

# MPI 3.0 Neighbourhood Collectives

---

Advanced Parallel Programming



# Overview

- Review of topologies in MPI
- MPI 3.0 includes new neighbourhood collective operations:
  - MPI\_Neighbor\_allgather[v]
  - MPI\_Neighbor\_alltoall[v|w]
- Example usage:
  - Halo-exchange can be done with a single MPI communication call
- Practical tomorrow:
  - Replace all point-to-point halo-exchange communication with a single neighbourhood collective in your MPP coursework code

# Topology communicators (review 1)

- Regular n-dimensional grid or torus topology
  - MPI\_CART\_CREATE
- General graph topology
  - MPI\_GRAPH\_CREATE
    - All processes specify all edges in the graph (not scalable)
- General graph topology (distributed version)
  - MPI\_DIST\_GRAPH\_CREATE\_ADJACENT
    - All processes specify their incoming and outgoing neighbours
  - MPI\_DIST\_GRAPH\_CREATE
    - Any process can specify any edge in the graph (too general?)

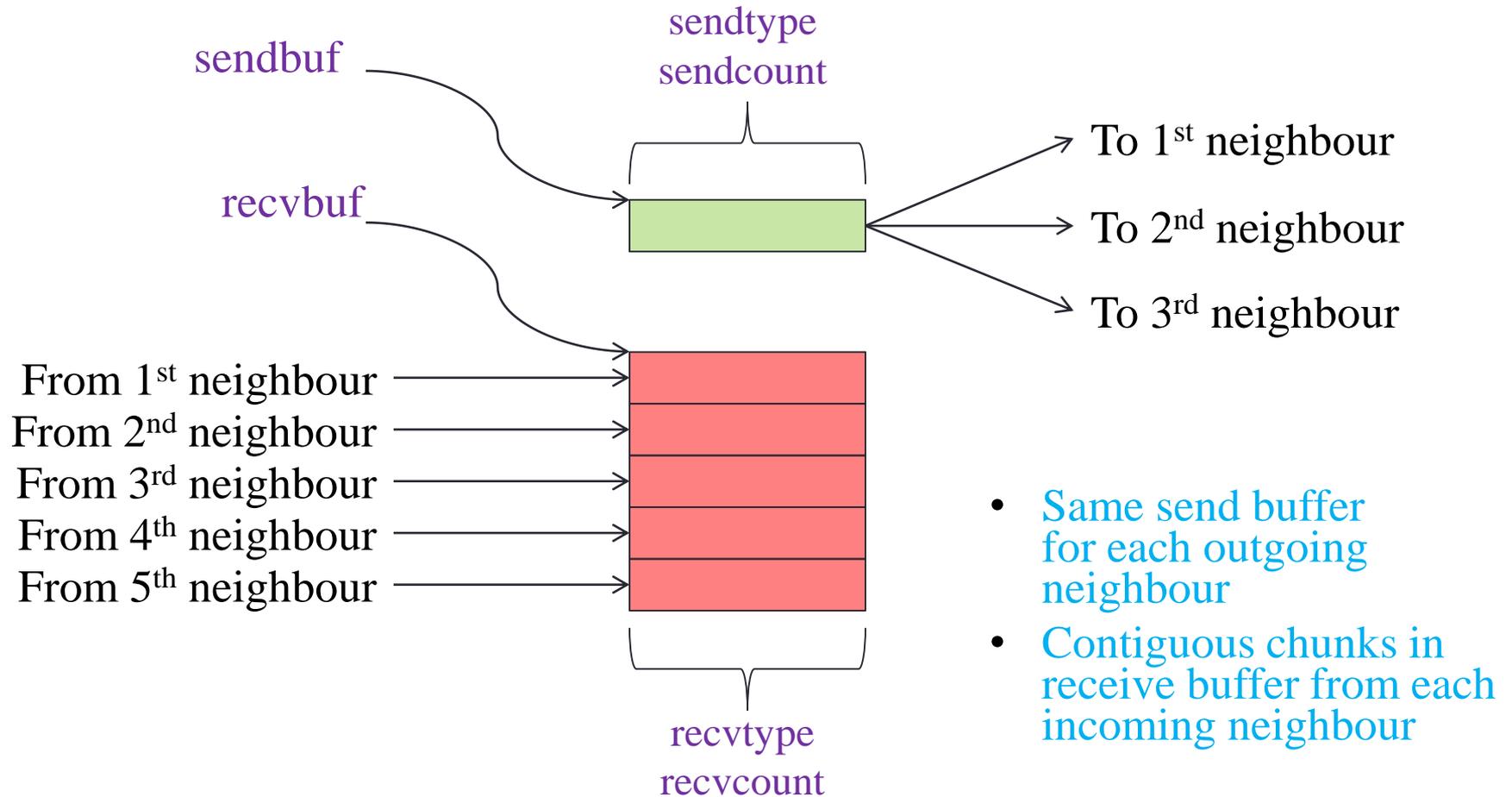
# Topology communicators (review 2)

- Testing the topology type associated with a communicator
  - MPI\_TOPO\_TEST
- Finding the neighbours for a process
  - MPI\_CART\_SHIFT
  - Find out how many neighbours there are:
    - MPI\_GRAPH\_NEIGHBORS\_COUNT
  - Get the ranks of all neighbours:
    - MPI\_GRAPH\_NEIGHBORS
  - Find out how many neighbours there are:
    - MPI\_DIST\_GRAPH\_NEIGHBORS\_COUNT
  - Get the ranks of all neighbours:
    - MPI\_DIST\_GRAPH\_NEIGHBORS

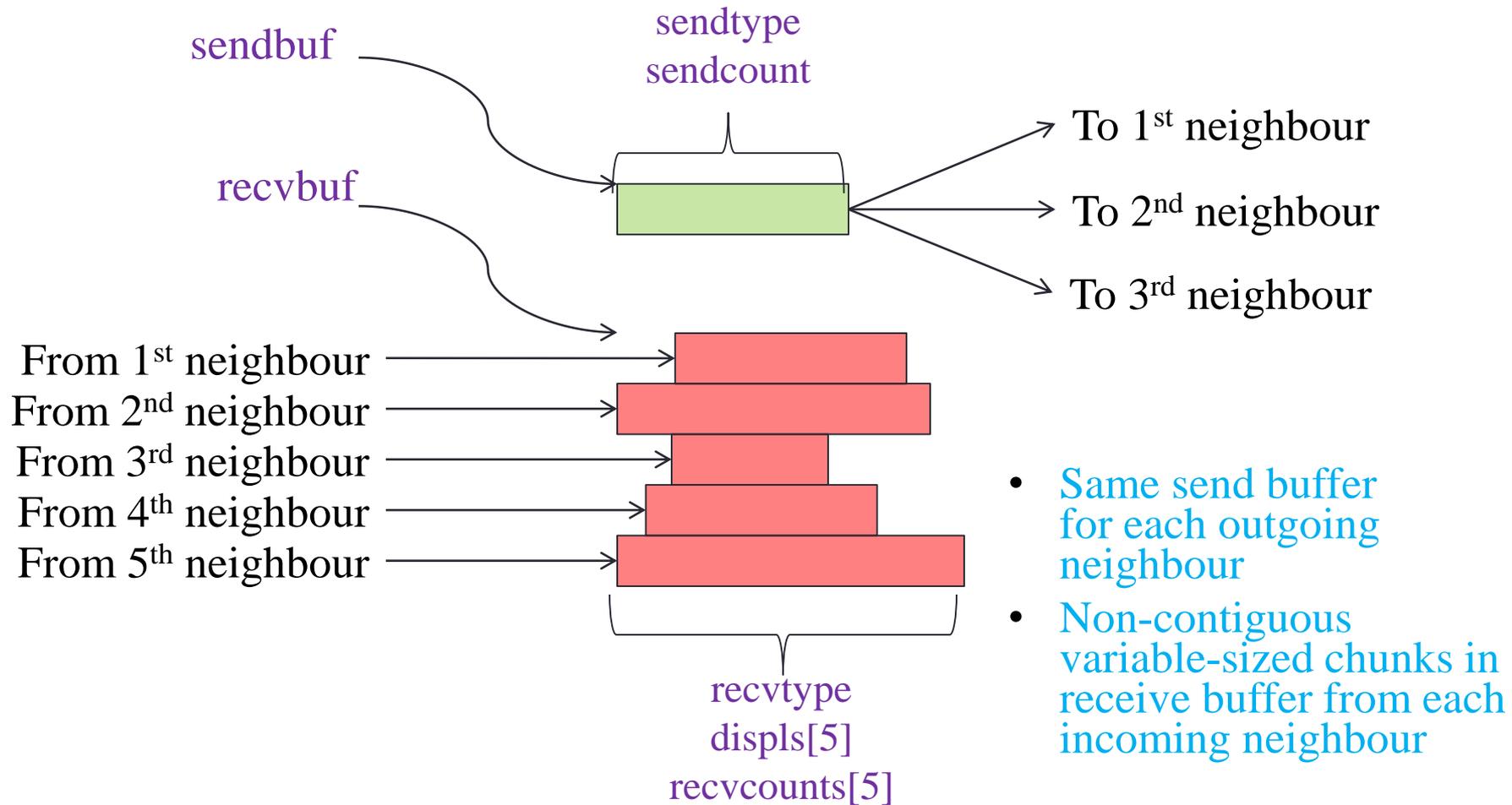
# Neighbourhood collective operations

- See section 7.6 in MPI 3.0 for blocking functions
  - See section 7.7 in MPI 3.0 for non-blocking functions
  - See section 7.8 in MPI 3.0 for an example application
    - But beware of the mistake(s) in the example code!
- MPI\_[N|In]ighbor\_allgather[v]
  - Send one piece of data to all neighbours
  - Gather one piece of data from each neighbour
- MPI\_[N|In]ighbor\_alltoall[v|w]
  - Send different data to each neighbour
  - Receive different data from each neighbour
- Use-case: regular or irregular domain decompositions
  - Where the decomposition is static or changes infrequently
  - Because creating a topology communicator takes time

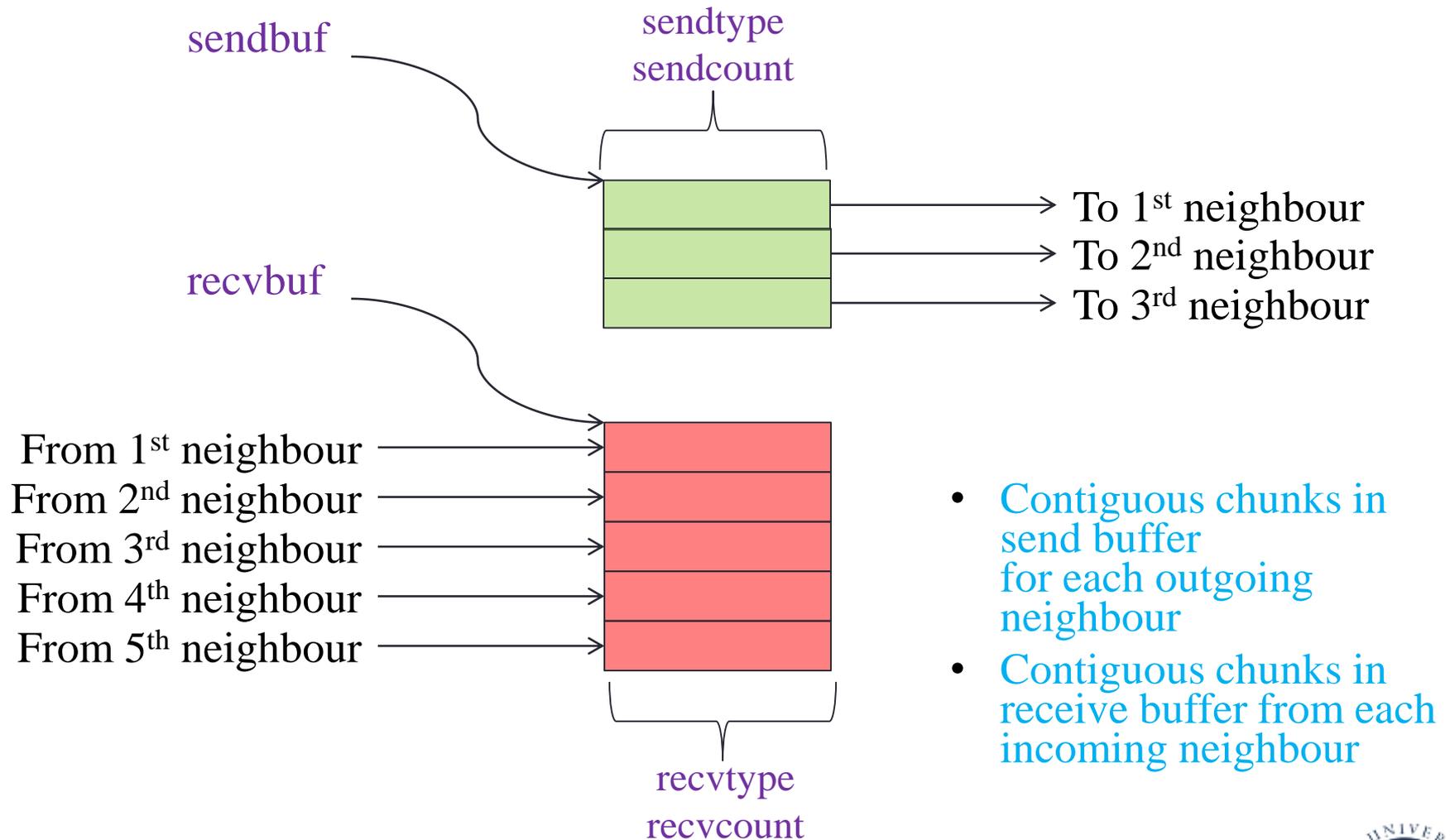
# MPI\_Neighbor\_allgather



# MPI\_Neighbor\_allgatherv

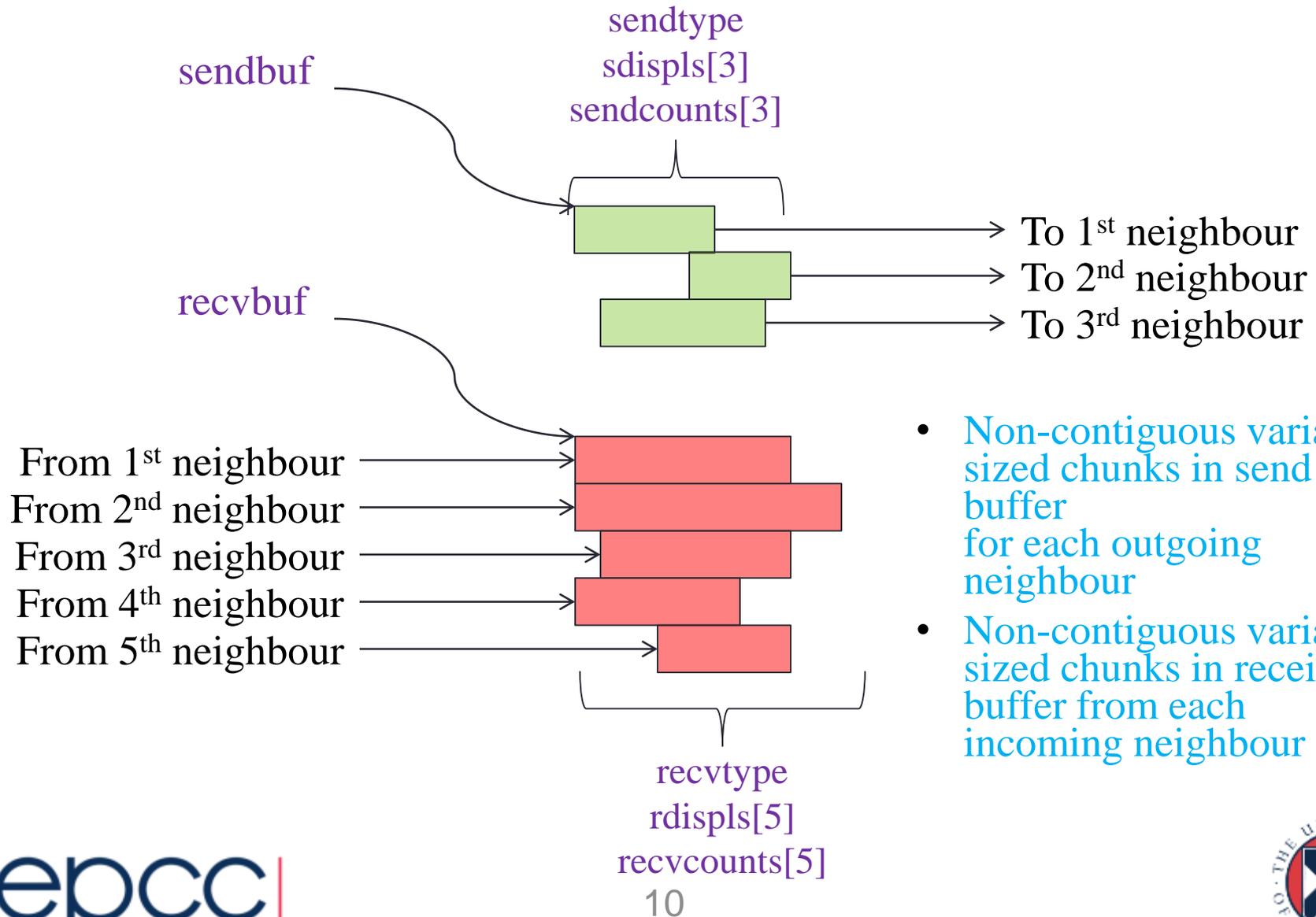


# MPI\_Neighbor\_alltoall



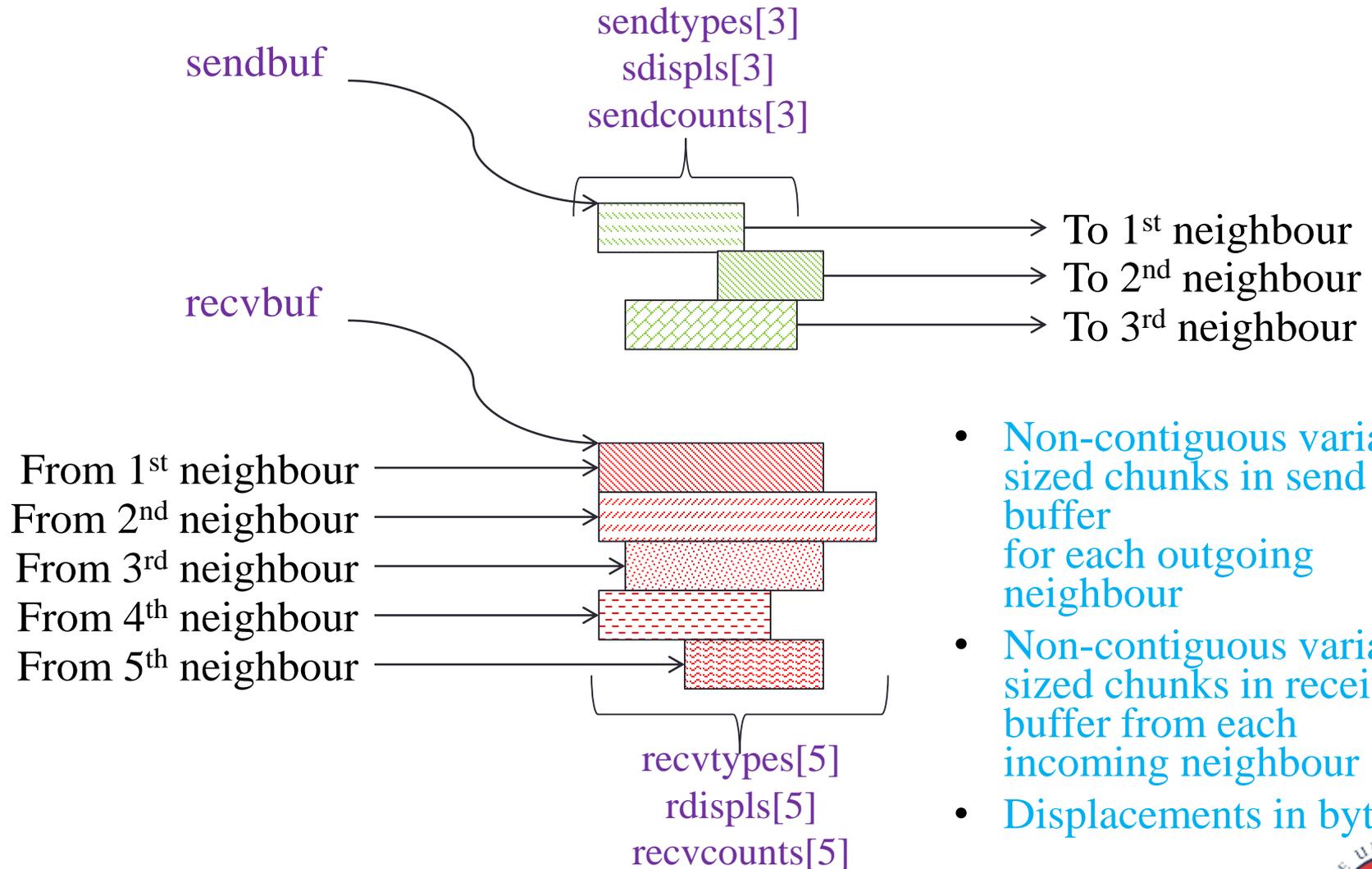
- Contiguous chunks in send buffer for each outgoing neighbour
- Contiguous chunks in receive buffer from each incoming neighbour

# MPI\_Neighbor\_alltoallv



- Non-contiguous variable-sized chunks in send buffer for each outgoing neighbour
- Non-contiguous variable-sized chunks in receive buffer from each incoming neighbour

# MPI\_Neighbor\_alltoallw



- Non-contiguous variable-sized chunks in send buffer for each outgoing neighbour
- Non-contiguous variable-sized chunks in receive buffer from each incoming neighbour
- Displacements in bytes

# MPI\_Neighbor\_alltoallw

```
for (int i=0;i<4;++i) {  
    sendcounts[i] = 1;  
    recvcounst[i]=1; }
```

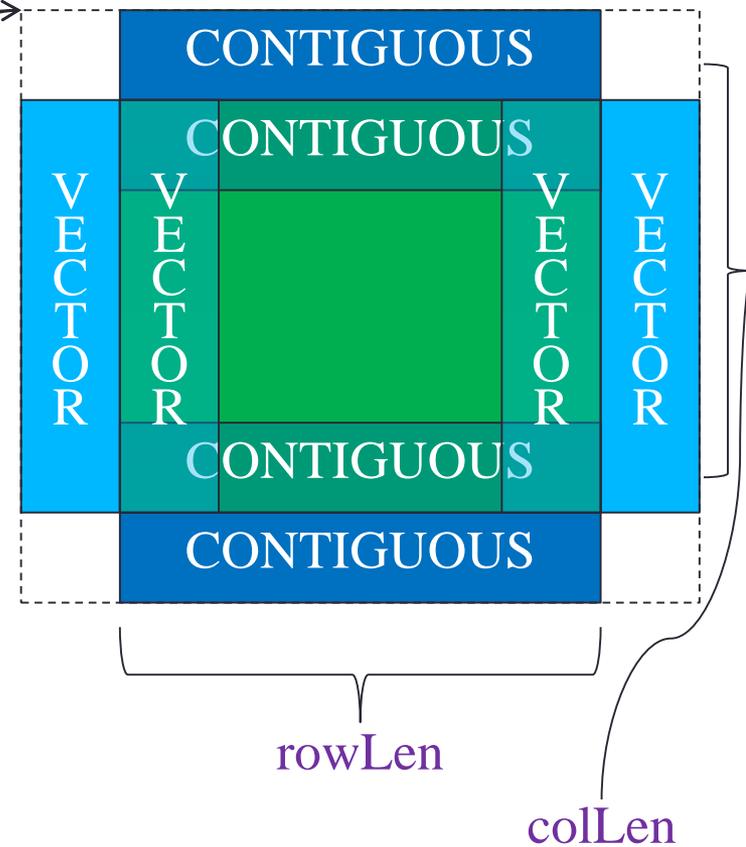
sendbuf

recvbuf

```
sendtypes[0] = contigType;  
senddispls[0] = (colLen*(rowLen+2)+1)*dblesize;  
sendtypes[1] = contigType;  
senddispls[1] = (1*(rowLen+2)+1) *dblesize;  
sendtypes[2] = vectorType;  
senddispls[2] = (1*(rowLen+2)+1)*dblesize;  
sendtypes[3] = vectorType;  
senddispls[3] = (2*(rowLen+2)-2) *dblesize;
```

// similarly for recvtypes and recvdispls

```
MPI_Neighbor_alltoallw(sendbuf, sendcounts, senddispls, sendtypes,  
    recvbuf, recvcounst, recvdispls, recvtypes,  
    comm);
```



# Summary

- Useful for regular or irregular domain decomposition
  - Where the decomposition is static or changes infrequently
- Should investigate replacing point-to-point communication
  - E.g. halo-exchange communication
- With neighbourhood collective communication
  - Probably `MPI_Neighbor_alltoallw` / `MPI_Ineighbor_alltoallw`
- So that MPI can optimise the whole pattern of messages
  - Rather than trying to optimise each message individually
- And so your application code is simpler and easier to read