

به نام خدا

## برنامه‌نویسی چندهسته‌ای

### دستور کار آزمایشگاه ۷



دانشگاه صنعتی امیرکبیر  
( پلی تکنیک تهران )

### مقدمه

در آزمایش‌های قبلی با هیستوگرام و انواع پیاده‌سازی آن آشنا شدید. این پیاده‌سازی‌ها بر اساس داده‌های ورودی، سطل‌ها<sup>۱</sup> خروجی و بر اساس هیستوگرام‌های محلی و روی پردازنده‌ی اصلی بوده است. همان‌طور که می‌دانید هیستوگرام یک برنامه‌ی حافظه<sup>۲</sup> محور است و نیاز به دسترسی‌های زیادی به حافظه دارد. به همین جهت بهینه‌سازی و پوشش تأخیر دسترسی‌ها به حافظه یکی از مهم‌ترین رویکردها جهت بهبود سرعت انجام کار است. از جهتی دیگر، در آزمایش قبل با اثر اشغال<sup>۳</sup> روی سرعت انجام کار آشنا شده‌اید. در این آزمایش علاوه بر آشنایی با جریان<sup>۴</sup>‌ها، با بهینه‌سازی بیشتر دسترسی‌های حافظه و اثر منفی افزایش تصرف روی سرعت انجام کارهای حافظه محور آشنا خواهید شد.

### اهداف آزمایش

- آشنایی با اثر منفی تصرف در کارهای حافظه محور
- آشنایی با حافظه‌ی پین‌شده<sup>۵</sup>
- آشنایی با جریان
- اثر تصرف در اجرای هم‌زمان جریان‌ها

جهت انجام این آزمایش فایل histogram.cu در اختیار شما قرار گرفته است. در این فایل کد به دست آوردن هیستوگرام روی پردازنده‌ی گرافیکی و قسمت موردنیاز برای پر کردن آرایه‌ی موردنظر نوشته شده است و شامل متغیرهای زیر می‌شود و استفاده‌ی هر یک به شرح زیر است

```
#define MAX_HISTOGRAM_NUMBER 10000
#define ARRAY_SIZE 102400000
#define CHUNK_SIZE 100
#define SCALER 80
#define THREAD_COUNT 512
```

- MAX\_HISTOGRAM\_NUMBER: برای مشخص کردن محدوده‌ی اعداد موجود یا به عبارت دیگر، تعیین تعداد سطل‌های هیستوگرام استفاده می‌شود
- ARRAY\_SIZE: تعداد کل عناصر مورد بررسی را مشخص می‌کند
- CHUNK\_SIZE: تعداد اعضای موردبررسی توسط یک نخ را مشخص می‌کند
- SCALER: افزایش‌دهنده‌ی میزان کار هر نخ با انجام کارهای تکراری. از این متغیر برای طولانی کردن زمان اجرای تابع هسته به اندازه‌ای که مقادیر زمانی به دست آمده قابل استناد و برای نمونه‌گیری Visual profiler طولانی باشد. برای معماری پاسکال عدد مناسب ۸۰ است.
- THREAD\_COUNT: تعداد نخ‌های هر بلوک را مشخص می‌کند

<sup>1</sup> bucket

<sup>2</sup> Memory intensive

<sup>3</sup> Occupancy

<sup>4</sup> Streams

<sup>5</sup> Pinned memory

## قسمت اول: بررسی اثر تصرف بر میزان تسريع

تابع هسته همراه کد سمت پردازنده‌ی اصلی در فایل histogram.cu موجود است. ابتدا کد را مطالعه کنید و پس از اطمینان از درک درست آزمایش با توجه به اندازه‌ی آرایه‌ی ورودی و تغییر THREAD\_COUNT و CHUNK\_SIZE تعداد بلوک‌های موردنیاز برای فراخوانی تابع را به دست آورید و جدول را کامل کنید و سؤالات زیر را پاسخ دهید. اندازه‌ی آرایه ورودی و بازه‌ی اعداد موردنظر در قسمت بالا معرفی شده است.

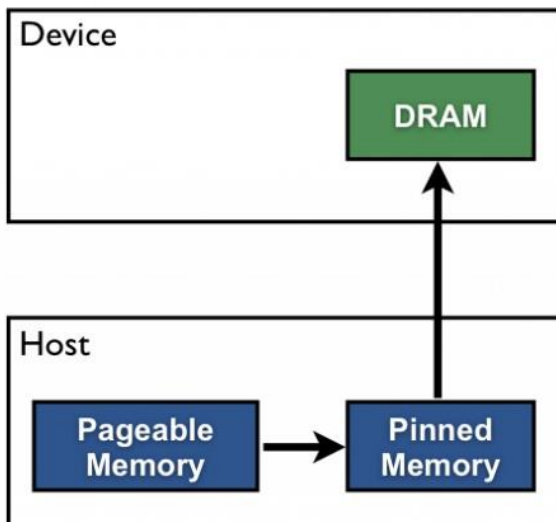
۱. دلیل اثر منفی افزایش اشغال در زمان اجرا چیست؟
۲. چگونه می‌توان با کم نکردن اشغال زمان اجرا را بهبود داد؟

تعداد نخ	۸	۱۶	۳۲	۲۵۶	۱۰۲۴
متغیر SCALER	80	80	80	80	80
تعداد بلوک					
اشغال نظری					
اشغال به‌دست‌آمده					
زمان اجرای تابع هسته					

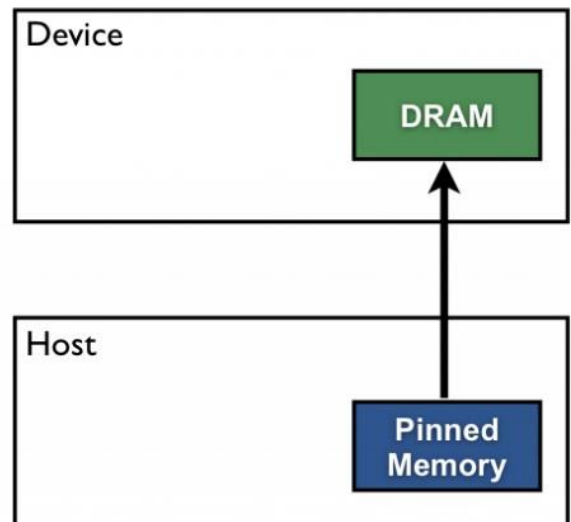
## قسمت دوم: استفاده از حافظه‌ی پین شده

در برنامه‌نویسی کودا به‌طور پیش‌فرض تخصیص حافظه برای پردازنده به‌صورت pageable است و پردازنده‌ی گرافیکی امکان دسترسی مستقیم به این نوع حافظه را ندارد. پس هنگام انتقال داده از حافظه‌ی pageable، ابتدا درایور CUDA حافظه‌ی page-locked یا پین شده را موقتاً تخصیص می‌دهد سپس آرایه‌ی موجود در حافظه‌ی pageable را به این حافظه منتقل می‌کند و سپس از این حافظه به حافظه‌ی پردازنده‌ی گرافیکی منتقل می‌کند. در شکل زیر این فرایند را می‌توانید ببینید.

### Pageable Data Transfer



### Pinned Data Transfer



شکل (۱): تفاوت حافظه پین شده و pageable

برای تخصیص حافظه به صورت پین شده، میانای برنامه‌ی کاربردی<sup>۶</sup> کودا توابعی را در اختیار برنامه‌نویس قرار می‌دهد. یکی از این توابع، تابع `cudaMallocHost` است؛ که از نظر ساختاری بسیار شبیه `cudaMalloc` است اما حافظه تخصیص یافته روی حافظه‌ی پردازنده است.

در این قسمت آزمایش، با انجام تغییرات مناسب در کد حافظه را به صورت پین شده تخصیص دهید و زمان اجرای برنامه را در حالت زیر با حالت متناظر در قسمت قبل مقایسه کنید. سپس جواب سؤالات زیر را بدهید

۱. تفاوت حافظه `Pageable` با حافظه پین شده در چیست؟
۲. روش‌های دیگر تخصیص حافظه‌ی پین شده را معرفی کنید و تفاوت آن‌ها را با یکدیگر بیان کنید.
۳. دلیل دیگر بهبود زمان اجرا در استفاده از این روش چیست؟ (به جز حذف حافظه‌ی پین شده موقت و زمان انتقال داده به آن)

## قسمت سوم: استفاده از جریان‌ها و اثر اشغال در اجرای هم‌زمان آن‌ها

در این قسمت با توجه به آشنایی با محدودیت موجود در قسمت اول، به رویکرد جدیدی برای بهبود زمان اجرا نیاز داریم. یکی از این رویکردها، استفاده از جریان‌ها است. با انجام تغییرات مناسب حلقه‌ی خارجی که مربوط به `SCALER` است را بین جریان‌های مختلف تقسیم کنید و زمان اجرای تابع هسته را به دست آورید و جدول زیر را کامل کنید. سپس به سؤالات زیر پاسخ دهید. برای مشاهده‌ی اجرای هم‌زمان جریان‌ها از `Visual Profiler` استفاده کنید.

تعداد نخ	۸	۱۶	۳۲	۲۵۶	۱۰۲۴	۸	۱۶	۳۲	۲۵۶	۱۰۲۴
تعداد جریان	۲	۲	۲	۲	۲	۴	۴	۴	۴	۴
متغیر <code>SCALER</code>	۲۰	۲۰	۲۰	۲۰	۲۰	۲۰	۲۰	۲۰	۲۰	۲۰
تعداد بلوک										
اندازه‌ی تکه										
اشغال نظری										
جریان‌های هم‌زمان										
اشغال به دست آمده										
زمان اجرای تابع هسته										

۱. به نظر شما مسائل مناسب برای استفاده از جریان‌ها چه ساختاری دارند؟
۲. پردازنده‌ی گرافیکی اولویت را به حفظ حداکثر اشغال می‌دهد یا اجرای هم‌زمان حداکثری جریان‌ها؟
۳. آیا با تغییر کامل رویکرد و اختصاص هر جریان به شمردن یک یا چند نوع عضو از آرایه‌ی اصلی می‌توان زمان اجرا را بهبود داد؟ جواب خود را با دلیل توضیح دهید. (پیاده‌سازی مورد نیاز برای این سؤال امتیازی است)

<sup>6</sup> Application Programming Interface (API)