

# Out of Order Execution

Idea  $\Rightarrow$  independent first, dependent last  $\Rightarrow$  Rest areas for dep. instr

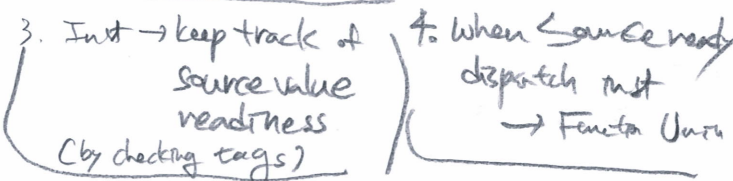
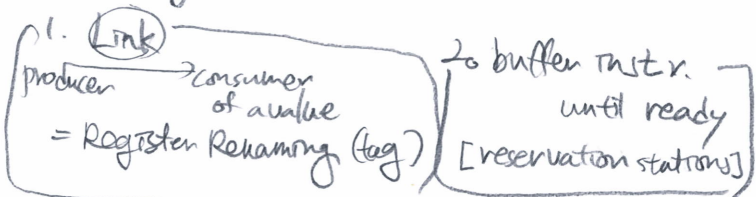
• Monitor the source "values"  $\Rightarrow$  Reservation stations in the resting area

• when all source "values" are available  $\Rightarrow$  "fire" the instruction.

$\therefore$  Instr. dispatched in dataflow order

In order dispatch + precise exceptions vs out of order dispatch + precise exceptn.

## Enabling OOO Execution



## • Tomasulo's Algorithm

## • Two Hump in a Modern Pipeline



# Lecture 7: Out of Order Dataflow + Superscalar Execution.

10/10

\* Most modern OOO Execution w/ precise execution

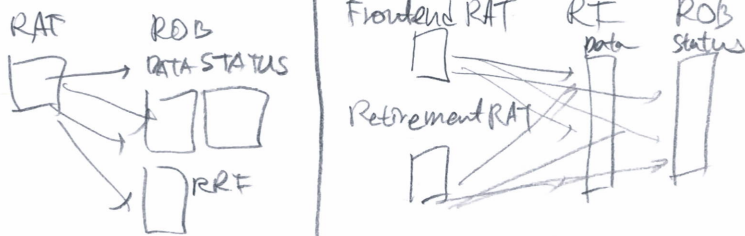
- $\hookrightarrow$  ① Reorder buffer to support in-order retirement
- ② A single <sup>physical</sup> register file to store all register values
- ③ 2 register maps

Future / Front reg. map  $\rightarrow$  used for renaming

Architectural reg. map  $\rightarrow$  used for maintaining precise state.

$\Rightarrow$  OOO Ex = "Restricted Dataflow" (limited to instruction window)

Pentium III | Pentium IV

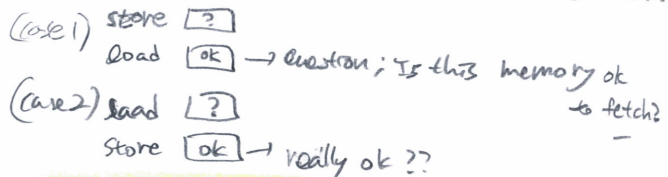


## Q. Memory Dependence Handling?

$\Rightarrow$  Memaddr is not known until a load/store executes

$\hookrightarrow$  renaming?  $\rightarrow$  difficult

- Determining dep?  $\rightarrow$  after (partial) execution
- When a load/store has its address ready, there may be younger/older loads/stores with undetermined addresses in the machine



$\Rightarrow$  "memory disambiguation" or "unknown address problem"

- Solution:
- ① Stall
  - ② Aggressive; Just do it! Assume load is independent
  - ③ Intelligent; Predict if the load is dependent on the/any unknown address store

How to detect?

Option 1 Wait until all previous stores committed.

Option 2 Keep a list of pending stores in a "store buffer" & check whether "load address" matches a prev. "store address"



"Store buffer" = list of stores that are pending in the machine

load address

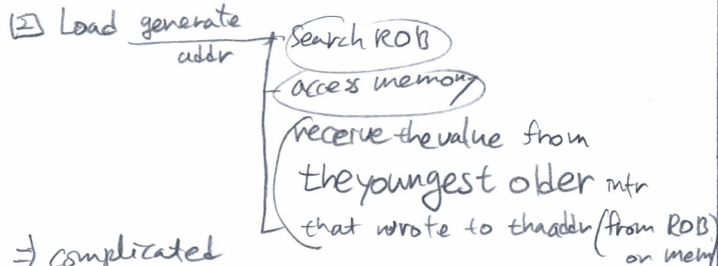
How to schedule? ⇒ Option 1 Assume dependent → simple & slow  
 Option 2 " independent → simple  
 Option 3 Predict → more accurate still need recovery & need recovery

### \* Data Forwarding Between Store & Load

↳ Modern processors use LQ (load queue) & SQ (Store queue) age-based comparison !!

### \* Out of Order Completion of Memops

Store execution finish → ROB



### \* Store-Load Forwarding Complexity

Option 1 CAM  
 Option 2 Range Search  
 Option 3 Age-based Search  
 also Option 4 Load data from [SQ mem/cache]

## Other Approaches to Concurrency (MDC 15)

(on Instruction Level Parallelism)

- 1 Super-scalar execution
- 2 VLIW
- 3 Fine-grain multi-threading
- 4 SIMD Processing Vector & array processors, GPUs
- 5 Decoupled Access Execute
- 6 Systolic Arrays

### \* Superscalar Execution

Idea: Fetch/Decode/Execute/Retire multiple instr per cycle.  
 • N-wide superscalar → N instr per cycle.

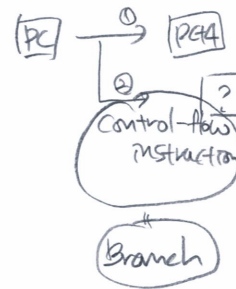
• In-Order superscalar processor

- Copies of datapath.
- Dependency make it tricky.

⇒ Adv: High IPC  
 Disadv: Dependency, more hardware.

## Lecture 18. Branch Prediction

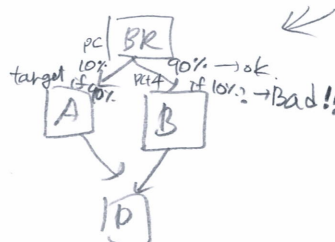
### • Control Dependence Handling



Branch types	Direction at fetch time	# of possible next PC	
Conditional	unknown	2	E
Unconditional	Always taken	1	D
Call	"	1	D
Return	"	many	E
Indirect	"	many	E

### \* Branch Prediction \*

if (X = 0) { A } else { B } ~ { D }



if (pointer != NULL) { A } else { B }  
 more probable!  
 → programmer needs to try maximize probability.



## ◦ Guessing Next PC+4

Idea: Get rid of control flow instr.

- ① Predicate combining
- ② Predicated execution.

## ◦ Branch Prediction

↳ fetch instr == Branch?

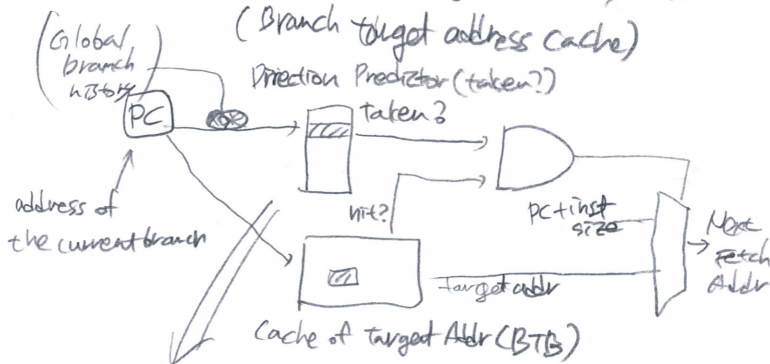
↳ Conditional Branch?

↳ Branch target addr?

→ store the target addr. (if taken)

from prev. instance & access it with the PC.

◦ Called Branch Target Buffer, BTB



So, How can we predict the branch? direction ..

- Compile time (static)
    - Always Not Taken
  - Run time (dynamic)
    - Always Taken
    - BTBN (Back Taken Forward Not Taken)
    - Profile based
    - Program analysis based
    - Hybrid (Dynamic + Static)
    - Advanced algorithms (perceptrons)
- pragmas → #(intel) if(unbrdpr)

## Intel Pentium Pro Branch Predictor

⇒ Two level global branch predictor

- ↳ 4-bit global history register
- ↳ Multiple pattern history tables (of 26bit counters)

## Lecture 19: Branch Prediction

MOULS

+ VLIW + Fine-Grained Multithreading

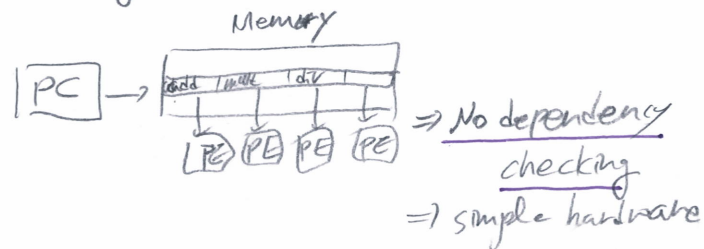
Other control dependence handling

- ↳ Stall
- ↳ Branch Prediction
- ↳ Branch Delay Slot
- ↳ Fine-grained multithreading
- ↳ predicated execution
- ↳ multipath execution

**VLIW** very large instruction word

↳ superscalar; multiple instructions & check dependencies between them.

Software packs independent instructions in a larger "instruction bundle"



(But complex compiler)

Q. what about variable latency operation?

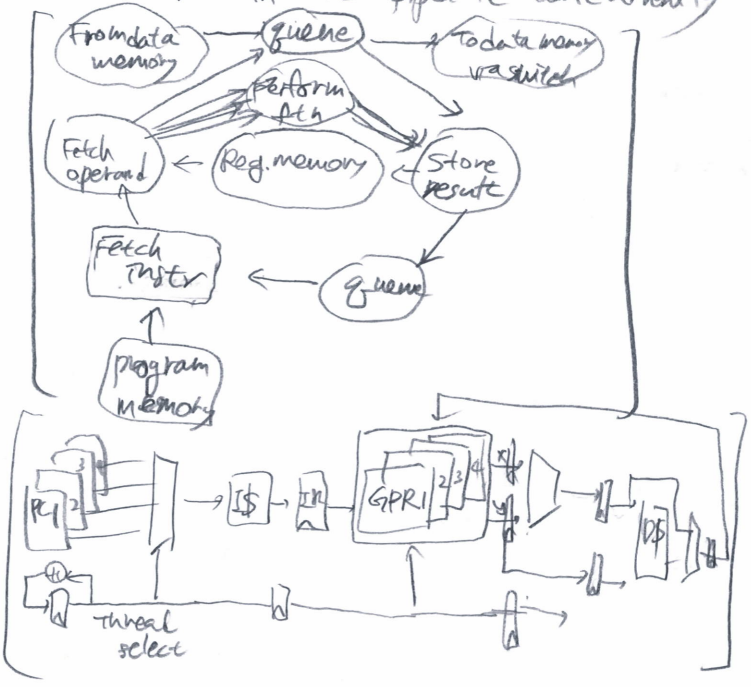
- VLIW philosophy
  - ⇒ RISC, Intel IA-64, Superblock
  - ✓ Impact.

- Trade-off - Adv: No need for dynamic scheduling hardware
  - (dependency checking instr. align)
- Disadv: Compiler more complex !!

# [Fine-Grained Multithreading.]

idea 1 Fetch engine fetches from a diff. thread with hardware having multiple thread contexts

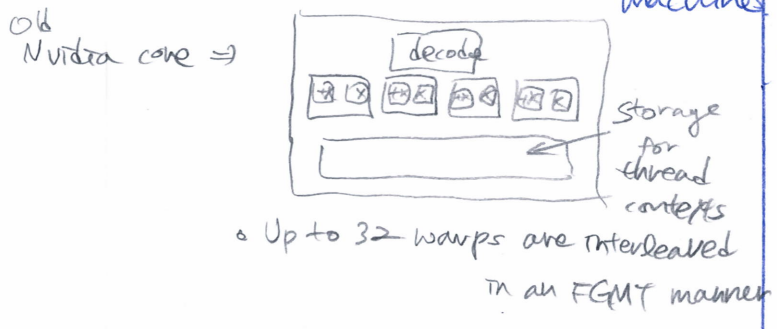
idea 2: Switch to another thread every cycle s.t no two instr from a thread are in the pipeline concurrently



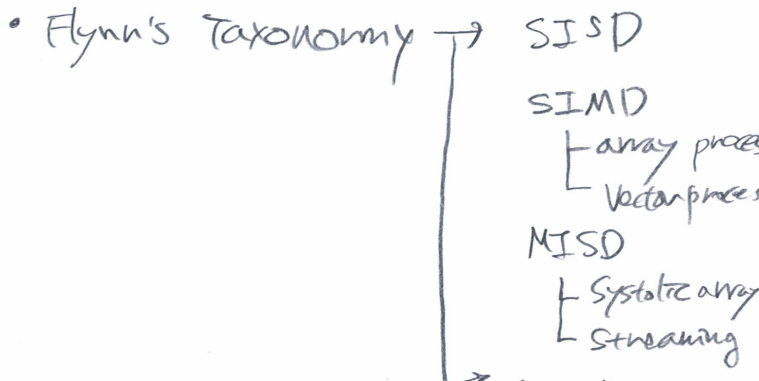
• Sun Niagara ; The first multicore machine

- Adv: No need for dep check, branch prediction
- Otherwise bubble cycles used for useful instr. from different threads
- Disadv: Extra hardware complexity
- Reduced single thread performance

• Modern GPU's are fine grained multicore machines



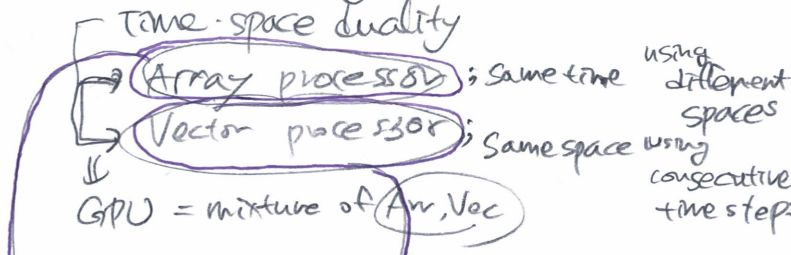
# Lecture 20: SIMD Processors



• Data parallelism.

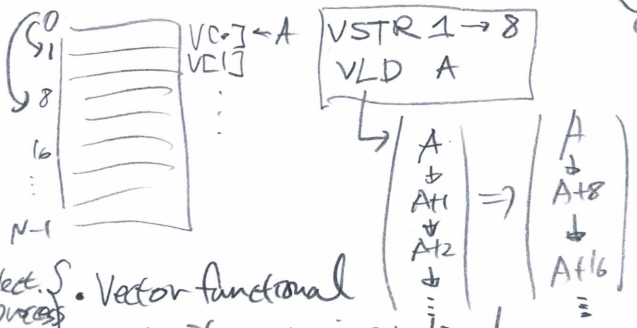
ex) dot product of two vectors  
 Same operation on different pieces of data  
 (≠ Data flow)

• SIMD processing



multiple processing element.  
 (easy..)

- sequential execution by
- 1 Vector registers
  - 2 Vector length register (VLEN)
  - 3 Vector stride register (VSTR)



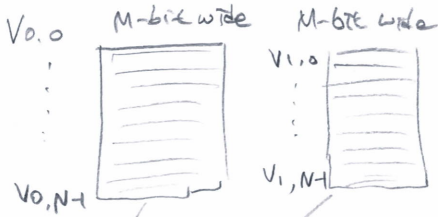
- Vector functional units are pipelined
- Each pipeline stage operates on diff data element
- Deeper pipeline ⇒ No intra-vector dep.
- Must be regular
  - No control flow.
  - Known stride → easy addr calc ⇒ prefetching



# \* Vector Processor Limitation

↳ Data should be regular!  
 very inefficient if irregular

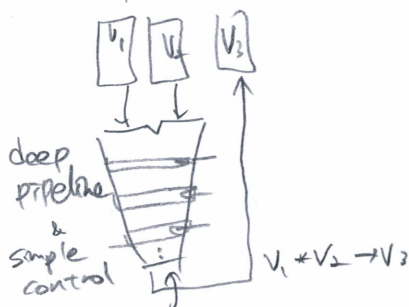
- Memory (Bandwidth) can be a bottleneck
- Vector registers



- vector control register = VLEN, VSTR, VMASK

$$VMASK[i] = (VLEN[i] == 0)$$

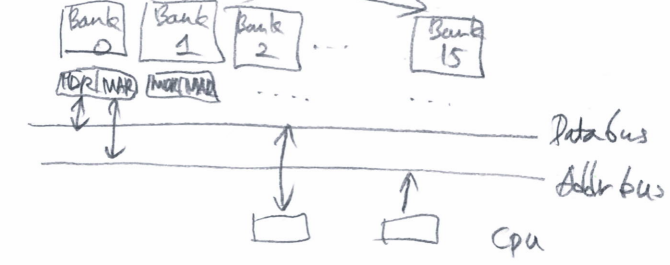
## • Vector Functional Units



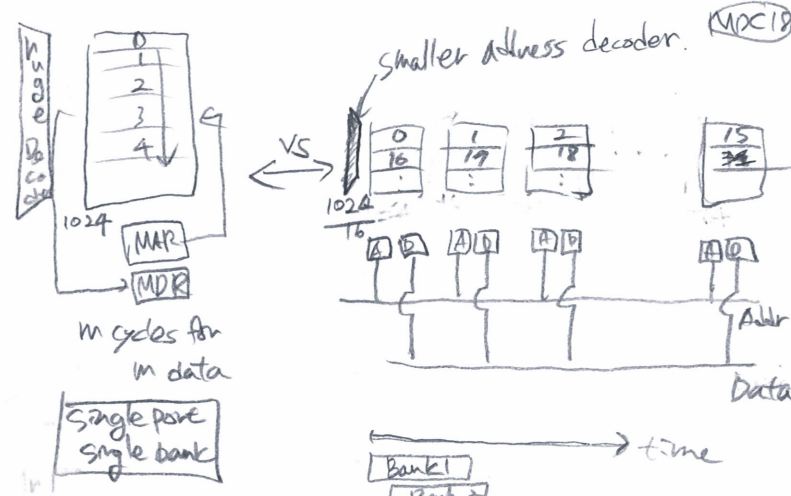
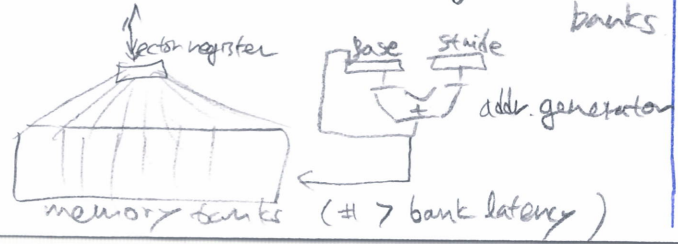
6 stage multiply pipeline.

• memory banks → loading/storing vectors from/to memory

Multiple (16) ⇒ 1 element for 1 cycle!  
 ⇒ Need multiple mem. banks!!

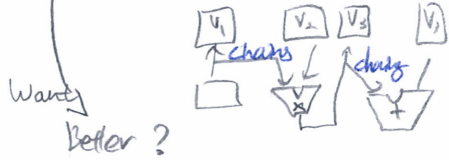


⇒ sustain N parallel access if all N go to different banks



• Vectorizable Loop. (depends on the stride)

• Vector chaining; data forwarding... from one vector functional unit to another



Vector chaining w/ 2 ports per bank

QI: what if # of data element > # elements in a vector register?

AI: Break loops? ↳ Vector Stripmining

QII: If irregular data? AI: Use indirection to combine/pack elements into vector register

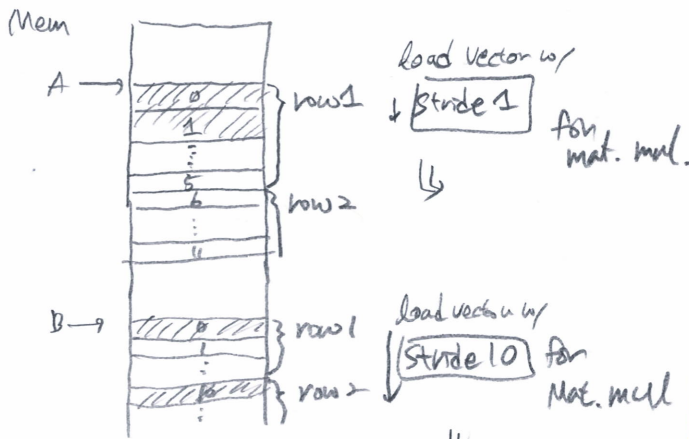
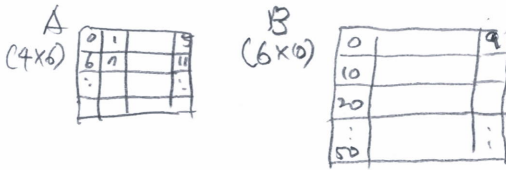
↳ Scatter/Gather operators

QIII: If Conditional Operations in a loop (Forbidden Operation for some element)

AI: Masked operation (VMASK) → "predicated execution"

# Lect 21 SIMD & GPUs

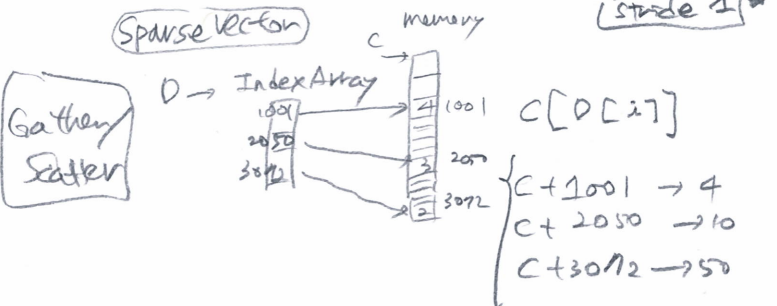
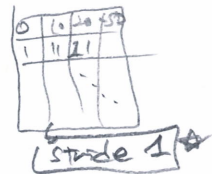
- Stride & banking issues (should be relatively prime)
- Storage of a matrix.
  - Row major: store data row by row
  - Column major: store data column by column.
- Matrix multiply  $A \times B$



Different strides can lead to "bank conflicts" to "bank conflicts"

- 1) more banks
- 2) Better Data layout  $\rightarrow$  natural  $A \times B^T$  (transpose) (i.e., column major)
- 3) Better mapping of address to bank.

E.g.) randomized mapping



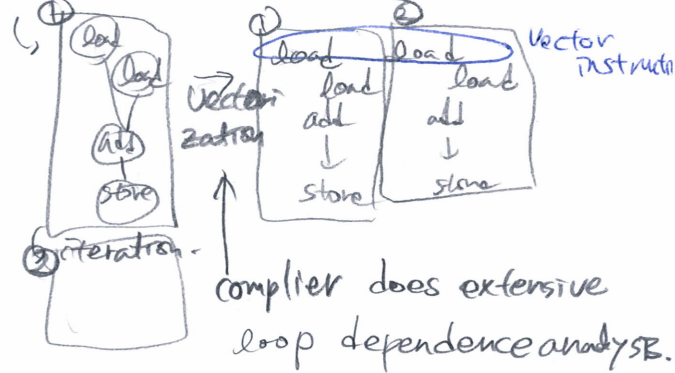
- Modern SIMD = mixture of array processor & vector processor.

Ex) Vector Instr Level Parallelism.



```
for (i=0; i<N; i++)
```

```
  C[i] = A[i] + B[i];
```



- SIMD ISA Extensions (graphics)
- Intel Pentium MMX Operations.
  - $\hookrightarrow$  Image Overlaying (w/ bit mask)

## GPU (Graphic Processing Units)

- programming is done using threads, NOT SIMD instructions.
- Programming Model vs. Hardware Execution Model.

$\hookrightarrow$  next page!

Q. How to calculate

```
for (i=0; i<N; i++)
  C[i] = A[i] + B[i].
```