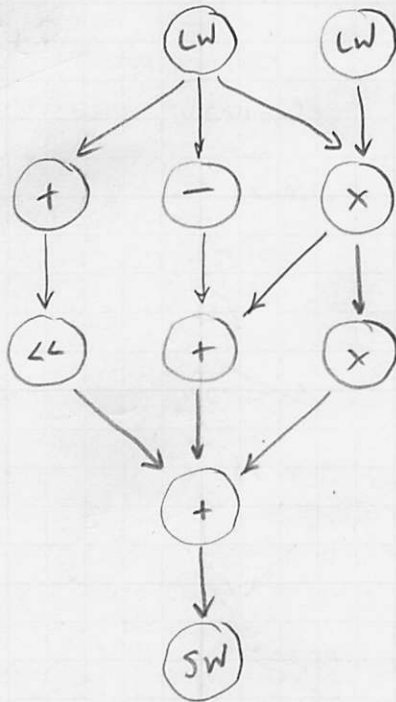


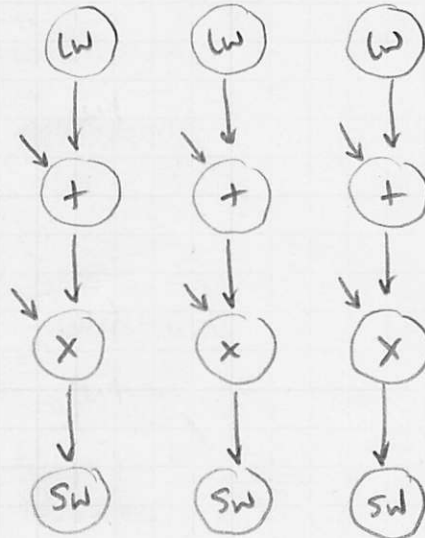
# ILP vs DLP

ILP = INSTRUCTION LEVEL PARALLELISM

DLP = DATA LEVEL PARALLELISM



ILP



DLP

OOO AND VLIW EXPLOIT ILP WHEN ARBITRARY INSTRUCTIONS CAN BE EXECUTED IN PARALLEL

VECTOR PROC CAN ALSO EXPLOIT ILP TO A LIMITED EXTENT, IT IS JUST MORE EFFECTIVE TO EXPLOIT ARBITRARY ILP ON OOO AND VLIW PROC

VECTOR EXPLOITS DLP WHEN THE SAME INSTRUCTION OPERATING ON DIFFERENT DATA IS EXECUTED IN PARALLEL

OOO AND VLIW CAN DEFINITELY STILL EXPLOIT DLP, IT IS JUST MORE EFFICIENT TO EXPLOIT DLP ON VECTOR PROC

# SUBWORD - SIMD

ADD A NEW INSTRUCTION THAT DOES FOUR 8b ADDITIONS IN PARALLEL

• `ADDU.qb rd, rt, rs`

$$\begin{aligned}
 R[rd][7:0] &= R[rt][7:0] + R[rs][7:0] \\
 R[rd][15:8] &= R[rt][15:8] + R[rs][15:8] \\
 R[rd][23:16] &= R[rt][23:16] + R[rs][23:16] \\
 R[rd][31:24] &= R[rt][31:24] + R[rs][31:24]
 \end{aligned}$$

SEMANTICS

for ( $i=0; i < n; i++$ )

$C[i] = A[i] + B[i]$

← ASSUME A, B, C ARE ARRAYS OF BYTES NOT ARRAYS OF WORDS

## SCALAR ASSEMBLY

```

loop:
  lbu  r1, 0(r2)
  lbu  r3, 0(r4)
  ADDU r5, r1, r3
  sb   r5, 0(r6)
  ADDIU r2, r2, 1
  ADDIU r4, r4, 1
  ADDIU r6, r6, 1
  ADDIU r7, r7, -1
  bne  r7, r0, loop
    
```

LOAD/STORE BYTE

Pointer bump by 1 BYTE

## SUBWORD-SIMD ASSEMBLY

```

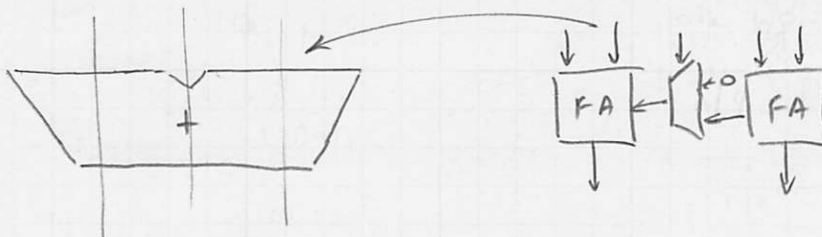
loop:
  lw   r1, 0(r2)
  lw   r3, 0(r4)
  ADDU.qb r5, r1, r3
  sw   r5, 0(r6)
  ADDIU r2, r2, 4
  ADDIU r4, r4, 4
  ADDIU r6, r6, 4
  ADDIU r7, r7, -1
  bne  r7, r0, loop
    
```

LOAD/STORE WORD

ONLY n/4 iterations

Four ADDS in parallel

## SUBWORD - SIMD MICROARCHITECTURE



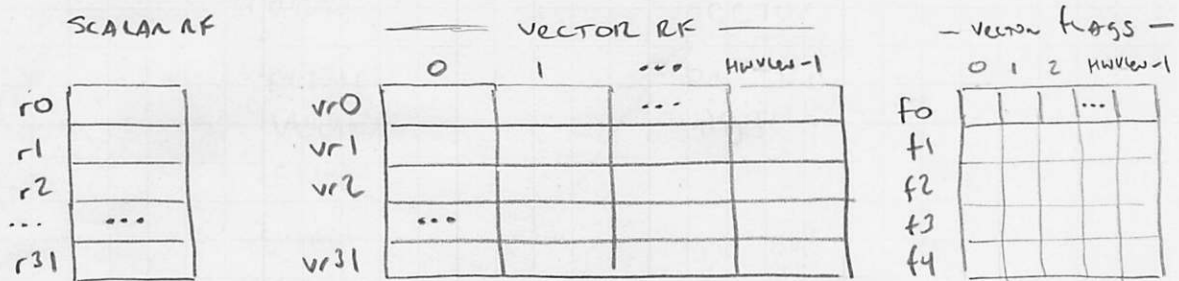
Simply Break CARRY CHAIN

CAN ALSO DO SUBWORD MULTIPLY AND FLOATING POINT

## VECTOR INSTRUCTION SET ARCHITECTURE

- ADD A NEW "VECTOR REGISTER FILE" THAT HOLDS VECTOR INSTEAD OF SCALAR VALUES
- ADD NEW "VECTOR INSTRUCTIONS" THAT OPERATE ON THE VECTOR REGISTERS
- NUMBER OF ELEMENTS IN A VECTOR REGISTER IS CALLED THE "HARDWARE VECTOR LENGTH"

### VECTOR ARCHITECTURAL STATE



### VECTOR INSTRUCTIONS

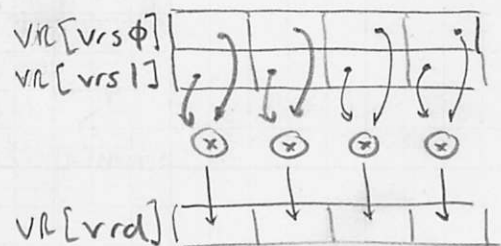
$\swarrow$  VLR   
 Number of Active elements

- VECTOR-VECTOR MULTIPLY

MUL.VV vrd, vrs0, vrs1

for  $i = 0$  to  $VLR-1$

$$VR[vrd][i] = VR[vrs0][i] \times VR[vrs1][i]$$

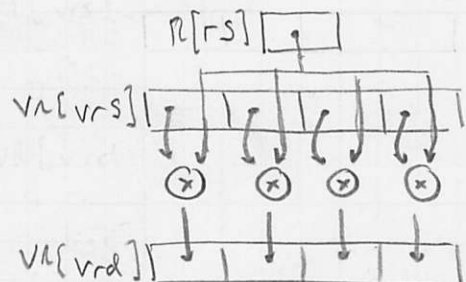


- VECTOR-SCALAR MULTIPLY

MUL.VS vrd, vrs, rs

for  $i = 0$  to  $VLR-1$

$$VR[vrd][i] = VR[vrs][i] \times R[rs]$$

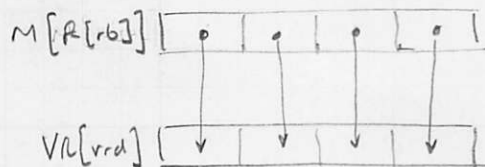


- VECTOR UNIT-STRADE LOAD

LW.V vrd, rb

for  $i=0$  to  $vcr-1$

$$VR[vrd][i] = M[R[rb] + 4 \times i]$$

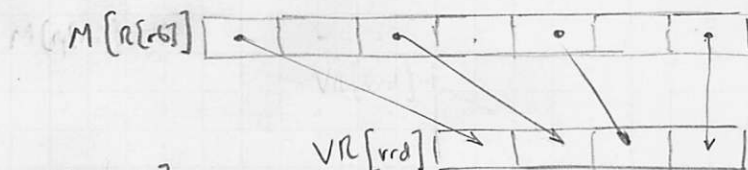


- VECTOR STRIDED LOAD

LWST.V vrd, rb, rs

for  $i=0$  to  $vcr-1$

$$VR[vrd][i] = M[R[rb] + R[rs] \times i]$$

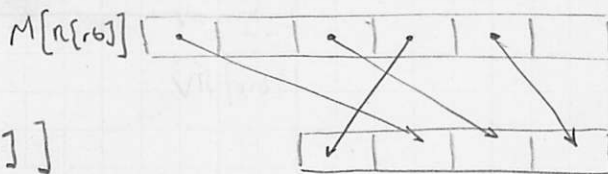


- VECTOR INDEXED LOAD

LWX.V vrd, rb, vrs

for  $i=0$  to  $vcr-1$

$$VR[vrd][i] = M[R[rb] + VR[vrs][i]]$$

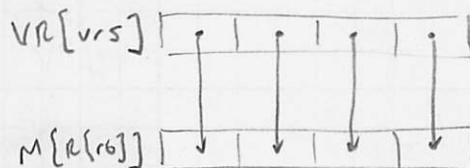


- VECTOR UNIT-STRADE STORE

SW.V vrs, rb

for  $i=0$  to  $vcr-1$

$$M[R[rb] + 4 \times i] = VR[vrs][i]$$



- SET LESS THAN (FLAG REG)

SET.f.vv fd, vrsφ, vrs1

for  $i=0$  to  $vcr-1$

$$FR[fd][i] = (VR[vrsφ][i] < VR[vrs1][i])$$

- VECTOR-VECTOR MULTIPLY (UNDER MASK)

MUL.vv vrd, vrsφ, vrs1, fs

for  $i=0$  to  $vcr-1$

if (FR[fs][i] == 1)

$$VR[vrd][i] = VR[vrsφ][i] \times VR[vrs1][i]$$

- SET ACTIVE VECTOR LENGTH

```
SETVL ra, rb
```

```
temp = min(hwvlen, R[rb]) ; vlen = min(hwvlen, n)
R[ra] = temp ; vlen = temp
```

## VECTOR CODE EXAMPLES

```
for (int i = 0 ; i < n ; i++)
  B[i] = A[i] * C
```

### SCALAR Assembly

loop:

```
lw r1, 0(r2)
mul r3, r1, r4
sw r3, 0(r5)
addiu r2, r2, 4
addiu r5, r5, 4
addiu r7, r7, -1
bne r7, r0, loop
```

### VECTOR Assembly

```
lw.v vr1, r2
mul.v vr2, vr1, r4
sw.v vr2, r5
```

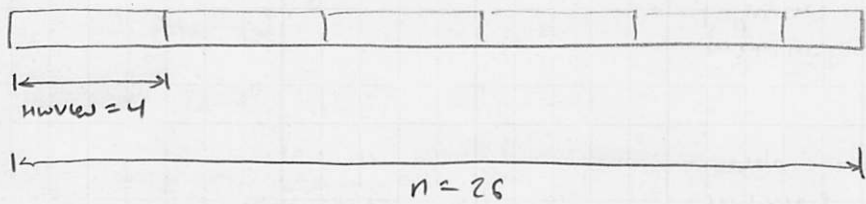
THIS ASSUMES  $n == hwvlen$ !  
WHAT IF  $hwvlen != n$ ?

if  $n = 2$  AND  $hwvlen = 4$

```
addiu r8, r0, 2
setvl r9, r8
lw.v vr1, r2
mul.v vr2, vr1, r4
sw.v vr2, r5
```

} EVEN THOUGH hardware vector length is four, because we used setvl to set the ACTIVE vector length we will only process two elements per vector instruction

if  $n$  or  $hwvlen$  is unknown at compile time MUST USE STRAPPING, execute in "hwvlen" "CHUNKS"



```
for ( int i = 0 ; i < n ; i ++ )
    B[i] = A[i] * C
```

SCALAR ASSEMBLY

```
loop:
lw   r1, 0(r2)
mul  r3, r1, r4
sw   r3, 0(r5)
addi r2, r2, 4
addi r5, r5, 4
addi r7, r7, -1
bne  r7, r0, loop
```

VECTOR ASSEMBLY

```
loop:
setvl r8, r7 ← INITIALIZES TO n
lw.v  vr1, r2
mul.v vr2, vr1, r4
sw.v  vr2, r5
sll   r9, r8, 2 # MULT VLR x 4
addu  r2, r2, r9 # pointers + VLR
addu  r5, r5, r9 #
subu  r7, r7, r8 # i = i - VLR
bne   r7, r0, loop
```

VECTOR CODE

EXECUTING MORE IRREGULAR DATA-LEVEL PARALLELISM

LESS STRUCTURE IN DATA ACCESS →

LESS STRUCTURE  
IN CONTROL FLOW

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

```
for ( i = 0 ; i < n ; i ++ )
    C[i] = x * A[i] + B[2*i]
```

```
for ( i = 0 ; i < n ; i ++ )
    if ( A[i] < 0 )
        x = y
    else
        x = z
    C[i] = x * A[i] + D[i]
```

```
for ( i = 0 ; i < n ; i ++ )
    E[C[i]] = D[A[i]] + B[i]
```

```
for ( i = 0 ; i < n ; i ++ )
    if ( A[B[i]] < 0 )
        C[i] = D[i] + E[i]
```

USING VECTOR INDEXED ACCESSES (SCATTER/GATHER)  
TO IMPLEMENT IRREGULAR MEMORY ACCESSES

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

lw.v vr1, r2 ← c ptr  
lw.v vr2, r3, vr1 ← b ptr  
sw.v vr2, r4 ← a ptr

USING VECTOR FLAGS/MASKS TO IMPLEMENT  
IRREGULAR CONTROL FLOW

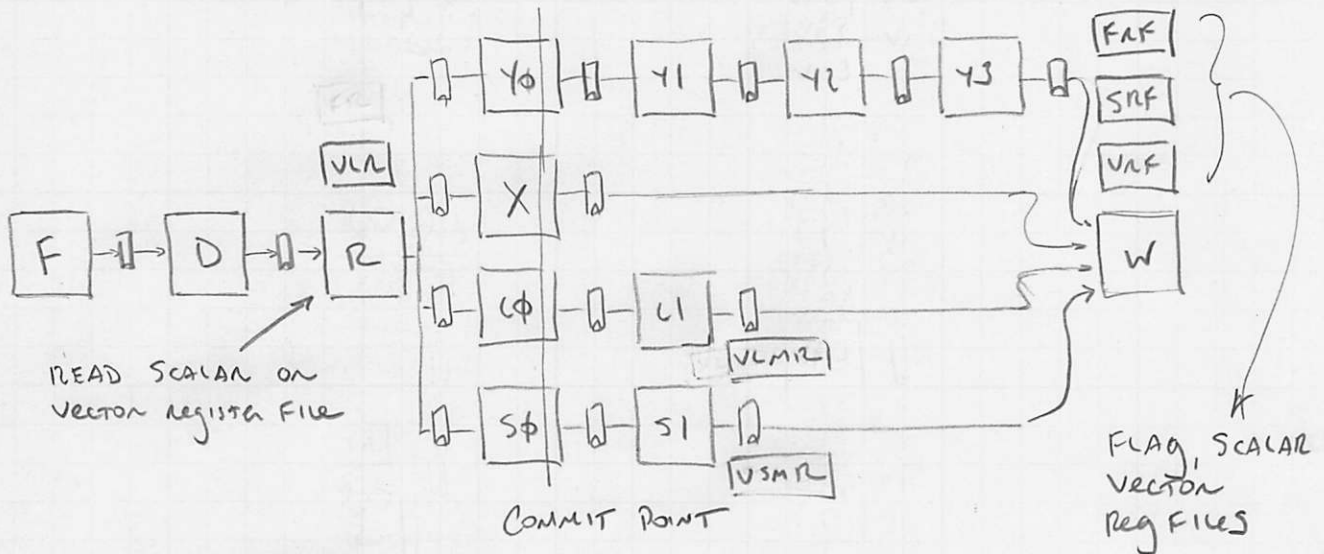
```
for (i=0; i<n; i++)
    if (C[i] < 0)
        B[i] = A[i] * C
```

lw.v vr1, r2  
sh.fw f1, vr1, vr0  
lw.v vr2, r3, f1  
mul.vs vr3, vr2, r4, f1  
sw.v vr3, r5, f1

f1  
f1  
f1

UNDER MASK

### VECTOR MICROARCHITECTURE



IN LI AND SI STAGES WE CAN NOW ALSO REQUEST TO READ OR WRITE A VECTOR OF DATA. SEPARATE LOAD/STORE PIPES WILL EVENTUALLY ALLOW MULTIPLE VECTOR INSTRUCTIONS TO EXECUTE IN PARALLEL.

ASSUME 32 VECTOR REGISTERS, VECTOR REGISTER FILE HAS A TOTAL OF 32 x HWLEN ELEMENTS

VLMBR = VLOAD MEM RESPONSE BUF VSMBR = VSTORE MEM REQ BUF

## EXECUTING VECTOR INSTRUCTIONS

MVLEN = 2 MUL.VS F D R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y1 Y3 W

MVLEN = 4 MUL.VS F D R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W

MVLEN = 8 MUL.VS F D R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W  
R Y0 Y1 Y2 Y3 W

MVLEN = 4 LW.V F D R I0 L1 W  
R I0 L1 W  
R I0 L1 W  
R I0 L1 W

IMPORTANT: WE READ BASE ADDRESS, GENERATE EFFECTIVE ADDRESS, AND SEND OUT A REQUEST TO READ A VECTOR OF DATA INTO VLMR WHEN HANDLING FIRST ELEMENT.

2ND-4TH ELEMENTS DO NOTHING IN R AND L0 STAGES. WE JUST WRITE BACK THE APPROPRIATE ELEMENT IN THE DESTINATION VECTOR REGISTER.

MVLEN = 4 SW.V F D R S0 S1 W  
R S0 S1 W  
R S0 S1 W  
R S0 S1 W

IMPORTANT: WE READ BASE ADDRESS, GENERATE EFFECTIVE ADDRESS, AND STORE THIS ADDRESS IN THE VSMR WHEN HANDLING THE FIRST ELEMENT.

All elements READ STORE DATA AND WRITE VSMR. LAST ELEMENT SENDS OUT MEMORY REQUEST TO WRITE A VECTOR OF DATA.



LET'S STUDY A SIMPLE LOOP EXECUTING ON THIS VECTOR MICROARCHITECTURE

SCALAR ASSEMBLY

```

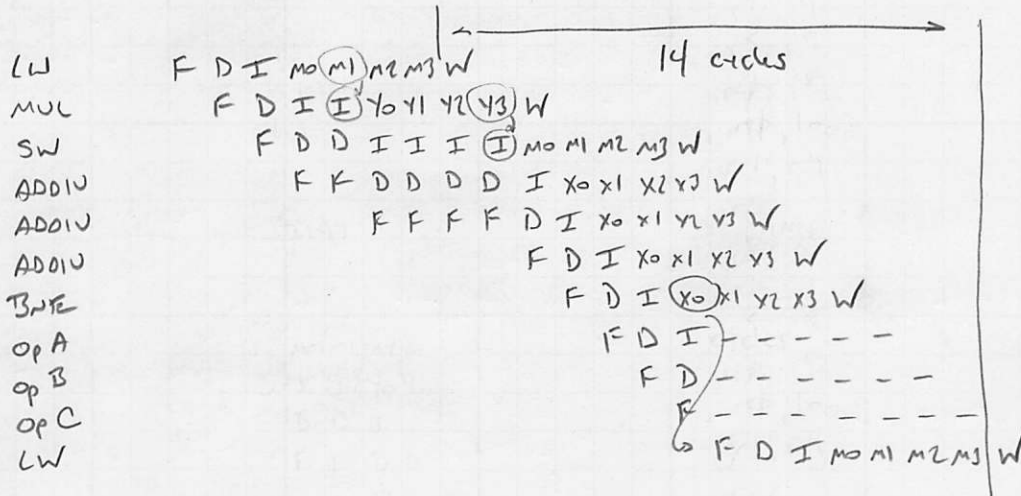
loop:
LW    r1, 0(r2)
MUL   r3, r1, r4
SW    r3, 0(r5)
ADDIU r2, r2, 4
ADDIU r5, r5, 4
ADDIU r7, r7, -1
BNE   r7, r0, loop
    
```

VECTOR ASSEMBLY

```

loop:
SETV1  r0, r7
LW.V   vr1, r2
MUL.VS vr3, vr1, r4
SW.V   vr3, r5
SLL    r9, r0, 2
ADDU   r2, r2, r9
ADDU   r5, r5, r9
SUBU   r7, r7, r0
DNE    r7, r0, loop
    
```

FOR REFERENCE WHAT IS THE EXECUTION TIME OF SCALAR CODE?



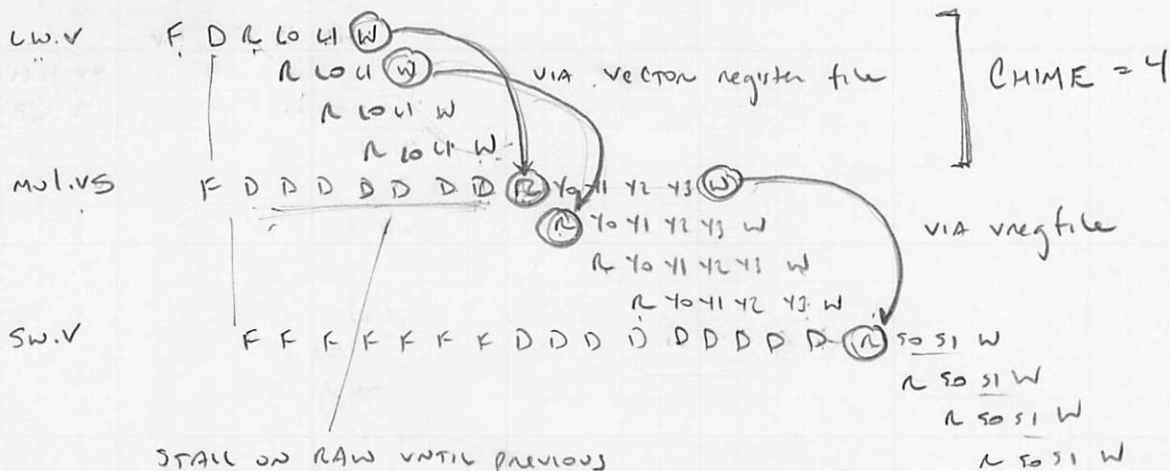
ASSUME 64 ITERATIONS

#INSTN = 7 x 64 = 448  
 Cycles/INSTN = (14 x 64) / 448 = 2  
 Cycle time = 1ns

TOTAL EXECUTION TIME = 896 ns

40 SHEETS EYE-GLASS, 4 SQUARES  
 42-382 100 SHEETS EYE-GLASS, 4 SQUARES  
 42-389 200 SHEETS EYE-GLASS, 4 SQUARES  
 National Brand

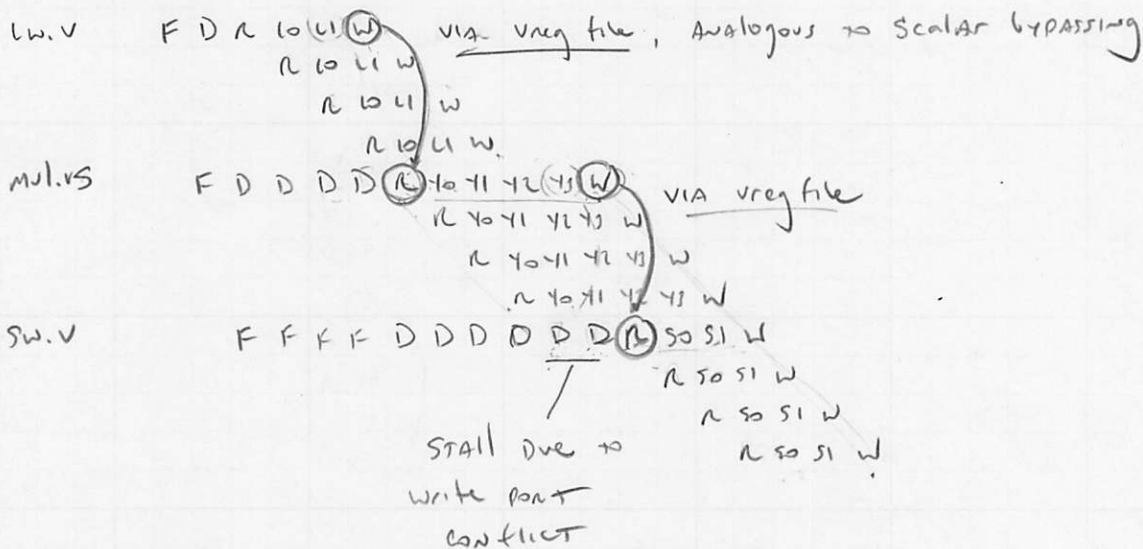
BASIC EXECUTION, HWLEN = 4



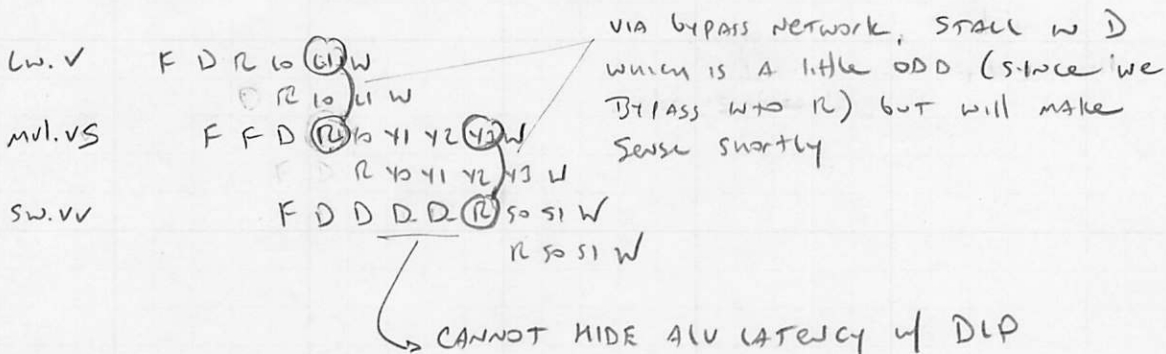
STALL ON RAW UNTIL PREVIOUS VECTOR INSTRUCTION HAS WRITTEN BACK RESULT.

\* DEP CHECKING ONCE PER VECTOR INSTR!

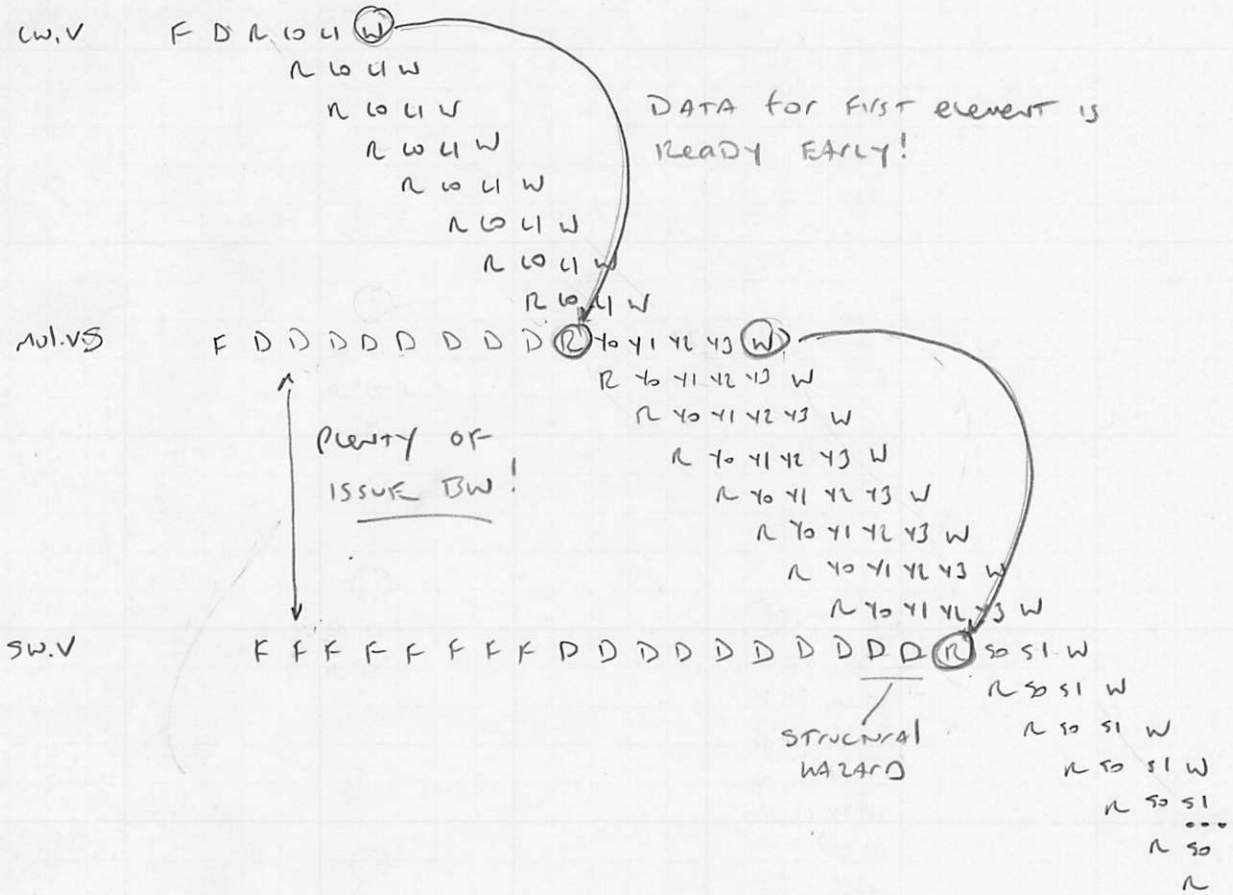
CHAINING, HWLEN = 4



CHAINING, HWLEN = 2 (ASSUME n ALSO = 2)



Chaining, HWUes = 8 (ASSUME n = 8)



KEY OBSERVATION: WITH RELATIVELY LITTLE FETCH/ISSUE BANDWIDTH + INORDER FETCH/ISSUE, WE CAN STILL KEEP FUNCTIONAL UNITS BUSY FOR MANY CYCLES.

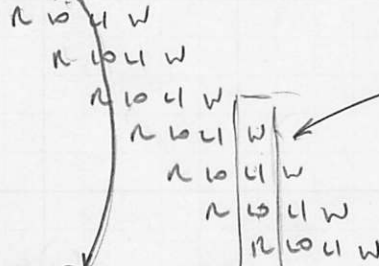
\* ADD MORE VREG FILE PORTS TO EXPLOIT VECTOR INSTR. LEVEL PARALLELISM

TO KEEP ALL FUNCTIONAL UNITS BUSY AT SAME TIME, NEEDS:

Y-PIPE	2r	1w
X-PIPE	2r	1w
L-PIPE	1r	1w
S-PIPE	2r	0w
<b>TOTAL</b>	<b>7r</b>	<b>3w</b>

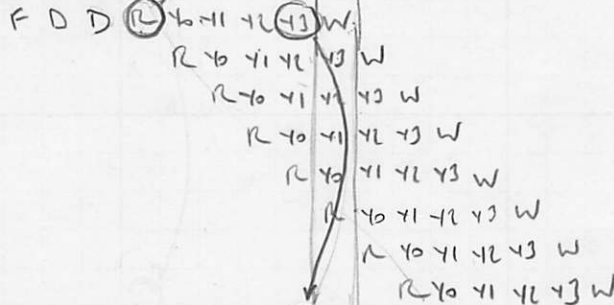
# CHAINING, HWLEN = B (ASSUME $n=B$ ) w/ $F, 3W$ VECTORS

CW.V F D R L O U W

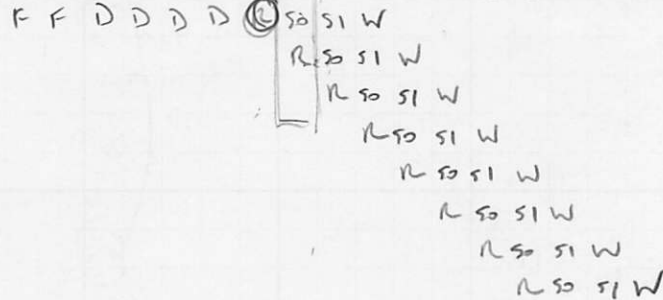


3 vector w/str exec AT SAME TIME, EQUIVALENT TO 12 SCALAR OPERATIONS IN FLIGHT AT ONCE w/ SINGLE ISSUE!

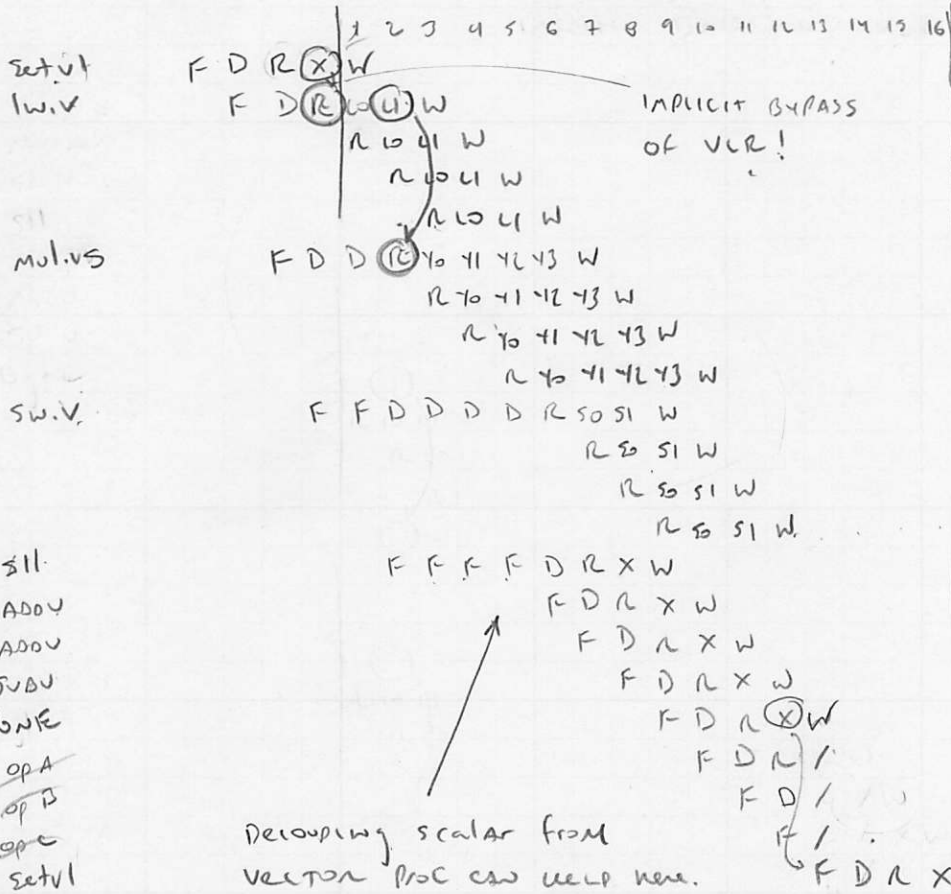
MULVS F D D



SW.V F F D D D D



30 SHEETS EYE-EASE  
100 SHEETS EYE-EASE  
200 SHEETS EYE-EASE  
National Brand



16 cycles/1In

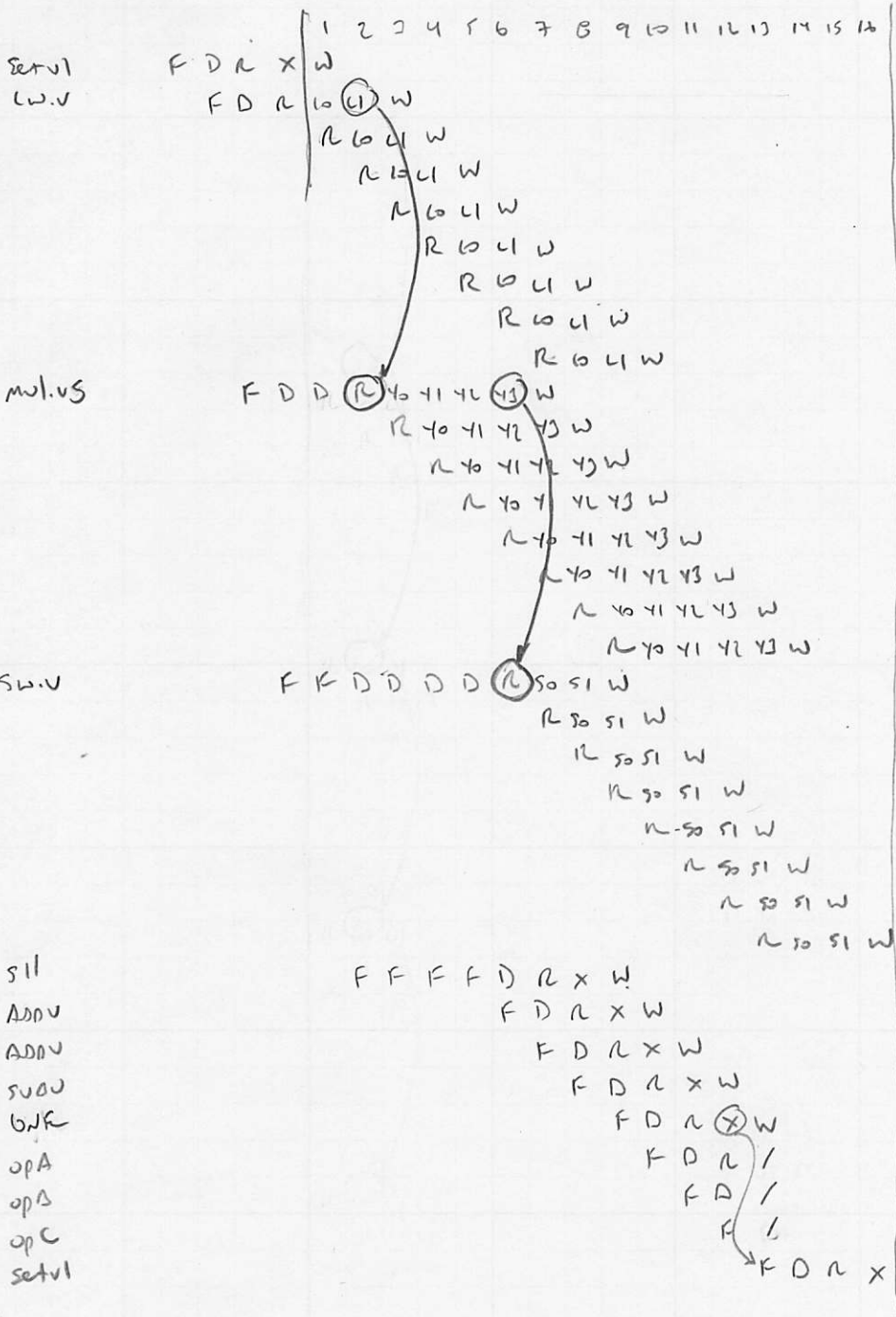
# instr =  $9 \times 1.6 = 144$   
 cycles/instr =  $(16 \times 16) / 144 = 1.78$   
 cycle time = 1ns

TOTAL EXECUTION TIME = 256 cycles

SPEEDUP OF 3.5 BUT NEEDS TO ADD 4x REGISTERS AND 7-3w PORTS

WHAT IF HUBLES = 8?

National Brand



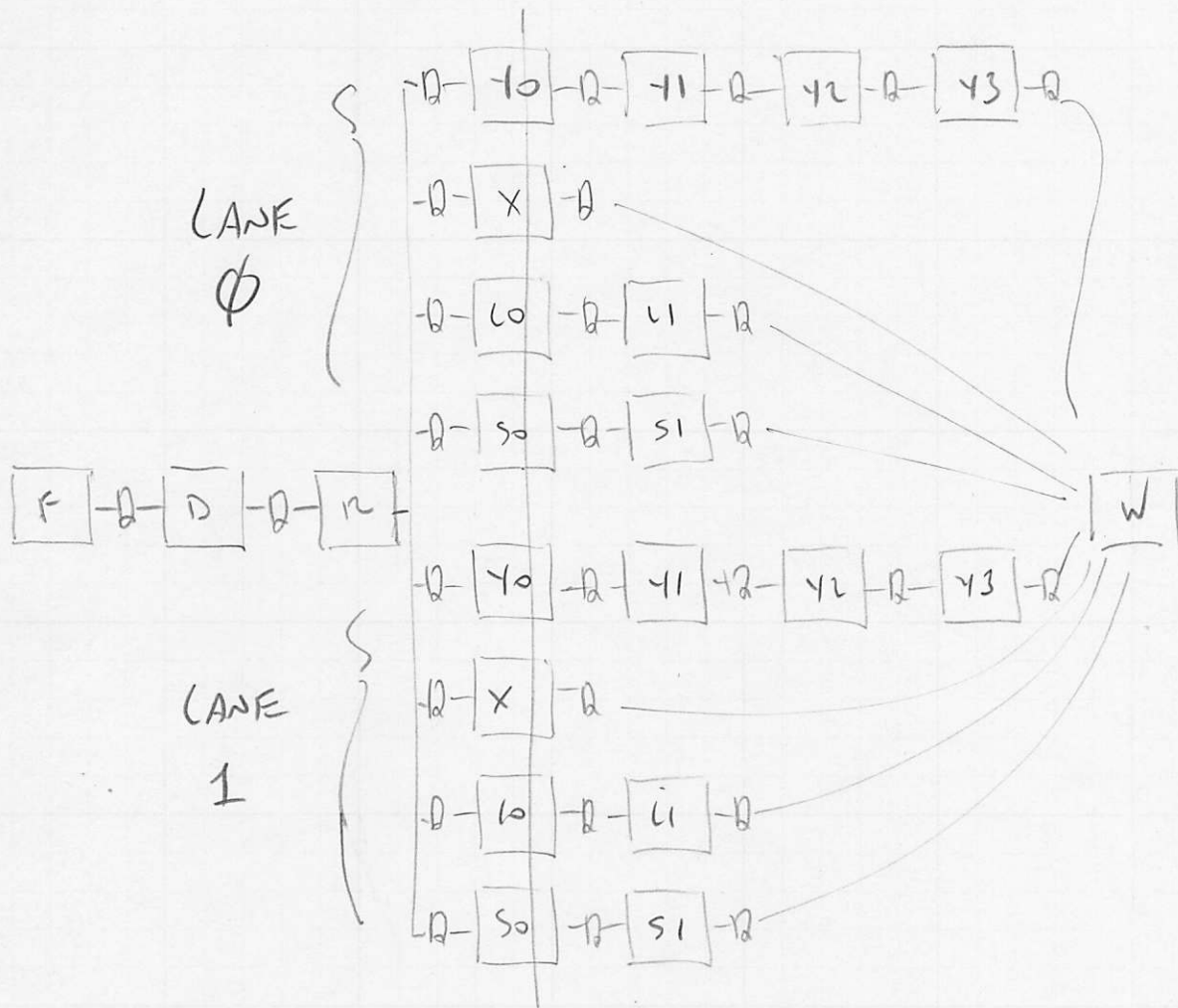
# INSTN = 9 x 8 = 72  
 CPI = (16 x 8) / 72 = 1.78  
 cycle time = 1ns

TOTAL EXECUTION TIME = 128.16

SPEEDUP OF 7x WITH 8x MORE REGISTERS + STILL 7x W PORTS

301 50 SHEETS EYE-EASE 10 SQUARES  
 42-382 100 SHEETS EYE-EASE 5 SQUARES  
 42-389 200 SHEETS EYE-EASE .5 SQUARES  
 National Brand

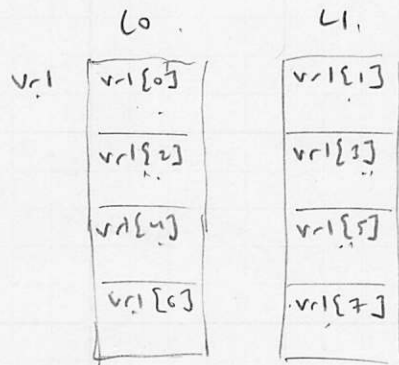
# MULTIPLE LANES



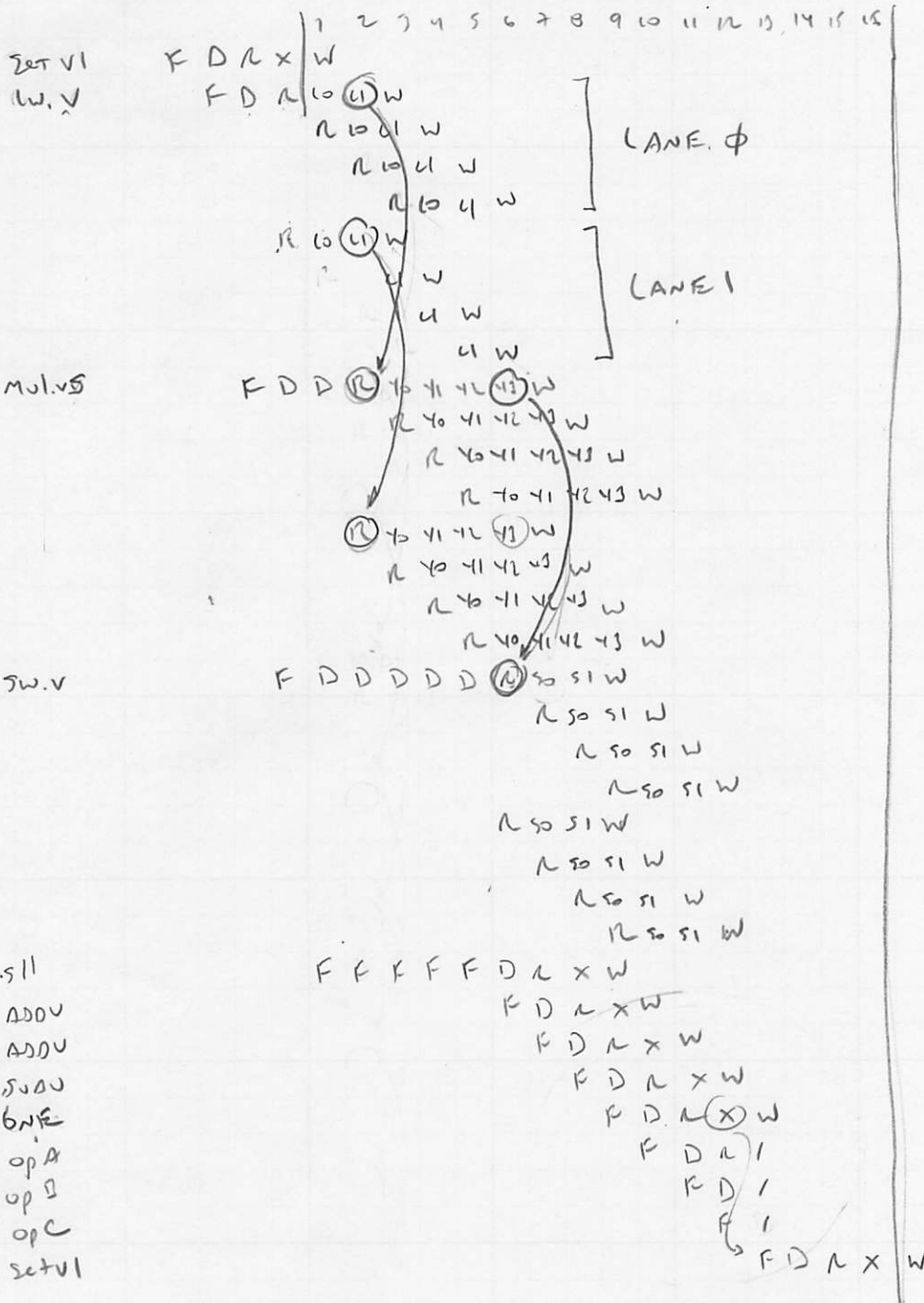
Commit Point

VRegist File IS PARTITIONED BY LANE.

ASSUME HWLEN = 8 w/ 2 LANES



only needs 7r3w / lane

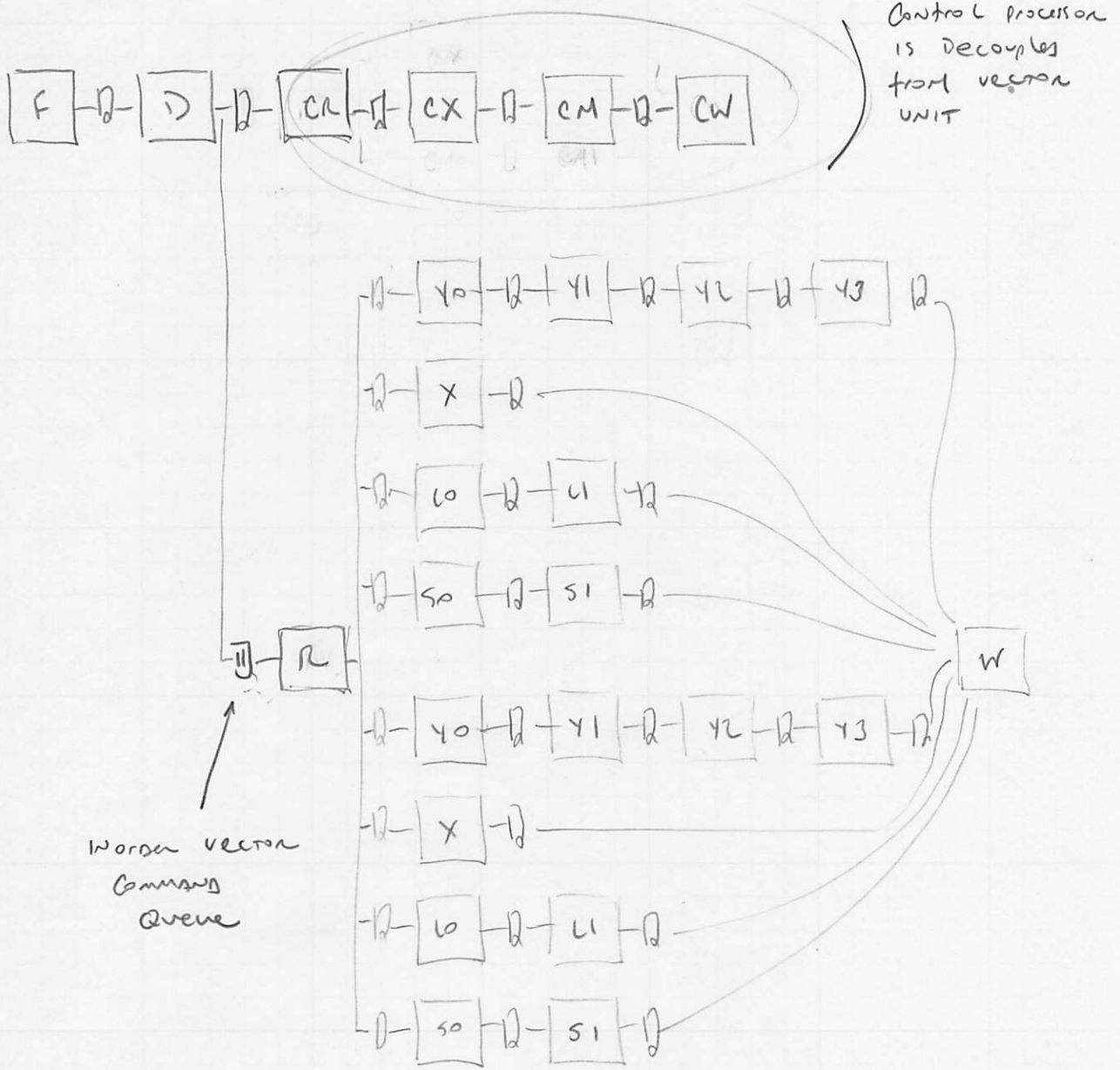


16 cycles / instr, SAME AS w/ 1 lane w/ view = 4

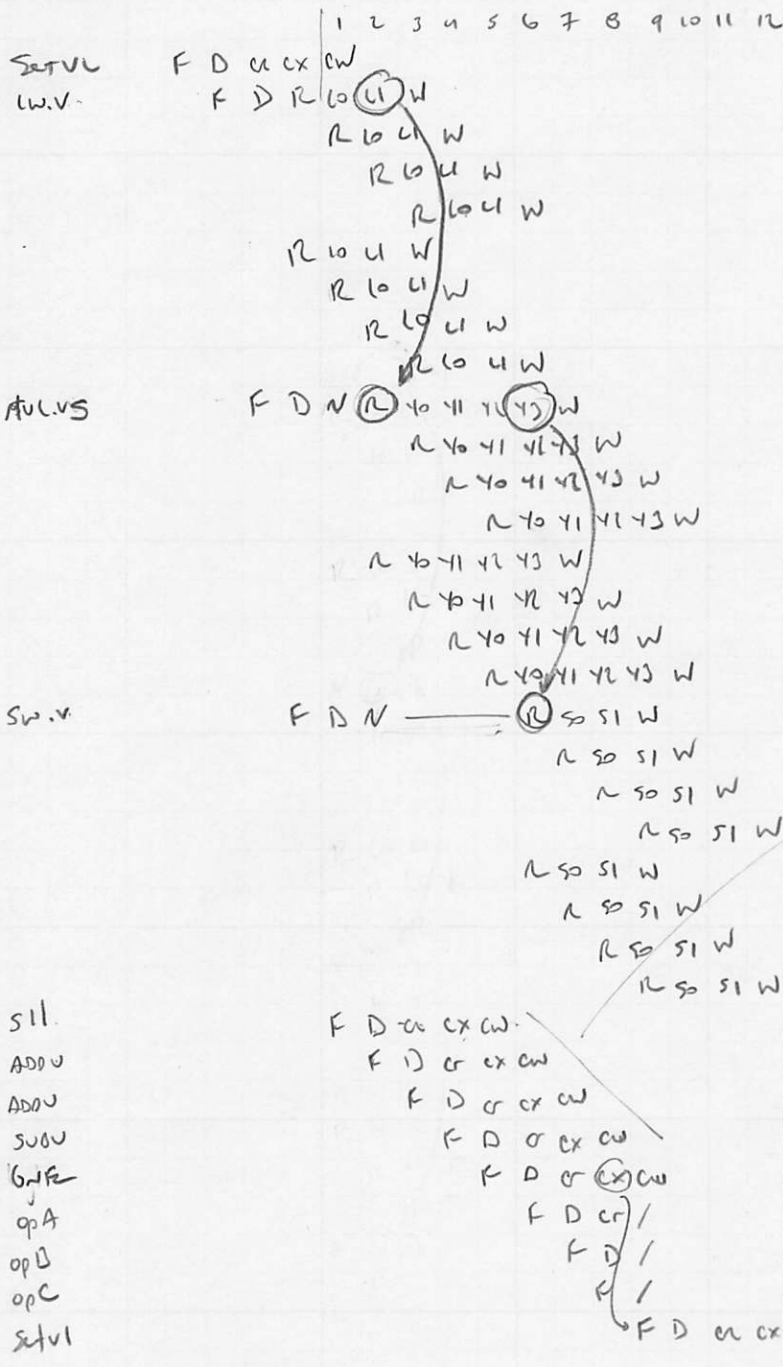
- SCALAR CODE GETTING IN THE WAY
- ↑ view
- Decouple scalar processor



MULTIPLE LANES w/ Decoupled CONTROL PROC



381 50 SHEETS EYE-EASE SQUARES  
42-362 100 SHEETS EYE-EASE 5 SQUARES  
42-369 200 SHEETS EYE-EASE 5 SQUARES  
National Brand



Decouples control proc  
 MEANS CANNOT SUPPORT  
 precise exceptions w  
 vector unit, although  
 could add ROB  
 trace.

Complexity issue  
 this execution  
 latency!

# instr =  $9 \times 8 = 72$   
 CPI =  $(12 \times 8) / 72 = 1.33$   
 cycle time = 1 ns

TOTAL execution TIME = 96 cycles!  
 ↳ with worse, single issue!