

# Framework Design

2019/09/07

# Horovod/Bluefog Consideration

Comm.
MPI
NCCL
Gloo

Framework
TensorFlow
PyTorch
MXNet

Device
CPU
GPU

Type
Int
Double
Float

Alg. Style
Async
Sync

Operation
Send/Recv
Broadcast
Reduce

Comm Style
All
Neighbor
Point to Point

Fault Tolerance
Timeout
Backup
Self-healing?

# Bluefog Consideration (First Priority)

Comm.
MPI
NCCL
<del>Gloo</del>

Framework
<del>TensorFlow</del>
PyTorch
<del>MXNet</del>

Device
CPU
GPU

Type
<del>Int</del>
<del>Double</del>
Float

Alg. Style
Async
Sync

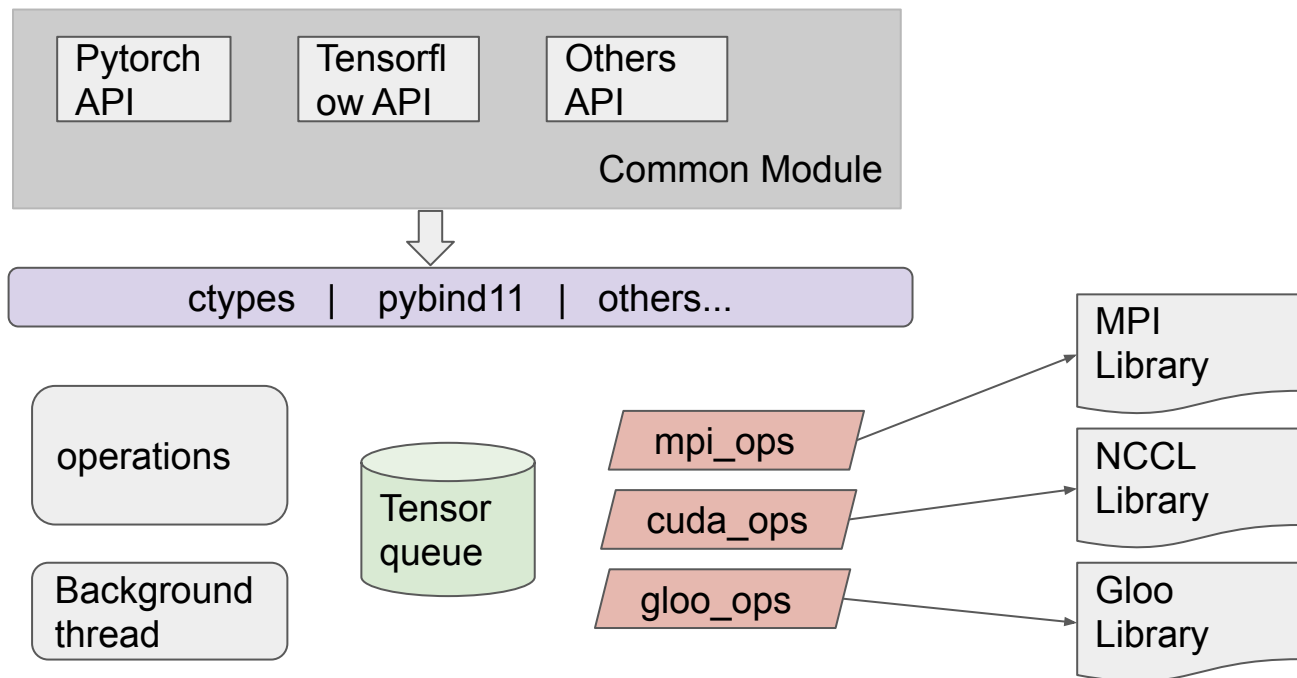
Operation
Send/Recv
Broadcast
Reduce

Comm Style
<del>All</del>
Neighbor
<del>Point to Point</del>

Fault Tolerance
<del>Timeout</del>
<del>Backup</del>
<del>Self-healing?</del>

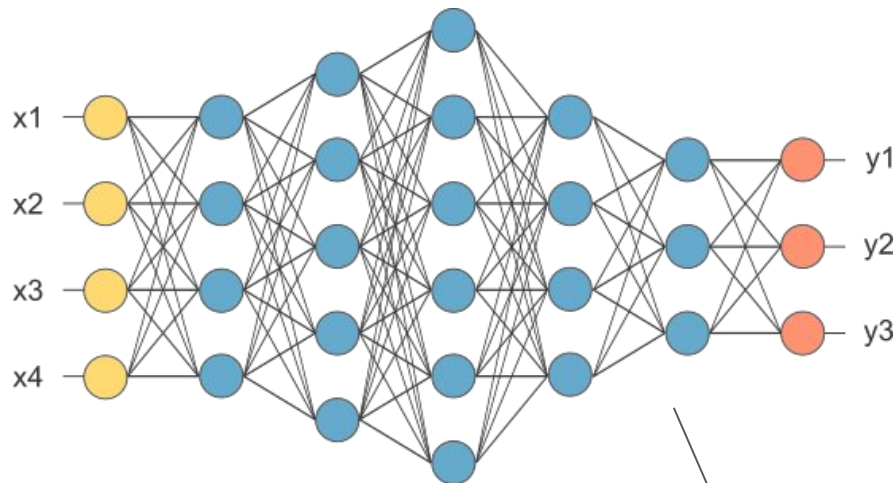
# Code Design

In order to manage and be able to easily extensibility for code, we have to decouple the logic well. A big picture of horovod code structure:



# Horovod main idea

## Paralleling communication with computation



Computation Thread

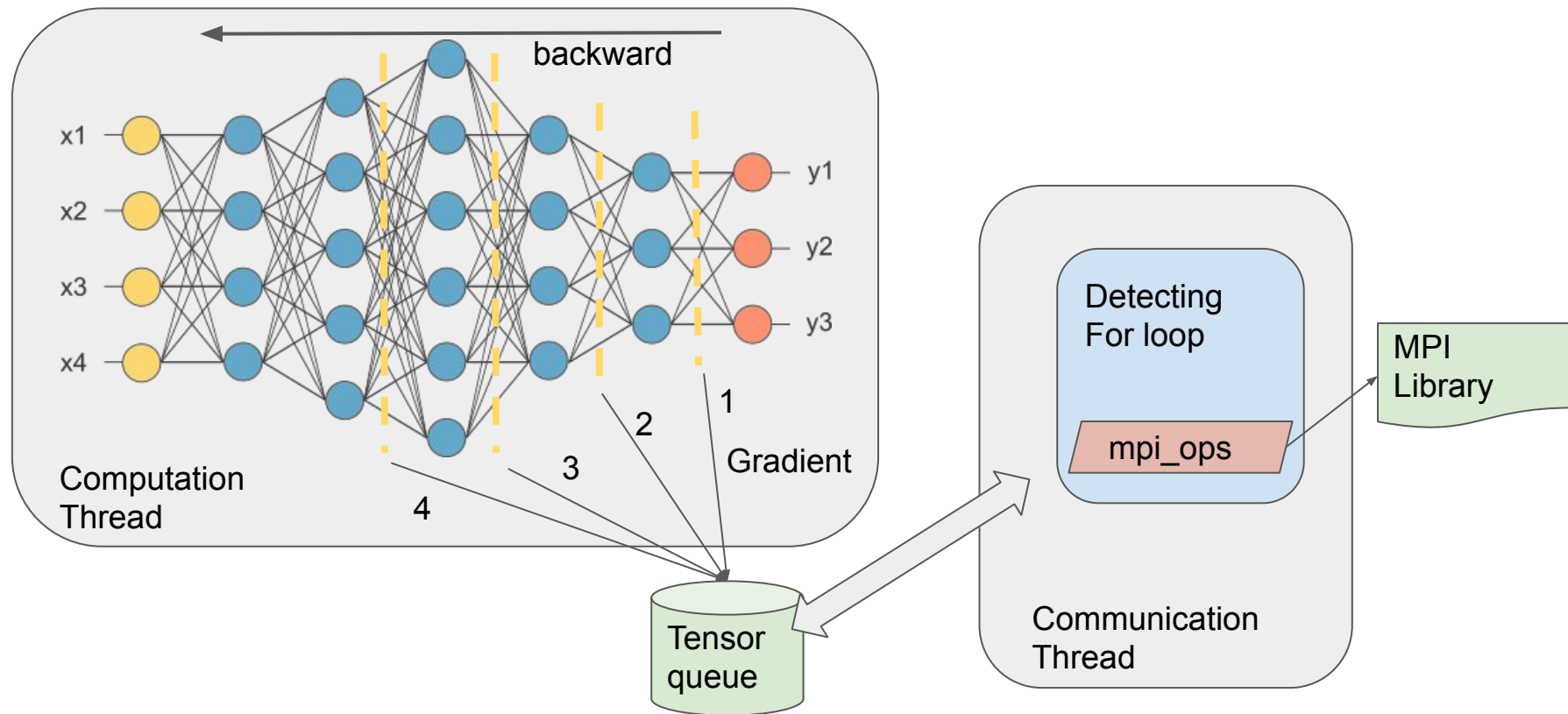


Communication Thread

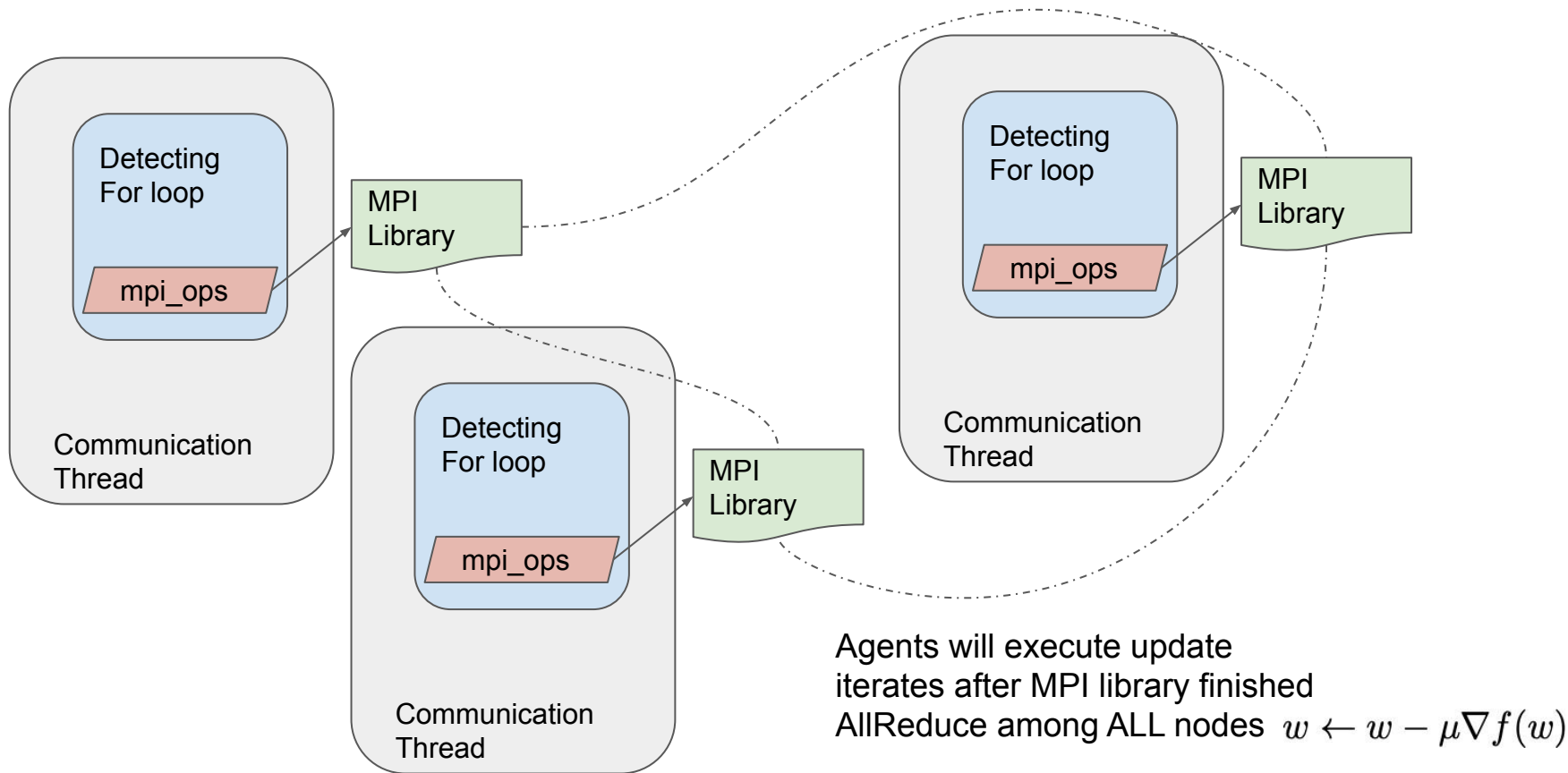
1. Forward ----  $y = f(w)$
2. Backward ---  $\nabla f(w)$
3. Step ----  $w \leftarrow w - \mu \nabla f(w)$

Each line represents the parameters  $w_{\{i,m\}^k}$  in the forward pass and its corresponding gradient in the backward pass

# Paralleling communication with computation



# Ring All-Reduce and Synchronized Update

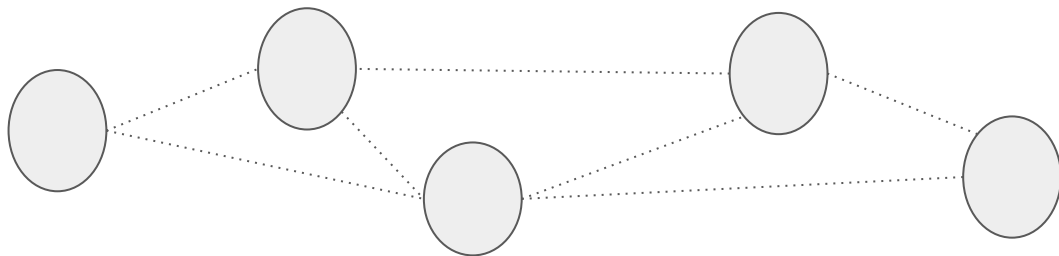


# Bluefog Three Main Differences

1. Computing average on iterate  $w$  instead of gradient

$$w_{i,k} \leftarrow \sum_{\ell \in \mathcal{N}_k} a_{k\ell} w_{i-1,\ell}$$

2. Local/neighbor communication (need topology)

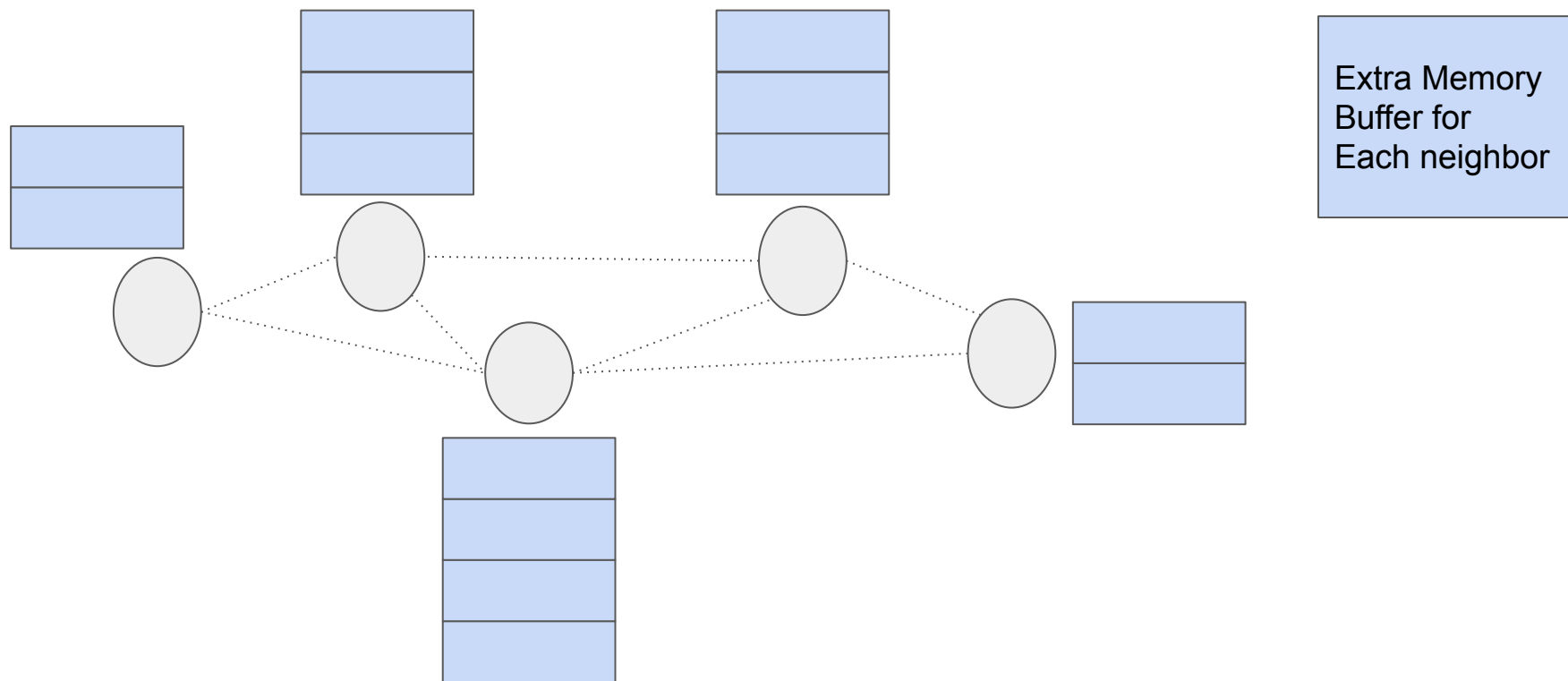


`file:///Users/ybc/Documents/dev/picpic-blue-fog/python/gossip_simulations/simulations_result/push_sum.html`

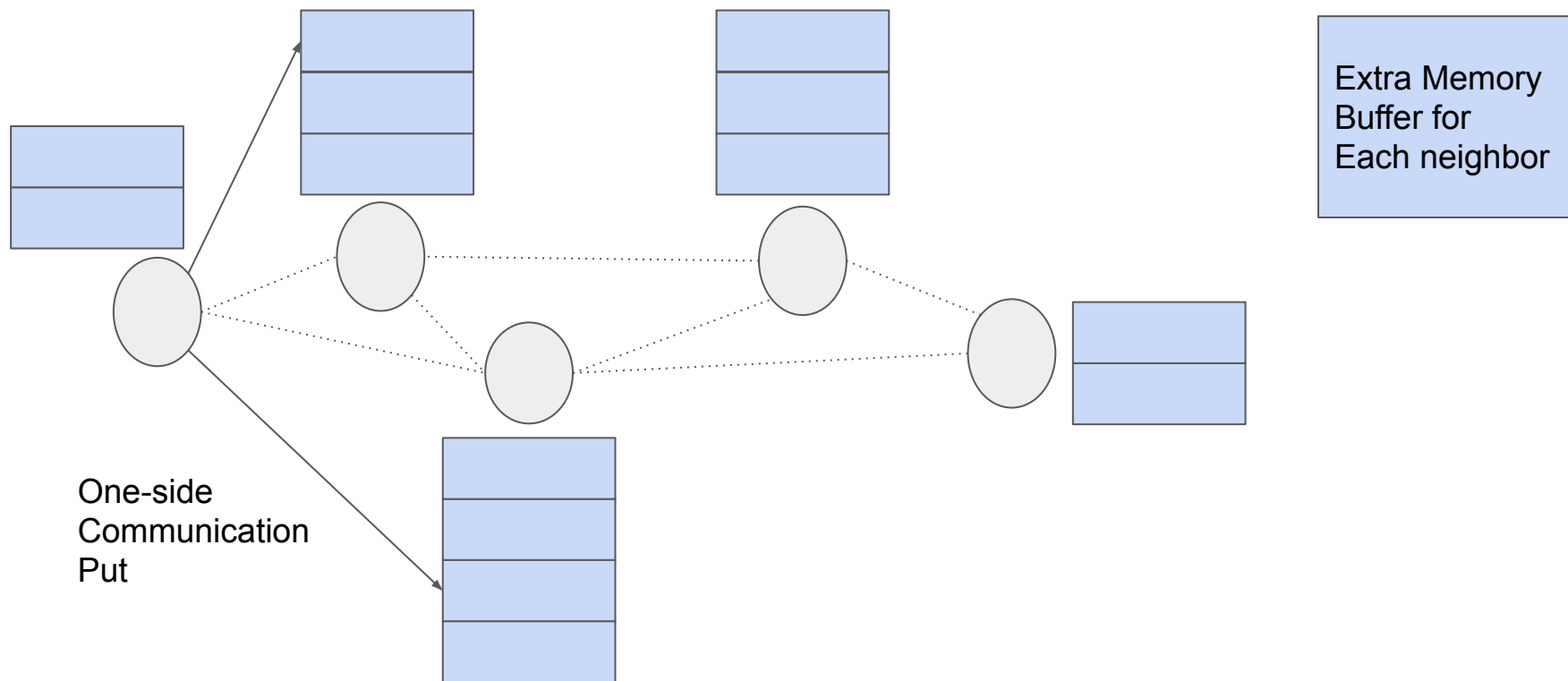
3. Asynchronous algorithm



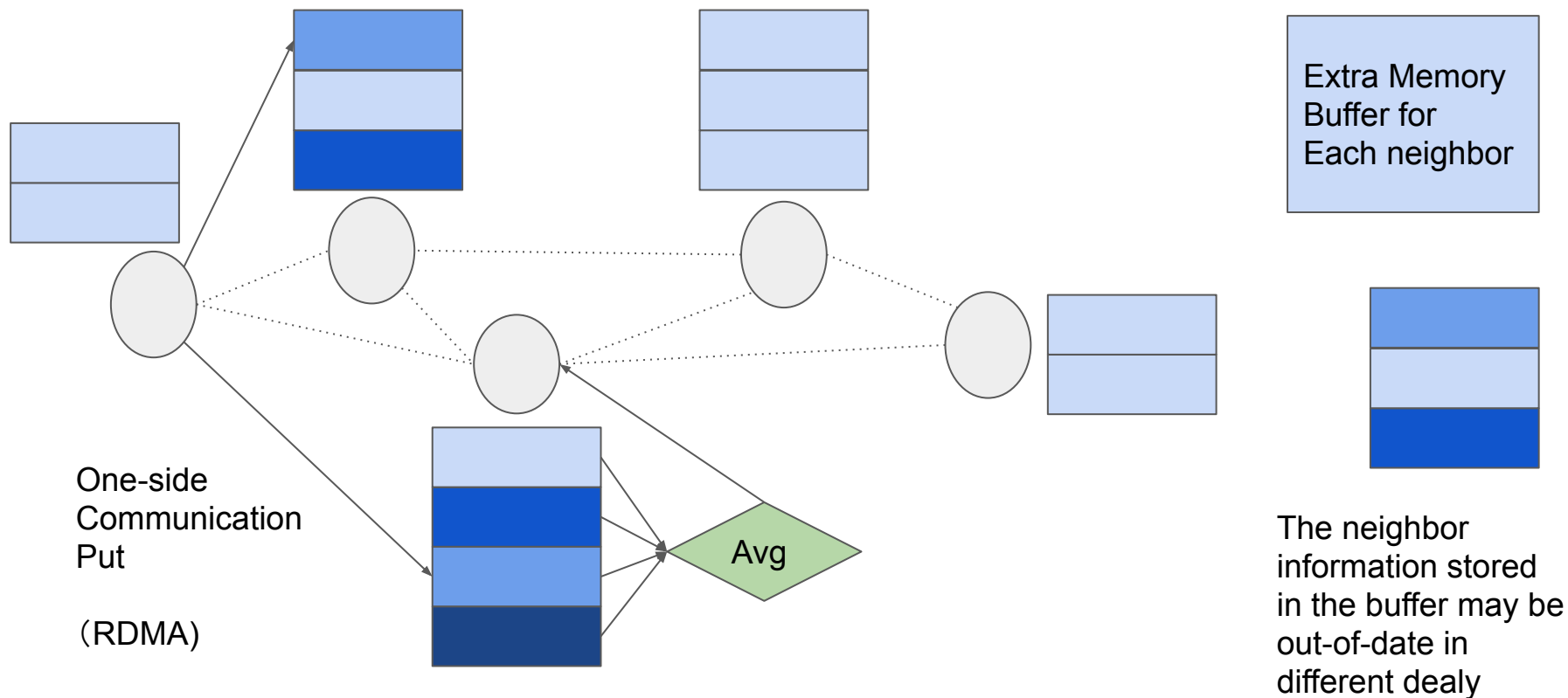
# Our asynchronization



# Our asynchronization



# Our asynchronization



# Our asynchronization

