



NEW YORK UNIVERSITY

# Energy-Based Models (part 4)

<http://bit.ly/DLSP20>

Yann LeCun

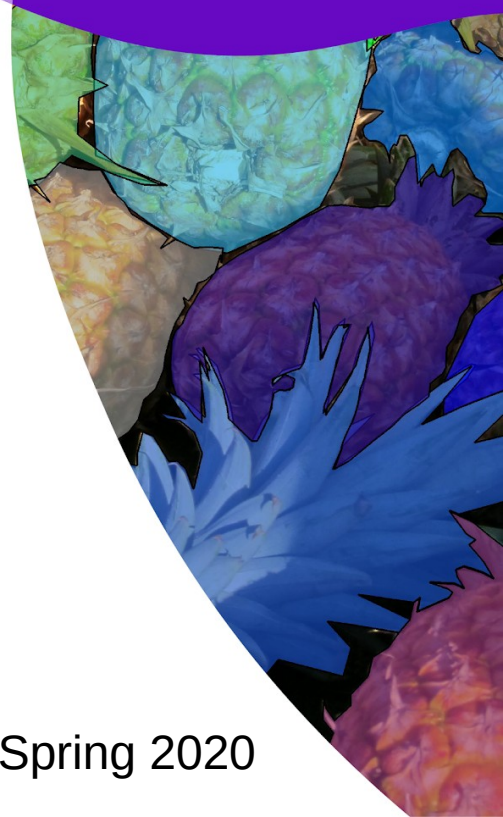
NYU - Courant Institute & Center for Data Science

Facebook AI Research

<http://yann.lecun.com>

TAs: Alfredo Canziani, Mark Goldstein

Deep Learning, NYU, Spring 2020



# Architecture and Loss Function

► **Family of energy functions**  $\mathcal{E} = \{E(W, Y, X) : W \in \mathcal{W}\}.$

► **Training set**  $\hat{\mathcal{S}} = \{(X^i, Y^i) : i = 1 \dots P\}.$

► **Loss functional / Loss function**  $\mathcal{L}(E, \mathcal{S}) \quad \mathcal{L}(W, \mathcal{S})$

► Measures the quality of an energy function on training set

► **Training**  $W^* = \min_{W \in \mathcal{W}} \mathcal{L}(W, \mathcal{S}).$

► **Form of the loss functional**

► invariant under permutations and repetitions of the samples

$$\mathcal{L}(E, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^P L(Y^i, E(W, \mathcal{Y}, X^i)) + R(W).$$

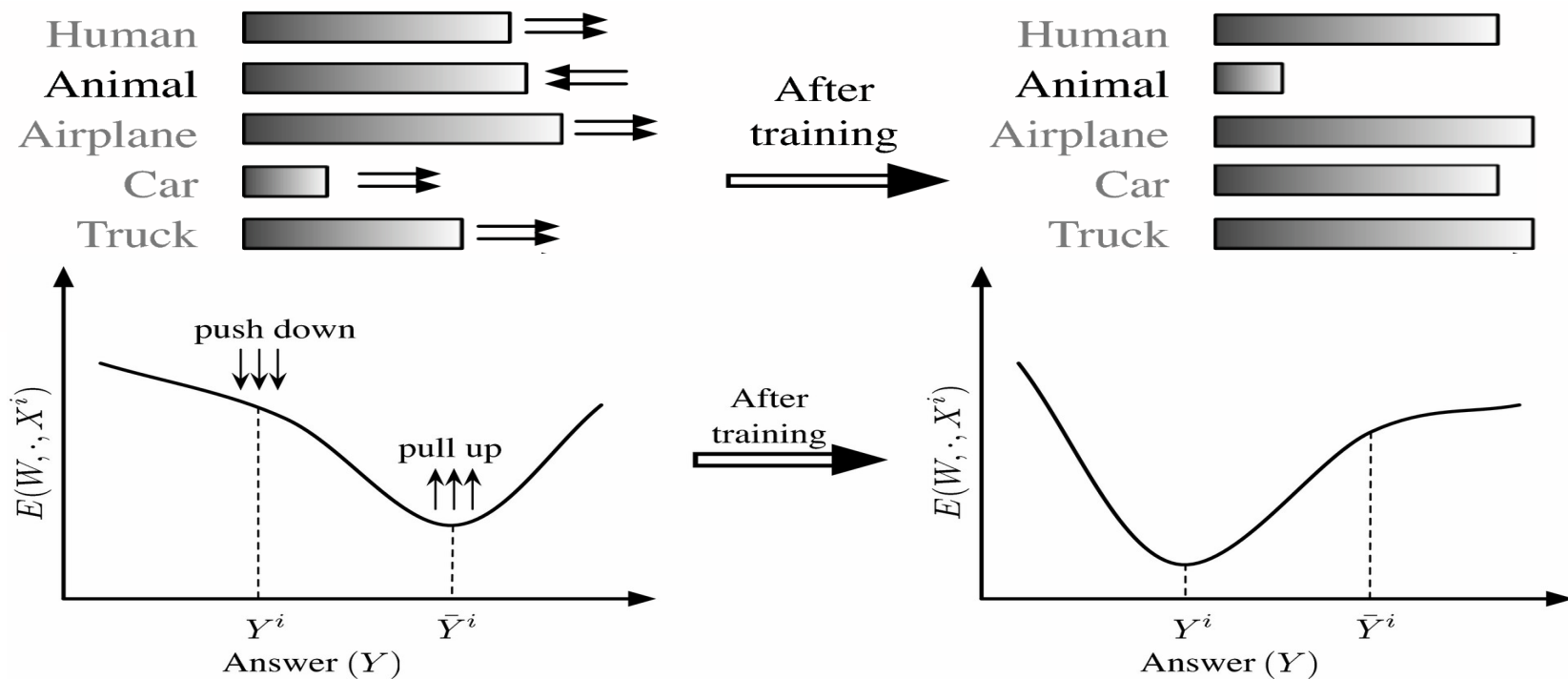
Per-sample  
loss

Desired  
answer

Energy surface  
for a given  $X_i$   
as  $Y$  varies

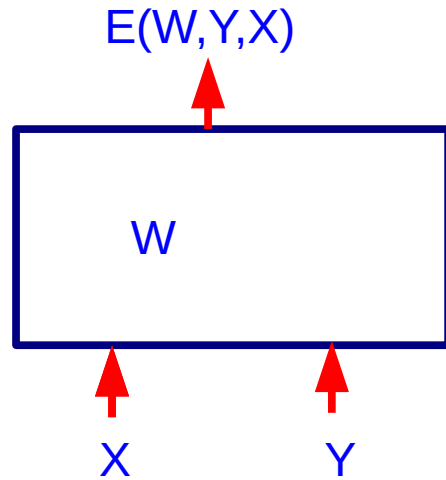
Regularizer

# Designing a Loss Function



- ▶ **Push down** on the energy of the correct answer
- ▶ **Pull up** on the energies of the incorrect answers, particularly if they are smaller than the correct one

# Architecture + Inference Algo + Loss Function = Model



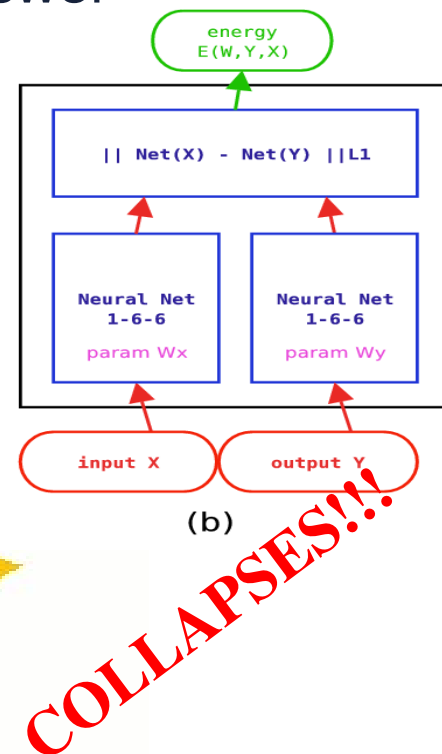
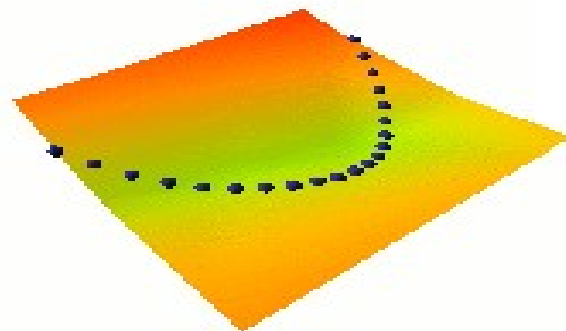
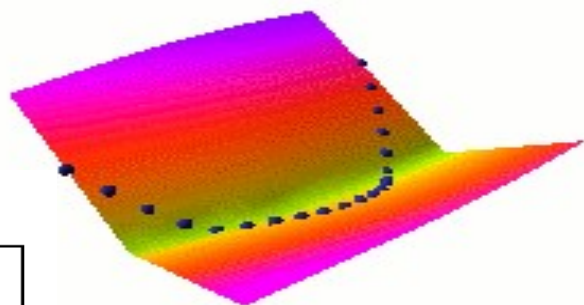
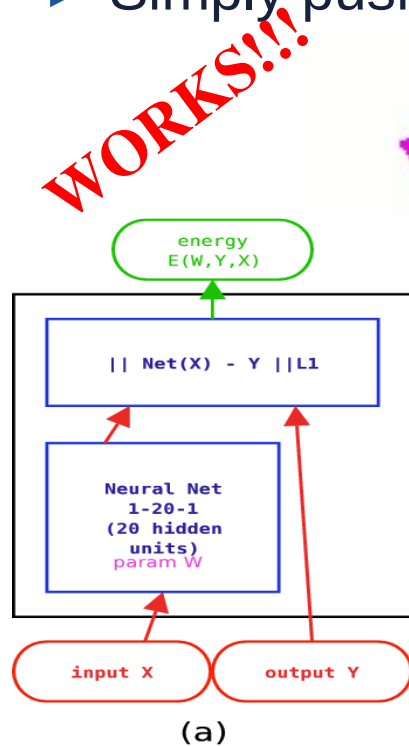
1. **Design an architecture:** a particular form for  $E(W, Y, X)$ .
2. **Pick an inference algorithm for Y:** MAP or conditional distribution, belief prop, min cut, variational methods, gradient descent, MCMC, HMC.....
3. **Pick a loss function:** in such a way that minimizing it with respect to  $W$  over a training set will make the inference algorithm find the correct  $Y$  for a given  $X$ .
4. **Pick an optimization method.**

❏ **PROBLEM:** What loss functions will make the machine approach the desired behavior?

# Examples of Loss Functions: Energy Loss

► **Energy Loss**  $L_{energy}(Y^i, E(W, \mathcal{Y}, X^i)) = E(W, Y^i, X^i)$ .

► Simply pushes down on the energy of the correct answer



# Negative Log-Likelihood Loss

- **Conditional probability of the samples (assuming independence)**

$$P(Y^1, \dots, Y^P | X^1, \dots, X^P, W) = \prod_{i=1}^P P(Y^i | X^i, W).$$

$$-\log \prod_{i=1}^P P(Y^i | X^i, W) = \sum_{i=1}^P -\log P(Y^i | X^i, W).$$

- **Gibbs distribution:**  $P(Y | X^i, W) = \frac{e^{-\beta E(W, Y, X^i)}}{\int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)}}.$

$$-\log \prod_{i=1}^P P(Y^i | X^i, W) = \sum_{i=1}^P \beta E(W, Y^i, X^i) + \log \int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)}.$$

- **We get the NLL loss by dividing by P and Beta:**

$$\mathcal{L}_{\text{nll}}(W, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^P \left( E(W, Y^i, X^i) + \frac{1}{\beta} \log \int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)} \right).$$

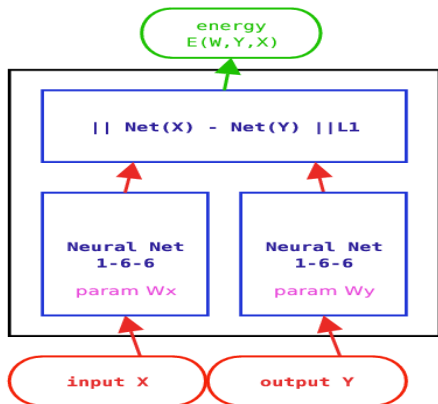
- **Reduces to the perceptron loss when Beta->infinity**

# Negative Log-Likelihood Loss

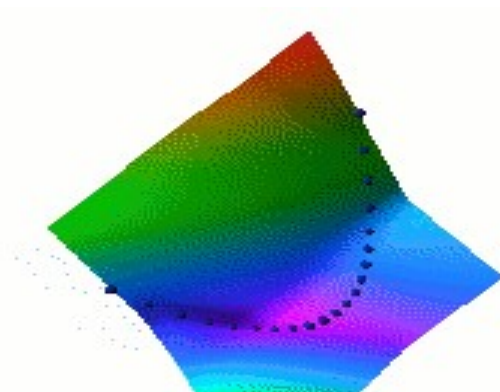
- ▶ Pushes down on the energy of the correct answer
- ▶ Pulls up on the energies of all answers in proportion to their probability

$$\mathcal{L}_{\text{nll}}(W, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^P \left( E(W, Y^i, X^i) + \frac{1}{\beta} \log \int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)} \right).$$

$$\frac{\partial \mathcal{L}_{\text{nll}}(W, Y^i, X^i)}{\partial W} = \frac{\partial E(W, Y^i, X^i)}{\partial W} - \int_{Y \in \mathcal{Y}} \frac{\partial E(W, Y, X^i)}{\partial W} P(Y|X^i, W),$$



(b)



# Negative Log-Likelihood Loss

- ▶ **A probabilistic model is an EBM in which:**
  - ▶ The energy can be integrated over  $Y$  (the variable to be predicted)
  - ▶ The loss function is the negative log-likelihood
- ▶ **Negative Log Likelihood Loss has been used for a long time in many communities for discriminative learning with structured outputs**
  - ▶ Speech recognition: many papers going back to the early 90's [Bengio 92], [Bourlard 94]. They call “Maximum Mutual Information”
  - ▶ Handwriting recognition [Bengio LeCun 94], [LeCun et al. 98]
  - ▶ Bio-informatics [Haussler]
  - ▶ Conditional Random Fields [Lafferty et al. 2001]
  - ▶ Lots more.....
  - ▶ In all the above cases, **it was used with non-linearly parameterized energies.**



# A Simpler Loss Functions: Perceptron Loss

$$L_{\text{perceptron}}(Y^i, E(W, \mathcal{Y}, X^i)) = E(W, Y^i, X^i) - \min_{Y \in \mathcal{Y}} E(W, Y, X^i).$$

## ► Perceptron Loss [LeCun et al. 1998], [Collins 2002]

- Pushes down on the energy of the correct answer
- Pulls up on the energy of the machine's answer
- Always positive. Zero when answer is correct
- No “margin”: technically does not prevent the energy surface from being almost flat.
- Works pretty well in practice, particularly if the energy parameterization does not allow flat surfaces.
- This is often called “**discriminative Viterbi training**” in the speech and handwriting literature

# Perceptron Loss for Binary Classification

$$L_{\text{perceptron}}(Y^i, E(W, \mathcal{Y}, X^i)) = E(W, Y^i, X^i) - \min_{Y \in \mathcal{Y}} E(W, Y, X^i).$$

► **Energy:**  $E(W, Y, X) = -Y G_W(X),$

► **Inference:**  $Y^* = \operatorname{argmin}_{Y \in \{-1, 1\}} -Y G_W(X) = \operatorname{sign}(G_W(X)).$

► **Loss:**  $\mathcal{L}_{\text{perceptron}}(W, \mathcal{S}) = \frac{1}{P} \sum_{i=1}^P (\operatorname{sign}(G_W(X^i)) - Y^i) G_W(X^i).$

► **Learning Rule:**  $W \leftarrow W + \eta (Y^i - \operatorname{sign}(G_W(X^i))) \frac{\partial G_W(X^i)}{\partial W},$

► **If  $G_W(X)$  is linear in  $W$ :**  $E(W, Y, X) = -Y W^T \Phi(X)$

$$W \leftarrow W + \eta (Y^i - \operatorname{sign}(W^T \Phi(X^i))) \Phi(X^i)$$

# A Better Loss Function: Generalized Margin Losses

► First, we need to define the **Most Offending Incorrect Answer**

► **Most Offending Incorrect Answer: discrete case**

**Definition 1** Let  $Y$  be a discrete variable. Then for a training sample  $(X^i, Y^i)$ , the **most offending incorrect answer**  $\bar{Y}^i$  is the answer that has the lowest energy among all answers that are incorrect:

$$\bar{Y}^i = \operatorname{argmin}_{Y \in \mathcal{Y} \text{ and } Y \neq Y^i} E(W, Y, X^i). \quad (8)$$

► **Most Offending Incorrect Answer: continuous case**

**Definition 2** Let  $Y$  be a continuous variable. Then for a training sample  $(X^i, Y^i)$ , the **most offending incorrect answer**  $\bar{Y}^i$  is the answer that has the lowest energy among all answers that are at least  $\epsilon$  away from the correct answer:

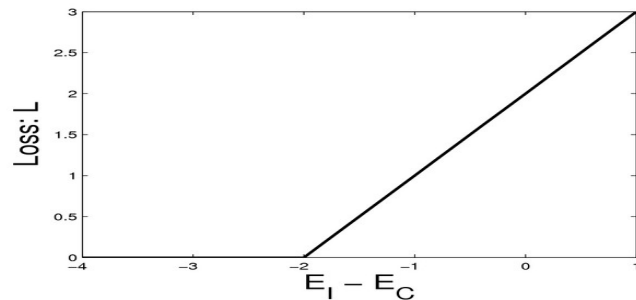
$$\bar{Y}^i = \operatorname{argmin}_{Y \in \mathcal{Y}, \|Y - Y^i\| > \epsilon} E(W, Y, X^i). \quad (9)$$

# Examples of Generalized Margin Losses

$$L_{\text{hinge}}(W, Y^i, X^i) = \max(0, m + E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)) ,$$

## ► Hinge Loss

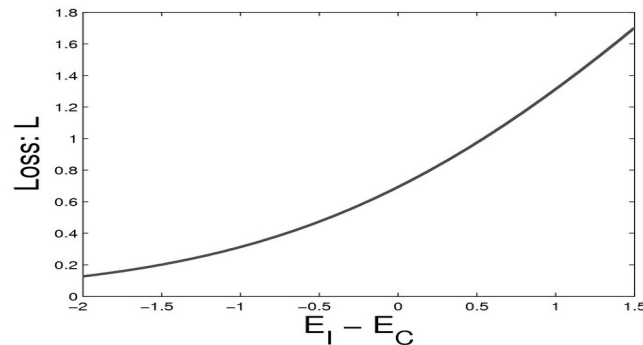
- [Altun et al. 2003], [Taskar et al. 2003]
- With the linearly-parameterized binary classifier architecture, we get linear SVMs



$$L_{\text{log}}(W, Y^i, X^i) = \log \left( 1 + e^{E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)} \right) .$$

## ► Log Loss

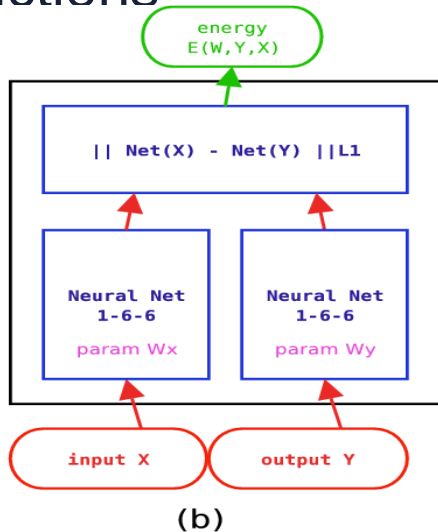
- “soft hinge” loss
- With the linearly-parameterized binary classifier architecture, we get linear Logistic Regression



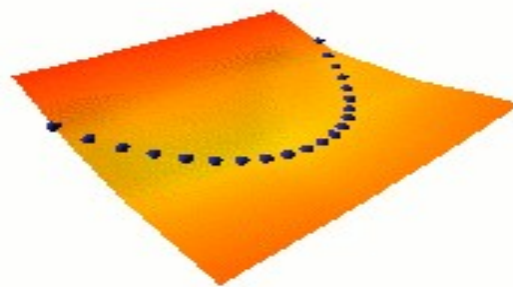
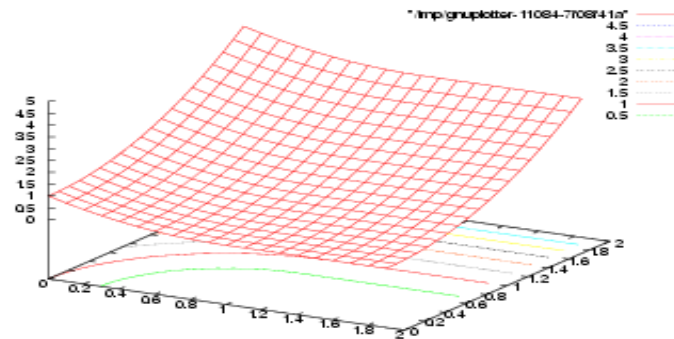
# Examples of Margin Losses: Square-Square Loss

$$L_{\text{sq-sq}}(W, Y^i, X^i) = E(W, Y^i, X^i)^2 + (\max(0, m - E(W, \bar{Y}^i, X^i)))^2.$$

- **Square-Square Loss**
- [LeCun-Huang 2005]
- Appropriate for positive energy functions



Learning  $Y = X^2$



**NO COLLAPSE!!!**

# Other Margin-Like Losses

► **LVQ2 Loss** [Kohonen, Oja], Driancourt-Bottou 1991]

$$L_{\text{lvq2}}(W, Y^i, X^i) = \min \left( 1, \max \left( 0, \frac{E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)}{\delta E(W, \bar{Y}^i, X^i)} \right) \right),$$

► **Minimum Classification Error Loss** [Juang, Chou, Lee 1997]

$$L_{\text{mce}}(W, Y^i, X^i) = \sigma \left( E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i) \right),$$
$$\sigma(x) = (1 + e^{-x})^{-1}$$

► **Square-Exponential Loss** [Osadchy, Miller, LeCun 2004]

$$L_{\text{sq-exp}}(W, Y^i, X^i) = E(W, Y^i, X^i)^2 + \gamma e^{-E(W, \bar{Y}^i, X^i)},$$

# What Make a “Good” Loss Function

## ► Good and bad loss functions

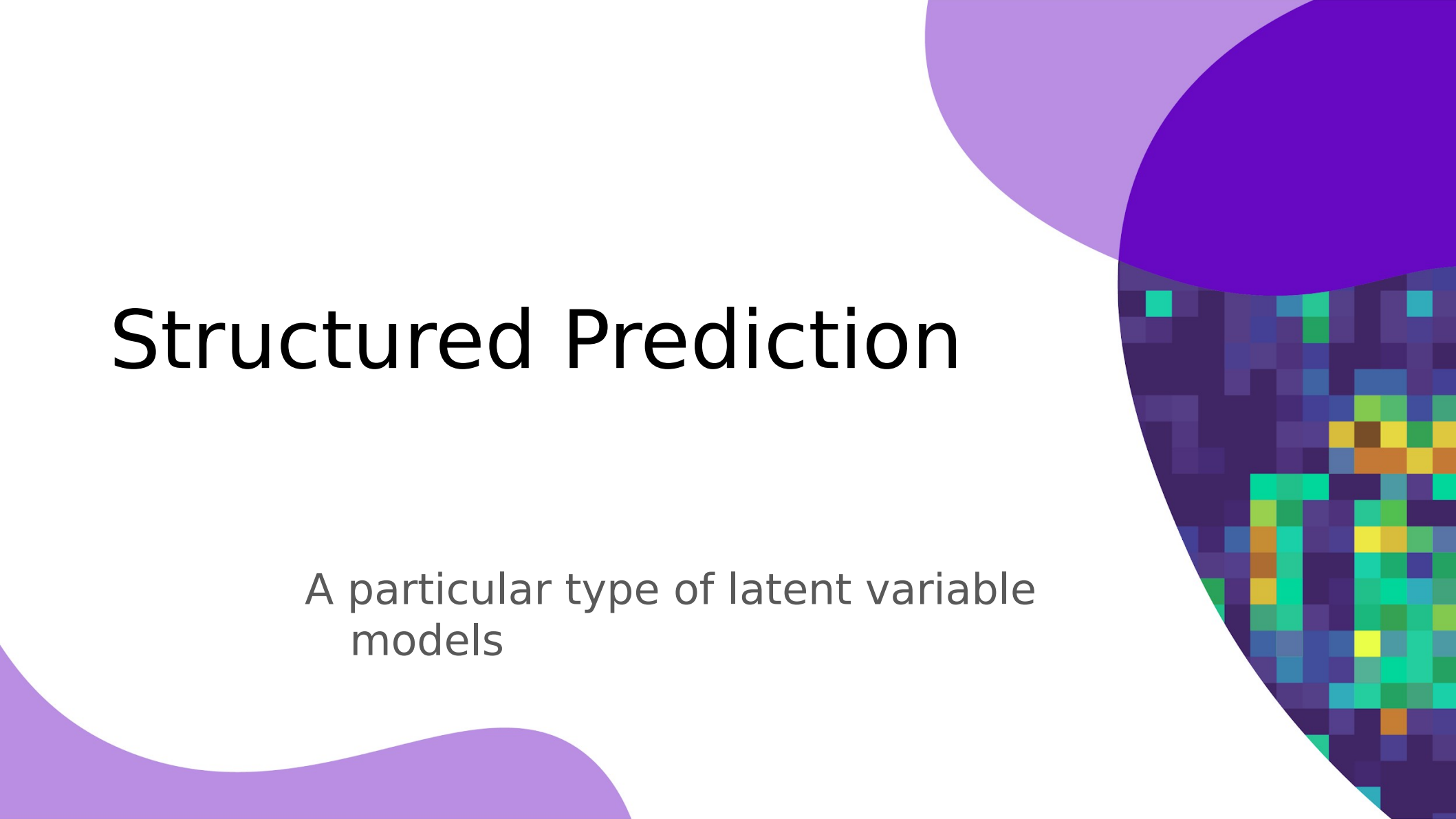
Loss (equation #)	Formula	Margin
energy loss	$E(W, Y^i, X^i)$	none
perceptron	$E(W, Y^i, X^i) - \min_{Y \in \mathcal{Y}} E(W, Y, X^i)$	0
hinge	$\max(0, m + E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i))$	$m$
log	$\log \left( 1 + e^{E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)} \right)$	$> 0$
LVQ2	$\min(M, \max(0, E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)))$	0
MCE	$\left( 1 + e^{-(E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i))} \right)^{-1}$	$> 0$
square-square	$E(W, Y^i, X^i)^2 - (\max(0, m - E(W, \bar{Y}^i, X^i)))^2$	$m$
square-exp	$E(W, Y^i, X^i)^2 + \beta e^{-E(W, \bar{Y}^i, X^i)}$	$> 0$
NLL/MMI	$E(W, Y^i, X^i) + \frac{1}{\beta} \log \int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)}$	$> 0$
MEE	$1 - e^{-\beta E(W, Y^i, X^i)} / \int_{y \in \mathcal{Y}} e^{-\beta E(W, y, X^i)}$	$> 0$

## ► Slightly more general form:

$$L(W, X^i, Y^i) = \sum_y H(E(W, Y^i, X^i) - E(W, y, X^i) + C(Y^i, y))$$

# Structured Prediction

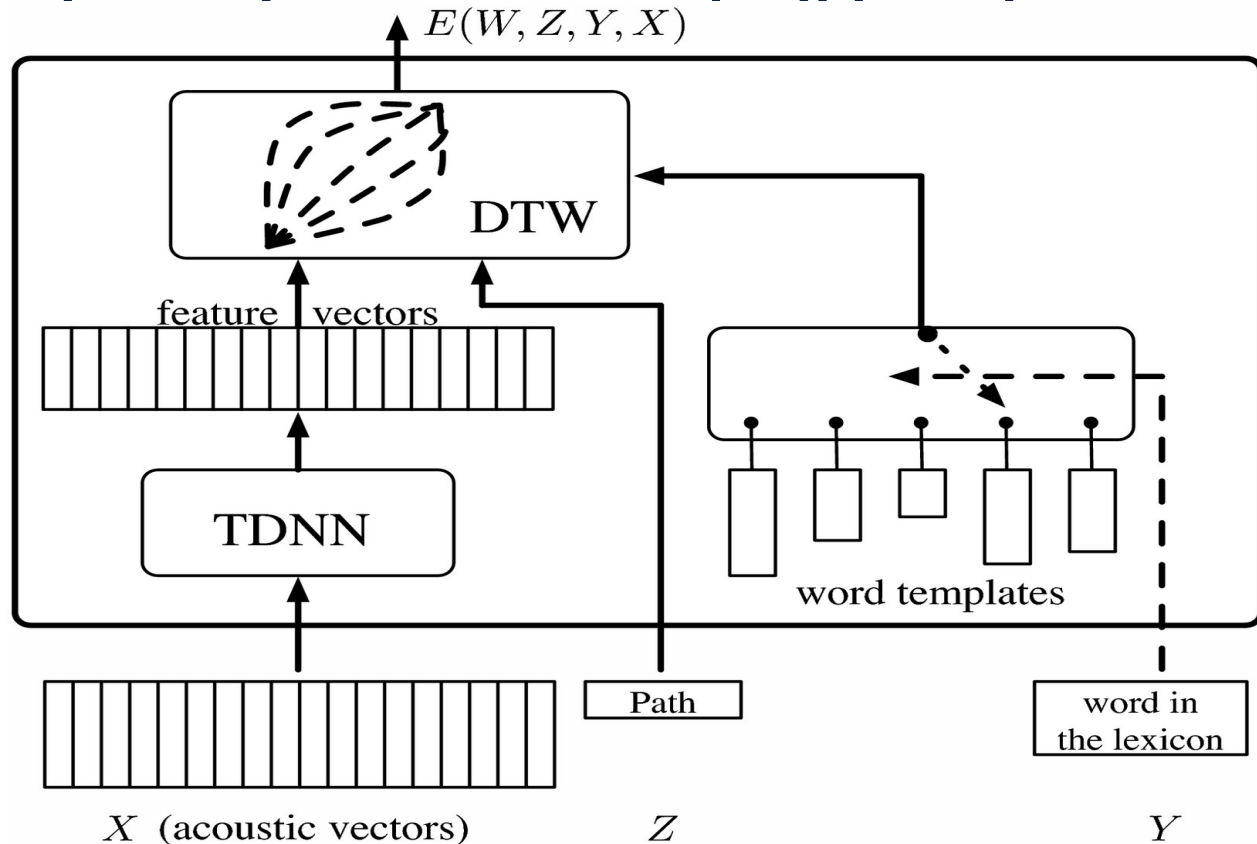
A particular type of latent variable models





# The Oldest Example of Structured Prediction

- ▶ Trainable Automatic Speech Recognition system with a **convolutional net** (TDNN) and dynamic time warping (DTW)
- ▶ The feature extractor and the structured classifier are trained simultaneously in an integrated fashion.
- ▶ with the LVQ2 Loss :
  - ▶ Driancourt and Bottou's speech recognizer (1991)
- ▶ with NLL:
  - ▶ Bengio's speech recognizer (1992)

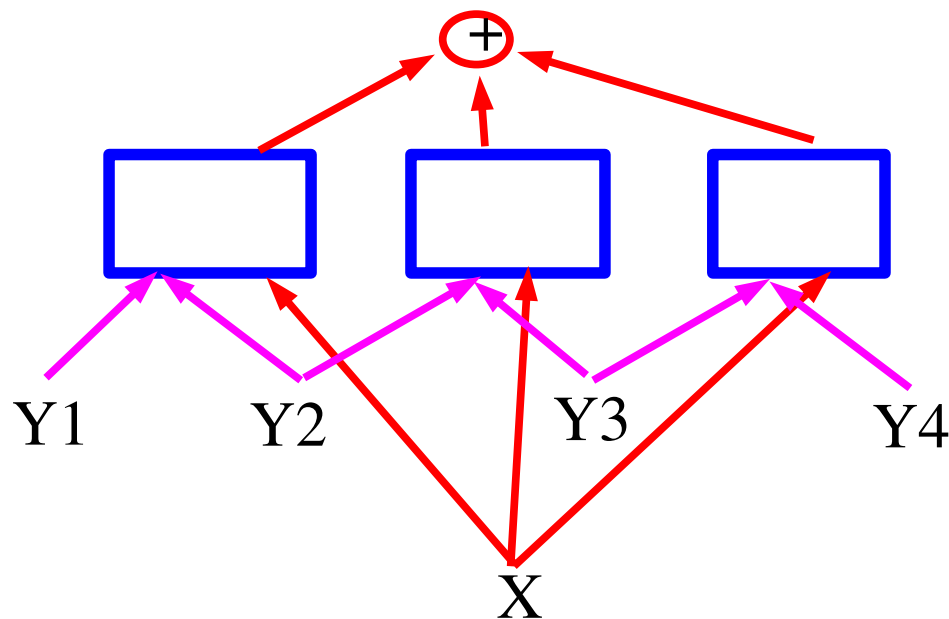


# Energy-Based Factor Graphs: Energy = Sum of “factors”

## ► Sequence Labeling

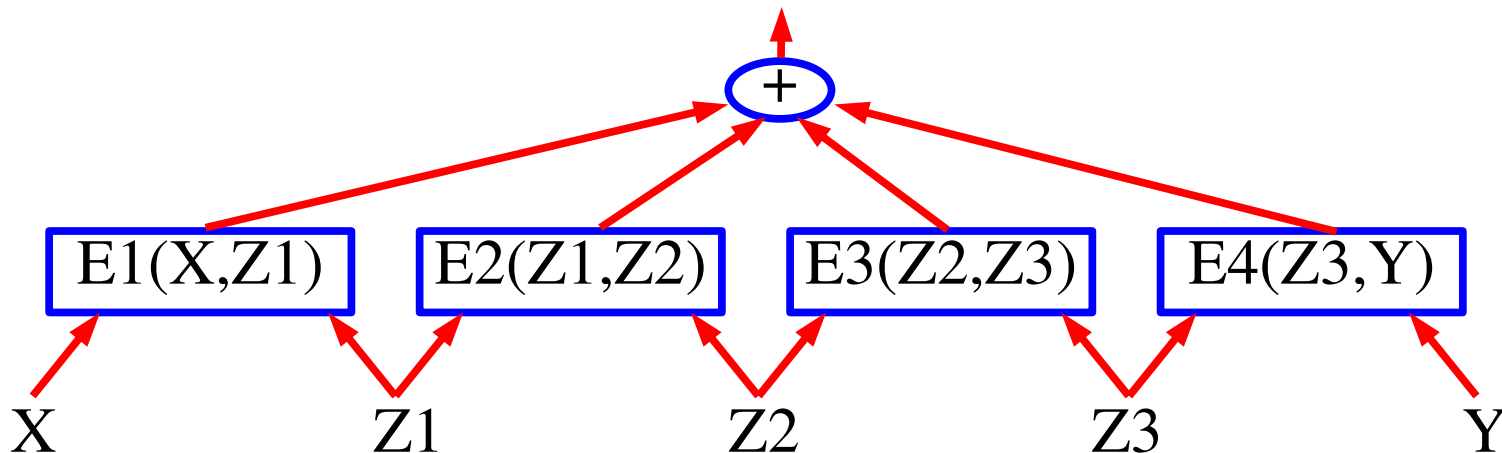
- Output is a sequence  $Y_1, Y_2, Y_3, Y_4, \dots$
- NLP parsing, MT, speech/handwriting recognition, biological sequence analysis
- The factors ensure grammatical consistency
- They give low energy to consistent sub-sequences of output symbols
- The graph is generally simple (chain or tree)
- Inference is easy (dynamic programming, min-sum)

$$Y^* = \operatorname{argmin}_{Y \in \mathcal{Y}, Z \in \mathcal{Z}} E(Z, Y, X).$$



# Energy-Based Factor Graphs

- ▶ When the energy is a sum of partial energy functions (or when the probability is a product of factors):
  - ▶ Efficient inference algorithms can be used for inference (without the normalization step).

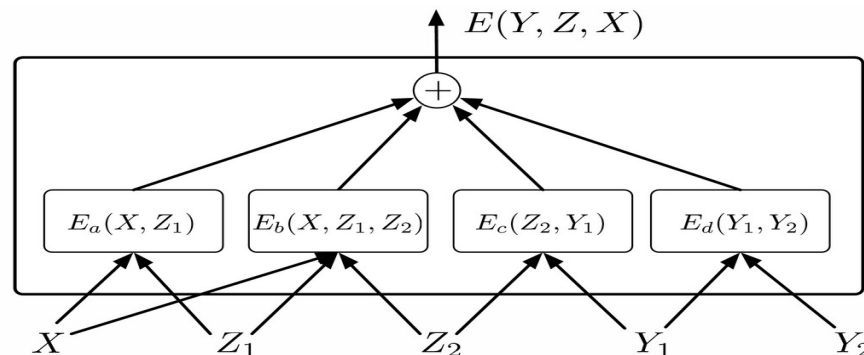


# Efficient Inference: Energy-Based Factor Graphs

## Example:

- ▶  $Z_1, Z_2, Y_1$  are binary
- ▶  $Z_2$  is ternary
- ▶ A naïve exhaustive inference would require  $2 \times 2 \times 2 \times 3 = 24$  energy evaluations (= 96 factor evaluations)
- ▶ BUT:  $E_a$  only has 2 possible input configurations,  $E_b$  and  $E_c$  have 4, and  $E_d$  6.

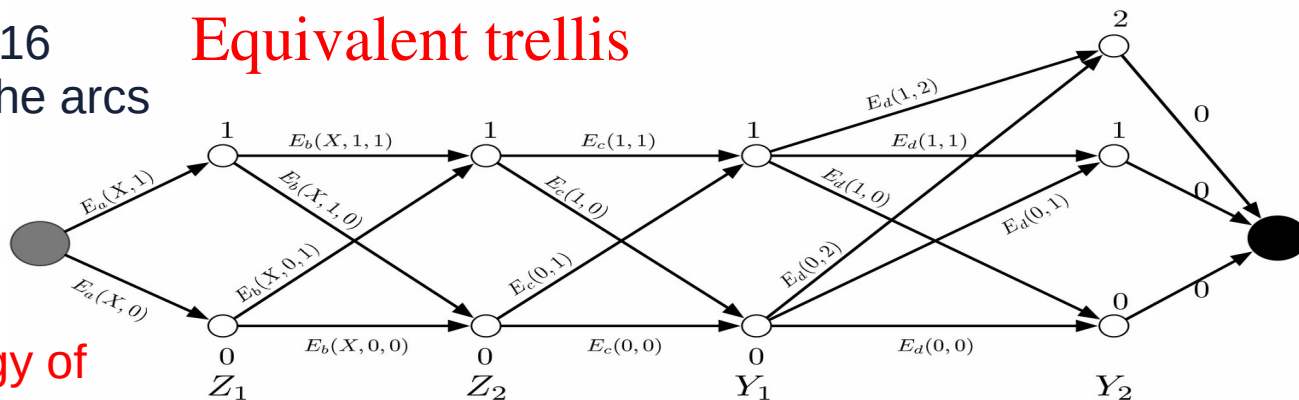
- ▶ The energy is a sum of “factor” functions
- Factor graph



- ▶ Hence, we can precompute the 16 factor values, and put them on the arcs in a trellis.

## Equivalent trellis

- ▶ A path in the trellis is a config of variable
- ▶ The cost of the path is the energy of the config



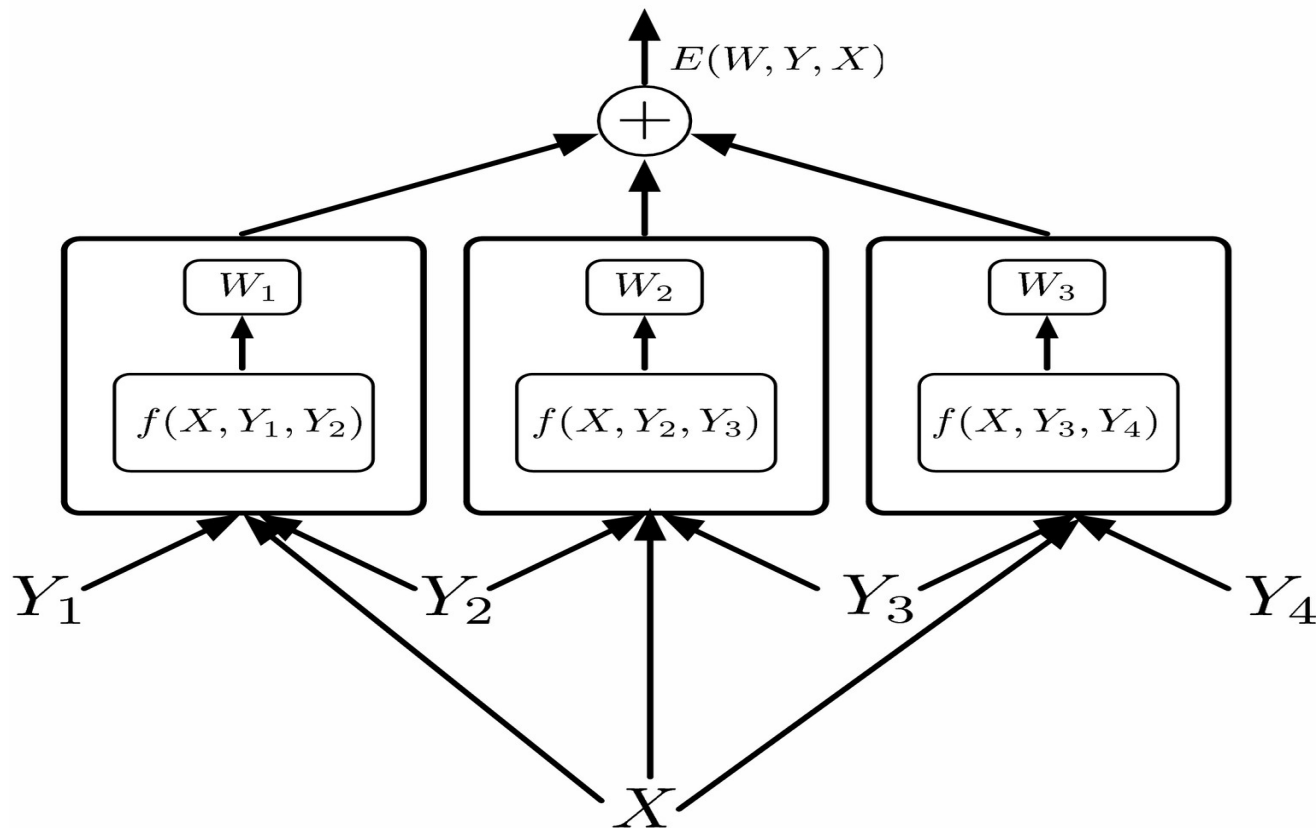
# Energy-Based Belief Prop

- ▶ The previous picture shows a chain graph of factors with 2 inputs.
- ▶ The extension of this procedure to trees, with factors that can have more than 2 inputs the “min-sum” algorithm (a non-probabilistic form of belief propagation)
- ▶ Basically, it is the sum-product algorithm with a different semi-ring algebra (min instead of sum, sum instead of product), and no normalization step.
- ▶ [Kschischang, Frey, Loeliger, 2001][McKay's book]

# Simple Energy-Based Factor Graphs with “Shallow” Factors

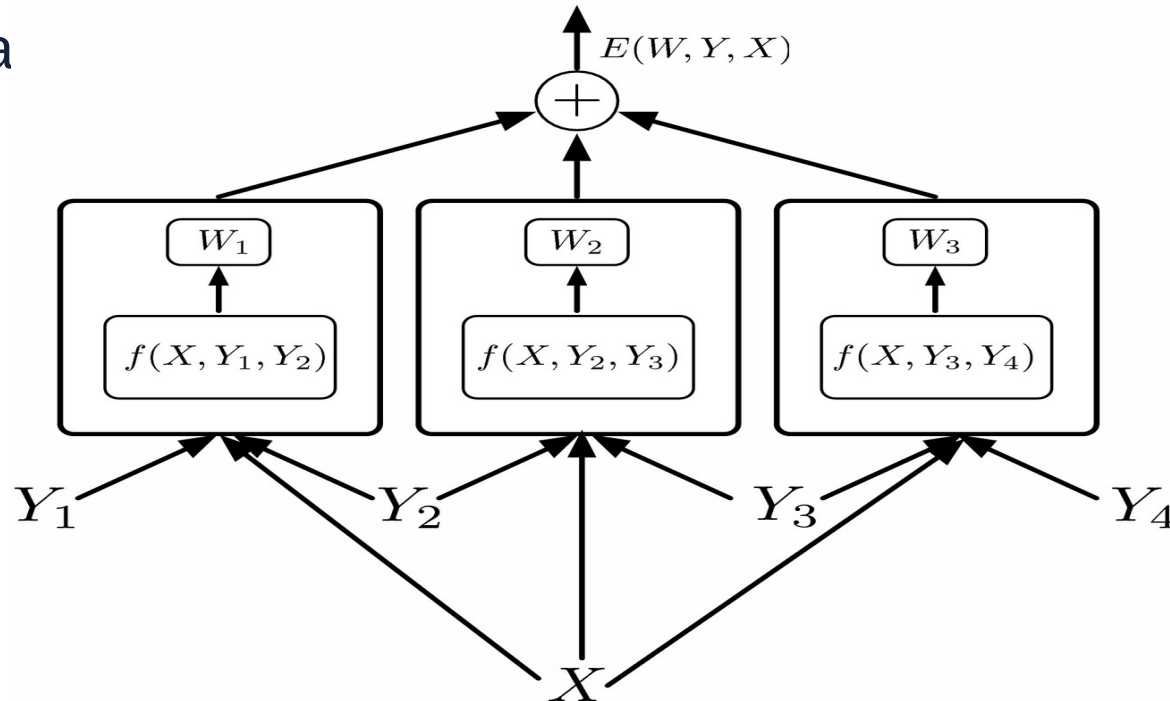
## ► Linearly Parameterized Factors

- with the NLL Loss :
  - Lafferty's **Conditional Random Field**
- with Hinge Loss:
  - Taskar and Altun/Hofmann's **Max Margin Markov Nets** and **Latent SVM**
- with Perceptron Loss
  - Collins's **Structured Perceptron** model



# Example : The Conditional Random Field Architecture

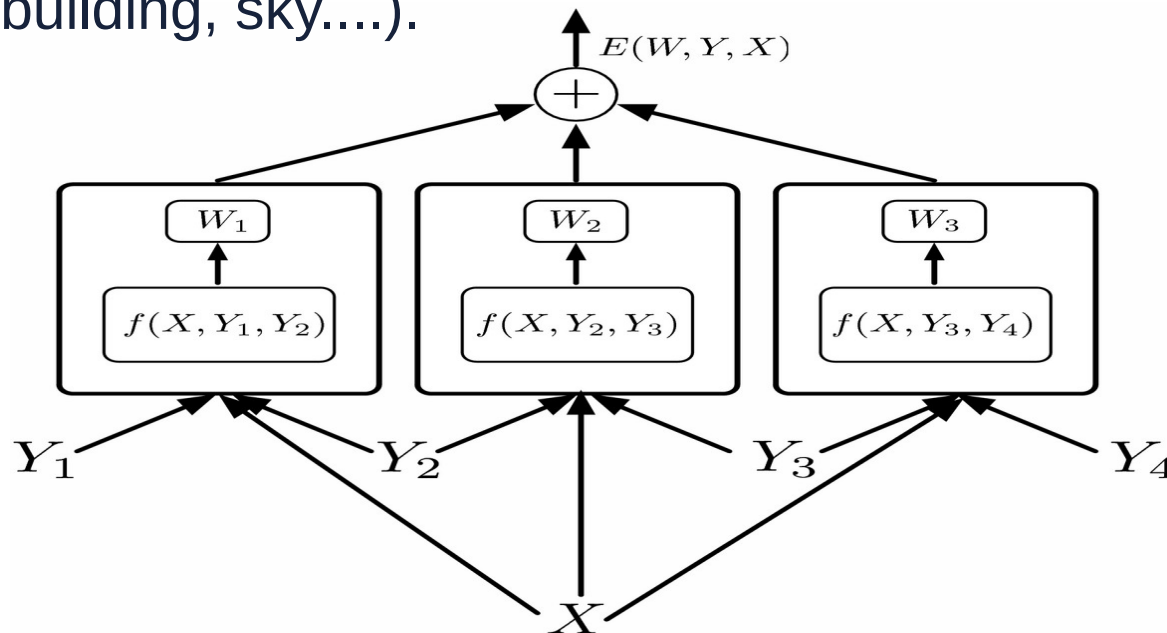
- ▶ A CRF is an energy-based factor graph in which:
  - ▶ the factors are **linear in the parameters** (shallow factors)
  - ▶ The factors take neighboring output variables as inputs
  - ▶ The factors a



# Example : The Conditional Random Field Architecture

## ► Applications:

- $X$  is a sentence,  $Y$  is a sequence of Parts of Speech Tags (there is one  $Y_i$  for each possible group of words).
- $X$  is an image,  $Y$  is a set of labels for each window in the image (vegetation, building, sky....).



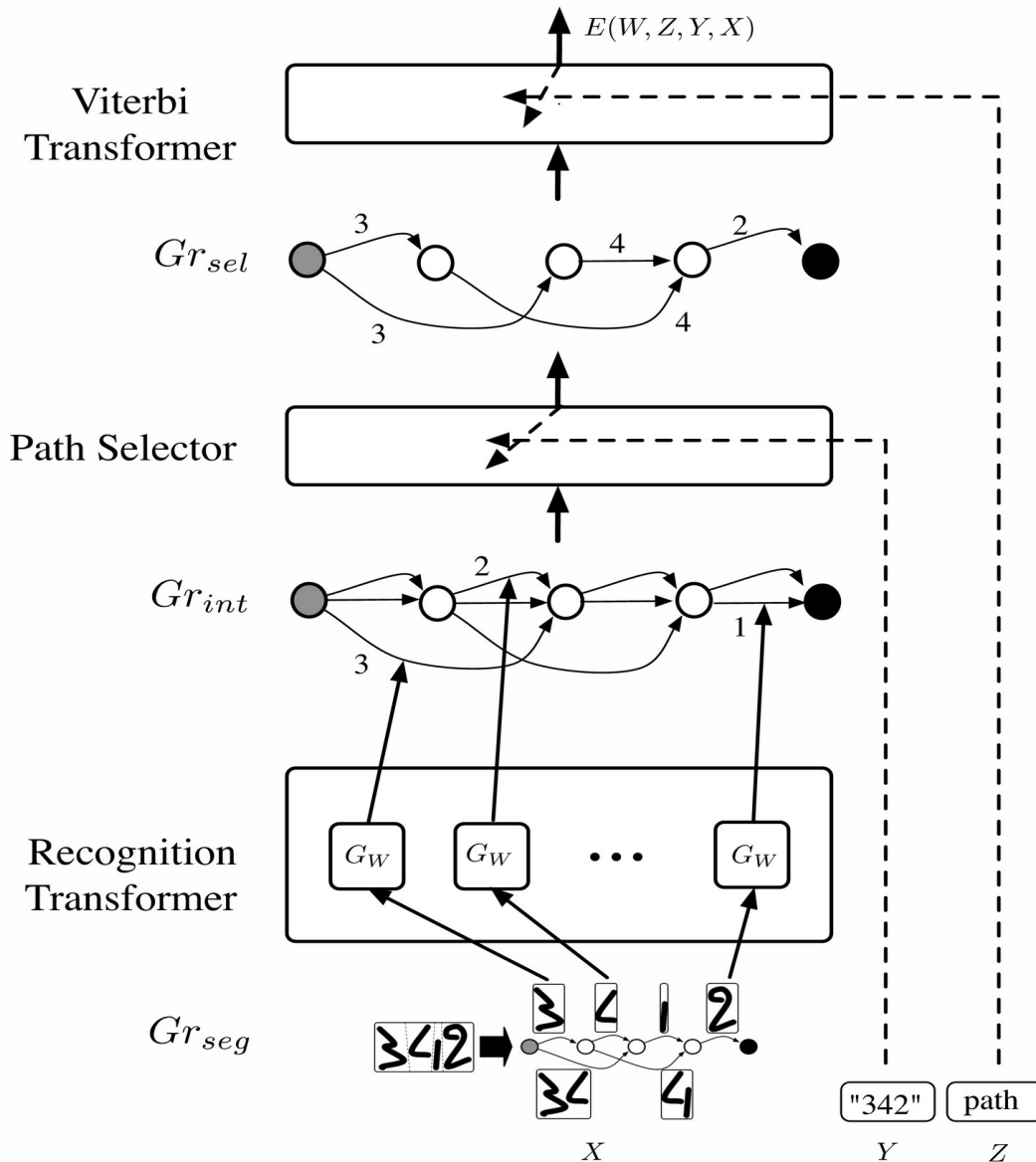


# Deep/non-linear Factors for Speech and Handwriting

- ▶ **Trainable Speech/Handwriting Recognition systems that integrate Neural Nets (or other “deep” classifiers) with dynamic time warping, Hidden Markov Models, or other graph-based hypothesis representations**
- ▶ **Training the feature extractor as part of the whole process.**
  - ▶ **With Minimum Empirical Error loss**
    - ▶ Ljolje and Rabiner (1990)
  - ▶ **with NLL:**
    - ▶ Bengio (1992), Haffner (1993), Bourlard (1994)
- ▶ **with the LVQ2 Loss :**
  - ▶ Driancourt and Bottou's speech recognizer (1991)
- ▶ **with NLL:**
  - ▶ Bengio's speech recognizer (1992)
  - ▶ Haffner's speech recognizer (1993)
- ▶ **With MCE**
  - ▶ Juang et al. (1997)
- ▶ **Late normalization scheme (un-normalized HMM)**
  - ▶ Bottou pointed out the **label bias problem** (1991)
  - ▶ Denker and Burges proposed a solution (1995)

# Deep Factors & implicit graphs: GTN

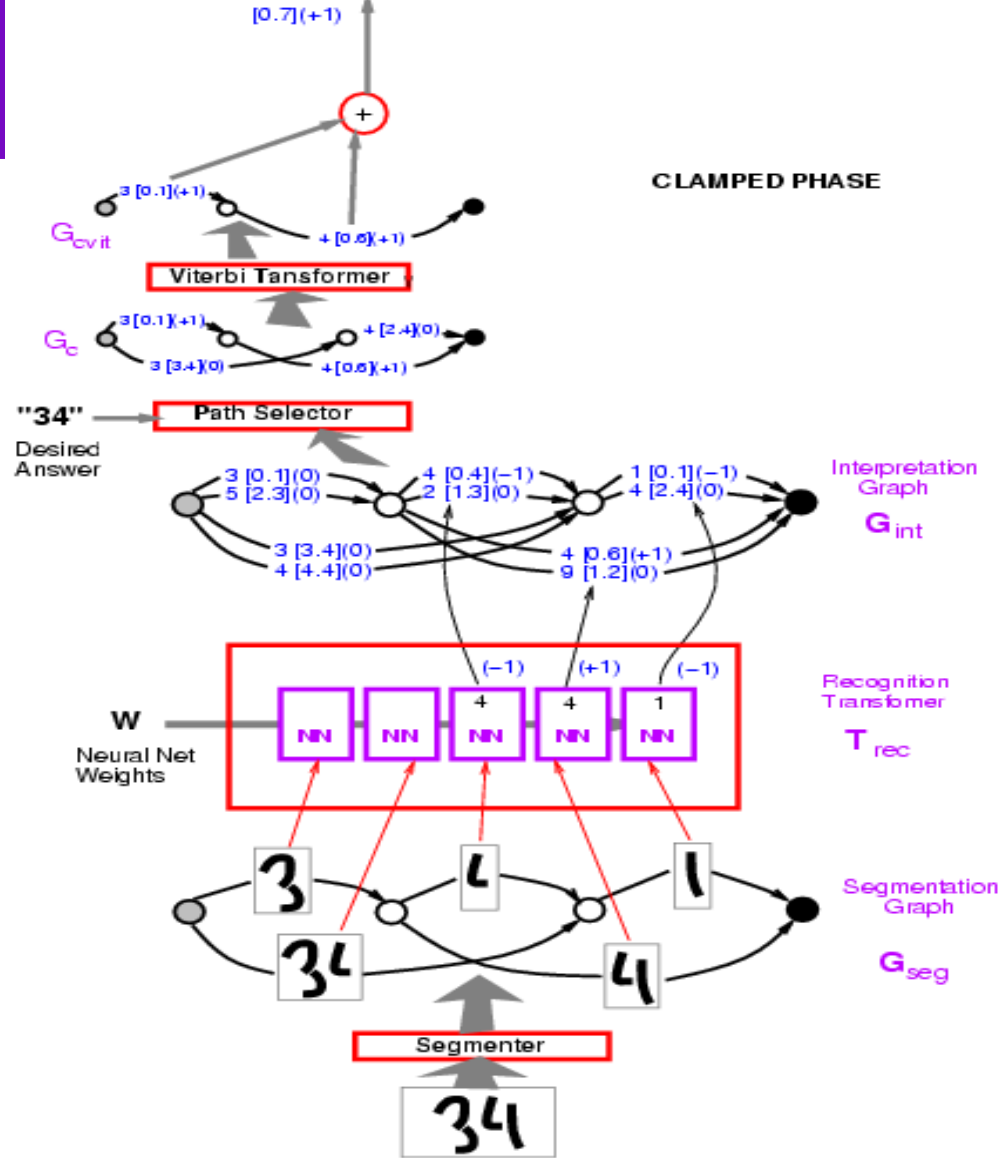
- ▶ Handwriting Recognition with **Graph Transformer Networks**
- ▶ Un-normalized hierarchical HMMs
  - ▶ Trained with Perceptron loss [LeCun, Bottou, Bengio, Haffner 1998]
  - ▶ Trained with NLL loss [Bengio, LeCun 1994], [LeCun, Bottou, Bengio, Haffner 1998]
- ▶ Answer = sequence of symbols
- ▶ Latent variable = segmentation



# Graph Transformer Networks

- **Variables:**
  - X: input image
  - Z: path in the interpretation graph/segmentation
  - Y: sequence of labels on a path
- **Loss function: computing the energy of the desired answer:**

$$E(W, Y, X)$$



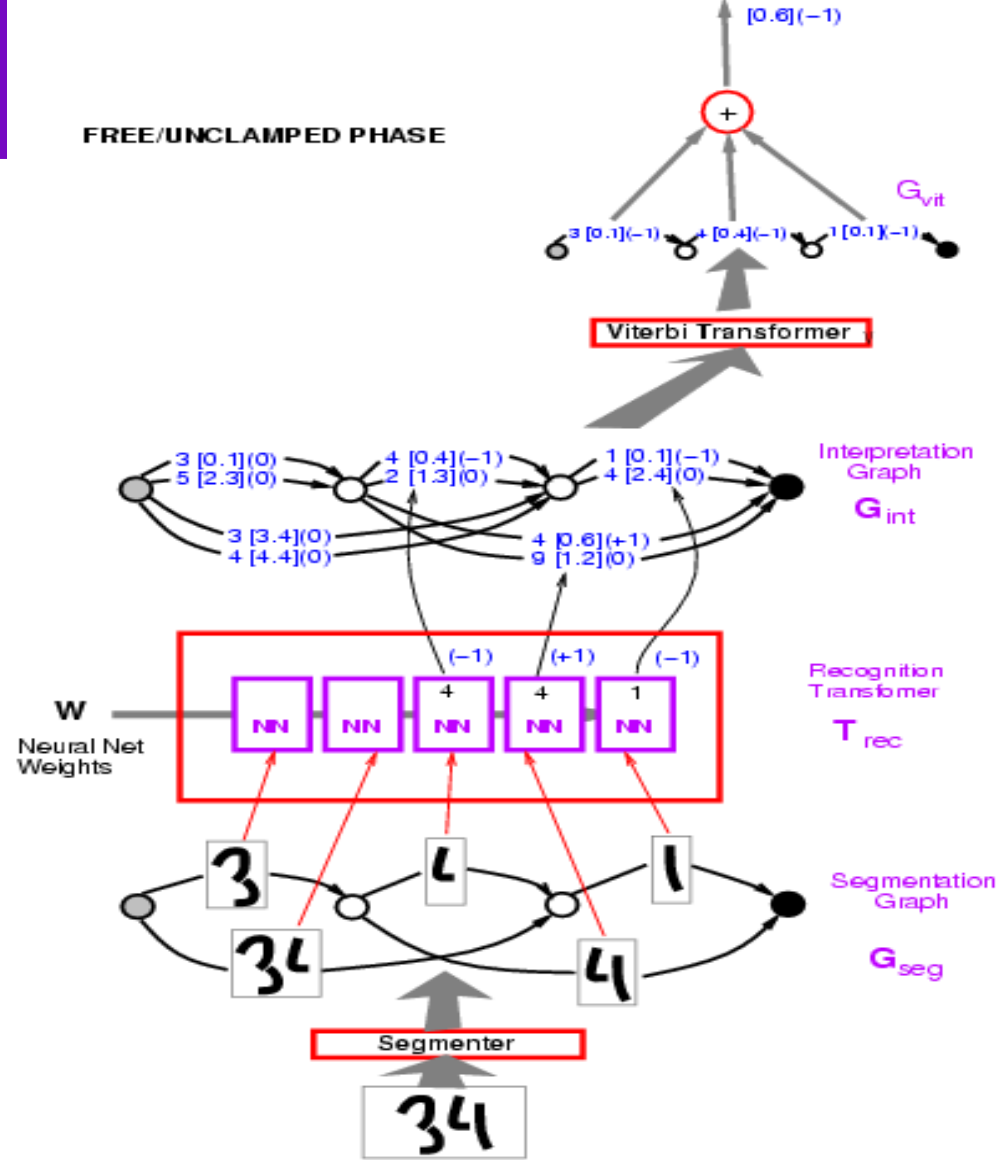
# Graph Transformer Networks

## ► Variables:

- X: input image
- Z: path in the interpretation graph/segmentation
- Y: sequence of labels on a path

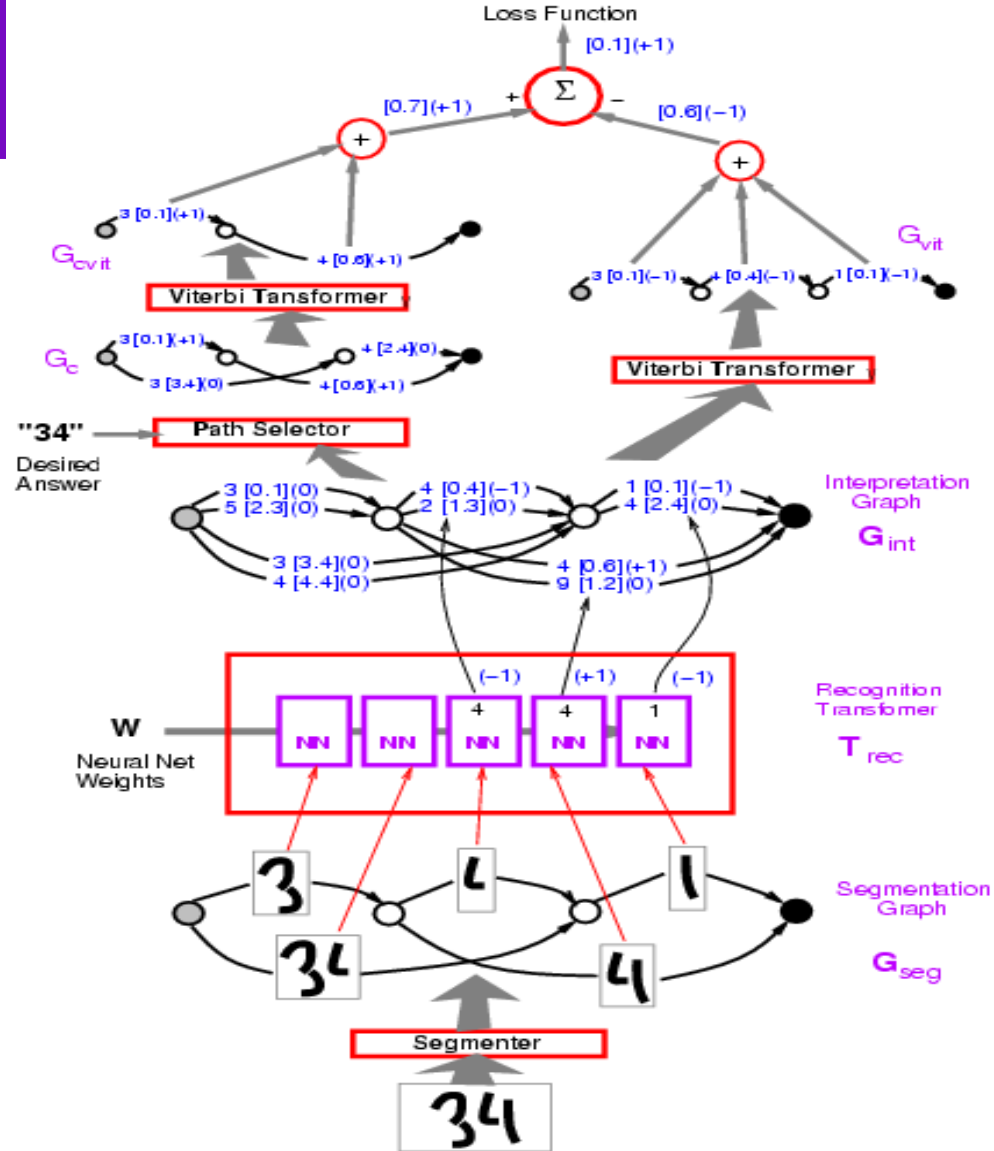
## ► Loss function: computing the constrastive term:

$$E(W, \check{Y}, X)$$



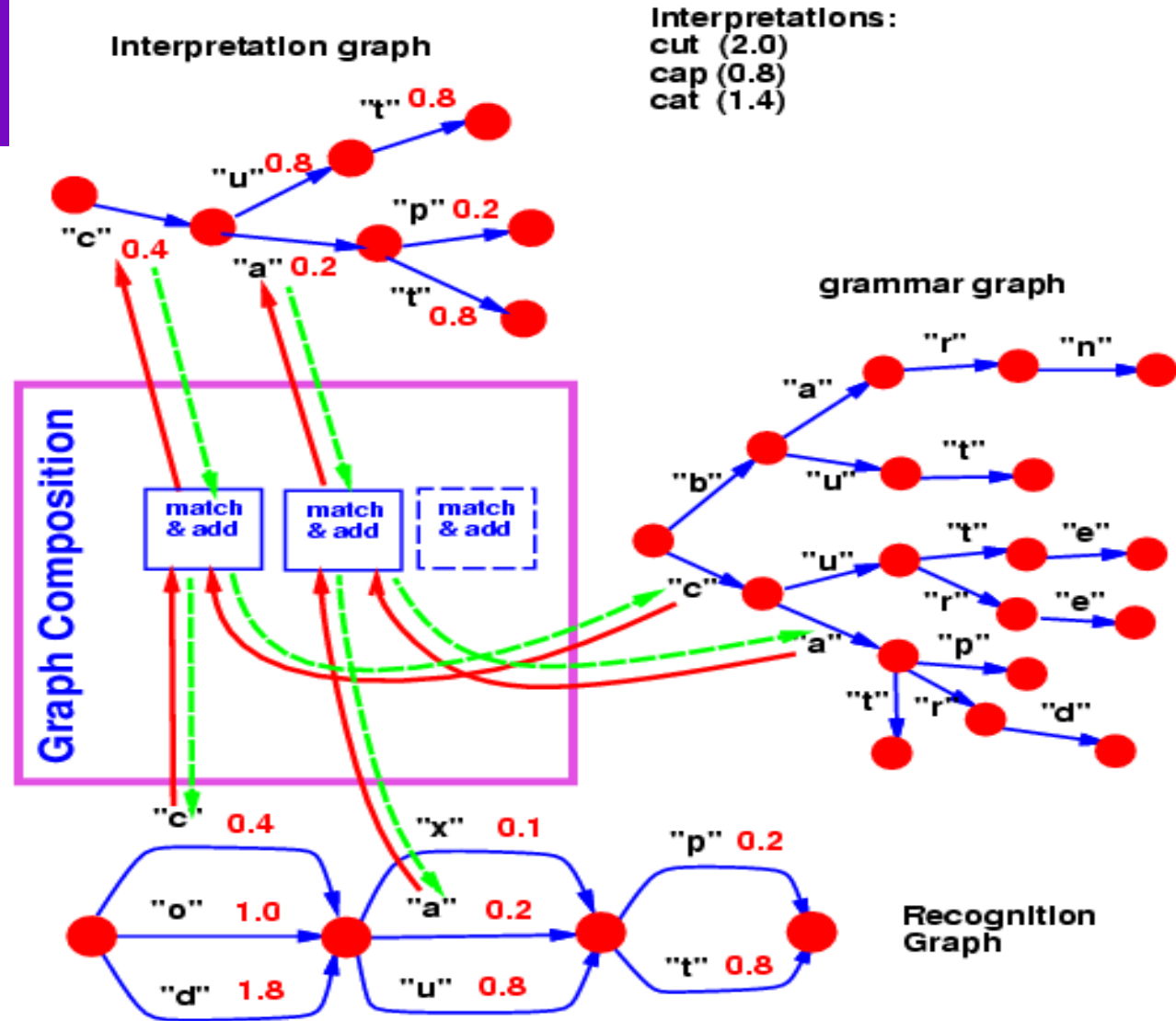
# Graph Transformer Networks

- ▶ **Example: Perceptron loss**
- ▶ **Loss = Energy of desired answer – Energy of best answer.**
- ▶ (no margin)



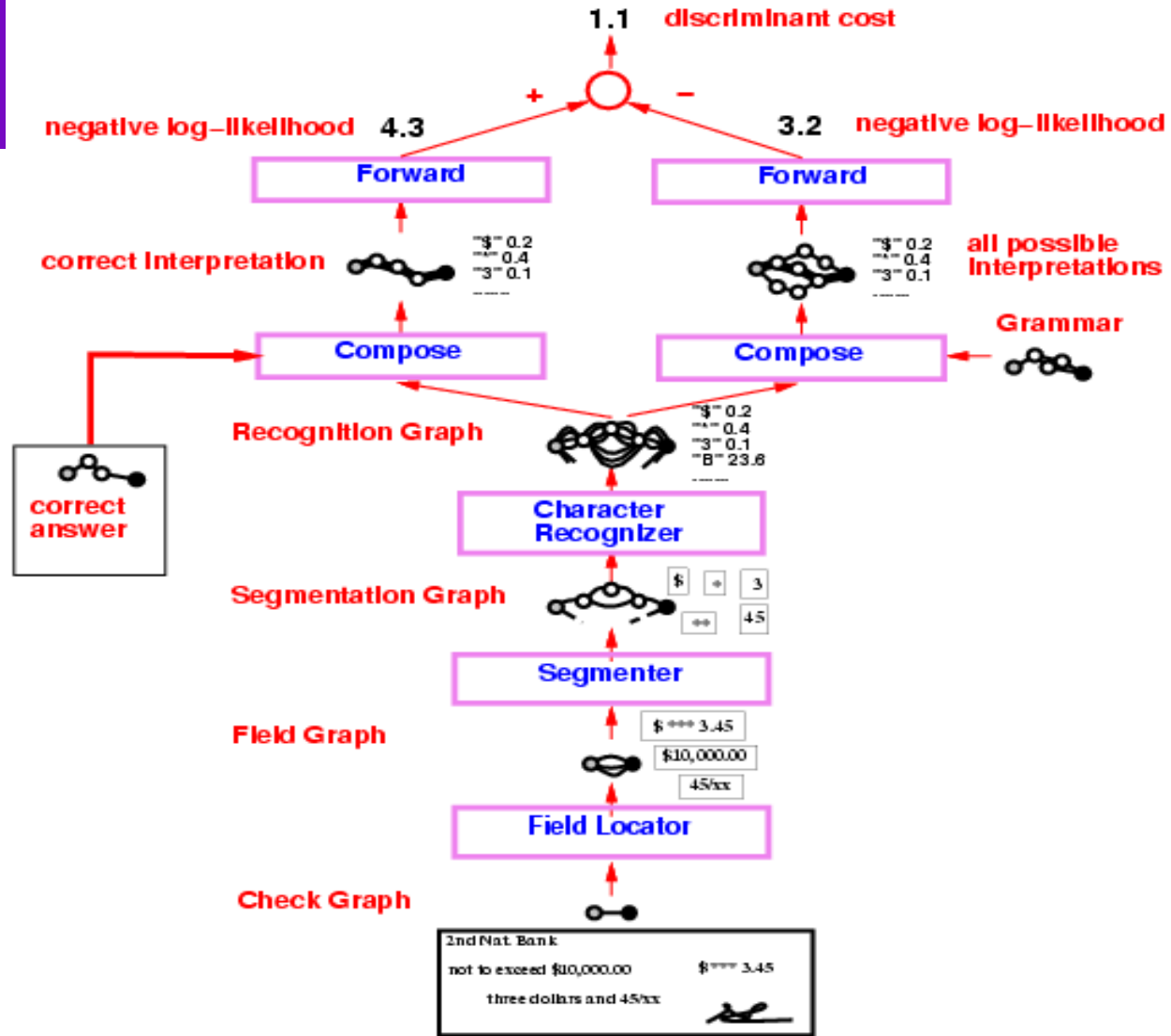
# Graph Composition, Transducers.

- ▶ The composition of two graphs can be computed, the same way the dot product between two vectors can be computed.
- ▶ General theory: semi-ring algebra on weighted finite-state transducers and acceptors.



# Check Reader

- ▶ Graph transformer network trained to read **check amounts**.
- ▶ Trained globally with Negative-Log-Likelihood loss.
- ▶ 50% percent correct, 49% reject, 1% error (detectable later in the process).
- ▶ **Fielded in 1996**, used in many banks in the US and Europe.
- ▶ Processes an estimated **10% of all the checks written in the US**.



# Deep Factors / Deep Graph: ASR with TDNN/HMM

- ▶ **Discriminative Automatic Speech Recognition system with HMM and various acoustic models**
  - ▶ Training the acoustic model (feature extractor) and a (normalized) HMM in an integrated fashion.
- ▶ **With Minimum Empirical Error loss**
  - ▶ Ljolje and Rabiner (1990)
- ▶ **with NLL:**
  - ▶ Bengio (1992)
  - ▶ Haffner (1993)
  - ▶ Bourlard (1994)
- ▶ **With MCE**
  - ▶ Juang et al. (1997)
- ▶ **Late normalization scheme (un-normalized HMM)**