

Bc. Milan Laslop, Bc. Tomáš Hurban

DISTRIBUOVANÉ VYHLADÁVANIE

Distribúované programové systémy

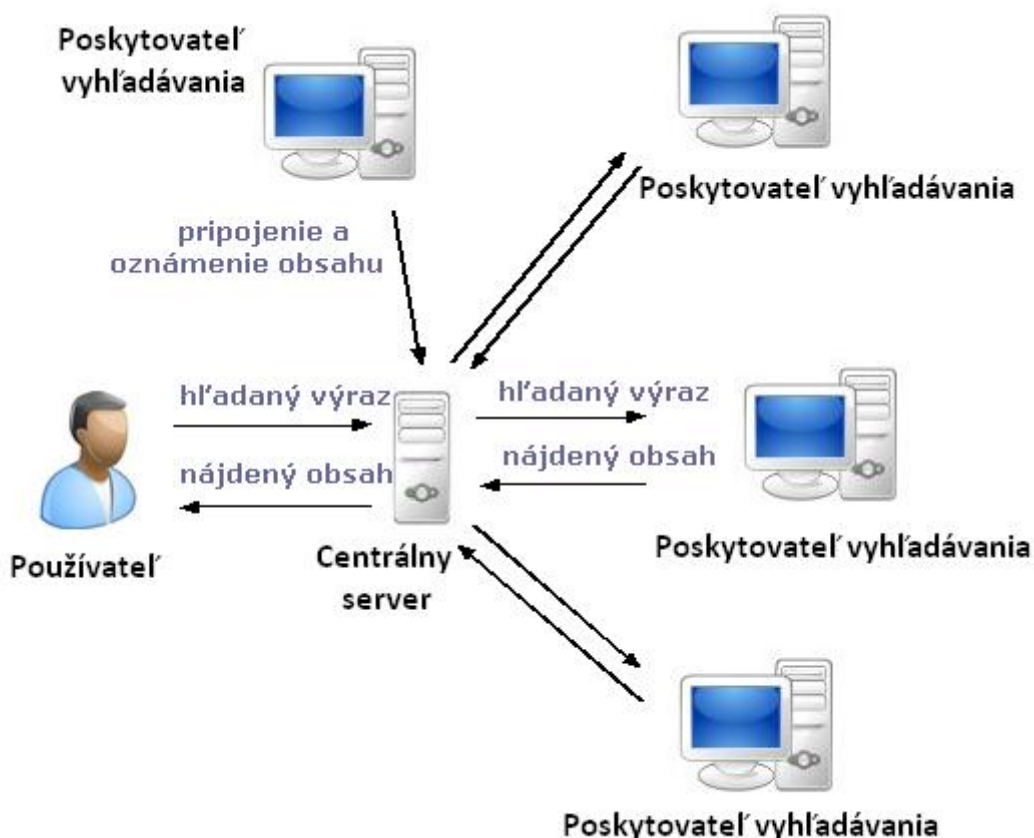
Obsah

1 Opis zadania	1
2 Architektúra systému.....	1
2.1 Dátové štruktúry.....	2
2.2 Implementované moduly.....	3
2.3 Rozhrania modulov	4
3 Implementácia	4
3.1 Vyhľadávanie v systéme	4
4 Riešené distribuované problémy	5
4.1 Riešené problémy.....	5
4.2 Nedoriešené problémy.....	6
4.3 Neriešené problémy.....	6

1 Opis zadania

Vyhľadávanie vo veľkom počte súborov alebo vo veľkých databázach býva často pomalé. Jedným z riešení tohto problému môže byť aj rozdistribuovanie obsahu (či už celého alebo len jeho častí) medzi viaceré vyhľadávacie jednotky. Cieľom tohto projektu bolo vytvoriť takýto distribuovaný systém, ktorý by celý obsah rozdelil medzi viacero samostatných vyhľadávacích jednotiek. Každá jednotka by tak mala lokálne uloženú len svoju časť databázy, v ktorej by vyhľadávala. Táto časť obsahu nemusí byť unikátna a môže ju mať uloženú viacero jednotiek. Keďže obsah sa môže časom meniť, je potrebné zabezpečiť synchronizáciu medzi jednotlivými vyhľadávacími jednotkami, aby sa zmenený obsah aktualizoval vo všetkých vyhľadávacích jednotkách, ktoré túto časť obsahu majú lokálne uloženú. Pri vyhľadávaní sa potom hľadaný výraz pošle viacerým vyhľadávacím jednotkám a z ich odpovedí sa potom vyskladá konečná odpoveď pre zadávateľa (klienta).

2 Architektúra systému



Obr. 1. Architektúra systému

Architektúra systému sa skladá z 3 častí:

- *Používateľ (Klient)* – chce vyhľadávať zadané slová alebo texty. Komunikuje len so serverom.
- *Poskytovatelia vyhľadávania* – majú pridelenú určitú časť obsahu, v ktorom vyhľadávajú. Komunikujú so serverom od ktorého získavajú hľadané výrazy a po prehladaní obsahu mu vrátia výsledky vyhľadávania.

Inicializácia poskytovateľa:

1. Pri pripojení poskytovateľa vyhľadávania k serveru, poskytovateľ oznámi serveru obsah lokálnej časti obsahu, v ktorom môže vyhľadávať. Server mu oznámi, ktoré časti obsahu by ešte mal mať a od ktorých iných poskytovateľov si ich má vypýtať.
 2. Poskytovateľ vyhľadávania získa časti obsahu od iných poskytovateľov.
- *Centrálny server* – komunikuje s klientom aj s poskytovateľmi vyhľadávania. Zabezpečuje vyberanie vhodných poskytovateľov, ktorým posiela požiadavky na vyhľadávanie a po vrátení jednotlivých výsledkov vyhľadávania z týchto vyskladá kompletnú odpoveď, ktorú pošle klientovi.

Všetok obsah sa pri spustení systému nachádza priamo na centrálnom serveri, no po pripojení poskytovateľov vyhľadávania sa tento obsah rozdistribuuje a počas vyhľadávania sa žiadny obsah na centrálnom serveri nenachádza.

2.1 Dátové štruktúry

Centrálny server:

```
-record(server_state, {providers, parts, connected_providers,  
waiting_parts}).
```

```
-record(part_info, {current_version, providers}).
```

```
-record(waiting_part_info, {part_data, copy_requests_to_providers}).
```

```
providers = dict()
parts = dict()
connected_providers = dict()
waiting_parts = dict()

current_version = number()
providers = dict()

part_data = term()
copy_requests_to_providers = dict()
```

Poskytovateľ vyhľadávania:

```
-record(provider_state, {id, parts}).

-record(part_info, {part_version, part_data}).
```

```
id = term()
parts = dict()

part_version = number()
part_data = term()
```

2.2 Implementované moduly

- central_server.erl – modul obsahuje implementáciu centrálného servera.
- search_provider_end_notifier.erl
- search_provider_supervisor.erl – modul obsahuje implementáciu supervisor pre poskytovateľov vyhľadávania.
- search_provider.erl – modul obsahuje implementáciu poskytovateľa vyhľadávania.
- tests.erl – modul s ukážkami testovacích funkcií.
- util.erl – modul s implementáciou funkcií na generovanie náhodných ID.

2.3 Rozhrania modulov

Tab. 1. Rozhrania modulov

	rozhranie pre správcu	rozhranie pre klienta	rozhranie pre centrálny server	rozhranie pre poskytovateľa
Centrálny server	<i>update</i> (PartName, PartData)	<i>search</i> (What) : results		<i>connect</i> (ProviderId, StateDiff) <i>disconnect</i> (ProviderId)
Poskytovateľ pripojenia	<i>create_and_activate</i> ()		<i>invalidate</i> (Pid) <i>search</i> (What, In, Pid) : results, search_times	<i>get</i> (PartName) : PartData

- **update** označuje aj vloženie novej časti obsahu, aj zmenu časti na vyššiu verziu (update teda označuje aj proces vkladania nových častí, aj počas počiatočnej inicializácie systému).

3 Implementácia

Centrálny server aj jednotlivý poskytovatelia vyhľadávania využívajú správanie *gen_server* a *supervisor*.

Na ukladanie údajov sa používajú dátové typy *dict* a *sets*.

3.1 Vyhľadávanie v systéme

Proces vyhľadávania prebieha nasledovne:

1. Klient zadá hľadaný výraz.
2. Server pošle požiadavky jednotlivým pripojeným poskytovateľom vyhľadávania na základe obsahu v ktorom môžu vyhľadávať.
3. Poskytovatelia obsahu vrátia nájdené texty.
4. Server tieto čiastočné výsledky poskladá do finálnej odpovede klientovi.

V systéme paralelne beží viacero procesov:

- 1 x central_server (gen_server)
- n x search_provider (gen_server)
- n x search_provider_supervisor (supervisor)
- n x search_provider_end_notifier:notifier_impl
- m x central_server:invalidate
- p x central_server:search_using_provider

4 Riešené distribuované problémy

4.1 Riešené problémy

- *Distribovanie obsahu* – server oznamuje poskytovateľom vyhľadávania po ich pripojení, ktoré časti obsahu majú mať a tým zabezpečuje rovnomerné distribuovanie a zálohovanie obsahu. Navyše treba riešiť inicializáciu celého systému (musí napríklad existovať jeden poskytovateľ, ktorý na začiatku má celý obsah, v ktorom sa bude hľadať, a tak sa tento obsah dá od neho distribuovať ďalším poskytovateľom, ktorí môžu vzniknúť neskôr).
- *Update starých častí obsahu* – centrálny server vie, ktorý poskytovateľ vyhľadávania má uložené ktoré časti obsahu, a preto pri zmene jednej časti obsahu vie poslať informáciu len tým poskytovateľom, ktorí túto časť obsahu majú, aby si mohli novú časť obsahu stiahnuť.
- *Rýchle vyhľadávanie vo veľkom obsahu* – centrálny server rozdelí celý obsah medzi poskytovateľov vyhľadávania a tým umožní paralelné prehľadávanie obsahu.
- *Riešenie výpadku poskytovateľa* – ak sa poskytovateľ odpojí od servera počas vyhľadávania zavolá sa funkcia *disconnect()* a tým sa poskytovateľ vyhodí zo zoznamu dostupných poskytovateľov vyhľadávania.

4.2 Nedoriešené problémy

- *Optimalizované vyhľadavanie* – server si pamätá, ktorému poskytovateľovi koľko trvalo vyhľadavanie v danej časti a snaží sa vyberať najrýchlejších poskytovateľov. Zaznamenávanie časov implementované je, no nie je implementované rozhodovanie o výbere najlepších poskytovateľov podľa týchto časov.

4.3 Neriešené problémy

- *Riešenie výpadku centrálného servera* – ak sa stane, že spadne centrálny server, tak spadne celý systém, pretože používateľ sa nemá inak ako cez centrálny server dostať k výsledkom vyhľadávania.
- *Synchronizácia viacerých centrálnych serverov bežiacich súčasne* – systém je navrhnutý na beh len s jedným centrálnym serverom.