

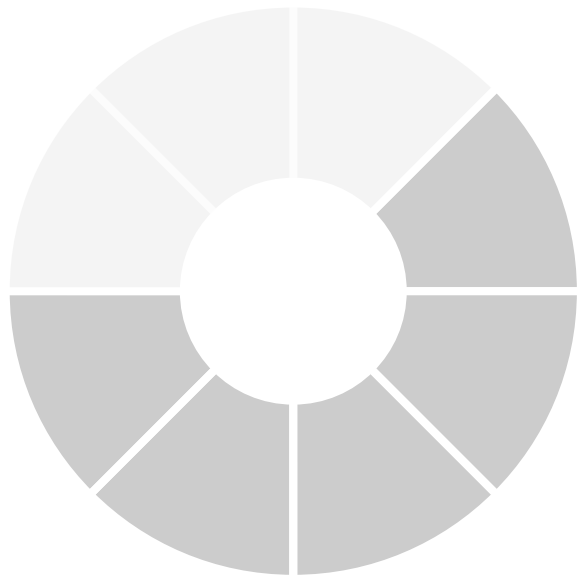
# Tackling Societal Challenges with Style Analysis

---

Martin Potthast  
Leipzig University  
[www.temir.org](http://www.temir.org)

joint work with the  
Webis Group  
[www.webis.de](http://www.webis.de)

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

1

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

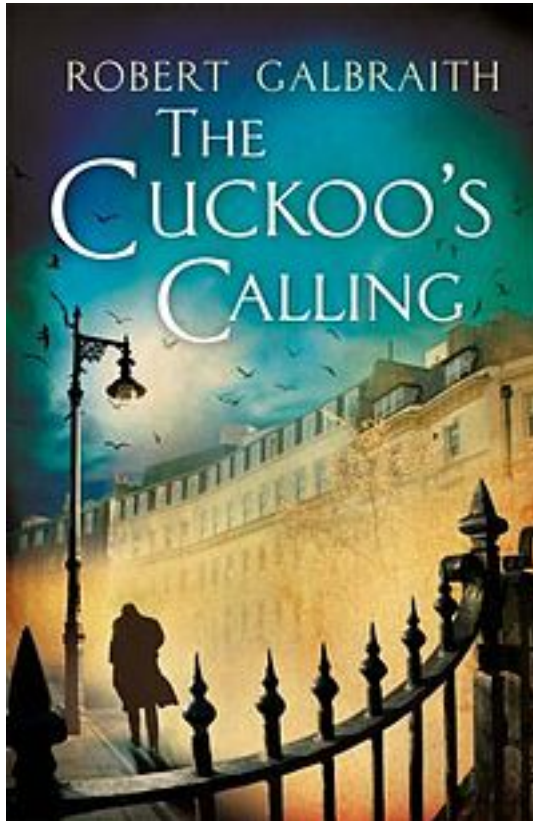
Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship



Fake likes

Fake news

Fake clicks

Fake users

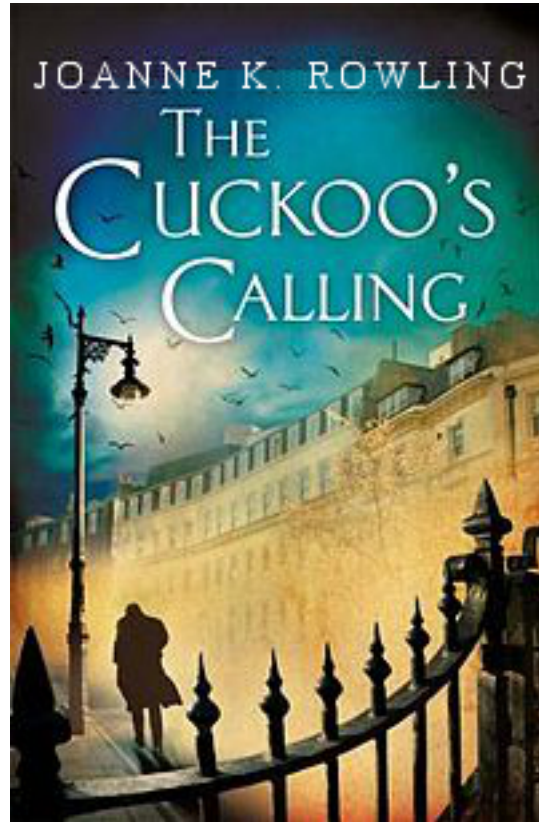
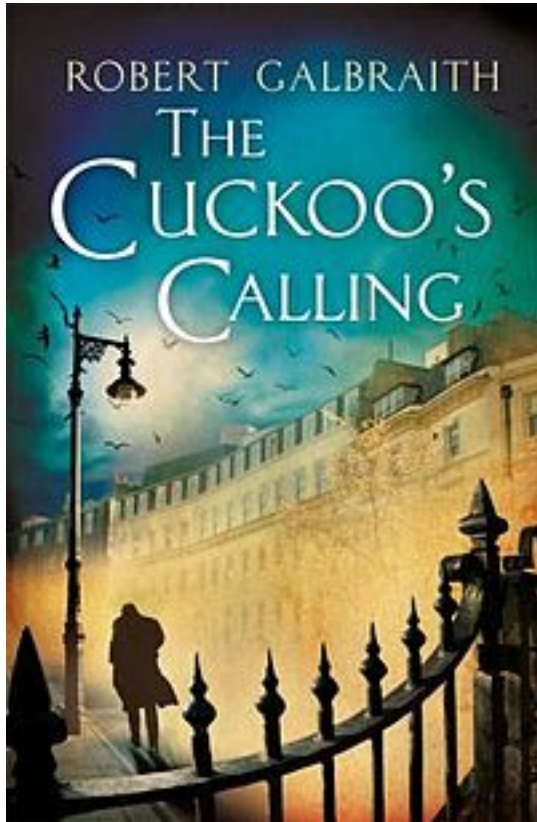
Fake reviews

Fake comments

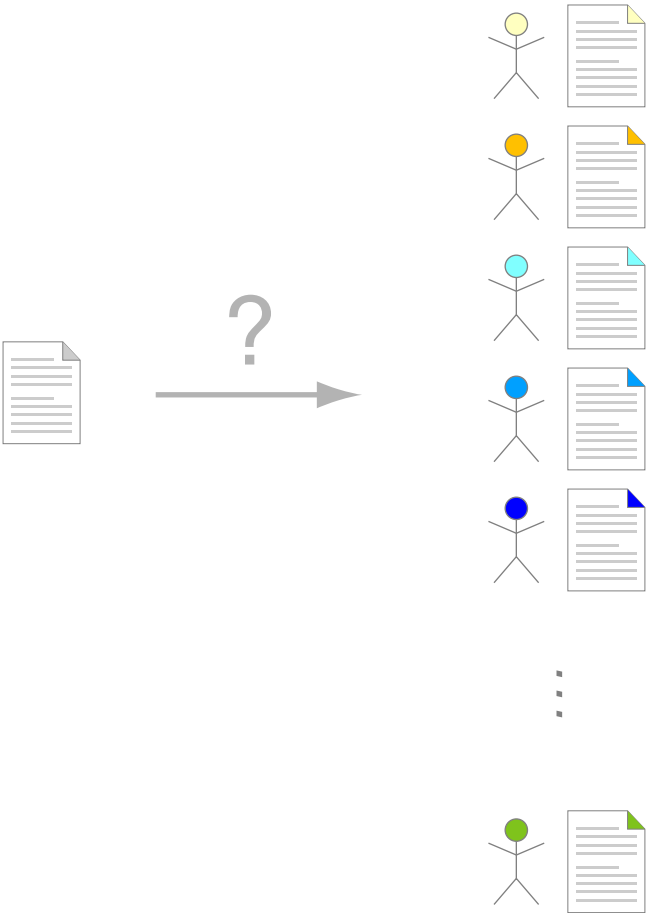
⋮

Fake identities (pseudonyms)



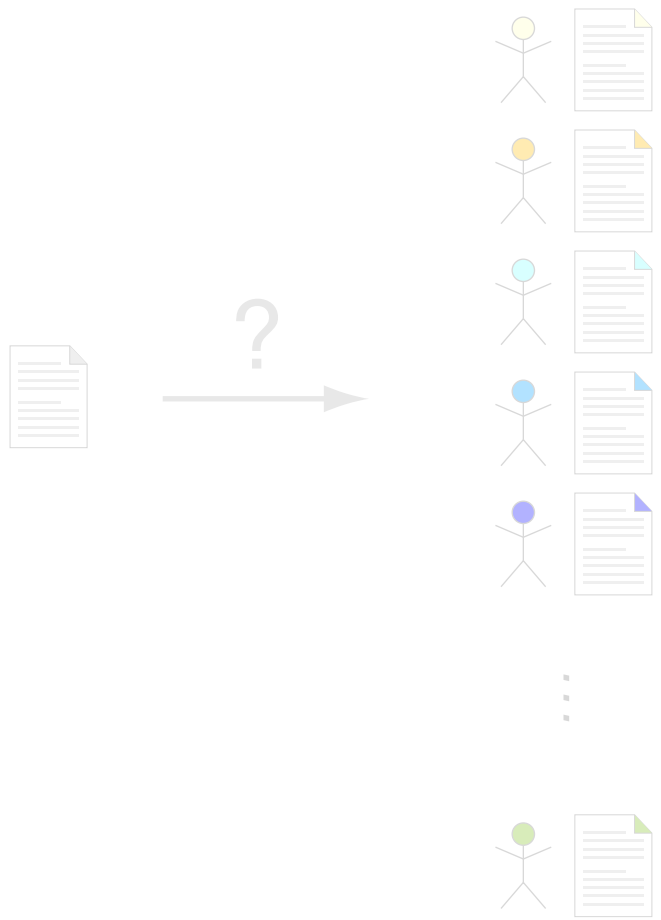


# Authorship Attribution



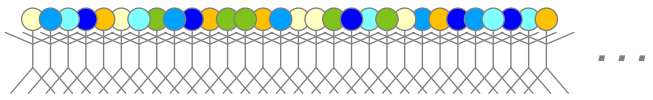
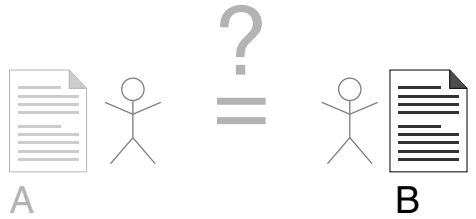
To which author does a text belong?

# Authorship Attribution



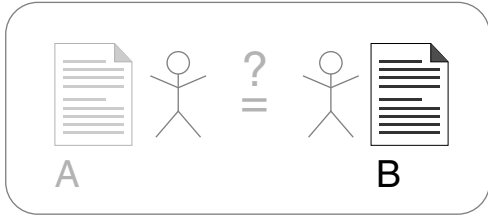
To which author does a text belong?

# Authorship Verification

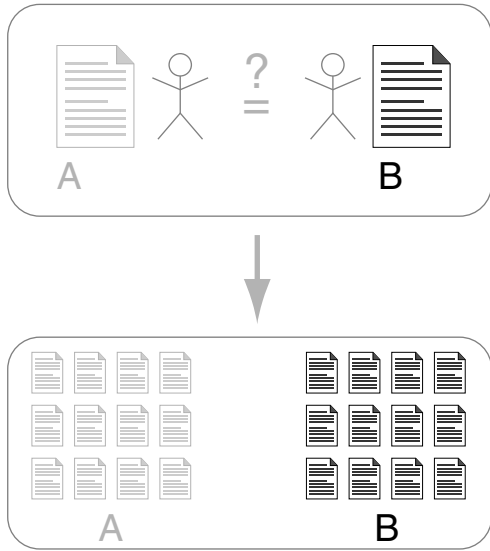


Originate two texts from the same author?

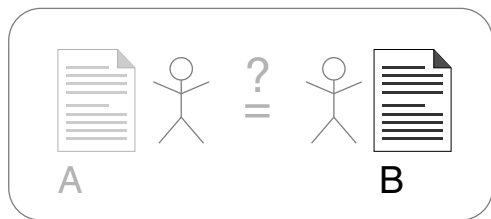
# Authorship Verification “Unmasking” [Koppel/Schler 2004]



# Authorship Verification "Unmasking" [Koppel/Schler 2004]



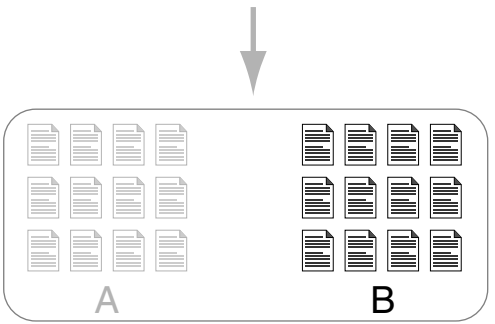
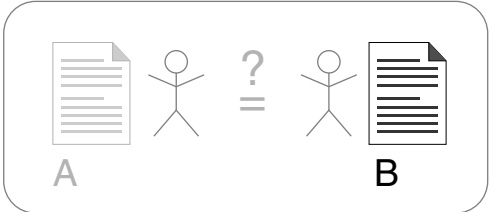
# Authorship Verification “Unmasking” [Koppel/Schler 2004]



0.3	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.3	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.6	0.3	0.2	0.2	0.2
0.0	0.0	0.3	0.1	0.6	0.0	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.4	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0

**A** **B**

# Authorship Verification "Unmasking" [Koppel/Schler 2004]

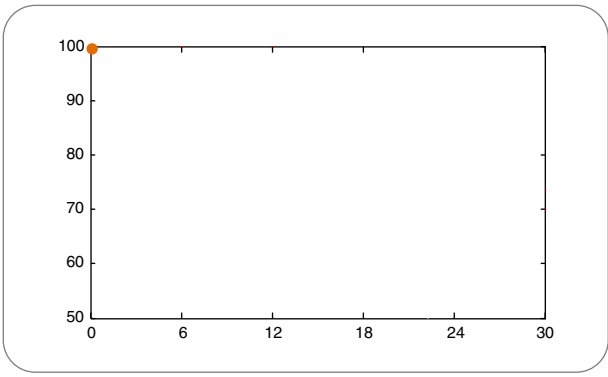
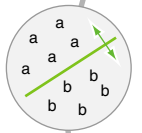


0.3	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.3	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.6	0.3	0.2	0.2	0.2
0.0	0.0	0.3	0.1	0.6	0.0	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.4	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0

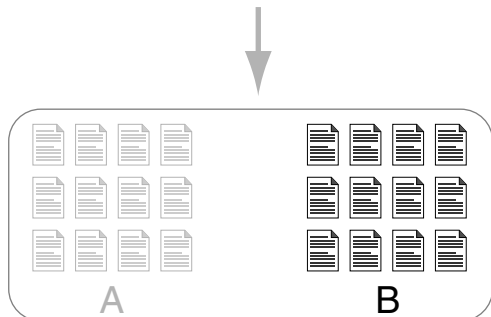
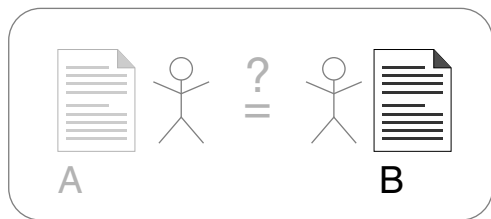
A

0.0	0.2	0.5	0.3
0.1	0.3	0.1	0.4
0.2	0.1	0.0	0.1
0.2	0.1	0.4	0.5
0.3	0.2	0.2	0.2
0.6	0.0	0.4	0.2
0.1	0.3	0.2	0.2
0.1	0.2	0.2	0.1
0.2	0.1	0.2	0.0
0.1	0.2	0.4	0.6
0.5	0.1	0.4	0.2
0.5	0.2	0.2	0.5
0.0	0.3	0.1	0.2
0.2	0.1	0.0	0.3
0.2	0.1	0.0	0.0

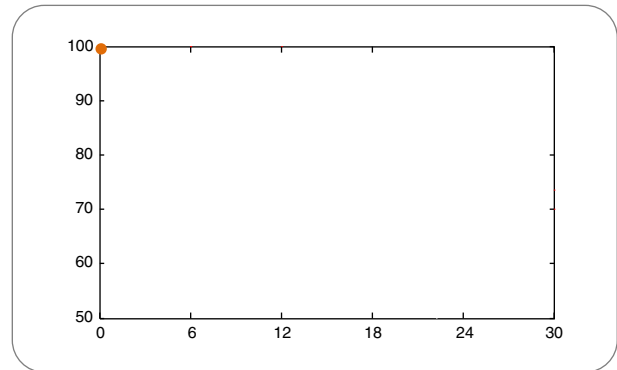
B



# Authorship Verification "Unmasking" [Koppel/Schler 2004]

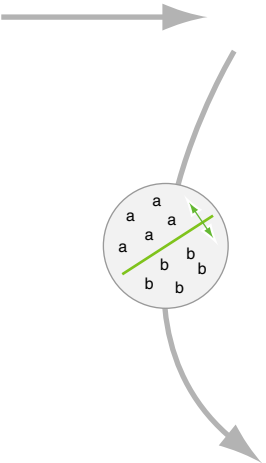
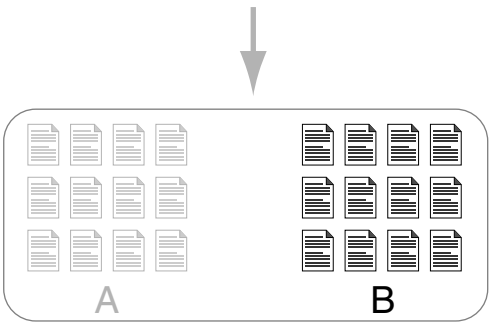
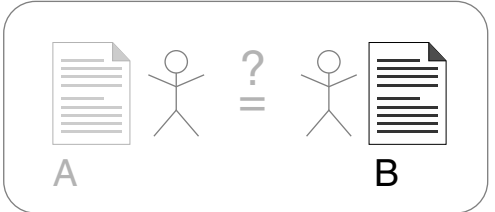


0.3	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.3	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.6	0.3	0.2	0.2	0.2
0.6	0.6	0.3	0.1	0.6	0.6	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.1	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0

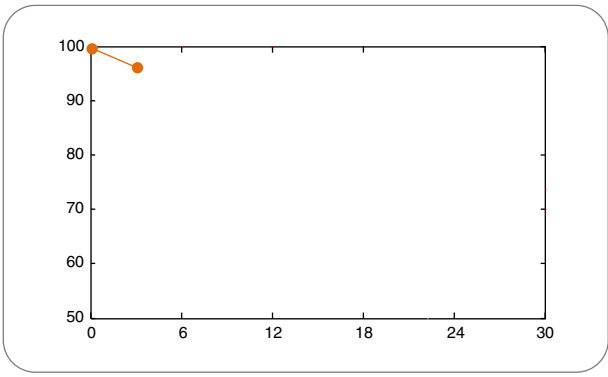




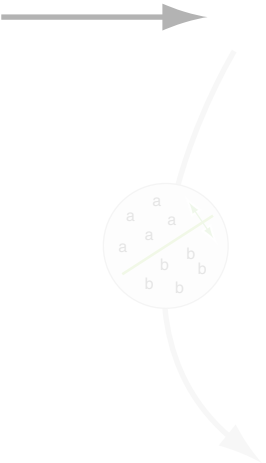
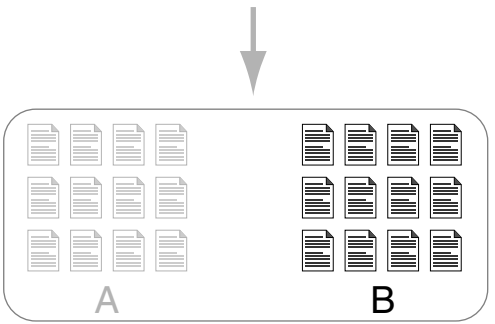
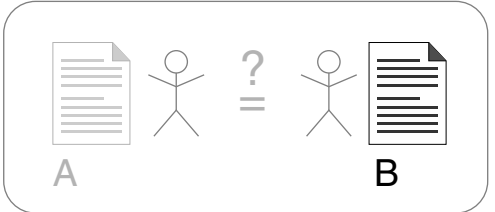
# Authorship Verification "Unmasking" [Koppel/Schler 2004]



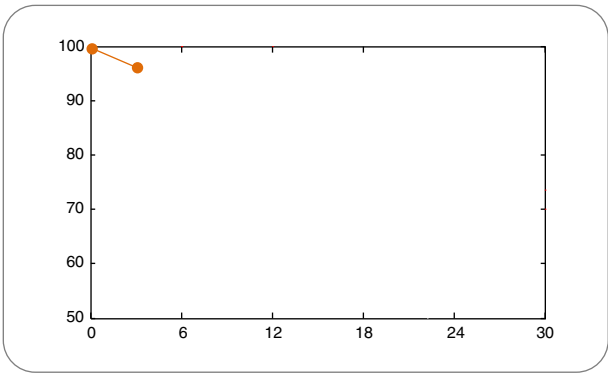
0.3	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.3	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.6	0.3	0.2	0.2	0.2
0.6	0.6	0.3	0.1	0.6	0.6	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.1	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0



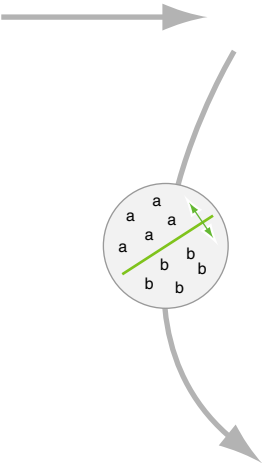
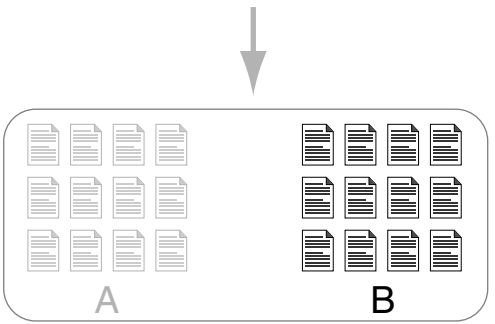
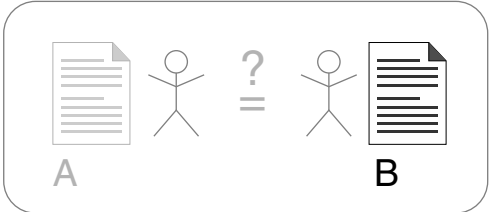
# Authorship Verification "Unmasking" [Koppel/Schler 2004]



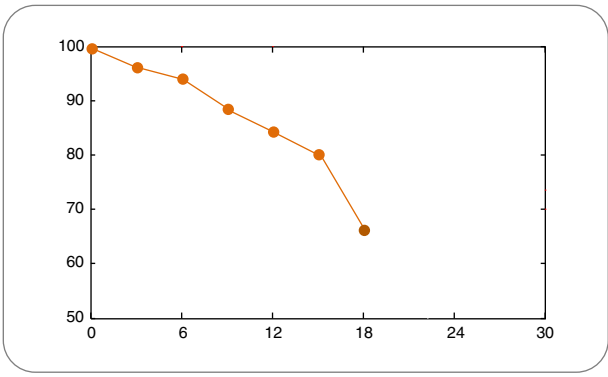
0.3	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.1	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.6	0.3	0.2	0.2	0.2
0.6	0.6	0.3	0.1	0.0	0.0	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.4	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0



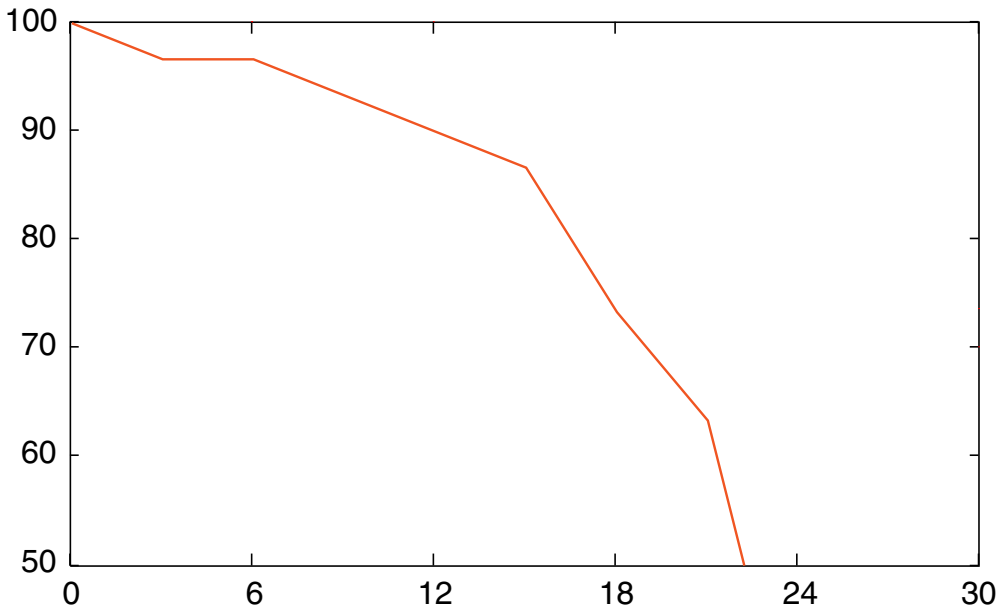
# Authorship Verification "Unmasking" [Koppel/Schler 2004]



0.0	0.1	0.1	0.2	0.0	0.2	0.5	0.3
0.2	0.2	0.4	0.1	0.1	0.3	0.1	0.4
0.0	0.4	0.1	0.2	0.2	0.1	0.0	0.1
0.1	0.0	0.2	0.3	0.2	0.1	0.4	0.5
0.2	0.2	0.1	0.0	0.3	0.2	0.2	0.2
0.0	0.0	0.3	0.1	0.0	0.0	0.4	0.2
0.3	0.3	0.4	0.2	0.1	0.3	0.2	0.2
0.1	0.1	0.1	0.3	0.1	0.2	0.2	0.1
0.0	0.1	0.2	0.2	0.2	0.1	0.2	0.0
0.1	0.2	0.6	0.1	0.1	0.2	0.4	0.6
0.4	0.4	0.2	0.3	0.5	0.1	0.4	0.2
0.5	0.1	0.3	0.2	0.5	0.2	0.2	0.5
0.2	0.2	0.1	0.1	0.0	0.3	0.1	0.2
0.1	0.3	0.2	0.4	0.2	0.1	0.0	0.3
0.0	0.2	0.0	0.1	0.2	0.1	0.0	0.0

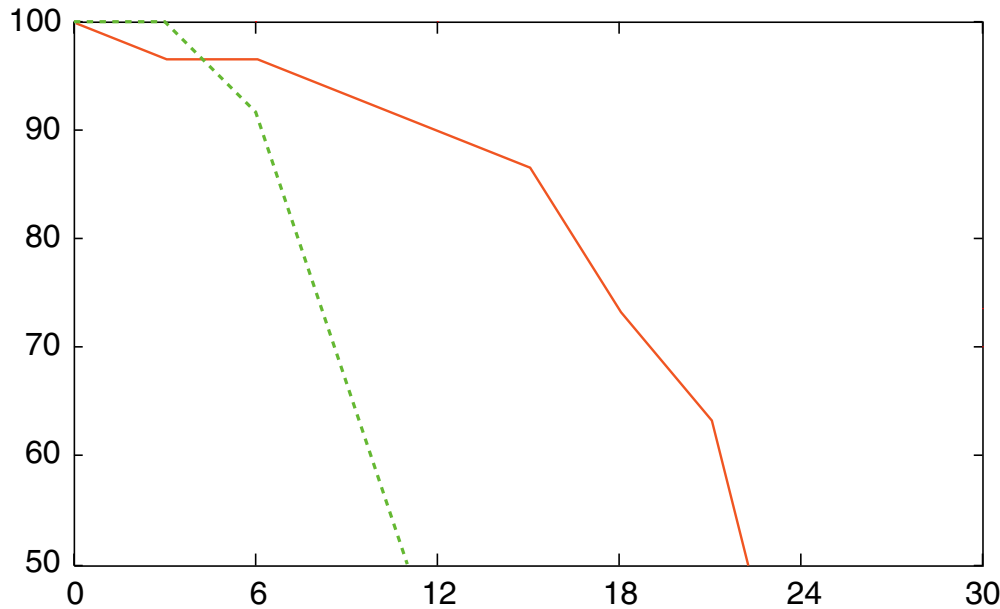


Typical learning characteristic for ...



different authors ( $A \neq B$ )

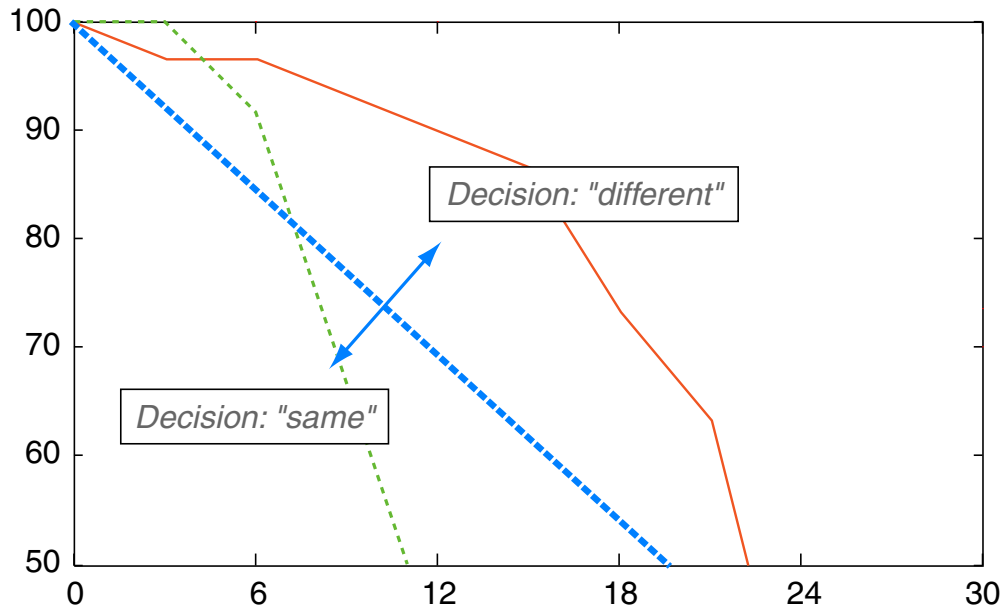
Typical learning characteristic for ...



different authors ( $A \neq B$ )  
same author ( $A = B$ )

# Authorship Verification Unmasking at Work

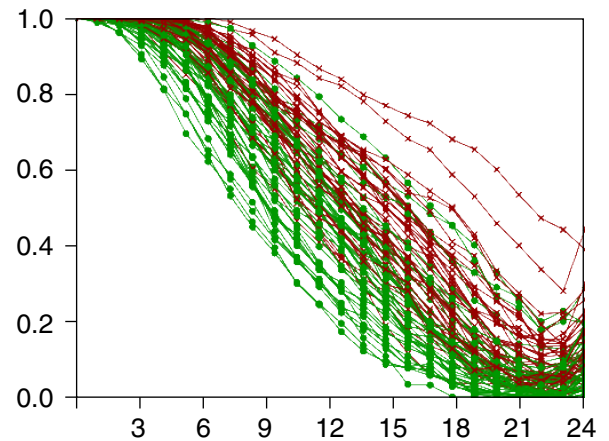
Typical learning characteristic for ...



different authors ( $A \neq B$ )  
same author ( $A = B$ )

The typical learning characteristic can be learned.

Experiment	I	II
<i>Performance</i>		
Precision	0.96	1.00
Accuracy	0.63	0.91
Classified	100%	26%
Omitted	0%	74%
<i>Configuration</i>		
Number of cases	180 training / 78 test	
Size of each case	4 000 words	
Number of authors	135	
Number of chunks	25	
Size of each chunk	600 words	
Vocabulary	250 words	
Removed per round	10 words	
Smoothing	no	

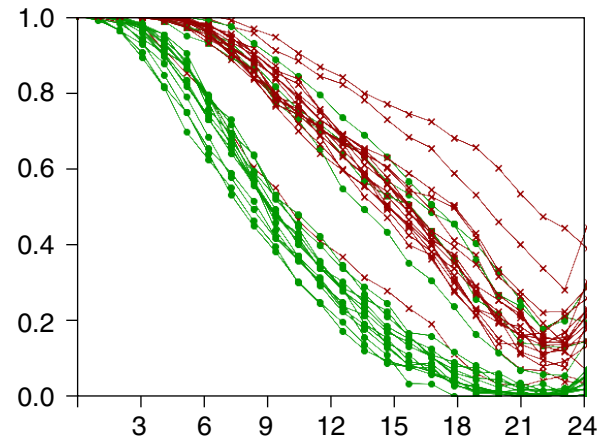
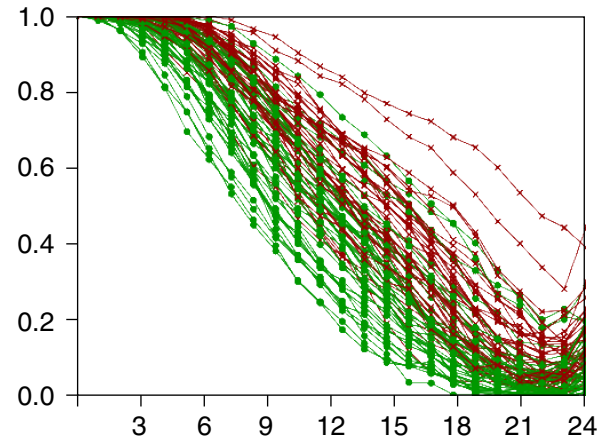


# Authorship Verification

Recent Results

[Bevendorff et al., 2019]

Experiment	I	II
<i>Performance</i>		
Precision	0.96	1.00
Accuracy	0.63	0.91
Classified	100%	26%
Omitted	0%	74%
<i>Configuration</i>		
Number of cases	180 training / 78 test	
Size of each case	4 000 words	
Number of authors	135	
Number of chunks	25	
Size of each chunk	600 words	
Vocabulary	250 words	
Removed per round	10 words	
Smoothing	no	





②

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

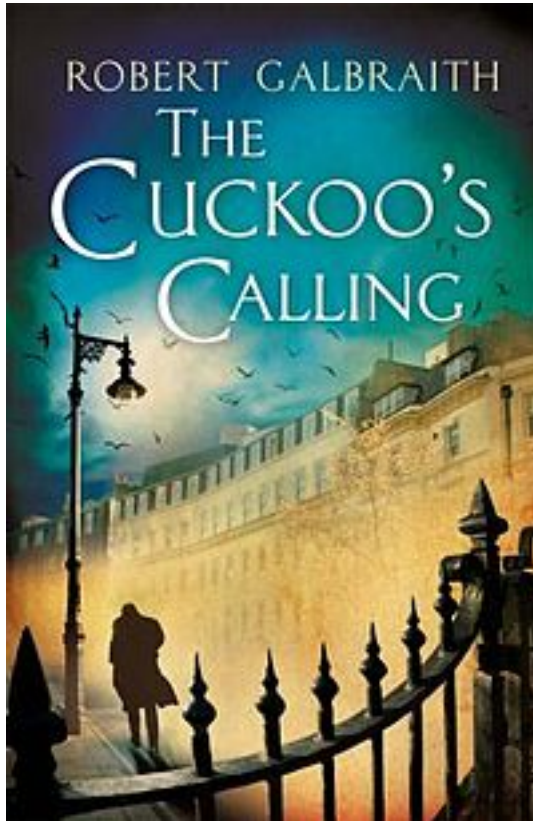
Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and  
Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship



Fake likes

Fake news

Fake clicks

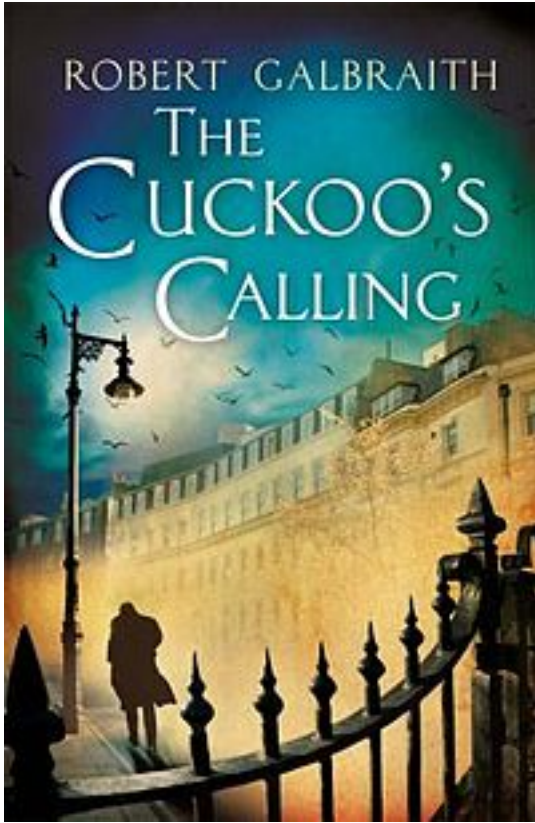
Fake users

Fake reviews

Fake comments

⋮

Fake identities (pseudonyms)



**DAILY PROPHET**  
THE WIZARD WORLD'S BEGILLING ENGAGEMENT OF CHOICE

OFFICIAL GUIDES TO ELEMENTARY HOME & PERSONAL DEFENCE WILL BE DELIVERED TO ALL WIZARDING HOME

National Weather  
Zodiac • Aspects

FIRST-SECOND EDITION

**SPECIAL EDITION**

# HE WHO MUST NOT BE NAMED RETURNS

HE WHO MUST NOT BE NAMED HAS RETURNED TO THIS COUNTRY AND IS ONCE MORE ACTIVE

— spells 2 — M. OF MAGIC AFFAIRS 3 — potions 6 — health 7 — RAD NEWS

# Countermeasure: Obfuscation

[Bevendorff et al., 2019]

A beautiful\_χchristmas you know jesus our saviour w patiently stooping to hunger and pain, so he mig ones\_χfrom shame; now if we love him, he bids u brothers and sisters who need. blessed old nick! i it, you would remember and certainly do it; this you empty your pack, pray give a portion to all w there's anything left and you can bring a small gi wasn't that dandy? sure, little mary\_χann has a wo she has! she\_χtakes after her own mother. i was jus that age. and you're just like\_χher still, mollie mullig

B sure, little mary\_χann has a wonderful education, s s after her own mother. i was just like her when i wa e just like\_χher still, mollie mulligan. sure you're\_χth an alley and the belle of shantytown. whist now! it lushes. but, hush! i think the show is about to begi ο, samson symbolical! come and see slivers\_χclow me and see zip, the foremost of freaks! come an ister sheiks! eager equestriennes\_χeach unexcelle enagerie ever beheld, the giant, the fat girl, the lion artists from far-off japan, audacious acrobats sho

Idea: Obfuscate by increasing the Kullback-Leibler Divergence (KLD).

Desired: A “minimally invasive” procedure.

Strategy: Determine “high-impact” n-grams.  $\frac{\partial}{\partial q} \left( p \log_2 \frac{p}{q} \right) \rightarrow \max,$

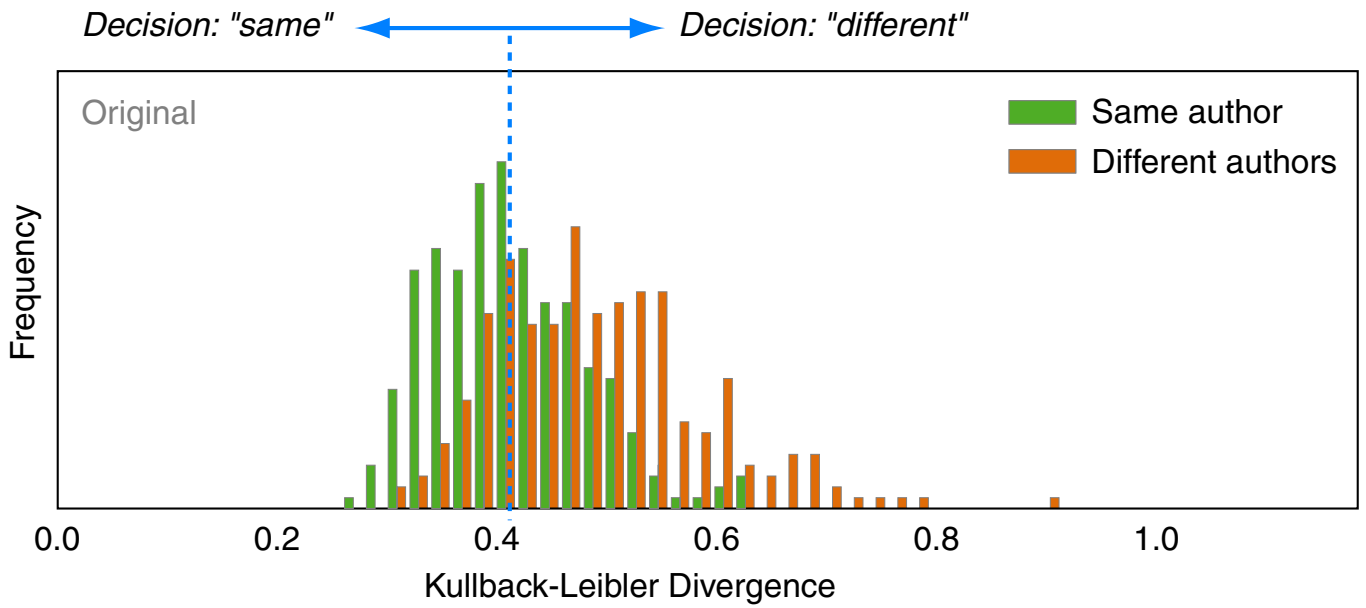
where  $p$  and  $q$  denote the n-gram-specific occurrence probabilities in texts A and B respectively.

# Countermeasure: Obfuscation

[Bevendorff et al., 2019]

A beautiful\_χchristmas you know jesus our saviour w patiently stooping to hunger and pain, so he mig ones\_s,χfrom shame; now if we love him, he bids u brothers and sisters who need. blessed old nick! i it, you would remember and certainly do it; this you empty your pack, pray give a portion to all wh there's anything left and you can bring a small gi wasn't that dandy? sure, little mary\_χann has a wor she has! she\_eχtakes after her own mother. i was jus that age. and you're just like\_χher still, mollie mullig

sure, little mary\_χann has a wonderful education, s s after her own mother. i was just like her when i wa e just like\_χher still, mollie mulligan. sure you're\_eχth an alley and the belle of shantytown. whist now! it luses. but, hush! i think the show is about to begi ο, samson symbolical! come and see slivers\_s,χclow me and see zip, the foremost of freaks! come an ister sheiks! eager equestriennes\_s,χeach unexcelle enagerie ever beheld, the giant, the fat girl, the lion artists from far-off japan, audacious acrobats sho

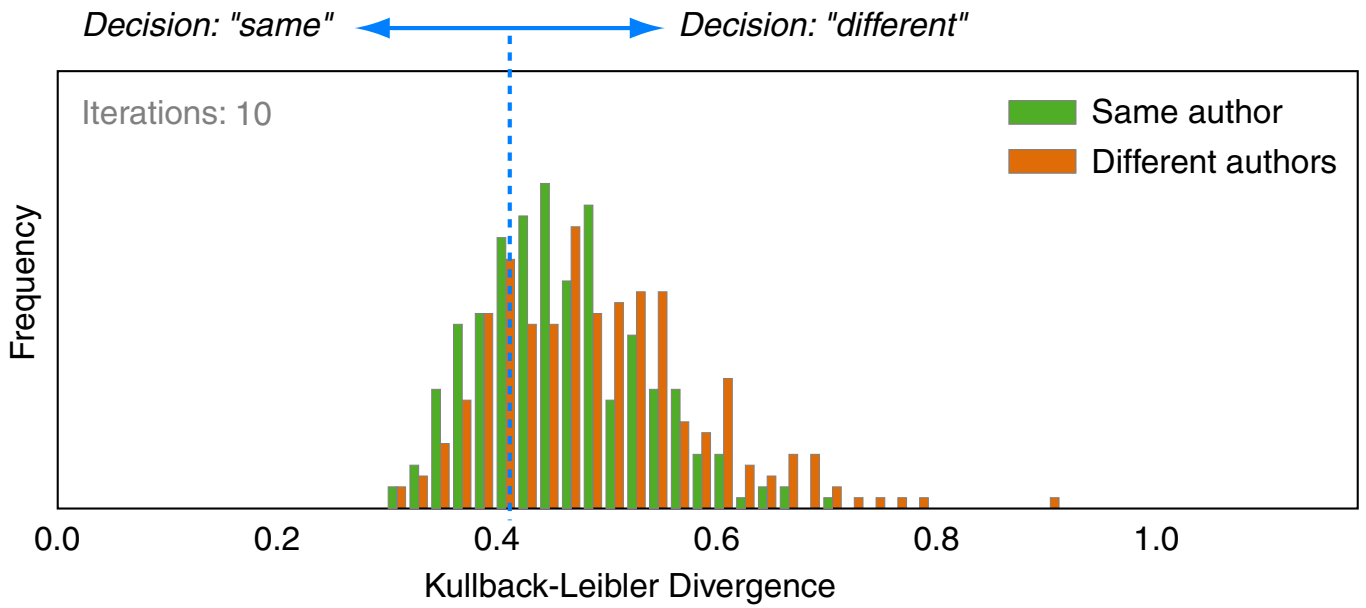


# Countermeasure: Obfuscation

[Bevendorff et al., 2019]

A beautiful\_χchristmas you know jesus our saviour w patiently stooping to hunger and pain, so he mig ones\_ϑfrom shame; now if we love him, he bids u brothers and sisters who need. blessed old nick! i it, you would remember and certainly do it; this you empty your pack, pray give a portion to all w there's anything left and you can bring a small gi wasn't that dandy? sure, little mary\_ϑann has a wor she has! she\_ϑtakes after her own mother. i was jus that age. and you're just like\_ϑher still, mollie mullig

sure, little mary\_ϑann has a wonderful education, s s after her own mother. i was just like her when i wa e just like\_ϑher still, mollie mulligan. sure you're\_ϑth an alley and the belle of shantytown. whist now! it lishes. but, hush! i think the show is about to begi ϑo, samson symbolical! come and see slivers\_ϑclow me and see zip, the foremost of freaks! come an ister sheiks! eager equestriennes\_ϑeach unexcelle enagerie ever beheld, the giant, the fat girl, the lion artists from far-off japan, audacious acrobats sho

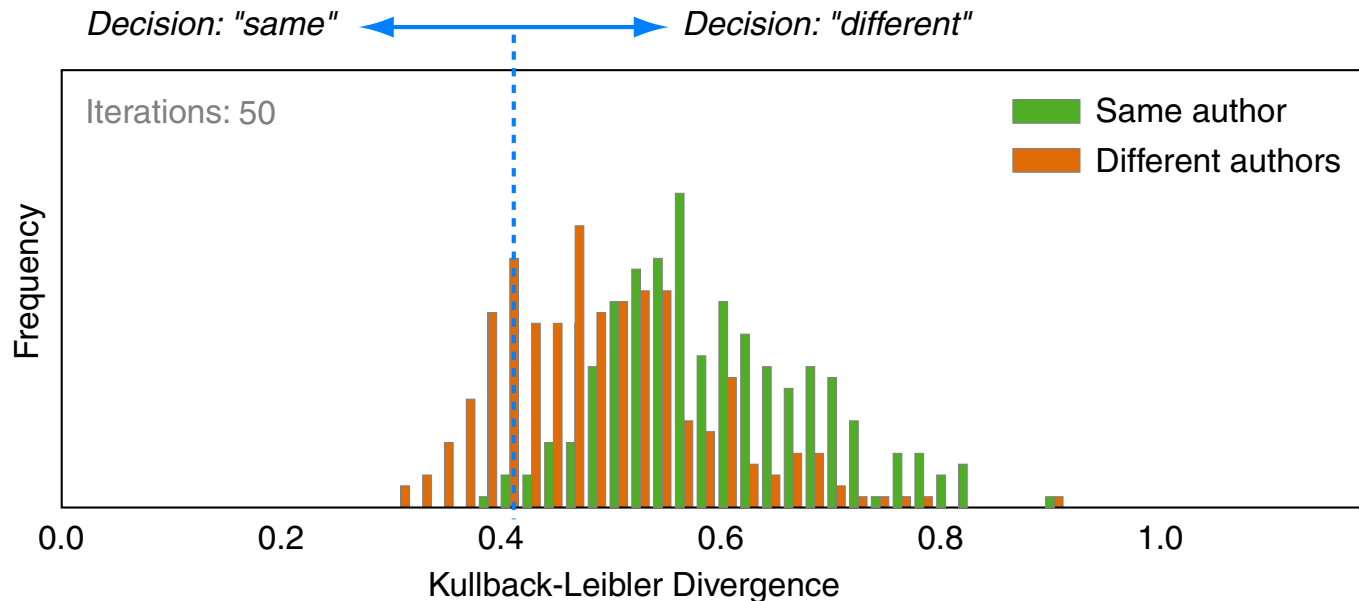


# Countermeasure: Obfuscation

[Bevendorff et al., 2019]

A beautiful\_χchristmas you know jesus our saviour w patiently stooping to hunger and pain, so he mig ones\_s,χfrom shame; now if we love him, he bids u brothers and sisters who need. blessed old nick! i it, you would remember and certainly do it; this you empty your pack, pray give a portion to all w there's anything left and you can bring a small gi wasn't that dandy? sure, little mary\_χann has a wor she has! she\_eχtakes after her own mother. i was jus that age. and you're just like\_χher still, mollie mullig

sure, little mary\_χann has a wonderful education, s s after her own mother. i was just like her when i wa e just like\_χher still, mollie mulligan. sure you're\_eχth an alley and the belle of shantytown. whist now! it lishes. but, hush! i think the show is about to begi ο, samson symbolical! come and see slivers\_s,χclow me and see zip, the foremost of freaks! come an ister sheiks! eager equestriennes\_s,χeach unexcelle enagerie ever beheld, the giant, the fat girl, the lion artists from far-off japan, audacious acrobats sho



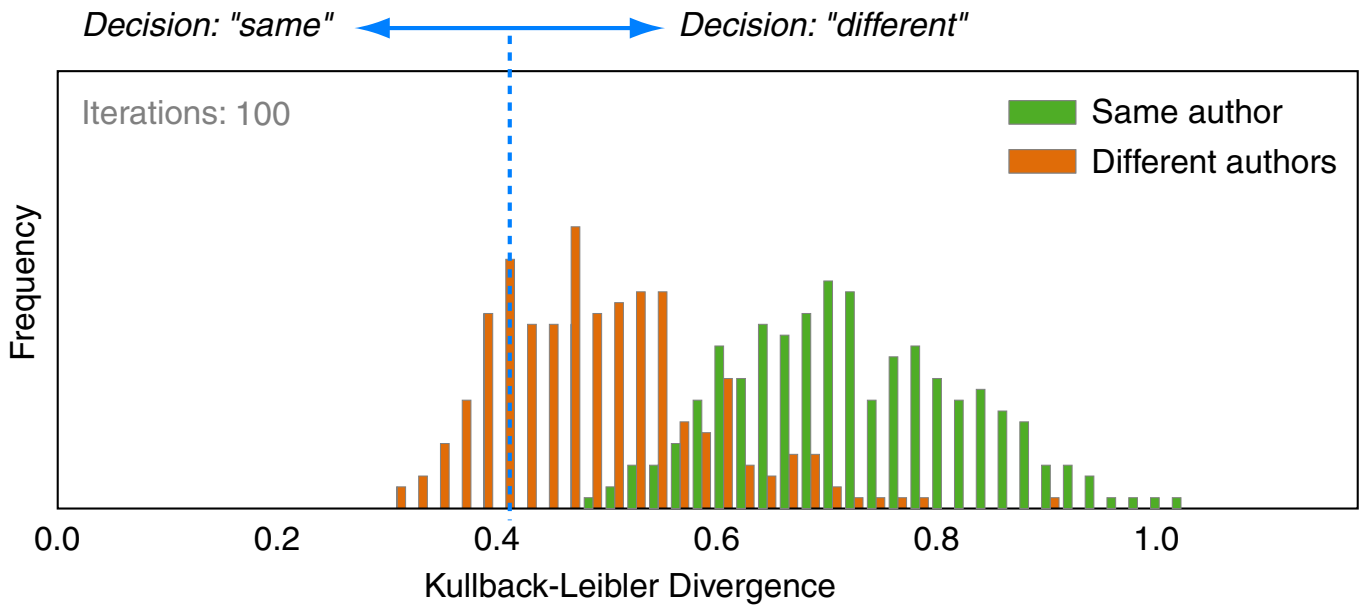


# Countermeasure: Obfuscation

[Bevendorff et al., 2019]

A beautiful\_χchristmas you know jesus our saviour w patiently stooping to hunger and pain, so he mig ones\_χfrom shame; now if we love him, he bids u brothers and sisters who need. blessed old nick! i it, you would remember and certainly do it; this you empty your pack, pray give a portion to all w there's anything left and you can bring a small gi wasn't that dandy? sure, little mary\_χann has a wor she has! she\_χtakes after her own mother. i was jus that age. and you're just like\_χher still, mollie mullig

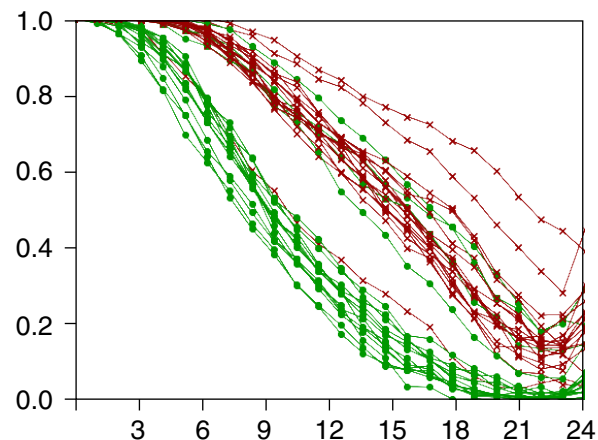
sure, little mary\_χann has a wonderful education, s s after her own mother. i was just like her when i wa e just like\_χher still, mollie mulligan. sure you're\_χth an alley and the belle of shantytown. whist now! it luses. but, hush! i think the show is about to begi ο, samson symbolical! come and see slivers\_χclow me and see zip, the foremost of freaks! come an ister sheiks! eager equestriennes\_χeach unexcelle enagerie ever beheld, the giant, the fat girl, the lion artists from far-off japan, audacious acrobats sho



# Countermeasure: Obfuscation

## Recent Results

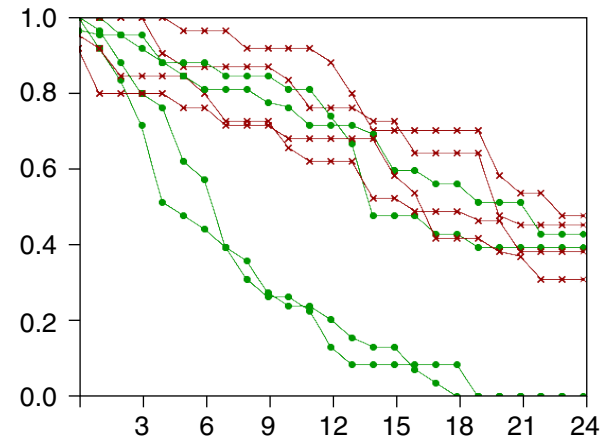
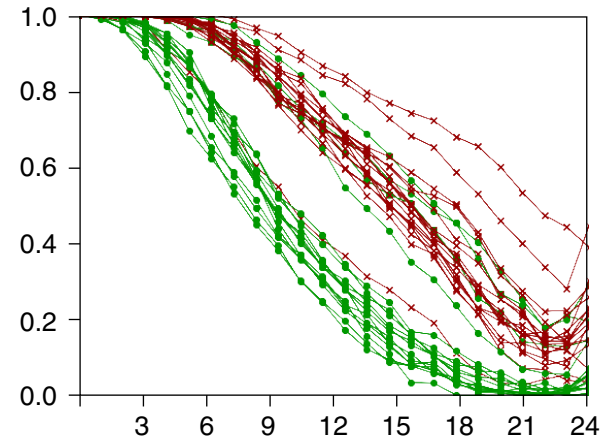
Experiment	I	II	III
<i>Performance</i>			
Precision	1.00	1.00	1.00
Accuracy	0.91	0.75	0.83
n-gram removals	0	80	200
Text coverage	0%	1%	2.6%
Classified	26%	10%	7%
Omitted	74%	90%	93%
<i>Configuration</i>			
Number of cases	180 training / 78 test		
Size of each case	4 000 words		
Number of authors	135		
Number of chunks	25		
Size of each chunk	600 words		
Vocabulary	250 words		
Removed per round	10 words		
Smoothing	no		



# Countermeasure: Obfuscation

## Recent Results

Experiment	I	II	III
<i>Performance</i>			
Precision	1.00	1.00	1.00
Accuracy	0.91	0.75	0.83
n-gram removals	0	80	200
Text coverage	0%	1%	2.6%
Classified	26%	10%	7%
Omitted	74%	90%	93%
<i>Configuration</i>			
Number of cases	180 training / 78 test		
Size of each case	4 000 words		
Number of authors	135		
Number of chunks	25		
Size of each chunk	600 words		
Vocabulary	250 words		
Removed per round	10 words		
Smoothing	no		



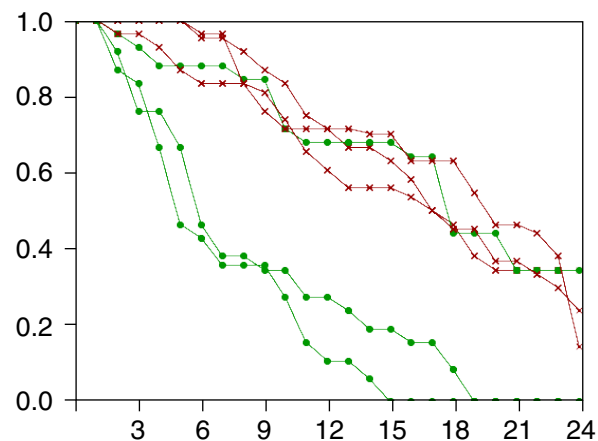
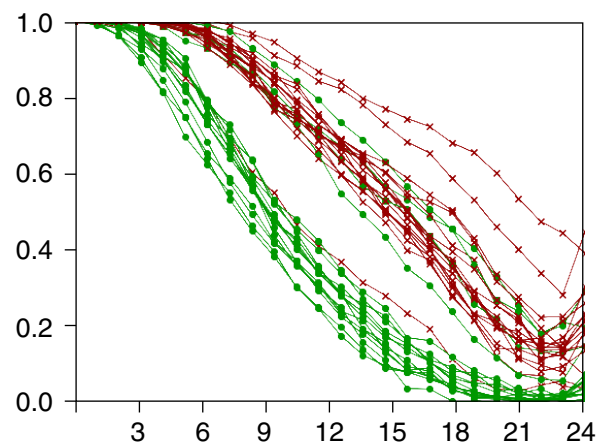
# Countermeasure: Obfuscation

## Recent Results

Experiment	I	II	III
<i>Performance</i>			
Precision	1.00	1.00	1.00
Accuracy	0.91	0.75	0.83
n-gram removals	0	80	200
Text coverage	0%	1%	2.6%
Classified	26%	10%	7%
Omitted	74%	90%	93%

### *Configuration*

Number of cases	180 training / 78 test	
Size of each case	4 000 words	
Number of authors	135	
Number of chunks	25	
Size of each chunk	600 words	
Vocabulary	250 words	
Removed per round	10 words	
Smoothing	no	



③

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship



To the Members of the California State Assembly:

I am returning Assembly Bill 1176 without my signature.

For some time now I have lamented the fact that major issues are overlooked while many unnecessary bills come to me for consideration. Water reform, prison reform, and health care are major issues my Administration has brought to the table, but the Legislature just kicks the can down the alley.

Yet another legislative year has come and gone without the major reforms Californians overwhelmingly deserve. In light of this, and after careful consideration, I believe it is unnecessary to sign this measure at this time.

Sincerely,

Arnold Schwarzenegger

To the Members of the California State Assembly:



I am returning Assembly Bill 1176 without my signature.

For some time now I have lamented the fact that major issues are overlooked while many unnecessary bills come to me for consideration. Water reform, prison reform, and health care are major issues my Administration has brought to the table, but the Legislature just kicks the can down the alley.

Yet another legislative year has come and gone without the major reforms Californians overwhelmingly deserve. In light of this, and after careful consideration, I believe it is unnecessary to sign this measure at this time.

Sincerely,

Arnold Schwarzenegger



To the Members of the California State Assembly:



I am returning Assembly Bill 1176 without my signature.

For some time now I have lamented the fact that major issues are overlooked while many unnecessary bills come to me for consideration. Water reform, prison reform, and health care are major issues my Administration has brought to the table, but the Legislature just kicks the can down the alley.

Yet another legislative year has come and gone without the major reforms Californians overwhelmingly deserve. In light of this, and after careful consideration, I believe it is unnecessary to sign this measure at this time.

Sincerely,

*“My goodness. What a coincidence [...]”*

Arnold

[Aron McLearn, Schwarzenegger spokesman, Oct. 2009]

# Paraphrasing “The Acrostify Benchmark”

*An acrostic is a poem or other form of writing in which the first letter, syllable or word of each line, paragraph or other recurring feature in the text spells out a word or a message.*

[Wikipedia]

**A poem** [Kuperavage 2000] :

**H** He broke my heart

**E** Every piece, shattered

**A** All I wanted was his love

**R** Real, as he promised

**T** True, as mine for him

...

# Paraphrasing “The Acrostify Benchmark”

*An acrostic is a poem or other form of writing in which the first letter, syllable or word of each line, paragraph or other recurring feature in the text spells out a word or a message.*

[Wikipedia]

A poem [Kuperavage 2000] :

**H** He broke my heart  
**E** Every piece, shattered  
**A** All I wanted was his love  
**R** Real, as he promised  
**T** True, as mine for him  
...

## Task

Given: (1) A text  $T$  and an acrostic  $x$ .

(2) Lower and upper bounds on the desired line lengths.

Task: Find a paraphrased version  $T^*$  of  $T$  in monospaced font that encodes  $x$  in some consecutive lines, if possible. Each line of  $T^*$  has to meet the length constraints.

# Paraphrasing “The Acrostify Benchmark”

*An acrostic is a poem or other form of writing in which the first letter, syllable or word of each line, paragraph or other recurring feature in the text spells out a word or a message.*

[Wikipedia]

A poem [Kuperavage 2000] :

**H** He broke my heart

**E** Every piece, shattered

**A complex search problem.**

(that may be tackled with AI technologies)

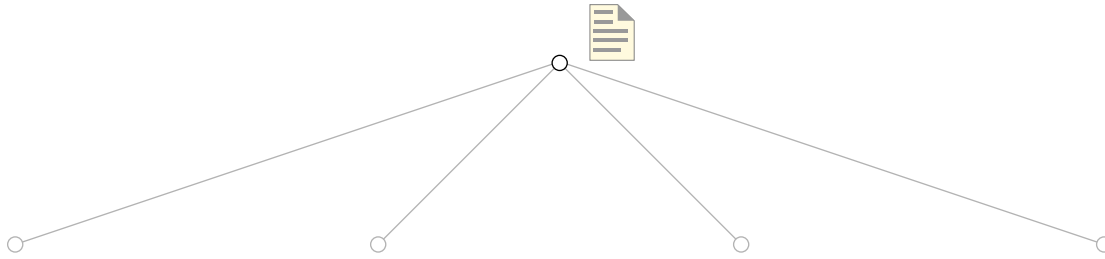
e

## Task

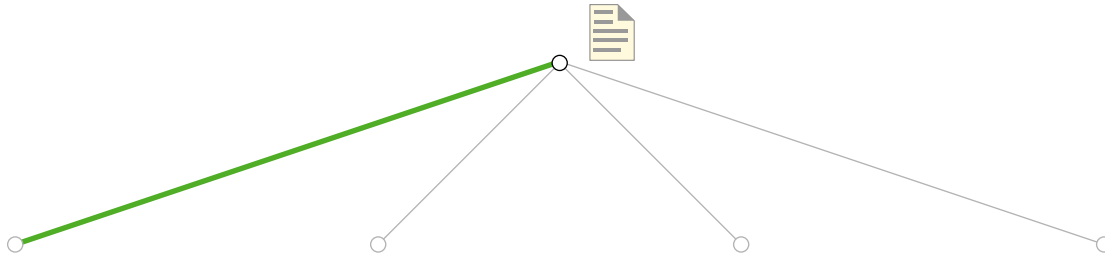
Given: (1) A text  $T$  and an acrostic  $x$ .

(2) Lower and upper bounds on the desired line lengths.

Task: Find a paraphrased version  $T^*$  of  $T$  in monospaced font that encodes  $x$  in some consecutive lines, if possible. Each line of  $T^*$  has to meet the length constraints.



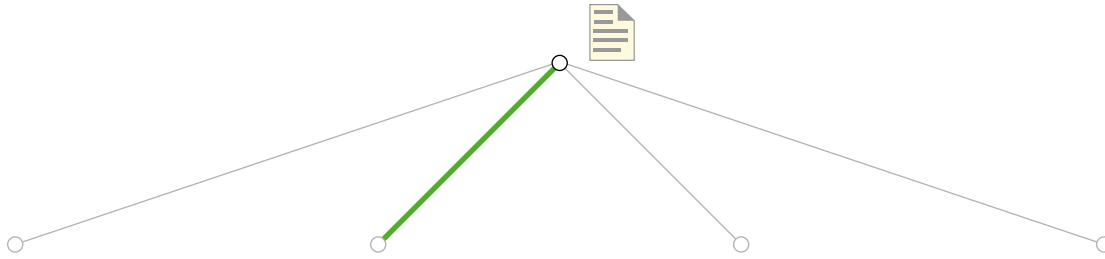
Subtask: Create the character **b**auhaus



**Be**fore some time  
~~now~~ I have  
lamented the  
fact that major  
issues are  
overlooked while  
many bills come  
to

«Preposition»

Subtask: Create the character **b**auhaus



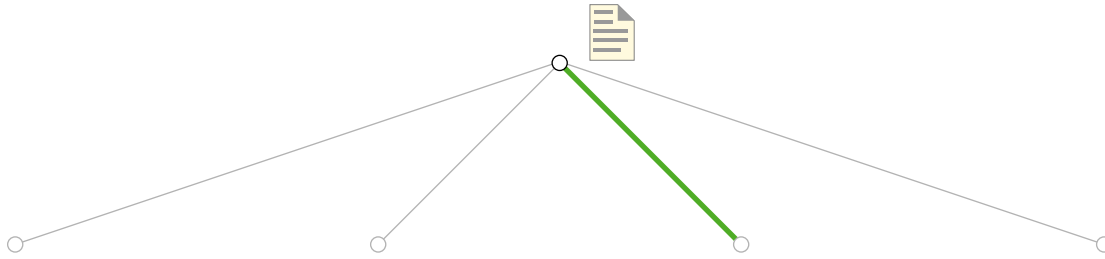
Before some time  
~~now~~ I have  
lamented the  
fact that major  
issues are  
overlooked while  
many bills come  
to

«Preposition»

For some time now  
I have lamented |  
**but** the fact that  
major issues are  
overlooked while  
many bills

«Add Connective»

Subtask: Create the character **b**auhaus



**Before** some time  
~~now~~ I have  
 lamented the  
 fact that major  
 issues are  
 overlooked while  
 many bills come  
 to

«Preposition»

For some time now  
 I have lamented |  
**but** the fact that  
 major issues are  
 overlooked while  
 many bills

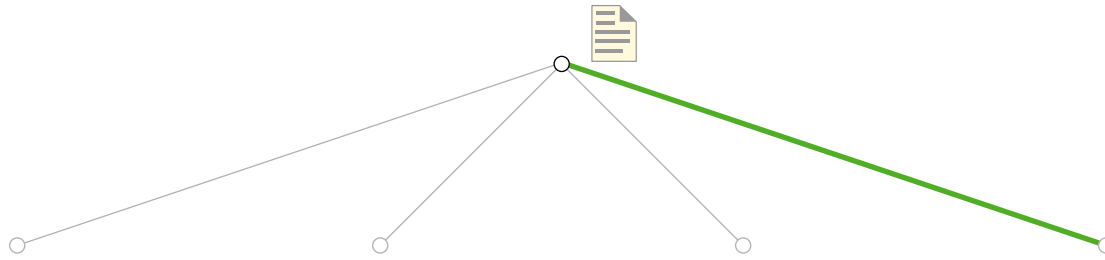
«Add Connective»

**Been** for some  
 time now I have  
 lamented the  
 fact that major  
 issues are  
 overlooked while  
 many bills come  
 to

«Change Tense»

Subtask: Create the character **b**auhaus





**Before** some time  
~~now~~ I have  
 lamented the  
 fact that major  
 issues are  
 overlooked while  
 many bills come  
 to

«Preposition»

For some time now  
 I have lamented |  
**but** the fact that  
 major issues are  
 overlooked while  
 many bills

«Add Connective»

**Been** for some  
 time now I have  
 lamented the  
 fact that major  
 issues are  
 overlooked while  
 many bills come  
 to

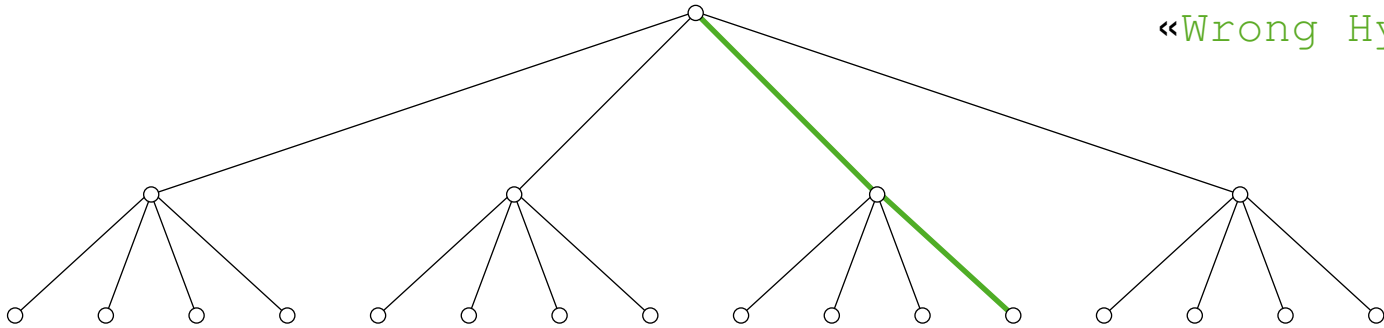
«Change Tense»

For some time now  
 I have lamented  
 the fact that  
 major issues are  
 overlooked while  
 many |  
**b**ills come to

«Linebreak»

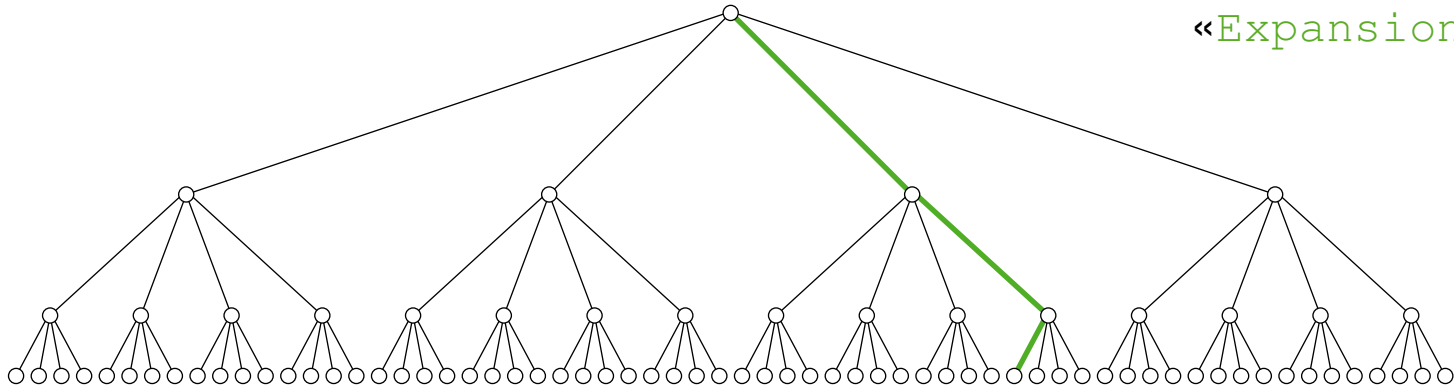
Subtask: Create the character **b**auhaus

«Wrong Hyphen»



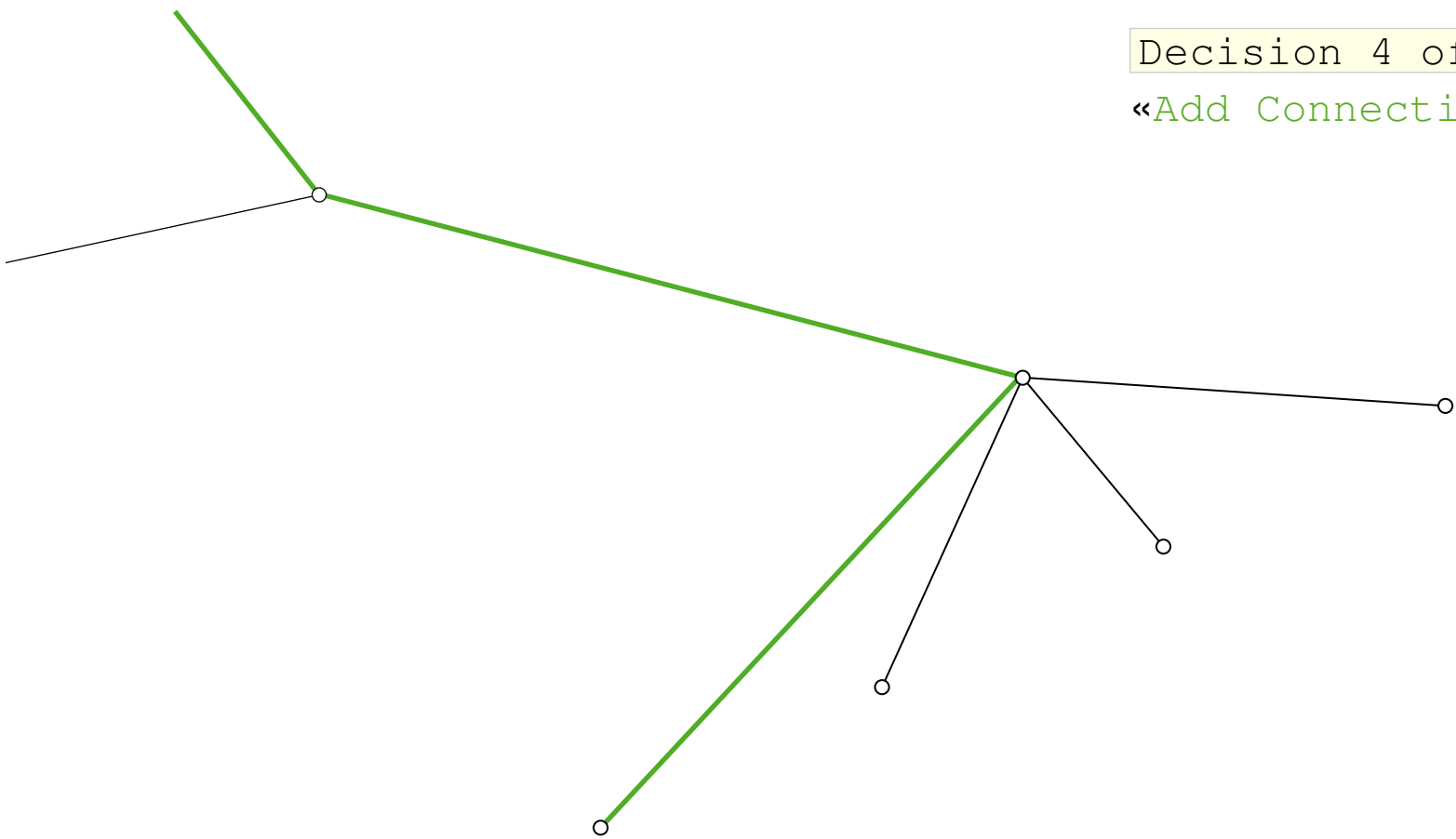
Subtask: Create the character **b**auhaus

«Expansion»



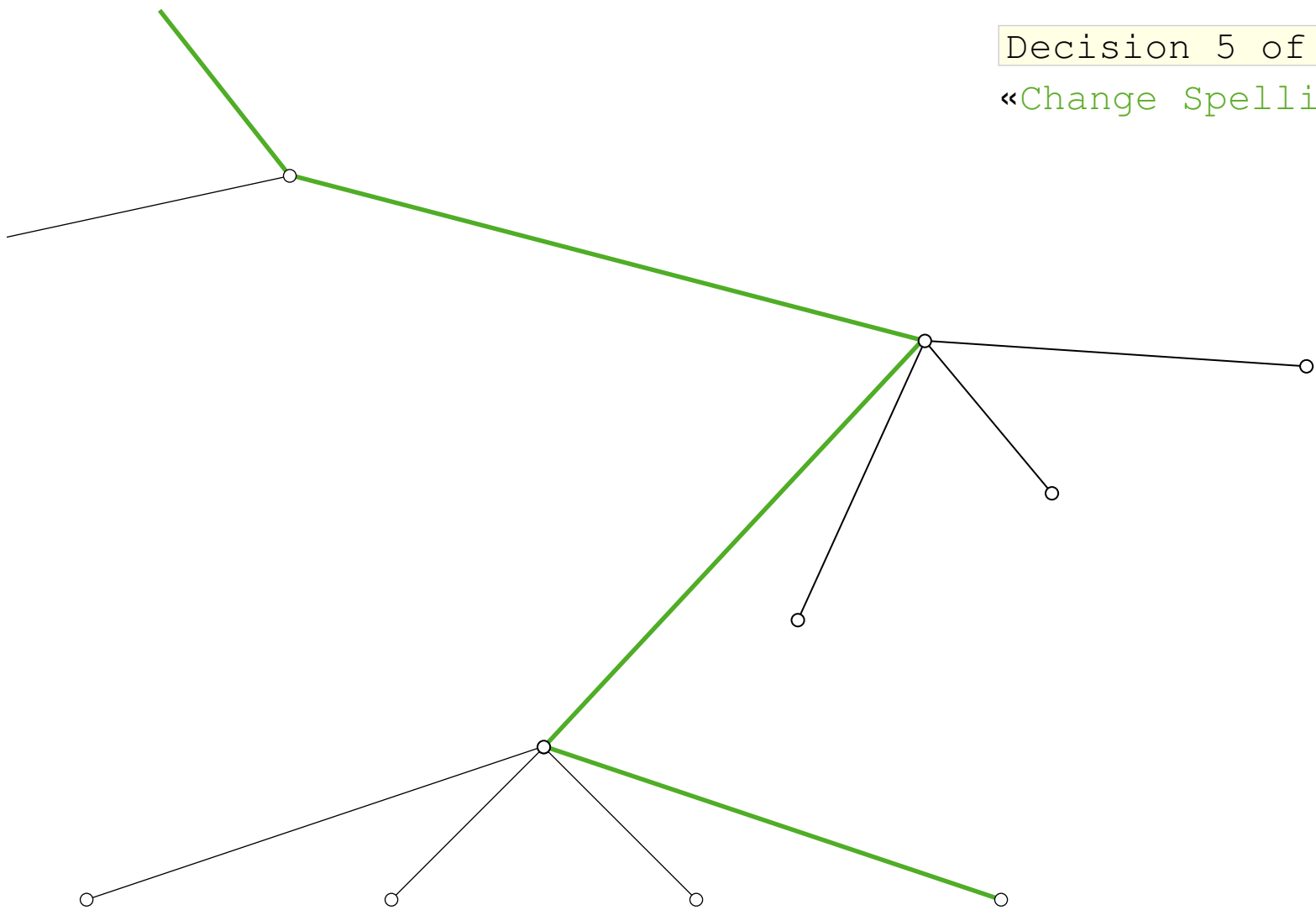
Subtask: Create the character bauhaus

«Add Connective»



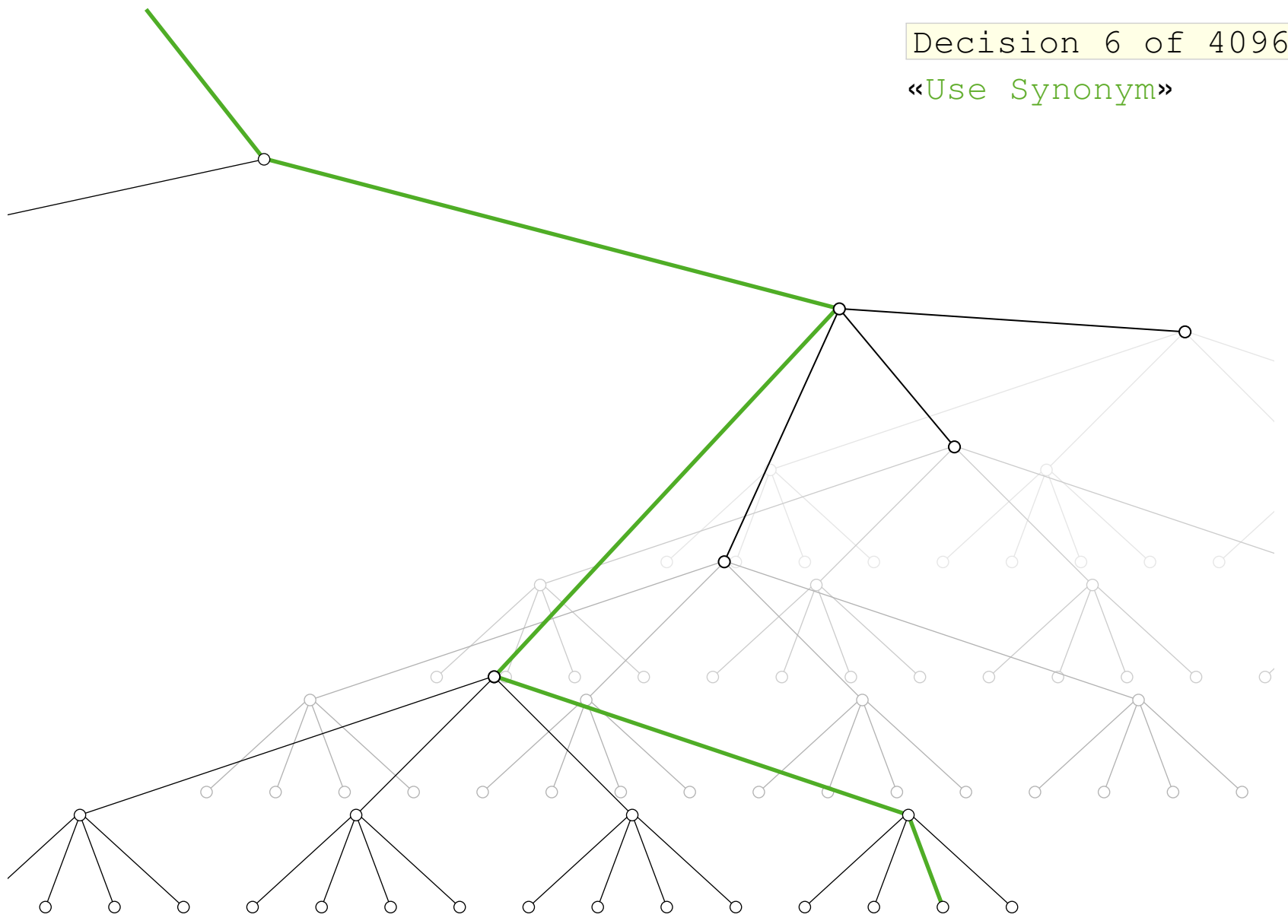
Subtask: Create the character bauhaus

«Change Spelling»



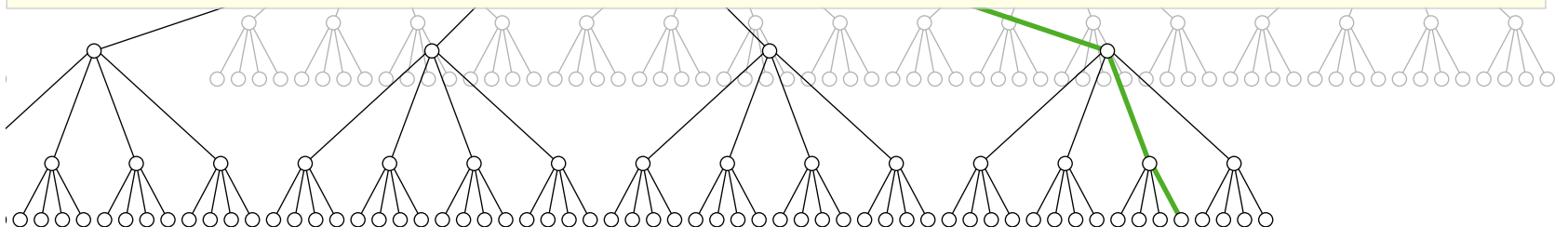
Subtask: Create the character bauhaus

«Use Synonym»



## «Hyphenation»

**B** Been for some time now I have lamented the fact th-  
**a** at major issues are overlooked while many  
**u** unnecessary bills come to me for consideration. [...]  
**h** health care are major issues my Administration [...]  
**a** ature just kicks the can down the alley. Yet [...]  
**u** ut the major reforms Californians overwhelmingly de-  
**s** serve. In light of this, and after careful [...]



# Paraphrasing

Searchspace Facts

Consider a text with a length of 100 words (the Schwarzenegger Letter) . . .

- ≈ 10 · 3 possibilities to change tense
- ≈ 100 possibilities to break a line
- ≈ 100 · 3 possibilities to introduce a synonym
- ≈ 100 · 3 possibilities to introduce filler words
- ≈ 100 · 5 possibilities to hyphenate a word
- ≫ 100 possibilities to introduce tautologies
- . . .



# Paraphrasing

Searchspace Facts

Consider a text with a length of 100 words (the Schwarzenegger Letter) . . .

≈ 10 · 3 possibilities to change tense

≈ 100 possibilities to break a line

≈ 100 · 3 possibilities to introduce a synonym

≈ 100 · 3 possibilities to introduce filler words

≈ 100 · 5 possibilities to hyphenate a word

≫ 100 possibilities to introduce tautologies

. . .

→ > 1 000 possible operations to generate a **single letter** of an acrostic

→  $O(10^{3n})$  possibilities to synthesize an  $n = 7$  letter word like '**Bauhaus**'

Compare the following numbers:

$10^{80}$  atoms in the observable universe

$10^{123}$  game-tree complexity of chess

---

Acrostic type	Length	Runtime	Nodes	Quality-related measures		
	(in letters)	(total in s)	(total)	△ WFC	△ ARI	△ SMOG

---

*Common English words*

Adjective

Noun

Verb

---

*Common US first names*

Male

Female

---

*Self-referential*

First words

---

**Average**

---

Setup details:

- Text genres: Reuters newspaper articles, Enron emails, English Wikipedia articles
- Hardware: standard quad-core PC with 16GB RAM

# Paraphrasing

Selected Results

---

Acrostic type	Length (in letters)	Runtime (total in s)	Nodes (total)	Quality-related measures		
				△ WFC	△ ARI	△ SMOG
<i>Common English words</i>						
Adjective	4.4	3.3	287 000			
Noun	4.8	3.4	285 000			
Verb	3.6	2.8	251 000			
<hr/>						
<i>Common US first names</i>						
Male	6.0	9.3	852 000			
Female	6.1	7.8	740 000			
<hr/>						
<i>Self-referential</i>						
First words	10.3	36.1	3 165 000			
<hr/>						
<b>Average</b>	<b>5.2</b>	<b>8.5</b>	<b>760 000</b>			

---

Setup details:

- Text genres: Reuters newspaper articles, Enron emails, English Wikipedia articles
- Hardware: standard quad-core PC with 16GB RAM

# Paraphrasing

## Selected Results

Acrostic type	Length (in letters)	Runtime (total in s)	Nodes (total)	Quality-related measures		
				$\Delta$ WFC	$\Delta$ ARI	$\Delta$ SMOG
<i>Common English words</i>						
Adjective	4.4	3.3	287 000	-1.0	-1.6	-0.9
Noun	4.8	3.4	285 000	-0.4	-1.0	-0.5
Verb	3.6	2.8	251 000	-1.0	-1.6	-0.9
<i>Common US first names</i>						
Male	6.0	9.3	852 000	-0.7	-1.9	-0.9
Female	6.1	7.8	740 000	-0.6	-1.8	-0.9
<i>Self-referential</i>						
First words	10.3	36.1	3 165 000	-0.3	-0.1	0.2
<b>Average</b>	<b>5.2</b>	<b>8.5</b>	<b>760 000</b>	<b>-0.8</b>	<b>-1.5</b>	<b>-0.8</b>

### Setup details:

- Text genres: Reuters newspaper articles, Enron emails, English Wikipedia articles
- Hardware: standard quad-core PC with 16GB RAM

4

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

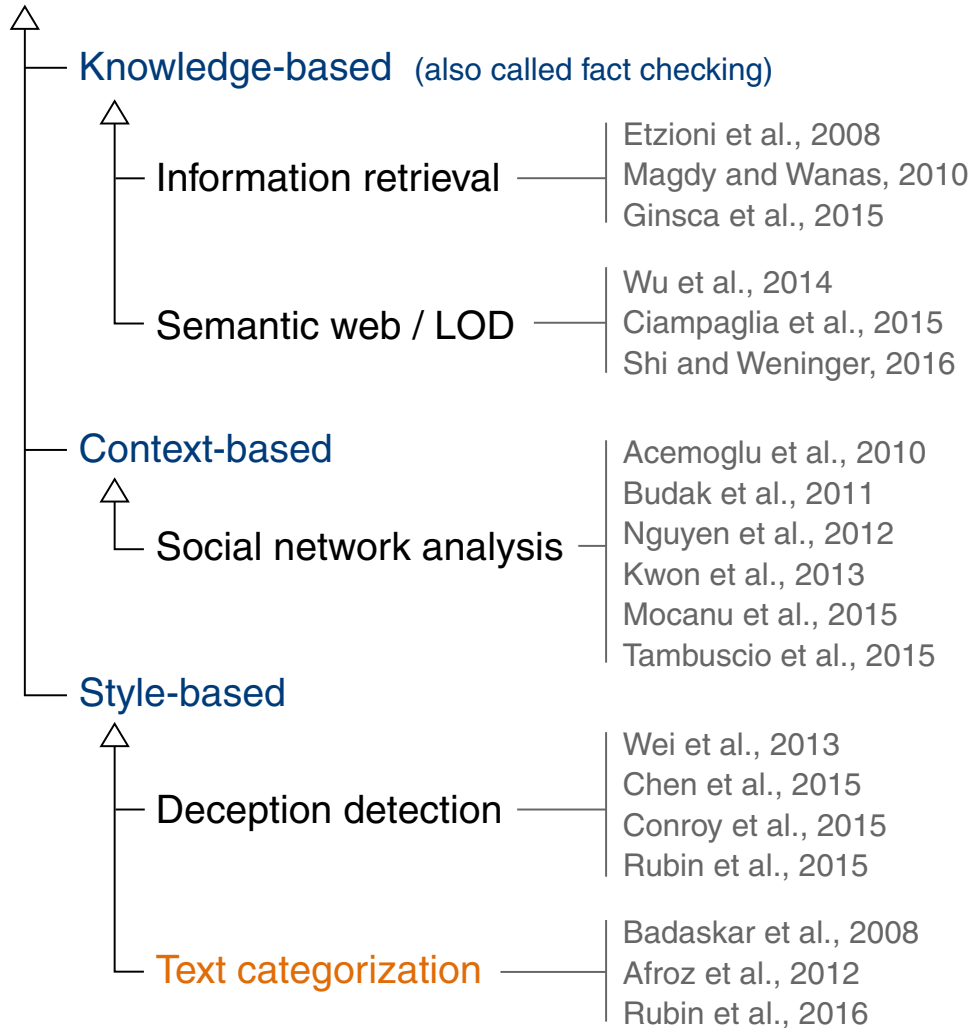
Clickbait  
Text quality  
**Fake News and Hyperpartisanship**  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Fake News and Hyperpartisanship Taxonomy of Approaches

## Fake news detection







# Fake News and Hyperpartisanship Corpus Construction

<b><i>Orientation</i></b>	<b>Fact-checking results</b>				$\Sigma$
	true	mix	false	n/a	
<i>Publisher</i>					
<i>Mainstream</i>	806	8	0	12	826
ABC News	90	2	0	3	95
CNN	295	4	0	8	307
Politico	421	2	0	1	424
<i>Left-wing</i>	182	51	15	8	256
Addicting Info	95	25	8	7	135
Occupy Democrats	59	25	7	0	91
The Other 98%	28	1	0	1	30
<i>Right-wing</i>	276	153	72	44	545
Eagle Rising	106	47	25	36	214
Freedom Daily	49	24	22	4	99
Right Wing News	121	82	25	4	232
$\Sigma$	1264	212	87	64	1627

Annotations provided by journalists at BuzzFeed

# Fake News and Hyperpartisanship Selected Results

<i><b>Orientation</b></i>	<b>Fact-checking results</b>				
	true	mix	false	n/a	$\Sigma$
<i>Mainstream</i>	806	8	0	12	826
ABC News				3	95
CNN				8	307
Politico				1	424
<i>Left-wing</i>				8	256
Addicting Info				7	135
Occupy Democr				0	91
The Other 98%				1	30
<i>Right-wing</i>	276	153	72	44	545
Eagle Rising	106	47	25	36	214
Freedom Daily	49	24	22	4	99
Right Wing News	121	82	25	4	232
$\Sigma$	1264	212	87	64	1627

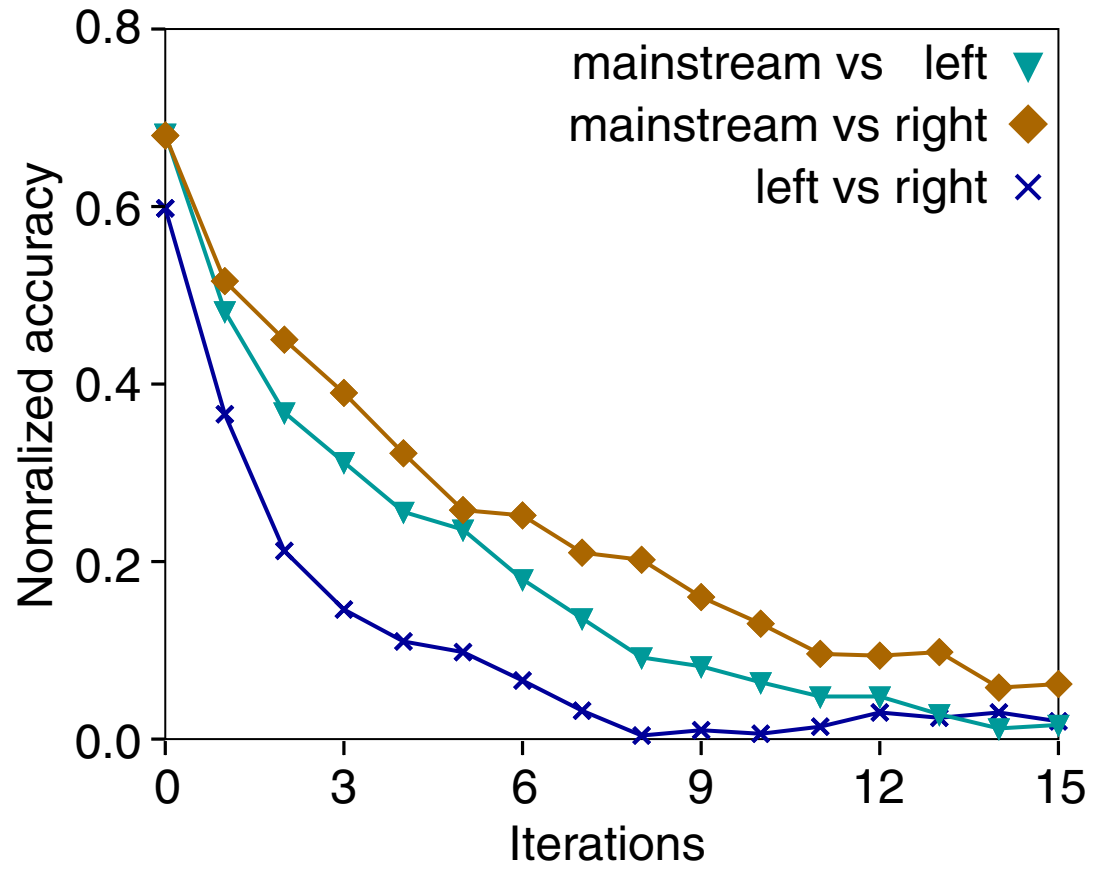
**Fake News Detection**  
 Precision  $\approx$  42%  
 Recall  $\approx$  41%

# Fake News and Hyperpartisanship Selected Results

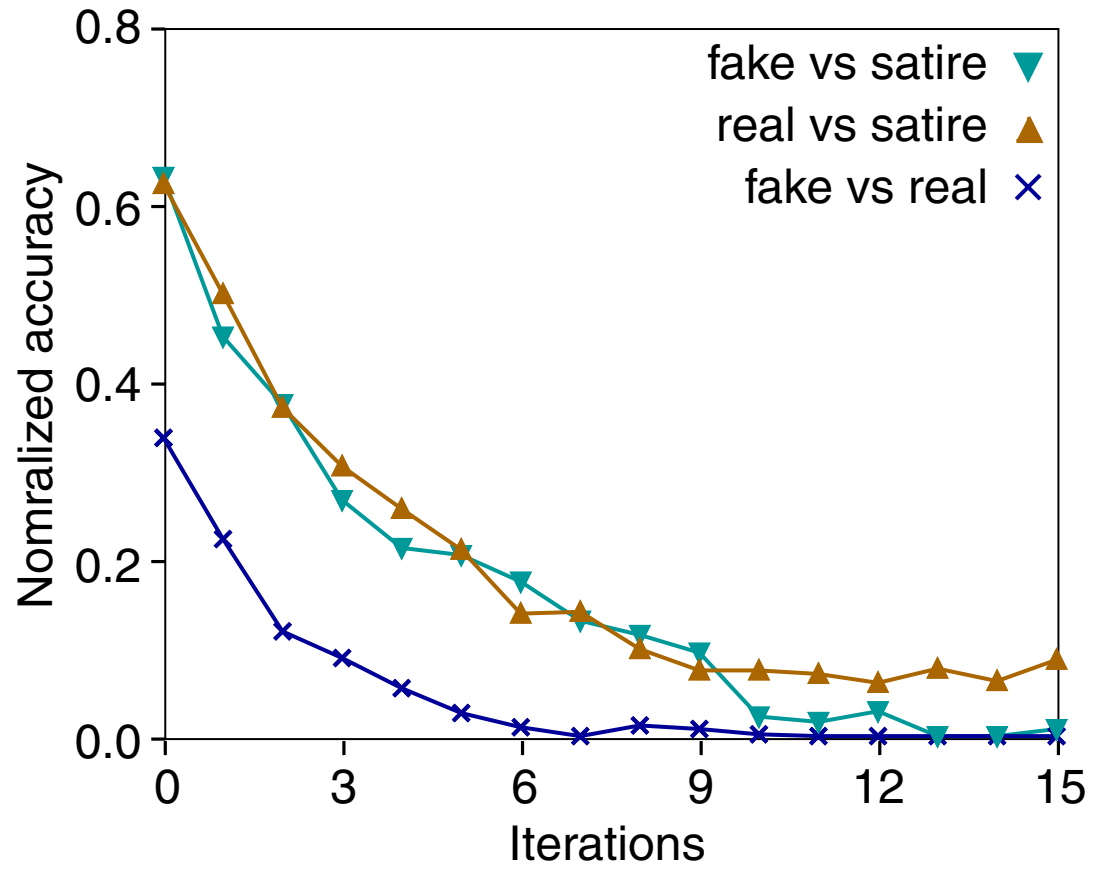
<b>Orientation</b>	<b>Fact-checking results</b>				
	true	mix	false	n/a	$\Sigma$
<i>Mainstream</i>	806	8	0	12	826
ABC News					95
CNN					307
Politico					424
<i>Left-wing</i>					256
Addicting In					135
Occupy Der					91
The Other 9					30
<i>Right-wing</i>	276	153	72	44	545
Eagle Rising	106	47	25	36	214
Freedom Daily	49	24	22	4	99
Right Wing News	121	82	25	4	232
$\Sigma$	1264	212	87	64	1627

**Hyperpartisanship Detection**  
 Precision  $\approx$  69%  
 Recall  $\approx$  89%

# Fake News and Hyperpartisanship Unmasking Orientation



# Fake News and Hyperpartisanship Unmasking Satire



5

# Challenges



## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation

## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship



“You won’t believe how far cats can count.”



welt

WELT

6. Oktober um 23:13 · 🌐

So hoch soll die natürliche Obergrenze für euer Alter liegen.



Wie alt Menschen höchstens werden können

WELT.DE



Huffington Post Deutschland

23 Std. · 🌐

Ob sie den Job des Geldes wegen macht?



## Nach Gehaltserhöhung: So viel verdient Angela Merkel jetzt

Es gibt wohl begehrtere Jobs als den des Bundeskanzlers.

HUFFINGTONPOST.DE

Source: BildBlog (<http://www.bildblog.de/ressort/fuer-sie-geklickt/>); Spoiler: 226.000 Euro p.a.



**BuzzFeed** ✓  
@BuzzFeed

 Follow



Hier ist der absolut genialste Weg, ein schlechtes Pokémon-Tattoo zu retten!

[bzfd.it/1C3yToz](https://bzfd.it/1C3yToz)



Source: Twitter @BuzzFeed (<http://www.twitter.com/buzzfeed/>); Spoiler: No clickbait

Register	Publisher*	Impact (retweets in 2015)	Tweets (in week 24)	Clickbait probability
Print + online	New York Times	23.8 · 10 <sup>6</sup>	875	21%
	The Guardian	14.0 · 10 <sup>6</sup>	744	15%
	Forbes	11.5 · 10 <sup>6</sup>	721	38%
	Daily Mail	6.9 · 10 <sup>6</sup>	516	22%
	Wall Street Journal	6.5 · 10 <sup>6</sup>	747	19%
Online only	Mashable	20.6 · 10 <sup>6</sup>	803	33%
	Huffington Post	11.6 · 10 <sup>6</sup>	770	46%
	Bleacher Report	10.2 · 10 <sup>6</sup>	196	9%
	BuzzFeed	10.0 · 10 <sup>6</sup>	695	42%
	Yahoo!	8.2 · 10 <sup>6</sup>	195	23%
Television	BBC News	39.6 · 10 <sup>6</sup>	694	17%
	ABC News	17.6 · 10 <sup>6</sup>	279	9%
	CNN	15.0 · 10 <sup>6</sup>	345	17%
	Fox News	10.2 · 10 <sup>6</sup>	378	8%
	NBC News	9.7 · 10 <sup>6</sup>	408	14%

Average: 28%

\* Top publishers on Twitter in 2014.

Register	Publisher*	Impact (retweets in 2015)	Tweets (in week 24)	Clickbait probability
Print + online	New York Times	$23.8 \cdot 10^6$	875	21%
	The Guardian	$14.0 \cdot 10^6$	744	15%
	Forbes	$11.5 \cdot 10^6$	721	38%
	Daily Mail	$6.9 \cdot 10^6$	516	22%
	We			19%
Online only	Ma			33%
	Hu			46%
	Ble			9%
	Bu			42%
	Yal			23%
Television	BB			17%
	AB			9%
	CNN	$15.0 \cdot 10^6$	345	17%
	Fox News	$10.2 \cdot 10^6$	378	8%
	NBC News	$9.7 \cdot 10^6$	408	14%

**Clickbait detection\*\***  
  
 Precision  $\approx$  71%  
 Recall  $\approx$  73%

Average: 28%

\* Top publishers on Twitter in 2014.

\*\* [Potthast et al., ECIR'16]

# Clickbait A Challenge!



<http://www.clickbait-challenge.org/>

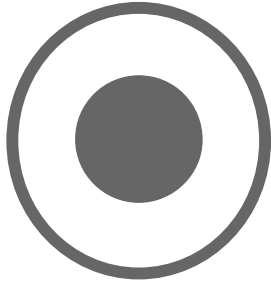
A challenge is organized to encourage researchers in this.

---

Corpus size	40 000 tweets
Votes per tweet	5
Votes per “check instance”	> 60
Number of AMT workers	3 500

---

Dedicated acquisition technology and statistics.



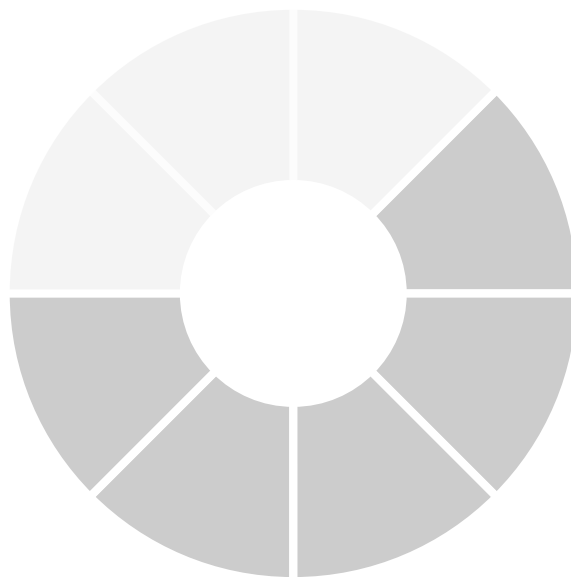
# Summary

## Web as Corpus

Mnemonic passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation



## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and  
Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship



# Summary

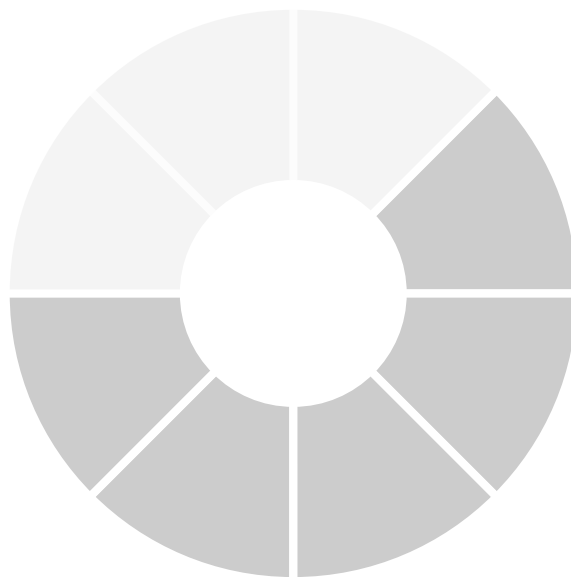
## Web as Corpus

Mnemonic  
passwords

## Synthesis

### Summarization

Paraphrasing  
Obfuscation



## Search

### Question queries

Axiomatic re-ranking  
Argument search

## Assessment

### Clickbait

Text quality  
Fake News and  
Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

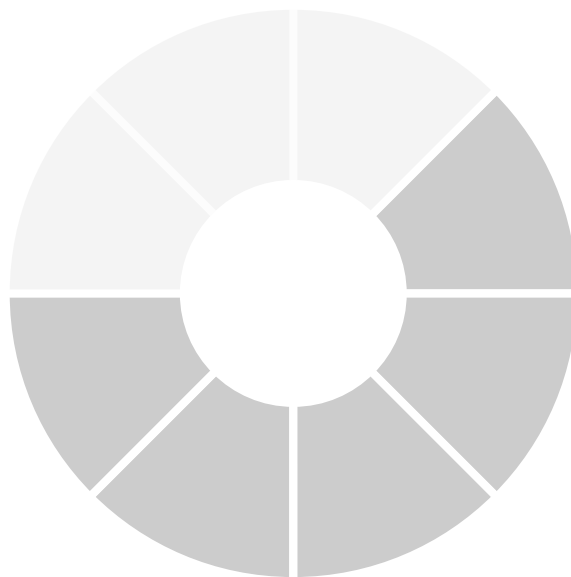
# Summary

## Web as Corpus

Mnemonic  
passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation



## Search

Question queries  
Axiomatic re-ranking  
**Argument search**

## Assessment

Clickbait  
Text quality  
**Fake News and  
Hyperpartisanship**  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Summary

## Web as Corpus

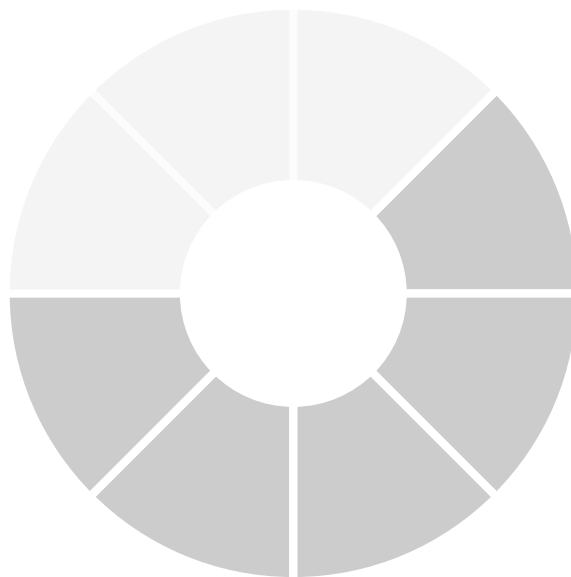
Mnemonic  
passwords

## Synthesis

Summarization

Paraphrasing

**Obfuscation**



## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and  
Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism

**Authorship**

# Summary

## Web as Corpus

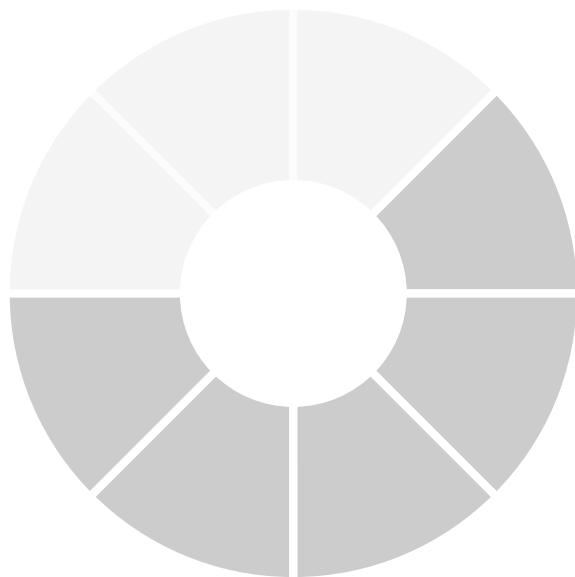
Mnemonic  
passwords

## Synthesis

Summarization

**Paraphrasing**

Obfuscation



## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and  
Hyperpartisanship

**Offensive language**

## Detection

Vandalism  
Plagiarism  
Authorship

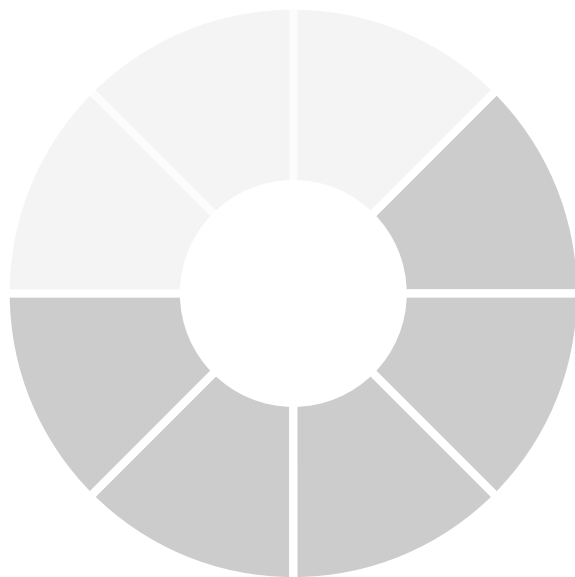
# Summary

## Web as Corpus

Mnemonic  
passwords

## Synthesis

Summarization  
Paraphrasing  
Obfuscation



## Search

Question queries  
Axiomatic re-ranking  
Argument search

## Assessment

Clickbait  
Text quality  
Fake News and  
Hyperpartisanship  
Offensive language

## Detection

Vandalism  
Plagiarism  
Authorship

# Thank you!



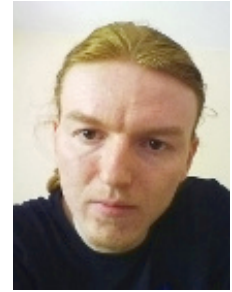
Benno Stein



Matthias Hagen



Henning Wachsmuth



Janek Bevendorff



Michael Völske

Thank you!



Johannes Kiesel



Khalid Al-Khatib



Tim Gollub



Shahbaz Syed



Yamen Ajjour