

سوال ۱. در این تمرین لازم است تا برنامه‌ی ضرب دو ماتریس دوبعدی را با استفاده از CUDA در زبان برنامه‌نویسی C/C++ بنویسید. داده‌ها را از نوع float و ابعاد آن‌ها را  $10240 * 10240$  فرض کنید. لطفا گزارشی از آن چه کردید و تحلیل آن و همچنین سورس کد پاسخ خود را آپلود کنید.

الف) ابتدا این برنامه را به C/C++ عادی بنویسید و زمان اجرای برنامه در CPU را محاسبه کنید. در صورتی که ابعاد  $10240$  بیش از حد بزرگ بود و اجرای آن زیاد طول می‌کشید مجازید برنامه‌ی خود را در ابعاد  $1024$  اجرا کنید اما در تحلیل‌های بعدی خود این مسأله را در نظر بگیرید.

ب) محاسبات برنامه را در CUDA بازنویسی کنید. در این مرحله نیازی به استفاده از امکانات GPU مانند Shared Memory ندارید.

ج) برنامه‌ی پیشین را با استفاده از Shared Memory بازنویسی کنید. بررسی کنید که Shared Memory حدوداً در کدام ابعاد ماتریس‌ها پر می‌شود.

د) راجع به تکنیک tiling تحقیق کنید. با استفاده از tiling برنامه‌ی پیشین خود را بهبود ببخشید و توضیح بدهید که چرا این برنامه نسبت به برنامه‌ی قبلی عملکرد بهتری دارد.

ه) راجع به memory coalescing در GPU تحقیق کنید. اکنون کدام دستورات load در برنامه‌ی شما از memory coalescing استفاده می‌کنند و کدام‌ها نه؟ چه روشی برای استفاده از memory coalescing در این دستورات پیشنهاد می‌دهید؟ آن را پیاده‌سازی کنید و overhead پیاده‌سازی و همچنین بهبود عملکرد حاصل از آن را گزارش دهید.

راهنمایی: می‌توانید راجع به نمایش‌های row major order و column major order در آرایه‌های دوبعدی تحقیق کنید.

موفق باشید.