

## ProArc

### – open source řešení pro produkci a archivaci digitálních dokumentů

**Martina NEZBEDOVÁ**

Knihovna AV ČR, v. v. i., Praha

[nezbedova@knav.cz](mailto:nezbedova@knav.cz)

INFORUM 2015: 21. ročník konference o profesionálních informačních zdrojích

Praha, 26. – 27. 5. 2015

#### **Abstrakt**

*Produkční a archivační systém ProArc je volně dostupný nástroj na výrobu a editaci popisných, technických a administrativních metadat k digitalizovaným i born digital dokumentům. Systém umožňuje vytvoření importních dat pro systém Kramerius. Jedním z výstupů jsou PSP balíčky vytvořené podle NDK standardů Národní knihovny. Na vývoji systému ProArc spolupracuje Knihovna AV ČR, v. v. i. s firmou INCAD, která zajišťuje analytické a programátorské práce. Od roku 2013, kdy byl tento systém plně nasazen v Digitalizačním centru KNAV, v něm bylo zpracováno více než 400 tisíc stran různých typů dokumentů. Systém začínají v současné době využívat i další instituce. Ve vývoji je archivační část, která umožní propojení na systém Archivematica. Příspěvek představí ProArc z pohledu zpracovatele.*

*Systém ProArc je jedním z výstupů projektu Česká digitální knihovna. Cílem celého projektu je vytvoření jednotného rozhraní nad digitálními knihovnami v ČR pro koncové uživatele a poskytování dat pro mezinárodní projekty. Projekt je financován z programu NAKI Ministerstva kultury ČR.*

Produkční a archivační systém ProArc je volně dostupný nástroj na výrobu a editaci popisných, technických a administrativních metadat k digitalizovaným i born digital dokumentům.

Na vývoji systému ProArc úzce spolupracuje Knihovna AV ČR, v. v. i., hlavní koordinátor vývoje, s firmou INCAD, která zajišťuje analytické a programátorské práce. Systém ProArc je jedním z výstupů projektu Česká digitální knihovna. Cílem celého projektu je vytvoření jednotného rozhraní nad digitálními knihovnami v ČR pro koncové uživatele a poskytování dat pro mezinárodní projekty. Projekt je financován z programu NAKI Ministerstva kultury ČR.

Od roku 2013, kdy byl tento systém plně nasazen v Digitalizačním centru KNAV, v něm bylo zpracováno více než 400 tisíc stran různých typů dokumentů. Systém začínají v současné době využívat i další instituce např. Studijní a vědecká knihovna v Hradci Králové.

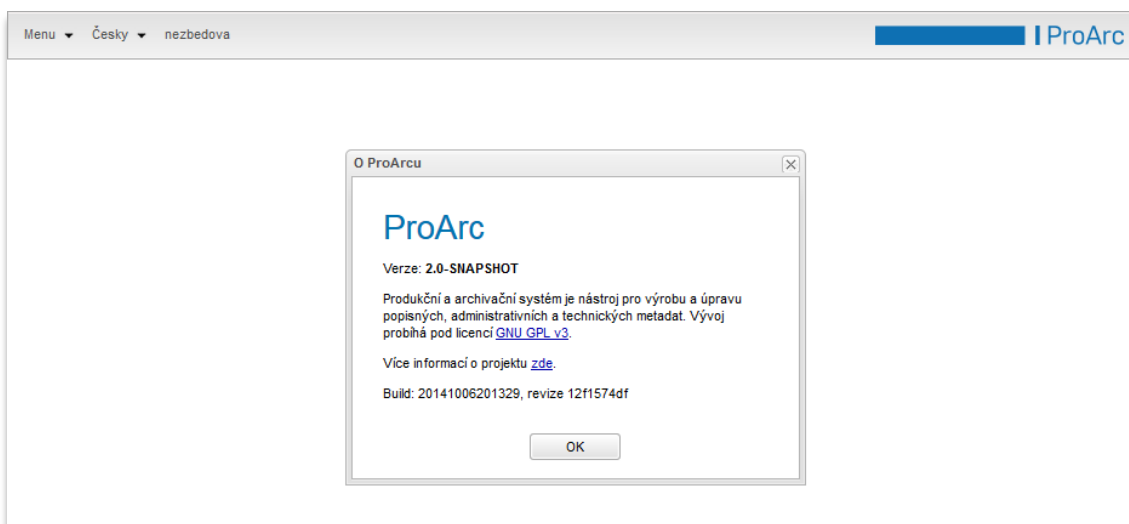
Produkční část systému ProArc je nástroj na výrobu a úpravu popisných, administrativních, technických metadat a jejich editaci. V tomto systému je možné nejen zakládat zcela nové objekty, ale i používat metadata z externích systémů (např. Aleph, Registr digitalizace). Systém ProArc podporuje dávkové a hromadné úpravy a také částečně automatizovanou produkci digitálních dokumentů. Systém ProArc automaticky generuje UUID pro jednotlivé objekty a umožňuje vyhledávání v metadatech. Rovněž je možné následné rozšíření o nové datové modely.

Systém ProArc podporuje standardy Národní knihovny ČR pro digitalizaci (plná podpora MODS, administrativních a technických metadat) a přidělování URN:NBN. Dodržování těchto standardů umožní spolupracovat s Národní knihovnou při sdílení dat. Plnění těchto standardů je požadováno v dotačním programu VISK 7 a v krajských digitalizačních projektech. Typy těchto dokumentů jsou zakládány ve formulářích označených jako NDK (Národní digitální knihovna). Zároveň je pro ně připraven export v podobě PSP balíčku dle standardů pro NDK. Systém ProArc je plně kompatibilní se systémem Kramerius.

Systém ProArc je určen i pro zpracování born digital dokumentů. Tyto dokumenty je po zpracování možné exportovat jak pro systém Kramerius, tak i jako metadata do společné bibliografické databáze akademií věd Visegradské čtyřky CEJSH.

Systém ProArc je open source, který je vystavěn na volně dostupných řešeních Fedora Commons repozitory, Java a PostgreSQL. Systém ProArc je webová aplikace, která využívá lokální server. Pro generování grafického formátu JPEG2000 využívá v rámci standardů NDK program Kakadu, ale podporuje i užití jiných programů. Pro potřeby OCR je využit komerční ABBYY Recognition Server, který umožňuje generovat formát ALTO XML. Technickou podporu zajišťuje firma INCAD.

Projektová dokumentace, informace o aktuálním stavu vývoje a řešených issues jsou spolu s instalačním balíčkem verze ProArc 2.0 NDK umístěny na adrese <https://code.google.com/p/archivacni-system/wiki/ProArc>.



V úvodní obrazovce se v horní liště nachází rozbalovací roletka Menu, která skrývá základní složky Import a Úložiště s jednotlivými pořadači a složku Zařízení obsahující formuláře odpovídající NDK standardům pro výrobu technických metadat (informace o skenerech a parametrech skenování).

Tvorba metadat začíná vytvořením databáze možností používaných technických metadat vyplněním různých typů používaných skenerů a aktuálních parametrů skenování. Tyto údaje lze doplňovat nebo měnit podle potřeby digitalizace. Nastavení je společné pro všechny pracovníky na jedné instalaci ProArcu.

Dalším krokem je založení objektu, ke kterému budou následně připojeny soubory s hotovými popisnými metadaty. Jednotlivé objekty, které lze založit jsou:

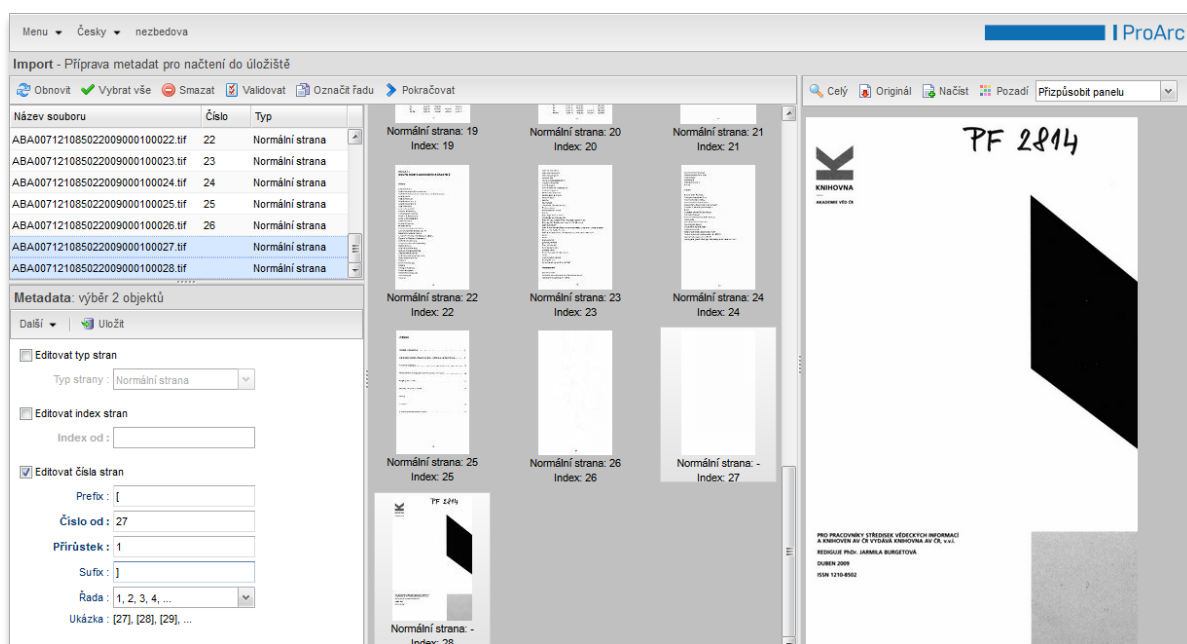
NDK Periodikum	K4 Periodikum
NDK Ročník	K4 Ročník
NDK Číslo	K4 Výtisk
NDK Příloha periodika	K4 Monografie
NDK Článek	K4 Monografie – volná část
NDK Obrázek/Mapa – vnitřní část	Strana
NDK Vicedílná monografie	E-článek
NDK Svazek monografie	
NDK Příloha monografie	
NDK Kapitola	
NDK Kartografický dokument	
NDK Hudebnina	

Formuláře s označením NDK obsahují položky pro vyplnění odpovídající NDK standardům a nad poli s povinností vyplnění (Mandatory) je prováděna validace vyplnění těchto polí. Tyto objekty jsou exportovány jako PSP balíčky, ale je u nich možný i export pro K4 v xml. Formuláře s označením K4 obsahují vybraná pole MODS (lze si přepnout i na plný MODS). I zde probíhá validace povinných polí. Exportem je pouze xml.

Založit lze jak prázdný formulář pro následné ruční vyplnění, tak i formulář s poli vyplněnými ve vybrané databázi (např. Registr digitalizace, Aleph). Pro výběr slouží jeden z identifikátorů (např. čárový kód, ISSN, číslo ČNB).

Dalším krokem je výběr adresáře s digitalizovanými daty pro přípravu metadat z importního adresáře na serveru, ke kterému se připojí technická metadata a vybrané soubory jsou načteny do sestavy určené pro tvorbu popisných metadat.

Obrazovka s názvem Tvorba metadat je velice variabilní a lze si ji přizpůsobit podle typu zpracovávaných dokumentů nebo pracovních zvyklostí. Popisná metadata lze tvořit po jednotlivých stranách nebo výběrem souvislé řady stran nebo výběrem různých jednotlivých stran. Popisná metadata je možné přiřazovat po řádcích obsahujících název souboru, v malých náhledech, ve velkých náhledech nebo jejich kombinací.



Pro výrobu popisných metadat jsou připraveny typy stran odpovídající standardům NDK. Pro usnadnění tohoto popisu je automaticky přednastaven typ strany – normální strana. K doplnění ostatních typů stran slouží roletka s předepsanými typy stran. Pro doplňování čísel jednotlivých stran jsou připravené číselné řady arabských i římských číslic a písmenná řada, doplněné možností prefixu a sufixu. Před připojením souborů k nadřazenému objektu prochází všechny soubory validací na vyplnění čísla strany. I po připojení k nadřazenému objektu lze typy i čísla stran editovat, prohlížet si náhledy jednotlivých stran nebo tyto strany přesouvat.

Užitečnou funkcí je možnost Označit řadu, tj. řadu po sobě jdoucích stran s vyplněnými popisnými metadaty a tyto následně aplikovat na další řadu objektů. Při popisu výtisků, které mají vždy stejné typy stran a jsou stejně číslované, to velice zjednoduší a zrychlí jejich popis, protože jednou vytvořená popisná metadata se jen zkopírují na další výtisky.

Důležitou součástí systému ProArc je správa digitálních objektů, která obsahuje všechny založené typy objektů, ve kterých lze vyhledávat. U jednotlivých objektů jsou informace o datu založení, vlastníku záznamů a případném exportu. V této databázi lze vyhledávat podle různých kritérií. Rovněž je zde možná editace ať už jednotlivých objektů nebo hromadná přepnutím do obrazovky vazby.

Ze systému ProArc lze po přidělení URN:NBN exportovat NDK PSP balíčky plně odpovídající standardům NDK. Dalšími možnostmi jsou export pro K4, export skenů, export původních skenů a pro born digital dokumenty export pro CEJSH.

The screenshot shows the ProArc web interface. At the top, there is a navigation bar with 'Menu', 'Česky', and 'nezbedova'. Below it, the main heading is 'Správa digitálních objektů'. There are several radio buttons for filtering: 'Naposledy vytvořené' (selected), 'Naposledy editované', 'Dotaz', and 'Rozšířený dotaz'. A dropdown menu shows 'Model: NDK Periodikum'. Below this is a search bar labeled 'Hledat'. A toolbar contains icons for 'Filtr', 'Obnovit', 'Metadata', 'Komentář', 'Vazby', 'ATM', 'FOXML', 'Exporty', 'Smazat', and 'URN:NBN'. The main table lists objects with columns: 'Popis', 'Model', 'PID', 'Vytvořen', 'Změněn', 'Vlastník', and 'Exp'. The selected object is 'Zprávy Slezského ústavu ČSAV v Opavě'. A detailed view of this object is shown below, with a tree structure of folders and files. A context menu is open over the 'Exporty' icon, showing options: 'Export pro Kramerius 4', 'NDK PSP Export', 'CEJSH Export', 'Export skenů', and 'Export původních skenů'.

Popis	Model	PID	Vytvořen	Změněn	Vlastník	Exp
Právník: časopis věnovaný vědě právní i státní, jež vydává ...	NDK Periodikum	uuid:10917b4e-be20-4931-9ba6-0abd19960d89	08.04.2015 14:14	08.04.2015 16:31	melichova	Ano
Soudobé dějiny	NDK Periodikum	uuid:43792fe2-7390-41fa-97d5-22b902e54cf4	20.03.2015 09:35	20.03.2015 09:48	melichova	Ne
Česká literatura: časopis pro literární vědu	NDK Periodikum	uuid:646fb512-92d9-4d74-843c-d799b15535a6	02.03.2015 11:08	02.03.2015 11:10	melichova	Ne
Archeologické rozhledy: informační orgán archeologických a...	NDK Periodikum	uuid:365b1357-50ff-4c23-aeb8-e22090693324	02.03.2015 10:38	02.03.2015 11:03	melichova	Ne
Preslia: věstník Československé botanické společnosti	NDK Periodikum	uuid:42896897-0d7f-11e3-993b-005056ae0003	28.01.2015 11:42	23.02.2015 23:36	proarc	Ano
Zprávy Slezského ústavu ČSAV v Opavě	NDK Periodikum	uuid:02da2e84-f187-4c92-ae7f-deba6e912882	08.07.2014 14:29	06.03.2015 10:51	nezbedova	Ano

Popis	Model	PID	Vytvořen	Změněn	Vlastník	Export
Zprávy Slezského ústavu ČSAV v Opavě	NDK Periodikum	uuid:02da2e84-f187-4c92-ae7f-deba6e912882	08.07.2014 14:29	06.03.2015 10:51	nezbedova	Ano
1959	NDK Ročník	uuid:77fb6271-1497-4c92-ae7f-deba6e912882	07.2014 14:44	19.09.2014 13:10	nezbedova	Ano
1960	NDK Ročník	uuid:5e6d01a3-c60d-4c92-ae7f-deba6e912882	07.2014 14:44	19.09.2014 13:10	nezbedova	Ano
1961	NDK Ročník	uuid:9194245f-7d64-4c92-ae7f-deba6e912882	07.2014 14:44	19.09.2014 13:10	nezbedova	Ano
116	NDK Číslo	uuid:61d4677d-c6df-410f-a577-5bae9e115336	08.07.2014 15:16	19.09.2014 13:10	nezbedova	Ano
117	NDK Číslo	uuid:a352ad31-8486-4720-ab60-15e6f3ebc80a	08.07.2014 15:20	19.09.2014 13:10	nezbedova	Ano
[1], Titulní strana	Strana	uuid:d2451087-ee7b-4d84-b12f-bb91c6d5c2c	09.07.2014 09:23	19.09.2014 13:10	nezbedova	Ano
2	Strana	uuid:39eba823-d0db-4ec4-aaef-b003d7030396	09.07.2014 09:23	19.09.2014 13:10	nezbedova	Ano
3	Strana	uuid:58bcae57-c3ec-48b6-a2a6-a2c0d57213f7	09.07.2014 09:23	19.09.2014 13:10	nezbedova	Ano
4	Strana	uuid:f87fd2c1-a295-47af-bbc6-43eb393932a7	09.07.2014 09:23	19.09.2014 13:10	nezbedova	Ano

Prohlížení a editace již hotových objektů umožňuje obrazovka Vazby, ve které je také možné přesouvat jednotlivé soubory i celé řady souborů jak v jednotlivých objektech, tak i mezi objekty navzájem. Tím je zajištěna náprava možných chyb.

V systému ProArc je také možné zpracovávat born digital dokumenty. Po založení e-článku se připojí již hotová metadata z připojené databáze Knihovna AV ČR Analytika a v dalším kroku se přidá plný text v pdf formátu. Exportovat lze jak plný text s přidanými metadaty do Krameria, tak jen metadata formou Export CEJSH.

Produkční část systému ProArc je velice vhodným a uživatelsky příjemným nástrojem pro výrobu metadat a jejich následného využití jak v systému Kramerius, tak i ke vzájemnému sdílení s ostatními institucemi, které dodržují standardy NDK. Systém se stále vyvíjí a vylepšuje pro širší využití v praxi (např. popis starých tisků).

Informace o systému ProArc jsou na <https://code.google.com/p/archivacni-system/>

K dispozici je i diskusní skupina na <http://groups.google.com/group/proarc-users>