

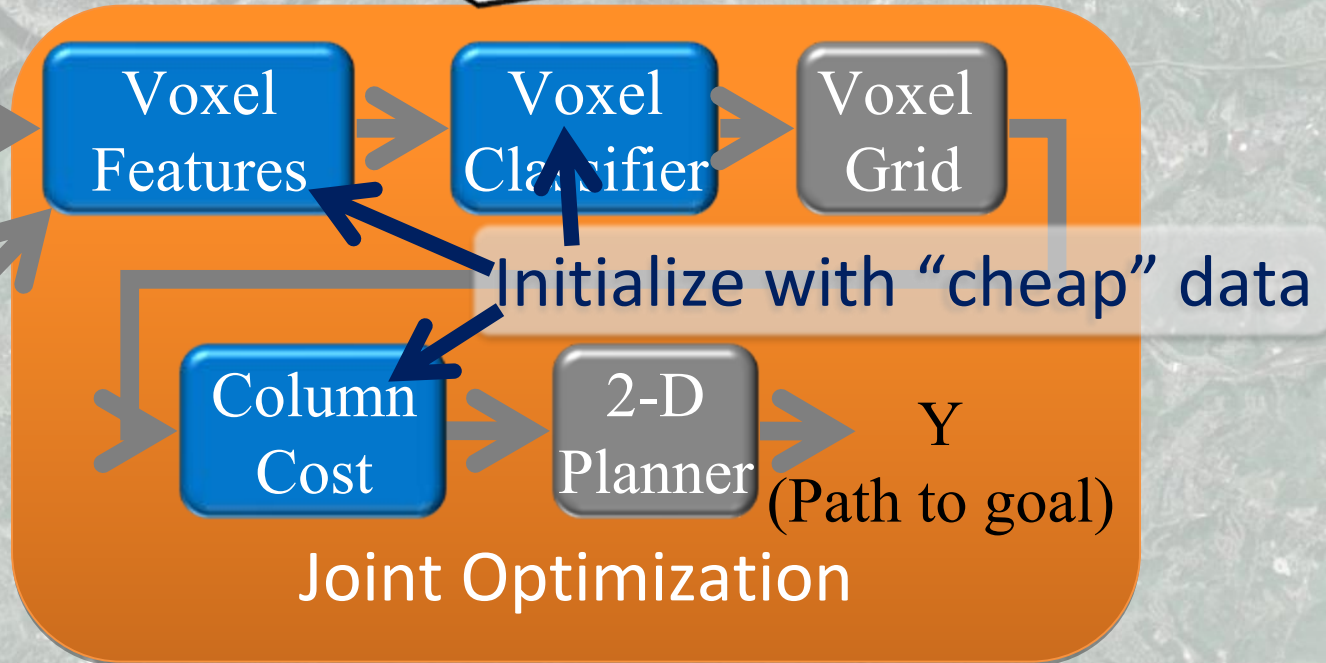
# Differentiable Sparse Coding

David Bradley and J. Andrew Bagnell

NIPS 2008

# 100,000 ft View

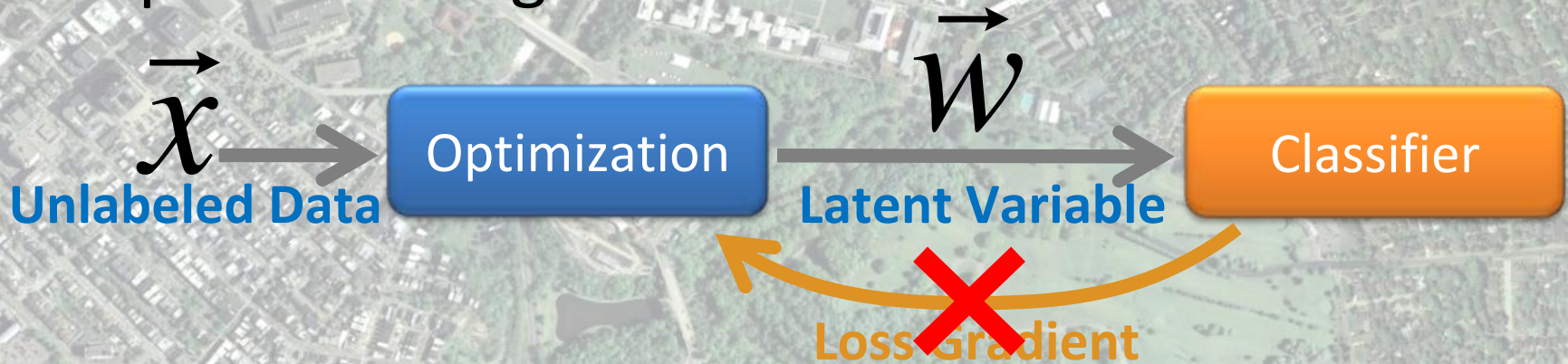
- Complex Systems





# 10,000 ft view

- Sparse Coding = Generative Model



- Semi-supervised learning
- **KL-Divergence Regularization**
- **Implicit Differentiation**

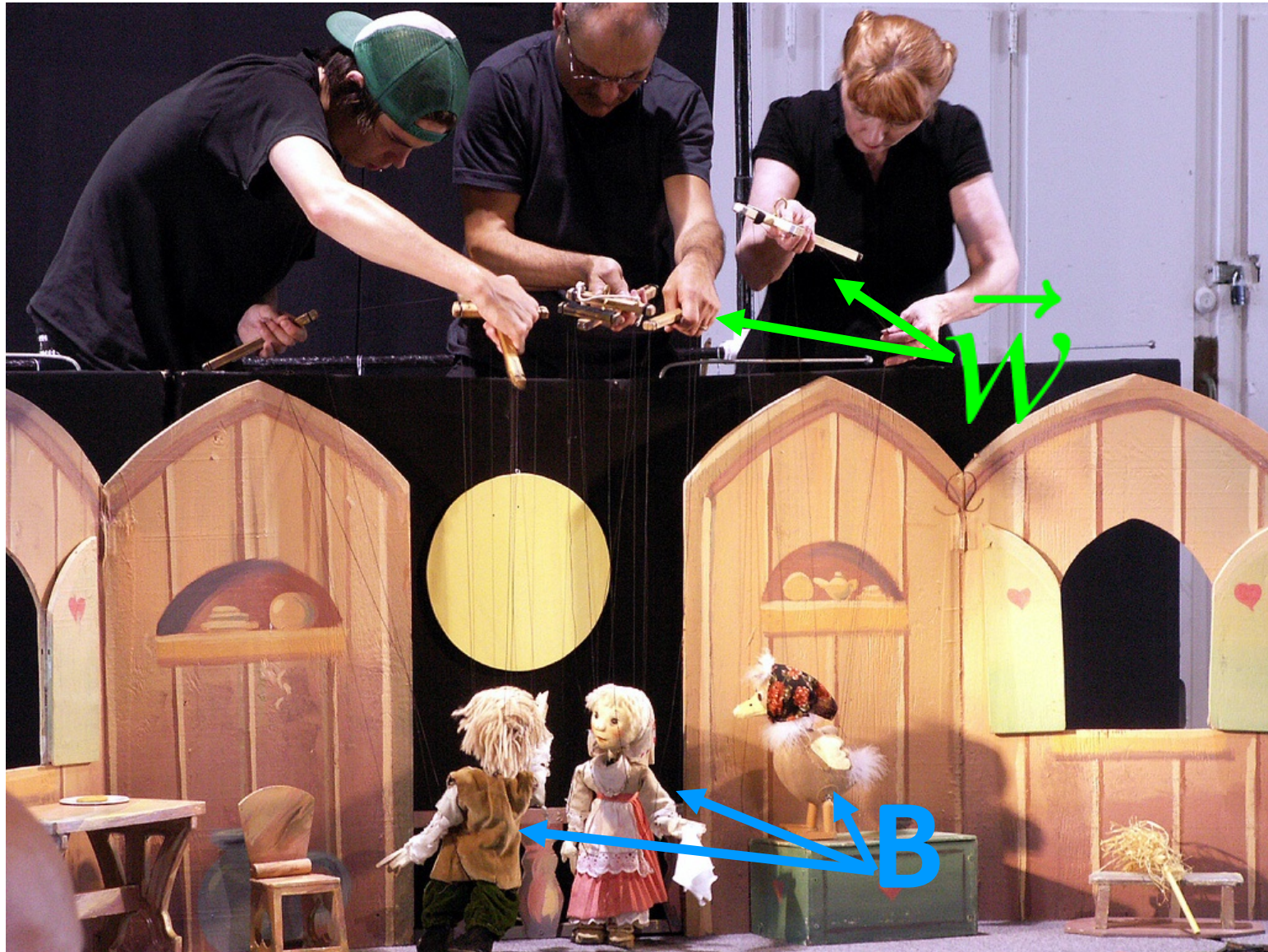
# Sparse Coding



Understand X



# As a combination of factors



# Sparse coding uses optimization

Reconstruction

loss function

$$\vec{x} \approx f(B\vec{w})$$

Some vector

Want to use  
to classify  $x$

Projection (feed-forward):

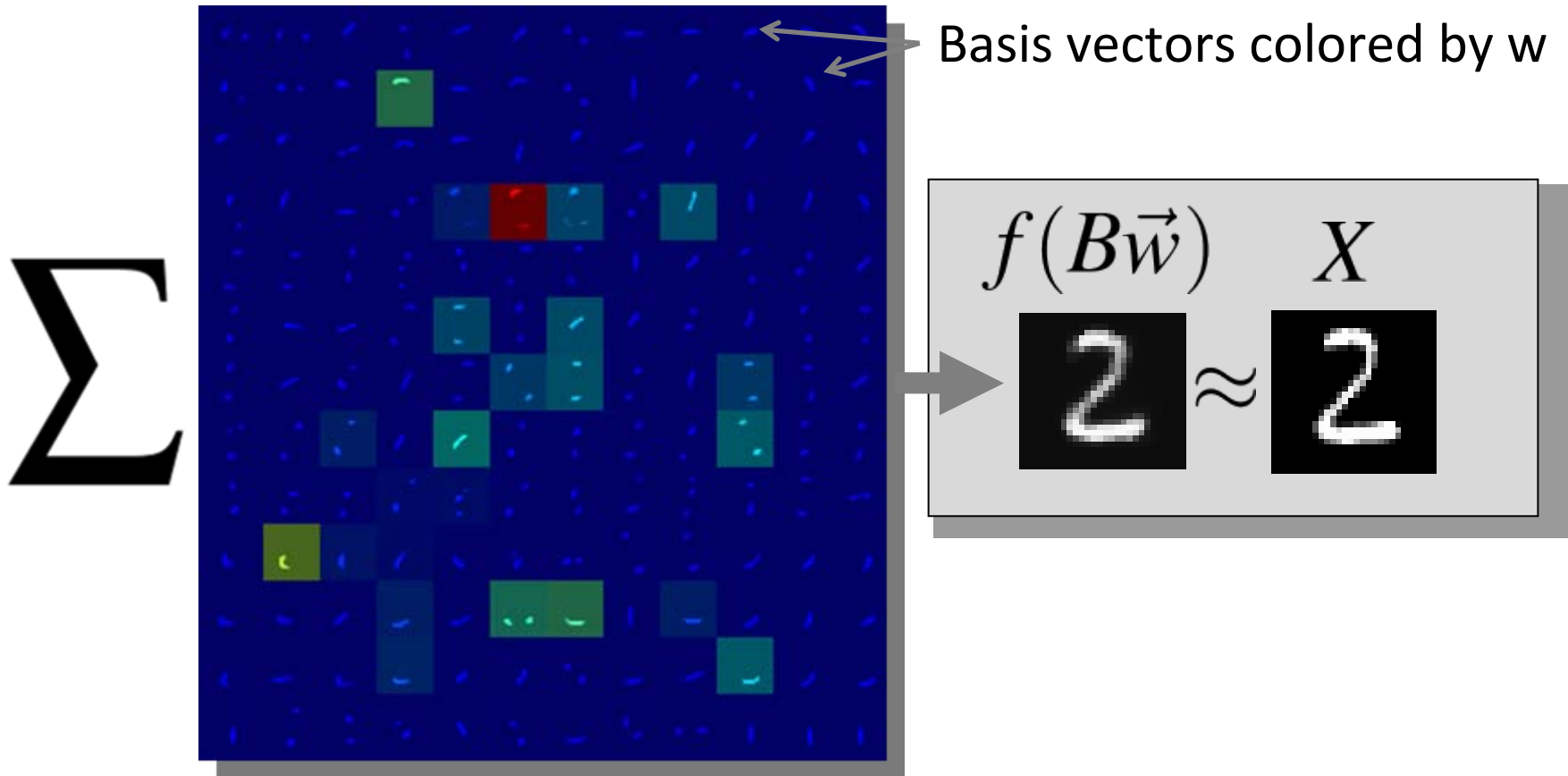
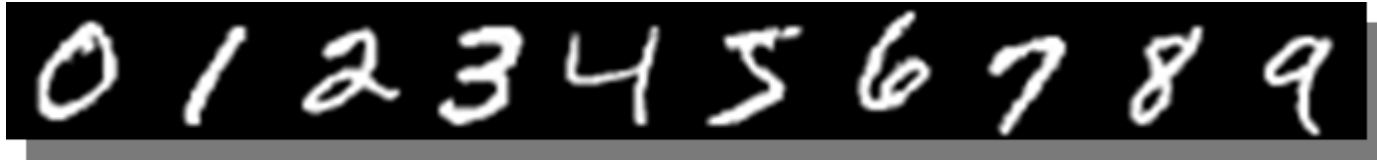
$$\vec{w} = f(\vec{x}^T B)$$



Sparse vectors:  
Only a few significant elements



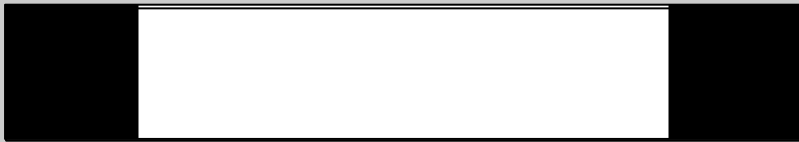
# Example: X=Handwritten Digits



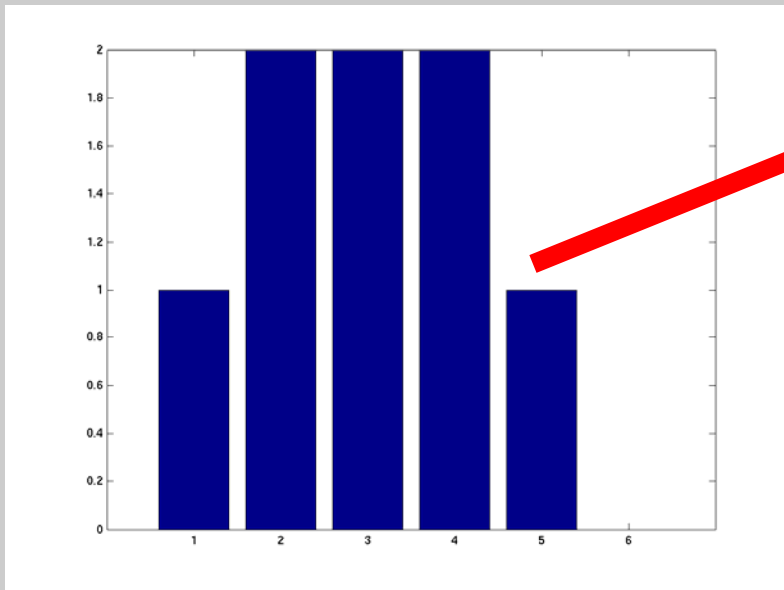


# Optimization vs. Projection

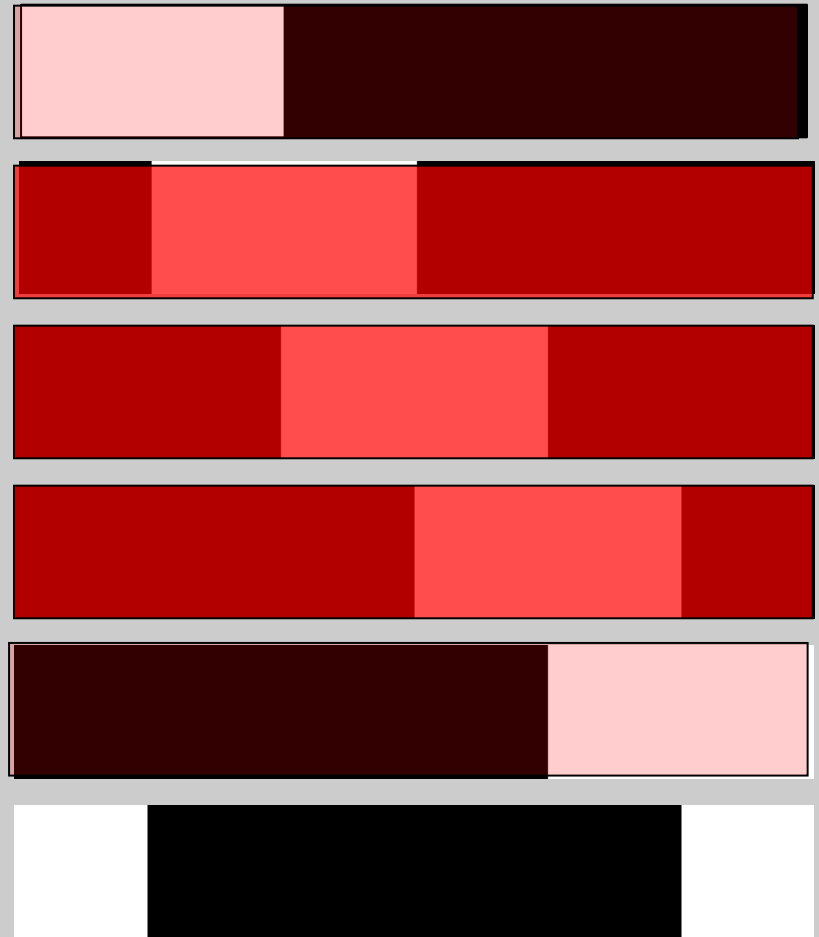
Input



Projection

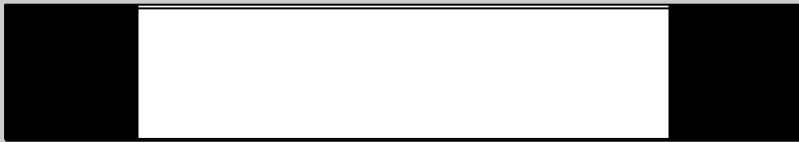


Basis



# Optimization vs. Projection

Input

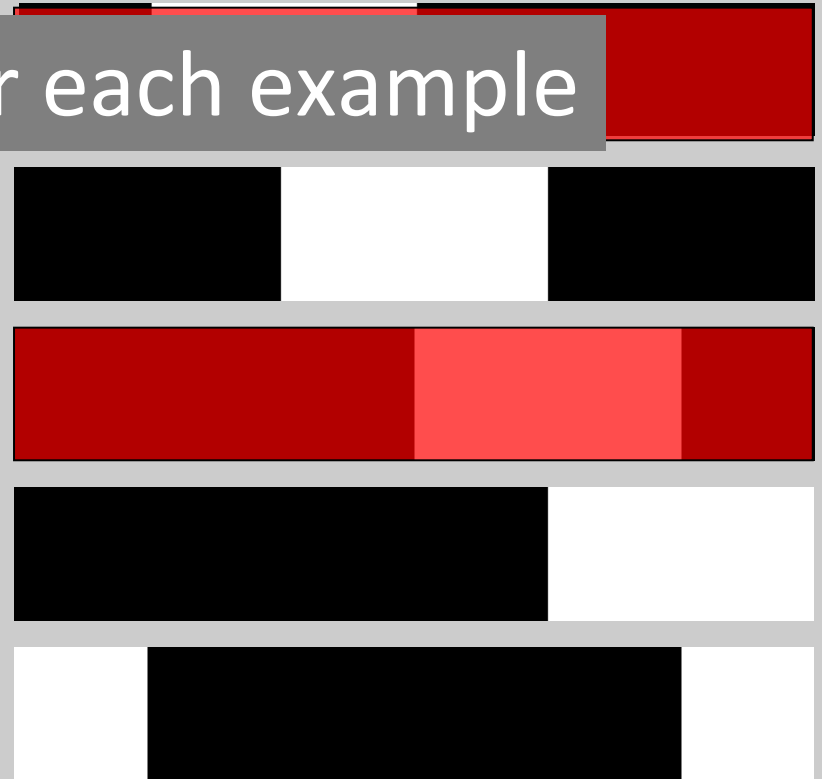
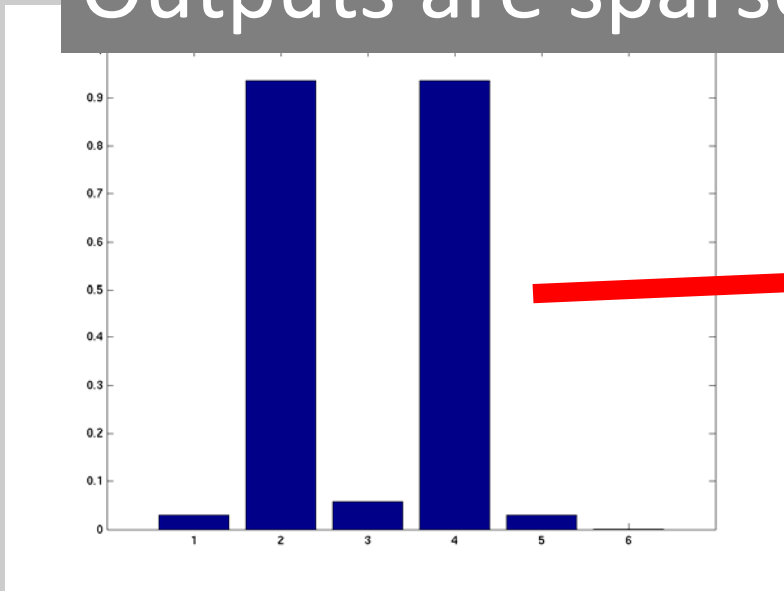


Basis



KL

Outputs are sparse for each example



# Generative Model

$$X_i = \text{2}$$

Latent variables are Independent

$$P(X) = \int_B \int_W P(X|W, B) P(W) P(B) dW dB$$

$$P(X) = \int_B P(B) \int_W \prod_i \text{Likelihood} P(X_i|W_i, B) \text{Prior} P(W_i) dW dB$$

Examples are Independent



# Sparse Approximation

$$X = \text{2}$$

$$P(\vec{x}) \approx \max_{\vec{w}} \overset{\text{Likelihood}}{P(\vec{x}|\vec{w}, B)} \overset{\text{Prior}}{P(\vec{w})}$$

$$\hat{w} = \arg \max_{\vec{w}} P(\vec{x}|\vec{w}, B) P(\vec{w})$$

MAP Estimate

# Sparse Approximation

$$\arg \max_{\vec{w}} P(\vec{x}) = \arg \min_{\vec{w}} (-\log P(\vec{x}))$$

Distance between  
reconstruction and input

Distance between weight  
vector and prior mean

$$-\log P(\vec{x}) \propto \text{Loss}(\vec{x} \parallel f(B\vec{w})) + \lambda \text{Prior}(\vec{w} \parallel \vec{p})$$

Regularization Constant

# Example: Squared Loss + L1

$$\hat{w} = \arg \min_{\vec{w}} \sum_i (x_i - B^i \vec{w})^2 + \lambda \sum_i |\vec{w}_i|$$

- Convex + sparse (widely studied in engineering)
- Sparse coding solves for B as well (non-convex for now...)
- Shown to generate useful features on diverse problems

Tropp, *Signal Processing*, 2006

Donoho and Elad, *Proceedings of the National Academy of Sciences*, 2002

Raina, Battle, Lee, Packer, Ng, ICML, 2007



# L1 Sparse Coding

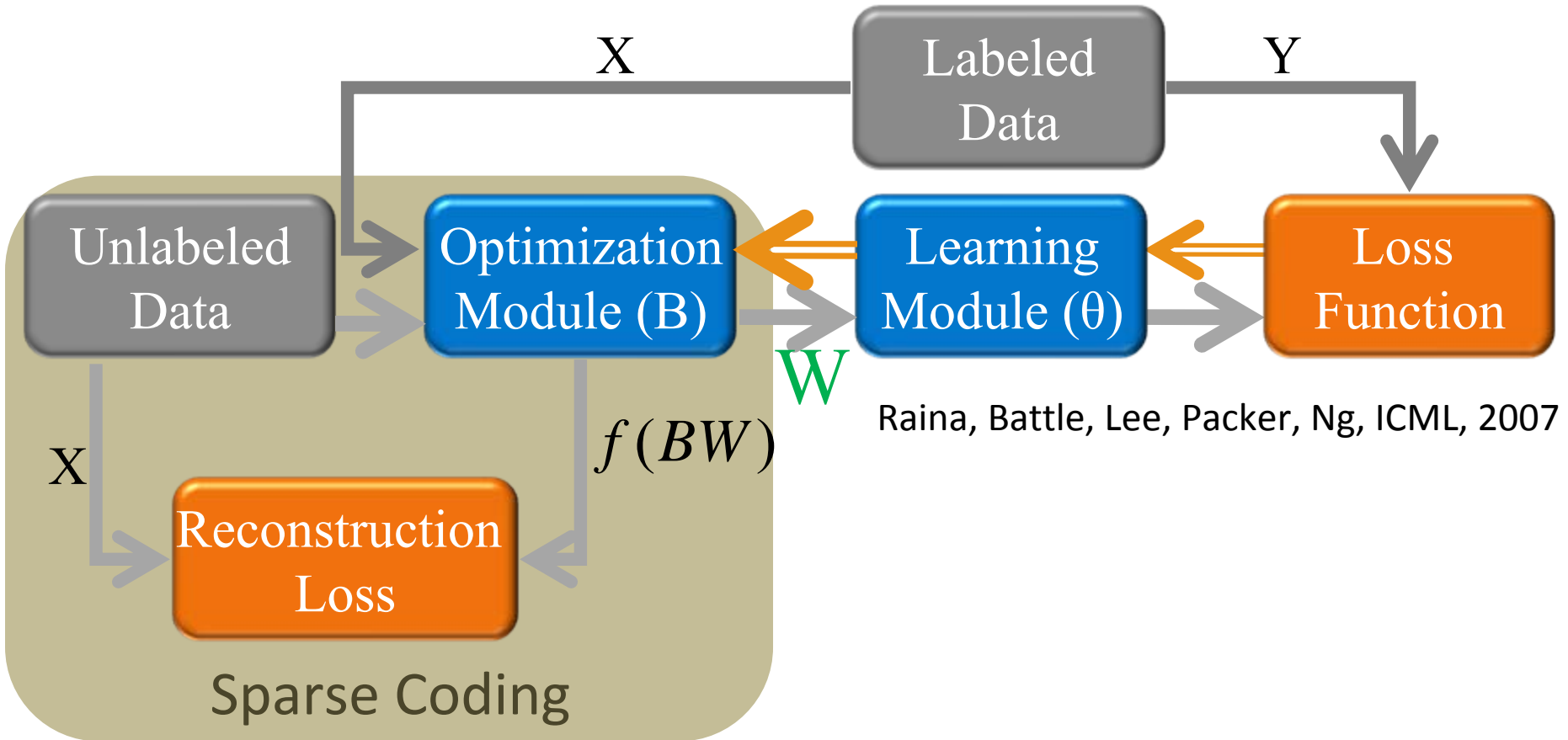
$$\arg \min_{W, B} \sum_j \frac{1}{2} \|B w_j - x\|_2^2 + \lambda \|w_j\|_1$$

Optimize B over  
all examples

$$s.t. \quad \|b_j\|_2 = 1$$

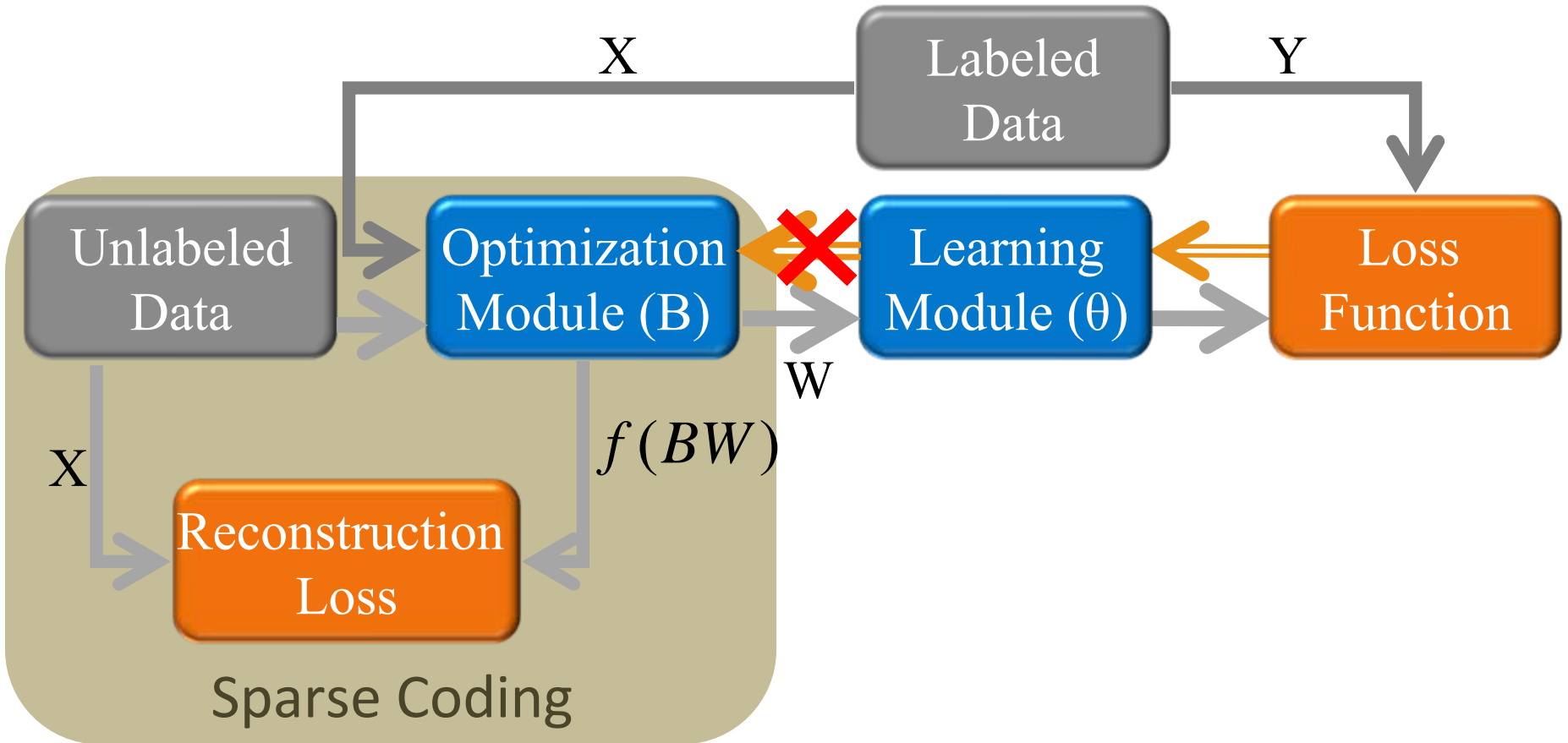
Shown to generate useful features on diverse  
problems

# Differentiable Sparse Coding



Bradley & Bagnell, "Differentiable Sparse Coding", NIPS 2008

# L1 Regularization is Not Differentiable



Bradley & Bagnell, "Differentiable Sparse Coding", NIPS 2008



Why is this unsatisfying?



# Problem #1: Instability

- L1 Map Estimates are discontinuous
- Outputs are not differentiable
- Instead use KL-divergence

Proven to compete with  
L1 in online learning



## Problem #2:

No closed-form Equation

$$\hat{w} = \arg \min_{\vec{w}} \text{Loss}(\vec{x} \parallel f(B\vec{w})) + \lambda \text{Prior}(\vec{w} \parallel \vec{p})$$

At the MAP estimate:

$$\nabla_{\vec{w}} \text{Loss}(\vec{x} \parallel f(B\hat{w})) + \lambda \nabla_{\vec{w}} \text{Prior}(\hat{w} \parallel \vec{p}) = 0$$

# Solution: Implicit Differentiation

Differentiate both sides with respect to an element of B:

$$\frac{\partial}{\partial B_j^i} \left( \nabla_{\vec{w}} \text{Loss}(\vec{x} \parallel f(B\hat{w})) + \lambda \nabla_{\vec{w}} \text{Prior}(\hat{w} \parallel \vec{p}) \right) = 0$$

Since  $\hat{w}$  is a function of B:

$$\frac{\partial}{\partial B_j^i} (\hat{w}) = \frac{\partial \hat{w}}{\partial B_j^i}$$

Solve for this



# Example: Squared Loss, KL prior

$$\hat{w} = \arg \min_w \frac{1}{2} \|x - Bw\|_2^2 + \lambda \sum_i w_i \log \frac{w_i}{p_i} - w_i + p_i$$

KL-Divergence

$$\frac{\partial \hat{w}}{\partial B_i^k} = - \left( B^T B + \text{diag} \left( \frac{\lambda}{\hat{w}} \right) \right)^{-1} \left( (B^k \hat{w}_i)^T + \vec{e}_i (f(B^k \hat{w}) - x_k) \right)$$



# Handwritten Digit Recognition

50,000 digit training set

10,000 digit validation set

10,000 digit test set



# Handwritten Digit Recognition

Step #1:

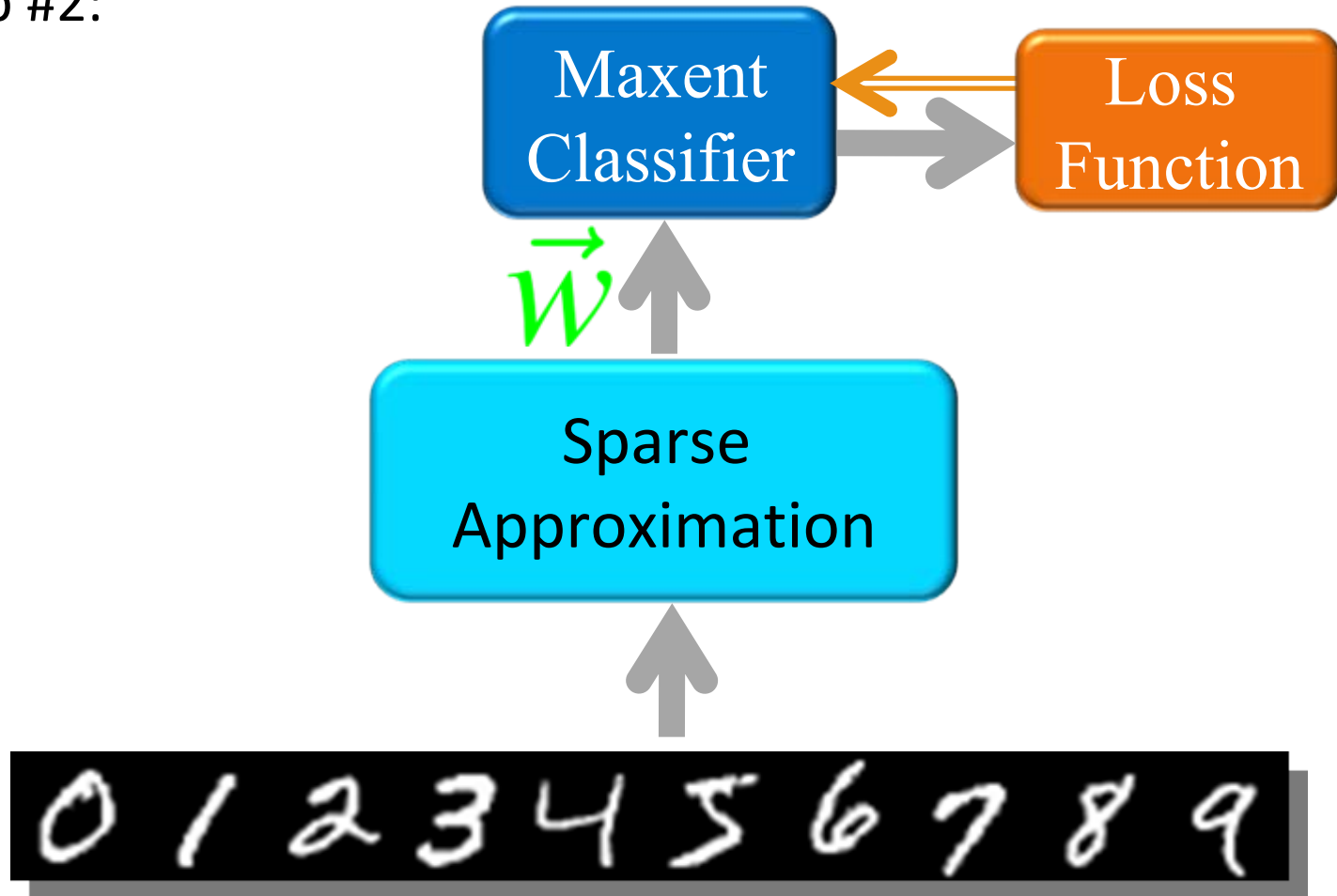
Unsupervised  
Sparse Coding  
 $L_2$  loss and  $L_1$  prior



Training Set

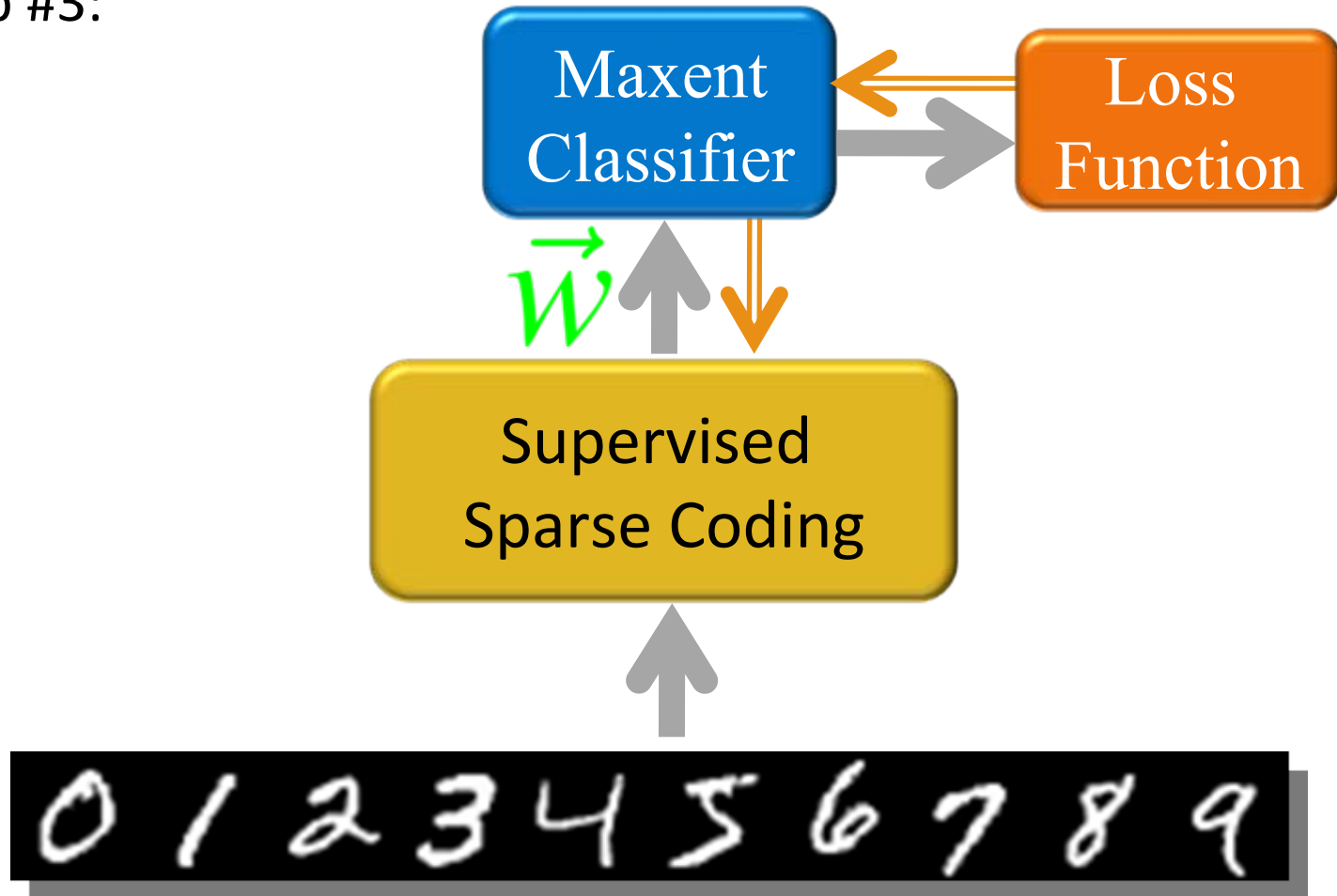
# Handwritten Digit Recognition

Step #2:

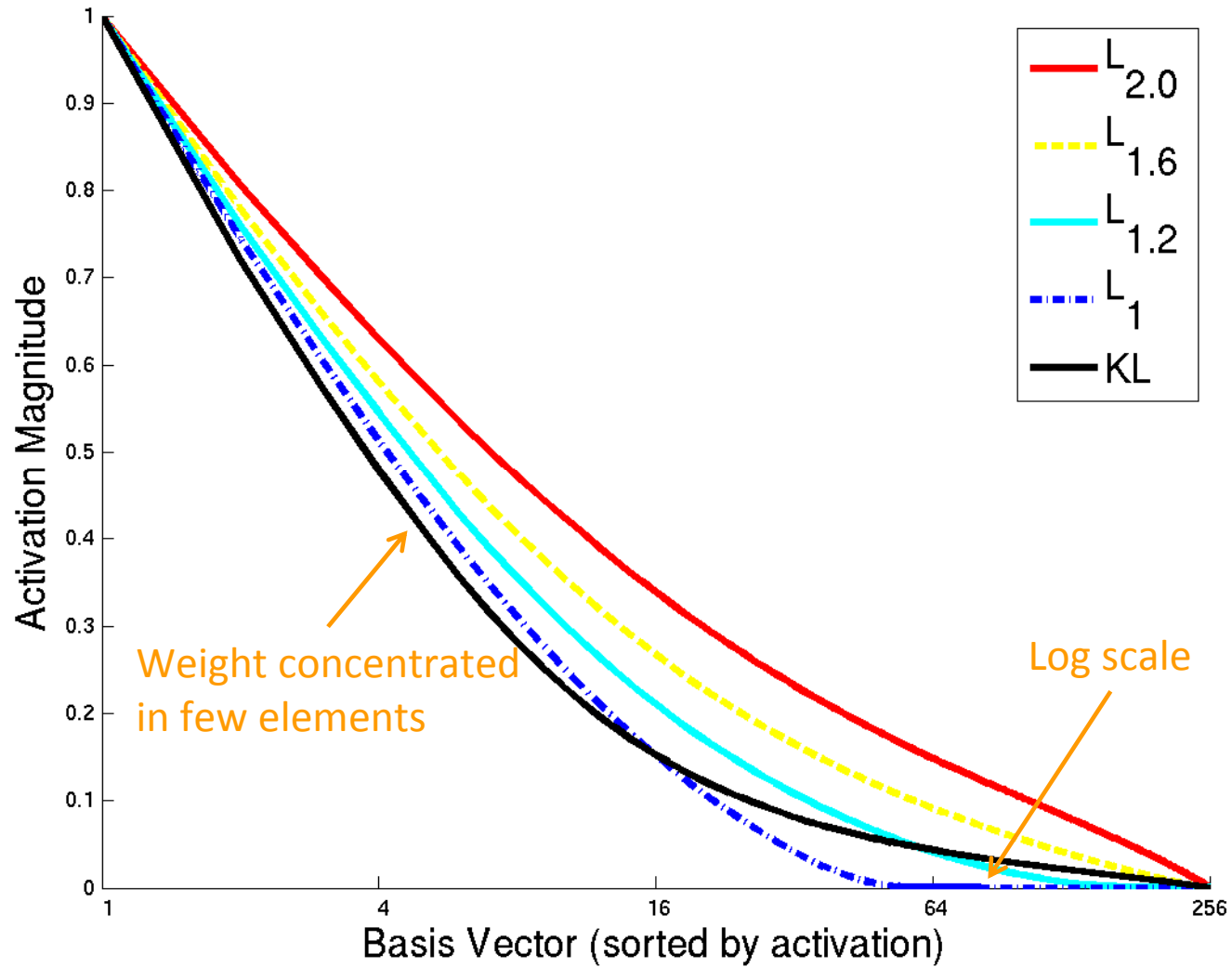


# Handwritten Digit Recognition

Step #3:

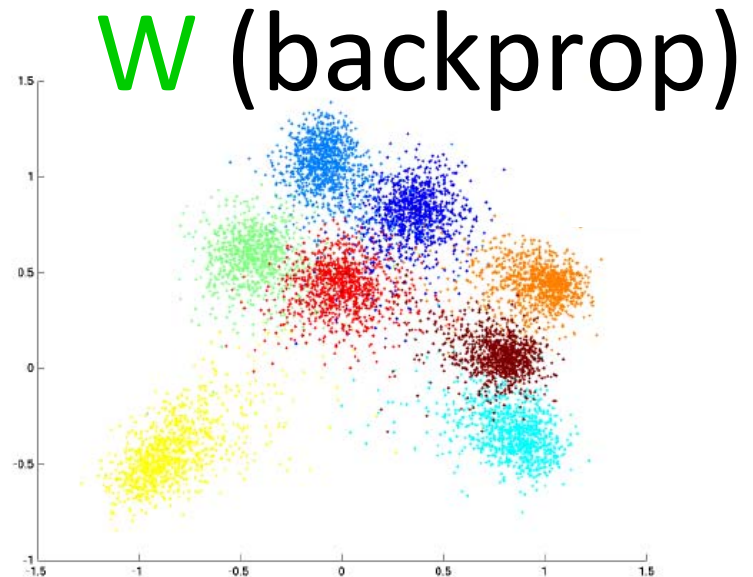
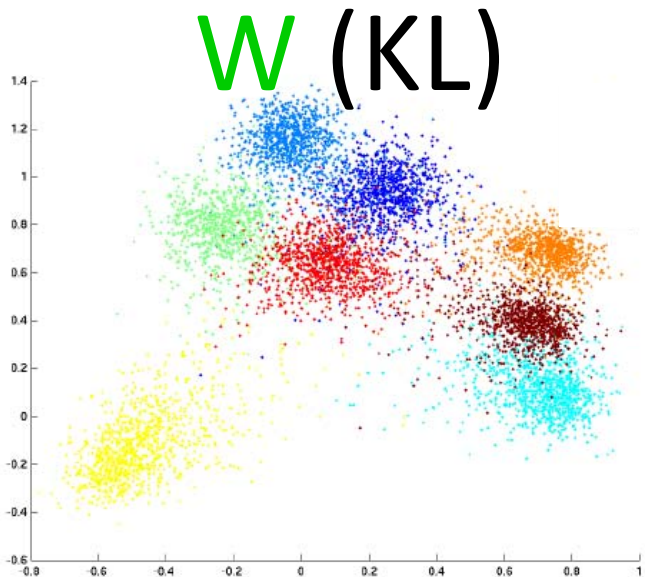
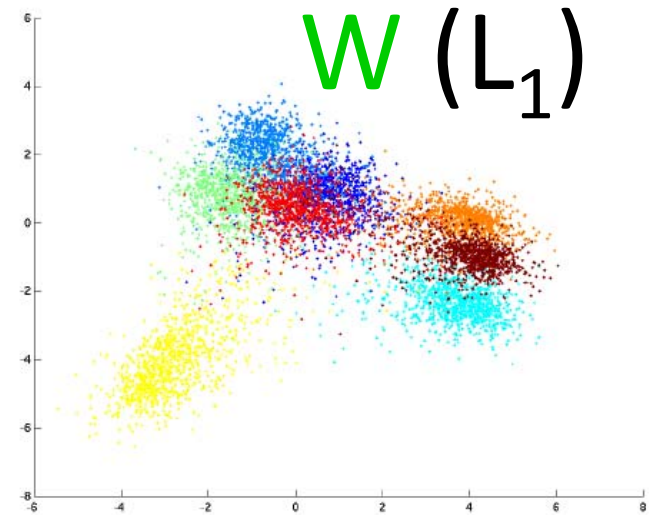
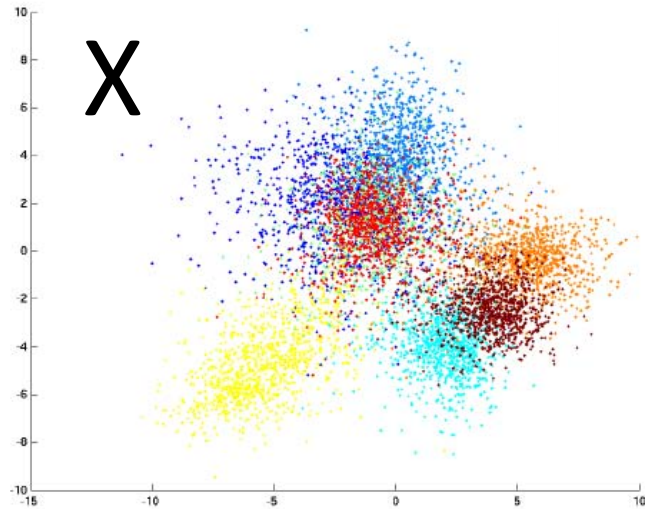


# KL Maintains Sparsity

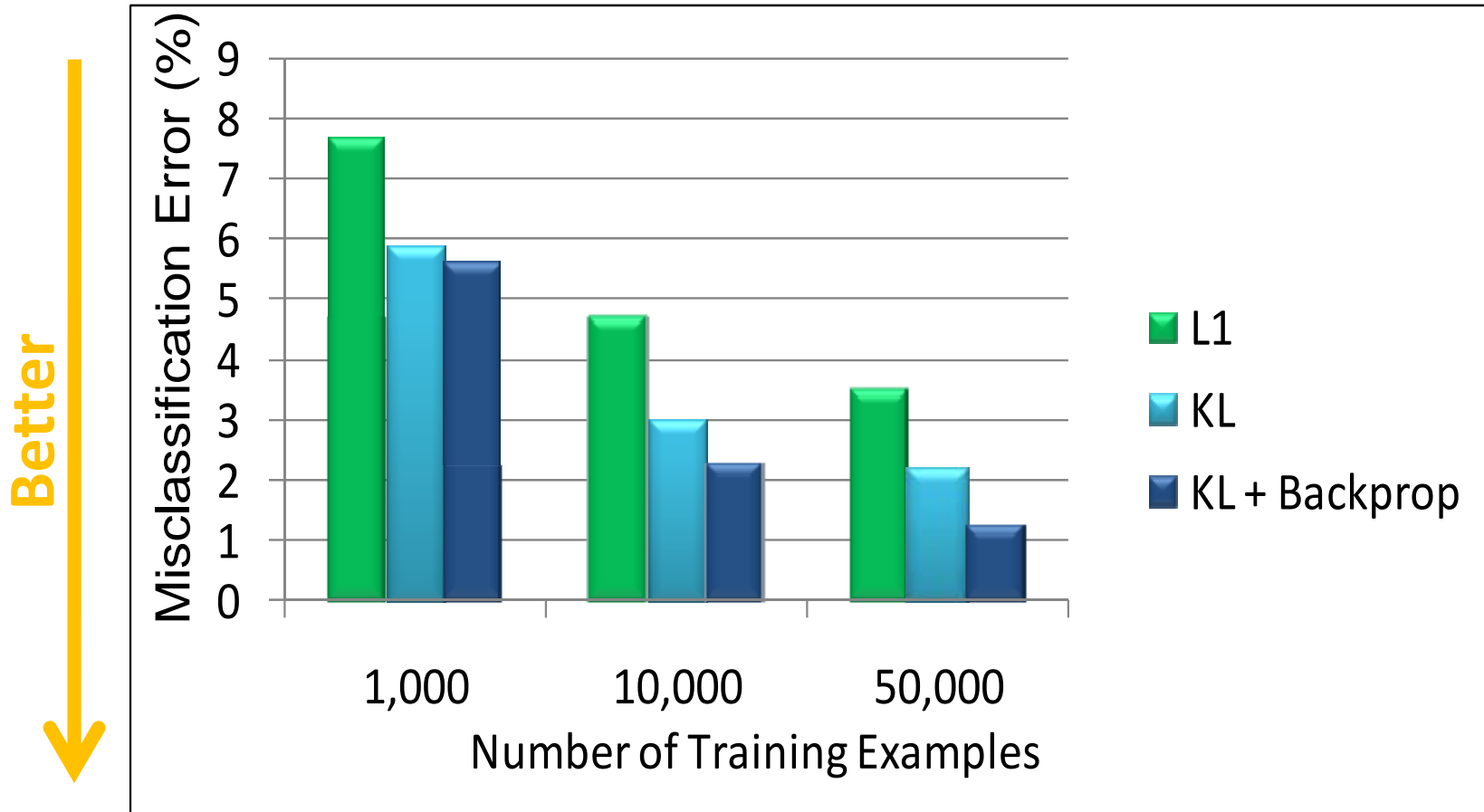




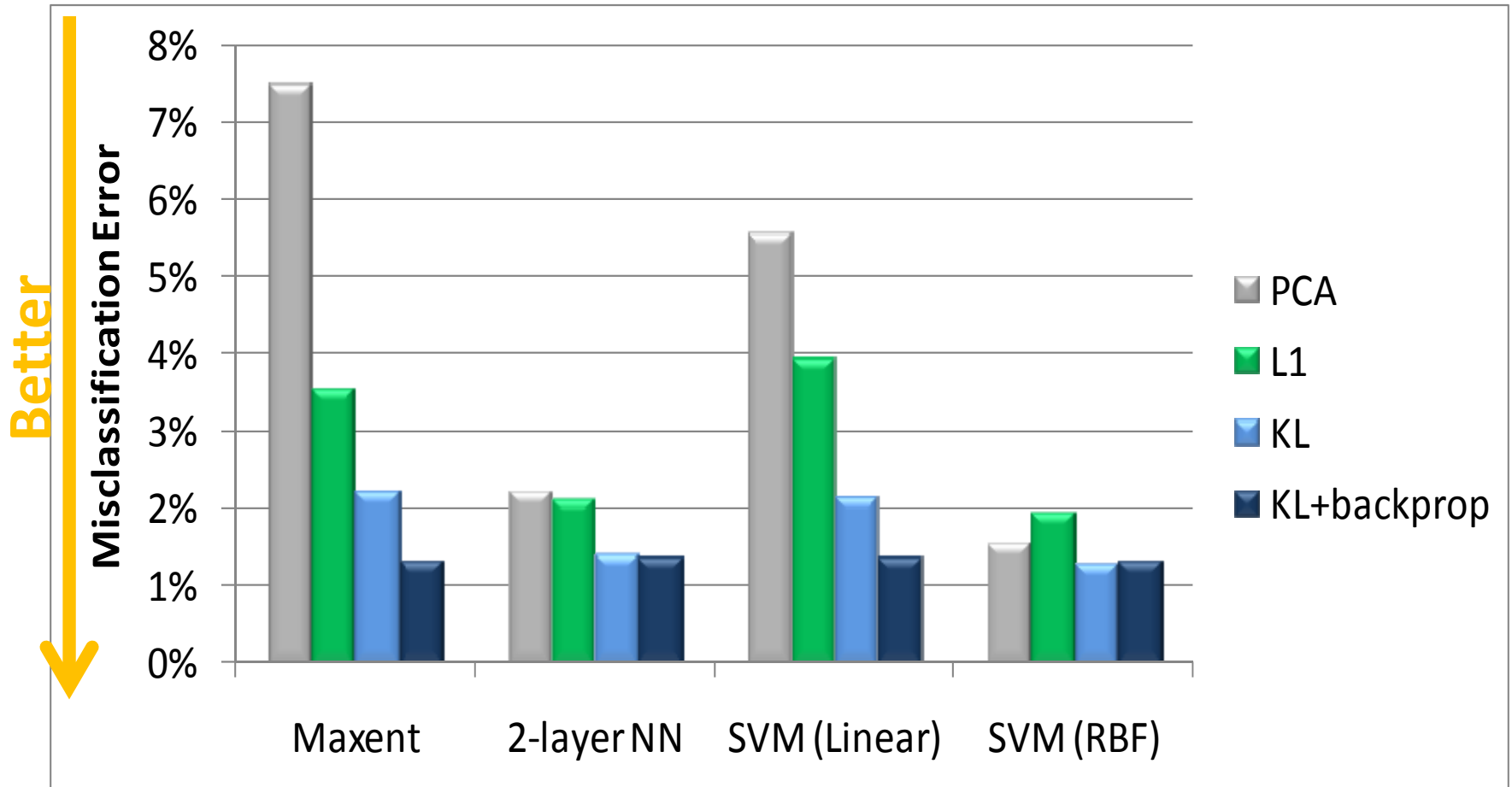
# KL adds Stability



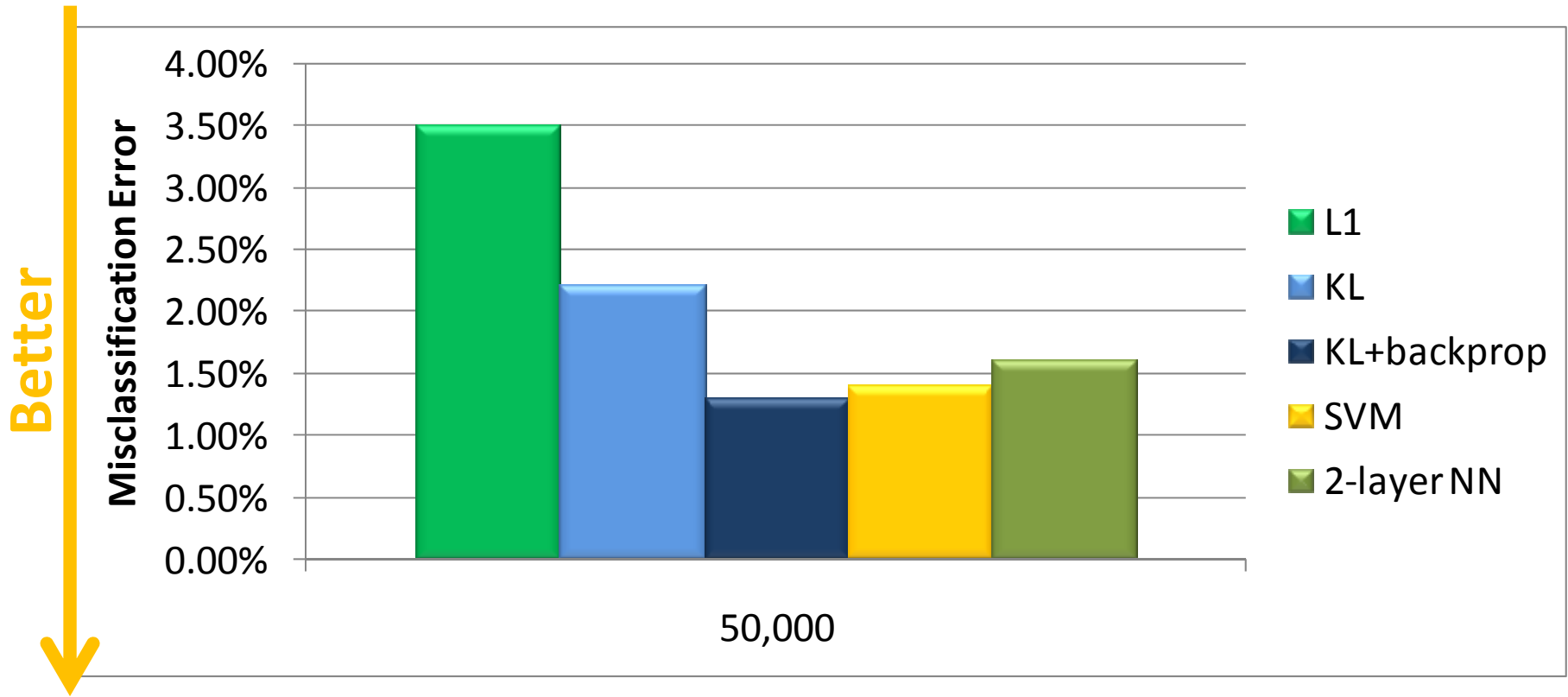
# Performance vs. Prior



# Classifier Comparison



# Comparison to other algorithms



Algorithm	L1	KL	KL+backprop	SVM	2-layer NN [15]
Test Set Error	3.53%	2.21%	<b>1.30%</b>	1.4%	1.6%

# Transfer to English Characters

24,000 character training set

12,000 character validation set

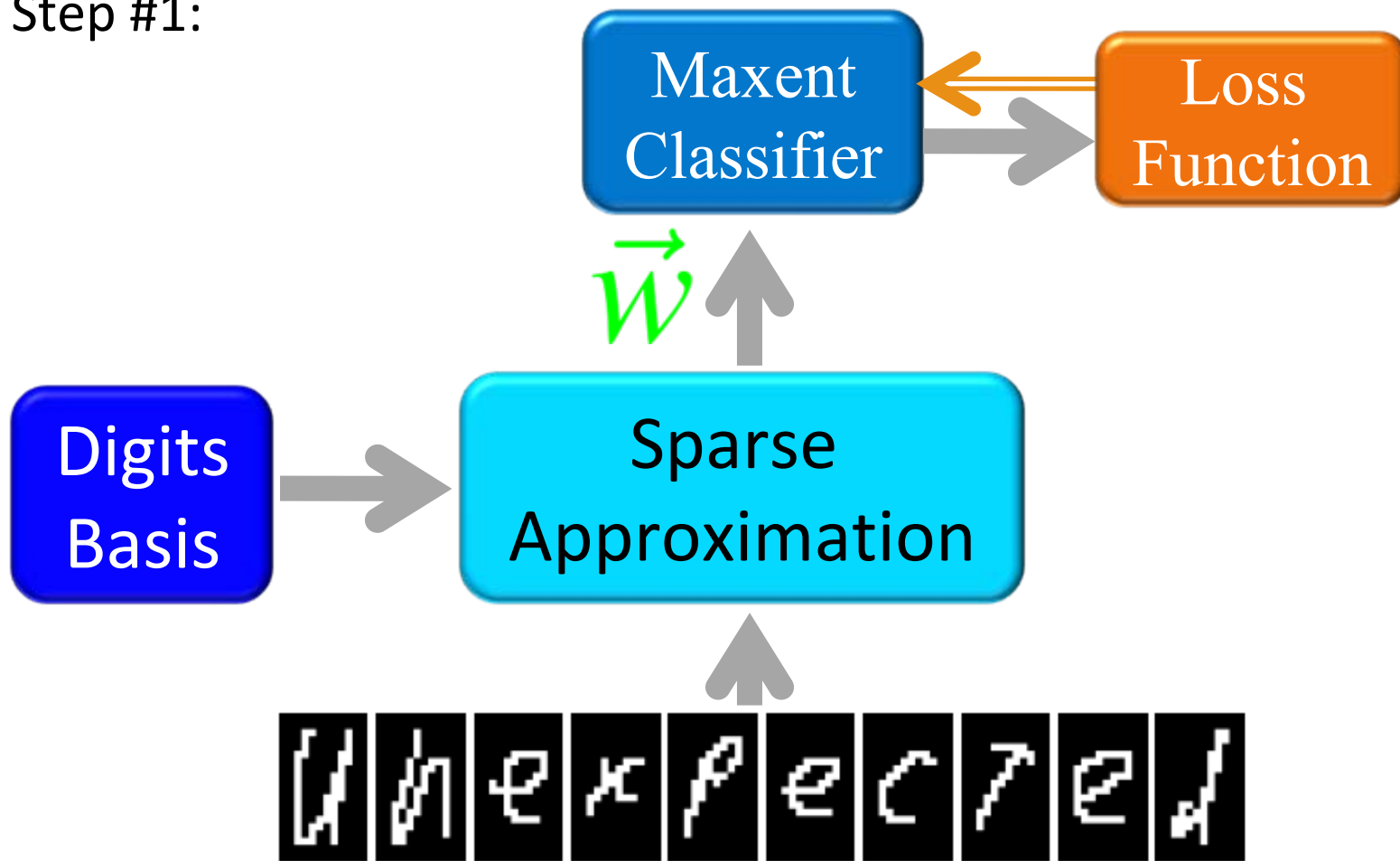
12,000 character test set





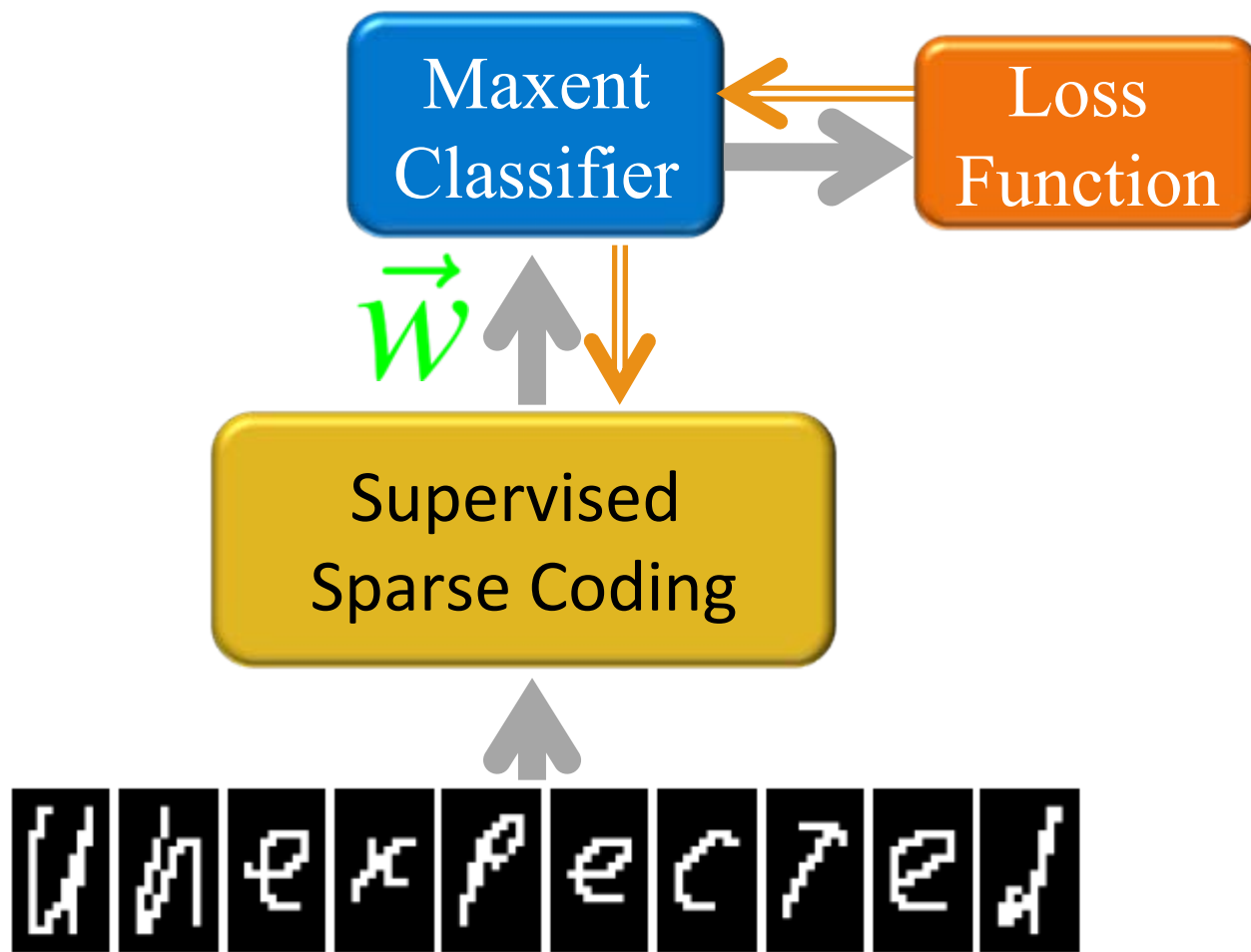
# Transfer to English Characters

Step #1:

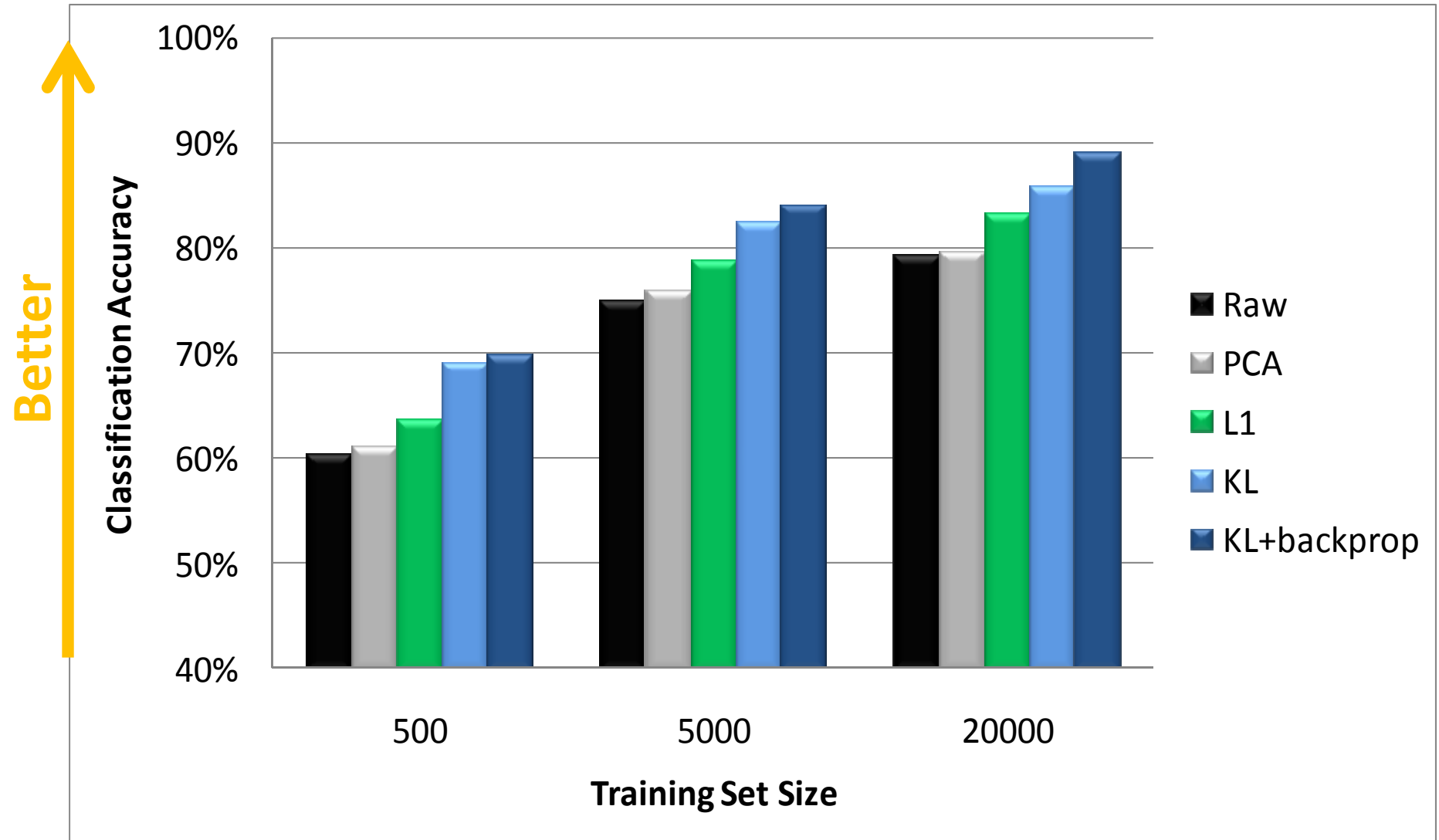


# Transfer to English Characters

Step #2:



# Transfer to English Characters



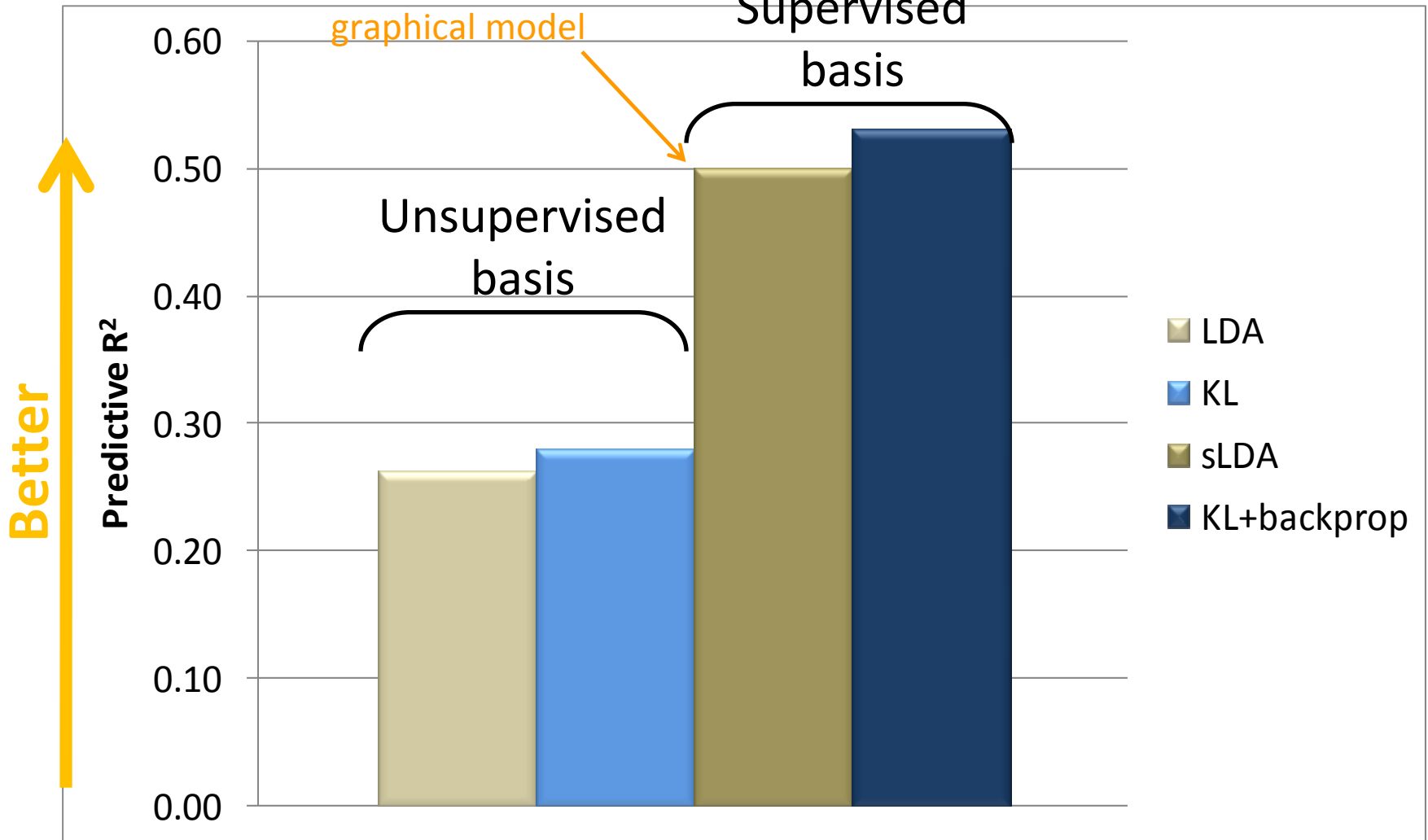






# Movie Review Sentiment

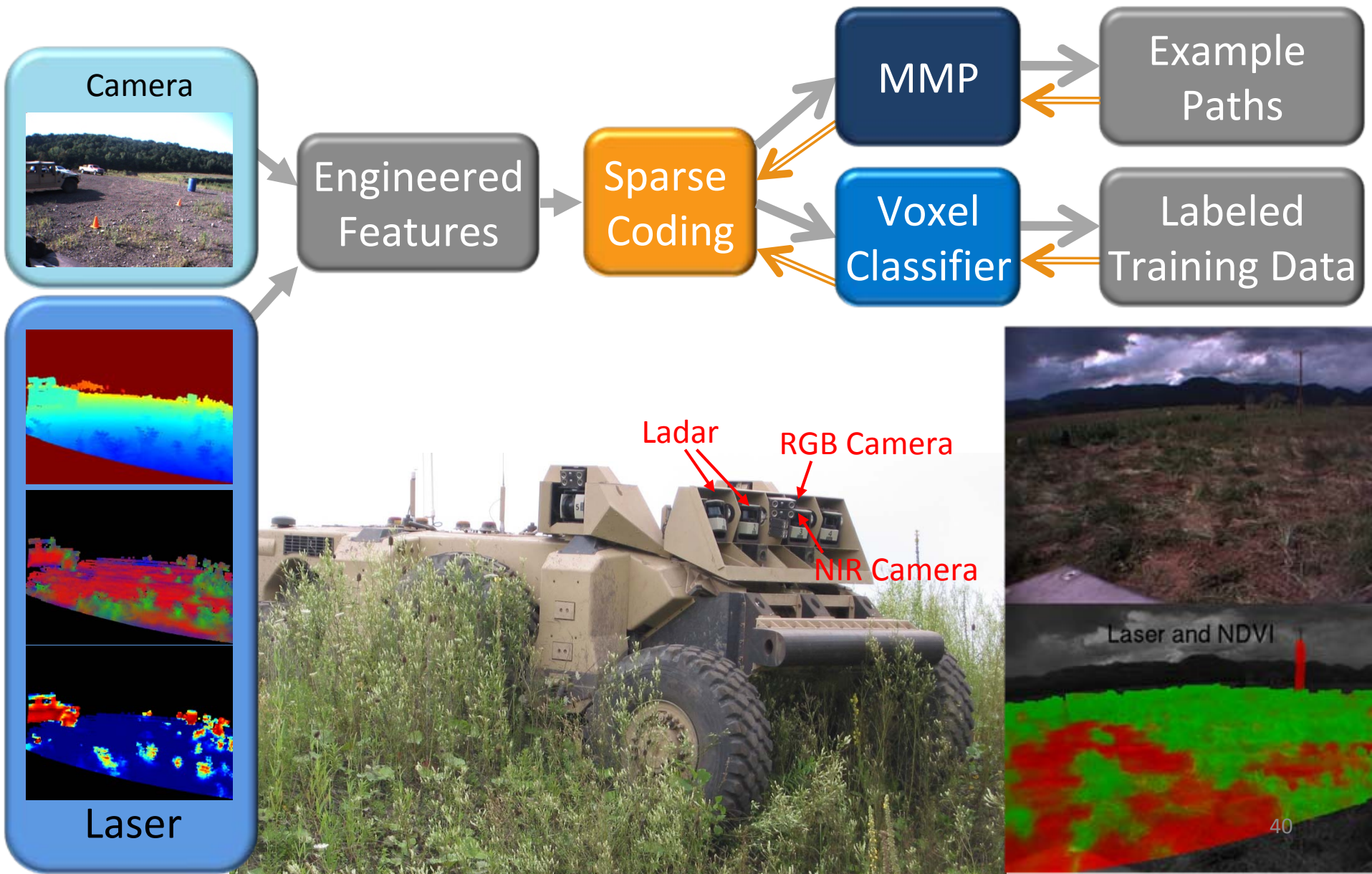
State of the art  
graphical model



Blei, McAuliffe, NIPS, 2007



# Future Work



# Future Work:

## Convex Sparse Coding

- Sparse approximation is convex
- Sparse coding is not because fixed-size basis is a non-convex constraint
- Sparse coding  $\leftrightarrow$  sparse approximation on infinitely large basis + non-convex rank constraint
  - Relax to a convex  $L_1$  rank constraint
- Use boosting for sparse approximation directly on infinitely large basis

Bengio, Le Roux, Vincent, Dellalleau, Marcotte, NIPS, 2005

Zhao, Yu. *Feature Selection for Data Mining*, 2005

Rifkin, Lippert. *Journal of Machine Learning Research*, 2007

# Questions?