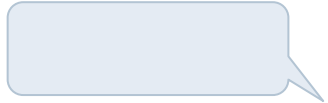# Conversations Gone Awry:

## *Detecting Early Signs of Conversational Failure*

**Justine Zhang**, Jonathan P. Chang, Cristian Danescu-Niculescu-Mizil, Lucas Dixon, Yiqing Hua, Dario Taraborelli, Nithum Thain
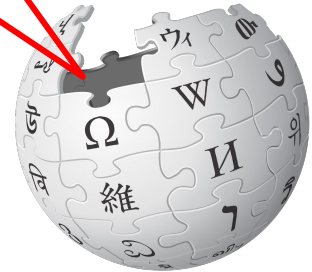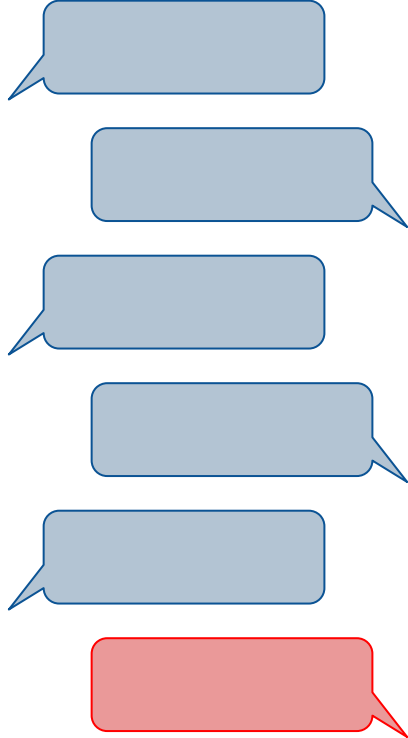
**Cornell University**, Jigsaw, Wikimedia Foundation

# Motivating problem: toxic behaviour

Yin et al., 2009
Sood et al., 2012
Nobata et al., 2016
Cheng et al., 2017
Gambäck & Sikdar, 2017
Pavlopoulos et al., 2017
Wulczyn et al., 2017

"Wow, you're coming off as a total dick...what the hell is wrong with you?"

2

# Testing intuitions of early signals: a guessing game

**A**

**B**



*"Wow, you're coming off as a total d\*\*k...what the hell is wrong with you?"*

4

**A**

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I would assume that it's as reliable as any other mainstream news source.

[...]

**B**

Why's there no mention of it here? Namely, an altercation with a foreign intelligence group? True, by *some* standards the source has some weak points, but that doesn't mean it shouldn't exist.

So what you're saying is we should put a bad source in the article because it exists?

[...]

*"Wow, you're coming off as a total d\*\*k…what the hell is wrong with you?"*

**Other humans: 72% Accuracy**

**Can we reconstruct some of this intuition?**

# Identifying awry-turning conversations

**Civil Start:**



← *Toxic Comment*

# Identifying awry-turning conversations

**Civil Start:**

**(no *toxic* behaviour)**

*rude, insulting, or disrespectful towards a person/group or that person/group's actions, comments, or work*

Wulczyn et al., 2017

← *Toxic Comment*

# Identifying awry-turning conversations

**Civil Start:**

*(one of the initial commenters eventually attacks the other)*

Arazy et al., 2013

← *Personal Attack*

**Civil Start:** → *early signals*

← *Personal Attack*

**Civil Start:**

Coser, 1956
De Dreu & Weingart, 2003
Galley et al., 2004
Andreas et al., 2012
Hillard et al., 2012
Allen et al., 2014
Wang & Cardie, 2014
Rosenthal & McKeown, 2015
*inter alia*

*(different from disagreement)*

← *Personal Attack*

# Data: Wikipedia Talk Pages

*(same **talk page**)*



**635 pairs**

*(details in paper; data available!)*

# Detecting early signals: *(im)politeness*

**A**

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I would assume that it's as reliable as any other mainstream news source.

**B**

Why's there no mention of it here? Namely, an altercation with a foreign intelligence group? True, by *some* standards the source has some weak points, but that doesn't mean it shouldn't exist.

So what you're saying is we should put a bad source in the article because it exists?

*direct questioning*
*hedging*

# Detecting early signals: *(im)politeness*

**A**

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I would assume that it's as reliable as any other mainstream news source.

**B**

Why's there no mention of it here? Namely, an altercation with a foreign intelligence group? True, by *some* standards the source has some weak points, but that doesn't mean it shouldn't exist.

So what you're saying is we should put a bad source in the article because it exists?

Goffman, 1955
Fraser, 1980
Brown & Levinson, 1987

*"building rapport, softening potential conflicts"*

16

# Detecting early signals: *(im)politeness*

A
> maybe
> don't think
>
> seems

> would assume

B
> Why's there no mention of it here? Namely, an altercation with a foreign intelligence group?

> So what you're saying is
>
> ?

*38 politeness strategies features*

Goffman, 1955
Fraser, 1980
Brown & Levinson, 1987

D-N-M et al., 2013
politeness.cornell.edu

# Detecting early signals

> Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

> I have been working on creating a new section…

> Let me know if you agree with this…

## Domain-specific intuition: *"let's coordinate"*

Kittur & Kraut, 2008

# Detecting early signals: *rhetorical intentions*

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I have been working on creating a new section…

Let me know if you agree with this…

**Automatically inferring**
*"let's coordinate"*

*idea:*
*similar rhetorical intentions prompt*
*similar responses*

Zhang, Spirling & D-N-M, 2017

# Detecting early signals: *rhetorical intentions*

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I have been working on creating a new section...

Let me know if you agree with this...

**typical responses**

Nice work, thanks for helping.

If you can do it I would appreciate it.

OK, I'll take a look later.

*idea:*
*similar rhetorical intentions prompt similar responses*

Zhang, Spirling & D-N-M, 2017

# Detecting early signals: *rhetorical intentions*

Is the St. Petersberg Times considered a reliable source by Wikipedia? I'm going to maybe do a rewrite of the article. I don't think the bulk of this article should rely on this one so-so source, which seems to speculate about UFOs.

I have been working on creating a new section…

Let me know if you agree with this…

**typical responses**

Nice work, thanks for helping.

If you can do it I would appreciate it.

OK, I'll take a look later.

*idea:*
*similar rhetorical intentions prompt*
*similar responses*

Zhang, Spirling & D-N-M, 2017

# Detecting early signals: *rhetorical intentions*

**typical** **responses**

going to

do a rewrite

have been working

Let me know

*Coordination*

Nice work, thanks for helping.

If you can do it I would appreciate it.

OK, I'll take a look later.

# Detecting early signals: *rhetorical intentions*

**Coordination**
Moderation
**Factual Check**
Casual remark
Action statement
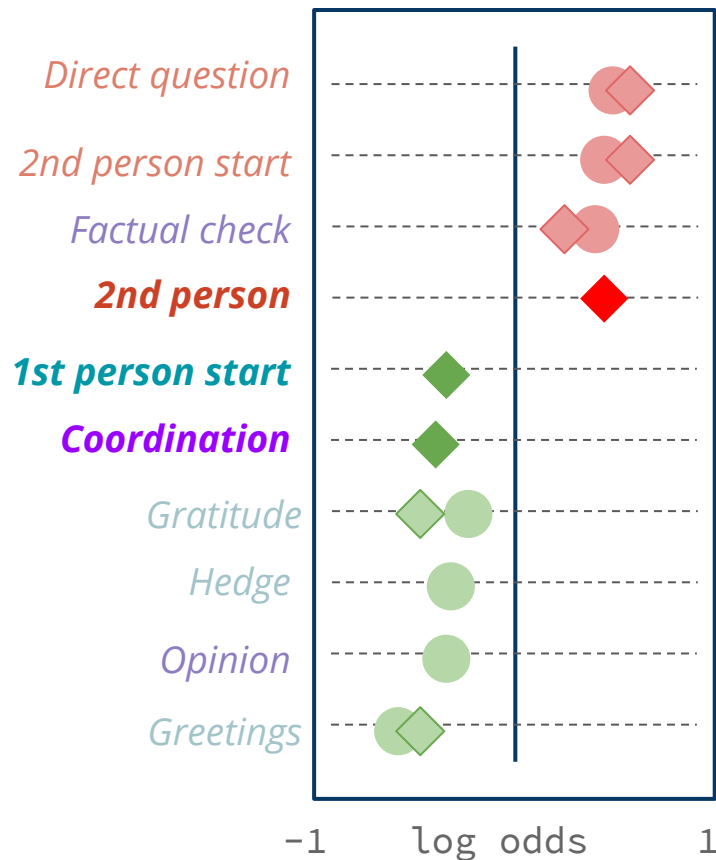**Opinion**

I'm going to maybe do a rewrite of the article...

The census is not talking about families here.

It is hard to address both of these issues.

*12 rhetorical intention features*
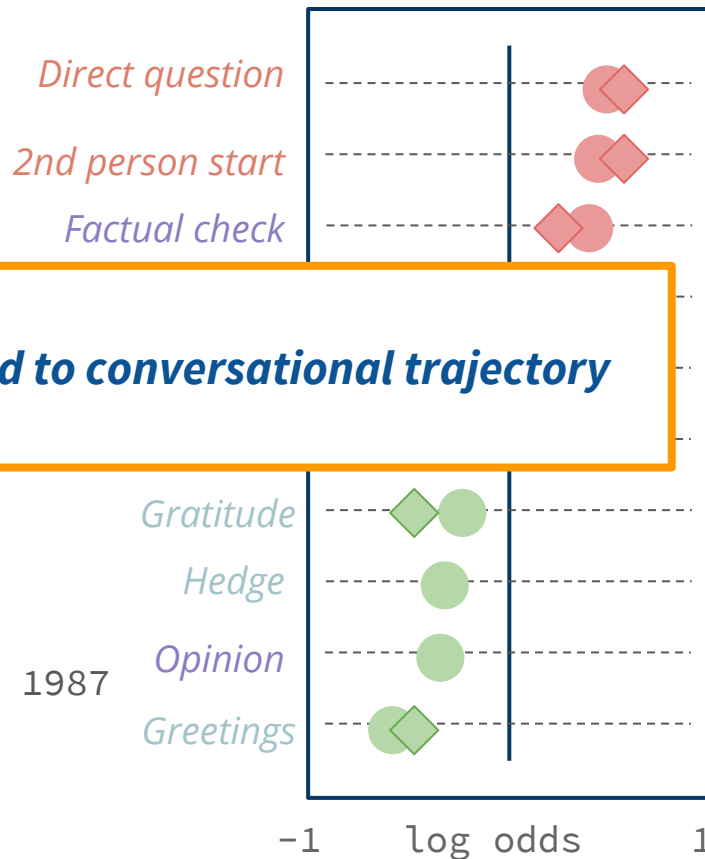
# Comparing awry-turning and on-track conversations



occurrence in
● **first comment**

◆ **second comment**

binomial $p < 0.05$

More likely to occur in
awry-turning conversation

# Comparing awry-turning and on-track conversations



Direct question
2nd person start
Factual check

**politeness is tied to conversational trajectory**

Gratitude
Hedge
Opinion
Greetings

Goffman, 1955
Fraser, 1980
Brown & Levinson, 1987

-1    log odds    1

*occurrence in*
● *first comment*

◆ *second comment*

*More likely to occur in awry-turning conversation*

# Measuring the strength of our early signals



A

B

How well do our signals
(politeness strategies,
rhetorical intentions)
perform at the guessing game?

[...]          [...]

# Measuring the strength of our early signals

**Humans: 72%**

**Our extracted signals: 65%**

**Can we reconstruct human intuition?**

**80% accuracy over cases that humans got right**

BOW/sentiment: 57%

random: 50%

100%

# Conclusion

**Conversations are <span style="color:orange">dynamic</span>:**
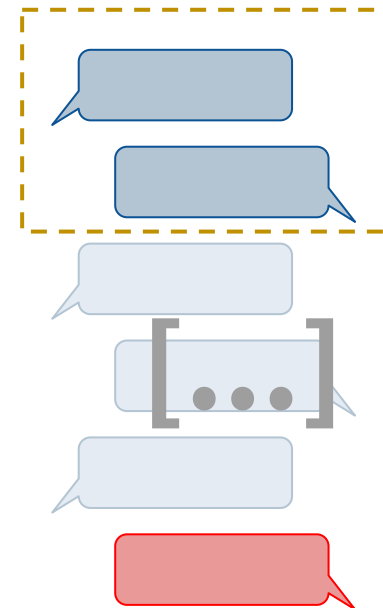**They can start civil…**
**but later turn <span style="color:darkred">awry</span>.**

**We detect <span style="color:goldenrod">early signals</span>**
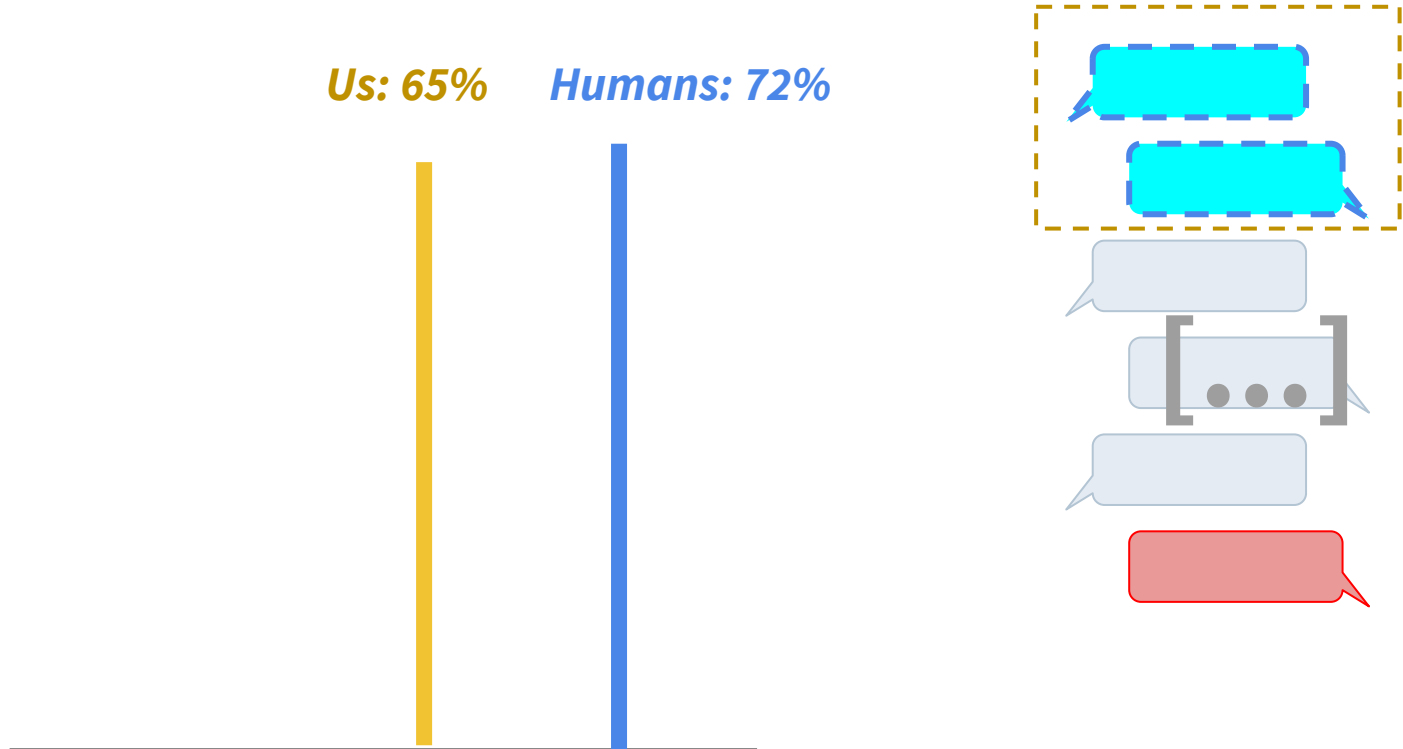**(<span style="color:teal">politeness strategies</span>,**
**<span style="color:purple">rhetorical intentions</span>)**
**of <span style="color:darkred">eventual derailment</span>,**

*towards* **the level of**      *Us: 65%*      *Humans: 72%*
**human intuition…**
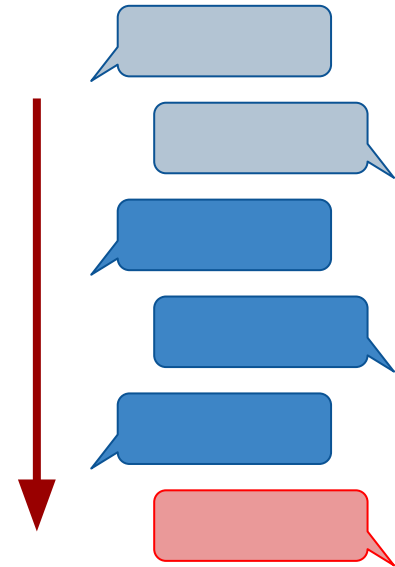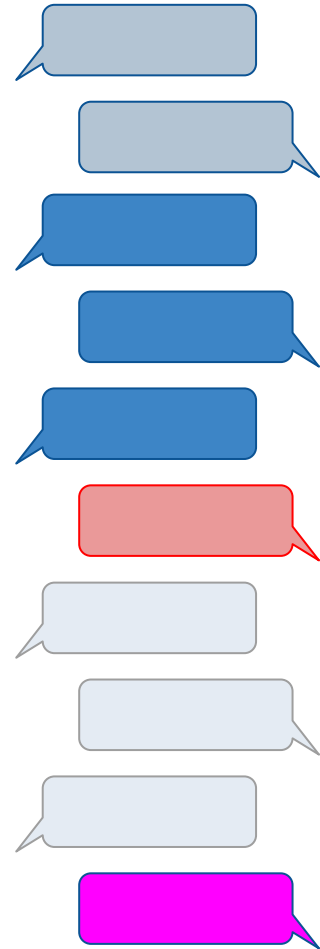
# Future work: early signals we might have missed

# Future work: understanding the derailment *process*

# Future work: other trajectories

**how might derailed conversations recover?**

# Questions?

**Data and code:**
**Cornell Conversational Analysis Toolkit**
convokit.infosci.cornell.edu

**Online guessing game:**
awry.infosci.cornell.edu