

Processing Elements

①

Road map:

- models of computation (system spec.)
 - KPN, SDF, State Charts
- performance estimation (system design and validation)
 - SystemC
- system architecture (system design)
 - PEs, Memory, On-chip interconnect

PE taxonomy

- ① IP instruction processors
 - execute program according to the ISA
- ② HW accelerators
 - execute one (or more) fixed functions
- ③ HW Coprocessors
 - extends ISA

Instruction Processors Types

① General Purpose

- designed for a variety of uses
- measured against benchmarks

e.g. SPEC CPU2017 (sequential) www.spec.org

e.g. PARSEC (parallel) cs.princeton.edu

e.g. EEMBC (embedded) www.eembc.org

② Application-Specific

- DSP: digital signal processors
- GPU: image processing / numerical computing

Instruction Processor Performance

- improve performance by increasing:
 - clock rate
 - instruction parallelism

① increasing clock rate:

option i) increase voltage

- clock frequency is roughly linear w.r.t. voltage

- dynamic (switching) power: $P = CV^2f$

↳ of transistors & wires

⇒ switching power $\propto f^3$

- challenges: distributing power, removing heat

option ii) deeper pipeline

- deeper pipeline \Rightarrow shorter stages \Rightarrow faster clock without increasing voltage
- need good branch prediction and out-of-order execution \Rightarrow consume power
- pipeline depths

ARM 7 (1993)	ARM 11 (2002)	Cortex A9 (2007)	Cortex A15 (2010)	Cortex A77 (2019)
3	8	8	15	13
Pentium (1993)	Pentium 3 (1999)	Pentium 4	Core (2006)	Atom (2008)
5	10	20-31	14-19	12-14

↳ hit the "power wall"

② increasing instruction parallelism

option i) superscalar datapath

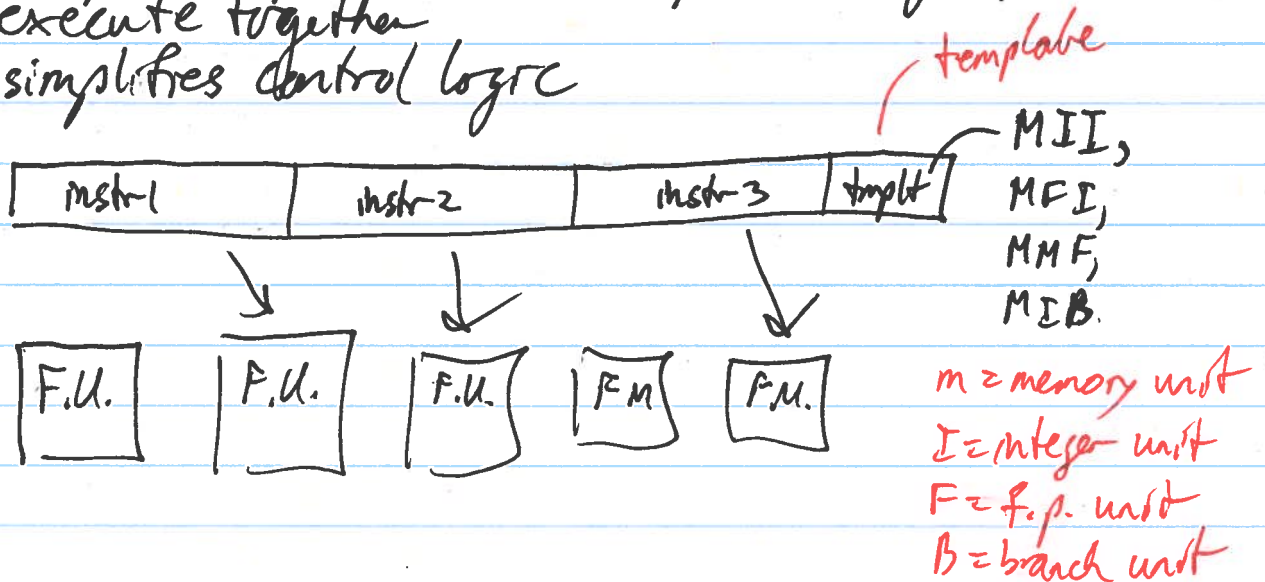
- fetch, decode, dispatch, execute, commit 2/4/6/8 instr. per cycle
- extracts instruction-level parallelism (ILP) from a thread of execution (transparent to the software)
- decode widths:

ARM 11	Cortex A9	Cortex A15	Cortex A77
1 (scalar)	2	3	4
Pentium	Core	Apple M1	IBM Power
2	4	8	8

(4)

option ii) VLIW

- very long instruction word
- compiler specifies ILP by bundling operations to execute together
- simplifies control logic



- VLIW processors:

DSPs (TMS320)

Itanium (Intel IA-64 arch.)

Nvidia Bluefield-2 accelerators

Qualcomm Hexagon DSP

option iii) parallel programming

- simultaneous multi-threading (thread-level parallelism)
- single core has state for 2+ threads (registers, PC, status register)
- the threads share the functional units
- better utilization of resources
- not popular in embedded systems - thread performance depends on other threads on the core
- multiprocessing (program-level parallelism)
 - use multiple PEs
 - memory/interconnect can be a bottleneck

(5)

DSP Digital Signal Processors

- processors optimized for sequential math operations on streaming data

e.g. TI TMS320 (1983-present)

- fixed-point and floating-point variants

e.g. Qualcomm Hexagon (2009-present)

- VLIW

- was fixed-point, now floating-point

- in Qualcomm Snapdragon SoC (CPU, DSP, GPU, GPS, cellular radio, NFC, image processor)

- Features:

- multiply-accumulate:

MAC $a \quad R_i, R_j \quad ; \quad R_a \leftarrow R_a + R_i * R_j$

- circular address registers

- increment modulo defined buffer size

- single loop instruction

RPT i - loop next instruction i times

6

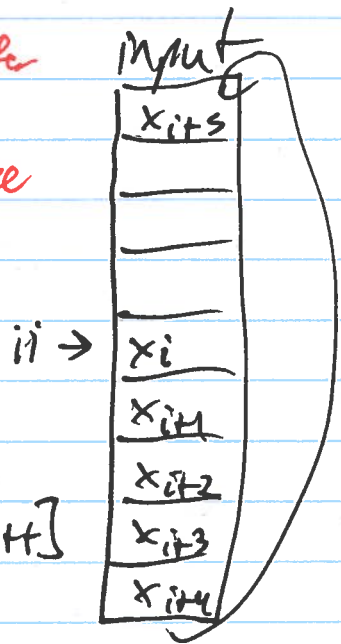
float *inputs = SOME_BUFFER, coeff[LARGE_VALUE] = CONSTANTS;
res = 0;

int ii = START_VALUE; *initial index into circular input buffer*
for (int c=0; c < LARGE_VALUE; c++)
res += input[(ii+c)%LARGE_VALUE] * coeff[c];

index to input { MOV Rca0, STARTVALUE
coeff { MOV Rcs0, LARGE_VALUE
accum { MOV Rca1, 0
MOV Rcb1, LARGE_VALUE
MOV Ra, 0

circular
circular buffer index
circ. buff. size

RPT LARGE_VALUE
MACa inputs[Rca0++], coeff[Rca1++]
STR res, Ra



! instr. loop
vs multiple instructions on
a general-purpose processor

⑦

GP GPU: General Purpose GPU Programming

top500.org

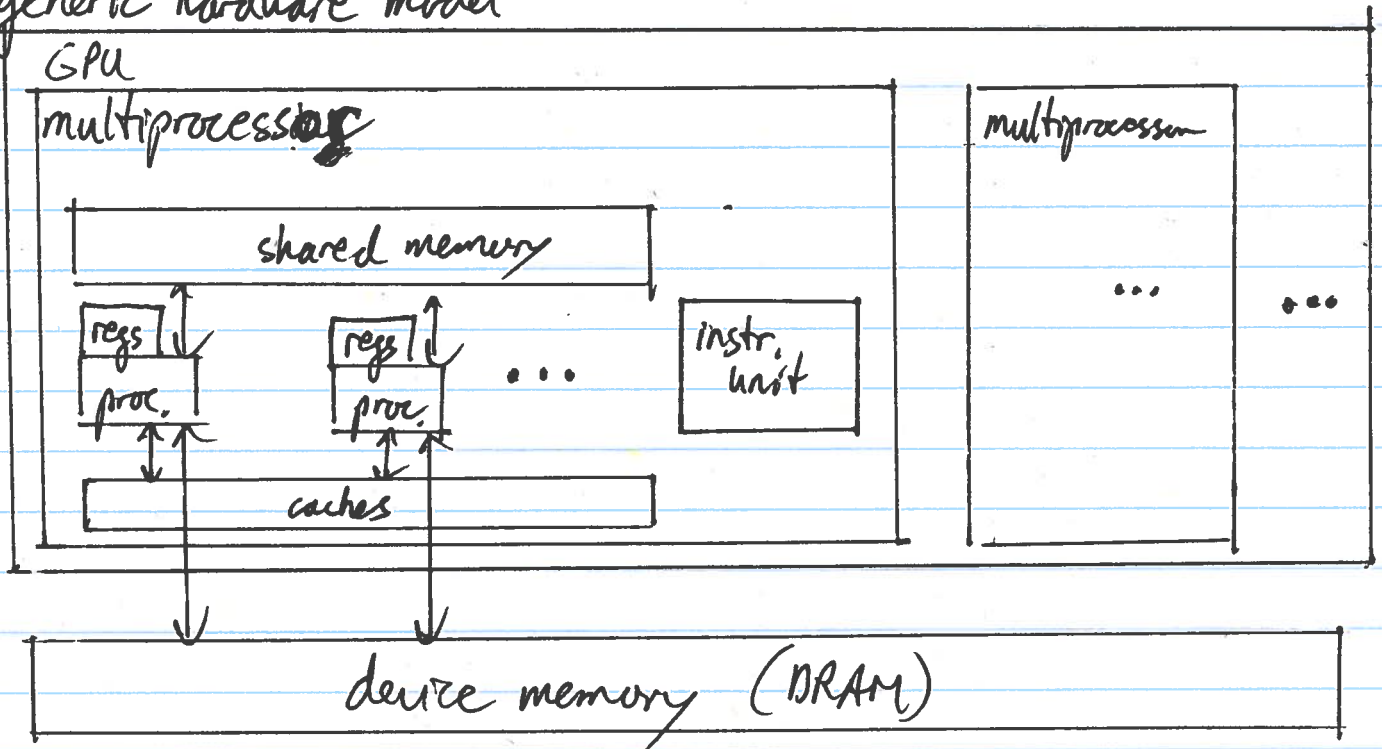
- GPUs can do thousands of floating-point operations in parallel

- gpgpu frameworks:


CUDA - proprietary (Nvidia)

OpenCL - open (AMD, IBM, Intel, Nvidia, Qualcomm, ...)

- generic hardware model

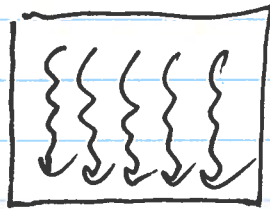


- programming model

- threads 

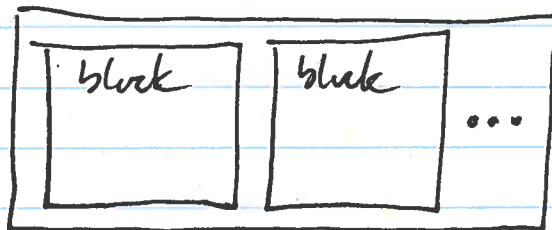
- executes an instance of a kernel
- unique thread id

- blocks



- set of threads executing the same code
- unique block ids
- blocks can be multi-dimensional (x, y, z)

- grid



- array of blocks running same kernel
- can be multi-dimensional (x, y, z)

- execution model

- blocks get assigned to multiprocessors
- threads are executed in warps of (32, 64) executing the instructions in lock step
- they can use predication for conditional execution

e.g. `var = 0;`
`if(condition)`
`var = 1;`
`else`
`var = 2;`

```
mov.s32 var, 0
setp    p, condition
@p mov.s32 var, 1
@!p mov.s32 var, 2
```

Handouts > cuda Saxpy

CUDA SAXPY Example

Source: <https://devblogs.nvidia.com/easy-introduction-cuda-c-and-c/>

SAXPY: Single-precision A * X Plus Y

```
#include <stdio.h>

__global__
void saxpy(int n, float a, float *x, float *y)
{
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    if (i < n) y[i] = a*x[i] + y[i];
}

int main(void)
{
    int N = 1<<20;
    float *x, *y, *d_x, *d_y;
    x = (float*)malloc(N*sizeof(float));
    y = (float*)malloc(N*sizeof(float));

    cudaMalloc(&d_x, N*sizeof(float));
    cudaMalloc(&d_y, N*sizeof(float));

    for (int i = 0; i < N; i++) {
        x[i] = 1.0f;
        y[i] = 2.0f;
    }

    cudaMemcpy(d_x, x, N*sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(d_y, y, N*sizeof(float), cudaMemcpyHostToDevice);

    // Perform SAXPY on 1M elements
    saxpy<<<(N+255)/256, 256>>>(N, 2.0f, d_x, d_y);

    cudaMemcpy(y, d_y, N*sizeof(float), cudaMemcpyDeviceToHost);

    float maxError = 0.0f;
    for (int i = 0; i < N; i++)
        maxError = max(maxError, abs(y[i]-4.0f));
    printf("Max error: %f\n", maxError);

    cudaFree(d_x);
    cudaFree(d_y);
    free(x);
    free(y);
}
```

Handwritten annotations:

- array size* (pointing to `n`)
- coefficient* (pointing to `a`)
- "kernel"* (pointing to the `saxpy` function)
- 1M* (pointing to `N = 1<<20`)
- in CPU's memory* (pointing to `x` and `y` mallocs)
- device memory* (pointing to `cudaMalloc` calls)
- initialize in CPU memory* (pointing to the initialization loop)
- copy from CPU to GPU (DMA)* (pointing to `cudaMemcpyHostToDevice` calls)
- copy result from GPU to CPU* (pointing to `cudaMemcpyDeviceToHost`)
- free device mem.* (pointing to `cudaFree` calls)
- free CPU memory* (pointing to `free` calls)

Additional notes:

- # blocks = $\lceil N / \text{block size} \rceil$* (pointing to `(N+255)/256`)
- block size* (pointing to `256`)

ASIP: Application-Specific Instruction-Set Processor

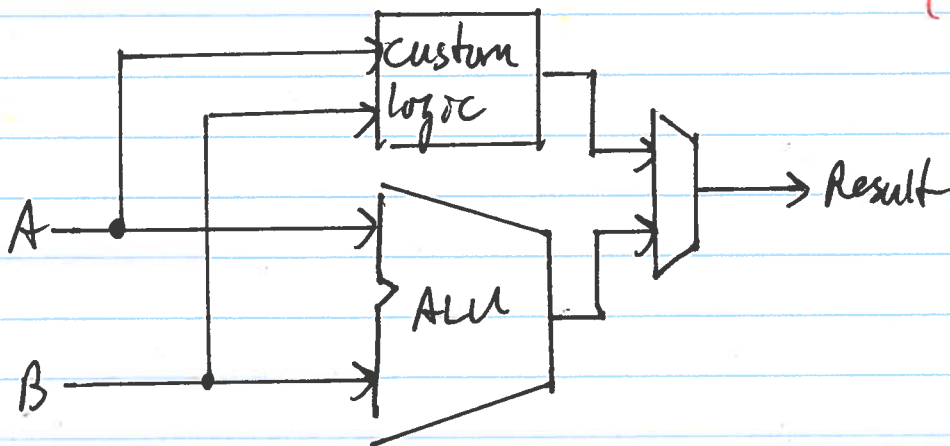
Types:

- configurable architecture
 - eg. vary #/type of functional units
 - eg. configure memory system (eg. Nios II - MMU is optional, Caches are optional)
- ISA optimization (Instruction Set Architecture)
 - eg. add custom instructions e.g. CRC instructions
 - eg. remove unused ~~functions~~ instructions
- tools generate RTL and configure compilers

Extensible Processors

- Nios II (user-defined instructions)

(soft-core CPU)



- up to 256 custom instructions
 - software macros are used to invoke it
 - eg. `int y = ALT_CI_CRC(x);`
- ↑
integer

- Xtensa (synthesized instructions)

e.g.

```
; use generator
; generator:: regfile(name => "L",
                      cname => "long128",
                      sname => "s",
                      width => 128);
```

- specifies a
16 x 128-bit
reg. file

```
opcode add128
add128 { assign sr = st + ss};
```

- given the above, the processor generator generates RTL, testbench, etc.

- the RTL is automatically pipelined
- also generates processor pipeline interlocks (hazard detection and stalling) and result bypass logic (data forwarding)

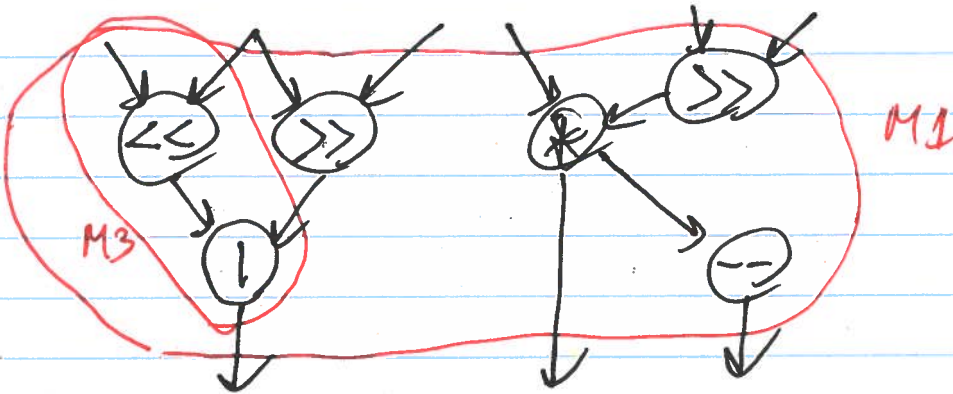
- to use the above in C:

```
int main() {
    long128 src1[N], src2[N], dst[N];
    for (int i=0; i<N; i++)
        dst[i] = add128(src1[i], src2[i]);
}
```

Instruction Set Synthesis

- ① extract CFG for the application
- ② analyze the DFG of each basic block and extract "cuts" (synthesized instructions)
- ③ estimate area and clock cycles for each cut
- ④ evaluate application speedup
- ⑤ pick the best set of cuts based on constraints (area)

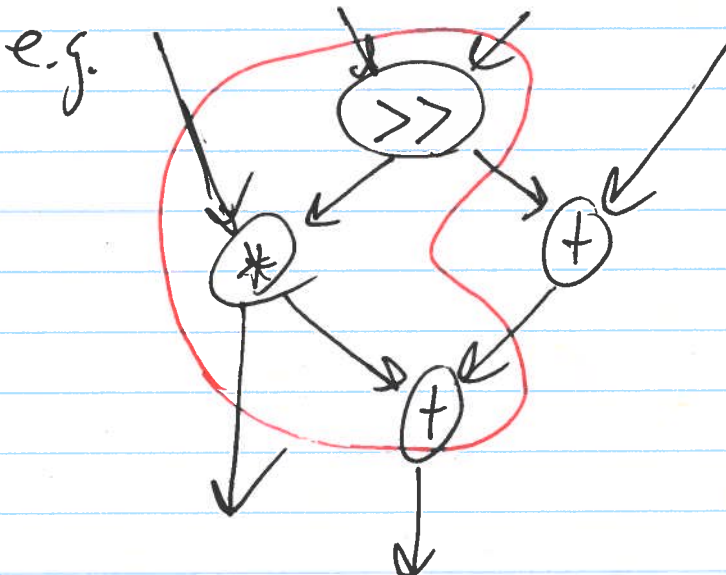
e.g. basic block



M3 : 3 inputs, 1 output

M1: 6 inputs, 3 outputs

might not be feasible based on reg file ports



- illegal: can't be done in one instruction

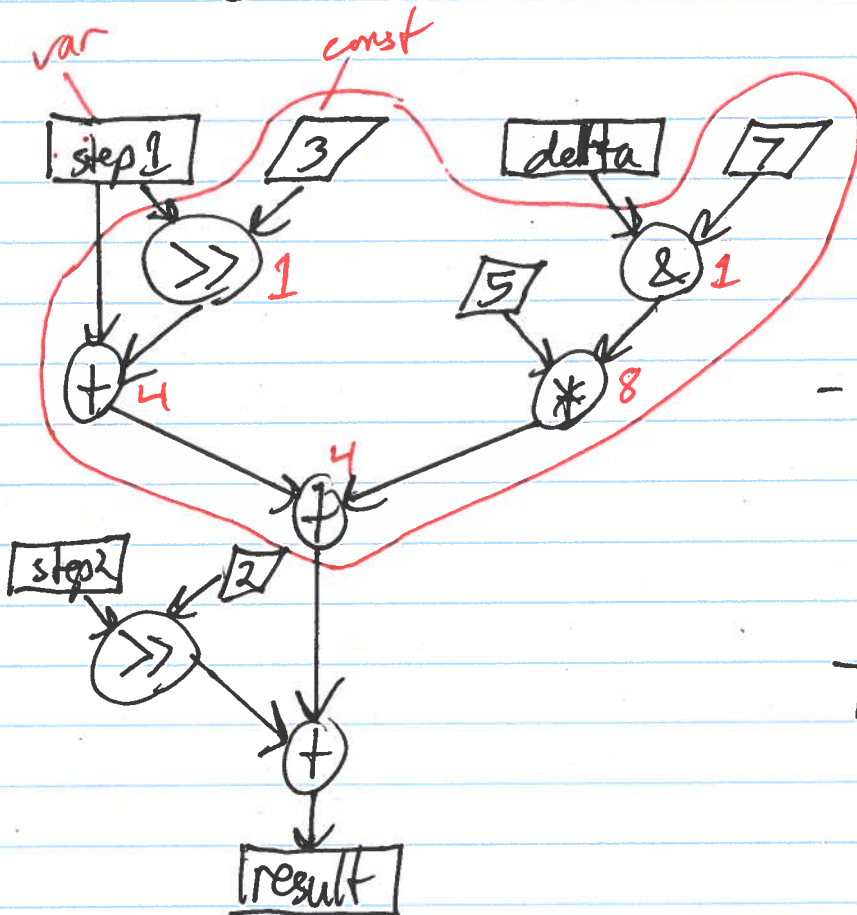
- speedup of a cut, C , can be computed as:

$$S = \frac{T_{\text{tot}}}{T_{\text{tot}} - n_c (T_{\text{sw}}^c - T_{\text{hw}}^c)}$$

n_c times cut
is executed

difference in time
of the cut

e.g. "PROC" block



← largest cut with
2 inputs / 1 output?

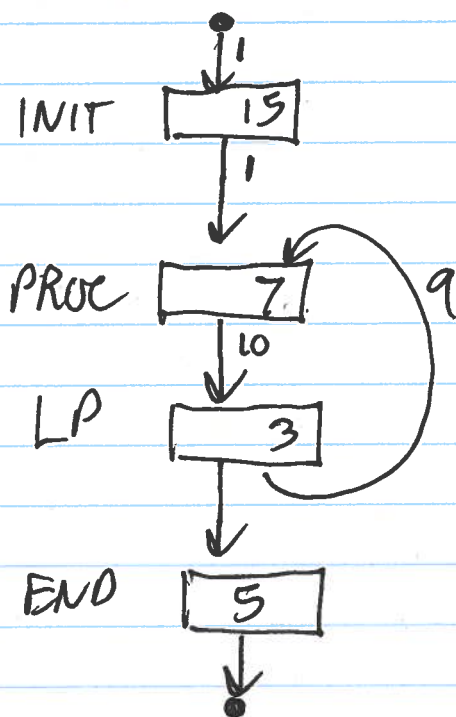
T_{sw}^c :

- 5 operations/instructions
- processor clock 100MHz
- $5 \times 10\text{ns} = \underline{50\text{ns}}$

T_{hw}^c :

- logic = 1 ns
- adder = 4 ns
- multiplier = 8 ns
- critical path = 13 ns
- $\lceil 13\text{ns} / 10\text{ns} \rceil = 2 \text{ cycles}$
- $2 \times 10\text{ns} = \underline{20\text{ns}}$

CFG (application)



$$T_{\text{tot}} = 15 + (7+3) * 10 + 5 = 120 \text{ cycles}$$

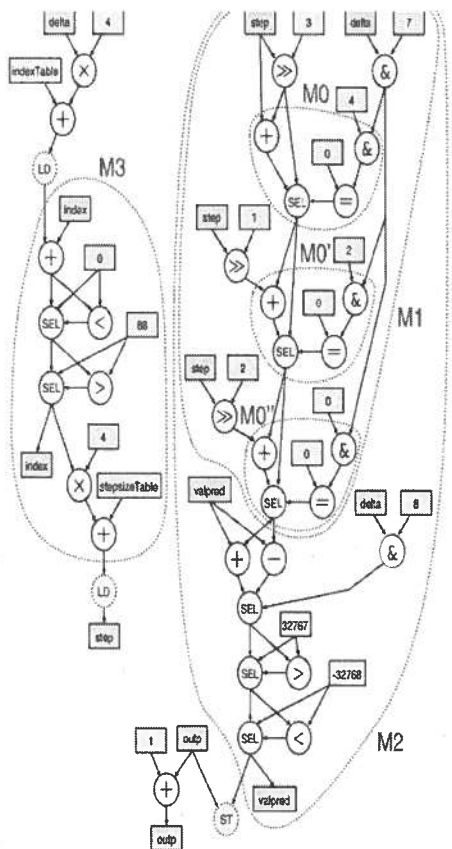
$$n_c = 10$$

$$T_{\text{sw}}^c = 5, T_{\text{hw}}^c = 2$$

$$S = \frac{120 \text{ cycles}}{(120 - 10(5-2))} = \frac{120}{90} \approx 1.33$$

- choose a set of cuts that maximizes total speedup while staying within area constraints and register file limitations

Handouts > instr Synth



adpcm example

- SEL is a selector node for conditional execution
- can execute M3 and M2 in parallel
- M0 is executed multiple times but it is small and have a lot of inputs
- M1 is a 16x3-bit multiply with 2 inputs/1 output
- depending on hardware M1, M2, or M3 provide better speedup

Atasu, Pozzi, lenne. Automatic Application-Specific Instruction-Set Extensions under Microarchitectural Constraints, Proceedings 40th Design Automation Conference, 2003.

Input/Output Multiplexing

- if #inputs/outputs > register file ports
- use mux/demux – takes multiple cycles

e.g. 4 inputs, 2 outputs (reg file 2 read ports, 1 write port)

