

Provozně  
ekonomická  
fakulta

Diplomová práce

Autor práce:

Bc. Jiří Švec

Vedoucí práce:

Ing. Jan Přichystal, Ph.D.

# Aplikace pro získávání názorů z uživatelských recenzí



## Cíl práce

- Webová aplikace
- Automatizované stahování a zpracování recenzí
- Rozšiřitelná o další datové zdroje
- Graficky přehledná a konfigurovatelná
- Uživatel získá přehled o probíraných tématech

## Současný stav

- Neexistuje konkurenční aplikace pro hromadné zpracování českých uživatelských recenzí
- Nutnost manuálního vyhledávání a vyhodnocení
- Časově náročný proces

# Požadavky

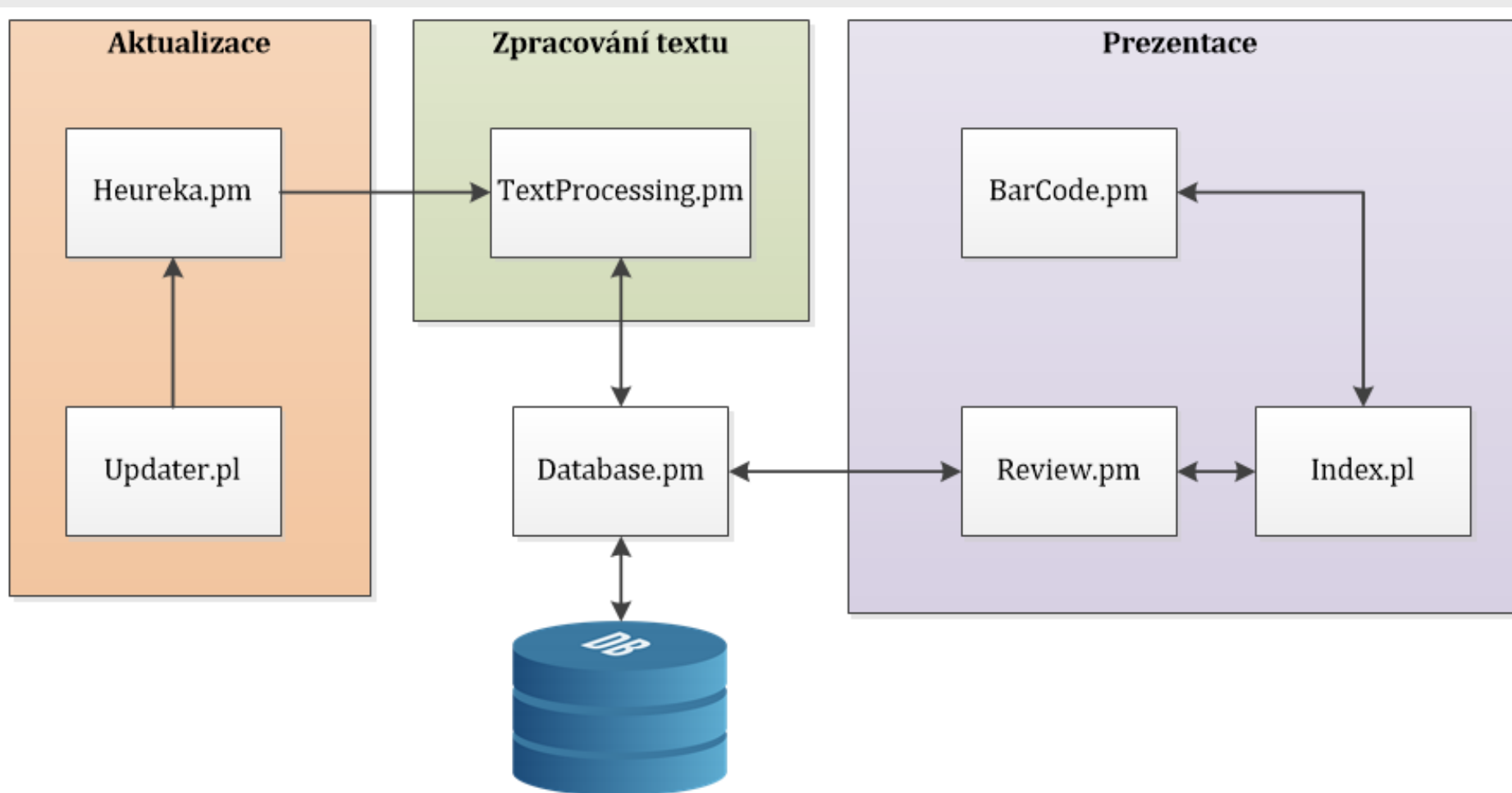
## Funkční

- Stahovat recenze
- Zpracovávat stažené recenze
- Ukládat výsledky do databáze
- Nahrávat snímky čárových kódu prostřednictvím kamery mobilního telefonu
- Zobrazovat zpracované recenze

## Nefunkční

- Webová aplikace
- Přívětivý vzhled
- Jednoduchá manipulace s aplikací
- Snadná integrace dalších zdrojů recenzí
- GUI v češtině
- Konfigurovatelnost
- Rychlá odezva

# Architektura



# Konfigurace

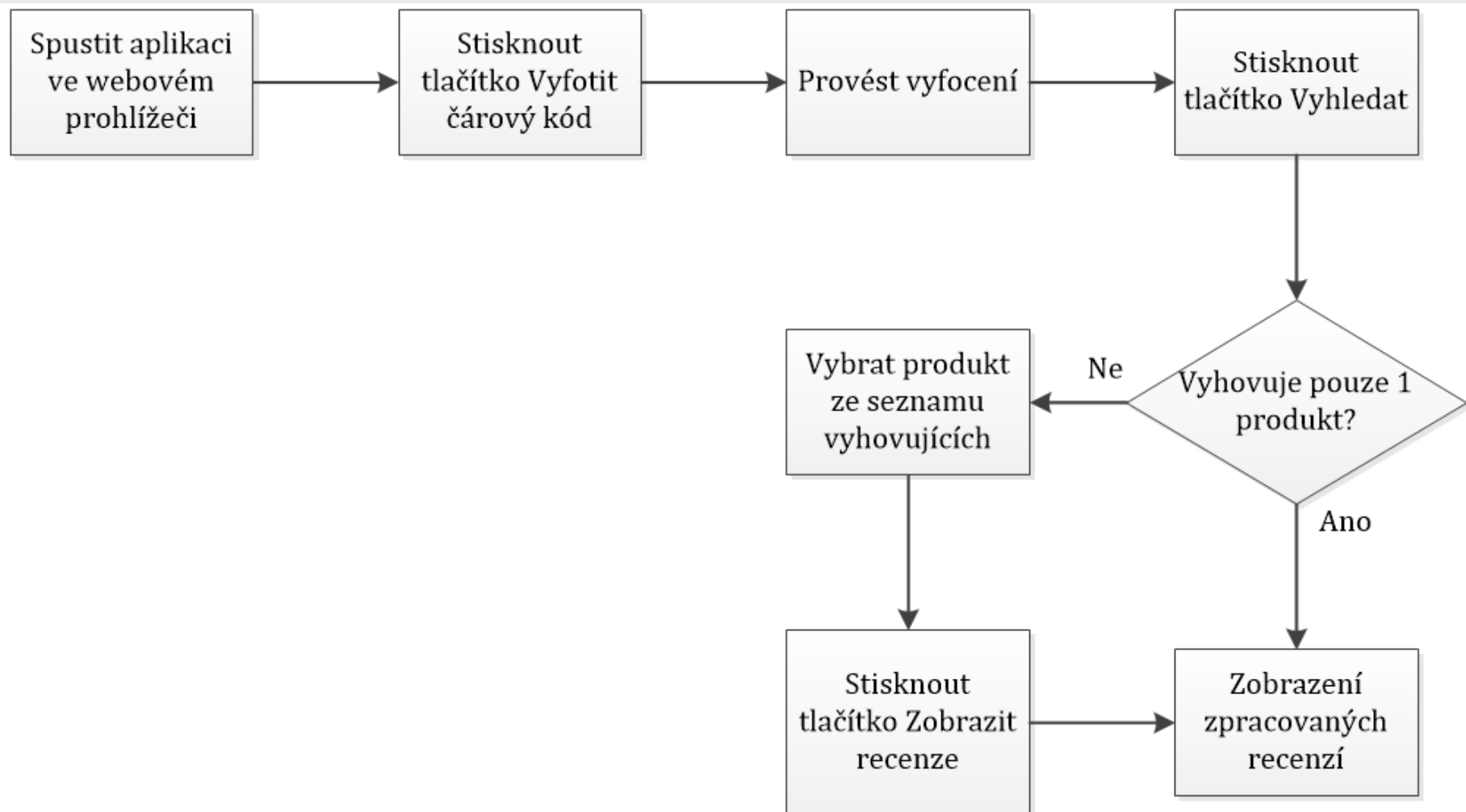
- Nastavení výkonu
- Nastavení aktualizací
- Párování názvů
- Připojení k databázi
- Cesty ke slovníkům
- Nastavení word cloudů
- Nastavení a optimalizace vendora

```
1 [GLOBAL]
2 DEBUG=0
3 THREAD_COUNT=8
4 THREAD_SLEEP_INT=3
5 PRODUCTS_PER_FILE_LIMIT=3000
6 PRODUCT_NAME_PERCENTAGE_MATCH=70
7 WORD_CLOUD_MIN_WORD_FREQUENCY=2
8 SHOW_WORD_CLOUD=1
9 SHOW_WORD_CLOUD_NEIGHBORS=1
10 UPDATE_EVERY=30
11 DELETE_ALL_ON_UPDATE=1
12
13 # downloaded from https://sites.google.com/site/kevinl
14 STOPWORDS_PATH=conf/stopwords_cz.txt
15
16 # downloaded from http://www.lexiconista.com/datasets
17 LEMMATIZATION_PATH=conf/lemmatization_cs.txt
18
19 # path to zbarimg tool, empty if installed on linux
20 ZBARIMG_HOME=
21
22
23 #-----
24 [DATABASE]
25 HOST=localhost
26 PORT=3306
27 USER=root
28 PASSWORD=root
29 SCHEMA=dp
30
31 #-----
32 [HEUREKA]
33 URL=http://www.heureka.cz
34 POPULATE_PRODUCTS_LIST=0
35 POPULATE_PRODUCTS=0
36 POPULATE_REVIEWS=0
37 GET_ERRORS_LIMIT=5
38 GET_ERRORS_SLEEP=5
39 MAX_CATEGORY_PAGE=15000
```

# Použité metody a nástroje

- Web scraping
- Synchronizační vzor fronta
- Normalizace
- Stopwords
- Tokenizace
- Stemizace
- Lemmatizace
- Shlukování
- Sousední dvojice
- Asociační pravidla
- Levenshteinova vzdálenost
- Responzivní design

# Prezentace názorů 1 / 2







# Klady a zápory

## Klady

- Autonomnost
- Škálovatelnost
- Konfigurovatelnost
- Jedinečnost
- Rozšiřitelnost
- Inkrementální stahování produktů
- Paralelní zpracování
- Responzivní GUI

## Zápory

- Zdlouhavý proces pořízení a zpracování dat
- Závislost na slovnících
- Nelze zajistit párování názvů produktu

**DĚKUJI ZA POZORNOST**