

# Nested Dissection Based Parallel Transient Power Grid Analysis on Public Cloud Virtual Machines

## Jiawen Cheng, Zhiqiang Liu, Lingjie Li, Wenjian Yu Department of Computer Science and Technology, Tsinghua University, Beijing 100084, China

2024-01-24

# Outline

#### Background

- Nested Dissection Based Parallel Transient Power Grid Analysis with Cloud Computing
- Experimental Results
- Conclusion

# Integrated Circuit (IC) Design Flow

- Backend related EDA technology
  - Physical design optimization (P&R, power delivery network)
  - Physical modelling analysis (RC extraction, delay analysis)
  - Signoff verification (signal integrity, circuit simulation)
  - Layout verification
  - Mask optimization
- Challenge of accurate circuit simulation
  - Parasitic effect from interconnects with widths < 0.25µm</p>
  - Extremely **time consuming** due to large scale



# **Power Grid Analysis**

- Compute the IR drop under various current load
- Critical for analyzing the robustness of the circuit
- Challenge
  - circuit scale /, IR drop margin /, demand for accuracy /
- Efficient, accurate, memory-saving software
- Classification
  - Direct current: only consider resistors
  - Transient: consider resistors, capacitors and inductors (RLC) accurate  $\sqrt{}$



# **Transient Power Grid Analysis**

- Modeled as an RLC network, a first-order ODE system is derived  $C\dot{x}(t) + Gx(t) = Bu(t)$  (1)
  - C: capacitance and inductance matrix
  - G: conductance matrix u(t): input sources
  - B: input selector matrix x(t): node voltages branch currents

• Apply trapezoidal rule with time step h, (1) becomes Ax = b problem!  $\left(\frac{C}{2} + \frac{G}{2}\right)x(t+h) = \left(\frac{C}{2} - \frac{G}{2}\right)x(t) + B\frac{u(t+h) + u(t)}{2}\int$ (2)

# Solving Linear Equation Systems

	Direct: LU/Cholesky + Fwd/Bwd Substitution	Iterative: Preconditioning + CG/GMRES
Pros	Fast substitution, accurate	Variable time step, memory saving
Cons	Fixed time step, fill-ins	Accuracy loss, convergence issue

#### Motivation

- Design an efficient, accurate and robust solver
- Accelerate solving and reduce memory cost
- Public cloud computing platforms make distributed computing feasible
- Internal structure of the circuit won't be exposed ⇒ Public cloud computing

 $\Rightarrow$  Direct method

 $\Rightarrow$  Distributed computing

# Domain Decomposition Method (DDM)

■ Symmetric matrix⇔Undirected graph

D

- Use the edge/vertex separator to partition the graph
- Arranging subdomains and separator together

$$\begin{pmatrix} D_1 & F_1 \\ D_2 & F_2 \\ \vdots \\ D_m & F_m \\ E_1 & E_2 & \cdots & E_m & S \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \\ y \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_m \\ g \end{pmatrix} \implies \frac{\text{Schur complement}}{[S - ED^{-1}F]y} = g - ED^{-1}f \\ D_i x_i = f - F_i y$$



Edge separator



Vertex separator

- Solve  $x_i$  in parallel (3) Solving y is bottleneck when m is large (3)
- Learn from nested dissection to reduce the size of Schur complement

# Our contribution

#### Efficient, accurate and robust transient PG analysis

- Multi-level distributed parallel LU factorization and forward/backward substitution approach based on nested dissection
- Scheme for efficient RHS updating to further exploit parallelism

#### Leverage cloud computing for distributed parallelism

- Use public cloud computing while ensuring security
- Ideal for deployment on cloud computing platforms without InfiniBand

# Outline

- Background
- Nested Dissection Based Parallel Transient Power Grid Analysis with Cloud Computing
- Experimental Results
- Conclusion

# **Overall Idea – LU Factorization**

• Top level of partition:  $\begin{pmatrix} D_1 \\ T \end{pmatrix}$ 

$$egin{pmatrix} egin{pmatrix} egi$$

• Solve the submatrices in *L*, *U* 

$$L_{31} = E_1 U_{11}^{-1}, L_{32} = E_2 U_{22}^{-1}$$
(3)

$$\boldsymbol{U}_{13} = \boldsymbol{L}_{11}^{-1} \boldsymbol{F}_1, \boldsymbol{U}_{23} = \boldsymbol{L}_{22}^{-1} \boldsymbol{F}_2 \tag{4}$$

$$L_{33}U_{33} = S - L_{31}U_{13} - L_{32}U_{23}$$
(5)

- Factorize  $D_1, D_2$  recursively in the same manner  $\Rightarrow$  nested dissection
- Opportunities for parallelization: (3), (4) and factorization of  $D_1$ ,  $D_2$
- Acceptable cost of solving (5) on a single process

# **Overall Idea – Forward/Backward Substitution**

Forward substitution

$$\begin{pmatrix} L_{11} \\ L_{22} \\ L_{31} & L_{32} & L_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ y \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ g \end{pmatrix} \implies L_{11}x_1 = f_1, L_{22}x_2 = f_2, L_{33}y = g - L_{31}x_1 - L_{32}x_2$$

- Parallel solving  $x_1, x_2$  in the same manner recursively
- Acceptable cost of solving y on a single process
- Backward substitution is similarly done from the top down

# Parallel Reordering into Nested Dissection Form

#### Notations

- $l_m$ : max nesting level  $2^{l_m}$ : number of processes
- Superscript (*r*): data resides in process *r*
- Subscript *l*: data describes information of level *l*
- Task assignment for process *r* in level *l* 
  - Level  $l_m$ : partition the full matrix if r = 0
  - Level 0: compute fill-in reducing reordering for  $D^{(r)}$
  - Others: partition the subdomain from level l 1 if r divides  $2^l$
- Record additional information for subsequent solving
  - *f*, *g*, *h*: offset of subdomain 1, subdomain 2, separator in *A*
  - e: end of submatrix



#### Level 0

• 
$$D^{(r)} = L_0^{(r)} U_0^{(r)} \Rightarrow$$
 Each process r factorizes  $D^{(r)}$  to get  $L_0^{(r)}, U_0^{(r)}$ 



#### Level 1

•  $E_1^{(r)} = L_1^{(r)} U_0^{(r)}$ ,  $F_1^{(r)} = L_0^{(r)} U_1^{(r)} \Rightarrow$  Each process r solves triangular equations to get  $L_1^{(r)}$ ,  $U_1^{(r)}$ 



#### Level 1

- $E_1^{(r)} = L_1^{(r)} U_0^{(r)}$ ,  $F_1^{(r)} = L_0^{(r)} U_1^{(r)} \Rightarrow$  Each process r solves triangular equations to get  $L_1^{(r)}$ ,  $U_1^{(r)}$
- $S_1^{(r)} = L_1^{(r+1)} U_1^{(r+1)} + L_1^{(r)} U_1^{(r)} + \tilde{L}_1^{(r)} \tilde{U}_1^{(r)} \Rightarrow$  Even-numbered process r factorizes the Schur Matrices sent from r + 1 to r complement to get  $\tilde{L}_1^{(r)}$ ,  $\tilde{U}_1^{(r)}$



#### • Level l > 1

•  $E = L_l^{(r)} U, F = L U_l^{(r)} \Rightarrow$  Each process r dividing  $2^{l-1}$  solves triangular equations to get  $L_l^{(r)}, U_l^{(r)}$ 



#### • Level l > 1

- $E = L_l^{(r)} U, F = L U_l^{(r)} \Rightarrow$  Each process r dividing  $2^{l-1}$  solves triangular equations to get  $L_l^{(r)}, U_l^{(r)}$
- $S_l^{(r)} = L_l^{(r+2^{l-1})} U_l^{(r+2^{l-1})} + L_l^{(r)} U_l^{(r)} + \tilde{L}_l^{(r)} \tilde{U}_l^{(r)} \Rightarrow$  Each process r dividing  $2^l$  factorizes the Schur Matrices sent from  $r + 2^{l-1}$  to r complement to get  $\tilde{L}_l^{(r)}, \tilde{U}_l^{(r)}$

### Parallel Forward/Backward Substitution

- Level 0
  - $L_0^{(r)}y_0^{(r)} = b_0^{(r)} \Rightarrow$  Each process *r* solves triangular equations to get  $y_0^{(r)}$

$$\begin{pmatrix} \begin{pmatrix} \mathbf{L}_{0}^{(3)} & & \\ \mathbf{L}_{1}^{(3)} & \mathbf{L}_{1}^{(2)} & \\ \mathbf{L}_{1}^{(3)} & \mathbf{L}_{1}^{(2)} & \tilde{\mathbf{L}}_{1}^{(2)} \end{pmatrix} \\ & & \begin{pmatrix} \mathbf{L}_{0}^{(1)} & & \\ & \mathbf{L}_{0}^{(1)} & \\ & \mathbf{L}_{1}^{(1)} & \mathbf{L}_{1}^{(0)} & \tilde{\mathbf{L}}_{1}^{(0)} \end{pmatrix} \\ & & \mathbf{L}_{2}^{(2)} & & \mathbf{L}_{2}^{(0)} \end{pmatrix} \begin{pmatrix} \mathbf{y}_{0}^{(3)} & \\ \mathbf{y}_{0}^{(2)} & \\ \mathbf{y}_{1}^{(1)} & \\ \mathbf{y}_{0}^{(0)} & \\ \mathbf{y}_{0}^{(0)} & \\ \mathbf{y}_{2}^{(0)} \end{pmatrix} = \begin{pmatrix} \mathbf{b}_{0}^{(3)} & \\ \mathbf{b}_{0}^{(2)} & \\ \mathbf{b}_{0}^{(1)} & \\ \mathbf{b}_{0}^{(0)} & \\ \mathbf{b}_{0}^{(0)} & \\ \mathbf{b}_{2}^{(0)} \end{pmatrix} \\ & & L & * \mathbf{y} = \mathbf{b}$$

## Parallel Forward/Backward Substitution

Level 
$$l > 0$$
Level  $l > 0$ 
 $L_l^{(r+2^{l-1})} y_l^{(r+2^{l-1})} + L_l^{(r)} y_l^{(r)} + \tilde{L}_l^{(r)} y_l^{(r)} = b_l^{(r)} \Rightarrow$  Each process  $r$  dividing  $2^l$  solves triangular equations to get  $y_l^{(r)}$ 

$$\begin{pmatrix} \begin{pmatrix} L_{0}^{(3)} & & \\ & L_{0}^{(2)} \\ & & L_{1}^{(2)} & \tilde{L}_{1}^{(2)} \end{pmatrix} \\ & & \\$$

## Parallel Forward/Backward Substitution

- Backward substitution is similarly done from the top down
- Properties of the processes
  - Own a **continuous** segment of solution vector ⇒ **Communication friendly**
  - Read and write the **same** segment of solution vector ⇒ **RHS updating friendly**

$$\begin{pmatrix} \begin{pmatrix} L_{0}^{(3)} & & \\ & L_{0}^{(2)} \\ & & L_{1}^{(3)} & L_{1}^{(2)} & \tilde{L}_{1}^{(2)} \end{pmatrix} & & \\ & & \begin{pmatrix} L_{0}^{(1)} & & \\ & & L_{0}^{(0)} \\ & & L_{1}^{(1)} & L_{1}^{(0)} & \tilde{L}_{1}^{(0)} \end{pmatrix} & & \\ & & & L_{2}^{(2)} & & L_{2}^{(0)} \end{pmatrix} \begin{pmatrix} y_{0}^{(3)} \\ & y_{0}^{(2)} \\ & y_{1}^{(1)} \\ & y_{0}^{(0)} \\ & y_{0}^{(0)} \\ & y_{1}^{(0)} \\ & y_{2}^{(0)} \end{pmatrix} = \begin{pmatrix} b_{0}^{(3)} \\ & b_{0}^{(2)} \\ & b_{0}^{(1)} \\ & b_{0}^{(0)} \\ & b_{0}^{(0)} \\ & b_{1}^{(0)} \\ & b_{2}^{(0)} \end{pmatrix} \\ & & L & * & \mathcal{Y} = b$$

# Scheme for Efficient Right-Hand Side Updating

• 
$$b = \left(\frac{C}{2} - \frac{G}{2}\right)x(t) + B\frac{u(t+h) + u(t)}{2} \triangleq Mx + Bt$$

- Fwd/bwd substitution read and write the same segment of solution vector
- Optimization strategy
  - No need to store the full *M*, *B*
  - Smaller SpGEMV and only 1 Allgatherv ⇒ Increase parallelism
  - Allocate shared memory for solution vectors on the same computing node

 $\Rightarrow$  Reduce communication overhead

Process 3

:

Process p - 2

Process p-1

Process p

 $\boldsymbol{v}^{(r)}$ 

\_

# $\frac{u(t)}{\Delta Mx + Bu}$ and write $\frac{time t \quad Allgatherv \quad time t + 1}{Process 1}$ $\frac{Process 1}{Process 2}$

Process 3

Process p - 2

Process p-1

Process p

 $M^{(r)}$ 

Process 3

Process p-2

Process p-1

Process p

 $\boldsymbol{u}^{(r)}$ 

 $\Rightarrow$  Reduce memory overhead

×

Identical on

all processes

u

# Outline

- Background
- LU Factorization & Forward/Backward Substitution Based on Nested Dissection
- Experimental Results
- Conclusion

# **Experimental Configurations**

#### Environment

- 4-node Huawei Cloud cluster with Ethernet connectivity
- Each node with a 32-core Intel Xeon Platinum 8378A CPU and 256GB RAM
- Test data: ibmpg3-6t, thupg1,3,5,8,10t (500 time points of 10ps each)
- Implementation: OpenMPI + Eigen + MKL + METIS
- Methods compared to: NICSLU, edge separator based DDM
- Allocate 2 threads to each process for DDM and our method

# Results on a Single Computing Node

#### Notations

- $N_{th}$ : # of threads  $T_s$ : total time of serial solving
- $T_f$ , T: time of factorization, total time
- $N_s$  : size of Schur complement
- $\overline{N}_s$ : average size of Schur complement in level C
- Sp1,2: speedup compared to NICSLU, DDM
- Analysis
  - All fail to complete the simulation of thupg8,10t
  - 2.06X speedup over NICSLU
  - 1.43X speedup over DDM

	Case	$N_{th}$	NICSLU		DDM			Proposed				
	Case		$T_s$	T	$T_{f}$	Т	$N_S$	$T_{f}$	Т	$\overline{N}_S$	Sp1	Sp2
	ibmpg3t N = 1.0E6	4	74.4	47.8	9.07	24.8	0	8.88	25.9	0	1.84	0.96
		8		40.7	5.60	18.5	2325	6.08	16.3	561	2.49	1.13
		16		32.2	5.43	23.1	6288	4.25	13.1	377	2.45	1.76
	nnz = 4.5 E0	32		34.8	5.47	27.3	10860	3.83	10.6	261	3.30	2.58
	ibmna4t	4	94.5	56.2	12.6	34.4	0	12.3	34.5	0	1.63	1.00
	$N = 1.2 \Sigma c$	8		43.1	7.82	24.9	3334	9.13	25.0	723	1.72	0.99
	N = 1.2E0	16		37.3	6.51	27.2	6410	5.55	16.4	335	2.28	1.66
	nnz = 0.8E0	32		41.9	8.43	39.1	13112	4.79	14.6	353	2.86	2.67
	ibmna5t	4	75.5	50.1	11.1	32.5	540	11.43	39.2	0	1.28	0.83
	N = 1  GFG	8		40.1	6.87	23.4	2469	6.87	23.6	330	1.70	0.99
	N = 1.0E0	16		36.0	5.78	23.7	4670	5.78	16.7	204	2.15	1.41
	nnz = 0.4E0	32		43.8	5.25	26.9	9388	5.25	13.9	294	3.15	1.93
	il	4	103	68.9	12.8	40.0	0	12.8	39.2	0	1.76	1.02
	Iompgot N 2 4 FC	8		53.2	7.64	27.8	1864	7.64	30.4	266	1.75	0.91
	N = 2.4 E0 nnz = 9.7 E6	16		48.3	5.23	23.7	4178	5.23	22.9	296	2.11	1.04
		32		61.6	4.51	22.6	7309	4.51	20.6	167	2.99	1.10
	thupg1t N = 5.3E6 nnz = 2.3E7	4	458	276	66	176	2953	92.5	193	1401	1.43	0.92
		8		221	87.1	185	8792	66.0	142	1404	1.55	1.30
		16		187	69.4	166	15064	44.6	95.4	661	1.96	1.74
		32		200	51.6	170	25981	36.5	80.5	697	2.49	2.12
	thupg3t N = 1.2E7	4	1243	761	213	487	4646	323	610	2270	1.25	0.80
		8		587	374	611	14110	316	493	2268	1.19	1.24
		16		539	218	453	23615	168	300	994	1.80	1.51
	1112-0.411	32		599	188	496	41340	138	252	1034	2.38	1.97
	thupg5t N = 2.0F7	4		1932	533	1013	6777	626	1061	3234	1.82	0.95
		8	2102	1203	511	915	17513	425	695	2445	1.73	1.32
	10 = 2.0E1	16	2193	1058	482	891	30558	327	539	1576	1.96	1.65
	nn2 - 0.0E1	32		1211	522	1068	51244	263	445	1210	2.72	2.40
	Average	-	-	-	-	-	-	-	-	-	2.06	1.43

# **Results on Multiple Computing Nodes**

- 64, 128 processes assigned to 2, 4 computing nodes respectively
- $V_f$ ,  $V_{tr}$ : communication volume of factorization, fwd&bwd substitution

#### Analysis

- Size of Schur complement effectively controlled
- Efficient
  - 2.84X speedup over DDM > 1.43X on a single node
  - High parallel scalability
- Accurate and robust
  - Relative error of solution compared to NICSLU < 1e-7</p>
- Low communication volume

				- D D I								
	Case	$N_{th}$	DDM			Proposed						
			$T_f$	T	$N_S$	$T_f$	T	$\overline{N}_{S}$	$V_f$	$V_{tr}$	Sp	
	ibmpg3t	64	6.67	36.0	17508	3.96	10.9	162	0.50	0.06	3.29	
		128	10.2	56.1	26646	4.16	11.5	126	0.93	0.07	4.88	
lled	ibmpg4t	64	8.15	45.5	19966	4.47	13.6	196	0.65	0.07	3.35	
		128	13.1	64.8	32215	5.01	14.0	172	1.19	0.09	4.63	
	ibmpg5t	64	5.57	30.8	14865	5.57	13.9	159	0.76	0.09	2.22	
		128	10.5	51.5	24295	5.84	14.4	138	1.11	0.12	3.57	
	ibmpg6t	64 128	5.73	30.4	14441	5.75	18.2	163	0.77	0.13	1.6/	
		128	9.98	48.4	22794	0.15	18.9	224	1.50	0.18	2.50	
	thupg1t	128	40.2	212	38348 59929	31.8	75.0	324	5.14	0.31	2.33	
е	thupg3t	120	47.0	472	50050	35.2	10.5	515	3.33	0.56	$\frac{2.71}{2.10}$	
		128	197	541	01613	108	208	489	13 5	0.08	2.10	
		64	201	863	77062	224	301	763	12.9	113	2.00	
	thupg5t	128	371	987	117742	209	368	594	22.3	1.47	2.68	
	thupg8t	()	950	0714	10(470	000	1175	1040	07.4	0.00	0.01	
	N = 4.1 E7	64 128	859	2/14	106472	806	11/5	1049	27.4	2.33	2.31	
	nnz = 1.8 E8	128	801	2875	142819	/50	1089	8/4	4/.1	2.97	2.04	
	thupg10t											
~ 7	N = 6.3 E7	128	1239	4375	181934	1058	1695	1034	73.0	4.57	2.58	
8-1	nnz = 2.7E8											
	Average	-	-	-	-	-	-	-	-	-	2.84	
	The notations of $N_{th}$ , $T_{f}$ , $T$ , $N_{S}$ , $\overline{N}_{S}$ and $S_{F}$ are the same as Table I. V.											
	denotes the communication volume of matrix reordering and factorization											
	and $V_{tr}$ denotes the average communication volume of forward/backward											
	substitution and RHS updating per time point.								05			

# Parallel scalability

- X-axis: # of threads Y-axis: speedup compared to serial NICSLU
- No improvement of speedup of NICSLU and DDM
- Increasing speedup of our method, achieving up to 6.0X using 4 computing nodes
- Cloud computing allows for faster simulation
- Cost of adding computing nodes may be less than the benefit of saving time



# Conclusion

- Leverage the public cloud computing for transient PG analysis while ensuring security
- Multi-level distributed parallel LU factorization and forward/backward substitution approach based on nested dissection
- Scheme for efficient RHS updating to further exploit parallelism
- Ideal for deployment on cloud computing platforms without InfiniBand
- 2.06X speedup compared to NICSLU on a single computing node
- 2.84X speedup compared to DDM on multiple computing nodes
- Good parallel scalability up to 6.0X speedup over serial simulation



# Thank you! Q & A

#### Jiawen Cheng 2024-01-24