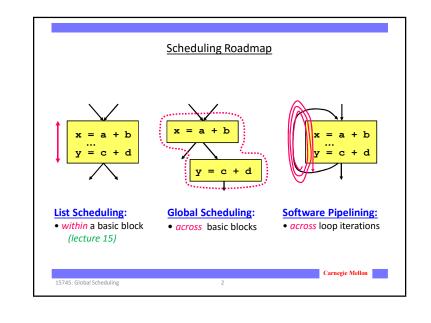
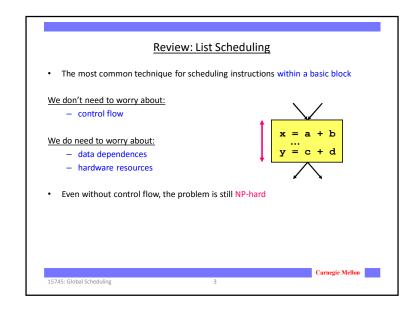
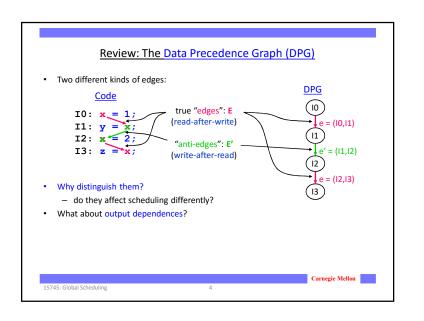
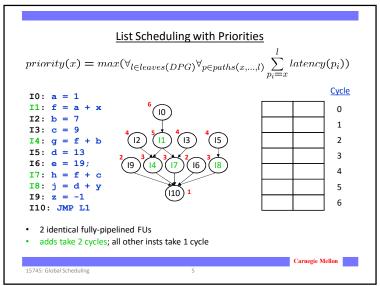
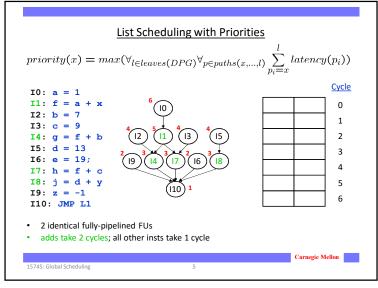
#### Lecture 20 Global Scheduling & Software Pipelining [ALSU 10.4-10.5] Phillip B. Gibbons 15745: Global Scheduling & Software Pipelining 1

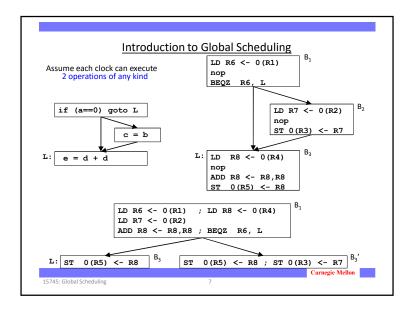


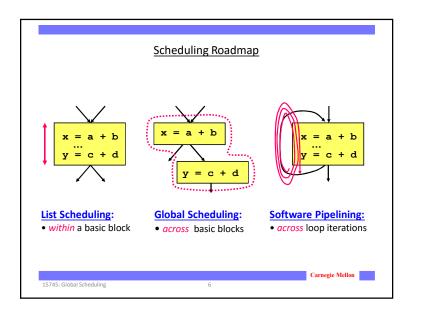


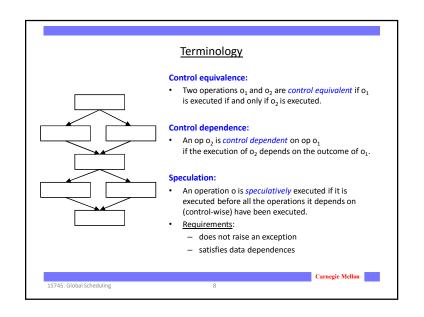


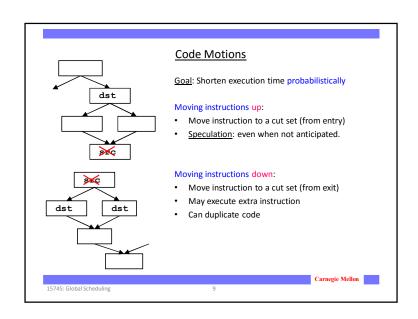




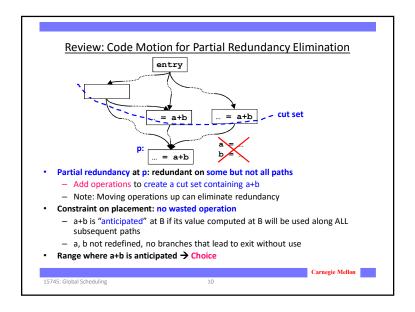


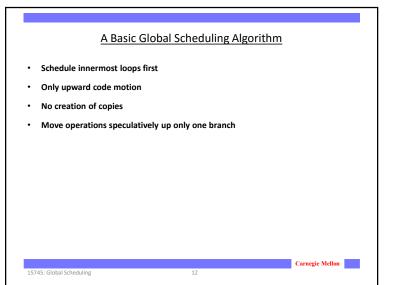






# General-Purpose Applications • Lots of data dependences • Key performance factor: memory latencies • Move memory fetches up - Speculative memory fetches can be expensive • Control-intensive: get execution profile - Static estimation • Innermost loops are frequently executed - back edges are likely to be taken • Edges that branch to exit and exception routines are not likely to be taken - Dynamic profiling • Instrument code and measure using representative data

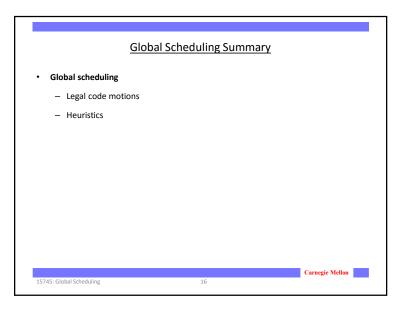


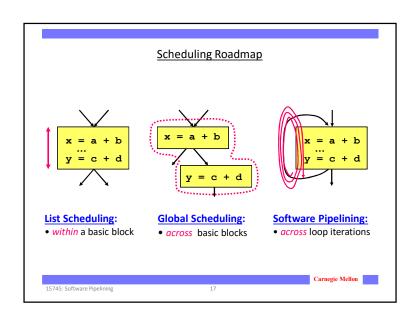


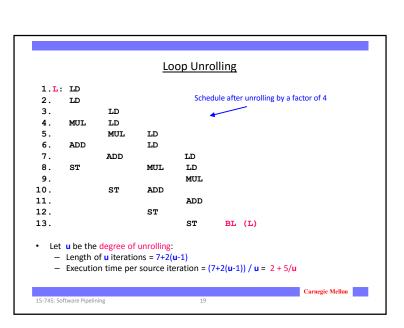
## Program Representation • Recall: A region in a control flow graph is: — a set of basic blocks and all the edges connecting these blocks, — such that control from outside the region must enter through a single entry block • A procedure is represented as a hierarchy of regions — The whole control flow graph is a region — Each natural loop (single entry with back edge to it) in the flow graph is a region — Natural loops are hierarchically nested • Schedule regions from inner to outer — treat inner loop as a black box unit • can schedule around it but not into it — ignore all the loop back edges → get an acyclic graph

## Extensions • Prepass before scheduling: loop unrolling • Especially important to move operation up loop back edges Carnegie Mellon

#### Algorithm Compute data dependences; For each region from inner to outer { For each basic block B in prioritized topological order { CandBlocks = ControlEquiv{B} ∪ Dominated-Successors{ControlEquiv{B}}; CandInsts = ready operations in CandBlocks; For (t = 0, 1, ... until all operations from B are scheduled) { For (n in CandInst in priority order) { if (n has no resource conflicts at time t) { S(n) = < B, t > Update resource commitments Update data dependences Update CandInsts; Priority functions: non-speculative before speculative Carnegie Mellon 15745: Global Scheduling







```
Example of DoAll Loops

    Machine:

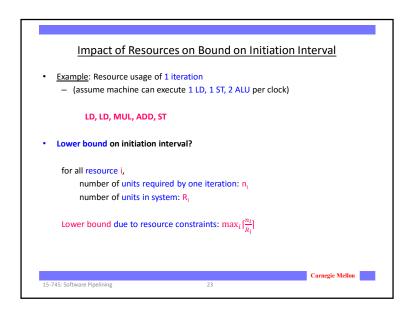
    - Per clock: 1 read, 1 write, 1 (2-stage) arithmetic op, with hardware loop op
      and auto-incrementing addressing mode.
· Source code:
        For i = 1 to n
             D[i] = A[i]*B[i] + c
· Code for one iteration:
      1. LD R5,0(R1++)
      2. LD R6,0(R2++)
      3. MUL R7,R5,R6
      5. ADD R8,R7,R4
      7. ST 0(R3++),R8
· Little or no parallelism within basic block
                                                           Carnegie Mellon
15-745: Software Pipelining
```

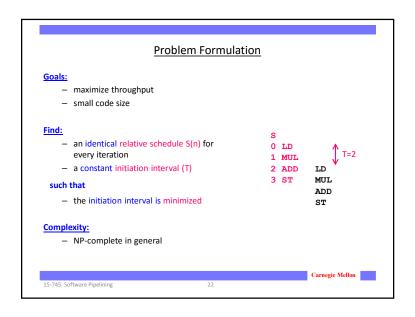
```
Software Pipelined Code
       LD
 1.
       LD
 2.
 3.
       MUL
                LD
 4.
                LD
 5.
                MUL
                        LD
 6.
       ADD
                         LD
 7. L:
                         MUL
                                 LD
                                          BL (L)
      ST
                ADD
                                 LD
 8.
 9.
                                 MUL
10.
                ST
                        ADD
11.
12.
                         ST
                                 ADD
13.
                                 ST
• Unlike unrolling, software pipelining can give optimal result

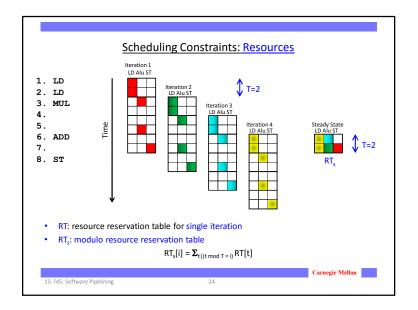
    Locally compacted code may not be globally optimal

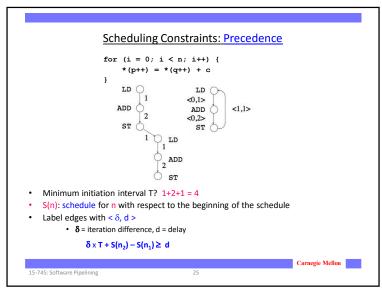
• DOALL: Can fill arbitrarily long pipelines with infinitely many iterations
                                                             Carnegie Mellon
15-745: Software Pipelining
```

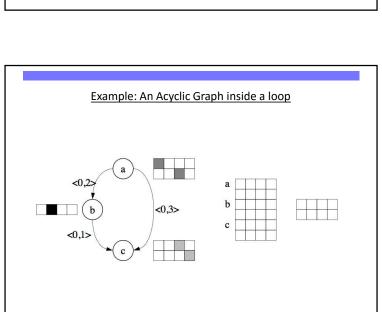
```
Example of DoAcross Loop
Loop:
                                          1. LD // A[i]
     Sum = Sum + A[i];
                                          2. MUL // A[i]*c
     B[i] = A[i] * c;
                                          3. ADD // Sum += A[i]
                                          4. ST // B[i]
Software Pipelined Code
  1. LD
   2. MUL
   3. ADD
             LD
   4. ST
              MUL
              ADD
              ST
Doacross loops
· Recurrences can be parallelized
· Harder to fully utilize hardware with large degrees of parallelism
15-745: Software Pipelining
```











15-745: Software Pipelining

Carnegie Mellon

