

Digitální knihovny a využití umělé intelligence

Petr Žabička

Moravská zemská knihovna v Brně

Komunita AI4LAM



- [AI4LAM.org](https://ai4lam.org) (AI pro knihovny, archivy a muzea)
- [slack](#), [twitter](#), [google groups](#), [github](#)
- otevřená platforma pro spolupráce paměťových organizací v oblasti AI
- zlepšuje povědomí a sdílí zkušenosti ve využití nástrojů umělé inteligence

Pravidelná setkání

- Zoom online jednou za ca. 6 týdnů, obvykle v 17 nebo 18 hod. našeho času, k dispozici zápisy a nahrávky jednání
- globální + pracovní skupiny + neformální skupina koordinátorů pro AU a NZ
- Konference Fantastic Futures

Metadata v knihovnách

- Bibliografické záznamy - formát MARC -> MODS
 - až 30 let staré - různá pravidla, konverze
 - neúplné
 - retrokonverze - přepis z lístkových katalogů různého stáří a kvality
 - jazyk záznamu, zkratky
 - nekonzistence uvnitř záznamu
 - duplicitní záznamy
 - nespolehlivé/nestandardní vzájemné provazby
 - popis shora vs. zdola
 - popis na úrovni titulu vs. exempláře
 - přítomnost / absence věcného popisu



Metadata v knihovnách

- Analytické záznamy (články)
 - nejednoznačné linkování do digitálních verzí
- Autoritní záznamy
 - pravděpodobně nejkonzistentnější
- TEI (bibliografický záznam, někdy fulltext)
- Speciální databáze (adresář knihoven apod.)
- Speciální datasety (georeferencování apod.)
- Obohacování (obalkyknih.cz)

<https://knihovny.cz/api>



Data v digitálních knihovnách

Obrazová data

- naskenovaný text
 - více sloupců
 - charakter databáze (slovníky, tabulky,...)
 - speciální objekty (vzorce atd.)
- mapy + mapová díla
- grafika
- noty
- kombinované předlohy (anotace - výřezy)

Textová data

- OCR
- e-born: pdf / epub / ...

Zvuk, video, 3D objekty, datová média, web



Využití strojového učení

- nasazení při digitalizaci
 - náhrada nebo zefektivnění lidské práce
 - zkvalitnění výsledků práce
- zpřístupnění
 - zkvalitnění dat / metadat pro indexaci
 - zkvalitnění dat pro prezentaci (obraz / zvuk / text)
 - zlepšení procesů indexace / vyhledávání
 - přídatné funkce při zpřístupnění
- aplikace pro výzkum
- specifické skupiny dokumentů



Strojové učení a proces digitalizace

Úpravy skenů

- Ořez, vyrovnání stránky, barevné podání
- Exon (Kaitos)

Scelování dokumentů

- rozměrné předlohy skenované po částech
- využití 3D snímání

Struktura dokumentu

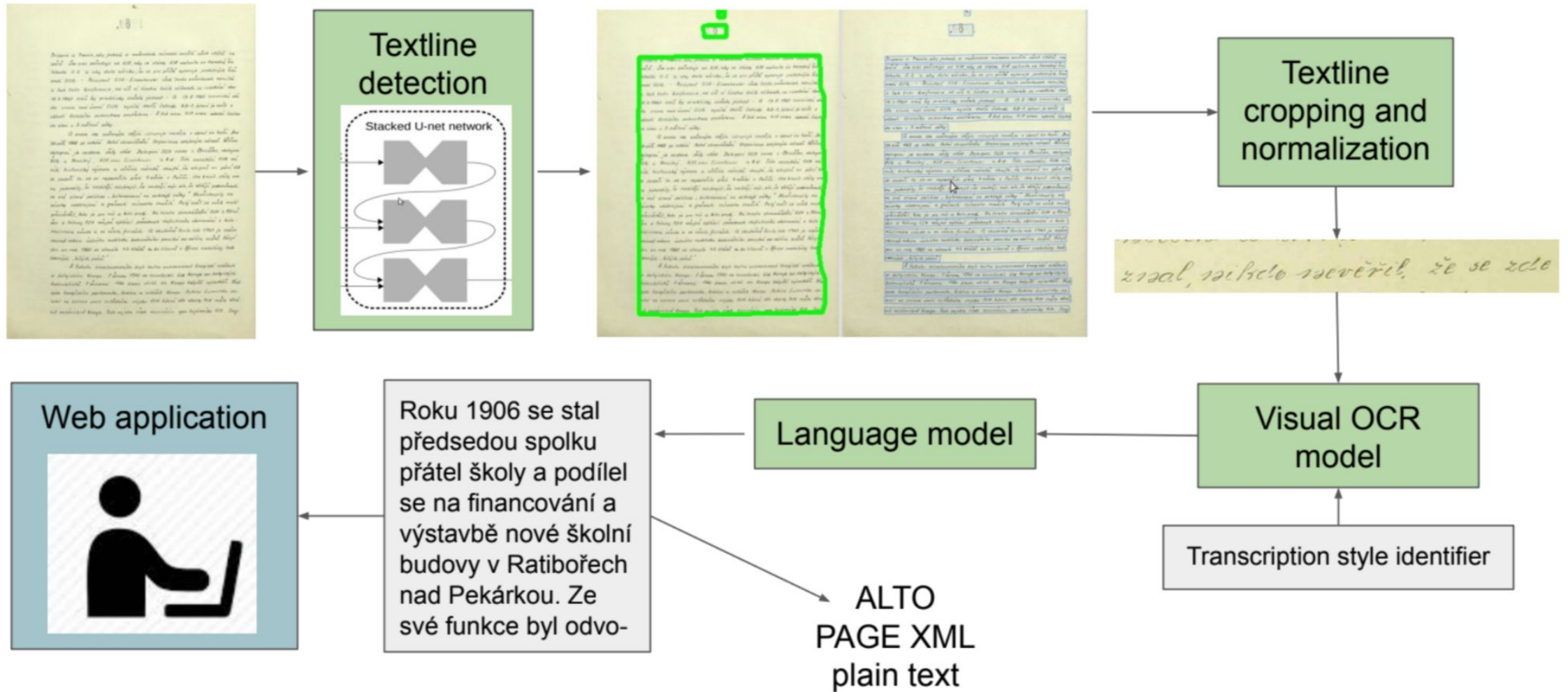
- Identifikace čísel stran, vydání periodik, kapitol, čísel a článků v periodiku, ...

OCR

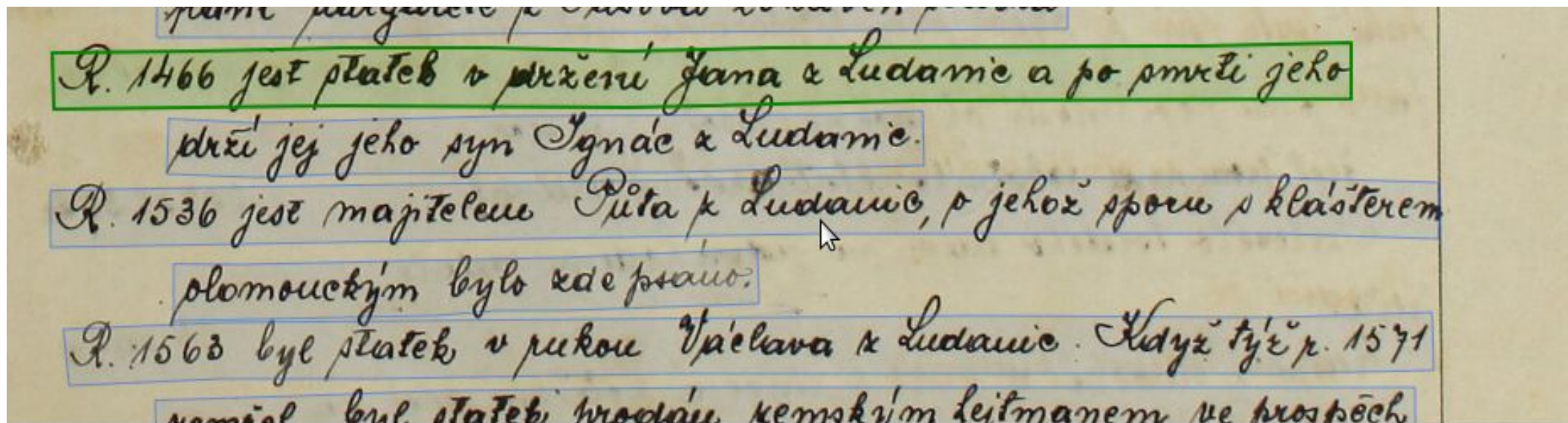
- Tesseract
- PERO



OCR - projekt PERO



České kroniky 20. století

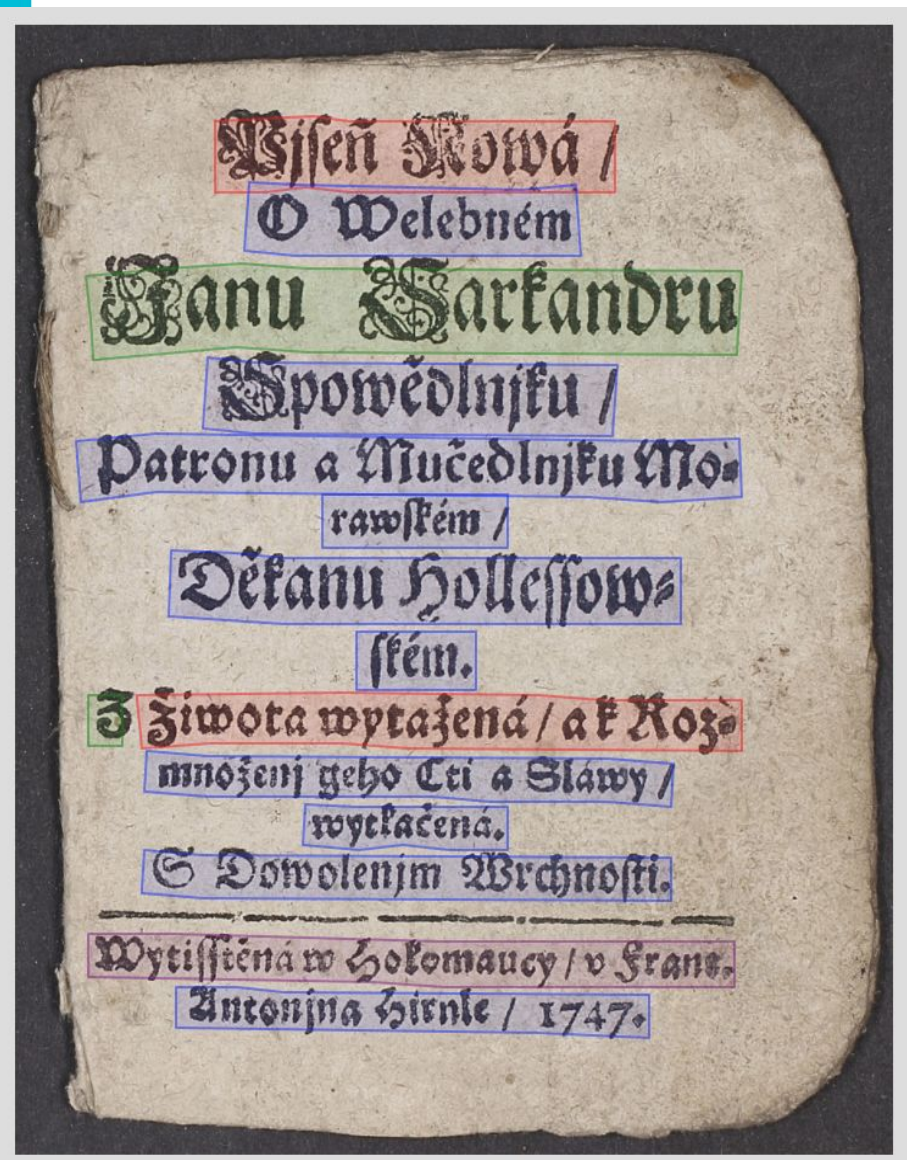


R. 1466 jest statek v držení Jana z Ludanic a po smrti jeho
drží jej jeho syn Ignác z Ludanic.

R. 1536 jest majitelem Půta z Ludanic, o jehož sporu s klášteřem
olomouckým bylo zde psáno.

R. 1563 byl statek v rukou Václava z Ludanic. Když týž r. 1571

Kramářské tisky



Pjšeň **N**owá/

O Welebném

Janu Sarkandru

Spowědnjku/

Patronu a Mučedlnjku Mo=

rawském/

Děkanu Holleffow=

ském.

Žiwota wytažená/ a k Roz=

množenj geho Cti a Sláwy/

wytačena.

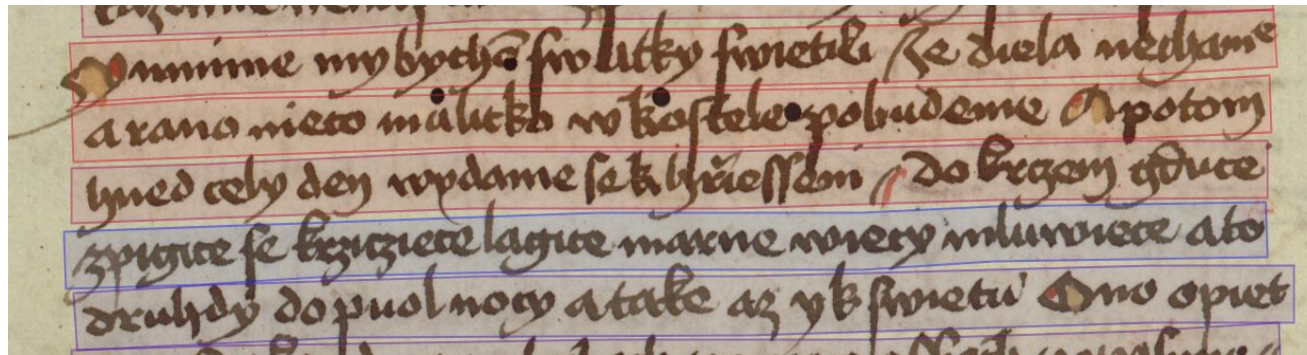
S Dowolenjm Wrchnosti.

Z

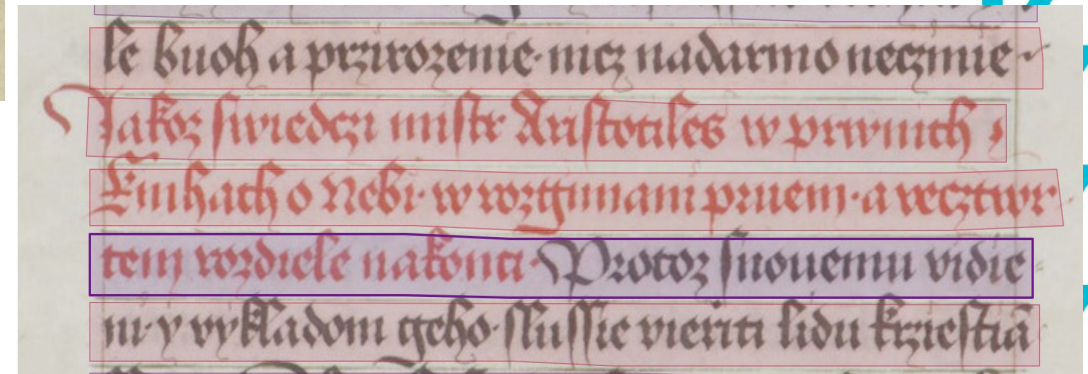
Wytisštěná w Holomaucy/ v Frane.

Antonjna Hirnle/ 1747.

Jan Hus, Vavřinec z Březové

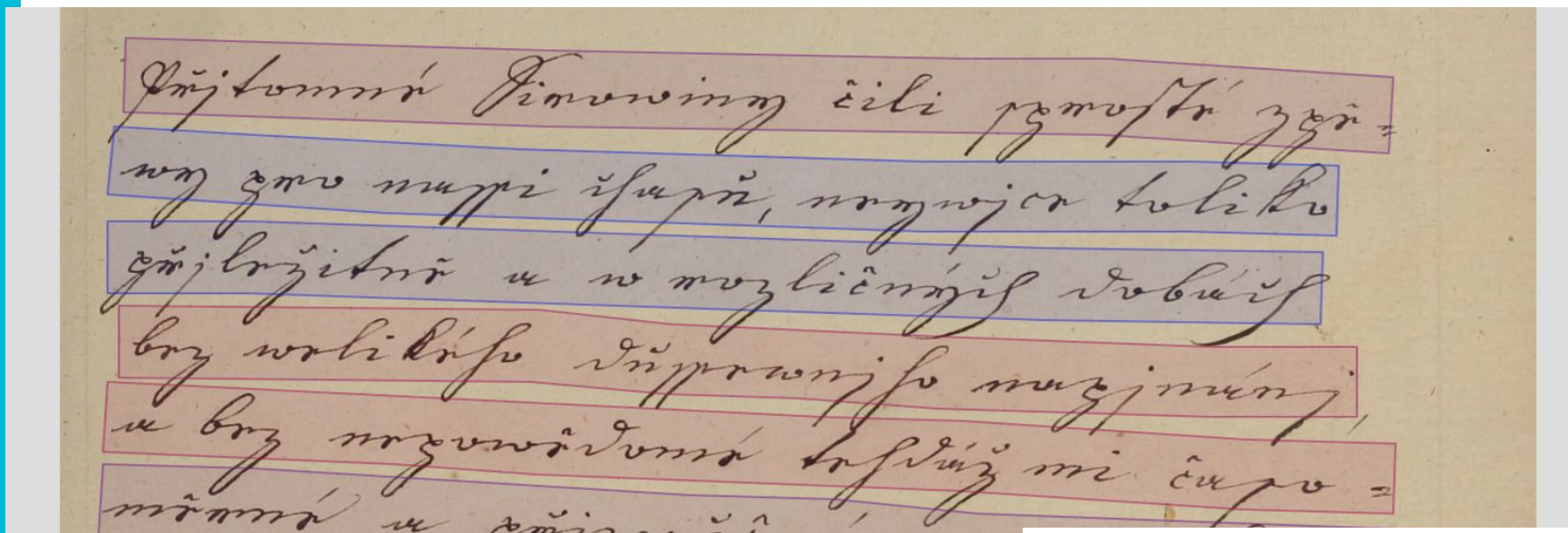


Imnime mybycha swatky swietili se diela nedani
a rano niesto malitko w kostele pobudeme A potom
hned cely den wydame sek hriessom Do krczem gduce
zpigice se krziciece agite marne wiery mluwiece a to
druhdy do puo nocy a take az ys swietu Ono opiet



le buoh a przirozenie niez nadarmo necziie
Jakoz swiedczy mistr Aristociles w prwnich
Eiihach o Nebi w rozgiani pruem a vecztyr
tem rozdziele nakon a Protoz Inouemu vidie
ni y vykladom geho sluffie vieriti lidu krziestia

Kurent



Přjtomné Sirowiny čili sprofté zpě,,
wy pro naši chasu, neywjce toliko
přjležitně a w rozličných dobách
bez welikého duffewnjho napjmáej
a bez nepowědomé tehďáž mi čapo=



PERO - důležité odkazy

- Jádru OCR - pero-ocr python balíček <https://github.com/DCGM/pero-ocr>
- Webová aplikace pro kontrolu a opravy - pero_ocr_web
 - Běží na <https://pero-ocr.fit.vutbr.cz>
 - Zdrojové kódy https://github.com/DCGM/pero_ocr_web
- OCR API pro hromadné zpracování
 - <https://pero-ocr.fit.vutbr.cz/api>
- Informace o projektu - <https://pero.fit.vutbr.cz/>



Strojové učení a zpřístupňování

- analýza+zlepšení dat / metadat pro indexaci
 - identifikace pojmenovaných entit a zvýšení jejich relevance
 - automatická předmětová kategorizace
 - identifikace a kategorizace obrázků
 - identifikace osob na obrázcích
 - speech2text
- zkvalitnění dat pro prezentaci (obraz / zvuk / text)
- zlepšení procesů indexace / vyhledávání
 - např. Solr + LearningToRank / RankNet
- vyhledávání obrázků / částí obrázku
- přídatné funkce při zpřístupnění
 - text2speech
 - překlad do jiného jazyka
 - doporučování (na základě dokumentu, na základě historie uživatele)





KNIHOVNY.CZ a strojové učení

- první experimenty - rozdělení na beletrii a odbornou literaturu
 - problém byly např. čítanky
- přiřazení třídy konspektu
 - bylo identifikováno několik nejednoznačných tříd = tříd, u kterých váhá i knihovník
 - došlo k redukci (zjednoznačnění) tříd konspektu pro tento účel
 - bylo nutné opravit v trénovacích datech neplatná MDT a nesprávně přidělené znaky konspektu
- pro publikace v češtině
- původně s využitím fulltextu, později jen na základě bibliografického záznamu
- pro některé třídy konspektu chybí dostatek trénovacích dat
- velmi přesné tam kde bylo v záznamu MDT
- konspektem obohaceno téměř 300.000 záznamů monografií
- využití zejména při filtrování (faseta Obor = konspekt)





KNIHOVNY.CZ a strojové učení

Nejčastější záměny při klasifikaci s využitím fulltextu (na základě testovací datové sady)

Originální kategorie	Kategorie přiřazená klasifikátorem
Dějiny zemí střední Evropy	Dějiny Česka a Slovenska
Malířství	Výtvarné umění
Řízení a správa podniku	Management. Řízení
Literatura. Literární život	Česká literatura (o ní)
Výtvarné umění	Malířství
Geografie Česka a Slovenska, reálie, cestování	Dějiny Česka a Slovenska
Umění	Výtvarné umění
Biografie	Vnitropolitický vývoj, politický život
Biografie	Film. Cirkus. Lidová zábava
Dějiny Česka a Slovenska	Dějiny zemí střední Evropy



Zlepšování obrazu

- Během digitalizace
- Pro zpřístupnění

Jak si radnice počíná proti úřadům, svědčí to, že nechala reklamace české okresní školní rady, jež podány byly na vyloučení 600 českých dětí z německých škol zcela klidně ležet a věc tu nevyřizuje. Radnice je zde podporována vládou a proti takovému nepříteli nezbývá žádný jiný prostředek, než jen jeden — volební lístek. (Bouřlivý souhlas.)

Za zřízení českých měšťanských škol v Brně bojují Čechové již od r. 1881, leč

before

Jak si radnice počíná proti úřadům, svědčí to, že nechala reklamace české okresní školní rady, jež podány byly na vyloučení 600 českých dětí z německých škol zcela klidně ležet a věc tu nevyřizuje. Radnice je zde podporována vládou a proti takovému nepříteli nezbývá žádný jiný prostředek, než jen jeden — volební lístek. (Bouřlivý souhlas.)

Za zřízení českých měšťanských škol v Brně bojují Čechové již od r. 1881, leč

after

“Vylepšování” obrazu

- na přání

Бѣара іаѢас нѣѢрыѣѢ ѢѢ	Dobrupospol.	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	[REDACTED]	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	Bratislava	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	Knižovna	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	Polokvium	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	lenochod	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	monitor	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	basšytara	temu af. ѢѢ.	Dobrze b
Бѣара іаѢас нѣѢрыѣѢ ѢѢ	gentleman	temu af. ѢѢ.	Dobrze b

Identifikace obrázků na stránce

- Cílem je automaticky detekovat pozici obrázků v naskenovaných dokumentech
- Proč je to těžké?



Co není text ještě nemusí být obrázek

- Grafické elementy, artefakty, pozadí

Žízeň kasi
CITRONKA.
Náhrada čerstvých citronů spřirod. ovoc.
šťavami. Láhev 7/10 lit. po K 8-50,
stačí na 25 litrů
jemného lihuprostého nápoje. Med, medo-
vina, ov. malaga s med., jitrocel s med. a j.
Dlouhý-Med-Soběslav
Filial.: Praha II. Vodičkova, palác České
banky. Kr. Vinohrady, Jungmannova 31.
4181

Koncerty a zábavy.
JOKOHAMA. Úspěšný
program.
9403
Již jen několik dní.

účetní kapitál K 13,000,000--
Telefon 4
šéfredakta Rozumová hani

Ročník IV. (1922). Str. 3.

Sbírka
rozhodnutí nejvyšších stolic
soudních republiky
československé.

Rozhodnutí nejvyššího správního soudu
ve věcech správních.

Z příkazu prezidia nejvyššího správního soudu pořádl JUDr. JOS. V. BOHUSLAV
senátní prezident nejvyššího správního soudu

Vydání druhé měkké
Obsah na straně druhé.

V PRAZE 1923.
Vydavatelství nejvyššího správního soudu. — Administrace Král. Vinohrady čp. 1234.
Nakladatel a výtiskárna: Právnická tiskárna v Praze, spol. s r. o. — Zodpovědný redaktor
JUDr. Václav Tomáš advokát Král. Vinohrady, u divadla č. 1.

Dnešní noviny

Užijte si
www.IHNED.cz

STAVEBNÍ
V Česku ub
nejméně v l
letruhu v Br
méně než v
z dění na ve
ru, který bo

VANHAR
Rožený Brň
dal majetek
podnikat. P
jichž obrat
maji malý r
ničí," říká o

ELEKTRIK
New York p
to dostalo
v náročném
Bloomberg
vozů elektr

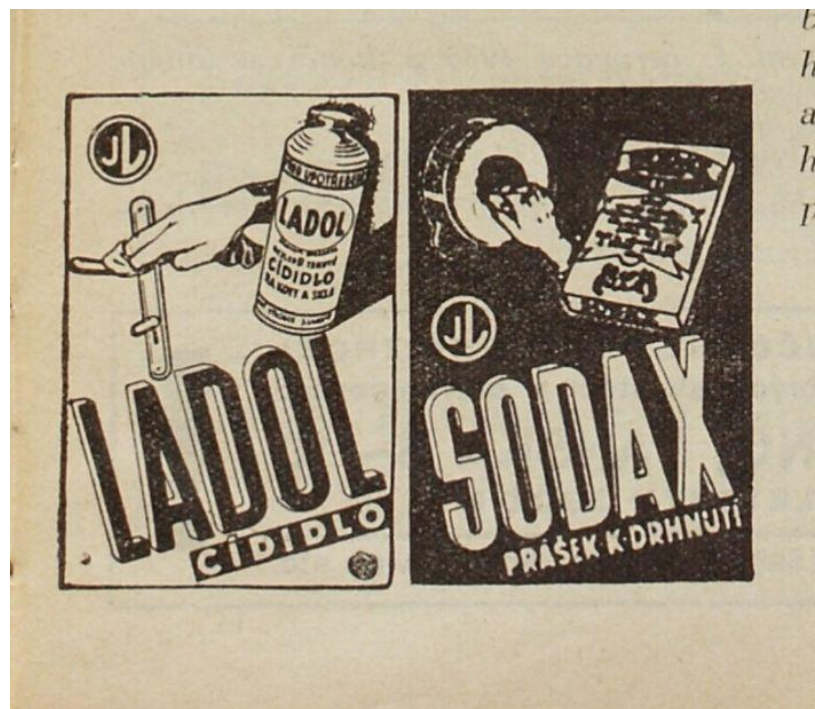
RECENZE
Povedená D
vě dobrým
litním vide
i citelný ne
sionálního.

HN Exkluz



Co obsahuje text stále může být obrázek

- text v obrázcích, překryv



Pokud vám rodinný rozpočet právě teď dovoluje koupit si jen jedno zkrášlovadlo, vyberte si voděodolnou řasenku. Nalíčené řasy dodají tváři výraz a upravený vzhled a je to otázka dvou minut, které se dají ukrást i v tom nejhroším ranním sklužu. Složení, odolávající dešti, vám dodá jistotu, že se ani po srážce s podzimní plískanicí neproměníte ve smutnou pandu.

ZKUSTE:
Voděodolnou objemovou řasenku Oriflame Beauty (199 Kč) se silikonovými složkami, které nabídnou voděodolávající vlastnosti a prodlouží trvanlivost nalíčení.

Oriflame BEAUTY
WATERPROOF mascara

(ten) / Foto archiv Irem

Obrázky nemusí být ohraničeny



vodafone

Žádná bouda.

Ale skutečně neomezené volání a SMS až 4 kamarádům.

Ve Vodafone si myslíme, že by člověk člověku měl být přítelem, ne výem. Proto s vámi jednáme na rovnou. Žádné uzavírání smlouvou. Žádné skryté podmínky. Žádné hvězdičky ani žádné poznámky malé jako psi blechy. Jen přátelské nabídky. Jako třeba Program kamarádů. Tady stačí podle počtu kamarádů zaplatit měsíční paušál od 179 do 286 Kč (vč. DPH) a můžete jít v síti Vodafone neomezeně volat i posílat SMS. Nabídka platí pro vybrané tarify v síti Vodafone.

Nyní 50% sleva na 2 měsíce při aktivaci do konce dubna.

Aktivace a více informací na **800 777 777** nebo na www.vodafone.cz

Teď je to ve vašich rukou.

50% SLEVA PŘI AKTIVACI DO KONCE DUBNA

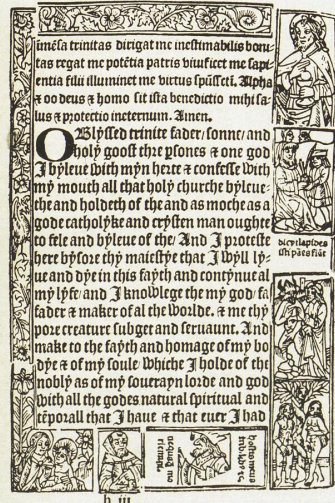
Kategorie nejsou jasné ani lidem

- Sken textu, typografie, grafy a tabulky

1481 překlad z holandštiny The Historye of Reynart the Foxe. Největší kniha a záslužným překladem bylo patrně vydání Caxtonovy verze knihy Zlatá legenda (Golden Legend) podle knihy ze 13. století od Jacobuse de Voragise Legendae aureae, přičemž však Caxton, který byl hotov s překladem v listopadu 1483, použil pro svou kompilaci dvou nových verzí, francouzské a anglické. Kniha je krásně vypravená

téhož roku, Blanchardyn and Eglantine, The Foure Sonnes of Aymon a Fayette of Armes and of Chyualrye Christy de Pisan roku 1489, Morale Proverbs atd. Z těchto knih byly mnohé znovu vydány v dnešní době a jsou běžně k dostání.

Knihy Willama Caxtona nemají titulní stránky a od roku 1487 jsou většinou zdobeny zajímavým štitkem o velikosti 5,5 x 4,5 palce (14 x



Stránka z knihy vytištěná roku 1499
Výsledkem
de Worde
ve Westminsteru

Q uo quisque in se ipso non habet vitam eternam, sed solummodo vitam in seculo. Quia ergo quisque in se ipso non habet vitam eternam, sed solummodo vitam in seculo. Quia ergo quisque in se ipso non habet vitam eternam, sed solummodo vitam in seculo.



o čem se nemluví

HLAVNÍ TÉMA: RODINNÉ VZTAHY BY MĚLY BYT TY NEJPEVNĚJŠÍ. JENŽE MÍSTO TOHO JE RODINA ČASTO MÍSTEM NEPOCHOPENÍ.

Ten dům je naše prokletí

ADRIANA (45): „ASI DOJDE AŽ K SOUDU.“

CO TOMU ŘÍKÁTE: VĚRA (54)
Adriana má tři děti. Protože tasy byli na venecii, umi si představit, kolik práce a starostí se vším mělo. Jedine řešení by asi bylo najmout si nějakého právníka, který by uměl poradit, jestli máš dědit, nebo ne. Ale vždyť se souzence se už asi nikdy nenapraví.

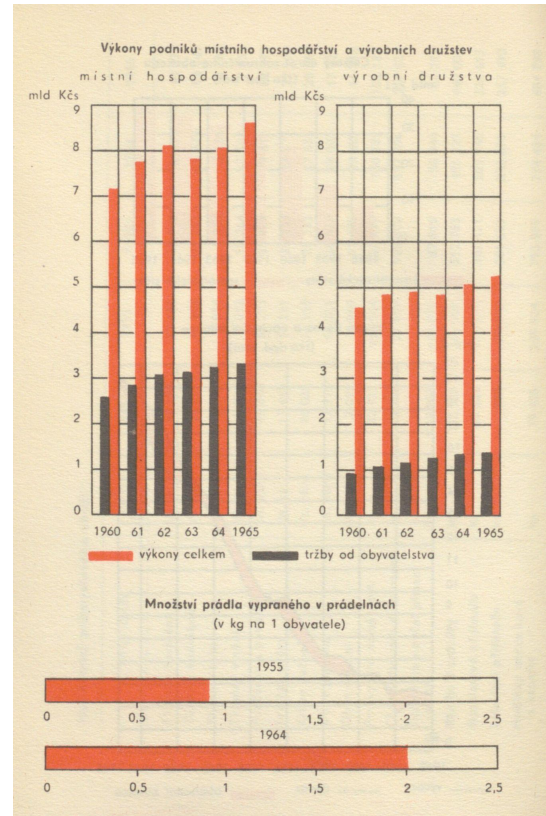
Něde rodina nikdy nebyla vzorem pospolitosti, kde by všechny nezájímající pomohli a byli si oporou. Nevím, jak se to stalo, asi to byla nějaká výchovná chyba rodičů nebo sp, ale zkrátka já, sestra a bratr nejsme žádní dokonalí souzence. Mezi mnou a sestrou je rozdíl po věku, bratr je ještě o dva roky mladší než Pepina. Jako malí jsme si spolu samozřejmě hráli, ale jakmile jsme šli do puberty, přestali jsme si rozumět.

Každý si našel kamarády, a i když jsme se nikdy nijak zvlášť neshádali, začali jsme si být cizí. Navíc sestra na mě začala žádat, možná proto, že jsem starší. Já jsem se taky brzo vdala a měla děti a jsem s mužem dodnes šťastná, zatímco ona se sice taky brzo vdala, ale stejně brzo se rozvedla. Od té doby, a to už je přes dvacet let, se plací mezi mužskými, ale žádný s ní nevydrží víc než pár roků. Žáá se mi taky zapokká, takže kolo by s ní byl... Bratr zase nemá žádný cíl, žije se tím, co se zrovna namane, a taky dost pije. Já se svým spořádaným životem jím možná píju krev, možná jsem vzor toho, co oni by sice rádi, ale nedokážou to. Rozhodně, jak jsem se dozvěděla, pro mě nemají moc dobrých slov, dokonce prý sestra říkala našim známým, že mi je manžel nevěrný a já se tvářím, že o tom nevím.

Nehořší situace nastala po smrti našich rodičů. Táta umřel už dávno, ale maminka ani ne před půlrokem. Požít je v tom, že my jsme byli s mužem a dětmi

a ními v jejich původním domě. Ne že bych o to nějak záležela, ale Pepina se odebírala se svým mužem do města a pak už tam jenom strádala bydlíště. A bratr zmizel v osmnácti... Nikdy nevím, kdy se náhodou objevil. Takže naše rodina zůstala v domě, což sice zní jako výhra, ale kdo někdy měl něco takového jako vesnický dům, ví, že to je hrůzná práce a stojí to spousta peněz. Těch hodin, co manžel strávil při opravách! V životě jsme nebyli na pořádné dovolené, protože všechno, co jsme ušetřili, šlo na opravy domu, topení a tak. Maminka říkala, že ona nám s dítětem přispívat nemohla, ale že by si přála, aby za naši práci a investice a také za to, že jsme se o ni starali, když už byla nemocná, připsali dům nám. Proto taky naptsala závěť, ve které to píše. Potom umřela a začalo peklo. Sestra si napden vzpomněla, kolik asi takový dům stojí a že my nemáme peníze na něj dělat náklady. Podle ní by se měl dům prodat a peníze by se měly rozdělit. Nebo, když my tam chceme bydlet, bychom měli jít i bratra vyplatit. Bratr a ni samozřejmě souhlasil, přestože nevěděl do oprav ani korunu. A že by se třeba staral jeden nebo druhý o dmu, když už nemohla, chudák, ani sama dojit na zachod, o tom ani nemluhám.

Jsem nešťastná. Prodat dům, to si neumím představit, protože kdy bychom bydleli? A kdy bychom měli peníze na to, abych souzence vyplatila, taky nevím, protože všechny peníze, co jsme měli navíc, šly do materiálu a na dělníky. Sestra říká, že jestli okamžitě nenavrhu nějaké řešení, dá vše k soudu. Jestli jsme se ani nenasázili zjistit, jak to právně všechno je. Vím jenom, že mě mrzí nejenom to, jak se teď souzencei handrkují o něco, na čem jim nikdy nezáleželo, ale hlavně to, že se o maminku pořádně nezajímali, když byla ještě na světě. Podle mě se musí v hrůbě oloupat. Manžel mě uklidňuje, že vše nějak dopadne, ale samozřejmě nemá prý svagra se svaarovou dobré slovo. Takže, když jsem plakala, se rozčilil a řval, že už nikdy nepřekročí náš práh. Copak takhle se chováš lidé v rodině? Je mi jasné, že takovéto spory se už nikdy nenapraví, protože jsme už příliš rozdílní.



Dostupná řešení - OCR od ABBYY (v ALTO)

- Dataset: 500 náhodných, manuálně oannotovaných stránek
- Výsledek:
 - Celkové IOU*: 0.69
 - Sensitivity**: 0.69
 - Precision**: 0.36
 - 18% obrázků není detekováno vůbec
 - 57% detekovaných obrázků jsou falešná pozitiva***

* Intersection over union

** Vypočteno pro práh citlivosti IOU=0.5

*** Neobsahují obrázek, nebo obsahují obrázek který už byl detekován, tzn. obrázky se překrývají.



OCR od ABBYY Ruční anotace

Neuronová síť dhSegment

- Celkové IOU*: **0.65**
- Sensitivity**: **0.65**
- Precision**: **0.4**
- **24%** obrázků není detekováno vůbec
- **53%** detekovaných obrázků jsou falešná pozitiva***
- Rychlost (notebook bez GPU):
 - 2.87s na stránku
 - 33 dní na milion stránek
 -

Při vynaložení relativně malého množství práce je kvalita prakticky identická s OCR od ABBYY

- Celkové IOU*: **0.69**
- Sensitivity**: **0.69**
- Precision**: **0.36**
- **18%** obrázků není detekováno vůbec
- **57%** detekovaných obrázků jsou falešná pozitiva***

* Intersection over union

**Vypočteno pro práh citlivosti IOU=0.5

***Neobsahují obrázek, nebo obsahují obrázek který už byl detekován, tzn. obrázky se překrývají.



Slepé uličky

- **Newspaper Navigator model**
 - Vyvinut v Library of Congress na detekci obrázků, nadpisů a dalších objektů ve starých novinách.
 - Založen na síti Faster-RCNN od Facebooku.
 - Při evaluaci (bez našeho trénování) měl mnohem horší výkon než OCR od ABBYY.
 - Domníváme se, že model je náchylný k overfittingu a špatně generalizuje na náš dataset.
- **Natrénování UNet sítě (bez předtrénovaných vah)**
 - Architektura využívaná k segmentaci v medicíně a při zpracování satelitních snímků.
 - Nepodařilo se nám dosáhnout dostatečné přesnosti, náš dataset byl zřejmě příliš malý na to, aby se model naučil všechny potřebné vizuální znaky (features).



Vyhledávání obrázků podle podobnosti - VISE

- Vyhledávání příbuzných obrázků podle prostorové podobnosti (výřez)
- Periodika, staré ilustrace, grafiky, loga...
- Přesný i pro malé rozlišení (testováno na 1024x1024)
- Možnost kombinace daty identifikujícími obrázky na stránce
- 13 000 obrázků zaindexováno za 7 hodin



Datová sada obrázků (jpg, png)

Indexace / trénování
vizuálních deskriptorů..



Zaindexovaná datová sada

Search ready



Demonstrace na inzerátech deníku Svobodné slovo

ARMIN **HYDRO**

jest nejlepší a nejlevější domácí výrobek továrny v Praze
Jan Hubinek a syn, Praha-VII, Libeň.
Dělní zábratka za pět set Kč a v celku s daní za zboží.
Všechny práce. Město 10. května 1918.
Kd' Píseň obecního úřadu v celku s daní za zboží.

BOH. KREJČÍK,
velkoobchodní úřad, Lázeň - státní - pražská, ul. Na Příkopě 101/102.

Romány z českých dějin od Jos. Svátka

První díl: "Týden prázdný" - Praha a Fress, 100 Kč. Druhý díl: "První světová válka" - Praha a Fress, 100 Kč. Třetí díl: "První světová válka" - Praha a Fress, 100 Kč. Čtvrtý díl: "První světová válka" - Praha a Fress, 100 Kč. Pátý díl: "První světová válka" - Praha a Fress, 100 Kč. Šestý díl: "První světová válka" - Praha a Fress, 100 Kč. Sedmý díl: "První světová válka" - Praha a Fress, 100 Kč. Osmý díl: "První světová válka" - Praha a Fress, 100 Kč. Devátý díl: "První světová válka" - Praha a Fress, 100 Kč. Desátý díl: "První světová válka" - Praha a Fress, 100 Kč.

Svatovítského pojítkového závodu
na Svatovítské, Svatovítské náměstí č. 102,
jednotlivě každé číslo od 1000 Kč.

Americká PULTY
Jos. Jiroušek,
Praha 1, Svatovítské náměstí 102.

Vazby skvostné
Prostřední ul. č. 102, Praha 1.

Japonská restaurace v Praze a Fress
Královské náměstí 102, Praha 1.

PROČ KUPUJÍ NAŠE DÁMY TAK RÁDY

Dr. Frant. Bačkovský,
Lázeňská ul. č. 102, Praha 1.

Výborný prádlo pro nevěsty.
J. NOVÁK, Tuřínova ul. č. 102, Praha 1.

Hledat

Výhody VISE

- Velmi dobrá přesnost i pro obrázky s malým rozlišením
- Vyhledávání je rychlé
- Není příliš mnoho konkurenčních systémů
- Není potřeba grafické karty
- Systém je dále rozvíjen

Nevýhody

- Nefunguje pro běžné bloky textů, vhodné spíše pro obrázky, ilustrace nebo větší texty jako tituly, nadpisy atd.
- Může nastat nepřesnost ve vyhledávání, např. pokud je výřez málo detailní nebo je špatná trénovací sada
- Nelze použít obrázky s vysokým rozlišením
- Nelze přidávat nové obrázky již k natrénované datové sadě
- Nelze bez konverze použít JPEG2000



VISE - Odkazy

- Abhishek Dutta, Relja Arandjelović, and Andrew Zisserman. 2021. VGG Image Search Engine. from <https://www.robots.ox.ac.uk/~vgg/software/vise/>
- Gitlab: <https://gitlab.com/vgg/vise>
- Oxford Visual geometry group: <https://www.robots.ox.ac.uk/~vgg/>



Další projekty

- Webarchiv.cz: strojová analýza typu stránky, tématu, sentimentu
- [Projekt Impresso](#):
 - zpracování a analýza 200 let historických novin v různých jazycích
 - zpracování přirozeného jazyka (NLP) pro postkorekci OCR, indexování n-gramů, distribuční sémantické indexování, zpracování pojmenovaných entit a kategorizaci textu a shlukování
- [Gallicapix.bnf.fr](#)
 - vyhledávání obrázků z digitální knihovny Gallica
- [The Newspaper and Photos project](#)



Závěrem

- Existuje řada zajímavých projektů, aplikací, modelů
- Není snadné najít hotové řešení
- Specifika reálných datových sad
 - rozsah, variabilita
- Velký prostor pro další rozvoj



Děkuji za pozornost!

Dotazy?

Petr Žabička

Moravská zemská knihovna v Brně