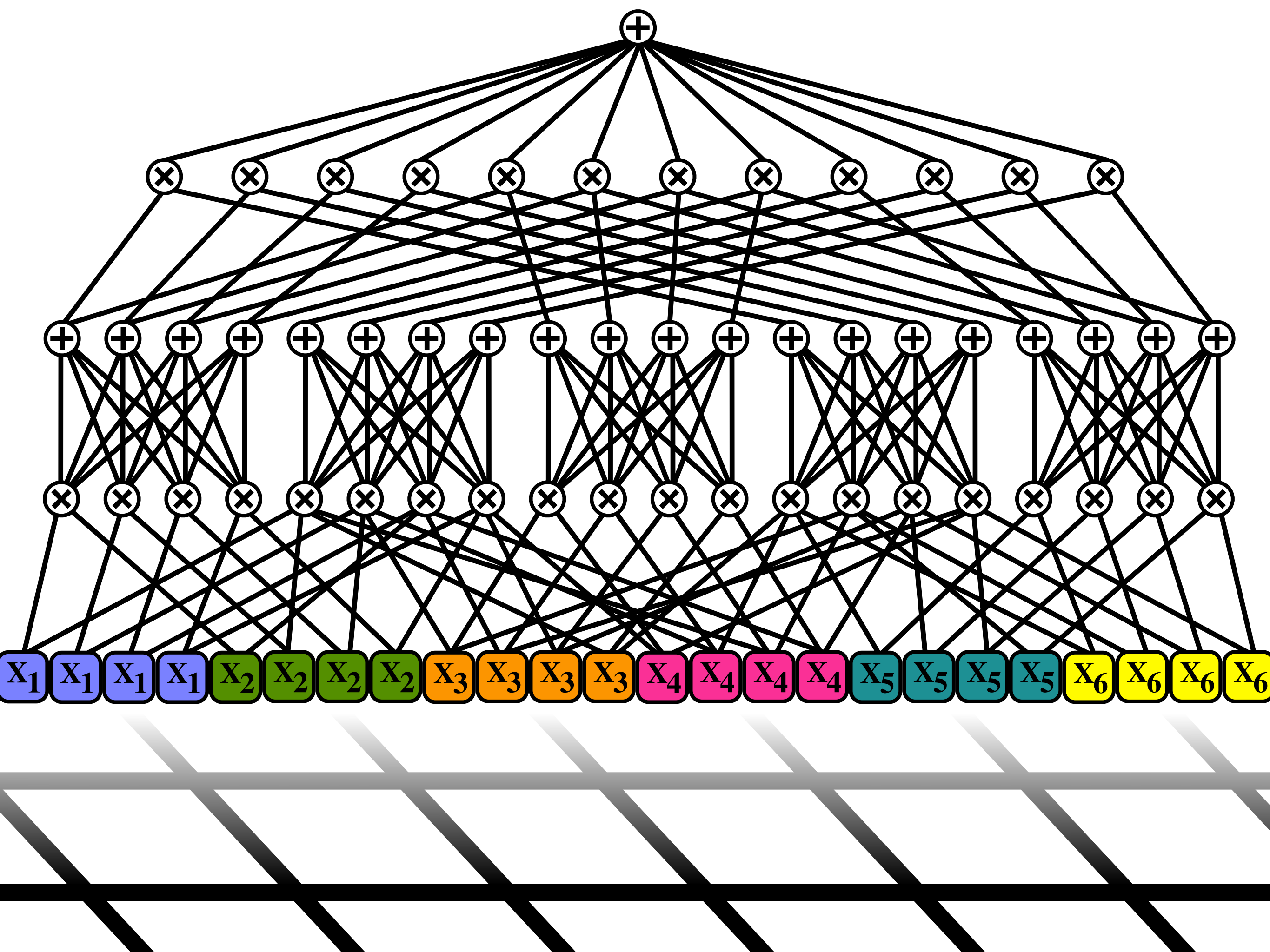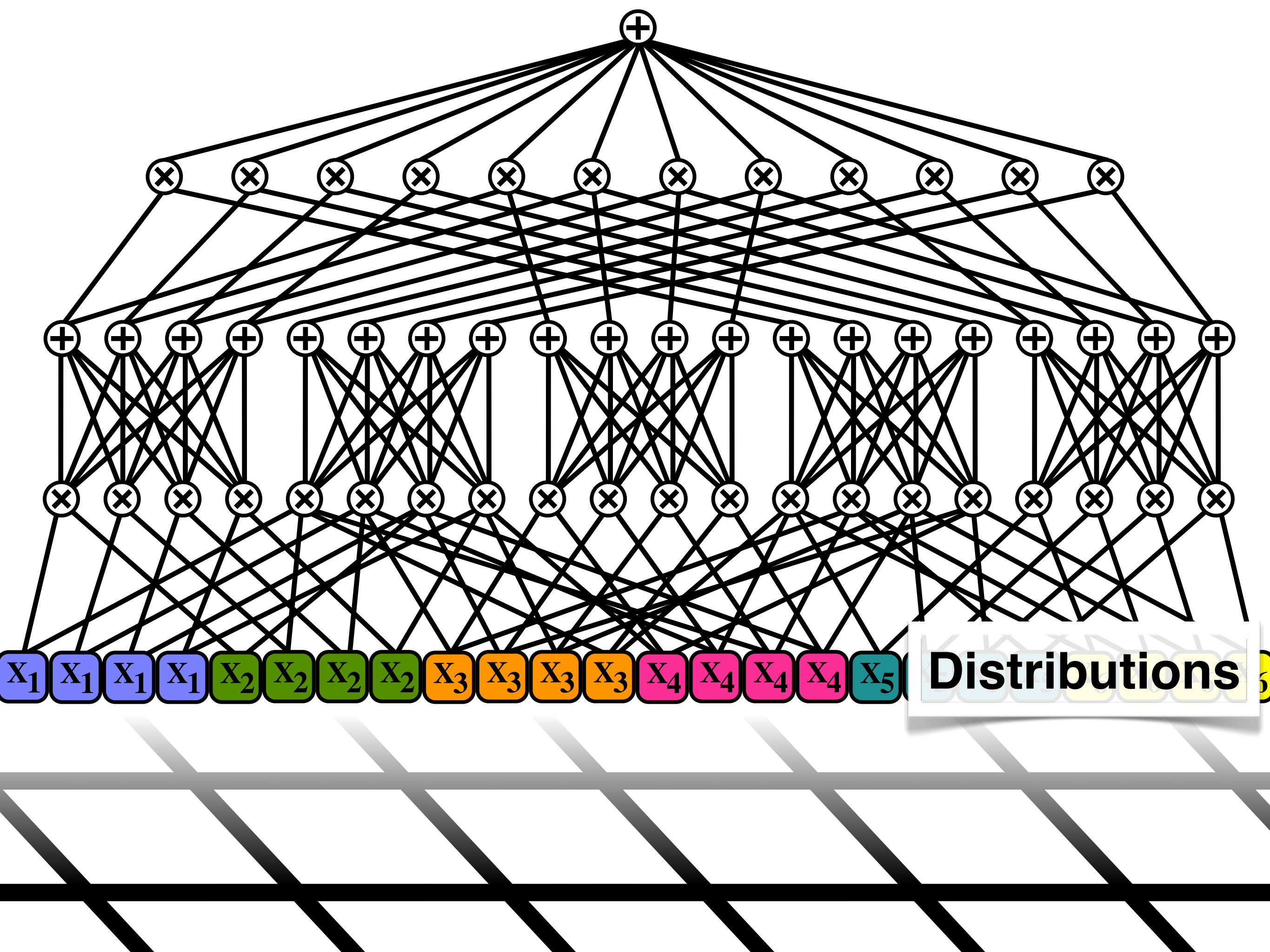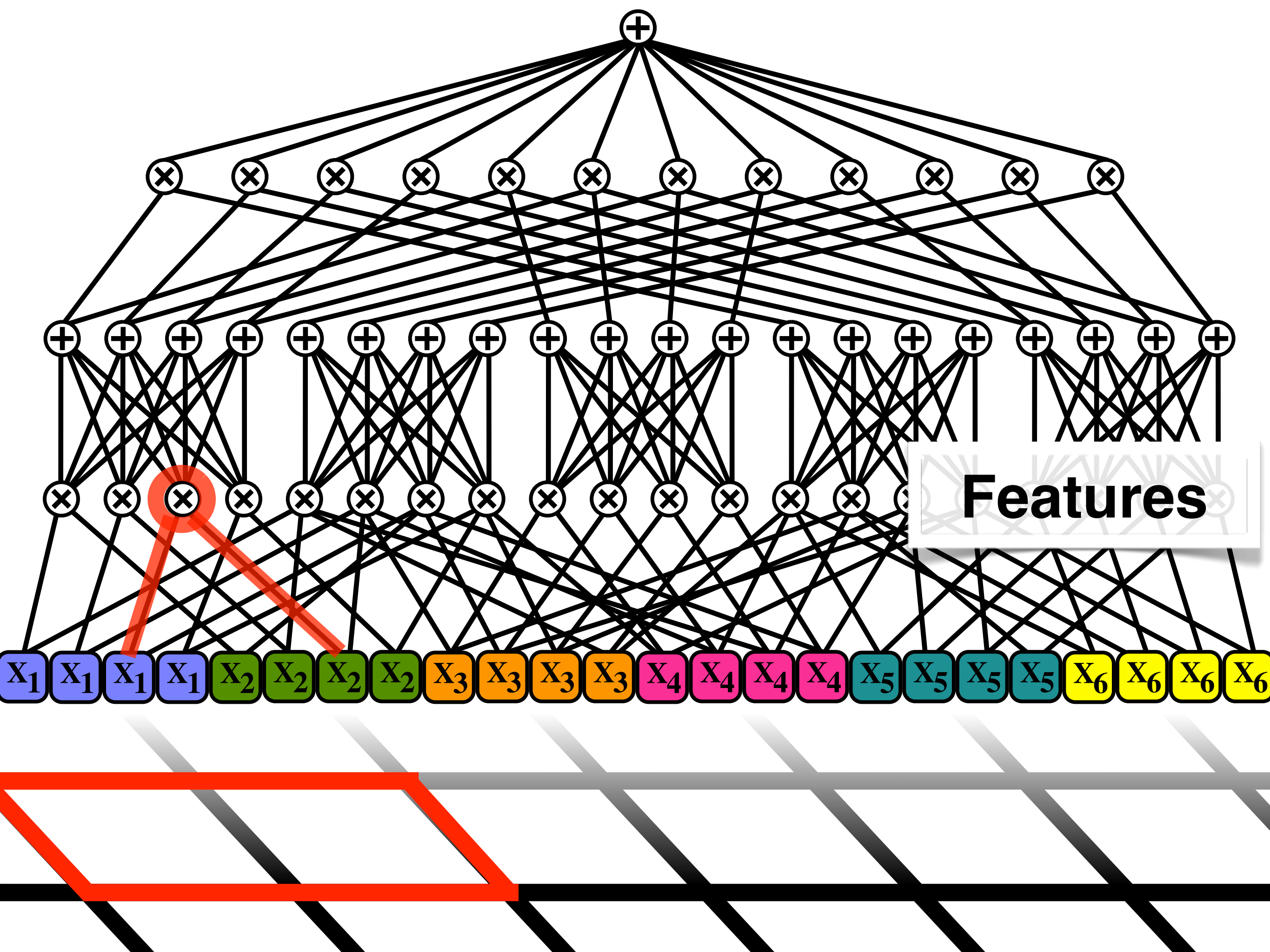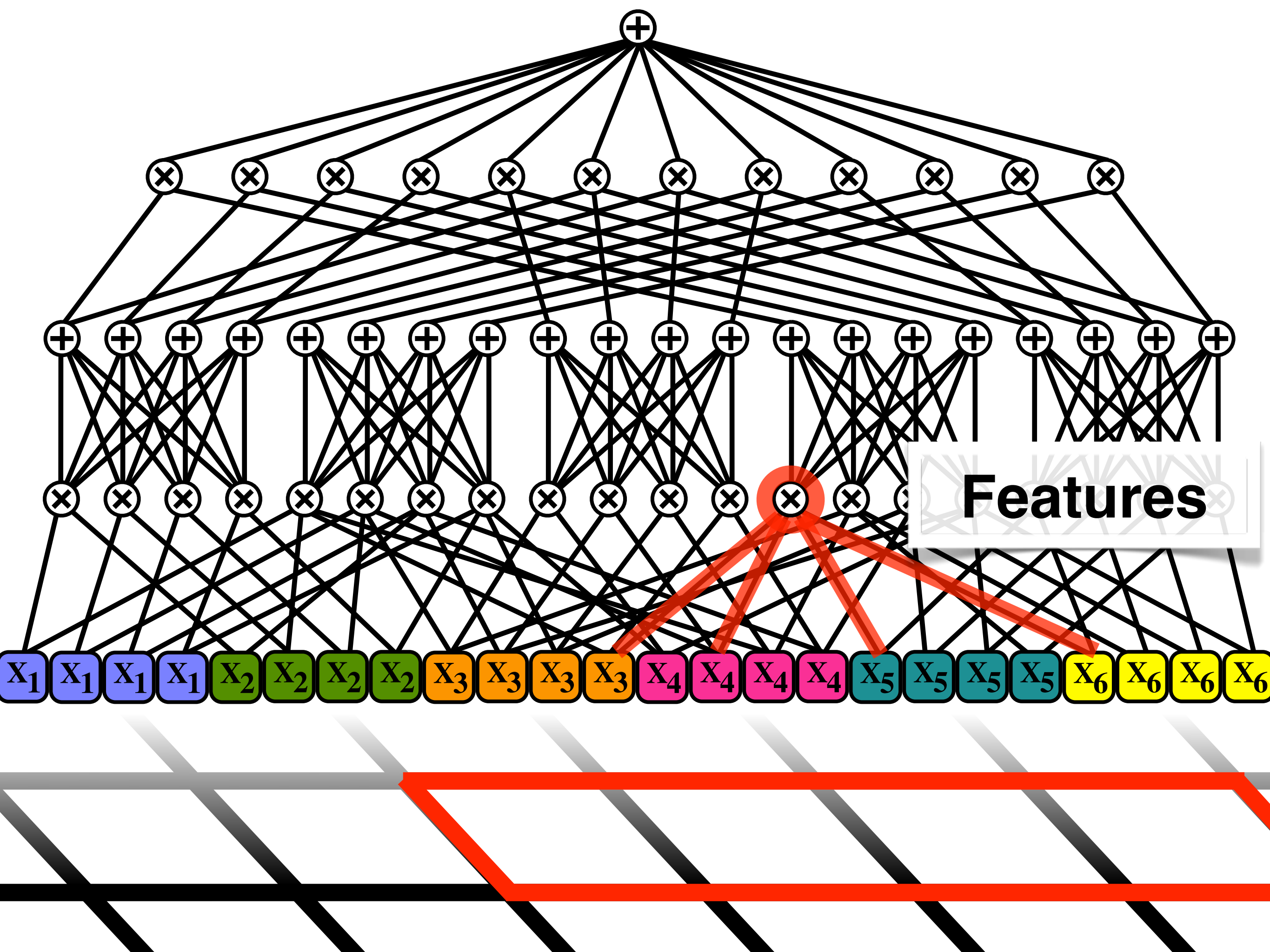# Discriminative Learning of Sum-Product Networks

Robert Gens
Pedro Domingos

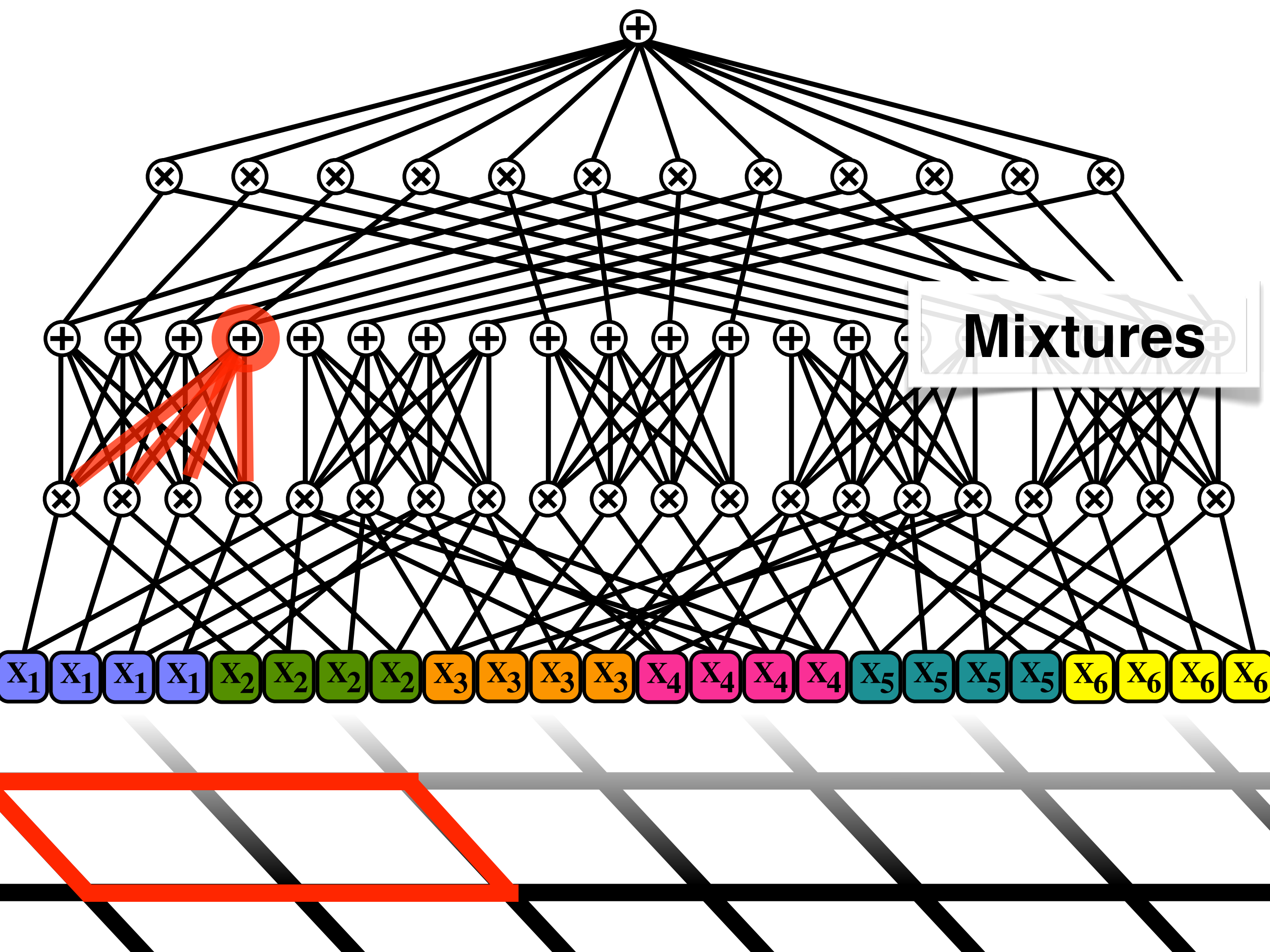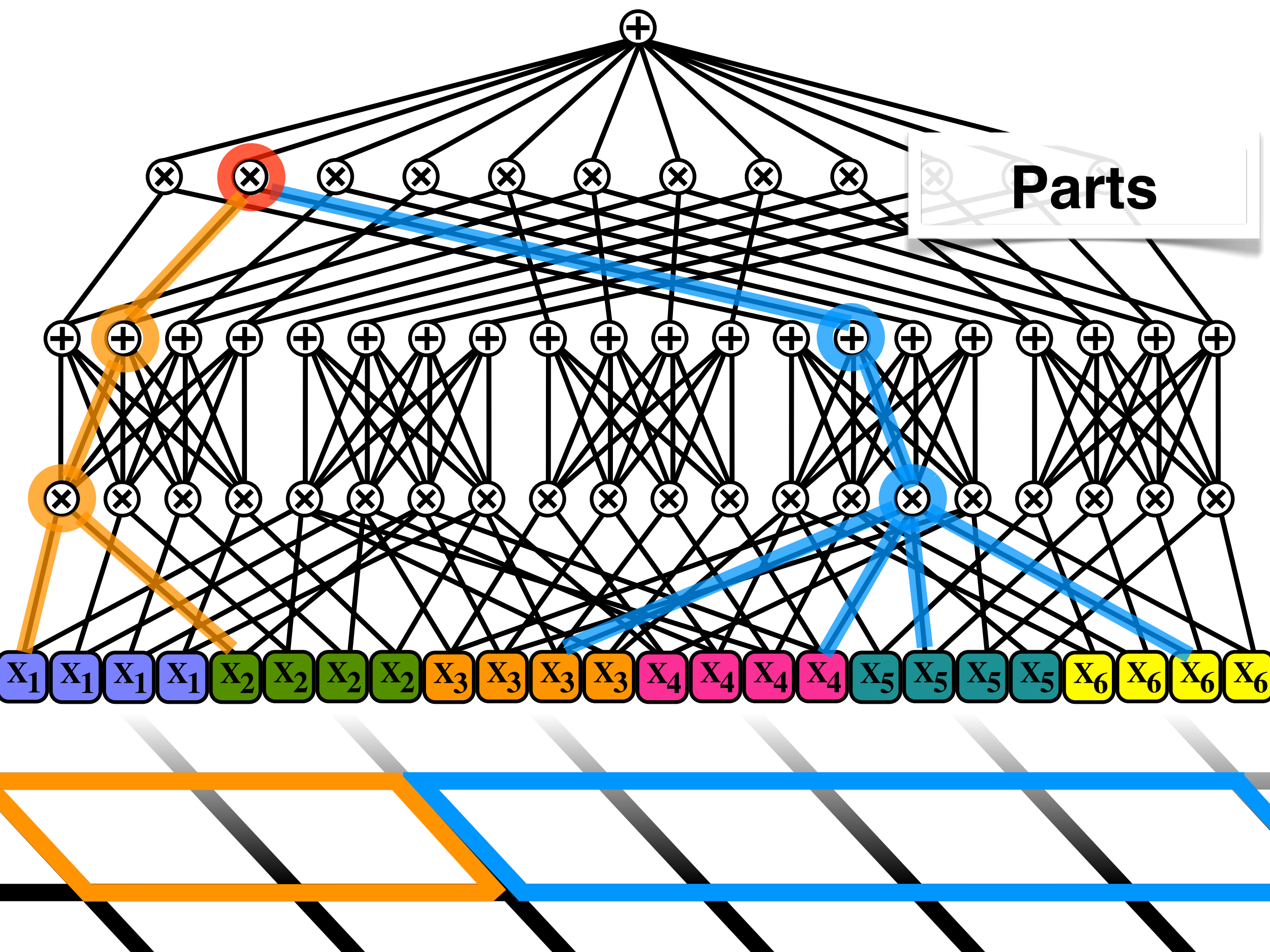UNIVERSITY *of* WASHINGTON

University of Washington
Computer Science & Engineering

Distributions

Features

Features

**Mixtures**

$X_1$ $X_1$ $X_1$ $X_1$ $X_2$ $X_2$ $X_2$ $X_2$ $X_3$ $X_3$ $X_3$ $X_3$ $X_4$ $X_4$ $X_4$ $X_4$ $X_5$ $X_5$ $X_5$ $X_5$ $X_6$ $X_6$ $X_6$ $X_6$

Parts

Parts

Pooling

**Motivation**   **SPN Review**   **Discriminative Training**   **Experiments**

**Motivation**

SPN Review

Discriminative Training

Experiments

# Graphical Models



➡️ SPNs perform fast, exact inference on high treewidth models
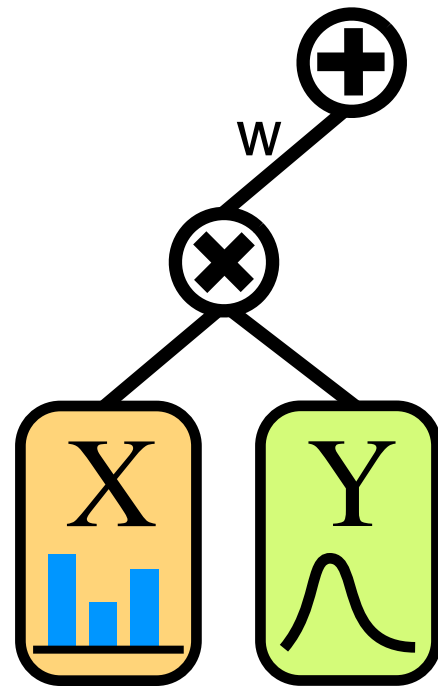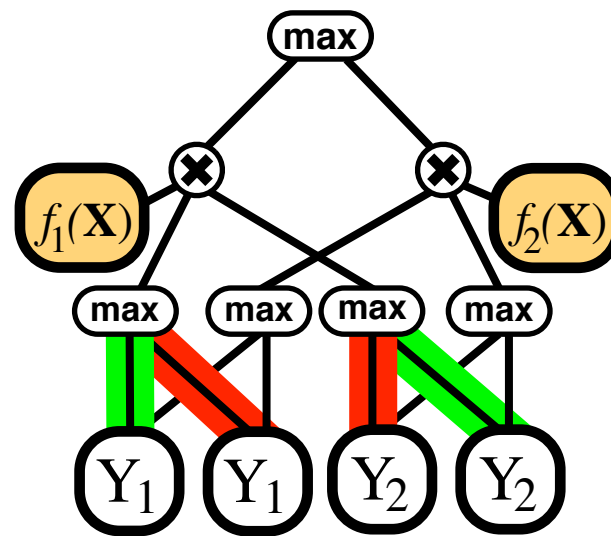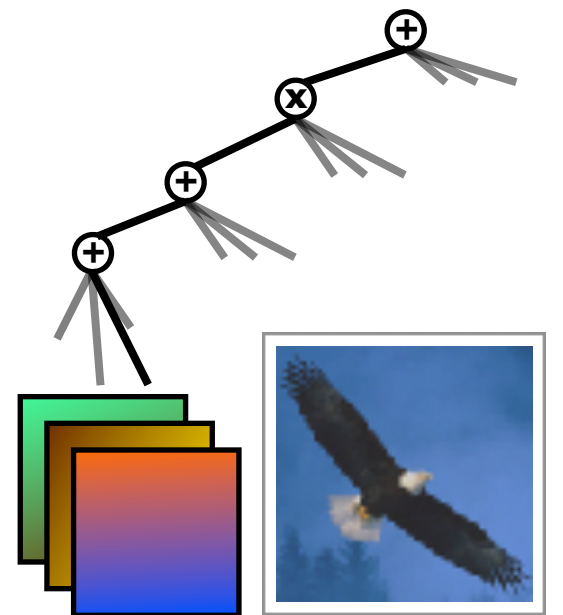
# Deep Architectures



➡️SPNs have full probabilistic semantics and tractable inference over many layers

# Discriminative Learning

SPNs combine features with fast, exact inference over high treewidth models

NIPS'12

Motivation

**SPN Review**

Discriminative Training

Experiments

A Product of SPNs over Disjoint Variables Is an SPN.

# A Weighted Sum of SPNs over the Same Variables Is an SPN.
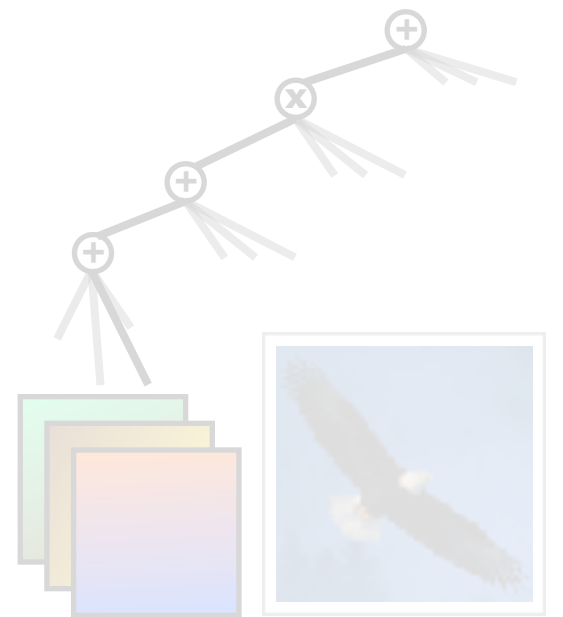


Sums out a mixture variable

# All Marginals Are Computable in Linear Time

# All Marginals Are Computable in Linear Time

# All Marginals Are Computable in Linear Time

$P(X{=}0)$ ?

# All Marginals Are Computable in Linear Time

$$P(X{=}0)\ ?$$

# All Marginals Are Computable in Linear Time

# All MAP States Are Computable in Linear Time

$$\max_y P(X{=}0, Y{=}y) \ ?$$

# All MAP States Are Computable in Linear Time

# All MAP States Are Computable in Linear Time

$$\max_{y} P(X{=}0, Y{=}y) \text{ ?}$$

# All MAP States Are Computable in Linear Time

$$\max_{y} P(X{=}0, Y{=}y) \text{ ?}$$

# All MAP States Are Computable in Linear Time

$$\max_{y} P(X{=}0, Y{=}y) = \boxed{0.12}$$

# All MAP States Are Computable in Linear Time

$$\max_y P(X{=}0, Y{=}y) = \boxed{0.12}$$

# Special Cases of SPNs

- Junction trees

- Hierarchical mixture models

- Non-recursive probabilistic context-free grammars

- Models with context-specific independence

- Models with determinism

- Other high-treewidth models

# Learning SPNs

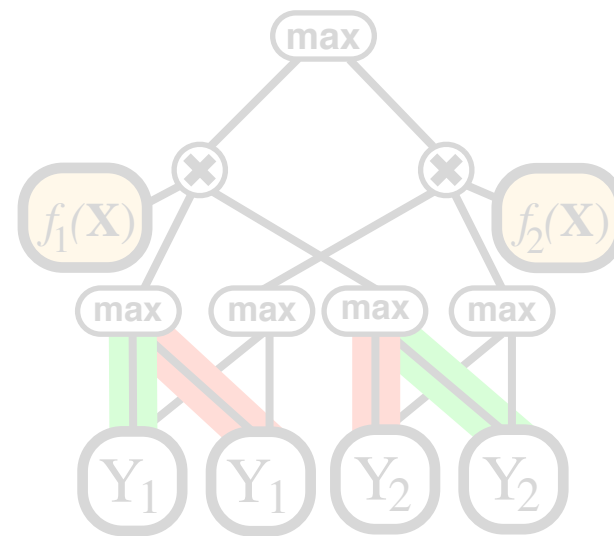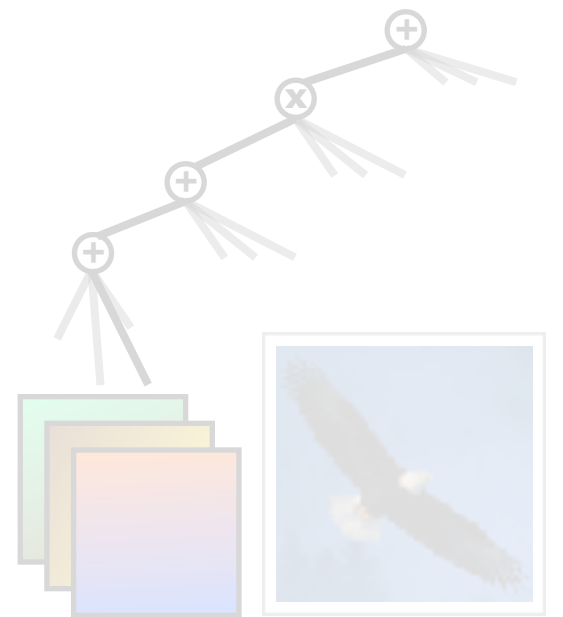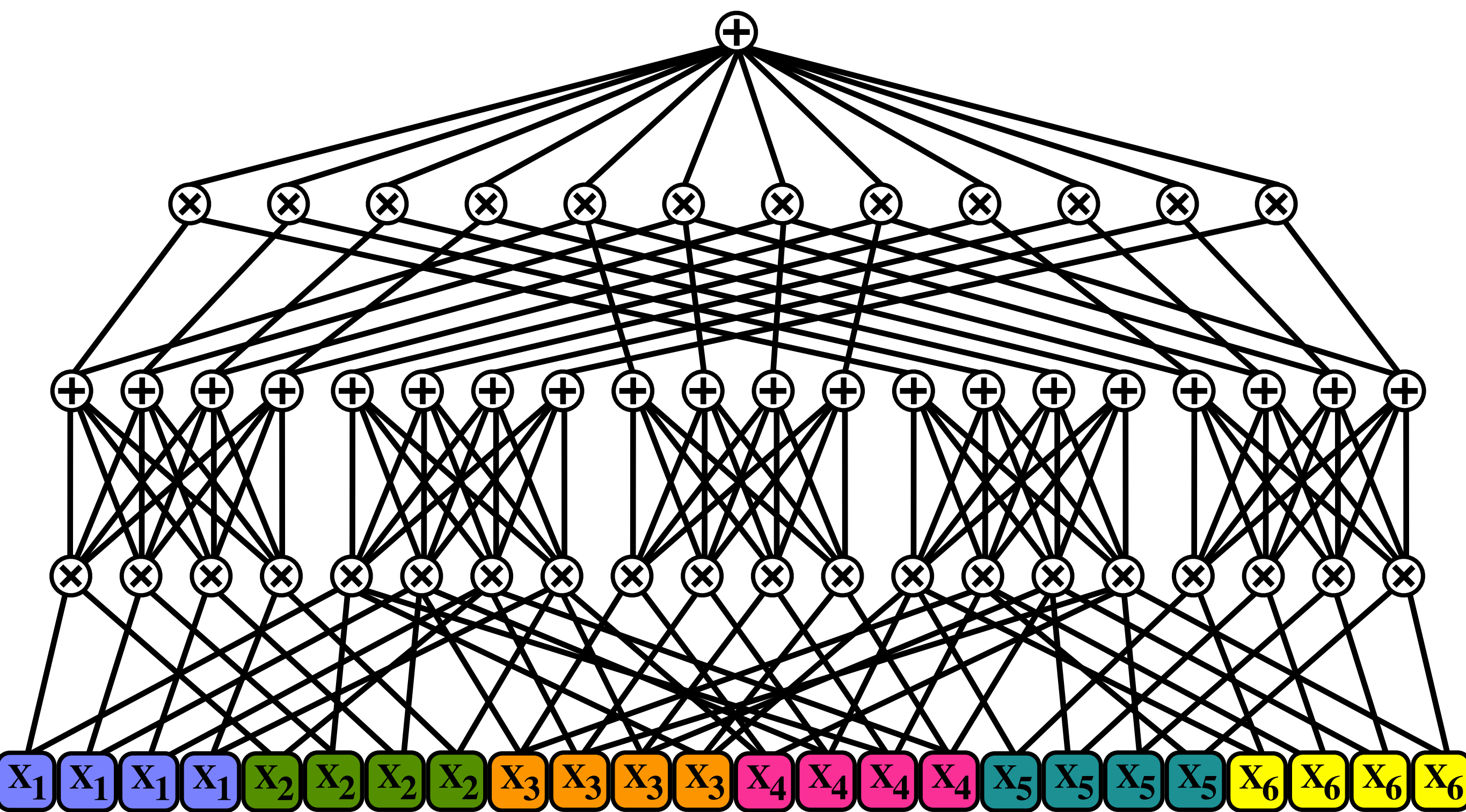| Update | Soft Inference (Marginals) | Hard Inference (MAP States) |
|---|---|---|
| Gen. EM | ✅ | ✅ |
| Gen. Gradient | ✅ | |
| Disc. Gradient | | |

Poon & Domingos, UAI 2011

Motivation

SPN
Review

**Discriminative
Training**

Experiments
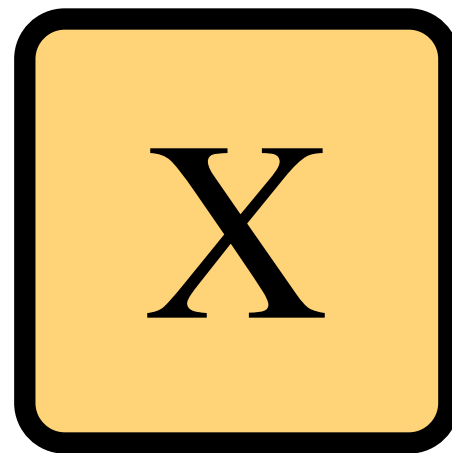
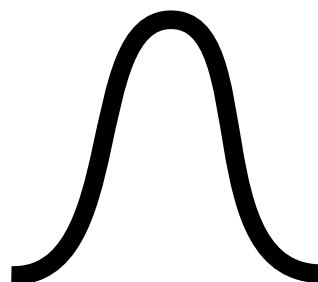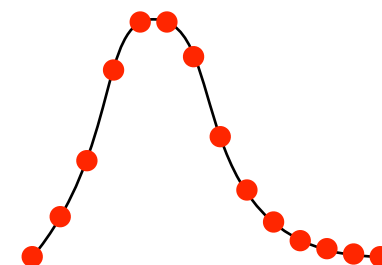# Discriminative SPNs

$$P(\mathbf{Y}|\mathbf{X})$$

$\mathbf{Y}$    Query

$\mathbf{H}$    Hidden

$\mathbf{X}$    Evidence

Treat as constants

# Discriminative SPNs

# Discriminative SPNs

$\mathbb{H}$  Hidden

$\mathbb{Y}$  Query

$f(\mathbf{X})$  Features (non-negative)

$\mathbb{X}$  Evidence

# Discriminative SPNs

**0.8**     0.1     **1.5**

$Y_1$  $Y_1$  $Y_2$  $Y_2$

$f_1(\mathbf{X})$     $f_2(\mathbf{X})$

$\mathbb{H}$  Hidden

$\mathbb{Y}$  Query

$f(\mathbf{X})$  Features
(non-negative)

$\mathbb{X}$  Evidence

# Discriminative SPNs

Greater variety than generative SPNs

Y    Query

$f(\mathbf{X})$    Features (non-negative)

X    Evidence

# Discriminative Training

$$\nabla \log P(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{P(\mathbf{y}, \mathbf{x})}{P(\mathbf{x})} =$$

$$\nabla \log \sum_{\mathbf{h}} P(\mathbf{y}, \mathbf{h}, \mathbf{x}) - \nabla \log \sum_{\mathbf{y}', \mathbf{h}} P(\mathbf{y}', \mathbf{h}, \mathbf{x})$$

**Tractable!**

Correct label

Best guess

# SPN Backpropagation

# SPN Backpropagation

# SPN Backpropagation

# SPN Backpropagation

# SPN Backpropagation

# SPN Backpropagation



For each child $j$:

$$\frac{\partial S}{\partial S_j} \leftarrow \frac{\partial S}{\partial S_j} + w_{n,j}\frac{\partial S}{\partial S_n}$$

$$\frac{\partial S}{\partial w_{n,j}} \leftarrow S_j\frac{\partial S}{\partial S_n}$$

# SPN Backpropagation



0.3

.24

.3

?

0.9

$f_1(\mathbf{X})$

0.6

.45

For each child $j$:
$$\frac{\partial S}{\partial S_j} \leftarrow \frac{\partial S}{\partial S_j} + \frac{\partial S}{\partial S_n} \prod_{k \in Ch(n) \setminus \{j\}} S_k$$

# SPN Backpropagation

0.3

.3

.24

.1

0.9

$f_1(\mathbf{X})$

0.6

.1

.2

.45

For each child $j$:

$$\frac{\partial S}{\partial S_j} \leftarrow \frac{\partial S}{\partial S_j} + \frac{\partial S}{\partial S_n} \prod_{k \in Ch(n) \setminus \{j\}} S_k$$

# Problem with Backpropagation



Gradient diffusion

# Hard Inference Overcomes Gradient Diffusion



**Soft Inference**
(Marginals)

**Hard Inference**
(MAP States)

# Reasons to Use Hard Inference

- To overcome gradient diffusion

- When goal is to predict most probable structure

- For speed or tractability

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$

$$\nabla \log \max_{\mathbf{h}} P(\mathbf{y}, \mathbf{h}, \mathbf{x}) - \nabla \log \max_{\mathbf{y}', \mathbf{h}} P(\mathbf{y}', \mathbf{h}, \mathbf{x})$$

Correct label

Best guess

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$

$$\nabla \log \max_{\mathbf{h}} P(\mathbf{y}, \mathbf{h}, \mathbf{x}) \ \text{-} \ \nabla \log \max_{\mathbf{y}', \mathbf{h}} P(\mathbf{y}', \mathbf{h}, \mathbf{x})$$

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$

$$\nabla \log \left( \max_{\mathbf{h}} P(\mathbf{y}, \mathbf{h}, \mathbf{x}) \right) - \nabla \log \left( \max_{\mathbf{y}', \mathbf{h}} P(\mathbf{y}', \mathbf{h}, \mathbf{x}) \right)$$

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y},\mathbf{x})}{\tilde{P}(\mathbf{x})} =$$



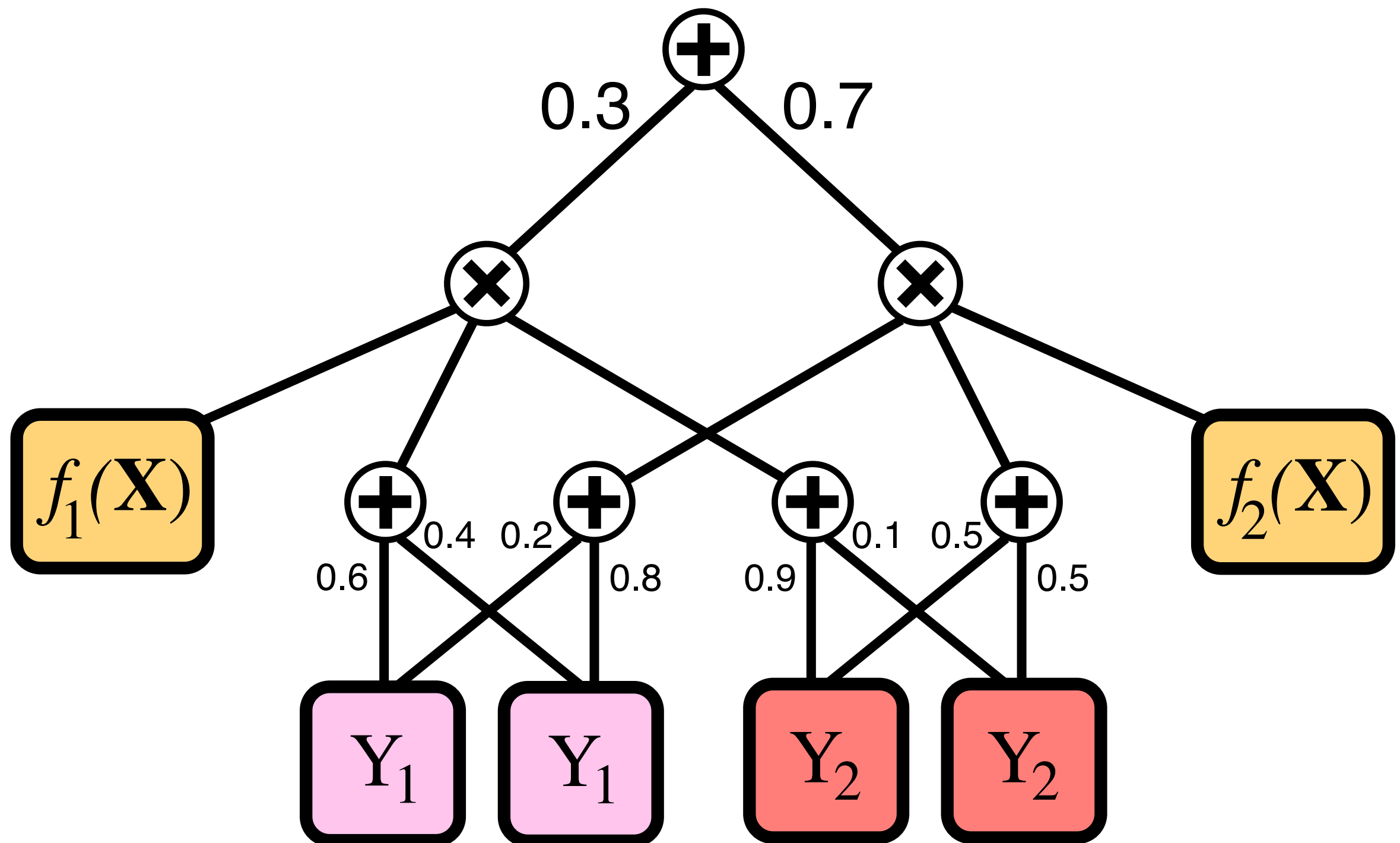$$\nabla \log \left( \max_{\mathbf{h}} P(\mathbf{y},\mathbf{h},\mathbf{x}) \right) - \nabla \log \left( \max_{\mathbf{y}',\mathbf{h}} P(\mathbf{y}',\mathbf{h},\mathbf{x}) \right)$$

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$



$$\nabla \log \left( \begin{array}{c} \max \\ \mathbf{h} \end{array} P(\mathbf{y}, \mathbf{h}, \mathbf{x}) \right) - \nabla \log \left( \begin{array}{c} \max \\ \mathbf{y}', \mathbf{h} \end{array} P(\mathbf{y}', \mathbf{h}, \mathbf{x}) \right)$$

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$

# Hard Gradient

$$\nabla \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \nabla \log \frac{\tilde{P}(\mathbf{y}, \mathbf{x})}{\tilde{P}(\mathbf{x})} =$$



# w/ correct label $-$ # w/ model guess

$$\frac{\partial}{\partial w_i} \log \tilde{P}(\mathbf{y}|\mathbf{x}) = \frac{\Delta c_i}{w_i}$$

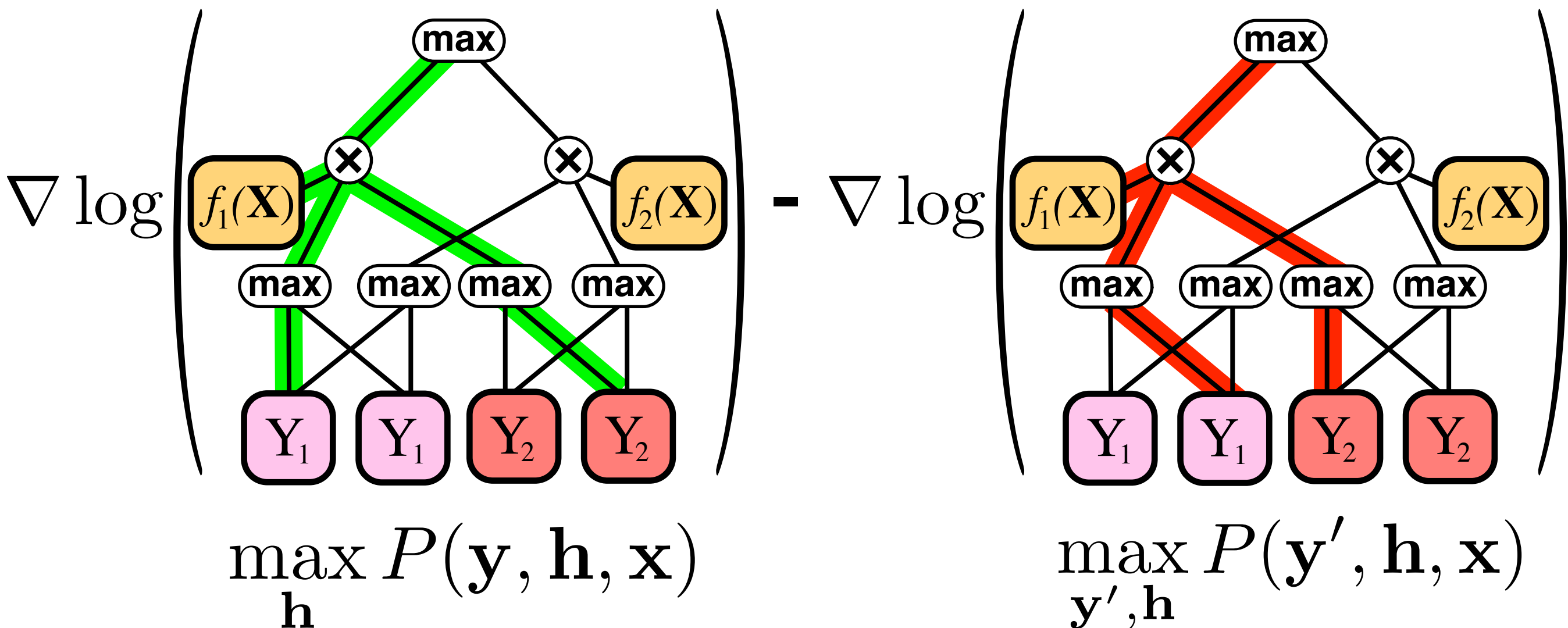# Learning SPNs: Summary



| Update | Soft Inference (Marginals) | Hard Inference (MAP States) |
|---|---|---|
| Gen. EM | $\Delta w_i \propto w_i \dfrac{\partial S}{\partial S_k}$ | $\Delta w_i = c_i$ |
| Gen. Gradient | $\Delta w_i = \eta \dfrac{\partial S}{\partial S_k} S_i$ | $\Delta w_i = \eta \dfrac{c_i}{w_i}$ |
| Disc. Gradient | $\Delta w_i = \eta \left( \overbrace{\dfrac{S_i}{S} \dfrac{\partial S}{\partial S_k}}^{\text{true label}} - \overbrace{\dfrac{S_i}{S} \dfrac{\partial S}{\partial S_k}}^{\text{exp. label}} \right)$ | $\Delta w_i = \dfrac{\eta}{w_i} ( \overbrace{c_i}^{\text{true}} - \overbrace{c_i}^{\text{test}} )$ |

# Learning SPNs: Summary



| Update | Soft Inference (Marginals) | Hard Inference (MAP States) |
|---|---|---|
| Gen. EM | $\Delta w_i \propto w_i \dfrac{\partial S}{\partial S_k}$ | $\Delta w_i = c_i$ |
| Gen. Gradient | $\Delta w_i = \eta \dfrac{\partial S}{\partial S_k} S_i$ | $\Delta w_i = \eta \dfrac{c_i}{w_i}$ |
| Disc. Gradient | $\Delta w_i = \eta \left( \overbrace{\dfrac{S_i}{S} \dfrac{\partial S}{\partial S_k}}^{\text{true label}} - \overbrace{\dfrac{S_i}{S} \dfrac{\partial S}{\partial S_k}}^{\text{exp. label}} \right)$ | $\Delta w_i = \dfrac{\eta}{w_i} \left( \overbrace{c_i}^{\text{true}} - \overbrace{c_i}^{\text{test}} \right)$ |

# Learning SPNs: Summary
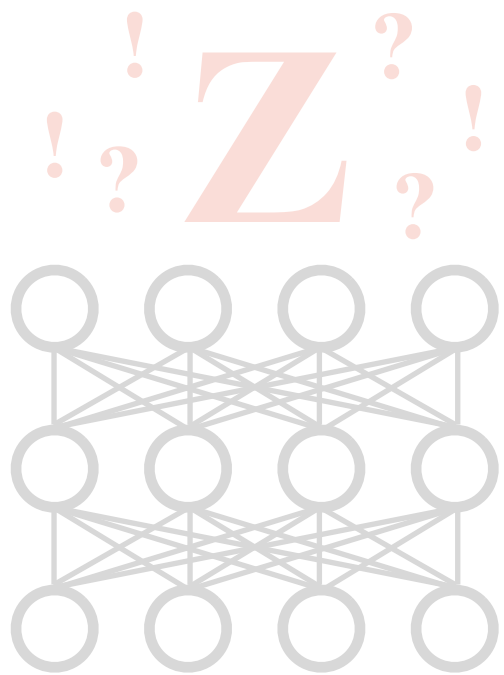


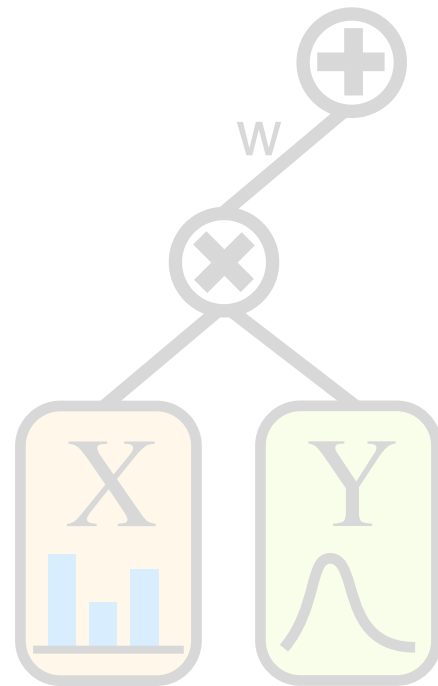| **Update** | **Soft Inference** (Marginals) | **Hard Inference** (MAP States) |
|---|---|---|
| Gen. EM | $\Delta w_i \propto w_i \dfrac{\partial S}{\partial S_k}$ | $\Delta w_i = c_i$ |
| Gen. Gradient | $\Delta w_i = \eta \dfrac{\partial S}{\partial S_k} S_i$ | $\Delta w_i = \eta \dfrac{c_i}{w_i}$ |
| Disc. Gradient | $\Delta w_i = \eta \left( \overbrace{\dfrac{S_i}{S}\dfrac{\partial S}{\partial S_k}}^{\text{true label}} - \overbrace{\dfrac{S_i}{S}\dfrac{\partial S}{\partial S_k}}^{\text{exp. label}} \right)$ | $\Delta w_i = \dfrac{\eta}{w_i}\left( \overbrace{c_i}^{\text{true}} - \overbrace{c_i}^{\text{test}} \right)$ |

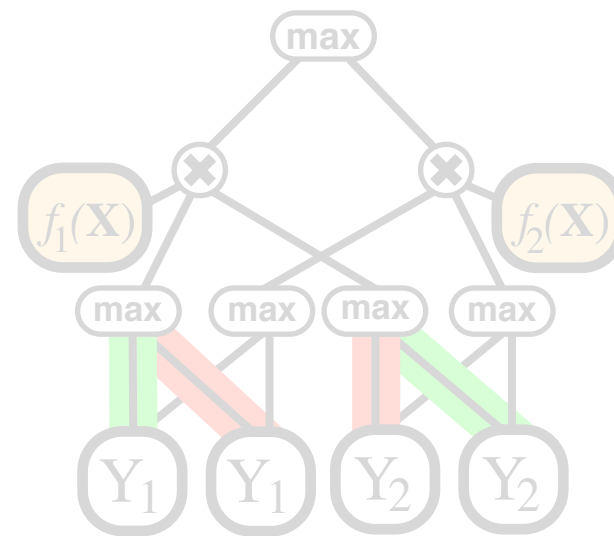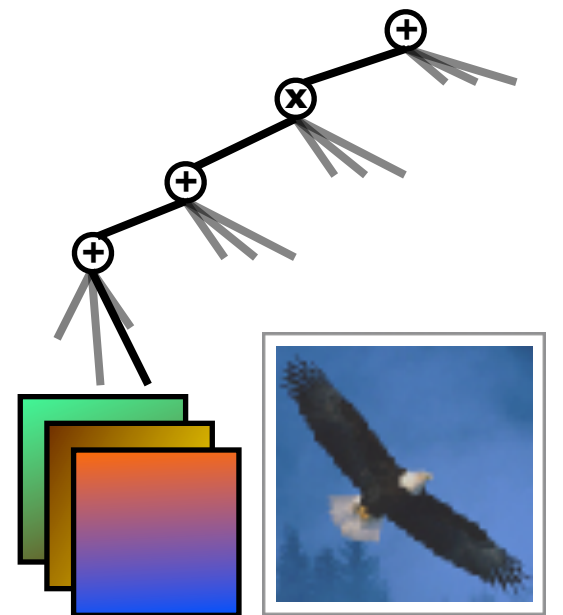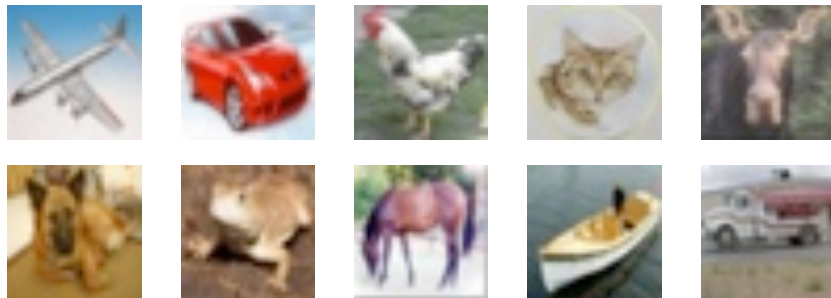**Motivation**  **SPN Review**  **Discriminative Training**  **Experiments**
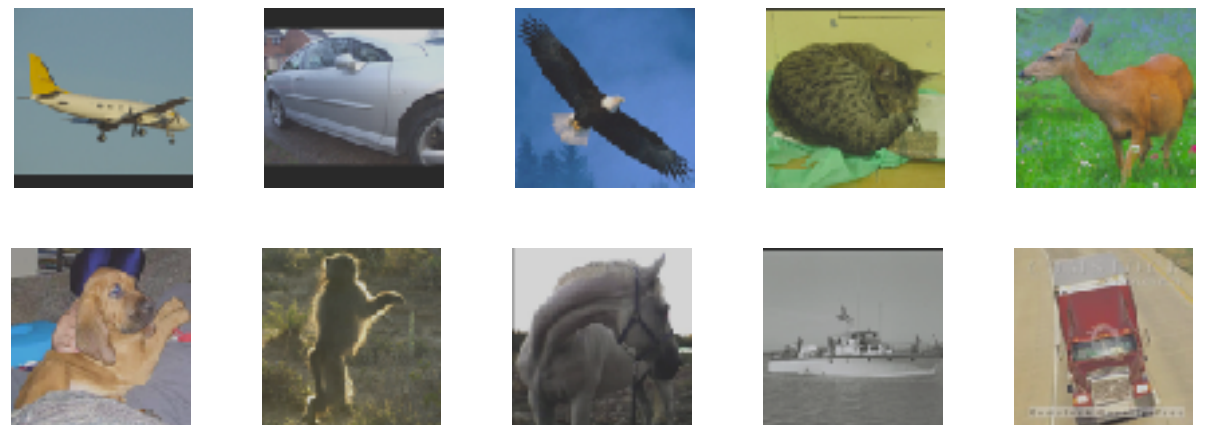
# Image Classification



## CIFAR-10

32x32px
50k train
10k test

## STL-10

96x96px
5k train
8k test  } 10 folds
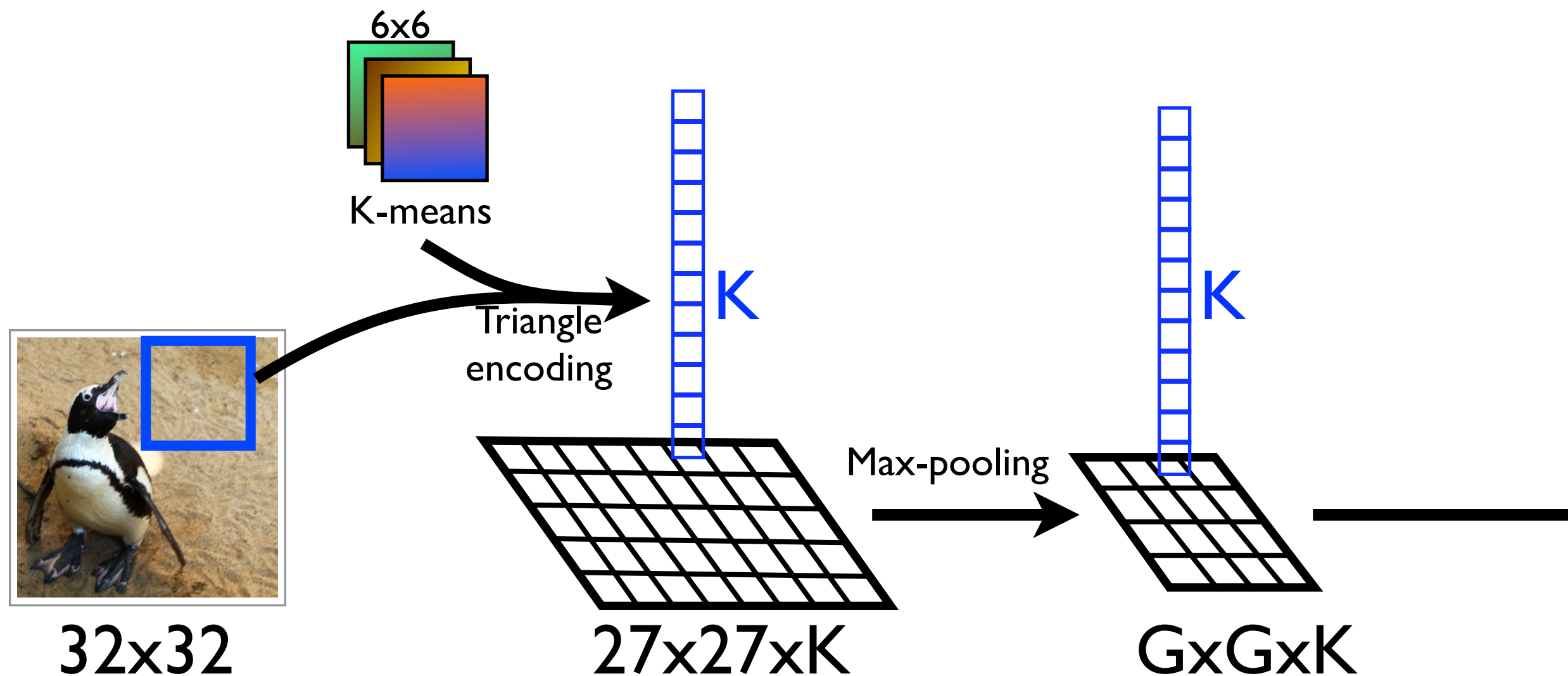~~100k unlabeled~~
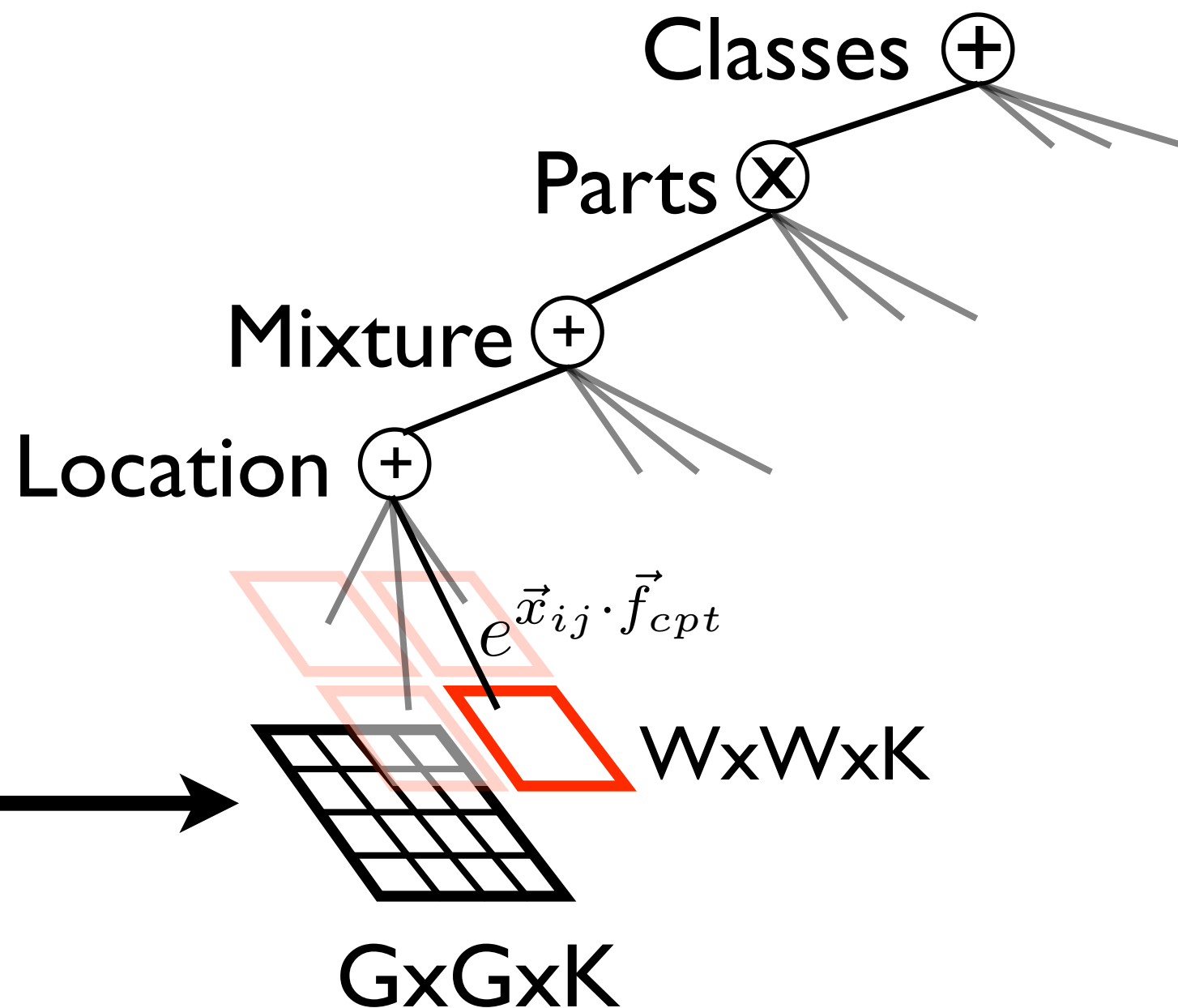
# Feature Extraction

Coates et al., AISTATS 2011

# SPN Architecture

Classes ⊕

Parts ⊗

Mixture ⊕

Location ⊕

$e^{\vec{x}_{ij} \cdot \vec{f}_{cpt}}$

WxWxK

GxGxK
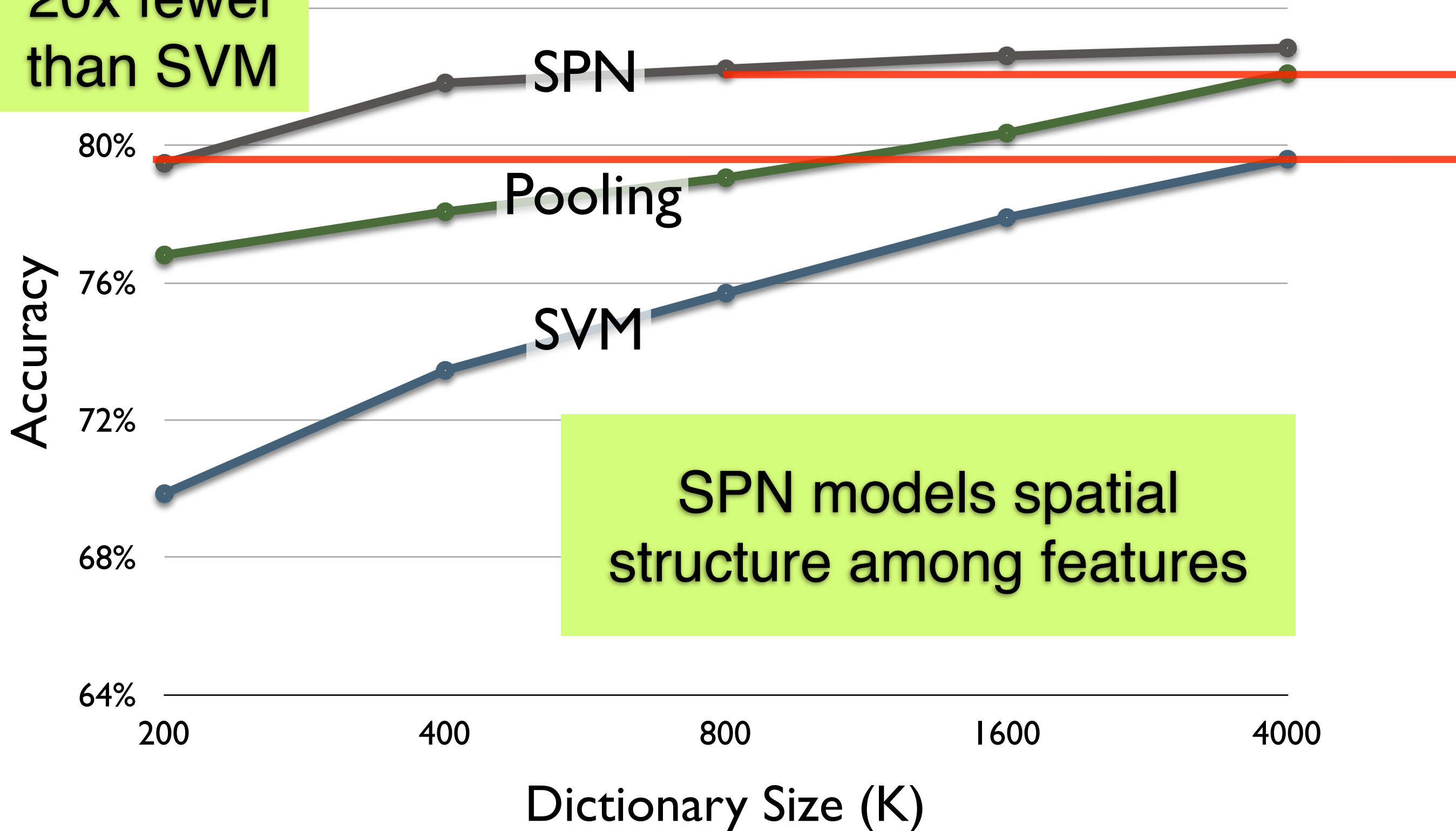
CIFAR-10 Results

# CIFAR-10 Results

20x fewer than SVM

SPN

Pooling

SVM

SPN models spatial structure among features

Accuracy

80%

76%

72%

68%

64%

200    400    800    1600    4000

Dictionary Size (K)

# CIFAR-10 Results

7x7x400

SPN

Best published

Learned Pooling

3-Layer Learned RF

SVM

Accuracy

84%

83%

82%

81%

80%

79%

0k    38k    75k    113k    150k

#Features

# STL-10 results



| | |
|---|---|
| 1-layer Vector Quantization | 54.9% |
| 1-layer Sparse Coding | 59.0% |
| 3-layer Learned Receptive Field | 60.1% |
| Discriminative SPN | 62.3% |

without unlabeled data

52%    55%    58%    61%    64%

# Future Work

- Max-margin SPNs

- Learning SPN structure

- Applying discriminative SPNs to structured prediction

- Approximate inference using SPNs

# Summary

- Discriminative SPNs combine the advantages of
  - Tractable inference
  - Deep architectures
  - Discriminative learning

- Hard gradient combats diffusion in deep models

- Discriminative SPNs outperform SVMs and deep models on image classification benchmarks