

SPRÍSTUPŇOVANIE A UCHOVÁVANIE VEDECKÝCH DÁT

Jitka Dobbersteinová, vzdelávací program Otvorená veda v praxi – Otvorené vedecké dáta, 17.10. – 19. 10. 2023



EURÓPSKA ÚNIA
Európske štrukturálne a investičné fondy
OP Integrovaná infraštruktúra 2014 – 2020

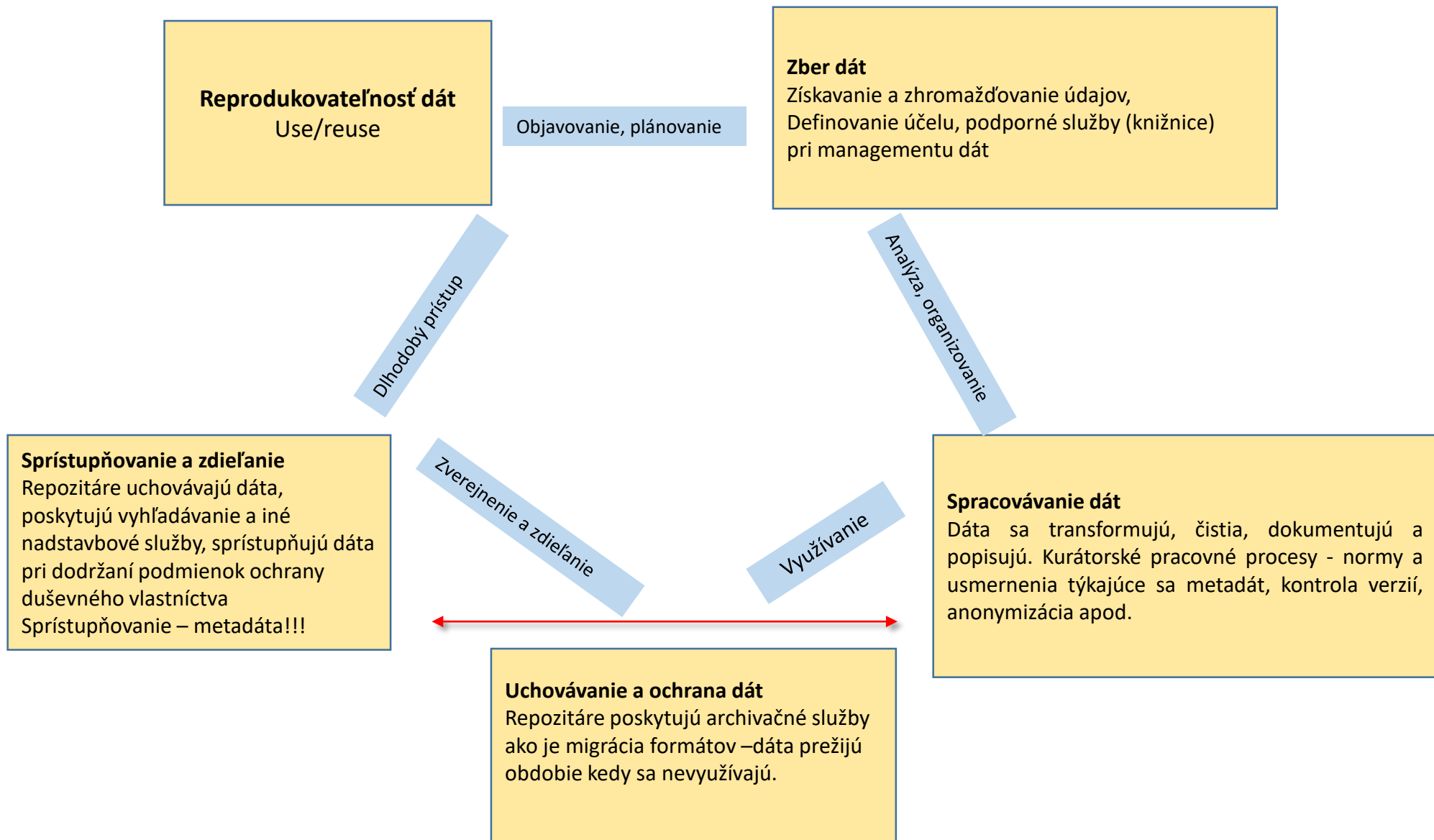


MINISTERSTVO
DOPRAVY A VÝSTAVBY
SLOVENSKEJ REPUBLIKY



MINISTERSTVO
ŠKOLSTVA, VEDY,
VÝSKUMU A ŠPORTU
SLOVENSKEJ REPUBLIKY





Legenda:

Žlté polia: kurátorský cyklus dát

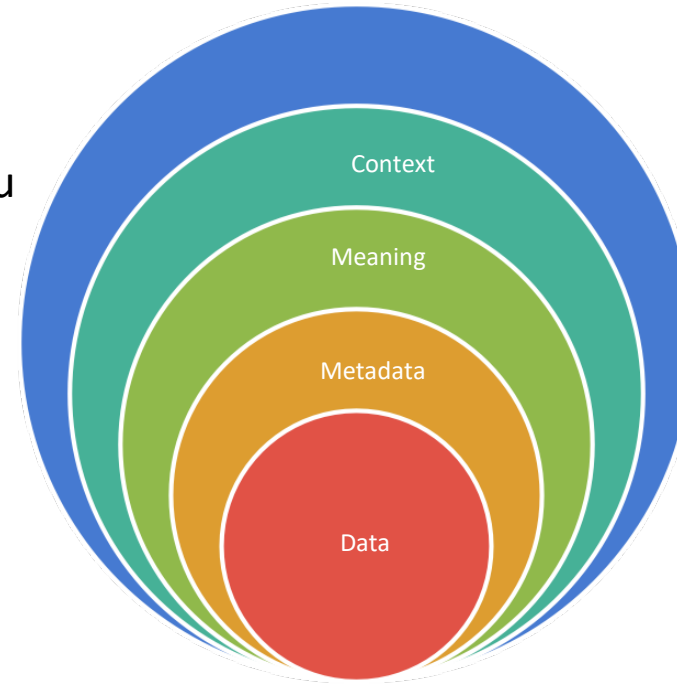
Modré šípky: životný cyklus managementu vedeckých dát

1

Získanie, analýza
Spracovanie

DÁTA vs METADÁTA?

- Digitálny objekt môže odvodzovať veľkú čas svojho významu z atribútov, ktoré nie sú samotnou jeho súčasťou
 - Kontext – okolnosti, za ktorých vznikol
 - Cesta – ako vznikol a dostal sa tam, kde je
 - To platí najmä pre dáta
- Metadáta môžu niesť viac významu ako samotné dáta
- **Vieme vôbec jednoznačne definovať metadáta?**



AKO ZABEZPEČÍTE ZROZUMITEĽNOSŤ DÁT PRE INÝCH VEDCOV...?

Opis použitých štandardov a metód pomôže pochopiť iným, ako ste realizovali projekt

(dokumentácia – kontextové informácie, na úrovni dát, na úrovni štúdie, nástroje; readme súbor)

AKO BUDETE DÁTA DOKUMENTOVAŤ, OPISOVAŤ?

Metadát je veľa – ktoré budú kľúčové pre projekt?

Metadáta môžu byť:

Opisné metadáta – informácie o obsahu a kontexte dát

Například: názov, tvorca, kľúčové slová predmetu a opis (abstrakt)

Štrukturálne metadáta – opisujú fyzickú štruktúru zložených dát

Například: použitý fotoaparát, clona, expozícia, formát súboru a vzťah k iným údajom alebo súborom

Administratívne metadáta – informácie používané na manažment dát

Například: kedy a ako boli vytvorené, kto k nim môže pristupovať, softvér potrebný na ich použitie a autorské práva

zdroj: <https://guides.lib.unc.edu/metadata/definition>



METADÁTA = DÁTA O DÁTACH

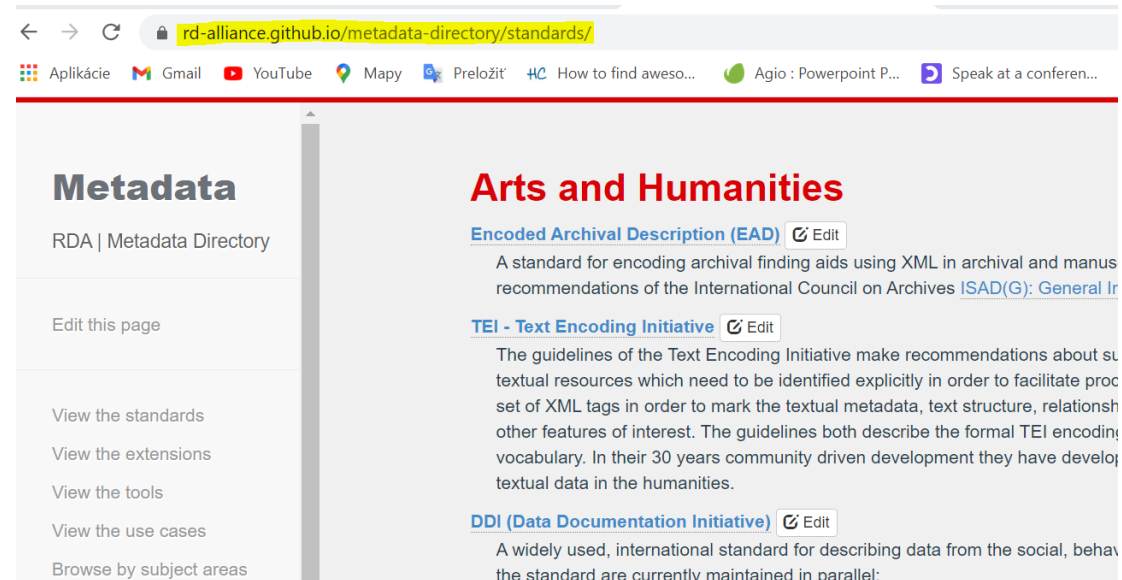
- Popisujú **charakteristiky** objektu/zdroja : obsah, formát, lokalita, prístupové práva...
- Popisujú **fyzické** alebo **digitálne** objekty
- **Forma:**
 - **voľný text** (readme file)
 - **štruktúrovaná a štandardizovaná** forma (MARC, DublinCore...)
- **Výskumné dáta** : nájditel'nosť, citovanosť a opätovné použitie
- **Spojené s dátami**, ktoré popisujú:
 - spolu v dátovom súbore
 - separátny súbor
 - záznam v katalógu/repozitári

Metadátové štandardy

<https://rd-alliance.github.io/metadata-directory/standards>

Štandardy:

- Dáta sú rovnako štruktúrované
- Štandard podľa disciplíny: **odborové metadátové štandardy**, napr. [DDI](#), [Darwin Core](#), [EML](#), [VRA Core](#)...
- Repozitáre si môžu určovať konkrétne štandardy
- ISO, Dublin Core Metadata Initiative....



<https://www.dcc.ac.uk/guidance/standards/metadata>

Sample Dublin Core metadata for an image

Below is an example of Dublin Core metadata for this photograph of Howard Walter Florey. [Generic examples](#) for the Dublin Core metadata schema are available from the [Dublin Core website](#).



- **Title** = Howard Walter Florey
- **Author/Creator** = Not known
- **Subject and Keywords** = The Australian National University photographs
- **Subject and Keywords** = Nobel Laureates
- **Description** = This photograph was taken at the 1948 Easter Conference, which considered the early planning of the Australian National University, at the Institute of Anatomy building. A keen amateur photographer, Florey carries his camera over his shoulder.
- **Publisher** = The Australian National University Digital Collections
- **Date** = 1948
- **Type** = image
- **Format** = image/tif
- **Resource identifier** = <http://hdl.handle.net/1885/7491>
- **Resource identifier** = ANUA225-405-3
- **Source** = ANU Archives
- **Rights** = This image is provided for research purposes only and must not be reproduced without the prior permission of the Archives Program, Australian National University

Metadátové schémy

- **štruktúra** metadát
- **špecifikuje elementy** = povinné, voliteľné, opakovateľné
- vychádza zo špecifickej komunity (knihovníci - MARC), popisuje špecifický formát al. doménu (ISO 19115)
- môže špecifikovať: obsahové pravidlá (formáty, riadené slovníky), syntax (napr. XML)

Example of MARC 21 record

```
=LDR 00000nam\a2200000\a\4500
=005 20111115115722.0
=008 111126s2008\\i\\a\\b\\001\0\eng\\
=020 \\$a9780073046235
=040 \\$aBD-DhIUB$cBD-DhIUB
=082 00$a332.1068$222
=100 1\$aRose, Peter S.
=245 10$aBank management & financial services /$cPeter S. Rose, Sylvia C. Hudgins.
=246 3\$aBank management and financial services
=250 \\$a7th ed.
=260 \\$aNew York :$bMacmillan,$c2008.
=300 \\$axxv, 722 p. :$bill. ;$c25 cm.
=504 \\$aIncludes bibliographical references and index.
=650 \0$aFinancial institutions$zUnited States.
=650 \0$aBank management$zUnited States.
=700 1\$aHudgins, Sylvia Conway,$d1956-
```

```
<head profile="http://dublincore.org/documents/dcq-html/">
  <title>Dublin Core</title>
  <link rel="schema.DC" href="http://purl.org/dc/elements/1.1/" />
  <link rel="schema.DCTERMS" href="http://purl.org/dc/terms/" />
  <meta name="DC.format" scheme="DCTERMS.IMT" content="text/html" />
  <meta name="DC.type" scheme="DCTERMS.DCMIType" content="Text" />
  <meta name="DC.title" content="Dublin Core" />
  <meta name="DC.publisher" content="Jimmy Wales" />
  <meta name="DC.subject" content="Dublin Core Metadaten-Elemente, Anwendungen" />
  <meta name="DC.creator" content="Björn G. Kulms" />
  <meta name="DCTERMS.license" scheme="DCTERMS.URI" content="http://www.gnu.org/copyleft/fdl.html" />
  <meta name="DCTERMS.rightsHolder" content="Wikimedia Foundation Inc." />
  <meta name="DCTERMS.modified" scheme="DCTERMS.W3CDTF" content="2006-03-08" />
</head>
```

dc.contributor.author	SCHÜTZE, Fabian	dc.rights	info:eu-repo/semantics/openAccess
dc.coverage.spatial	USA	dc.rights.uri	http://creativecommons.org/licenses/by/4.0/
dc.coverage.temporal	1959-2013	dc.subject	Statistical data
dc.date.accessioned	2019-10-16T15:43:18Z	dc.subject	Consumption
dc.date.available	2019-10-16T15:43:18Z	dc.subject	Equity prices
dc.date.created	2013-2014	dc.subject	E20
dc.date.issued	2014	dc.subject	E21
dc.identifier.other	EUI_ResData_00001_ECO	dc.subject	G12
dc.identifier.uri	http://hdl.handle.net/1814/64586	dc.subject.classification	YP-CS
dc.description	1 data file. Dataset elaborated from Robert Shiller stock market database, 1992-2013; and Federal Reserve Economic Data (FRED) on aggregate consumption.	dc.subject.ddc	317.3
dc.description.abstract	Dataset of adjusted monthly data for aggregate U.S. consumption expenditures and equity returns, 1959-2013, sourced from Robert Shiller stock market database, 1992-2013; and Federal Reserve Economic Data (FRED) on aggregate consumption - generated in the context of a research project illustrating an econometric technique, Generalised Method of Moments (GMM).	dc.subject.lcsh	Consumption (Economics) -- United States -- Statistics -- Databases
dc.format	Excel file	dc.subject.lcsh	Stocks -- Prices -- United States
dc.format.mimetype	application/vnd.openxmlformats-officedocument.spreadsheetml.sheet	dc.title	Dataset of adjusted monthly data for aggregate U.S. consumption expenditures and equity returns, 1959-2013
dc.language.iso	en	dc.type	Dataset
dc.publisher	European University Institute, ECO	dc.identifier.doi	10.2870/342675
dc.relation.ispartofseries	EUI Research Data	eui.subscribe.skip	true
dc.relation.ispartofseries	2014	dc.rights.license	Creative Commons Attribution 4.0 International
dc.relation.ispartofseries	Department of Economics		

	A	B	C	D	E	F
1	ID	author	date	subject.age	subject.gender	DataFile
2	sample 1	Ben Kok	-10-2015		male	FLProj_int_subj1.txt
3	sample 2		12-10-2015		male	FLProj_int_subj2.txt
4	sample 3		12-10-2015		female	FLProj_int_subj23.txt
5	sample 4		12-10-2015		female	FLProj_int_subj4.txt

<https://www.dcc.ac.uk/guidance/standards/metadata/tools>

List of Metadata Tools

AgriMetamaker

A service to facilitate the publication of metadata in the AGRIS Repository; it conforms to the the AGRIS Application Profile, which draws from the Dublin Core and AgMES standards.

ANZ-MEST - Metadata Entry and Search Tool

A GeoNetwork web application for metadata management and searching, with profiles available for two extensions of ISO 19115: ANZLIC and the Marine Community Profile.

AVM Adobe Metadata Panels

A set of metadata panels that can be added to Adobe Creative Suite 4 applications to allow AVM-compliant metadata to be entered directly into images.

AVM Web Tool

A web-based tool for assembling an AVM-compliant XMP packet for insertion into an image file.

Bio-Formats

Bio-Formats reads proprietary microscopy image data and metadata, and converts them to OME-TIFF, a combination of TIFF and OME-XML.

zdieľanie a
zverejnenie
dát

Princípy FAIR a vedecké dáta

AKO BUDETE DÁTA DOKUMENTOVAŤ, OPISOVAŤ?

Metadát je veľa – ktoré budú kľúčové pre váš projekt?

RIADENÉ SLOVNÍKY

Príklady a odkazy na riadené slovníky

<https://guides.lib.utexas.edu/metadata-basics/controlled-vocabs>

[Library of Congress
Authorities for subject
headings](https://authorities.loc.gov/)

<https://authorities.loc.gov/>

[ICD-11: International Classification of Diseases 11th
Revision](https://www.who.int/standards/classifications/classification-of-diseases)

[The global standard for diagnostic health information](https://www.who.int/standards/classifications/classification-of-diseases)

<https://www.who.int/standards/classifications/classification-of-diseases>

The Astronomy Thesaurus

<https://www.mso.anu.edu.au/library/thesaurus/>

[Art & Architecture Thesaurus® Online](http://www.getty.edu/research/tools/vocabularies/aat/)

Getty Research Institute

<http://www.getty.edu/research/tools/vocabularies/aat/>



2

Uchovávanie Ochrana

DLHODOBÉ UCHOVÁVANIE VÝSKUMNÝCH DÁT PO UKONČENÍ PROJEKTU

Ktoré z množstva získaných dát majú byť uchovávané dlhodobo?

Posúdiť hodnotu svojich dát, dlhodobo uchovávať len dáta s vysokou hodnotou:

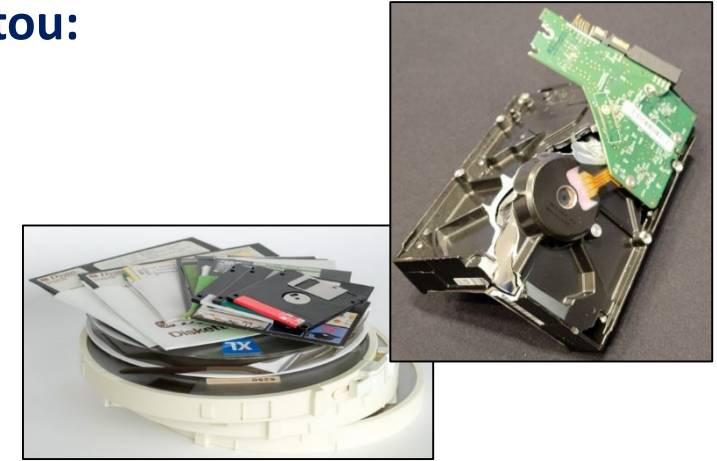
- všetky podkladové dáta súvisiace s publikáciou
- dáta, ktorých zber je náročný alebo nákladný
- dáta, ktoré sa nedajú nahradiť/znovu vygenerovať

KTORÉ DÁTA MAJÚ BYŤ ZNIČENÉ/VYMAZANÉ? AKO?

- biologický materiál, ktorý sa časom znehodnotí
- dôverné údaje alebo osobné údaje, ktoré nemožno anonymizovať

AKO DLHO PO UKONČENÍ PROJEKTU DÁTA UCHOVÁVAŤ?

- minimálne obdobie pre dáta, ktoré sú podkladom pre publikácie
- dodržiavať právne predpisy, nariadenia platné pre váš typ dát, pravidlá či usmernenia inštitúcie...



Nástroj DCC '[Five steps to decide what data to keep: a checklist for appraising research data](#)'

- Zabezpečenie dát v dôveryhodnom repozitári/úložisku
- Uchovávanie dát počas a po projekte
- Plány dlhodobej ochrany, archivácie (veľkosť dát, ako dlho, zničenie dát, finančné náklady, stratégie proti strate dát)

Výber riešenia pre uchovávanie dát <https://www.re3data.org/>

Typy riešení		
inštitucionálny repozitár	publikácie, ale aj dáta	
disciplinárny repozitár		https://www.re3data.org/
inštitucionálne úložisko	úložiská prispôbené na uchovávanie dátových súborov	
národná/nadnárodná infraštruktúra		národné repozitáre, Zenodo, World Data System
komerčné úložiská	úložiská prispôbené na uchovávanie dátových súborov	Google Drive, Figshare, Dropbox...

Bezpečnosť, formáty dát

- Nekompatibilita
- Univerzálné formáty: ASCII, Unicode
- „Bezpečné“ formáty
- Zálohovanie

Type	Preferred format(s)	Non-preferred format(s)
Text documents	<ul style="list-style-type: none"> • PDF/A (.pdf) • ODT (.odt) 	<ul style="list-style-type: none"> • Microsoft Word (.doc) • Office Open XML (.docx) • Rich Text File (.rtf) • PDF other than PDF/A (.pdf)
Markup language	<ul style="list-style-type: none"> • XML (.xml) • HTML (.html) • Related files: .css, .xslt, .js, .es 	<ul style="list-style-type: none"> • SGML (.sgml) • Markdown (.md)
Programming languages	<ul style="list-style-type: none"> • MATLAB • NetCDF • Text-Fabric 	
Spreadsheets	<ul style="list-style-type: none"> • ODS (.ods) • CSV (.csv) 	<ul style="list-style-type: none"> • Microsoft Excel (.xls) • Office Open XML Workbook (.xlsx) • PDF/A (.pdf)
Databases	<ul style="list-style-type: none"> • SQL (.sql) • SIARD (.siard) • CSV (.csv) 	<ul style="list-style-type: none"> • Microsoft Access (.mdb, .accdb) • dBase (.dbf) • HDF5 (.hdf5, .he5, .h5)
Raster images	<ul style="list-style-type: none"> • JPEG (.jpg, .jpeg) • TIFF (.tif, .tiff) • PNG (.png) • JPEG 2000 (.jp2) • DICOM (.dcm) 	
Vector images	<ul style="list-style-type: none"> • SVG (.svg) 	<ul style="list-style-type: none"> • Adobe Illustrator (.ai) • EPS (.eps) • WMF/EMF (.wmf, .emf)

DANS Dutch national centre of expertise and repository for research data

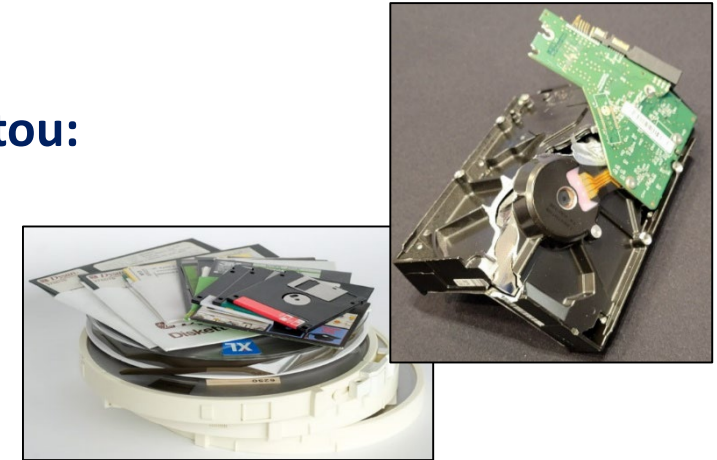
Uchovávanie dát

DLHODOBÉ UCHOVÁVANIE VÝSKUMNÝCH DÁT PO UKONČENÍ PROJEKTU

Ktoré z množstva získaných dát majú byť uchovávané dlhodobo?

Posúdiť hodnotu svojich dát, dlhodobo uchovávať len dáta s vysokou hodnotou:

- všetky podkladové dáta súvisiace s publikáciou
- dáta, ktorých zber je náročný alebo nákladný
- dáta, ktoré sa nedajú nahradiť/znovu vygenerovať



Nástroj DCC ['Five steps to decide what data to keep: a checklist for appraising research data,](#)

Nástroj na vyhodnotenie nákladov na uchovávanie dát [Data Management costing tool and checklist \(UK Data Service\)](#)

AKO DLHO PO UKONČENÍ PROJEKTU DÁTA UCHOVÁVAŤ?

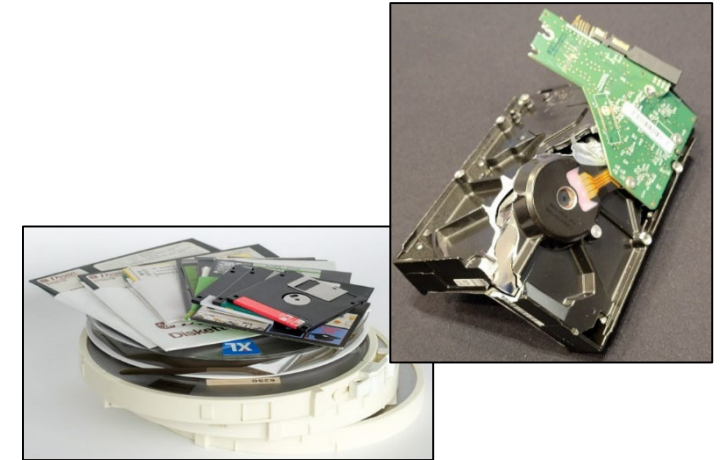
- minimálne obdobie pre dáta, ktoré sú podkladom pre publikácie
- dodržiavať právne predpisy, nariadenia platné pre váš typ dát, pravidlá či usmernenia inštitúcie...

Uchovávanie dát

KTORÉ DÁTA MAJÚ BYŤ ZNIČENÉ/VYMAZANÉ? AKO?

- biologický materiál, ktorý sa časom znehodnotí
- dôverné údaje alebo osobné údaje, ktoré nemožno anonymizovať
- usb kľúče
- papier, optické disky

Softvér na vymazanie súborov z pevných diskov: [BCWipe](#), [WipeFile](#), [DeleteOnClick](#) a [Eraser](#) pre platformy Windows; a [Permanent Eraser](#) pre platformy MacOS.





Ďakujem
Otázky?

