# A shared memory implementation of pipelined Parareal

# Daniel Ruprechta,b

<sup>a</sup>School of Mechanical Engineering, University of Leeds, Woodhouse Lane, Leeds LS2 9JT, UK
<sup>b</sup>Centre for Computational Medicine in Cardiology, Institute of Computational Science, Università della Svizzera italiana, Via Giuseppe Buffi 13,

CH-6900 Lugano, Switzerland

#### **Abstract**

The paper introduces an OpenMP implementation of pipelined Parareal and compares it to a standard MPI-based implementation. Both versions yield essentially identical runtimes, but, depending on the compiler, the OpenMP variant consumes about 7% less energy. However, its key advantage is a significantly smaller memory footprint. The higher implementation complexity, including manual control of locks, might make it difficult to use in legacy codes, though.

*Keywords:* Parareal, parallel-in-time integration, pipelining, shared memory, memory footprint, energy-to-solution 2010 MSC: 68W10, 65Y05, 68N19

#### 1. Introduction

Computational science faces a variety of challenges stemming from the massive increase in parallelism in state-of-the-art high-performance computing systems: projections suggest that exascale machines will require 100-million way concurrency [1]. This development mandates rethinking algorithms with respect to concurrency, fault-tolerance and energy efficiency but also the development of new and inherently parallel numerical methods. As a novel direction for the parallelisation of the solution of initial value problems, parallel-in-time (or time-parallel) integration schemes are attracting a quickly growing amount of attention [2]. One widely studied and used parallel-in-time method is *Parareal* [3], but other methods are frequently being introduced: RIDC [4], PFASST [5], MGRIT [6] or a DG-based time multi-grid method [7]. Parareal is based on ideas from multiple shooting for boundary value problems, which have been adopted for the time dimension [8, 9]. It has been applied problems from finance [10] over plasma physics [11] to fluid flow [12] and neutron kinetics [13]. The principle efficacy of parallelisation-in-time for extreme-scale parallel computations has been demonstrated [14] and recently the first software packages have also started to emerge.

Parareal's iterates between an expensive fine integrator run in parallel and a cheap coarse method which runs in serial and propagates corrections forward in time. While the unavoidable serial part limits parallel efficiency according to Amdahl's law, some improvements are possible by using a so-called *pipelined* implementation [15, 16]. Pipelining reduces the effective cost of the serial correction step in Parareal and therefore improves speedup. Even more optimisation is possible by using an event-based approach [17], but this requires a suitable execution framework that is not available on all machines. In contrast, pipelining comes naturally when implementing Parareal in MPI. It is, however, not straightforward in OpenMP and so far no shared memory version of Parareal with pipelining has been described. Studies therefore use almost exclusively MPI to implement Parareal and the very few using OpenMP implementations of Parareal use only the non-pipelined version [18, 19]. However, using shared memory can have advantages, since it avoids e.g. the need to allocate buffers for message passing. The disadvantage is that naturally OpenMP is limited to a shared memory unit, typically a single compute node. Since Parareal's convergence tends to deteriorate of too many parallel time slices are computed [20] and given the trend to compute nodes with large

Email address: d.ruprecht@leeds.ac.uk (Daniel Ruprecht)

numbers of cores, shared memory implementations might nevertheless be an attractive choice - as long as the benefit from pipelining does not have to be relinquished.

This paper introduces an OpenMP-based version of pipelined Parareal and compares it to a standard MPI-based implementation. For the comparison, a special-purpose Fortran code is used [21], which is freely available under a BSD license. The code is deliberately designed to not use external libraries other than MPI and OpenMP to avoid interference and facilitate rerunning the benchmarks on other architectures: since it only requires an MPI FORTRAN compiler built with thread support, it can be compiled and run on almost any platform.

#### 2. The method: Parareal

The starting point for Parareal is an initial value problem of the form

$$\dot{\mathbf{q}} = \mathbf{f}(\mathbf{q}(t), t), \ \mathbf{q}(0) = \mathbf{q}_0, \ t \in (0, T],$$
 (1)

which, in the numerical examples below, arises from the spatial discretisation of a PDE ("method-of-lines") with  $\mathbf{q} \in \mathbb{R}^{N_{\text{dof}}}$  being a vector containing all degrees-of-freedom. Let  $\mathcal{F}_{\delta t}$  denote a numerical procedure for the approximate solution of (1), for example a Runge-Kutta method. Denote further by

$$\mathbf{q} = \mathcal{F}_{\delta t}(\tilde{\mathbf{q}}, t_2, t_1) \tag{2}$$

the result of approximately integrating (1) forward in time from some starting value  $\tilde{\mathbf{q}}$  at a time  $t_1$  to a time  $t_2 > t_1$  using  $\mathcal{F}_{\delta t}$ . That is,

$$\mathbf{q} = \mathcal{F}_{\delta t}(\mathbf{q}_0, T, 0) \tag{3}$$

would indicate sequentially solving the full initial value problem in serial using the time stepper denoted by  $\mathcal{F}_{\delta t}$ . In order to parallelise the numerical integration of (1), Parareal and other so-called *parallel-across-the-steps* methods [22] introduce a decomposition of the time interval [0,T] into time-slices  $[t_p,t_{p+1}], p=0,\ldots,P-1$  where P is the number of cores, equal to the number of processes or threads, to be used in time. In the implementations sketched in Section 3, time slice  $[t_p,t_{p+1}]$  is assigned to the core running either thread number p (in the OpenMP variant) or the process with rank p (in the MPI version). For simplicity, assume here that all time slices have the same length and that the whole interval [0,T] is covered with a single set of time slices - otherwise, some form of restarting or *moving window* [23] has to be used.

Introduce now a second time integrator denoted  $\mathcal{G}_{\Delta t}$ , which has to be much cheaper to compute but can also be much less accurate (commonly referred to as the "coarse propagator"). Typically, a lower order method with a larger time step is used here, but a lower order spatial discretisation can be used as well. Using fewer degrees of freedom on the coarse level is also possible, but necessitates the introduction of suitable transfer operators (interpolation and restriction) and can degrade convergence. Parareal starts off with a *prediction step*, computing a rough guess of the starting value  $\mathbf{q}_p^0$  at the beginning of each time slice by

$$\mathbf{q}_p^0 = \mathcal{G}_{\Delta t}(\mathbf{q}_0, t_p, 0), \quad p = 0, \dots, P - 1.$$

$$\tag{4}$$

Here, subscript p indicates an approximation of the solution at time  $t_p$ . These approximate starting values are used to start the following iteration, run concurrently on each time slice,

$$\mathbf{q}_{p+1}^{k} = \mathcal{G}_{\Delta t}(\mathbf{q}_{p}^{k}, t_{p+1}, t_{p}) + \mathcal{F}_{\delta t}(\mathbf{q}_{p}^{k-1}, t_{p+1}, t_{p}) - \mathcal{G}_{\Delta t}(\mathbf{q}_{p}^{k-1}, t_{p+1}, t_{p}), \quad p = 0, \dots, P-1, \quad k = 1, \dots, K.$$
 (5)

As k increases, this iteration converges at the endpoints of the slices to the same solution provided by (3), that is  $\mathbf{q}_p^k \to \mathcal{F}_{\delta t}(\mathbf{q}_0, t_p, 0)$  for  $p = 0, \dots, P$ . Because the computationally expensive evaluation of the fine propagator (referred to as the *fine integrator step* in Section 3) can be parallelised across time slices, iteration (5) can run in less wall clock time than the direct time-serial integration (3) – provided the coarse method is cheap enough and the number of required iterations K is small. Computation of the fine values is followed by the *correction step*: the updated value

<sup>&</sup>lt;sup>1</sup>Here, the prescribed value  $\mathbf{q}_0$  at t = 0 is referred to as initial value, while the value  $\mathbf{q}_p$  the beginning of a time slice is called starting value.

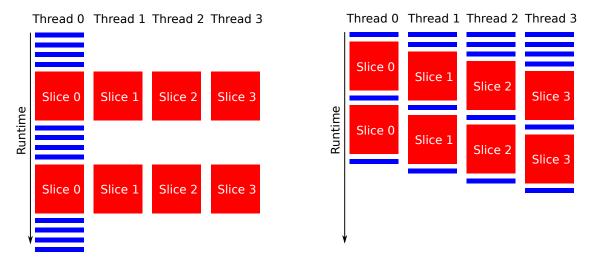


Figure 1: Execution diagram for Parareal without (left) and with (right) pipelining. Red blocks correspond to  $c_f$ , the time needed to run the fine integrator  $\mathcal{F}_{\delta t}$  over one time slice  $[t_p, t_{p+1}]$  while blue blocks correspond to  $c_c$ , the time required by the coarse method  $\mathcal{G}_{\Delta t}$ . Pipelining allows to hide some of the cost of the coarse method. While it comes naturally when using MPI to implement Parareal (there, Thread would refer to a Process), using simple loop-based parallelism with OpenMP results in the non-pipelined version shown on the left. A pipelined OpenMP Parareal is introduced in Section 3.

 $\mathbf{q}_p^k$  from the previous time slice is obtained, the coarse integrator is applied to it and (5) is evaluated to obtain  $\mathbf{q}_{p+1}^k$  which is then used to correct the value on the next time slice and so on. Note that, in contrast to the fine integrator, the correction step has to be performed in correct order, going step by step from the first time slice to the last. Also, when using a distributed memory parallelisation, the value  $\mathbf{q}_p^k$ , required in (5) to compute  $\mathcal{G}_{\Delta t}(\mathbf{q}_p^k, t_{p+1}, t_p)$ , has to be communicated from the process handling time slice  $[t_{p-1}, t_p]$  to the one handling  $[t_p, t_{p+1}]$ .

#### 2.1. Performance model

The expected performance of Parareal, referred to here as *projected speedup*, can be described by a simple theoretical model [16]. Here, the model is briefly repeated to more clearly illustrate the effect and importance of pipelining Parareal, see Subsection 2.2. It is also used as a baseline to compare against the measured speedups reported in Section 4.

Denote, as above, by P the number of cores in time, equal to the number of time slices. Further, denote by  $c_c$  the cost of integrating over one time slice using  $\mathcal{G}_{\Delta t}$  and by  $c_f$  the cost when using  $\mathcal{F}_{\delta t}$ . Because all time slices are assumed to consist of the same number of steps and an explicit method is used here, it can be assumed that  $c_f$  and  $c_c$  are identical for all time slices. This does not necessarily hold for an implicit method where differences in the number of iterations required by the nonlinear solver to converge can introduce load imbalances. A simple theoretical model for the speedup of Parareal using K iterations against running the fine method in serial now reads

$$s_{\rm np}(P) = \frac{Pc_{\rm f}}{(1+K)Pc_{\rm c} + Kc_{\rm f}} = \frac{1}{(1+K)\frac{c_{\rm c}}{c_{\rm f}} + \frac{K}{P}}.$$
 (6)

Equation (6) illustrates the necessity for a cheap coarse method to minimise the ratio  $c_c/c_f$  while still guaranteeing rapid convergence to minimise the terms K/P and 1 + K.

### 2.2. Pipelining Parareal

It has been pointed out that a proper implementation of Parareal allows to hide some of the cost of the coarse propagator and the name *pipelining* has been coined for this approach [16]. Figure 1 sketches the execution diagrams of both a non-pipelined (left) and pipelined (right) implementation for four time slices. As can be seen, pipelining

reduces the effective cost of the coarse correction step in each iteration from  $P \times c_c$  to  $c_c$  – but note that the initial prediction step still has cost  $P \times c_c$  as before. For pipelined Parareal, estimate (6) changes to

$$s_{\rm p}(P) = \frac{Pc_{\rm f}}{Pc_{\rm c} + Kc_{\rm c} + Kc_{\rm f}} = \frac{1}{\left(1 + \frac{K}{P}\right)\frac{c_{\rm c}}{c_{\rm f}} + \frac{K}{P}}.$$
 (7)

Because  $K/P \ll K$ , the pipelined version allows for better speedup, that is  $s_{np}(P) \leq s_p(P)$ . However, because pipelining only hides cost from the coarse integrator, the effect is smaller when the coarse method is very cheap and  $c_c/c_s \ll 1$ . In that case, the term K/P dominates the estimate which is not affected by pipelining.

### 3. PararealF90: Pipelined Parareal in MPI and OpenMP

The code used here for benchmarking is written in Fortran 90 and available under an open-source BSD license [21]. It is special-purpose and tailored to solve a single benchmark problem, 3D Burgers' equation

$$u_t + u \cdot \nabla u = \nu \Delta u \tag{8}$$

on  $[0,1]^3 \subset \mathbb{R}^3$  with periodic boundary conditions. While this a rather specific setup, finite difference stencils are a widely used motif in computational science: even though specifics and runtimes will be different if more complex problems are solved, the general concepts are used in complex application codes, e.g. by means of domain specific embedded languages (DSEL) [24] and the possibility to use Parareal with such a DSEL has been shown [25]. Thus, conclusions in terms of performance of different implementations of Parareal can be generalised to some extent, in particular because the Parareal routines only operate on linear arrays and do not see any specifics of the underlying time steppers or spatial discretisation. By exchanging the modules implementing the spatial discretisation, PararealF90 could be used to solve and benchmark different equations using stencil-based discretisations.

The time integration module provides two different methods, a strong stability preserving Runge-Kutta method (RK3-SSP) [26] and a first order forward Euler. The module for the spatial discretisation provides a fifth order WENO finite difference discretisation [26] and a simple first order upwind stencil for the advection term as well as a second and fourth order centred stencil for the diffusive term. Two modules provide the different implementations of Parareal, either using MPI as described in Subsection 3.1 or OpenMP as introduced in Subsection 3.2.

For  $\mathcal{G}_{\Delta t}$ , the forward Euler with upwind and second order centred stencils is used, while  $\mathcal{F}_{\delta t}$  uses the RK3-SSP integrator, a 5th order WENO for advection and a fourth order centred stencil for diffusion. Note that both implementations of Parareal use the same modules to provide the coarse and fine integrator and spatial discretisation. Since all three versions rely on the same implementation for the actual integrators and spatial discretisation, differences in performance should therefore solely emerge from the different Parareal routines wrapped around the compute routines. To run coarse and fine serial reference simulations, driver functions are used that directly call the same time stepper routines that are used within Parareal with the same configurations.

The code also comes with a test harness: in particular, the tests guarantee that all three implementations of Parareal produce results that are identical up to a tolerance of  $\varepsilon = 10^{-14}$  and thus essentially to round-off error. To detect possible race conditions, the comparison test can be performed multiple times: up to 100 instances of the test were run and passed on both used architectures described in Subsection 4.1. The tests also use various randomised parameters (randomised in a small range to avoid too large problems with very long runtimes): diffusion parameter  $\nu$ , number of time slices and processors P, number of finite difference nodes N (each direction has a different value in the tests) and number of fine and coarse steps per time slice. Furthermore, both implementations of Parareal use three auxiliary buffers per time-slice: q to store the fine value and for communication,  $\delta q$  to store the difference  $\mathcal{F}_{\delta l}(q) - \mathcal{G}_{\Delta l}(q)$  needed in the correction step and  $q_c$  to store the coarse value from the previous iteration. It seems that without introducing additional communication, three copies per time slice is the minimum storage required for Parareal.

Both used systems have nodes with two multi-core CPUs, see the description in Subsection 4.1. For the OpenMP version to be efficient, it has take care not only of thread safety but also take into account *non-uniform memory access* (NUMA): fetching data from memory associated with the other CPU requires communication through the intra-node interconnect and incurs overhead. Thus, in Parareal F90, solution buffers are extended by one dimension, e.g. Q(i,j,k,p), where i,j,k refer to the three spatial coordinates and p to the thread number. In the discussion

# Algorithm 1: Parareal using MPI

```
input: Initial value q_0; number of iterations K
 1.1 q \leftarrow q_0
 1.2 p = MPI\_COMM\_RANK()
 1.3 q \leftarrow \mathcal{G}_{\Lambda t}(q, t_p, 0)
 1.4 q_c \leftarrow \mathcal{G}_{\Delta t}(q, t_{p+1}, t_p)
 1.5 for k = 1, K do
            q \leftarrow \mathcal{F}_{\delta t}(q, t_{p+1}, t_p)
 1.6
            \delta q \leftarrow q - q_c
 1.7
            if Process not first then
 1.8
                MPI_RECV(q, source = p - 1)
 1.9
1.10
            end
            else
1.11
            q \leftarrow q_0
1.12
            end
1.13
1.14
            q_c \leftarrow \mathcal{G}_{\Delta t}(q, t_{p+1}, t_p)
            q \leftarrow q_c + \delta q
1.15
            if Process not last then
1.16
                 MPI\_SEND(q, target = p + 1)
1.17
1.18
            end
1.19 end
```

below, the first three dimensions are omitted and the buffer corresponding to thread p is indicated simply by  $\mathbb{Q}(p)$ . By using first touch initialisation, each thread ensures its respective part of the solution data is allocated locally. This strategy provides thread safety, as long as each thread is only accessing its own part of the buffers, as well as data affinity.

### 3.1. Parareal in MPI

The implementation of Parareal with MPI is straightforward and has been illustrated before [25]. However, it is repeated here for the sake of completeness and sketched in Algorithm 1.

- **Prediction phase:** lines 1.3 and 1.4. Every process generates its own coarse starting value  $\mathbf{q}_p^0 = \mathcal{G}_{\Delta t}(\mathbf{q}_0, t_p, 0)$  and computes  $\mathcal{G}_{\Delta t}(\mathbf{q}_p^0, t_{p+1}, t_p)$  for use in the correction. Later processes have to perform more time steps and thus take longer to complete the prediction phase, leading to the pipelined execution model sketched in Figure 1, since no global synchronisation points exist.
- **Fine integrator:** lines 1.6 and 1.7. Each process computes the fine value  $\mathcal{F}_{\delta t}(\mathbf{q}_p^k, t_{p+1}, t_p)$  in line 1.6 and, in order to free the buffer q for reuse in the receive, computes the difference  $\delta q := \mathcal{F}_{\delta t}(\mathbf{q}_p^k, t_{p+1}, t_p) \mathcal{G}_{\Delta t}(\mathbf{q}_p^k, t_{p+1}, t_p)$  between coarse and fine value in line 1.7 and stores it in  $\delta q$ .
- **Update phase:** lines 1.8–1.18. To receive the updated value from the previous time slice, every process but the first posts an MPI\_RECV. The first process simply copies the initial value  $\mathbf{q}_0$  into the buffer q instead. After the receive is completed, the coarse value  $\mathcal{G}_{\Delta t}(\mathbf{q}_p^{k+1},t_{p+1},t_p)$  of the new starting value is computed and the updated end value  $\mathbf{q}_{p+1}^{k+1} = \mathcal{G}_{\Delta t}(\mathbf{q}_p^{k+1},t_{p+1},t_p) + \delta q$  is computed and send to the process handling the following time slice using MPI\_SEND.

Note how this implementation naturally features pipelining as sketched in Figure 1 (right). In order to obtain the non-pipelined execution flow sketched in Figure 1 (left), MPI\_BARRIER directives would have to be added artificially after the prediction step, that is after line 1.4, and before the correction, that is before line 1.8. Parareal without pipelining in MPI would therefore be more complex and deliver reduced performance, so that this strategy seems rather senseless. Previous benchmarks did not show a significant difference between blocking and non-blocking communication [25] and for the sake of simplicity, blocking communication is used here.

Algorithm 2: Parareal with pipelining using OpenMP

```
input: Initial value q_0; number of iterations K
 2.1 q(0) \leftarrow q_0
 2.2 OMP PARALLEL
 2.3 N = OMP\_GET\_MAX\_THREADS()
 2.4 OMP DO
 2.5 for p = 0, N - 1 do
         q(p) \leftarrow q(0)
 2.6
         if Thread not first then
 2.7
             q(p) \leftarrow \mathcal{G}(q(p), t_p, 0)
 2.8
 2.9
2.10
         g_c(p) \leftarrow \mathcal{G}(q(p), t_{p+1}, t_p)
2.11 end
2.12 OMP END DO NOWAIT
2.13 for k = 1, K do
         OMP DO ORDERED
         for p = 0, N - 1 do
2.15
              OMP_SET_LOCK(p)
2.16
              q(p) \leftarrow \mathcal{F}_{\delta t}(q(p), t_{p+1}, t_p)
2.17
2.18
              \delta q(p) \leftarrow q(p) - q_c(p)
              OMP_UNSET_LOCK(p)
2.19
              OMP\_ORDERERD
2.20
              if Thread is first then
                   OMP_SET_LOCK(0)
2.22
                  q(0) \leftarrow q_0
2.23
                  OMP_UNSET_LOCK(0)
2.24
2.25
              q_c(p) \leftarrow \mathcal{G}_{\Delta t}(q(p), t_{p+1}, t_p)
2.26
              if Thread not last then
                   OMP_SET_LOCK(p+1)
2.28
                   q(p+1) \leftarrow q_c(p) + \delta q(p)
2.29
                   OMP_UNSET_LOCK(p+1)
2.30
              end
2.31
              OMP END ORDERED
2.32
2.33
         OMP END DO NOWAIT
2.34
2.35 end
2.36 OMP END PARALLEL
```

#### 3.2. Parareal in OpenMp

The implementation of Parareal with pipelining in OpenMP introduced here is sketched in Algorithm 2. To implement pipelined Parareal with OpenMP, essentially the whole algorithm is enclosed in one parallel region: threads are spawned by the OMP PARALLEL directive in line 2.2 and terminated by OMP END PARALLEL in line 2.36. Because this version requires manual synchronisation, a number of OpenMP locks is created using OMP\_INIT\_LOCK (not shown), one for each thread. During the fine integrator and update step, these locks are set and unset using OMP\_SET\_LOCK and OMP\_UNSET\_LOCK to protect buffers during writes and avoid race conditions.

• **Prediction step:** lines 2.4–2.12. Just as in the MPI example, each thread is computing its own coarse prediction of its starting value  $\mathbf{q}_p^0$  in a parallelised loop. The coarse value  $\mathcal{G}_{\Delta t}(\mathbf{q}_p^0, t_{p+1}, t_p)$  is also computed and stored for use in the first iteration. The later the time slice (indicated by a higher thread number p), the more steps the

thread must compute and thus the larger its workload. Therefore, at the end of the coarse prediction loop, the NOWAIT clause is required to avoid implicit synchronisation and enable pipelining. <sup>2</sup>

- **Parareal iteration:** lines 2.13–2.35. Here, both the fine integrator and update step are performed inside a single loop over all time slices, parallelised by OMP DO directives. Because parts of the loop (the update step) have to be executed in serialised order, the ORDERERD directive has to be used in line 2.14. Again, to avoid implicit synchronisation at the end of the loop, the NOWAIT clause is required in line 2.34.
  - Fine integrator: lines 2.16–2.19. Before the fine integrator is executed, an OMP\_LOCK is set to indicate that the thread will start writing into buffer q(p). Because thread p-1 accesses this buffer in its update step, locks are necessary to prevent race conditions and incorrect solutions. After the lock is set, the thread proceeds with the computation of  $\mathcal{F}_{\delta l}(\mathbf{q}_p^k, t_{p+1}, t_p)$  and computation of the difference between coarse and fine value  $\delta q$ . Then, since q(p) is now up to date and  $\delta q$  ready, the lock can be released.
  - **Update step:** lines 2.20-2.32. The update step has to be performed in proper order, from first to last time slice. Therefore, it is enclosed in ORDERED directives, indicating that this part of the loop is to be executed in serial order. Then, as in the two other versions, the update step is initialised with  $\mathbf{q}_0^{k+1} = \mathbf{q}_0$ . For every time slice, the coarse value of the updated initial guess is computed and the update performed. The updated end value is written into buffer q(p+1) to serve as the new starting value for the following time slice. However, to prevent thread p from writing into q(p+1) while thread p+1 is still running the fine integrator, thread p sets OMP\_LOCK number p+1 while performing the update.

This implements a pipelined Parareal in OpenMP and achieves runtimes that are essentially identical to the one provided by the MPI version. The necessity to manually control synchronisation between threads, however, makes it more complex. As shown in Section 4, it can outperform the MPI implementation in terms of memory requirements and energy consumption.

#### 4. Numerical results

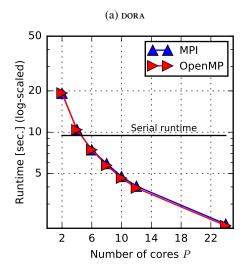
In this section, the OpenMP and MPI implementation are compared with respect to runtime in Subsection 4.2, memory footprint in Subsection 4.3 and energy consumption in Subsection 4.4. The parameters for the simulation are a viscosity parameter of v = 0.02 and a spatial discretisation on both levels with  $N_x = N_y = N_z = 40$  grid points in every direction. The simulation is run until T = 1.0 with a coarse time step of  $\Delta t = 1/192$  and a fine step of  $\Delta t = 1/240$ . Because of the quite high computational cost of the WENO-5 method in comparison to a cheap first order upwind scheme and the fact that RK3SSP needs three evaluations of the right hand side per step while the Euler method needs only one, the coarse propagator is about a factor of forty faster, despite the fact that the coarse step is only a factor of 1.25 larger than the fine.

To fix the number of iterations to a meaningful value which guarantees comparable accuracy from Parareal and serial fine integrator, we estimate the discretisation error of  $\mathcal{F}_{\delta t}$  by comparing against a reference solution with time step  $\delta t/10$ . This gives estimates for the fine relative error at T=1 of about  $e_{\text{fine}} \approx 5.9 \times 10^{-5}$  and for the coarse error of about  $e_{\text{coarse}} \approx 7.3 \times 10^{-2}$ . For P=24 time slices, after three iterations, the defect between Parareal and the fine solution is approximately  $1.4 \times 10^{-4}$ , after four iterations  $1.5 \times 10^{-5}$ . We therefore fix the number of iterations to K=4 so that for all values of P Parareal produces a solution with the same accuracy as the fine integrator.

#### 4.1. Hardware

Benchmarks are run on a single node of two different systems, respectively. The first one is cub, a commodity Linux cluster at the Institute of Computational Science in Lugano, consisting of 3x14 IBM Blade nodes. Each node has two quad-core Opteron (Barcelona) CPUs, for a total of 8 cores per node, and 16 GigaByte main memory. Nodes

<sup>&</sup>lt;sup>2</sup>Note that the prediction step could be implemented without a loop, by just using the thread number. However, within the Parareal iteration the ORDERED directive is required to guarantee a serialised execution of the update step and this directive is only available as part of a parallelised loop. Therefore, to keep the code consistent, the prediction step is also written as a parallel loop.



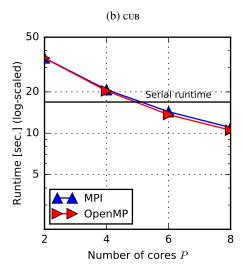


Figure 2: Five run averages of runtime with relative standard deviation below 0.01 on DORA (2a) and 0.025 on CUB (2b).

are connected through an Infiniband interconnect. The used MPI implementation is OpenMPI-1.4.2<sup>3</sup>, compiled with GCC-4.5.0<sup>4</sup> with flags -03 for maximum optimisation and -fopenmp to enable OpenMP. Note that, as the code is stand alone, no external libraries have to be linked.

The second system is PIZ DORA at the Swiss National Supercomputing Centre. DORA is a Cray XC40 with a total of 1,256 compute nodes. Each node contains two 12-core Intel Haswell CPUs and has 64 or 128 GigaByte of RAM and nodes are connected through a Cray Aries interconnect, using a dragonfly network topology. Two compilers are tested, the GCC-4.9.2 and the Cray Fortan compiler version 8.3.12. Both use the MPICH MPI library version 7.2.2. As on cub, compiler flags -03 and -fopenmp (GCC) or omp (Cray compiler) are used. Performance data for each completed job is generated using the Cray *Resource Utilisation Reporting* tool RUR [27]. RUR collects compute node statistics before and after each job and provides data on user and system time, maximum memory used, amount of I/O operations, consumed energy and other metrics. However, it only collects data for a full node and not for individual CPUs or cores.

## 4.2. Wall clock time and speedup

At first, runtime and speedup compared to the serial execution of the fine integrator are assessed. On both Cub and Dora, for each variant of Parareal and each value of P, five runs are performed and the average runtime is reported here. Measured runtimes are quite stable across different runs: the largest relative standard deviation of all performed five-run ensembles is smaller than 0.025 on cub and smaller than 0.01 on dora. Therefore, plots show only the average values without error bars, because those are hardly recognisable and clutter the figure.

Figure 2 shows runtimes in seconds depending on the number of cores on Dora (left) and cub (right). For Dora, only results from the Cray compiler are shown, which tends to generate slightly faster code than the GCC. The runtime of the serial fine integrator is indicated by a horizontal black line. Note that the y axis is scaled logarithmically, so the distance from 5 to 10 seconds is the same as from 10 to 20. Because the more modern CPUs on Dora are faster, runtimes are generally smaller on Dora than on cub. For P = 8, for example, Parareal runtimes on cub are around 10 seconds but only around 7 seconds on Dora. Both versions give almost identical performance: on both Dora and cub, OpenMP is marginally faster than the MPI version.

In addition, Figure 3 shows the speedup relative to the fine integrator run serially. Both versions fall short of the theoretically possible speedup indicated by the black line, but differences between MPI and OpenMP are small. As far

<sup>3</sup>http://www.open-mpi.org/

<sup>4</sup>https://gcc.gnu.org/

<sup>&</sup>lt;sup>5</sup>Detailed specification can be found here http://www.cscs.ch/computers/piz\_daint/index.html

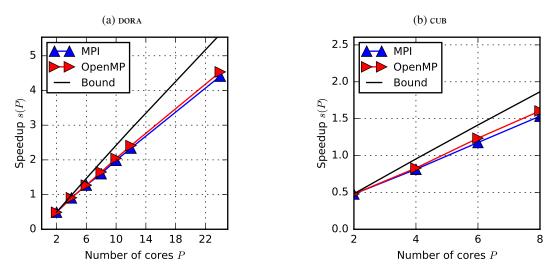


Figure 3: Speedup computed from average runtimes shown in Figure 2 on DORA (3a) and CUB (3b).

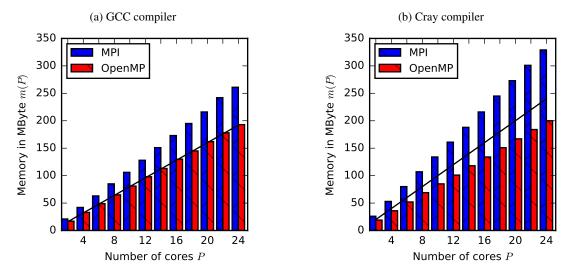


Figure 4: Maximum memory allocated in MegaByte for GCC (4a) and Cray compiler (4b) for the three different versions of Parareal depending on the number of used cores *P*. The black line indicates expected memory consumption computed as number of cores time memory footprint of serial fine integrator.

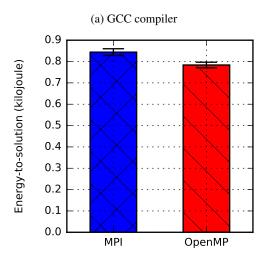
as runtimes and speedup are concerned, there is no indication that using the more complex OpenMP version provides benefits.

# 4.3. Memory footprint

The memory footprint of the code is measured only on dora where RUR is available. In contrast to runtime and energy, the memory footprint, as expected, does not vary between runs. Therefore Figure 4 shows a visualisation of the data from a single run with no averaging. The bars indicate the maximum required memory in MegaByte (MB) while the black line indicates the expected memory consumption using *P* cores computed as

$$m(P) = P \times m_{\text{serial}}$$
 (9)

where  $m_{\text{serial}}$  is the value measured for a reference run of the fine integrator. Because copies of the solution have to be stored for every time slice, the total memory required for Parareal can be expected to increase linearly with the



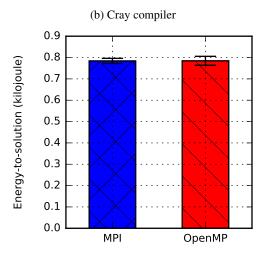


Figure 5: Energy-to-solution in Kilojoule for GCC in (5a) and Cray compiler (5b) for the three versions of Parareal, averaged across 50 samples, each using all 24 cores of the node. The largest relative standard deviation is 0.07, error bars indicate 95% confidence intervals.

number of cores in time. Note, however, that memory required per core stays constant if it follows (9).

For the OpenMP variant compiled with GCC, the memory footprint shown in Figure 4a exactly matches the expected values. The Cray compiler, shown in Figure 4b, leads to a smaller than expected memory footprint, but memory requirements still increase linearly with the number of time slices. For both compilers, the MPI version causes a noticeable overhead in terms of memory footprint, most likely because of internal allocation of additional buffers for sending and receiving [28]. For both OpenMP and MPI, the total memory footprint is *larger* for the Cray than for the GCC compiler, but the effect is much more pronounced for the MPI implementation (329 MB versus 261MB) than for the OpenMP variant (200MB versus 193MB).

It is important to note that both implementations allocate three auxiliary buffers per core. The overhead in terms of memory in MPI does thus not simply stem from allocating an additional buffer for communication, but comes from within the MPI library. The OpenMP implementation avoids this overhead. Given that memory will be a much more precious resource on future supercomputers and that the memory requirements of "across-the-steps" parallel-in-time methods have raised concerns [29], these savings might be important. Since the memory overhead grows as more and more cores are used in time, savings from OpenMP will be especially pronounced when Parareal is used with many time slices on nodes with a large numbers of cores.

#### 4.4. Energy-to-solution

The RUR tool reports the energy-to-solution for every completed job. Because RUR can only measure energy usage for a full node, results are reported here for runs using the full number of P=24 cores available on a dora node. In contrast to runtimes and memory footprint, energy measurements show significant variations between runs due to random fluctuations: thus, the presented values are averages over ensembles of 50 runs for each version of Parareal. This number of runs has been sufficient to reduce the relative standard deviation to below 0.09 in both configurations and therefore gives a robust indication of actual energy requirements.

When using the GCC compiler, the MPI version consumes more energy than OpenMP. The resulting 95% confidence intervals (assuming energy-to-solution is normally distributed), are  $844.04 \pm 15.89$ J for MPI and  $783.72 \pm 11.53$ J for OpenMP. Figure 5 gives a graphical representation of these values including the confidence intervals. The average for OpenMP is well outside the confidence interval for the MPI version, so this is very unlikely just a chance result. Moreover, because runtimes are almost identical, the differences in energy-to-solution cannot simply be attributed to differences in time-to-solution. It is also noteworthy that for the GCC compiler OpenMP without pipelining (results not shown here), even though it is *slower* than the MPI version, still requires less energy ( $844.04 \pm 15.89$ J versus  $801.14 \pm 13.18$ J). Tracking down the precise reason for the differences in energy-to-solution and power requirement

will require detailed tracing of power uptake, which is only possible on specially prepared machines [30] and thus left for future work.

Note that the energy consumption of Parareal has previously been studied [25]. By comparing against a simple theoretical model, it has been shown that the energy overhead of Parareal (defined as energy-to-solution of Parareal divided by energy-to-solution of the fine serial integrator), is mostly due to Parareal's intrinsic suboptimal parallel efficiency. While improving parallel efficiency of parallel-in-time integration remains the main avenue for improving energy efficiency, the results here suggest that in some cases a shared memory approach can provide non-trivial additional savings.

For code generated with the Cray compiler, both OpenMP and MPI lead to almost identical energy requirements: here, confidence intervals are  $784.24 \pm 11.93$ J for MPI and  $784.24 \pm 12.18$ J for OpenMP. It seems likely that the compiler optimises the message passing to take advantage of the shared memory on the single node. Supposedly, the MPI version handles communication in a way that is similar to what is explicitly coded in the OpenMP version. However, as shown in Subsection 4.3, this automatic optimisation comes at the expense of a significantly larger memory footprint.

#### 5. Summary

The paper introduces and analyses an OpenMP implementation of the parallel-in-time method Parareal with pipelining. Pipelining allows to hide some of the cost of the serial coarse correction step in Parareal and is important to optimise its efficiency (even though it cannot relax the inherent limit on parallel efficiency given by the inverse of the number of required iterations). Pipelining comes naturally in a distributed memory MPI implementation, but is not straightforward when using OpenMP. The new OpenMP implementation is compared to a standard MPI variant in terms of runtime, memory footprint and energy consumption for both a Cray compiler and the GCC. Both versions produce essentially identical runtimes. For both compilers, using OpenMP leads to significant reductions in memory footprint, but the effect is more pronounced for the Cray compiler. In terms of energy-to-solution, the results strongly depend on the compiler: while for GCC the OpenMP version is more energy efficient than the MPI version, there is no difference for the Cray compiler.

The results show that contemplating a shared memory strategy to implement "parallel-across-the-steps" methods like Parareal can be worthwhile. Even though it is more complicated, it can reduce memory requirements. A potential caveat is whether the benefits carry over to the full space-time parallel case, where a parallel-in-time method is combined with spatial decomposition. For Parareal without pipelining the potential of such a hybrid space-time parallel approach has been illustrated [19] but whether applies to the pipelined version introduced here remains to be seen.

# Acknowledgments

I would like to thank the Centre for Interdisciplinary Research (ZIF) at the University of Bielefeld, Germany, for inviting me for a research visit in August 2015. The writing of this article greatly benefited from the tranquil and productive atmosphere at ZIF. I gratefully acknowledge Andrea Arteaga's support with the energy measurements.

#### References

- [1] J. Dongarra, P. Beckman, al., The international exascale software roadmap, Int. Journal of High Performance Computer Applications 25 (1). doi:10.1177/1094342010391989.
- [2] M. J. Gander, 50 years of Time Parallel Time Integration, in: Multiple Shooting and Time Domain Decomposition, Springer, 2015. URL http://www.unige.ch/%7Egander/Preprints/50YearsTimeParallel.pdf
- [3] J.-L. Lions, Y. Maday, G. Turinici, A "parareal" in time discretization of PDE's, Comptes Rendus de l'Acadmie des Sciences Series I Mathematics 332 (2001) 661–668. doi:10.1016/S0764-4442(00)01793-6.
- [4] A. J. Christlieb, C. B. Macdonald, B. W. Ong, Parallel high-order integrators, SIAM Journal on Scientific Computing 32 (2) (2010) 818–835. doi:10.1137/09075740X.
- [5] M. Emmett, M. L. Minion, Toward an Efficient Parallel in Time Method for Partial Differential Equations, Communications in Applied Mathematics and Computational Science 7 (2012) 105–132. doi:10.2140/camcos.2012.7.105.
- [6] R. D. Falgout, S. Friedhoff, T. V. Kolev, S. P. MacLachlan, J. B. Schroder, Parallel time integration with multigrid, SIAM Journal on Scientific Computing 36 (2014) C635–C661. doi:10.1137/130944230.

- [7] M. J. Gander, M. Neumueller, Analysis of a Time Multigrid Algorithm for DG-Discretizations in Time (2014). URL http://arxiv.org/abs/1409.5254
- [8] P. Chartier, B. Philippe, A parallel shooting technique for solving dissipative ODE's, Computing 51 (3-4) (1993) 209–236. doi:10.1007/BF02238534.
- [9] M. Kiehl, Parallel multiple shooting for the solution of initial value problems, Parallel Computing 20 (3) (1994) 275–295. doi:10.1016/ S0167-8191(06)80013-X.
- [10] G. Bal, Y. Maday, A "Parareal" time discretization for non-linear PDE's with application to the pricing of an American Put, in: L. Pavarino, A. Toselli (Eds.), Recent Developments in Domain Decomposition Methods, Vol. 23 of Lecture Notes in Computational Science and Engineering, Springer Berlin, 2002, pp. 189–202. doi:10.1007/978-3-642-56118-4\_12.
- [11] D. Samaddar, D. E. Newman, R. S. Snchez, Parallelization in time of numerical simulations of fully-developed plasma turbulence using the parareal algorithm, Journal of Computational Physics 229 (2010) 6558–6573. doi:10.1016/j.jcp.2010.05.012.
- [12] E. Celledoni, T. Kvamsdal, Parallelization in time for thermo-viscoplastic problems in extrusion of aluminium, International Journal for Numerical Methods in Engineering 79 (5) (2009) 576–598. doi:10.1002/nme.2585.
- [13] A.-M. Baudron, J.-J. Lautard, Y. Maday, O. Mula, The parareal in time algorithm applied to the kinetic neutron diffusion equation, in: Domain Decomposition Methods in Science and Engineering XXI, Lecture Notes in Computational Science and Engineering, Springer International Publishing, 2014, pp. 437–445. doi:10.1007/978-3-319-05789-7\_41.
- [14] R. Speck, D. Ruprecht, R. Krause, M. Emmett, M. L. Minion, M. Winkel, P. Gibbon, A massively space-time parallel N-body solver, in: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC '12, IEEE Computer Society Press, Los Alamitos, CA, USA, 2012, pp. 92:1–92:11. doi:10.1109/SC.2012.6.
- [15] E. Aubanel, Scheduling of Tasks in the Parareal Algorithm, Parallel Computing 37 (2011) 172-182. doi:10.1016/j.parco.2010.10.004.
- [16] M. L. Minion, A Hybrid Parareal Spectral Deferred Corrections Method, Communications in Applied Mathematics and Computational Science 5 (2) (2010) 265–301. doi:10.2140/camcos.2010.5.265.
- [17] L. A. Berry, W. R. Elwasif, J. M. Reynolds-Barredo, D. Samaddar, R. S. Snchez, D. E. Newman, Event-based parareal: A data-flow based implementation of parareal, Journal of Computational Physics 231 (17) (2012) 5945–5954. doi:10.1016/j.jcp.2012.05.016.
- [18] D. Ruprecht, R. Krause, Explicit parallel-in-time integration of a linear acoustic-advection system, Computers & Fluids 59 (0) (2012) 72–83. doi:10.1016/j.compfluid.2012.02.015.
- [19] R. Krause, D. Ruprecht, Hybrid Space-Time Parallel Solution of Burgers' Equation, in: Domain Decomposition Methods in Science and Engineering XXI, Vol. 98 of Lecture Notes in Computational Science and Engineering, Springer International Publishing, 2014, pp. 647–655. doi:10.1007/978-3-319-05789-7\_62.
- [20] M. J. Gander, S. Vandewalle, Analysis of the Parareal Time-Parallel Time-Integration Method, SIAM Journal on Scientific Computing 29 (2) (2007) 556–578. doi:10.1137/05064607X.
- [21] D. Ruprecht, PararealF90: Implementing Parareal OpenMP or MPI?, release v1.0 (2015). doi:10.5281/zenodo.31288.
- [22] K. Burrage, Parallel methods for ODEs, Advances in Computational Mathematics 7 (1997) 1-3. doi:10.1023/A:1018997130884.
- [23] M. Schreiber, A. Peddle, T. Haut, B. Wingate, A decentralized parallelization-in-time approach with parareal (2015). URL http://arxiv.org/abs/1506.05157
- [24] T. Gysi, C. Osuna, O. Fuhrer, M. Bianco, T. C. Schulthess, Stella: A domain-specific tool for structured grid methods in weather and climate models, in: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, 2015. doi:10.1145/2807591.2807627.
- [25] A. Arteaga, D. Ruprecht, R. Krause, A stencil-based implementation of Parareal in the C++ domain specific embedded language STELLA, Applied Mathematics and Computation 267 (2015) 727–741. doi:10.1016/j.amc.2014.12.055.
- [26] C.-W. Shu, S. Osher, Efficient implementation of essentially non-oscillatory shock-capturing schemes II, Journal of Computational Physics 83 (1989) 32–78. doi:10.1016/0021-9991(89)90222-2.
- [27] A. Barry, Resource utilization reporting: Gathering and evaluating HPC system usage, in: CUG2013 Proceedings, 2013. URL https://cug.org/proceedings/cug2013\_proceedings/includes/files/pap103-file2.pdf
- [28] R. Rabenseifner, G. Hager, G. Jost, Hybrid MPI/OpenMP parallel programming on clusters of multi-core SMP nodes, in: 17th Euromicro International Conference on Parallel, Distributed and Network-based processing, 2009, pp. 427–436. doi:10.1109/PDP.2009.43.
- [29] D. E. Keyes, Exaflop/s: The why and the how, Comptes Rendus Mécanique 339 (23) (2011) 70 77. doi:10.1016/j.crme.2010.11.002.
- [30] C. Isci, M. Martonosi, Runtime power monitoring in high-end processors: Methodology and empirical data, in: Proceedings of the 36th Annual IEEE/ACM International Symposium on Microarchitecture, MICRO 36, IEEE Computer Society, Washington, DC, USA, 2003, pp. 93-.
  - URL http://dl.acm.org/citation.cfm?id=956417.956567