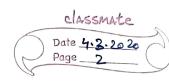
## SIMD

NUDIA tech reports VECTOR ARCHITECTURE formi replan. vector achitecture grab sets of data elements scattered arch. about memory, place them into large sequential register files, operate on data, in those register files and then dispense the results back to memory single instruction works on vector of data, which results in dozens of operations on individual claments. these large registers out as compiler controlled buffers, both to hide memory latency and to leverage memory bandwidth, vector loads and stones are deeply populined. RVV -> vector extensions bon RISC-V, loosely based on to years-old-cray, +c RV64V) the primary components of ISA of RV64V rector registers & each vector vegister holds a single vector, and RVB4V has 22 of them, each 64-bits wide, (1) xxx atleast 16 read ports and 8 write ports Weston functional wonits! deach unit wix fully pipelined and vit can stant a new operation of every clack cycle. a control

unit is needed to detect hazards (structural), for



functional units and data hazards on register accesses. We have a third unit, load & store unit the vector memory unit loads on stores a vector to or from memory

the vector loads and stores are folly pipelined, this unit also normally handles scalar loads and stores.

scalar registers likewike provide data as input
to vector functional units as well as computed addresses
to pass to the vector load and store unit.

there are normal 31 general purpose registers, and 32 floating point registers for the RV64V.

fld for a Hload scalar a man

addi 1028, 18285, # 256 # local address to load

fruid 61, 61, 60 th anxis

fed 62,0 (x6) # load Y[i]

fadded 62,62,61 # a+x[i]+2/[i]

folds f2, 0(x6) # Store into Y[i]
addi x5, x5, #8 # cincrement indento x

andi X6, x6, #8 Charement characto Y.

bne 1828 x5, loop # check of done

| - 1 |          |                |                               |  |
|-----|----------|----------------|-------------------------------|--|
| 1   | vsetdcfg | 4 * FP64       | # enable a DPFP negs          |  |
|     | jid      | fo, a          | # load scalar.                |  |
|     | vold     | vo, x5         | # load Veopor x               |  |
|     | Vmo      | VI, vo, to     | # Vector scalar mult.         |  |
|     | vid      | V2, <b>x</b> 6 | # lead vectory                |  |
|     | vadd     | V3, V1, V2     | # Vector-vectoradd,           |  |
|     | ust.     | V3.X6.         | # Store the som,              |  |
|     |          | . n. 10 - 11   | # disable vector a registers. |  |
|     | •        |                | 10/13/4/                      |  |
| - 1 |          |                |                               |  |

8 incts. Vs 258 for RV64V

vectore execution time depends upon

@ length of operand vectors.

& structural hazards, among the operations

3 data dependencies

load and add could happen concurrently.

convoy: set of vector instructions that could potentially execute together, the instruction in a convoy must not contain any structural hazands, chaining allows a vector operation to start as soon as the individual elements of its vector source operands becomes available.

blekible chaining allows a vector ainstruction to chargein to essentially any other active vector the truction.

Chime: - simply a unit of time taken to execute one convoy.

DMPS

chime approximation ignores some processon Specific overheads.

vld vo, x5 # load vector x vmub VI, vo, 60 # vector sealon moltiply

vld v2, x6 & Load vector y.

vooled V3 V1, V2 # Vector - vector rada VSL V3, X6 # Store the cum.

vld vadel

vld, mul

fp-add: 6 cycles.

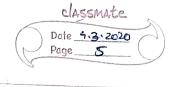
fp-moltiply; 7 cycles 6p-divides: 20 cycles.

vector-load : 12 cycles

optimization that either improve the performance on morease the the type of programs that can

run well on vector arch.

how can a vector processon exercite a single vector. forter than one element per clock cycle, multiple elements per clock cycle to improve the perf.



how does a vector processor hardle programs where
the vector lengths are not the same as maximum
vector length, we need an efficient soln in this
Common case,

Strip mining 
what happens when there is an if statement inside

what does a vector processor need from memory system?

how does a vector processor handle multiple dimensional matrices (stride)

of there are sparse matrices, how does a vector

(base rector, index vector)

S. E. C. C. No. of Co.

how do you program a vector computer,