

MAX PLANCK INSTITUTE FOR DYNAMICS OF COMPLEX TECHNICAL SYSTEMS MAGDEBURG



COMPUTATIONAL METHODS IN SYSTEMS AND CONTROL THEORY

Mixed precision HODLR matrices Erin Carson¹, Xinye Chen², Xiaobo Liu³

¹Department of Numerical Mathematics, Charles University, Czech Republic (carson@karlin.mff.cuni.cz). ²LIP6, Sorbonne University, CNRS, France (Xinye.Chen@lip6.fr). ³Max Planck Institute for Dynamics of Complex Technical Systems, Germany (xliu@mpi-magdeburg.mpg.de).

Hierarchical Off-Diagonal Low-Rank matrix



HODLR matrix-vector product

Core idea: balance the approximation error in \hat{H} (mixed-precision HODLR) representation) and the finite-precision computation error in $b \leftarrow Hx$.

Theorem (Backward error). If b = Hx is computed recursively, exploiting the HODLR structure, in working precision $u \leq \varepsilon/n$, then the computed b satisfies

 $\widehat{b} = \mathrm{fl}(\widehat{H}x) = (H + \Delta H)x, \quad \|\Delta H\|_F \leq 2(\sqrt{2} + 1)\sqrt{2^{\ell+1} + 2^{\ell-1}\varepsilon} \|H\|_F.$

• Intuition: more inexact the low-rank representation, lower the precision one can safely use.

■ Numerical experiments on three *kernel matrices* of size *n* = 2000,

Off-diagonal blocks at each level stored as low-rank ($\leq p$) matrix outer product UV^{T} , to achieve

- 1. $\mathcal{O}(pn \log n)$ storage requirement: O(np) each level $\times O(\log n)$ levels;
- 2. arithmetic operation cost of $\mathcal{O}(p^{\alpha}n\log^{\beta}n)$, $\alpha, \beta \in \{1, 2\}$, for matrix summation, multiplication, inversion, factorizations, etc.

Applications: fast solvers for elliptic PDEs, kernel methods in machine learning, integral equations in computational physics, etc.

New definition: (p, ε) **-HODLR matrix**

Definition. The off-diagonal blocks of the (p, ε) -HODLR matrix H (associated) with a HODLR matrix H) satisfy $\|\widetilde{H}_{ii}^{(k)} - H_{ii}^{(k)}\| \le \varepsilon \|H_{ii}^{(k)}\|$ at any level k, where $0 \le \varepsilon < 1$ is the low-rank approximation threshold.

• Rank-constrained \rightarrow tolerance-constrained HODLR format, to facilitate error analysis of the low-rank approximations.

Mixed-precision storage scheme



$$K_{ij} = \begin{cases} \frac{1}{x-y}, & \text{if } x \neq y; \\ 1, & \text{otherwise.} \end{cases}, \quad K_{ij} = \begin{cases} \log \|x_i - x_j\|_2, & \text{if } x \neq y; \\ 0, & \text{otherwise.} \end{cases}, \quad K_{ij} = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2}\right)$$

evaluated at point sets s_1 (1D uniform grid points in [0, 1]) and $s_2 = s_3$ (2D) uniform grid points in $[-1, 1] \times [-1, 1]$:



Backward error in different working precisions. Depth $\ell = 8$. The *x*-axis indicates the value of ε .

HODLR LU factorization

- HODLR LU factorization $H^{(k)} \approx LU$ schematic: \approx Key steps in the recursive algorithm: 1: Partition $H^{(k)}$ into $\begin{vmatrix} H_{11}^{(k+1)} & H_{12}^{(k+1)} \\ H_{21}^{(k+1)} & H_{22}^{(k+1)} \end{vmatrix}$ 2: $L_{11}, U_{11} \leftarrow \text{HODLR_LU}(H_{11}^{(k+1)}, k+1)$ 3: $U_{12} \leftarrow$ Solve triangular system $L_{11}U_{12} = H_{12}^{(k+1)}$ 4: $L_{21} \leftarrow$ Solve triangular system $L_{21}U_{11} = H_{21}^{(k+1)}$ 5: $H_{22}^{\varepsilon} \leftarrow H_{22}^{(k+1)} - L_{21}U_{12}$ (possible rank-*p* truncation)

Norm distributions among layers for two HODLR matrices with depth $\ell = 6$.

Core idea: reduced precision $u_k \ge u$ for storing the low-rank factors U and V^{T} of off-diagonal blocks in the kth layer.

- Choose $u_k \leq \varepsilon/(2^{k/2}\xi_k)$, where $\xi_k := \max_{i \neq j} \|\widetilde{H}_{ii}^{(k)}\|_F / \|\widetilde{H}\|_F$ characterizes the relative importance of the off-diagonal blocks in magnitude.
- With $\varepsilon = 10^{-4}$, Use {bf16, fp16, fp16, fp32, fp32, fp32} for sayIr3 and {q52, fp32, fp32, fp32, fp32, fp32} for LeGresley 2508.

Numerical experiments on Schur complements of *SuiteSparse* matrices:



6: $L_{22}, U_{22} \leftarrow \text{HODLR}_{LU}(H_{22}^{\varepsilon}, k+1)$ 7: $L \leftarrow \begin{bmatrix} L_{11} \\ L_{21} & L_{22} \end{bmatrix}, U \leftarrow \begin{bmatrix} U_{11} & U_{12} \\ U_{22} \end{bmatrix}$

• At the bottom level $k = \ell$, dense LU factorization of $H_{ii}^{(\ell)}$ is computed. **Theorem** (Backward error). If the LU decomposition of \hat{H} is computed in a working precision $u \leq \varepsilon/n$, then

 $\widehat{L}\widehat{U} = H + \Delta H, \quad \|\Delta H\|_F \lesssim 2(2^\ell - 1)\varepsilon \|H\|_F + 11(2^\ell - 1)\varepsilon \|\widehat{L}\|_F \|\widehat{U}\|_F.$

• For larger ε (a coarser approximation), the LU factorization can be computed in a lower precision without affecting the backward error.

• Bound remains in the same form for LU factorization of one-precision HODLR matrices, albeit with smaller constant factors.

Numerical experiments on the Schur complement of ex37:



Backward error in different working precisions. The x-axis indicates the value of ε .

Further reading

E. Carson, X. Chen, and X. Liu. Mixed precision HODLR matrices.







