2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)
**Tutorial**

**Energy-Based Models with Applications to Speech and Language Processing**

**Zhijian Ou**

Speech Processing and Machine Intelligence (SPMI) Lab, Tsinghua University
http://oa.ee.tsinghua.edu.cn/ouzhijian/

May 22, 2022

# Introduction

- **Energy-Based Models (EBMs)** are an important class of probabilistic models, also known as **random fields (RFs)** and **undirected graphical models (UGMs)**.

- EBMs have **unique properties** and are radically different from some other popular probabilistic models such as
  - hidden Markov models (HMMs), auto-regressive models, Generative Adversarial Nets (GANs) and Variational Auto-encoders (VAEs), which are self-normalized (i.e., sum to one).

- EBMs have attracted **increasing interests** not only from core machine learning but also from application domains such as vision, speech, natural language processing
  - with significant theoretical and algorithmic progress

- **The sequential nature of speech and language** also presents special challenges and needs treatment different from processing fix-dimensional data (e.g., images).

The purpose of this tutorial is to present a **systematic** introduction to energy-based models, including both algorithmic progress and applications in speech and language processing.

# Content

I. **Basics for EBMs (45 min)**
   1. Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
   2. Learning EBMs by Monte Carlo methods
   3. Learning EBMs by noise-contrastive estimation (NCE)

II. **EBMs for language modeling (45 min)**
   1. Trans-dimensional random field (TRF) LMs for speech recognition
   2. Residual energy-based models for text generation
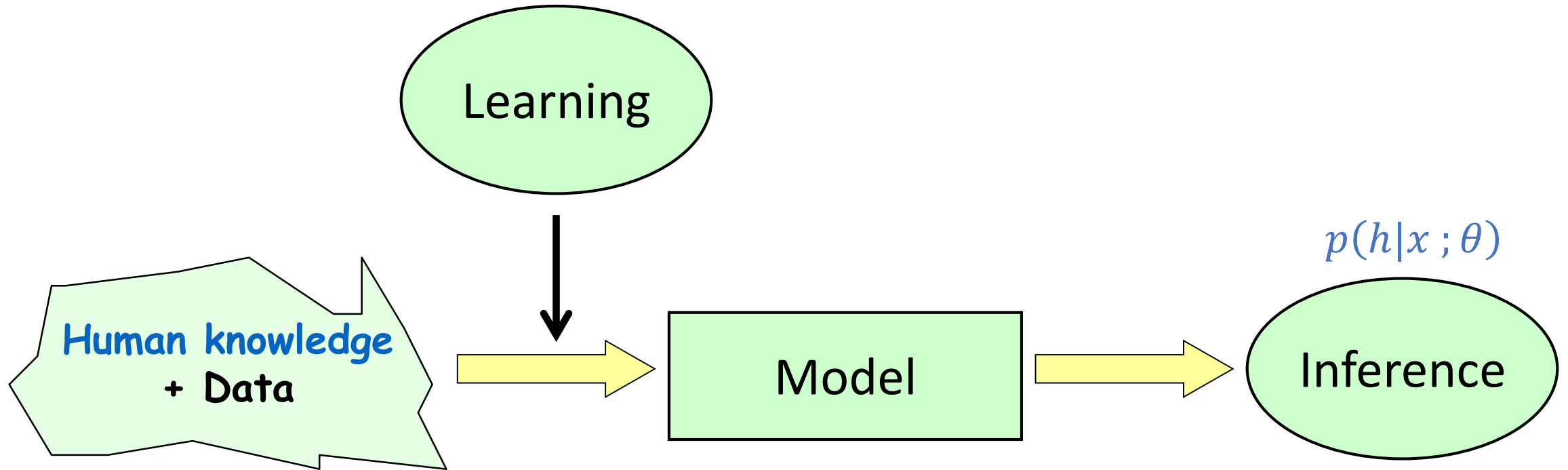   3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for speech recognition and natural language labeling(45 min)**
   1. CRFs as conditional EBMs
   2. CRFs for speech recognition
   3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
   1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
   2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
   3. JRFs for semi-supervised natural language labeling

# Probabilistic Framework



$p(x, h; \theta)$: Generative model, e.g., Hidden Markov Model (HMM)

$p(h|x; \theta)$: Discriminative model, e.g., Conditonal Random Field (CRF)

We need probabilistic models, besides neural nets.

# Roadmap



II. EBMs for language modeling
$p_\theta(x)$

I. Basics for EBMs

$p_\theta(h|x)$

$p_\theta(x, h)$

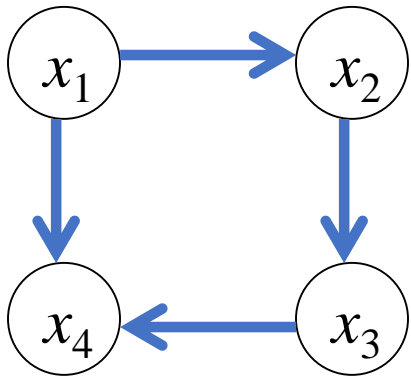III. EBMs for speech recognition and natural language labeling

IV. EBMs for semi-supervised natural language labeling

# Probabilistic Graphical Modeling (PGM) Framework

A general framework for describing and applying probabilistic models

- A graphical model is a family of probability distributions defined in terms of a directed or undirected graph.
- Semantics: how the family of distributions is defined.

- Semantics of Directed Graphical Models (DGMs)

Consider a directed acyclic graph (DAG)

- $x_V$ : a collection of random variables indexed by the nodes
- $pa(v)$ : the parent nodes of $v$

$$p(x_V) = \prod_{v \in V} p(x_v | x_{pa(v)})$$

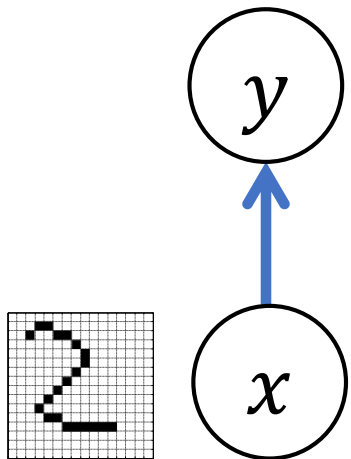$$p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2|x_1)p(x_3|x_2)p(x_4|x_1, x_3)$$

# DGM example - Neural Net (NN) based classifier

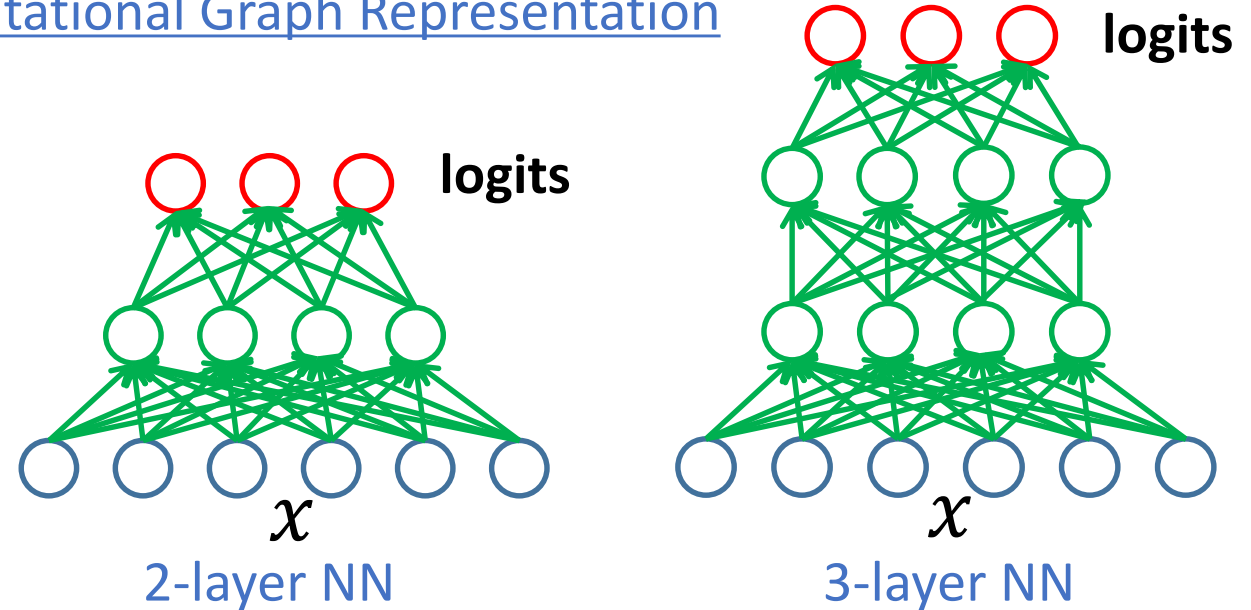- Multi-class logistic regression

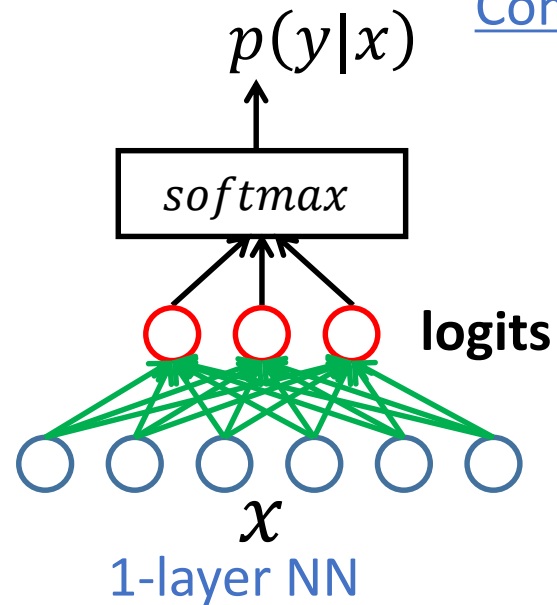Consider observation/features $x \in \mathbb{R}^d$, class label $y \in \{1, \cdots, K\}$

$$p(y = k|x) = \frac{exp(z_k)}{\sum_{j=1}^{K} exp(z_j)} \triangleq softmax(z_k)$$

where $z_k = w_k^T x + b_k, k = 1, \cdots, K,$ often called **logits**

GM Representation

Computational Graph Representation



$p(y|x)$

$softmax$

**logits**

1-layer NN

**logits**

2-layer NN

**logits**

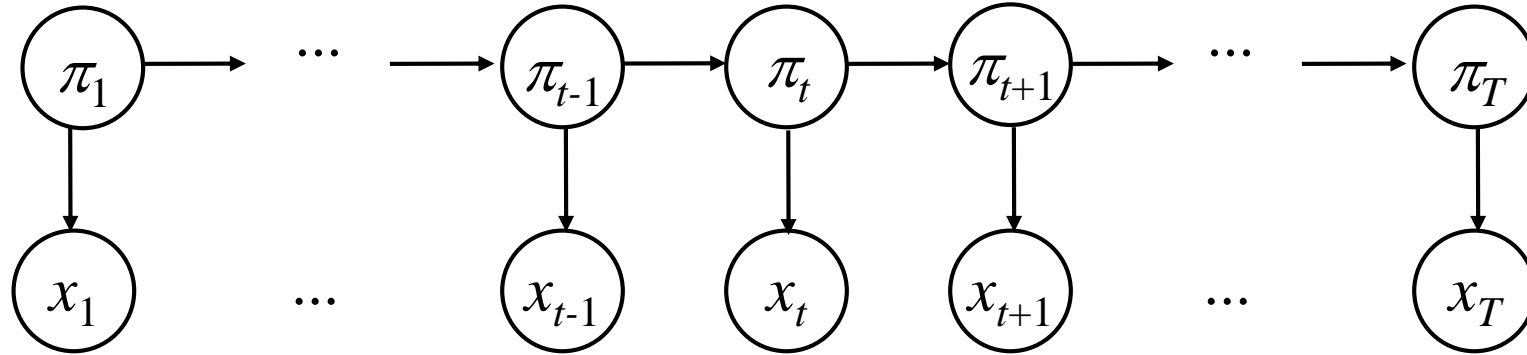3-layer NN

(NNs as feature extractors)

# HMM Viewed as Directed Graphical Model



The joint probability distribution of a hidden Markov model (HMM) :

$$p(\pi_{1:T}, x_{1:T}) = p(\pi_1) \prod_{t=1}^{T-1} p(\pi_{t+1}|\pi_t) \prod_{t=1}^{T} p(x_t|\pi_t)$$

State Initial Distr.    State Transition Distr.    State Observation Distr.
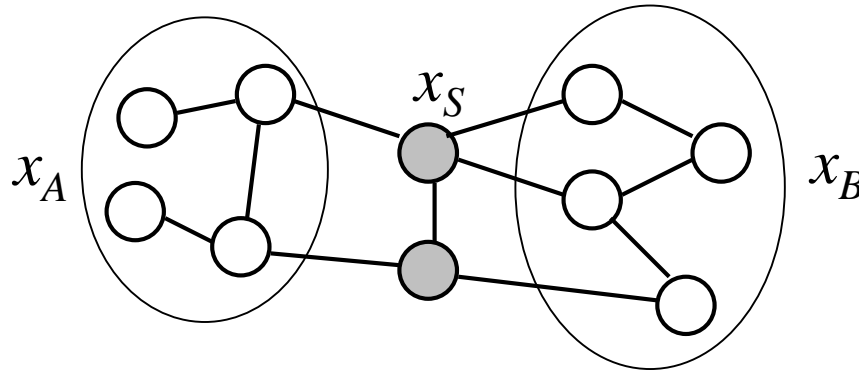
# Undirected Graphical Model (UGM) Semantics - G property

❖ A probability distribution $p(x_V)$ is said to obey the Global Markov property, relative to a undirected graph $g$ , if for any triple $(A, B, S)$ of disjoint subsets of $V$ such that $S$ *separates A from B,*

$$x_A \perp x_B \mid x_S$$



*S separates A from B* : if all trails from $A$ to $B$ intersect $S$

# UGM Semantics - Factorization property

**Hammersley-Clifford Theorem: If $p$ is strictly positive, (F)$\Longleftrightarrow$(G).**

❖ A probability distribution $p(x_V)$ is said to factorize according to $g$, if there exist non-negative functions (called potential functions) $\phi_C(x_C)$ for all cliques $C$ such that

$$p(x_V) = \frac{1}{Z}\prod_{C\in\mathcal{C}}\phi_C(x_C) \quad \text{or} \quad p(x_V) \propto \prod_{C\in\mathcal{C}}\phi_C(x_C)$$

where $Z$ is the normalizing constant (partition function)

$$Z = \sum_{x_V}\prod_{C\in\mathcal{C}}\phi_C(x_C)$$

$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z}\phi(x_1, x_2)\phi(x_2, x_3)\phi(x_3, x_4)\phi(x_1, x_4)$$

A subset of nodes $C$ is called a clique, if every pair of nodes in $C$ is joined.

# UGMs and Energy-Based Models (EBMs)

❖ Let every clique potential be associated with a clique energy $E(x_C)$

$$E_C(x_C) = -log\phi_C(x_C)$$

energy $= -log$ potential

❖ The resulting distribution is known as the Gibbs (or Boltzmann) distribution, originating from statistical physics

$$p(x_V) \propto exp\left[-\sum_{C \in \mathcal{C}} E_C(x_C)\right]$$

High probability states correspond to low energy configurations.

# UGMs and Log-Linear Models

❖ Let each clique potential be a log-linear function

$$log\phi_C(x_C) = \theta_C^T f_C(x_C)$$

where $f_C(x_C)$ is a <u>feature</u> vector derived from (the values of) the variables $x_C$, $\theta_C$ is the associated <u>feature weight</u> vector.

❖ The resulting distribution has the form

$$p(x_V) = \frac{1}{Z(\theta)} exp\left[\sum_C \theta_C^T f_C(x_C)\right]$$

i.e., the distribution with the maximum entropy
s.t. empirical expectation of $f_C$
= model expectation $f_C$

This is known as a Log-Linear Model or a Maximum Entropy Model.

It can be proved that the maxent distribution is the same as
the maximum likelihood distribution from the closure of the set of log-linear distributions.

S. D. Pietra, V. D. Pietra, and J. Lafferty, "Inducing features of random fields", IEEE PAMI, 1997.

# UGM Example - Ising model



- Consider a lattice of binary RV's, $x_i \in \{-1, 1\}$

$$p(x_{1:N^2}) \propto exp\left\{ -\sum_{i \sim j} E(x_i, x_j) \right\} = exp\left\{ \beta \sum_{i \sim j} x_i x_j \right\} \quad \beta > 0$$

- $\beta$ : how much neighboring variables take identical values is favored.

- Samples of Ising models on a lattice with different $\beta$ :



$\beta$ = 0.1 or 10 ?       $\beta$ = 1       $\beta$ = 0.1 or 10 ?

EBMs are natural for modeling interactions (mutual influences), where the directions of edges cannot be clearly defined.

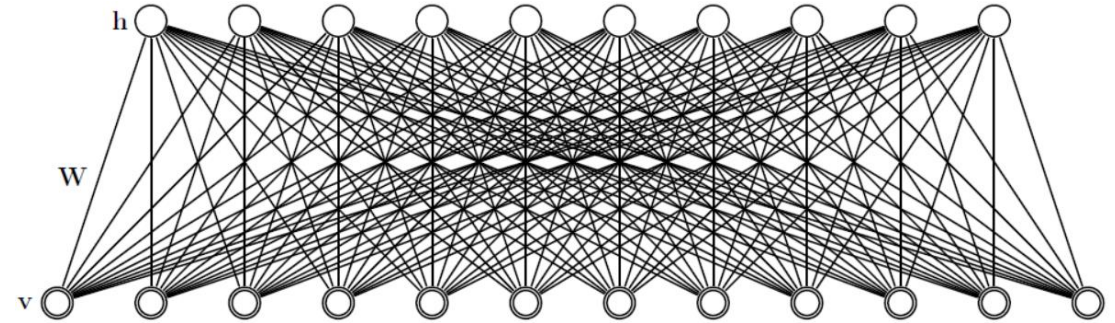# UGM Example - Restricted Boltzmann Machines (RBMs)

- RBM is the main building block of Deep Belief Network, which ignites Deep Learning

- RBM is a UGM over a bipartite graph
  - Binary visible variables $v \in \{0,1\}^D$
  - Binary hidden variables $h \in \{0,1\}^F$
  - $\theta = \{W, b, a\}$



RBM: a stochastic version of a NN

$$p(v, h; \theta) = \frac{1}{Z(\theta)} exp[-E(v, h; \theta)]$$

$$E(v, h; \theta) = -v^T W h - b^T v - a^T h$$

$$= -\sum_{i=1}^{D}\sum_{j=1}^{F} v_i W_{ij} h_j - \sum_{i=1}^{D} b_i v_i - \sum_{j=1}^{F} a_j h_j$$

$$p(h|v; \theta) = \prod_j p(h_j|v), \qquad p(h_j = 1|v) = \sigma\left(\sum_i W_{ij} v_i + a_j\right)$$

$$p(v|h; \theta) = \prod_i p(v_i|h), \qquad p(v_i = 1|h) = \sigma\left(\sum_j W_{ij} h_j + b_i\right)$$

Sigmoid function : $\sigma(x) = \frac{1}{1+e^{-x}}$

Hinton and Salakhutdinov. "Reducing the dimensionality of data with neural networks." Science, 2006.

14

# Learned features $W_{*j}$

Learned receptive fields for unit $h_j$

h

W

v

Observed Data
Subset of 25,000 characters

Learned W: "edges"/"parts"
Subset of 1000 features

$$p(v|h; \theta) = \prod_i p(v_i|h),$$

$$p(v_i = 1|h) = \sigma\left(\sum_j W_{ij}h_j + b_i\right)$$

$W_{*7}$     $W_{*29}$     $W_{*3}$

$$\sim \sigma\left( h_7 \times \qquad + h_{29} \times \qquad + h_3 \times \qquad \cdots\right)$$

$$v|h \sim \sigma(h_1 \cdot W_{*1} + h_2 \cdot W_{*2} + \cdots + b)$$

15

# Neural Random Fields (NRFs) - Basics

$u_\theta(x)$

- NRFs are defined by using NNs to implement $u_\theta(x): \mathbb{R}^d \to \mathbb{R}$

$$p_\theta(x) = \frac{1}{Z(\theta)} exp[u_\theta(x)]$$

Potential function

  - $u_\theta(x)$ can be very flexibly defined

- This type of EBMs has been studied several times in different contexts

  - Deep energy models (DEMs)

$x \in \mathbb{R}^d$

    - Ngiam et al., 2012
    - Kim & Bengio, 2016 - includes linear and squared terms in $u_\theta(x)$
  - Descriptive models / Generative ConvNet
    - Xie et al., 2016 / Dai et al., 2014 - defines in the form of exponential tilting of a reference distribution (Gaussian white noise)
  - Neural random field language models

See references in: Yunfu Song, Zhijian Ou. Learning Neural Random Fields with Inclusive Auxiliary Generators. arXiv:1806.00271, 2018.

    - Wang & Ou, 2017 - defines over sequences

# Probabilistic Graphical Modeling (PGM) Framework - Summary

- **Directed Graphical Models / Bayesian Networks (BNs)**
  - Self-normalized/Local-normalized
  - e.g. Hidden Markov Models (HMMs), Neural network (NN) based classifiers, Variational AutoEncoders (VAEs), Generative Adversarial Networks (GANs), auto-regressive models (e.g. RNNs/LSTMs)

$$p(x_1, x_2, x_3, x_4) = p(x_1)p(x_2|x_1)p(x_3|x_2)p(x_4|x_1, x_3)$$

- **Undirected Graphical Models / Random Fields (RFs) / Energy-based models**
  - Involves the normalizing constant $Z$ / Globally-normalized
  - e.g. Ising model, Conditional Random Fields (CRFs)

$$p(x_1, x_2, x_3, x_4) = \frac{1}{Z}\Phi(x_1, x_2)\Phi(x_2, x_3)\Phi(x_3, x_4)\Phi(x_1, x_4)$$

# Content

# MCMC example: the Metropolis–Hastings algorithm

Problem: We want to draw samples from a target distribution $p(x)$ ?

Solution: Construct a Markov chain that has $p(x)$ as the stationary distribution.

1. Randomly initialize $x_0$

2. For $t = 1, \cdots$

     Generates $x^*$ from a proposal $q(x^*|x_{t-1})$,

     Accept $x_t = x^*$ with probability $min\left\{1, \; \dfrac{p(x^*)q(x_{t-1}|x^*)}{p(x_{t-1})q(x^*|x_{t-1})}\right\}$,

     otherwise set $x_t = x_{t-1}$

  Burn-in : first few samples are discarded.

- For an irreducible & ergodic Markov chain, there exist stationary distribution $\pi$, which satisfies equation $\pi = \pi P$.
- A sufficient condition: satisfy the detailed balance equation

$$\pi_i \, P_{ij} = \pi_j \, P_{ji}$$

# Metropolis–Hastings example

e.g. $q(x^*|x_{t-1}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^*-x_{t-1})^2}{2\sigma^2}}$

# Training of EBMs in general

$$p(x; \theta) = \frac{1}{Z(\theta)} \exp[u(x; \theta)]$$

Normalizing constant:

$$Z(\theta) = \sum_x \exp[u(x; \theta)]$$

- Maximum likelihood (ML) training

The scaled log-likelihood of observations $\{x_i, i = 1, \cdots, N\}$

$$L(\theta) \triangleq \frac{1}{N} \sum_{i=1}^{N} \log p(x_i; \theta) = \left[\frac{1}{N} \sum_{i=1}^{N} u(x_i; \theta)\right] - \log Z(\theta)$$

$$\frac{\partial L(\theta)}{\partial \theta} = E_{\tilde{p}(x)}\left[\frac{\partial u(x; \theta)}{\partial \theta}\right] - E_{p(x; \theta)}\left[\frac{\partial u(x; \theta)}{\partial \theta}\right] = 0 \qquad \text{Maximum Entropy}$$

Expectation under
empirical distribution $\tilde{p}(x) = \frac{1}{N} \sum_{i=1}^{N} 1(x = x_i)$

Expectation under
model distribution $p(x; \theta)$

21

# Learning EBMs by Monte Carlo methods

gradient = empirical expectation – model expectation

1. Approximate the model expectations using Monte Carlo sampling.
   - We can use MCMC to generate the samples, but running MCMC to convergence at the inner loop would be extremely slow.
   - Fortunately, it was shown by Younes (1989) that we can start the MCMC chain at its previous value from the outer loop, and just take a few steps in the inner loop.

> *For mini-batch iterations*
>    Obtain empirical expectations;
>    *For Monte Carlo iterations*
>        Obtain model expectations;
>    End
>    update parameters;
> End

2. We can combine this with stochastic gradient descent (SGD), which takes mini-batches of samples from the empirical distribution.

   Both ideas are applications of the Stochastic Approximation (SA) methodology!

- Robbins and Monro. A stochastic approximation method. The annals of mathematical statistics, 1951.
- L. Younes, "Parametric inference for imperfectly observed Gibbsian fields," Probability Theory and Related Fields, 1989.

# Stochastic Approximation (SA)

Problem: find a solution $\theta$ to $f(\theta) \triangleq E_{z \sim p(\cdot;\theta)}[F(z;\theta)] = 0$, where $\theta \in R^d$, noisy measurement $H(z;\theta) \in R^d$

Method:

(1) <u>Sampling</u>: Generate $z_t \sim K(z_{t-1}, \cdot\;; \theta_{t-1})$, a Markov transition kernel that admits $p(\cdot;\theta_{t-1})$ as the invariant distribution.

(2) <u>Updating</u>: Set $\theta_t = \theta_{t-1} + \gamma_t H(z_t; \theta_{t-1})$

When $f(\theta)$ corresponds to the gradient of some objective function, then under certain regularity conditions, $\theta_t$ will converge to a optimal solution.

$E_{z \sim p(\cdot;\theta)}[H(z;\theta)]$

$\theta$

- Robbins and Monro. A stochastic approximation method. The annals of mathematical statistics, 1951.
- Chen (2002), Stochastic Approximation and Its Applications, Kluwer Academic Publishers.

# Connection between existing RF training methods

- Stochastic Approximation (SA), Robbins and Monro 1951.

- aka Stochastic Maximum Likelihood (SML), Younes 1989.

- This was independently discovered by Tieleman in 2008, who called it persistent contrastive divergence (PCD).

- In regular contrastive divergence (CD), proposed by Hinton 2002, we restart the Markov chain at the training data rather than at the previous state. This will not converge to the MLE.

- "Clearly, the widely used practice of CD1 learning is a rather poor "substitute" for maximum likelihood learning. " (Salakhutdinov phd thesis 2009).

# Content

I.    **Basics for EBMs (45 min)**
   1. Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
   2. Learning EBMs by Monte Carlo methods
   3. Learning EBMs by noise-contrastive estimation (NCE)

II. **EBMs for language modeling (45 min)**
   1. Trans-dimensional random field (TRF) LMs for speech recognition
   2. Residual energy-based models for text generation
   3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for speech recognition and natural language labeling(45 min)**
   1. CRFs as conditional EBMs
   2. CRFs for speech recognition
   3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
   1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
   2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
   3. JRFs for semi-supervised natural language labeling

# Learning EBMs by noise-contrastive estimation (NCE)

- The target RF model $\qquad p_\theta(x) = \dfrac{1}{Z(\theta)}\, e^{u_\theta(x)}$

- Treat $\log Z(\theta)$ as a parameter $\zeta$ and rewrite $\qquad p_{\theta,\zeta}(x) \propto e^{u_\theta(x) - \zeta}$

- Introduce a **noise distribution** $q(x)$, and consider a binary classification

$$x \sim p_0(x)$$
$$\boxed{\text{Binary discriminator}}$$
$$x \sim q(x)$$
$$C = 0/1$$

$$P(C = 0|x) = \frac{p_{\theta,\zeta}(x)}{p_{\theta,\zeta}(x) + \nu q(x)}, where\ \nu = \frac{P(C = 1)}{P(C = 0)}$$

$$P(C = 1|x) = 1 - P(C = 0|x)$$

- Noise Contrastive Estimation (NCE):

$$\max_{\theta,\zeta} E_{x \sim p_0(x)}[\log P(C = 0|x)] + \nu E_{x \sim q(x)}[\log P(C = 1|x)]$$

Consistency: $p_\theta \to p_0$ (oracle), under infinite amount of data and infinite capacity of $p_\theta$.

Michael Gutmann, Aapo Hyvarinen. "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," AIStat, 2010.

# NCE discussion

- In applying NCE, a natural question to ask is: from a statistical standpoint, what the best choice of $q$ and $\nu$ would be?

To get estimates with a small estimation error, the foregoing discussion suggests the following

1. Choose noise for which an analytical expression for $\ln p_n$ is available.

2. Choose noise that can be sampled easily.

3. Choose noise that is in some aspect, for example with respect to its covariance structure, similar to the data.

4. Make the noise sample size as large as computationally possible.

Michael Gutmann, Aapo Hyvarinen. "Noise-contrastive estimation: A new estimation principle for unnormalized statistical models," AIStat, 2010.

# Dynamic noise-contrastive estimation (DNCE)

- Reliable NCE needs a large $\nu \approx 20$, which almost linearly increases the training cost.

- The model estimated by NCE could be overfitted to the empirical distribution.

$$x \sim \alpha p_0(x) + (1-\alpha)q_\phi(x)$$

$$x \sim q_\phi(x)$$

Binary discriminator

$$C = 0/1$$

Dynamic noises help optimization, by gradually increasing the difficulty of the two-class discrimination task

Introduce a dynamic noise distribution, simultaneously optimized to be close to the data distribution, so that small $\nu$ could be used.

$$\left\{ \begin{array}{l} \min_{\phi} KL(p_0 || q_\phi) \\ \max_{\theta,\zeta} E_{x \sim \alpha p_0(x) + (1-\alpha)q_\phi(x)}[\log P(C = 0|x)] + \nu E_{x \sim q_\phi(x)}[\log P(C = 1|x)] \end{array} \right.$$

Using interpolation to alleviate overfitting

Consistency under infinite amount of data and infinite capacity of $p_\theta$ and $q_\phi$

Bin Wang, Zhijian Ou. Improved training of neural trans-dimensional random field language models with dynamic noise-contrastive estimation. SLT, 2018.

# Content

I. **Basics for EBMs (45 min)**
1. Probabilistic graphical modeling (PGM) framework and EBM model examples
2. Learning EBMs by Monte Carlo methods
3. Learning EBMs by noise-contrastive estimation (NCE)

II. **EBMs for language modeling (45 min)**
→ 1. Trans-dimensional random field (TRF) LMs for speech recognition
2. Residual energy-based models for text generation
3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for natural language labeling and speech recognition (45 min)**
1. CRFs as conditional EBMs
2. CRFs for speech recognition
3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
3. JRFs for semi-supervised natural language labeling

$p_\theta(x)$

I. Basics for EBMs

$p_\theta(h|x)$     $p_\theta(x, h)$

29

# Section Content

1. Motivation

2. Related work

3. Method: **RF LMs**

4. Experiments

5. Conclusion

- Bin Wang, Zhijian Ou, Zhiqiang Tan. Trans-dimensional Random Fields for Language Modeling. ACL Long Paper, 2015.
- Bin Wang, Zhijian Ou, Zhiqiang Tan. Learning Trans-dimensional Random Fields with Applications to Language Modeling. TPAMI, 2018.
- Bin Wang, Zhijian Ou. Language modeling with neural trans-dimensional random fields. ASRU, 2017.
- Bin Wang, Zhijian Ou. Learning neural trans-dimensional random field language models with noise-contrastive estimation. ICASSP, 2018.
- Bin Wang, Zhijian Ou. Improved training of neural trans-dimensional random field language models with dynamic noise-contrastive estimation. SLT, 2018.
- Silin Gao, Zhijian Ou, Wei Yang, Huifang Xu. Integrating discrete and neural features via mixed-feature trans-dimensional random field language models. ICASSP, 2020. [Oral]

# N-gram LMs

- Language modeling (LM) is to determine the joint probability of a sentence, i.e. a word sequence.

- Dominant: Directed modeling approach

Current word        All previous words/history

$$p(x_1, x_2, \cdots, x_l) = \prod_{i=1}^{l} p(x_i | x_1, \cdots, x_{i-1})$$

Previous $n-1$ words

$$\approx \prod_{i=1}^{l} p(x_i | x_{i-n+1}, \cdots, x_{i-1})$$

- Using Markov assumption leads to the N-gram LMs
  – One of the state-of-the-art LMs

Chen and Goodman. "An empirical study of smoothing techniques for language modeling", Computer Speech & Language, 1999.

# Recurrent Neural Nets (RNNs)/LSTM/Transformer LMs



$$p(x_i|x_1,\cdots,x_{i-1}) \approx p\big(x_i|h_{i-1}(x_1,\cdots,x_{i-1})\big) \approx \frac{h_{i-1}^T w_k}{\sum_{k=1}^{V} h_{i-1}^T w_k}$$

☹.1 Computational expensive in both training and testing [1]

    e.g. $V = 10^4 \sim 10^6, w_k \in \mathbb{R}^{250 \sim 1024}$

☹.2 As directed sequential model /Auto-regressive model, potentially suffers from Exposure Bias and Label Bias

[1] Partly alleviated by using un-normalized models (e.g., through NCE) or a small set of tokens (e.g., BPE).

# Trans-dimensional Random Field (TRF) LM: motivation

☺.1 Avoid local normalization

$$p_\theta(x^l) = \frac{1}{Z_l(\theta)} e^{u_\theta(x^l)}, x^l \triangleq x_1, x_2, \cdots, x_l$$

☺.2 Flexible

| Type | Features |
|------|----------|
| w | $(w_{-3}w_{-2}w_{-1}w_0)(w_{-2}w_{-1}w_0)(w_{-1}w_0)(w_0)$ |
| c | $(c_{-3}c_{-2}c_{-1}c_0)(c_{-2}c_{-1}c_0)(c_{-1}c_0)(c_0)$ |
| ws | $(w_{-3}w_0)(w_{-3}w_{-2}w_0)(w_{-3}w_{-1}w_0)(w_{-2}w_0)$ |
| cs | $(c_{-3}c_0)(c_{-3}c_{-2}c_0)(c_{-3}c_{-1}c_0)(c_{-2}c_0)$ |
| wsh | $(w_{-4}w_0)\ (w_{-5}w_0)$ |
| csh | $(c_{-4}c_0)\ (c_{-5}c_0)$ |
| cpw | $(c_{-3}c_{-2}c_{-1}w_0)\ (c_{-2}c_{-1}w_0)(c_{-1}w_0)$ |
| tied | $(c_{-9:-6}, c_0)\ (w_{-9:-6}, w_0)$ |

Discrete features



CNN features



BLSTM features

# Trans-dimensional random fields (TRFs): model

- Assume the sentences of length $l$ are distributed as follows:

$$p_l(x^l; \lambda) = \frac{1}{Z_l(\lambda)} e^{\lambda^T f(x^l)} \qquad l = 1, \dots, l_{max}$$

$x^l \triangleq x_1, x_2, \cdots, x_l$ is a word sequence with length $l$;

$f(x^l) = \left( f_1(x^l), \dots, f_d(x^l) \right)^T$ is the feature vector;

$\lambda = (\lambda_1, \dots, \lambda_d)^T$ is the parameter vector;

$Z_l(\lambda) = \sum_{x^l} e^{\lambda^T f(x^l)}$ is the normalization constant.

Needed to be estimated

- Assume length $l$ is associated with priori probability $\pi_l$. Therefore the pair $(l, x^l)$ is jointly distributed as:

$$p(l, x^l; \lambda) = \pi_l \cdot p_l(x^l; \lambda)$$

34

# Feature definition

$$p_l(x^l; \lambda) = \frac{1}{Z_l(\lambda)} e^{\lambda^T f(x^l)}$$

- $f_i(x^l)$ returns the count of a specific phrase observed in the input sentence $x^l$

    $\underline{x^l = \textit{he is a teacher and he is also a good father.}}$

    $f_{he\ is}(x^l) =$ count of "*he is*" observed in $x^l = 2$

    $f_{a\ teacher}(x^l) =$ count of "*a teacher*" observed in $x^l = 1$

    $f_{she\ is}(x^l) =$ count of "*she is*" observed in $x^l = 0$

    ... ...

- For example, n-grams and skip n-grams (tied or not) of orders ranging from 1 to 10, observed in the training set are added to the features.

# Review the development of TRF LMs

| | |
|---|---|
| ACL-2015<br>TPAMI-2018 | • Discrete features<br>• Augmented stochastic approximation (**AugSA**) for model training |
| ASRU-2017 | • Potential function as a deep CNN.<br>• Model training by **AugSA plus JSA** (joint stochastic approximation) |
| ICASSP-2018 | • Potential function in the form of **exponential tilting** (revisited in residual EBMs)<br>• Use LSTM on top of CNN<br>• **NCE** is introduced to train TRF LMs |
| SLT-2018 | • Simplify the potential definition by using only Bidirectional LSTM<br>• Propose **Dynamic NCE** for improved model training |
| ICASSP-2020 | • **Mixed-feature** TRFs, by integrating discrete and neural features |

# WSME - Introduction

- Whole-sentence maximum entropy (WSME)

  - Rosenfeld, Chen, Zhu. "Whole-sentence exponential language models: a vehicle for linguistic-statistical integration". Computer Speech & Language, 2001.

$$p(x; \lambda) = \frac{1}{Z(\lambda)} \exp[\lambda^T f(x)]$$

- The empirical results of previous WSME models are not satisfactory

  - After incorporating lexical and syntactic information, 1% and 0.4% respectively in perplexity and in WER is reported for the resulting WSME (Rosenfeld et al., 2001).

  - Amaya and Benedi. "Improvement of a whole sentence maximum entropy language model using grammatical features", ACL 2001.

  - Ruokolainen, Alumae, Dobrinkat. "Using dependency grammar features in whole sentence maximum entropy language model for speech recognition". HLT 2010.

# RFLMs vs WSME

- Whole-sentence maximum entropy (WSME)

$$p(l, x^l; \lambda) = \frac{1}{Z(\lambda)} \exp[\lambda^T f(x^l)], \qquad x \triangleq (l, x^l), \qquad x^l \triangleq (x_1, x_2, \cdots, x_l)$$

Essentially a mixture distribution with unknown weights (differ from each other greatly, $10^{40}$) !
Poor sampling → poor estimate of gradient → poor fitting

$$p(l, x^l; \lambda) = \frac{Z_l(\lambda)}{Z(\lambda)} \cdot \frac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)], \, Z_l(\lambda) = \sum_{x^l} \exp[\lambda^T f(x^l)]$$

# RFLMs vs WSME

- Whole-sentence maximum entropy (WSME)

$$p(l, x^l; \lambda) = \frac{1}{Z(\lambda)} \exp[\lambda^T f(x^l)], \qquad x \triangleq (l, x^l), \qquad x^l \triangleq (x_1, x_2, \cdots, x_l)$$

Essentially a mixture distribution with unknown weights (differ from each other greatly, $10^{40}$) !
Poor sampling → poor estimate of gradient → poor fitting

$$p(l, x^l; \lambda) = \frac{Z_l(\lambda)}{Z(\lambda)} \cdot \frac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)], \; Z_l(\lambda) = \sum_{x^l} \exp[\lambda^T f(x^l)]$$

- We propose a trans-dimensional RF model

$$p(l, x^l; \lambda) = \pi_l \cdot \frac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)], \qquad l = 1, \cdots, m$$

Empirical length probabilities in the training data
Serve as a control device to improve sampling from multiple distributions !

# Apply SA to RFLM training

- The trans-dimensional RF model $\qquad p(l, x^l; \lambda) = \pi_l \cdot \dfrac{1}{Z_l(\lambda)} \cdot \exp[\lambda^T f(x^l)]$ (1)

$$E_{\tilde{p}(x)}[f_i(x)] - E_{p(x;\lambda)}[f_i(x)] = 0, \qquad x \triangleq (l, x^l)$$

- Consider the joint distribution of the pair $(l, x^l)$ $\qquad p(l, x^l; \lambda, \zeta) \propto \pi_l \cdot \dfrac{1}{e^{\zeta_l}} \cdot \exp[\lambda^T f(x^l)]$ (2)

where $\zeta_l$ is hypothesized values of the true $\zeta_l^*(\lambda) = log Z_l(\lambda)$

The marginal probability of length $l$ is: $p(l; \lambda, \zeta) = \dfrac{\pi_l e^{-\zeta_l + \zeta_l^*(\lambda)}}{\sum_j \pi_l e^{-\zeta_j + \zeta_j^*(\lambda)}}$

- SA is used to find $\zeta_l^* = \zeta_l^*(\lambda^*)$ and $\lambda^*$ that solves

$$\begin{cases} \pi_l = p(l; \lambda, \zeta), & l = 1, \cdots, m \\ 0 = E_{\tilde{p}(x)}[f_i(x)] - E_{p(l, x^l; \lambda, \zeta)}[f_i(x)] \end{cases}$$

Zhiqiang Tan. 2015. Optimally adjusted mixture sampling and locally weighted histogram. In Technical Report, Dept. of Statistics, Rutgers Univ.

# RFLMs – Breakthrough in training (1)

- Propose Augmented Stochastic Approximation (AugSA) Training Algorithm
  - Simultaneously updates the model parameters and normalizing constants

**Input:** training set
1: set initial values $\lambda^{(0)} = (0, \ldots, 0)^T$ and
   $$\zeta^{(0)} = \zeta^*(\lambda^{(0)}) - \zeta_1^*(\lambda^{(0)})$$
2: **for** $t = 1, 2, \ldots, t_{max}$ **do**
3:     set $B^{(t)} = \emptyset$
4:     set $(j^{(t,0)}, x^{(t,0)}) = (j^{(t-1,K)}, x^{(t-1,K)})$
    ***Step I: MCMC sampling***
5:     **for** $k = 1 \rightarrow K$ **do**
6:       $(j^{(t,k)}, x^{(t,k)}) = \text{sample}((j^{(t,k-1)}, x^{(t,k-1)}))$ (See Tab.1)
7:       set $B^{(t)} = B^{(t)} \cup \{(j^{(t,k)}, x^{(t,k)})\}$
8:     **end for**
    ***Step II: SA updating***
9:     Compute $\lambda^{(t)}$ based on (16)
10:     Compute $\zeta^{(t)}$ based on (17) and (18)
11: **end for**

Fig. 1. Augmented stochastic approximation (AugSA)

Bin Wang, Zhijian Ou, Zhiqiang Tan. Learning Trans-dimensional Random Fields with Applications to Language Modeling. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 2018.

# RFLMs – Breakthrough in training (2)

- Propose Trans-dimensional mixture sampling
  - Sampling from $p(l, x^l; \lambda, \zeta)$, a mixture of RFs on subspaces of different dimensions.
  - Formally like RJ-MCMC (Green, 1995).

```
 1: function SAMPLING((L^(t-1), X^(t-1)))
 2:      set k = L^(t-1)
 3:      set L^(t) = k
 4:      set X^(t) = X^(t-1)
              Stage I: Local jump
 5:      generate j ~ Γ(k, ·)
 6:      if j = k + 1 then
 7:
 8:              generate Y ~ g_{k+1}(y|X^(t-1)) (equ.24)
 9:              set L^(t) = j and X^(t) = {X^(t-1), Y} with
         probability equ.22
10:      end if
11:      if j = k - 1 then
12:              set L^(t) = j and X^(t) = X^(t-1)_{1:k-1} with prob-
         ability equ.23
13:      end if
              Stage II: Markov move
14:      for i = 1 → L^(t)  do
15:
16:
17:              a ~ p(L^(t), {X^(t)_{1:i-1}, ·, X^(t)_{i+1:L^(t)}}; Λ, ζ)
18:          X_i^(t) ← a
19:      end for
20:      return (L^(t), X^(t))
21: end function
```

42

# Motivation: Integrating discrete and neural features

- Language models using discrete features (N-gram LMs, Discrete TRF LMs)
  - Mainly capture local lower-order interactions between words
  - Better suited to handling symbolic knowledges

- Language models using neural features (LSTM LMs, Neural TRF LMs)
  - Able to learn higher-order interactions between words
  - Good at learning smoothed regularities due to word embeddings

- Interpolation of LMs[1, 2]: usually achieves further improvement
  - Discrete and neural features have complementary strength. ☺
  - Two-step model training is sub-optimal. ☹

[1]Xie Chen, Xunying Liu, Yu Wang, Anton Ragni, Jeremy HM Wong, and Mark JF Gales, "Exploiting future word contexts in neural network language models for speech recognition," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 9, pp. 1444–1454, 2019.

[2]Bin Wang, Zhijian Ou, Yong He, and Akinori Kawamura, "Model interpolation with trans-dimensional random field language models for speech recognition," *arXiv preprint arXiv:1603.09170*, 2016.

# Mixed-feature TRF LMs: Definition

- Mixed-feature TRF LMs:

    ▪ $p(l, x^l; \eta) = \dfrac{\pi_l}{Z_l(\eta)} e^{V(x^l, \eta)}, \ V(x^l, \eta) = \lambda^T f(x^l) + \phi(x^l; \theta), \ \eta = (\lambda, \theta)$

Discrete n-gram features, with parameter $\lambda$:

$$f(x^l) = (f_1(x^l), f_2(x^l), \cdots, f_d(x^l))$$

$d$: the total number of types of n-grams

$$f_k(x^l) = c$$

where $c$ is the count of the $k$th n-gram type in $x^l$

$\underline{x^l = \textit{he is a teacher and he is also a good father.}}$

$f_{he\ is}(x^l) = $ count of "$he\ is$" in $x^l = 2$

$f_{a\ teacher}(x^l) = $ count of "$a\ teacher$" in $x^l = 1$

Neural network features, with parameter $\theta$



Forward LSTM

Word embedding

Backward LSTM

$$\phi(x^l; \theta) = \sum_{i=1}^{l-1} h_{f,i}{}^T e_{i+1} + \sum_{i=2}^{l} h_{b,i}{}^T e_{i-1}$$

44

# Mixed-feature TRF LMs: Training

- Treat $\log Z_l(\eta)$ as a parameter $\zeta_l$ and rewrite

$$p(l, x^l; \eta) = \frac{\pi_l}{Z_l(\eta)} e^{V(x^l, \eta)} \longrightarrow p(x; \xi) = \pi_l e^{V(x^l, \eta) - \zeta_l}, x = (l, x^l), \xi = (\eta, \zeta)$$

- Introduce a **noise distribution** $q_\phi(x)$, and consider a binary classification

$$q_\phi(x) = \pi_l \times NNLM(x^l), NNLM \text{ could be LSTM/Transformer autoregressive LM}$$

$x \sim p_0(x)$ → Binary discriminator → $C = 0/1$

$x \sim q_\phi(x)$

$$P(C = 0|x) = \frac{p(x; \xi)}{p(x; \xi) + \nu q_\phi(x)}, where \ \nu = \frac{P(C = 1)}{P(C = 0)}$$

$$P(C = 1|x) = 1 - P(C = 0|x)$$

- Noise Contrastive Estimation (NCE):

$$\max_{\xi} E_{x \sim p_0(x)}[\log P(C = 0|x)] + \nu E_{x \sim q_\phi(x)}[\log P(C = 1|x)]$$

Reliable NCE needs a large $\nu \approx 20$; Dynamic-NCE works well with $\nu = 1$

# Experiments: PTB dataset

WER curves of the three TRF LMs during the first 100 training epochs:



- **Mixed TRF converges faster than the state-of-the-art Neural TRF, using only 58% training epochs.**

☺ **The discrete features in Mixed TRF lower the non-convexity of the optimal problem, and reduce the amount of patterns for neural features to capture.**

46

# On Google one-billion word benchmark

**Training**: Google One-Billion word benchmark, 0.8 billion words,  568K vocabulary
**Testing**: WSJ'92 test data, 330 utterances, rescoring 1000-best lists

| Model | WER (%) | #Param (M) | Training time | Inference Time |
|---|---|---|---|---|
| KN5 | 6.13 | 133 | 2.5  h (1 CPU) | 0.491 s (1 CPU) |
| LSTM-2x1024 | 5.55 | 191 | 144 h (2 GPUs) | 0.909 s (2 GPUs) |
| discrete-TRF basic | 6.04 | 102 | 131 h (8 cores and 2 GPUs) | 0.022 s (1 CPU) |
| neural-TRF | 5.47 | 114 | 336 h (2 GPUs) | 0.017 s (2 GPUs) |
| mix-TRF | 5.28 | 216 | 297 h (8 cores and 2 GPUs) | 0.024 s (1 core and 2 GPUs) |
| LSTM-2x1024+KN5 | 5.38 | 324 | | |

**5%**  **33%**  **38x**

## Open-source LM toolkit

https://github.com/thu-spmi/SPMILM

# Section Conclusion

- Language models play an important role for ASR and more applications!

- Random Field language models

  - Avoid local normalization
  - Being flexible to integrate rich features (both discrete and neural)
  - Overcome "label bias" and "Exposure bias"

- More related work

  - Residual energy-based models for text generation
  - Electric: an energy-based cloze model for representation learning over text

- Yuntian Deng, Anton Bakhtin, Myle Ott, Arthur Szlam, and Marc'Aurelio Ranzato. Residual energy-based models for text generation, ICLR 2020.
- Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Pre-training transformers as energy-based cloze models, EMNLP 2020.

# Content

I.  Basics for EBMs (45 min)
    1.  Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
    2.  Learning EBMs by Monte Carlo methods
    3.  Learning EBMs by noise-contrastive estimation (NCE)

II. EBMs for language modeling (45 min)
    1.  Trans-dimensional random field (TRF) LMs for speech recognition
    2.  Residual energy-based models for text generation
    3.  Electric: an energy-based cloze model for representation learning over text

III. EBMs for speech recognition and natural language labeling (45 min)
    1.  CRFs as conditional EBMs
    2.  CRFs for speech recognition
    3.  CRFs for sequence labeling in NLP

IV. EBMs for semi-supervised natural language labeling (45 min)
    1.  Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
    2.  Upgrading CRFs to Joint random fields (JRFs) for sequential data
    3.  JRFs for semi-supervised natural language labeling

# Motivation

- Text generation is ubiquitous in many NLP tasks, from summarization, to dialogue and machine translation.

- Locally normalized LMs are plagued by exposure bias and label bias.

- EBMs are ideal for modeling text as ... , but seldom explored
  - they can score the whole input at once,
  - they are not prone to exposure bias and label bias ,
  - they may enable generation of large chunks of text, which should help improve coherency (e.g., resampling from the large set of candidates produced by the base locally normalized LM).

Yuntian Deng, Anton Bakhtin, Myle Ott, Arthur Szlam, and Marc'Aurelio Ranzato. Residual energy-based models for text generation, ICLR 2020.

# Model - Residual EBMs

- Deng et al. investigate an EBM trained on the residual of a pretrained autoregressive LM [a, b]

$$P_\theta(x) \propto P_{LM}(x) exp[-E_\theta(x)]$$

- $P_{LM}(x)$: a pretrained autoregressive LM and fixed
- $E_\theta(x)$: the residual energy function parameterized by $\theta$
- Call $P_\theta$ the joint model

- Consider the problem of conditional generation of discrete sequences

$$P_\theta(x_{p+1}, \cdots, x_T | x_1, \cdots, x_p) = \frac{P_{LM}(x_{p+1}, \cdots, x_T | x_1, \cdots, x_p) \exp(-E_\theta(x_1, \cdots, x_T))}{Z_\theta(x_1, \cdots, x_p)} \quad (2)$$

- Given a prefix $x_1, \cdots, x_p$, generate a sequence of total length $T > p$

a. Bin Wang, Zhijian Ou. Learning neural trans-dimensional random field language models with noise-contrastive estimation. ICASSP, 2018.
b. Tetiana Parshakova, Jean-Marc Andreoli, and Marc Dymetman. Global autoregressive models for data-efficient sequence learning. CoNLL, 2019.

# Training of Residual EBMs

- Trained using NCE, and more specifically its conditional version
  - Using $P_{LM}(x)$ as the noise distribution
  - $E_\theta(x)$ initialized with BERT/RoBERTa; in the final layer we project the mean-pooled hidden states to a scalar energy value.
  - $x_+$ : positive sentence taken from the training data
  - $x_-$ : negative sentence drawn from $P_{LM}(x)$ (for a given ground truth prefix)

$$\max \mathbb{E}_{x_+ \sim P_{data}} \log \frac{1}{1 + \exp(E_\theta(x_+))} + \mathbb{E}_{x_- \sim P_{LM}} \log \frac{1}{1 + \exp(-E_\theta(x_-))}$$

In all our experiments we use a prefix of size 120 tokens and we generate the following 40 tokens; with the notation of Eq. 2, $p = 120$ and $T = 160$. For training the joint models, for efficiency we generated 16/128 samples per prefix for CC-News/Book Corpus offline, and sample uniformly from those samples at training time.

# Generation by Residual EBMs

- Use self-normalizing importance sampling

    a) Sampling from the proposal - the auto-regressive language model $P_{LM}(x)$

    b) Resampling according to the energy function.

---

**Algorithm 1:** Top-k Joint Sampling

**Input:** number of samples $n$ drawn from $P_{LM}$, value of $k$ in top-k

```
// Get a set of samples from P_LM
```
sample $n$ samples $\{x^1, \cdots, x^n\}$ from $P_{LM}$ with top-k sampling
calculate energies $s^i = E_\theta(x^i)$ for each $x^i \in \{x^1, \cdots, x^n\}$
```
// Resample from the set of LM samples
```
sample $x = x^i$ with probability $\dfrac{\exp(-s^i)}{\sum_{j=1}^{n} \exp(-s^j)}$

**return** $x$

---

# Evaluation

- PPL (perplexity) evaluation needs to estimate the partition function

$$Z_\theta = \sum_x P_{LM}(x) exp[-E_\theta(x)] = \mathbb{E}_{x \sim P_{LM}(x)}\{exp[-E_\theta(x)]\}$$

- Two estimators for the lower and upper bounds of the partition function

**Theorem 2.** *Denote $T_n$ as the empirical estimate of $\log \mathbb{E}_{x \sim P_{LM}} \exp(-E(x))$ with $n$ samples $x_i \sim P_{LM} (i = 1, \cdots, n)$: $T_n = \log \frac{1}{n} \sum_{i=1}^n \exp(-E(x_i))$, then $\forall \epsilon > 0, \exists N > 0$ such that $\forall n > N$ we have*

$$Z_\theta - \epsilon < \mathbb{E}[T_n] < Z_\theta < \mathbb{E}[(2n-1)T_n - 2(n-1)T_{n-1}] < Z_\theta + \epsilon \qquad (4)$$

- Per-step probabilities

$$P(x_t|x_{<t}) = P_{LM}(x_t|x_{<t}) \frac{\mathbb{E}_{x'_{t+1}, \cdots, x'_T \sim P_{LM}(\cdot|x_{\leq t})}[\exp(-E_\theta(x_{\leq t}, x'_{t+1}, \cdots, x'_T))]}{\mathbb{E}_{x'_t, \cdots, x'_T \sim P_{LM}(\cdot|x_{\leq t-1})}[\exp(-E_\theta(x_{\leq t-1}, x'_t, \cdots, x'_T))]}. \qquad (5)$$

  - The basic PLM distribution is adjusted

A central technical contribution of this paper!

# Experiments

- On two large datasets, <u>residual EBMs</u> have demonstrated improved generation ability against very strong auto-regressive baselines, both in terms of estimated perplexity and through human evaluation.



Figure 1: Perplexity gain of **JOINT BIT-MED** and JOINT BIT-LARGE∗ (using BASE LM-24L) at each position relative to **BASE LM-24L** on the test set of CC-News.

# Content

I. **Basics for EBMs (45 min)**
1. Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
2. Learning EBMs by Monte Carlo methods
3. Learning EBMs by noise-contrastive estimation (NCE)

II. **EBMs for language modeling (45 min)**
1. Trans-dimensional random field (TRF) LMs for speech recognition
2. Residual energy-based models for text generation
3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for speech recognition and natural language labeling (45 min)**
1. CRFs as conditional EBMs
2. CRFs for speech recognition
3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
3. JRFs for semi-supervised natural language labeling

# Motivation

- The cloze task of predicting the identity of a token given its surrounding context has proven highly effective for representation learning over text.

- BERT implements the cloze task by replacing input tokens with [MASK], but
  - Drawback in efficiency (only 15% of tokens are masked out at a time)
  - A pre-train/fine-tune mismatch where BERT sees [MASK] tokens in training but not in fine-tuning
  - Less concerned with producing (pseudo-)likelihood score for text

$$
\text{the} \rightarrow \boxed{\text{BERT}} \rightarrow \begin{cases} p_\theta(\text{dog}\,|\,\text{the }\underline{\quad}\text{ barked}) \\ p_\theta(\text{cat}\,|\,\text{the }\underline{\quad}\text{ barked}) \\ p_\theta(\text{the}\,|\,\text{the }\underline{\quad}\text{ barked}) \\ \quad\cdots \end{cases}
$$

the → [MASK] → barked

Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Pre-training transformers as energy-based cloze models, EMNLP 2020.

# The Electric Model

- Learn $p_{data}(x_t|x_{\backslash t})$ of a token $x_t$ occurring in the surrounding context

$$x_{\backslash t} = [x_1, \cdots, x_{t-1}, x_{t+1}, \cdots, x_n]$$

  - Maps the unmasked input $x = [x_1, \cdots, x_n]$ into contextualized vector representations $h(x) = [h_1, \cdots, h_n]$ using a Transformer net.

  - Define a conditional EBM, with a learned weight vector $w$

$$P_\theta(x_t|x_{\backslash t}) \propto exp[-w^T h_t], \qquad \text{where energy function } E(x)_t \triangleq w^T h_t$$

**BERT**
- Mask 15% of the input sequence
- Calculate the distribution over the vocabulary using a softmax layer for each masked position

**Electric**
- Using unmasked input
- Likelihood scores $P_\theta(x_t|x_{\backslash t})$ can be computed simultaneously for all positions rather than only for a small masked-out subset.

# Training of the Electric Model

- Trained using NCE, and more specifically its conditional version

  - Define the un-normalized output

  $$\hat{p}_\theta(x_t|x_{\backslash t}) = exp[-w^T h_t]$$

  - A binary classifier is trained to distinguish positive $x_t$ vs negative $\hat{x}_t$, with $k$ negatives sampled for every $n$ positive data points.

Formally, the NCE loss $\mathcal{L}(\theta)$ is

$$n \cdot \underset{x,t}{\mathbb{E}} \left[ -\log \frac{n \cdot \hat{p}_\theta(x_t|\boldsymbol{x}_{\backslash t})}{n \cdot \hat{p}_\theta(x_t|\boldsymbol{x}_{\backslash t}) + k \cdot q(x_t|\boldsymbol{x}_{\backslash t})} \right] + k \cdot \underset{\substack{x,t \\ \hat{x}_t \sim q}}{\mathbb{E}} \left[ -\log \frac{k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})}{n \cdot \hat{p}_\theta(\hat{x}_t|\boldsymbol{x}_{\backslash t}) + k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})} \right]$$

  - Noise distribution: a two-tower cloze model

  $$\overrightarrow{\boldsymbol{h}} = T_{\text{LTR}}(\boldsymbol{x}), \quad \overleftarrow{\boldsymbol{h}} = T_{\text{RTL}}(\boldsymbol{x})$$

  $$q(x_t|\boldsymbol{x}_{\backslash t}) = \text{softmax}(\boldsymbol{W}[\overrightarrow{\boldsymbol{h}}_{t-1}, \overleftarrow{\boldsymbol{h}}_{t+1}])_{x_t}$$

  Two causal transformers
  (left-to-right & right-to-left)

# Training of the Electric Model

$$n \cdot \underset{\boldsymbol{x},t}{\mathbb{E}} \left[ -\log \frac{n \cdot \hat{p}_\theta(x_t|\boldsymbol{x}_{\backslash t})}{n \cdot \hat{p}_\theta(x_t|\boldsymbol{x}_{\backslash t}) + k \cdot q(x_t|\boldsymbol{x}_{\backslash t})} \right] + k \cdot \underset{\substack{\boldsymbol{x},t \\ \hat{x}_t \sim q}}{\mathbb{E}} \left[ -\log \frac{k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})}{n \cdot \hat{p}_\theta(\hat{x}_t|\boldsymbol{x}_{\backslash t}) + k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})} \right]$$

- Naïve calculation is expensive: $k + 1$ forward passes through the transformer to compute the $\hat{p}_\theta$s (once for the positive samples $x_t|x_{\backslash t}$ and once for each negative sample $\hat{x}_t|x_{\backslash t}$)

**Algorithm 2** Efficient NCE loss estimation

**Given:** Input sequence $\boldsymbol{x}$, number of negative samples $k$, noise distribution $q$, model $\hat{p}_\theta$.

**1.** Pick $k$ unique random positions $R = \{r_1, ..., r_k\}$ where each $r_i$ is $1 \leq r_i \leq n$.

**2.** Replace the $k$ random positions with negative samples: $\hat{x}_i \sim q(x_i|\boldsymbol{x}_{\backslash i})$ for $i \in R$,
$x^{\text{noised}} = \text{REPLACE}(\hat{\boldsymbol{x}}, R, \hat{X})$.

**3.** For each position $t = 1$ to $n$: add to the loss

$-\log \frac{k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})}{(n-k) \cdot \hat{p}_\theta(\hat{x}_t|\boldsymbol{x}^{\text{noised}}_{\backslash t}) + k \cdot q(\hat{x}_t|\boldsymbol{x}_{\backslash t})}$    if $t \in R$

$-\log \frac{(n-k) \cdot \hat{p}_\theta(x_t|\boldsymbol{x}^{\text{noised}}_{\backslash t})}{(n-k) \cdot \hat{p}_\theta(x_t|\boldsymbol{x}^{\text{noised}}_{\backslash t}) + k \cdot q(x_t|\boldsymbol{x}_{\backslash t})}$    otherwise

**Trick:**

- Simultaneously choose $k = \lceil 0.15n \rceil$ random positions

- One pass through the transformer over $x^{\text{noised}}$

- Assume $\hat{p}_\theta(\cdot|x_{\backslash t}) \approx \hat{p}_\theta(\cdot|x^{\text{noised}}_{\backslash t})$

# Experiments

| Model | MultiNLI | SQuAD 2.0 | GLUE Avg. |
|---|---|---|---|
| BERT | 84.3 | 73.7 | 82.2 |
| XLNet | 85.8 | 78.5 | – |
| ELECTRA | 86.2 | 80.5 | 85.1 |
| Electric | 85.7 | 80.1 | 84.1 |

Table 1: Dev-set scores of <mark>pre-trained models on down-stream tasks</mark>. To provide direct comparisons, we only show base-sized models pre-trained on WikiBooks.

| Model | Pre-train Data | Clean WER | Other WER | Runtime |
|---|---|---|---|---|
| None | – | 7.26 | 20.37 | 0 |
| BERT | WikiBooks | 5.41 | 17.41 | $n$ |
| Electric | WikiBooks | 5.65 | 17.42 | 1 |
| GPT-2 | OWT | 5.64 | 17.60 | 1 |
| TwoTower | OWT* | 5.32 | 17.25 | 1 |
| ELECTRA-TT | OWT* | 5.22 | 17.01 | 1 |
| Electric | OWT* | 5.18 | 16.93 | 1 |

Table 2: <mark>Test-set word error rates on LibriSpeech</mark> after rescoring with base-sized models. None, GPT-2, and BERT results are from Salazar et al. (2020). <mark>Runtime</mark> is measured in passes through the transformer. "Clean" and "other" are easier and harder splits of the data. *We use a public re-implementation of OpenWebText.

- Electric slightly underperforms ELECTRA on downstream tasks, better than BERT.

- Pseudo-log-likelihood (PLL) scores, are used to re-rank the outputs of a speech recognition system, perform better and faster than masked models

$$\text{PLL}(\boldsymbol{x}) = \sum_{t=1}^{n} \log(\hat{p}_\theta(x_t|\boldsymbol{x}_{\backslash t})) = \sum_{t=1}^{n} -E(\boldsymbol{x})_t$$

# Content

## I. Basics for EBMs (45 min)

1. Probabilistic graphical modeling (PGM) framework and EBM model examples
2. Learning EBMs by Monte Carlo methods
3. Learning EBMs by noise-contrastive estimation (NCE)

$p_\theta(x)$

$p_\theta(h|x)$

I. Basics for EBMs

$p_\theta(x,h)$

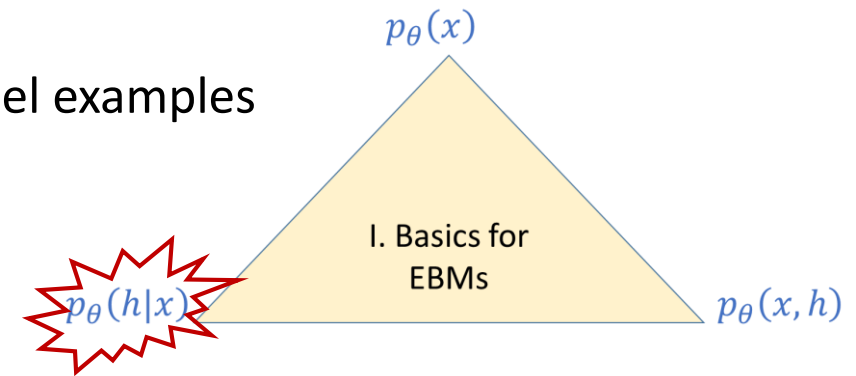## II. EBMs for language modeling (45 min)

1. Trans-dimensional random field (TRF) LMs for speech recognition
2. Residual energy-based models for text generation
3. Electric: an energy-based cloze model for representation learning over text

## III. EBMs for speech recognition and natural language labeling (45 min)

1. **CRFs as conditional EBMs**
2. CRFs for speech recognition
3. CRFs for sequence labeling in NLP

## IV. EBMs for semi-supervised natural language labeling (45 min)
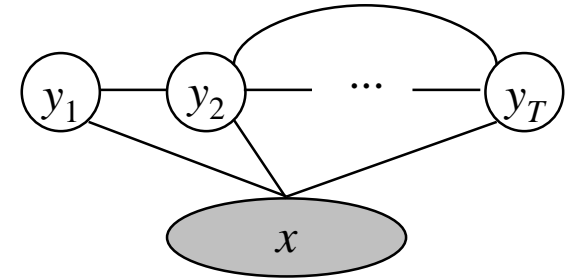
1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
3. JRFs for semi-supervised natural language labeling

# UGM Example - Conditional Random Fields (CRFs)

- A CRF is a conditional distribution $p(y|x)$ defined as a UGM/RF/EBM

$$p(y|x) = \frac{1}{Z(x)} exp \left[ \sum_C \phi_C(y_C, x) \right]$$



- $x$ is observed sequence, which is always given;

- $y$ is hidden sequence;

- $\phi_C(y_C, x)$ : Clique (log-)potential function over clique $C$ in the subgraph induced by $y$
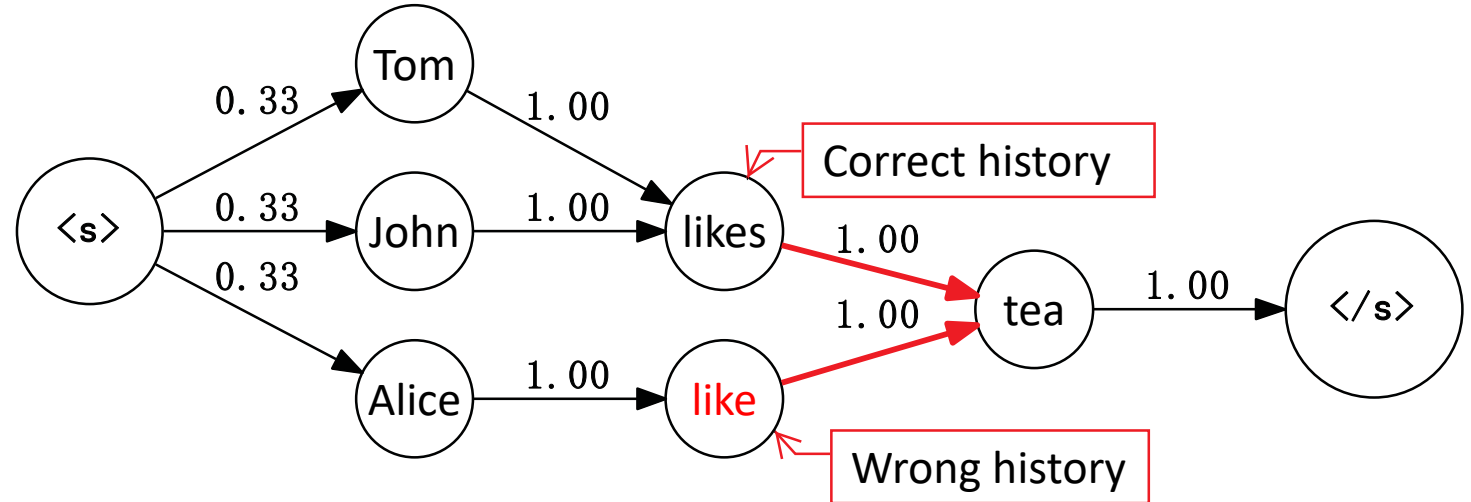
☺ CRFs can overcome "label bias" and "exposure bias" suffered by locally-normalized models

J. Lafferty, A. McCallum, and F. Pereira, "Conditional random fields: Probabilistic models for segmenting and labeling sequence data", ICML 2001.

# Label bias

▶ Word probabilities at each time-step are locally normalized, so successors of incorrect histories receive the same mass as do the successors of the true history. [Wiseman, et al., 2016]

**Training data**

Tom likes tea
John likes tea
Alice like tea



▶ [Andor, et al., 2016]

- "Intuitively, we would like the model to be able to revise an earlier decision made during search, when later evidence becomes available that rules out the earlier decision as incorrect."
- "the label bias problem means that locally normalized models often have a very weak ability to revise earlier decisions."
- A proof that globally normalized models are strictly more expressive than locally normalized models.

- Wiseman, et al., "Sequence-to-sequence Learning as Beam-Search Optimization", EMNLP, 2016.
- Andor, et al., "Globally Normalized Transition-Based Neural Networks", ACL, 2016.
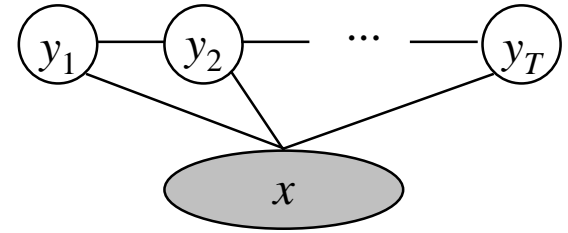
# Exposure bias

▶ Mismatch between **training** (teacher forcing) and **testing** (prediction) of locally-normalized sequence models:

- **Training**: maximize the likelihood of each successive target word, conditioned on the gold history of the target word.
- **Testing**: the model predict the next step, using its own predicted samples in testing.

▶ The model is never exposed to its own errors during training, and so the inferred histories at test-time do not resemble the gold training histories. [Wiseman, et al., 2016]

▶ **Exposure bias** results from training in a certain way (maybe alleviated by scheduled sampling), **Label bias** results from properties of the model itself.

- Wiseman, et al., "Sequence-to-sequence Learning as Beam-Search Optimization", EMNLP, 2016.
- Andor, et al., "Globally Normalized Transition-Based Neural Networks", ACL, 2016.

# Linear-chain CRFs

for sequence tagging, e.g. POS tagging, shallow parser, Chinese word segmentation, …

$$p(y_{1:T}|x) \propto exp\left\{\sum_{t=1}^{T-1}\phi_t(y_t,x)+\sum_{t=1}^{T-1}\psi_t(y_{t-1},y_t,x)\right\}$$



**Traditional**: **log-linear models** with discrete features

$$p(y_{1:T}|x) \propto exp\left\{\sum_{t=1}^{T-1}\sum_i \lambda_i f_i(y_t,x,t)+\sum_{t=1}^{T-1}\sum_j \mu_j f_j(y_{t-1},y_t,x,t)\right\}$$

- Node potential

$$\lambda_1 f_1(y_t,x,t) = \lambda_1 \cdot 1(y_t = \text{prep}, x_t = \text{on})$$

$$\lambda_2 f_i(y_t,x,t) = \lambda_2 \cdot 1(y_t = \text{adv}, x_t \text{ ends in ly})$$

- Edge potential

$$\mu_1 f_1(y_{t-1},y_t,x,t) = \mu_1 \cdot 1(y_t = \text{prep}, y_{t+1} = \text{non})$$

**Recently**: **neural CRFs**

Use NN to extract features
$$\text{LSTM}(x_{1:T}): x_{1:T} \to h_{1:T}$$

- Node potential, calculated via a linear layer

$$\phi_t(y_t = k, x) = w_k^{\mathrm{T}} h_t \triangleq \phi_t^k$$

$w_k$ is the weight vector for label $k$

- Edge potential, mostly implemented as a matrix $A$

66

# Training of CRFs in general

$$p_\theta(y|x) = \frac{1}{Z_\theta(x)} \exp[u_\theta(x,y)]$$

Normalizing constant:

$$Z_\theta(x) = \sum_y \exp[u_\theta(x,y)]$$

- Conditional Maximum likelihood (CML) training

The scaled log conditional likelihood of training data $\{(x_i, y_i), i = 1, \cdots, N\}$

$$L(\theta) \triangleq \frac{1}{N}\sum_{i=1}^N log p_\theta(y_i|x_i) = \frac{1}{N}\sum_{i=1}^N \{u_\theta(x_i, y_i) - \log Z_\theta(x_i)\}$$

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{1}{N}\sum_{i=1}^N \left\{ \frac{\partial u_\theta(x_i, y_i)}{\partial \theta} - E_{p_\theta(y|x_i)}\left[\frac{\partial u_\theta(x_i, y)}{\partial \theta}\right] \right\}$$

Expectation under empirical distribution

Expectation under model distribution $p_\theta(y|x_i)$

For linear-chain CRFs, this expectation can be exactly calculated.

67

# Training of Neural CRFs

Model $\quad p_\theta(y|x) = \dfrac{1}{Z_\theta(x)} \exp[u_\theta(x,y)],$ where $u_\theta(x,y) = \sum_t \phi_t^{y_t} + \sum_t A_{y_{t-1},y_t}$

For potential value $\phi_t^k, 1 \le t \le T, 1 \le k \le K$

$$\frac{\partial \log p(y|x)}{\partial \phi_t^k} = \delta(y_t = k) - E_{p(y|x)}[\delta(y_t = k)]$$

$$= \delta(y_t = k) - p(y_t = k|x)$$

i.e., the **error signal** received by the NN feature extractor during training

i.e., $\gamma_t^k$, the posterior **state occupation probability**, calculated using the alpha-beta variables from the forward-backward algorithm [Rabiner, 1989]
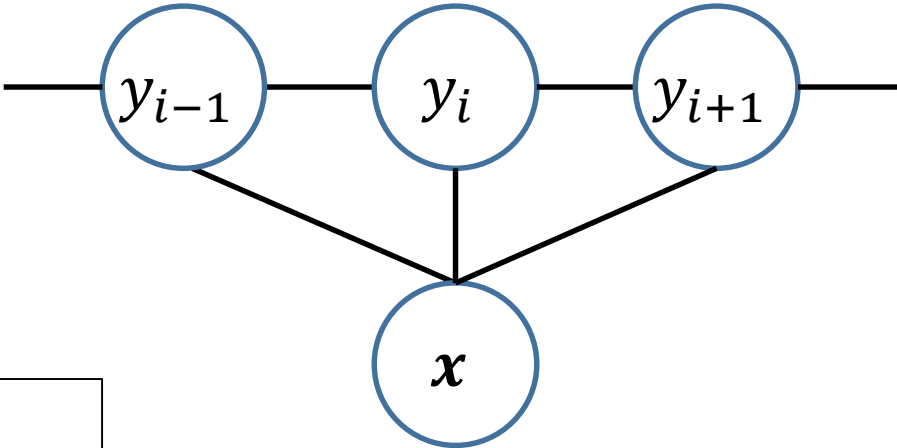
# Conditional random field (CRF) - Summary

A CRF define a conditional distribution over output sequence $y^l$ given input sequence $x^l$ of length $l$ :

$$p_\theta(y^l|x^l) = \frac{1}{Z_\theta(x^l)} \exp(u_\theta(x^l, y^l)) \qquad Z_\theta(x^l) = \sum_{y^l} \exp(u_\theta(x^l, y^l))$$

Potential for linear-chain:  Node potential  Edge potential

$$u_\theta(x^l, y^l) = \sum_{i=1}^{l} \phi_i(y_i, x^l) + \sum_{i=1}^{l} \psi_i(y_{i-1}, y_i, x^l)$$

☺ CRFs can overcome "label bias" and "exposure bias".

$y_{i-1}$ — $y_i$ — $y_{i+1}$

$x$

Example of a linear-chain CRF

Successfully applied for sequence labeling in NLP, less so for ASR

▶ CRFs was explored for phone classification, using zero, first and second order features [Gunawardana, et al., 2005].

▶ CTC-CRF: the first CRF successfully developed for end-to-end ASR

A. Gunawardana, et al.,"Hidden conditional random fields for phone classification", Europspeech, 2005.

# Content

I. **Basics for EBMs (45 min)**
1. Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
2. Learning EBMs by Monte Carlo methods
3. Learning EBMs by noise-contrastive estimation (NCE)

II. **EBMs for language modeling (45 min)**
1. Trans-dimensional random field (TRF) LMs for speech recognition
2. Residual energy-based models for text generation
3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for speech recognition and natural language labeling (45 min)**
1. CRFs as conditional EBMs
2. CRFs for speech recognition
3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
3. JRFs for semi-supervised natural language labeling

# End-to-end ASR: Basics

ASR is a *sequence discriminative* problem

- For acoustic observations $\boldsymbol{x} \triangleq x_1, \cdots, x_T$, find the most likely labels $\boldsymbol{y} \triangleq y_1, \cdots, y_L$

1. How to obtain $p(\boldsymbol{y} \mid \boldsymbol{x})$

2. How to handle alignment, since $L \neq T$

- Need a differentiable sequence-level loss of mapping acoustic sequence $\boldsymbol{y}$ to label sequence $\boldsymbol{x}$

- **Explicitly**: introduce hidden state sequence $\boldsymbol{\pi}$, as in Connectionist Temporal Classification (CTC), RNN Transducer (RNNT), CRF
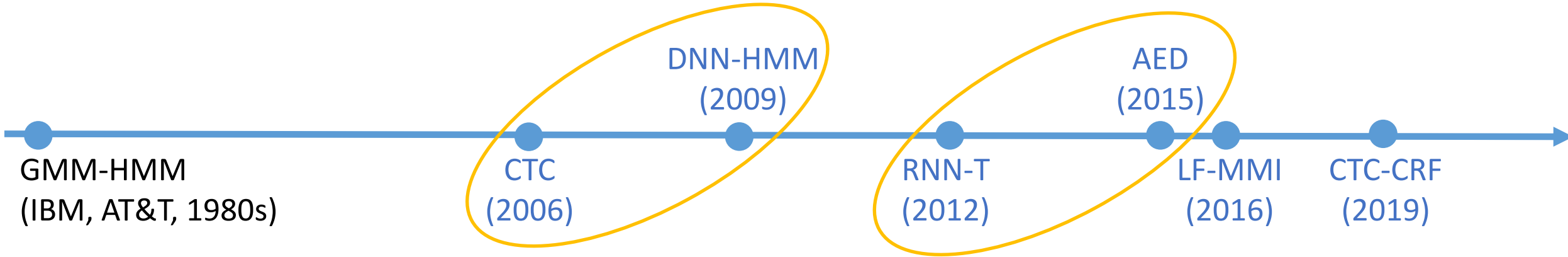- **Implicitly**: as in Attention based Encoder-Decoder (AED)

Labels

$\boldsymbol{y}$
$\parallel$
$y_1$
$\vdots$
$y_L$

$L \neq T$



Observations $\boldsymbol{x} = x_1 \cdots x_T$

Example of explicit alignment

# Brief ASR History



- [CTC] Graves, et al., "Connectionist Temporal Classification: Labelling unsegmented sequence data with RNNs", ICML 2006.
- [DNN-HMM] A. Mohamed, et al., "Deep belief networks for phone recognition", NIPS Workshop Deep Learning for Speech Recognition and Related Applications, 2009.
- [RNNT] A. Graves, "Sequence transduction with recurrent neural networks", ICML 2012 Workshop on Representation Learning.
- [AED] D. Bahdanau, et al., "Neural machine translation by jointly learning to align and translate", ICLR 2015.
- [LF-MMI] D. Povey, et al., "Purely sequence-trained neural networks for ASR based on lattice-free MMI", INTERSPEECH 2016.
- [CTC-CRF] Xiang&Ou. "CRF-based Single-stage Acoustic Modeling with CTC Topology", ICASSP, 2019.

# CTC: introducing blank symbol

- Motivation: training $p(\boldsymbol{y} \mid \boldsymbol{x})$ without the need for frame-level alignments between the acoustics $\boldsymbol{x}$ and the transcripts $\boldsymbol{y}$

  - Introduce a state sequence $\boldsymbol{\pi} \triangleq \pi_1, \cdots, \pi_T$, where $\pi_t \in$ the-alphabet-of-labels $\cup$ <b>



Path posterior
$$p(\boldsymbol{\pi} \mid \boldsymbol{x}) = \prod_{t=1}^{T} \underset{\text{State posterior}}{p(\pi_t \mid \boldsymbol{x})}$$

Linear&Softmax Layer
$$z_t = W h_t \in \mathbb{R}^{K+1}$$

$$p(\pi_t = k \mid \boldsymbol{x}) = \frac{exp(z_t^k)}{\sum_i exp(z_t^i)} \triangleq p_t^k : \text{the}$$

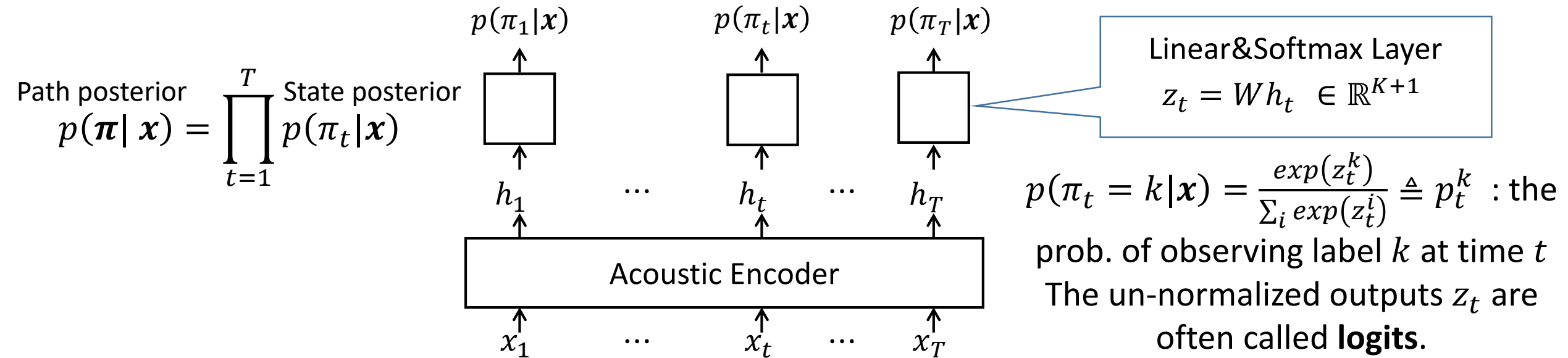prob. of observing label $k$ at time $t$
The un-normalized outputs $z_t$ are often called **logits**.

Graves, et al., "Connectionist Temporal Classification: Labelling unsegmented sequence data with RNNs", ICML 2006.

# CTC topology

- **State topology** refers to the state transition structure in $\boldsymbol{\pi}$, which basically determines the mapping $\mathcal{B}_{CTC}$ from $\boldsymbol{\pi}$ to $\boldsymbol{y}$

CTC topology : a mapping $\mathcal{B}_{CTC}$ maps $\boldsymbol{\pi}$ to $\boldsymbol{y}$ by
1. reducing repetitive symbols to a single symbol;
2. removing all blank symbols.

$$\mathcal{B}(-CC--AA-T-) = CAT$$

Path posterior

$$p(\boldsymbol{\pi}|\boldsymbol{x}) = \prod_{t=1}^{T} p(\pi_t|\boldsymbol{x})$$

Label-seq posterior

$$p(\boldsymbol{y}|\boldsymbol{x}) = \sum_{\boldsymbol{\pi}:\,\mathcal{B}_{CTC}(\boldsymbol{\pi})=\boldsymbol{y}} p(\boldsymbol{\pi}|\boldsymbol{x})$$

Summing over all possible paths, which map to $\boldsymbol{y}$



B B **c** B B **a** **a** B B **t**

B **c** **c** B **a** B B B B **t**

· · ·

B **c** B B **a** B B **t** **t** B

# CTC: shortcoming

- Conditional independence assumption

Overcome

RNN-T

CTC-CRF

$$p(\boldsymbol{\pi}| \boldsymbol{x}) = \prod_{t=1}^{T} p(\pi_t|\boldsymbol{x})$$

$p(\pi_1|\boldsymbol{x})$  $p(\pi_t|\boldsymbol{x})$  $p(\pi_T|\boldsymbol{x})$

$h_1$  $\cdots$  $h_t$  $\cdots$  $h_T$

Acoustic Encoder

$x_1$  $\cdots$  $x_t$  $\cdots$  $x_T$

Computational flow

$\pi_{t-1}$  $\pi_t$  $\pi_{t+1}$

$\boldsymbol{x}$

Graphical Model Representation

# Section Content

1. Motivation

2. Related work

3. Method: **CTC-CRF**

4. Experiments

5. Conclusion

- H. Xiang, Z. Ou. "CRF-based Single-stage Acoustic Modeling with CTC Topology", ICASSP, 2019.
- K. An, H. Xiang, Z. Ou. "CAT: A CTC-CRF based ASR Toolkit Bridging the Hybrid and the End-to-end Approaches towards Data Efficiency and Low Latency", INTERSPEECH, 2020.
- Fan, et al., "The SLT 2021 children speech recognition challenge: Open datasets, rules and baselines", SLT, 2021.
- H. Zheng, W. Peng, Z. Ou, J. Zhang. "Advancing CTC-CRF Based End-to-End Speech Recognition with Wordpieces and Conformers", arXiv:2107.03007, 2021.

# Motivation: data-efficient end2end

- End-to-end system:
  - Eliminate the construction of GMM-HMMs and phonetic decision-trees, and can be trained from scratch (flat-start or single-stage)

- In a more strict/ambitious sense:
  - Remove the need for a pronunciation lexicon and, even further, train the acoustic and language models jointly rather than separately
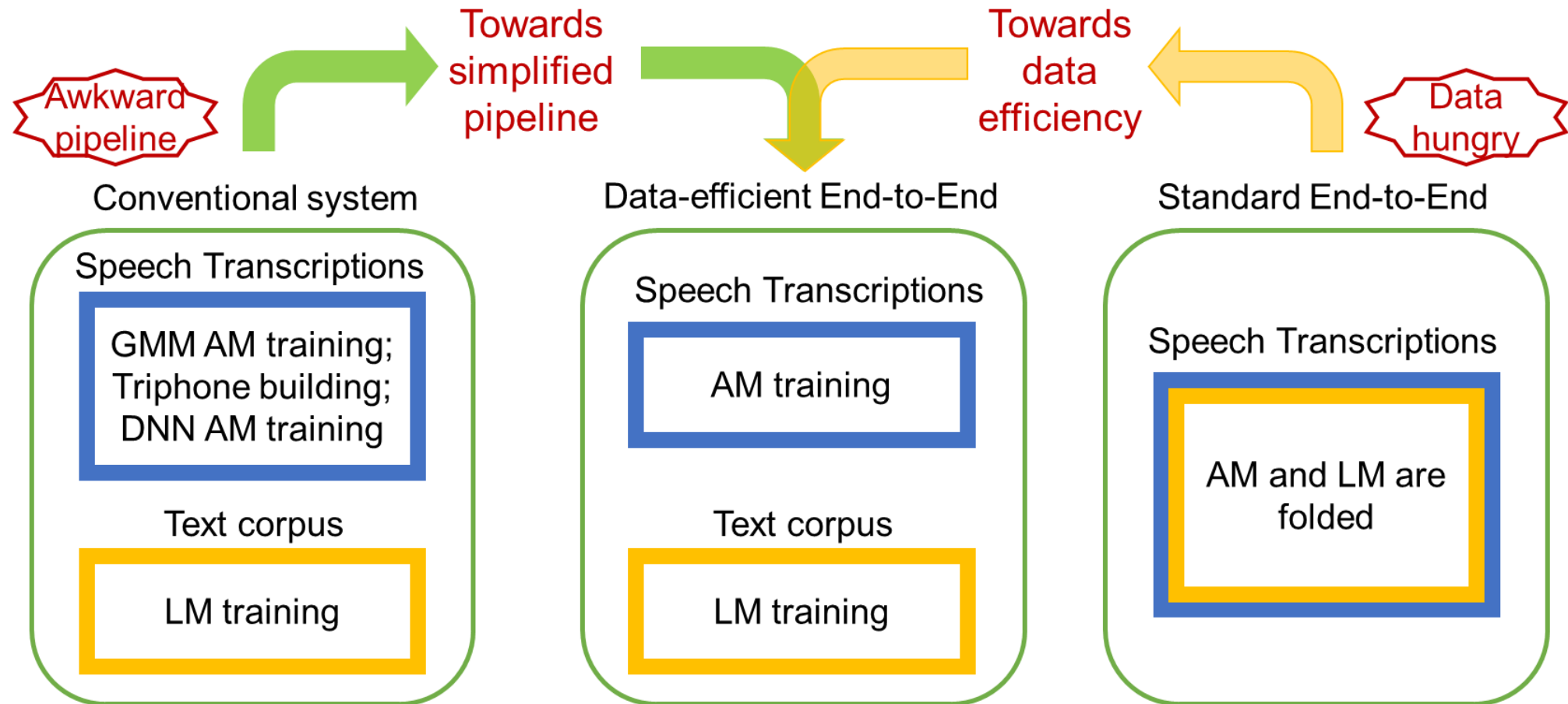  - Data-hungry

We need data-efficient end2end speech recognition, which can flexibly use a separate language model (LM) with or without a pronunciation lexicon.

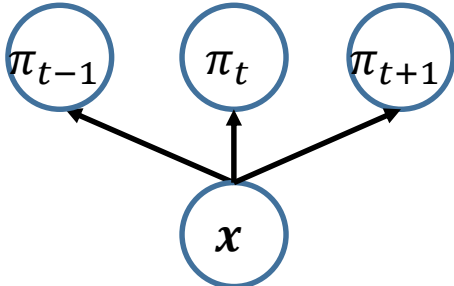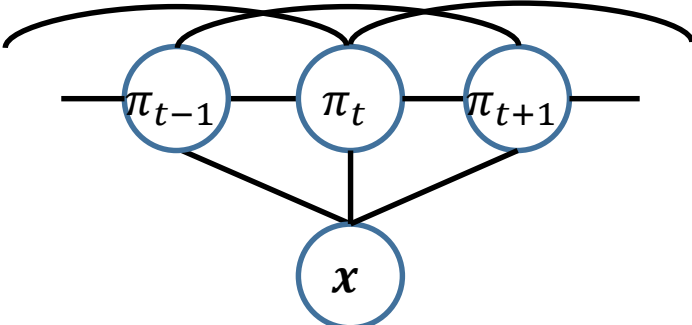- Text corpus for language modeling are cheaply available.
- Data-efficient

# Motivation: bridging

**Modularization promote Data-efficiency**

✓ Keep necessary factorization of AM and LM

# CTC vs CTC-CRF

| CTC | CTC-CRF |
|---|---|
| $p(\boldsymbol{y}|\boldsymbol{x}) = \sum_{\boldsymbol{\pi}:\mathcal{B}(\boldsymbol{\pi})=\boldsymbol{y}} p(\boldsymbol{\pi}|\boldsymbol{x})$, using CTC topology $\mathcal{B}$ ||

State Independence

$$p(\boldsymbol{\pi}|\boldsymbol{x};\boldsymbol{\theta}) = \prod_{t=1}^{T} p(\pi_t|\boldsymbol{x})$$

$$p(\boldsymbol{\pi}|\boldsymbol{x};\boldsymbol{\theta}) = \frac{e^{\phi(\boldsymbol{\pi},\boldsymbol{x};\boldsymbol{\theta})}}{\sum_{\boldsymbol{\pi}'} e^{\phi(\boldsymbol{\pi}',\boldsymbol{x};\boldsymbol{\theta})}}$$

Node potential, by NN

$$\phi(\boldsymbol{\pi},\boldsymbol{x};\boldsymbol{\theta}) = \begin{pmatrix} \Sigma_{t=1}^{T}\log p(\pi_t|\boldsymbol{x}) \\ + \log p_{LM}(\mathcal{B}(\boldsymbol{\pi})) \end{pmatrix}$$

Edge potential,

by n-gram denominator LM of labels, like in LF-MMI

$$\frac{\partial \log p(\boldsymbol{y}|\boldsymbol{x};\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\boldsymbol{\pi}|\boldsymbol{y},\boldsymbol{x};\boldsymbol{\theta})}\left[\frac{\partial \log p(\boldsymbol{\pi}|\boldsymbol{x};\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right]$$

$$\frac{\partial \log p(\boldsymbol{y}|\boldsymbol{x};\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbb{E}_{p(\boldsymbol{\pi}|\boldsymbol{x},\boldsymbol{y};\boldsymbol{\theta})}\left[\frac{\partial \phi(\boldsymbol{\pi},\boldsymbol{x};\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right] - \mathbb{E}_{p(\boldsymbol{\pi}'|\boldsymbol{x};\boldsymbol{\theta})}\left[\frac{\partial \phi(\boldsymbol{\pi}',\boldsymbol{x};\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}\right]$$

# Training of CTC-CRFs

Model $\quad p(\boldsymbol{y}|\boldsymbol{x}) = \displaystyle\sum_{\boldsymbol{\pi}:\mathcal{B}(\boldsymbol{\pi})=\boldsymbol{y}} p(\boldsymbol{\pi}|\boldsymbol{x})$, where $p(\boldsymbol{\pi}|\boldsymbol{x}) = \dfrac{e^{\phi(\boldsymbol{\pi},\boldsymbol{x})}}{\sum_{\boldsymbol{\pi'}} e^{\phi(\boldsymbol{\pi'},\boldsymbol{x})}}$

Node potential $\qquad$ Edge potential

$$\phi(\boldsymbol{\pi}, \boldsymbol{x}) = \Sigma_{t=1}^{T}\log p(\pi_t|\boldsymbol{x}) + \log p_{LM}(\mathcal{B}(\boldsymbol{\pi}))$$

Denote $p(\pi_t = k|\boldsymbol{x}) = \phi_t^k$, then for potential value $\phi_t^k, 1 \leq t \leq T, 1 \leq k \leq K+1$

$$\frac{\partial \log p(\boldsymbol{y}|\boldsymbol{x})}{\partial \phi_t^k} = \mathbb{E}_{p(\boldsymbol{\pi}|\boldsymbol{x},\boldsymbol{y})}\left[\frac{\partial \phi(\boldsymbol{\pi},\boldsymbol{x})}{\partial \phi_t^k}\right] - \mathbb{E}_{p(\boldsymbol{\pi'}|\boldsymbol{x})}\left[\frac{\partial \phi(\boldsymbol{\pi'},\boldsymbol{x})}{\partial \phi_t^k}\right]$$

$$= E_{p(\boldsymbol{\pi}|\boldsymbol{x},\boldsymbol{y})}[\delta(\boldsymbol{\pi}_t = k)] - E_{p(\boldsymbol{\pi'}|\boldsymbol{x})}[\delta(\boldsymbol{\pi'}_t = k)]$$

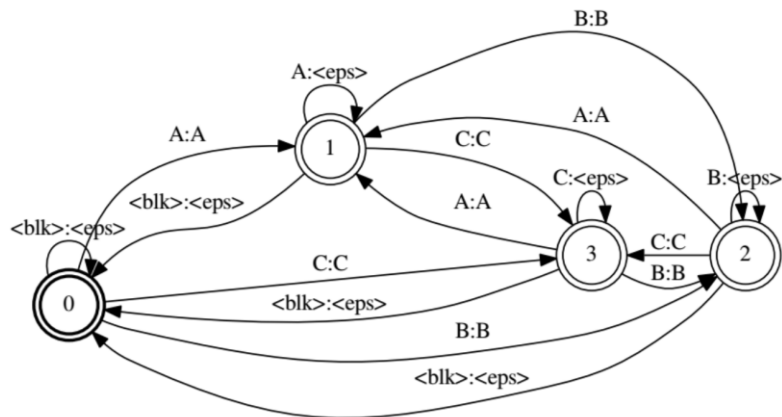i.e., the **error signal** received the acoustic encoder NN during training

i.e., the posterior **state occupation probability**, by running the FB algorithm over the WFST determined by $\boldsymbol{y}$

i.e., the posterior **state occupation probability**, by running the FB algorithm over the WFST determined by n-gram LM of labels
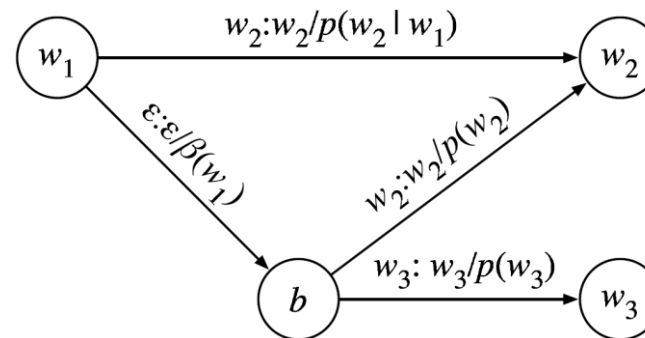
# Edge potential in CTC-CRF

$$\log p_{LM}(\mathcal{B}(\boldsymbol{\pi}))$$

**WFST representing CTC topology: T**

**n-gram denominator LM of labels: G**



$T \circ G$: Composed into a single WFST

The arc in the WFST is treated as the state



$$p(\pi_t = k | \boldsymbol{x}) = \sum_{\text{input symbol of } a \text{ is } k} p(a_t = a | \boldsymbol{x})$$

# Related work

- **Directed Graphical Model/Locally normalized**

  - DNN-HMM : Model $p(\boldsymbol{\pi}, \boldsymbol{x})$ as an HMM, could be discriminatively trained, e.g. by $\max_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\boldsymbol{y} \mid \boldsymbol{x})$

  - CTC : $p(\boldsymbol{\pi} \mid \boldsymbol{x}) = \prod_{t=1}^{T} p(\pi_t \mid \boldsymbol{x})$

  - RNNT : $p(\pi_{1:T+U} \mid x_{1:T}) = \prod_{j=1}^{T+U} p(\pi_j \mid \pi_{1:j-1})$

  - AED : $p(\boldsymbol{y} \mid \boldsymbol{x}) = \prod_{i=1}^{L} p(y_i \mid y_1, \cdots, y_{i-1}, \boldsymbol{x})$

- **Undirected Graphical Model/Globally normalized**

  - CRF : $p(\boldsymbol{\pi} \mid \boldsymbol{x}) \propto exp[\phi(\boldsymbol{\pi}, \boldsymbol{x})]$

  CTC-CRF is fundamentally different from all history models!

# Related work (SS-LF-MMI/EE-LF-MMI)

- ## Single-Stage (SS) Lattice-Free Maximum-Mutual-Information (LF-MMI)

  - 10 - 25% relative WER reduction on 80-h WSJ, 300-h Switchboard and 2000-h Fisher+Switchboard datasets, compared to CTC, Seq2Seq, RNN-T.

  - Cast as MMI-based discriminative training of an HMM (generative model) with

    *Pseudo state-likelihoods calculated by the bottom DNN,*

    *Fixed state-transition probabilities.*

  - 2-state HMM topology

  - Including a silence label

CTC-CRF
- Cast as a CRF;
- CTC topology;
- No silence label.

Hadian, et al., "Flat-start single-stage discriminatively trained HMM-based models for ASR", T-ASLP 2018.

# SS-LF-MMI vs CTC-CRF

| | SS-LF-MMI | CTC-CRF |
|---|---|---|
| State topology | HMM topology with two states | CTC topology |
| Silence label | Using silence labels.<br>Silence labels are randomly inserted when estimating denominator LM. | No silence labels. Use <blk> to absorb silence.<br>☺ No need to insert silence labels to transcripts. |
| Decoding | No spikes. | The posterior is dominated by <blk> and non-blank symbols occur in spikes.<br>☺ Speedup decoding by skipping blanks. |
| Implementation | Modify the utterance length to one of 30 lengths; use leaky HMM. | ☺ No length modification; no leaky HMM. |

# Experiments

- We conduct our experiments on three benchmark datasets:
  - WSJ 80 hours
  - Switchboard 300 hours
  - Librispeech 1000 hours

- Acoustic model: 6 layer BLSTM with 320 hidden dim, 13M parameters

- Adam optimizer with an initial learning rate of 0.001, decreased to 0.0001 when cv loss does not decrease

- Implemented with Pytorch.

- Objective function (use the CTC objective function to help convergences):
$$\mathcal{J}_{CTC-CRF} + \alpha \mathcal{J}_{CTC}$$

- Decoding score function (use word-based language models, WFST based decoding):
$$\log p(\boldsymbol{l}|\boldsymbol{x}) + \beta \log p_{LM}(\boldsymbol{l})$$

H. Xiang, Z. Ou. "CRF-based Single-stage Acoustic Modeling with CTC Topology", ICASSP, 2019.

# Experiments (Comparison with CTC, phone based)

**WSJ 80h**

| Model | Unit | LM | SP | dev93 | eval92 |
|-------|------|-----|-----|--------|--------|
| CTC | Mono-phone | 4-gram | N | 10.81% | 7.02% |
| CTC-CRF | Mono-phone | 4-gram | N | 6.24% | 3.90% |

44.4%

**Switchboard 300h**

| Model | Unit | LM | SP | SW | CH |
|-------|------|-----|-----|-----|-----|
| CTC | Mono-phone | 4-gram | N | 12.9% | 23.6% |
| CTC-CRF | Mono-phone | 4-gram | N | 11.0% | 21.0% |

14.7%   11%

**Librispeech 1000h**

| Model | Unit | LM | SP | Dev Clean | Dev Other | Test Clean | Test Other |
|-------|------|-----|-----|-----------|-----------|------------|------------|
| CTC | Mono-phone | 4-gram | N | 4.64% | 13.23% | 5.06% | 13.68% |
| CTC-CRF | Mono-phone | 4-gram | N | 3.87% | 10.28% | 4.09% | 10.65% |

19.1%   22.1%

SP: speed perturbation for 3-fold data augmentation.

# Experiments (Comparison with STOA)

**Switchboard 300h**

| Model | SW | CH | Average | Source |
|---|---|---|---|---|
| Kaldi chain triphone | 9.6 | 19.3 | 14.5 | IS 2016 |
| Kaldi e2e chain monophone | 11.0 | 20.7 | 15.9 | ASLP 2018, 26M |
| Kaldi e2e chain biphone | 9.8 | 19.3 | 14.6 | ASLP 2018, 26M |
| CTC-CRF monophone | 10.3 | 19.7 | 15.0 | ICASSP 2019, BLSTM, 13M |
| **CTC-CRF monophone** | **9.8** | **18.8** | **14.3** | **IS 2020, VGG BLSTM, 16M** |

**10%**

RWTH IS 2018, "Improved training of end-to-end attention models for speech recognition".
RWTH IS 2019, "RWTH ASR Systems for LibriSpeech Hybrid vs Attention -- Data Augmentation".
IBM IS19, "Forget a Bit to Learn Better Soft Forgetting for CTC-based Automatic Speech Recognition".
Espnet ASRU19, "Espresso: A Fast End-to-end Neural Speech Recognition Toolkit".
Google IS19, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition".

# Experiments (Comparison with STOA)

**Librispeech 1000h**

| Model | Test Clean | Test Other | Source |
|---|---|---|---|
| Kaldi chain triphone | 4.28 | - | IS 2016 |
| **CTC-CRF monophone** | **4.0** | **10.6** | **ICASSP 2019, BLSTM (6,320), 13M** |

RWTH IS 2018, "Improved training of end-to-end attention models for speech recognition".
RWTH IS 2019, "RWTH ASR Systems for LibriSpeech Hybrid vs Attention -- Data Augmentation".
IBM IS19, "Forget a Bit to Learn Better Soft Forgetting for CTC-based Automatic Speech Recognition".
Espnet ASRU19, "Espresso: A Fast End-to-end Neural Speech Recognition Toolkit".
Google IS19, "SpecAugment: A Simple Data Augmentation Method for Automatic Speech Recognition".

# Mandarin Aishell-1 results

- 170 hours mandarin speech corpus
- 400 speakers from different accent areas
- 15% CER reduction compared with LF-MMI
- 5% CER reduction compared with end-to-end transformer

| Model | %CER |
|---|---|
| LF-MMI with i-vector [1] | 7.43 |
| Transformer [2] | 6.7 |
| CTC-CRF [3] | 6.34 |

[1] D. Povey, A. Ghoshal, and et al, "The Kaldi speech recognition toolkit," ASRU 2011.
[2] S. Karita, N. Chen, and et al, "A comparative study on transformer vs RNN in speech applications," ASRU 2019.
[3] Keyu An, Hongyu Xiang, and **Zhijian Ou**, "CAT: A CTC-CRF based ASR toolkit bridging the hybrid and the end-to-end approaches towards data efficiency and low latency," INTERSPEECH 2020.

2021 SLT CHILDREN
SPEECH RECOGNITION CHALLENGE (CSRC)

ORGANIZER：西北工业大学　清华大学　厦门大學　Databaker technology　CCF

- 400 hours of data, targeting to boost children speech recognition research.
- Evaluated on 10 hours of children's reading and conversational speech.
- 3 baselines (Chain model, Transformer and CTC-CRF) are provided.

| model | Kaldi Chain model | Espnet Transformer | CTC-CRF |
|-------|-------------------|--------------------|---------| 
| CER% | 28.75 | 27.28 | **25.34** |

Fan Yu, Zhuoyuan Yao, Xiong Wang, Keyu An, Lei Xie, **Zhijian Ou**, Bo Liu, Xiulin Li, Guanqiong Miao. The SLT 2021 children speech recognition challenge: Open datasets, rules and baselines. SLT 2021.

# Advancing CTC-CRF Based End-to-End Speech Recognition with Wordpieces and Conformers

Huahuan Zheng, Wenjie Peng, **Zhijian Ou** and Jinsong Zhang, arXiv:2107.03007

| Basic Units of Labels | Label Sequence |
|---|---|
| phoneme | DH AE1 T N IY1 DH ER0 AH1 V DH EH1 M HH AE1 D K R AO1 S T DH AH0 TH R EH1 SH OW2 L D S IH1 N S DH AH0 D AA1 R K D EY1 |
| character /grapheme | t h a t _ n e i t h e r _ o f _ t h e m _ h a d _ c r o s s e d _ t h e _ t h r e s h o l d _ s i n c e _ t h e _ d a r k _ d a y _ |
| subword /wordpiece | that_ ne i ther_ of_ them_ had_ cro s sed_ the_ th re sh old_ sin ce_ the_ d ar k_ day_ |
| word | that neither of them had crossed the threshold since the dark day |

# Experiments (Comparison between different units, WER%)

**Switchboard 300h**

| Model | Unit | LM | Augmentation | Eval2000 | SW | CH |
|---|---|---|---|---|---|---|
| Conformer (this work) | monophone | 4-gram | SP, SA | 12.1 | 7.9 | 16.1 |
| | monophone | Trans.* | SP, SA | 10.7 | 6.9 | 14.5 |
| | wordpiece | 4-gram | SP, SA | 12.7 | 8.7 | 16.5 |
| | wordpiece | Trans.* | SP, SA | 11.1 | 7.2 | 14.8 |

**Librispeech 1000h**

| Model | Unit | LM | Augmentation | Test Clean | Test Other |
|---|---|---|---|---|---|
| Conformer (this work) | monophone | 4-gram | SA | 3.61 | 8.10 |
| | monophone | Trans.** | SA | 2.51 | 5.95 |
| | wordpiece | 4-gram | SA | 3.59 | 8.37 |
| | wordpiece | Trans.** | SA | 2.54 | 6.33 |

SP: speed perturbation for 3-fold data augmentation.

SA: our implementation of SpecAug with ratio

* Latest Kaldi Transformer LM rescoring

** RWTH 42-layer Transformer

English: a low degree of grapheme-phoneme correspondence

# Experiments (Comparison between different units, WER%)

**CommonVoice German 700h**

| Model | #params | unit | LM | Augmentation | Test |
|---|---|---|---|---|---|
| Conformer (This work) | 25.03 | char | 4-gram | SP, SA | 12.7 |
| | 25.03 | char | Trans. | SP, SA | 11.6 |
| | 25.03 | monophone | 4-gram | SP, SA | 10.7 |
| | 25.03 | monophone | Trans. | SP, SA | 10.0 |
| | 25.06 | wordpiece | 4-gram | SP, SA | 10.5 |
| | 25.06 | wordpiece | Trans. | SP, SA | 9.8 |

German: a high degree of grapheme-phoneme correspondence

# Experiments (Comparison with STOA)

**Switchboard 300h**

| Model | #params | LM | unit | SW | CH | Eval2000 |
|---|---|---|---|---|---|---|
| RNN-T, 2021 [10] | 57 | RNN LM | char | 6.4 | 13.4 | 9.9 |
| Conformer [9] | 44.6 | Trans. | bpe | 6.8 | 14.0 | 10.4 |
| TDNN-F [11] | - | Trans.* | triphone | 7.2 | 14.4 | 10.8 |
| TDNN-F [11] | - | Trans.** | triphone | 6.5 | 13.9 | 10.2 |
| VGGBLSTM [2] | 39.15 | RNN LM | monophone | 8.8 | 17.4 | [13.0] |
| Conformer (This work) | 51.82 | Trans. | monophone | 6.9 | 14.5 | 10.7 |
| | 51.85 | Trans. | wordpiece | 7.2 | 14.8 | 11.1 |

* N-best rescoring, ** Iterative lattice rescoring

[2] "CAT: A CTC-CRF based ASR toolkit bridging the hybrid and the end-to-end approaches towards data efficiency and low latency," INTERSPEECH 2020.
[9] "Conformer: Convolution-augmented Transformer for Speech Recognition", Interspeech 2020.
[10] "Advancing RNN transducer technology for speech recognition," ICASSP 2021.
[11] "A paralleliz- able lattice rescoring strategy with neural language models," ICASSP, 2021

# Section Conclusion

- The CTC-CRF framework inherits the data-efficiency of the hybrid approach and the simplicity of the end-to-end approach.

- CTC-CRF significantly outperforms regular CTC on a wide range of benchmarks, and is on par with other state-of-the-art end-to-end models.
  - English WSJ-80h, Switchboard-300h, Librispeech-1000h; Mandarin Aishell-170h; …

- Flexibility
  - Streaming ASR <- INTRESPEECH 2020
  - Neural Architecture Search <- SLT 2021
  - Children Speech Recognition <- SLT 2021
  - Wordpieces, Conformer architectures
  - Multilingual and Crosslingual <- ASRU2021
  - …

https://github.com/thu-spmi/cat

# Content

I.    Basics for EBMs (45 min)
   1. Probabilistic graphical modeling (PGM) framework and EBM model examples (classic & modern)
   2. Learning EBMs by Monte Carlo methods
   3. Learning EBMs by noise-contrastive estimation (NCE)

II. EBMs for language modeling (45 min)
   1. Trans-dimensional random field (TRF) LMs for speech recognition
   2. Residual energy-based models for text generation
   3. Electric: an energy-based cloze model for representation learning over text

III. EBMs for speech recognition and natural language labeling (45 min)
   1. CRFs as conditional EBMs
   2. CRFs for speech recognition
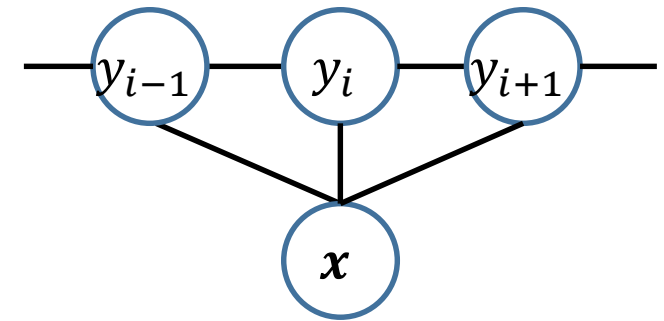→ 3. CRFs for sequence labeling in NLP

IV. EBMs for semi-supervised natural language labeling (45 min)
   1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
   2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
   3. JRFs for semi-supervised natural language labeling

# Motivation

- Conditional random fields (CRFs) have been shown to be one of the most successful approaches to sequence labeling.

- Various linear-chain neural CRFs (NCRFs) have been developed
  - Node potential modeling is improved by using NNs
  - But the linear-chain structure is still kept, i.e. using a bigram table as the edge potential

  - Linear-chain NCRFs capture only first-order [1] interactions and neglect higher-order dependencies between labels, which can be potentially useful in real-world sequence labeling applications

[1] *Fixed n-th order can be cast as first-order.*

Linear-chain CRF

How can we improve CRFs to capture long-range dependencies in the label sequence (preferably non-Markovian)?

# Related work

- Infinite-Order CRFs (based on the Hierarchical Pitman-Yor Process) [a], Semi-Markov CRFs [b], Latent-dynamic CRFs [c], but not by using NNs

- Attention-based encoder-decoder (AED) and RNNT exploit non-Markovian dependences between labels, but both are locally normalized models and thus suffer from the label bias and exposure bias problems

- [d] extends AED, by removing the final softmax in the RNN decoder to learn global sequence scores, but cast as a non-probabilistic variant of the seq2seq model

- [e] proposes **neural CRF transducers**: RNNT+CRF

a. Sotirios P. Chatzis and Yiannis Demiris, "The Infinite-Order Conditional Random Field Model for Sequential Data Modeling," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 35, pp. 1523–1534, 2013.
b. Sunita Sarawagi and William W. Cohen, "Semi-Markov Conditional Random Fields for Information Extraction," in NIPS, 2004.
c. Louis-Philippe Morency, Ariadna Quattoni, and Trevor Darrell, "Latent-Dynamic Discriminative Models for Continuous Gesture Recognition," in CVPR, 2007.
d. Wiseman, et al., "Sequence-to-sequence Learning as Beam-Search Optimization", EMNLP, 2016.
e. Kai Hu, Zhijian Ou, et al. Neural CRF Transducers for Sequence Labeling. ICASSP, 2019.
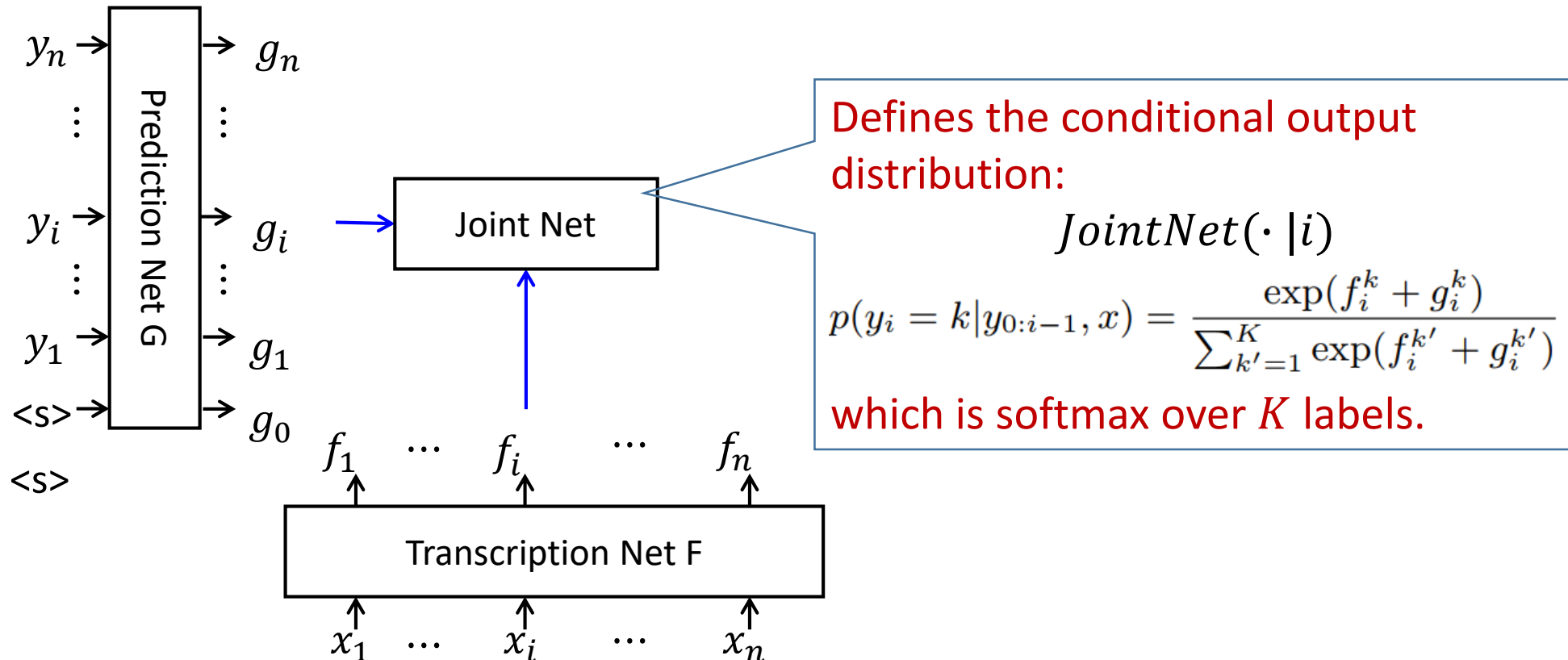
# Related work - RNNT

- ## RNN Transducers (RNNT)

  - Originally developed for general sequence-to-sequence learning, which do not assume that the input and output sequences are of equal lengths and aligned, e.g., in speech recognition

  - In the following, we introduce RNN transducers in a simple form for applications in sequence labeling: i.e., for the aligned setting: one label for one observation in each position

  - Similar idea can be used to revise general RNNT

A. Graves, "Sequence transduction with recurrent neural networks," ICML 2012 Workshop on Representation Learning.

# RNNT: introducing prediction network for labels

- Motivation: extending CTC by considering output-output dependencies

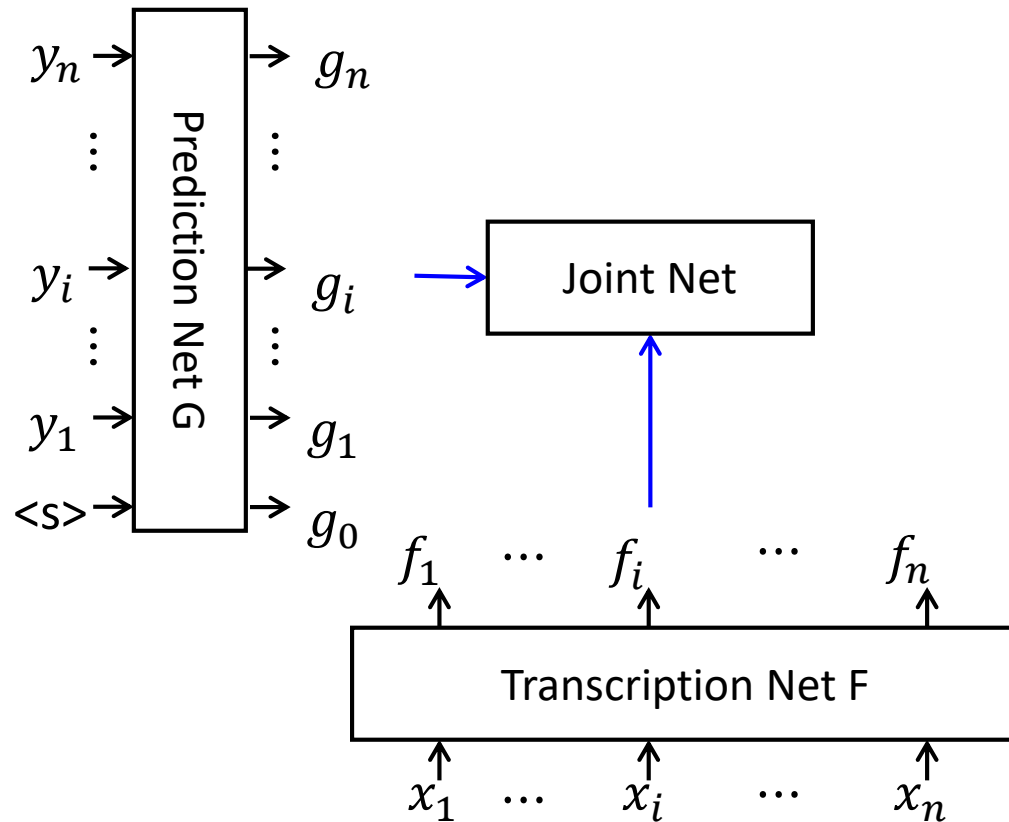- Introduce the prediction network, which attempts to model each output in $y_{1:n}$ given the previous ones
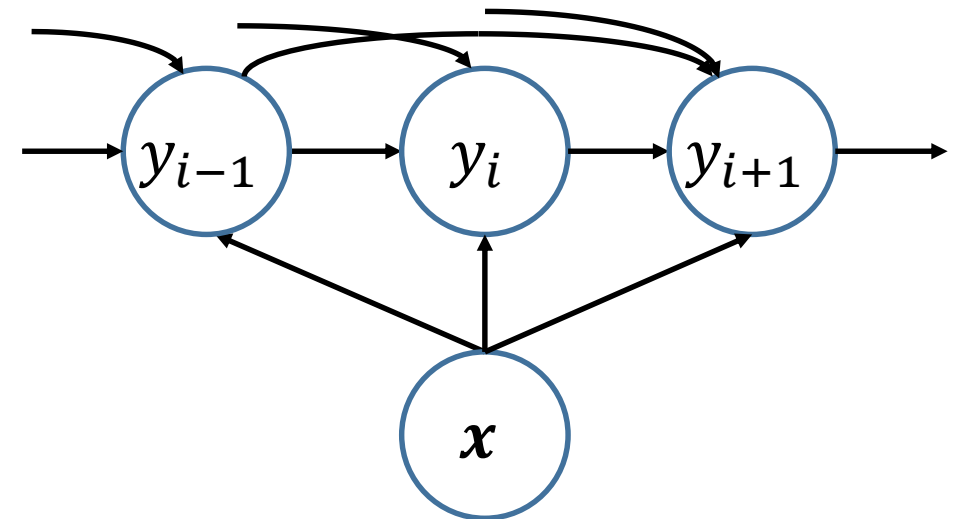
$$p(y|x) = \prod_{i=1}^{n} p(y_i|y_{0:i-1}, x)$$



Defines the conditional output distribution:

$$JointNet(\cdot \,|i)$$

$$p(y_i = k|y_{0:i-1}, x) = \frac{\exp(f_i^k + g_i^k)}{\sum_{k'=1}^{K} \exp(f_i^{k'} + g_i^{k'})}$$

which is softmax over $K$ labels.

$y_0$ is the special token <s>

$f_i, g_i \in \mathbb{R}^K$

A. Graves, "Sequence transduction with recurrent neural networks," ICML 2012 Workshop on Representation Learning.

# RNNT: shortcoming

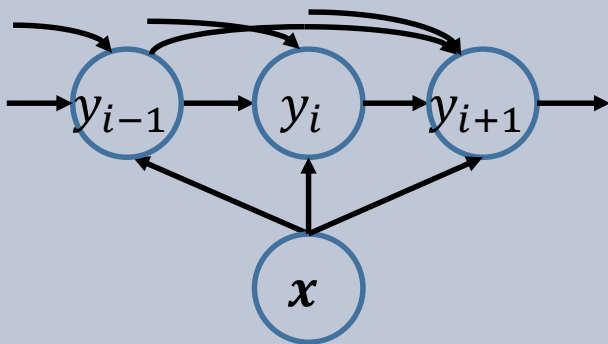

$$p(y|x) = \prod_{i=1}^{n} p(y_i|y_{0:i-1}, x)$$

Graphical Model Representation

- As directed sequential model /Auto-regressive model, RNNT potentially suffers from Exposure Bias and Label Bias. A recent effort in [Cui, et al., 2021].

X. Cui, et al., "Reducing Exposure Bias in Training Recurrent Neural Network Transducers", INTERSPEECH 2021.

# Neural CRF Transducer

$$p(y|x) = \prod_{i=1}^{n} p(y_i|y_{0:i-1}, x)$$

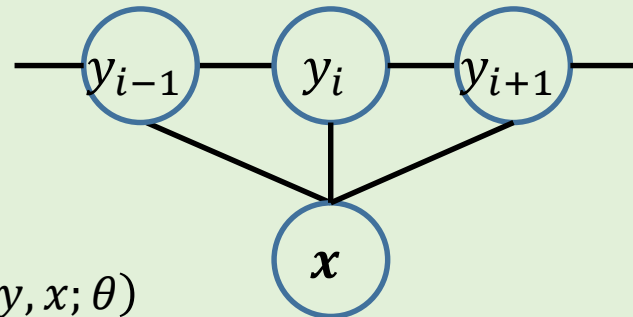$$p(y|x;\theta) = \frac{\exp\{u(y,x;\theta)\}}{Z(x;\theta)}, \text{where } Z(x;\theta) = \sum_{y'} \exp\{u(y',x;\theta)\}$$

### RNNT



$$p(y_i = k|y_{0:i-1}, x) = \frac{\exp(f_i^k + g_i^k)}{\sum_{k'=1}^{K} \exp(f_i^{k'} + g_i^{k'})}$$
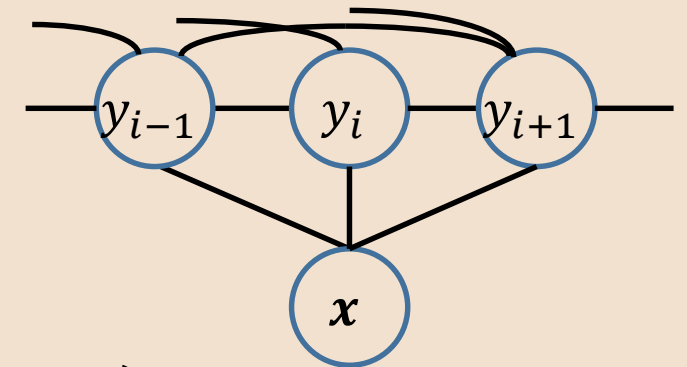
Local conditional

### Linear-chain CRF



$$u(y,x;\theta) = \sum_{i=1}^{n} \{\phi_i(y_i,x;\theta) + \psi_i(y_{i-1}, y_i;\theta)\}$$

Local log-potential: $f_i^k + A_{jk}$
for $y_{i-1}=j$, $y_i=k$

### Neural CRF Transducer



$$u(y,x;\theta) = \sum_{i=1}^{n} \{\phi_i(y_i,x;\theta) + \psi_i(y_{0:i-1}, y_i;\theta)\}$$

Local log-potential: $f_i^k + g_i^k$ for $y_i=k$

Kai Hu, Zhijian Ou, et al. Neural CRF Transducers for Sequence Labeling. ICASSP, 2019.

# Neural CRF Transducer: training and decoding

▶ **Training**

- Negative log-likelihood over input seq. $x$ and oracle label seq. $y^*$

$$L(y^*; \theta) = -u(y^*, x; \theta) + \textcolor{red}{logZ(x; \theta)}$$

- Monte Carlo Method

$$\nabla_\theta logZ(x; \theta) = \textcolor{red}{E_{p(y'|x;\theta)}}[\nabla_\theta u(y', x; \theta)]$$

- Beam search with early updates [b]

$$L(y^*_{1:j}; \theta) = -u(y^*_{1:j}; \theta) + log \sum_{\textcolor{red}{y' \epsilon \mathcal{B}_j}} \exp\{u(y'_{1:j}; \theta)\}$$

$\textcolor{red}{\mathcal{B}_j}$ contains all paths in the beam at step $j$, together with the oracle path prefix $y^*_{1:j}$

▶ **Decoding: beam search**

a. Kai Hu, Zhijian Ou, et al. Neural CRF Transducers for Sequence Labeling. ICASSP, 2019.
b. Andor, Alberti, et al., "Globally Normalized Transition-Based Neural Networks", ACL 2016.

# Experiment and Conclusion

| Model \ Task | POS (Accuracy) | Chunking (F1 score) | English NER (F1 score) | Dutch NER (F1 score) | Globally normalized | Long-range dependencies |
|---|---|---|---|---|---|---|
| Linear-chain CRF | 97.52 | 95.01 | 91.11 | 81.53 | √ | ✕ |
| RNN Transducer | 97.50 | 95.02 | 91.02 | 81.59 | ✕ | √ |
| CRF Transducer | 97.52 | 95.14 | 91.40 | 81.84 | √ | √ |

Experiment results show that CRF transducers achieve consistent improvements over linear-chain CRFs and RNN transducers across four sequence labeling tasks, and obtain state-of-theart results.

Reproducible code is at https://github.com/thu-spmi/SPMISeq

Kai Hu, Zhijian Ou, et al. Neural CRF Transducers for Sequence Labeling. ICASSP, 2019.

# Content

I. **Basics for EBMs (45 min)**
  1. Probabilistic graphical modeling (PGM) framework and EBM model examples
  2. Learning EBMs by Monte Carlo methods
  3. Learning EBMs by noise-contrastive estimation (NCE)

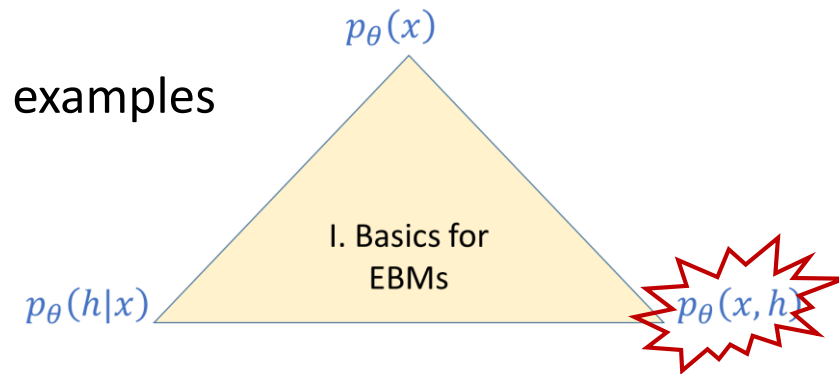II. **EBMs for language modeling (45 min)**
  1. Trans-dimensional random field (TRF) LMs for speech recognition
  2. Residual energy-based models for text generation
  3. Electric: an energy-based cloze model for representation learning over text

III. **EBMs for speech recognition and natural language labeling (45 min)**
  1. CRFs as conditional EBMs
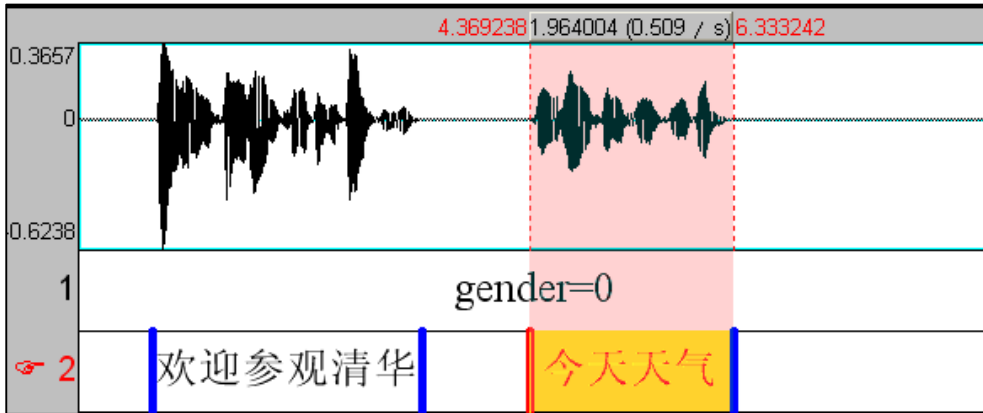  2. CRFs for speech recognition
  3. CRFs for sequence labeling in NLP

IV. **EBMs for semi-supervised natural language labeling (45 min)**
  1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
  2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
  3. Comparison of joint-training and pre-training for semi-supervised learning via EBMs

$p_\theta(x)$

I. Basics for EBMs

$p_\theta(h|x)$

$p_\theta(x, h)$

105

# Supervised learning from Labeled data $\{(x_j, y_j), j = 1, 2, \cdots\}$
## Tremendous Success!



**Speech Recognition**

ATIS UTTERANCE EXAMPLE IOB REPRESENTATION

| Sentence | show | flights | from | Boston | To | New | York | today |
|---|---|---|---|---|---|---|---|---|
| Slots/Concepts | O | O | O | B-dept | O | B-arr | I-arr | B-date |
| Named Entity | O | O | O | B-city | O | B-city | I-city | O |
| Intent | | | | Find_Flight | | | | |
| Domain | | | | Airline Travel | | | | |

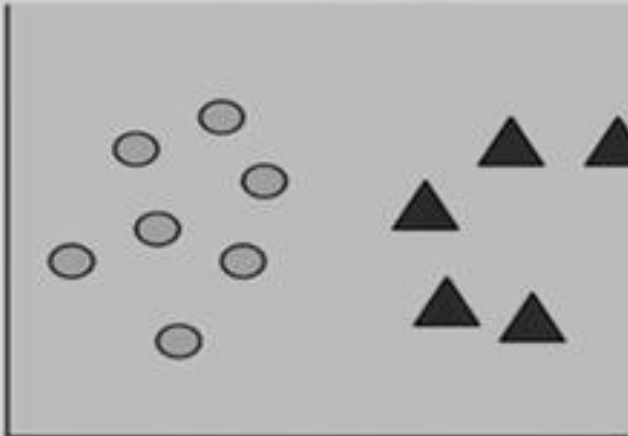**Intent Detection, Slot Filling, Named Entity Recognition**



**Object Detection and Tracking**



**Syntactic Parsing**

# Semi-supervised learning (SSL)



Labeled Data
(a)

Labeled and Unlabeled Data
(b)

$p_\theta(y|x)$

Classification plane

Supervised Learning
(c)

Semi-Supervised Learning
(d)

**The key to designing SSL methods is:**
**How to effectively exploit <span style="color:#4472C4">the information contained in the unlabeled data $\{x\}$</span>,**
**which can provide**
**priors/regularizations/inductive biases**
**for finding the posterior $p_\theta(y|x)$.**

Xiaojin Zhu, "Semi-supervised learning literature survey," Technical report, University of Wisconsin-Madison, 2006.

# SSL methods (for using DNNs)

- Recent SSL methods with DNNs can be distinguished by the priors they adopt, and, can be divided into two classes.

  - **Generative SSL**
  - **Discriminative SSL**: The outputs from the discriminative classifier are smooth with respect to local and random perturbations of the inputs [1-5].

[1] Takeru Miyato, et al, "Virtual adversarial training: a regularization method for supervised and semi-supervised learning," TPAMI, 2018.
[2] Samuli Laine and Timo Aila, "Temporal ensembling for semisupervised learning," ICLR, 2017.
[3] Antti Tarvainen and Harri Valpola, "Mean teachers are better role models: Weight-averaged consistency targets improve semi-supervised deep learning results," NIPS, 2017.
[4] Kihyuk Sohn, David Berthelot, Chun-Liang Li, and et al, "FixMatch: Simplifying semi-supervised learning with consistency and confidence," arXiv:2001.07685, 2020.
[5] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, "A simple framework for contrastive learning of visual representations," arXiv:2002.05709, 2020.

# Discriminative SSL

- Recent SSL methods with DNNs can be distinguished by the priors they adopt, and, can be divided into two classes.

  - **Generative SSL**
  - **Discriminative SSL**: The outputs from the discriminative classifier are smooth with respect to local and random perturbations of the inputs.

  ☹ heavily rely on domain-specific data augmentations, which are tuned intensively for images leading to impressive performance in some image domains
  ☹ less successful for other domains where these augmentations are less effective (e.g., medical images and text). For instance, random input perturbations are more difficult to apply to discrete data like text [6].

[6] Kevin Clark, Minh-Thang Luong, Christopher D Manning, and Quoc Le, "Semi-supervised sequence modeling with cross-view training," in EMNLP, 2018

# Generative SSL - Basics

- Exploit unsupervised learning of generative models over unlabeled data, blend unsupervised learning and supervised learning.

☺ inherently not require data augmentations and generally can be applied to a wider range of domains.
☺ make fewer domain-specific assumptions and tend to be domain-agnostic.
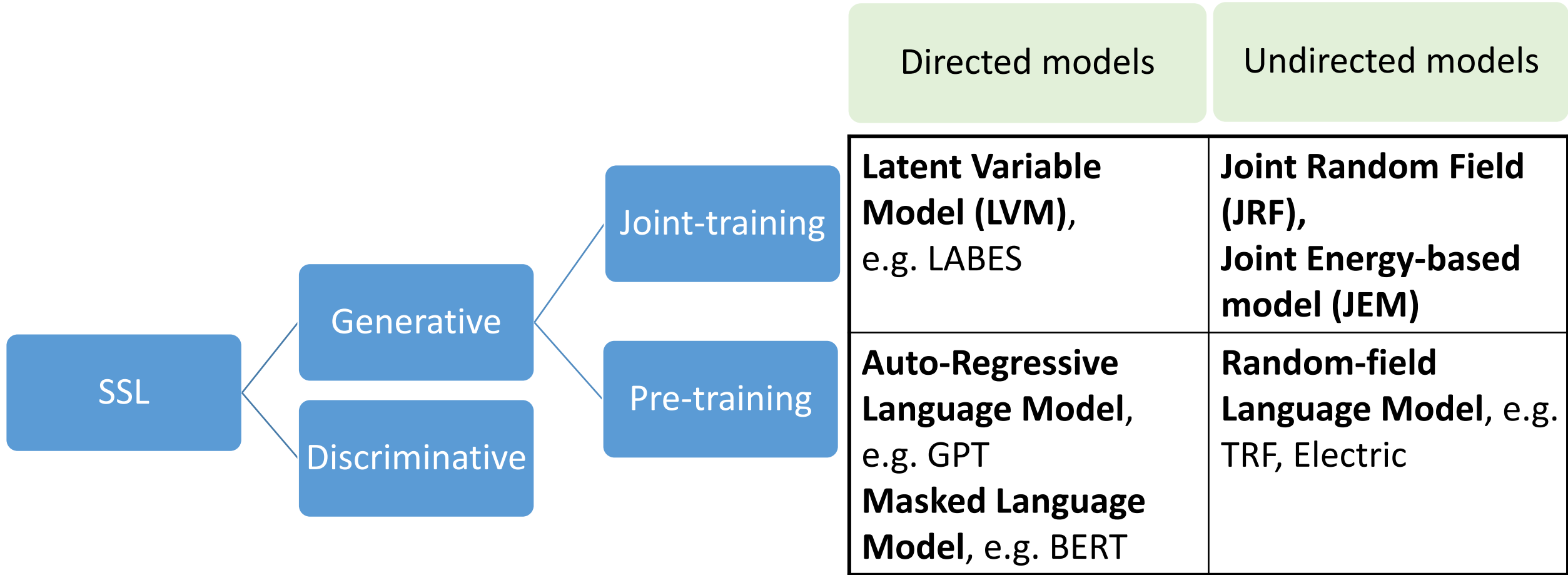
# Generative SSL - Two Different Approaches

- **Joint-training**
  - A joint model of p(x,y) is defined.
  - When we have label y, we maximize p(y|x) (the supervised objective), and when the label is unobserved, we marginalize it out and maximize p(x) (the unsupervised objective).
  - Semi-supervised learning over a mix of labeled and unlabeled data is formulated as maximizing the (weighted) sum of log p(y|x) and log p(x).

- **Pre-training**
  - Only defines p(x) without y.
  - Perform unsupervised representation learning (called pre-training) on unlabeled data, followed by supervised training (called fine-tuning) on labeled data.
  - This manner of pre-training followed by fine-tuning has received increasing application in Natural Language Processing.

# There are many open questions in designing semi-supervised methods for particular tasks !



| | Directed models | Undirected models |
|---|---|---|
| Joint-training | **Latent Variable Model (LVM)**, e.g. LABES | **Joint Random Field (JRF)**, **Joint Energy-based model (JEM)** |
| Pre-training | **Auto-Regressive Language Model**, e.g. GPT **Masked Language Model**, e.g. BERT | **Random-field Language Model**, e.g. TRF, Electric |

[LABES] Y. Zhang, Z. Ou, et al. A Probabilistic End-To-End Task-Oriented Dialog Model with Latent Belief States towards Semi-Supervised Learning. EMNLP, 2020.
[JRF] Y. Song, Z. Ou, et al. Upgrading CRFs to JRFs and its benefits to sequence modeling and labeling. ICASSP, 2020.
[JEM] S. Zhao, J.H. Jacobsen, et al. Joint energy-based models for semi-supervised classification. ICML Workshop on Uncertainty and Robustness in Deep Learning, 2020.
[TRF] B. Wang, Z. Ou. Improved training of neural trans-dimensional random field language models with dynamic noise-contrastive estimation. SLT, 2018.
[Electric] K. Clark, M.T. Luong, et al. Pre-Training Transformers as Energy-Based Cloze Models. EMNLP, 2020.

**Table 1.** Applications of EBMs across different domains: comparison and connection (See text for details).

| | Image classification | Natural language labeling |
|---|---|---|
| Observation | $x \in \mathbb{R}^D$ <br> continuous, fixed-dimensional | $x \in \bigcup_l \mathbb{V}^l$ <br> discrete, sequence |
| Label | $y \in \{1, 2, \cdots, K\}$ | $y \in \bigcup_l \{1, 2, \cdots, K\}^l$ |
| Pre-training | ① $\quad u_\theta(x) = w^T h$ | ③ $\quad u_\theta(x)$ in Eq.(3) |
| Joint-training | ② $u_\theta(x, y) = \Psi_\theta(x)[y]$ | ④ $\quad u_\theta(x, y)$ in Eq.(6) |

Yunfu Song, Huahuan Zheng, Zhijian Ou. An empirical comparison of joint-training and pre-training for domain-agnostic semi-supervised learning via energy-based models. MLSP, 2021.

# Content

$p_\theta(x)$

I. Basics for EBMs

$p_\theta(h|x)$     $p_\theta(x, h)$

# ① Pre-training of an EBM for semi-supervised image classification

1) **Pre-training**: estimate $p_\theta(x)$ over unlabeled images

$$p_\theta(x) = \frac{1}{Z(\theta)} exp[u_\theta(x)]$$

Use a feedforward NN to implement $u_\theta(x) \colon \mathbb{R}^d \to \mathbb{R}$

which, in the final layer, calculates $u_\theta(x) = w^T h$ via a linear layer.

$u_\theta(x)$ $K$

$h$

Potential function

$x$

2) **Fine-tuning**: throw $w$ and fed $h$ into an new linear output layer, followed by softmax$(Wh)$, to predict $y \in \{1, \cdots, K\}$, where $W \in \mathbb{R}^{K \times H}$

Yunfu Song, Zhijian Ou. Learning Neural Random Fields with Inclusive Auxiliary Generators. arXiv:1806.00271, 2018.

# ② Joint-training of an EBM for semi-supervised image classification

$u_\theta(x, y)$

- **Joint modeling** of observation $x \in \mathbb{R}^d$ and class label $y \in \{1, \cdots, K\}$:

$$p_\theta(x, y) = \frac{1}{Z(\theta)} exp[u_\theta(x, y)]$$

Potential function

- Consider a NN $\Psi_\theta(x): \mathbb{R}^d \to \mathbb{R}^K$ and define:

$$u_\theta(x, y) = \Psi_\theta(x)[y]$$

- **Classifier**: $p_\theta(y|x) = \frac{p_\theta(x,y)}{p_\theta(x)} = \frac{exp[u_\theta(x,y)]}{\sum_y exp[u_\theta(x,y)]}$, like a $K$-class logistic regression

$x$

**Marginal density**: $p_\theta(x) = \frac{1}{Z(\theta)} exp[u_\theta(x)]$, where $u_\theta(x) \triangleq log \sum_y exp[u_\theta(x, y)]$

$$\begin{cases} \min_\theta KL[\tilde{p}(\tilde{x})||p_\theta(\tilde{x})] -\alpha \sum_{(\tilde{x},\tilde{y})\sim\mathcal{L}} log p_\theta(\tilde{y}|\tilde{x}) \\ \min_\phi KL[p_\theta(x)||q_\phi(x)] \end{cases}$$

Yunfu Song, Zhijian Ou. Learning Neural Random Fields with Inclusive Auxiliary Generators. arXiv:1806.00271, 2018.

# Learning Neural Random Fields with Inclusive Auxiliary Generators

Yunfu Song, Zhijian Ou

In this paper we develop Neural Random Field learning with Inclusive-divergence minimized Auxiliary Generators (NRF-IAG), which is under-appreciated in the literature. The contributions are two-fold. First, we rigorously apply the stochastic approximation algorithm to solve the joint optimization and provide theoretical justification. The new approach of learning NRF-IAG achieves superior unsupervised learning performance competitive with state-of-the-art deep generative models (DGMs) in terms of sample generation quality. Second, semi-supervised learning (SSL) with NRF-IAG gives rise to strong classification results comparable to state-of-art DGM-based SSL methods, and simultaneously achieves superior generation. This is in contrast to the conflict of good classification and good generation, as observed in GAN-based SSL.



Published as a conference paper at ICLR 2020

## YOUR CLASSIFIER IS SECRETLY AN ENERGY BASED MODEL AND YOU SHOULD TREAT IT LIKE ONE

**Will Grathwohl**
University of Toronto & Vector Institute
Google Research
wgrathwohl@cs.toronto.edu

**Kuan-Chieh Wang*& Jörn-Henrik Jacobsen***
University of Toronto & Vector Institute
wangkua1@cs.toronto.edu
j.jacobsen@vectorinstitute.ai

**David Duvenaud**
University of Toronto & Vector Institute
duvenaud@cs.toronto.edu

**Kevin Swersky & Mohammad Norouzi**
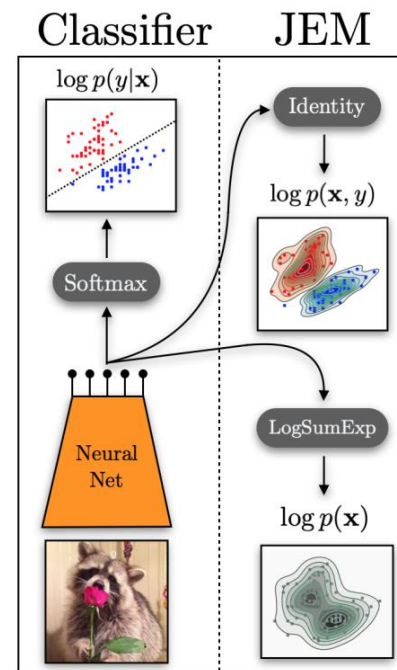Google Research
{kswersky, mnorouzi}@google.com



Figure 1: Visualization of our method, JEM, which defines a joint EBM from classifier architectures.

# Inclusive-NRF algo. for learning from continuous data, e.g., Images.
## simultaneously training a random field and a generator.



$u_\theta(x)$

$h$

Potential function

Update $(x, h)$ Update

Generator

$x$

Revise

$(x', h')$ Propose $g_\phi(h)$ $\oplus$ $\epsilon$

$x$

$$\begin{cases} \min_\theta KL[\tilde{p}(\tilde{x})||p_\theta(\tilde{x})] \\ \min_\phi KL[p_\theta(x)||q_\phi(x)] \end{cases} \Rightarrow \begin{cases} \nabla_\theta = E_{\tilde{p}(\tilde{x})}[\nabla_\theta log p_\theta(\tilde{x})] = E_{\tilde{p}(\tilde{x})}[\nabla_\theta u_\theta(\tilde{x})] - E_{p_\theta(x)}[\nabla_\theta u_\theta(x)] \\ \nabla_\phi = E_{p_\theta(x)}[\nabla_\phi log q_\phi(x)] = E_{p_\theta(x)q_\phi(h|x)}[\nabla_\phi log q_\phi(x, h)] \end{cases}$$

Yunfu Song, Zhijian Ou. Learning Neural Random Fields with Inclusive Auxiliary Generators. arXiv:1806.00271, 2018.

# Content

# ③ Pre-training of an EBM for semi-supervised natural language labeling

1) **Pre-training**: estimate $p_\theta(x)$ over unlabeled sentences $x = (x_1, \cdots, x_l)$

$$p_\theta(x) = \frac{1}{Z(\theta)} exp[u_\theta(x)]$$

Use a B-LSTM to implement $u_\theta(x) : \mathbb{V}^l \to \mathbb{R}$

$$u_\theta(x) = \sum_{i=1}^{l-1} h_{f,i}^T e_{i+1} + \sum_{i=2}^{l} h_{b,i}^T e_{i-1}$$



2) **Fine-tuning**: we add a CRF on top of the extracted representations $\{(h_{f,i}, h_{b,i}), i = 1, \cdots, l\}$ to predict label sequence $y = (y_1, \cdots, y_l)$.

Bin Wang, Zhijian Ou. Improved training of neural trans-dimensional random field language models with dynamic noise-contrastive estimation. IEEE Workshop on Spoken Language Technology (SLT), 2018.

# ④ Joint-training of an EBM for semi-supervised natural language labeling

- JRF: Define a joint distribution over $x = (x_1, \cdots, x_l)$ and $y = (y_1, \cdots, y_l)$

$$p_\theta(l, x^l, y^l) = \pi_l p_\theta(x^l, y^l; l) = \frac{\pi_l}{Z_\theta(l)} \exp\left(u_\theta(x^l, y^l)\right)$$

- Consider a NN $\Psi_\theta(x): \mathbb{V}^l \rightarrow \mathbb{R}^{l \times K}$ and define:

$$u_\theta(x, y) = \sum_{i=1}^{l} \Psi_\theta(x)[i, y_i] + \sum_{i=1}^{l} A[y_{i-1}, y_i]$$

- From JRF we have:

$$p_\theta(y^l | x^l) = \frac{1}{\sum_{y^l} \exp\left(u_\theta(x^l, y^l)\right)} \exp\left(u_\theta(x^l, y^l)\right)$$

which is a CRF

- From JRF we have

$$p_\theta(l, x^l) = \frac{1}{Z_\theta}$$

$$= \frac{\pi_l}{Z_\theta(l)} \exp\left(\right.$$

where $u_\theta(x^l) = \log \sum$

which is a trans-dim

**Recently: neural CRFs**

Use NN to extract features
$$\mathrm{LSTM}(x_{1:l}): x_{1:l} \rightarrow h_{1:l}$$

- Node potential, calculated via a linear layer

$$\phi_t(y_t = k, x) = w_k^{\mathrm{T}} h_t \triangleq \phi_t^k$$

$w_k$ is the weight vector for label $k$

- Edge potential, mostly implemented as a matrix $A$

Yunfu Song, Zhijian Ou, et al. Upgrading CRFs to JRFs and its benefits to sequence mode

# Upgrading CRFs to Joint random fields (JRFs) for sequential data



CRF
$$p_\theta(y^l \mid x^l)$$

TRF
$$p_\theta(l, x^l)$$

Supervised Learning

Unsupervised Learning

JRF
$$p_\theta(l, x^l, y^l)$$

Labeled Data

Unlabeled Data

Edge Potentials

Node Potentials

$y_1 \quad\quad y_2 \quad\quad y_3$

$o_1 \quad\quad o_2 \quad\quad o_3$

Bi-LSTM

$h_1 \quad\quad h_2 \quad\quad h_3$

$x_1 \quad\quad x_2 \quad\quad x_3$

123

# Dynamic NCE algo. for learning from discrete data, e.g., texts.
## Simultaneously train a random field and a generator.

- The target RF model $\quad p_\theta(x) = \dfrac{1}{Z(\theta)} e^{u_\theta(x)}$

- Treat $\log Z(\theta)$ as a parameter $\zeta$ and rewrite $\quad p_{\theta,\zeta}(x) \propto e^{u_\theta(x) - \zeta}$

- Introduce a **noise distribution** $q(x)$, and consider a binary classification

$x \sim p_0(x)$ → Binary discriminator → $C = 0/1$
$x \sim q(x)$

$$P(C = 0|x) = \frac{p_{\theta,\zeta}(x)}{p_{\theta,\zeta}(x) + \nu q(x)}, where \ \nu = \frac{P(C = 1)}{P(C = 0)}$$

$$P(C = 1|x) = 1 - P(C = 0|x)$$

- Noise Contrastive Estimation (NCE):

$$\max_{\theta,\zeta} E_{x \sim p_0(x)}[\log P(C = 0|x)] + E_{x \sim q(x)}[\log P(C = 1|x)]$$

- Consistency: $p_\theta \to p_0$ (oracle), under infinite amount of data and infinite capacity of $p_\theta$.
- Reliable NCE needs a large $\nu \approx 20$; Dynamic-NCE works well with $\nu = 1$.

# Content

I.    Basics for EBMs (45 min)
   1. Probabilistic graphical modeling (PGM) framework and EBM model examples  (classic & modern)
   2. Learning EBMs by Monte Carlo methods
   3. Learning EBMs by noise-contrastive estimation (NCE)

II. EBMs for language modeling (45 min)
   1. Trans-dimensional random field (TRF) LMs for speech recognition
   2. Residual energy-based models for text generation
   3. Electric: an energy-based cloze model for representation learning over text

III. EBMs for speech recognition and natural language labeling (45 min)
   1. CRFs as conditional EBMs
   2. CRFs for speech recognition
   3. CRFs for sequence labeling in NLP

IV. EBMs for semi-supervised natural language labeling (45 min)
   1. Upgrading EBMs to Joint EBMs (JEMs) for fixed-dimensional data
   2. Upgrading CRFs to Joint random fields (JRFs) for sequential data
➡ 3. Comparison of joint-training and pre-training for semi-supervised learning via EBMs

**EBM models can be very flexibly defined for SSL, by either of joint-training and pre-training.**

**… previously known in the literature†, but it is unclear which is better when evaluated in a common experimental setup.**

**To the best of our knowledge, this paper‡ is the first to systematically compare joint-training and pre-training for EBM-based for SSL, across domains (image classification and natural language labeling).**

† EBM based SSL results have been reported across different data modalities (images, natural languages, an protein structure prediction and year prediction from the UCI dataset repository) [12,13,14].
‡ Yunfu Song, Huahuan Zheng, Zhijian Ou. An empirical comparison of joint-training and pre-training for domain-agnostic semi-supervised learning via energy-based models. MLSP, 2021.

**Table 1.** Applications of EBMs across different domains: comparison and connection (See text for details).

| | Image classification | Natural language labeling |
|---|---|---|
| Observation | $x \in \mathbb{R}^D$ <br> continuous, fixed-dimensional | $x \in \bigcup_l \mathbb{V}^l$ <br> discrete, sequence |
| Label | $y \in \{1, 2, \cdots, K\}$ | $y \in \bigcup_l \{1, 2, \cdots, K\}^l$ |
| Pre-training | ① $\quad u_\theta(x) = w^T h$ | ③ $\quad u_\theta(x)$ in Eq.(3) |
| Joint-training | ② $\quad u_\theta(x, y) = \Psi_\theta(x)[y]$ | ④ $\quad u_\theta(x, y)$ in Eq.(6) |

Pre-training aims to learn representations that may be useful for multiple downstream tasks, and any information about the labels is not utilized until the fine-tuning stage.

Yunfu Song, Huahuan Zheng, Zhijian Ou. An empirical comparison of joint-training and pre-training for domain-agnostic semi-supervised learning via energy-based models. MLSP, 2021.

**Table 2.** SSL for image classification over CIFAR-10 with 4,000 labels. The upper/lower blocks show generative/discriminative SSL methods respectively. The means and standard deviations are calculated over ten independent runs with randomly sampled labels.

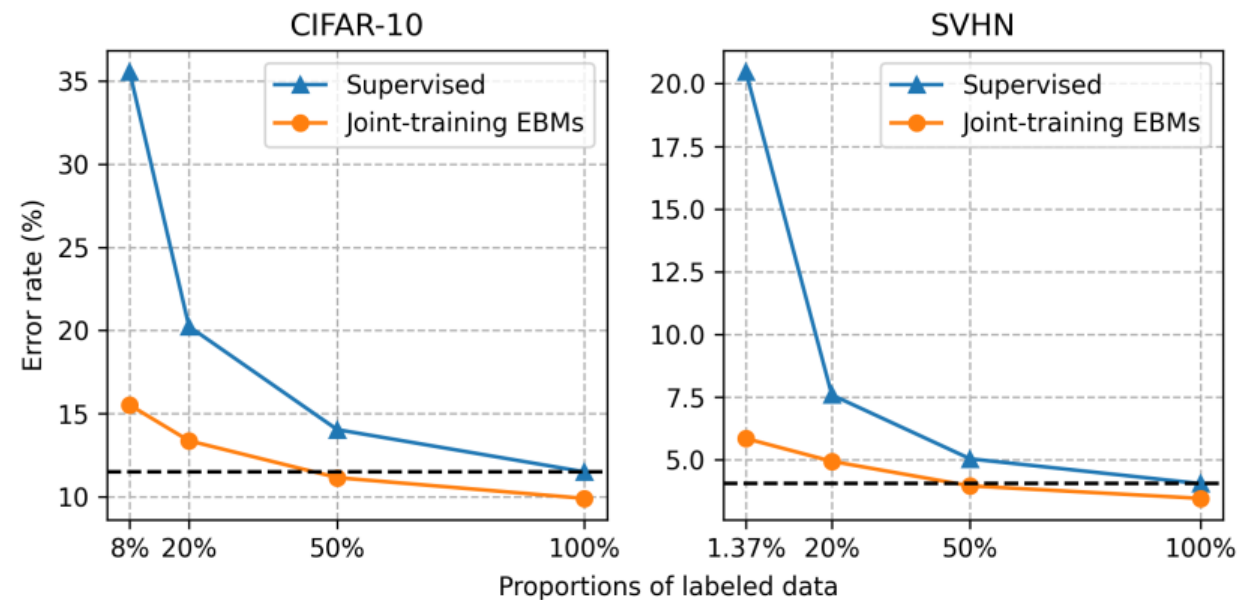| Methods | error (%) |
|---|---|
| CatGAN [30] | 19.58±0.46 |
| Ladder network [31] | 20.40±0.47 |
| Improved-GAN [32] | 18.63±2.32 |
| BadGAN [33] | 14.41±0.30 |
| Sobolev-GAN [34] | 15.77±0.19 |
| **Supervised baseline** | 25.72±0.44 |
| **Pre-training+fine-tuning EBM** | 21.40±0.38 |
| **Joint-training EBM** | 15.12±0.36 |
| Results below this line cannot be directly compared to those above. | |
| VAT small [1] | 14.87 |
| Temporal Ensembling [2] | 12.16±0.31 |
| Mean Teacher [3] | 12.31±0.28 |



**Fig. 1.** Error rates of supervised baseline and joint-training EBMs as the amount of labels varies on SVHN and CIFAR-10 datasets. The dash line is the supervised result trained with 100% labeled data.

Joint-training EBMs outperform pre-training+fine-tuning EBMs by a large margin in this task.

Can reduce 50% of labels without losing performance.

**Table 3.** Natural language labeling results. The evaluation metric is accuracy for POS and $F_1$ for chunking and NER. "Labeled" denotes the amount of labels in terms of the proportions w.r.t. the full set of labels. "U/L" denotes the ratio between the amount of unlabeled and labeled data. "U/L=0" denotes the supervised baseline. "pre." and "joint" denote the results by pre-training+fine-tuning EBMs and joint-training EBMs, respectively.

| Labeled | U/L | POS tagging pre. | POS tagging joint | Chunking pre. | Chunking joint | NER pre. | NER joint |
|---|---|---|---|---|---|---|---|
| 2% | 0 | 95.57 | | 78.73 | | 78.19 | |
| | 50 | 95.72 | 95.92 | 81.62 | 82.24 | 76.74 | 77.61 |
| | 250 | 95.96 | 96.13 | 82.10 | 82.26 | 78.49 | 78.51 |
| | 500 | 96.08 | 96.24 | 83.10 | 83.05 | 79.47 | 79.17 |
| 10% | 0 | 96.81 | | 90.06 | | 86.93 | |
| | 50 | 96.87 | 96.99 | 91.60 | 91.85 | 86.37 | 87.05 |
| | 250 | 96.88 | 97.00 | 91.09 | 91.93 | 86.86 | 86.77 |
| | 500 | 96.92 | 97.08 | 91.93 | 92.23 | 87.57 | 87.06 |
| 100% | 0 | 97.41 | | 94.77 | | 90.74 | |
| | 50 | 97.40 | 97.49 | 95.05 | 95.31 | 91.24 | 91.34 |
| | 250 | 97.45 | 97.54 | 95.12 | 95.48 | 91.19 | 91.51 |
| | 500 | 97.46 | 97.57 | 95.19 | 95.50 | 91.30 | 91.52 |

**Table 4.** Relative improvements by joint-training EBMs compared to the supervised baseline (abbreviated as sup.) and pretraining+fine-tuning EBMs respectively. Refer to Table 3 for notations.

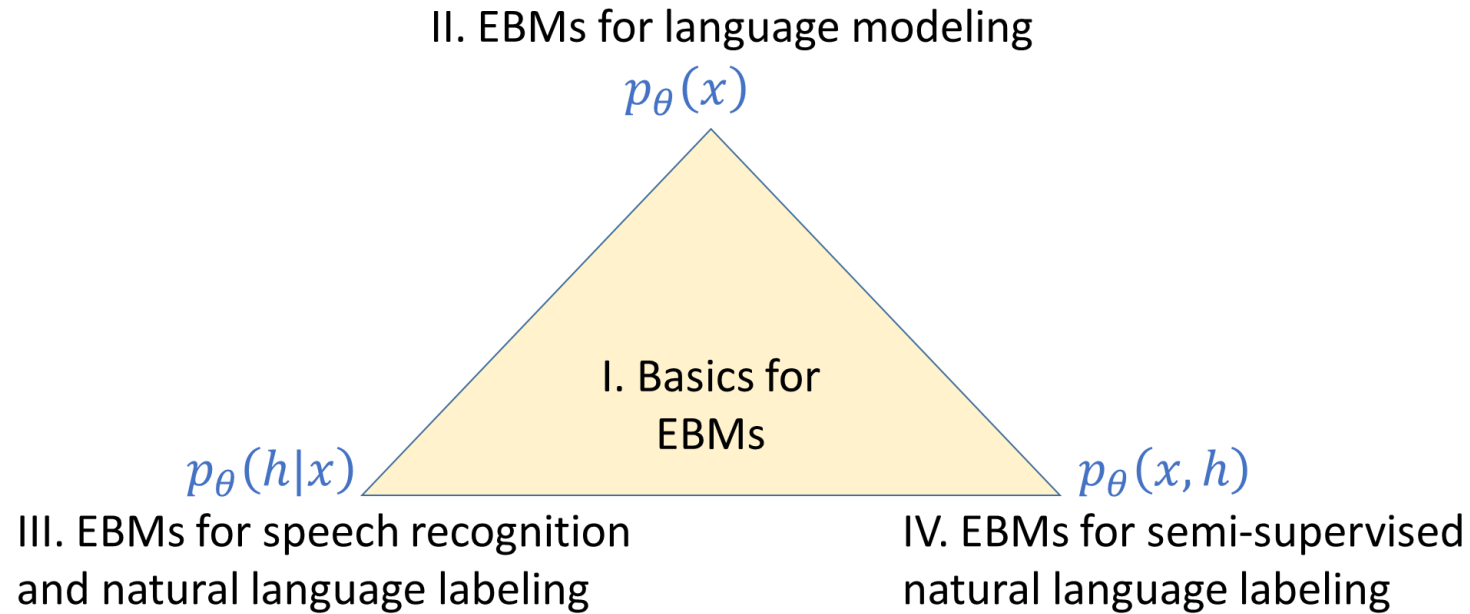| Labeled | U/L | joint over sup. POS | joint over sup. Chunking | joint over sup. NER | joint over pre. POS | joint over pre. Chunking | joint over pre. NER |
|---|---|---|---|---|---|---|---|
| 2% | 50 | 7.9 | 16.5 | -2.7 | 4.7 | 3.4 | 3.7 |
| | 250 | 12.6 | 16.6 | 1.5 | 4.2 | 0.9 | 0.1 |
| | 500 | 15.1 | 20.3 | 4.5 | 4.1 | -0.3 | -1.5 |
| 10% | 50 | 5.6 | 18.0 | 0.9 | 3.8 | 3.0 | 5.0 |
| | 250 | 6.0 | 18.3 | -1.2 | 3.8 | 9.4 | -0.7 |
| | 500 | 8.5 | 21.8 | 1.0 | 5.2 | 3.7 | -4.1 |
| 100% | 50 | 3.1 | 10.3 | 6.5 | 3.5 | 5.3 | 1.1 |
| | 250 | 5.0 | 13.6 | 8.3 | 3.5 | 7.4 | 3.6 |
| | 500 | 6.2 | 14.0 | 8.4 | 4.3 | 6.4 | 2.5 |

- Joint-training EBMs outperform pre-training EBMs in 23 out of the 27 settings marginally but nearly consistently.
- A possible explanation is that pretraining is not aware of the labels for the targeted task and is thus weakened for representation learning.

# Section Conclusion

- We systematically evaluate and compare joint-training and pre-training for EBM-based domain-agnostic SSL, through a suite of experiments across a variety of domains such as image classification and natural language labeling.

- **Joint-training EBMs outperform pre-training EBMs marginally but nearly consistently.**

  ▶ Presumably, this is because that the optimization of joint-training is directly related to the targeted task, but pre-training is not aware of the labels for the targeted task.

- This new finding would be helpful for future work to further explore better methods to leverage unlabeled data.

Reproducible code is at https://github.com/thu-spmi/semi-EBM

# Summary



II. EBMs for language modeling
$p_\theta(x)$

I. Basics for EBMs

$p_\theta(h|x)$

III. EBMs for speech recognition and natural language labeling

$p_\theta(x, h)$

IV. EBMs for semi-supervised natural language labeling

**Take-home messages for EBMs/UBMs/RFs**

1. Flexibility in modeling
2. Computation efficiency in inference
3. Overcome label bias and exposure bias suffered by locally-normalized models
4. Joint EBMs for generative semi-supervised learning
5. Difficult in model training

We are making progress, and there are many interesting open questions...

# Main References

1. D. Koller and N. Friedman. Probabilistic graphical models: principles and techniques. MIT press, 2009.

2. Eric Fosler-Lussier, et al. Conditional random fields in speech, audio, and language processing. Proceedings of the IEEE, 2013.

3. Zhijian Ou. A Review of Learning with Deep Generative Models from Perspective of Graphical Modeling. arXiv:1808.01630, 2018.

4. Bin Wang, Zhijian Ou, Zhiqiang Tan. Trans-dimensional Random Fields for Language Modeling. Annual Meeting of the Association for Computational Linguistics (ACL Long Paper), 2015.

5. Bin Wang, Zhijian Ou, Zhiqiang Tan. Learning Trans-dimensional Random Fields with Applications to Language Modeling. IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), 2018, vol.40, no.4, pp.876-890. https://github.com/thu-spmi/SPMILM

6. Bin Wang, Zhijian Ou. Language modeling with neural trans-dimensional random fields. ASRU, 2017.

7. Bin Wang, Zhijian Ou. Learning neural trans-dimensional random field language models with noise-contrastive estimation. ICASSP, 2018.

8. Bin Wang, Zhijian Ou. Improved training of neural trans-dimensional random field language models with dynamic noise-contrastive estimation. SLT, 2018.

9. Silin Gao, Zhijian Ou, et al. Integrating discrete and neural features via mixed-feature trans-dimensional random field language models. ICASSP (Oral), 2020.

10. Yunfu Song, Zhijian Ou. Learning Neural Random Fields with Inclusive Auxiliary Generators. arXiv:1806.00271, 2018. https://github.com/thu-spmi/Inclusive-NRF

11. Hongyu Xiang, Zhijian Ou. CRF-based Single-stage Acoustic Modeling with CTC Topology. ICASSP (Oral), 2019. https://github.com/thu-spmi/CAT

12. Kai Hu, Zhijian Ou, Min Hu, Junlan Feng. Neural CRF Transducers for Sequence Labeling. ICASSP, 2019. https://github.com/thu-spmi/SPMISeq

13. Yunfu Song, Zhijian Ou, Zitao Liu, Songfan Yang. Upgrading CRFs to JRFs and its benefits to sequence modeling and labeling. ICASSP, 2020. https://github.com/thu-spmi/semi-EBM

14. Yunfu Song, Huahuan Zheng, Zhijian Ou. An empirical comparison of joint-training and pre-training for domain-agnostic semi-supervised learning via energy-based models. IEEE Machine Learning for Signal Processing Workshop (MLSP), 2021. https://github.com/thu-spmi/semi-EBM

15. Will Grathwohl, Kuan-Chieh Wang, Jörn-Henrik Jacobsen, David Duvenaud, Mohammad Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should treat it like one. ICLR 2020.

16. Yuntian Deng, Anton Bakhtin, Myle Ott, Arthur Szlam, and Marc'Aurelio Ranzato. Residual energy-based models for text generation, ICLR 2020.

17. Kevin Clark, Minh-Thang Luong, Quoc V. Le, and Christopher D. Manning. Pre-training transformers as energy-based cloze models, EMNLP 2020.

# Thanks for your attention !

Thanks to my collaborators and students :

Zhiqiang Tan, Bin Wang, Hongyu Xiang, Yunfu Song, Kai Hu,

Keyu An, Huahuan Zheng, Silin Gao, Wenjie Peng