

A photograph of a cityscape across a body of water, likely Reykjavik, Iceland. The buildings are multi-story and have various colors, including red roofs and white walls. A church with a green steeple is visible on the right. The text 'LREC 2014' is written in a large, white, sans-serif font, and 'Reykjavik' is written below it in a white, cursive script.

LREC 2014  
Reykjavik

# LQVSumm:

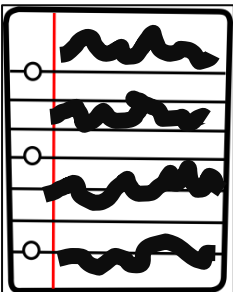
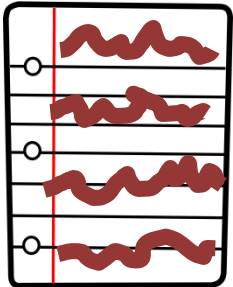
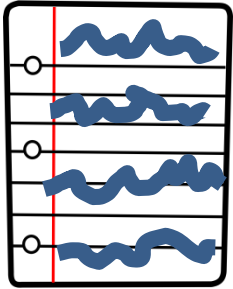
## A Corpus of Linguistic Quality Violations in Multi-Document Summarization

Annemarie Friedrich, Marina Valeeva and Alexis Palmer

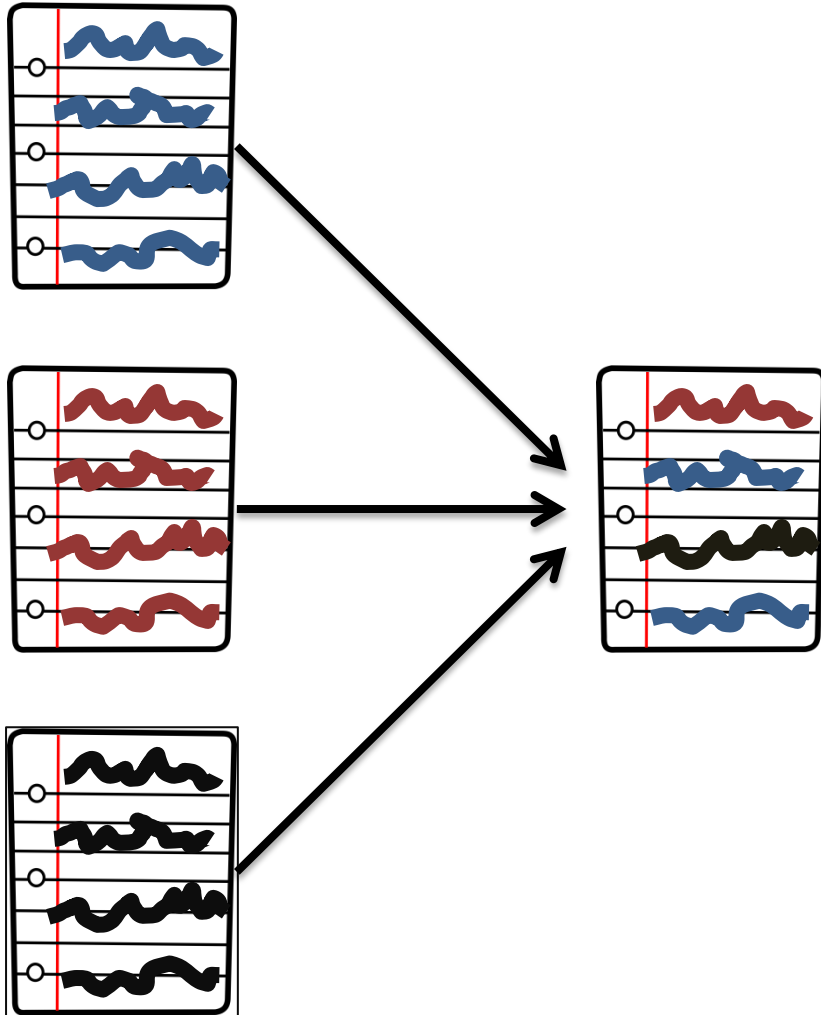


COMPUTATIONAL LINGUISTICS & PHONETICS  
SAARLAND UNIVERSITY, GERMANY

# Extractive Multi-Document Summarization



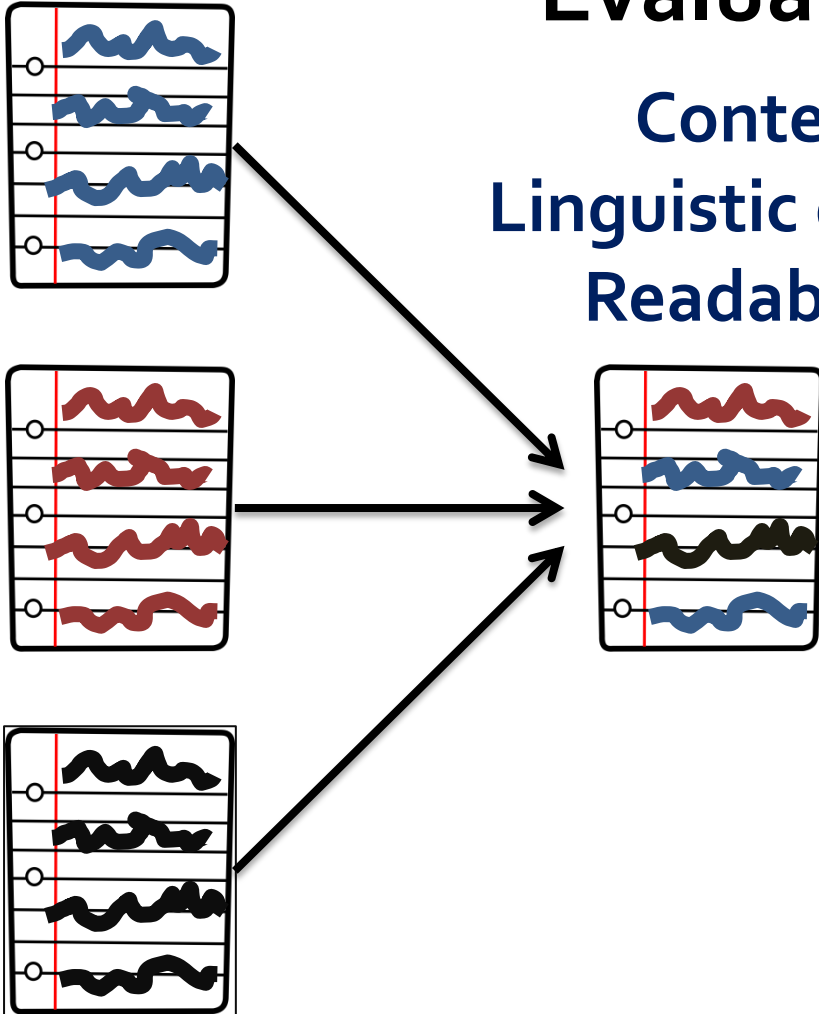
# Extractive Multi-Document Summarization



# Extractive Multi-Document Summarization

## Evaluation

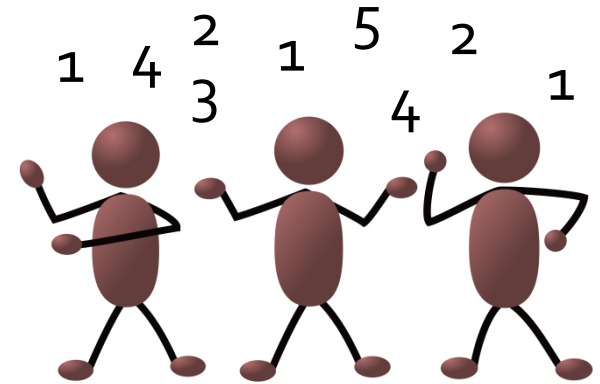
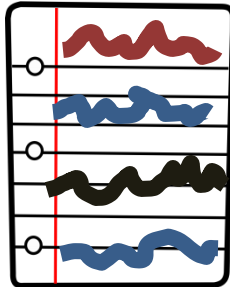
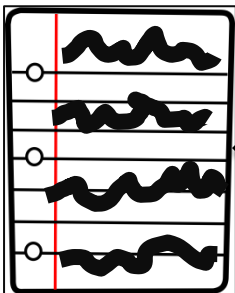
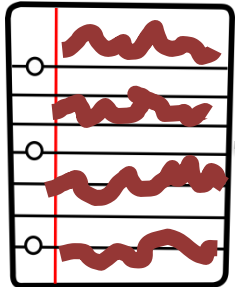
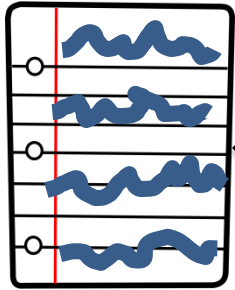
Content?  
Linguistic quality /  
Readability?



# Extractive Multi-Document Summarization

## Evaluation

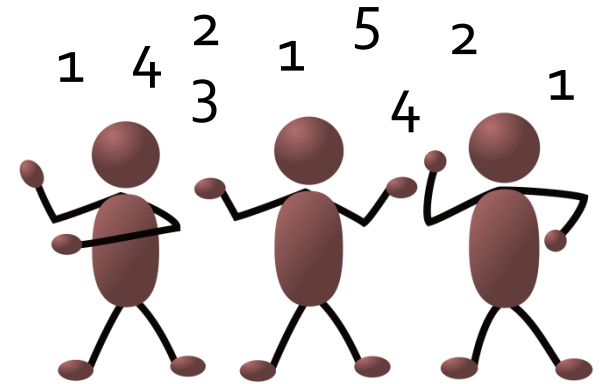
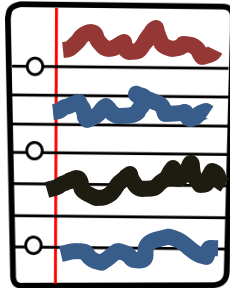
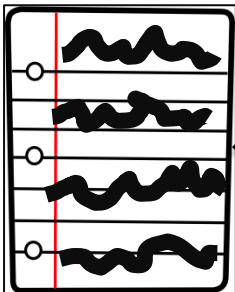
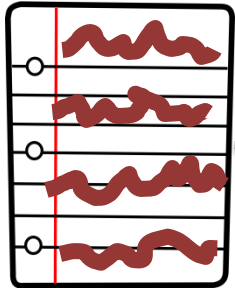
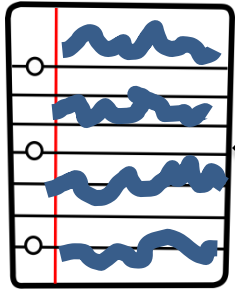
Content?  
Linguistic quality /  
Readability?



# Extractive Multi-Document Summarization

## Evaluation

Content?  
Linguistic quality /  
Readability?

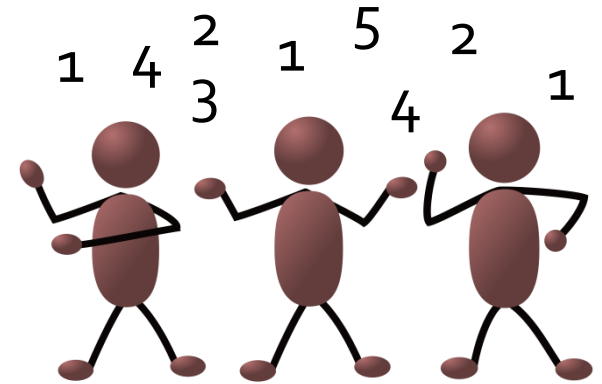
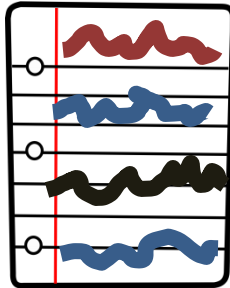
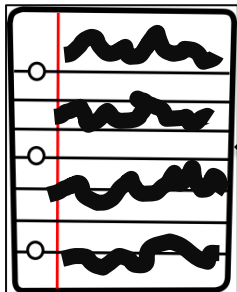
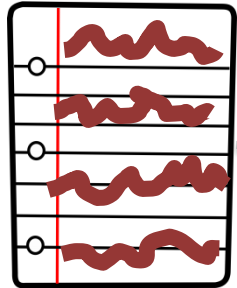
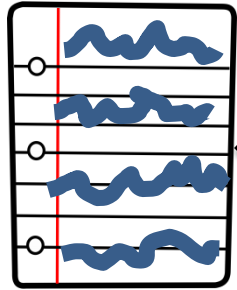


Automatic Evaluation Methods

# Extractive Multi-Document Summarization

## Evaluation

Content?  
Linguistic quality /  
Readability?



Automatic Evaluation Methods

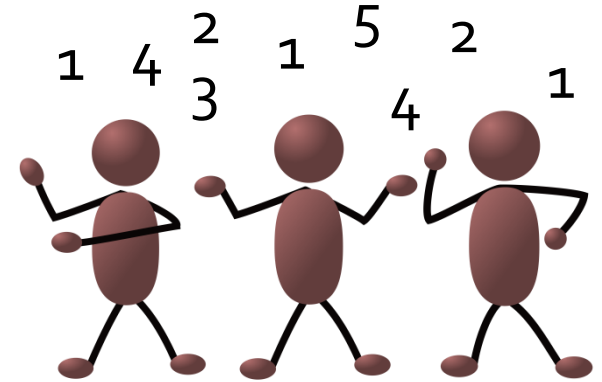
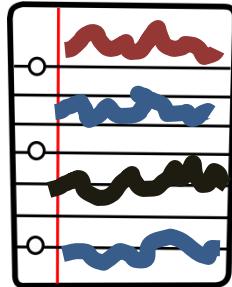
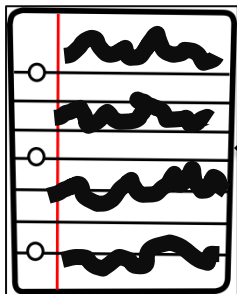
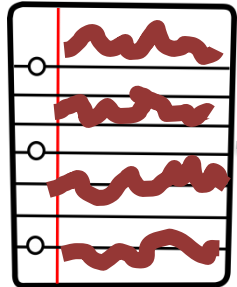
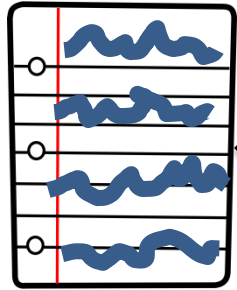
Automatic Content Evaluation



# Extractive Multi-Document Summarization

## Evaluation

Content?  
Linguistic quality /  
Readability?



Automatic Evaluation Methods

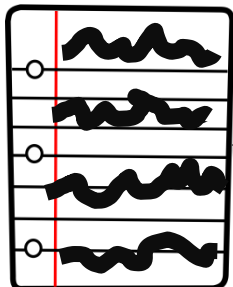
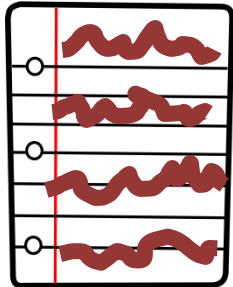
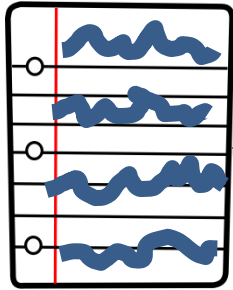
Automatic Content Evaluation



Automatic Linguistic Quality Evaluation ?



# Violations of Linguistic Quality



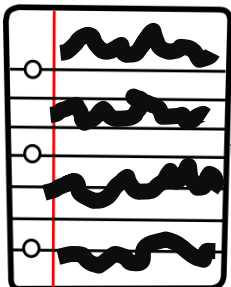
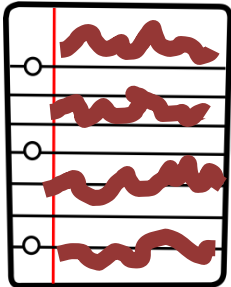
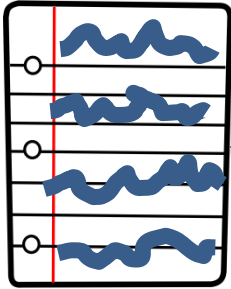
entity mentions:  
reference unclear

**The suspect** apparently called **her** from a cell phone shortly before the shooting began, saying he was “**acting out in revenge for something that happened 20 years ago**”, Miller said. **The gunman, a local truck driver Charles Roberts**, was apparently **acting in “revenge” for an incident that happened to him 20 years ago**. Charles Carl Roberts IV may have planned to

# Violations of Linguistic Quality

subsequent mention of  
entity too specific

entity mentions:  
reference unclear

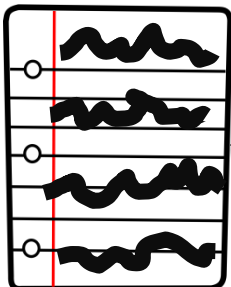
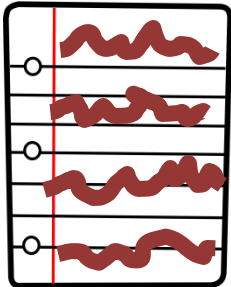
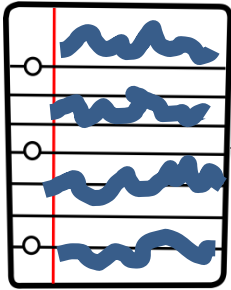


**The suspect** apparently called **her** from a cell phone shortly before the shooting began, saying he was “**acting out in revenge for something that happened 20 years ago**, Miller said. **The gunman, a local truck driver Charles Roberts**, was apparently **acting in “revenge” for an incident that happened to him 20 years ago**. Charles Carl Roberts IV may have planned to

# Violations of Linguistic Quality

subsequent mention of  
entity too specific

entity mentions:  
reference unclear



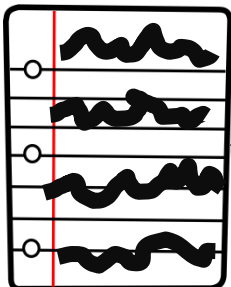
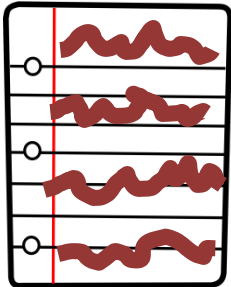
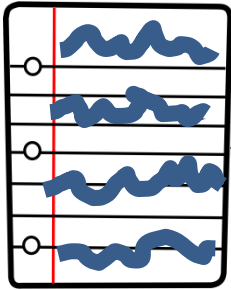
**The suspect** apparently called **her** from a cell phone shortly before the shooting began, saying he was “**acting out in revenge for something that happened 20 years ago**, Miller said. **The gunman, a local truck driver Charles Roberts**, was apparently **acting in “revenge” for an incident that happened to him 20 years ago**. Charles Carl Roberts IV may have planned to

redundant  
information

# Violations of Linguistic Quality

subsequent mention of entity too specific

entity mentions: reference unclear

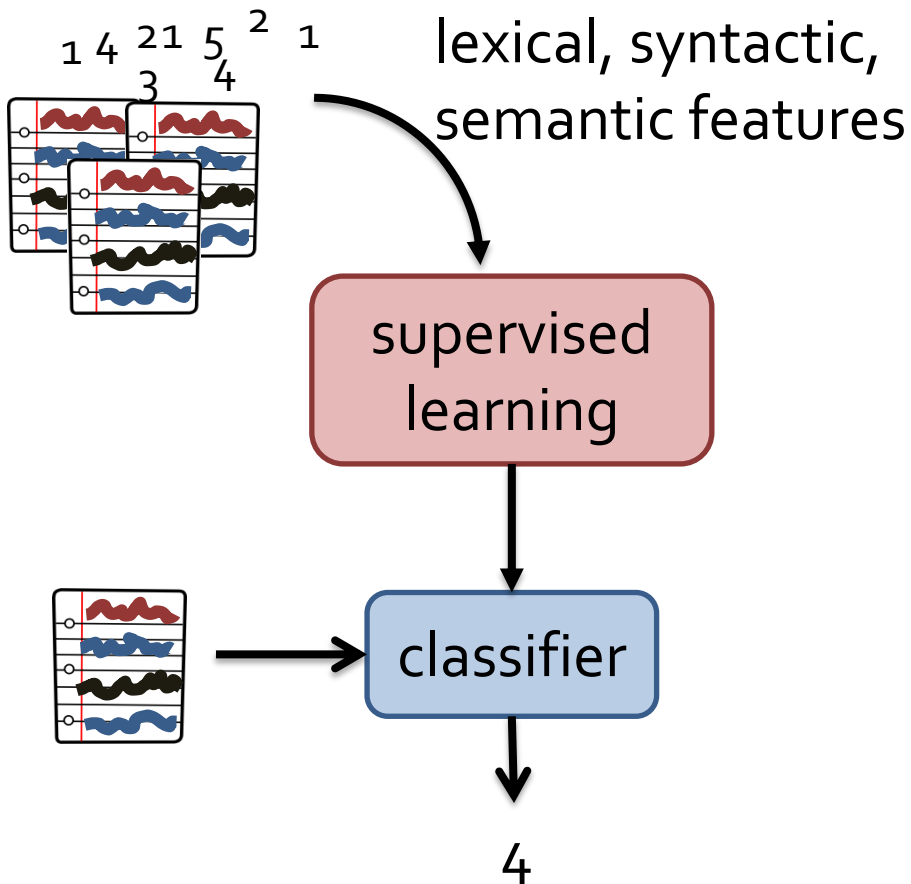


**The suspect** apparently called **her** from a cell phone shortly before the shooting began, saying he was “**acting out in revenge for something that happened 20 years ago**, Miller said. **The gunman, a local truck driver Charles Roberts**, was apparently **acting in “revenge” for an incident that happened to him 20 years ago**. Charles Carl Roberts IV may have planned to

redundant information

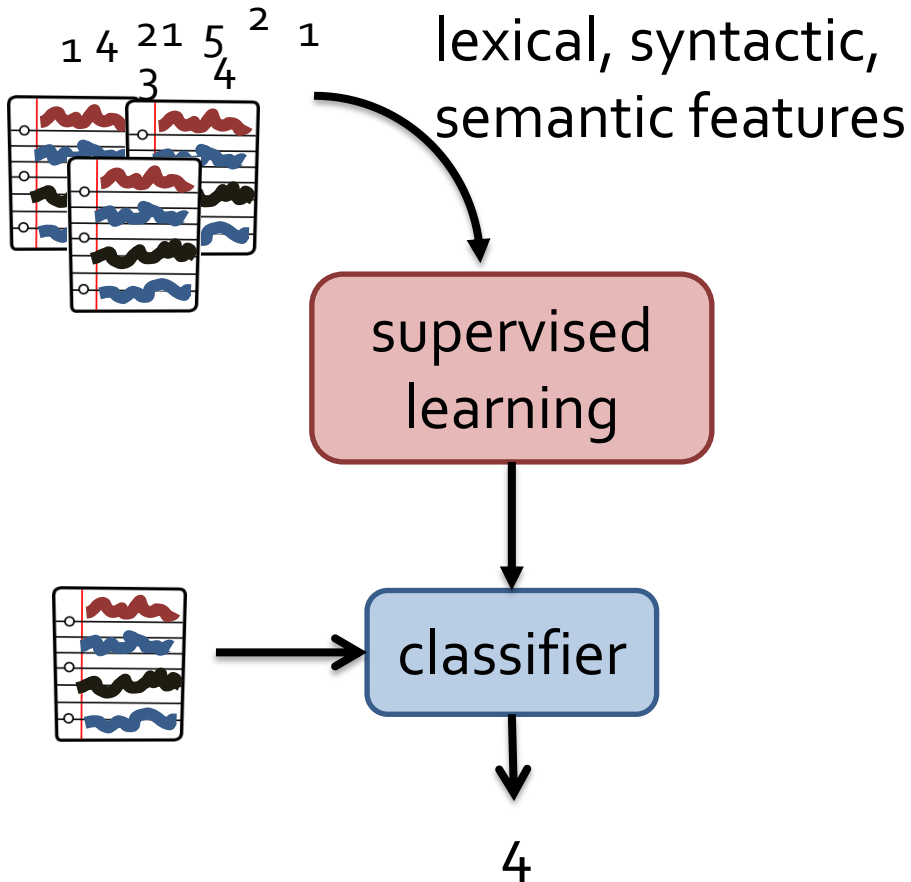
incomplete sentence

# Automatic Evaluation of Linguistic Quality for Automatic Summarization



[Pitler et al., 2010; Conroy et al., 2011;  
Giannakopoulos and Karkaletsis, 2011;  
de Oliveira, 2011; Lin et al., 2012]

# Automatic Evaluation of Linguistic Quality for Automatic Summarization



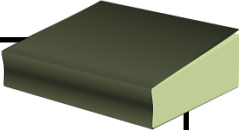
## Revision-based approach



[Pitler et al., 2010; Conroy et al., 2011; Giannakopoulos and Karkaletsis, 2011; de Oliveira, 2011; Lin et al., 2012]

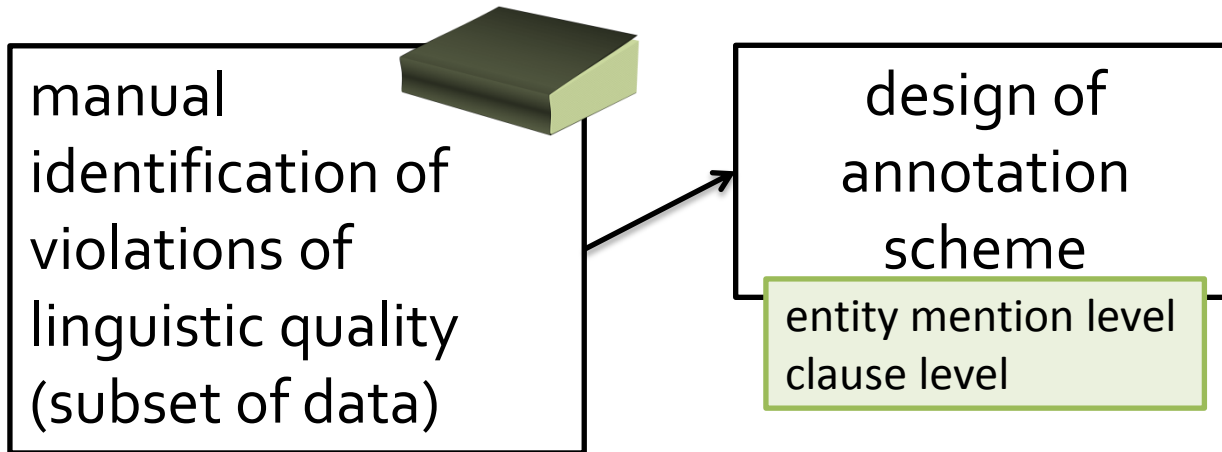
[Mani et al. 1999, Jing & McKeown 2000, Otterbacher et al. 2002]

# LQVSumm corpus



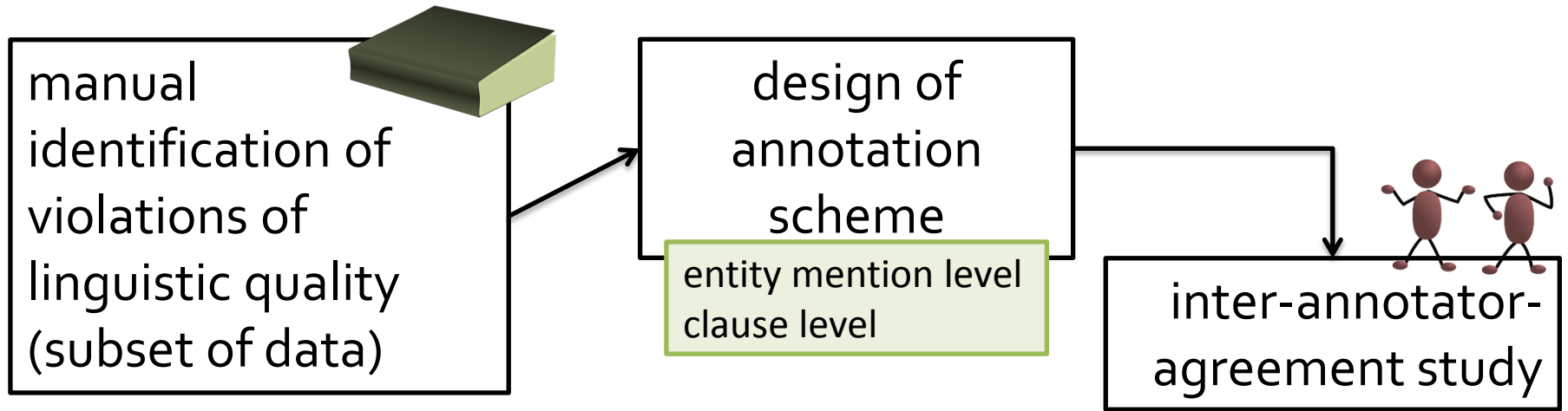
manual  
identification of  
violations of  
linguistic quality  
(subset of data)

# LQVSumm corpus

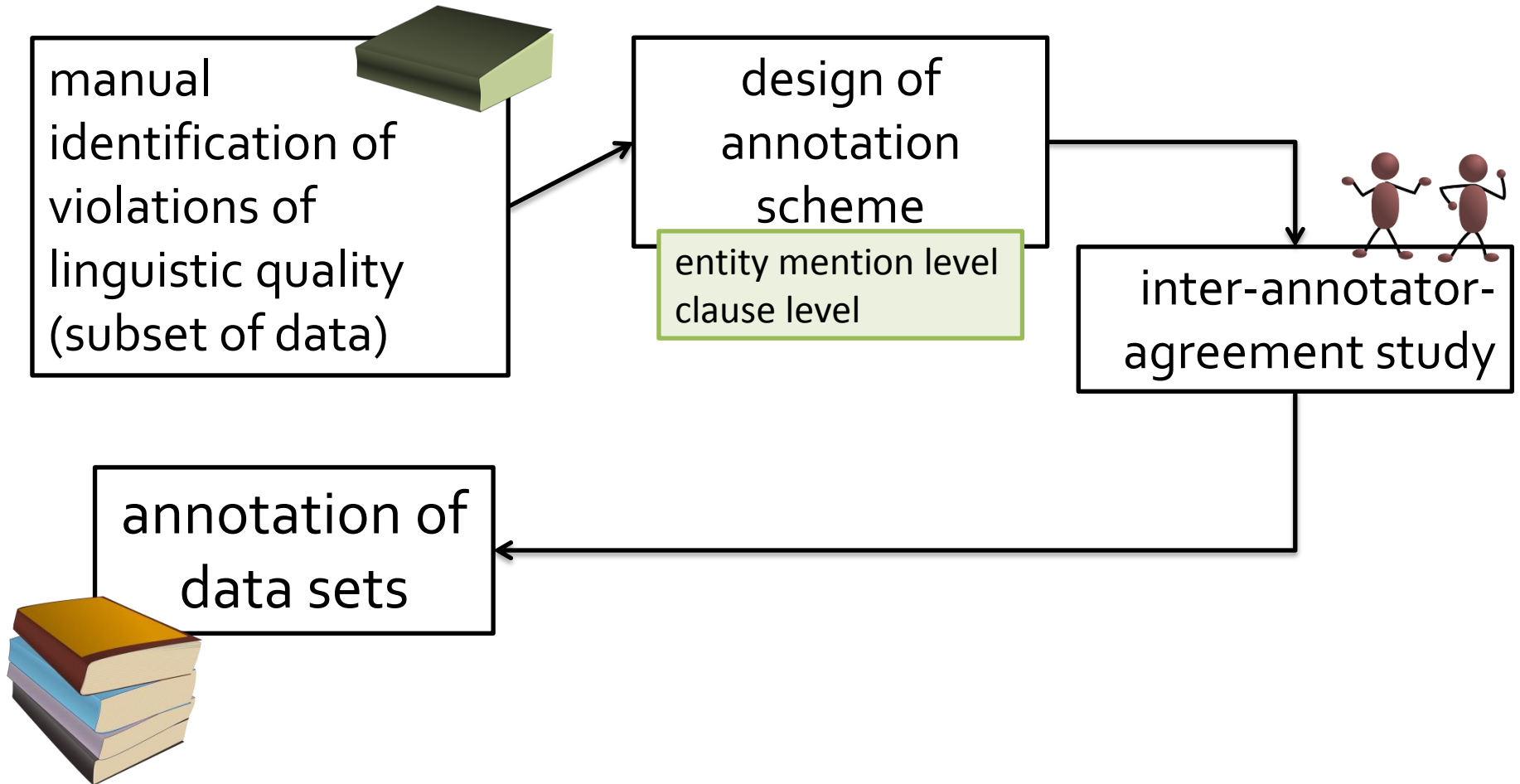




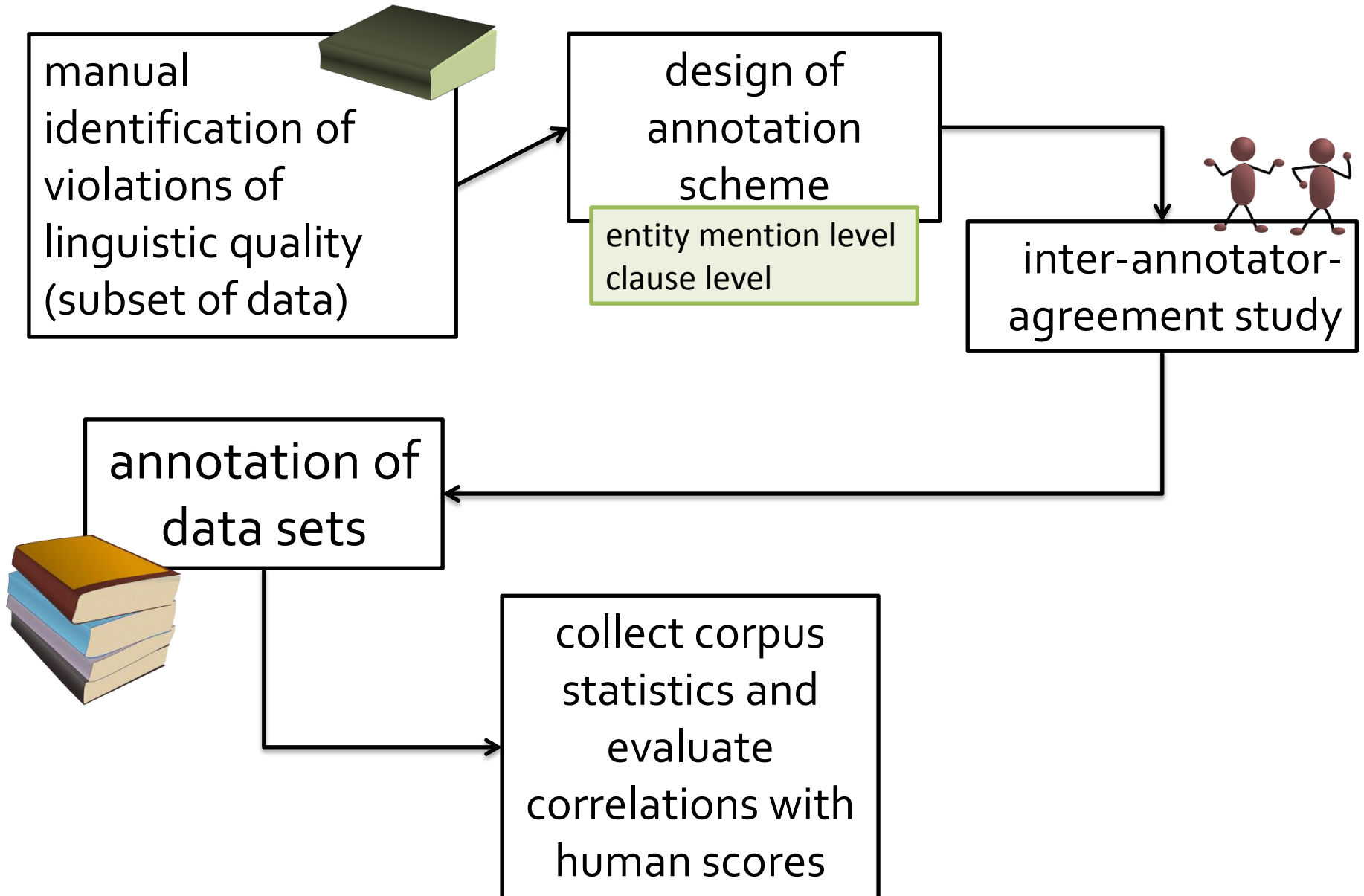
# LQVSumm corpus



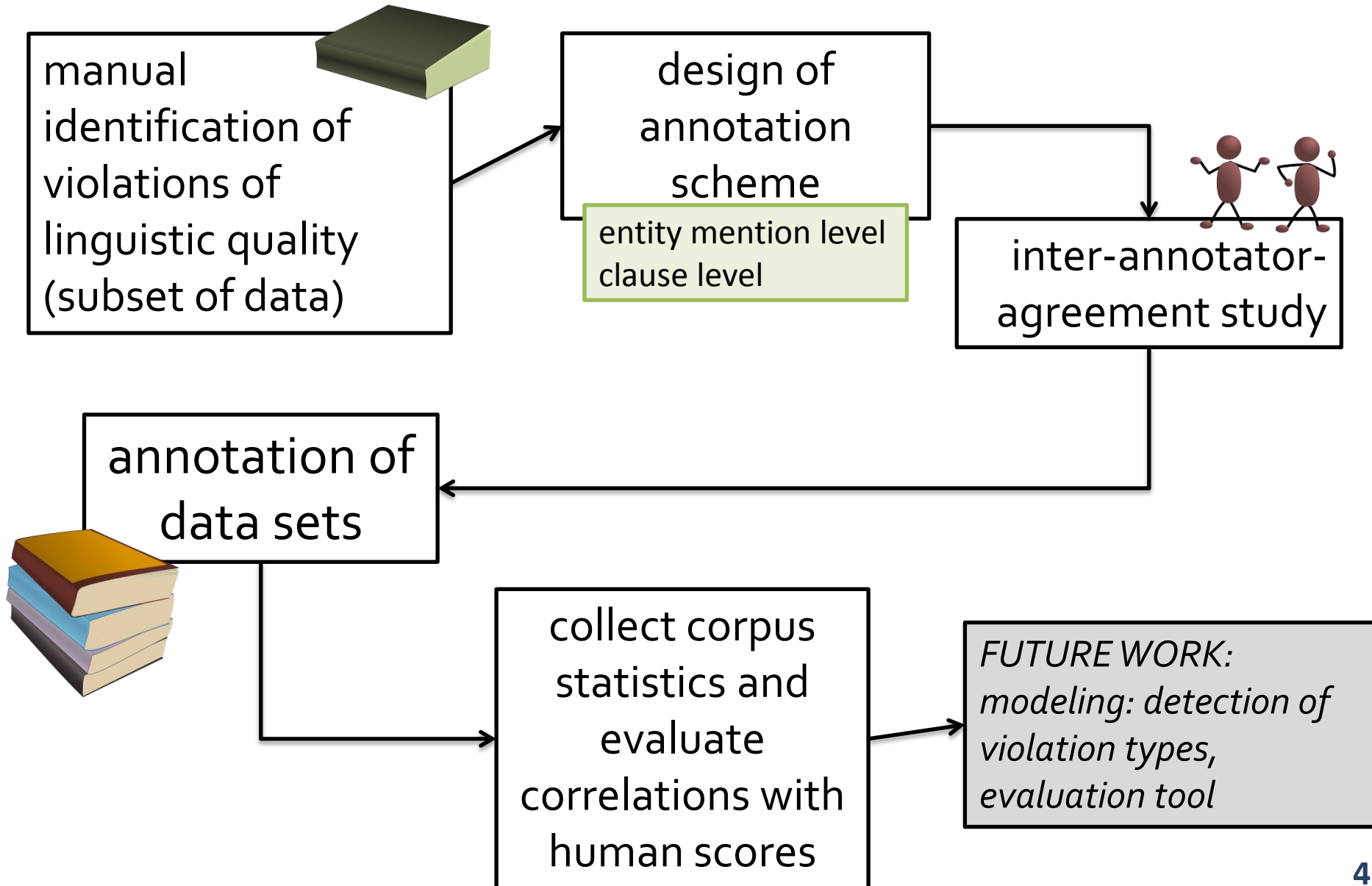
# LQVSumm corpus



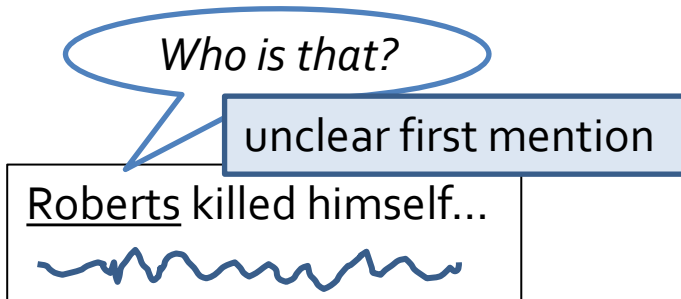
# LQVSumm corpus



# LQVSumm corpus



# Annotation Scheme: **Entity Mention level**



# Annotation Scheme: **Entity Mention level**

*Who is that?*

unclear first mention

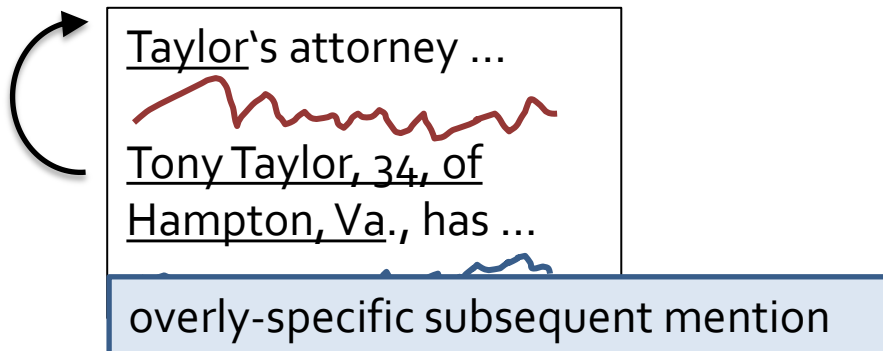
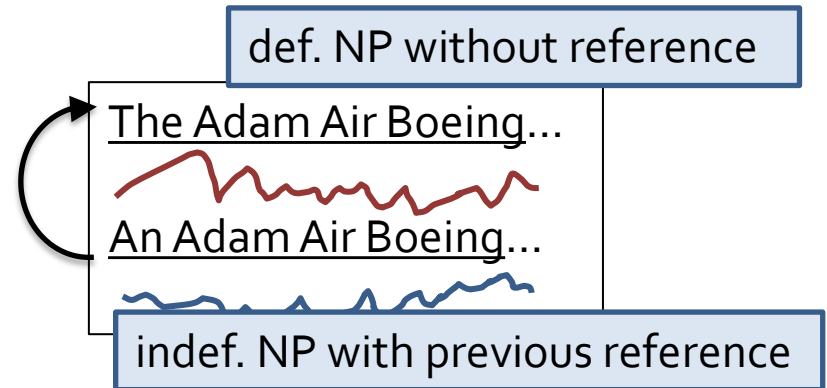
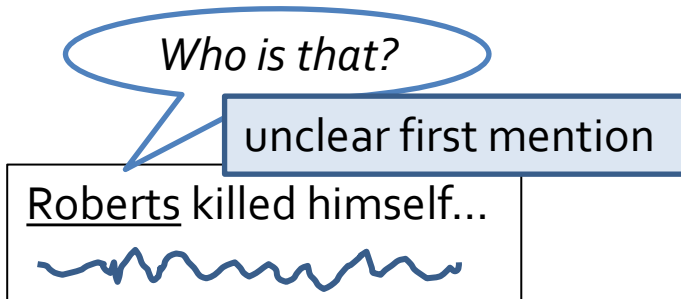
Roberts killed himself...

Taylor's attorney ...

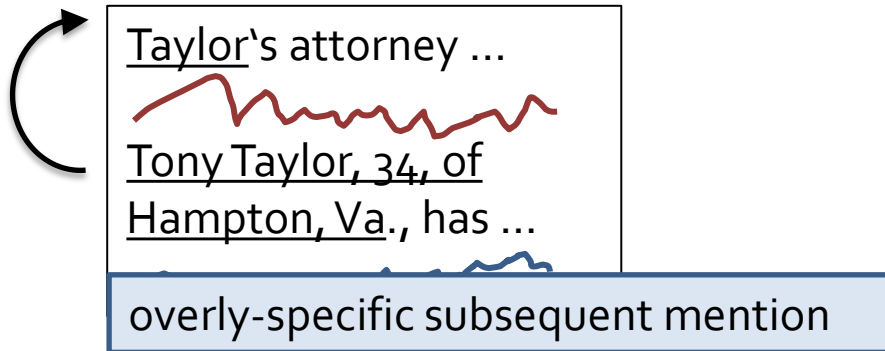
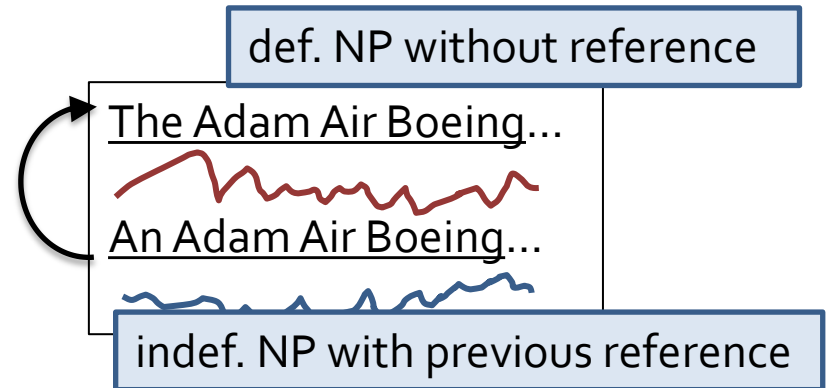
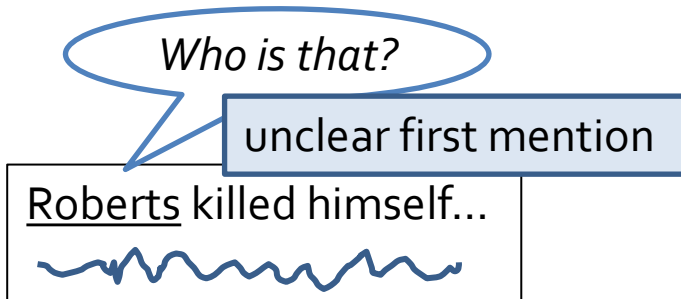
Tony Taylor, 34, of Hampton, Va., has ...

overly-specific subsequent mention

# Annotation Scheme: Entity Mention level



# Annotation Scheme: Entity Mention level



pronouns without antecedents

pronouns with misleading antecedents

unclear acronyms



# Annotation Scheme: **Clause level**

(sentence, phrase, sequence of tokens)

ungrammaticality

incomplete sentence

# Annotation Scheme: **Clause level**

(sentence, phrase, sequence of tokens)

ungrammaticality


incomplete sentence

dateline included



GEORGETOWN, Pennsylvania

2006-10-05 16:53:53 UTC



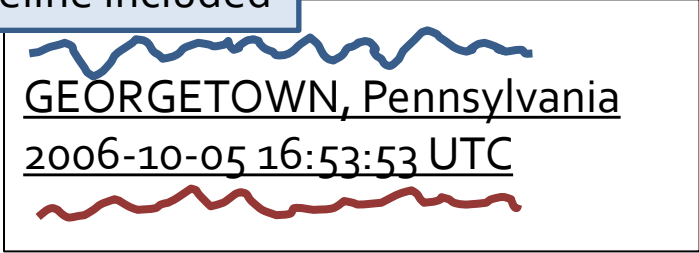
# Annotation Scheme: Clause level

(sentence, phrase, sequence of tokens)

ungrammaticality

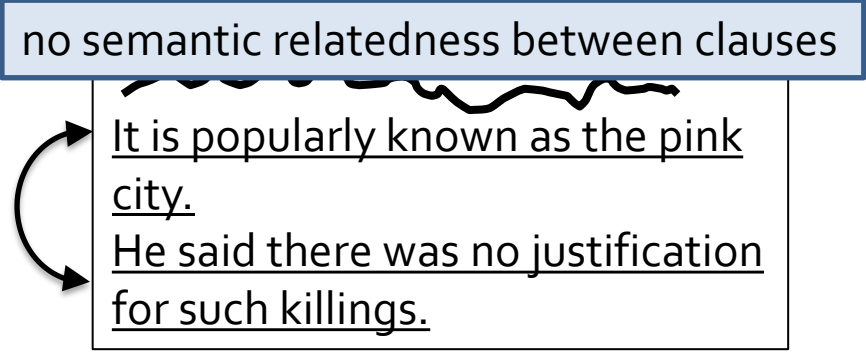
incomplete sentence

dateline included



GEORGETOWN, Pennsylvania  
2006-10-05 16:53:53 UTC

no semantic relatedness between clauses



It is popularly known as the pink city.  
He said there was no justification for such killings.

# Annotation Scheme: Clause level

(sentence, phrase, sequence of tokens)

ungrammaticality

incomplete sentence

dateline included

GEORGETOWN, Pennsylvania  
2006-10-05 16:53:53 UTC

redundant information

He was acting out in revenge  
for something that happened  
20 years ago....

...was apparently acting in  
revenge for an incident that  
happened to him 20 years ago.

no semantic relatedness between clauses

It is popularly known as the pink  
city.  
He said there was no justification  
for such killings.

# Annotation Scheme: Clause level

(sentence, phrase, sequence of tokens)

ungrammaticality

incomplete sentence

dateline included

GEORGETOWN, Pennsylvania  
2006-10-05 16:53:53 UTC

redundant information

He was acting out in revenge  
for something that happened  
20 years ago....

...was apparently acting in  
revenge for an incident that  
happened to him 20 years ago.

no semantic relatedness between clauses

It is popularly known as the pink  
city.  
He said there was no justification  
for such killings.

inappropriate use of  
discourse connective



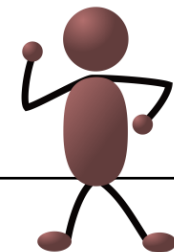
# LQVSumm: Annotated Data

	TAC
data source	1935 summaries, TAC 2011 (initial summaries), generated by 44 different <u>extractive</u> summarization systems
input to systems	sets of 10 news articles
Output	100-word summaries
summarization approaches	sentence selection + compression



# LQVSumm: Annotated Data

	TAC
data source	1935 summaries, TAC 2011 (initial summaries), generated by 44 different <u>extractive</u> summarization systems
input to systems	sets of 10 news articles
Output	100-word summaries
summarization approaches	sentence selection + compression
manual scores for summaries	Readability (1-5), Pyramid (content), Responsiveness (1-5)



# Inter-annotator agreement

- 100 randomly chosen summaries
- two annotators (A) and (B)
- annotations match if same type & overlapping span



# Inter-annotator agreement

- 100 randomly chosen summaries
- two annotators (A) and (B)
- annotations match if same type & overlapping span

level	Precision(B:A)	Recall(B:A)	F1
entity mention	90.4	54.5	67.5
clause	84.1	83.3	83.6

# Inter-annotator agreement

- 100 randomly chosen summaries
- two annotators (A) and (B)
- annotations match if same type & overlapping span

level	Precision(B:A)	Recall(B:A)	F1
entity mention	90.4	54.5	67.5
clause	84.1	83.3	83.6

A creates twice as many annotations,  
B's annotations are a  
subset of A's

# Inter-annotator agreement

- 100 randomly chosen summaries
- two annotators (A) and (B)
- annotations match if same type & overlapping span

level	Precision(B:A)	Recall(B:A)	F1
entity mention	90.4	54.5	67.5
clause	84.1	82.5	83.6

Agreement higher on clause level than on entity mention level

# Inter-annotator agreement

- 100 randomly chosen summaries
- two annotators (A) and (B)
- annotations match if same type & overlapping span

level	Precision(B:A)	Recall(B:A)	F1
entity mention	90.4	54.5	67.5
clause	84.1	83.3	83.6



degree of subjectivity is manageable

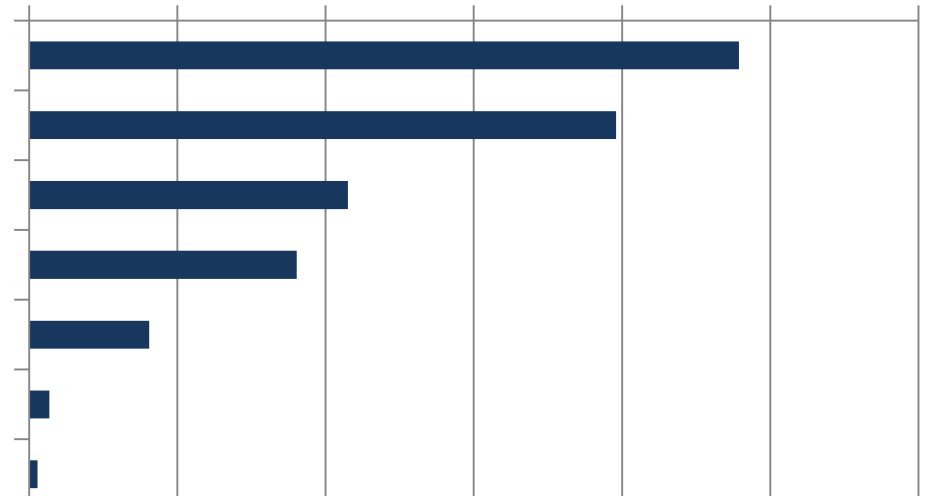
# Absolute Frequencies of LQVs by type

total: 1935 summaries

## Entity mention level

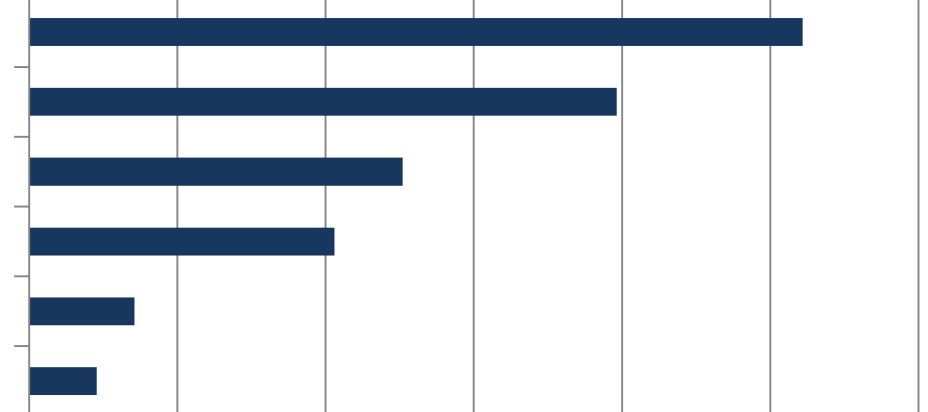
0 200 400 600 800 1000 1200

def. NP without reference  
unclear first mention  
indef. NP with previous reference  
pronoun without antecedent  
overly-specific subsequent mention  
pronoun with misleading antecedent  
unclear acronym



## Clause level

incomplete sentence  
ungrammaticality  
redundant information  
dateline included  
no semantic relatedness between clauses  
inappropriate discourse connective



# Ranking systems: average number of violations per summary

- compare rankings with TAC 2011 rankings
- draw conclusions about strengths/weaknesses of systems

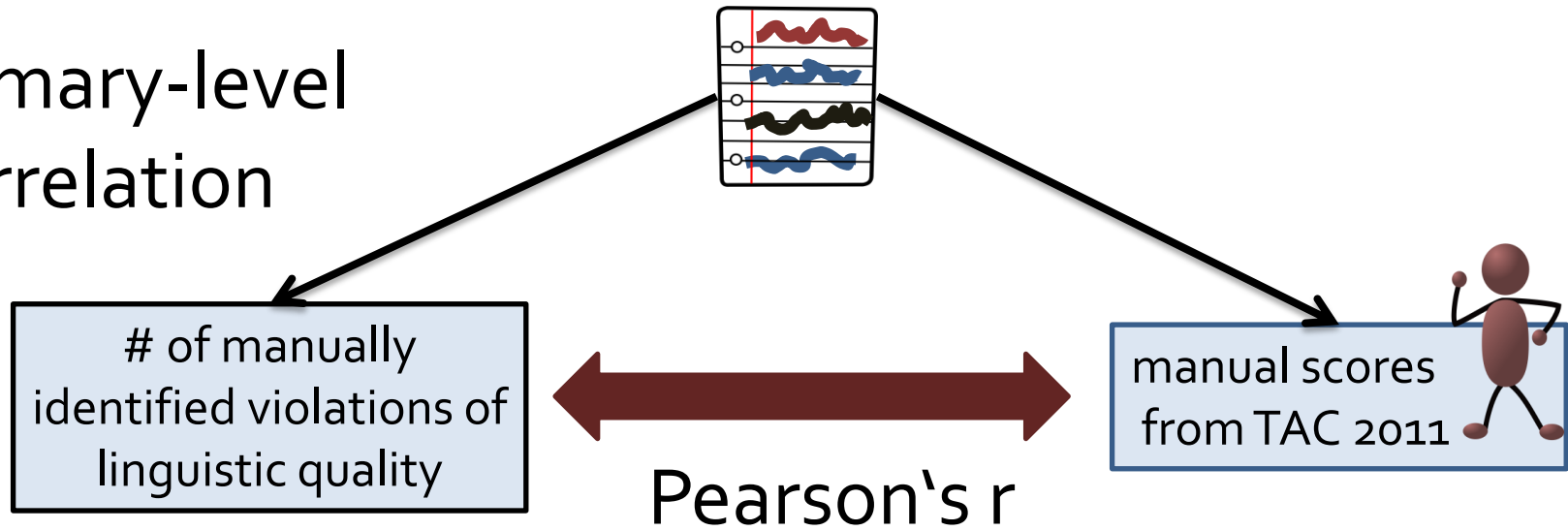
System	Entity mention level	Clause level	All LQV types
<b>1</b> (baseline using first 100 words as summary)	0.34	1	1.34
<b>21</b>	0.84	0.45	1.3
...	...	...	...
<b>7</b>	1.14	4.63	5.77

# Ranking systems: average number of violations per summary

- compare rankings with TAC 2011 rankings
- draw conclusions about strengths/weaknesses of systems

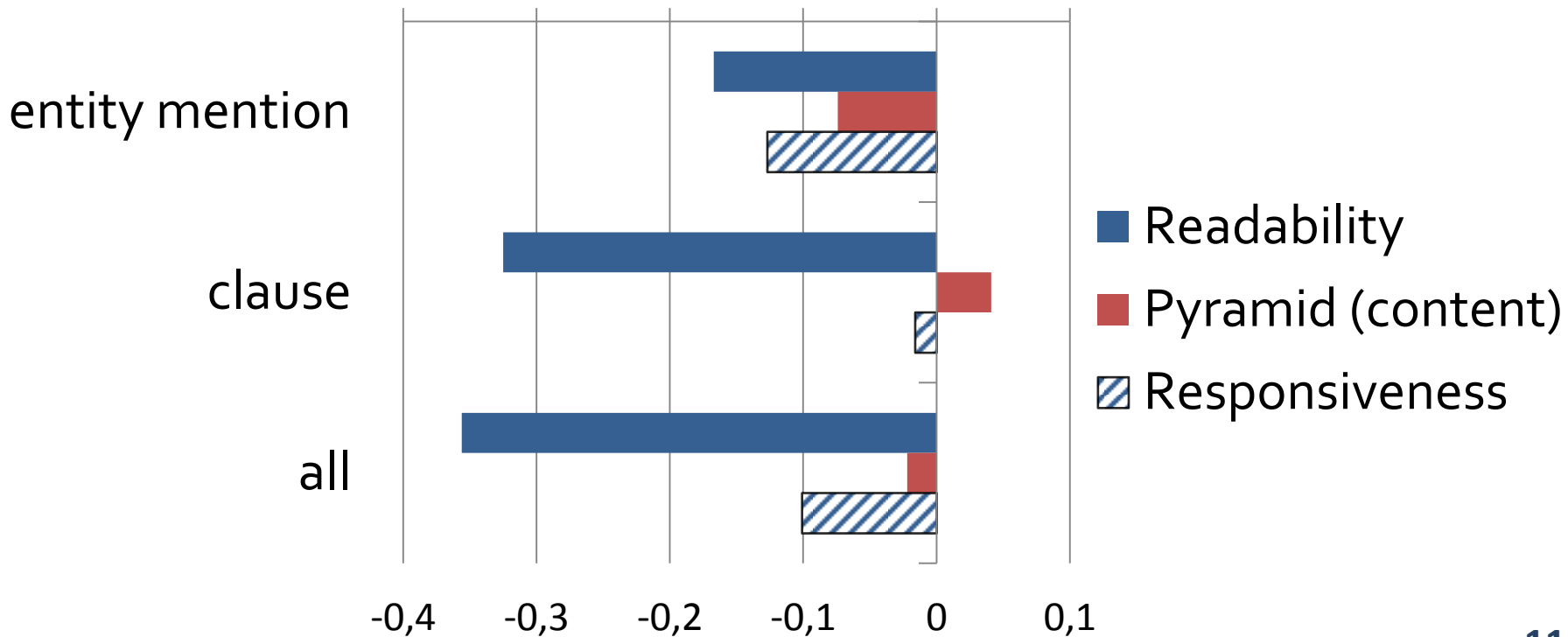
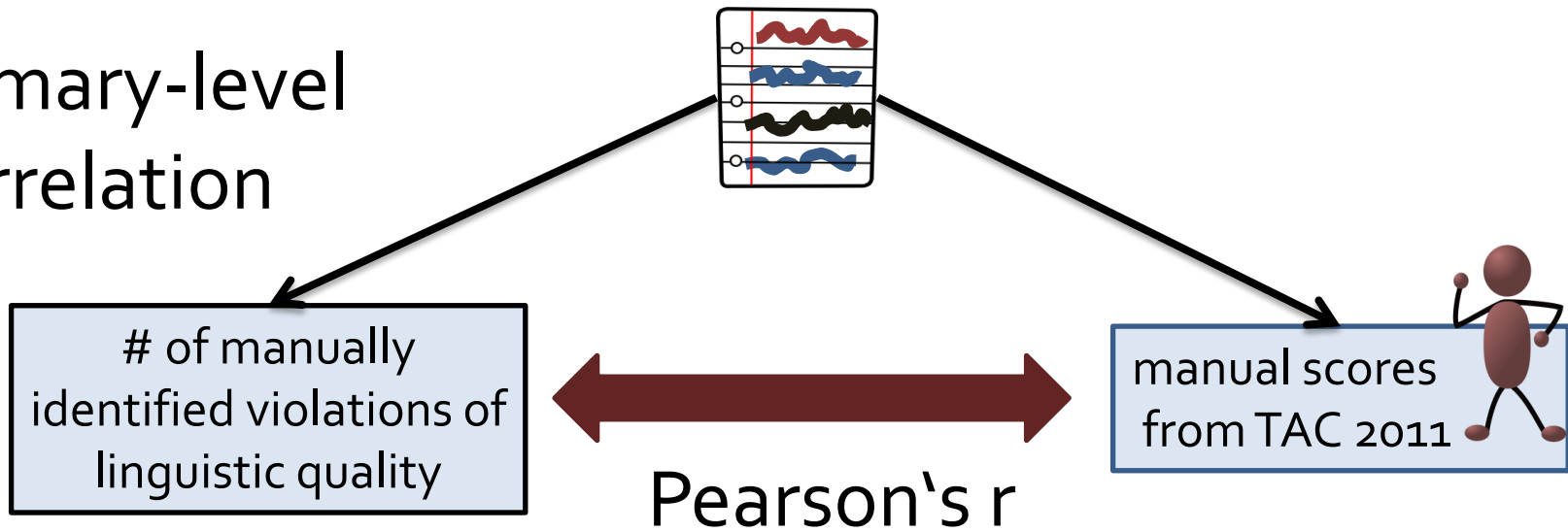
System	Entity mention level	Clause level	All LQV types
<b>1</b> (baseline using first 100 words as summary)	0.34	1	1.34
<b>21</b>	0.84	0.45	1.3
...	...	...	...
<b>7</b>	1.14	4.63	5.77
<b>Best TAC system</b> (differs for each column, TAC 2011)	(System 1) 0.34	(System 16) 0.23	(System 21) 1.30
<b>Average of systems in TAC</b>	1.42	1.54	2.96

# Summary-level correlation





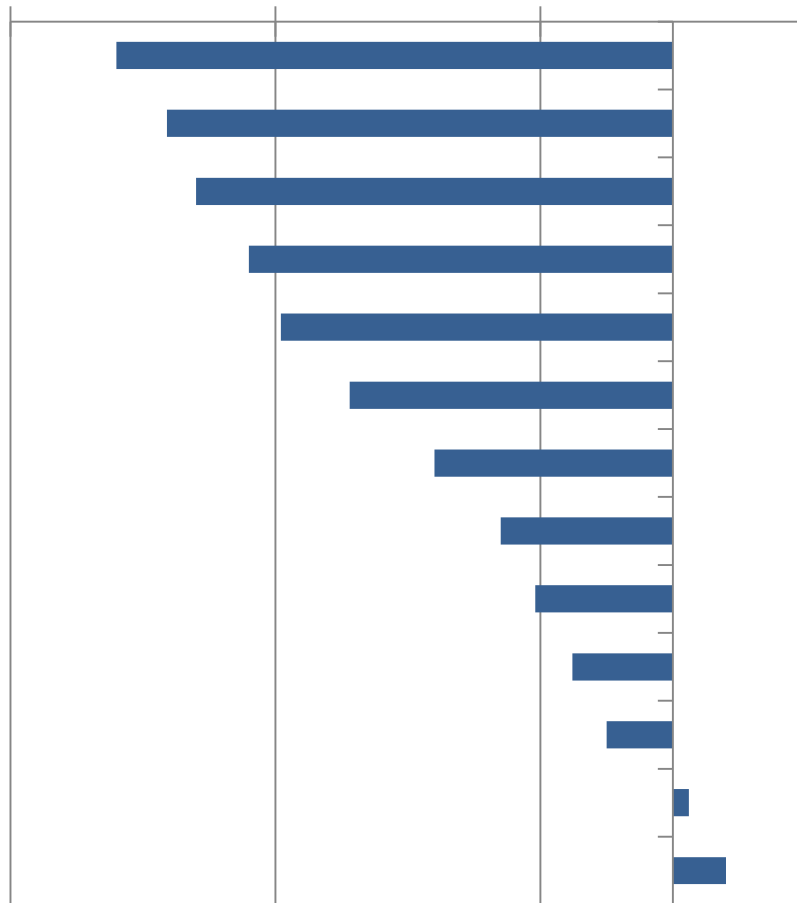
# Summary-level correlation



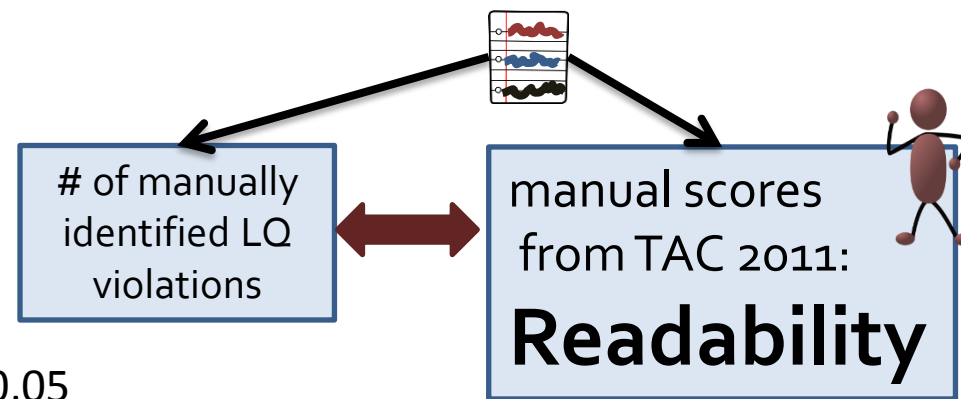
# Summary-level correlation

Pearsons's r

-0,25      -0,15      -0,05      0,05

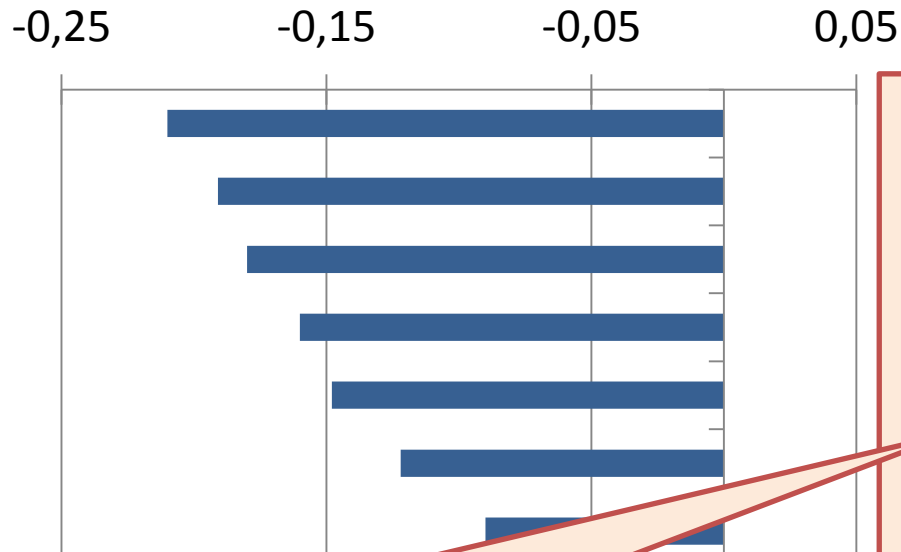


- incomplete sentence
- pronoun without antecedent
- ungrammaticality
- redundant information
- no semantic relatedness between clauses
- def. NP without referent
- dateline included
- pronoun with misleading antecedent
- indef. NP with previous referent
- unclear acronym
- inappropriate discourse connective
- unclear first mention
- overly specific subsequent mention



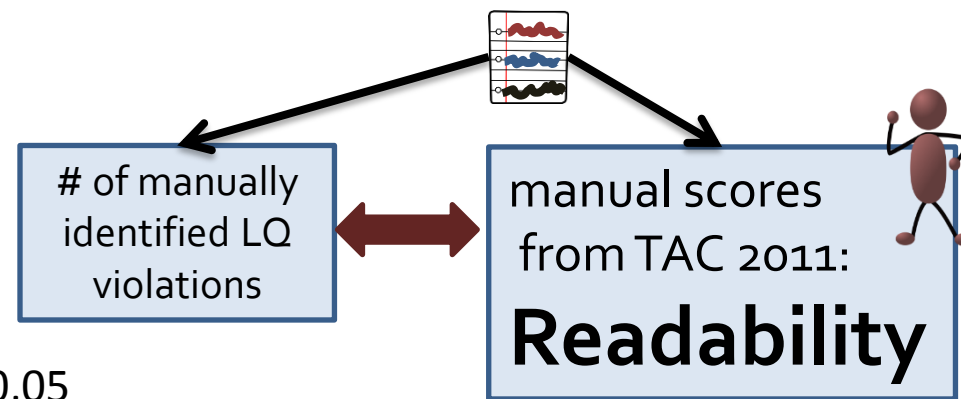
# Summary-level correlation

Pearsons's r

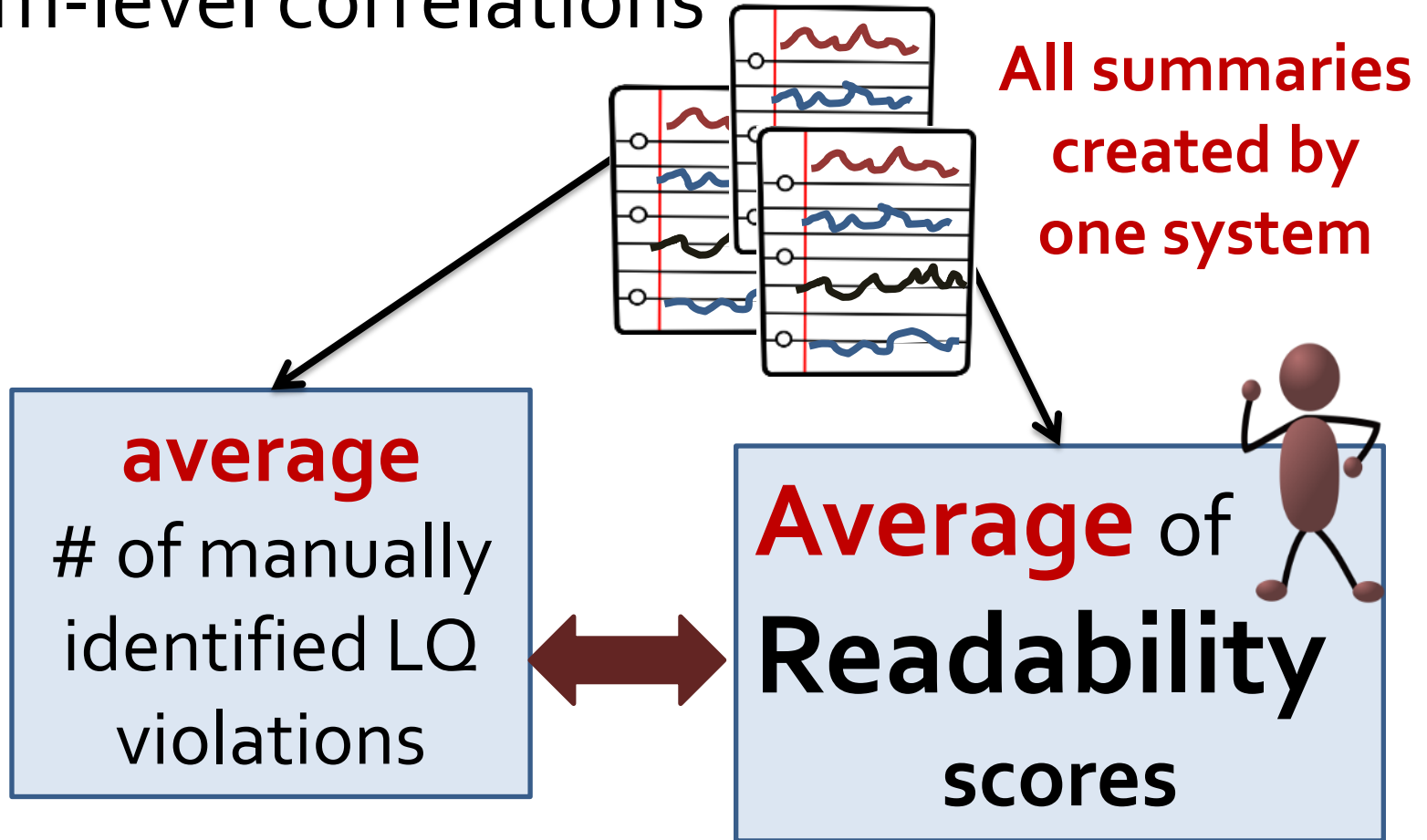


Significantly correlated to intuitively assigned Readability scores  
→ **play a role for judgment**

- incomplete sentence
- pronoun without antecedent
- ungrammaticality
- redundant information
- no semantic relatedness between clauses
- def. NP without referent
- dateline included
- pronoun with misleading antecedent
- indef. NP with previous referent
- unclear acronym
- inappropriate discourse connective
- unclear first mention
- overly specific subsequent mention

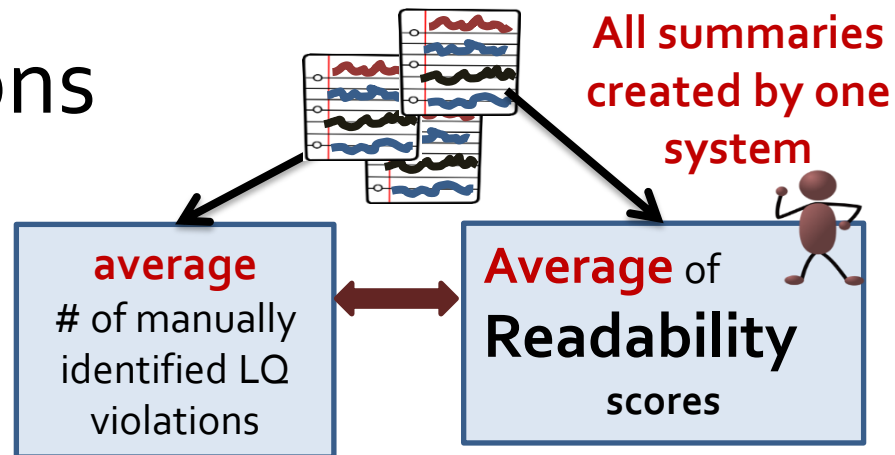


# System-level correlations



<i>System 21</i>	1.30		3.75
<i>System 2</i>	1.74	←→	3.34
<i>System 7</i>	5.77		2.09
...	...		...

# System-level correlations

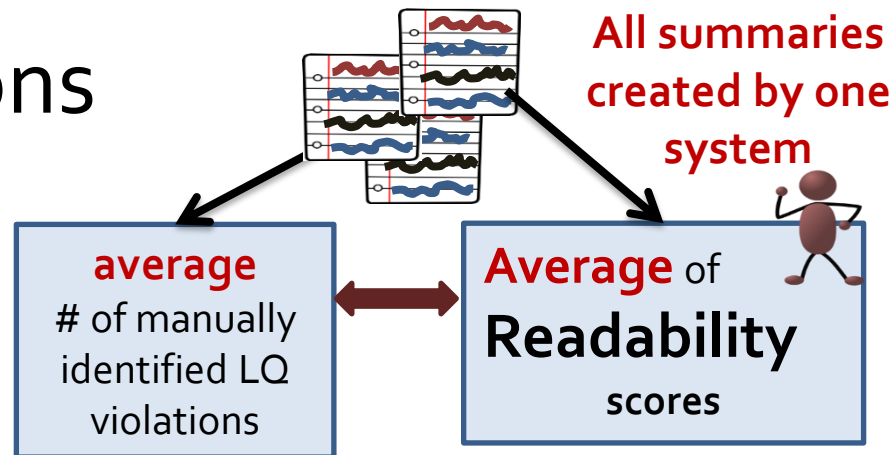


**DICOMER:** features from Penn Discourse TreeBank-style discourse parser

*higher absolute correlation*  
→ *better ranking*

Method	Ranking of	Pearson's $r$	Spearman's $\rho$	Kendall's $\tau$
DICOMER [Lin et al. 2012]	<i>all 50 systems</i>	<b>0.867</b>	0.712	0.535
LQVSumm <i>sum(violations)</i>	<i>44 systems</i>	-0.820	<b>-0.858</b>	<b>-0.713</b>

# System-level correlations



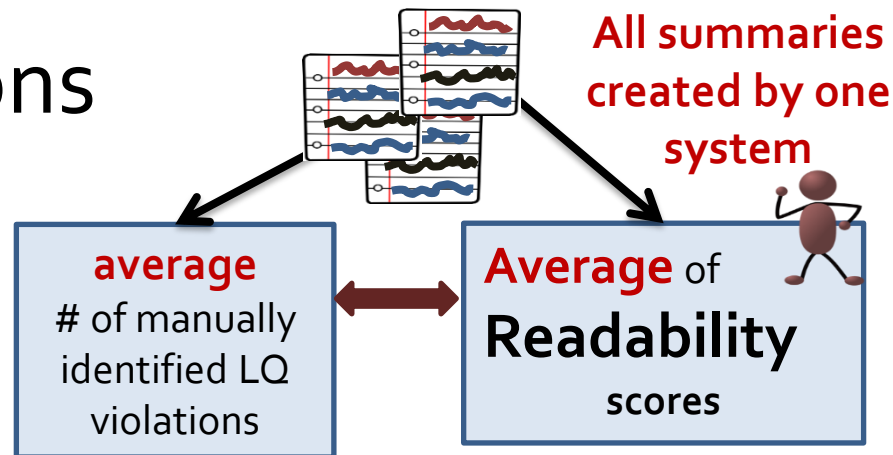
**DICOMER:** features from Penn Discourse TreeBank-style discourse parser

*higher absolute correlation*  
 → *better ranking*

Method	Ranking of	Pearson's $r$	Spearman's $\rho$	Kendall's $\tau$
DICOMER [Lin et al. 2012]	<i>all 50 systems</i>	<b>0.867</b>	0.712	0.535
LQVSumm <i>sum(violations)</i>	<i>44 systems</i>	-0.820	<b>-0.858</b>	<b>-0.713</b>

Pearson's $r$	Spearman's $\rho$ , Kendall's $\tau$
actual scores	ranking only

# System-level correlations



**DICOMER:** features from Penn Discourse TreeBank-style discourse parser

*higher absolute correlation*  
 → *better ranking*

Method	Ranking of	Pearson's $r$	Spearman's $\rho$	Kendall's $\tau$
DICOMER [Lin et al. 2012]	<i>all 50 systems</i>	<b>0.867</b>	0.712	0.535
LQVSumm <i>sum(violations)</i>	<i>44 systems</i>	-0.820	<b>-0.858</b>	<b>-0.713</b>

Pearson's $r$	Spearman's $\rho$ , Kendall's $\tau$
actual scores	ranking only
DICOMER is better (trained on TAC 2009 & TAC 2010)	<b>counting the number of violations works better than a supervised system.</b>

# Conclusions

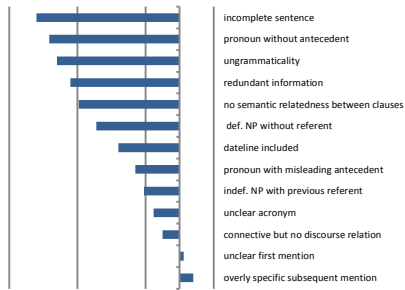
**LQVSumm:**

≈ 2000 summaries  
marked with LQV types



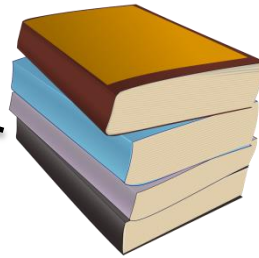


# Conclusions

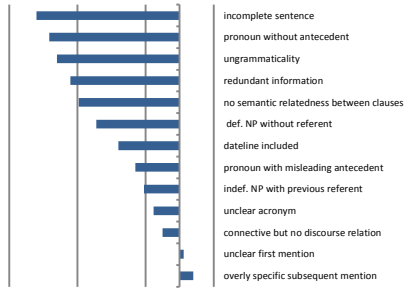


most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types

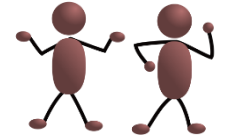
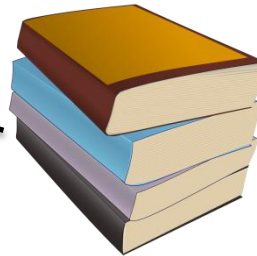


# Conclusions



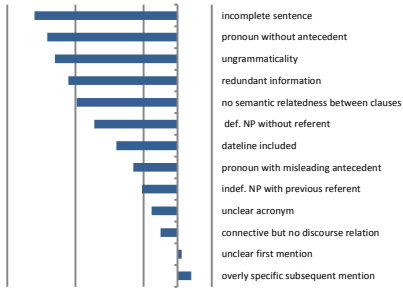
most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types



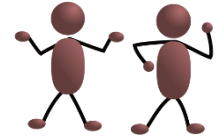
good inter-annotator agreement

# Conclusions

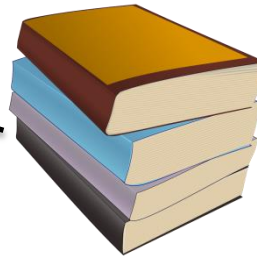


most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types

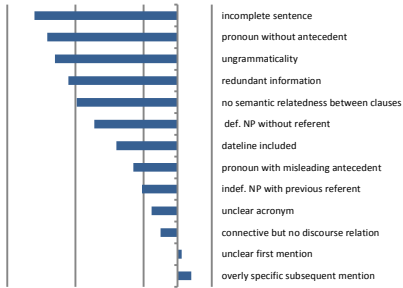


good inter-annotator agreement



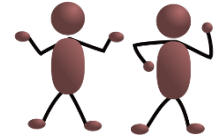
counts and marked instances of linguistic quality violations allow for:

# Conclusions

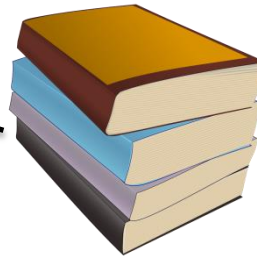


most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types



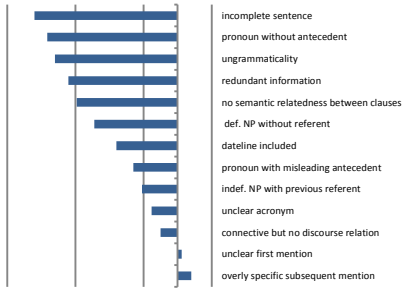
good inter-annotator agreement



counts and marked instances of linguistic quality violations allow for:

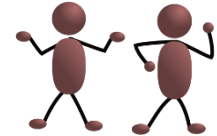
analyzing what a particular system is good/bad at (rather than just obtaining a numeric score)

# Conclusions

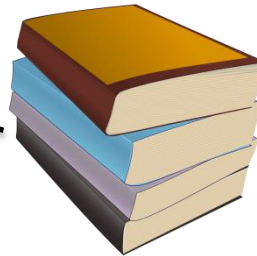


most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types



good inter-annotator agreement

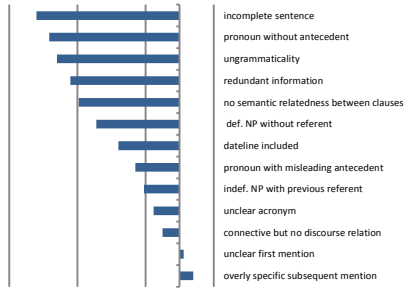


counts and marked instances of linguistic quality violations allow for:

analyzing what a particular system is good/bad at (rather than just obtaining a numeric score)

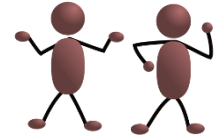
developing automatic methods to detect LQVs (future work)

# Conclusions

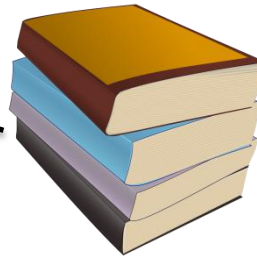


most types correlated to human judgments; others are infrequent

**LQVSumm:**  
≈ 2000 summaries  
marked with LQV types



good inter-annotator agreement



counts and marked instances of linguistic quality violations allow for:

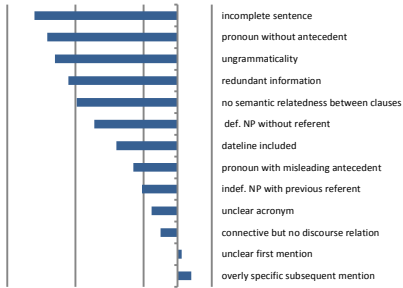
analyzing what a particular system is good/bad at (rather than just obtaining a numeric score)

developing automatic methods to detect LQVs (future work)

Available in stand-off format at: [www.coli.uni-saarland.de/~afried](http://www.coli.uni-saarland.de/~afried)

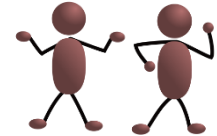
# Conclusions

# Thanks!

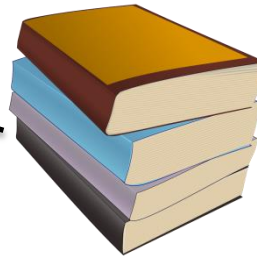


most types correlated to human judgments; others are infrequent

LQVSumm:  
≈ 2000 summaries  
marked with LQV types



good inter-annotator agreement



counts and marked instances of linguistic quality violations allow for:

analyzing what a particular system is good/bad at (rather than just obtaining a numeric score)

developing automatic methods to detect LQVs (future work)

Available in stand-off format at: [www.coli.uni-saarland.de/~afried](http://www.coli.uni-saarland.de/~afried)

# Backup Slides



# Annotation Scheme: Overview

## entity mention level

- pronouns without antecedents
- indefinite NPs with a previous mention
- ...

## clause level

(sentence, phrase,  
sequence of tokens)

- ungrammatical sentences
- no semantic relatedness
- ...

# Performance of the G-Flow summarization system


- G-Flow system: Christensen et al. (NAACL 2013): Towards Coherent Multi-Document Summarization
- system incorporates coherence information into sentence extraction
- marked 50 summaries provided on the web site of the authors

System	Entity mention level	Clause level	All LQV types
<b>Best TAC system</b> (differs for each column, TAC 2011)	(System 1) 0.34	(System 16) 0.23	(System 21) 1.30
<b>G-Flow</b> (DUC 2004 data)	<b>0.30</b>	<b>0.20</b>	<b>0.50</b>



G-Flow succeeds in producing more coherent / readable summaries

# inappropriate use of discourse connective



Taylor's attorney could not be reached for comment Friday night.

**And** the person who cooperates first gets the biggest reward.

