

Toruń, 20 września 2023 r.

dr hab. Piotr Wiśniewski, prof. UMK  
Wydział Matematyki i Informatyki  
Uniwersytet Mikołaja Kopernika  
Chopina 12/18  
87-100 Toruń

## Recenzja rozprawy doktorskiej

mgra inż. Michała Kukowskiego

*Indeksowanie baz danych na nowoczesnych typach pamięci*

### 1. Tematyka rozprawy

Rozprawa doktorska magistra inżyniera Michała Kukowskiego dotyczy konstruowania optymalnych indeksów składowanych na nowoczesnych nośnikach pamięci trwałych. Temat badań jest ważny, gdyż obecnie wykorzystywane nośniki pamięci trwałych mają inne charakterystyki niż dawniej stosowane dyski magnetyczne. W szczególności lepiej radzą sobie z odczytem losowym danych, z drugiej strony stawiają znaczne większe wymagania w organizacji zapisu danych. W rozprawie zaprezentowano autorskie algorytmy dla pamięci typu flash, dysków SSD oraz pamięci zmiennofazowej PCM.

### 2. Ocena merytoryczna

Rozprawa skomponowana jest nasepująco: We wstępie autor wprowadza czytelnika w tematykę badań, rzeczowo uzasadniając ich celowość, a następnie prezentuje pokrótce zawartość poszczególnych rozdziałów. Rozdział drugi prezentuje metody indeksowania danych oraz formułuje problemy badawcze wynikające z przechowywania indeksów na nowoczesnych nośnikach pamięci. W rozdziale trzecim autor prezentuje nowoczesne nośniki pamięci wraz z porównaniem ich do dysków magnetycznych. Rozdział czwarty zawiera omówienie platformy testowej zaprojektowanej na potrzeby rzetelnej analizy wydajności opracowywanych algorytmów. Następne trzy rozdziały prezentują algorytmy opracowane w ramach prezentowanych badań. W piątym omówiony jest algorytm Flash Aware Tree – adaptacja drzewa na potrzeby pamięci flash w urządzeniach wbudowanych. Rozdział szósty jest najobszerniejszy, prezentuje kilka algorytmów przewidzianych dla dysków SSD, które z racji swojej olbrzymiej popularności wymagają najszerszego opracowania. Rozdział siódmy skupia się na algorytmach dla pamięci PCM, które jeszcze niedawno wydawały się istotnym przełomem. Każdy algorytm omawiany w tych trzech rozdziałach został gruntownie przetestowany na platformie testowej. Wyniki testów zostały zaprezentowane i przeanalizowane w rozprawie.

Rozdział ósmy jest podsumowaniem pracy oraz zapowiedzią dalszych kierunków badań.

Przechodząc do omawiania poszczególnych rozdziałów pominę rozdział wstępny.

WPLYNEŁO  
2.1-09-2023

20N-III/174 a/2023

## **Rozdział drugi**

Rozdział drugi prezentuje czytelnikowi ideę oraz metody indeksacji. Sam w sobie nie prezentuje nowych wyników badań. Jednakże porządkuje wiedzę czytelnika, co istotnie jest przydatne przy lekturze kolejnych rozdziałów.

## **Rozdział trzeci**

O ile w poprzednim rozdziale mówimy jedynie o porządkowaniu wiedzy, to ten rozdział wnosi czytelnikowi wiedzę, która nie jest aż tak powszechna. Omówione są w nim typy pamięci trwałych poczynając od przypomnienia budowy i zasad działania dysków magnetycznych, następnie przechodząc do pamięci Flash, dysków SSD, kończąc na pamięciach PCM. Autor klarownie prezentuje wszystkie aspekty budowy nowoczesnych pamięci, ze szczególnym uwzględnieniem zawiłości związanych z zapisem. Wiedza ta jest punktem wyjścia do opracowania algorytmów prezentowanych w rozdziałach 5 – 7.

## **Rozdział czwarty**

Platforma testowa omówiona w rozdziale czwartym opiera się o opracowany na potrzeby badań symulator SIPS oraz adaptację testów TPC w szczególności TPC-H i TPC-C. Możliwości prezentowanego symulatora są ciekawe również w oderwaniu od zastosowania w recenzowanej rozprawie, szkoda, że nie doczekał się samodzielnej publikacji. Konieczność opracowania symulatora jest podyktowana z jednej strony ograniczoną informacją pochodzącą od producentów dysków, z drugiej w dobie kryzysu pozwala zaoszczędzić środki wydatkowane na dyski testowe. W szczególności pamięci masowe PCM są bardzo drogie i trudno dostępne.

## **Rozdział piąty**

Pierwszy z trzech głównych rozdziałów pracy. Autor omawia w nim adaptację indeksowania za pomocą drzew dostosowaną do pamięci Flash typu NAND wykorzystywaną w tzw. urządzeniach wbudowanych. Opracowana struktura została nazwana Flash Aware Tree. Kluczowym aspektem modyfikacji jest efektywne wykorzystanie zapisu pełnych stron pamięci. Autor precyzyjnie opisuje metody działania z wprowadzonym mechanizmem reorganizacji poziomów drzewa. Po zaprezentowaniu algorytmu autor prezentuje wyniki testów w oparciu omówioną wcześniej platformę testową. Wyniki prezentowane w rozdziale piątym autor opublikował w [1].

## **Rozdział szósty**

Najobszerniejszy rozdział. Prezentuje on wyniki opublikowane przez autora w publikacjach [2,3,4]. Każda z tych prac wprowadza modyfikacje innego rodzaju indeksu. Podobnie każdemu indeksowi w rozprawie autor poświęcił osobny podrozdział. Pierwszym jest FA-LSM [2]. Zmodyfikowana względem klasycznego LSM struktura została zaadaptowana do potrzeb dysków SSD. W wyniku modyfikacji uzyskano kilkukrotne przyspieszenie. Nie do końca precyzyjna jest uwaga o 6 krotnie mniejszym wykorzystaniu pamięci flash w ostatnim zdaniu pierwszego akapitu podrozdziału 6.3. Sam algorytm został precyzyjnie wyjaśniony, a wyniki testów na symulatorze rzetelnie dowodzą tezy o przyspieszeniu.

W podrozdziałach 6.4 i 6.5 prezentowane są wyniki pracy [3] dotyczącej indeksu CF-Tree do indeksacji kolumnowych baz danych. Efektywność zaproponowanej struktury istotnie zależy od scenariusza użycia. Autor dokonuje gruntownej analizy scenariuszy wskazując kiedy proponowana przez niego struktura jest efektywna, a kiedy skuteczniejsze są struktury znane z literatury.

Podrozdziały 6.6 i 6.7 omawiają zagadnienie indeksowania częściowego [4]. Omówiony jest w nich mechanizm Lazy Adaptive Merging, będący adaptacją znanego z literatury Adaptive Merging. Główny pomysł dotyczy leniwego usuwania, które znacznie zmniejsza liczbę dokonywanych zapisów, jednocześnie oszczędzając dysk SSD i czas. Zaprezentowane przyspieszenie nie jest tak spektakularne jak w uzyskane w poprzednich metodach, jednakże oszczędność dysku jest istotnym aspektem.

### **Rozdział siódmy**

Pamięci PCM mają zupełnie inną charakterystykę niż pamięci Flash. Ich zapis i odczyt bardziej przypominają pamięci RAM - adresowane są bajtowo, nie ma potrzeby zapisu całych stron, jednakże wciąż zapis podobnie jak w przypadku pamięci flash jest nawet 10krotnie wolniejszy niż odczyt. Podobnie jak w poprzednich rozdziałach autor prezentuje najpierw metody indeksacji wierszowe, potem kolumnowe, a na końcu indeksy częściowe. Przy czym pierwsze dwa aspekty opracowane są na podstawie wiedzy zastanej, natomiast indeksowanie częściowe prezentuje nowe wyniki opublikowane w [5]. Jak poprzednio wprowadzone rozwiązania są dokładnie wyjaśniane oraz gruntownie przetestowane.

### **3. Wnioski końcowe**

Autor recenzowanej rozprawy zaproponował istotne modyfikacje algorytmów i struktur indeksujących, które pozwalają przyspieszyć obsługę baz danych jednocześnie istotnie oszczędzając zużycie nowoczesnych nośników pamięci. Wyniki te zostały zaprezentowane w 5 recenzowanych publikacjach naukowych. W rozprawie widać ogrom pracy jaką autor musiał wykonać przy opracowaniu symulatora i testach algorytmów.

Rozprawa została zredagowana w sposób czytelny i spójny. Stanowi doskonały przegląd metod indeksowania i problematyki nowoczesnych nośników danych. W opinii recenzenta rozprawa powinna zostać wydana jako całość w formie książki.

**Niniejszym stwierdzam, że recenzowana rozprawa spełnia wszystkie zwyczajowe i formalne normy stawiane rozprawom doktorskim w naukach technicznych, w dyscyplinie informatyka techniczna i telekomunikacja. Wnoszę o dopuszczenie mgr inż. Michała Kukowskiego do dalszych etapów przewodu doktorskiego.**

*Piotr Wiśniewski*