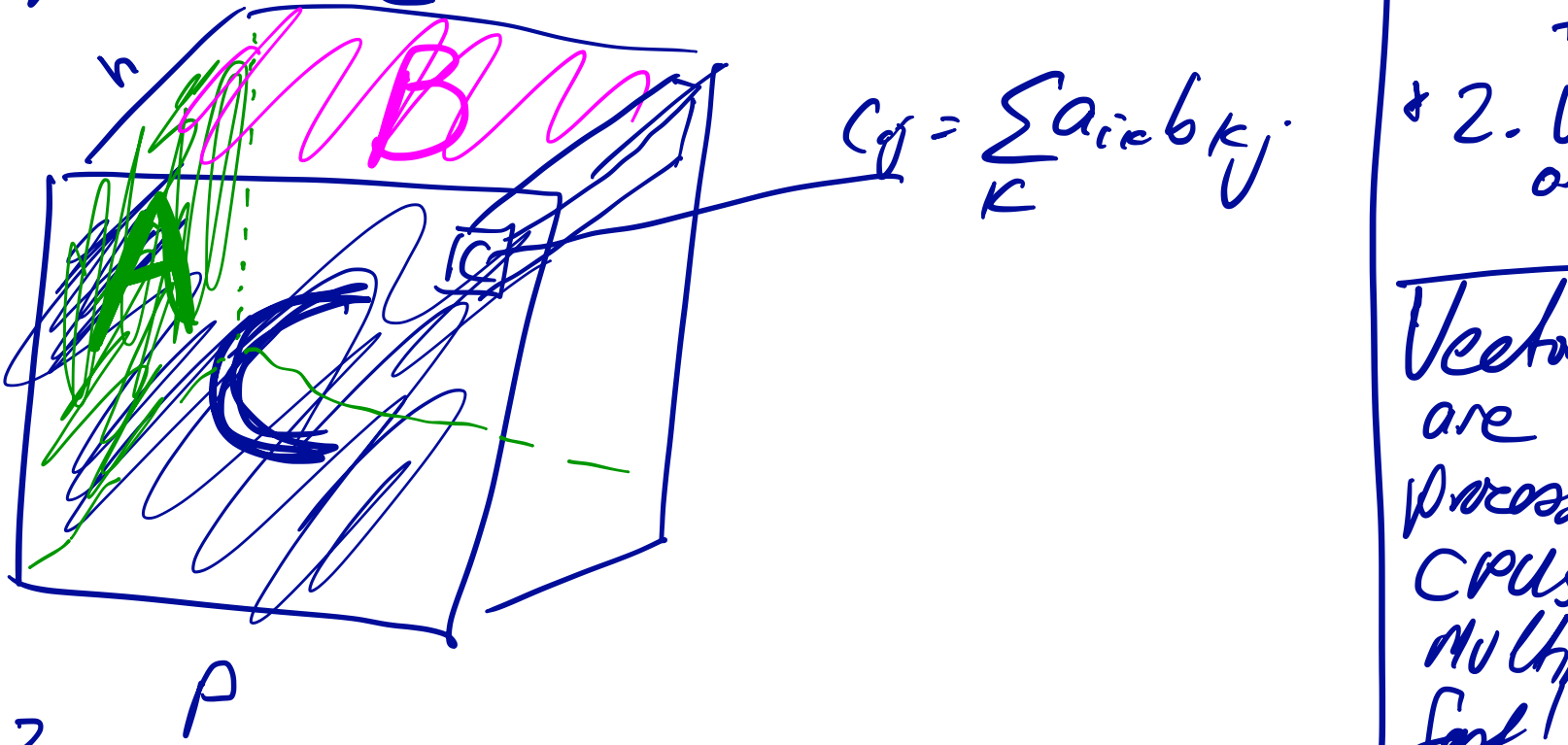


# Lecture 15

## The Mat Mul Cube

A: M x n  
B: n x p



Why the cube?  
It describes memory locality

Why does this matter?

- 1. Memory Locality + Memory Hierarchy
- 2. Vector Units on CPUs

Vector Units are specialized processing units on CPUs to do multiple float ops fast

A common Vector op:

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \\ a_4 \end{bmatrix} * \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix}$$

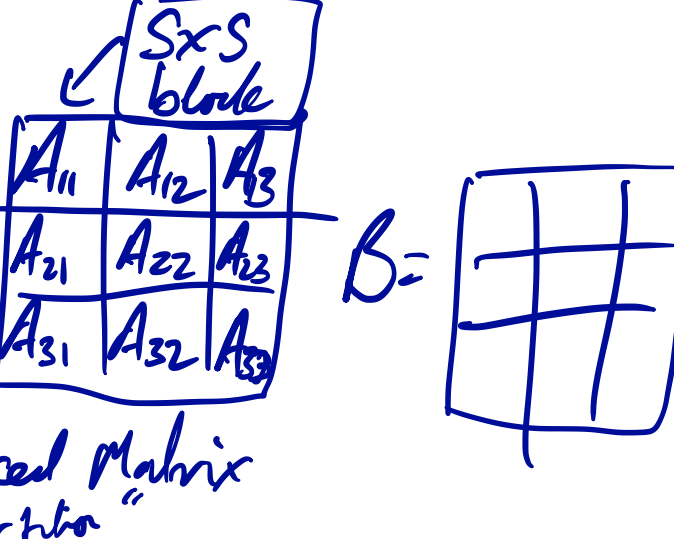
a vector unit can do all 4 @ once!

In taking @simd for ...  
it will try and use the vector unit.

function dot(a,b)

@simd for i = 1 : length(a)  
C[i] = a[i] \* b[i]  
end

Another View on Mat Mul



Then SxS blocks of C

$$C_{21} = A_{21} B_{11} + A_{22} B_{21} + A_{23} B_{31}$$

SxS - MatMul

(16x16) 256 "p"

$$C_{ij} = \sum_{k=1}^p A_{ik} B_{kj}$$

SxS matrix

This is called the Blocked Approach

