# Low Rank Adaptation (LoRA)
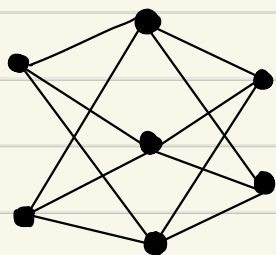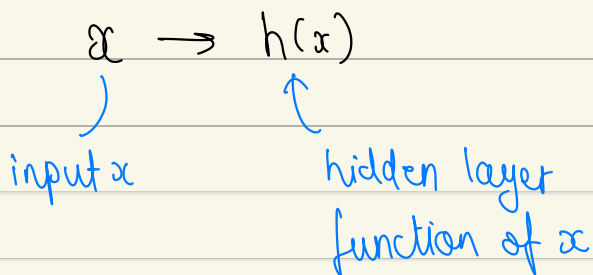
LoRA fine tunes a model by adding new trainable parameters



consider a neural network

$x \rightarrow h(x)$

input $x$

hidden layer
function of $x$

$h(x) = W_0 x$

$W_0$ = weight matrix
$x$ = input vector

$W_0 \in R^{d \times k}$
$x \in R^{k \times 1}$
$h(x) \in R^{d \times 1}$

→ so what happens if we retrain all parameters?

Suppose our weight matrix has shape 1000×1000
so this would mean changing <u>1000000</u> parameters!

so how does LoRA help us?

$$h(x) = W_0 x + \Delta W x$$

old term                    new term!

where $\Delta W x = B A$          where B & A are matrices

$$\therefore h(x) = W_0 x + B A x$$
$$h(x) = (W_0 + B A) x$$

r is called the intrinsic rank
of the model

$W_0, \Delta W \in R^{d \times k}$
$\begin{cases} B \in R^{d \times r} \\ A \in R^{r \times k} \end{cases}$
$h(x) \in R^{d \times l}$

the reason this works & we get efficiency gains because this
r is a lot smaller than d and k

$$\left( \overset{\text{frozen}}{\widetilde{W_0}} + \overset{\text{trainable}}{\widetilde{BA}} \right) x = h(x)$$

B and A contain
far fewer terms
than $W_0$

so eg if $d = 1000$
$\quad\quad\quad k = 1000 \quad\quad \rightarrow (d \times r) + (k \times r)$
$\quad\quad\quad r = 2 \quad\quad\quad\quad = (1000 \times 2) + (1000 \times 2)$
$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad = \underline{4000 \text{ trainable parameters}}$

significantly reduces the
parameters we have to train