# Social Media & Text Analysis

lecture 7 - learn large-scale paraphrase from Twitter


Follow @cocoweixu

**Instructor: Wei Xu**
**Website: socialmedia-class.org**

# Paraphrase

| | | |
|---|---|---|
| *wealthy* | **word** | *rich* |
| *the king's speech* | **phrase** | *His Majesty's address* |
| *… the forced resignation of the CEO of Boeing, Harry Stonecipher, for …* | **sentence** | *… after Boeing Co. Chief Executive Harry Stonecipher was ousted from …* |

# Application

**Information Extraction**

end_job (Harry Stonecipher, Boeing)

↑ **extract**

| … the <u>forced resignation</u> of the CEO of Boeing, Harry Stonecipher, for … |
|---|

| … after Boeing Co. Chief Executive Harry Stonecipher was <u>ousted</u> from … |
|---|

# Application

**Question Answering**

Who is the CEO <u>stepping down</u> from Boeing?

**match**

*… the forced <u>resignation</u> of the CEO of Boeing, Harry Stonecipher, for …*

*… after Boeing Co. Chief Executive Harry Stonecipher was <u>ousted</u> from …*
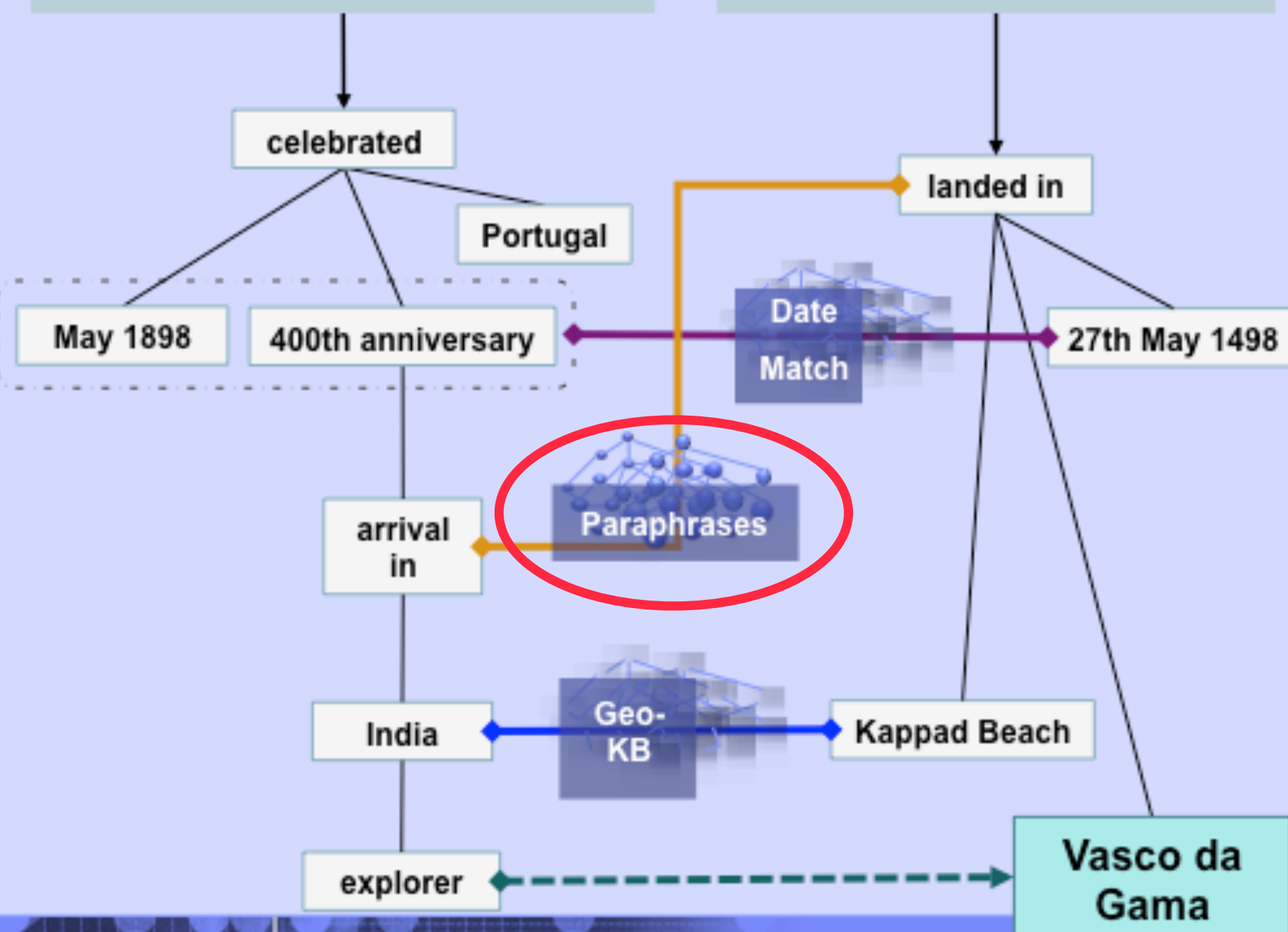
# Application

## Text Simplification

_They are culturally akin to_ the coastal peoples of Papua New Guinea.

↓

_Their culture is like that of_ the coastal peoples of Papua New Guinea.

Wei Xu ○ socialmedia-class.org

# Application

## Stylistic Rewriting



Palpatine:
*If you will not be turned, you will be destroyed!*

↓

*If you will not be turn'd, you will be undone!*

Luke:
*Father, please! Help me!*

↓

*Father, I pray you! Help me!*

Wei Xu, Alan Ritter, Bill Dolan, Ralph Grishman, Colin Cherry. "Paraphrasing for Style" In COLING (2012)

# Paraphrase Data

| 80s<br>WordNet | '01<br>Novels | '04<br>News | '05 '13<br>Bi-Text | '11<br>Video | '12<br>Shakespeare | (This Work)<br>Twitter |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| ● | ● | ● | ● | ● | ● | |

Barzilay
McKeown

Dolan
Quirk
Brockett

Callison-Burch
Ganitkevitch
Van Durme
Bernard

Chen
Dolan

Xu
Ritter
Dolan
Grishman
Cherry

Xu
Ritter
Callison-Burch
Dolan
Ji

# Paraphrase Research

**WordNet**   **Novel**   **News**   **Bi-Text**   **Video** **Shakespeare**              **Twitter**

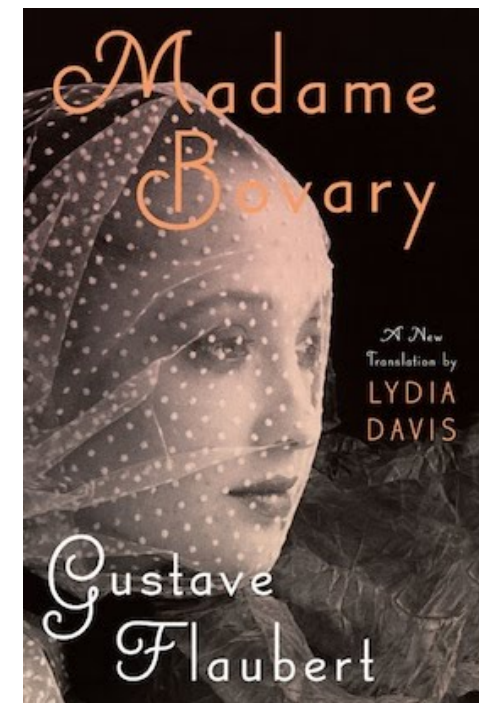●             ●           ●          ●              ●          ●

information extraction
question answering
semantic similarity
semantic parsing
text-to-text generation

…

but, primarily for formal language usage and
well-edited text

# Previous Work



multiple English translations of novels

(Barzilay and McKeown, 2001)

# Previous Work



only a few hundreds news agencies
report big events
using formal language

(Dolan, Quirk and Brockett, 2004; Dolan and Brockett, 2005; Brockett and Dolan, 2005)

# Previous Work



ask dozens of annotators to write
one sentence for a short video (<= 10 seconds)

(Chen and Dolan, 2011)

# Previous Work

... 5  farmers were    **thrown  into  jail**    in  Ireland ...

... fünf Landwirte    **festgenommen**    , weil  ...

... oder  wurden    **festgenommen**    , gefoltert ...

... or  have  been    **imprisoned**    ,  tortured ...

pivoting through bilingual text
from European Parliament proceedings,
multilingual websites etc.

(Bannard and Callison-Burch, 2005; Ganitkevitch,  Van Durme and Callison-Burch, 2013)

# 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

Wei Xu, Alan Ritter, Ralph Grishman.
"A Preliminary Study of Tweet Summarization using Information Extraction" in LASM (2014)

# Twitter as a powerful resource

thousands of users
talk about both big and micro events
using formal, informal, erroneous language

**Very diverse!**

# Enables new applications

**Noisy Text Normalization**

| oscar nom'd doc |
| --- |

↓

| Oscar-nominated documentary |
| --- |

| don't want for |
| --- |

↓

| don't wait for |
| --- |

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)

# Enables new applications

## Human-computer Interaction



Apple Siri     Google Now     Windows Cortana

who wants to get a beer?  →

*want to get a beer?*

*who else wants to get a beer?*

*who wants to go get a beer?*

*who wants to buy a beer?*

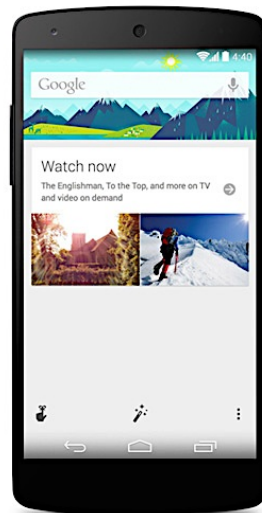*who else wants to get a beer?*

*trying to get a beer?*

*… (21 different ways)*

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

# Enables new applications

Listen & Speak
Like a Native Speaker



**Language Education**

Aaaaaaaaand stephen curry _is on fire_

What a incredible performance from Stephen Curry

# Enables new applications

**Sentiment Analysis**   🙂 or ☹ ?

| |
|---|
| *This nets vs bulls game is <u>great</u>* |

| |
|---|
| *This Nets vs Bulls game is <u>nuts</u>* |

| |
|---|
| *<u>Wowsers</u> to this nets bulls game* |

| |
|---|
| *this Nets vs Bulls game is <u>too live</u>* |

| |
|---|
| *This Nets and Bulls game is a <u>good</u> game* |

| |
|---|
| *This netsbulls game is <u>too good</u>* |

| |
|---|
| *This NetsBulls series is <u>intense</u>* |

# Learn Paraphrases

# Learn Paraphrases

**identify parallel sentences automatically from Twitter's big data stream**

| |
|---|
| *Mancini has been sacked by Manchester City* |
| *Mancini gets the boot from Man City* |

Yes!

| |
|---|
| *WORLD OF JENKS IS ON AT 11* |
| *World of Jenks is my favorite show on tv* |

No!

# Early Attempts on Twitter

- 1242 tweet pairs, tracking celebrity & hashtags
(Zanzotto, Pennacchiotti, Tsioutsiouliklis, 2011)

- named entity + date
(Xu, Ritter, Grishman, 2013)

- bilingual posts
(Ling, Dyer, Black, Trancoso, 2013)

# Named Entity + Time

**Tyler Anderson**
@tylerjanderson

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

**Jeff Clutter**
@Pibbbs

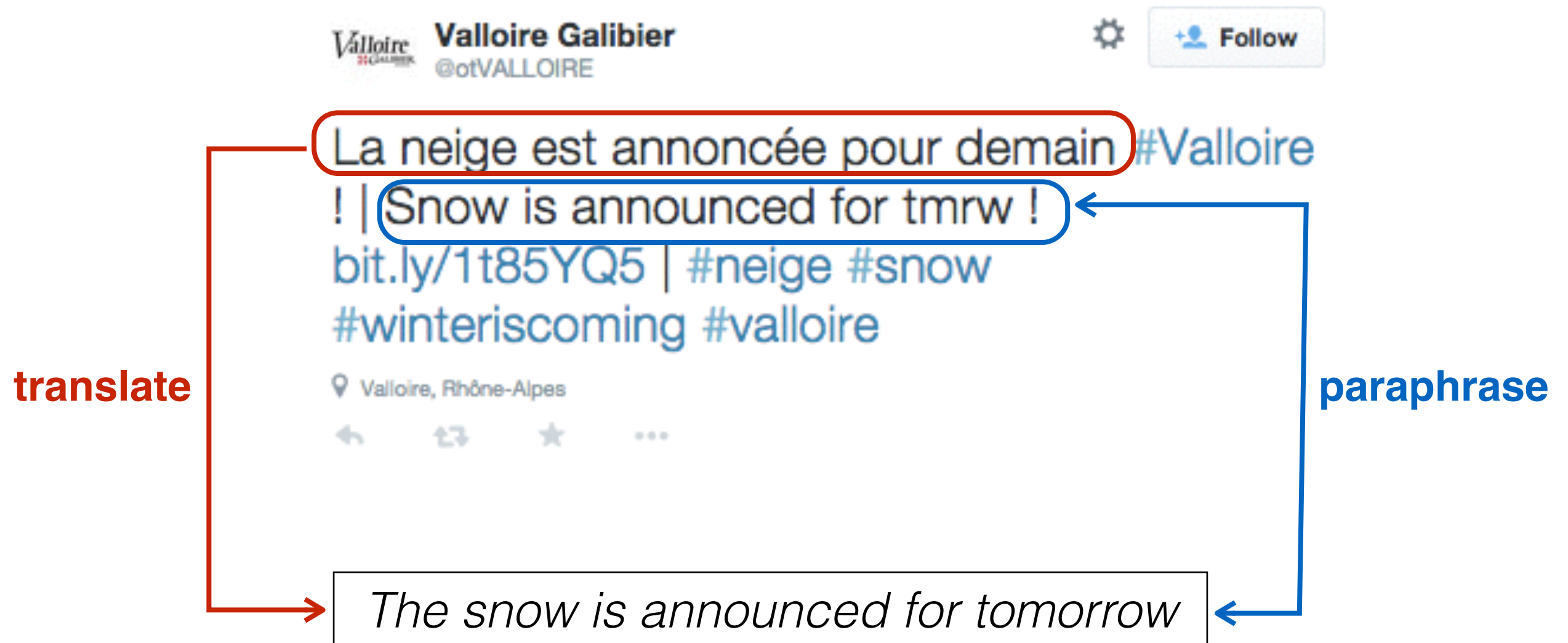Instagram can sell your photos without
consent starting January 16th,

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

# Self-translation



**translate**

**paraphrase**

*The snow is announced for tomorrow*

(Ling, Dyer, Black and Trancoso, 2013)

# Design a Model

# Train it on data

# A Challenge

| Mancini has been sacked by Manchester City |
|---|

| Mancini gets the boot from Man City |
|---|

very short
lexically divergent

(less word overlap, even in high-dimensional space)

# Previous Methods

| Algorithm | Reference | Description | Supervision |
|-----------|-----------|-------------|-------------|
| Vector Based | Mihalcea et al. (2006) | cosine similarity with tf-idf weighting | unsupervised |
| ESA | Hassan (2011) | explicit semantic space | unsupervised |
| KM | Kozareva and Montoyo (2006) | combination of lexical and semantic features | supervised |
| LSA | Hassan (2011) | latent semantic space | unsupervised |
| RMLMG | Rus et al. (2008) | graph subsumption | unsupervised |
| MCS | Mihalcea et al. (2006) | combination of several word similarity measures | unsupervised |
| WTMF | Guo and Diab (2012) | latent semantics model of missing words | unsupervised |
| STS | Islam and Inkpen (2007) | combination of semantic and string similarity | unsupervised |
| SSA | Hassan (2011) | salient semantic space | unsupervised |
| QKC | Qiu et al. (2006) | sentence dissimilarity classification | supervised |
| ParaDetect | Zia and Wasif (2012) | PI using semantic heuristic features | superviseded |
| SDS | Blacoe and Lapata (2012) | simple distributional semantic space | supervised |
| matrixJcn | Fernando and Stevenson (2008) | JCN WordNet similarity with matrix | unsupervised |
| FHS | Finch et al. (2005) | combination of MT evaluation measures as features | supervised |
| PE | Das and Smith (2009) | product of experts | supervised |
| WDDP | Wan et al. (2006) | dependency-based features | supervised |
| SHPNM | Socher et al. (2011) | recursive autoencoder with dynamic pooling | supervised |
| MTMETRICS | Madnani et al. (2012) | combination of eight machine translation metrics | supervised |
| LEXLATENT | Ji and Eienstein (2013) | combination of latent space and lexical features | supervised |

mostly based on sentence similarity of surface words or latent semantics

# Design a Model

**At-least-one-anchor Assumption**

two sentences about the same <u>topic</u> are paraphrases
if and only if
they contain at least one word pair that is a paraphrase **anchor**

| |
|---|
| *That boy <u>Brook Lopez</u> with a deep **3*** |
| *<u>brook lopez</u> hit a **3*** |

Yes!

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

# Another Challenge

not every word pair of similar meaning indicates
sentence-level paraphrase

| |
|---|
| *Iron Man **3** was brilliant fun* |

| |
|---|
| *Iron Man **3** tonight see what this is like* |

← No!

Solution:
   a discriminative model using features at word-level

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

# Features

- String Features:

  - words, stemmed forms, normalized forms

  - same, similar or dissimilar

- POS Features:

  - fine grained tags:
        "a", "be", "do", "have",
        "get", "go", "follow", "please"

- Topical Features:

  - word significantly associated with each topic

  - e.g. "3" for basketball; "RIP" for death events
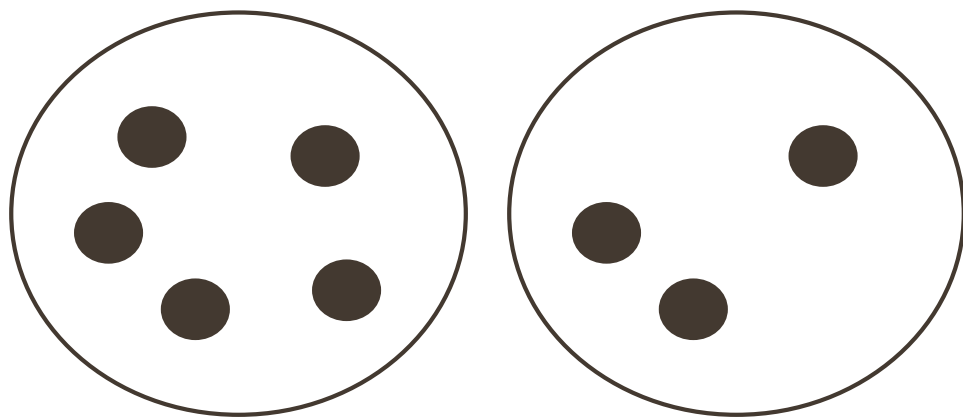
# Multi-instance Learning Paraphrase Model

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

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

$\gamma$ paraphrase

$\gamma$ non-paraphrase

sentence pair

word pair

$Z_1$ 0

$Z_2$ 0

$Z_3$ 1

$Z_4$ 0

...

**manti** | teo

**be** | is

**next** | new

**manti** | little

...

diff_word
same_pos_nn
both_sig
...

same_stem
same_pos_be
not_both_sig
...

diff_word
same_pos_jj
both_sig
...

diff_word
diff_pos_nn
diff_pos_jj
not_both_sig
...

features

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

# [Mini Tutorial]
# Multi-instance Learning

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

Negative Bags

Positive Bags



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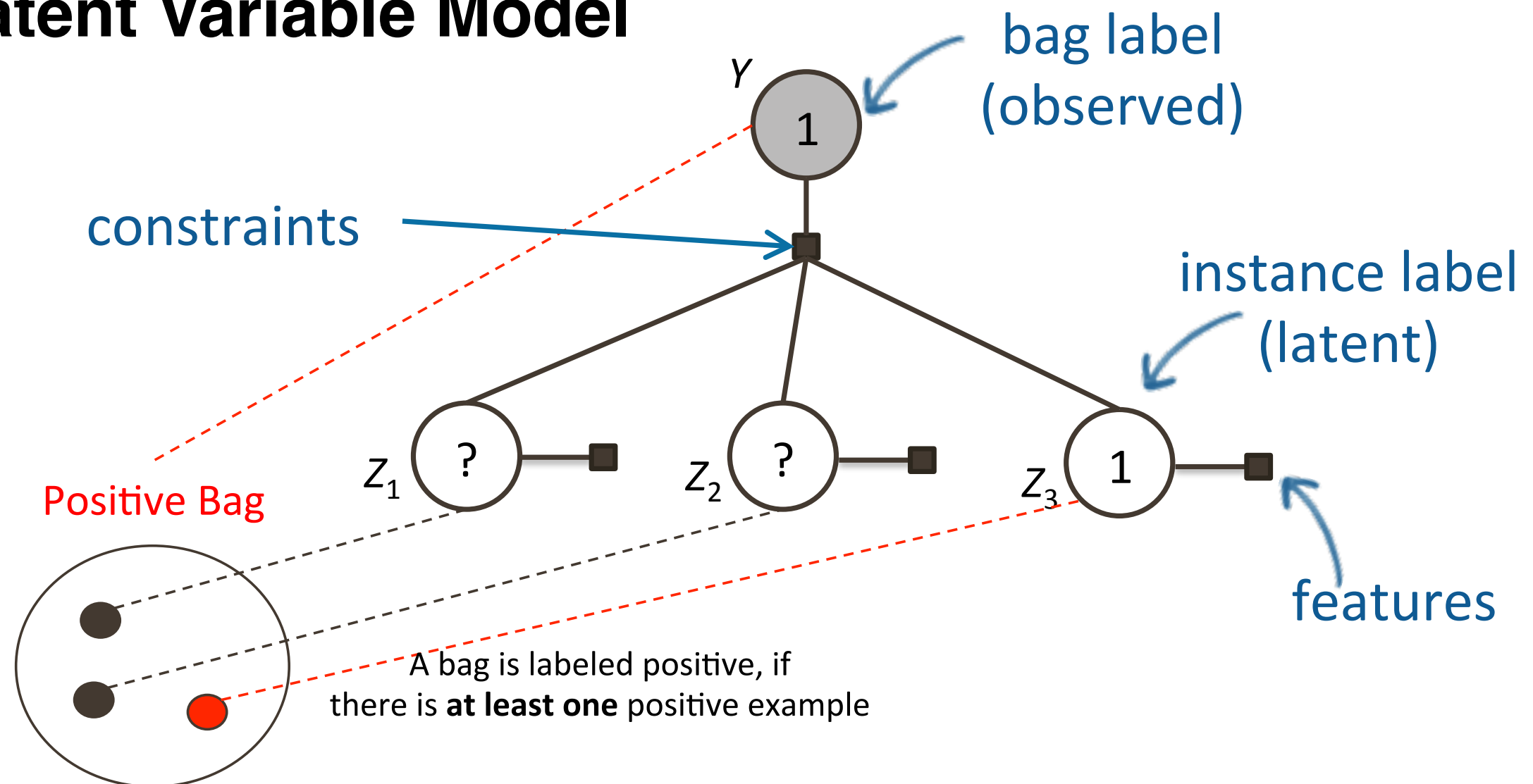**all** the examples in it are negative

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

(Dietterich et al., 1997)

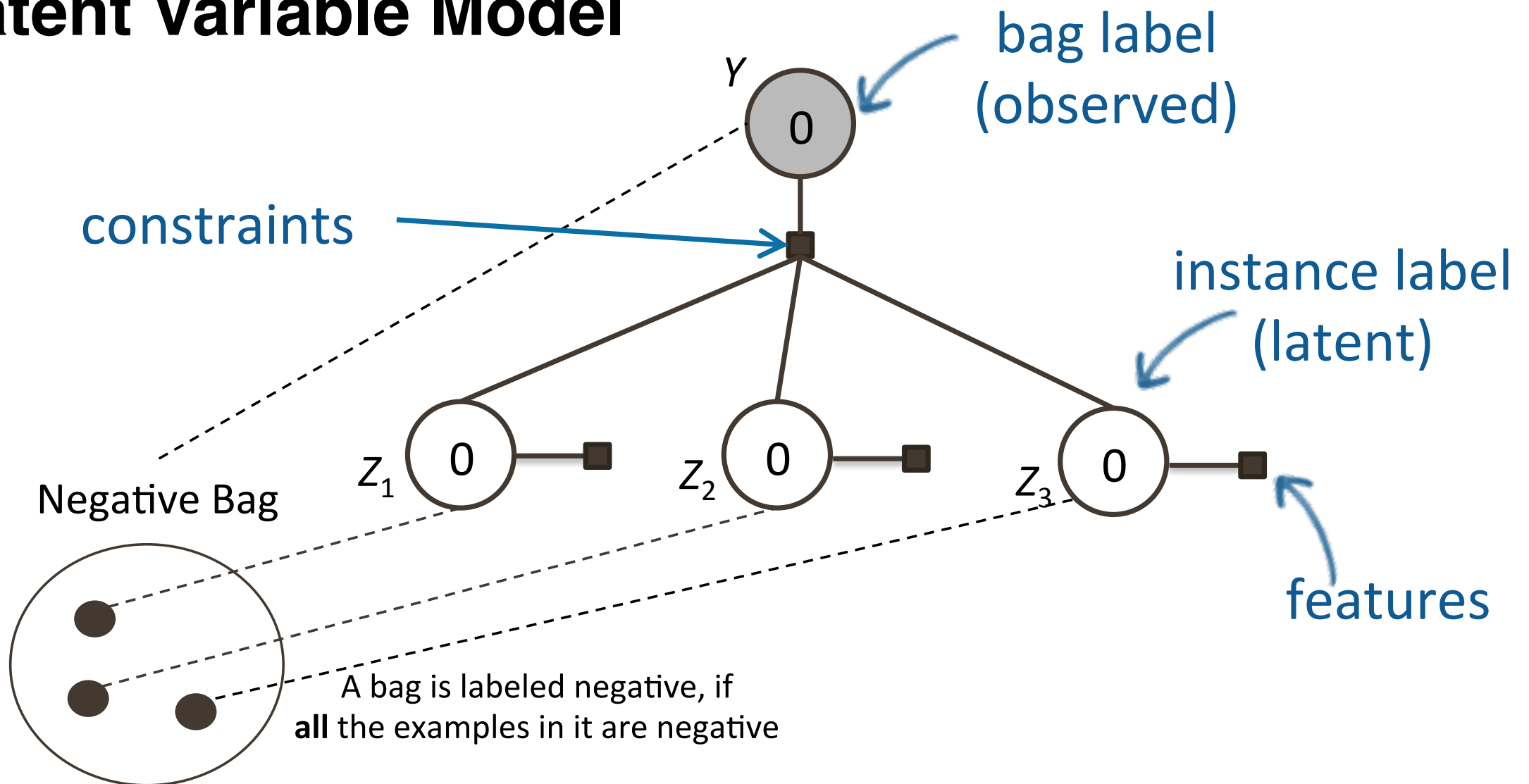# [Mini Tutorial] Multi-instance Learning

**Latent Variable Model**



bag label
(observed)

*Y*

constraints

instance label
(latent)

Positive Bag

$Z_1$   ?   $Z_2$   ?   $Z_3$   1

features

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

# [Mini Tutorial]
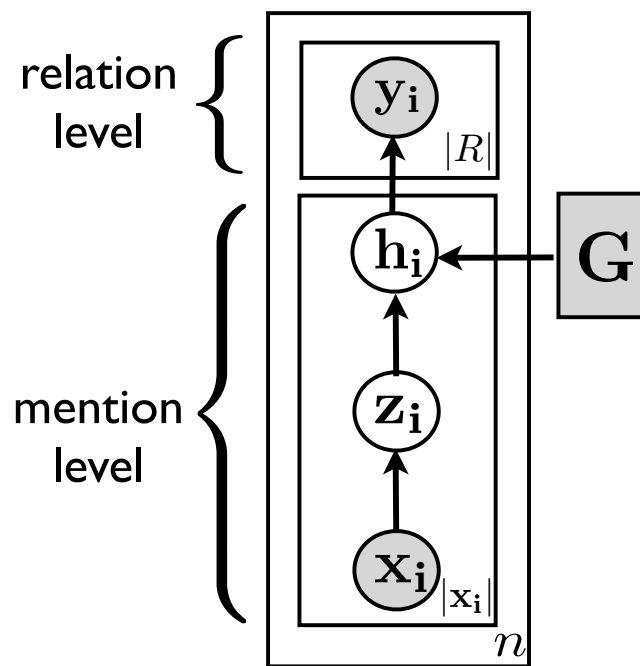## Multi-instance Learning

**Latent Variable Model**

bag label
(observed)

$Y$

constraints

instance label
(latent)

$Z_1$   $Z_2$   $Z_3$

Negative Bag

features

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

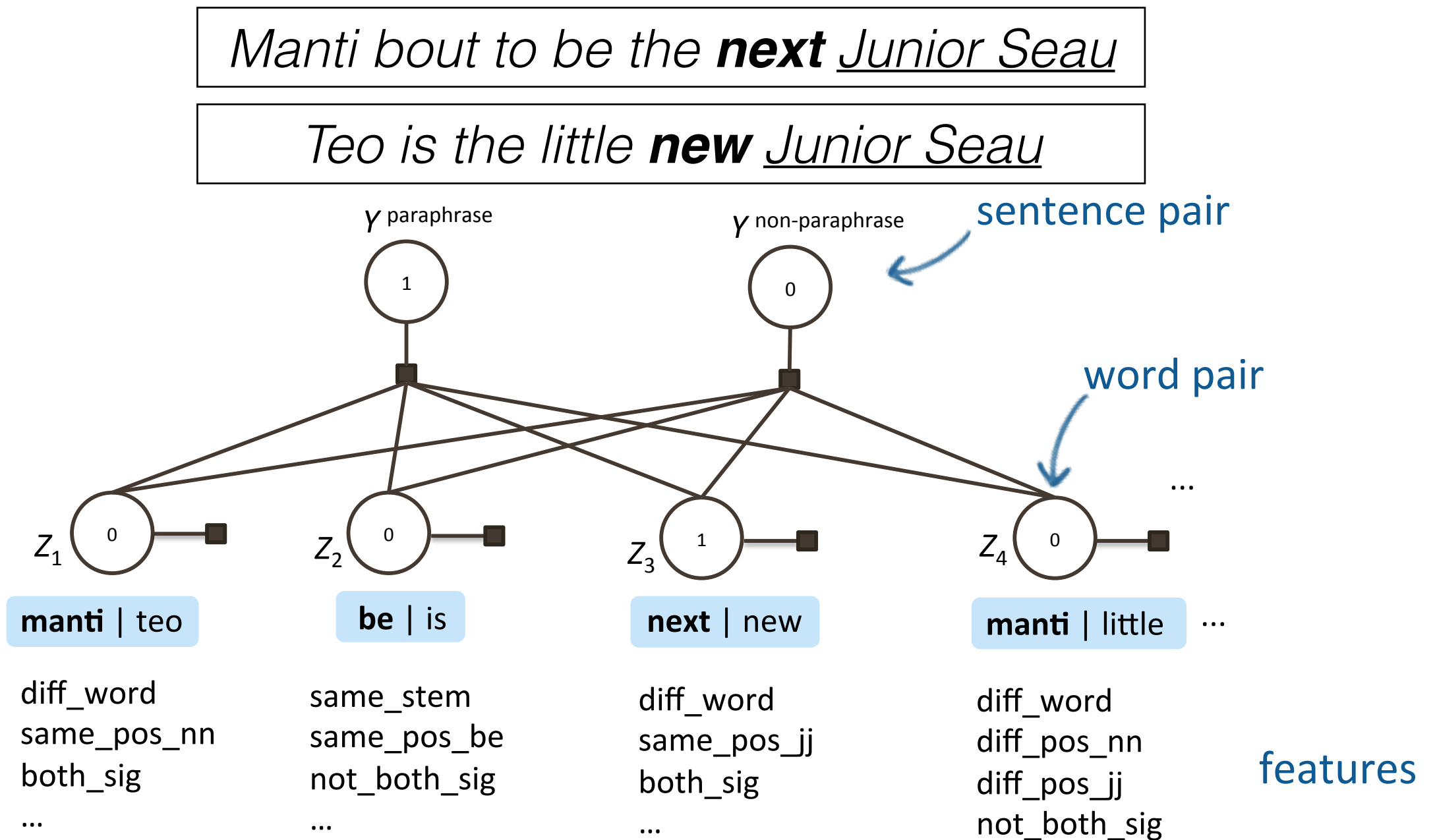# [Mini Tutorial]
# Multi-instance Learning

**Distantly Supervised Information Extraction**



1. incomplete knowledge base problem

2. distant supervision + human-labeled data

3. IE + IR

Wei Xu, Ralph Grishman, Le Zhao. "Passage Retrieval for Information Extraction using Distant Supervision" In IJCNLP (2011)

Wei Xu, Raphael Hoffmann, Le Zhao, Ralph Grishman. "Filling Knowledge Base Gaps for Distant Supervision of Relation Extraction" In ACL (2013)

Maria Pershina, Bonan Min, Wei Xu, Ralph Grishman. "Infusion of Labeled Data into Distant Supervision for Relation Extraction" In ACL (2014)
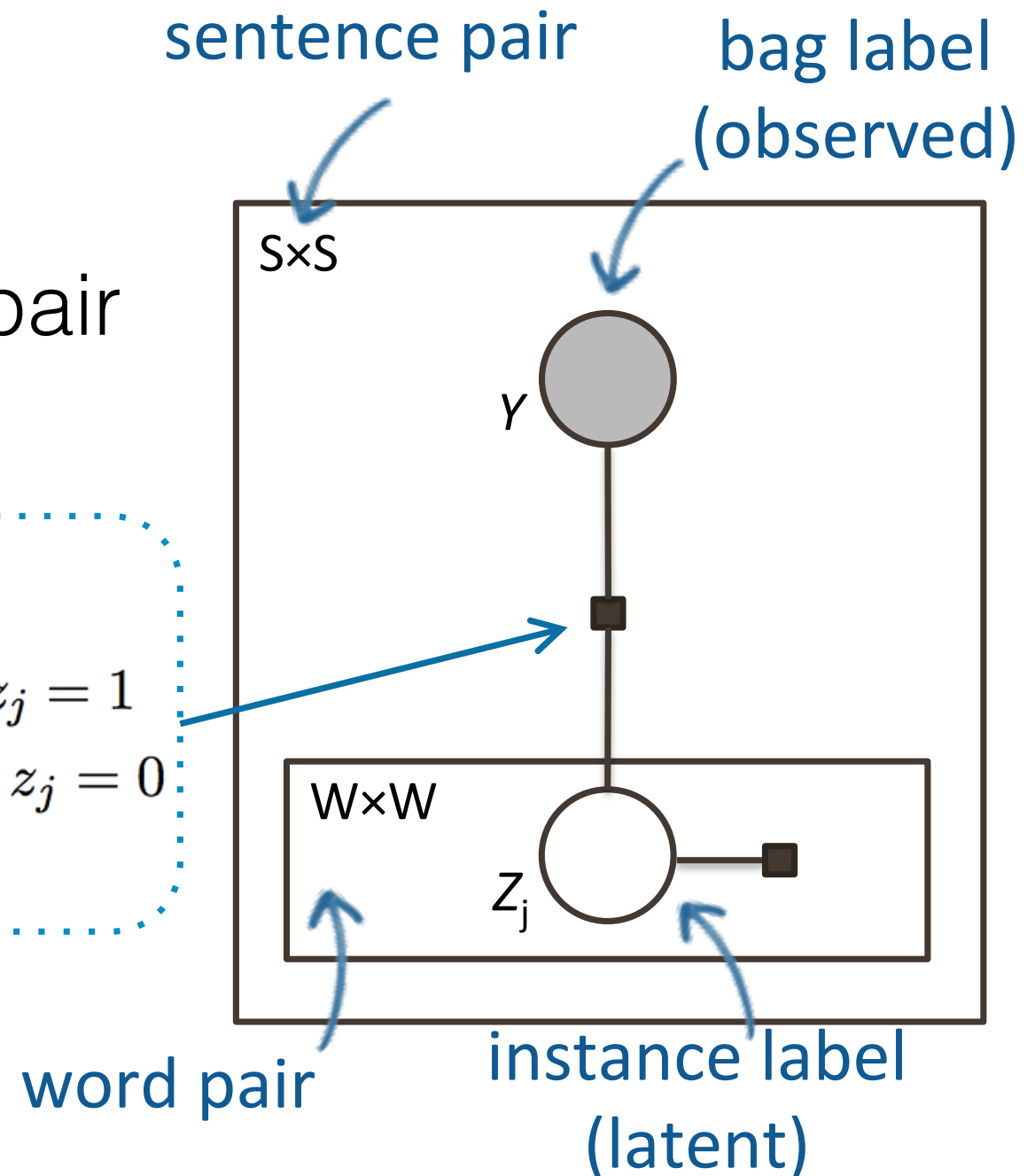
# [Recap] Multi-instance Learning Paraphrase Model

Manti bout to be the **next** _Junior Seau_

Teo is the little **new** _Junior Seau_

$\gamma$ paraphrase

$\gamma$ non-paraphrase

sentence pair

1

0

word pair

$Z_1$ 0

$Z_2$ 0

$Z_3$ 1

$Z_4$ 0

...

**manti** | teo

**be** | is

**next** | new

**manti** | little ...

diff_word
same_pos_nn
both_sig
...

same_stem
same_pos_be
not_both_sig
...

diff_word
same_pos_jj
both_sig
...

diff_word
diff_pos_nn
diff_pos_jj
not_both_sig

features

...

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

# Joint Word-Sentence Model

**Model the assumption:**
sentence-level paraphrase
is anchored by at-least-one word pair

deterministic OR

$$\sigma(\mathbf{z}_i, y_i) = \begin{cases} 1 & \text{if } y_i = true \wedge \exists j : z_j = 1 \\ 1 & \text{if } y_i = false \wedge \forall j : z_j = 0 \\ 0 & \text{otherwise} \end{cases}$$

sentence pair

bag label
(observed)

S×S

$Y$

W×W

$Z_j$

word pair

instance label
(latent)

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

# Joint Word-Sentence Model



$i$th sentence pair's label
(observed or to be predicated)

$j$th word pair

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

parameters       features       deterministic OR

latent labels for all word pairs
in the $i$th sentence pair

# Learning Algorithm

**Objective:**
learn the parameters that maximize
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

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

# Learning Algorithm

**Perceptron-style Update:**

Viterbi approximation + online learning
O(# word pairs)



$$\frac{\partial \log P(\mathbf{y}|\mathbf{w};\theta)}{\partial \theta} \approx \sum_i f(\mathbf{z}_i^*, \mathbf{w}_i) - \sum_i f(\mathbf{z}_i', \mathbf{w}_i)$$

**reward correct
(conditioned on labels)**

**penalize wrong
(ignoring labels)**

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

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

# Training Data

# Twitter Trends

# Annotation

**Crowdsourcing**

(Courtesy: The Sheep Market by Aaron Koblin)

# Annotation

## Crowdsourcing

**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

amazon mechanical turk™
Artificial Artificial Intelligence

Wei Xu. "Data-driven Approaches for Paraphrasing Across Language Variations" PhD Thesis, New York University. (2014)

# A Problem

only **8%** sentence pairs about the same topic
have similar meaning

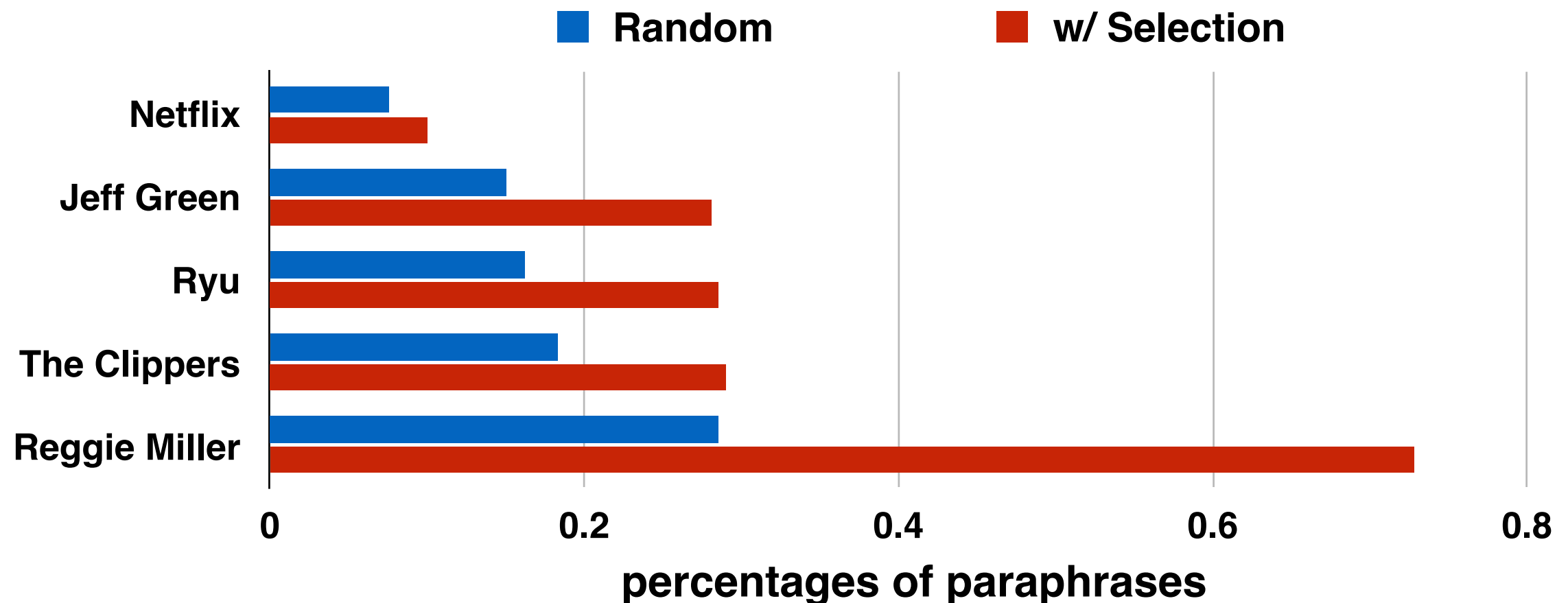hurts both quantity and quality

non-experts lower their bars

Wei Xu. "Data-driven Approaches for Paraphrasing Across Language Variations" PhD Thesis, New York University. (2014)

# Sentence Selection

## SumBasic Algorithm

**8%** → **16%**

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



■ **Random**　　　■ **w/ Selection**

percentages of paraphrases

Wei Xu. "Data-driven Approaches for Paraphrasing Across Language Variations" PhD Thesis, New York University. (2014)

# Topic Selection

**Multi-Armed Bandits**

**16%** ➞ **34%**

$$\max \sum_{\{i \mid r_i(t_0) > 0\}} \hat{\mu}_i(t_0) r_i(t_1)$$

$$\text{s.t.} \quad \sum_i c_i r_i(t_1) \leq (1 - \epsilon)B, \forall i : 0 \leq r_i(t_1) \leq l - r_i(t_0).$$

Wei Xu. "Data-driven Approaches for Paraphrasing Across Language Variations" PhD Thesis, New York University. (2014)

# Twitter Paraphrase Dataset
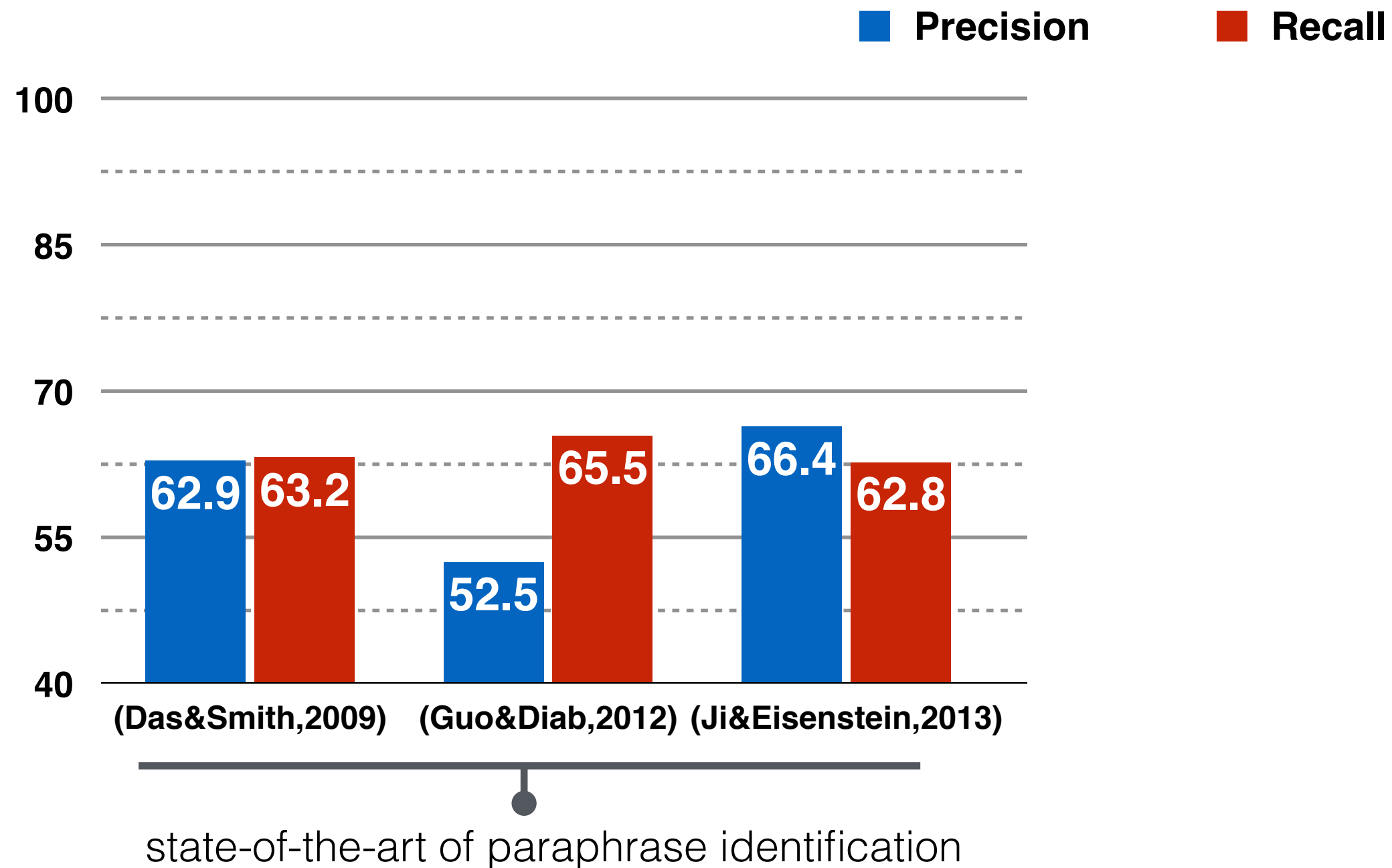
18,762 sentence pairs labeled
cost only $200

important but difficult to obtain
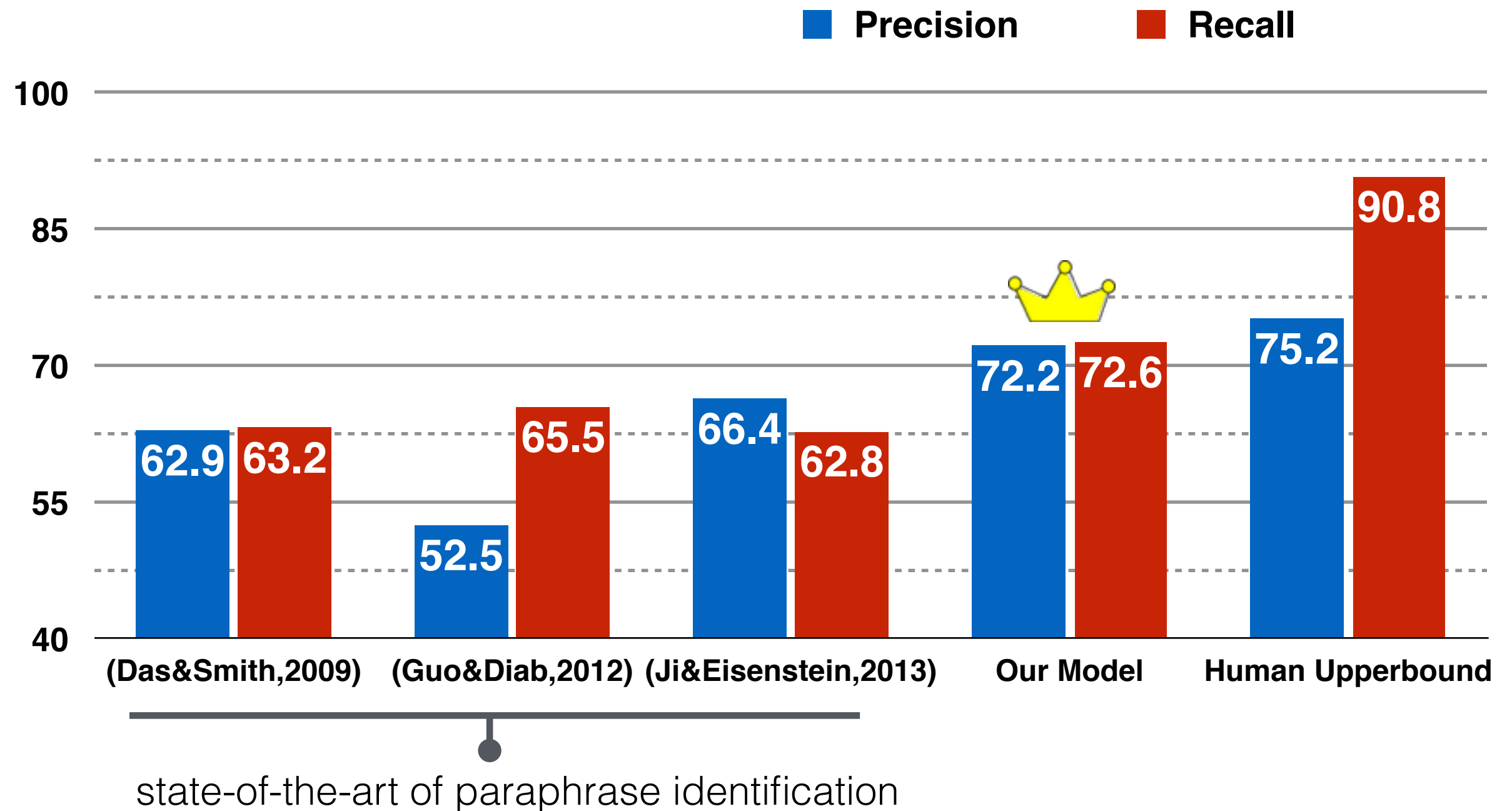
1/3 paraphrase, 2/3 non-paraphrase (very balanced)

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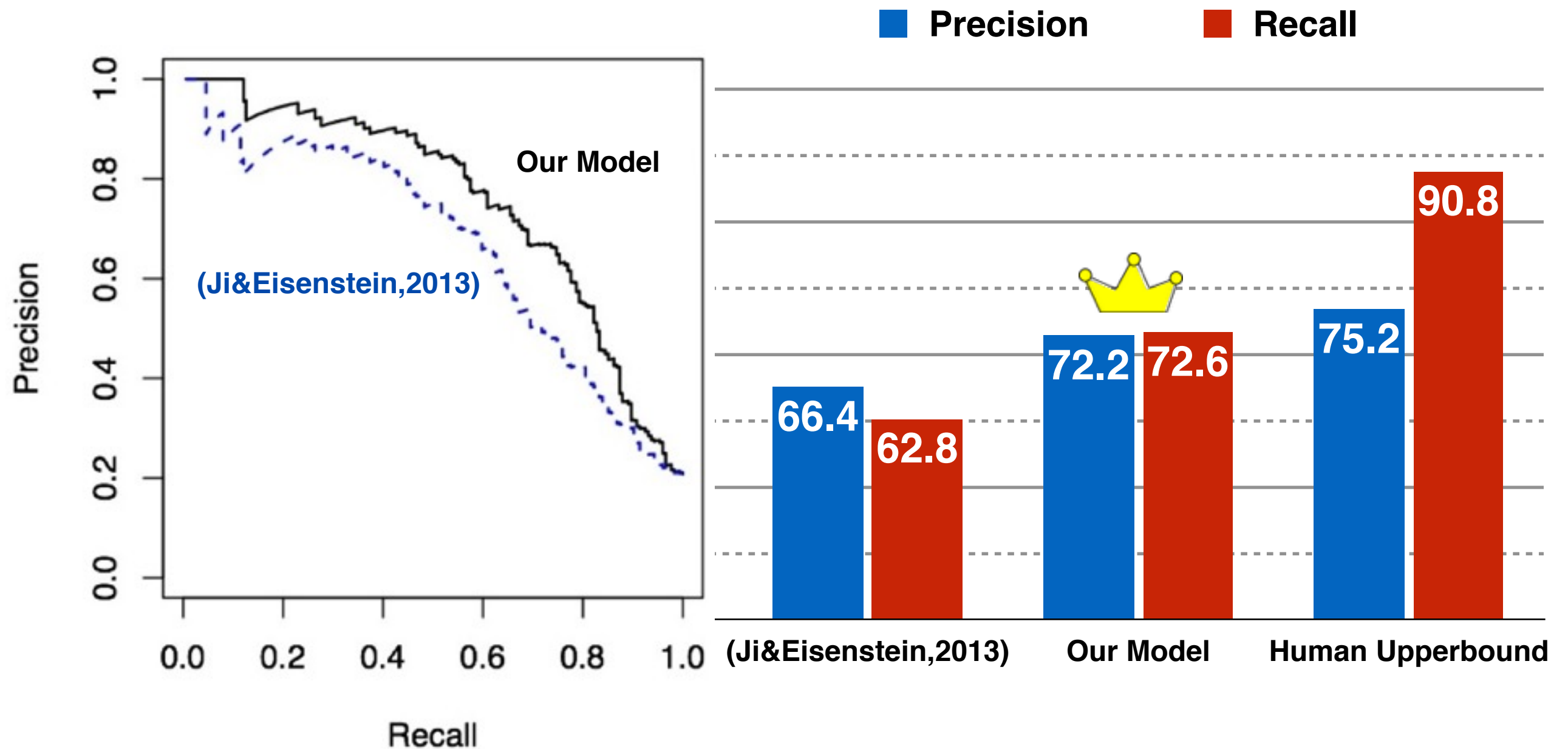
# Performance

# Performance

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

# Performance

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

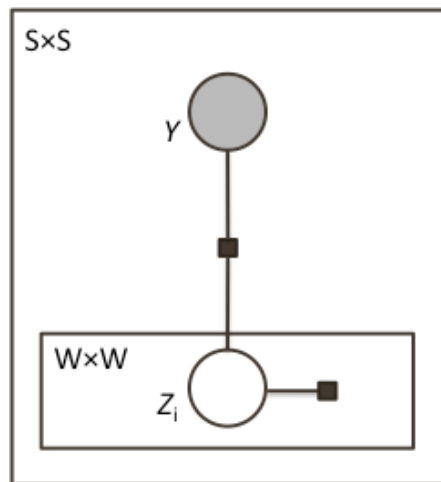# Performance

# Innovations

That boy <u>Brook Lopez</u> with a deep **3**

<u>brook lopez</u> hit a **3**

Yes!

## Multi-instance Learning Paraphrase Model (MultiP)



- Twitter's big data stream
- potential beyond Twitter and English
- joint sentence-word alignment
- extensible latent variable model

(a lot of space for future work)

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

# Impact

Shared Task: Paraphrase and Semantic Similarity in Twitter
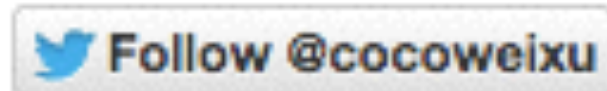19 research groups participated (100+ requested the data)

TU Munich
FBK

MITRE
Stanford
UMBC
UMD
Columbia

U Groningen
U Zagreb
U Edinburgh
U Sussex
Dublin City U
MTA

East China U
Wuhan U
HK UST

U Tokyo

Masaryk U
Amrita U

Wei Xu, Chris Callison-Burch, Bill Dolan.
"SemEval-2015 Task 1: Paraphrase and Semantic Similarity in Twitter" In SemEval (2015)

thank you very much

thank u 4 ur time

# Thank you

thanks

thanking you

appreciate it

gratitude

thx

3x

Follow @cocoweixu

say thanks

**Instructor: Wei Xu**

**www.cis.upenn.edu/~xwe/**

tyvm

thnx

**Course Website: socialmedia-class.org**

wawwww thankkkkkkkkkkk you alottttttttttt!

thanks a lot

am grateful