



Vladimír Míč

Zápočtová úloha: životopis

Vzdělání

- 2009–2012 **Bakalářské studium**, *Fakulta informatiky Masarykovy univerzity, Brno, titul Bc.*
obor: Matematická informatika
- 2001–2009 **Střední škola**, *Gymnázium Křenová 36, Brno, maturita s vyznamenáním.*
osmileté studium

Závěrečná práce

- téma *Hierarchické shlukování rozsáhlých množin vektorových dat*
vedoucí RNDr. Tomáš Homola
popis Bakalářská práce vypracovaná na FI MU.

Praxe

Informatická praxe

- 2012–dosud **Externí spolupráce**, *Fakulta informatiky Masarykovy univerzity, Brno.*
Aplikovaný výzkum zaměřený na detekci a rozpoznávání obličejů v obrázcích.

Fyzikální praxe

- 2009 **Termoluminiscenční dozimetrie**, *Fakulta jaderná a fyzikálně inženýrská ČVUT, Praha.*
Článek a prezentace vytvořené v rámci Fyzikálního týdne 2009
- 2008 **Rentgenfluorescenční analýza - pomocník nejen při studiu památek**, *Fakulta jaderná a fyzikálně inženýrská ČVUT, Praha.*
Článek a prezentace vytvořené v rámci Fyzikálního týdne 2008
- 2007 **Stanovení délky a útlumu optického vlákna metodou optické reflektometrie**, *Fakulta jaderná a fyzikálně inženýrská ČVUT, Praha.*
Článek a prezentace vytvořené v rámci Fyzikálního týdne 2007

Hudební úspěchy

- 1996–2009 **Hra na klavír**, *ZUŠ Došlíkova 48, Brno.*

Hudební úspěchy

- 2009 **1. místo v krajském kole Národní soutěže ZUŠ**, Brno.
obor Komorní smyčcový soubor

Bzenecká 7, Brno, 628 00

☎ +420 728 237 xxx • ☎ 539 012 xxx • ✉ v.mic@mail.muni.cz

🌐 www.fi.muni.cz/~xmich/PB029/

- 2008 **2. Místo v okresním kole Národní soutěže ZUŠ**, Brno.
obor Hra na klavír
- 2005 **2. Místo v okresním kole Národní soutěže ZUŠ**, Brno.
obor Hra na klavír
- 2002 **2. Místo v okresním kole Národní soutěže ZUŠ**, Brno.
obor Hra na klavír

Jazyky

Angličtina **Upper intermediate**

Počítačové znalosti

- Programování
 - jazyk Java – pokročilý,
 - jazyk C – základy,
 - jazyk Haskell – základy.
- Znalosti softwaru
 - Microsoft office – pokročilý,
 - Open office – pokročilý,
 - komprimace videa:
 - MCE Buddy service,
 - Avidemux.

Zájmy

Věda zejména matematika,
technika počítače, mobilní telefony, fotoaparáty,
vážná hudba baroko, romantismus a další.

Bzenecká 7, Brno, 628 00

+420 728 237 xxx • 539 012 xxx • v.mic@mail.muni.cz
www.fi.muni.cz/xmic/PB029/

Úvod k bakalářské práci

Vladimír Míč

11. prosince 2012

S rozvojem informačních technologií dochází k obrovskému nárůstu velikosti ukládaných a zpracovávaných dat. Společně s tím dochází i k rozvoji sítí, a tím se zlepšuje datová dostupnost a možnosti sdílení. Hlavně v síťové oblasti se dnes běžně setkáváme se servery, které obsahují řádově desítky terabytů dat. Tento objem zaplní zejména multimediální soubory, konkrétně obrázky a videa, která se dnes po síti velmi rychle šíří. Pro zachování praktické použitelnosti, tedy dostatečně rychlé reakce na běžné příkazy, je nutné vylepšovat hardware a zároveň zavádět do praxe nové metody pro práci s daty [1].

Tato bakalářská práce vznikla pro potřeby vyhledávacího systému. Vyhledávání v rozsáhlých množinách dat už není možné realizovat lineárním průchodem jejich celého objemu, jako je tomu u menších datových úložišť, ani indexováním veškerého obsahu [2]. Data proto musí být předzpracována. Vhodnou volbou je použití některého z algoritmů pro shlukování. Jedná se o algoritmy, které rozdělí objekty do množin (dále třídy nebo shluky) a to tak, že objekty ve stejné třídě si jsou co nejvíce podobné a objekty v rozdílných třídách jsou co nejvíce odlišné [3]. Pro potřeby vyhledávání je dále vybrán pro každou třídu reprezentant. Ten se v rámci dané podobnosti nachází blízko středu, což znamená, že vystihuje danou třídu nejlépe. Reprezentantem může být i fiktivní objekt, který leží přesně uprostřed, a vystihuje tak danou třídu lépe, než kterýkoli její reálný objekt. Vyhledávání se realizuje pouze nad jednotlivými reprezentanty, čímž dochází k urychlení procesu. Algoritmy provádějící hierarchické shlukování zavádí více úrovní shluků. Počet úrovní může a nemusí být omezen.

Cílem bakalářské práce bylo implementovat v jazyce Java vybrané algoritmy hierarchického shlukování množiny obecných bodů a porovnat jejich výsledky, zejména rychlost a kvalitu zpracování. Požadované bylo využití knihovny Messif¹, což je knihovna jazyka Java vyvíjená na Fakultě informatiky Masarykovy univerzity v Brně (dále FI MU) pro práci nad metrickými prostory. Program měl být schopen zpracovávat různé typy dat. Podmínkou pro ně byla pouze jejich reprezentace pomocí třídy v jazyce Java, a možnost jejich porovnávání s jiným objektem stejného typu pomocí metody `getDistance()`. Výsledný program měl být schopen rozdělit zpracování dat paralelně na více počítačů tak, aby bylo možné konečný výsledek jednoduše získat z dílčích výsledků. Každý počítač měl také zpracovávat úlohu paralelně v rámci vícevláknové aplikace. Tím mělo být dosaženo co nejvyšší efektivity a rychlého zpracování při použití v praxi. Při psaní programu bylo využito třídy `AbstractObject` a dalších tříd knihovny Messif, čímž byly splněny první dva požadavky. Program využívá paralelismu na úrovni vláken a knihovnu Hadoop² společnosti Apache pro rozdělení výpočtů

¹Stránky knihovny Messif: <http://isd.fi.muni.cz/trac/messif>

²Stránky knihovny Hadoop: <http://hadoop.apache.org>

mezi více počítačů. Tím byly splněny zbývající požadavky.

Program byl začleněn do vyhledávacího systému Mufin³, rovněž vyvíjeného na FI MU, který slouží jako univerzální vyhledávací systém. V době vzniku bakalářské práce byl zaměřen zejména na obrázky, ale počítalo se s jeho použitím v rámci podobnostního vyhledávání v biologii, geografii, při vyhledávání plagiátů a v dalších oblastech. Matematický popis konkrétních dat, například obrázků, a tvorba metody `getDistance()` nebyly úkolem bakalářské práce.

Práce obsahuje i kapitolu věnující se testování konkrétních algoritmů hierarchického shlukování. V této kapitole jsou vidět rozdíly mezi časovou náročností algoritmů a kvalitou shluků při aplikaci na konkrétní soubor jednoho milionu objektů, jimiž byly textové popisy obrázků z galerie Flickr.com. Tyto popisy byly vytvořeny v rámci projektu Mufin. V praxi to může být i tisíckrát tolik. Po dohodě s vedoucím práce bylo usouzeno, že pro potřeby systému Mufin bude nevhodnější MacNaughton–Smithova metoda hierarchického shlukování.

Bakalářská práce obsahuje přehled základních metod hierarchického shlukování a jejich popis. Dále jsou uvedeny a rozebrány konkrétní algoritmy hierarchického shlukování včetně jejich složitostí, návrh programu, jeho implementace a testování pro použití pěti různých shlukovacích algoritmů. Diskutován je problém zastavení algoritmů.

První kapitola je věnována metodám hierarchického shlukování. Druhá kapitola se zabývá konkrétními algoritmy, třetí pak návrhem výsledného programu, čtvrtá jeho implementaci a pátá testováním.

Rejstřík

Java, 1

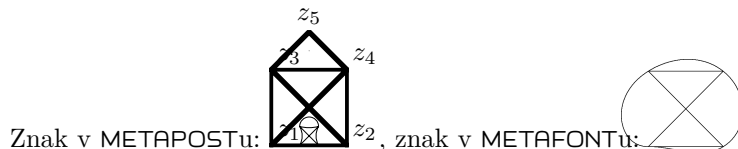
Vyhledávání, 1

Shluk, 1

Reference

- [1] Výzkum a vývoj v oblasti flexibilních infrastruktur.
- [2] Mark Baggesen. A review of windows search 4.0 from microsoft, 2011.
- [3] Jiří Kučera. Shluková analýza, 2008.

Grafika



³Stránky projektu Mufin: <http://mufin.fi.muni.cz/tiki-index.php>