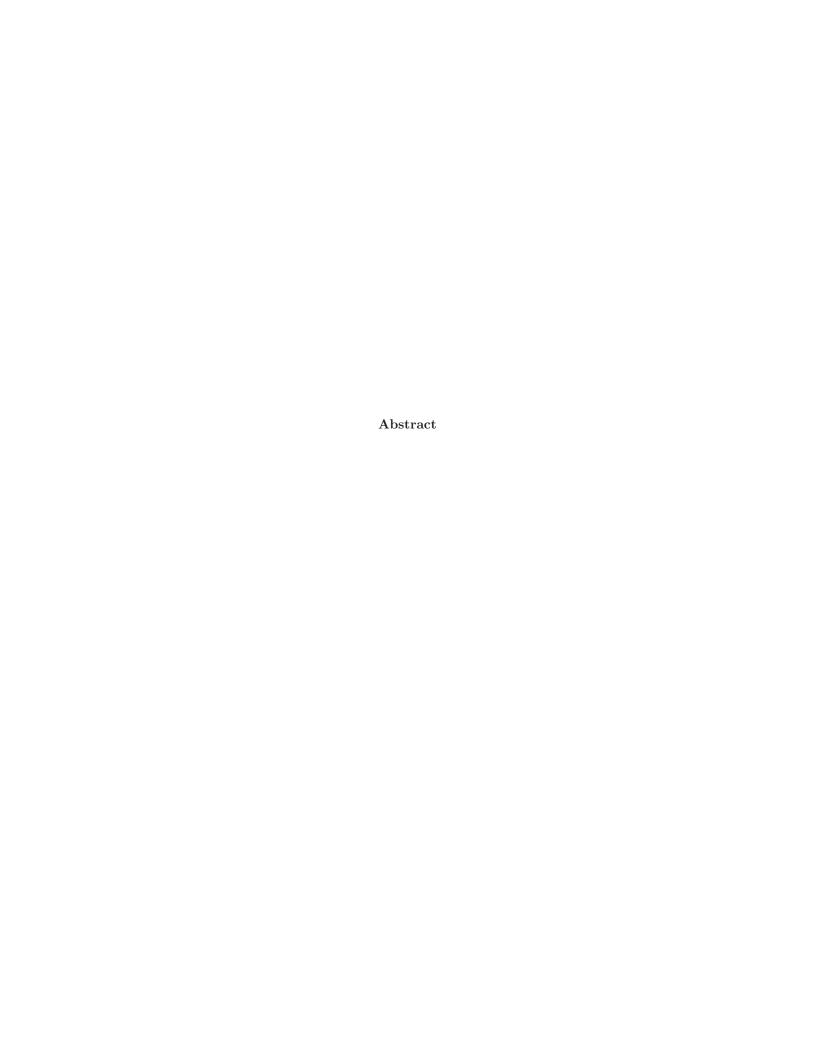
## Interface for Model Parallelism

**Author** Dr. Martin Kühn

May 14, 2020



# Contents

1	Definition	<b>2</b>
	1.1 Input parameter	2
2	Tarantella	3
3	Convolutionial layer	4
	3.1 Callback functions	5
	3.2 Implementation	5
	3.2.1 Forward pass	5
	3.2.2 Backward pass	5
4	Instance Normalisation	7
	4.1 Callback functions	7
	4.2 Implementation	8
	4.2.1 Forward pass	8
	4.2.2 Backward pass	8
5	Concat	9

# Definition

### 1.1 Input parameter

r is the rank of the current process.

 $M=r_0,r_1,\ldots$  is the full machine consisting of all the ranks. The list is ordered. |M| is the number of ranks.

D is the decompostion of a tensor over the ranks of the machine as follows  $\mathcal{N}\times\mathcal{C}\times\mathcal{H}\times\mathcal{W}.$ 

# Tarantella

Tarantella should provide a function that returns a pointer to a buffer in a GPI segment.

get\_segment\_buffer(size)

# Convolutionial layer

The layer must have global parameters describing the global parameters. These are N, C, F, H, W. N is the batch size, C the number of input channels, F the number of output channels or filters, H the height and W the width of the images.

Additionally for a distributed convolutional layer we need to know the following.

The own rank r and the machine M.

The distribution of the input tensor over the dimensions  $\bar{\mathcal{N}} \times \bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}}$ . Here  $\bar{\mathcal{N}}$  is not the number of samples but the number of instances this dimension is distributed over and so on. So  $\bar{\mathcal{N}} = 1$  means no distribution.

The distribution of the output tensor over the dimensions  $\hat{\mathcal{N}} \times \hat{\mathcal{C}} \times \hat{\mathcal{H}} \times \hat{\mathcal{W}}$ . Maybe we want to give strides for distribution of the dimensions over the ranks.  $\tilde{\mathcal{N}} \times \tilde{\mathcal{C}} \times \tilde{\mathcal{P}}$  or  $\tilde{\mathcal{N}} \times \tilde{\mathcal{C}} \times \tilde{\mathcal{P}}$ . Here stands  $\tilde{\mathcal{P}}$  for  $\tilde{\mathcal{H}} \times \tilde{\mathcal{W}}$ .

Further we allow

$$\bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{\mathcal{P}}$$
 (3.1)

$$\bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{\mathcal{C}} \times \tilde{\mathcal{P}} \tag{3.2}$$

$$\bar{\mathcal{N}} \times \bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{\mathcal{N}} \times \tilde{\mathcal{C}} \times \tilde{\mathcal{P}} \leq |M| \tag{3.3}$$

$$\hat{\mathcal{H}} \times \hat{\mathcal{W}} \leq \breve{\mathcal{P}} \tag{3.4}$$

$$\hat{C} \times \hat{\mathcal{H}} \times \hat{\mathcal{W}} \leq \tilde{C} \times \tilde{\mathcal{P}} \tag{3.5}$$

$$\hat{\mathcal{N}} \times \hat{\mathcal{C}} \times \hat{\mathcal{H}} \times \hat{\mathcal{W}} \leq \breve{\mathcal{N}} \times \breve{\mathcal{C}} \times \breve{\mathcal{P}} \leq |M| \tag{3.6}$$

Later on we would need a parameter that defines distribution of the weights. Currently we are distributing the images only, so the weights are copied and synchronized always.

Further we would need resources like GPI segments and allreduces.

In summary we have

```
\begin{array}{l} (N,C,F,H,W) \\ r \\ M \\ \bar{\mathcal{N}} \times \bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}} \\ \tilde{\mathcal{N}} \times \hat{\mathcal{C}} \times \hat{\mathcal{P}} \\ \hat{\mathcal{N}} \times \hat{\mathcal{C}} \times \hat{\mathcal{H}} \times \hat{\mathcal{W}} \\ \tilde{\mathcal{N}} \times \dot{\mathcal{C}} \times \tilde{\mathcal{P}} \\ (\text{distribution pattern weights}) \\ \text{GPI segmnts all reduce double buffer} \end{array}
```

#### 3.1 Callback functions

The ranks to synchronize the weights is returned by a callback function.

get\_ranks\_synch\_weights().

The size of the gradients is returned.

get\_sizeof\_gradients()

Tarantella has to create an allreduce (double buffer) for this layer if the return of get\_ranks\_synch\_weights() is not empty. The size is given by get\_sizeof\_gradients(). get\_sizeof\_input\_tensor() get\_sizeof\_output\_tensor()

### 3.2 Implementation

Tarantella creates an allreduce of size get\_sizeof\_gradients() on the ranks given by get\_ranks\_synch\_weights(). Tarantella delivers pointers within GPI segments to store the local gradients in.

The implementation puts the local gradients directly into the segments.

The Tarantella function get\_segment\_buffer(size) is used to allocate a buffer in a GPI segment. On the first call they exchange a pointer to the exchange buffers with all their neighbors.

#### 3.2.1 Forward pass

- start collecting all the input data needed
- calculate convolution with local data
- calculate convolution with remote data
- fuse both data

### 3.2.2 Backward pass

- start collecting all the input data needed
- calculate gradient weights with local data
- calculate error signal with local data

- calculate gradient weights with remote data
- $\bullet\,$  calculate error signal with remote data
- fuse data

### Instance Normalisation

The layer must have global parameters describing the global parameters. These are N, C, H, W. N is the batch size, C the number of input channels, H the height and W the width of the images.

Additionally for a distributed convolutional layer we need to know the following.

The own rank r and the machine M.

The distribution of the input tensor over the dimensions  $\bar{\mathcal{N}} \times \bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}}$ . The dimension of the input tensor is the dimension of the output tensor.

Maybe we want to give strides for distribution of the dimensions over the ranks.  $\tilde{\mathcal{N}} \times \tilde{\mathcal{C}} \times \tilde{\mathcal{P}}$ . Here stands  $\tilde{\mathcal{P}}$  for  $\tilde{\mathcal{H}} \times \tilde{\mathcal{W}}$ .

Further we allow

$$\bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{\mathcal{P}}$$
 (4.1)

$$\bar{C} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{C} \times \tilde{\mathcal{P}} \tag{4.2}$$

$$\bar{\mathcal{N}} \times \bar{\mathcal{C}} \times \bar{\mathcal{H}} \times \bar{\mathcal{W}} \leq \tilde{\mathcal{N}} \times \tilde{\mathcal{C}} \times \tilde{\mathcal{P}} \leq |M| \tag{4.3}$$

Further we would need resources like allreduces.

In summary we have

r

M

 $\bar{\mathcal{N}}\times\bar{\mathcal{C}}\times\bar{\mathcal{H}}\times\bar{\mathcal{W}}$ 

$$ilde{\mathcal{N}} imes ilde{\mathcal{C}} imes ilde{\mathcal{P}}$$

allreduce double buffer weights allreduce double buffer norm

The allreduce double buffer norm is not part of the gradients update but is used inside the layer to reduce the norm. Latency is critical here.

### 4.1 Callback functions

The ranks to synchronize the weights is returned by a callback function.

```
get_ranks_synch_weights().
```

The size of the gradients is returned.

get\_sizeof\_gradients().

The ranks to synchronize the norm is returned by a callback function get\_ranks\_synch\_norm().

The size of the norm is returned by get\_sizeof\_norm().

Tarantella has to create two all reduces (double buffer) for this layer if the return of get\_ranks\_synch\_weights () andget\_ranks\_synch\_norm() are not empty. The size is given by get\_size of\_gradients() and get\_ranks\_synch\_norm(). get\_size of\_tensor()

### 4.2 Implementation

#### 4.2.1 Forward pass

- compute local norm
- allreduce norm
- perform scaling

#### 4.2.2 Backward pass

Not sure if it works like this:

- calculate local gradient of norm
- calculate local error signal of norm
- allreduce local data
- perform backward gradient of norm.
- perform backward gradient of error signal

Concat