# Harnessing Web Archives
# to Tackle Selected Societal Challenges

The Oral Exam of
**Johannes Kiesel**

To Obtain the Academic Degree of
**Dr. rer. nat.**

Web Technology & Information Systems Group
Bauhaus-Universität Weimar

www.uni-weimar.de                    www.webis.de

# Harnessing Web Archives to Tackle Selected Societal Challenges

## Societal challenges

*Issues that concern most if not all members of a society, either now or in a likely future.*

Well-known challenges:*

❑ Critical assessment of information

❑ Protection of the environment

❑ Preservation of culture

❑ Ensuring public health

❑ Security and privacy

# Harnessing Web Archives to Tackle Selected Societal Challenges
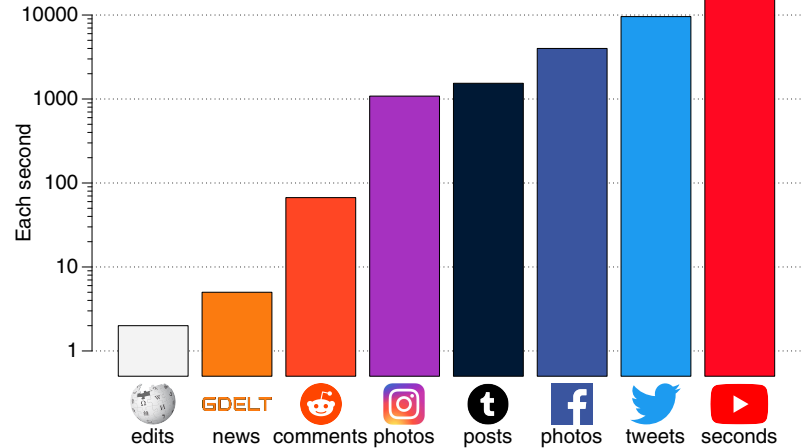
## Societal challenges

*Issues that concern most if not all members of a society, either now or in a likely future.*

Well-known challenges:*

- ❑ Critical assessment of information
- ❑ Protection of the environment
- ❑ Preservation of culture
- ❑ Ensuring public health
- ❑ Security and privacy

*Taken from European Commission (Horizon 2020), World Economic Forum, Gesellschaft für Informatik

## Society → Web



Source: DOMO, Reddit, GDELT, Wikipedia

## Web archives

- ❑ Allow for large-scale analyses
- ❑ Allow to trace changes
- ❑ Allow to replicate analyses

# Harnessing Web Archives to Tackle Selected Societal Challenges

## Main contributions

### 1. Preservation of digital culture

❑ 10K pages high-fidelity archive (FAIRest dataset award)

❑ Reproduction assessment task

❑ 9K pages segmentation dataset

❑ Segmentation evaluation measures

### 2. Critical assessment of information

❑ Revert-based vandalism detection

❑ 30K edits Wiki vandalism dataset

❑ 1M hyperpartisan news dataset

❑ Style-based polarity detection

❑ Hyperpartisan news challenge (SemEval, 42 teams)

### 3. Online security and privacy

❑ 3B web sentences dataset

❑ Position-dependent language model

❑ Security estimate: mnemonic passwords

❑ Personal archiving tool

**Tailored web archiving technology**  (Webis Web Archiver)

# Harnessing Web Archives to Tackle Selected Societal Challenges
## Main contributions

## 1. Preservation of digital culture

- ❑ 10K pages high-fidelity archive (FAIRest dataset award)

- ❑ **Reproduction assessment task**

- ❑ 9K pages segmentation dataset

- ❑ Segmentation evaluation measures

## 2. Critical assessment of information

- ❑ Revert-based vandalism detection

- ❑ 30K edits Wiki vandalism dataset

- ❑ 1M hyperpartisan news dataset

- ❑ Style-based polarity detection

- ❑ **Hyperpartisan news challenge (SemEval, 42 teams)**

## 3. Online security and privacy

- ❑ 3B web sentences dataset

- ❑ Position-dependent language model

- ❑ Security estimate: mnemonic passwords

- ❑ Personal archiving tool

**Tailored web archiving technology**  (Webis Web Archiver)

→ **New tasks**

# Harnessing Web Archives to Tackle Selected Societal Challenges
## Main contributions

**1. Preservation of digital culture**

❏ 10K pages high-fidelity archive (FAIRest dataset award)

❏ Reproduction assessment task

❏ 9K pages segmentation dataset

❏ Segmentation evaluation measures

**2. Critical assessment of information**

❏ Revert-based vandalism detection

❏ 30K edits Wiki vandalism dataset

❏ 1M hyperpartisan news dataset

❏ Style-based polarity detection

❏ Hyperpartisan news challenge (SemEval, 42 teams)

**3. Online security and privacy**

❏ 3B web sentences dataset

❏ Position-dependent language model

❏ Security estimate: mnemonic passwords

❏ Personal archiving tool

**Tailored web archiving technology**  (Webis Web Archiver)

→ New tasks      → New or improved algorithms

# Harnessing Web Archives to Tackle Selected Societal Challenges
## Main contributions

### 1. Preservation of digital culture

- ❑ 10K pages high-fidelity archive (FAIRest dataset award)

- ❑ Reproduction assessment task

- ❑ 9K pages segmentation dataset

- ❑ **Segmentation evaluation measures**

### 2. Critical assessment of information

- ❑ Revert-based vandalism detection

- ❑ 30K edits Wiki vandalism dataset

- ❑ 1M hyperpartisan news dataset

- ❑ Style-based polarity detection

- ❑ Hyperpartisan news challenge (SemEval, 42 teams)

### 3. Online security and privacy

- ❑ 3B web sentences dataset

- ❑ Position-dependent language model

- ❑ Security estimate: mnemonic passwords

- ❑ Personal archiving tool

**Tailored web archiving technology**  (Webis Web Archiver)

→ New tasks     → New or improved algorithms     → More adequate evaluation measures
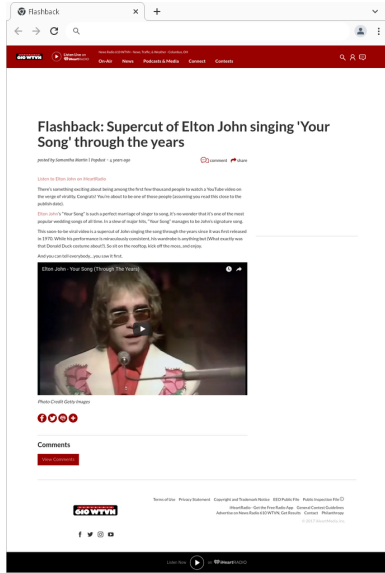
# Harnessing Web Archives to Tackle Selected Societal Challenges
## Main contributions

### 1. Preservation of digital culture

- ❑ 10K pages high-fidelity archive (FAIRest dataset award)
- ❑ Reproduction assessment task
- ❑ 9K pages segmentation dataset
- ❑ Segmentation evaluation measures

### 2. Critical assessment of information

- ❑ Revert-based vandalism detection
- ❑ 30K edits Wiki vandalism dataset
- ❑ 1M hyperpartisan news dataset
- ❑ Style-based polarity detection
- ❑ Hyperpartisan news challenge (SemEval, 42 teams)

### 3. Online security and privacy

- ❑ 3B web sentences dataset
- ❑ Position-dependent language model
- ❑ Security estimate: mnemonic passwords
- ❑ Personal archiving tool

## Tailored web archiving technology  (Webis Web Archiver)

→ New tasks        → New or improved algorithms        → More adequate evaluation measures        → Larger and more accurate datasets

Challenge 1
Preservation of Digital Culture

# Web Page Segmentation

(highlighting reproducibility)

# Web Page Segmentation

# Web Page Segmentation

# Web Page Segmentation
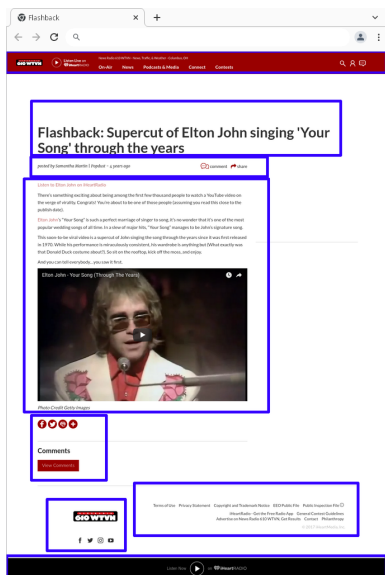
# Web Page Segmentation
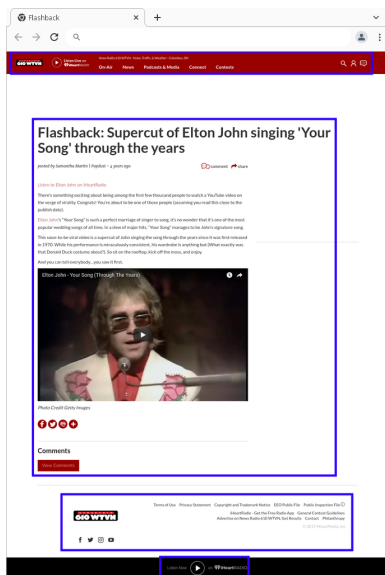


Visually distinct segments



Self-contained segments

# Web Page Segmentation
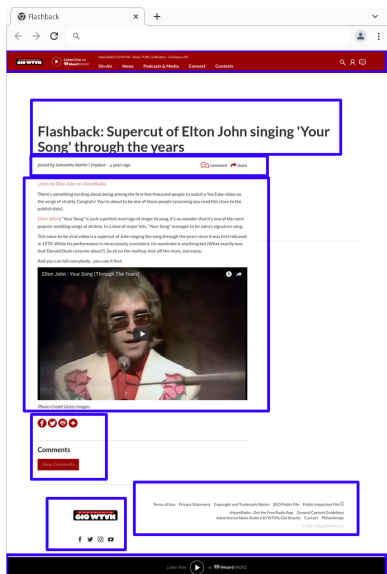


Visually distinct segments
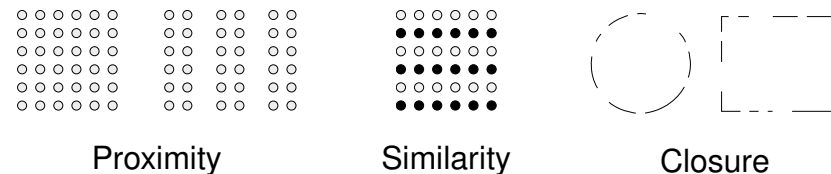


Self-contained segments

Existing definitions (9): biased towards downstream tasks

- ❑ Segments are visual blocks (4), edge-delineated (2), visually distinct (1), self-contained (1), have a heading (1)

- → Problem: inconsistent evaluation methodology

- → No reliable benchmark of algorithms

Existing datasets (20): not re-usable

- ❑ The 12 with human annotations are small (max 1000 pages)

- ❑ Only 3 of these allow for algorithms based on computer vision

- ❑ None allow to reproduce page for browser-based algorithms

# Web Page Segmentation



Visually distinct segments



Self-contained segments

## Existing definitions (9): biased towards downstream tasks

- ❑ Segments are visual blocks (4), edge-delineated (2), visually distinct (1), self-contained (1), have a heading (1)

- → Problem: inconsistent evaluation methodology

- → No reliable benchmark of algorithms

## Existing datasets (20): not re-usable

- ❑ The 12 with human annotations are small (max 1000 pages)

- ❑ Only 3 of these allow for algorithms based on computer vision

- ❑ None allow to reproduce page for browser-based algorithms

## Solution

- ❑ Segment concept based on human viewer (Gestalt principles)

- ❑ Dataset of 8490 archived web pages
  (5 segmentations each; reproducible in browser)

- ❑ Segmentation fusion method

- ❑ Evaluation measure, tweakable towards downstream tasks

## Gestalt principles (selection)



Proximity          Similarity          Closure

A web page segment is a part of a web page
containing those elements that belong together
as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

Large-scale human annotation (8490 pages $\times$ 5)



Example: Drawing rectangles on screenshot. Usage: 1. Click 2. Draw rectangle 3. Click.

$\rightarrow$ Annotation of 600,000 segments in 4 months of full-time work

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

Ground-truth fusion: hierarchical clustering (UPGMA)

Large-scale human annotation (8490 pages $\times$ 5)



$\rightarrow$   Annotation of 600,000 segments in 4 months of full-time work

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

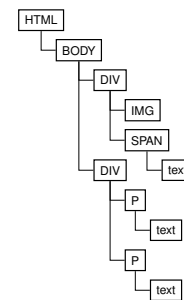Ground-truth fusion: hierarchical clustering (UPGMA)

Evaluation: $F_{B^3} \in [0, 1]$ (from clustering evaluation)

$\rightarrow$ Decomposition into $P_{B^3}$, $R_{B^3}$
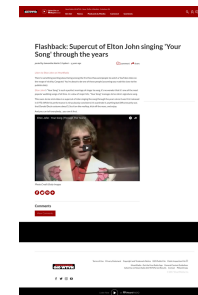  $\approx$ errors of oversegmentation, undersegmentation

Large-scale human annotation (8490 pages $\times$ 5)



$\rightarrow$ Annotation of 600,000 segments in 4 months of full-time work

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

Ground-truth fusion: hierarchical clustering (UPGMA)

Evaluation: $F_{B^3} \in [0, 1]$ (from clustering evaluation)

$\rightarrow$ Decomposition into $P_{B^3}$, $R_{B^3}$
   $\approx$ errors of oversegmentation, undersegmentation

Elements of downstream tasks



Listen Live on iHeartRADIO News Radio 610WTVN-News, Traffic, Weather - Columbus, OH On-Air News Podcasts Media Connect Contests Flashback: Supercur of Elton John singing 'Your Song' through the years posted by Samantha Martin | Popdust - 4 years ago comment share Listen to Elton John on iHeartRadio The

Characters      DOM nodes      Pixels      Edge pixels

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

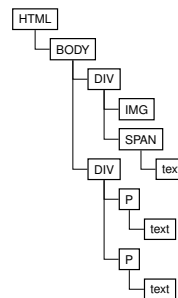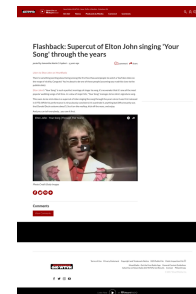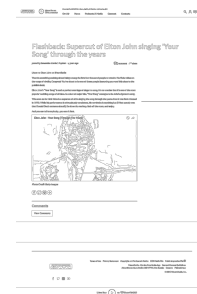Ground-truth fusion: hierarchical clustering (UPGMA)

Evaluation: $F_{B^3} \in [0, 1]$ (from clustering evaluation)

$\rightarrow$ Decomposition into $P_{B^3}$, $R_{B^3}$
   $\approx$ errors of oversegmentation, undersegmentation

Elements of downstream tasks



Characters          DOM nodes          Pixels          Edge pixels

High agreement for all tasks

| Agreement measure | Characters | Nodes | Pixels | Edge pixels |
| --- | --- | --- | --- | --- |
| $F_{B^3}$ | 0.78 | 0.74 | 0.65 | 0.73 |
| $\max(P_{B^3}, R_{B^3})$ | 0.97 | 0.95 | 0.94 | 0.96 |

A web page segment is a part of a web page containing those elements that belong together as per agreement among a majority of viewers.

Elements $E = \{e_1, \ldots, e_n\}$

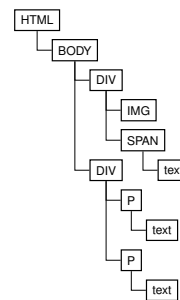Segmentation $S = \{s_1, \ldots, s_m\}$ with segments $s_i \subseteq E$

Ground-truth fusion: hierarchical clustering (UPGMA)

Evaluation: $F_{B^3} \in [0, 1]$ (from clustering evaluation)

$\rightarrow$ Decomposition into $P_{B^3}$, $R_{B^3}$
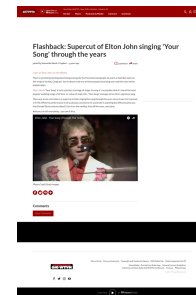$\approx$ errors of oversegmentation, undersegmentation

## Elements of downstream tasks



| Characters | DOM nodes | Pixels | Edge pixels |

## High agreement for all tasks

| Agreement measure | Characters | Nodes | Pixels | Edge pixels |
|---|---|---|---|---|
| $F_{B^3}$ | 0.78 | 0.74 | 0.65 | 0.73 |
| $\max(P_{B^3}, R_{B^3})$ | 0.97 | 0.95 | 0.94 | 0.96 |

## Insights into segmentation technology ($F_{B^3}$)

| Elements/task | 1Seg | VIPS | HEPS | Cor. | MMD. | Meier | MV@2 |
|---|---|---|---|---|---|---|---|
| Characters | 0.52 | **0.67** | 0.50 | 0.61 | 0.61 | 0.50 | 0.62 |
| Pixels | 0.24 | 0.38 | 0.33 | 0.36 | **0.42** | 0.32 | 0.39 |

Challenge 2
Critical Assessment of Information

# Spatio-Temporal Analysis of Vandalism in Wikipedia

(highlighting temporal dynamics)

# Wikipedia Vandalism
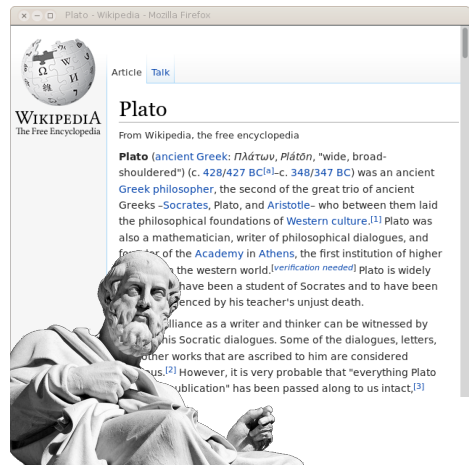
# Wikipedia Vandalism



Vandalism is a problem for Wikipedia

- ❑  470 million edits to the English Wikipedia (14 years)
- ❑  40 million (9.5%) are vandalism
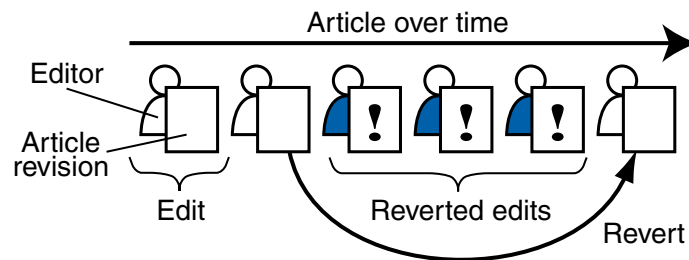- →  Rate of today: a vandalism case every 5 seconds

How to fight vandalism?

- ❑  Explain **why** people vandalize
- ❑  Analyze **when** people vandalize
- ❑  Analyze **where** these people are

# Wikipedia Vandalism



## Vandalism is a problem for Wikipedia

- ❏ 470 million edits to the English Wikipedia (14 years)
- ❏ 40 million (9.5%) are vandalism
- → Rate of today: a vandalism case every 5 seconds

## How to fight vandalism?

- ❏ Explain **why** people vandalize
- ❏ Analyze **when** people vandalize
- ❏ Analyze **where** these people are

## Language-independent detection approach

- ❏ Take all 1.2 billion edits to the 7 most-edited Wikipedias
  (english, german, french, spanish, russian, italian, japanese)
- ❏ Historical geolocation of anonymous editors
  (77% of edits by cross-checking RIR, IPligence, and IP2Location)
- ❏ Vandalism detector based on revert patterns (community behavior)
- → Spatio-temporal analysis per local time of anonymous editors

## Reverts (supported by Wiki interface)
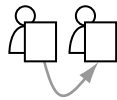
Not all reverts indicate vandalism

- ❏ Prior work: use only reverts whose comment indicates vandalism
- → Underestimates vandalism; language-dependent

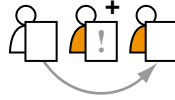- ❏ Our approach: identify revert patterns indicating non-vandalism
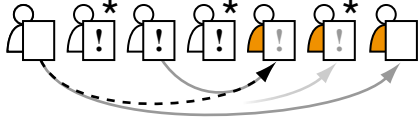


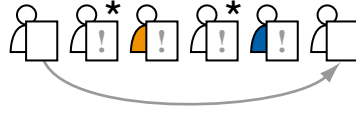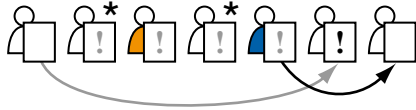Revert to blank page

Empty revert

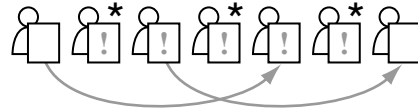Self-revert

Revert correction (enlargement)

Revert reverting more than one editor

Reverted revert

Interleaved reverts (edit war)
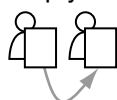
Not all reverts indicate vandalism

- ❑ Prior work: use only reverts whose comment indicates vandalism
- → Underestimates vandalism; language-dependent

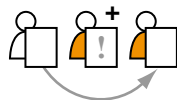- ❑ Our approach: identify revert patterns indicating non-vandalism

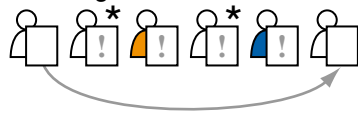Revert to blank page

Empty revert

Self-revert

Revert correction (enlargement)

Revert reverting more than one editor

Reverted revert

Interleaved reverts (edit war)

- ❑ Only 46% of reverted edits are vandalism
- ❑ Human evaluation: precision 82.8%, recall 84.7%
  (4 times the recall of prior work)

# Not all reverts indicate vandalism

❑ Prior work: use only reverts whose comment indicates vandalism

→ Underestimates vandalism; language-dependent

❑ Our approach: identify revert patterns indicating non-vandalism

# Vandalism over time of day



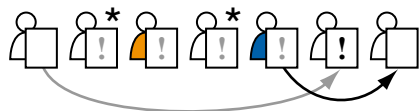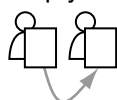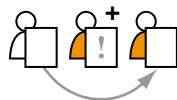Revert to blank page

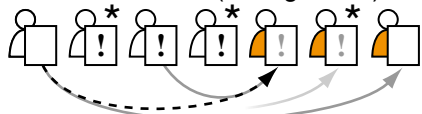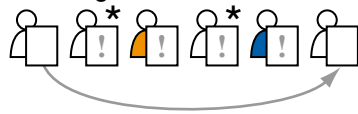

Empty revert



Self-revert



Revert correction (enlargement)



Revert reverting more than one editor



Reverted revert



Interleaved reverts (edit war)





❑ Only 46% of reverted edits are vandalism

❑ Human evaluation: precision 82.8%, recall 84.7%
(4 times the recall of prior work)

# Spatio-temporal vandalism analysis

# Security Estimate for Mnemonic Passwords

(highlighting volume)

# The mnemonic password advice

1. Create a sentence
2. Memorize it
3. Concatenate the first characters of each word
4. Use the string as password

*When I walked to the grocery store, there were camels flying overhead!*

Password: | wiwttgstwcfo |

☑ Show password

# The mnemonic password advice

1. Create a sentence
2. Memorize it
3. Concatenate the first characters of each word
4. Use the string as password

*When I walked to the grocery store, there were camels flying overhead!*

Password: wiwttgstwcfo

☑ Show password

Passwords that require a botnet ($H_1 \approx 65$ Bit):

☐ 14 random lowercase letters (out of 26)
☐ 10 random ASCII characters (out of 96)
☐ 5 random words (out of 7776)

And for mnemonic passwords?



Bar chart — x-axis: Word initials (t a o s i w h c b f m p d r e l n g u y v j k q z x), y-axis: Probability (0.00 to 0.20)

# The mnemonic password advice

1. Create a sentence
2. Memorize it
3. Concatenate the first characters of each word
4. Use the string as password

*When I walked to the grocery store,*
*there were camels flying overhead!*

Password: | wiwttgstwcfo |

☑ Show password

Passwords that require a botnet ($H_1 \approx 65$ Bit):

❑ 14 random lowercase letters (out of 26)
❑ 10 random ASCII characters (out of 96)
❑ 5 random words (out of 7776)

And for mnemonic passwords?



Probability — Word initials
t a o s i w h c b f m p d r e l n g u y v j k q z x

Depends on password distribution (Kerckhoffs' principle) $\rightarrow$ model distribution from a billion passwords

## The mnemonic password advice

(as per German BSI, Google, etc.)

1. Create a sentence
2. Memorize it
3. Concatenate the first characters of each word
4. Use the string as password

*_When _I _walked _to _the _grocery _store,*
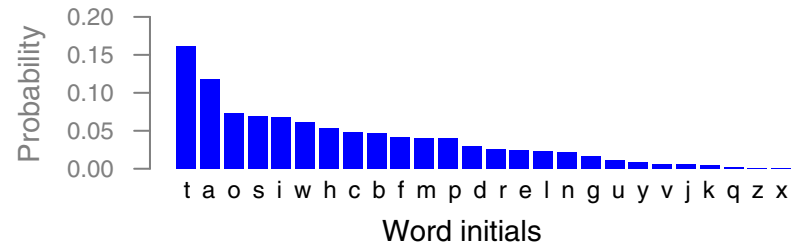*_there _were _camels _flying _overhead!*
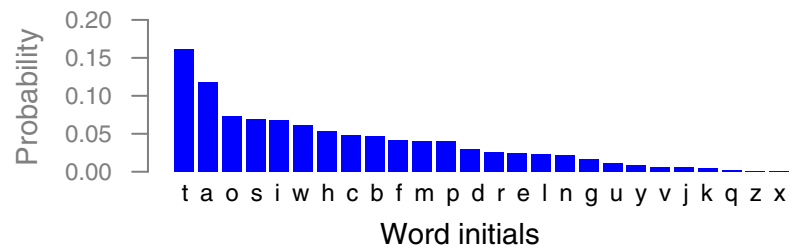
Password: | wiwttgstwcfo |

☑ Show password

Passwords that require a botnet ($H_1 \approx 65$ Bit):

❑ 14 random lowercase letters (out of 26)
❑ 10 random ASCII characters (out of 96)
❑ 5 random words (out of 7776)

And for mnemonic passwords?



Depends on password distribution (Kerckhoffs' principle) → model distribution from a billion passwords

## Approach: substitute mnemonics by web sentences

❑ 3 billion web sentences corpus from a standard web archive
❑ Statistically align the sentence corpus to mnemonics
❑ Estimate password distribution using position-dependent language models
→ Security estimates against offline ($H_1$) and online attacks ($H_0$, $\lambda_n$)

Sentence acquisition for password distribution estimate

| | | |
|---:|:---|---:|
| 5,000 | Mnemonics | Study by Yang et al., 2016 |
| 80,000 | Sentences | The Bible |
| 5,000,000 | Sentences | Encyclopedia Britannica |
| 70,000,000 | Passwords | Largest password corpus |
| 730,000,000 | Web pages | ClueWeb12, 27.3 TB |
| 3,400,000,000 | Sentences | Extracted and filtered |
| 500,000,000 | Sentences | And aligned to mnemonics |

Alignment in sentence complexity (≈ readability)



Sentences of length 12

# Security estimates (per character)

| Language model | Lowercase letters | | ASCII | |
|---|---|---|---|---|
| | $H_1$ | Ppl. | $H_1$ | Ppl. |
| Uniform | 4.70 | 26.0 | 6.55 | 94.0 |
| Order 0 | 4.15 | 17.8 | 5.09 | 34.1 |
| Order 8 | 3.71 | 13.1 | 3.98 | 15.8 |
| Order 8, position-dependent | 3.65 | 12.6 | 3.70 | 13.0 |

Security estimates (per character)

| Language model | Lowercase letters | | ASCII | |
| --- | --- | --- | --- | --- |
| | $H_1$ | Ppl. | $H_1$ | Ppl. |
| Uniform | 4.70 | 26.0 | 6.55 | 94.0 |
| Order 0 | 4.15 | 17.8 | 5.09 | 34.1 |
| Order 8 | 3.71 | 13.1 | 3.98 | 15.8 |
| Order 8, position-dependent | 3.65 | 12.6 | 3.70 | 13.0 |

Reaching $H_1$ = 65 Bit with mnemonic passwords

❑ Lowercase letters from 13+ words sentence          54 Bit

❑ 7-bit visible ASCII  (incl. %, !, @, #, etc.)          8 Bit
  (adds on average 2 characters $\approx$ 6.4 Bit)

❑ Word replacements  (and $\rightarrow$ &, to $\rightarrow$ 2, etc.)          2 Bit

❑ Different characters  (last of each word)          0 Bit

❑ Complex sentences  (rich vocabulary)          $+$          2 Bit

          66 Bit

# Harnessing Web Archives to Tackle Selected Societal Challenges
## Summary

### 1. Preservation of digital culture

- ❏ 10K pages high-fidelity archive (FAIRest dataset award)

- ❏ Reproduction assessment task

- ❏ 9K pages segmentation dataset

- ❏ Segmentation evaluation measures

### 2. Critical assessment of information

- ❏ Revert-based vandalism detection

- ❏ 30K edits Wiki vandalism dataset

- ❏ 1M hyperpartisan news dataset

- ❏ Style-based polarity detection

- ❏ Hyperpartisan news challenge (SemEval, 42 teams)

### 3. Online security and privacy

- ❏ 3B web sentences dataset

- ❏ Position-dependent language model

- ❏ Security estimate: mnemonic passwords

- ❏ Personal archiving tool

### Tailored web archiving technology  (Webis Web Archiver)

Highlighted aspects:

- ❏ Reproducibility
- ❏ Temporal dynamics
- ❏ Volume