

MUNI  
FI



# Towards General Document Understanding through QA

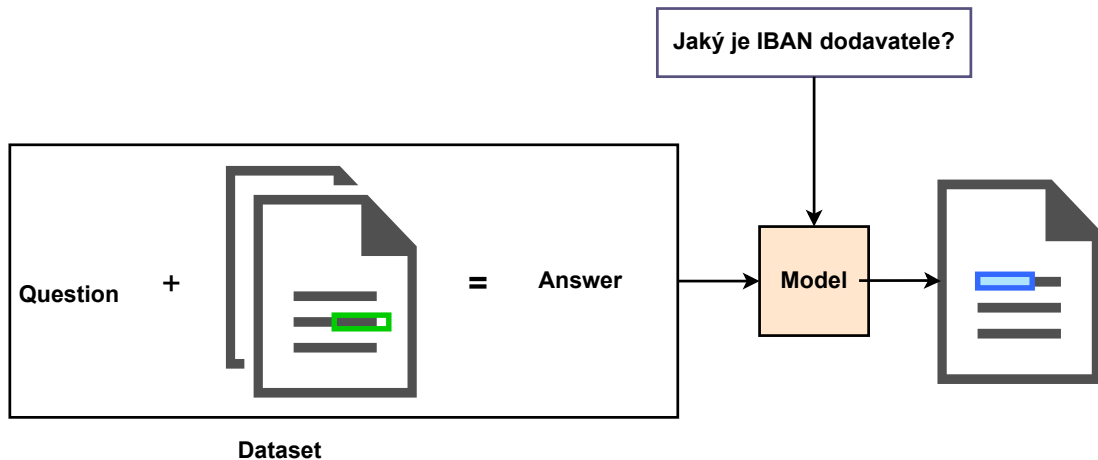
**Šárka Ščavnická, Michal Štefánik, Marek Kadlčík,  
Martin Geletka, Petr Sojka**  
**527352@mail.muni.cz, sojka@fi.muni.cz**

Math Information Retrieval Research Group,  
Faculty of Informatics, Masaryk University

December 9th, 2022

# Motivation

# Model Description



# Competitive system for the English language

## 1. Select a file

Clear

Gallery

## INVOICE

**East Repair Inc.**  
1912 Harvest Lane  
New York, NY 12210

**BILL TO**  
John Smith  
2 Court Square  
New York, NY 12210

LOGO

**SHIP TO**  
John Smith  
3787 Franklin Drive  
Cambridge, MA 12210

**INVOICE #** us-001  
**INVOICE DATE** 11/02/2019  
**P.O.#** 2312/2019  
**DUE DATE** 20/02/2019

QTY	DESCRIPTION	UNIT PRICE	AMOUNT
1	Front and rear brake cables	100.00	100.00
2	New set of pedal arms	15.00	30.00
3	Labor 3hrs	5.00	15.00
		Subtotal	145.00
		Sales Tax 6.25%	9.06
		<b>TOTAL</b>	<b>\$154.06</b>

John Smith

Thank you

**TERMS & CONDITIONS**

Payment is due within 15 days  
Please make checks payable to: East Repair Inc.

## 2. Ask a question

Question

Model

 LayoutLMv1 🐱  LayoutLMv1 for Invoices 🐼  Donut 🍩

Clear

Submit

Top Answer

us-001

## Textual Representation of the Invoice

- FAKTURA - DAŇOVÝ DOKLAD © číslo: (25 Dodavatel | IČO 42993865 DIČ [CZ6804111963] Variabilní symbol 25 Josef Sedlář Konstantní symbol [0308 Josef Sedlář Specifický symbol , Podolí 88 Částka 78438.00 Kč 789 85 Mohelnice Peněžní ústav KB Mohelnice Objednávka Číslo účtu (100 Číslo fixmy 00024 Příjemce Odběratel IČO Drč Kamil Sedlář Kamil Sedlář Staškova 55 Staškova 55 789 85 MOHELNICE 789 85. MOHELNICE Platební Den splatnosti 24.10.2009
  - Who is the supplier, and what is his address?

## Textual Representation of the Invoice

- FAKTURA - DAŇOVÝ DOKLAD © číslo: (25 **Dodavatel** | IČO 42993865 DIČ [CZ6804111963] Variabilní symbol 25 Josef Sedlář Konstantní symbol [0308 Josef Sedlář Specifický symbol , Podolí 88 Částka 78438.00 Kč 789 85 Mohelnice Peněžní ústav KB Mohelnice Objednávka Číslo účtu (100 Číslo fixmy 00024 Příjemce Odběratel IČO Drč Kamil Sedlář Kamil Sedlář Staškova 55 Staškova 55 789 85 MOHELNICE 789 85. MOHELNICE Platební Den splatnosti 24.10.2009
  - Who is the supplier, and what is his address?
- We need more context

# Visual Representation of the Invoice

FAKTURA - DAŇOVÝ DOKLAD číslo: 25 9

Dodávateľ IČO 42993865 DIČ CZ6804111963 Josef Sedlář Josef Sedlář Podolí 88 789 85 Mohelnice Peněžní ústav KB Mohelnice Číslo účtu 43-3975450267/0100	Variabilní symbol 25 Konstantní symbol 0308 Specifický symbol Částka =78438.00 Kč Objednávka Číslo firmy 00024																				
Příjemce Kamil Sedlář Staňkova 55 789 85 MOHELNICE	Odběratel IČO DIČ Kamil Sedlář Staňkova 55 789 85 MOHELNICE																				
Platební podmínky Den splatnosti <b>24.10.2009</b> Způsob úhrady v hotovosti Datum vystavení dokladu 14.10.2009 Datum uskutečnění zdanitelného plnění 14.10.2009																					
Položka (ceny v Kč bez daně) JednCena Množství CelkCena DPH																					
Fakturujeme Vám stavební práce																					
<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="width: 50%;">stavební práce</td> <td style="width: 10%;">150.00</td> <td style="width: 10%;">24.000 hrd</td> <td style="width: 10%;">3600.00</td> <td style="width: 10%;">9%</td> </tr> <tr> <td>materiál f. 2913000010</td> <td>61374.00</td> <td>1.000 Kč</td> <td>61374.00</td> <td>9%</td> </tr> <tr> <td>materiál f. 3001537373</td> <td>6987.39</td> <td>1.000 Kč</td> <td>6987.39</td> <td>9%</td> </tr> <tr> <td>celkem .....</td> <td></td> <td></td> <td>71961.39</td> <td>Kč</td> </tr> </table>		stavební práce	150.00	24.000 hrd	3600.00	9%	materiál f. 2913000010	61374.00	1.000 Kč	61374.00	9%	materiál f. 3001537373	6987.39	1.000 Kč	6987.39	9%	celkem .....			71961.39	Kč
stavební práce	150.00	24.000 hrd	3600.00	9%																	
materiál f. 2913000010	61374.00	1.000 Kč	61374.00	9%																	
materiál f. 3001537373	6987.39	1.000 Kč	6987.39	9%																	
celkem .....			71961.39	Kč																	
<table border="1" style="margin: auto; border-collapse: collapse;"> <tr> <td style="padding: 2px;">sazba</td> <td style="padding: 2px;">bez daně</td> <td style="padding: 2px;">DPH</td> <td style="padding: 2px;">s daní</td> </tr> <tr> <td style="padding: 2px;">snižovaná</td> <td style="padding: 2px;">71961.39</td> <td style="padding: 2px;">6476.53</td> <td style="padding: 2px;">78437.92</td> </tr> </table>		sazba	bez daně	DPH	s daní	snižovaná	71961.39	6476.53	78437.92												
sazba	bez daně	DPH	s daní																		
snižovaná	71961.39	6476.53	78437.92																		
k úhradě ..... (zakrouhlení +0.08) 78438.00 Kč *Živn. list 11.12.2008, č.j.:ZJV/642/2008/ÚJ/11																					
Počet stran 1 Vystavil Kateřina Pačková	<b>Josef SEDLÁŘ</b> 789 85 Mohelnice, Podolí 88 292 02 Lottice, U Mlýna 97																				

## Text Representation of the Invoice with the Visual Knowledge

- Who is the supplier, and what is his address?
- FAKTURA - DAŇOVÝ DOKLAD © číslo: (25 **Dodavatel** | IČO 42993865 DIČ [CZ6804111963] Variabilní symbol 25 **Josef Sedlář** Konstantní symbol [0308 **Josef Sedlář** Specifický symbol , **Podolí 88** Částka 78438.00 Kč **789 85 Mohelnice** Peněžní ústav KB Mohelnice Objednávka Číslo účtu (100 Číslo fixmy 00024 Příjemce Odběratel IČO Drč Kamil Sedlář Kamil Sedlář Staškova 55 Staškova 55 789 85 MOHELNICE 789 85. MOHELNICE Platební Den splatnosti 24.10.2009
- Better results if we use the visual layout of the document



# Methods and Datasets

# Question Answering

- Two types:
  - Extractive QA
  - Generative QA
- Multilingual datasets, primarily in English
- Datasets:
  - Stanford Question Answering Dataset (SQuAD)
    - Wikipedia articles
  - French Question Answering Dataset (FQuAD)
    - Wikipedia articles

# Visual Question Answering

- Focusing only on the visual part of an image
- We can use the same images for different languages
- Datasets:
  - Visual Question Answering (VQA) dataset
  - The Image-Set Visual Question Answering (ISVQA)

# Named Entity Recognition

- Aims to locate and identify entities in the given text
- Datasets:
  - CoNLL-2003 dataset
    - English and German
    - Newspaper articles
  - Few-NERD
    - Wikipedia articles
  - Dataset of German Legal Documents for NER
    - Non-English
    - Legal documents

# Document Visual Question Answering

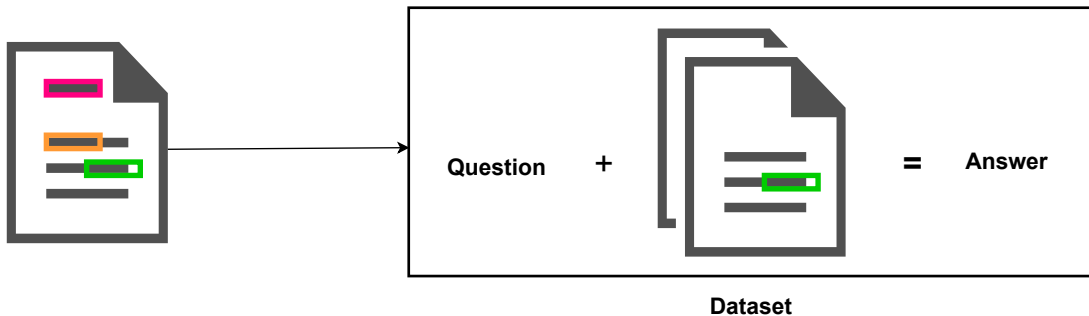
- Combines the visual and textual sides of a document
- Lacking coverage of non-English models for Document Visual Question Answering
- Datasets:
  - DocVQA
    - From the UCSF Industry Documents Library
    - Contains invoices
    - Documents are from the period between 1960 and 2000
  - Visual Machine Reading Comprehension (VisualMRC)
    - Consisting of screenshots of web pages

# Dataset creation

# Invoices - 6,849 documents and 15 entities

FAKTURA - DAŇOVÝ DOKLAD		číslo: 25			
Dodavatel IČO 42993865 DIČ CZ6804111963 Josef Sedlář Josef Sedlář Podolí 88 789 85 Mohelnice Peněžní ústav KB Mohelnice Číslo účtu 43-3975450267/0100		Variabilní symbol 25 Konstantní symbol 0308 Specifický symbol Číska +78438.00 Kč Objednávka Číslo firmy 00024			
Příjemce Kamil Sedlář Staškova 55 789 85 MOHELNICE	Odběratel IČO DIČ Kamil Sedlář Staškova 55 789 85 MOHELNICE				
		Platební podmínky Den splatnosti 24.10.2009 Způsob úhrady v hotovosti Datum vystavení dokladu 14.10.2009 Datum uskutečnění zdanitelného plnění 14.10.2009			
Položka	(ceny v Kč bez daně)	JednCena	Množství	CelkCena	DPH
Fakturuje Vám stavební práce					
stavební práce	150.00	24.000 hod	3600.00	9%	
materiál ř. 2913000010	61374.00	1.000 Kč	61374.00	9%	
materiál f. 3001537373	6987.39	1.000 Kč	6987.39	9%	
celkem .....			71961.39 Kč		
	sazba	bez daně	DPH	s daní	9%
	smíšená	71961.39	6476.53	78437.92	
k úhradě .....	(zakrouhlení +0.08)			78438.00 Kč	
*živn. list 11.12.2008, č.j.:ZV/642/2008/UJ/11					
Počet stran 1	Josef SEDLÁŘ 789 85 Mohelnice, Podolí, 88 ..... 11. srpna 07				

# Dataset Description





# Research Proposal

**How can Textual QA and Visual QA datasets  
improve Document Understanding?**

## First Question

- Our model will be trained on the Czech Document VQA dataset
  - A large number of relevant but narrow questions
- Hypothesis:
  - If we first introduce our model to much more general and comprehensive QA datasets
    - More robust in the unseen evaluation scenarios
    - Perform better

**How well can Document VQA models generalize beyond training entity types?**

## Second Question

- Conventional Named Entity Recognition models
  - A closed set of entity types present in the training set
  - Must be retrained when a new entity needs to be recognized
- Document VQA
  - A closed set of questions
  - Retraining?
    - How well our model can generalize to unseen entities?

**How much can Document VQA in non-English languages benefit from English datasets?**

## Third Question

- Compare the model performance on a target language
- Train on:
  - Only English data
  - A mixture of English and the target-language data
  - Only the target-language data
- The evaluation of the approach will assess the applicability of our models to unseen language

**MUNI**

FACULTY

OF INFORMATICS