

Digitaalisten aineistojen pitkäaikaissäilytys - Tiedostformaattien standardointi

Valtakunnalliset kuva-arkistopäivät
4.-5.11.2013
Vesa Hongisto



Miksi tiedostformaatteja standardoidaan

Digitaalisen pitkäaikaissäilytyksen tasot

- **Bittien säilyttäminen:**
 - Alkuperäinen bittijono voidaan toistaa nykyaikaisilla laitteilla, mutta ei takeita siitä voidaanko sisältöä tulkita.
 - Säilyttämisen perusta ja optio tulevaisuuteen: antaa mahdollisuuden nostaa vaatimustasoa myöhemmin jos tarvetta ja rahaa on.
- **Sisällön ymmärrettävyyden säilyttäminen:**
 - Teksti, kuvat ja muu olennainen sisältö voidaan tulkita.
 - Ymmärrettävyys riippuu myös kohdeyleisöstä.
 - Käytön ja uuden tiedon tuottamisen kannalta yleensä tärkein taso.
- **Alkuperäisen käyttökokemuksen säilyttäminen:**
 - Kalvoesityksen tehosteet toistuvat alkuperäisessä muodossaan.
 - Vanha tietokonepeli näyttää ja kuulostaa samalta kuin 30 vuotta sitten.
 - **Vaativin**, mutta tutkimuksen kannalta mielenkiintoinen taso.

5.11.2013 / Vesa Hongisto



Miksi tiedostoformaatteja standardoidaan

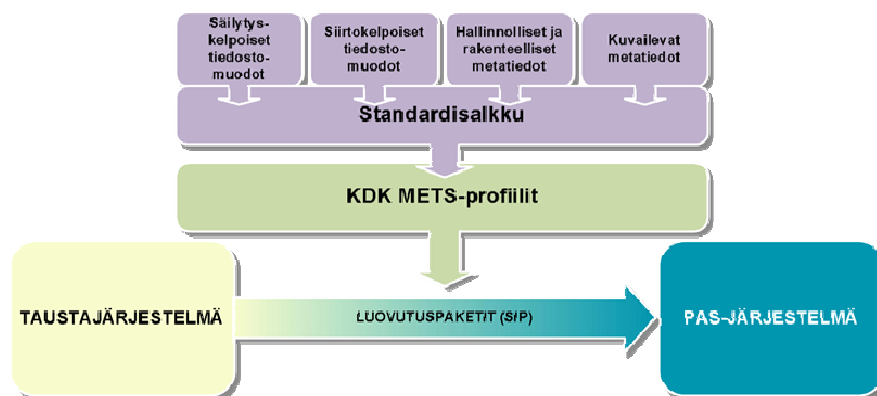
Digitaalisen pitkäaikaissäilytys

- Bittien säilyttäminen:
 - Alkuperäinen bittijono voi sisältää tulkittavaa sisältöä tulkita.
 - Säilyttämisen perusta ja myöhemmin jos tarve on.
- Sisällön ymmärrettävyyden säilyttäminen:
 - Teksti, kuvat ja muu olennainen sisältö voidaan tulkita.
 - Ymmärrettävyys riippuu myös kohdeyleisöstä.
 - Käytön ja uuden tiedon tuottamisen kannalta yleensä tärkein taso.
- Alkuperäisen käyttökokemuksen säilyttäminen:
 - Kalvoesityksen tehosteet toistuvat alkuperäisessä muodossaan.
 - Vanha tietokonepeli näyttää ja kuulostaa samalta kuin 30 vuotta sitten.
 - Vaativin, mutta tutkimuksen kannalta mielenkiintoinen taso.

Sisällön ymmärrettävyyden säilyttämiseksi tiedostot on muutettava kunakin ajanhetkenä käytössä oleviin formaatteihin. Tiedostoformaattien muunnossa häviää aina jonkin verran informaatiota.

5.11.2013 / Vesa Hongisto

HKDK:n tuottamat ohjaavat dokumentit



5.11.2013 / Vesa Hongisto

Säilytyssuunnitelma

- Jokaiselle pitkäaikaissäilytykseen tulevalle aineistolle tehdään säilytyssuunnitelma
- Säilytyssuunnitelmassa määritellään mitkä ovat ne ominaisuudet jotka eivät saa muuttua migraatiossa
- Jos on käytettävissä vaihtoehtoisia tiedostoformaatteja valitaan se, joka parhaiten takaa halutun ominaisuuden säilymisen
- Esimerkkejä:
 - Asiakirjassa voi riittää, että informaatio sisältö, teksti, säilyy, asettelu voi muuttua
 - Taidevalokuvassa on oleellista, että myös värisävyt säilyvät

5.11.2013 / Vesa Hongisto

Tiedostoformaattien arviointikriteerit

- **Avoimuus**
 - Kuinka helppoa tiedostomuodosta on saada tietoja?
- **Käyttö PAS-standardina**
 - Missä määrin tiedostomuoto on muodollisesti hyväksytty pitkäaikaissäilytyksen välineeksi kansalliskirjastoissa, kansallisarkistoissa ja muissa alan laitoksissa?
- **Vakaus / yhteensopivuus**
 - (a) Missä määrin tiedostomuoto on eteen- ja taaksepäin yhteensopiva?
 - (b) Missä määrin tiedostomuoto on suojattu tiedoston korruptoitumista vastaan?
 - (c) Kuinka usein tiedostomuodosta julkaistaan korvaavia versioita?
- **Riippuvuudet / yhteentoimivuus**
 - Missä määrin tiedostomuoto on sidottu esimerkiksi tiettyyn laitteistoon tai ohjelmistoon?
- **Standardisuus**
 - Missä määrin tiedostomuoto on käynyt läpi perusteellisen standardointiprosessin?

5.11.2013 / Vesa Hongisto

Tiedostiformaattien arviointi

Kriteeri	Arviointiohje	Arvio
1. Avoimuus	Määritykset saatavissa yhdeltä tai useammalta seuraavista:	A
	(a) avoimen jäsenyyden järjestö [kuten W3C (World Wide Web Consortium) tai OMG (Object Management Group)]	
	(b) Kansainvälinen standardointijärjestö (esim. ISO)	
	(c) Tuotannonalan avoimen jäsenyyden järjestö	
	Määritykset saatavilla vain maksusta	A [€]
	Määritykset mahdollisesti saatavilla useista lähteistä (ei saatu vahvistettua)	B
	Määrityksiä jakelee tai niiden jakelua valvoo yksi kaupallinen toimija tai pieni kaupallisten toimijoiden joukko.	C

5.11.2013 / Vesa Hongisto

Tiedostiformaattien arviointi

Kriteeri	Arviointiohje	Arvio
2. Käyttö PAS-standardina	Merkittävä osa (yli puolet) organisaatioista käyttää tai aikoo käyttää tiedostomuotoa säilytyskelpoisena tiedostomuotona	A
	Jotkut organisaatioista (alle puolet) käyttävät tai aikovat käyttää tiedostomuotoa säilytyskelpoisena tiedostomuotona	B
	Yksikään organisaatio ei käytä eikä aio käyttää tiedostomuotoa säilytyskelpoisena tiedostomuotona	C

5.11.2013 / Vesa Hongisto

Tiedostoformaattien arviointi

Kriteeri	Arviointiohje	Arvio
3. Vakaus / yhteen-sopivuus		
(a) Alas-/ ylöspäin yhteenso- pivuuden taso	Tiedostomuodon versio on alaspäin yhteensopiva, jos se sisältää kaikki aiemman version toiminnallisuudet	
	Tiedostomuodon versio on ylöspäin yhteensopiva, jos sen avulla voi ongelmitta tallentaa sisältöä, joka on tarkoitettu tiedostomuodon uudemmalle versiolle (kääntäen: ohjelmisto, joka on suunniteltu tulkitsemaan tai näyttämään tiedostomuodon aiempaa versiota pystyy tulkitsemaan tai näyttämään myös tiedostomuodon nykyistä versiota).	
	Hyvä yhteensopivuus – tiedostomuoto on yhteensopiva niin ylös- kuin alaspäinkin.	A
	Kohtalainen yhteensopivuus – tiedostomuoto on vain alaspäin yhteensopiva.	B
	Huono yhteensopivuus – tiedostomuoto ei ole yhteensopiva ylös- eikä alaspäin.	C

5.11.2013 / Vesa Hongisto

MUSEOVIRASTO

9

Tiedostoformaattien arviointi

Kriteeri	Arviointiohje	Arvio
3. Vakaus / yhteen-sopivuus		
(b) Korruptoitumisen sieto	Korruptoitumisen sieto tarkoittaa sitä, että tiedostomuoto sietää sisällön bitti- tai tavutason satunnaisia muutoksia.	
	Hyvä sietokyky – Muutokset eivät juurikaan tai lainkaan vaikuta tiedostomuodon A näyttämiseen tai tulkintaan; tai tiedostomuoto sisältää menetelmiä, joilla muutokset havaitaan ja/tai korjataan	
	Kohtalainen sietokyky – Muutokset vaikuttavat tiedostomuodon näytettävyyteen B mutteivät tulkittavuuteen; tiedostomuoto voi jossakin määrin palautua muutoksista.	
	Huono sietokyky – Kaikki muutokset vaikuttavat näytettävyyteen ja tulkittavuuteen.	C

5.11.2013 / Vesa Hongisto

MUSEOVIRASTO

10

Tiedostoformaattien arviointi

Kriteeri	Arviointiohje	Arvio
3. Vakaus / yhteen-sopivuus		
(c) Versio- päivitysten määrä	Tiedostomuodon vakaus ilmaistuna uusien versioiden tai laajennusten lukumääränä; tiedostomuodon käyttö johdannaisissa ja/tai tuotannonalakohtaisissa sovelluksissa.	
	Suuri vakaus	A
	Keskitason vakaus	B
	Epävakaus	C

5.11.2013 / Vesa Hongisto

Tiedostoformaattien arviointi

Kriteeri	Arviointiohje	Arvio
4. Riippumat- tomuus / yhteentoimivu- us	Korkea riippumattomuus & yhteentoimivuus	A
	Korkea riippumattomuus & keskitason yhteentoimivuus	
	Keskitason riippumattomuus & korkea yhteentoimivuus	
	Korkea riippumattomuus & alhainen yhteentoimivuus	B
	Keskitason riippumattomuus & yhteentoimivuus	
	Keskitason riippumattomuus & alhainen yhteentoimivuus	
	Alhainen riippumattomuus & yhteentoimivuus	C
	Alhainen riippumattomuus & keskitason yhteentoimivuus	
Alhainen riippumattomuus & korkea yhteentoimivuus		

5.11.2013 / Vesa Hongisto

Tiedostoformaattien arviointi

Kriteeri	Arviointiohje	Arvio
5. Standardisuus	Tiedostomuotoa sääntelee jonkin seuraavista hallitsema muodollinen prosessi: <ul style="list-style-type: none">• Avoimen jäsenyyden organisaatio (W3C, OMG yms.)• Kansainvälinen standardiorganisaatio (esim. ISO)• Tuotannonalakohtainen avoimen jäsenyyden organisaatio	A
	Tiedostomuotoa sääntelee dokumentoitu prosessi, jonka on luonut yksittäinen C kaupallinen toimija tai pieni kaupallisten toimijoiden joukko; tai tiedostomuotoa koskevia prosesseja ei ole dokumentoitu.	C

Tiedostoformaatit

- Säilytyskelpoiset tiedostomuodot
 - Säilytyskelpoiseksi hyväksytään sellaiset tiedostomuodot, joissa tietosisällön säilyminen ja ymmärrettävyys voidaan taata pidemmällä aikavälillä.
- Siirtokelpoiset tiedostomuodot
 - Siirtokelpoiseksi hyväksytään vain sellaisia tiedostomuotoja, joita käytetään useassa KDK:n PAS-järjestelmää hyödyntävässä organisaatiossa ja joissa pitkäaikaissäilytettävää aineistoa on runsaasti tallennettu. KDK:n PAS-ratkaisu muuntaa säilytysuunnitelman asettamien vaatimusten ja ehtojen mukaisesti siirtokelpoisessa muodossa vastaanotetut tiedostot säilytyskelpoiseen tiedostomuotoon ennen säilyttämisen aloittamista.

Tiedostoformaatit

Sisältö	Tiedostomuoto	Avoimuus	Käyttö PAS standardina	Vakaus / yhteensopivuus			Riippumattomuus / yhteensopivuus	Standardisuus
				Alas /löspäin yhteensopivuuden taso	Korruptoitumisen sieto	Verkopäivitysten määrä		
TEKSTI	Electronic Publications (EPUB)	A	B	B		A	A	A
	<u>Extensible Hypertext Markup Language (XHTML)</u>	A	B	B		A	A	A
	<u>Extensible Markup Language (XML)</u>	A	A	A		A	A	A
	<u>Hypertext Markup Language (HTML)</u>	A	A	B		A	A	A
	<u>Open Document Format (ODF)</u>	A	A	B		B	A	A
	<u>PDF for long-term preservation (PDF/A)</u>	A ^c	A	B		A	A	A
	Tekstitiedosto (<u>plain text</u>)	A ^c	A	B		A	A	A
ÄÄNI	Audio <u>Interchange File Format (AIFF)</u> , <u>PCMkoodattu</u>	A	A	A		A	A	A
	<u>Broadcast Wave Format (BWF)</u>	A	A	A		A	A	A
	<u>Free Lossless Audio Codec (FLAC)</u>	A	B	A	A	A	A	A
	<u>MPEG-4 AAC – Advanced Audio Coding (AAC)</u>	A	B	A			A	A
	<u>Waveform Audio Format (WAV)</u>	A	A	A		A	A	A
	<u>Motion JPEG 2000</u>	A ^c	A	A	A	A	A	A
ELÄVÄ KUVA	<u>Joint photographic experts group (JPEG)</u>	A ^c	A			A	A	A
KUVA	<u>Joint photographic experts group JPEG 2000 (JP2)</u>	A ^c	A			A	A	A
	<u>Tagged Image File Format (TIFF)</u>	A	A			A	A	A
VERKKOARKISTO	<u>Web ARChive Format (WARC)</u>	A	A		B	A	A	A
TIETOKANNAT	Määritellään myöhemmin							

5.11.201

ASTO

15

Tiedostoformaattien pakolliset metatiedot

- Aineistoja digitaaliseen pitkäaikaissäilytykseen vietäessä on tiedostoista annettava riittävän tarkka kuvaus niiden teknisistä ominaisuuksista, jotta myöhempi muunto uusiin formaatteihin on mahdollista
- KDK-hankkeessa on määritelty joidenkin aineistotyyppien osalta pakolliset tekniset metatiedot
- Määrittely on tehty muutoksina suhteessa USA:n kongressin kirjaston eri tiedostotyyppien skeemoihin
- Esimerkki: tekstitiedostoissa täytyy mainita käytetty merkkistö
- Teknisissä metatiedoissa on myös käytettävä samoja yksikäsitteisiä ilmaisuja esimerkiksi mittayksiköissä

Formaattikirjasto

- KDK:ssa käytetään pääsääntöisesti tiedostomuodoille PRPNOM formaattikirjaston tunnisteita
- KDK ylläpitää kontrolloitua sanastoa jossa määritetään miten kukin formaatti ja sen versio ilmaistaan



5.11.2013 / Vesa Hongisto

17

Tiedostomuotojen migraatiossa huomioitavat ominaisuudet

- Esimerkki tiedostomuunnoksessa säilyvistä, muuttuvista ja katoavista ominaisuuksista

Lähdemuoto	Kohdemuoto				Lisätiedot
	PNG	JPEG	JPEG 2000	TIFF	
GIF					
Datan häviötön säilyminen	X		X	X	PNG: muuttuu Deflate-muotoon, JPEG 2000: muuttuu Wavelet-muotoon, JPEG: muuttuu häviölliseen DCT-muotoon
Kuvatasot				(X)	Tarvittaessa eri tasoista voidaan tehdä useita kuvia. TIFF-muodossa ne saa kokoelmaksi.
Maksimidimensio	X	X	X	X	
Dimensio	X	X	X	X	
Esitysresoluutio (89a)	X		X	X	Konvertoituu uuteen muotoon. JPEG 2000-muodossa määritettävä kuvan fyysinen koko.
Kuvan kanavat	X	X	X	X	
Bittisyvyys	X	X	X	X	
Värikoordinaatisto	X	(X)	X	X	JPEG-muodossa konvertoitava YCbCr-koordinaatistoon.
Paletti	X		(X)	(X)	PNG-muodossa lokaalit paletit muuttuvat nk. ehdotetuiksi paleteiksi. JPEG 2000- ja TIFF-muodossa yhden paletin tuki.
Paletin bittisyvyys	X		X	X	TIFF-muodossa paletin bittisyvyys kaksinkertaistuu
Taustaväri	X				
Läpinäkyvyys (89a)	X		X	X	PNG-muodossa muuttuu läpinäkyvyysväriksi. JPEG 2000 ja TIFF-muodossa muuttuu kuvakanavaksi.
Lomitus	X	(X)	(X)		JPEG- ja JPEG 2000 -muodoissa muuttuu progressiiviseksi.
Kommenttikenttä (89a)	X	X	X	X	
Kuvaan upotettava teksti (89a)					
GIF-laajennusosio					
Animaatio				(X)	Animaation kuvat voidaan TIFF-muodossa koostaa kuvakokoelmaksi, mutta ilman animaation toiminnallisuutta.

5.11.2013 / Vesa Hongisto



Kiitos!

vesa.hongisto@nba.fi

www.nba.fi

<http://www.kdk.fi/fi/pitkaaikaissailytys/maeaerittely-ja-dokumentit/5-suomi/pitkaeikaissaailytys/141-kdk/sailytys-ja-siirtokelpoiset-tiedostomuodot>

