HOKUSAI BigWaterfall2 利用手引書

Version 1.12 2025 年 5 月 9 日 理化学研究所 情報セキュリティ・システム部

改版履歴

版	発行年月日	章	改版理由
1.0	2023.12.4	-	新規作成
1.1	2023.12.7	1.2	超並列演算システム、大容量メモリ演算サーバの演算性能の
			計算式を修正
1.2	2023.12.8	5.7.5	説明文とコマンド実行例を修正
1.3	2023.12.11	5.6	ジョブスクリプト例を修正
1.4	2023.12.12	2.2	接続先アドレスを追加
		5.6	ジョブスクリプト例を修正
1.5	2023.12.19	3.2	Lfs quota コマンドを修正
1.6	2023.12.21	5.6	ジョブスクリプト例を修正
1.7	2023.12.22	5.2	ジョブ実行制限を修正
1.8	2023.12.27	5.2	メモリの単位を GiB 表記に修正
		5.3	表 5-10 を修正
1.9	2024.4.2	3.2	HPCI 課題の説明を追加
		4.2	GPU 向けのコンパイル/リンク例を追記
		5.2	mpc_ht の記述を削除
			GPU サーバのジョブ実行リソースを追記
		5.4	GPU でインタラクティブジョブを実行する際の注釈を追記
		5.3.2	ジョブ投入オプションを追記
		5.6.1	GPU サーバ(GPU)向けジョブスクリプトの例の記載を追記
1.10	2024.7.1	全体	用語統一
		1.1	システム構成図の修正
		1.2	メモリの単位変更
		1.6	パーティションの優先度の説明を追記
		3.1	
		3.2	GID の確認方法、lfs コマンドの-h オブションのコメント追記
		4.1	(default)表記の説明、intel ロード時の補足を追記
		4.5	開発ツールの記載を6章から4.5に移動
		5.2	ハーティションにおける課題ことの設定についての記載内容を
			ハーテインヨンについての記載から「mpc_a」と「Imc_a」の記載
			そりは
		5.2	
		0.0 6 空	ypus オフションを迫記 Singularity コンニナ 泊記
1 11	2024 40 25	0 早 皮染	Singulanty コノノノ 迫記 問い合わせてはの記載を修正
1.11	2024.10.20	/ナ 証冊	同じロクビクズの記載を修正
1 12	2025 5 9	16	ジョブ実行順序の設定変更に伴う記載内容の修正

はじめに	1
1. HOKUSAI BigWaterfall2 システムの概要	2
1.1 システム構成	2
1.2 ハードウェア概要	4
1.3 ソフトウェア構成	5
1.4 サービス時間	5
1.5 利用規約について	6
1.6 ジョブの実行順序について	7
2. HOKUSAI BigWaterfall2 システムへのログイン	8
2.1 システムログインまでの流れ	8
2.2 接続先アドレス	8
3. ファイルシステム環境	9
3.1 ストレージ領域	9
3.2 ディスク使用状況と上限値の表示	11
3.3 一時領域	11
4. コンパイル/リンク	12
4.1 環境設定	12
4.2 コンパイラ	14
4.3 コンパイル/リンク概要	14
4.4 数値計算ライブラリ	21
4.5 開発ツール	22
5. バッチジョブ・インタラクティブジョブ	23
5.1 ジョブ管理システム概要	23
5.2 ジョブ実行リソース	24
5.3 ジョブ投入オプション	25
5.4 インタラクティブジョブ実行	28
5.5 バッチジョブ投入	29
5.6 バッチジョブスクリプト例	34
5.7 ジョフ状態表示	41
5.8 ジョフキャンセル	47
5.9	47
6. Singularity コンテナ	49
6.1 Singularity イメージの取得	49
6.2 Singularity の"Command Line Interface(CLI)"について	49
6.3 Ubuntu の sif イメージを利用したインタラクティブジョブ実行	50

Copyright (C) RIKEN, Japan. All rights reserved.

6.4 Ubuntu の sif イメージを利用したバッチジョブ実行	50
6.5 CPU 版(MPC) Singularity Ubuntu 環境構築	51
6.6 GPU ノード向け Singularity コンテナの取得	52
6.7 インタラクティブジョブ Singularity コンテナ内で Shell 起動	53
6.8 インタラクティブジョブ PyTorch コンテナで CuPy を使用した実行例	54
6.9 バッチジョブ PyTorch コンテナの GPU 実行例	55

Copyright (C) RIKEN, Japan. All rights reserved.

はじめに

本利用手引書は、理化学研究所が導入したスーパーコンピュータ・システム(HOKUSAI BigWaterfall2 システム)の利用方法について説明した資料です。システムを利用する方は必ずお読みください。

本利用手引書の内容は不定期に更新いたします。最新の利用手引書は、以下のHBW2 ポータルより入手することが可能です。

https://hokusai.riken.jp/

HBW2 ポータルでは、利用可能なアプリケーションのバージョンや実行方法がまとめて あります。また、ssh 公開鍵の登録もここで行います。詳細につきましては「HOKUSAI BigWaterfall2(HBW2) Web ポータル 利用手引書」をご参照ください。

システムの運用状況は HBW2 ポータルや利用者向けメーリングリストでアナウンスしま す。また、HOKUSAI BigWaterfall2システムについてのあらゆるご質問は、以下の Webペ ージをご参照頂き、お問い合わせください。

https://i.riken.jp/supercom/contact/

本書の一部、または全部を無断で複製、転載、再配布することを禁じます。

Copyright (C) RIKEN, Japan. All rights reserved.

1. HOKUSAI BigWaterfall2 システムの概要

1.1 システム構成

HOKUSAI BigWaterfall2 システム(以下、HBW2)は、超並列演算システム(MPC)、大容 量メモリ演算サーバ(LMC)と、システムの利用入口となるフロントエンドサーバ、共有ストレ ージから構成されるシステムです。



図 1-1 システム構成図

HBW2 の超並列演算システム(MPC)は、PRIMERGY CX2550 M7 を 312 ノードで構成 します。1 ノードの理論演算性能は 6.8TFLOPS、主記憶容量は 119GiB です。各ノードは、 InfiniBand NDR(400GB/s)で接続され、高速なノード間通信とファイル共有を実現します。

大容量メモリ演算サーバ(LMC)は、PRIMERGY RX4770 M7 を 2 ノードで構成します。1 ノードの理論演算性能は 6.4TFLOPS、主記憶容量は 2.7TiB です。各ノードは、InfiniBand NDR(400.0GB/s)で接続され、高速なノード間通信とファイル共有を実現します。

Copyright (C) RIKEN, Japan. All rights reserved.

ストレージ環境は、共有ストレージ(SFS)で構成します。

共有ストレージ(SFS)は、各ユーザーのホームディレクトリや課題グループ用の共有ディ レクトリなど、広帯域でオンライン性のあるファイルシステムであり、超並列演算システム (MPC)、大容量メモリ演算サーバ(LMC)およびフロントエンドサーバから参照可能です。利 用可能容量は 9.8PB です。

HBW2 へのアクセスは、ssh/scp によるアクセスと HTTPS アクセス(HBW2 ポータル)が 可能です。ユーザーはフロントエンドサーバ上にて、プログラムの編集、コンパイル/リンク、 バッチジョブの操作、インタラクティブジョブの実行、チューニング、デバッグ等の作業を行う ことが可能です。

Copyright (C) RIKEN, Japan. All rights reserved.

1.2 ハードウェア概要

1.2.1 超並列演算システム(MPC)

● 演算性能

CPU: Intel Xeon Max 9480 (1.9GHz) 312 台(624CPU, 34,944 コア) 理論ピーク性能: 2.12 PFLOPS (1.9GHz × 56 コア × 32FMA × 2CPU × 312 ノー ド)

- 主記憶 メモリ容量: 36TiB(119GiB × 312 台) メモリバンド幅: 3.26TB/s
- 内蔵ディスク
 ディスク容量: 300TB (480GB × 2 × 312 台)
- インターコネクト InfiniBand NDR 通信性能:ノード間 400GB/s

1.2.2 大容量メモリ演算サーバ(LMC)

- 演算性能 CPU: Intel Xeon Gold 6418H (2.1GHz) 2 台(8CPU, 192 コア) 理論ピーク性能: 6.1TFLOPS (2.1GHz × 24 コア × 4CPU × 32FMA)
 メモリ
 - メモリ容量: 5.4TiB(2.7TiB × 2台) メモリバンド幅: 1126.4GB/s
- 内蔵ディスク
 ディスク容量: 2.4TB (600GB × 2 × 2 台)
- インターコネクト InfiniBand NDR 通信性能:ノード間 400GB/s

Copyright (C) RIKEN, Japan. All rights reserved.

1.3 ソフトウェア構成

HBW2 で利用可能なソフトウェアー覧を以下に示します。

	F •				
項目	超並列演算システム(MPC)	大容量メモリ演算サーバ	フロントエンドサーバ		
		(LMC)			
OS	Rocky Linux 8(x 312 台)	Red Hat Enterprise Linux 8	Red Hat Enterprise Linux 8		
	(Linux kernel version 4.18)	(Linux kernel version 4.18)	(Linux kernel version 4.18)		
コンパ	インテル oneAPI ベース & HPC ツー	ルキット マルチノード			
イラ	インテル C/C++および Fortran コンパ	パイラ			
	インテル TBB				
	インテル Distribution for Python				
ライブ	インテル oneAPI ベース & HPC ツ	ールキット マルチノード			
ラリ	インテル MKL				
	インテル MPI ライブラリ				
	インテル IPP				
	インテル DAL				
ツール	インテル oneAPI ベース & HPC ツールキット マルチノード				
	インテル VTune				
	インテル Advisor				
	インテル Inspector				
	インテル Trace Analyzer & Collecto	r			
アプリ	Gaussian, ADF, AMBER,	Gaussian, ADF, AMBER,	GaussView,		
ケーシ	GAMESS,GROMACS,	GAMESS, GROMACS,	VMD, ROOT		
ョン	NAMD, ROOT	NAMD, ROOT			

表 1-1 ソフトウェア一覧

HBW2 で利用可能なアプリケーションやライブラリなどの最新情報、利用方法は以下の HBW2 ポータルに掲載します。

https://hokusai.riken.jp/

1.4 サービス時間

HBW2 は原則 24 時間サービスを行いますが、システムの定期保守作業、緊急保守作 業、設備保守作業(空調機メンテナンス・電源設備点検)を行う必要がある場合など、サービ スを休止します。HBW2 の運用状況は HBW2 ポータルや利用者向けメーリングリストでア ナウンスします。

Copyright (C) RIKEN, Japan. All rights reserved.

1.5 利用規約について

詳細は「スーパーコンピュータ・システム利用規約」をご参照ください。「スーパーコンピュ ータ・システム利用規約」は以下の URL から参照できます。

https://i.riken.jp/supercom/

1.5.1 利用可能演算時間

課題によって、利用できる演算時間が異なります。課題番号(project ID)、割り当てられた演算時間、使用済み演算時間、割り当てられた演算時間の有効期限は、*listcpu* コマンドで確認できます。

使用済み演算時間が100%に達した場合、ジョブを新規に投入/実行できなくなります。

[username@h	[username@hokusai1~]\$ listcpu -p 課題番号(project ID)					
[RB999999]						
Node/User	Used(h)	Limit(h)	Used(%)	Expiry date		
mpc	0	5115801.6	0.0	2024-03-31		
lmc	0	28108.8	0.0	2024-03-31		
userA	0	-	-	2024-03-31		
userB	0	-	-	2024-03-31		
:						

表 1-2 listcpu コマンドの表示項目

項目	説明
Used (h)	使用済み演算時間 (単位: 時間)
Limit (h)	割り当て演算時間 (単位:時間)
Use (%)	使用済み演算時間 / 割り当て演算時間 (単位: %)
Expiry date	割り当て演算時間の有効期限

Copyright (C) RIKEN, Japan. All rights reserved.

1.6 ジョブの実行順序について

本システムでは、ジョブの優先度をリアルタイムに評価し実行順序を制御しています。優 先度の評価は以下の項目を基に行っています。

表 1-3 優先度制御における評価項目

評価順序	評価項目	概要
1	課題毎のフェアシェア値	利用状況によって算出される課題毎のフェアシ
		ェア値が高い順に実行

フェアシェアによる優先度制御は、課題毎に割り当てられたリソースに対して、公平にリ ソースを利用できるようにするための仕組みです。

課題毎にフェアシェア値を持ち、フェアシェア値が一番高い課題から優先的にジョブを実行します。直近の利用量が多いとフェアシェア値が下がり、時間が経つと回復します。その ため割り当てられたリソースを有効に活用するためには、一年を通して継続的に利用する ことが必要です。

2. HOKUSAI BigWaterfall2 システムへのログイン

2.1 システムログインまでの流れ

HBW2 にログインするためにはアカウント申請、公開鍵登録を実施する必要がありま す。アカウント申請や公開鍵登録は HBW2 ポータルから実施してください。詳細の手順に つきましては「HOKUSAI BigWaterfall2(HBW2) Web ポータル 利用手引書」をご参照くだ さい。

2.2 接続先アドレス

HBW2 に SSH/HTTPS アクセス可能なサーバはフロントエンドサーバです。接続先ア ドレスを表 2-1 に示します。

フロントエンドサーバは合計 4 台のサーバで構成されており、hokusai.riken.jp にアクセ スすると、DNS ラウンドロビンにより、いずれかのログインノードにアクセスすることができ ます。HBW2 ポータルは、HTTPS アクセスにより、SSH 公開鍵登録やマニュアル参照等 が可能です。

ホスト名	プロトコル	アクセス用途
hokusai.riken.jp	ssh	● 仮想端末
hokusai1.riken.jp		● ファイル転送
hokusai2.riken.jp		
hokusai3.riken.jp		
hokusai4.riken.jp		
hokusai.riken.jp	https	● SSH 公開鍵登録
		• マニュアル参照
		● 各種申請

表 2-1 接続先アドレス一覧

3. ファイルシステム環境

3.1 ストレージ領域

HBW2 では、以下のストレージ領域が利用可能です。

表	3-1	スト	レージ	·領域-	-覧
---	-----	----	-----	------	----

ディレクトリ	サイズ	目的
/home		ホーム領域
/data	9.6PB	「データ領域(課題単位、利用負担金の対象)
/tmp_work		一時領域
/arc	申請量	階層型ストレージ領域(課題単位、新規受付停止)

各サーバから各ストレージへのアクセス可否を以下に示します。

衣 3-2 ストレーン視域へのアクセス可容一	表	3-2	ストレージ領域へのアクセス可否-	-覧
------------------------	---	-----	------------------	----

ストレージ種別	フロントエンドサーバ	MPC	LMC
共有ストレージ	0	0	0
階層型ストレージ	0	×	×

O:利用可能 ×:利用不能

3.1.1 共有ストレージ

共有ストレージはユーザーのホーム領域および課題単位のデータ領域となります。ホーム領域は、利用者あたり4TB ずつ割り当てられます。データ領域は利用負担金の対象となります。データ領域を利用されたい場合は、HBW2 ポータルから申請を行ってください。

	衣 5-5 六有ヘトレ Ju	の前の当て前限	
ディレクトリ	割当ディレクトリ	使用量制限	Inode 数制限
/home	/home/ユーザ名	4TB	1000 万
/data	/data/課題の GID(GID of project)	1TB~	1TB あたり 60 万

表 3-3 共有ストレージの割り当て制限

Copyright (C) RIKEN, Japan. All rights reserved.

3.1.2 階層型ストレージ

階層型ストレージ領域(テープ領域)はキャッシュストレージとテープライブラリ装置からな る長期保存用のストレージとなります。階層型ストレージはフロントエンドサーバから参照可 能であり、課題単位での利用となりますが、現在は新規の受付は行っていません。

表 3-4 階層型ストレージの割り当て制限

領域名	使用量制限	Inode 数制限
/arc	テープ 2 巻/8TB~テープ 13 巻/52TB)	テープ 1 巻あたり 40,000 ファイ ル

主な使用目的は以下のとおりです。

- 大量データの長期保存
- バックアップ



階層型ストレージはテープにデータを保管しますので、サイズの小さなファイル(100MB 🕂 未満)の格納には不向きです。多くのデータを格納する場合、サイズの小さな複数のフ アイルを1ファイルにアーカイブして格納してください。

3.2 ディスク使用状況と上限値の表示

ディスク使用状況と上限値を表示する場合は、Ifs quota コマンドを実行してください。 /data 領域の使用状況を表示する場合、-p オプションにて課題の GID(GID of project)を指 定してください。理研内課題の GID は申請時に通知されました"RB"から始まる課題番号 (project ID)の GID になります。HPCI 課題の GID は申請時に通知されました"HP"から始 まる課題番号の GID の先頭に 1 を付与し桁上げした数値となります。 課題の GID は id コマンドで確認することができます。

[username@hokusai1 ~]\$ Ifs quota -p \${UID} \${HOME}						
Disk quotas for usr username (uid XXXX	():					
Filesystem kbytes quota limi	t grace	files	quota	limit	grace	
/home/username 8637844 0	0 –	35783	0	0	-	
[username@hokusai1 ~]\$ Ifs quota -p 課題の GID/data/課題番号						
Disk quotas for prj XXXXX (pid XXXXX):						
Filesystem kbytes quota limi	t grace	files	quota	limit	grace	
/data/RB999999 8637844 0 200000	- 00000	35783	0	0	-	

表 3-5 ディスク使用状況と上限値の表示項目

項目	説明
[kbytes]	ディスクの使用容量と使用上限、使用率をサイズ表示(単位:KB)
	※-h オプションを指定してコマンド実行すると見やすい単位に
	変換して表示されます
[Files]	ファイル(inode)数の使用容量と使用上限、使用率を数量表示(単
	位:K)

3.3 一時領域

テンポラリファイルを格納する一時領域として、/tmp_work が利用できます。

/tmp_work 配下のデータについては、更新日が 1 週間以前のデータは自動的に削 除されます。テンポラリファイルだけを格納する領域として使用してください。

Copyright (C) RIKEN, Japan. All rights reserved.

4. コンパイル/リンク

4.1 環境設定

コンパイラ、ライブラリ、アプリケーション、ツールの利用時に必要となる環境変数は、 module コマンドを使用して容易に設定することができます。

\$ module <subcommand> <subcommand-args>

moduleコマンドのサブコマンドは以下のとおりです。

耒	4-1	module \exists	マンドの	サブコマ	マンドー	·眥
~						-

オプション	説明
avail	利用可能な設定一覧を表示します。
list	利用中の設定一覧を表示します。
load module	設定を読み込みます。
unload module	設定を解除します。
purge	全ての設定を解除します。
switch module1 module2	設定を切り替えます。

例)利用可能な設定一覧を表示する。

[username@hokusai1~];	\$ module avail	
/lustre/opt/modulefile	es/hokusai/apps	
adf/2022.103(default) gaussview/6.0.16 amber/22(default) gold/2023.2(default)	gamess/2023R2 gromacs/2023.3(default) gaussian/16_C.02(default) vmd/1.9.4	
/lustre/opt/modulefile	es/hokusai/compilers	
boost/1.83.0 gcc/13.2.0(default)	fftw-mpi-double/3.3.10 intelmpi/impi_23.2.0	

※(default)表示のバージョンは、HBW2 での推奨バージョンです。 バージョン指定無しでモジュールを load した場合は、(default)表示のバージョンの 設定を読み込みます。

例) Intel コンパイラを利用する場合

[username@hokusai1 ~]\$ module load intel

※intelを指定すると、intelmpiも合わせて設定を読み込みます。

Copyright (C) RIKEN, Japan. All rights reserved.

例)利用中の設定一覧を表示する。

[username@hokusai1 ~]\$ module list Currently Loaded Modulefiles: 1) intelmpi/impi_23.2.0 2) intel/23.02.1

排他関係にある設定は同時に利用することはできません。利用しようとした場合、以下の エラーメッセージが出力され、設定に失敗します。

[username@hokusai1 ~]\$ module load openmpi/4.1.6 Loading openmpi/4.1.6 ERROR: openmpi/4.1.6 cannot be loaded due to a conflict. HINT: Might try "module unload intelmpi" first.

例)利用中の設定を解除する。

[username@hokusai1 ~]**\$ module purge** [username@hokusai1 ~]**\$** module list No Modulefiles Currently Loaded.

Copyright (C) RIKEN, Japan. All rights reserved.

4.2 **コンパイラ**

HBW2 のフロントエンドサーバでは、超並列演算システム(MPC)と大容量メモリ演算サ ーバ(LMC)で実行可能なロードモジュールを作成するためのコンパイラが利用できます。

例) 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けのコンパイル/リン クを行う場合

[username@hokusai1 ~]\$ module load intel Loading intel/23.02.1 Loading requirement: intelmpi/impi_23.2.0 [username@hokusai1 ~]\$ module list Currently Loaded Modulefiles: 1) intelmpi/impi_23.2.0 2) intel/23.02.1

GPU サーバでは、GPU サーバ(GPU)で実行可能なロードモジュールを作成するための コンパイラが利用できます。

例) GPU サーバ(GPU)向けのコンパイル/リンクを行う場合

[username@hokusai1 ~]\$ srun --partition gpu_i --account RB999999 --pty \$SHELL
[username@bwgpu4 ~]\$ module load cuda/12.4
[username@bwgpu4 ~]\$ module list
Currently Loaded Modulefiles:
1) cuda/12.4

4.3 コンパイル/リンク概要

コンパイル/リンクのコマンド一覧は以下のとおりです。

表 4-2 超並列演算システム(MPC)、大容量メモリ演算サーバ(LMC) コンパイル/リンクコマンドー覧

種別	言語処理系	コマンド	自動並列 ^{注1}	OpenMP ^{注1}
逐次	Fortran	ifort		
(非 MPI)	С	icc		
	C++	ісрс	norollol	achonmh
MPI 並列	Fortran	mpiifort	-parallel	-dobeninb
	С	mpiicc		
	C++	mpiicpc		

注 1: 自動並列、OpenMP オプションはデフォルトでは無効です。

4.3.1 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けコンパイル/リンク 逐次プログラムのコンパイル/リンク

逐次プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けにコンパイル/リンクする場合、*ifort/icc/icpc*コマンドを使用します。

ifort/icc/icpc[option] file [...]

例1)逐次 Fortran プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC) 向けにコンパイル/リンク

[username@hokusai1 ~]\$ ifort seq.f

例2)逐次 C プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け にコンパイル/リンク

[username@hokusai1 ~]\$ icc seq.c

スレッド並列プログラムのコンパイル/リンク

スレッド並列プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けにコンパイル/リンクする場合、*ifort/icc/icpc*コマンドを使用します。

ifort/icc/icpc thread-option [option] file [...]

例 1) Fortran プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC) 向けに自動並列を有効にしてコンパイル/リンク

[username@hokusai1 ~]\$ ifort -parallel para.f

例 2) C プログラムをアプリケーション超並列演算システム(MPC)/大容量メモリ演算サー バ(LMC)向けに自動並列を有効にしてコンパイル/リンク

[username@hokusai1 ~]\$ icc -parallel para.c

例 3) OpenMP Fortran プログラムを超並列演算システム(MPC)/大容量メモリ演算サー バ(LMC)向けにコンパイル/リンク

[username@hokusai1 ~]\$ ifort -qopenmp omp.f

例 4) OpenMP C プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ (LMC)向けにコンパイル/リンク

[username@hokusai1 ~]\$ icc -qopenmp omp.c

Copyright (C) RIKEN, Japan. All rights reserved.

例 5) OpenMP Fortran プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けに自動並列も有効にしてコンパイル/リンク

[username@hokusai1 ~]\$ ifort -parallel -qopenmp omp_para.f

例 6) OpenMP C プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ (LMC)向けに自動並列も有効にしてコンパイル/リンク

[username@hokusai1 ~]\$ icc -parallel -qopenmp omp_para.c

MPI 並列プログラムのコンパイル/リンク

MPI 並列プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け にコンパイル/リンクする場合は、mpiifort/mpiicc/mpiicpc コマンドを使用します。

mpiifort/mpiicc/mpiicpc [option] file [...]

例 1) MPI Fortran プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ (LMC)向けにコンパイル/リンク

[username@hokusai1 ~]\$ mpiifort mpi.f

例 2) MPI C プログラムを超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向 けにコンパイル/リンク

[username@hokusai1 ~]\$ mpiicc mpi.c

例 3) ハイブリッド並列(MPI+OpenMP)Fortran プログラムを超並列演算システム (MPC)/大容量メモリ演算サーバ(LMC)向けにコンパイル/リンク

[username@hokusai1 ~]\$ mpiifort -qopenmp mpi_omp.f

例 4) ハイブリッド並列(MPI+OpenMP)C プログラムを超並列演算システム(MPC)/大容 量メモリ演算サーバ(LMC)向けにコンパイル/リンク

[username@hokusai1 ~]\$ mpiicc -qopenmp mpi_omp.c

オプションによるプログラムの最適化(Fortran/C/C++)

<u>最適化は演算結果に影響を与える場合があります。最適化を施した場合は実行結</u> <u>果が正しいかどうか検証してください。以下に記載するオプションが全てではありま</u> <u>せん。詳細は各コンパイルコマンドのマニュアル(man コマンド)をご参照ください。</u>

コンパイル オプション	説明
-O0	最適化は行われません。このオプションは、アプリケーション 開発の初期段階およびデバッグ時に使用します。アプリケー ションが正常に動作することを確認した後は、より高度なオプ ションを使用してください。
-01	サイズを考慮した最適化を行います。オブジェクトのサイズを 増やす傾向が ある最適化を省略します。多くの場合、最小 限のサイズで最適化されたコー ドが作成されます。 コードサイズが大きいため、メモリーページングが問題になる 巨大なサーバー / データベース・アプリケーションにおいて、 このオプションは効果的です。
-02	最速化します(デフォルト設定)。ベクトル化と実行速度を改善 する多くの 最適化を有効にします。多くの場合、/O1 (-O1) よりも速いコードを作成します。
-O3	-O2 の最適化に加えて、スカラー置換、ループアンロール、 分岐を除去するコード反復、効率的にキャッシュを使用す るループ・ブロッキング、データ・プリフェッチ機能など、強力 なループ最適化およびメモリアクセス最適化を行います。 -O3 オプションは、特に浮動小数点演算を多用するループや 大きなデータセットを処理するループを含むアプリケーション に推奨します。この最適化は、場合によって-O2 の最適化よ りもアプリケーションの実行が遅くなることがあります。
-fast	プログラム全体の速度を最大限にします。次のオプションを 設定します。 -ipo, -O3, -no-prec-div, -static, -fp-model fast=2, -xHost * -static は MPI プログラムのリンク時には指定できません。
-qopt-report[= <i>n</i>]	最適化レポートを作成し、標準エラー出力します。n には、 0(レポートなし)から 5(最大限の情報)の範囲で詳細レベルを 指定します。 デフォルトは 2 です。
-qopt-report- phase[= <i>list</i>]	最適化フェーズごとのレポートを生成します。指定可能な list は man を参照してください。
-qopt-report-help	上記の-qopt-report-phase で利用可能なすべてのフェーズ オプションを表示します。コンパイルは行いません。
-qopt-report- routine=string	string で指定された文字列を含む関数またはサブルーチンに 関するレポートを生成します。デフォルトでは、すべての関数 およびサブルーチンのレポートが生成されます。

表 4-3 一般的な最適化オプション

Copyright (C) RIKEN, Japan. All rights reserved.

コンパイル オプション	説明
-qopenmp	OpenMP 指示句がある場合、その指示によるマルチスレッ ド・コードが生成されます。スタックのサイズを増やさなければ ならないことがあります。
-parallel	自動パラレライザーは、安全に並列実行可能な構造のルー プを検出し、そのループに対するマルチスレッド・コードを自動 生成します。
-par-threshold[<i>n</i>]	パフォーマンス向上の可能性に基づいて、ループの自動並列 化のしきい値を設定します(n=0 から n=100。デフォルト: n=100)。 0 - 計算量にかかわらず並列化を行います。 100 - パフォーマンス向上が確実に見込める場合にのみ、 ループを並列化します。 -parallel とともに使用する必要があります。
-guide[= <i>n</i>]	ガイド付き自動並列化。ループのベクトル化や自動並列化を 促進するアドバイスを生成します。オブジェクト・ファイルや実 行ファイルは生成されません。自動並列化のアドバイスは、- parallel オプションが同時に指定された場合に提供されます。 n には 1 から 4 までの値で、アドバイスのレベルを指定しま す。4 は最も詳細なレベルのアドバイスです。n が省略された 場合、デフォルトは 4 です。
-qopt-matmul	コンパイラが生成する行列乗算(matmul)ライブラリ呼び出し を有効にします。行列乗算におけるループの入れ子(該当す る場 合)を特定し、ループを matmul ライブラリ呼び出しに置 換して、パフォーマンスを向上します。このオプションは、-O3 と-parallel が指定されると、デフォルトで有効になります。ま た、このオプションは、-O2 オプション以上が指定されない限 り、効果はありません。
-coarray=shared	共有メモリーシステムにおいて Fortran 2008 標準の Co- Array を有効にします(Fortran のみ)。

表 4-4 並列パフォーマンスオプション

表 4-5 プロセッサ専用最適化オプション

コンパイル オプション	説明
-xtarget	target で指定された命令セットをサポートするインテルプロセッサ向けの専用コードを生成します。 超並列演算システム(MPC)向けに最適化する場合、オプション-xCORE-AVX512を指定してください。 GPU 演算サーバ向けに最適化する場合、オプション- xCORE-AVX2 を指定してください。
-xhost	コンパイルを行うホスト・プロセッサで利用可能な最上位の命 令セッ トを利用したコードを生成します。これは、最適な-x オ プションに対応します。

Copyright (C) RIKEN, Japan. All rights reserved.

コンパイル オプション	説明
-ip	単ーファイルの最適化を行います。現在のソースファイルを 対象にしたインライン展開を含むプロシージャ間の最適化で す。
-ipo[= <i>n</i>]	インライン展開およびその他のプロシージャ間の最適化を複 数のソースファイルに対して行います。オプションの n 引数に は、コンパイル時に生成するオブジェクトファイルの最大数を 指定します。デフォルトの n は 0 です(コンパイラが最適なフ ァイル数を自動選択)。
-ipo-jobs[<i>n</i>]	プロシージャ間の最適化(IPO)のリンクフェーズで、同時に実 行するジョブ数を指定します。デフォルトは1ジョブです。
-finline-functions - inline-level=2	インライン展開をコンパイラの判断に任せます。このオプションは、-O2 および-O3 を指定すると有効になります。
-inline-factor=n	インライン化される関数の合計サイズと最大サイズを指定し ます。n のデフォルト値は 100(100%、スケール係数 1)です (Fortran のみ)。
-prof-gen	プロファイル最適化で参照する動的なパフォーマンスデータを 生成するため、プログラムにインストルメント・コードを埋め込 みます。
-prof-use	最適化中に-prof-gen オブションで生成した実行ファイルのプロファイリング情報を参照します。

表 4-6 プロシージャ間最適(IPO)オプションとプロファイルに基づく最適化オプション

表 4-7 浮動小数点演算最適化オプション

コンパイル オプション	説明
-fp-model <i>name</i>	浮動小数点演算における演算モデルを制御します。特定の 最適化を制限して浮動小数点結果の一貫性を強化します。
-ftz[-]	メインプログラムをこのオプションでコンパイルすると、プログ ラム全体で SSE 命令によるデノーマル結果をゼロにフラッシ ュします。-O0 を除き、デフォルトはオンです。
-fimf-precision: <i>name</i>	算術ライブラリ関数の精度を制御します。デフォルトはオフで す(コンパイラは、デフォルトのヒューリスティックを使用)。 nameの値は、high、medium、lowのいずれかです。精度を 下げるとパフォーマンスの向上につながります。また、その逆 についても同じことが言えます。算術ライブラリのルーチンの 多くは、インテル製マイクロプロセッサ向けにより高度に最 適化されています。
-fimf-arch- consistency= <i>true</i>	算術ライブラリ関数が、同じアーキテクチャの異なるインテル プロセッサにおいて一貫 した結果を生 成します。ランタイ ム・パフォーマンスが低下することがあります。デフォルトは "false"(オフ)です。
-prec-div	浮動小数点除算の精度を上げます。精度を上げるとパフォー マンスに影響します。
-prec-sqrt	平方根計算の精度を上げます。精度を上げるとパフォーマン スに多少影響します。

コンパイル オプション	説明
-unroll[<i>n</i>]	ループをアンロールする最大回数を設定します。-unroll0 は ループアンロールを無効にします。デフォルトは、-unroll で、 コンパイラの判断に任せます。
-qopt-prefetch[=n]	ソフトウェア・プリフェッチのレベルを制御します。n には 0(プ リフェッチなし)から 4(強力なプリフェッチ)までの値です。高レ ベルの最適化が有効な場合、デフォルトは 2 です。
-qopt-block-factor= <i>n</i>	ループのフロッキング最適化におけるフロッキング係数を指定します。 ブロッキング係数は、1 ブロックのループ反復数です。 ループ・ブロッキングは-O で有効になり、キャッシュ中の データ利用率を高めるように設計されています。
-qopt-streaming-stores =mode	ストリーミング・ストアを生成するかどうかを指定します。 modeの値は次のとおりです。 always: アプリケーションがメモリ上のデータ再利用がほとん どないことを仮定して、キャッシュをバイパスするストリーミン グ・ストアの生成をコンパイラに指示します。 never: ストリーミング・ストアの生成を無効にします。 auto: ストリーミング・ストアの生成をコンパイラの判断に任せ ます。
-fno-alias	ブログラムでエイリアシングしないことを前提に最適化しま す。デフォルトではオフです。
-fno-fnalias	関数内でエイリアシングしないことを前提に最適化します。デ フォルトではオフです。
-fexceptions	-fexceptions: C++のデフォルトです。例外処理テーブルの生 成を有効にします。 -fno-exceptions: C、Fortran のデフォルトです。コードサイズ が小さくなります。C++では、例外指定は解析されますが、無 視されます。構造化例外処理(try ブロックや throw 文)を使用 していると、一連の呼び出し中の関数が-fno-exceptions でコ ンパイルされている場合、エラーが発生します。
-vec-threshold n	パフォーマンス向上の可能性に基づいて、ループのベクトル 化のしきい値 n を設定します(n=0 から n=100。デフォルト : n =100)。 0 - 計算量にかかわらずベクトル化されます。 100 - パフォーマンス向上が確実に見込める場合にのみル ープはベクトル化されます。

表 4-8 きめ細かなチューニングオプション

4.4 数値計算ライブラリ

数値計算ライブラリとして BLAS/LAPACK/ScaLAPACK が利用可能です。

4.4.1 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)の数値計算ライブラリ

超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け BLAS/LAPACK/ScaL APACK を利用する場合、以下のコンパイル/リンクオプションが利用可能です。

及 + DEAD/EAT ADIVOCALAT ADIV J ノノヨノ 見			
ライブラリ	並列性	指定オプション	備考
BLAS	逐次	-qmkl=sequential	
	スレッド並列	-qmkl=parallel	
LAPACK	逐次	-qmkl=sequential	
	スレッド並列	-qmkl=parallel	
ScaLAPACK	MPI 並列	-qmkl=cluster	

表 4-9 BLAS/LAPACK/ScaLAPACK オプション一覧

- 例 1) 逐次版 BLAS/LAPACK を利用する。
 - **\$** ifort <u>-qmkl=sequential</u> blas.f
- 例 2) スレッド並列版 BLAS/LAPACK を利用する。
- \$ ifort <u>-qmkl=parallel</u> blas.f
- 例 3) ScaLAPACK を利用する(逐次版 BLAS/LAPACK をリンク)。
 - \$ mpiifort __qmkl=cluster scalapack.f

上記以外の組み合わせについては、以下 URL を参照してください。

https://software.intel.com/en-us/articles/intel-mkl-link-line-advisor

Copyright (C) RIKEN, Japan. All rights reserved.

4.5 開発ツール

超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けに以下のツールが利用できます。

- インテル VTune Amplifier XE(パフォーマンス分析ツール)
- インテル Inspector XE(エラー検出ツール)
- インテル Advisor XE(スレッド化アシスタント)
- インテル Trace Analyzer & Collector(MPI 通信パフォーマンス分析ツール)

5. バッチジョブ・インタラクティブジョブ

5.1 ジョブ管理システム概要

HBW2 ではジョブ管理システムによりバッチジョブおよびインタラクティブジョブの実行が 制御されます。ユーザーはジョブ実行に必要なパーティション名、ノード数、コア数、経過時 間などを指定し、ジョブ管理システムにジョブ実行を依頼します。

HBW2 で実行可能なジョブは以下の2種類です。

ジョブ形式	用途
バッチジョブ	バッチ形式でジョブを実行する。
	ノードダウンなどの異常発生時、オプションによりジョブの自動再実
	行が可能。
インタラクティブジ	ユーザー端末からのデータ入力により、会話形式でジョブを実行す
ョブ	る。
	主にデバッグ等での利用を想定。
	ノードダウンなどの異常発生時、ジョブは再実行されない。

表 5-1 ジョブの種類

バッチジョブは、以下の3種類に分類されます。

表 5-2 バッチジョブの種類

バッチジョブ 種別	用途	投入形式
通常ジョブ	スクリプト単位でジョブを実行する	「5.5.1 通常ジョブ投入」参 照
ステップジョブ	投入した複数のジョブを 1 つのまと まりとして扱い、その中で実行順 序、依存関係をもつジョブ	「5.5.2 ステップジョブ投入」 参照
アレイジョブ	1 度のジョブ投入で、複数の同じ通 常ジョブを投入し実行するジョブ	「5.5.3 アレイジョブ投入」参照

ユーザーがジョブ操作に用いるコマンドは、以下のとおりです。

表 5-3 ジョブ操作コマンド一覧

機能	コマンド名	
ジョブ投入	sbatch/srun	
ジョブ参照	squeue/sacct	
ジョブ削除	scancel	

5.2 ジョブ実行リソース

ジョブを投入する際は、ジョブを実行するハードウェアを指定する「パーティション」を選択 します。

5.2.1 パーティション

パーティションはジョブを実行するハードウェアを指定するオプションです。

パーティションにおける課題ごとの設定

各課題は、以下の計算資源を同時に使用することができます。

表 5-4 各サブシステムにおける課題ごとの設定

ジョブ実行先:サブシステム		実行コア数	投入ジョブ本数
超並列演算システム	(MPC)	7168	5000
大容量メモリ演算サーバ	(LMC)	192	100
GPU サーバ	(GPU)	448	100

パーティション	ノード数	最大 経過時間	最大 コア数	最大 ノード数	ノードあたりの メモリ量
mpc	312	24h	7168	64	112GiB
mpc_l	156	72h	896	8	112GiB
Imc	2	24h	192	2	2880GiB
gpu	3	72h	112	1	448GiB
gpu_i	1	24h	112	1	448GiB

表 5-5 パーティションごとの設定

*デフォルト経過時間は 1h です。

5.3 ジョブ投入オプション

ジョブ投入時は、ジョブの実行に応じて、以下の3つのオプションを指定します。

- 基本オプション
- ジョブ資源オプション

5.3.1 基本オプション

ジョブ投入時に指定する基本オプションは以下のとおりです。

衣 5-6 ジョノ技入基本オノジョノ		
オプション	説明	
account <project id=""></project>	ジョブ実行に利用する課題番号(project ID)を指定 (複数課題を持つユーザーは指定必須)	
mail-user	メールの送信先を指定	
mail-type	メール通知を指定	
BEGIN	ジョブ開始時にメール通知を指定	
END	ジョブ終了時にメール通知を指定	
REQUEUE	ジョブ再実行時にメール通知を指定	
job-name <name></name>	ジョブ名を指定	
output <filename></filename>	標準出力を指定されたファイルに出力 (-e 指定無しの場合、標準エラー出力も同じファイル に出力)	
error <filename></filename>	標準エラー出力を指定されたファイルに出力	
requeue	障害発生時ジョブを再実行することを指定 (デフォルト:no-requeue)	
dependency	ステップジョブとしてジョブを投入	
afterany: <jobid></jobid>	指定したジョブ(jobid)終了後に投入	
afterok: <jobid></jobid>	指定したジョブ(jobid)が正常終了後に投入	
afternotok: <jobid></jobid>	指定したジョブ(jobid)が異常終了後に投入	
after: <jobid[+minute]></jobid[+minute]>	指定したジョブ(jobid)開始後時間の経過で投入	
singleton	同ーユーザー・同ージョブ名のジョブを逐次実行	
array <start-end></start-end>	アレイジョブとしてジョブを投入	
export= <env></env>	ジョブ投入時の環境変数をジョブ実行環境に継承	
begin= <time></time>	ジョブ開始時刻を指定	

表 5-6 ジョブ投入基本オプション

5.3.2 ジョブ資源オプション

ジョブが利用する資源に関する主要オプションは以下のとおりです。

<u>д</u> ,	
オプション	説明
partition <partition_name></partition_name>	キュー(パーティション)を指定
nodes <num></num>	ノード数を指定
ntasks <num></num>	プロセス数を指定
ntasks-per-node <num></num>	ノードあたりのプロセス数を指定
cpus-per-task <num></num>	プロセスあたり確保する論理 CPU コア数を指定
ntasks-per-core= <num></num>	CPU コアあたりのプロセス数を指定
mem= <num[m g t]></num[m g t]>	ノードあたりのメモリ量を指定
mem-per- cpu= <num[k m g t]></num[k m g t]>	CPU コアあたりのメモリ量を指定 超並列演算システム(MPC)の最大: 112G 大容量メモリ演算サーバ(LMC)の最大: 2880G
cpus-per-gpu= <num></num>	GPU あたりの CPU コア数を指定
mem-per- gpu= <num[k m g t]></num[k m g t]>	GPU あたりのメモリ量を指定 GPU サーバ(GPU)の最大: 448G
gpus= <num></num>	GPU の使用枚数を指定 ※GPU あたり CPU:28 コア、メモリ:112G
time	経過時間を指定
<mm></mm>	分で指定
<mm:ss></mm:ss>	分、秒で指定
<hh:mm:ss></hh:mm:ss>	時間、分、秒で指定
<days-hh></days-hh>	日数、時間で指定
<days-hh:mm></days-hh:mm>	日数、時間、分で指定
<days-hh:mm:ss></days-hh:mm:ss>	日数、時間、分、秒で指定

表 5-7 ジョブ資源オプション(共通)

メモリ量を指定する場合、以下の単位が使用可能です。

表 5-8 メモリ量で使用可能な単位

単位	説明
KiB	キビバイト
MiB	メビバイト
GiB	ギビバイト
TiB	ティビバイト

コアあたりのメモリ割当量のデフォルトは、以下のとおりです。

表 5-9 コアあたりのメモリ割当量のデフォルト

システム	コアあたりのメモリ割当量	
超並列演算システム(MPC)	1G	
「大容量メモリ演算サーバ(LMC)	30G	
GPU サーバ(GPU)	4G	

Copyright (C) RIKEN, Japan. All rights reserved.



コアあたりのメモリ割当量より多いメモリ量を要求した場合、要求メモリ量に応じて演 🔼 算時間が計算されるので注意してください。HBW2システムよりCPUコアあたりのメ モリ量は少なくなっています。

生成するプロセス数やスレッド数は MPI オプションや環境変数 OMP NUM THREADS を使用して明に指定するようにしてください。省略すると意図しないプロセス数、スレッド数 でプログラムが実行される可能性があり、性能劣化を招く恐れがあります。

5.4 インタラクティブジョブ実行

インタラクティブジョブを実行するためには、srun コマンドを実行します。ジョブ管理シス テムがインタラクティブジョブに対して計算資源を割り当てることでジョブの実行が開始され ます。

計算ノードにログインするインタラクティブジョブの投入は、ジョブ投入オプションは引数と して以下の通り指定します。

srun --partition パーティション名 --account projectname --pty \$SHELL

ログインノードではプロセス数制限があるため、多数のプロセス等のリソースが必要なコンパイルを実行する場合は、インタラクティブジョブを利用されることを推奨します。

例1) 超並列演算システム(MPC)へログインしプログラムを実行する例

```
[username@hokusai1 ~]$ srun --account RB9999999 --pty $SHELL
[username@bwmpc001 ~]$ ./a.out
Hello world
[username@bwmpc001 ~]$ exit
exit
[username@hokusai1 ~]$
```

フロントエンドサーバから計算ノードに直接プログラムを実行する場合、ジョブ投入オプションは引数として以下の通り指定します。ジョブ投入オプションは実行するプログラムに合わせ、「5.3 ジョブ投入オプション」を参照し指定してください。

srun --partition=mpc --account projectname [option] 実行プログラム

例2) 超並列演算システム(MPC)へ直接プログラムを実行する例

[username@hokusai1 ~]\$ module load intel [username@hokusai1 ~]\$ srun --partition=mpc --account RB999999 --nodes=1 -ntasks=2 ./a.out Hello world [username@hokusai1 ~]\$



srun コマンド実行前に実行するプログラムに合わせ、環境設定を実施してください。

なお、GPU サーバでインタラクティブジョブを実行する際は、 partition=gpu_i を指定してください。

Copyright (C) RIKEN, Japan. All rights reserved.

5.5 バッチジョブ投入

バッチジョブを実行するためには、実行するプログラムとは別に「ジョブスクリプト」を作成 し、利用するパーティション、経過時間、ノード数などの資源を記載したオプションを記述し た上で、実行するコマンド列を記載します。ユーザーはジョブスクリプトを sbatch コマンドで 投入します。投入されたジョブはジョブ管理システムにより資源の空き状況やジョブ間の優 先度などを元に自動で実行開始されます。

5.5.1 通常ジョブ投入

通常ジョブを投入する場合、sbatchコマンドの引数にバッチジョブとして実行するジョブス クリプトを指定します。

sbatch [option] [job-script]

- ジョブスクリプトを指定しない場合、標準入力から実行命令を読み込みます。
- ジョブ投入オプションは、ジョブスクリプト内もしくは標準入力にディレクティブを用いて 指定することが可能です。
- ジョブ投入が完了後、ジョブに対して識別用 ID(ジョブ ID)が割り当てられます。
- 例) 通常ジョブ投入例

[username@hokusai1~]\$ **sbatch run.sh** Submitted batch job 1234.

5.5.2 ステップジョブ投入

ステップジョブは、複数のバッチジョブを1つのまとまりとして扱い、その中で実行の順序 関係や依存関係を指定することで、ジョブチェイン機能を実現するジョブモデルです。ステッ プジョブは複数サブジョブから構成され、各サブジョブは同時に実行されることはありませ ん。ステップジョブの動作イメージを以下に示します。



図 5-1 ステップジョブイメージ

Copyright (C) RIKEN, Japan. All rights reserved.

ステップジョブの投入形式は以下のとおりです。

sbatch --dependency=[option] [job-script]

例 1-1) 指定したジョブ(jobid)が終了後にジョブを投入

[username@hokusai1 ~]\$ sbatch step1.sh Submitted batch job 1234 [username@hokusai1 ~]\$ sbatch --dependency=afterany:1234 step2.sh Submitted batch job 1235 [username@hokusai1 ~]\$ sbatch --dependency=afterany:1235 step3.sh Submitted batch job 1236

例 1-2) 同ーユーザー・同ージョブ名のジョブを逐次実行(同時に1本のみ実行可能)

[username@hokusai1 ~]\$ sbatch step1.sh
Submitted batch job 1234
[username@hokusai1 ~]\$ sbatch --dependency=singleton step2.sh
Submitted batch job 1235

例 1-3) 指定したジョブ(jobid)が正常終了後にジョブを投入

```
[username@hokusai1 ~]$ sbatch step1.sh
Submitted batch job 1234
[username@hokusai1 ~]$ sbatch ---dependency=afterok:1234 step2.sh
Submitted batch job 1235
```

例 1-4) 指定したジョブ(jobid)が異常終了後にジョブを投入

[username@hokusai1 ~]\$ sbatch step1.sh
Submitted batch job 1234
[username@hokusai1~]\$ sbatchdependency=afternotok:1234 step2 .sh

例 2) 複数の条件を指定してステップジョブを投入

[username@hokusai1 ~]\$ sbatch step1.sh Submitted batch job 1234 [username@hokusai1 ~]\$ sbatch ---dependency=afterok:1234 step2.sh Submitted batch job 1235 [username@hokusai1 ~]\$ sbatch ---dependency=afterok:1234, afternotok:1236 step3.sh Submitted batch job 1236 [username@hokusai1 ~]\$ sbatch ---dependency=afterany:1234:1235, after:1236+1 step4.sh Submitted batch job 1237

オプション	説明
afterany: <jobid></jobid>	指定したジョブが終了後に投入
afterok: <jobid></jobid>	指定したジョブが正常終了後(終了ステータ
	スが 0)に投入
afternotok: <jobid></jobid>	指定したジョブが異常終了後(終了ステータ
	スが0以外)に投入
after: <jobid[+minute]></jobid[+minute]>	指定したジョブの開始後、指定した時間の
	経過後に投入
singleton	同ーユーザー・同ージョブ名のジョブを逐次
	実行

表 5-10 ステップジョブオプション

Copyright (C) RIKEN, Japan. All rights reserved.

5.5.3 アレイジョブ投入

アレイジョブは、複数の同じバッチジョブを同時に投入し、実行するジョブです。 例えば、ジョブのパラメータを変えて、それぞれの実行結果を確認したい場合、通常のバ ッチジョブであれば、ジョブを1つ1つ投入する必要がありますが、アレイジョブを使えば、1 度に複数パターンの投入ができます。

アレイジョブの投入形式は以下のとおりです。

sbatch --array 0-2:1 jobscript

アレイジョブでは、サブジョブごとにジョブの入出力ファイルを変えられるようにジョブスク リプトを作成する必要があります。このために、サブジョブごとに設定されるアレイ番号を使 います。アレイ番号は環境変数 SLURM_ARRAY_TASK_ID で参照することが可能です。 各アレイジョブには、ジョブ ID_タスク ID という形式のジョブ ID に加え、通常のジョブと 同様のジョブ ID も付与されます。ただし、通常のジョブ ID はジョブの実行が開始され

た順に付与されるため、タスク ID の順番とジョブ ID の順番が一致しないことがあります。

5.5.4 バッチジョブの標準出力と標準エラー出力

バッチジョブの標準出力ファイルと標準エラー出力ファイルは、ジョブ投入ディレクトリもし くは、ジョブ投入時に指定されたファイルに出力されます。

標準出力ファイルにはジョブ実行中の標準出力、標準エラー出力ファイルにはジョブ実 行中のエラーメッセージが出力されます。ジョブ投入時に、標準出力ファイル、標準エラー 出力ファイルを指定しなかった場合は、以下のファイルに出力されます。

slurm-XXXXX.out --- 標準出力・エラーファイル (XXXXX はジョブ投入時に割り当てられたジョブ ID)

Copyright (C) RIKEN, Japan. All rights reserved.

5.5.5 ジョブスクリプト記述

バッチジョブを投入するためには、vi コマンドや emacs コマンドにてジョブスクリプトを作成します。

(1) 先頭行は "#!" に続けて、シェルのパスを指定してください。

[記述例]

#!/bin/bash



ログインシェルが bash 以外の場合、ジョブスクリプト内で module コマンドを実行す る場合は、「#!/bin/sh -l」を指定する必要があります。

(2) 2 行目以降にジョブ投入オプションをディレクティブ"#SBATCH"を用いて指定します。[記述例]

#SBATCHnodes=1	ノード数を指定
#SBATCHtime=1:00:00	経過時間を指定

(3) ジョブ投入オプションに続けて、実行時の環境変数設定と、プログラム実行を指定します。Slurmでは srun コマンドで実行プログラム(task)を指定することにより、ジョブ投入オプションに従ってノードやプロセス等のリソース割当が行われます。
[記述例]

export OMP_NUM_THREADS=20	環境変数を設定
srun ./a.out	プログラムを実行

5.5.6 MPI プログラムの実行

srun オプション

ジョブスクリプトファイルに記載しているジョブ投入オプションは引き継がれます。そのため、srun コマンドも sbatch 同様のオプションがありますが、MPI プログラム実行時にオプションを記載する必要はありません。

	オプション	説明
cpu	-bind= <type></type>	タスクを CPU にバインドします
	verbose	実行する前にバインディングを詳細に報告する
	rank	タスクランクごとに自動的にバインドする
	rank_ldom	ランクごとに NUMA ローカリティドメインにバインドする ※ノード内の全コア使用時のみ指定可能

表 5-11 srun オプション

5.6 バッチジョブスクリプト例

5.6.1 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向けジョブスクリプト

5.6.1.1 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け単一ノード逐次 ジョブ用スクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: mpc
・ノード数	: 1 ノード
・プロセス数(スレッド数)	: 1 プロセス(1 スレッド)
·経過時間	: 60 分
▪課題番号(project ID)	: RB999999

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.2 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)単一ノードスレッド並 列スクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: mpc
・ノード数	: 1 ノード
・プロセス数(スレッド数)	: 1 プロセス(10 スレッド)
・経過時間	: 60 分
▪課題番号(project ID)	: RB999999

[username@hokusai1 ~]\$ vi mpc-para.sh
#!/bin/sh
slurm option#
#SBATCHpartition=mpc
#SBATCHaccount=RB999999
#SBATCHnodes=1
#SBATCHcpus-per-task=10
<pre>#SBATCHtime=1:00:00</pre>
Program execution#
module load intel
export OMP_NUM_THREADS=\${SLURM_CPUS_PER_TASK} sruncpus-per-task=\${SLURM_CPUS_PER_TASK} ./a.out



超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)でスレッド並列ジョブを実 行する場合、環境変数 OMP_NUM_THREADS の指定は必須です。指定するメモリ 量により、確保するコア数が変わる場合があり、省略すると意図しないスレッド数で実 行される可能性があります。

OMP_NUM_THREADS :スレッド数(環境変数\${SLURM_CPUS_PER_TASK} を指定、または cpus-per-task で指定した値を記載する)

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.3 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け単一ノード MPI 並列ジョブスクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: mpc
・ノード数	: 1
・プロセス数(スレッド数)	: 20 プロセス(1 スレッド)
・経過時間	: 60 分
▪課題番号(project ID)	: RB999999

[username@hokusai1~]\$ vi mpc-single-mpi.sh
#!/bin/bash
slurm option#
#SBATCHpartition=mpc
#SBATCHaccount=RB999999
#SBATCHnodes=1
#SBATCHntasks-per-node=20
#SBATCH

#----- Program execution -----#

module load intel

srun ./a.out



超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)で MPI 並列ジョブを実 🔼 行する場合、指定するメモリ量により、確保するコア数が変わる場合があり、省略す ると意図しないプロセス数で実行される可能性があります。

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.4 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け単一ノードハイ ブリッド並列ジョブスクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: mpc
・ノード数	: 1
・プロセス数(スレッド数)	: 2 プロセス(10 スレッド)
・コア数	: 20 コア(2 x 10)
・経過時間	: 60 分
▪課題番号(project ID)	: RB999999

[username@hokusai1 ~]\$ vi mpc-single-hybrid.sh #!/bin/bash #----- slurm option ------# #SBATCH --partition=mpc #SBATCH --account=RB999999 #SBATCH --nodes=1 #SBATCH --ntasks-per-node=2 #SBATCH --ntasks-per-node=2 #SBATCH --cpus-per-task=10 #SBATCH --time=1:00:00 #----- Program execution -----# module load intel export OMP_NUM_THREADS=\${SLURM_CPUS_PER_TASK} srun --cpus-per-task=\${SLURM_CPUS_PER_TASK} ./a.out



超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)でハイブリッド並列ジョブ を実行する場合、環境変数 OMP_NUM_THREADS は指定必須です。指定するメモ リ量により、確保するコア数が変わる場合があり、省略すると意図しないプロセス数/ スレッド数で実行される可能性があります。

OMP_NUM_THREADS:スレッド数(環境変数\${SLURM_CPUS_PER_TASK}を 指定、または cpus-per-task で指定した値を記載する)

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.5 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け複数ノード MPI 並列ジョブスクリプト

以下のジョブ実行を想定した記述方法を説明します。

: mpc
: 2
: 80 プロセス(1 スレッド)
: 40 プロセス
:90 分
: RB999999

[username@hokusai1~]\$ vi mpc-multi-mpi.sh
#!/bin/sh
slurm option#
#SBATCHpartition=mpc
#SBATCHaccount=RB999999
#SBATCHnodes=2
#SBATCHntasks-per-node=40
#SBATCHtime=1:30:00
Program execution#
module load intel
srun ./a.out



超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)で MPI 並列ジョブを実行 する場合、指定するメモリ量により、確保するコア数が変わる場合があり、省略すると 意図しないプロセス数で実行される可能性があります。

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.6 超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)向け複数ノードハイ ブリッド並列ジョブスクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: mpc
・ノード数	: 2
・プロセス数(スレッド数)	: 4 プロセス(20 スレッド)
・ノードあたりのプロセス数	:2 プロセス
•経過時間	:1 時間 30 分
▪課題番号(project ID)	: RB999999

[username@hokusai1 ~]\$ vi mpc-multi-hybrid.sh #!/bin/sh #----- slurm option -----# #SBATCH --partition=mpc #SBATCH --partition=mpc #SBATCH --account=RB999999 #SBATCH --nodes=2 #SBATCH --ntasks-per-node=2 #SBATCH --ntasks-per-node=2 #SBATCH --cpus-per-task=5 #SBATCH --time=1:30:00 #----- Program execution -----# module load intel export OMP_NUM_THREADS=\${SLURM_CPUS_PER_TASK} srun --cpus-per-task=\${SLURM_CPUS_PER_TASK} ./a.out



超並列演算システム(MPC)/大容量メモリ演算サーバ(LMC)でハイブリッド並列ジョブ を実行する場合、環境変数 OMP_NUM_THREADS は指定必須です。指定するメモ リ量により、確保するコア数が変わる場合があり、省略すると意図しないプロセス数/ スレッド数で実行される可能性があります。

OMP_NUM_THREADS :スレッド数(環境変数\${SLURM_CPUS_PER_TASK}を 指定、または cpus-per-task で指定した値を記載する)

Copyright (C) RIKEN, Japan. All rights reserved.

5.6.1.7 GPU サーバ(GPU)向け単一ノード逐次ジョブ用スクリプト

以下のジョブ実行を想定した記述方法を説明します。

・パーティション	: gpu
・ノード数	: 1 ノード
・プロセス数(スレッド数)	: 1 プロセス(1 スレッド)
·経過時間	:60 分
・課題番号(project ID)	: RB999999

[username@hokusai1~]\$ vi gpu-seq.sh
#!/bin/bash
slurm option#
#SBATCHpartition=gpu
#SBATCHgpus=1
#SBATCHaccount=RB999999
#SBATCHnodes 1
#SBATCHtime 01:00:00
Program execution#
module load intel
srun ./a.out

Copyright (C) RIKEN, Japan. All rights reserved.

5.7 ジョブ状態表示

投入したジョブ状態やリソース情報を確認する場合、squeue コマンドを使用します。

	squeue [option]	[JOBID [JOBID …]]	
--	-----------------	-------------------	--

表 5-12 squeue コマンドオプション一覧

オプション	説明
なし	実行待ち、実行中のジョブ情報を表示
long	実行待ち、実行中の詳細なジョブ情報を表示
account <account></account>	指定した課題のジョブ情報を表示
jobs <jobid></jobid>	指定したジョブ ID のジョブ情報を表示
partition <partitionname></partitionname>	指定したパーティションのジョブ情報を表示
states <states></states>	指定したジョブステータスのジョブ情報を表示
nodelist <nodename></nodename>	指定したノードに割り当てられたジョブ情報を表示
iterate <seconds></seconds>	指定した間隔(秒単位)でジョブ情報を表示

Copyright (C) RIKEN, Japan. All rights reserved.

5.7.1 ジョブ状態表示

squeue コマンドを実行すると、現在実行中もしくは実行待ちのジョブ状態を表示します。

[username@hokusai1	~]\$ squeue						
JOBID	PARTITION	NAME	USER	ST	TIME	NODES	NODELIST (REASON)
1234	mpc	job. sh	username	R	0:04	1	bwmpc185

表 5-13 ジョブ情報の表示項目

項目	説明
JOBID	ジョブ ID
PARTITION	ジョブのキュー名
NAME	ジョブ名
USER	ユーザー名
ST	ジョブの現在の状態(表 5-14 ジョブの状態一覧参照)
TIME	ジョブの実行時間
NODES	ジョブ実行に使用するノード数
NODELIST	ジョブが実行されるホスト名のリスト

表 5-14 ジョブの状態一覧

状態	説明
PD	ジョブ実行待ち
CF	ジョブ実行に必要な資源を獲得中
R	ジョブ実行中
CG	ジョブの終了待ち
CD	ジョブ終了
CA	ジョブキャンセルによる終了
F	ゼロ以外の終了コードまたはその他の障害状態で終了
NF	割り当てられたノードのいずれかの故障による終了
ТО	実行時間制限による終了
RJT	投入が受付けられなかった状態

Copyright (C) RIKEN, Japan. All rights reserved.

5.7.2 詳細ジョブ情報の表示(sacct --format オプション)

sacct コマンドの--format オプションを指定すると、指定した項目のジョブ情報を表示します。表示する項目は--helpformat オプションで確認することができます。

[username@h	okusai1	~]\$		sacct	-X	
format=JobI	D, Partition, Acc	ount,Submit,S	start, End, CPUTim	eRAW, Alloc [°]	「res%40,State −j 1440	
JobID	Partition	Account	Subm	it	Start	End
CPUTimeRAW			AllocTRES	State		
1234	mpc	RB999999	2023-11-16T13:0	4:05 2023-	11-16T13:04:18 2023-11-1	16T13:09:19
301	billing=8, cpu=	1, mem=8800M, n	ode=1 COMPLETE	D		

表 5-15 sacct コマンドオプション一覧

オプション	説明
-X	ジョブ割り当て自体に関連する統計のみを表示
accounts	指定した課題の情報を表示
jobs	指定したジョブ ID の情報を表示
helpformat	fomat オプションで表示する項目の確認
fomat	指定した情報を表示
starttime	指定した日時(YYYY-MM-DD[THH:MM[:SS]])からのジョブ情報を表示

表 5-16	詳細ジョブ	「情報の	表示項目
--------	-------	------	------

項目	説明
Account	課題番号(project ID)
Partition	パーティション
Submit	ジョブの投入日時
Start	ジョブの開始時間
End	ジョブ終了時間
Elapsed	実行経過時間
CPUTimeRAW	経過時間(CPU 秒)
AllocTres	割り当てられたリソース
billing	演算時間計算に用いられるコア数(実際に使用されたコア数)
сри	割り当てコア数
mem	割り当てメモリ量
node	割り当てられたノード数
State	ジョブの状態

5.7.3 終了ジョブ情報の表示(sacct コマンド)

sacct コマンドを実行すると、投入済みジョブに加え、既に実行が終了したジョブも表示されます。--starttime オプションを指定すると、コマンド実行日以前の情報を表示することが可能です。

[username@hokusai1~]\$ sacct						
JobID	JobName	Partition	Account	AllocCPUS	State	ExitCode
 12345	test_job	mpc	username	4	COMPLETED	0:0

表 5-17 終了ジョブ†	青報の表示項目
---------------	---------

項目	説明
JobID	ジョブ ID
JobName	ジョブ名
Partition	ジョブ実行したパーティション
Account	課題番号(project ID)
AllocCPUS	割り当て CPU 数
State	ジョブの終了状態
ExitCode	ジョブの終了コード

Copyright (C) RIKEN, Japan. All rights reserved.

5.7.4 パーティションの表示

sinfo コマンドを実行すると、ユーザーが利用可能なパーティションを表示します。

[username@) hokusai	1 ~]\$ sinfo	I		
PARTITION	AVAIL	TIMELIMIT	NODES $(A/I/O/T)$	NODELIST	
mpc*	up 1·	-00:00:00	0/312/0/312	bwmpc[001-312]	
mpc_l	up 3-	-00:00:00	0/156/0/156	bwmpc[001-156]	
lmc	up 1·	-00:00:00	0/2/0/2	bwlmc[1-2]	

表 5-18 /	パーティ	ションの	表示項目
----------	------	------	------

項目	説明
PARTITION	キュー名(パーティション名)
AVAIL	キューの状態
TIMELIMIT	最大実行時間
NODES	ノードの状態を表示します。
	A:allocated
	l:idle
	O:other
	T:total
NODELIST	キュー(パーティション)に割り当てられたノード

Copyright (C) RIKEN, Japan. All rights reserved.

5.7.5 ジョブ投入時に指定可能なジョブ資源の上限表示(sinfo コマンド)

sinfo コマンドの「--format」オプションを指定すると、パーティション毎に利用可能なノード数、ノード当たりの CPU 数、最大経過時間を表示します。

以下のように表示された場合、超並列演算システム(MPC)の mpc パーティションでは 1 から 12 ノードが利用でき、ジョブは最大 24 時間(1 日)、1 ノード当たりの CPU 数は 112 コ アで実行できることを示しています。

<u>表示された上限を超えた値を指定したジョブも投入することが可能ですが、ジョブ実</u> 行はされません。そのため、scancelコマンドを実行しジョブを削除してください。

PARTIT	AVAIL	TIMELIMIT	JOB_SIZE	NODES	MAX_CPUS_PER_NODE	MEMORY	
mpc*	up	1-00:00:00	1-12	260	112	128300	
mpc_l	up	3-00:00:00	1-6	156	112	128300	
lmc	up	1-00:00:00	1	2	96	3095859	

項目	説明
PARTIT	パーティション
AVAIL	パーティションの状態表示
TIMELIMIT	最大経過時間
JOB_SIZE	割り当てノード数
NODES	ノード数
MAX_CPUS_PER_NODE	ノードあたりの最大 CPU 数
MEMORY	ノード当たりのメモリ量(MB)

表 5-19 ジョブ資源の上限表示項目

Copyright (C) RIKEN, Japan. All rights reserved.

5.8 ジョブキャンセル

投入済みのジョブをキャンセルする場合、scancel コマンドを実行します。

scancel JOBID [JOBID…]

キャンセルするジョブのジョブ ID を scancel の引数に指定します。

[username@hokusai1 ~]\$ scancel 12345

5.9 環境変数

5.9.1 スレッド並列、MPI 実行に関する環境変数

スレッド並列および MPI プログラム実行時に指定可能な主な環境変数について説明します。

	説明
OMP_NUM_THREADS	OpenMP もしくは自動並列によりスレッド並列化された プログラムを実行する場合は、環境変数 OMP_NUM_THREADS にスレッド数を指定します。
KMP_AFFINITY	スレッドのコアバインドを制御します。
I_MPI_PIN_DOMAIN	MPI プロセスのコアバインドを制御します。

表 5-20 実行時環境変数

5.9.2 ジョブに関する環境変数

ジョブ内では、ジョブ管理システムにより、以下の環境変数が設定されます。

環境変数	説明
SLURM_JOB_ID	ジョブ ID
SLURM_JOB_NAME	ジョブ名
SLURM_SUBMIT_DIR	sbatch コマンド実行時のカレントディレ クトリ
SLURM_ARRAY_JOB_ID	アレイジョブ ID(アレイジョブのみ設定)
SLURM_ARRAY_TASK_ID	アレイジョブのタスク番号(アレイジョブ のみ設定)

表 5-21 ジョブ内で参照可能な環境変数

また、ジョブ投入時にオプションで指定する資源量の一部は、ジョブ内で以下の環境変数 に設定されます。ジョブ投入時に指定されなかった項目は、自動で適切な値が設定されま す。

環境変数	説明					
SBATCH_PARTITION	パーティション名 (-p の値)					
SLURM_JOB_NUM_NODES	ノード数 (-N の値)					
SLURM_JOB_CPUS_PER_NODE	1ノードあたりのコア数 カンマ区切りでノード毎のコア数					
SLURM_MEM_PER_NODE	1ノードあたりのメモリ量 mem の値					
SLURM_MEM_PER_CPU	1コアあたりのメモリ量 mem-per-cpu の値					

表 5-22 ジョブ内で参照可能な環境変数

ログインノード上でジョブ投入前に設定した環境変数は、指示がない限り、ジョブ内には 引き継がれません。ジョブ内に引き継ぐためには、*sbatch* コマンドの--export オプションを 指定します。

Copyright (C) RIKEN, Japan. All rights reserved.

6. Singularity コンテナ

6.1 Singularity イメージの取得

Singularity イメージは利用者自身で用意する必要があります。 dockerhub や Singularity ライブラリなどから取得できます。 Dockerhub:https://hub.docker.com/ Singularity ライブラリ:<u>https://cloud.sylabs.io/library</u> 以下は dokcerhub から ubuntu の sif イメージを取得します。 [username@hokusai1 ~]\$ singularity pull docker://ubuntu:24.04

[username@hokusai1 ~]\$ Is -I ubuntu_24.04.sif

6.2 Singularity の"Command Line Interface(CLI)"について

Command Line Interface には Run, Exec, Shell が指定できます。各 CLI については 以下を参照ください。

Run

コンテナ内でユーザ定義のデフォルトコマンドを実行することができます。sif ファイルの 作成時に指定できる %runscript が実行されます。Docker や Podman でいう CMD に近 いものです。%runscript が指定されていない場合は Shell が立ち上がります。

```
[username@bwmpcxxx ~]$ singularity run ubuntu_24.04.sif << EOF
hostname
date
EOF
bwgpu4
Thu Jun 27 11:59:12 JST 2024
[username@bwmpcxxx ~]$
```

Exec

コンテナ内でコマンドを実行することができます。Docker や Podman のようにコンテナを 作成してから exec するのではなく、sif ファイルに直接コマンドを投げることが可能です。

Copyright (C) RIKEN, Japan. All rights reserved.

[username@bwmpcxxx ~]\$ singularity exec lolcow.sif cowsay "Fresh from the library!"



Shell

コンテナ内で Shell を実行する. %runscript などは実行されず Shell が起動する。デフォルトの Shell は Bash が起動するが, --shell で Shell を指定することも可能。

[username@bwmpcxxx ~]\$ singularity shell lolcow.sif Singularity>

6.3 Ubuntu の sif イメージを利用したインタラクティブジョブ実行

[username@hokusai1 ~]\$ srun --partition=mpc --account=RBxxxxxx --nodes=1 -ntasks=1 --pty bash
[username@bwmpcxxx ~]\$ singularity run ubuntu_24.04.sif
Singularity> cat /etc/os-release
PRETTY_NAME="Ubuntu 24.04 LTS"
NAME="Ubuntu"
(snip)
Singularity> exit

6.4 Ubuntu の sif イメージを利用したバッチジョブ実行

バッチジョブで Ubuntu sif イメージを MPC 計算ノードで実行。

Copyright (C) RIKEN, Japan. All rights reserved.

```
[username@hokusai1 ~]$ cat go.sh
#!/bin/sh
#SBATCH --partition=mpc
#SBATCH ---nodes=1
#SBATCH --ntasks=1
#SBATCH -J JOB
#SBATCH -o %x_%j.log
#SBATCH -e %x_%j.err
#SBATCH --time=1:00:00
#SBATCH --account=RBxxxxxx
singularity run ubuntu_24.04.sif << EOS
hostname
echo
cat /etc/os-release
EOS
[username@hokusai1 ~]$ sbatch go.sh
```

6.5 CPU 版(MPC) Singularity Ubuntu 環境構築

Singularity イメージは利用者の環境でも構築できますが、ログインノードでも作成できます。

6.5.1 Ubuntu 構築設定ファイル作成

```
[username@hokusai1 ~]$ cat ubuntu.def
BootStrap: docker
From: ubuntu:24.04
%post
   apt-get update
   apt-get -y upgrade
   apt-get -y install build-essential vim
```

6.5.2 Ubuntu 環境イメージ作成

```
[username@hokusai1~]$ singularity build --fakeroot ubuntu.img ubuntu.def
INF0:
         Starting build...
Getting image source signatures
Copying blob d4c3c94e5e10 done
Copying config 2abc4dfd83 done
Writing manifest to image destination
Storing signatures
2024/05/17 11:22:20 info unpack layer:
sha256:d4c3c94e5e10ed15503bda7e145a3652ee935c0b2e9de9b5c98df7ec0a0cd925
INF0:
         Running post scriptlet
+ apt-get update
Get:1 http://security.ubuntu.com/ubuntu focal-security InRelease [114 kB]
(snip)
Processing triggers for libc-bin (2.31-Oubuntu9.15) ...
INFO:
        Creating SIF file...
INF0:
         Build complete: ubuntu.img
[username@hokusai1~]$ ls -l ubuntu.img
-rwxr-x--- 1 username username 160432128 May 17 11:26 ubuntu.img
```

6.6 GPU ノード向け Singularity コンテナの取得

GPU ノードにインタラクティブジョブでログインし、最新の PyTorch コンテナ SIF ファイルを 取得します。 NGC や dokcerhub、Singularity のライブラリなどから SIF ファイルを取得できます。 NGC カタログ:

https://catalog.ngc.nvidia.com/?filters=&orderBy=weightPopularDESC&query=&pa ge=&pageSize=

GPU ノードのコンテナ向けバッチジョブでは、「--gpus」で GPU の使用枚数を指定します。 「--gpus=1」は 1GPU、「--gpus=2」は 2GPU、「--gpus=4」は 4GPU を使用できます。

Copyright (C) RIKEN, Japan. All rights reserved.

「--gpus=1」を指定した場合、CPU は 28 コア確保されます。

NGC から取得する場合は docker://nvcr.io/{イメージ名}を指定 [username@hokusai1 ~]\$ srun -p gpu_i --account=RBxxxxxx --nodes=1 --gpus=1 -time=5:00:00 --pty bash [username@bwgpu4 ~]\$ singularity pull pytorch:24.01-py3.sif docker://nvcr.io/nvidia/pytorch:24.01-py3

6.7 インタラクティブジョブ Singularity コンテナ内で Shell 起動

singularityコマンド実行時「--nv」オプションを指定するとGPU環境が設定されます。

[userna	ame@hok	kusai1	~]\$ srun -p	gpu_iaccou	unt=RBxxxxxxnodes	s=1gpu	us=1time=5:00:00
pty l	bash						
[userna	ame@bwg	gpu4 ~]	\$ singular	ity runnv p	oytorch:24.01-py3.si	if	
INF0:	Conv	verting	;SIF file 1	to temporary a	sandbox		
(snip)							
NOTE: (CUDA Fo	orward	Compatibil	ity mode ENABL	_ED.		
Usin	g CUDA	12.3 d	river vers	ion 545.23.08	with kernel driver	version	535. 129. 03.
See I	nttps:/	//docs.	nvidia.com/	/deploy/cuda-c	compatibility/ for o	details.	
Singula Thu Apr + NVID	arity> r 18 16 IA-SMI	nvidia 5:26:02 535.12	-smi 2024 	Driver Ve	ersion: 535.129.03	CUDA Vo	+ ersion: 12.2
 GPU	Name		P(ersistence-M	⊦ Bus−Id Disp.A	+ Volat	+ ile Uncorr. ECC
Fan	Temp	Perf	Pv	wr:Usage/Cap	Memory-Usage	GPU–U⁺	til Compute M.
						1	MIG M.
=====			=======================================		+======================================	=+=======	
0	NVIDIA	A H100	80GB HBM3	0n	00000000:10:00.0)ff	0
N/A	28C	P0		72W / 700W	4MiB / 81559M	MiB∣ (0% Default
						I	Disabled
+					+	+	+
(snip)							
Singula	arity>						

6.8 インタラクティブジョブ PyTorch コンテナで CuPy を使用した実行例

[username@hokusai1~]\$ srun -p gpu_i --account=RBxxxxxx --nodes=1 --gpus=1 -time=5:00:00 --pty bash \$ singularity pull pytorch:24.01-py3.sif docker://nvcr.io/nvidia/pytorch:24.01py3 コンテナ内で Shell を起動する。 [username@bwgpu4~]\$ singularity run --nv pytorch:24.01-py3.sif (snip) Singularity> pip freeze|grep cupy cupy==13.1.0 Singularity> pip freeze|grep Cython Cython==0. 29. 37 Singularity> cat cupy_ex.py import time import cupy as cp start_time = time.time() a = cp. random. rand (50000, 50000) b = cp. random. rand (50000, 50000) result = cp. dot(a, b)end_time = time.time() print(f"CuPy Time: {end_time - start_time} seconds") Singularity> python cupy_ex.py CuPy Time: 331.78111362457275 seconds

6.9 バッチジョブ PyTorch コンテナの GPU 実行例

```
[username@hokusai1 ~]$ cat go.sh
#!/bin/sh
#SBATCH --partition=gpu
#SBATCH ---nodes=1
#SBATCH --gpus=1
#SBATCH --ntasks-per-gpu=1
#SBATCH -J GPU
#SBATCH -o %x_%j.log
#SBATCH -e %x_%j.err
#SBATCH --time=1:00:00
#SBATCH --account=RBxxxxxx
singularity run --nv pytorch:24.01-py3.sif << EOS
python cupy_ex.py
echo
cat /etc/os-release
EOS
[username@hokusai1 ~]$ sbatch go.sh
```