ECE 6775 High-Level Digital Design Automation Fall 2024

Midterm Review HLS Design Practice



Cornell University



Announcements

- Lab 4 released
 - Group with your teammate on CMS by Monday
- Midterm on Tue 10/22 at 11:40am
 - In class, 75 mins
 - Open book, open notes, closed Internet
 - Please arrive 5 mins earlier
- Links to past quizzes are posted on Ed (thread #1)
- Another paper reading session on Tuesday 10/29
 - C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. Cong,
 <u>"Optimizing FPGA-based Accelerator Design for Deep</u> <u>Convolutional Neural Networks</u>", FPGA 2015

Agenda

- Midterm review
 - Overview of key lecture topics
 - Homework question discussion
- HLS design practice
 - Lab design review
 - Additional case studies

Midterm next Tuesday (20%)

- ► Topics covered: lectures 01~11, 13, 14
 - Hardware specialization
 - Algorithm basics
 - FPGA
 - C-based synthesis
 - Control flow graph and SSA
 - Scheduling
 - Resource sharing
 - Pipelining

Key Topics (1)

- Algorithm basics
 - Time complexity, esp. big-O notation
 - Graphs
 - Trees, DAGs, topological sort
 - BDDs, timing analysis
- FPGAs
 - LUTs and LUT mapping
- C-based synthesis
 - Arbitrary precision and fixed-point types
 - Key HLS optimizations to improve design performance

Key Topics (2)

- Control data flow graph
 - Dominance relation
 - Loops
 - SSA
- Scheduling
 - TCS & RCS algorithms: ILP, list scheduling, SDC
 - Operation chaining, frequency/latency/resource constraints
 - Ability to devise a simple scheduling algorithm to optimize a design metric

Key Topics (3)

- Resource sharing
 - Conflict and compatibility graphs
 - Ability to determine minimum resource usage in # of functional units and/or registers, given a fixed schedule
- Pipelining
 - Dependence types
 - Ability to determine minimum II given a code snippet
 - Modulo scheduling concepts: MII, RecMII, ResMII

HW1 Q5

Q5. Implement the circuit shown below using a **minimum** number of 3-input lookup tables (LUTs). Please do NOT simplify the gate-level circuit before mapping it to LUTs. Indicate your answer by creating a table similar to Table 1. Please add a new row per LUT, and use the signal names from the circuit to label the inputs (i.e., the associated cut) and output of each LUT. Mark an unused LUT input with a "-".



HW2 Q3

Q3. With the code snippet listed below, what is the minimum achievable II if we pipeline the loop? Please (1) draw the dependence graph for the loop, and (2) calculate both ResMII and RecMII.

```
for (int i = 2; i < N; i++) {
    #pragma pipeline II=?
    mem[i] += (mem[i-2] >> i);
}
```

In this problem, we assume that all operations take a full cycle to execute and no combinational chaining is allowed. We also assume that the 'mem' array will be mapped to a single-ported memory block (i.e., one read/write port) without any write-through support (i.e., a memory read must happen one cycle after a write).



HW2 Q5

Q5. To form an integer linear programming (ILP) model for the constrained scheduling problem, we have learned to use a binary decision variable X(i,k) to indicate if operation i starts at cycle k, where $1 \le k \le L$ with L being the maximum schedule latency.

In this problem, we will make use of these schedule variables to test if the output value of a given operation v is live at a specific cycle s $(1 \le s \le L)$. Here we make the following assumptions:



Output value of V_2 is only live at cycle 3

Figure 3: A scheduled graph with three operations.

Can you introduce a **derived** binary variable Y(v, s) to indicate if an operation v is (\mathbf{a}) scheduled before cycle s (excluding s)? Hint: Make use of the relevant X's.

 (\mathbf{b}) Can you introduce another derived binary variable Z(v,s) to indicate if an operation v is scheduled at or after cycle s? Hint: Make use of Y.

 (\mathbf{c}) Finally, can you make use of the variables introduced from (a) and (b) to derive binary variable R(v, s), which indicates if the output value of v is live at cycle s. Hint: You may need to introduce additional linear constraint(s) between Y, Z, and/or R. You may also make use of USE[v] based on the assumptions stated.

HW2 Q6

Q6. Given a chain of N operations $\{O_1, O_2, ..., O_N\}$ (N > 1) where each operation O_i is associated with a positive integral delay value of d_i $(d_i \le 10ns)$, our goal is to automatically place (or insert) one or multiple registers into this chain to minimize the clock period of the circuit (i.e., maximize the operating frequency).



A (Simplified) Loopy Form of Processor Pipeline

```
for ( pc++ ) {
    inst = fetch( pc )
    { rs, rt, rd } = decode( inst )
    result = execute( reg[rs], reg[rt] )
    mem( result, reg[rt] )
    reg[rd] = writeback( result )
}
```

Case Study: CORDIC



Case Study: Digit Recognition

- Use a simple machine learning algorithm to recognize handwritten digits
 - 2000 training instances per digit
 - Each training/test instance is a 7x7 bitmap after downsampling



MNIST dataset: http://yann.lecun.com/exdb/mnist/

K-Nearest-Neighbor (KNN) Implementation

```
bit4 digitrec( digit input )
Ł
 #include "training_data.h"
  // This array stores K minimum distances per training set
  bit6 knn_set[10][K_CONST];
  // Initialize the knn set
  for ( int i = 0; i < 10; ++i )
    for ( int k = 0; k < K_CONST; ++k )
      // Note that the max distance is 49
      knn_set[i][k] = 50;
```

Main compute loop

```
L2000: for ( int i = 0; i < TRAINING SIZE; ++i ) {
  L10: for (int j = 0; j < 10; j++) {
    // Read a new instance from the training set
    digit training instance = training data[j * TRAINING SIZE + i];
    // Update the KNN set
    update_knn( input, training_instance, knn_set[j] );
 }
}
```

Assuming 10 cycles per innermost loop (L10) ~200K cycles by default without optimizations

10x Speedup through Parallel Processing

```
bit4 digitrec( digit input )
ł
 #include "training_data.h"
  // This array stores K minimum distances per training set
  bit6 knn_set[10][K_CONST];
  // Initialize the knn set
  for ( int i = 0; i < 10; ++i )
    for ( int k = 0; k < K CONST; ++k )
      // Note that the max distance is 49
      knn set[i][k] = 50;
 Unroll inner loop completely
  L2000: for ( int i = 0; i < TRAINING_SIZE; ++i )
                                                    {
Partition training
   L10: for ( int j = 0; j < 10; j++ ) {
                                                     set into 10 banks
      // Read a new instance from the training set
      digit training_instance = training_data[j * TRAINING_SIZE + i];
      // Update the KNN set
      update_knn( input, training_instance, knn_set[j] );
   }
  }
               10 instances of "update_knn" running in parallel
```

~20K cycles after parallelization

Further Speedup through Pipelining

```
bit4 digitrec( digit input )
{
 #include "training_data.h"
  // This array stores K minimum distances per training set
  bit6 knn_set[10][K_CONST];
  // Initialize the knn set
  for ( int i = 0; i < 10; ++i )
    for ( int k = 0; k < K CONST; ++k )
      // Note that the max distance is 49
      knn set[i][k] = 50;
                                            Pipeline outer loop
  L2000: for ( int i = 0; i < TRAINING_SIZE; ++i ) {
    L10: for ( int j = 0; j < 10; j++ ) {
      // Read a new instance from the training set
      digit training instance = training data[j * TRAINING SIZE + i];
      // Update the KNN set
      update_knn( input, training_instance, knn_set[j] );
   }
  }
                    Outer loop (L2000) pipelined to II=1
                       ~2K cycles after pipelining
                                                                   16
```

Array Partitioning Caveats

Example: Array partitioning through *constant indices* (after unrolling)



Array Partitioning Caveats

Example: Array partitioning through *constant indices* (after unrolling)



Case Study: Revisiting 3x3 Convolution

```
for (r = 1; r < H; r++)
  for (c = 1; c < W; c++) {
    #pragma HLS pipeline II=?
    for (i = 0; i < 3; i++)
      for (j = 0; j < 3; j++)
        out += img[r+i-1][c+j-1] * f[i][j];
      out[r][c] = out;
    }
}</pre>
```

- Inner loops "i" & j are automatically unrolled
- The 3x3 filter array "f" is automatically partitioned into 9 registers after unrolling
- The entire input image "img" is stored in an on-chip buffer with two read ports

ResMII = ? RecMII = ?

Achieving II=1 for 3x3 Convolution using a Line Buffer and Shift Registers





- New pixel read from frame buffer in main memory (DRAM)
- 2. Update line buffer by removing the oldest pixel and shifting in the new one

Line Buffer + Shift Registers:

a custom "cache" + a custom "register file"



Resulting Specialized Memory Hierarchy

Memory architecture customized for convolution



HLS Code Snippet

```
LineBuffer<2,C,pixel_t> linebuf;
 1
 \mathbf{2}
     Window<3,3,pixel_t> window;
     for (int r = 1; r < R+1; r++) {
 3
       for (int c = 1; c < C+1; c++) {
 4
         #pragma HLS pipeline II=1
 5
         pixel_t new_pixel = img[r][c];
6
         // Update shift window
7
8
         window.shift_left();
9
         if (r < R \&\& c < C) {
           for (int i = 0; i < 2; i++ )
10
             window.insert(buf[i][c]);
11
         }
12
13
         else { // zero padding
           for (int i = 0; i < 2; i++)
14
             window.insert(0);
15
         }
16
         window.insert(new_pixel);
17
18
         // Update line buffer
         linebuf.shift_up(c);
19
20
         if (r < R \&\& c < C)
           linebuf[1].insert(c, new_pixel);
21
         else // Zero padding
22
23
           linebuf[1].insert(c, 0);
         // Perform 3x3 convolution
24
         out[r-1][c-1] = convolve(window, weights);
25
26
       }
     }
27
```

Next Class

Midterm exam