



# Aineistot kohti pitkäaikaissäilytystä

## Digitaalisten aineistojen säilyttäminen Kulttuuriperintö-PAS-palvelussa

4.10.2022

Johan Kylander



# Digitaalisen aineiston hallinta



# Digitaalisten aineistojen jäsentäminen

- Digitaalisen aineiston hallinta lähtee siitä, että tunnistetaan mitä digitaalista aineistoa omistetaan ja missä muodossa se on
- Aineistoa pitää dokumentoida, varustaa metatiedoilla, jotka kertovat mistä on kyse
- Digitaalista aineistoa pitää osata jäsentää, esimerkiksi aineistokokonaisuus voi koostua useasta erityyppisestä tiedostosta
- Esimerkiksi oheisdokumentaatioita, joka kuvaa tai täydentää digitaalista aineistokokonaisuutta, on syytä säilyttää



## Resurssin käyttötarkoitus

fi ▼

Lataa ▼

Luonnos Rekisteri: Tutkimusaineistojen koodistot Tietoalue: Koulutus

Organisaatio: CSC - Tieteen tietotekniikan keskus

KOODIT

TIEDOT

Hae koodia



7 koodia

source - Lähdeaineisto

Luonnos

outcome - Tulosaaineisto

Luonnos

publication - Julkaisu

Luonnos

documentation - Dokumentaatio

Luonnos

configuration - Konfiguraatiotiedosto

Luonnos

method - Metodi

Luonnos

rights - Oikeuksien kuvaus

Luonnos



# Aineiston synty

- Aineiston syntyhistoria kertoo paljon siitä, miksi digitaalisella aineistolla on tiettyjä piirteitä
  - Esim. skannattu kuva on erilainen verrattuna kameralla otettuun kuvaan
  - Eri PDF-ohjelmistot tuottavat piirteitään erilaisia PDF-dokumentteja
- Syntyhistoriaa ei välttämättä voida rekonstruoida myöhemmin
- Digitaalisen aineiston laatua ja kestävyyttä pitää miettiä jo alusta saakka
  - Tiedostomuodot hallintaan jo luontivaiheessa
- Korkealaatuinen digitaalinen aineisto alkaa jo suunnittelusta

# Pitkäaikaissäilyttäminen osana digitointivaihetta

Aineiston säilyttämisen suunnittelu osaksi digitaalisen aineiston syntyä:

1. Laatu
  - Korkeampi laatu kestää pidempään
2. Tiedostomuodot
  - "Turhia" normalisointeja kannattaa välttää
3. Metatiedot
  - Aineiston synnyn ja digitointiprosessin dokumentointi
4. Käyttötarkoitus
  - Soveltamisohjeilla voidaan päästä haluttuun lopputulokseen







# Tiedostomuodot

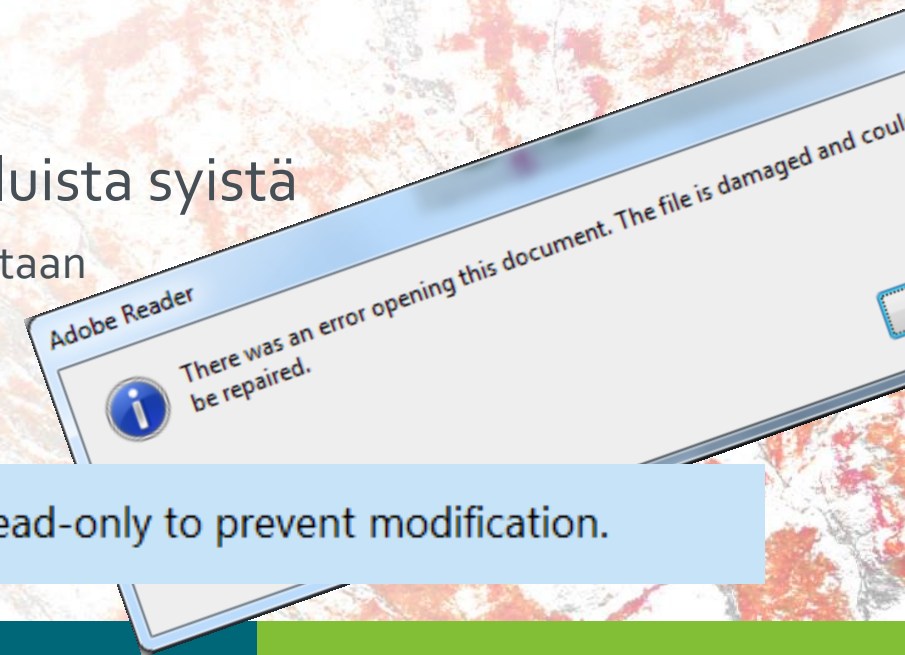
Meitä on enemmän kuin tarpeeksi...




- Mime-tyypit (~2000 rekisteröityä tiedostomuotoa (IANA))
  - Tunniste on merkkijono: application/pdf, text/plain, audio/x-wav ...
  - Ei versioita, joten siltä osin kattaa huomattavasti suuremman joukon tiedostomuotoja
  - Omia mime-tyyppejä on mahdollista käyttää: application/x-mun-oma-formaatti
- PRONOM tiedostomuotokirjasto (~1400 rekisteröityä tiedostomuotoa)
  - Pysyvät tunnisteet tiedostomuodoille: fmt/431, fmt/569, ...
  - Tiedostomuodon eri versioille omat tunnisteet
  - Pysyvät tunnisteet UK:n kansallisarkiston (TNA) myöntämiä; omien käyttäminen ei mahdollista
- Lisäksi lukuisia tiedostomuotoja joille ei ole mime-tyyppiä eikä PRONOM tunnistetta
  - Erityisesti laitevalmistajien omat eksoottiset suljetut tiedostomuodot
  - Osittain myös uudemmilta tiedostomuodoilta puuttuu toinen tai molemmat

# Tunnistaminen & validointi

- Tiedostomuodon tunnistaminen ei normaalisti ole riittävä toimenpide
  - On myös varmistuttava, että tiedosto on ko. tiedostomuodon määrittelyn mukainen
  - Jotkut ohjelmistot avaavat tiedoston normaalisti vaikka se olisi pilkulleen oikein
- PAS-palvelu validoi kaikki palveluun siirrettävät tiedostot...
  - Ja pienenkin virheen löytyessä PAS-palvelu ei ota säilytysvastuuta tiedostosta
  - Mahdollistetaan automaattiset massamigraatiot tulevaisuudessa
  - Validointi ei aina ole täydellistä (puutteet validointityökaluissa)
- Tietyissä tapauksissa validointi voidaan ohittaa perustelluista syistä
  - Korjaaminen voi olla mahdotonta, mutta tiedoston säilyttäminen katsotaan välttämättömäksi/tarpeelliseksi
  - Tällöin tiedosto otetaan vain bittitason säilyttämiseen

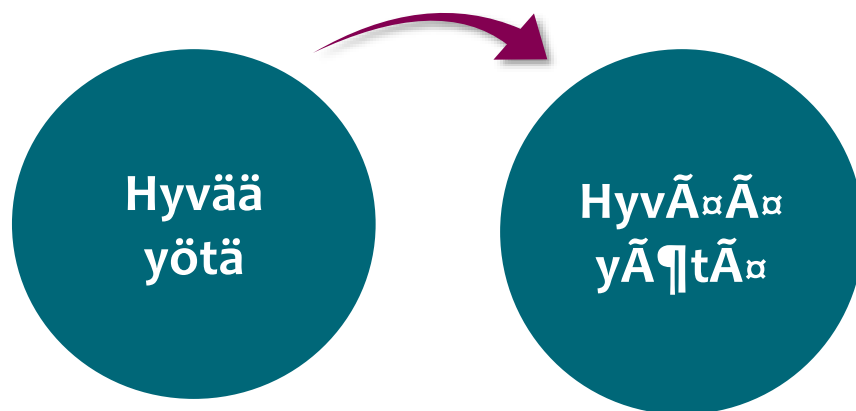


 This file claims compliance with the PDF/A standard and has been opened read-only to prevent modification.



# Tekstitiedostojen merkistöistä

- Monet tiedostomuodot sisältävät metatietoja miten tiedosto tulee tulkita
- Tekstitiedostoissa (plain text) ei mitään metatietoja ole, mutta käytettävissä olevia erilaisia merkistöjä on paljon
  - ANSI, UTF-8, ISO-8859-#, Windows-1257, ...
- Käytettävä merkistö on tallennettava johonkin muualle, jotta tiedosto voidaan tulkita oikein



# Tarkistussummat

Tiedostomuodosta riippumaton sormenjälki



- Tiedoston tarkistussummalla voidaan varmistua, että tiedosto ei ole ajansaatoissa (tahattomasti) muuttunut
  - Mutta ei siis takaa etteikö se voisi muuttua, mutta mahdolliset muutokset voidaan huomata koneellisesti
  - Jos muutoksia on tapahtunut, niin muuttunut tiedosto voidaan korvata eheällä kopiolla varmuuskopiosta (...jos sellainen on...)
- Tarkistussumma tulee laskea tiedostolle mahdollisimman aikaisessa vaiheessa
  - Aikaleima myös talteen
- PAS-palvelun tukemat tarkistussummat
  - md5, sha1, sha224, sha256, sha384 ja sha512
  - Kaikille käyttöjärjestelmille löytyy näiden laskemiseen valmiit ohjelmistot, joten ei tarvitse ymmärtää algoritmien toimintaa



# Tarkistussummat

Näin...

sha-1:310570

sha-1:445566 != sha-1:310570



PAS-palvelu



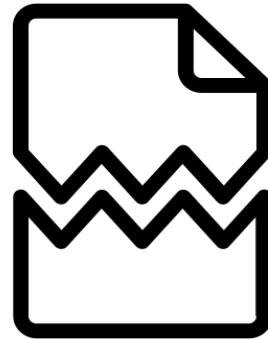


# Tarkistussummat

Ei näin...



sha-1:445566



sha-1:445566 == sha-1:445566  
=> Tiedosto "OK"



PAS-palvelu

# Tekninen metatieto

- Tiedostomuoto (mimetyyppi ja versio)
- Tiedoston eheystieto (tarkistussumma)
- Aineistotyyppikohtaiset tiedot (merkistö, kuvan korkeus, äänen näytteenottotaajuus jne.)
- (Lähes) kaikki tekninen metatieto on luettavissa tiedostoista
- Tarkistussumma on kuitenkin eheyden kannalta hyvä ottaa talteen mahdollisimman aikaisessa vaiheessa

# Tapahtumahistoria

- Aineiston tapahtumahistorialla kerrotaan mitä aineistolle on tapahtunut ja milloin
- Tapahtumahistoria esitetään tyypillisesti tapahtumina
- Tapahtumahistoria selittää miksi digitaalinen aineisto on juuri tämänkaltainen (onko se digitoitu, alkujaan digitaalinen, onko aineistoa muokattu, millä ohjelmistoilla tiedosto on käsitelty)
- Tapahtumahistoriaa on käytännössä mahdotonta luoda retroaktiivisesti
- Varsinkin aineiston synty (laite, ohjelmisto) on yleensä vaikeaa luoda jälkikäteen



# Hyödyntäviä organisaatioita ohjaavat PAS-määrittely



- Yksi PAS-palvelun näkyvimpiä osia hyödyntäville organisaatioille
- PAS-määrittelyt on tehty tiiviissä yhteistyössä hyödyntävien organisaatioiden kanssa

# Kansallinen tiedostomuotojen määrittely

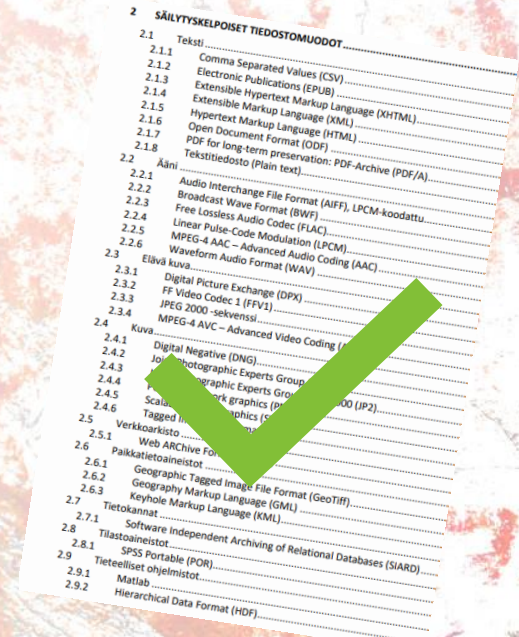
<https://urn.fi/urn:nbn:fi-fe2020100578095>



- Säilytys- ja siirtokelpoiset tiedostomuodot määrittely
  - "Recommended" vs. "acceptable for transfer"
  - Yleisesti hyväksytty jako kansainvälisesti
    - Toki muitakin lähestymistapoja on olemassa
- Tietyin ehdoin myös muita tiedostomuotoja voidaan ottaa säilytykseen
  - Organisaation pitää itse pystyä arvioimaan säilyttämisen riskit pitkällä tähtäimellä
  - Asiasta on kuitenkin aina sovittava erikseen
- Päivittyvä määrittely
  - Viime vuosina päivittynyt varsin maltillisesti
  - Määrittelyä päivitetään yhteistyössä hyödyntävien organisaatioiden kanssa
    - Mutta päivitykset aina tarve edellä

# Säilytyskelpoiset tiedostomuodot

- Säilytyskelpoinen tiedostomuoto on sellainen, jonka tietosisällön säilyminen ja ymmärrettävyys voidaan taata pidemmällä aikavälillä
  - Tällä hetkellä yhteensä 32 tiedostomuotoa
- PAS-palvelu vastaanottaa säilytyskelpoisia tiedostomuotoja siirtopaketeissa
- PAS-palvelu lupaa voivansa migroida hallitusti ja suunnitellusti
- Tiedostomuodot on jaettu aineistotyyppittäin, esimerkiksi
  - Kuva: JPEG, PNG, SVG, TIFF ...
  - Teksti: plain text, CSV, ODF, PDF ...
  - Tietokannat: SIARD, POR
  - Ääni: BWF, FLAC ...
  - Elävä kuva: FFV1, AVC ...

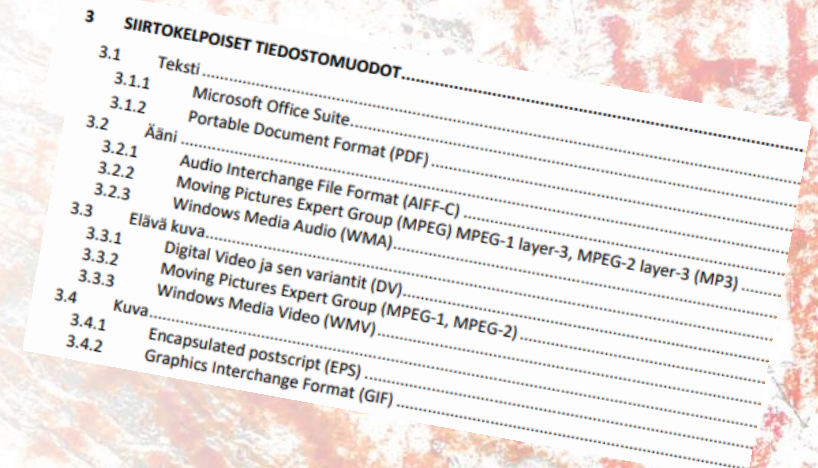


2	SÄILYTYSKELPOISET TIEDOSTOMUODOT
2.1	Teksti
2.1.1	Comma Separated Values (CSV)
2.1.2	Electronic Publications (EPUB)
2.1.3	Extensible Hypertext Markup Language (XHTML)
2.1.4	Extensible Markup Language (XML)
2.1.5	Hypertext Markup Language (HTML)
2.1.6	Open Document Format (ODF)
2.1.7	PDF for long-term preservation: PDF-Archive (PDF/A)
2.1.8	Tekstitiedosto (Plain text)
2.2	Ääni
2.2.1	Audio Interchange File Format (AIFF), LPCM-koodattu
2.2.2	Broadcast Wave Format (BWF)
2.2.3	Free Lossless Audio Codec (FLAC)
2.2.4	Linear Pulse-Code Modulation (LPCM)
2.2.5	MPEG-4 AAC - Advanced Audio Coding (AAC)
2.2.6	Waveform Audio Format (WAV)
2.3	Elävä kuva
2.3.1	Digital Picture Exchange (DPX)
2.3.2	FF Video Codec 1 (FFV1)
2.3.3	JPEG 2000 -sekvenssi
2.3.4	MPEG-4 AVC - Advanced Video Coding (AVC)
2.4	Kuva
2.4.1	Digital Negative (DNG)
2.4.2	Joint Photographic Experts Group
2.4.3	Joint Photographic Experts Group
2.4.4	Joint Photographic Experts Group
2.4.5	Scalable Vector Graphics (SVG)
2.4.6	Tagged Image File Format (TIFF)
2.5	Verkoarkisto
2.5.1	Web ARChive For
2.6	Paikkatietoaineistot
2.6.1	Geographic Tagged Image File Format (GeoTIFF)
2.6.2	Geography Markup Language (GML)
2.6.3	Keyhole Markup Language (KML)
2.7	Tietokannat
2.7.1	Software Independent Archiving of Relational Databases (SIARD)
2.8	Tilastoaineistot
2.8.1	SPSS Portable (POR)
2.9	Tieteelliset ohjelmistot
2.9.1	Matlab
2.9.2	Hierarchical Data Format (HDF)



# Siirtokelpoiset tiedostomuodot

- Siirtokelpoiset tiedostomuodot ovat sellaisia, joiden tietosisällön säilymistä tai ymmärrettävyyttä ei voida taata pidemmällä aikavälillä
  - Mutta joita on kuitenkin huomattavasti organisaatioilla
  - Siirtokelpoisia tiedostomuotoja ei pitäisi enää tuottaa
- Käytännössä siirtokelpoinen on pakko migroida (joskus) ensiksi säilytyskelpoiseen, eli se on askeleen säilytyskelpoisen jäljessä
- Käytössä vastaava jako aineistotyyppittäin kuin säilytyskelpoisissa muodoissa
  - Kaikilla aineistotyyppiryhmillä ei kuitenkaan ole siirtokelpoisia muotoja



3	SIIRTOKELPOISET TIEDOSTOMUODOT.....
3.1	Teksti.....
3.1.1	Microsoft Office Suite.....
3.1.2	Portable Document Format (PDF).....
3.2	Ääni.....
3.2.1	Audio Interchange File Format (AIFF-C).....
3.2.2	Moving Pictures Expert Group (MPEG) MPEG-1 layer-3, MPEG-2 layer-3 (MP3).....
3.2.3	Windows Media Audio (WMA).....
3.3	Elävä kuva.....
3.3.1	Digital Video ja sen variantit (DV).....
3.3.2	Moving Pictures Expert Group (MPEG-1, MPEG-2).....
3.3.3	Windows Media Video (WMV).....
3.4	Kuva.....
3.4.1	Encapsulated postscript (EPS).....
3.4.2	Graphics Interchange Format (GIF).....

# Yleiset rajoitukset tiedostomuodoille

- Tiedostoissa ei saa käyttää salasanasuojauksia eikä mitään muita salaustekniikoita
- Tiedostoissa ei saa käyttää DRM (Digital Rights Management) -tekniikoita
- Tiedostoja ei saa allekirjoittaa digitaalisesti, jos se estää tiedoston käsittelyn
- Tiedostoja ei saa (tarpeettomasti) pakata
- Tiedostosta ei saa puuttua sen esittämiseen tarvittavia ulkoisia komponentteja



# Tiedostomuoto vs. sen soveltaminen

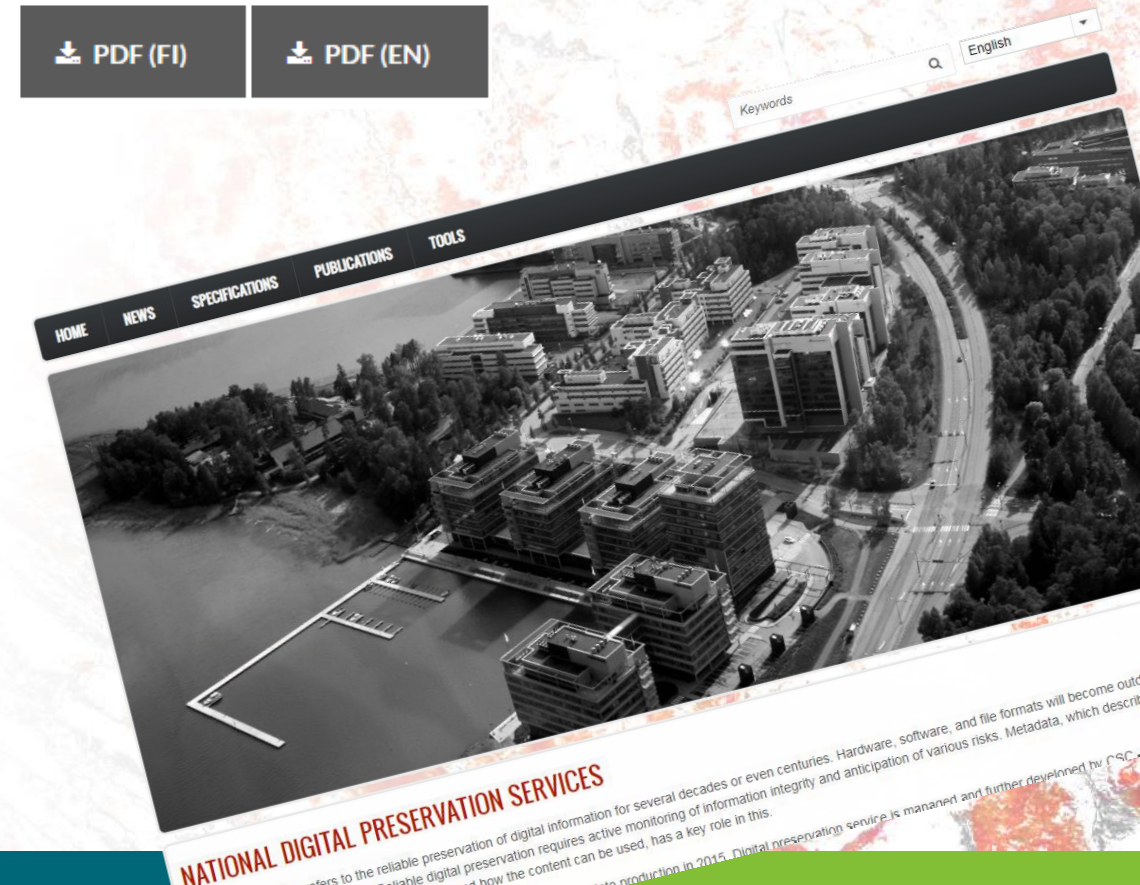
- Pelkkä tiedostomuoto ei takaa, että säilytettävä sisältö säilyy tarkoituksen mukaisesti, vaan sitä osattava soveltaa käyttötarkoituksen mukaisesti
  - Skannatussa asiakirjassa värisävyjen laadulla ei välttämättä ole suurta merkitystä, jos kyseessä on esimerkiksi konekirjoitettu asiakirja (kuvan pakkausmenetelmä voi kuitenkin vaikuttaa automaattiseen tekstin tunnistamiseen!)
  - Valokuvassa värien laadulla on yleensä aivan toisenlainen merkitys
- Pakkaamaton vs. pakattu
  - Muutamit tiedostomuodot mahdollistavat pakkaamisen (esim. jpeg)
  - Kannattaa aina käyttää mahdollisimman hyvää laatua ja pakata tiedostoa mahdollisimman vähän, erityisesti jos kyseessä on häviöllinen pakkaaminen (kuten jpeg)
  - Hyvä laatuinen voidaan aina myöhemmin muuntaa heikompi laatuiseksi, mutta muunnos toiseen suuntaan ei yleensä ole mahdollista



# Määrittelyt saatavilla

Myös englanniksi

- <http://digitalpreservation.fi/specifications>
  - Aineistojen ja niiden metatietojen paketointi pitkäaikaissäilytykseen
  - Säilytys- ja siirtokelpoiset tiedostomuodot
  - PAS-palveluiden rajapinnat
- Vuosittaiset laaturaportit ja muita julkaisuja
- Ohjelmistokoodia, skeemat & schematron säännöt saatavilla @GitHub
  - <https://github.com/Digital-Preservation-Finland/>



# Mitä pitää ottaa huomioon jo ennen PAS-palveluun siirtymistä?

- ✓ Laske tarkistussummat tiedostoille mahdollisimman aikaisessa vaiheessa
- ✓ Huolehdi laadukkaista metatiedoista
- ✓ Huomioi tiedostomuodot-määrittely ja mahdolliset soveltamisohjeet
- ✓ Dokumentoi aineiston syntyhistoria (digitointi, ...)
- ⊘ Älä muodosta siirtopaketteja "varastoon"

Kysyvä ei  
tieltä eksy

# Muuta huomioitavaa

- Kaikki toimenpiteet tiedostoille pitää dokumentoida koko elinkaaren ajan
  - Mitä tehtiin, koska, miksi, kuka teki,...
  - Myös ennen säilyttämisen aloittamista (ja erityisesti silloin)
  - Tapahtumat sanasto: <https://www.digitalpreservation.fi/specifications/vocabularies>
- Tiedostojen nimeämiseen kannattaa kiinnittää huomiota
  - "presentation1.pptx" vs "pas-museot-webinaari-2022-10-04-jkyl.pptx"
  - Paljon muita konkreettisia ohjeita webinaarissa: <https://youtu.be/Xkqkg1oiUOQ>



# Mitä taustajärjestelmään?

- Linkki aineiston kuvailun ja tiedostojen välillä
  - Mahdollisimman tarkka polku (ei pelkkä tiedostonnimi)
- Tiedoston tiedot:
  - Tarkistussumma (sekä käytetty algoritmi että laskettu summa)
  - Tiedostomuoto ja -versio
  - Tiedostokohtaiset tunnisteet
- Tapahtumahistoria
  - Tärkeimmät tapahtumat ja linkitykset tiedostoihin
  - Suorittaja (ohjelmisto) mukaan
- Kirjanpito siitä, mitä aineistoja on viety PAS-palveluun!



**YOU ARE NOT ALONE...**

pas-support@csc.fi  
digitalpreservation.fi  
@dpres\_fi