



národní
úložiště
šedé
literatury

Metodika pro tvorbu balíčků SIP se zaměřením na digitalizáty tištěných dokumentů

Cubr, Ladislav,; Fremrová, Květa; Jiroušek, Václav; Kočišová, Pavlína; Kopský, Vojtěch;
Miláček, Ivo; Ostráková, Natalie; Pavčík, Filip
2023

Dostupný z <http://www.nusl.cz/ntk/nusl-538239>

Dílo je chráněno podle autorského zákona č. 121/2000 Sb.

Licence Creative Commons Uveďte autora-Neužívejte dílo komerčně-Nezasahujte do díla 3.0 Česko

Tento dokument byl stažen z Národního úložiště šedé literatury (NUŠL).

Datum stažení: 20.05.2024

Další dokumenty můžete najít prostřednictvím vyhledávacího rozhraní nusl.cz .



Metodika pro tvorbu balíčků SIP se zaměřením na digitalizáty tištěných dokumentů

VERZE 2.0 (2023)

Autoři: Ladislav Cubr, Květa Fremrová, Václav Jiroušek,
Pavlína Kočišová, Vojtěch Kopský, Ivo Miláček, Natalie
Ostráková, Filip Pavčík

Realizováno v rámci institucionálního výzkumu Národní knihovny České republiky
financovaného Ministerstvem kultury ČR v rámci Dlouhodobého koncepčního rozvoje
výzkumné organizace

Oponenti:

- 1) Mgr. Pavlína Nimrichtrová, Národní archiv
- 2) Ing. Jiří Pavlík, Národní technická knihovna

Obsah

I. TEORETICKÁ ČÁST	6
ÚČEL METODIKY	6
URČENÍ	9
1 OBECNÁ ČÁST	10
1.1 Model OAIS	10
1.1.1 Koncept archivu	10
1.1.2 Prostředí archivu OAIS	12
1.1.3 Informační model OAIS	14
1.1.3.1 Informační obsah a interpretační informace	16
1.1.3.2 Archivační informace	19
1.1.3.3 Informační balíček	20
1.1.4 Přístupy k uchování	22
1.1.5 Specifická standardizace informačních balíčků	25
1.2 Formát objektu CDO	25
1.2.1 Roviny a aspekty užití formátu	27
1.2.2 Výběr archivačního formátu	28
1.2.2.1 Archivační formáty pro obrazovou komponentu	30
1.2.2.2 Archivační formát pro OCR komponentu	32
1.2.3 Prezentací formáty	33
1.2.3.1 Prezentací formáty pro obrazovou komponentu	33
1.2.4 Formátové registry	33
1.2.5 Specifikace obrazových dat	35
1.2.5.1 Generace obrazových dat	35
1.2.5.2 Strukturální model obrazové reprezentace	36
1.2.5.3 Formátový profil	37
1.2.5.4 Typy komprese	37
1.2.5.5 Obecné obrazové vlastnosti	39
1.2.5.6 Prezentací varianty	41
1.3 Metadatové standardy	42
1.3.1 Přehled metadatových standardů pro digitalizáty tištěných dokumentů	44
1.3.2 PREMIS	45
1.3.3 METS	46
1.3.4 MODS	47
1.3.5 Dublin Core	48
1.3.6 MIX	49
1.3.7 Metadata v obrazových souborech	49
2 SPECIFICKÁ ČÁST	51
2.1 Standard NDK a související předpisy	51
2.1.1 Metadatový aplikační profil	52
2.1.2 Standardy pro metadata	53
2.1.3 Standardy pro formáty	55
2.1.4 Standardy pro obrazová data	56
2.1.5 Podmínky VISK 7 pro rok 2023	57

II. IMPLEMENTAČNÍ ČÁST	59
Úvod	59
3 PLÁNOVÁNÍ DIGITALIZAČNÍHO PROJEKTU	60
3.1 Digitalizační projekt	60
3.2 Technické zajištění	61
3.2.1 Snímací zařízení	61
3.2.2 Softwarové nástroje pro tvorbu digitalizátů	62
3.2.3 Validační nástroje	62
3.2.4 Kontrola předloh	63
3.3 Základní standardizační doporučení	63
3.3.1 Stanovení základní intelektuální entity a granularity	63
3.3.2 Složení SIP balíku podle Standardu NDK	64
3.3.3 Perzistentní identifikátory tištěné předlohy	65
3.3.4 Perzistentní identifikátory digitalizátu	66
3.3.4.1 URN:NBN	66
3.3.4.2 UUID	68
3.3.5 Projektová dokumentace	68
4 DIGITALIZACE	70
4.1 Příprava bibliografických záznamů	70
4.2 Snímání předloh	70
4.2.1 Věrnost reprodukce tištěné předloze	70
4.2.2 Základní parametry pro skenování	71
4.2.3 Závislost rozlišení obrazu a typografických bodů písma	72
4.2.4 Snímkový formát	73
4.2.5 Zabudování EXIF metadat do souborů	73
4.2.6 Barevný profil	73
4.3 Zpracování dat	75
4.3.1 Zpracování obrazové komponenty	75
4.3.2 Archivační formát (formát pro archivní kopie)	76
4.3.2.1 Druh komprese a transformace	76
4.3.2.2 Kompresní poměr	76
4.3.2.3 Dlaždice	77
4.3.2.4 Průběh zobrazení	77
4.3.2.5 Dekompoziční úrovně	77
4.3.2.6 Vrstvy kvality	78
4.3.2.7 Regiony	78
4.3.2.8 Zájmové oblasti	78
4.3.2.9 Velikost bloků	79
4.3.2.10 Lokalizace dlaždice	79
4.3.2.11 Přemostění	79
4.3.2.12 ICC profily	79
4.3.2.13 Hlavička segmentu packetů	80
4.3.2.14 Vložená metadata	80
4.3.3 Prezentační formát (formát pro uživatelské kopie)	80
4.3.4 Vytváření OCR komponenty	81

4.4 Vytváření metadat	82
4.4.1 Převod bibliografických metadat	82
4.4.2 Získávání technických metadat	82
4.4.3 Nástroje pro formátovou identifikaci	83
4.4.4 Propojování událostí s objektem a agentem	84
4.4.5 Záznam událostí a nástrojů	85
4.4.5.1 Snímání (skenování, fotografování)	86
4.4.5.2 Formátová konverze z TIFF do JP2 archivní kopie	86
4.4.5.3 Vytvoření ALTO XML z OCR	87
4.4.5.4 Formátová identifikace	87
4.4.5.5 Formátová validace	88
5 KONTROLA KVALITY	89
5.1 Digitální otisk	89
5.2 Validace metadat	90
5.3 Formátová validace	91
5.4 Datová (obrazová) validace	92
5.5 Validace balíčku SIP	93
TERMINOLOGIE	94
SEZNAM ZKRATEK	102
PŘÍLOHA – MAPOVÁNÍ VÝSTUPŮ METADATOVÝCH EXTRAKTORŮ DO METADAT BALÍČKŮ SIP	106
O AUTORECH	112
CITOVANÁ LITERATURA	113

I. TEORETICKÁ ČÁST

Účel metodiky

Metodika předkládá postup pro vytváření balíčků SIP pro digitalizáty tištěných dokumentů, popisuje procedurální postup pro užití Standardu NDK pro metadata, formáty a obrazová data užívaná při digitalizaci tištěných dokumentů v českých knihovnách. Standard NDK je závazný pro Národní knihovnu ČR a Moravskou zemskou knihovnu v Brně pro digitalizaci v rámci pokračující spolupráce, navazující na projekt Vytvoření Národní digitální knihovny, a dále pro knihovny digitalizující na základě podpory získané z podprogramu VISK 7, případně pro další organizace, které budou odevzdávat svoje digitalizáty do LTP úložiště NK ČR k dlouhodobému uchovávání. Také díky široké podpoře Standardu NDK ze strany open-source nástrojů pro digitalizaci (ProARC), zpřístupňování (Kramerius) i archivaci (ARCLib) je Standard NDK využíván již přibližně 140 knihovnami a jinými paměťovými institucemi. Metodika pro balíčky SIP obsahuje i některá obecnější doporučení, která lze vztahovat i na jiné typy dokumentů než digitalizáty tištěných dokumentů (monografií a periodik), primárně se však zaměřuje na tento typ dokumentů.

Předkládaná metodika je druhou verzí původní metodiky se stejným názvem. První verze přinášela v době svého vydání v roce 2018 pro oblast digitalizace knihovních fondů metodiku, která v českém prostředí dosud chyběla. V následujících pěti letech byla tato metodika využívána při digitalizaci v Národní knihovně i v dalších knihovních institucích, zapojených do digitalizace pod programem VISK 7. Na základě nasbíraných zkušeností jasněji vyvstaly otázky, které bylo třeba v metodice vyjasnit, zpřesnit či aktualizovat. Taktéž bylo třeba reflektovat i změny v mezinárodních standardech. Nová verze metodiky tedy na původní verzi staví, ale navíc zohledňuje všechny tyto potřeby. Pro oblast digitalizátů tištěných dokumentů tak nadále zůstává jedinečnou.

V českém prostředí dále existuje pouze Metodika pro vytváření bezpečnostních kopií archiválií v digitální podobě (Dvořák a kol., 2015), která byla vytvořena Národním archivem ČR. Tato metodika se však vztahuje na jiný typ dokumentů (digitalizáty archiválií) a pokrývá celý životní cyklus dokumentu. V oblasti digitalizace podrobněji popisuje některé technologické otázky digitalizace (výběr a kalibrace skeneru, barevná specifikace apod.), ale nikoliv procedurální postup v takové podobě jako tato metodika.

V zahraničí vznikly metodiky pro vytváření digitalizátů zaměřené na otázku obrazových dat a částečně metadat. Za nejvýznamnější v tomto směru lze považovat směrnice „Technical Guidelines for Digital Cultural Content Creation Programmes“ z roku 2008 a směrnice americké iniciativy FADGI (Federal Agencies Digital Guidelines Initiative), které byly vytvořeny pro paměťové instituce v USA a nesou název „Technical Guidelines for Digitizing Cultural Heritage Materials: Creation of Raster Image Files“ (FADGI, 2010; FADGI, 2016; FADGI 2023). Směrnice vyšla v roce 2010, v roce 2016 a 2023 pak byly vydány její revidované verze. Poskytují zejména technologické specifikace z hlediska vlastního snímání a kontroly kvality, nejsou však specifické pro užití konkrétního metadatového profilu a formátu. Obě směrnice lze doporučit pro potřeby výběru snímacího zařízení, jeho kalibrace, volbu prostorového rozlišení, způsob zpracování obrazových dat a práci s barevným prostorem. Oba dokumenty jsou také využity v teoretické a praktické části předkládané metodiky. Metodika pro balíčky SIP obsahuje specifická procedurální doporučení pro vytváření konkrétního typu digitálního dokumentu s užitím konkrétního metadatového a formátového profilu (které jsou popsány ve Standardu NDK) a konkrétních nástrojů v českém prostředí digitalizačních projektů knihoven.

Mimo záběr metodiky je dlouhodobé uchovávání i zpřístupňování digitalizátů. Pro dlouhodobé uchovávání lze užit Metodiku logické ochrany digitálních dat (Hutař a kol., 2018), která vznikla v projektu NAKI ARCLib. Tato metodika obsahuje v teoretické a praktické části pasáže, které lze do určité míry užit pro dlouhodobé uchovávání digitálních dokumentů v různých organizacích užívajících různá softwarová řešení. Její implementační část je pak zaměřena specificky na konkrétní softwarový systém, systém ARCLib.

Specifičnost pro určitý typ dokumentu, konkretizace metadatových, formátových a datových (obrazových, textových ve smyslu OCR) otázek je to, co činí Metodiku pro tvorbu balíčků SIP novou v českém prostředí. Zkušenosti z jejího využívání byly v roce 2021 vyhodnoceny a výsledkem bylo zjištění potřeby některé části metodiky doplnit, zpřesnit a po pěti letech aktualizovat tak, aby byly uživatelsky jasnější a využitelnější. Nová verze metodiky reaguje na dotazy uživatelské komunity, zkušenosti z provozu Komplexního validátoru NDK či nedostatky v datech, odevzdávaných v uplynulých letech do LTP úložiště NK ČR. Potřebným způsobem rovněž reaguje na rozvoj mezinárodních standardů, zahraničních směrnic pro oblast digitalizace, dostupných softwarových nástrojů, ale také změny v samotném Standardu NDK.

Doplněno bylo několik nových kapitol a podkapitol, nově byly přidány například části

obsahující podrobnější doporučení pro pořizování obrazových kopií týkající se správy barev, určení vhodného rozlišení skenování či práce s identifikátory. Metodika jako celek prošla celkovou aktualizací, rozšířením o nové celky, v částech, které se změnilo jen částečně, pak došlo k přepracování struktury informací, a doplněním dalších vazeb.

Konkrétně byl v teoretické části metodiky doplněn popis informačních balíčků převzatých z rámce OAIS o definici a způsob vzniku produkčního (PSP) balíčku, který, ač je v české digitalizační praxi používán, dosud v popisu chyběl. Došlo též k přepracování a rozšíření kapitoly 1.3 Metadatové standardy, kde byly kompletně doplněny podkapitoly pro formáty Dublin Core a MODS, a návazných kapitol ve Specifické části metodiky, věnujících se aplikačnímu profilu Národní knihovny a jejím standardům v digitalizaci NDK.

Značných změn doznala Specifická část metodiky (Část 2). Některé její podkapitoly (např. o prezentačních a archivačních formátech, tj. současné kapitoly 1.2.2.1, 1.2.2.2 a 1.2.3.1) byly jako celek přesunuty na relevantní místa v první kapitole. Jiné byly upraveny a staly se součástí kapitoly 3. Stávající druhá kapitola tak nyní tvoří koherentní celek zaměřený na způsob převedení obecných metadatových a formátových standardů na český metadatový aplikační profil zohledňující národní specifika. Nové uspořádání Metodiky by mělo zajistit lepší přehlednost pro její uživatele, a tedy zvýšit její potenciál jakožto praktické příručky pro instituce, které při digitalizaci využívají Standard NDK, případně také pro žadatele o participaci v programu VISK 7.

Ve Specifické části metodiky bylo dále doplněno a aktualizováno použití identifikátorů a jejich přidělování pomocí systému ČIDLO. Nově prošly aktualizací a doplněním veškeré požadavky na standardy pro obrazová data. V nové verzi metodiky byl uveřejněn aktuální soupis pravidel a požadavků pro žádost o grant z programu VISK 7, k některým původním požadavkům byl dodán kontext.

Zcela nová je podkapitola věnovaná UUID, která tento identifikátor jednak obecně představuje, jednak blíže popisuje způsob jeho využívání v rámci NDK. V technických kapitolách metodiky kompletně prošly revizí části, věnující se konceptu formátových politik (v oblasti souborových formátů), správy barevných profilů a byly aktualizovány veškeré pokyny, které souvisí s metodikou digitalizací obrazů FADGI, aby korespondovaly s pokyny, uvedenými v současné třetí edici.

Metodika byla v novém vydání opatřena podrobným slovníkem pojmů, vycházejícím z původního krátkého slovníku na začátku implementační části, a doplněna seznamem používaných zkratk.

Určení

Metodika rozšiřuje standardizaci v oblasti digitalizace textových dokumentů v knihovnách a jiných paměťových institucích v ČR. Je určena všem organizacím, které využívají Standard NDK, ať již povinně, neboť pak odevzdávají své balíčky SIP do LTP úložiště NK ČR k dlouhodobému uchovávání, nebo využívají Standard NDK dobrovolně, jako cesty ke standardizaci výstupů digitalizace. V obou případech je účelem Standardu NDK i Metodiky pro tvorbu balíčků SIP vytvářet data a metadata při digitalizaci v takové podobě, aby při tvorbě balíčku AIP nebylo třeba užít formátovou normalizaci (čímž dochází k minimalizaci nákladů pro archiv) a aby byly během produkce dokumentů zachyceny všechny důležité informace (metadata) pro následné dlouhodobé uchovávání a dodrženy postupy, přičemž obojí by mělo zaručit zvýšení kvality vytvářených digitalizátů, jejich autenticity a předjímaných požadavků dlouhodobého uchovávání i zpřístupňování. Metodika v neposlední řadě představuje vhodné postupy tvorby dat tak, aby byla minimalizována potřeba nákladných dodatečných datových oprav či aktualizací metadat.

Metodika je určena všem uživatelům Standardu NDK, kteří jsou v roli producenta digitálních dat. Cílovou skupinou metodiky jsou knihovny všech úrovní nebo obdobné instituce (např. muzea bez registrované knihovny, ale s knihovní sbírkou) a externí producenti dat zajišťující digitalizaci. Celkový počet uživatelů Standardu NDK je možné dle počtu registrátorů v systému ČIDLO (Český systém pro identifikaci a lokalizaci dokumentů digitálního kulturního dědictví) nepřímě odhadovat na přibližně 140 knihoven a jiných paměťových institucí vlastnících knihovní fondy.

Základní vymezení uživatelské komunity je následující:

- Národní knihovna České republiky, Moravská zemská knihovna v Brně
- krajské knihovny,
- specializované knihovny,
- knihovny vědeckých ústavů,
- vysokoškolské knihovny,
- knihovny archivů, muzeí a galerií,
- producenti digitalizovaných dat kulturní povahy.

1 Obecná část

1.1 Model OAIS

Základní konceptuálním rámcem pro oblast dlouhodobého uchovávání digitálních dokumentů (digitální archivace) je model OAIS (Open Archival Information System) obsažený v normě ISO 14721 (česky ČSN ISO 14721). Tento model je již od počátku milénia (kdy byl již znám jako návrh normy) základním a obecně přijímaným referenčním rámcem pro řízení životního cyklu digitálních dokumentů z hlediska digitální archivace. Model OAIS se primárně zabývá archivací a zpřístupněním digitálních dokumentů, ale do určité míry vymezuje i otázku produkce digitálních dokumentů. Činí tak v podobě konceptu dohody o dodávání dat (*submission agreement*), která je vyjádřením toho, v jaké podobě budou digitální objekty dodávány do archivu, a je výsledkem dohody mezi producentem (vkladatelem) a archivem. Tato podoba je v normě ČSN ISO 14721 specifikována jako tzv. vstupní informační balíček (*submission information package*, balíček SIP). V praxi může docházet k tomu, že archiv nemá možnost ovlivnit tuto podobu způsobem, který odpovídá požadavkům digitální archivace, například tehdy, pokud má legislativou stanovené povinnosti přijímat a archivovat digitální publikace od vydavatelů a současně platí, že a) legislativa vydavatelům neukládá povinnost dodávat digitální publikace do archivu ve formátech, které archiv preferuje (z důvodů jejich vhodnosti pro archivaci), b) archiv nemá možnost provádět formátovou konverzi publikací do požadované podoby (např. z důvodu neexistence adekvátních nástrojů pro takovou konverzi). Optimální situací je, pokud archiv může specifikovat podobu balíčku SIP tak, aby při převodu do balíčku AIP nebylo potřeba provádět formátovou normalizaci. To je možné v případě, kdy producent i archiv jsou jedna a tatáž organizace (a tudíž za cílové formáty digitalizace lze zvolit formáty vhodné pro archivaci) a kdy jsou dostupné technologie pro danou oblast (tj. existují vhodné formáty a adekvátní konverzní nástroje), nebo kdy je archiv taková organizace, která má pravomoc nebo autoritu (jako metodické centrum) vydávat pokyny pro standardizaci pro producenty v nějaké oblasti (např. knihovnictví).

1.1.1 Koncept archivu

Norma ISO 14721 pojednává o specifickém modelu archivu, který označuje jako „otevřený archivační informační systém“ (*open archival information system*), zkráceně

archiv OAIS a definuje jej jako „archiv (*archive*), který sestává z organizace tvořené lidmi a systémy, jež přijala odpovědnost uchovávat informace a zpřístupňovat je cílové komunitě (*designated community*), přičemž tento archiv může být součástí větší organizace.“ (ISO 14721:2012, s. 24)

Norma ISO 14721 vymezuje šest obecných závazných povinností, které musí daná organizace plnit, aby mohla provozovat archiv OAIS. Tyto povinnosti odlišují pojem „archiv OAIS“ od jiných užití pojmu „archiv“. Jde o následující povinnosti (ISO 14721:2012, s. 39):

1. Vyjednávat s producenty informací a přijímat od nich příslušné informace.
2. Získávat možnost s poskytnutými informacemi dostatečně nakládat, aby bylo možné zajistit jejich dlouhodobé uchování.
3. Určit, ať již samostatně nebo ve spolupráci s dalšími stranami, které komunity by se měly stát cílovými komunitami, a tudíž by měly být schopny porozumět poskytovaným informacím. Tím je vymezena znalostní základna dané skupiny.
4. Zajistit, aby informace určené k uchování byly pro cílovou komunitu srozumitelné samy o sobě. Cílová komunita by měla být schopna informacím porozumět bez využití odborných zdrojů, například bez rady odborníků, kteří informace vytvořili.
5. Dodržovat zdokumentovaná pravidla a postupy, které zajistí, že informace budou chráněny před všemi možnými nepředvídatelnými událostmi (včetně zániku archivu), a zajistit, že informace nebudou nikdy smazány (s výjimkou případu, kdy jejich smazání bude součástí schválené strategie). Nemělo by docházet k žádnému jednorázovému mazání dat.
6. Zpřístupňovat uchovávané informace cílové komunitě a umožňovat šíření informací v podobě kopií původně dodaných datových objektů nebo v takové podobě, aby bylo možné zpětně dohledat (*as traceable to*), ke kterým původně dodaným datovým objektům se vztahují, a to společně s doklady o jejich autenticitě.

První povinnost zahrnuje zejména dohodu o dodávání dat. Druhá povinnost zahrnuje vyřešení otázky práv duševního vlastnictví (pokud archiv nemá právo dodaná data měnit, například je převést do jiného formátu, je jeho činnost *de facto* znemožněna). Třetí povinnost se týká cílové komunity z hlediska jejího vymezení. Čtvrtá povinnost se týká cílové komunity

s ohledem na související koncepty nezávislé srozumitelnosti (tj. informace srozumitelné samy o sobě) a interpretačních informací. Pátá povinnost specificky pojednává o problematice dlouhodobého uchovávání navzdory změnám technologií (datových nosičů, formátů apod.). Šestá povinnost se týká zpřístupňování informací a specificky zohledňuje otázku autenticity.

Za účelem plnění těchto funkcí norma přináší tzv. funkční model. Funkční model OAIS vymezuje šest základních funkčních celků (*functional entities*) archivu OAIS. Funkční celek je konceptuální model představující kategorii činností (funkcí a služeb), které musejí být archivem vykonávány, aby byl archivem OAIS. Výkon těchto činností pro digitální archivaci zajišťují počítačové technologie stejně jako lidské pracovníky (nejen jako obsluha, ale i jako tvůrci návrhu softwaru apod.). Těmito funkčními celky jsou: Příjem, Archivní úložiště, Správa dat, Administrace, Plánování uchovávání a Zpřístupnění (ISO 14721:2012, s. 45). Sedmým funkčním celkem jsou základní služby (*common services*), jimiž jsou „podpůrné služby nezbytné pro provoz archivu OAIS“ (ISO 14721:2012, s. 20). Funkční model není pro tuto metodiku podstatný vzhledem k tomu, že se zabývá doporučeními pro producenta. Důležitější je vymezení prostředí, ve kterém se archiv nachází.

1.1.2 Prostředí archivu OAIS

Norma ISO 14721 vymezuje okolí, ve kterém se archiv OAIS nachází a které určuje jeho vstupy a výstupy, zejména z hlediska informací.

Producent (*producer*) je „úloha vykonávaná osobami nebo klientskými systémy poskytujícími informace určené k uchovávání; může se jednat o další archivy OAIS nebo také o osoby či systémy v daném archivu OAIS“ (ISO 14721:2012, s. 25). Producent je podle normy ISO 14721 role spočívající v dodávání informací do archivu OAIS. V praxi masové digitalizace knihoven vykonává úlohu producenta, archivu i managementu často jedna a tatáž knihovna.

Koncový uživatel (*consumer*) je „úloha vykonávaná osobami nebo klientskými systémy, které využívají služeb archivu OAIS za účelem nalezení a vlastního zpřístupnění uchovávaných informací; tuto úlohu mohou vykonávat další archivy OAIS nebo též osoby nebo systémy z daného archivu OAIS.“ (ISO 14721:2012, s. 21).

Management (*management*) je „úloha vykonávaná těmi, kdo určují celková pravidla archivu OAIS jako součást širších pravidel, například v rámci větší organizace“ (ISO

14721:2012, s. 24). Management může například ve vztahu k archivu schvalovat zřizovací listinu, určovat rozsah působnosti, být hlavním zdrojem financování nebo vyhodnocovat výkon.

Z definice normy ISO 14721 vyplývá, že hlavní činností producenta je dodávat informace do archivu. Tato norma však dále vymezuje také některé další dílčí činnosti, ze kterých vyplývá, že producent je často též samotným producentem informací (a samozřejmě sám výraz „producent“ tuto funkci jasně implikuje).

Archiv OAIS podle normy ISO 14721 uzavírá s vkladatelem dohodu o dodávání dat (*submission agreement*), což je „dohoda uzavřená mezi archivem OAIS a producentem, která stanovuje datový model a další potřebná nastavení pro relaci dodávání dat (*data submission session*); datový model určuje formát/obsah a logické konstrukty užívané producentem a způsob, jakým jsou reprezentovány na všech dodaných datových nosičích nebo při všech telekomunikačních spojeních“ (ISO 14721:2012, s. 26). Relace dodávání dat je „jednotlivá dodávka datového nosiče nebo telekomunikační spojení, kterými jsou archivu OAIS poskytována data“ (ISO 14721:2012, s. 22).

Tato dohoda by podle normy ISO 14721 měla vždy v nějaké podobě existovat, ale nemusí jít vždy o formální podobu (smlouvu). Uváděným příkladem je webarchiv (jako typ archivu OAIS, který uchovává sklizený webový obsah), kde dohoda o dodávání dat nabývá podobu nastavení sklízecího robota (ISO 14721:2012, s. 36).

S koncovým uživatelem pak archiv OAIS uzavírá dohodu o objednavce (*order agreement*), což je „dohoda mezi archivem a koncovým uživatelem, v níž jsou stanoveny údaje o dodání, například typ datového nosiče a formát dat“ (ISO 14721:2012, s. 24). Tato dohoda opět nemusí být formální a ustanovení normy lze interpretovat tak, že dohodou o objednavce může být jednoduše to, že uživatel v digitální knihovně (jako součásti archivu OAIS) vyhledá požadovaný dokument.²¹ Rozdíl mezi koncovým uživatelem a cílovou komunitou spočívá v tom, že koncový uživatel je jakýkoliv subjekt, který interaguje s archivem OAIS s cílem získání informací (tedy i softwarový systém). Člen cílové komunity je takový koncový uživatel, na základě jehož znalostní základny se udržují informace tak, aby byly srozumitelné samy o sobě.

¹ Zejména viz ISO 14721:2012, s. 38.

1.1.3 Informační model OAIS

Z hlediska této metodiky je klíčový informační model OAIS, protože ten se vztahuje i na balíčky SIP.

Pojem „informace“ (*information*) definuje norma ISO 14721 jako „jakékoliv znalosti (*knowledge*), které mohou být předmětem výměny (*exchange*)“ a udává, že při výměně jsou informace „vždy vyjádřeny (tj. reprezentovány) určitým typem dat“ (ISO 14721:2012, s. 30). Data jsou definována jako „opakovaně interpretovatelná reprezentace informací ve formalizované podobě vhodné pro komunikaci, interpretaci nebo zpracování.“ (ISO 14721:2012, s. 21). Pojmy informace a data jsou kategorie. Jednotlivý objekt spadající do první kategorie nazývá norma ISO 14721 informační objekt (*information object*), objekt z druhé kategorie pak datový objekt (*data object*).

Při výměně tedy příjemce získává informace vždy z dat, v tom smyslu, že převádí datový objekt na informační objekt. Aby se tento proces mohl uskutečnit, musí příjemce disponovat odpovídající znalostní základnou (*knowledge base*), což je „množina informací, které si osvojila osoba nebo systém a které této osobě nebo tomuto systému umožňuje porozumět přijímaným informacím“ (ISO 14721:2012, s. 23).

Porozumění informací příjemcem je tedy chápáno jako převod datového objektu na informační objekt užitím znalostní základny příjemce. Pokud si například čtenář české národnosti zapůjčí knihu „Information Science in Theory and Practice“ Briana Vickeryho, pak pozorovatelné znaky (vytištěná slova) představují datový objekt. Aby příjemce knize rozuměl, (tj. mohl převést tyto znaky na informace), musí rozumět jazyku, ve kterém je napsaná, tj. angličtině. Pokud tomu tak je, znamená to, že znalost angličtiny tvoří součást jeho znalostní základny. Čtenář, který anglicky neumí, bude potřebovat získat dodatečné informace, aby knize rozuměl. Tento typ informací norma OAIS označuje jako interpretační informace (*representation information*) a definuje je jako „informace, které převádějí datový objekt do smyslupnějších významových celků (*the information that maps a data object into more meaningful concepts*)“ (ISO 14721:2012 s. 25). Tato definice znamená, že srozumitelnost informací lze posuzovat stupňovitě.

Ve výše uvedeném příkladu i čtenář, který anglicky neumí, rozumí tomu, že text je v angličtině; pokud by to nevěděl, mohl by rozumět alespoň tomu, že text je v cizím jazyce; na nejnižším stupni (např. pokud je negramotný) pak může stále rozumět alespoň tomu, že datový objekt je kniha.

Interpretační informace jsou tedy informace, které popisují formu reprezentace informací v datech (v našem příkladu je jí anglický jazyk). Jinými slovy, interpretační informace slouží k doplnění znalostní základny potřebné k rozklíčování informací ve formě datového objektu. Čtenář, který nedisponuje angličtinou, může interpretační informace získat z jiného datového objektu - např. z anglicko-českého slovníku. Náš český překlad „interpretační informace“ není doslovný (ten by byl „informace o reprezentaci“). Domníváme se ale, že lépe vystihuje skutečnost, že veškerá data musejí být předmětem interpretace.

Datový objekt norma ISO 14721 definuje výčtem – datový objekt je „buď fyzický objekt, nebo digitální objekt“ (ISO 14721:2012 s. 21). Fyzický objekt (*physical object*) je „objekt s fyzicky pozorovatelnými vlastnostmi, které reprezentují informace, jež je pro účely uchovávání, šíření a samostatného využívání vhodné patřičně zdokumentovat“, například vzorek horniny, biologického materiálu či archeologický nález (ISO 14721:2012 s. 24). Digitální objekt (*digital object*) je „objekt složený z množiny bitových posloupností (*a set of bit sequences*)“ (ISO 14721:2012 s. 22). Informační objekt (*information object*) je pak „datový objekt se svými interpretačními informacemi“ (ISO 14721:2012 s. 23).

Norma ISO 14721 nedefinuje slovní spojení „digitální informace“ (*digital information*), ačkoliv jej často užívá. Z jejího textu lze vyvodit, že digitální informace jsou informační objekty, které vznikly interpretací digitálních objektů. Pro označení informací vzniklých interpretací fyzických objektů užívá norma termín „nedigitální“ (*non-digital*).²

Informační model přináší typologii informačních objektů, které se vyskytují v průběhu životního cyklu informací. Z nich musíme vyčlenit jako klíčový informační objekt ten, který norma ISO 14721 nazývá informační obsah a který je podle normy hlavním předmětem dlouhodobého uchovávání v archivu OAIS; dále pak zmíněné interpretační informace. Kategorie informací, které slouží jako nezbytné dodatečné informace, jež musejí být uchovávány spolu s informačním obsahem, nazývá norma archivační informace. Ty obsahují pět podkategorií informací. Dalším důležitým prvkem informačního modelu je informační balíček jako logická schránka, která zpravidla obsahuje informační obsah a archivační informace.

² Již tyto úvodní definice ukazují, že pojetí informací v normě ISO 14721 se liší od řady jiných modelů. Například podle modelu projektu InterPARES 2 jsou informace (*information*) „sdělení zamýšlené pro komunikaci napříč časem a prostorem“, data (*data*) jsou „nejmenší smysluplná část informace“ a dokument (*document*) je „zaznamenaná informace“ (Duranti a Thibodeau, 2006, s. 15). Data jsou tedy podle modelu podkategorií informací (zatímco v normě ISO 14721 jsou „data“ a „informace“ dvě odlišné kategorie stejné úrovně) a definice termínu „dokument“ odpovídá „datovému objektu“ v normě OAIS.

1.1.3.1 Informační obsah a interpretační informace

Informační obsah (*content information*) je podle normy ISO 14721 „množina informací, která je původním předmětem uchování, nebo která obsahuje část těchto informací či všechny tyto informace; informační obsah je informační objekt složený ze svého datového objektu s obsahem a svých interpretačních informací“ (ISO 14721:2012, s. 21).

Oproti první verzi normy ISO 14721 nynější definice informačního obsahu připouští, že některé části informačního obsahu nemusejí nebo nemohou být zachovány (příklad z prvního vydání normy ISO 14721: Máme tabulku, obsahující čísla, reprezentující teplotní údaje; respektive, obsahující údaje, které pozorovatel vykládá jako záznam teplot, bez toho, aby k tabulce měl dokumentaci, osvětlující původ, souvislosti s jiným měřením či historií dat).

Datový objekt s obsahem (dále zkráceně jako objekt CDO, *content data object*) tedy nese informační obsah, pokud je příjemcem spojen s odpovídajícími interpretačními informacemi (které mohou, ale nemusejí být součástí jeho znalostní základny). Je třeba upozornit, že digitálním objektem podle normy ISO 14721 není nutně soubor – může jím být, stejně jako jím může být například množina souborů. V případě objektu CDO může jít o jeden soubor ve formátu PDF, který reprezentuje digitální knihu, nebo o stovku souborů ve formátu JPEG, které dohromady reprezentují sto stran knihy – záleží na konkrétní implementaci. Jednotlivý digitální objekt CDO v archivu je tedy taková podmnožina bitových posloupností z celkové množiny všech bitů uložených na datových nosičích archivu, které jsou potřebné k transformaci bitů do daného konkrétního informačního obsahu.

Interpretační informace člení norma ISO 14721 na tři druhy:

Strukturální interpretační informace (*structure information*) jsou takové „interpretační informace, které udávají, jak jsou další informace složeny; mohou například převádět bitové toky (*bit streams*) na základní typy dat, jako jsou znaky, čísla a pixely, a na seskupení těchto typů dat, jako znakové řetězce a pole“ (ISO 14721:2012, s. 26).

Sémantické interpretační informace (*semantic information*) jsou takové „interpretační informace, které podrobněji popisují význam nesený strukturálními interpretačními informacemi“ (ISO 14721:2012, s. 26).

Třetí druh interpretačních informací zavedlo až druhé vydání normy – jde o **ostatní**

interpretační informace (*other representation information*), které jsou definovány jako takové „interpretační informace, které nelze jednoduše zařadit mezi strukturální interpretační informace nebo sémantické interpretační informace. Například k porozumění datovému objektu s obsahem mohou být potřeba software, algoritmy, šifrování, psané pokyny atd., přičemž všechny tyto informace budou odpovídat definici interpretačních informací. Nemusí být zřejmé, zda se vztahují ke struktuře nebo významu interpretované informace. Dále se mezi ně mohou řadit informace udávající vztah mezi strukturálními interpretačními informacemi a sémantickými interpretačními informacemi nebo popisující software potřebný pro zpracování databázového souboru“ (ISO 14721:2012, s. 24).

Pokud se vrátíme k případu uvedenému výše, interpretační informace potřebné k převedení fyzického objektu ve formě tištěné knihy na informační obsah (tj. člověku srozumitelný intelektuální obsah) představují česko-anglický slovník a schopnost čtenáře rozeznat lineární text. Konkrétněji jde v případě slovníku o sémantické interpretační informace ve formě dalšího datového objektu a ve druhém případě o strukturální interpretační informace.

Předmětem této metodiky je především otázka interpretace digitálních objektů. Lze říci, že problematika interpretačních informací tvoří jádro digitální archivace. Digitální objekt musí být v daný okamžik vždy uložen na nějakém datovém nosiči, který je sice fyzickým objektem (pevný disk, magnetický pásek apod.), ale data na něm zapsaná člověk nedokáže vnímat, natož převést do srozumitelné podoby. Z tohoto hlediska je nutné, aby interpretaci digitálního objektu zprostředkovala počítačová technologie.

V praxi digitální archivace se klade největší důraz na strukturální interpretační informace, zejména na formát souboru. Formáty patrně nejviditelněji odrážejí problém zastarávání digitálních technologií.

Z hlediska své činnosti musí archiv OAIS shromáždit dostatečné interpretační informace k objektu CDO. Způsoby konkrétního řešení se mohou lišit podle posouzení situace archivem. V případě formátu to může znamenat získat formátovou specifikaci a uložit ji spolu s objektem CDO nebo jen zaznamenat technické informace o formátu, na základě kterých je možno tyto interpretační informace dohledat (tj. předpoklad, že jsou tyto informace běžně dostupné, a tedy samy o sobě srozumitelné pro cílovou komunitu). Tyto technické informace jsou v praxi zaznamenávány do datových objektů, které se označují jako technická metadata. Z informačního modelu OAIS je zřejmé, že i k tomuto typu metadat (tak jako k jakýmkoliv jiným typům dat), je nutné získat dostatečné interpretační informace. Tento

proces podle normy ISO 14721 v digitálním světě znamená rekurzivitu – interpretační informace o technických metadatech budou zaznamenány v dalším digitálním objektu (např. v PDF popisujícím metadatovým standard). Tato rekurzivita podle normy končí tehdy, když jsou interpretační informace zaznamenány v podobě fyzického objektu (tedy např. vytištěný standard). V praxi to v současnosti není považováno za větší problém vzhledem k tomu, že metadata jsou zpravidla ukládána v textových formátech (zejména v XML), jejichž interpretace se nepovažuje za problematickou.

Norma ISO 14721 uvádí dva typy specializovaných softwarových nástrojů pro interpretační informace (tyto nástroje samy patří mezi ostatní interpretační informace). **Software pro zobrazení interpretačních informací** (*representation rendering software*)³ je software, který umožňuje interpretační informace reprodukovat v podobě srozumitelné lidem. Příkladem je prohlížeč Adobe Acrobat, který dokáže zobrazit formátovou specifikaci uloženou v PDF. **Zpřístupňovací software** (*access software*) je software, který dokáže prezentovat samotný informační obsah nebo jeho část. Jinými slovy, software pro zobrazení interpretačních informací reprodukuje interpretační informace a slouží jako pomůcka pro uchovávání informačního obsahu, zatímco zpřístupňovací software reprodukuje informační obsah a slouží jako prostředek pro zpřístupňování informačního obsahu cílové komunitě.

Shrnutí: Norma ISO 14721 nepopisuje odlišnosti užití různých typů informačního obsahu. Archiv OAIS musí uchovávat objekt CDO spolu se zvolenou množinou interpretačních informací (uchované v podobě dalšího digitálního objektu, zpravidla označovaného jako “metadata”). Některé typy informačního obsahu mohou být určeny k reprodukci lidským uživatelům, jako tomu bývá v případě dokumentů digitálního dědictví (např. zobrazení v digitální knihovně), zatímco jiné nemusí být vůbec určeny pro vnímání člověkem, ale pro zpracování jinými systémy (to je příklad rozsáhlých datových sad získaných z pozorování, které mohou být reprezentovány v podobě tabulek, obsahovat numerické údaje a sloužit jako zdroj pro automatizované analýzy).

Za jeden ze způsobů, jak řešit otázku interpretačních informací, bylo určeno vytvoření mezinárodních registrů interpretačních informací (zejména pro formáty). Otázce těchto registrů a bližšímu přiblížení problematiky formátů se věnuje kapitola 1.2.4.

Je třeba rozlišovat interpretační informace a metadata, která je reprezentují. Metadata

³ Anglický název je zavádějící, měl by být spíše „representation information rendering software“.

jsou datové objekty a potřebují rovněž své interpretační informace. Rekurzivitu v praxi ukončuje užití souborů v textových formátech typu XML, jejichž interpretace se z hlediska digitální archivace považuje (vzhledem k jednoduchosti souborového formátu) za relativně bezproblémovou.

1.1.3.2 Archivační informace

Informační model OAIS definuje archivační informace (*preservation description information*) jako „informace, které jsou nezbytné k adekvátnímu uchovávání jednotlivého informačního obsahu“ (ISO 16363:2012, s. 25). Jde o kategorii informací, která obsahuje pět podkategorií.

Identifikační informace (*reference information*) jsou „informace, které plní funkci identifikátoru informačního obsahu“ (ISO 14721:2012, s. 25). Tyto informace mohou také zahrnovat „identifikátory, které vnějším systémům umožňují jednoznačně odkazovat na konkrétní informační obsah“ a udávat a popisovat způsob jejich přidělování (ISO 14721:2012, s. 74). Uváděným příkladem z oblasti knihoven je perzistentní identifikátor (konkrétním příkladem pak ISBN) a bibliografický popis. Tyto informace tedy mohou zahrnovat i globální perzistentní identifikátory, na základě nichž mohou uživatelé vyhledávat konkrétní informační obsah.

Provenienční informace (*provenance information*) jsou „informace, které dokumentují historii informačního obsahu; tyto informace vypovídají o původu nebo zdroji informačního obsahu, o veškerých změnách, které mohly od doby jeho vzniku nastat, a o tom, kdo o něj od doby jeho vzniku pečoval“ (ISO 14721:2012, s. 25). Jak dále norma uvádí, archiv nese odpovědnost za vytváření a uchovávání těchto informací až od okamžiku jejich příjmu do archivu; provenienční informace z dřívější doby by měl poskytnout vkladatel. Uváděnými příklady z oblasti knihoven jsou metadata o procesu uchovávání (ukazatele k předchozím verzím jednotky, historie změn) a metadata o procesu digitalizace.

Informace o neporušenosti (*fixity information*) jsou „informace, které udávají, jak je zajištěno, aby objekt s informačním obsahem nebyl nezdokumentovaným způsobem změněn“ (ISO 14721:2012, s. 22). Tyto informace „fungují jako obal nebo ochranný štít, který chrání informační obsah“ (ISO 14721:2012, s. 34). Příkladem je digitální otisk.

Kontextuální informace (*context information*) jsou „informace, které dokládají vztah informačního obsahu k jeho okolí; patří mezi ně důvod vytvoření informačního obsahu a jeho

vztah k dalším objektům s informačním obsahem“ (ISO 14721:2012, s. 21).

Informace o přístupových právech (*access rights information*) jsou „informace, které udávají omezení týkající se přístupu k informačnímu obsahu, a to včetně právního rámce, licenčních podmínek a řízení přístupu“ (ISO 14721:2012, s. 19).

1.1.3.3 Informační balíček

Informační balíček je „logická schránka, která může obsahovat informační obsah a archivační informace; k tomuto informačnímu balíčku jsou připojeny informace o zabalení (*packaging information*), které vymezují a určují informační obsah, a informace o popisu balíčku (*package description*), které usnadňují vyhledání informačního obsahu“ (ISO 14721:2012, s. 23).

Norma dále odlišuje tři typy informačních balíčků:

- **Archivní informační balíček** (*archival information package*) je „informační balíček, který je složen z informačního obsahu a přidružených archivačních informací a je uchováván v archivu OAIS“ (ISO 14721:2012, s. 23). Dále bude nazýván jen jako balíček AIP.
- **Vstupní informační balíček** (*submission information package*) je „informační balíček, který dodává producent do archivu OAIS tak, aby mohl být využit při sestavení nebo aktualizaci jednoho nebo více AIP a/nebo přidružených popisných informací (*descriptive information*)“ (ISO 14721:2012, s. 26). Dále bude nazýván jako balíček SIP.
- **Výstupní informační balíček** (*dissemination information package*) je „informační balíček odvozený z jednoho nebo více balíčků AIP a zasláný archivem OAIS koncovému uživateli jako odpověď na jeho požadavek vůči tomuto archivu“ (ISO 14721:2012, s. 22). Dále bude nazýván jako balíček DIP.

Klíčový informační balíček z hlediska archivu OAIS je balíček AIP. Ten musí obsahovat informační obsah a přidružené archivační informace. Definice nadřazené kategorie (tj. informačního balíčku) uvádí, že tyto dvě složky může obsahovat. To znamená, že při dodávání balíčků SIP do archivu nebo vydávání balíčků DIP koncovým uživatelům nemusí každý jednotlivý balíček SIP nebo DIP vždy obsahovat informační obsah a archivační informace. Do archivu mohou být například dodávány odděleně balíčky SIP s informačním obsahem a balíčky SIP s archivačními informacemi. Podobně mohou být archivem vydávány,

v závislosti na požadavcích koncových uživatelů, například jen balíčky DIP obsahující informační obsah, ale již nikoliv archivační informace.

Informace o zabalení (*packaging information*) jsou „informace, které slouží k propojení a popisu součástí informačního balíčku“ (ISO 14721:2012, s. 24). Popis balíčku (*package description*) jsou „informace určené pomůckám pro zpřístupnění“ (ISO 14721:2012, s. 24). Těmito pomůckami norma míní „softwarový program nebo dokument, který koncovým uživatelům umožňuje najít, analyzovat, objednat nebo získat informace z archivu OAIS“ (ISO 14721:2012, s. 19). Popisné informace (*descriptive information*) jsou „množina informací, která je složena především z popisů balíčků a je poskytována správě dat za účelem podpory koncových uživatelů při objednávání a získávání informačních jednotek z archivu OAIS“ (ISO 14721:2012, s. 22).

Hlavním účelem informací o zabalení je vymezit, které části balíčku AIP jsou informační obsah a které archivační informace. To znamená popsat jednak to, které soubory v balíčku tvoří objekt CDO, jenž reprezentuje informační obsah (např. pět souborů ve formátu JPEG obsahujících články) a které technická metadata reprezentující interpretační informace (např. jeden soubor v XML popisující formát JPEG). **Objekt CDO a technická metadata dohromady tvoří informační obsah.** Dále je třeba popsat, jaké soubory jsou archivační metadata reprezentující archivační informace (např. druhý soubor v XML obsahující bibliografický popis). V popisu balíčků a v popisných informacích se pak podle normy obvykle opakují identifikační informace. Pokud budeme uvažovat digitální knihu v PDF, pak informační obsah může tvořit jeden soubor ve formátu PDF (objekt CDO) a jeden soubor ve formátu XML obsahující technická metadata (např. informace o verzi a typu, příkladem je PDF/A-2u); archivační informace jeden soubor v XML s archivačními metadaty (blíže o metadatach bude pojednávat Kapitola 1.3), přičemž jejich součástí bude identifikátor ISBN (který bude současně obsažen i v samotném informačním obsahu, jelikož je vydavatelskou praxí uvádět jej přímo v knize) a tento identifikátor bude také obsažen v popisu balíčku a popisných informacích. Na základě tohoto identifikátoru pak čtenář jakožto koncový uživatel může knihu vyhledat v archivu.

Předmětem dlouhodobého uchování je informační obsah (např. konkrétní kniha v PDF) i přidružené archivační informace uložené v balíčku AIP spolu s informačním obsahem. Archiv vytváří balíček AIP z balíčku SIP (nebo z více balíčků SIP), který mu dodá vkladatel. Na základě dohody o dodávání dat oba subjekty (archiv a vkladatel) definují datový model, tedy podobu dodávaných informací, jehož součástí by měl být datový slovník (popis

všech typů dat, což zahrnuje i metadata). Osvědčeným postupem v komunitě digitální archivace je dohoda o užití takových formátů objektu CDO, které podporují dlouhodobé uchovávání, a užití mezinárodně rozšířených metadatových schémat.

V rámci digitalizace NDK kromě tří výše zmíněných typů balíčků rozeznáváme tzv. **produkční balíček** (*Producer submission package*) (Hutař, 2012, s.70.). Tento druh balíčku není součástí modelu OAIS, ale byl dodatečně vytvořen z praktických důvodů mezinárodní pracovní skupinou z prostředí LTP systémů. Tato pracovní skupina stávající model OAIS archivu rozšířila v oblasti modulu Ingest⁴, který původně rovnou pracoval s již hotovým balíčkem SIP, ovšem nijak neřešil proces, při kterém se data transformují do podoby, která splňuje archivační požadavky. V realitě LTP systémů dochází k úpravám surových dat podle požadavků příslušné archivní instituce do archivem požadovaných souborových formátů a struktury. PSP balík tedy představuje chybějící krok v procesu před vznikem SIP balíku, se kterým se dále v konceptu archivu OAIS pracuje. Toto předpole pro modul Ingest, ve kterém se odehrává změna z PSP balíku na balík SIP, se nazývá Pre-Ingest.

1.1.4 Přístupy k uchovávání

Norma ISO 14721 popisuje dvě hlavní kategorie opatření pro uchovávání informačního obsahu vykonávaná v archivu OAIS nad balíčkem AIP jako digitální migraci a emulaci. Emulaci se blíže věnovat nebudeme (a norma se jí také detailněji nezabývá).

Digitální migrace (*digital migration*) je podle normy „přesun (*transfer*) digitálních informací v rámci archivu OAIS se záměrem tyto informace uchovat“ (ISO 14721:2012, s. 22). Podle normy tento přesun od jiných typů přesunů odlišují následující tři atributy (ISO 14721:2012, s. 22):

- je zaměřen na uchování celého informačního obsahu,
- nová podoba informací v archivu nahrazuje podobu předchozí
- a archiv OAIS řídí všechny stránky přesunu a nese za ně plnou odpovědnost.

Norma ISO 14721 rozlišuje čtyři typy digitálních migrací. Renovace (*refreshment*) a replikace (*replication*) jsou migrace, při kterých nedochází ke změně datových objektů (tedy bitových posloupností), balíčkovací migrace (*repackaging*) a transformace (*transformation*) jsou migrace, při nichž ke změně datových objektů dochází. Definice těchto

⁴ Jedná se o část procesu, během které se data přijímají do úložiště LTP.

typů jsou následující (ISO 14721:2012, s. 103-104):

Renovace: Digitální migrace, při níž je jedna instance datového nosiče, která obsahuje balíček AIP, více balíčků AIP nebo části balíčků AIP, nahrazena jinou instancí datového nosiče stejného typu, a to zkopírováním bitů na datový nosič, využitý k umístění balíčků AIP a ke správě a přístupu k datovému nosiči. Díky tomu dokáže stávající mapovací infrastruktura archivního úložiště beze změny stále nalézat balíček AIP a přistupovat k němu.

Replikace: Digitální migrace, při níž nedochází k žádným změnám balíčkovacích informací, informačního obsahu ani archivačních informací. Bity nesoucí tyto informační objekty jsou při přesunu na novou instanci stejného nebo nového typu datového nosiče zachovány. Renovace je také replikací. Replikace však může vyžadovat změny mapovací infrastruktury archivního úložiště.

Balíčkovací migrace: Digitální migrace, při níž dochází k změně bitů balíčkovacích informací.

Transformace: Digitální migrace, při níž dochází k změnám bitů informačního obsahu nebo archivačních informací, přičemž je současně vyvinuta snaha uchovat informační obsah v úplnosti.

Renovaci a replikaci se nebudeme dále věnovat. V praxi se souhrnně označuje jako **bitová ochrana** a znamená zkopírování dat na jiný datový nosič. Rozdíl spočívá v tom, že replikace může vyžadovat úpravu mapovací infrastruktury (tato změna postihuje hardware a software, nikoliv data). Norma ISO 14721 tento rozdíl vysvětluje nepříliš srozumitelně, ale příklad z běžné praxe pomůže lepšímu porozumění. Zkopírování dat z jednoho nosiče CD na jiný je příkladem renovace. Příkladem replikace je zkopírování dat z deseti nosičů CD na pevný disk počítače.

Žádný z těchto dvou typů migrací nepředstavuje větší intelektuální problém, riziko spočívá především v nedostatku financí a nedostatečné kontrole stavu datových nosičů. Archiv musí mít nastaveny mechanismy pro monitorování stavu datových nosičů a plány pro migraci v případě, že je zaznamenáno riziko jejich degradace nebo zastarávání (tj. konec jejich hardwarové podpory). Rizikem samozřejmě může být i sama volba datových nosičů, které jsou zcela nevhodné svou povahou (názorným příkladem z minulosti je disketa, která nikdy nebyla bezpečným nosičem) nebo náročností údržby (např. z hlediska nedostatku finančních prostředků).

Ani balíčkovací migrace nepředstavuje větší riziko. Jedná se o změnu informací

o zabalení. Příkladem může být jiné uspořádání v rámci adresáře. Riziko, které přináší balíčkovací migrace, spočívá spíše v neudržení odlišení informačního obsahu a archivačních informací.

Transformace je nejdůležitějším a nejsložitějším typem digitální migrace. Pouze transformace podle normy zakládá novou verzi balíčku AIP. To znamená, že verze balíčku AIP, která prošla renovací, replikací nebo balíčkovací migrací, zůstává nezměněna.

Původně uložený balíček AIP (tj. informační balíček vytvořený z balíčku SIP dodaného vkladatelem) má být považován za první verzi balíčku AIP a norma jej označuje jako originál, resp. původní balíček AIP (*original AIP*). Tento původní balíček „může být udržován pro ověření uchovávání informací“ (ISO 14721:2012, s. 105). Norma takovýto postup (tj. uchovávat první verzi balíčku AIP i tehdy, když je vytvořena novější verze) tedy nepředepisuje, pouze uvádí jako možnost.

Transformaci norma dále dělí na dva dílčí typy. **Vratná transformace** (*reversible transformation*) je taková transformace, kdy je možná bezeztrátová zpětná transformace. Uváděným příkladem vratné transformace je migrace textového souboru (obsahujícího písmena anglické abecedy) v kódování ASCII do kódování UNICODE UTF-16 (ISO 14721:2012, s. 106).

Nevratná transformace (*non-reversible transformation*) je taková „transformace, u které nemůže být zaručeno, že se jedná o vratný převod“ (ISO 14721:2012, s. 24). Uváděný popis může být nesrozumitelný, proto uveďme jednoduchý příklad z praxe – formátovou konverzi z formátu TIFF (nekomprimovaná verze) do formátu JPEG (formát JPEG nepodporuje jinou než matematicky ztrátovou kompresi). Nevratná transformace je podle normy rizikem pro zachování autenticity. V případě, že jde o transformaci archivačních informací, je možné učinit zobecnění, že v tomto případě by nikdy nemělo jít o nevratnou transformaci.

Norma systematicky nepopisuje rozdíly mezi transformací, která mění informační obsah, a transformací, která mění archivační informace. Její koncepty však lze specifikovat následujícím způsobem. V případě transformace informačního obsahu mohou nastat dvě základní varianty. Zprvu může být potřeba změnit pouze interpretační informace o formátu. Například může dojít k tomu, že archivem doporučovaný nástroj pro zobrazení souboru ve formátu PDF zastaral a nebude adekvátně reprodukovat formát PDF. Pak bude muset archiv provést průzkum a vybrat jiný nástroj, a informace o tomto nově doporučeném

nástroji budou zdokumentovány a uchovávány spolu s objektem CDO. Tato transformace tedy zahrnuje pouze interpretační informace. Za druhé může nastat situace, kdy dochází ke změně bitů objektu CDO, typicky v případě formátové konverze. Pokud půjde o konverzi z formátu TIFF (nekomprimovaná verze) do formátu JP2 (matematicky bezeztrátová komprese), půjde o vratnou transformaci, nicméně nová reprezentace (objekt CDO ve formátu JP2) bude vyžadovat jiné interpretační informace, neboť se jedná o jiný formát.

Jinou situací je proces obohacování archivačních metadat. V případě, že se mění objekt CDO, pak se vždy také mění archivační informace v tom smyslu, že se doplní o záznamy těchto změn (zejména jde o provenienční informace). Proto bude v průběhu uchovávání objem archivačních informací narůstat. Rovněž může dojít k tomu, že se v průběhu uchovávání obohatí archivační informace o nové informace (např. další perzistentní identifikátor).

1.1.5 Specifická standardizace informačních balíčků

Pro specifické typy digitálních dokumentů existují specifické metadatové standardy a doporučení pro výběr vhodných formátů. Datový objekt s obsahem může být tvořen jedním nebo více soubory, které mohou být uloženy v jednom nebo více formátech (souborových formátech). Formát je typem strukturálních interpretačních informací.

1.2 Formát objektu CDO

Jedna z definic pojmu formát je “Vnitřní struktura a kódování digitálního objektu, která umožňuje jeho zpracování nebo zobrazení ve formě přístupné lidskému uživateli. Digitální objekt může být soubor, nebo datový tok do souboru vložený” (Brown, 2006, s. 5). Americká studie o formátech souborů uvádí: „Většina souborů – s výjimkou souborů, které jsou jednoduchými datovými toky – obsahuje dvě základní komponenty: strukturální prvky a datové prvky. Formát souboru reprezentuje jedinečné a specifické uspořádání těchto strukturálních a datových prvků“ (Lawrence, 2000, s. 2).

Formát je jedním z typů interpretačních informací modelu OAIS, konkrétně informací o způsobu, jakým je potřeba datový objekt uložený v konkrétním formátu interpretovat (softwarovou aplikací), a v případě digitalizátů knih také o způsobu, jak datový objekt v daném formátu reprodukovat, tj. adekvátně zobrazit. Digitalizát tištěného dokumentu lze pojímat jako objekt CDO tvořený dvěma datovými komponentami: obrazová komponenta

(archivní kopie) a OCR komponenta. V případě zvukových dokumentů je kromě obrazové (která je v tomto případě jenom doplňující) zastoupena také komponenta zvuková. Třetí komponentou je strukturální komponenta, popsána strukturální mapou v METS, obecně označovaná jako „strukturální metadata“.

V kontextu digitální archivace by v ideálním případě:

a) měl vždy být k dispozici dostatek volně dostupných (specializovaných) softwarových aplikací, které znalostí daného formátu disponují (tj. aplikace, které mají dostatečné strukturální interpretační informace) a dokáží s ním pracovat;

b) informace o formátu by měly být obsaženy v dostatečně dobře popsané dokumentaci (formátové specifikaci), která by měla být dostupná;

c) formát by neměl být zatížen patenty.

Počítačová realita má však k tomuto ideálu daleko a právě formáty a aplikace pro práci s nimi jsou předmětem častých změn, a tedy rizikem pro dlouhodobé uchovávání. Jeden z odhadů průměrné délky zastarání formátu (od doby uvedení) na trh se pohybuje v rozmezí osmi až dvaceti let (Kejser, Nielsen, Thirifays, 2011). Zastarávání se projevuje dvěma způsoby. Jednak narůstajícím rizikem ztráty dostupnosti softwarových aplikací, které dokáží formáty adekvátně reprodukovat (tj. zobrazit, přehrát nebo jiným způsobem prezentovat smyslům lidského uživatele). Zadruhé pak ztrátou schopnosti nové generace softwarových aplikací formát zpracovávat (upravovat data v daném formátu, provádět konverzi do jiného formátu apod.). V řadě případů také není dostupná formátová specifikace nebo se k formátům váží licenční omezení, což rovněž představuje velká rizika pro uchovávání informací. Připomeňme si druhou povinnost archivu OAIS (získat možnost s informacemi dostatečně nakládat), kterou lze zřejmě aplikovat i na tento případ, neboť nad informacemi uloženými ve formátu, jehož specifikace je nedostupná (uzavřená) nebo zatížená patenty, z principu nelze získat plnou kontrolu.

Pokud nedojde k včasné formátové konverzi v době, kdy ještě existují vhodné nástroje pro tuto transformaci, může dojít k nevratné ztrátě informačního obsahu (objekt CDO může existovat jako uložený objekt, ale nebude z něj možné získat informační obsah). Archiv by mohl teoreticky vytvořit nový nástroj na základě uložené formátové specifikace zastaralého formátu, prakticky je však takovou alternativu obtížné ověřit.⁵

⁵ Vzhledem k tomu, že je obtížné předvídat vlastnosti budoucího technologického prostředí, které mohou vytvoření takového nástroje znemožňovat, nebo vzhledem k tomu, že již nebudou dostupné další interpretační informace, které formátová

1.2.1 Roviny a aspekty užití formátu

Z hlediska popisu formátu musíme odlišovat několik rovin, přičemž v případě rastrových formátů jde zejména o tyto roviny:

- Rodina formátů
- Konkrétní formát
- Verze formátu
- Komprese
- Profil

Příkladem rodiny formátů je RAW. Existuje celá řada konkrétních formátů RAW této rodiny, které vytvořili výrobci zařízení. Například fotoaparáty Canon užívají formát Canon RAW. Verze je dána historicky a odlišné verze formátu mohou znamenat odlišné požadavky na zobrazení. Různé formáty nabízejí odlišné možnosti komprese (např. TIFF verze 6 může být v nekomprimované podobě, zatímco JP2 je vždy komprimovaný). Profil znamená specifické nastavení v rámci formátu při jeho vytváření (např. volba ztrátové nebo bezztrátové komprese), u některých formátů nastavení profilu vyžaduje specializované znalosti (což je zejména případ formátu JP2).

Z hlediska užití musíme u rastrových formátů odlišovat nejméně tyto aspekty:

a) archivační formát; b) prezentační formát; c) prezentační meziformát.

Archivační formát je takový, který je aktuálně vhodný z hlediska potřeb dlouhodobého uchování. Někdy je nazýván jako archivní obrazová matrice (*archival master*). Koncept obrazové matrice (*master*) byl do digitální archivace převzat z komerčního sektoru. Obrazová matrice jsou v obrazovém průmyslu obvykle komprimovaná data, která slouží jako zdroj pro vytváření obrazových dat v různé kvalitě, v různých formátech a pro různé účely a nabízí nejvyšší možnou kvalitu dané produkce (např. fotografa). Archivační formát je volen mj. právě s ohledem na to, aby měl tuto funkci obrazové matrice, přičemž však jsou na jeho výběr kladeny další omezující podmínky.

Prezentační formát je takový, který je aktuálně vhodný pro zpřístupňování z hlediska potřeb cílové komunity, v kontextu současné praxe formou prezentace v digitální knihovně.

specifikace předpokládala, ale nezaznamenala, neboť v době vytvoření formátové specifikace byly tyto další interpretační informace běžně dostupné.

Prezentační formát můžeme rozdělit na hlavní prezentační formát (formát představující nejvyšší možnou kvalitu) a doplňkové prezentační formáty (např. malé náhledy obrázků).

Prezentační meziformát je meziformát, ze kterého digitální knihovna generuje cílový prezentační formát. Nejčastějším případem současné praxe je formát JP2 ve ztrátové kompresi, ze kterého se v digitální knihovně generuje formát JPEG jako výsledný hlavní prezentační formát.

1.2.2 Výběr archivačního formátu

Široce známým doporučeným postupem pro řízení životního cyklu digitálních dokumentů je, že prvním a zásadním krokem je výběr vhodného archivačního formátu, tj. formátu, ve kterém budou dokumenty uchovávány v archivu (a v případě, kdy dokumenty vytváří i uchovává jedna a tatáž organizace, vytváření finálních dat přímo v tomto formátu). Archivační formát je volen z hlediska jeho (aktuální) vhodnosti pro uložení v balíčku AIP v digitálním archivu. Cílem je uložit obsah v takovém formátu, o kterém se předpokládá, že jeho užití v současnosti a blízké budoucnosti nebude představovat větší riziko. Pokud vkladatel do archivu dodá data v jiném formátu než archivačním, je doporučeným postupem, aby archiv provedl normalizaci do archivačního formátu (Cubr, 2010, s. 83-86).

Volba archivačního formátu není jednoduchá záležitost, nicméně dnes již existuje řada doporučení od uznávaných organizací, kterými se lze řídit. Za jeden z nejvýznamnějších zdrojů pro výběr archivačního formátu lze v současnosti označit dokument Deklarace doporučených formátů (Recommended Formats Statement), který od roku 2014 vydává Kongresová knihovna (Library of Congress). Tento dokument obsahuje konkrétní seznam doporučených formátů (pro digitální i fyzické objekty), který je určen jak pro vnitřní potřeby Kongresové knihovny, tak i pro jiné archivy. Zahrnuje seznamy preferovaných a akceptovaných formátů. Dokument je každý rok aktualizován, poslední seznam byl vydán v roce 2022 (LOC, 2022). Formátová doporučení nalezneme i ve směrnících, které vydává Iniciativa amerických federálních agentur pro zásady digitalizace (Federal Agencies Digital Guidelines Initiative - FADGI).⁶

Pro představu o rozšíření formátů v paměťových institucích je užitečné i studium tzv. formátových politik lokálních institucí. Řada z nich nese v názvu slovo „recommendation“, ale od výše uvedených doporučení (LOC a FADGI) se liší tím, že uvádějí

⁶ V Evropě seznam vlastních formátových doporučení spravuje například Katalog archivischer Dateiformate Švýcarského spolkového archivu, viz https://kost-ceco.ch/cms/kad_main_de.html.

formáty, které tato instituce přijímá k uložení, zatímco ty předchozí doporučují formáty pro archivaci obecně. První formátovou politikou citovanou v odborné literatuře je dokument Floridského digitálního archivu (Florida Center for Library Automation, 2012). Tento dokument obsahuje výčet nejběžnějších formátů s hodnocením jejich spolehlivosti na třístupňové škále (vysoká, střední a nízká spolehlivost) pro potřeby uchovávání v tomto archivu. V současnosti zveřejňuje své formátové politiky několik desítek institucí.

Řada institucí zveřejňuje též formátová hodnocení, jejichž příkladem může být studie americké iniciativy FADGI, která byla vydána v roce 2014 a zaměřuje se specificky na rastrové formáty (tedy formáty relevantní pro digitalizaci knih) (FADGI, 2014). Studie obsahuje zhodnocení nejběžnějších formátů (TIFF, JP2, PDF, PNG, JPEG) z hlediska digitální archivace a uvádí sadu podrobných srovnávacích kritérií seskupených do čtyř hlavních kategorií (udržitelnost, ekonomické faktory, požadavky na implementaci, nastavení a možnosti).

Kritéria, která se v různých doporučeních objevují nejčastěji, lze zobecnit do těchto kategorií: podpora, otevřenost a nezatíženost patenty.⁷ Podporou formátu se zde rozumí míra jeho užívanosti v dané komunitě a dostupnost nástrojů pro vytváření, zpracování a reprodukci. Nízká úroveň podpory znamená, že formát zastarává nebo že se vůbec neujal.⁸ Otevřeností formátu se rozumí skutečnost, že je dostupná dokumentace formátu (tj. formátová specifikace). Směrnice MINERVA uvádějí, že užití otevřených formátů „napomůže interoperabilitě a zajistí, že zdroje lze opětovně využívat, vytvářet a upravovat celou řadou aplikací. Rovněž napomůže vyhnout se závislosti na konkrétním dodavateli“ (Ferne, 2008, s. 32).

Za dostupnost se v praxi zpravidla považuje to, že formátová specifikace je buď volně dostupná online, nebo že je dostupná v nějaké normalizační organizaci, která ji udržuje. Například specifikace formátu TIFF je dostupná na webu jeho vlastníka (firmy Adobe)⁹ a je průmyslovým standardem. Specifikace formátu JP2 je mezinárodní normou, která je dostupná k zakoupení v normalizační organizaci ISO.¹⁰

Nezatíženost patenty nutně neznamená, že formát není nikým vlastněn, pouze to, že

⁷ Srv. např. Cubr, 2010, s. 83; Ferne, 2008, s. 11-14.

⁸ V tomto smyslu může být přínosný například pravidelně aktualizovaný seznam „Ohrožených digitálních druhů“ (*The Bit List of Digitally Endangered Species*), spravovaný koalici DPC, viz <https://www.dpconline.org/digipres/champion-digital-preservation/bit-list>.

⁹ <https://developer.adobe.com/content/dam/udp/en/open/standards/tiff/TIFF6.pdf>

¹⁰ Citace specifikace: ISO/IEC 15444-1:2004. Information technology - JPEG 2000 image coding system: core coding system. 2nd ed. Geneva: ISO, 2004.

výkon práv duševního vlastnictví není uplatňován.

Dalšími uváděnými kritérii pro výběr formátu jsou například: míra zpětné kompatibility; možnosti exportu do jiných formátů; míra nezávislosti na specifických hardwarových a softwarových platformách; rozumná rovnováha mezi nabídkou funkcí formátu na straně jedné a přiměřenou komplexitou na straně druhé (Cubr, 2010, s. 85).

Obecně lze říci, že archivační formát by měl splňovat mj. také funkce obrazové matrice. Mimo oblast digitální archivace je častou volbou (zejména profesionálních fotografů) uložení obrazové matrice ve formátu DNG. Takovéto užití je v digitální archivaci problematické jednak kvůli svázanosti formátu DNG s aplikacemi firmy Adobe (která je původcem tohoto formátu), jednak také k nemalé ceně těchto aplikací, která může být pro řadu paměťových organizací zásadní provozní překážkou. To ukazuje, jak již bylo uvedeno výše, že požadavky na archivační formát jdou nad rámec požadavků na obrazovou matici.

1.2.2.1 Archivační formáty pro obrazovou komponentu

Směrodatné zahraniční zdroje uvádějí formáty TIFF a JP2 jako dva hlavní archivační formáty pro obrazová (rastrová) data. Tyto formáty splňují v dostatečné míře uvedené tři základní podmínky kladené na archivační formát (podpora, otevřenost, nezátíženost patenty). Je nepochybně výhodou, že v oblasti rastrových formátů existuje možnost takové volby vzhledem k tomu, že existují jiné typy dat, pro které otevřené formáty zatím nejsou k dispozici nebo se neužívají.

Směrnice FADGI doporučují jako archivační formát pro všechny typy obsahu TIFF, verze 6 a JPEG 2000, u některých typů obsahu pak ještě PNG a PDF/A (FADGI, 2023, s. 23). Doporučení Floridského digitálního archivu označovalo za rastrové formáty s nejvyšší spolehlivostí nekomprimovaný TIFF, JP2 v bezztrátové kompresi a PNG (Florida Center for Library Automation, 2012). V Deklaraci doporučených formátů Kongresové knihovny pro 2022-2023 jsou pro digitální obrazová data (digitální fotografie a další typy digitálních obrazových dat) uvedeny jako preferované archivační formáty: TIFF, JP2, PNG a JPEG/JFIF (Library of Congress, 2022). Směrnice pro budování kvalitních digitálních sbírek (Framework of Guidance for Building Good Digital Collections) vydané americkou normalizační organizací NISO doporučují jako archivační formáty nekomprimovaný TIFF a bezztrátově komprimovaný JPEG 2000 (NISO Framework working group, 2007, s. 28). Podle rozsáhlého průzkumu provedeného v letech 2012-2013 považují severoamerické knihovny formát TIFF za nejspolehlivější archivační formát vůbec (bez ohledu na typ

dokumentu) (Rimkus et al., 2014).

V současné praxi je nejužívanějším archivačním formátem TIFF.¹¹ Nekomprimovaný TIFF jako archivační formát využívají pro své digitalizační projekty například Kongresová knihovna (FADGI, 2014, s. 5), pro některé projekty i Národní knihovna Švédska¹². Do roku 2014 jej jako archivační formát využívala i Národní knihovna Francie, od roku 2014 je jejich archivačním formátem JPEG 2000 (Duploy, 2017). Figuruje ve všech formátových politikách lokálních institucí zveřejněných online (Ostráková, Kopský, 2020, tabulka 2) a ve třech případech jako jediný důvěryhodný formát. Převaha formátu TIFF jako archivačního formátu digitalizačních projektů souvisí nejen s důvěrou, které se těšil a stále těší, ale také s jeho dlouhou historií (specifikace aktuální, tj. šesté verze byla vydána v roce 1992) a širokou podporou v technologickém prostředí od 90. let 20. století do současnosti.

Přibližně posledních deset let se začíná v novějších digitalizačních projektech (zejména evropské provenience) rozšiřovat užití formátu JP2 jakožto archivačního formátu (Van der Knijff, 2011; FADGI, 2014). Jistý vliv na to může mít i příznivá studie italských odborníků, která byla v roce 2008 vydána v periodiku D-Lib Magazine (Buonora, Liberati, 2008). Tato studie srovnávala tři rastrové formáty (TIFF, JP2 a JPEG) z hlediska jejich vhodnosti pro archivaci. Za nejvhodnější archivační formát označila právě JP2, mj. díky jeho největší robustnosti (odolnosti vůči menšímu poškození bitů). Formát JP2, který byl původně v knihovnách užíván spíše jako formát pro zpřístupnění, se stal archivačním formátem například pro digitální obrazové fondy knihovny Wellcome Library (Buckley, 2009) nebo pro masovou digitalizaci norské národní knihovny (Brygfjeld, 2010). Prvně jmenovaná knihovna užívá profil formátu JP2 ve ztrátové kompresi, druhá v bezztrátové. Užití ztrátové komprese není výjimečné, nicméně bezztrátové je častější. V praxi jsou voleny různé profily formátu JP2, zahrnující nejen volbu typu komprese, ale některé další parametry specifické pro tento formát (JPEG 2000 profiles, 2010). Vhodné nastavení těchto parametrů vyžaduje specialistu; z toho důvodu si některé knihovny nechaly vypracovat profily na zakázku u externích odborníků.¹³

K formátu JP2 se vyskytly i kritické názory (Van der Knijff, 2011) a o možnosti jeho využití jako archivačního formátu se vedla intenzivní diskuze (FADGI, 2014, s. 3). Obecně

¹¹ Srv. např. FADGI, 2014, s. 3-4; Rimkus et al., 2014; Van der Knijff, 2011.

¹² Digitalizované noviny jsou archivovány ve formátu JPEG 2000, ostatní dokumenty většinou ve formátu TIFF (Neiss, 2017).

¹³ Například profil formátu JP2 pro Wellcome Library vytvořil Robert Buckley, jeden z autorů specifikace formátu JP2 (viz Buckley, 2009). Buckley vytvořil specifikaci formátového profilu JP2 i pro řadu dalších knihoven.

lze říci, že v USA převažuje volba formátu TIFF a v Evropě se rozšiřuje užití formátu JP2. Užití formátu PNG jako archivačního formátu je v současné praxi řídké, což může být vzhledem k jeho kvalitám překvapivé.

1.2.2.2 Archivační formát pro OCR komponentu

Optické rozpoznávání znaků (OCR) je metoda získávání textu z obrazu. V současné praxi se užívá k tomu, aby se z rastrových dat vzniklých digitalizací vytěžil textový obsah. Výstup z OCR se ukládá v podobě strukturovaného textového formátu, který obsahuje informace o pozici (obrazem vyjádřených) konkrétních písmen (slov) v obrazovém souboru, z něhož byl vytvořen, aby bylo zajištěno namapování textu na obraz. Tímto formátem je v současné praxi převážně ALTO. Formát ALTO sám sebe popisuje jako „standardizovaný formát XML k ukládání informací o rozložení (layout) a obsahu“ (ALTO Principles, 2016). Formát ALTO XML je navržen jako externí schéma pro standard METS. Jde však především o datový formát (obsahuje vlastní text předlohy a jeho strukturu), částečně slouží i jako metadatový formát (např. popis informací o zdrojovém obrázku). V případě, že se vytvářejí archivační formát i prezentační varianty, je nutné, aby měly stejnou pixelovou velikost (počet pixelů na šířku a výšku obrázku, jinak text nebude správně namapován).

Míra efektivity OCR závisí na několika faktorech: stav předlohy, kvalita softwaru a jeho slovníků a komprese. Podle některých výzkumů¹⁴ přinášejí ztrátově komprimovaná obrazová data mírně lepší výsledky procesu OCR v porovnání s bezztrátovou kompresí (Chapman 2007, s. 39; Buckley, 2008, s. 23). Záleží však také na tom, v jakém formátu je ztrátová komprese – různé nástroje pro OCR mohou s různými formáty pracovat odlišně (např. Tesseract v minulosti nepodporoval formát JP2). V praxi platí, že pro novodobé fondy v dobrém stavu předloh (a pro běžné jazyky) je efektivita OCR velmi vysoká.

Prezentace digitalizátu knihy v digitální knihovně zahrnuje nejen zobrazení, ale i reprodukci těchto strukturovaných textových informací získaných z OCR, která umožňuje čtenářům plnotextové prohledávání obrazových dat. Přidáním této strukturované textové složky je digitalizát knihy obohacen o funkci, kterou jeho tištěná předloha nikdy neměla.

¹⁴ Národní archivy Velké Británie v roce 2017 naopak uvádějí že ztrátová komprese optické rozpoznávání textů ovlivňuje negativně (The National Archives, 2017, s. 3).

1.2.3 Prezentační formáty

Volba prezentačního formátu závisí na požadavcích cílové komunity. Jelikož v současné praxi je hlavní formou zpřístupnění digitalizátů jejich zobrazení v digitální knihovně, tyto požadavky se řídí především podporou formátu v internetových prohlížečích. Směrnice MINERVA doporučují zpřístupňovat digitální kopie dokumentů v různých velikostech nebo formátech, aby zacílení na uživatele bylo co nejširší (Ferne, 2008, s. 73). Směrnice pro budování kvalitních digitálních sbírek (Framework of Guidance for Building Good Digital Collections) vydané americkou normalizační organizací NISO doporučují pro zpřístupňování digitalizovaných dokumentů užití formátů JPEG a PDF.

1.2.3.1 Prezentační formáty pro obrazovou komponentu

V praxi je nejužívanějším hlavním prezentačním formátem digitálních knihoven JPEG (FADGI, 2014). Tato volba je logická vzhledem k jeho vysoké podpoře v internetových prohlížečích nebo zobrazovacích aplikacích. Formát JPEG jako prezentační formát užívá například Gallica, jedna z největších digitálních knihoven světa, provozovaná francouzskou národní knihovnou, v níž je tento prezentační formát vytvářen konverzí z formátu TIFF, který je archivačním formátem této knihovny (Bruys et al., 2019). Formát JP2 není podporován běžnými internetovými prohlížeči, pro jeho užití jako prezentačního formátu je nutno nainstalovat plugin. Běžným způsobem zpřístupnění, který řeší tento problém z hlediska komfortu cílové komunity, je využití formátu JP2 jako prezentačního meziformátu, ze kterého digitální knihovna generuje formát JPEG, v němž jsou obrazová data prezentovaná čtenářům internetovým prohlížečem (Buckley, 2009, s. 11). Jako doplňkové vedlejší formáty se užívají například formát PDF (pro možnost stažení digitalizátu knihy čtenářem v podobě jednoho souboru) nebo formát GIF (pro náhledy obrázků).

1.2.4 Formátové registry

V souvislosti s přijetím modelu OAIS se v komunitě paměťových institucí objevil navazující koncept globálního registru interpretačních informací. Idea takového registru spočívá v tom, že bude sloužit jako zásobník interpretačních informací, které jsou potřebné pro digitální objekty (zejména informací o formátech, souvisejících aplikacích a všech ostatních prvcích počítačového prostředí, jež podporuje adekvátní reprodukci digitálních objektů a jejich zpracování). Za tímto návrhem stála pragmatická úvaha, podle níž nejsou jednotlivé instituce schopny všechny potřebné interpretační informace shromažďovat

vlastními silami.

Jako první vznikl registr PRONOM (Brown, 2006). Byl založen a je již patnáct let provozován britskými Národními archivy (The National Archives). Ačkoliv PRONOM zdaleka nesplňuje své původní ambice (obsahuje poměrně rozsáhlou databázi formátů, ale jejich popis je většinou minimální), jde v současnosti o nejvýznamnější důvěryhodný projekt, který se alespoň snaží uskutečňovat původní vizi směrodatného globálního registru interpretačních informací. Později sice vznikly dva další registry, ale ty již zanikly. Prvním z nich byl Global Digital Format Registry (GDFR), který skončil již před několika lety.¹⁵ Registr Unified Digital Format Registry (UDFR) provozovaný Kalifornskou digitální knihovnou (California Digital Library) měl spojit registry GDFR a PRONOM, ale byl ukončen v dubnu 2016 kvůli nedostatku financí.¹⁶

Klíčovou funkcí registru PRONOM je to, že (jako jediný registr vůbec) nabízí jednoznačný a jedinečný identifikátor formátu, přesněji řečeno, jak uvádí sám registr:

„rozšiřitelné schéma pro poskytování perzistentních, jedinečných a jednoznačných identifikátorů pro jednotky interpretačních informací zaznamenané v registru PRONOM“ (Brown, 2006, s. 4). Formát je tedy pouze jedním z typů interpretačních informací, o nichž registr vede údaje, nicméně nejrozšířenějším. Funkce identifikátoru PUID jsou dvě: propojení se záznamem jednotky interpretačních informací v registru PRONOM (tj. způsob identifikace záznamu, přičemž tento záznam by ideálně měl obsahovat co nejpodrobnější informace o formátu nebo jiné jednotce interpretačních informací) a jedinečný perzistentní identifikátor, který odlišuje v maximální možné míře jeden formát od druhého (odlišuje se nejen typ formátu, ale často i verze¹⁷). Například PUID pro JPEG verze 1.00 je „fmt/42“, pro verzi 1.01 „fmt/43“ a pro verzi 1.02 „fmt/44“. Registr MIME,¹⁸ který je nejužívanějším obecným registrem formátů (sloužícím i pro účely mimo kontext digitální archivace), odlišuje formáty jen na základě typu a názvu, například formát JPEG všech verzí má označení „image/jpeg“.

Za druhý významný zdroj informací o formátech lze považovat registr Kongresové knihovny.¹⁹ Ten obsahuje nejen základní interpretační informace o formátech, ale také o rizicích s nimi spojených (Library of Congress, 2013).

¹⁵ https://web.archive.org/web/20171009011353/https://library.harvard.edu/preservation/digital-preservation_gdfr.html

¹⁶ <https://web.archive.org/web/20220511203412/http://udfr.org/>

¹⁷ Například u formátu EPUB se jeho jednotlivé verze v registru PRONOM nerozlišují.

¹⁸ <http://www.iana.org/assignments/media-types/media-types.xhtml>

¹⁹ <https://www.loc.gov/preservation/digital/formats/index.shtml>

Formátový registr PRONOM nabízí záznamy o formátech a dalších jednotkách interpretačních informací v podobě volně dostupných webových stránek. Současně poskytuje identifikační mechanismus, který obsahuje popis toho, kde se v daném formátu nachází tzv. „magické číslo“ (údaj o verzi formátu), a identifikátory PUID. „Magické číslo“ je interní mechanismus označení konkrétního formátu daný formátovou specifikací, zatímco identifikátor PUID je jedinečný externí identifikátor, který je přidělován registrem PRONOM. Registr PRONOM také nabízí svůj vlastní nástroj, DROID, který provádí formátovou identifikaci užitím uvedeného mechanismus na jednotlivé soubory.

1.2.5 Specifikace obrazových dat

V digitalizační praxi je třeba u obrazových dat rozlišovat několik dalších vlastností, které jsou nad rámec formátu: generace dat; strukturální model obrazové reprezentace; formátový profil; typy komprese; obecné obrazové vlastnosti; prezentační varianty obrazových dat.

1.2.5.1 Generace obrazových dat

Prvotní výstup snímání skenerem nebo fotoaparátem budeme nazývat původní snímky. Původní snímky jsou soubory v rastrových formátech uložené po snímání na datový nosič pracovní stanice digitalizační linky. Tyto původní snímky zpravidla procházejí dalším zpracováním do doby, než je vytvořen konečný obrazový výstup digitalizace (finální produkční data). Každá transformace, kterou prošla obrazová data od původních snímků po vytvoření finálních produkčních dat, znamená vytvoření nové generace obrazových dat. K vytváření dalších generací dochází také v pozdějších etapách (archivace a zpřístupňování), nicméně cíle vytváření těchto generací jsou odlišné. Cílem digitalizačních transformací je vytvořit konečný digitální produkt a teprve tento konečný produkt lze považovat za plnohodnotnou obrazovou složku digitalizátu knihy. Cílem archivace a zpřístupnění je zachovat tento digitalizační produkt v požadované kvalitě, v případě archivace z hlediska uchování informačního obsahu navzdory rizikům technologického zastarávání, v případě zpřístupnění vytvoření takové podoby digitalizátu dat, která je vhodná pro aktuální potřeby cílové komunity a která se může lišit od dat uložených v balíčku AIP. Původní snímky v současné praxi mohou být v odlišném formátu než finální produkční data. Generace se od sebe liší změnami, které lze zaznamenat na bitové úrovni, přičemž však může platit, že

některé transformace je možno vykonat společně před tím, než bude uložena nová generace obrazových dat.

1.2.5.2 *Strukturální model obrazové reprezentace*

Modelem obrazové reprezentace se zde rozumí vztah mezi strukturou tištěné knihy (posloupností stran, fyzická mapa) a způsobem digitální reprezentace této předlohy v digitalizátu knihy (logická mapa). Převažující strukturální model současné praxe lze charakterizovat následovně:

- Základní předmět reformátování = jedna stránka knihy.
- Základní objekt uložení = jeden soubor reprezentující tuto jednu stránku.

Tento model znamená, že jeden soubor v rastrovém formátu reprezentuje jednu stránku knižního bloku, přičemž všechny rastrové soubory mají stejnou velikost. Pro reprezentaci jiných částí knih, než jsou stránky (např. přebal), nebo nestandardních částí (např. stránka s mapou, kterou lze rozložit, takže její velikost bude jiná než ostatních stran), se v praxi užívají různé postupy. Běžnou současnou praxí je snaha zachytit všechny části knihy, včetně vakátů a prázdných přídeští (jedinou výjimkou v tomto směru je zadní část přebalu, která bývá bílá a nepovažuje se za smysluplné ji digitalizovat).

V rámci Standardu NDK jsou strukturální mapy dokumentu tvořeny pomocí formátu METS, konkrétně oddílu <structMap>.

Fyzická strukturální mapa reprezentuje dokument, otištěný do digitální reprezentace. Kromě fyzické posloupnosti jednotlivých naskenovaných stran (reprezentace stran) jsou do každého <div>, reprezentujícího jednu stranu, vloženy také informace o typu strany a odkazy na jednotlivé komponenty v jednotlivých složkách SIP balíku, ze kterých se digitální objekt strany skládá. Tyto komponenty tvoří: *mastercopy* (archivní verze) obrazu, *usercopy* (uživatelská kopie) obrazu, *amd_mets* (tzv. „vedlejší mets“ s technickými a administrativními metadaty), *alto* a *txt* (ocr data).

Logická strukturální mapa potom reprezentuje hierarchický model reprezentace dokumentu, založený na jeho fyzické předloze. Pomocí jednotlivých zanořených <div> pro každou úroveň popisu umožňuje interpretovat posloupnost celků, které tvoří zdigitalizovaný dokument od např. svazků ročníků periodik (VOLUME) přes jednotlivá čísla periodika (ISSUE) po jednotlivé články (ARTICLE).

1.2.5.3 *Formátový profil*

Formátovým profilem se rozumí nastavení v rámci konkrétního formátu. Mezi hlavní prvky profilu rastrových formátů patří komprese (tj. její nastavení v rámci možností formátu, přičemž ne všechny rastrové formáty umožňují všechny typy kompresí). Nejčastějším profilem pro formát TIFF je nekomprimovaná varianta. Formátový profil pro JP2 zahrnuje (nutnou) volbu mezi ztrátovou a bezztrátovou kompresí a dále několik dalších specifických nastavení typických pro tento formát.

1.2.5.4 *Typy komprese*

Klíčovým aspektem rastrových obrazových dat je komprese. Pro potřeby této práce postačí následující klasifikace čtyř základních typů:

1. Nekomprimovaná varianta (tj. komprese není užitá)
2. Matematicky bezztrátová komprese
3. Vizuálně bezztrátová komprese
4. Vizuálně ztrátová komprese

Nekomprimovaná varianta může být teoreticky nejlepší možnou volbou pro archivační formát. Studie iniciativy FADGI uvádí: „Nekomprimovaná datová struktura má jednu velkou výhodu: je relativně transparentní. Transparentnost souvisí s ukazatelem udržitelnosti: nemělo by být složité vytvořit nástroj, který dokáže přečíst informaci o obalu (wrapper) a rozbalit rastrová data tak, aby je bylo možno zobrazit“ (FADGI, 2014, s. 4). Jediným problémem s užitím nekomprimovaných dat je jejich velikost. V současné praxi může tento problém představovat velkou překážkou vzhledem k omezenému rozpočtu na úložné kapacity.

Matematicky bezztrátová komprese obrazových dat v principu odpovídá konceptu vratné transformace modelu OAIS. Kompresní algoritmus snižuje transparentnost formátu, ale současně umožňuje snížení požadavků na úložné kapacity.

Vizuálně bezztrátovou kompresí se myslí taková matematicky ztrátová komprese, která na základě nějakého přijatého psychofyziologického modelu stanoví kompresní poměr, jehož výstupem (při zobrazení) má být informační obsah, který by běžný pozorovatel neměl rozeznat od výstupu matematicky bezztrátové komprese, nebo jsou viděné rozdíly nepodstatné (Buckley, 2009, s. 4).

Vizuálně ztrátová komprese je taková matematicky ztrátová komprese, která přináší vizuálně patrné změny obrazové kvality. Míra komprese v rámci tohoto typu se může dále výrazně lišit, přičemž nejvyšší možnou vizuálně ztrátovou kompresi lze pochopitelně užít pouze pro náhledy obrázku (tedy pro doplňkové prezentační formáty).

Užití ztrátové komprese se v archivaci obecně nedoporučuje, protože je spojeno s rizikem generační ztráty informace v průběhu migrace z jednoho ztrátově komprimovaného formátu do jiného.²⁰ Podobné riziko (zde pod označením “*cumulative loss*”) je zohledňováno i u archivace filmových materiálů (Blood, 2011), přestože zde je z důvodů datové náročnosti motivace použít kompresní formát ještě vyšší. Výjimkou je samozřejmě situace, kdy je obsah do úložiště přijat již ve ztrátové kompresi. Nicméně i zde již u obrazových materiálů probíhají normalizace do bezztrátových formátů a u videoformátů v tomto směru běží diskuze a proběhlo testování konverzí mezi některými formáty.

Ne všechny rastrové formáty umožňují všechny výše uvedené volby. Formát TIFF umožňuje všechny varianty: nekomprimovanou variantu (která je také nejčastěji užívaná pro TIFF jako archivační formát), matematicky bezztrátovou kompresi (algoritmus LZW nebo ZIP) a ztrátovou kompresi (algoritmus JPEG). Formát JPEG nabízí pouze matematicky ztrátovou kompresi (typy 3-4).²¹ Formát PNG nabízí pouze bezztrátovou kompresi, podle studie FADGI s vynikajícími výsledky (FADGI, 2014, s. 3). Formát JP2 nabízí pouze komprimované varianty (typy 2-4). Koncept vizuálně bezztrátové komprese (typ 3) je v současné praxi spojován právě s tímto formátem. Skupina výzkumníků z několika významných knihoven světa provedla mezi čtenáři rozsáhlý test vnímání (různě vysoké) ztrátové komprese formátu JP2 s cílem navrhnout vhodné profily formátu JP2 ve vizuálně bezztrátové kompresi (Chapman, 2006). Některé knihovny pak začaly v praxi využívat vizuálně bezztrátovou kompresi pro archivační formáty.²² Pro matematicky bezztrátovou kompresi se udává, že kompresní poměr je obvykle zhruba 2:1 (Buckley, 2008, s. 6).

1.2.5.5 *Obecné obrazové vlastnosti*

Směrnice FADGI určují čtyři základní obrazové vlastnosti rastrových dat prostorové rozlišení (*spatial resolution*), bitovou hloubku (*bit-depth*), barevný prostor (*color space*) a barevný mód (*color mode*). Pro tyto vlastnosti udává doporučení odpovídající čtyřem

²⁰ https://en.wikipedia.org/wiki/Generation_loss

²¹ Srv. Buonora, 2008.

²² Například Wellcome Library (Buckley, 2009).

stupňům kvality, označené jednou až čtyřmi hvězdičkami (FADGI, 2023).

Prostorové rozlišení určuje množství informací v rastrovém souboru z hlediska počtu pixelů (obrazových prvků) na jednotku měření, obvykle palec (odtud zkratka PPI),²³ tj. „stanovuje, jak blízko od sebe jsou jednotlivé pixely umístěny“; bitová hloubka „stanovuje maximální počet odstínů nebo barev v digitálním obrazovém souboru“ (FADGI, 2010, s. 4). Prostorové rozlišení a rozměry digitálního obrazu určují celkový počet pixelů v souboru; při určení požadované velikosti souboru je nutno zadat hodnotu prostorového rozlišení a rozměry (např. 300 PPI + 8x10 palců). Ve směrnici FADGI je bitová hloubka 8 bitů na kanál u *general collections* (běžné knihy novodobých fondů) doporučena v nižším standardu 1 a 2 hvězdičky, zatímco ve vyšším (3 a 4 hvězdičky) je možné si zvolit mezi 8 a 16 bity na kanál. Obrazy ve stupních šedi mají jeden kanál, barevné obrazy tři a více kanálů. Barevný model je způsob číselné specifikace barev s užitím tří nebo více kanálů. Například barevný model RGB obsahuje tři kanály o bitové hloubce 8 nebo 16 bitů.

Prostorové rozlišení tedy vymezuje, jak detailně může být převedena tištěná kniha (např. čitelnost písma), a bitová hloubka určuje, jak věrně mohou být zachyceny její barvy nebo odstíny šedi (tedy barevná věrnost) ve výsledném souboru. Prostorové rozlišení, ani bitová hloubka logicky nemohou udávat ani zaručovat kvalitu uložených informací, pouze vymezují rozsah možné kvality digitalizačního převodu. Kvalita digitalizátu se odvíjí od míry detailnosti (např. velikost písma) a barevnosti předlohy, která je digitálně zachycována, a nastavených hodnot prostorového rozlišení a bitové hloubky. Obecně platí, že rozmezí prostorového rozlišení a bitové hloubky je na jedné straně ovlivněno minimální hranicí (tj. aby kniha ještě byla čitelná a její obrazové prvky do určité míry rozpoznatelné), na druhé straně maximální hranicí rozlišení, tj. takové, nad jejíž rámec snímání nemůže v principu přinést již žádný pozorovatelný rozdíl, a digitalizace ve vyšším rozlišení by byla neekonomická či jinak neúčelná. Existující doporučení se tedy mohou pohybovat pouze mezi těmito krajními případy. Doporučení současné praxe však také zohledňují skutečnost, že při masové digitalizaci je prakticky nemožné vytvářet tato nastavení pro každou knihu zvlášť. Z tohoto důvodu se vydávají plošná doporučení pro minimální rozlišení a bitovou hloubku pro různě definované kategorie knih z hlediska jejich předpokládané velikosti písma a barevnosti.

Pro běžné knihy novodobých fondů obsahující barevné prvky se obecně doporučuje

²³ Zkratka pro počet pixelů na palec (pixel per inch), někdy se užívá též počet pixelů na milimetr nebo centimetr.

barevný model RGB, bitová hloubka 8 bitů na kanál a minimální prostorové rozlišení 300-400 PPI (FADGI, 2016; The Association for library collections and technical services, 2013). Digitalizační projekt může vytvořit odlišné pracovní postupy například pro tištěné knihy obsahující barevné ilustrace nebo fotografie a tištěné knihy obsahující pouze text (a druhé snímat jen ve stupních šedi). Při masové digitalizace však může být z hlediska přípravy obtížné kontrolovat, zda kniha neobsahuje barevné prvky. Z tohoto důvodu se často vytváří jedno plošné nastavení obrazových vlastností pro všechny knihy.²⁴

V případě prostorového rozlišení nemohou být doporučení minimálního rozlišení ničím jiným než predikcí očekávané nejmenší velikosti písma u určité skupiny předloh. V případě, že se ve skupině vyskytne anomálie, pak tato plošná doporučení přirozeně nemohou zaručit kvalitní výsledek snímání. Hodnota prostorového rozlišení se nastavuje pouze na skeneru, pro fotografování tento údaj přirozeně nemá smysl. Na výsledné rozlišení fotografie má vliv kvalita senzoru fotoaparátu a objektivu, vzdálenost knihy od objektivu a zaostření.

Zatímco nastavení výše uvedených vlastností (prostorové rozlišení, bitová hloubka a barevný prostor) v současné praxi nepředstavuje větší problém, barevný profil je složitější problematika. Barevný profil „určuje interpretaci číselných hodnot popisujících pixely v obrázku tím, že popisuje chování zařízení nebo rozsah barevného prostoru“ (FADGI, 2010, s. 45). Barevný prostor je „geometrická reprezentace barev v prostoru, který lze vizuálně vnímat nebo vytvářet užitím konkrétního barevného modelu“ (FADGI, 2017). Barevný prostor je například vyžadován aplikacemi pro zobrazení. Snímací zařízení zpravidla užívají vlastní barevný profil, který je závislý na konkrétním zařízení nebo výrobci. Tento technologicky závislý profil lze převést do ICC profilu (standardu pro univerzální barevnou specifikaci) a uložit do obrázku.

Směrnice FADGI doporučuje barevné prostory v závislosti na originálu – ale i u general collections je od dvou hvězdiček výše i ProPhoto RGB (spolu s sRGB, Adobe RGB 1998 a ECIRGBv2), naopak u manuskriptů už vůbec chybí sRGB (FADGI, 2023). Knihovna LOC pak v doporučeních pro obrazové vlastnosti zdůrazňuje, že barevný prostor v němž byl digitalizát vytvořen musí být zachován, transkódování barevného prostoru se nedoporučuje (Library of Congress, 2006). Otázce správy barev a barevným profilům se podrobněji věnuje kapitola 4.2.6.

²⁴ Příkladem je digitalizace novodobých dokumentů v Národní knihovně ČR, kde bylo do roku 2021 plošně pro všechny digitalizované dokumenty nastaveno rozlišení 300 PPI, barevný model RGB, bitová hloubka 8 bitů na kanál. V průběhu roku 2021 došlo ke zvýšení rozlišení u všech dokumentů na 400 PPI.

1.2.5.6 *Prezentační varianty*

Hlavním cílem současné praxe je vytvoření finálních produkčních dat v archivačním formátu. Obrazová data v archivačním formátu se následně v nezměněné podobě uchovávají v archivu do doby, než bude z důvodů zastarávání technologií nutno přistoupit k archivačním opatřením. Běžnou praxí je, že se v digitální knihovně čtenářům nezpřístupňují obrazová data v archivačním formátu, ale vytváří se jejich prezentační varianta v prezentačním formátu. V případě užití archivačního formátu TIFF bývají prezentační variantou zpravidla obrazová data ve formátu JPEG jako hlavním prezentačním formátu,²⁵ v případě archivačního formátu JP2 bývá vytvořen prezentační meziformát ve formátu JP2 (ve ztrátové kompresi), ze kterého systém digitální knihovny vytváří za chodu prezentační varianty ve formátu JPEG.

Důvody vytváření prezentační varianty pro zpřístupnění jsou různé a zpravidla jsou kombinací více faktorů. Zaprvé jsou dány historicky. Digitalizace se prováděla již v dobách, kdy bylo internetové připojení ještě pomalé nebo nákladné, a tudíž prezentační varianta tvořená obrazovými soubory menší velikosti (zejména ve ztrátové kompresi) byla uživatelsky vhodnou formou prezentace. Druhým důvodem je, že formát TIFF ani JPEG 2000, nejčastější formáty digitalizačních projektů, nejsou podporovány internetovými prohlížeči. Třetím důvodem může být následování praxe obrazového průmyslu, kdy je obvyklým způsobem ukládat obrazové matrice v nejvyšší možné kvalitě a užívat je jako zdroj pro generování obrazových dat v různé kvalitě pro různé účely (mj. také prezentace na webu, např. v online periodících). Odůvodněním také mohou být výzkumy stanovující psychofyziologický model, podle něhož je informační obsah reprodukováný z obrazových dat ve ztrátové kompresi určité úrovně čtenářem vizuálně nerozeznatelný od informačního obsahu reprodukováného z obrazových dat v bezeztrátové kompresi nebo nekomprimované podobě.²⁶ Z toho se vyvozuje, že zpřístupňování obrazových matic v archivačním formátu je neúčelné. Nikoliv výjimečným případem současné praxe je, že se prezentační varianty vytvářejí již při produkci, jako součást finálních produkčních dat (a tedy nikoliv až v archivu), a to z důvodu jednoduššího zpracování.²⁷

²⁵ Srv. např. Smith, 2006, s. 10 a Vychodil, 2010, s. 64.

²⁶ Viz např. Chapman, 2007.

²⁷ Viz např. Standardy pro obrazová data na <https://www.ndk.cz/standardy-digitalizace/standardy-pro-obrazova-data>.

1.3 Metadatové standardy

Metadata lze ve zkratce definovat jako “data o jiných datech” (ISO 14721; s. 24). Výstižnější je ovšem definice, která metadata popisuje jako “informace, které vytváříme, ukládáme a sdílíme za účelem popisovat věci tak, abychom s nimi mohli interagovat a získávat z nich vědomosti, které potřebujeme” nebo ve specifickém případě digitálních knihoven “strukturované informace, které popisují, vysvětlují, lokalizují nebo jinak usnadňují vyhledávání, používání nebo správu informačního zdroje” a zároveň “data, spojovaná s informačním systémem nebo informačním objektem za účelem popisu, administrace, správy právních požadavků a technických funkcionalit, jejich použití a využití”. Metadata jsou stejně jako datový objekt s obsahem předmětem interpretace softwarovými nástroji, dalšími prvky počítačových systémů a také lidskými uživateli. Jsou zároveň součástí klíčových přidružených informací, které se zaznamenávají a udržují spolu s datovým objektem a obsahem, což je důvodem, proč jsou současné metadatové standardy, užívané pro specifikaci ISO 14721 v komunitě paměťových institucí pro záznam specifických typů informací využívány.

Ačkoliv se pojem metadata může vztahovat na širokou škálu typů (např. i na metadata uložená v obrazových datech nebo v databázi), standardy, užívané v paměťovými institucemi jsou specifickým typem metadat, která se vyznačují následujícími charakteristikami:

- jsou vysoce strukturovaná;
- ukládají se do samostatných textových souborů ve formátu XML;
- uchovávají se spolu s datovým objektem s obsahem v informačním balíčku;
- jsou pečlivě zdokumentovaná mezinárodními standardy;
- užívá je většina organizací dané oblasti digitálního dědictví (např. knihovnictví).

Tento typ metadatových formátů byl vyvinut mezinárodním společenstvím informačních profesionálů na základě katalogizačních standardů, původně sloužících pro popis tištěných a dalších dokumentů v knihovních fondech. Tyto katalogizační standardy již umožňovaly automatické zpracování záznamů a strojové čtení, což bylo pro metadatové formáty klíčové. Předchozí vysoká standardizace katalogizačních záznamů také umožnila široké přijetí metadatových standardů napříč knihovními institucemi.

Obecně se v knihovní digitalizaci používají metadatové formáty několika typů: popisné, strukturální, technické, administrativní a právní, přičemž každý z těchto typů popisuje některý aspekt uchovávaného digitálního objektu.

Všechny standardy se skládají z vlastní sady elementů (*element set*), což je kontrolovaná množina prvků, které se užívají buď samostatně nebo v kombinacích, které standard předepisuje. Standardy obecně předepisují možnosti popisu specifických typů zdrojů, konkrétní účel užití, definují význam jednotlivých elementů, jejich vztahy a hierarchii, a v neposlední řadě poskytují návod, jaké hodnoty a jak by měly být v konkrétních situacích nebo u konkrétních zdrojů užívány.

Sada elementů je obvykle doplněna sadou atributů, které dále dodefinovávají jednotlivé elementy ze sady. Každý element jich může mít přiřazených několik, které popisují jeho různé vlastnosti (název, typ, identifikátor a další).²⁸

Vzhledem k častému používání několika metadatových standardů v rámci jednoho metadatového zápisu je důležité, aby všechny byly formulovány v jednom obecném rámci, který umožňuje interoperabilitu, konzistenci, srozumitelnost a zároveň je komunikovatelný mezi různými komunitami jak knihovních, tak jiných paměťových institucí (Zeng, 2016, s. 38).

Z tohoto důvodu bylo obecně adaptováno schéma XML, které všechny tyto požadavky splňuje (Zeng, 2016 s. 131).

Pro ilustraci aktuálně vnímané důležitosti úlohy metadat (výše uvedeného typu) v současné praxi poslouží následující citace, která pochází přímo z jednoho z takových standardů: „Bez strukturálních metadat jsou obrázek stránky nebo textový soubor tvořící digitální dílo téměř k ničemu a bez technických metadat zohledňujících digitalizační proces si badatelé nemohou být jisti, jak přesný odraz originálu digitální verze poskytuje. Pro účely vnitřní správy musí mít knihovna přístup k náležitým technickým metadatům, aby mohla pravidelně obnovovat a migrovat data, a tak zajistit zachování cenných zdrojů“ (METS, 2016).

²⁸ Například PREMIS má u každého elementu uveden jedinečný identifikátor založený na hierarchické řazení, např. 1.2, 3.1.

1.3.1 Přehled metadatových standardů pro digitalizáty tištěných dokumentů

Již před érou digitálních metadat byla knihovnická komunita známá vysokou mírou jednotné standardizace na mezinárodní úrovni. Tak tomu je i v případě metadatových standardů pro digitální data. Knihovny pro správu digitalizátů nejčastěji užívají tyto mezinárodní metadatové standardy: METS,²⁹ PREMIS,³⁰ MODS,³¹ Dublin Core³² (dále jako “DC”), MIX³³ a ALTO³⁴.

Tyto standardy lze z hlediska rozsahu jejich možné aplikace rozdělit do dvou skupin:

a) omezené užití (MIX; ALTO), b) univerzální užití (METS, PREMIS, MODS, DC). První skupina je určena specificky pro digitalizáty knih (a některé další typy digitálních dokumentů), druhou skupinu lze aplikovat na většinu typů digitálních dokumentů od digitalizátů tištěných dokumentů přes digitalizáty audiovizuálních dokumentů po e-born dokumenty. Všechny tyto standardy, s výjimkou PREMIS, jsou definovány specificky jako XML schémata. Standard PREMIS není svázán se schématem XML, ale lze jej jako XML vyjádřit (existuje oficiální XML schéma k tomuto standardu) a XML schéma se také užívá pro balíčky SIP a AIP.

Standard ALTO je primárně datový formát (zaznamenávající text získaný procesem OCR a jeho souřadnicové umístění vzhledem k obrazu), ale obsahuje některé metadatové prvky (např. informace o obrazovém zdroji pro OCR). Standard MODS je užíván pro zápis bibliografických informací, a je tedy určen k naplňování role identifikačních informací OAIS (a případně popisných informací OAIS). Standard DC je vzhledem k omezenému rozsahu sady elementů využíván pouze pro základní zápis bibliografických informací. Vybrané mezinárodní standardy popíšeme detailněji v následujících oddílech.

1.3.2 PREMIS

Standard PREMIS sám sebe označuje jako standard pro archivační metadata (*preservation metadata*), který „podporuje životaschopnost, reprodukovatelnost,

²⁹ <http://www.loc.gov/standards/mets>

³⁰ <http://www.loc.gov/standards/premis/v3/premis-3-0-final.pdf>

³¹ <http://www.loc.gov/standards/mods/mods-overview.html>

³² <http://dublincore.org>

³³ <http://www.loc.gov/standards/mix/>

³⁴ <https://www.loc.gov/standards/alto/description.html>

srozumitelnost, autenticitu a identitu digitálních objektů v archivačním kontextu“ (PREMIS, 2015, s. 1). Slouží však nejen pro zápis archivačních informací, ale také interpretačních informací. Aktuální třetí verze PREMIS vyšla v roce 2015. Standard neobsahuje pouze sadu elementů, ale také vlastní komplexní datový model, terminologický slovník a podrobný text vysvětlující logiku a možnosti užití standardu v archivu. Ve své sebedefinici klade standard PREMIS také důraz na to, aby jeho elementy byly implementovatelné, což podle něj znamená, že hodnoty většiny elementů musí být možné automatizovaně vyplňovat a zpracovávat archivem (Premis, 2015, s. 3). K takovému cíli ostatně směřují všechny zde uvedené metadatové standardy.³⁵

Datový model v PREMIS definuje **čtyři základní entity**: objekt (*object*), činitel (*agent*), událost (*event*) a právní deklarace (*rights statement*). **Objekt** dále člení na čtyři úrovně:

1. intelektuální entita (*intellectual entity*) je „jednotlivý intelektuální nebo umělecký výtvar (*creation*), který je považován za relevantní pro cílovou komunitu v kontextu digitální archivace“;
2. reprezentace (*representation*) „množina souborů (včetně strukturálních metadat) potřebná pro úplnou reprodukci intelektuální entity“;
3. soubor (*file*) „pojmenovaná a uspořádaná posloupnost bajtů, kterou dokáže rozeznat operační systém“ a která je v určitém formátu;
4. bitový tok (*bitstream*) „data v rámci jednoho souboru, která mají smysluplné společné vlastnosti pro archivační účely“ (Premis, 2015, s. 8).

Všechny úrovně (vyjma intelektuální entity) odpovídají pojmu „digitální objekt“ v modelu OAIS, přičemž reprezentace v PREMIS odpovídá pojmu „objekt CDO“. Intelektuální entita odpovídá informačnímu obsahu modelu OAIS s tím rozdílem, že ve standardu PREMIS jde specificky o reprodukováný informační obsah (tj. obsah, který může vnímat člověk).

Intelektuální entitu je možné podle modelu PREMIS také dále specifikovat podle úrovní abstrakce popsaných ve známém knihovnickém modelu FRBR. Model FRBR stanovuje tyto čtyři úrovně: dílo (*work*), vyjádření (*expression*), manifestace (*manifestation*)

³⁵ Týká se to i metadat ve schématu MODS, které obvykle vznikají konverzí bibliografických záznamů ve formátu MARC.

a exemplář (*unit*).³⁶

PREMIS obsahuje elementy, které odpovídají všem typům archivačních informací. Klíčové jsou zejména elementy pro zápis provenienčních informací. V tomto ohledu PREMIS vhodně předepisuje logiku metadatového zápisu: „metadata, soubory, bitové toky a reprezentace uchovávané v archivu se popisují jako statické množiny bitů. Není možné změnit soubor (nebo bitový tok nebo reprezentaci); lze pouze vytvořit nový soubor (nebo bitový tok nebo reprezentaci), který se vztahuje k zdrojovému objektu“ (Premis, 2015, s. 22). Tento vztah mezi novým a předchozím objektem definuje jako vztah odvození (*derivation relationship*), u něhož musí být zaznamenán specifický typ události, odlišný od událostí, které nevytvářejí nový objekt. Standard odlišuje dva typy odvození ze zdrojového digitálního objektu do nového objektu: replikace (*replication*) a transformace (*transformation*) (Premis, 2015, s. 19). Replikace znamená vytvoření digitální kopie, která je bitově identická se zdrojovým digitálním objektem (Premis, 2015, s. 272), transformace má za výsledek vytvoření jednoho nebo více digitálních objektů, které nejsou bitově identické se zdrojovým objektem (Premis, 2015, s. 273).

Pro strukturální interpretační informace slouží sekce elementů popisujících formát (název formátu; verze formátu; název formátového registru; identifikátor záznamu formátu v tomto registru; role registru). Pro podrobnější popis interpretačních informací je ve standardu PREMIS vyčleněna možnost vnořit externí schéma.³⁷ Pro digitalizáty knih je za tímto účelem užíván standard MIX. PREMIS obsahuje i sekci signifikantních vlastností, která však není v praxi příliš užívána.

1.3.3 METS

Standard METS (*Metadata Encoding & Transmission Standard*) slouží primárně jako metadata zaznamenávající informace o zabalení modelu OAIS (tedy o zabalení balíčků SIP, AIP a DIP). Především umožňuje vnoření dalších metadatových schémat pro popis archivačních a interpretačních informací (a tím jejich identifikaci). Dále obsahuje sekci určenou pro zápis provenienčních informací (formou vnoření externího schématu) a zápis některých interpretačních informací (např. o chování objektu). V praxi se METS užívá

³⁶ https://www.ifla.org/files/assets/cataloguing/frbr/frbr_2008.pdf

³⁷ V rámci elementu „objectCharacteristicsExtension“.

zejména pro první funkci (záznam informací o zabalení) a také jako datový formát. Jeho sekce strukturálních map se využívá pro záznam informací o všech obrazových souborech (fyzická mapa) a jejich posloupnosti (logická mapa).

Tyto informace tedy netvoří strukturální interpretační informace, ale vlastní datovou součást digitalizátu knihy, bez níž by objekt CDO byl neúplný. V praxi je METS možno užít v kombinaci se standardem PREMIS, přičemž lze zvolit několik způsobů implementace. Americké směrnice NISO doporučují zaznamenat PREMIS do sekce METS pro zápis provenienčních informací (NISO, 2007, s.55).

1.3.4 MODS

MODS, neboli (*Metadata Object Description Schema*), je metadatový standard široce užívaný pro zápis deskriptivních metadat. Jde o výsledek projektu zaměřeného na vývoj standardu pro popis jakéhokoliv typu dokumentu a správu digitálních objektů v jazyce XML, který vedlo oddělení Network Development and MARC Standards Office, jež je součástí Kongresové knihovny. MODS vychází z katalogizačního standardu MARC 21, je ale jednodušší a snadno čitelný pro lidského uživatele. První verze MODS 1.2 byla zveřejněna v lednu 2002, aktuální verze 3.8 vyšla v září 2022. Vývoj standardu stále řídí Network Development and MARC Standards Office ve spolupráci s mezinárodní komunitou vývojářů, přičemž návrhy na vylepšení a nové prvky standardu probíhají i pomocí celosvětové emailové konference.

Vzhledem k historii vzniku umožňuje MODS téměř plnou konverzi záznamů v MARC 21 do metadatového zápisu v MODS pomocí převodní šablony MARCXML s minimální ztrátou informací. Formáty se však nepřevádí v poměru 1:1, jelikož slovník MODS je vyjádřen slovními značkami (elementy), na rozdíl od číselných značek MARC. Sada elementů MODS (*MODS Element Set*) ovšem umožňuje i vytváření kompletních originálních záznamů, nikoliv pouze konverze z existujících katalogizačních záznamů. Provázanost s MARC 21 je pak i v rozdílném vytváření záznamů a užitých hodnot v návaznosti na katalogizační pravidla, ve kterých byla zpracována předloha (AACR2 nebo RDA).

Elementová sada MODS obsahuje dvacet kontejnerových elementů (*top elements*), které jsou zpravidla dále každý doplněny sadou vlastních podřízených elementů (*subelements*). Každý element může být specifikován atributy, které konkretizují typ

vyplněné hodnoty. Pokud se MODS použije v kombinaci s METS, je možné zápis rozčlenit do hierarchických úrovní, odpovídajících např. vnitřnímu členění dokumentu. Schéma MODS také počítá s doplněním elementů ze sad jiných standardů pomocí kontejnerového top elementu <mods:extension>, díky čemuž je možné v rámci jednoho zápisu popsat i specifické dokumenty nebo specifické informace (např. technického rázu), pro které MODS nemá ve vlastním element setu vhodné vyjádření.

1.3.5 Dublin Core

Dublin Core Metadata Element Set vznikl původně jako soupis patnácti klíčových vlastností, kterými lze popsat libovolný digitální objekt včetně webových stránek. Těmito klíčovými vlastnostmi byly *contributor* (příspěvatel), *coverage* (pokrytí/rozsah), *creator* (tvůrce), *date* (datum), *description* (popis), *format* (formát), *identifier* (identifikátor), *language* (jazyk), *relation* (vztah), *publisher* (vydavatel), *rights* (práva), *source* (zdroj), *subject* (předmět), *title* (název) a *type* (typ).

Název standardu je odvozen od města Dublin ve státě Ohio, ve kterém se v roce 1995 konal OCLC/NCSA Metadata Workshop, na kterém bylo schéma vytvořeno.

Dublin Core (dále jako „DC“) byl formálně standardizován normami ISO 15836, ANSI/NISO Z39.85, a IETF RFC 5013. Jako metadatový standard je uznáván od roku 2002, kdy vznikla formální dokumentace DCMI Metadata Terms.³⁸ Standard v současné době spravuje iniciativa Dublin Core Metadata Initiative, která funguje na principu placeného členství.

Standard DC bylo původně možné rozdělit do dvou verzí. Základní sadu Jednoduchého DC (*Simple DC*) o patnácti elementech doplňuje rozšíření Kvalifikovaného DC (*Qualified DC*). Kvalifikovaná verze obsahuje navíc tři další elementy - *audience*, *provenance* a *rightsHolder*. Od r. 2012 Byly tyto dvě verze sjednoceny do jednotného slovníku DCMI Metadata Terms.

Pomocí elementů DC lze univerzálně popsat široké množství digitálních objektů od textových, přes obrazové, zvukové, audiovizuální až po webové stránky. Kromě popisu lze formát využít také jako klíč k propojení jiných metadatových standardů, respektive jejich

³⁸ <https://www.dublincore.org/specifications/dublin-core/dcmi-terms/>

elementových sad. Na rozdíl od jiných metadatových standardů nemá předepsaný syntax ani hierarchii mezi elementy. Díky svojí jednoduchosti slouží jako nástroj interoperability mezi standardy např. v oblasti linked data nebo v rámci sémantického webu. Některé metadatové standardy (konkrétně MODS) mají pro konverzi z jednoho formátu do druhého vlastní mapování.

1.3.6 MIX

Standard MIX (*Metadata for Images in XML Standard*) je XML schéma, které je založeno na americké normě ANSI/NISO Z39.87-2006. Podle vlastního popisu je účelem normy „standardizovaná sada metadatových elementů pro rastrová obrazová data“, přičemž tyto elementy „dokumentují digitální obrazová data vytvořená digitální fotografií nebo skenováním a též data, která byla pozměněna editováním nebo obrazovým převodem“ (ANSI/NISO, 2006, s. 1). Standard MIX obsahuje elementy této normy, přidává několik dalších (např. rozděluje prostorové rozlišení do dvou elementů) a snižuje povinnost vyplnění elementů. Podle své vlastní definice MIX vznikl jako formát pro výměnu nebo uložení dat specifikovaných v uvedené normě NISO. Standard MIX se v praxi užívá pro záznam obrazových vlastností digitalizátů (tedy dalších typů interpretačních informací), a to jako externí schéma pro PREMIS.

Norma ANSI/NISO Z39.87-2006 uvádí, že není určena pro záznam provenience (ANSI/NISO, 2006, s. 1). Kupodivu to není tak docela pravda vzhledem k tomu, že jedna její sekce elementů („Change History“) je určena pro záznam provenienčních informací z doby produkce (pro záznam generací dat vzniklých při vytváření finálních produkčních dat i užitých procesů), ale v praxi se za tímto účelem užívají spíše elementy standardu PREMIS, ačkoliv je PREMIS primárně určen pro záznamy operací v archivu, nikoliv pro digitalizaci.

1.3.7 Metadata v obrazových souborech

Speciální oblastí je zabudování metadat do obrazových souborů. Široce rozšířeným a obrazovým průmyslem podporovaným standardem pro zabudovaná metadata je EXIF.³⁹ Dva hlavní snímkové formáty (TIFF a RAW) podporují záznam metadat ve formátu EXIF.

³⁹ http://www.cipa.jp/std/documents/e/DC-008-2012_E.pdf

Tato metadata do souborů zapíše snímací zařízení (ať již skener, nebo fotoaparát) automaticky. Obsahují velké množství elementů popisujících mj. snímací zařízení (včetně sériového čísla), způsob snímání nebo obrazové nastavení. Formát JP2 záznam EXIF metadat neumožňuje. Umožňuje však vnořit libovolná metadata v XML (do strukturálního prvku „XML Box“, jehož volba je součástí profilu tohoto formátu). V praxi se do něj zapisují například bibliografické údaje (Library of Congress, 2006).

2 Specifická část

2.1 Standard NDK a související předpisy

Národní knihovna ČR byla spolu s Moravskou zemskou knihovnou v letech 2009-2014 řešitelem projektu Vytvoření Národní digitální knihovny (projekt NDK), který byl financován z Integrovaného operačního programu EU programového období 2007-2014. NK ČR se v souvislosti s tímto projektem začala poprvé soustavněji věnovat digitální archivaci. Výstupy projektu byly tři: digitalizáty tištěných dokumentů, jejich archivace v archivu a zpřístupnění uživatelům. V rámci projektu byl vybudován archiv pod názvem LTP úložiště NK ČR. V listopadu 2011 byla vydána nová zřizovací listina NK ČR, v jejímž článku II. ods. 2 d)⁴⁰ je již explicitně uveden závazek digitální archivace i odkaz na koncept důvěryhodného digitálního repozitáře. Fáze udržitelnosti projektu NDK skončila formálně v roce 2019, společná digitalizační linka NDK i LTP úložiště NK ČR jsou však ve spolupráci obou institucí nadále využívány a rozvíjeny. Současně v rámci projektu NDK došlo k adopci Standardu NDK na celorepublikové úrovni pro digitalizační projekty realizované z dotačního programu VISK 7 (od roku 2012) i v řadě dalších digitalizačních projektů a aktivit realizovaných zejména specializovanými a krajskými knihovnami.

Standard NDK byl vytvořen v souvislosti s realizací projektu NDK. Stal se závaznou specifikací pro digitalizáty tištěných monografií a periodik vytvořených Národní knihovnou ČR a Moravskou zemskou knihovnou v rámci tohoto projektu a je nadále rozvíjen i po skončení fáze udržitelnosti. Účelem Standardu NDK bylo vytvářet balíčky SIP tak, aby objekt CDO byl již v archivačním formátu a LTP úložiště NK ČR nemuselo provádět formátovou normalizaci. Metadatový profil byl navržen tak, aby umožnil zaznamenání všech důležitých informací, kterou popisují proces produkce digitalizátů.

V rámci Národní knihovny má odbornou a kurátorskou stránku digitální archivace na starosti specializovaný odbor - ONDS (Odbor novodobých digitálních sbírek, dříve pod názvem ODIF, Odbor digitálních fondů).⁴¹ Tehdejší ODIF při přípravě projektu NDK zavedl nové postupy ve třech klíčových oblastech digitální archivace (metadata, trvalé identifikátory a datové formáty). Pro metadata vytvořil Standard NDK (tj. metadatový aplikační profil pro

⁴⁰ Ve zřizovací listině je přesně uvedeno, že NK ČR: „Formuluje strategie a postupy dlouhodobé ochrany elektronických dokumentů a provozuje důvěryhodné digitální úložiště.“ Viz: <https://www.nkp.cz/soubory/ostatni/zrizovaci-listina-nk.pdf>

⁴¹ Na počátku nesl název Odbor digitální ochrany.

digitalizaci).⁴²

Vývoj a údržbu Standardu NDK má v gesci Oddělení standardů digitálních sbírek (dříve Oddělení pro standardy), které jej vyvíjí ve spolupráci s odbornou komunitou, sdruženou do Pracovních skupin pro digitalizaci textových dokumentů, audiovizuální dokumenty a elektronické dokumenty. Jako poradní orgán při vývoji Standardu NDK působí Formátový výbor, složený ze zástupců Národní knihovny a reprezentantů tří stávajících Pracovních skupin. Každý cyklus vývoje a aktualizace jednotlivých Definic metadatových formátů⁴³ je spoluvyvíjen a připomínkován členy Pracovních skupin, a následně prohlasován jak interně příslušnou Pracovní skupinou, tak Formátovým výborem při NDK. Návrhy k aktualizaci Standardu NDK vycházejí nejen z potřeby držet použité mezinárodní metadatové formáty v aktuální podobě v souladu s celosvětově platným standardem, ale také z potřeb digitalizační linky NDK a jiných digitalizačních institucí.

Návrhy na úpravu tak vycházejí nejen z práce jednotlivých útvarů Národní knihovny (např. Odboru digitalizace a knihovních systémů, Odboru zpracování fondů, Odboru historických a hudebních fondů, Slovanské knihovny), ale také prostřednictvím členů Pracovních skupin a též velkou měrou vycházejí z potřeb a požadavků uživatelů Standardu NDK z jiných paměťových institucí. Při vývoji Definic metadatových formátů pro speciální typy dokumentů pak Oddělení standardů digitálních sbírek spolupracuje kromě útvarů Národní knihovny i s jinými paměťovými institucemi, které s daným typem dokumentu pracují, či se na něj specializují (např. České muzeum hudby a Národní technická knihovna pro zvukové dokumenty, Knihovna Akademie věd pro e-born dokumenty, atd.).

Standard NDK slouží jako archivační formát metadat nejen v LTP systému digitalizace NDK, ale je také základním podporovaným formátem pro ukládání dat do open source LTP systému ARCLib⁴⁴, pro tvorbu dat v produkčním digitalizačním nástroji ProArc⁴⁵ a pro zpřístupňování dat v open source digitální knihovně Kramerius, na jejichž rozvoji se Národní knihovna spolupodílí.

⁴² <http://www.ndk.cz/standards-digitalizace/metadata>

⁴³ Jedná se o konkrétní metadatové standardy, které jsou vydávány pro určitý typ digitalizovaného dokumentu.

⁴⁴ <https://arclib.cz/>

⁴⁵ <https://github.com/proarc>

2.1.1 Metadatový aplikační profil

V digitalizační praxi si paměťové instituce pro konkrétní projekt vytvářejí tzv. metadatový aplikační profil. Koncept tohoto profilu je založen na ideji, že pro konkrétní kontext je nutno metadatové standardy lokalizovat a optimalizovat (Zeng, 2016, s. 54). Metadatový aplikační profil je soubor metadatových elementů, které jsou vybrány z jednoho nebo více mezinárodních metadatových standardů a jsou spojeny do jednoho sloučeného schématu, který je uzpůsoben na míru funkčním požadavkům konkrétního užití, zatímco je zachována interoperabilita s původními mezinárodními standardy (Duval, 2002). Součástí profilu může být i vypracování vlastních metadatových elementů, v praxi digitalizátů tištěných dokumentů to však není obvyklé. Na webové stránce Kongresové knihovny je možno nalézt ukázky metadatových profilů řady světových paměťových institucí (včetně NK ČR).⁴⁶ Metadatový profil se může stát národním standardem, jak se tomu stalo v případě Standardu NDK pro oblast digitalizace knihovnických dokumentů v ČR.

Metadatový aplikační profil pro digitalizaci slouží k tomu, aby byly zachyceny zejména provenienční informace o původu a historii změn digitalizátu v průběhu jeho vytváření, počínaje snímáním a konče finalizací balíčku SIP. Jde zejména o operace, které byly vykonány (události standardu PREMIS) a informace o všech generacích obrazových dat. Pokud nejsou tyto informace zachyceny v průběhu digitalizace, jejich pozdější zjišťování archivem může být obtížné nebo přímo nemožné. Rovněž je důležité, aby balíček SIP obsahoval kvalitní identifikační informace, v praxi jde především o bibliografická metadata a perzistentní identifikátory. Garantem kvality těchto metadat by měla být vždy digitalizující knihovna – předlohy jsou popsány v jejím katalogizačním systému. Tyto záznamy jsou předmětem konverze do formátu MODS.

2.1.2 Standardy pro metadata

Standard NKD je zastřešující pojem, který zahrnuje všechny jednotlivé metadatové standardy, které vydává Národní knihovna ČR. Jejich oficiální název zní Definice metadatových formátů (dále DMF). Historicky nejstarší jsou DMF pro digitalizaci textových monografií a periodik. Jako součást projektu NDK byly sice vytvořeny DMF i pro další typy dokumentů (např. pro zvukové dokumenty nebo e-born dokumenty), tato metodika se však zaměřuje jen na textové digitalizáty, které činí v české praxi největší objem digitalizovaných

⁴⁶ <http://www.loc.gov/standards/mets/mets-registered-profiles.html>

dokumentů.

Hlavním obsahem Standardu NDK je metadatový profil, jehož základem jsou mezinárodní standardy METS a PREMIS, doplněné o MODS, DC, MIX, ALTO a CopyrightMD (a v případě zvukových dokumentů také o specifický standard AES57). Označení “Standard NDK” zde znamená, že se jedná o metadatový profil, který je českým národním standardem pro knihovny.

Tento metadatový profil je založen na užití mezinárodních metadatových standardů, které jsou mu nadřazené. Do metadatového profilu Standardu NDK byly vybrány převážně mezinárodně široce aplikované formáty, které jsou dále vyvíjeny a aktualizovány mezinárodní uživatelskou komunitou z paměťových institucí a které disponují detailní specifikací a metadatovým schématem. Některé (MODS, METS, ALTO) jsou zaštitěny Kongresovou knihovnou. Tyto vlastnosti zajišťují aktuálnost metadatových profilů a interoperabilitu dat.

Do českého národního standardu (tj. do Standardu NDK) byly z každého metadatového schématu vybrány jejich všeobecné principy a elementová sada, která odpovídá katalogizačnímu popisu v ČR (v případě popisných metadat) a do které lze zaznamenat klíčové informace z tvorby digitálních objektů při vytváření SIP balíků. Mezinárodní standardy tedy byly adaptovány v míře, která odpovídá potřebám a realitám českého knihovního prostředí.

Dodržování Standardu NDK není povinné všechny knihovny a paměťové instituce v ČR, povinnost je dána pouze pro ty, které svá data ukládají do LTP úložiště NK ČR, tedy především instituce, které digitalizaci realizují pomocí dotačního programu VISK 7.

Jednotlivé DMF jsou metadatové profily založené na více souborech METS XML v balíčku SIP (jeden pro bibliografická metadata, ostatní vytvářené zvlášť pro každý soubor obrazové matrice v archivním formátu) (Národní knihovna, 2023b). V jiných aplikacích v zahraničí tomu bývá často jinak. V SIP balíčku vždy existuje pouze jeden soubor METS XML (také označovaný jako „hlavní mets“) (Národní knihovna, 2023b).

Součástí Standardu NDK je i užívání perzistentních identifikátorů založených na mezinárodním standardu URN:NBN (Uniform Resource Name: National Bibliography Number). Českou národní implementací standardu URN:NBN je identifikační systém ČIDLO (CZIDLO – CZech IDentification and LOcalization tool), který byl vytvořen v Odboru digitálních fondů NK ČR v letech 2011-2013 a který je dále vyvíjen. Hlavní funkcí

systemu ČIDLO je přidělování identifikátorů URN:NBN a zabezpečování trvalé identifikace dokumentů českého kulturního dědictví. V současnosti systém ČIDLO přidělil digitalizovaným knihám a číslům periodik více než 1,5 milionu identifikátorů URN:NBN.⁴⁷ ČIDLO však slouží nejen knihovnám a paměťovým institucím v ČR, ale i samotným uživatelům, a to jako prostředek pro zajištění trvalého zpřístupňování digitálních dokumentů v internetové síti (řeší problém s nestabilitou URL adres) a pro zajištění důvěryhodnosti citační praxe (řeší problém s ověřováním autenticity citovaných dokumentů).

V roce 2014 byla existující pravidla systému ČIDLO zapracována do „Metodiky pro přidělování a správu životního cyklu unikátních perzistentních identifikátorů digitálních dokumentů podle standardu URN:NBN“, která byla 5. června 2015 uznána Ministerstvem kultury ČR jako certifikovaná metodika. Vývoj systému ČIDLO pokračoval dál, podobně i nové technologické postupy a softwarová řešení, a v roce 2018 byla vydána nová aktualizovaná Metodika, certifikovaná Ministerstvem kultury ČR. Tato metodika spolu s metodikou pro vytváření balíčků SIP představují ucelený komplexní návod, jakým způsobem využívat systém ČIDLO.

Metadatové standardy DMF jsou pak doplněny Pravidly popisu pro digitalizaci, které přibližují specifika praktického popisu digitalizátů a jejich pravidel. Vychází z pravidel knihovnického popisu dokumentů, katalogizačních pravidel a příkladů dobré praxe. Zabývají se také pravidly pro doplňování různých popisných hodnot v rámci metadatového popisu (například číslování stran, datování výtisků či doplňování typů stran). V současné době NK spravuje Pravidla pro popis digitalizátů monografií, periodik a zvukových dokumentů.

2.1.3 Standardy pro formáty

Hlavním typem digitálních dat, se kterými knihovní digitalizace dlouhodobě pracuje, jsou rastrová obrazová data (tvořená především digitalizací tištěných dokumentů). Jako archivační formát byl před projektem NDK užíván JPEG a jako prezentační formát DjVu. Při přípravě projektu bylo třeba vhodný formát znovu zvážit s ohledem na výši investice, masový záběr plánované digitalizace a vývoj v oblasti digitální archivace. Již tehdy bylo také zřejmé, že formát DjVu jako prezentační formát (jako archivační nebyl zvažován nikdy) již zastaral.⁴⁸ Po interní analýze byl zvolen JP2 jako nový archivační i prezentační formát pro digitalizované (novodobé) tištěné dokumenty NK ČR. Jako hlavní výhody byly označeny

⁴⁷ https://standardy.ndk.cz/ndk/archivace/Certifik_metodika_urnnbn_2018.pdf

⁴⁸ Jedním z důvodů byla prohra souborů s jeho hlavním konkurentem, formátem PDF.

otevřená dokumentace, neproprietárnost, kompresní možnosti a využitelnost pro archivaci i zpřístupnění. Specifikem projektu NDK byla možnost využít matematicky bezztrátovou kompresi pro archivační formát. Celá řada institucí, která v současnosti přechází nebo se chystá přecházet na JP2 jako archivační formát (konverzí z formátu TIFF), totiž volí matematicky ztrátovou kompresi, a to z důvodu významné úspory úložných kapacit. I když jde o vizuálně bezztrátovou kompresi (tedy uživatel nic nepozná), ztrátově komprimovaný formát je vždy rizikem pro budoucí migrace. Pro prezentaci byla v Národní knihovně ČR i Moravské zemské knihovně zvolena implementační varianta s image serverem. Profil JP2 pro zpřístupnění je prezentační meziformát pro vytváření dočasných obrázků ve formátu JPEG pro uživatele při jeho procházení digitální knihovnou. Tím byl vyřešen problém s nutností pluginu. Archivační formát je v JP2 s profilem v bezztrátové kompresi, prezentační meziformát v JP2 s profilem ve ztrátové kompresi.

Pro OCR komponentu (*Optical Character Recognition*), tedy pro optické rozpoznávání znaků digitalizovaných tištěných textů, byl zvolen standardizovaný formát ALTO XML. Standard ALTO umožňuje ukládat informace potřebné k popisu vzhledu a obsahu textového dokumentu, které byly získány pomocí aplikace OCR. V metadatovém profilu NDK se doporučuje, aby všechny úpravy obrazů (které vedou ke změně rozměrů obrazů nebo rozlišení apod.), se udělaly ještě předtím, než se vytvoří OCR. Rovněž je třeba zachovat velikost obrazu uživatelských i archivních kopií stejnou (tj. počet pixelů, rozlišení), tak aby ALTO XML odpovídalo. Souběžně s ALTO XML je výstup z OCR ukládán zároveň ve formě jednoduchého neformátovaného ASCII textu ve formátu TXT.

2.1.4 Standardy pro obrazová data

Vedle předpisu konkrétního formátu a jeho profilu, obsahuje Standard NDK pro data vzniklá v podprogramu VISK 7 další doporučení pro tvorbu obrazových dat, kterými jsou:⁴⁹

- ořez dokumentů cca 1 mm vně okraje dokumentu
- narovnání podle řádků textu
- veškeré úpravy obrazů je třeba provádět na archivních souborech
- uživatelský soubor se bude generovat až po všech úpravách
- uživatelská i archivní kopie musí mít stejný rozměr (v pixelech) a stejné

⁴⁹ <http://www.ndk.cz/standardy-digitalizace/standardy-pro-obrazova-data>

rozlišení (v DPI)

- skenování v rozlišení minimálně 300 PPI
- barevná hloubka 24 bitů
- barevný model RGB

2.1.5 Podmínky VISK 7 pro rok 2023

Program VISK 7 je určen pro digitalizaci originálních tištěných novodobých dokumentů za účelem jejich ochrany a zpřístupnění širokému okruhu uživatelů napříč českými knihovnami. Zpřístupnění dokumentů je realizováno pomocí digitální knihovny Kramerius.

Zaměření programu se váže k ochraně a zpřístupnění dokumentů tištěných na kyselém papíru, dokumentů ve špatném fyzickém stavu a dokumentů jinak fyzicky ohrožených (viz Koncepce rozvoje knihoven v ČR na léta 2021-2027). Do programu digitalizace lze zahrnout převážně bohemikální dokumenty vyprodukované v českých zemích, případně zahraniční publikace bohemikálního charakteru od roku 1800 do současnosti.

Podmínky pro VISK 7 jsou vydávány každý rok. Na jejich vydávání se podílí Odbor novodobých digitálních sbírek NK ČR (dříve odbor Digitálních fondů). Žadatelé o VISK 7 se musejí řídit Standardem NDK, který je pro VISK 7 závazný, a dalšími doporučeními, která se týkají zejména způsobu dodávání balíčků SIP.

V podmínkách VISK 7 pro rok 2023 jsou specificky zdůrazněny následující podmínky, které souvisí s vyplňováním metadat podle Standardu NDK. Žadatel musí pro digitalizované dokumenty získat tyto identifikátory:

- platné číslo České národní bibliografie
- identifikátor URN:NBN (přes systém ČIDLO)
- u periodických dokumentů číslo ISSN (pokud u starších periodických dokumentů dosud nebylo ISSN přiděleno, musí zajistit jeho přidělení u České národní agentury ISSN při Národní technické knihovně)
- sigla instituce (pokud dosud nebyla přidělena, musí o ní požádat prostřednictvím Národní knihovny ČR)

Kromě získání výše uvedených identifikátorů je žadatel povinen splnit i následující

požadavky:

- musí být specifikován formát digitalizovaných titulů a úroveň zpracování metadat
- popis způsobu zpřístupňování digitalizovaných kopií
- popis způsobu uložení digitálních dat v instituci a předávání dat do úložiště NK ČR
- u dokumentů musí být ověřeno pomocí Registru digitalizace, zda již nebyly zpracovány, pokud ne, je nutné je v Registru zaevidovat ještě před podáním projektu

Do roku 2016 byl v podmínkách VISK 7 uveden i požadavek: „Konverze obrazových souborů pomocí OCR do textového formátu s úspěšností rozpoznávání min. 95%“ (Česko, 2015). Od roku 2017 již není uplatňován, neboť z objektivních důvodů nelze dané úspěšnosti dosáhnout u některých typů písma či dokumentů ve špatném fyzickém stavu.

II. IMPLEMENTAČNÍ ČÁST

Úvod

Uvedená doporučení jsou určena pro knihovny, případně jiné paměťové instituce, které vytvářejí balíčky SIP obsahující digitalizáty tištěných dokumentů, které jsou určeny k dlouhodobému uchovávání (digitální archivaci) v repozitáři. Doporučení se zaměřují především na obecnější procedurální postupy, konkrétnější postupy jsou uvedeny pouze v souvislosti s plněním technických metadat.

Metodika obsahuje pouze doporučení, která nejsou závazná, a snaží se navrhnout optimální postup, kterého nemusí být v současné praxi vždy možné dosáhnout, nicméně předpokladem je, že v budoucnosti to možné bude.

Metodika je určena jak pro knihovny, které digitalizují podle aktuálních DMF pro digitalizáty periodik a monografií (tj. metadatových profilů vydávaných Národní knihovnou ČR), tak pro digitalizace podle jiných metadatových profilů, pokud zachovávají některé základní rysy Standardu NDK.

Ve vztahu ke Standardu NDK tato metodika navrhuje postupy, které jsou nad rámec toho, co je povinné v rámci konkrétních DMF, resp. specifikuje postupy, které jednotlivá DMF nepopisují, ale předpokládají. V případě rozporu této metodiky s podmínkami VISK 7 platí, že podmínky tohoto podprogramu jsou závazné pro knihovny, které z něj získávají dotaci na svou digitalizaci.

3 Plánování digitalizačního projektu

3.1 Digitalizační projekt

Digitalizací se zde rozumí specificky digitalizace tištěných dokumentů do podoby rastrových dat. Tento typ digitalizace směrnice FADGI definuje jako „konverze analogových barevných a jasových hodnot do nespojitých číselných hodnot. Číslo nebo množina čísel označuje barvu a jas každého pixelu v rastrovém obrázku“ (FADGI, 2010, s. 44). Výstupem této digitalizace a předmětem následného uchování a zpřístupňování je digitalizát tištěného dokumentu, který je tvořen množinou souborů rastrových dat, jež reprezentují vizuální vlastnosti částí tištěné knihy (stránek, přebalu atd.) a jejich posloupnost, a který obvykle také zahrnuje textová data, jež jsou výstupem optického rozpoznávání znaků (OCR) v rastrových obrazech a která umožňují do určité míry pracovat s obrazy rovněž jako s textem (viz funkce plnotextového prohledávání digitalizátu).

Digitalizační projekt může mít různé cíle. Jedním z nejčastějších důvodů je záchrana fyzicky poškozených dokumentů a dokumentů na nestabilním médiu (typicky kyselý papír užívaný k tisku novin), kdy digitální reprezentace nahrazuje původní dokument. Dále je cílem ochrana původních dokumentů; tím, že uživatel knihovny používá ke studiu digitalizát, je jeho předloha ušetřena namáhání a poškozování. Cílem projektu může být i lepší zpřístupnění dokumentů ze sbírek širšímu okruhu uživatelů pomocí aplikace ke čtení nebo například zpřístupnění děl již nedostupných na trhu. Digitalizační projekty paměťových institucí mohou digitalizáty pouze shromažďovat a uchovávat ve formě *dark archive* (nepřístupného archivu), mohou je zpřístupňovat čtenářům v digitálních knihovnách, nebo mohou dělat obojí: uchovávat i zpřístupňovat.

Pokud je cílem výsledné digitalizáty dlouhodobě uchovávat i zpřístupňovat čtenářům, je optimálním řešením takové, aby byla digitalizace nastavena z hlediska potřeb digitální archivace. Pro digitální archivaci je klíčové, aby zahájení digitalizačního projektu předcházela prvotní specifikace balíčku SIP a aby byly stanoveny vhodné transparentní postupy (včetně specifických nástrojů) pro vytváření dat a metadat, které budou určitou zárukou toho, že vytváření balíčku SIP proběhlo tak, aby nebyla narušena autenticita informačního obsahu v balíčku SIP.

Specifikaci balíčku SIP pro kontext této metodiky poskytuje Standard NDK (resp. jednotlivá DMF) a určení vhodných postupů je pak vlastním cílem této metodiky. V praxi NDK jsou pak jednotlivé DMF doplněny o Pravidla popisu, která sumarizují doporučené postupy pro popisování jednotlivých typů dokumentů a výběr optimálních hodnot do zápisu deskriptivních metadat. Digitalizační projekt by měl být sepsán po obeznámení se se Standardem NDK a s touto metodikou, a teprve na základě těchto požadavků by měla být navržena jeho realizace.

Digitalizační projekt, po stanovení podoby balíčku SIP a postupů, musí zahrnovat řadu obecných specifikací technické a organizační povahy. Mezi ně patří: výběr vhodného snímacího zařízení, vytvoření digitalizačního pracoviště (včetně stanovení rolí a jejich obsazení lidskými zdroji), výběr testování softwarových nástrojů pro vytvoření digitalizátů v požadované podobě (včetně metadat), volba / vývoj digitalizačního systému pro řízení celého průběhu digitalizace, testování zařízení a softwaru a kontrola kvality (včetně postupů řešení chyb). Tyto otázky jsou již dobře popsány v existující odborné literatuře. Jako základní zdroj pro tuto problematiku lze doporučit směrnice americké iniciativy FADGI (Federal Agencies Digital Guidelines Initiative), které byly vytvořeny pro paměťové instituce v USA a které nesou název „Technical Guidelines for Digitizing Cultural Heritage Materials: Creation of Raster Image Files“. Směrnice vyšla v roce 2010,⁵⁰ v letech 2016⁵¹ a 2023⁵² pak byly vydány její revidované verze. Jako úvodní zdroj lze současně využít také českou „Metodiku pro vytváření bezpečnostních kopií archiválií v digitální podobě“ vytvořenou Národním archivem.⁵³

3.2 Technické zajištění

3.2.1 Snímací zařízení

Mezi prvotní rozhodnutí při přípravě digitalizačního projektu patří určení toho, zda bude snímání provedeno formou skenování (skener), nebo fotografování (digitální fotoaparát).⁵⁴ Následně je potřeba věnovat dostatečné úsilí výběru adekvátního skeneru nebo fotoaparátu, který dokáže vytvářet digitalizáty v předepsané podobě. Optimální variantou je vybraná snímací zařízení před jejich pořízením přímo otestovat.

⁵⁰ http://www.digitizationguidelines.gov/guidelines/FADGI_Still_Image-Tech_Guidelines_2010-08-24.pdf

⁵¹ http://www.digitizationguidelines.gov/guidelines/FADGI_Still_Image_Tech_Guidelines_2016.pdf

⁵²

https://www.digitizationguidelines.gov/guidelines/FADGI%20Technical%20Guidelines%20for%20Digitizing%20Cultural%20Heritage%20Materials_3rd%20Edition_05092023.pdf

⁵³ <https://www.nacr.cz/wp-content/uploads/2019/05/metodika2015.pdf>

⁵⁴ Tato metodika se blíže zaměřuje na skenování.

Mezi základní požadavky na výběr skeneru patří dostatečné prostorové rozlišení skeneru a možnost ukládat původní snímky v nekomprimované podobě.

Jako základní zdroj pro otázky výběru vhodných snímacích zařízení lze opět doporučit výše uvedenou směrnici „Technical Guidelines for Digitizing Cultural Heritage Materials: Creation of Raster Image Files“⁵⁵, a českou „Metodiku pro vytváření bezpečnostních kopií archiválií v digitální podobě“ (zejména kapitoly 4.1-4.7).

3.2.2 Softwarové nástroje pro tvorbu digitalizátů

Pro vytváření digitalizátů tištěných dokumentů je potřeba využívat řadu vysoce specializovaných nástrojů, zejména pro zpracování obrazových dat, formátové konverze nebo optické rozpoznávání znaků. Specializované nástroje pro tvorbu digitalizátů je potřeba vybrat již ve fázi přípravy digitalizačního projektu a provést testování, zda nástroje skutečně dokáží vytvářet digitalizáty v předepsané podobě. Toto testování se může opakovat, dokud nebude nalezen vhodný nástroj. Vzhledem k tomu, že některé kvalitní nástroje jsou komerční, musí výběr zohledňovat nejen kvalitu výstupů softwarových nástrojů, ale také finanční náklady spojené s jejich pořízením a užitím. Finanční náklady spojené se softwarovými nástroji musejí být důkladně propočítány a zahrnuty do rozpočtu digitalizačního projektu. Je špatnou praxí neověřovat si náklady spojené se specializovaným softwarem, který je nezbytný pro vytváření digitalizátů tištěných dokumentů, a teprve při zahájení digitalizace zjistit, že rozpočet projektu není dostatečný.

3.2.3 Validační nástroje

Pro kvalitní digitalizaci, která zohledňuje požadavky dlouhodobého uchování, jsou nezbytné nejen produkční softwarové nástroje, ale také specializované validační nástroje. Zejména důležitá je validace metadat, souborových formátů a balíčku SIP.

Ve fázi přípravy digitalizace je zapotřebí, aby zvolené validátory byly otestovány, zejména jejich zapojení do celkového digitalizačního systému. Je nutné, aby digitalizační systém dokázal integrovat tyto specializované nástroje. Pokud nedojde k testování validačních nástrojů, může se stát, že po zahájení digitalizace v daném projektu budou zjištěny dodatečné náklady na integraci validátorů, se kterými původní rozpočet projektu nepočítal.

⁵⁵ <http://www.digitizationguidelines.gov/guidelines/digitize-technical.html>

3.2.4 Kontrola předloh

Tištěné předlohy, které mají být zdigitalizovány v daném projektu, by měly být důkladně zkontrolovány z hlediska kvality a úplnosti před zahájením digitalizace. V případě, že součástí projektu musí být seznam konkrétních předloh, které mají být zdigitalizovány (a tedy existuje závazek uvedené tituly zdigitalizovat), měly by být tyto předlohy zkontrolovány ještě před podáním projektového záměru.

Pro digitalizaci by měl být vždy vybírán takový exemplář, který je úplný a nejzachovalejší (v případě digitalizace periodik by měl být kompletní minimálně ročník). Pokud knihovna nedisponuje žádným úplným exemplářem, případně jsou všechny její exempláře ve špatné kvalitě, měla by zvážit, zda je nutné, aby takovou předlohu digitalizovala – zpravidla existuje možnost přenechat digitalizaci daného titulu monografie nebo ročníku periodika jiné knihovně. Další variantou je zapůjčit úplný a kvalitní exemplář z jiné knihovny. Pokud to není nezbytné, nedoporučuje se kompletovat jeden titul monografie nebo čísla periodika z více exemplářů.

V případě ověření nedostupnosti úplného exempláře je kompletace z více nekompletních svazků postupem možným a žádoucím. V takovém případě je nicméně potřeba velmi pečlivě ověřovat shodnost vydání všech použitých předloh. Zejména u periodik deníkové povahy je navíc vhodné ověřit, zda náhradní svazek neobsahuje jinou regionální mutaci či mutační vydání. Skutečnost o použití více exemplářů je zároveň doporučeno zaznamenat v poznámce o fyzickém stavu exempláře v poli <mods:physicalDescription> a konkrétní použité předlohy identifikovat v poli <mods:location>.

3.3 Základní standardizační doporučení

3.3.1 Stanovení základní intelektuální entity a granularity

V praxi NDK se základní intelektuální entitou myslí smysluplný celek, jehož reprezentace je obsahem balíčku SIP. V případě monografií se jedná o svazek, u periodik o jednotlivé číslo, a například u zvukových dokumentů se za intelektuální entitu považuje zvukový dokument (což může být jak jeden nosič, obsahující skladbu/skladby, tak soubor nosičů, které byly vydány jako celek).

Standard NDK potom u každého typu dokumentu v jednotlivých Definicích metadatových formátů určuje jeho granularitu, tj. dělení základní intelektuální entity na nižší části, které zároveň slouží v metadatovém záznamu hlavního mets jako úrovně popisu. Tyto

úrovně popisu potom figurují jako jednotlivé položky hierarchického záznamu v logické strukturální mapě.

3.3.2 Složení SIP balíku podle Standardu NDK

Podle Standardu NDK se SIP balíček pro digitalizáty textových dokumentů stabilně skládá z následujících položek:

- 1) složka mastercopy, která obsahuje archivní snímky předlohy ve formátu TIFF
- 2) složka usercopy, která obsahuje uživatelskou kopii snímků ve formátu JPEG 2000 (tj. tu, která se uživateli zobrazí při čtení v digitální knihovně)
- 3) složka alto pro OCR soubory
- 4) složka txt pro jednoduchý přepis (plain text) obsahu digitalizátu
- 5) soubor hlavní METS, ve kterém se nachází bibliografické údaje o digitalizátu, strukturální mapy a další údaje o objektu CDO
- 6) složka vedlejších METS (amdsec) pro každou naskenovanou stranu, která obsahuje administrativní a technická metadata
- 7) soubor info, který slouží jako základní popis obsahu SIP balíku
- 8) soubor MD5 pro kontrolní součet

V případě jiných typů digitálních dokumentů, než jsou naskenované textové dokumenty se obsah SIP balíků liší a pro jejich vytváření je nutné se řídit příslušným dokumentem Definice metadatových formátů, které přihlíží k jejich specifikům.

Pro zvukové dokumenty SIP balík obsahuje navíc složky pro originální neupravené verze zvukových nahrávek (source_audio), dále archivní (mastercopy_audio) a uživatelskou (usercopy_audio) kopii zvukových nahrávek a složku pro katalogizační záznam (catalog_entry).

Pro digital-born dokumenty je SIP balík naopak stručnější. Obsahuje pouze složku original pro archivovaný dokument (nebo dokumenty), hlavní mets, soubor info a soubor MD5 pro kontrolní součet. V případě e-born skládaných periodik obsahuje SIP balík navíc ještě složku amdsec pro vedlejší mets záznamy jeho uživatelských kopií.

3.3.3 Perzistentní identifikátory tištěné předlohy

Pro propojení tištěných předloh s digitalizáty je klíčové užití perzistentních identifikátorů tištěné předlohy na úrovni titulu. Odpovídající identifikátory tohoto typu jsou ISBN (*International Standard Book Number*), ISMN (*International Standard Music Number*), ISSN (*International Standard Serials Number*) a čČNB (číslo České národní bibliografie), které by měly být zaznamenány do metadat digitalizátu tištěné předlohy, pokud to pravidla těchto identifikačních systémů umožňují. V případě ISBN, ISMN a ISSN jde o identifikátory řízené mezinárodními organizacemi prostřednictvím národních agentur.

Identifikátor ISBN je přidělován českým knihám od roku 1989. Řídí jej Národní agentura ISBN v ČR při Národní knihovně ČR.⁵⁶

Identifikátor ISMN je přidělován českým hudebninám od roku 1996.⁵⁷ Řídí jej Národní agentura ISMN v ČR při Národní knihovně ČR.⁵⁸

Identifikátorem ISSN jsou označována periodika v ČR (ČSSR) od 70. let 20. století. Přiděluje jej České národní středisko ISSN při Národní technické knihovně.⁵⁹

Ve všech případech je účast vydavatelů dobrovolná a neplatí, že veškerá česká produkce daných typů publikací od doby zavedení těchto identifikačních systémů do českého prostředí obsahuje některý z těchto identifikátorů. Identifikátory ISBN a ISMN jsou většinou uvedeny v samotném dokumentu (knize, hudebnině) a nelze je přidělovat zpětně (tj. po vydání).⁶⁰ Identifikátor ISSN lze přidělovat zpětně a o jeho přidělení může požádat nejen nakladatel periodika, ale i knihovna.

Identifikátor čČNB je užíván v ČR od roku 2010.⁶¹ Jedná se o kód národní bibliografie, pro který je ve formátu MARC 21 vyhrazeno pole 015 a který některé jiné národní knihovny přidělují již desítky let záznamům svých národních bibliografií. Je přidělován popisné jednotce České národní bibliografie, nelze jej tedy přidělit jednotlivým číslům periodika, jednotlivým svazkům vícesvazkového díla bez významných názvů části a monografické řadě/edici, která zahrnuje díla s vlastními názvy.⁶² Aby byl perzistentní, musejí být zachovány i neplatné

⁵⁶ <https://www.nkp.cz/sluzby/sluzby-pro/isbn-ismn-issn>

⁵⁷ <https://www.nkp.cz/soubory/ostatni/prirucka-ismn.pdf>

⁵⁸ <https://www.nkp.cz/sluzby/sluzby-pro/isbn-ismn-issn/oma#ISMN>

⁵⁹ <https://www.techlib.cz/cs/2844-ceske-narodni-stredisko-issn>

⁶⁰ Je-li dokument zaveden do národní bibliografie, pak mu jsou i zpětně přiděleny tyto identifikátory a jsou tedy součástí katalogizačních záznamů knihovny NK ČR.

⁶¹ <https://www.caslin.cz/caslin/spoluprace/sluzby-souborneho-katalogu-cr/cislo-cnb-v-sk-cr>

⁶² <http://www.registrdigitalizace.cz/rdcz/info/data/cnb>

identifikátory čČNB (zneplatněné např. po sloučení záznamu). Identifikátor čČNB by podle současných pravidel měly mít „veškeré publikované dokumenty vydané na území ČR od roku 1801 do současnosti.“ Na stránkách Souborného katalogu je dále uveden výčet dokumentů, které nemají nárok na přidělení. Jedná se hlavně o tzv. šedou literaturu, která není běžně dostupná na knižním trhu: zejména nepublikované vysokoškolské a habilitační práce, výroční zprávy institucí a škol, divadelní programy, propagační materiály, pozvánky, interní dokumenty apod. V případě pochybností se tento identifikátor nepřidělí.⁶³

Pokud tištěné periodikum nemá přidělený identifikátor ISSN, měla by knihovna požádat o jeho přidělení České národní středisko ISSN, a to s dostatečným předstihem před zahájením digitalizace. Skutečnost, zda má dané periodikum přidělen identifikátor ISSN, lze ověřit v národní databázi ISSN, kterou spravuje Národní technická knihovna.⁶⁴

Pokud pravidla pro přidělování identifikátoru čČNB umožňují přidělení tohoto identifikátoru⁶⁵ a tištěná předloha na úrovni záznamu titulu jej zatím nemá, měla by knihovna požádat o jeho přidělení Souborný katalog ČR.⁶⁶ I zde je třeba, aby se tak stalo s dostatečným předstihem před zahájením digitalizace. Přítomnost platného identifikátoru čČNB významně pomáhá předcházet duplicitní digitalizaci a zároveň poskytuje určitou garanci správně zpracovaného katalogizačního záznamu.⁶⁷

V případě digitalizací financovaných z podprogramu VISK 7 je přítomnost identifikátorů čČNB a ISSN povinná.

3.3.4 Perzistentní identifikátory digitalizátu

3.3.4.1 URN:NBN

Jako hlavní perzistentní identifikátor digitalizátu tištěných dokumentů na úrovni intelektuální entity se doporučuje užívat identifikátor URN:NBN. Přidělování daného identifikátoru v českém prostředí zajišťuje služba ČIDLO, která slouží knihovnám a dalším

⁶³ <https://www.caslin.cz/caslin/spoluprace/sluzby-souborneho-katalogu-cr/cislo-cnb-v-sk-cr>

⁶⁴ http://aleph.ntkcz.cz/F/?func=find-b-0&local_base=stk02

⁶⁵ <https://www.caslin.cz/caslin/spoluprace/sluzby-souborneho-katalogu-cr/cislo-cnb-v-sk-cr>

⁶⁶ <https://www.caslin.cz/caslin/spoluprace/sluzby-souborneho-katalogu-cr/jak-spravne-postupovat-nez-zacne-knihovna-digitalizovat-dokument>; viz také: <https://registrdigitalizace.cz/rdcz/info/data/ccnb>

⁶⁷ Ve společném digitalizačním projektu Národní knihovny ČR a Moravské zemské knihovny se ukázalo, že nepřítomnost identifikátoru čČNB či přítomnost neplatného identifikátoru čČNB může způsobovat problémy například při zpracování vícesvazkových monografií. Zatímco v jedné knihovně byl v čase digitalizace stejný titul zpracován "zdola" (každý díl měl vlastní záznam) v druhé byl popsán "shora" (existoval jeden společný záznam pro všechny části). Výsledné digitální kopie pak nebylo možné v digitální knihovně jednotně zobrazit a bylo nutné přistoupit k dodatečné opravě.

paměťovým institucím pro potřeby trvalé identifikace českého kulturního dědictví. Systém ČIDLO se skládá ze souboru pravidel, která jsou účastníci systému (tj. zapojené instituce a jejich zaměstnanci) povinni dodržovat, dále z identifikátorů URN:NBN platných pro český jmenný prostor (tj. z identifikátorů začínajících řetězcem „urn:nbn:cz“) a z několika technických podsystémů. Jedním z technických podsystémů je i online aplikace resolver.⁶⁸ Hlavní funkcí, kterou resolver poskytuje, je přesměrovávací služba, která zajistí přesměrování identifikátoru URN:NBN na aktuální URL adresu do digitální knihovny, ve které je identifikovaný dokument zpřístupňován, případně na metadatový záznam dokumentu v systému.⁶⁹

Identifikátor URN:NBN lze přidělit v případě digitalizovaných periodik článku, číslu a ročníku (titulu digitalizovaného periodika URN:NBN přidělit nelze), v případě digitalizovaných monografií svazku monografie (jednodílové nebo vícesvazkové) a vnitřní části monografie (kapitole). Také je možné přidělit identifikátor URN:NBN i vysokoškolským závěrečným pracím i dalším typům dokumentů, jako jsou kartografické dokumenty, grafiky, hudebniny nebo zvukové dokumenty. Pro odlišnou úroveň granularity musí být přidělen jiný identifikátor URN:NBN.

Podle aktuálních DMF pro digitalizáty tištěných dokumentů (verze 2.0 pro digitalizovaná periodika a verze 2.1 pro digitalizované monografie)⁷⁰ musí být identifikátor URN:NBN povinně přidělen základní intelektuální entitě (tj. číslu nebo ročníku periodika, resp. svazku monografie).

Podrobná pravidla pro přidělování identifikátorů URN:NBN i řízení celého jejich životního cyklu jsou obsažena v certifikované metodice popisující pravidla systému ČIDLO, která byla vydána v roce 2018.⁷¹ Metodika definuje nezbytné kroky k získání perzistentního identifikátoru URN:NBN a taktéž podrobná pravidla pro přidělení identifikátoru a následné dlouhodobé uložení a zpřístupnění digitálního dokumentu, kterému byl identifikátor přidělen. Účastníkem (registrátorem) systému ČIDLO se může stát jakákoliv knihovna nebo jiná instituce, která má v systému ADR (Centrální adresář knihoven a informačních institucí v ČR)⁷² Národní knihovny ČR přidělenou siglu. Podmínkou je dodržování pravidel systému. Registraci zařizuje kurátor systému ČIDLO. V případě, že pro knihovnu provádí digitalizaci externí

⁶⁸ <https://resolver.nkp.cz/>

⁶⁹ <https://standardy.ndk.cz/ndk/archivace/resolver-urn-nbn-sluzba-cidlo>

⁷⁰ <https://www.ndk.cz/standardy-digitalizace/metadata>

⁷¹ https://standardy.ndk.cz/ndk/archivace/Certifik_metodika_urnnbn_2018.pdf

⁷² <https://aleph.nkp.cz/cze/adr>

dodavatel, musí mít souhlas dané knihovny se zastupováním v systému ČIDLLO.

3.3.4.2 UUID

Kromě identifikátoru URN:NBN se v české digitalizační praxi užívá také další mezinárodní perzistentní identifikátor, UUID (*Universally unique identifier*, známý také jako GUID, *Globally unique identifier*). UUID má standardní délku 128 bitů, a je zpravidla reprezentován textovým řetězcem tvořeným 32 hexadecimálními znaky a 4 spojovníky (nejčastěji ve formátu 8-4-4-4-12, tedy xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxx). Původně byl vyvinut pro účely identifikace digitálních objektů v operačním systému Domain/OS počítačů Apollo v 80. letech. Později se uplatnil v Distribuovaném výpočetním prostředí (*Distributed computing environment, DCE*) projektu Open software foundation (OSF), který se inspiroval architekturou operačního systému počítačů Apollo. Později DCE adoptovala platforma Windows. Dnes je UUID součástí standardu Distribuovaného výpočetního prostředí Open software foundation pod ISO/IEC 11578:1996 (*Information technology – Open Systems Interconnection – Remote Procedure Call*) a také pod ITU-T Rec. X.667 | ISO/IEC 9834-8:2014.⁷³

V rámci Standardu NDK je identifikátor UUID povinně přidělován každé existující úrovni bibliografických metadat (od titulové úrovně až po úroveň strany či vnitřní části). Standard NDK nepředepisuje využití konkrétní verze UUID (existují verze 1 až 5), které se mezi sebou ostatně liší pouze způsobem, jakým se identifikátor vygeneruje, nikoliv výslednou podobu. V praxi digitalizace Národní knihovny je využíváno UUID v1, které spolehlivě zajišťuje unikátnost přiděleného identifikátoru využitím MAC adresy v kombinaci s časovým údajem v okamžiku přidělení. Jako vhodné lze doporučit rovněž UUID v4, kde dochází k vygenerování zcela náhodného řetězce znaků v celém rozsahu identifikátoru, pravděpodobnost přidělení duplicitního identifikátoru je rovněž v tomto případě zcela zanedbatelná.

3.3.5 Projektová dokumentace

Knihovna by měla k digitalizačnímu projektu vytvořit podrobnou dokumentaci, kterou by optimálně měla zpřístupňovat spolu s digitalizáty tištěných dokumentů. Tato dokumentace by měla přinejmenším obsahovat informace o užitých metadatových standardech (např. konkrétních metadatových standardech podle Standardu NDK včetně jejich verze) a zvolených

⁷³ <https://www.ietf.org/rfc/rfc4122.txt>

archivačních a prezentačních formátech a optimálně též informace o užitých validačních nástrojích.

V případě Národní knihovny je dokumentace z webových stránek⁷⁴ doplněna o repozitáře na vývojářské platformě GitHub,⁷⁵ které používá pro vydávání nových verzí validačních nástrojů a řešení uživatelských problémů. Archiv vláken zároveň slouží jako jednoduchá znalostní báze řešení a příkladů dobré praxe.

⁷⁴ <https://standards.ndk.cz/>

⁷⁵ <https://github.com/NLCR>

4 Digitalizace

4.1 Příprava bibliografických záznamů

Před zahájením vlastního snímání předlohy musí existovat její bibliografický záznam v katalogu knihovny, která digitalizuje. Tento záznam by také měl být důkladně zkontrolován. V případě nekvalitního nebo neúplného záznamu by měla být provedena rekatalogizace. Kontrolu záznamů (resp. rekatalogizaci) předloh určených k digitalizaci by měl vždy provádět školený a zkušený katalogizátor. Z toho důvodu by digitalizační projekt měl počítat s vyčleněním odpovídající pracovní síly (katalogizátorem na alespoň částečný úvazek). Nedoporučuje se kontrola pracovníky, kteří nejsou katalogizátory. Kvalita katalogizačního záznamu má zásadní vliv na kvalitu digitálních bibliografických metadat uložených v balíčků SIP a následně metadat zobrazených v digitální knihovně čtenářům. Šetření prostředků na katalogizátory při digitalizaci může významným způsobem snižovat kvalitu výstupů digitalizačního projektu. Také platí, že pozdější opravy bibliografických metadat (tj. opravy v repozitáři a digitální knihovně) jsou mnohem komplikovanější a technicky, personálně i finančně nákladnější. Případnou opravu je navíc vždy nutné provést současně nad archivním i uživatelským balíčkem.⁷⁶

Bibliografický záznam by již měl obsahovat identifikátory ISBN, resp. ISMN, ISSN nebo ČČNB, pokud lze některý z těchto identifikátorů na základě pravidel daných identifikačních systémů přidělit. Zejména u identifikátoru ČČNB je potřeba zkontrolovat také jeho platnost.

Pro převod bibliografických údajů do digitálních metadat (standarty MODS a DC) by měl být užit vlastní katalog knihovny, nikoliv jiné katalogy (souborný katalog, báze ČNB apod.), byť by tyto jiné katalogy obsahovaly agregované záznamy knihoven zapojených do digitalizace.

4.2 Snímání předloh

4.2.1 Věrnost reprodukce tištěné předloze

Digitalizát by měl být co nejvěrnější digitální reprodukcí tištěné předlohy. Pro účely specifikace věrnosti doporučujeme využít směrnici DLF (Digital Library Federation) nazvanou

⁷⁶ https://standarty.ndk.cz/ndk/archivace/Certifik_metodika_urnnbn_2018.pdf

„Benchmark for Faithful Digital Reproductions of Monographs and Serials“.⁷⁷ Směrnice uvádí, že cílem věrnosti je „přesně reprodukovat výchozí zdrojový dokument, s ohledem na jeho úplnost, vzhled původních stránek (včetně tonality a barvy) a správnou (tj. původní) posloupnost stránek“ (The digital library federation benchmark working group, 2002, s. 2). Směrnice DLF také uvádí několik zásad, které musejí splňovat obrazová data reprezentující tištěnou předlohu. Tyto zásady doporučujeme využít v maximální možné míře, pokud směřovatelné podmínky pro digitalizační projekt (např. pravidla podprogramu VISK 7 nebo Standard NDK) neuvádějí jinak.

4.2.2 Základní parametry pro skenování

Jako minimální parametry pro skenování doporučujeme prostorové rozlišení nejméně 300 PPI,⁷⁸ barevnou hloubku nejméně 24 bitů (tj. 3 x 8 bitů) a barevný model RGB. Tato doporučení jsou v souladu s aktuálním standardem NDK pro obrazová data.⁷⁹

Pokud je to časově možné, doporučujeme provést zběžný průzkum míry detailnosti předloh určených pro digitalizační projekt a v případě potřeby rozlišení zvýšit paušálně na 400 PPI,⁸⁰ případně i vyšší. Za tímto účelem lze užít tzv. Quality Index (QI) (viz doporučení AIIM TR26- 1993 Resolution as it Relates to Photographic and Electronic Imaging, původně určené pro mikrofilmy). Tento výpočet bere v úvahu velikost písmen⁸¹ předlohy a plánovanou kvalitu výsledného obrazu na stupnici od špatné kvality po kvalitu excelentní. Index udává kolik obrazových bodů (pixelů) je potřeba pro reprezentaci nejmenšího písmene ve zdrojovém textu. Barevné a šedé obrazy vyžadují nejméně 16 obrazových bodů (Quality Index =8) pro excelentní, detailní zobrazení nejmenšího písmene zdrojového textu, bitonální obrazy potřebují 24 obrazových bodů. Pro písmeno o velikosti 1 mm je tak ideální snímací rozlišení 400 PPI, písmeno pak bude reprezentováno 16 obrazovými body u obrazů v plných barvách a v odstínech šedi (podrobněji kapitola 4.2.3).

Pokud takový průzkum není časově možný, doporučujeme řídit se obecnými doporučeními, které jsou obsaženy ve směrnici FADGI.⁸² Tato směrnice pro různé kategorie dokumentů (rukopisy, další vzácné dokumenty, knihy, noviny, fotografie aj.) doporučuje

⁷⁷ <http://old.diglib.org/standards/bmarkfin.pdf>

⁷⁸ Rozlišení 300 PPI postačuje pro text s velikostí písmen 1,4mm.

⁷⁹ <https://www.ndk.cz/standardy-digitalizace/standardy-pro-obrazova-data>

⁸⁰ Dle směrnice FADGI je rozlišení 400 PPI použitelné i pro vzácné knihy.

⁸¹ V případě netextových materiálů (mapy, nákresy apod.) se výpočet odvíjí od šířky nejtenčí čáry, tahu.

⁸² <http://www.digitizationguidelines.gov/guidelines/digitize-technical.html>

optimální hodnotu rozlišení v závislosti na požadované výsledné kvalitě. Doporučujeme využít uváděné hodnoty rozlišení pro nejvyšší kategorii kvality („4 Star“).

4.2.3 Závislost rozlišení obrazu a typografických bodů písma

Typografický bod má jednotku označovanou v ČR malým písmenem „b“ nebo tečkou za číslem, v současné době však převládá označování z anglického point – pt. V různých zemích má typografický bod různou velikost. Přebírají dvě velikosti, a to přibližně 0,350 a 0,376 mm. Rovněž počítačové vyjádření jednoho bodu (*desktop publishing point*), vycházející z angloamerické měrné soustavy, je rovno 0,3528 mm. Pro účely stanovení optimálního rozlišení v závislosti na velikosti bodů písma tak můžeme drobné odchylky ignorovat.

V typografické praxi se nejčastěji využívají písma od velikosti 6 pt (*nonpareille*), určenou zejména pro rozsáhlé slovníky a vzorce, knižní a časopisecká produkce využívá nejčastěji písma o velikosti 8 pt (*petit*), 9 pt (*borgis*) nebo 10 pt (*garmond*). Zejména u dětských knih se můžeme setkat také s písmy ve velikostech 12 (*cicero*) nebo 14 (*střední*).

Společnost ABBYY, jejíž OCR software je využíván v digitalizační lince NDK, doporučuje optimální rozlišení obrazu v závislosti na velikosti písma. Podobná doporučení můžeme najít i u dalších výrobců OCR programů. Pro ideální rozpoznání znaků je pro písmo o velikosti 8 až 10 pt doporučeno skenovat alespoň na 300 DPI. Písmo velikosti 8 pt a menší je doporučeno skenovat na 400 až 600 DPI. Pro složitá písma jako jsou japonské či čínské je doporučeno rozlišení přibližně o 25 % větší než pro jednoduché jazyky (ruština, hebrejština, arabština aj.).

Kvalita OCR by měla být jedním z kritérií pro cílové rozlišení obrazu, ale nikoliv kritériem jediným. Pro vytvoření kvalitní digitální kopie nás zajímá především původní rozlišení, ve kterém byl originál knihy vytištěn. Respektovat přitom musíme optické rozlišení skenovací soustavy a fyzické rozlišení prvku, na kterém zaznamenáme obraz knihy. Nikdy bychom pro archivní účely neměli rozlišení obrazu softwarově interpolovat a pokud dojde k interpolaci pro potřeby OCR, neměly by takové soubory být archivovány, ale použity maximálně jako dočasné pracovní kopie právě pro vytvoření OCR.

Projekt Národní digitální knihovny má jako minimální standard uvedeno skenování na 300 PPI (pro účely skenování vycházíme ze zjednodušeného převodu 1 PPI = 1 DPI). Vzhledem k uvedenému by měly být všechny knihy s významnou částí textu písma 8 pt a menšího skenovány na skenerech nebo fotoaparátech s optickým a fyzickým rozlišením minimálně 400

PPI. U ostatních knih s textem s jednoznačně převládajícím písmem nad 8 pt postačuje minimálně 300 PPI. U fotografických předloh a map je sledováno rozlišení originálu knihy a digitální kopie by měly být vytvářeny ideálně ve stejném rozlišení. Vzhledem k technickým možnostem skenovacích zařízení toho není možné vždy dosáhnout. Proto jsou fotografie vytištěné na křídovém papíře a mapy s vrstevnicemi nebo jinými podrobnostmi digitalizovány alespoň v rozlišení 400 PPI. Pokud má instituce možnost využít k digitalizaci zařízení s fyzickým rozlišením 600 PPI, doporučujeme jej pro skenování takových předloh využít.

4.2.4 Snímkový formát

Jako snímkový formát obrazových dat (formát původních snímků) by měl být vždy zvolen nekomprimovaný formát. V případě skenování by měl být užit formát TIFF, verze 6, bez komprese; v případě fotografování formát rodiny RAW (například fotoaparáty firmy Canon využívají formát Canon RAW). Tento požadavek je obzvláště důležitý vzhledem k tomu, že jako archivační formát (výstup digitalizace určený k dlouhodobému uchovávání) by měl vždy vybrán formát s bezztrátovou kompresí. Užití ztrátové komprese pro snímkový formát znamená, že archivační formát bude pouze formálně bezztrátový, ale reálně bude obsahovat ztrátová obrazová data.

4.2.5 Zabudování EXIF metadat do souborů

Pro digitalizaci je vhodné využívat EXIF metadata zabudovaná do obrazových dat. EXIF je široce rozšířeným a obrazovým průmyslem podporovaným standardem pro zabudovaná metadata pro rastrová data. Tato metadata jsou obsažena přímo v souboru a obsahují informace týkající se procesu snímání. EXIF není podporován ve všech obrazových formátech, ale jen pro JPEG, TIFF a RAW. Formáty TIFF a RAW jsou doporučenými formáty pro původní snímky. Je vhodné, aby tyto původní snímky tato metadata obsahovaly. Zpravidla je tato funkce již nastavena na snímacích zařízeních (skenery, fotoaparáty), případně je třeba zkontrolovat, zda funguje, nebo tuto funkci na snímacím zařízení nastavit. Tyto údaje nejsou sice zachovány v doporučeném formátu pro archivní kopie (JP2), který standard EXIF nepodporuje, ale jsou doporučeným zdrojem pro plnění technických metadat o formátu TIFF, resp. RAW.

4.2.6 Barevný profil

Na základě směrnice FADGI doporučujeme využít pro kopie digitalizovaných dokumentů standardizované barevné profily, a to pro obrazy bez barevně významné informace

(např. noviny) profil sRGB a pro archivní obrazy s barevně významnou informací profil Adobe RGB 1998. Snímek musí vždy obsahovat informaci o barevném prostoru. Buď tuto informaci obsahuje snímek z procesu skenování, nebo je v následných procesech nutné snímku barevný prostor přiřadit či původní snímek převést do vhodného barevného profilu (FADGI, 2016, s. 61).

Barevný profil neboli ICC profil je soubor dat, která charakterizují vstupní zařízení pořizující obrazový snímek nebo výstupní barevný prostor podle standardů International Color Consortium. Skenovací zařízení používaná pro vytváření digitálních kopií by měla mít od výrobce daný barevný profil zařízení, který je možné použít jako vstupní barevný profil pro přiřazení barevného prostoru, ve kterém skenujeme. Druhou a lepší možností je nechat si barevný profil zařízení vyrobit za pomoci standardizovaného kalibračního terče a kalibračního softwaru. Pomocí získaného barevného profilu potom provedeme přepočítání z barevného prostoru skeneru na některý ze standardizovaných barevných prostorů. Minimem by měl být barevný prostor definovaný ICC profilem sRGB, který podporují běžná zobrazovací zařízení. Pro archivní snímky, kde chceme zachovat přesnější barevnou informaci, může být vhodné využít rozšířený barevný prostor Adobe RGB 1998.

Pořizovatel digitální kopie má 3 možnosti, jak uložit snímek s barevnou informací:

1. Pořídí digitální obraz s vestavěným ICC profilem zařízení, který uloží do snímku nebo do samostatného souboru a ten uloží do archivu. Kdykoli pak může převést snímek do libovolného standardizovaného barevného prostoru a získat tak například uživatelskou kopii snímku v barevném prostoru sRGB vhodnou pro zobrazovací zařízení uživatelů.
2. Pořídí digitální obraz s vestavěným profilem zařízení a ten převede do standardizovaného ICC profilu (např. Adobe RGB 1998). V archivu tak uchovává obraz s připojeným standardizovaným barevným profilem s jednoznačně definovanými barvami. Pro uživatelské kopie může archivní snímek kdykoli převést na běžně používaný sRGB barevný prostor, pokud je archivní snímek uložen v prostoru s větším rozsahem barev než je sRGB.
3. Pořídí digitální obraz s vestavěným profilem zařízení, který převede přes standardizovaný ICC profil do digitálního obrazu ve standardizovaném barevném prostoru. U obrazu není připojen žádný ICC profil, protože obraz sám je uložen ve standardizovaném barevném prostoru. Pro uživatelské kopie je opět možné

obraz ve standardizovaném prostoru převádět na jiný, obvykle sRGB barevný prostor, pokud je archivní snímek uložen v prostoru s větším rozsahem barev než je sRGB.

4.3 Zpracování dat

Data, která tvoří konečný digitalizát, lze rozdělit do tří skupin – obrazová komponenta (archivní a uživatelské kopie), OCR komponenta (textová data vzniklá optickým rozpoznáváním znaků, která souřadnicově sedí na obrazová data) a textová komponenta (prostý text pro vyhledávání).

4.3.1 Zpracování obrazové komponenty

Zpracování obrazové komponenty začíná zpracováním původních snímků a končí vytvořením archivních a uživatelských kopií. Každý krok zpracování, který končí uložením souboru, vytváří novou generaci obrazové komponenty.

Ořezy by se měly provádět na původním snímku ve formátu TIFF a měly by být provedeny cca 1 mm vně okraje stránky.⁸³ Vyrovnávání zešikmení by se mělo provádět rovněž na původním snímku, a to podle prostředního řádku textu stránky.⁸⁴ Pokud je text v předloze rovnoběžný s okrajem stránky, je vhodné usilovat o pořízení snímků s vodorovně nasnímaným textem tak, aby nemuselo docházet k jejich dodatečnému vyrovnávání.

V případě fotografování doporučujeme provádět uvedené úpravy na původním snímku RAW v aplikaci výrobce fotoaparátu a po jejich skončení provést konverzi do formátu TIFF (verze 6, bez komprese).

Po skončení ořezů by měly být soubory ve formátu TIFF převedeny do archivních a uživatelských kopií, přičemž soubor archivní i uživatelské kopie by měl být vytvořen z téhož souboru TIFF předchozí generace. Soubor archivní i uživatelské kopie musí mít stejnou pixelovou velikost i stejné rozlišení. Jedním z důvodů je namapování OCR komponenty pro současná i budoucí (vzniklá budoucí formátovou konverzí v repozitáři) obrazová data.

Jako formát archivních i uživatelských kopií doporučujeme využít formát JP2 (obrazový

⁸³ Ořezy vně uchovají hrany stránek, které mohou nést potenciálně zajímavé informace pro koncové uživatele, badatele. Je však zpravidla nutné takový ořez provádět manuálně nebo poloautomaticky. U novodobé literatury, kde se informace o hraně považuje za nevýznamovou, je tedy akceptovatelný i ořez dovnitř. Takový ořez by pak měl být minimální, tj. neořezávat víc než je nutné (Vychodil, 2012).

⁸⁴ Případně dle převažující části textu, není-li stránka tištěna rovně; cílem je aby strana působila rovně.

rastrový formát specifikovaný v první části standardu ISO/IEC 15444-1:2004) v souladu s aktuálními požadavky Standardu NDK pro obrazová data (blíže viz další oddíly).

Všechny generace obrazové komponenty by měly být uchovány nejméně do doby úspěšného dodání balíčku SIP do repozitáře. Důvodem zachování je možnost případných oprav (chyb, které vznikly až ve zpracování), kontrola kvality (dohledatelnost postupu) a vytváření a validace technických metadat.

V řetězci datových transformací obrazové komponenty od původních snímků do archivních kopií nesmí být nikdy užitá ztrátová komprese. Změny ve zpracování od původních snímků do archivních kopií budou nutně zahrnovat několik datových transformací (formátové konverze, ořezy apod.). Žádná z nich nesmí zahrnovat ztrátovou komprimaci obrazových dat.

4.3.2 Archivační formát (formát pro archivní kopie)

Archivační formát (formát archivních kopií) je v projektu Národní digitální knihovna bezztrátový JP2. Následující text popisuje nastavení parametrů pro formátový profil JP2 pro archivní kopie.

4.3.2.1 Druh komprese a transformace

Druh komprese (<i>Compression</i>)	Bezeztrátová
Transformace (<i>Transformation</i>)	5-3 reversible filter

Obrazy budou komprimovány bezztrátově, aplikací filtrů 5-3.⁸⁵

4.3.2.2 Kompresní poměr

Výsledný kompresní poměr (<i>Compression ratio</i>)	Záleží dle charakteristik obrazu.
-------------------------------------------------------	-----------------------------------

Kompresní poměr se u migrace do bezztrátového JP2 explicitně neudává, nástroj je na základě vlastností obrazů spočítá sám.⁸⁶ Barevné obrazy mohou mít při bezztrátové kompresi

⁸⁵ Aplikace filtrů 9-7 na řádky a sloupce obrazu vede k výsledné ztrátové kompresi.

⁸⁶ U zpřístupňujících kopií ve formátu JP2 je kompresní poměr explicitně zadán, jedná se však o hodnotu aplikovanou na vstupní data nikoliv o výslednou hodnotu komprese obrazu. Ta se může, a je to pro projekt NDK přijatelné, od předepsané aplikované hodnoty mírně lišit, vliv na konečný kompresní poměr mají totiž vlastnosti konkrétního obrazu. Aby výsledná hodnota kompresního poměru odpovídala předepsané hodnotě, musel by se pro každý jednotlivý obraz spočítat bitrate, což je

kompresní poměr okolo 1:2, strana s textem okolo 1:5 a prázdná, bílá strana i 1:400.

4.3.2.3 Dlaždice

Dlaždice (<i>Tiling</i>)	4096x4096
----------------------------	-----------

Vstupní obraz může být v začátku transformačního procesu zpracováván buď jako celek (tj. jeden obraz odpovídá jedné dlaždici), nebo může být rozdělen na dlaždice čtvercového tvaru. Každá dlaždice je pak zpracovávána separátně. Rozdělení obrazu na dlaždice urychluje proces komprese a dekomprese obrazu. Minimální povolená velikost jedné dlaždice je 128x128 obrazových bodů, tato velikost vede ke generaci velkého počtu dlaždic a naopak ke zpomalení procesů. Pro archivační obrazy není důležitá doba dekomprese obrazu a není tedy nutné obraz dělit na menší dlaždice, jako je tomu například u zpřístupňujících kopií JP2 v NDK.⁸⁷ Zároveň také dnes používané systémy velmi rychle dlaždice o velikosti 4096x4096 dekomprimují, není tedy nutné datový tok archivních kopií dále zesložit'ovat rozdělením obrazu na menší dlaždice.

4.3.2.4 Průběh zobrazení

Průběh zobrazení (<i>Progression order</i>)	RPCL
-----------------------------------------------	------

Parametr průběh zobrazení (*Progression order*) udává, jak budou posílány pakety při přenosu dat a při jejich dekompresi. Standard definuje 5 způsobů zobrazení: LRCP, RLCP, RPCL, PCRL a CPRL, kde jednotlivá písmena odpovídají přenášeným datům v paketu (L=vrstvě kvality, R= rozlišení, C= barevná komponenta a P= pozici). V případě hodnoty zvolené pro projekt Národní digitální knihovny, tj. RPCL, dochází k seskupování paketů dle rozlišení, tj. nejdříve jsou shromážděna všechna data odpovídající první vrstvě rozlišení, která je tak zobrazena ve své maximální možné kvalitě, následně se nahrává další vrstva rozlišení a jí odpovídající barevnost, kvalita apod. V případě tohoto pořadí zobrazení uživatel vidí postupně se zvětšující obraz.

4.3.2.5 Dekompoziční úroveň

proveditelné, nikoliv však pro projekt NDK nutné.

⁸⁷ Doba dekomprese dlaždic o velikosti 1024x1024 je rychlejší než u velikosti dlaždic 4096x4096.

Počet dekompozičních úrovní (<i>Decomposition level</i>)	5 nebo 6
---------------------------------------------------------------	----------

Počet dekompozičních úrovní se odvíjí od velikosti vstupního obrazu a od požadavků organizace na zpřístupnění (velikost náhledu, počet „obrazů“ mezi náhledem a plnou velikostí obrazu). Přítomnost několika dekompozičních úrovní má příznivý vliv na výsledný obraz a obecně se doporučuje do obrazu jich několik dát, standardem se zdá být minimální počet 5 dekompozičních vrstev.

4.3.2.6 Vrstvy kvality

Počet vrstev kvality (<i>Quality layers</i>)	1
------------------------------------------------	---

Počet vrstev kvality zjednodušeně určuje, kolik obrazů s různou úrovní kvality (s různým kompresním poměrem) je možné z jednoho datového proudu extrahovat. Protože archivní kopie JP2 nejsou primárně určeny pro zpřístupnění, je dostatečná jedna úroveň kvality.

4.3.2.7 Regiony

Velikost regionů (<i>Precinct size</i>)	256x256 pro první dvě dekompoziční úrovně, 128x128 pro nejnižší úrovně
-------------------------------------------	------------------------------------------------------------------------

Okrsky (příp. regiony, angl. *precincts*) sdružují bloky kódů souvisejících dat, jež se ukládají do jednotlivých balíčků a slouží k tomu, aby místně související data byla umístěna v jednom balíčku, přičemž jeden okrsek může být rozdělen do více paketů. Okrsky, stejně tak jako dlaždice, umožňují přístup k vybraným částem obrazu, tj. například načtení jednotlivých vybraných okrsků. Použitím okrsků lze tedy pro zpřístupnění rozdělit dlaždice na menší části.

4.3.2.8 Zájmové oblasti

Zájmové oblasti (<i>Regions of Interests</i>)	Ne
-------------------------------------------------	----

Zájmové oblasti jsou části obrazu, které jsou dekomprimovány prioritně vůči ostatním částem obrazu, jež jsou pro uživatele méně významné (např. pozadí). Při dekompresi se tedy tato část zobrazuje obvykle nejdříve a ve vyšší kvalitě (Vrtělová, 2017). V projektu NDK se funkce zájmových oblastí nevyužívá.

4.3.2.9 Velikost bloků

Velikost bloků (<i>Code block size</i>)	64 x 64
-------------------------------------------	---------

Bloky kódu (*codeblock*) jsou čtvercového tvaru a jsou na sobě nezávisle kódovány do výsledného datového toku. Velikost bloků může nabývat hodnot od 4x4 pixelů do 1024x1024 pixelů. Nejobvyklejší velikost těchto bloků v paměťových institucích je 64x64 obrazových bodů.

4.3.2.10 Lokalizace dlaždice

Značka lokalizující dlaždice TLM ⁸⁸ (<i>Tile Length Markers</i>)	Ano
----------------------------------------------------------------------------------	-----

Značkovací segment TLM nese informaci o délce dlaždic, resp. částí dlaždic v celém toku dat jednoho souboru. Tato informace může následně posloužit k rychlejší lokalizaci a orientaci v proudu dat při použití více dlaždic.

4.3.2.11 Přemostění

Přemostění (<i>Bypass</i>)	Ano
------------------------------	-----

Parametr přemostění, tj. BYPASS, se týká procesu komprese a dekomprese, jedná se o režim, v jakém bude obraz zpracován. Hodnota BYPASS znamená, že kodér při kompresi vynechá kompresi některých, méně významných dat, čímž se urychlí proces komprese i následné dekomprese (i o 20%).⁸⁸ Výsledná komprese obrazu je pak o něco menší.

4.3.2.12 ICC profily

ICC profily (<i>ICC Profile</i>)	Ano
------------------------------------	-----

Obrazy by měly vždy obsahovat informaci o svém barevném profilu, což zajistí, že barvy budou v následných zobrazovacích aplikacích správně interpretovány. Formát JP2 (tj. Part 1

⁸⁸ U nástroje OpenJpeg je až do verze 2.2.0 tento režim nefunkční.

standardu JPEG 2000) podporuje barevný prostor sRGB a vybrané ICC profily. Pro účely digitální archivace se kvůli co nejvěrnějšímu zachování barev doporučují profily Adobe RGB a ProPhotoRGB, které mají větší rozsah barev (gamut).

4.3.2.13 Hlavička segmentu paketů

Značka začátku hlavičky segmentu paketů SOP (<i>Start of Packet Header</i>)	Ano (Cuse_sop=yes)
Značka konce hlavičky segmentu paketů EPH (<i>End of Packet Header</i>)	Ano (Cuse_eph=yes)

Značky SOP a EPH označují začátek a konec paketů tvořících tok jednoho souboru. Zvyšují odolnost souboru proti přenosovým chybám, přijímající systém (protokol) díky nim dokáže rozpoznat, že mu nějaká data chybí.

4.3.2.14 Vložená metadata

Vložená metadata (<i>Embedded Metadata</i>)	Ne
-----------------------------------------------	----

Do obrazů ve formátu JPEG 2000 je možné vložit související metadata,⁸⁹ například identifikátor, informace o souvisejících právech apod. Takové soubory pak mohou obsahovat dostatečné informace a nehrozí jejich ztráta v systému. Aktuálně se tato funkcionality formátu v projektu NDK nevyužívá, obrazy jsou ale pojmenovány s využitím jednoznačného identifikátoru související intelektuální entity.

4.3.3 Prezentací formát (formát pro uživatelské kopie)

Pro zpřístupňování je doporučeno užít formát JP2 ve ztrátové kompresi, jež umožní rychlejší přenos a zobrazení dokumentu uživateli.⁹⁰ Tento JP2 je prezentačním meziformátem (zvaný někdy jako *production master copy*), tj. v prezentačním systému z něj aplikace (tzv. *image server*) generuje obrázky ve formátu JPEG.⁹¹

⁸⁹ Nejedná se o metadata EXIF, formát JP2 nepodporuje EXIF standard.

⁹⁰ Tento kompresní poměr je ideální zvolit dle skenovaných dokumentů a ověřit experimentálně. V projektu NDK bylo zjištěno, že i při kompresním poměru 1:20 nese obraz významné informace, z důvodu výskytu kompresních artefaktů u kompresních poměrů 1:20 až 1:30 byl však zvolen jako vhodnější kompresní poměr 1:8 až 1:10.

⁹¹ Tímto odpadá na straně uživatelů nutnost obrazy stahovat nebo instalovat plugin pro zobrazení obrazů JP2. Prohlížeče nativně formát JPEG podporují.

PARAMETRY PRO FORMÁTOVÝ PROFIL PRO UŽIVATELSKÉ KOPIE	
Druh komprese (<i>Compression</i>)	Ztrátová
Transformace (<i>Transformation</i>)	9-7 reversible filter
Výsledný kompresní poměr (<i>Compression ratio</i>)	1:8 až 1:30
Dlaždice (<i>Tiling</i>)	1024x1024
Průběh zobrazení (<i>Progression order</i>)	RPCL
Počet dekompozičních úrovní (<i>Decomposition level</i>)	5 nebo 6
Počet vrstev kvality (<i>Quality layers</i>)	12 (logaritmicky)
Velikost regionů (<i>Precinct size</i>)	256x256 pro první dekompoziční úroveň, 128x128 pro nejnižší dekompoziční úroveň
Zájmové oblasti (<i>Regions of Interests</i>)	Ne
Velikost bloků (<i>Code block size</i>)	64x64
Značka lokalizující dlaždice TLM (<i>Tile Length Markers</i>)	Ano („R“)
Přemostění (<i>Bypass</i>)	Ano
ICC profily (<i>ICC Profile</i>)	Ano
Značka začátku hlavičky segmentu paketů SOP (<i>Start of Packet Header</i>)	Volitelné
Značka lokace hlavičky segmentu paketů EPH (<i>End of Packet Header</i>)	Volitelné
Vložená metadata (<i>Embedded Metadata</i>)	Ne

4.3.4 Vytváření OCR komponenty

OCR komponenta (textová data získaná nástrojem pro optické rozpoznávání znaků) musí být zapsána ve formátu ALTO XML.⁹² OCR komponenta (výstup OCR) by měla být vytvářena až z archivních kopií ve formátu JP2. Důvodem je zachování mapování textu na pixelovou pozici v obraze. V metadatové části souboru s OCR by mělo být uvedeno, z jakého souboru byl daný soubor s OCR vytvořen (element `<sourceImageInformation>`).

4.4 Vytváření metadat

4.4.1 Převod bibliografických metadat

Bibliografická metadata (MODS, DC) by měla být automaticky přebírána z katalogizačního záznamu z katalogu knihovny, která je vlastníkem předlohy. Údaje na vyšší míře granularity, než kterou popisuje katalogizační záznam (ročník a číslo periodika, nebo svazek vícesvazkové monografie), musejí být zpravidla vytvářeny manuálně pracovníkem v digitalizaci. Při popisu periodik by měl mít pracovník k dispozici kompletní celý ročník fyzické předlohy.

Převod katalogizačního záznamu do formátu MODS zajišťuje série transformačních šablon. Tyto šablony využívají jazyk XSLT a pomocí série příkazů a podmínek jsou schopny z jazyka katalogizačních záznamů MARC 21 vytvořit záznam v MODS. Digitalizace NDK využívá upravené transformační šablony Kongresové knihovny Spojených států amerických.⁹³ V praxi jiných digitalizací je ale možné využívat vlastní verzi této šablony, nebo šablony, vytvořené pro vlastní specifické potřeby instituce. Jedinou podmínkou je, že vytvořená data musí odpovídat Standardu NDK.

Nejprve se katalogizační záznam převede z MARC 21 do MARCXML, kdy již dostane podobu XML jazyka, ovšem se zachovanou informací o polích a indikátorech, které slouží k další transformaci z MARCXML do MODS a Dublin Core.

Dále by hlavní METS dokument měl obsahovat fyzickou a logickou mapu záznamu. Fyzická mapa slouží k identifikaci každé jednotlivé naskenované strany v souborech master copy, user copy, alto, ocr, txt a vedlejšího METS souboru amd_mets. Logická mapa naproti

⁹² <https://www.loc.gov/standards/alto/>

⁹³ Ukázky šablon využívaných v projektu NDK a šablon využívaných v rámci digitalizačního nástroje ProArc jsou k dispozici zde: <https://standardy.ndk.cz/ndk/standardy-digitalizace/metadata>.

tomu mapuje hierarchické pořadí jednotlivých popisovaných úrovní dokumentu pomocí zanořených <div> a slouží jako podpora správného zobrazování digitalizátu v aplikacích pro čtení.

4.4.2 Získávání technických metadat

Technická metadata by měla být v maximální možné míře převzata z metadatových extraktorů, které tyto údaje získávají přímo ze souborů. Postup získávání metadat by tedy primárně neměl být založen na údajích, které se přednastaví do digitalizačního systému (např. název skeneru nebo skenovací aplikace pro jednu linku) a systém je pak jen automaticky přiděluje do metadat všech dokumentů dané linky. Pokud jsou tyto informace zaznamenávány do metadat z přednastavených hodnot uložených v digitalizačním systému, hrozí vždy lidská chyba (např. při výměně skeneru nebo aktualizaci skenovací aplikace se neprovede změna přednastavených hodnot). Dále platí, že pro získávání informací o formátu nelze užít prostý opis koncovky souboru.

V případě technických metadat, která se plní o původních snímcích (ve formátu TIFF nebo RAW) je vhodné využívat EXIF metadata, která jsou obsažena přímo v souboru. EXIF metadata jsou relevantní pouze pro popis formátu původního snímku (tedy TIFF nebo RAW) a některých dalších vlastností původních snímků, nikoliv pro další generace dat. Rovněž se doporučuje užít EXIF metadata pro informace o snímacím zařízení a způsobu snímání. Iniciativa FADGI vydala doporučení, jaké minimální elementy EXIF metadat ve formátu TIFF je vhodné zaznamenat do metadat (Embedded metadata working group-Smithsonian institution, 2010).

Pro převod ze schématu metadatového extraktoru (např. JHOVE) do technických metadat v MIX a PREMIS je potřeba předem ověřit možnosti namapování, a poté toto namapování nastavit v digitalizačním systému. Pro různé metadatové extraktory, včetně různých verzí téhož nástroje, se může namapování lišit. V některých případech nemusí být zcela jednoznačné.

Doporučené minimální metadatové extraktory jsou JHOVE a jpylyzer. Příloha této metodiky obsahuje doporučené namapování pro aktuální verze nástroje jpylyzer a JHOVE.

4.4.3 Nástroje pro formátovou identifikaci

Jako základní formátové identifikační nástroje musí být užity nástroje DROID⁹⁴ a JHOVE.

Nástroj DROID čerpá informace z formátového registru PRONOM. Tento registr obsahuje identifikátor PUID, který je klíčový pro jednoznačnou identifikaci formátu (Brown, 2006, s. 4). Název formátu není pro identifikaci formátu dostatečný. PUID je „rozšiřitelné schéma pro poskytování perzistentních, jedinečných a jednoznačných identifikátorů pro jednotky interpretačních informací zaznamenané v registru PRONOM“ (Brown, 2006, s. 4). Formát je tedy pouze jedním z typů interpretačních informací, o nichž registr vede údaje, nicméně nejrozšířenějším. Funkce identifikátoru PUID jsou dvě: propojení se záznamem jednotky interpretačních informací v registru PRONOM (tj. způsob identifikace záznamu, přičemž tento záznam by ideálně měl obsahovat co nepodrobnější informace o formátu nebo jiné jednotce interpretačních informací) a jedinečný perzistentní identifikátor, který odlišuje v maximální možné míře jeden formát od druhého (odlišuje se nejen typ formátu, ale i verze). Například PUID pro JPEG verze 1.00 je „fmt/42“, pro verzi 1.01 „fmt/43“ a pro verzi 1.02 „fmt/44“. Registr MIME,⁹⁵ který je nejužívanějším obecným registrem formátů (sloužícím i pro účely mimo kontext digitální archivace), odlišuje formáty jen na základě typu a názvu, například formát JPEG všech verzí má označení „image/jpeg“.

Nástroj DROID využívá ke své činnosti metadatové soubory obsahující záznamy z registru PRONOM, které se nazývají „signature files“. Z důvodu průběžné aktualizace registru PRONOM se pro účely identifikace doporučuje použít vždy nejnovější verzi nástroje DROID a nejnovější verzi „signature files“.

Standard NDK doporučuje využít formátový identifikační nástroj, který pracuje s registrem PRONOM, a tedy dokáže souboru přidělit identifikátor PUID. Tomu pak musí odpovídat záznam v metadatech PREMIS. Aktuální doporučovaný nástroj je uveden v příloze. Níže jsou uvedeny doporučené způsoby získávání a záznamu informací o formátové identifikaci.

Elementy PREMIS <formatName> a <formatVersion> by měly být vyplněny užitím nástroje JHOVE.

⁹⁴ Aktuálně je možné využít i nástroj FIDO (<http://openpreservation.org/technology/products/fido/>), jenž je též založen na registru PRONOM. Preference nástroje DROID v této metodice se odvíjí od jeho delší existence a faktu, že je vyvíjen stejnou institucí jako registr PRONOM. Mezi další důvěryhodné implementace registru PRONOM lze v současné době řadit například nástroj Siegfried (<https://itforarchivists.com/siegfried>). Využití kombinace více nástrojů založených na registru PRONOM lze obecně považovat za dobrou praxi, neboť kombinace více přístupů při identifikaci formátu posiluje její spolehlivost.

⁹⁵ <http://www.iana.org/assignments/media-types/media-types.xhtml>

Elementy PREMIS v metadatovém kontejneru <formatRegistry> by měly být vyplněny následujícím způsobem:

- element <formatRegistryName> by měl vždy obsahovat hodnotu „PRONOM“;
- element <formatRegistryKey> by měl obsahovat identifikátor PUID, který by měl být získán nástrojem DROID

4.4.4 Propojování událostí s objektem a agentem

Při propojování událostí s objektem je třeba v souladu se standardem PREMIS důsledně odlišovat vztah odvození od ostatních vztahů. Standard PREMIS předepisuje následující logiku metadatového zápisu pro změny objektu: „metadata, soubory, bitové toky a reprezentace se popisují jako statické množiny bitů. Není možné změnit soubor (nebo bitový tok nebo reprezentaci); lze pouze vytvořit nový soubor (nebo bitový tok nebo reprezentaci), který se vztahuje k zdrojovému objektu“ (Premis, 2015, s. 22). Tento vztah mezi novým a předchozím objektem definuje jako vztah odvození (*derivation relationship*). Standard odlišuje dva typy odvození ze zdrojového digitálního objektu do nového objektu: replikace (*replication*) a transformace (*transformation*) (Premis, 2015, s. 19). Replikace znamená vytvoření digitální kopie, která je bitově identická se zdrojovým digitálním objektem (Premis, 2015, s. 272), transformace má za výsledek vytvoření jednoho nebo více digitálních objektů, které nejsou bitově identické se zdrojovým objektem (Premis, 2015, s. 273).

Na událost, která je odvozením, je třeba ze záznamu pro dotčené objekty (zdrojový a výsledný) odkázat pomocí kontejnerového elementu <relatedEventIdentification> a v něm vnořených subelementů. Na jiné události (tj. na ty, které odvozením nejsou) je třeba ze záznamu pro příslušný objekt odkázat pomocí kontejnerového elementu <linkingEventIdentifier> a jeho subelementů. V obou případech (tj. bez ohledu na to, zda jde o odvození, či nikoli) se propojují i naopak událost s objektem. Za tímto účelem se v záznamu pro událost využije kontejnerový element <linkingObjectIdentifier>.

Jako agenti musejí být vždy uvedeny všechny nezávisle existující nástroje, které jsou bezprostředním původcem události. V praxi bývá někdy při digitalizaci užit komplexní digitalizační systém, který řídí operace a integruje různé jiné nástroje pro dílčí operace. Pokud je pro událost užit nástroj, který existuje nezávisle na tomto systému, je dobrou praxí tento nezávislý nástroj uvést jako samostatného agenta, tj. činitele události. Například nástroje pro formátovou identifikaci nebo validaci, které reálně vykonávají událost, digitalizační systém je

pouze využívá.

4.4.5 Záznam událostí a nástrojů

Kongresová knihovna udržuje řízený slovník pro PREMIS, který obsahuje zejména typy událostí.⁹⁶ Doporučení níže uvádějí jak události povinné z hlediska Standardu NDK, tak doporučené nad rámec Standardu NDK. Objektem událostí, uvedených v tomto oddíle, je vždy soubor (úroveň „file“ v modelu PREMIS): k jednomu souboru se váže jedna nebo více událostí.

<eventType>	<eventDetail>	popis
capture	digitization	Vytvoření původního skenu
capture	XML_creation	Vytvoření souboru ALTO
capture	TXT_creation	Vytvoření souboru textu
migration	MC_creation	Vytvoření archivní kopie v JP2 z původního skenu
deletion	PS_deletion	Smazání původního skenu

4.4.5.1 Snímání (skenování, fotografování)

Snímáním se rozumí vlastní proces skenování nebo fotografování, jehož bezprostředním výstupem jsou původní snímky. Tato událost se vztahuje k objektům, jimiž jsou soubory původních snímků (TIFF nebo RAW).

Podle Standardu NDK se zapíše typ události (PREMIS: <eventType>) s hodnotou „capture“ a detail události (PREMIS: <eventDetail>) s hodnotou „digitization“; snímací zařízení (MIX: <captureDevice>) se v případě skeneru zapíše s hodnotou buď „reflection print scanner“ (nejčastěji používaný typ skeneru), nebo „transmission scanner“, v případě fotografování s hodnotou „digital still camera“. Propojení objektu s událostí se provede elementy metadatového kontejneru PREMIS <relatedEventIdentification>.

4.4.5.2 Formátová konverze z TIFF do JP2 archivní kopie

Tato událost se váže k archivním kopiím (tj. souborům ve formátu JP2 v bezztrátové kompresi, které byly vytvořeny ze souborů předchozí generace ve formátu TIFF).

⁹⁶ <http://id.loc.gov/vocabulary/preservation/eventType.html>

Podle Standardu NDK se zapíše typ události (PREMIS: <eventType>) s hodnotou „migration“ a detail události (PREMIS: <eventDetail>) s hodnotou „MC_creation“; kodek, který byl užít k vytvoření souborů ve formátu JP2 (PREMIS: <creatingApplicationName>) se zapíše celým názvem (doporučené hodnoty pro nejčastěji užívané aplikace jsou „Kakadu“ a „OpenJPEG“) a zvlášť se zapíše i verze kodeku (PREMIS: <creatingApplicationVersion>). Údaje o kodeku se zapíší i do části PREMIS Agent; název agenta (PREMIS: <agentName>) se zapíše sloučením elementů PREMIS <creatingApplicationName> a <creatingApplicationVersion> (mezi nimi musí být mezera); typ agenta (PREMIS: <agentType>) se zapíše s hodnotou „software“; poznámka o agentovi (PREMIS: <agentNote>) by měla začínat hodnotou „command line: “ (za dvojtečkou je mezera) a následovat bude konkrétní příkazový řádek užitý v daném kodeku. Propojení objektu s událostí se provede elementy metadatového kontejneru PREMIS <relatedEventIdentification>.

4.4.5.3 Vytvoření ALTO XML z OCR

Tato událost se váže k souborům OCR komponenty, které vytvořil software pro optické rozpoznávání znaků z obrazových dat a které jsou ve formátu ALTO XML. Podle Standardu NDK se zapíše typ události (PREMIS: <eventType>) s hodnotou „capture“ a detail události (PREMIS: <eventDetail>) s hodnotou „XML_creation“; software, který byl použit k vytvoření souborů OCR komponenty se zapíše celým názvem (ALTO:<OCRProcessing>:<ocrProcessingStep>:<processingSoftware>:<softwareName>). Také se zapíše verze softwaru: (ALTO:<OCRProcessing>:<ocrProcessingStep>:<processingSoftware>:<softwareVersion>).

Nad rámec Standardu NDK doporučujeme zapsat název agenta (PREMIS: <agentName>) pomocí element ALTO <processingSoftware>:<softwareName>; typ agenta (PREMIS: <agentType>) zapsat s hodnotou „software“. Propojení objektu s událostí se provede elementy metadatového kontejneru PREMIS <relatedEventIdentification>.

4.4.5.4 Formátová identifikace

Formátová identifikace je jeden z klíčových procesů digitální archivace. Doporučujeme proto nad rámec povinností stanovených Standardem NDK zapsat do metadat událost formátové identifikace, a to jako opakovanou událost, při níž byly užity nejméně dva nástroje (JHOVE a DROID, viz 4.4.3). Tato událost se váže k souborům původních snímků (TIFF nebo RAW) a souborům archivních kopií (bezeztrátový JP2).

Jako typ události (PREMIS: <eventType>) se zapíše hodnota „format identification“. Pro užití nástroje JHOVE se zapíše název agenta (PREMIS: <agentName>) s hodnotou obsahující „JHOVE“ + číslo verze, například „JHOVE v1.20“. Pro užití nástroje DROID se zapíše název agenta (PREMIS: <agentName>) s hodnotou obsahující „DROID: version“ + číslo verze + verze souboru signature files, například „DROID: version: 6.4, Signature files: 1. Type: Container Version: 20171130 File name: container-signature-20171130.xml 2. Type: Binary Version: 93 File name: DROID_SignatureFile_V93.xml“.

Propojení objektu s událostí se provede elementem PREMIS <linkingEventIdentifier>.

4.4.5.5 Formátová validace

Formátová validace je další z klíčových procesů digitální archivace. Doporučujeme proto nad rámec povinností stanovených Standardem NDK zapsat do metadat událost formátové validace. Tato událost se váže k souborům původních snímků (TIFF nebo RAW), kdy doporučujeme užít nástroj JHOVE jako formátový validátor, a souborům archivních kopií (bezeztrátový JP2), kdy doporučujeme užít nástroje JHOVE a jpylyzer jako formátové validátory (tato událost bude tedy opakovatelná).

Jako typ události (PREMIS: <eventType>) se zapíše hodnota „validation“ a jako detail události (PREMIS: <eventDetail>) hodnota „format validation“.

Pro užití nástroje JHOVE se zapíše název agenta (PREMIS: <agentName>) s hodnotou obsahující „JHOVE“ + číslo verze, například „JHOVE v1.20“. Pro užití nástroje jpylyzer se zapíše název agenta (PREMIS: <agentName>) s hodnotou obsahující „jpylyzer“ + číslo verze, například „jpylyzer v1.18.0“. Dále doporučujeme vyplnit element (PREMIS:<eventOutcomeInformation>) hodnotou obsahující textové výstupy validátorů JHOVE a jpylyzer popisující výsledek formátové validace. V případě aktuální verze nástroje JHOVE doporučujeme převzít hodnotu obsaženou v elementu schématu JhoveView <status> (např. „Well-Formed and valid“). V případě aktuální verze nástroje jpylyzer, pokud jeho metadatové schéma obsahuje v elementu <isValidJP2> hodnotu „True“, pak se zapíše hodnota „valid“, pokud jpylyzer vypíše hodnotu „False“ pak by se v metadatech měla objevit hodnota „not valid“. Propojení objektu s událostí se provede elementem PREMIS <linkingEventIdentifier>.

5 Kontrola kvality

Z hlediska kontroly kvality je pro oblast dlouhodobého uchovávání klíčová validace následujících čtyř oblastí:

- Validace metadat
- Formátová validace
- Datová validace
- Validace balíčku SIP

5.1 Digitální otisk

Podle Standardu NDK je nutný digitální otisk (konkrétně MD5) pro soubory v konečném balíčku SIP. Digitální otisky pro jednotlivé soubory (vyjma souboru info.xml a souboru obsahujícího MD5) jsou uloženy v metadatech, pro celý balíček SIP (vyjma souboru info.xml a souboru obsahujícího MD5) pak rovněž v podobě samostatnému souboru MD5 v kořenovém adresáři balíčku SIP. Tento soubor MD5 má svůj digitální otisk v souboru info.xml.

Postup, který musí být dodržen pro vytvoření digitálních otisků souborů balíčku SIP v souladu se Standardem NDK, je následující: musejí již existovat soubory v následujících podadresářích: adresář se soubory OCR komponenty v ALTO XML („alto“), adresář se soubory archivních kopií („masterCopy“), adresář se soubory uživatelských kopií („userCopy“) a adresář se soubory obsahujícími prostý text („txt“). Digitální otisk k souborům v těchto podadresářích musí být zapsán do vedlejších metadatových souborů (podadresář „amdsec“). Hodnoty těchto digitálních otisků musejí být následně zapsány také do hlavního souboru v METS XML v kořenovém adresáři balíčků a k nim musejí být přidány digitální otisky vedlejších metadatových souborů. Teprve poté může být vytvořen soubor MD5 v kořenovém adresáři balíčku SIP (tento soubor obsahuje i digitální otisk hlavního souboru v METS XML). Na závěr musí být vytvořen soubor info.xml v kořenovém adresáři, který obsahuje výčet souborů včetně samostatného souboru MD5. Soubor MD5 v kořenovém adresáři tedy logicky nemůže obsahovat digitální otisk k samostatnému souboru MD5, ani k info.xml. Digitální otisk samostatného souboru MD5 se nakonec umístí do elementu <checksum> v souboru info.xml.

Balíčky SIP jsou však vytvářeny v několika fázích, ve kterých vzniká několik odlišných generací obrazových dat (minimálně tři základní generace původní snímky v TIFF / RAW,

ořezané snímky v TIFF, archivní a uživatelské kopie v JP2). Tyto generace by měly být podle doporučení z kapitoly 4.3.1 této metodiky zachovány do doby úspěšného dodání balíčku SIP do repozitáře.

Pro bezproblémový proces vytváření souborů (všech typů) doporučujeme, aby bezprostředně po každém uložení nového souboru byl vytvořen digitální otisk (a to i k předchozím generacím souborů, které nejsou obsaženy v balíčku SIP). Tyto digitální otisky by měly být dočasně uchovávány v digitalizačním systému a při finalizaci balíčku SIP zapsány do metadat a samostatného souboru MD5 podle postupu uvedeného v předchozím odstavci. Důvod, proč je důležité, aby byl digitální otisk vytvořen bezprostředně po vytvoření souboru, je ten, že později vytvořený digitální otisk již může být otiskem porušených souborů (v době mezi vytvořením souboru a finalizací balíčku SIP může z různých důvodů dojít k porušení souboru, například vadou datového nosiče dočasného digitalizačního úložiště), a tedy by mohl pozbýt své kontrolní funkce. Po vytvoření digitálního otisku by také měla být okamžitě spuštěna kontrola neporušenosti (jako kontrola toho, že samotný digitální otisk byl vytvořen správně). Dále je vhodné spustit kontrolu neporušenosti vždy před tím, než vznikne nová generace souboru (jako ověření toho, že zdrojový soubor není poškozen). To platí zejména pro formátovou migraci z TIFF do JP2.

Do elementu PREMIS <messageDigestOriginator> doporučujeme zapsat vždy software, který MD5 vytvořil.

5.2 Validace metadat

Validace metadat je kontrola souladu metadat s předepsanou podobou metadatového profilu v balíčku SIP. Metadatový profil je soubor metadatových elementů, které jsou vybrány z jednoho nebo více metadatových standardů a jsou spojeny do jednoho sloučeného schématu, který je uzpůsoben na míru funkčním požadavkům konkrétního užití, zatímco je zachována interoperabilita s původními standardy (Duval, 2002). Definice metadatových formátů ve Standardu NDK jsou metadatové profily, vytvořené z mezinárodních metadatových standardů, užívaných knihovnami i jinými institucemi pro potřeby digitální archivace (základem jsou mezinárodní standardy METS a PREMIS, k nimž se podle typu dokumentu/dat poji mezinárodní standardy pro popisná a technická metadata). Podobně je tomu s metadatovými profily, které vydávají jiné národní knihovny i další instituce. Validace metadat v případě metadatového profilu zahrnuje jednak kontrolu užitím oficiálních validačních schémat mezinárodních standardů (zaznamenaných v XSD), jednak kontrolu podle profilu. Validaci

metadat lze rozdělit na dva typy. Základní validace metadat je kontrola přítomnosti elementů, strukturálních vztahů mezi nimi a dodržení obecného omezení pro hodnoty elementů (např. pouze znakový řetězec). Rozšířená validace je specifitější sémantická kontrola (např. konkrétní počet znaků určitého typu).

Pro účely validace metadat balíčků SIP vytvořených podle Standardu NDK lze využít nástroj Komplexní validátor, který vyvíjí Národní knihovna ČR. Komplexní validátor je lokálně instalovatelná aplikace určená ke kontrole SIP balíčků vytvořených podle aktuálně platných Definic metadatových formátů Standardu NDK. Tento nástroj však nepokrývá všechny historické verze Standardu NDK a v případě vydání nových Definic metadatových formátů bude zřejmě vždy existovat určitá prodleva, než bude Komplexní validátor aktualizován o možnost validace podle nejnovějších Definic metadatových formátů. Z tohoto důvodu lze doporučit využívat jen takové verze Definic metadatových formátů Standardu NDK, které dokáže validovat Komplexní validátor, a pokud to není možné, pak alespoň provádět validaci mezinárodních metadatových standardů, obsažených ve Standardu NDK, dle jejich oficiálních validačních schémat (XSD). Kontrola pomocí validátoru pak postupuje podle několik dílčích kroků. Jedná se o kontrolu struktury a integrity SIP balíčku, o validaci bibliografických, technických i administrativních metadat pomocí šablon a validaci obrazových dat. Výsledek kontroly validátor vypíše formou reportu obsahujícího informaci, zda balíček prošel validací úspěšně, v opačném případě obsahuje report počet chyb a varování i jejich detailní popis. Nad rámec kontroly pomocí validátoru by měla být rovněž provedena manuální validace metadat na úrovni popisu částí základní intelektuální entity (desky, přideštiny, titulní strana apod.) a dodržení původní posloupnosti stran, optimálně s tištěnou předlohou v ruce.

5.3 Formátová validace

Formátová validace je kontrola, zda je souborový formát vytvořen dle požadavků daných jeho oficiální dokumentací, případně dalších požadavků (např. nastavení formátového profilu, tedy nastavení parametrů v rámci formátu).

Formátová validace by měla být provedena bezprostředně po vytvoření nové generace obrazových dat. Doporučené nástroje pro formátovou validaci jsou: pro formáty TIFF a JP2 nástroj JHOVE vždy nejnovější verze a pro formát JP2 dále ještě nástroj jpylyzer. Validace formátu TIFF (tedy souborů předchozích generací obrazových dat), které nebudou dlouhodobě uchovávány, je důležitá pro to, aby bylo ověřeno, že pro následnou formátovou migraci do formátu JP2 byly užity zdrojové soubory, které byly vytvořeny korektně (tj. ve validním

formátu TIFF).

Pro kontrolu formátového profilu JP2 (lokální nastavení formátu v rámci možností oficiální specifikace, která umožňuje volbu různých parametrů) doporučujeme využívat nástroj Komplexní validátor. Případně je možné, aby instituce využila jiný nástroj, který dokáže využívat nástroj jpylyzer a výstupy srovnávat s oficiálním profilem JP2 doporučeným ve Standardu NDK.

5.4 Datová (obrazová) validace

Datová validace je validace datových prvků souboru, která je nad rámec formátové validace. Datová validace, tak jak ji chápeme v této metodice, ověřuje, zda je soubor neporušený a zda jej lze v odpovídající aplikaci otevřít. Tyto vlastnosti někdy formátové validátory neodhalí. Tuto validaci je možné provádět automatizovaně pomocí specializovaných nástrojů či manuálně (ověření otevřítelnosti). Pro obrazy ve formátu JP2 doporučujeme využít nástroj ImageMagick, případně Kakadu. Tyto nástroje pomáhají odhalit chyby v obrazovém datovém toku, které formátové validátory nemohou odhalit (např. otevřítelnost souboru). Optimálně pro datovou obrazovou validaci lze rovněž využít Komplexní validátor, který umožňuje oba uvedené nástroje zapojit do validačního procesu.

Rovněž by měla být provedena manuální kontrola kvality obrazů při dodatečném popisu (označení desek, přídeští apod.) digitalizačními pracovníky, případně již pracovníky skenování. Součástí této kontroly je vizuální inspekce obrazu (příp. poslech záznamu). Měly by být prohlédnuty náhledy všech obrazů, čímž se zkontroluje kompletnost a vizuální konzistence skenování. Z těchto obrazů by následně měl být vybrán vzorek a ten by měl projít důkladnější vizuální kontrolou.⁹⁷

Iniciativa FADGI doporučuje kontrolovat následující oblasti pro shodu s projektovou specifikací a pro detekci defektů:

- kontrola otevřítelnosti souboru
- kontrola vlastností souboru (komprese, barevný prostor, bitová hloubka) zda odpovídají zadání
- kontrola informací o barevném prostoru (zda jsou správné a kompletní)

⁹⁷ Iniciativa FADGI doporučuje takto důkladněji zkontrolovat alespoň 10 obrazů nebo 10% z každé vzniklé dávky obrazů.

- kontrola obrazu vůči analogové předloze (rozměry, rozlišení, orientace, dopad ořezů, kompletnost dokumentu)
- kontrola kvality obrazu (jas, kontrast, barevná věrnost, šum, artefakty, míra detailu apod.); (FADGI, 2016, s. 87-89).

5.5 Validace balíčku SIP

Validace balíčku SIP je komplexní validace, která může zahrnovat všechny výše uvedené typy validací, a dále také kontrolu struktury balíčků (např. přítomnost předepsaných souborů a adresářů).

Základním prvkem validace je kontrola úplnosti a neporušenosti souborů před odesláním balíčku SIP do repozitáře. Tu lze jednoduše provést užitím souboru MD5 v balíčku SIP, který obsahuje digitální otisky i výčet všech souborů).

Důležité je však užit komplexní nástroj – pro Standard NDK jím může být Komplexní validátor, s omezením, které byly uvedeny výše (viz Metadatová validace). Optimálně by měl producent digitalizátů mít k dispozici také vlastní nástroj, který dokáže provést základní kontrolu struktury balíčku SIP.

Terminologie

AIP balíček - archivní informační balíček (*archival information package*); informační balíček složený z informačního obsahu a přidružených archivních informací a uchovávaný v archivu OAIS (dle normy ISO 14721:2012)

Archivační formát - formát digitalizátu, který je považován za aktuálně vhodný pro zajištění dlouhodobého uchování

Archivní kopie - konečný digitalizát v archivačním formátu

Autenticita digitálního dokumentu - vlastnost digitálního dokumentu, podle které může osoba nebo systém posoudit, že je digitální dokument tím, za co se vydává

Balíčkovací migrace - digitální migrace, při níž dochází ke změně bitů balíčkovacích informací (dle normy ISO 14721:2012)

Barevný model - způsob číselné specifikace barev za užití tří nebo více kanálů (např. RGB, CMYK nebo YCbCr)

Barevný profil (ICC profil) - soubor dat, která charakterizují zařízení pořizující obrazový snímek nebo výstupní barevný prostor podle standardů International Color Consortium

Barevný prostor - geometrická reprezentace barev v prostoru, který lze vizuálně vnímat nebo vytvářet užitím konkrétního barevného modelu - např. sRGB nebo Adobe RGB (dle FADGI, 2017)

Bitová hloubka - stanovuje maximální počet odstínů nebo barev v digitálním obrazovém souboru (dle FADGI, 2010)

Bitová ochrana - principem bitové ochrany je ochrana a zachování digitálního objektu v takovém stavu, v jakém byl uložen

Bitový tok - data v rámci jednoho souboru, která mají smysluplné společné vlastnosti pro archivační účely (dle PREMIS)

ČIDLO/CZIDLO - Český systém pro IDentifikaci a LOkalizaci dokumentů českého kulturního dědictví (*CZech IDentification and LOcalization tool*); systém, který zabezpečuje trvalou identifikaci digitálních dokumentů českého národního dědictví za využití mezinárodního standardu URN:NBN

Data - opakovaně interpretovatelná reprezentace informací ve formalizované podobě vhodné pro komunikaci, interpretaci nebo zpracování (dle ISO 14721:2012); jednotlivý objekt spadající to kategorie dat nazýváme “datový objekt”

Datový objekt s obsahem (*Content Data Object, CDO*) - datový objekt, který je spojený s relevantními interpretačními informacemi a nese tak informační obsah; může být tvořen jedním nebo více soubory

Definice metadatových formátů - jednotlivé metadatové specifikace vydávané a dále udržované Národní knihovnou ČR, které v souhrnu tvoří Standard NDK (viz níže); v době vydání této metodiky existují Definice metadatových formátů samostatně pro monografie, pro periodika, pro gramodesky, pro fonoválčky a pro 3 typy e-born dokumentů (e-monografie, e-periodika a skládaná periodika)

Digitalizace - převod fyzické předlohy do digitální podoby

Digitalizační zařízení - fyzické zařízení zajišťující digitalizaci předlohy

Digitalizát - digitální dokument vzniklý digitalizací fyzické předlohy (tištěné, zvukové)

Digitální migrace - přesun digitálních informací v rámci archivu OAIS s cílem tyto informace uchovat; dále je možné digitální migraci rozdělit na čtyři typy - renovace, replikace, transformace a balíčková migrace (dle normy ISO 14721:2012)

Digitální otisk - mechanismus pro kontrolu neporušenosti digitálních objektů, např. MD5, SHA-1

DIP balíček - výstupní informační balíček (*dissemination information package*); informační balíček odvozený z jednoho nebo více balíčků AIP a zasláný archivem OAIS koncovému uživateli jako odpověď na jeho požadavek vůči tomuto archivu (dle normy ISO 14721:2012)

DROID (*Digital Record Object Identification*) - softwarový nástroj pro formátovou identifikaci souborů, je provázaný s formátovým registrem PRONOM

Emulace - napodobení činnosti jednoho zařízení nebo systému jiným zařízením (systémem); z hlediska digitální archivace jde o ochranné opatření, pomocí kterého je možné znovuvytvořit softwarové a hardwarové prostředí, které je nutné k obnovení digitálního dokumentu

Formát - specifické uspořádání datových a strukturních složek souboru, které umožňuje jeho zpracování a zobrazení jeho obsahu

Formátová identifikace - jednoznačné určení formátu, optimálně prostřednictvím

identifikátoru PUID

Formátový identifikační nástroj - nástroj pro formátovou identifikaci

Formátová normalizace - jedna ze základních strategií dlouhodobé archivace; pokud se pod pojmem formátová migrace tradičně rozumí převod digitálního objektu z jednoho formátu do jiného za účelem zajištění jeho použitelnosti a uchování jeho obsahu, pak pojem formátová normalizace znamená, že jde o dlouhodobou a plánovanou aktivitu, která slouží k tomu, aby všechny digitální objekty byly uloženy v otevřených a dlouhodobě udržitelných formátech

Formátová politika - soubor pravidel, kterými se řídí určitá instituce při přijímání digitálních dokumentů k dlouhodobému uložení a kterými stanovuje formáty a případné další parametry souborů, které k uložení přijímá

Formátový profil - nastavení v rámci konkrétního formátu (zahrnuje např. výběr typu komprese a další specifická nastavení)

Formátový registr - registr interpretačních informací potřebných pro reprodukci a zpracování digitálních objektů; schraňuje zejm. informace o formátech, souvisejících aplikacích a dalších prvcích počítačového prostředí

Fyzická strukturální mapa - reprezentuje dokument otištěný do digitální reprezentace, zachycuje fyzickou posloupnost jednotlivých stran a odkazuje na jednotlivé komponenty v SIP balíčku

Generace (obrazových) dat - množiny (obrazových) dat, které vznikají během digitalizace, následného zpracování i archivace; každá transformace, kterou data prochází od vytvoření původních snímků po finální data, znamená vytvoření nové generace dat

Granularita - dělení základní intelektuální entity na dílčí části, které slouží jako úroveň popisu v metadatovém záznamu

Charakterizace - proces zjištění informací (zejména technických informací, např. informace o formátu) o souboru, které jsou extrahovány přímo ze souboru

Identifikační informace - informace, které plní funkci identifikátoru informačního obsahu (dle normy ISO 14721:2012); v prostředí knihoven jde např. o perzistentní identifikátory a bibliografický popis

Informace - jakékoli znalosti, které mohou být předmětem výměny, při této výměně jsou informace vždy reprezentovány určitým typem dat (dle normy ISO 14721:2012); jednotlivý

objekt spadající do kategorie informací nazýváme “informační objekt”

Informace o neporušenosti - informace, které udávají, jak je zajištěno, aby objekt s informačním obsahem nebyl změněn nezdokumentovaným způsobem (dle normy ISO 14721:2012); jde např. o digitální otisk

Informace o přístupových právech - informace udávající omezení, která se týkají přístupu k informačnímu obsahu včetně právního rámce, licenčních podmínek a řízení přístupu (dle normy ISO 14721:2012)

Informace o zabalení - slouží k propojení a popisu součástí informačního balíčku (dle normy ISO 14721:2012)

Informační balíček - logická schránka, která může obsahovat informační obsah a archivační informace; dle normy ISO 14721:2012 je možné rozlišit tři typy informačního balíčku: vstupní informační balíček (SIP), archivní informační balíček (AIP) a výstupní archivační balíček (DIP)

Informační obsah - množina informací, která je původním předmětem uchování, nebo která obsahuje část těchto informací nebo všechny tyto informace; informační obsah je informační objekt složený z datového objektu s obsahem a ze svých interpretačních informací (dle normy ISO 14721:2012)

Intelektuální entita - jednotlivý intelektuální nebo umělecký výtvar, který je považován za relevantní pro cílovou komunitu v kontextu digitální archivace (dle PREMIS); v praxi Standardu NDK je základní intelektuální entitou smysluplný celek, jehož reprezentace je obsahem balíčku SIP, tj. v případě monografií je to svazek, v případě periodik jednotlivé číslo a v případě zvukových dokumentů jeden nosič, resp. soubor nosičů, pokud tvoří jeden celek

Interpretační informace - informace, které převádějí datový objekt do smysluplnějších významových celků (dle normy ISO 14721:2012)

JP2 - rastrový formát rodiny JPEG 2000 podle Part I. specifikace JPEG 2000

Kontextuální informace - informace, které dokládají vztah informačního obsahu k jeho okolí; jde např. o důvod vytvoření informačního obsahu a o jeho vztah k dalším objektům s informačním obsahem (dle normy ISO 14721:2012)

Kontrola neporušenosti - kontrola neporušenosti souboru užitím digitálního otisku

Logická strukturální mapa - představuje hierarchický model reprezentace dokumentu, založený na fyzické předloze; umožňuje interpretovat posloupnost celků zdigitalizovaného

dokumentu např. od svazku až po jednotlivý článek

LTP úložiště - systém pro dlouhodobé uložení (*Long-Term Preservation*) digitálních dokumentů; jeho cílem by měla být ochrana digitálních dokumentů, která zajistí jejich bezpečné uložení, integritu a autenticitu v dlouhodobém horizontu

Matematicky bezeztrátová komprese - komprese odpovídající konceptu vratné transformace; kompresní algoritmus umožňuje snížení požadavků na úložnou kapacitu, ale zároveň jsou po dekomprimaci data totožná s daty před kompresí

Metadatový extraktor - nástroj pro charakterizaci, jenž ze souboru získá především jeho technické vlastnosti

Metadatový aplikační profil - soubor metadatových elementů, které byly vybrány z jednoho nebo více mezinárodních metadatových standardů a sloučeny do jednoho schématu, které zachovává interoperabilitu s původními mezinárodními standardy, ale zároveň vyhovuje specifickým požadavkům pro konkrétní užití (nebo např. reflektuje národní specifika v případě národního standardu); součástí metadatového profilu může být i vytvoření vlastní sady metadatových elementů

Nevratná transformace - digitální migrace, u které nelze zaručit zachování informačního obsahu v úplnosti; podrobněji viz heslo transformace

Obrazová matrice - komprimovaná data v nejvyšší možné kvalitě, která slouží jako zdroj pro vytváření obrazových dat pro různé účely, v různých formátech a různé kvalitě; v oblasti digitalizace dokumentů plní funkci obrazové matrice archivační formát

OCR - optické rozpoznávání znaků (*Optical Character Recognition*, OCR) je metoda pro získávání textového obsahu z obrazových souborů

Perzistentní identifikátor - znakový řetězec, který jednoznačným způsobem označuje nějaký (digitální) objekt, a který je trvale užíván k identifikaci právě toho (digitálního) objektu, kterému byl na počátku přidělen; perzistentní identifikátory mají většinou standardizovanou syntax a jsou přidělovány na základě jasně stanovených pravidel

Popis balíčku (*package description*) - informace určené pomůckám pro zpřístupnění (dle normy ISO 14721:2012)

Popisné informace (*descriptive information*) - množina informací, která je složena především z popisů balíčků a je poskytována správě dat za účelem podpory koncových uživatelů při objednávání a získávání informačních jednotek z archivu OAIS (dle normy ISO 14721:2012)

Prostorové rozlišení - určuje množství informací v rastrovém souboru z hlediska počtu pixelů na jednotku měření, kterou obvykle bývá palec, tj. určuje, jak jsou jednotlivé pixely blízko u sebe (dle FADGI, 2010)

Prezentační formát - formát, ve kterém je digitalizát zpřístupňován uživatelům (např. v digitální knihovně)

Prezentační meziformát - meziformát, ze kterého digitální knihovna generuje cílový prezentační formát

Provenienční informace - informace, které dokumentují historii informačního obsahu; vypovídají o jeho původu nebo zdroji, stejně tak jako o veškerých změnách, které mohly od jeho vzniku nastat, a o tom, kdo o něj od jeho vzniku pečoval (dle normy ISO 14721:2012)

Prvotní digitalizát - data, která jsou bezprostředním výstupem digitalizačního zařízení, tj. data, která digitalizační zařízení uloží do souboru nebo souborů na datový nosič po skončení procesu snímání

Předloha - fyzický objekt, ze kterého je digitalizačním zařízením vytvářen digitalizát

PSP balíček – produkční balíček (*Producer Submission Package*); balíček dat a metadat, který přichází od producenta dat, po konverzích a kontrolách dat a metadat z něj vzniká balíček SIP; pokud jsou při tvorbě PSP balíčku dodržovány aktuální standardy DMF, je PSP balíček shodný se SIP balíčkem

PUID - identifikátor registru PRONOM jednoznačně označující formát i jeho jednotlivé verze případně jiné parametry formátu (typ komprese apod.)

Původní snímek - prvotní digitalizát, který vznikl snímáním

RAW - třída obrazových formátů, které obsahují minimálně zpracovaná data ze senzoru; jedná se o proprietární formáty vlastněné výrobcem fotoaparátů

Renovace - digitální migrace, při níž je jedna instance datového nosiče nahrazena jinou instancí datového nosiče stejného typu, a to zkopírováním bitů na datový nosič, využitý k umístění balíčků a ke správě a přístupu k datovému nosiči; jde o specifický případ replikace; v případě renovace ale není třeba měnit mapovací strukturu úložiště a balíček lze stále v úložišti nalézt (dle normy ISO 14721:2012)

Replikace - digitální migrace, při níž nedochází k žádným změnám balíčkovacích informací, informačního obsahu ani archivačních informací; na rozdíl od renovace může vyžadovat změny

mapovací struktury archivního úložiště (dle normy ISO 14721:2012)

Reprezentace - množina souborů potřebná pro úplnou reprodukci intelektuální entity (dle PREMIS)

Resolver - on-line aplikace, která zajišťuje přidělování a správu perzistentních identifikátorů založených na standardu URN:NBN; přičemž resolver má také funkci přesměrovávací služby, která zajistí přesměrování webového prohlížeče z URN:NBN na aktuální URL adresu digitálního dokumentu

Signifikantní vlastnosti - vlastnosti datového objektu s obsahem (CDO), které je třeba v průběhu času udržovat, aby byla zabezpečena trvalá přístupnost, použitelnost a význam objektu

SIP balíček - vstupní informační balíček (*submission information package*); informační balíček, který dodává producent do archivu OAIS tak, aby mohl být využit při sestavení nebo aktualizaci jednoho nebo více AIP balíčků a/nebo přidružených popisných informací (dle normy ISO 14721:2012)

Snímací zařízení - digitalizační zařízení pro digitalizaci tištěných předloh (skener, fotoaparát)

Snímání - digitalizace užitím snímacího zařízení

Snímkový formát - formát, ve kterém je uložen původní snímek

Snímek - digitalizát vzniklý snímáním

Software pro zobrazení interpretačních informací - software, který reprodukuje interpretační informace v podobě, jež je srozumitelná lidem (dle normy ISO 14721:2012)

Soubor - pojmenovaná a uspořádaná posloupnost bajtů, kterou operační systém dokáže rozpoznat a která má určitý formát (dle PREMIS)

Standard NDK - zastřešující pojem, který zahrnuje standardy pro digitální dokumenty (nebo také Definice metadatových formátů; tzv. DMF), které vydává Národní knihovna ČR; Standard NDK slouží jako jednotný formát pro paměťové instituce, které chtějí svá data dlouhodobě uchovávat v úložišti NK ČR.

Technická metadata - informace o vlastnostech souboru, jako je velikost, formát, komprese, zařízení, kterým byl vytvořen, obrazové a další specifické vlastnosti

Transformace - digitální migrace, při které dochází ke změnám bitů informačního obsahu nebo archivačních informací, přičemž je ale snaha zachovat informační obsah v úplnosti

Transformační šablona - šablona (příp. serie šablon), které umožňují převod a zpracování záznamů z jednoho formátu do jiného; v praxi NDK jde zejm. o převody katalogizačních záznamů zpracovaných podle pravidel AACR2 nebo RDA do formátu MODS

UNICODE - mezinárodní standard pro kódování znaků

Uživatelská kopie - konečný digitalizát v prezentačním formátu

Validace - automatická kontrola toho, zda jsou data nebo metadata vytvořena v souladu s deklarovanými specifikacemi

Validátor - softwarový nástroj, který je vytvořený za účelem kontroly dat, metadatových záznamů nebo souborových formátů; v prostředí NDK se používá volně dostupný open-source nástroj Komplexní validátor, který je určený pro kontrolu SIP/PSP balíčků (tj. zda balíčky odpovídají stanoveným předpisům podle Standardu NDK); pro validaci formátů, tj. pro kontrolu zda daný souborový formát odpovídá specifikaci formátu se doporučuje užívat volně dostupné open source nástroje jako JHOVE nebo JPYLYZER

Vizuálně bezztrátová komprese - matematicky ztrátová komprese, u které běžný pozorovatel při zobrazení výstupu nerozezná změny oproti výstupu matematicky bezztrátové komprese

Vizuálně ztrátová komprese - matematicky ztrátová komprese, která přináší vizuálně patrné změny obrazové kvality

Vratná transformace - digitální migrace, u které je možná bezztrátová zpětná transformace (dle normy ISO 14721:2012); podrobněji viz heslo transformace

ZIP - formát pro archivaci a sdílení bezztrátově komprimovaných dat

Znalostní základna - množina informací, které si osvojila osoba nebo systém a která této osobě nebo systému umožňuje porozumět přijímaným informacím (dle normy ISO 14721:2012)

Zpřístupňovací software - software, který dokáže prezentovat informační obsah nebo jeho část, slouží tedy jako prostředek pro zpřístupnění informačního obsahu cílové komunitě (dle normy ISO 14721:2012)

Seznam zkratk

AACR2	Anglo-American Cataloguing Rules, 2nd Edition
ADR	Centrální adresář knihoven a informačních institucí v ČR
AIIM	Association for Intelligent Information Management
AIP	Archival Information Package
ALA	American Library Association
ALTO	Analyzed Layout and Text Object
ANSI	American National Standards Institute
ASCII	American Standard Code for Information Interchange
CD	Compact Disc
CDO	Content Data Object
CIPA	Camera & Imaging Products Association
CZIDLO	Czech Identification and Localization Tool
čČNB	číslo České národní bibliografie
ČIDLO	Český identifikační a lokalizační systém
ČNB	Česká národní bibliografie
ČR	Česká republika
ČSN	Česká technická norma (původně Československá státní norma)
ČSSR	Československá socialistická republika
DC	Dublin Core
DCE	Distributed Computing Environment
DCMI	Dublin Core Metadata Initiative
DIP	Dissemination Information Package
DLF	Digital Library Federation
DMF	Definice metadatových formátů

DNG	Digital Negative
DPI	dots per inch
DROID	Digital Record Object Identification
ECI	European Color Initiative
EPH	End of Packet Header
EPUB	electronic publication
EU	Evropská unie
EXIF	Exchangeable Image File Format
FADGI	Federal Agencies Digital Guidelines Initiative
FRBR	Functional Requirements for Bibliographic Records
GDFR	Global Digital Format Registry
GIF	Graphics Interchange Format
GUID	Globally Unique Identifier
ICC	International Color Consortium
IEC	International Electrotechnical Commission
IETF	Internet Engineering Task Force
ISBN	International Standard Book Number
ISMN	International Standard Music Number
ISO	International Organisation for Standardization
ISSN	International Standard Serial Number
ITU	International Telecommunication Union
ITU-T	ITU Telecommunication Standardization Sector
JFIF	JPEG File Interchange Format
JHOVE	JSTOR/Harvard Object Validation Environment
JP2	JPEG 2000
JPEG	Joint Photographic Experts Group

JSTOR	Journal Storage
LOC	Library of Congress
LTP	Long-term Preservation
LZW	Lempel-Ziv-Welch (kompresní algoritmus)
MAC	Media Access Control
MARC	Machine-readable Cataloging
MC	Master Copy
MD5	Message Digest Algorithm 5
METS	Metadata Encoding and Transmission Standard
MIME	Multipurpose Internet Mail Extension
MIX	Metadata for Images in XML Standard
MODS	Metadata Object Description Schema
NCSA	National Center for Supercomputing Applications
NDK	Národní digitální knihovna
NISO	National Information Standards Organization
NK ČR	Národní knihovna České republiky
OAIS	Open Archival Information System
OCLC	Online Computer Library Center
OCR	Optical Character Recognition
ODIF	Odbor digitálních fondů
ONDS	Odbor novodobých digitálních sbírek
OS	operační systém
OSF	Open Software Foundation
PDF	Portable Document Format
PDF/A	PDF pro Archivaci
PNG	Portable Network Graphics

PPI	pixels per inch
PREMIS	Preservation Metadata Implementation Strategies
PRONOM	Public Record Office and Nôm (formátový registr)
PS	původní sken
PUID	PRONOM's Persistent Unique Identifier
QI	Quality Index
RDA	Resource Description and Access
RFC	Request for Comments
RGB	red-green-blue
SIP	Submission Information Package
SOP	Start of Packet Header
sRGB	standardní RGB
TIFF	Tagged Image File Format
TLM	Tile Length Markers
UDFR	Unified Digital Format Registry
URL	Universal Resource Locator
URN:NBN	Uniform Resource Name: National Bibliography Number
USA	United States of America
UTF	Unicode Transformation Format
UUID	Universally Unique Identifier
VISK	Veřejné informační služby knihoven
XML	eXtensible Markup Language
XSD	XML Schema Definition
XSLT	Extensible Stylesheet Language Transformations

Příloha – mapování výstupů metadatových extraktorů do metadat balíčků SIP

Mapování výstupů nástrojů

Doporučení pro mapování výstupů nástroje jpylyzer, JHOVE do metadat k souborům původních snímků (TIFF) a archivním kopiím v bezztrátovém formátu JP2.

Nástroj jpylyzer

	Mapování
metadata MIX	
<fileSize>	fileInfo/fileSizeInBytes
<imageWidth>	properties/jp2HeaderBox/imageHeaderBox/width
<imageHeight>	properties/jp2HeaderBox/imageHeaderBox/height
<colorSpace>	properties/ jp2HeaderBox/colourSpecificationBox/enumCS (u properties/ jp2HeaderBox/colourSpecificationBox/meth =Enumerated) nebo properties/jp2HeaderBox/colourSpecificationBox/icc/colourSpace u properties/jp2HeaderBox/colourSpecificationBox/meth =Restricted)
<iccProfileName>	properties/ jp2HeaderBox/colourSpecificationBox/description
<iccProfileVersion>	properties/ jp2HeaderBox/colourSpecificationBox/icc/profileVersion
<codec>	properties/contiguousCodestreamBox/com/comment
<codecVersion>	properties/contiguousCodestreamBox/com/comment
<codestreamProfile>	properties/contiguousCodestreamBox/siz/rsiz ⁹⁸

⁹⁸ Pokud hodnota zde odpovídá číslu 1 pak se jedná o Profile 0 a do metadat se vyplní P0, pokud rsiz obsahuje hodnotu 2, pak se jedná o profil 1 a do metadat se vyplní hodnota P1, pokud rsiz obsahuje hodnotu "ISO/IEC 15444-1" pak se nejedná o žádný profil a žádné omezení, jedná se o profil 2 a vyplní se hodnota P2.

<tileWidth>	properties/ contiguousCodestreamBox/siz/xTsiz
<tileHeight>	properties/ contiguousCodestreamBox/siz/yTsiz
<qualityLayers>	properties/ contiguousCodestreamBox/cod/layers
<resolutionLevels>	properties/ contiguousCodestreamBox/cod/levels
<samplingFrequencyUnit>	properties/jp2HeaderBox/resolutionBox/ <i>hodnota podle toho jaké rozlišení se plní, nejsnazší je nastavit in a brát rozlišení z vRescInPixelsPerInch a hRescInPixelsPerInch</i>
<xSamplingFrequency>	Kontejnerový element, neobsahuje konkrétní hodnotu ale níže uvedené elementy, ty lze plnit z vícero elementů viz níže.
Varianta zápisu 1:	
<numerator>	properties/jp2HeaderBox/resolutionBox/captureResolutionBox/hResdInPixelsPerInch
<denominator>	Vždy 1
Varianta zápisu 2:	
<numerator>	properties/jp2HeaderBox/resolutionBox/captureResolutionBox/hRcN
<denominator>	Hodnota elementu <denominator> se získá výpočtem mezi hodnotami elementů v kontejnerovém element properties/jp2HeaderBox/resolutionBox/captureResolutionBox: hRcN / hRcD x 10 ^{hRcE} x 0,0254 (v palcích)
<ySamplingFrequency>	Kontejnerový element, neobsahuje konkrétní hodnotu ale níže uvedené elementy, ty lze plnit z vícero elementů viz níže.
Varianta zápisu 1:	
<numerator>	properties/jp2HeaderBox/resolutionBox/captureResolutionBox/vResdInPixelsPerInch

<denominator>	Vždy 1
Varianta zápisu 2:	
<numerator>	properties/jp2HeaderBox/resolutionBox/captureResolutionBox/vRcN
<denominator>	<i>Hodnota elementu <denominator> se získá výpočtem mezi hodnotami elementů v kontejnerovém elementu</i> properties/jp2HeaderBox/resolutionBox/captureResolutionBox: $vRcN / vRcD \times 10^{vRcE} \times 0,0254$ (v palcích)
<bitsPerSampleValue >	properties/contiguousCodestreamBox/siz/ssizDepth
<bitsPerSampleUnit>	properties/ jp2HeaderBox/imageHeaderBox/bPCDepth
<samplesPerPixel>	properties/ jp2HeaderBox/imageHeaderBox/nC

Nástroj JHOVE

	Mapování	
MIX		
<fileSize>	jhove/repInfo/size	JP2 TIFF
<formatName>	jhove/repInfo/format	JP2 TIFF
<formatVersion>	jhove/repInfo/version	JP2 TIFF
<byteOrder>	mix/byteOrder	JP2 TIFF
<compressionScheme>	properties/property/name=Transformation/value ⁹⁹ <i>pro formát TIFF:</i> mix/BasicDigitalObjectInformation/compressionScheme	JP2 TIFF
<imageWidth>	properties/property/name=XSize/value	JP2 TIFF
<imageHeight>	properties/property/name=YSize/value	JP2 TIFF
<colorSpace>	properties/property/property/name=EnumCS/value <i>pro formát TIFF:</i> mix/BasicImageInformation/BasicImageCharacteristics/PhotometricInterpretation/colorSpace	JP2 TIFF
<tileWidth>	properties/property/name=Codestream/property/name=XTSize/value	JP2
<tileHeight>	properties/property/name=Codestream/property/name=YTSize/value	JP2
<qualityLayers>	properties/property/name=NumberOfLayers/value	JP2
<resolutionLevels>	properties/property/name=NumberDecompositionLevels/value	JP2
<scannerManufacturer>	mix/ImageCaptureMetadata/ScannerCapture/scannerManufacturer	TIFF

⁹⁹ Value je číslo, kde 1=lossless, tj. 5-3 reversible a 0=lossy, tj. 9-7 irreversible.

<scannerModelName>	mix/ImageCaptureMetadata/ScannerCapture/ScannerModel/scannerModelName	TIFF
<scannerModelSerialNo>	mix/ImageCaptureMetadata/ScannerCapture/ScannerModel/scannerModelName	TIFF
<samplingFrequencyUnit>	property/name=NisoImageMetadata/mix/ImageAssessmentMetadata/SpatialMetrics/samplingFrequencyUnit	TIFF
<scanningSoftwareName>	mix/ImageCaptureMetadata/ScannerCapture/ScanningSystemSoftware/scanningSoftwareName	TIFF
<scanningSoftwareVersionNo>	mix/ImageCaptureMetadata/ScannerCapture/ScanningSystemSoftware/scanningSoftwareName <i>číslo verze je součástí tohoto elementu</i>	TIFF
<orientation>	mix/ImageCaptureMetadata/orientation	TIFF
<samplingFrequencyUnit>	mix/ImageAssessmentMetadata/SpatialMetrics/samplingFrequencyUnit	JP2 TIFF
<xSamplingFrequency>		JP2 TIFF
<numerator>	Properties/property/name=VertResolution/property/name=Numerator/value <i>nebo</i> property/name=NisoImageMetadata/mix/ImageAssessmentMetadata/SpatialMetrics/xSamplingFrequency/numerator ¹⁰⁰ <i>Pro TIFF:</i> mix/ImageAssessmentMetadata/SpatialMetrics/xSamplingFrequency/numerator	JP2 TIFF
<denominator>	Properties/property/name=VertResolution/property/name=Denominator/value <i>nebo</i> property/name=NisoImageMetadata/mix/ImageAssessmentMetadata/SpatialMetrics/xSamplingFrequency/denominator <i>Pro TIFF:</i> mix/ImageAssessmentMetadata/SpatialMetrics/xSamplingFrequency/denominator	JP2 TIFF

¹⁰⁰ Informace o rozlišení obrazu se může nacházet na dvou místech, dle toho, jaký kodek obraz konvertoval.

<ySamplingFrequency>		
<numerator>	Properties/property/name=HorizResolution/property/name=Numerator/value <i>nebo</i> property/name=NisoImageMetadata/mix/ImageAssessmentMetadata/SpatialMetrics/ySamplingFrequency/numerator <i>Pro TIFF:</i> mix/ImageAssessmentMetadata/SpatialMetrics/ySamplingFrequency/numerator	JP2 TIFF
<denominator>	Properties/property/name=HorizResolution/property/name=Denominator/value <i>nebo</i> property/name=NisoImageMetadata/mix/ImageAssessmentMetadata/SpatialMetrics/ySamplingFrequency/denominator <i>Pro TIFF:</i> mix/ImageAssessmentMetadata/SpatialMetrics/ySamplingFrequency/denominator	JP2 TIFF
<bitsPerSampleValue>	mix:ImageColorEncoding/mix:BitsPerSample/mix:bitsPerSample Value	JP2 TIFF

Je možné využít i nástroje další. Například nástroj FITS agreguje několik nástrojů, mezi nimi i nástroje JHOVE a jpylyzer, stejně tak nástroj Kost-Val. Průběžně vznikají nové nástroje a inovují se existující. Výše uvedené nástroje JHOVE a jpylyzer jsou komunitou okolo digitální archivace hojně používané a průběžně aktualizované (zvl. JHOVE), považujeme je tedy za nástroje spolehlivé.

O autorech

1) Metodika verze 1.0 (2019)

Hlavním autorem této metodiky je Ladislav Cubr za spolupráce s Natalií Ostrákovou a Pavlínou Kočišovou.

2) Metodika verze 2.0 (2023)

Autory aktualizace metodiky jsou Květa Fremrová, Václav Jiroušek, Pavlína Kočišová, Vojtěch Kopský, Ivo Miláček, a Filip Pavčík. Všechny využití texty vznikly v rámci institucionálního výzkumu Národní knihovny České republiky.

Citovaná literatura

ADOBE RGB (1998) *Color Image Encoding*, 2005 [online]. Version 2005-05. San Jose (CA): Adobe Systems Inc. [cit. 2023-06-26]. Dostupné z:

<https://www.adobe.com/digitalimag/pdfs/AdobeRGB1998.pdf>

ALTO Principles, 2016. The Library of Congress [online]. Washington (DC): The Library of Congress [cit. 2023-06-22]. Dostupné z: <https://www.loc.gov/standards/alto/description.html>

ALTO: *Technical Metadata for Layout and Text Objects*, 2022. The Library of Congress [online]. Washington (DC): The Library of Congress [cit. 2023-06-22]. Dostupné z: <https://www.loc.gov/standards/alto/>

ANSI/NISO Z39.87-2006, 2006. *Data Dictionary – Technical Metadata for Digital Still Images*. Bethesda (MD): NISO Press, xiv, 107 s. ISBN 978-1-937522-37-7. ISSN 1041-5653.

ARCLib - *komplexní řešení pro dlouhodobou archivaci digitálních (knihovných) sbírek* [online], 2015-2020. Praha: Knihovna AV ČR [cit. 2023-06-23]. Dostupné z: <https://arclib.cz/>

BLOOD, George, 2011. *Refining Conversion Contract Specifications: Determining Suitable Digital Video Formats for Medium-term Storage* [online]. Washington (DC): FADGI [cit. 2023-06-23]. Dostupné z: https://www.digitizationguidelines.gov/audio-visual/documents/IntrmMastVidFormatRecs_20111001.pdf

BROWN, Adrian, 2006. *The PRONOM PUID Scheme: a scheme of persistent unique identifiers for representation information* [online]. London: National Archives, 9 s. [cit. 2023-06-23]. Digital Preservation Technical Paper, issue 2. Dostupné z: https://www.nationalarchives.gov.uk/aboutapps/pronom/pdf/pronom_unique_identifier_scheme.pdf

BRUYS, Alix, Bertrand CARON, GRANDCOLAS, Yannick a LEDOUX, Thomas, 2019. *JPEG Got the Blues: Properly Rendering 32-bits JPEG*. *Open Preservation Foundation* [online]. The Open Preservation Foundation [cit. 2023-06-21]. Dostupné z: <https://openpreservation.org/blog/2019/11/07/jpeg-got-the-blues/>

BRYGFJELD, Svein Arne, 2010. *JP2K for preservation and access, experiences from the National Library of Norway* [online]. The National Library of Norway, 16 Nov 2010 [cit. 2023-06-23]. Dostupné z: <http://www.dpconline.org/docman/miscellaneous/events/521-jp2knov2010brygfjeld/file>.

BUCKLEY, Robert, 2008. *JPEG 2000 - a Practical Digital Preservation Standard?* [online]. Glasgow: Digital Preservation Coalition, 21 s. [cit. 2023-06-21]. DPC Technology Watch Series

Report, 08-01. Dostupné z: <http://www.dpconline.org/docman/technology-watchreports/87-jpeg-2000-a-practical-digital-preservation-standard/file>.

BUCKLEY, Robert, 2009. *JPEG 2000 as a Preservation and Access Format for the Wellcome Trust Digital Library* [online]. Edited by Simon Tanner. London: King's College London, 17 s. [cit. 2023-06-15]. Dostupné také z:

<https://web.archive.org/web/20150319030938/http://wellcomelibrary.org:80/content/documents/22082/JPEG2000-preservation-format.pdf>

BUONORA, Paolo a LIBERATI, Franco, 2008. A Format for Digital Preservation of Images. *D-Lib Magazine* [online]. **14**(7/8), [cit. 2023-06-23]. DOI: 10.1045/july2008-buonora. Dostupné z: <http://www.dlib.org/dlib/july08/buonora/07buonora.html>.

CAPLAN, Priscilla, 2018. *PREMIS a jak mu porozumět*. Praha: Univerzita Karlova, 26 s.

CIPA, 2012. *Exchangeable image file format for digital still cameras, Exif Version 2.3* [online]. Tokyo: CIPA DC-008-2012, 185 s. [cit. 2023-06-23]. Dostupné z:

https://www.cipa.jp/std/documents/e/DC-008-2012_E.pdf

CUBR, Ladislav, 2010. *Dlouhodobá ochrana digitálních dokumentů*. Praha: Národní knihovna ČR, 154 s. ISBN 978-80-7050-588-5.

CUBR, Ladislav, 2017. *Autenticita a digitální informace* [online]. Praha. Disertační práce. Univerzita Karlova, Filozofická fakulta, Ústav informačních studií a knihovnictví, 227 s.

Vedoucí práce Jiří IVÁNEK [cit. 2023-06-23]. Dostupné z:

<https://dspace.cuni.cz/bitstream/handle/20.500.11956/94159/140057523.pdf?sequence=1&isAllowed=y>

CUBR, Ladislav, LODROVÁ, Iveta, ŘEHÁNEK, Martin a VAŠEK, Zdeněk, 2016. Srovnání vybraných národních identifikačních systémů užívajících identifikátory URN:NBN. *ProInflow: časopis pro informační vědy* [online]. Brno: Masarykova univerzita, Filozofická fakulta, **8**(1) [cit. 2023-06-15]. DOI: <https://doi.org/10.5817/ProIn2016-1-3>. ISSN 1804–2406. Dostupné z:

<http://www.phil.muni.cz/journals/index.php/proinflow/article/view/2016-1-3>

CUBR, Ladislav, OSTRÁKOVÁ, Natalie a KOČIŠOVÁ, Pavlína, 2019. *Metodika pro tvorbu balíčků SIP se zaměřením na digitalizáty tištěných dokumentů* [online]. Praha: Národní knihovna ČR [cit. 2023-06-26]. Dostupné z: <http://www.nusl.cz/ntk/nusl-432324>

ČESKO, 2015. Ministerstvo kultury, Odbor umění, literatury a knihoven. Veřejné informační služby knihoven (VISK): podprogram č. 7: národní program ochrany a digitalizace dokumentů ohrožených degradací kyselého papíru - KRAMERIUS. *Veřejné informační služby knihoven* [online]. Praha: Národní knihovna ČR, 16. 9. 2015, 8 s. [cit. 2023-06-20]. Dostupné z:

<http://visk.nkp.cz/dokumenty/visk7/2016/VISK7-podm2016.doc>

DCMI USAGE BOARD, 2020. DCMI Metadata Terms [online]. *The Dublin Core Metadata Initiative* [cit. 2023-06-23]. Dostupné z: <http://dublincore.org/specifications/dublin-core/dcmi-terms/2020-01-20/>

DUPLOY, Laurent, 2017. JPEG 2000 as a preservation format for digitization: lessons learned from a library. In: *Archiving2017: Final Program and Proceedings : May 15–18, 2017, Riga, Latvia*. Riga: Society for Imaging Science and Technology, s. 157–159. ISBN 978-0- 89208-326-8.

DURANTI, Luciana a THIBODEAU, Kenneth, 2006. The concept of record in Interactive, Experiential and Dynamic Environments: the view of InterPARES. *Archival Science*. **6**(1), s. 13 - 68. DOI: <https://doi.org/10.1007/s10502-006-9021-7>

DUVAL, Erik, HODGINS, Wayne, SUTTON, Stuart a WEIBEL, Stuart L., 2002. Metadata Principles and Practicalities. *D-Lib Magazine* [online]. **8**(4) [cit. 2023-06-20]. DOI: 10.1045/april2002-weibel. Dostupné z: <http://www.dlib.org/dlib/april02/weibel/04weibel.html>

DVOŘÁK, Tomáš, KOUCKÝ, Karel, ŠULC, Jaroslav a kol, 2015. *Metodika pro vytváření bezpečnostních kopií v digitální podobě, Verze 1.0* [online]. Praha: Národní archiv, Státní oblastní archiv v Praze [cit. 2023-06-26]. Dostupné z: <https://www.nacr.cz/wp-content/uploads/2019/05/metodika2015.pdf>

EMBEDDED METADATA WORKING GROUP – SMITHSONIAN INSTITUTION, 2010. *Basic Guidelines for Minimal Descriptive Embedded Metadata in Digital Images. Federal agencies digital guidelines initiative* [online]. Smithsonian Libraries & Archives [cit. 2023-06-23]. Dostupné z: <http://www.digitizationguidelines.gov/guidelines/GuidelinesEmbeddedMetadata.pdf>

FEDERAL AGENCIES DIGITIZATION GUIDELINES INITIATIVE, 2010. Still Image Working Group. *Technical Guidelines for Digitizing Cultural Heritage Materials: creation of raster image files* [online]. Washington (DC): FADGI, 96 s. [Cit. 2023-06-23]. Dostupné z: http://www.digitizationguidelines.gov/guidelines/FADGI_Still_Image_Tech_Guidelines_2010-08-24.pdf

FEDERAL AGENCIES DIGITIZATION GUIDELINES INITIATIVE, 2014. Still Image Working Group. *Raster Still Images for Digitization: A Comparison of File Formats . Part 3. Narrative and Summary Table* [online]. Washington (DC): FADGI, Revised 29th Aug 2014, 9 s. [cit. 2023-06-15]. Dostupné z: http://www.digitizationguidelines.gov/guidelines/FADGI_RasterFormatCompare_p3_20140829_r.pdf

FEDERAL AGENCIES DIGITIZATION GUIDELINES INITIATIVE, 2016. *Technical Guidelines for Digitizing Cultural Heritage Materials: creation of raster image files* [online].

Washington (DC): FADGI, September 2016 [Cit. 2023-06-24]. Dostupné z:

http://www.digitizationguidelines.gov/guidelines/FADGI_Still_Image_Tech_Guidelines_2016.pdf

FEDERAL AGENCIES DIGITIZATION GUIDELINES INITIATIVE, 2017. *Color space - Glossary. Federal Agencies Digital Guidelines Initiative* [online]. Washington (DC): FADGI

[cit. 2023-06-16]. Dostupné z: <http://www.digitizationguidelines.gov/term.php?term=colospace>

FEDERAL AGENCIES DIGITIZATION GUIDELINES INITIATIVE, 2023. *Technical Guidelines for Digitizing Cultural Heritage Materials, 3rd. edition* [online]. Washington (DC): FADGI, May 2023 [Cit. 2023-06-24]. Dostupné z:

https://www.digitizationguidelines.gov/guidelines/FADGI%20Technical%20Guidelines%20for%20Digitizing%20Cultural%20Heritage%20Materials_3rd%20Edition_05092023.pdf

FERNIE, Kate, ed., 2008. *Technical Guidelines for Digital Cultural Content Creation*

Programmes. Version 2.0 [online]. [Rome]: Minerva EC, 92 s. [cit. 2023-06-20] Dostupné z:

<http://www.minervaeurope.org/publications/MINERVA%20TG%202.0.pdf>.

FLORIDA CENTER FOR LIBRARY AUTOMATION, 2012. *Recommended Data Formats for Preservation Purposes in the Florida Digital Archive* [online]. FCLA, Mar. 2012 [cit. 2023-06-08]. Dostupné z:

<https://libraries.flvc.org/documents/181844/502298/Recommended+Data+Formats/>

Generation loss, 2023 [online], poslední aktualizace 5. června 2023 [cit. 2023-06-23],

Wikipedia. Dostupné z: https://en.wikipedia.org/wiki/Generation_loss

Global digital format registry, 2016 [online]. Harvard Library [cit. 2023-06-23]. Dostupné z:

https://web.archive.org/web/20171009011353/https://library.harvard.edu/preservation/digital-preservation_gdfr.html

HUTAŘ, Jan, 2012. *Digitalizace, popis pomocí metadat a jejich formáty* [online]. Praha.

Disertační práce. Univerzita Karlova. Filozofická fakulta. Ústav informačních studií a knihovnictví, 244 s. Vedoucí práce Stanislav KALKUS. [cit. 2023-06-26] Dostupné z:

<https://dspace.cuni.cz/bitstream/handle/20.500.11956/44181/140015545.pdf?sequence=1&isAllowed=y>

HUTAŘ, Jan, MIRANDA, Andrea, PAVLÁSKOVÁ, Eliška, et al. 2018. *Metodika logické ochrany digitálních dat* [online]. Praha: Knihovna AV ČR [cit. 2023-06-23]. Dostupné z:

<http://www.nusl.cz/ntk/nusl-371612>

CHAPMAN, Stephen et al., 2007. *Page Image Compression for Mass Digitization*. 2007.

Arlington (VA): Proc. IS&T Archiving Conference, s. 37-42. [cit. 2023-06-26] Dostupné také z:

https://web.archive.org/web/20140705020820/http://library.harvard.edu/sites/default/files/IST_PageImageCompression_preprint.pdf

IFLA STUDY GROUP ON THE FUNCTIONAL REQUIREMENTS FOR BIBLIOGRAPHIC RECORDS, 2009. *Functional Requirements for Bibliographic Records: final report* [online]. Haag: IFLA, September 1997, as amended and corrected through Feb 2009, v, 137 s. [cit. 2023-06-27]. Dostupné z: https://www.ifla.org/files/assets/cataloguing/frbr/frbr_2008.pdf.

ISO/IEC 11578:1996, 1996. *Information technology. Open Systems Interconnection. Remote Procedure Call (RPC)*. Geneva: ISO, 570 s.

ISO/IEC 15444-1:2004, 2004. *Information technology - JPEG 2000 image coding system: core coding system*. 2nd ed. Geneva: ISO, 194 s.

ISO 14721:2012, 2012. *Space data and information transfer systems - Open archival information system (OAIS) - Reference model*. 2nd ed. Geneva: ISO, 126 s.

ISO 16363:2012, 2012. *Space data and information transfer systems - Audit and certification of trustworthy digital repositories*. Geneva: ISO, 70 s.

ISO/IEC 9834-8:2014, 2014. *Information technology - Procedures for the operation of object identifier registration authorities. Part 8: Generation of universally unique identifiers (UUIDs) and their use in object identifiers*. 3rd ed. Geneva: ISO, 23 s.

JPEG 2000 profiles – examples from a range of institutions (footnotes on reverse), [2010]. Digital Preservation Coalition, [cit. 2023-06-23]. Dostupné z: <http://www.dpconline.org/docman/miscellaneous/events/529-jp2knov2010parametercomparisonchart/file>

KEJSER, Ulla Bøgvad, NIELSEN, Anders Bo a THIRIFAYS, Alex, 2011. Cost Aspects of Ingest and Normalization. In: *BORBINHA, José et al., ed. iPRES 2011: 8th International Conference on Preservation of Digital Objects* [online]. Singapore: National Library Board & Nanyang Technological University, s. 107-115 [cit. 2023-06-14]. Dostupné z: <https://services.phaidra.univie.ac.at/api/object/o:294222/download>

Koncepce rozvoje knihoven v České republice na léta 2021 - 2027 s výhledem do roku 2030, 2020 [online]. Praha: Národní knihovna ČR, 70 s [cit. 2023-06-26]. Dostupné z: <https://ipk.nkp.cz/docs/koncepce-rozvoje-2021-2027/koncepce-rozvoje-knihoven-2021-2027/view>

LAWRENCE, Gregory W. et al, 2000. *Risk management of digital information: a file format investigation*. Washington (DC): Council on Library and Information Resources, vi, 75 s. ISBN 18-873-3478-5. Dostupné také z: <https://www.clir.org/pubs/reports/pub93/pub93.pdf>

LEACH, Paul et al., 2005. *A Universally Unique Identifier (UUID) URN Namespace* [online]. The internet society [cit. 2023-06-23]. Dostupné z: <https://datatracker.ietf.org/doc/html/rfc4122>

LIBRARY OF CONGRESS, 2006. Office of Strategic Initiatives. *JPEG 2000 Profile for the National Digital Newspaper Program* [online]. Prepared by: Robert Buckley a Roger Sam. Washington (DC): The Library of Congress, 24 s. [cit. 2023-06-23]. Dostupné z: https://www.loc.gov/ndnp/guidelines/docs/NDNP_JP2HistNewsProfile.pdf

LIBRARY OF CONGRESS, 2017. *Sustainability of Digital Formats: Planning for Library of Congress Collections*. [online]. Washington (DC): The Library of Congress, Last updated 3rd October 2017 [cit. 2023-06-21]. Dostupné z: <http://www.digitalpreservation.gov/formats/index.shtml>

LIBRARY OF CONGRESS, 2022-2023. *Recommended Formats Statement 2022-2023* [online]. Washington (DC): The Library of Congress, [cit. 2023-21-06]. Dostupné z: <https://www.loc.gov/preservation/resources/rfs/RFS%202022-2023.pdf>

Media types, 2023 [online]. Los Angeles (CA): Internet Assigned Numbers Authority, [cit. 2023-06-23]. Dostupné z: <http://www.iana.org/assignments/media-types/media-types.xhtml>

METS: An Overview & Tutorial. Metadata Encoding and Transmission Standard (METS), 2022 [online]. Washington (DC): The Library of Congress, Last updated 28th March 2022 [cit. 2023-06-13]. Dostupné z: <http://www.loc.gov/standards/mets/METSOverview.v2.html>

Metadata Object Description Schema. MODS: Uses and Features, 2022 [online]. Washington (DC): The Library of Congress, Last updated 4th February 2022 [cit. 2023-06-10]. Dostupné z: <https://www.loc.gov/standards/mods/mods-overview.html>

MILLER, Steven J., 2011. *Metadata for digital collections: a how to do it manual*. London: Facet, 343 s. ISBN 978-1-85604-771-5

MORAVSKÁ ZEMSKÁ KNIHOVNA, [2015]. *LTP-portál.cz. Web o digitální archivaci* [online]. Brno: Moravská zemská knihovna [cit. 2023-06-23]. Dostupné z: <https://ltp-portal.mzk.cz/>

Národní centrum ISSN [online]. Praha: Centrum PID, [cit. 2023-06-19]. Dostupné z: <https://identifikatory.cz/cs/sluzby/nc-issn/>

NÁRODNÍ KNIHOVNA ČR, 2016. *Průručka uživatele systému ISMN*. 2. české vydání. Z anglického originálu přeložil a upravil Antonín JEŘÁBEK [online]. Praha: Národní agentura ISMN v ČR, Národní knihovna ČR, původní vydání 2008, revize březen 2016. [cit. 2023-06-20]. Dostupné z: <https://www.nkp.cz/soubory/ostatni/prirucka-ismn.pdf>

NÁRODNÍ KNIHOVNA ČR, 2019. *Číslo ČNB v SK ČR*. [online] Praha: Národní knihovna ČR, poslední aktualizace 17.10.2019 [cit. 2023-06-26]. Dostupné z:

<https://www.caslin.cz/caslin/spoluprace/sluzby-souborneho-katalogu-cr/cislo-cnb-v-sk-cr>

NÁRODNÍ KNIHOVNA ČR, 2021. *Standardy pro obrazová data. Národní digitální knihovna* [online]. Praha: Národní knihovna ČR, poslední aktualizace 27.10.2021 [cit. 2023-06-23].

Dostupné z: <https://standardy.ndk.cz/ndk/standardy-digitalizace/standardy-pro-obrazova-data>

NÁRODNÍ KNIHOVNA ČR, 2023a. *Mezinárodní registrační systémy* [online.]. Praha: Národní knihovna ČR, poslední aktualizace 25.01.2023 [cit. 2023-06-26]. Dostupné z:

<https://www.nkp.cz/sluzby/sluzby-pro/isbn-ismn-issn>

NÁRODNÍ KNIHOVNA ČR, 2023b. *Standardy pro metadata. Národní digitální knihovna* [online]. Praha: Národní knihovna ČR, poslední aktualizace 15.06.2023 [cit. 2023-06-26].

Dostupné z: <http://www.ndk.cz/standardy-digitalizace/metadata>

NEISS, Bengt. *File format for still images* [elektronická pošta]. Message to: natalie.ostrakova@nkp.cz. 9. listopadu 2017 [cit. 2023-06-10]. Osobní komunikace.

NISO FRAMEWORK WORKING GROUP, 2007. *A framework of guidance for building good digital collections: a NISO recommended practice* [online]. 3rd ed. Baltimore (MD): National Information Standards Organization (NISO), iii, 95 s. [cit. 2023-06-23]. ISBN 978-1- 880124-74-1. Dostupné z: <https://www.niso.org/sites/default/files/2017-08/framework3.pdf>

NISO Metadata for Images in XML Schema. Technical Metadata for Digital Still Image Standard, 2021. Washington (DC): Library of Congress, Last updated 21th October 2021 [cit. 2023-06-20]. Dostupné z: <https://www.loc.gov/standards/mix/>

OSTRÁKOVÁ, Natalie, KOČIŠOVÁ, Pavlína a BEŇAČKOVÁ, Miroslava, 2019. Vývoj standardu PREMIS a možnosti jeho dalšího využití ve standardech NDK. *ProInFlow* [online] **11**(2), 72-85 [cit. 2023-06-20]. DOI: <https://doi.org/10.5817/ProIn2019-2-6>. Dostupné z:

<https://journals.phil.muni.cz/proinflow/article/view/2019-2-6>

OSTRÁKOVÁ, Natalie a KOPSKÝ, Vojtěch, 2020. Posuzování souborových formátů z hlediska dlouhodobého uchování a návrh metodiky pro Národní knihovnu České republiky. *Knihovna: knihovnická revue*. **31**(2), 83–105. ISSN 1801-3252. [cit. 2023-06-20] Dostupné z:

<https://knihovnav revue.nkp.cz/archiv/2020-2/recenzovane-prispevky/posuzovani-souborovych-formatu-z-hlediska-dlouhodobeho-uchovavani-a-navrh-metodiky-pro-narodni-knihovnu-ceske-republiky>

POMERANTZ, Jeffrey, 2015. *Metadata*. Cambridge, Massachusetts: MIT Press, 252 s. ISBN 978-0-262-52851-1

PREMIS EDITORIAL COMMITTEE, 2015. *PREMIS Data Dictionary for Preservation Metadata* [online]. Version 3.0. Washington (DC): Library of Congress, rev. Nov 2015, viii, 273 s. [cit. 2023-06-23]. Dostupné z: <http://www.loc.gov/standards/premis/v3/premis-3-0-final.pdf>

RIMKUS, Kyle, Thomas PADILLA, Tracy POPP a MARTIN, Greer, 2014. Digital Preservation File Format Policies of ARL Member Libraries: An Analysis. *D-Lib Magazine* [online]. **20**(3/4) [cit. 2023-06-23]. Dostupné z: <http://www.dlib.org/dlib/march14/rimkus/03rimkus.html>

RŮŽIČKA, Michal, MIRANDA, Andrea, HEJTMÁNEK, Lukáš et al. 2019. *Metodika bitové ochrany digitálních dat* [online]. Praha: Knihovna AV ČR [cit. 2023-06-23]. Dostupné z: <http://www.nusl.cz/ntk/nusl-393240>

SMITH, Neil, 2006. Digitising Documents for Public Access. In: MACDONALD, Lindsay, ed. *Digital heritage: applying digital imaging to cultural heritage*. Amsterdam: Elsevier, s. 3-32. ISBN 0-75-066183-6.

THE ASSOCIATION FOR LIBRARY COLLECTIONS AND TECHNICAL SERVICES, 2013. Preservation & Reformatting Section. *Minimum Digitization Capture Recommendations*. In: Association for Library Collections & Technical Services (ALCTS) [online]. Chicago: ALA, June 2013 [cit. 2023-06-21]. Dostupné z: <https://www.ala.org/alcts/resources/preserv/minimum-digitization-capture-recommendations>

THE DIGITAL LIBRARY FEDERATION BENCHMARK WORKING GROUP, 2002. *Benchmark for Faithful Digital Reproductions of Monographs and Serials* [online]. Version 1. Washington (DC): Digital Library Federation, 6 s. [cit. 2023-06-20]. Dostupné z: <http://old.diglib.org/standards/bmarkfin.pdf>.

THE NATIONAL ARCHIVES, 2017. *General hints and tips for digitisation for business use* [online]. [cit. 2023-06-15]. Dostupné z: <https://www.nationalarchives.gov.uk/documents/information-management/hints-tips-digitisation-for-business-use.pdf>

TIFF, revision 6.0, 3rd June 1992, Mountain View (CA): Adobe Systems Inc. [cit. 2023-06-24]. Dostupné z: <https://developer.adobe.com/content/dam/udp/en/open/standards/tiff/TIFF6.pdf>

Unified digital format registry, 2012-2016 [online]. University of California [cit. 2023-06-23]. Dostupné z: <https://web.archive.org/web/20220511203412/http://udfr.org/>

VAN DER KNIJFF, Johan, 2011. JPEG 2000 for Long-term Preservation: JP2 as a Preservation Format. *D-Lib Magazine* [online]. **17**(5/6) [cit. 2023-06-13]. Dostupné z: <http://www.dlib.org/dlib/may11/vanderknijff/05vanderknijff.html>

VAŠEK, Zdeněk, CUBR, Ladislav, ŘEHÁNEK, Martin, 2018. *Metodika pro přidělování a správu životního cyklu unikátních perzistentních identifikátorů digitálních dokumentů podle*

standardu URN:NBN, Verze 2.0. In: Národní digitální knihovna [online]. Praha: Národní knihovna ČR [cit. 2023-06-23]. Dostupné z:
https://standardy.ndk.cz/ndk/archivace/Certifik_metodika_urnbn_2018.pdf

VRTĚLOVÁ, Lucie, 2017. *Analýza nastavení formátu JPEG 2000*. Brno. Bakalářská práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Vedoucí práce David BAŘINA.

VYCHODIL, Bedřich, 2010. JPEG2000 - Aneb nemyslete si, že vás mine!. *Knihovna* [online]. **21**(2), s. 53-68 [cit. 2023-06-23]. Dostupné z: <http://oldknihovna.nkp.cz/knihovna102/10253.htm>

ZENG, Marcia Lei a QIN, Jian, 2016. *Metadata*. 2nd edition. Chicago: Neal-Schuman, an imprint of the American Library Association. xxvii, 555 stran. ISBN 978-1-55570-965-5.

Zřizovací listina Národní knihovny České republiky, 2011. Praha: Ministerstvo kultury ČR, vydaná 30. listopadu 2011 [cit. 2023-06-23]. Dostupné z:
<https://text.nkp.cz/soubory/ostatni/zrizovaci-listina-nk.pdf>