



D1

D2

D3

D4

D5

C1

ChilledDoor
Chilled

Chilled

Chilled

Chilled

Chilled

MATLAB HPC Mentors

June 25, 2020

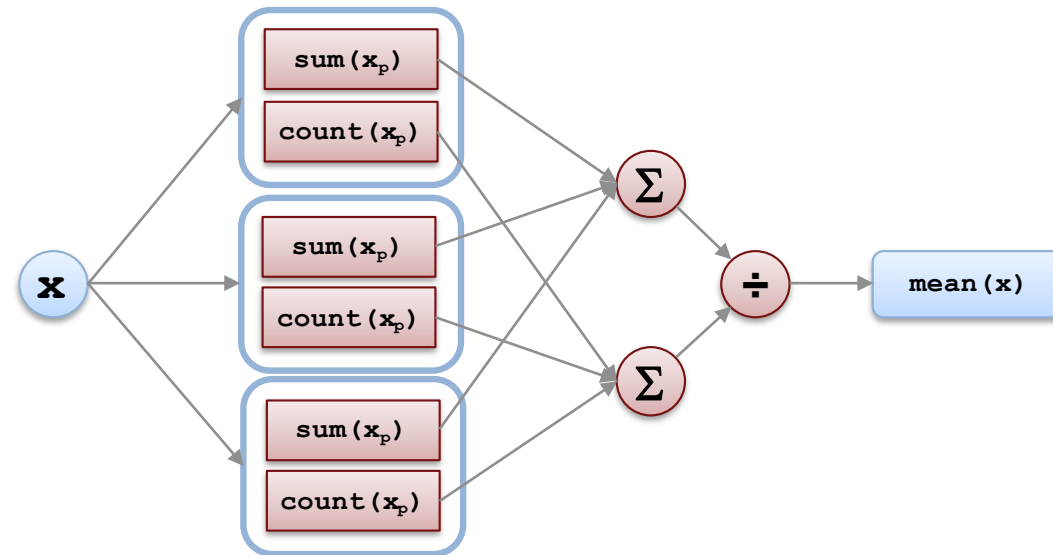
Host: **Nicholas Ide**

Guest Speakers: **Ben Tordoff & Oli Tissot**, Parallel Computing Development

Quick Updates

- No meeting in July or August – enjoy the summer, stay healthy!
- If you have issues with recent SSH and MATLAB Parallel Server, let me know.
 - Have been a few reports
 - Members of the community have found temporary solutions

Distributed Arrays: techniques and best practices for handling very large calculations



Ben Tordoff & Oli Tissot,
Parallel Computing Development

Agenda

- ➔ What are distributed arrays?
 - What can you do with distributed arrays?
 - Building a distributed array
 - Advanced manoeuvres
 - Debugging methods and tips

Remote arrays in MATLAB

Rule: take the calculation to where the data is

Normal array – calculation happens in main memory:



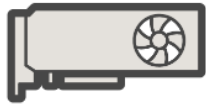
```
x = rand(...)
```

```
x_norm = (x - mean(x)) ./ std(x)
```

Remote arrays in MATLAB

Rule: take the calculation to where the data is

gpuArray – all calculation happens on the GPU:



```
x = gpuArray(...)
```

```
x_norm = (x - mean(x)) ./ std(x)
```

distributed – calculation is spread across the memory of a cluster:



```
x = distributed(...)
```

```
x_norm = (x - mean(x)) ./ std(x)
```

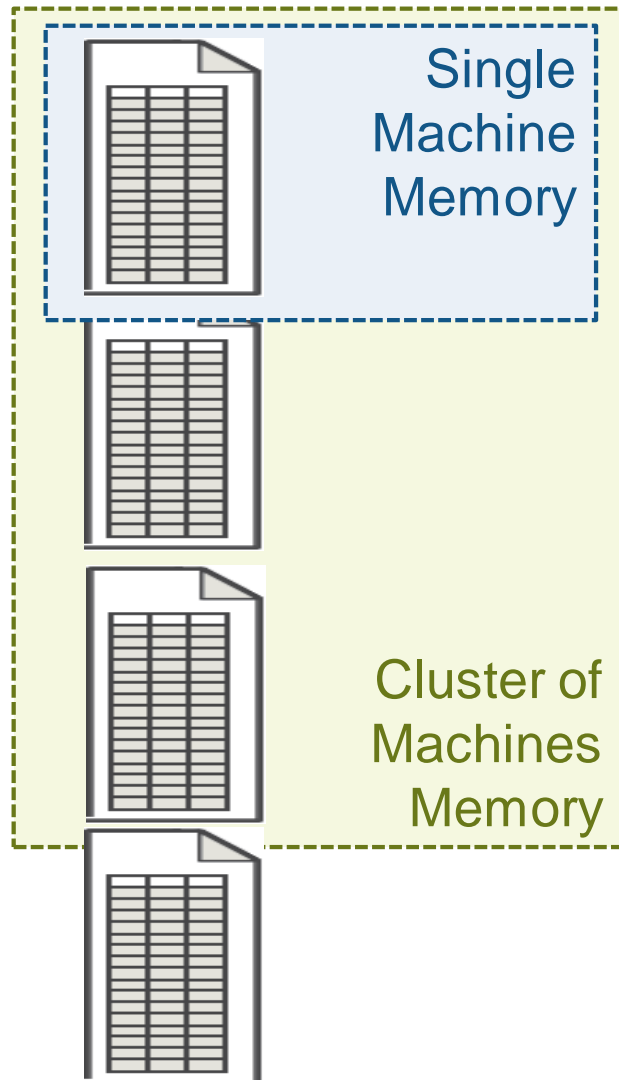
tall – calculation is performed by stepping through files:



```
x = tall(...)
```

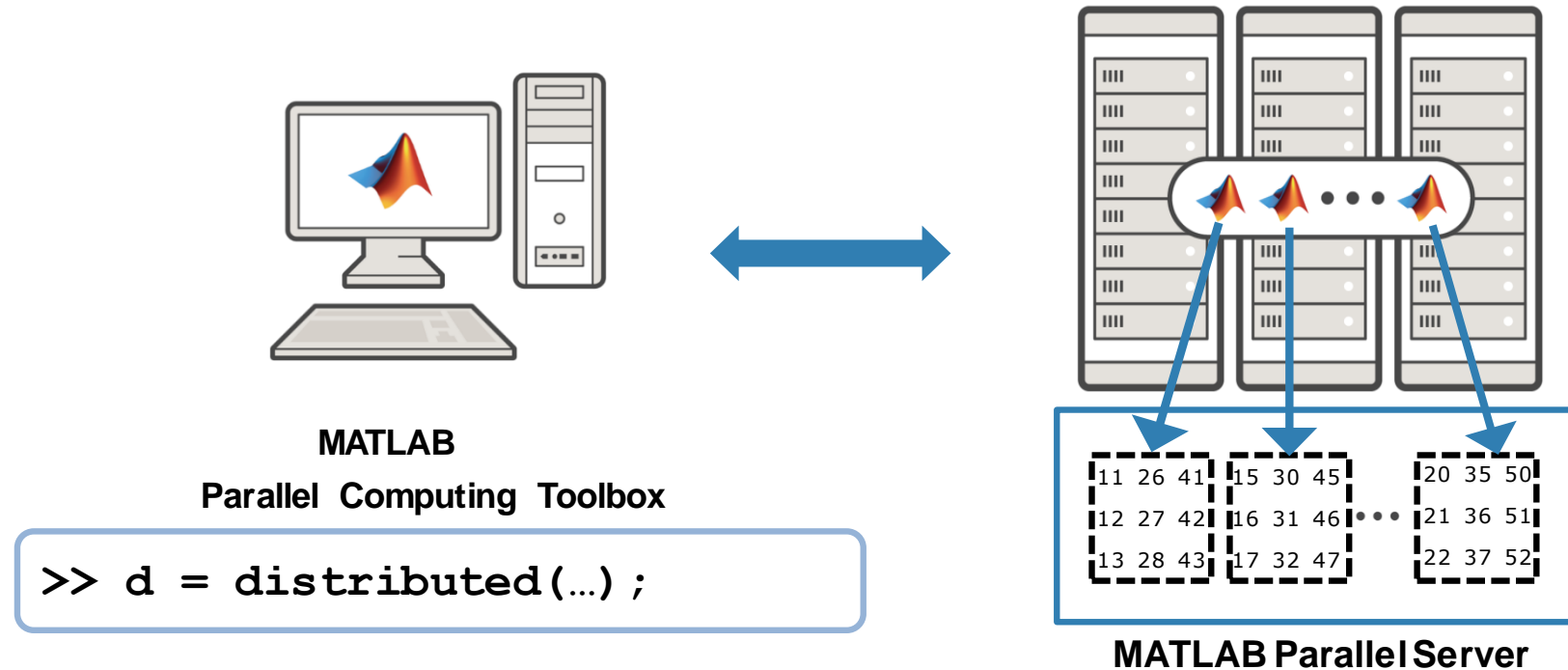
```
x_norm = (x - mean(x)) ./ std(x)
```

When should I reach for distributed?



- If your data fits in memory, just use MATLAB normally
- If it won't fit on one machine, maybe it can be split across the combined memory of a cluster of machines? Use `distributed` arrays
- If it won't fit in the combined memory of a cluster of machines then use `tall` arrays

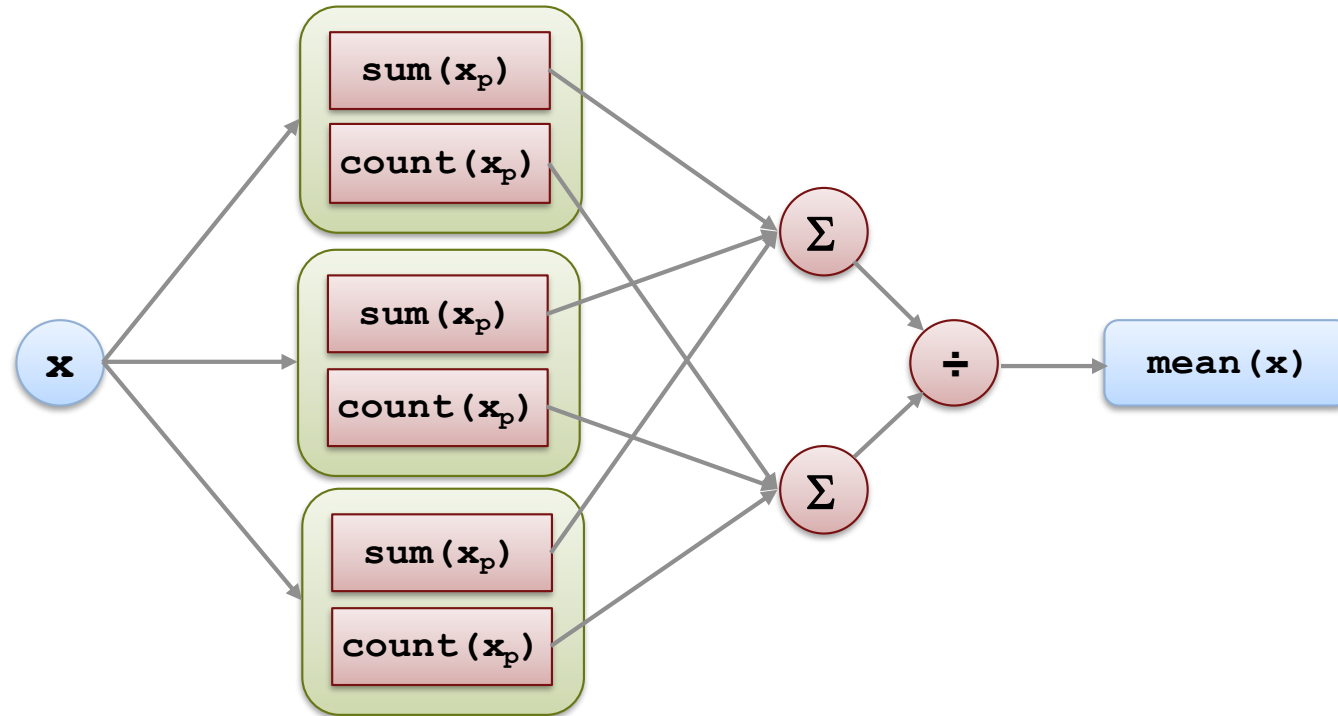
Using distributed arrays



- Use the memory of multiple machines as though it was your own
- Client sees a “normal” MATLAB variable
- Work happens on cluster

Distributed Algorithms

```
>> m = mean (d) ;
```



Agenda

- What are distributed arrays?
- ➔ What can you do with distributed arrays?
- Building a distributed array
- Advanced manoeuvres
- Debugging methods and tips

What can you do with distributed arrays?

Distributed array functions

The screenshot shows the MathWorks Help Center interface. The main article is titled "Run MATLAB Functions with Distributed Arrays". It includes a code snippet:

```
D = distributed/gallery('lehmer',n);
e = eig(D);
```

Below the code, there is a "Tip" box with the text: "For a filtered list of all MATLAB functions that support distributed arrays, see [Function List \(Distributed Arrays\)](#)." A red box highlights this link, and a red arrow points from it to the right-hand screenshot.

The right-hand screenshot shows the "Help Center" search results for "Distributed Arrays". The "Extended Capability" section is checked, showing 518 functions. The "Operators and Elementary Operations" section lists various arithmetic operations:

Operator	Description
+	Addition
sum	Sum of array elements
cumsum	Cumulative sum
-	Subtraction
diff	Differences and approximate derivatives
.*	Multiplication
*	Matrix multiplication
prod	Product of array elements
cumprod	Cumulative product
/	Right array division

- Extensive support:
 - Includes most linear algebra
 - Scale up mathematical operations

[Run MATLAB functions with distributed arrays](#)

What can you do with distributed arrays?

- >500 functions supported (as of R2020a)
- support for dense and sparse linear algebra
- support for numerics (double, single, logical, etc.)
- support for datetimes, durations, categoricals, tables, ...

- focus is on data preparation and large system solve

What can you do with distributed arrays?

Code written for distributed arrays looks like normal MATLAB code

```
% Create / read distributed data
A = distributed(...)

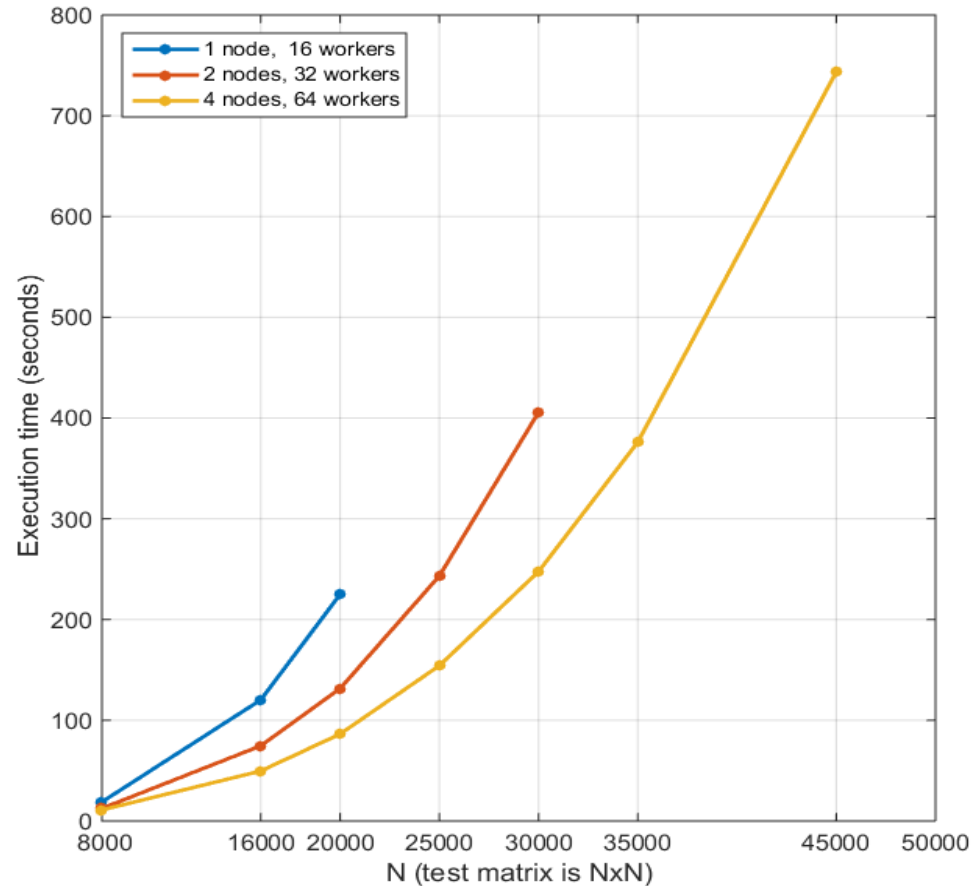
% Same code as for in-memory
b = sum(A, 2);
x1 = A\b; % direct solution
x2 = pcg(A, b); % iterative solution

% Bring back from cluster
[x1, x2] = gather(x1, x2);
```

Processing quite big data

Multiplication of 2 NxN matrices

$$\gg C = A * B$$



N	Execution time (seconds)		
	1 node, 16 workers	2 nodes, 32 workers	4 nodes, 64 workers
8000	19	13	11
16000	120	75	50
20000	225	132	86
25000	-	243	154
30000	-	406	248
35000	-	-	376
45000	-	-	743
50000	-	-	-

Processor: Intel Xeon E5-class v2
16 cores, 60 GB RAM per compute node, 10 Gb Ethernet

Agenda

- What are distributed arrays?
- What can you do with distributed arrays?
- ➔ Building a distributed array
 - Advanced manoeuvres
 - Debugging methods and tips

Building a distributed array

There are four main ways to build a distributed array:

1. Create from in-memory data
2. Build functions
3. Read from datastore
4. Construct from local parts

Building a distributed array: from in-memory

```
>> x = gallery("poisson", 10000);  
>> dx = distributed(x); % Data sent to workers
```

- All data is sent from client to workers
- Useful for debugging before scaling up
- Useful for data that is close to filling local memory (i.e. can be created but not operated on due to fill-in etc.)

Building a distributed array: build functions

```
>> dx = distributed.ones(1e9, 100);  
>> dx = distributed.rand(1e9, 100);  
>> dx = distributed.speye(1e9);  
>> dx = distributed.sprand(1e9, 1e9, 1e-8);  
... etc.
```

- No data is sent from client to workers
- Useful for creating test data and examples

Building a distributed array: build from datastore

Datastore:

- Simple interface for data in multiple files/folders
- Presents data a piece at a time
- Access pieces in serial (desktop) or in parallel (cluster)
- Back-ends for tabular text, images, databases and more
- Data always stacked vertically

```
>> ds = datastore("data/*.csv")
>> preview(ds)
  Rows      Cols      Vals
  _____  _____  _____
           1         1         4
           2         1        -1
    10001         1        -1
           :         :         :
           :         :         :
```

Building a distributed array: build from datastore

Read datastore into distributed

- each worker reads its own part of the data
- data is distributed vertically (workers have blocks of rows)

```
>> dt = distributed(ds); % distributed table of [Rows, Cols, Vals]
>> d = sparse(dt.Rows, dt.Cols, dt.Vals);

>> size(d)
ans = 100000000    100000000

>> nnz(d)
ans = 499960000
```

Agenda

- What are distributed arrays?
- What can you do with distributed arrays?
- Building a distributed array
- ➔ Advanced manoeuvres
- Debugging methods and tips

Writing your own algorithms

SPMD let's you craft parallel algorithms:

- Inside SPMD `distributed` -> `codistributed`
- `getLocalPart` extracts the data for this worker
- Use `gop` to reduce across workers (binary tree)
- Use `gcat` to concatenate results across workers

```
>> d = distributed(...)  
>> spmd  
    S = sum(getLocalPart(d), "all"); % Sum values on this worker  
    totalSum = gop(@plus, S); % Add results across all workers  
end
```

Writing your own algorithms

SPMD let's you craft parallel algorithms:

- **MPI-like** `labSend`, `labReceive`, `labSendReceive` for low level communication
- `labindex` for working out which worker you are

```
>> d = distributed(...)  
>> spmd  
    mydata = getLocalPart(d);  
    % Cycle data to the right (wrapping round on last worker)  
    prevWorker = mod(labindex-2, numlabs)+1; % worker to the left  
    nextWorker = mod(labindex, numlabs)+1; % worker to the right  
    mydata = labSendReceive(nextWorker, prevWorker, mydata)  
  
end
```


Modifying the data distribution

SPMD also gives control over distribution:

- `redistribute` for changing the distribution of an existing array
- `codistributorXX` for creating new (co)distributed arrays

```
>> d = distributed.ones(1e5); % Default is 1D distribution in dim 2
>> spmd
    % Switch to having whole rows per worker
d2 = redistribute(d, codistributor1d(1));
    % Use a custom (uneven) distribution of the data
partition = 1e4 * [1 2 3 4];
d3 = redistribute(d, codistributor1d(1, partition));
    % Switch to block-cyclic
d4 = redistribute(d, codistributor2dbc());

end
```

Agenda

- What are distributed arrays?
- What can you do with distributed arrays?
- Building a distributed array
- Advanced manoeuvres
- ➔ Debugging methods and tips

Debugging

Some rules of thumb:

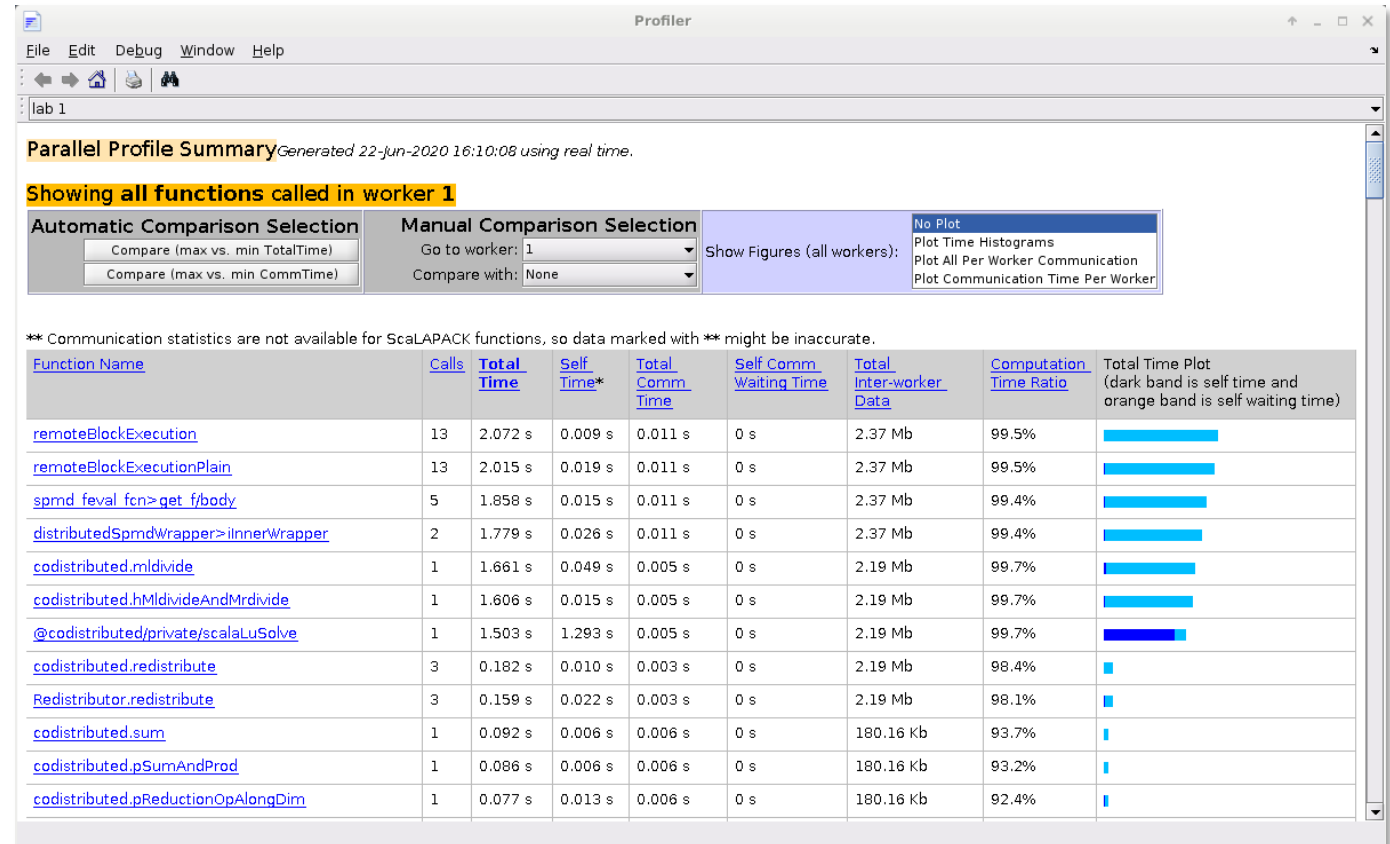
- Start with “local” pool and small data
 - faster iteration time
 - proves correctness
- Use pre-built algorithms if you can
- If writing your own, prefer higher-level gop, gcat over labSend/Receive
 - It’s easy to create mismatched communications with labSend/Receive!
 - Using unique tags for each communication helps to spot stray messages
 - Keeping computation and communication code separate helps in debugging both
- Run for a few different sizes to understand how duration scales with size
- Think about where your data lives – minimize transmission of files across networks

Debugging

Use parallel profiler (mpiprofiler):

- Shows what code ran, how much data transferred and allows comparison between workers (spot uneven loading)

```
>> mpiprofiler on
>> % Lots of parallel code ...
>> mpiprofview
```



Debugging

Use `spmd` or `parfevalonall` to interact with workers and see local state

```
>> spmd, d, end
```

```
Lab 1: This worker stores d(:,1:1667).
```

```
    LocalPart: [10000x1667 double]
```

```
    Codistributor: [1x1 codistributor1d]
```

```
Lab 2: This worker stores d(:,1668:3334).
```

```
    LocalPart: [10000x1667 double]
```

```
    Codistributor: [1x1 codistributor1d]
```

```
Lab 3:
```

```
    etc.
```

Questions

