

Rozšíření systému pro získávání, zpracování a analýzu rozsáhlých kolekcí textů z webu

Matějka Jiří

30. 1. 2018

Co umíme:

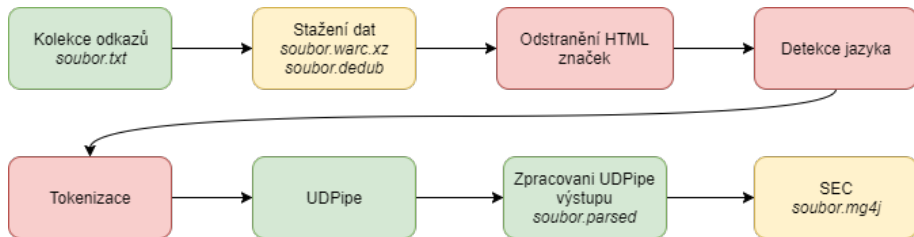
- Umíme vytvořit jednorázově korpus ze stažených.
- Sémanticky označit korpus.

Kde je tedy problém:

- Neumíme tvorbu korpusu provádět automaticky.
- Sémanticky označený korpus neumíme aktualizovat.
- Špatná detekce jazyka a deduplikace dat.
- Výpočetní náročnost celého procesu zpracování.

Odkud brát nová data:

- Atom a RSS zdroje,
- blogy,
- komentáře u článků,
- sociální sítě,
- ...



- Spousta skriptů (často i nečitelných),
- vertikalizace není úplně vychytaná (např. neošetřené chybové stavy),
- zpracování není automatizované,
- spousta logů, u chyb často není možné zjistit přesně příčinu.

- Sjednocení skriptů do větších celků (použití přepínačů, sdílené moduly) – během února,
- podrobně analyzovat vertikalizátor a buď opravit chyby nebo najít vhodnější řešení – během března (možná i duben),
- zautomatizovat zpracování – duben,
- sledování zpracování, měření náročnosti na procesor a paměť, tvorba statistik – duben, květen.

- Logy snadno parsovatelné regulárními výrazy nebo skripty,
- vedení podrobných statistik už během zpracování, možnost zasílat je pravidelně emailem,
- využití emailové notifikace na upozornění při kritických chybách.