

TRT

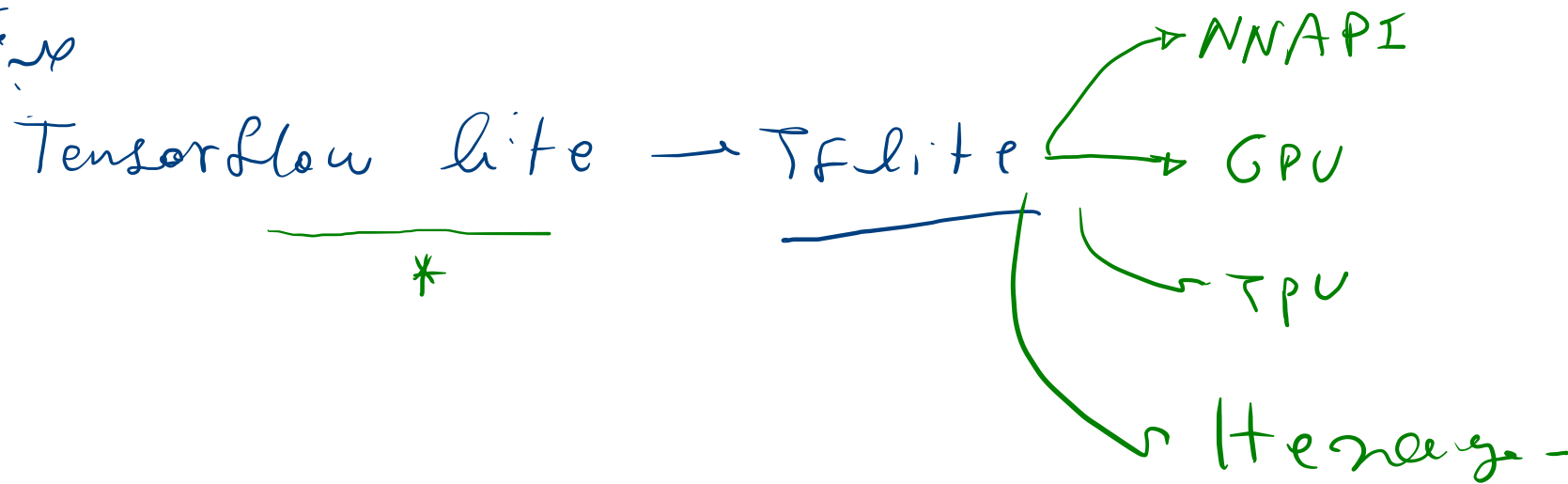


NVIDIA

/ purpose: NVIDIA
GPU

Tensor Run Time (TRT)

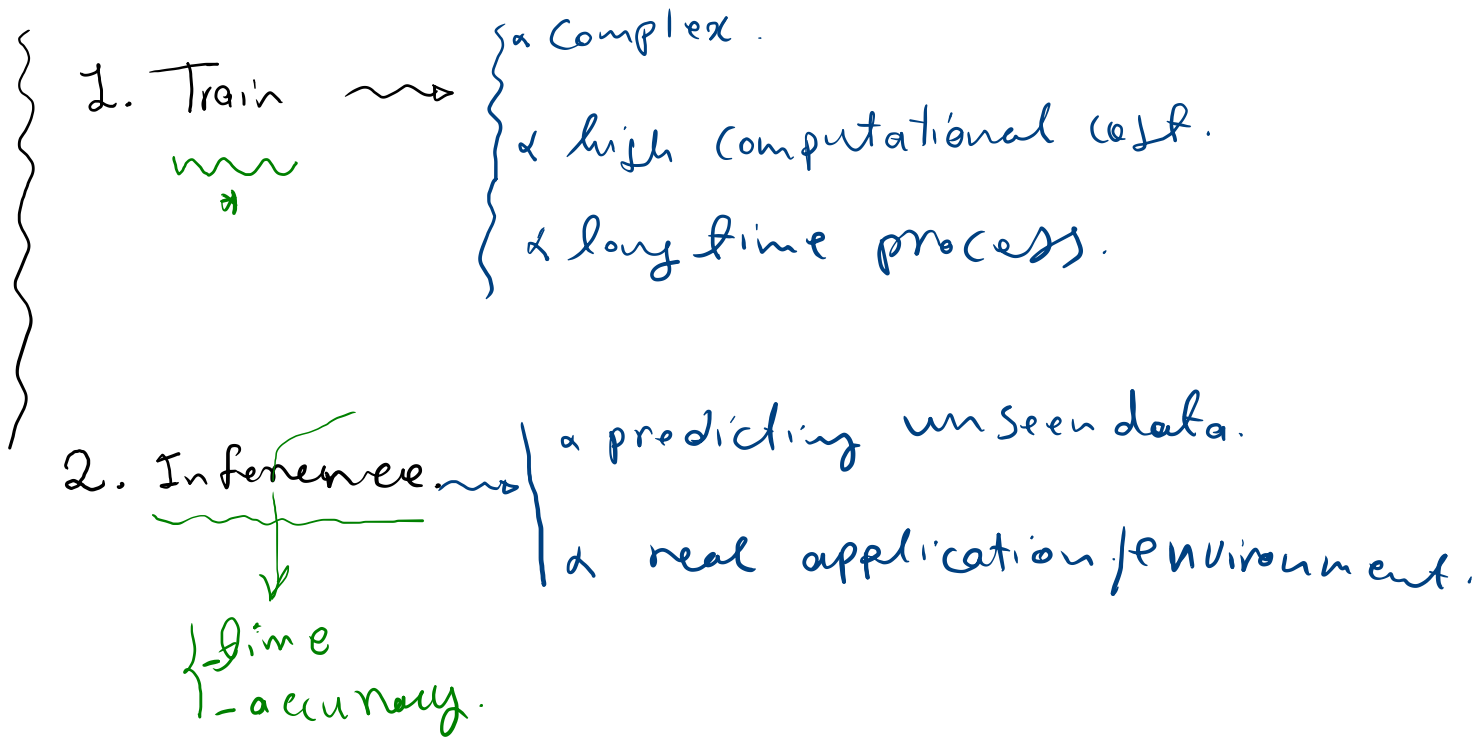
في



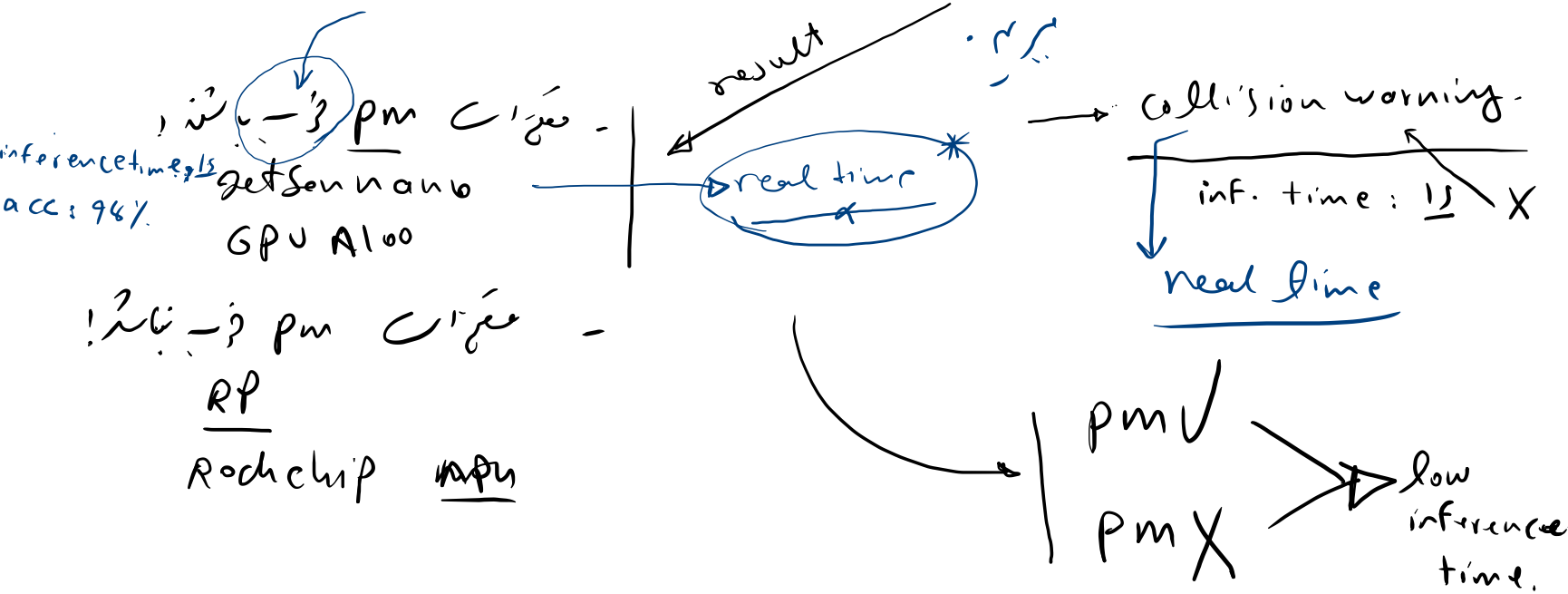
ORT
↓

optimization ~> real time applications.

DL →



آزمون Train شده را از رویکرد و مدل نتایج Inference



نتایج pm - 3 - 1

RP

Rochchip mpu

real time \longrightarrow response : $< \underline{40ms}$
 $< 10ms$
 $< 5ms$

limited resource $\xrightarrow{\text{must}}$ optimization.

imp. optimization - libs \rightarrow TKLITE
 \swarrow TRT
x

Fp32



int8 dynamic range.

-127 & +127

Fp
~> 19,23487 → 4 x 19,23487 = 76,8 → cat

↓
int

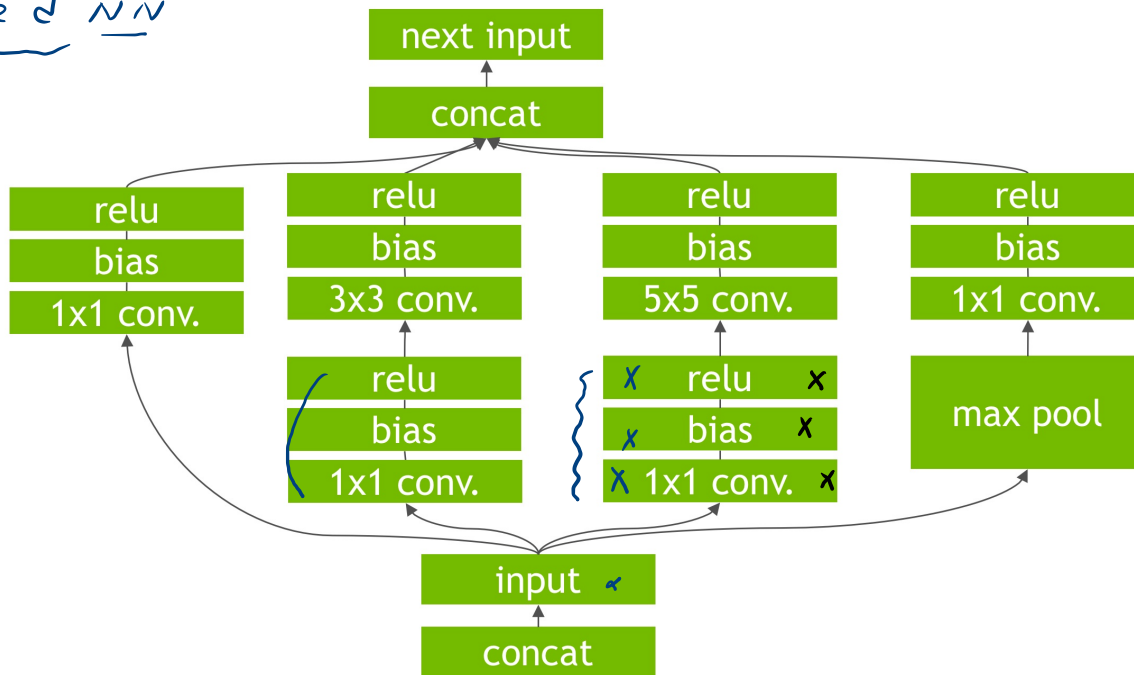
19

→ 4 x 19 → 76 → dog

1. Layers Fusion.

TRT

unoptimized nn



$$p1 \rightarrow x = 8$$

$$p2 \rightarrow y = \underset{x}{x}^2 + 12$$

$$p3 \rightarrow \underset{x}{x} = \underset{y}{y} + \underset{x}{x} + 1$$

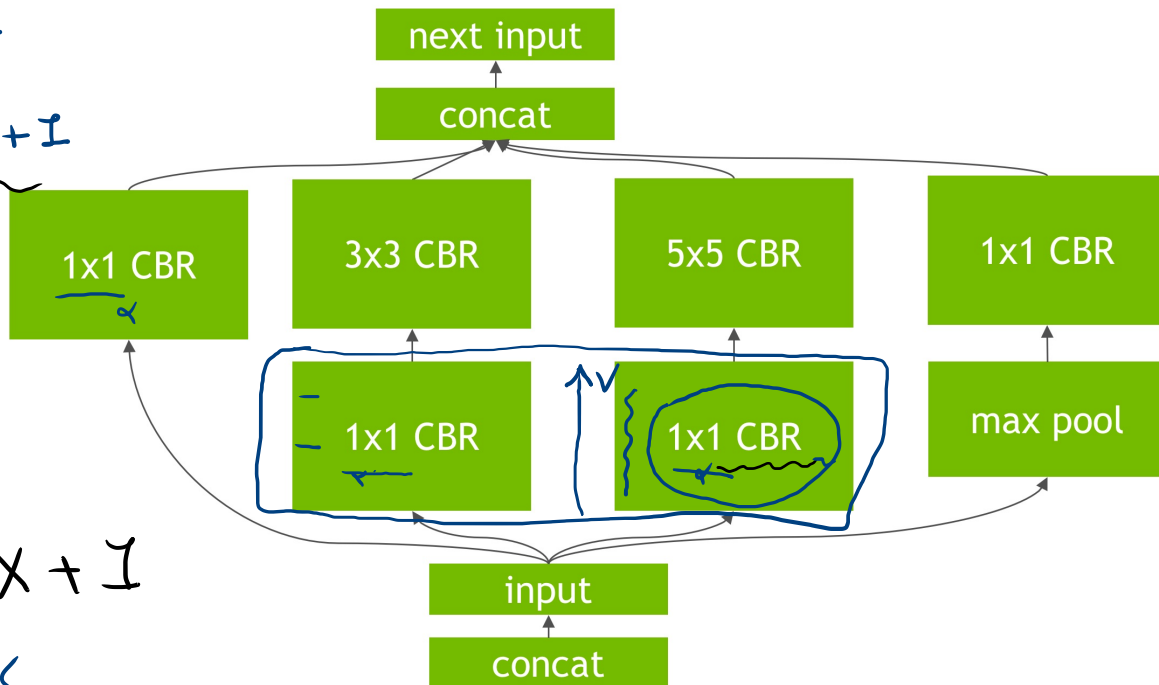
حفظ x

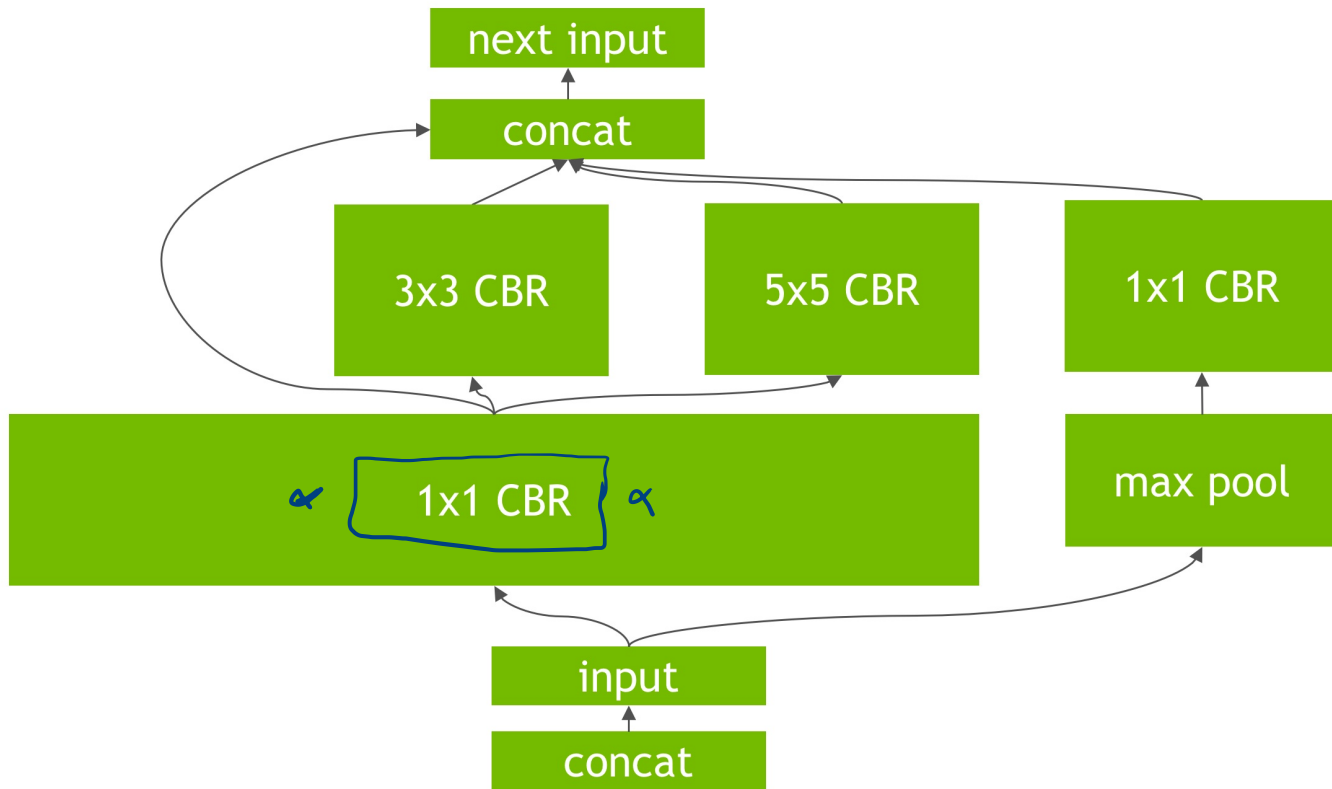
حفظ x ، y ، z

$$z = x^2 + 12 + x + 1$$

حفظ x

حفظ x ، y ، z





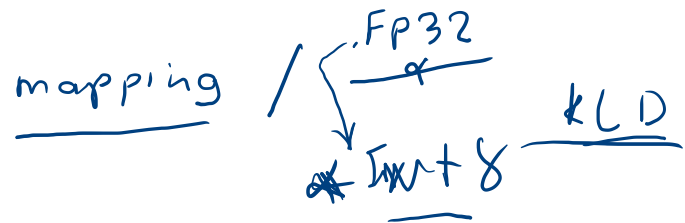
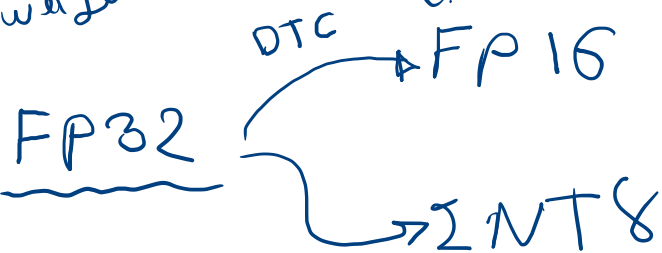
layer fusion

network	before	after
VGG19	<u>43</u>	27
Inception V3	309	113
Resnet 152	670	<u>159</u>

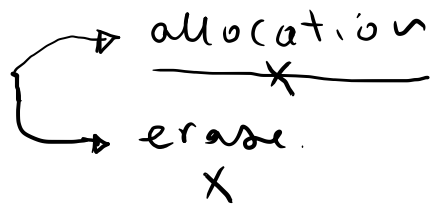
2. pre cision calibration

↓ LPR: low precision inference.

weights, bias, activation
Func.

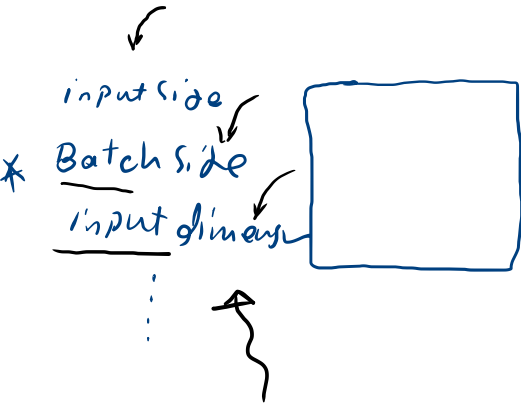


3. Dynamic Tensor memory allocation!



low inference time

4. Kernel auto tuning!



TRT \rightsquigarrow other ~~parameters~~ parameters
opt.

× ResNet 18

1.69
s



PT :	$\frac{1X}{5.5}$	$\frac{FP32}{\times}$
TRT	$\frac{FP32}{\times}$	→ <u>1.5X</u>
TRT	FP16	→ 5.5X

× Real Time

* Tensor R T α $\left\{ \begin{array}{l} \rightarrow \text{C++} \quad \alpha \\ \rightarrow \text{python} \quad \alpha \end{array} \right.$

{ pytorch \leftarrow
Tensorflow \leftarrow

*
Resnet 50.

α



image Net

