

On the optimization of imaging pipelines Schoonhoven, R.A.

Citation

Schoonhoven, R. A. (2024, June 11). On the optimization of imaging pipelines. Retrieved from https://hdl.handle.net/1887/3762676

Version:	Publisher's Version
License:	Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden
Downloaded from:	<u>https://hdl.handle.net/1887/3762676</u>

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

CURRICULUM VITAE

Richard Schoonhoven was born in 1995 in Utrecht, the Netherlands. He completed his secondary education at RSG Broklede in Breukelen, the Netherlands. Afterwards, he completed bachelor's degrees (both cum laude) in mathematics and physics at Utrecht University in 2016, and obtained master's degrees (both cum laude) in mathematics and computer science in 2019 from Utrecht University. The master's thesis with the title "Improving cryo-ET reconstructions of ER-associated ribosomes with tomographic reconstruction methods and deep learning" was supervised by Dr. Tristan van Leeuwen. In 2019, he started as a PhD candidate at Centrum Wiskunde & Informatica (the national research institute for mathematics and computer science in Amsterdam) under the supervision of Prof.dr. K.J. Batenburg.

ACKNOWLEDGMENTS

First, I thank my advisors, Prof. Joost Batenburg and Dr. Daniël Pelt, for the time and energy they have put into providing motivating, educational, and fun supervision of our research projects. In particular, I would like to thank them for providing me with large amounts of freedom to pursue different research directions, which has made the past four years a thoroughly enjoyable part of my career.

In addition, I would like to thank Dr. Ben van Werkhoven for his energetic and stimulating approach to research and supervision, which has made several of our projects a great pleasure to collaborate on.

I would also like to thank my colleague Dr. Alexander Skorikov for his companionable collaboration, and his tireless energy for dealing with my ceaseless stream of thoughts, and at times, tumultuous working style.

A special thank you goes out to my co-authors at the ESRF who have hosted me for several weeks on many occasions, and for their efforts to make me feel welcome and show me around their impressive facilities. I have to mention in particular Dr. Alexandra Pacureanu who was kind enough to allow me to stay in their city apartment in Grenoble during these trips.

I would like to thank my co-authors Allard Hendriksen, Bram Veenboer, and Jan-Willem Buurlage, who have greatly helped me with my research projects and spent time and effort teaching me more about new topics.

In particular, I would like to thank Willem Jan Palenstijn for his expert and patient help on countless problems I encountered.

I thank the colleagues whom I shared an office with, Vladyslav Andriiashen, Mathé Zeegers, Adriaan Graas, Poulami Ganguly, Francien Bossema, Jordi Minnema, Dirk Schut, Maximilian Kiss, Tianyuan Wang, Rien Lagerwerf for providing stimulating conversations and a pleasant work environment.

Many others contributed to a great research environment, among which Floris-Jan, Jiayang, Serban, Alex, Nicola, Henri, Maureen, Dzemila, Georgios, Sophia, Felix, Rob, Robert, Hamid, Roozbeh, Ajynkya, Jan, and Tristan.

I would like to thank my family and friends for their support, and camaraderie over the years.

Finally, I would like to thank Sasha for her love, infinite patience, and wholehearted support.

A Appendices

A.1 Appendix: (LEAN) graph-based pruning for convolutional neural networks by extracting longest chains

A.1.1 Datasets

In this appendix we discuss some more details on the datasets used for experimentation.



Figure A1: Example input and target images of the (top left) Circle-Square (CS), (top right) CamVid, (bottom) real-world dynamic CT datasets.

Simulated Circle-Square (CS) dataset: We used a simulated high-noise 5-class segmentation dataset containing 256×256 images of randomly placed squares and circles (CS dataset) [183] (see Figure A1). The objects were assigned a random grey value and Gaussian noise was added to the images. In total, we generated 1000 training, 250 validation, and 100 test images. Experimental results on the CS dataset are quantified using global accuracy, i.e., the ratio of correctly classified pixels, regardless of class, to the total number of pixels.

CamVid: The Cambridge-driving Labeled Video Database (CamVid) [24, 25] is a collection of videos with labels, captured from the perspective of a driving automobile. In total, 700 labeled frames are split into 367 training, 100 validation, and 233 test images. As there are few training images, we combined the training and validation datasets and trained for a fixed 500 epochs. Similar to other papers that apply CNNs to CamVid [9, 180], we use 11 classes, and a single class representing unlabeled pixels (see Figure A1).

We used median frequency balancing [56] to balance classes for training, and set the unlabeled class weights to zero. During training, we used data augmentation by cropping and (horizontally) flipping input images.

A.1. APPENDIX: (LEAN) GRAPH-BASED PRUNING FOR CONVOLUTIONAL NEURAL NETWORKS BY EXTRACTING LONGEST CHAINS 167



Figure A2: Adjacency matrices of active convolutions (in white) after pruning. All pruned network were pruned to a ratio of 10%. From left to right, we have the unpruned matrix of a 100-layer MS-D network trained on the real-world dynamic CT dataset, randomly pruned convolutions, structured magnitude pruning, structured operator norm pruning, and LEAN.

Real-world dynamic CT dataset: The real-time dynamic X-ray CT dataset contains images of a dissolving tablet suspended in gel [38, 39]. The bubbles are to be segmented within a glass container filled with gel [209] (see Figure A1). The dataset consists of 512×512 images, split into 9216 training images, 2048 validation images, and 1536 test images. As in [209], we use the F1-score because the large amount of background pixels make global accuracy an unsuitable metric.

A.1.2 Reducing the size of the pruning graph

The procedure outlined in Section 4.4 can lead to large pruning graphs, but the size of the graph can be reduced. First, according to Equation 4.3, the operator norm of ReLU is 1. Therefore, the combination of a convolution followed by a ReLU can be combined into a single edge whose weight equals the norm of the convolution.

Batch normalization often succeeds a convolution. Batch-normalization scaling is applied with different learned parameters per input channel, and output a single channel. Therefore, the input convolution edge and the following batch normalization edges can be combined. The edges can be combined into a single edge whose weight is the product of the two edge weights, preserving the path length.

A.1.3 Structure of pruned MS-D networks

To investigate the structure of pruned networks we plotted the adjacency matrices of pruned networks where an entry is 0 if it is pruned (black) and 1 if it is still active (white). Here, we show the adjacency matrices of MS-D networks pruned to a ratio of 10% in Figure A2. After pruning, LEAN retains only connections linked to nearby layers in the densely connected MS-D network. Compared to individual filter pruning, LEAN exposes a distinct structure which may suggest that LEAN could be used for architecture discovery.

A.2 Appendix: Benchmarking optimization algorithms for auto-tuning GPU kernels

A.2.1 Tunable parameters per GPU kernel

In Table A.1 we show the tunable parameters per kernel, and the values each parameter could take. For the convolution kernel, the MI50 GPU (the only AMD model) required a different problem setup due to hardware constraints.

Kernel	parameter to tune	list of values	number of
			possible values
Convolution	block_size_x	1, 2, 4, 8, 16, 32, 48,	12
(except MI50)		64, 80, 96, 112, 128	
	block_size_y	1, 2, 4, 8, 16, 32	6
	tile_size_x	1, 2, 3, 4, 5, 6, 7, 8	8
	tile_size_y	1, 2, 3, 4, 5, 6, 7, 8	8
	use_padding	0, 1	2
	read_only	0, 1	2
Convolution	block_size_x	16, 32, 48, 64,	8
(MI50)		80, 96, 112, 128	
	block_size_y	1, 2, 4, 8, 16, 32	6
	tile_size_x	1, 2, 4	3
	tile_size_y	1, 2, 4	3
	use_padding	0, 1	2
GEMM	MWG	16, 32, 64, 128	4
	NWG	16, 32, 64, 128	4
	MDIMC	8, 16, 32	3
	NDIMC	8, 16, 32	3
	MDIMA	8, 16, 32	3
	NDIMB	8, 16, 32	3
	VWM	1, 2, 4, 8	4
	VWN	1, 2, 4, 8	4
	SA	0, 1	2
	SB	0, 1	2
Point-in-polygon	block_size_x	32, 64, 96, 128, 160, 192, 224,	31
		256, 288, 320, 352, 384, 416,	
		448, 480, 512, 544, 576, 608,	
		640, 672, 704, 736, 768, 800,	
		832, 864, 896, 928, 960, 992	
	tile_size	1, 2, 4, 6, 8, 10,	11
		12, 14, 16, 18, 20	
	between_method	0, 1, 2, 3	4
	use_precomputed_slopes	0, 1	2
	use_method	0, 1, 2	3

Table A.1: Tunable parameters per kernel, and list of possible values for each parameter.

A.2.2 Alternative splits for competition heatmaps

In Figures A1 and A6 we show the algorithm competition heatmaps such as in Figures 5.1, 5.2 and 5.3, but when split at 100 and 400 budgets instead of 200.



Algorithm Column beats Row - convolution feval <= 100



A.2. APPENDIX: BENCHMARKING OPTIMIZATION ALGORITHMS FOR AUTO-TUNING GPU KERNELS 171

	'			11 01			eut.	5 R0	vv -	GEN	1111	eva	<=	- 10	0
BasinHopping	0	0	0	0	1	4	1	2	0	2	0	0	0	0	1
BestILS	15	0	4	0	16	18	14	12	1	10	14	12	11	9	15
BestMLS	15	1	0	0	15	16	8	4	1	7	12	11	7	7	12
BestTabu	18	8	13	0	18	18	16	16	4	13	18	15	14	16	17
DifferentialEvolution	12	0	0	0	0	10	3	1	0	1	0	0	0	0	2
DualAnnealing	6	0	0	0	0	0	1	0	0	0	0	0	0	0	0
FirstILS	11	0	0	0	8	13	0	0	0	0	5	8	6	5	6
FirstMLS	13	0	0	0	11	14	2	0	0	1	6	9	7	5	7
FirstTabu	17	5	9	0	17	18	12	11	0	11	16	13	14	13	16
GLS	14	0	2	0	12	14	4	0	0	0	8	10	9	5	9
GeneticAlgorithm	13	0	0	0	7	14	3	2	0	2	0	2	1	0	2
ParticleSwarm	15	1	2	0	12	15	6	6	0	5	5	0	1	1	5
RandomSampling	17	5	5	0	15	16	8	6	0	5	9	6	0	3	8
SMAC4BB	17	3	4	0	17	17	9	6	0	6	11	9	8	0	12
SimulatedAnnealing	13	0	1	0	9	12	3	1	0	1	2	2	2	1	0
	}asinHopping	BestILS	BestMLS	BestTabu	itialEvolution	ualAnnealing	FirstILS	FirstMLS	FirstTabu	GLS	sticAlgorithm	articleSwarm	lomSampling	SMAC4BB	tedAnnealing
	,	Algo	rith	m C	olur	nn k	beat	s Ro	w -	GEI	MM	feva	>	100	
BasinHopping	0	16	17	11	13	16	18	18	14	18	13	13	0	0	20
BestILS	2	0	1	6	1	3	16	9	6	6	0	1	0	0	13
BestMLS	3	6	0	7	2	5	14	9	9	٩	1	-	0	0	19
BestTabu	2	7	•				14			9		1			
DifferentialEvolution			8	0	6	7	14	11	5	8	7	1 5	0	0	11
	3	16	8 18	0 9	6 0	7 13	14 10 22	11 20	5 13	8 21	7 13	1 5 11	0	0 0	11 21
DualAnnealing	3 1	16 14	8 18 15	0 9 10	6 0 1	7 13 0	14 10 22 19	11 20 17	5 13 12	8 21 17	1 7 13 5	1 5 11 6	0 0 0	0 0 0	11 21 20
DualAnnealing FirstILS	3 1 0	16 14 1	8 18 15 0	0 9 10 2	6 0 1 0	7 13 0 0	14 10 22 19 0	11 20 17 2	5 13 12 4	8 21 17 2	1 7 13 5 0	1 5 11 6 1	0 0 0 0	0 0 0 0	11 21 20 6
DualAnnealing FirstILS FirstMLS	3 1 0 1	16 14 1 6	8 18 15 0 2	0 9 10 2 5	6 0 1 0 0	7 13 0 0 1	14 10 22 19 0 9	11 20 17 2 0	5 13 12 4 7	9 8 21 17 2 2	1 7 13 5 0 0	1 5 11 6 1 0	0 0 0 0 0	0 0 0 0 0	111 21 20 6 111
DualAnnealing FirstILS FirstMLS FirstTabu	3 1 0 1 4	16 14 1 6 5	8 18 15 0 2 6	0 9 10 2 5 2	6 0 1 0 0 4	7 13 0 1 1 7	14 10 22 19 0 9 14	11 20 17 2 0 9	5 13 12 4 7 0	8 21 17 2 2 8	1 7 13 5 0 0 0 6	1 5 11 6 1 0 4	0 0 0 0 0 0	0 0 0 0 0	11 21 20 6 11 10
DualAnnealing FirstILS FirstMLS FirstTabu GLS	3 1 0 1 4 1	16 14 1 6 5 3	8 18 15 0 2 6 0	0 9 10 2 5 2 5	6 0 1 0 0 4 0	7 13 0 1 1 7 1	14 10 22 19 0 9 14 10	11 20 17 2 0 9 1	5 13 12 4 7 0 7	9 8 21 17 2 2 8 8 0	1 7 13 5 0 0 6 0	1 5 11 6 1 0 4 0	0 0 0 0 0 0 0 0	0 0 0 0 0 0	11 21 20 11 10 11
DualAnnealing FirstILS FirstMLS FirstTabu GLS GeneticAlgorithm	3 1 0 1 4 1 2	16 14 1 6 5 3 13	8 18 15 0 2 6 0 13	0 9 10 2 5 2 5 9	6 0 1 0 4 0 1	7 13 0 1 1 7 1 1 6	14 10 22 19 0 9 14 10 20	11 20 17 2 0 9 1 20	5 13 12 4 7 0 7 11	9 8 21 17 2 2 8 0 17	1 7 13 5 0 0 6 0 0 0	1 5 11 6 1 0 4 0 4	0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0	11 21 20 6 11 10 11 22
DualAnnealing FirstILS FirstMLS FirstTabu GLS GeneticAlgorithm ParticleSwarm	3 1 0 1 4 1 2 4	16 14 1 6 5 3 13 13	8 18 15 2 6 6 13 13	0 9 10 2 5 2 5 9 8	6 0 1 0 4 0 1 3	7 13 0 1 1 7 1 1 6 9	14 10 22 19 0 9 14 10 20 21	11 20 17 2 0 9 1 20 21	5 13 12 4 7 0 7 11 11	8 21 17 2 2 8 0 17 20	1 7 13 5 0 0 6 0 0 0 0 9	1 5 11 6 1 0 4 0 4 0 4 0	0 0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0	11 21 20 11 10 11 22 20
DualAnnealing FirstILS FirstMLS FirstTabu GLS GeneticAlgorithm ParticleSwarm RandomSampling	3 1 0 1 4 1 2 4 24	16 14 1 6 5 3 13 13 17 24	8 18 15 2 6 0 13 13 17 24	0 9 10 2 5 2 5 9 8 8 18	6 0 1 0 4 0 1 3 24	7 13 0 1 1 7 1 6 9 24	14 10 22 19 0 9 14 10 20 21 21 24	111 20 177 2 0 9 1 20 21 21	5 13 12 4 7 0 7 11 11 11	8 21 17 2 2 8 0 17 20 24	1 7 13 5 0 0 6 0 6 0 0 9 24	1 5 11 6 1 0 4 0 4 0 4 0 23	0 0 0 0 0 0 0 0 0 0 0 0	0 0 0 0 0 0 0 0 0 0 0 0	11 21 20 11 10 11 22 20 24
DualAnnealing FirstILS FirstMLS FirstTabu GLS GeneticAlgorithm ParticleSwarm RandomSampling SMAC4BB	3 1 0 1 4 1 2 4 24 11	16 14 1 6 5 3 13 13 17 24 12	8 18 15 0 2 6 0 13 13 17 24 12	0 9 10 2 5 2 9 8 8 18 18	6 0 1 0 4 0 1 3 24 12	7 13 0 1 1 7 1 3 4 9 24 12	14 10 22 19 9 14 10 20 21 21 24 12	11 20 17 2 0 9 1 20 21 21 21 24	5 13 12 4 7 0 7 11 11 11 11 19 7	 8 21 17 2 8 0 17 20 24 12 	1 7 13 5 0 0 0 0 0 0 0 0 12	1 5 11 6 1 1 0 4 0 4 0 4 0 2 3 1 1	0 0 0 0 0 0 0 0 0 0 0 0 0 2	0 0 0 0 0 0 0 0 0 0 0 0 0 0	11 21 20 6 11 10 11 22 20 24 12
DualAnnealing FirstILS FirstMLS FirstTabu GLS GeneticAlgorithm ParticleSwarm RandomSampling SMAC4BB SimulatedAnnealing	3 1 0 1 4 1 2 4 24 11 0	16 14 1 6 5 3 13 13 17 24 12 0	 8 18 15 0 2 6 0 13 17 24 12 0 	 9 10 2 5 4 5 9 8 18 6 3 	6 0 1 0 4 0 1 3 24 12 0	7 13 0 1 1 7 1 1 6 9 24 12 24	14 10 22 19 0 9 14 10 20 21 24 12 24 12	11 20 17 2 0 9 1 20 21 20 21 24 12 0	5 13 12 4 7 0 7 11 11 11 19 7 5	 8 21 17 2 8 0 17 20 24 12 1 	1 7 3 5 0 0 0 6 0 0 0 0 0 0 9 24 12 0	1 5 11 6 1 1 0 4 0 4 0 4 0 2 3 11 0	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 2 0	0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	11 20 6 11 10 11 22 20 24 12 0

Figure A2: (GEMM:) Occurrences when the column algorithm found better solutions than the row algorithm. An occurrence is counted when 50 runs for a budget are statistically significantly better according to a two-sample independent t-test ($\alpha = 0.05$). (Top): Heatmap for low ≤ 100 budgets (25, 50, 100). (Bottom): Heatmap for mid and high > 100 budgets (200, 400, 800, 1600). Algorithms with low values (blue) in their rows were not often beaten for those budgets, and algorithms with high values in their column (red) often beat other algorithms.



Algorithm Column beats Row - pnpoly feval <= 100

Figure A3: (PnPoly:) Occurrences when the column algorithm found better solutions than the row algorithm. An occurrence is counted when 50 runs for a budget are statistically significantly better according to a two-sample independent t-test ($\alpha = 0.05$). (Top): Heatmap for low ≤ 100 budgets (25, 50, 100). (Bottom): Heatmap for mid and high > 100 budgets (200, 400, 800, 1600). Algorithms with low values (blue) in their rows were not often beaten for those budgets, and algorithms with high values in their column (red) often beat other algorithms.

A.2. APPENDIX: BENCHMARKING OPTIMIZATION ALGORITHMS FOR AUTO-TUNING GPU KERNELS 173



Figure A4: (Convolution:) Occurrences when the column algorithm found better solutions than the row algorithm. An occurrence is counted when 50 runs for a budget are statistically significantly better according to a two-sample independent t-test ($\alpha = 0.05$). (Top): Heatmap for low ≤ 400 budgets (25, 50, 100, 200, 400). (Bottom): Heatmap for mid and high > 400 budgets (800, 1600). Algorithms with low values (blue) in their rows were not often beaten for those budgets, and algorithms with high values in their column (red) often beat other algorithms.



Algorithm Column beats Row - GEMM feval <= 400

Figure A5: (GEMM:) Occurrences when the column algorithm found better solutions than the row algorithm. An occurrence is counted when 50 runs for a budget are statistically significantly better according to a two-sample independent t-test ($\alpha = 0.05$). (Top): Heatmap for low ≤ 400 budgets (25, 50, 100, 200, 400). (Bottom): Heatmap for mid and high > 400 budgets (800, 1600). Algorithms with low values (blue) in their rows were not often beaten for those budgets, and algorithms with high values in their column (red) often beat other algorithms.

A.2. APPENDIX: BENCHMARKING OPTIMIZATION ALGORITHMS FOR AUTO-TUNING GPU KERNELS 175



Figure A6: (PnPoly:) Occurrences when the column algorithm found better solutions than the row algorithm. An occurrence is counted when 50 runs for a budget are statistically significantly better according to a two-sample independent t-test ($\alpha = 0.05$). (Top): Heatmap for low ≤ 400 budgets (25, 50, 100, 200, 400). (Bottom): Heatmap for mid and high > 400 budgets (800, 1600). Algorithms with low values (blue) in their rows were not often beaten for those budgets, and algorithms with high values in their column (red) often beat other algorithms.

A.2.3 Per kernel graphs of experimental results

In Figures A7 to A15 we show plots of algorithm performance in terms of fraction of optimal fitness found for certain budget used (per GPU).



Figure A7: **Convolution:** Fraction of optimal runtime per GPU for FirstILS, FirstMLS, dual annealing, simulated annealing, and GLS over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A8: **Convolution:** Fraction of optimal runtime per GPU for GA, BestMLS, BestILS, basin hopping, and differential evolution over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A9: **Convolution:** Fraction of optimal runtime per GPU for SMAC, FirstTabu, BestTabu, PSO, and random sampling over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic *x*-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A10: **GEMM:** Fraction of optimal runtime per GPU for FirstILS, FirstMLS, dual annealing, simulated annealing, and GLS over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A11: **GEMM:** Fraction of optimal runtime per GPU for GA, BestMLS, BestILS, basin hopping, and differential evolution over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A12: **GEMM:** Fraction of optimal runtime per GPU for SMAC, FirstTabu, BestTabu, PSO, and random sampling over 50 runs. Each point is the mean fraction of optimal runtime found (*y*-axis) for mean budget used (logarithmic *x*-axis), with error bars indicating the standard deviation in fraction of optimum.



Figure A13: **Point-in-polygon:** Fraction of optimal runtime per GPU for FirstILS, FirstMLS, dual annealing, simulated annealing, and GLS over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum. The point-in-polygon kernel was not implemented for the MI50 GPU.



Figure A14: **Point-in-polygon:** Fraction of optimal runtime per GPU for GA, BestMLS, BestILS, basin hopping, and differential evolution over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum. The point-in-polygon kernel was not implemented for the MI50 GPU.



Figure A15: **Point-in-polygon:** Fraction of optimal runtime per GPU for SMAC, FirstTabu, BestTabu, PSO, and random sampling over 50 runs. Each point is the mean fraction of optimal runtime found (y-axis) for mean budget used (logarithmic x-axis), with error bars indicating the standard deviation in fraction of optimum. The point-in-polygon kernel was not implemented for the MI50 GPU.