

Q3. Part 1

$$1) \text{ minibatches per epoch} = \frac{76800}{1536} = 50$$

$$\text{total mini batches} = 50 \times 10 = 500$$

$$\text{how many times is it distributed among 6 machines} = \frac{500}{6} = 83.3$$

Typical training round with an even load

$$= 1.5s \text{ (comms from master)} + 4s \text{ (computation)} + 1.5s \text{ (comms to master)} + 6 \times 0.015s = 7.09s$$

$$\begin{aligned} \text{Training with the last } \frac{1}{3} \text{ load} &= 1.5s + 2 \times 0.015s \\ &= 7.03s \end{aligned}$$

$$\text{Total time} = 7.09 \times 83 + 7.03 = 595.5s$$

$$2) \text{ Comm and agg time} = 1.5s + 1.5s + 0.015m$$

$$\text{Training time} = 4s$$

$$3 + 0.015m > 4$$

$$0.015m > 1$$

$$m > \frac{1}{0.015}$$

$$m > 66.6$$

$$m \geq 67$$



Q3. Part 2

$$\text{Data per device} = \frac{204800}{8} = 25600$$

$$\text{Processed data per device after 10 epochs} = 256000$$

$$\# \text{ batches processed} = \frac{256000}{32} = 800$$

Assuming the updates are sent back after all 10 epochs are processed.

$$\text{Total time} = 7s (\text{receiving the model}) + 4 \times 800 (\text{updates}) + 7s (\text{sending the model back}) = 32014s$$



## Q2. Part 1

Assuming "the same" means equal number of operations

### Asynchronous

Core 1:	starts at $t = 0.15s$	Update intervals = $0.15 - 0.35$
Core 2:	at $t = 0.20s$	= $0.2 - 0.4$
Core 3:	at $t = 0.25s$	= $0.25 - 0.45$
Core 4:	at $t = 0.50s$	= doesn't update

### Synchronous

Core 1:  $0.0 - 0.2s$   $0.2s - 0.4s$   $0.4s - 0.6s$  with no delay in between

$$x_{\text{async}} = x_0 - \eta \nabla f_1(x_0) - \eta \nabla f_2(x_0) - \eta \nabla f_3(x_0)$$

$$x_{\text{seq}} = x_0 - \eta \nabla f_1(x_0) - \eta \nabla f_2(x_1) - \eta \nabla f_3(x_2)$$

$$\text{where } x_1 = x_0 - \eta \nabla f_1(x_0) \text{ and } x_2 = x_1 - \eta \nabla f_2(x_1)$$

$$\|x_{\text{async}} - x_{\text{seq}}\| = \| -\eta \nabla f_2(x_0) - \eta \nabla f_3(x_0) - (-\eta \nabla f_2(x_1) - \eta \nabla f_3(x_2)) \|$$

$$= \| -\eta \nabla f_2(x_0) - \eta \nabla f_3(x_0) + \eta \nabla f_2(x_1) + \eta \nabla f_3(x_2) \|$$

$$= \| (\eta \nabla f_2(x_1) - \eta \nabla f_2(x_0)) + (\eta \nabla f_3(x_2) - \eta \nabla f_3(x_0)) \|$$

$$\text{triangle inequality} \leq \| \eta \nabla f_2(x_1) - \eta \nabla f_2(x_0) \| + \| \eta \nabla f_3(x_2) - \eta \nabla f_3(x_0) \|$$

$$\text{triangle inequality} \leq \| \eta \nabla f_2(x_1) \| + \| \eta \nabla f_2(x_0) \| + \| \eta \nabla f_3(x_2) \| + \| \eta \nabla f_3(x_0) \|$$

$$\begin{aligned} \text{(via } 0 < \|\nabla f_i\| \leq \Delta) &\leq \eta \Delta + \eta \Delta + \eta \Delta + \eta \Delta \\ &\leq 4\eta \Delta \end{aligned}$$



Assuming "the same" means executed for the same amount of time 0.6s

### Asynchronous

Execution times: Core 1: 0.15 - 0.35    0.35 - 0.55  
Core 2: 0.2 - 0.4    0.4 - 0.6  
Core 3: 0.25 - 0.45  
Core 4: Doesn't update

### Synchronous

Same as before.

~~Unsure of the notation but basically two more operations that don't cancel out with other operations so  $4\eta\Delta + 2\eta\Delta = 6\eta\Delta$~~

Unsure of the notation but basically two more operations that don't cancel out with other operations so  $4\eta\Delta + 2\eta\Delta = 6\eta\Delta$

### Q2. Part 2

All gradient calcs finish at 0.5s

$$t = 0.5s + 0.3s \text{ (aggregation)} + 0.25s \text{ (update)} \\ = 1.05s$$