

Vorlesung Semantic Web



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Vorlesung im Wintersemester 2012/2013

Dr. Heiko Paulheim

Fachgebiet Knowledge Engineering

Einführung

- Was ist das Semantic Web?
- Bausteine des Semantic Web
- Grundlagen: URIs, Unicode, XML

Was ist das Semantic Web?

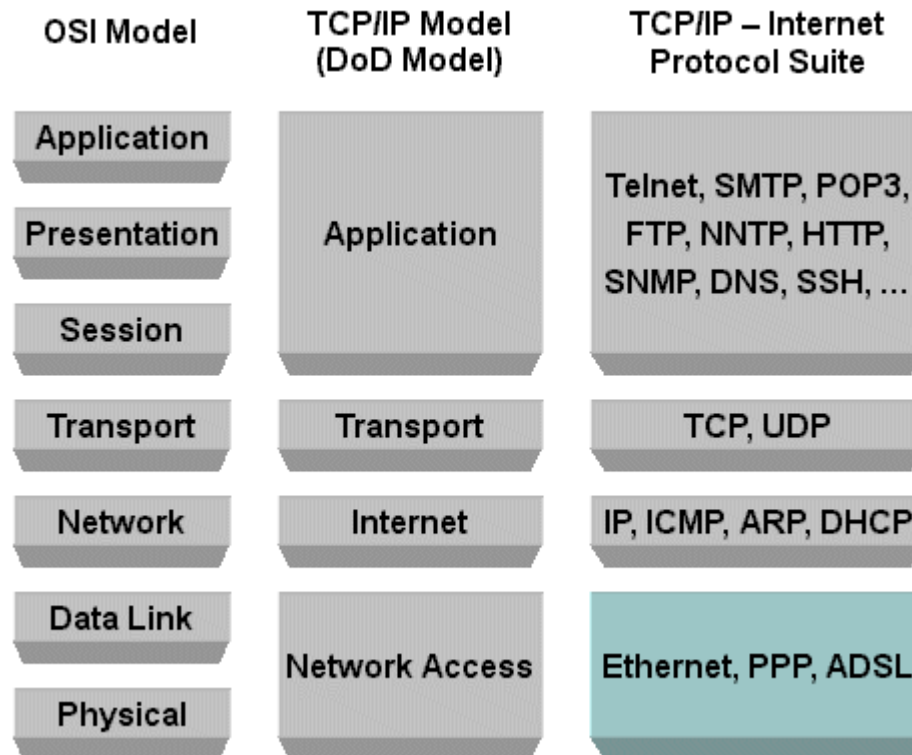
- Artikel von Tim Berners-Lee, Jim Hendler, und Ora Lassila:

„The Web is the killer app of the Internet.
The Semantic Web is another killer app
of that magnitude.“



Berners-Lee et al. (2001): *The Semantic Web*. In: Scientific American, Mai 2001.

Web vs. Internet?



Chin-Shiuh Shieh (2000): *TCP/IP - Internet Protocol Suite and Ethernet*.
<http://bit.kuas.edu.tw/~csshie/teach/np/tcpip/index.html>

Das „klassische“ Web

- HTTP-Protokoll
- URLs
- HTML als Auszeichnungssprache
 - plus CSS, JavaScript, ...
 - plus weitere mehr oder weniger standardisierte Formate (GIF, JPEG, Flash, ...)
- Browser als universeller Client

Das „klassische“ Web

- Hypertext: verlinkte Dokumente

Das World Wide Web

Der Grundstein für das World Wide Web wurde in den 90er-Jahren durch [Tim Berners-Lee](#) am [CERN](#) gelegt.

...

Tim Berners-Lee

Tim Berners-Lee (geboren 1955) gilt als einer der Erfinder des [World Wide Web](#).

...

CERN

Das CERN ist ein europäisches Forschungszentrum in der Nähe von Genf.

...

Eine kurze Geschichte des Webs



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Was wissen Sie über die Geschichte des Webs?
- Versuchen Sie, die folgenden Ereignisse zu ordnen:
 1. Erste HTML-Version
 2. Wikipedia geht online
 3. Gründung von Skype
 4. Erster Web-Katalog
 5. Gründung des W3C
 6. Erste Volltext-Suchmaschine
 7. Gründung von Twitter
 8. HTTP-Standard
 9. 500 Webserver am Netz
 10. Gründung von Facebook
 11. Dotcom-Blase und Börsencrash
 12. Erste Version von Internet Explorer
 13. Gründung von Google
 14. Erste Domain registriert
 15. Erste Version von Firefox
 16. TCP/IP-Standard
 17. 1.000 Computer am Netz
 18. 1.000.000 Computer am Netz
 19. 1.000.000.000 Computer am Netz
 20. Erstes Multi-User-Game im Netz



Eine kurze Geschichte des Webs



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- 1974: TCP/IP-Standard
- 1979: Erstes Multi-User-Game
- 1985: Erste Domains registriert, ~1.000 Computer am Netz
- 1989: Hypertext-Konzept von Tim Berners-Lee am CERN
- 1991: Erste HTML-Version (20 Elemente)
- 1992: ~1.000.000 Computer am Netz
- 1993: Mosaic-Browser, ca. 500 Webserver (weltweit)
- 1994: Volltext-Suchmaschinen (WebCrawler, Lycos),
Web-Kataloge (Yahoo, AltaVista), Gründung des W3C



Eine kurze Geschichte des Webs



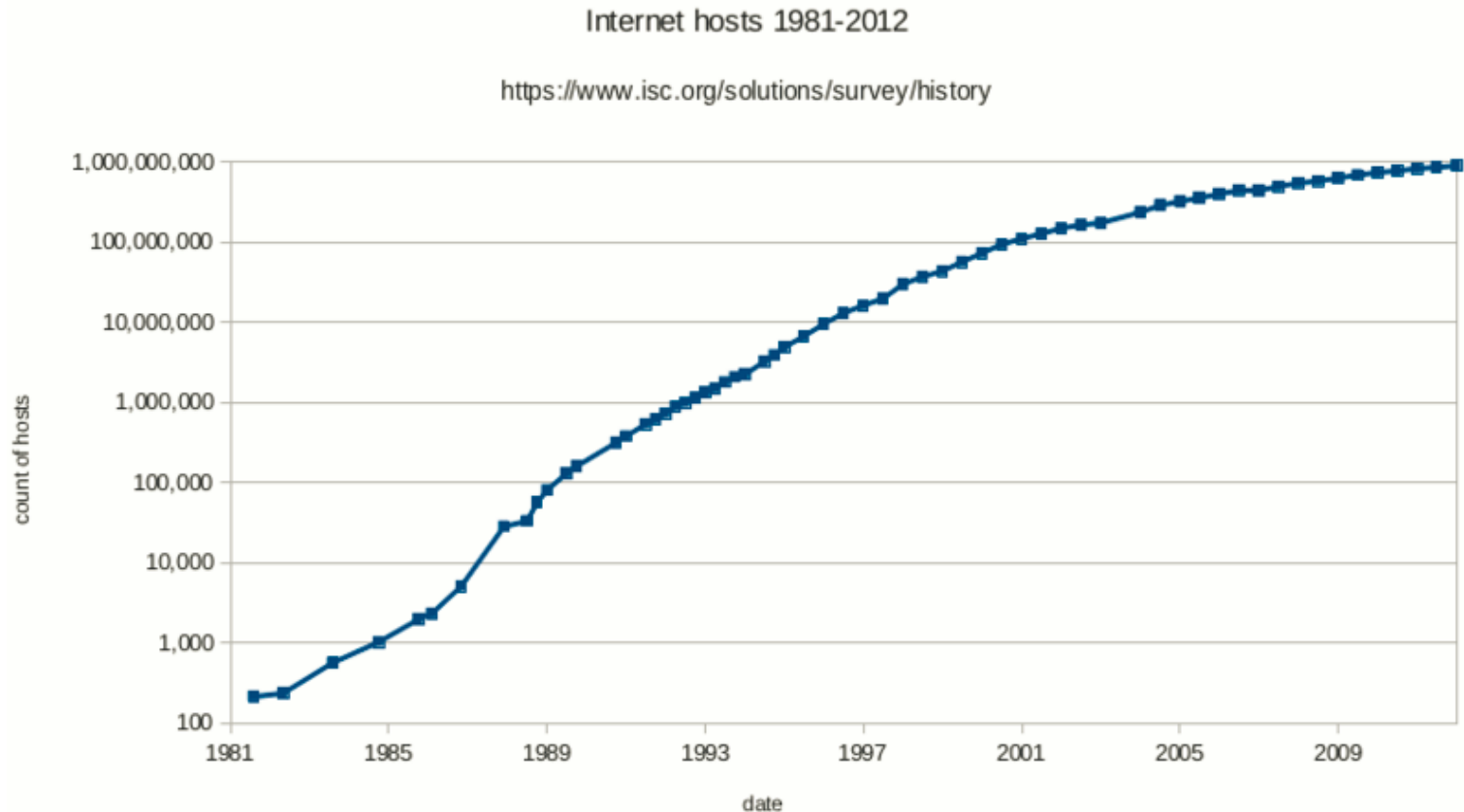
TECHNISCHE
UNIVERSITÄT
DARMSTADT

- 1995: Internet Explorer
- 1996: HTTP-Standard
- 1998: Google
- 2000: Dotcom-Börsencrash
- 2001: Wikipedia
- 2003: Skype
- 2004: Facebook, Firefox
- 2006: Twitter, WikiLeaks
- ...

...1.000.000.000 Computer sind bis jetzt nicht am Netz!



Wachstum des Webs



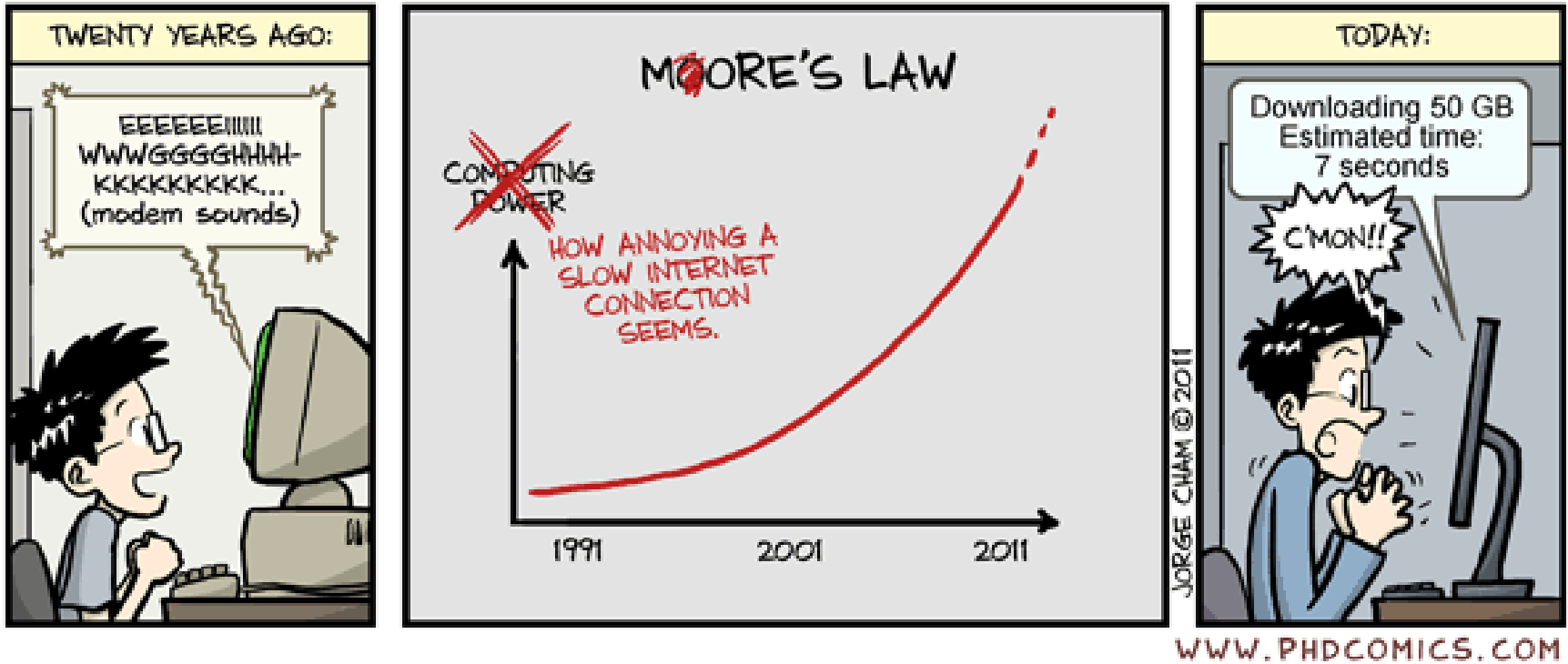
http://de.wikipedia.org/w/index.php?title=Datei:NASDAQ_IXIC_-_dot-com_bubble.png&filetimestamp=20050426161953

Der Dotcom-Börsencrash



http://de.wikipedia.org/w/index.php?title=Datei:NASDAQ_IXIC_-_dot-com_bubble.png&filetimestamp=20050426161953

Das "klassische" Web

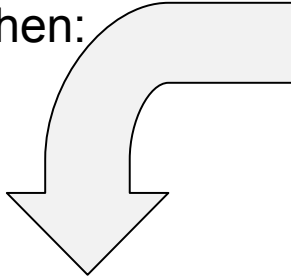


<http://www.phdcomics.com/comics.php?n=1456>

Das „klassische“ Web

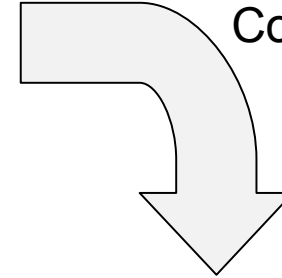


aus Sicht des
Menschen:



```
<html>
...
<b>Dr. Mark Smith</b>
<i>Physician</i>
Main St. 14
Smalltown
Mon-Fri 9-11 am
Wed 3-6 pm
...
</html>
```

aus Sicht des
Computers:



Dr. Mark Smith
Physician
Main St. 14
Smalltown
Mon-Fri 9-11 am
Wed 3-6 pm

Print in bold: „hmf298hmmhudsa“
Print in italics: „mj2i9ji0“
Print normal: „fdsah
02hfadsh0um2m0adsmf0ihm
asdfjköfdsa298ndsfmij32mio
lk2mjpoimjiofdpmsajiomjm“



Informationssuche im Web



Volltextsuche nach Stichwörtern (z.B. Google):

- „Mark Smith“
- „Physician in Smalltown“
- „Doctor in Smalltown“
- „Physician in Smalltown with opening hours on Wednesday afternoon“
- „Somebody in Smalltown who can fix a broken leg“

```
<html>
...
<b>Dr. Mark Smith</b>
<i>Physician</i>
Main St. 14
Smalltown
Mon-Fri 9-11 am
Wed 3-6 pm
...
</html>
```

=> „klassisches“ Web für Suche zu unflexibel

=> intelligente Agenten können es nicht nutzen

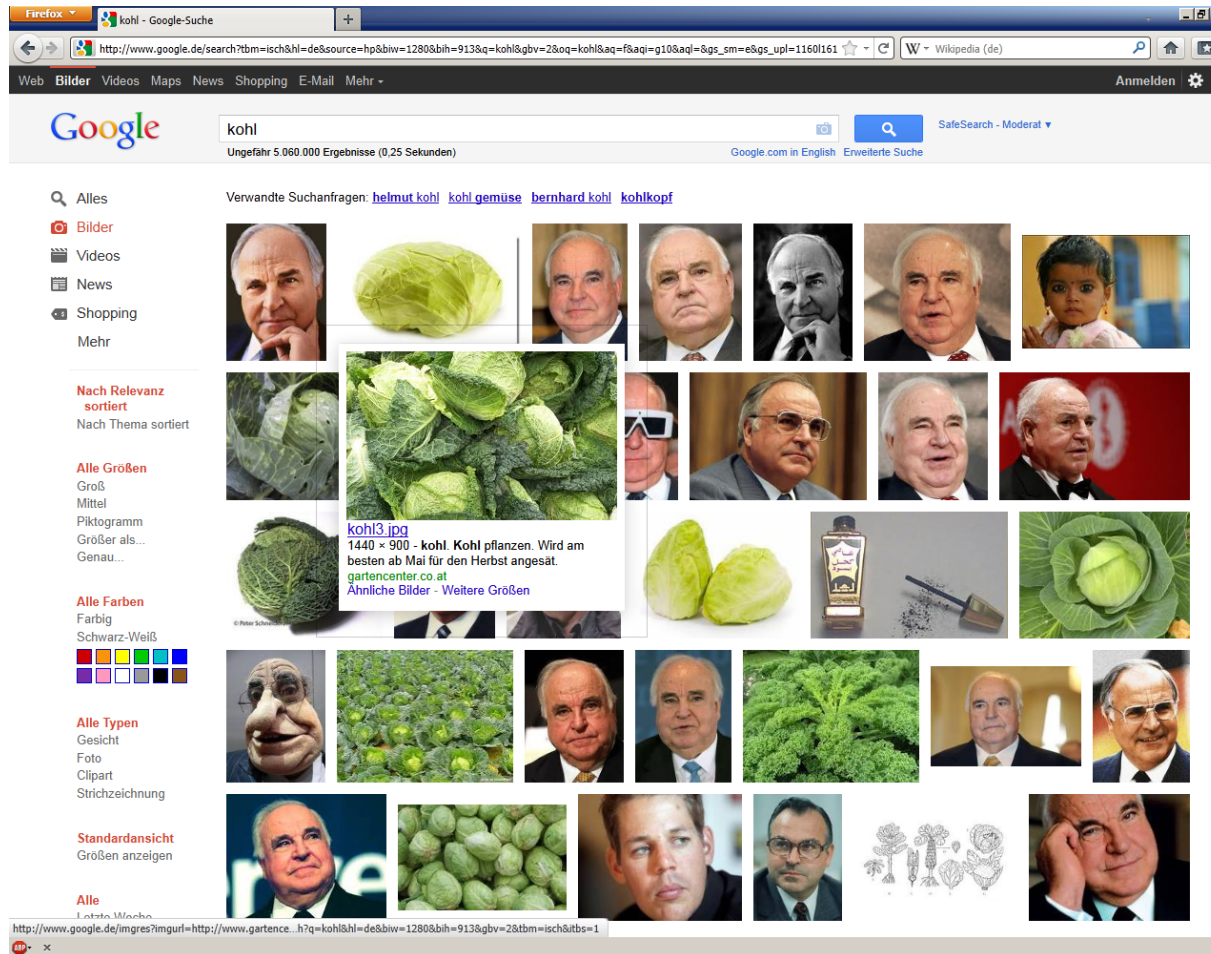
Probleme des Web

- Informationen finden
 - Stichwortbasierte Textsuche statt echter Fragen
 - verschiedene natürliche Sprachen
 - Homonyme/Polyseme
 - Synonyme
- Information verarbeiten
 - Formate (Encodings, Bilder, Videos, PDFs, ...)
- Information verwerten
 - Verteilt auf verschiedene Seiten
 - Bsp.: Information zu Buchautor auf Verlagsseite, Adresse auf Uni-Seite



<http://geekandpoke.typepad.com/geekandpoke/2011/08/coders-love-unicode.html>

Probleme des Web



Probleme des Web

Pfahls-Prozess: Richter rollen
die Skandale der Ära Kohl
neu auf

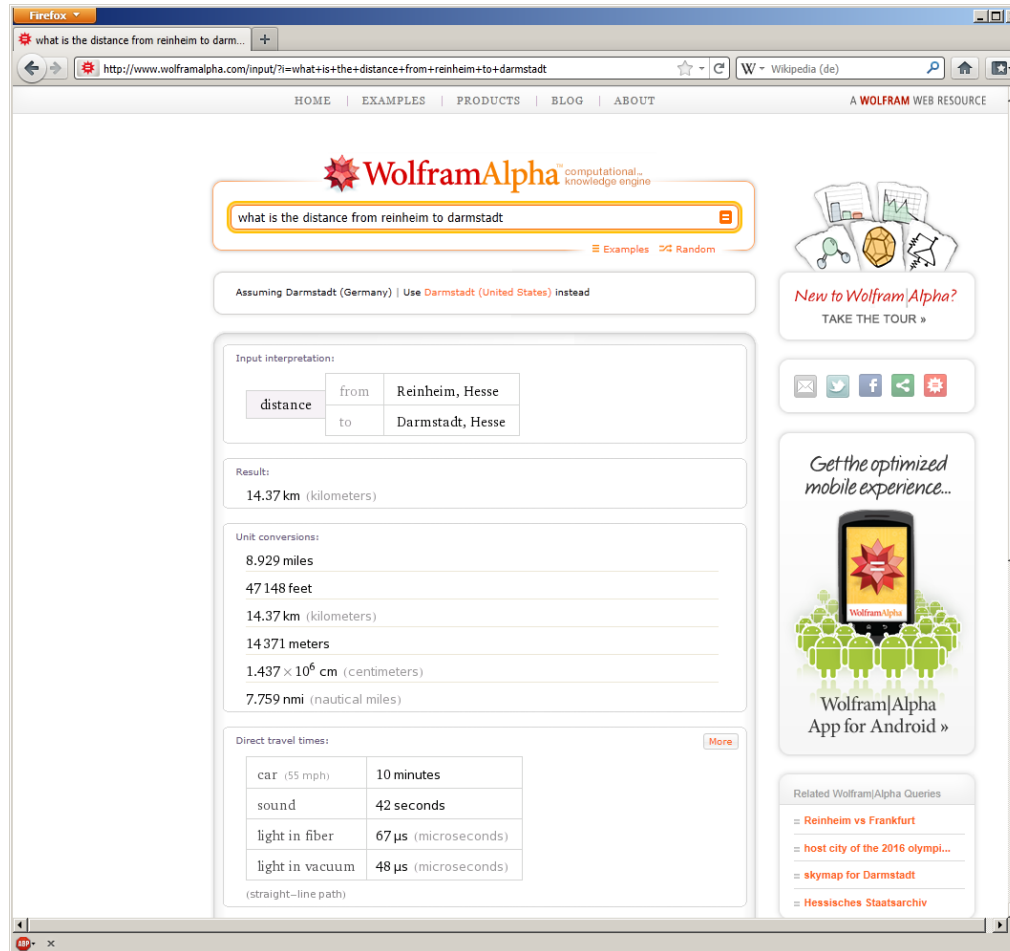
?

?

?



Beispiel: Wolfram Alpha



Firefox

what is the distance from reinheim to darm...

http://www.wolframalpha.com/input/?i=what+is+the+distance+from+reinheim+to+darmstadt

W - Wikipedia (de)

HOME | EXAMPLES | PRODUCTS | BLOG | ABOUT

A WOLFRAM WEB RESOURCE

WolframAlpha computational knowledge engine

what is the distance from reinheim to darmstadt

Examples Random

Assuming Darmstadt (Germany) | Use Darmstadt (United States) instead

New to Wolfram Alpha?
TAKE THE TOUR »

Input interpretation:

distance	from	Reinheim, Hesse
	to	Darmstadt, Hesse

Result:

14.37 km (kilometers)

Unit conversions:

8.929 miles

47 148 feet

14.37 km (kilometers)

14 371 meters

1.437 × 10⁶ cm (centimeters)

7.759 nmi (nautical miles)

Direct travel times:

car (55 mph)	10 minutes
sound	42 seconds
light in fiber	67 μs (microseconds)
light in vacuum	48 μs (microseconds)

(straight-line path)

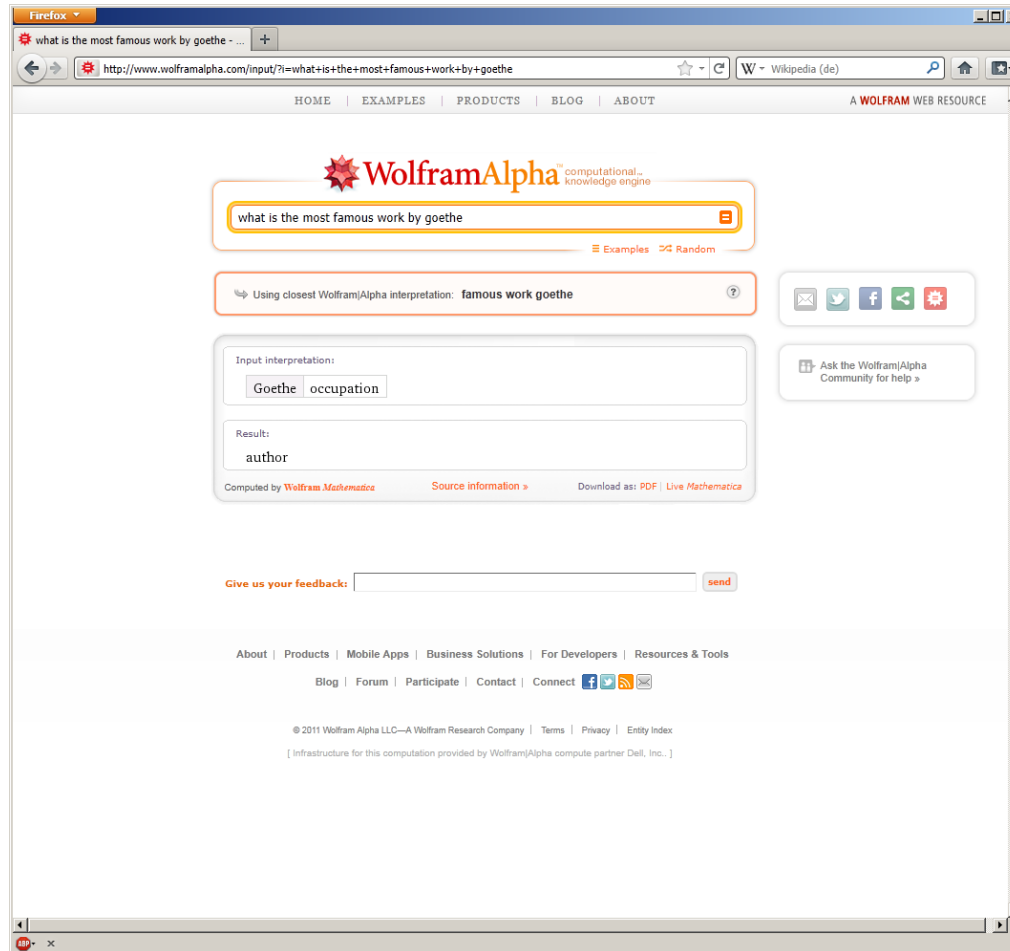
Get the optimized mobile experience...

Wolfram|Alpha App for Android »

Related Wolfram|Alpha Queries

- Reinheim vs Frankfurt
- host city of the 2016 olympi...
- skymap for Darmstadt
- Hessisches Staatsarchiv

Beispiel: Wolfram Alpha

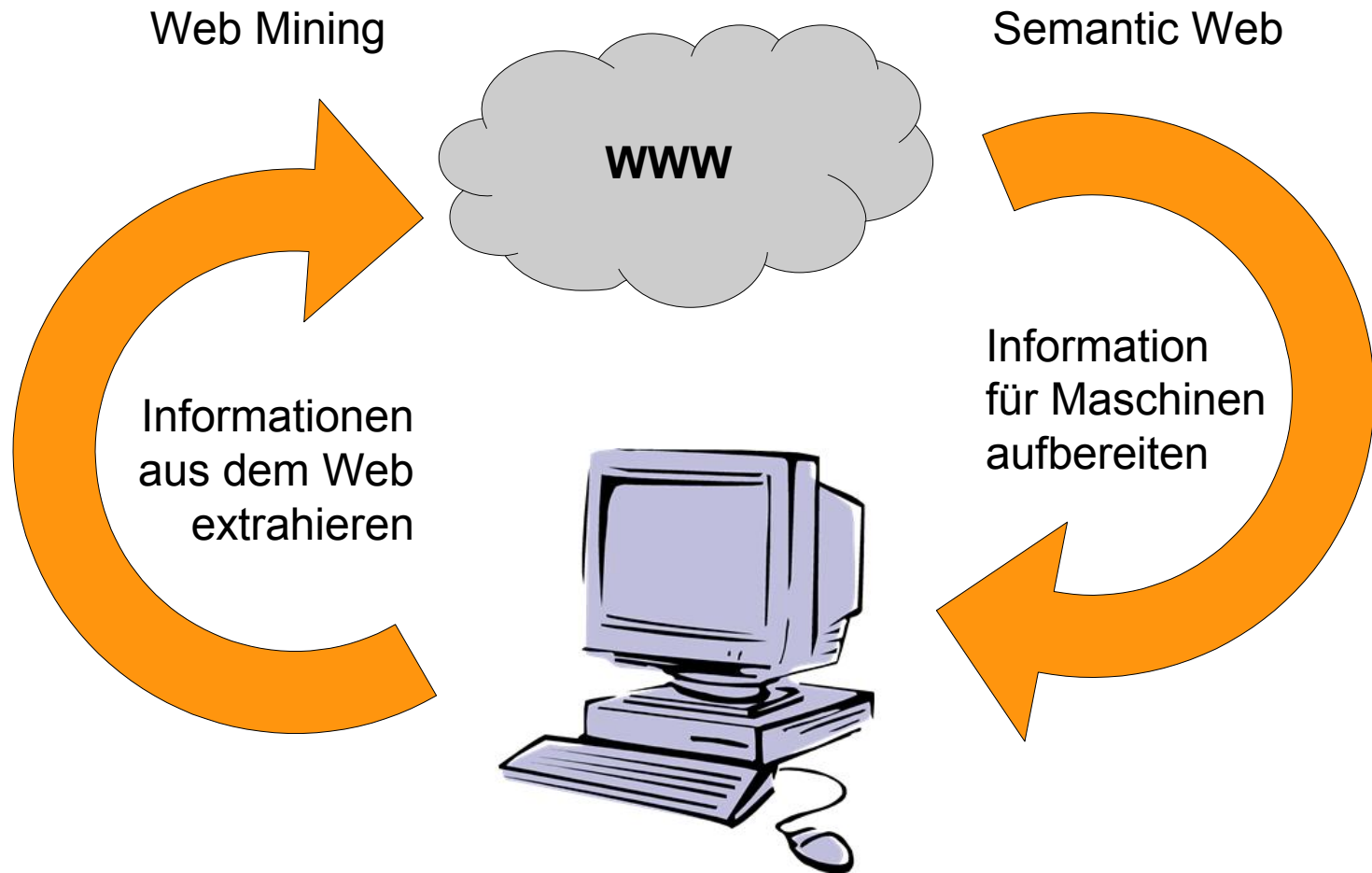


The screenshot shows a Firefox browser window displaying the Wolfram Alpha website. The address bar shows the URL <http://www.wolframalpha.com/input/?i=what+is+the+most+famous+work+by+goethe>. The search input field contains the text "what is the most famous work by goethe". Below the input field, the result is displayed as "Using closest Wolfram|Alpha interpretation: famous work goethe". The input interpretation is "Goethe" and the result is "author". The page includes navigation links (HOME, EXAMPLES, PRODUCTS, BLOG, ABOUT) and a footer with copyright information: © 2011 Wolfram Alpha LLC—A Wolfram Research Company | Terms | Privacy | Entity Index. [Infrastructure for this computation provided by Wolfram|Alpha compute partner Dell, Inc.]

Lösungsansätze



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Idee des Semantic Web

- Information in maschinenlesbarer Form bereitstellen
- (Semantische) Verweise zwischen Seiten nutzbar machen
- Schlussfolgern ermöglichen
- Komplexe (nützliche!) Abfragen ermöglichen

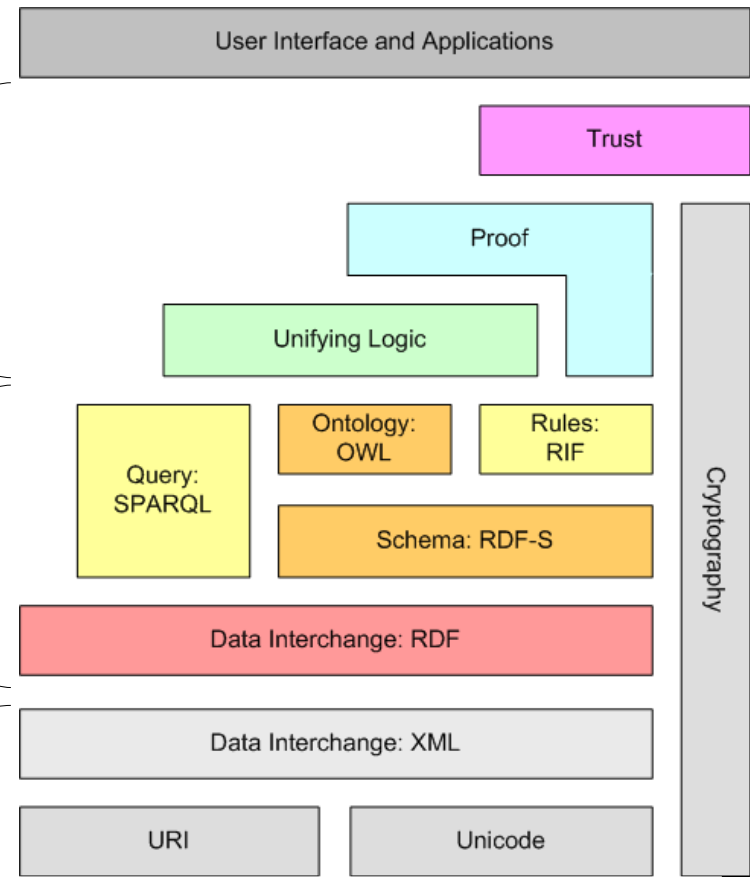
Semantic Web – Aufbau



here be dragons...

Semantic-Web-
Technologie
(Fokus der Vorlesung)

Technische
Grundlagen

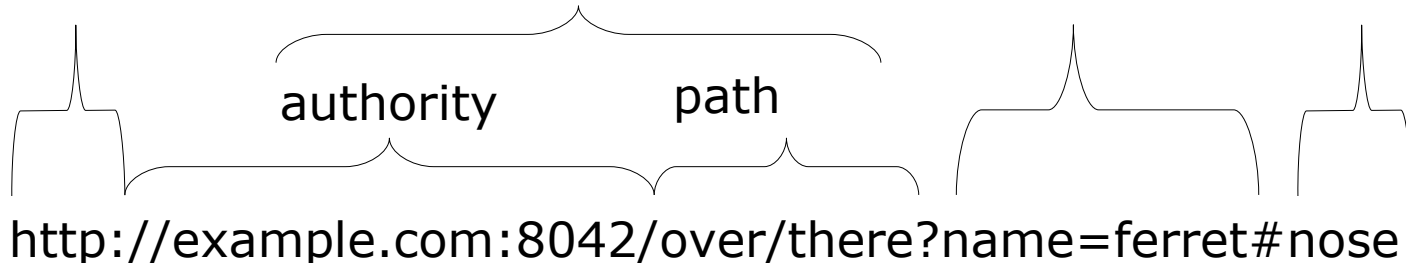


Berners-Lee (2009): *Semantic Web and Linked Data*
<http://www.w3.org/2009/Talks/0120-campus-party-tbl/>

Uniform Resource Identifiers (URIs)

- Als „Universal Resource Identifier“ erstmals vorgeschlagen von Tim Berners-Lee (1994): IETF RFC 1630
- Standard: IETF RFC 3986 (2005)
- Dienen zur Benennung und zum Auffinden von Ressourcen im Internet

URI = scheme ":" hier-part ["?" query] ["#" fragment]



URIs und URLs

- Uniform Resource Locators (IETF RFC 1738, 1994) sind eine Untermenge von URIs
- URIs können *beliebige* Dinge mit *beliebigen* Namen versehen
- ein URL ist der Name einer Resource [im Internet]
- URL-typische Präfixe:
 - http
 - ftp
 - mailto
 - telnet
 - file
 - ...

URLs im Web

- Häufigste Verwendung:
(HTTP-)Links
- Links haben in der Regel
keine Metainformation!

Tim Berners-Lee

Tim Berners-Lee (geboren
1955) gilt als einer der Erfinder
des [World Wide Web](#).

...

Das World Wide Web

Der Grundstein für das World
Wide Web wurde in den 90er-
Jahren durch [Tim Berners-Lee](#)
am [CERN](#) gelegt.

...

Zeichensätze im Web

- ASCII („American Standard Code for Information Interchange“) ISO 646 (1963), 127 Zeichen, davon 95 druckbar:

```
!"#$%&'()*+,-./0123456789:;<=>?  
@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_  
`abcdefghijklmnopqrstuvwxyz{|}~
```

- Erweiterung auf 8 Bit: ISO 8859-1 bis -16 (1998)
 - deckt Zeichen der europäischen Sprachen ab
 - bekannt ist insbesondere 8859-1 („Latin-1“)
- Das Web spricht aber viel mehr Sprachen...

وللعبّ علامات يقفوها الف
فأولها إدمان النظر، والعج
سراؤها، والمعبرة لضمائرها
بر لا يطرف، يتنقل بتنقل
ن مال، كالحرباء مع الشمس

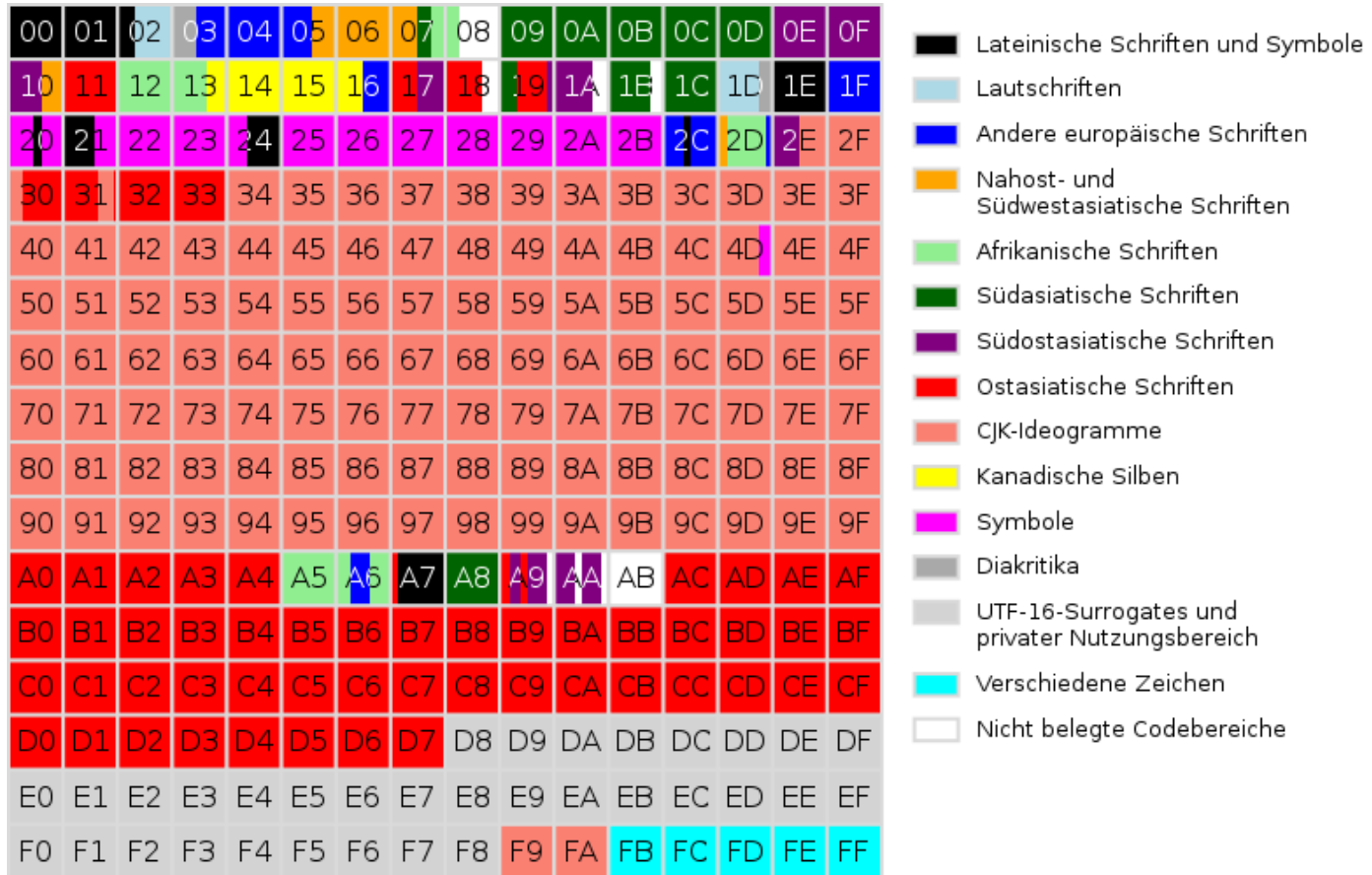
我爱中国
国中爱我

Unicode

- ISO 10646
 - erste Version 1991 (Europa, Nahost, Indien)
 - Version 6.0 im Oktober 2010
 - 17 Codebereiche à 16 Bit
 - deckt selbst exotischste Sprachen ab



Unicode



Quelle: Wikimedia Commons

Informationsrepräsentation in XML

XML (eXtensible Markup Language)

- Standardisiert vom W3C (1998)
- Universelles Datenaustauschformat



```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

XML: Konzepte



- Tags (beliebig definierbar):
 - Paare:
`<physician> ... </physician>`
 - Empty-Element-Tags:
`<young />`
- Attribute:
`<physician location="Smalltown">`
- Schachtelung (genau ein Root-Element!):
`<physician>`
`<address> ... </address>`
`</physician>`

XML: wohlgeformte Dokumente



```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  <telephone>
    <number>+44 123 456789</number>
  </address>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

HTML und XML

- HTML-Dokumente sind i.d.R. keine wohlgeformten XML-Dokumente!

```
<p> Look at this! <img src=smiley.gif> <br>
```

- XHTML: HTML als wohlgeformte XHTML-Dateien
- Vom W3C standardisiert (seit 2000)

```
<p> Look at this!  <br/> </p>
```



XPath: Zugriff auf XML-Daten



- Abfragesprache für XML
- Standardisiert vom W3C (1999, Version 2.0: 2010)

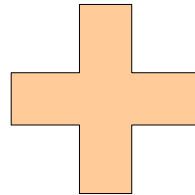
`/physician[name='Dr. Mark Smith']/telephone/number`

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

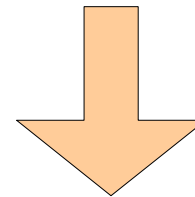
XSLT: Transformation von XML-Dokumenten

- Verarbeitung von XML-Dokumenten basierend auf Stylesheets
- Standardisiert vom W3C (1999)
- Verwendet XPath

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```



```
<xsl:template match="/pyhsician">
  <b>
    <xsl:value-of select="name"/>
  </b>
</xsl:template/>
```



```
<b>Dr. Mark Smith</b>
```

Namensräume in XML

- Elemente gleichen Namens können mehrfach vorkommen
- ...aber mit anderem Inhalt (und anderer Semantik!)
- Wie können wir solche Elemente unterscheiden?

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

Namensräume in XML



- Mit Präfix unterscheidbar (Schreibweise: `prefix:name`)
- Ein Namensraum ist ein URI
- Default-Namensraum möglich

```
<physician xmlns      ="http://www.med.com/physician"
           xmlns:addr="http://www.med.com/addr">
  <name>Dr. Mark Smith</name>
  <addr:address>
    <addr:street>Main St.</addr:street>
    <addr:number>14</addr:number>
    <addr:city>Smalltown</addr:city>
  </addr:address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

XML: Document Type Definition (DTD)



- Definiert gültige Elemente für einen XML-Dokumenttyp
 - Name
 - zulässige Attribute
 - zulässige Kind-Elemente
- DTD ist Teil der W3C-XML-Spezifikation
- XML-Dokumente, die auf eine DTD passen, heißen „gültig“

XML: Document Type Definition (DTD)

```
<!DOCTYPE physician [  
  
<!ELEMENT physician (  
  name,  
  address*,  
  telephone?,  
  fax?,  
  hours)>  
  
<!ELEMENT address (  
  street,  
  number,  
  city)>  
  
<!ELEMENT street (#PCDATA)>  
  
  ...  
  
>
```

```
<!DOCTYPE physician SYSTEM  
  "physician.dtd">  
  
<physician>  
  <name>Dr. Mark Smith</name>  
  <address>  
    <street>Main St.</street>  
    <number>14</number>  
    <city>Smalltown</city>  
  </address>  
  <telephone>  
    <number>+44 123 456789</number>  
  </telephone>  
  <hours>  
    <monday>9-11 am</monday>  
    <tuesday>9-11 am</tuesday>  
    ...  
  </hours>  
</physician>
```

XML: Document Type Definition (DTD)

- Definition von Kind-Elementen und deren Reihenfolge:
`<!ELEMENT address(street,nr,addtl*,zip,city,state?) >`
 - ? und * markieren optionale und wiederholbare Elemente
- Definition von Attributlisten:
`<!ATTLIST person title CDATA>`
 - mögliche Zusätze: #REQUIRED, #FIXED, #IMPLIED, "..."
 - Aufzählung zulässiger Werte: (dr|prof)
- Definition von Entitäten:
`<!ENTITY sw "Semantic Web">`
 - Können als Abkürzung verwendet werden: `&sw;`

- W3C-Standard (seit 2004)
- XML-Schema-Dateien sind selbst XML-Dateien
- Flexibler als DTDs:
 - Minimale und maximale Anzahl von Elementen
 - Kombinationen von Elementen (entweder oder, Auswahl ohne Reihenfolge, ...)
 - Datentypen (Zahlen, Daten, ...), eigene Definitionen möglich
 - Unterstützung von Namespaces
 - Modulare Schemata möglich

XML Schema



```
<xs:schema elementFormDefault="qualified"
xmlns:xs="http://www.w3.org/2001/XMLSchema">
```

```
<xs:element name="physician">
  <xs:complexType>
    <xs:sequence>
      <xs:element name="name"
        type="xs:string">
      <xs:element name="address">
        <xs:complexType>
          <xs:sequence>
            <xs:element name="street"
              type="xs:string">
            ...
          </xs:sequence>
        </xs:complexType>
      </xs:element>
      ...
    </xs:sequence>
  </xs:complexType>
</xs:element>
</xs:schema>
```

```
<physician xmlns:xsi=
"http://www.w3.org/2001/XMLSchema-instance"
xsi:noNamespaceSchemaLocation=
"physician.xsd">
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

XML Schema – modulare Schemata



```
<xs:schema elementFormDefault="qualified"
xmlns:xs="http://www.w3.org/2001/XMLSchema"
xmlns:addr="http://www.address.com/">

<xs:import
  namespace="http://www.address.com/"
  schemaLocation="address.xsd"/>
<xs:element name="physician">
  <xs:complexType>
    <xs:sequence>
      <xs:element name="name"
        type="xs:string">
      <xs:element ref="addr:address" />
      ...
    </xs:sequence>
  </xs:complexType>
</xs:element>
</xs:schema>
```

```
<xs:schema elementFormDefault="qualified"
xmlns:xs="http://www.w3.org/2001/XMLSchema">

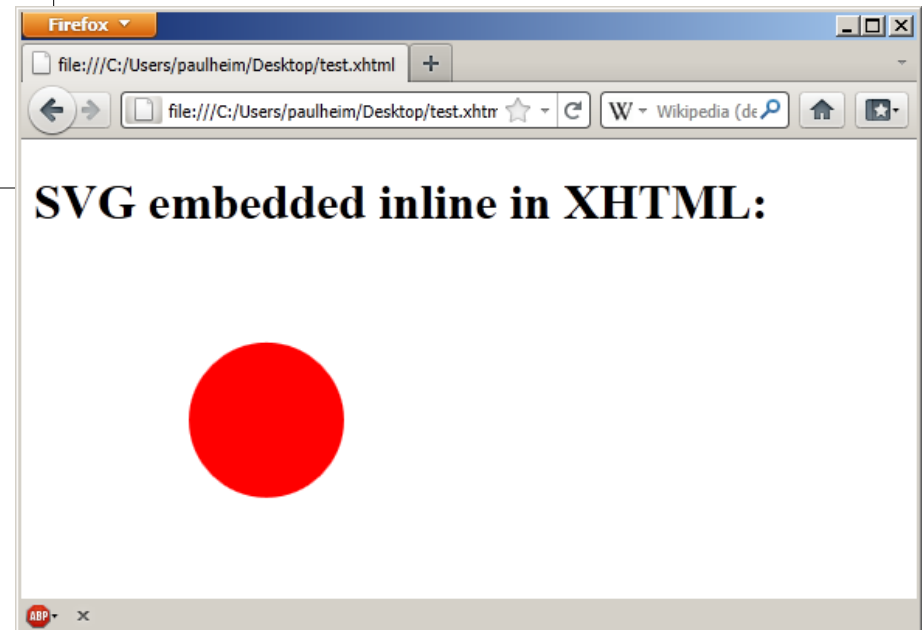
  <xs:element name="address">
    <xs:complexType>
      <xs:sequence>
        <xs:element name="street"
          type="xs:string">
          ...
        </xs:sequence>
      </xs:complexType>
    </xs:element>
  </xs:schema>
```

Beispiel: Modulare Schemata in XHTML



TECHNISCHE
UNIVERSITÄT
DARMSTADT

```
<html xmlns="http://www.w3.org/1999/xhtml"
      xmlns:svg="http://www.w3.org/2000/svg">
  <body>
    <h1>SVG embedded inline in XHTML:</h1>
    <svg:svg width="300px" height="200px">
      <svg:circle cx="150" cy="100" r="50"
        fill="#ff0000"/>
    </svg:svg>
  </body>
</html>
```



https://developer.mozilla.org/En/SVG:Namespaces_Crash_Course



RELAX NG



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Alternative zu XML Schema
- Standardisiert (ISO/IEC 19757-2, 2003)



- Unterschiede zu XML Schema:
 - weniger flexibel in Kardinalitäten
 - kein expliziter Link zwischen XML-Dokument und Schema
 - kein eigenes Typsystem (kann XML Schema Datatypes verwenden)
 - bessere Unterstützung von schwach strukturierten Inhalten
 - XML-basierte und kompakte Darstellung möglich



```
<element name="physician"
xmlns="http://relaxng.org/ns/structure/1.0">
  <element name="address">
    <group>
      <element name="street">
        <text/>
      </element>
      <element name="number">
        <text/>
      </element>
      <element name="city">
        <text/>
      </element>
    </group>
  </element>
  ...
</element>
```

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

Schematron



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Fokus: *Validierung* von XML-Dokumenten
- Standardisiert (ISO/IEC 19757-3, 2006)
- Verwendet Regeln zur Validierung
- Regeln basieren auf XPath-Ausdrücken
- Ausführung der Validierung mit XSLT möglich
- Fehlermeldungen werden direkt im Schema definiert



Schematron



```
<?xml version="1.0" encoding="utf-8"?>
<schema
  xmlns="http://purl.oclc.org/dsdl/schematron"
  <title>Physician validation schema</iso:title>

  <pattern>
    <rule context="physician">
      <assert test="address">A physician must have an
        address</assert>
      ...
      <assert test="hours/*">Hours must not be empty
        </assert>
    </rule>
  </pattern>
</schema>
```

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

Schematron – Ausführung

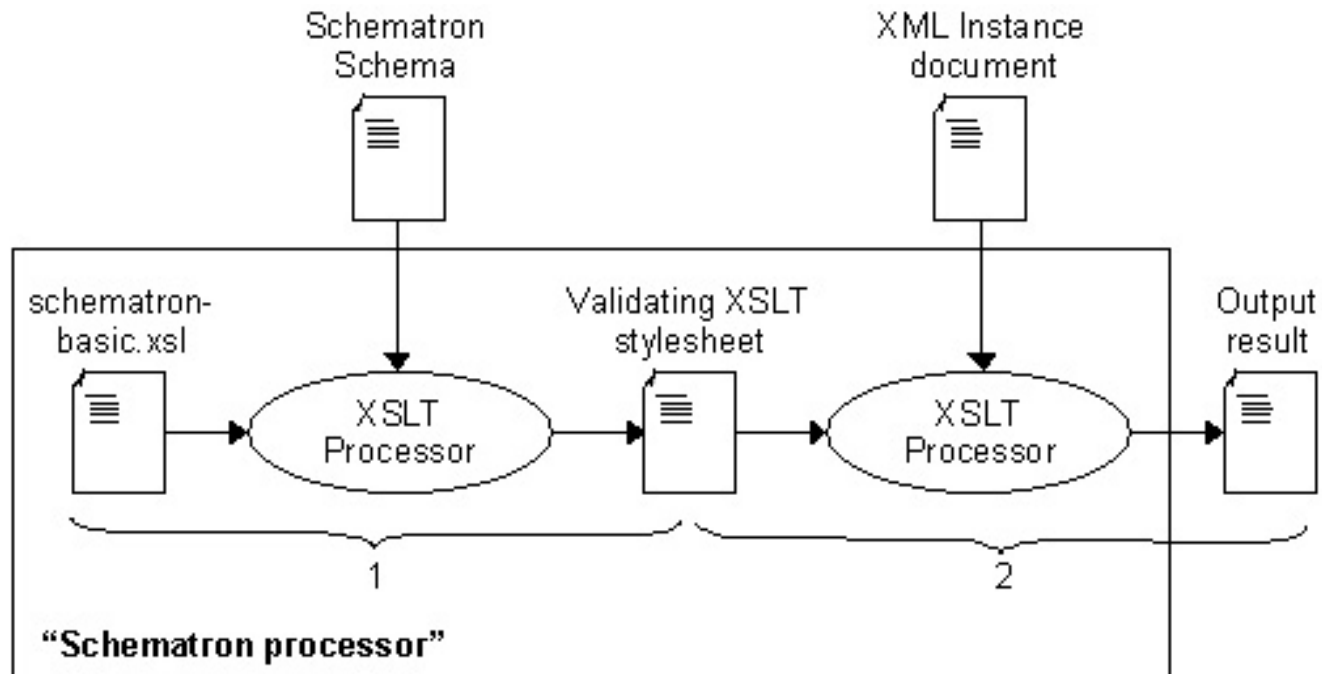


Figure 1: Schematron processing

Eddie Robertson: Combining Schematron with other XML Schema Languages.
http://www.topologi.com/resources/schtrn_xsd_paper.html

XML Schema, DTD & Co – Was wird hier definiert?



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Syntax – σύνταξις („zusammen“ + „Ordnung“)
 - Welche Elemente gibt es?
 - Wie werden sie angeordnet?
 - Welche Kombinationen sind zulässig?
- Gegensatz: Semantik - σημαίνει („bezeichnen“)
 - Wie interpretiert man den Inhalt von Elementen?
 - In welchem *inhaltlichen* Zusammenhang stehen Elemente?



Syntax und Semantik: Exkurs in die Sprachwissenschaft



TECHNISCHE
UNIVERSITÄT
DARMSTADT

- Syntax: wie bildet man korrekte Wörter und Sätze?

„Dieser Satz kein Verb.“

„Die träumende Lampe ~~schenkst~~ schenkt dem müden Wasserhahn unaufmerksam ~~eine~~ einen abgesägten Saft.“

- Semantik: was bedeutet ein Wort/Satz/Text?

Syntax und Semantik: Exkurs in die Sprachwissenschaft



TECHNISCHE
UNIVERSITÄT
DARMSTADT

The screenshot shows the Duden online interface for the word 'Lampe, die'. The page includes a search bar, navigation links, and a detailed entry. The 'Bedeutungsübersicht' section is circled in red, containing the following text:

Bedeutungsübersicht

1. als Träger einer künstlichen Lichtquelle (besonders von Glühbirnen) dienendes, je nach Zweck sehr unterschiedlich gestaltetes, hängendes, stehendes oder auch frei bewegliches Gerät
2. (besonders Fachsprache) künstliche Lichtquelle (z. B. Glühlampe)

The 'Inhalte' sidebar on the right lists various linguistic aspects: Rechtschreibung, Bedeutungsübersicht, Wussten Sie schon?, Synonyme zu Lampe, Aussprache, Herkunft, Grammatik, Typische Verbindungen (computergeneriert), Bedeutungen, Beispiele und Wendungen, and Blättern. Below this, there are links for Drucken, Zitieren, Wortvorschlag, Hilfe zum Wörterbuch, and Weitersagen.



XML Schema, DTD & Co – Was wird hier definiert?

Personalverzeichnis
des Krankenhauses:

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

(wahrscheinlich)
die Privatadresse

?
=

Gelbe Seiten:

```
<physician>
  <name>Dr. Mark Smith</name>
  <address>
    <street>Main St.</street>
    <number>14</number>
    <city>Smalltown</city>
  </address>
  <telephone>
    <number>+44 123 456789</number>
  </telephone>
  <hours>
    <monday>9-11 am</monday>
    <tuesday>9-11 am</tuesday>
    ...
  </hours>
</physician>
```

(wahrscheinlich)
die Adresse der Praxis

XML Schema, DTD & Co – Was wird hier definiert?



- XML Schema / DTD definiert die *Syntax* eines XML-Dokuments, nicht die *Semantik*
- Tag-Namen sind für Maschinen nicht a priori interpretierbar
 - das macht die Informationssuche nicht leichter...
 - Semantik der Daten (ver-)steckt hart verdrahtet in der Anwendung
- Das Semantic Web soll hier Abhilfe schaffen
 - *Semantic Web ist/kann mehr als XML!*

```
<2nf3oiü*>
  <34f0>Dr. Mark Smith</34f0>
  <rmd4935r>
    <e2m4>Main St.</e2m4>
    <dur3>14</dur3>
    <jfa34>Smalltown</jfa34>
  </rmd4935r>
  <d24r3fmö>
    <deß5>+44 123 456789</deß5>
  </d24r3fmö>
  <vsfif>
    <f02>9-11 am</f02>
    <fj9>9-11 am</fj9>
    ...
  </vsfif>
</2nf3oiü*>
```

Zusammenfassung

- Probleme des klassischen Web
 - Nicht für Maschinen nutzbar
- URIs
 - eindeutige Identifier für Ressourcen
 - URL = dereferenzierbarer URI
- Unicode
 - ein einheitlicher Zeichensatz für alle
- XML
 - XPath
 - XSLT
 - Schemasprachen

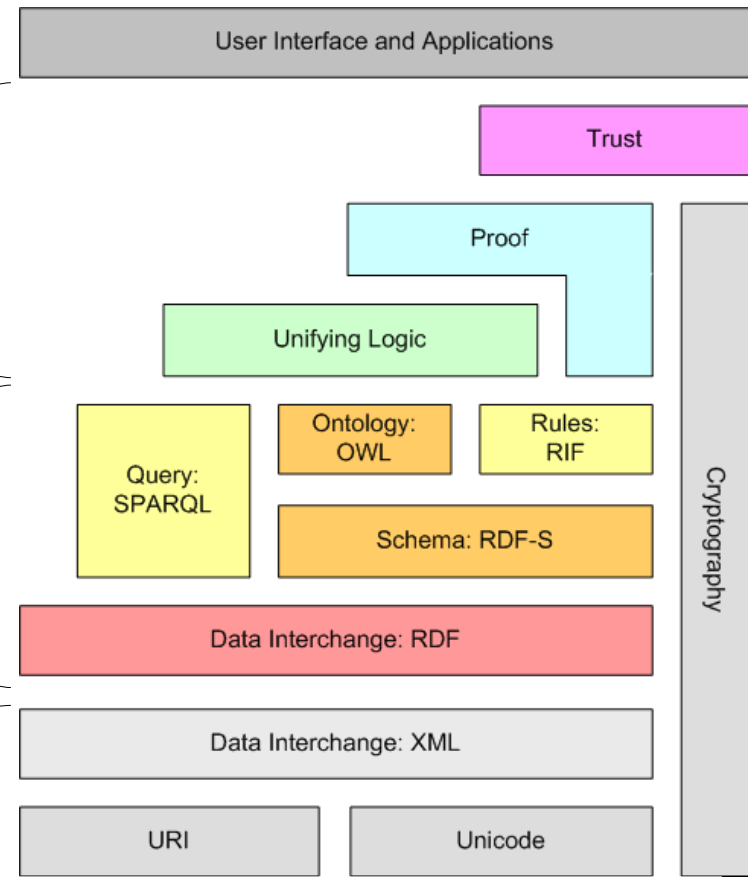
Semantic Web – Aufbau



here be dragons...

Semantic-Web-
Technologie
(Fokus der Vorlesung)

Technische
Grundlagen



Berners-Lee (2009): *Semantic Web and Linked Data*
<http://www.w3.org/2009/Talks/0120-campus-party-tbl/>

Vorlesung Semantic Web



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Vorlesung im Wintersemester 2012/2013

Dr. Heiko Paulheim

Fachgebiet Knowledge Engineering