

Fast Online Training with Frequency-Adaptive Learning Rates for Chinese Word Segmentation and New Word Detection

A decorative graphic consisting of three overlapping chevrons pointing to the right. The leftmost chevron is green, the middle one is blue, and the rightmost one is yellow. The authors' names are written in black text on the yellow chevron.

Xu SUN, Houfeng WANG, Wenjie LI

The Hong Kong Polytechnic University
Peking University

Outline

◆ Introduction

◆ Method

- Joint modeling: word segmentation + new word detection
- New features

◆ A new online training method

- Feature-frequency adaptive online training
- Finish the training in 10 passes

◆ Experiments

◆ Conclusions

Introduction

◆ 3 proposals in this work

Introduction

◆ 3 proposals in this work

◆ 1) Joint modeling

- word segmentation + new word detection

◆ 2) New features

- high dimensional edge features on CRFs

Introduction

◆ 3 proposals in this work

◆ 1) Joint modeling

- word segmentation + new word detection

◆ 2) New features

- high dimensional edge features on CRFs

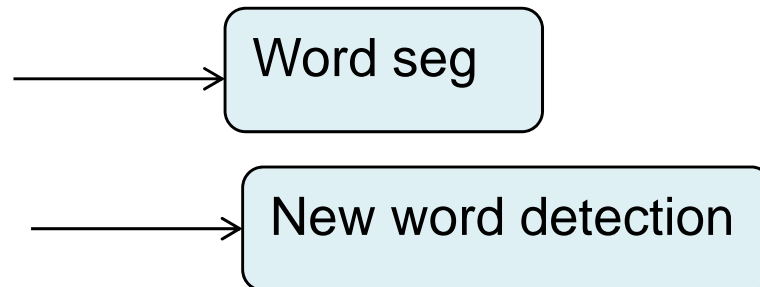
◆ 3) **Very Fast online training **major proposal****

- finish the training in 10 passes

1) Joint modeling

◆ Prior work on word seg & new word detection

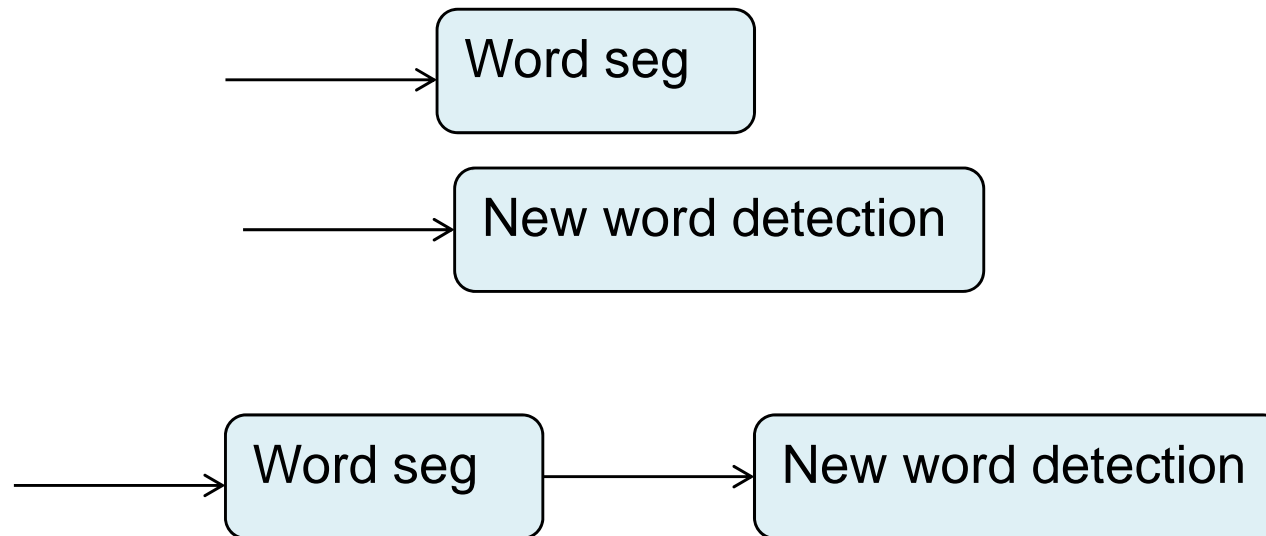
- Separate: two tasks are independent



1) Joint modeling

◆ Prior work on word seg & new word detection

- Separate: two tasks are independent
- Pipeline: first word seg, then new word detection



1) Joint modeling

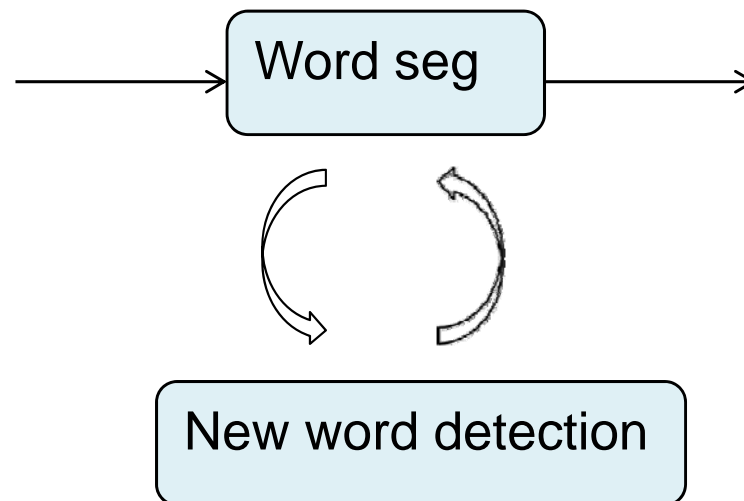
◆ Our work

- **Joint modeling**: word seg + new word detection
- For simplicity, we use a single CRF model

1) Joint modeling

◆ Our work

- **Joint modeling**: word seg + new word detection
- For simplicity, we use a single CRF model



2) New features

- ◆ **Traditional features are character ngrams**

2) New features

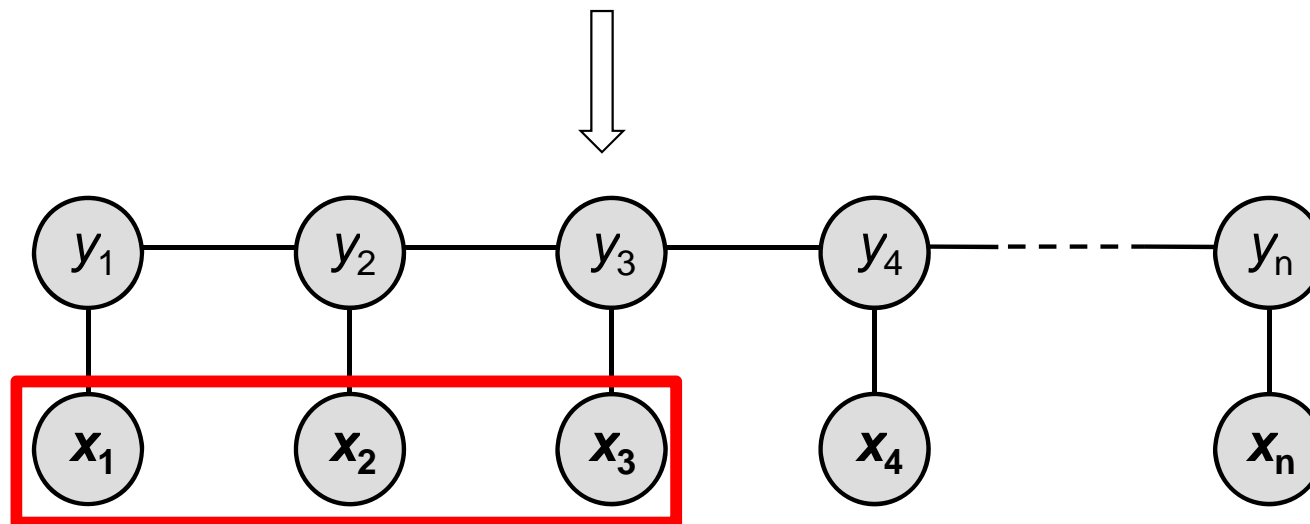
- ◆ **Traditional features are character ngrams**
- ◆ **New features: word-based ngrams**
 - Word lexicon is collected from training data
 - Word unigram features
 - Word bigram features

2) New features

◆ Traditional features are character ngrams

◆ New features: word-based ngrams

- Word lexicon is collected from training data
- Word unigram features
- Word bigram features



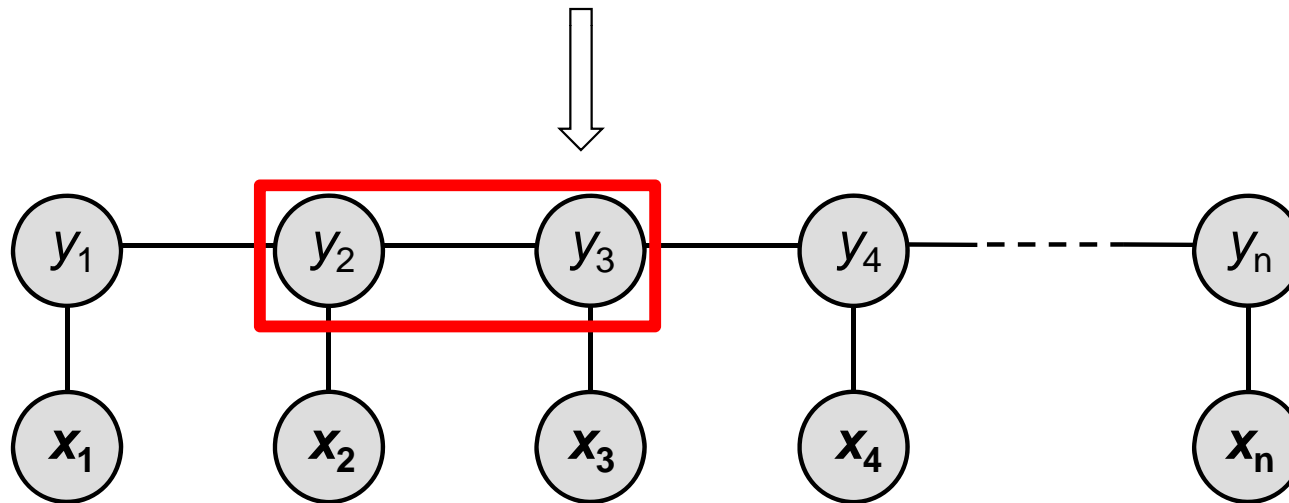
2) New edge (transition) features

- ◆ **Traditional edge features avoid observations**
 - to avoid feature explosion

2) New edge (transition) features

◆ Traditional edge features avoid observations

- to avoid feature explosion



2) New edge (transition) features

◆ Traditional edge features avoid observations

- to avoid feature explosion

◆ **Proposal:** New edge features involving observations

- Tractable feature explosion: #feature increased 9 times
- → 10 millions features

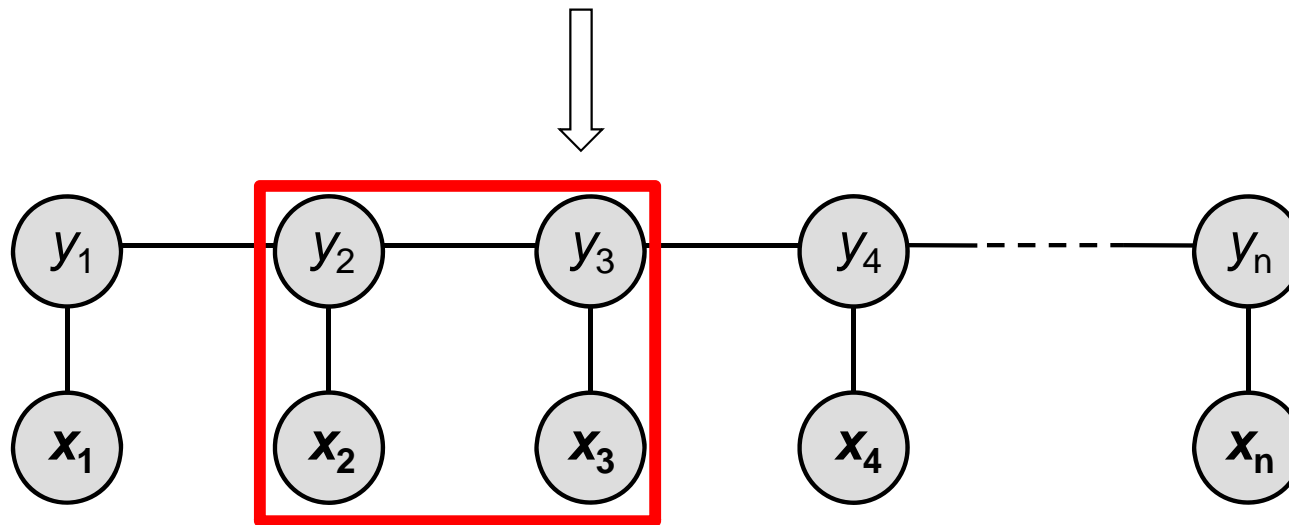
2) New edge (transition) features

◆ Traditional edge features avoid observations

- to avoid feature explosion

◆ **Proposal:** New edge features involving observations

- Tractable feature explosion: #feature increased 9 times
- → 10 millions features



3) Very fast online learning ***

- ◆ Large-scale word segmentation is already slow

3) Very fast online learning ***

- ◆ Large-scale word segmentation is already slow
 - ◆ With high dimensional new features (10 millions)
- Even more slower training speed than before

3) Very fast online learning ***

- ◆ Large-scale word segmentation is already slow
- ◆ With high dimensional new features (10 millions)
 - Even more slower training speed than before
 - Need a fast training method, even with 10 millions of features

3) Very fast online learning ***

◆ Batch training

- Limited memory BFGS (LBFGS)
- Too slow, need 300 passes for training

3) Very fast online learning ***

◆ Batch training

- Limited memory BFGS (LBFGS)
- Too slow, need 300 passes for training

◆ Existing online training methods

- Stochastic gradient descent (SGD)
- Moderately fast, need 50 passes for training

3) Very fast online learning ***

◆ Batch training

- Limited memory BFGS (LBFGS)
- Too slow, need 300 passes for training

◆ Existing online training methods

- Stochastic gradient descent (SGD)
- Moderately fast, need 50 passes for training

◆ **Main proposal:** a very fast online training method

- Finish the training in 10 passes

3) Very fast online learning ***


◆ Potential problem of existing online training

- Over-simplified setting
- Learning rate: use only 1 scalar for all weights

3) Very fast online learning ***

◆ Potential problem of existing online training

- Over-simplified setting
- Learning rate: use only 1 scalar for all weights

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \boxed{\gamma_t} \nabla_{\mathbf{w}_t} \mathcal{L}_{stoch}(\mathbf{z}_i, \mathbf{w}_t)$$


3) Very fast online learning ***


- ◆ better to use different learning rates for different features?

3) Very fast online learning ***

- ◆ better to use different learning rates for different features?
- ◆ → YES. We use **a vector of learning rates**

3) Very fast online learning ***

- ◆ better to use different learning rates for different features?
- ◆ → YES. We use a vector of learning rates

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \boxed{\boldsymbol{\gamma}_t} \cdot \mathbf{g}_t$$


$$\boldsymbol{\gamma}_t \in \mathbb{R}_+^f$$

$$\mathbf{g}_t = \nabla_{\mathbf{w}_t} \mathcal{L}_{stoch}(\mathbf{z}_i, \mathbf{w}_t) = \nabla_{\mathbf{w}_t} \left\{ \ell(\mathbf{z}_i, \mathbf{w}_t) - \frac{\|\mathbf{w}_t\|^2}{2n\sigma^2} \right\}$$

3) Very fast online learning ***

- ◆ better to use different learning rates for different features?
- ◆ → YES. We use **a vector of learning rates**

3) Very fast online learning ***

- ◆ better to use different learning rates for different features?
- ◆ → YES. We use **a vector of learning rates**
- ◆ → **Learning of Learning rates**: How to learn the values of the learning rates?

3) Very fast online learning ***

- ◆ better to use different learning rates for different features?
- ◆ → YES. We use **a vector of learning rates**
- ◆ → **Learning of Learning rates**: How to learn the values of the learning rates?

*Learning rates are learned from **feature frequency information**

***Higher** frequency feature → **lower** learning rate

3) Very fast online learning ***

ADF learning algorithm

```
1: procedure ADF( $q, c, \alpha, \beta$ )
2:    $w \leftarrow 0, t \leftarrow 0, v \leftarrow 0, \gamma \leftarrow c$ 
3:   repeat until convergence
4:     . Draw a sample  $z_i$  at random
5:     .  $v \leftarrow \text{UPDATE}(v, z_i)$ 
6:     . if  $t > 0$  and  $t \bmod q = 0$ 
7:       . .  $\gamma \leftarrow \text{UPDATE}(\gamma, v)$ 
8:       . .  $v \leftarrow 0$ 
9:       .  $g \leftarrow \nabla_w \mathcal{L}_{stoch}(z_i, w)$ 
10:      .  $w \leftarrow w + \gamma \cdot g$ 
11:      .  $t \leftarrow t + 1$ 
12:   return  $w$ 
```

3) Very fast online learning ***

ADF learning algorithm

```
1: procedure ADF( $q, c, \alpha, \beta$ )
2:    $w \leftarrow 0, t \leftarrow 0, v \leftarrow 0, \gamma \leftarrow c$ 
3:   repeat until convergence
4:     . Draw a sample  $z_i$  at random
5:     .  $v \leftarrow \text{UPDATE}(v, z_i)$ 
6:     . if  $t > 0$  and  $t \bmod q = 0$ 
7:       . .  $\gamma \leftarrow \text{UPDATE}(\gamma, v)$ 
8:       . .  $v \leftarrow 0$ 
9:       .  $g \leftarrow \nabla_w \mathcal{L}_{stoch}(z_i, w)$ 
10:      .  $w \leftarrow w + \gamma \cdot g$ 
11:      .  $t \leftarrow t + 1$ 
12:   return  $w$ 
```


3) Very fast online learning ***

ADF learning algorithm

```
1: procedure ADF( $q, c, \alpha, \beta$ )
2:    $w \leftarrow 0, t \leftarrow 0, v \leftarrow 0, \gamma \leftarrow c$ 
3:   repeat until convergence
4:     . Draw a sample  $z_i$  at random
5:     .  $v \leftarrow \text{UPDATE}(v, z_i)$ 
6:     . if  $t > 0$  and  $t \bmod q = 0$ 
7:       . .  $\gamma \leftarrow \text{UPDATE}(\gamma, v)$ 
8:       . .  $v \leftarrow 0$ 
9:       .  $g \leftarrow \nabla_w \mathcal{L}_{stoch}(z_i, w)$ 
10:      .  $w \leftarrow w + \boxed{\gamma} g$ 
11:      .  $t \leftarrow t + 1$ 
12:   return  $w$ 
```

Convergence Analysis

◆ Good convergence properties of the proposed method

The ADF training is convergent

Theorem 1 Assume ϕ is the largest eigenvalue of the function $\mathbf{C}_t = \prod_{m=1}^t (\mathbf{I} - \gamma_0 \beta^m \mathbf{H}(\mathbf{w}^*))$. For the proposed ADF training, its convergence rate is bounded by ϕ , and we have

$$\phi \leq \exp \left\{ \frac{\gamma_0 \lambda \beta}{\beta - 1} \right\},$$

where λ is the minimum eigenvalue of $\mathbf{H}(\mathbf{w}^*)$.

Experiments

◆ Data

◆ Sighan bakeoff 2004

- Microsoft Research data (MSR)
- Peking University data (PKU)
- City University of Hongkong data (CU)

	#W.T.	#Word	#C.T.	#Char
MSR	8.8×10^4	2.4×10^6	5×10^3	4.1×10^6
CU	6.9×10^4	1.5×10^6	5×10^3	2.4×10^6
PKU	5.5×10^4	1.1×10^6	5×10^3	1.8×10^6

Experiments

◆ Results

- Baseline: CRF with SGD training

Data	Method	Passes	Train-Time (sec)	NWD Rec	Pre	Rec	CWS F-score
MSR	Baseline	50	4.7e3	72.6	96.3	95.9	96.1
	+ New features	50	1.2e4	75.3	97.2	97.0	97.1
	+ New word detection	50	1.2e4	78.2	97.5	96.9	97.2
	+ ADF training	10	2.3e3	77.5	97.6	97.2	97.4
CU	Baseline	50	2.9e3	68.5	94.0	93.9	93.9
	+ New features	50	7.5e3	68.0	94.4	94.5	94.4
	+ New word detection	50	7.5e3	68.8	94.8	94.5	94.7
	+ ADF training	10	1.5e3	68.8	94.8	94.7	94.8
PKU	Baseline	50	2.2e3	77.2	95.0	94.0	94.5
	+ New features	50	5.2e3	78.4	95.5	94.9	95.2
	+ New word detection	50	5.2e3	79.1	95.8	94.9	95.3
	+ ADF training	10	1.2e3	78.4	95.8	94.9	95.4

Experiments

◆ Results

- New feature (word feature + high dimensional edge features) helps

Data	Method	Passes	Train-Time (sec)	NWD Rec	Pre	Rec	CWS F-score
MSR	Baseline	50	4.7e3	72.6	96.3	95.9	96.1
	+ New features	50	1.2e4	75.3	97.2	97.0	97.1
	+ New word detection	50	1.2e4	78.2	97.5	96.9	97.2
	+ ADF training	10	2.3e3	77.5	97.6	97.2	97.4
CU	Baseline	50	2.9e3	68.5	94.0	93.9	93.9
	+ New features	50	7.5e3	68.0	94.4	94.5	94.4
	+ New word detection	50	7.5e3	68.8	94.8	94.5	94.7
	+ ADF training	10	1.5e3	68.8	94.8	94.7	94.8
PKU	Baseline	50	2.2e3	77.2	95.0	94.0	94.5
	+ New features	50	5.2e3	78.4	95.5	94.9	95.2
	+ New word detection	50	5.2e3	79.1	95.8	94.9	95.3
	+ ADF training	10	1.2e3	78.4	95.8	94.9	95.4

Experiments

◆ Results

- Joint modeling (word seg + new word detection) helps slightly on word segmentation, but largely on NWD rec.

Data	Method	Passes	Train-Time (sec)	NWD Rec	Pre	Rec	CWS F-score
MSR	Baseline	50	4.7e3	72.6	96.3	95.9	96.1
	+ New features	50	1.2e4	75.3	97.2	97.0	97.1
	+ New word detection	50	1.2e4	78.2	97.5	96.9	97.2
	+ ADF training	10	2.3e3	77.5	97.6	97.2	97.4
CU	Baseline	50	2.9e3	68.5	94.0	93.9	93.9
	+ New features	50	7.5e3	68.0	94.4	94.5	94.4
	+ New word detection	50	7.5e3	68.8	94.8	94.5	94.7
	+ ADF training	10	1.5e3	68.8	94.8	94.7	94.8
PKU	Baseline	50	2.2e3	77.2	95.0	94.0	94.5
	+ New features	50	5.2e3	78.4	95.5	94.9	95.2
	+ New word detection	50	5.2e3	79.1	95.8	94.9	95.3
	+ ADF training	10	1.2e3	78.4	95.8	94.9	95.4

Experiments

◆ Results

- ADF training greatly reduced training time, yet with even higher F-score!

Data	Method	Passes	Train-Time (sec)	NWD Rec	Pre	Rec	CWS F-score
MSR	Baseline	50	4.7e3	72.6	96.3	95.9	96.1
	+ New features	50	1.2e4	75.3	97.2	97.0	97.1
	+ New word detection	50	1.2e4	78.2	97.5	96.9	97.2
	+ ADF training	10	2.3e3	77.5	97.6	97.2	97.4
CU	Baseline	50	2.9e3	68.5	94.0	93.9	93.9
	+ New features	50	7.5e3	68.0	94.4	94.5	94.4
	+ New word detection	50	7.5e3	68.8	94.8	94.5	94.7
	+ ADF training	10	1.5e3	68.8	94.8	94.7	94.8
PKU	Baseline	50	2.2e3	77.2	95.0	94.0	94.5
	+ New features	50	5.2e3	78.4	95.5	94.9	95.2
	+ New word detection	50	5.2e3	79.1	95.8	94.9	95.3
	+ ADF training	10	1.2e3	78.4	95.8	94.9	95.4

Experiments

◆ Results

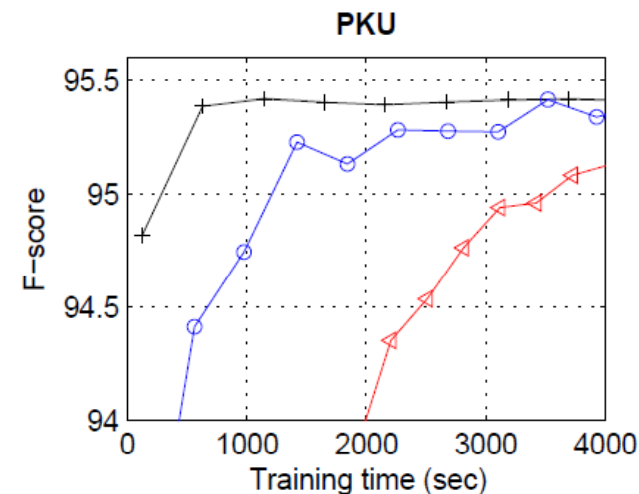
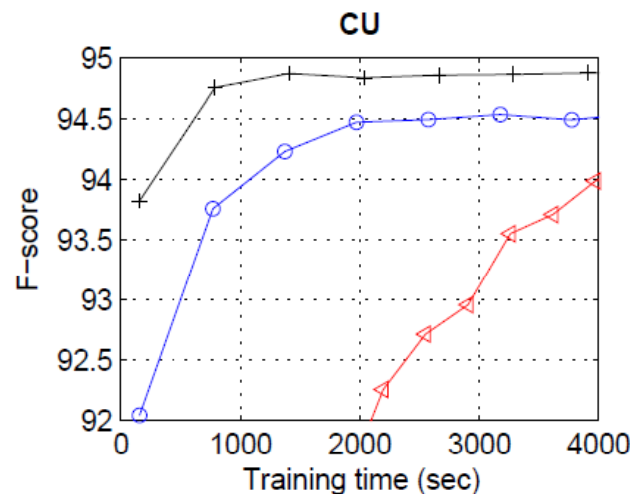
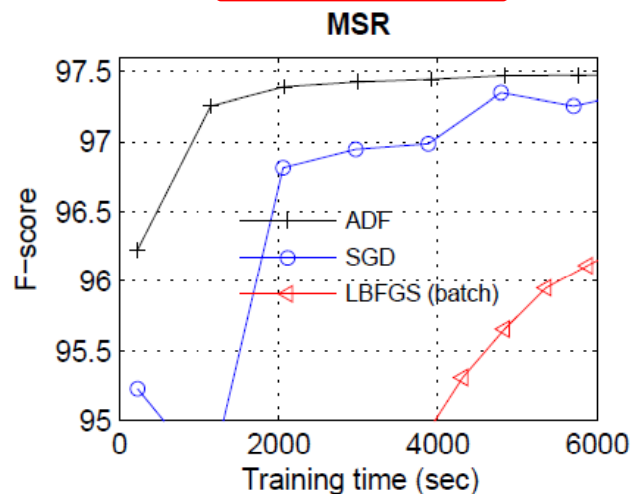
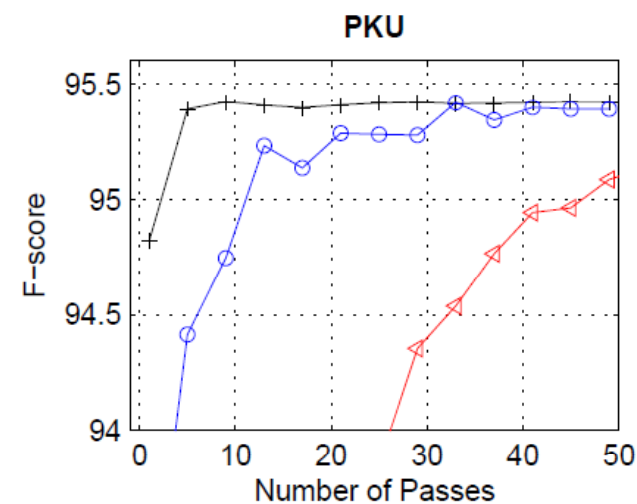
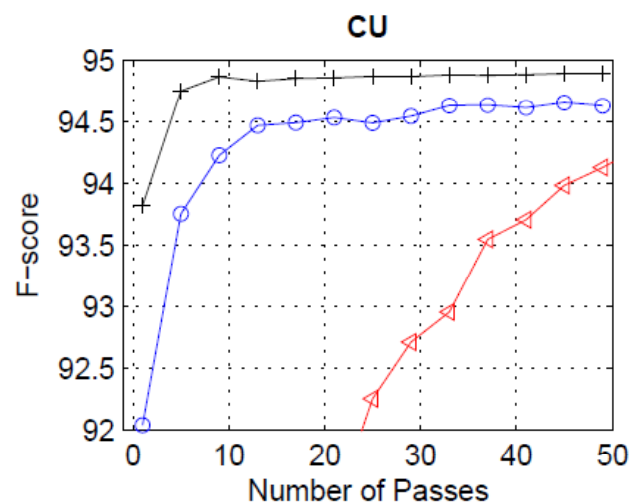
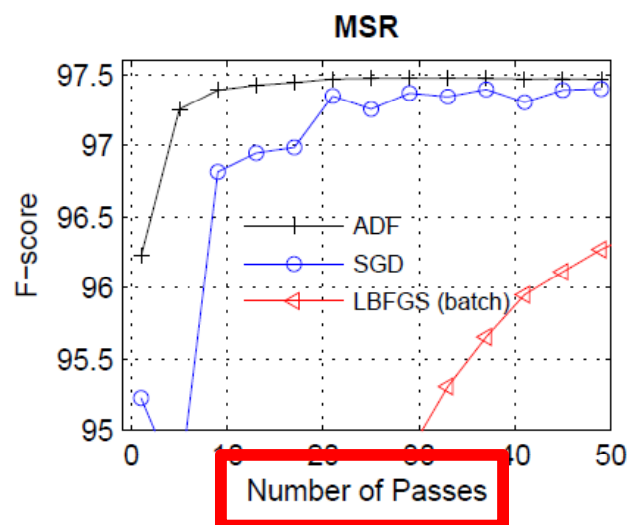
- ADF training greatly reduced training time, yet with even higher F-score!

Data	Method	Passes	Train-Time (sec)	NWD Rec	Pre	Rec	CWS F-score
MSR	Baseline	50	4.7e3	72.6	96.3	95.9	96.1
	+ New features	50	1.2e4	75.3	97.2	97.0	97.1
	+ New word detection	50	1.2e4	78.2	97.5	96.9	97.2
	+ ADF training	10	2.3e3	77.5	97.6	97.2	97.4
CU	Baseline	50	2.9e3	68.5	94.0	93.9	93.9
	+ New features	50	7.5e3	68.0	94.4	94.5	94.4
	+ New word detection	50	7.5e3	68.8	94.8	94.5	94.7
	+ ADF training	10	1.5e3	68.8	94.8	94.7	94.8
PKU	Baseline	50	2.2e3	77.2	95.0	94.0	94.5
	+ New features	50	5.2e3	78.4	95.5	94.9	95.2
	+ New word detection						
	+ ADF training						

We achieved best accuracy reports on the MSR & PKU datasets

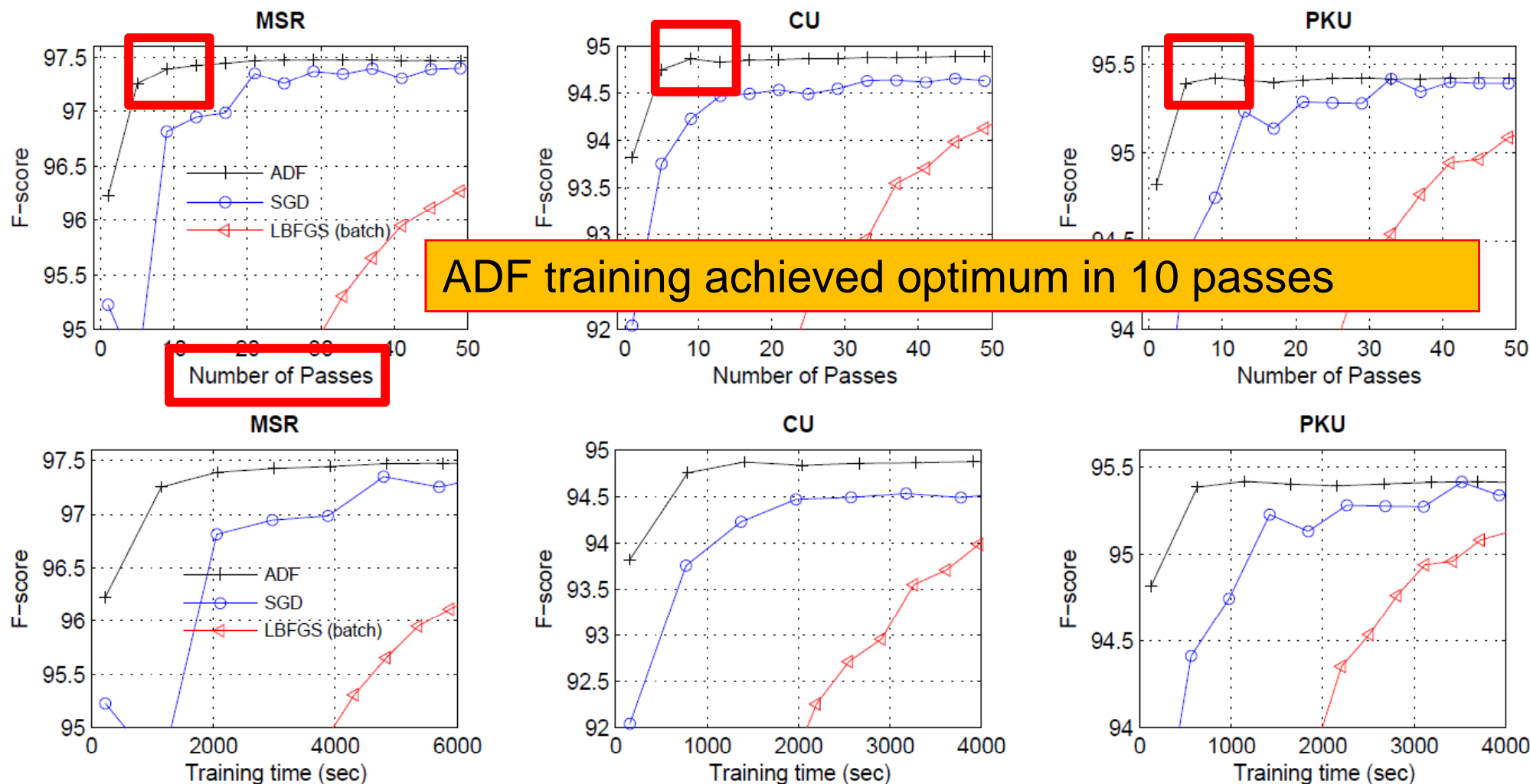
Experiments

◆ New training ADF vs. SGD vs. LBFGS



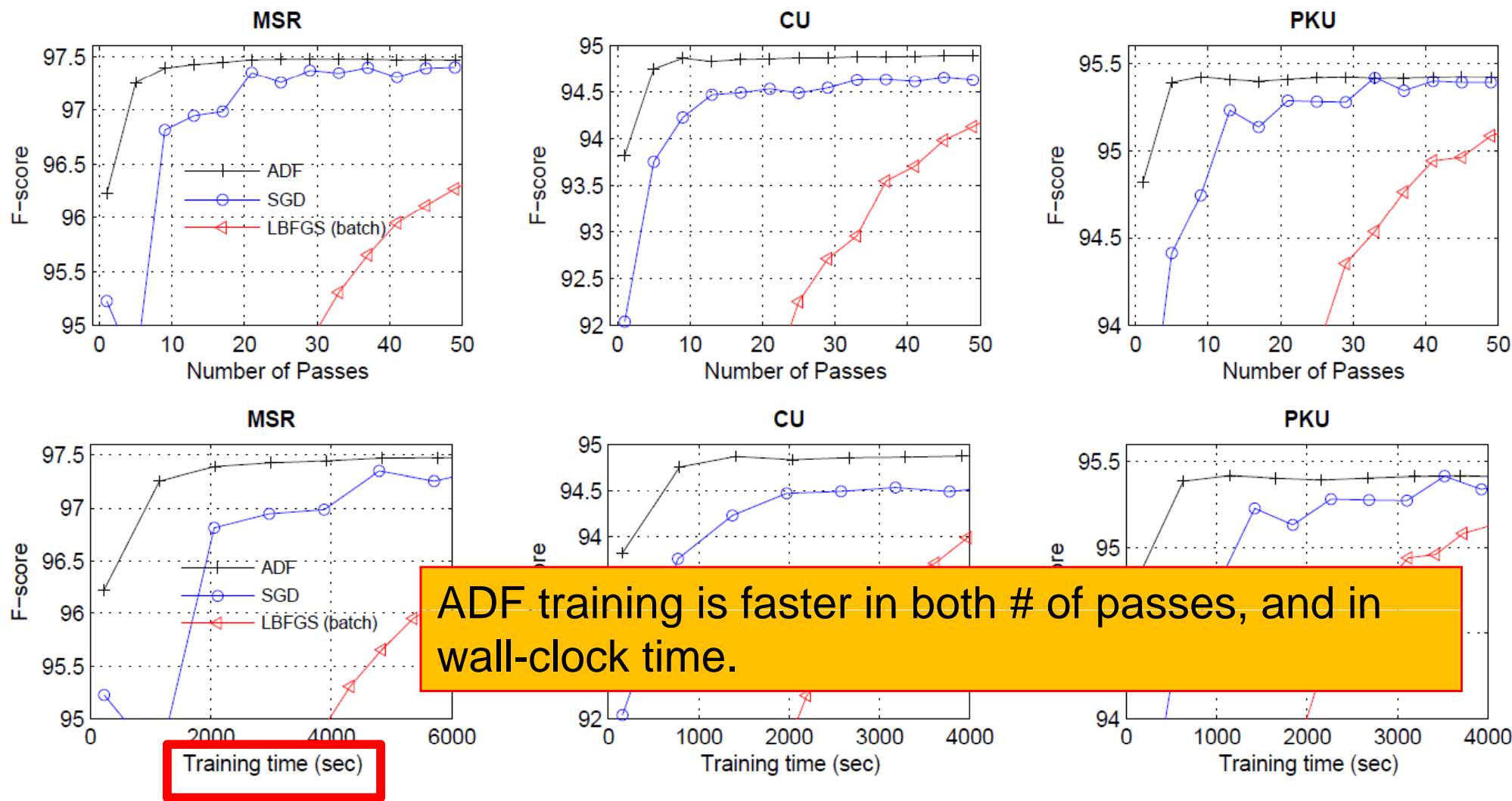
Experiments

◆ New training ADF vs. SGD vs. LBFGS



Experiments

◆ New training ADF vs. SGD vs. LBFGS



ADF training is faster in both # of passes, and in wall-clock time.

Conclusions

◆ 3 proposals

- Joint modeling: segmentation + new word detection
- New features: word features + high dimensional edge features
- A very fast online training method (main proposal)

Conclusions

◆ 3 proposals

- Joint modeling: segmentation + new word detection
- New features: word features + high dimensional edge features
- A very fast online training method (main proposal)

◆ Experiments

- Joint modeling helps

Conclusions

◆ 3 proposals

- Joint modeling: segmentation + new word detection
- New features: word features + high dimensional edge features
- A very fast online training method (main proposal)

◆ Experiments

- Joint modeling helps
- New features helps a lot on model accuracy

Conclusions

◆ 3 proposals

- Joint modeling: segmentation + new word detection
- New features: word features + high dimensional edge features
- A very fast online training method (main proposal)

◆ Experiments

- Joint modeling helps
- New features helps a lot on model accuracy
- The new training method finishes the training in 10 passes


Conclusions

◆ 3 proposals

- Joint modeling: segmentation + new word detection
- New features: word features + high dimensional edge features
- A very fast online training method (main proposal)

◆ Experiments

- Joint modeling helps
- New features helps a lot on model accuracy
- The new training method finishes the training in 10 passes
- Final results beat the existing best reports on the datasets
 - More accurate & faster!

- 
- ◆ **Thanks!**
 - ◆ **Any question?**