

Acquiring Web Site Metadata by Heterogeneous Information Sources Processing

What is the problem?

- too many web documents emerge
- in december 2011 there were about 556 000 000 websites
- web documents structure differ from site to site
- contemporary approaches do not achieve satisfactory results

What to do about it?

- we need to utilize **features** of today's web pages
- web documents provide **implicit semantics** in style
- **keyword paradigm** seems appropriate to describe web documents

Why keyword paradigm?

- widespread in many disciplines (summarization and categorization)
- „de facto“ standard in searching (google search)
- used in social media (e.g. del.icio.us)
- used in advertising (e.g. AdWords)

Our approach to keywords extraction

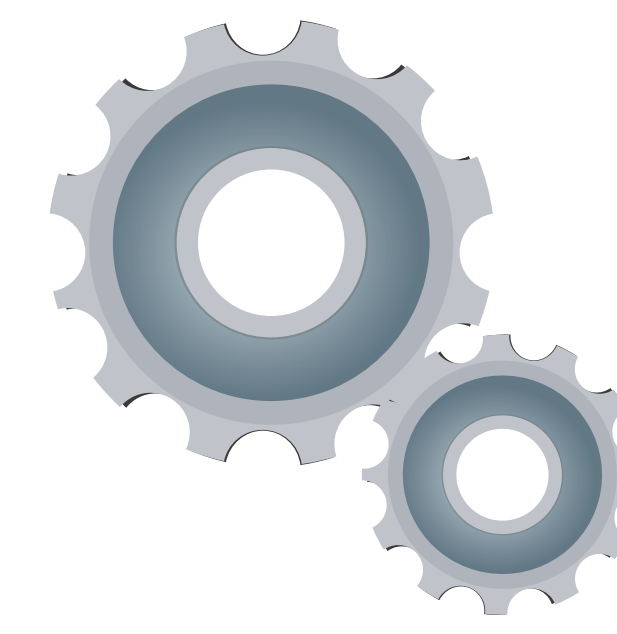
3. Find emphasized phrases

- 1. emphasized phrase + CSS features
- 2. emphasized phrase + CSS features
- ...

6. Get relevant keywords

- ✓ Keyword 1
- ✓ Keyword 2
- ✓ Keyword 3
- ✓ Keyword 4
- ✓ Keyword 5
- ...

5. ATR + some magic

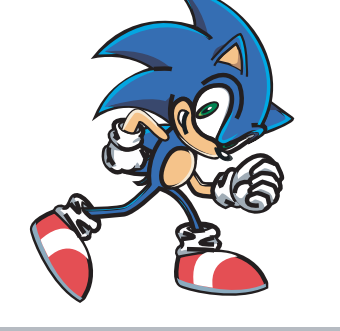


2. Extract content

Lorem ipsum dolor sit amet, consectetur adipiscing elit

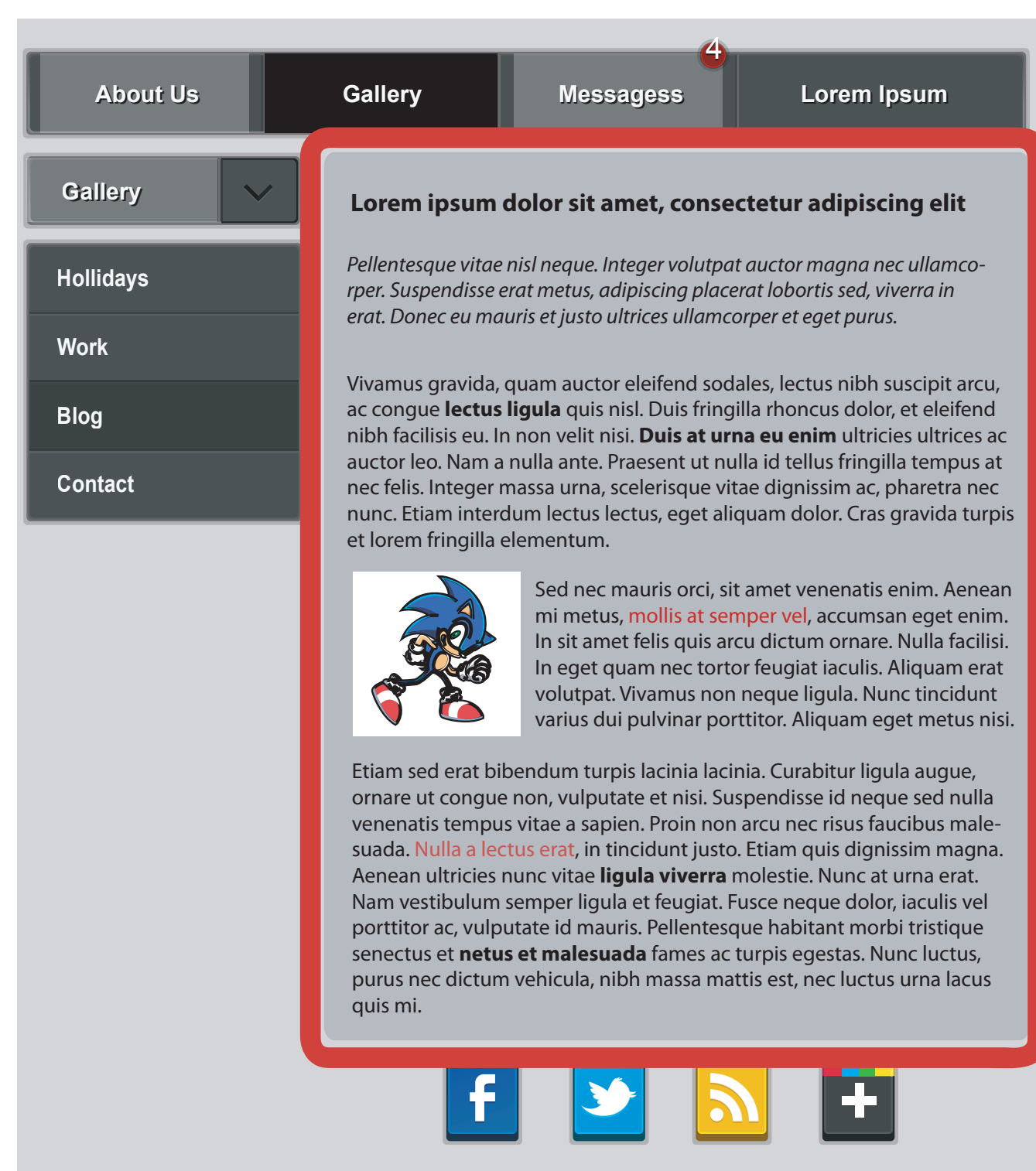
Pellentesque vitae nisi neque. Integer volutpat auctor magna nec ullamcorper. Suspendisse erat metus, adipiscing placerat lobortis sed, viverra in erat. Donec eu mauris et justo ultrices ullamcorper et eget purus.

Vivamus gravida, quam auctor eleifend sodales, lectus nibh suscipit arcu, ac congue **lectus ligula** quis nisi. Duis fringilla rhoncus dolor, et eleifend nibh facilisis eu. In non velit nisi. **Duis at urna eu enim** ultricies ultrices ac auctor leo. Nam a nulla ante. Praesent ut nulla id tellus fringilla tempus at nec felis. Integer massa urna, scelerisque vitae dignissim ac, pharetra nec nunc. Etiam interdum lectus lectus, eget aliquam dolor. Cras gravida turpis et lorem fringilla elementum.



Sed nec mauris orci, sit amet venenatis enim. Aenean mi metus, **mollis at semper vel**, accumsan eget enim. In sit amet felis quis arcu dictum ornare. Nulla facilisi. In eget quam nec tortor feugiat iaculis. Aliquam erat volutpat. Vivamus non neque ligula. Nunc tincidunt varius dui pulvinar porttitor. Aliquam eget metus nisi.

Etiam sed erat bibendum turpis lacinia lacinia. Curabitur ligula augue, ornare ut congue non, vulputate et nisi. Suspendisse id neque sed nulla venenatis tempus vitae a sapien. Proin non arcu nec risus faucibus malesuada. **Nulla a lectus erat**, in tincidunt justo. Etiam quis dignissim magna. Aenean ultricies nunc vitae **ligula viverra** molestie. Nunc at urna erat. Nam vestibulum semper ligula et feugiat. Fusce neque dolor, iaculis vel porttitor ac, vulputate id mauris. Pellentesque habitant morbi tristique senectus et **netus et malesuada** fames ac turpis egestas. Nunc luctus, purus nec dictum vehicula, nibh massa mattis est, nec luctus urna lacus quis mi.



1. Find all CSS styles



4. Extract plain text

Lorem ipsum dolor sit amet, consectetur adipiscing elit Pellentesque vitae nisi neque. Integer volutpat auctor magna nec ullamcorper. Suspendisse erat metus, adipiscing placerat lobortis sed, viverra in erat. Donec eu mauris et justo ultrices ullamcorper et eget purus. Vivamus gravida, quam auctor eleifend sodales, lectus nibh suscipit arcu, ac congue lectus ligula quis nisi. Duis fringilla rhoncus dolor, et eleifend nibh facilisis eu. In non velit nisi. Duis at urna eu enim ultricies ultrices ac auctor leo. Nam a nulla ante. Praesent ut nulla id tellus fringilla tempus at nec felis. Integer massa urna, scelerisque vitae dignissim ac, pharetra nec nunc. Etiam interdum lectus lectus, eget aliquam dolor. Cras gravida turpis et lorem fringilla elementum...

Magic explanation

- extract emphasized words and phrases from web document
- compute the *visibility factor*
- extract phrases and words from plain text
- improve weight of phrases, which are emphasized by CSS styles

Visibility factor

- add points for every style feature, which is different from main textual content.

Experimental evaluation

- So far we have evaluated
- 60 web documents
- from 11 different web sites
- together 2048 phrases

Further work

- determine important style features
- train SVM to automatically detect keywords

Base algorithms	Average relevance %	After utilizig CSS features	AVG relevance %	Delta %
<i>CValue</i>	72,98	<i>Cvalue</i>	74,71	1,73
<i>TermEx</i>	71,44	<i>TermEx</i>	73,64	2,2
<i>Weirdness</i>	63,61	<i>Weirdness</i>	70,26	6,65
<i>GlossEx</i>	55,37	<i>GlossEx</i>	63,1	7,73
<i>TF-IDF</i>	48,27	<i>TF-IDF</i>	69,98	21,71

