

Interview Session

CNNs

(Q)

CNN



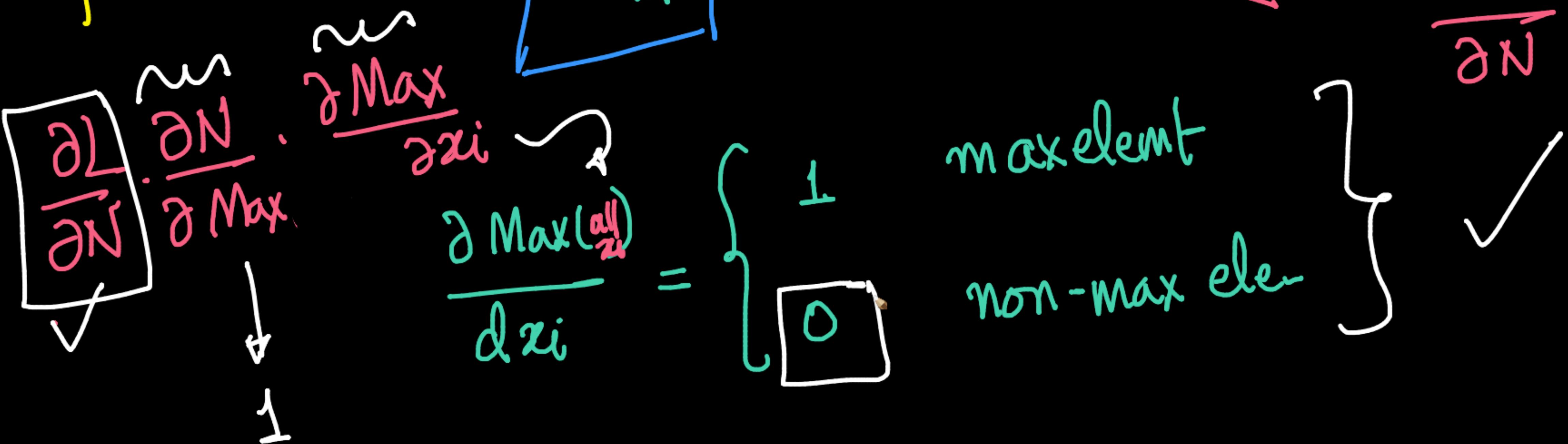
#params (straight forward)

(Q1)

Maxpool  $3 \times 3 \rightarrow$  no params

How does back-prop over  
on Maxpool layer

Input - . . .



\* \* \*  
Lesson!  
==

operation

conv

Max pool

Dropout

Residual --

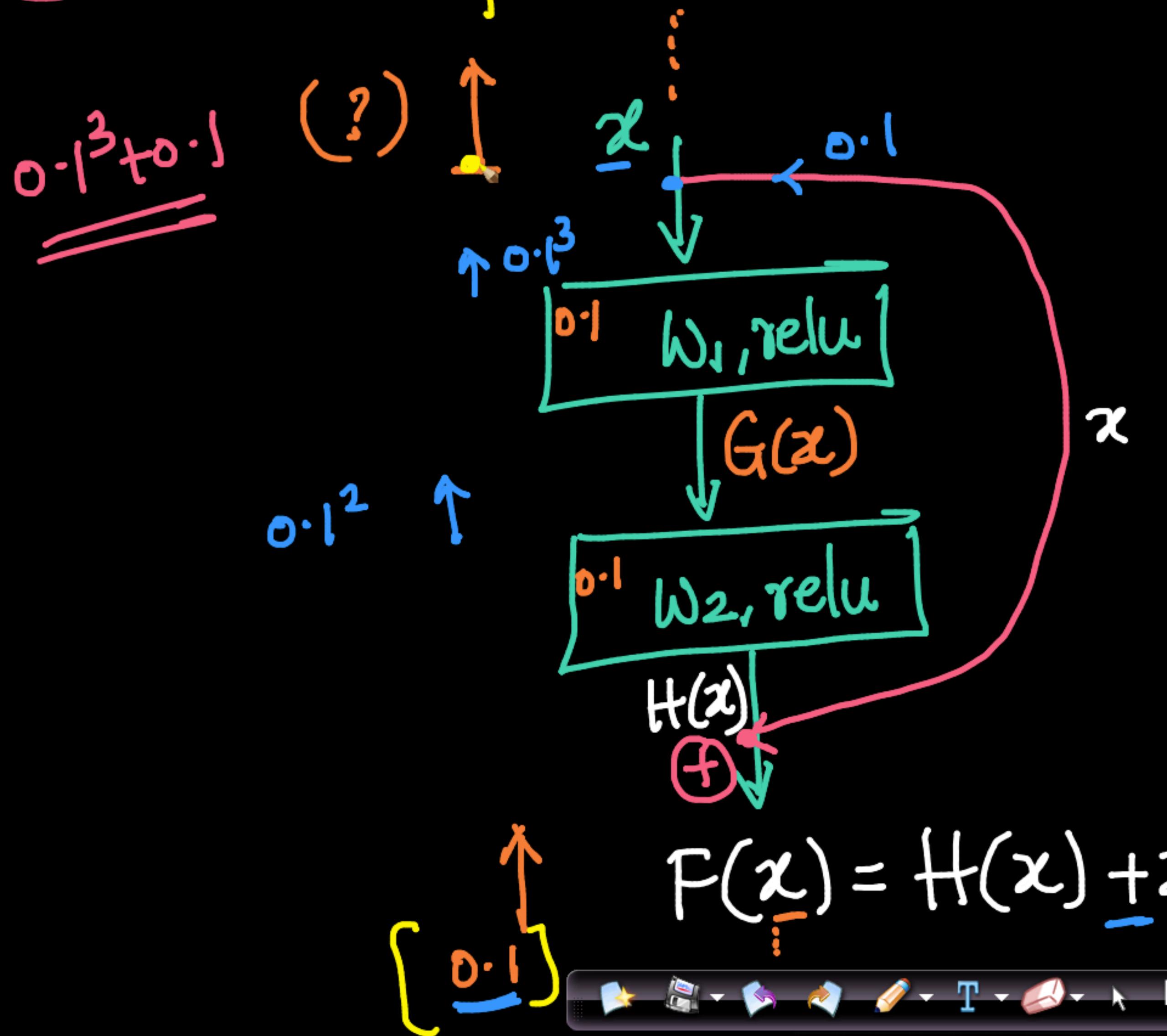
backprop ...

Transformer. Module/block

Q

skip-connection:

back-prop



a

$$0.1 \times 0.1 \times 0.1$$

b

$$1 + 0.1 \times 0.1 \times 0.1$$

c

$$0.3$$

d

$$0$$

e random value

f

$$0.2$$

g

$$\underline{0.101}$$

g 0.11  
0.001

h 0.1

$$\frac{\partial H}{\partial G} \cdot \frac{\partial G}{\partial x}$$

$$\frac{\partial L}{\partial f}$$



0.1  
(given).

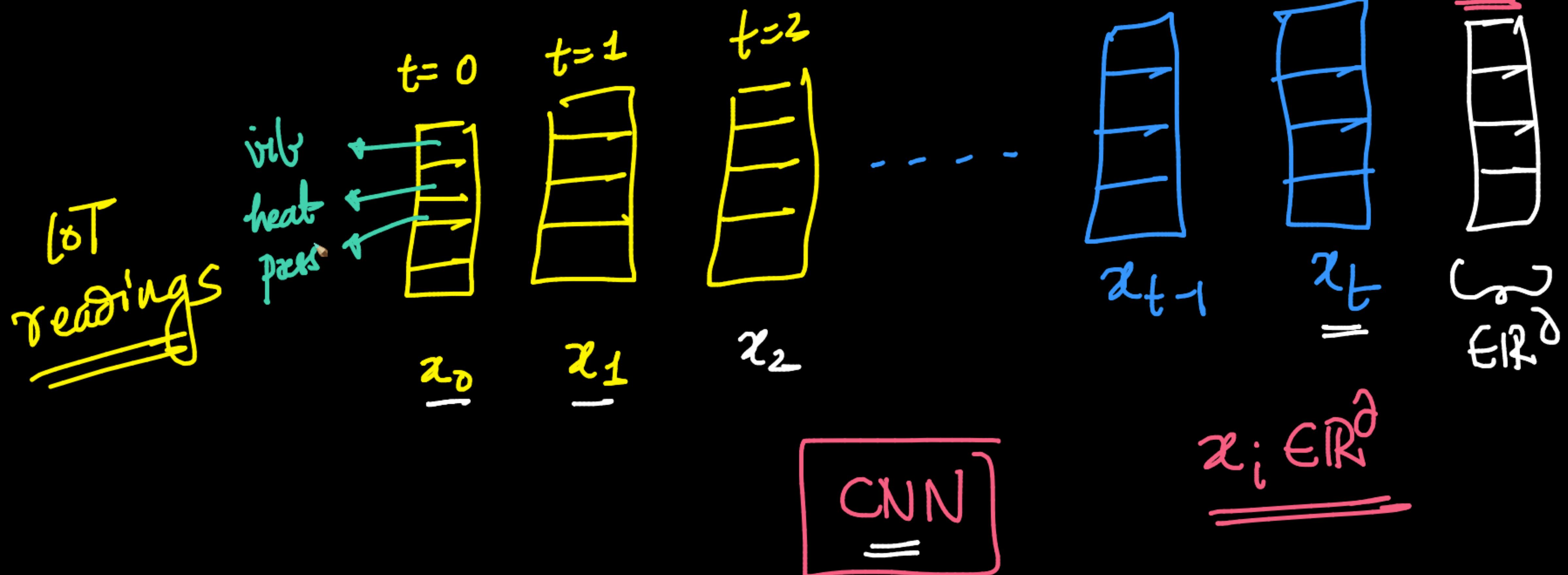
$$\frac{\partial F}{\partial x} = \left[ \frac{\partial H}{\partial \bar{a}} + \frac{\partial x}{\partial x} \right] \frac{\partial L}{\partial F}$$

$$(0.1^2 + 1) \stackrel{0.1}{=}$$

$$0.1^3 + 0.1 = 0.10$$

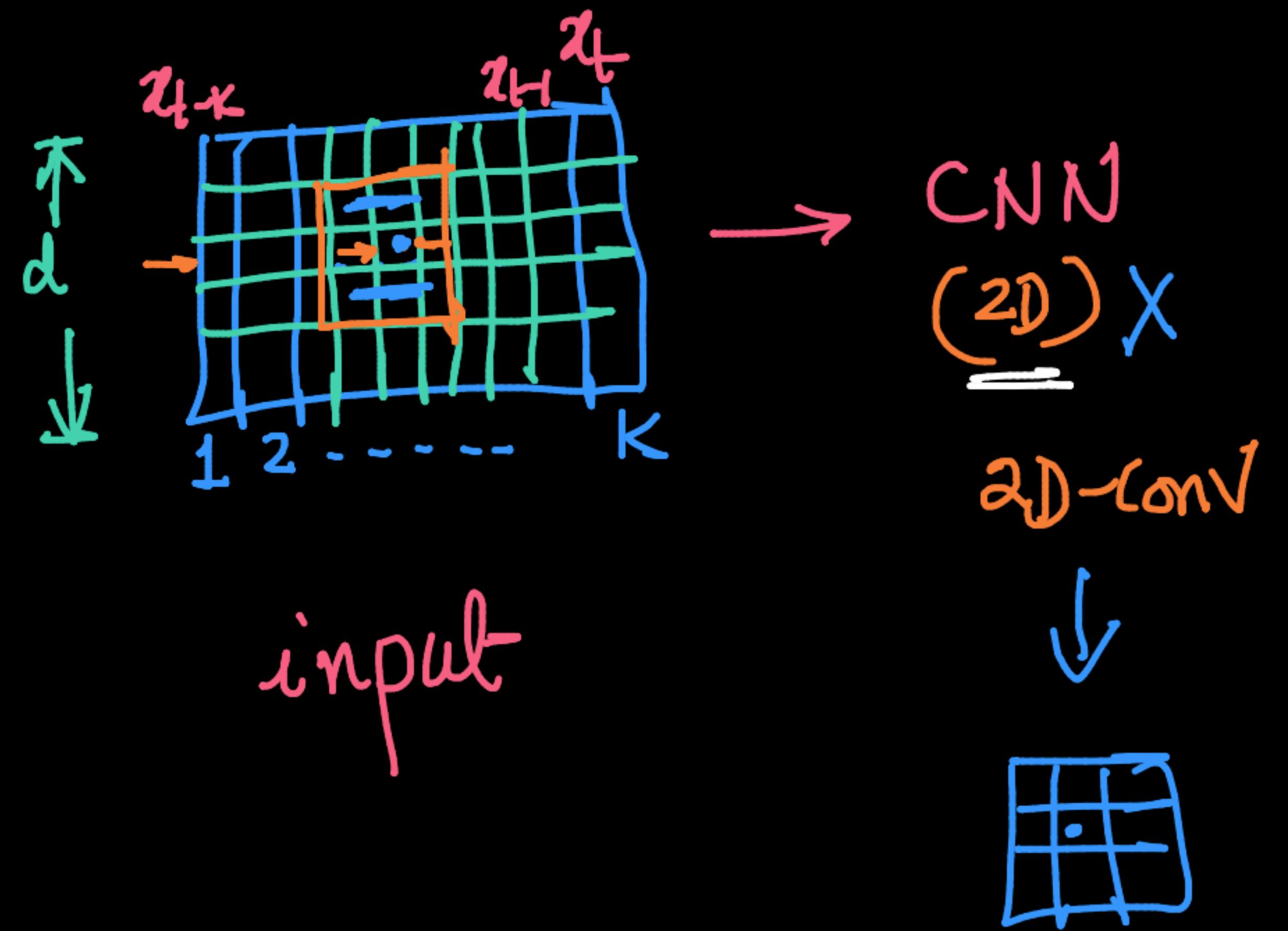
(Q3)

## Time-Series - forecasting

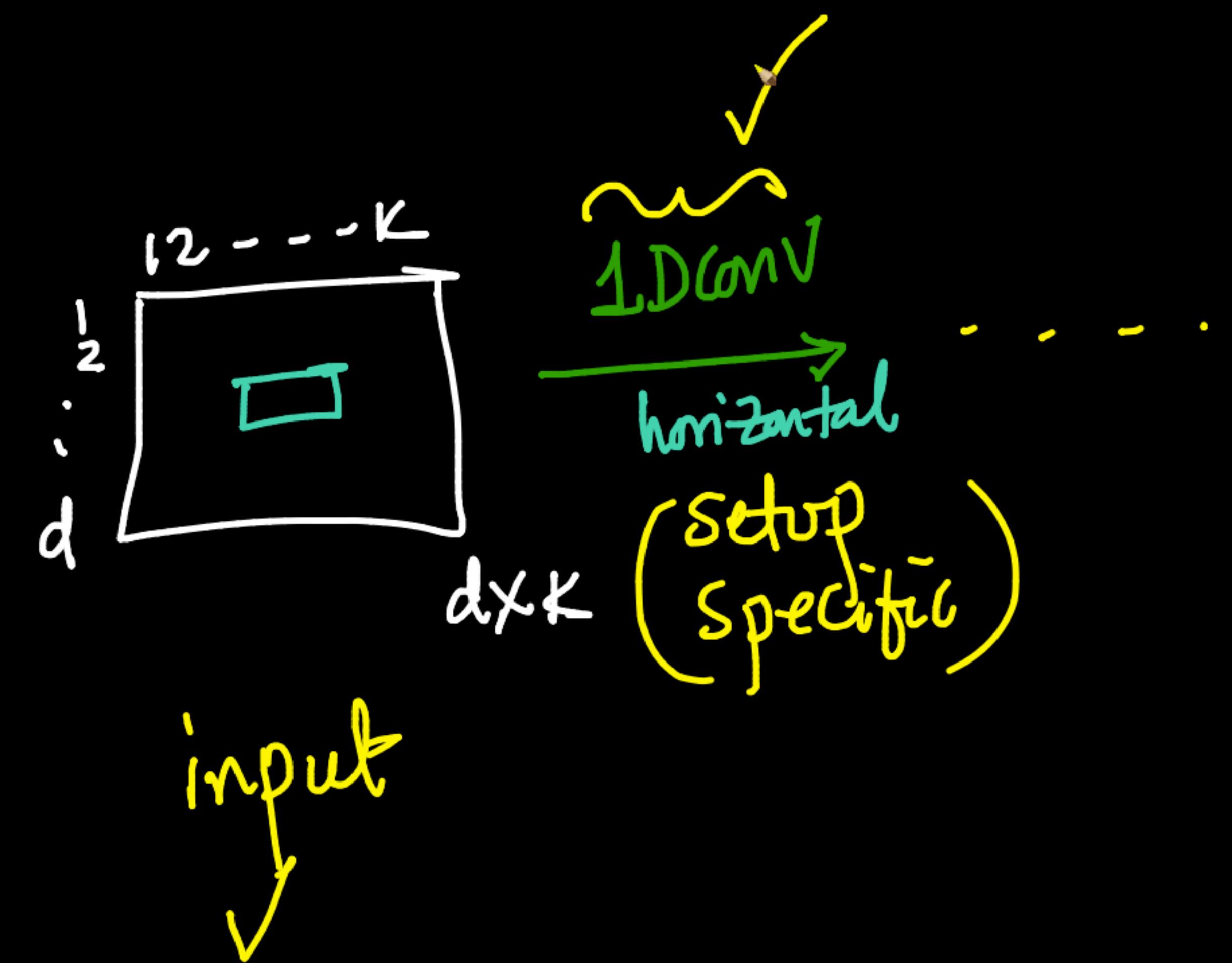


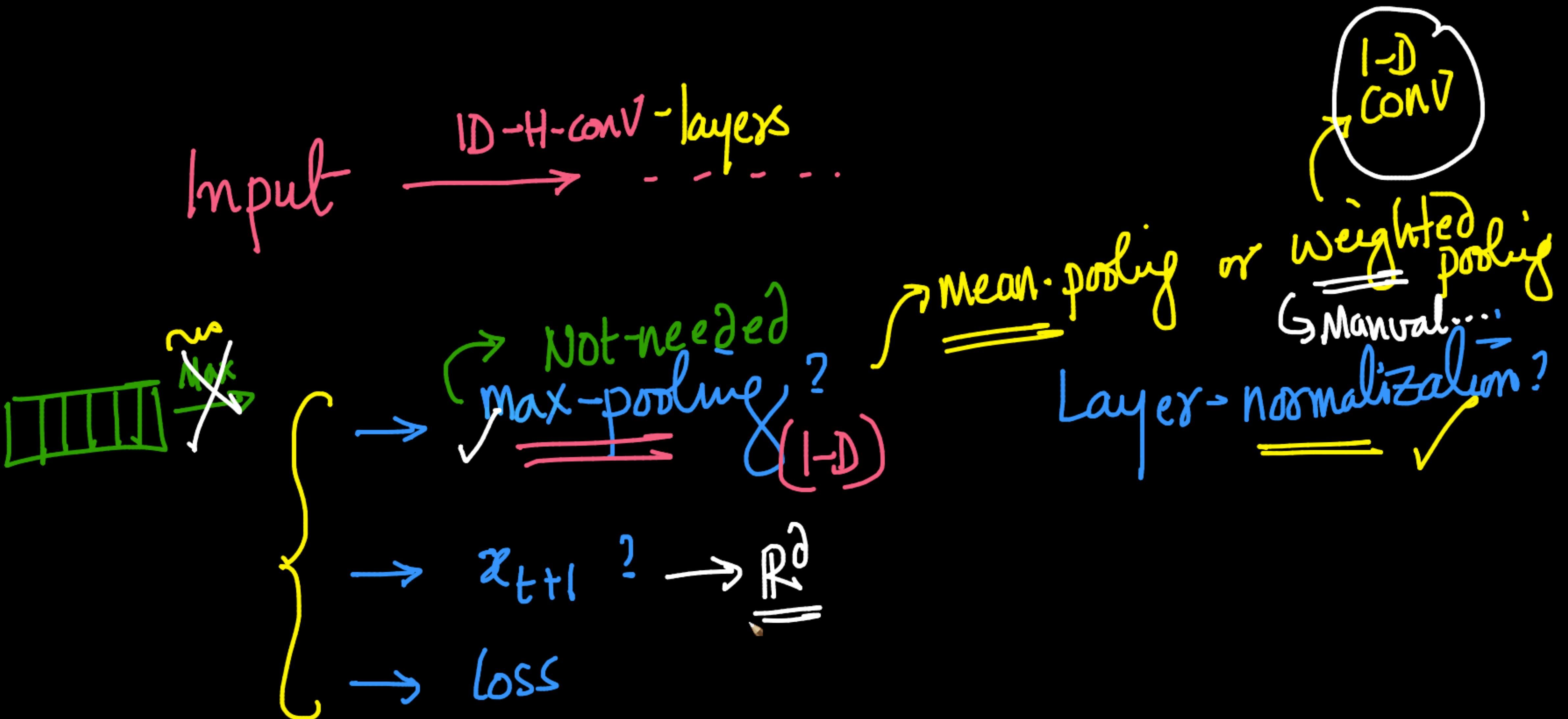
AI

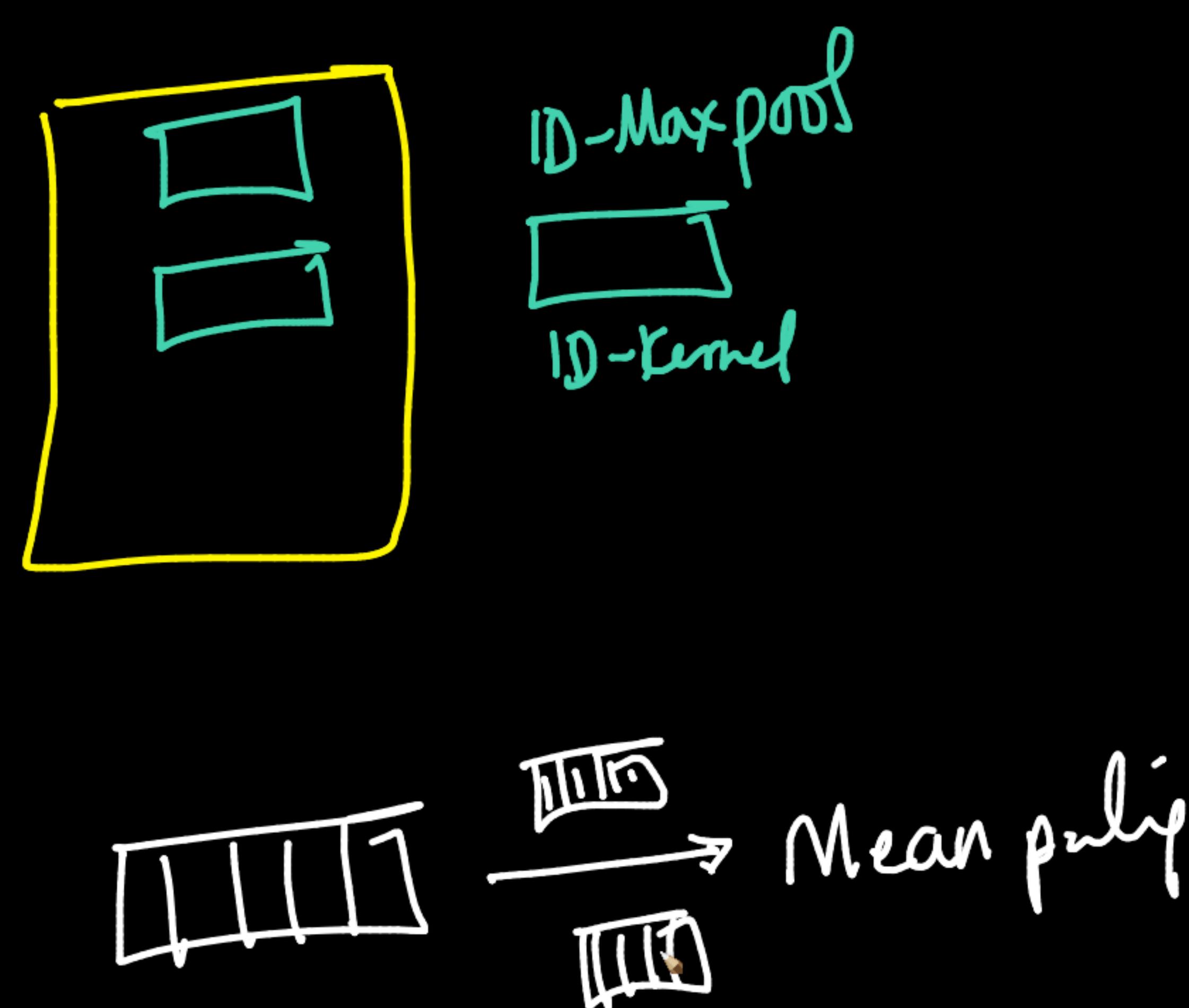
@ time  $t$ , last  $K$  time stamps → move window

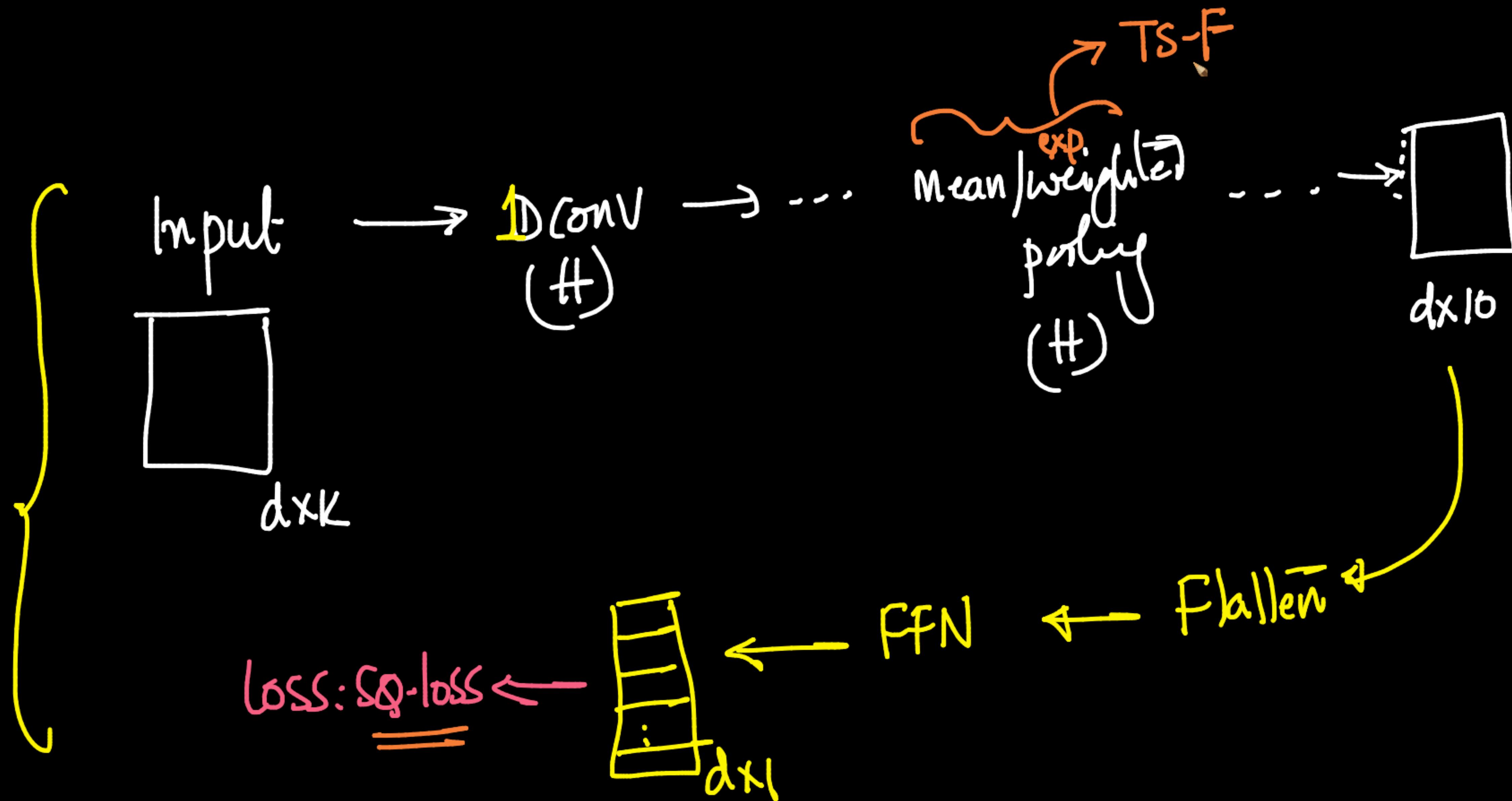


2D-conv



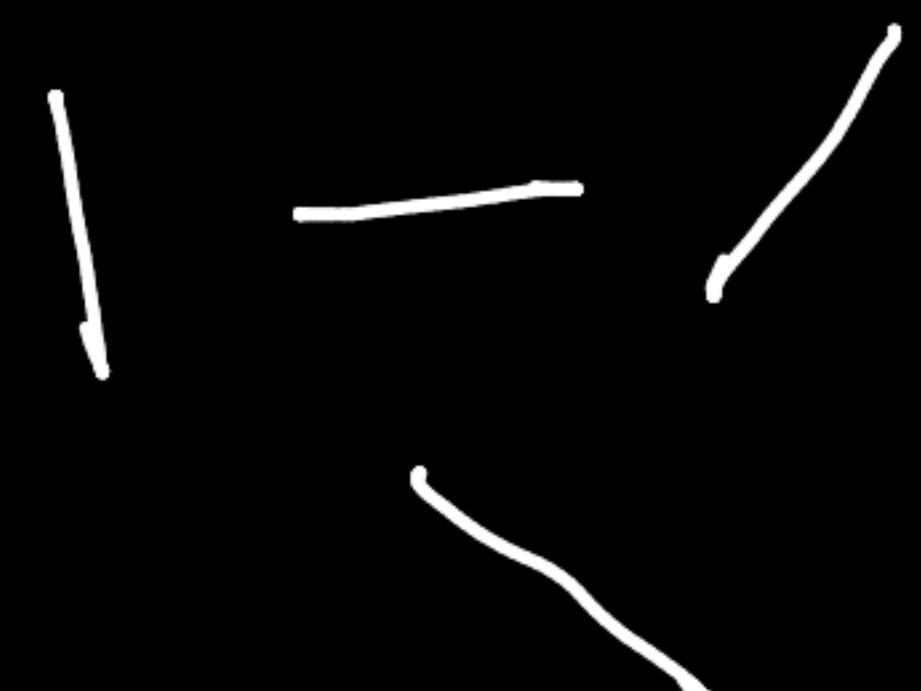
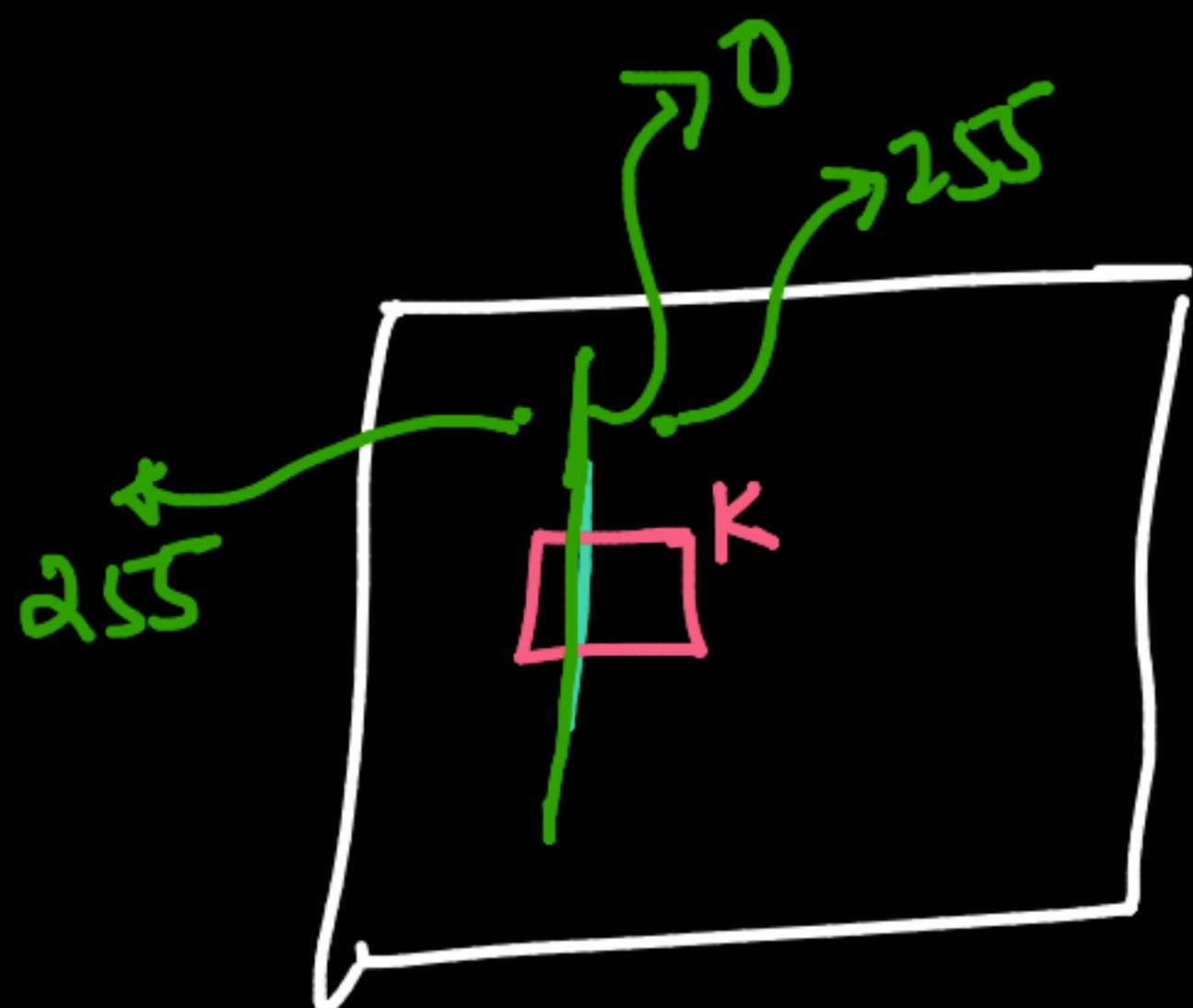
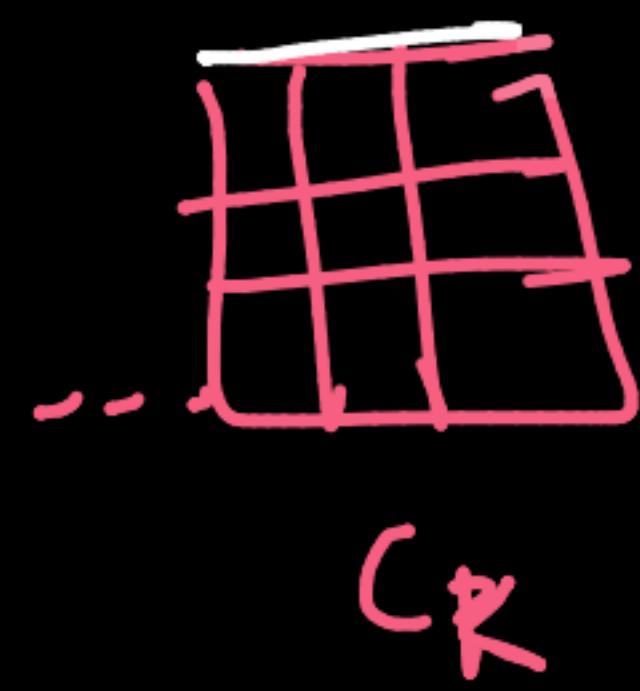
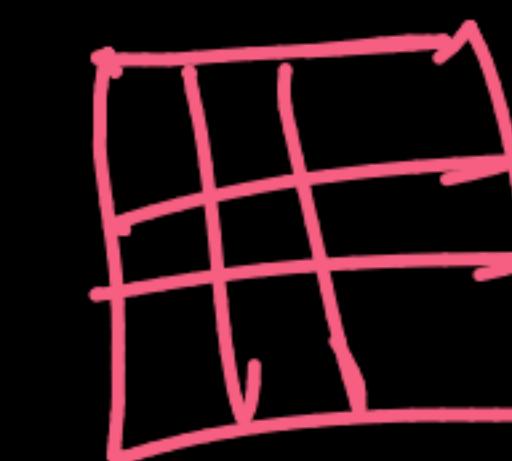
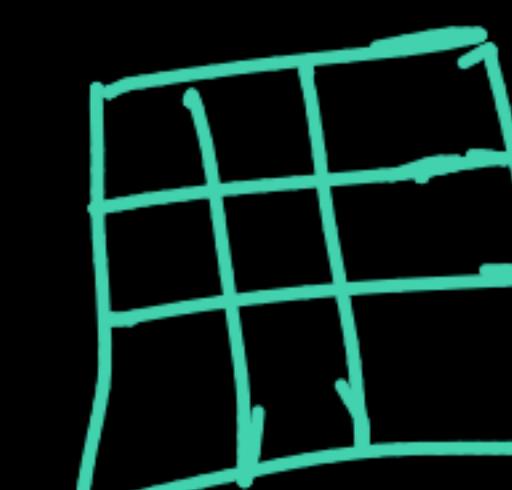






(Q)

Image - CNN

1st-layer - convolutions  $3 \times 3$ Specific  
shapeto look like edge  
detectors

built  
CNN model

Task:

s.t. 1st layer conv

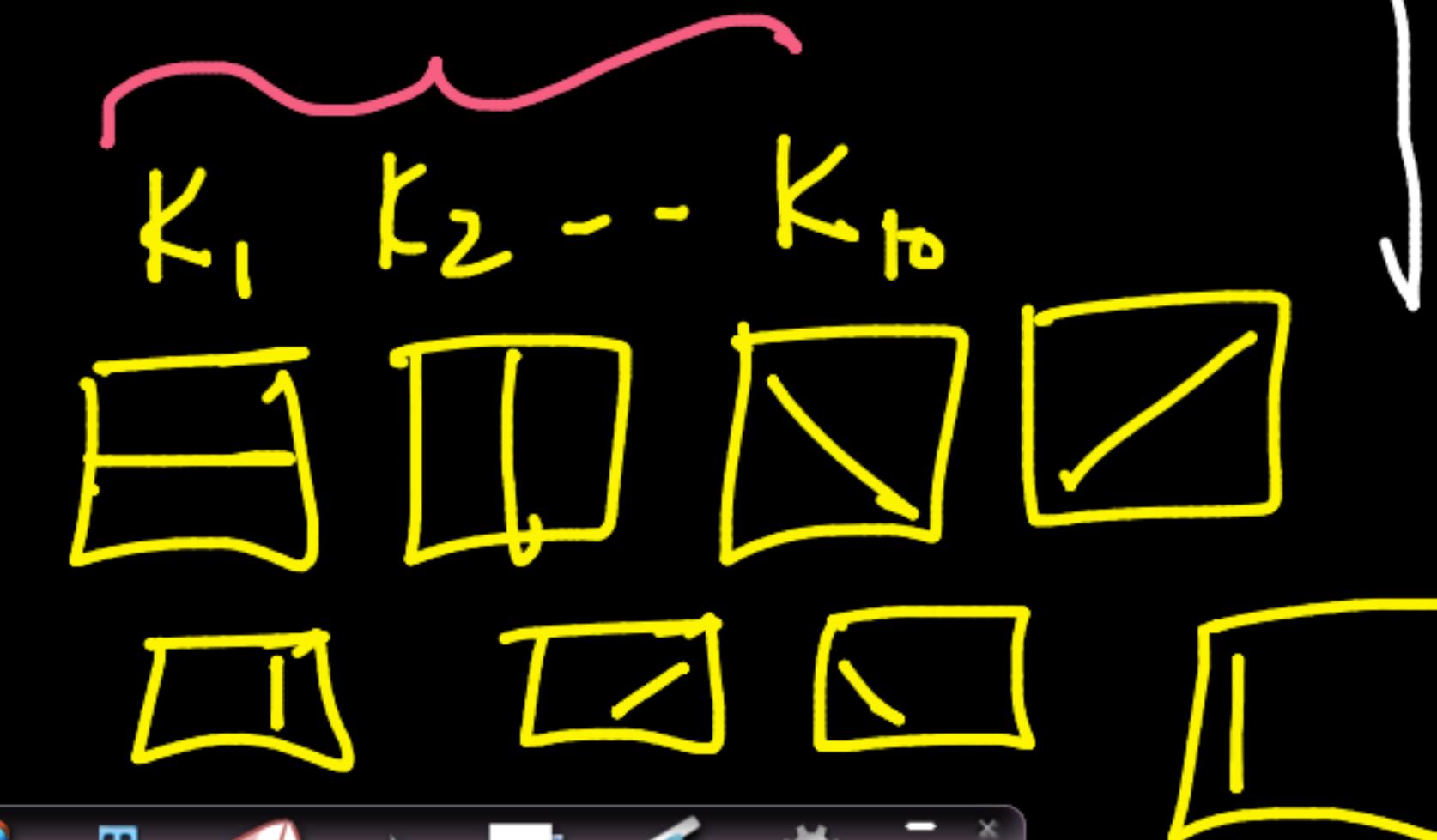
look like edge  
detectors always

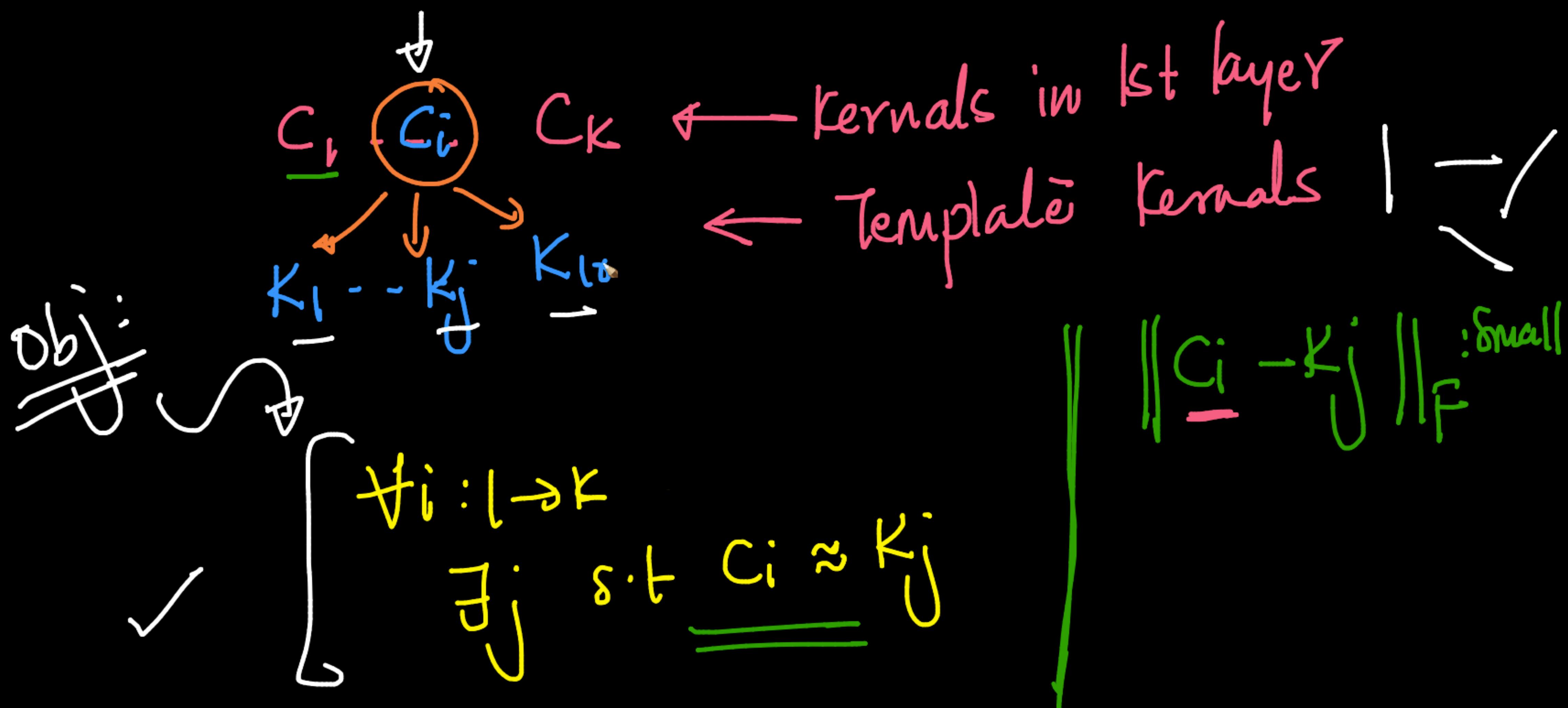
X (a) inc size of kernel

? (b) remove noise in image

Good Idea (c) initialize kernels to be  
edge-detectors

(d) template kernels:





min  $\theta$

all Kernels

Loss +  $\lambda$  reg

$\sum_i \sum_j [ \|c_i - k_j\|_F^2 ]$

min

Kernels in 1st layer

Back Prop

$L_M$

Diagram illustrating the optimization process for learning kernels. The overall loss function is given by  $\text{Loss} + \lambda \text{reg}$ , where  $\lambda$  is a regularization parameter. The first term,  $\sum_i \sum_j [ \|c_i - k_j\|_F^2 ]$ , represents the Frobenius norm of the difference between feature vectors  $c_i$  and  $k_j$ . This term is highlighted with a green bracket and labeled "Kernels in 1st layer". The second term,  $\lambda \text{reg}$ , is highlighted with a red bracket and labeled "Back Prop". A green bracket labeled "Back Prop" also encloses the entire loss function. The optimization is performed over all kernels  $\theta$ .



nearest to it. These  $\mathbf{x}_i$  are called *support vectors*.

## Soft-margin [edit]

To extend SVM to cases in which the data are not linearly separable, the *hinge loss* function is helpful

$$\max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)).$$

Note that  $y_i$  is the  $i$ -th target (i.e., in this case, 1 or -1), and  $\mathbf{w}^T \mathbf{x}_i - b$  is the  $i$ -th output.

This function is zero if the constraint in (1) is satisfied, in other words, if  $\mathbf{x}_i$  lies on the correct side of the margin. For data on the wrong side of the margin, the function's value is proportional to the distance from the margin.

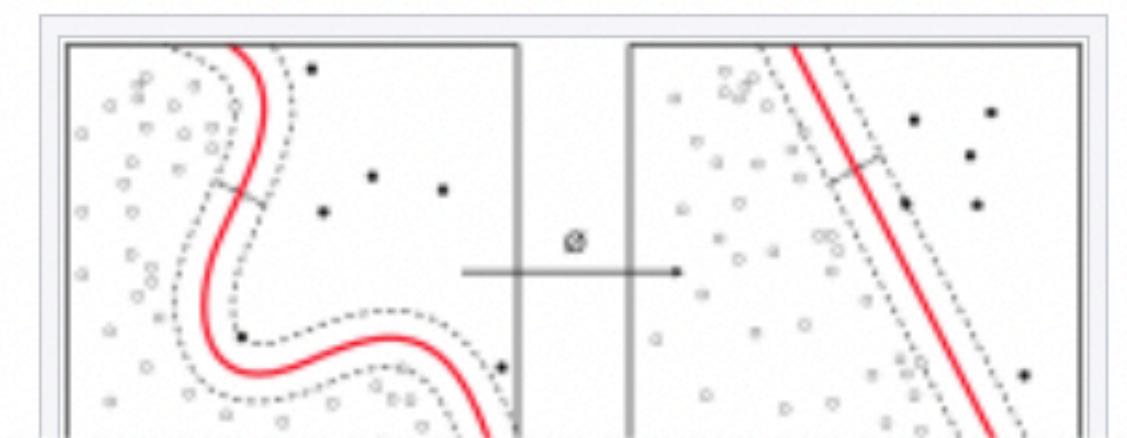
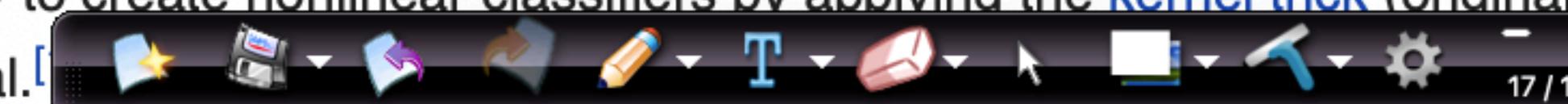
The goal of the optimization then is to minimize

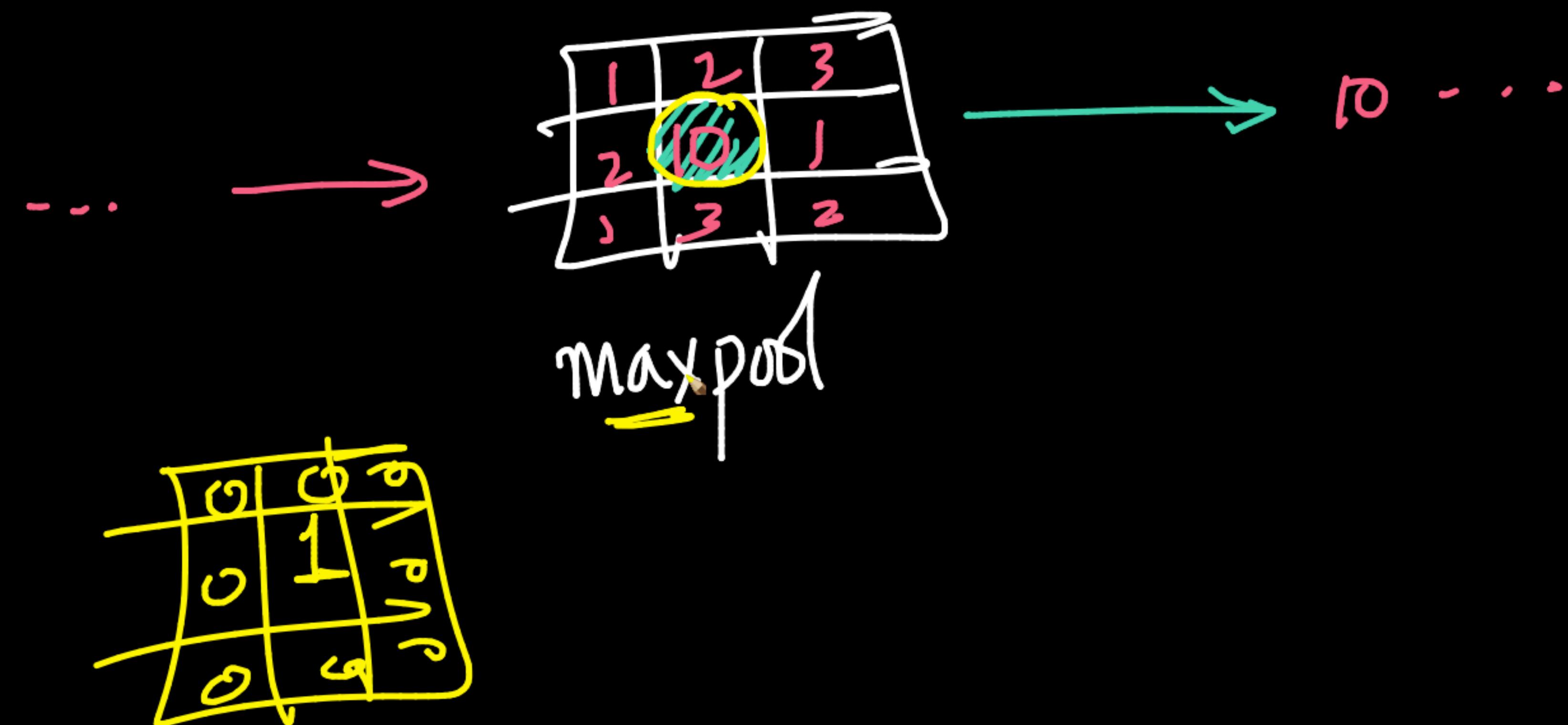
$$\lambda \|\mathbf{w}\|^2 + \left[ \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)) \right],$$

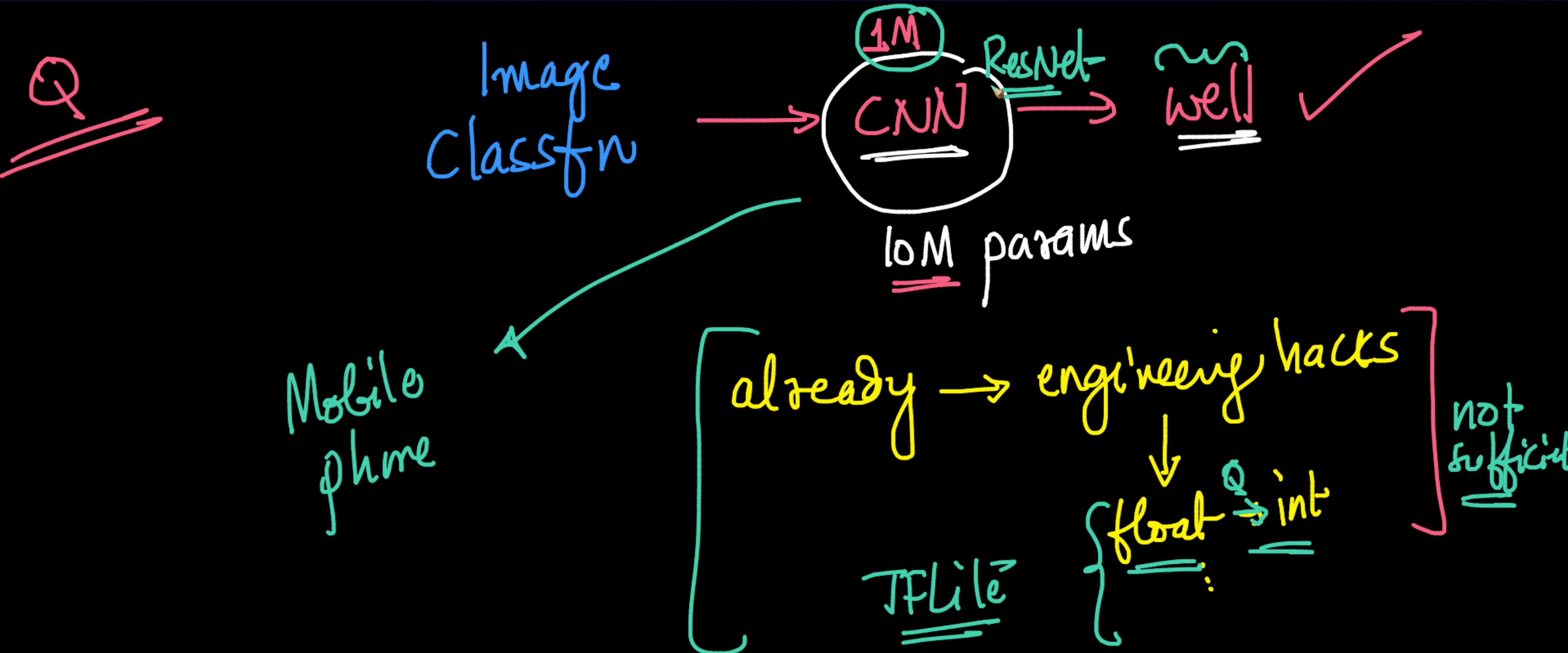
where the parameter  $\lambda > 0$  determines the trade-off between increasing the margin size and ensuring that the  $\mathbf{x}_i$  lie on the correct side of the margin. Thus, for sufficiently small values of  $\lambda$ , it will behave similar to the hard-margin SVM, if the input data are linearly classifiable, but will still learn if a classification rule is viable or not. (This parameter  $\lambda$  is often also called  $C$ , e.g. in LIBSVM, but usually refers to the inverse of  $\lambda$ .)

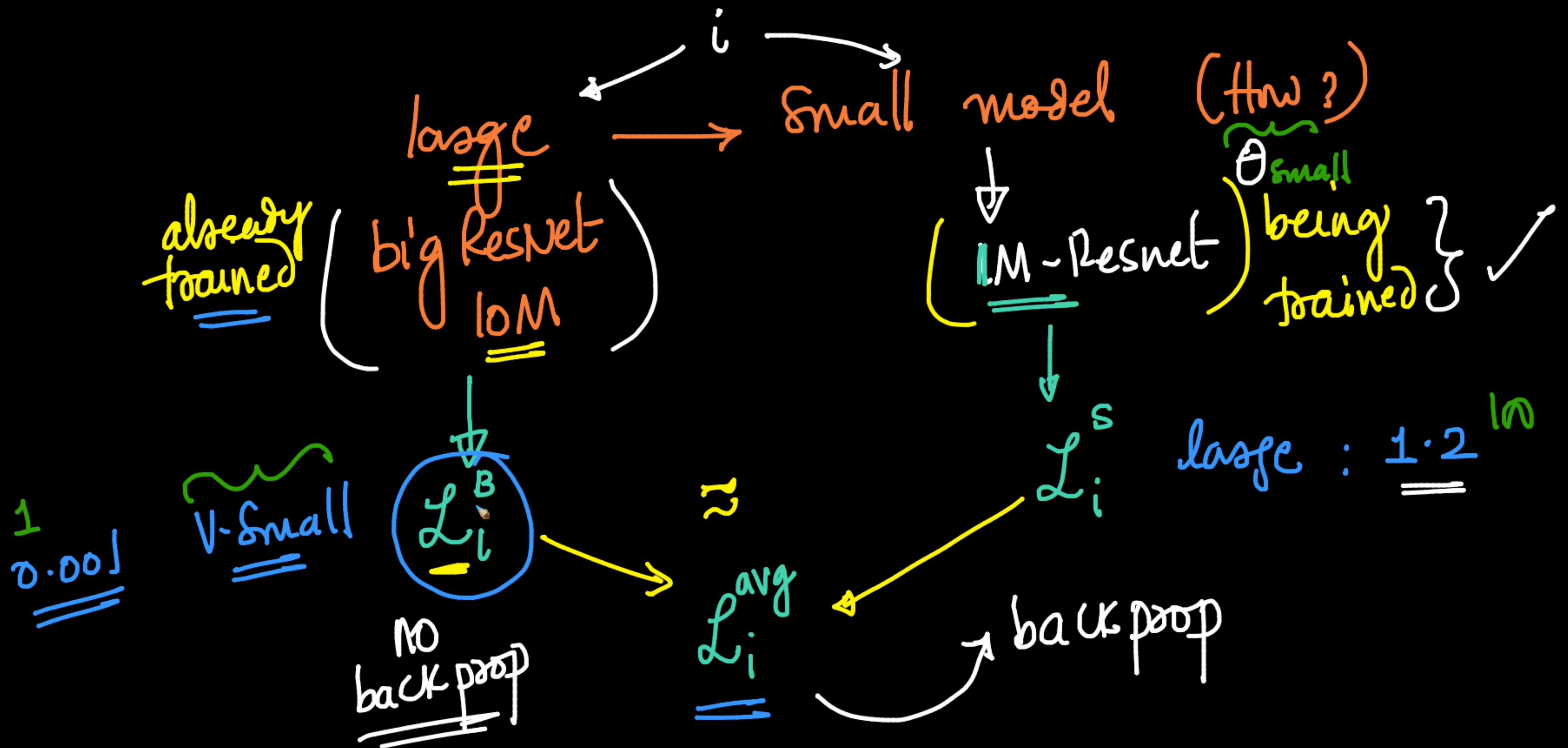
## Nonlinear Kernels [edit]

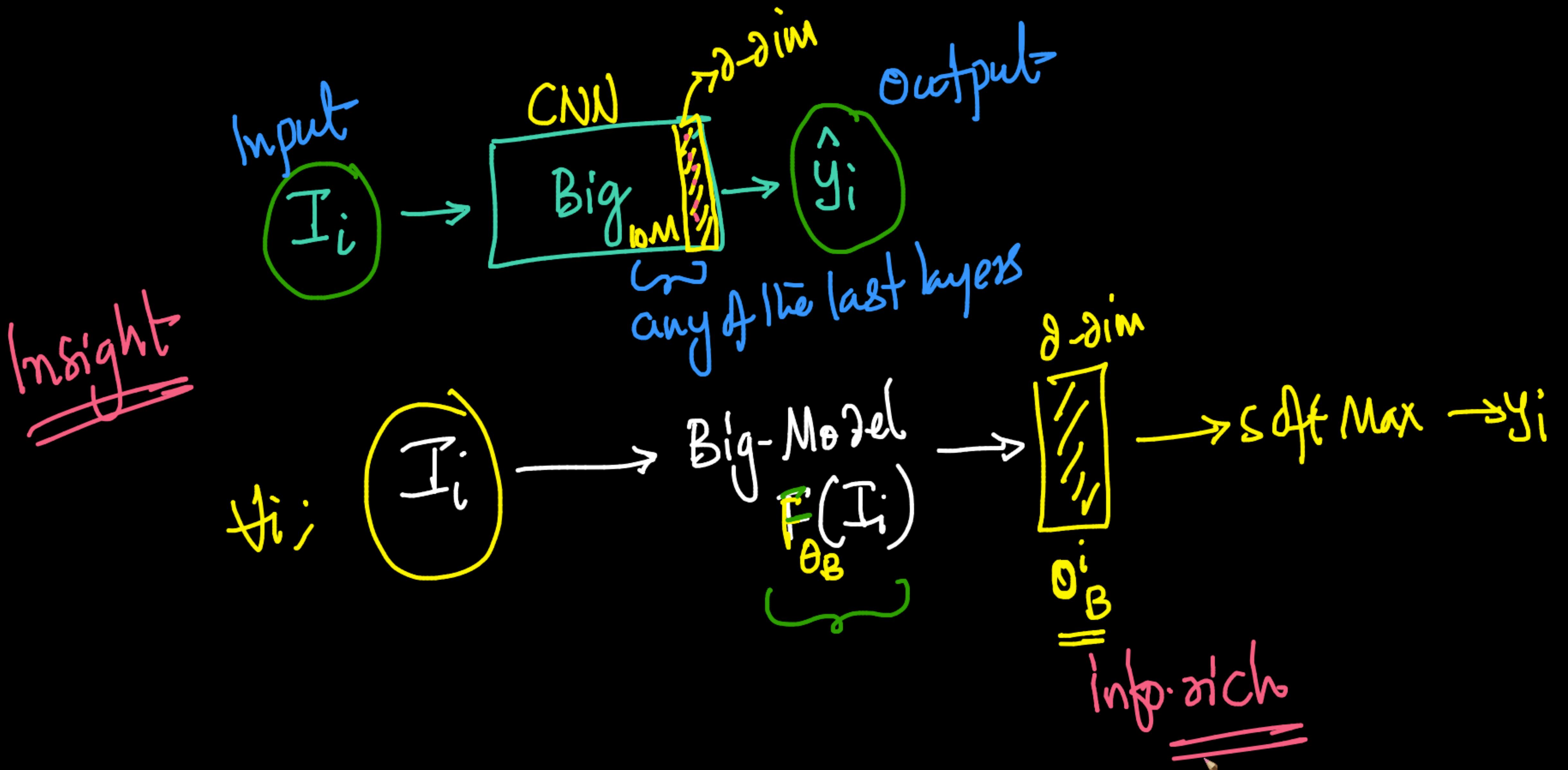
The original maximum-margin hyperplane algorithm proposed by Vapnik in 1963 constructed a [linear classifier](#). However, in 1992, Bernhard Boser, Isabelle Guyon and Vladimir Vapnik suggested a way to create nonlinear classifiers by applying the [kernel trick](#) (originally proposed by Aizerman et al.).

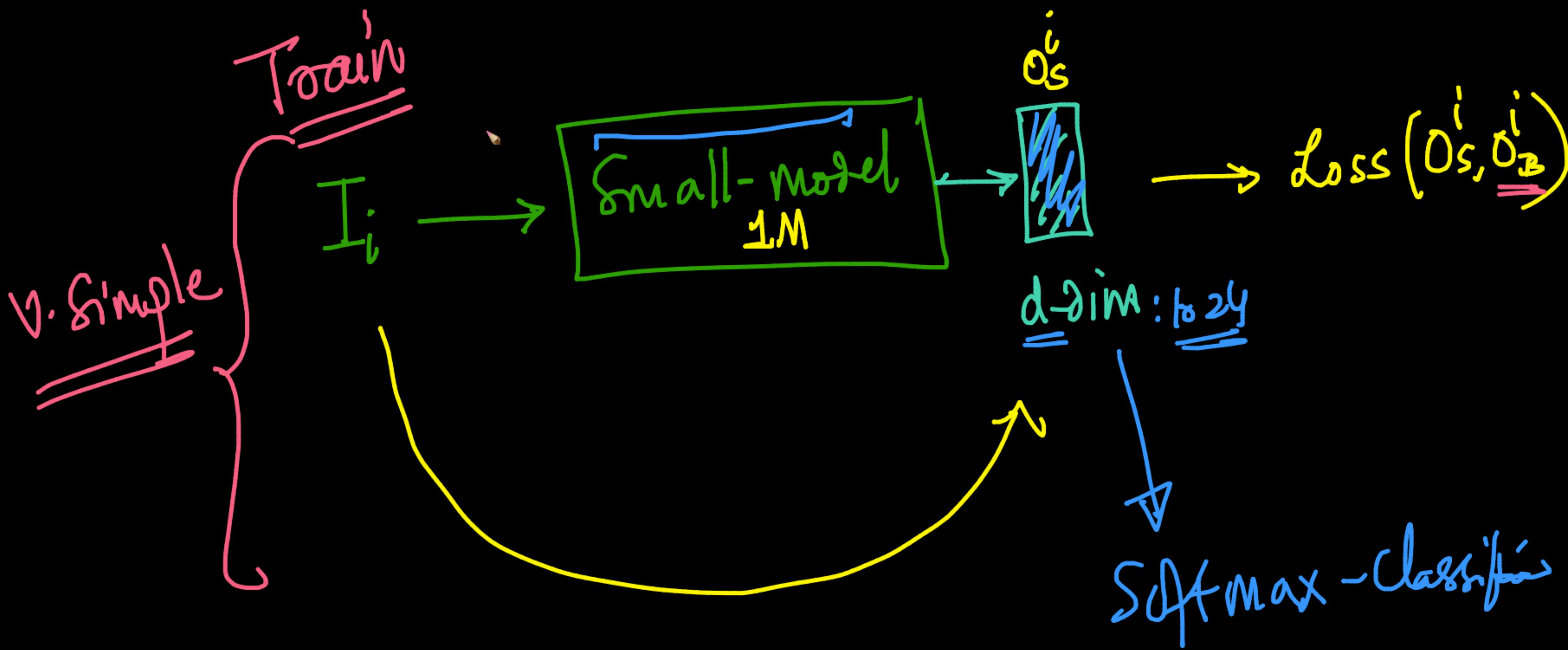














[Deep Learning: Convolutional Neural Nets. - Google Docs](#)[en.wikipedia.org/wiki/Support-vector\\_machine](#)[Support-vector machine - Wikipedia](#)

nearest to it. These  $\mathbf{x}_i$  are called *support vectors*.

## Soft-margin [edit]

To extend SVM to cases in which the data are not linearly separable, the *hinge loss* function is helpful

$$\max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)).$$

Note that  $y_i$  is the  $i$ -th target (i.e., in this case, 1 or  $-1$ ), and  $\mathbf{w}^T \mathbf{x}_i - b$  is the  $i$ -th output.

This function is zero if the constraint in (1) is satisfied, in other words, if  $\mathbf{x}_i$  lies on the correct side of the margin. For data on the wrong side of the margin, the function's value is proportional to the distance from the margin.

The goal of the optimization then is to minimize

$$\lambda \|\mathbf{w}\|^2 + \left[ \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)) \right],$$

where the parameter  $\lambda > 0$  determines the trade-off between increasing the margin size and ensuring that the  $\mathbf{x}_i$  lie on the correct side of the margin. Thus, for sufficiently small values of  $\lambda$ , it will behave similar to the hard-margin SVM, if the input data are linearly classifiable, but will still learn if a classification rule is viable or not. (This parameter  $\lambda$  is often also called  $C$ , e.g. in *LIBSVM*, but usually refers to the inverse of  $\lambda$ .)

## Nonlinear Kernels [edit]

The original maximum-margin hyperplane algorithm proposed by Vapnik in 1963 constructed a [linear classifier](#). However, in 1992, [Bernhard Boser](#), [Isabelle Guyon](#) and [Vladimir Vapnik](#) suggested a way to create nonlinear classifiers by applying the [kernel trick](#) (originally proposed by Aizerman et al. [formally](#))

