

# A Model For the Time Complexity of the Time Stepper w/ MPI

Max Ruth

```
$Assumptions = { $\tau \in \text{Integers}$ ,  $\tau > 0$ };  
(* Assumptions about the variables we are working with. In particular,  
we want any Sums to evaluate correctly*)
```

## Description of the problem

There are a few parameters that we have to keep track of in this problem. They are:

```
Clear[Nx, Ny, L, M, T]  
Nx; Ny; (* Number of processors in the x and y directions *)  
L; (* Length of the x or y direction of the domain (it's square!) *)  
M; (* Number of grid points the x or y direction *)  
T; (* Total time for which the problem is run *)  
 $\tau$ ; (* Number of time steps taken per MPI Communication *)  
CFL; (* The CFL Condition *)  
h; u; v; (*The things we are solving for: height, x velocity, y velocity *)  
 $\alpha$ ;  $\beta$ ; (* Coefficients in the  $\alpha$ -  
     $\beta$  model of communication in MPI. Assume  $\beta$  has units that work with floats naturally *)  
 $\gamma$ ; (* The time required for one predict-correct time step at a single point *)  
 $\delta$ ; (* The time to calculate the speed at a single point *)
```

Derived from these parameters, we can find a few important things:

```
dx =  $\frac{L}{M}$ ; (* The distance between grid points *)  
dt =  $\frac{\text{CFL}}{\text{Max}[u, v]}$  dx; (* The time step (I think. It's something like this) *)  
mx =  $\frac{M}{Nx}$ ; my =  $\frac{M}{Ny}$ ;  
(* Number of grid points belonging to each cell (ignoring any integer nonsense) *)  
mg =  $4\tau$ ; (* Number of ghost grid cells on one side in one dimension *)  
Nt; (* Number of time steps taken. This isn't know a priori,  
but it is independent of the number of processors (assuming dx is fixed). However,  
it does depend on the CFL condition and dx, so it is in the "derived" category. *)
```

The algorithm that we are performing goes like this:

```

for i = 1 : Nt/τ
  Communicate ghost cells
  Find local time step
  Allreduce for global time step
  Perform the time step operation 2 τ times

```

To analyze the complexity of the algorithm, we need to think about the cost of each of these operations individually

## Finding Time Step Costs

Here, we calculate the cost of doing a block of time steps **per processor**.

### 1. Ghost Cell Communication Cost

The variable we are interested in is:

```
cg; (* Cost of ghost cell communication *)
```

To perform the ghost cell communication, we have to communicate with all four neighbors. The first two communications are to the left and right, and communicate blocks of size  $ng \cdot mx$ . The last two communications to the top and bottom have to include the corners as well, so they communicate blocks of size  $ng \cdot (my + ng)$ . Plugging this into the  $\alpha$ - $\beta$  model, we have

```
In[27]:= cg = 2 (α + β mg mx) + 2 (α + β mg (mg + my)) // Collect[#, {α, β}, FullSimplify] &
```

```
Out[27]:= 4 α + 8 β τ (M (1/Nx + 1/Ny) + 4 τ)
```

### 2. Find Local Time Step

The cost of this is simply

```
In[43]:= clts = δ mx my;
```

### 3. Allreduce for Global Time Step

This is the term that has the worst scaling with the number of processors. Assuming a reasonable algorithm, this should be done via some sort of divide-and-conquer, giving an approximate cost

```
In[45]:= car = α Log[Nx Ny];
```

### 4. Perform the time-step operation $2\tau$ times

Using the cost to predict-correct for a single node  $\gamma$ , we can find the cost to calculate a layer to be

```
In[37]:= clayer =  $\gamma$  (mx + 2 k) (my + 2 k);
```

where  $k$  is the number of ghost cells that are being calculated for the next step. That is, the first step after communicating, we have to calculate the value at  $mg - 2$  ghost cells in each direction. After that step, we calculate  $mg - 4$ , and so on until the final step where we calculate 0 ghost cells. In total, this becomes the sum

```
In[38]:= cts = Sum[clayer, {k, 0, mg - 2, 2}] // Collect[#,  $\tau$ , Simplify] &
```

$$\text{Out[38]= } \frac{2 \left( 3 M^2 + 8 N_x N_y - 6 M (N_x + N_y) \right) \gamma \tau}{3 N_x N_y} + 8 \left( -4 + M \left( \frac{1}{N_x} + \frac{1}{N_y} \right) \right) \gamma \tau^2 + \frac{128 \gamma \tau^3}{3}$$

As an example, if we communicate every time step ( $\tau = 1$ ), we have:

```
In[39]:= cts /.  $\tau \rightarrow 1$  // FullSimplify
```

$$\text{Out[39]= } \frac{2 \left( M^2 + 8 N_x N_y + 2 M (N_x + N_y) \right) \gamma}{N_x N_y}$$

## Finding Total Cost and Average Cost per Time Step

The total cost of the algorithm comes to multiplying the cost per time step block by the number of blocks, i.e.

```
In[46]:= Ctotal =  $\frac{Nt}{\tau}$  (cg + clts + car + cts) // Collect[#,  $\tau$ , FullSimplify] &
```

$$\text{Out[46]= } \frac{2 Nt \left( 6 M (N_x + N_y) \left( 2 \beta - \gamma \right) + 3 M^2 \gamma + 8 N_x N_y \gamma \right)}{3 N_x N_y} +$$

$$Nt \left( 32 \beta + 8 \left( -4 + M \left( \frac{1}{N_x} + \frac{1}{N_y} \right) \right) \gamma \right) \tau + \frac{128}{3} Nt \gamma \tau^2 + \frac{Nt \left( 4 \alpha + \frac{M^2 \delta}{N_x N_y} + \alpha \text{Log}[N_x N_y] \right)}{\tau}$$

The average cost per time step would then simply be

```
In[47]:= cavg =  $\frac{Ctotal}{Nt}$  // Collect[#,  $\tau$ , FullSimplify] &
```

$$\text{Out[47]= } \frac{2 \left( 6 M (N_x + N_y) \left( 2 \beta - \gamma \right) + 3 M^2 \gamma + 8 N_x N_y \gamma \right)}{3 N_x N_y} +$$

$$\left( 32 \beta + 8 \left( -4 + M \left( \frac{1}{N_x} + \frac{1}{N_y} \right) \right) \gamma \right) \tau + \frac{128 \gamma \tau^2}{3} + \frac{4 \alpha + \frac{M^2 \delta}{N_x N_y} + \alpha \text{Log}[N_x N_y]}{\tau}$$

## Analysis

## What is the best $\tau$ ?

It seems that we should probably ignore terms that are constant in  $\tau$ , as these cannot be tuned away via  $\tau$ . So,

```
In[48]:= cavganalysis =
  cavg - ((D[tau] (tau cavg // FullSimplify)) /. tau -> 0) // Collect[#, tau, FullSimplify] &
```

$$\text{Out[48]} = \left( 32\beta + 8 \left( -4 + M \left( \frac{1}{N_x} + \frac{1}{N_y} \right) \right) \gamma \right) \tau + \frac{128\gamma\tau^2}{3} + \frac{4\alpha + \frac{M^2\delta}{N_x N_y} + \alpha \log[N_x N_y]}{\tau}$$

Translating this big-O notation, we have

```
In[50]:= cavgtBigO = O[tau^2] + O[1/tau];
```

Because the coefficients of both of these terms are positive, there must be some optimal  $\tau$  that minimizes cavganalysis. The balance is essentially against latency type costs (i.e. terms with  $\alpha$  in them and the time to calculate the time step everywhere in the block) and the extra work required to calculate in the ghost cells (scaling for big  $\tau$  in the  $\tau^2$  term). The actual minimum will depend on all of these parameters however, so it is not so pretty to calculate.

## Weak Scaling

For the scaling problems, let's assume that  $N_x = N_y$ ,  $\tau$  is fixed, and that we only care about the cost of a time step (and not of the total cost). That is,

```
In[52]:= cscaling = cavg /. Ny -> Nx
```

$$\text{Out[52]} = \frac{2 \left( 12 M N_x (2\beta - \gamma) + 3 M^2 \gamma + 8 N_x^2 \gamma \right)}{3 N_x^2} + \left( 32\beta + 8 \left( -4 + \frac{2M}{N_x} \right) \gamma \right) \tau + \frac{128\gamma\tau^2}{3} + \frac{4\alpha + \frac{M^2\delta}{N_x^2} + \alpha \log[N_x^2]}{\tau}$$

For weak scaling, we go one step farther, and assume that  $M = M_0 N_x$ :

```
In[54]:= cweak = cscaling /. M -> M0 Nx // FullSimplify
```

$$\text{Out[54]} = \frac{4\alpha}{\tau} + 8M_0 \left( 2\beta - \gamma + 2\gamma\tau \right) + \frac{M_0^2 (\delta + 2\gamma\tau)}{\tau} + \frac{16}{3} (\gamma + 6\beta\tau - 6\gamma\tau + 8\gamma\tau^2) + \frac{\alpha \log[N_x^2]}{\tau}$$

We see that the cost per processor actually scales with the number of processors, but only logarithmically due to the Allreduce. The “flop rate” for all processors ( $\frac{\text{Total work of time step}}{\text{compute time per time step}} = K \frac{M^2}{c_{\text{weak}}}$ ) has the form

```
In[63]:= FRweakBigO = \frac{O[Nx^2]}{1 + O[\log[Nx^2]]};
```

This is increasing with  $N_x$  for all  $N_x$  still, so the optimal number of processors in weak scaling is still unbounded.

## Strong Scaling

In this case, we assume instead that  $M$  is independent of  $N_x$ . Collecting by  $N_x$ , the strong scaling cost for a single processor is

In[56]:= **cscaling** // **Collect**[#,  $N_x$ , **FullSimplify**] &

$$\text{Out[56]} = \frac{4\alpha}{\tau} + 32\beta\tau + \frac{M^2(\delta + 2\gamma\tau)}{N_x^2\tau} + \frac{16}{3}\gamma(1 - 6\tau + 8\tau^2) + \frac{8M(2\beta + \gamma(-1 + 2\tau))}{N_x} + \frac{\alpha \text{Log}[N_x^2]}{\tau}$$

In Big-O notation, again, we see that, so there is a decent chance

In[65]:= **cscalingBigO** =  $O\left[\frac{1}{N_x^2}\right] + O[\text{Log}[N_x^2]]$ ;

That is, for small  $N_x$ , the cost is dominated by time stepping, whereas for large (huge)  $N_x$ , the cost is dominated by Allreduce communication costs. Similarly, for the “flop rate” with strong scaling, we find that

In[67]:= **FRstrongBigO** =  $\frac{O[N_x^2]}{1 + O[\text{Log}[N_x^2] N_x^2]}$ ;

This expression will have a minimum - i.e. there is some amount of processors at which the code will actually get *slower*. However, even in the absence of the Allreduce step, this will still converge to a constant flop rate for large  $N_x$ , dominated by the constant term in cscaling above. This is because, regardless of the number of points belonging to a given processor, it still has to calculate the ghost cells, which has a finite lower bound of required work per processor per time step.