

# Semanttinen Web ja Web palvelut

## Sisältötuotannon ja palveluiden tulevaisuus

Internetissä  
Prof. Eero Hyvönen

Helsingin yliopisto ja Tietotekniikan tutkimuslaitos HIIT  
Semantic Computing Research Group  
<http://cs.helsinki.fi/group/seco/>

20-Oct-04

1

## Sisältö

- WWW tänään
  - Palvelut ja tiedonhaku: ongelmia
  - Tiedon esitys: merkkausekset
  - Toiminnan esitys: ohjelmointi
- Tulevaisuuden WWW
  - Semantic Web -visio
  - Tasot ja standardit
  - Sovellusalueita
  - Demonstraatio semanttisesta portaalista

20-Oct-04

2

## WWW tänään

20-Oct-04

3

## WWW:n perusta

- URI osoitteet: resurssit
  - Sivustot, dokumentit, kuvat jne.
- HTML-kieli
  - WWW-sivujen julkaiseminen
  - Hyperlinkit
- HTTP ym. protokollat
  - Hyperteksti
  - Sähköposti
  - Keskusteluryhmät

20-Oct-04

4

## WWW tänään: palvelut ja tiedonhaku

- Toiminnalliset palvelut
  - Pankit, kaupat, virastot jne.
- Tiedonhaun palvelut
  - Hakukoneet (Google, AltaVista jne.)
  - Portaalit, hakemistot
    - Yahoo! jne.
    - Kallis ja hankala ylläpito
  - Arkistot, museot, kirjastot jne.
  - Tietokannat eri sovelluksissa

20-Oct-04

5

## Toiminnalliset palvelut: ongelmia

- Yhteisten sanastojen ja standardien puute
  - Tiedot: esim. eri yritysten tuotekatalogit
  - Prosessit: esim. ostoprosessin kuvaaminen
- Järjestelmien monimutkaisuus ja kalleus
  - Esim. EDIFACT
- Toimintojen/palveluiden monimutkaisuus
  - Esim. liikematkan osapalveluiden yhdistäminen

20-Oct-04

6

## Tiedonhaku: ongelmia

- Laadun mittaaminen
  - Recall: Kuinka monta % relevanteista löytyy
  - Precision: Kuinka monta % löytyneistä relevantteja
  - Relevance: Vastaako tulos haluttua
- Suomessa sanamuodot ja johdokset ongelmana

20-Oct-04

7

- Hakusana voi esiintyä epärelevantissa dokumentissa
  - "This page is *not* about *politics*"
- Synonyymien tunnistaminen
  - Venus  $\neq$  Aamutähti  $\neq$  Iltatähti
  - => low recall
- Homonyymien tunnistaminen
  - Nokia -> firman ja kaupungin nimi
  - => low precision

20-Oct-04

8

- Yleistermien käyttö vaikeaa
  - Esim: Pohjoismaat -> Suomi, Ruotsi, ...
  - Käyttäjän tunnettava eritystermit
- Fraasien käyttö
  - Esimerkiksi "Helsingin yliopisto"
- Relevanssi
  - Haun tuloksena yleensä paljon osumia
  - Dokumenttien järjestys "hyvyyden" mukaan
    - vrt. Googlen innovaatio

20-Oct-04

9

- Implisiittinen tieto
  - Joulus-sanalla ei välttämättä löydy pukki-sivua
- Hajautunut tieto
  - Esim. hae tutkimusryhmään kuuluvien julkaisut

20-Oct-04

10

## Tiedon esitys: merkkaukielet (markup languages)

20-Oct-04

11

## Merkkauskielten idea

- Ympäristöriippumattomia standardeja dokumenttien
  - luomiseen
  - hallitsemiseen
  - siirtämiseen
- Dokumentit tekstitiedostoja
  - Avoin yksinkertainen formaatti
  - Käytössä kaikilla HW/SW-alustoilla
  - Helppo muokata, tallentaa, lukea, siirtää
  - Käytettävissä tulevaisuudessakin

20-Oct-04

12

- Ideana erottaa *rakenne, sisältö ja ulkoasu*
  - Kuvataan rakenne yleisesti merkkauksilla (ohjelmoija)
    - Esim. HTML: <H1> Otsikko </H1>
  - Kuvataan sisältö (ohjelmoija)
    - Esim. XML: <OSOITE> Tietotie 3 </OSOITE>
  - Ulkoasusta päättää lukija (selain)
    - Esim. PC, kännykkä tms.

20-Oct-04

13

- Käytännössä työnjako menee helposti sekaisin
  - Esim. tekstin korostus tai koko on selaimen asia:
    - <EM> Korostettu teksti </EM>
      - *Loogisesti oikein; ei oteta kantaa siihen miten korostus tehdään*
    - <I> Korostettu teksti (kursiivi) </I>
      - *Loogisesti väärin, jos ajatuksena on vain korostaa tekstiä*

20-Oct-04

14

## HTML

- **H**yper **T**ext **M**arkup **L**anguage
- WWW-sivujen kirjoittamiskieli
- Kaikkien selaimien tukema
  - Tuettu versio kuitenkin vaihtelee!
    - Esim. HTML:n uudet ominaisuudet, Java-tuki jne.
  - Selaimet voivat näyttää sivuja hieman eri lailla
    - Esim. puutteellisten kuvausten oletusarvoiset täydentämiset, fonttivalikoimat jne.
  - Sivut on siksi aina hyvä testata eri selaimilla!

20-Oct-04

15

## SGML

- Standard Generalized Markup Language
- Metakieli, jonka avulla voidaan määrittellä merkkaukiskieliä
  - Data Type Definition (DTD) -määrittelyt
  - Monimutkainen
- ISO standardi 1986
- Esimerkiksi HTML on yksi SGML:n määrittely

20-Oct-04

16

## XML

- **E**xtensible **M**arkup **L**anguage
- SGML:n yksinkertaisempi osajoukko (20%)
  - XML korvannut käytännössä SGML:n
- Voidaan määrittellä sovelluskohtaisia markup-kieliä
  - <HENKILO>
    - <NIMI> Onni Opiskelija </NIMI>
    - <PUHELIN> 123456 </PUHELIN>
  - </HENKILO>

20-Oct-04

17

- Uuden kielen syntaksin määrittely
  - DTD-kuvaus: oma määrittelykieli
  - XML Schema: määrittely XML-perustaisesti
- Ulkoasu: eXtensible Style Language (XSL)
  - XSLT: muunnokset, esim. HTML-sivuksi
  - XSL-FO: oma layout-kieli

20-Oct-04

18

## XML-esimerkki (DTD)

```
<?xml version="1.0"?>
<!DOCTYPE CONTACTS (CONTACT+)>
[
<ELEMENT CONTACT (NAME, PHONE+, ADDRESS, IMG)?>
<ELEMENT NAME (#PCDATA)>
<ELEMENT PHONE (#PCDATA)>
<ATTLIST PHONE type (home | work | gsm) "work">
<ELEMENT ADDRESS (#PCDATA)>
<ELEMENT IMG EMPTY>
<ATTLIST IMG src CDATA #REQUIRED>
<ENTITY Uni "University of Helsinki"> ]
<CONTACTS>
<CONTACT>
<NAME>Mika Klemettinen</NAME>
<PHONE type="work">191 44159</PHONE>
<PHONE type="gsm">050-54 78 595</PHONE>
<ADDRESS-&Uni;</ADDRESS>
<IMG src="mika.jpg"/>
</CONTACT>
</CONTACTS>
```

20-Oct-04

19

## Miksi XML?

- Samalle sisällölle eri ulkoasuja
  - Eri laitteet (PC, kännykkä, ...)
  - Eri sovellukset (WWW-sivu, painettu kirja, ...)
- Sisällön/rakenteen hyödyntäminen
  - Esim. parempi osumatarkkuus hakukoneissa
- Laadun kontrollointi
  - Syntaksin tarkistus mahdollista

20-Oct-04

20

- XML on Webin perusta jatkossa
  - Tietojen koodaus *avoimessa* muodossa
    - Runsaasti standardeja eri aloille
  - *Avoimet* rajapinnat Java ym. kieliin
    - Ohjelmallinen sivujen käsittely

20-Oct-04

21

## Merkkauskielten merkitys

- Muodostavat WWW:n perustan
  - Helppokäyttöisiä näytettäviä sivuja käyttäjille
  - Helppoja tehdä toteuttajan näkökulmista
  - Avoimet yhteiset standardit
- Valmistajariippumattomuus
- Stabiilisuus tiedostoformaattien muutoksia vastaan
  - Sivut ovat yksinkertaisia tekstitiedostoja
- Sovellusaluekohtaiset standardikiel

20-Oct-04

22

## Standardointi

- WWW-kehityksen yleiskoordinointi
  - WWW Consortium (www.w3.org)
    - Valmistajien, operaattoreiden jne. yhteistyöelin
    - Laatii WWW-suosituksia
- Sovellusaluekohtaiset organisaatiot
  - ISO: Eri alat paitsi sähkö/elektroniikka
  - IEC, CEN, UN/CEFACT, OASIS, ...
  - Loputtomasti työryhmiä eri aloilla

20-Oct-04

23

## Toiminnan esitys: ohjelmointi

20-Oct-04

24

## Ohjelmointi

- WWW-selaimen sovellusohjelmointi
  - Hajautettu toiminnallisuus
- WWW-palvelimen sovellusohjelmointi
  - Keskitetty toiminnallisuus

20-Oct-04

25

## WWW-selaimen sovellusohjelmointi

- Java-appletit (asiakaspää)
  - Java-ohjelma siirtyy palvelimelta selaimeseen
  - Ohjelma ajetaan asiakaskoneessa
- Dynamic HTML (asiakaspää)
  - ECMAScript (JavaScript, J Script)
    - HTML-koodin sekaan ajettavia ohjelmia (script)
  - Cascading Style Sheets (CSS)
    - Yleisiä tyyliäärittelyjä HTML-kielen elementeille
  - Domain Object Model (DOM)
    - Sivun oliomalli skriptejä varten

20-Oct-04

26

## WWW-palvelimen sovellusohjelmointi

- Server Side Includes (SSI)
  - HTML-koodilla korvattavia koodeja HTML-sivulla
    - Esim. päivitys tai muu dokumentin osa
    - Palvelin hoitaa korvaamisen ennen sivun lähettämistä

20-Oct-04

27

- Server Pages -systemit (ASP, JSP)
  - HTML-sivulla myös Javaa tms. ohjelmointikieltä
  - Koodit ajetaan ja korvataan HTML-tuloksella
  - Palvelimella ohjelma luo HTML-sivut
    - Esim. tietojen haku tietokannasta
  - Tulos lähetetään selaimelle
- TAG Libraries
  - Koodit korvataan omilla merkkauksilla

20-Oct-04

28

- CGI-skriptit ja servletit
  - Palvelimen ohjelma
  - Saa tiedot selaimelta esim. lomakkeella
    - PUT ja GET metodit
  - Palauttaa selaimelle HTML-tuloksen

20-Oct-04

29

## Semantic Web & Web Services:

Visio

20-Oct-04

30

## Mitä hyötyä on Webistä?

- Keskeistä WWW:n tarjoamat palvelut
  - Viestintä (email, puhe, kuva, video,...)
  - Tiedonhaku (hakukoneet, portaalit,...)
  - Toimenpiteiden suorittaminen
    - Sähköinen liiketoiminta
    - Sähköinen asiointi
    - Ym.

20-Oct-04

31

## Kehityksen este Webissä?

- WWW-palvelu  $\approx$  kone auttaa ihmistä
  - Edellyttää sisältöjen koneellista "ymmärtämistä"
- WWW:n sisällöt ovat ihmislukijaa varten
  - HTML, PDF, JPEG, ...
- Kone ei ymmärrä WWW:n sisältöjä
  - Hakukoneet, ostoagentit, verkkomönkijät jne.
  - Periaatteessa kaikki WWW-sovellukset
- => Perustavaa laatua oleva ristiriita

20-Oct-04

32

## Miten Webistä tulee älykkäämpi?

1. Älykkäämmät sovellukset
  - Sisältö pysyy samana
  - Koneesta tehdään ihmismäisempi
2. Älykkäämmin esitetty sisältö
  - Sisältö helpommin ymmärrettäväksi
  - Kone pysyy tyhempänä
- Käytännössä molempia tapoja tarvitaan
  - Yhä älykkäämmät järjestelmät käsittelevät yhä älykkäämmin esitettyjä tietoja

20-Oct-04

33

## Ratkaisumalli 1: Älykkäämmät sovellukset

- Kielen automaattisen tulkinnan vaikeus
  - Dokumenttien vapaamuotoisuus
  - Sisällön semantiikka
- Ei-tekstuaaliset sisällöt
  - Kuva, ääni, musiikki, video, ohjelmisto,...
  - Miten tulkita algoritmisesti?
- Tulkintaan ei riitä itse dokumentti
  - Tarvitaan konteksti, common sense
  - Tekoälyn perusvaikeuksia, ihmiselle helppoa!
- *Suuria tieteellisteknisiä haasteita*

20-Oct-04

34

## Ratkaisumalli 2: Älykkäämmin esitetty sisältö

- Semantic Webin lähtökohta
  - Talletetaan tieto niin, että tyhempikin sen ymmärtää!
  - Ihminen tulee konetta vastaan
  - Kone voi auttaa ihmistä itsensä auttamisessa
- Kiihkeä kehitystyö käynnistynyt
  - W3C:n Semantic Web Activity 2001
  - W3C:n Web Services Activity 2002

20-Oct-04

35

## Webin sukupolvia

- 1G WWW:
  - WWW-sivut ihmisen tulkittavaksi
  - HTML-kieli
- 2G WWW:
  - Rakenteet ihmisen/koneen tulkittavaksi
  - XML-kieli
- 3G WWW: Semantic Web
  - Merkitykset ihmisen/koneen käytettäväksi
  - RDF(S)-kieli
- => Uusi perusta älykkäille WWW-palveluille
  - Kansainvälinen yhteistyö (W3C, ISO, FIPA, ym....)

20-Oct-04

36

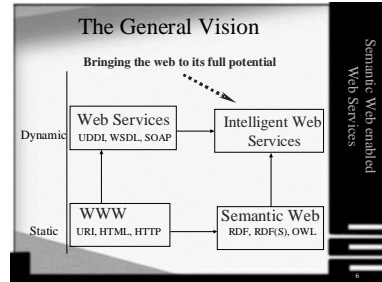
## WWW:n kaksi kehityksakselia

- 1. *Sisällön* rikastuminen (semantiikka)
  - Semanttinen Web
- 2. *Dynaamisuuden* lisääntyminen
  - Toiminnallisuuden lisääntyminen
    - Agent technologies, web services, grid computing
    - Kohti *aktiivisten* palveluiden verkkoa
  - Adaptiivisuuden lisääntyminen
    - Verkon rakenne ja yhteydet muuttuvat lennossa
    - Mobile systems, ambient computing

20-Oct-04

37

## Älykkäät verkkopalvelut



(Dieter Fensel, 2002)

20-Oct-04

38

## Semantic Web & Web Services:

Teknologioita

20-Oct-04

39

## Semantic Web: Technology push

<b>Trust level</b> Digital signature, annotations...	<b>Planning</b> CPR, SPAR, PDDL, ... <b>Processes</b> BPML, WFDL, PSL, ... <b>Services</b> UDDI, WSDL, DAML-S, ... <b>Transactions</b> XML/EDI, KQML, ... <b>Communication</b> TCP/IP, HTTP, SOAP, ...
<b>Logic level</b> KR, RuleML, ...	
<b>Ontology level</b> OWL, RDFS, ...	
<b>Metadata level</b> RDF, RDFS, Topic Maps, ...	
<b>Structure level</b> XML, XML DTD, Schema, XSL, ...	
<b>Internet level</b> Unicode, URL, ...	

20-Oct-04

Laajennettu Tim Berners-Leen (W3C) "teknologiakakku"

## Metadata level

20-Oct-04

41

## Miksi XML ei ole "semanttinen"?

```

- <OSOITE>
  <NIMI>Onni Ohjelmoija</NIMI>
  <PUHELIN> 123 456 </PUHELIN>
</OSOITE>
    
```

```

& vHXQTωΣξ
  vφTφTξXφφx XθωΩsYωxωYvHφTφTξ
  vχΑςΣυTφξ θθI ιKκ vHχΑςΣυTφξ
  vHXQTωΣξ
    
```

- Semantiikka on vain nahkakansissa, ei peltikuoressa

20-Oct-04

42

## Semanttiset metakuvaukset

- Idea: rakenteelle on annettava merkitys (semantiikka) toisella tasolla
  - WWW-resurssien metakuvaukset
  - Käsitteiden loogiset kuvaukset
- Tärkeimmät Semantic Web -standardit
  - W3C: RDF(S)
  - ISO: Topic Map, XTM

20-Oct-04

43

## RDF(S)

- RDF Resource Description Framework (1999)
  - Yleinen verkkoresurssien kuvaamiskieli
  - Relaatiotietomalli, *ei* syntaksi kuten XML
- RDF Schema (2000)
  - RDF-terminologian määrittely
  - Olioajattelu WWW-kuvauksiin
    - Käsitehierarkiat, periytyminen (Class/subClass/type)

20-Oct-04

44

## Ontology level

20-Oct-04

45

## Ontologian käsite

- “Ontologia on formaali, eksplisiittinen määrittely yhteisestä käsitteistöstä” (Gruber, 1993)
  - Formaali: jämpä
  - Eksplisiittinen: konekin ymmärtää
  - Yhteinen: kommunikaatio mahdollista
- Kuvaa sovellusmaailmassa olevat käsitteet/oliot
- Ensimmäinen edellytys sille, että ihmiset ja koneet voivat ymmärtää toisiaan

20-Oct-04

46

## Ontologiat käytännössä

- Yleisiä hierarkkisia sanastoja
  - Esim. YSA, WordNet
- Yleisiä maailmanmalleja
  - Esim. CYC
- Spesifejä ammattitermistöjä
  - Esim. RosettaNet Dictionaries
- Luokittelujärjestelmiä
  - Esim. tuotteet/palvelut UN/SPSC
  - Esim. kirjastojen UDK
- Loogisia terminologisia malleja
  - Olioperustaisissa ohjelmistoissa

20-Oct-04

47

## Esimerkki: Museoalan ontologia

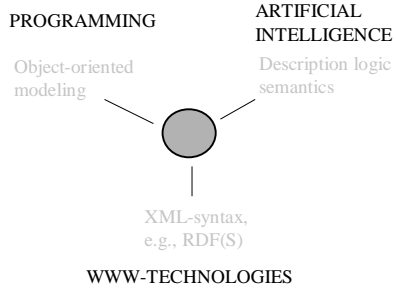
- Yli 6000:n käsitteen hierarkkinen taksonomia
- Kehittiin HY/HIIT:n Semanttisen laskennan tutkimusryhmässä
- Sovelletaan mm. MuseoSuomi -portaalissa (Finnish Museums on the Semantic Web)

20-Oct-04

48



## WHAT IS NEW?



20-Oct-04

49

## Semantic Web & Web Services:

sovelluksia, tutkimusta

20-Oct-04

50

## Sovellusalueita

- Interoperability
- Informaation haku (information retrieval)
- Tietämyksen hallinta (knowledge management)
- Sähköinen liiketoiminta, Web Services
- Profilointi ja kustomointi

20-Oct-04

51

## Interoperability

- XML-perustaiset lukemattomat ”standardit” kommunikointikieliksi (syntaksi)
- Tiedot semanttisesti yhteismitallisiksi metakuvausten + ontologioiden avulla
- Järjestelmien yhteiskäyttö: web services
  - Avoimet WWW standardit
    - SOAP, WSDL, UDDI, WSMF, DAML-S, ...
  - Legacy- ym. järjestelmien yhdistäminen web service –teknologioilla

20-Oct-04

52

## Informaation (täsmä)haku

- Seuraavan polven hakurobotit
  - Metatietojen hyödyntäminen
- Älykkäät hakemistot
- Matchmaker-sovellukset
- Semanttiset portaalit

20-Oct-04

53



20-Oct-04

54

## What is MuseumFinland?



- A semantic portal for Finnish museums to publish their collections together on the Semantic Web
- Case study in a research project on semantic interoperability on the Web
- Duration: 3/2002-3/2004

20-Oct-04

55

## Research Consortium



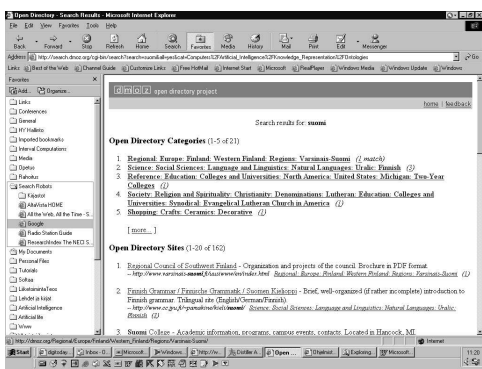
ANTIKVARIA GROUP



Co-operation with:  
Finnish National Gallery



## Open Directory Project: Koko WWW RDF(S) ontologia!



57

## Case: Promoottori

- Photo repository of the Promotion Ceremonies at the University of Helsinki
- Exhibition to be opened at the Senate Square at Helsinki city center 10/2003
  - The Helsinki University Museum
- An interesting view into the life of the university
- Complicated semantics of obscure events
- Novice users who do not know the contents before
- Team: Eero Hyvönen, Samppa Saarela, Kim Viljanen, Arvil Styrman, Jaana Tegelberg, ...



20-Oct-04

58

## Tietämyksen hallinta (knowledge management)

- Ongelmia
  - Dokumenttien monimuotoisuus
  - Maapalloistuminen -> sisältöjen hajautus
  - Tietämiskannan komplisoituminen
- SemWeb-tekniikat antavat uusia työkaluja
  - Liima heterogeenisten hajautettujen dokumenttien hallintaan
- Adoben XMP
  - Kaikki WWW-julkaisut tukevat RDF-metadattaa

20-Oct-04

59

## Sähköinen kaupankäynti: Web Services

- Miten tarjoan oman tuotteen/palvelun kansainvälisille markkinoille?
  - Visio: globaalit rekisterit ja sanastot käytettävissä (UDDI-rekisterit, RosettaNet,...)
- Miten hoidan transaktiot ja prosessit kumppaneiden kanssa?
  - Esim. tarjouspyyntöön vastaaminen
  - EDI-XML, SOAP, WSDL
- Toimintakehykset: ebXML, ...

20-Oct-04

60

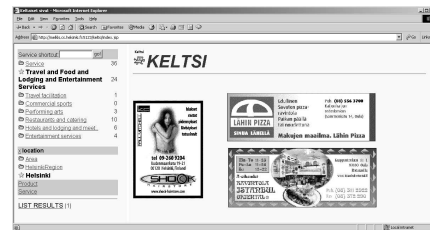
- Semantic Webin mahdollisuuksia
  - Sisältöjen ja tietojärjestelmien yhteentoimivuus
  - Eri standardien yhdistäminen
    - Esim. toimialaportaalit
  - Sisältöjen rikastaminen
  - Liiketoimintaprosessien automatisointi
    - Semanttiset kuvaukset

20-Oct-04

61

## Case: älykkäät Keltaiset sivut

- Ontologiat haun perustana
- Toisiaan tukevien palveluiden suositus



20-Oct-04

62

## Adaptiivisuus: profilointi ja kustomointi

- Ihmiset ja palvelut (P3P)
  - Omat preferenssit, tietosuoja,...
- Laitteet (CC/PP, FIPA Device Ontology)
  - Esim. MV-kännykälle ei värejä
- Dokumentit (transcoding)
  - Sisältöjen kustomointi eri laitteille ja tarpeisiin

20-Oct-04

63

## Semantic Web & Web Services

- Juna lähti jo
    - XML-standardointihankkeet, 90-luvun loppu
    - W3C Semantic Web Activity, 2001/kevät
    - EU:n OntoWeb 2001/kesä
    - W3C Web Services Activity, 2002/kevät
  - Suomessa tapahtunutta
    - "Semantic Web Kick-Off in Finland", 2001/syyskuu
    - "Towards the Semantic Web and Web Services", 2002/syyskuu
    - Tim Berners-Lee saa Millennium-palkinnon 2004/kesä
    - "Web Intelligence – Älyä verkossa", 2004/syyskuu
  - Tärkeä ala jatkossa monessa mielessä
    - Teollinen intressi
    - Tekninen mahdollisuus
    - Tieteellinen haaste
- 20-Oct-04 • Kansallinen intressi

20-Oct-04

64