9th International Conference in Deep Learning in Computational Physics

ML-Based Optimum Sub-system Size Heuristic for the GPU Implementation of the Tridiagonal Partition Method

Milena Veneva

 $\checkmark \boxtimes$ milena.p.veneva@gmail.com \mathbb{R} Milena_Veneva in mveneva

 \checkmark RIKEN Center for Computational Science, supervisor **Dr. Toshiyuki Imamura**

 \checkmark in collaboration with the Joint Institute for Nuclear Research, supervisor $\mbox{Dr.}$ Alexander Ayriyan

July 2025, Moscow, Russia.

Tridiagonal partition method (0/3)

 $0. \ \mbox{Starting from the initial SLAE}$ with the following coefficient matrix.

$(\times$	\times														(x_0)		(y_0)
×	\times	×													x_1		y_1
	\times	×	×												x_2		y_2
		×	×	×											x_3		<i>y</i> 3
			×	\times	\times										X 4		y 4
				\times	\times	×									x_5		y_5
					\times	×	×								x_6		y 6
						\times	×	\times							<i>x</i> ₇	=	<i>y</i> 7
							×	\times	\times						x_8		<i>y</i> 8
								\times	\times	\times					X 9		y 9
									\times	\times	\times				<i>x</i> ₁₀		y ₁₀
										\times	\times	\times			<i>x</i> ₁₁		y ₁₁
											\times	\times	×		x_{12}		y ₁₂
												\times	×	×	x ₁₃		y 13
$\langle \rangle$													×	×/	$\langle x_{14} \rangle$		y_{14}

[1] Austin T. M., Berndt M., and Moulton J. D. A memory efficient parallel tridiagonal solver, pp.1–13, Preprint

Tridiagonal partition method (1a/3)

1a. Partitioning.

$(\times$	×														(x_0)		(y_0)
×	\times	×													x_1		y 1
	\times	×	×												x_2		y_2
		×	×	\times											x_3		<i>y</i> 3
			×	\times	\times										x_4		y ₄
				×	×	×									x_5		y 5
					×	\times	×								x_6		<i>у</i> 6
						\times	×	\times							<i>x</i> ₇	=	y 7
							×	\times	\times						x_8		y 8
								\times	\times	×					X 9		y 9
									×	×	×				<i>x</i> ₁₀		<i>y</i> ₁₀
										×	\times	\times			<i>x</i> ₁₁		y ₁₁
											×	×	\times		x_{12}		y ₁₂
												\times	\times	×	<i>x</i> ₁₃		y 13
													\times	×/	$\langle x_{14} \rangle$		y_{14}

Tridiagonal partition method (1b/3)

1b. Obtaining the interface equations (on the device).



2. Assembling the interface system and solving it with the Thomas method (on the host).



3. Substituting the already found unknowns corresponding to the \boxtimes coefficients, and solving the rest of the tridiagonal sub-systems (on the device).



- on the basis of NVIDIA GPU RTX 2080 Ti;
- SLAE sizes: 10^i , 2×10^i , 4×10^i , 5×10^i , 8×10^i , i = 2, 3, ..., 7, and 4.5×10^3 , 2.5×10^4 , 3×10^4 , 6×10^4 , 7×10^4 , 7.5×10^4 ;
- Parameters: 256 CUDA threads within a block were used; number of CUDA streams: according to [2]; precision FP64, _NO_WD interface option; no recursions; no overwriting of RHS;
- very many different sub-system sizes were tested;
- event synchronization when collecting the times.

[2] Veneva M., Imamura T., ML-Based Optimum Number of CUDA Streams for the GPU Implementation of the Tridiagonal Partition Method, arXiv:2501.05938 [cs.DC] (2025), to be published.

Considerations

- Consideration about blockSize: there should be enough warps in a block so as to keep the SM busy while one warp is waiting for resources (e.g. loading data from memory). Choosing the blockSize to be 256 gives the best performance. BlockSize smaller than 128 is not enough to supply the SM with enough active warps, and blockSize bigger than 256 requires too much resources (registers and shared memory), therefore the performance starts deteriorating.
- Balance between latency hiding, resource utilization, and occupancy.
- Within the partition method, the gridSize depends on the sub-system size *m*.

[3] NVIDIA, CUDA C++ Best Practices Guide, Release 12.9,

https://docs.nvidia.com/cuda/cuda-c-best-practices-guide/ (2025).

Existing heuristics

- Exhaustive search over all possible combinations of parameters (e.g. QUDA), but this requires additional runs of the application and is energy-costly. This approach is beneficial for applications that are run many times with the same parameters.
- Some promote a performance characteristic hoping that in the worst case scenario this approach is not going to lead to pathologically bad performance (e.g. occupancy in the case of Thrust).
- Some use two-step approach so as to gather information about the GPU limitations, and the kernel that needs to be launched (e.g. KLARAPTOR).
- Some use machine learning techniques to predict the combinations of parameters that would lead to optimal or near optimal performance (e.g. [4]).

[4] Sato K., Takizawa H., Komatsu K., and Kobayashi H.. Automatic tuning of CUDA execution parameters for

stencil processing, In: Software Automatic Tuning, From Concepts to State-of-the-Art Results (2010).

Comparison between the achieved and the theoretical occupancy



Thus, the occupancy cannot be our main reference point.

SLAE size	opt m	#str	comments	corrected opt m
10 ²	4	1	_	4
2×10^{2}	4	1	-	4
4×10^{2}	4	1	-	4
5×10^{2}	4	1	-	4
8×10^{2}	4	1	-	4
10^{3}	4	1	-	4
2×10^{3}	4	1	-	4
4×10^{3}	4	1	-	4
4.5×10^{3}	4	1	-	4
5×10^{3}	8	1	-	8
8×10^{3}	8	1	-	8
10^{4}	8	1	-	8
2×10^{4}	8	1	-	8
2.5×10^{4}	8	1	-	8
3×10^{4}	16	1	-	16
4×10^{4}	16	1	-	16
5×10^{4}	16	1	-	16
6×10^{4}	20	1	-	20
7×10^{4}	35	1	small diff with the time when $m = 20$: 0.0008099	20
7.5×10^{4}	40	1	small diff with the time when $m = 20$: 0.00719	20
8×10^{4}	32	1	-	32
105	40	1	small diff with the time when $m = 32$: 0.000621	32
2×10^{5}	64	2	small diff with the time when $m = 32$: 0.068788	32
4×10^{5}	64	4	small diff with the time when $m = 32$: 0.073638	32
5×10^{5}	40	8	small diff with the time when $m = 32$: 0.045666	32
8×10^{5}	64	8	diff with the time when $m = 32$: 0.182118	32
100	32	8	-	32
2×10^{6}	32	16	-	32
4×10^{6}	32	32	-	32
5×10^{6}	32	32	-	32
8×10^{6}	64	32	small diff with the time when $m = 32$: 0.44834	32
107	32	32	-	32
2×10^{7}	64	32	=	64
4×10^{7}	64	32	=	64
$5 \times 10^{\prime}$	64	32	=	64
$8 \times 10'$	64	32	=	64
108	64	32	-	64

ML model for the optimum sub-system size (1/2)

- k-nearest neighbors (kNN) classification (independent variable: SLAE size; dependent (target variable): the optimum sub-system size).
- k was found empirically by using scikit-learn tool GridSearchCV (which looks for the best combination of hyper-parameters, in particular we looked for the best k (between 1, and the number of unique sub-system sizes), and best weight – 'uniform' or 'distance') to be equal to 1 (and 'uniform' weight).
- The data was split into two data sets training (that is, neighbours), and test (for which we make predictions) using the Python scikit-learn routine train_test_split with shuffle option turned on, and splitting ratio 3 : 1.

ML model for the optimum sub-system size (2/2)

- Two approaches when fitting the data: (1) using the experimentally found optimum sub-system sizes, and (2) using the corrected sub-system sizes which takes into account the trend we have noticed.
- The former approach gave us 0.7 normalised accuracy score for the test set, which means that it manages to find the expected sub-system size in 7 out of 10 cases.
- On the other hand, using the latter idea, we get 1.0 normalised accuracy score, that is, 100% accuracy.
- The null accuracy (accuracy that could be achieved by always predicting the most frequently met sub-system size) was found to be 0.4.
- The results are in a good correspondence with the memory alignments considerations (aligned to at least 256 bytes).

Results from the kNN classification model



(a) Results from the kNN sub-system size using the corrected data for m.

(b) Results from the kNN classification model for the optimum classification model for the optimum sub-system size using the observed data for *m*.

Changing from sub-optimal to optimal (bigger) sub-system size might give us speed-up up to 1.70 (for number_of lines $= 8 imes 10^7$ compare m = 64 and m = 4).

Results for FP32



(a) Results from the kNN classification model for the optimum sub-system size using the corrected data for *m*.

(b) Results from the kNN classification model for the optimum sub-system size using the observed data for *m*.

- 0.8 normalised accuracy score for the test set for the observed data,
- 1.0 normalised accuracy score for the test set for the corrected data,
- null accuracy of 0.4.

SLAE size	#str	opt m	heuristic	opt m	diff with	opt m	diff with
		2080 Ti	on 2080 Ti	A5000	the heuristic	4080	the heuristic
10^{2}	1	4	4	4	-	4	-
2×10^2	1	4	4	4	-	4	-
4×10^2	1	4	4	4	-	4	-
5×10^2	1	4	4	4	- 1	4	-
8×10^2	1	4	4	4	- 1	8	small
10^{3}	1	4	4	4		4	
2×10^3	1	4	4	4		4	
4×10^3	1	4	4	8	small	8	small
4.5×10^3	1	4	4	4		4	
5×10^3	1	8	8	4	small	4	small
8×10^3	1	8	8	8		4	small
10^{4}	1	8	8	8		8	-
2×10^4	1	8	8	8		16	small
2.5×10^4	1	8	8	8		8	
3×10^4	1	16	16	16	-	16	-
4×10^4	1	16	16	16		16	
5×10^4	1	16	16	16	-	16	-
6×10^{4}	1	20	20	32	2.65%	40	small
7×10^4	1	35	20	20	-	20	-
7.5×10^4	1	40	20	20	-	40	small
8×10^4	1	32	32	40	small	32	-
10^{5}	1	40	32	32	-	32	-
2×10^5	2	64	32	64	6.26%	64	4.59%
4×10^{5}	3	64	32	64	3.54%	64	small
5×10^5	8	40	32	40	2.38%	40	4.19%
8×10^5	8	64	32	64	6.03%	64	2.50%
10^{6}	8	32	32	64	9.44%	64	7.13%
2×10^{6}	16	32	32	64	8.15%	64	6.00%
4×10^{6}	32	32	32	64	5.60%	64	6.90%
5×10^{6}	32	32	32	64	3.65%	64	5.66%
8×10^{6}	32	64	32	64	5.63%	64	7.09%
107	32	32	32	64	6.06%	64	6.75%
2×10^{7}	32	64	64	64	-	64	-
4×10^7	32	64	64	64	-	64	-
5×10^7	32	64	64	64	-	64	-
8×10^7	32	64	64	64	-	64	-
10^{8}	32	64	64	64		64	

Thank you for your attention!

Acknowledgements: R-CCS team (RIKEN), Dr. Toshiyuki Imamura (RIKEN), Dr. Alexander Ayryan (JINR).