

**ZDOKONALENÍ VIRTUÁLNÍHO BADATELSKÉHO
PROSTŘEDÍ MANUSCRIPTORIA
ROZŠÍŘENÍ SLUŽEB MANUSCRIPTORIA ROZŠÍŘENÍM
NORMALIZOVANÉ INFORMACE O TYPU OBSAHU
OBRAZŮ**

Zpráva

AiP Beroun, autor: Ing. Tomáš Psohlavec

Obsah

1	Úvod o dokumentu	3
1.1	Účel	3
1.2	Předpokládaný čtenář	3
1.3	Termíny a konvence	3
1.4	Reference	3
2	Úvod	4
3	Analytická fáze	4
4	Implementační fáze.....	5
4.1	Cíle implementace	5
4.2	Příprava implementace.....	6
4.3	Workflow zpracování metadat	6
4.4	Implementace do koncového rozhraní.....	7
4.5	Pilotní řešení	8
5	Závěr	8
6	Příloha 1: Seznam dokumentů cíleně vybraných do testovacího vzorku.....	9

1 Úvod o dokumentu

AiP Beroun uzavřela s Národní knihovnou České republiky Smlouvu o spolupráci ve výzkumu a vývoji: Zdokonalení virtuálního badatelského prostředí Manuscriptoria - rozšíření služeb Manuscriptoria rozšířením normalizované informace o typu obsahu obrazů.

1.1 Účel

Tento dokument tvoří zprávu k rozšíření aktuálních služeb Manuscriptoria začleněním normalizované informace o typu obsahu obrazů (text, iluminace, hudební notace, bordura, tabulka, diagram) a doplněním hledání dle této položky dle poptávky zadavatele č.j. 2295/KGR/11 ze dne 23.8.2011.

1.2 Předpokládaný čtenář

Tento dokument je určen především pro Zadavatele (NK ČR) a pro Řešitele úkolu (AiP Beroun) jako popis pilotního řešení. Dále je tento dokument určen všem, kteří se podílejí na rozvoji projektu Manuscriptorium jako uživatelé.

1.3 Termíny a konvence

Termíny a konvence použité v tomto dokumentu, pokud zde nejsou přímo vysvětleny, jsou popsány a definovány v dokumentu [2].

- MnS Manuscriptorium
- IIR Rozpoznání informací v obraze (Image Information Recognition)
- IIR Aplikace Aplikace pro rozpoznávání informací v obraze

1.4 Reference

V dokumentu se odkazujeme na následující zdroje:

- [1] Guidelines for Electronic Text Encoding and Interchange, Representation of Primary Sources
<http://www.tei-c.org/release/doc/tei-p5-doc/en/html/PH.html#PHFAX>
- [2] Manuscriptorium v.2.0 – analýza systému, AiP Beroun 2004

2 Úvod

NK ČR (Zadavatel) vyvíjí se svým subdodavatelem desktopovou aplikaci pro rozpoznávání informací v obrazech (IIR Aplikace). Zadavatel IIR Aplikaci poskytl Řešiteli. Řešitel zanalyzoval možnosti využití výstupů této aplikace v Manuscriptoriu a připravil pilotní řešení, na jehož základě lze posuzovat účelnost rozšíření služeb, které Manuscriptorium poskytuje koncovým uživatelům, případně Content Providerům.

3 Analytická fáze

V souvislosti se zamýšlenou implementací IIR do Manuscriptoria bylo nejprve nutno zanalyzovat možnosti zachycení požadovaných informací v metadatech. Vzhledem k tomu, že MnS pracuje s formátem TEI P5, byly přirozeně prozkoumány možnosti zápisu informací v TEI P5 přímo v elementu `<facsimile>`: využity byly podelementy `<zone>` způsobem, který popisuje následující příklad.



```
<surface xml:id="ID0011R">
  <desc>
    <label>11r</label>
  </desc>
  <zone xml:id="ID0001v__graphic" ulx="0" uly="0" lrx="1082" lry="1555">
    <graphic url="I/MVCHK_HR_13_II_A_13__2FMI/N1/HR_13_II_A_13__2FMIN10011R.JPG"/>
  </zone>
  <zone xml:id="ID0011R__pic_1" ulx="63" uly="51" lrx="681" lry="681">
    <desc>Libovolný volitelný popis (červená oblast).</desc>
  </zone>
  <zone xml:id="ID0011R__pic_2" points="21,1044 822,1044 822,783 1059,783 1059,1536
21,1536">
    <desc> Libovolný volitelný popis (zelená oblast).</desc>
  </zone>
</surface>
```

Postup při generaci atributu `xml:id` pro element `<zone>` byl navržen tak, že

```
@xml:id = ../surface/@xml:id+"__"+pic+"_"+N
```

Kde `pic` je zkratka typu rozezaného grafického elementu a `N` je pořadové číslo výskytu daného typu grafického elementu v daném obraze.

Aby bylo možné vztáhnout umístění grafického elementu relativně k danému obrazu i zrcadlu stránky, je vždy součástí záznamu o daném `<surface>` i informace o rozměrech obrazu a o zrcadle, například:

```
<zone xml:id="ID0001v__graphic" ulx="0" uly="0" lrx="995" lry="1120"/>
<zone xml:id="ID0001v__mirror" ulx="164" uly="32" lrx="849" lry="1031"/>
```

Tak je zajištěno, že mění-li se rozlišení obrazu (například použije-li se jiná uživatelská kvalita), budou informace zanesené v metadatech i nadále platné.

Poznámka: povšimněte si v příkladu uvedeného zapouzdření elementu `<graphic>` s odkazem na odpovídající datový soubor do elementu `<zone>` - tento přístup lze aplikovat i u ostatních `<zone>` elementů v případě existence doplňkových obrazových dat (například výřez iniciály ve vyšším rozlišení dokumentů). Bližší informace o aparátu k zachycení nově generovaných metadat uvádí [1].

4 Implementační fáze

4.1 Cíle implementace

Cílem řešení je obohatit uživatelské rozhraní Manuscriptoria o funkce založené na existenci nových metadat generovaných IIR Aplikací. Za tím účelem AIP vytvořila pilotní řešení, umožňující Zadavateli zhodnotit kvalitu generovaných informací a zvážit možnosti jejího nasazení do ostrého provozu MnS.

V součinnosti s odpovědnými pracovníky Zadavatele byly popsány tyto dvě nové funkce:

1. Vyhledávání dokumentů dle obsahu faksimile:
vyhledávací formulář bude obohacen o checkboxy, které omezí hledání například jen na dokumenty s hudební notací, či dokumenty s iniciálami.
2. Rychlý přístup na relevantní stránky:
z detailu záznamu bude možno otevřít faksimile na konkrétní stránce, například na stránce s iluminací.

Seznam rozeznávaných grafických elementů, které budou využity k dosažení výše uvedeného:

- Ilustrace
- Iniciála
- Notace
- Marginélie
- Tabulka

4.2 Příprava implementace

Součástí přípravných fází řešení byla formulace požadavků na IIR Aplikaci dodávané Zadavatelem, které se týkaly požadovaných vstupů a výstupů.

Cílem úkolu bylo umožnit efektivní hromadné zpracování obrazů. Proto bylo požadováno, aby IIR Aplikace mohla pracovat dvěma způsoby, které se liší v závislosti na použitém druhu vstupů:

1. Neexistují metadata, vstupem je adresářová struktura s obrazy:
IIR Aplikace předepsaným způsobem generuje nová metadata ve formátu TEI P5 a ukládá jako samostatné soubory (jeden pro každý adresář s obrazy)
2. Existují metadata, vstupem je adresářová struktura s XML soubory:
IIR Aplikace zpracovává obrazové soubory z umístění odkazovaného v metadatach a doplňuje do existujících metadat nové informace.

Kromě výstupu ve formátu TEI P5 XML bylo též požadováno, aby aplikace generovala log soubory s informací o míře jistoty, s jakou byly jednotlivé výskyty informací ve stránce rozeznány. Tento log soubor může být mimo jiné využit při analýze úspěšnosti.

4.3 Workflow zpracování metadat

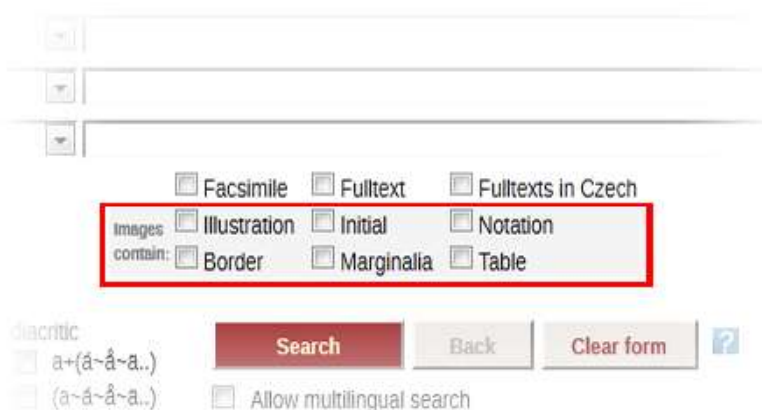
IIR Aplikace byla začleněna do workflow pilotního prostředí MnS. Vstupem ke zpracování jsou existující XML soubory. Data ke zpracování musí být IIR Aplikaci

dostupná lokálně, čili je nutno centralizovat zpracovávané obrazové soubory i odpovídající metadata (pracovní kopie) a ta aktualizovat podle aktuálního místa uložení dat. Toto řešení obsahuje výrazný podíl ruční práce při přípravě a zpracování dat a metadat, což je při realizaci pilotního řešení akceptovatelné. Pro případné ostré nasazení bude nutné jednotlivé procesy automatizovat, aby bylo možno bezpečně provádět dávkové hromadné zpracování.

Začlenění výstupů IIR Aplikace do workflow zpracování metadat se samozřejmě neobešlo bez změny rutin pro generování FRT/FDM souborů – datových souborů katalogu MnS tak, aby bylo možno vyhledávat dle typu obrazu.

4.4 Implementace do koncového rozhraní

V souladu s cílem pilotního řešení byly provedeny změny v pokročilém formuláři, jež umožňují žádoucím způsobem omezit výsledky vyhledávání.



The screenshot shows a search interface with several dropdown menus at the top. Below them is a section for filtering search results, with a red box highlighting the 'images contain:' section. The filters are as follows:

- Facsimile
- Fulltext
- Fulltexts in Czech
- images contain:
 - Illustration
 - Initial
 - Notation
 - Border
 - Marginalia
 - Table

Below the filters are buttons for 'Search', 'Back', and 'Clear form', along with a help icon. There are also checkboxes for 'diacritic' and 'Allow multilingual search'.

Dále bylo upraveno zobrazení detailu záznamu tak, že při výpisu informací se generují seznamy stránek s odkazy na stránky obsahující dané typy informací.

Image Contents Summary
Illustration: FC FS 2 73 95 125 172 175 186 191 193 194 197 198 200 230 254 265 270 271 273 274 294 296 306 312 314 315 320 324 341 342 347 348 354 364 367 380 389 372 380 384 388 389 bis 397 398 405 406 422 431 432 441 454 456 457 459 460 BC
Marginalia: 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 22 33 41 42 44 46 52 54 58 59 61 83 103 109 113 135 157 159 186 187 188 189 190 196 203 214 232 248 251 257 258 263 265 312 329 331 332 334 339 346 348 351 359 365 371 372 374 377 385 387 396 408 410 412 434 444 449 450
Border: FC 113 131 169 202 244 bis 264 269 312 317 318 319 348 364 370 389 406 432 442 458 460 SP
Notation: 453
Initial: FS 5 10 75 170 171 172 173 174 176 191 192 194 195 197 198 200 221 222 227 259 260 262 266 269 271 272 273 274 282 300 314 340 342 369 371 372 389 390 391 402 403 404 406 408 409 415 421 422 427 432 434 452 453 460 SP
Table: 73 75 77 79 81 83 85 87 89 91 93 95 97 99 101 103 105 107 109 111 113 115 117 119 121 123 125 127 129 131 133 135 139 141 143 145 147 149 151 153 409 432

4.5 Pilotní řešení

Pro pilotní řešení byla zvolena data dokumentů české provenience. Z části se jedná o náhodný výběr a z části o sadu zvláštních dokumentů různých typů a vlastností, jež AIP vybrala nad rámec smlouvy a které považuje za důležité pro testování.

Sadu dokumentů je možno změnit dle požadavků Zadavatele v závislosti na aktuálních potřebách a výsledcích testování.

V době kompletace této zprávy se jednalo o 88 dokumentů náhodně vybraných a 25 dokumentů cíleně vybraných [viz dokumentace].

Pilotní řešení je volně dostupné na adrese

<http://www.manuscriptorium.com/apps/pilot/iir/>

a aby bylo možno testovat jej v mezinárodním měřítku, je k dispozici i anglická verze

<http://www.manuscriptorium.com/apps/pilot/iir/en>

Některé pro testování nedůležité koncovo-uživatelské funkce byly v pilotním řešení zakázány.

5 Závěr

Pilotní řešení je nyní plně k dispozici na výše uvedené adrese. Pracuje s aktuálními výsledky poskytovanými IIR Aplikací.

Manuscriptorium aktuálně zpřístupňuje cca 5 000 000 obrazů. Měření rychlosti při zpracování vzorků dat ukazuje, že lze IIR Aplikaci využívat i v takto masovém měřítku. Pokud tedy Zadavatel dojde po testování k rozhodnutí, že je účelné IIR Aplikaci a jí generovaná metadata využívat v MnS, lze dále řešit její nasazení v ostrém provozu, mimo jiné:

- řešit nasazení do workflow ostrého MnS,
- vyřešit otázku aplikace nad externě uloženými daty (například zvážit možnost centralizace obrazů a spojení IIR s generováním standardizované Normal sady dat),
- zapojení technologie do aplikace M-Tool a usnadnit popis oblastí obrazu,
- dále rozšířit koncovo-uživatelské funkce při prohlížení faksimile atp.

Lze také uvažovat o tom, že kvalitu rozeznávaných informací je možné dále vylepšovat v poloautomatickém provozu. S vhodně navrženým nástrojem pro hromadnou poloruční práci (například rychlou vizuální kontrolu a odmazávání nesprávně rozeznávaných oblastí – bordury, marginálie). Rychlost práce by mohlo navýšit využití informací z logu IIR Aplikace.

6 Příloha 1: Seznam dokumentů cíleně vybraných do testovacího vzorku

Seznam je tvořen FyzId jednotlivých digitalizovaných exemplářů. Hledání v pilotu lze provést podle signatury (např. CO.X.13 pro KKPS_CO_X_13_____332I). Prvních 6 znaků FyzId je interní kód knihovny dle seznamu MnS.

KKPS_CO_X_13_____332I
 KNMP_MS_F_1_____181I
 MZM_A_7077_3_____1FAH
 NKCR_XIII_A_6_____1IWS
 NKCR_XXIII_C_124____1MZN
 NMP_25_A_13_____2INJ
 NMP_25_A_5_____18TJ
 NMP_26_B_8_____2TGS
 PNP_TR_I_27_____24NN
 SK_BE_VIII_66_____0QEO
 SK_DE_III_22_____1POI
 SOAZ_FOND_VELKOSTATE101P
 VKO_M_I_306_____05WR
 VMO_K_14905_____01RS
 ZK_20_H_3_14131____03WM
 ZK_20_K_21_14261_H0YIH
 ZK_27_C_18_1_1875516VH
 ZK_7_D_8_5093_____16PJ
 ZMP_502_F_13_____27MI
 ZMP_510_A_011_____0R9N
 ZMP_510_A_016_____2C4O
 ZMP_510_B_003_____0QLJ
 ZMP_510_B_024_____0OSH
 ZMP_510_CH_1_____2SXL
 ZMP_513_A_4_____0UBQ