Towards Exascale FFT Computation A survey, benchmark, and MPI Challenges

<u>A. Ayala¹</u> <u>S. Tomov¹</u> M. Stoyanov² J. Dongarra^{1,2}

¹Innovative Computing Laboratory University of Tennessee at Knoxville

²Oak Ridge National Laboratory

The 10th Annual MVAPICH User Group Meeting







Section summary

Background **Discrete and Fast Fourier Transform** FFT for Exascale Libraries for single-device systems Distributed and Multi-GPU FFT Libraries for parallel systems heFFTe Acceleration Batched FFT on GPUs Distributed FFT Convolution Mixed Precision FFT **3** Experiments & Profiling Scalability Communication Bottleneck FFT (Auto) Tuning Profiling Multi-GPU FFT

- The FFT is an algorithm developed by Cooley-Tukey in 1965
- Considered one of the top 10 algorithms of the 20th century.

Discrete Fourier Transform (DFT)

Let x be an m-dimensional array of size $N := N_1 \times N_2 \times \cdots \times N_m$. Its DFT is defined by y := DFT(x), obtained as:

$$y(k_1, k_2, \dots, k_m) := \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \cdots \sum_{n_m=0}^{N_m-1} \tilde{x} \cdot e^{-2\pi i \left(\frac{k_1 n_1}{N_1} + \frac{k_2 n_2}{N_2} \cdots + \frac{k_m n_m}{N_m}\right)},$$

where $\tilde{x} := x(n_1, n_2, ..., n_m)$.

- A naive DFT costs $\mathcal{O}(N^2)$
- Using the FFT, the cost can be reduced to $\mathcal{O}(N \log_2 N)$.

- The FFT is an algorithm developed by Cooley-Tukey in 1965
- Considered one of the top 10 algorithms of the 20th century.

Discrete Fourier Transform (DFT)

Let x be an m-dimensional array of size $N := N_1 \times N_2 \times \cdots \times N_m$. Its DFT is defined by y := DFT(x), obtained as:

$$y(k_1, k_2, \dots, k_m) := \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \cdots \sum_{n_m=0}^{N_m-1} \tilde{x} \cdot e^{-2\pi i \left(\frac{k_1 n_1}{N_1} + \frac{k_2 n_2}{N_2} \cdots + \frac{k_m n_m}{N_m}\right)},$$

where $\tilde{x} := x(n_1, n_2, ..., n_m)$.

- A naive DFT costs $\mathcal{O}(N^2)$
- Using the FFT, the cost can be reduced to $\mathcal{O}(N \log_2 N)$.

- The FFT is an algorithm developed by Cooley-Tukey in 1965
- Considered one of the top 10 algorithms of the 20th century.

Discrete Fourier Transform (DFT)

Let x be an m-dimensional array of size $N := N_1 \times N_2 \times \cdots \times N_m$. Its DFT is defined by y := DFT(x), obtained as:

$$y(k_1, k_2, \dots, k_m) := \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \cdots \sum_{n_m=0}^{N_m-1} \tilde{x} \cdot e^{-2\pi i \left(\frac{k_1 n_1}{N_1} + \frac{k_2 n_2}{N_2} \dots + \frac{k_m n_m}{N_m}\right)},$$

where $\tilde{x} := x(n_1, n_2, ..., n_m)$.

- A naive DFT costs $\mathcal{O}(N^2)$
- Using the FFT, the cost can be reduced to $\mathcal{O}(N \log_2 N)$.

- The FFT is an algorithm developed by Cooley-Tukey in 1965
- Considered one of the top 10 algorithms of the 20th century.

Discrete Fourier Transform (DFT)

Let x be an m-dimensional array of size $N := N_1 \times N_2 \times \cdots \times N_m$. Its DFT is defined by y := DFT(x), obtained as:

$$y(k_1, k_2, \dots, k_m) := \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \cdots \sum_{n_m=0}^{N_m-1} \tilde{x} \cdot e^{-2\pi i \left(\frac{k_1 n_1}{N_1} + \frac{k_2 n_2}{N_2} \cdots + \frac{k_m n_m}{N_m}\right)},$$

where $\tilde{x} := x(n_1, n_2, ..., n_m)$.

• A naive DFT costs $\mathcal{O}(N^2)$

• Using the FFT, the cost can be reduced to $\mathcal{O}(N \log_2 N)$.

- The FFT is an algorithm developed by Cooley-Tukey in 1965
- Considered one of the top 10 algorithms of the 20th century.

Discrete Fourier Transform (DFT)

Let x be an m-dimensional array of size $N := N_1 \times N_2 \times \cdots \times N_m$. Its DFT is defined by y := DFT(x), obtained as:

$$y(k_1, k_2, \dots, k_m) := \sum_{n_1=0}^{N_1-1} \sum_{n_2=0}^{N_2-1} \cdots \sum_{n_m=0}^{N_m-1} \tilde{x} \cdot e^{-2\pi i \left(\frac{k_1 n_1}{N_1} + \frac{k_2 n_2}{N_2} \dots + \frac{k_m n_m}{N_m}\right)},$$

where $\tilde{x} := x(n_1, n_2, ..., n_m).$

• A naive DFT costs $\mathcal{O}(N^2)$

• Using the FFT, the cost can be reduced to $\mathcal{O}(N \log_2 N)$.

und FFT for Exascale

Applications Relying on Parallel FFTs



Figure: Several applications from the U.S. ECP project heavily rely on FFTs.

MUG 2022

Single-Device FFT Libraries

Library	Language	Developer	GPU support	Open Source	2D & 3D support	Stride data support
CUFFT	С	NVIDIA	\checkmark		\checkmark	~
ESSL	C++	IBM			\checkmark	√
FFTE	Fortran	Riken		\checkmark	\checkmark	√
FFTPACK	Fortran	NCAR		\checkmark		
FFTS	С	U. Waikato		\checkmark		
FFTW3	С	MIT		\checkmark	\checkmark	√
FFTX	С	LBNL	\checkmark	\checkmark	\checkmark	√
KFR	C++	KFR		\checkmark		√
KISS	C++	Sandia		\checkmark	\checkmark	√
OneMKL	С	Intel	\checkmark		\checkmark	√
ROCM	C++	AMD	\checkmark	\checkmark	\checkmark	√
VkFFT	C++	D. Tolmachev	\checkmark	\checkmark	\checkmark	√

Figure: State-of-the-art of FFT libraries targeting a single-device unit.

Ref.: Interim Report on Benchmarking FFT Libraries on High Performance Systems Ayala et al., ICL Tech Report 2021.

Single-Device FFT Comparison

- Useful when input data is small or can be batched.
- heFFTe provides portability to run FFT experiment on different devices.



Figure: Comparison of single-device performance for a 512^3 FFT.

- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: Transfer computed data to neighbor processes
- 9: end for
- 10: Unpack data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}

- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid \cdot
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: Transfer computed data to neighbor processes
- 9: end for
- 10: Unpack data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}



- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: Transfer computed data to neighbor processes
- 9: end for
- 10: Unpack data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}



- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: Transfer computed data to neighbor processes
- 9: end for
- 10: Unpack data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}





- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: *Transfer* computed data to neighbor processes
- 9: end for
- 10: Unpack data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}





- 1: Input: 3-D array, processor grids: Pin, Pout
- 2: Transfer data from P_{in} to a pencil or slab grid
- 3: Define processor grids (MPI groups) for each direction
- 4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do
- 5: Compute local 1-D or 2-D FFTs on the GPUs
- 6: *Pack* data in contiguous memory
- 7: for P on my MPI group do
- 8: Transfer computed data to neighbor processes
- 9: end for
- 10: *Unpack* data in contiguous memory
- 11: end for
- 12: Transfer data from the pencil or slab grid to P_{out}





Algorithm 1 Parallel 3-D FFT computation on GPUs			
1: Input: 3-D array, processor grids: P_{in} , P_{out}			
2: Transfer data from P_{in} to a pencil or slab grid			
: Define processor grids (MPI groups) for each direction			
4: for $r \leftarrow 1, \cdots, n_{\text{exchanges}}$ do			
5: Compute local 1-D or 2-D FFTs on the GPUs			
6: <i>Pack</i> data in contiguous memory			
7: for P on my MPI group do			
8: <i>Transfer</i> computed data to neighbor processes			
9: end for			
0: Unpack data in contiguous memory			
11: end for			
Transfer data from the pencil or slab grid to P_{out}			





These 3 tasks can be replaced by 1 via MPI Alltoallw

Communication can be accelerated by enabling Mixed-Precision, c.f., Advances in Mixed Precision Algorithms: 2021 Edition. *Abdelfattah et al.*, *LLNL-TR-825909*

Parallel FFT Libraries

Library	Developer	Language	CPU Backend	GPU Backend	Real-to- Complex	Slab Decomp.	Brick Decomp.
2DECOMP&FFT	NAG	Fortran	FFTW3, ESSL	-	\checkmark	\checkmark	
AccFFT	Georgia Tech	C++	FFTW3	CUFFT	\checkmark		
Cluster FFT	Intel	Fortran	MKL	-			
CRAFFT	Cray	Fortran	FFTW3	-	\checkmark		
cuFFTMp	NVIDIA	С	-	CUFFT	\checkmark		
FFTE	U. Tsukuba / Riken	Fortran	FFTE	CUFFT	\checkmark	\checkmark	
fftMPI	Sandia	C++	FFTW3, MKL, KISS	-			\checkmark
FFTW3	MIT	С	FFTW3	-	\checkmark	\checkmark	
heFFTe	ICL - UTK	C++	FFTW3, MKL, Stock	CUFFT, ROCM, OneMKL	\checkmark	\checkmark	\checkmark
nb3DFFT	RWTH Aachen	Fortran	ESSL	-	\checkmark		
P3DFFT	UC San Diego	C++	FFTW3	-	\checkmark	\checkmark	
spFFT	ETH	C++	FFTW3	CUFFT, ROCM	\checkmark	\checkmark	
SWFFT	Argonne	C++	FFTW3	-			\checkmark

Figure: State-of-the-art of FFT libraries targeting parallel systems.

Ref.: Interim Report on Benchmarking FFT Libraries on High Performance Systems Ayala et al., ICL Tech Report 2021.

The Highly Efficient FFT for Exascale (heFFTe)

- heFFTe is part of the US Exascale Project (ECP).
- Funded by DoE, it aims to provide reliable FFT computation on exascale systems.
- Integrated to ECP projects such as Copa-Cabana.
- Unique features: Batched FFT, Convolution, Sine/Cosine transform, Mixed Precision FFT.



Figure: heFFTe interfaces and architecture support.

Moving FFT Kernels to GPUs

- Moving local kernels to GPUs yields high speedups.
- In the following example, we show this effect on Summit.



heFFTe Acceleration

Figure: Performance comparison between fftMPI, using 1280 IBM Power9 Cores, 40 cores per node (left) and heFFte, using 192 NVIDIA V100-GPUs, 6 per node (right).

Ref.: heFFTe: Highly Efficient FFT for Exascale. Ayala_et al., ICCS 2020.

LAMMPS Rhodopsin Benchmark using heFFTe

• Molecular dynamics apps heavily rely on FFTs, and often have their own parallel FFT implementation (e.g., fftMPI, SWFFT).

heFFTe Acceleration

LAMMPS Rhodopsin Benchmark using heFFTe

- Molecular dynamics apps heavily rely on FFTs, and often have their own parallel FFT implementation (e.g., fftMPI, SWFFT).
- Using heFFTe real-to-complex accelerates LAMMPS K space kernel around $1.76 \times$.



Figure: Breakdown for the LAMMPS Rhodopsin experiment. Using 32 Summit nodes, 6 V-100 GPUs per node, and 1 MPI per GPU.

Ref.: Performance Analysis of Parallel FFT on Large Multi-GPU Systems. Ayala et al., IEEE IPDPS 2022.

3-D Batched FFT experiments

• Batched FFTs are needed in applications such image processing, filtering, particle energy computations, among others.



Figure: Batch of 3-D FFT of size 64^3 on NVIDIA and AMD GPUs, 1 MPI per GPU. Speedups of over $2\times$ with respect to the not batched version.

3-D Batched FFT experiments

• Batched FFTs are needed in applications such image processing, filtering, particle energy computations, among others.



Figure: Batch of 3-D FFT of size 64^3 on NVIDIA and AMD GPUs, 1 MPI per GPU. Speedups of over $2\times$ with respect to the not batched version.

3-D Convolution

- Convolutions are widely used for solving PDEs, e.g., in Earthquake simulation.
- They are also useful for convolutional networks.



Figure: Convolution of a 3-D FFT of size 64^3 , using 40 Power9 cores and 6 V-100 GPUs per node.

Enabling Mixed Precision in heFFTe

- Can be achieved using third party lossless and lossy compression libraries (NVOMP, ZFP).
- Currently, heFFTe uses casting compression + our own implementation of MPLOSC_Alltoall



Figure: Speedup obtained reducing Floating Point (FP).

Ref.: Mixed precision and approximate 3D FFTs. Cayrols et al., IEEE Cluster 2022.

Enabling Mixed Precision in heFFTe

- Can be achieved using third party lossless and lossy compression libraries (NVOMP, ZFP).
- Currently, heFFTe uses casting compression + our own implementation of MPLOSC_Alltoall



Figure: Accuracy obtained reducing Floating Point (FP).

Ref.: Mixed precision and approximate 3D FFTs. Cayrols et al., IEEE Cluster=2022,

• Parallel FFT scales at the same rate of the underlying MPI framework, until a breakdown point.



Figure: Strong Scalability on up to 6144 V-100 GPUs on Summit (left), and Weak Scalability on up to ~ 1.2 M Fugaku core, 48 cores per node (right). Using heFFTe.

Ref.: Scalability Issues in FFT Computation. Ayala et al., ACM PACT 2021.

• Parallel FFT scales at the same rate of the underlying MPI framework, until a breakdown point.



Figure: Strong Scalability on up to 6144 V-100 GPUs on Summit (left), and Weak Scalability on up to $\sim 1.2M$ Fugaku core, 48 cores per node (right). Using heFFTe.

Ref.: Scalability Issues in FFT Computation. Ayala et al., ACM PACT 2021.

• Similar behavior is observed for state-of-the-art FFT libraries.



Figure: Strong Scalability on 32K Power9 cores for CPU-based libraries (left), and 4096 V-100 for GPU-based libraries (right).

Ref.: FFT Benchmark Performance Experiments on Systems Targeting Exascale. Ayala et al., ICL Tech Report 2022.

• Similar behavior is observed for state-of-the-art FFT libraries.



Figure: Strong Scalability on 32K Power9 cores for CPU-based libraries (left), and 4096 V-100 for GPU-based libraries (right).

Ref.: FFT Benchmark Performance Experiments on Systems Targeting Exascale. Ayala et al., ICL Tech Report 2022.

3-D FFT with All-to-All Communication



Figure: Vampir trace of back-to-back 3-D FFTs of size 1024^3 (5 forward + 5 backward), using 4 Summit nodes with 16 NVIDIA GPUs, 4 MPIs per node.

3-D FFT with Point-to-Point Communication



Figure: Vampir trace of back-to-back 3-D FFTs of size 1024^3 (5 forward + 5 backward), using 4 Summit nodes with 16 NVIDIA GPUs, 4 MPIs per node.

Scaling Communication

- While messages become smaller, latency effects become significant.
- Next figure shows variability of time spent on MPI_Alltoall.



Figure: Strong Scalability on up to 6144 V-100 GPUs on Summit (left), and Weak Scalability on up to $\sim 1.2M$ Fugaku core, 48 cores per node (right). Using heFFTe.

 Ref.: Accelerating MultProcess Communication for Parallel 3-D FFT. Ayala et al.,

 ExaMPI SC 2021.

Tuning Processor grid

- At every transposition step a processor grid $P \times Q$ is defined
- The selection of this grid highly impacts scalability, in the next figure we set P = 5 and let Q vary.



Figure: Strong Scalability on 40960 Power9 cores setting different grids $5 \times Q$. Using P3DFFT. Auto-tuning also available with XTune.

Ref.: Interim Report on Benchmarking FFT Libraries on High Performance Systems Ayala et al., ICL Tech Report 2021.

Tuning Processor grid

- At every transposition step a processor grid $P \times Q$ is defined
- This grid can be found using autotuning, e.g., using GPU-Tune

# Task parameters				
dimx	= Integer	<pre>(64, 128, transform="normalize", name="dimx")</pre>		
dimy	= Integer	<pre>(64, 128, transform="normalize", name="dimy")</pre>		
dimz	= Integer	<pre>(64, 128, transform="normalize", name="dimz")</pre>		
# Inpu	t/ <mark>tuning</mark> paramet	ters		
#### note that the px_i , py_i , px_o , py_o are in log scale				
px_i	= Integer	<pre>(0, np.log2(nodes*npernode), transform="normalize", name="px_i")</pre>		
py_i	= Integer	<pre>(0, np.log2(nodes*npernode), transform="normalize", name="py_i")</pre>		
px_o	= Integer	<pre>(0, np.log2(nodes*npernode), transform="normalize", name="px_o")</pre>		
py_o	= Integer	<pre>(0, np.log2(nodes*npernode), transform="normalize", name="py_o")</pre>		
<pre>comm_type = Categoricalnorm (['a2a', 'p2p'], transform="onehot", name="comm_type")</pre>				
# Tuning Objective				
time	= Real	<pre>float("-Inf") , float("Inf"), name="time")</pre>		

Figure: GPU-Tune for autotuning heFFTe processor grids.

Ref.: Autotuning heFFTe with GPU-Tune. Sherry Li's Group at LBNL.

MUG 2022

Tuning Algorithm

- Choosing between Pencil or Slab can lead to considerable speedups, $\sim 30\%$.
- The choice of Binary or Collective MPI is important at large-scale.



Figure: Phase Diagram (left) and best settings for a 512^3 FFT with 6 GPUs per node and manual tuning (right). Using heFFTe on Summit.

Ref.: Impacts of Multi-GPU MPI collective communications on large FFT computation Ayala et al., ExaMPI SC 2019.

Profiling Communication of Distributed FFTs

- We have used the following tools:
 - NVVP, rocProf
 - Score-P, Vampir, Cube
- Available ongoing developments include: OSU-INAM
- Challenges of the state-of-the-art:
 - Limited support for modern C++
 - Lack of support for all MPI distributions
 - Some do not support GPUs yet
 - If any, support for GPUs is vendor-specific
 - Hard to perform user-specific requests
 - In general, output does not reflect FFT workloads
 - Not mapping to the architecture and network
 - Developing a profiler tool requires lots of effort

• heFFTe tracing provides a detailed timing for tasks

- heFFTe tracing provides a detailed timing for tasks
- We need an adaptable visualization tool

- heFFTe tracing provides a detailed timing for tasks
- We need an adaptable visualization tool
- We can take advantage of the following analogies:
 - Network Traffic \longleftrightarrow Vehicular Traffic
 - Memory unit \longleftrightarrow Vehicle
 - Practical Bandwidth \leftrightarrow Average Velocity
 - Peak Bandwidth \longleftrightarrow Road Capacity

- heFFTe tracing provides a detailed timing for tasks
- We need an adaptable visualization tool
- We can take advantage of the following analogies:
 - Network Traffic \longleftrightarrow Vehicular Traffic
 - Memory unit \longleftrightarrow Vehicle
 - Practical Bandwidth \leftrightarrow Average Velocity
 - Peak Bandwidth \longleftrightarrow Road Capacity
- Using the power of heFFTe tracing + Python + TransCAD, we can build an architecture-aware visualization of network congestion.

- heFFTe tracing provides a detailed timing for tasks
- We need an adaptable visualization tool
- We can take advantage of the following analogies:
 - Network Traffic \longleftrightarrow Vehicular Traffic
 - Memory unit \longleftrightarrow Vehicle
 - Practical Bandwidth \leftrightarrow Average Velocity
 - Peak Bandwidth \longleftrightarrow Road Capacity
- Using the power of heFFTe tracing + Python + TransCAD, we can build an architecture-aware visualization of network congestion.
- Defining network peak values, TransCAD can also provide optimized paths

Ref.: heFFTe profiler Ayala, 2022.

Profiling Communication with heFFTe

• We can obtain a video of data-volume exchange over time.

Profiling with heFFTe

• We can obtain a video of network traffic over time.



Acknowledgments

- heFFTe is funded by the Department of Energy (DoE) Exascale Project WBS 2.3.3.13.
- Collaborators:
 - A. Haidar (NVIDIA)
 - ICL OpenMPI Team (UTK)
 - ICL FIBER Team (UTK)
 - Network-Based Computing Research (DK. Panda's group, OSU)
 - ECP X-Tune (Sherry Li's group, LBNL)
 - D. Takahashi (U. Tsukuba)
 - D. Pekurovsky (SDSC)