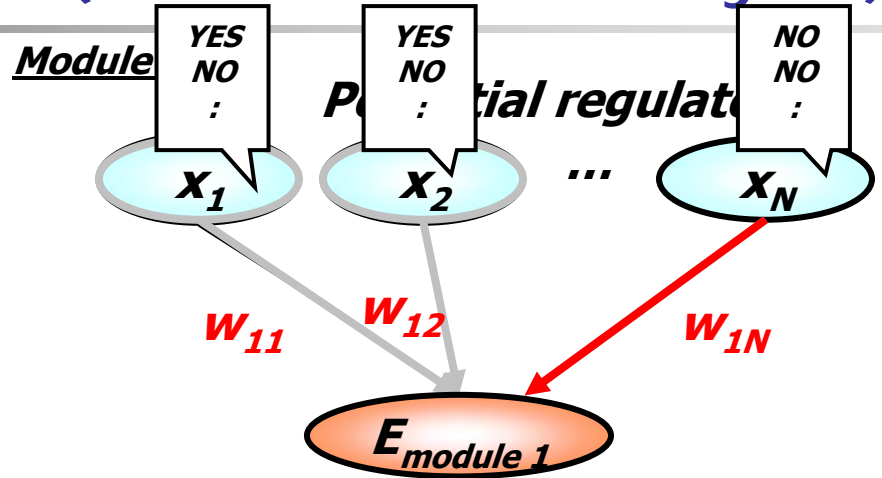
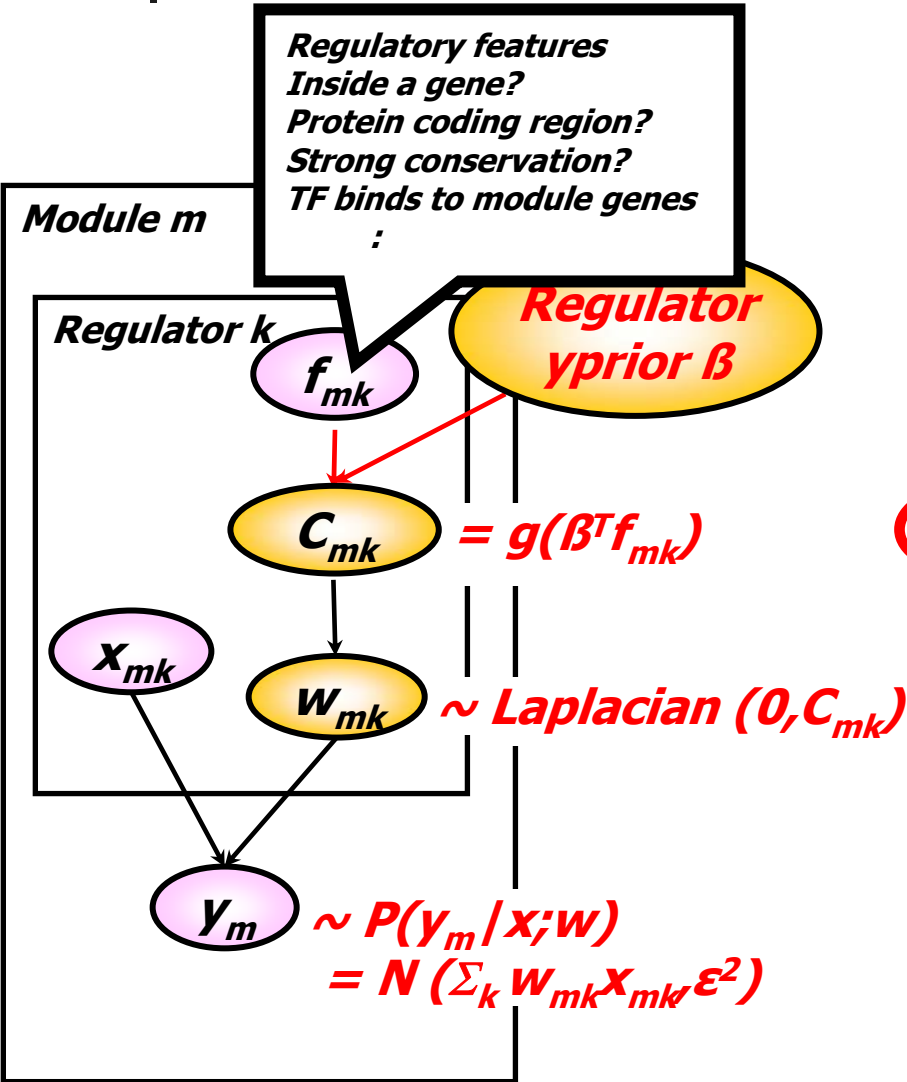


$\sim P(w)$
 $= \text{Laplacian}(0, C)$

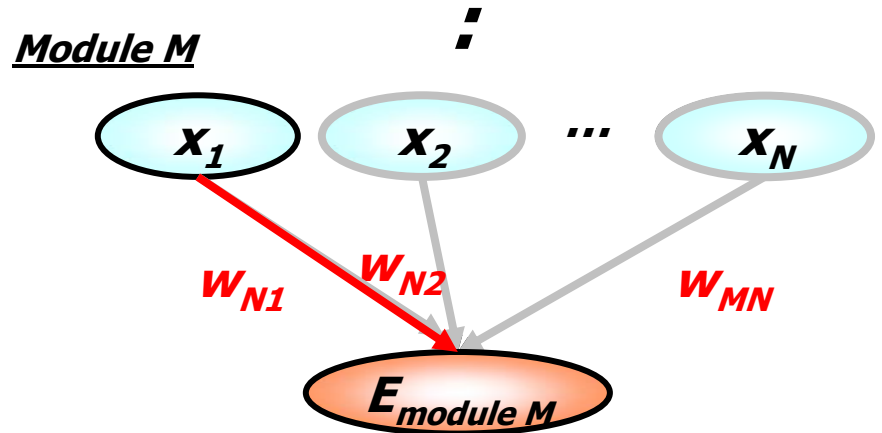
$\sim P(y_m | \mathbf{x}_m; \mathbf{w})$
 $= N(\sum_k w_{mk} x_{mk}, \epsilon^2)$



Metaprior Model (Hierarchical Bayes)



"Regulatory potential" =
 $\beta_1 \times$ Inside a gene? + $\beta_2 \times$ Protein coding region? + $\beta_3 \times$ Conserved? ...

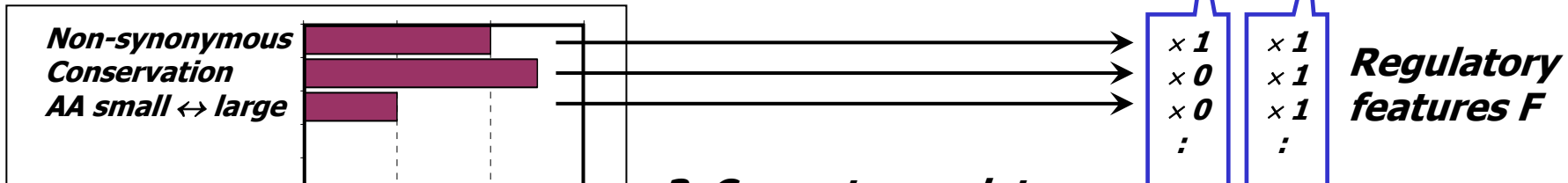




Metaprior Method

BY(lab) MVLIT ELVAQ VSDASKQLW

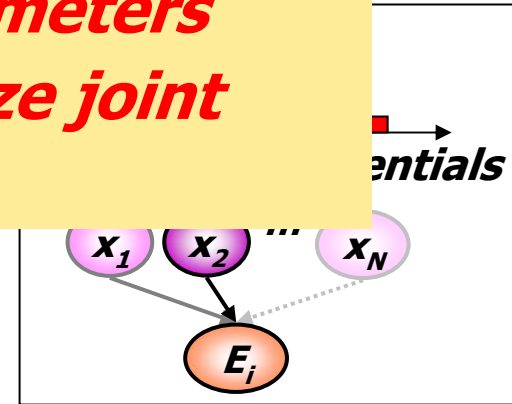
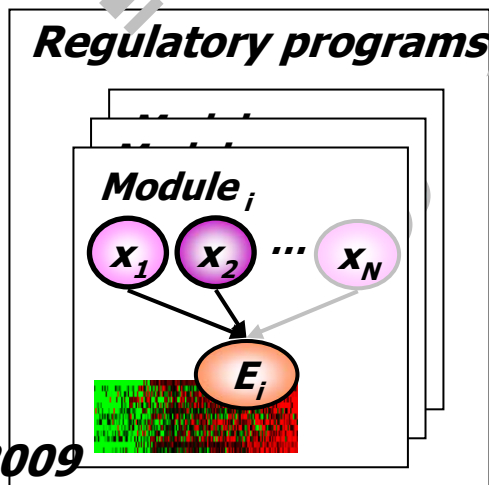
RM(wild) MVDIT ELVQ VSDASKQLW



Empirical hierarchical Bayes

- Use point estimate of model parameters
- Learn priors from data to maximize joint posterior

Maximize $P(E, \beta, W | X)$



Maximize $P(E, \beta, W | X)$



Transfer Learning

- What do regulatory potentials do?
 - They do **not** change selection of “strong” regulators – those where prediction of targets is clear
 - They only help disambiguate between weak ones
- Strong regulators help teach us what to look for in other regulators

***Transfer of knowledge
between different prediction tasks***

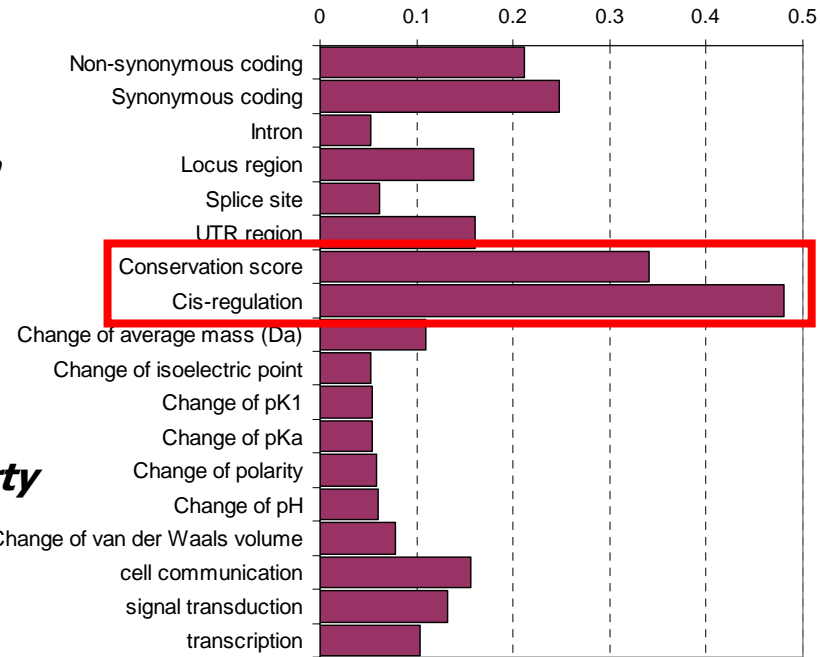
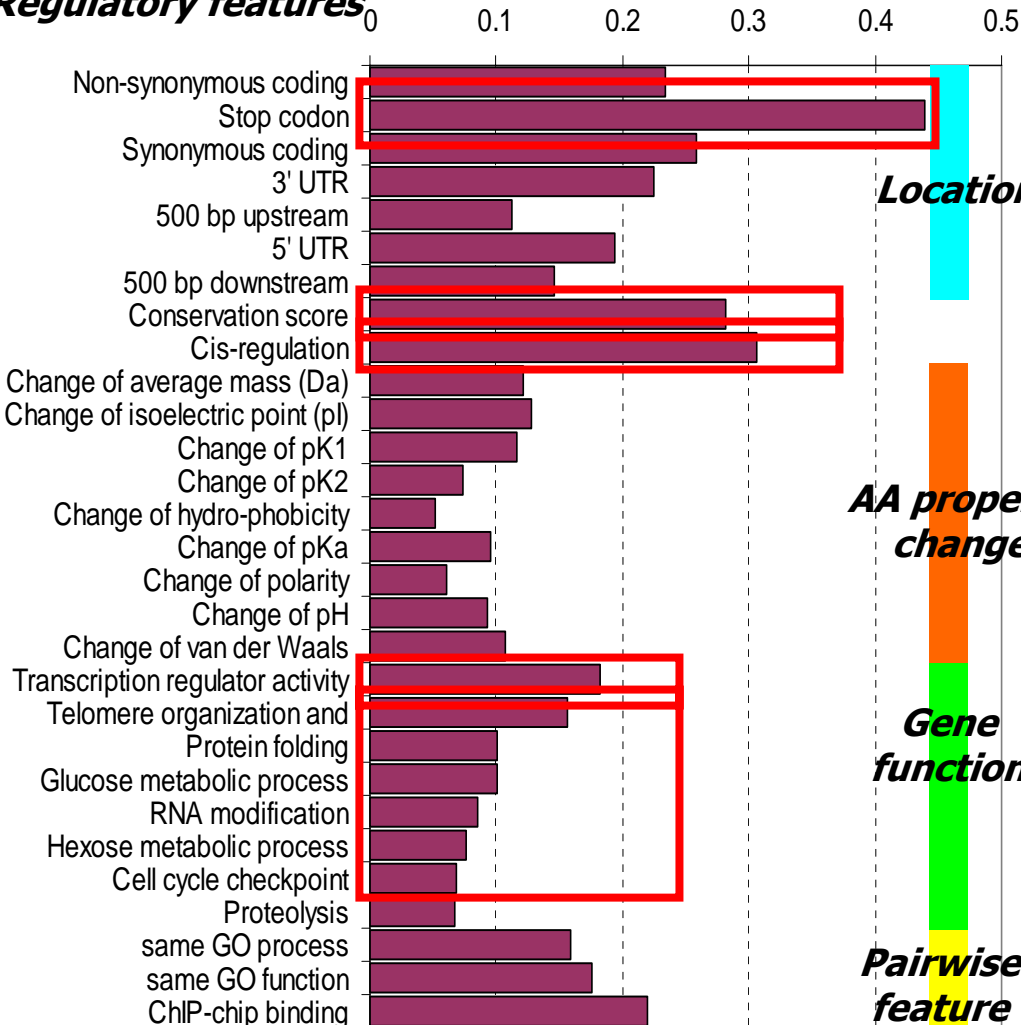


Learned regulatory weights

Yeast regulatory weights

Human regulatory weights

Regulatory features

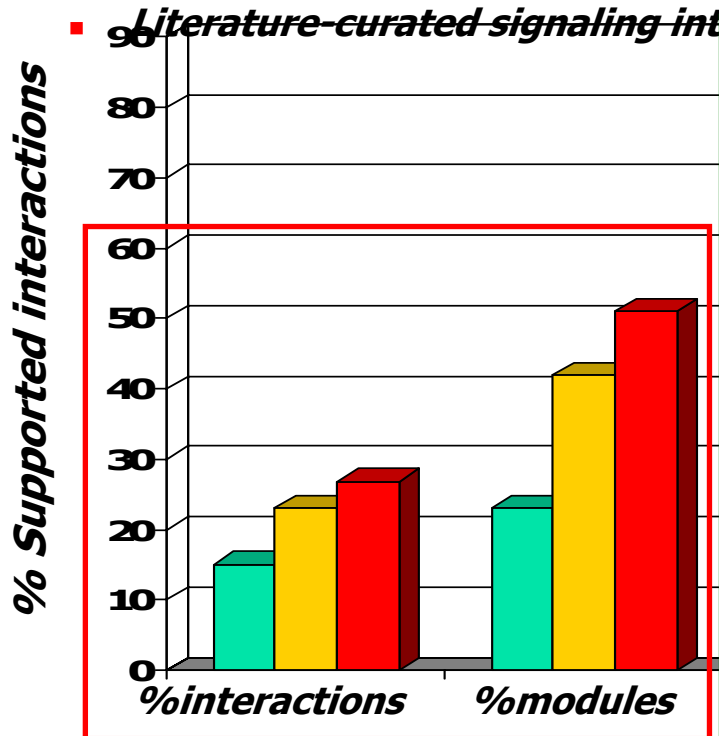




Biological evaluation I

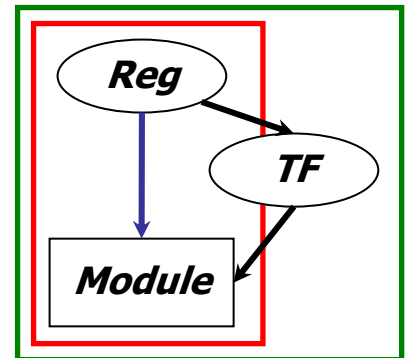
■ *How many predicted interactions have support in other data?*

- *Deletion/ over-expression microarrays [Hughes et al. 2000; Chua et al. 2006]*
- *ChIP-chip binding experiments [Harbison et al. 2004]*
- *Transcription factor binding sites [Maclsaac et al. 2006]*
- *mRNA binding pull-down experiments [Gerber et al. 2004]*
- *Literature-curated signaling int*



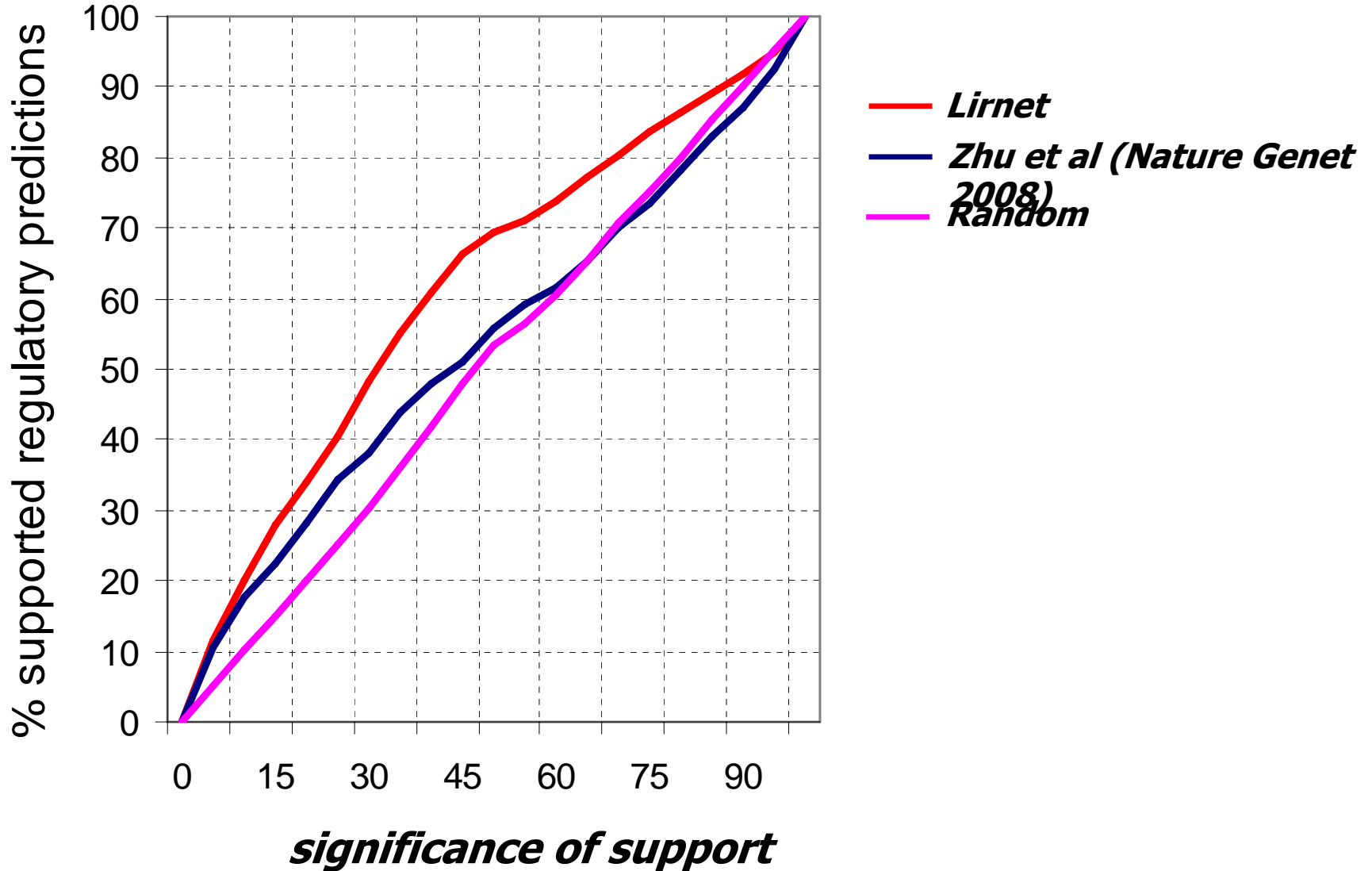
Lirnet without regulatory features

- Geronemo
- Flat Limet
- Limet





Biological Evaluation II



Lee et al., PLOS Genetics 2009



Preferential Chromosomal Regulation

8 validated regulators in 7 regions

14 validated regulators in 11 regions

- Finding causal regulators for 13 “chromosomal hot spots”

Region	Zhu et al [Nat Genet 08]	Lirnet (top 3 are considered)		
1	None	SEC18	RDH54	SPT7
2	TBS1, TOS1, ARA1, CSH1, SUP45, CNS1, AMN1	AMN1	CNS1	TOS1
3	None	TRS20	ABD1	PRP5
4	LEU2 , ILV6 , NFS1, CIT2, MATALPHA1	LEU2	PGS1	ILV6
5	MATALPHA1	MATALPHA1	MATALPHA2	RBK1
6	URA3	URA3	NPP2	PAC2
7	GPA1	STP2	GPA1	NEM1
8	HAP1	HAP1	NEJ1	GSY2
9	YRF1-4, YRF1-5, YLR464W	SIR3	HMG2	ECM7
10	None	ARG81	TAF13	CAC2
11	SAL1, TOP2	MKT1	TOP2	MSK1
12	PHM7	PHM7	ATG19	BRX1
13	None	ADE2	ORT1	CAT5



Current & Future Directions

- Understand mechanism by which individual genotype leads to changes in phenotype
 - Genotype & copy number changes (e.g., in cancer)
 - First step to personalized medicine
- Analysis and reconstruction of cellular pathways
- Understand immune response and how it is affected by aging
- ...



The Computer Science Inside

- **Computational Issues:** Huge graphical models require development of new algorithms
 - Convex optimization methods for learning network structure
 - Learning using MAP inference
 - Using combinatorial optimization within standard inference
- **Statistical issues:** Sparse data in high dimension
 - “Holistic models” to exploit correlations between different labels
 - Transfer learning between related problems
 - New algorithms for feature selection



Thank You!

<http://ai.stanford.edu/~koller/>

<http://dags.stanford.edu/>