

An MPI-based 3D FFT Implementation on CPU+GPU Heterogeneous Clusters

Kun Li

Yan Li, Ting Cao, Haipeng Jia, Yunquan Zhang SKL of Computer Architecture, ICT, CAS



Basics

* FFT(Fast Fourier Transform)

A fast algorithm to compute DFT, widely applied

An important algorithm of data processing in SKA project

Bandwidth-intensive task

Considering the characteristic of matrix W_n caused by the periodicity of \mathcal{O}_N , the basic idea of FFT algorithm is to decompose matrix into many small matrix to calculate



3

Background

(a):
$$y_k = \sum_{j=0}^{N-1} \omega_n^{jk} \cdot x_j \ (0 \le k < n) \quad \omega_n = e^{-\frac{2\pi i}{n}} (i = \sqrt{-1})$$

$$DFT_n = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_n^1 & \omega_n^2 & \dots & \omega_n^{n-1} \\ 1 & \omega_n^2 & \omega_n^4 & & \omega_n^{2(n-1)} \\ \vdots & \ddots & \vdots \\ 1 & \omega_n^{n-1} & \omega_n^{2(n-1)} & \dots & \omega_n^{(n-1)(n-1)} \end{bmatrix}, y = DFT_n \cdot x$$
(b): $j = (j_1, j_2) = j_1 \cdot r + j_2 \quad (0 \le j_1 < q, \ 0 \le j_2 < r)$
 $k = (k_2, k_1) = k_1 + k_2 \cdot q \quad (0 \le k_1 < q, \ 0 \le k_2 < r)$

 $k = (k_2, k_1) = k_1 + k_2 \cdot q \ (0 \le k_1 < q, 0 \le k_2)$ $n = q \cdot r$

And Person .

(c):
$$y_k = \sum_{j_2=0}^{r-1} \left[\left(\sum_{j_1=0}^{q-1} \omega_q^{j_1k_1} \cdot x_{j_1 \cdot r + j_2} \right) \cdot \omega_n^{j_2k_1} \right] \cdot \omega_r^{j_2k_2}$$



(c):
$$y_k = \sum_{j_2=0}^{r-1} \left[\left(\sum_{j_1=0}^{q-1} \omega_q^{j_1k_1} \cdot x_{j_1 \cdot r + j_2} \right) \cdot \omega_n^{j_2k_1} \right] \cdot \omega_r^{j_2k_2}$$

$$A \otimes B = \begin{pmatrix} a_{0,0}B & \cdots & a_{0,m-1}B \\ \vdots & \ddots & \vdots \\ a_{m-1,0}B & \cdots & a_{m-1,m-1}B \end{pmatrix}$$

 $A \otimes B = [a_k, B]$, for $A = [a_k, P]$.

(d): $DFT_n = (DFT_r \otimes I_q)T_q^n (I_r \otimes DFT_q)L_r^n$



(a) Matrix factorization



- General factorizations
 - * FFT(Fast Fourier Transform)

Example: $N=N_1 \times N_2 = 4 \times 4$





Salar Chart



FFT on CPU+GPU Clusters

- CPU+GPU heterogeneous computing system has gradually become a hot research direction in the field of high performance computing.
 - Tianhe-IA
 - Dawning Nebulae
- Possible platform candidate for SDP
- Case study: 3D FFT on heterogeneous clusters
 - Application scenarios: astrophysical N-body simulations, blood flow simulations, etc.
 - Parallel algorithm: based on 1D and 2D decomposition



Algorithm Strategy

- The research of 3D FFT on clusters is mainly based on 1D and 2D decomposition
 - * $n_0 \times n_1 \times n_2$ ($n_0 \ge n_1 \ge n_2$), Computing sequence: n_2 , n_1 , n_0
 - 1D decomposition
 - Strategy
 - 1. Along n_0 dimension, divide data into $P(P \le n)$ processes.
 - 2. Secondly, *P* processes computes 2D FFT(size= $n_1 \times n_2$) of n_0/p batches parallelly.

3. Because the data needed for the n_0 dimension are distributed in all processes, We need to do a full exchange of data on all processes to complete the computation of n_0 dimension.

bottleneck

 $_{n_0}$ The scale of parallel computing is limited to n_1 or n_2





Algorithm Strategy

- The research of 3D FFT on clusters is mainly based on 1D and 2D decomposition
 - $\bullet \quad n_0 \times n_1 \times n_2 \ (n_0 \ge n_1 \ge n_2)$
 - 2D decomposition
 - The number of processes allowed: $n_0 \times n_1$





Parallel Algorithm

According to the ratio of communication, algorithm is divided into 3 stages :

a: Initialize data distribution, and compute local multi-dimension FFT

for $t \leftarrow 0, ..., d-r-2$ do $h_0 \leftarrow X_{s=0}^{r-1} \frac{N_s}{P_s} \times X_{s=r}^{d-2-t} N_s$ $N \leftarrow N_{d-1-t}$ $h_1 \leftarrow X_{s=d-t}^{d-1} \stackrel{\circ}{N_s}$ $h_0 \times N \times h_1 \stackrel{\text{FFT}}{\rightarrow} h_0 \times \stackrel{\circ}{N} \times h_1$ end for (a) Scale: $N = N_0 \times ... \times N_{d-1}$ Processes: $P = P_0 \times ... \times P_{r-1}$ r < d $\times_{s=l}^{u} N_s = N_l \times ... \times N_u$

local FFT computation

• Call existing 1D FFT library

Example(using 2D process grid to compute 3D FFT)

a.
$$\frac{n_0}{p_0} \times \frac{n_1}{p_1} \times n_2 \xrightarrow{FFT} \frac{n_0}{p_0} \times \frac{n_1}{p_1} \times \hat{n}_2$$



Parallel Algorithm

According to the ratio of communication, algorithm is divided into 3 stages :

b: 1) local transposition2) global transposition

3)local computation



 T_0 :local transposition T_1 :global transposition MPI_AlltoAll MPI_Alltoallv

Example(using 2D process grid to compute 3D FFT)



Parallel Algorithm

According to the ratio of communication, algorithm is divided into 3 stages :

c: Do *r* local transposition and global transposition on the outputs of B to get the correct outputs.

$$\begin{aligned} h_{1} \leftarrow \times_{s=r+1}^{d-1} \stackrel{\wedge}{N} \\ h_{0} \leftarrow \times_{s=0}^{r-1} \frac{\stackrel{\wedge}{N}_{s+1}}{P_{s}} \\ N \leftarrow N_{0} \\ h_{0} \times N \times h_{1} \stackrel{\text{FFT}}{\rightarrow} h_{0} \times \stackrel{\wedge}{N} \times h_{1} \end{aligned}$$
for $t \leftarrow 0, \dots, r-1$ do
$$L_{1} = N_{r-t} \\ h_{1} \leftarrow \times_{s=r-t}^{r-1} \frac{N_{s+1}}{P_{s}} \times \times_{s=0}^{r-t-2} \frac{N_{s}}{P_{s}} \\ L_{0} \leftarrow N_{r-t-1} \\ h_{0} \leftarrow \times_{s=r+1}^{d-1} N_{s} \\ P \leftarrow P_{r-t-1} \\ \frac{L_{1}}{P} \times h_{1} \times L_{0} \times h_{0} \stackrel{\text{T}}{\rightarrow} L_{1} \times h_{1} \times \frac{L_{0}}{P} \times h_{0} \\ N \leftarrow N_{r-t} \\ N \times h_{0} \times h_{1} \stackrel{\text{T}}{\rightarrow} h_{0} \times N \times h_{1} \end{aligned}$$
end for
(c)

$$\frac{n_1}{p_0} \times \frac{n_2}{p_1} \times n_0 \xrightarrow{T_1} n_1 \times \frac{n_2}{p_1} \times \frac{n_0}{p_1} \xrightarrow{T_0} \frac{n_2}{p_1} \times \frac{n_0}{p_1} \times n_1 \xrightarrow{T_1} n_2 \times \frac{n_0}{p_0} \times \frac{n_1}{p_1} \xrightarrow{T_0} \frac{n_0}{p_0} \times \frac{n_1}{p_1} \times n_2$$



Experimental Result

Implementation Strategy

- Local computation calls NVIDIA CUFFT.
- Local transposition computes on GPU.

When GPU computation completed, do transposition immediately to make use of high bandwidth of GPU effectively.

Global transposition uses encapsulated FFTW transposition functions.



Experimental Result

Experimental Platform

	Platform	nodes	CPU	RAM	GCC	GPU	GPU SDK
	А	2	Dual Intel E5430 2.66GHz	24G	4.4.3	Tesla C1060	CUDA 4.2
	В	16	Dual Intel E5430 2.66GHz	24G	4.4.3	Tesla C1060	CUDA 4.2
	network	A: QDR Infiniband			B: Intel 82573 Gigabit Ethernet		

- san chart

Experimental Result



A: QDR Infiniband



B: Intel 82573 Gigabit Ethernet



15

Summary

✤ The basics of FFT

- FFT algorithm strategy: 1D decomposition and 2D decomposition, and the analysis of scalability
- Illuminate the MPI parallel FFT algorithm
- * FFT algorithm performance evaluation



16

Summary

