# CS 6230
## Term Project
### Due 11:59pm, Friday, 12/8/2023

The term project can be done either alone or in teams of two. When grading, some consideration will be given to whether it was a solo effort or a two-person effort, with higher expectations for projects done by two-person teams.

Numerical libraries implement efficient routines for matrix-matrix multiplication (commonly called GEMM – GEneral Matrix Multiplication), where four variants are implemented:

i) **C=AB**, ii) **C=A$^T$B**, iii) **C=AB$^T$**, and iv) **C=A$^T$B$^T$**.

For the term project, these four variants of GEMM are to be optimized by you: 1) Using OpenMP for multi-core CPUs, and 2) Using CUDA for GPUs.

A primary difference between previous optimization exercises for programming assignments and the term project is that of addressing *generality*. For the programming exercises with variants of matrix multiplication, the exercises only involved the "square" case (equal values for all three size parameters for matrix multiplication) and the evaluation involved a single problem size. The term project highlights the challenge of developing high-performance numerical libraries: high performance is desired across a range of different problem sizes and the specific matrix sizes used by an application program invoking the library routine is not known a priori. Thus, adaptivity with respect to the problem size parameters in the call can be beneficial. The adaptivity may involve dynamically choosing parameters such as tile sizes, shapes/sizes of thread blocks, or even different code versions, e.g., representing different loop permutations and/or degrees of loop unrolling. If multiple code versions are used, an automatic dynamic selection mechanism must be implemented that chooses one of the versions based on problem size parameters.

The code versions to be optimized are:

1) **Base Matrix-Matrix Multiplication (C=AB):**
```
   for (i=0;i<Ni;i++)
    for (j=0;j<Nj;j++)
     for (k=0;k<Nk;k++)
   // C[i][j] += A[i][k]*B[k][j];
      C[i*Nj+j] += A[i*Nk+k]*B[k*Nj+j];
```

2) **Transposed-Nontransposed Matrix-Matrix Multiplication (C=A$^T$B):**
```
   for (i=0;i<Ni;i++)
    for (j=0;j<Nj;j++)
     for (k=0;k<Nk;k++)
   // C[i][j] += A[k][i]*B[k][j];
      C[i*Nj+j] += A[k*Ni+i]*B[k*Nj+j];
```

3) **Nontransposed-Transposed Matrix-Matrix Multiplication (C=AB$^T$):**
```
   for (i=0;i<Ni;i++)
    for (j=0;j<Nj;j++)
     for (k=0;k<Nk;k++)
   // C[i][j] += A[i][k]*B[j][k];
      C[i*Nj+j] += A[i*Nk+k]*B[j*Nk+k];
```

4) **Transposed-Transposed Matrix-Matrix Multiplication (C=A$^T$B$^T$):**
```
   for (i=0;i<Ni;i++)
    for (k=0;k<Nk;k++)
     for (j=0;j<Nj;j++)
   // C[i][j] += A[k][i]*B[j][k];
      C[i*Nj+j] += A[k*Ni+i]*B[j*Nk+k];
```

| 32 B | 0.06 | 32 KB | 6 |
| 64 B | 0.12 | 64 KB | 9 |
| 128 B | 0.24 | 128 KB | 13.5 |
| 256 B | 0.48 | 256 KB | 20.25 |
| 512 B | 0.96 | 512 KB | 30.375 |
| MAC | 0.075 | DRAM | 200 |
| Hop | 0.035 | | |

Like the programming assignments, template codes will be provided, along with a number of test cases with varying sizes of the matrices. Here are examples of problem-size variation for test cases:

| Ni | Nj | Nk |
|---|---|---|
| 8*1024 | 8*1024 | 16 |
| 4096 | 4096 | 64 |
| 2048 | 2048 | 256 |
| 1024 | 1024 | 1024 |
| 256 | 256 | 16*1024 |
| 64 | 64 | 256*1024 |
| 16 | 16 | 4*1024*1024 |

| Ni | Nj | Nk |
|---|---|---|
| 9*999 | 9*999 | 37 |
| 3*999 | 3*999 | 111 |
| 999 | 999 | 999 |
| 333 | 333 | 9*999 |
| 111 | 111 | 81*999 |
| 37 | 37 | 729*999 |

more advanced technologies. Also, many of our observations in Section 6 are technology-independent.

To capture data reuse opportunity in an arch, we follow the approach used in [30]:

$$e_i = \prod_{j=i}^{L} R_j$$

Here $e_i$ is the per-access energy for level with level total number of data reuse differently. The number of times the data is accessed from the lower-cost (child) level during its lifetime in this level. To support direct inter-PE communication in systolic arrays, we treat neighbor PEs as an additional level in the hierarchy. We distinguish the cost for different communication distances (Figure 3) which is an improvement over [8].

Since all designs today are power constrained, finding the optimal accelerator design now becomes an optimization problem of minimizing $E$ over the 3D design space, similar to Yang's work [44]. $e_i$ is determined by the resource allocation (Table 3), and $RT_i$ can be directly calculated from the dataflow and loop blocking schemes. In this simply perform a conservatively pruned search design space guided by domain-specific knowledge. An artifact model is available at https://github.com/Interstellar-CNN-scheduler.

**Framework validation:** We have thoroughly our model by comparing its results generated by our synthesis tools simple designs we have generated. Figure 7 shows the energy consumption model and post-synthesis average less than 2%. Furthermore, we are able to reproduce the results from [30] operations each cycle. However, it needs to store the intermediate results, which used flow in wires from one logic unit to the next, so it can access them when needed. Thus, the size of the memory needed is related to the complexity of the communication in the algorithm, so the energy increases with communication complexity.

# 6   Results

Using our dataflow taxonomy and the ability to generate and evaluate large numbers of accelerator Halide, we first explore different dataflow and loop choices, and then consider hardware resource constraints. At the end we leverage the characteristics of the memory to introduce an efficient optimizer for DNN accelerators many factors, to first order it grows as the square root of its size, which roughly corresponds to the length of the wires that need to transport the address and data values across the memory array. This was noted long ago by Amrutur and Horowitz [8], citing even earlier work by Evans and Franzon [9]. The memory energy also depends on the fetch width, but this dependence is much weaker than you might expect. For example, moving from 16- to 64-b fetches only changes the energy by 1.5x, so wider fetches are generally more efficient in terms of energy per byte. This means that for a 16-b machine, a fetch from eve

---

[2] Energy for memory and integer ops come from Verilog, placed and routed using commercial tools. Energy for floating-point ops come from the Galal thesis which also used data from placed and routed designs.

[3] Internally, most SRAMs fetch 64–256 b on each access, so returning a small number of bits increases the effective energy cost per bit. To address this issue, you could create a SIMD machine and fetch the 16-b data for four lanes from a single SRAM. While this is more efficient, it also makes the memory four times larger, since it now needs to hold four lanes' worth of working set, so the benefit is modest.