So if we do Batch norm then pass to a sigmoid we could constraine the none linearity to it's Linear region. $\underline{\quad\int}$ this region

So then we add $(\gamma,\beta)$ $\gamma\hat{x} + \beta$

these are Learned

↑ normalize

& can be addjusted, so we could recover the original activations pre normalization. [thats all from paper] so now why even bother doing this work Just to have the possibility to undo it? Well the original outputs are Stongly dependent on the initutization of the weights In the layer Which means that the activations could Be large doto the fact that the weights Were randomly initalized some way.

Now we are putting a prior

that says at the start the activations will have the same shape But the mean and sd will be 0 & 1 respectivly + we can change that even going Back to the original activations however now it must be Justified, Doing so (changing $\beta$ & $\gamma$) must be Justified by addiqetly minimizing the loss. meaning the Magnitude of the activation is no longer Largly determined by some arbetrary initial weights but rather must be Justified by a minimization of loss

So Why exploding & Vanishing
gradient in Basic Recurrent Neural
nets?

It's pretty damn simple.
Think of it ~~like~~ Like this in
a network with no rolutent connections
~~our~~ our gradient is calculated
based only on the values of
the current training example
But in a RNN

you get some thing like

$$\frac{dL}{dW_i} = \ldots\ldots\ldots \cdot O_{t-1}$$

times ↙    ← activation
output at
some previous
time period

So if that was a huge output
or a very small one it will
still effect us.

but not only there

You will also get terms like

$$\frac{dL}{dw_i} = \ldots \ldots \cdot O_{O_t}$$

some activation output at the curr time ←

↓ tiny

?

this looks like normal but $O_{O_t}$ can be calc like this if it's the output of a recurrent connection.

$$RELU\left(X w_0 + b + w_i O_{O_{t-1}}\right) = O_{O_t}$$

So it too is effected by pass output outputs. all this means if we had an output that was HUGE or tiny we continue to have it effet out weight & bias updates even after the particular example has gone.

This is why RNN are more prone to these issue (exploding/vanishing grad)