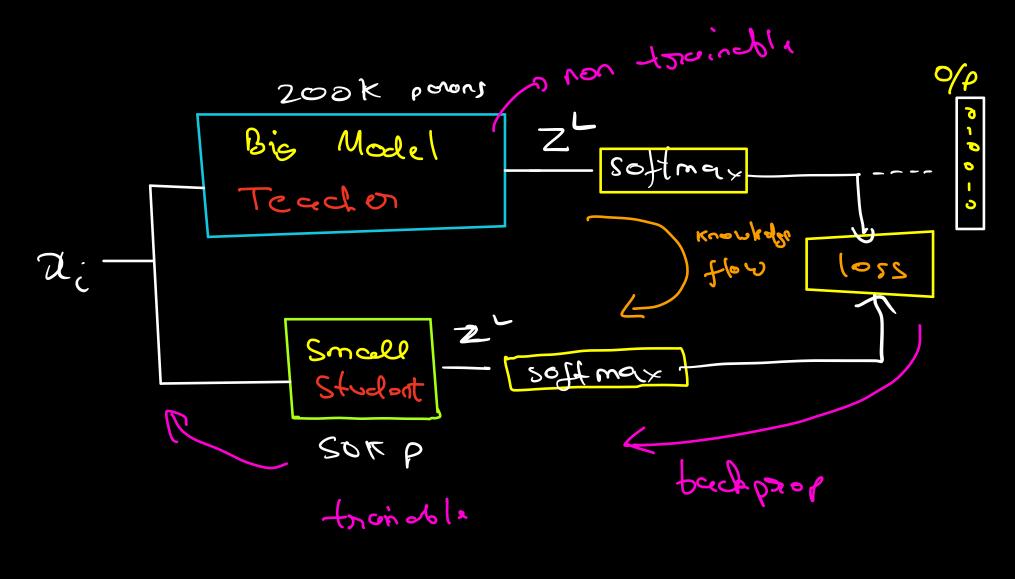
## Knowledge Distilation (2015) -> papor

Issues with DL models

-) Most nodern models are very large -> There is high lateray / prediction time

> (of af monning them in high

Distillation réfers to the toronsfer of information from a big model to a small one.



Knowledge can blow from larger model to smarller model

Which loss? -> Loss foncé con be chosen deponding on the application. However this setup is for classification -> After softmax, we get probability O: What is a way to compare prob. clist? 3 KL divorgence loss Z Pi los (Pi)

## Temporadure Controlled Softman

fon distillation, in above setup un hour "softnex". One issue with it i.  $\frac{2^{1}}{2^{1}} = \frac{20}{100} = \frac{e^{20}}{e^{20} + e^{40} + e^{60}} = \frac{4 \times 10^{18}}{2^{1}}$   $\frac{2^{1}}{2^{1}} = \frac{20}{100} = \frac{4 \times 10^{18}}{100}$   $\frac{2^{1}}{2^{1}} = \frac{e^{20}}{100} = \frac{4 \times 10^{18}}{100}$ Es! = 0,999

Squashing effet is to. (carge! 20 40 60 2 = ~2×~7~~1.5×~~ 4×10 2×10 0.979 So after softmax, most of the time, 0/p looks like this -, 0.00--- 0.00--- 0.979... 0.000 

So softmax orases a lot of lecerning by squashing! Temperadure-> S:2 E ezil+ (T) is just a factor. (sceling) 8g: T=10  $\frac{e^{20/10}}{2}$ ,  $\frac{e^{50/10}}{2}$ ,  $\frac{e^{50/10}$ 

## Proctical Setup

