Johannes Kiesel    Fabienne Hubricht    Benno Stein    Martin Potthast
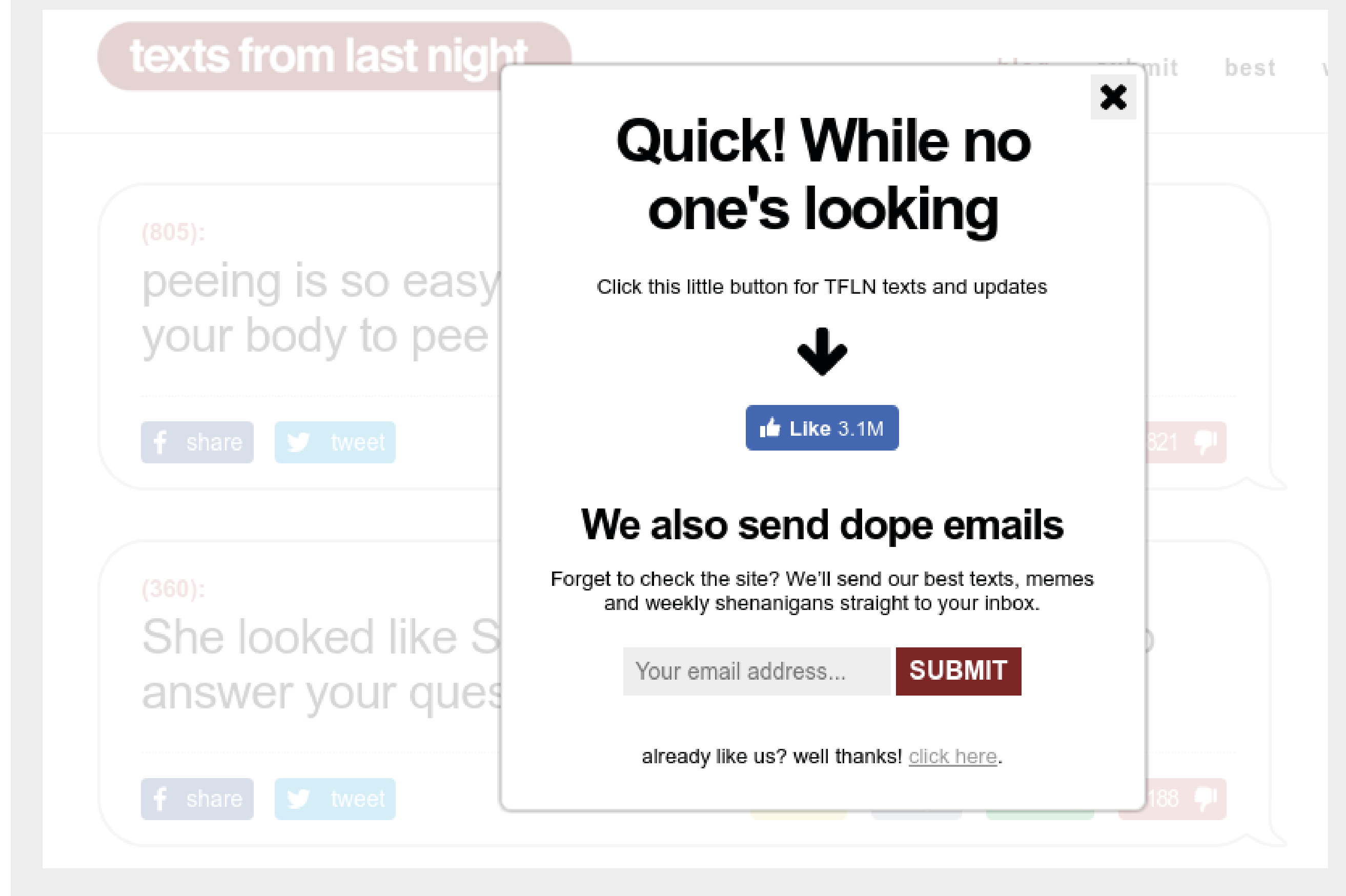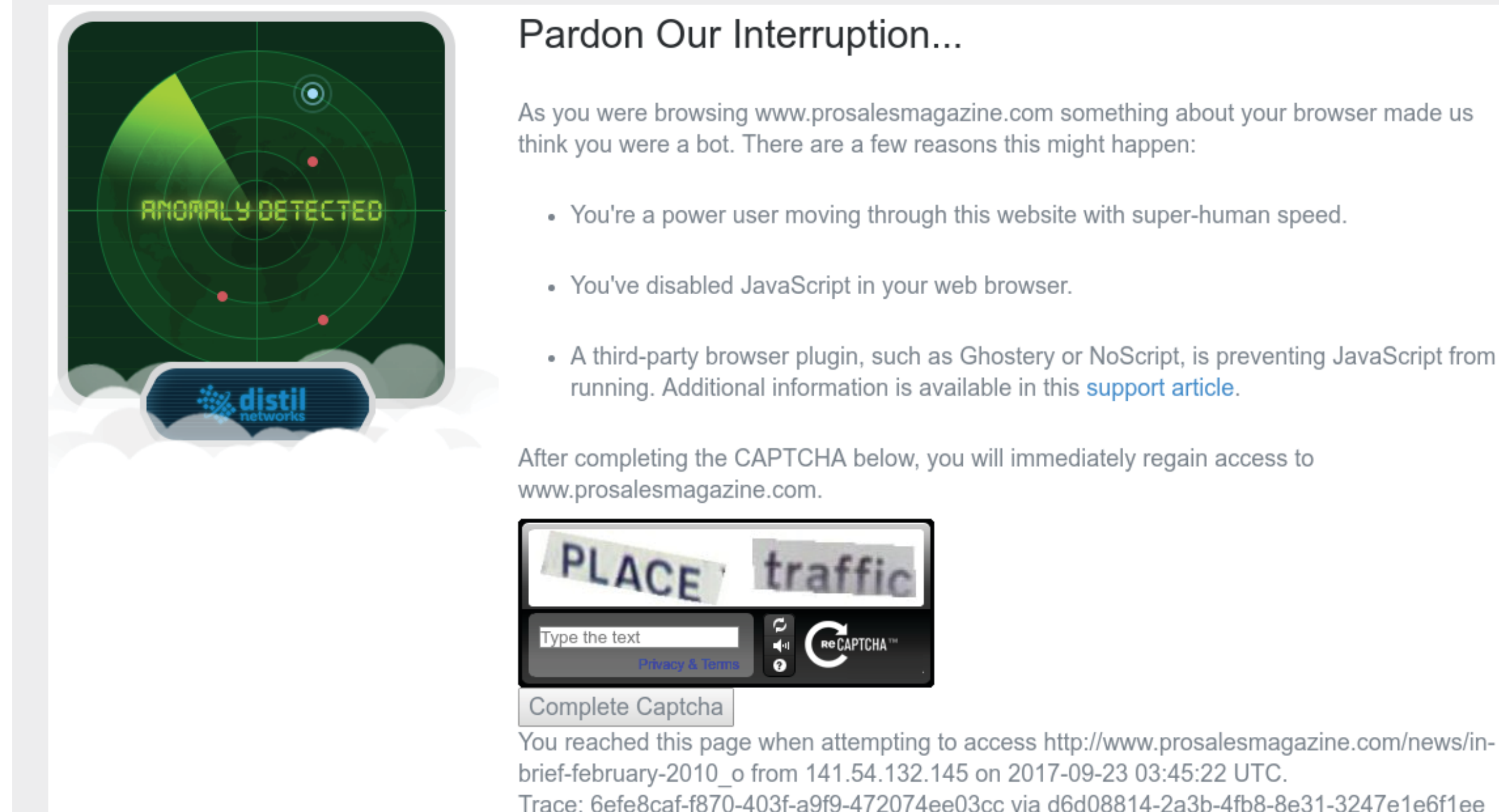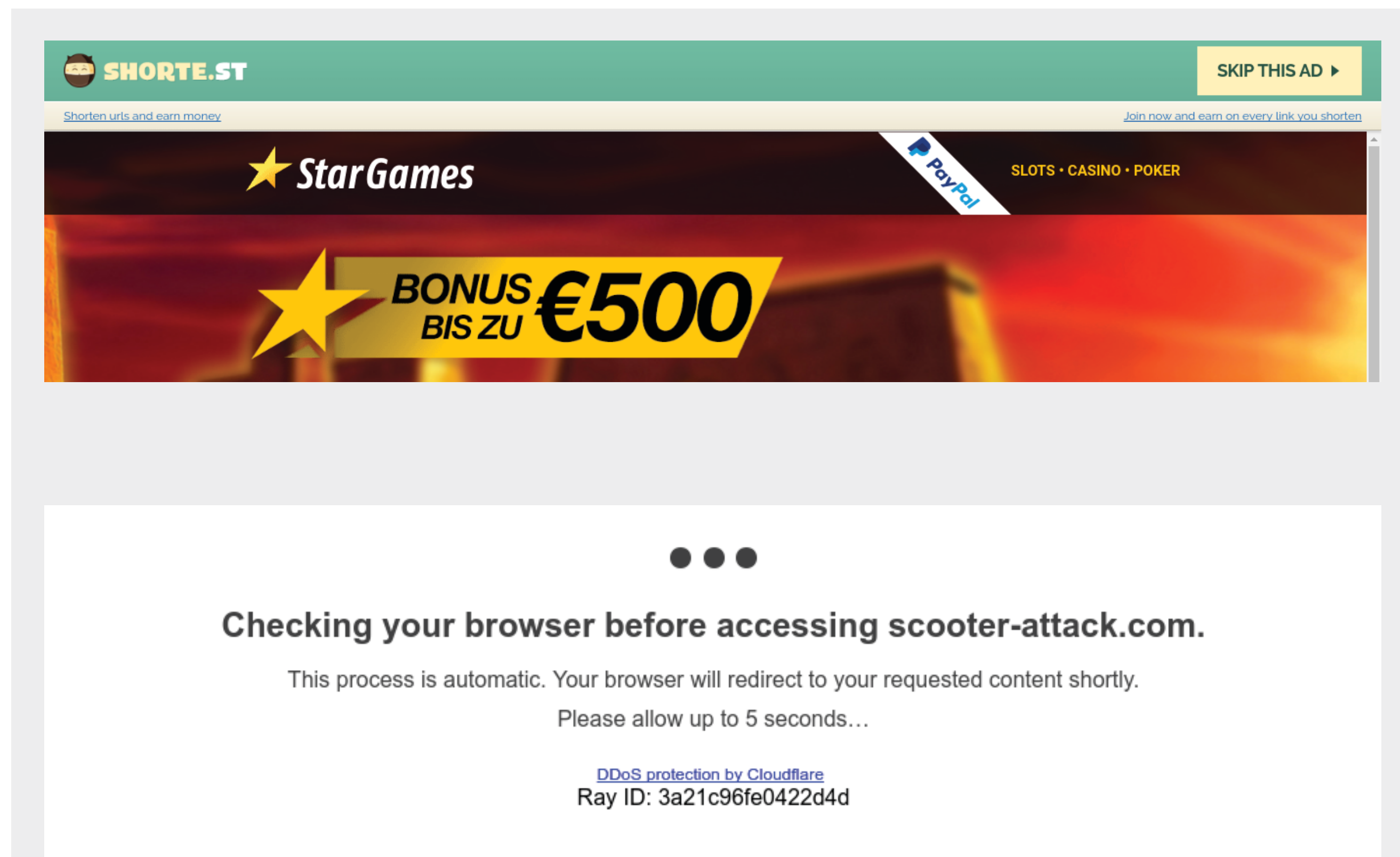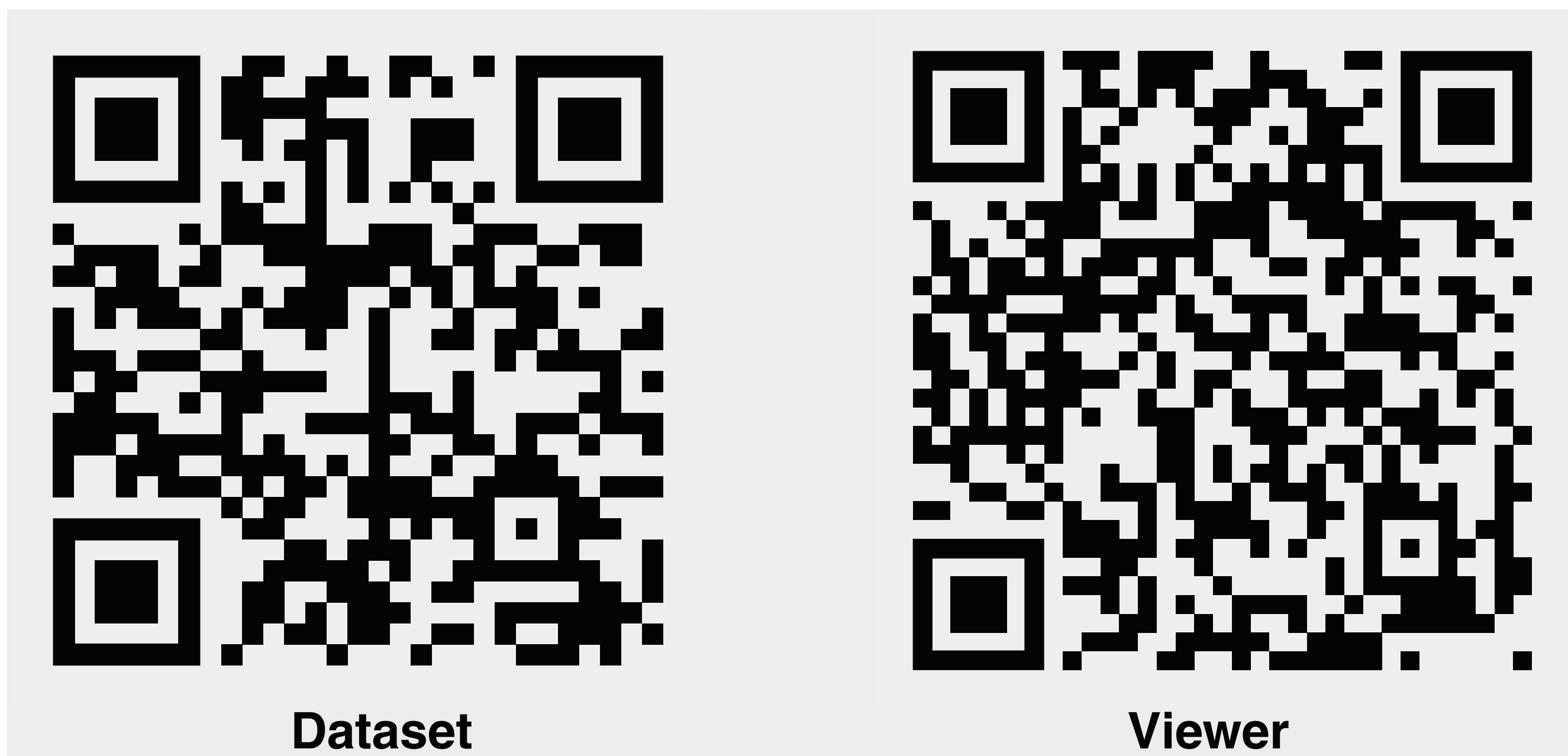
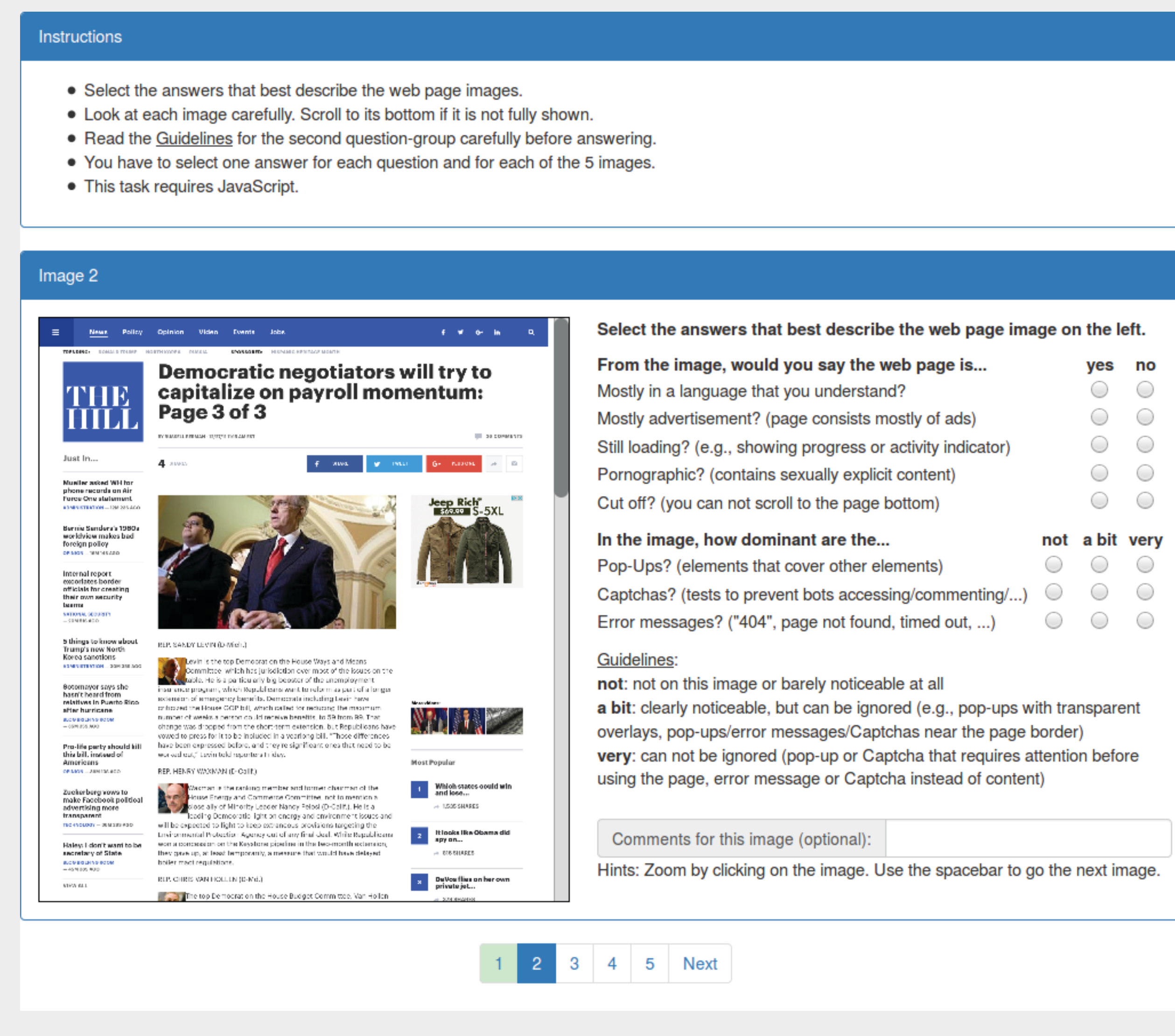# A Dataset for Content Error Detection in Web Archives



## A Dataset of 10,000 Web Pages

**Content errors** if a web page is different than one expects.
We define 5 types of content errors.

**Webis-Web-Archive-17** sampled from Common Crawl.
Annotated for reproduction errors in previous work

**Annotation** using Mechanical Turk with 5 annotators per page.
Distinction whether errors can be ignored when using the page.

**Curation** by two of the authors if annotators disagreed.
Assuring consistency in the labels.



**Dataset**          **Viewer**

| Content error | Agreement | Corrections | Distribution | | %Error |
|---|---|---|---|---|---|
| | | | No | Yes | |
| Ad page | 0.65 | 329 | 9865 | 105 | 1.1 |
| Loading indicators | 0.89 | 48 | 9950 | 50 | 0.5 |
| | | | Not | A bit | Very | |
| Pop-ups | 0.82 | 394 | 9297 | 315 | 388 | 3.9 |
| CAPTCHAs | 0.91 | 124 | 9865 | 60 | 75 | 0.8 |
| Error messages | 0.89 | 331 | 9554 | 83 | 363 | 4.5 |