Agenda:

1. Previous Session $\rightarrow$ Multiprocessing & Multi-threading

2. Design parallel algorithms
   for training
   $\rightarrow$ Data science
   $\rightarrow$ ML
   $\rightarrow$ DL

3. Parallel processing for productionization

Parallel processing is a vast - subject

introductory - session

Introduce ideas
as we solve
problems

many models
of computation
$\begin{bmatrix} \text{CPUs,} \\ \text{GPUs, FPGAs} \\ \text{clusters} \end{bmatrix}$

Parallel        vs        Distributed    computing

multi-core                              Spark / Hadoop

multi-threaded                          [ Disb. memory ]

GPU

[ Shared memory ]

# Parallel Matrix Addition

assume: A & B fit into RAM

data parallelism



$$n_1 \quad n_2 \quad n_3$$

$A$

$n \times M$

$B$

$n \times M$

[ROW-MAJOR]

$$A + B = C$$

Shared memory: A, B & C

CODE: https://docs.python.org/3/library/multiprocessing.shared_memory.html

# Parallel matrix Multiplication

$\longrightarrow$ most ML/DL algos' update rules

$-$ many variations of parallel algos.

$$C = \begin{bmatrix} \boxed{\phantom{xxxxxxx}} & n_1 \\ \boxed{\phantom{xxxxxxx}} & n_2 \\ \vdots \\ \boxed{\phantom{xxxxxxx}} \end{bmatrix} \times \begin{bmatrix} \phantom{xxxxxxx} \end{bmatrix}$$

$C = $
$n \times K$

$A$
$n \times M$

$B$
$m \times K$

$\hookrightarrow$ Data Parallelism

APPLIED COURSE

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \times \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix}$$
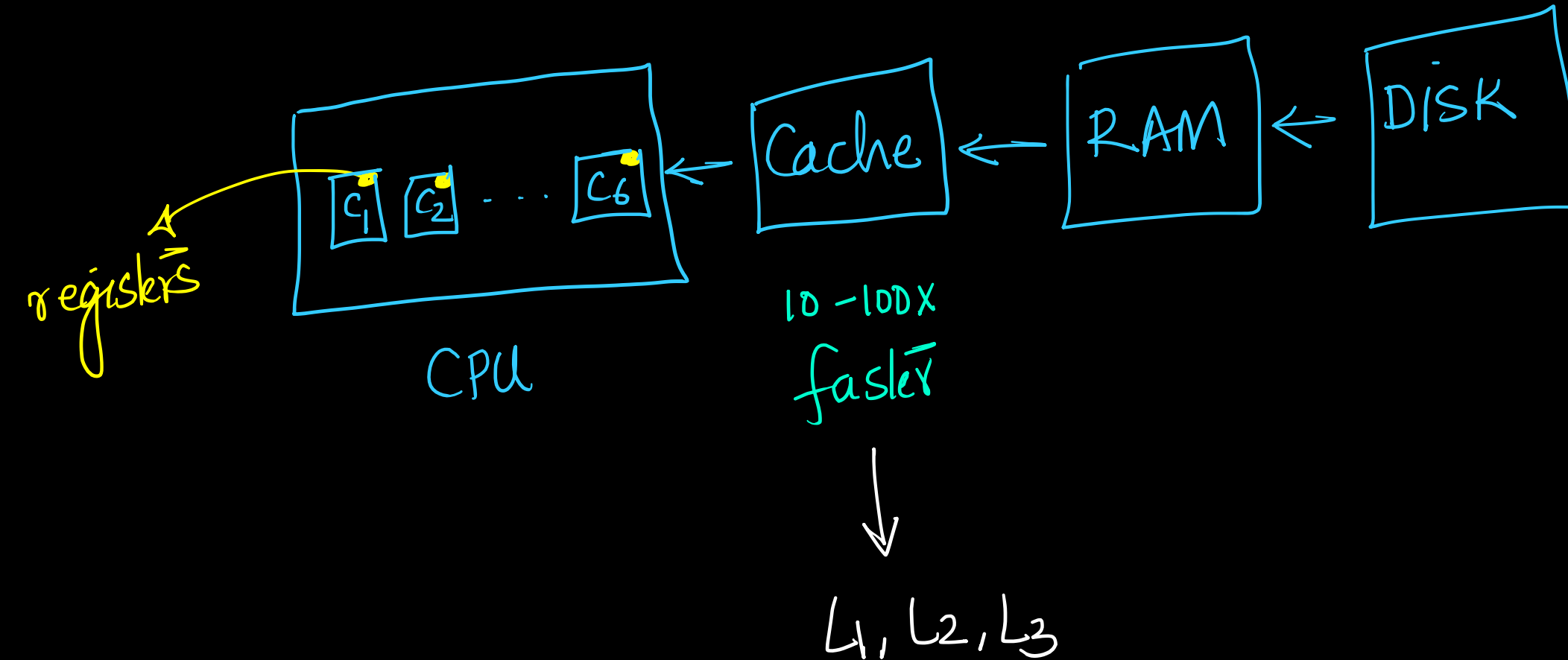
$$A \qquad\qquad B$$

Sub-matrix

$$C = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}$$

NOTE: Cache vs RAM                    [Cache-aware algos]



registers

CPU

$10-100X$
faster

$L_1, L_2, L_3$

# Vectors & Tensors

$n \times 1$
matrices

$n \times m \times K$

$\Downarrow$

$K (n \times m)$ matrices

# Logistic (or) Linear regression

loss

$$W^i_{new} = W^i_{old} - \eta \frac{\partial L}{\partial W^i}\bigg|_{W_{old}}$$

$i^{th}$ component of $W$
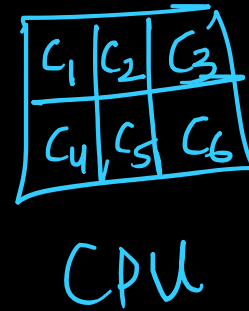
learning rate

summation over $j: 1 \to n$ train points

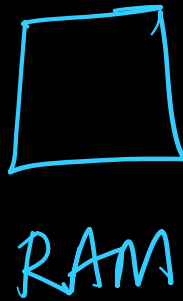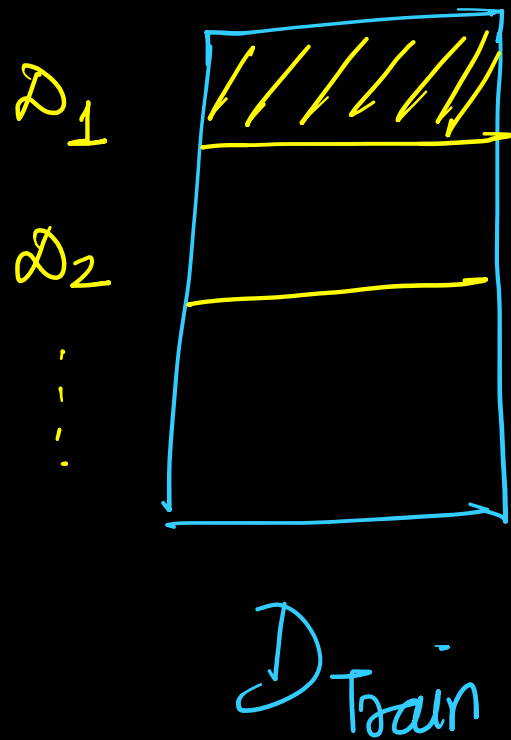# (Q) what if $D_{Train}$ doesnt fit into RAM



$D_1$

$\partial_2$

$D_{Train}$

RAM

| $C_1$ | $C_2$ | $C_3$ |
|---|---|---|
| $C_4$ | $C_5$ | $C_6$ |

CPU

USE distributed computing

$\hookrightarrow$ SPARK
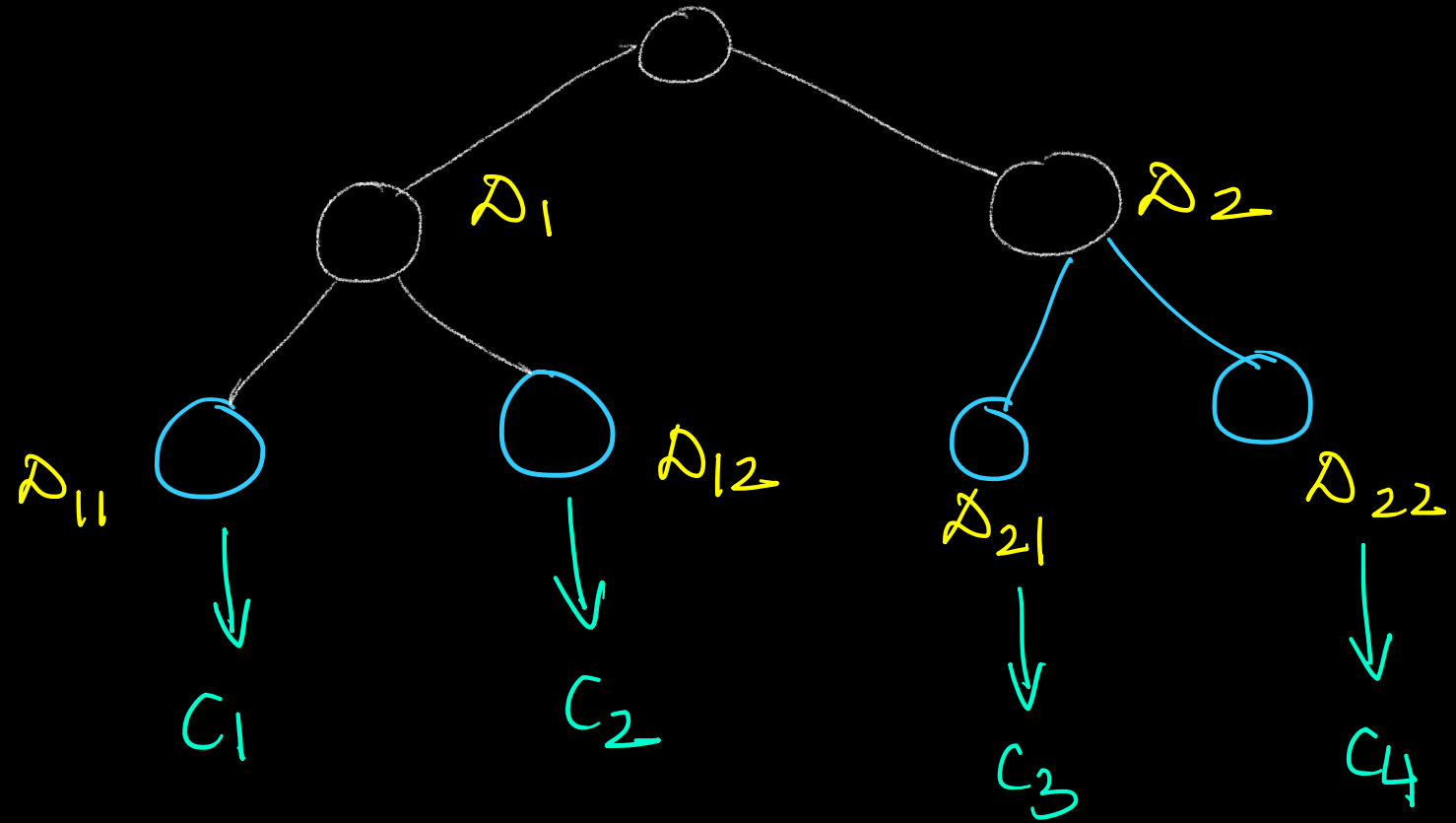
# Decision Trees:

① 

$f_1$ } 
$f_2$ } $\rightarrow C_1$

$f_3$ } $\rightarrow C_2$

$f_d$

$D_1$

Task - parallelism

best split & entropy gain
per feature

② Data parallelism

# Random Forest

$C_1$    $C_2$    $C_3$ --- --- $C_6$

$\downarrow$    $\downarrow$    $\downarrow$       $\downarrow$

$T_1$    $T_2$    $T_3$ --- --- $T_6$

$T_7$    $T_8$    $T_9$ --- --- $T_{12}$

$\vdots$

Trivially-paralledizable

# Pool of processes

CODE: https://docs.python.org/3.8/library/multiprocessing.html#using-a-pool-of-workers

# GBDT

$$F_M(x) = F_{M-1}(x) + \gamma_m h_m(x)$$

model with m-1 base learners

$m^{th}$ base learner

inherently serial

soln: biuld each base learner (DT) parallely.

# Deep Learning



$W_{5\times3}$

matrix - mul : forward pass

backward pass:

$$W_{5\times3} = \begin{bmatrix} c_1 & c_2 \\ c_3 & c_4 \end{bmatrix}$$

$5 \times 3$

updates

# Deep - Learning :

- matrix - multiplication

- Update weights/params per layer parallely

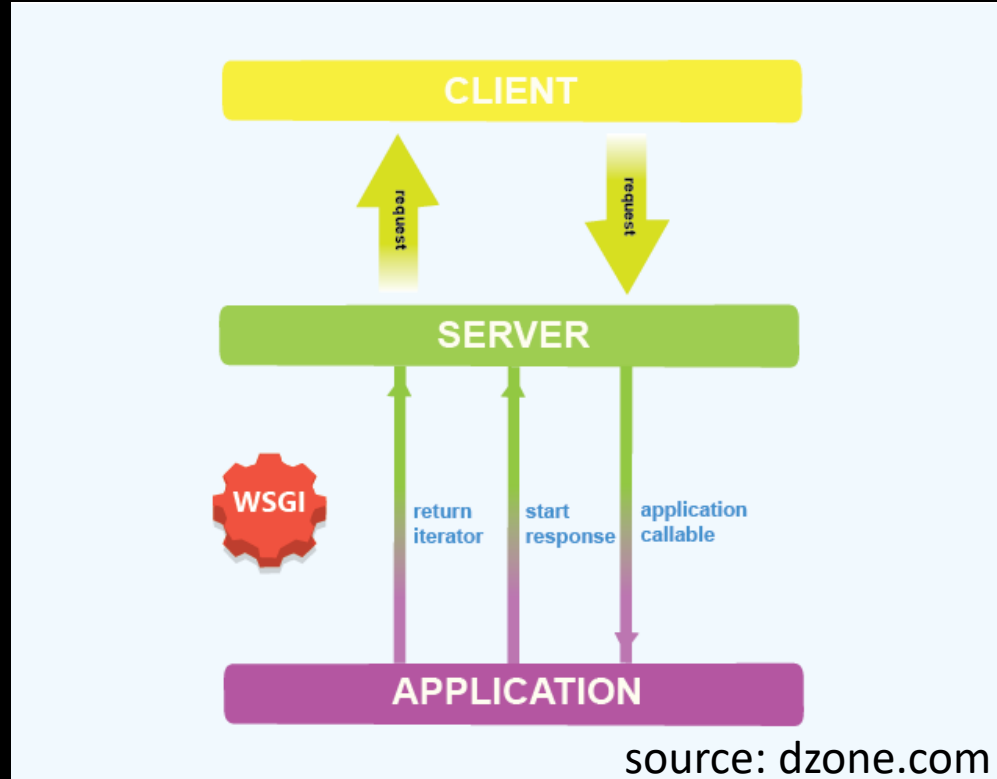$$\downarrow$$

independent

# Productionization :

Flask
API $\longleftarrow x$

$\quad\quad \llcorner \longrightarrow y$

$\quad\quad f(x)$

simultaneous requests

# WSGI : Web-server Gateway Interface



CLIENT

SERVER

WSGI | return iterator | start response | application callable

APPLICATION

source: dzone.com

Gunicorn
https://gunicorn.org

↓

Green Unicorn

↓

WSGI HTTP Server
for unix

```
$pip install gunicorn

$ cat myapp.py
from flask import Flask
app = Flask(__name__)


@app.route('/')
def hello_world():
    return 'Hello, World!'


if __name__ == "__main__":
    app.run()


$ python3 myapp
Use a production WSGI server instead.
 * Debug mode: off
 * Running on http://127.0.0.1:5000/ (Press CTRL+C to quit)
```

```
$cat wsgi.py
from myapp import app


if __name__ == "__main__":
    app.run()


$ gunicorn -w 4 -b 127.0.0.1:4000 wsgi:app
[2020-08-30 17:53:47 +0530] [13230] [INFO] Starting gunicorn 20.0.4
[2020-08-30 17:53:47 +0530] [13230] [INFO] Listening at: http://127.0.0.1:4000
(13230)
[2020-08-30 17:53:47 +0530] [13230] [INFO] Using worker: sync
[2020-08-30 17:53:47 +0530] [13233] [INFO] Booting worker with pid: 13233
[2020-08-30 17:53:48 +0530] [13234] [INFO] Booting worker with pid: 13234
[2020-08-30 17:53:48 +0530] [13235] [INFO] Booting worker with pid: 13235
[2020-08-30 17:53:48 +0530] [13236] [INFO] Booting worker with pid: 13236
```