

Why does $\text{AllGather}_x(A[I_x]) \rightarrow A[I]$ in forward
 imply $\text{ReduceScatter}_x(A'[I]\{U_x\}) \rightarrow A'[I_x]$ in backward?

Imagine: $A[I_x]$ represents model parameters sharded
 along X axis
 and data $[I_x]$ is also sharded.

Forward pass:

$$\begin{aligned} & \text{data}[I_x] \cdot A[I_x] \\ & \quad \downarrow \text{AllGather} \\ \rightarrow & \text{data}[I_x] \cdot A[I] \\ \rightarrow & \text{loss}[I_x] = (\text{value} - \text{label}) \end{aligned}$$

$$\rightarrow \text{loss} = \sum_{I_x} \text{loss}[I_x]$$

$$\frac{\partial \text{loss}}{\partial A} = \sum_{I_x} \frac{\partial \text{loss}[I_x]}{\partial A}$$

↳ This means if we just compute $\frac{\partial \text{loss}(I_x)}{\partial A}$ on each device
 it will result in $\underbrace{A'[I]\{U_x\}}_{\substack{\text{unreduced} \\ \frac{\partial \text{loss}}{\partial A}}}$

To obtain the true gradient $\frac{\partial \text{loss}}{\partial A}$, we could

$$\text{AllReduce}(A'[I]\{U_x\}) \rightarrow A'[I]$$

But the goal is to keep $A[I_x]$ on each device.

We can just do:

$$\text{ReduceScatter}_x(A'[I]\{U_x\}) \rightarrow A'[I_x]$$

$A'[I_x]$ already gives us enough info to do

$$A[I_x] \leftarrow A[I_x] - \text{lr} \cdot A'[I_x]$$

Why does $\text{ReduceScatter}_x(A[I]\{U_x\}) \rightarrow A[I_v]$
in forward pass imply

$$\text{AllGather}_x(A'[I_x]) \rightarrow A'[I]$$

Scenario:

$$\text{Input}[B, D], W_1[D, F_x], W_2[F_x, G]$$

$$\text{Input}[B, D] @ W_1[D, F_x]$$

$$= \text{out}_1[B, F_x]$$

$$\Rightarrow \text{out}_1[B, F_x] @ W_2[F_x, G] \xrightarrow{\text{locally}} \text{out}_2[B, G]\{U_x\}$$

In theory, if we want to make $\text{out}_2[B, G]$
fully replicated, we need to do

$$\text{AllReduce} =^{\text{which}} \text{is} \text{ReduceScatter} \rightarrow \text{AllGather}$$

In practice, we don't have to do that, we can

$$\text{ReduceScatter}_x(\text{out}_2[B, G]\{U_x\}) = \text{out}_2[B, G_x]$$

And don't do AllGather because not necessary needed

In this case, when we receive $\text{out}_2'(B, G_x)$,
to compute $W_2'[F_x, G]$

By chain-rule:
$$\underset{\substack{\uparrow \\ \mathbb{R}^{1 \times G}}}{Z} = \underset{\substack{\uparrow \\ \mathbb{R}^{1 \times F}}}{X} \underset{\substack{\uparrow \\ \mathbb{R}^{F \times G}}}{W}$$

If we already have: $\frac{\partial L}{\partial Z} \in \mathbb{R}^{1 \times G}$

$$\frac{\partial L}{\partial W} = \underset{\substack{\uparrow \\ \mathbb{R}^{F \times 1}}}{X^T} \underset{\substack{\uparrow \\ \mathbb{R}^{1 \times G}}}{\frac{\partial L}{\partial Z}}$$

So to compute: $W_2'[F_x, G]$

we need $\text{out}_1[B, F]$ and $\text{out}_2'(B, G)$

More precisely: $\frac{1}{B} (\text{out}_1^T[F, B] \odot \text{out}_2'[B, G])$

In this case, $\text{out}_1[B, \underbrace{F_x}_{\text{sharded}}]$ already matches $W_2'[\underbrace{F_x}_{\text{sharded}}, \underbrace{G}_{\text{match}}]$

but we need G fully sharded.

Therefore, during the backward pass,

We need to $\text{AllGather}_x(\text{out}_2'[B, G_x])$

and then multiply this with $\text{out}_1[B, F_x]$