

国立大学法人 大阪大学 サイバーメディアセンター 御中

スーパーコンピュータ バッチシステム入門/応用

2022年6月22日 日本電気株式会社 第一官公ソリューション事業部

Orchestrating a brighter world

NECは、安全・安心・公平・効率という社会価値を創造し、 誰もが人間性を十分に発揮できる持続可能な社会の実現を目指します。

バッチシステム入門/応用

- 1. バッチシステムの概要
- 2. ジョブの投入方法
- 3. 各種計算資源の利用方法
- 4. ジョブスクリプトのテクニック
- 5. 高度なジョブ投入・実行方法
- 6. SX-Aurora TSUBASA 公開情報



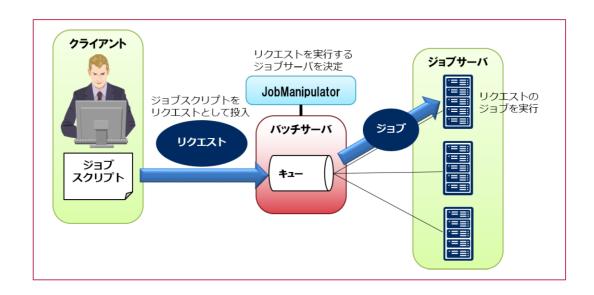
1. バッチシステムの概要

1. バッチシステムの概要

- ▶ スーパーコンピュータは、多くの利用者で共同利用するため、利用者の計算資 源要求(ジョブ)を受け付け、計算資源が空いたタイミングでジョブを順次実行す るバッチシステムを採用しています。
- ▶ バッチシステム内で、ジョブの実行を制御・管理するものが、ジョブ管理シス テムです。SQUIDでは「NEC NQSV」、OCTOPUSでは「NEC NQSII」を採用 しています。
- ▶ 「NEC NQSV」と「NEC NQSII」の基本的な利用方法は同じであるため、本講 習会では「NEC NQSV」を基に説明します。
- ※「NEC NQSV」、「NEC NQSII」における"リクエスト"と"ジョブ" 「NEC NQSV」、「NEC NQSII」の製品マニュアルでは以下の用語が使われます。
 - ・リクエスト:ユーザがフロントエンドから要求するバッチ利用の単位
 - ·**ジョブ**:個々の計算ノード上で実行される、リクエスト内の実行単位 本講習会では、他システムも含めた理解のしやすさを考慮し、"ジョブ"という用語を、 フロントエンドから要求するバッチ利用の単位(NQSV、NQSIIのリクエストに相当)の意味で使用します。

1. バッチシステムの概要

- ▶ 利用方法には、「**会話ジョブ**」と「バッチジョブ」の2種類があります。
- フロントエンドからジョブ管理システムにジョブ実行を要求(ジョブを投入) します。
- ▶ 要求されたジョブ要求は、他のジョブ要求の優先度や要求資源量、ジョブサー バの利用状況などがジョブ管理システムによって加味・判断され、実行順序が 決定されます。





Orchestrating a brighter world

2. ジョブの投入方法

- 2.1. 会話ジョブ投入方法
- 2.2. バッチジョブ投入方法
- 2.3. ジョブ管理システムのコマンドの使い方

2.1. 会話ジョブ投入方法

- ◆ 会話ジョブは、会話的(インタラクティブ)に計算ノードを利用するジョブで す。会話ジョブを使用する場合は**qloginコマンド**を使用します。
- ◆ 会話ジョブ投入コマンドのサンプル(汎用CPUノードで実行する例)を以下に記載します。

\$ qlogin -q INTC --group=G01234 [オプション]

-q キュー名…会話キューを指定

--group=グループ名…課金対象にするグループ名を指定

◆ qloginコマンドを実行すると、リクエストIDが採番され、以下の例のように標準出力に表示されます。

(出力例)

Request **123.sqd** submitted to queue: INTC.

Waiting for **123.sqd** to start.

2.2. バッチジョブ投入方法(1)

- ◆ ジョブスクリプトファイルの作成
 - ジョブスクリプトファイルは、バッチジョブ内での実行内容、キューの指定、計算資源の要求量やその他オプションを記述するファイルです。
 - このジョブスクリプト内に、**計算を行うための実行コマンド**を書くとともに、ジョブ 管理システムに対して**要求する資源量などをオプション(#PBS)として記述**していきま す。なお、オプションを指定しない場合は**既定値が適用**されます。
 - ・ジョブスクリプトのサンプル(汎用CPUノードでシリアル実行する例)を以下に記載します。

```
#!/bin/bash
#----- asub option -----
                                                          要求する資源量な
#PBS -q SQUID
                         → バッチジョブを投入するキュー名の指定
                                                          どを記載
                         → 課金対象にするグループ名
#PBS --group=G01234
#PBS -I elapstim reg=01:00:00
                         → ジョブの最大実行時間の要求値 1時間の例
#PBS -I cpunum job=76
                         → 使用するCPUコア数の要求値
#PBS -m b
                         → バッチジョブ実行開始時にメールを送信
#PBS -M user@hpc.cmc.osaka-u.ac.jp → 送信先アドレス
#----- Program execution -----
                                                          計算を行うための
module load BaseCPU/2021
                         → ベース環境をロードします
                                                          実行コマンドを記
                                                         載
cd $PBS O WORKDIR
                         → gsub実行時のカレントディレクトリへ移動
                         →プログラムの実行
./a.out
```

2.2. バッチジョブ投入方法(2)

- ◆ ジョブ管理システムへのバッチジョブ要求
 - qsub バッチジョブの投入

\$ qsub [オプション] [ジョブスクリプトファイル名]

※[オプション]については、頁11-14にて記載

qsubコマンドを実行すると、リクエストIDが採番され、以下の例のように標準出力に表示されます。

(出力例)

Request 1182. sqd submitted to queue: SQUID.

計算ノードのファイルシステム

計算ノードでは以下ファイルシステムを利用することが可能です。フロントエンドからもアクセス可能です。

- ・home領域(/sqfs/home/(利用者番号)) ※ジョブ実行開始時のカレントディレクトリ
- ・拡張領域(/sqfs/work/(グループ番号)/(利用者番号)) ※5TB(購入により拡張可能)
- ・高速領域(/sqfs/ssd/(グループ番号)/(利用者番号)) ※要購入

home領域の容量は10GBであるため、拡張領域または高速領域の利用を推奨します。

・環境変数PBS_O_WORKDIRの利用

ジョブスクリプトでは、環境変数PBS_O_WORKDIRにqsub実行時のカレントディレクトリのパスが自動的に 代入されます。バッチスクリプト内で、以下のコマンドを実行すると、qsub実行時のカレントディレクトリへ 簡単に移動できますので、**実行ディレクトリでジョブを投入することを推奨します**。

\$ cd \$PBS_O_WORKDIR

2.2. バッチジョブ投入方法(3)

◆ ノード共通オプション(必須)

オプション	機能
-q [バッチキュー名]	バッチジョブを投入するキューを指定します。 キューの詳細は「3. 各主計算資源環境の使い方」で説明します。
group=[グループ名]	指定グループでジョブを実行します。 指定グループのポイントが消費されます。



2.2. バッチジョブ投入方法(4)

▶ 汎用CPUノード向けオプション

オプション	機能
-b [ノード数]	ジョブを実行するノード数を指定します。
-T [使用MPIライブラリ名]	MPI実行を行う場合に指定が必要です。 Intelコンパイラを利用している場合、IntelMPIが利用可能です。 -T intmpi と指定してください。 NVIDIA HPC SDKコンパイラを利用している場合、OpenMPIが 利用可能です。-T openmpiを指定してください。

※MPIライブラリとモジュール、コンパイラの関係 本システムではコンパイラ、ライブラリ、アプリケーションの環境変数設定を「Environment modules | で管理しています。MPIライブラリとモジュール、コンパイラの関係は以下のとおりです。

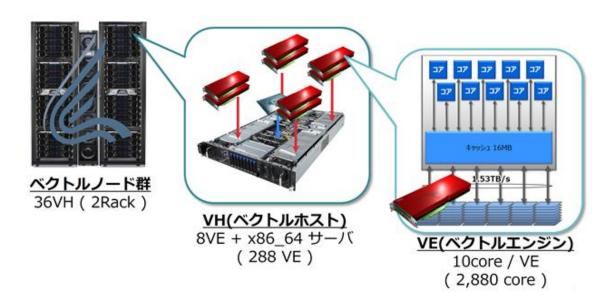
計算環境	MPIライブラリ	コンパイラ	推奨環境 モジュール名
汎用CPUノード	Intel MPI	Intel Parallel Studio	BaseCPU
ベクトルノード	NEC MPI	NEC SDK for VE	BaseVEC
GPUノード	Open MPI	NVIDIA HPC SDK CUDA	BaseGPU
(なし)	Open MPI	GNU Compiler	BaseGCC

2.2. バッチジョブ投入方法(5)

◆ ベクトルノード向けオプション

オプション	機能
venode=[総VE数]	利用する総VE数を指定します。 ベクトルノードを使用する場合、必須オプションです。
venum-lhost=[論理ホストあたりのVE数]	論理ホストを構成するVE数を指定します。 (大まかには、[1VH内で確保するVE数]とご理解ください)
-T necmpi	SX-Aurora TSUBASAにてMPI実行を行う場合に指定が必要です。

【参考】ベクトル計算ノード $x86_64$ サーバである VH (ベクトルホスト) に 8VE を搭載し、全 36VH で構成されます。



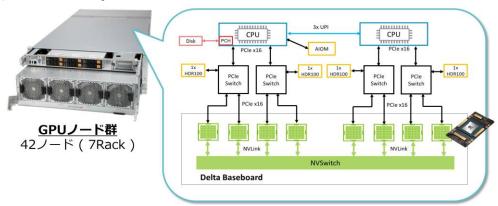
2.2. バッチジョブ投入方法(6)

◆ GPUノード向けオプション

オプション	機能
-l gpunum_job=[ノードあたりのGPU数]	ノードあたりの利用するGPU数を指定します。 GPUノードを使用する場合、必須オプションです。
gpunum-lhost=[ノードあたりのGPU数]	上記「-l gpunum_job」と同様です。
-b [ノード数]	ジョブを実行するノード数を指定します。
-T [使用MPIライブラリ名]	MPI実行を行う場合に指定が必要です。 Intelコンパイラを利用している場合、IntelMPIが利用可能です。 -T intmpi と指定してください。 NVIDIA HPC SDKコンパイラを利用している場合、OpenMPIが 利用可能です。-T openmpiを指定してください。

【参考】GPUノード

CPUとして「Intel Xeon Platinum 8368」、GPUとして「NVIDIA A100」を有し、4Uの筐体に8GPUを搭載しています。



GPUノード CPU 2.4GHz/38core x2 GPU 9.7TFlops x8 + NVSwitch



2.3. ジョブ管理システムのコマンドの使い方(1)

- ◆ バッチジョブの確認
 - qstat 投入したジョブの状況を確認

\$ qstat

ユーザ自身が投入したジョブの一覧を表示

\$ qstat -I

省略なく一覧を表示

上記のオプションはアルファベットのLの小文字です。

「--adjust-column」を同時に指定することで幅の最適化が行われます。

\$ qstat -f [リクエストID]

特定ジョブの詳細を表示

2.3. ジョブ管理システムのコマンドの使い方(2)

- ◆ バッチジョブの確認
 - qstatgroup ユーザが所属するグループが投入したジョブを確認

\$ qstatgroup

ユーザが所属するグループが投入したジョブの一覧を表示

■ sstat – 投入リクエストの実行開始時刻を確認

\$ sstat

投入リクエストの実行開始時刻を表示 実行開始時刻が決まっていない場合は、時刻は表示されない。



\Orchestrating a brighter world

2.3. ジョブ管理システムのコマンドの使い方(3)

- ◆ バッチジョブの標準出力/標準エラー出力表示
 - qcat ジョブ実行中に、標準出力/標準エラー出力を表示
 - \$ qcat -e [リクエストID]

標準エラー出力を表示する場合

\$ qcat -o [リクエストID]

標準出力を表示する場合

以下のオプションを組み合わせることも可能です。

- -f ファイルの内容が増え続けるとき、追加されたデータを出力します。
- -n 指定した行数分表示します。(無指定時は10行分です。)
- -b ファイルの先頭から表示します。 (無指定時は最終行から表示されます。)



2.3. ジョブ管理システムのコマンドの使い方(4)

- ◆ バッチジョブの保留・解除
 - qhold 投入したジョブを保留

\$ qhold [リクエストID]

ホールドすることでスケジューリング対象から外れ、ジョブを実行開始されない状態にする。

■ qrls - 保留にしたジョブを保留解除

\$ qrls [リクエストID]

ホールド状態が解除され、再度スケジューリング対象の状態にする。



2.3. ジョブ管理システムのコマンドの使い方(5)

- ◆ ジョブ情報の表示
 - acstat ユーザ自身が過去に投入したジョブ情報を表示

\$ acstat コマンド実行時から24時間以内のジョブの情報を表示

\$ acstat -A コマンド実行年度のジョブの情報を表示

■ acstatgroup - ユーザが所属するグループが過去に投入したジョブ情報を表示

\$ acstatgroup コマンド実行時から24時間以内のジョブの情報を表示

\$ acstatgroup - A コマンド実行年度のジョブの情報を表示

- ◆ バッチジョブの削除
 - qdel 投入したリクエストを削除

\$ qdel [リクエストID]

バッチジョブが実行状態(RUN)の場合は、最初にSIGTERMを送り、その後SIGKILLが送られます。



2.3. ジョブ管理システムのコマンドの使い方(6)

- ▶ キューの状態確認
 - qstat -Q バッチジョブが実行されるキューの状態を表示

\$ qstat -Q

[EXECUTION QUEUE] Batch Server Host: sqd

SCH JSVs ENA STS PRI TOT ARR WAI QUE PRR RUN POR EXT HLD HOL RST SUS MIG STG CHK QueueName

1 1518 ENA ACT 10 169 0 0 5 0 164 0 0 0 0 0 0 SC1

主な項目の説明は以下のとおりです。

ENA: ENA → ジョブ投入可能、DIS → ジョブ投入不可

STS: ACT → 実行可能、INA → 実行不可

QUE:実行待ちジョブ件数 RUN: 実行中ジョブ件数

\Orchestrating a brighter world

3. 各種計算資源環境の使い方

3.1. SQUIDシステム

3.2. OCTOPUSシステム

3.1.1. 【SQUID】ジョブクラス

- ◆ それぞれのジョブクラスは、ジョブ管理システム上のキューに対応しており、利用者は キューにジョブを投入することで計算環境の利用が可能です。
- ◆ キューは、利用者がジョブを直接投入する投入キューと、実行順を待ち合わせる実行 キューに分かれます。
- ・汎用CPUノード向けジョブクラス

利用方法	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	SQUID	24時間	38,912core (76c/ノード)	124TiB (248GB/ノード)	512	
	SQUID-R	24時間	38,912core (76c/ノード)	124TiB (248GB/ノード)	512	※ 1
	SQUID-H	24時間	38,912core (76c/ノード)	124TiB (248GB/ノード)	512	※ 2
	SQUID-S	24時間	38core	124GiB	1	※ 3
	DBG	10分	152core (76c/ノード)	496GiB (248GB/ノード)	2	デバッグ用
	INTC	10分	152core (76c/ノード)	496GiB (248GB/ノード)	2	会話ジョブ利用
占有利用	mySQUID	無制限	76core × 占有ノード数	248GiB × 占有ノード数	占有数	

- ※1 NW帯域が狭い経路の利用を許容して、実行待ち時間を短縮されたい方向け
- ※2 ポイント消費を多くして高優先度ジョブを投入し、実行待ち時間を短縮されたい方向け
- ※3 他のジョブとのノード内の共有を許容して、ポイント消費を抑えたい方向け
- ※利用可能経過時間を指定しない場合、利用経過時間の規定値はDBGが10分、その他は1時間となります
- ※Core数を指定しない場合、利用Core数の規定値はSQUID-Sが38core、その他は76coreとなります
- ※メモリを指定しない場合、利用メモリの規定値はSQUID-Sが124GB、その他は248GBとなります



3.1.1. 【SQUID】ジョブクラス

ベクトルノード向けジョブクラス

利用方法	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 VE数	備考
共有利用	SQUID	24時間	2,560core (10c/VE)	12TiB (48GB/VE)	256	
	SQUID-H	24時間	2,560core (10c/VE)	12TiB (48GB/VE)	256	※ 1
	SQUID-S	24時間	40core (10c/VE)	192GiB (48GB/VE)	4	※ 2
	DBG	10分	40core (10c/VE)	192GiB (48GB/VE)	4	デバッグ用
	INTV	10分	40core (10c/VE)	192GiB (48GB/VE)	4	会話ジョブ利用
占有利用	mySQUID	無制限	10core × 占有VE数	48GiB × 占有VE数	占有数	

※1 ポイント消費を多くして高優先度ジョブを投入し、実行待ち時間を短縮されたい方向け

- ※2他のジョブとのノード内の共有を許容して、ポイント消費を抑えたい方向け
- ※利用可能経過時間を指定しない場合、利用経過時間の規定値はDBGが10分、その他は1時間となります

3.1.1. 【SQUID】ジョブクラス

・GPUノード向けジョブクラス

利用方法	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	SQUID	24時間	2,432core (76c/ノード)	15.75TiB (504GB/ノード)	32	
	SQUID-H	24時間	2,432core (76c/ノード)	15.75TiB (504GB/ノード)	32	※ 1
	SQUID-S	24時間	38core	252GiB	1	※ 2
	DBG	10分	152core (76c/ノード)	1,008GiB (504GB/VE)	2	デバッグ用
	INTG	10分	152core (76c/ノード)	1,008GiB (504GB/VE)	2	会話ジョブ利用
占有利用	mySQUID	無制限	76core × 占有ノード数	504GiB × 占有ノード数	占有数	

※1 ポイント消費を多くして高優先度ジョブを投入し、実行待ち時間を短縮されたい方向け

- ※2 他のジョブとのノード内の共有を許容して、ポイント消費を抑えたい方向け
- ※利用経過時間を指定しない場合、利用経過時間の規定値はDBGが10分、その他は1時間となります
- ※Core数を指定しない場合、利用Core数の規定値はSQUID-Sが8core、その他は76coreとなります
- ※メモリを指定しない場合、利用メモリの規定値はSQUID-Sが252GB、その他は504GBとなります



3.1.2. 【SQUID】汎用CPUノードの使い方(1)

◆ シリアル実行利用方法 1ノード内でのプログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -q SQUID
#PBS --group=G01234
#PBS - l elapstim_req=00:30:00
#---- Program execution -----
module load BaseCPU/2021
                          → コンパイル時にmodule loadしていたものを記載
module load xxx/xxx
cd $PBS O WORKDIR
. /a. out
```

3.1.2. 【SQUID】汎用CPUノードの使い方(2)

◆ スレッド並列利用方法

1ノード内でのスレッド並列プログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -q SQUID
#PBS --group=G01234
#PBS -I elapstim_req=00:30:00
#PBS -v OMP NUM THREADS=76
                       → 並列実行数を指定
#---- Program execution -----
module load BaseCPU/2021
                   → コンパイル時にmodule loadしていたものを記載
module load xxx/xxx
cd $PBS O WORKDIR
. /a. out
```

利用時の注意点

実行スクリプトにて「OMP NUM THREADS」を忘れずに指定してください。 「OMP NUM THREADS」を指定しない場合、あるいは間違った値を指定してしまった場合、 意図しない並列数での実行となる可能性があります。



Orchestrating a brighter world

3.1.2. 【SQUID】汎用CPUノードの使い方(3)

◆ MPI利用方法(IntelMPI)

4ノードで合計304プロセスを生成するIntelMPIのプログラム実行を想定したジョブスクリプトのサンプルです。

3.1.2. 【SQUID】汎用CPUノードの使い方(4)

◆ MPI+ノード内並列利用方法 4ノードで各ノード上に1プロセスを生成し、各プロセスで76個のスレッドを生成する IntelMPIプログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- asub option -----
#PBS -q SQUID
#PBS --group=G01234
#PBS - I elapstim_req=00:30:00
#PBS -b 4
                               → ノード数を指定
                              → IntelMPIを指定
#PBS -T intmpi
#PBS -v OMP NUM THREADS=76
                              → スレッド並列実行数を指定
#---- Program execution -----
module load BaseCPU/2021
module load xxx/xxx
                              → コンパイル時にmodule loadしていたものを記載
cd $PBS 0 WORKDIR
mpirun ${NQSV_MPIOPTS} -np 4 ./a.out → mpirunの引数に${NQSV_MPIOPTS}を指定
                                  1プロセスを1CPUに割り当て、4コアに76スレッドを割り当てる
```

3.1.3. 【SQUID】ベクトルノードの使い方(1)

◆ シリアル実行利用方法 1VE内でのプログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#------ qsub option ------
#PBS -q SQUID-S → 1~7VE使用する場合は、SQUID-Sキューを指定
#PBS --group=G01234
#PBS -I elapstim_req=00:30:00
#PBS --venode=1 → 全体で1VE使用することを指定
#----- Program execution ------
module load BaseVEC/2021
module load xxx/xxx → コンパイル時にmodule loadしていたものを記載

cd $PBS_0_WORKDIR
./a. out
```

3.1.3. 【SQUID】ベクトルノードの使い方(2)

◆ スレッド並列利用方法

1VE内でのスレッド並列プログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option --
                           → 1~7VE使用する場合は、SQUID-Sキューを指定
#PBS -a SQUID-S
#PBS --group=G01234
#PBS - I elapstim_req=00:30:00
                      → 並列実行数を指定
#PBS -v OMP NUM THREADS=10
#PBS --venode=1
                           → 全体で1VE使用することを指定
#---- Program execution -----
module load BaseVEC/2021
module load xxx/xxx
                 → コンパイル時にmodule loadしていたものを記載
cd $PBS_O_WORKDIR
. /a. out
```

利用時の注意点

実行スクリプトにて「OMP NUM THREADS」を忘れずに指定してください。 「OMP NUM THREADS」を指定しない場合、あるいは間違った値を指定してしまった場合、 意図しない並列数での実行となる可能性があります。



3.1.3. 【SQUID】ベクトルノードの使い方(3)

◆ MPI利用方法(NEC MPI)

40VEで合計400プロセスを生成するNEC MPIのプログラム実行を想定したジョブスクリプトのサンプルです。

3.1.3. 【SQUID】ベクトルノードの使い方(4)

◆ MPI+ノード内並列利用方法

40VEでVEあたり1プロセスを生成し、各プロセスで10個のスレッドを生成するNEC MPI プログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -a SQUID
#PBS --group=G01234
\#PBS - I = lapstim req = 00:30:00
#PBS --venode=40
                              → 全体で40VE使用することを指定
#PBS -T necmpi
                              → NEC MPIを指定
#PBS -v OMP_NUM_THREADS=10
                              → スレッド並列実行数を指定
#---- Program execution -----
module load BaseGPU/2021
module load xxx/xxx
                            → コンパイル時にmodule loadしていたものを記載
cd $PBS O WORKDIR
mpirun -venode -np 40 ./a.out
                          → MPIプロセスのVHおよびVEへの割り当てはNQSVが自動で実施
```

利用時の注意点

「OMP NUM THREADS」で指定するVEあたりのスレッド数は、コア数である10を超えないよ うにしてください。



3.1.4. 【SQUID】GPUノードの使い方(1)

◆ シリアル実行利用方法 1ノード内でのプログラム実行を想定したジョブスクリプトのサンプルです。



3.1.4. 【SQUID】GPUノードの使い方(2)

◆ スレッド並列利用方法

1ノード内でのスレッド並列プログラム実行を想定したジョブスクリプトのサンプルです。

利用時の注意点

実行スクリプトにて「OMP_NUM_THREADS」を忘れずに指定してください。 「OMP_NUM_THREADS」を指定しない場合、あるいは間違った値を指定してしまった場合、 意図しない並列数での実行となる可能性があります。



3.1.4. 【SQUID】GPUノードの使い方(3)

◆ MPI利用方法(OpenMPI)

2ノードで合計4プロセスを生成するOpenMPIのプログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -a SQUID
#PBS --group=G01234
#PBS - I elapstim_req=00:30:00
#PBS -b 2
                                → ノード数を指定
#PBS - I gpunum_job=8
                                → ノードあたり8GPU使用することを指定
#PBS -T openmpi
                                → OpenMPIを指定
#PBS -v NQSV MPI MODULE=BaseGPU/2021 → OpenMPIの起動に必要となるモジュール名を指定
#---- Program execution -----
module load BaseGPU/2021
module load xxx/xxx
                                → コンパイル時にmodule loadしていたものを記載
cd $PBS O WORKDIR
mpirun ${NQSV_MPIOPTS} -np 4 -npernode 2 ${PBS_0_WORKDIR}/mpi_prog
                                → mpirunの引数に${NQSV MPIOPTS}を指定
```

利用時の注意点

NQSV_MPI_MODULEに複数のモジュールを指定する場合は、コロン(:)区切りで指定することが可能です。 (例) BaseGPU/2021とcuda/11.3を指定する #PBS -v NQSV_MPI_MODULE=BaseGPU/2021:cuda/11.3



3.1.4. 【SQUID】GPUノードの使い方(4)

◆ MPI+ノード内並列利用方法(OpenMPI)

2ノードで各ノード上に2プロセスを生成し、OpenMPにより各プロセスで12個のスレッド を生成するOpenMPIプログラム実行を想定したジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -a SQUID
#PBS --group=G01234
#PBS - I elapstim_req=00:30:00
#PBS -b 2
                               → ノード数を指定
#PBS -I gpunum_job=8
                               → ノードあたり8GPU使用することを指定
#PBS -T openmpi
                               → OpenMPIを指定
#PBS -v OMP_NUM_THREADS=12
                               → スレッド並列実行数を指定
#PBS -v NQSV MPI MODULE=BaseGPU/2021 → OpenMPIの起動に必要となるモジュール名を指定
#---- Program execution -
module load BaseGPU/2021
module load xxx/xxx
                               → コンパイル時にmodule loadしていたものを記載
cd $PBS O WORKDIR
mpirun ${NQSV MPIOPTS} -np 4 -npernode 2 --bind-to socket ¥
                                 --report-bindings ${PBS_0_WORKDIR}/mpi_prog
                               → 1プロセスを1CPUに割り当て、12コアに12スレッドを割り当てる
                                  ため、「--bind-to-socket」等を指定
```

利用時の注意点

NQSV MPI MODULEに複数のモジュールを指定する場合は、コロン(:)区切りで指定することが可能です。

(例) BaseGPU/2021とcuda/11.3を指定する

#PBS -v NOSV MPI MODULE=BaseGPU/2021:cuda/11.3

3.2.1. **【OCTOPUS】**ジョブクラス

・汎用CPUノード向けジョブクラス

利用方法	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	OCTOPUS	120時間	3,072 core (24c/ノード)	23,680 GB (185GB/ノード)	128	
	DBG	10 分	24 core (24c/ノード)	185 GB (185GB/ノード)	1	デバッグ用
	INTC	120時間	3,072 core (24c/ノード)	23,680 GB (185GB/ノード)	128	インタラクティブ利用
占有利用	myOCTOPUS	無制限	24core× 占有ノード数	185GB× 占有ノード数	占有ノード	

※利用可能経過時間を指定しない場合、利用経過時間の規定値はINTCが20分、その他は1分となります

※Core数を指定しない場合、利用Core数の規定値は24coreとなります

※メモリを指定しない場合、利用メモリの規定値は185GBとなります

※OCTOPUSのジョブスクリプトファイルのサンプルは以下を参照ください。

参考: OCTOPUSサンプルスクリプト集

http://www.hpc.cmc.osaka-u.ac.ip/system/manual/octopus-use/cpu-script collection/



3.2.1. **【OCTOPUS】**ジョブクラス

・GPUノード向けジョブクラス

利用方法	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	OCTOPUS	120時間	768 core (24c/ノード)	5,920 GB (185GB/ノード)	32	1ノードあたり GPU4枚まで利用可能
	DBG	10 分	24 core (24c/ノード)	185 GB (185GB/ノード)	1	デバッグ用
	INTG	120時間	768 core (24c/ノード)	5,920 GB (185GB/ノード)	32	会話ジョブ利用 1ノードあたり GPU4枚まで利用可能
占有利用	myOCTOPUS	無制限	24core× 占有ノード数	185GB× 占有ノード数	占有ノード	

※利用可能経過時間を指定しない場合、利用経過時間の規定値はINTGが20分、その他は1分となります

※Core数を指定しない場合、利用core数の規定値は24coreとなります

※メモリを指定しない場合、利用メモリの規定値は185GBとなります



Orchestrating a brighter world

© NEC Corporation 2021

3.2.1. **【OCTOPUS】**ジョブクラス

・Xeon Phiノード向けジョブクラス

利用方法/	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	ОСТРНІ	120時間	2,048 core (64c/ノード)	5,920 GB (185GB/ノード)	32	
	INTP	120時間	2,048 core (64c/ノード)	5,920 GB (185GB/ノード)	32	会話ジョブ利用
占有利用	myOCTPHI	無制限	24core× 占有ノード数	185GB× 占有ノード数	占有ノード	

※利用可能経過時間を指定しない場合、利用経過時間の規定値はINTPが20分、その他は1分となります

※Core数を指定しない場合、利用core数の規定値は64coreとなります

※メモリを指定しない場合、利用メモリの規定値は185GBとなります

・大容量主記憶搭載ノード向けジョブクラス

利用方法/	ジョブ クラス	利用可能 経過時間	利用可能 最大Core数	利用可能 メモリ	同時利用可能 ノード数	備考
共有利用	OCTMEM	120時間	256 core (128c/ノード)	12TB (6TB/ノード)	2	
	INTM	120時間	256 core (128c/ノード)	12TB (6TB/ノード)	2	会話ジョブ利用

※利用可能経過時間を指定しない場合、利用経過時間の規定値はINTMが20分、その他は1分となります

※Core数を指定しない場合、利用core数の規定値は128coreとなります

※メモリを指定しない場合、利用メモリの規定値は6TBとなります



3.2.2. 【OCTOPUS】SQUIDとの違い

◆ SQUIDシステムとOCTOPUSシステムでは主に以下の違いがあります。

システム名	SQUID	OCTOPUS	補足
MPIライブラリ	Intel MPI NEC MPI Open MPI	Intel MPI Open MPI NVAPICH2	使用可能なライブラリが異なります
環境変数	NQSV_MPIOPTS	NQSII_MPIOPTS	ジョブスクリプト内で指定する 環境変数名が異なります
GPUノード オプション	gpunum_job gpunum-lhost	gpunum_job	OCTOPUSではNQSVにて gpunum_lhostオプションを使用できません
moduleコマンド	対応	非対応	OCTOPUSではmoduleコマンドを 使用できません

参考: OCTOPUSの利用方法

(CPU) http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/cpu-use/

(GPU) http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/gpu-use/

(Phi) http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/phi-use/

(大容量メモリノード)

http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/mem-use/



4. ジョブスクリプトのテクニック

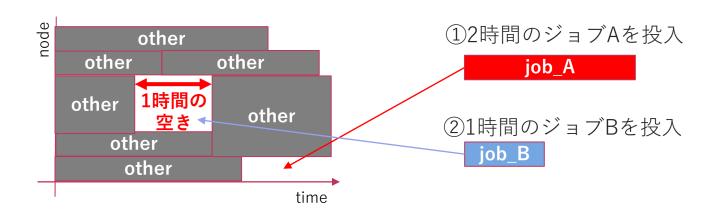
4.1. ジョブスクリプトの推奨オプション(1)

◆ 経過時間制限値の指定

ジョブ管理システムでは要求資源に合わせて最短で実行可能なノードにジョブをアサインします。

ジョブスクリプトにて、

経過時間制限値(# PBS -I elapstim_req=[経過時間制限値])に 適切な時間(プログラムの実行時間 + α)を指定することで、 アサインされやすくなり、実行開始時間が早くなる可能性があります。



4.1. ジョブスクリプトの推奨オプション(2)

◆ ノード共通おすすめオプション 必須ではありませんが、便利なオプションを紹介します。

オプション	機能
-N [ジョブ名]	ジョブの名前を指定します。 指定がなければ、バッチスクリプト名がジョブ名になります。
-o [標準出力ファイル名]	バッチジョブの標準出力の出力ファイル名を指定します。 指定がなければ、ジョブ投入時のディレクトリに「ジョブ名.oリクエストID」の ファイル名で出力されます。
-e [標準エラー出力ファイル名]	バッチジョブの標準エラー出力の出力ファイル名を指定します。 指定がなければ、リクエスト投入時のディレクトリに「ジョブ名.eリクエストID」 のファイル名で出力されます。
-j [o,e]	バッチジョブの標準出力と標準エラー出力をマージします。 o:マージした結果を標準出力に出力します。 e:マージした結果を標準エラー出力に出力します。 (-jとoもしくはeの間にはスペースが必要です)



4.1. ジョブスクリプトの推奨オプション(3)

◆ ノード共通おすすめオプション

オプション	機能
-M [メールアドレス]	メールの送信先を指定します。複数指定する場合は -M メールアドレス1,メールアドレス2 のように「,」で区切ってください。
-m [b,e,a]	バッチジョブの状態の変化についてのメールを送ります。 b:ジョブが開始したときにメールを送信 e:ジョブが終了したときにメールを送信 a:ジョブが異常終了したときにメールを送信 (-mとbやeの間にはスペースが必要です) 複数を指定することができます。 例) 開始時と終了時にメール通知 -m be
-v 環境変数	バッチジョブを実行するときに使用する環境変数を指定します。 複数ノードで実行する場合、全てのノードで指定した環境変数の値が設定されます。
-r { y n }	バッチジョブのリラン可否を指定します。 y:リラン可能 n:リラン不可



4.2. 環境変数の活用(1)

◆ ジョブスクリプトの作成では、環境変数を活用すると便利です。 以下は、コア数、ノード数を環境変数として設定するジョブスクリプトのサンプルです。

2ノードで76並列(38コア/ノード)のプログラム実行を想定しています。

```
#!/bin/bash
#---- gsub option -----
#PBS -a SQUID
#PBS --group=G01234
#PBS - I elapstim reg=00:30:00
#PBS -v NODE=2
#PBS -v CORE=38
#PBS -b ${NODE}
                               → ノード数を指定
                            → ノード当たりのコア数を指定
#PBS - I cpunum job=${CORE}
#PBS -T intmpi
#---- Program execution -----
module load BaseCPU/2021
module load xxx/xxx
cd $PBS O WORKDIR
                               → qsub実行時のカレントディレクトリへ移動
NP=`expr ${NODE} "*" ${CORE}` → 総並列数を計算
mpirun ${NQSV MPIOPTS} -np ${NP} -ppn ${CORE} ./a.out
```

実行するノード数を変更する場合、ジョブスクリプトの変更箇所を1箇所にすることができます。



4.2. 環境変数の活用(2)

◆ ジョブ管理システムの既定値として設定される環境変数 以下は、リクエスト実行時に既定値として設定される環境変数です。

環境変数名	説明	-vオプションによる 変更可否
PBS_ENVIRONMENT	バッチリクエストがバッチ処理であることを示す、"PBS_BATCH"が設定 されます。	×
PBS_JOBID	バッチリクエストのジョブIDが設定されます。	×
PBS_SUBREQNO	サブリクエスト番号が設定されます(パラメトリックリクエストの場合) パラメトリックリクエストは「5.1. パラメトリックリクエスト機能」で説明します。	0
PBS_JOBNAME	バッチリクエスト名が設定されます。	×
PBS_NODEFILE	リクエストを構成するジョブが実行されているホスト一覧が記載されたファイルへのパスが設定されます。	0
PBS_O_HOME	クライアントホスト上の環境変数 HOME が設定されます。	\circ
PBS_O_HOST	クライアントホスト名が設定されます。	0
PBS_O_LANG	クライアントホスト上の環境変数 LANG が設定されます。	\circ
PBS_O_LOGNAME	クライアントホスト上の環境変数 LOGNAME が設定されます。	0
PBS_O_MAIL	クライアントホスト上の環境変数 MAIL が設定されます。	\circ
PBS_O_PATH	クライアントホスト上の環境変数 PATH が設定されます。	0
PBS_O_SHELL	クライアントホスト上の環境変数 SHELL が設定されます。	0
PBS_O_TZ	クライアントホスト上の環境変数 SHELL が設定されます。	0
PBS_O_WORKDIR	クライアントホスト上の作業用ディレクトリを設定します。	\circ

※OCTOPUSでは、「PBS_SUBREQNO」、「PBS_NODEFILE」は使用できません。

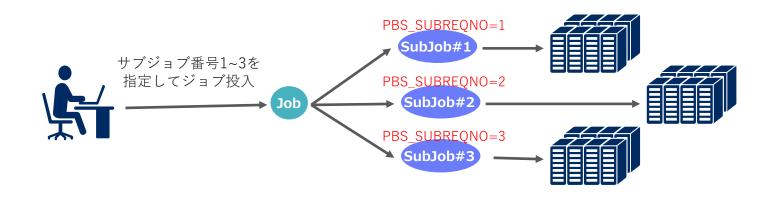


5. 高度なジョブ投入・実行方法

5.1. パラメトリックリクエスト機能(1)

◆ パラメトリックリクエストとは

入力となるパラメータや入力ファイルを変えながら、同一のジョ ブスクリプト をサブジョブとして複数回実行するジョブです。



利用シーン

同一ソルバーを用い、複数の入力パターンのあるアンサンブル実行



Orchestrating a brighter world

5.1. パラメトリックリクエスト機能(2)

◆ パラメトリックリクエストの投入方法 サブジョブ番号を指定したジョブスクリプトを作成し、qsubコマンドにて投入 します。

ジョブスクリプトのサンプル

```
#!/bin/bash
#PBS -q SQUID
#PBS --group=G01234
#PBS -I elapstim_req=01:00:00
#PBS -t 1, 2, 3
# ↑ パラメトリックリクエストのサブジョブ番号として1から3を指定(-t 1-3 と指定することもできます)
echo $PBS_SUBREQNO
# ↑ 環境変数PBS_SUBREQNOにサブジョブごとのサブジョブ番号が設定されています
```



5.1. パラメトリックリクエスト機能(3)

◆ パラメトリックリクエストの状態確認

投入したパラメトリックリクエストの状態は、qstatコマンドで確認することが できます。

\$ <mark>qstat</mark> RequestID	ReqName	UserName	Queue	Pri	STT	S	Memory	C	PU	Elapse	R	Н	M	Jobs
 1 <mark>234[]. sqd</mark> ↑ qstatコマン	-										Y	Y	Y	1
\$ <mark>qstat -s</mark> RequestID	ReqName	UserName	Queue	Pri	STT	S	Memory	C	PU	Elapse	R	Н	M	Jobs
1234[1]. sqd										0				
1234[2]. sqd	qsub. sh	user	SC1	0	QUE	_	0.00B	0.	00	0	Υ	Υ	Υ	1
1234[3]. sqd														
↑qstatコマン []内の数字に \$qstat -R RequestID		ブ番号です	-								ONE	<u> </u>		
1234[]. sqd ↑ qstatコマン TOTAL に、≦ DONE に、終	ドに-Rオフ 全サブジョ	プションを ブ数、ACT	指定した IVE に、	場合、 バッチ	サマ サー	リカ バゴ	ヾ表示され ニで保持し	ます ている			 ()		

50

5.1. パラメトリックリクエスト機能(4)

◆ パラメトリックリクエストの削除

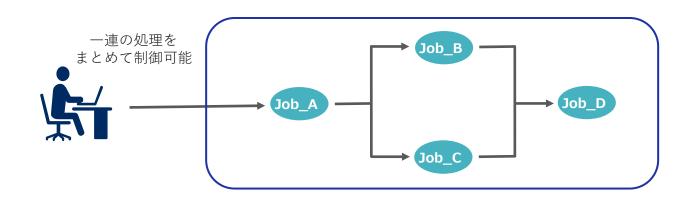
パラメトリックリクエストは全サブジョブを含む親ジョブに対する削除と、サブスクリプト単位での削除を行うことができます。

\$ qdel 1234[]. sqd
Request 1234[]. sqd was deleted.
 ↑全サブジョブを削除します
\$ qdel 1234[2]. sqd
Request 1234[2]. sqd was deleted.
 ↑指定されたサブジョブのみ削除します

5.2. ワークフロー機能(1)

◆ ワークフローとは

ジョブの投入・実行を含む一連の処理をシェルスクリプトとして記述し、ジョブ の実行順序等を制御する機能です。



利用シーン

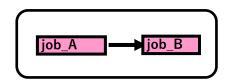
(例1) 最初に前処理用の1本のジョブを実行 その後、前処理で生成したファイルを利用する複数のジョブを実行

(例2) ジョブA、B、…、Zを順次実行 途中で正常終了しないジョブがあった場合、以降の処理をすべて自動削除

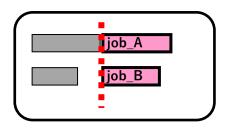


5.2. ワークフロー機能(2)

- ◆ 指定可能な依存関係、制御
 - 順次実行(--afterオプション) 「--after job_A ./job_B」のように指定することで、 job_Aの実行終了後、job_Bの実行が開始されます。

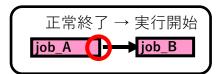


■ 同時実行(--parallelオプション) 「--parallel ./job_A ./job_B」のように指定することで、 job_Aとjob_Bが同時に実行開始されます。

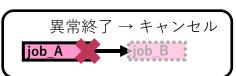


■ エラー発生時の後続ジョブのキャンセル(--cancel-afterオプション) 「--cancel-after ./job A」のように指定した上で、

「--after job_A ./job_B」のように指定することで、 job_Aが異常終了した場合、job_Bはキャンセル(qdel)されます。



ジョブAが正常終了した場合、 ジョブBの実行が開始される



ジョブAが異常終了した場合、 ジョブBはキャンセルされる



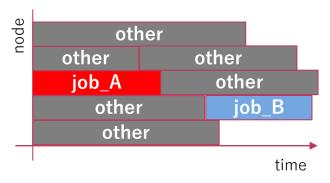
5.2. ワークフロー機能(3)

- ◆ ワークフローのジョブ実行制御を使うメリット
 - ワークフローを使わずに順次実行制御 (例:gwaitの利用)

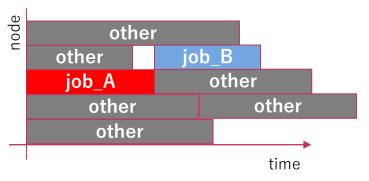
フロントエンドサーバ上でqwaitを利用し、1本目のジョブが終了したことを確認したら次のジョブをqsubする。

```
#!/bin/bash
qwait <reqID_A> ←1本目のジョブが終了するまでここで待つ
qsub./B_req ←2本目のジョブを投入
```

■ ワークフローを使うと隙間なくスケジューリングされる



qwaitを利用する場合 (job_Aの実行終了後、job_Bを投入)



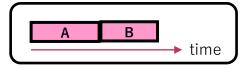
ワークフローを利用する場合 (job_Aの実行終了後、job_Bを投入)

⇒ job_Bの実行開始時刻が早い

5.2. ワークフロー機能(4)

- ◆ ワークフロースクリプトの作成
 - ワークフロースクリプトとは?
 - ワークフローの一連の処理を記述したシェルスクリプト
 - 後述するwstartコマンドを利用して、ワークフロースクリプトを投入します。
 - スクリプト例(Aの終了後、Bの実行開始)

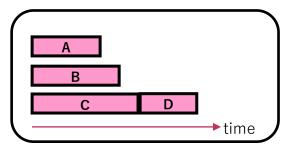
#!/bin/bash qsub -N A A.sh qsub --after A -N B B.sh



※-Nオプションによりジョブに 名前を付けることができる。

■ スクリプト例(A、B、Cのすべてが終了後、Dの実行開始)

#!/bin/bash qsub -N A A.sh qsub -N B B.sh qsub -N C C.sh qsub --after A,B,C -N D D.sh

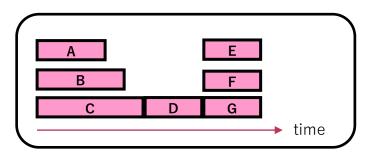




5.2. ワークフロー機能(5)

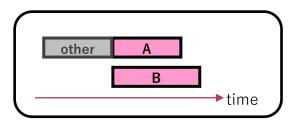
■ スクリプト例(A、B、Cのすべてが終了後、Dの実行開始、Dの終了後、E、F、Gの開始)

#!/bin/bash
qsub -N A A.sh
qsub -N B B.sh
qsub -N C C.sh
qsub --after A,B,C -N D D.sh
qsub --after D -N E E.sh
qsub --after D -N F F.sh
qsub --after D -N G G.sh



■ スクリプト例(A、Bを実行開始)

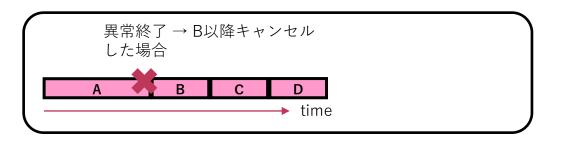
#!/bin/bash qsub --parallel A.csh B.sh



5.2. ワークフロー機能(6)

■ スクリプト例(Aが異常終了した場合、B以降をすべてキャンセル)

#!/bin/bash qsub --cancel-after -N A A.sh qsub --after A -N B B.sh qsub --after B -N C C.sh qsub --after B -N D D.sh



- ※以下の場合が異常終了とみなされます
 - ・終了ステータスが0以外
 - ・qdelで削除された
 - ・システム障害によって終了した

5.2. ワークフロー機能(7)

▶ ワークフローの投入 wstartコマンドの引数にワークフロースクリプトを指定します。

\$ wstart wfl.sh

※ワークフロースクリプトの名前をwfl.shとしています。

▶ ワークフローの状態確認 wstatコマンドを利用します。引数にワークフローIDを指定してください。

\$ wstat 123

ワークフローIDは投入時に表示されます。もしくは、引数なしでwstatコマンドを実行することで確認可能です。

<出力例>

\$ wstat 123 WFL-ID	RequestID	Reques	tName Stat Ex	it
123_sqd 123_sqd	12345.sqd 12346.sqd	А В	RUN QUE	 - -



5.2. ワークフロー機能(8)

◆ ワークフローの削除 wdelコマンドの引数にワークフローIDを指定します。

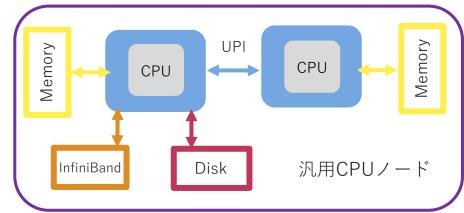
\$ wdel 123

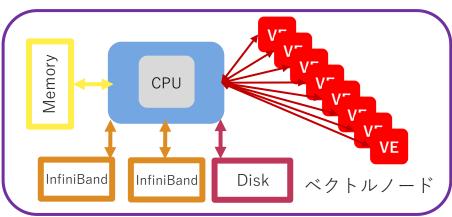
これにより、ワークフロー内で投入されたリクエストがすべて削除されます。

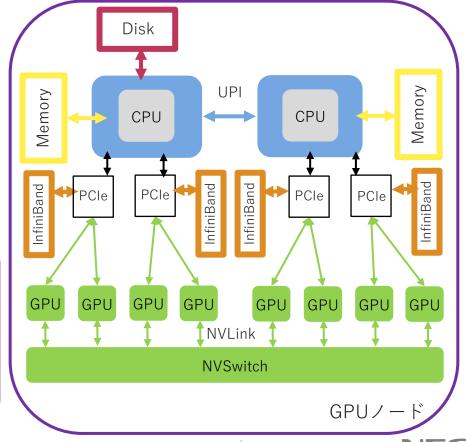


5.3. 高度なMPIの実行[アフィニティ] (1)

◆ プロセスのコアへの割り当て 各計算ノードは、以下のような構成をしています。計算ノードの構成に合わせて、MPIプログラムの各プロセスをコアへ割り当てることで、性能向上する場合があります。







5.3. 高度なMPIの実行[アフィニティ] (2)

◆ プロセス割り当ての確認 各MPIライブラリは、実行時に割り当てたCPU(VE)を出力することが可能です。 採取方法は以下のとおりです。

■ Intel MPIの場合

実行時に「I_MPI_DEBUG=4」を指定します。 標準出力に以下のように出力されます。

```
[0] MPI startup(): Rank
                      Pid
                            Node name Pin cpu
[0] MPI startup(): 0
                    446143
                            cpu0101
                                      0 ←MPIのRank#0プロセスは
[0] MPI startup(): 1
                                         cpu0101ノードの#0コアに割り当てられた
                    446144
                            cpu0101
[0] MPI startup(): 2
                    446145
                            cpu0101
[0] MPI startup(): 3
                    446146 cpu0101
```

5.3. 高度なMPIの実行[アフィニティ] (3)

■ NECMPIの場合

mpirunのオプションに「-v」オプションを指定します。

mpirun -v -venode -np 3 ./.out

標準出力に以下のように出力されます。

/opt/nec/ve/mpi/libexec/mpid: Creating 1 processes of './a.out' on VE 0 of local host vec0101 /opt/nec/ve/mpi/libexec/mpid: Creating 1 processes of './a.out' on VE 1 of local host vec0101 /opt/nec/ve/mpi/libexec/mpid: Creating 1 processes of './a.out' on VE 2 of local host vec0101 ↑ MPIの1プロセスがvec0101ノードのVE#2に割り当てられた

■ OpenMPIの場合

mpirunのオプションに「-display-map」オプションを指定します。

mpirun \${NQSV_MPIOPTS} -display-map -np 2 ./a.out

標準出力に以下のように出力されます。



5.3. 高度なMPIの実行[アフィニティ] (4)

◆ アフィニティ

各MPIライブラリは、実行時に割り当てるコアを制御することが可能です。 ここでは、ソケット分散、ソケット集中を例に説明します。

ソケット集中

- ソケットに詰めてプロセスを割り当てます。
- 近接プロセスがソケットに集中し、同一メモリを使用するため、ソケット分 散に比べてUPI間の転送時間が短縮されます。

ソケット分散

- ソケットに対して交互にプロセスを割り当てます。
- コアを余らせて実行する場合、ソケット集中に比べ1コア当たりが使用す るメモリバンド幅(メモリ転送速度)を確保できるため、メモリへのデータ アクセスが高速になります。



5.3. 高度なMPIの実行[アフィニティ] (5)

◆ ソケット分散の指定方法 各MPIライブラリは、実行時に割り当てるコアを制御することが可能です。 ここでは、ソケット分散、ソケット集中を例に説明します。

■ ソケット分散

- Intel MPIの場合 実行時に環境変数「I_MPI_PIN_PROCESSOR_LIST="all:map=scatter"」を指定する。
- OpenMPIの場合 mpirunのオプションに「--map-by socket」を指定する。

mpirun \${NQSV_MPIOPTS} -display-map -np 2 --map-by socket ./a.out

■ ソケット集中

- Intel MPIの場合 実行時に環境変数「I_MPI_PIN_PROCESSOR_LIST="all:map=bunch"」を指定する。
- OpenMPIの場合 mpirunのオプションに「--map-by core」を指定する。

mpirun \${NQSV_MPIOPTS} -display-map -np 2 --map-by core ./a.out



5.3. 高度なMPIの実行[アフィニティ] (6)

▶ Intel MPIにおけるノード内プロセス数のマニュアル指定

Intel MPIでは、-ppn、-rr、-perhostオプション(I MPI PERHOST 環境変数含む)を指定することで、 ノード内のプロセスを限定することが可能です。ただし、-machinefile オプションを併用した場合、machinefile オプションが優先されます。

SQUIDのジョブ管理システムでは、NQSV MPIOPTS環境変数の中に、-machinefile オプションを含ん でいるため、-ppn、-rr、-perhost オプションは無効となります。ノード内のプロセス数を限定する場 合は、-l cpunum_job オプションに必要コア数を指定することで可能です。ただし、コアのピニング等 を併用するなどより細かい指定をする場合は、PBS NODEFILE環境変数を利用します。

以下は、PBS NODEFILE環境変数を利用して、2ノードで128プロセス(ノードあたり64プロセス)を生 成するジョブスクリプトのサンプルです。

```
#!/bin/bash
#---- gsub option -----
#PBS -a SQUID
#PBS --group=G01234
#PBS - I elapstim_req=00:30:00
#PBS -b 2
#PBS -T intmpi
#---- Program execution -----
module load BaseCPU/2021
module load xxx/xxx
cd $PBS 0 WORKDIR
mpirun -hostfile ${PBS_NODEFILE} -np 128 -ppn 64 ./a.out
```

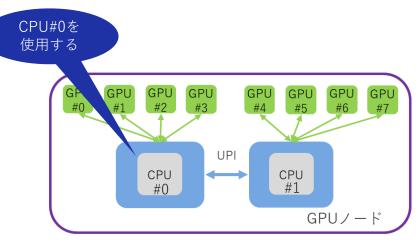
5.3. 高度なMPIの実行[アフィニティ] (7)

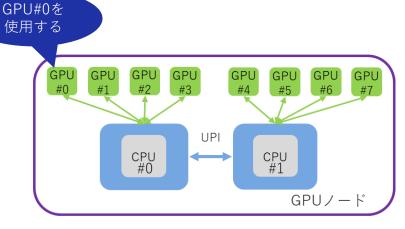
- ◆ GPUノードにおける使用CPU、GPUの指定方法 GPUノードにて使用するCPU、GPUは、NQSVでは制御しておりません。 コマンドや環境変数を利用することで、任意のCPU、GPUに割り当てることが可能です。
 - numactlコマンド numactlコマンドを使用することでCPU番号を指定 することができます。 以下は、CPU#0を使用することを指定しています。

numact| -N 0 -| \$@

■ 環境変数CUDA_VISIBLE_DEVICES 使用するGPU番号を指定することができます。 以下は、GPU#0を使用することを指定しています。

export CUDA VISIBLE DEVICES=0





5.3. 高度なMPIの実行[アフィニティ] (8)

■ MPIランクごとに使用するGPUを指定する方法 ラッパースクリプトを使用することで、MPIランクごとに使用するGPUを指定することができます。

以下は、1ノード内で2並列(1プロセス/ソケット)で実行し、それぞれのソケットにつながるGPUを使用すること想定した実行方法です。

```
WRAP=mpiwrap-single.sh mpirun ${NQSV_MPIOPTS} -mca pml ucx -np 2 -npernode 2 $WRAP ${EXECDIR}/a.out
```

上記で使用した、ラッパースクリプトmpiwrap-single.shの内容は以下のとおりです。

```
#!/bin/bash
echo LANK=${OMPI COMM WORLD LOCAL RANK}
                                                                                             GPU
if [ ${OMPI_COMM_WORLD_LOCAL_RANK} -eq 0 ] ; then
                                                                                #4
        export CUDA VISIBLE DEVICES=0
        numactl -N 0 -I $@
lrank#0
                                                                                   CPU
#1
e se
                                                                    CPU
                                                                     #0
        export CUDA VISIBLE DEVICES=4
                                                                                      GPUノード
        numactl -N 1 -I $@
rank#1
fi
```



5.3. 高度なMPIの実行[アフィニティ] (9)

◆ GPUノードにおけるNVLinkの利用方法

NVLinkとは、GPU間で直接データ転送を行う通信プロトコルです。CPUを介さないため、 高速な通信を行うことが可能です。

シングルノードにて、MPI APIでNVLinkを利用し、転送レイヤにてUCXを使用する (CUDA 対応 UCX + GPUDirect)場合に、指定する環境変数は以下のとおりです。

UCX_TLS=sm,cuda_copy,gdr_copy,cuda_ipc

※UCX(Unified Communication X)とは

HPC 用の通信 API フレームワークです。

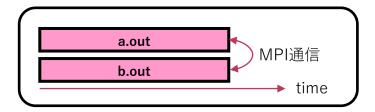
InfiniBand経由のMPI通信のために最適化されています。

Orchestrating a brighter world

5.4. 高度なMPIの実行[MPMD] (1)

- ◆ MPMDジョブ
 - MPMDジョブとは Multi Program Multi Dataの略です。

シミュレーション分野において、異なる2つ(もしくは3つ以上)の実行モジュールを同時に実行し、相互にMPI通信を行いながら連携実行するモデルをMPMDモデルと呼びます。



- 利用可能なMPIライブラリ 以下3つのMPIライブラリは、いずれもMPMDに対応しています。
 - Intel MPI
 - NEC MPI
 - OpenMPI



5.4. 高度なMPIの実行[MPMD] (2)

◆ MPMDジョブの実行方法汎用CPUノードにてa.outをノード#1で76並列、b.outをノード#2,#3で152並列実行する場合

```
#!/bin/bash
#PBS -q SQUID
#PBS --group=G01234
#PBS -T intmpi
#PBS -b 3
: (中略)

module load BaseCPU
export I_MPI_DEBUG=4
cd $PBS_0_WORKDIR

mpirun ${NQSV_MPIOPTS} -np 76 ./a.out : ¥
-np 152 ./b.out
```



5.5. コンテナの実行方法(1)

SQUIDシステムでは、コンテナを利用してジョブを実行することができま す。コンテナを利用してジョブを実行する場合の、ジョブスクリプトにつ いて説明します。(※OCTOPUSシステムでは利用できません)

- ▶ コンテナ実行
 - コンテナの実行は、exec サブコマンドを指定して実施します。 (実行例)centos.sif コンテナイメージ内のhostname コマンドを実行する場合
 - \$ singularity exec centos.sif hostname
 - コマンドの書式は以下のとおりです。
 - \$ singularity exec <イメージファイル名> <コンテナ内の実行コマンド>

指定するコマンドは、**コンテナ内の実行コマンドとなっている点にご注意くだ** さい。コマンドがパス指定なしであれば、コンテナ内のPATH環境変数で探索さ れたコマンドが実行されます。パス指定有りの場合も、絶対パスはコンテナ内 のファイル構造に従います。

5.5. コンテナの実行方法(2)

◆ 環境変数

コンテナ外で定義した環境変数は、基本的にはコンテナ内にも引き継がれます。

ただし、build時などにコンテナ側で明示的に定義されている環境変数は、コンテナ側の定義に従います。

コンテナ側で定義されている環境変数を上書きする場合には、--env オプションによる個別の指定や、--env-fileオプションによる一括の指定でコンテナ内に渡すことが可能です。

• --env オプションによる個別指定

\$ singularity exec --env MYVAR="My Value!" centos.sif myprog.exe

• --env-file オプションによる一括指定

\$ cat myenvfile

MYVAR="My Value!"

\$ singularity exec --env-file myenvfile centos.sif myprog.exe

環境変数に関する詳細な内容は、下記の公式ドキュメントを参照ください。 https://sylabs.io/guides/3.7/user-guide/environment_and_metadata.html



5.5. コンテナの実行方法(3)

◆ ホストOSのマウント コンテナ内からホストOSのファイルシステムのread/writeを行いたい場合には、 ホストOSの特定のディレクトリをbindマウントすることで利用可能です。

■ ホストOSの標準マウントディレクトリ 下記のディレクトリは標準でマウントされており、コンテナ内でも同じパスで利用可能

- ホームディレクトリ:/sqfs/home/(利用者番号)
- テンポラリ領域:/tmp

(例)コンテナ内からホストOSのhomeディレクトリに置かれたプログラム (/sqfs/home/(利用者番号)/a.out)を実行

\$ singularity exec centos.sif ./a.out

上記例では、カレントディレクトリが、コンテナ内でhomeに移動しています。



5.5. コンテナの実行方法(4)

■ ホストOSの特定のディレクトリ

ホストOSの特定のディレクトリをマウントする場合には、--bindオプションを利用します。--bindオプションの書式は以下となります。

--bind <ホストOSのパス>:<コンテナ内のパス>:<モード>

コンテナ内のパス並びにモード(ro/rw)は省略可能です。省略した場合は、コンテナ内のパスは、ホストOSのパスと同じパスでread/writeでマウントされます。

例(拡張領域上のディレクトリに置かれたプログラム(a.out)を実行)

- \$ cd /sqfs/work/(グループ名)/(利用者番号)
- \$ singularity exec --bind `pwd` centos.sif ./a.out

上記例では、コンテナ外での環境変数PWDが、コンテナ内へも引き継がれており、 コンテナ内のカレントディレクトリが、拡張領域上のディレクトリとなっています。

ホストOSのマウントに関する詳細な内容は、下記の公式ドキュメントを参照ください。

https://sylabs.io/guides/3.7/user-guide/bind_paths_and_mounts.html



5.5. コンテナの実行方法(5)

◆ 汎用CPUノードでの実行方法

1ノード内でスレッド並列プログラム実行を想定したジョブスクリプトのサンプルです。

利用時の注意点

ホストOSの特定のディレクトリをマウントする場合には、--bindオプションを利用してください。



5.5. コンテナの実行方法(6)

◆ GPUノードでの実行方法

1ノード内でスレッド並列プログラム実行を想定したジョブスクリプトのサンプルです。

利用時の注意点

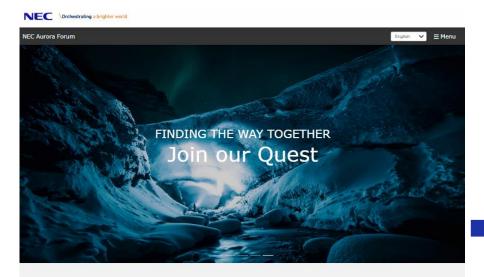
ホストOSの特定のディレクトリをマウントする場合には、--bindオプションを利用してください。



6. SX-Aurora TSUBASA 公開情報

6. SX-Aurora TSUBASA 公開情報

- NEC Aurora Forum
 - https://www.hpc.nec/ 各種ドキュメント、コミュニティ



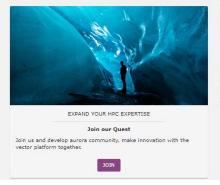




Table 5: NQSV



基本的な操作に関するマニュアルは以下の2点です

- ・NQSV 利用の手引き 操作編 バッチリクエストの作成・投入・操作などの一連の 操作に関するマニュアルです
- ・NQS利用の手引き リファレンス編 バッチリクエスト操作関連のコマンドに関する リファレンスマニュアルです。

フリック

Orchestrating a brighter world

