# Web Page Segmentation Revisited: Evaluation Framework and Dataset

## CIKM 2020

**Johannes Kiesel**[1]

Florian Kneist[1]

Lars Meyer[1]

Kristof Komlossy[1]

Benno Stein[1]

Martin Potthast[2]

[1] Bauhaus-Universität Weimar
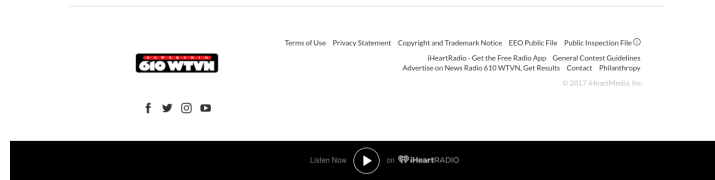
[2] UNIVERSITÄT LEIPZIG

[1,2] Webis

# Web Page Segmentation

# Web Page Segmentation

# Web Page Segmentation

# Web Page Segmentation

# Web Page Segmentation: Downstream Tasks (Examples)



❑ **Content Extraction**
Image: *Language Independent Content Extraction from Web Pages.* Javier et al., DIR'09.

# Web Page Segmentation: Downstream Tasks (Examples)



- ❑ **Content Extraction**
  Image: *Language Independent Content Extraction from Web Pages.* Javier et al., DIR'09.

- ❑ **Template Detection**
  Image: *Automatic Data Extraction From Template Generated Web Pages.* Ma et al., PDPTA'03.

# Web Page Segmentation: Downstream Tasks (Examples)





❑ **Content Extraction**
Image: *Language Independent Content Extraction from Web Pages.* Javier et al., DIR'09.

❑ **Template Detection**
Image: *Automatic Data Extraction From Template Generated Web Pages.* Ma et al., PDPTA'03.

❑ **Design Mining**
Image: *Webzeitgeist: Design Mining the Web.* Kumar et al., CHI'13.



**LAYOUT QUERY**

# Concept Formation: Web Page Segment

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.

Rationale: Web pages are created for human viewers, and so are segments

Gestalt Principles provide common ground

Proximity

Similarity

Closure

Symmetry

# Evaluation Framework for Web Page Segmentation

> 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 $s_i \subseteq E$

# Evaluation Framework for Web Page Segmentation

> 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 $s_i \subseteq E$

Suggested sets of elements:



Characters



DOM nodes



Pixels



Edges

# Evaluation Framework for Web Page Segmentation

> 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 $s_i \subseteq E$

Precision

$$P_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S|} \right)$$

Recall

$$R_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S^*|} \right)$$

F-Measure, $F_{B^3}$, is defined as the harmonic mean of precision and recall as usual

Note: $P_{B^3}(S, S') = R_{B^3}(S', S) \quad \Rightarrow \quad F_{B^3}(S, S') = F_{B^3}(S', S)$

# Evaluation Framework for Web Page Segmentation

# Evaluation Framework for Web Page Segmentation



$$S \qquad \text{Characters of } S \text{ in } S^* \qquad S^*$$

$$P_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S|} \right)$$

# Evaluation Framework for Web Page Segmentation



$S$      Pixels of $S$ in $S^*$      $S^*$

$$P_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S|} \right)$$

# Evaluation Framework for Web Page Segmentation



$$S \qquad\qquad \text{Pixels of } S^* \text{ in } S \qquad\qquad S^*$$

$$P_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S|} \right)$$

$$R_{B^3}(S, S^*) = \text{avg}_e \left( \frac{|\text{elements in same segment as } e \text{ in both } S \text{ and } S^*|}{|\text{elements in same segment as } e \text{ in } S^*|} \right)$$

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Remove 644 error pages (JCDL'19)

Preprocessing and filtering
(from 10,000 to 8,490 pages)

Draw 5 segmentations for each web page

Curate crowdsourced segmentations

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

Fit drawn segments to DOM nodes

Fuse into one coherent ground truth

Post-processing
(77,017 segments on 8,490 pages)

Webis-WebSeg-20 (CIKM'20)

# The Webis-WebSeg-20 Dataset

| Webis-Web-Archive-17 (JDIQ'18) |
|---|

Extract node locations and text content

Remove 866 overly simple web pages

Remove 644 error pages (JCDL'19)

Preprocessing and filtering
(from 10,000 to 8,490 pages)

Draw 5 segmentations for each web page

Curate crowdsourced segmentations

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

Fit drawn segments to DOM nodes

Fuse into one coherent ground truth

Post-processing
(77,017 segments on 8,490 pages)

| Webis-WebSeg-20 (CIKM'20) |
|---|

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Remove 644 error pages (JCDL'19)

> Preprocessing and filtering
> (from 10,000 to 8,490 pages)

Draw 5 segmentations for each web page

Curate crowdsourced segmentations

> Crowdsourcing on Mechanical Turk
> ($\approx$ half a year of full-time work)

Fit drawn segments to DOM nodes

Fuse into one coherent ground truth

> Post-processing
> (77,017 segments on 8,490 pages)

Webis-WebSeg-20 (CIKM'20)

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Remove 644 error pages (JCDL'19)

Preprocessing and filtering
(from 10,000 to 8,490 pages)

**Draw 5 segmentations for each web page**

**Curate crowdsourced segmentations**

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

Fit drawn segments to DOM nodes

Fuse into one coherent ground truth

Post-processing
(77,017 segments on 8,490 pages)

Webis-WebSeg-20 (CIKM'20)

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Remove 644 error pages (JCDL'19)

Preprocessing and filtering
(from 10,000 to 8,490 pages)

Draw 5 segmentations for each web page

Curate crowdsourced segmentations

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

| Agreement measure | Elements | | | |
|---|---|---|---|---|
| | Characters | Nodes | Pixels | Edges |
| $F_{B^3}$ | 0.78 | 0.74 | 0.65 | 0.73 |

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Preprocessing and filtering
(from 10,000 to 8,490 pages)

Remove 644 error pages (JCDL'19)

Draw 5 segmentations for each web page

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

Curate crowdsourced segmentations

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

# The Webis-WebSeg-20 Dataset

Webis-Web-Archive-17 (JDIQ'18)

Extract node locations and text content

Remove 866 overly simple web pages

Preprocessing and filtering
(from 10,000 to 8,490 pages)

Remove 644 error pages (JCDL'19)

Draw 5 segmentations for each web page

Crowdsourcing on Mechanical Turk
($\approx$ half a year of full-time work)

Curate crowdsourced segmentations

Fit drawn segments to DOM nodes

Post-processing
(77,017 segments on 8,490 pages)

Fuse into one coherent ground truth

Webis-WebSeg-20 (CIKM'20)

# Web Page Segmentation Revisited: Evaluation Framework and Dataset

Evaluation Framework for Web Page Segmentation

- ❑ Segmentation similarity, quality, and fusion
- ❑ Comparison of level of detail of segmentations
- ❑ Adjustable for different downstream tasks

Webis-WebSeg-20

- ❑ 8,490 pages from 4,824 sites
- ❑ 5 human annotators each page
- ❑ Segments in "Simple Feature Access" standard
- ❑ Web pages provided in several representations:
  - HTML file
  - Screenshot
  - Screenshot coordinates of DOM nodes
  - Webis-Web-Archive-17 WARC file



https://webis.de/publications.html?q=
johannes+kiesel+web+archive#stein_2020w

Paper, browser, code, data