





Machine Learning and GPU Accelerated Sparse Linear Solvers for Transistor-Level Circuit Simulation: A Perspective Survey

ASP-DAC 2024

Zhou Jin, Wenhao Li, Yinuo Bai, Tengcheng Wang, Yicheng Lu, Weifeng Liu Super Scientific Software Laboratory, China University of Petroleum-Beijing, China Email: jinzhou@cup.edu.cn

Contents

- Background
- Sparse Linear Solvers for Circuit Simulation
 - SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
 - Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
 - PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems
- Conclusions

Contents

Background

Sparse Linear Solvers for Circuit Simulation

- SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
- Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
- PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems

Conclusions

Transistor-Level Circuit Simulation

Transistor-level circuit simulation (SPICE simulation) plays a crucial role in verifying circuit performance, and serving as the basis of timing, yield and reliability analysis, etc.



Design process of analog circuits

Challenges in SPICE simulation: the continuous expansion of scale, making it extremely time-consuming during the IC design process.



Transistor-Level Circuit Simulation

> SPICE workflow example

Read the netlist and

establish the circuit

equation

 G_{\cdot}

-G

D4 0 2 1k

 $-G_{1}-k$

 $G_1 + G_2 + G_3$

 $\left[e_{1} \right]$

 $||e_2| = ||0|$

-G.

 $-G_3 + k$ $G_3 + G_4 - k \parallel e_3$

The most timeconsuming step in SPICE simulation is solving linear equations Ax = b.



LU Factorization

LU factorization is part of the direct method to solve the linear equations and is also one of the most important tasks in many scientific computing applications. The following figure shows an example of LU factorization:

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ l_{21} & 1 & 0 & 0 \\ l_{31} & l_{32} & 1 & 0 \\ l_{41} & l_{42} & l_{43} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} & u_{14} \\ 0 & u_{22} & u_{23} & u_{24} \\ 0 & 0 & u_{33} & u_{34} \\ 0 & 0 & 0 & u_{44} \end{bmatrix}$$

The elements of *L* and *U* are calculated from the following two equations:

$$u_{ij} = a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \qquad \begin{cases} i = 1, 2, \dots, N\\ j = i, i+1, \dots, N \end{cases} \qquad l_{ij} = \frac{1}{u_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} u_{kj} \right) \begin{cases} i = 1, 2, \dots, N\\ j = 1, 2, \dots, i-1 \end{cases}$$

Sparse LU Factorization

When the input matrix is sparse, the sparse LU factorization can be used to reduce unnecessary computation and storage, which is called sparse LU factorization.



consuming

Development of Direct Sparse Solver

MUMPS is a distributed direct solver using the multifrontal method. It was developed by Patrick and Duff et al. (PARA00) SuperLU using supernodal method is developed by Li and Demmel et al. SuperLU_DIST is a distributed version. (TOMS03) SuperLU_DIST has been continuously updated and added versions for GPU.(JPDC19) Only on share memory CPUs

On Single GPU

On share memory CPUs and single GPU

On distributed systems



Development of Direct Sparse Solver



Circuit Matrix Properties



Circuit sparse matrix properties.

			Entri	es per row		
Circuit Matrix	Ν	max	min	average	variation	Symmetry
ASIC_320k memchip Freescale1 circuit5M G2_circuit transient	321,821 2,707,524 3,428,755 3,523,317 150,102 178,866	203,800 27 27 27 4 60423	1 2 1 1 1 1 1	8.2 5.5 5.5 10.7 2.9 5.4	502.95 2.06 2.07 1356.61 0.52 147.2	100.00% 0.32% 7.67% 55.99% 0.0005% 68.99%

In circuit simulation, the matrix usually

has the following properties:

Development of Direct Sparse Solver



PanguLU is developed by SSSLab. It uses strategies such as regular two-dimensional blocking, synchronization-free and sparse kernels on distributed heterogeneous platform. (SC23 Best paper)



Contents

Background

Sparse Linear Solvers for Circuit Simulation

- SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
- Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
- PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems
- Conclusions

Using GPU for general computing has become a research hotspot.

Name	Strategy	Characteristic	
SuperLU ^[1]	Running GPU-version dense matrix-matrix multiplication on (often small) matrices generated from supernodal pattern.	Can not use GPU for scheduling to save the	
GLU ^[2]	Finding dependencies between columns, generating levels, and using level scheduling method for running multiple columns in the same level on GPUs concurrently.	cost for kernel relaunching and cache data flushing.	
?	?	Can use GPU for scheduling.	

TABLE I: There are two ways to utilize GPUs for sparse LU.

We need to design a method that can use GPU for scheduling.

[1] Demmel, James W, Eisenstat, et al, A supernodal approach to sparse partial pivoting. SIAM Journal on Matrix Analysis & Applications, 1999.

[2] Wai-Kong L, Ramachandra A, Nakhla M S, Dynamic GPU Parallel Sparse LU Factorization for Fast Circuit Simulation, IEEE Transactions on Very Large Scale Integration (VLSI) Systems, 2018



Due to the logic between symbolic and numeric factorization are similar. we only use the symbolic factorization to introduce our algorithm.

Specifically:





We test 1309 matrices from the SuiteSparse matrix set on Titan RTX (Turing) GPUs, Intel 20-core CPUs, comparing the performance of SFLU with SuperLU and GLU.

In the symbolic factorization, the acceleration is up to 205.14x and 701x compared to SuperLU and GLU, respectively.

In the numeric factorization, the acceleration is up to 3585.2x and 252.2x compared to SuperLU and GLU, respectively.

Contents

Background

Sparse Linear Solvers for Circuit Simulation

- SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
- Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
- PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems
- Conclusions

The numeric factorization in supernodal LU factorization follows four steps, where K signifies the K-th iteration and N denotes the number of matrix blocks on the diagonal.

- Factorize the diagonal block.
 - contains a large number of GEMN

The Schur-complement phase

- Factorize the sub-matrices in L panel: L(K : N, K) operations.
- Factorize the sub-matrices in U panel: U(K, K +1 : N)
- Perform the Schur-complement for all the tailing sub-matrices by using A = A-L × U.

GEMM takes up much of the time for numeric factoriztaion.

We tested some circuit matrices, such as the G3_circuit matrix in the figure, which accounts for as much as 73.4% of the time, and most of the other matrices tend to account for 40%-60% of the GEMM time.

In order to better **utilize the dense GEMM of Level-3 BLAS**, the traditional approach is to form dense blocks using supernodal (shown in Figure 1) or multifrontal methods. However, for matrices without a regular sparse structure, **supernodes** formed through relaxation **may not always be dense** (shown in Figure 2).

(a) Original supernodes (b) Supernodes after merging

Schur-complement

 $A(I,J) \leftarrow A(I,J) - L(I,K) \times U(K,J);$

Figure 2. Example of Schurcomplement computation.

We select three representative circuit matrices and can see that the input matrix (L and U blocks) involved in GEMM are not always dense.

Figure 1. Density distribution of matrix factors (L and U blocks) participating in GEMM.

The introduction of sparse matrix multiplication (SpMM) shows great potential for further accelerating the LU factorization.

After introducing SpMM, we test the data in Table 1. The "Oracle" in the table indicates the optimal total time computed by picking the optimal kernel in SpMM and GEMM for each matrix multiplication in the ideal case.

Circuit matrix	nnz (A)	GEMM (s)	SpMM (s)	Oracle (s)	Speedup1	Speedup2
ASIC_320k	1,931,828	3.5809	0.4543	0.3653	9.80x	1.24x
Freescale1	17,052,626	8.8363	35.9429	8.5388	1.03x	4.21x
ckt11752_dc_1	333,029	0.0279	0.0433	0.0225	1.24x	1.92x
pre2	5,834,044	57.2954	10.9046	7.2531	7.90x	1.50x
meg4	58,142	0.0037	0.0027	0.0022	1.68x	1.23x
G2_circuit	726,674	- 1	\sim	115	1.16x	4.25x
Freescale2	14,313,7	Which m	ethod is		1.19x	1.95x
FullChip	26,62				1.48x	1.40x
ASIC_320ks	1,	best? G	EMM?		10.24x	1.13x
ASIC_680ks	1,69	SoM	N/2		2.81x	1.44x
circuit5M_dc	14,80	Spivi	IVI:	0230	6.04x	1.21x
transient	961,368		.0	0.2100	2.69x	1.29x

TABLE I: Analysis of the performance improvement space of matrix multiplication.

Compared with GEMM and SpMM, the performance improvement potential of "Oracle" is 1.03x-10.24x and 1.13x-4.25x, respectively.

We need to pick the optimal kernel for each execution of matrix multiplication!

It is difficult to define thresholds for either density or matrix size. The use of **density-aware adaptive** matrix multiplication equipped with random forest^[1] provides an effective solution to the issues mentioned above.

Figure 2. Density-aware LU factorization.

- Dataset Source: The large-scale circuit matrices in SuiteSparse and exported by SuperLU_DIST 8.0.0.
- Samples in the dataset: Each sample contains 15 matrix features (F1-F15) and a label P, where P=1 indicates that using GEMM is better than that of SpMM, and P=0 is vice versa.
- Data preprocessing:
 - Z-score normalization
 - Sample equalization
- Number of samples in the dataset :

Dataset	GEMM	SpMM
Total sample set(D1)	456,000	383,000
Training set(D2)	150,000	150,000
Testing set(D3)	306000	233000

An example of 15 matrix features (F1-F15).

We test several matrix in a 2 * Intel Xeon Silver 4210 CPU @ 2.20GHz, 512GB DDR4 platform configuration. We test 11 test matrix from the SuiteSparse, including 6 circuit matrix and 5 non-circuit matrix.

Figure 1. Model accuracy

For circuit matrix, up to **9.35x** acceleration for matrix mutiplication and **1.76x** for LU factorization;

For other irregular matrix, up to 4.32x and 2.08x, respectively.

Performance evaluation on circuit matrices.										
				Matrix	multiplicat	ion time (s)		Numeric factorization time (s)		
circuit matrix	nnz	GEMM	SpMM	Oracle	AI	Speedup (AI vs GEMM)	Speedup (AI vs SpMM)	SuperLU	Our work	Speedup
ASIC_320k	1,931,828	3.5809	0.4543	0.3653	0.5045	7.10x	0.90x	7.5201	4.5770	1.64x
ASIC_320ks	1,316,085	2.9155	0.3223	0.2846	0.3117	9.35x	1.03x	6.0602	3.5947	1.69x
ASIC_680ks	1,693,767	2.6196	1.3393	0.9323	1.0085	2.60x	1.33x	5.3502	3.8139	1.40x
circuit5M_dc	14,865,409	7.5022	1.2418	1.0231	2.0084	6.04x	0.62x	19.2301	14.5605	1.32x
pre2	5,834,044	57.2954	10.9046	7.2531	8.1021	7.07x	1.35x	103.5721	58.7290	1.76x
transient	961,368	0.5721	0.2709	0.2121	0.2532	2.26x	1.07x	1.7202	1.4517	1.18x
Average	-	-	-	-	-	5.35x	1.05x	-	-	1.50x

Performance evaluation on non-circuit matrices.

		Matrix multiplication time (s)						Numeric factorization time (s)		
non-circuit matrix	nnz	GEMM	SpMM	Oracle	AI	Speedup (AI vs GEMM)	Speedup (AI vs SpMM)	SuperLU	Our work	Speedup
sinc12	283,992	9.7211	2.7060	2.0120	2.2491	4.32x	1.20x	14.7541	7.4277	2.04x
psmigr_3	543,160	16.3840	7.8951	5.8241	6.2611	2.62x	1.26x	28.4921	19.8387	1.54x
psmigr_2	540,022	30.8151	12.0731	8.6251	9.1721	3.36x	1.32x	45.9061	24.5057	2.08x
epb2	175,027	0.1291	0.10021	0.05631	0.07621	1.69x	1.31x	0.4351	0.3937	1.08x
benzene	242,669	8.1778	9.6211	6.2315	7.2724	1.13x	1.32 x	17.6351	16.9551	1.04x
Average	-	-	-	-	-	2.62x	1.28x	-	-	1.55x

Contents

Background

Sparse Linear Solvers for Circuit Simulation

- SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
- Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
- PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems
- Conclusions

Motivation 1: Uneven Block Sizes

When the supernodal method aggregates columns of the same structure, there are different sizes of supernodes.

It can be seen that the matrix blocks generated by the supernodes can be very irregular, which often affects the computational efficiency, and their irregular structure makes it difficult to optimise the performance at the kernel level.

PanguLU: Sparse Solver on Distributed Heterogeneous Systems (SC23)

Motivation 2: Redundant Zero Fill-ins

The supernodal method divides the matrix into uneven dense blocks to be computed with dense BLAS, which may cause two questions.

1. Redundant zero fill-ins can occur when forming dense blocks that add extra floating-point operations.

2. The local sparsity of the matrix cannot be exploited during GEMM

compared with dense BLAS.

PanguLU: Sparse Solver on Distributed Heterogeneous Systems (SC23)

Motivation 3: High Synchronisation Costs

SuperLU_DIST uses the level-set method to generate an elimination tree with tree nodes as the minimum scheduling unit in distributed systems, which can lead to high synchronization costs.

- We test six real-world matrices from SuiteSparse Matrix Collection on 64 NVIDIA A100 GPUs.
- The synchronization cost would be more than 50% at 64 processes on irregular matrices.
- We need a new method to reduce synchronization costs.

Data Layout of PanguLU

Our approach is completely different from the classic supernodal method.

PanguLU splits the matrix into multiple blocks of equal size and uses sparse kernel to calculate.

Each process uses a two-layer sparse structure to store the matrix.

First layer sparse structure stores the positional information of sparse matrix blocks. We will get the position of sub-matrix in this layer sparse format.

Second layer sparse structure stores the non-zero elements of sub-matrix block with sparse format. We will get internal information about the block, including rows, non-zeros, and sub-matrix structure.

Sparse Kernels of PanguLU

In the numeric factorization of PanguLU, since the matrix blocks are sparse, we develop four dedicated sparse kernels: sparse general triangular factorization (GETRF), sparse upper triangular solve (TSTRF), sparse lower triangular solve (GESSM), and Schur complement with sparse-sparse matrix multiplication (SSSSM).

L * X = A - L * U + C = C

All four sparse kernels use sparse sub-matrix blocks as inputs and outputs.

Sparse Kernels of PanguLU

There are several critical factors for the performance of numeric factorzation, such as density, structure and size of the matrix. To exploit the best performance, we implement 17 sparse kernels in PanguLU.

Kernel	Versi on	Addressin g Method	Paralleling Method	Dense Mappi ng	4 $ \begin{array}{c} 4 \\ \bullet \\ C_V1 \\ \bullet \\ G_V1 \\ \bullet \\ \bullet \\ G_V2 \end{array} $ 4 4 $ \begin{array}{c} C_V1 \\ \bullet \\ C_V1 \\ \bullet \\ C_V2 \\ \bullet \\ G_V3 \\ \bullet \\ G_V1 \end{array} $
	C_V1	Direct	Row	\checkmark	
GETR F	G_V1	Bin-search	Un-sync SFLU	-	
	G_V2	Direct	Un-sync SFLU	\checkmark	
	C_V1	Merge	Column	-	
	C_V2	Direct	Column	\checkmark	
TSTR F/ GESS	G_V1	Bin-search	Warp-level column	-	nnz(A) (log ₁₀ scale) nnz(B) (log ₁₀ scale) (a) GETRF (c) GESSM
M	G_V2	Bin-search	Un-sync warp-level row	-	$4 C_V 1 G_V 2 \qquad 4 C_V 1 G_V 1 \qquad 4 C_V 2 $
	G_V3	Direct	Warp-level column	\checkmark	
	C_V1	Direct	Approximate equal load column block	√	
SSSS M	C_V2	Bin-search	Adaptive split-bin type	-	o million o
	G_V1	Bin-search	Adaptive multi-level	-	
	G_V2	Direct	Warp-level column	\checkmark	$ \frac{1}{2} - \frac{2}{2} + \frac{4}{6} = \frac{6}{2} - \frac{2}{2} + \frac{6}{6} = \frac{6}{2} + 6$

Different sparse kernels have different

(d) SSSSM

How to select one of these sparse kernels at runtime?

The "Addressing Method" indicates how the calculation value is located and updated. In addition, the "Dense Mapping" represents that the sparse structure is mapped to a dense space.

ASP-DAC 2024

(b) TSTRF

Sparse Kernels of PanguLU

We propose an algorithm selection strategy for constructing sparse kernels based on large amounts of performance data to select better sparse kernels.

We develop four decision trees according to #non-zeros of matrix A/B and the FLOPs involved in the computation from the performance of 17 sparse kernels to guide our algorithm selection.

Mapping of PanguLU

PanguLU develops a mapping method and redistributes the computational load between processes. We map sparse block of heavy-load process to less-load process according to kernel types and time slices.

Synchronisation-Free Scheduling

PanguLU are computed sequentially in the order of time slices, to maintain the correctness of the sparse LU factorisation as show in figure.

Synchronisation-Free Scheduling

PanguLU are computed sequentially in the order of time slices, to maintain the correctness of the sparse LU factorisation as show in figure.

But obviously these time slices can be computed in parallel on different processes.

Executive direction

The synchronisation between each time slice is required to ensure that the correctness of the calculations.

Is there another method of calculation that can guarantee correctness while keeping the processes in a working state as much as possible?

Synchronisation-Free Scheduling

So we researched the order of computation of sparse LU factorisation using sparse kernels.

In particular, we do not use any synchronisation in here.

Experimental Setup

Two experimental platforms	Two heterogeneous distributed Solvers
4 * NVIDIA A100 GPUs (40 GB, B/W 1555GB/s) 2 * Intel Xeon 8180 CPUs @ 2.5 GHz (512GB DDR4)	SuperLU_DIST 8.1.2 PanguLU 3.5
4 * AMD MI50 GPUs (16 GB, B/W 1024GB/s) 1 * AMD Epyc 7601 CPU @ 2.2 GHz (128GB DDR4)	SuperLU_DIST 8.1.2 PanguLU 3.5

Compared with SuperLU_DIST, PanguLU achieves an average speedup of 2.53x and 2.79x on the NVIDIA GPU platform and the AMD GPU platform, with speedups ranging from 1.10x - 11.70x and 1.12x - 17.97x

In this figure, the red lines and blue lines the performance of SuperLU_DIST and respectively, while the solid lines and do indicate the performance on the A100 G MI50 GPU platforms, respectively. Especially, PanguLU achieves significant performance advantages for irregular matrices such as ASIC_680k (circuit matrix) with speedups of up to 11.70x and 17.97x on the two GPU platforms, respectively.

PanguLU achieves better performance by using sparse kernel on a single GPU. In the 16 matrices, PanguLU has a geometric mean of 6.54x compared with SuperLU.

Matrix	Panel Fa	ict Time (s)	Schur	Time (s)	All Time (s)			
IVIATIX	SuperLU	PanguLU	SuperLU	PanguLU	SuperLU	PanguLU	Speedup	
apache2	9.4	6.32	27.36	3.86	36.76	10.18	3.61x	
ASIC_680k	3.39	3.21	982.36	17.8	985.76	21.01	46.92x	
audikw 1	103.44	47.94	764.83	159.35	868.27	207.29	4.19x	
cage12	15.8	13	3570.76	70.65	3586.56	83.65	42.88x	
CoupCons3D	15.69	11.11	55.55	8.5	71.24	19.61	3.63x	
dielFilterV3real	41.62	24.63	127.92	25.27	169.54	49.9	3.40x	
ecology1	2.26	1.8	2.88	0.57	5.14	2.37	2.17x	
G3_circuit	5.69	4.98	14.71	1.35	20.4	6.34	3.22x	
Ga41As41H72	979.85	80.91	21893.96	1275.41	22873.81	1356.32	16.86x	
Hook_1498	117.77	66.3	1208.79	230.17	1326.56	296.47	4.47x	
inline_1	10.15	6.15	19.66	2.04	29.82	8.19	3.64x	
ldoor	7.2	4.55	13.87	1.05	21.07	5.59	3.77x	
nlpkkt80	173.69	80.48	1816.94	537.19	1990.63	617.68	3.22x	
Serena	238.6	115.98	3551.09	899.36	3789.69	1015.35	3.73x	
Si87H76	152.21	74.88	17245.73	1034.96	17397.94	1109.84	15.68x	
SiO2	77.78	38.57	9762.91	410.6	9840.69	449.17	21.91x	
Geometric Mean		-		-		-	6.54x	

Kernel time comparison of SuperLU_DIST and PanguLU for 16 test matrices on a signle A100 GPU.

Contents

Background

Sparse Linear Solvers for Circuit Simulation

- SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs
- Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation
- PanguLU: A Scalable Regular Two-Dimensional Block-Cyclic Sparse Direct Solver on Distributed Heterogeneous Systems

Conclusions

Sparse Linear Solvers for Circuit Simulation

Machine Learning and GPU Accelerated Sparse Linear Solver (on CPU, GPU and Distributed Heterogeneous Systems):

- Jianqi Zhao, Yao Wen, Yuchen Luo, Zhou Jin, Weifeng Liu, Zhenya Zhou. SFLU: Synchronization-Free Sparse LU Factorization for Fast Circuit Simulation on GPUs. The 58th ACM/IEEE Design Automation Conference (DAC' 21). 2021.
- https://gitee.com/ssslab/sflu
- Tengcheng Wang, Wenhao Li, Haojie Pei, Yuying Sun, Zhou Jin, Weifeng Liu. Accelerating Sparse LU Factorization with Density-Aware Adaptive Matrix Multiplication for Circuit Simulation. The 60th ACM/IEEE Design Automation Con ference (DAC' 23). 2023
- <u>https://github.com/SuperScientificSoftwareLaboratory/DALU</u>
- Xu Fu, Bingbin Zhang, Tengcheng Wang, Wenhao Li, Yuechen Lu, Enxin Yi, Jianqi Zhao, Xiaohan Gen g, Fangying Li, Jingwen Zhang, Zhou Jin, Weifeng Liu. PanguLU: A Scalable Regular Two-Dimensional BlockCyclic Sparse Direct Solver on Distributed Heterogeneous Systems. 36th ACM/IEEE International Conference for High Performance Computing, Networking, Storage, and Analysis (SC' 23). 2023. (Best paper award)
- https://www.ssslab.cn/software.html

An overview of surveyed sparse direct solver

	Solver	Left/Right looking	Blocking method	Kernel			Parallelism		
					Level	Method	Distributed	Multi- thread	GPU
	MUMPS	Left	Multifrontal	Level-3 BLAS	Tree/Node	Level-set	\checkmark	\checkmark	—
	UMFPAC K	Left	Multifrontal	Level-3 BLAS	—	—		—	_
Gene-	SuperLU_ DIST	Right	Supernode	Level-3 BLAS	Supernode	Level-set	\checkmark	\checkmark	\checkmark
Tai	PARDISO	Left & Right	Supernode	Level-3 BLAS	Supernode	Pipeline	\checkmark	\checkmark	—
	PanguLU	Right	Regular 2D block	Adaptive sparse kernel	2D block	Synchronizat ion-free	\checkmark		\checkmark
	KLU	Left	Block diagonal	—	—	—	—	—	—
Dedi-	NICSLU	Left	Supernode	Level-3 BLAS	Supernode	Level-set	—	\checkmark	\checkmark
cated to	FLU	Left	Supernode	Level-3 BLAS	Supernode	Register Level	—	\checkmark	—
circuit	GLU	Left & Right	—	—	Element	Level-set	—	—	\checkmark
matrix	SFLU	Left & Right	—	—	Element	Synchronizat ion-free	—	—	\checkmark
	Density- aware LU	Right	Supernode	Adaptive kernel	Supernode	Level-set	\checkmark	\checkmark	—

Some Observations for Future

The future opportunities in this evolving landscape:

- Taking into account various factors such as the matrix and hardware platform characteristics in algorithm design.
 - How to introduce sparse computations and strike a trade-off between dense and sparse kernels for different computational tasks remains a pivotal challenge.
 - Better task scheduling strategies and computational kernels need to be designed on GPUs to improve the parallel computing performance.

Welcome to collaborate with us!

Thanks!

Zhou Jin

jinzhou@cup.edu.cn