



NEW YORK UNIVERSITY

Structured Prediction

<http://bit.ly/DLSP20>

Yann LeCun

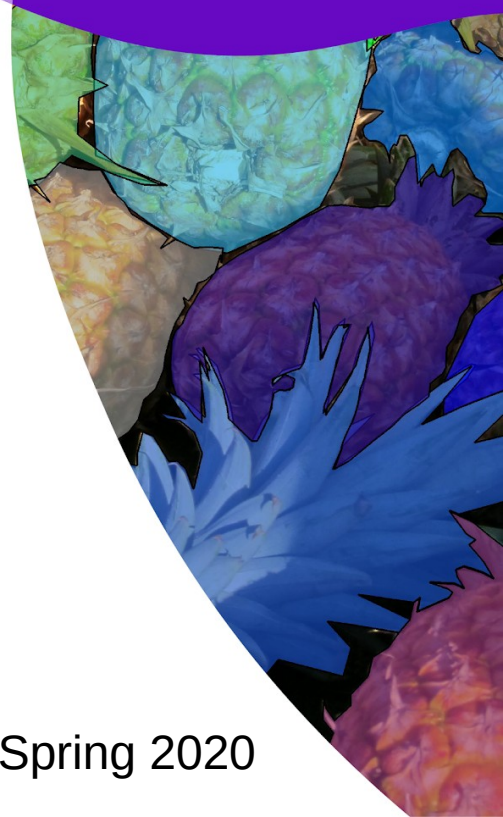
NYU - Courant Institute & Center for Data Science

Facebook AI Research

<http://yann.lecun.com>

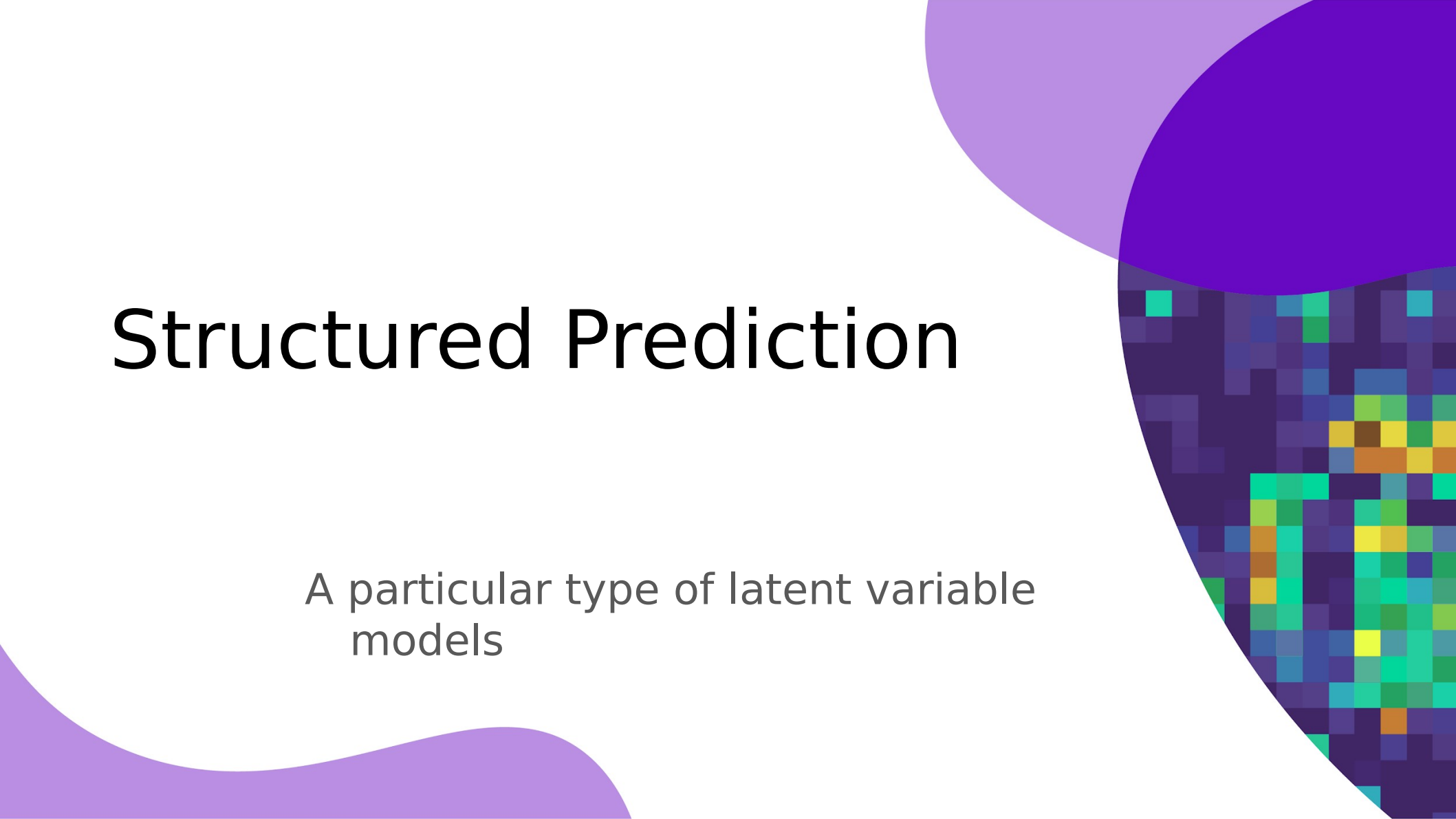
TAs: Alfredo Canziani, Mark Goldstein

Deep Learning, NYU, Spring 2020



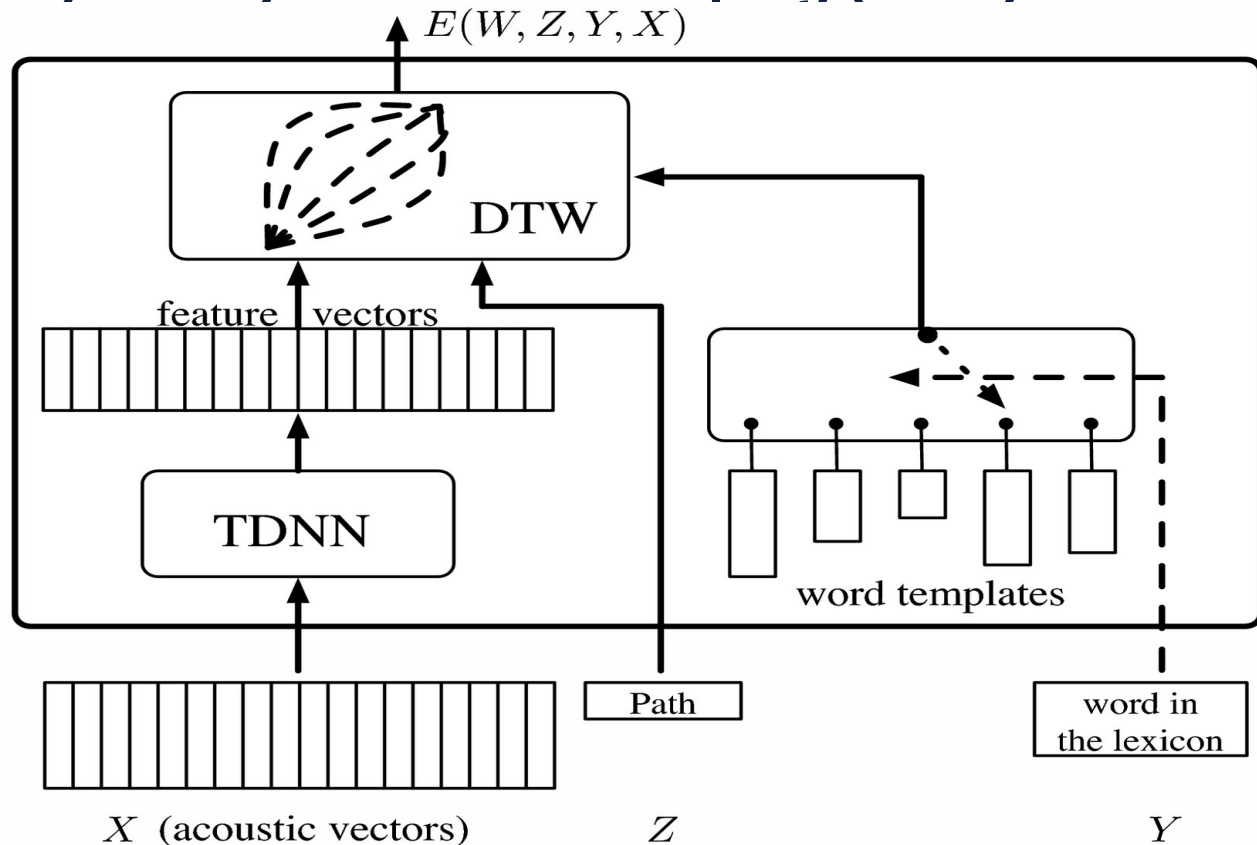
Structured Prediction

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)
 - ▶ Haffner's speech recognizer (1993)

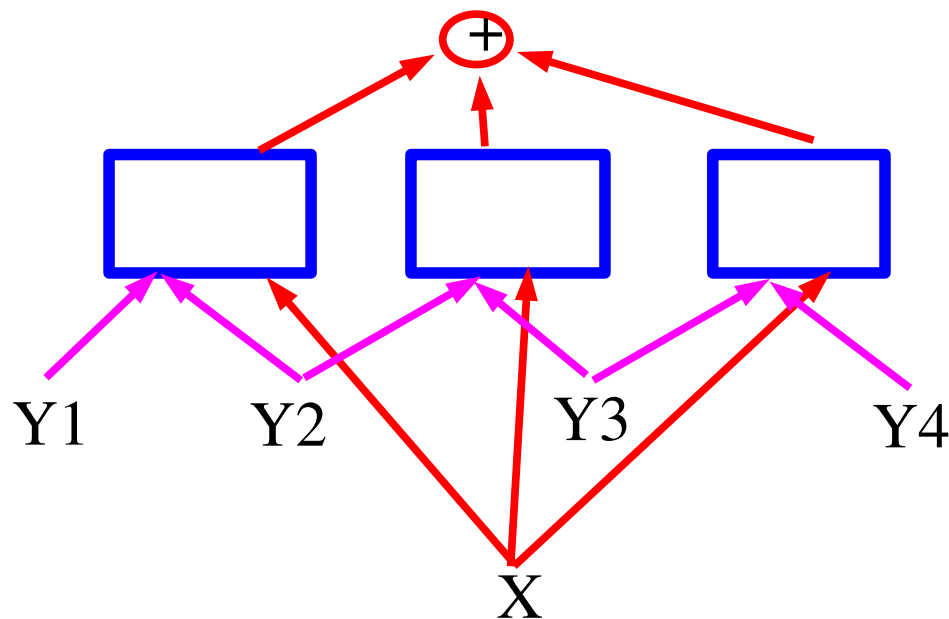


Energy-Based Factor Graphs: Energy = Sum of “factors”

► Sequence Labeling

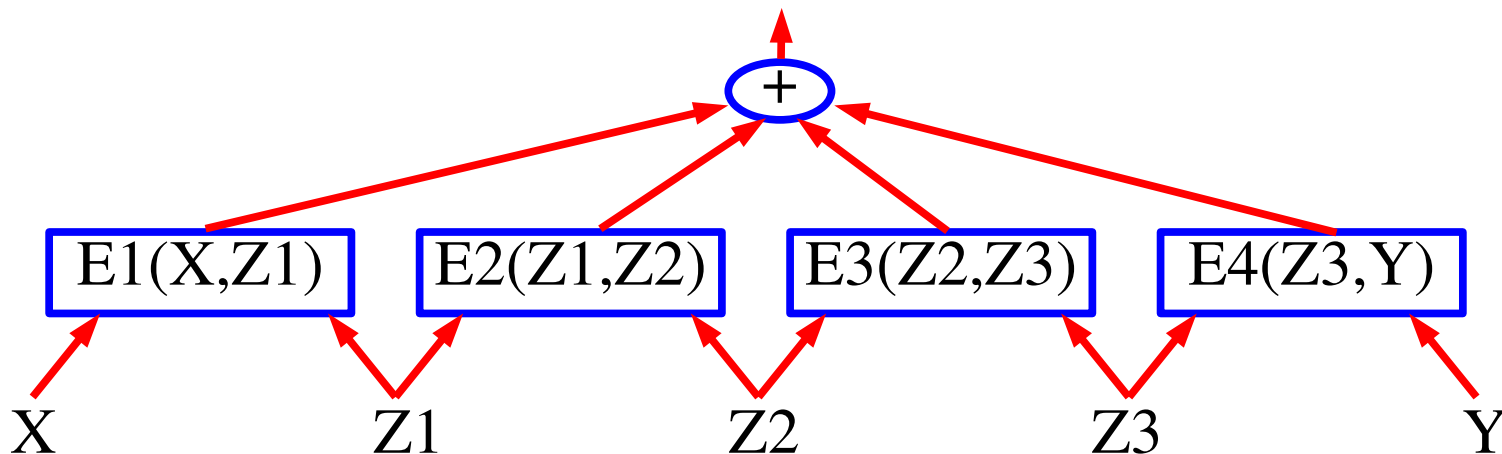
- Output is a sequence $Y_1, Y_2, Y_3, Y_4, \dots$
- NLP parsing, MT, speech/handwriting recognition, biological sequence analysis
- The factors ensure grammatical consistency
- They give low energy to consistent sub-sequences of output symbols
- The graph is generally simple (chain or tree)
- Inference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



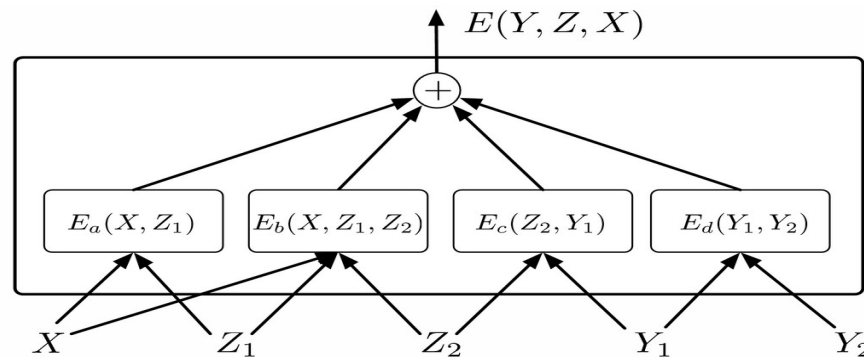
Efficient Inference: Energy-Based Factor Graphs

Example:

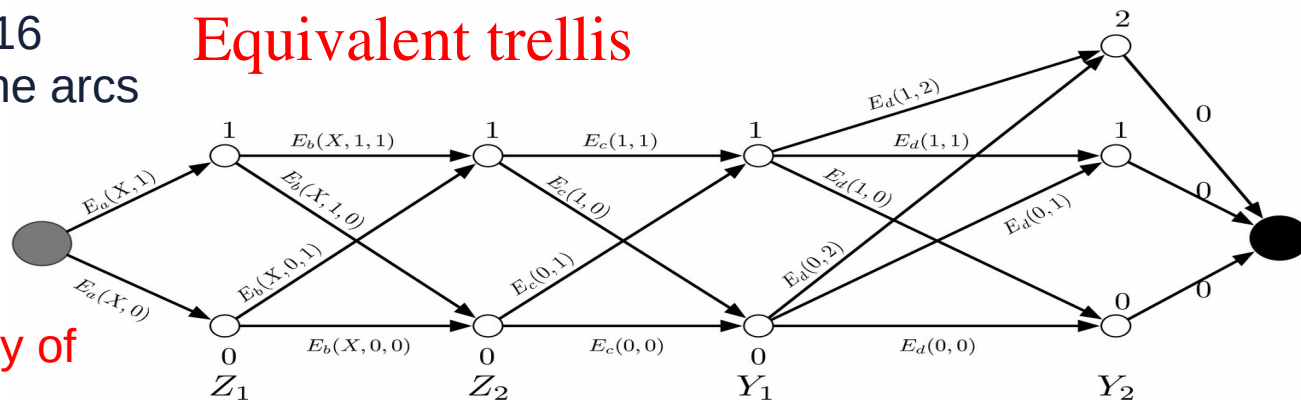
- ▶ Z_1, Z_2, Y_1 are binary
- ▶ Z_2 is ternary
- ▶ A naïve exhaustive inference would require $2 \times 2 \times 2 \times 3 = 24$ energy evaluations (= 96 factor evaluations)
- ▶ BUT: E_a only has 2 possible input configurations, E_b and E_c have 4, and E_d 6.
- ▶ Hence, we can precompute the 16 factor values, and put them on the arcs in a trellis.
- ▶ A path in the trellis is a config of variable
- ▶ The cost of the path is the energy of the config

▶ The energy is a sum of “factor” functions

Factor graph



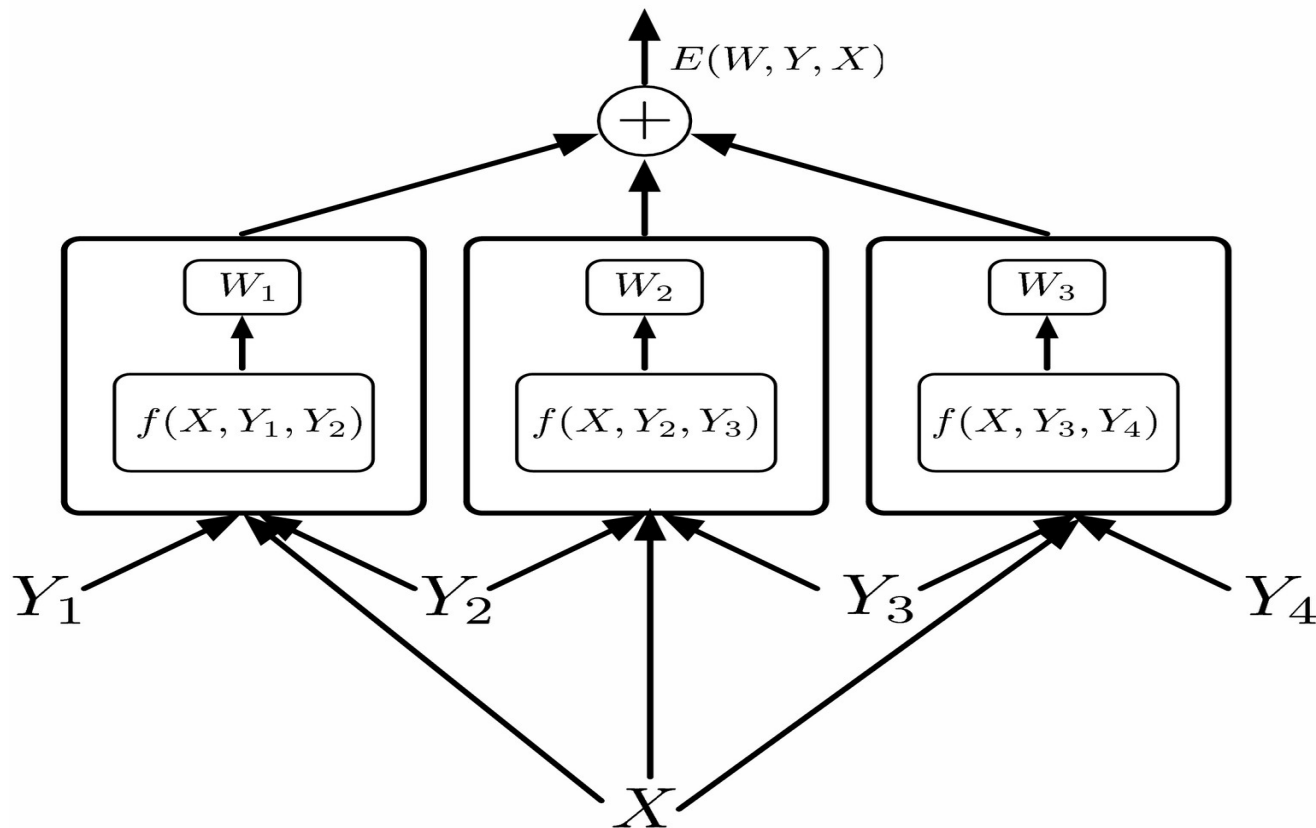
Equivalent trellis



Simple Energy-Based Factor Graphs with “Shallow” Factors

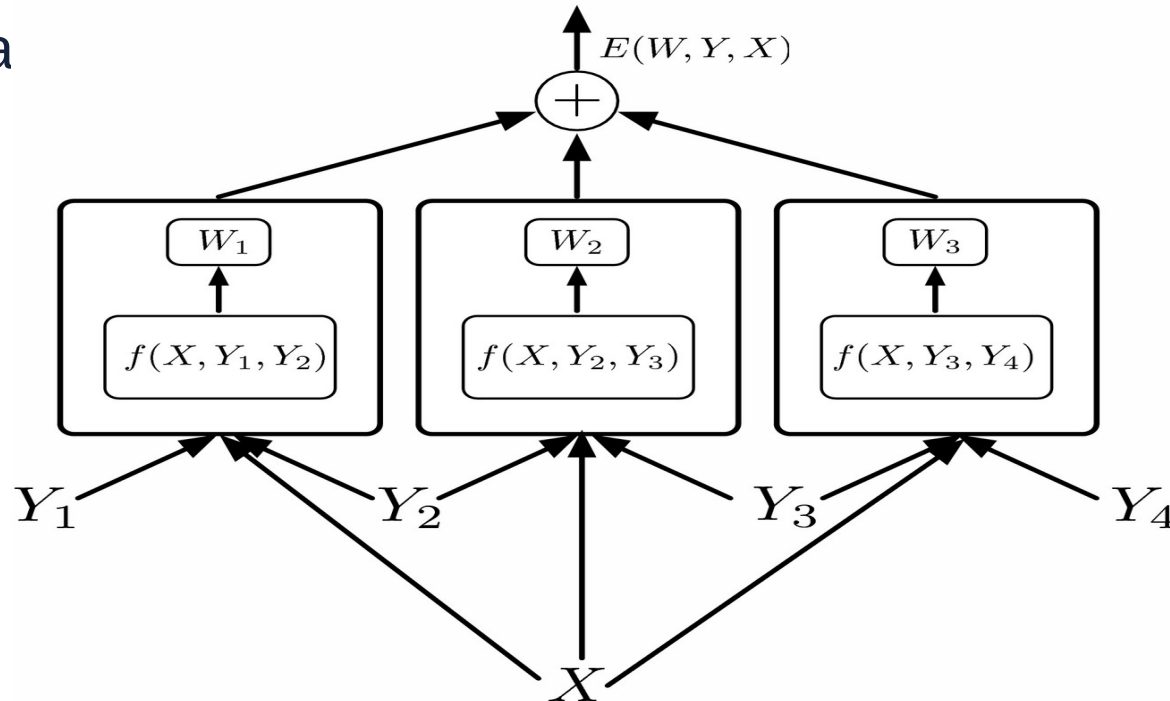
► Linearly Parameterized Factors

- 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

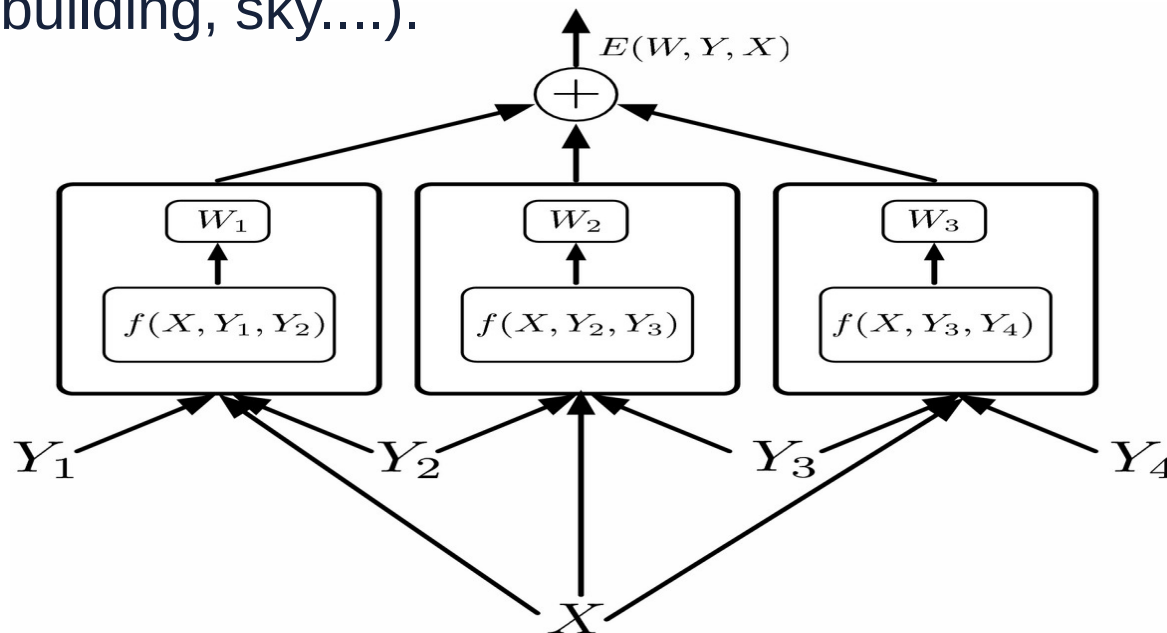
- ▶ A CRF is an energy-based factor graph in which:
 - ▶ the factors are **linear in the parameters** (shallow factors)
 - ▶ The factors take neighboring output variables as inputs
 - ▶ The factors a



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....).

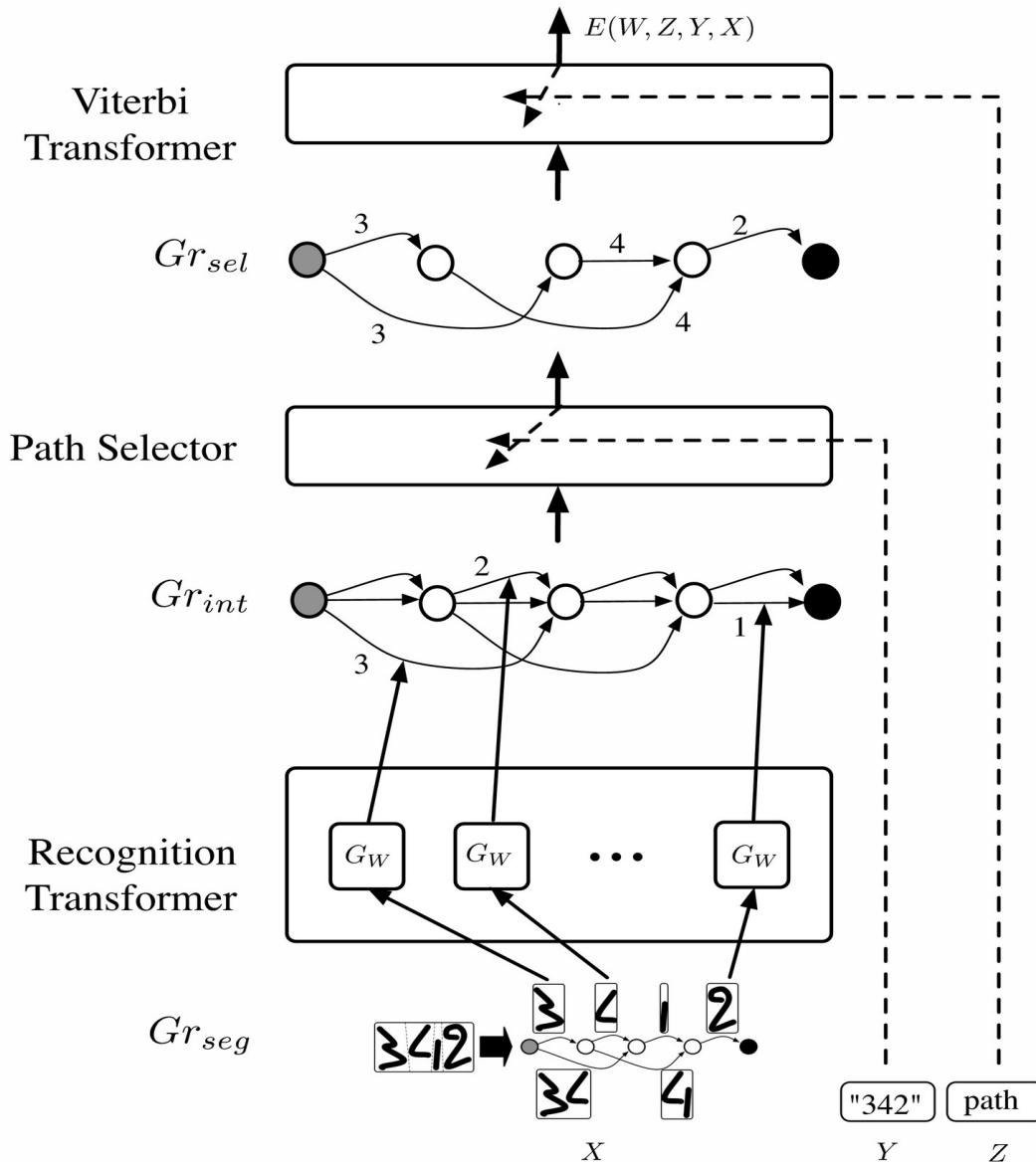


Deep/non-linear Factors for Speech and Handwriting

- ▶ **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

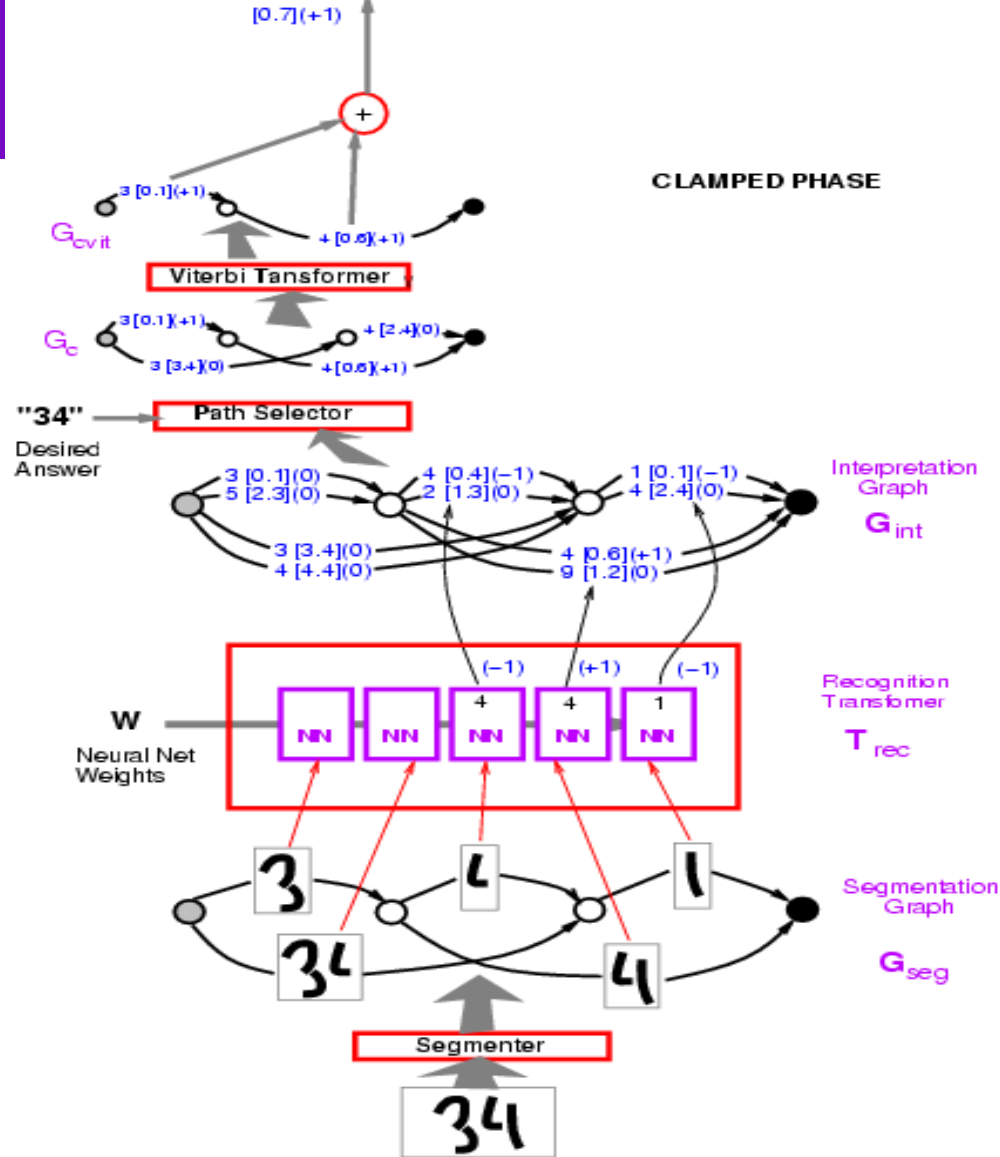
- ▶ Handwriting Recognition with **Graph Transformer Networks**
- ▶ Un-normalized hierarchical HMMs
- ▶ Trained with Perceptron loss [LeCun, Bottou, Bengio, Haffner 1998]
- ▶ Trained with NLL loss [Bengio, LeCun 1994], [LeCun, Bottou, Bengio, Haffner 1998]
- ▶ Answer = sequence of symbols
- ▶ Latent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



Graph Transformer Networks

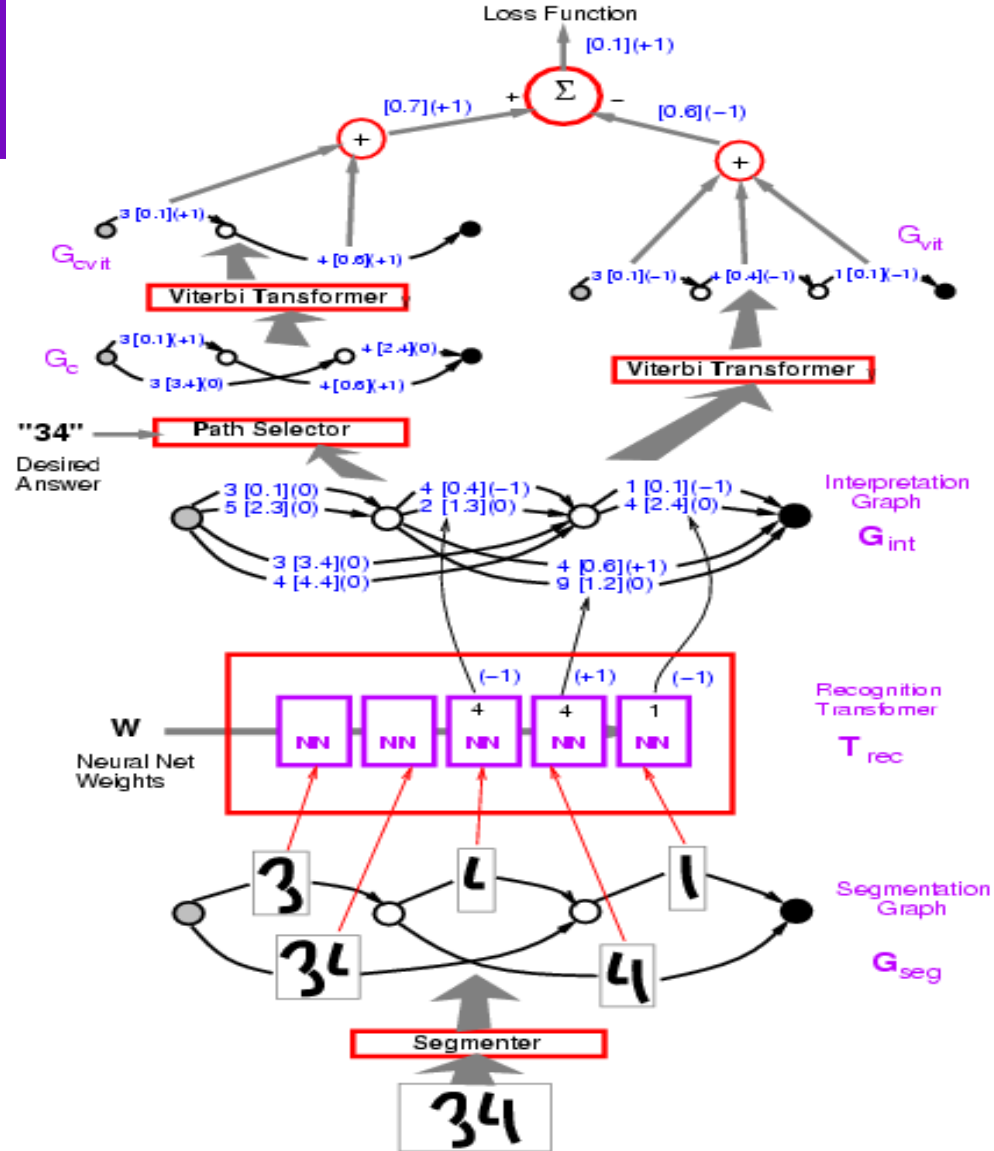
- **Variables:**
 - X: input image
 - Z: path in the interpretation graph/segmentation
 - Y: sequence of labels on a path
- **Loss function: computing the energy of the desired answer:**

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



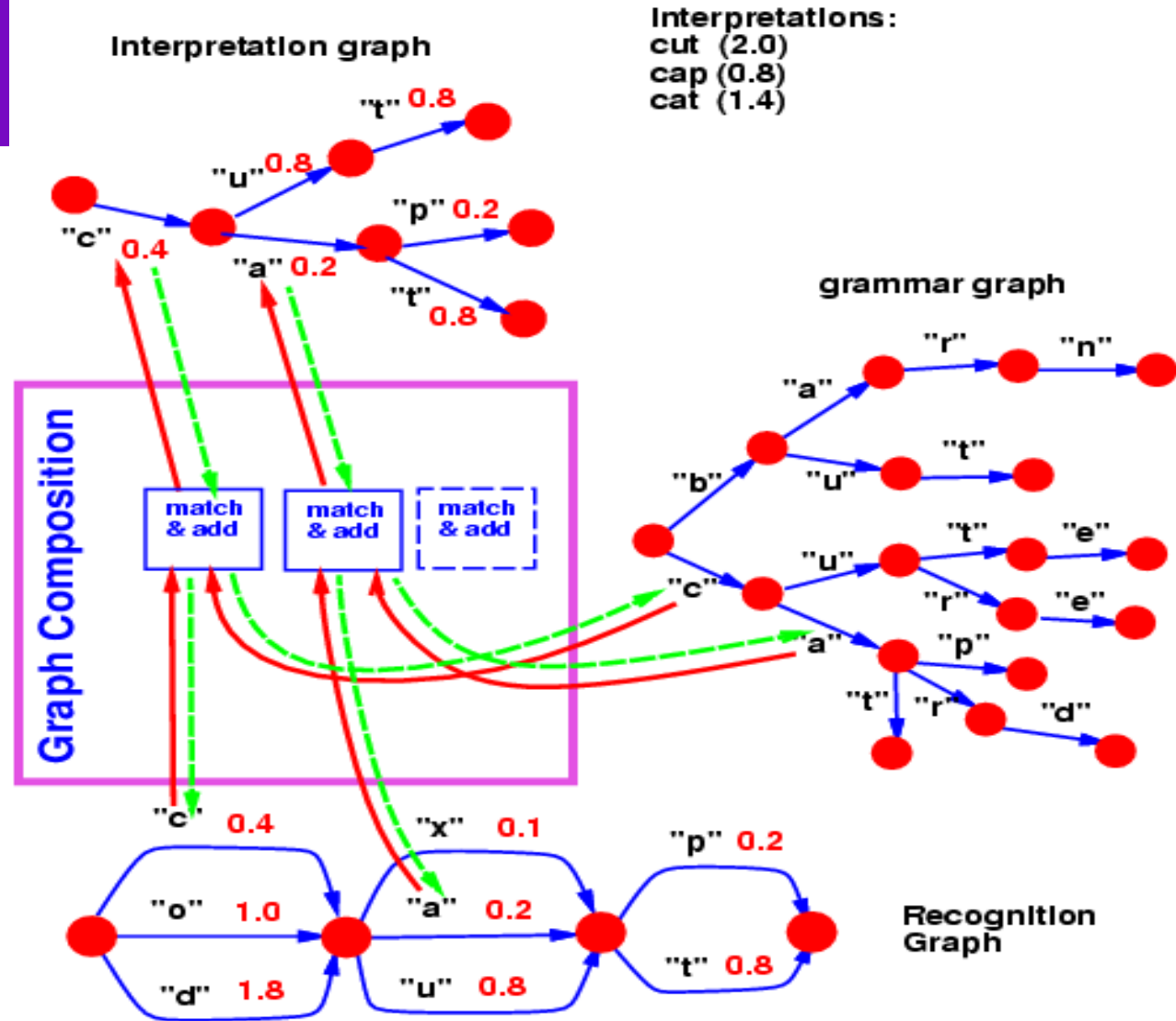
Graph Transformer Networks

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



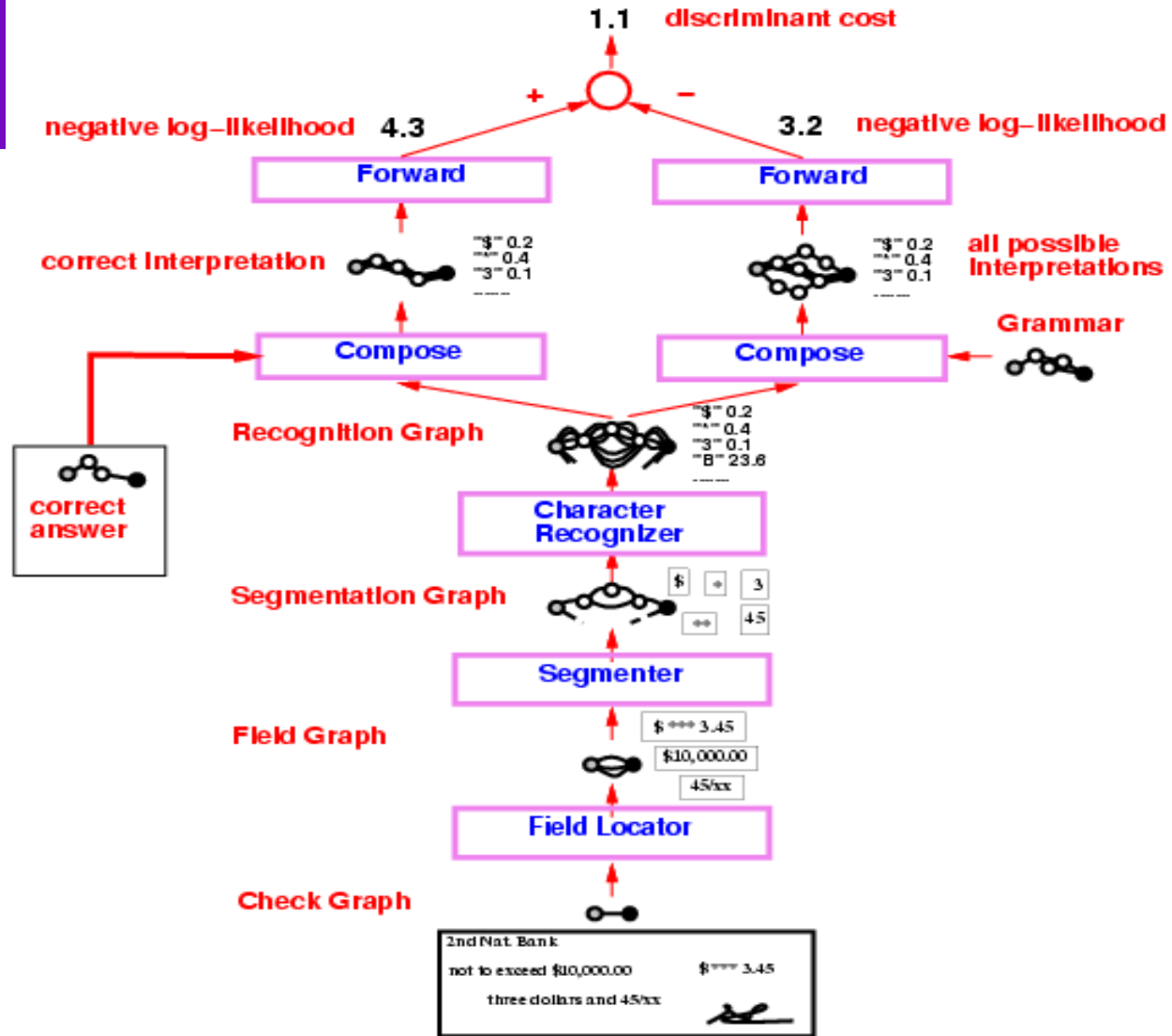
Graph Composition, Transducers.

- ▶ 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)**
 - ▶ Bottou pointed out the **label bias problem** (1991)
 - ▶ Denker and Burges proposed a solution (1995)

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	$\min \left(M, \max(0, E(W, Y^i, X^i) - E(W, \bar{Y}^i, X^i)) \right)$	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