

Tähtitieteen data ja sen linkitykset julkaisuihin

Eva Isaksson

Suomi lähti kunnianhimoisesti mukaan suureen kansainväliseen tähtitieteen dataprojektiin. Kuulostaako tämä tuoreelta uutiselta? Kovin tuoreesta tapah- tumasta ei ollut kyse, sillä vuosiluku oli 1890. Historioitsijat pohtivat vieläkin, söikö tämä vuosikymmeniä kestänyt suurhanke suomalaisen tähtitieteen voi- mavarat. Kaikki tutkimusdataan liittyvät elementit olivat jo tuolloin olemassa. Tämä prosessoitu data päätyi julkaisuihin, pinoon kookkaita katalogeja.

Jo 1800-luvun lopussa tarjolla oli raakadataa, eli suuren valokuvakartoituksen Helsingin tai- vaanviipale, jota kuvattiin kylminä talviöinä la- silevyille. Tutkijoiden lisäksi tuli laskijoita, jouk- ko ahkeria naisia jotka mittasivat dataa levyiltä ja redusoidivat sitä taulukoiksi. Painetut katalogit löytyvät kirjastoista.

Entä jos joku haluaisi päästä käsiksi alkuperäi- seen raakadataan? Aivan ensiksi pitäisi selvittää, missä se on. Mikään luettelotieto ei asiaa paljas- ta, ja kaikkein vähiten nämä astrofotograafisen kartoituksen painetut luettelot.

Helsingin Observatorion entiset työntekijät eh- kä sentään tietävät, että lasilevyt siirrettiin joulu- kuussa 2009 Kumpulaan, Fysiikan laitoksen kel- lariin. Seuraava haaste olisi löytää vanhoista puu- laatikosta jokin ”datasetti” eli yksittäistä taivaan- kohtaa esittävä valokuvauslevy, joka toivottavasti saattaa vielä olla käyttökelpoinen.

Taivaan data on avointa kaikille

Yleensä kun puhutaan datasta, ajattelemme digi- taalisia tietomassoja. Tähtitieteessä syntyy paljon havaintodataa, joten sitä alettiin tallentaa biteik- si heti tilaisuuden tullen. NASAn kuullennoilla syntynyttä kuvamateriaalia alettiin 1970-luvulla tallentaa uudelleenlaiseen, avoimen lähdekoodin for- maattiin, joka julkistettiin 1981 FITS-formaatti- na (Flexible Image Transport System).

FITS vakiintui nopeasti tähtitieteellisen datan perusformaatiksi, ja on niin vakaa, että nyky-

tähtitieteilijä pystyy huoletta käsittelemään kol- men vuosikymmenen ikäisiä FITS-tiedostoja. Sen etuihin voidaan lukea mm. se, että tiedostot voivat olla todella suuria kooltaan, ja havainto- laitteet voivat kirjoittaa metadataa suoraan näi- hin tiedostoihin. FITS Liberator –ohjelmalla jopa pelkkä näppärä harrastelija voi muokata itselleen kuvia kaukoputken tuottamasta raakadatasta.

Tähtitieteen datan avoimuus juontaa juurensa siitä syvään juurtuneesta näkemyksestä, että tai- vas on avoin kaikille, eivätkä tähdet ole kenen- kään omaisuutta. Tutkijat on ollut helppo saada jakamaan dataa keskenään. Toisaalta ne laitteet joilla tämä avoin data on saatu kuuluvat kaik- kein kalleimpiin tutkimuslaitteisiin, joita ei lä- hetetä avaruuteen pikkurahalla.

Tähtitiede tuottaa suuria määriä dataa

Yleensä laitehankkeiden takana ovat isot yhtei- siä varoja kanavoivat organisaatiot kuten NASA, ESA (Euroopan avaruusjärjestö) tai ESO (Euro- pean Southern Observatory). Satojen tuhansien tai miljardien eurojen hintaiset laitteet ovat ku- kin ainutkertaisia lajissaan. Niillä tehtävät tai- vaankartoitukset kattavat usein tietyn aallonpi- tuuden ja koko havaittavan taivaan ja tuottavat datamääriä, joita tyyppillisesti mitataan petatavuus- sa. Mikäli havaintolaitteen suuntaamiseen haluaa vaikuttaa, tutkija joutuu anomaan ankarasti kil- pailtua havaintoaikaa. Alan käytännön mukaan tällaiselle havaintodatalle saa vuoden mittaisen



suoja-ajan, sitten data siirtyy kaikkien saataville.

Lähitulevaisuudessa toteutetaan yhä valtavampia tähtitieteen havaintoprojekteja. Etelä-Afrikkaan ja Australiaan on rakenteilla mittavia radioteleskooppeja. Valmistuttuaan tämä SKA-hanke tuottaa enemmän dataa kuin mitä koko nykyisessä internet-tietovuossa on liikenteessä.

Vielä 2000-luvun alussa tähtitieteen tutkimuksessa uskottiin ns. virtuaaliobservatorioihin. Kaikki taivaalta tallennettu data olisi tietoverkon kautta jokaisen tähtitieteilijän hyppysissä. Virtuaaliobservatorioprojekteja käynnistettiin eri puolilla maailmaa. Sitten Yhdysvaltain kansallinen tiedesäätiö NSF leikkasi 2012 kansallisen virtuaaliobservatorion rahoituksen viidesosaan. Vaik-

ka tutkimusdata on päivän kuuma sana, sen saatavuus voi olla milloin hyvänsä vaarassa, kun jokin keskeinen rahoittaja sulkee rahahanat. Alalle on yllättäen ilmaantunut yksityisiä toimijoita. Microsoft käynnisti 2008 ”World Wide Telescope”-hankkeen, jossa hyödynnetään monia virtuaaliobservatoriota varten kehitettyjä työkaluja.

Tähtitieteen datakeskukset syntyvät

Kirjastonäkökulmasta meitä kiinnostaa, miten kaikki tämä data on järjestetty, ja miten se linkittyy julkaisuihin. Euroopassa ongelmaa on pohtinut vuodesta 1971 lähtien Strasburgissa toimiva tähtitieteen datakeskus CDS (Centre des Données Astronomiques). Sen piirissä toimii tähtitieteilijöitä, IT-asiantuntijoita ja informaatikkoja, jotka kaikki tuovat peliin oman erityisosaamisensa.

CDS:n tietoranteiden suunnittelun lähtökohdanna on ollut, että sekä julkaisuissa että havaintodatassa puhutaan kohteista. Taivaalta löyty-

vät valopisteet tai muunlaiset, himmeämmät läiskät saavat yleensä jonkin tunnisteiden tai useampia. Strasburgissa alettiin pitää kirjaa kohteiden nimistä ja kerätä niitä tietokannaksi. Näitä tunnisteita haravoidaan tähtitieteen julkaisujen koktekteistä DJIN-nimisellä ohjelmalla, ja tarvittaessa voidaan lennosta luoda kokonaan uusia tunnisteita. Kaukaiset kohteet eivät meidän aikaskaalamme mukaan juurikaan liiku paikoiltaan, joten tietyissä tähtitaivaan koordinaateissa sijaitseva kohde erilaisine havaintodatoinen voidaan yhdistää tiettyyn tunnisteeseen.

CDS on tuottanut tähtitieteilijöiden käyttöön monia datatyökaluja: kohteet kattava SIMBAD, havaintotaulukoita keräävä Vizier ja vuorovai-

kutteisesti käytettävä Aladin-taivaankartasto. Nämä kuuluvat alan tutkijoiden perustyökaluihin. Uuden artikkelin ilmestyessä sen tekstissä esiintyvien kohteiden nimet tunnistetaan ja artikkelit linkittyvät jatkossa suoraan näihin kohteisiin. Myös suurimmat tähtitieteen kustantajat osallistuvat CDS:n työhön mm. huolehtimalla siitä että kohteiden ja taulukoiden tiedot on mahdollisimman helppo poimia teksteistä.

ADS kokoaa kaikki maailman tähtitieteen julkaisut

Atlantin toisella puolen NASAn rahoittama Astrophysics Data System (ADS) puolestaan on rakentanut viitetietokannan, joka sisältäisi kaikki maailmassa ilmestyneet tähtitieteen julkaisut. Näitä oli 2013 noin 10 miljoonaa viitettä. ADS käynnistettiin 1991. Se on pyrkinyt alusta asti rakentamaan linkkejä muihin tietokantoihin. Muut toimijat (CDS, kustantajat, preprint-tietokanta arXiv, tähtitieteen suuret havainto-ohjelmat) ovat sen yhteistyökumppaneita. Niinpä SIMBADiin haravoitujen kohdetunnisteiden avulla voidaan suoraan hakea kirjallisuusviitteitä ADSista. Tämä hienosti hiottu kokonaisuus on sitonut tähtitieteen tutkijat niin tiukasti verkkoonsa, ettei tähtitieteilijä yleensä tietoa hakiesaan harhaannu muualle ADSin ääreltä.

ADSiin alettiin lisätä datalinkkejä vuodesta 1997 lähtien. Tämä tarkoittaa sitä, että kyseisissä artikkeleissa on ollut linkki verkossa löytyvään dataan. Aluksi kyseessä oli vain muutamia kymmeniä linkkejä, mutta muutamassa vuodessa datalinkkien määrä alkoi lisääntyä. Nykyään niiden määräksi on tasaantunut n. 1500 linkiksi vuosittain (Pepe et al. 2013).

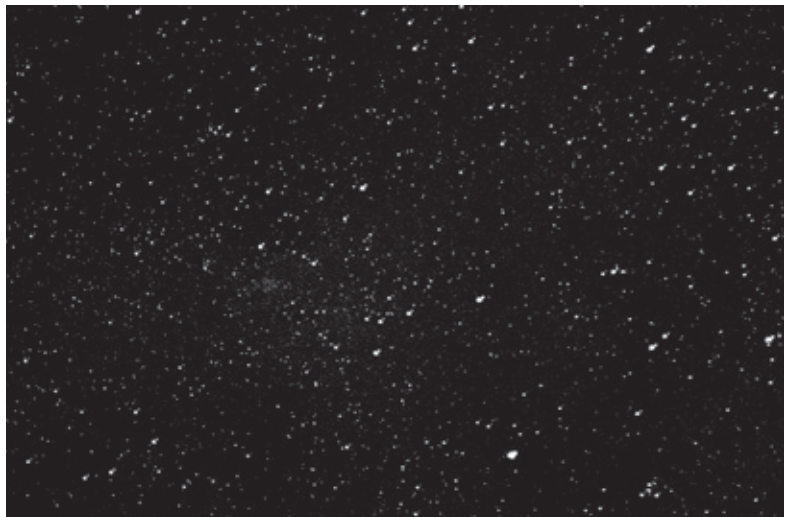
Käytännössä kyse on kahdenlaisista linkeistä. Data-

linkki voi sijaita keskitetysti ylläpidetyllä palvelimella, josta sen voi olettaa löytyvän jatkossakin. Näin on varsinkin silloin, kun kyse on käsittelemättömästä datasta ja koko taivaan kattavista kartoituksista. Data voi myöskin olla jo kertaalleen käsiteltyä. Tällöin se päättyy herkästi yksittäisten tutkijoiden kotisivuille tai usein lyhytikäisiin ftp-arkistoihin. ADSIn vanhimmista datalinkeistä jopa yli puolet oli jossain vaiheessa ehtinyt muuttua kuolleiksi linkeiksi. Suurin osa näistä oli nimenomaan yksittäisten tutkijoiden ylläpitämiä linkkejä. Sen sijaan keskitetyissä palveluissa säilytetyn datan jatkosaatavuus oli melko vakaata vuosienkin jälkeen.

Datalinkkejä metsästämissä

Alkuperäisen, käsittelemättömän datan käyttäminen ei tähtitieteessä aina ole mahdollista tai käytännöllistä. Sitä voi olla yksinkertaisesti niin paljon, ettei sitä kannata säilyttää sellaisenaan saatavilla. Monelle riittää jo kertaalleen käsitelty data. Erään tähtitieteilijän sanoin: ”Moniko on koskaan käsitellyt SDSS-kuvatiedostoja uudestaan? Antakaa kun arvaan: maa kantaa pinnallaan enintään kymmenen ihmistä jotka ovat ikinä käsitelleet Sloan [Digital Sky Survey] –kuvia uudelleen.”

Eräs keskeinen syy dataviittauksiin on niiden artikkelille tuoma lisäarvo. Henneken ja Accomazzi (2011) tarkastelivat 3814 datalinkattua artik-




Kuva: Morguefile.com / seriousfun

kelia, jotka oli julkaistu 1995-2000 ja vertasivat niiden saamia viittauksia mahdollisimman hyvin näitä vastaaviin linkkaamattomiin artikkeleihin. Kymmenessä vuodessa datalinkatut artikkelit saivat noin 20% enemmän viittauksia kuin verrokkit.

Dataviittausten seuranta kiinnostaa erityisesti suuria laiterahoittajia. Lähes kaikilla suurilla tähtitieteen laitteistoilla on nykyään oma kura-toitu viitetietokanta, johon listataan ne artikkelit, joissa on käytetty näillä laitteilla kerättyä dataa. Artikkelit linkitetään muihin viitetietokantoihin (ADS tai Web of Science) joista saadaan viittauserät. Näin yksittäisten havaintolaitteistojen tuottavuutta voidaan seurata.

Myös laitteita käyttäneistä tutkijoista voidaan hakuprosesseissa saada kiinnostavia tietoja. Jos vaikkapa suomalaiselta hakijalta toivotaan näyttöä ESO:n havaintolaitteiden käyttökokemuksesta, ESO Telbib-tietokanta kertoo, paljonko ESO:n havaintolaitteilla saatua dataa tutkija on käyttänyt, ja onko se raakadataa vai ehkä muiden käsittelemää.

Seuraava askel tähtitieteen datapalvelun palapellissä on semanttisten teknologioiden tehostettu

soveltaminen yhä kasvavaan määrään kokotekstejä. Tietomassojen kasvaessa on viisasta sijoittaa datan hallintaan mahdollisimman aikaisessa vaiheessa, jotta kalliiden havaintolaitteiden tuottama tieto voidaan hyödyntää mahdollisimman hyvin. Kallista dataa ei kannata unohtaa kellareihin! 

Lähteet

E. Henneken & A. Accomazzi (2011), Linking to Data - Effect on Citation Rates in Astronomy. arXiv:1111.3618[cs.DL]

A. Pepe et al. (2013), Sharing, archiving, and citing data in astronomy. <http://authorea.com/288>

Linkkejä

ADS: <http://labs.adsabs.harvard.edu/adsabs/>

CDS: <http://cdsweb.u-strasbg.fr/>

World Wide Telescope: <http://www.worldwidetelescope.org>

ESO Telbib: <http://telbib.eso.org/>

Tietoa kirjoittajasta

*Eva Isaksson, kirjastonhoitaja
Helsingin yliopiston kirjasto
Email. Eva.isaksson@helsinki.fi*