

Scalable Tensor Algorithms for Modern Computing Systems

Suraj KUMAR

CNRS LIP/LaBRI Applicant

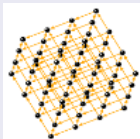
March 3, 2022

My Past Experience

Parallelization in Polyhedral Model

(IISc, 2012)

- Linked-list operations
- Improved spatial locality
- Parallelization using OpenMP



Seismic Imaging on GPU (IBM, 2013)

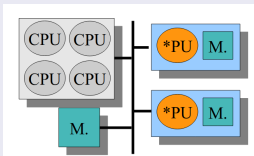
$$H_1 = \sin^2 \theta \cos^2 \phi \frac{\partial^2}{\partial x^2} + \sin^2 \theta \sin^2 \phi \frac{\partial^2}{\partial y^2} + \cos^2 \theta \frac{\partial^2}{\partial z^2} + \sin^2 \theta \sin 2\phi \frac{\partial^2}{\partial x \partial y} + \sin 2\theta \sin \phi \frac{\partial^2}{\partial y \partial z} + \sin 2\theta \cos \phi \frac{\partial^2}{\partial x \partial z}$$
$$H_2 = \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} + \frac{\partial^2}{\partial z^2} - H_1$$

Schedulers for Blue Gene Supercomputers (IBM, 2013)

- GASNET API
- Unbalanced Tree Search benchmark
- Comparison to Charm++

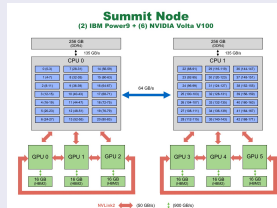


Scheduling on Heterogeneous Platforms (Inria, 2017)



Molecular Simulations on Supercomputers (PNNL, 2019)

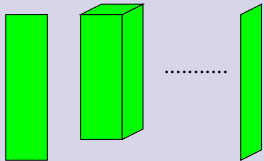
- NWChemEx Project
- Tamm library
- Hartree Fock and CCSD applications



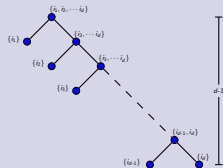
My Past Experience

Parallel Tensor Train Approximation (Inria, current)

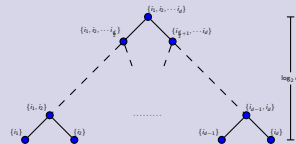
Small object representation



Sequential algorithm

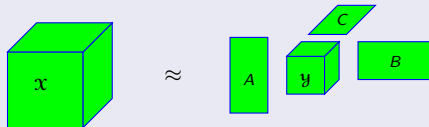


For better parallelization



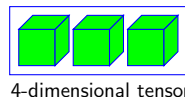
Communication Optimal Parallel Algorithms for Tensor Computations (Inria, current)

- Obtain \mathcal{Y} from \mathcal{X}, A, B, C
- Obtain \mathcal{X} from \mathcal{Y}, A, B, C



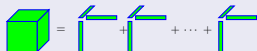
Tensors and Importance of Communication

- **Neuroscience:** Neuron \times Time \times Trial
- **Media:** User \times Movie \times Time
- **Ecommerce:** User \times Product \times Time
- **Social-Network:** Person \times Person \times Time \times Type



- High dimensional tensors: Neural network, Molecular simulation, Quantum computing
- People work with low dimensional structure (decomposition) of the tensors

Canonical decomposition



Tucker decomposition



Tensor Train decomposition

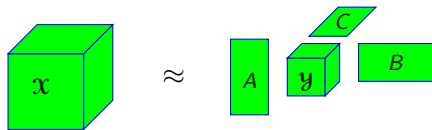


- Gaps between computation and communication growing exponentially on a parallel computer

	time-per-operation	Network-bandwidth	Network-latency
Annual improvements	59 %	26 %	15 %

Source: Getting up to speed: The future of supercomputing

Higher-order SVD (HOSVD) to compute Tucker decomposition



Algorithm 1 HOSVD Algorithm(\mathcal{X} , R_1 , R_2 , R_3)

- 1: $A \leftarrow R_1$ left singular vectors of $\mathcal{X}_{(1)}$
 - 2: $B \leftarrow R_2$ left singular vectors of $\mathcal{X}_{(2)}$
 - 3: $C \leftarrow R_3$ left singular vectors of $\mathcal{X}_{(3)}$
 - 4: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
 - 5: Return \mathcal{Y} , A , B , C
-

- \mathcal{X} , \mathcal{Y} : 3-dimensional input and output tensors (or arrays) & A , B , C : matrices
- $\mathcal{X}_{(i)}$: matricization of \mathcal{X} (i th dimension represents rows and remaining dimensions represent columns)
- Multiple Tensor-Times-Matrix (Multi-TTM) computation: $\mathcal{Y} = \mathcal{X} \times_1 A^T \times_2 B^T \times_3 C^T$
 - To obtain full tensor, $\mathcal{X} = \mathcal{Y} \times_1 A \times_2 B \times_3 C$

Communication Lower Bounds and Communication Optimal Algorithms

- 1 For Matrix Matrix Multiplications
- 2 For Multi-TTM Computation

Assumptions

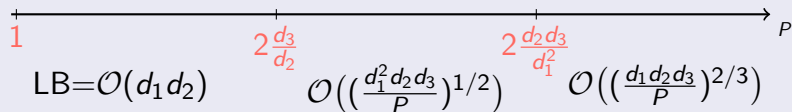
- P number of processors
- Each processor performs (asymptotically) equal amount of operations
- No redundant operations
- One copy of data is in the system
 - $1/P$ th amount of inputs (before the computation) and output (after the computation) on each processor
- Each processor has enough memory

This is joint work with Laura Grigori (Inria Paris, France), Grey Ballard (Wake Forest University, USA), Kathryn Rouse (Inmar Intelligence, USA), and Hussam Al Daas (Rutherford Appleton Laboratory, UK).

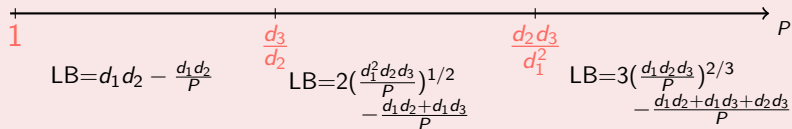
Existing Lower Bounds for Matrix Matrix Multiplications

- $C = AB$, where $A \in \mathbb{R}^{n_1 \times n_2}$, $B \in \mathbb{R}^{n_2 \times n_3}$, and $C \in \mathbb{R}^{n_1 \times n_3}$
- Let $d_1 = \min(n_1, n_2, n_3) \leq d_2 = \text{median}(n_1, n_2, n_3) \leq d_3 = \max(n_1, n_2, n_3)$

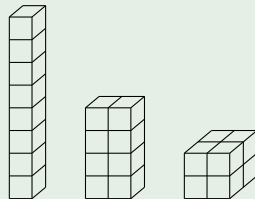
Existing Communication Lower Bounds (CARMA [IPDPS 2013])



Our Communication Lower Bounds



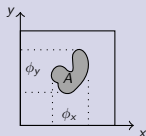
Arrangements of 8 processors



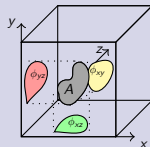
Loomis-Whitney & Hölder-Brascamp-Lieb inequalities

Size of $d - 1$ dimensional projections (Loomis-Whitney inequality)

- 2-dimensional object A and its 1-dimensional projections ϕ_x, ϕ_y
- $\phi_x \phi_y \geq \text{Area}(A)$



- 3-dimensional object A and its 2-dimensional projections: $\phi_{xy}, \phi_{yz}, \phi_{xz}$
- $(\phi_{xy} \phi_{yz} \phi_{xz})^{\frac{1}{3-1}} \geq \text{Volume}(A)$



Hölder-Brascamp-Lieb (HBL) inequality – Generalization of Loomis-Whitney inequality

$$\Delta = \begin{matrix} & A & B & C \\ \begin{matrix} i \\ j \\ k \end{matrix} & \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 0 \end{pmatrix} \end{matrix}$$

for $i = 1:n_1$, for $k = 1:n_2$, for $j = 1:n_3$

$$C[i][j] + = A[i][k] * B[k][j]$$

- Find $\mathbf{x} = [x_1 \ x_2 \ x_3]^T$ such that $\Delta \cdot \mathbf{x} \geq \mathbf{1}$, $\mathbf{1}$ is vector of all ones
- ϕ_A, ϕ_B, ϕ_C : projections of computations on arrays A, B, C
- HBL inequality: $\phi_A^{x_1} \phi_B^{x_2} \phi_C^{x_3} \geq \text{Amount of computations}$
- To make inequality tight select \mathbf{x} such that $\mathbf{1}^T \mathbf{x}$ is minimum $\Rightarrow x_1 = x_2 = x_3 = \frac{1}{2}$

Constraints for Matrix Multiplications

for $i = 1:n_1$, for $k = 1:n_2$, for $j = 1:n_3$

$$C[i][j] += A[i][k] * B[k][j]$$

- Total number of multiplications = $n_1 n_2 n_3$
- Consider a processor which performs $\frac{n_1 n_2 n_3}{P}$ amount of multiplications
- Optimization problem:

Minimize $\phi_A + \phi_B + \phi_C$ s.t.

$$\phi_A^{\frac{1}{2}} \phi_B^{\frac{1}{2}} \phi_C^{\frac{1}{2}} \geq \frac{n_1 n_2 n_3}{P}$$

Extra constraints (our contributions)

- Each element of A (resp. B) is involved in n_3 (resp. n_1) multiplications
 - To perform at least $\frac{n_1 n_2 n_3}{P}$ multiplications: $\phi_A \geq \frac{n_1 n_2}{P}, \phi_B \geq \frac{n_2 n_3}{P}$
- Each element of C is the sum of n_2 multiplications, therefore $\phi_C \geq \frac{n_1 n_3}{P}$
- Projections can be at max the size of the arrays: $\phi_A \leq n_1 n_2, \phi_B \leq n_2 n_3, \phi_C \leq n_1 n_3$

Optimization Problem to Compute Communication Lower Bounds

- Projections (ϕ_A, ϕ_B, ϕ_C) indicate the amount of array access
- Communication lower bound = $\phi_A + \phi_B + \phi_C$ – data owned by the processor

Minimize $\phi_A + \phi_B + \phi_C$ s.t.

$$\phi_A^{\frac{1}{2}} \phi_B^{\frac{1}{2}} \phi_C^{\frac{1}{2}} \geq \frac{n_1 n_2 n_3}{P}$$

$$\frac{n_1 n_2}{P} \leq \phi_A \leq n_1 n_2$$

$$\frac{n_2 n_3}{P} \leq \phi_B \leq n_2 n_3$$

$$\frac{n_1 n_3}{P} \leq \phi_C \leq n_1 n_3$$

Generalized version (in terms of d_1, d_2, d_3)

Minimize $\phi_1 + \phi_2 + \phi_3$ s.t.

$$\phi_1^{\frac{1}{2}} \phi_2^{\frac{1}{2}} \phi_3^{\frac{1}{2}} \geq \frac{d_1 d_2 d_3}{P}$$

$$\frac{d_1 d_2}{P} \leq \phi_1 \leq d_1 d_2$$

$$\frac{d_1 d_3}{P} \leq \phi_2 \leq d_1 d_3$$

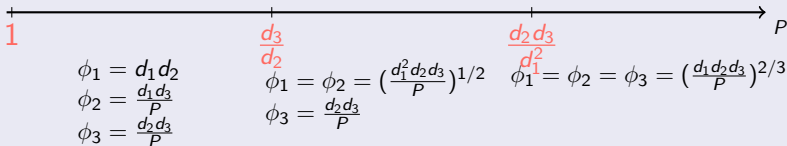
$$\frac{d_2 d_3}{P} \leq \phi_3 \leq d_2 d_3$$

$$d_1 \leq d_2 \leq d_3$$

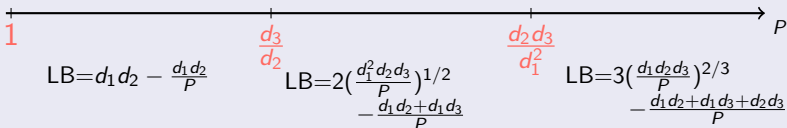
Amount of Accesses and Communication Lower bounds

- Estimate the solution based on Lagrange multipliers
- Prove optimality using all KarushKuhnTucker (KKT) conditions are satisfied

Amount of accesses $= \phi_1 + \phi_2 + \phi_3$



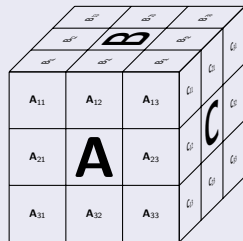
Communication Lower Bounds (Amount of Data Transfers)



Design of Communication Optimal Algorithms

Data Distribution (P is organized into a $p_1 \times p_2 \times p_3$ grid)

- Select p_1, p_2 , and p_3 based on the lower bounds
- Each processor has $\frac{1}{p}$ th amount of A, B and C
- $A_{11} = A(1 : \frac{n_1}{p_1}, 1 : \frac{n_2}{p_2})$ is evenly distributed among $(1, 1, *)$ processors
- Similar data distributions for B and C



Algorithm 2 $C = AB$ Matrix Multiplication Algorithm

- 1: (p'_1, p'_2, p'_3) is my processor id
- 2: //All-gather input matrices A and B
- 3: $A_{p'_1 p'_2} = \text{All-Gather}(A, (p'_1, p'_2, *))$
- 4: $B_{p'_2 p'_3} = \text{All-Gather}(B, (*, p'_2, p'_3))$
- 5: $T = \text{Local-Matrix-Multiplication}(A_{p'_1 p'_2}, B_{p'_2 p'_3})$ // Local matrix multiplication in a temporary
- 6: $\text{Reduce-Scatter}(C_{p'_1 p'_3}, T, (p'_1, *, p'_3))$ // Reduce-scatter the output

3-dimensional Multi-TTM ($\mathcal{Y} = \mathcal{X} \times_1 \mathbf{A}^{(1)\top} \times_2 \mathbf{A}^{(2)\top} \times_3 \mathbf{A}^{(3)\top}$)

- TTM-in-Sequence approach (used in Tucker-MPI)
 - For 2-dimensional computation, $\mathbf{Y} = \mathbf{A}^{(1)\top} \mathbf{X} \mathbf{A}^{(2)}$
- $\mathcal{X} : n_1 \times n_2 \times n_3$, $\mathcal{Y} : r_1 \times r_2 \times r_3$, $\mathbf{A}^{(k)} : n_k \times r_k$

All-at-Once approach (our contribution)

for $n'_1 = 1:n_1$, for $n'_2 = 1:n_2$, for $n'_3 = 1:n_3$

for $r'_1 = 1:r_1$, for $r'_2 = 1:r_2$, for $r'_3 = 1:r_3$

$\mathcal{Y}(r'_1, r'_2, r'_3) = \mathcal{Y}(r'_1, r'_2, r'_3)$

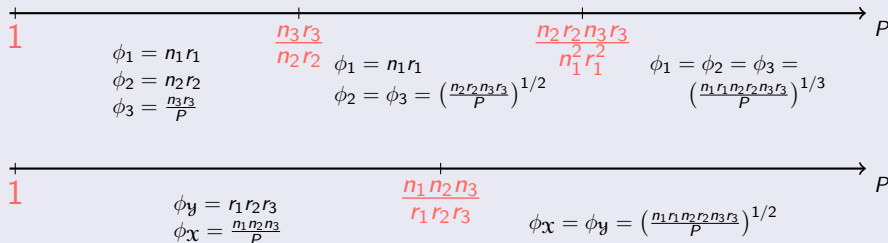
$+ \left(\mathcal{X}(n'_1, n'_2, n'_3) * \mathbf{A}^{(1)}(n'_1, r'_1) * \mathbf{A}^{(2)}(n'_2, r'_2) * \mathbf{A}^{(3)}(n'_3, r'_3) \right)$

$$\Delta = \begin{bmatrix} \mathbf{I}_{3 \times 3} & \mathbf{1}_3 & \mathbf{0}_3 \\ \mathbf{I}_{3 \times 3} & \mathbf{0}_3 & \mathbf{1}_3 \end{bmatrix}$$

- Total number of inner (4 – array) operations = $n_1 r_1 n_2 r_2 n_3 r_3$
- Δ is not full rank: allows us to get multiple constraints related to computations
- Possible to solve matrix and tensor optimization problems separately

Amount of Accesses and Lower bounds

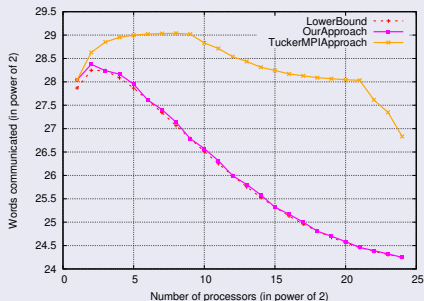
Amount of accesses = $\phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3$



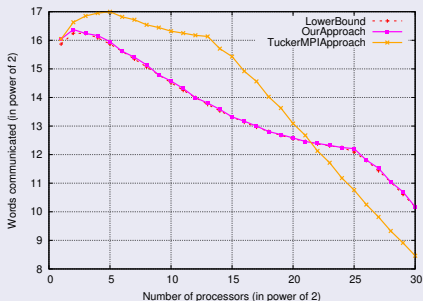
- We assume $n_1 r_1 \leq n_2 r_2 \leq n_3 r_3$ and $r_1 r_2 r_3 \leq n_1 n_2 n_3$
- Communication lower bound = $\phi_x + \phi_y + \phi_1 + \phi_2 + \phi_3 - \frac{n_1 n_2 n_3 + r_1 r_2 r_3 + n_1 r_1 + n_2 r_2 + n_3 r_3}{P}$
- Can design similar communicational optimal algorithm (though 6-dimensional) for this
- Selection of optimal processor grid dimensions based on the lower bound requires some adaption

Simulated Performance Comparison of Our Algorithm

$$n_1 = n_2 = n_3 = 2^{20}, r_1 = r_2 = r_3 = 2^8$$



$$n_1 = n_2 = n_3 = 2^{12}, r_1 = r_2 = r_3 = 2^4$$



- Typical scenarios in data compression problems
- Lower Bound is only valid for our approach
- For $P \ll \frac{n_1 n_2 n_3}{r_1 r_2 r_3}$, our approach communicates much less than the state-of-the-art approach (TuckerMPI)

Project: Scalable Tensor Algorithms for Modern Computing Systems

- 1 Design of Scalable Communication Optimal Algorithms for Tensors (Main Focus)
- 2 Extension of Existing Approaches/Algorithms (Short/Mid Term Research Plans)
- 3 Exploratory Topics (Mid/Long Term Research Plans)

Scalable communication optimal algorithms for tensors

- Analyze existing algorithms
- Determine communication lower bounds
- Propose communication optimal algorithms
- Implement the proposed algorithms

Main focus

Extension of existing approaches

- Strassen's concepts to tensors
- Concepts of hierarchical matrices to tensors
- Separation order of dimensions in tensor train

Short/Mid term plans

Exploratory topics

- New tensor representations
- Architecture aware algorithms
- Randomization in tensors
- Factorizations of tensors

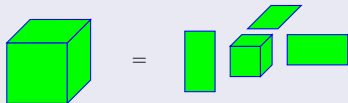
Mid/Long term plans

Scalable algorithms for popular tensor operations

- Determine the communication lower bounds for tensor decompositions
- Analyse the popular decomposition algorithms and communications performed by them
- Propose new scalable communication optimal algorithms
 - If possible design tiles/tasks based algorithms
- Implement the proposed algorithms
 - Handle performance issues for homogeneous systems
 - Load balancing
 - Memory aware approaches
 - scheduling strategies
- Same for manipulation operations of popular tensor representations
- Extend implementation for heterogeneous systems (start with Nvidia GPUs based heterogeneous systems)
- Create a tensor library

Popular tensor decompositions

Tucker decomposition



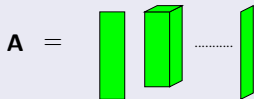
- Determine communication lower bounds for this operation
- Analyse communications performed by state of the art algorithms
- Propose and implement new scalable communication algorithms

Canonical decomposition



- No deterministic algorithm to find the decomposition
- Analyse one iteration of the popular existing algorithms
- Propose and implement scalable algorithms for one iteration

Tensor Train decomposition



- Determine communication lower bounds for this operation
- Analyse communication performed by popular algorithm
- Propose and implement new scalable communication algorithms

Proving communication lower bounds for parallel computations

How people did it for linear algebra operations?

- People obtain results for matrix multiplication operations
- Same lower bounds apply to almost all direct linear algebra operations using reduction [Ballard et. al., 09] , for instance, bound for LU factorization

$$\begin{pmatrix} I & & -B \\ A & I & \\ & & I \end{pmatrix} = \begin{pmatrix} I & & \\ A & I & \\ & & I \end{pmatrix} \begin{pmatrix} I & -B \\ & I & AB \\ & & I \end{pmatrix}$$

Approach to compute lower bounds for tensor computations

Notation: Tensors are denoted by solid shapes and number of lines denote the dimensions of the tensors. Connecting two lines implies summation (or contraction) over the connected dimensions.

- Obtain bounds for basic tensor operations: Tensor times matrix (TTM), Multiple tensor times matrix (Multi-TTM), Tensor contraction



- Express decompositions and manipulations in terms of these basis operations

Strassen's concepts to tensors

Matrix multiplication of $n \times n$ square matrices

- Complexity of traditional matrix multiplication is $\mathcal{O}(n^3)$
- Strassen's matrix multiplication
 - Expressed matrix multiplication as a tensor computation
 - Canonical rank of the tensor determines the complexity of the computation
 - Complexity is $\mathcal{O}(n^{2.81})$
- Plan to extend Strassen's concepts to tensor contractions

Contraction of a 3-dimensional tensor with a matrix

```
for  $i_1 = 1 : n$  do
  for  $i_2 = 1 : n$  do
    for  $i_3 = 1 : n$  do
      for  $j_2 = 1 : n$  do
         $\mathcal{G}(i_1, i_2, j_2) = \mathcal{G}(i_1, i_2, j_2) + \mathcal{A}(i_1, i_2, i_3) * B(i_3, j_2)$ 
      end for
    end for
  end for
end for
```

- Total $\mathcal{O}(n^4)$ operations
- Apply Strassen's algorithm for each i_1 , total $\mathcal{O}(n^{3.81})$ operations
- Expressing as a canonical decomposition of $8 \times 8 \times 4$ tensor can further reduce the number of operations

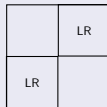
Hierarchical Matrix concepts to Tensors

Hierarchical Matrices

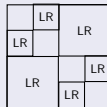
- Data sparse approximation of non-sparse matrices



Original Matrix



Step 1

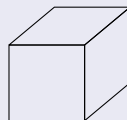


Step 2

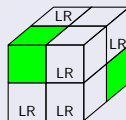
LR: low rank block

Tensors

- $f(i, j, k) = \frac{1}{|i-j|+|j-k|+|k-i|}$
- Value is small if difference of any pair is large
- Formalize and evaluate this approach for tensors



Original Tensor



Step 1

High Dimensional Tensor Representations and Randomization

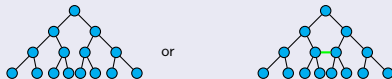
- Tensor Train is a popular representation to work with high dimensional tensors
- Adding tensors and applying an operator in this representation



- Requires a truncation process which iterates over cores one by one
- This representation is not much suited to work in parallel

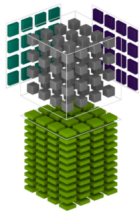
New tensor representations

- Look at new representations in tree format – suitable for parallelization



- Data will be stored at the leaf nodes
- Apply randomization to tensor operations

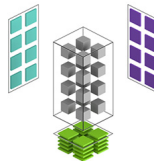
Architecture Aware Algorithms



$$D = \begin{bmatrix} A_{00} & A_{01} & A_{02} & A_{03} \\ A_{10} & A_{11} & A_{12} & A_{13} \\ A_{20} & A_{21} & A_{22} & A_{23} \\ A_{30} & A_{31} & A_{32} & A_{33} \end{bmatrix} + \begin{bmatrix} B_{00} & B_{01} & B_{02} & B_{03} \\ B_{10} & B_{11} & B_{12} & B_{13} \\ B_{20} & B_{21} & B_{22} & B_{23} \\ B_{30} & B_{31} & B_{32} & B_{33} \end{bmatrix} + \begin{bmatrix} C_{00} & C_{01} & C_{02} & C_{03} \\ C_{10} & C_{11} & C_{12} & C_{13} \\ C_{20} & C_{21} & C_{22} & C_{23} \\ C_{30} & C_{31} & C_{32} & C_{33} \end{bmatrix}$$

FP16 or FP32 FP16 FP16 FP16 or FP32

NVIDIA A100 Tensor Core FP64



- Recent Nvidia GPUs have tensor cores to accelerate AI computations
- Most linear algebra computations do not take advantages of these units
- Design algorithms which take architecture details into account

Fig source: www.nvidia.com

Integration in the LIP/LaBRI laboratory

Communication optimal tensor algorithms

- Analyze existing algorithms
- Determine communication lower bounds
- Propose communication optimal algorithms
- Implement the proposed algorithms

Main focus

Extension of existing approaches

- Strassen's concepts to tensors
- Concepts of hierarchical matrices to tensors
- Separation order of dimensions in tensor train

Short/Mid term plans

Exploratory topics

- New tensor representations
- Architecture aware algorithms
- Randomization in tensors
- Factorizations of tensors

Mid/Long term plans

ROMA team (LIP laboratory)

- *Bora Ucar*: design of tensor compression and manipulation algorithms
- *Gregoire Pichon*: low-rank based algorithms
- *Anne Benoit, Loris Marchal, Yves Robert and Frederic Vivien*: scalability and scheduling aspects in the long term

SATANAS team (LaBRI laboratory)

- *Olivier Beaumont and Lionel Eyraud-Dubois*: tensor train based neural network models
- *Mathieu Faverge*: low-rank based methods
- *Abdou Guermouche, Samuel Thibault*: exploitation of maximum potential of HPC systems in the long term

Bringing additional skills in the team

- High dimensional dense tensor computations, use of tensors in molecular simulations
- Communication lower bounds for linear algebra computations
- Scalable approaches for large HPC systems