# Embedding Knowledge in HTML

## HTML is Everywhere

- We usually think of HTML as the language of Web pages
- But it's also widely used on/for mobile devices and tablets
  - It readily adapts for different screen sizes/orientations
- And is the basis of many ebook formats
  - E.g. Kindle's formats, mobi, epub
- How can we add knowledge to HTML pages?

## Adding RDF-like data to HTML

- We'd like to add semi-structured know-ledge to a conventional HTML document
  - Humans can see and understand the regular HTML content (text, images, videos, audio)
  - Machines can see and understand the data markup in XML, RDF or some other format
- Possibilities include
  - Add a link to a separate document with the knowledge
  - Embed the knowledge as comments, javascript, etc.
  - Distribute the knowledge markup throughout the HTML as attributes of existing HTML tags

## One page, not two

- Content providers prefer not to generate multiple pages, one for humans (HTML) and another for machines (RDF)
  - RDF serializations are complex
  - Requires a separate storage, generation, etc. mechanism
  - Introduces redundancy, which can lead to errors if we change one page but not the other
- Simplifies the job of search engines as well

## General approach

- Provide or reuse tag *attributes* to encode the metadata
  - Browsers and web apps ignore attributes they don't understand
- Three approaches have been developed
  - Microformats (~ 2005)
  - RDFa (~ 2007)
  - Microdata (aka schema.org) (~ 2012)
- Status 2014/5 (IMHO)
  - *Microformats* used but future is limited
  - *RDFa* becoming the encoding of choice
  - Schema.org vocabularies getting large uptake
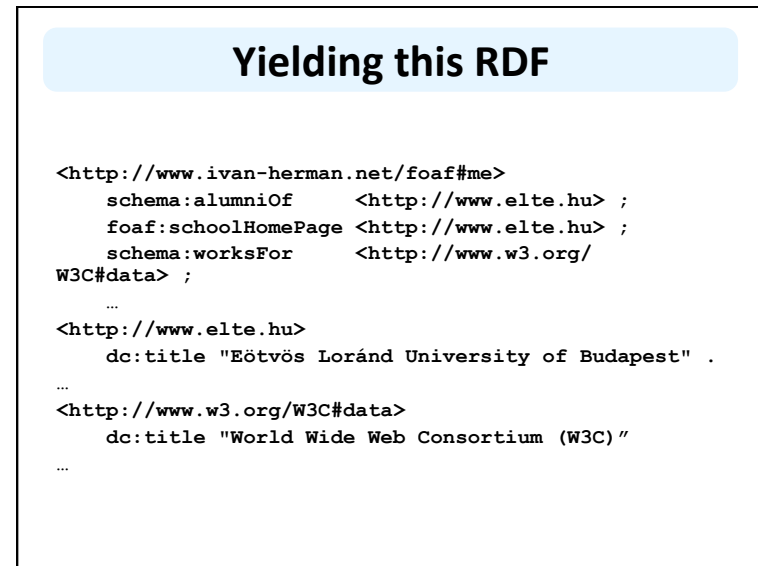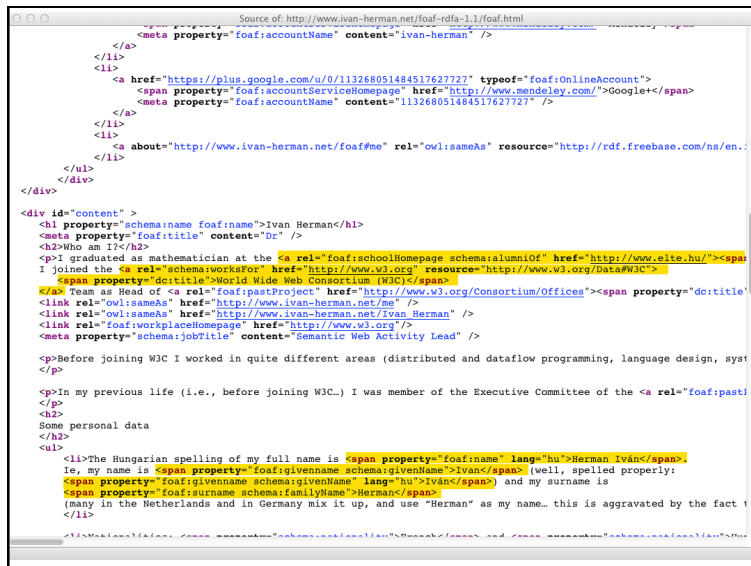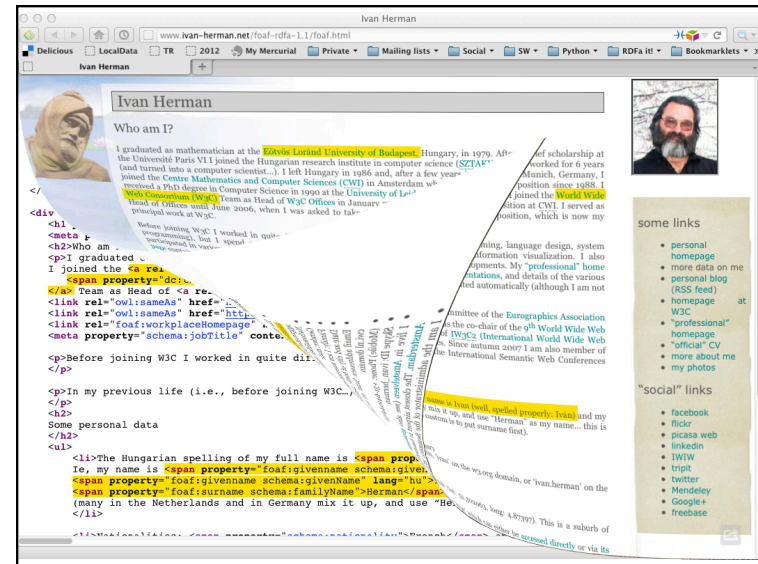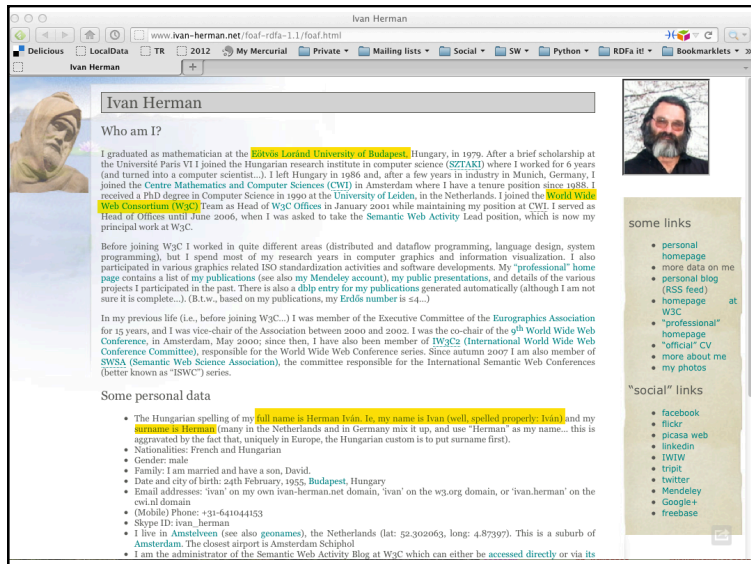
## Microformats approach

- Reuses HTML attributes like @class, @title
- Separate vocabularies (address, CV, …)
- Difficult to mix microformats (no concept of namespaces)
- Does not, inherently, define an RDF representation

  possible to transform via, e.g., XSLT + GRDDL, but transformations are vocabulary dependent
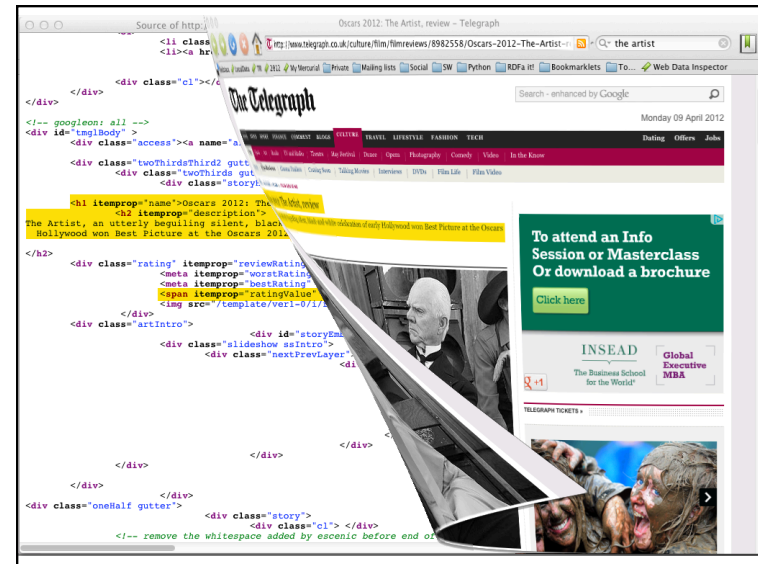
## Microdata approach

- Defined and supported by Google, Bing, Yahoo and Yandex
- Adds new attributes to HTML5 to express metadata
- Works well for simpler "single-vocabulary" cases, but not well suited for mixing vocabularies or for complex vocabularies
- No notion of datatypes or namespaces
- Defines a generic mapping to RDF

## RDFa approach

- Adds new (X)HTML/XML attributes
- Has namespaces and URIs at its core
  - So mixing vocabulary is easy, as in RDF
- Complete flexibility for using literals or URI resources
- Is a complete serialization of RDF

## Yielding this RDF

```
<http://www.ivan-herman.net/foaf#me>
    schema:alumniOf       <http://www.elte.hu> ;
    foaf:schoolHomePage   <http://www.elte.hu> ;
    schema:worksFor       <http://www.w3.org/
W3C#data> ;
    …
<http://www.elte.hu>
    dc:title "Eötvös Loránd University of Budapest" .
…
<http://www.w3.org/W3C#data>
    dc:title "World Wide Web Consortium (W3C)"
…
```

## Yielding this RDF

```
[ rdf:type schema:Review ;
  schema:name "Oscars 2012: The Artist, review" ;
  schema:description "The Artist, an utterly
beguiling…" ;
  schema:ratingValue "5" ;
  …
]
```

4

## Rich Snippets

- Search engines add text under results to preview what's on page and why it's relevant
- Text ften extracted from structured data embedded on the page
- See http://bit.ly/RichSN for more information



## RDFa and Microdata: similarities

- RDFa and Microdata are modern options
  - Microformats is another
- Both have similar approaches
  - Structured data encoded in *HTML attributes only* – no new elements
  - Define some special *attributes*
    - e.g., `itemscope` for microdata, `resource` for RDFa
  - Reuse *some* HTML core attributes (e.g., `href`)
  - Use textual content of HTML source, if needed
- RDF data can be extracted from both

## RDFa and microdata: differences

- Microdata *optimized* for simpler use cases:
  - One vocabulary at a time
  - Tree shaped data
  - No datatypes
- RDFa provides full serialization of RDF in XML or HTML
  - Price is extra complexity over Microdata
- RDFa 1.1 Lite is a simplified authoring profile of RDFa, very similar to microdata
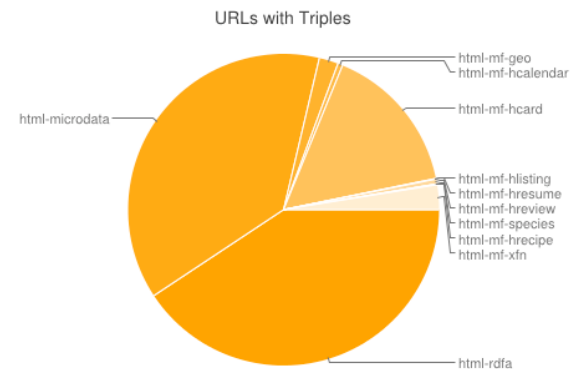
## Amount of structured data on Web?

- Web Data Commons project uses Common Crawl data to estimate how much structured data is on the Web
- Looked for Microdata, RDFa, and nine common Microdata formats (e.g., hCalendar, hCard) in URLs parsable as HTML
- Nov. 2013 crawl:
  - 44TB (compressed) data from 2.2B URLs from 13M domains
  - 14% of domains, 26% of URLs had semantic data
- Processing 40TB (compressed) of the 2012 crawl took 5.6K machine hours on 100 machines and cost ~$400

## What formats were found?

- Microdata use up (140K->463K sites form 2012->13)
- See here for details on 2013 crawl



URLs with Triples

html-mf-geo
html-mf-hcalendar
html-mf-hcard
html-microdata
html-mf-hlisting
html-mf-hresume
html-mf-hreview
html-mf-species
html-mf-hrecipe
html-mf-xfn
html-rdfa

## Conclusions

- The amount of structured data on the web is growing steadily
- Microdata shows the strongest growth
- RDFa also common
- Microformat data is probably not growing as much