

Julkaisuarkistot – pysyvän saatavuuden infrastruktuuri?

Jyrki Ilva (jyrki.ilva@helsinki.fi)

Missä mennään -webinaari, 16.11.2021

Pysyvä osoite: <https://urn.fi/URN:NBN:fi-fe2021111655567>

Verkkoaineistojen pysyvyys usein häilyvää

- Verkkoaineistojen pysyvyydessä monenlaisia haasteita
- Aineistoja katoaa saatavilta tai niiden sisältö saattaa muuttua
- Ongelma koskettaa myös tieteellisiä julkaisuja ja muita julkaisuarkistoihin tallennettavia aineistoja



Jonathan Zittrain: [The Internet Is Rotting](#). Atlantic, June 30, 2021.

Webinaari julkaisuarkistoista ja pitkäaikaissäilytyksestä

- Kansalliskirjasto järjestää [webinaarin julkaisuarkistoaineistojen pitkäaikaissäilytyksestä](#) 9.12. klo 9-11
 - Ilmoittautuminen auki 7.12. asti
- Tässä esityksessä keskitytään ensisijaisesti aineistojen pysyvään saatavuuteen

Webinaari: Julkaisuarkistojen aineistot pitkäaikaissäilytykseen?

Created by Jyrki Ilva, last modified on Nov 12, 2021

Aika: **Torstai 9.12.2021 klo 9-11**

Paikka: Zoom-etäkokous, [osallistumisohjeet](#)

Ilmoittautuminen: [viimeistään tiistaina 7.12.](#)

Julkaisuarkistoihin tallennetaan vuosittain kymmeniä tuhansia julkaisuja, sekä opinnäytteitä, rinnakkaistallenteita, sarjajulkaisuja että muitakin aineistoja. Kansalliskirjasto kerää julkaisuarkistoissa julkaistuja aineistoja talteen osana kulttuuriaineistolain mukaista toimintaansa. CSC ylläpitää kansallista pitkäaikaissäilytyspalvelua. Millä tavalla ja missä laajuudessa julkaisuarkistojen aineistot päätyvät pitkäaikaissäilytykseen?

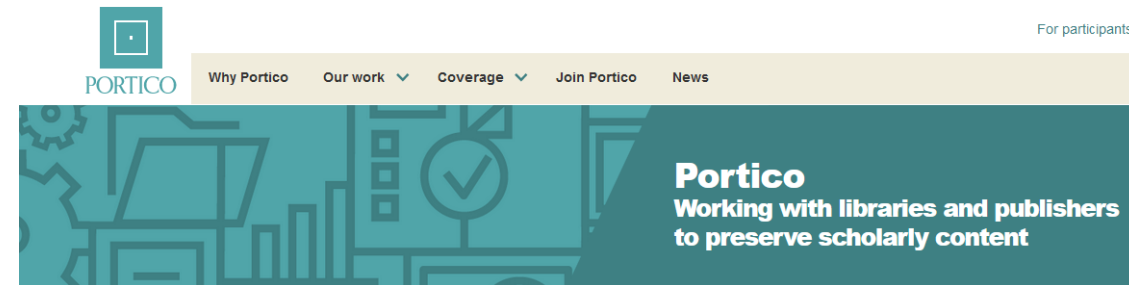
Kansalliskirjaston järjestämässä seminaarissa keskustellaan julkaisuarkistoaineistojen pitkäaikaissäilytykseen (PAS) liittyvistä kysymyksistä. Tilaisuus on suunnattu sekä julkaisuarkistojen parissa työskenteleville että niiden toiminnasta vastaaville henkilöille.

Tilaisuus on ilmainen, ja siihen voi **ilmoittautua viimeistään tiistaina 7.12. tällä verkkolomakkeella**. Etäkokouksen osoite lähetetään ilmoittautuneille tilaisuutta edeltävänä päivänä.

Tervetuloa mukaan!

Pysyvä saatavuus ja pitkäaikaissäilytys

- Pysyvä saatavuus ja pitkäaikaissäilytys kytköksissä toisiinsa, mutta kuitenkin eri asioita
 - Pitkäaikaissäilytys ei sinällään tarkoita sitä, että säilytetyt aineistot olisivat laajasti saatavilla
- Pysyvä saatavuus ei sekään ole aivan uusi idea
 - Käytännön esimerkkinä tieteellisten kustantajien käyttämä [Portico-palvelu](#)
 - Pysyvyyden pituus ja saatavuuden laajuus enemmän tai vähemmän suhteellisia asioita



Portico is a community-supported preservation archive that safeguards access to e-journals, e-books, and digital collections. Our unique, trusted process ensures that the content we preserve will remain accessible and usable for researchers, scholars, and students in the future.

Resources for our community and business continuity

[Read ITHAKA's COVID-19 response](#)

NEWS

[All News](#)

September 30, 2021

Access alert: former De Gruyter journals

August 13, 2021

Portico to participate in Embedding Preservability project

August 5, 2021

Access alert: OA journals from De Gruyter, Dedicated Juncture Researcher's Association, and Manchester University Press

Search for titles in the Portico archive

Search titles by keyword:

E-journals E-books Collections

Search now

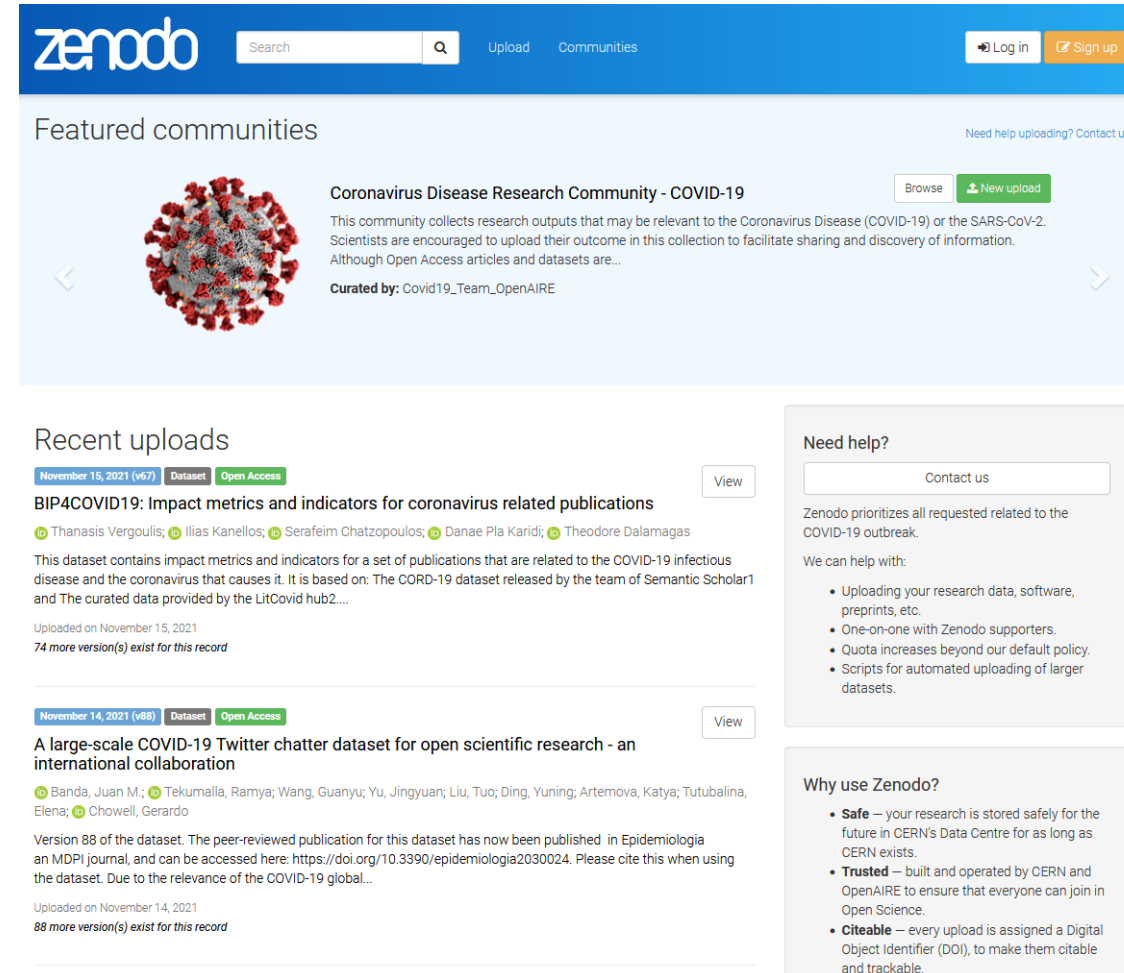
Request a free holdings comparison

- Inform your decision about joining Portico
- Compare costs to other preservation solutions
- Manage the preservation status of your collection

Get started

Pysyvä saatavuus ja julkaisuarkistot

- Maailmalla julkaisuarkistoilla, pitkäaikaissäilytyksellä ja pysyvällä saatavuudella on usein läheinen yhteys
 - PAS-järjestelmiä saatetaan rakentaa myös organisaatiotasolla
- CERN:in ja OpenAIRE:n [Zenodo-julkaisuarkistoon](#) viitataan usein pitkäaikaissäilytyksen ja pysyvän saatavuuden takaavana ratkaisuna
 - "Tutkimuksesi on tallessa niin pitkään kuin CERN on olemassa"



The screenshot shows the Zenodo website interface. At the top, there is a blue navigation bar with the Zenodo logo, a search bar, and links for 'Upload' and 'Communities'. On the right side of the navigation bar, there are 'Log in' and 'Sign up' buttons. Below the navigation bar, the main content area is divided into two sections: 'Featured communities' and 'Recent uploads'.

Featured communities: The first featured community is the 'Coronavirus Disease Research Community - COVID-19'. It features a 3D model of a coronavirus particle. The description states: 'This community collects research outputs that may be relevant to the Coronavirus Disease (COVID-19) or the SARS-CoV-2. Scientists are encouraged to upload their outcome in this collection to facilitate sharing and discovery of information. Although Open Access articles and datasets are...'. It is curated by 'Covid19_Team_OpenAIRE'. There are 'Browse' and 'New upload' buttons.

Recent uploads: Two recent uploads are listed. The first is 'BIP4COVID19: Impact metrics and indicators for coronavirus related publications', uploaded on November 15, 2021. It is a dataset and is open access. The authors listed are Thanasis Vergoulis, Ilias Kanellos, Serafeim Chatzopoulos, Danae Pla Karidi, and Theodore Dalamagas. The description mentions that the dataset contains impact metrics and indicators for a set of publications related to COVID-19. The second upload is 'A large-scale COVID-19 Twitter chatter dataset for open scientific research - an international collaboration', uploaded on November 14, 2021. It is also a dataset and open access. The authors listed are Banda, Juan M.; Tekumalla, Ramya; Wang, Guanyu; Yu, Jingyuan; Liu, Tuo; Ding, Yuning; Artemova, Katya; Tutubalina, Elena; and Chowell, Gerardo. The description mentions that version 88 of the dataset has been published in the journal Epidemiologia and MDPI, and provides a DOI link.

Need help?: A section titled 'Need help?' with a 'Contact us' button. Below it, it states: 'Zenodo prioritizes all requested related to the COVID-19 outbreak. We can help with:' followed by a list of services: 'Uploading your research data, software, preprints, etc.', 'One-on-one with Zenodo supporters.', 'Quota increases beyond our default policy.', and 'Scripts for automated uploading of larger datasets.'

Why use Zenodo?: A section titled 'Why use Zenodo?' with a list of benefits: 'Safe' (research is stored safely in CERN's Data Centre), 'Trusted' (built and operated by CERN and OpenAIRE), and 'Citeable' (every upload is assigned a Digital Object Identifier (DOI)).

Entä Suomessa?

- OKM:n rahoittamat ja CSC:n rakentamat kansalliset pitkäaikaissäilytyspalvelut, Kulttuuriperintö-PAS ja Fairdata-PAS
- Pysyvästä saatavuudesta on puhuttu aika ajoin kunnianhimoisesti pitkäaikaissäilytyksen päälle rakennettavana lisäpalveluna
 - Joskus tulevaisuudessa, kenties – mutta miten eletään siihen asti?
- Missä määrin julkaisuarkistot vastaavat – tai voisivat vastata – aineistojen pysyvää saatavuutta koskeviin tarpeisiin?



KANSALLISET PITKÄAIKAISÄILYTYSPALVELUT

PAS-palveluilla tarkoitetaan kulttuuriperintöaineistojen ja tutkimusaineistojen pitkäaikaissäilyttämiseen tuotettuja palveluita yhdessä.

Pitkäaikaissäilytys tarkoittaa digitaalisen informaation säilyttämistä ymmärrettävänä ja käytettävänä useiden kymmenien ja jopa satojen vuosien ajan. Laitteet, ohjelmistot ja tiedostomuodot vanhenevat ajan myötä, mutta informaation täytyy säilyä. Luotettava pitkäaikaissäilyttäminen edellyttää sisällön eheyden aktiivista valvontaa ja monenlaisiin riskeihin varautumista. Tässä ovat keskeisessä asemassa metatiedot, jotka kuvailevat mm. aineiston sisältöä, historian ja aikuperän sekä tiedot siitä, miten informaatiota voidaan käyttää.

Kulttuuriperintöaineistojen pitkäaikaissäilyttämiseen tuotettu PAS-palvelu (Kulttuuriperintö-PAS-palvelu) takaa kirjastojen, arkistojen ja museoiden keskeisten kansallisten digitaalisten tietovarantojen pitkäaikaissäilyttämisen. Digitaalisilla kulttuuri-perintöaineistoilla tarkoitetaan sekä digitoituja että digitaaliseen muotoon tuotettuja kulttuuriperintöaineistoja: lakisääteisen säilyttämisen piiriin kuuluvia digitaalisia kulttuuriaineistoja, kansalliseen kulttuuriperintöön kuuluvaa digitaalista asiakirjallista aineistoa sekä aineellisen ja henkisen kulttuuriperinnön säilyttämisestä vastaavien, opetus- ja kulttuuriministeriön hallinnonalalla toimivien organisaatioiden muita pitkäaikaissäilytyksen piiriin kuuluvia digitaalisia tietovarantoja.

Tutkimusaineistojen pitkäaikaissäilyttämiseen tuotettu PAS-palvelu (Fairdata PAS-palvelu) varmistaa tutkimuksen digitaalisten aineistojen saatavuuden ja pitkäaikaisen säilyvyyden. Tämä PAS-palvelu tukee osaltaan pysyvää ja koordinoitua toimintamallia tutkimusaineistojen hallinnan tueksi. Pyrkimyksenä on, että tutkimuksen todennettavuus ja toistettavuus elinkaaren eri vaiheissa onnistuu ja tulosten hyödyntäminen on helppoa. Tällöin tutkimustuloksia voidaan käyttää yhä uudelleen, arvioida, hyödyntää päätöksenteossa ja turvata digitalisoitumisen myötä yhä nopeammin kasvavat tietomäärät tulevien tutkijakupolvien käyttöön.

Organisaatiot ja vastuu

- Pysyvä saatavuus ei toteudu itsestään, vaan vaatii työtä ja riittäviä resursointia
- Edellyttää käytännössä sitä, että taustalla on jokin vakaa ja pysyvä organisaatio
 - Organisaatiolla pitää olla valmiudet ottaa vastuu aineistoista ja palveluista myös pitkällä tähtäimellä
 - Organisaation pitää olla valmis sitoutumaan siihen, että se pitää huolta aineistoista
 - Jos tämä ei ole jossain vaiheessa mahdollista, vastuu pitää välittää eteenpäin jollekin muulle taholle
 - Käytännössä fuusiot, organisaatiomuutokset, rahoitustilanteen vaihtelu ja uudet strategisen tason priorisoinnit saattavat aiheuttaa haasteita vakainakin pidetyille organisaatioille

Palvelut ja tekninen infrastruktuuri

- Tällä hetkellä kaikilla suomalaisilla korkeakouluilla ja monilla tutkimuslaitoksilla on käytössään julkaisuarkisto
 - Palvelukonseptit saattavat ajan myötä elää – esim. julkaisuarkistojen ja tutkimustietojärjestelmien välinen työnjako selkiytynyt vähitellen
 - Miten pitkäikäinen konsepti julkaisuarkisto on?
- Palveluita voidaan tuottaa joko organisaation omin voimin tai yhteistyössä
 - Suurin osa suomalaisista organisaatioista käyttää Kansalliskirjaston tarjoamaa keskitettyä palvelua
 - Pysyvän saatavuuden takaamiseksi teknisen infrastruktuurin ylläpitoon ja kehittämiseen pitää olla riittävät resurssit

Ohjelmistojen elinkaari

- Avoimen lähdekoodin julkaisuarkisto-ohjelmistot
 - Ohjelmistot vaativat ylläpitoa ja kehitystyötä
 - Koodi ei synny itsestään vaan kehitystyöhön pitää löytää resursseja tavalla tai toisella
- Suomessa käytössä DSpace
 - Julkaistiin alun perin vuonna 2002
 - Käytössä tuhansissa organisaatioissa eri puolilla maailmaa
 - Kesällä 2021 julkaistu versio 7 ajanmukaisti ohjelmiston arkkitehtuuria ja teknisiä ratkaisuja
- Kansainvälisen DSpace-yhteisön toimintaa rahoitetaan nykyään pääosin jäsenmaksuilla
 - Kansalliskirjasto on DSpace-yhteisön jäsen
 - DSpace nyt myös [SCOSS:in](#) rahoituskohteena



WHAT IS SCOSS?

The Global Sustainability Coalition for Open Science Services (SCOSS) is a network of influential organisations committed to helping secure OA and OS infrastructure well into the future. Officially formed in early 2017, SCOSS's purpose is to provide a new co-ordinated cost-sharing framework that will ultimately enable the broader OA and OS community to support the non-commercial services on which it depends. [READ MORE >>](#)

[KEY FACTS & FIGURES](#)

[GOVERNANCE](#)

[DEFINING OPEN INFRASTRUCTURE](#)

[PAST SCOSS AWARDEES](#)

[WHO IS BEHIND SCOSS](#)

Pysyvät osoitteet

- Pysyviin tunnisteisiin perustuvat verkko-osoitteet yksi keino helpottaa organisaatioiden, palveluiden ja ohjelmistoversioiden muutoksista aiheutuvia ongelmia
 - Vaikka julkaisun verkko-osoite muuttuu, se on edelleen löydettävissä pysyvän osoitteen avulla
- Suomalaisissa julkaisuarkistoissa käytetään yleisesti URN:ejä, joissakin käytössä myös Handle ja DOI
 - Järjestelmästä riippumatta pysyviä osoitteita käyttävät organisaatiot sitoutuvat ylläpitämään tietoa siitä, missä osoitteessa tunnistetta vastaava dokumentti sijaitsee
 - Tiedot toimitetaan resolverille, joka ohjaa käyttäjät oikeaan osoitteeseen
 - DOI:n etuna kansainvälisten palveluntarjoajien (Crossref, DataCite) rakentama tieteellisille aineistoille optimoitu infrastruktuuri
 - Handle- ja DOI-pohjaisten osoitteiden uudelleenohjaus saattaa olla organisaatiomuutoksissa jonkin verran mutkikkaampaa kuin URN:in

Julkaisuoikeudet

- Julkaisun pysyvän saatavuuden kannalta on eduksi, jos siihen liittyvistä oikeuksista on sovittu ja ne on määritelty selkeästi
 - Creative Commons -lisenssit yleinen keino jatkokäyttöoikeuksien ilmaisemiseen
 - Esim. CC BY -lisenssi mahdollistaa aineiston tallentamisen ja julkaisemisen muissa palveluissa, ml. julkaisuarkistot
- Muussa tapauksessa julkaisun avoin saatavuus saattaa rajoittua palveluun johon se on alun perin tallennettu
 - Aiheuttaisi ongelmia myös mahdolliselle kansalliselle pysyvän saatavuuden palvelulle
- Julkaisuarkistoissa CC-lisenssoitua aineistoa edelleen suhteellisen vähän
 - Rinnakkaistallenteiden osalta kustantajan kanssa solmitut sopimukset rajoittaneet; tallennus perustunut usein kustantajan politiikkaan, ei lisenssiin
 - Opinnäytteissä lisenssin käyttö kiinni tekijän päätöksestä

Kuvailutiedot ja löydettävyys

- Julkaisuarkistoaineistoilla usein hyvä näkyvyys Googlen ja muiden hakukoneiden tuloksissa
 - DSpace-kehittäjät tehneet yhteistyötä Google Scholarin kanssa indeksoitumisen parantamiseksi
 - Julkaisuihin kohdistuvat linkit ja viittaukset parantavat näkyvyyttä
 - Pysyvien osoitteiden käyttö linkityksissä suositeltavaa, mutta saattaa vaikuttaa hakukoneiden tuloksiin
- Myös näkyvyydellä kirjastojen hakukäyttöliittymissä merkitystä
 - Julkaisuarkistojen Dublin Core -muotoista metadataa haravoidaan jatkuvasti mm. Finnan indeksiin
 - Metadatan jatkokäyttö on joitakin poikkeuksia lukuunottamatta vapaata (CC0)
 - Julkaisuarkistojen metadataa saatetaan kirjastoissa konvertoida myös MARC-formaattiin, missä on omat haasteensa
 - Kansalliskirjaston metadataprosessi ja pitkäaikaissäilytys?

Tallennusformaatit

- Tallennusformaattien elinkaari yksi mahdollinen pullonkaula pysyvän saatavuuden toteutumisessa
 - Pysyvätkö dokumentit luettavina teknisen ympäristön muutoksista huolimatta?
 - Pitkäaikaissäilytysjärjestelmissä tehdään tarvittaessa migraatioita formaatista toiseen; julkaisuarkistoissa ei tällaiseen pääsääntöisesti mahdollisuuksia
- Julkaisuarkistojen aineistot pääosin suhteellisen yksinkertaisia PDF-formaattiin tallennettuja tekstimuotoisia dokumentteja
 - Yleisesti käytetty ja suhteellisen hyvin dokumentoitu formaatti
 - Kasvava osa aineistoista tallennettu jo valmiiksi PDF/A-formaatissa
 - XML-pohjaiset tallennusformaatit (esim. tieteellisten kustantajien käyttämä JATS-XML) eivät ole toistaiseksi yleistyneet prosesseihin ja kustannuksiin liittyvistä syistä johtuen
 - Audiovisuaalisten aineistojen kanssa saattaa olla enemmän ongelmia

Sisällölliset muutokset ja versiointi

- Käytännön elämässä aineistoja saatetaan joutua korjaamaan ja päivittämään
 - Julkaisu on saatettu tallentaa vääränä versiona tai väärässä muodossa
 - Julkaisussa saattaa olla jokin korjaamista vaativa virhe
 - Tieteellisestä artikkelista saatetaan tallentaa ensin tekijän oma versio ja sitten kustantajan versio
- Päivitysprosessien käytännöissä vielä parannettavaa
 - DSpace-julkaisuarkistoissa mahdollisuus tallentaa julkaisuille perustason provenienssitietoja, joihin rekisteröidään myös tiedostojen tarkistussummat
 - Julkaisujen versiointi parantaisi tilannetta mm. viittaamisen kannalta
 - Esim. Zenodossa tähän jo valmiimpia välineitä
 - Kysymyksiä: Miten käyttäjät ohjataan julkaisun uusimpaan versioon? Miten pysyvien tunnisteiden kanssa toimitaan, kun julkaisusta on erilaisia versioita?

Julkaisujen poistaminen saatavilta

- Yksittäisiä aineistoja saatetaan kaikesta huolimatta välillä joutua poistamaan saatavilla
 - Esim. opinnäytteen tekijä saattaa peruuttaa julkaisuluvan
 - Julkaisu saattaa sisältää vanhentunutta tietoa, josta voi olla haittaa
- Poistamiseen liittyvissä käytännöissä vielä monin paikoin kehittämistä
 - Poistetulle julkaisulle on suositeltavaa jättää julkaisuarkistoon "hautakivisivu"
 - Poistoprosessi on hyvä dokumentoida ja samalla on syytä pitää huolta, että poistettu julkaisu on arkistoitu

Julkaisuarkistot ja pysyvä saatavuus?

- Julkaisuarkistot takaavat kohtuullisen hyvin aineistojen saatavuuden lyhyellä ja keskipitkällä aikavälillä
 - Pitkän aikavälin pysyvyyden osalta ennusteiden tekeminen on vaikeampaa
- Mitä jos pysyvyys ei toteudukaan?
 - Julkaisun tallentaminen moneen eri paikkaan parantaa sen säilymismahdollisuuksia, vaikka jokin yksittäinen palvelu poistuisi käytöstä
 - Vaikka aineisto ei enää olisi saatavilla julkaisuarkistossa, se saattaa löytyä verkkoarkistoista (mm. [Internet Archive](#) ja [Suomalainen verkkoarkisto](#)) tai vapaakappaleena Kansalliskirjastosta ja muista vapaakappalekirjastoista



www.kansalliskirjasto.fi