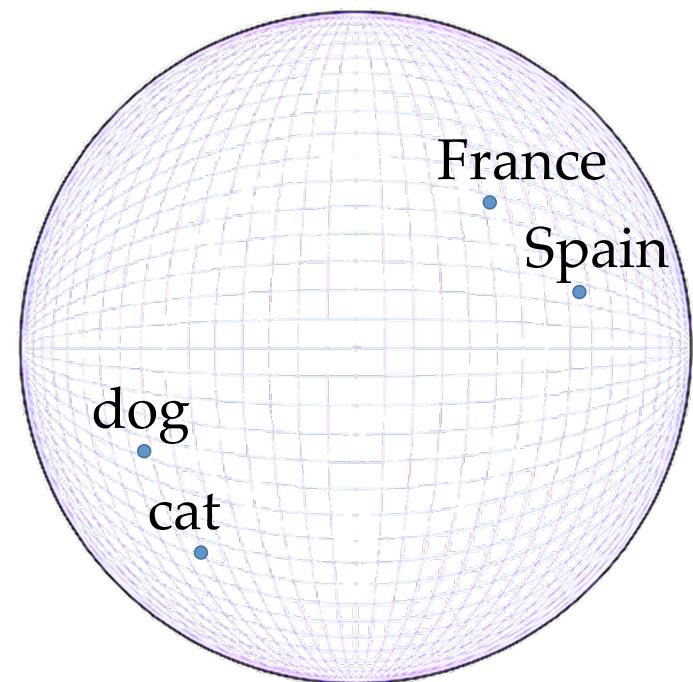


# Learning to translate with neural networks

Michael Auli

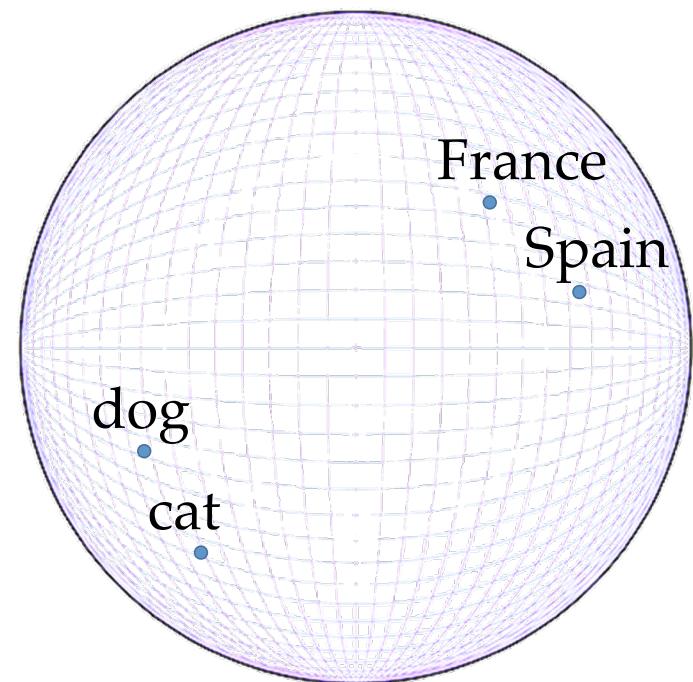
# Neural networks for text processing

- Similar words near each other



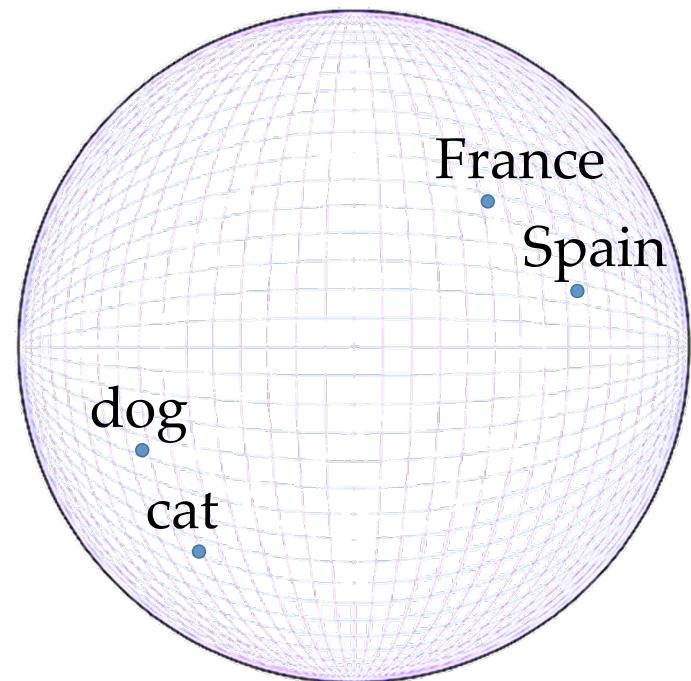
# Neural networks for text processing

- Similar words near each other
- Changing model parameters for one example effects similar words in similar contexts



# Neural networks for text processing

- Similar words near each other
- Changing model parameters for one example effects similar words in similar contexts
- Traditional discrete models treat each word separately



# Neural networks



Error 26% → 15% (Krizhevsky 2012)



Error 27% → 18 % (Hinton 2012)

Language Modeling

PPLX 141 → 101 (Mikolov 2011)

# Neural networks



Error 26% → 15% (Krizhevsky 2012)



Error 27% → 18 % (Hinton 2012)

Language Modeling    PPLX 141 → 101 (Mikolov 2011)

→ Machine Translation    This talk  
Le (2012), Kalchbrenner (2013),  
Devlin (2014), Sutskever (2014),  
Cho (2014), ...

# What happened in MT over the past 10 years?

The screenshot shows the English version of the United Nations website. At the top, there's a banner with the UN logo and the text "Welcome to the United Nations. It's your world." Below the banner, the "UNITED NATIONS" logo is displayed with the tagline "We the peoples... A stronger UN for a better world." A navigation bar at the top includes links for "Peace and Security", "Development", "Human Rights", "Humanitarian Affairs", and "International Law". On the left, a sidebar features a photo of a man speaking at a podium labeled "PRESIDENT". The main content area has sections like "Your United Nations" (with links to "UN at a Glance", "UN Charter", etc.) and "In Focus" (with links to "Syria | Chemical Weapons Report | Factsheet | Final report", "Central African Republic", "Ukraine", etc.).

The screenshot shows the Chinese version of the United Nations website. The header is in Chinese, with the UN logo and the text "欢迎来到联合国". The main title is "联合国" (United Nations) with the subtitle "我联合国人民, 团结起来, 追求更美好的世界!". The navigation bar includes "和平与安全", "发展", "人权", "人道主义事务", and "国际法". The sidebar on the left shows a photo of children in a camp. The "In Focus" section contains links to "叙利亚危机 | 化学武器调查报告 | 方案 | 最终报告", "中非共和国局势", "乌克兰", etc.

# What happened in MT over the past 10 years?



“Learning simple models from large bi-texts is a solved problem”

(Lopez & Post, 2013)

# What happened in MT over the past 10 years?



“Learning simple models from large bi-texts is a solved problem”

(Lopez & Post, 2013)



WMT 2013



9 / 10 times

# Machine translation

本 地 区 的 发 展 和 进 步 。

development and progress of the region .

# Machine translation

本 地 区 的 发 展 和 进 步 。

Translation modeling



development and progress of the region .

# Machine translation

本 地 区 的 发 展 和 进 步 。

Translation modeling



Language modeling

development and progress of the region .



# Machine translation

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling

Language modeling

Optimization

# Machine translation

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling

Language modeling

Optimization

Reordering

# Machine translation

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling

Auli et al., EMNLP 2013; Hu et al., EACL 2014



Language modeling



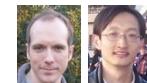
Optimization

Auli & Gao, ACL 2014



Reordering

Auli et al., EMNLP 2014



Koehn et al. (2003)

## Discrete phrase-based translation

本 地区 的 发展 和 进步 。

development and progress of the region .

Koehn et al. (2003)

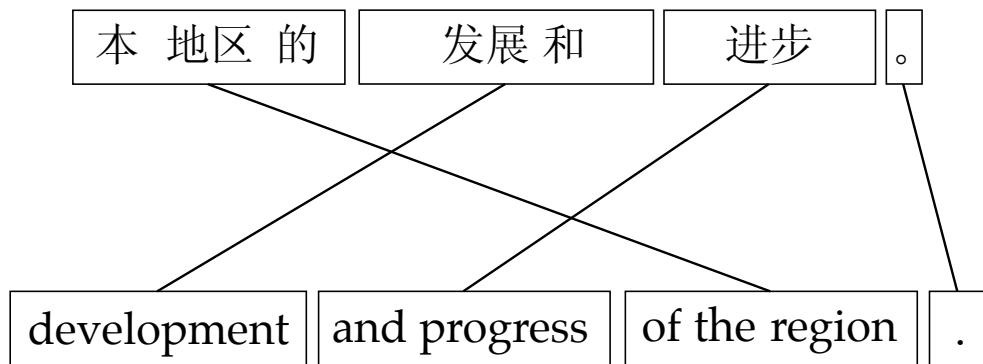
# Discrete phrase-based translation

本 地区 的 发展 和 进步 。

development and progress of the region .

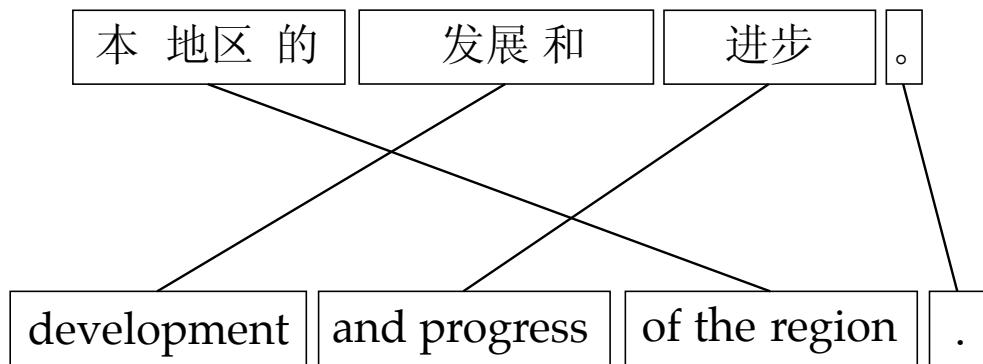
Koehn et al. (2003)

# Discrete phrase-based translation



Koehn et al. (2003)

# Discrete phrase-based translation



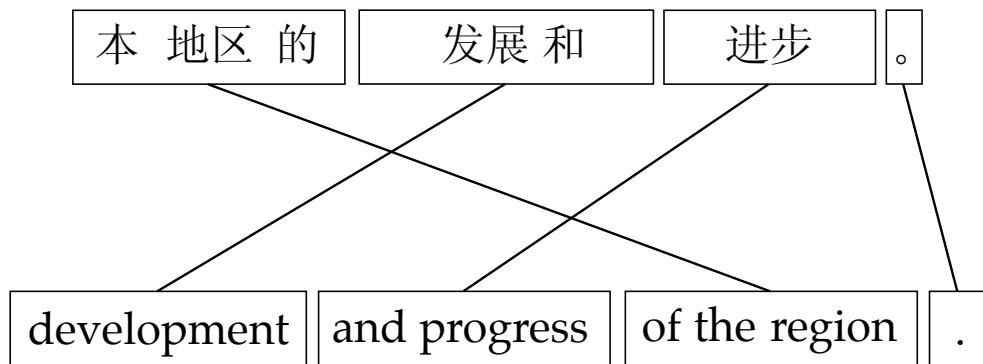
本 地区 的 → of the region

发展 → development

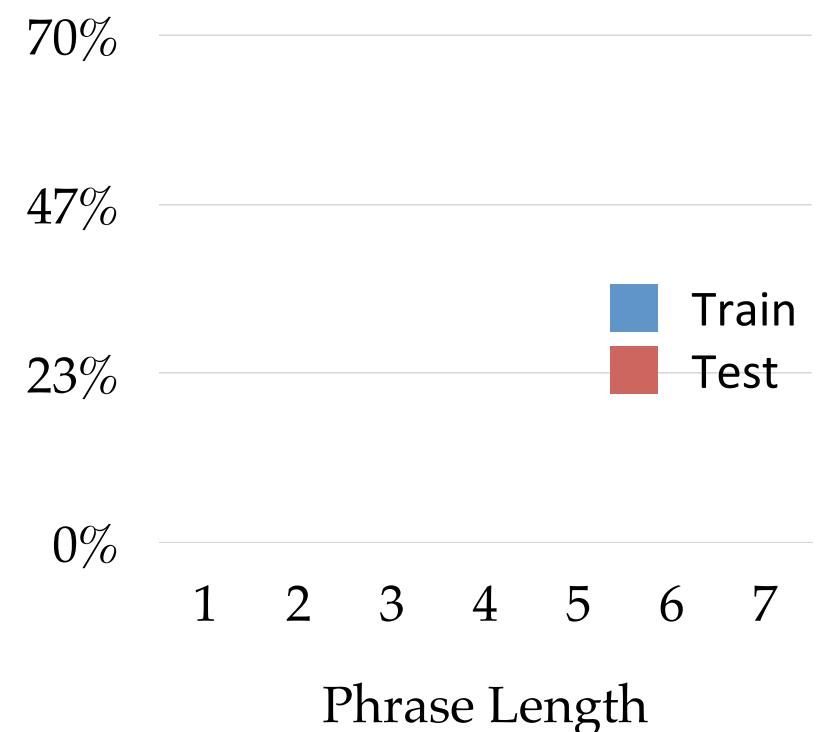
和 进步 → and progress

Koehn et al. (2003)

# Discrete phrase-based translation

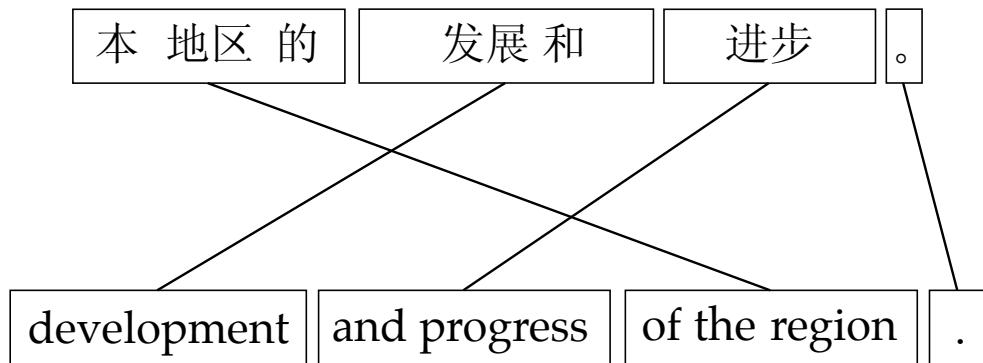


本 地区 的 → of the region  
发展 → development  
和 进步 → and progress

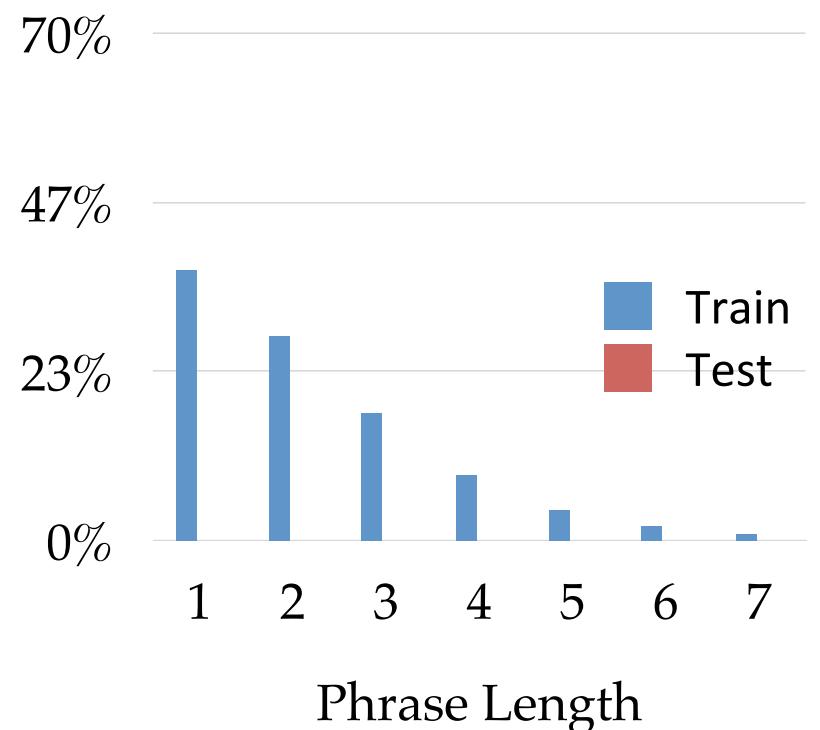


Koehn et al. (2003)

# Discrete phrase-based translation

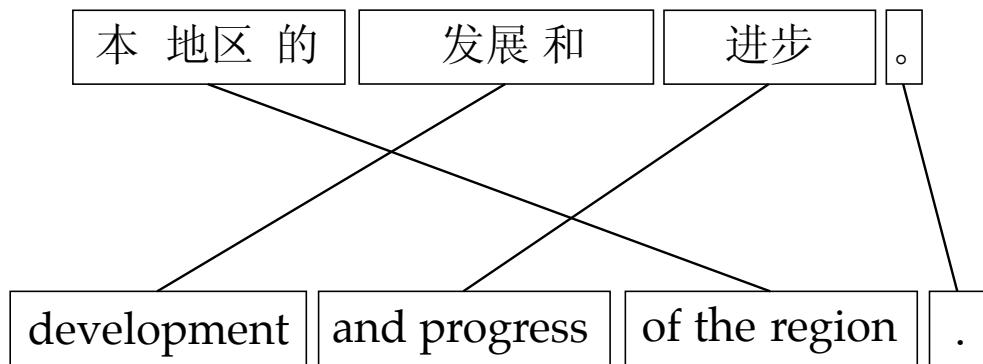


本地区的 → of the region  
发展 → development  
和进步 → and progress

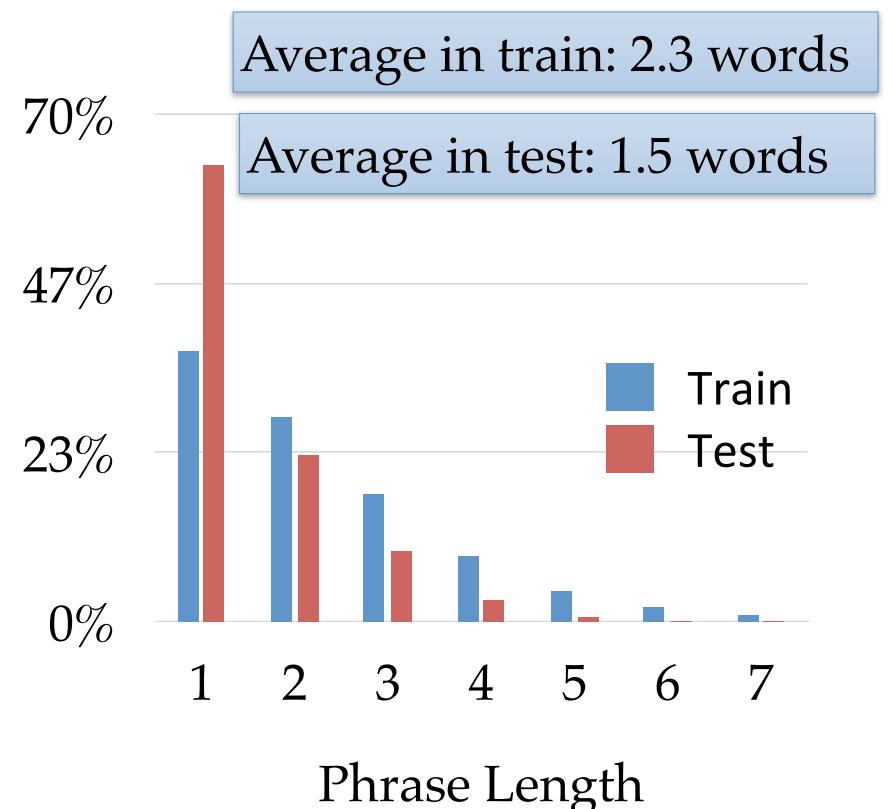


Koehn et al. (2003)

# Discrete phrase-based translation



本 地区 的 → of the region  
发展 → development  
和 进步 → and progress



Kneser & Ney (1996)

# Discrete n-gram language modeling

$p(\text{progress in the region}) =$

Kneser & Ney (1996)

# Discrete n-gram language modeling

$p(\text{progress in the region}) =$

Train data:

...  
development and progress of  
**the region. in ...**

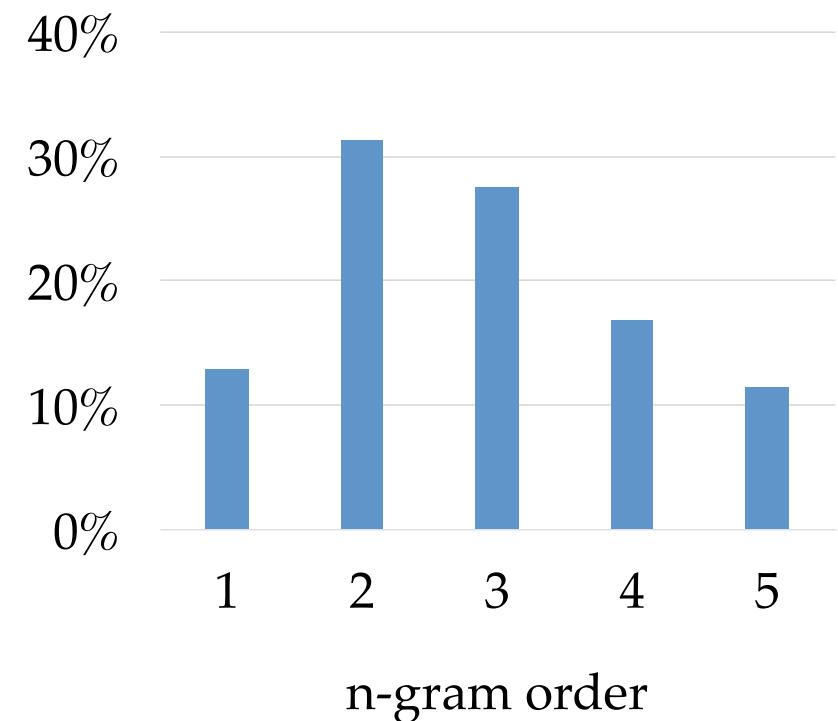
Kneser & Ney (1996)

# Discrete n-gram language modeling

$$\begin{aligned} p(\text{progress in the region}) = \\ p(\text{progress}) p(\text{in}) \\ p(\text{the}) p(\text{region} \mid \text{the}) \end{aligned}$$

Train data:

...  
development and progress of  
**the region.** in ...  
...



Does not include out-of-vocabulary tokens

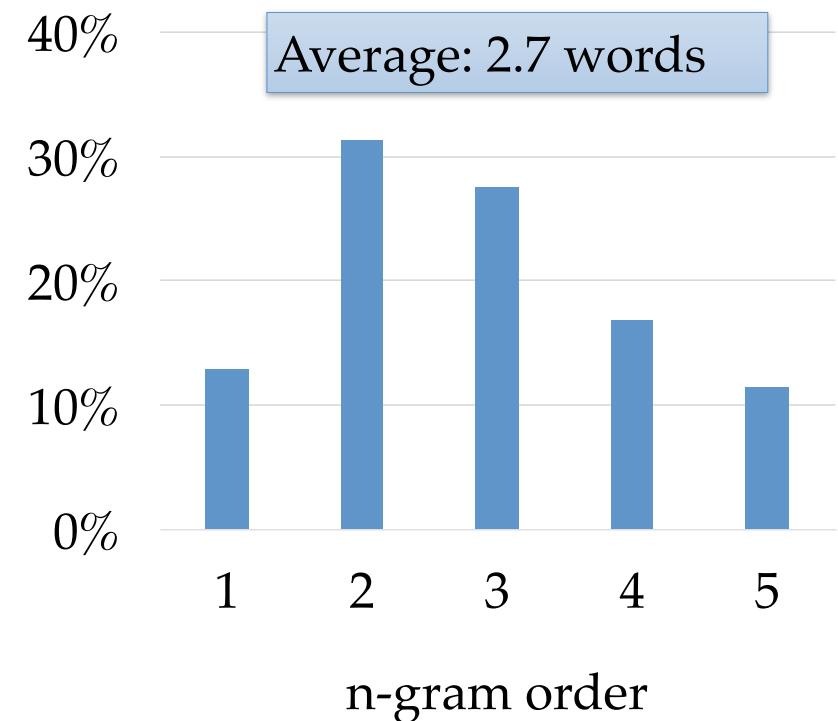
Kneser & Ney (1996)

# Discrete n-gram language modeling

$$p(\text{progress in the region}) = p(\text{progress}) p(\text{in}) \\ p(\text{the}) p(\text{region} \mid \text{the})$$

Train data:

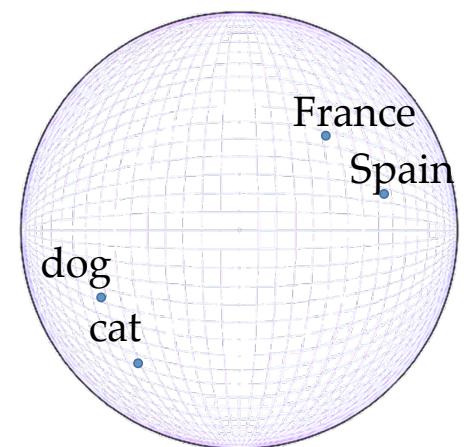
...  
development and progress of  
**the region.** in ...  
...



Does not include out-of-vocabulary tokens

# How can we improve this?

- Or: how to capture relationships beyond 1.5 to 2.7 words?
- Neural nets: distributional representations make it easier to capture relationships
- Recurrent nets: easy to model variable-length sequences



# This talk

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling

Auli et al., EMNLP 2013; Hu et al., EACL 2014



Language modeling



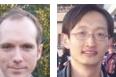
Optimization

Auli & Gao, ACL 2014



Reordering

Auli et al., EMNLP 2014



# This talk

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling



Auli et al., EMNLP 2013; Hu et al., EACL 2014

Language modeling



Optimization

Auli & Gao, ACL 2014



Reordering

Auli et al., EMNLP 2014

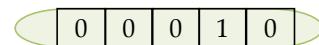


# Feed-forward network

e.g. bigram LM

# Feed-forward network

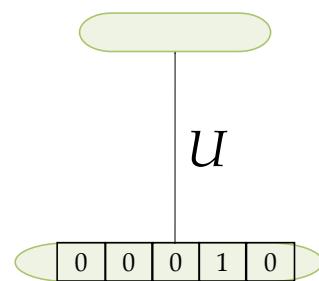
e.g. bigram LM



and

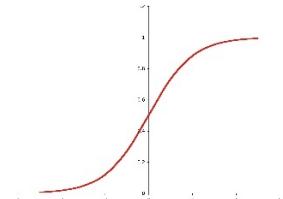
# Feed-forward network

e.g. bigram LM



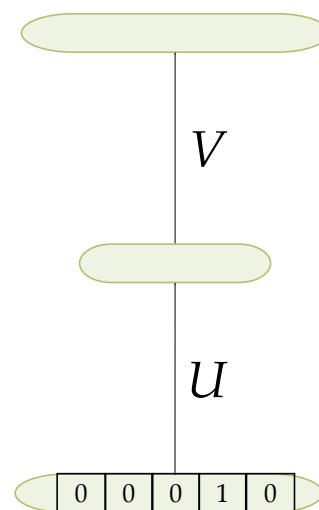
and

$$h_t = \sigma(Ux_t)$$
$$\sigma(z) = \frac{1}{1 + \exp\{-z\}}$$



# Feed-forward network

e.g. bigram LM

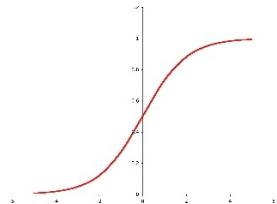


$$y_t = s(V h_t)$$

$$s(z) = \frac{\exp\{z\}}{\sum_{z'} \exp\{z'\}}$$

$$h_t = \sigma(U x_t)$$

$$\sigma(z) = \frac{1}{1 + \exp\{-z\}}$$



and

# Feed-forward network

e.g. bigram LM

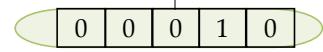
$p(\text{progress} \mid \text{and})$



$V$



$U$



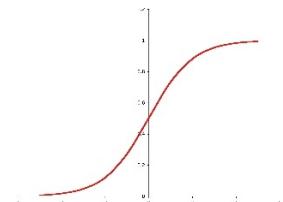
and

$$y_t = s(V h_t)$$

$$s(z) = \frac{\exp\{z\}}{\sum_{z'} \exp\{z'\}}$$

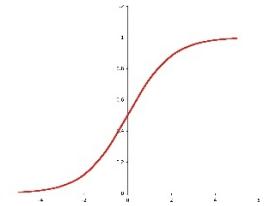
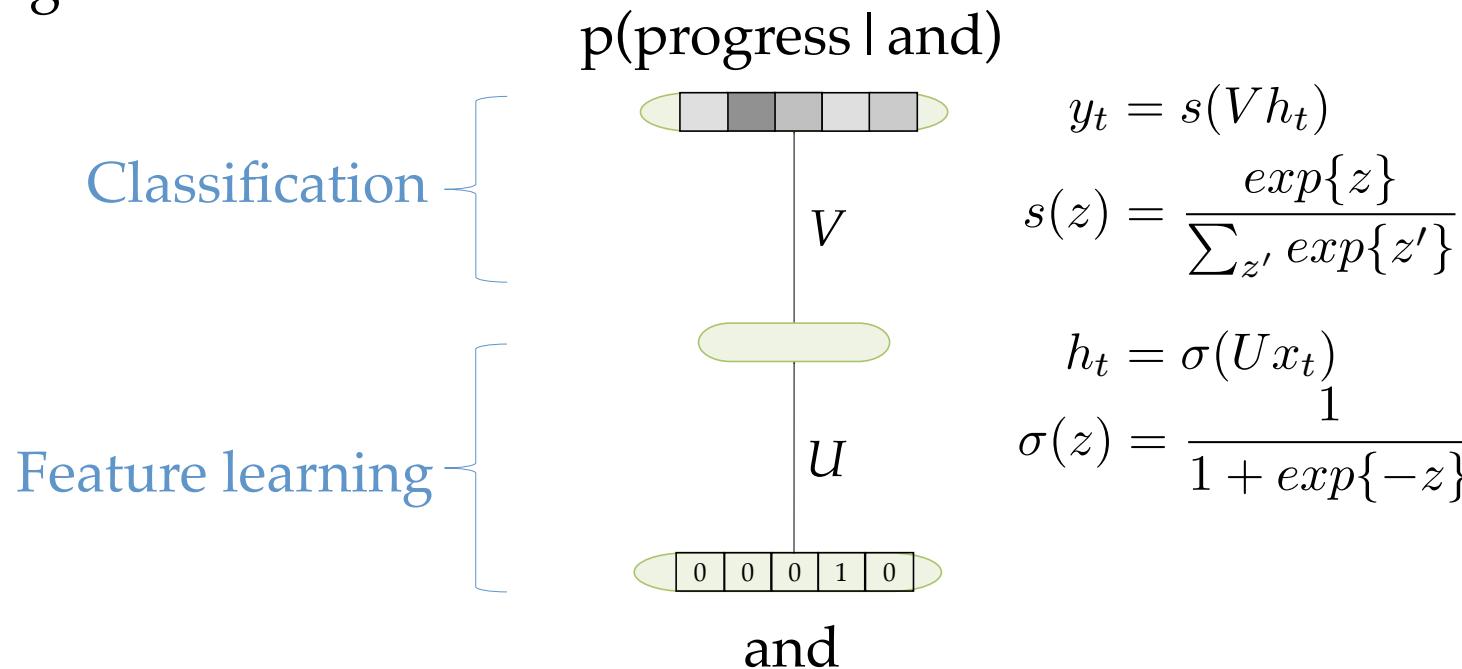
$$h_t = \sigma(U x_t)$$

$$\sigma(z) = \frac{1}{1 + \exp\{-z\}}$$



# Feed-forward network

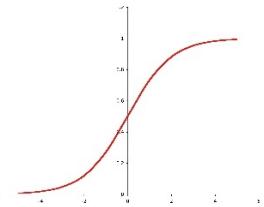
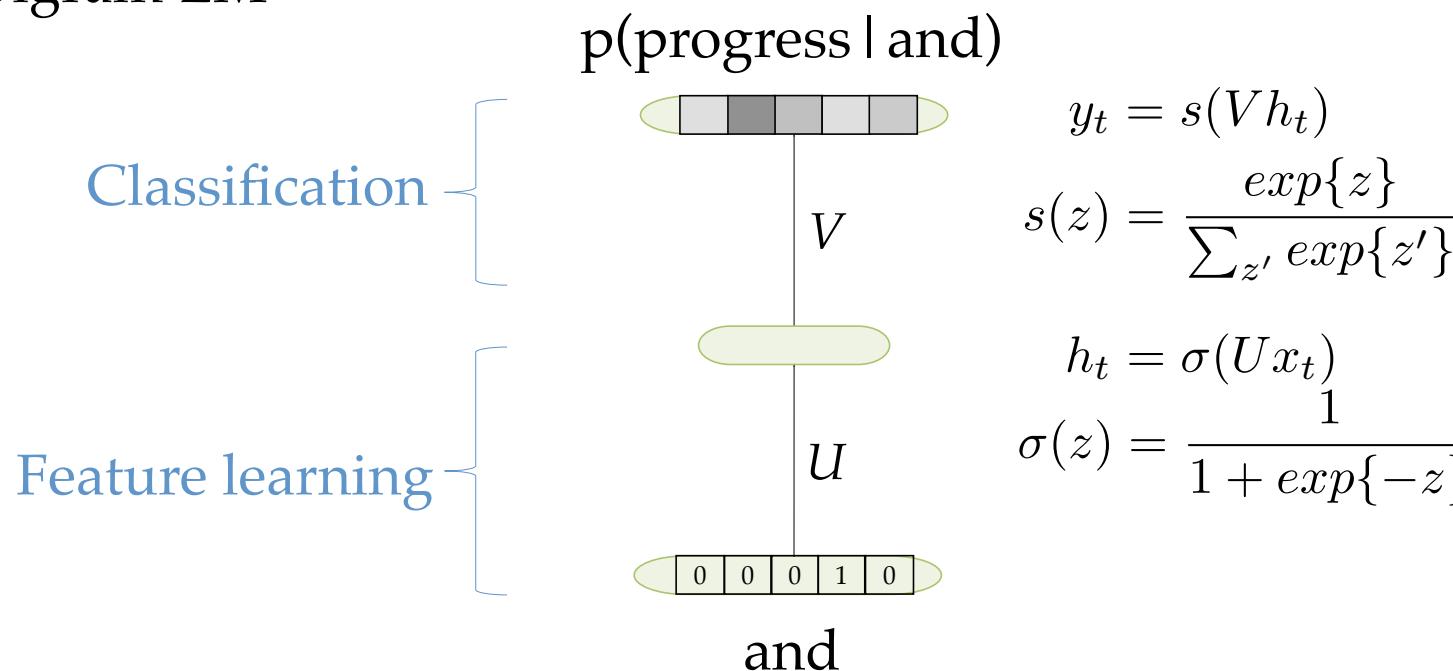
e.g. bigram LM



# Feed-forward network

e.g. bigram LM

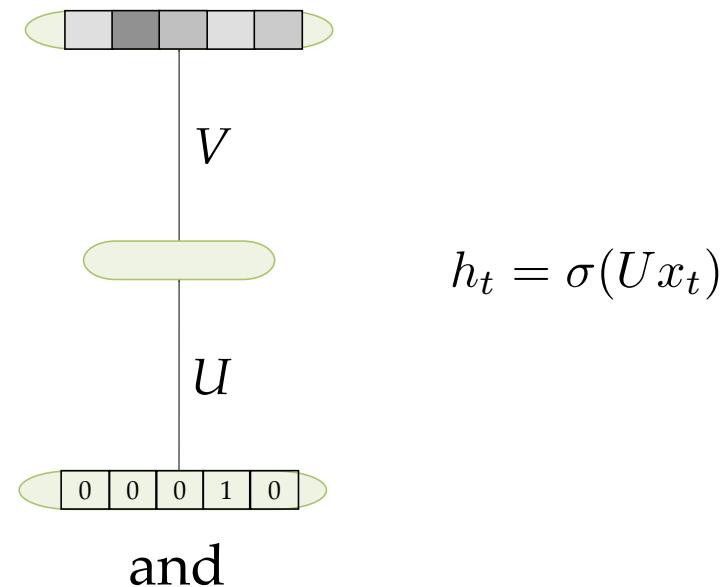
Still based on  
limited context!



# Recurrent network

e.g. bigram LM

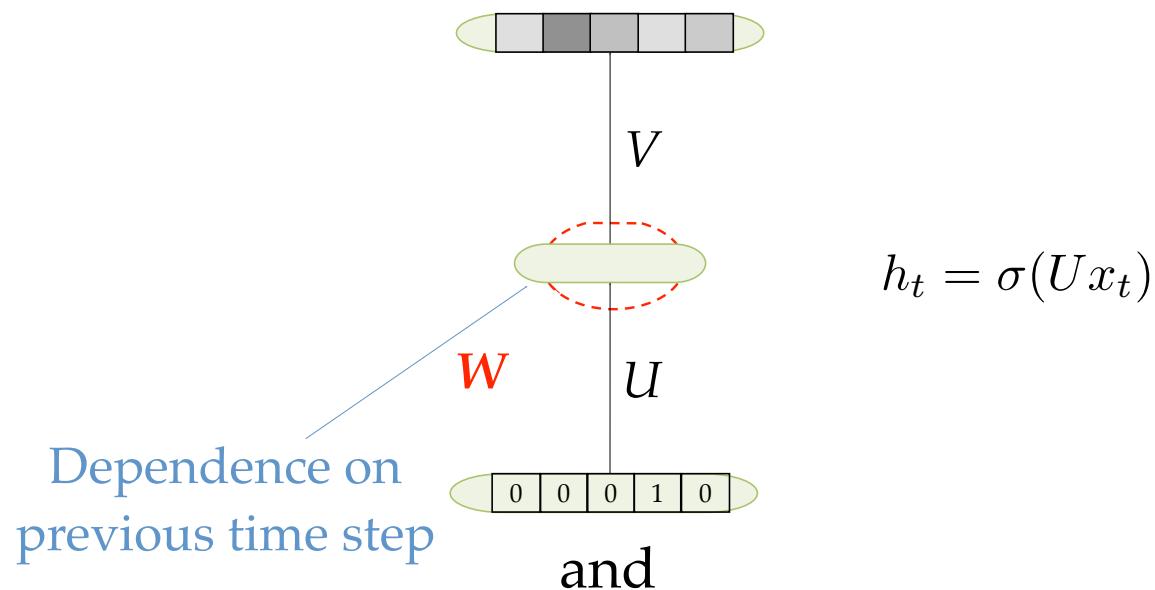
$p(\text{progress} \mid \text{and})$



# Recurrent network

e.g. bigram LM

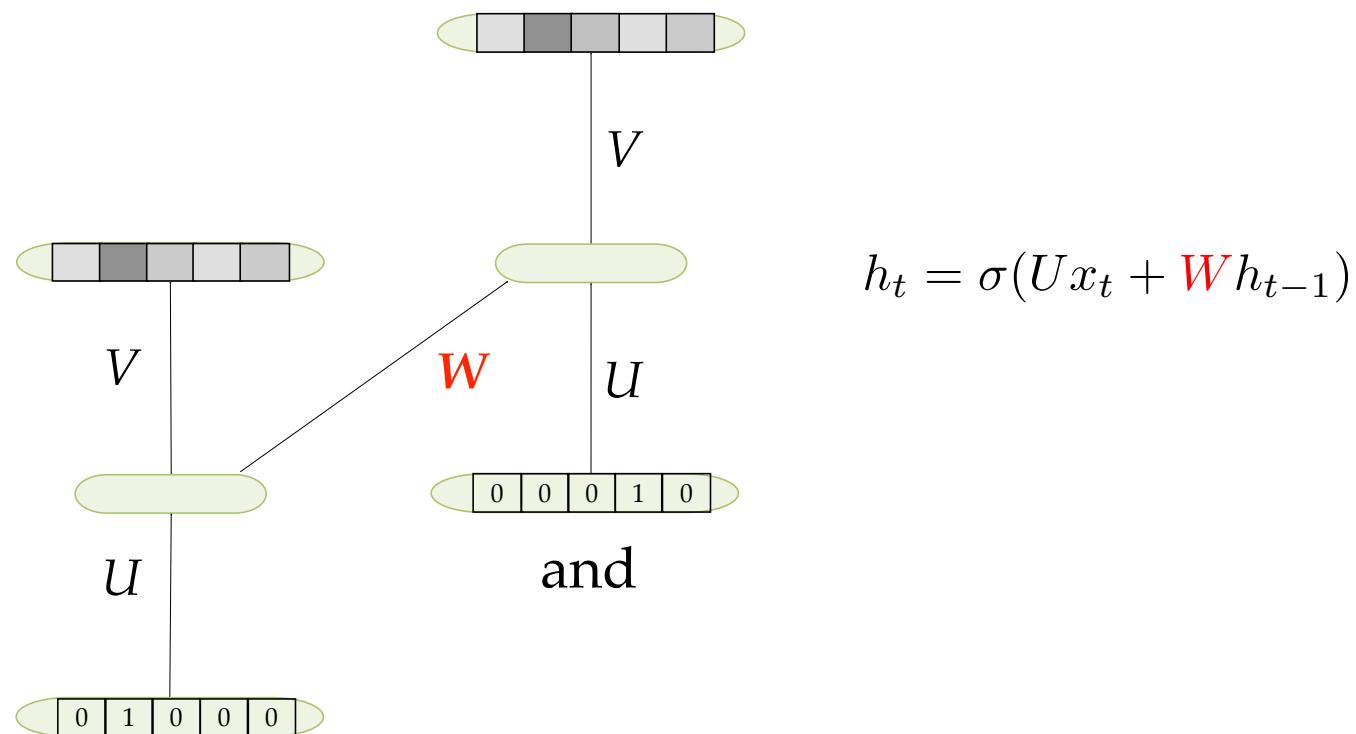
$p(\text{progress} \mid \text{and})$



# Recurrent network

e.g. bigram LM

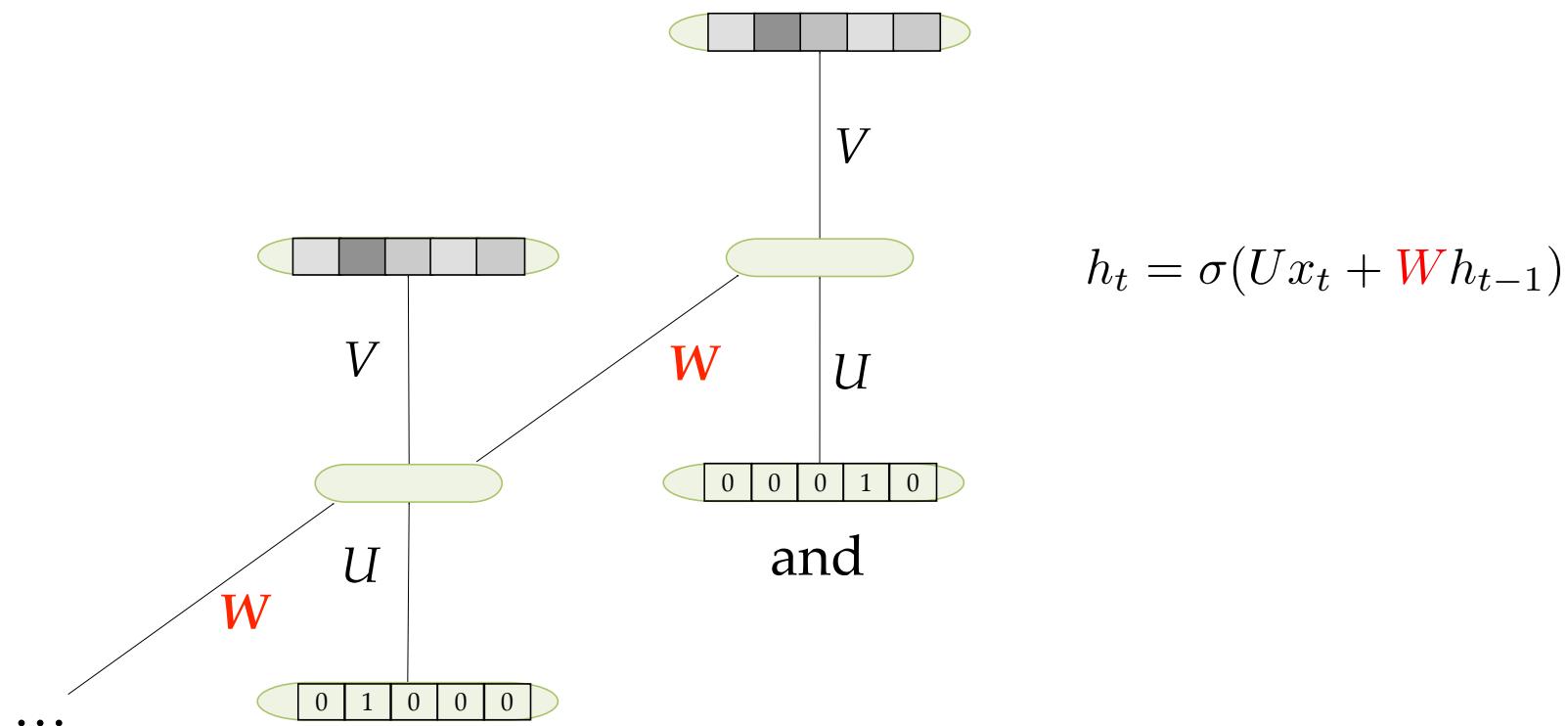
$p(\text{progress} \mid \text{and})$



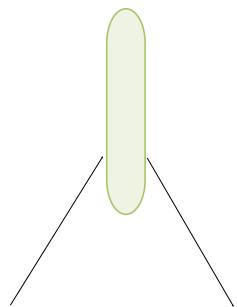
# Recurrent network

e.g. bigram LM

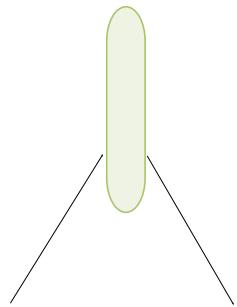
p(progress | and)



# Recurrent network

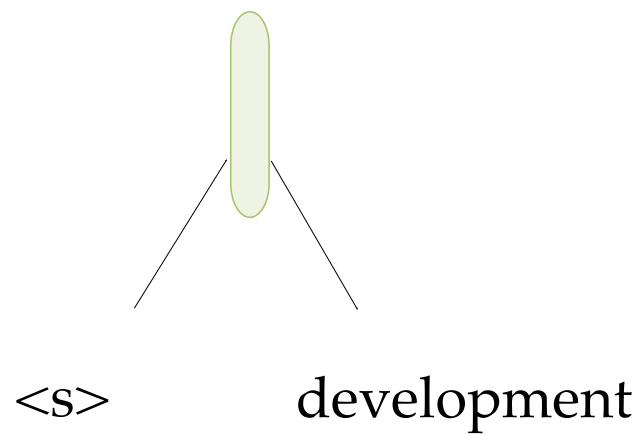


# Recurrent network

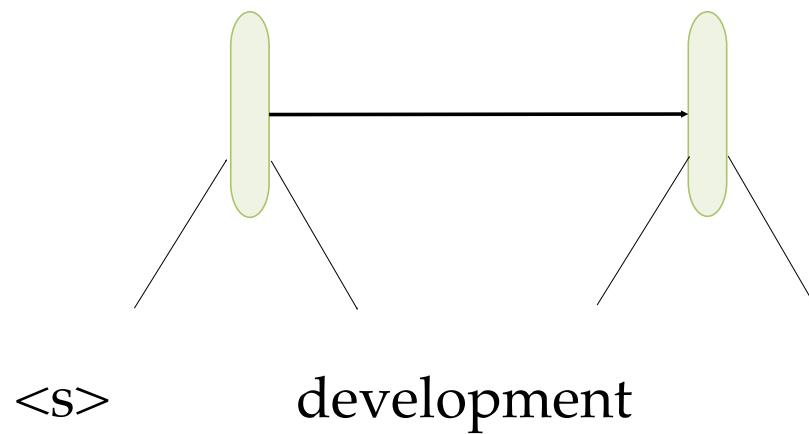


$\langle s \rangle \quad p(\text{development} | \langle s \rangle)$

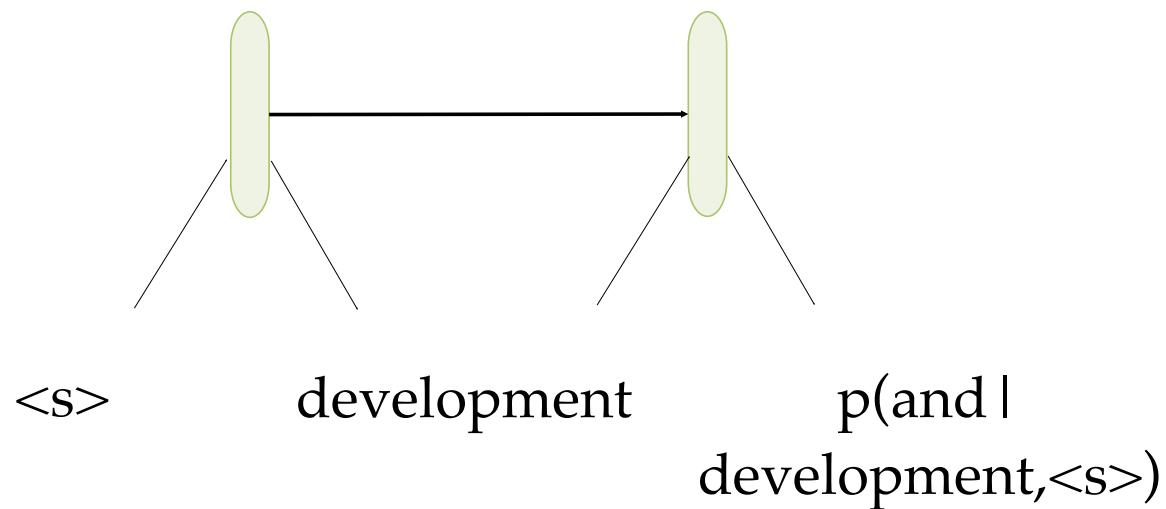
# Recurrent network



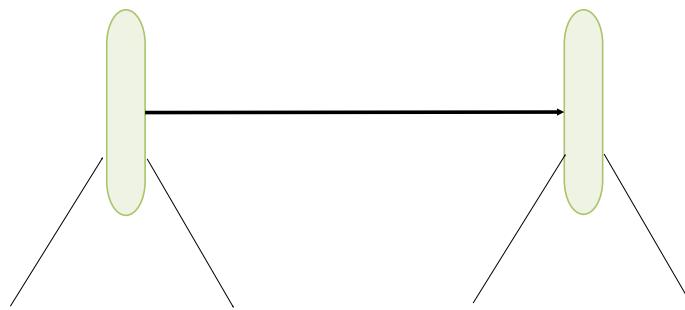
# Recurrent network



# Recurrent network



# Recurrent network

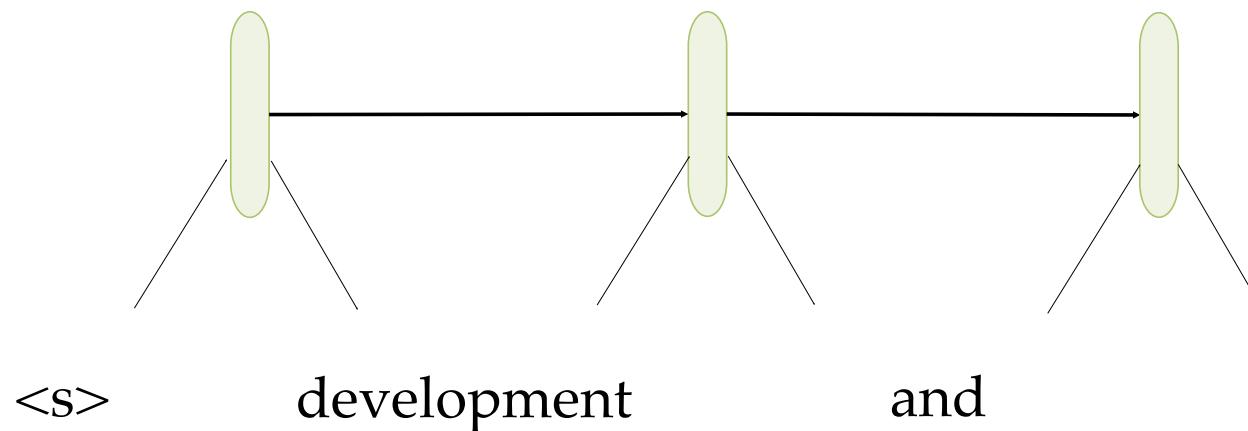


<S>

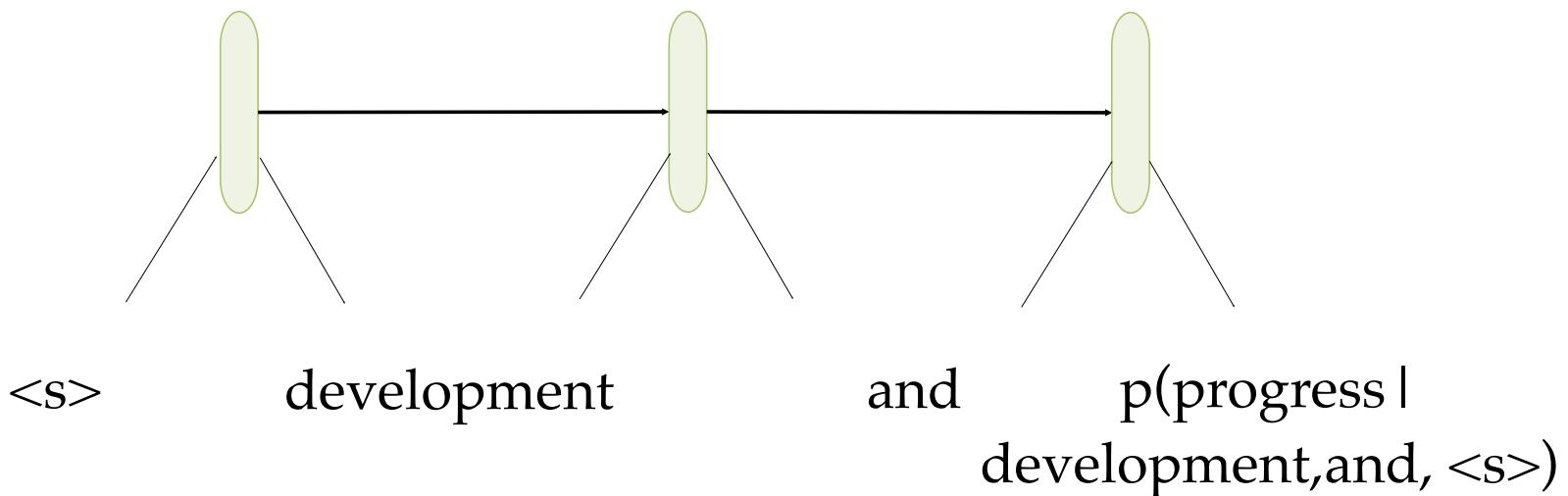
development

and

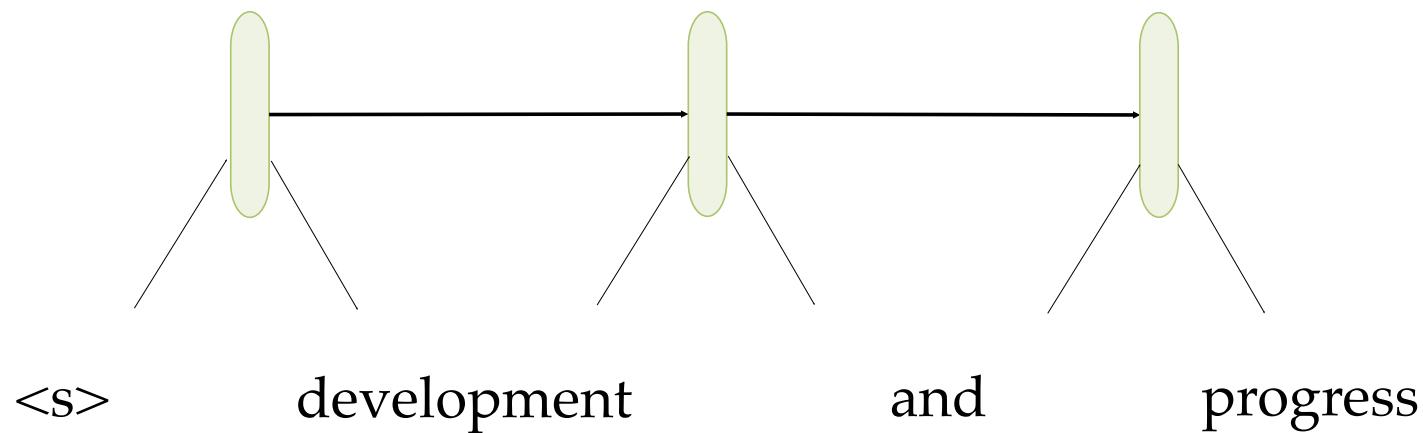
# Recurrent network



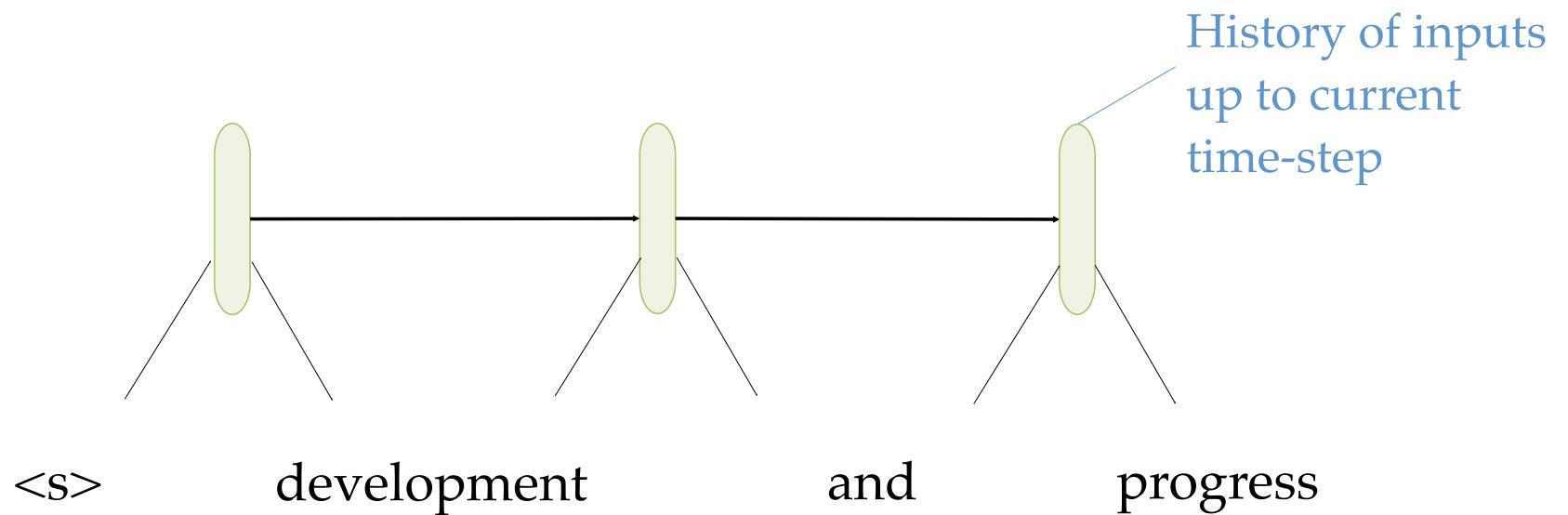
# Recurrent network



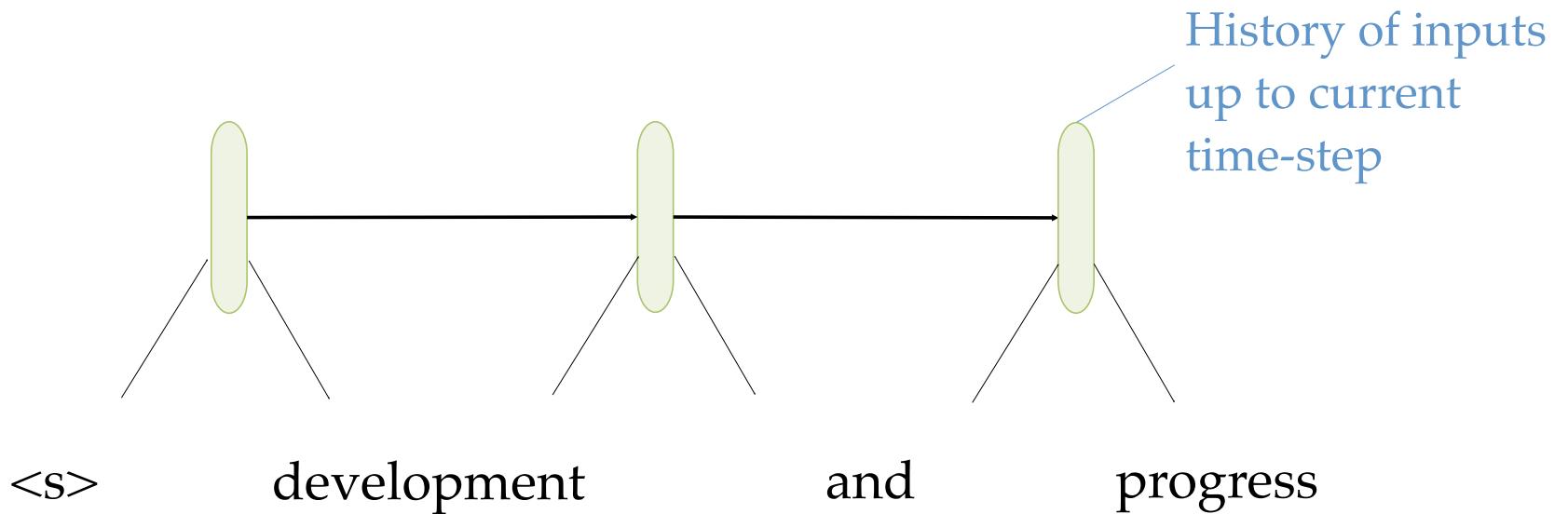
# Recurrent network



# Recurrent network



# Recurrent network



State of the art in language modeling (Mikolov 2011)

More accurate than feed-forward nets (Sundermeyer 2013)

# Combined language and translation model

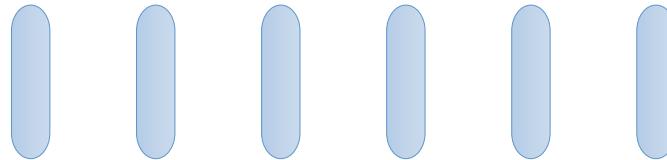
本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

# Combined language and translation model

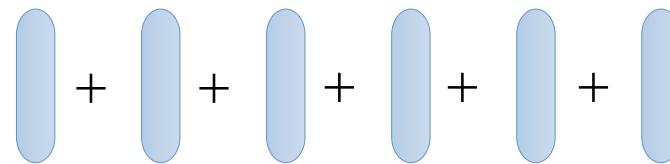
本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

# Combined language and translation model

本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

# Combined language and translation model

本 地 区 的 发 展 和 进 步



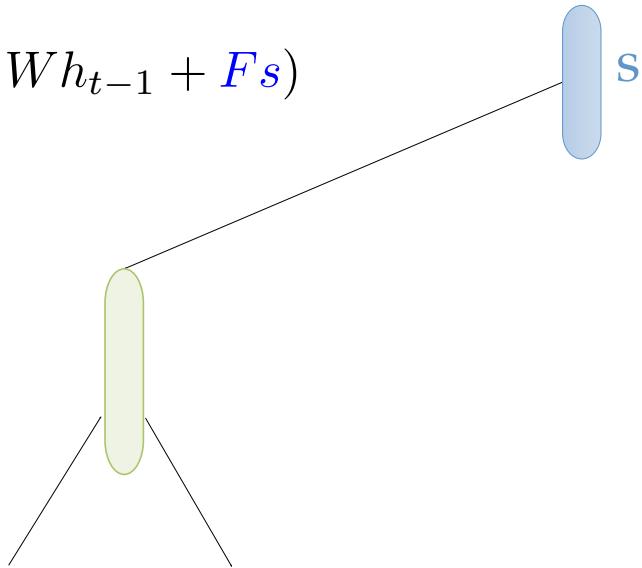
Auli et al., EMNLP 2013

Entire source sentence  
representation

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Auli et al., EMNLP 2013

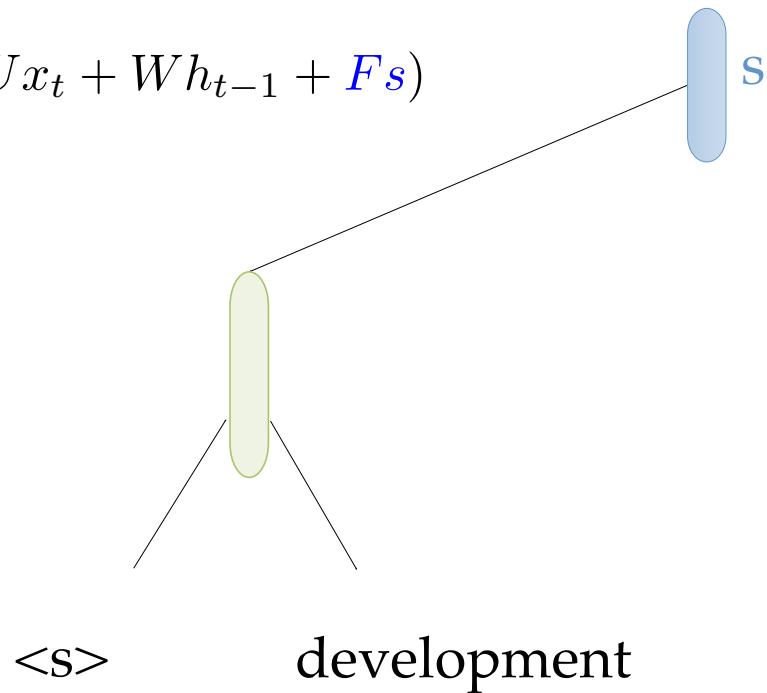
Entire source sentence  
representation

<S>

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Auli et al., EMNLP 2013

Entire source sentence  
representation

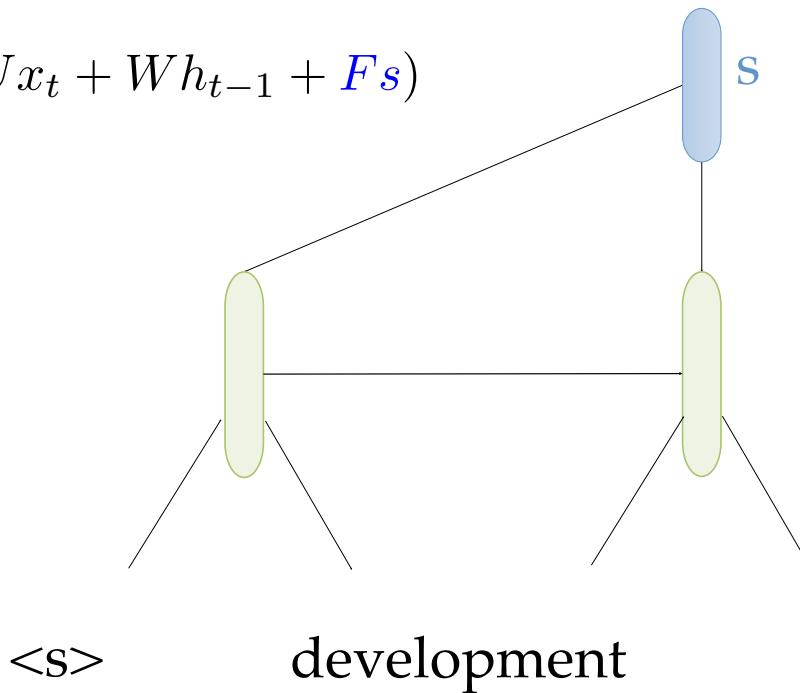
# Combined language and translation model

本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

$$h_t = \sigma(Ux_t + Wh_{t-1} + \mathbf{F}s)$$



Entire source sentence  
representation

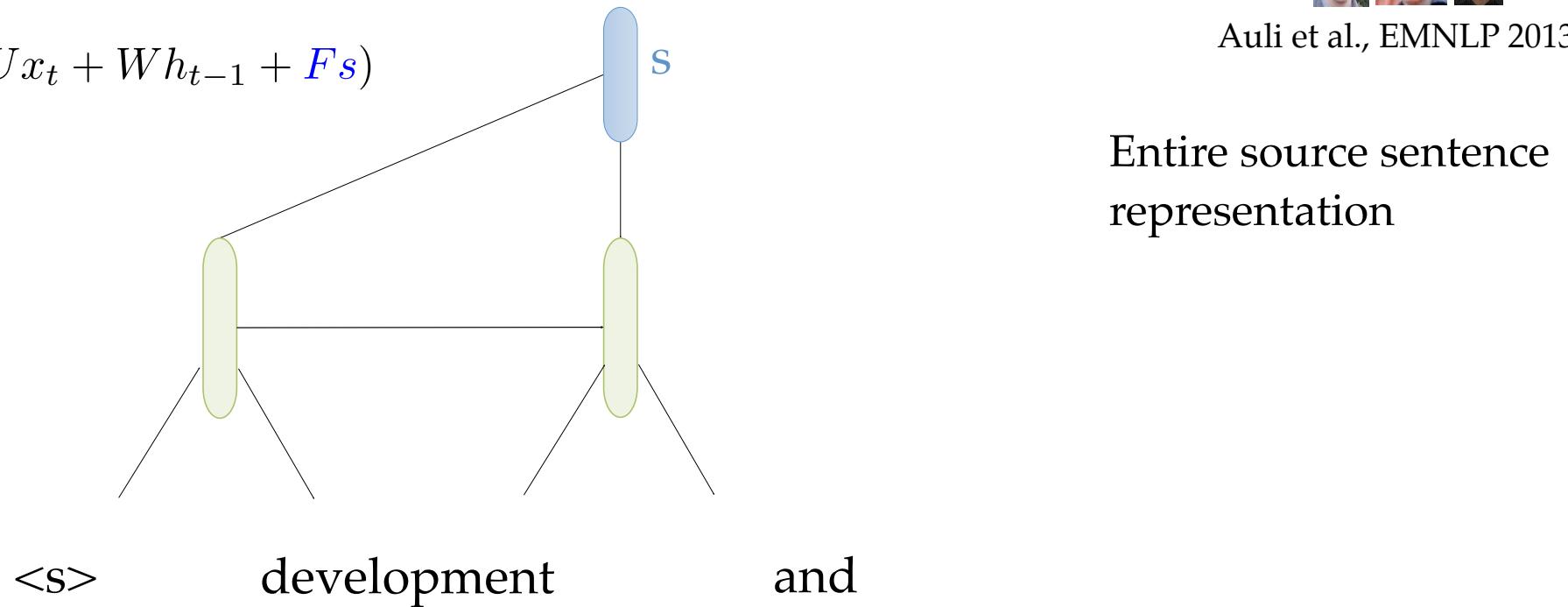
# Combined language and translation model

本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

$$h_t = \sigma(Ux_t + Wh_{t-1} + \mathbf{F}s)$$



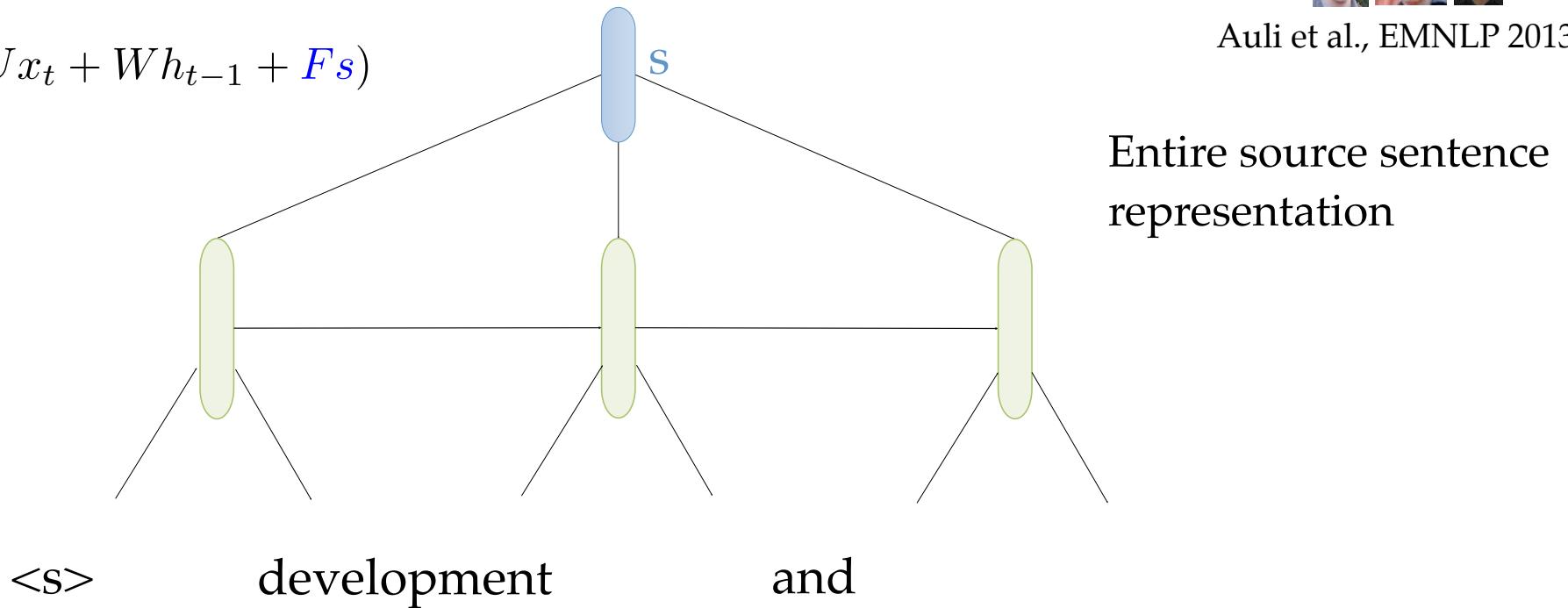
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \mathbf{F}s)$$



Auli et al., EMNLP 2013



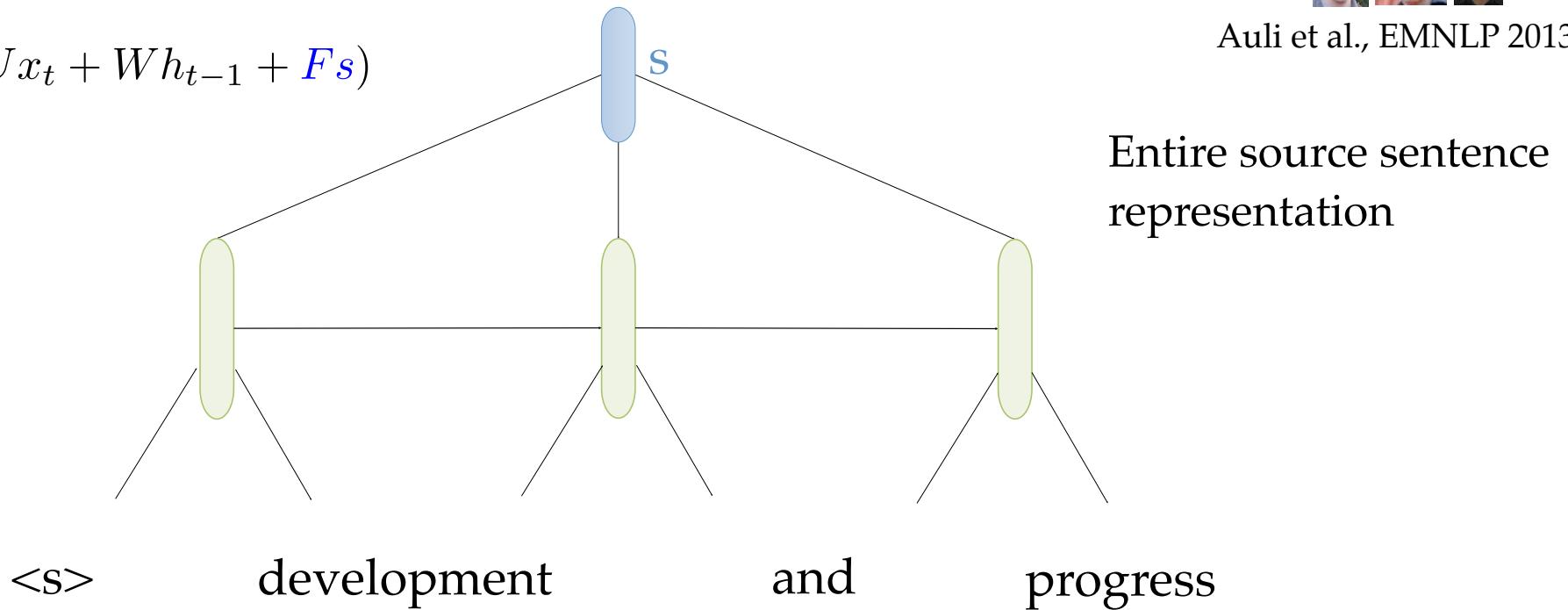
# Combined language and translation model

本 地 区 的 发 展 和 进 步



Auli et al., EMNLP 2013

$$h_t = \sigma(Ux_t + Wh_{t-1} + \mathbf{F}s)$$



# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs}) \left( \begin{array}{c|c|c} \textcolor{blue}{\text{---}} & \textcolor{blue}{\text{---}} & \textcolor{blue}{\text{---}} \\ \hline \textcolor{blue}{\text{---}} & \textcolor{blue}{\text{---}} & \textcolor{blue}{\text{---}} \end{array} \right)$$

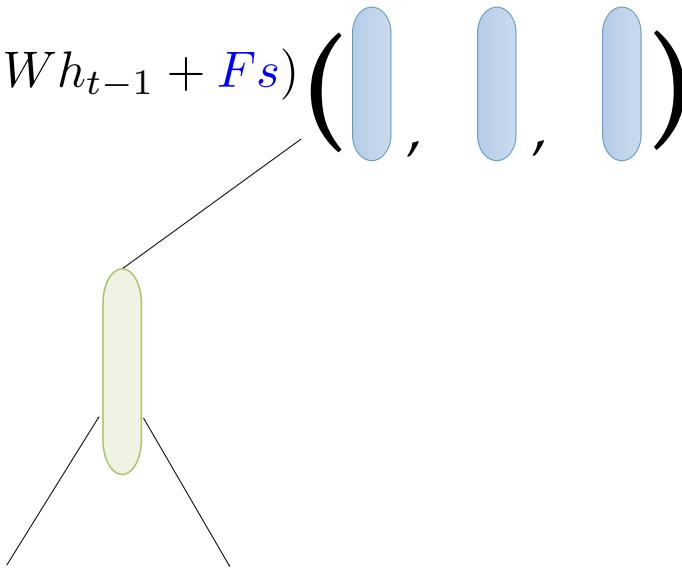
Source word-window

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs}) (, , )$$

Source word-window



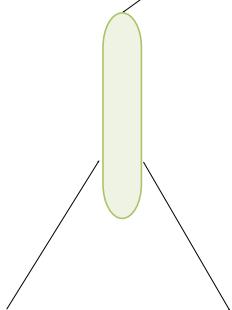
<S>

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs}) (, , )$$

Source word-window



< s >

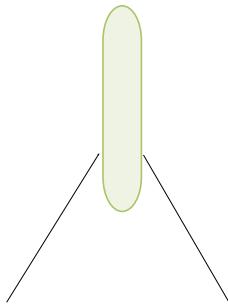
development

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$

Source word-window

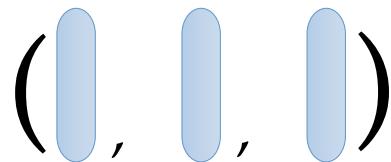


<s> development

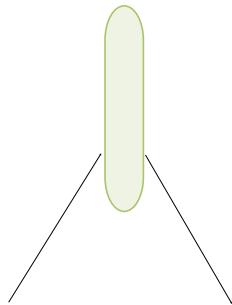
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



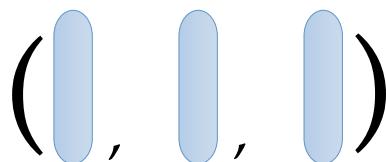
< s >

development

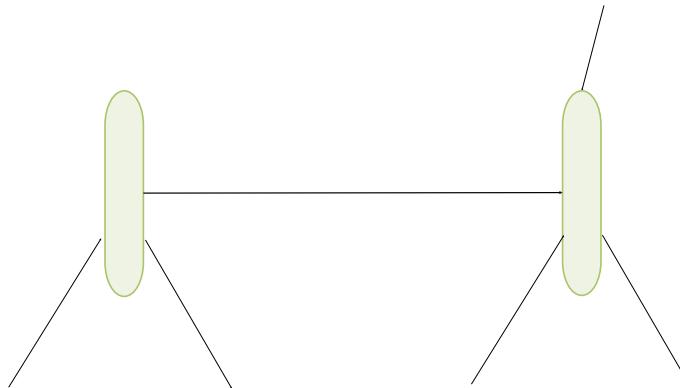
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



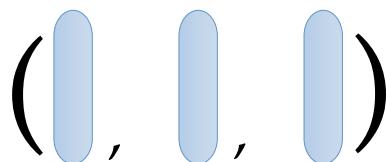
< s >

development

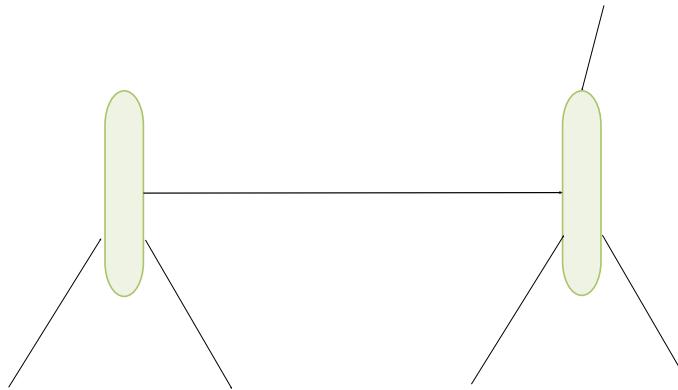
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



< s >

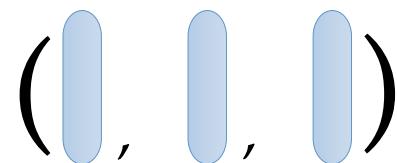
development

and

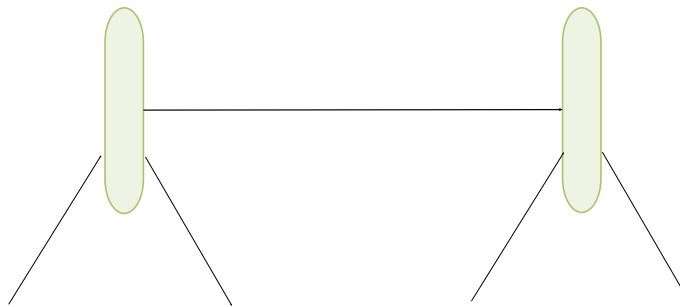
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



< s >

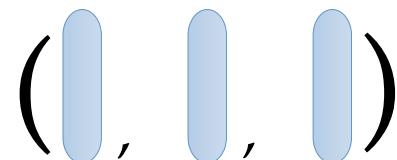
development

and

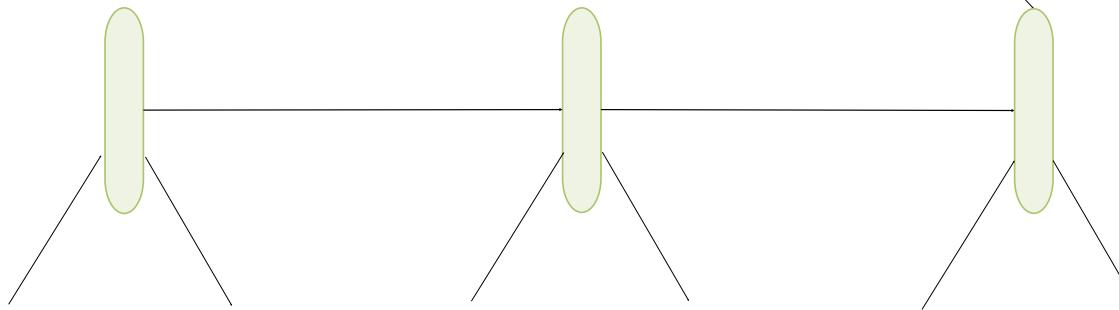
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



< s >

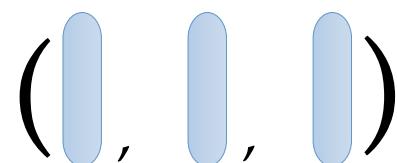
development

and

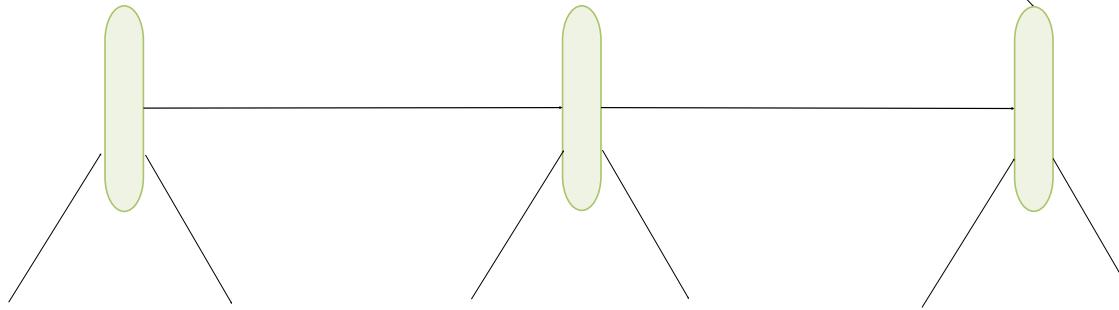
# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Source word-window



$\langle s \rangle$

development

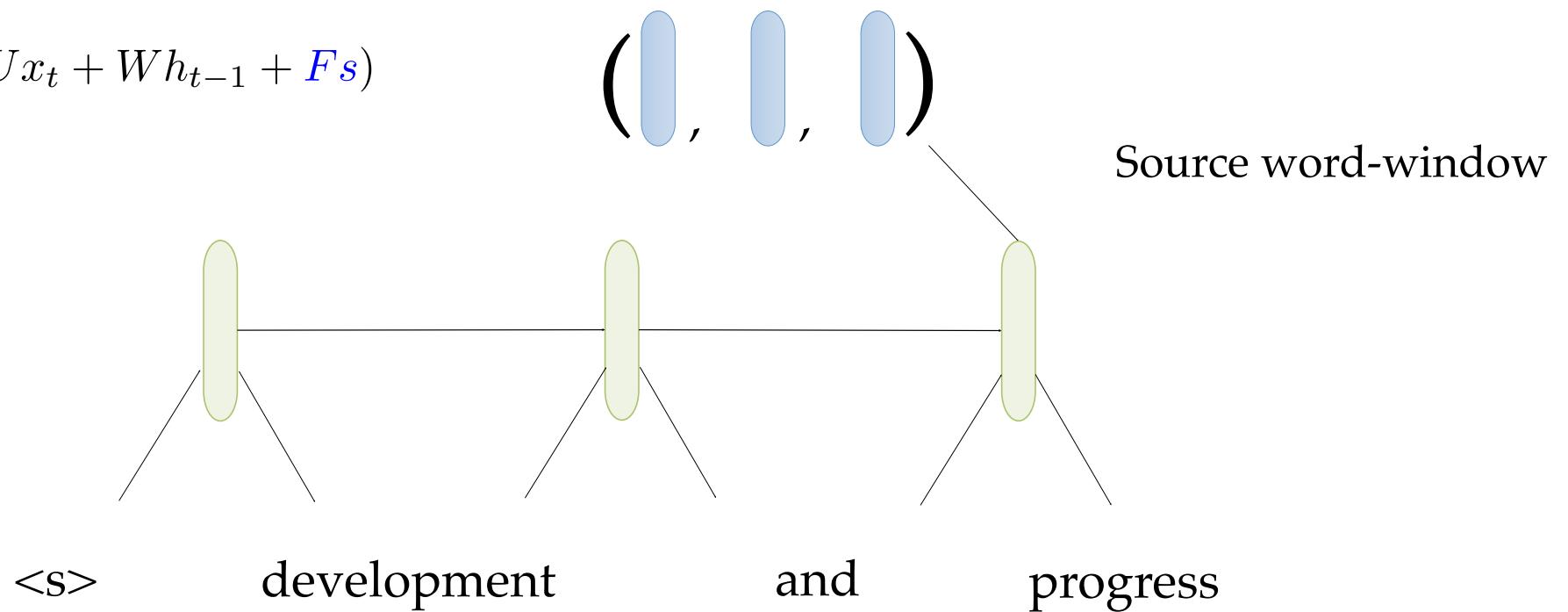
and

progress

# Combined language and translation model

本 地 区 的 发 展 和 进 步

$$h_t = \sigma(Ux_t + Wh_{t-1} + \textcolor{blue}{Fs})$$



Le (2012), Kalchbrenner (2013), Devlin (2014), Cho (2014), Sutskever (2014)

# Experimental setup

- WMT 2012 French-English translation task
- Data: 100M words
- Baseline: Phrase-based model similar to Moses
- Rescoring
- Mini-batch gradient descent
- Class-structured output layer (Goodman, 1996)

# Does the neural model learn to translate?

-Discrete translation model

WMT 2012 French-English, 100M words, phrase-based baseline, n-best rescoring

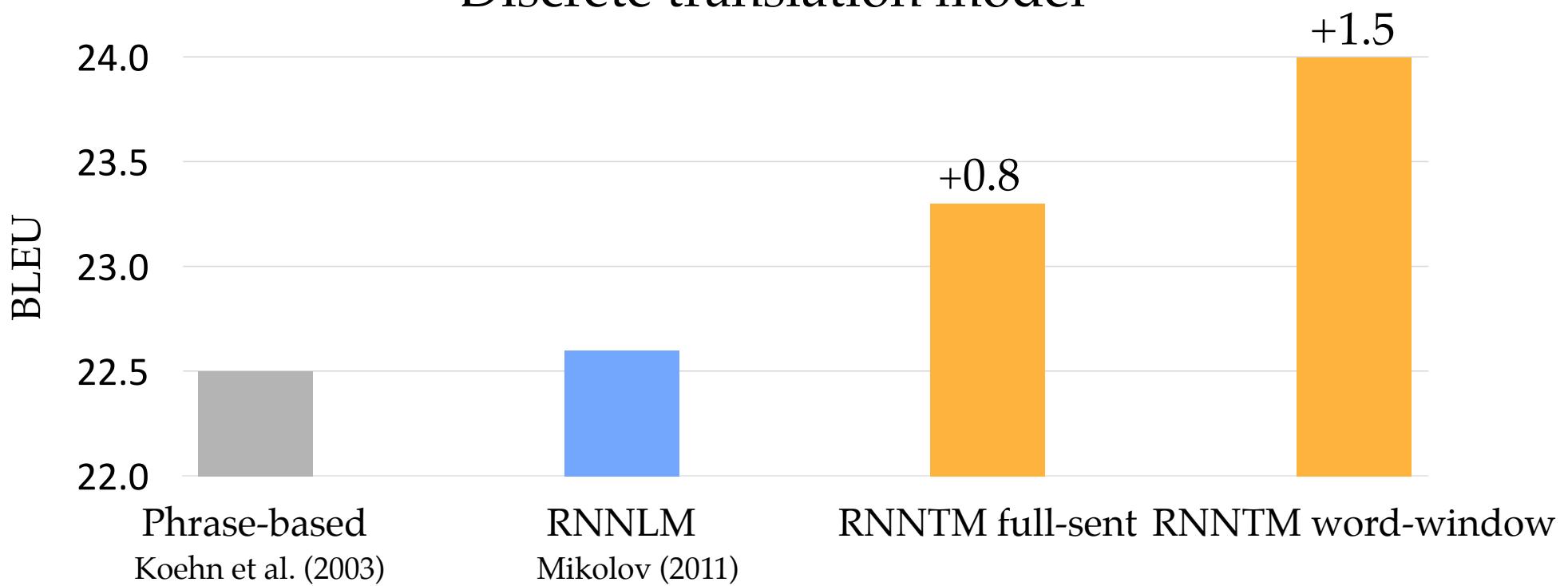
# Does the neural model learn to translate?

-Discrete translation model



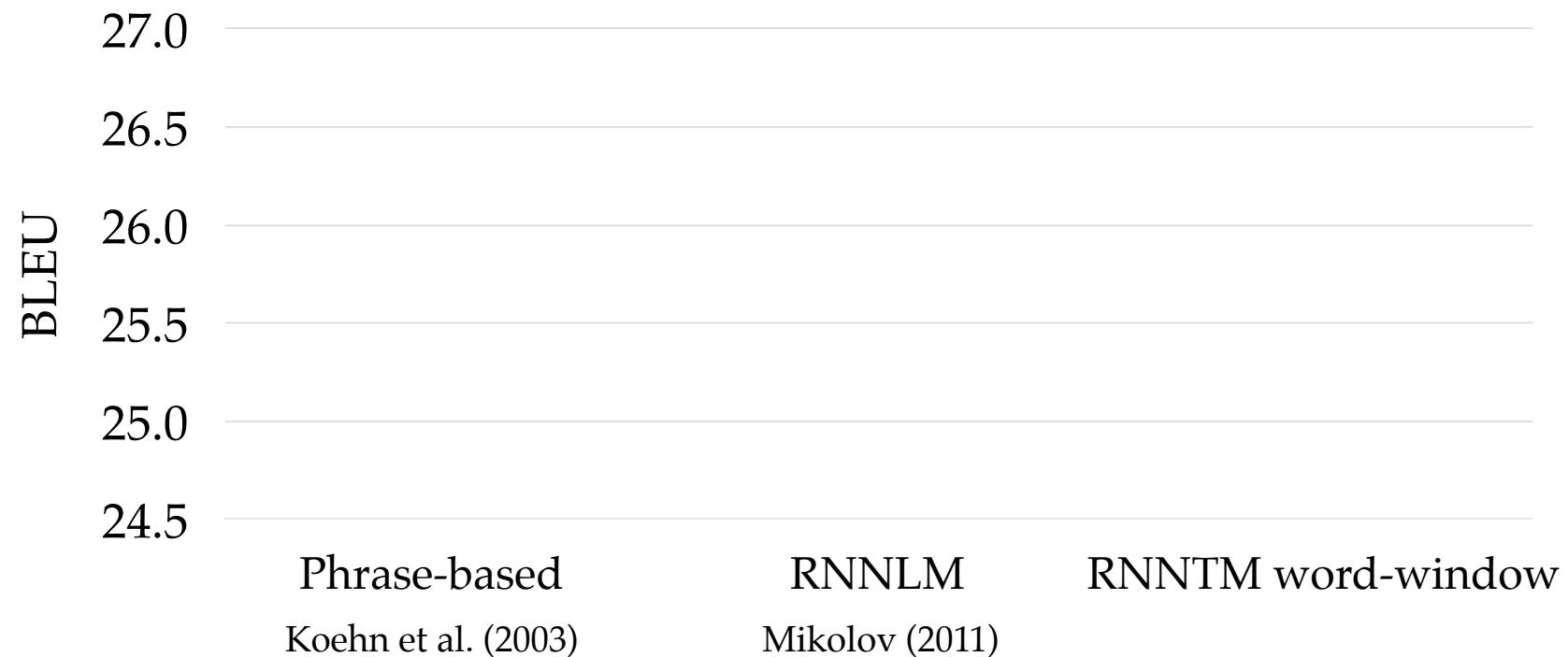
WMT 2012 French-English, 100M words, phrase-based baseline, n-best rescoring

# Does the neural model learn to translate? -Discrete translation model



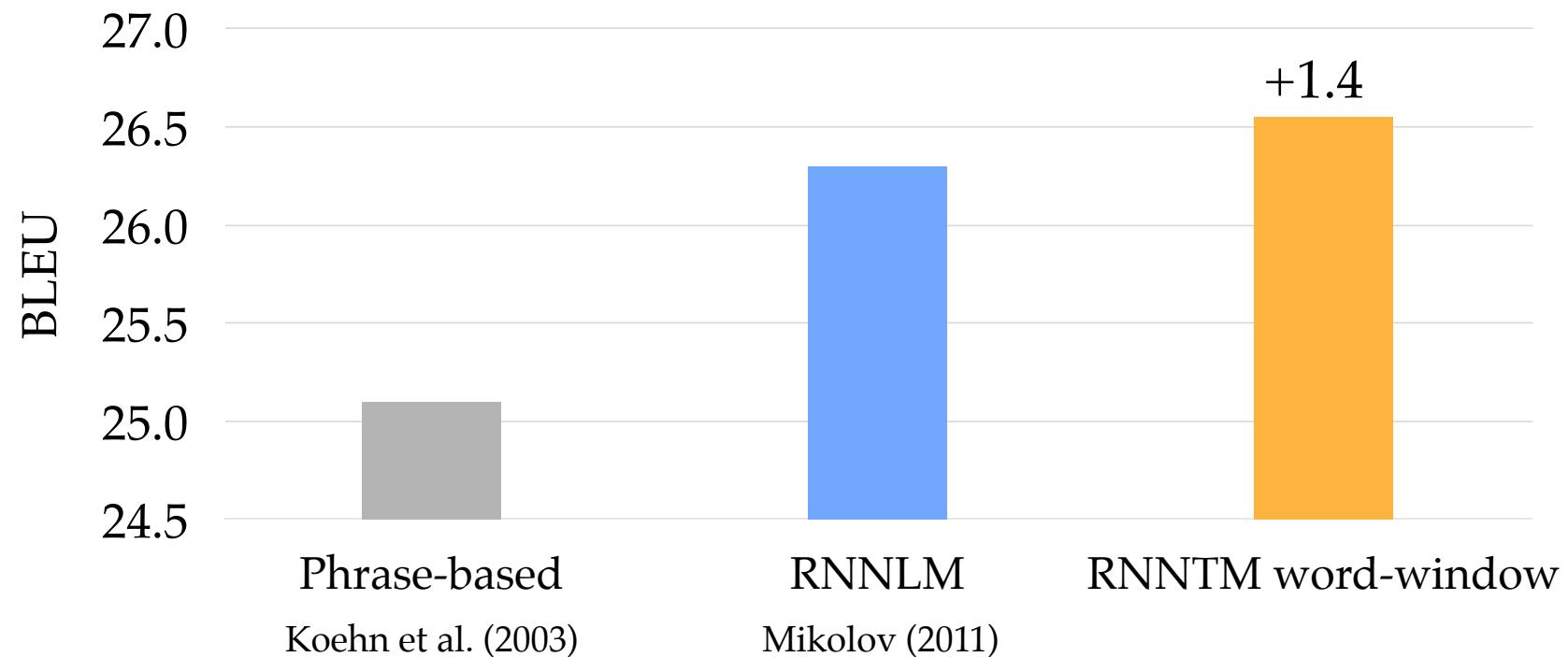
WMT 2012 French-English, 100M words, phrase-based baseline, n-best rescoring

# Improving a phrase-based baseline +Discrete translation model



WMT 2012 French-English, 100M words, phrase-based baseline, lattice rescoring

# Improving a phrase-based baseline +Discrete translation model



WMT 2012 French-English, 100M words, phrase-based baseline, lattice rescoring

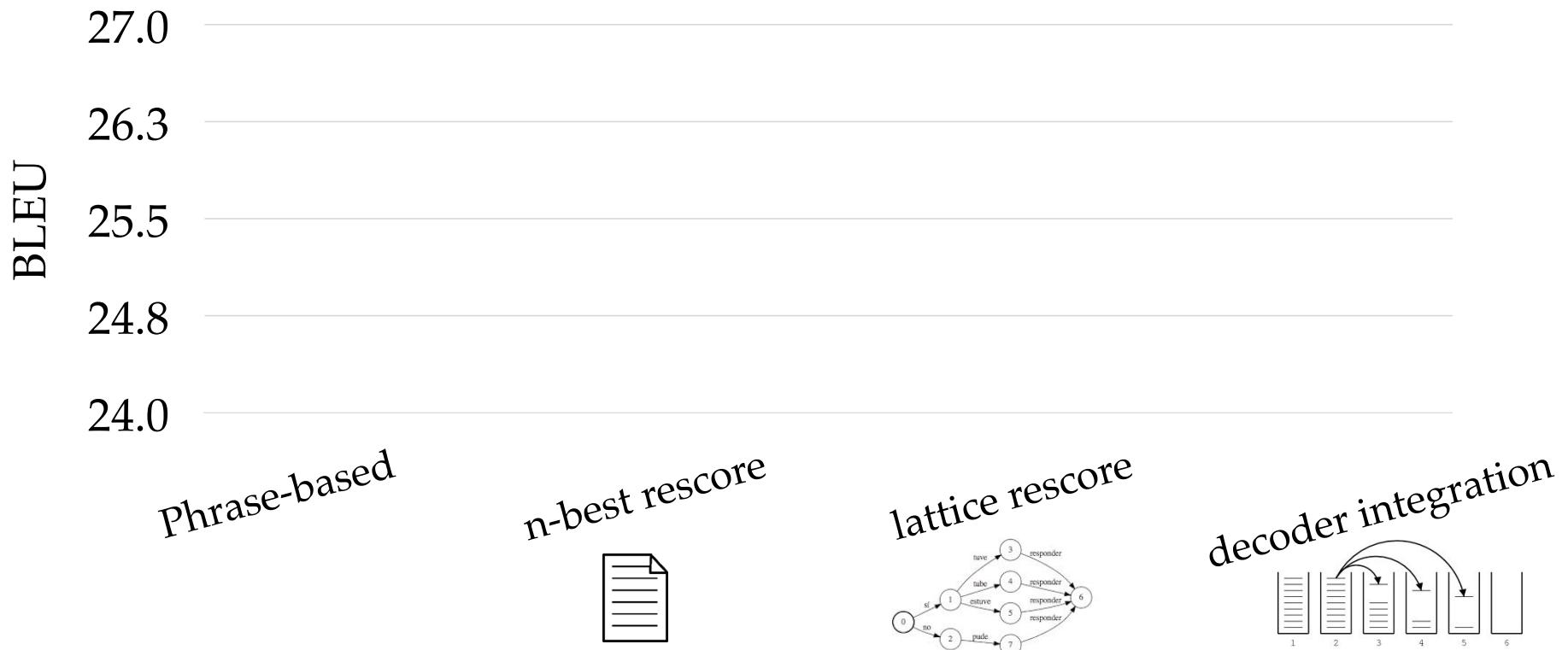
# Combining recurrent nets with discrete models

Auli & Gao, ACL 2014



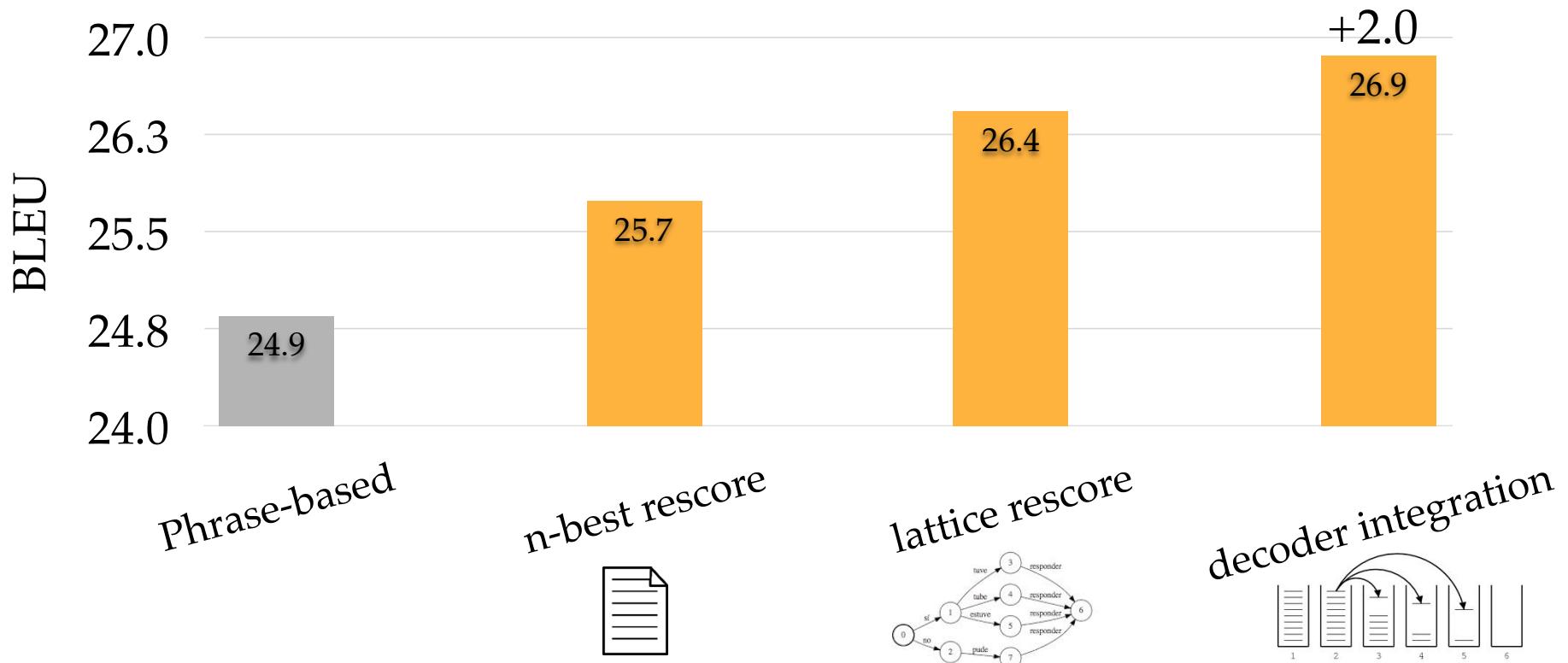
# Combining recurrent nets with discrete models

Auli & Gao, ACL 2014

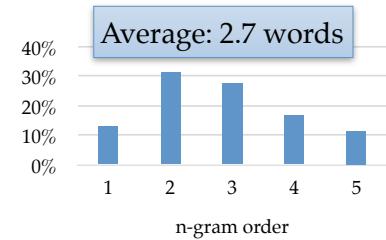


# Combining recurrent nets with discrete models

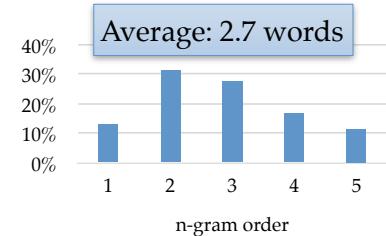
Auli & Gao, ACL 2014



# Neural nets vs discrete models

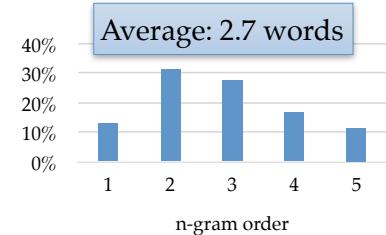


# Neural nets vs discrete models



- Neural models more robust when discrete models rely on sparse estimates?

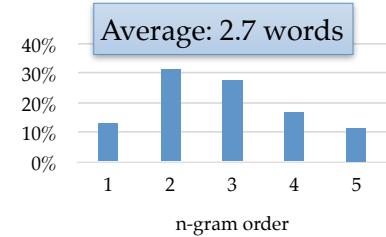
# Neural nets vs discrete models



- Neural models more robust when discrete models rely on sparse estimates?
- Language modeling: **n-gram LM** and **neural net LM** two components of log-linear model of translation

$$\hat{e} = \operatorname{argmax}_e \sum_i \lambda_i h_i(f, e) \quad h_1(f, e) = \log p_{lm}(e) \quad h_2(f, e) = \log p_{rnn}(e)$$

# Neural nets vs discrete models



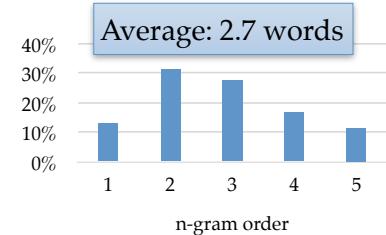
- Neural models more robust when discrete models rely on sparse estimates?
- Language modeling: **n-gram LM** and **neural net LM** two components of log-linear model of translation

$$\hat{e} = \operatorname{argmax}_e \sum_i \lambda_i h_i(f, e) \quad h_1(f, e) = \log p_{lm}(e) \quad h_2(f, e) = \log p_{rnn}(e)$$

- Split each model into five features, one for each n-gram order s.t.

$$\log p_{lm}(e) = \sum_{i=1}^5 h_i(f, e) \quad \log p_{rnn}(e) = \sum_{i=6}^{10} h_i(f, e)$$

# Neural nets vs discrete models



- Neural models more robust when discrete models rely on sparse estimates?
- Language modeling: **n-gram LM** and **neural net LM** two components of log-linear model of translation

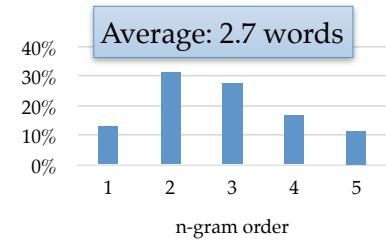
$$\hat{e} = \operatorname{argmax}_e \sum_i \lambda_i h_i(f, e) \quad h_1(f, e) = \log p_{lm}(e) \quad h_2(f, e) = \log p_{rnn}(e)$$

- Split each model into five features, one for each n-gram order s.t.

$$\log p_{lm}(e) = \sum_{i=1}^5 h_i(f, e) \quad \log p_{rnn}(e) = \sum_{i=6}^{10} h_i(f, e)$$

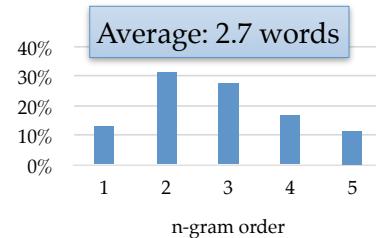
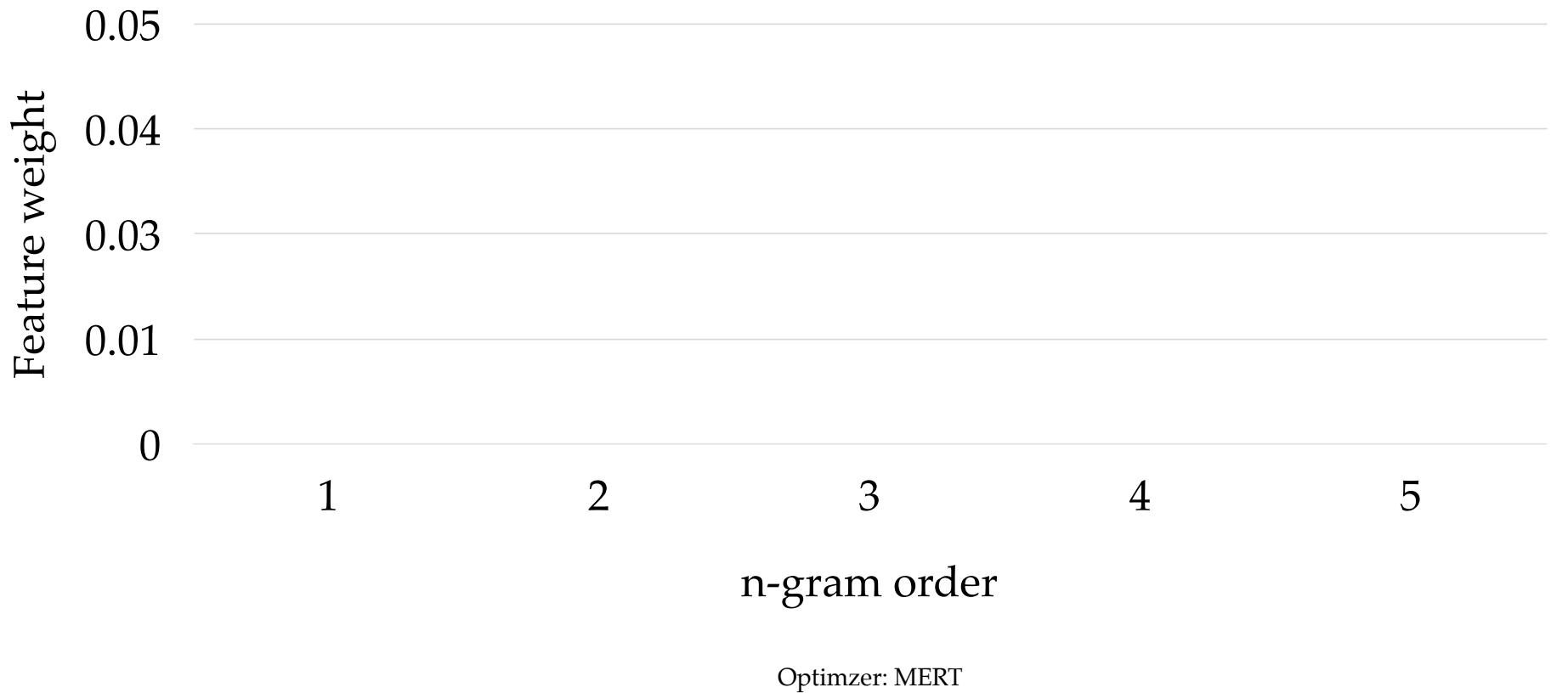
- Standard optimizer (MERT) to find weights for each n-gram order

# Neural nets vs discrete models



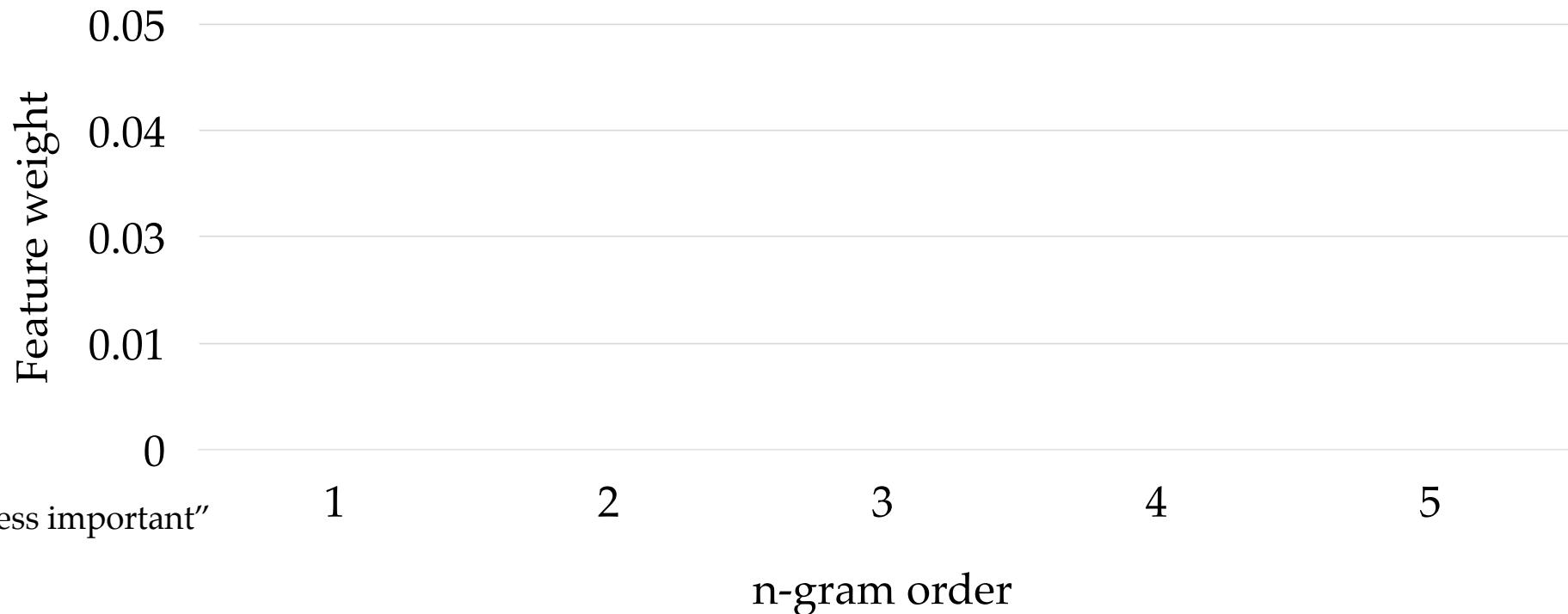
Optimizer: MERT

# Neural nets vs discrete models



# Neural nets vs discrete models

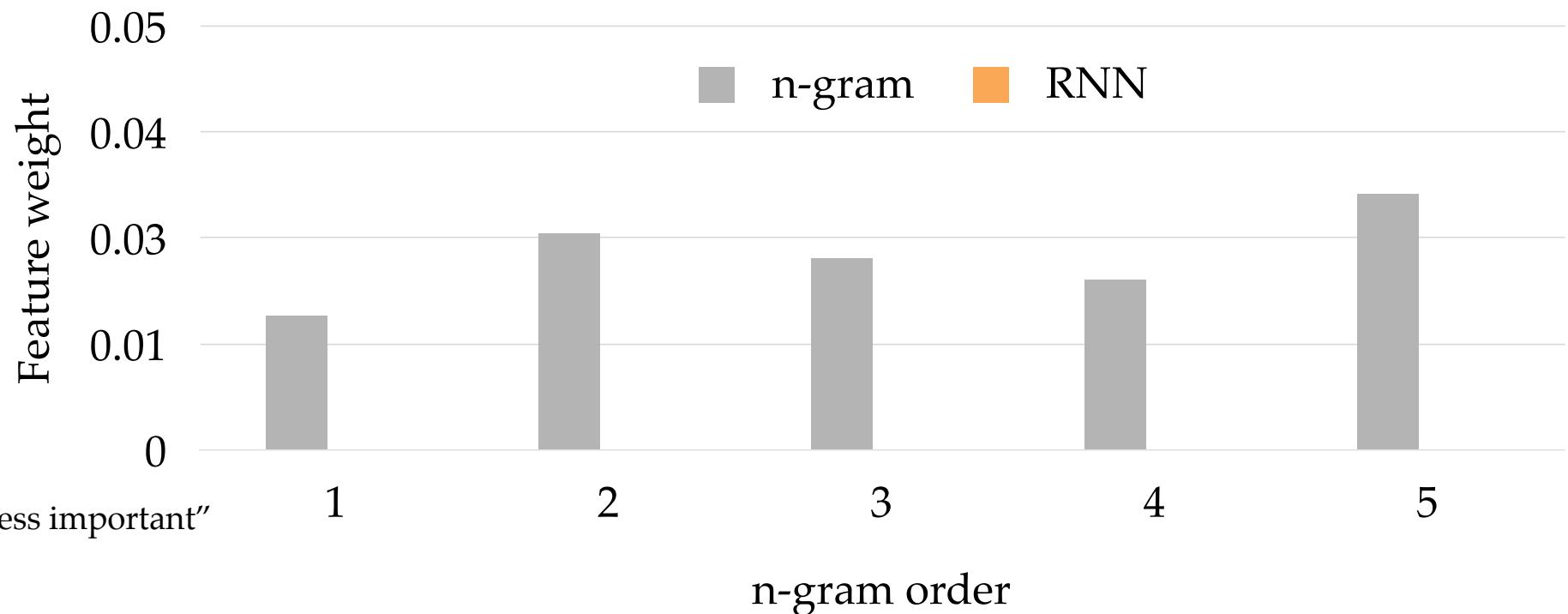
“more important”



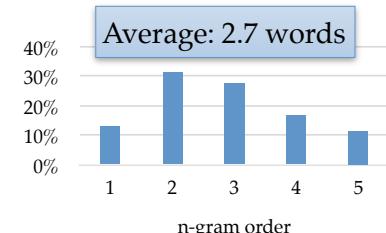
Optimizer: MERT

# Neural nets vs discrete models

“more important”

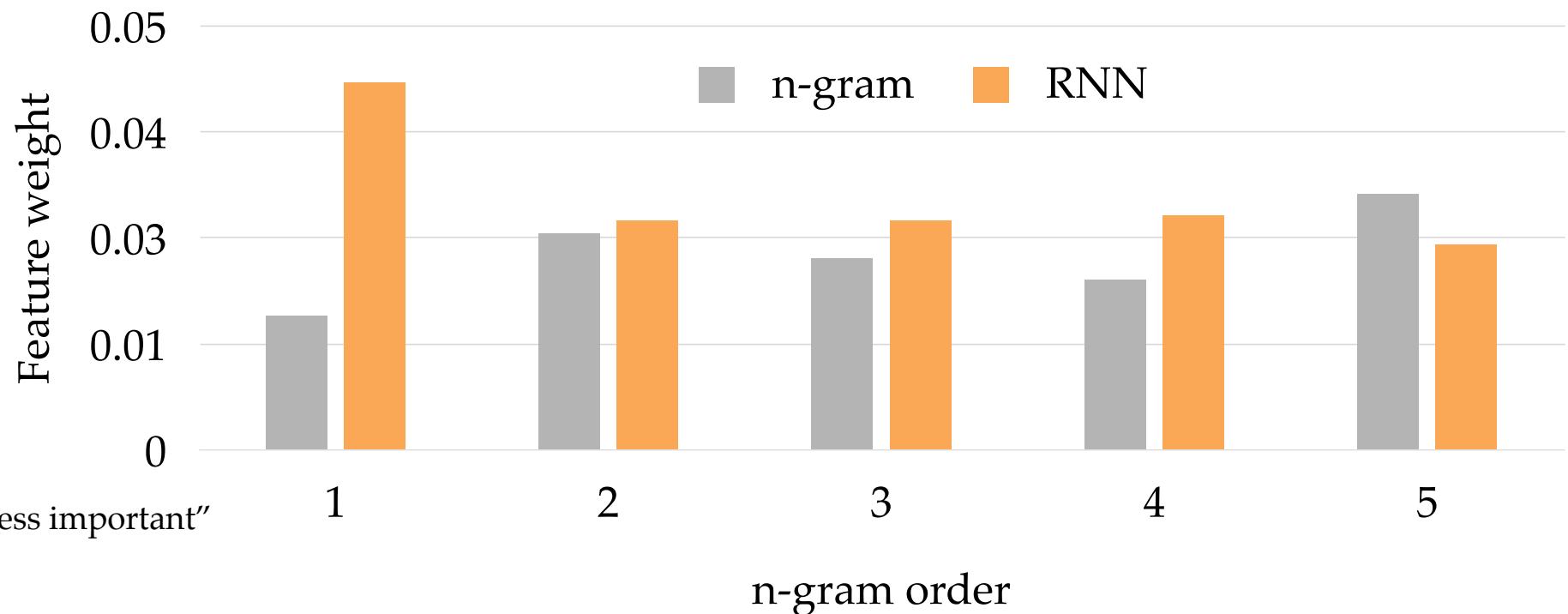


Optimizer: MERT

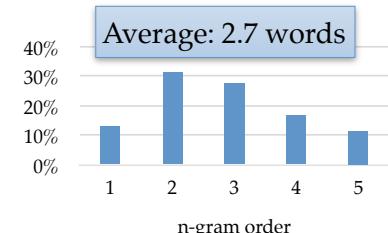


# Neural nets vs discrete models

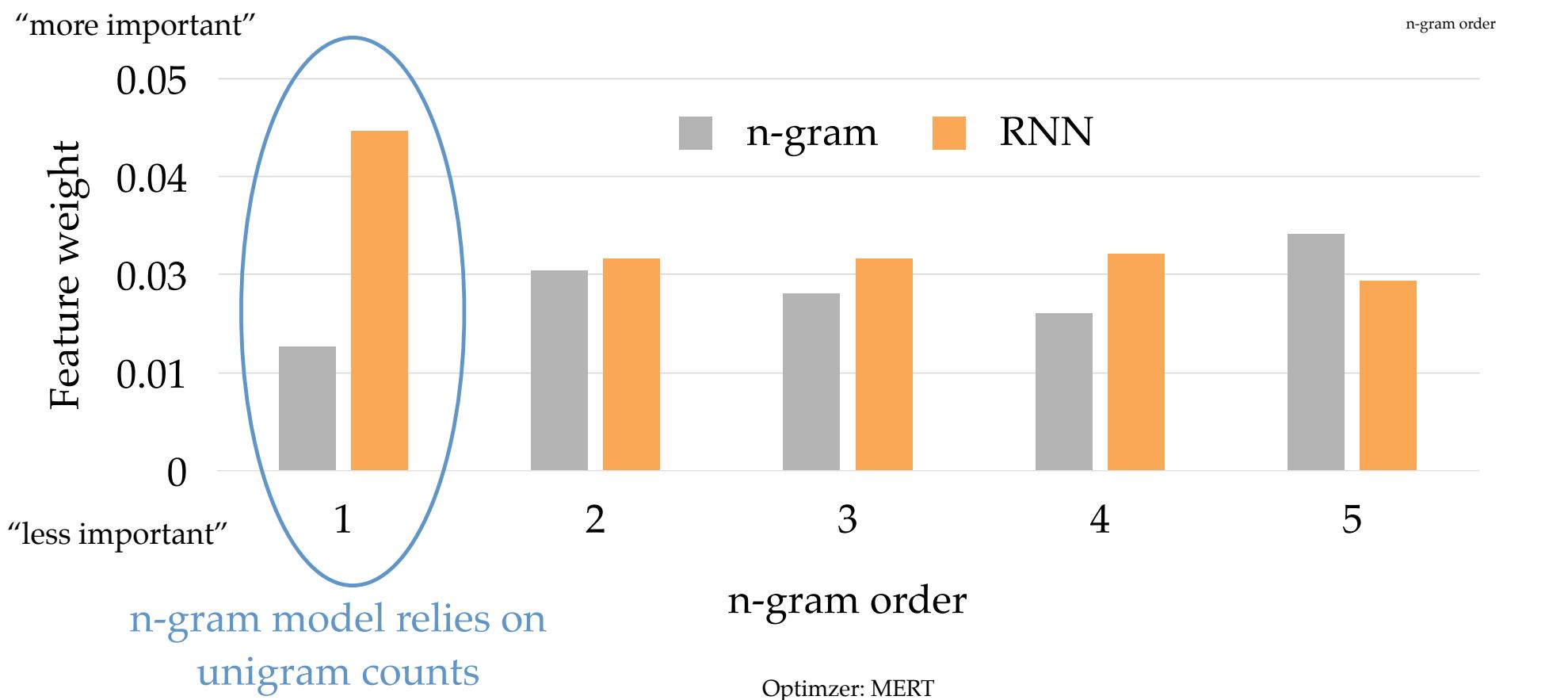
“more important”



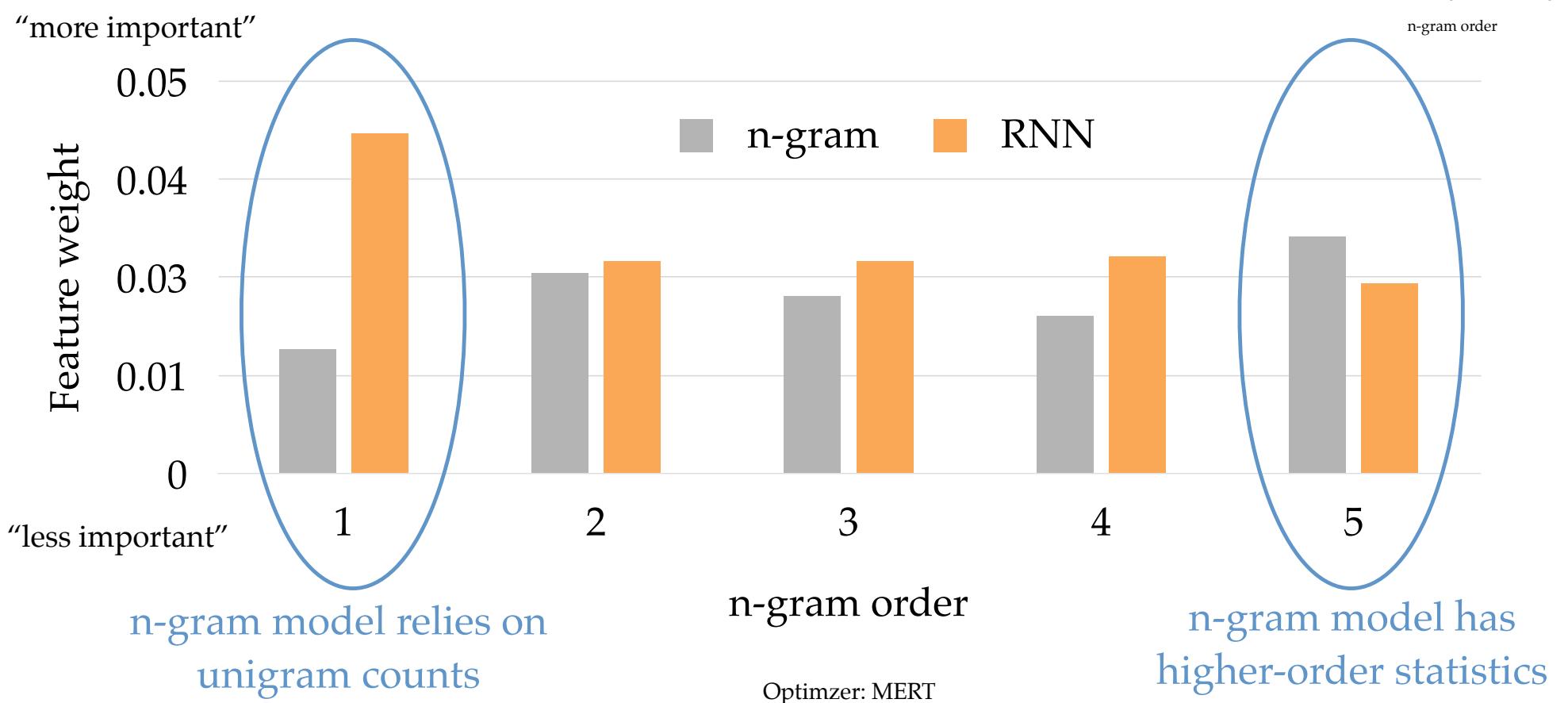
Optimizer: MERT



# Neural nets vs discrete models



# Neural nets vs discrete models

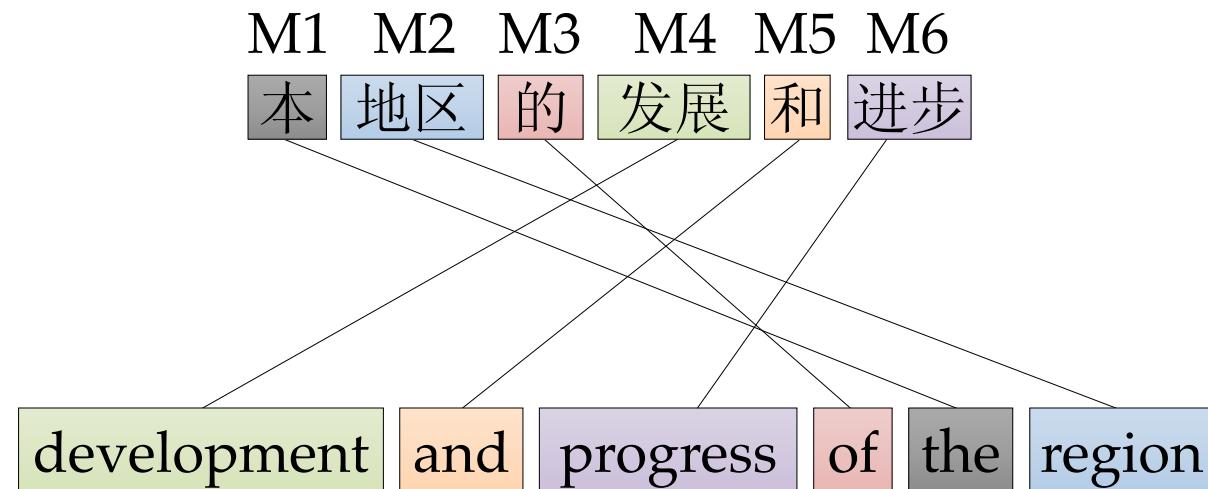


# Recurrent Minimum Translation Unit Models

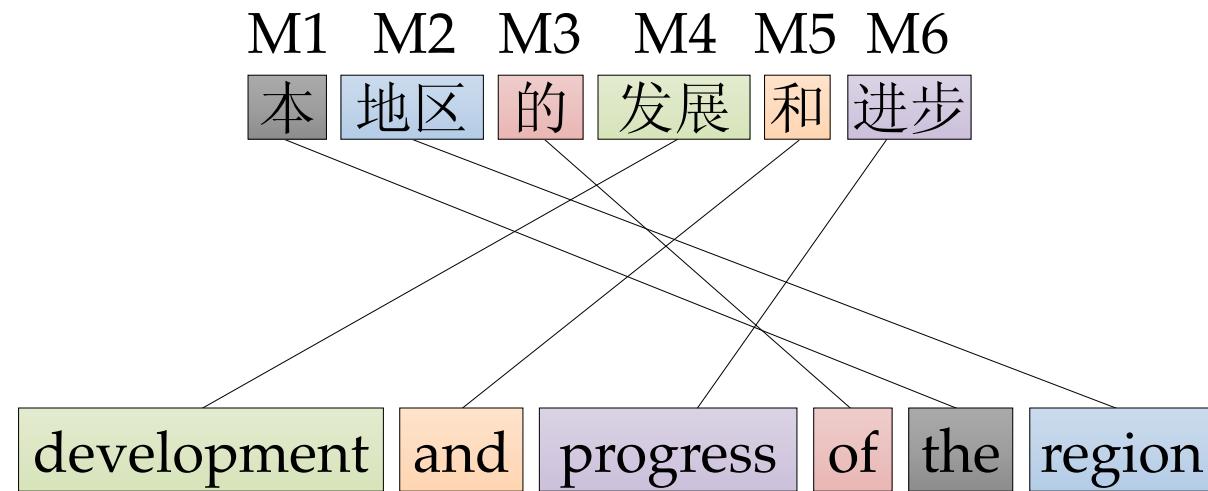
本 地 区 的 发 展 和 进 步

development and progress of the region

# Recurrent Minimum Translation Unit Models

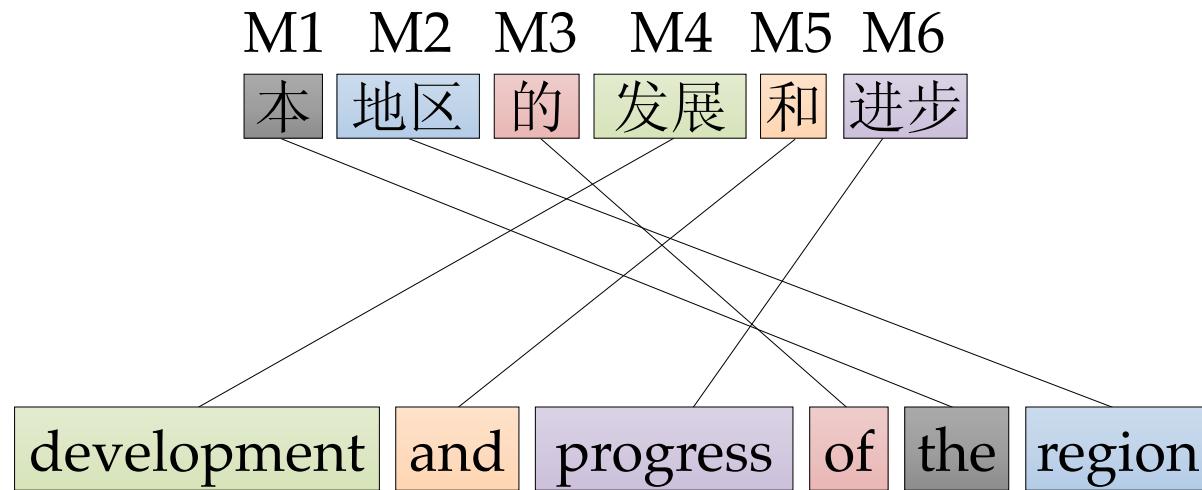


# Recurrent Minimum Translation Unit Models



M1      M2      M3  
本      地区      的  
the      region      of      ...

# Recurrent Minimum Translation Unit Models



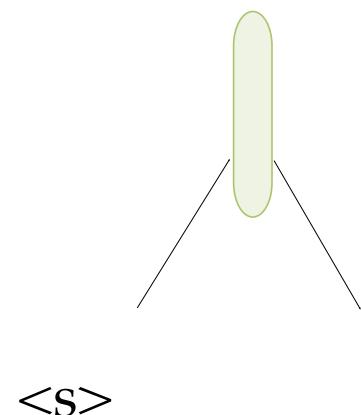
M1      M2      M3  
本      地区      的  
the      region      of      ...

n-gram models over MTUs:  
 $p(M1) \ p(M2 | M1) \ p(M3 | M1, M2) \dots$

Banchs et al. (2005), Quirk & Menezes (2006)

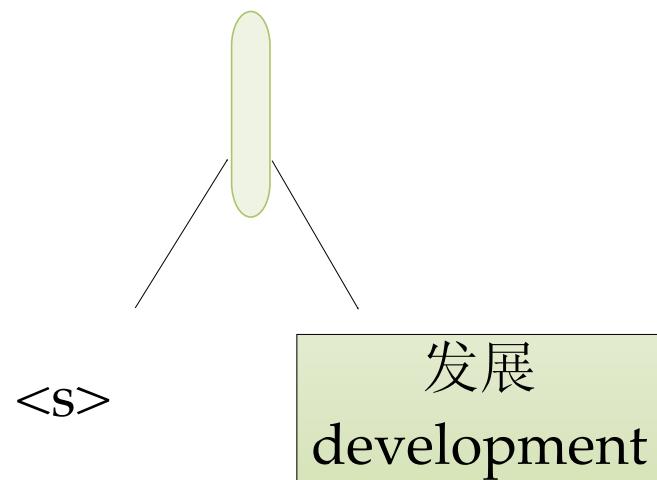
# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014



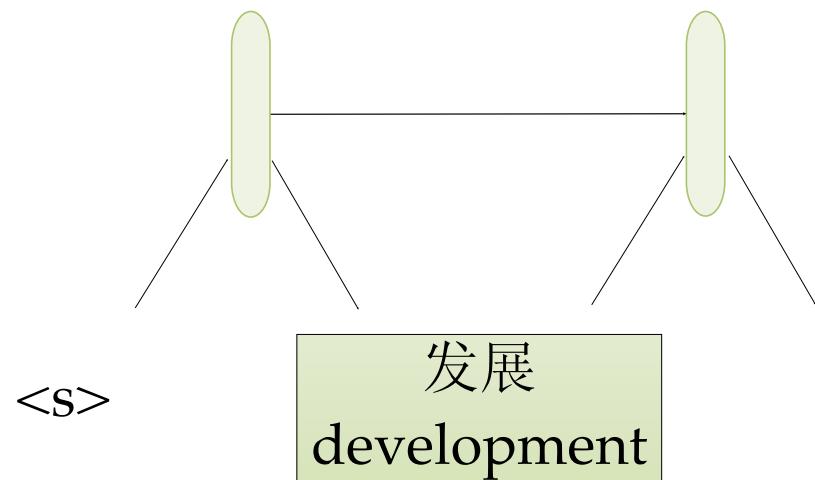
# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014



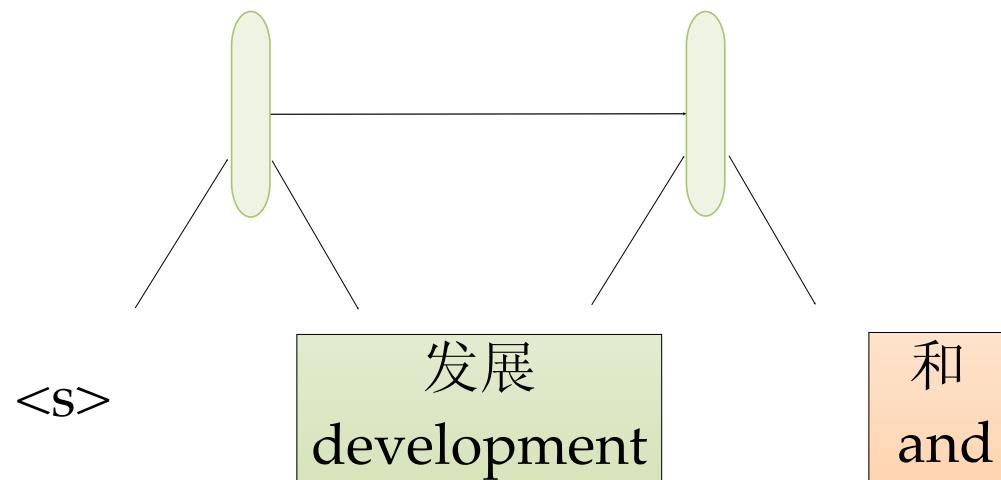
# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014



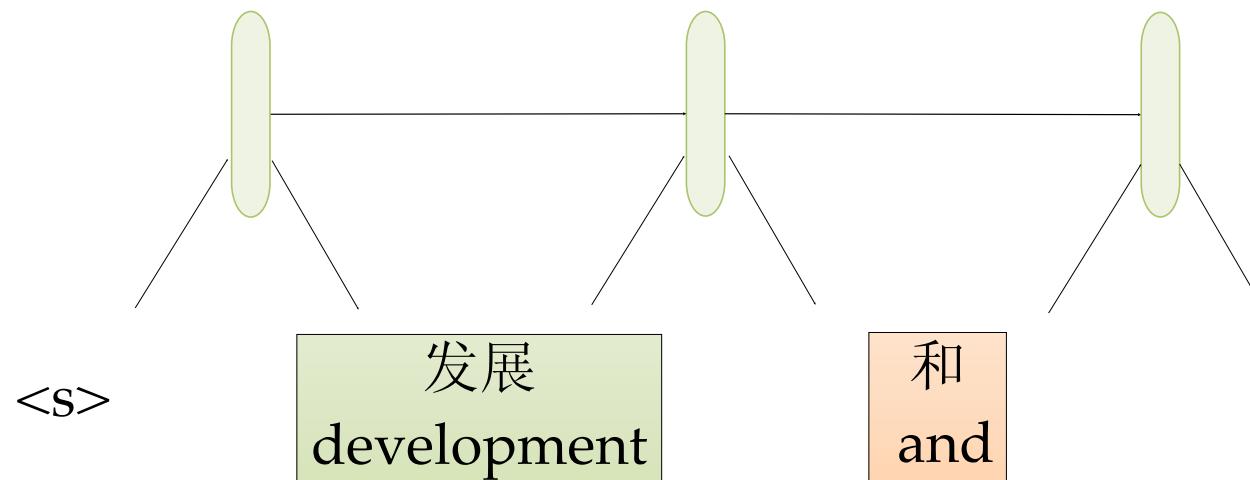
# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014



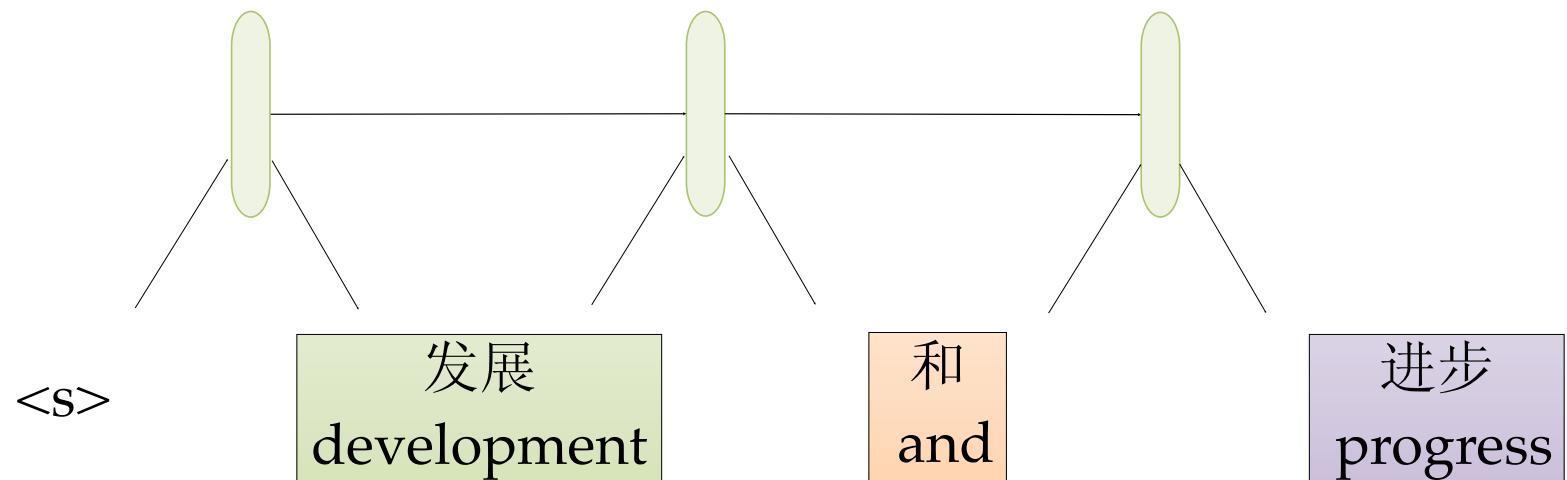
# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014

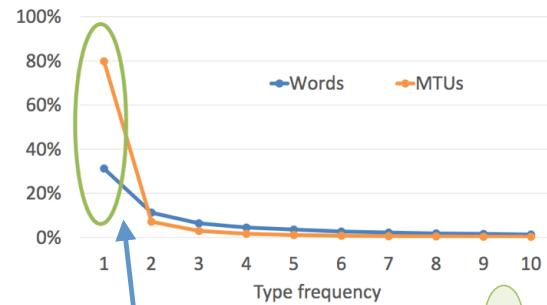


# Recurrent Minimum Translation Unit Models

Hu et al., EACL 2014



# Recurrent Minimum Translation Unit Models



Singletons

<S>

发展  
development

和  
and

进步  
progress

Hu et al., EACL 2014



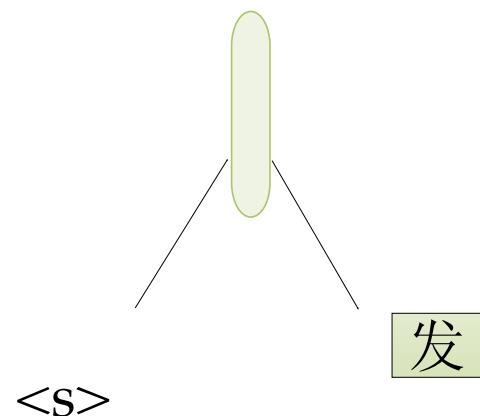
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



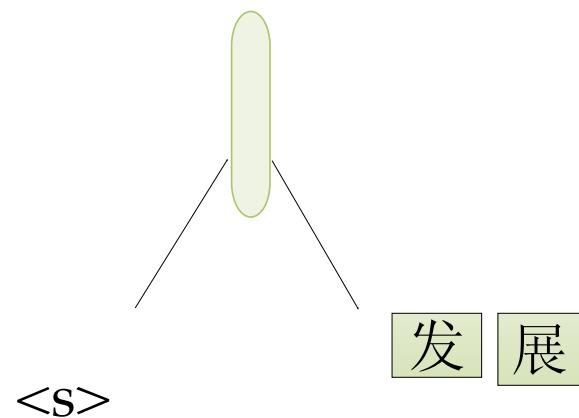
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



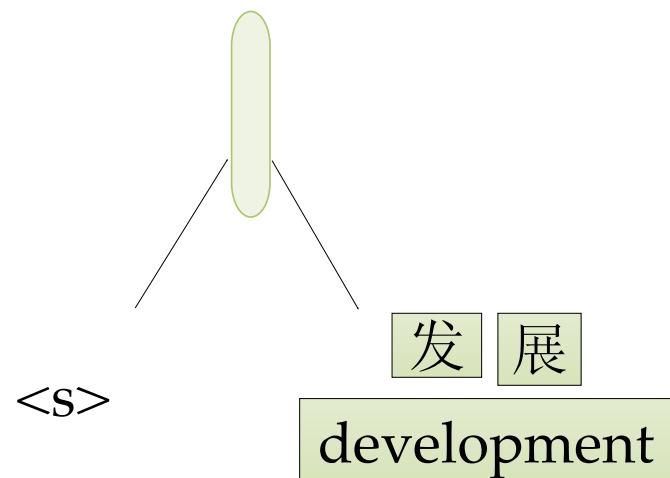
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



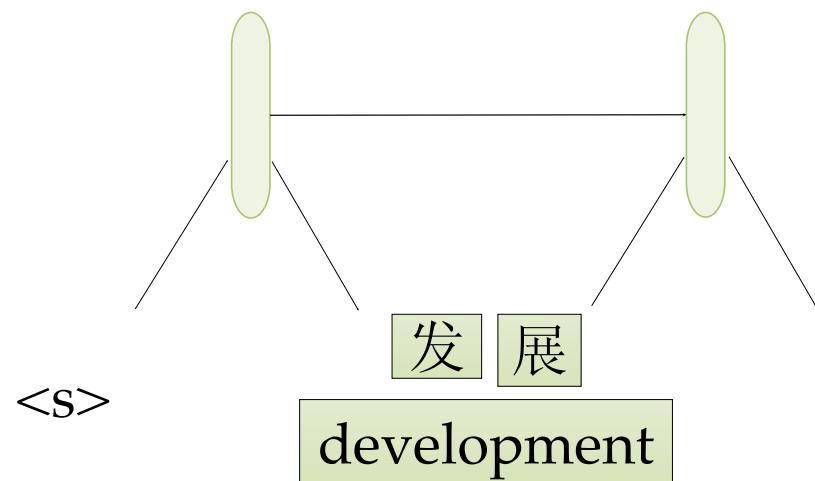
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



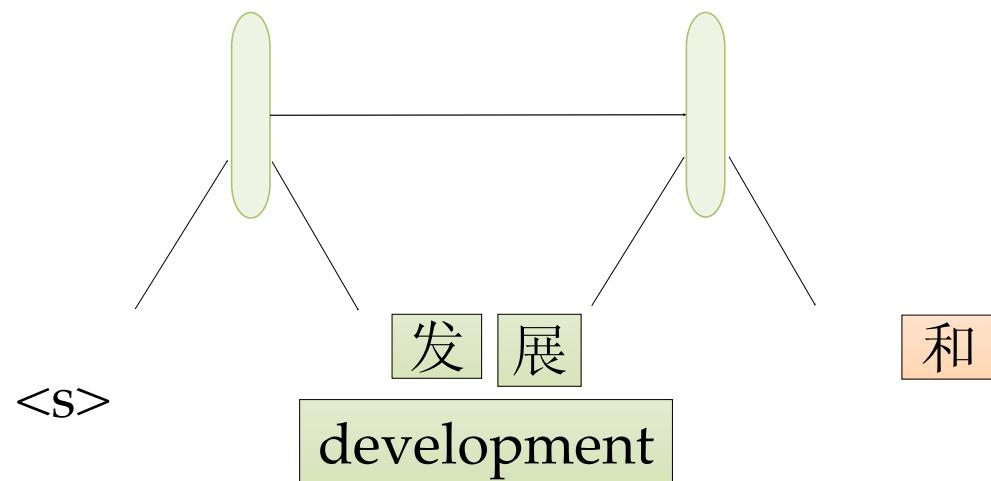
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



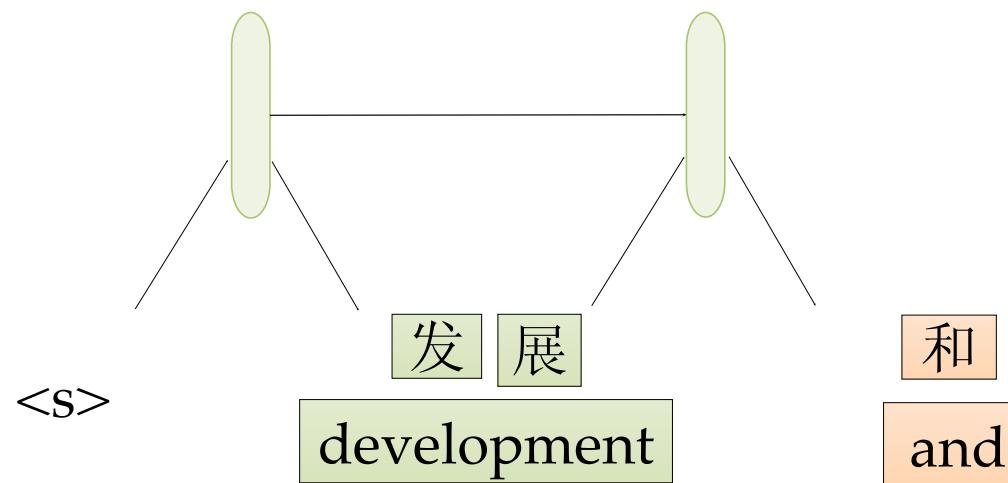
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



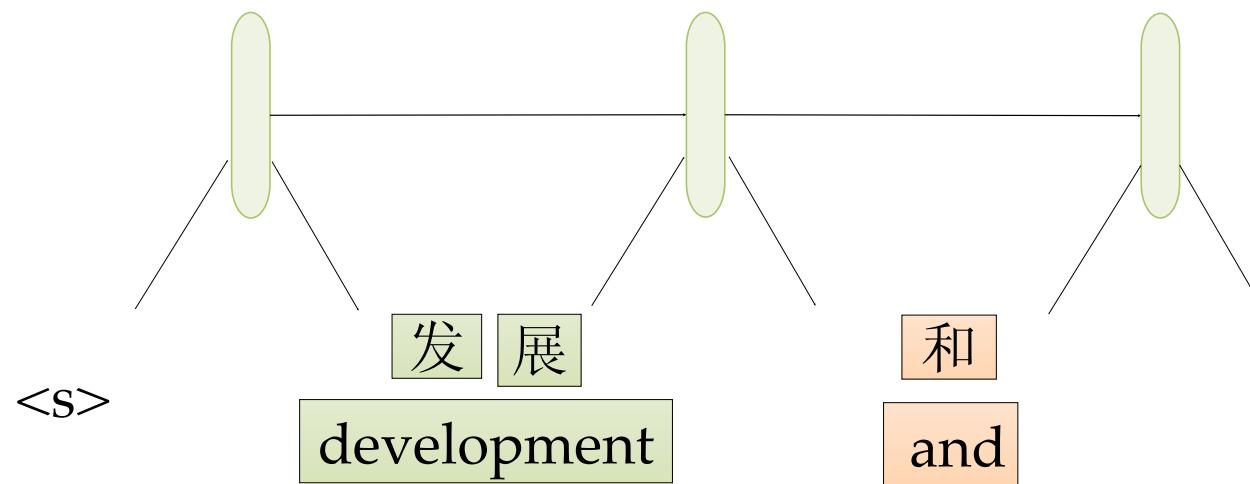
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



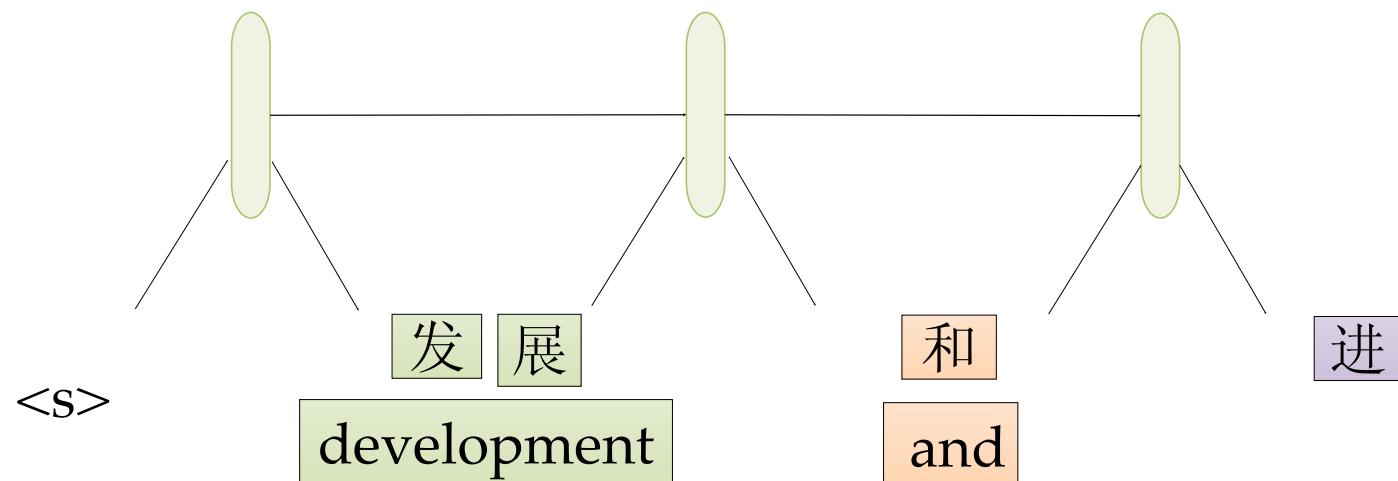
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



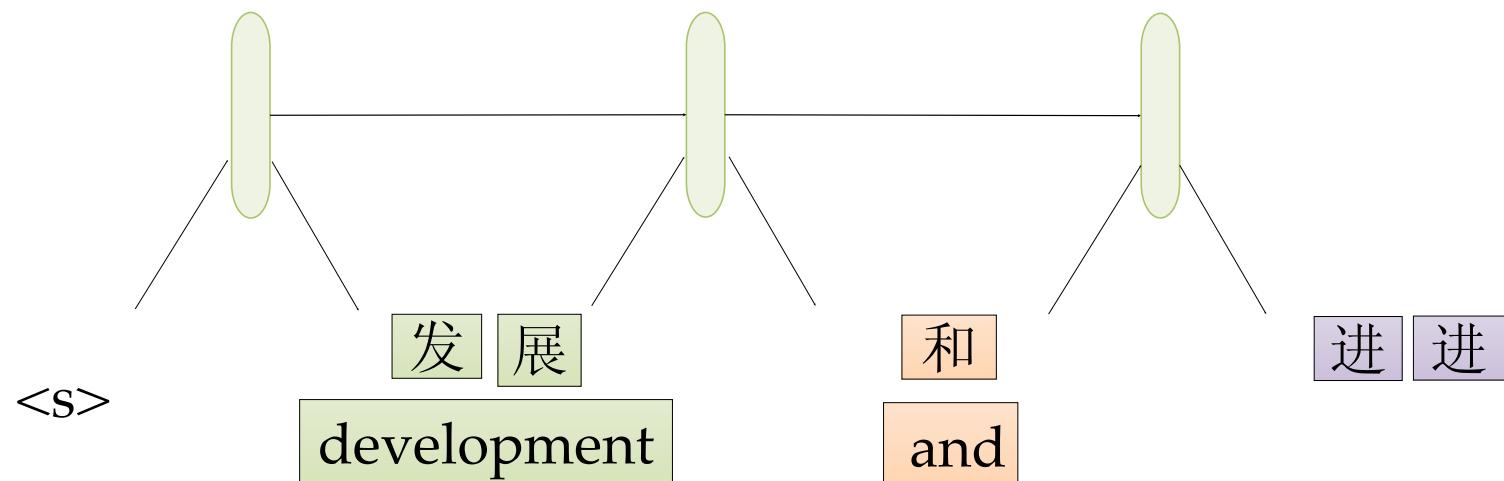
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



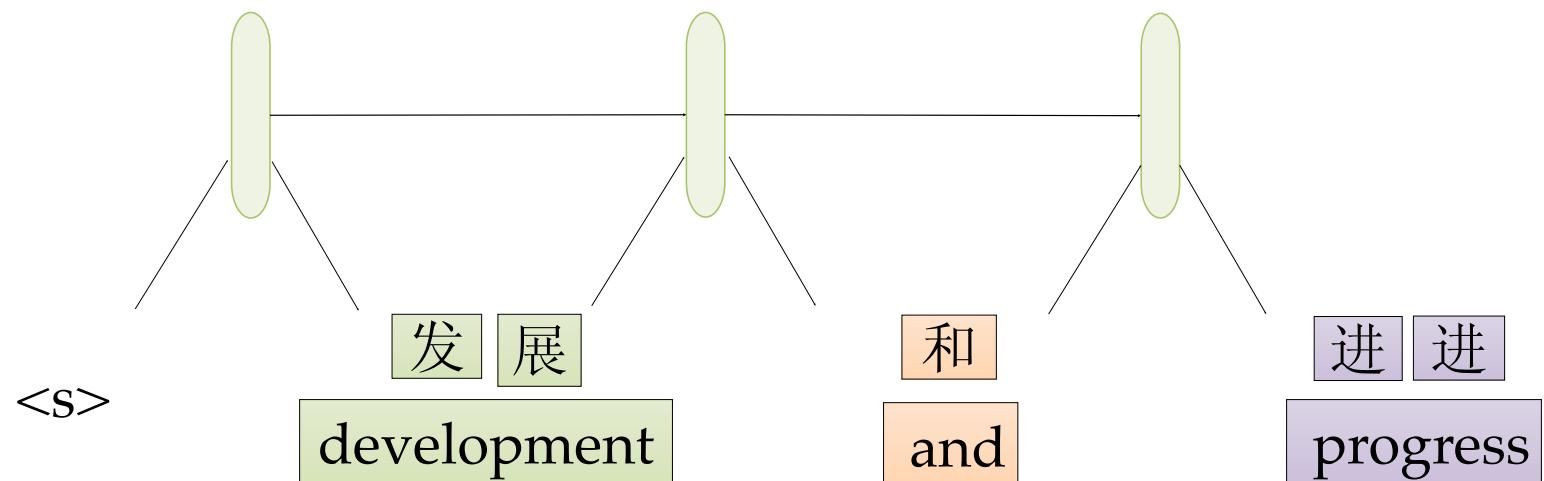
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation



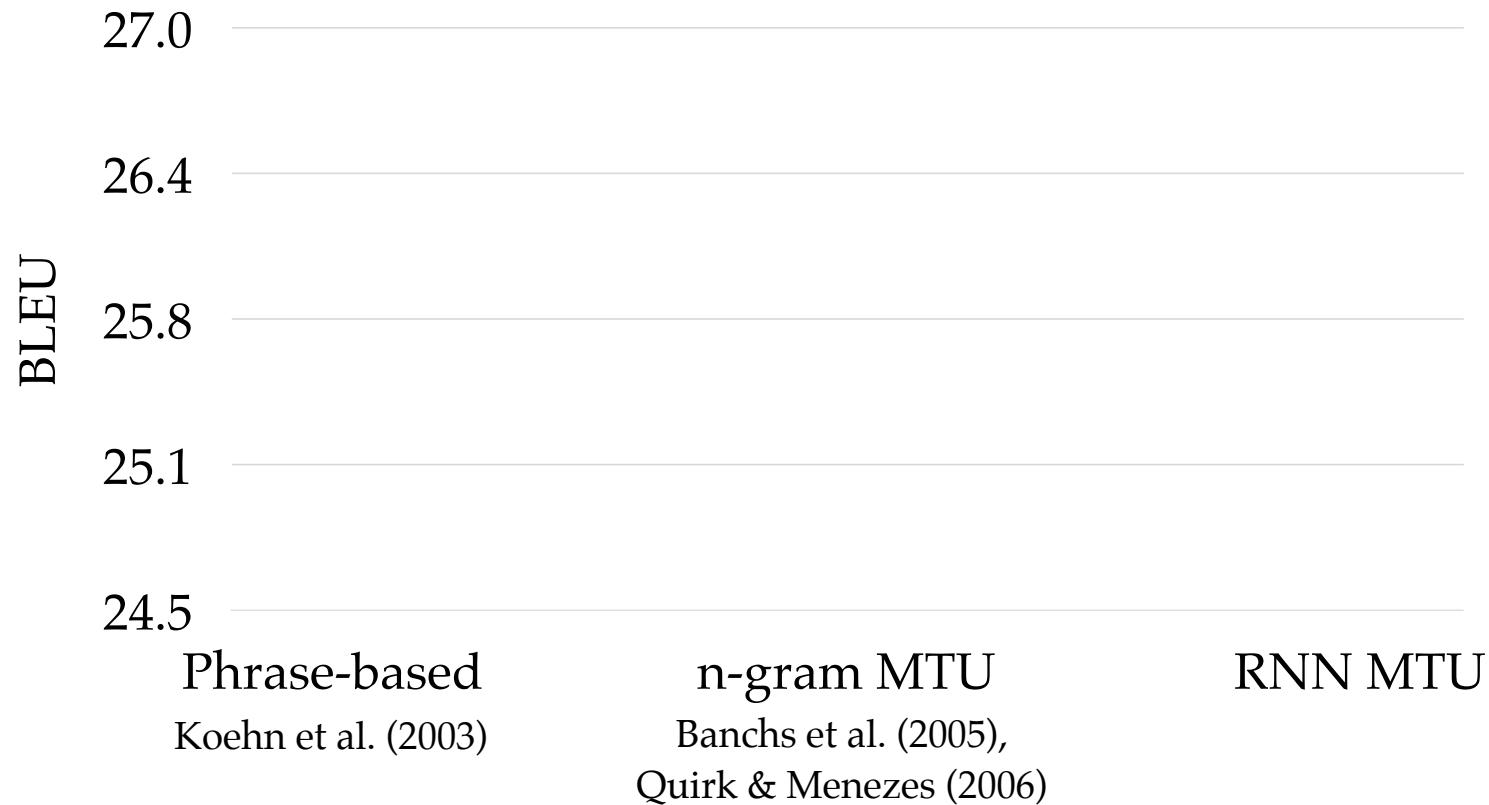
# Recurrent Minimum Translation Unit Models

Reduce sparsity by bag of words representation

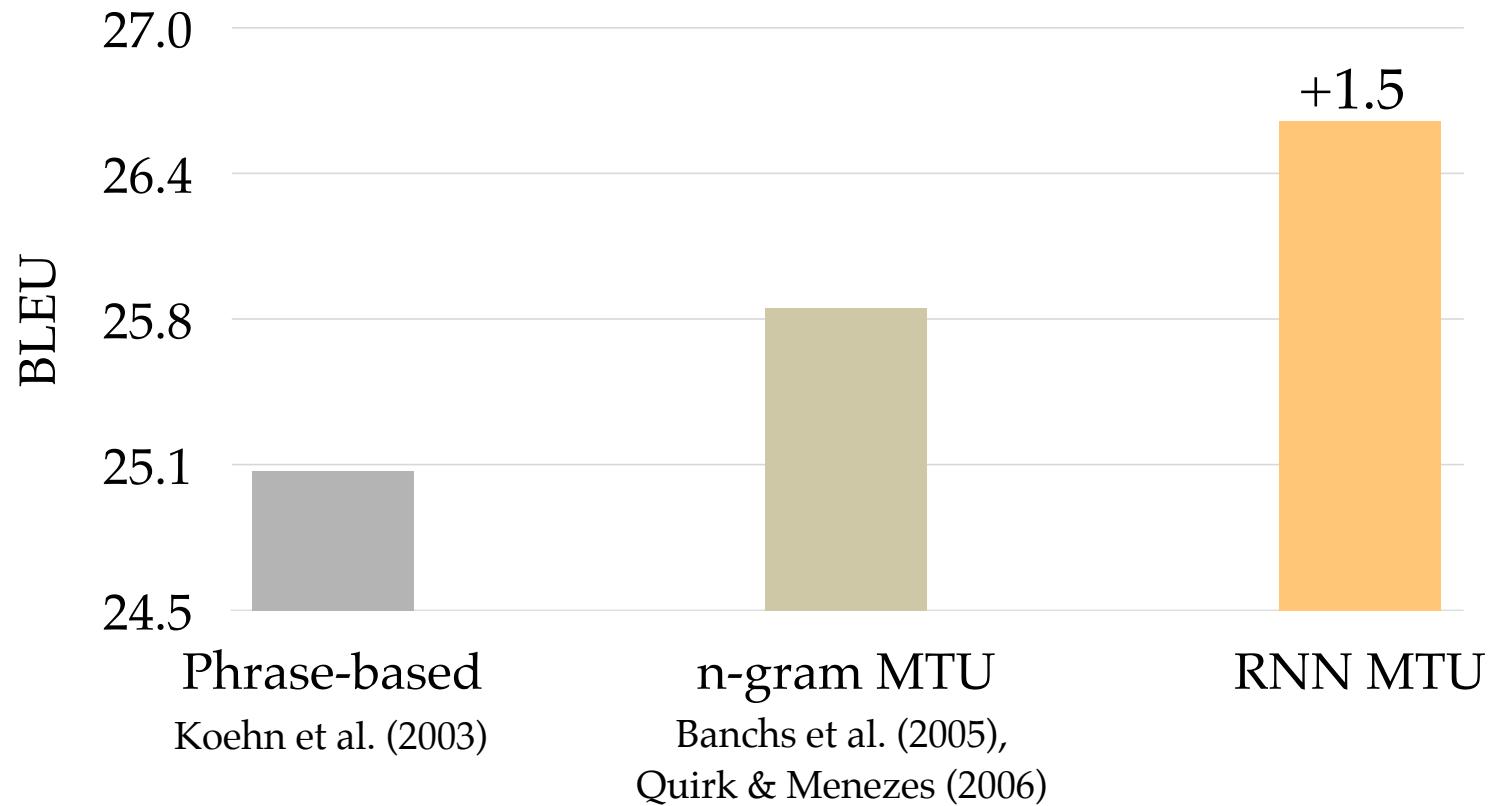


# Recurrent Minimum Translation Unit Models

# Recurrent Minimum Translation Unit Models



# Recurrent Minimum Translation Unit Models



# This talk

本 地 区 的 发 展 和 进 步 。



development and progress of the region .



Translation modeling

Auli et al., EMNLP 2013; Hu et al., EACL 2014



Language modeling



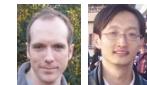
Optimization

Auli & Gao, ACL 2014

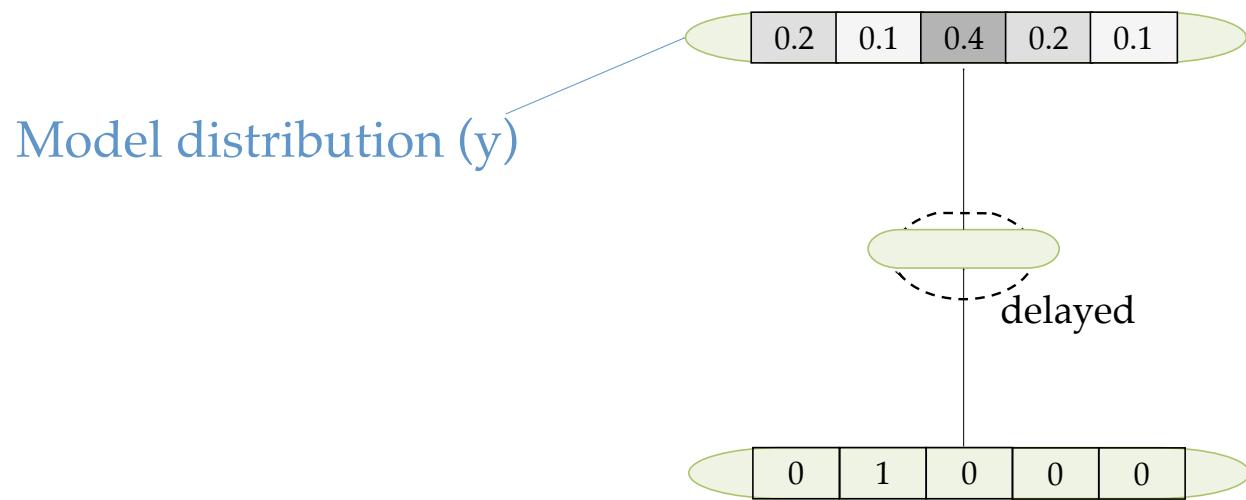


Reordering

Auli et al., EMNLP 2014



# Back propagation with cross entropy error

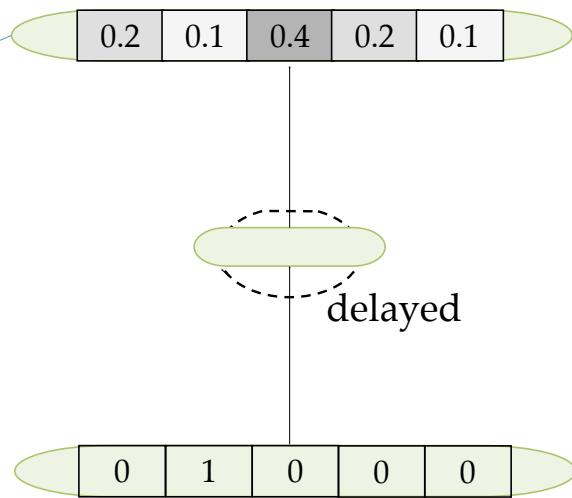


# Back propagation with cross entropy error

True distribution ( $t$ ) —————

0	0	1	0	0
---	---	---	---	---

Model distribution ( $y$ )

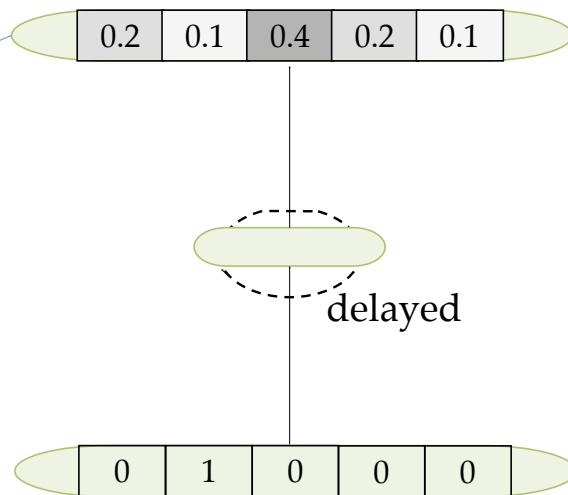


# Back propagation with cross entropy error

True distribution ( $t$ ) —————

0	0	1	0	0
---	---	---	---	---

Model distribution ( $y$ )



$$J = - \sum_i t^i \log y^i$$

Goal: Make correct outputs most likely

# Back propagation with cross entropy error

True distribution ( $t$ )

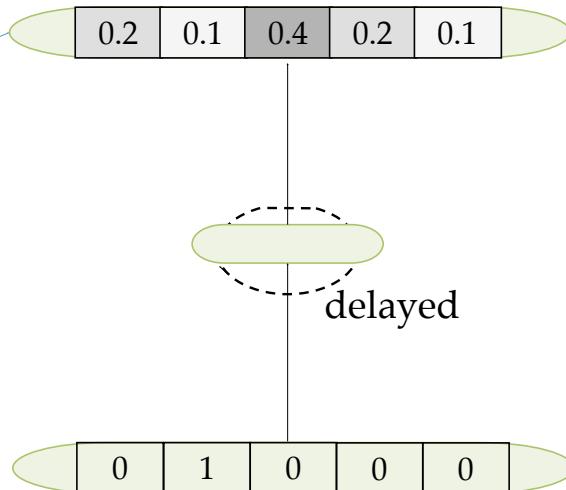
0	0	1	0	0
---	---	---	---	---

Error vector ( $\delta$ )

0.2	0.1	-0.6	0.2	0.1
-----	-----	------	-----	-----

$$\delta^i = y^i - t^i$$

Model distribution ( $y$ )



$$J = - \sum_i t^i \log y^i$$

Goal: Make correct outputs most likely

# Back propagation with cross entropy error

True distribution ( $t$ )

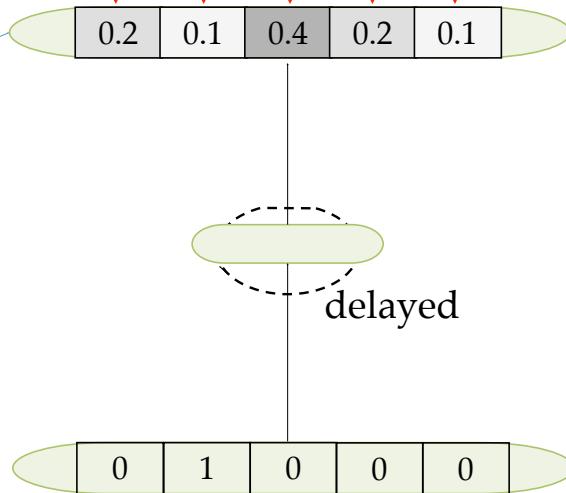
0	0	1	0	0
---	---	---	---	---

Error vector ( $\delta$ )

0.2	0.1	-0.6	0.2	0.1
-----	-----	------	-----	-----

$$\delta^i = y^i - t^i$$

Model distribution ( $y$ )



$$J = - \sum_i t^i \log y^i$$

Goal: Make correct outputs most likely

# Back propagation with cross entropy error

True distribution ( $t$ )

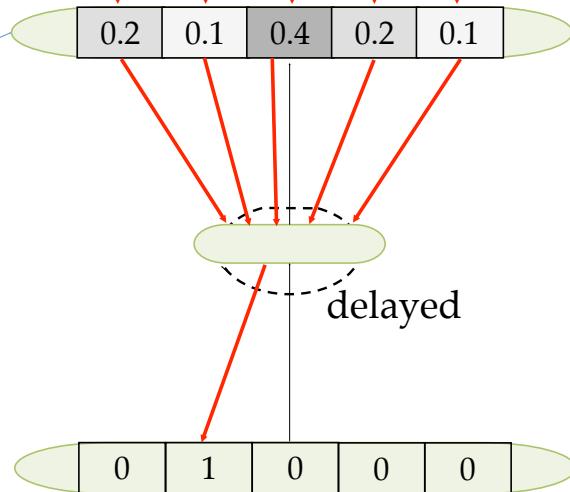
0	0	1	0	0
---	---	---	---	---

Error vector ( $\delta$ )

0.2	0.1	-0.6	0.2	0.1
-----	-----	------	-----	-----

$$\delta^i = y^i - t^i$$

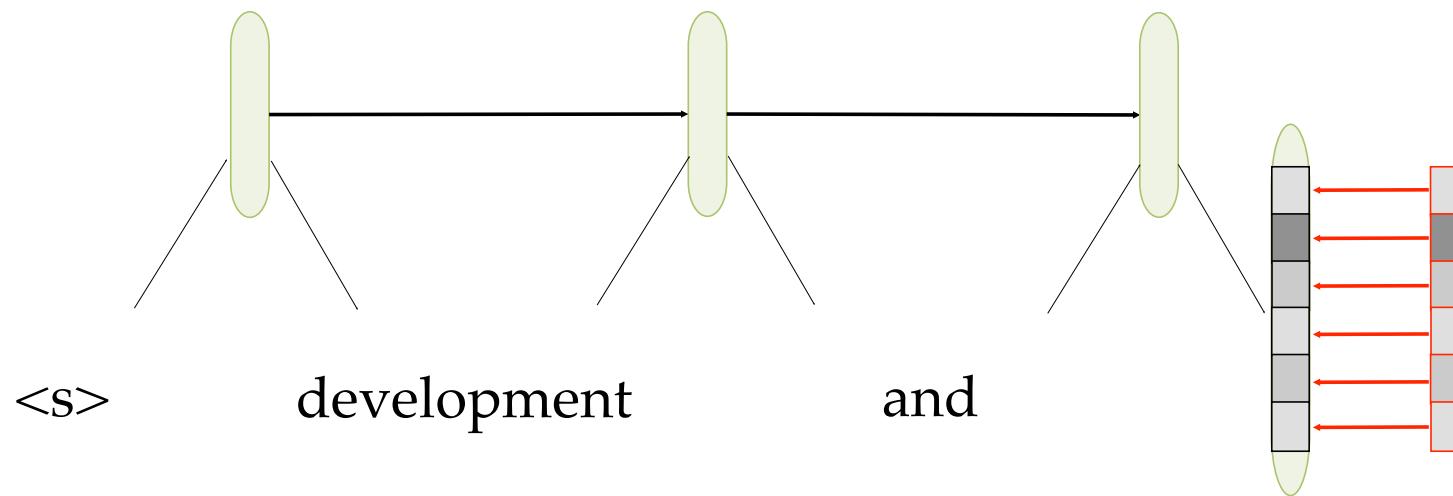
Model distribution ( $y$ )



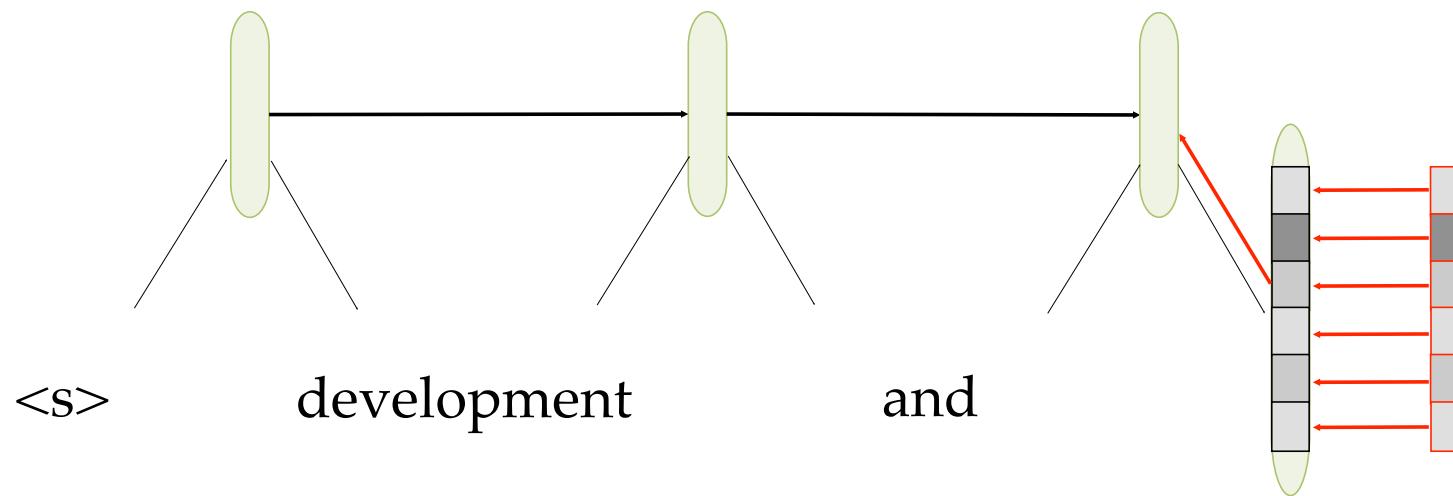
$$J = - \sum_i t^i \log y^i$$

Goal: Make correct outputs most likely

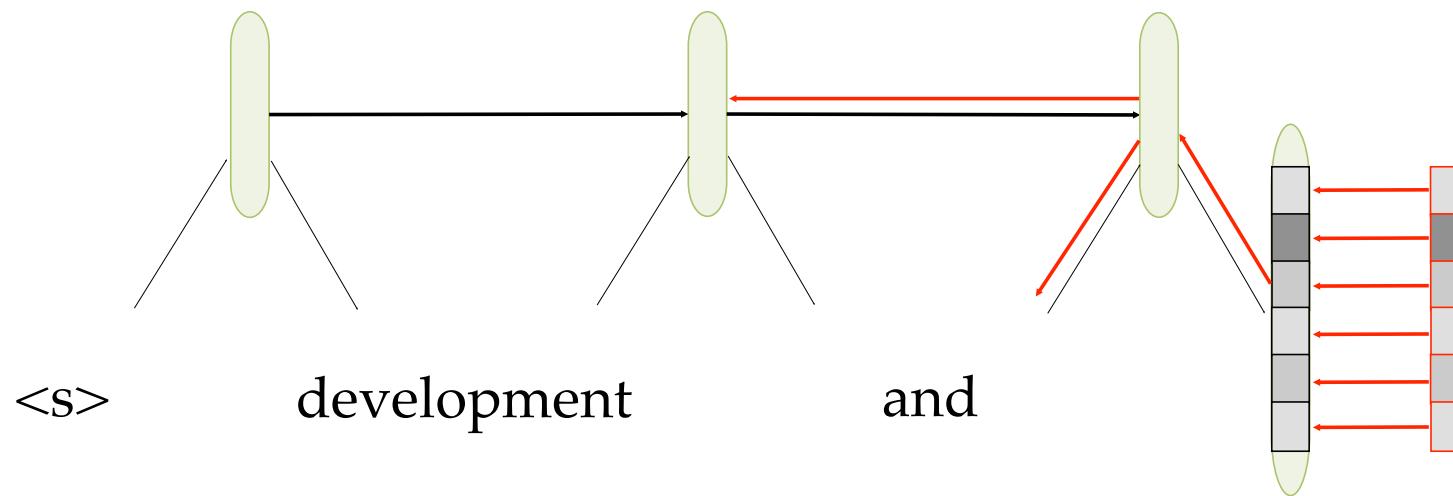
# Back propagation through time



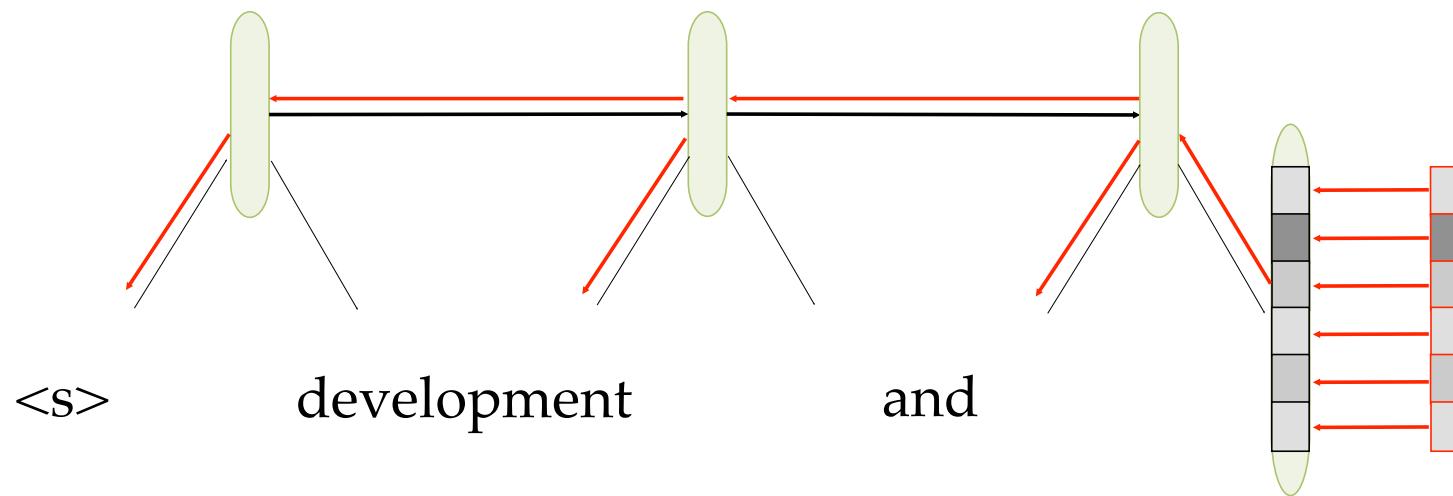
# Back propagation through time



# Back propagation through time



# Back propagation through time



# Optimization

Auli & Gao, ACL 2014



- Likelihood training very common
- Optimizing for evaluation metrics difficult, but empirically successful (Och 2003, Smith 2006, Chiang 2009, Gimpel 2010, Hopkins 2011)

# Optimization

Auli & Gao, ACL 2014



- Likelihood training very common
- Optimizing for evaluation metrics difficult, but empirically successful (Och 2003, Smith 2006, Chiang 2009, Gimpel 2010, Hopkins 2011)
- **Next:** Task-specific training of neural nets for translation

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

precision scores

brevity penalty

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

precision scores

brevity penalty

Human: development and progress of the region

System: advance and progress of region

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

precision scores

brevity penalty

Human: development and progress of the region

System: advance and progress of region

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

precision scores

brevity penalty

Human: development and progress of the region

System: advance and progress of region

# BLEU Metric

(Bilingual Evaluation Understudy; Papineni 2002)

$$\text{BLEU} = \exp \left( \sum_{n=1}^4 \frac{1}{4} \log p_n \right) \text{BP}$$

precision scores

brevity penalty

Human: development and progress of the region

System: advance and progress of region

# Expected BLEU Training

(Smith 2006, He 2012, Gao 2014)

L:

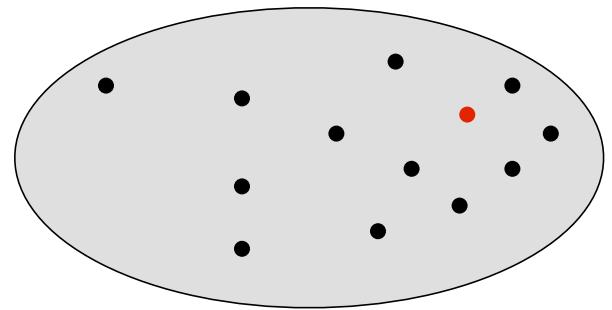
$$\max_{\theta} p(\tilde{e}|f; \theta)$$

# Expected BLEU Training

(Smith 2006, He 2012, Gao 2014)

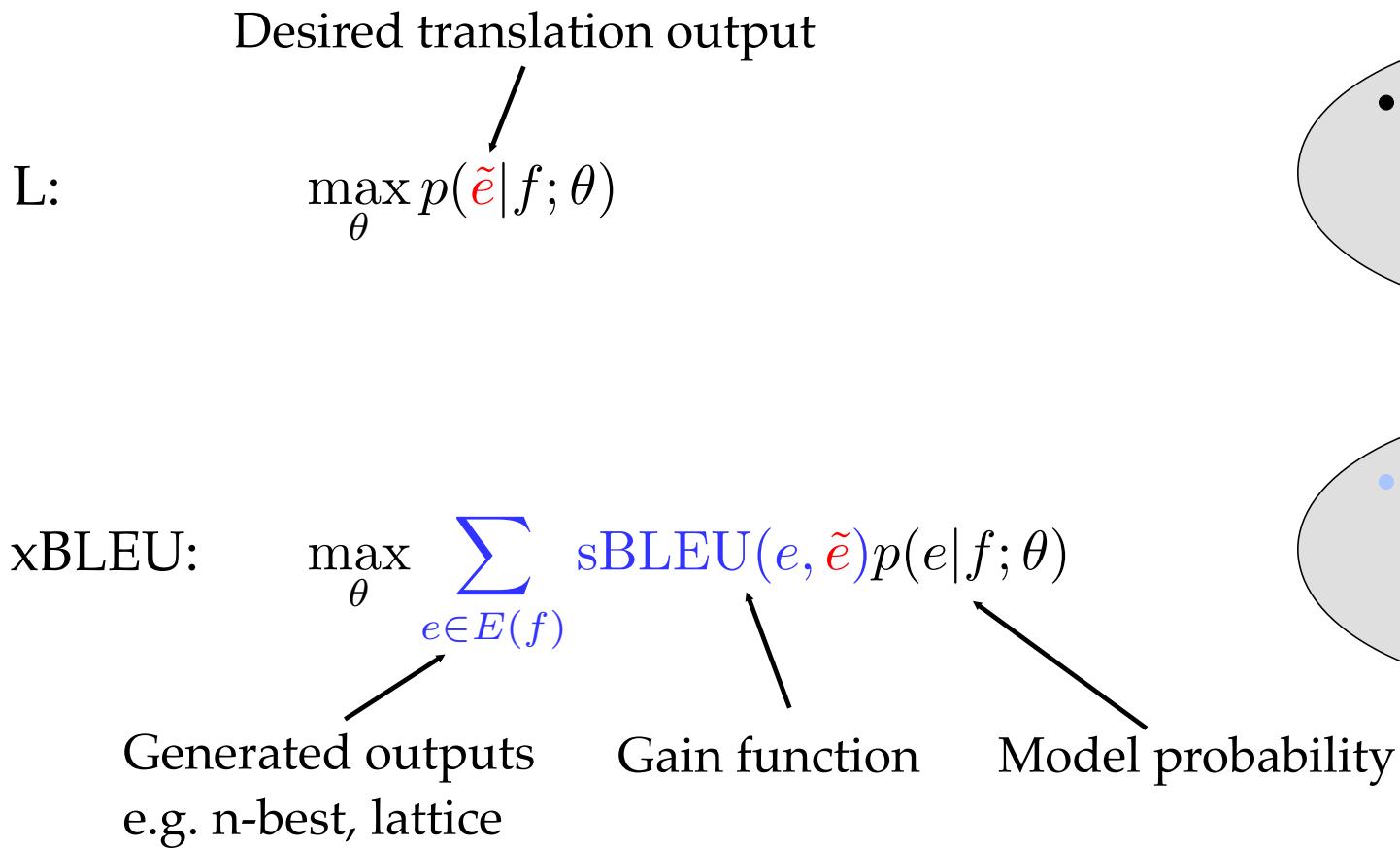
Desired translation output

$$\text{L: } \max_{\theta} p(\tilde{e}|f; \theta)$$

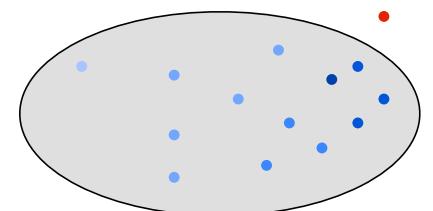


# Expected BLEU Training

(Smith 2006, He 2012, Gao 2014)



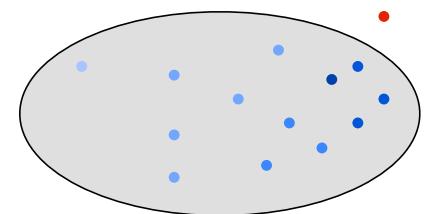
# Expected BLEU Training



本 地 区 的 发 展 和 进 步

Human: development and progress of the region

# Expected BLEU Training



本 地 区 的 发 展 和 进 步

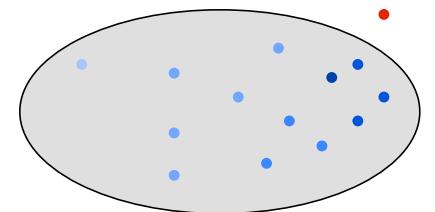
Human: development and progress of the region

advance and progress of the region

development and progress of this province

progress of this region

# Expected BLEU Training

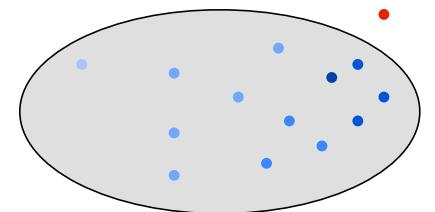


本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU
advance and progress of the region	0.8
development and progress of this province	0.5
progress of this region	0.3

# Expected BLEU Training

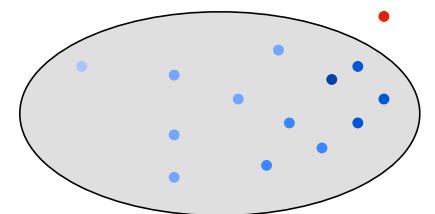


本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$
advance and progress of the region	0.8	0.2
development and progress of this province	0.5	0.3
progress of this region	0.3	0.5

# Expected BLEU Training

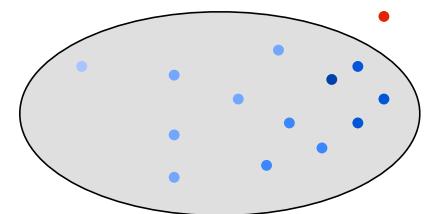


本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$
advance and progress of the region	0.8	0.2
development and progress of this province	0.5	0.3
progress of this region	0.3	0.5

# Expected BLEU Training

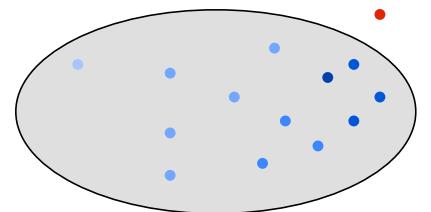


本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$
advance and progress of the region	0.8	0.2
development and progress of this province	0.5	0.3
progress of this region	0.3	0.5

# Expected BLEU Training



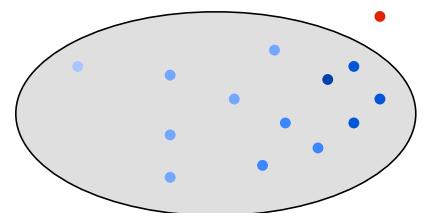
本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$
advance and progress of the region	0.8	0.2
development and progress of this province	0.5	0.3
progress of this region	0.3	0.5

$$\text{xBLEU} = \sum_{e \in E(f)} \text{sBLEU}(e, \tilde{e}) p(e|f; \theta) = 0.5$$

# Expected BLEU Training



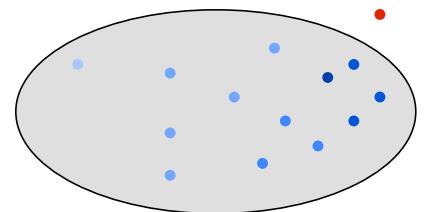
本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$	$\delta_t$
advance and progress of the region	0.8	0.2	0.3
development and progress of this province	0.5	0.3	0
progress of this region	0.3	0.5	-0.2

$$\text{xBLEU} = \sum_{e \in E(f)} \text{sBLEU}(e, \tilde{e}) p(e|f; \theta) = 0.5$$

# Expected BLEU Training



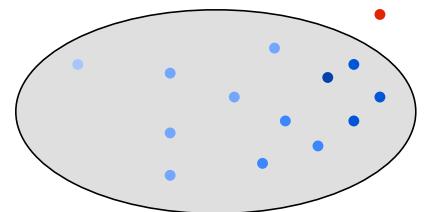
本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$	$\delta_t$	$p_{t+1}(e f; \theta)$
advance and progress of the region	0.8	0.2	0.3	0.5
development and progress of this province	0.5	0.3	0	0.3
progress of this region	0.3	0.5	-0.2	0.2

$$\text{xBLEU} = \sum_{e \in E(f)} \text{sBLEU}(e, \tilde{e}) p(e|f; \theta) = 0.5$$

# Expected BLEU Training



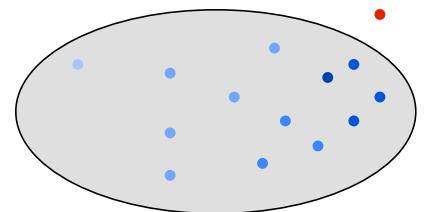
本 地 区 的 发 展 和 进 步

Human: development and progress of the region

	sBLEU	$p_t(e f; \theta)$	$\delta_t$	$p_{t+1}(e f; \theta)$
advance and progress of the region	0.8	0.2	0.3	0.5
development and progress of this province	0.5	0.3	0	0.3
progress of this region	0.3	0.5	-0.2	0.2

$$\text{xBLEU} = \sum_{e \in E(f)} \text{sBLEU}(e, \tilde{e}) p(e|f; \theta) = 0.5$$

# Expected BLEU Training



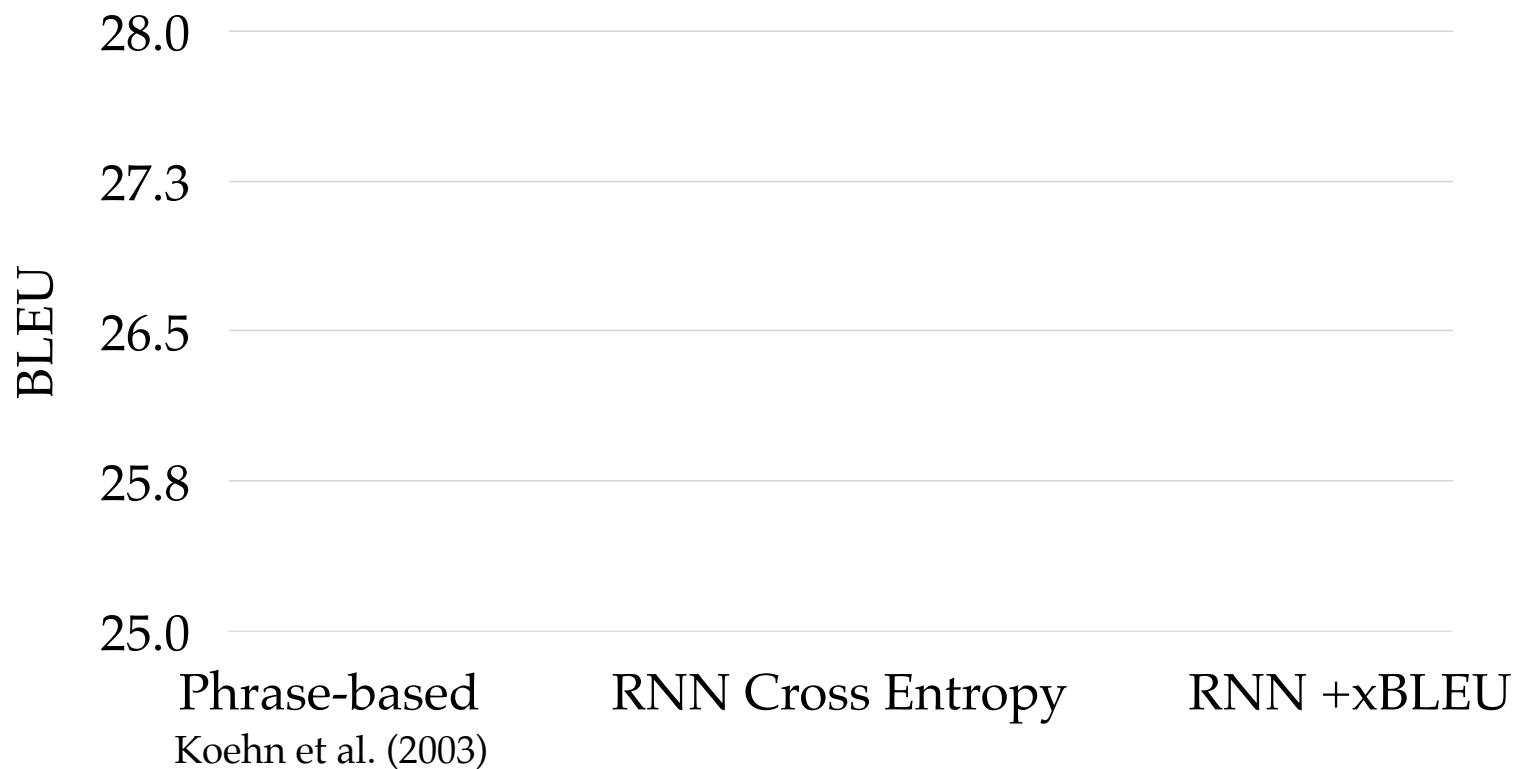
本 地 区 的 发 展 和 进 步

Human: development and progress of the region

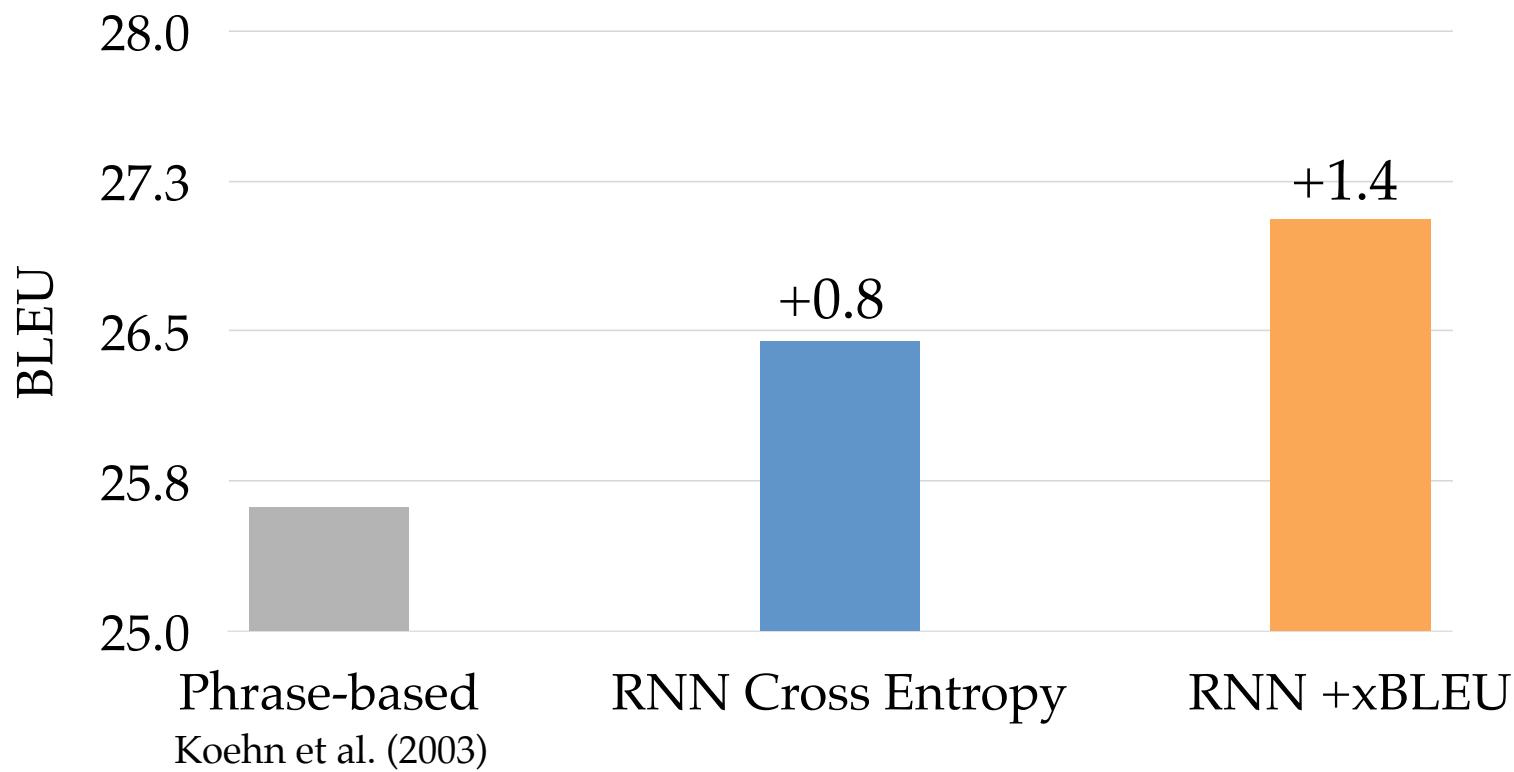
	sBLEU	$p_t(e f; \theta)$	$\delta_t$	$p_{t+1}(e f; \theta)$
advance and progress of the region	0.8	0.2	0.3	0.5
development and progress of this province	0.5	0.3	0	0.3
progress of this region	0.3	0.5	-0.2	0.2

$$\text{xBLEU} = \sum_{e \in E(f)} \text{sBLEU}(e, \tilde{e}) p(e|f; \theta) = 0.5 \rightarrow \textbf{0.6}$$

# Results



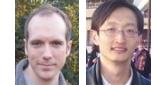
# Results



# Scaling linear reordering models

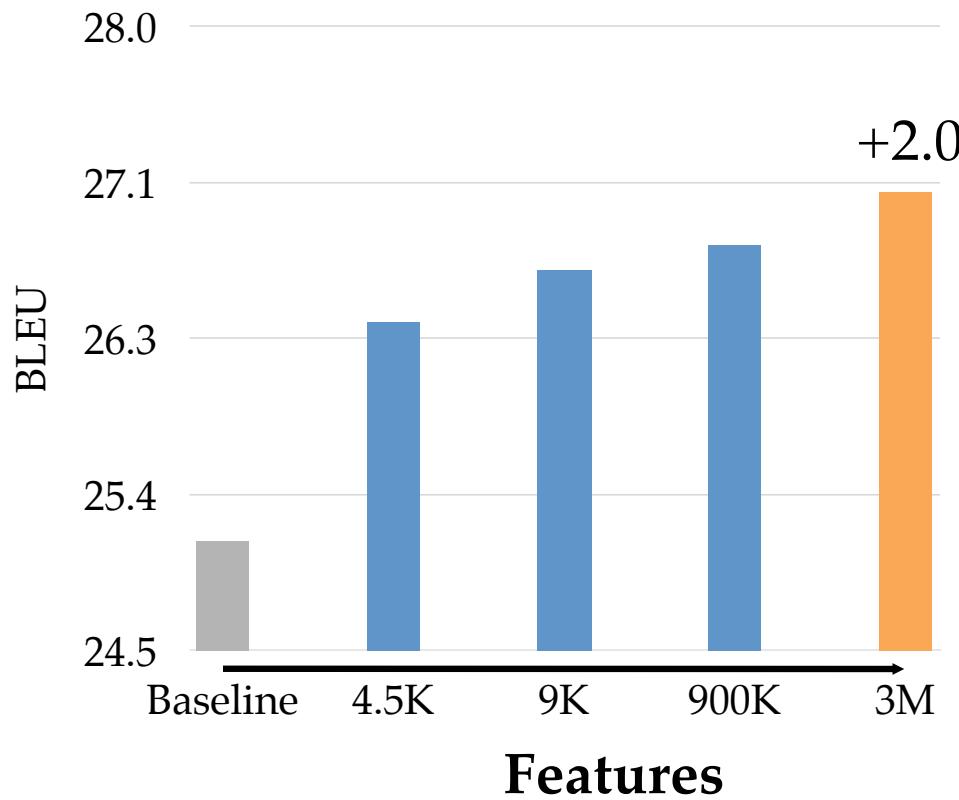
xBLEU training of millions of linear features

Auli et al., EMNLP 2014

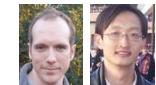


# Scaling linear reordering models

xBLEU training of millions of linear features



Auli et al., EMNLP 2014



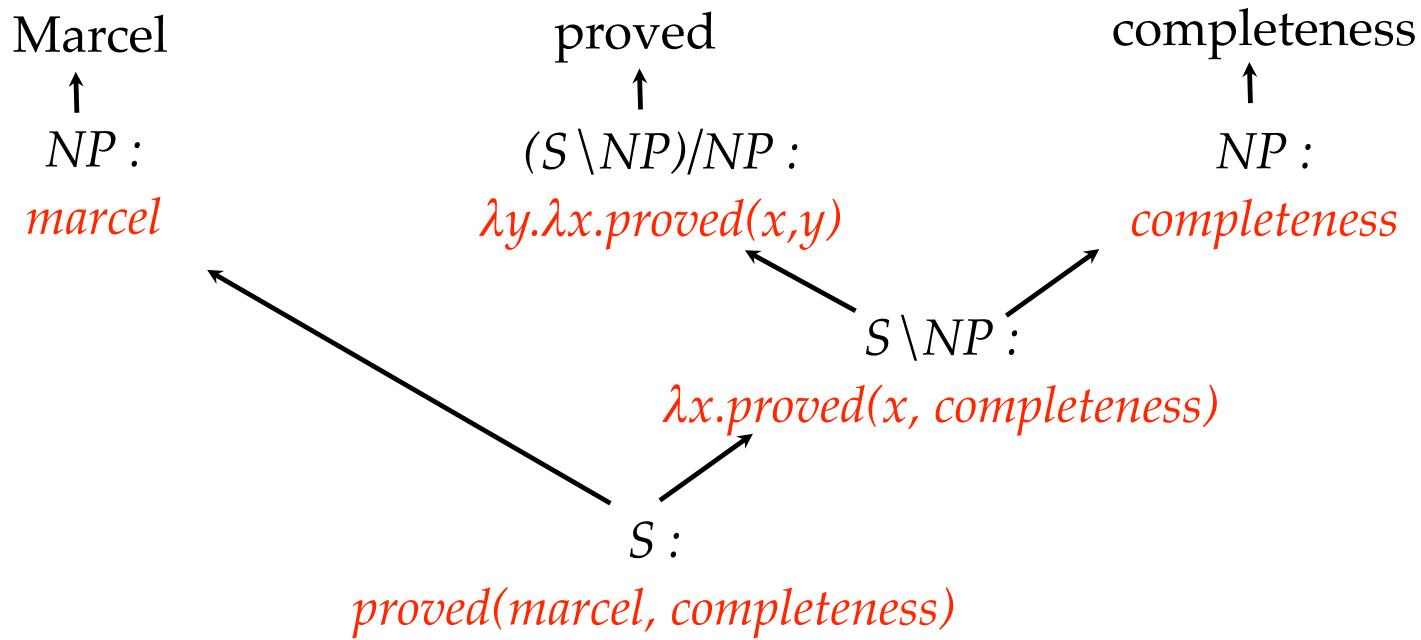
# My other neural network projects

- Social media response generation with RNNs  
building neural net-based conversational agents based on twitter conversations
- Semi-supervised phrase table expansion with word embeddings  
using distributional word and phrase representations and by mapping between distributional source and target spaces with RBVs
- CCG parsing & tagging with RNNs



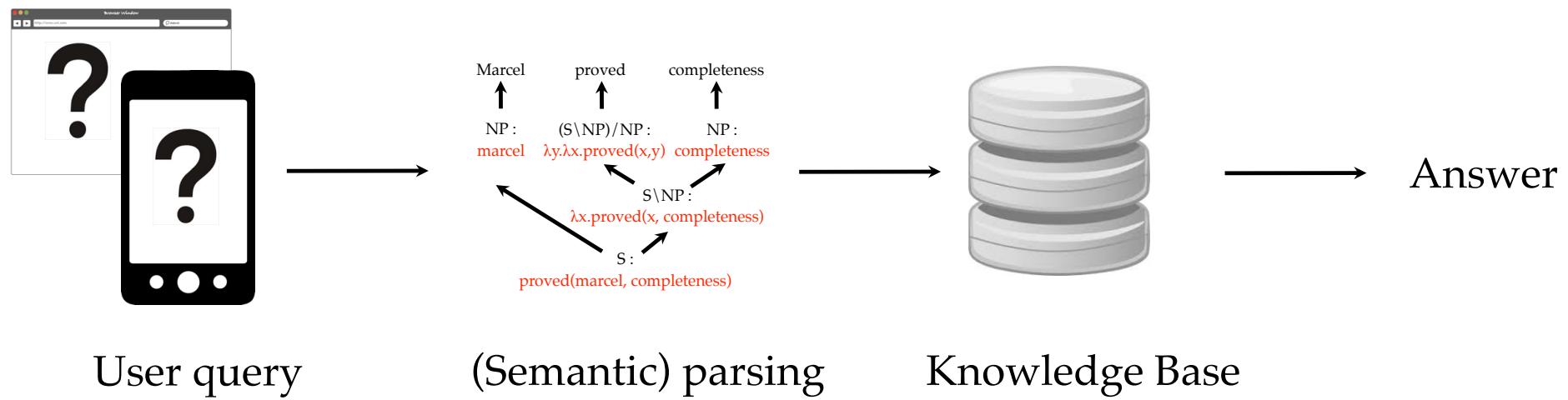
# Semantic CCG Parsing

Zettlemoyer (2005, 2007), Bos (2008), Kwiatkowski (2010, 2013) Krishnamurthy (2012), Lewis (2013a,b)



Combinatory Categorial Grammar (CCG; Steedman 2000)

# How is this useful?



# How is this useful?



User query

(Semantic) parsing

Knowledge Base

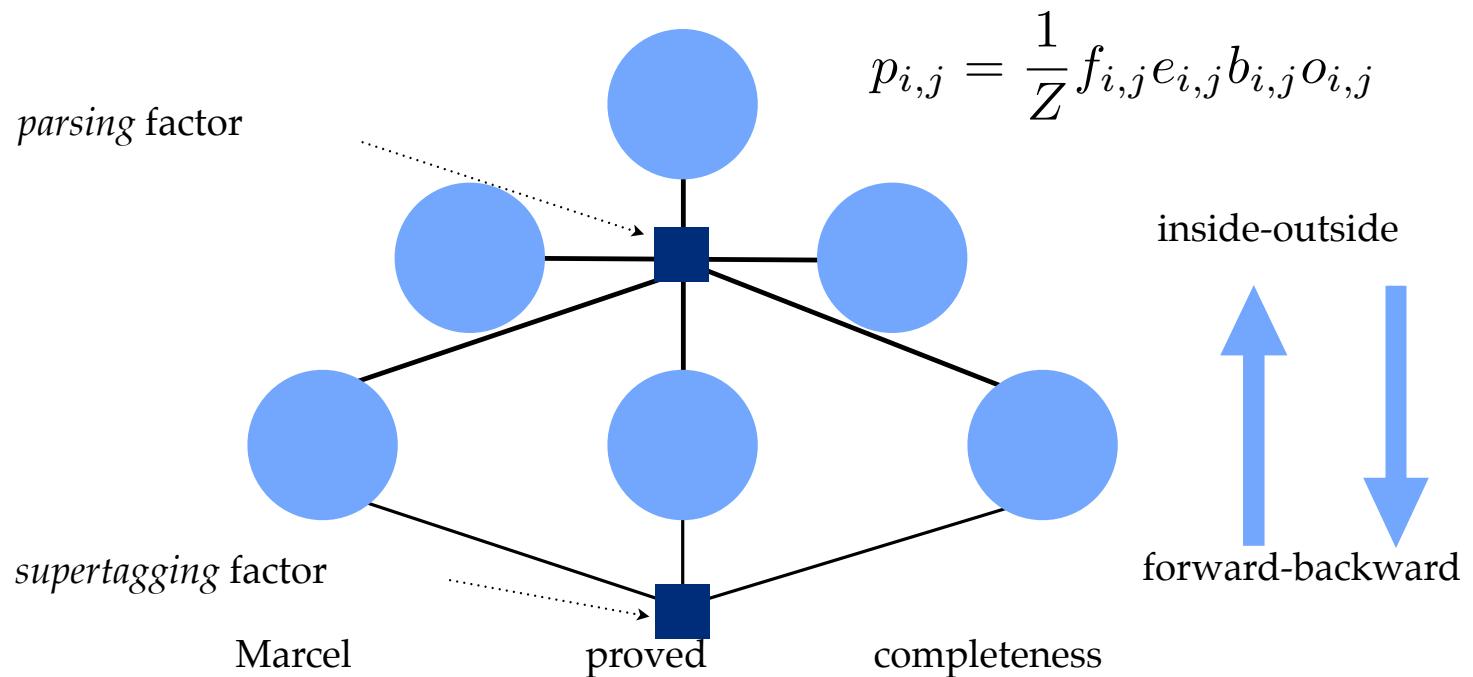
# Integrated Parsing & Tagging

with belief propagation, dual decomposition and softmax-margin training



Auli & Lopez ACL 2011a,b

Auli & Lopez EMNLP 2011



# Integrated Parsing & Tagging

- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).

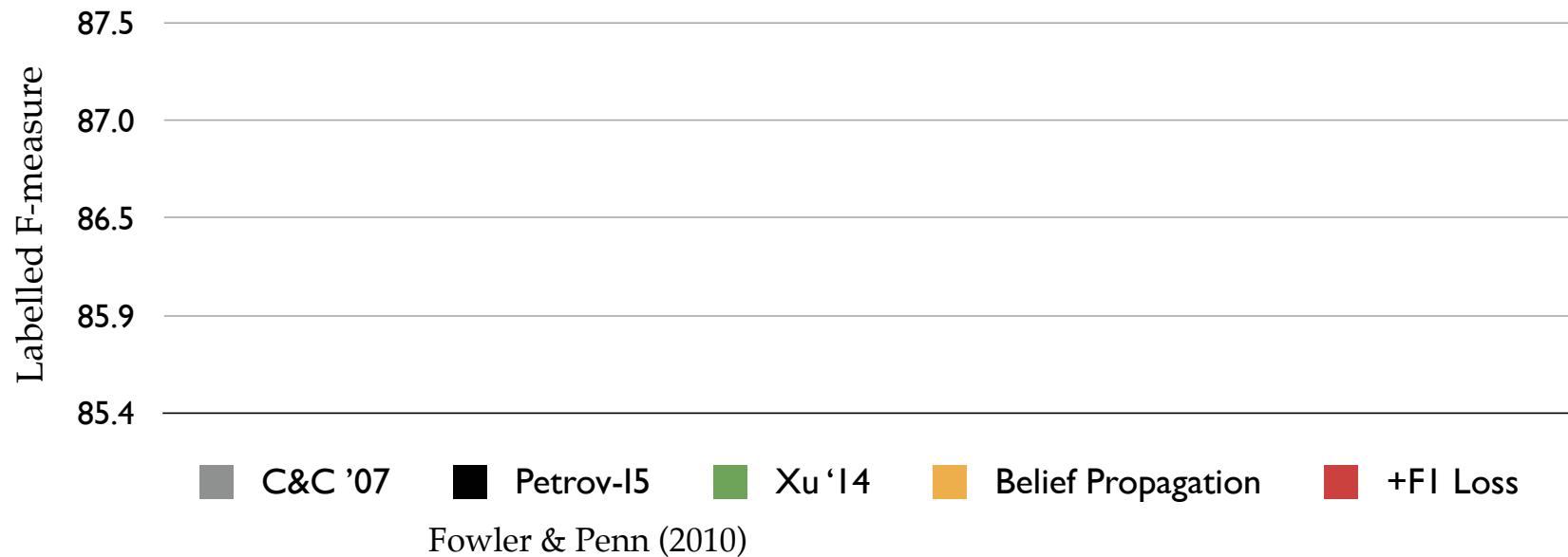


Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).

best CCG parsing  
results to date

- Belief propagation for inference.



# Integrated Parsing & Tagging

- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).

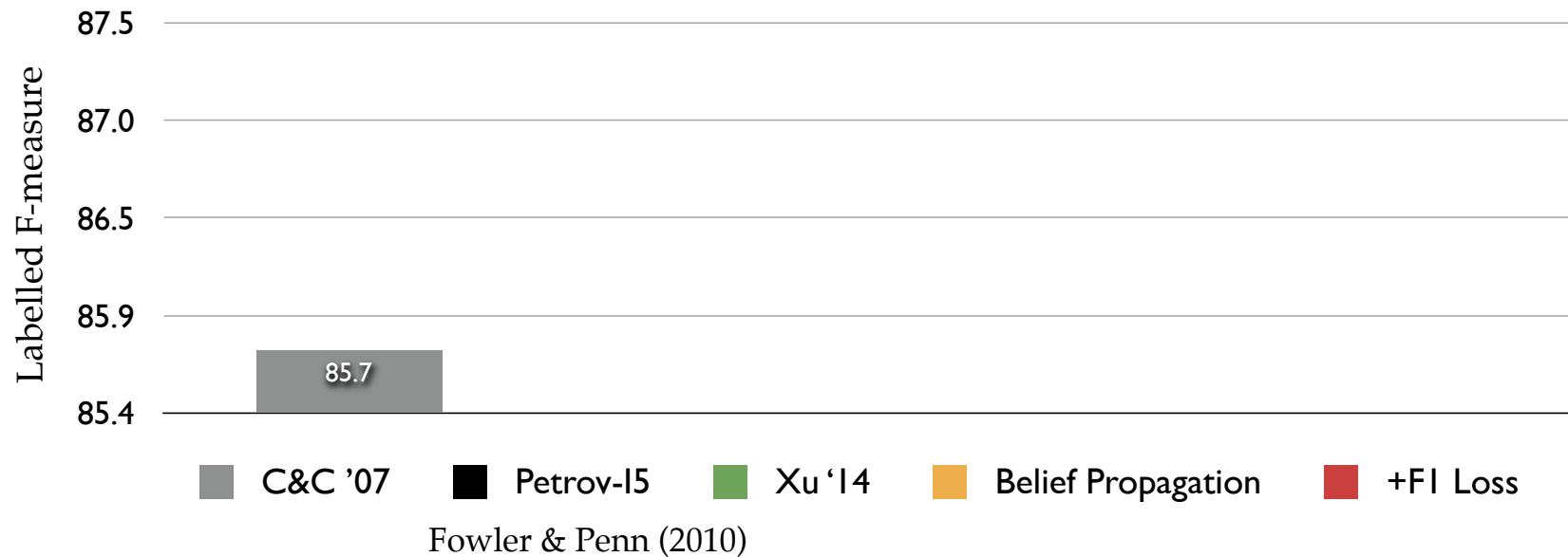


Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).

best CCG parsing  
results to date

- Belief propagation for inference.



# Integrated Parsing & Tagging

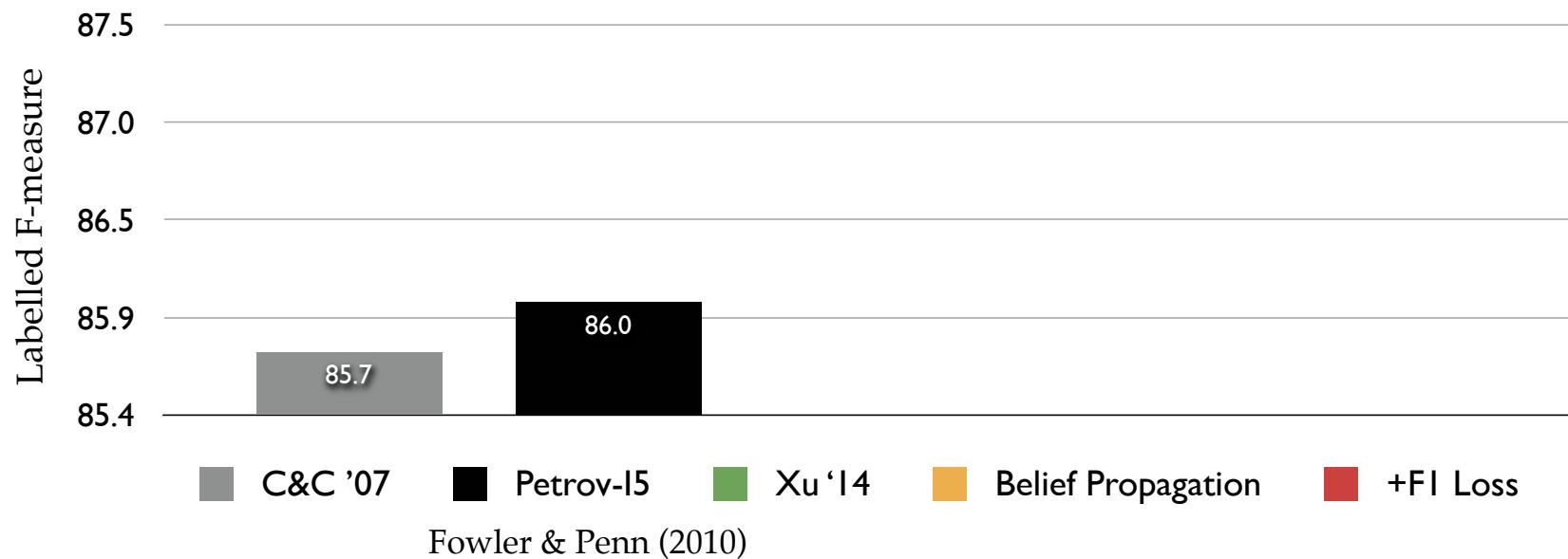
- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).



Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).
- Belief propagation for inference.

best CCG parsing  
results to date



# Integrated Parsing & Tagging

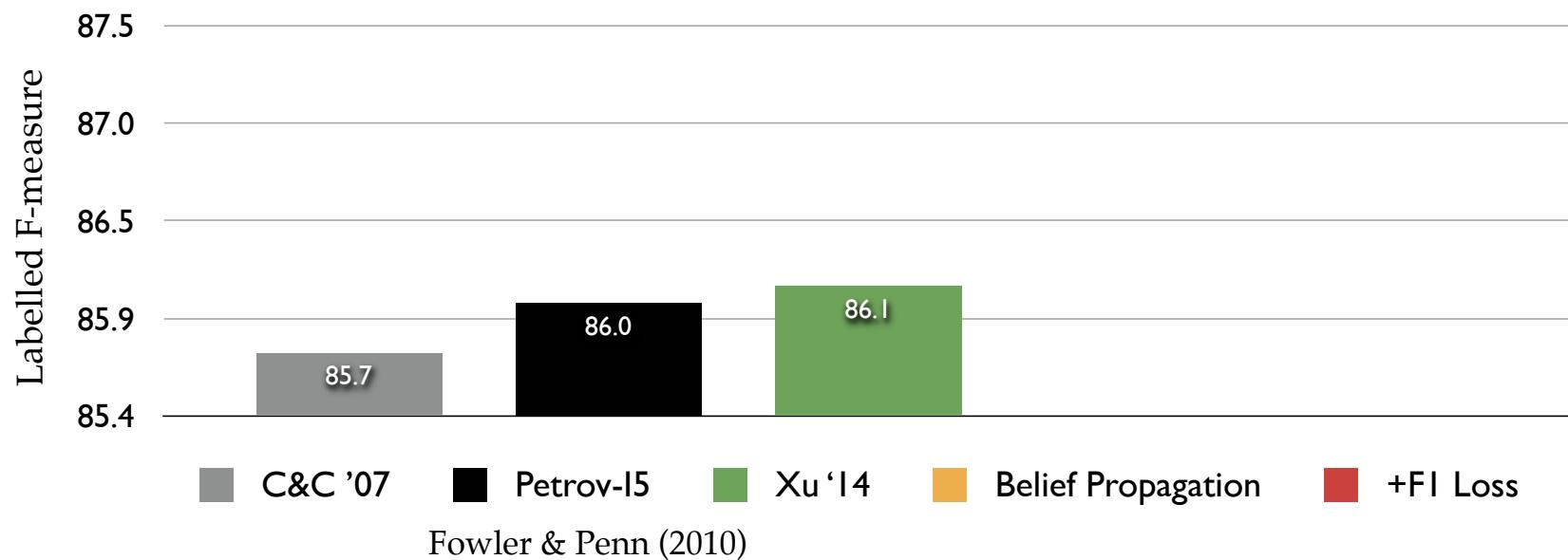
- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).



Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).
- Belief propagation for inference.

best CCG parsing  
results to date



# Integrated Parsing & Tagging

- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).

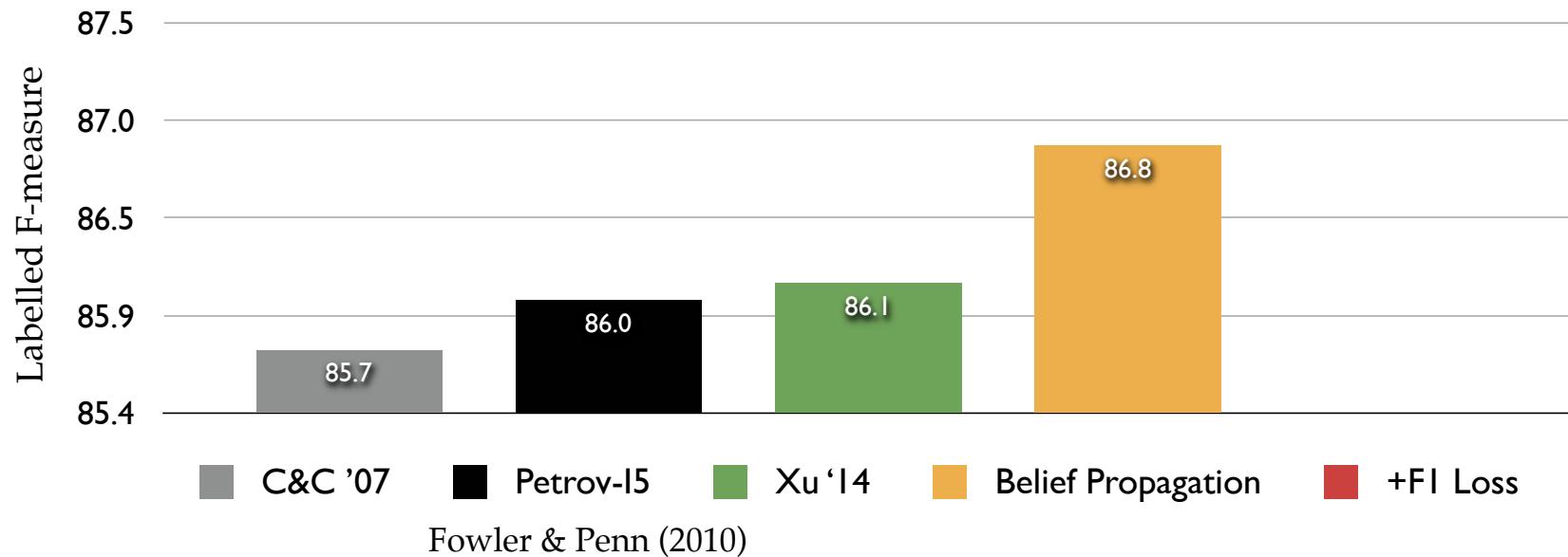


Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).

best CCG parsing  
results to date

- Belief propagation for inference.



# Integrated Parsing & Tagging

- F-measure loss for parsing sub-model (+DecF<sub>1</sub>).

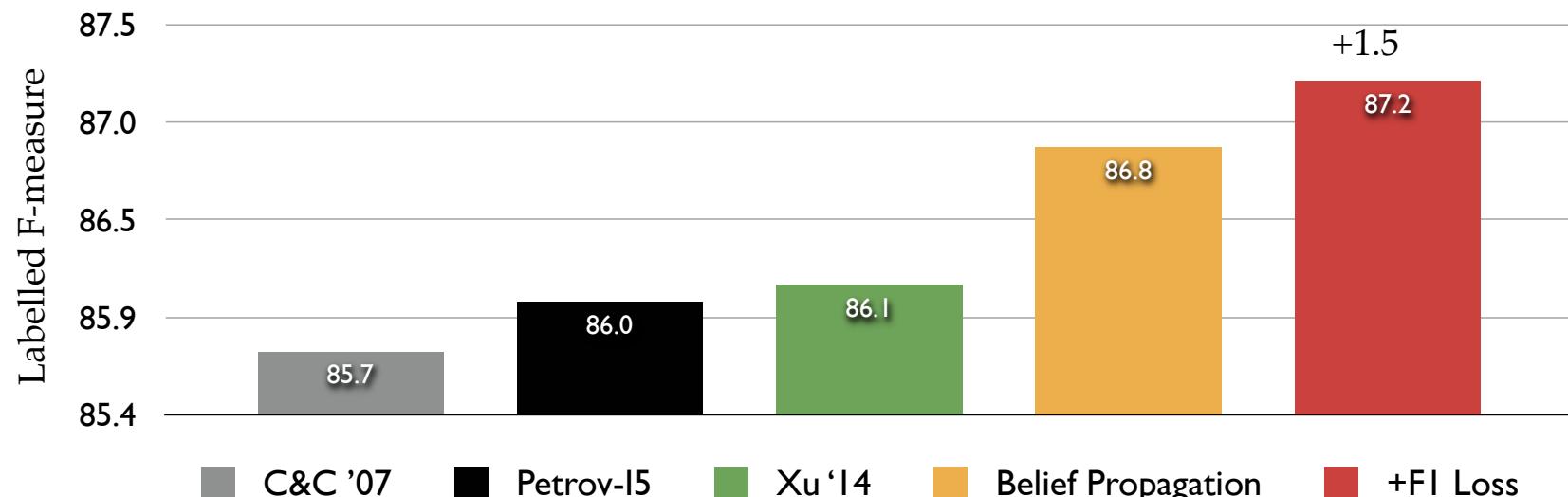


Auli & Lopez EMNLP 2011

- Hamming loss for supertagging sub-model (+Tagger).

best CCG parsing  
results to date

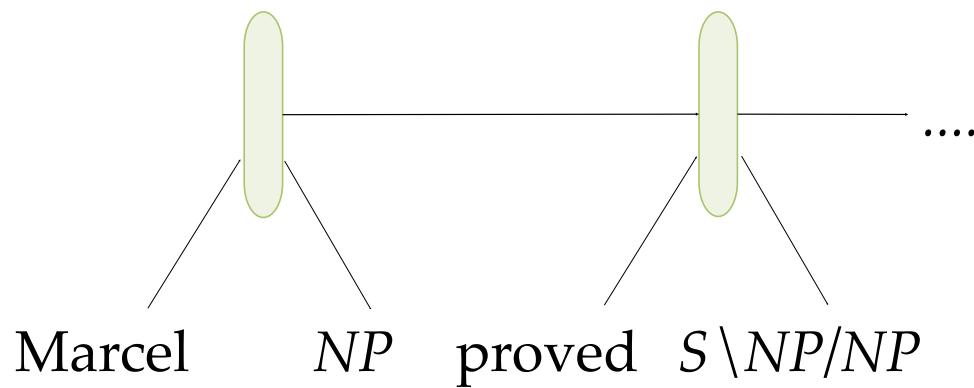
- Belief propagation for inference.



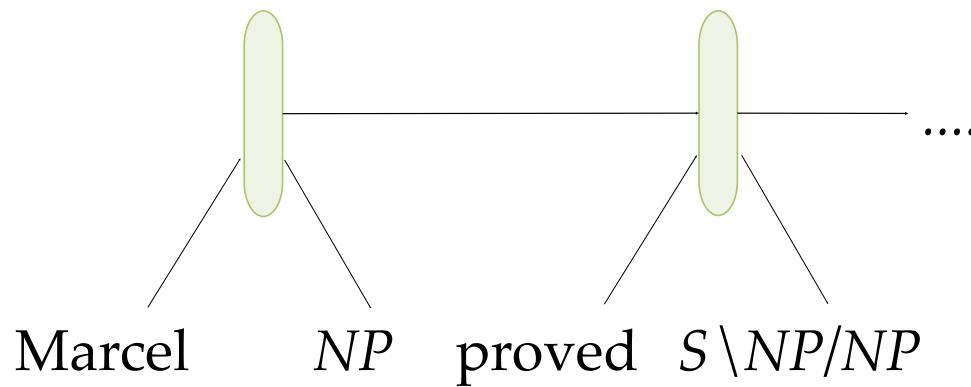
Fowler & Penn (2010)

# Recurrent nets for CCG supertagging & parsing

with Wenduan Xu



# Recurrent nets for CCG supertagging & parsing



with Wenduan Xu



	tagging	parsing
CRF (Clark & Curran '07)	91.5	85.3
FFN (Lewis & Steedman, '14)	91.5	86.0
RNN	92.3	86.5

# Summary

- Two RNN translation models
- Neural nets help most when discrete models sparse
- Task-specific objective gives best performance
- Next: Better modeling of source-side, e.g., bi-directional RNNs, different architectures