

groupnorm chronicles

latentCall145

April 2024

1 Introduction

Normalization is important because it smoothens out loss landscapes which makes training easier. Image models often use batch normalization but some other models like VQGAN (and VQGAN-based models like Stable Diffusion 1, 2, and SDXL) use group normalization (GN) because GN normalizes each element within the batch separately, resulting in stable training even with small batch sizes.

This writeup describes the GN forward + backward pass with a focus on efficient GPU implementation on the backward pass. I got interested in this project in the first place because mixed-precision convolutions work faster in NHWC format (as opposed to NCHW which is PyTorch's default). The problem is that PyTorch's GN GPU implementation doesn't work for NHWC tensors, so I added NHWC support to GN myself to get the NHWC convolution speedup¹.

2 Forward Pass

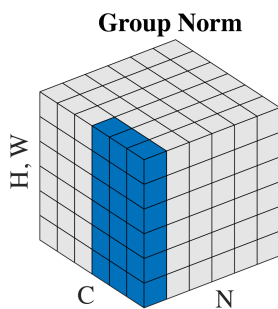


Figure 1: Diagram showing values in image (shaded in blue) that are reduced to one mean/variance during normalization. i drew this cube so many times while coding

¹Code here

Let the GN forward activation go as the following: $X \in \mathbb{R}^{N \times H \times W \times C}$ is the input. In GN, X is reshaped to $x \in \mathbb{R}^{N \times R \times G \times D}$ (where $R = HW$ and represents the resolution dimension) since that better represents the dimension which is being reduced in the normalization process (along the D and R dimension for x), and thus x will be mentioned instead of X for the remainder of this description. y is the output and is of same shape as x , $\gamma \in \mathbb{R}^C$ is the weights, and $\beta \in \mathbb{R}^C$ is the biases. In a similar way to x , γ, β will be reshaped to $\mathbb{R}^{G \times D}$.

$$\mu_{ng} = \frac{1}{RD} \sum_{r,d} x_{nrzd} \quad (1)$$

$$\sigma_{ng} = \sqrt{\frac{1}{RD} \sum_{r,d} (x_{nrzd} - \mu_{ng})^2} \quad (2)$$

$$\hat{x}_{nrzd} = \frac{x_{nrzd} - \mu_{ng}}{\sigma_{ng}} \quad (3)$$

$$y_{nrzd} = \gamma_{gd} \hat{x}_{nrzd} + \beta_{gd} \quad (4)$$

3 Backward Pass

3.1 Chain Rule Refresher

To preface the backwards derivation, here's a quick reminder of the chain rule: For function $f(x_1(t), x_2(t), \dots, x_n(t))$ (note how t is involved in the computation of multiple intermediate functions x_i which then affects f),

$$\frac{\partial f}{\partial t} = \sum_i \frac{\partial f}{\partial x_i} \frac{\partial x_i}{\partial t}$$

3.2 Weight Partial

The backwards derivation for the weight and bias is pretty straightforward. Note that in the below equations, we are looping over n, r (instead of every element of y) since a single element of γ or β affects the output of $n \times r$ elements of y :

$$\frac{\partial f}{\partial \gamma_{gd}} = \sum_{n,r} \frac{\partial f}{\partial y_{nrzd}} \frac{\partial y_{nrzd}}{\partial \gamma_{gd}} \quad (5)$$

$$= \sum_{n,r} \frac{\partial f}{\partial y_{nrzd}} \hat{x}_{nrzd} \quad (6)$$

$$\frac{\partial f}{\partial \beta_{gd}} = \sum_{n,r} \frac{\partial f}{\partial y_{nrzd}} \frac{\partial y_{nrzd}}{\partial \beta_{gd}} \quad (7)$$

$$= \sum_{n,r} \frac{\partial f}{\partial y_{nrzd}} \quad (8)$$

3.3 Activation Partialals

The backwards derivation for the input is not as straightforward since nudging any element $x_{nr'gd}$ affects every element within its group (because the nudging affects the mean/variance of the group), which in turn affects the output:

$$\frac{\partial f}{\partial x_{nr'gd}} = \sum_{r',d'} \frac{\partial f}{\partial y_{nr'gd'}} \frac{\partial y_{nr'gd'}}{\partial \hat{x}_{nr'gd'}} \frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr'gd}} \quad (9)$$

$$= \sum_{r',d'} \frac{\partial f}{\partial y_{nr'gd'}} \gamma_{gd'} \frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr'gd}} \quad (10)$$

Note that the only reason I'm using r', d' for the sums instead of r, d is because r, d specifies the activation input whose partial I want to calculate. We're still summing through the R and D dimension, just under different labels. Focusing on the partials for the normalized variable \hat{x} , let's first apply the quotient rule:

$$\frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr'gd}} = \frac{1}{\sigma_{ng}^2} \left(\sigma_{ng} \frac{\partial (x_{nr'gd'} - \mu_{ng})}{\partial x_{nr'gd}} - (x_{nr'gd'} - \mu_{ng}) \frac{\partial \sigma_{ng}}{\partial x_{nr'gd}} \right) \quad (11)$$

$$= \frac{1}{\sigma_{ng}} \left(\frac{\partial x_{nr'gd'}}{\partial x_{nr'gd}} - \frac{\partial \mu_{ng}}{\partial x_{nr'gd}} \right) - \frac{1}{\sigma_{ng}^2} \left((x_{nr'gd'} - \mu_{ng}) \frac{\partial \sigma_{ng}}{\partial x_{nr'gd}} \right) \quad (12)$$

$$= \frac{1}{\sigma_{ng}} (\delta_{r'r} \delta_{d'd} - \frac{\partial \mu_{ng}}{\partial x_{nr'gd}}) - \frac{1}{\sigma_{ng}^2} \left((x_{nr'gd'} - \mu_{ng}) \frac{\partial \sigma_{ng}}{\partial x_{nr'gd}} \right) \quad (13)$$

If you haven't seen the δ_{ij} symbol yet, it's called the Kronecker delta and equals 1 if $i = j$, and 0 otherwise. Thus, $\delta_{r'r} \delta_{d'd}$ equals 1 if $(r', d') = (r, d)$ and 0

otherwise. Now let's unpack the partials for the mean and standard deviation:

$$\frac{\partial \mu_{ng}}{\partial x_{nr gd}} = \frac{1}{RD} \quad (14)$$

$$\frac{\partial \sigma_{ng}}{\partial x_{nr gd}} = \frac{1}{2RD\sigma_{ng}} \sum_{r',d'} 2(x_{nr'gd'} - \mu_{ng}) \frac{\partial (x_{nr'gd'} - \mu_{ng})}{\partial x_{nr gd}} \quad (15)$$

$$= \frac{1}{RD\sigma_{ng}} \sum_{r',d'} (x_{nr'gd'} - \mu_{ng}) \left(\frac{\partial x_{nr'gd'}}{\partial x_{nr gd}} - \frac{\partial \mu_{ng}}{\partial x_{nr gd}} \right) \quad (16)$$

$$= \frac{1}{RD\sigma_{ng}} \sum_{r',d'} (x_{nr'gd'} - \mu_{ng}) (\delta_{r'r} \delta_{d'd} - \frac{1}{RD}) \quad (17)$$

$$= \frac{1}{RD\sigma_{ng}} \left((x_{nr gd} - \mu_{ng}) - \frac{1}{RD} \sum_{r',d'} (x_{nr'gd'} - \mu_{ng}) \right) \quad (18)$$

$$= \frac{1}{RD\sigma_{ng}} \left((x_{nr gd} - \mu_{ng}) - \frac{1}{RD} \sum_{r',d'} x_{nr'gd'} + \frac{1}{RD} \sum_{r',d'} \mu_{ng} \right) \quad (19)$$

$$= \frac{1}{RD\sigma_{ng}} \left((x_{nr gd} - \mu_{ng}) - \mu_{ng} + \mu_{ng} \right) \quad (20)$$

$$= \frac{x_{nr gd} - \mu_{ng}}{RD\sigma_{ng}} \quad (21)$$

We're not going to write $\frac{\partial \sigma_{ng}}{\partial x_{nr gd}} = \frac{\hat{x}_{nr gd}}{RD}$ even though it's mathematically equivalent because we don't want to store/load \hat{x} when performing the backward pass on the GPU. Doing so wastes memory and would actually worsen performance as loading/storing values in GPUs (specifically to global memory, which is the GPU equivalent to RAM) is much slower than math operations. Plugging the above partials to $\frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr gd}}$:

$$\frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr gd}} = \frac{1}{\sigma_{ng}} \left(\delta_{r'r} \delta_{d'd} - \frac{\partial \mu_{ng}}{\partial x_{nr gd}} \right) - \frac{1}{\sigma_{ng}^2} \left((x_{nr'gd'} - \mu_{ng}) \frac{\partial \sigma_{ng}}{\partial x_{nr gd}} \right) \quad (22)$$

$$= \frac{\delta_{r'r} \delta_{d'd}}{\sigma_{ng}} - \frac{1}{RD\sigma_{ng}} - \frac{1}{RD\sigma_{ng}^3} (x_{nr'gd'} - \mu_{ng})(x_{nr gd} - \mu_{ng}) \quad (23)$$

And now plugging into $\frac{\partial f}{\partial x_{nr gd}}$:

$$\frac{\partial f}{\partial x_{nr gd}} = \sum_{r',d'} \frac{\partial f}{\partial y_{nr'gd'}} \gamma_{gd'} \frac{\partial \hat{x}_{nr'gd'}}{\partial x_{nr gd}} \quad (24)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nr gd}} &= \sum_{r',d'} \frac{\partial f}{\partial y_{nr'gd'}} \gamma_{gd'} \left(\frac{\delta_{r'r} \delta_{d'd}}{\sigma_{ng}} \right. \\ &\quad \left. - \frac{1}{RD\sigma_{ng}} - \frac{1}{RD\sigma_{ng}^3} (x_{nr'gd'} - \mu_{ng})(x_{nr gd} - \mu_{ng}) \right) \quad (25) \end{aligned}$$

$$\frac{\partial f}{\partial x_{nr gd}} = \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \left(\frac{\gamma_{gd'} \delta_{r'r} \delta_{d'd}}{\sigma_{ng}} - \frac{\gamma_{gd'}}{RD \sigma_{ng}} - \frac{\gamma_{gd'}}{RD \sigma_{ng}^3} (x_{nr' gd'} - \mu_{ng})(x_{nr gd} - \mu_{ng}) \right) \quad (26)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nr gd}} &= \frac{1}{\sigma_{ng}} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} \delta_{r'r} \delta_{d'd} \\ &+ \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \left(-\frac{\gamma_{gd'}}{RD \sigma_{ng}} - \frac{\gamma_{gd'}}{RD \sigma_{ng}^3} (x_{nr' gd'} - \mu_{ng})(x_{nr gd} - \mu_{ng}) \right) \end{aligned} \quad (27)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nr gd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} \\ &+ \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \left(-\frac{\gamma_{gd'}}{RD \sigma_{ng}} - \frac{\gamma_{gd'}}{RD \sigma_{ng}^3} (x_{nr' gd'} - \mu_{ng})(x_{nr gd} - \mu_{ng}) \right) \end{aligned} \quad (28)$$

3.4 Optimizations

At this point, you can start implementing this on a GPU, but it's going to be relatively slow because each element in Equation 28's sum takes an unnecessary number of operations (RD elements * (3 subtractions and 5 products per element) + $(RD - 1)$ additions for the actual sum). To speed this sum up, let's factor out the terms in the loops as much as possible:

$$\begin{aligned} \frac{\partial f}{\partial x_{nr gd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} - \frac{1}{RD \sigma_{ng}} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} \\ &+ \frac{\mu_{ng} - x_{nr gd}}{RD \sigma_{ng}^3} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} (x_{nr' gd'} - \mu_{ng}) \end{aligned} \quad (29)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nr gd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} - \frac{1}{RD \sigma_{ng}} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} \\ &+ \frac{\mu_{ng} - x_{nr gd}}{RD \sigma_{ng}^3} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} x_{nr' gd'} \\ &+ \frac{\mu_{ng}(x_{nr gd} - \mu_{ng})}{RD \sigma_{ng}^3} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} \end{aligned} \quad (30)$$

$$\begin{aligned}
\frac{\partial f}{\partial x_{nr gd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} + \left(\frac{\mu_{ng}(x_{nr gd} - \mu_{ng})}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} \\
&\quad + \frac{\mu_{ng} - x_{nr gd}}{RD\sigma_{ng}^3} \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \gamma_{gd'} x_{nr' gd'}
\end{aligned} \tag{31}$$

$$\begin{aligned}
\frac{\partial f}{\partial x_{nr gd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} + \left(\frac{\mu_{ng}(x_{nr gd} - \mu_{ng})}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) \sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr' gd'}} \\
&\quad + \frac{\mu_{ng} - x_{nr gd}}{RD\sigma_{ng}^3} \sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr' gd'}} x_{nr' gd'}
\end{aligned} \tag{32}$$

Now the loops are really simple. The $\sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr' gd'}}$ loops have $D(R-1)$ adds for the inner sum, then D mults, then $(D-1)$ adds for the outer sum (so $(RD-1)$ adds and D mults total). The $\sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr' gd'}} x_{nr' gd'}$ loops have $RD-1$ adds and $RD+D$ mults total.

Even better, the inner loops can be reused for the weight partial calculation. Let $S_{ngd}^{(y)} = \sum_r \frac{\partial f}{\partial y_{nr gd}}$ and $S_{ngd}^{(xy)} = \sum_r \frac{\partial f}{\partial y_{nr gd}} x_{nr gd}$. We can sub these sums in Equation 6:

$$\frac{\partial f}{\partial \gamma_{gd}} = \sum_{n,r} \frac{\partial f}{\partial y_{nr gd}} \hat{x}_{nr gd} \tag{33}$$

$$= \sum_n \sum_r \frac{\partial f}{\partial y_{nr gd}} \frac{(x_{nr gd} - \mu_{ng})}{\sigma_{ng}} \tag{34}$$

$$= \sum_n \frac{1}{\sigma_{ng}} \sum_r \frac{\partial f}{\partial y_{nr gd}} (x_{nr gd} - \mu_{ng}) \tag{35}$$

$$= \sum_n \frac{1}{\sigma_{ng}} \left(\sum_r \frac{\partial f}{\partial y_{nr gd}} x_{nr gd} - \mu_{ng} \sum_r \frac{\partial f}{\partial y_{nr gd}} \right) \tag{36}$$

$$\frac{\partial f}{\partial \gamma_{gd}} = \sum_n \frac{S_{ngd}^{(xy)} - \mu_{ng} S_{ngd}^{(y)}}{\sigma_{ng}} \tag{37}$$

Doing the same for Equation 8:

$$\frac{\partial f}{\partial \beta_{gd}} = \sum_{n,r} \frac{\partial f}{\partial y_{nr gd}} \tag{38}$$

$$= \sum_n \sum_r \frac{\partial f}{\partial y_{nr gd}} \tag{39}$$

$$\frac{\partial f}{\partial \beta_{gd}} = \sum_n S_{ngd}^{(y)} \tag{40}$$

And because we haven't done so yet, let's also do Equation 32:

$$\begin{aligned} \frac{\partial f}{\partial x_{nrzd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + \left(\frac{\mu_{ng}(x_{nrzd} - \mu_{ng})}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) \sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr'gd'}} \\ &\quad + \frac{\mu_{ng} - x_{nrzd}}{RD\sigma_{ng}^3} \sum_{d'} \gamma_{gd'} \sum_{r'} \frac{\partial f}{\partial y_{nr'gd'}} x_{nr'gd'} \end{aligned} \quad (41)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nrzd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + \left(\frac{\mu_{ng}(x_{nrzd} - \mu_{ng})}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) \sum_{d'} \gamma_{gd'} S_{ngd'}^{(y)} \\ &\quad + \frac{\mu_{ng} - x_{nrzd}}{RD\sigma_{ng}^3} \sum_{d'} \gamma_{gd'} S_{ngd'}^{(xy)} \end{aligned} \quad (42)$$

There are a lot of terms in Equation 42 meaning lots of operations per partial, which is bad because we have to run Equation 42 many times ($NRGD$ times, once for each element of x). By rearranging some terms, we can rewrite Equation 42 in the form:

$$\frac{\partial f}{\partial x_{nrzd}} = \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + c^{(1)} x_{nrzd} + c^{(2)}$$

This'll greatly reduce the number of operations per element. Here it is:

$$S_{ng}^{(y\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(y)} \quad (43)$$

$$S_{ng}^{(xy\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(xy)} \quad (44)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nrzd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + \left(\frac{\mu_{ng}(x_{nrzd} - \mu_{ng})}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) \sum_{d'} \gamma_{gd'} S_{ng}^{(y\gamma)} \\ &\quad + \frac{\mu_{ng} - x_{nrzd}}{RD\sigma_{ng}^3} \sum_{d'} \gamma_{gd'} S_{ng}^{(xy\gamma)} \end{aligned} \quad (45)$$

$$\begin{aligned} \frac{\partial f}{\partial x_{nrzd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + \left(\frac{\mu_{ng} x_{nrzd}}{RD\sigma_{ng}^3} - \frac{\mu_{ng}^2}{RD\sigma_{ng}^3} - \frac{1}{RD\sigma_{ng}} \right) S_{ng}^{(y\gamma)} \\ &\quad + \left(\frac{\mu_{ng}}{RD\sigma_{ng}^3} - \frac{x_{nrzd}}{RD\sigma_{ng}^3} \right) S_{ng}^{(xy\gamma)} \end{aligned} \quad (46)$$

$$\begin{aligned}
\frac{\partial f}{\partial x_{nrgd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrgd}} + \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} x_{nrgd} \\
&\quad - \left(\frac{\mu_{ng}^2}{RD\sigma_{ng}^3} + \frac{1}{RD\sigma_{ng}} \right) S_{ng}^{(y\gamma)} + \frac{\mu_{ng}}{RD\sigma_{ng}^3} S_{ng}^{(xy\gamma)}
\end{aligned} \tag{47}$$

$$\begin{aligned}
\frac{\partial f}{\partial x_{nrgd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrgd}} + \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} x_{nrgd} \\
&\quad + \frac{-\mu_{ng}^2 S_{ng}^{(y\gamma)} + \mu_{ng} S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} - \frac{S_{ng}^{(y\gamma)}}{RD\sigma_{ng}}
\end{aligned} \tag{48}$$

$$\begin{aligned}
\frac{\partial f}{\partial x_{nrgd}} &= \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrgd}} + \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} x_{nrgd} \\
&\quad - \mu_{ng} \frac{(\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)})}{RD\sigma_{ng}^3} - \frac{S_{ng}^{(y\gamma)}}{RD\sigma_{ng}}
\end{aligned} \tag{49}$$

From here, you can see that $c^{(1)}, c^{(2)}$ have shape $N \times G$ where:

$$c_{ng}^{(1)} = \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} \tag{50}$$

$$c_{ng}^{(2)} = -\mu_{ng} c_{ng}^{(1)} - \frac{S_{ng}^{(y\gamma)}}{RD\sigma_{ng}} \tag{51}$$

Putting it all together, here are the partials for the backward pass:

$$S_{ngd}^{(y)} = \sum_r \frac{\partial f}{\partial y_{nrzd}} \quad (52)$$

$$S_{ngd}^{(xy)} = \sum_r \frac{\partial f}{\partial y_{nrzd}} x_{nrzd} \quad (53)$$

$$S_{ng}^{(y\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(y)} \quad (54)$$

$$S_{ng}^{(xy\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(xy)} \quad (55)$$

$$\frac{\partial f}{\partial \gamma_{gd}} = \sum_n \frac{S_{ngd}^{(xy)} - \mu_{ng} S_{ngd}^{(y)}}{\sigma_{ng}} \quad (56)$$

$$\frac{\partial f}{\partial \beta_{gd}} = \sum_n S_{ngd}^{(y)} \quad (57)$$

$$c_{ng}^{(1)} = \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD\sigma_{ng}^3} \quad (58)$$

$$c_{ng}^{(2)} = -\mu_{ng} c_{ng}^{(1)} - \frac{S_{ng}^{(y\gamma)}}{RD\sigma_{ng}} \quad (59)$$

$$\frac{\partial f}{\partial x_{nrzd}} = \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nrzd}} + c_{ng}^{(1)} x_{nrzd} + c_{ng}^{(2)} \quad (60)$$

4 Fusing Normalization with Activation

One more thing. A lot of models that GN (such as VQGAN-based models) use an activation directly after the normalization. If you implement this with a GN layer followed by an activation layer, values are stored in GPU global memory after normalization, which are then re-loaded to run the activation layer.

As I stated earlier, loading/storing from GPU global memory is really slow, so we can speed up normalization and activation by running both operations in one layer; this technique is called operator fusion². This removes the slow GPU global memory transfer between layers. Operator fusion also saves memory as ML libraries need to store intermediate activations for each layer to calculate gradients, so fusing layers reduces the number of intermediate activations to store.

For the forward pass, we add an activation function $\phi(x)$ to $y_{nrzd} = \gamma_{gd} \hat{x}_{nrzd} + \beta_{gd}$, so y_{nrzd} now becomes:

$$y_{nrzd} = \phi(\gamma_{gd} \hat{x}_{nrzd} + \beta_{gd}) \quad (61)$$

This is trivial to implement on a GPU. For the backward pass:

²Operator fusion is one of the most important features (if not THE most important feature) of ML compilers like Triton or XLA.

$$\frac{\partial f}{\partial x_{nr gd}} = \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \frac{\partial y_{nr' gd'}}{\partial \hat{x}_{nr' gd'}} \frac{\partial \hat{x}_{nr' gd'}}{\partial x_{nr gd}} \quad (62)$$

$$= \sum_{r', d'} \frac{\partial f}{\partial y_{nr' gd'}} \phi'(\gamma_{gd'} \hat{x}_{nr' gd'} + \beta_{gd'}) \gamma_{gd'} \frac{\partial \hat{x}_{nr' gd'}}{\partial x_{nr gd}} \quad (63)$$

Essentially, for every $\frac{\partial f}{\partial y_{nr gd}}$ term you see, you also multiply it with $\phi'(\gamma_{gd} \hat{x}_{nr gd} + \beta_{gd})$, so Equations 52 - 60 become:

$$S_{ngd}^{(y)} = \sum_r \frac{\partial f}{\partial y_{nr gd}} \phi'(\gamma_{gd} \hat{x}_{nr gd} + \beta_{gd}) \quad (64)$$

$$S_{ngd}^{(xy)} = \sum_r \frac{\partial f}{\partial y_{nr gd}} \phi'(\gamma_{gd} \hat{x}_{nr gd} + \beta_{gd}) x_{nr gd} \quad (65)$$

$$S_{ng}^{(y\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(y)} \quad (66)$$

$$S_{ng}^{(xy\gamma)} = \sum_{d'} \gamma_{gd'} S_{ngd'}^{(xy)} \quad (67)$$

$$\frac{\partial f}{\partial \gamma_{gd}} = \sum_n \frac{S_{ngd}^{(xy)} - \mu_{ng} S_{ngd}^{(y)}}{\sigma_{ng}} \quad (68)$$

$$\frac{\partial f}{\partial \beta_{gd}} = \sum_n S_{ngd}^{(y)} \quad (69)$$

$$c_{ng}^{(1)} = \frac{\mu_{ng} S_{ng}^{(y\gamma)} - S_{ng}^{(xy\gamma)}}{RD \sigma_{ng}^3} \quad (70)$$

$$c_{ng}^{(2)} = -\mu_{ng} c_{ng}^{(1)} - \frac{S_{ng}^{(y\gamma)}}{RD \sigma_{ng}} \quad (71)$$

$$\frac{\partial f}{\partial x_{nr gd}} = \frac{\gamma_{gd}}{\sigma_{ng}} \frac{\partial f}{\partial y_{nr gd}} \phi'(\gamma_{gd} \hat{x}_{nr gd} + \beta_{gd}) + c_{ng}^{(1)} x_{nr gd} + c_{ng}^{(2)} \quad (72)$$