

Generative Models for Discriminative Problems

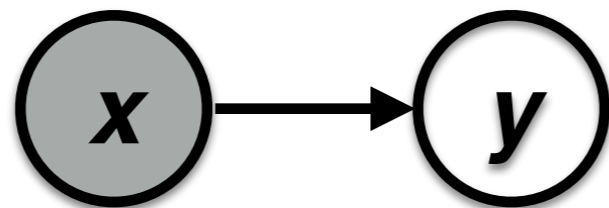
Chris Dyer
DeepMind

Terminological clarification

- A **discriminative problem**: for some input \mathbf{x} , find the most likely \mathbf{y} in a set $\mathcal{Y}(\mathbf{x})$

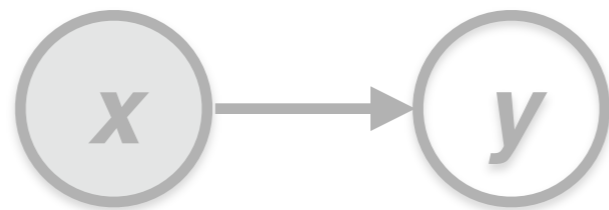
Terminological clarification

- A **discriminative problem**: for some input \mathbf{x} , find the most likely \mathbf{y} in a set $\mathcal{Y}(\mathbf{x})$
- A **discriminative model** directly models $p(\mathbf{y} | \mathbf{x})$
logistic/linear/... regressions, MLPs, CRFs, MEMMs, seq2seq(+att)

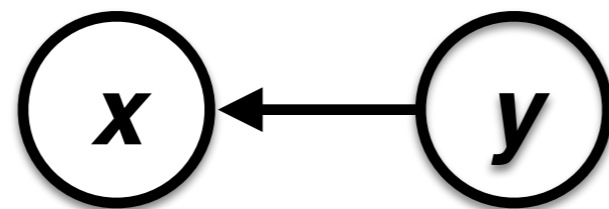


Terminological clarification

- A **discriminative problem**: for some input \mathbf{x} , find the most likely \mathbf{y} in a set $\mathcal{Y}(\mathbf{x})$
- A **discriminative model** directly models $p(\mathbf{y} | \mathbf{x})$
logistic/linear/... regressions, MLPs, CRFs, MEMMs, seq2seq(+att)

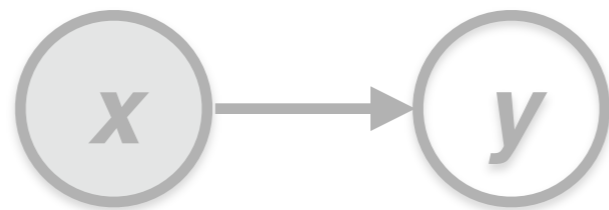


- A **generative model for a discriminative problem** models $p(\mathbf{x}, \mathbf{y})$, often by breaking it into $p(\mathbf{y})p(\mathbf{x} | \mathbf{y})$
Naive Bayes, GMMs, HMMs, PCFGs, the IBM translation models

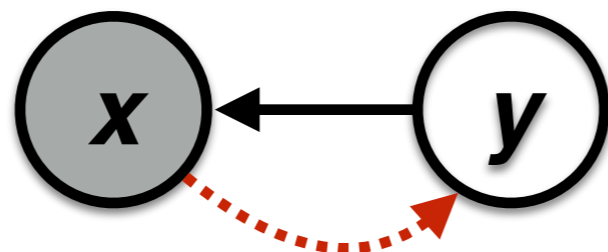


Terminological clarification

- A **discriminative problem**: for some input \mathbf{x} , find the most likely \mathbf{y} in a set $\mathcal{Y}(\mathbf{x})$
- A **discriminative model** directly models $p(\mathbf{y} | \mathbf{x})$
logistic/linear/... regressions, MLPs, CRFs, MEMMs, seq2seq(+att)

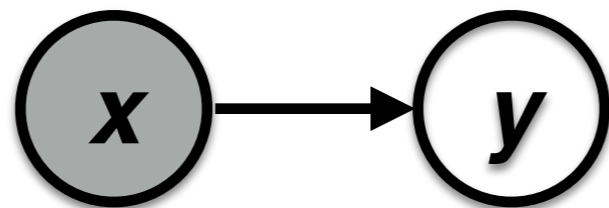


- A **generative model for a discriminative problem** models $p(\mathbf{x}, \mathbf{y})$, often by breaking it into $p(\mathbf{y})p(\mathbf{x} | \mathbf{y})$
Naive Bayes, GMMs, HMMs, PCFGs, the IBM translation models

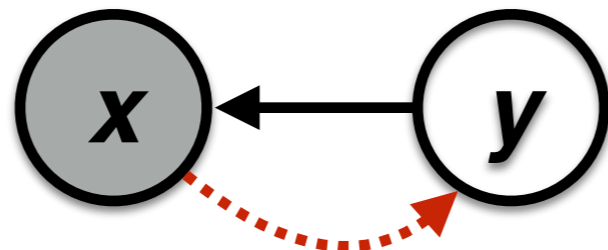


Terminological clarification

- A **discriminative problem**: for some input \mathbf{x} , find the most likely \mathbf{y} in a set $\mathcal{Y}(\mathbf{x})$
- A **discriminative model** directly models $p(\mathbf{y} | \mathbf{x})$
logistic/linear/... regressions, MLPs, CRFs, MEMMs, seq2seq(+att)



- A **generative model for a discriminative problem** models $p(\mathbf{x}, \mathbf{y})$, often by breaking it into $p(\mathbf{y})p(\mathbf{x} | \mathbf{y})$
Naive Bayes, GMMs, HMMs, PCFGs, the IBM translation models



system	BLEU	HTER	mTER
PBSY	25.3	28.0	21.8
HPB	24.6	29.9	23.4
SPB	25.8	29.0	22.7
NMT	31.1*	21.1*	16.2*

Table 2: Overall results on the HE Set: BLEU, computed against the original reference translation, and TER, computed with respect to the targeted post-edit (HTER) and multiple post-edits (mTER).

(Bentivogli et al., 2016)

But why?

Exp-ID	Model	Unidi	1st pass Model Size
E8	Proposed	5.6	0.4 GB
E9	Conventional LFR system	6.7	0.1 GB (AM) + 2.2 GB (PM) + 4.9 GB (LM) = 7.2GB

Table 5: The improved LAS outperforms the conventional LFR system while being more compact. Both models use second-pass rescoring.

(Chiu et al., last week)

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** ([Lake et al., 2015](#))

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** ([Lake et al., 2015](#))
- Generative models may be more **sample efficient** than equivalent discriminative models ([Ng & Jordan, 2001](#))
 - In some domains, we can build (relatively) accurate models of data generation → even better sample efficiency

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** (Lake et al., 2015)
- Generative models may be more **sample efficient** than equivalent discriminative models (Ng & Jordan, 2001)
 - In some domains, we can build (relatively) accurate models of data generation → even better sample efficiency
- **Exploit alternative data/variables**: zero shot learning, learning from unpaired samples, semisupervised learning, *exploit natural conditional independencies*

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** (Lake et al., 2015)
- Generative models may be more **sample efficient** than equivalent discriminative models (Ng & Jordan, 2001)
 - In some domains, we can build (relatively) accurate models of data generation → even better sample efficiency
- **Exploit alternative data/variables**: zero shot learning, learning from unpaired samples, semisupervised learning, *exploit natural conditional independencies*
- **Reduce label bias** when producing sequential outputs

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** (Lake et al., 2015)
- Generative models may be more **sample efficient** than equivalent discriminative models (Ng & Jordan, 2001)
 - In some domains, we can build (relatively) accurate models of data generation → even better sample efficiency
- **Exploit alternative data/variables**: zero shot learning, learning from unpaired samples, semisupervised learning, *exploit natural conditional independencies*
- **Reduce label bias** when producing sequential outputs
- **Safety considerations**: model introspection by sampling, generative models “know what they know”

Why generative models?

Five reasons

- “Human-like learning” looks more like **model building+inference** than **optimizing pattern recognition functions** ([Lake et al., 2015](#))
- Generative models may be more **sample efficient** than equivalent discriminative models ([Ng & Jordan, 2001](#))
 - In some domains, we can build (relatively) accurate models of data generation → even better sample efficiency
- **Exploit alternative data/variables**: zero shot learning, learning from unpaired samples, semisupervised learning, *exploit natural conditional independencies*
- **Reduce label bias** when producing sequential outputs
- **Safety considerations**: model introspection by sampling, generative models “know what they know”

**But didn't we use generative models
and give them up for some reason?**

Why not generative models?

- To use “generative models for discriminative problems” we must **model complex distributions** (sentences, documents, speech, images)
 - Complex distributions → lots of bad independence assumptions
(*naive Bayes, n-grams, HMMs, statistical translation models*)

Why not generative models?

- To use “generative models for discriminative problems” we must **model complex distributions** (sentences, documents, speech, images)
 - Complex distributions → lots of bad independence assumptions (*naive Bayes, n-grams, HMMs, statistical translation models*)
 - **But:** neural networks let the learner figure out their own independence assumptions!

Why not generative models?

- To use “generative models for discriminative problems” we must **model complex distributions** (sentences, documents, speech, images)
 - Complex distributions → lots of bad independence assumptions (*naive Bayes, n-grams, HMMs, statistical translation models*)
 - **But:** neural networks let the learner figure out their own independence assumptions!
- Using generative models require solving **difficult inference problems**
 - Inference problems are especially difficult when you get rid of the “bad independence assumptions”!

Why not generative models?

- To use “generative models for discriminative problems” we must **model complex distributions** (sentences, documents, speech, images)
 - Complex distributions → lots of bad independence assumptions (*naive Bayes, n-grams, HMMs, statistical translation models*)
 - **But:** neural networks let the learner figure out their own independence assumptions!
- Using generative models require solving **difficult inference problems**
 - Inference problems are especially difficult when you get rid of the “bad independence assumptions”!
- You aren’t “**optimizing the task**”!

Why not generative models?

- To use “generative models for discriminative problems” we must **model complex distributions** (sentences, documents, speech, images)
 - Complex distributions → lots of bad independence assumptions (*naive Bayes, n-grams, HMMs, statistical translation models*)
 - **But:** neural networks let the learner figure out their own independence assumptions!
- Using generative models require solving **difficult inference problems**
 - Inference problems are especially difficult when you get rid of the “bad independence assumptions”!
- You aren’t “**optimizing the task**”!

Case studies

- **Text categorization**

$x =$

US surrounds new London embassy with a moat

Heavily defended and dubbed a glass Berlin city, the new moated building also featured a

It is central London's first new moated building since the medieval era. The new US embassy in Nine Elms is a paradox, a heavily defended delicate glass box. Its architect calls it a "crystaline radiant beacon"; in fact, it resembles a corporate cube.

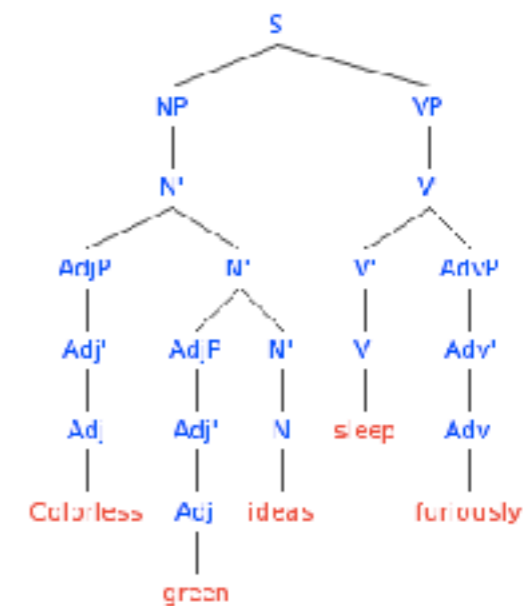
It is also one of the world's most expensive embassies, costing a cool \$400m. Remarkably, not a cent of US taxpayer money has been spent. Speaking at the press launch on Wednesday, Ambassador William Moran, principal deputy director of the Bureau of US Overseas Buildings Operations, confirmed that the new building "was entirely funded from the proceeds of real estate sales".

$y =$ POLITICS

- **Syntactic parsing**

$x =$ Colorless green ideas
sleep furiously

$y =$



- **Sequence to sequence transduction**

$x =$ Welcome to Okinawa

$y =$ 沖縄へようこそ。

Case studies

- **Text categorization**

$x =$

US surrounds new London embassy with a moat

Heavily defended and dubbed a glass Berlin city, the new moated building also featured a

It is central London's first new moated building since the medieval era. The new US embassy in Nine Elms is a paradox, a heavily defended delicate glass box. Its architect calls it a "crystaline radiant beacon"; in fact, it resembles a corporate cube.

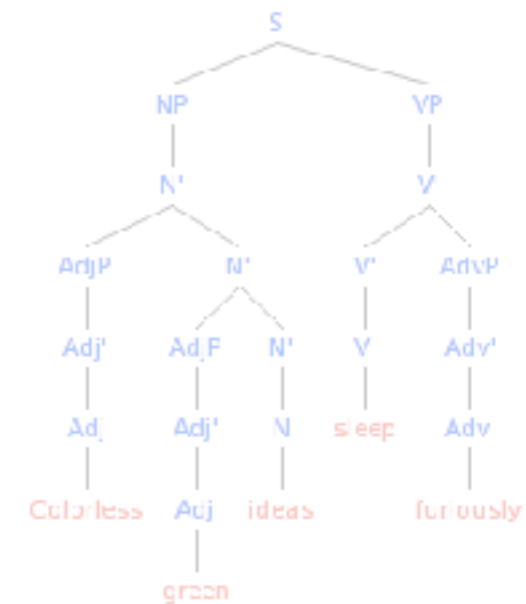
It is also one of the world's most expensive embassies, costing a cool \$400m. Remarkably, not a cent of US taxpayer money has been spent. Speaking at the press launch on Wednesday, Ambassador William Moran, principal deputy director of the Bureau of US Overseas Buildings Operations, confirmed that the new building "was entirely funded from the proceeds of real estate sales".

$y =$ POLITICS

- **Syntactic parsing**

$x =$ Colorless green ideas
sleep furiously

$y =$



- **Sequence to sequence transduction**

$x =$ Welcome to Okinawa

$y =$ 沖縄へようこそ。

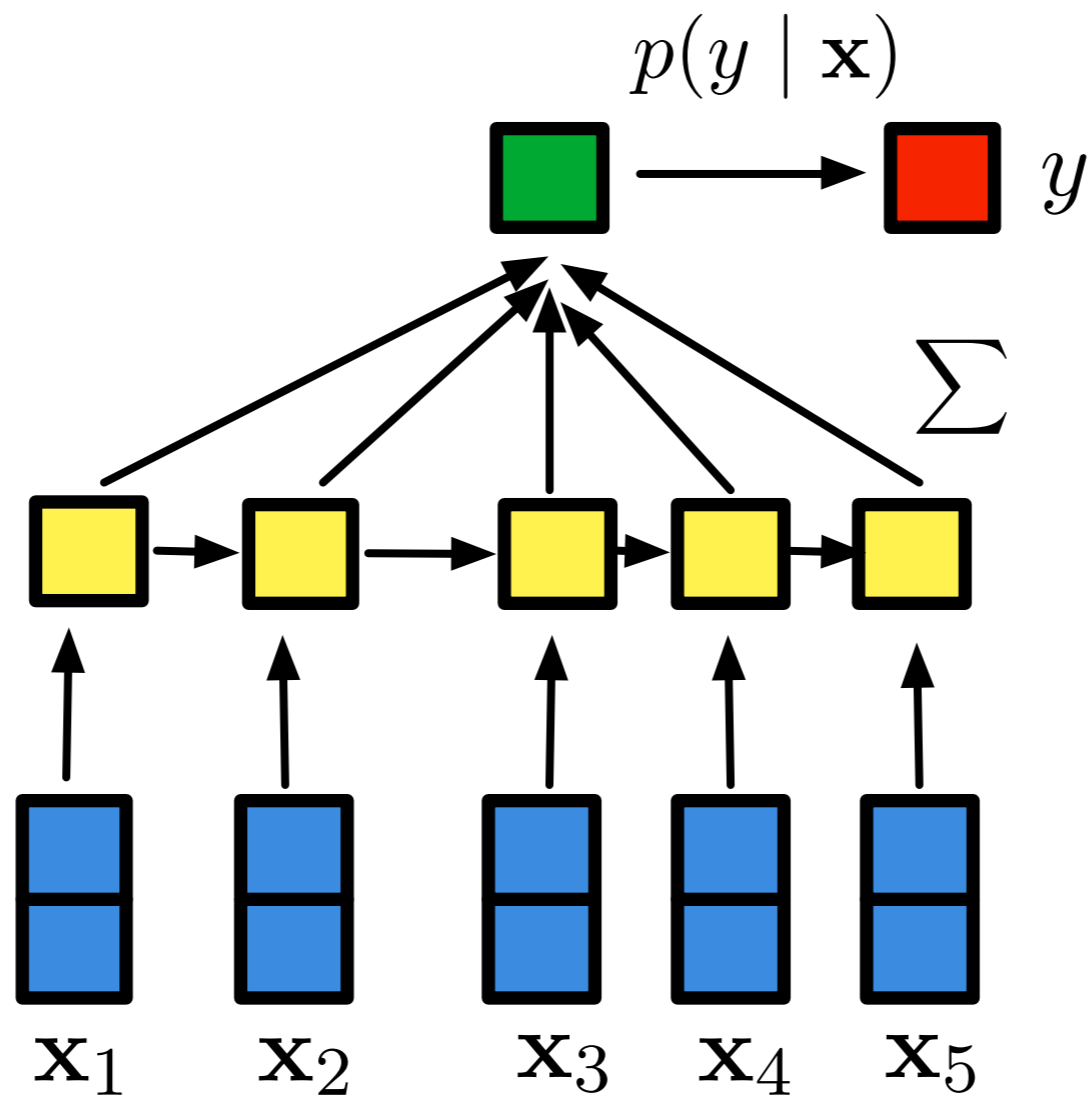
Experimental setup

Text categorization

- **Supervised classification**
 - Sample efficiency of a generative-discriminative pair ([Ng and Jordan, 2001](#))
 - How well do generative models do on standard datasets “at scale”?
 - How well do generative models do across a range of data conditions?

Discriminative model

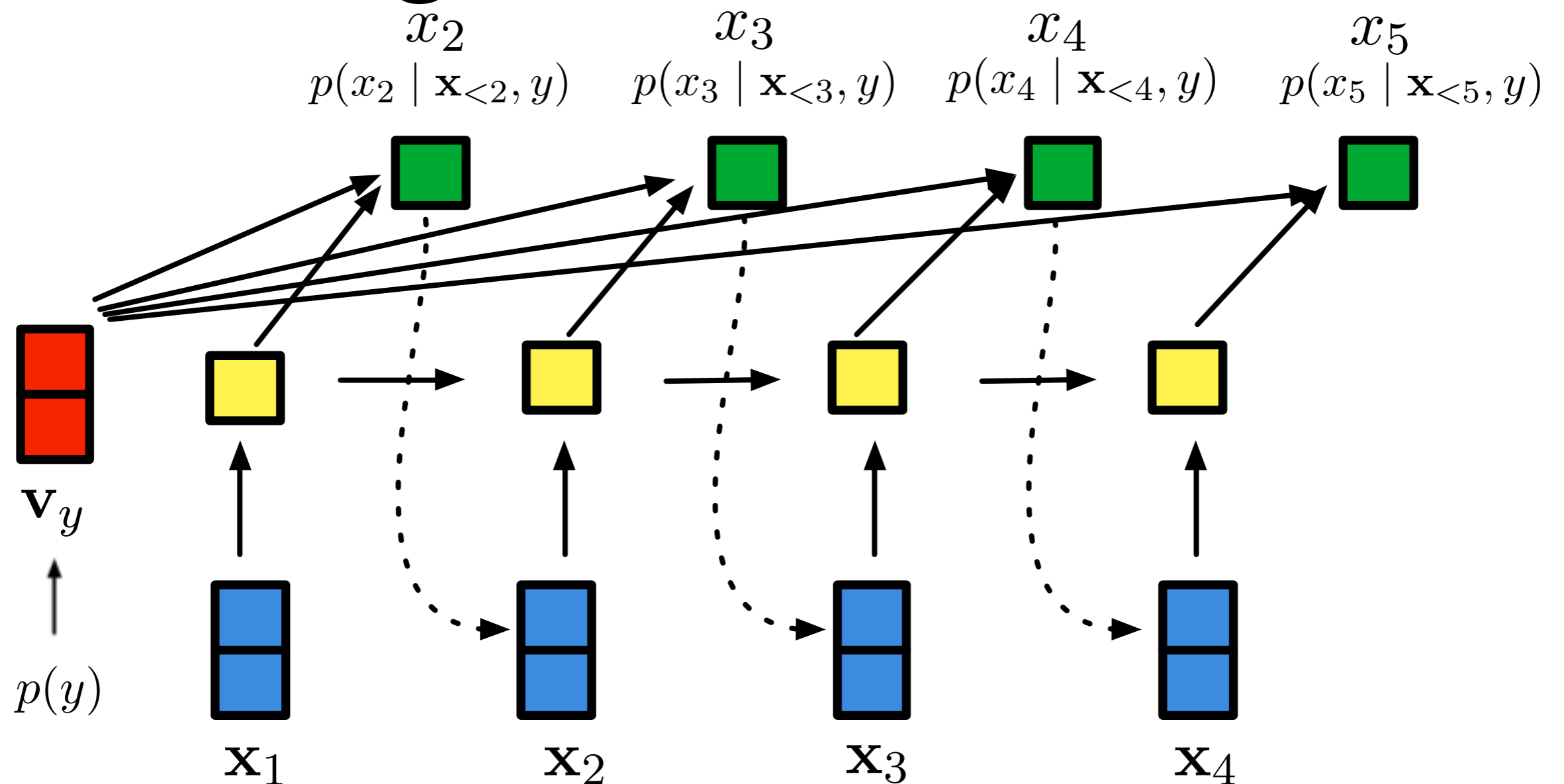
Text categorization



$$\mathcal{L}(\mathbf{W}) = \sum_i \log p(y_i | \mathbf{x}_i; \mathbf{W})$$

Generative model

Text categorization



$$\mathcal{L}(\mathbf{W}) = \sum_i \log p(\mathbf{x}_i | y_i) p(y_i)$$

Supervised text categorization

	AGNews	DBPedia	Yahoo	Yelp Binary
Bag of Words (Zhang et al., 2015)	88.8	96.6	68.9	92.2
char-CRNN (Xiao and Cho, 2016)	91.4	98.6	71.7	94.5
very deep CNN (Conneau et al., 2016)	91.3	98.7	73.4	95.7

Supervised text categorization

	AGNews	DBPedia	Yahoo	Yelp Binary
Bag of Words (Zhang et al., 2015)	88.8	96.6	68.9	92.2
char-CRNN (Xiao and Cho, 2016)	91.4	98.6	71.7	94.5
very deep CNN (Conneau et al., 2016)	91.3	98.7	73.4	95.7
Discriminative LSTM	92.1	98.7	73.7	92.6

Supervised text categorization

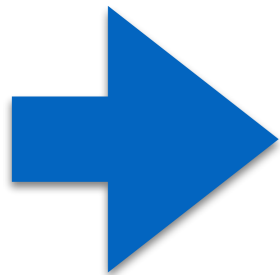
	AGNews	DBPedia	Yahoo	Yelp Binary
Bag of Words (Zhang et al., 2015)	88.8	96.6	68.9	92.2
char-CRNN (Xiao and Cho, 2016)	91.4	98.6	71.7	94.5
very deep CNN (Conneau et al., 2016)	91.3	98.7	73.4	95.7
Discriminative LSTM	92.1	98.7	73.7	92.6
Naive Bayes	90.0	96.0	68.7	86.0
Kneser-Ney Bayes	89.3	95.4	69.3	81.8

Supervised text categorization

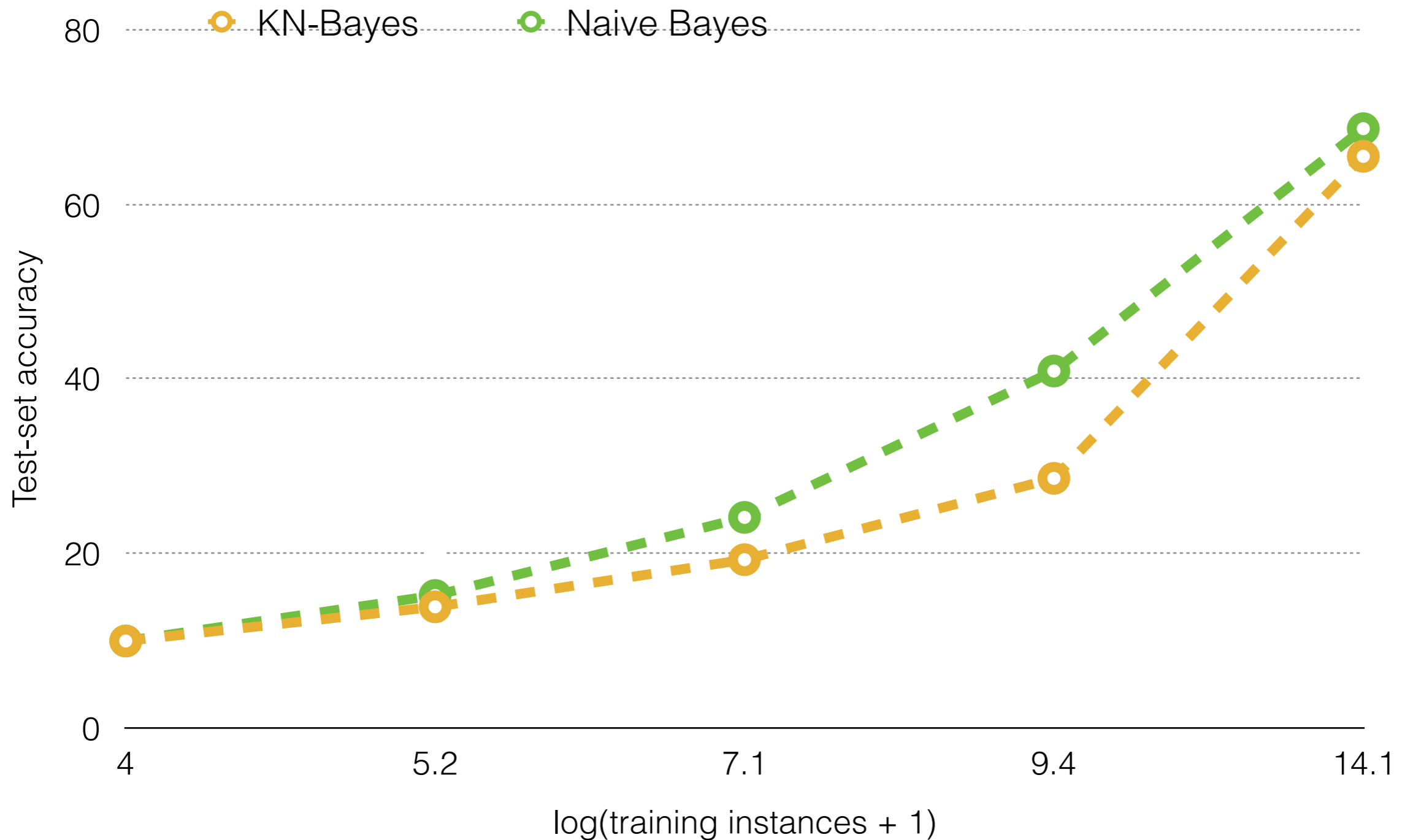
	AGNews	DBPedia	Yahoo	Yelp Binary
Bag of Words (Zhang et al., 2015)	88.8	96.6	68.9	92.2
char-CRNN (Xiao and Cho, 2016)	91.4	98.6	71.7	94.5
very deep CNN (Conneau et al., 2016)	91.3	98.7	73.4	95.7
Discriminative LSTM	92.1	98.7	73.7	92.6
Naive Bayes	90.0	96.0	68.7	86.0
Kneser-Ney Bayes	89.3	95.4	69.3	81.8
Generative LSTM	90.7	94.8	70.5	90.0

Supervised text categorization

	AGNews	DBPedia	Yahoo	Yelp Binary
Bag of Words (Zhang et al., 2015)	88.8	96.6	68.9	92.2
char-CRNN (Xiao and Cho, 2016)	91.4	98.6	71.7	94.5
very deep CNN (Conneau et al., 2016)	91.3	98.7	73.4	95.7
Discriminative LSTM	92.1	98.7	73.7	92.6
Naive Bayes	90.0	96.0	68.7	86.0
Kneser-Ney Bayes	89.3	95.4	69.3	81.8
Generative LSTM	90.7	94.8	70.5	90.0

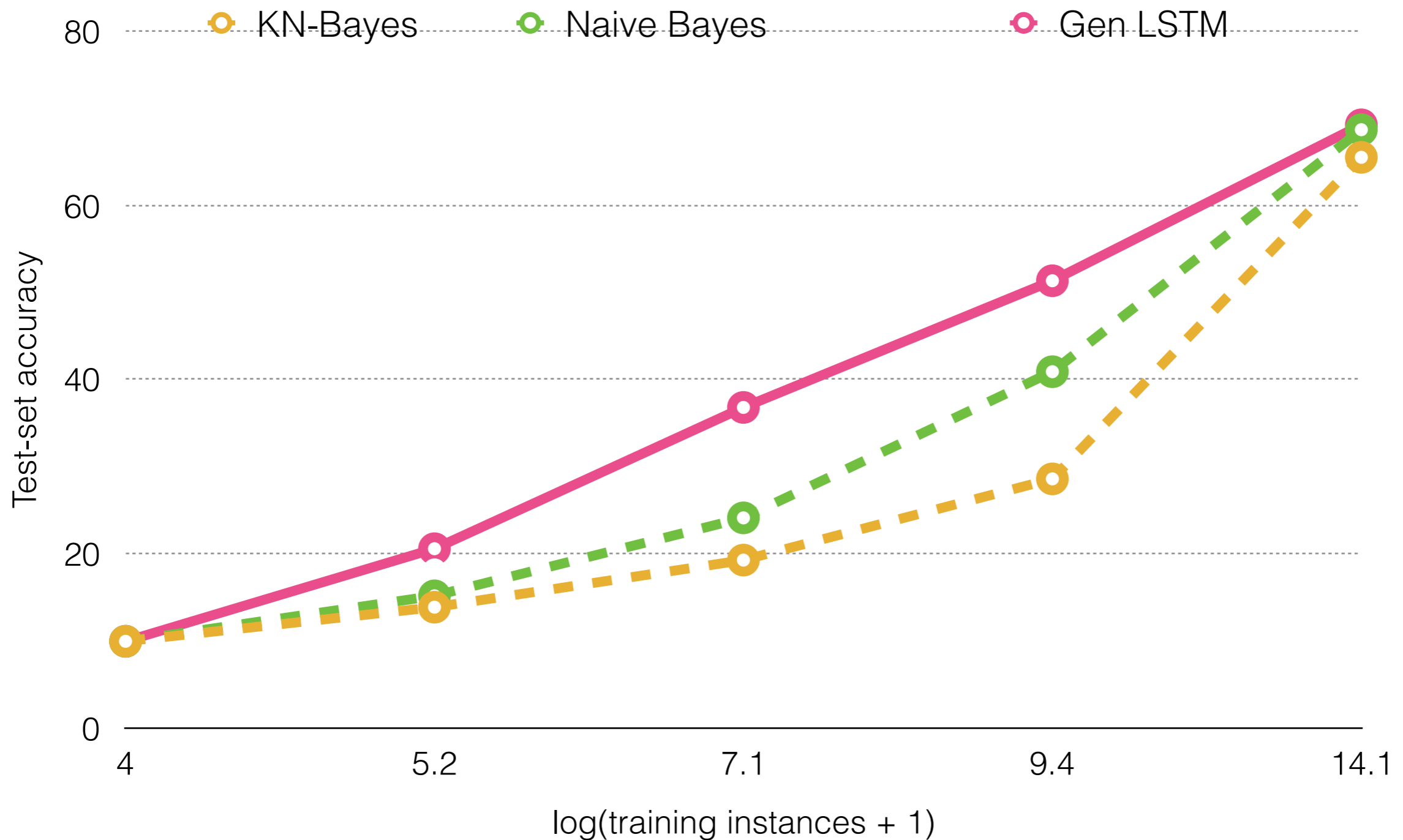


Sample efficiency



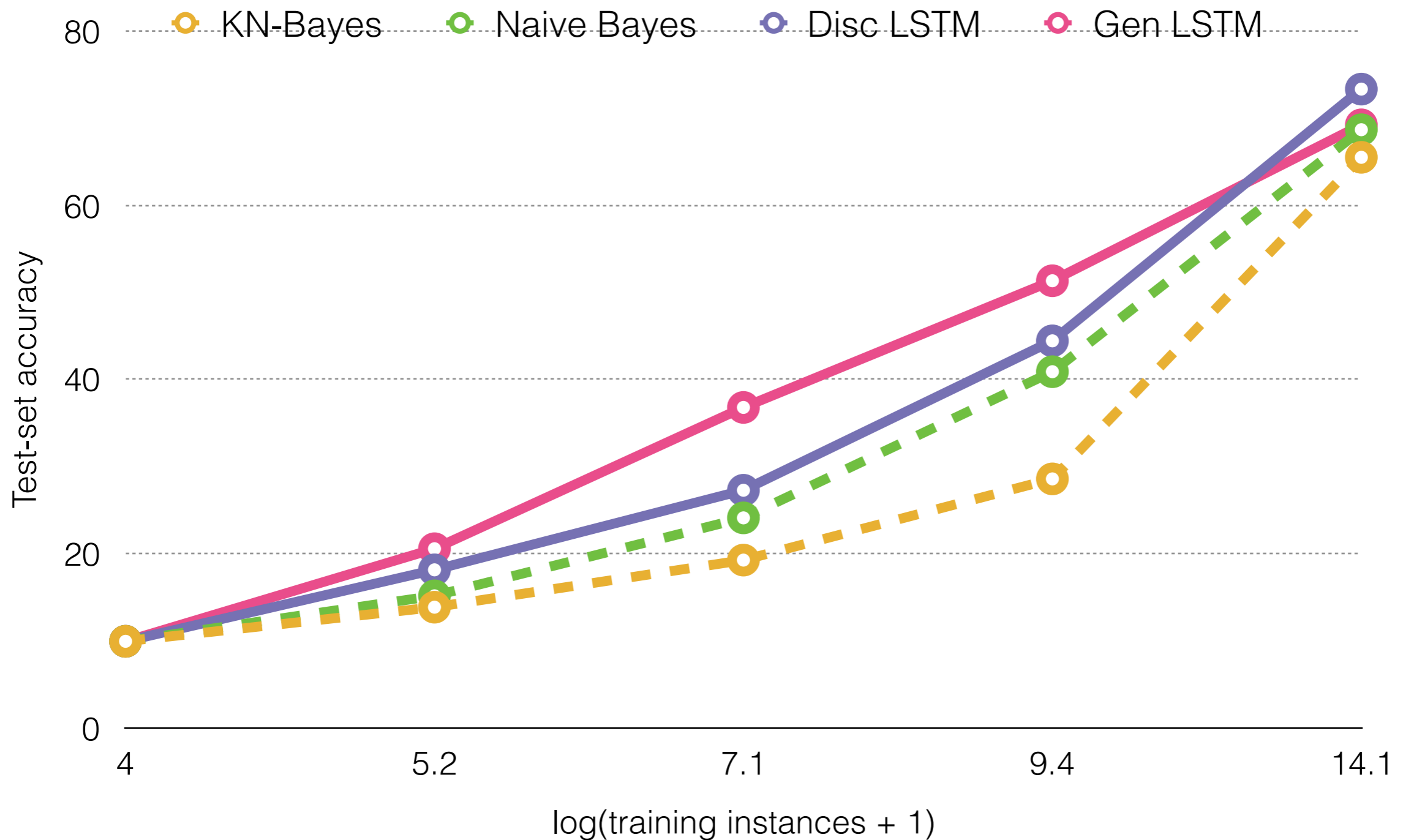
Yahoo! Answers data: 1,395,000 instances / 10 classes

Sample efficiency



Yahoo! Answers data: 1,395,000 instances / 10 classes

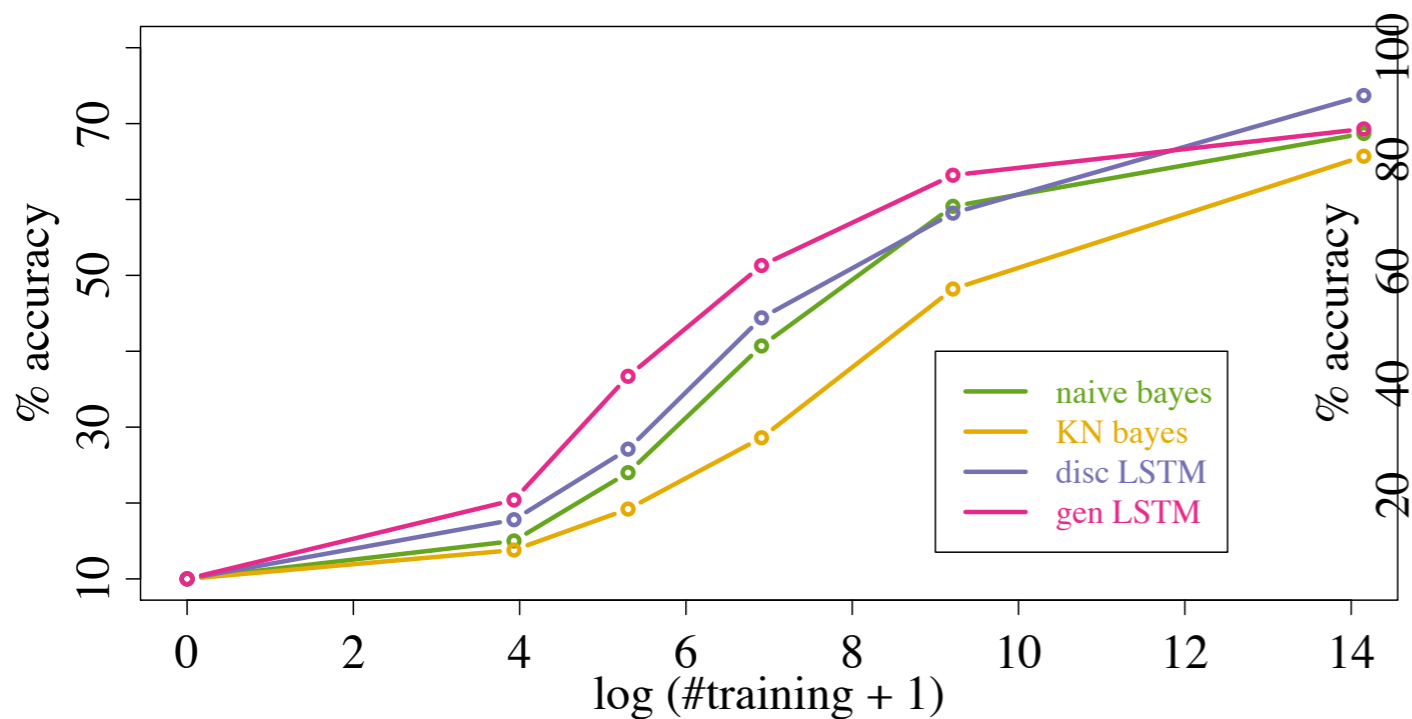
Sample efficiency



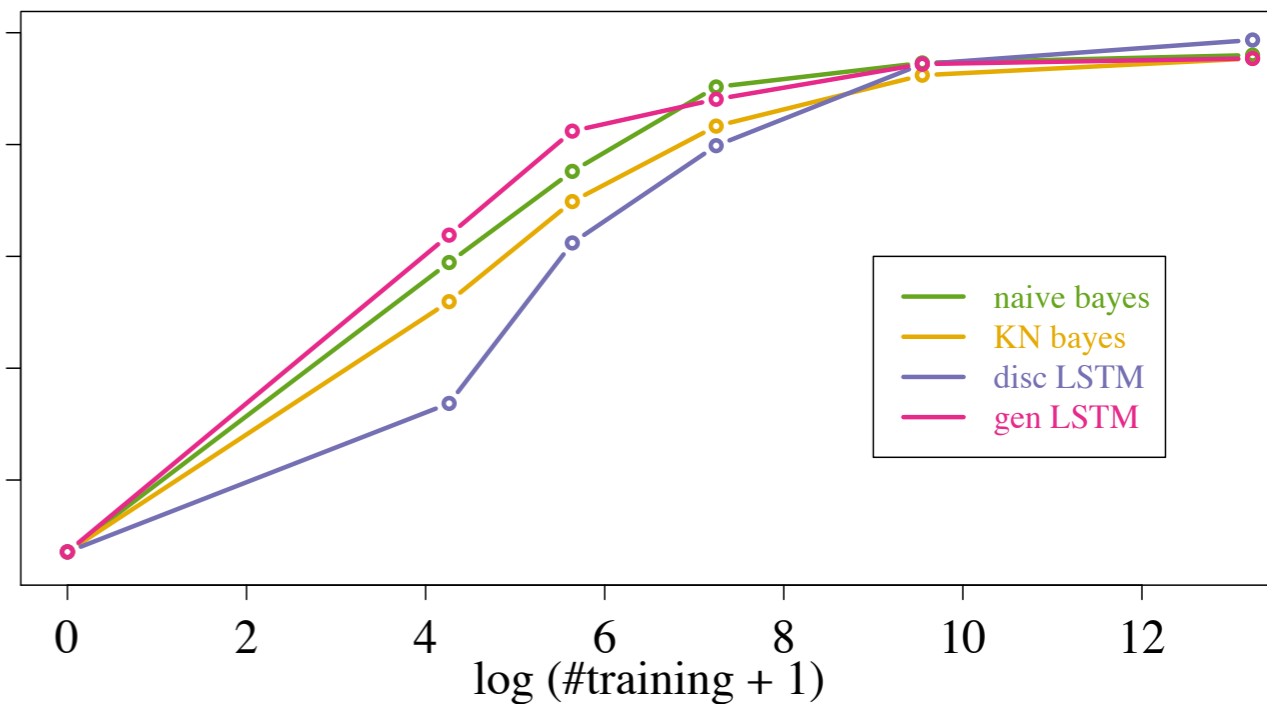
Yahoo! Answers data: 1,395,000 instances / 10 classes

Sample efficiency

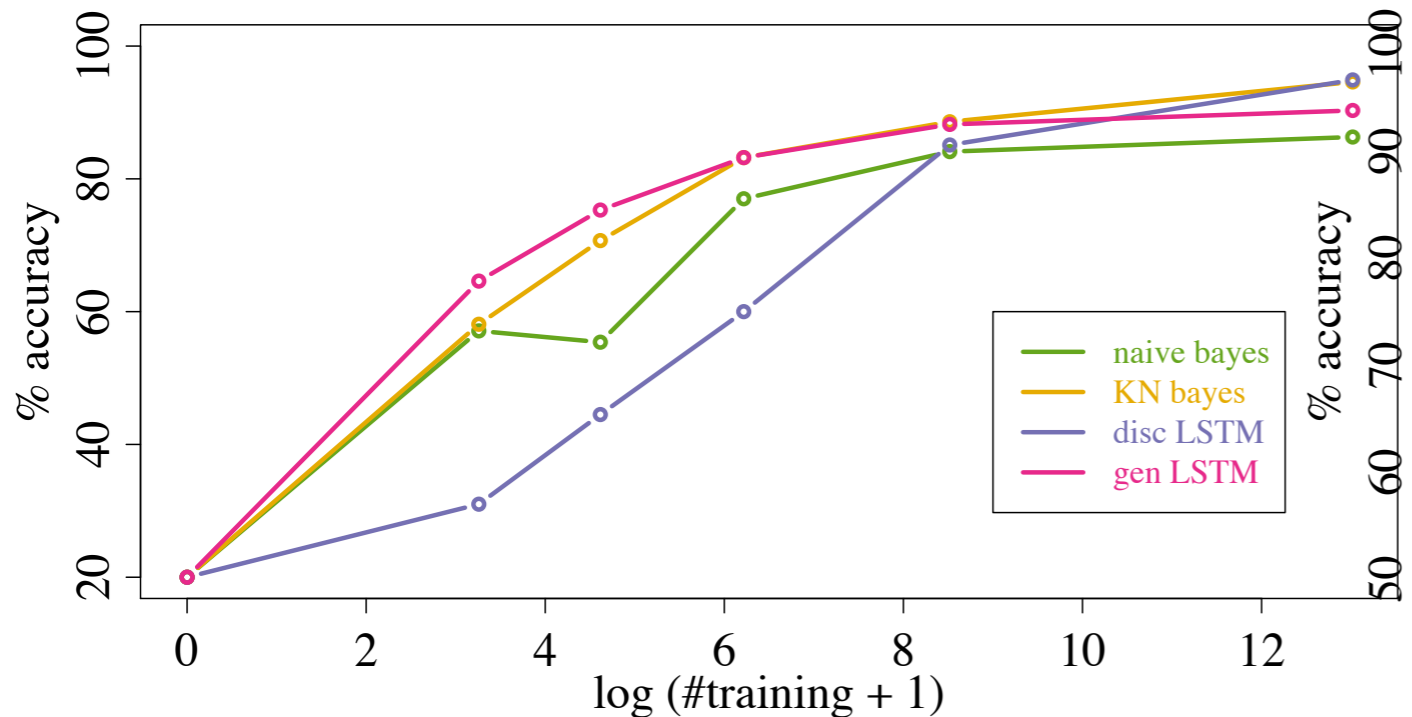
Yahoo



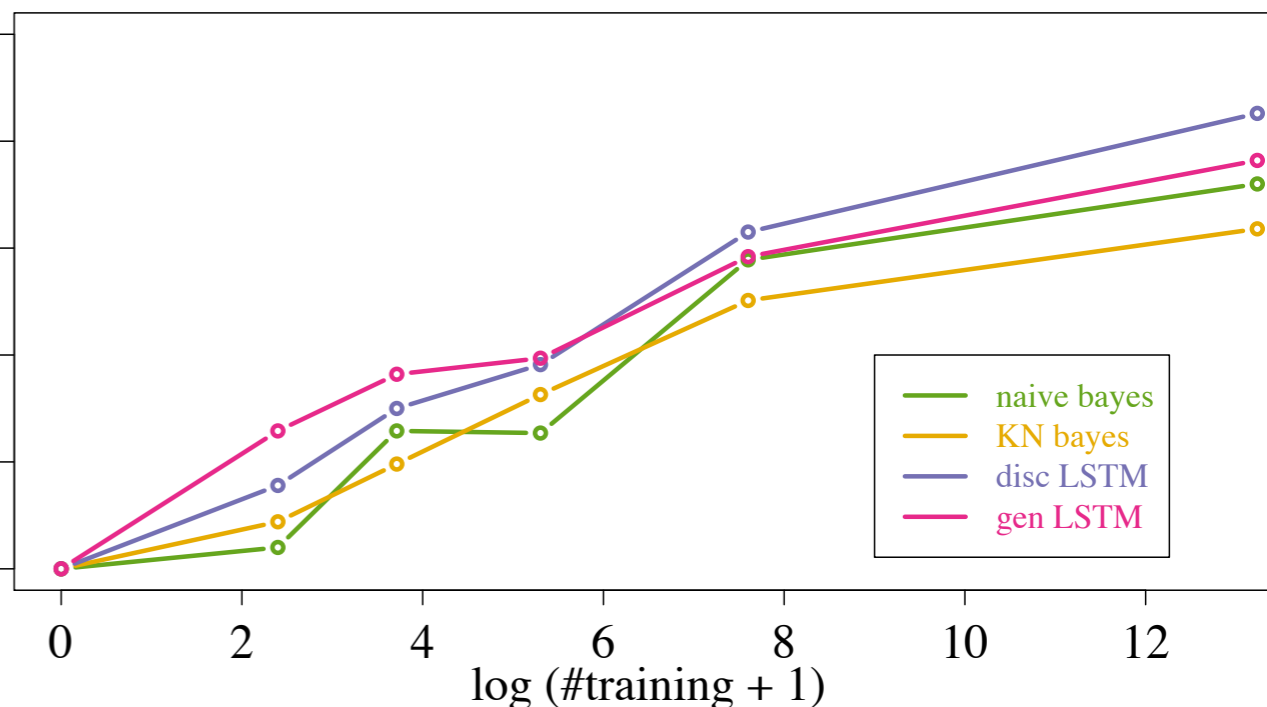
DBPedia



Sogou



Yelp Binary



Discussion

- **Generative** models of text **approach** their **asymptotic errors** more **rapidly** (better in small-data regime).
- **Discriminative** models of text have **lower asymptotic errors, faster training and inference time, and a good estimate of $p(\mathbf{x})$**
- The downside is **inference is expensive**. We have to evaluate the likelihood of the document for every class!

Case studies

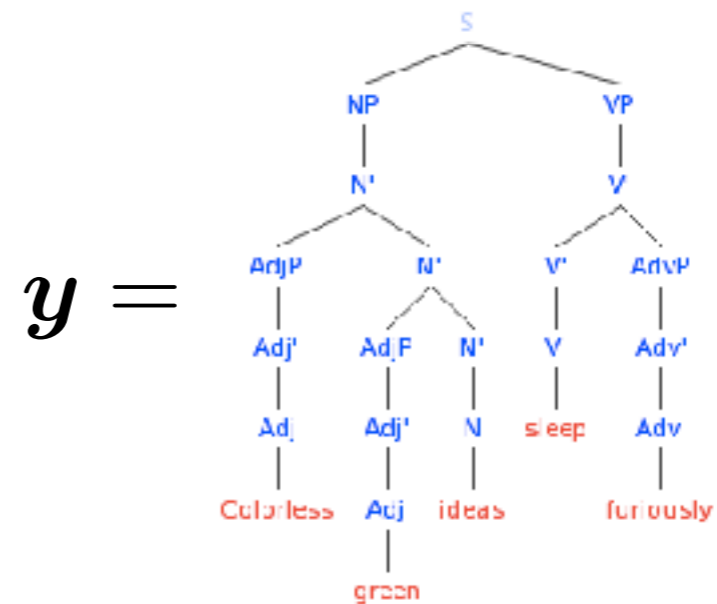
- **Text categorization**

$x =$ The image shows a screenshot of a news article snippet. The title is "US surrounds new London embassy with a moat". The text below the title describes the new US embassy in Nine Elms as a "heavily defended delicate glass box" and mentions that it is one of the world's most expensive embassies, costing a cool \$4bn. It also notes that remarkably, not a cent of US taxpayer money has been spent. The article is attributed to Ambassador William Moran, principal deputy director of the Bureau of US Overseas Buildings Operations.

$y =$ POLITICS

- **Syntactic parsing**

$x =$ Colorless green ideas
sleep furiously



- **Sequence to sequence transduction**

$x =$ Welcome to Okinawa

$y =$ 沖縄へようこそ。

Syntactic parsing

Recurrent Neural Net **Grammars**

Syntactic parsing

Recurrent Neural Net **Grammars**

- Generate **symbols** sequentially using an **RNN**

Syntactic parsing

Recurrent Neural Net **Grammars**

- Generate **symbols** sequentially using an **RNN**
- Add some **control symbols** to rewrite the history occasionally
 - Occasionally **compress** a sequence into a **constituent**
 - RNN predicts next terminal/control symbol based on the history of compressed elements and non-compressed terminals

Syntactic parsing

Recurrent Neural Net **Grammars**

- Generate **symbols** sequentially using an **RNN**
- Add some **control symbols** to rewrite the history occasionally
 - Occasionally **compress** a sequence into a **constituent**
 - RNN predicts next terminal/control symbol based on the history of compressed elements and non-compressed terminals
- This is a **top-down, left-to-right generation** of a tree+sequence (other traversal orders are possible)

(**D**, et al., ACL 2016; Kuncoro, **D**, et al., EACL 2017)

Example derivation



The hungry cat meows loudly

stack	action	probability
	NT(S)	$p(\text{NT}(\text{S}) \mid \text{TOP})$
(S	NT(NP)	$p(\text{NT}(\text{NP}) \mid (\text{S}))$
(S (NP	GEN(<i>The</i>)	$p(\text{GEN}(\textit{The}) \mid (\text{S}, (\text{NP}))$
(S (NP <i>The</i>	GEN(<i>hungry</i>)	$p(\text{GEN}(\textit{hungry}) \mid (\text{S}, (\text{NP}, \textit{The}))$
(S (NP <i>The hungry</i>	GEN(<i>cat</i>)	$p(\text{GEN}(\textit{cat}) \mid \dots)$
(S (NP <i>The hungry cat</i>	REDUCE	$p(\text{REDUCE} \mid \dots)$
(S (NP <i>The hungry cat</i>)		
(S (NP <i>The hungry cat</i>)		

Compress “The hungry cat”
into a single composite symbol

stack	action	probability
	NT(S)	$p(\text{NT}(\text{S}) \mid \text{TOP})$
(S	NT(NP)	$p(\text{NT}(\text{NP}) \mid (\text{S}))$
(S (NP	GEN(<i>The</i>)	$p(\text{GEN}(\textit{The}) \mid (\text{S}, (\text{NP}))$
(S (NP <i>The</i>	GEN(<i>hungry</i>)	$p(\text{GEN}(\textit{hungry}) \mid (\text{S}, (\text{NP}, \textit{The}))$
(S (NP <i>The hungry</i>	GEN(<i>cat</i>)	$p(\text{GEN}(\textit{cat}) \mid \dots)$
(S (NP <i>The hungry cat</i>	REDUCE	$p(\text{REDUCE} \mid \dots)$
(S (NP <i>The hungry cat</i>)	NT(VP)	$p(\text{NT}(\text{VP}) \mid (\text{S}, (\text{NP } \textit{The hungry cat}))$
(S (NP <i>The hungry cat</i>) (VP	GEN(<i>meows</i>)	
(S (NP <i>The hungry cat</i>) (VP <i>meows</i>	REDUCE	
(S (NP <i>The hungry cat</i>) (VP <i>meows</i>)	GEN(.)	
(S (NP <i>The hungry cat</i>) (VP <i>meows</i>) .	REDUCE	
(S (NP <i>The hungry cat</i>) (VP <i>meows</i>) .)		

Deriving the model

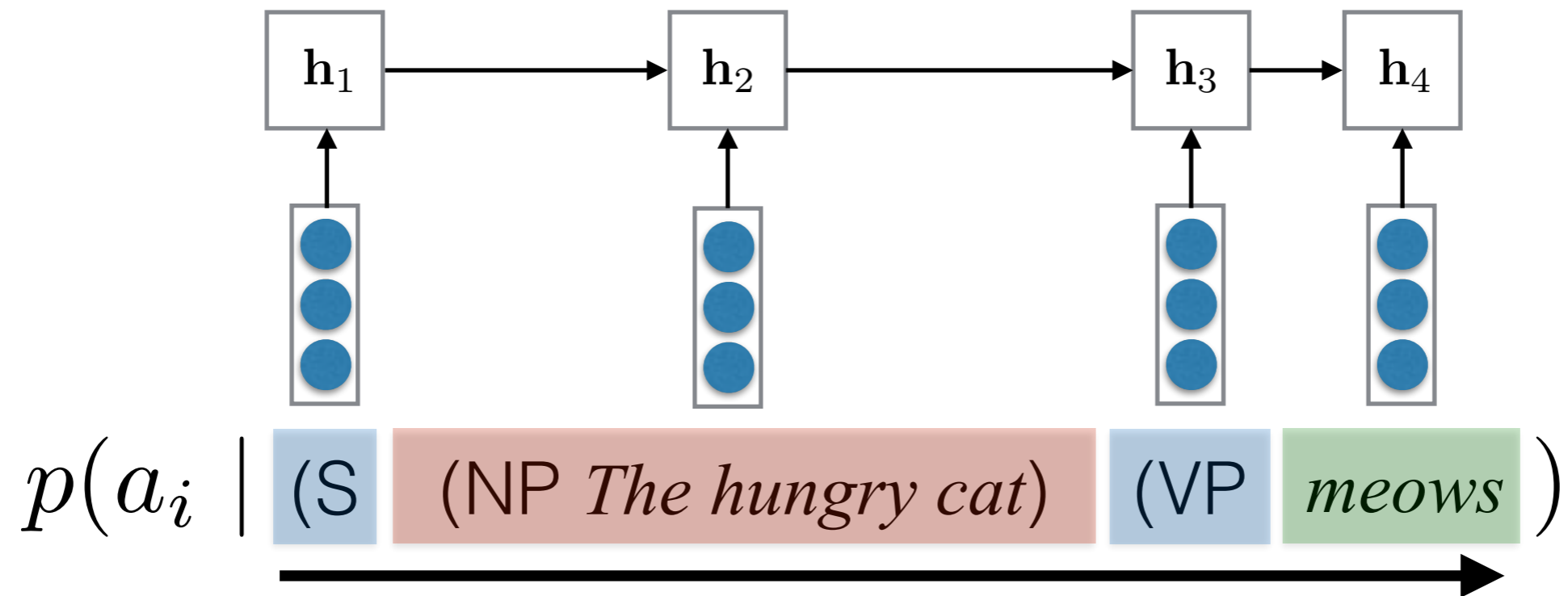
- Valid (***tree***, ***string***) pairs are in bijection to valid sequences of actions (specifically, the DFS, left-to-right traversal of the trees)
- Every stack configuration perfectly encodes the complete history of actions.
- Therefore, the probability decomposition is justified by the chain rule, i.e.

$$p(\mathbf{x}, \mathbf{y}) = p(\mathit{actions}(\mathbf{x}, \mathbf{y})) \quad (\text{prop 1})$$

$$p(\mathit{actions}(\mathbf{x}, \mathbf{y})) = \prod_i p(a_i \mid \mathbf{a}_{<i}) \quad (\text{chain rule})$$

$$= \prod_i p(a_i \mid \mathit{stack}(\mathbf{a}_{<i})) \quad (\text{prop 2})$$

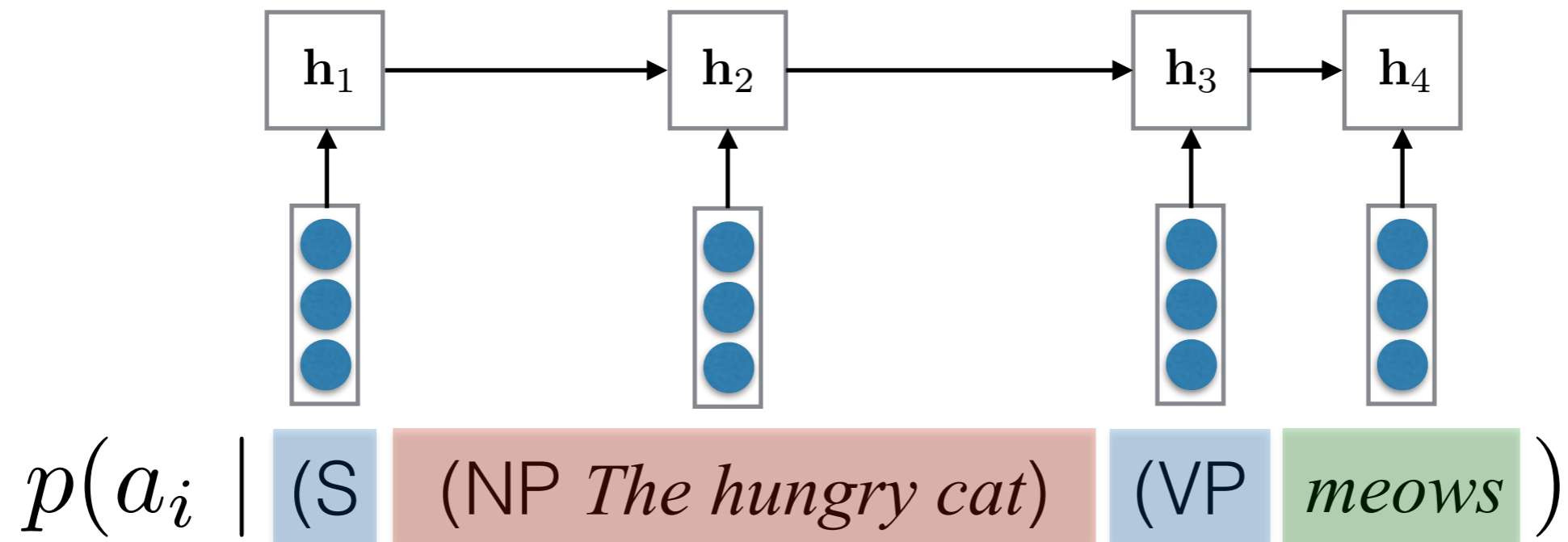
Modeling the next action



1. unbounded depth

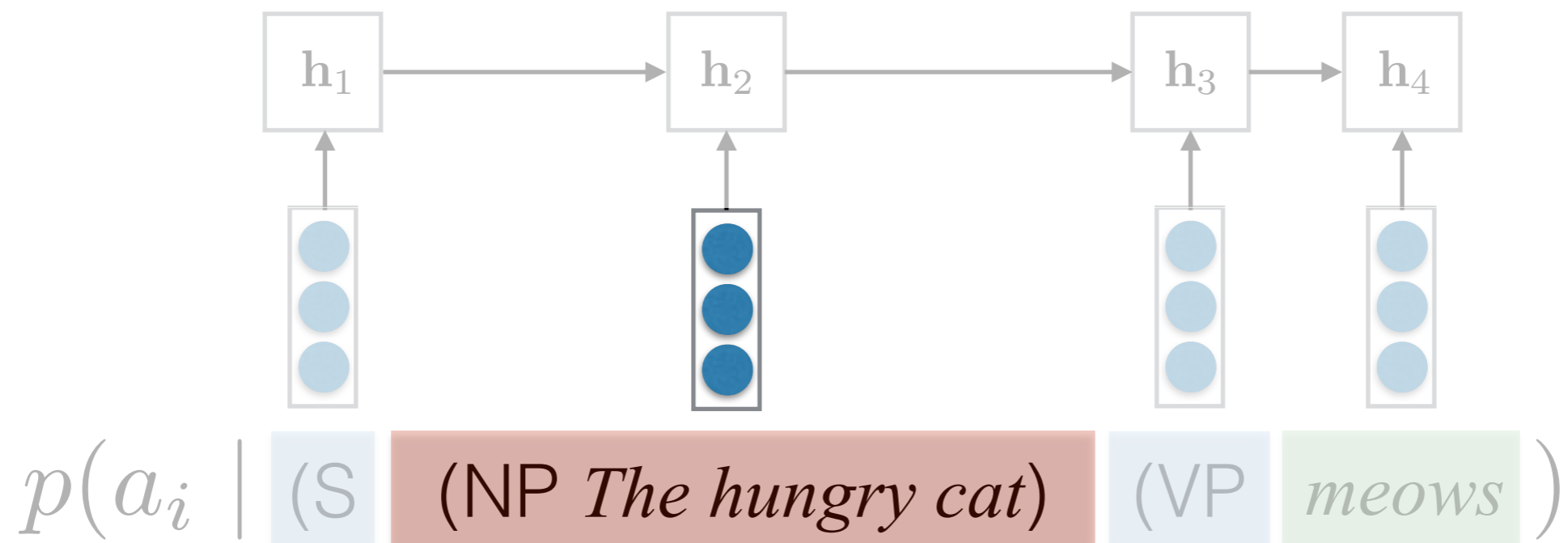
1. Unbounded depth \rightarrow recurrent neural nets

Modeling the next action



1. Unbounded depth \rightarrow recurrent neural nets

Modeling the next action

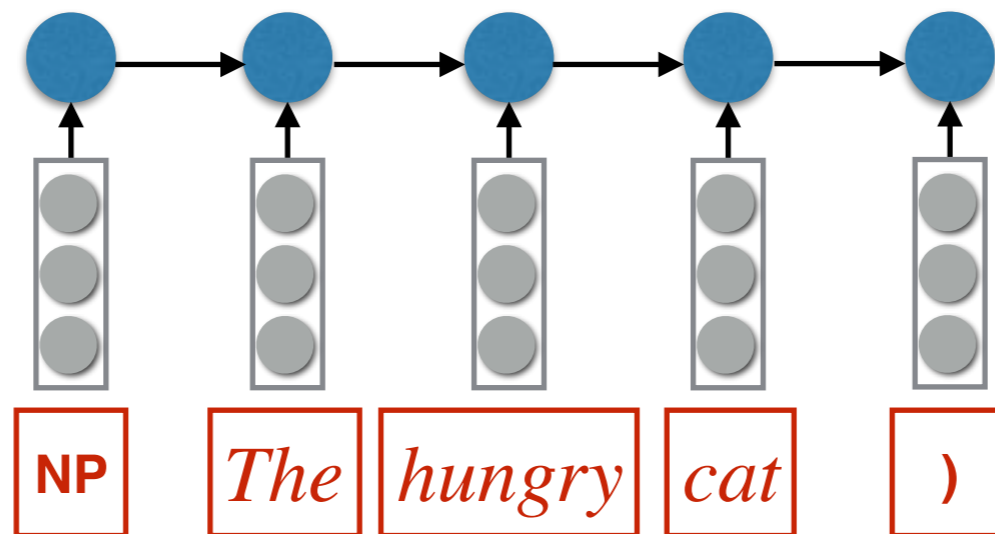


2. arbitrarily complex trees

1. Unbounded depth \rightarrow recurrent neural nets
2. Arbitrarily complex trees \rightarrow recursive neural nets

Syntactic composition

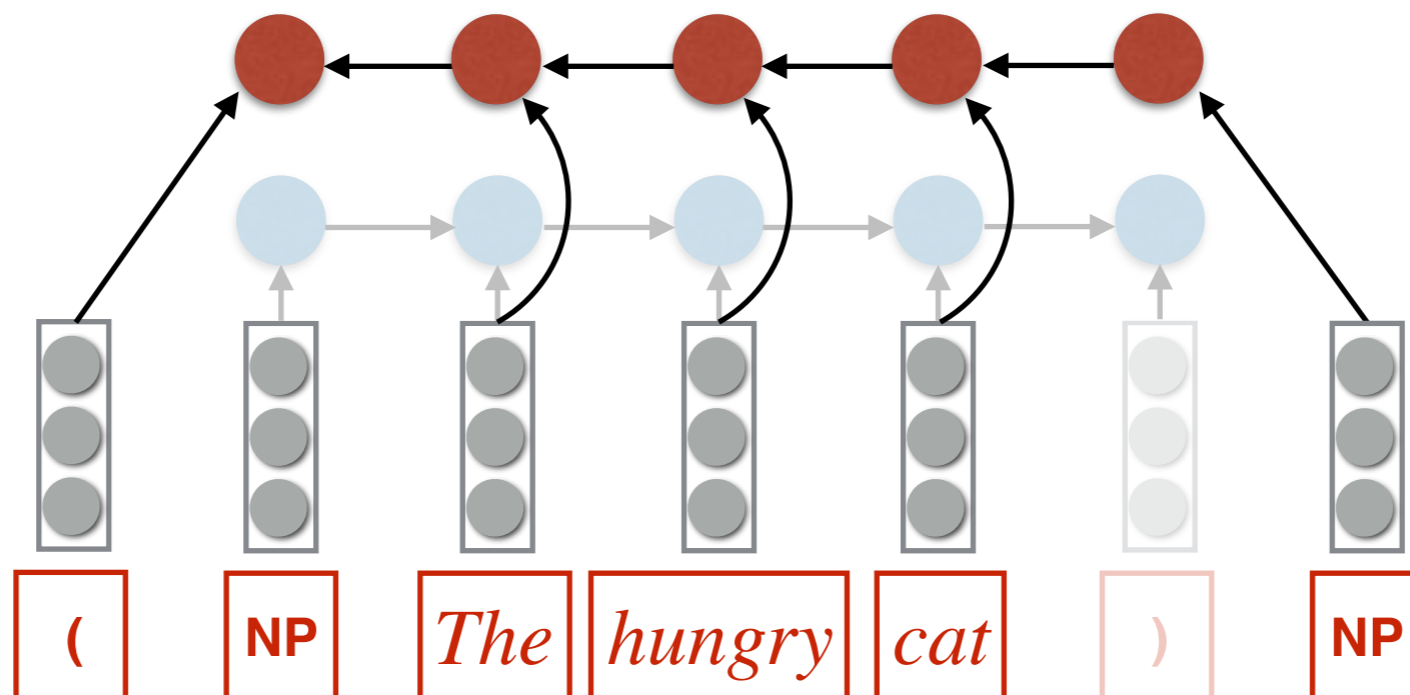
Need representation for: (NP *The hungry cat*)



What head type? ↗

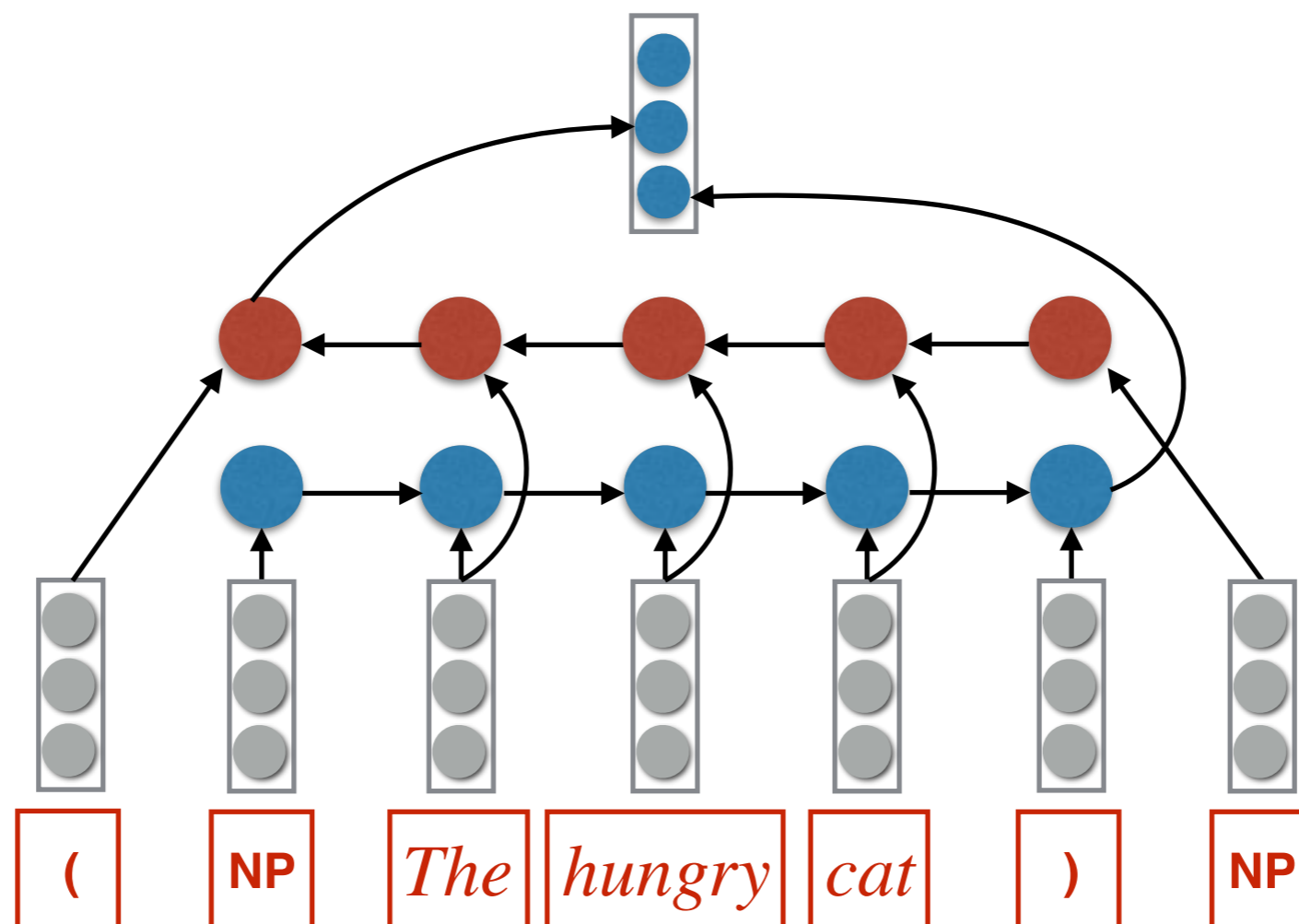
Syntactic composition

Need representation for: (NP *The hungry cat*)



Syntactic composition

Need representation for: (NP *The hungry cat*)



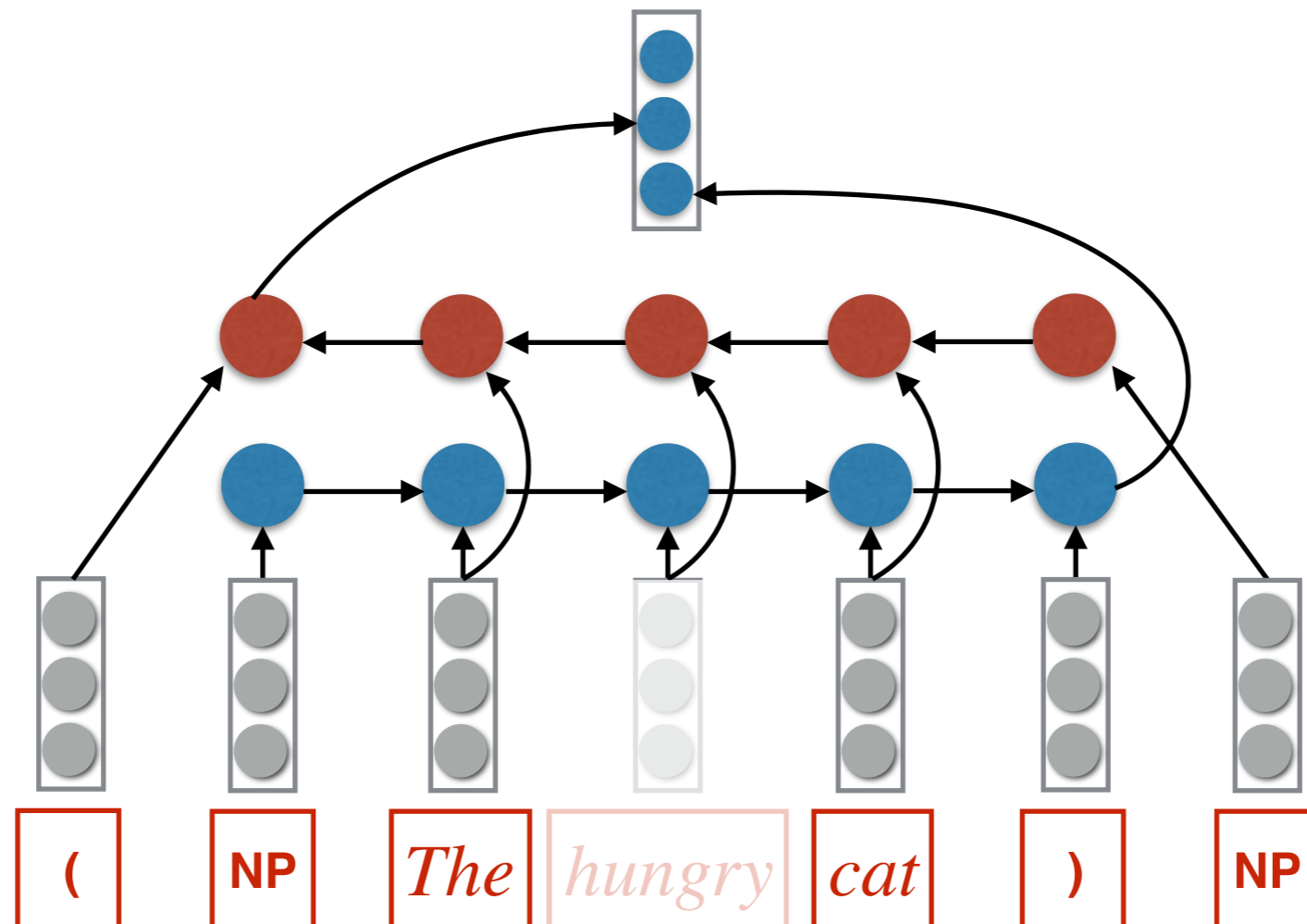
Syntactic composition

Recursion

Need representation for:

(NP *The hungry cat*)

(NP *The (ADJP very hungry) cat*)



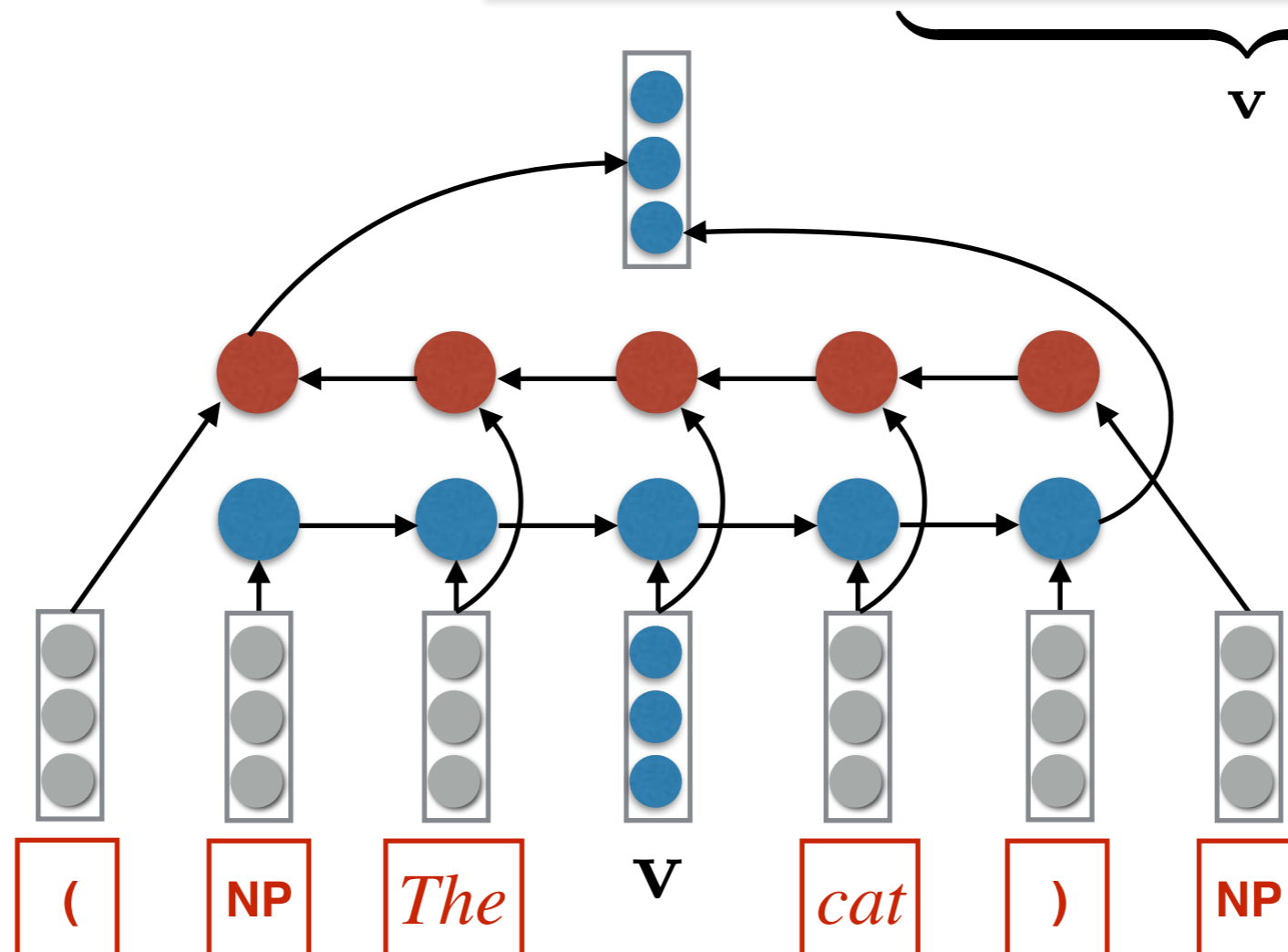
Syntactic composition

Recursion

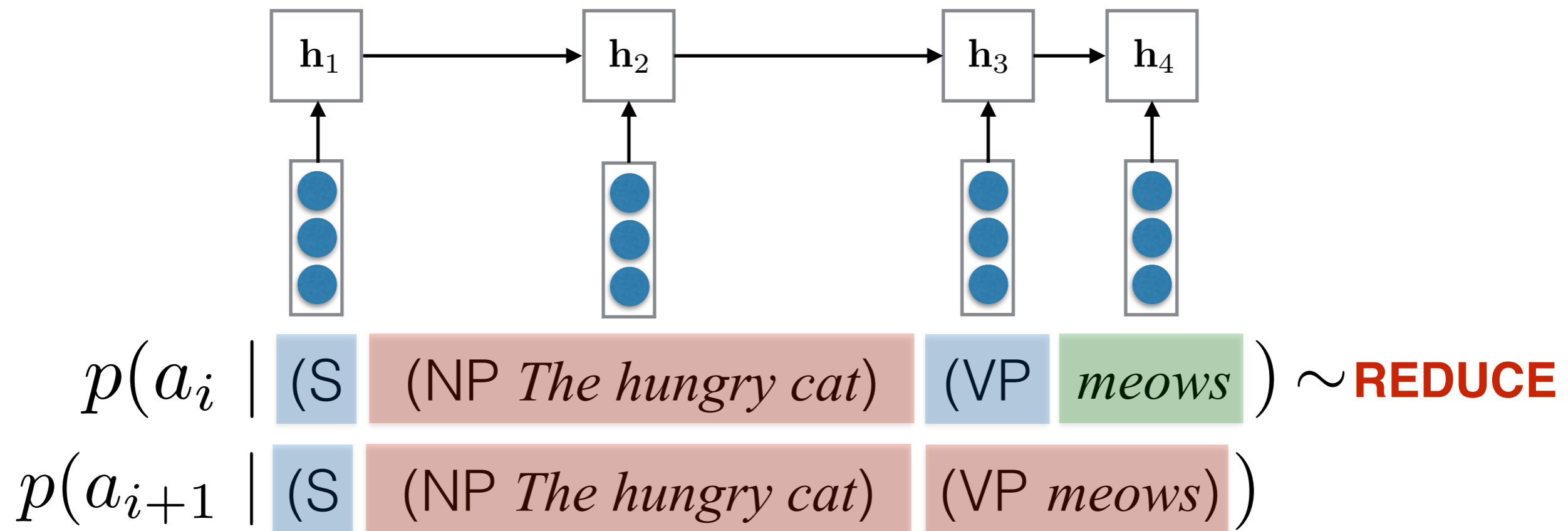
Need representation for:

(NP *The hungry cat*)

(NP *The (ADJP *very hungry*) cat*)



Modeling the next action




3. limited updates

1. Unbounded depth \rightarrow recurrent neural nets
2. Arbitrarily complex trees \rightarrow recursive neural nets
3. Limited updates to state \rightarrow stack RNNs

Inference

- In text categorization, it was not really a problem to exhaustively evaluate all candidate y 's.
- Here, we can't do that — we have $O(2^{|\mathbf{x}|})$ candidates!
- Outline of the solution
 - Learn a tractable instrumental distribution, $q(\mathbf{y} | \mathbf{x})$, which approximates the posterior over trees
 - Use **importance sampling** to solve the inference problems (maximization, marginalization) we care about

Results: Parsing

	Type	F1
Petrov and Klein (2007)	Gen	90.1
Shindo et al (2012) Single model	Gen	91.1
Vinyals et al (2015) PTB only	Disc	90.5
Shindo et al (2012) Ensemble	<i>Gen+Ensemble</i>	92.4
Vinyals et al (2015) Semisupervised	<i>Disc+SemiSup</i>	92.8
 Discriminative PTB only	Disc	91.7
Generative PTB only	Gen	93.6

Results: Parsing

	Type	F1
Petrov and Klein (2007)	Gen	90.1
Shindo et al (2012) Single model	Gen	91.1
Vinyals et al (2015) PTB only	Disc	90.5
Shindo et al (2012) Ensemble	<i>Gen+Ensemble</i>	92.4
Vinyals et al (2015) Semisupervised	<i>Disc+SemiSup</i>	92.8
Discriminative PTB only	Disc	91.7
Generative PTB only	Gen	93.6
Choe and Charniak (2016) Semisupervised	Gen <i>+SemiSup</i>	93.8
Fried et al. (2017)	<i>Gen+Semi +Ensemble</i>	94.7

Discussion

- RNNs are effective both for **modeling language** and **parsing**
- Generative parser outperforms discriminative parser
- Expectation: the discriminative model would do better with more data
- We are in the “generative” regime!

Case studies

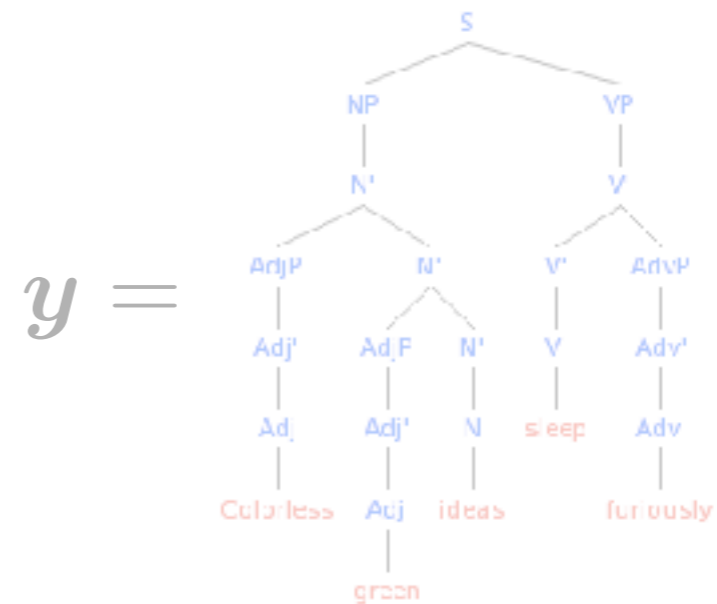
- **Text categorization**

$x =$ The image shows a screenshot of a news article snippet. The title is "US surrounds new London embassy with a moat". The text below the title describes the new US embassy in Nine Elms as a heavily defended, delicate glass box, a "crystaline radiant beacon" in fact, that resembles a corporate cube. It is also noted as one of the world's most expensive embassies, costing a cool \$4bn. Remarkably, not a cent of US taxpayer money has been spent. Speaking at the press launch on Wednesday, Ambassador William Miller, principal deputy director of the Bureau of US Overseas Buildings Operations, confirmed that the new building "has actively funded from the proceeds of real estate sales".

$y =$ POLITICS

- **Syntactic parsing**

$x =$ Colorless green ideas
sleep furiously



- **Sequence to sequence transduction**

$x =$ Welcome to Okinawa

$y =$ 沖縄へようこそ。

Seq2Seq Modeling

Direct model

$$\begin{aligned} p(\mathbf{y} \mid \mathbf{x}) &= \text{ConditionalRNNLM}(\mathbf{x}) \\ &= \prod_i p(y_i \mid \mathbf{x}, \mathbf{y}_{<i}) \end{aligned}$$

- State of the art performance in most applications
- Two serious problems that concern us:
 - Nontrivial to use “unpaired” samples of \mathbf{x} or \mathbf{y} to train the model
 - “Explaining away effects” - models like this learn to ignore “inconvenient” inputs (i.e., \mathbf{x}), in favor of high probability continuations of an output prefix ($\mathbf{y}_{<i}$)

Seq2Seq Modeling

What is label bias?

Label bias is a species of “explaining away” that causes trouble in directed (locally normalized) models.

a b c → x y z

a b c' → x y z

a b' c → x y z

d → w

a b' d → x y z

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} \mid \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} \mid \mathbf{y})$$

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model” “Channel model”

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model” “Channel model”

$$\mathbf{y} \sim p(\mathbf{y})$$

The world is colorful because of the Internet.

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model”

“Channel model”

$$\mathbf{y} \sim p(\mathbf{y})$$

The world is colorful because of the Internet.



Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model”

“Channel model”

$$\mathbf{y} \sim p(\mathbf{y})$$

The world is colorful because of the Internet.



$$\mathbf{x} \sim p(\mathbf{x} | \mathbf{y})$$

世界はインターネットのためにカラフルです。

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model”

“Channel model”

$$\mathbf{y} \sim p(\mathbf{y})$$

The world is colorful because of the Internet.



$$\mathbf{x} \sim p(\mathbf{x} | \mathbf{y})$$

世界はインターネットのためにカラフルです。



Source model can be estimated from unpaired \mathbf{y} 's

Seq2Seq Modeling

Generative model

$$p(\mathbf{y} | \mathbf{x}) \propto p(\mathbf{y}) \times p(\mathbf{x} | \mathbf{y})$$

“Source model”

“Channel model”

$$\mathbf{y} \sim p(\mathbf{y})$$

The world is colorful because of the Internet.



$$\mathbf{x} \sim p(\mathbf{x} | \mathbf{y})$$

世界はインターネットのためにカラフルです。

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} p(\mathbf{y})p(\mathbf{x} | \mathbf{y})$$

👍 Inference model form avoids explaining away of inputs (“label bias”).

Seq2Seq Modeling

Generative model

- Question: Can we use **neural network component models** without bad independence assumptions?
 - **Training — straightforward**
 - **Decoding — challenging**

Decoding

- Some bad initial results
 - The IS algorithm we proposed hurt us unless the number of samples (k) was massive
 - Reranking an k -best list from a direct model didn't help unless k was even bigger
- Question: **can we develop a left-to-right decoder for a noisy channel MT model?**

Decoding

Direct vs. generative

Direct model:

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \mathbf{x}, \hat{\mathbf{y}}_{<i})$$
$$i \leftarrow i + 1$$

Decoding

Direct vs. generative

Direct model:

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \mathbf{x}, \hat{\mathbf{y}}_{<i})$$

$$i \leftarrow i + 1$$



Chain rule!

Decoding

Direct vs. generative

Direct model:

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \mathbf{x}, \hat{\mathbf{y}}_{<i})$$

$$i \leftarrow i + 1 \quad \text{👍 Chain rule!}$$

Not perfect, but $\hat{\mathbf{y}} \approx \arg \max_{\mathbf{y}} p(\mathbf{y} \mid \mathbf{x})$

(Compare to using greedy decoding with MEMMs)

Decoding

Direct vs. generative

Generative model (naive):

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \hat{\mathbf{y}}_{<i}) p(\mathbf{x} \mid \hat{\mathbf{y}}_{<i}, y)$$

$$i \leftarrow i + 1$$

Decoding

Direct vs. generative

Generative model (naive):

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \hat{\mathbf{y}}_{<i}) p(\mathbf{x} \mid \hat{\mathbf{y}}_{<i}, y)$$

$$i \leftarrow i + 1$$



Chain rule!

Decoding

Direct vs. generative

Generative model (naive):

while $y_i \neq \text{STOP}$:

$$\hat{y}_i = \arg \max_y p(y \mid \hat{\mathbf{y}}_{<i}) p(\mathbf{x} \mid \hat{\mathbf{y}}_{<i}, y)$$

$$i \leftarrow i + 1$$



Probability doesn't work like this.

Decoding

Direct vs. generative

Outline of solution:

Introduce a latent variable \mathbf{z} that determines when enough of the conditioning context has been read to generate another symbol

$$p(\mathbf{x} | \mathbf{y}) = \sum_{\mathbf{z}} p(\mathbf{x}, \mathbf{z} | \mathbf{y})$$

$$p(\mathbf{x}, \mathbf{z} | \mathbf{y}) \approx \prod_{j=1}^{|\mathbf{x}|} \underbrace{p(z_j | z_{j-1}, \mathbf{y}_1^{z_j}, \mathbf{x}_1^{j-1})}_{\text{alignment probability}} \underbrace{p(x_j | \mathbf{y}_1^{z_j}, \mathbf{x}_1^{j-1})}_{\text{word probability}}$$

How much of \mathbf{y} do we need to read to model the j^{th} token of \mathbf{x} ?

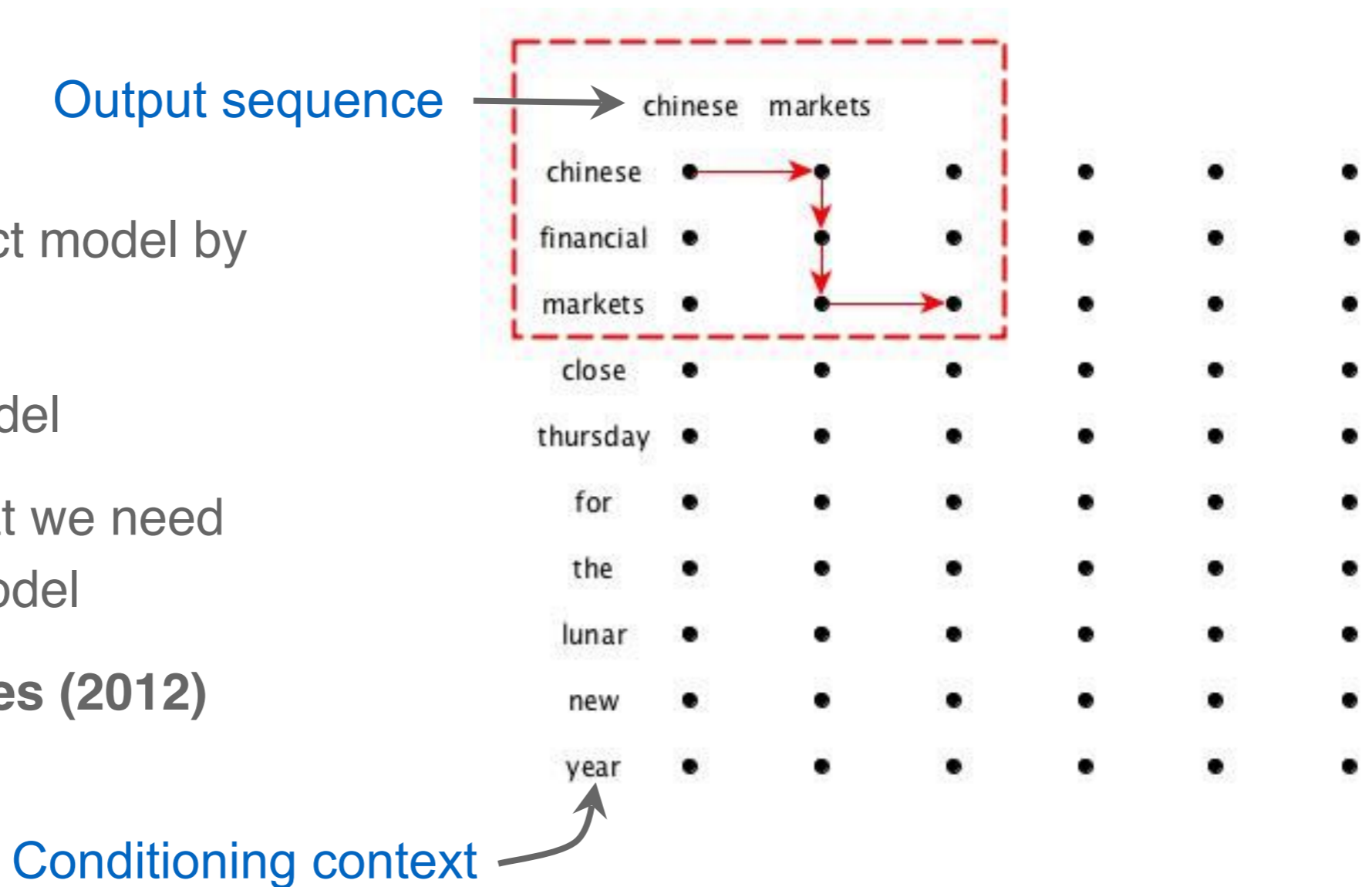
The Segment to Segment Model

Introduced as a direct model by
Yu et al. (2016)

It's a good direct model

It also is exactly what we need
for the channel model

Similar model: **Graves (2012)**



Decoding with an auxiliary model

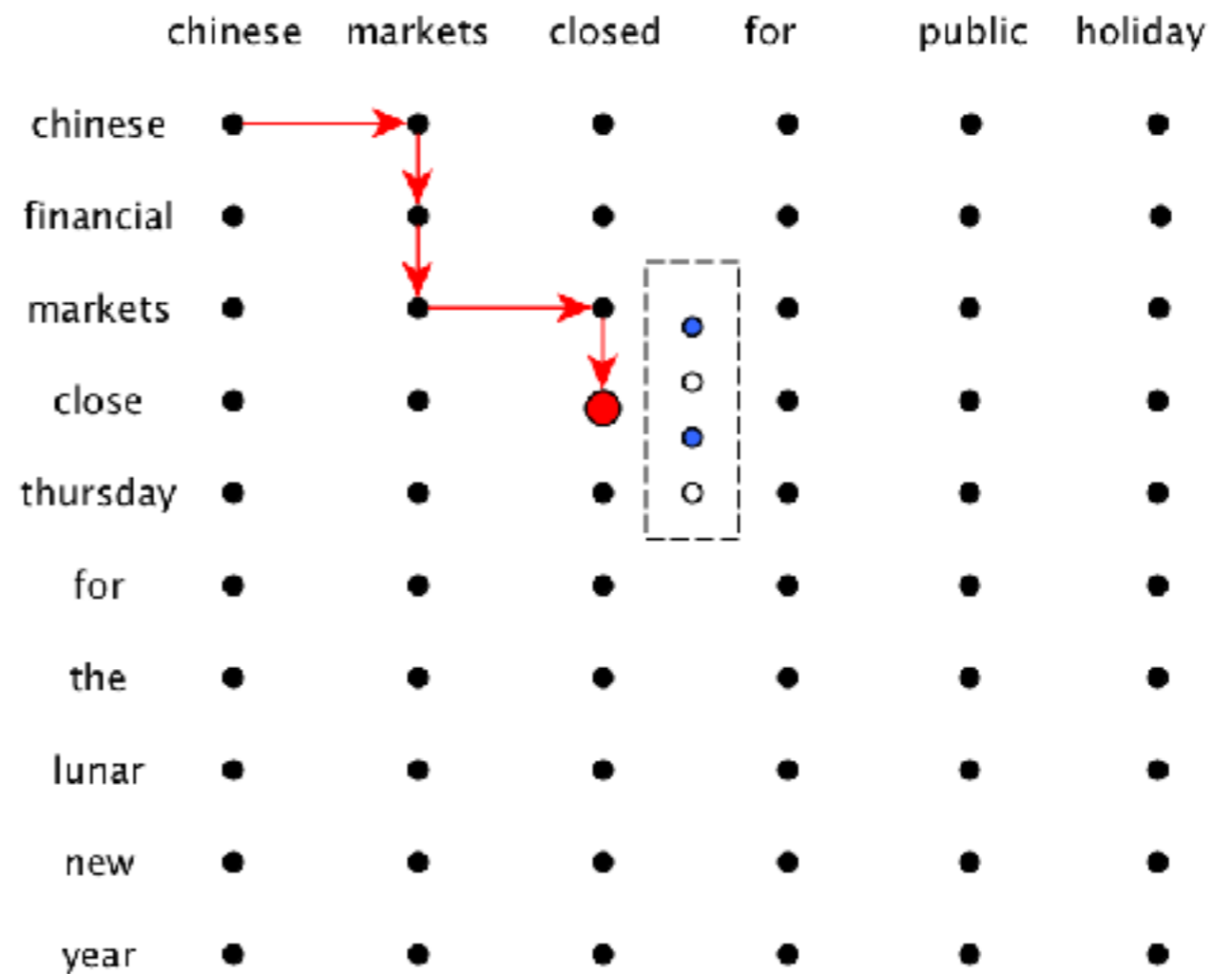
Possible proposals:

Chinese markets open

Chinese markets closed

Market close

Financial markets



Decoding with an auxiliary model

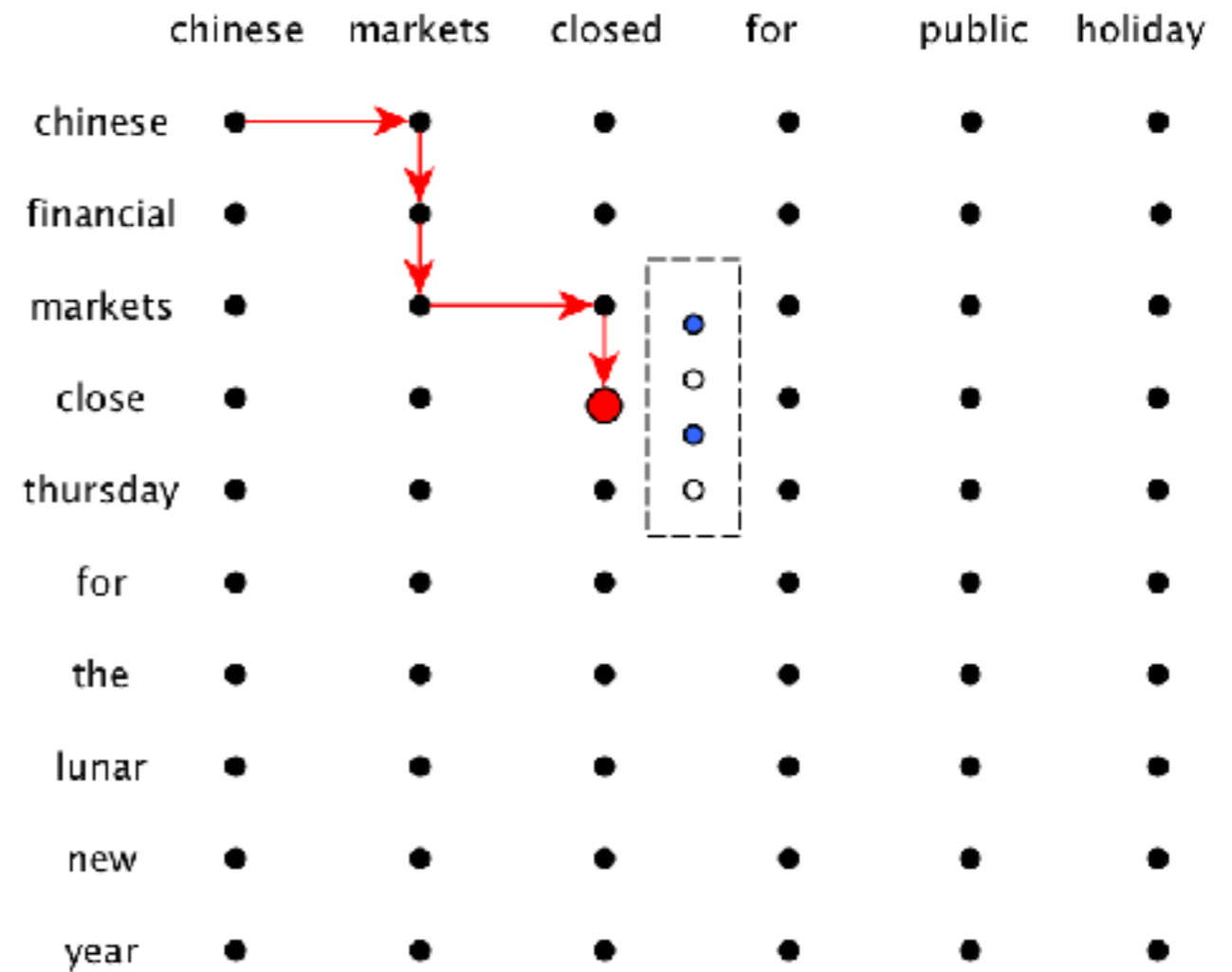
Possible proposals:

Chinese markets open

Chinese markets closed

Market close

Financial markets



Expanded objective

$$O_{\mathbf{x}_1^i, \mathbf{y}_1^j} = \lambda_1 \log p(\mathbf{y}_1^j | \mathbf{x}_1^i) + \lambda_2 \log p(\mathbf{x}_1^i | \mathbf{y}_1^j) + \lambda_3 \log p(\mathbf{y}_1^j) + \lambda_4 |\mathbf{y}_1^j|.$$

Experiments

Machine translation

- Medium-sized Chinese-English news parallel data
- Large LSTM language model trained on English news + target side of parallel data
- Evaluation using BLEU-4 (higher is better)

Experiments

Machine translation

Gen Discriminative

Model	BLEU
Seq2seq with attention	25.27
Direct model (q by itself)	23.33
Direct + LM + bias	23.33
Channel + LM + bias	26.28
Direct + channel + LM + bias	26.44

Experiments

Machine translation

	Model	BLEU
Discriminative	Seq2seq with attention	25.27
	Direct model (q by itself)	23.33
	Direct + LM + bias	23.33
Gen	Channel + LM + bias	26.28
	Direct + channel + LM + bias	26.44

Experiments

Machine translation

Gen Discriminative

Model	BLEU
Seq2seq with attention	25.27
Direct model (q by itself)	23.33
Direct + LM + bias	23.33
Channel + LM + bias	26.28
Direct + channel + LM + bias	26.44

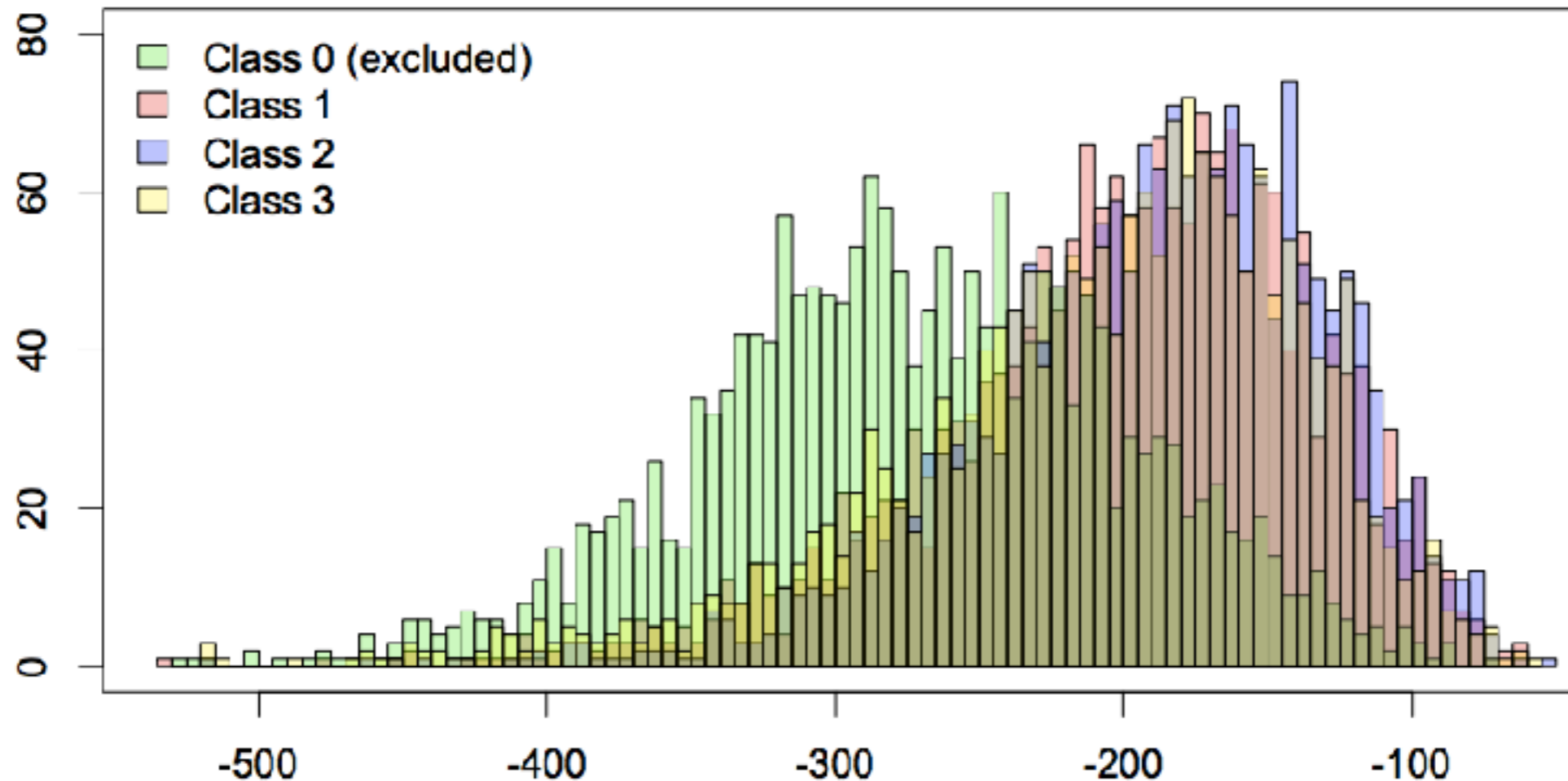
Conclusions

- **Generative can be used well for “discriminative problems”**
 - Especially in data-restricted scenarios
 - Especially with neural nets, which let us define great generative models
- **Open questions**
 - Inference is hard, but there are lots of exciting possibilities for learning to do inference
 - Is there a theoretical account for when a particular dataset is in the “generative” vs. “discriminative” regime and where the crossover point is?

Thank you!

Outlier detection

- Generative models also provide an estimate of $p(\mathbf{x})$
- The likelihood of the input is a good estimate of “what the model knows”. Examples that fall out of this are a good indication that the model should stop what it’s doing and get help.



Zero-shot learning

- Train on $n - 1$ classes
- Predict for all classes
- Learn **(label) concepts**, to be used as class embeddings \mathbf{v}_y from an **auxiliary task**
 - For example, from a large unannotated corpus, learn standard **word embeddings** and use them **as class embeddings**
- **Fix** the class embeddings **during training**
- When we see a **new class**, use the word embedding for the class

Zero-shot learning

Class	Precision	Recall	Accuracy
company	98.9	46.6	93.3
educational institution	99.2	49.5	92.8
athlete	96.5	90.1	94.6
means of transportation	96.5	74.3	94.2
building	99.9	37.7	92.1
natural place	98.9	88.2	95.4
village	99.9	68.1	93.8
animal	99.7	68.1	93.8
plant	99.2	76.9	94.3
film	99.4	73.3	94.5
written work	93.8	26.5	91.3
AVERAGE	98.3	63.6	93.6

Inference

Importance sampling

Assume we've got a conditional distribution $q(\mathbf{y} \mid \mathbf{x})$

- s.t.
- (i) $p(\mathbf{x}, \mathbf{y}) > 0 \implies q(\mathbf{y} \mid \mathbf{x}) > 0$
 - (ii) $\mathbf{y} \sim q(\mathbf{y} \mid \mathbf{x})$ is tractable and
 - (iii) $q(\mathbf{y} \mid \mathbf{x})$ is tractable

Let the importance weights $w(\mathbf{x}, \mathbf{y}) = \frac{p(\mathbf{x}, \mathbf{y})}{q(\mathbf{y} \mid \mathbf{x})}$

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} p(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} w(\mathbf{x}, \mathbf{y}) q(\mathbf{y} \mid \mathbf{x}) \\ &= \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} \mid \mathbf{x})} w(\mathbf{x}, \mathbf{y}) \end{aligned}$$

Inference

Importance sampling

$$\begin{aligned} p(\mathbf{x}) &= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} p(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{x})} w(\mathbf{x}, \mathbf{y}) q(\mathbf{y} | \mathbf{x}) \\ &= \mathbb{E}_{\mathbf{y} \sim q(\mathbf{y} | \mathbf{x})} w(\mathbf{x}, \mathbf{y}) \end{aligned}$$

Replace this expectation with its Monte Carlo estimate.

$$\mathbf{y}^{(i)} \sim q(\mathbf{y} | \mathbf{x}) \quad \text{for } i \in \{1, 2, \dots, N\}$$

$$\mathbb{E}_{q(\mathbf{y} | \mathbf{x})} w(\mathbf{x}, \mathbf{y}) \stackrel{\text{MC}}{\approx} \frac{1}{N} \sum_{i=1}^N w(\mathbf{x}, \mathbf{y}^{(i)})$$

Results: **Language modeling**

	Perplexity
5-gram IKN	169.3
LSTM LM	113.4
Generative (IS)	102.4