

ベータ版 oneAPI for AMD* GPU 2023.2.1 ガイド

この記事は、Codeplay 社の許可を得て iSUS (IA Software User Society) が作成した 2023 年 8 月 18 日時点の『oneAPI for AMD* GPUs (beta) 2023.2.1』の日本語参考訳です。原文は更新される可能性があります。原文と翻訳文の内容が異なる場合は原文を優先してください。

[バージョン 2023.0.0 のガイドはこちら。](#)

[バージョン 2023.1.0 のガイドはこちら。](#)

[『oneAPI for NVIDIA* GPU 2023.2.1 ガイド』はこちら。](#)



ベータ版 oneAPI for AMD* GPU は、開発者が DPC++/SYCL* を利用して oneAPI アプリケーションを作成し、それらを AMD* GPU 上で実行できるようにするインテル® oneAPI ツールキット向けのプラグインです。

注意: これはベータ品質のソフトウェアであり、主要機能のほとんどが含まれていますが、まだ完全ではなく、既知および未確認のバグがあることに留意してください。サポートされる機能の詳細については、「[機能](#)」を参照してください。

このプラグインは、HIP バックエンドを DPC++ 環境に追加します。このドキュメントでは、「ベータ版 oneAPI for AMD* GPU」と「DPC++ HIP プラグイン」が同じ意味で使われています。

oneAPI の詳細については、[インテル® oneAPI の概要](#) (英語) を参照してください。

ベータ版 oneAPI for AMD* GPU の使用を開始するには、「[導入ガイド](#)」を参照してください。

導入ガイド

- [ベータ版 oneAPI for AMD* GPU のインストール](#)
- [DPC++ を使用して AMD* GPU をターゲットにする](#)
- [DPC++ のリソース](#)
- [SYCL* のリソース](#)
- [SYCL* アプリケーションのデバッグ](#)

パフォーマンス・ガイド

- [はじめに](#)
- [プログラミング・モデル](#)
- [最適化の目的](#)
- [パフォーマンス解析](#)
- [一般的な最適化](#)

サポート

- [機能](#)
- [更新履歴](#)
- [トラブルシューティング](#)
- [使用許諾契約書 \(英語\)](#)

導入ガイド

ベータ版 oneAPI for AMD* GPU のインストール

このガイドには、DPC++ と DPC++ HIP プラグインバージョン 2023.2.1 を使用して、AMD* GPU で SYCL* アプリケーションを実行する方法を説明します。

これはベータ品質のソフトウェアであり、主要機能のほとんどが含まれていますが、まだ完全ではなく、既知および未確認のバグがあることに留意してください。サポートされる機能の詳細については、「[機能](#)」を参照してください。

DPC++ に関連する一般的な情報は、「[DPC++ のリソース](#)」の節を参照してください。

サポートされるプラットフォーム

このリリースは、次のプラットフォームで検証されています。

GPU ハードウェア	アーキテクチャー	オペレーティング・システム	HIP	GPU ドライバー
AMD Radeon* Pro W6800	gfx1030	Ubuntu* 22.04.2 LTS	5.4.3	6.1.0-1006-oem
AMD Radeon* Pro W6800	gfx1030	Ubuntu* 22.04.2 LTS	4.5.2	6.1.0-1006-oem

- このリリースは HIP 5.x または HIP 4.5.x (5.4.x まで) で動作するはずですが、HIP 5.4.3 と 4.5.2 でのみテストされています。Codeplay は、HIP 5.x よりも古いバージョンでは正常な動作を保証いたしかねます。
 - HIP 5.4.3 は、既存の HIP インストールと共存できます。ROCm* インストール・ガイドの「[複数バージョンの ROCm* インストールでインストーラー・スクリプトを使用する](#)」(英語) の説明を参照してください。
 - HIP から配布されるビットコード・ライブラリーの変更により、HIP 5.5 以降は DPC++ ではサポートされません。これは次のリリースで修正される予定です。
- このリリースは各種 AMD* GPU と HIP バージョンで動作するはずですが、Codeplay は評価されていないプラットフォームでの正常な動作を保証するものではありません。
- このパッケージは Ubuntu* 22.04 でのみテストされていますが、一般的な Linux* システムにインストールできます。
- プラグインは、システムにインストールされている HIP のバージョンに依存します。HIP は Windows* と macOS* をサポートしていないため、これらのオペレーティング・システムではベータ版 oneAPI for AMD* GPU パッケージは利用できません。

要件

- C++ 開発ツールインストールします。

oneAPI アプリケーションをビルドして実行するには、C++ 開発ツールの `cmake`、`gcc`、`g++`、`make` および `pkg-config` をインストールする必要があります。

次のコンソールコマンドは、一般的な Linux* ディストリビューションに上記のツールをインストールします。

Ubuntu*

```
$ sudo apt update
$ sudo apt -y install cmake pkg-config build-essential
```

Red Hat* と Fedora*

```
$ sudo yum update
$ sudo yum -y install cmake pkgconfig
$ sudo yum groupinstall "Development Tools"
```

SUSE*

```
$ sudo zypper update
$ sudo zypper --non-interactive install cmake pkg-config
$ sudo zypper --non-interactive install pattern devel_C_C++
```

次のコマンドで、ツールがインストールされていることを確認します。

```
$ which cmake pkg-config make gcc g++
```

次のような出力が得られるはずです。

```
/usr/bin/cmake
/usr/bin/pkg-config
/usr/bin/make
/usr/bin/gcc
/usr/bin/g++
```

2. DPC++/C++ コンパイラーを含む [インテル® oneAPI ツールキット 2023.2.1](#) をインストールします。
 - インテル® oneAPI ベース・ツールキットは、多くの利用環境に適用できます。
 - oneAPI for AMD* GPU をインストールするには、インテル® oneAPI ツールキットのバージョン 2023.2.1 が必要です。これよりも古いバージョンにはインストールできません。
3. AMD* GPU 向けの GPU ドライバーと ROCm* ソフトウェア・スタックをインストールします。
 - 例えば、ROCm* 5.4.1 の場合、『[ROCm* インストール・ガイド v5.4.1](#)』（英語）の手順に従ってください。
 - `--usecase="dkms,graphics,opencl,hip,hiplibsdk` 引数を指定して `amdgpu-install` インストーラーを起動し、必要となるすべてのコンポーネントを確実にインストールすることを推奨します。

インストール

1. [ベータ版 oneAPI for AMD* GPU のインストーラー](#)（英語）をダウンロードします。
2. ダウンロードした自己展開型インストーラーを実行します。

```
$ sh oneapi-for-amd-gpus-2023.2.1-rocm-5.4.1-linux.sh
```

- インストーラーは、デフォルトの場所にあるインテル® oneAPI ツールキット 2023.2.1 のインストールを検索します。インテル® oneAPI ツールキットが独自の場所にインストールされている場合、`--install-dir /path/to/intel/oneapi` でパスを指定します。
- インテル® oneAPI ツールキットが home ディレクトリー外にある場合、`sudo` を使用してコマンドを実行する必要があります。

環境を設定

1. 実行中のセッションで oneAPI 環境を設定するには、インテルが提供する `setvars.sh` スクリプトを `source` します。

システム全体へのインストールの場合:

```
$ . /opt/intel/oneapi/setvars.sh --include-intel-llvm
```

プライベート・インストールの場合 (デフォルトの場所):

```
$ . ~/intel/oneapi/setvars.sh --include-intel-llvm
```

- `clang++` などの LLVM ツールにパスを追加するには、`--include-intel-llvm` オプションを使用します。
 - ターミナルを開くたびにこのスクリプトを実行する必要があります。セッションごとに設定を自動化する方法については、「[CLI 開発向けの環境変数を設定する](#)」(英語) など、関連するインテル® oneAPI ツールキットのドキュメントを参照してください。
2. HIP ライブラリーとツールが環境内にあることを確認します。
 - `rocminfo` を実行します。実行時の表示に明らかなエラーが認められなければ、環境は正しく設定されています。
 - 問題があれば、環境変数を手動で設定します。

```
$ export PATH=/PATH_TO_ROCM_ROOT/bin:$PATH
```

```
$ export LD_LIBRARY_PATH=/PATH_TO_ROCM_ROOT/lib:$LD_LIBRARY_PATH
```

ROCm* は通常 `/opt/rocm-x.x.x/` にインストールされます。

インストールの確認

DPC++ HIP プラグインのインストールを確認するには、DPC++ の `sycl-ls` ツールを使用して、SYCL* で利用可能な AMD* GPU があることを確認します。AMD* GPU が利用できる場合、`sycl-ls` の出力に次のような情報が表示されます。

```
[ext_oneapi_hip:gpu:0] AMD HIP BACKEND, AMD Radeon PRO W6800 0.0 [HIP 40421.43]
```

- 上記のように利用可能な AMD* GPU が表示されていれば、DPC++ HIP プラグインが適切にインストールされ、設定されていることが確認できます。
- インストールや設定に問題がある場合、「[トラブルシューティング](#)」の「`sycl-ls` の出力でデバイスが見つからない場合」を確認してください。
- 利用可能なハードウェアとインストールされている DPC++ プラグインに応じて、OpenCL* デバイス、インテル® GPU、または NVIDIA* GPU など、ほかのデバイスもリストされることがあります。

サンプルアプリケーションを実行

1. 次の C++/SYCL* コードで構成される simple-sycl-app.cpp ファイルを作成します。

```
#include <sycl/sycl.hpp>

int main() {
    // カーネルコード内で使用する 4 つの int バッファを作成
    sycl::buffer<sycl::cl_int, 1> Buffer(4);

    // SYCL* キューを作成
    sycl::queue Queue;

    // カーネルのインデックス空間サイズ
    sycl::range<1> NumOfWorkItems{Buffer.size()};

    // キューへコマンドグループ (ワーク) を送信
    Queue.submit([&](sycl::handler &cgh) {

        // デバイス上のバッファへの書き込み専用アクセサを作成
        auto Accessor = Buffer.get_access<sycl::access::mode::write>(cgh);

        // カーネルを実行
        cgh.parallel_for<class FillBuffer>(
            NumOfWorkItems, [=](sycl::id<1> WIid) {
                // インデックスでバッファを埋めます
                Accessor[WIid] = (sycl::cl_int)WIid.get(0);
            });
    });

    // ホスト上のバッファへの読み取り専用アクセサを作成。
    // キューのワークが完了するのを待機する暗黙のバリア
    const auto HostAccessor = Buffer.get_access<sycl::access::mode::read>();

    // 結果をチェック
    bool MismatchFound = false;
    for (size_t I = 0; I < Buffer.size(); ++I) {
        if (HostAccessor[I] != I) {
            std::cout << "The result is incorrect for element: " << I
                << " , expected: " << I << " , got: " << HostAccessor[I]
                << std::endl;
            MismatchFound = true;
        }
    }

    if (!MismatchFound) {
        std::cout << "The results are correct!" << std::endl;
    }

    return MismatchFound;
}
```

2. アプリケーションをコンパイルします。

```
$ icpx -fsycl -fsycl-targets=amdgcN-amd-amdhsa -Xsycl-target-backend --  
offload-arch=<ARCH> simple-sycl-app.cpp -o simple-sycl-app
```

ARCH には GPU のアーキテクチャー (例えば gfx1030) を指定します。次のコマンドで確認できます。

```
$ rocminfo | grep 'Name: *gfx.*'
```

出力に GPU アーキテクチャーが表示されます。例えば、次のようになります。

```
Name: gfx1030
```

3. アプリケーションを実行します。

```
$ ONEAPI_DEVICE_SELECTOR="hip:*" SYCL_PI_TRACE=1 ./simple-sycl-app
```

次のような出力が得られます。

```
Warning: ONEAPI_DEVICE_SELECTOR environment variable is set to hip:*.  
To see the correct device id, please unset ONEAPI_DEVICE_SELECTOR.  
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded: libpi_hip.so  
[ PluginVersion: 13.32.1 ]  
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded:  
libpi_unified_runtime.so [ PluginVersion: 13.32.1 ]  
SYCL_PI_TRACE[all]: Requested device_type: info::device_type::automatic  
SYCL_PI_TRACE[all]: Selected device: -&gt; final score = 1500  
SYCL_PI_TRACE[all]: platform: AMD HIP BACKEND  
SYCL_PI_TRACE[all]: device: AMD Radeon Graphics
```

これで、oneAPI for AMD* GPU の環境設定が確認でき、oneAPI アプリケーションの開発を開始できます。

以降では、AMD* GPU で oneAPI アプリケーションをコンパイルして実行するための一般的な情報を説明します。

DPC++ を使用して AMD* GPU をターゲットにする

AMD* GPU 向けのコンパイル

AMD* GPU 対応の SYCL* アプリケーションをコンパイルするには、DPC++ に含まれる clang++ コンパイラを使用します。

例:

```
$ icpx -fsycl -fsycl-targets=amdgcN-amd-amdhsa -Xsycl-target-backend=amdgcN-amd-  
amdhsa --offload-arch=gfx1030 -o sycl-app sycl-app.cpp
```

次のフラグが必要です。

- `-fsycl`: C++ ソースファイルを SYCL* モードでコンパイルするようにコンパイラーに指示します。このフラグは暗黙的に C++ 17 を有効にし、SYCL* ランタイム・ライブラリーを自動でリンクします。
- `-fsycl-targets=amdgcncn-amd-amdhsa`: AMD* GPU をターゲットに SYCL* カーネルをビルドすることをコンパイラーに指示します。
- `-Xsycl-target-backend=amdgcncn-amd-amdhsa --offload-arch=gfx1030:gfx1030 AMD*` GPU をターゲットに SYCL* カーネルをビルドすることをコンパイラーに指示します。

AMD* GPU をターゲットにする場合、GPU の特定のアーキテクチャーを指定する必要があることに注意してください。

利用できる SYCL* コンパイルフラグの詳細は、『[DPC++ コンパイラー・ユーザーズ・マニュアル](#)』(英語)を参照してください。すべての DPC++ コンパイラー・オプションの詳細は、『[インテル® oneAPI DPC++/C++ コンパイラー・デベロッパー・ガイドおよびリファレンス](#)』の「[コンパイラー・オプション](#)」(英語)を参照してください。

icpx コンパイラーを使用する

`icpx` コンパイラーは、デフォルトで `-O2` と `-ffast-math` オプションを有効にするため、通常の `clang++` ドライバーよりも積極的な最適化を行います。多くの場合、これによりパフォーマンスは向上しますが、一部のアプリケーションでは問題が生じる可能性があります。その場合、`-fno-fast-math` を使用して `-ffast-math` を無効にして、`-O2` 以外の `-O` オプションを指定することで最適化レベルを変更できます。`$releasedir/compiler/latest/linux/bin-llvm/clang++` にある `clang++` ドライバーを直接起動することで、通常の `clang++` の最適化レベルを適用できます。

現在、`icpx` コンパイラーには以下に示す既知の問題が報告されています。詳細については、「[トラブルシューティング](#)」を参照してください。

- 倍精度浮動小数点を使用すると、SYCL* グループ・アルゴリズム、`broadcast`、`joint_exclusive_scan`、`joint_inclusive_scan`、`exclusive_scan_over_group`、`inclusive_scan_over_group` がハングアップします。

複数ターゲット向けのコンパイル

AMD* GPU をターゲットにするだけでなく、一度のコンパイルで複数のハードウェア・ターゲットで実行できる SYCL* アプリケーションを生成できます。次の例は、AMD* GPU、NVIDIA* GPU、および SPIR* をサポートする任意のデバイス (インテル® GPU など) で実行できるコードを含む単一のバイナリーを生成する方法を示しています。

```
$ icpx -fsycl -fsycl-targets=amdgcncn-amd-amdhsa,nvptx64-nvidia-cuda,spir64 \  
-Xsycl-target-backend=amdgcncn-amd-amdhsa --offload-arch=gfx1030 \  
-Xsycl-target-backend=nvptx64-nvidia-cuda --offload-arch=sm_80 \  
-o sycl-app sycl-app.cpp
```

AMD* GPU でアプリケーションを実行

AMD ターゲットの SYCL* アプリケーションをコンパイルしたら、ランタイムが SYCL* デバイスとして AMD* GPU を選択しているか確認する必要があります。

通常、デフォルトのデバイスセクターを使用するだけで、利用可能な AMD* GPU の 1 つが選択されます。しかし、場合によっては、SYCL* アプリケーションを変更して、GPU セクターやカスタムセクターなど、より正確な SYCL* デバイスセクターを設定することもあります。

環境変数 `ONEAPI_DEVICE_SELECTOR` を設定して、利用可能なデバイスセットを限定することで SYCL* デバイスセクターを支援できます。例えば、DPC++ HIP プラグインでサポートされるデバイスのみを許可するには、次のように設定します。

```
$ export ONEAPI_DEVICE_SELECTOR="hip:*
```

この環境変数の詳細については、インテル® oneAPI DPC++ コンパイラーのドキュメントで「[環境変数](#)」参照してください。

DPC++ のリソース

- [インテル® DPC++ の概要 \(英語\)](#)
- [DPC++ 導入ガイド](#)
- [DPC++ コンパイラー・ユーザーズ・マニュアル \(英語\)](#)
- [DPC++ コンパイラーとランタイムのアーキテクチャー設計](#)
- [DPC++ 環境変数](#)

SYCL* のリソース

- [SYCL* 2020 仕様](#)
- [SYCL* アカデミー学習教材 \(英語\)](#)
- [Codingame インタラクティブ SYCL* チュートリアル \(英語\)](#)
- [IWOCL SYCL* トーク \(英語\)](#)
- [無料の DPC++ 電子書籍 \(英語\)](#)
- [SYCL* の最新ニュース、学習教材、プロジェクトの紹介 \(英語\)](#)

SYCL* アプリケーションのデバッグ

この節では、さまざまなデバイスで SYCL* アプリケーションをデバッグするための情報、ヒント、およびポイントについて説明します。

SYCL* アプリケーションのホストコードは、単純に C++ アプリケーションとしてデバッグできますが、カーネルデバッグのサポートやツールは、ターゲットデバイスによって異なる可能性があります。

注意: SYCL* アプリケーションに汎用性がある場合、実際のターゲットデバイスではなく、インテルの OpenCL* CPU デバイスなど、豊富なデバッグサポートとツールを備えたデバイスでデバッグしたほうが有用なことがあります。

インテルの OpenCL* CPU デバイスでのデバッグ

インテルの OpenCL* CPU デバイスを使用した DPC++ アプリケーションのデバッグについては、『インテル® oneAPI プログラミング・ガイド』の「DPC++ と OpenMP* オフロードプロセスのデバッグ」の節を参照してください。

ROCm* デバッガーのサポート

ROCm* SDK には `rocgdb` デバッガーが付属しており、HIP アプリケーションの AMD* GPU 上のカーネルをデバッグできます。

ただし、DPC++ では現在、AMD* GPU ターゲットの SYCL* カーネルに対し、適切なデバッグ情報を生成することができません。そのため、`rocgdb` を使用して SYCL* カーネルをデバッグすると、次のようなエラーが表示されることがあります。

```
Thread 5 "dbg" hit Breakpoint 1, with lanes [0-63],
main::{lambda(sycl::_V1::handler&)#1}::operator() (sycl::_V1::handler&)
const::{lambda(sycl::_V1::id<1>)#1}::operator() (sycl::_V1::id<1>) const
(/long_pathname_so_that_rpms_can_package_the_debug_info/src/rocm-
gdb/gdb/dwarf2/frame.c:1032: internal-error: Unknown CFA rule.
```

デバッグ情報を生成せずにアプリケーションをビルドしても、デバッガーは役立ちます。例えば、カーネルが無効なメモリアドレスなどのエラーをスローする場合、`rocgdb` を使用してプログラムを実行することができます。エラー発生時にブレークして、`disas` コマンドを使用してエラーを引き起こした場所のカーネル・アセンブリー行を確認できます。

パフォーマンス・ガイド

はじめに

このガイドは、SYCL* プログラミング・モデルと一般的な GPU におけるパフォーマンスの紹介から始まります。次に、GPU でのパフォーマンス解析の基本と、そこで使用される一般的なツールを紹介します。最後に、ベンダー固有の GPU と利用可能なツールについて紹介します。

GPU に適用される一般的な SYCL* 最適化については、「[一般的な最適化](#)」を参照してください。

NVIDIA* GPU をターゲットにする固有の最適化については、「[NVIDIA* GPU 上のパフォーマンス](#)」を参照してください。

プログラミング・モデル

グラフィックス処理ユニットは、超並列アーキテクチャーにより、CPU よりも 1 秒あたり多くの浮動小数点演算を実行でき、メモリー帯域幅も高くなっています。これらの機能は、コードの開発時点で GPU アーキテクチャーを使用することを選択した場合にのみ活用できます。

ここでは、GPU における大規模並列処理を表現するプログラミング・モデルが基本となります。SYCL* は OpenCL* や CUDA* と同様のプログラミング・モデルを採用しており、カーネル (GPU によって実行される関数) は work-item によって実行される操作で表現されます。

[SYCL* 仕様 \(Rev 6\) の 3.7.2 節](#)では次のように定義されています。

カーネルが実行のため送信されると、インデックス空間が定義されます。カーネルボディのインスタンスは、インデックス空間の各ポイントで実行されます。カーネル・インスタンスは work-item (ワーク項目) と呼ばれ、グローバル id を提供するインデックス空間内のポイントで識別されます。それぞれの work-item は同じコードを実行しますが、コードと操作されるデータの実行パスは、work-item のグローバル id を使用して計算を特殊化することで異なります。

SYCL* では、2 つの異なるカーネル実行モデルを利用できます。

[SYCL* 仕様 \(Rev 6\) の 3.7.2.1 節](#)では次のように記述されています。

`range<N>` (N は 1、2 または 3) で定義される N 次元のインデックス空間でカーネルを呼び出す単純な実行モデルをサポートします。この場合、カーネルの work-item は独立して実行されます。各 work-item は、タイプ `item<N>` の値によって識別されます。タイプ `item<N>` は、タイプ `id<N>` の work-item 識別子と、カーネルを実行する work-item の数を示す `range<N>` をカプセル化します。

[SYCL* 仕様 \(Rev 6\) の 3.7.2.2 節](#)では次のように記述されています。

work-item を work-group に編成できる ND-range の実行モデルは インデックス空間よりも粗い粒度の分解を提供します。それぞれの work-group には、work-item で使用できる

インデックス空間と同じ次元の work-group id が割り当てられます。work-item には、それぞれ work-group 内で一意のローカル id が割り当てられるため、単一 work-item は、グローバル id、またはローカル id と work-group id の組み合わせで識別できます。特定の work-group 内の work-item は、単一の計算ユニットの処理ユニットで同時に実行されます。SYCL* で使用される work-group は、ND-range と呼ばれます。ND-range は、N 次元のインデックス空間であり、N は 1、2 または 3 です。SYCL* では、ND-range は `nd_range<N>` クラスを介して表現されます。`nd_range<N>` は、グローバルレンジとローカルレンジで構成され、それぞれ `range<N>` タイプの値で表現されます。さらに、タイプ `id<N>` 値で表現されるグローバルオフセットが存在することもあります。これは SYCL* 2020 では非推奨です。タイプ `nd_range<N>` と `id<N>` は、それぞれ N 要素の整数配列です。`nd_range<N>` で定義される反復回数は、ND-range のグローバルオフセットで開始される N 次元のインデックス空間であり、サイズはグローバルレンジで、ローカル・レンジ・サイズの work-group に分割されます。ND-range の各 work-item は、タイプ `nd_range<N>` の値によって識別されます。タイプ `nd_range<N>` は、グローバル id、ローカル id、および work-group id をすべて `id<N>` (`id<N>` タイプの反復空間オフセットですが、SYCL* 2020 では非推奨) としてカプセル化し、グローバルとローカルレンジを同期して work-group を有効にします。work-group には、work-item のグローバル id と同様の方法で id が割り当てられます。work-item には work-group とゼロからその次元の work-group サイズから 1 を引いた範囲のコンポーネントを保持するローカル id が割り当てられます。つまり、work-group id と work-group 内のローカル id の組み合わせで work-item が一意に定義されます。

work-item は、次の OpenCL* メモリーモデルに従って 3 つの異なるメモリー領域にアクセスできます。

- **グローバルメモリー:** すべての work-group のすべての work-item 間で共有されます。
- **ローカルメモリー:** 同一 work-group のすべての work-item 間で共有されます。
- **プライベート・メモリー:** 各 work-item でプライベートです。

アーキテクチャー

SYCL* 仕様では、独立して動作する 1 つ以上の計算ユニット (CU) で構成されるデバイスを考慮することで、OpenCL* 1.2 の仕様に従います。NVIDIA では CU を ストリーミング・マルチプロセッサ (*streaming multiprocessor*) と呼び、AMD では単純に計算ユニット (*compute unit*) と呼んでいます。それぞれの CU は、1 つ以上の処理エレメント (PE) とローカルメモリーで構成されます。work-group は単一の CU で実行されますが、work-item は 1 つ以上の PE で実行されることがあります。一般に、CU は SIMD 形式で work-item の小さなセット (*sub-group* として定義) を実行します。sub-group は NVIDIA では ワープ (warp)、AMD では ウェーブフロント (wavefront) と呼ばれます。sub-group サイズは NVIDIA 向けには 32 で、AMD 向けには通常 64 (一部のアーキテクチャー向けには 32) です。

計算

カーネルを構成する work-group は、CU 全体にスケジューリングされます。この時点で、それぞれの CU は処理エレメントで 1 つ以上の *sub-group* を実行します。計算ユニットには、算術演算を実行する整数論理ユニットや浮動小数点ユニット、メモリー操作を行うロード/ストアユニット、超越関数 (正弦、余弦、逆数、平方根など) を実行する特別なユニット、AI で役立つ行列操作など、さまざまな種類の処理エレメントが含まれます。処理エレメントが操作を完了するのに要する時間 (クロックサイクルで測定) は、レイテンシーと呼ばれます。レイテ

ンシーは操作の種類によって異なります。例えば、グローバル・メモリー・トランザクションのレイテンシーは、レジスター呼び出しに比べ桁違いに大きく、これは各種算術演算でも同じことが当てはまります。

スループットは、実行された操作の数と、それらの完了に要する時間の比率です。この比率は、命令のレイテンシーを減らすか、同時に実行する命令数を増やすことで高めることができます。これまで、CPU はクロック周波数を上げて命令レイテンシーを最小化することでスループットを向上させてきました。一方、GPU はレイテンシーを隠匿することでスループットを向上させます。これにより、CU は sub-group 間で「コンテキスト」(レジスター、命令カウンタなど)をわずかな労力に変更できます。そのため、操作に多くのクロックサイクルを要する場合、CU は「コンテキスト」を変更し、別の sub-group の操作を実行することでそれらを隠匿できます。アーキテクチャーによって、同時に実行できる sub-group の最大数は異なります。実際に実行中の sub-group と実行中の sub-group の最大数の比率は「占有率」として定義されます。次の節で詳しく説明します。

GPU における work-item の同時実行は、複数レベルで実現されます。

1. 同一 sub-group 内の異なる work-item は SIMD 形式で同期実行されます。つまり、同じ操作が異なるデータを実行します。
2. 前述したように、CU はレイテンシーを隠匿するため、同一または異なる work-group から複数の sub-group を同時に実行します。
3. GPU を構成する CU は、異なる work-group に属する、異なる sub-group を同時に実行します。

これらの並列実行の機能は、起動されたカーネルが GPU 全体をビジー状態にする十分な大きさの work-item を持っている場合にフル活用されます。

メモリー

次の図は、ディスクリット GPU を搭載したシステムにおける一般的な接続方法を示しています。[1] ホストとデバイスを接続し、[2] CU をグローバルメモリーに接続します。例えば、NVIDIA* GA100 GPU の目安となる帯域幅は次のようになります。[1] PCIe* x16 4.0 では 31GB/秒、および [2] HBM2 では 1555GB/秒。

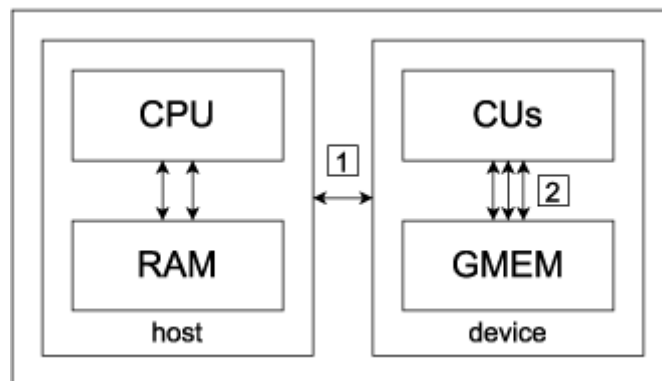


図 1

CPU と GPU 間の接続 [1] が大きなボトルネックになる可能性があります。そのため、ホストとデバイス間のデータ転送を慎重に検討し、GPU 上のデータの局所性を可能な限り維持することが重要です。ただし、カーネルの実行とオーバーラップすることで、PCIe* メモリーのトランザクションで生じるレイテンシーを隠匿することができます。

GPU の主要な特徴として、CU とグローバルメモリー間の高い帯域幅があります [3]。これは、それらを接続するメモリー・コントローラーの数と幅によるものです。例えば、NVIDIA* GA100 GPU には、12 個の 512 ビットの HBM メモリー・コントローラーがあります。これにより、クロックサイクルごとに大量のデータを転送できます。NVIDIA* GA100 GPU では、クロックごとに 6144 ビットです。ただし、この高帯域幅のメモリーを十分に活用するには、メモリーアクセスを結合する必要があります。つまり、work-item はキャッシュに最適な方法でメモリーアクセスする必要があります。

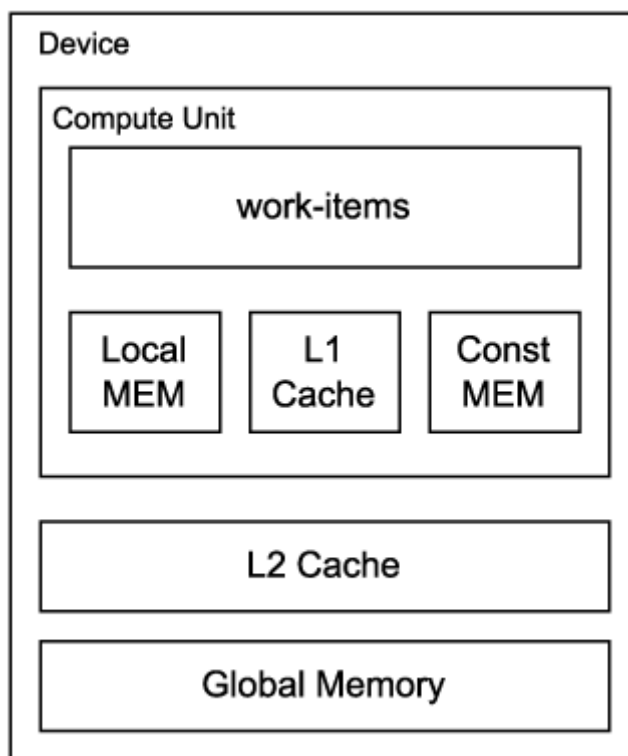


図 2

work-item とグローバルメモリー間にはいくつかのメモリー階層があります。以下に、それらをアクセス・レイテンシーが低い順に示します。

- **レジスター**は、ワークメモリーとして使用される work-item からは透過なデータを維持します。
- **コンスタント・メモリー**は、CU が使用する読み取り専用メモリーです。
- **ローカルメモリー**は CU ごとにあり、同一 work-group 内の work-item 間で共有されます。ローカルメモリーは、グローバルメモリーよりも高速であり、再利用されるグローバルメモリーのデータをキャッシュするために使用されます。
- グローバルメモリー (DDR または HBM) と CU を接続するメモリーシステムを構成する **L1** および **L2** キャッシュ。

最適化の目的

優先度

GPU コードのパフォーマンスに影響する主な要因を重要度の高い順に示します。

- **合成されていないグローバル・メモリー・アクセス。** キャッシュが完全に活用されると、メモリーアクセスは結合され、高い帯域幅を維持できます。結合の方法はアーキテクチャーによって異なりますが、一般に、同じ sub-group 内の work-item が連続したメモリー位置をアクセスすることで実現されます。
- **ローカルメモリーのバンク競合。** ローカルメモリーは複数のバンクに分割されており、異なる work-item から同時にアクセスできます。異なる work-item が同じメモリーバンクにアクセスすると、バンク競合が発生してトランザクションはシリアル化されます。
- **if 文などの条件式やループの反復回数は work-item によって異なるため、同じ sub-group に属する work-item が異なる命令を実行することで発散が発生します。** 近年のアーキテクチャーではこの事象が緩和され、パフォーマンスのペナルティーが軽減されています。

計算の種類が異なれば最適化の優先順位も変わってきます。例えば、メモリー・トランザクションに対し算術演算が少ないメモリー依存のタスクを考えてみます。この場合、GPU を十分に活用するには、メモリーアクセスを結合することが重要です。一方、メモリー・トランザクションに対し算術演算が多い計算依存タスクがあります。この場合、スレッドの発散を回避することが有用な場合があります。算術演算数とリード/ライトデータのバイト数の比率は、**演算強度**として定義されます。

$$I = (\text{浮動小数点操作数}) / (\text{リード/ライトデータのバイト数}) \quad [\text{FLOP/バイト}]$$

ループライン・モデルを利用して、カーネルの演算強度をハードウェア特性に関連付けることで、カーネルがメモリー依存であるか計算依存であるか確認できます。**ループライン・モデル**は 2 次元座標として表示され、x 軸には演算強度が、y 軸には浮動小数点演算のスループット (FLOPS: 1 秒あたりの浮動小数点演算) が示されます。

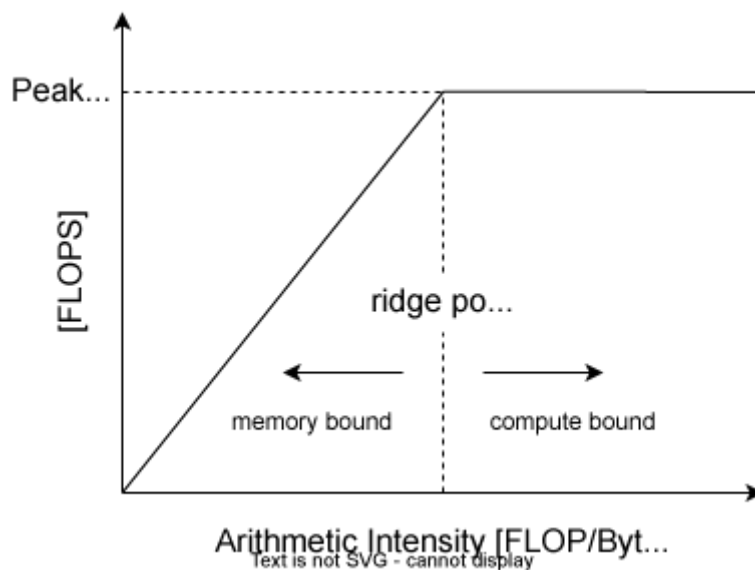


図 3

実際の**ループライン**を構成する最初のセグメントは、 $y = x * B$ を示します。ここで、 B はグローバル・メモリー・システムの帯域幅です。次に水平線 ($y = P_{max}$) は、FMA など特定の操作の最大浮動小数点スループット (P_{max}) に依存します。セグメントが遭遇するポイントは**リッジポイント**と呼ばれます。

カーネルのパフォーマンスは、**ループライン**にプロットされたポイント (点) で示されます。 x 軸はカーネルの演算強度を示し、 y 軸は計測されたカーネルの FLOPS を示します。このポイントがリッジポイントの左にある場合、そのカーネルは**メモリー依存**であり、右にある場合は**計算依存**です。

占有率

カーネルのパフォーマンスを評価するには、その**占有率**を考慮します。占有率は、次のように定義される計算ユニットの式で求められます。

占有率 = アクティブな sub-group 数 / アクティブな sub-group の最大数

アクティブな sub-group は、CU で実際に実行される sub-group です。アクティブな sub-group の最大数は、計算ユニットのアーキテクチャーによって異なります。例えば、NVIDIA* GA100 CU アーキテクチャーでは 64 です。

占有率を高めるにはアクティブな sub-group 数を最大化する必要があります。ただし、計算ユニットのアーキテクチャーによって制約は異なります。

- **work-group あたりの work-item の最大数**
- **CU で実行される work-group の最大数:** work-group サイズが小さすぎると、CU はアクティブな sub-group の最大数を実行することができません。
- **レジスター数の制限:** カーネルコードが複雑になるとレジスターの使用量が増加します。コードを簡素にすることでレジスターの使用量を軽減できます。これは、コードを複数のカーネルに分割することで実現することもできます。
- **ローカルメモリー量の制限:** work-group がローカルメモリーを消費しすぎると、同時に実行できる work-group 数が減少します。

work-group が使用するレジスターやローカルメモリーが多いと、占有率が制限される可能性があります。ユーザーは、work-group のサイズを変更することで占有率を改善できます。このサイズは、sub-group サイズの倍数で、アクティブな sub-group の最大数の除数である必要があります。

例えば、NVIDIA* GA100 GPU では、各 work-item は最大 32 個のレジスターを使用して完全な占有を実現できます。

$r_{max} = (\text{CU ごとのレジスター数}) / (\text{アクティブな sub-group の最大数}) * (\text{sub-group サイズ}) = 32$

work-item が 32 未満のレジスターを使用する場合、CU で同時に実行できる work-group の最大数 wg_{max} とすると、ローカルメモリーにも同じことが当てはまります。

$wg_{max} = (\text{アクティブな sub-group の最大数}) * (\text{sub-group サイズ}) / (\text{実際の work-group サイズ})$

各 work-group が最大 48Kb / wg_{max} のローカルメモリーを割り当てる場合、完全な占有率が得られます。

実効占有率は重要ですが、パフォーマンスにおける最重要のメトリックではありません。命令レベルの並列処理が十分にあり、同じ *sub-group* に属する独立した命令の同時実行が可能であれば、低い占有率でもレイテンシーを十分に隠匿できます。これについては、[こちら](#) (英語) をご覧ください。

さらに、GPU のすべての CU を利用するためカーネルで起動される work-item の最小数は、少なくとも次の `wi_min` でなければなりません。

$$wi_min = (\text{アクティブな sub-group の最大数}) * (\text{sub-group サイズ}) * (\text{CU 数}) = 262144$$

これらのパラメーターはすべて、特定ベンダーの各アーキテクチャーのドキュメントに記載されていますが、以下の表にいくつかの一般的な GPU アーキテクチャーの数値を示します。

一般的な GPU アーキテクチャーのリファレンス占有率					
アーキテクチャー	アクティブな sub-group の最大数	work-item の最大数	work-group の最大数	レジスター数	ローカルメモリー (バイト)
NVIDIA* S.M. 7.0	64	2048	32	65536	65536
NVIDIA* S.M. 7.5	32	1024	16	65536	65536
NVIDIA* S.M. 8.0	64	2048	32	65536	65536
AMD* GFX9xx	40 ¹	1024	16	29184 (?)	65536

- 1 この図は、AMD アーキテクチャー全般に適用されます。work-group が 1 つの sub-group (例えば 64 work-item) のみである場合、CU あたりの work-group の最大数は 40 です。

NVIDIA* GPU の場合、[NVIDIA* Nsight* Compute](#) (英語) は占有計算機を提供しており、理論上の占有率がどのように計算されるか判断するのに役立ちます。非推奨ですが、NVIDIA の [spreadsheet](#) (英語) でも同様の機能を提供することを示しています。

パフォーマンス解析

パフォーマンス解析と最適化は繰り返し作業です。開発者は、ツールを使用してアプリケーションのパフォーマンスを測定しボトルネックを特定して、それらを改善しながら、この手順を繰り返します。それぞれの反復作業で、以前は見つからなかったボトルネックが明らかになることがあります。

ある時点で、アプリケーションの制限要因となる部分で、可能な限り高いパフォーマンスを特定することが重要です。これは、光速またはルーフラインと呼ばれることもあり、アプリケーションの理論上のピーク・パフォーマンスを予測したり、そのパフォーマンスにどれだけ近づいているかを判断するのに役立ちます。

以降の節では、解析ツールと制限要因について詳しく説明します。

解析の方法論

パフォーマンス解析に使用されるツールはプロファイラーとも呼ばれます。プロファイルという用語はいろいろな意味で使用されます。ここでは、パフォーマンス解析に使用されるツールの総称という意味で使用します。特定のパフォーマンス・ツールの説明では、より具体的な意味で使用されることがあります。

パフォーマンス解析は、大きく分けてトレースとサンプリングに分類されます。トレースは、アプリケーションの実行中に 1 つ以上のイベントが発生するたびに記録します。サンプリングは、実行中のアプリケーションの状態を定期的に調査して、その状態を記録します。頻繁に発生するイベントでは、トレースで大量のデータが蓄積される可能性があります。サンプリングでは、サンプリング間隔を調整することでデータ量を制御できます。間隔を長くするとデータ量は減りますが、短い間隔の動作を記録できないことがあります。どちらも改良すべき点がありますが、トレースまたはサンプリングのいずれかを実行中のデータ軽減と組み合わせることができます。

どの解析ツールでも、考慮すべき 2 つのことがあります。

- **オーバーヘッド:** ツールが通常のプログラムの実行時間をどれくらい増加させるかを表わします。パフォーマンス・ツールは、オーバーヘッドを最小限に抑えることが求められます。ただし、オーバーヘッドの増加を十分に理解している場合は、データを解釈する際にこれを補正できます。例えば、あるツールはコードの GPU 実行領域では正確な結果を提供し、CPU 実行領域では実行時間が長くなります。
- **データ量:** 生成されるファイルの大きさを示します。データ量が多いと、オーバーヘッドも増加します。また、大きな出力データセットは管理が困難で、特に出力データセットを表示するためリモートマシンに移動する場合、後処理ツールの応答性にも問題があります。

システムレベルの解析

システムレベルの解析では、同一ノードまたは異なるノード上のプロセス間の相互作用、および CPU と GPU 間の相互作用を調査します。

複雑なワークロードにおける CPU と GPU 間の相互作用を解析するのは困難なことがあります。ベンダーは、このような解析を支援するためトレースツールを提供することがあります。それらは、メモリー割り当て、メモリー転送、カーネルの起動、同期など、GPU 間の API 呼び出しのタイムスタンプと期間を記録します。これらのツールには、シリアル化や過度のアイドル時間などのボトルネックを視覚的に特定するタイムライン表示が含まれます。

状況に応じて、OS のカーネルトレース (Linux* ftrace など) を使用して、それをアプリケーションの実行に関連付けると便利です。これには、root 権限が必要になります。パフォーマンスの問題に関連するカーネルのアクティビティーが理解できない場合は、循環バッファーを利用するすべての OS アクティビティーを記録し、パフォーマンスの問題が検出されたときにアプリケーションの制御下でバッファーをダンプすると便利です (例えば、タイムステップが平均時間や予測時間よりも大幅に長くかかる場合)。循環バッファーによる手法は、トレース・データ・ストリーム全体を記録する際にコストが高い場合に有効です。

分散アプリケーションのスケールリング (通常 **メッセージ・パッシング・インターフェイス** (英語) を使用) は特筆に値します。一般に使用されるスケールリングには 2 つの定義があります。**強力なスケールリング**は、問題のサイズを一定に保ち、MPI ランクの数が増加するのにしたがって経過時間を測定します。**弱いスケールリング**は、MPI ランクの数に比例して問題サイズを大きくします。

強力なスケーリングはより困難な問題です。多くの場合、すべての MPI ランクをビジーに保つのに十分なワークがありません。MPI プロファイル・ツールを利用して、異なる数のランクでスイープを実行し、MPI プロファイルを比較することが有用です。

特に大規模なスケーリングでは、そのほかの MPI の問題がしばしば発生します。リダクション操作は $\log(N)$ に反比例します。さらに、小さなリダクション操作 (スカラー値への MPI_Allreduce など) は、OS によるノイズの影響を受ける可能性があります。ネットワークが混雑する可能性があるため、大規模な共有クラスターではポイントツーポイント操作でも影響を受ける可能性があります。強力なスケーリングでは、メッセージサイズは通常、ランク数が多いほど小さくなるため、MPI レイテンシーがさらに重要になります。

開発者は、アプリケーションの動作が大規模なノードと小規模なノードで実行される際の違いを予測する必要があります。通常のように MPI プロファイル・ツールを使用すると、動作の違いを理解するのに役立ちます。オーバーヘッドの低いツールは、大規模なケースでは特に重要です。

カーネルレベルの解析

カーネルレベルの解析では、GPU カーネルの実行に費やされた時間と、個々の GPU カーネルのパフォーマンスに注目します。

前の節で説明したツールは、通常、起動パラメーター、起動回数、カーネルで消費された時間など、カーネル実行ごとのサマリーを示します。アプリケーションの合計時間は、CPU の経過時間と GPU の経過時間の合計として見積もられることが多く、GPU での経過時間はカーネルの実行時間の合計として概算されます。これにより、GPU カーネルの実行時間を改善することで、全体でどれだけ改善されるかが分かります。実行がオーバーラップしていたり、データの転送時間が長い場合は、常に正確であるとは限りませんが、経験則としては適切です。

カーネルのパフォーマンスを詳しく解析するには以下が必要です。

- カーネルのソースコードを調査
- カーネル向けにコンパイラーが生成するアセンブリー言語の調査
- カーネル実行中のハードウェア・パフォーマンス・メトリックの収集

アセンブリー言語を生成する方法は、コンパイラーと GPU によって異なります。詳細については以降で説明します。

ここからは、GPU (多くの場合 CPU にも該当) で利用可能なメトリックと、それらを解釈してパフォーマンスを改善する作業で導入できる一般的な手法について説明します。異なる GPU 向けの詳細については、このドキュメントの後半で説明します。

重要な GPU メトリック

レートメトリック

アプリケーションが GPU を使用するのには、利用可能な計算リソースを増やすためです。通常、計算スループットは、単位時間あたりに処理される演算数で示されます。例えば、倍精度浮動小数点演算の数/秒、32 ビット整数演算の数/秒などです。特定の GPU では、これらのピーク値が公開されています。

多くの場合、アプリケーションのピーク・パフォーマンスは、非計算リソース (特にメインメモリーやスクラッチパッド・メモリーなどさまざまなメモリー領域) へのアクセスによって制限されます。ここにもピーク値があります。例えば、メインメモリーの帯域幅は、単位時間あたりのバイト数で表現されます。

従来のルーファインのようなモデルでは、ほかのリソースによる制限 (一般的なものはメインメモリーの帯域幅) を考慮して、達成可能な計算パフォーマンスを定量化しようとしています。アプリケーションが計算以外のメトリックでピーク・パフォーマンスに達している場合、ピーク計算パフォーマンスを達成することはできません。これにより、開発者は、アプリケーションで達成可能なピーク・パフォーマンスに関する情報を得ることができます。

利用率メトリック

特定のリソースや機能ユニットがどれだけビジーであるかを知るのには有用です。この利用率メトリックは、レートメトリックとは異なります。リソースの利用率が高くても、ピーク・パフォーマンスにほど遠い場合があります。1つの例として、メモリー・アクセス・パターンが不均一なカーネルが挙げられます。この場合、メモリー帯域幅がピークから離れていても、メモリーユニットの利用率は非常に高くなる場合があります。利用率メトリックは、ルーファイン・モデルでは明らかにならないボトルネックを理解するのに役立ちます。

通常、利用率メトリックはメモリーユニットと計算ユニットで利用できます。また、キャッシュやローカルメモリーなど、各種マイクロアーキテクチャー・ブロックでも利用できる場合があります。

発散

前述のように、GPU は複数の work-item を SIMD (単一命令複数データ) 方式で同時に実行する複数の計算ユニット (CU) で構成されています。

開発者は、単一の work-item に対して実行する操作を記述します。コンパイラーは、このコードを複数の work-item を同時に処理する命令に変換します。各 GPU には、sub-group サイズと呼ばれる、同時に実行される work-item の最小数がネイティブに設定されています。

発散 (Divergence) は、異なる work-item が異なるパスをたどることで発生します。多くの work-item が特定の命令で実行される場合、コンパイラーは可能なすべてのパスの組み合わせを考慮して命令を生成する必要があります。特定の命令で非アクティブな work-item は無効になります。これにより SIMD レーンの一部しか利用されないため、利用率は低下します。

GPU は発散を測定するメトリックを提供し (通常、sub-group ごとにアクティブな work-item)、ネイティブの sub-group サイズと比較できます。

占有率

GPU の占有率については前述しましたが、これは簡単に言うと、特定のカーネルで実際にアクティブな sub-group の数と、アクティブな sub-group の理論上の最大数との比率です。占有率は、カーネルが利用可能な最大の並列性をどれくらい活用できているかを開発者に示すことから重要です。

一部の GPU には、実際の占有率を測定するハードウェア機能が備わっています。理論上の占有率は、コンパイルされたカーネルとハードウェアのプロパティから計算できます。

起動パラメーター

カーネルは、グローバルレンジとローカルレンジで起動されます。後者は work-group のサイズです。work-group のサイズは、sub-group サイズの倍数である必要があります。そのため、グローバル問題サイズを切り上げたり、グローバル問題サイズ外の work-item を処理しないようカーネルにコードを追加する必要があります。

占有率を改善するため、特定の GPU ハードウェアの work-group サイズに制約が課される場合があります。CU 数など、特定の GPU ハードウェアと何らかの関連性のあるグローバル問題サイズを選択することも有益な場合があります。グローバルとローカルの問題サイズは自然なサイズに合わせる必要がなく、ハードウェアに適合するように選択できます。

すべての GPU は、カーネルの起動ごとに実際の起動パラメーターを確認するメカニズムを提供しています。これには、グローバルおよびローカル問題サイズ、レジスター数、およびローカル・メモリー・サイズなどのカーネル・プロパティーが含まれます。

一般的な最適化

ここでは、DPC++ を使用する際の一般的なパフォーマンスの問題や落とし穴、そしてその対処方法について説明します。

インデックスの入れ替え

SYCL* 仕様の [4.9.1 節](#)では、次のことが規定されています。

整数から多次元 id やレンジを構成する場合、多次元空間の線形化において右端の要素が最も速く変化するように要素を記述します。

そのため、インテル® DPC++ コンパイラーでは、右端の次元が CUDA* または HIP の x 次元にマップされ、右から 2 つ目の次元が CUDA* や HIP の y 次元にマップされます。以下に例を示します。

```
cgh.parallel_for(sycl::nd_range{sycl::range(WG_X), sycl::range(WI_X)}, ...)
```

```
cgh.parallel_for(sycl::nd_range<2>{sycl::range<2>(WG_Y, WG_X),  
sycl::range<2>(WI_Y, WI_X)}, ...)
```

```
cgh.parallel_for(sycl::nd_range<3>{sycl::range<3>(WG_Z, WG_Y, WG_X),  
sycl::range<3>(WI_Z, WI_Y, WI_X)}, ...)
```

WG_X と WI_X は、x 次元の **work-group 数**と **work-group ごとの work-item 数** (CUDA* では、名前付きの **グリッドサイズ**と **ブロックあたりのスレッド数**) であり、_Y と _Z は y 次元と z 次元のものになります。

次の場合、2 次元または 3 次元のカーネルの parallel_for 実行では特に重要であることに注意してください。

- ローカルまたはグローバルメモリー内の (1-d) 配列には、手動で線形化されたアクセスがあります。結合されていないグローバル・メモリー・アクセスまたはローカルメモリー内のバンク競合によるパフォーマンスの問題を回避するため、これを考慮する必要があります。線形化の詳細については、SYCL* 仕様の [3.11 節](#)の多次元オブジェクトと線形化を参照してください。

- 次のエラー (または同等のエラー) が発生します。

```
Number of work-groups exceed limit for dimension 1 : 379957 > 65535
```

これは、CUDA* など一部のプラットフォームでは、x 次元が y および z 次元よりも多くの work-group をサポートするためです。

```
Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
```

このトピックの詳細については、[こちら](#) (英語) を参照してください。

インライン展開

場合によっては、DPC++ がインライン展開に対し保守的になり、パフォーマンスが低下することもあります。

これが発生する一般的なケースは、カーネルラムダがカーネル実装を含む大きな関数を単純に呼び出すような SYCL* アプリケーションの場合です。

この問題を回避するには、特定の関数に `always_inline` 属性を追加して、インライン展開を強制します。以下に例を示します。

```
__attribute__((always_inline)) void function(...) {
    ...
}

...

q.submit([&](sycl::handler &cgh) {
    cgh.parallel_for(..., [=](...) {
        function(...);
    });
});
```

高速数学ビルトイン

SYCL* 数学ビルトインは、同等の OpenCL* 1.2 数学ビルトインの精度要件と一致するように定義されていますが、一部のアプリケーションでは必要以上に精度が高くなり、パフォーマンスが低下する可能性があります。

これに対処するため、SYCL* 仕様では、数学関数のサブセットのネイティブバージョン (4.17.5 節の「[数学関数](#)」に完全なリストがあります) が提供されています。これには、精度とパフォーマンスのトレードオフがあります。これらは、ネイティブ名前空間で定義されています。例えば、`sycl::cos()` のネイティブバージョンは、`sycl::native::cos()` です。

一般に、精度が問題にならない場合、ネイティブバリエーションを使用すると大幅に改善できる可能性があります。すべてのバックエンドがすべてのビルトインに対し緩和された精度を使用するわけではないことに注意してください。

注意: `-ffast-math` コンパイルオプションは、標準の `sycl::` 数学関数を、対応する `sycl::native::` 関数に入れ替えます (利用可能であれば)。指定された数学関数のネイティブバージョンが存在しない場合、`-ffast-math` フラグは影響しません。

icpx コンパイラーでは `-ffast-math` がデフォルトです。icpx で `-ffast-math` を無効にするには、`-fno-fast-math` を使用します。

ループアンロール

コンパイラーは一部のループアンロールを自動的に行いますが、次のように `unrolling` プラグマを使用して、デバイスコード内の計算集約型ループのアンロールを手動でコンパイラーに指示することが有益な場合もあります。

```
#pragma unroll <unroll factor>
for( ... ) {
    ...
}
```

エイリアス解析

エイリアス解析では、2 つのメモリー参照が互いにエイリアスでないことが証明できます。これにより最適化が有効になることがあります。デフォルトでは、コンパイラーはエイリアス解析によって証明されない限り、メモリー参照はエイリアスであると想定する必要があります。ただし、デバイスコード内のメモリー参照がエイリアスではないことをコンパイラーに明示的に通知することもできます。これは、バッファー/アクセサーと USM モデルのそれぞれのキーワードを使用することで実現できます。

前者は、`oneapi` 拡張の `no_alias` プロパティをアクセサーに追加することができます。

```
q.submit([&](sycl::handler &cgh) {
    sycl::accessor acc{...,
    sycl::ext::oneapi::accessor_property_list{sycl::ext::oneapi::no_alias}};
    ...
});
```

後者の場合、`__restrict` 修飾子をポインターに追加できます。

`__restrict` は C++ では非標準であり、SYCL* 実装全体で一貫性がない可能性があることに注意してください。dpc++ では、`restrict` 修飾されたデバイス関数 (SYCL* カーネルから呼び出される関数) パラメーターのみが考慮されます。

例:

```
void function(int *__restrict__ ptr) {
    ...
}

...
int *ptr = sycl::malloc_device<int>(..., q);
...
q.submit([&](sycl::handler &cgh) {
    cgh.parallel_for(..., [=](...) {
        function(ptr);
    });
});
```

より強制的なアプローチは、`[[intel::kernel_args_restrict]]` 属性をカーネルに追加することです。これは、各 USM ポインター間、またはそのモデルがカーネル内で使用される場合はバッファアクセス間のすべてのエイリアス依存関係を無視するようにコンパイラーに指示します。

例 (バッファ/アクセサーモデル):

```
q.submit([&](handler& cgh) {
    accessor in_accessor(in_buf, cgh, read_only);
    accessor out_accessor(out_buf, cgh, write_only);
    cgh.single_task<NoAliases>([=]() [[intel::kernel_args_restrict]] {
        for (int i = 0; i < N; i++)
            out_accessor[i] = in_accessor[i];
    });
});
```

CUDA* プラットフォームでは、[sycl_ext_oneapi_cuda_tex_cache_read 拡張 \(英語\)](#) の `ldg` テンプレート関数を使用することも、メモリアクセスのパフォーマンス向上に役立ちます。

サポート

機能

コア機能

機能	サポート
コンテキスト内の複数デバイス	いいえ
サブグループ (sub-group)	部分的 ¹
グループ関数/アルゴリズム	部分的 ¹
整数関数	はい
数学関数 (スカラー)	はい
数学関数 (ベクトル)	はい
数学関数 (marray)	いいえ
共通関数	はい
ジオメトリー関数	はい
リレーショナル関数	はい
atomic ref	はい
オペレーティング・システム	Linux*
バッファの再解釈	はい
stream	いいえ
デバイスイベント	いいえ
グループの非同期コピー	はい
プラットフォームの get info	はい
カーネルの get info	はい
<code>sycl::nan</code> と <code>sycl::isnan</code>	はい
デバイスセレクター	いいえ
階層的並列化	はい
ホストタスク	はい
インオーダー・キュー	はい
リダクション	部分的 ¹
キューのショートカット	はい
vec	はい
marray	はい
errc	はい
匿名カーネルラムダ	はい
機能を評価するマクロ	はい

機能	サポート
sycl::span	はい
sycl::dynamic_extent	いいえ ²
sycl::bit_cast	はい
aspect_selector	いいえ
カーネルバンドル	いいえ
特殊化定数	いいえ

非コア機能

機能	サポート
image	いいえ
fp16 データタイプ	いいえ
fp64 データタイプ	はい
prefetch	はい
USM	ホスト、デバイス、共有
USM アトミックホスト割り当て	はい
USM アトミック共有割り当て	問題あり [^shared-usm-status]
USM システムに割り当て	はい
SYCL_EXTERNAL	いいえ
アトミックメモリーの順序付け	relaxed
アトミック・フェンス・メモリーの順序付け	いいえ
アトミック・メモリー・スコープ	work_group
アトミック・フェンス・メモリーのスコープ	いいえ
64 ビット・アトミック	いいえ
バイナリー形式	AMDGCN
デバイスのパーティション化	いいえ
ホストデバッグ可能デバイス	いいえ
オンラインコンパイラ	いいえ
オンラインリンカー	いいえ
キューのプロファイル	はい
mem_advise	いいえ
バックエンド仕様	いいえ
アプリケーション・バックエンドの相互運用	いいえ
カーネル・バックエンドの相互運用	いいえ
ホストタスク (ハンドルと相互運用)	いいえ
reqd_work_group_size	いいえ

機能	サポート
キャッシュビルド結果	いいえ
ビルドログ	いいえ
ビルトインカーネル関数	なし

拡張機能

機能	サポート
uniform	いいえ
USM アドレス空間 (デバイス、ホスト)	部分的 ³
固定ホストメモリの使用	はい
サブグループ・マスク (+ グループ投票)	いいえ
静的ローカルメモリ使用量照会	いいえ
sRGB イメージ	いいえ
デフォルト・プラットフォーム・コンテキスト	はい
メモリチャネル	いいえ
最大ワークグループ照会	一部分
結合行列	いいえ
すべて制限 (restrict all)	いいえ
プロパティ・リスト (property list)	いいえ
カーネル・プロパティ (kernel properties)	いいえ
SIMD 呼び出し	いいえ
低レベルデバイス情報	いいえ
カーネルキャッシュ設定	いいえ
FPGA lsu	いいえ
FPGA reg	いいえ
データ・フロー・パイプ	いいえ
キューに投入されたバリア	いいえ
フィルターセクター	はい
グループソート	はい
フリー関数の照会	はい
明示的な SIMD	いいえ
discard_queue_events	部分的 ¹
device_if	いいえ
device_global	いいえ
C と C++ 標準ライブラリーのサポート	はい?
カーネルでの assert	はい?

機能	サポート
buffer_location	いいえ
accessor_property_list (+ no_offset, no_alias)	はい
グループ・ローカル・メモリー	はい
printf	いいえ
ext_oneapi_bfloat16	いいえ
拡張デバイス情報	いいえ
sycl_ext_oneapi_cuda_tex_cache_read	はい ⁴
sycl_ext_oneapi_native_math	はい
sycl_ext_oneapi_bfloat16_math_functions	いいえ

1 (1、2、3、
4) 一部のテストで失敗

2 numeric_limits<size_t>::max() の使用

3 <https://github.com/intel/llvm/pull/6289> (英語) に追加 (未テスト)

4 技術的にサポートされますが、sycl_ext_oneapi_cuda_tex_cache_read は NVIDIA* GPU のみ有効です。

更新履歴

2023.2.0

改良点

SYCL* コンパイラー

- gfx9+ HIP atomic が追加されました [b13561c9]。
- basic HIP atomic が追加されました [c3c5e923]。
- AMD HIP プラットフォームで `__CUDA_ARCH__` マクロが無効化されました [8a7cf2b2]。

SYCL* ライブラリー

- AMD バックエンドで `sycl_ext_oneapi_memcpy2d` をサポートしました - [oneAPI memcpy2d](#) (英語) [9008a5d2]。
- PCI デバイス ID と UUID のサポートが追加されました [e09ff588]。
- `SYCL_PI_HIP_MAX_LOCAL_MEM_SIZE` 環境変数がサポートされました [92f6d688]。
- `-fsycl-targets` で `amd-gpu-gfx1034` を指定できるようになりました [5e86a41d]。

バグフィックス

- 無効な work-group サイズに関するエラーを `PI_ERROR_INVALID_WORK_GROUP_SIZE` に置き換えるようになりました [2357af0a]。

- `sycl::ctz` 関数からの間違っただ結果に対応しました [5a9f601e]。
- イベントが意図したとおりに待機しない原因となる問題に対処しました [1b225447]、[ce7c594f]。

2023.1.0

改良点

SYCL* コンパイラー

- `-fsycl-targets` の引数として AMD* アーキテクチャー (`amd_gpu_gfx1032` など) を指定できるようになりました [e5de913f]。

SYCL* ライブラリー

- デバイス拡張に `cl_khr_fp64` が追加されました [cd832bff]。
- HIP* バックエンドのゼロ・レンジ・カーネルをサポートしました [a3958865]。

バグフィックス

- ガードが正しく構築されない問題を修正しました [ce7c594f]。
- 相互運用ヘッダーとデバイスの特実化が追加されました [998fd91e]。

2023.0.0

oneAPI for AMD* GPU の最初のベータリリースです。

このリリースは、[intel/llvm repository at commit 0f579ba](#) (英語) から作成されました。

新機能

- HIP バックエンドのベータサポート

SYCL* コンパイラー

- デバイスでの `assert` をサポート
- ローカル・メモリー・アクセサーのサポート
- グループ集合関数のサポート
- `sycl::ext::oneapi::sub_group::get_local_id` のサポート

SYCL* ライブラリー

- `atomic64` デバイス機能の照会をサポート
- SYCL* キューごとに複数の HIP ストリームをサポート
- 相互運用のサポート
- `sycl::queue::submit_barrier` のサポート

トラブルシューティング

この節では、トラブルシューティングのヒントと一般的な問題の解決方法について説明します。ここで説明する方法で問題が解決しない場合は、[Codeplay のコミュニティー・サポート・ウェブサイト \(英語\)](#) からサポートリクエストをお送りください。完全なサポートは保証できませんが、できる限り支援させていただきます。サポートリクエストを送信する前に、ソフトウェアが最新の安定したバージョンであることを確認してください。

問題、パフォーマンス、機能要望は、[oneAPI DPC++ コンパイラーのオープンソース・リポジトリ \(英語\)](#) から報告できます。

sycl-ls の出力にデバイスが表示されない

sycl-ls がシステム上の期待されるデバイスを報告しない場合:

1. システムに互換性のあるバージョンの CUDA* または ROCm* ツールキット (それぞれ CUDA* と HIP プラグイン向け)、および互換性のあるドライバーがインストールされていることを確認してください。
2. nvidia-smi または rocm-smi がデバイスを正しく認識できることを確認します。
3. プラグインが正しくロードされていることを確認します。これは、環境変数 SYCL_PI_TRACE に 1 を設定して、sycl-ls を再度実行することで分かります。

例:

```
$ SYCL_PI_TRACE=1 sycl-ls
```

次のような出力が得られるはずです。

```
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded: libpi_opencl.so
[ PluginVersion: 11.15.1 ]
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded:
libpi_level_zero.so [ PluginVersion: 11.15.1 ]
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded: libpi_cuda.so
[ PluginVersion: 11.15.1 ]
[ext_oneapi_cuda:gpu:0] NVIDIA CUDA BACKEND, NVIDIA A100-PCIE-40GB 0.0
[CUDA 11.7]
```

インストールしたプラグインが sycl-ls の出力に表示されない場合、SYCL_PI_TRACE に -1 を設定して再度実行することで、詳細なエラー情報を取得できます。

```
$ SYCL_PI_TRACE=-1 sycl-ls
```

大量の出力が得られますが、次のようなエラーが表示されているか確認してください。

```
SYCL_PI_TRACE[-1]:
dlopen (/opt/intel/oneapi/compiler/2023.0.0/linux/lib/libpi_hip.so) failed
with <libamdhip64.so.4: cannot open shared object file: No such file or
directory>
SYCL_PI_TRACE[all]: Check if plugin is present.Failed to load plugin:
libpi_hip.so
```

- CUDA* プラグインには、CUDA* SDK で提供される libcuda.so と libcupti.so が必要です。

- HIP プラグインには、ROCm* の libamdhip64.so が必要です。

CUDA* または ROCm* のインストールと、環境が適切に設定されていることを確認してください。また、LD_LIBRARY_PATH が上記のライブラリーを検出できる場所を指しているか確認してください。

4. ONEAPI_DEVICE_SELECTOR または SYCL_DEVICE_ALLOWLIST などのデバイスフィルター環境変数が設定されていないことを確認します (ONEAPI_DEVICE_SELECTOR が設定されていると、sycl-ls は警告を表示します)。
5. 権限を確認します。POSIX* では、アクセラレーター・デバイスへのアクセスは、通常、適切なグループのメンバーであることを条件としています。例えば、Ubuntu* Linux* の場合、GPU へのアクセスには video グループと render グループのメンバーである必要がありますが、これは設定によって異なります。

不正バイナリーエラーの扱い

CUDA* または HIP をターゲットにする SYCL* アプリケーションを実行すると、特定の状況でアプリケーションが失敗し、無効なバイナリーであることを示すエラーが報告されることがあります。例えば、CUDA* の場合は CUDA_ERROR_NO_BINARY_FOR_GPU がレポートされる場合があります。

これは、選択された SYCL* デバイスに適切でないアーキテクチャーのバイナリーが送信されたことを意味します。この場合、次の点を確認してください。

1. アプリケーションが、利用するハードウェアのアーキテクチャーと一致するようにビルドされていることを確認してください。
 - CUDA* 向けのフラグ:
-Xsycl-target-backend=nvptx64-nvidia-cuda --cuda-gpu-arch=<arch>
 - HIP 向けのフラグ:
-Xsycl-target-backend=amdgcN-amd-amdhsa --offload-arch=<arch>
2. 実行時に適切な SYCL* デバイス (ビルドされたアプリケーションのアーキテクチャーに一致するもの) が選択されていることを確認します。環境変数 SYCL_PI_TRACE=1 を設定すると、選択されたデバイスに関連するトレース情報を表示できます。以下に例を示します。

```
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded: libpi_opencl.so
[ PluginVersion: 11.16.1 ]
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded:
libpi_level_zero.so [ PluginVersion: 11.16.1 ]
SYCL_PI_TRACE[basic]: Plugin found and successfully loaded: libpi_cuda.so
[ PluginVersion: 11.16.1 ]
SYCL_PI_TRACE[all]: Requested device_type: info::device_type::automatic
SYCL_PI_TRACE[all]: Requested device_type: info::device_type::automatic
SYCL_PI_TRACE[all]: Selected device: -> final score = 1500
SYCL_PI_TRACE[all]: platform: NVIDIA CUDA BACKEND
SYCL_PI_TRACE[all]: device: NVIDIA GeForce GTX 1050 Ti
```

3. 誤ったデバイスが選択されている場合、環境変数 ONEAPI_DEVICE_SELECTOR を使用して SYCL* デバイスセクターが選択するデバイスを変更できます。インテル® oneAPI DPC++/C++ コンパイラーのドキュメントにある「[環境変数](#)」の節を参照してください。

コードの実行中にハングアップする

次のグループ・アルゴリズムが倍精度浮動小数点数を使用する場合、`icpx` コンパイラーでコンパイルすると、CUDA* バックエンドでコードの実行がハングアップします。

- `broadcast`
- `joint_exclusive_scan`
- `joint_inclusive_scan`
- `exclusive_scan_over_group`
- `inclusive_scan_over_group`

グループ・アルゴリズムを使用する場合は、DPC++ `clang++` コンパイラー・ドライバーを使用する必要があります。

詳細は、「[ベータ版 oneAPI for AMD* GPU のインストール](#)」を参照してください。

外部参照関数「…」を解決できません/外部シンボル「…」が未定義です

これにはいくつかの原因が考えられます。

1. 現在 DPC++ では `std::complex` はサポートされていません。代わりに `sycl::complex` を使用してください。
2. DPC++ の AMD* GPU バックエンドのカーネルコードでは、`<cmath>` で宣言された C++ 標準ライブラリーの一部の数学機能 (`std::cos`、`logf`、`sinf` など) がサポートされていません。代わりに、同等の `sycl` 名前空間バージョンを使用してください。

詳細は、「[ベータ版 oneAPI for AMD* GPU のインストール](#)」を参照してください。

oneAPI for AMD* GPU 使用許諾契約書

重要 - ソフトウェアを複製、インストール、または使用する前に[使用許諾契約書 \(英語\)](#) をお読みになり、同意する必要があります。