

## Лабораторные курса ТПГУ.

1. Написать код, реализующий последовательный алгоритм умножения матриц и параллельный алгоритм для вычислений на GPU. Сравнить время выполнения.
2. Изменить код ядра в программе, вычисляющей произведение матриц на GPU таким образом, чтобы генерировать ошибку обращения к памяти определенными нитями.
  - a. Обработать ошибку при выполнении ядра.
  - b. Провести отладку программы с помощью `cuda-gdb` и выявить некорректно выполняемые нити.
3. Реализовать транспонирование матрицы размерностью  $N \times K$  без использования разделяемой памяти, с разделяемой памятью без разрешения конфликта банков и с разрешением конфликта банков. С помощью метрик псу:
  - a. Определить время выполнения соответствующих ядер на GPU.
  - b. Для всех трёх случаев определить эффективность использования разделяемой памяти.
  - c. Определить для всех трех случаев пропускную способность при загрузке из глобальной памяти и при сохранении в глобальной памяти.
4. Эмулировать недостаток регистров (большой размер локальных переменных в ядре) и, используя метрики псу, определить использование локальной памяти.
5. Провести сравнительный анализ производительности программ, реализующих произведение матриц с использованием *CUDA Runtime API* и *CUDA Driver API*.<sup>\*</sup>
6. Провести сравнительный анализ производительности программ, реализующих произведение матриц с использованием *PyCuda*, *Numba cuda* и *numpy.matmul*.
7. Провести сравнительный анализ производительности программ, реализующих произведение матриц с использованием библиотеки *cuBLAS* и интерфейса *wmma*, при различных типах данных.
8. Переписать код обучения нейросети для распознавания цифр, представленный в Лекции 8, на языке C/C++ (с использованием CUDA API).
9. Разработать простой фреймворк (Python или C/C++) для обучения нейронных сетей.