

# Halo Exchanges and Tausch

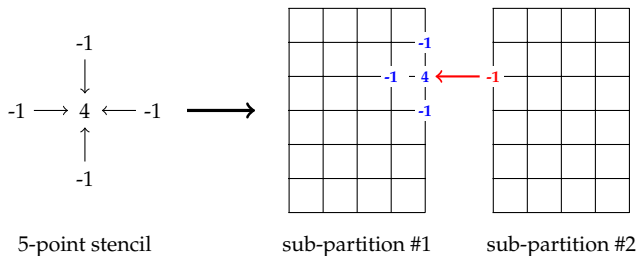
Lukas Spies

Scientific Computing Seminar  
University of Illinois

2<sup>nd</sup> October 2018

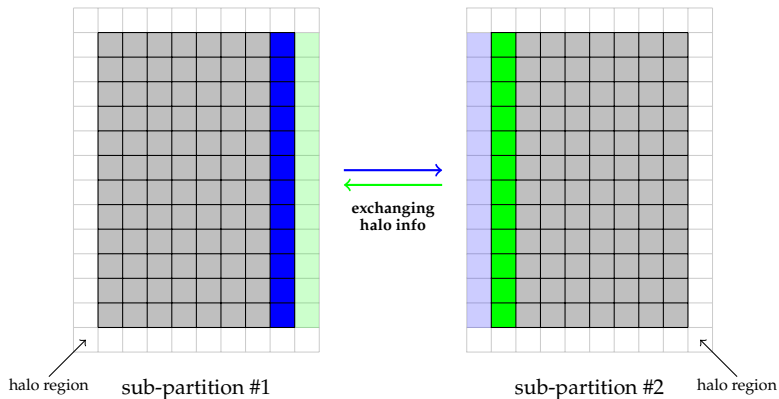
## Sample problem:

5-point stencil on two sub-partitions of some domain



Stencil weights known, node values unknown  
 $\Rightarrow$  Need for communication!

**Solution:** Store extra values as halo around domain



**New problem:** Need to keep updating halos

## Existing Solutions:

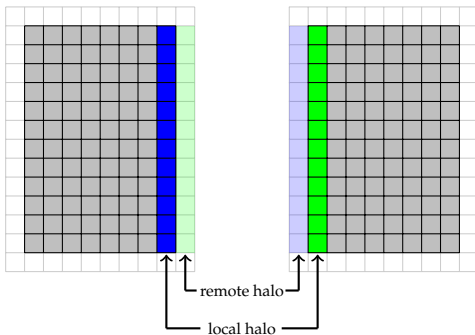
- ▶ Libraries like MSG (previously used in Cedar)  
→ complicated API, not designed for parallel computing
- ▶ Custom implementation  
→ development work needs to be repeated for every project, takes time away from actual development

⇒ Tausch to the rescue!



## Glossary:

- ▶ **remote halo:** values needed by one sub-partition that are computed on another sub-partition.
- ▶ **local halo:** values computed by one sub-partition that are needed by another sub-partition.



## Tausch overview:

1. **Defining halo info:**  
addLocalHaloInfo()  
addRemoteHaloInfo()
2. **Extract values of halo:**  
packSendBuffer()
3. **Send/Recv values:**  
send()  
recv()
4. **Move received values back into main buffer:**  
unpackRecvBuffer()

```
#include "tausch.h"

int main(int argc, char** argv) {

    MPI_Init(&argc, &argv);

    Tausch<double> *tausch = new Tausch<double>(MPI_DOUBLE, MPI_COMM_WORLD);

    (...)

    int haloidSend = tausch->addLocalHaloInfo(sendindices);
    int haloidRecv = tausch->addRemoteHaloInfo(recvindices);

    (...)

    tausch->packSendBuffer(haloidSend, 0, buf);
    tausch->send(haloidSend, 1, right);

    tausch->recv(haloidRecv, 1, left);
    tausch->unpackRecvBuffer(haloidRecv, 0, buf);

    (...)

    MPI_Finalize();

    return 0;

}
```

```
...  
Tausch<double> *tausch = new Tausch<double>(MPI_DOUBLE, MPI_COMM_WORLD);  
...
```

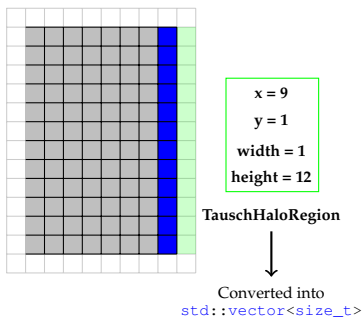
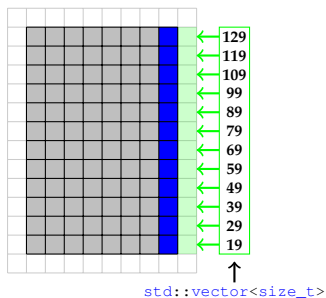
## Constructing new Tausch object:

- ▶ templated library: can be used with any data type supported by MPI
- ▶ To duplicate or not to duplicate communicator

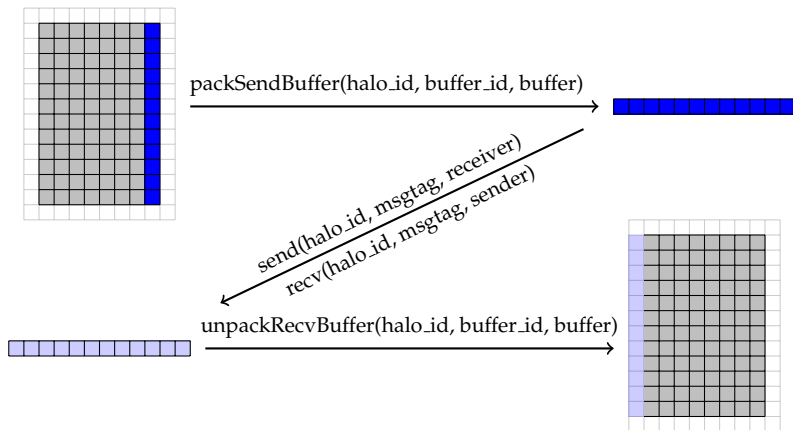


```
...  
int haloidSend = tausch->addLocalHaloInfo(sendindices);  
int haloidRecv = tausch->addRemoteHaloInfo(recvindices);  
...
```

## Two ways to specify halo indices:

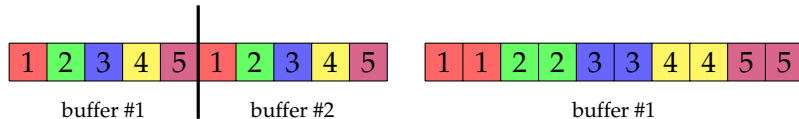


```
...  
tausch->packSendBuffer(haloidSend, 0, buf);  
tausch->send(haloidSend, 1, right);  
  
tausch->recv(haloidRecv, 1, left);  
tausch->unpackRecvBuffer(haloidRecv, 0, buf);  
...
```



## Some design considerations:

- ▶ Using MPI persistent communication decreases overhead on repeated calls to send/recv.
- ▶ Multiple buffers can be combined into one message
- ▶ Multiple values per point can be combined into one message
- ▶ Aligned memory buffers



Performance model in 2D:

$$T_{\text{perf}} = 2nT_{b2h,\text{nseq}} + 2nT_{b2h,\text{seq}} + 4nT_{\text{MPI,send/recv}} + 2nT_{h2b,\text{seq}} + 2nT_{h2b,\text{nseq}}$$

$T_{\text{perf}}$  := Performance prediction

$n$  := number of points in one dimension

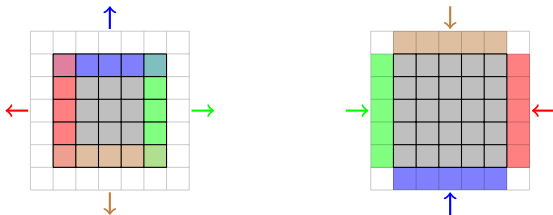
$T_{b2h,\text{seq}}$  := Time to copy one value from buffer to halo, sequential memory access

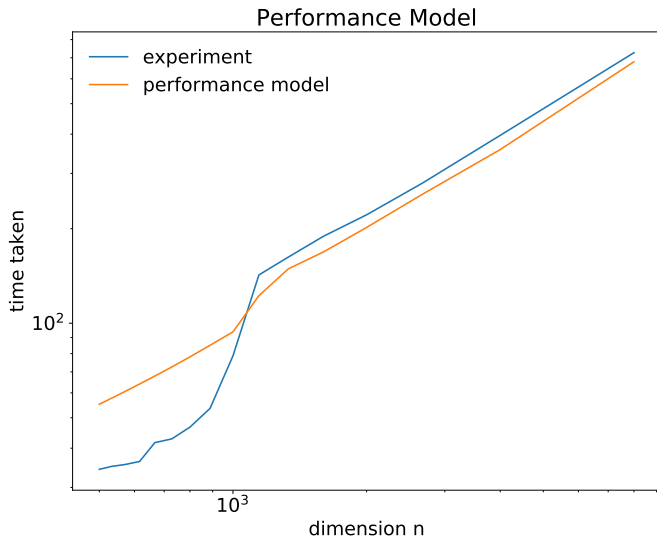
$T_{b2h,\text{nseq}}$  := Time to copy one value from buffer to halo, non-sequential memory access

$T_{h2b,\text{seq}}$  := Time to copy one value from halo to buffer, sequential memory access

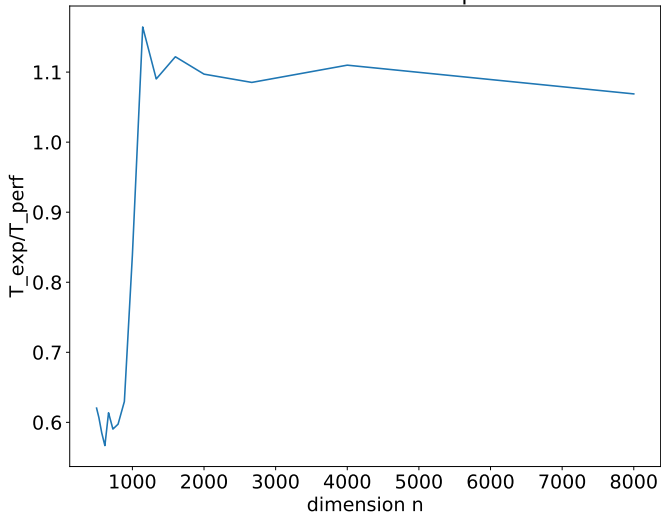
$T_{h2b,\text{nseq}}$  := Time to copy one value from halo to buffer, non-sequential memory access

$T_{\text{MPI,send/recv}}$  := Time to send one value using MPI communication





Performance Model - relative performance



## Current status and future work:

- ▶ Tausch is used already by Cedar, integration into PlasCom2 in progress
- ▶ Performance model demonstrates achievement of expected performance
- ▶ API very simple and straightforward → easy integration
- ▶ CPU/GPU communication still incomplete/unverified
- ▶ Small paper planned