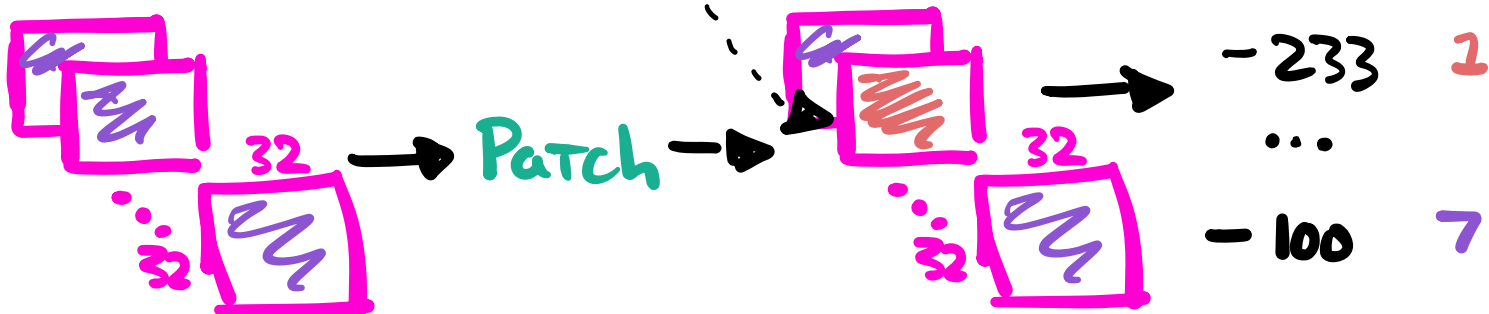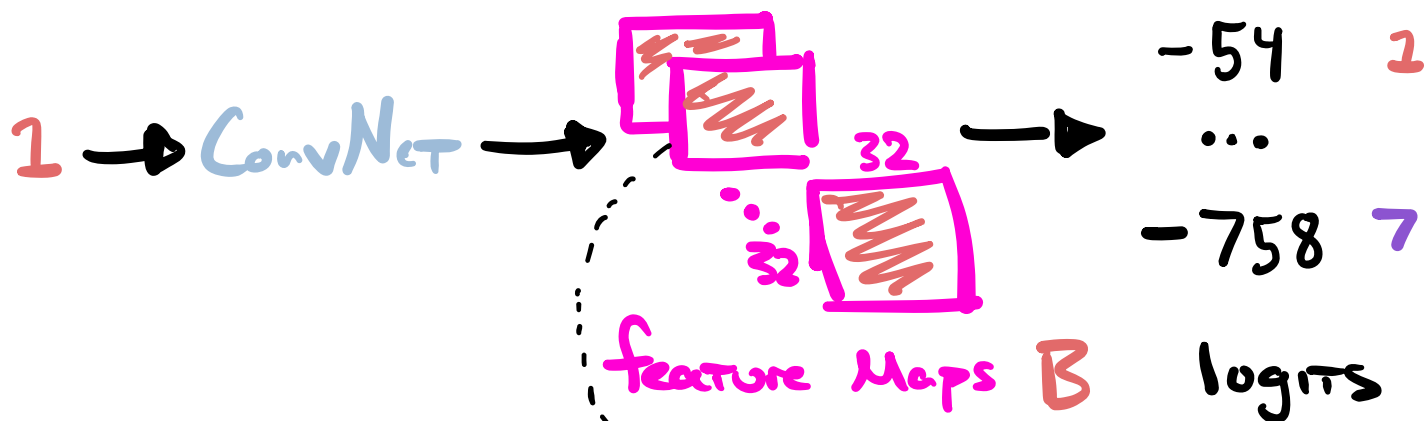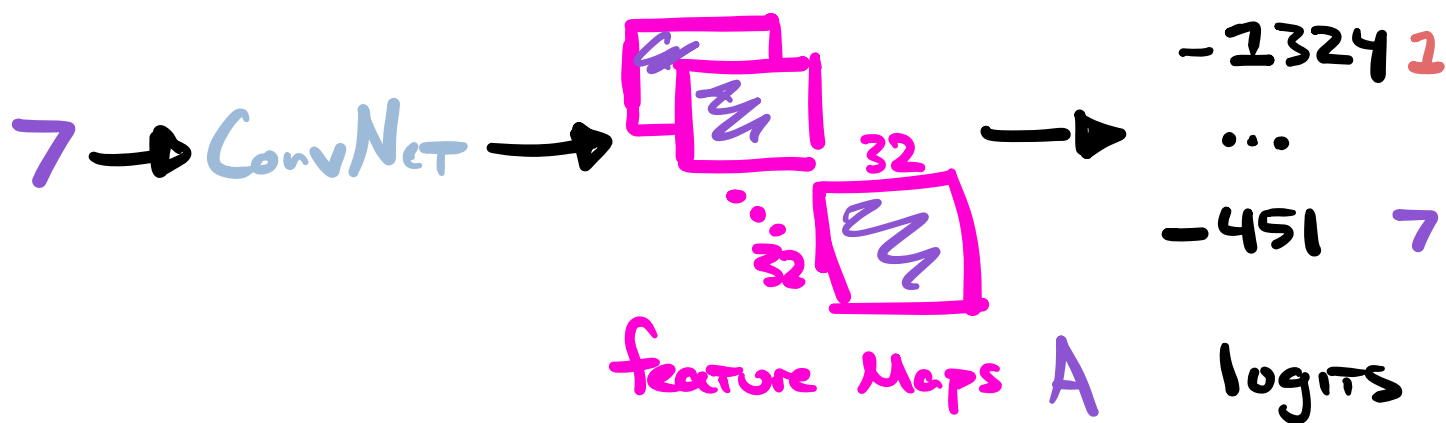# Activation Patching

> Which (learned) features led to prediction **A** rather than **B**?

## Approach (general)

(i) Run **A** and **B** fwd through Network

(ii) Collect **activations** of interest from **B** pass

(iii) <u>Patch</u> these in to same location in **A** pass

(iv) Observe change ($\Delta$) in output

For the homework, we consider ConvNet and specifically the last layer feature maps (32 × 32); We have **64** of them. Which is most important for telling 7s from 1s?

7 → ConvNet → feature Maps A → logits

-1324  1
...
-451  7

1 → ConvNet → feature Maps B → logits

-54  1
...
-758  7

→ Patch → → logits

-233  1
...
-100  7

$$\blacktriangle_{7,1} \overset{..}{=} \tilde{y}_7 - \tilde{y}_1 \; ; \; \text{measure for } \textcolor{green}{\text{patched}}$$

and $\textcolor{red}{\text{1s}}$

$$\textcolor{green}{\blacktriangle}_{7,1} - \textcolor{red}{\blacktriangle}_{7,1} = \underbrace{(\overset{\tilde{y}_7}{-100} - (-233))}_{} - \underbrace{(\overset{\tilde{y}_1}{-758} - (-54))}_{}$$

$$= 133 - 704 = 837$$

(if $\textcolor{green}{\tilde{y}_1}$ smaller say $-1000 \rightarrow 1604$ ↑

bigger say $+1000 \rightarrow -396$ ↓

$\tilde{y}_7$ smaller say $-1000 \rightarrow -63$ ↓

bigger say $+1000 \rightarrow 1937$ ↑ )

The range is funny so we

normalize (ish) by

$$\underbrace{\textcolor{blue}{\blacktriangle}_{7,1} - \textcolor{red}{\blacktriangle}_{7,1}}_{}$$

for run with $\textcolor{purple}{7s}$ ; sort of an

upper bound