



Universiteit
Leiden
The Netherlands

Progressive Indexes

Timbó Holanda, P.T.

Citation

Timbó Holanda, P. T. (2021, September 21). *Progressive Indexes. SIKS Dissertation Series*. Retrieved from <https://hdl.handle.net/1887/3212937>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/3212937>

Note: To cite this publication please use the final published version (if applicable).

Progressive Indexes

Pedro Holanda

Progressive Indexes

Proefschrift

ter verkrijging van
de graad van doctor aan de Universiteit Leiden,
op gezag van rector magnificus prof.dr.ir. H. Bijl,
volgens besluit van het college voor promoties
te verdedigen op dinsdag 21 september 2021
klokke 10:00 uur

door

Pedro Thiago Timbó Holanda
geboren te Fortaleza, Brazilië
in 1992

Promotiecommissie

Promotor:	prof. dr. Stefan Manegold	(CWI & Universiteit Leiden)
Copromotores:	dr. Hannes Mühleisen	(CWI & UvA)
	prof. dr. Peter Boncz	(CWI & VU)
Overige leden:	prof. dr. Aske Plaat	(Universiteit Leiden)
	prof. dr. Thomas Bäck	(Universiteit Leiden)
	prof. dr. Yanlei Diao	(École Polytechnique de Paris)
	dr. Stratos Idreos	(Harvard University)
	dr. Eduardo Cunha de Almeida	(UFPR)

The research reported in this thesis has been carried out within the Database Architectures group at Centrum Wiskunde & Informatica (CWI), the National Research Institute for Mathematics and Computer Science in the Netherlands.

SIKS Dissertation Series No. 2021-21 The research reported in this thesis has been carried out under the auspices of SIKS, the Dutch Research School for Information and Knowledge Systems.

This research is financially supported by the Dutch funding agency NWO, under project number 628.006.002 (the DAMIOSO project), in collaboration with Honda Research Institute, Leiden University and Centrum Wiskunde & Informatica (CWI).



Acknowledgments

I want to thank my direct supervisor, Stefan Manegold. I must say that I feel like I won the supervisor lottery. Stefan always gave me complete freedom to pursue any path I wanted during my studies while encouraging me to see the weak points of any ideas I had. To this day, I am still amazed by his "eagle eyes" skill. His attention to detail is awe-inspiring, and having the opportunity of having him review all my papers, really enabled me to take them to the next level.

I want to acknowledge the impact that my co-supervisor, Hannes Mühleisen had on my thesis and post-PhD career. He provided me with the opportunity to work with him on DuckDB. It was one of the most fun projects I had during my Ph.D. and helped me improve my programming skills. In the final months of my Ph.D., when learning I had the intention to continue in the Netherlands for a few more years, Hannes also moved mountains to secure a post-doc position for me (on record time, I must say).

I have a high debt with Mark Raasveldt. I was fortunate to have Mark as a friend (and as a roommate for almost two years), and although we've been Ph.D. students at the same time, the reality is that Mark was already on another level. Living with him was an excellent opportunity to expand all the skill-set needed to be a successful researcher. Not only was he always available for any presentation/writing and programming questions I would have, but working with him was truly fun. With Mark, not all were related to work. As my flatmate, he also introduced me to all the greatness of Dutch culture (e.g., Febo, Action, New-Kids Turbo).

In the office, I was lucky enough to have the greatest officemate ever, Tim Gubner. Tim, thanks for all the great times and for singing with me all the greatest hits of Jon Lajoie, Tenacious D, and Backstreet Boys. I'm pleased that we could continue our musical journey even after being asked if we were killing rats in the office (although we were just singing slightly out of tune). I also think we still have the hidden potential of revolutionizing the computer science field with all of our great research ideas (e.g., Commie-coin, a communistic crypto-currency; DataBreaks: Breaks for Fast Databases). Thank you for all the laughs, and I hope we get to share an office again in the future!

During the Corona Pandemic, I was very thankful I managed to bring a little "Brazilian Gang" into the Database Architectures group. Diego Tome and Matheus Nerone helped me maintain my sanity throughout the whole of 2020 by having daily coffee breaks, workouts, and video-game marathons. I think 2020 would be almost

unbearable without our activities. Besides all that, I also enjoyed working with Matheus on the first few months of the pandemic. I think we both did a great job motivating each other to finish off our multidimensional work, and I'm still impressed with what he accomplished with his mad plotting skills.

I would also like to thank my friends that made my years in Amsterdam some of the best years in my life. Especially Bianca Jabur (Kiki), Tijs Kramer, and George Anastasiou, thanks for all the fun!

Last but not least, I would like to thank my family for their support in all of my academic and life choices: my parents Tarcisio and Ana Holanda, and my sister Camila Holanda.

Contents

1	Introduction	11
1	Data Analysis	11
2	Interactive Data Analysis	12
2.1	Index Creation Problem	13
2.2	Research Questions	15
3	Our Contributions	16
4	Structure and Covered Publications	16
2	Background	19
1	Relational Database Systems	19
1.1	Physical Layout	20
2	Interactive Exploratory Data Analysis	21
3	Index Structures	22
4	Index Selection Problem	23
4.1	Automatic Index Selection	24
4.2	Adaptive Index Creation	24
3	Progressive Indexing	27
1	Introduction	27
1.1	Contributions	29
1.2	Outline	29
2	Related Work	29
2.1	Cracking Kernels	30

2.2	Adaptive Indexing for Robustness	31
3	Progressive Indexing	35
3.1	Progressive Quicksort	37
3.2	Progressive Radixsort (MSD)	38
3.3	Progressive Bucktersort	40
3.4	Progressive Radixsort (LSD)	42
4	Greedy Progressive Indexing	43
4.1	Greedy Progressive Quicksort	44
4.2	Greedy Progressive Radixsort (MSD)	46
4.3	Greedy Progressive Bucktersort	47
4.4	Greedy Progressive Radixsort(LSD)	47
5	Experimental Analysis	48
5.1	Setup.	48
5.2	Delta Impact	50
5.3	Cost Model Validation	53
5.4	Interactivity Threshold	56
5.5	Varying Interactivity	59
5.6	Adaptive Indexing Comparison	61
6	Summary	66
4	Multidimensional Progressive Indexing	69
1	Introduction	69
1.1	Contributions	70
1.2	Outline	70
2	Related Work	71
2.1	Multidimensional Data Structures	71
2.2	Adaptive/Progressive Index	73
3	Multidimensional Progressive Indexing	75
3.1	Data Structure	76
3.2	Creation Phase	77
3.3	Refinement Phase	80
3.4	Greedy Progressive Indexing	82
4	Experimental Analysis	85
4.1	Setup.	85
4.2	Data Sets & Workloads	86
4.3	Delta Impact	87

4.4	Performance Comparison	91
4.5	Impact of Dimensionality	96
4.6	Full Scan Exceeding the Interactivity Threshold	97
5	Summary	98
5	Progressive Merges	101
1	Introduction	101
1.1	Contributions	102
1.2	Outline	102
2	Related Work	102
2.1	Merge Complete (MC)	103
2.2	Merge Gradual (MG)	104
2.3	Merge Ripple (MR)	104
3	Progressive Mergesort	106
4	Experimental Analysis	110
4.1	Setup	111
4.2	Performance Comparison	112
4.3	Varying Data Sizes	113
4.4	Appends during Index Creation	115
5	Summary	116
6	Big Picture	117
1	The Elephant In The Room	117
2	Future Work	118
2.1	Progressive Indexes	118
2.2	Progressive Merges	119
	Summary	123
	Samenvatting	125
	Publications	127

