

Light-Supervision of Structured Prediction Energy Networks

Andrew McCallum

**Pedram
Rooshenas**

Oregon PhD → UMass Postdoc



**Aishwarya
Kamath**

UMass MS



SPENs
[2016]

David Belanger

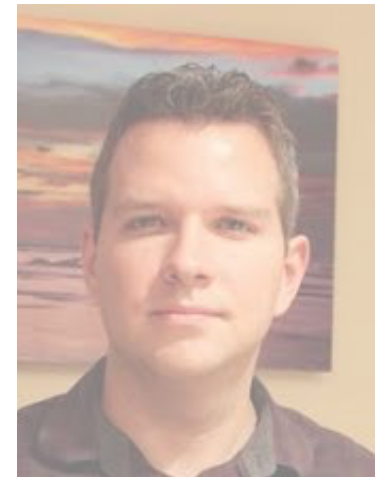
UMass PhD → Google Brain



Generalized Expectation
[Mann; Druck 2010-12]

Greg Druck

UMass PhD → Yummly



Light-Supervision

Prior Knowledge as *Generalized Expectation*

...induces extra structural dependencies...

Structured Prediction

Complex dependencies with *SPENs*

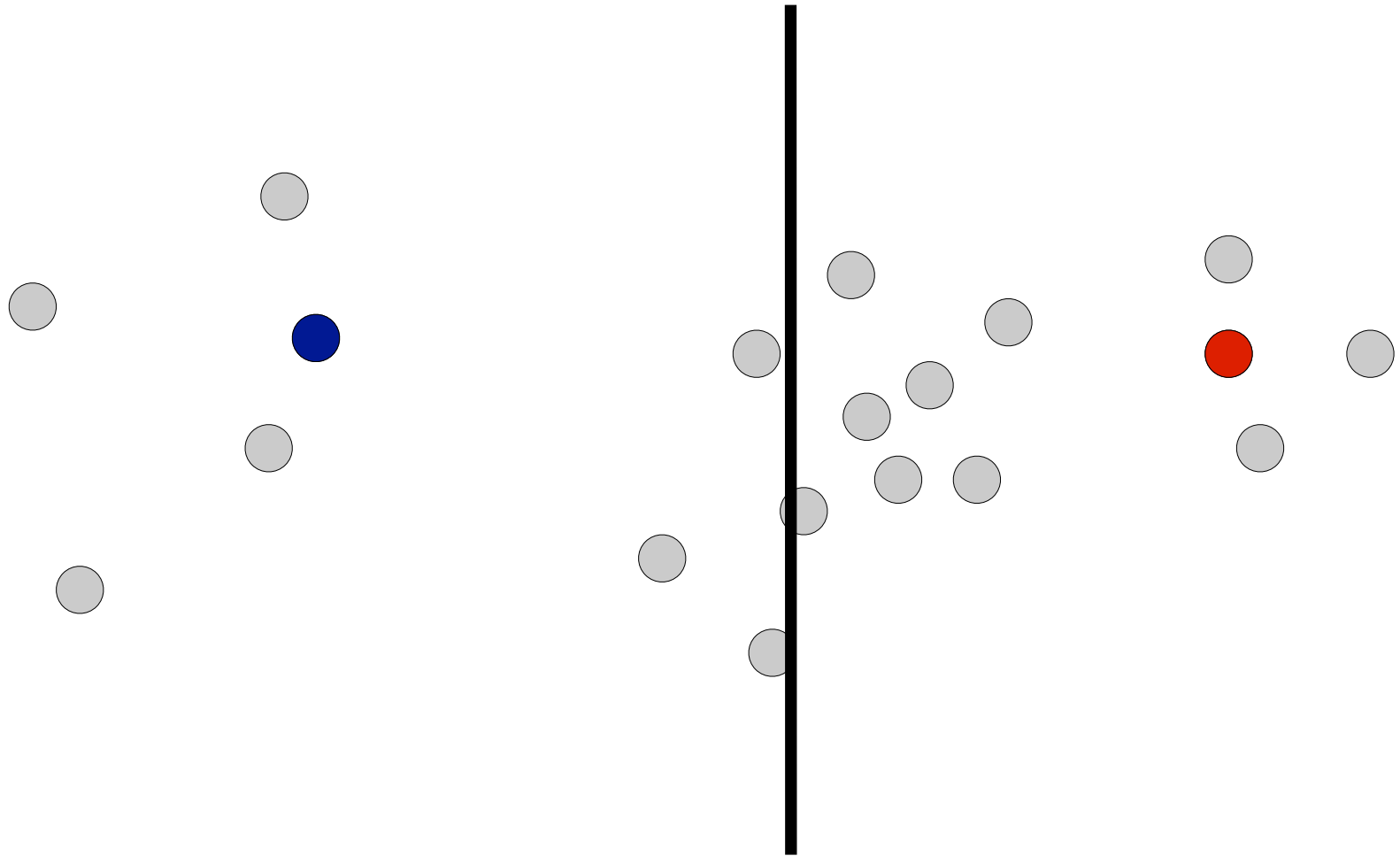
Chapter 1

Generalized Expectation

Learning from small labeled data

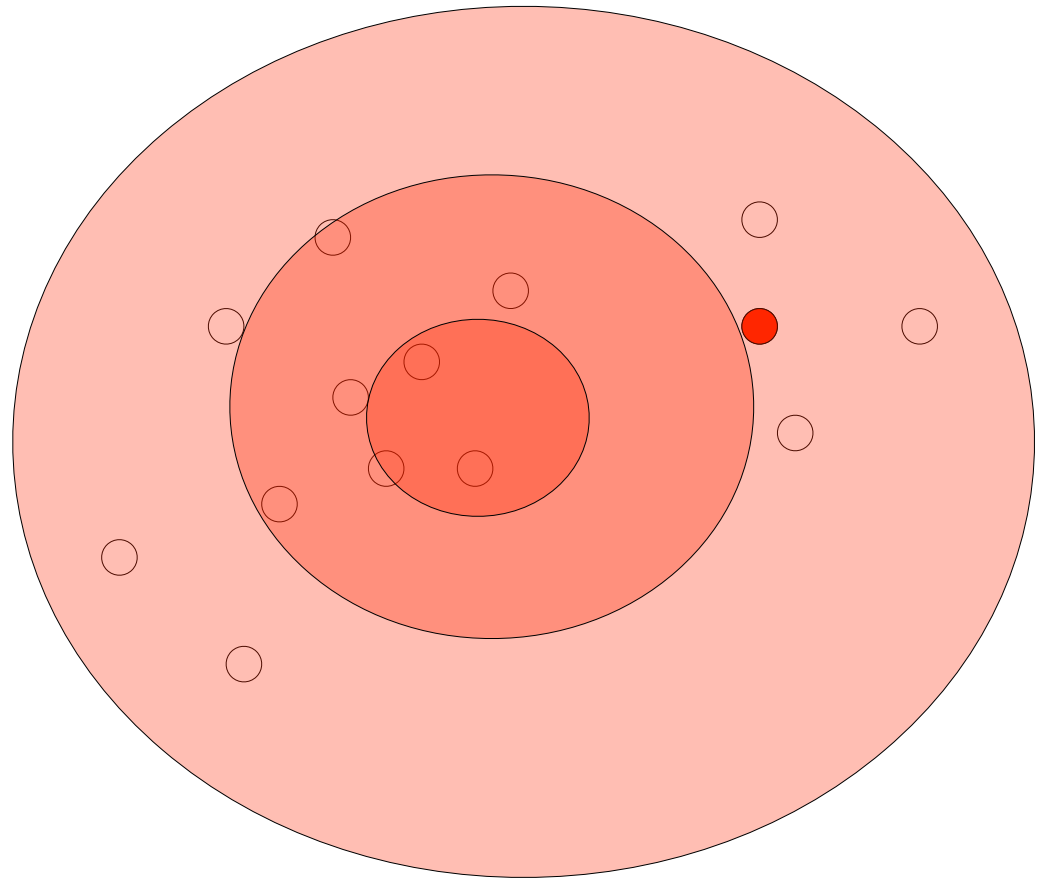
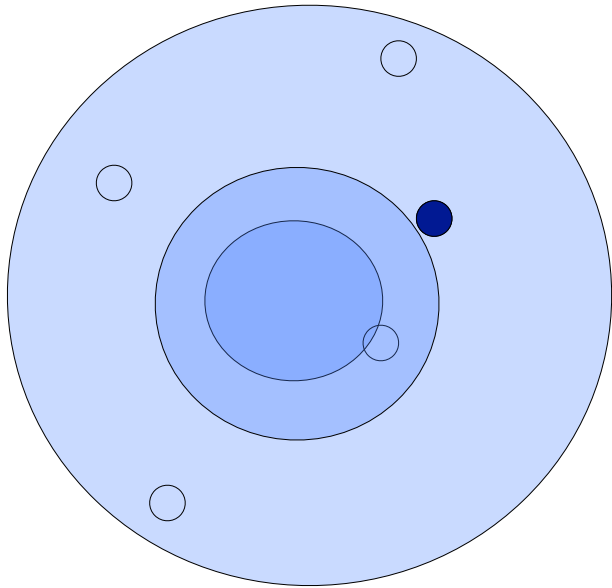


Leverage unlabeled data



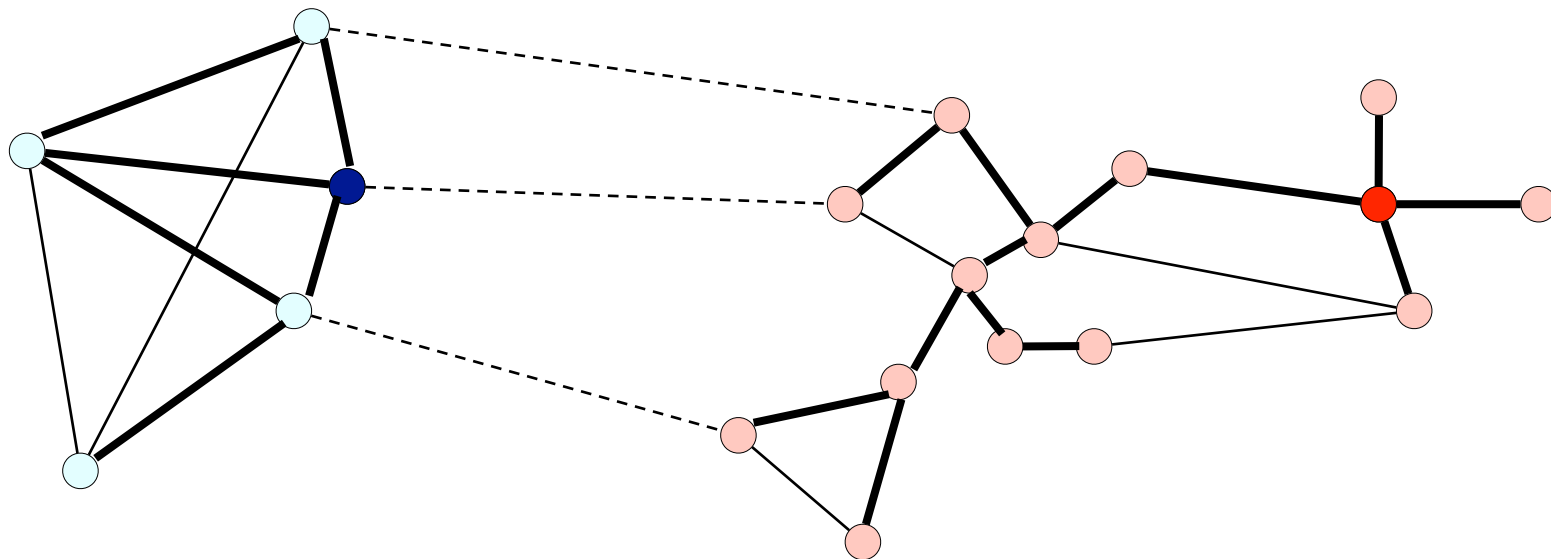
Family 1: Expectation Maximization

[Dempster, Laird, Rubin, 1977]



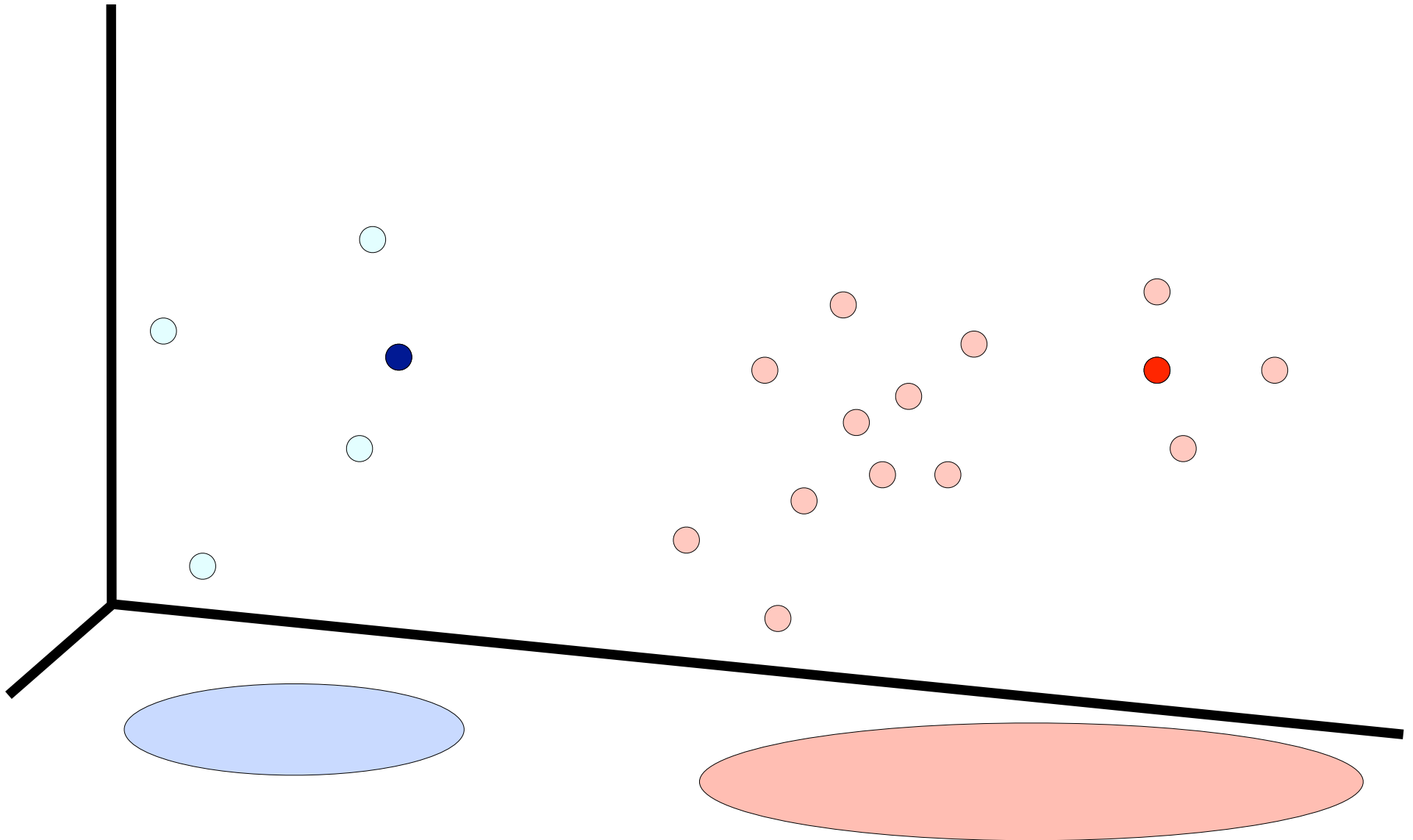
Family 2: Graph-Based Methods

[Szumner, Jaakkola, 2002] [Zhu, Ghahramani, 2002]



Family 3: Auxiliary-Task Methods

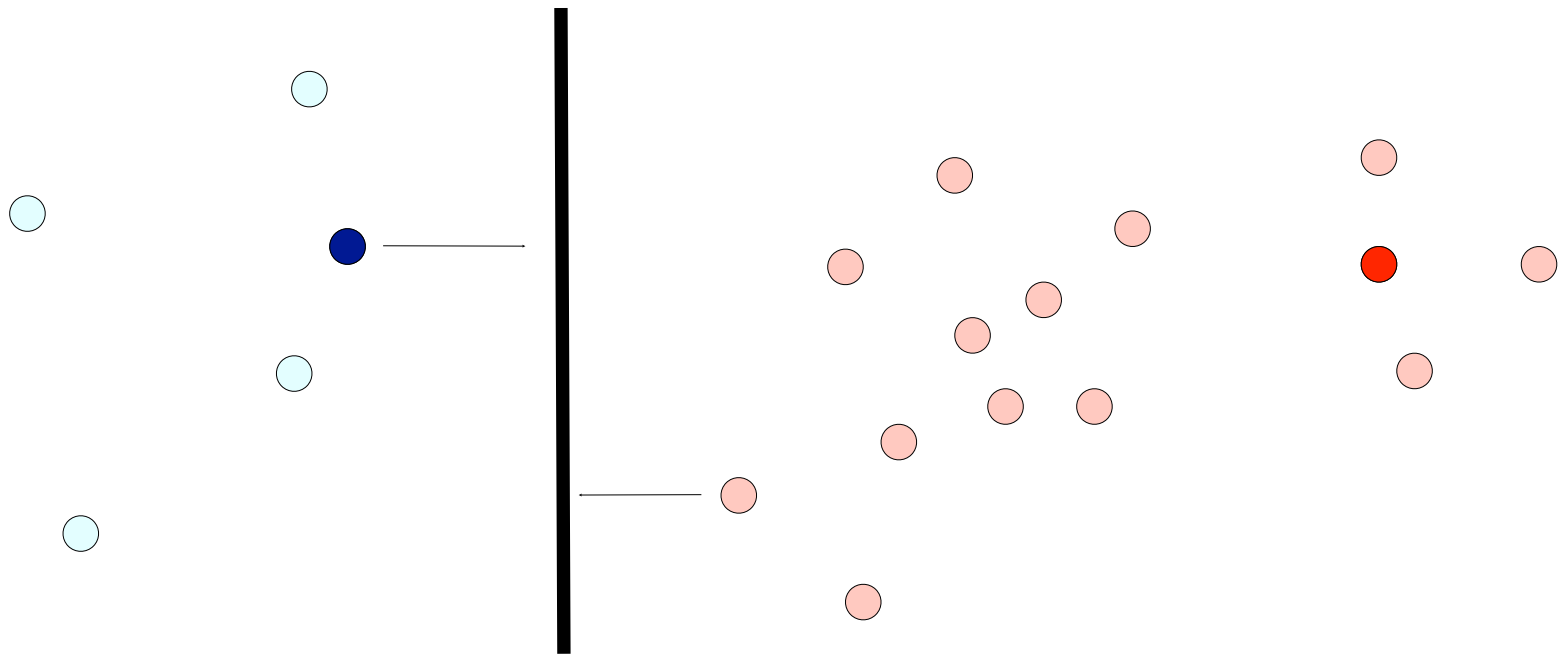
[Ando and Zhang, 2005]



Family 4: Boundary in Sparse Region

Transductive SVMs [Joachims, 1999]: Sparsity measured by margin

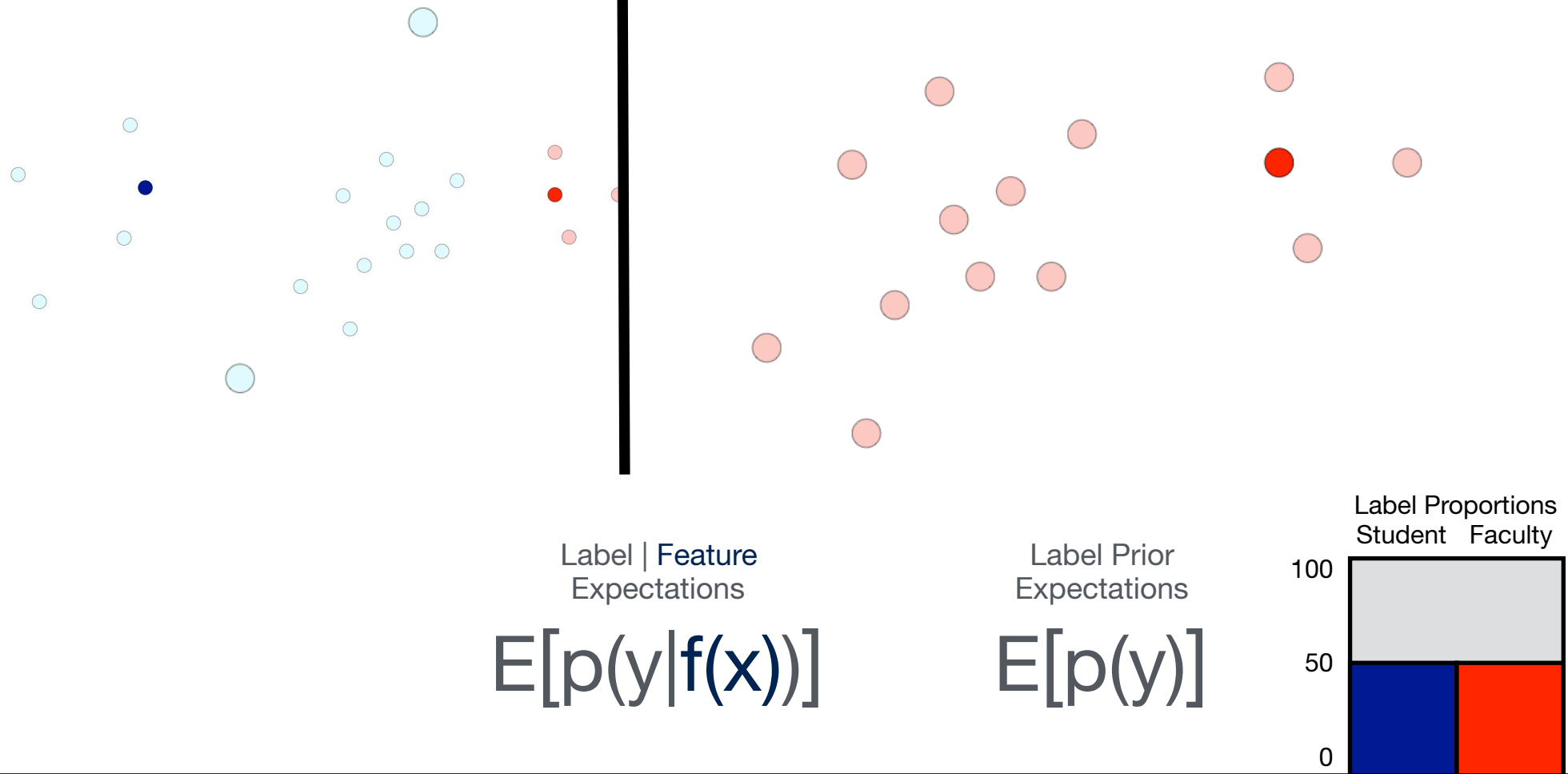
Entropy Regularization [Grandvalet & Bengio, 2005]: minimize label entropy



Family 54 GB General Library Expansion Reigedra

[Mann, McGarrigle 2018] / [McGarraigs, Mann 2019] / [Sparsity, Deak & Callaghan 2012]

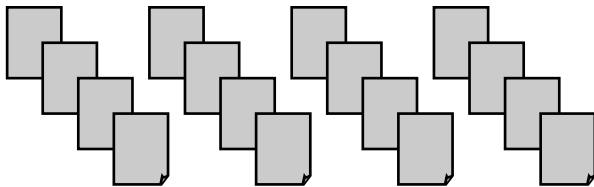
Entropy Regularization [Zadachynski, Leventhal, & Bengio, 2005]: minimize label entropy
best solution?



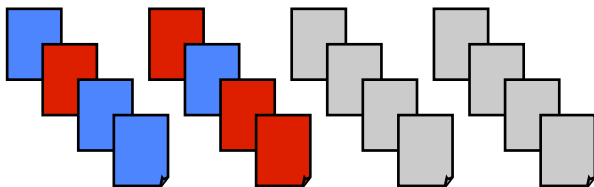
Expectations on Labels | Features

Classifying *Baseball* versus *Hockey*

Traditional

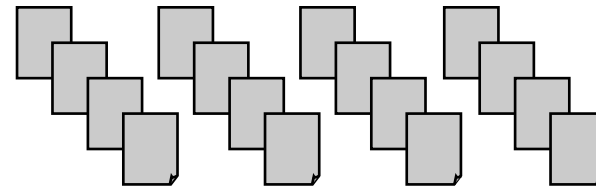


Human
Labeling
Effort



(Semi-)Supervised Training via
Maximum Likelihood

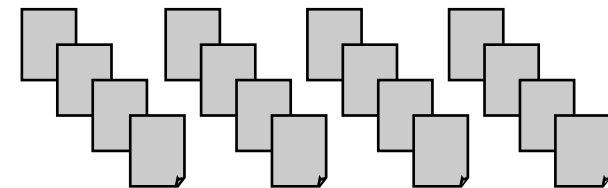
Generalized Expectation



Brainstorm
a few
Keywords

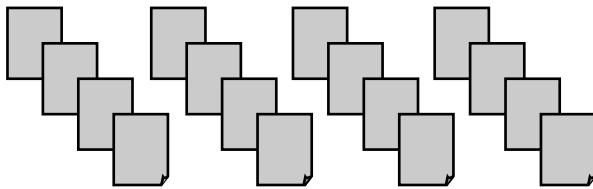


$p(\text{HOCKEY} \mid \text{"puck"}) = .9$



Semi-Supervised Training via
Generalized Expectation

Labeling Features



~1000 unlabeled examples

features labeled . . .

hockey
baseball
HR
Mets

goal
Buffalo
Leafs
puck
Lemieux

Toronto Maple
Leafs

ball
Oilers
Sox
Pens
runs

Edmonton Oilers

Pittsburgh
Penguins

batting
base
NHL
Bruins
Penguins

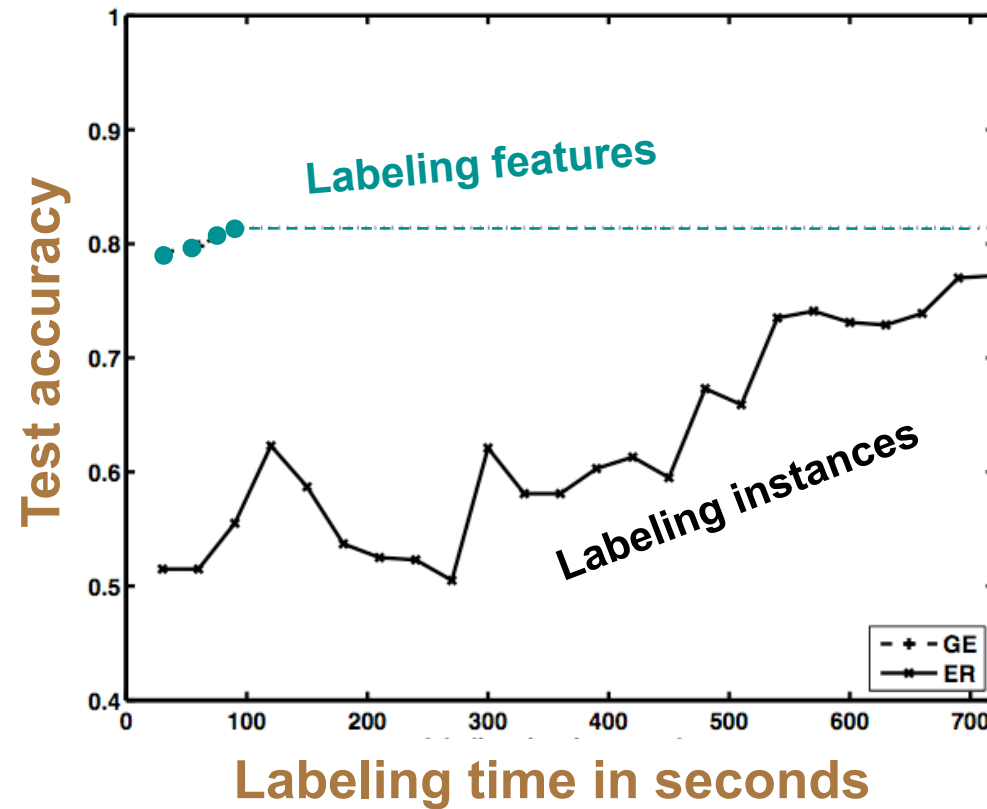
Accuracy **85%**

92%

94.5%

96%

Accuracy per Human Effort



Prior Knowledge

Feature labels from humans

baseball/hockey classification

baseball	hockey
hit	puck
braves	goal
runs	nhl

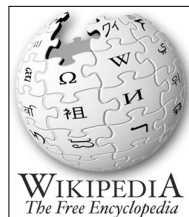
many other sources

resources on the web

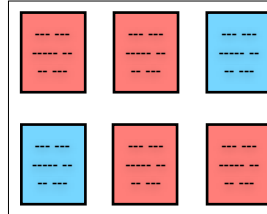
DBLP Record 'conf/aaai/MaierTOJ10'

BibTeX

```
@inproceedings{DBLP:conf/aaai/MaierTOJ10,  
  author = {Marc Maier and  
    Brian Taylor and  
    Russein Okuy and  
    David Jensen},  
  title = {Learning Causal Models of Relational Domains},  
  booktitle = {AAAI},  
  year = {2010},  
  url = {http://www.aaai.org/ocs/index.php/AAAI/AAAI10/paper/view/1919},  
  crossref = {DBLP:conf/aaai/2010},  
  bibsource = {DBLP, http://dblp.uni-trier.de}
```

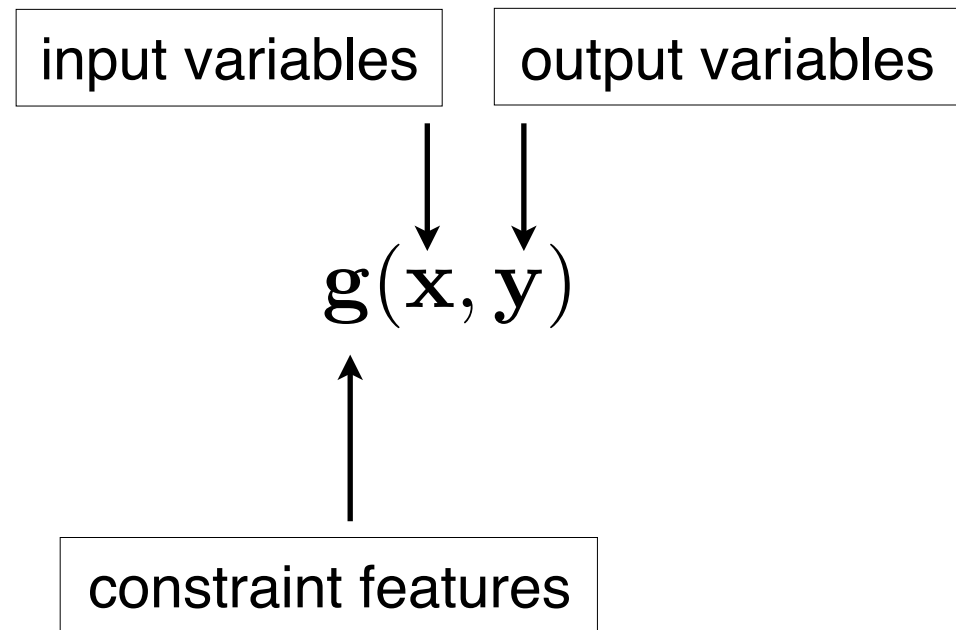


data from related tasks



W.H. Enright. Improving the efficiency of matrix operations in the numerical solution of stiff ordinary differential equations. *ACM Trans. Math. Softw.*, 4(2), 127-136, June 1978.

Generalized Expectation (GE)



returns 1 if \mathbf{x} contains “hit” and \mathbf{y} is **baseball**

Generalized Expectation (GE)

assume general CRF [Lafferty et al. 01]

$$p(\mathbf{y}|\mathbf{x}; \theta) = \frac{1}{Z_{\theta, \mathbf{x}}} \exp(\theta^\top \mathbf{f}(\mathbf{x}, \mathbf{y}))$$

$$\mathbb{E}_{p(\mathbf{y}|\mathbf{x}; \theta)} [\mathbf{g}(\mathbf{x}, \mathbf{y})]$$

model distribution

model features

model probability of **baseball** if \mathbf{x} contains “hit”

Generalized Expectation (GE)

$$E_{\tilde{p}(\mathbf{x})} [E_{p(\mathbf{y}|\mathbf{x};\theta)} [\mathbf{g}(\mathbf{x}, \mathbf{y})]]$$



empirical distribution

(can be defined as)
model's probability that
documents that contain "hit" are labeled **baseball**

Generalized Expectation (GE)

(soft) expectation constraint

$$S(\mathbb{E}_{\tilde{p}(\mathbf{x})} [\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)} [\mathbf{g}(\mathbf{x}, \mathbf{y})]])$$



score function

larger score if model expectation matches prior knowledge

Generalized Expectation (GE)

Objective Function

$$\mathcal{O}(\theta) = S(\mathbb{E}_{\tilde{p}(\mathbf{x})}[\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)}[\mathbf{g}(\mathbf{x}, \mathbf{y})]]) + r(\theta)$$

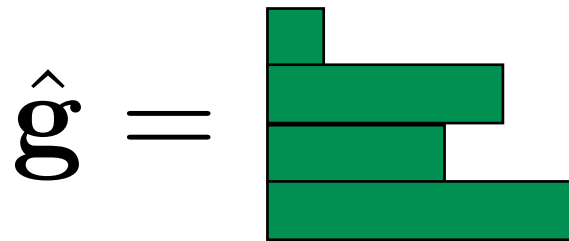
regularization



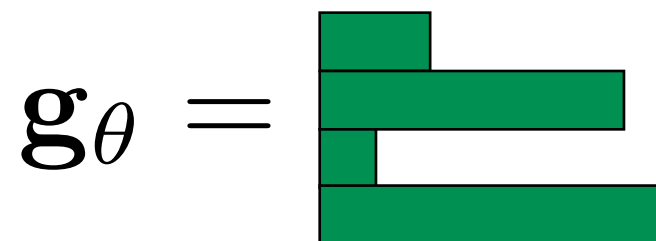
GE Score Functions

$$\mathcal{O}(\theta) = S(\mathbb{E}_{\tilde{p}(\mathbf{x})}[\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)}[\mathbf{g}(\mathbf{x}, \mathbf{y})]]) + r(\theta)$$

target expectations



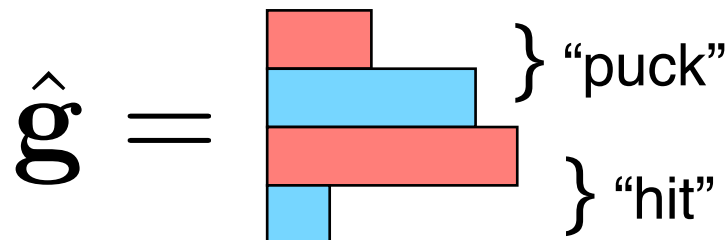
model expectations



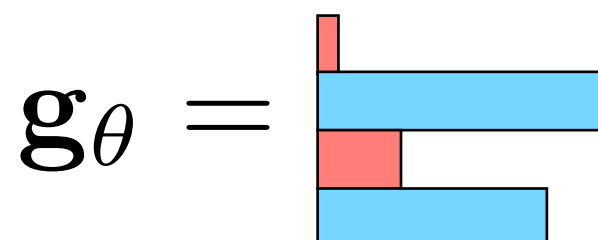
squared error:

$$S_{l_2^2}(\theta) = -\left\| \hat{\mathbf{g}} - \mathbf{g}_\theta \right\|_2^2$$

target expectations



model expectations



KL divergence:

$$S_{KL}(\theta) = -\sum_q \hat{\mathbf{g}}_q \log \frac{\hat{\mathbf{g}}_q}{\mathbf{g}_{\theta,q}}$$

Estimating Parameters with GE

$$\mathcal{O}(\theta) = S(\mathbb{E}_{\tilde{p}(\mathbf{x})}[\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)}[\mathbf{g}(\mathbf{x}, \mathbf{y})]]) + r(\theta)$$

violation term:

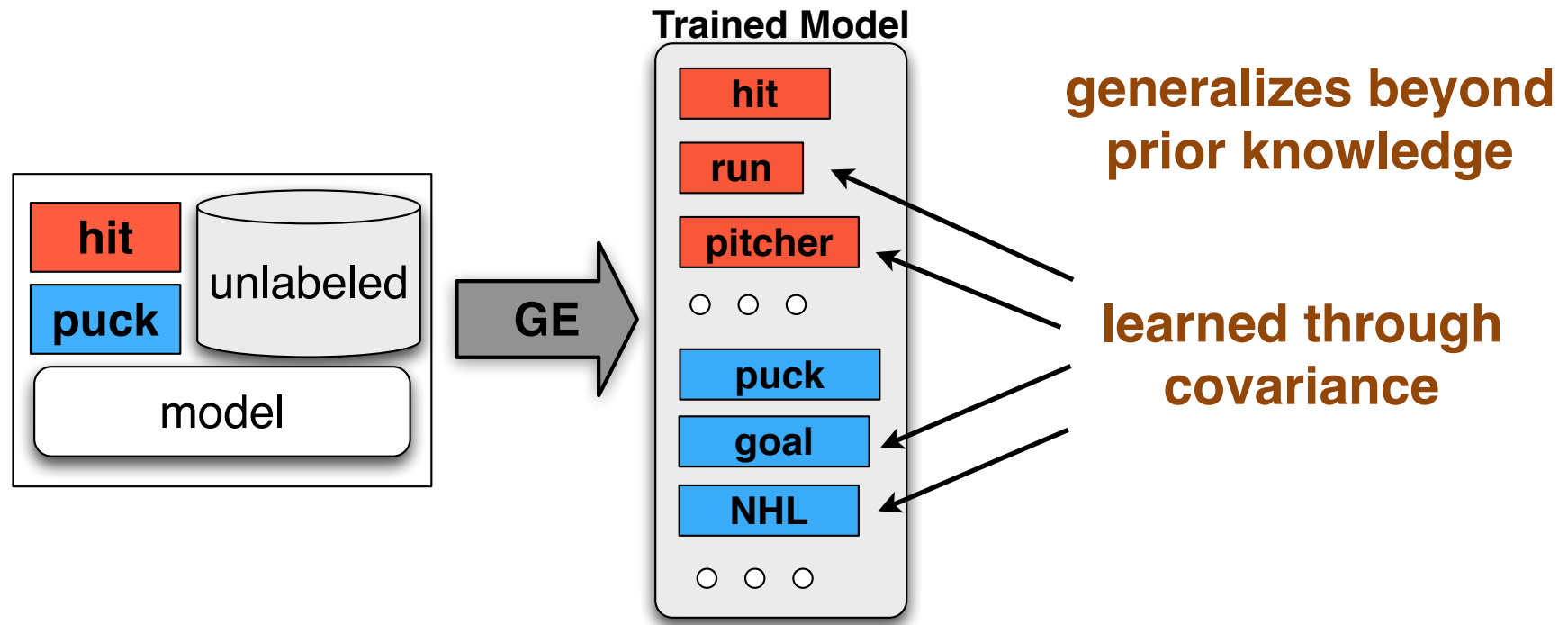
$$\text{KL: } v_i = \frac{\hat{g}_i}{g_{\theta i}}$$

$$\text{sq. error: } v_i = -2(\hat{g}_i - g_{\theta i})$$

$$\nabla_{\theta} \mathcal{O}(\theta) = \mathbf{v}^{\top} \left(\mathbb{E}_{\tilde{p}(\mathbf{x})} \left[\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)} [\mathbf{g}(\mathbf{x}, \mathbf{y}) \mathbf{f}(\mathbf{x}, \mathbf{y})^{\top}] \right. \right. \\ \left. \left. - \mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)} [\mathbf{g}(\mathbf{x}, \mathbf{y})] \mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)} [\mathbf{f}(\mathbf{x}, \mathbf{y})^{\top}] \right] \right) + \nabla_{\theta} r(\theta)$$

estimated covariance between model and constraint features

Learning About Unconstrained Features



Generalized Expectation criteria

Easy communication with domain experts

- Inject domain knowledge into parameter estimation
- Like “informative prior”...
- ...but rather than the “language of parameters”
(difficult for humans to understand)
- ...use the “language of expectations”
(natural for humans)

IID Prediction

“classification” e.g. *logistic regression*

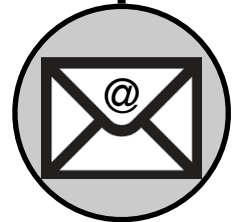
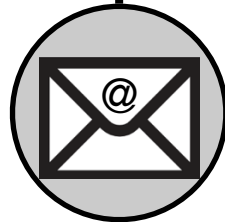
Example: Spam Filtering



Predicted
Y



Observed
X



Structured Prediction

e.g. “sequence labeling” Chinese Word Segmentation

$$\mathcal{O}(\theta) = S(\mathbb{E}_{\tilde{p}(\mathbf{x})}[\mathbb{E}_{p(\mathbf{y}|\mathbf{x};\theta)}[\mathbf{g}(\mathbf{x}, \mathbf{y})]]) + r(\theta)$$

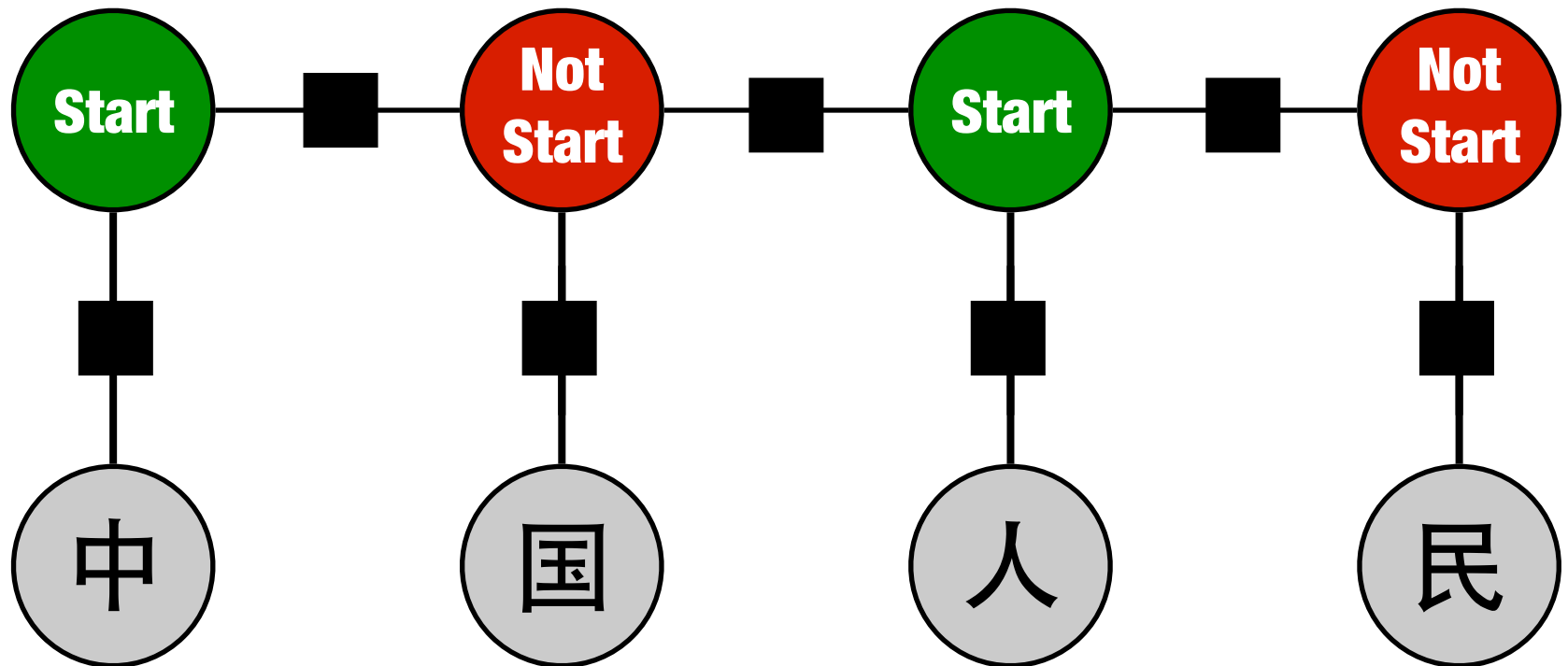
Linear-chain CRF

GE
Gradient

$$\mathbf{v}^\top \sum_{\mathbf{y}} \sum_i \sum_j p(y_{i-1}, y_i, y_j | \mathbf{x}; \theta) \mathbf{g}(\mathbf{x}, y_j, j) \mathbf{f}(\mathbf{x}, y_{i-1}, y_i, i)^\top$$

marginal over three, non-consecutive positions

Y



X

C h i n e s e P e o p l e

手勝只的面是，總統表
對州契州尼戰總均
號三瑞亞穆挑黨前
頭在金治羅大和目
尼倫而喬。一共和
穆托，得選的共人
羅桑選，贏初臨其參選

Natural Expectations lead to **Difficult Training Inference**

“AUTHOR field should be contiguous, only appearing once.”

AUTHOR Anna Popescu (2004), “Interactive Clustering,”
AUTHOR
EDITOR Wei Li (Ed.), Learning Handbook, Athos Press,
EDITOR
LOCATION Souroti.

$p(y_{i-1}, y_i, y_j, y_k)$

The downfall of GE.

Chapter 2

A framework providing easier inference for complex dependencies?

Structured Prediction Energy Networks

Deep Learning
+
Structured Prediction

Structured Prediction

“classification” e.g. logistic regression

Example: Spam Filtering



Predicted

Y

$= \operatorname{argmin}_Y$

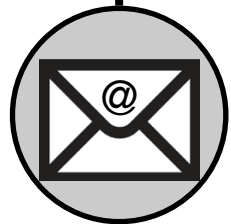
$E(Y;X) = \Sigma$

Observed

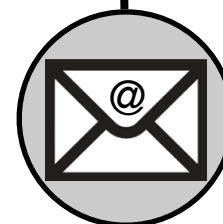
X



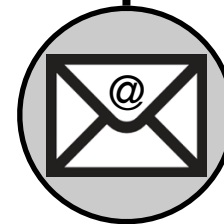
Factor



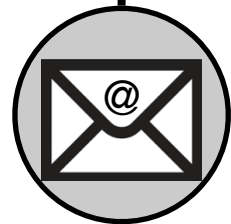
Factor



Factor



Factor

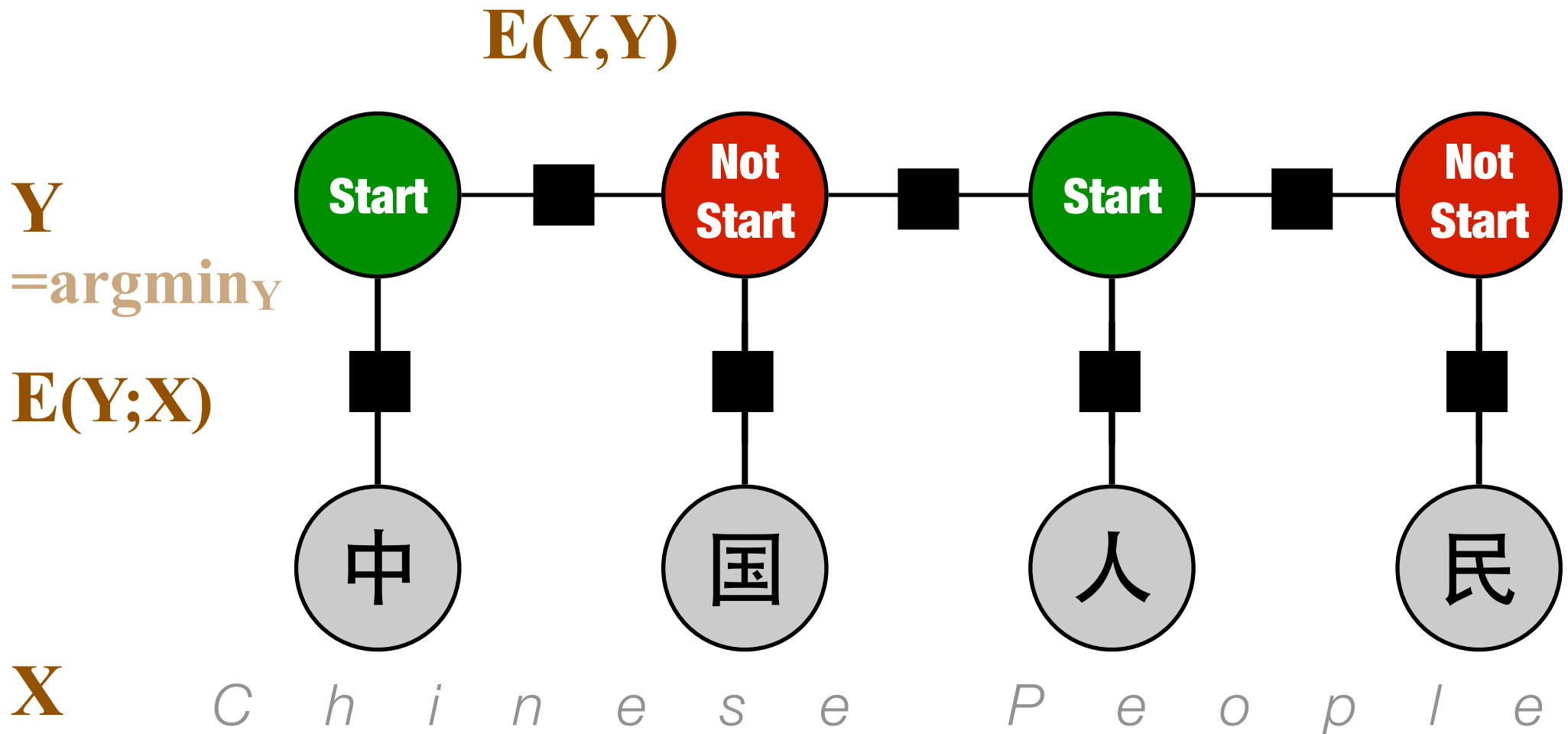


Structured Prediction

e.g. “sequence labeling”

Example: Chinese Word Segmentation

手勝只的面是，總統表
對州契州尼戰總均
號三瑞亞穆挑黨前
頭在金治羅大和目
尼倫而喬。一共和
穆托，得選的他選
羅桑選贏初臨其參

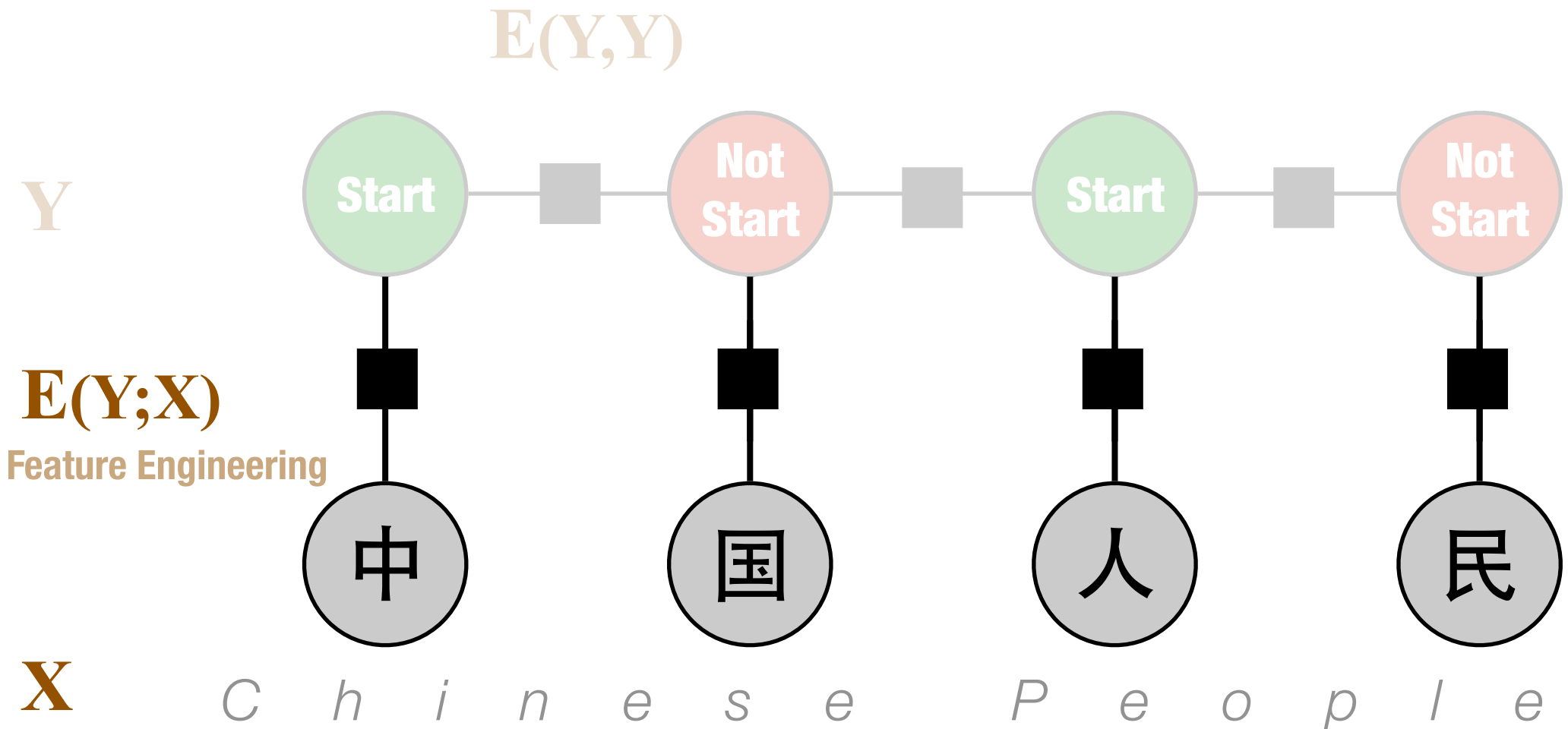


Structured Prediction

e.g. "sequence labeling"

Example: Chinese Word Segmentation

手勝只的面是，
對州契州尼戰總
號三瑞亞穆挑黨
頭在金治羅大和
尼倫而喬。一共
穆托，得選的。人
羅桑選，贏初臨其
參

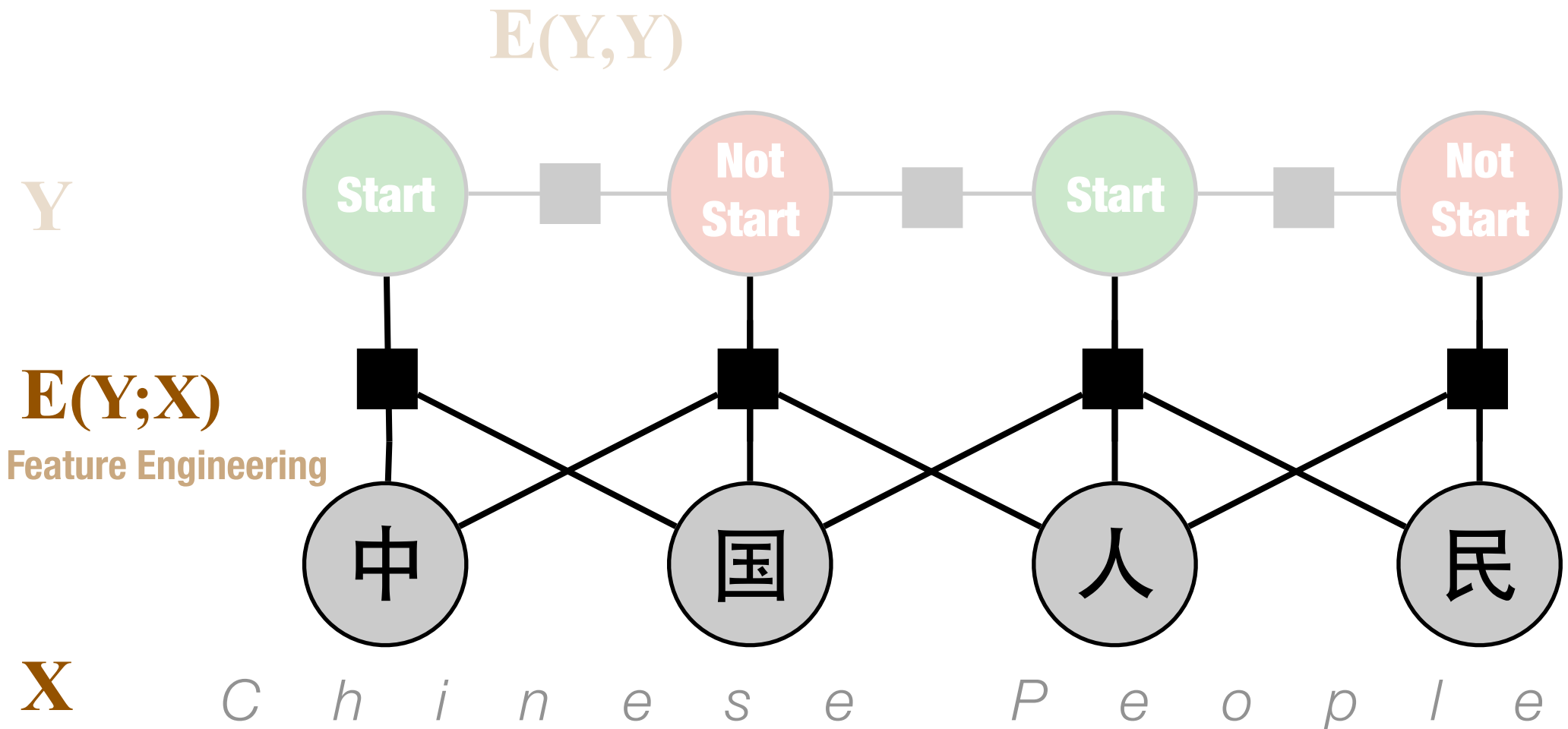


Structured Prediction

e.g. “sequence labeling”

Example: Chinese Word Segmentation

手勝只的面是，
對州契州尼戰總
號三瑞亞穆挑黨
頭在金治羅大和
尼倫而喬。一共
穆托，得選的。人
羅桑選，贏初臨其
參

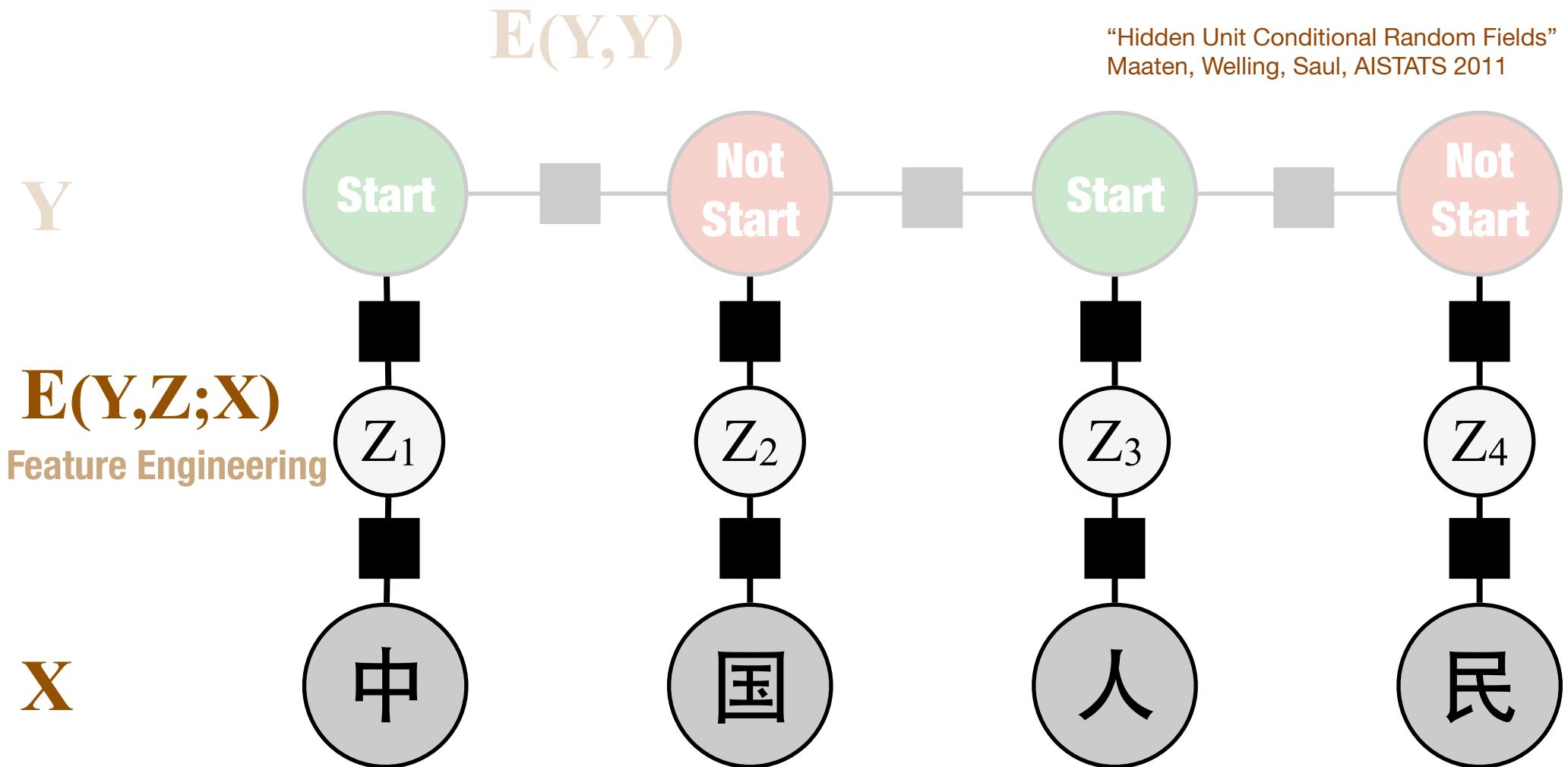


Structured Prediction

e.g. “sequence labeling”

Example: Chinese Word Segmentation

對手勝只的面是總統表
對三契州尼挑戰均
號瑞亞穆大和黨前
頭在金治羅挑目
尼倫而喬選的一共人
穆托，得選的他人
羅桑選，贏初臨其參選



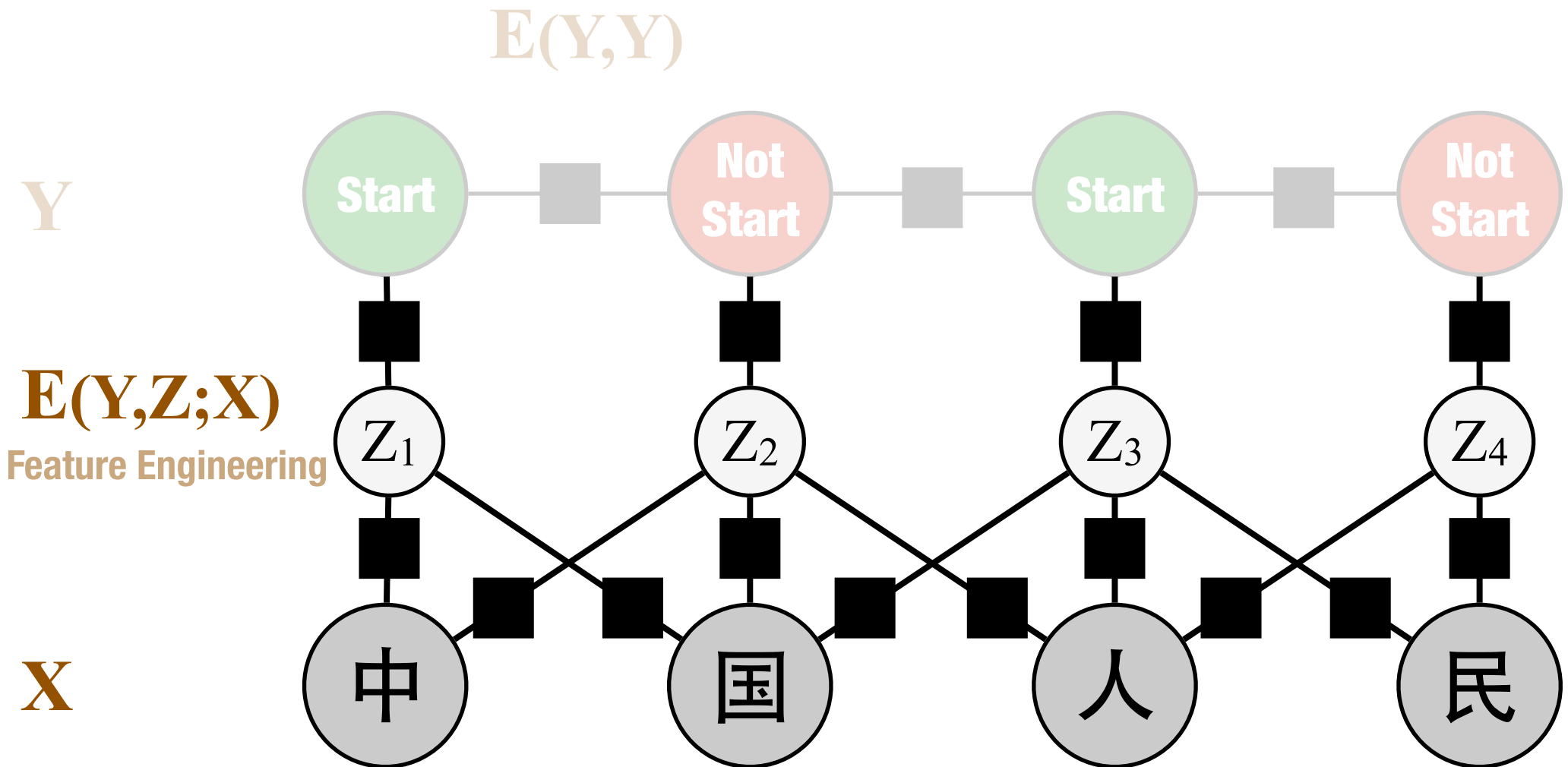
Chinese People

Structured Prediction

e.g. "sequence labeling"

Example: Chinese Word Segmentation

手勝只的面是，
對契州尼戰總
號三瑞亞挑黨
頭在金治羅大
尼倫而喬。一
穆托，得選的
羅桑選贏初臨
其參



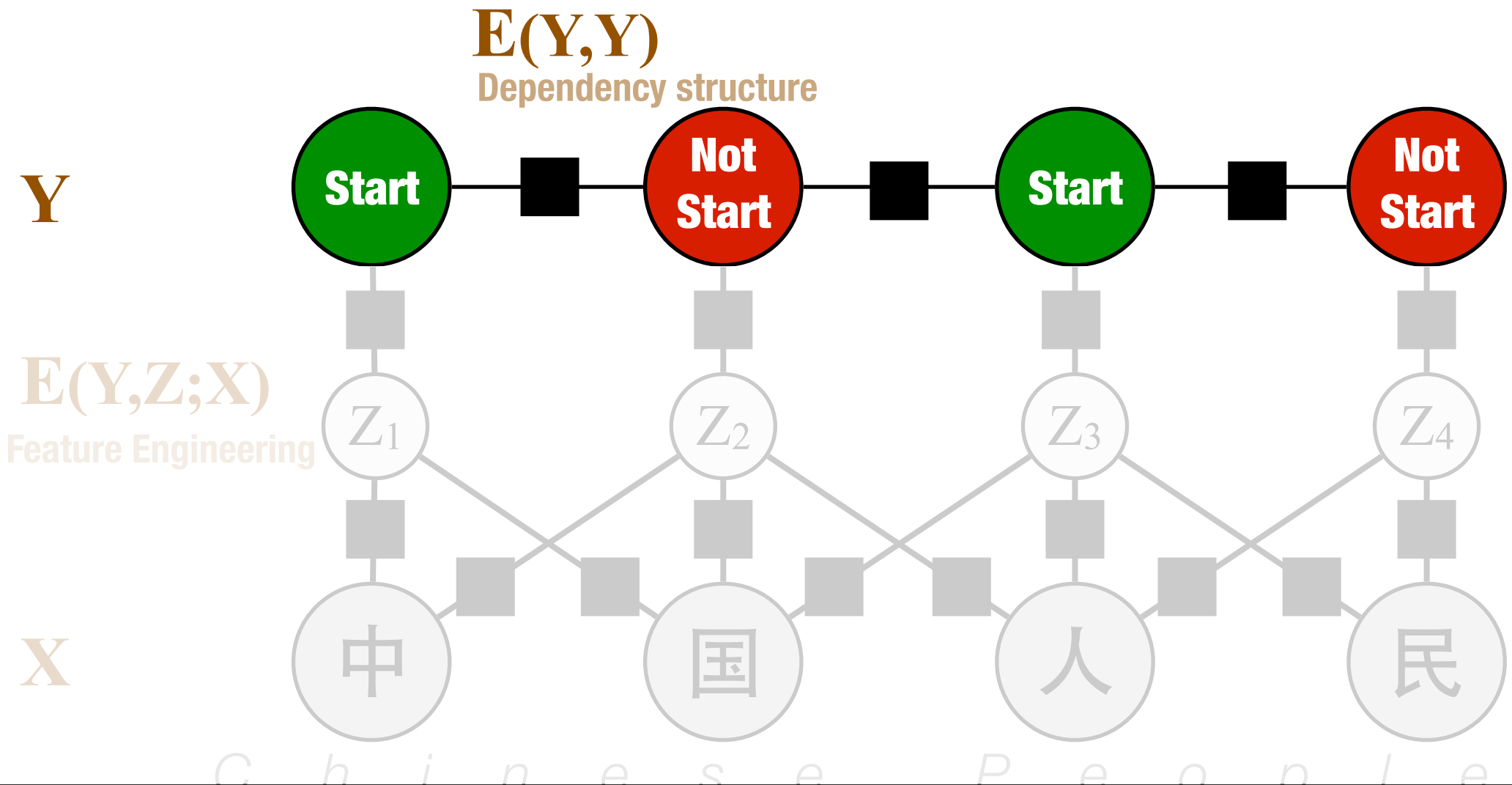
Chinese People

Structured Prediction

e.g. “sequence labeling”

Example: Chinese Word Segmentation

手勝只的面是，
對契州尼戰總
號三瑞亞挑黨
頭在金治羅大和
尼倫而喬。一
穆托，得選的
羅桑選，贏初
其參



Structured Prediction

e.g. "sequence labeling"

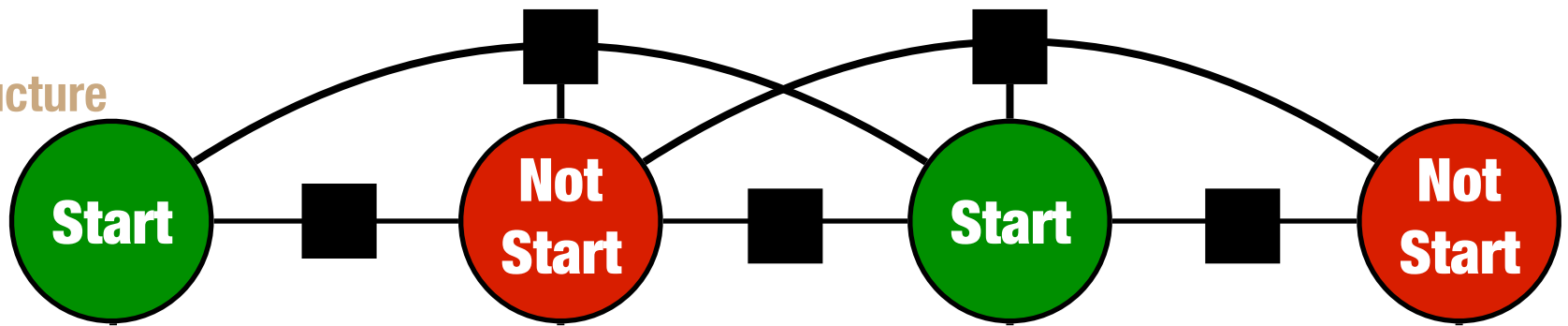
Example: Chinese Word Segmentation

手勝只的面是，
對契州尼戰總
號三瑞亞挑黨
頭在金治羅大
尼倫而喬。一
穆托，得選的
羅桑選贏初臨
其參選

$E(Y,Y)$

Dependency structure

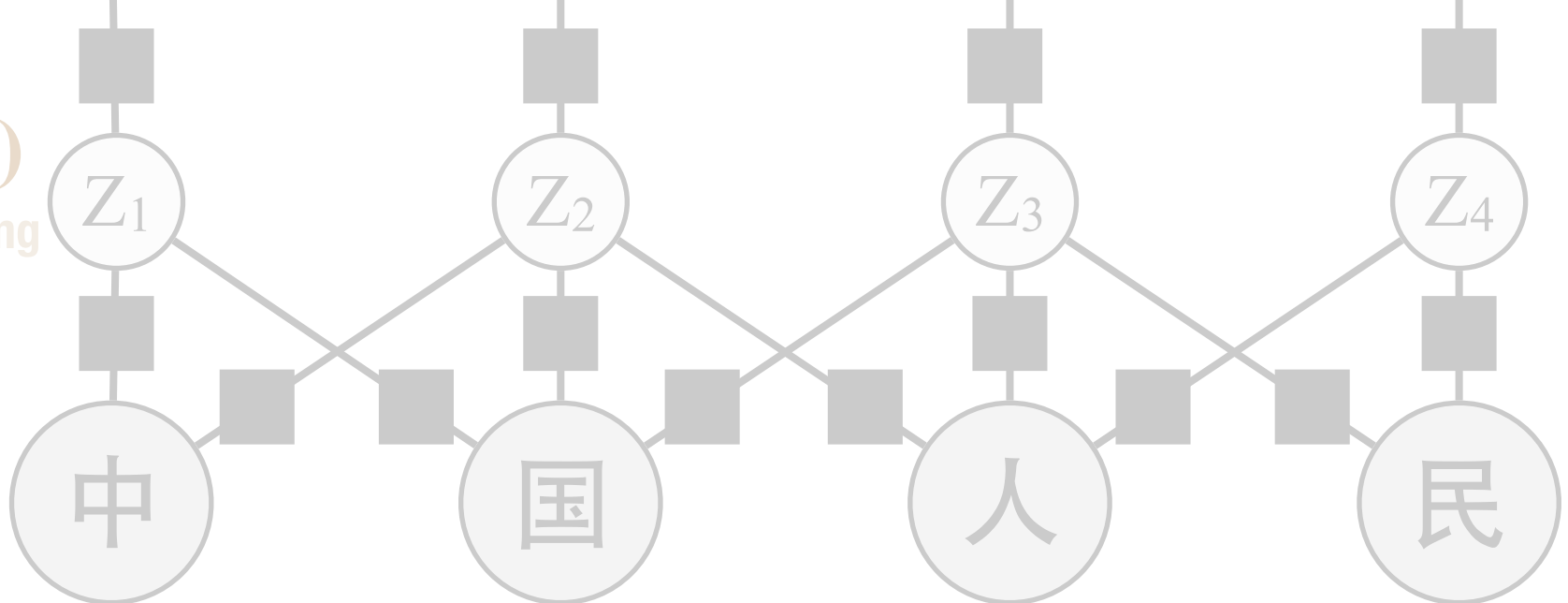
Y



$E(X,Z...,Y)$

Feature Engineering

X

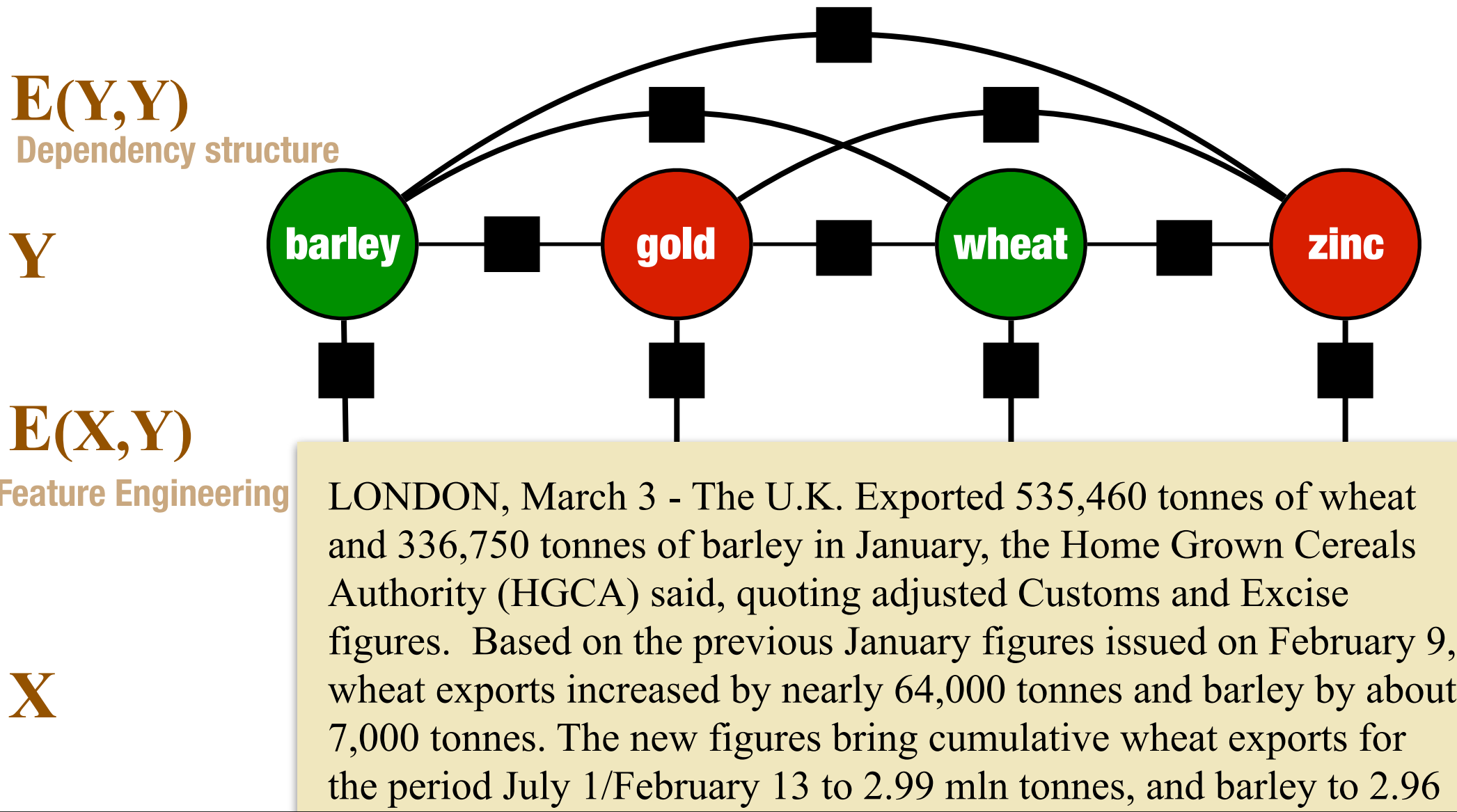


C h i n e s e P e o p l e

Structured Prediction

e.g. “multi-label classification”

Example: Multi-label Document Classification



Structured Prediction

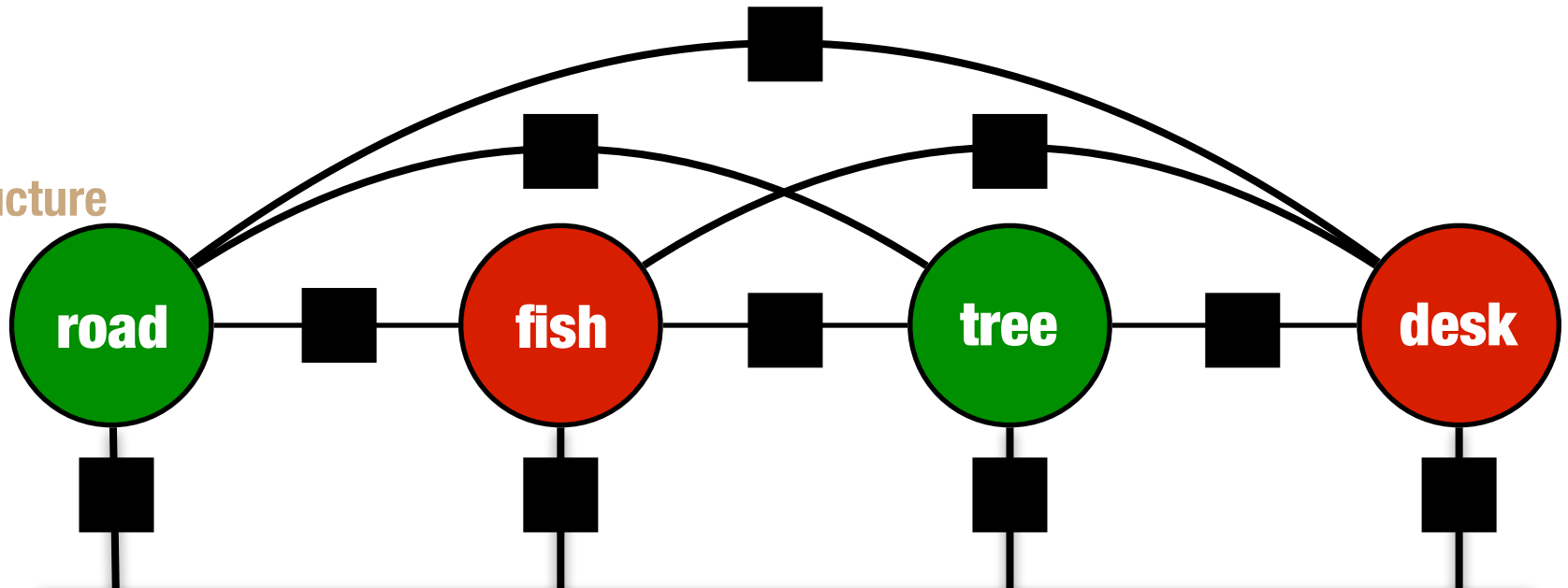
e.g. “multi-label classification”

Example: Multi-label Image Classification

$E(Y,Y)$

Dependency structure

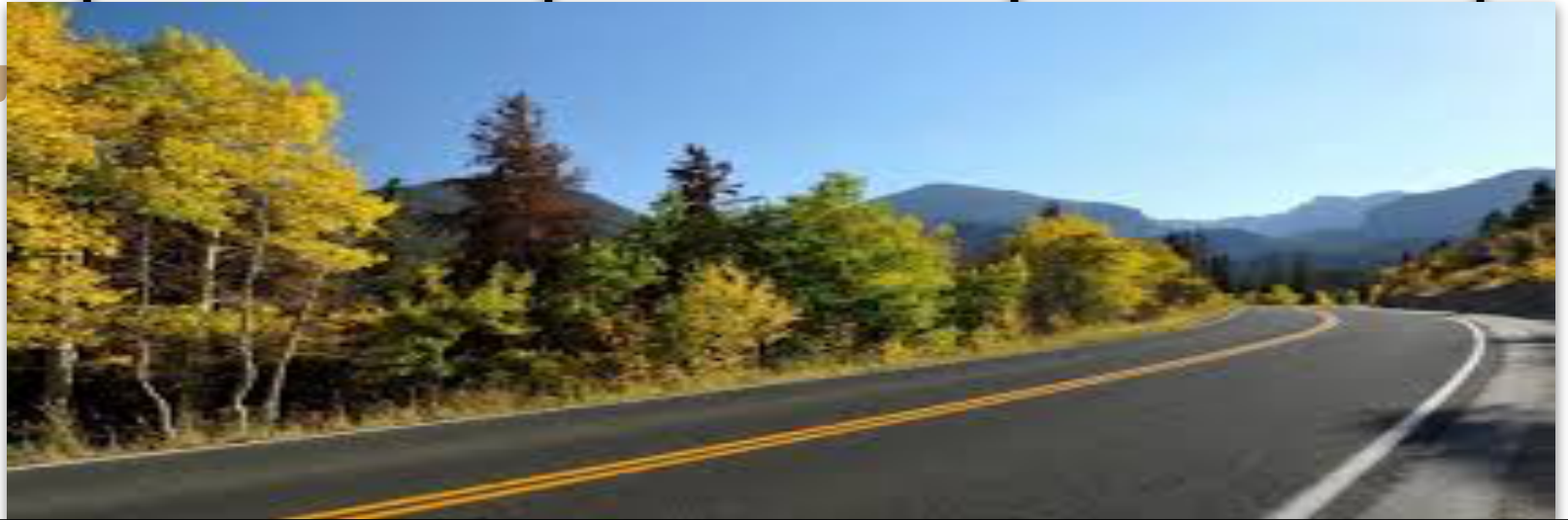
Y



$E(X,Y)$

Feature Engineering

X



Structured Prediction

Example:
Scene Understanding

$E(Y,Y)$

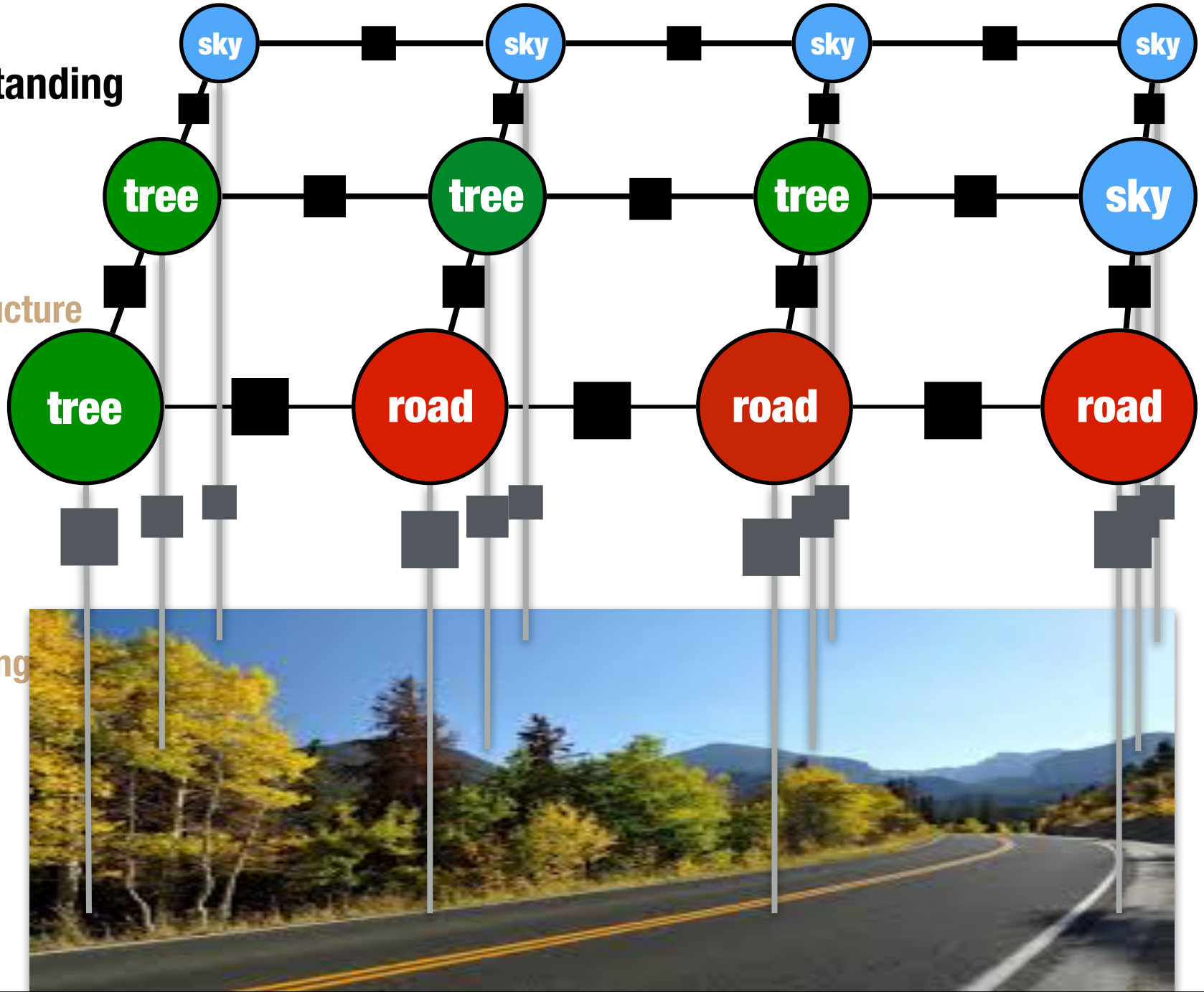
Dependency structure

Y

$E(X,Y)$

Feature Engineering

X



Structured Prediction

Example:
Scene Understanding

$E(Y, Y)$

Dependency structure

Y

$E(X, Y)$

Feature Engineering

X

- **Expressivity** of dependencies
- **Parsimony** of parameterization
- **Tractability** of inference



Structured Prediction *Sampling Inference*

Example:
Scene Understanding

$E(Y, Y)$

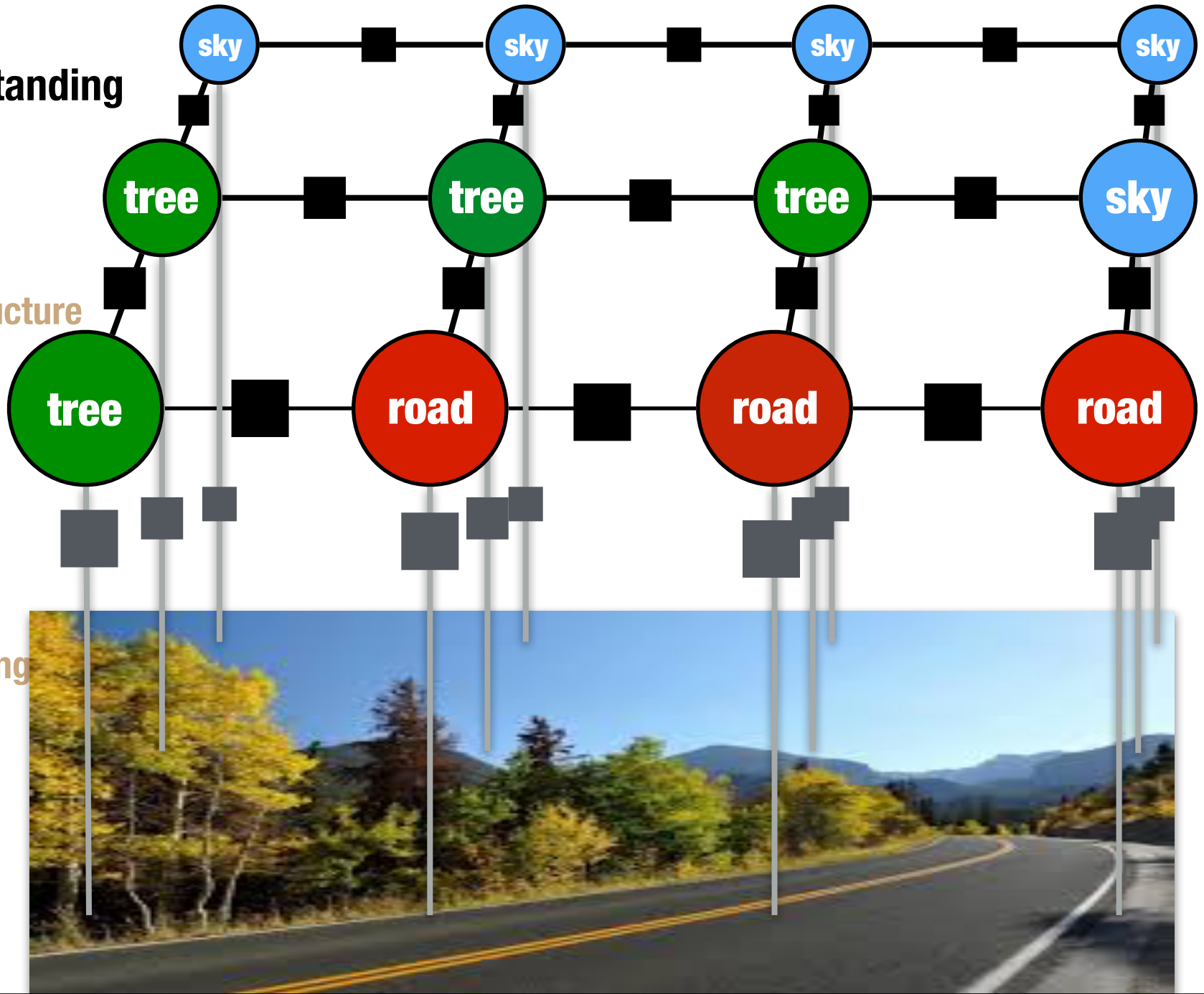
Dependency structure

Y

$E(X, Y)$

Feature Engineering

X



Structured Prediction *Sampling Inference*

Example:
Scene Understanding

$E(Y, Y)$

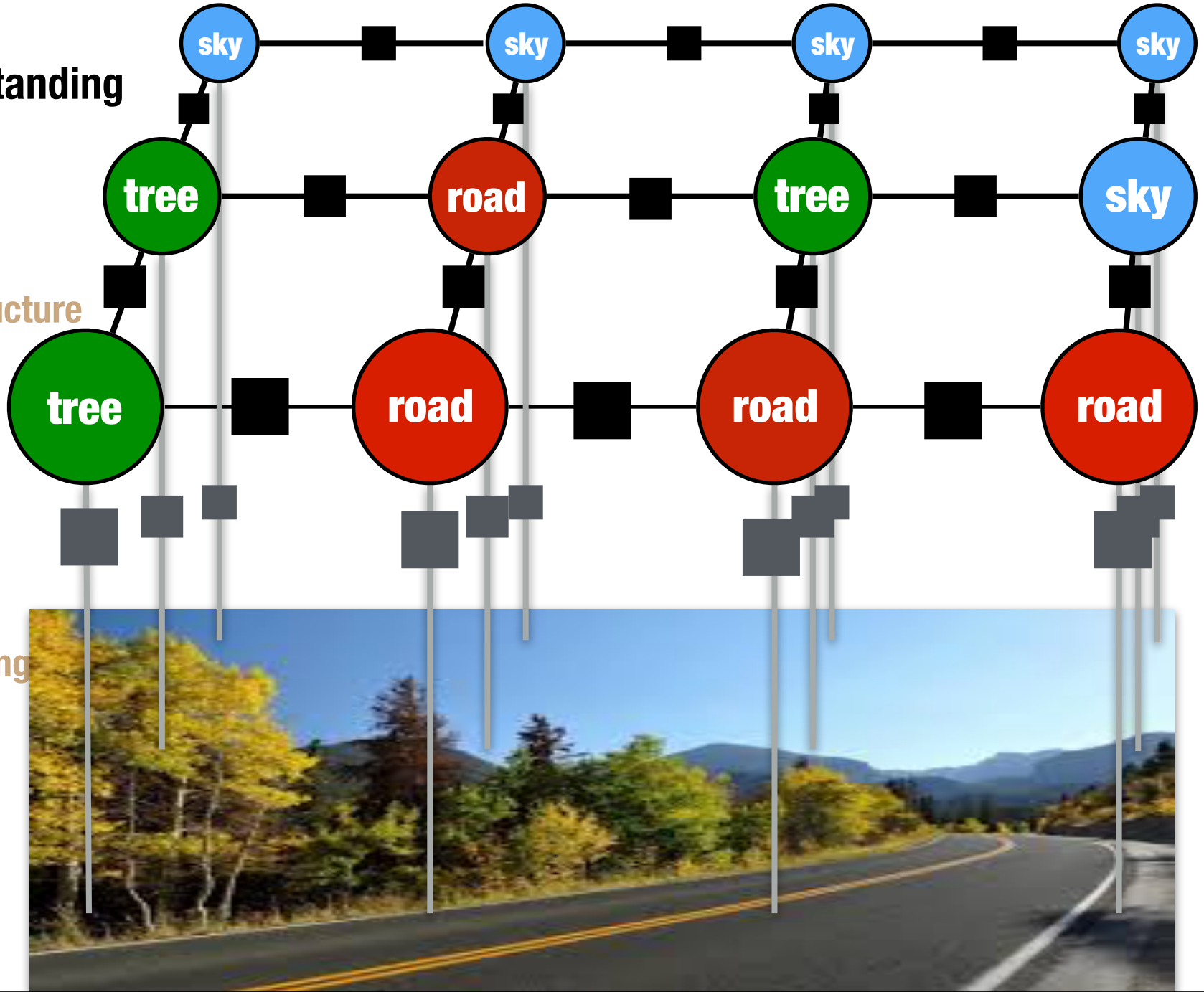
Dependency structure

Y

$E(X, Y)$

Feature Engineering

X



Structured Prediction *Sampling Inference*

Example:
Scene Understanding

$E(Y, Y)$

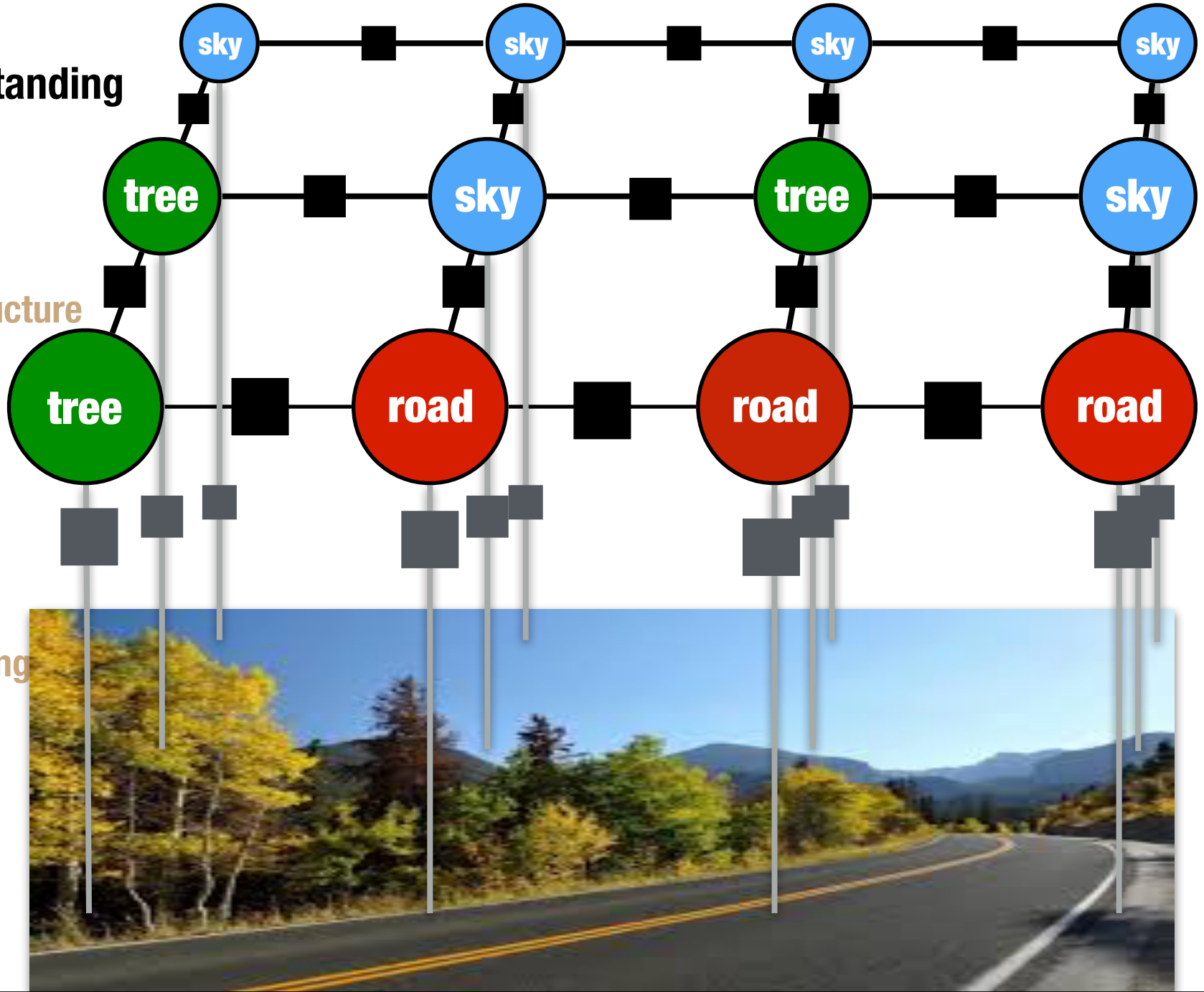
Dependency structure

Y

$E(X, Y)$

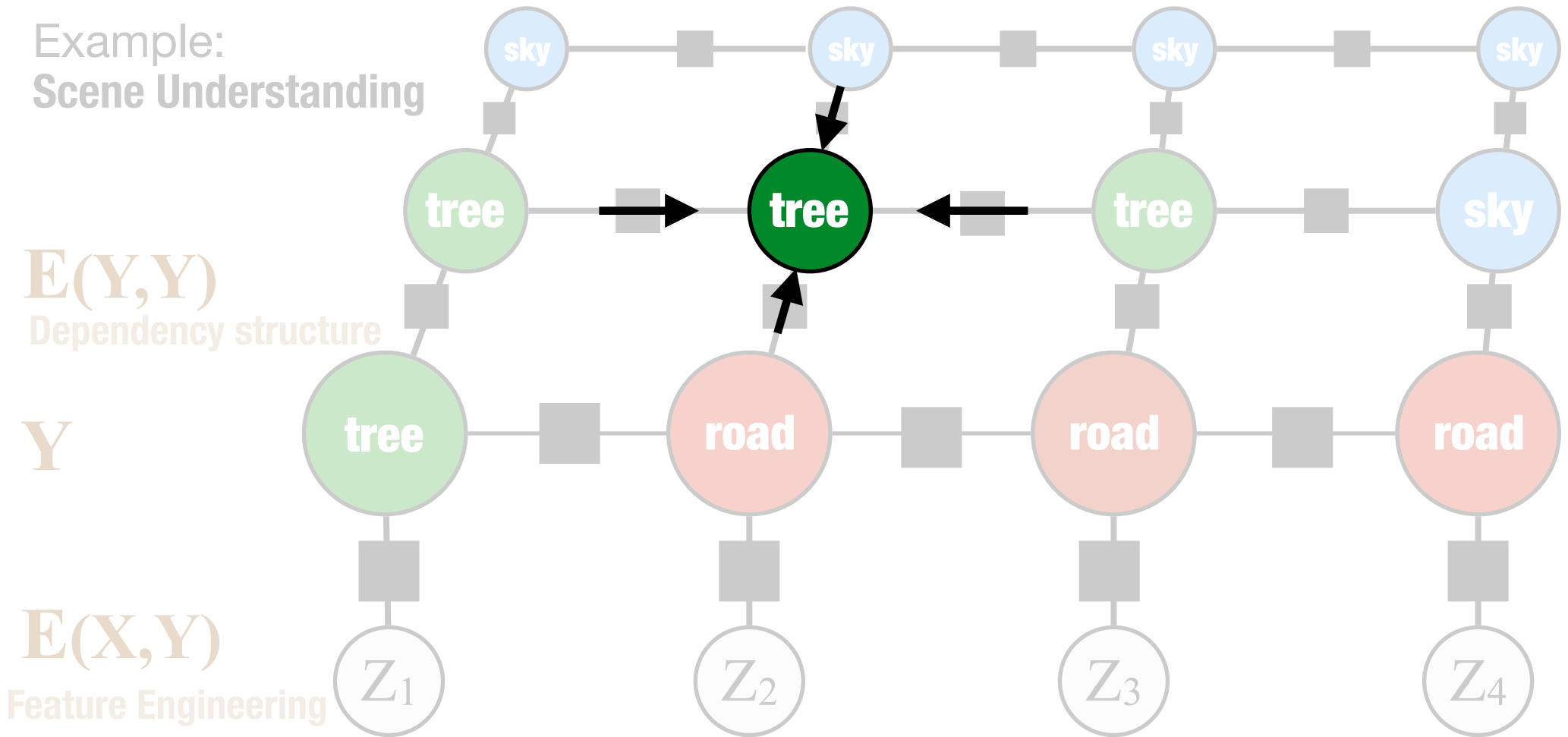
Feature Engineering

X



Structured Prediction *Variational Inference*

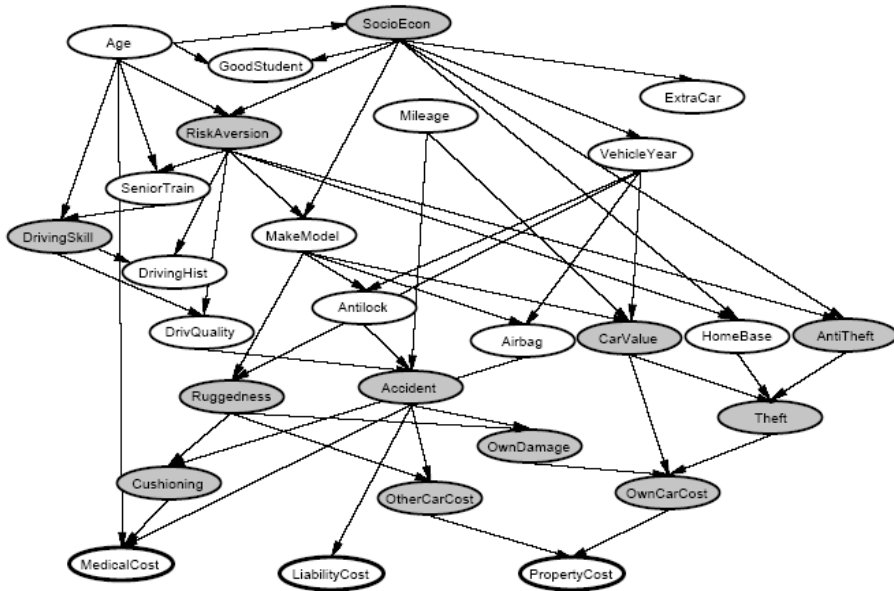
Example:
Scene Understanding



$$m_{i \rightarrow j}^{(t+1)}(x_j) = \sum_{x_i} \Phi_{ij}(x_i, x_j) \Phi_i(x_i) \prod_{k \in N(i)} m_{k \rightarrow i}^{(t)}(x_i)$$

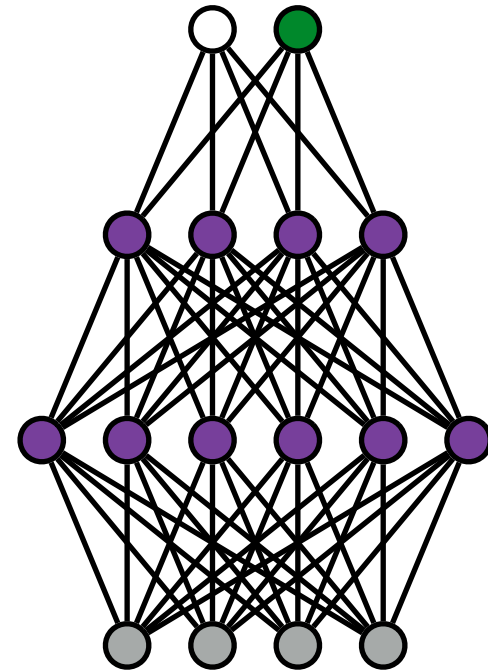
X

Bayesian Network



Sparsely connected
Hand-designed representations
Loopy/iterated inference (typically)
Cautious about capacity
“Statistically conscientious”

Deep Learning



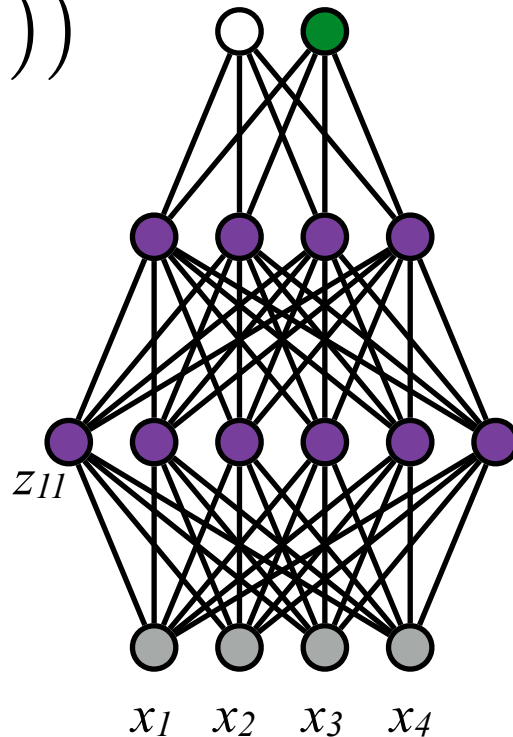
Densely connected (learn connectivity)
Learned, distributed representations
Feed-forward inference (typically)
Wild about high capacity
“Wild West” 😊

Deep Learning

$$y = \sigma(W_3 z_2) = \sigma(W_3 \sigma(W_2 \sigma(W_1 x)))$$

$$z_2 = \sigma(W_2 z_1)$$

$$z_{11} = \sigma \left(\sum_i z_{1i} w_{11i} \right) = \sigma(W_1 x)$$

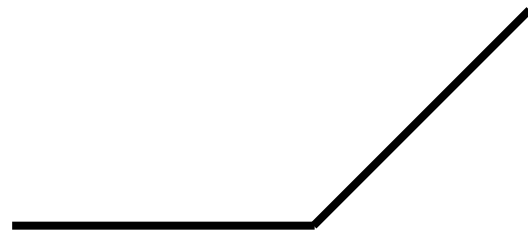


y
Predicted

z_2

z_1

x
Observed



$$\sigma(\cdot) = \max(\cdot, 0)$$

Deep Learning

$$y = F(\mathbf{x}; W)$$

Training Data $\left\{ \mathbf{x}^{(i)}, \mathbf{y}^{(i)} \right\}_{i=1}^N$

Loss

$$\mathcal{L} = \sum_i L \left(F(\mathbf{x}^{(i)}; W), \mathbf{y}^{(i)} \right)$$

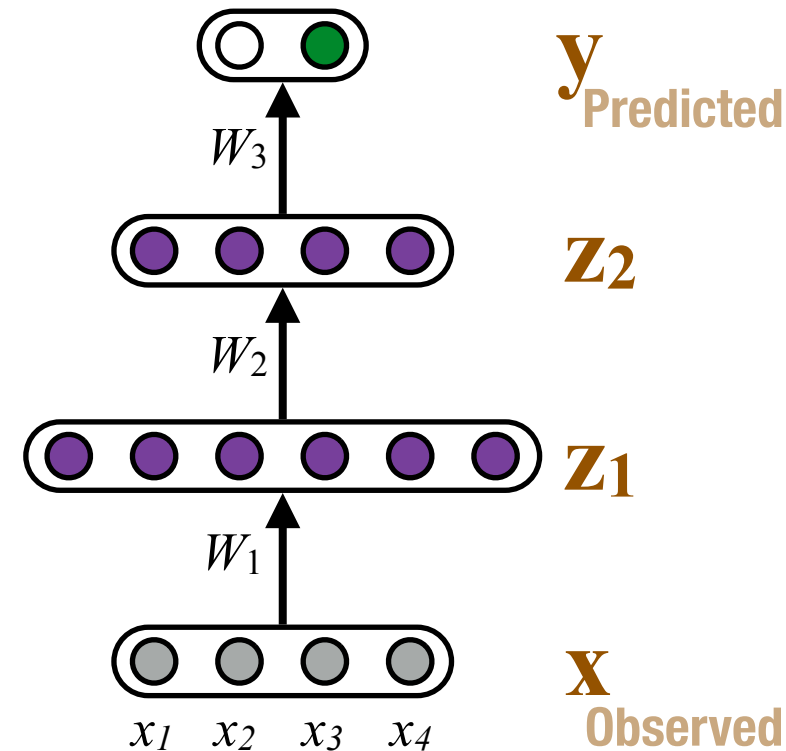
e.g. Squared error, Cross-entropy,...

Training

$$\arg \min_W \mathcal{L}$$

Gradient descent

$$W_{\text{new}} = W_{\text{old}} - \alpha \frac{\partial \mathcal{L}(W)}{\partial W}$$

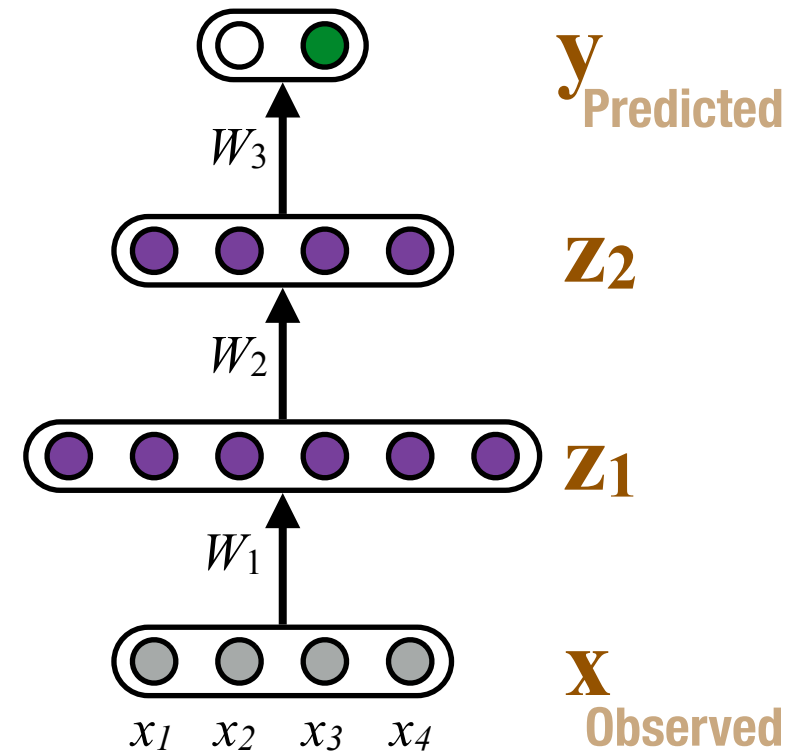


Key tools:

- (1) Back-propagation
- (2) Stochastic gradient descent

Deep Learning

$$y = F(\mathbf{x}; W)$$



Back-propagation

Deep Learning

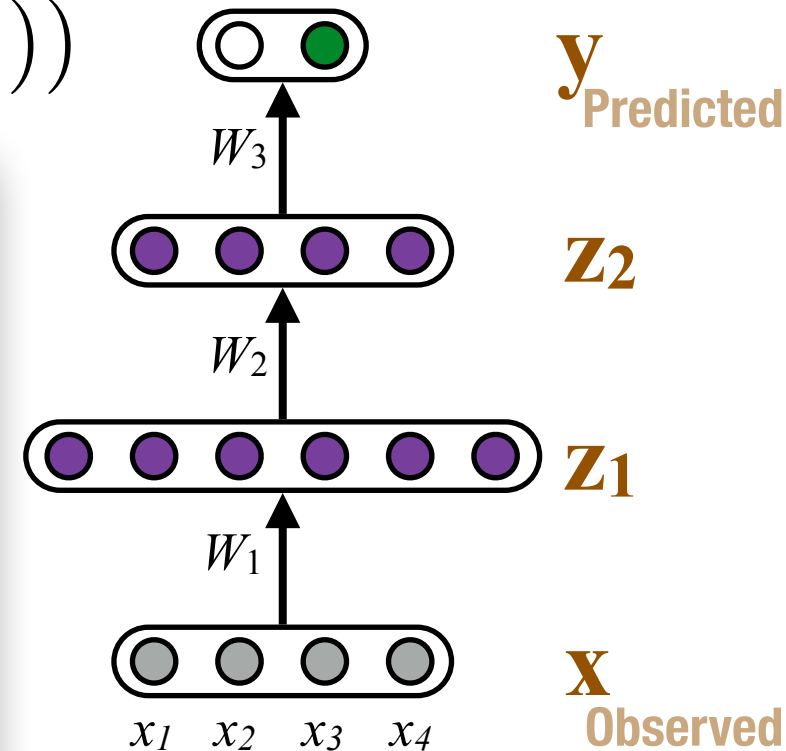
$$y = \sigma(W_3 \sigma(W_2 \sigma(W_1 \mathbf{x})))$$

The “chain rule”

$$g(f(x))' = g'(f(x)) \cdot f'(x)$$

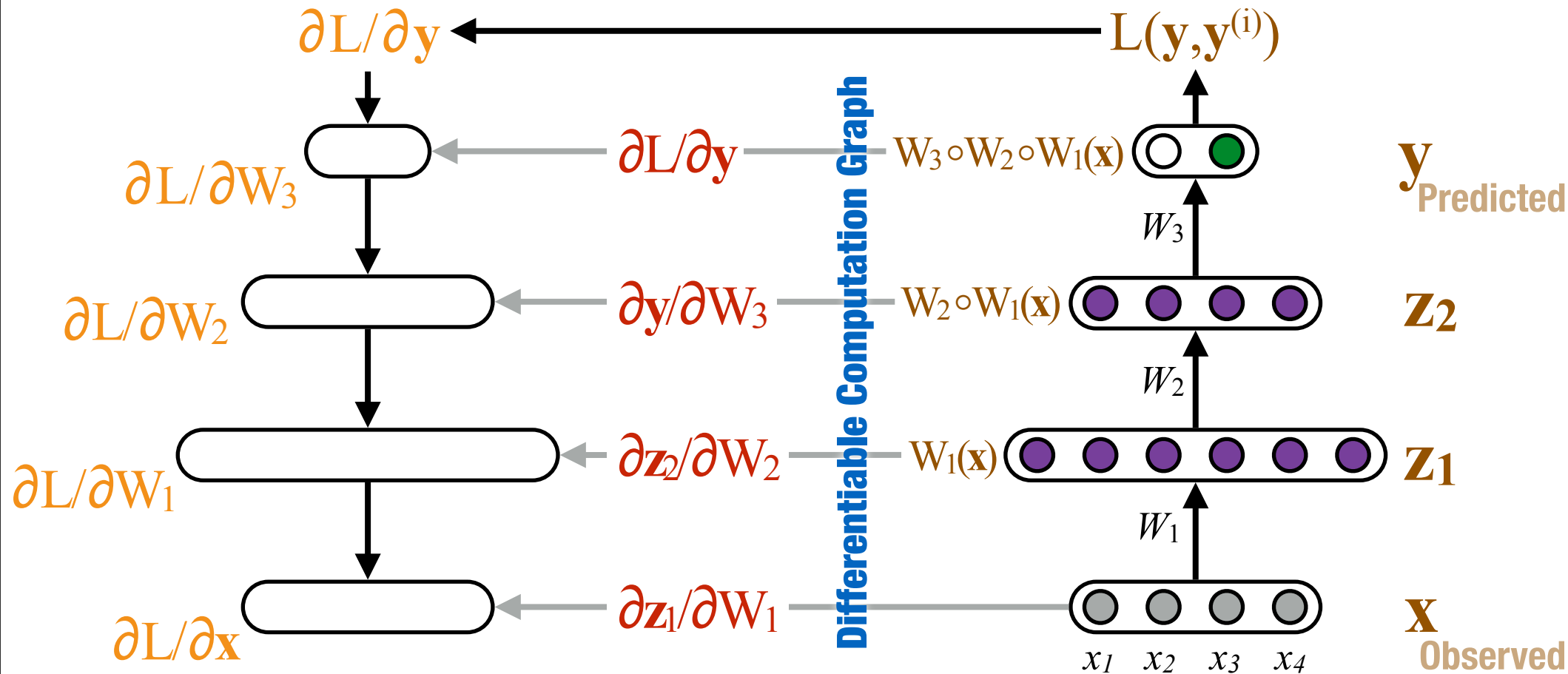
$$\frac{\partial g \circ f}{\partial x} = \frac{\partial g}{\partial f} \cdot \frac{\partial f}{\partial x}$$

$$\frac{\partial j \circ i \circ h \circ g \circ f}{\partial x} = \frac{\partial j}{\partial i} \frac{\partial i}{\partial h} \frac{\partial h}{\partial g} \frac{\partial g}{\partial f} \frac{\partial f}{\partial x}$$



Back-propagation

Deep Learning



Can get gradient of *Loss* wrt parameters at any depth from

- (1) local partial derivative functions
- (2) numeric gradient from above

Example: CNNs for Object Classification in Images

Representation Learning

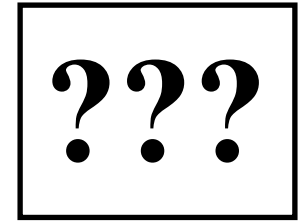


Motivation for SPENs

1. Use power of deep learning for *structure learning*

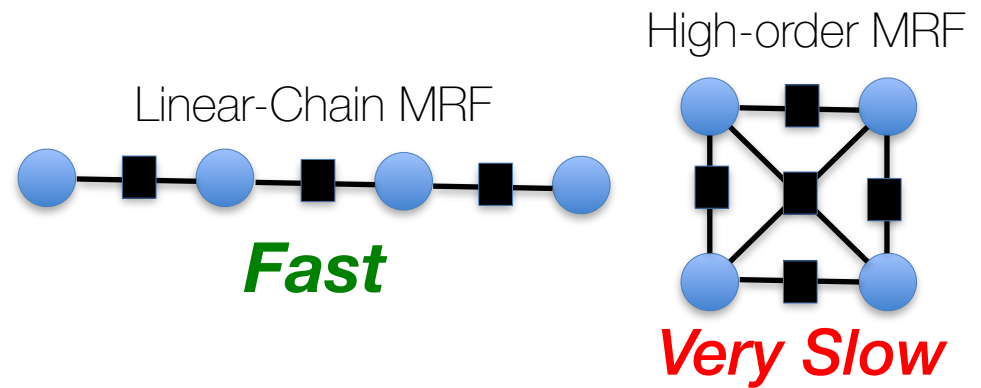


x

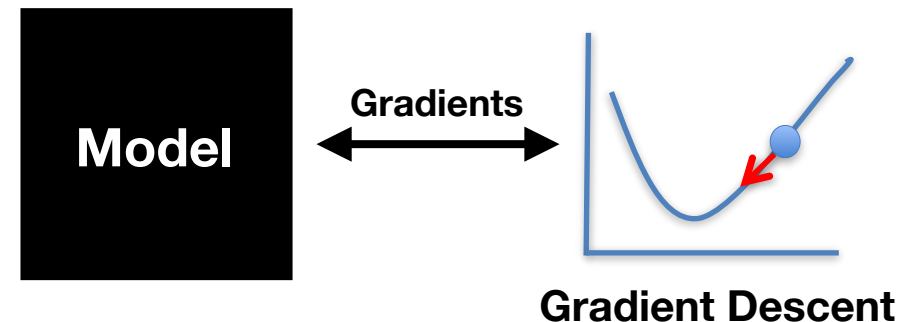


y

2. Provide an alternative to graphical models.



3. Black-box interaction with model.



Structured Prediction Energy Networks

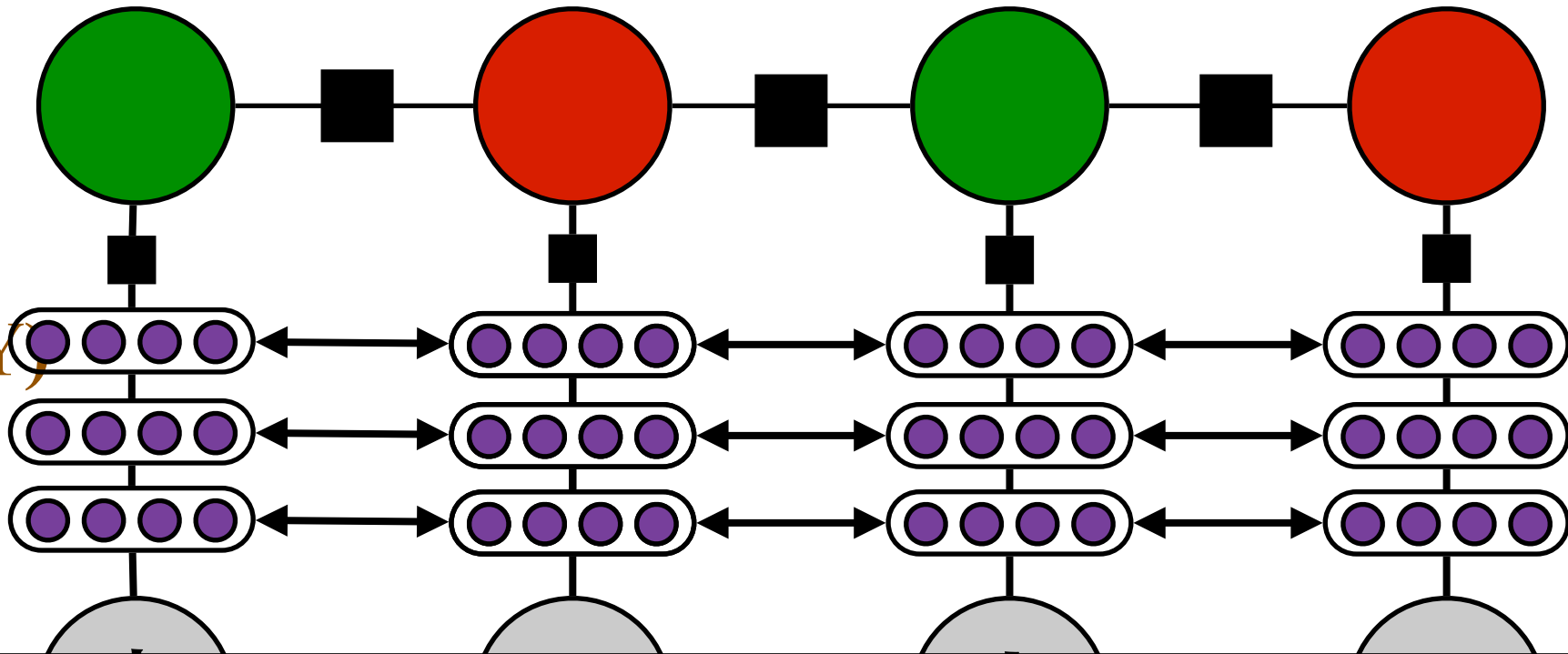
[Belanger, McCallum, ICML 2016]

$E(Y, Y)$

$$\Psi_0[y_0, y_1] + \Psi_1[y_1, y_2] + \Psi_2[y_2, y_3]$$

$Y \in \{0, 1\}$

$E(X, Z_{..}, Y)$



Structured Prediction Energy Networks

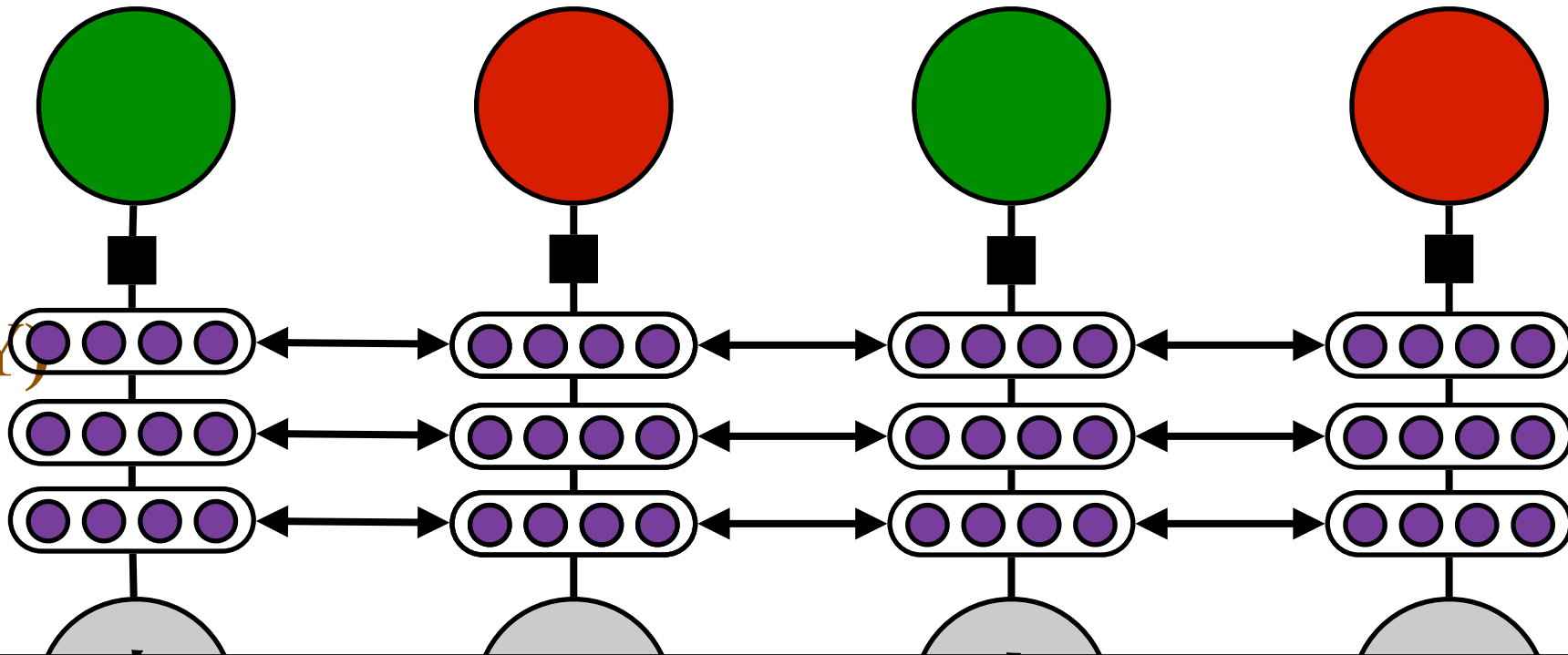
[Belanger, McCallum, ICML 2016]

$E(Y, Y)$

$Y \in \{0, 1\}$

$E(X, Z_{\dots}, Y)$

X



Structured Prediction Energy Networks

[Belanger, McCallum, ICML 2016]

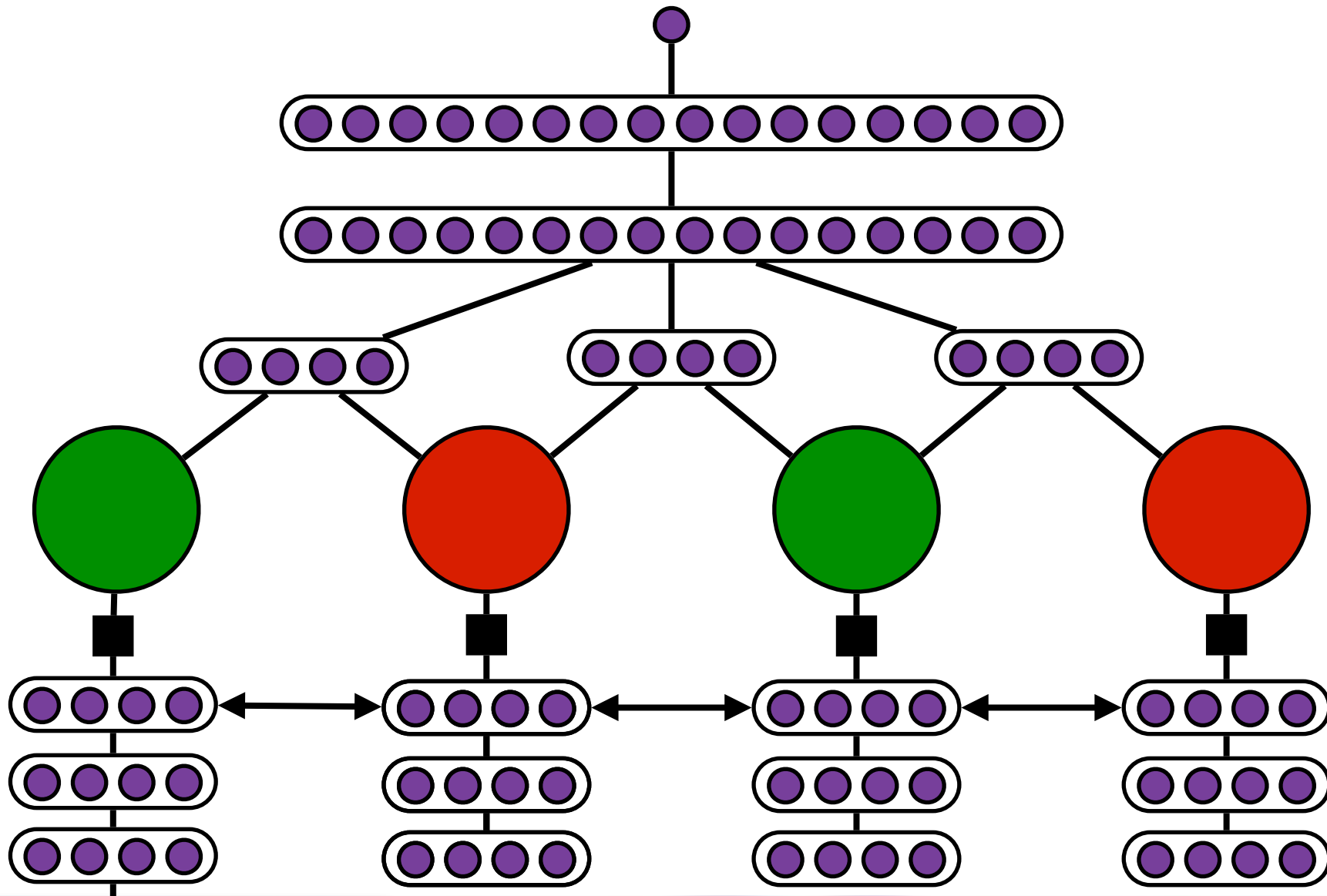
$E(y, y)$

y

$E(y, z; x)$

z

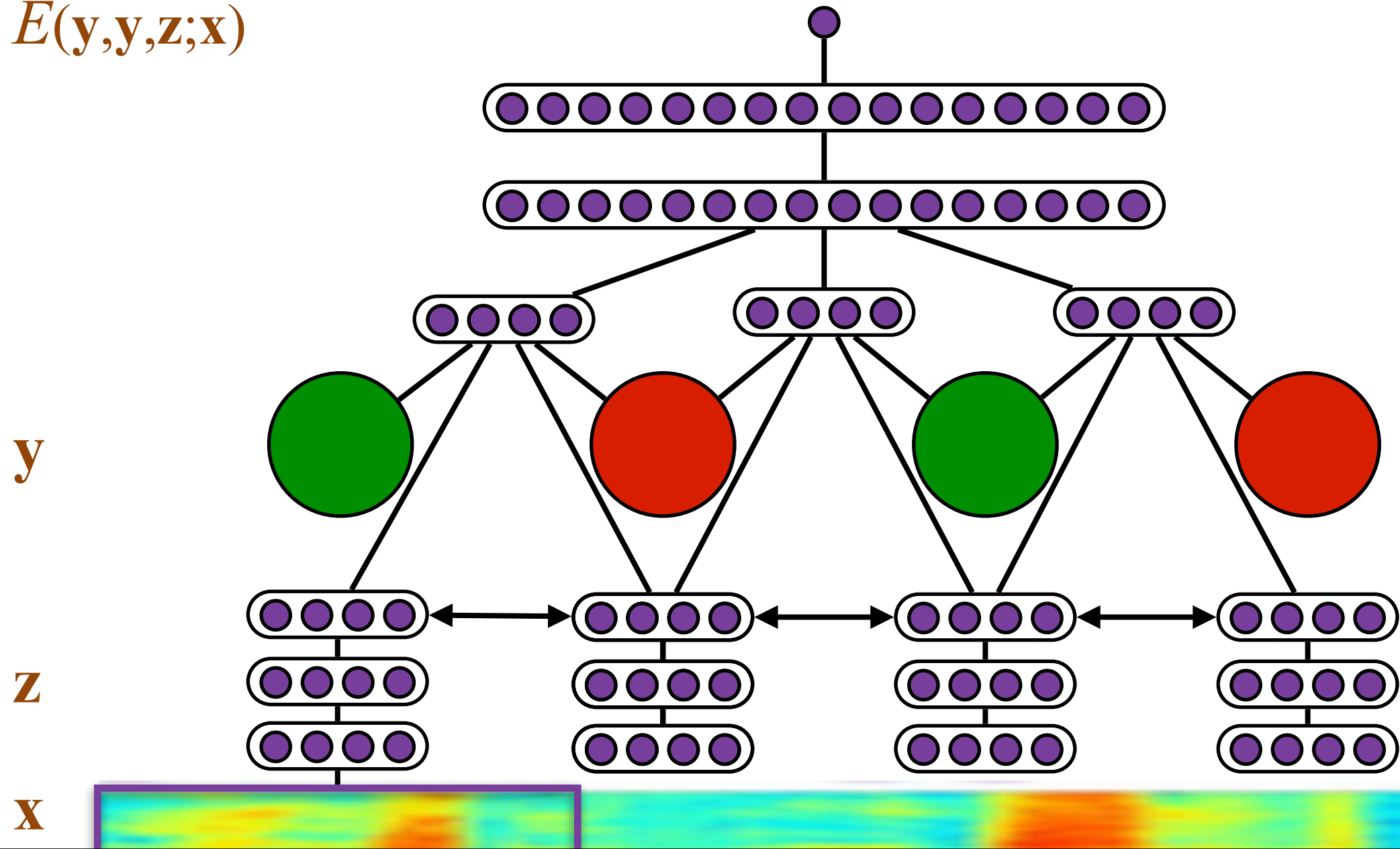
x



Structured Prediction Energy Networks

[Belanger, McCallum, ICML 2016]

$$E(y, y, z; x)$$



Structured Prediction Energy Networks

[Belanger, McCallum, ICML 2016]

Energy network

$$E(\bar{\mathbf{y}}; F(\mathbf{x}))$$

$$E(\mathbf{y}, \mathbf{y}, \mathbf{z}; \mathbf{x})$$

Soft prediction...

found by gradient descent

$$\bar{\mathbf{y}}^* = \arg \min_{\bar{\mathbf{y}} \in [0, 1]^L} E(\bar{\mathbf{y}}; F(\mathbf{x}))$$
$$\frac{\partial E(\bar{\mathbf{y}}; F(\mathbf{x}))}{\partial \bar{\mathbf{y}}}$$

Relax \mathbf{y} , to be continuous

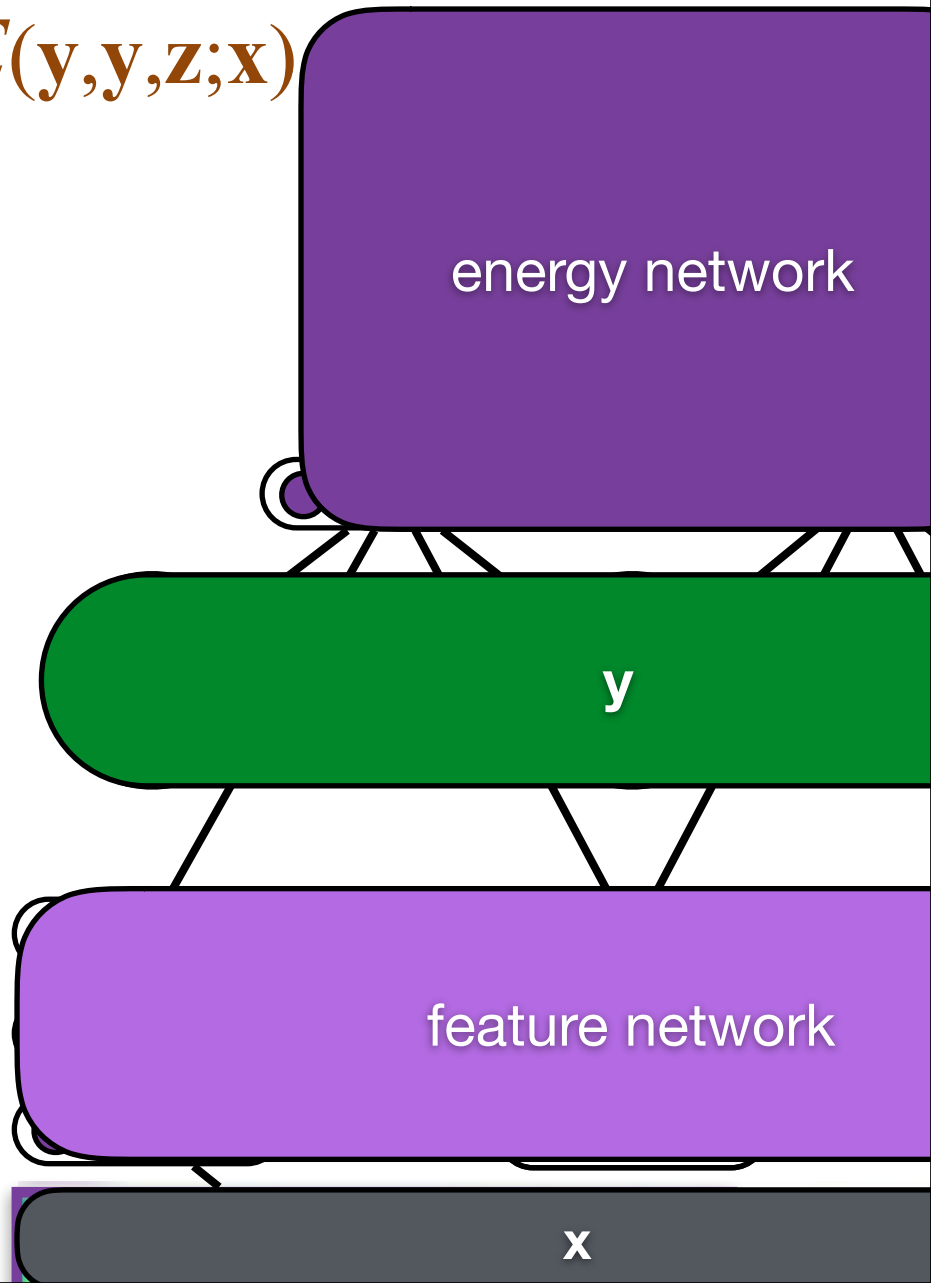
$$\mathbf{y} \in \{0, 1\}^L \rightarrow \bar{\mathbf{y}} \in [0, 1]^L$$

Feature Network

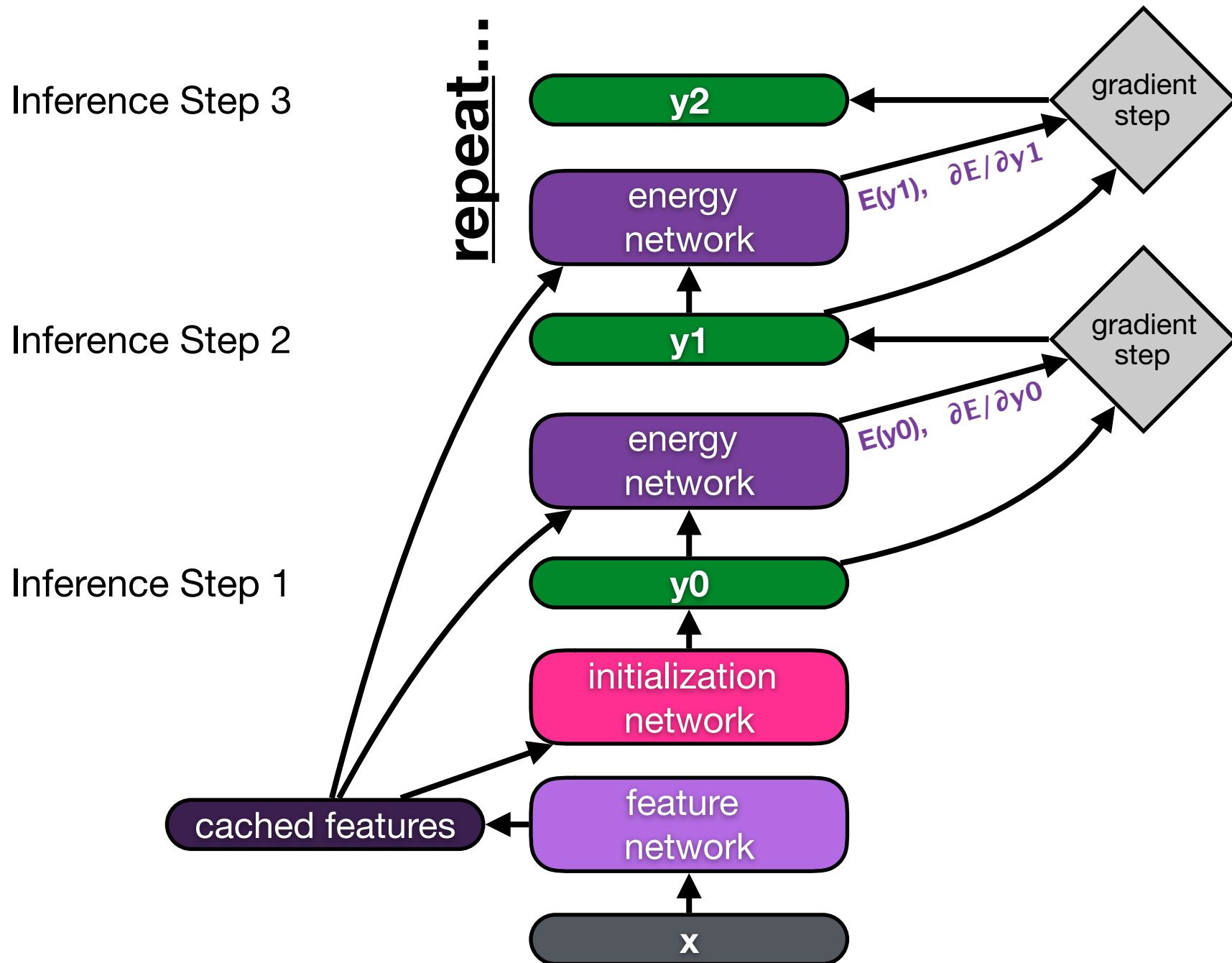
$$F(\mathbf{x})$$

\mathbf{z}

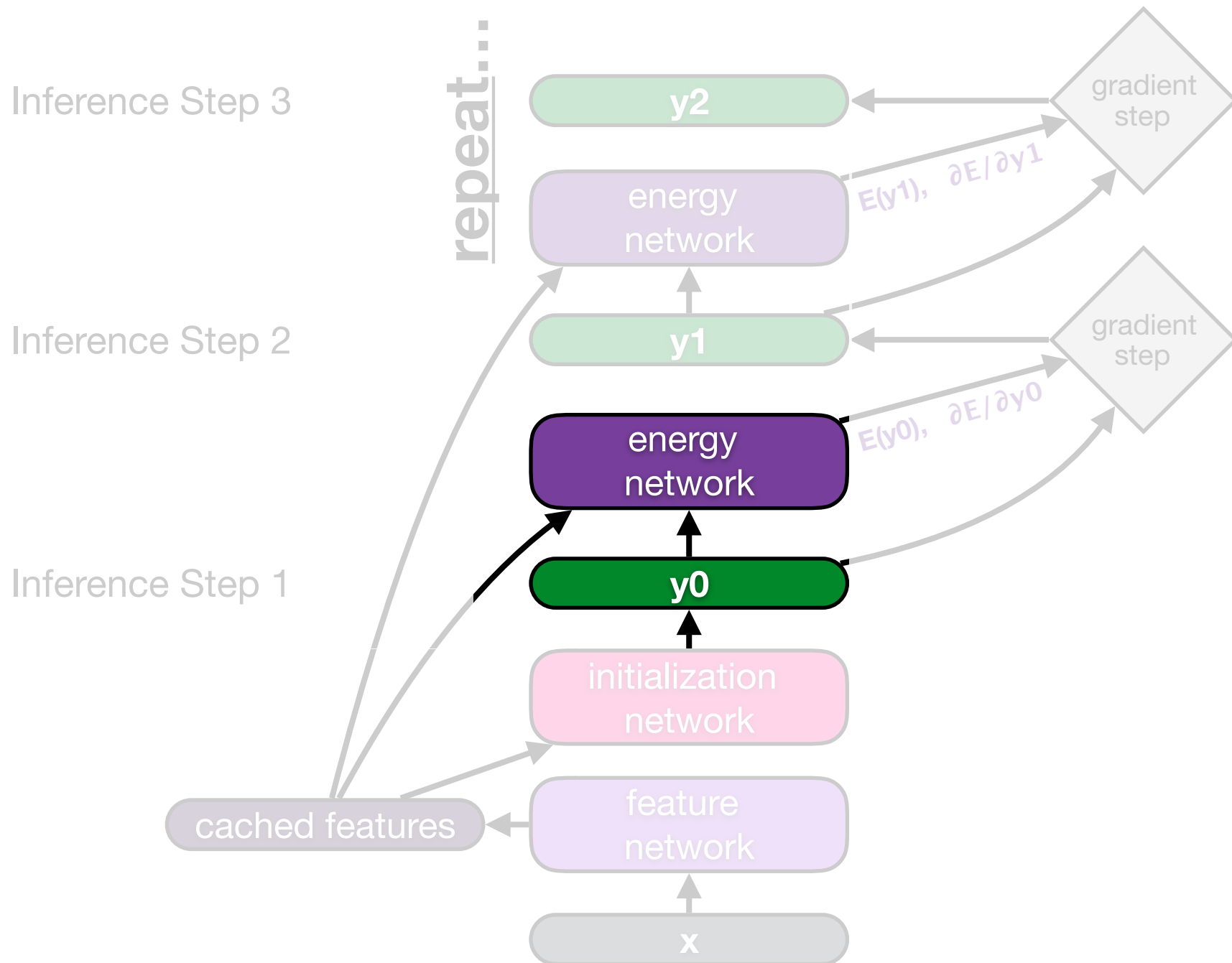
\mathbf{x}



SPEN Inference Graph



SPEN Inference Graph



Gradient used to Modify Inputs

“A Neural Algorithm for Artistic Style”

[Gatys et al. 2015]



SPENs use similar idea:

Optimize energy using backprop all the way down to the raw pixels.

Learning Algorithm 1: Structured SVM

Belanger, McCallum,
ICML 2016

Training Loss= $\mathcal{L} =$

(Taskar et al., 2004; Tsochantaridis et al., 2004)

$$\sum_{\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{\text{training data}}} \max_{\bar{\mathbf{y}}} \left[\underbrace{\Delta(\mathbf{y}^{(i)}, \bar{\mathbf{y}})}_{\substack{\text{Penalty} \\ \text{true} \quad \text{predicted}}} - \underbrace{\left(E_{\bar{W}}(\bar{\mathbf{y}}; \mathbf{x}^{(i)}) - E_W(\mathbf{y}^{(i)}; \mathbf{x}^{(i)}) \right)}_{\substack{\text{Model's energy difference} \\ \text{predicted} \quad \text{true}}} \right]_+$$

search requires **Loss-Augmented Inference**

$$\arg \min_{\bar{\mathbf{y}}} \left(\underbrace{-\Delta(\mathbf{y}^{(i)}, \bar{\mathbf{y}})}_{\substack{\text{Penalty must be} \\ \text{differentiable}}} + E_W(\bar{\mathbf{y}}; \mathbf{x}^{(i)}) \right)$$

Stochastic Gradient

$$\frac{\partial \mathcal{L}}{\partial W}$$

Learning Algorithm 2:

Belanger, McCallum,
ICML 2017

End-to-end “backprop through inference”

Training Loss= $\mathcal{L} =$

Direct Risk Minimization

$$\sum_i L \left(\mathbf{y}^{(i)}, \text{Algorithm}_{\mathbf{y}}(\mathbf{x}^{(i)}) \right)$$

training data

Algorithm for inference

$$\bar{\mathbf{y}}^* = \bar{\mathbf{y}}^{[0]} + \sum_{t=1}^T \alpha_t \frac{\partial}{\partial \bar{\mathbf{y}}} E_W(\mathbf{x}, \bar{\mathbf{y}}^{[t-1]})$$

sum over “time steps” of inference

Direct application of:
Justin Domke, AISTATS, 2012.

“Generic Methods for Optimization-Based Modeling”

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial L}{\partial \bar{\mathbf{y}}^*} \frac{\partial \bar{\mathbf{y}}^*}{\partial W} = \sum_{t=1}^T \alpha_t \frac{\partial L}{\partial \bar{\mathbf{y}}^*} \left(\frac{\partial}{\partial W} \frac{\partial}{\partial \mathbf{y}} E_W(\mathbf{x}, \bar{\mathbf{y}}^{[t-1]}) \right)$$

Hessian-Vector product can be approximated using one-dimensional finite differences

sum over “time steps” of inference

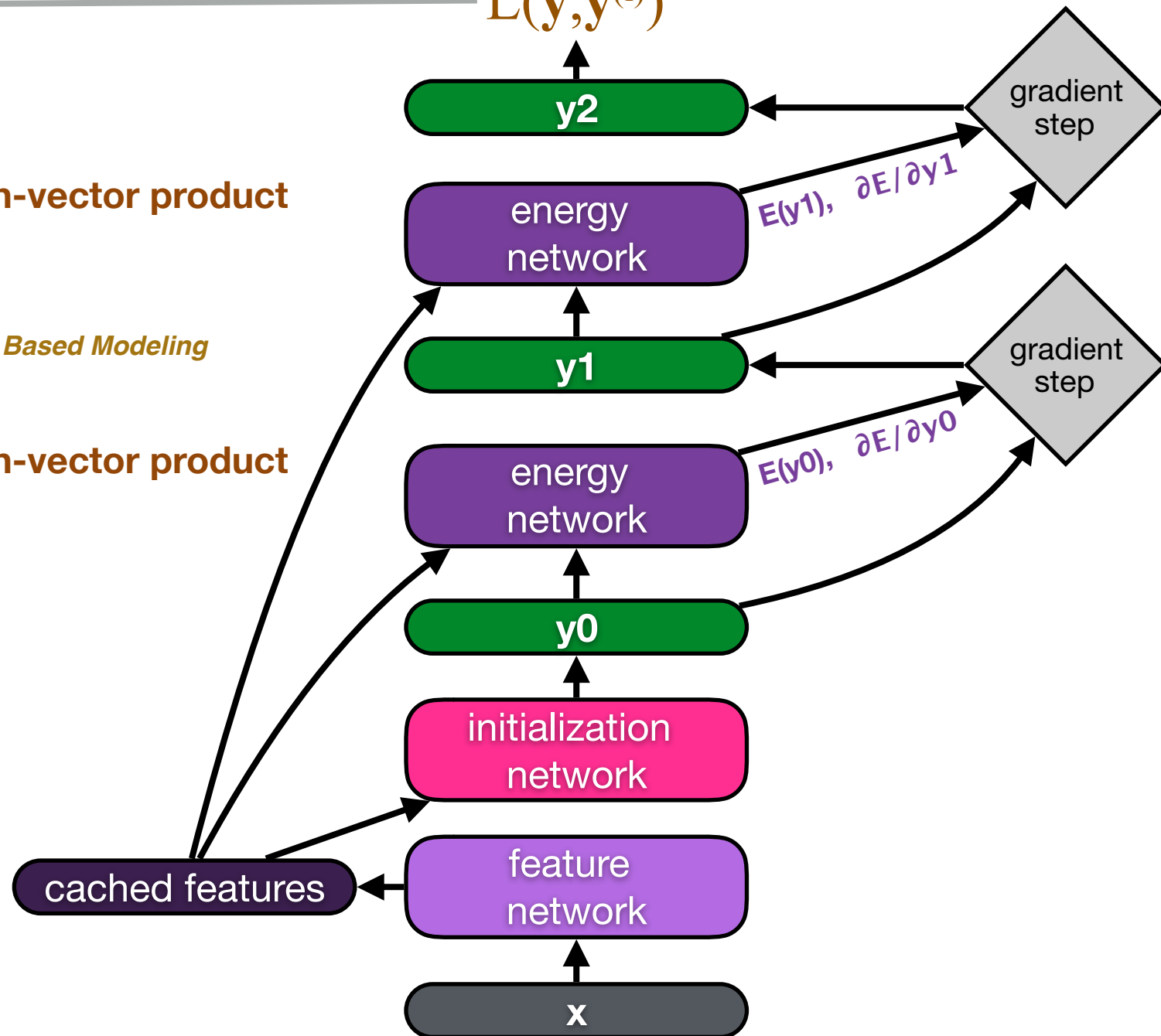
Learning Algorithm 2 Graph

$$\partial L / \partial \mathbf{y} \leftarrow L(\mathbf{y}, \mathbf{y}^{(i)})$$

Hessian-vector product

Domke, 2012.
Generic Methods for Optimization Based Modeling

Hessian-vector product



Chapter 3

Light Supervision training of Structured Prediction Energy Networks

(Turing complete!)

1. Human writes arbitrary prior knowledge
2. Learn model with arbitrary dependencies.
(SPEN)
3. Efficient inference by gradient descent.

AUTHOR Anna Popescu (2004), “Interactive Clustering,”
EDITOR Wei Li (EDITOR Ed.), Learning Handbook, Athos Press,
LOCATION Souroti.

Human writes arbitrary prior knowledge...

“AUTHOR field should be contiguous, only appearing once.”

...as a scoring function $V(x=\text{citation}, y=\text{labeling})$

```
score = 0
score -= 1 foreach AUTHOR non-contiguous
score -= 1 if has both JOURNAL & BOOKTITLE
score -= 1 foreach “using” not in TITLE
score -= 1 foreach [A-Z]\. not AUTHOR|EDITOR
score -= 1 if PUBLISHER before JOURNAL ...
```

(like rule-based AI before ML was popular)

Why use ML if we get a ruled-based scoring function?

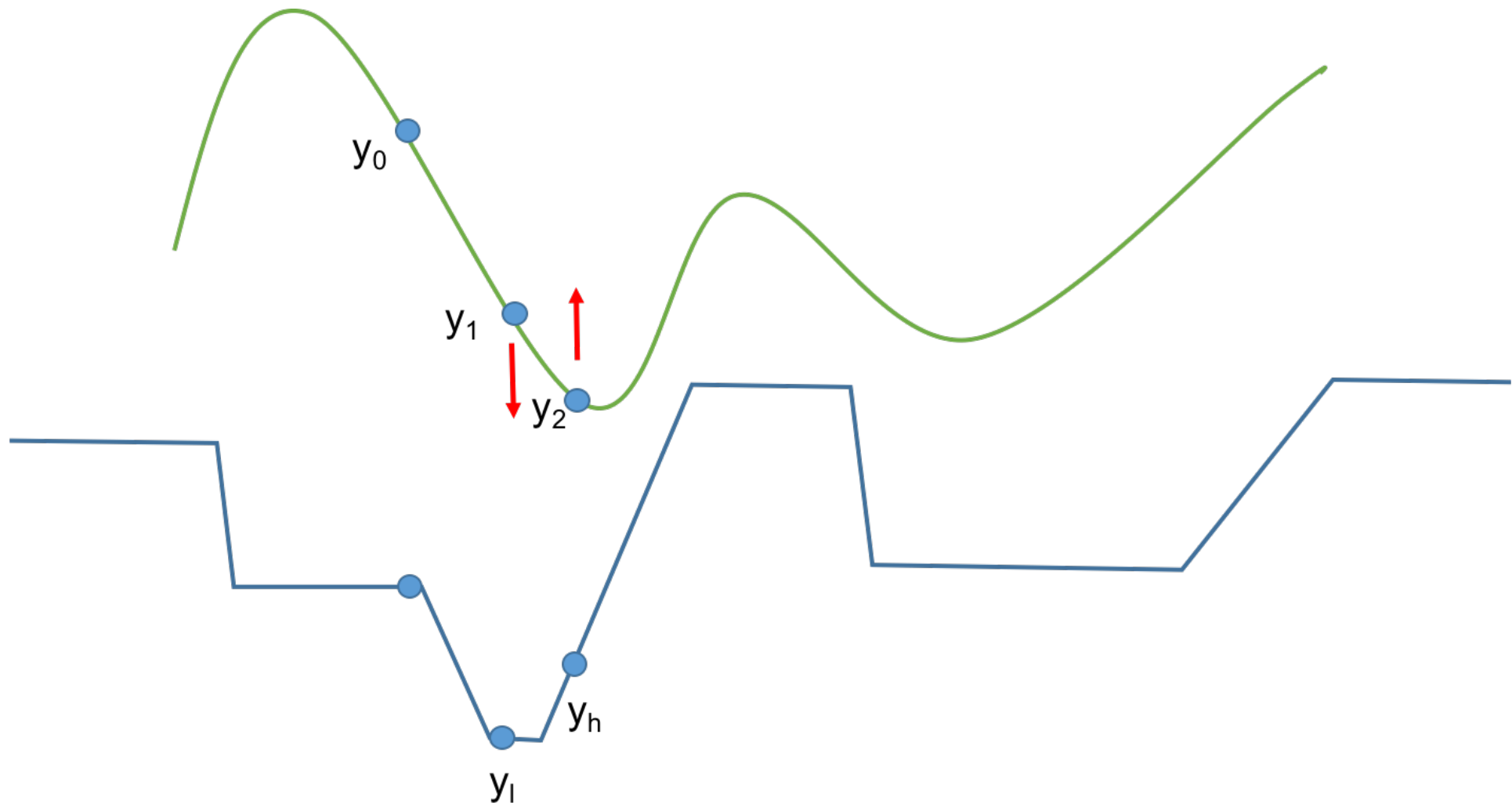
- **Doesn't generalize**
 - examines just a few features
 - SPENs will learn correlated features, labels.
- **No inference procedure** just scores for given (x,y)
 - stochastic optimization is slow
 - SPENs provide gradient-descent inference

Learning Algorithm 3:

Rooshenas, McCallum,...
forthcoming

“ranking successive gradient steps”

Training Loss = $\sum_{\mathbf{x} \in \mathcal{D}} [\alpha(V(\mathbf{y}_h, \mathbf{x}) - V(\mathbf{y}_l, \mathbf{x})) - E_{\mathbf{w}}(\mathbf{y}_h, \mathbf{x}) + E_{\mathbf{w}}(\mathbf{y}_l, \mathbf{x})]_+$



Preliminary Experiments

(...much more work and comparisons in future...)

Weak-Sup SPEN: simple test

Multi-label Document Classification

x = Medical bag-of-words

[amount, cystourethrogram, diagnosed, episode, evaluate, exam, fever, grade, growth, hematuria, infection, interval, kidney, left, lower, occurred, patient, pole, previously, purpose, reflux, renal, scar, scarring, small, study, tract, urinary, vesicoureteral, voiding, year]

y = multiple ICM-9-CD codes

[593-70, 599-00]

x = Human background knowledge

Keyword descriptions of ICM-9-CD codes. (Not gathering any labeled correlation knowledge.)

593-70: vesicoureteral, reflux, unspecified, nephropathy
V79-99: viral, chlamydial, infection, conditions, unspecified
753-00: renal, agenesis, dysgenesis

Scoring function gives +1 for each label:keyword cooccurrence.

$$V(y^i, x^i) = \sum_j I(l_j \in y^i) I(|x^i \cap w_j| > 0) - \gamma \max(|y^i| - 1, 0)$$

Label, Keyword matches Sparsity constraint

Does the SPEN generalize over the human scoring function?

ICM-9-CD code data set, evaluate **F1 of label set**

Human Scoring Function, Exhaustive Search						SPEN
$N \leq 1$	$N \leq 2$	$N \leq 3$	$N \leq 4$	$N \leq 5$	$N \leq 6$	
15.5	18.3	19.6	20.5	21.1	20.3	22.6

(~10x faster)

Weak-Sup SPEN: better test Citation Field Extraction

x = Citation Token Sequence

Anna Popescu (2004), “Interactive Clustering,”
Wei Li (Ed.), Learning Handbook,
Athos Press, Soutoti.

y = Seq. of Labels $\in |14|$

**AUTHOR AUTHOR YEAR TITLE TITLE
EDITOR, EDITOR EDITOR BOOKTITLE, BOOKTITLE
PUBLISHER PUBLISHER LOCATION**

x = Human background knowledge

Human-written scoring function. 50 lines of code. Written in ~1 hour.

```
score -= 1 foreach AUTHOR non-contiguous  
score -= 1 if has both JOURNAL & BOOKTITLE  
score -= 1 foreach “using” not in TITLE  
...
```

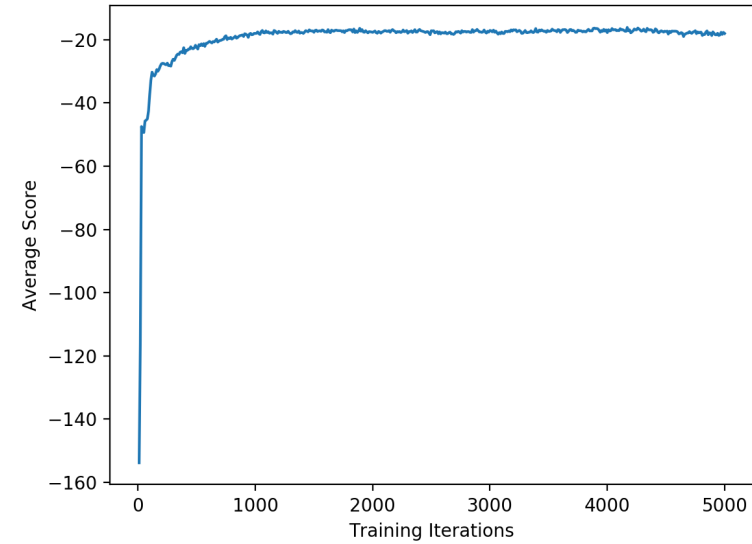
~4000 unlabeled examples, 0 labeled.

Scoring function advice:

- Penalties only, so 0 = best.
- Can use varying magnitudes, -1, -5, -10.
- Debug with some stochastic optimization.

Citation Field Extraction Accuracy

Method (no labeled data)	Token accuracy	Time sec/citation	Ave. V() score
GE [Mann & McCallum '10]	37%	?	N/A
V search 10	34%	14	-1.86
V search 100	39%	170	-0.98
V search 1000	42%	1240	-0.62
SPEN	52%	0.0008	~ -20



Example text

Wright, A. K. Simple imperative polymorphism. *Lisp and Symbolic Computation* 8, 4 (Dec. 1995), 343-356.

V search 100 output

AUTHOR	TITLE	AUTHOR	AUTHOR	AUTHOR	NOTE	NOTE	NOTE	NOTE	NOTE	NOTE	NOTE	DATE	DATE	PUB	PUB	PUB
--------	-------	--------	--------	--------	------	------	------	------	------	------	------	------	------	-----	-----	-----

SPEN output

AUTHOR	TITLE	TITLE	TITLE	TITLE	TITLE	TITLE	TITLE	TITLE	DATE	DATE	PAGES	PAGES	PAGES
--------	-------	-------	-------	-------	-------	-------	-------	-------	------	------	-------	-------	-------

Related Work

- **Deep Value Networks...**

[Gygli, Norouzi, Angelova 2017 ICML]

- Matching magnitude (**rather than just ranking**).
- Hurts accuracy? 5% vs SPEN's 52%

- **Constraint-Driven Learning**

[Chang, Ratnov, Roth 2007 ACL]

- Supervised training \Rightarrow Pseudo-label **data** w/ constraints \hookrightarrow

- **Snorkel: Rapid Training Data Creation with Weak Supervision**

[Ratner, Bach, Ehrenberg, Fries, Wu, Ré 2017 VLDB]

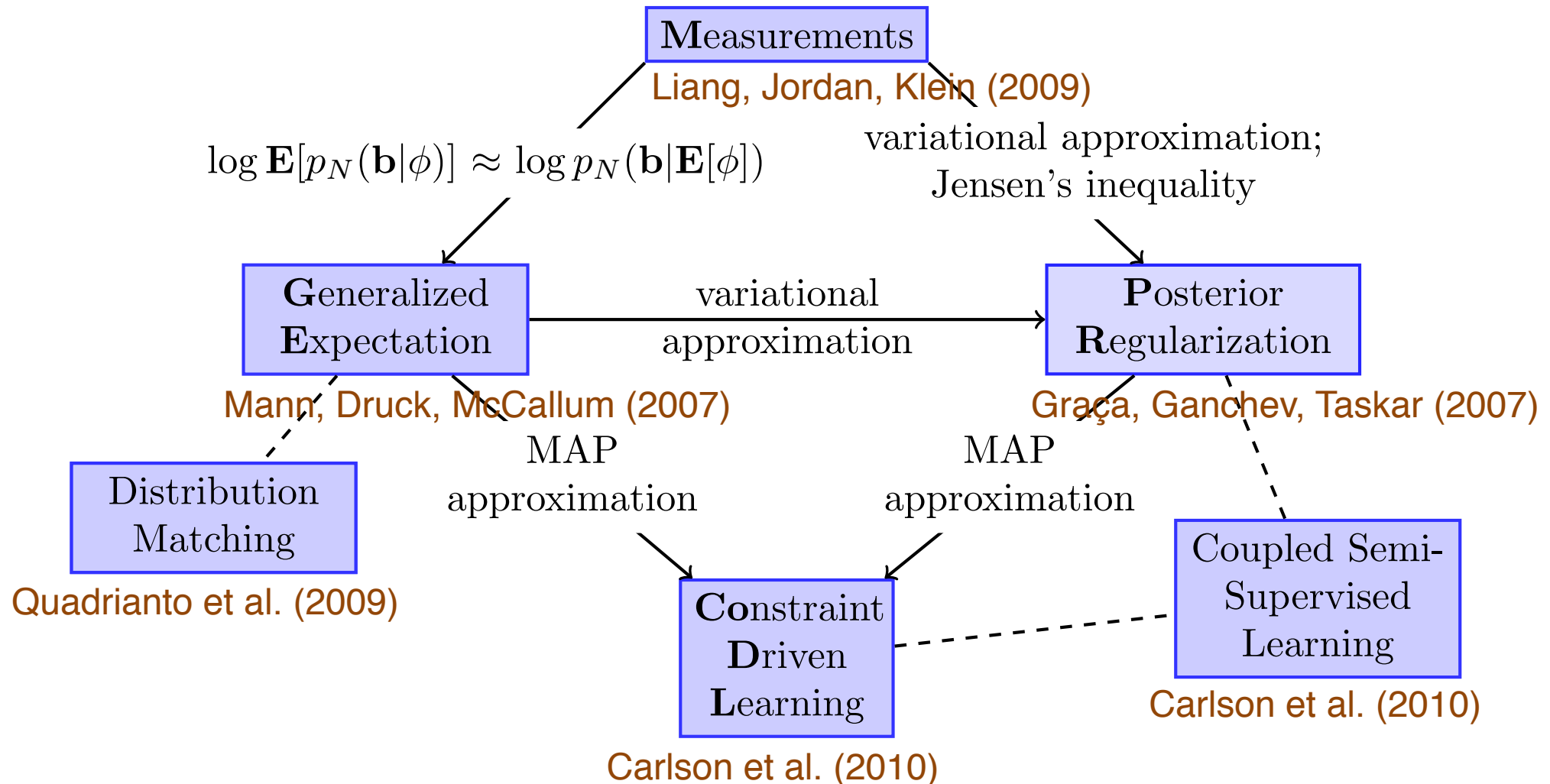
- Rules \Rightarrow Pseudo-labeled **data** \Rightarrow Supervised (self) training

- **Label-Free Supervision of NNs w/ ... Domain Knowledge**

[Stewart, Ermon 2017 AAAI]

- Constraints \Rightarrow Loss function \Rightarrow Train **feed-forward** NN.

GE Related Work



Summary

- ***Generalized Expectation***

- Learning from unlabeled data + “labeled features”
- Hard to do inference

- ***Structured Prediction Energy Networks***

- Representation learning for *output* variables
- Test-time inference by gradient descent
- New SPEN training method: Ranking

- **Experiments**

- Multi-label Classification: ICM-9
- Sequence labeling: Citation field extraction

- **Next**

- Training on corpus-wide expectations.
- Interactive tools for score function development.

