

# Obrazová data a strojové učení

Petr Žabička, Ján Bogár, Michal Tran - Moravská zemská knihovna v Brně

# Obsah

- Úvod
- Projekt PERO
- Identifikace obrázků na stránce
- Systém VISE



# Strojové učení a proces digitalizace

- Ořez, vyrovnání stránky, barevné podání
- Scelování dokumentů
- Struktura dokumentu
- OCR
- ...




# Strojové učení a zpřístupňování

- Vyhledávání obrázků
- Vyhledávání a podobnost částí stránek
- Identifikace obsahu obrázku
- ...



# OCR - projekt PERO



←  Digitální studovna Ministerstva obrany ČR

Hledat v celé digitální knihovně

Sbírky Procházet Informace English

🖨️ 📄 🖼️ T

Hledat v dokumentu


9 z 122 stránek

1 2 3 4 5 6 7 8 9 10 11 12

3

smu si naši byt na půdě  
malého domku a dva dny  
smu čekali co si bude dít!

Domek v  
Petrovaradině



Těti dny jsme se měli se

Muj Deňik od r 19 30/7 14

Nakladatelské údaje  
[S. I.], 1914-1915

Typ dokumentu  
Kniha

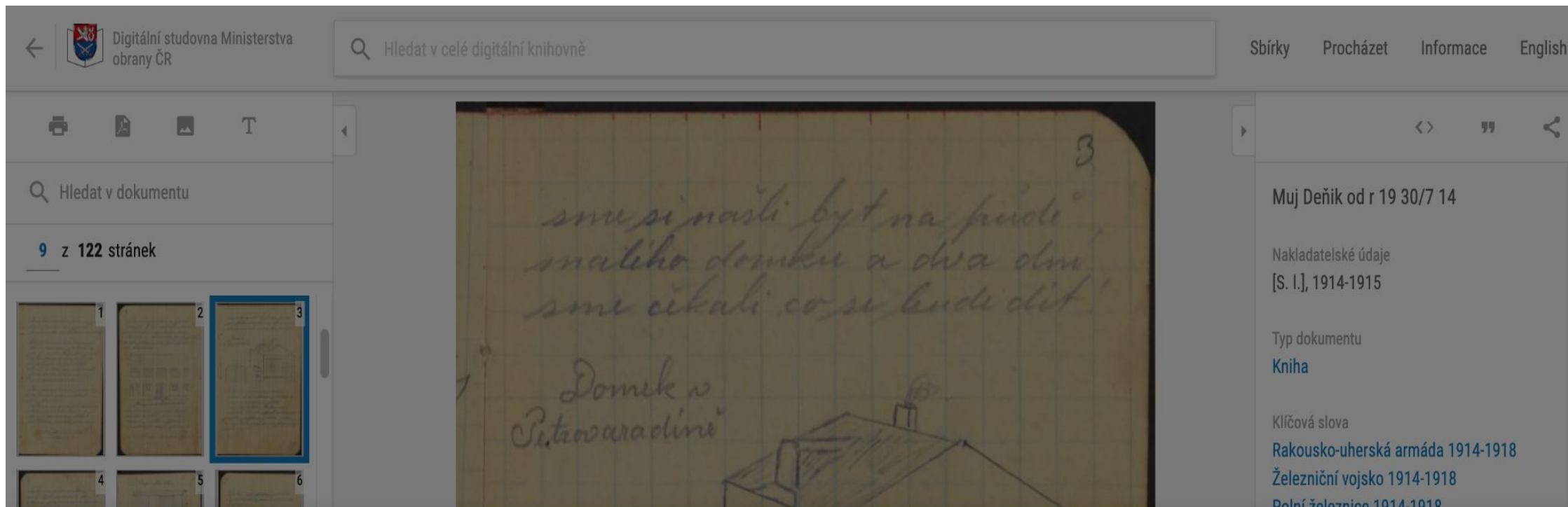
Klíčová slova  
Rakousko-uherská armáda 1914-1918  
Železniční vojsko 1914-1918  
Polní železnice 1914-1918  
Srbská fronta 1914-1915  
Východní fronta 1914-1918

Jazyk  
Čeština

Místo uložení  
Vojenský historický ústav Praha  
Signatura: XIII-10715 (př. č. 2862/2005)

Fyzický popis

# OCR - projekt PERO



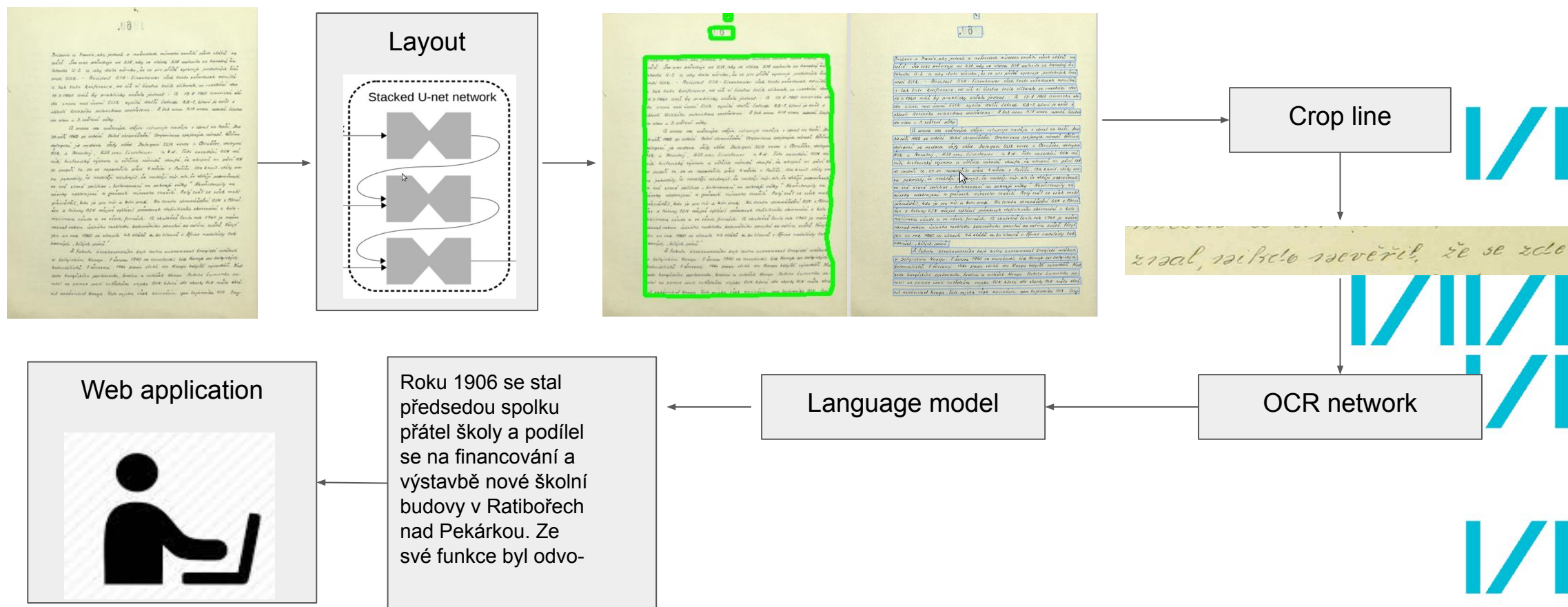
The screenshot shows a digital library interface. At the top left, there is a logo and the text "Digitální studovna Ministerstva obrany ČR". A search bar contains the text "Hledat v celé digitální knihovně". On the right, there are navigation links: "Sbírký", "Procházet", "Informace", and "English". Below the search bar, there are icons for print, download, and text. A search bar for the document is labeled "Hledat v dokumentu". Below that, it says "9 z 122 stránek". A grid of document thumbnails is shown, with the third one highlighted. The main content area displays a handwritten page with the text: "sme si našli byt na půdě malého domku a dva dny sme čekali co se bude dít!" and a drawing of a house. The page is numbered "3" in the top right corner. On the right side, there is a metadata panel with the following information: "Muj Deňik od r 19 30/7 14", "Nakladatelské údaje [S. I.], 1914-1915", "Typ dokumentu Kniha", and "Klíčová slova Rakousko-uherská armáda 1914-1918, Železniční vojsko 1914-1918, Polní železnice 1914-1918".

sme si našli byt na půdě malého domku a dva dny sme čekali co se bude dít!

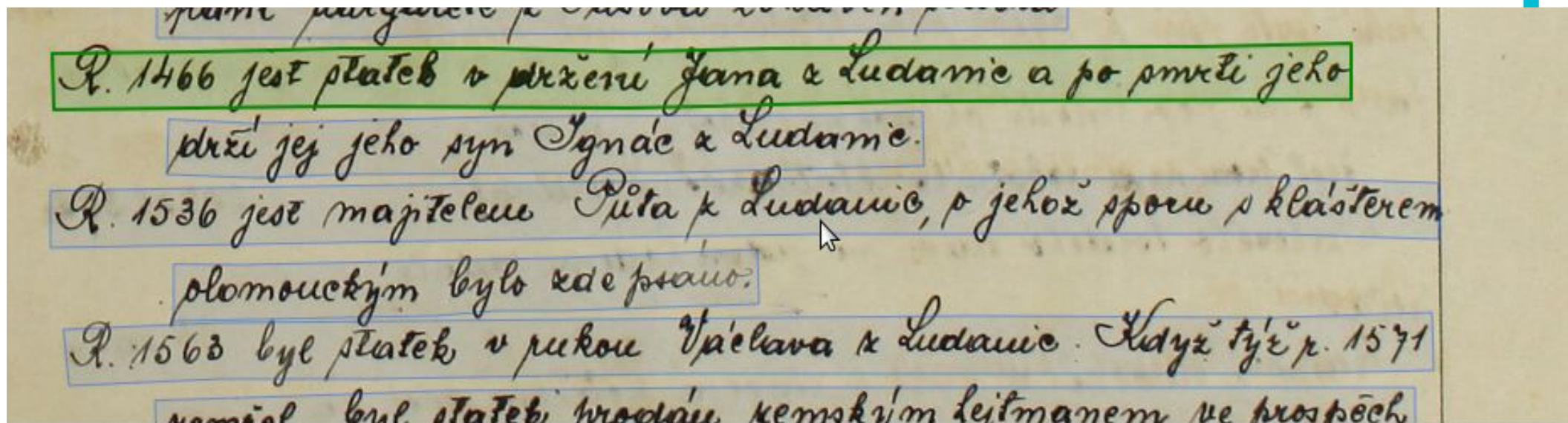
Muj Deňik od r 19 30/7 14. [S. I.], 1914-1915, s. 3. Dostupné také z: <http://www.digitalniknihovna.cz/dsmo/uuid/uuid:32f7ba5e-a323-11ea-95c0-001b63bd97ba>

# OCR - projekt PERO

## Řetězec procesu rozpoznání textu



## 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



## Periodika na mikrofilmech

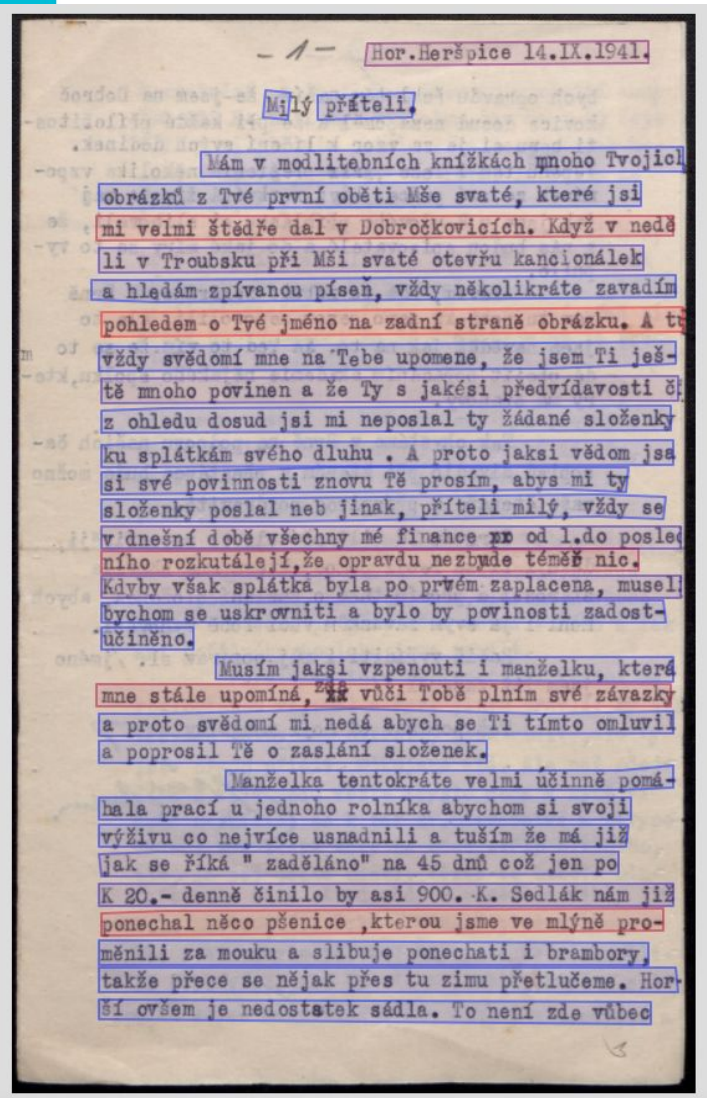
Socialisté italské mají za sebou veliký úspěch: vymožení amnestie pro odsouzené z posledních bouří. Co se nakřičely reakcionářské listy, že socialisté svou hlučnou agitací pro amnestii odsouzeným vlastně škodí, poněvadž prý vláda nemůže ustoupiti nátlaku z ulice -- události však daly socialistům za pravdu. Mocné hnutí v celé zemi přinutilo vládu k silným ústupkům. Nejsou sice amnestováni všichni političtí provinilci, ale téměř 3000 obětem vojenských i civilních soudů po milánských bouřích, kteří ni-

nakřičely

zence z posledních bouří. Co se nakřičely reakcionářské listy, že socialisté svou hluč-

reakcionářské

# Strojopis

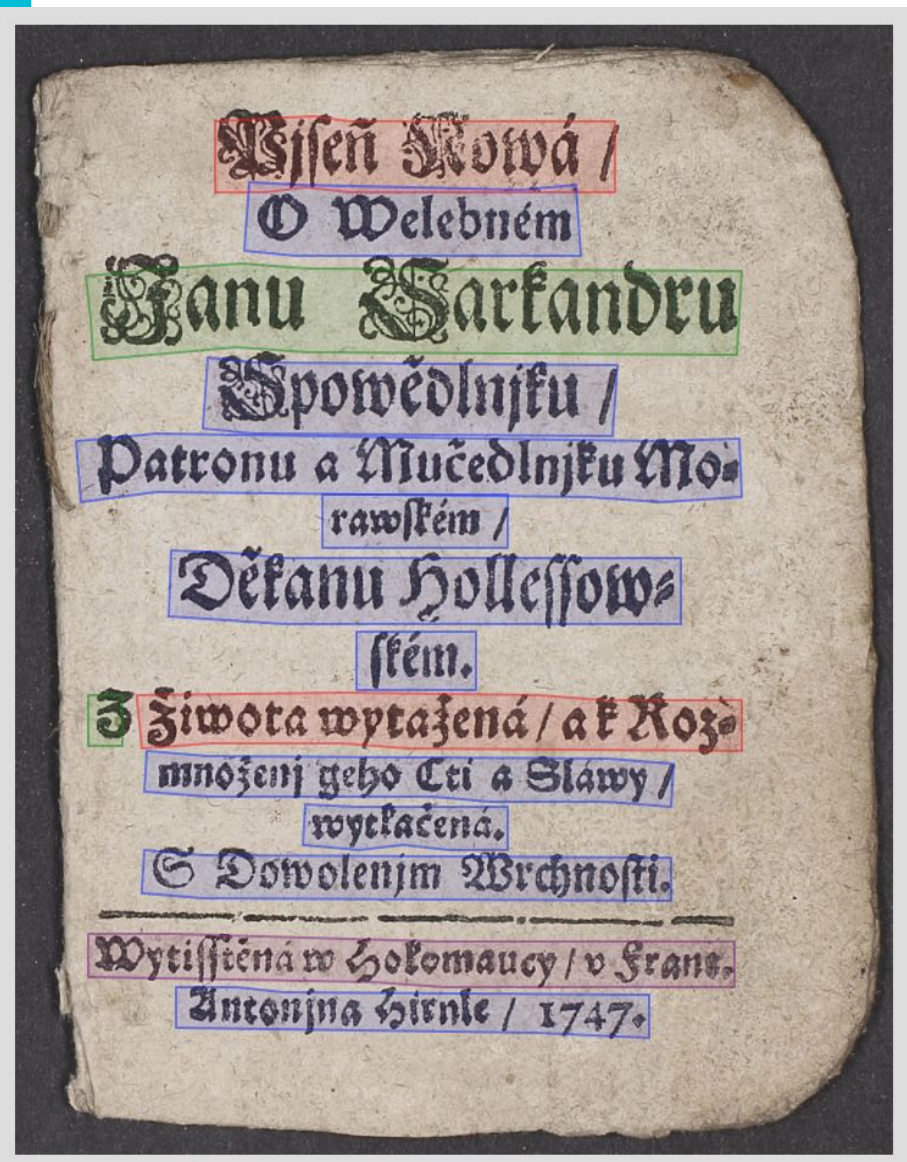


v dnešní době všechny mé finance od 1. do posledního rozkutálejí, že opravdu nezbyde téměř nic.

Kdyby však splátka byla po prvním zaplácena, museli

finance ~~xx~~ od 1. do posledního rozkutálejí, že opravdu nezbyde téměř nic.  
po prvním zaplácena, museli:

# Kramářské tisky



Pjšeň **N**owá/

O Welebném

Janu Sarkandru

Spowědljku/

Patronu a Mučedlnjku Mo=

rawském/

Děkanu Hollešow=

ském.

**Z**iwota wytažená/ a k Roz=

množenj geho Cti a Sláwy/

wytačena.

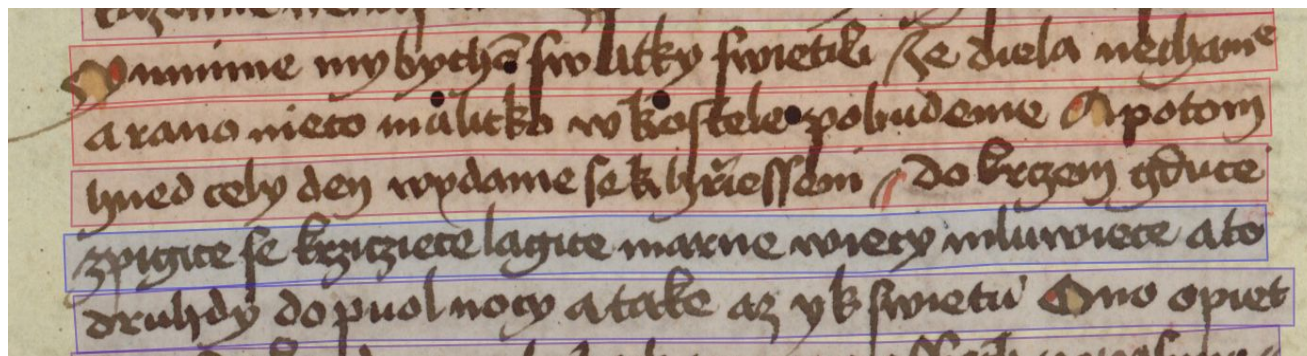
S Dowolenjm Wrchnosti.

**Z**

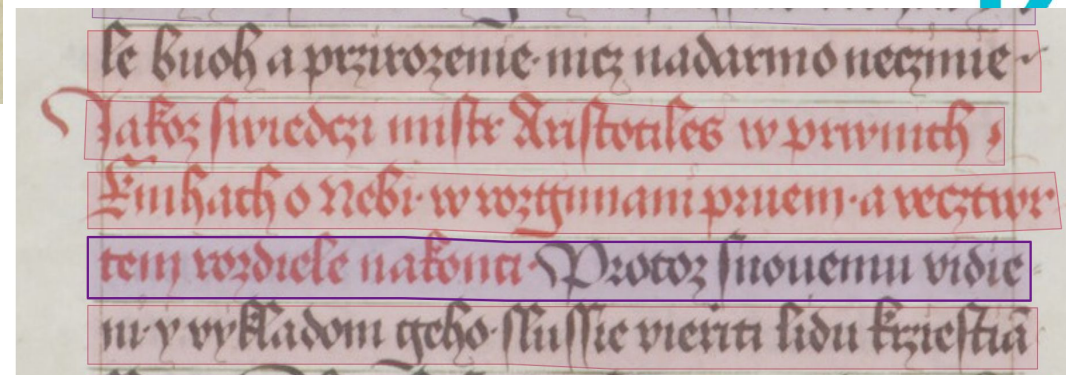
Wytisštěná w Hořomaucy/ v Frane**.**

Antonjna Hirnle/ 1747.

# Jan Hus, Vavřinec z Březové



I mnime mybycha swatky swietili se diela nedani  
a rano nieto malitko w kostee pobudeme A potom  
hned cely den wydame sek hriessem 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·  
Jakoż swiedczy mistr Aristocileš w prwnich  
Eiihach o Nebi w rozgiani pruem·a vecztryr  
tem rozdziele nakonā Protoz Inouemu vidie  
ni·y vykladom geho sluffie vieriti lidu krzieſtia

# České kroniky 20. století

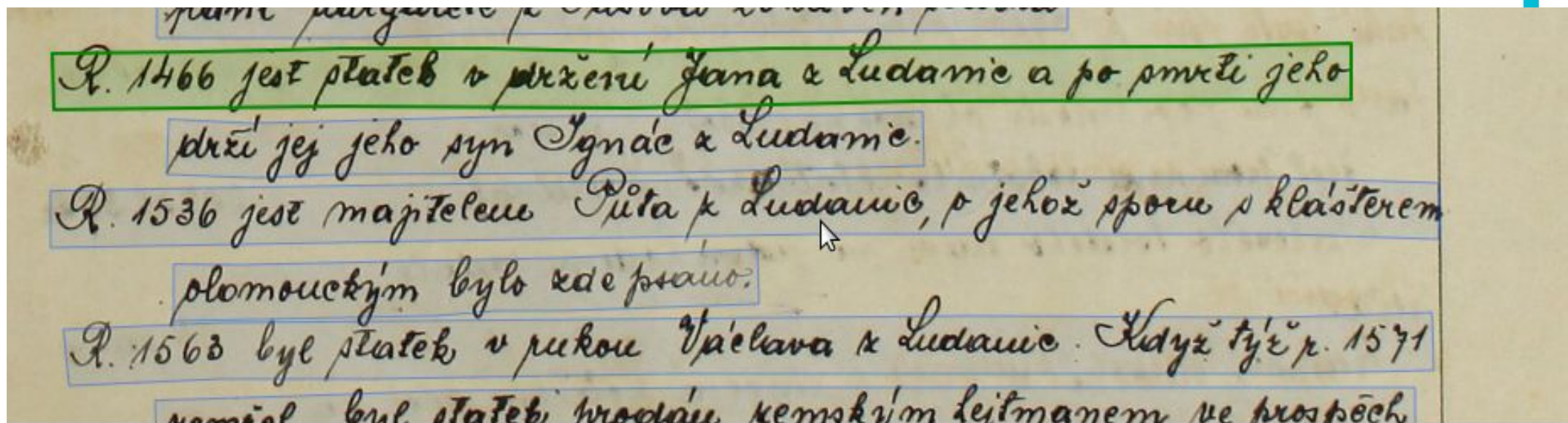
Moravoleň v horním závodě uvolnil této škole 2 nyní nepoužívané místnosti, které budou upraveny na dílnu pro práci s dřevem a s kovy.

ani projednávání komplexního plánu rozvoje školství v 3. pětiletce za účasti zástupců MNV a místních závodů však nevyřešilo ožehavý problém, to je zavedení 2 směnného vyučování. V příštích letech přibude na škole asi 100 dětí t. j. 3 třídy, pro které není učeben. Příslušné komise ONV a

Moravoleň v horním závodě uvolnil této škole 2 nyní nepoužívané místnosti, které budou upraveny na dílnu pro práci s dřevem a s kovy

ani projednávání komplexního plánu rozvoje školství v 3. pětiletce za účasti zástupců MNV a místních závodů však nevyřešilo ožehavý problém, to je zavedení 2 směnného vyučování. V příštích letech přibude na škole asi 100 dětí t. j. 3 třídy, pro které není učeben. Příslušné komise ONV a

## 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

# Rukopisy

195  
hudbného a slavného jako Innocenc III., jenž řídil osudy téměř celé Evropy. Přední jeho starostí bylo způsobiti novou výpravu křížovou na osvobození Palaestiny; doba byla příhodná, neboť Saladin právě byl zemřel. Proto rozeslal posly po veškerém světě křesťanském, by hlásali kříž a vybírali peníze na novou výpravu. Ale poměry v Evropě nebyly utěšeny, neboť

hudbného a slavného jako Innocenc III., jenž řídil osudy téměř celé Evropy. Přední jeho starostí bylo způsobiti novou výpravu křížovou na osvobození Palaestiny,

doba byla příhodná, neboť Saladin právě byl zemřel.

Proto rozeslal posly po veškerém světě křesťanském, by hlásali kříž a vybírali peníze na novou výpravu. Ale poměry v Evropě nebyly utěšené, neboť

# Rukopisy

7. 2

Reverendissime, Amplissime  
ac Graciosissime Domine  
Abbas!

Quum id munus mihi impositum caperet ratio, ut  
quos fecerim profectus infirmum, natum facis, me  
in subeundo et Religionis materiae examine, solum  
modo circumstantius, me licet involam, ad quoddam  
tempus detineret. Absolventur enim omni die Jovis  
qui rigoreo examine unicus destinari posset et debet  
concursum designate, quare ea deducor, ut exprobram  
donec in arenam descensus vocor. Operam omnem

impendit Dominus Professor  
concursum ea Historia Naturali,  
matica absolutis, trahatur).  
currentibus tempus reliquum

Quod studia concernit currenti  
eruditione ac fervore litteris  
Inter hos palmam fert Dom.  
Theologia, quem in cathedra  
In Oeconomia Dominus Prof.  
modo huius reddidit celeberrimam. Huius  
trahens pedagogicam et catech.

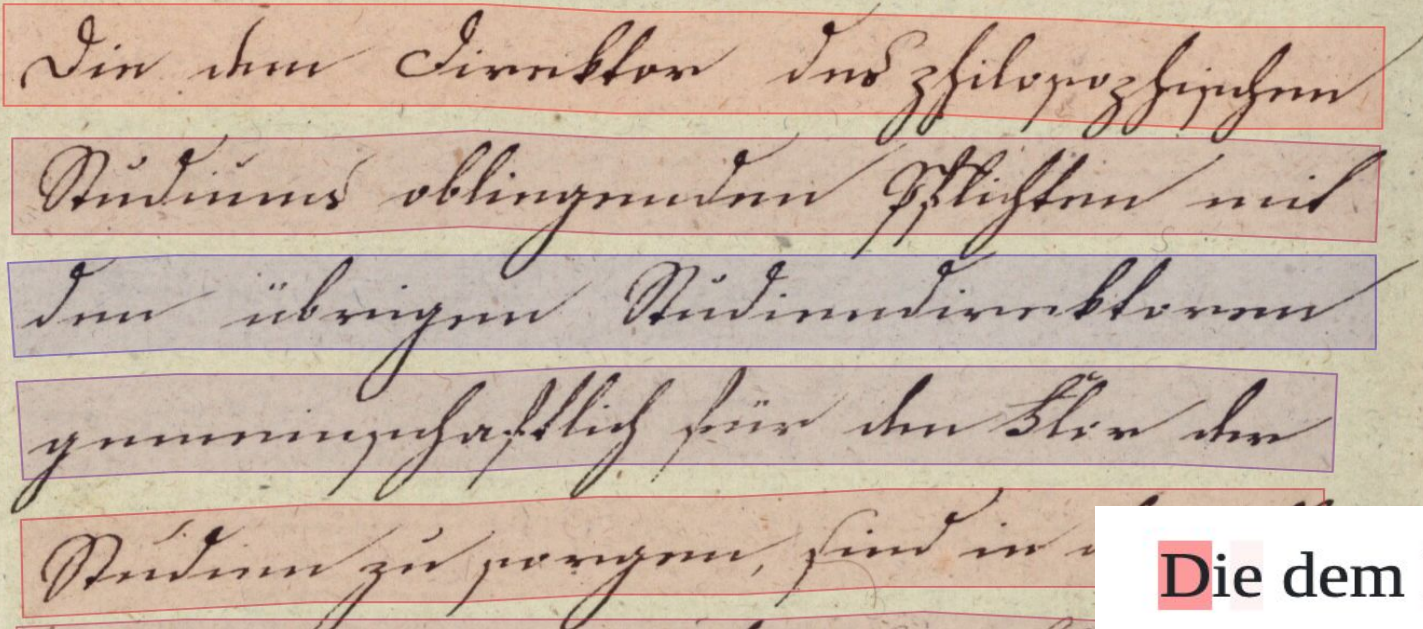
impendit Dominus Professor Löcher ne mora amplius,  
concursum ex Historia Naturali, Diplomatica, et Nemi-  
matica absolutis, trahatur. Interea omne a studiis  
currentibus tempus reliquum materia revolvenda conse-  
Quod, studia concernit currentia, nactus sum duces ingenio

impendit Dominus Professor Löcher ne mora amplius,  
concursum ea Historia Naturali, Diplomatica, et Nemi-  
matica absolutis, trahatur). Interea omne a studiis  
currentibus tempus reliquum materia revolvenda conse-  
Quod studia concernit currentia, nactus sum duces ingenio

eruditione ac fervore litteris rite tractandi



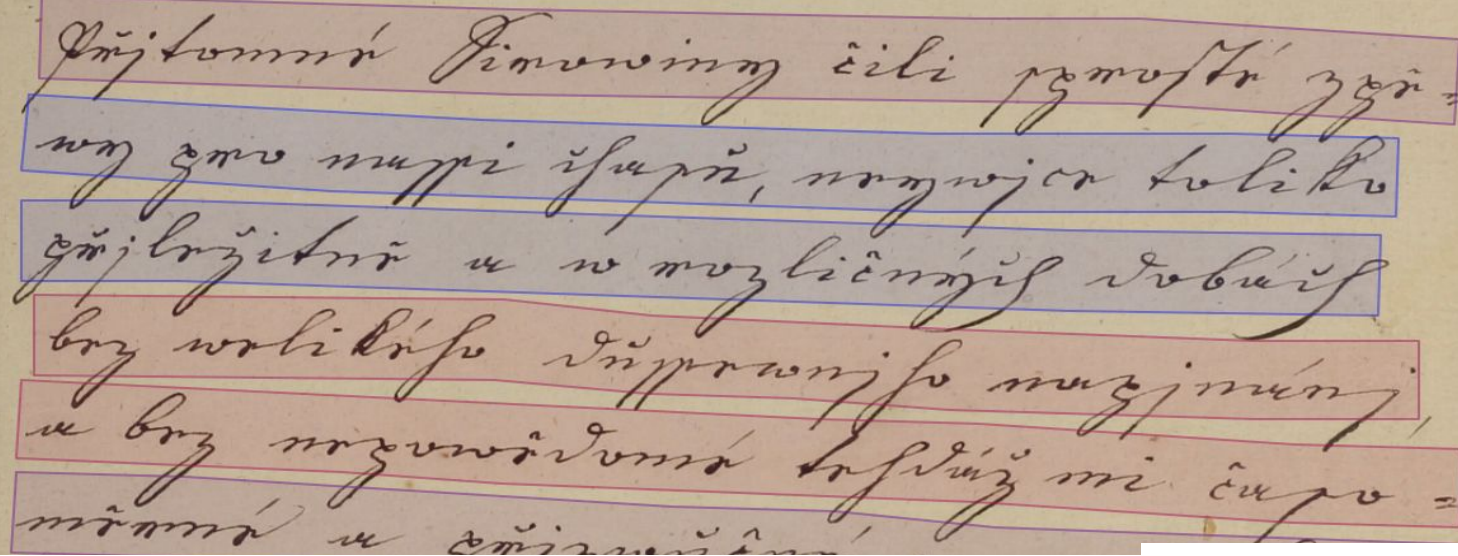
# Kurent



Die dem Direktor, des philosophischen  
Studiums obliegenden Pflichten mit  
den übrigen Studiendirektoren  
gemeinschaftlich für den Fler der  
Studien zu sorgen, sind in der all,

Die dem Direktor, des philosophischen  
Studiums obliegenden Pflichten mit  
den übrigen Studiendirektoren  
gemeinschaftlich für den Fler der  
Studien zu sorgen, sind in der all,

# Kurent



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

Přjtomné Sirowiny čili sprofté zpě,  
wy pro naši chasu, neywjce toliko  
přjležitně a w rozličných dobách  
bez welikého dušewnjho 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](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/>



# 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.**  
Již jen několik dní.

9403

ý kapitál K 13,000,000--  
Telefon 4  
šnéna Rozumné ber

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: Přírodně-tykářství 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](http://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

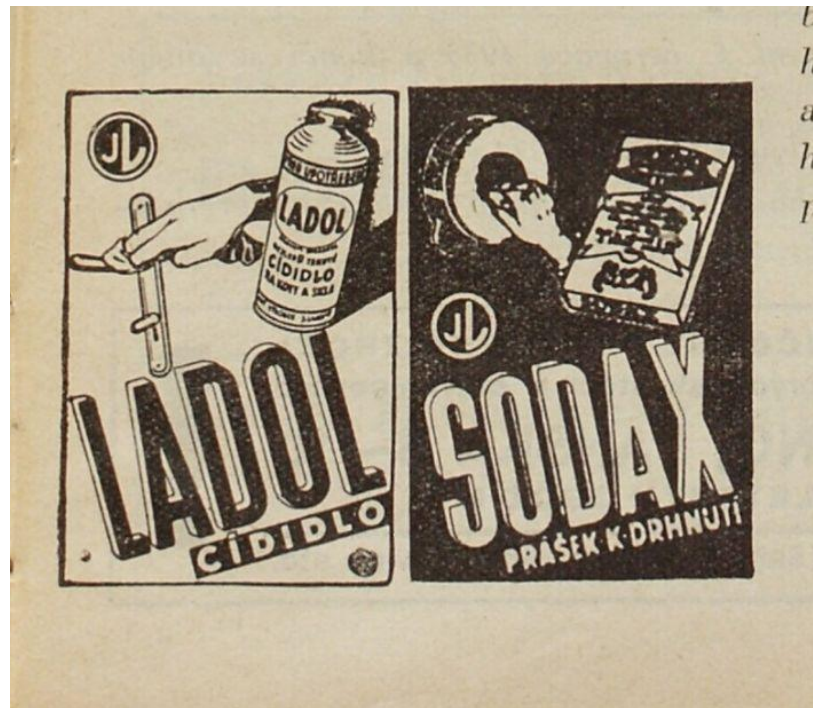
**RECENZE**  
Povedená D  
vě dobrým  
litním vide  
i citelný ne  
sionálního.

HN Exkluz



# Co obsahuje text stále může být obrázek

- text v obrázcích, překryv

An advertisement for Oriflame Beauty mascara. It features a large image of a mascara wand and tube. The text includes a paragraph about the benefits of waterproof mascara, a price list for Nivea and Oriflame products, and a 'ZKUSTE' (Try) section.

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 nejhorším ranním sklužu. Složení, odolávající dešti, vám dodá jistotu, že se ani po srážce s podzimní plísňovostí neproměníte ve smutnou pandu.

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

**NIVEA**  
KREMLIN  
SANTALIN  
ALZHA MANGEL  
ENTENSER  
104 Kč

**Oriflame**  
BEAUTY  
WATERPROOF  
MASCARA

(ten) / Foto archiv Iitem

# 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 vládnem. 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.

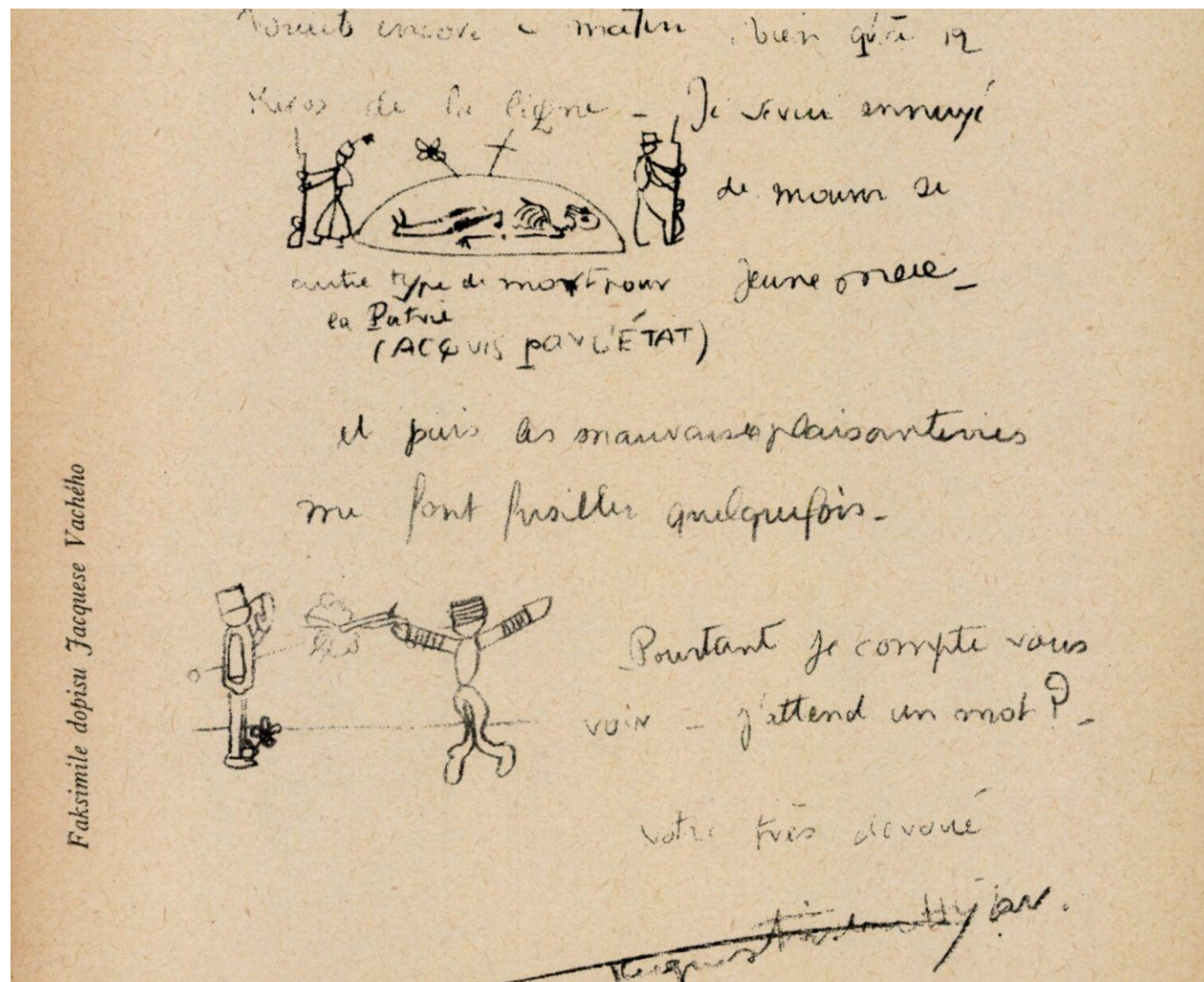
Nyní 50% sleva na 2 měsíce při aktivaci do konce dubna.

Aktivace a více informací na 800 777 777 nebo na [www.vodafone.cz](http://www.vodafone.cz)

Teď je to ve vašich rukou.

**50% SLEVA PŘI AKTIVACI DO KONCE DUBNA**

# Obrázek a text mohou být vizuálně podobné





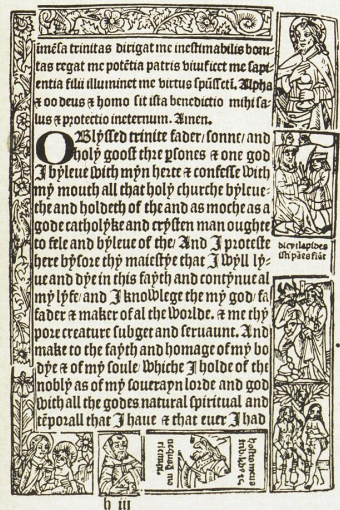
# Kategorie nejsou jasné ani lidem

- 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 Christiny 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. Wynuknem de Worde ve Westminsteru.

pende nec. Quid tep e alle other man Quid tho fo modis penat op to the fpe et alle my good e If I hade not god ano to alle; thepe fialo no t as foeth the holi the renuice that fere good hret geyt fepth caitus fo to an four maymore f Dpdy hlyp buce ga fo fo vrlamp



## 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 mec lituje. Protože tady bydlí na vesnici, umí si představit, kolik práce a starostí se vším mělo. Jedinec řešeni by asi bylo najmout si nějakého právníka, který by uměl poradit, jestli máš právo, nebo ne. Ale vzta- h se sourozenci se už asi nikdy nenapraví.

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

ROZPOČETÍ

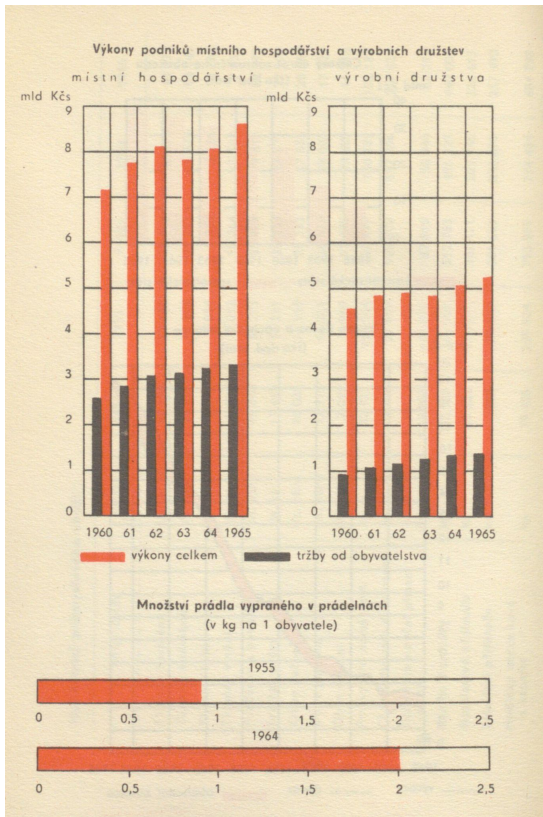
ROZPOČETÍ

ROZPOČETÍ

Náš rodina nikdy nebyla vzorem pospolitosti, kde by všichni nezájmně pomáhali a byli si oporou. Nevím, jak se to stalo, asi to byla nějaká výchovná chyba rodičů nebo ep, ale zkrátka já, sestra a bratr nejsme žádní dokonalí sourozenci. 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. Nejhorší situace nastala po smrti našich rodičů. Táta umřel už dávno, ale maminka ani ne před půlrokem. Pořád je v tom, že my jsme byli s mužem a dětmi

a ními v jejich původním domě. Ne že bych o to nějak zvlášť stála, 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 spoustu 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 přejívat 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, přišel dům nám. Proto taky napláta živéč, ve které to píše. Potom umřela a začalo peklo. Sestra si napel- non vzpomněla, kolik asi takový dům stojí a že my nemáme prvo si na něj dělat nároky. 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 si na samozřejmě souhlasil, přestože nevěděl do oprav ani korunu. A že by se třeba staral jeden nebo druhý o mědu, když už nemohla, chudák, ani sama dojit na záchod, o tom ani nemluhám.

Jsem nešťastná. Prodat dům, to si neumím představit, protože kdy bychom bydlili? A kdy bychom vřali peníze na to, abych osourouzene vypla- tila, taky nevím, protože všechny peníze, co jsme měli navíc, šly do materiálu a na dělníky. Sestra říká, že jestli okamžitě nenavrhnu nějaké řešení, dá vše k soudu. Jestli jsem se ani nemesadili zjistit, jak to právně všechno je. Vím jenom, že mě mrzí ne- jenom to, jak se teď sourozenci handrkují o něco, na čem jim nikdy nezaleželo, ale hlavně to, že se o maminku pořádně nezajímali, když byla ještě na světě. Podle mě se musí v hrůbě olouacet. Manžel mě uklidňuje, že vše nějak dopadne, ale samozřej- mě nemá prý švagrovu dobré slovo. Tuhle, když jsem plakala, se rozčílil 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 rozhádati.



# 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

# Detekce pomocí vlastního modelu strojového učení

Řešení: Segmentace pomocí konvoluční neuronové sítě

## Proč konvoluční sítě?

- Translační invariance: poloha obrázku na stránce není důležitá
- Množství volně dostupných implementací state of the art modelů (napr. ResNet, AlexNet, VGG,...)
- Možnost použít váhy předtrénované na velkém datasetu.



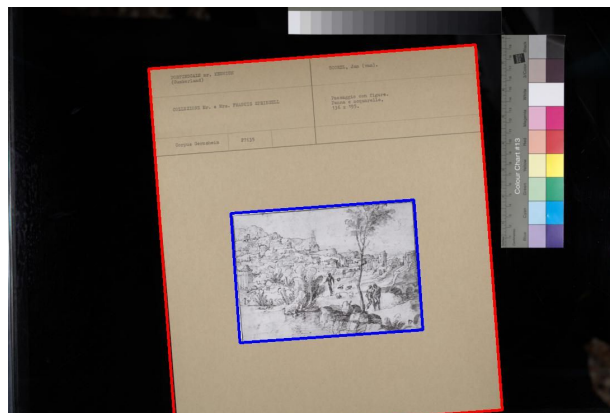
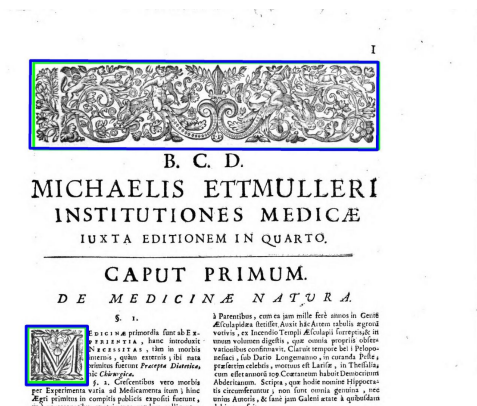
# Slepé uličky

- **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).



# Neuronová síť dhSegment

- Vyvinuta v 2018 specificky k zpracování skenovaných historických dokumentů
- Použita k detekci ornamentů a fotek na naskenovaných stránkách
- Inspirována segmentační sítí UNet
- Možnost využít předtrénované váhy z ResNet-50
- Implementována v Pythonu pomocí Tensorflow



# Neuronová síť dhSegment

- Natrénována na ~1000 manuálně anotovaných stránkách na notebooku bez GPU

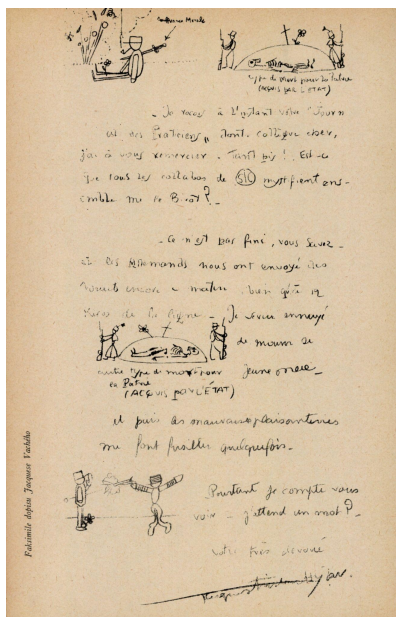
zpracování sítí



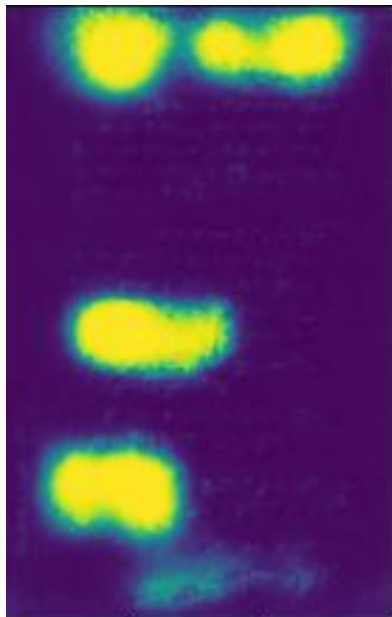
binarizace



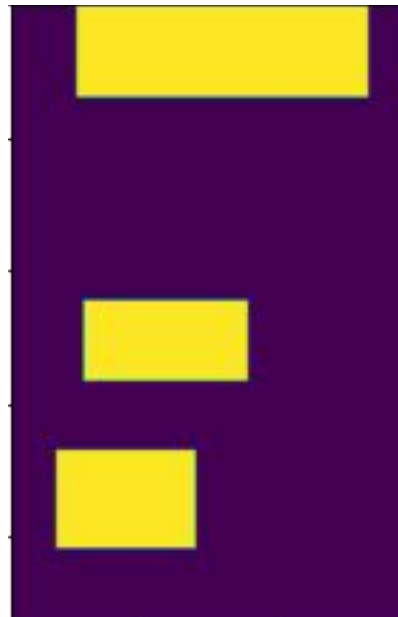
výstup



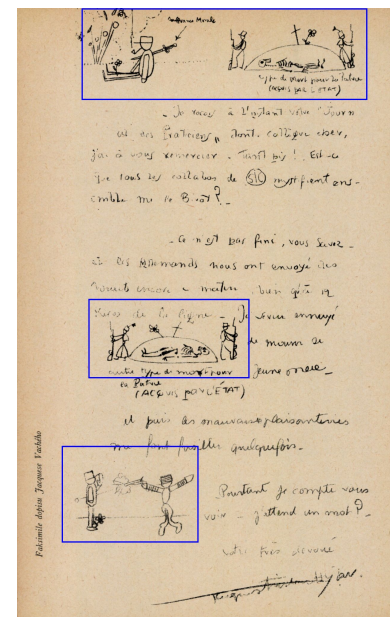
vstup



pravděpodobnost  
pro každý pixel



binární maska  
pro každý pixel



poloha obrázků



# Neuronová síť dhSegment

- Výsledek:
  - Celkové IOU: 0.65
  - Senzitivity\*: 0.65
  - Precission\*: 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**

\*Vypočteno pro práh citlivosti IOU=0.5

\*\*Neobsahují obrázek, nebo obsahují obrázek který už byl detekován



# Odkazy

- dhSegment
  - Github <https://github.com/dhlab-epfl/dhSegment>
  - Publikace: Oliveira, Sofia Ares, Benoit Seguin, and Frederic Kaplan. "dhSegment: A generic deep-learning approach for document segmentation." *2018 16th International Conference on Frontiers in Handwriting Recognition (ICFHR)*. IEEE, 2018, preprint at <https://arxiv.org/abs/1804.10371>
- NewspaperNavigator projekt: <https://github.com/LibraryOfCongress/newspaper-navigator>
- OCR od ABBYY: <https://www.abbyy.com/finereader-server/>
- Resnet: <https://arxiv.org/pdf/1512.03385.pdf>
- UNet: <https://arxiv.org/pdf/1505.04597.pdf>





# 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





Hledat

# Porovnání obrázků inzerce deníku Svobodné slovo

### ARMIN

Je to nejlepší a nejlevnější domácí výrobci továrna mýdel

**Jan Hubínek a syn, Praha-VIII. Libeň.**

Zkuste ARMIN ke praní neb mytí a uvidíte, že jeho útržkem.

V dvojnásobné ceně kus 6 krejcarů.

Pánům obchodníkům cenitky a vzorky zdarma vyžádání.

### BOH. KREJČÍK,

velkoobchod ublín, libeň - skid - písárna

Velkoobchod ublín, libeň - skid - písárna

Velkoobchod ublín, libeň - skid - písárna

### Romány z českých dějin od Jos. Svátka

romány z českých dějin od Jos. Svátka

romány z českých dějin od Jos. Svátka

### AMERIKO PULTY

americké skříně, skříně, skříně

americké skříně, skříně, skříně

americké skříně, skříně, skříně

### Vazby skvostné

vazby skvostné, vazby skvostné

vazby skvostné, vazby skvostné

vazby skvostné, vazby skvostné

### Dr. Frant. Bačkovský,

lékař v Praze, třída st. číslo 36.

lékař v Praze, třída st. číslo 36.

lékař v Praze, třída st. číslo 36.

### Výbavy prádla pro nevěsty.

výbavy prádla pro nevěsty, výbavy prádla pro nevěsty

výbavy prádla pro nevěsty, výbavy prádla pro nevěsty

výbavy prádla pro nevěsty, výbavy prádla pro nevěsty

veškeré - GISEVOE - POTREBY

veškeré - GISEVOE - POTREBY

veškeré - GISEVOE - POTREBY

### Bratři! První výrobní družstvo dělnictva krejčovského

Bratři! První výrobní družstvo dělnictva krejčovského

Bratři! První výrobní družstvo dělnictva krejčovského

### Jarní mody

Jarní mody, jarní mody, jarní mody

Jarní mody, jarní mody, jarní mody

Jarní mody, jarní mody, jarní mody

### Malý seznamovatel.

Malý seznamovatel, Malý seznamovatel

Malý seznamovatel, Malý seznamovatel

Malý seznamovatel, Malý seznamovatel

### Hynek Gotwald

Hynek Gotwald, Hynek Gotwald

Hynek Gotwald, Hynek Gotwald

Hynek Gotwald, Hynek Gotwald

### AMERIKO PULTY

americké skříně, skříně, skříně

americké skříně, skříně, skříně

americké skříně, skříně, skříně

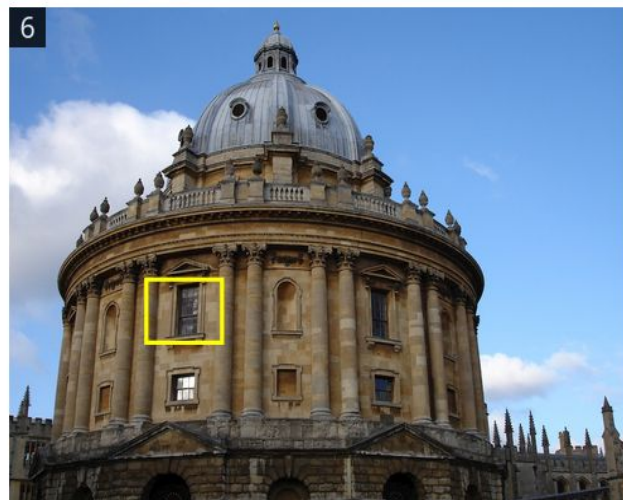
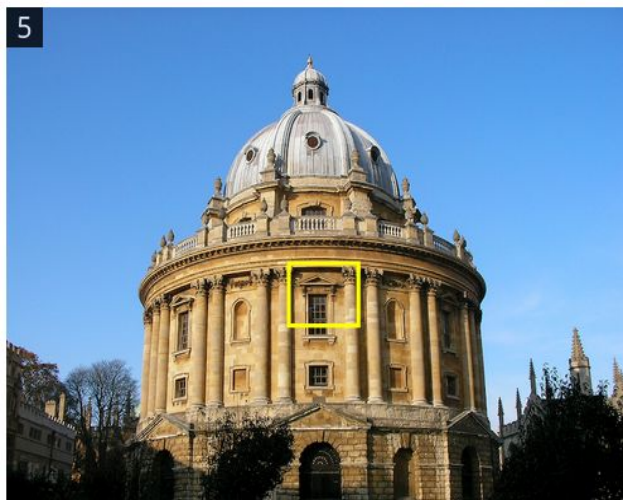
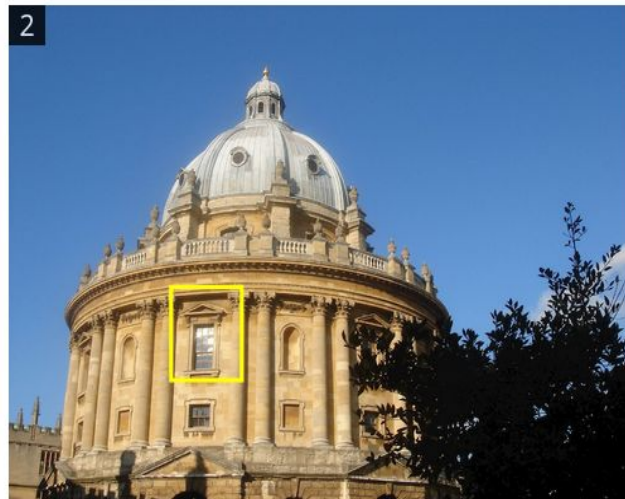
### AMERIKO PULTY

americké skříně, skříně, skříně

americké skříně, skříně, skříně

americké skříně, skříně, skříně

# Funkční i pro 3D prostorové obrázky



# 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 použít JPEG2000



# 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/>



# Co dál?

- 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, Ján Bogár, Michal Tran - Moravská zemská knihovna v Brně