

# Oponentský posudek diplomové práce

## Jindřich Šedek:

### Algoritmy nad rozšířeným sufixovým polem

Predkládaná práce zkoumá vhodnost náhrady sufixového stromu sufixovým polem v některých z jeho aplikací. Posuzována je praktická časová a prostorová složitost obou struktur a její závislost na menících se parametrech jednotlivých úloh.

Na práci je vidět, že se autor tématu dlouhodobě venoval. Vynaložené úsilí je dobře vidět na kvalitě kódu, který obsahuje rozsáhlé unittesty a je podrobně okomentován. Na druhou stranu, vlastní text práce má radu nedostatků, kterých bylo možné se snadno vyvarovat. Především jde o chybející implementační detaily a zdůvodnění uvedených implementačních rozhodnutí, a také o absenci pokusu zdůvodnit získané výsledky. I přes níže rozebírané problémy jde o dobrou práci, která by měla být obhájena.

Práce v první části uvádí do problematiky a podává přehled o diskutovaných strukturách, jejich rozšířeních a práci s nimi (kapitoly 1 až 4). V druhé části je pak popsána implementace jednotlivých struktur, metody měření, uvažované aplikace a výsledky experimentu.

Po krátkém úvodu v kapitole 1, přichází definice základních pojmů v kapitole 2. Zde je řada drobných chyb. Například není zřejmé co je „minimum z nejdelsích společných předpon“. Dále není definováno ani prázdné slovo ani  $S[n \dots n - 1]$ , ale je později použito v definici sufixového stromu i pole. V definici sufixového stromu chybí zmínka o orientaci hran směrem od kořene, ale implicitně se používá. Na obrázku sufixového stromu i CDAWGu chybí  $\$$  v nejdelsí příponě. Nadpis posledního sloupce v příkladu 2.3, i ve všech podobných uvedených později, má být  $S[suffix[i] \dots n - 1]\$$ .

Kapitola 3 popisuje známá rozšíření sufixového pole, která jej přibližují schopnostmi více k sufixovému stromu. U rozšíření je překvapivě rovnou diskutována implementace, ac se jí autor věnuje až později v kapitole 5, kde je pro ni lepší místo. Sekce 3.1 obsahuje v lematu 1 formulaci „předpona ... je menší nebo rovna ... předponě“, ale ve skutečnosti jde o srovnání délky obou předpon. Lemma 2 pak obsahuje formulaci „předpona se zmenší o jedničku“ místo „zkrátí se o jeden znak“. Na straně etnáct místo  $leptab - leptab[suffix[0]]$  má být  $leptab - leptab[ISA[0]]$ . Popis implementace  $leptab$  neobsahuje potřebné detaily navrhovaného řešení ani odkaz do literatury.

Sekce 3.2 popisuje  $lep$ -interval s hodnotou  $l$ , ale nezavádí pro něj vhodné značení, což později vede k těžko čitelným pasázím. Rovněž značení intervalu  $[i \dots j]$ , místo například  $[i, j]$ , není příliš šťastné. Ve zdůvodnění platnosti lematu 3 je znatelný argument zdůvodňující platnost třetí podmínky, jejíž platnost plyne rovnou z definice  $h$ . Rozbor věty 4 obsahuje chybu v indexu, kdy  $leptab[j]$  má být  $leptab[j + 1]$ . Část 3.3 popisuje rozšíření  $childtab$  a zavádí výšku vrcholu místo obvyklejší hloubky, která spíše měla být definována již u sufixového stromu a dále častěji používána. V příkladu  $childtab$  mají všechny 0 být 1. Oddíl 3.4 vysvětluje  $suffintab$ . Bohužel příklad 3.4 neodpovídá definici ani popsanému algoritmu. Řádky dvojsloupce  $suffintab$ , kde  $i$  i  $j$  jsou 0 nemají mít zápornou hodnotu. Také řádky s  $i = 0$  a  $j = 11$  nevyhovují definici ani algoritmu. Problém zřejmě pochází z článku [1], odkud autor převzal definici  $lep$ -intervalu, ale narozdíl od článku ji správně interpretoval a nevsílil si rozdíl.

jen jeden algoritmus. Sekce opět obsahuje řadu experimentálních výsledků.

V závěrečné kapitole 7 autor shrnuje obsah práce. I sem se vloudilo nedopatření v podobě poslední vety na straně sedesát, která zřejmě patří o dva odstavce výše.

V celé práci podle mne schází více odkazů do použité literatury. Především jde o první zmínky konceptu a postupu, kdy není na první pohled jasné, že nejde o autorův objev. Například sekce 3.2 neobsahuje žádnou citaci, ale jde o koncept převzatý z [1]. Nicméně, autor si rozhodně nepřivlastňuje cizí výsledky a cituje relevantní literaturu.

Priložené CD obsahuje vše potřebné a zejména umožňuje zhodnotit kvalitu kódu. Zde se dobře ukazuje vynaložené úsilí. Kód obsahuje rozsáhlé unit testy a je podrobně komentován. Navíc je k němu vygenerována i pěkná dokumentace. Kladem je rovněž bezproblémový preklad a fungování na 64 bitovém OS.

V Praze 15.5.2009

Martin Šenft

