# From Shakespeare to Twitter:
# What are Language Styles all about?

Wei Xu

Department of Computer Science and Engineering

THE OHIO STATE UNIVERSITY

Follow @cocoweixu

Sep-08-2017 @ EMNLP 2017

# My very first Shakespeare play



WILLIAM SHAKESPEARE'S
# THE TEMPEST
DIRECTED BY MADELINE SAYET

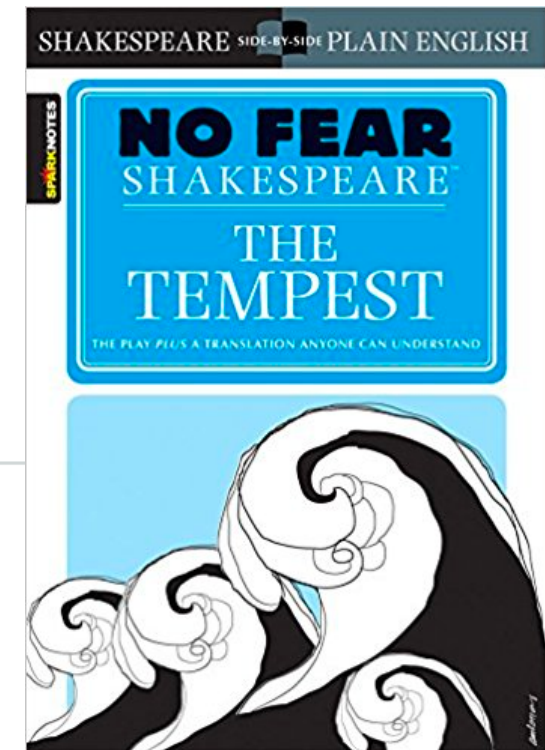*Shakespeare's The Tempest directed by Madeline Sayet*
**March 29-31**
Brooklyn Lyceum, 227 4th Ave.
www.brooklynlyceum.com

**(Brooklyn - March 31, 2012)**

# How can I learn Early Modern English in no time?

**ORIGINAL TEXT**

**BOATSWAIN**

Here, master. What cheer?

**MASTER**

Good, speak to th' mariners. Fall to 't yarely, or we run ourselves aground. Bestir, bestir.

**MODERN TEXT**

**BOATSWAIN**

I'm here, sir. How can I help you?

**MASTER**

My good boy, give the other sailors a pep talk—and do it fast, before we're shipwrecked. Hurry, hurry!

**I trained machine translation models …**

If you will not be turned, you will be destroyed!

Shakespearized (real system output):

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

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

# My very first paper on Stylistics

## Paraphrasing for Style

Wei Xu[1]  Alan Ritter[2]  William B. Dolan[3]  Ralph Grishman[1]  Colin Cherry[4]

(1) New York University
(2) University of Washington
(3) Microsoft Research
(4) National Research Council, Canada

xuwei@cims.nyu.edu, aritter@cs.washington.edu, billdol@microsoft.com,
grishman@cs.nyu.edu, Colin.Cherry@nrc-cnrc.gc.ca

ABSTRACT

We present initial investigation into the task of paraphrasing language while targeting a particular writing style. The plays of William Shakespeare and their modern translations are used as a testbed for evaluating paraphrase systems targeting a specific style of writing. We show that even with a relatively small amount of parallel training data, it is possible to learn paraphrase models which capture stylistic phenomena, and these models outperform baselines based on dictionaries and out-of-domain parallel text. In addition we present an initial investigation into automatic evaluation metrics for paraphrasing writing style. To the best of our knowledge this is the first work to investigate the task of paraphrasing text with the goal of targeting a specific style of writing.

KEYWORDS: Paraphrase, Writing Style.

**(COLING - June 15, 2012)**
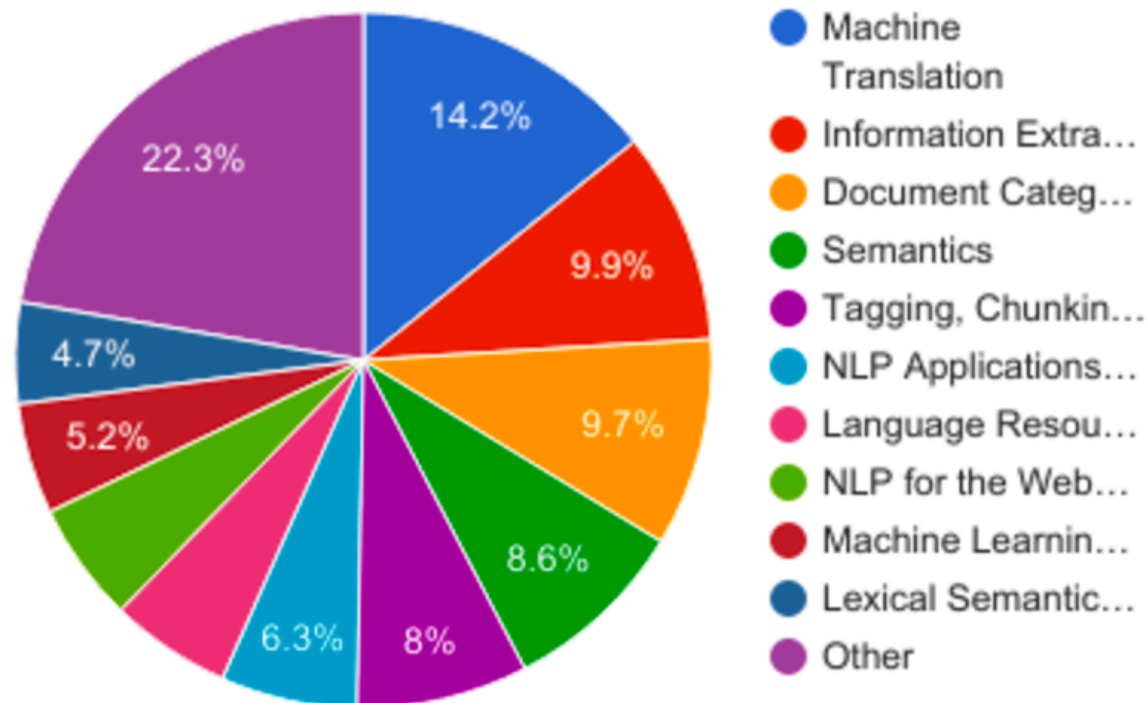
# Stylistic Paraphrase becomes a thing!



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

- Had Kabbara, Jackie Cheung. "Stylistic Transfer in Natural Language Generation Systems Using Recurrent Neural Networks" (EMNLP Uphill Battles 2016)

- Harsh Jhamtani, Varun Gangal, Eduard Hovy, Eric Nyberg "Shakespearizing Modern Language Using Copy-Enriched Sequence to Sequence Models" (**EMNLP StyVa 2017**)

- Se Won Jang, Jesik Min, Mark Kwon. "Writing Style Conversion using Neural Machine Translation" (Stanford CS224n 2017)
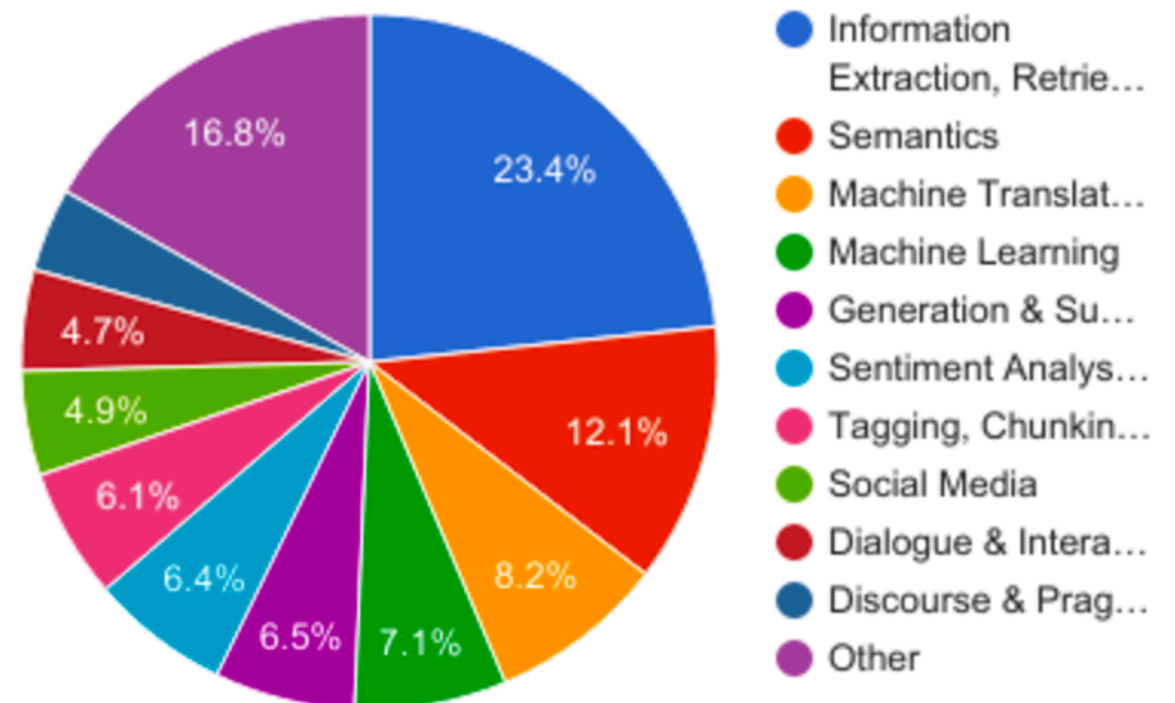
**Paraphrasing ≈ Monolingual Translation ≈ T2T Generation**

# Stylistic *text-to-text* Generation becomes a thing!



**ACL 2014 Submissions**

- Machine Translation
- Information Extra...
- Document Categ...
- Semantics
- Tagging, Chunkin...
- NLP Applications...
- Language Resou...
- NLP for the Web...
- Machine Learnin...
- Lexical Semantic...
- Other

14.2% / 9.9% / 9.7% / 8.6% / 8% / 6.3% / 5.2% / 4.7% / 22.3%

**ACL 2017 Submissions**

- Information Extraction, Retrie...
- Semantics
- Machine Translat...
- Machine Learning
- Generation & Su...
- Sentiment Analys...
- Tagging, Chunkin...
- Social Media
- Dialogue & Intera...
- Discourse & Prag...
- Other

23.4% / 16.8% / 12.1% / 8.2% / 7.1% / 6.5% / 6.4% / 6.1% / 4.9% / 4.7%
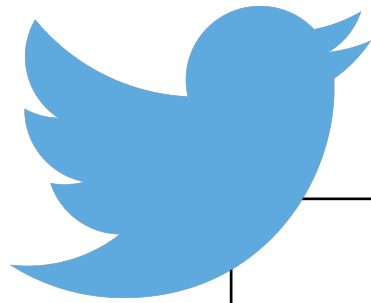
Generation & Summarization is now in top 5 areas,
while in 2014 it didn't even make top 10

# Can help User-generated Text

Hostes is going outta biz.

Normalized:

Hostess is going out of business.

Wei Xu, Joel Tetreault, Martin Chodorow, Ralph Grishman, Le Zhao.
"Exploiting Syntactic and Distributional Information for Spelling Correction with Web-Scale N-gram Models" (EMNLP 2011)

Wei Xu, Alan Ritter, Ralph Grishman. "Gathering and Generating Paraphrases from Twitter with Application to Normalization" (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" (WNUT 2015)

# Can help children read
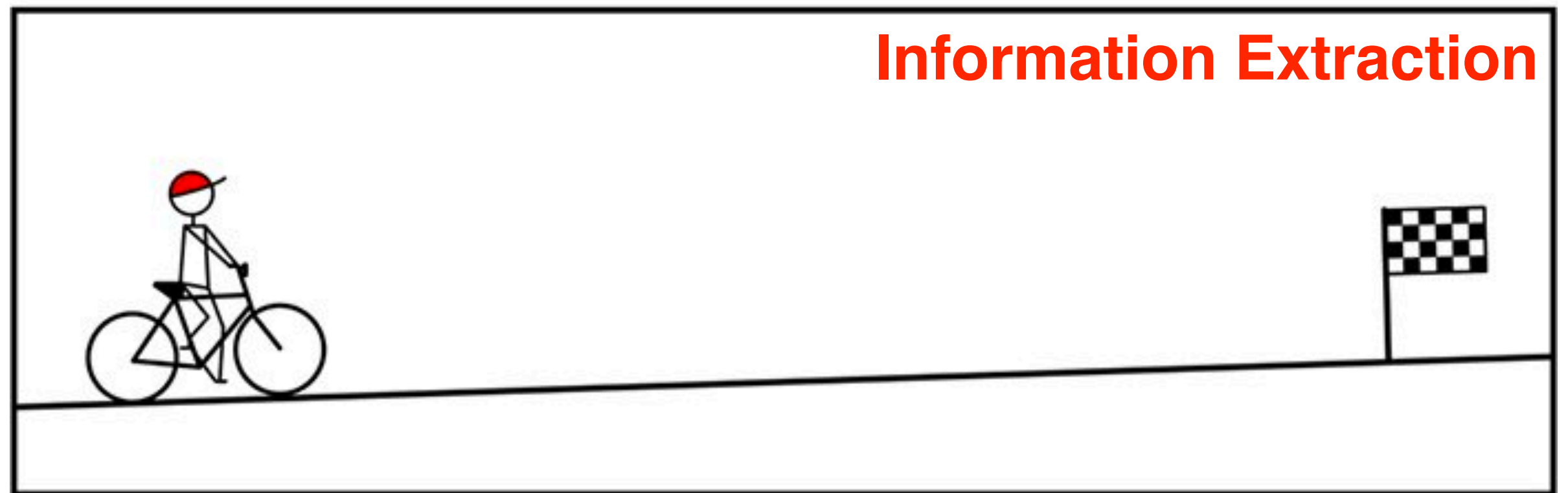
Slightly more fourth-graders nationwide are reading proficiently compared with a decade ago.

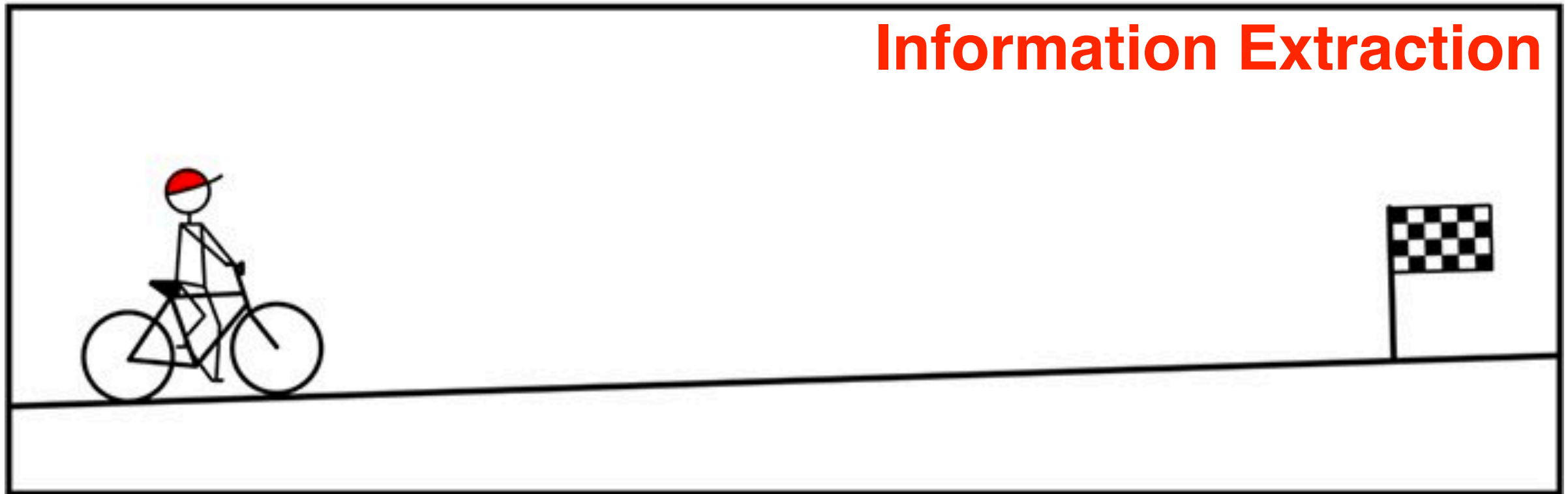Most fourth-graders are better readers than they were 10 years ago.

Wei Xu, Chris Callison-Burch, Courtney Napoles "Problems in Current Text Simplification Research: New Data Can Help" (TACL 2015)
Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

# My Research

My plan



**Information Extraction**

My plan

**Information Extraction**

Reality

**Paraphrase/Stylistics is exciting!**
**(with many pitfalls)**

Peak Experience

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

# My Research on Stylistics/Paraphrase

| | | | |
|---|---|---|---|
| erroneous | → | correct | **(Xu et al. EMNLP '11)** |
| writer style | ↔ | plain | **(Xu et al. COLING '12)** |
| complex | → | simple | **(Xu et al. TACL '15)** <br> **(Xu et al. TACL '16)** |
| noisy | ↔ | standard | **(Xu et al. BUCC '13)** <br> **(Xu et al. TACL '14)** <br> **(Xu et al. SemEval '15)** <br> **(Lan, Qiu, He, Xu EMNLP ' 17)** |
| feminine | ↔ | masculine | **(Preotiuc, Xu, Ungar AAAI '16)** |

# Common Pitfalls and Lessons I learned

**Lesson #1: Don't take data quality for granted.**

# Text Simplification Dataset

Parallel Simple-Normal Wikipedia Corpus was Benchmark since 2010.

Wei Xu, Chris Callison-Burch, Courtney Napoles. "Problems in Current Text Simplification Research: New Data Can Help" (TACL 2015)

# Text Simplification Dataset

**Pitfall**:  Simple Wikipedia is not simple.



Wei Xu, Chris Callison-Burch, Courtney Napoles. "Problems in Current Text Simplification Research: New Data Can Help"  (TACL 2015)

# Newsela Corpus

We proposed a new high-quality dataset of news articles simplified by professional editors.

Wei Xu, Chris Callison-Burch, Courtney Napoles. "Problems in Current Text Simplification Research: New Data Can Help" (TACL 2015)

# Wikipedia

# Newsela

**alignment error**

**real simplification**

**17%**

**50%**

**33%**

**not simpler**

**alignment error**

**not simpler**

**real simplification**

**2%**

**6%**

**92%**

manual inspection of aligned sentence pairs

see other analyses in the paper and EMNLP '15 talk (https://vimeo.com/150290363)

**Lesson #2: Evaluation is difficult (yet very important).**

# Human Evaluation

**Pitfall**: common practice is not necessarily good/correct

5-point Likert scale

- grammaticality

- meaning preservation

# Sentence Compression

Compression rate (CR) strongly correlates wit human judgements of meaning and grammaticality



Courtney Napoles, Benjamin Van Durme, and Chris Callison-Burch. "Evaluating Sentence Compression: Pitfalls and Suggested Remedies"  (T2T-G 2011)

# Sentence Simplification

**Pitfall**: Unfairly bias towards deletion over paraphrasing

5-point Likert scale



- 🟣 **grammaticality**
- 🔴 **meaning preservation**
- 🟢 **simplicity**

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

# Automatic Evaluation

**Pitfall:** outputs with no change get high **BLEU** and high meaning and grammaticality

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

# Automatic Evaluation

**SARI tunable metric:** Compare **S**ystem output **a**gainst **R**eference sentences and against the **I**nput sentence.



*keep*
O ∩ R ∩ I

*del*
I ∩ not O ∩ not R

Input
I

*add*
O ∩ R ∩ not I

System output
O

Human references
R

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

# Automatic Evaluation

**SARI** <span style="color:red">**tunable**</span> **metric:** Compare **S**ystem output **a**gainst **R**eference sentences and against the **I**nput sentence.



*keep*
$O \cap R \cap I$

*del*
$I \cap$ not $O \cap$ not $R$

Input
I

System
output
O

Human
references
R

*add*
$O \cap R \cap$ not $I$

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

# Automatic Simplification

Legend: Grammar (blue), Meaning (orange), Simplicity+ (red)

Categories: Simple Wikipedia, Turk, Moses + reranking (Wubben et al. 2012), Joshua -BLEU-, Joshua -SARI-

# Automatic Simplification



Legend: ■ Grammar (blue) ■ Meaning (orange) ■ Simplicity+ (red)

Categories: Simple Wikipedia, Turk, Moses + reranking (Wubben et al. 2012), Joshua -BLEU-, Joshua -SARI-

**tuning towards BLEU leaves input unchanged**

Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen, Chris Callison-Burch. "Optimizing Statistical Machine Translation for Text Simplification" (TACL 2016)

**Lesson #3: Linguistic styles often conflate w/ topics.**

# User Profiling



Delighted I kept my Xmas vouchers - Happy
Friday to me 😊 #shopping

# Paraphrase can help control topic and focus on style

wonderfully  delightfully  beautifully  fine  well  good  nicely  superbly

**she says**

**he says**

(also age & income)

Daniel Preotiuc, Wei Xu, Lyle Ungar. "Discovering User Attribute Stylistic Differences via Paraphrasing"  (AAAI  2016)

# Style is often subtle

**Good News**: "Wisdom of the Crowd"
aggregated human judgements are pretty good

Daniel Preotiuc, Wei Xu, Lyle Ungar. "Discovering User Attribute Stylistic Differences via Paraphrasing"  (AAAI  2016)

# Style is often subtle

**Good News**: "Wisdom of the Crowd"
aggregated human judgements are pretty good

**Bad News**: Text-to-Text Generation in famine or masculine
styles is much more difficult than Shakespearean.

**Lesson #4: Our language is ever evolving.**

Oxford Dictionaries revealed this week the earliest known usage of word "selfie" is from a 2002 online ABC forum post.

# My take on Evolving Language

**learn and model very-large-scale paraphrases**

|  | word |  |
|---|---|---|
| *selfie* |  | *photo* |

|  | phrase |  |
|---|---|---|
| *gets the boot from* |  | *has been sacked by* |

|  | sentence |  |
|---|---|---|
| *Mr Corbyn is actually a secret supporter of Brexit.* |  | *Jeremy Corbyn is a closest Brexiteer.* |

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

# Twitter is a powerful 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, Chris Callison-Burch, Bill Dolan, Yangfeng Ji. "Extracting Lexically Divergent Paraphrases from Twitter" (TACL 2014)
Wei Xu, Chris Callison-Burch, Bill Dolan. "SemEval-2015 Task 1: Paraphrase and Semantic Similarity in Twitter" (SemEval 2015)

# Twitter is a powerful resource

thousands of users
talk about both big/micro events daily



**Very diverse!**

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

Wei Xu, Alan Ritter, Chris Callison-Burch, Bill Dolan, Yangfeng Ji. "Extracting Lexically Divergent Paraphrases from Twitter" (TACL 2014)
Wei Xu, Chris Callison-Burch, Bill Dolan. "SemEval-2015 Task 1: Paraphrase and Semantic Similarity in Twitter" (SemEval 2015)

# Automatic Paraphrase Identification + Word Alignment

- Streaming data + Unsupervised model  (Xu et al. 2013)

- Topic detection + Multiple Instance Learning (Xu et al. 2014)

- URL linked data + Deep Pairwise Model (Lan et al. 2017)

- Ongoing work …

[Das & Smith 2014; Socher et al. 2011; Ling et al. 2013; Ji & Eisenstein 2013; Parikh et al. 2016; Witting & Gimpel 2017; and many others]

Wei Xu, Alan Ritter, Ralph Grishman. "Gathering and Generating Paraphrases from Twitter with Application to Normalization"  (BUCC 2013)
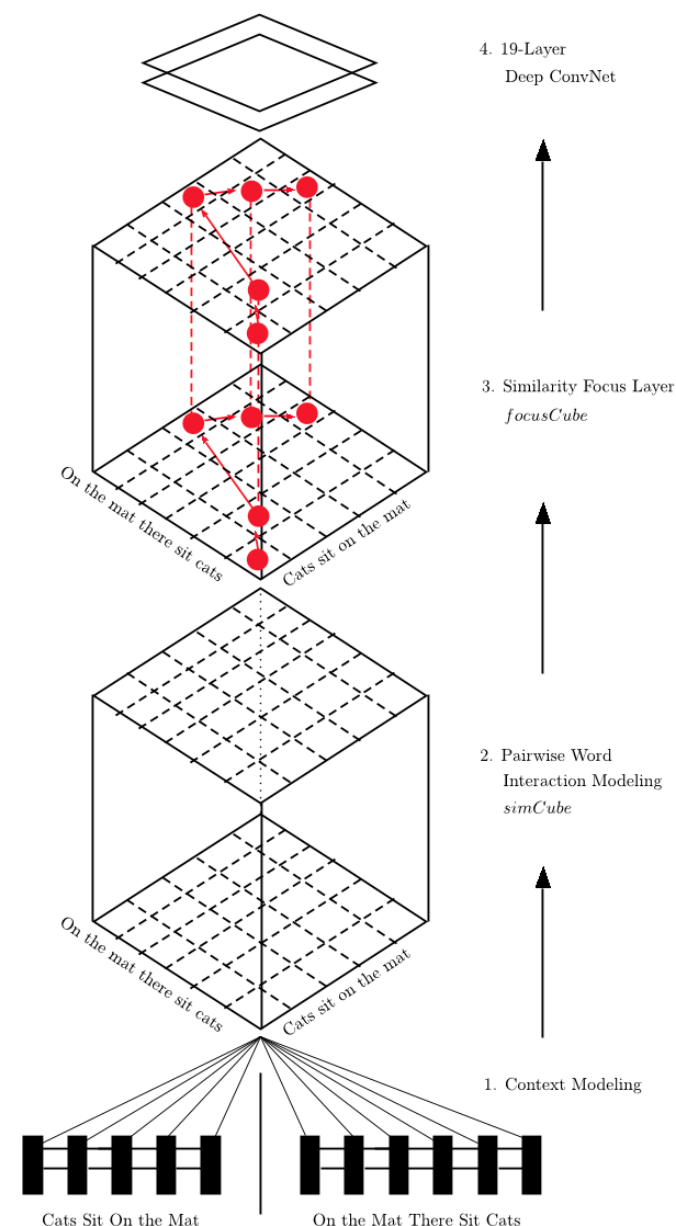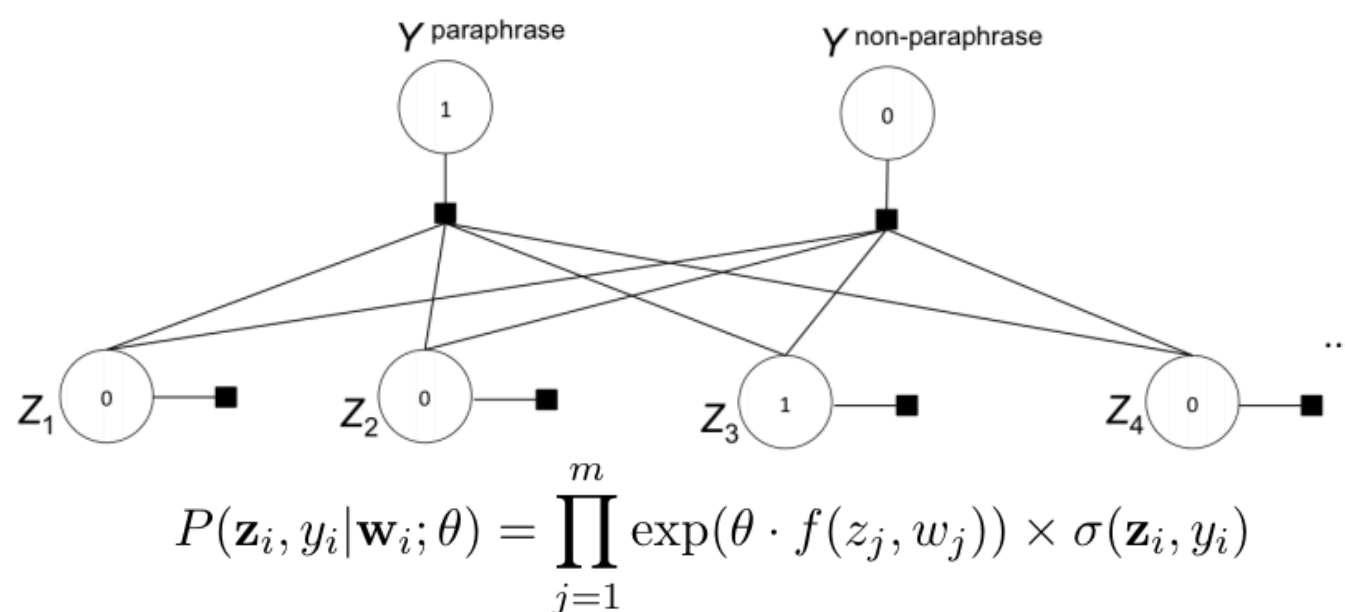
Wei Xu, Alan Ritter, Chris Callison-Burch, Bill Dolan, Yangfeng Ji. "Extracting Lexically Divergent Paraphrases from Twitter"  (TACL 2014)

Wuwei Lan, Siyu Qiu, Hua He, Wei Xu. "A Continuously Growing Dataset of Sentential Paraphrases"  (EMNLP  2017)

# Automatic Paraphrase Identification + Word Alignment



- **LEX-OrMF** (Orthogonal Matrix Factorization) [Guo and Diab 2012]
- **DeepPairwiseWord** (Deep Neural Networks) [He et al. 2015; Ongoing Work]
- **MultiP** (Multiple Instance Learning) [Xu et al. 2014; Ongoing Work]

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

Wuwei Lan, Siyu Qiu, Hua He, Wei Xu. "A Continuously Growing Dataset of Sentential Paraphrases"  (EMNLP 2017)
Wuwei Lan, Wei Xu. "A Better Pairwise Neural Model"  Ongoing Work

# [Twitter Paraphrase Corpus]

**51,524 sentence pairs
(manually annotated)**

**> 30,000 new sentential paraphrases
every month (automatically harvested)**

Wuwei Lan, Siyu Qiu, Hua He, Wei Xu. "A Continuously Growing Dataset of Sentential Paraphrases" (EMNLP 2017)

# Timely Paraphrases

Donald Trump, DJT, Drumpf, Mr Trump, Idiot Trump, Chump, Evil Donald, #OrangeHitler, Donald @realTrump, D*nald Tr*mp, Comrade #Trump, Crooked #Trump, CryBaby Trump, Daffy Trump, Donald KKKrump, Dumb Trump, GOPTrump, Incompetent Trump, He-Who-Must-Not-Be-Named, Pres-elect Trump, President-Elect Trump, President-elect Donald J . Trump, PEOTUS Trump, Emperor Trump

# Lesson #5: Data! Data! Data!

# Sentence Compression Dataset

**A Dataset and Evaluation Metrics for Abstractive Compression of Sentences and Short Paragraphs**

**Kristina Toutanova**
Microsoft Research
Redmond, WA, USA

**Chris Brockett**
Microsoft Research
Redmond, WA, USA

**Ke M. Tran**[*]
University of Amsterdam
Amsterdam, The Netherlands

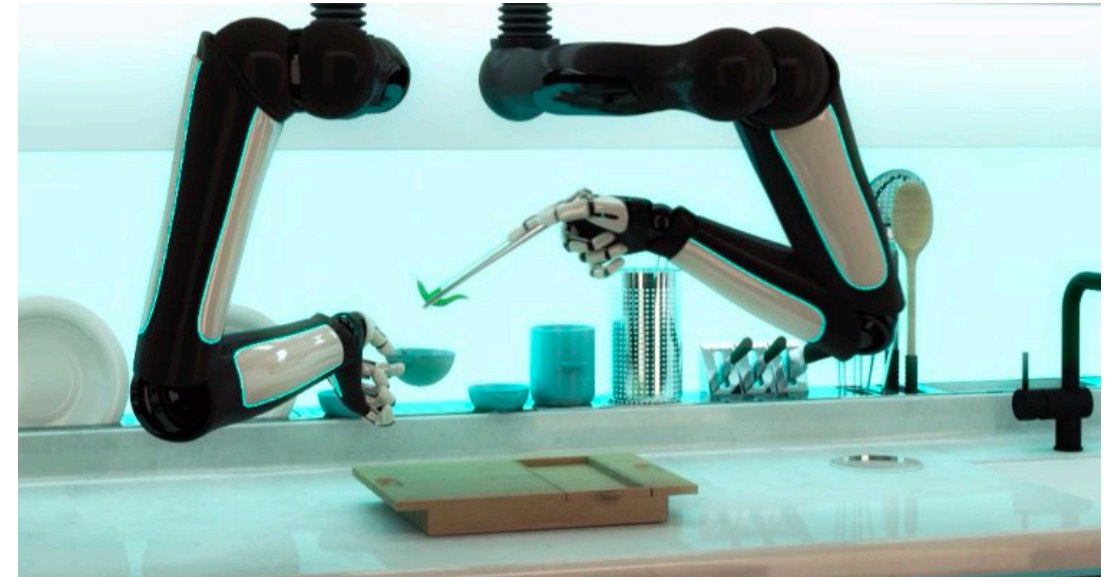**Saleema Amershi**
Microsoft Research
Redmond, WA, USA

## Abstract

We introduce a manually-created, multi-reference dataset for abstractive sentence and short paragraph compression. First, we examine the impact of single- and multi-sentence level editing operations on human compression quality as found in this corpus. We observe that substitution and rephrasing operations are more meaning preserving than other operations, and that compressing in context improves quality. Second, we systematically
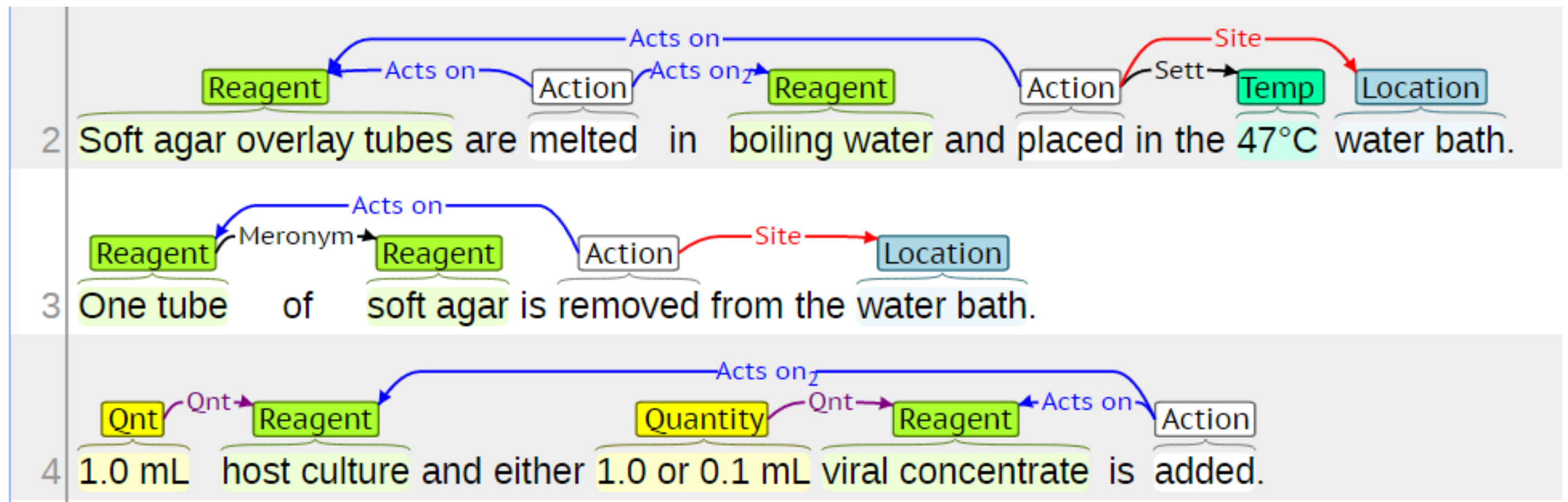
This paper has two parts. In the first half, we introduce a manually-created *multi-reference* dataset for *abstractive* compression of sentences and short paragraphs, with the following features:
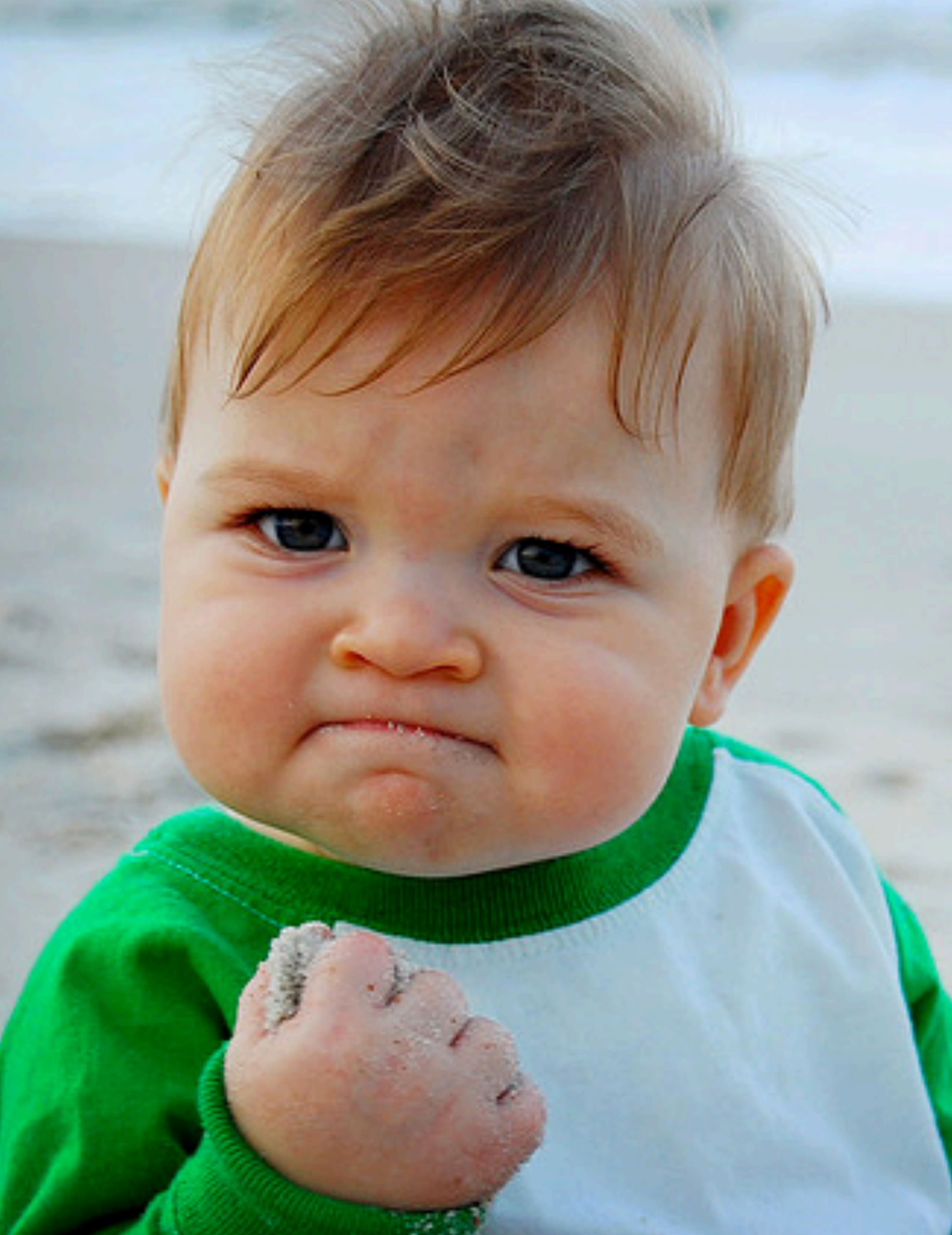
- It contains approximately 6,000 source texts with multiple compressions (about 26,000 pairs of source and compressed texts), representing business letters, newswire, journals, and technical documents sampled from the Open American National Corpus (OANC[1]).

# Instructional Language



Wet Lab Protocols:



Chaitanya Kulkarni, Wei Xu, Alan Ritter, Raghu Machiraju
"Taking the First Essential Steps in Automating the Wet Laboratory: Annotating a Corpus of Protocols for Reproducibility" (Ongoing Work)

We Can Do It!

as long as
we have data

thank u 4 ur time

thanku

# Thank You

thanking you

gramercies

thanks a lot

appreciate it

gratitude

thx

3x

tyvm

thanks

say thanks

thank you very much

thnx

thanks a ton

I can no other answer make but thanks,
And thanks, and ever thanks.

wawwww thankkkkkkkkkkk you alottttttttttttt!

I am grateful