

Kohti yhteisiä aineistokäytänteitä

Tutkimusaineistojen yhteiskäytöstä on digitaalisena aikana ja viimeistään 2000-luvulla tullut itsestään selvä tavoite kielentutkimuksen parissa, ja aineistoasiat ovat olleet yhä enemmän esillä niin Suomessa kuin kansainvälisestikin. *Virittäjässä* 3/2010 aineistoasioita käsiteltiin monipuolisesti Suomen kielen nauhoitearkiston 50-vuotisjuhlavuoden tapahtumien ja teemojen kautta. Lehden Suunvuoro-palstalla Toni Suutari (2010a) nosti esiin ajankohtaisen ja tärkeän kysymyksen kansallisen kieliaineistoinfrastruktuurin tarpeesta Suomessa ja kiinnitti samalla huomiota muun muassa tutkimusrahoittajien kiinnostukseen aineistojen jatkokäyttöä kohtaan. Erkki Lyytikäisen ja Jaakko Yli-Paavolan (2010) katsaus puolestaan valotti sitä, miten kielenäytteiden yhtenäisestä tallentamisesta on keskusteltu ja miten sitä on kehitetty nauhoitearkiston 50-vuotisen historian aikana. Tässä kirjoituksessa jatkan keskustelua kieliaineistojen arkistoinnista ja kiinnitän huomiota erityisesti kielentutkijoiden rooliin aineistokäytänteiden kehittäjinä.

Suomalaisen arkistoinfrastruktuurin on tulevaisuudessa toimittava yhä enemmän kansainvälisessä kontekstissa ja tuettava useiden eri kielten tutkimusta. Kansallisesti merkittävän aineistorakenteen profiili painottuu kuitenkin suomalais-ugrilaisiin ja Suomessa puhuttaviin kieliin ja niiden tutkimukseen. Kes-

tävän aineistoinfrastruktuurin rakentaminen edellyttää riittävien resurssien ja teknisen toteutuksen lisäksi aineistohallinnan käytänteiden suunnittelua ja lingvistiyhteisön, erityisesti fennistien ja fennougristien, vahvaa panosta. Uusien aineistojen tuottamisen, käsittelyn ja hallinnan yhteisiä käytänteitä ja suuntaviivoja tarvitaan, jotta yhteiskäyttöön soveltuvien aineistojen tuottaminen olisi tutkijoille vaivatonta ja luonteva osa tutkimusprosessia.

Aineistojen jakaminen säästää tutkimusyhteisön aikaa ja vaivaa, ja rinnakkain samanlaisten pienten aineistojen parissa tehtävän työn sijaan on mielekkäämpää keskittyä kartuttamaan eri aineistoista ja niiden analyyseistä muodostuvaa, edustavaa yhteiskorpusta. Suomessa etenkin pienempien suomalais-ugrilaisien kielten tutkimuksessa samojen klassikkoaineistojen hyödyntäminen aina uusissa tutkimuksissa on tuttua: esimerkiksi Suomalais-Ugrilaisen Seuran kentälle lähettämien stipendiaattien keräämät laadukkaat ja verrattain suuret aineistot ja niiden pohjalta 1900-luvun mittaan toimitetut kieliopit, tekstikoelmat ja sanakirjat ovat olleet monilta osin suomalaisen fennougristiikan perusta (ks. lähemmin esim. Saarinen 2007; Grünthal 2010). Näitä klassikkoaineistoja käytetään edelleen pääaineistoina yksittäisten suomalais-ugrilaisien kielten tutkimuksessa, ja erityisesti kielten muu-

tosta heijastavina vertailuaineistoina niiden rooli on merkittävä myös tulevaisuudessa (Saarinen 2007). Vaikka 1800-luvun lopulta alkaen kerättyjä aineistoja on julkaistu läpi 1900-luvun ja digitoitu yhä enemmän, valitettavan pieni osa erityisesti aineistoihin eri tutkimusten yhteydessä tehdyistä kieliopillisista analyseistä ja käänöksistä on muiden tutkijoiden ja kieliyhteisöjen saatavilla sähköisessä muodossa. Sukukielten lisäksi sama puute koskee monelta osin myös suomen kielen aineistoja.

Entistä kattavampien ja monipuolisemmin analysoitujen yhteisten aineistojen hyödyntäminen mahdollistaa lisäksi empiirisen tutkimuksen toistettavuuden. Tältä osin kielitiede lähestyy luonnontieteitä laadunvarmistuksen ja vertaisarvioinnin näkökulmasta. Vähemmän tutkittujen, pienten tai uhanalaisten kielten materiaalien pitkäjänteisesti suunniteltu arkistointi ja jakaminen voidaan nähdä myös työnä kielidiversiteetin tallentamiseksi ja esiin tuomiseksi, tutkimuksen lisäksi myös kielten revitalisaatio- ja opeutushankkeiden käyttöön. Kieliteknologian tarjoamien sovellusmahdollisuuksien rajana ovat vain mielikuvitus ja määrärahat sekä aineistojen käyttöoikeudet.

Suomalainen arkistoinfrastruktuuri

Kieliaineistojen hallintaan ja tallentamiseen on eri maissa, tutkimuslaitoksissa ja yliopistoissa tarjolla monenlaisia työkaluja ja sähköisiä arkistoja. Sovellusten ja palveluiden runsaan määrän vuoksi haasteena onkin kieliaineistojen sekä eri arkistojen formaattien ja rakenteiden yhteensopivuus aineistojen siirrettävyyden takaamiseksi, jotta tietyllä sovelluksella tallennettu ja käsitelty aineistokokonaisuus olisi käytettävissä ja analysoitavissa

myös muualla. Yhtä suuri haaste on taata tarjolla olevien palveluiden lähestyttävyyden ja käytettävyyden kieliaineistojen tallentajille ja tutkijoille, joista suurimmalla osalla ei ole kieliteknologista erityisosaamista.

Viime vuosina yhteiseurooppalainen Clarin-hanke (ks. <http://www.clarin.eu>) on rakentanut pohjaa kieliaineistojen yhteiskäytölle Euroopassa, ja sen kansallisen yhteistyötahona toimii Suomessa Fin-Clarin-konsortio. Yhteiseurooppalaisen hankkeen päämääränä on parantaa aineistojen ja niiden käyttöön ja käsittelyyn tarkoitettujen työkalujen saatavuutta yli instituutti- ja maarajojen esimerkiksi selkeyttämällä ja yhtenäistämällä aineistojen luettelointia ja käyttöluparatkaisuja. Suomalaisen konsortion ja siihen liittyvien hankkeiden tärkeimpänä tavoitteena on kehittää Suomeen Clarin-yhteensopiva tekninen infrastruktuuri ja kartuttaa sen puitteissa tarjottavia aineistoresursseja. Useiden suomalaisten yliopistojen lisäksi konsortioon kuuluvat Tieteen tietotekniikan keskus (CSC) ja Kotimaisten kielten tutkimuskeskus (Kotus) ovat avainasemassa teknisen infrastruktuurin tarjoajina. Niiden palveluiden piiriin on konsortion päämäärän mukaan tarkoitus koota kaikki olemassa olevat ja karttavat sähköisesti käytettävät kieliaineistot ja -työkalut erilaisista yksittäisistä arkistoista.

Kotus on viime vuosina panostanut aineistojensa käytettävyyteen esimerkiksi vuonna 2006 avatulla Kaino-aineistopalvelullaan, jonka www-pohjaisen käyttöliittymän kautta aineistot on pystytty tarjoamaan helposti ja nopeasti entistä suuremmalle käyttäjäkunnalle (Suutari 2010b). Samansuuntaista työtä tehdään myös CSC:n palvelimella sijaitseva Kielipankkia kehittävässä Helsingin yli-

opiston Fin-Clarín- ja Meta-Nord-hankkeissa, joiden yhtenä tavoitteena on saattaa nykyistä laajempi valikoima Kielipankin aineistoja myös Unix-käyttöön totuttomien tutkijoiden ulottuville. Nykyisin monet suomen ja muiden uralilaisten kielten erikokoisia korpuksia sisältävät kokoelmat ovat käyttöluvanhaltijoidensa hyödynnettävissä ainoastaan komentorivin käyttöosaamista vaativalla Kielipankin Unix-palvelimella, mutta tulevaisuudessa niistä voitaneen tehdä hakuja erilaisiin tarpeisiin verkon kautta toimivalla, käyttäjänsä eteenpäin ohjeistavalla korpushakuohjelmalla.

Aineistojen monipuolisemman käytettävyyden parantamisen ohella hanke työskentelee Kielipankin aineistojen käyttöluokitusten selkiyttämiseksi ja yhdenmukaistamiseksi. Aineistot voidaan ehdotetun yhteiseurooppalaisen standardin mukaisesti luokitella kokonaan julkiseksi tai tutkimuskäyttöön tarkoitetuiksi, tai niiden käyttöoikeus voidaan rajoittaa luvanvaraiseksi (Oksanen, Lindén & Westerlund 2010). Tämän lisäksi CSC ottaa parhaillaan käyttöön uutta sähköistä käyttöluopajärjestelmää Kielipankille. Aineistojen helppo käytettävyys, selkeä aineisto- ja käyttöoikeusluokitus ja käyttöluopien sujuva hallinta kuuluvatkin pitkäjänteisen ja turvallisen säilytyksen ohella tärkeimpiin arkistoilta vaadittaviin palveluihin aineistojen yhteiskäytön sujuvoittamiseksi.

Aineiston tuottajan velvollisuudet ja hyödyt

Aineistojen yhteiskäyttöön ja jakamiseen liittyy kiinteästi sekä tekijänoikeuden että yksityisyydensuojan näkökulma. Kieliaineiston tekijänoikeus on pääsääntöisesti sen tuottajalla tai kerääjällä, ellei oikeutta

ole luovutettu erikseen aineiston julkaisijalle tai muulle taholle. Aineiston osien, analyysien ja käännösten tekijänoikeus saattaa määrittä erikseen kunkin osion laatijalle, mikäli aineiston käsittelyyn osallistuu useampi henkilö. Tekijänoikeus tuo tutkijalle mukanaan velvollisuuden aineiston takana olevia tahoja, haastateltuja tai muita oikeudenomistajia sekä aineistoa myöhemmin käytettäviä tutkimus- ja muita yhteisöjä kohtaan. Aineiston tuottamiseen osallistuneiden yksityisyydensuojasta ja tekijänoikeuden kunnioittamisesta huolehtimisen lisäksi hyvin tutkimuseettisiin periaatteisiin kuuluu aineiston myöhemmän käytön määrittäminen mahdollisia eri käyttötarkeitua ajatellen. Tutkijan apuna myöhemmän käytön huomioon ottamisessa toimivat arkistojen tarjoamat valmiit, kansainvälisesti yhtenäiset käyttöoikeusluokitukset. Aineiston ja sen osien oikeudenomistajien vastuulle jää silti huolehtiminen myöhemmän käytön ja julkisuusasteen sopimisesta yhdessä informanttien ja muiden aineiston tuottamiseen osallistuneiden tahojen kanssa.

Aineiston tuottajana tutkijan tehtäväksi jää lisäksi huolehtia aineiston ja sen metatietojen työstämisestä arkiston ohjeistuksen mukaiseen muotoon siirrettävyyden, luokittelun ja käytettävyyden varmistamiseksi. Eri työvaiheissa käsiteltävän aineiston tulee taipua eri formaatteihin, ja esimerkiksi puheaineiston tie nauhurilta litteroiduksi ja mahdollisesti myös tarkempia analyysitasoja sisältäväksi aineistoksi kulkee monen välineen ja sovelluksen kautta. Valmis aineisto tulee samalla tavoin voida tarjota käyttäjille erilaisia tutkimustarpeita ajatellen ja mahdollisesti useamman eri ohjelman kautta käytettäväksi. Fin-Clarín-hanke on antanut Kielipankkiin tallennettavien

uusien aineistojen tuottajille ohjeeksi aineistojen toimittamisen xml-muodossa. Äänitallenteiden olisi hyvä lisäksi olla litteroituja ja aikakoodattuja siten, että litteraatio on automaattisesti yhdistetty äänitallenteen vastaavaan kohtaan, mikä helpottaa huomattavasti hakujen tekemistä aineistosta analysointivaiheessa. Puheaineiston litteroinnin lisäksi ei vaadita muita nimikointitasoja, mutta ihanneta-pauksessa toki aineistoa eri tutkimustar-koituksiin käyttäneet tutkijat oheistavat aineistoon omat analyysitasonsa ja jaka-vat ne edelleen uuteen tutkimuskäyttöön. Sen sijaan tiedot aineiston sisällöstä, me-tatiedot, laaditaan aineiston ohien kan-sainvälisten standardien ja arkiston oh-jeiden mukaisesti. (Ks. esim. CLARIN Metadata Now 2009.)

Myöhemmän tutkimuskäytön mah-dollistaminen on aineiston tuottajalle prosessi, johon aineiston varsinaisen ker-äämisen lisäksi kuluu aikaa ja vaivaa. Aineiston tuottaminen kestävään, myös muuta kuin omaa käsillä oleva tutki-musta palvelemaan muotoon, vaatii eri-tyistä työpanosta, jonka tulee näkyä myös tutkijan tieteellisenä ansiona. Asianmu-kaisesti toimitettu raaka-aineisto voi pal-vella myöhempää tutkimusta ja tieteen-alaa omalta osaltaan yhtäläisesti kuin sen keräämisen motivaationa oleva ja tutki-jan omana analyysinä syntyvä julkaisu. Myöhempää tutkimusta palveleva tie-deyhteisön käytettäväksi toimitettu ja jul-kaistu aineisto tulee näistä syistä nähdä tutkijalle ansioksi luettavana erillisenä tieteellisenä tuotoksena esitelmien, artik-keleiden, ohjelmistojen ja patenttien jou-kossa. Julkaisuksi luettavalle aineistolle asetettavat vaatimukset määrittää lopulta tiedeyhteisö, joka määrittelemällä aineis-ton tieteelliset laatuvaatimukset nostaa samalla aineiston tuottajan työn arvoa ja

tekee sen näkyväksi. Raaka-aineiston tai sen osan tuottajan kannalta sen jakami-nen ja käyttöoikeuksin rajattu tai vapaa julkaiseminen merkitsee yhä uusia viit-tauksia aineistoon, joka pystyy todennä-köisesti palvelemaan hyvinkin erityyppiä tutkimuksia tuottajansa oman tut-kimusintressin lisäksi. Aineiston arvos-taminen tieteellisenä julkaisuna motivoi sen tutkijaa käsittelemään aineistonsa myös kansainvälisen käyttäjäkunnan huomioon ottaen.

Tavoitteena aineistohallinnan yhteiset käytänteet

Aineistoa keräävän tutkijan kannalta suuri kysymys on, miten aineisto työste-tään julkaisuksi arkistointia ja myöhem-pää käyttöä varten joustavasti ja omaa tutkimustyötä samalla tukien. Tutki-jan on esimerkiksi ratkaistava, miten ai-neisto työstetään tiettyyn formaattiin, millä sovelluksilla litteroidaan, annotoi-daan, kirjoitetaan käännökset ja käsitel-lään sanastoa, mitä aineiston käytöstä sovitaan informanttien kanssa ja miten voi selvittää eri lähteistä poimitun ai-neistokokonaisuuden tekijänoikeudet. Suurin osa kysymyksistä vaatii selvit-tämistä pelkästään aineiston kerääjän omia tutkimuskysymyksiä ja -julkaisuja ajatellen, mutta niiden ratkaisutapa riip-puu samalla aineiston myöhemmästä ar-kistointi- ja käyttötavasta. Aineiston kä-sittelyyn ja tuottamiseen liittyviin kysy-myksiin tarvitaankin yhteistä keskuste-lua ja aineistohallinnan hyvien käytän-teiden jakamista. Kokemuksen ja tutki-mustapojen ja -tekniikoiden kehittymi-sen myötä päivitettävät aineiston työstä-mistä koskevat ohjeet ja suositukset tu-kevat yksittäisiä tutkijoita aineiston hal-linnassa, ja ne kuuluvat teknisen raken-

teen lisäksi erottamattomasti osaksi toimivaa aineistoinfrastruktuuria.

Arkistot ja aineistoihin keskittyvät verkostot kuten Clarin voivat tarjota aineiston työstämiseen erilaisia suunta- viivoja kuten metatietostandardeja ja käyttöoikeusluokituksia. Suomessa Fin-Clarin-konsortion perustaminen ja sen kautta tehtävä työ eri aineistokokonai- suuksien kokoamiseksi yhteen on hyvä alku. Konkreettiset aineiston työstämi- seen liittyvät toimet ja käytänteet on kui- tenkin valittava laajassa lingvistiyhtei- sössä. Aineistojen tärkeimmät tuotta- jat, tutkijat, ja aineistojen käyttäjinä tut- kijoiden ohella myös opiskelijat, tietä- vät parhaiten, millaisia aineistosisältöjä ja -muotoja tarvitaan ja voidaan käyttää. Myös aineiston käsittelyyn liittyvät käy- tännön kysymykset, kuten millainen ai- neisto kannattaa litteroida Elan- tai Tran- sana-ohjelmalla tai milloin informanttien kanssa tulisi käyttää kirjallista sopimus- pohjaa, jäävät tutkijoiden ratkaistavaksi. Samankaltaisia asioita ei kannata tehdä eri yliopistoissa ja tutkimusryhmissä eri tavoin, mikäli siihen ei ole tutkimuksen tavoitteiden määrittelemää sisällöllistä tarvetta, vaan kokemuksia on syytä ke- rätä ja jakaa tutkimusyhteisössä entistä aktiivisemmin. Toisaalta kieliaineistoja tallentava arkistotahokaan ei saa jäädä pelkäksi aineistoja vastaanottavaksi tek- niseksi rakenteeksi, vaan sen täytyy pys- tyä jatkuvaan vuoropuheluun aineistoa tuottavan ja käyttävän lingvistiyhteisön kanssa. Tutkijat ovat parhaita asiantun- tijoita kertomaan aineistoihin liittyvistä käytännön haasteista ja tarpeista; arkisto puolestaan pystyy luomaan ja kehittä- mään erilaisten aineistojen hallintaan ja käyttöön sovellettavia kansainvälisen inf- rastruktuurin mukaisia teknisiä ja sisäl- löllisiä ohjeita ja sovelluksia.

Aineistoja on yliopistoissa ja tutki- muslaitoksissa kartutettu pitkäjänteisesti erilaisissa tutkimusprojekteissa ja myös opiskelijavoimin (puheaineistojen keruu- hankkeista ks. Lyytikäinen & Yli-Paavola 2010; Siirainen 2010; Karttunen & Rou- hikoski 2010). Erilaisten hankkeiden yh- teistyön ja suunnitelmallisuuden lisäämi- nen eri tutkimusten tarpeita palvelevien aineistojen keruussa ja työstämisessä ko- rostuukin varmasti tulevaisuudessa. Ai- neistonhallinnan kysymykset liittynevät tulevaisuudessa yhä kiinteämmin myös opetukseen, ja aineistoja ja niiden ana- lyysejä voi olla järkevää kartuttaa sys- temaattisemmin myös kurssitöinä. Esi- merkkinä tästä on Helsingin yliopiston suomen kielen ja suomalais-ugrialaisten kielten oppiaineryhmien suunnitelma kytkeä aineistojen käsittelyyn tarkoitettu- jen sovellusten käyttö ja olemassa olevien aineistojen analyysitasojen kartuttami- nen yhä tiiviimmin osaksi opinnäytese- minaareja, pääaineeseen kytkettyjä tieto- ja viestintäteknikan opintoja, sukukiel- ten kursseja ja kenttätöyöpetusta.

Aineistonhallinnan käytänteiden luo- minen, kehittäminen ja vakiinnuttami- nen, erilaiset aineistotalkoot ja johdon- mukainen aineistonhallinta säästäne- vät tulevilta lingvistikupolvilta paljon työtä, jota muuten tehtäisiin nykyaineis- tojen formaattien muokkauksen ja niiden käyttöoikeuksien selvittämisen parissa – elleivät nykyaineistot onohdu ja jää hyö- dyntämättä uudelleen. Ylimääräistä vai- vaa on syytä välttää myös aineiston käsit- telemisessä arkistojulkaisuksi, eikä kyn- nys tuottaa myös myöhempään käyttöön sopivaa aineistoa saa nousta liian kor- keaksi tutkijalle. Kysymystä siitä, miten mahdollisimman vähällä vaivalla voitai- siin tuottaa mahdollisimman monikäyt- töisiä aineistoja, on pidettävä jatkuvasti

esillä aineistoja käyttävän ja tuottavan lingvistiyyhteisön ja tutkimushakkeiden piirissä, kieliarkistojen ja niiden kehittäjien työssä sekä näiden yhteisillä olemassa olevilla ja aktiivisesti rakennettavilla uusilla foorumeilla.

LOTTA JALAVA
etunimi.sukunimi@helsinki.fi

Lähteet

- CLARIN Metadata Now. Short Guide, March 2009. http://www.clarin.eu/system/files/Metadata_now-CLARIN-ShortGuide.pdf. (30.12.2010.)
- GRÜNTAL, RIHO 2010: Matkueita ja yksittäisiä tutkijoita. Suomalais-Ugrilaisen Seuran keruuretkien tausta ja tavoitteet. – *Uralica Helsingiensia* 4 s. 17–51.
- KARTTUNEN, MIIA – ROUHIKOSKI, ANU 2010: Kentällä kokeillen. Murresyntaxin tutkijat uusia aineistoja kokoamassa. – *Virittäjä* 114 s. 426–432.
- LYYTIKÄINEN, ERKKI – YLI-PAAVO-
LA, JAAKKO 2010: Suomen kielen nauhoitearkisto 50-vuotias. – *Virittäjä* 114 s. 411–419.
- OKSANEN, VILLE – LINDÉN, KRISTER – WESTERLUND, HANNA 2010: Laundry symbols and license management. Practical considerations for the distribution of LRs based on experiences from CLARIN. – *Proceedings of LREC 2010: Workshop on Language Resources: From Storyboard to Sustainability and LR Lifecycle Management. May 2010, Malta*. <http://www.lrec-conf.org/proceedings/lrec2010/workshops/W20.pdf>.
- SAARINEN, SIRKKA 2007: Fennougristinen kenttätyö. <http://www.kotus.fi/index.phtml?s=734>. (30.12.2010.)
- SIIROINEN, MARI 2010: Aikamatkoja ja nykypäivää äänitallenteiden maailmassa. Suomen kielen nauhoitearkiston 50-vuotisjuhlavuoden satoa. – *Virittäjä* 114 s. 420–422.
- SUUTARI, TONI 2010a: Suunvuoro. – *Virittäjä* 114 s. 323.
- 2010b: Suomen kielen nauhoitearkisto – vireä viisikymppinen. – *Virittäjä* 114 s. 423–426.