



# Semanttinen Web

## Sisältötuotannon ja palveluiden tulevaisuus Internetissä

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

24-Oct-02

1

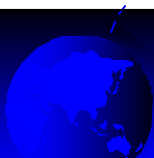
## Sisältö

- ☞ WWW tänään
  - Palvelut ja tiedonhaku: ongelmia
  - Tiedon esitys: merkkiauskielet
  - Toiminnan esitys: ohjelmointi
- ☞ Tulevaisuuden WWW
  - Semantic Web -visio
  - Tasot ja standardit
  - Sovellusalueita
  - Miksi aihepiiri on tärkeä?



24-Oct-02

2



## WWW tänään

24-Oct-02

3

## WWW:n perusta

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

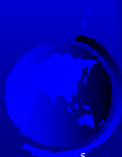


24-Oct-02

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



24-Oct-02

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



24-Oct-02

6

## Tiedonhaku: ongelmia

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



24-Oct-02

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



24-Oct-02

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



24-Oct-02

9

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



24-Oct-02

10

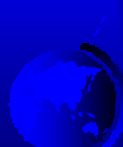
## Tiedon esitys: merkkaukielet (markup languages)

24-Oct-02

11

## Merkkaukielten idea

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



24-Oct-02

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



24-Oct-02

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*



24-Oct-02

14

## HTML

- ☞ Hyper Text Markup Language
- ☞ 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!



24-Oct-02

15

## SGML

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

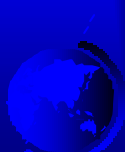


24-Oct-02

16

## XML

- ☞ Extensible Markup Language
- ☞ SGML:n yksinkertaisempi osajoukko (20%)
  - XML korvannut käytännössä SGML:n
- ☞ Voidaan määrittellä sovelluskohtaisia markup-kieliä
  - ◆ <HENKILO>  
<NIMI> Onni Opiskelija </NIMI>  
<PUHELIN> 123456 </PUHELIN>  
</HENKILO>



24-Oct-02

17

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



24-Oct-02

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 Uui "University of Helsinki" > ]
<CONTACTS>
<CONTACT>
<NAME>Mika Klemettinen</NAME>
<PHONE type="work">191 44159</PHONE>
<PHONE type="gsm">050-54 78 595</PHONE>
<ADDRESS>& Uui</ADDRESS>
<IMG src="mika.jpg"/>
</CONTACT>
</CONTACTS>
```

24-Oct-02

19

## Miksi XML?

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

24-Oct-02

20

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

24-Oct-02

21

## Merkkaukielten 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

24-Oct-02

22

## Standardointi

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

24-Oct-02

23

## Toiminnan esitys: ohjelmointi

24-Oct-02

24

## Ohjelmointi

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

24-Oct-02

25

## WWW-selaimen sovellusohjelmointi

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

24-Oct-02

26

## WWW-palvelimen sovellusohjelmointi

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

24-Oct-02

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

24-Oct-02

28

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

24-Oct-02

29

## Semantic Web & Web Services:

Visio

24-Oct-02

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.

24-Oct-02

31

## Kehityksen este Webissä?

- ☞ WWW-palvelu ≈ 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**

24-Oct-02

32

## Miten Webistä tulee älykkäämpi?

- ☞ 1. Älykkäämmät sovellukset
  - Sisältö pysyy samana
  - Koneesta tehdään ihmismäisempi
- ☞ 2. Älykkäämmiin esitetty sisältö
  - Sisältö helpommin ymmärrettäväksi
  - Kone pysyy tyhempänä
- ☞ Käytännössä molempia tapoja tarvitaan
  - Yhä älykkäämmät järjestelmät käsittelevät yhä älykkäämmiin esitettyä tietoa

24-Oct-02

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*

24-Oct-02

34

## Ratkaisumalli 2: Älykkäämmiin 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

24-Oct-02

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....)

24-Oct-02

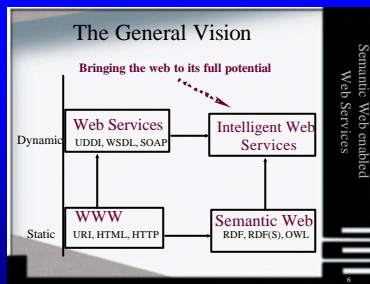
36

## WWW:n kaksi kehityksakselia

- ☞ 1. *Sisällön* rikastuminen (semantiikka)
  - Semanttinen Web
- ☞ 2. *Dynaamisuu*den 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



## Älykkäät verkkopalvelut



(Dieter Fensel, 2002)

## Semantic Web & Web Services:

Teknologioita



Laajennettu Tim Berners-Leen (W3C) "teknologiakakku"

## Metadata level

## Miksi XML ei ole "semanttinen"?

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

☞ Semantiikka on vain nahkakansissa, ei peltikuoressa



## 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

24-Oct-02

43

## RDF(S)

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

24-Oct-02

44

## Ontology level

24-Oct-02

45

## Ontologian käsite

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

24-Oct-02

46

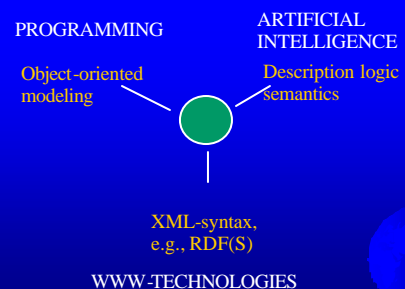
## Ontologiat käytännössä

- 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

24-Oct-02

47

## WHAT IS NEW?



24-Oct-02

48






# Semantic Web & Web Services:

sovelluksia, tutkimusta

24-Oct-02 49

# Sovellusalueita

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



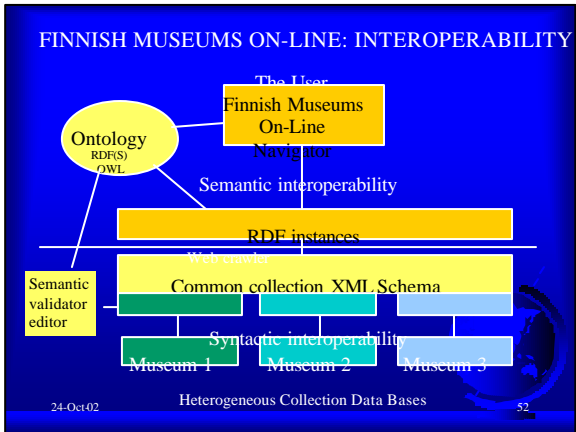
24-Oct-02 50

# Interoperability

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




24-Oct-02 51



# Informaation (täsmä)haku

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

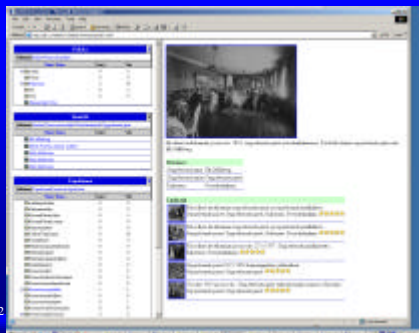


24-Oct-02 53

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

24-Oct-02 54

## Helsingin yliopisto: semanttinen kuvatietokanta



24-Oct-02

55

## 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

24-Oct-02

56

## Quid-tietosanakirja



24-Oct-02

57

## 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, ...
- Nokia: "40% alihankinnoista RosettaNetiin v. 2002"

24-Oct-02

58

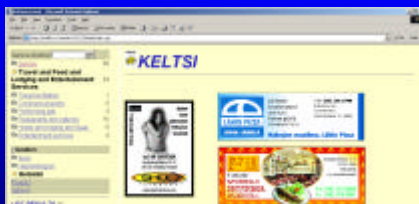
- ☞ Semantic Webin mahdollisuuksia
  - Sisältöjen rikastaminen
  - Eri standardien yhdistäminen
    - ◆ Esim. toimialaportaalit
  - Liiketoimintaprosessien automatisointi
    - ◆ Semanttiset kuvaukset

24-Oct-02

59

## Case: älykkäät Keltaiset sivut

- ☞ Ontologiat haun perustana
- ☞ Toisiaan tukevien palveluiden suositus



24-Oct-02

60

## Adaptiivisuus: profilointi ja kustomointi

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

24-Oct-02

61

## Semantic Web & Web Services

- ☞ Juna lähti jo
  - ◆ XML-standardointihankkeet, 90-luvun loppu
  - ◆ W3C Semantic Web Activity, 2001/kevät
  - ◆ EU:n OntoWeb 2001/kesä
  - ◆ "Semantic Web Kick-Off in Finland", 2001/syysy
  - ◆ W3C Web Services Activity, 2002/kevät
  - ◆ "Towards the Semantic Web and Web Services", 2002/syysy
- ☞ Tärkeä ala jatkossa monessa mielessä
  - ◆ Teollinen intressi
  - ◆ Tekninen mahdollisuus
  - ◆ Tieteellinen haaste
  - ◆ Kansallinen intressi

24-Oct-02

62

## Lisätietoja

- ☞ W3C:n Semantic Web/Web Services Activity
  - ◆ [www.w3.org](http://www.w3.org)
- ☞ Tutkimusmaailman portaali
  - ◆ [www.semanticweb.org](http://www.semanticweb.org)
- ☞ EU:n yhteistyöverkosto
  - ◆ [www.ontoweb.org](http://www.ontoweb.org)
- ☞ Semantic Web in Finland
  - ◆ [www.cs.helsinki.fi/u/eahyvone/stes/semanticweb](http://www.cs.helsinki.fi/u/eahyvone/stes/semanticweb)
  - ◆ SW Kick-Off in Finland -proceedings available (300 pp)
- ☞ Syksyn päätapahtuma Suomessa 20-21.10.
  - ◆ "Towards the Semantic Web and Web Services"
  - ◆ <http://www.xml-finland.org/events/xml2002/>

24-Oct-02

63