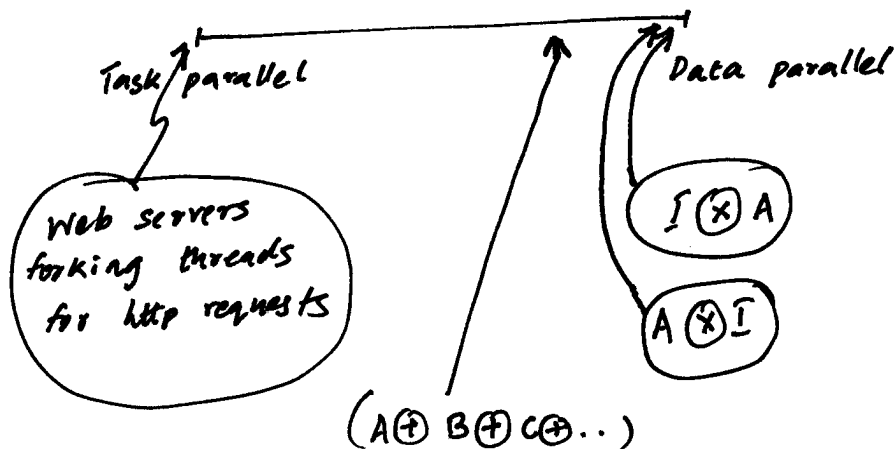


Task vs. Data Parallelism: continuum



Parallelism in mathematical constructs.

First, we look at different ways to visualize/represent constructs like $(I \otimes A)$ and $(A \otimes I)$:

Consider $(I_n \otimes F_2)$: $(F_2 = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix})$

Matrix representation looks like:

$$M = \begin{bmatrix} \begin{matrix} 1 & 1 \\ 1 & -1 \end{matrix} & & & \\ & \begin{matrix} 1 & 1 \\ 1 & -1 \end{matrix} & & \\ & & \begin{matrix} 1 & 1 \\ 1 & -1 \end{matrix} & \\ & & & \dots \\ & & & & \dots \end{bmatrix}$$

(n times)

Our transform matrix M is used to transform input

vector X into output vector Y : $Y = M \cdot X$

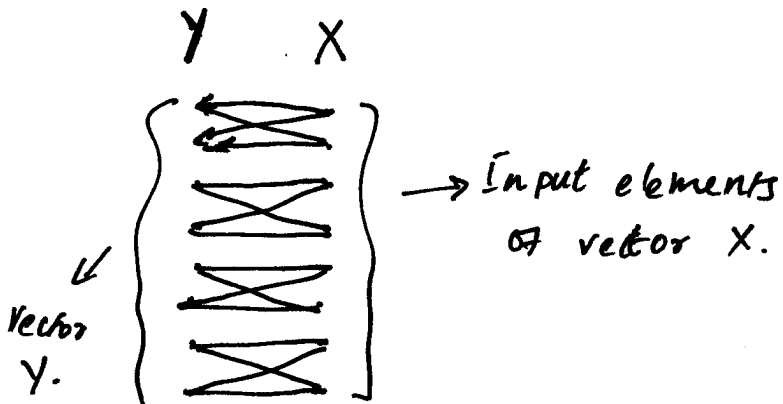
To see how the data flows from x to Y , we draw

a Data Flow Graph:

For eg., $(I_n \otimes F_2)$

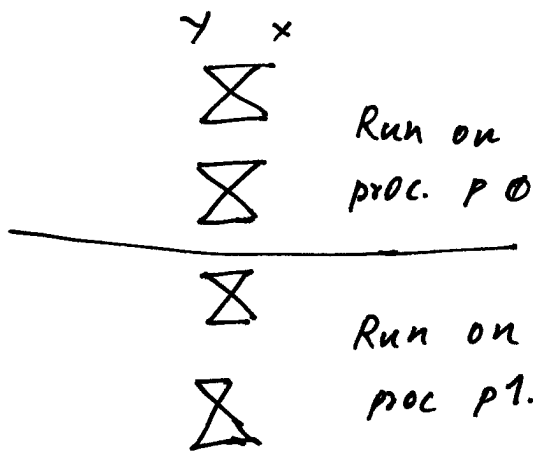
looks like:

(data flows from right to left)



PARALLELIZATION

This visualization shows us how to partition the transform to run in parallel on multiple processors.



In general, $(I_n \otimes A)$ is already (naturally) parallelized for execution on upto n processors

VECTORIZATION (form of parallelism):

Example:

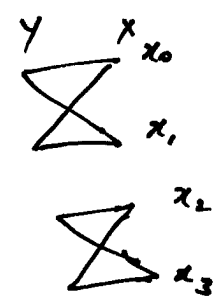
Goal: To vectorize $(I_2 \otimes F_2)$ on a 2-way vector machine.

- Constraints:
- ① Vector loads (loads to vector registers) must be done on consecutive memory elements.
 - ② loads might be expensive (cache misses etc.)
 - ③ Vector ops. use vertical parallelism.

Attempt #1: V_n : vector registers

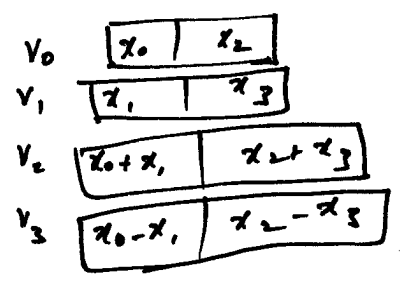
Load $V_0 \leftarrow x_0$
 Load $V_0 \leftarrow x_2$
 Load $V_1 \leftarrow x_1$
 Load $V_1 \leftarrow x_3$

} Might actually involve more operations (rotate, and etc.) on most machines)



Add $V_2 \leftarrow V_0 + V_1$
 Sub $V_3 \leftarrow V_0 - V_1$

Store $Y_0 \leftarrow V_{2,0}$
 Store $Y_0 \leftarrow V_{2,1}$
 Store $Y_2 \leftarrow V_{3,0}$
 Store $Y_3 \leftarrow V_{3,1}$



Problem with this: Too many loads stores.
 \therefore Might not get any vector speed up.
 (might get slowdown).

ATTEMPT #2

Load $V_0 \leftarrow x_0, x_1$
 Load $V_1 \leftarrow x_2, x_3$
 Perm $V_2 \leftarrow V_0, V_1 (0, 2)$
 Perm $V_3 \leftarrow V_0, V_1 (1, 3)$ } These are cheap!
 Add $V_4 \leftarrow V_2 + V_3$
 Sub $V_5 \leftarrow V_2 - V_3$
 Perm $V_6 \leftarrow V_4, V_5 (0, 2)$
 Perm $V_7 \leftarrow V_4, V_5 (1, 3)$ } cheap!
 Store $Y_0, Y_1 \leftarrow V_6$
 Store $Y_2, Y_3 \leftarrow V_7$

We replaced expensive loads/stores with register permutations.

With the tensor notation:

$$\begin{aligned}
 (\bar{I}_2 \otimes F_2) &= L_2^T (F_2 \otimes \bar{I}_2) L_2^T \quad \text{--- ①} \\
 \text{(Also, } \mathbf{A} (F_2 \otimes \bar{I}_2) &= L_2^T (\bar{I}_2 \otimes F_2) L_2^T \quad \text{--- ②}
 \end{aligned}$$

① Helps us go from $(\bar{I}_n \otimes A)$ to $(A \otimes \bar{I}_n)$. Note that $(A \otimes \bar{I}_n)$ is naturally vectorized for an n-way vector machine.

In general,

$$(I_n \otimes A_m) = L_n^{nm} (A_m \otimes I_n) L_m^{mn} \quad \text{--- (3)}$$

$$(A_m \otimes I_n) = L_m^{mn} (I_n \otimes A_m) L_n^{nm} \quad \text{--- (4)}$$

③ & ④ help us convert between parallelized and vectorized forms of the same transform.