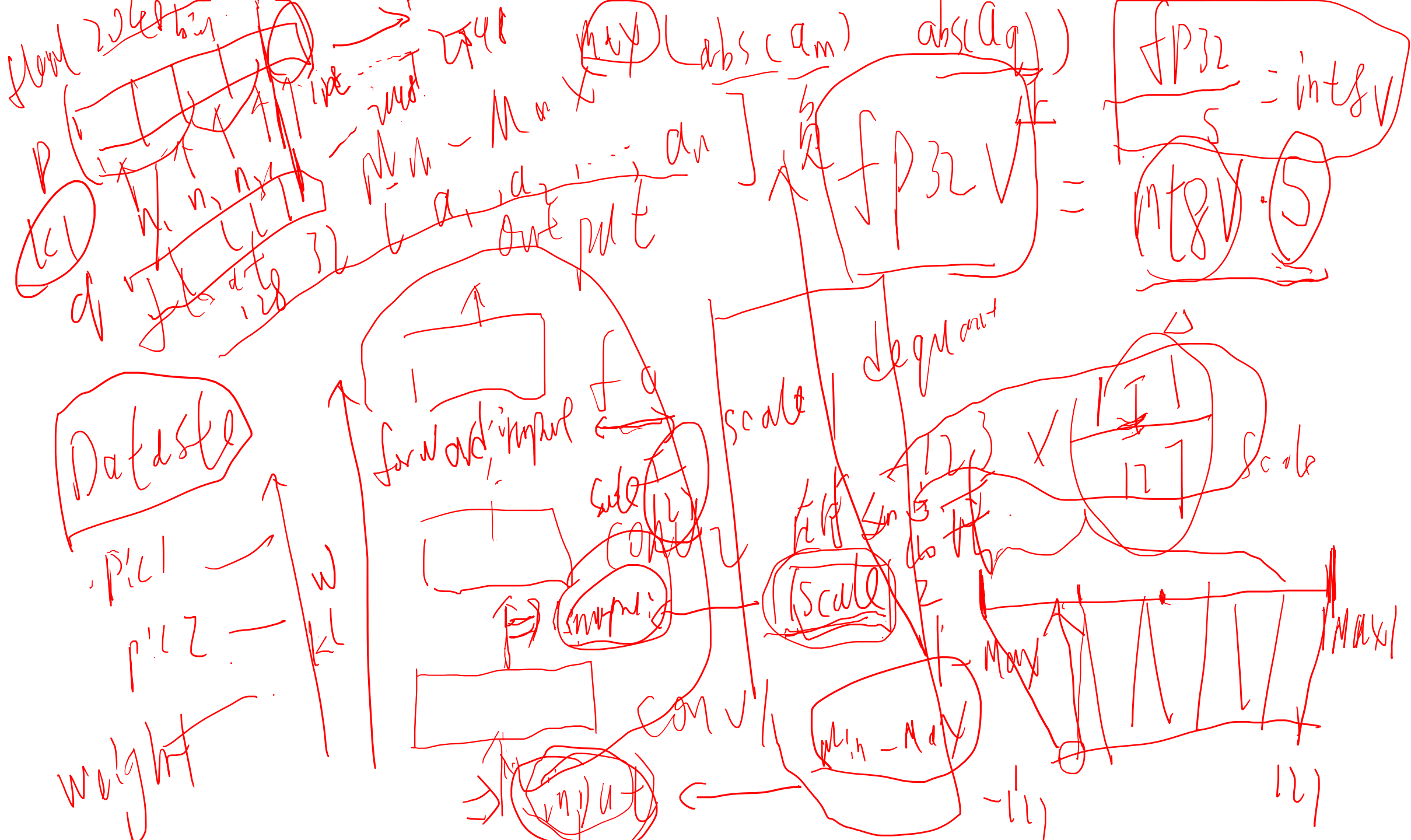


主要内容

- 量化方法MIN-MAX、ADMM量化
- TensorRT 的组网及使用



$$f_{PSLV} = \frac{\text{int of } V}{S} \quad (5)$$

Weight of PSLV

$$\frac{\text{int of } V}{S}$$

int of V

$$\text{int of } V = \frac{f_{PSLV} \cdot S}{S}$$

$$f(x) = \frac{x}{S}$$

$$f^{-1}(x) = S \cdot x$$

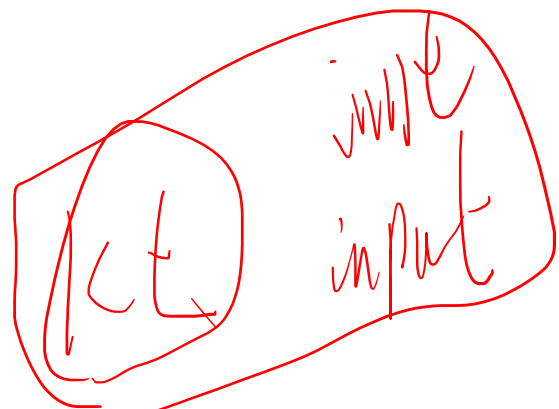
loss function

$$L = \left(S \cdot \bar{E}\left(\frac{x}{S}\right) - x \right)^2$$

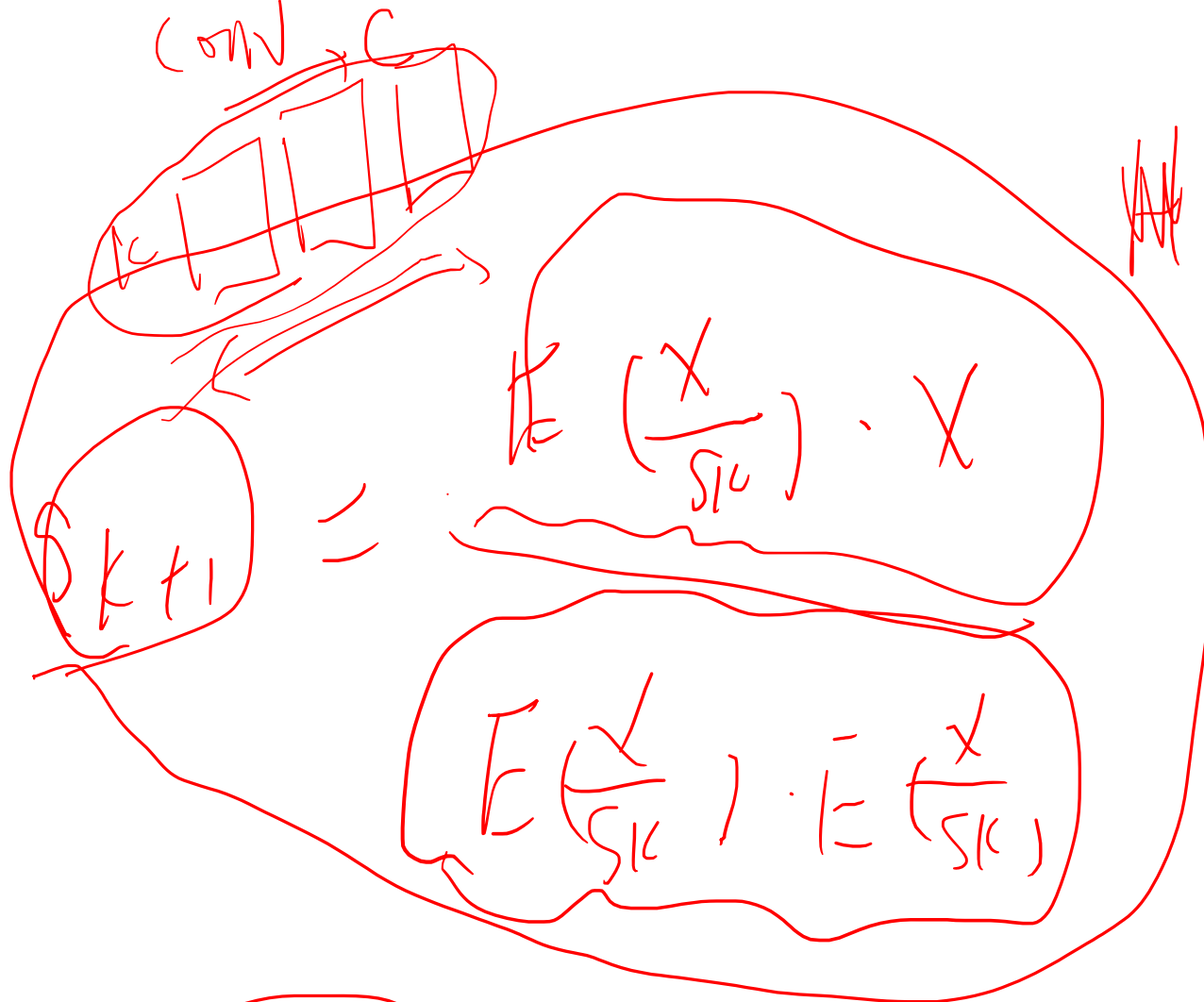
$$\frac{\partial L}{\partial S} = \bar{E}\left(\frac{x}{S}\right) \cdot \left[S \cdot \bar{E}\left(\frac{x}{S}\right) - x \right]$$

$$\bar{E}\left(\frac{x}{S}\right) \cdot \left(S_{k+1} \cdot \bar{E}\left(\frac{x}{S_{k+1}}\right) - x \right) = 0$$

A DMM

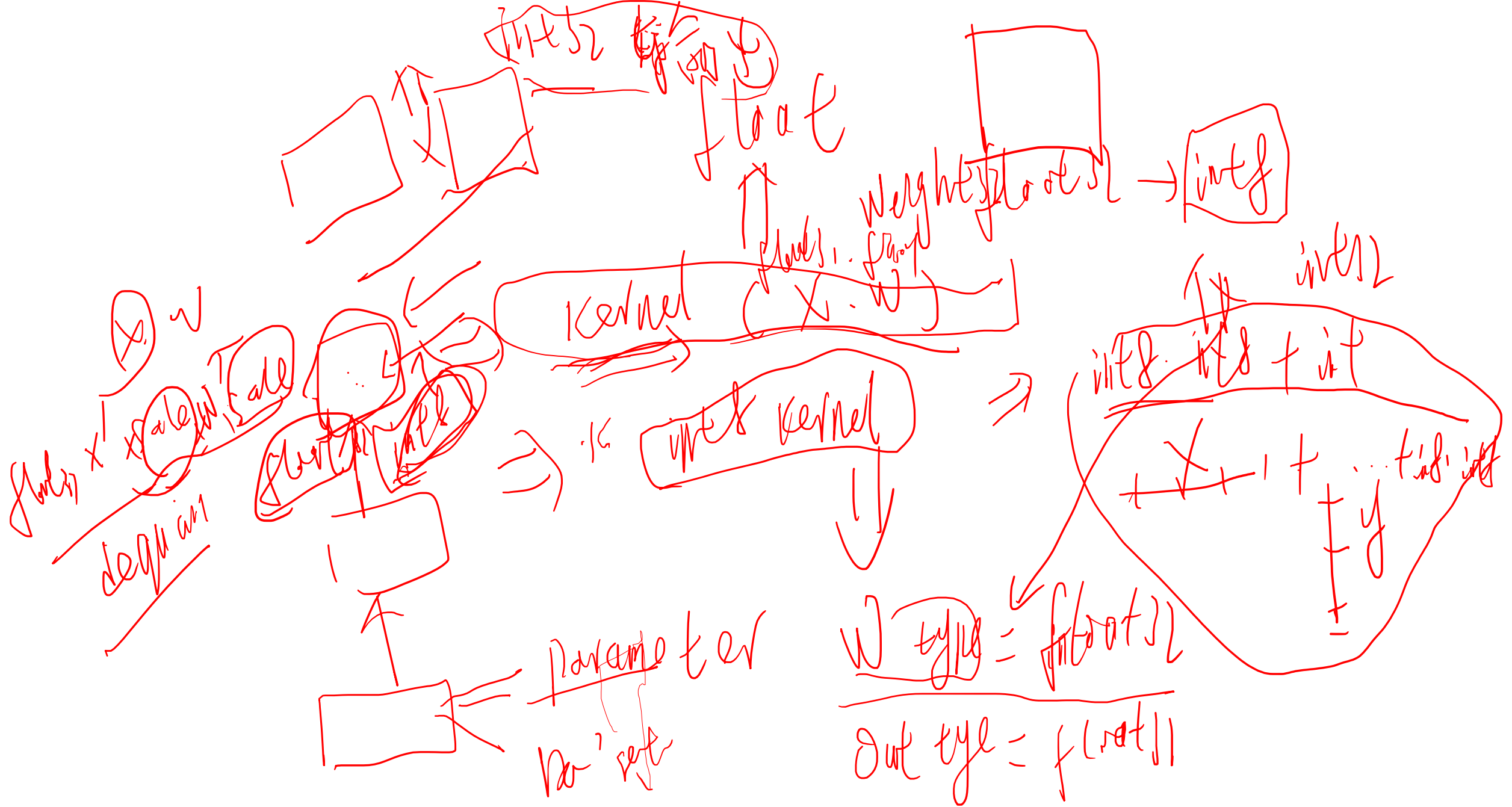


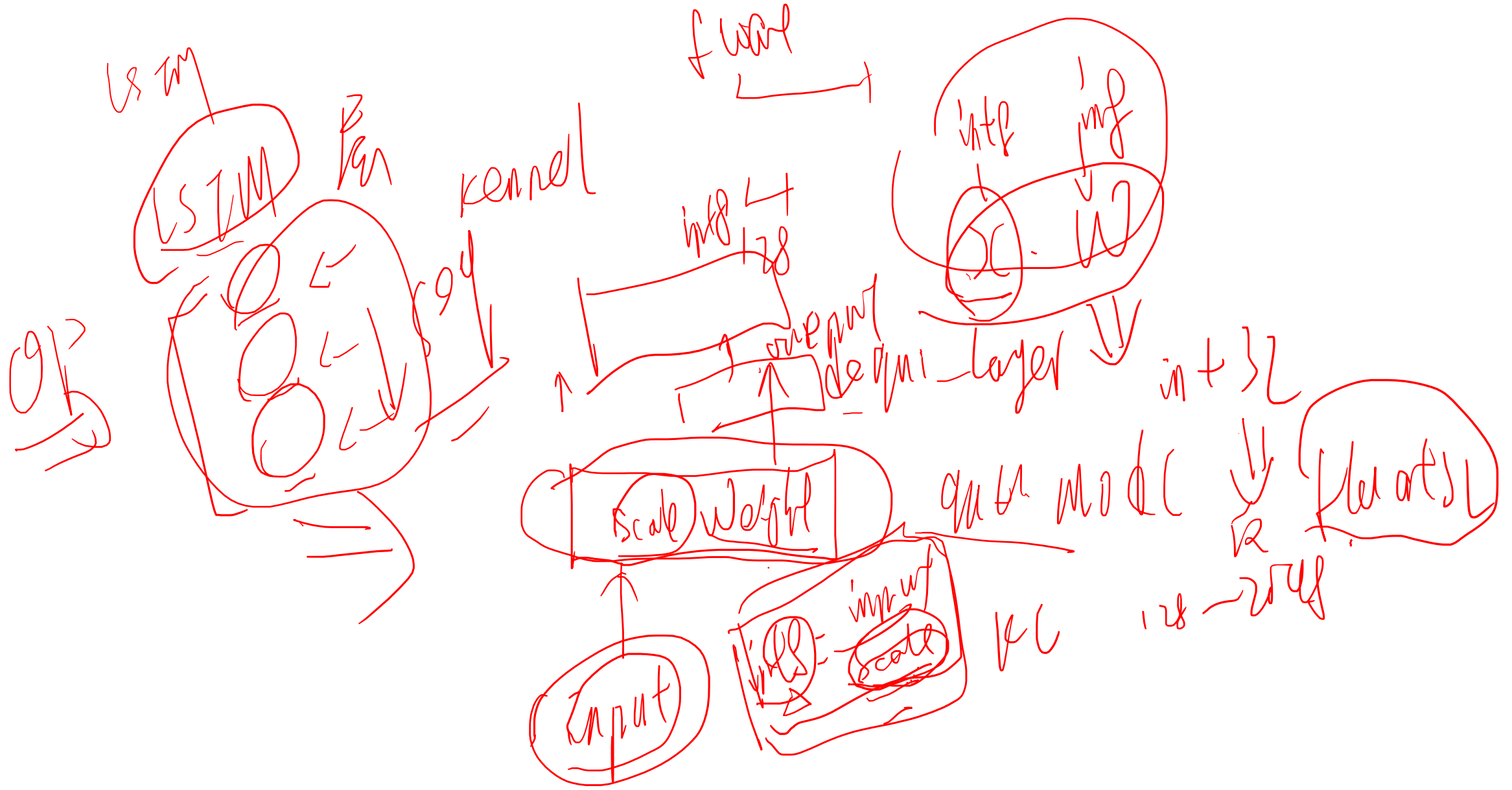
权重 w_{kt} :



S1k









build engine

low-level dev (philman)

Serial

onnx

uff

W / b
network



shared

input
output
W

shape

shape

shared mem

kernel

filter

动态核

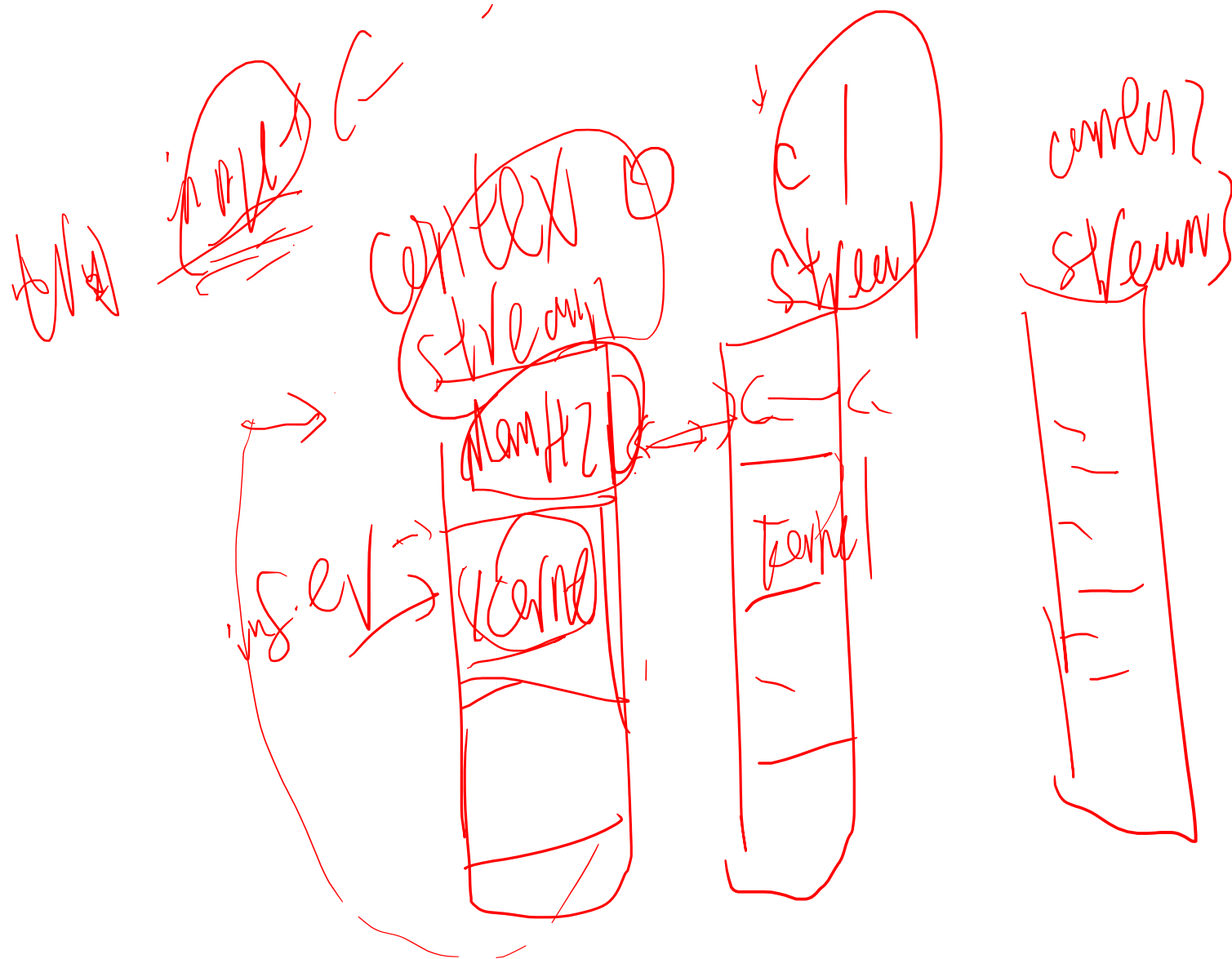
context

context

context

context

context



考勤&反馈

