

Social Media & Text Analysis

lecture 7 - Twitter Paraphrases and
Latent Variable Models



CSE 5539-0010 Ohio State University
Instructor: Wei Xu
Website: socialmedia-class.org

Homework #2 is out

Due in three weeks

Autumn 2016



AU16 5539 > Files

Home

Assignments

Grades

People

Modules

Files

Collaborations

Announcements

Search for files



0 items selected

▼ AU16 CSE 5539 - Interim: Artif Int

Name ▲

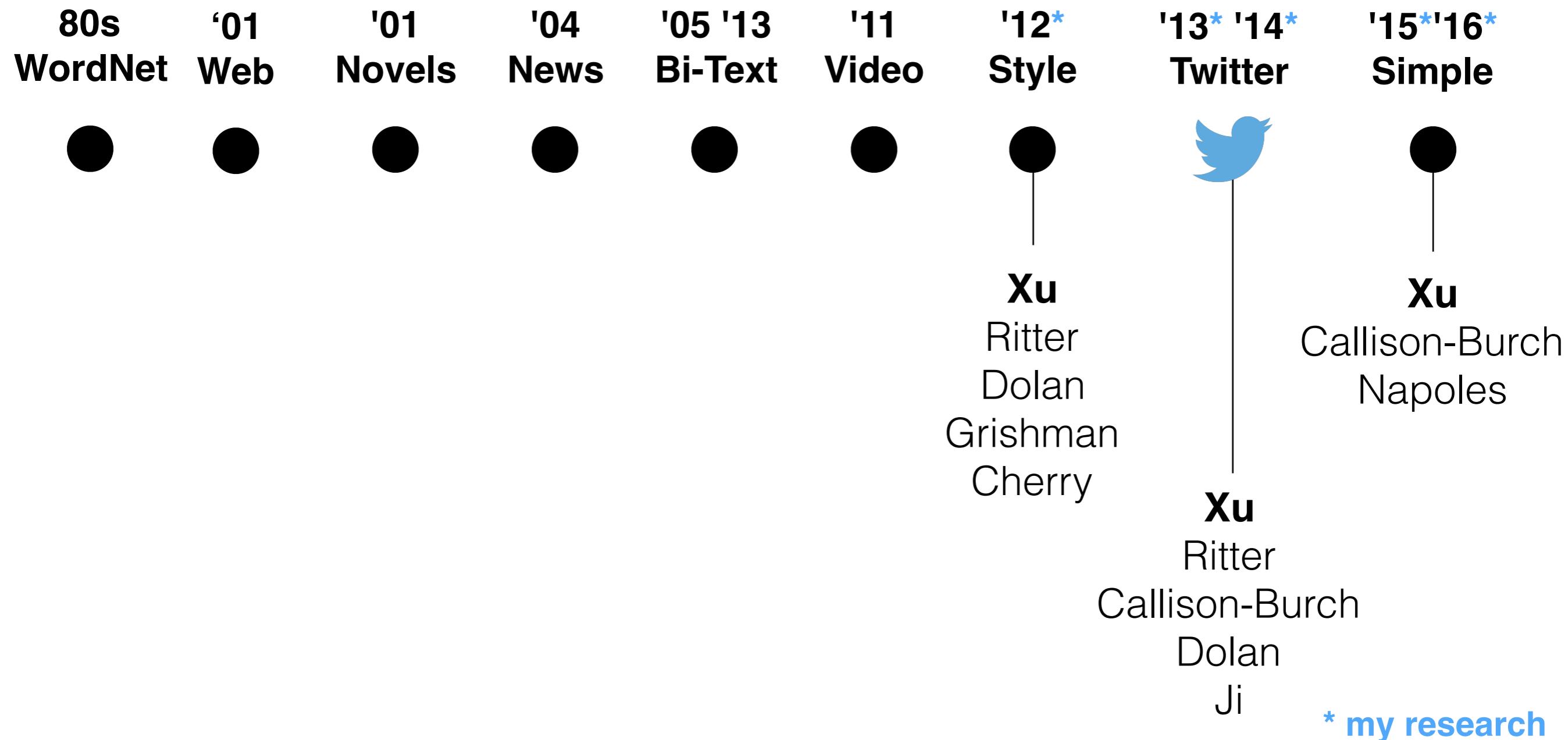


AU16_5539_0010_homework2.zip

Collaborators

Chris Callison-Burch	UPenn
Alan Ritter	UW / OSU
Bill Dolan	MSR
Yangfeng Ji	GaTech
Jeniya Tabassum	OSU
Wuwei Lan	OSU
Ralph Grishman	NYU
Raphael Hoffmann	UW / AI2 Incubator
Joel Tetreault	ETS / Yahoo!
Le Zhao	CMU / Google
Maria Pershina	NYU
Colin Cherry	NRC
Courtney Napolis	JHU
Lyle Ungar	UPenn
Daniel Preoțiuc-Pietro	UPenn
Ellie Pavlick	UPenn
Mingkun Gao	UPenn / UIUC
Quanze Chen	UPenn / UW
Martin Chodorow	CUNY

Paraphrase Research



Paraphrase Research



News



only a few hundreds news agencies
only big events
only well-edited text
(the MSR Paraphrase Corpus)

Twitter as a new resource



Rep. Stacey Newman @staceynewman · 5h

So sad to hear today of former WH Press Sec **James Brady's passing**.
@bradybuzz & family will carry on his legacy of #gunsense.



Jim Sciutto @jimsciutto · 4h

Breaking: Fmr. WH Press Sec. **James Brady** has died at 73, crusader for gun control after wounded in '81 Reagan assassination attempt



NBC News @NBCNews · 2h

James Brady, President Reagan's press secretary shot in 1981 assassination attempt, dead at 73 nbcnews.to/WX1Btq pic.twitter.com/1ZtuEakRd9



average sentence length: news ≈18.6 words Twitter ≈11.9 words

Twitter as a powerful resource

thousands of users
talk about both big/micro events daily



a very broad range of paraphrases:
synonyms, misspellings, slang, acronyms and colloquialisms

Paraphrase Model

obtain sentential paraphrases automatically

Mancini has been sacked by Manchester City

Yes!

Mancini gets the boot from Man City

WORLD OF JENKS IS ON AT 11

No!

World of Jenks is my favorite show on tv

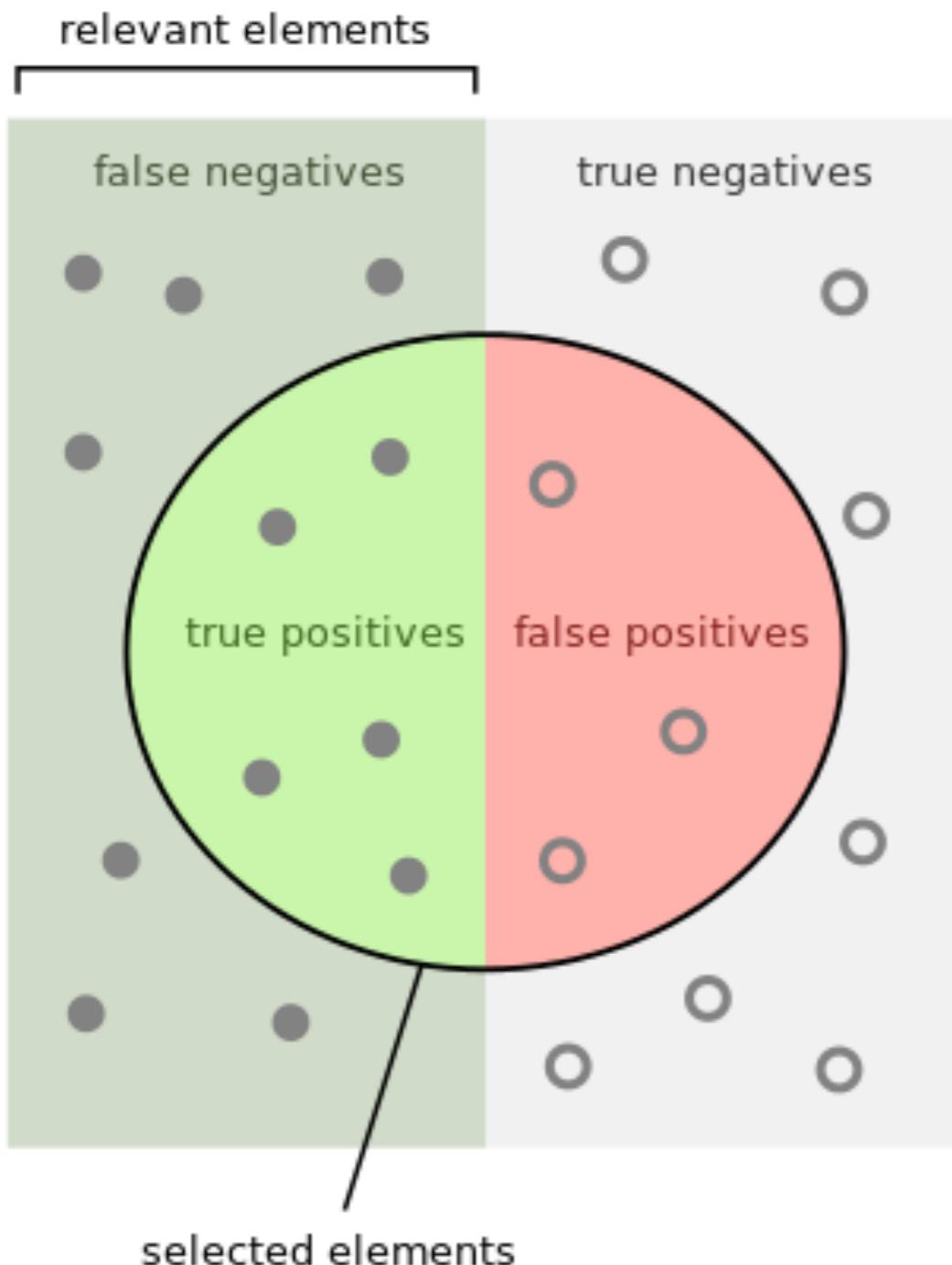
(On news data - MSR Paraphrase Corpus)



Paraphrase Identification

Algorithm	Reference	Description	Supervisi	Accurac	F1
Vector	Mihalcea et al. (2006)	cosine similarity with tf-idf weighting	unsupervised	65.40%	75.30%
ESA	Hassan (2011)	explicit semantic space	unsupervised	67.00%	79.30%
KM	Kozareva and Montoyo (2006)	combination of lexical and semantic features	supervised	76.60%	79.60%
LSA	Hassan (2011)	latent semantic space	unsupervised	68.80%	79.90%
RMLMG	Rus et al. (2008)	graph subsumption	unsupervised	70.60%	80.50%
MCS	Mihalcea et al. (2006)	combination of several word similarity measures	unsupervised	70.30%	81.30%
WTMF	Guo and Diab (2012)	latent space model for short text	unsupervised	71.51%	---
STS	Islam and Inkpen (2007)	combination of semantic and string similarity	unsupervised	72.60%	81.30%
SSA	Hassan (2011)	salient semantic space	unsupervised	72.50%	81.40%
QKC	Qiu et al. (2006)	sentence dissimilarity classification	supervised	72.00%	81.60%
ParaDetect	Zia and Wasif (2012)	PI using semantic heuristic features	supervised	74.70%	81.80%
SDS	Blacoe and Lapata (2012)	simple distributional semantic space	supervised	73.00%	82.30%
matrixJcn	Fernando and Stevenson (2008)	JCN WordNet similarity with matrix	unsupervised	74.10%	82.40%
FHS	Finch et al. (2005)	combination of MT evaluation measures as features	supervised	75.00%	82.70%
PE	Das and Smith (2009)	product of experts	supervised	76.10%	82.70%
WDDP	Wan et al. (2006)	dependency-based features	supervised	75.60%	83.00%
SHPNM	Socher et al. (2011)	recursive autoencoder with dynamic pooling	supervised	76.80%	83.60%
MTMETRICS	Madnani et al. (2012)	combination of eight machine translation metrics	supervised	77.40%	84.10%
Bi-CNN-MI	Yin and Schutze (2015)	convolutional neural network w/ multi-granular interaction	supervised	78.40%	84.60%
MP-CNN	He et al. (2015)	convolutional neural network w/ multiple perspectives	supervised	78.60%	84.73%
Recon	Cheng and Kartsaklis (2015)	Recursive NNs w/ syntax-aware multi-sense word embeddings	supervised	78.60%	85.30%
LEXDISCRIM	Ji and Eisenstein (2013)	combination of latent space and lexical features	supervised	80.41%	85.96%

Classification Evaluation: Precision, Recall, F-measure



How many selected items are relevant?

Precision =



How many relevant items are selected?

Recall =



$$F_1 = \frac{2 \cdot precision \cdot recall}{precision + recall}$$

Classification Evaluation: Accuracy

- In classification problems, the primary source for accuracy estimation is the **confusion matrix**

		True Class	
		Positive	Negative
Predicted Class	Positive	True Positive Count (TP)	False Positive Count (FP)
	Negative	False Negative Count (FN)	True Negative Count (TN)

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

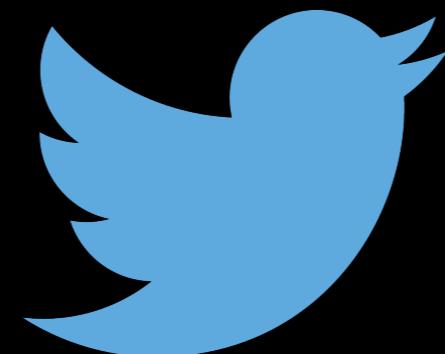
$$\text{True Positive Rate} = \frac{TP}{TP + FN}$$

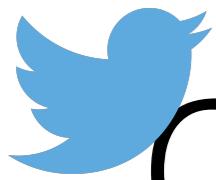
$$\text{True Negative Rate} = \frac{TN}{TN + FP}$$

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

Multi-instance Learning Paraphrase Model





Challenges of Twitter Data

Mancini has been sacked by Manchester City

Mancini gets the boot from Man City

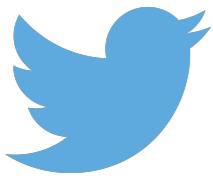
Yes!



very short, lexically divergent

Techniques

- Multiple Instance Learning
- **Probabilistic Graphical Models**
- Markov Networks with Latent Variables



Assumption

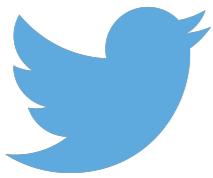
*Mancini has been sacked by Manchester **City***

*Mancini gets the boot from Man **City***

Yes!

At-least-one Paraphrase Anchor

two sentences about the same topic are paraphrases
if and only if
they contain at least one word pair that is
a paraphrase **anchor** (in the context)



Assumption

*The **new** Ciroc flavor has arrived*

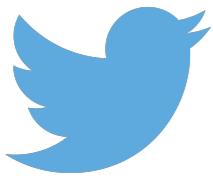
*Ciroc got a **new flavor** coming out*

Yes!

At-least-one Paraphrase Anchor

two sentences about the same topic are paraphrases
if and only if

they contain at least one word pair that is
a paraphrase **anchor** (in the context)



Assumption

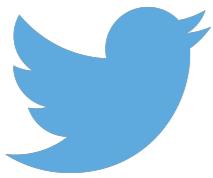
*Manti bout to be the **next** Junior Seau*

*Teo is the little **new** Junior Seau*

Yes!

At-least-one Paraphrase Anchor

two sentences about the same topic are paraphrases
if and only if
they contain at least one word pair that is
a paraphrase **anchor** (in the context)



Assumption

WORLD OF JENKS **IS ON AT 11**

World of Jenks **is my favorite show on tv**

No!

At-least-one Paraphrase Anchor

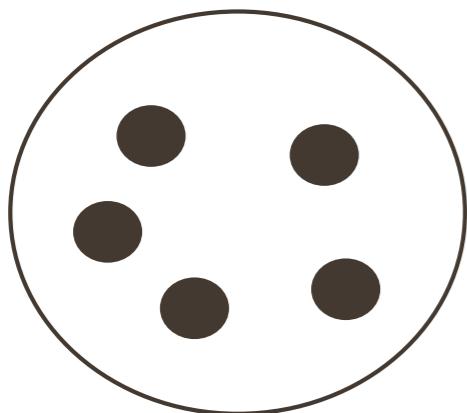
two sentences about the same topic are paraphrases
if and only if

they contain at least one word pair that is
a paraphrase **anchor** (in the context)

Multi-instance Learning

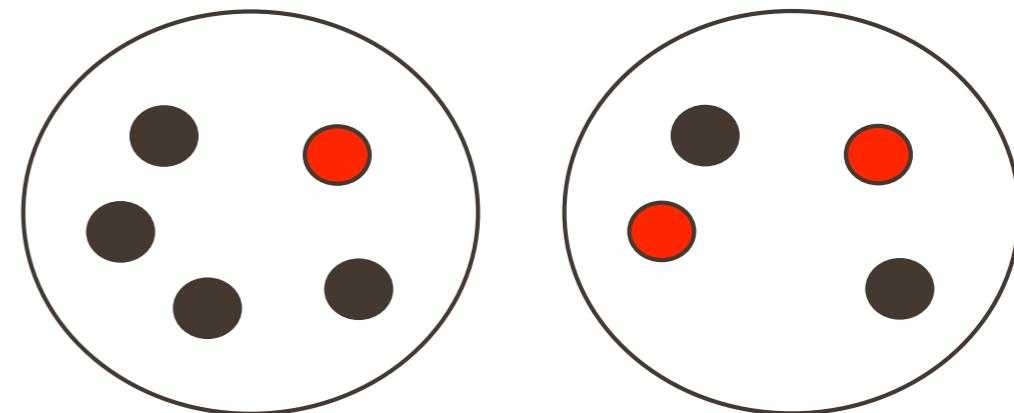
Instead of labels on each individual instance, the learner only observes labels on bags of instances.

Negative Bags



A bag is labeled negative, if **all** the examples in it are negative

Positive Bags

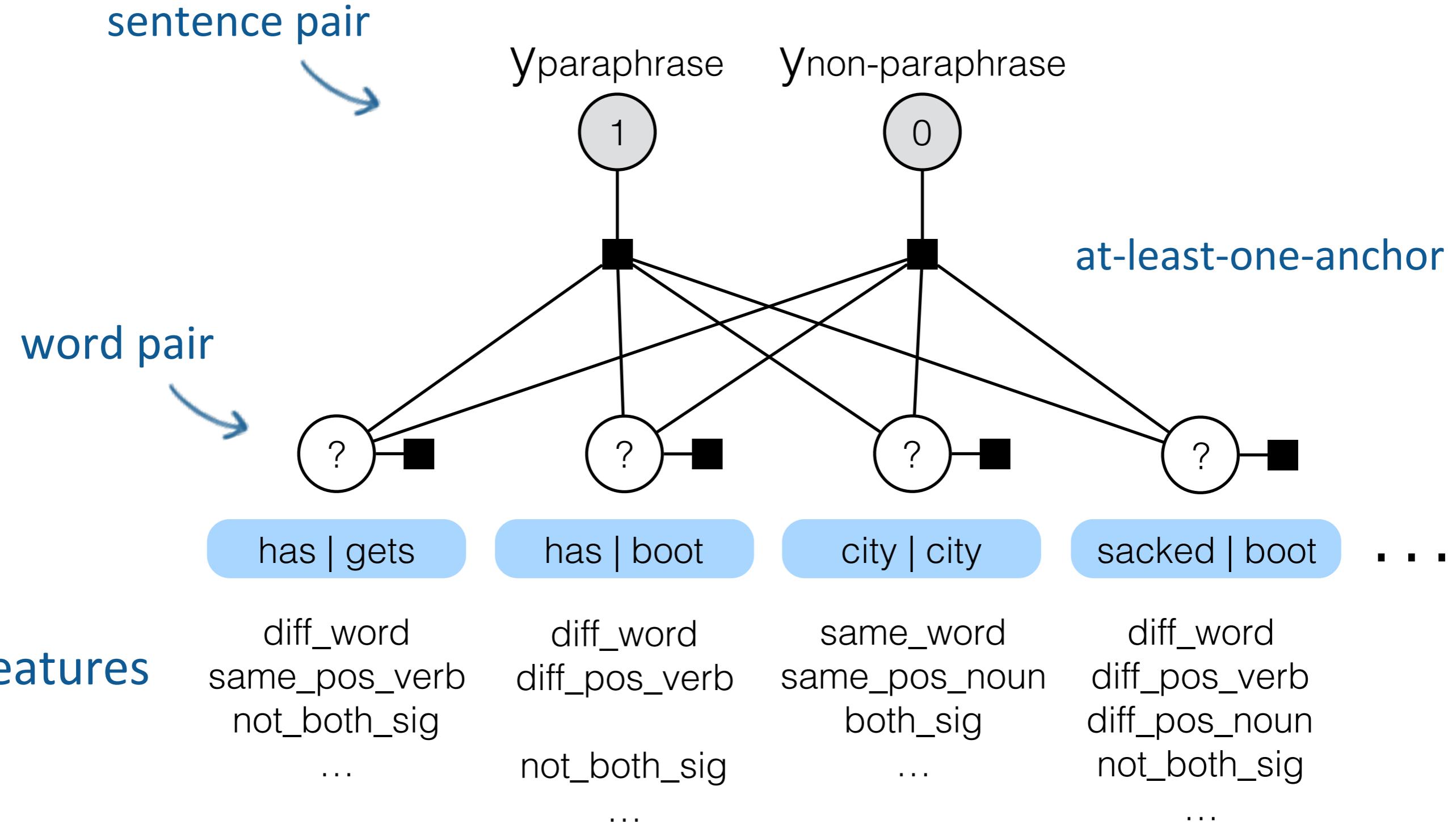


A bag is labeled positive, if there is **at least one** positive example

Multi-instance Learning Paraphrase Model

Mancini has been sacked by Manchester City

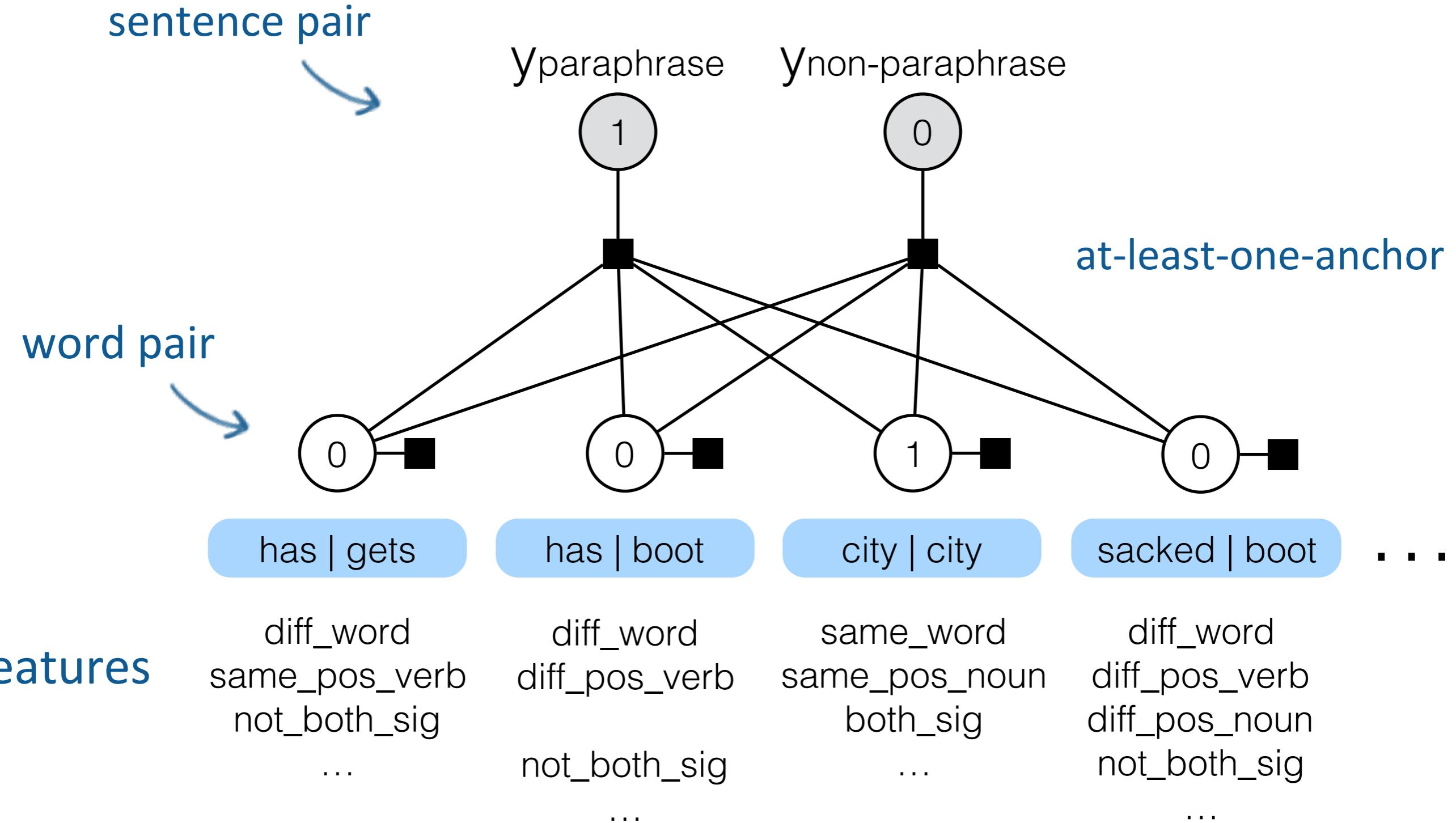
Mancini gets the boot from Man City



Multi-instance Learning Paraphrase Model

*Mancini has been sacked by Manchester **City***

*Mancini gets the boot from Man **City***



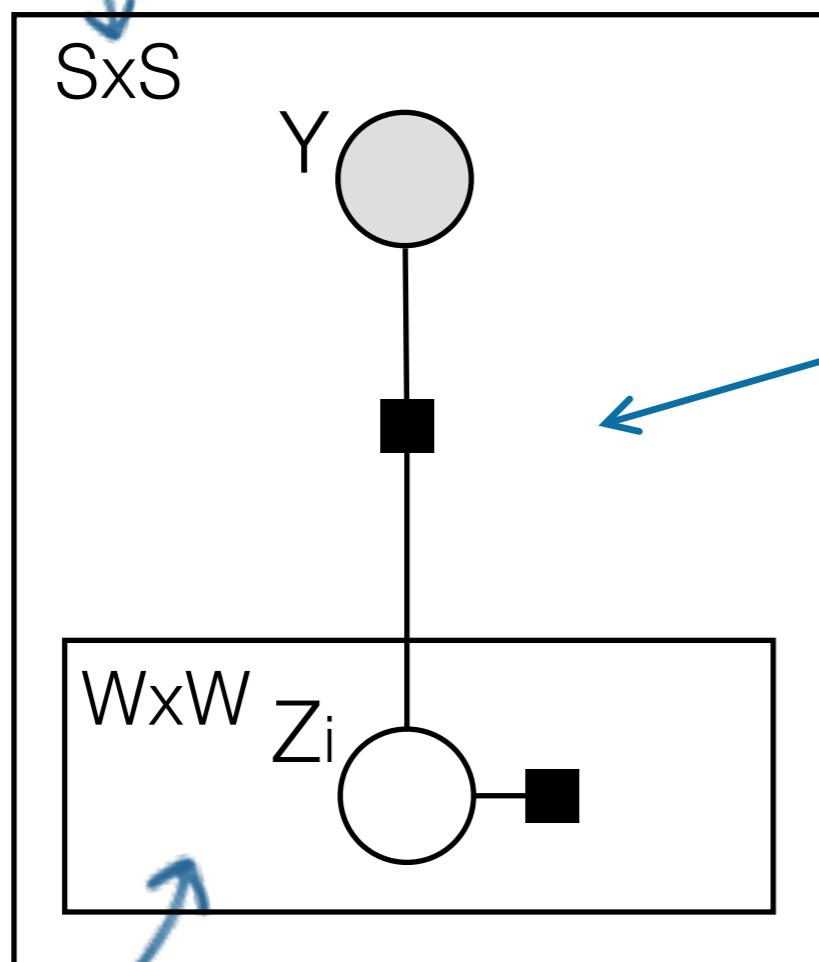
Multi-instance Learning

Model the assumption:

sentence-level paraphrase

is anchored by at-least-one word pair

sentence pair



word pair

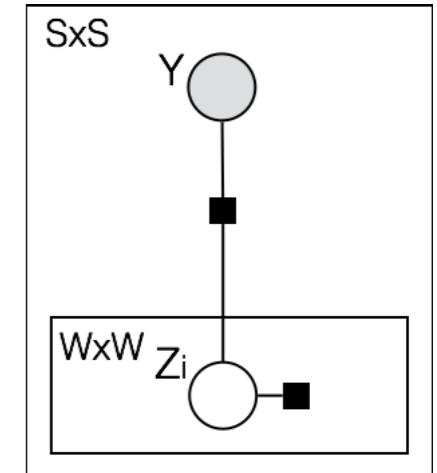
$$\sigma(\mathbf{z}, \mathbf{y}) = \begin{cases} 1 & \text{if } \mathbf{y} = \text{true} \wedge \exists i : z_i = 1 \\ 1 & \text{if } \mathbf{y} = \text{false} \wedge \forall i : z_i = 0 \\ 0 & \text{otherwise} \end{cases}$$

$$P(\mathbf{z}_i, y_i | \mathbf{w}_i; \theta) = \prod_{j=1}^m \exp(\theta \cdot f(z_j, w_j)) \times \sigma(\mathbf{z}_i, y_i)$$

Learning Algorithm

Objective:

learn the parameters that maximize conditional likelihood over the training corpus



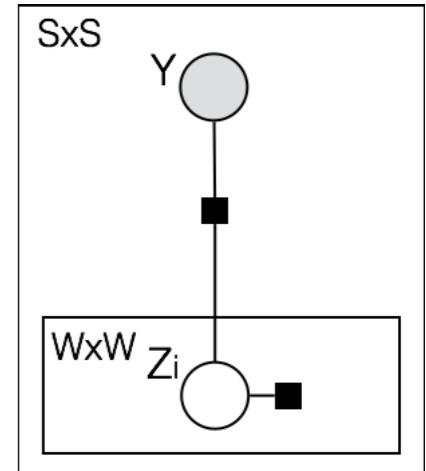
$$\theta^* = \arg \max_{\theta} P(\mathbf{y}|\mathbf{w}; \theta) = \arg \max_{\theta} \prod_i \sum_{\mathbf{z}_i} P(\mathbf{z}_i, y_i | \mathbf{w}_i; \theta)$$

***i*th training sentence pair**

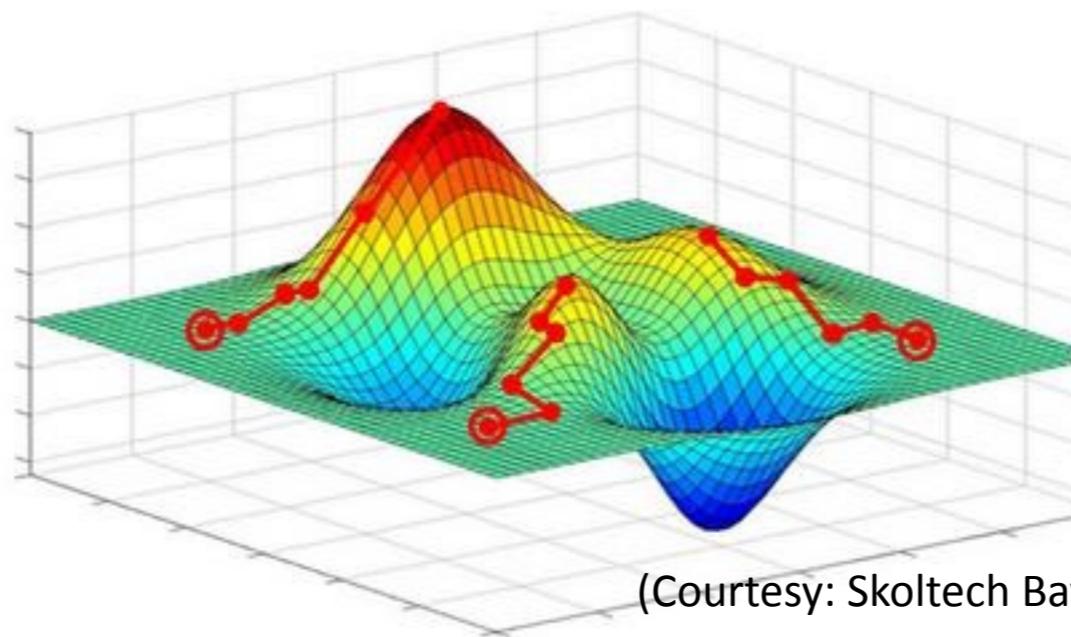
**all possible values
of the latent variables**

Learning Algorithm

Finding maximum likelihood estimate:
stochastic gradient ascent



$$\frac{\partial \log P(\mathbf{y}|\mathbf{w}; \theta)}{\partial \theta} = \mathbf{E}_{P(\mathbf{z}|\mathbf{w}, \mathbf{y}; \theta)} \left(\sum_i f(\mathbf{z}_i, \mathbf{w}_i) \right) - \mathbf{E}_{P(\mathbf{z}, \mathbf{y}|\mathbf{w}; \theta)} \left(\sum_i f(\mathbf{z}_i, \mathbf{w}_i) \right)$$

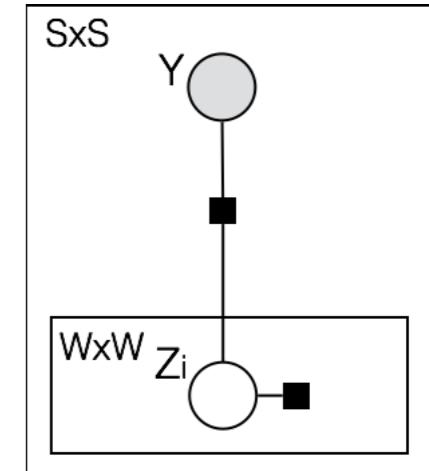


Learning Algorithm

Perceptron-style Update:

Viterbi approximation + online learning

$O(\# \text{ word pairs})$



$$\begin{aligned}\frac{\partial \log P(\mathbf{y}|\mathbf{w}; \theta)}{\partial \theta} &= \mathbf{E}_{P(\mathbf{z}|\mathbf{w}, \mathbf{y}; \theta)} \left(\sum_i f(\mathbf{z}_i, \mathbf{w}_i) \right) - \mathbf{E}_{P(\mathbf{z}, \mathbf{y}|\mathbf{w}; \theta)} \left(\sum_i f(\mathbf{z}_i, \mathbf{w}_i) \right) \\ &\approx \underbrace{\sum_i f(\mathbf{z}_i^*, \mathbf{w}_i)}_{\text{reward correct}} - \underbrace{\sum_i f(\mathbf{z}'_i, \mathbf{w}_i)}_{\text{penalize wrong}}\end{aligned}$$

reward correct
(conditioned on labels)

$$\mathbf{z}^* = \arg \max_{\mathbf{z}} P(\mathbf{z}|\mathbf{w}, \mathbf{y}; \theta)$$

penalize wrong
(ignoring labels)

$$\mathbf{y}', \mathbf{z}' = \arg \max_{\mathbf{y}, \mathbf{z}} P(\mathbf{z}, \mathbf{y}|\mathbf{w}; \theta)$$

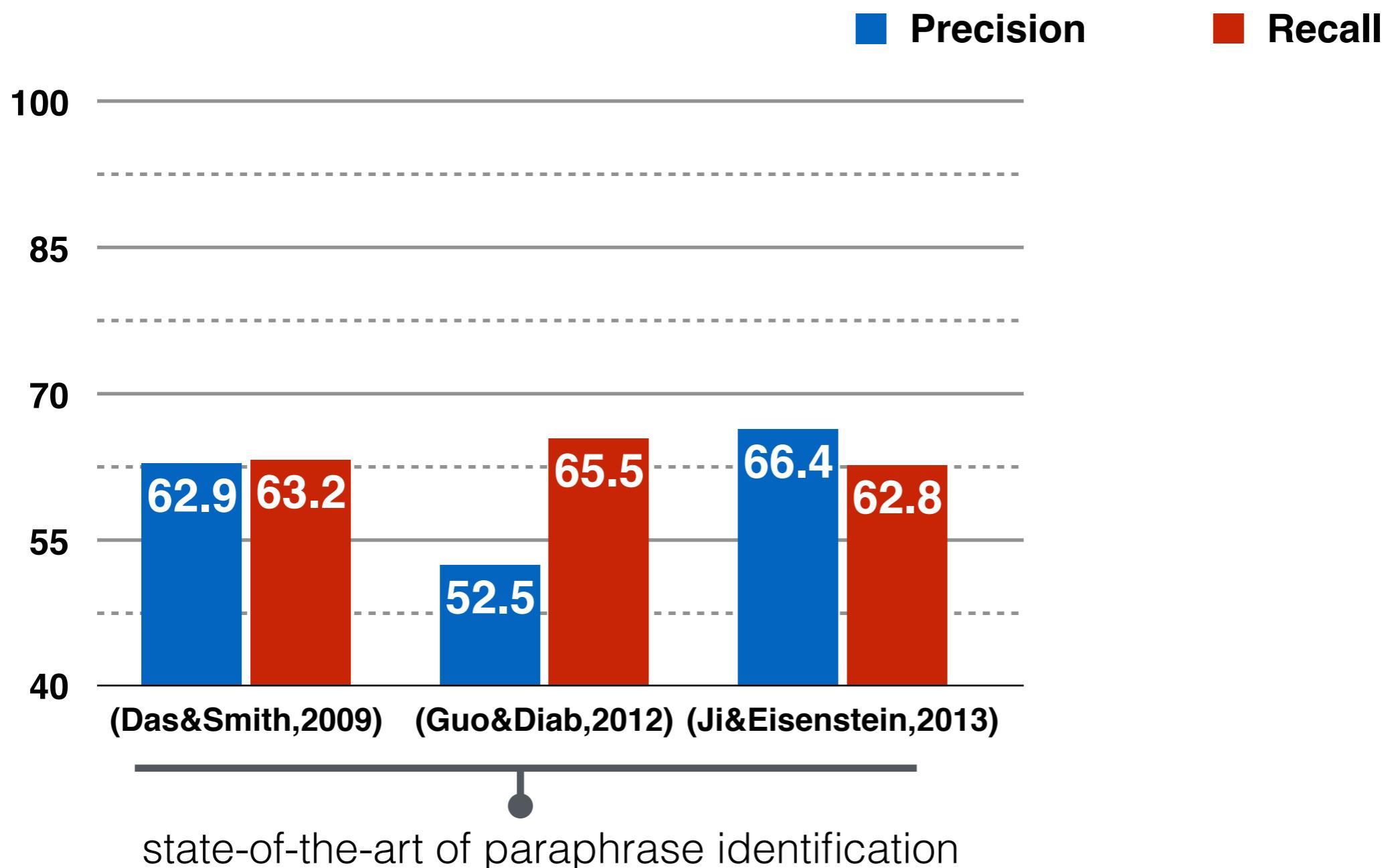
Twitter Paraphrase Dataset

18,762 labeled sentence pairs
1/3 paraphrase, 2/3 non-paraphrase

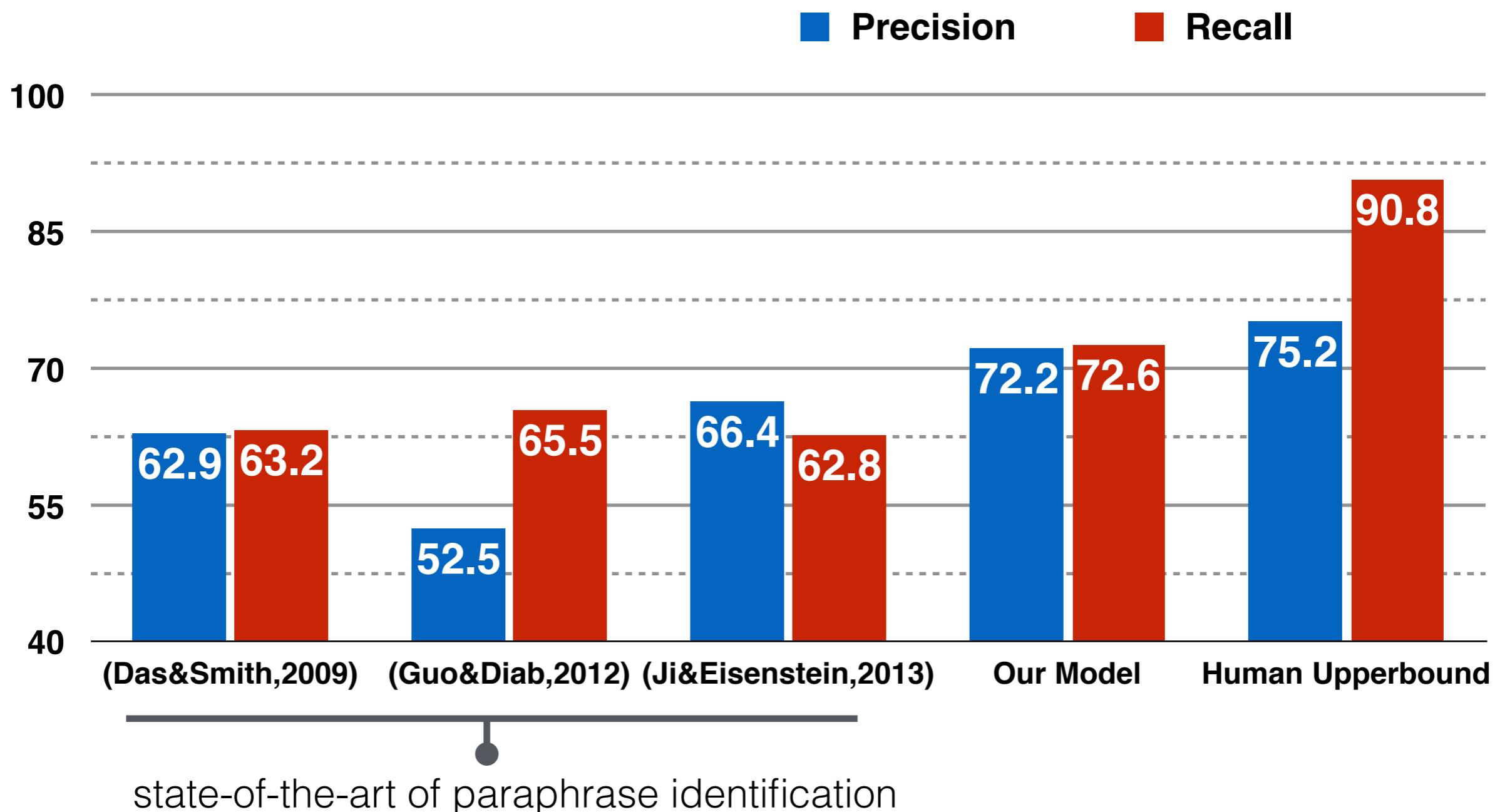
Techniques

- **Crowdsourcing** (Human-Computer Interaction)
- SumBasic algorithm for sentence filtering
- Multi-armed bandits algorithm for topic selection

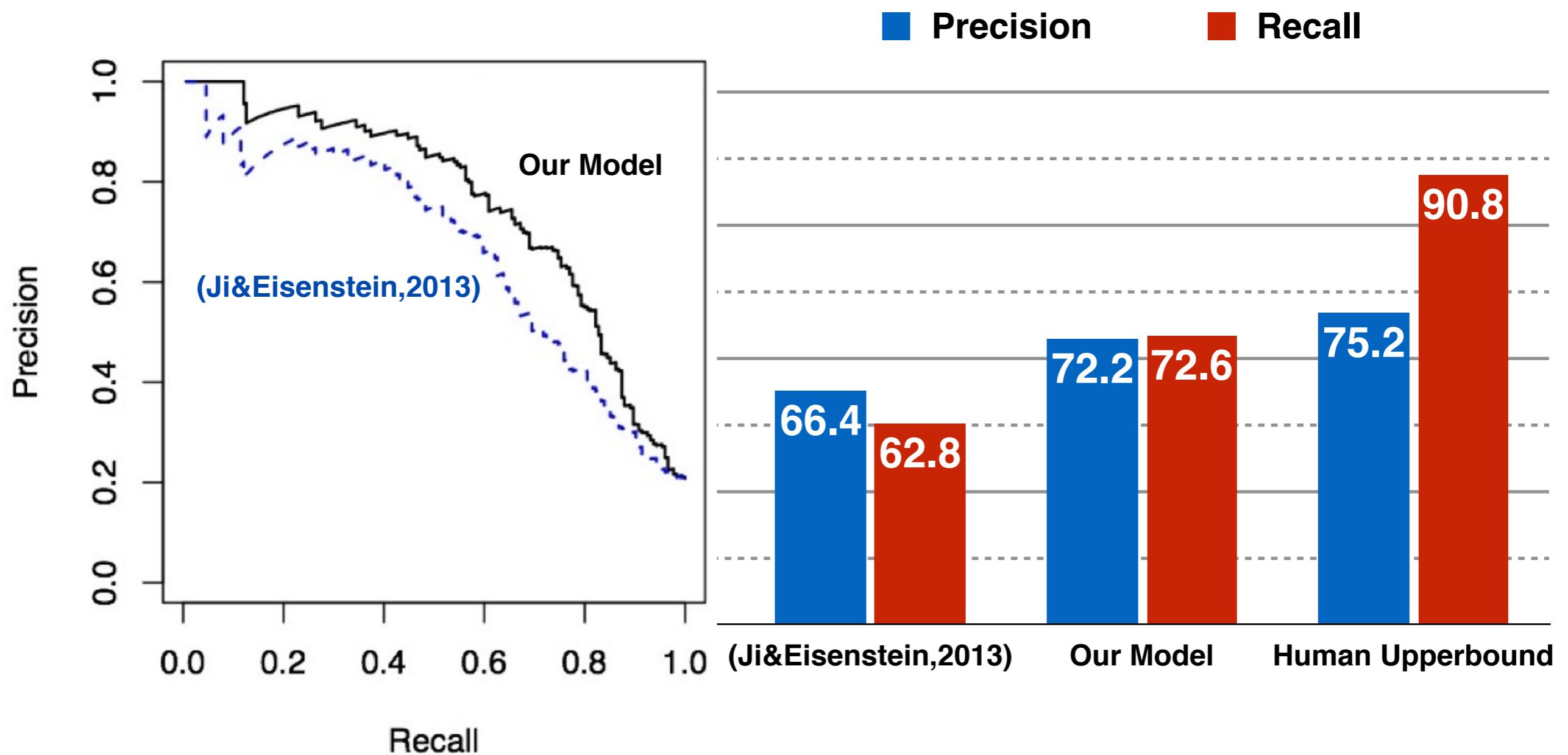
Performance



Performance



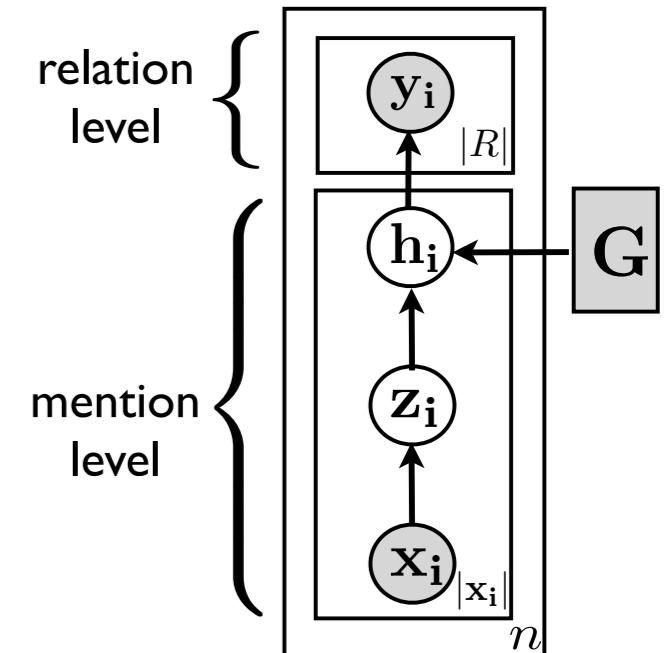
Performance



Multi-instance Learning

Other Application: Learning Knowledge Base

... the forced resignation of the CEO of Boeing, Harry Stonecipher, for ...



Multi-instance Learning

Other Application: Resolving Time Expressions

98.9 THE DRIVE @989THEDRIVE

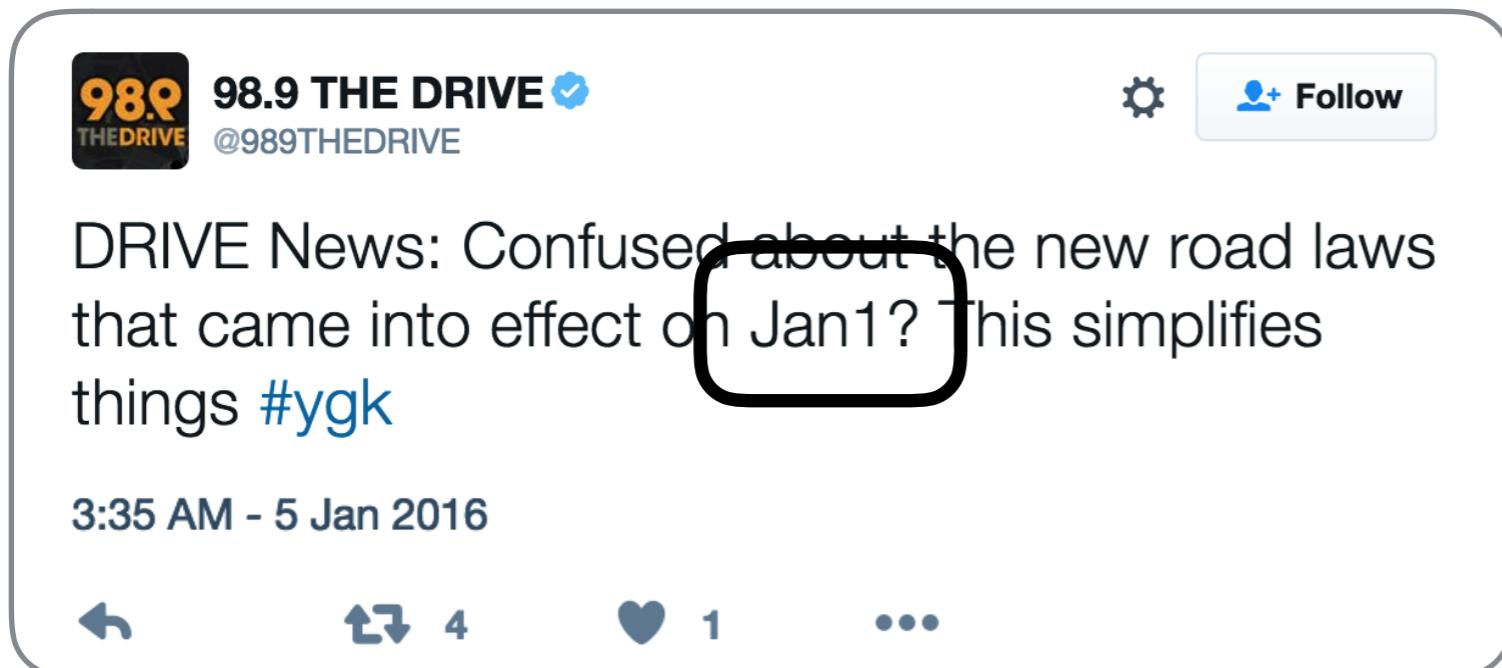
DRIVE News: Confused about the new road laws that came into effect on Jan1? This simplifies things [#ygk](#)

3:35 AM - 5 Jan 2016

4 1 ...

Multi-instance Learning

Other Application: Resolving Time Expressions



98.9 THE DRIVE 
@989THEDRIVE

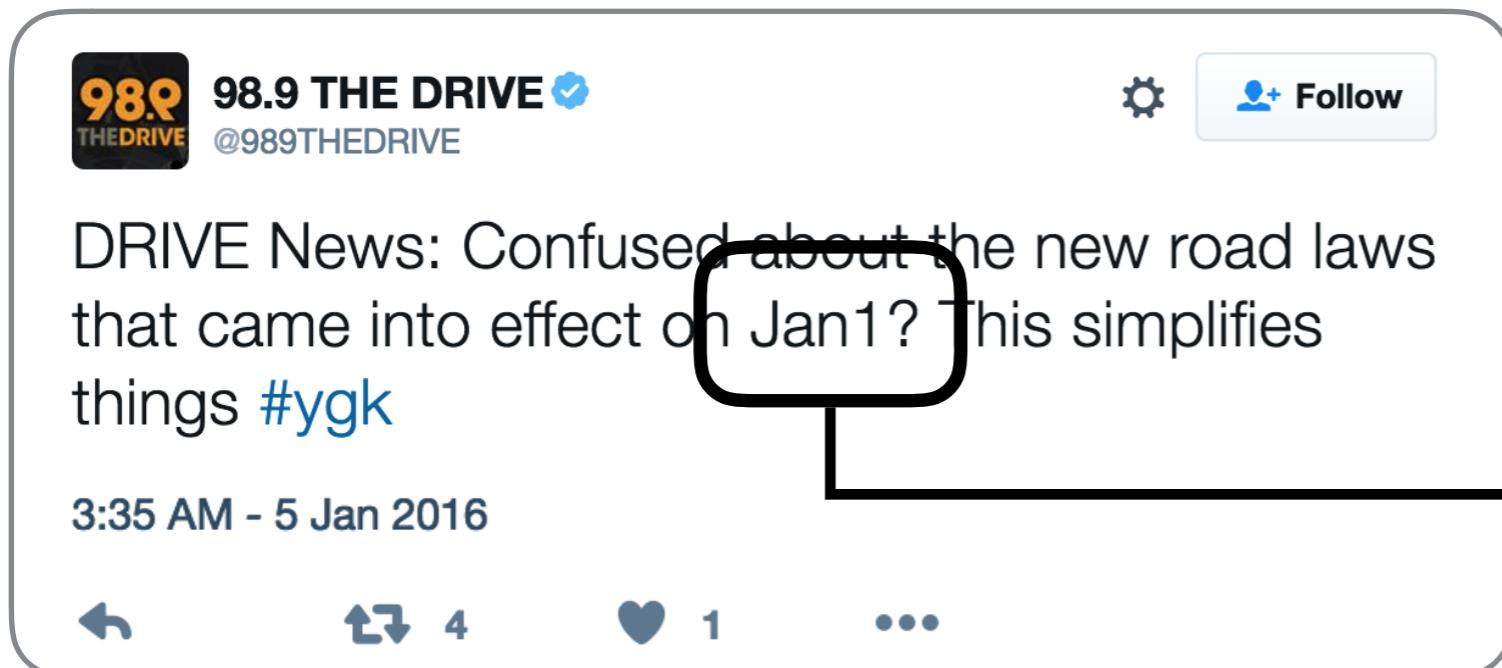
DRIVE News: Confused about the new road laws
that came into effect on Jan1? This simplifies
things #ygk

3:35 AM - 5 Jan 2016

4 1 ...

Multi-instance Learning

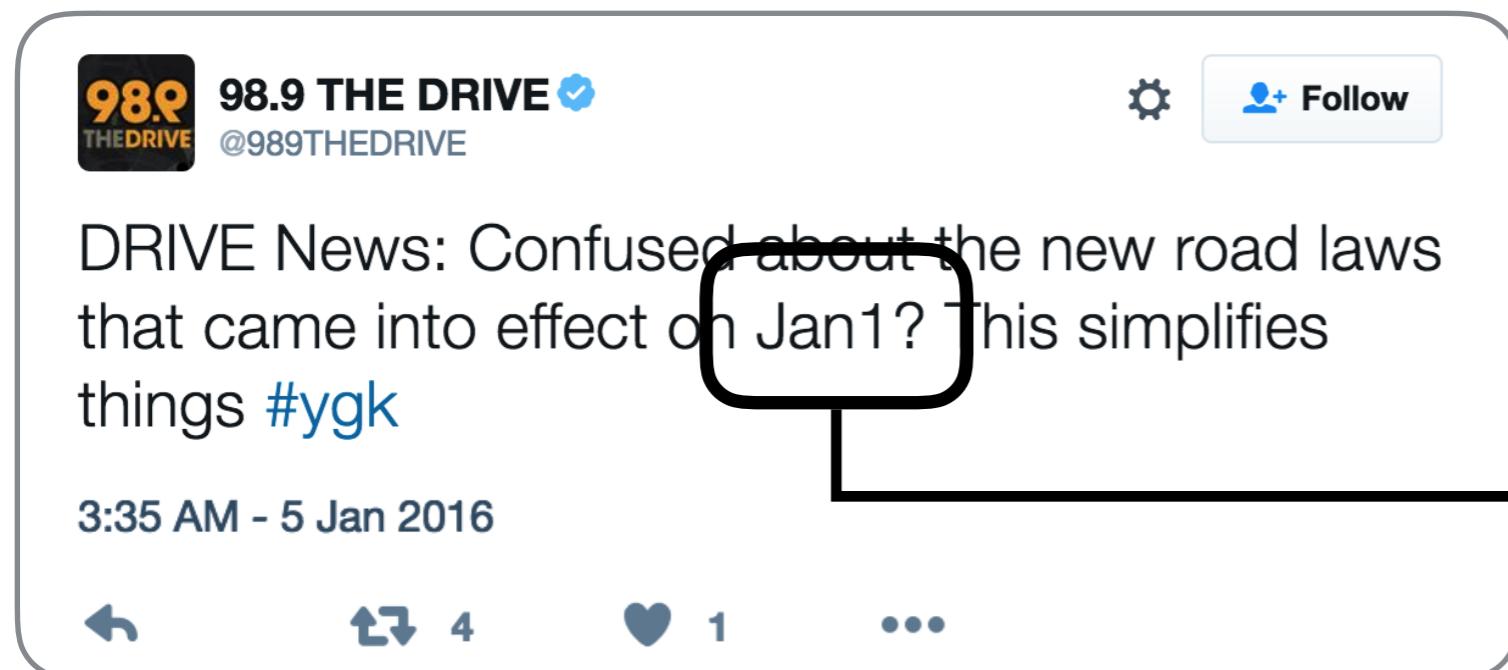
Other Application: Resolving Time Expressions



TweeTime
→ 1 Jan 2016

Multi-instance Learning

Other Application: Resolving Time Expressions



TweeTime

→ 1 Jan 2016

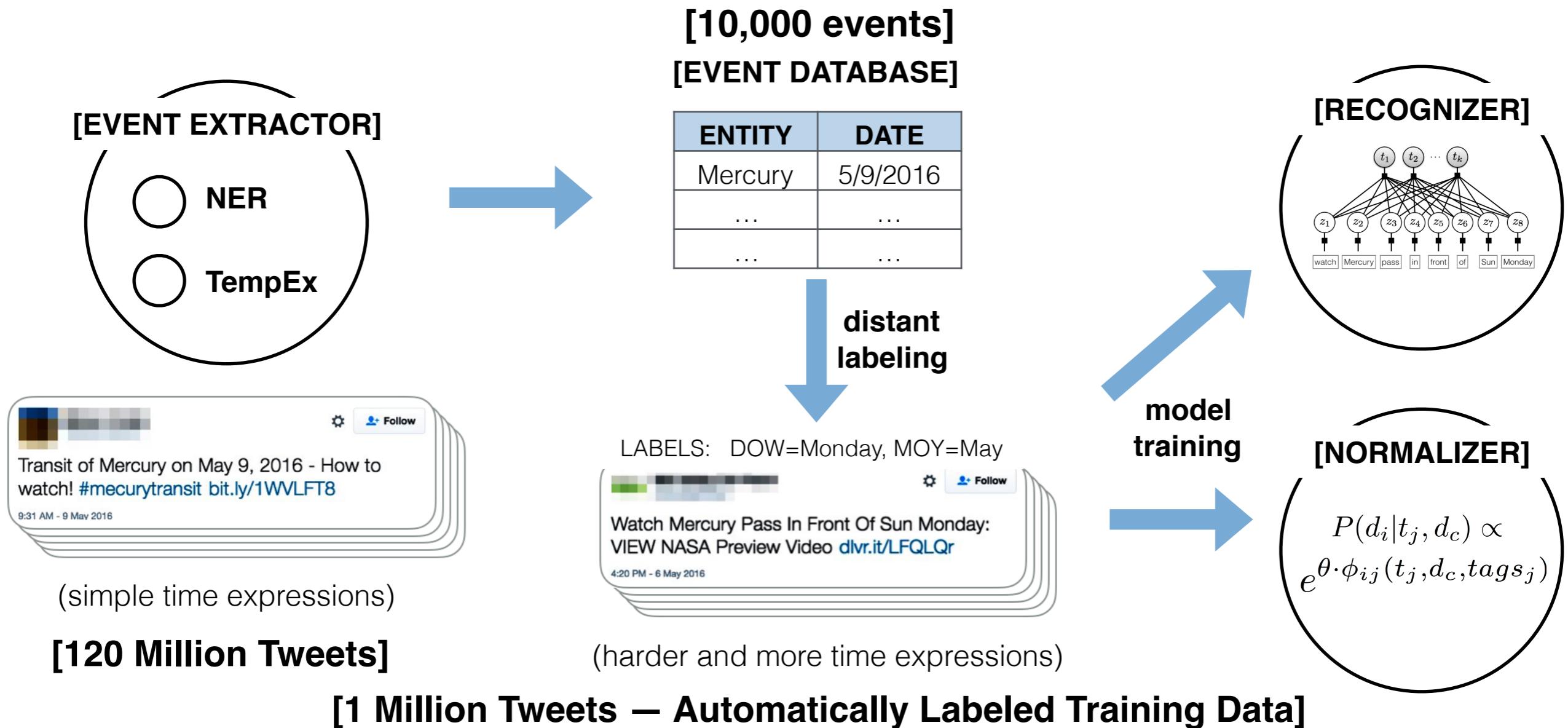
All other state-of-the-art time resolvers

{ TempEX
HeidelTime
SUTime
UWTime }

failed.

Multi-instance Learning

Other Application: Resolving Time Expressions



More about **TweeTime**



- **Jeniya Tabassum (OSU)**
- **Tuesday**, Oct 25, 12:45 pm, McPherson 2019
- A Minimally Supervised Method for Recognizing and Normalizing Time Expressions in Twitter

Crowdsourcing Training Data

Early Attempts on Twitter Paraphrase

- 1242 tweet pairs, tracking celebrity & hashtags
(Zanzotto, Pennacchiotti, Tsoutsouliklis, 2011)
- named entity + date
(Xu, Ritter, Grishman, 2013)
- bilingual posts, only phrases
(Ling, Dyer, Black, Trancoso, 2013)

Early Attempt:

Named Entity + Date

 **Tyler Anderson**
@tylerjanderson



From **January 16**, **Instagram** can sell your photos without permission
geek.com/articles/geek-...



 **Jeff Clutter**
@Pibbbs



Instagram can sell your photos without consent starting **January 16th**.



Early Attempt:

Self-translation

Valloire Galibier (@otVALLOIRE)

La neige est annoncée pour demain #Valloire
! | Snow is announced for tmrw !
bit.ly/1t85YQ5 | #neige #snow
#winteriscoming #valloire

translate

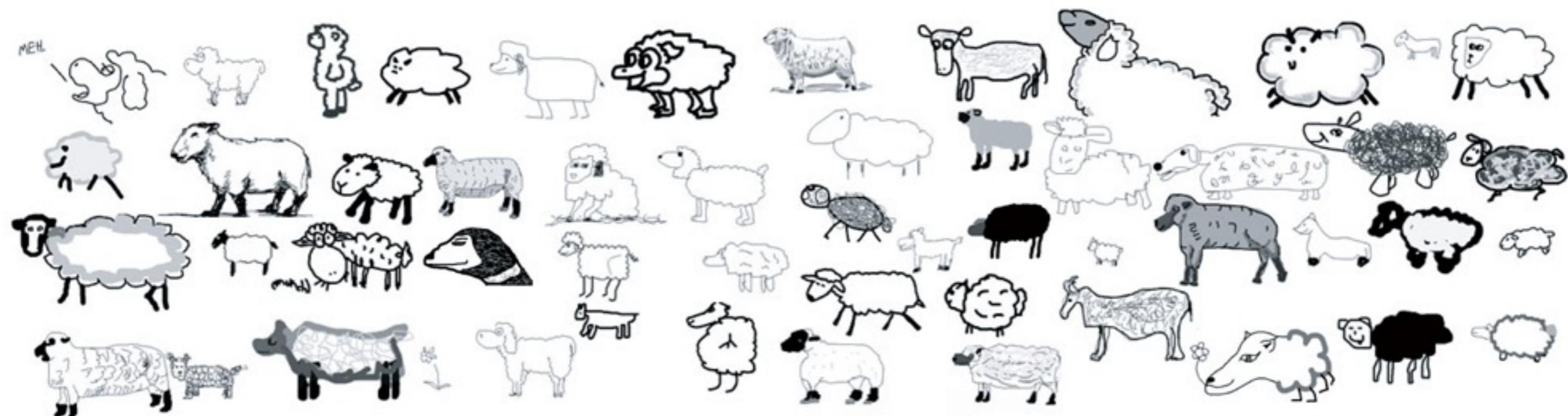
paraphrase

The snow is announced for tomorrow

Google Translate

Crowdsourcing

art



(Courtesy: The Sheep Market by Aaron Koblin)

Crowdsourcing

paraphrase

Here Is The Question To You:

Original Sentence: ***Borussia Dortmund advanced to the final***

Select ALL sentences that have similar meaning from below:

- Borussia Dortmund has clinched their Champions League final spot
- Real Madrid efforts are not enough as Cinderella Borussia Dortmund advances to the Champions League Final
- But it's Borussia Dortmund whose heading to Wembley Park
- Congratulations Borussia Dortmund's going to Wembley



Simple Quality Control

- Cohen's Kappa — used to remove bad workers

simple agreement between two raters



$$\kappa = \frac{p_0 - p_e}{1 - p_e}$$

Annotations: A blue arrow points from the term p_0 to the numerator, and another blue arrow points from the term p_e to the denominator.

likelihood that agreement
is attributed to chance

		B	
		Yes	No
A	Yes	45	15
	No	25	15

$$\kappa = \frac{0.60 - 0.54}{1 - 0.54} = 0.1304$$

		B	
		Yes	No
A	Yes	25	35
	No	5	35

$$\kappa = \frac{0.60 - 0.46}{1 - 0.46} = 0.2593$$

A Problem

only **8%** sentence pairs
about the same Twitter's trending topic
have similar meaning

hurts both quantity and quality

non-experts lower their bars

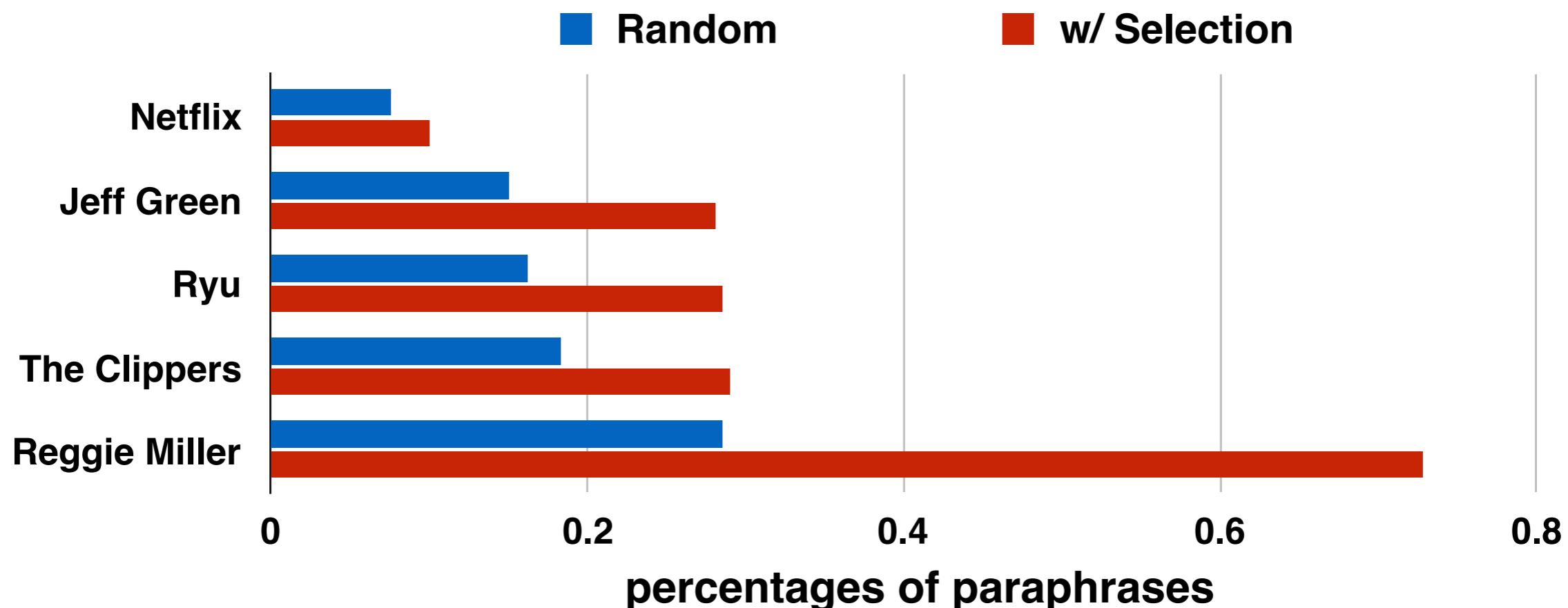


Sentence Selection

SumBasic Algorithm

8% → 16%

$$Salience(s) = \sum_{w_i \in s} \frac{P(w_i)}{|w_i| w_i \in s|}$$



Wei Xu, Alan Ritter, Ralph Grishman. "A Preliminary Study of Tweet Summarization using Information Extraction" in LASM (2013)

Wei Xu, Alan Ritter, Chris Callison-Burch, Bill Dolan, Yangfeng Ji. "Extracting Lexically Divergent Paraphrases from Twitter" In TACL (2014)

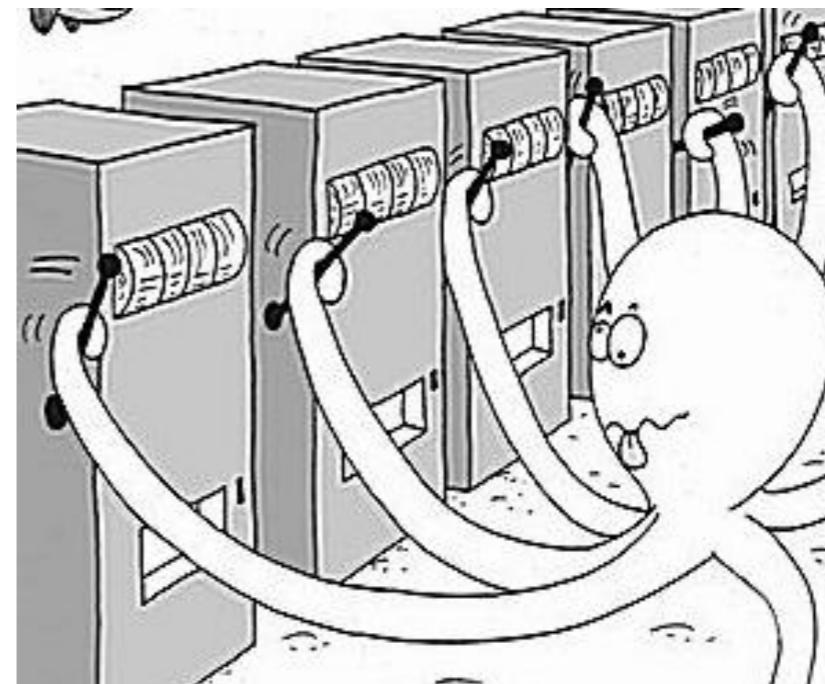
Topic Selection

Multi-Armed Bandits

16% → 34%

$$\max \sum_{i \in \{n | r_n(t_1) > 0\}} \hat{\mu}_i(t_0) r_i(t_1)$$

$$\text{s.t. } \sum_{i \in \{m | r_m(t_0) > 0\}} r_i(t_0) \leq (1 - \epsilon)B, \forall i : 0 \leq r_i(t_1) \leq l - r_i(t_0)$$



Innovations

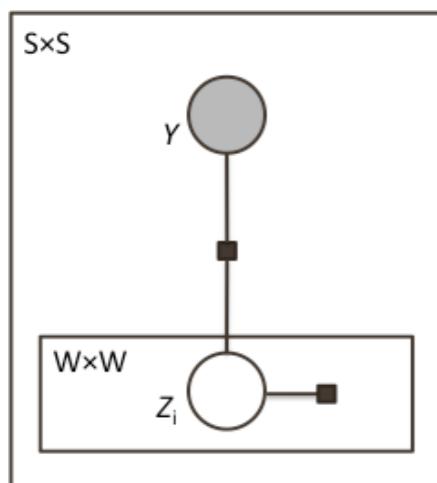
Web-scale Paraphrase from Twitter

Mancini has been sacked by Manchester City

Mancini gets the boot from Man City

Yes!

Multi-instance Learning Paraphrase Model



- Twitter's big data stream
- joint sentence-word alignment
- no word-level annotation needed
- extensible latent variable model

Impact & Future Work

[SemEval 2015 Shared-Task 1]

Paraphrase and Semantic Similarity in Twitter



Wei Xu

(University of Pennsylvania)



Chris Callison-Burch

(University of Pennsylvania)



Bill Dolan

(Microsoft Research)

Participation

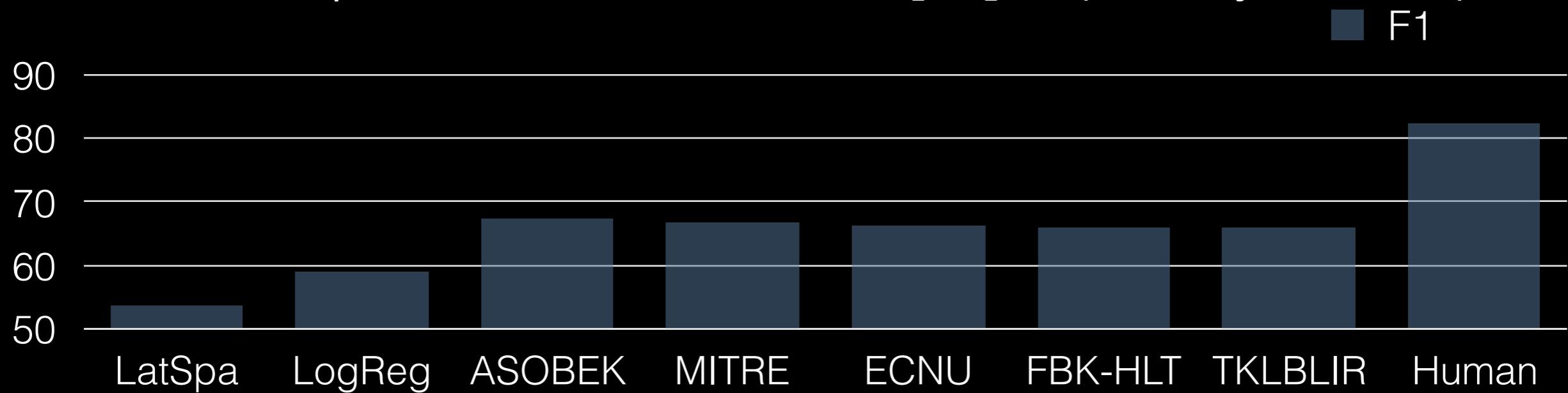


Popular Techniques

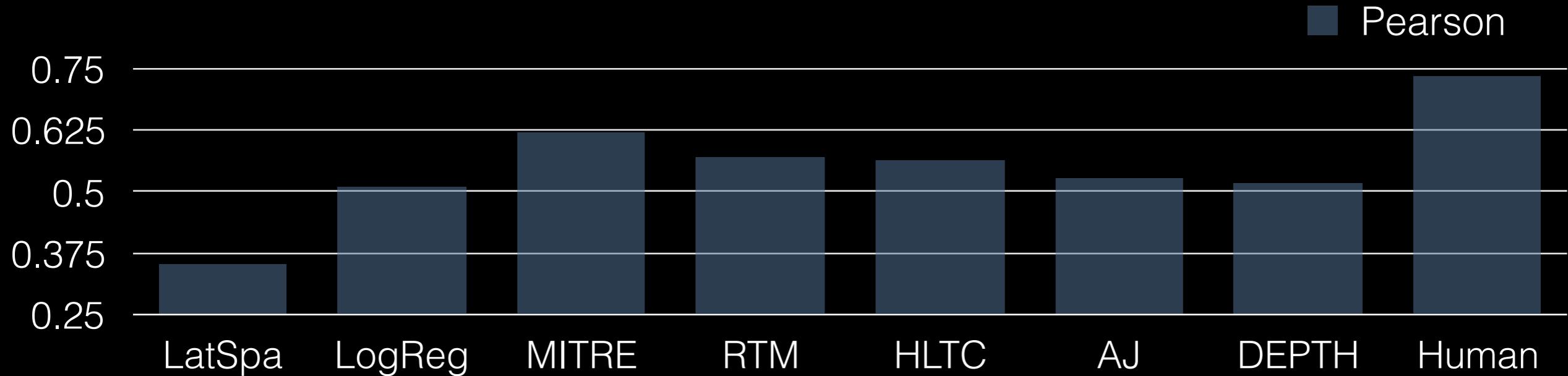
- Word Embeddings
 - String Similarities
 - Syntactic Similarities
 - Semantics
(LSA, CRM; WordNet; Entailment;
different parsers: constituency, logical, Boxer ...)
 - MT metrics
 - POS
 - Named Entity
 - various classifiers: SVM, MaxEnt, neural network
- ...

Performance

#1 Paraphrase Identification [PI] (binary 0 or 1)



#2 Semantic Similarity [SS] (degreeed 0~1)



Challenging Cases

paraphrase

(Mariano "Mo" Rivera is a baseball pitcher)

Classy gesture by the Mets for Mariano

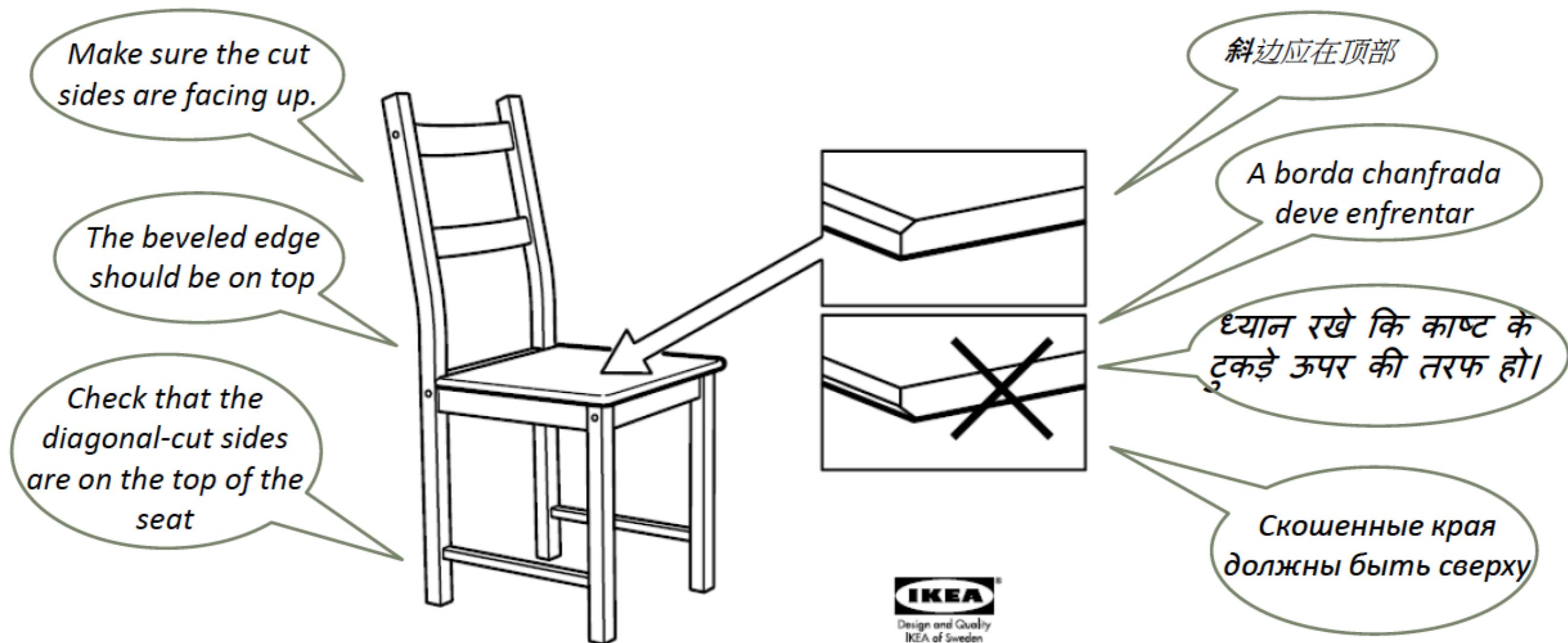
real class shown by the Mets Mo Rivera is a legend

non-paraphrase

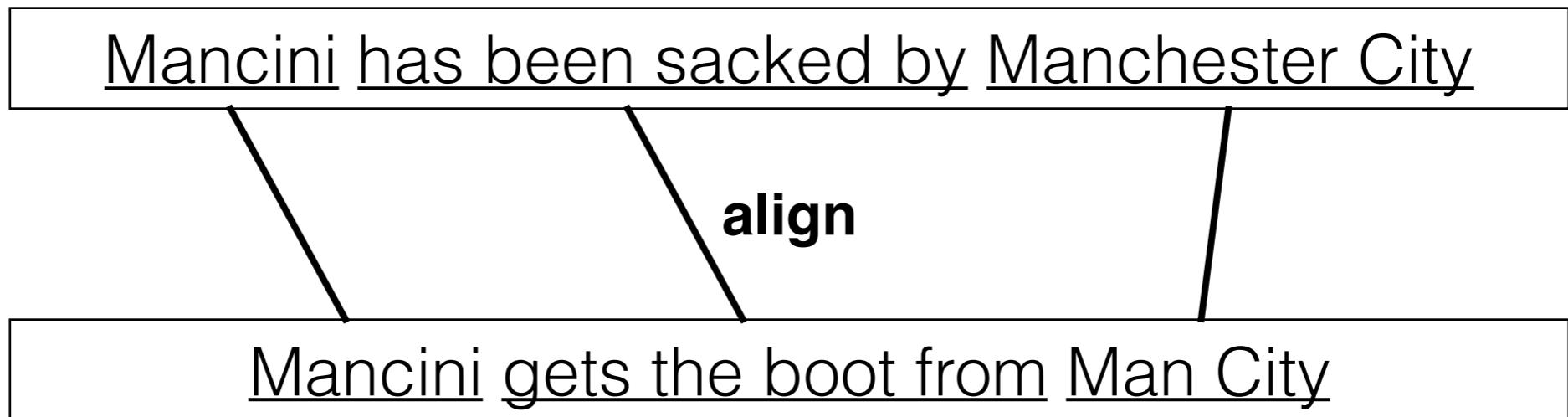
The world of jenks is such a real show

Jenks from the World of Jenks is such a good person

Multi-lingual Paraphrases



Extract Phrasal Paraphrases



Techniques

- Word/Phrase Alignment
- Probabilistic Graphical Models

Wei Xu, Joel Tetreault, Martin Chodorow, Ralph Grishman, Le Zhao.

"Exploiting Syntactic and Distributional Information for Spelling Correction with Web-Scale N-gram Models" In EMNLP (2011)

Wei Xu, Alan Ritter, Ralph Grishman. "Gathering and Generating Paraphrases from Twitter with Application to Normalization" In BUCC (2013)

Timothy Baldwin, Marie-Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, Wei Xu. "Shared Tasks of the

2015 ACL Workshop on Noisy User-generated Text: Twitter Lexical Normalization and Named Entity Recognition" In WNUT (2015)

Extract Phrasal Paraphrases



has been sacked by

gets the boot from

manchester city

man city

4

for

4

four

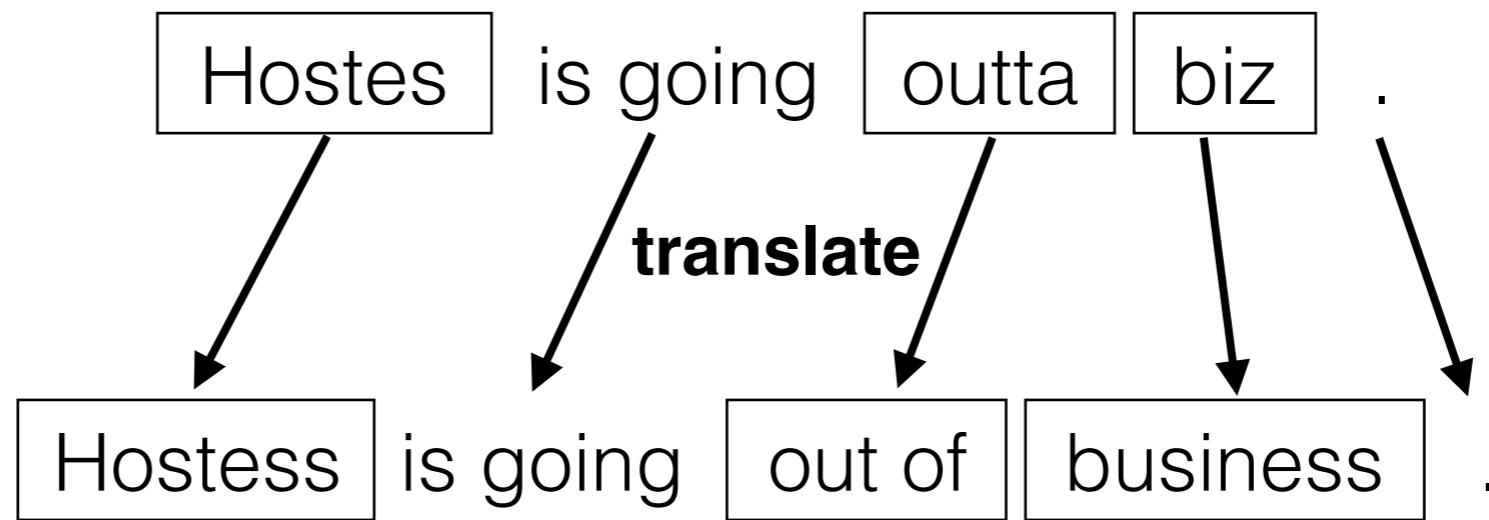
outta

out of

hostes

hostess

Noisy Text Normalization



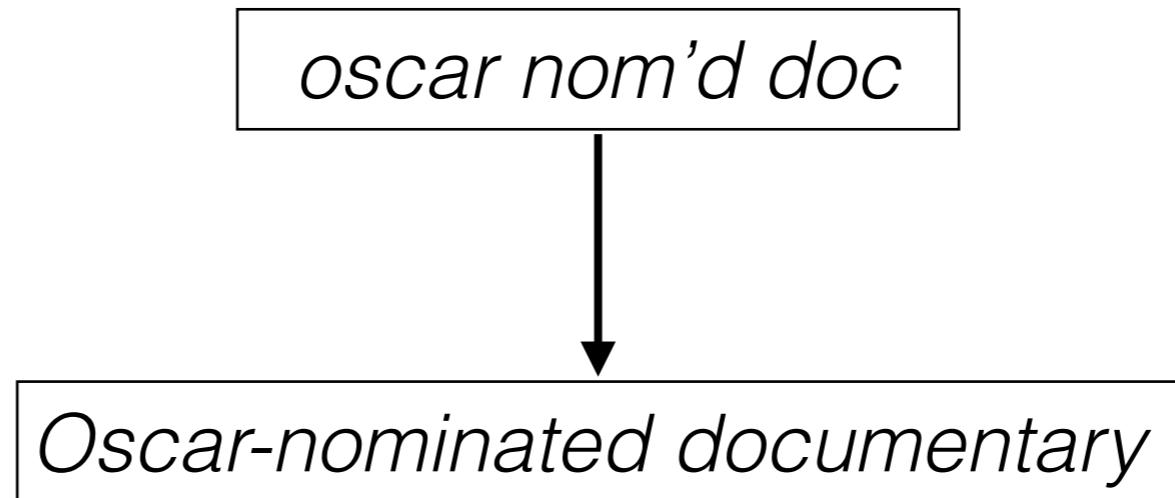
Wei Xu, Joel Tetreault, Martin Chodorow, Ralph Grishman, Le Zhao.

"Exploiting Syntactic and Distributional Information for Spelling Correction with Web-Scale N-gram Models" In EMNLP (2011)

Wei Xu, Alan Ritter, Ralph Grishman. "Gathering and Generating Paraphrases from Twitter with Application to Normalization" In BUCC (2013)

Timothy Baldwin, Marie-Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, Wei Xu. "Shared Tasks of the 2015 Workshop on Noisy User-generated Text: Twitter Lexical Normalization and Named Entity Recognition" In WNUT (2015)

Noisy Text Normalization



Wei Xu, Joel Tetreault, Martin Chodorow, Ralph Grishman, Le Zhao.

“Exploiting Syntactic and Distributional Information for Spelling Correction with Web-Scale N-gram Models” In EMNLP (2011)

Wei Xu, Alan Ritter, Ralph Grishman. “Gathering and Generating Paraphrases from Twitter with Application to Normalization” In BUCC (2013)

Timothy Baldwin, Marie-Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, Wei Xu. “Shared Tasks of the 2015 Workshop on Noisy User-generated Text: Twitter Lexical Normalization and Named Entity Recognition” In WNUT (2015)

Natural Language Generation



noisy



standard

'13* '14*



stylistic



plain

'12* '14*



complex



simple

'15* '16*

wen
or
when?

erroneous



correct

'11* '13*



feminine



masculine

'16*

and more (future work) ...



* my research

Voice Assistant



Unlimited Text in theory

“Almost any single (relatively complex) meaning can be implemented by an astonishingly high number of synonymous surface expressions.”

Meaning-Text Linguistic Theory (Žolkovskij & Mel’čuk, 1965; ~ now)

meaning = invariant of paraphrases

text = ‘virtual paraphrasing’

paraphrases = synonymous linguistic expressions

Unlimited Text in theory

“Almost any single (relatively complex) meaning can be implemented by an astonishingly high number of synonymous surface expressions.”

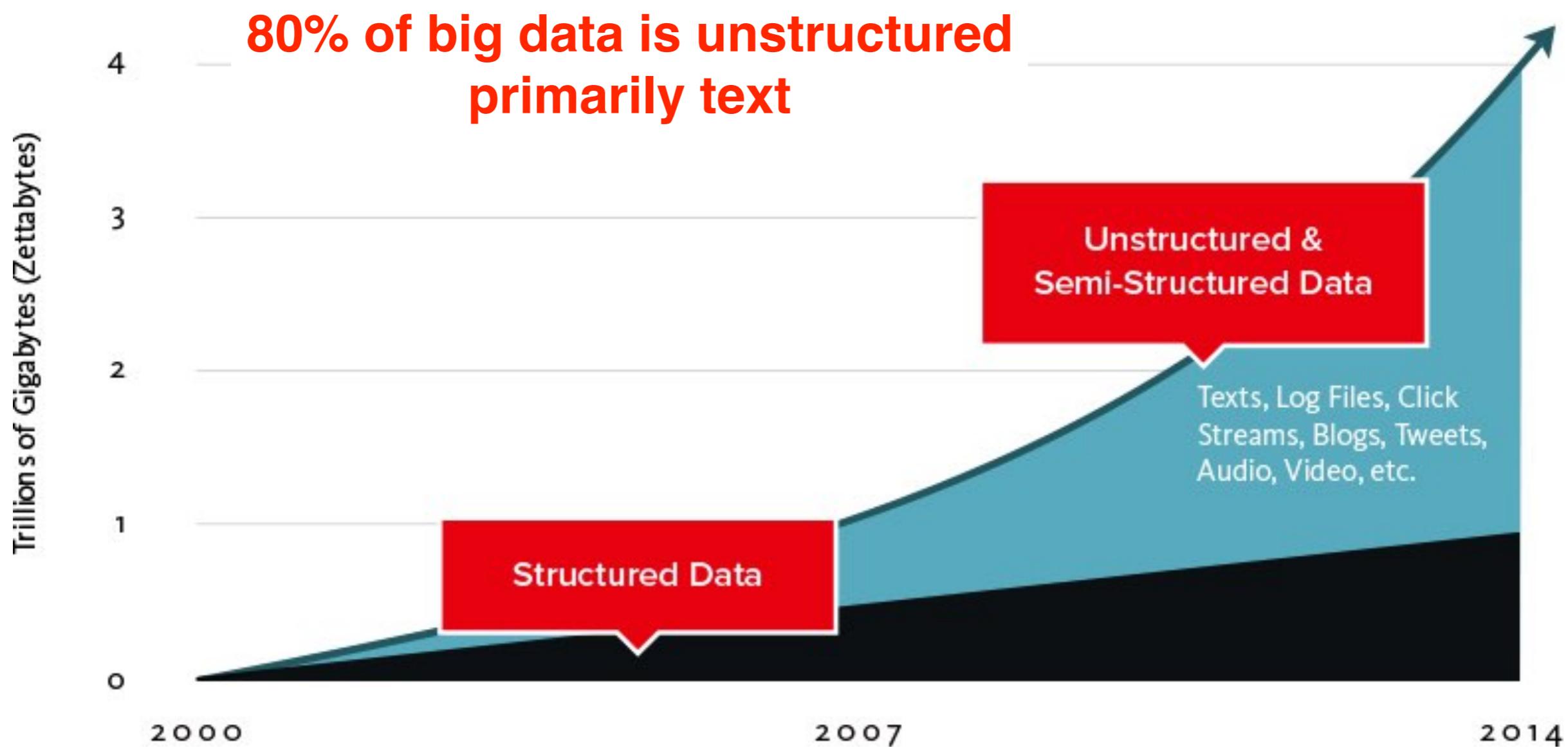
Meaning-Text Linguistic Theory (Žolkovskij & Mel’čuk, 1965; ~ now)

meaning = invariant of **paraphrases**

text = ‘virtual **paraphrasing**’

paraphrases = synonymous linguistic expressions

Unlimited Text in practice



(Source: IDC Research & Couchbase)

Social Science

“Parallel universe” experimental paradigm

Exploit situations with *many* instances of:

...the same speaker

...in the same situation, or

conveying the same info...

...varying their wording (beyond a fixed set of lexical choices)

and see the effects.



Relates to work on style (e.g., Annie Louis and Ani Nenkova, 2013) and paraphrasing (e.g., Wei Xu, Alan Ritter, Chris Callison-Burch, Bill Dolan, Yangfeng Ji, 2014)



Social Science

wonderfully delightfully beautifully fine well good nicely superbly



she says



he says

(also age & income)

Sentiment Analysis



*This nets vs bulls game is **great***

*This Nets vs Bulls game is **nuts***

Wowzers to this nets bulls game

*this Nets vs Bulls game is **too live***

*This Nets and Bulls game is a **good** game*

*This netsbulls game is **too good***

*This NetsBulls series is **intense***

Language Education

Aaaaaaaaaand Stephen Curry **is on fire**



What an incredible performance from Stephen Curry

Listen & Speak
Like a Native Speaker



thanku

thank u 4 ur time

Thank You

thankning you

gratitude

appreciate it

thx

3x

tyvm

thanks

say thanks

thank you very much

thnx

wawwww thankkkkkkkkkkk you alottttttttt!

thanks a lot

I am grateful

socialmedia-class.org