

Haskell vs. F# vs. Scala

A High-level Language Features and Parallelism Support Comparison

Prabhat Tooto Pantazis Deligiannis Hans-Wolfgang Loidl

School of Mathematical and Computer Sciences,
Heriot-Watt University, Edinburgh, EH14 4AS, U.K.
{pt114, pd85, H.W.Loidl}@hw.ac.uk

Abstract

This paper provides a performance and programmability comparison of high-level parallel programming support in Haskell, F# and Scala. Developing several parallel versions, we employ skeleton-based, semi-explicit and explicit approaches to parallelism. We focus on advanced language features for separating computational and coordination aspects of the code and tuning performance. We also assess the impact of functional purity and multi-paradigm design of the languages on program development and performance.

Basis for these comparisons are several Barnes-Hut implementations of the n-body problem in all three languages, on both Linux and Windows. Our performance measurements on state-of-the-art multi-cores achieve a speedup up to 5.62 (on 8 cores) with a highly-tuned Haskell version. For comparable implementations in Scala and F# we achieve speedups of 4.51 (on 8 cores) and 2.28 (on 4 cores), respectively. We observe that near best speedups are achieved using the highest level abstraction in these languages.

Categories and Subject Descriptors D.1.3 [Programming Techniques]: Concurrent Programming—Parallel programming; D.3.2 [Programming Languages]: Language Classifications—Applicative (functional) languages

General Terms Algorithms, Languages, Measurement, Performance

Keywords Haskell, F#, Scala, Parallelism, n-body, Barnes-Hut

1. Introduction

Because functional languages are not defined in terms of operations on a hidden global state, they avoid unnecessary sequentialisation and provide ample latent parallelism that can be exploited by compiler or runtime-system. This property makes them an attractive platform for exploiting common-place parallel hardware without imposing new concepts of explicit threads with explicit communication onto every parallel application.

Functional languages provide a high degree of abstractions and expressiveness, enabling the parallel programmer to only specify *what* value the program should compute instead of *how* to compute it. Managing parallelism is all about *how* and therefore largely

hidden from the programmer. However, for tuning parallel performance, some limited control of operational aspects is desirable.

The approaches to efficiently exploit such latent parallelism, provided in the latest implementations of state-of-the-art languages such as Haskell, F# and Scala, all aim to be minimally intrusive to the code, while giving the expert parallel programmer sufficient control to perform parallel performance tuning. In this paper, we evaluate language mechanisms provided in each of these three languages, covering several levels of abstraction as well as a range from purely declarative to mixed paradigm languages. We present a head-to-head comparison of the resulting parallel performance.

A new generation of programming languages, such as F# and Scala, often take a multi-paradigm approach, embedding the advantages of functional languages into a mainstream, object-oriented language. They use existing, highly-optimised VM technology, .NET and JVM respectively, to combine the ease of expressing parallelism with efficient sequential execution. In this paper, we perform a head-to-head programmability comparison between purely functional Haskell and multi-paradigm F# and Scala, all with semi-explicit control of parallelism that rule out the explicit use of threads. We assess the impact of key language design issues, in particular laziness and mutable data-structures, on sequential and parallel performance. Finally, we give a head-to-head parallel performance comparison of Haskell, F# and Scala, using different techniques to expose parallelism at different levels of abstraction. The results from our measurements of a Barnes-Hut implementation of the n-body algorithm show that we achieve respectable speedups in all languages. We achieve a speedup of 5.62 (on 8 cores) with a highly-tuned Haskell version. With the implementations in Scala and F# we achieve speedups of 4.51 (on 8 cores) and 2.28 (on 4 cores), respectively.

The remainder of the paper is organised as follows: Section 2 presents related work; Section 3 presents the background on Haskell, F# and Scala, and their support for parallelism; Section 4 presents our multiple Barnes-Hut algorithm implementations, both sequential and parallel; Section 5 presents results from our measurements on two different multi-core architectures, an 8-core Linux machine and a 4-core (with hyper-threading) Windows machine; and we summarise our findings in Section 6.

2. Related Work

We can broadly classify parallel declarative languages as implicitly parallel, without any explicit control of parallelism, semi-explicit, only exposing potential parallelism, and explicit, with constructs for the generation and handling of explicit threads. In this section we focus on and survey semi-explicit approaches, although the lower level constructs in Scala can also be classified as explicit. Another useful classification is by purity of the programming model,

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

FHPC'12, September 15, 2012, Copenhagen, Denmark.
Copyright © 2012 ACM 978-1-4503-1577-7/12/09...\$10.00

with Haskell and Clean [24] representing pure languages and F#, Scala, as well as most Lisp and ML dialects, such as OCaml [15], representing impure, mixed paradigm languages. The general area of parallel declarative programming languages is surveyed in depth in [35]. A more focused comparison of parallel Haskell variants is given in [34].

Crucial to a high-performance implementation of a declarative language is an adaptive runtime-system, that can make good decisions about the management of parallelism, usually deferred to the programmer in explicit languages. Important concepts are futures as handles for a data-structure, that might be evaluated in parallel and on which other threads should synchronise, first introduced in the Mul-T [12] variant of Lisp. Importantly for performance, this system introduced lazy task creation [21] as a technique, where one task can automatically subsume the computation of another task, thus increasing the granularity of the parallelism. Both, the language- and the system-level contributions have been picked up in recent implementations of parallel functional languages: F#, PolyML [20] and parallel CML [27], implemented in Manticore, provide language-level futures; *GpH*'s runtime system uses lazy task creation, by representing potential parallelism as "sparks" that can move freely and cheaply between processors and work represented by one spark can be subsumed by a running thread if no additional parallelism is required. Another key runtime-system design goal is to support light-weight threads, thus reducing the overhead for creating parallelism and encouraging a programming style that generates a massive amount of parallelism, giving the runtime-system the flexibility to arrange the parallelism in a way most suitable to the underlying hardware. Haskell/GHC excels at light-weight threads, as shown by the thread-ring benchmark of the Computer Language Shootout [29]. Filaments [17] and Cilk [10], now integrated in Intel's Cilk Plus compiler, are other examples of runtime-systems for light-weight threads.

Several experimental languages explored the use of high-level, parallelism language features in object-oriented languages: Fortress [30], X10 [6] and Chapel [5]. Of these, Chapel is currently best supported, in particular on massively parallel supercomputers. These languages introduce high-level constructs such as virtual shared memory (X10), structured programming constructs for parallel execution (Chapel), and software transactional memory (Fortress) to avoid a re-design of the software architecture due to specifics of the underlying, parallel architecture.

An increasingly important area of high-level abstractions for parallelism are parallel patterns or algorithmic skeletons [7], higher-order functions with pre-defined parallel computation structures. Because they can hide all complexities of the efficient, possibly hardware-dependent, handling of parallelism in a library, it is being picked up as technology of choice in mainstream languages without built-in high-level parallelism support. Prominent examples are Google's MapReduce [9] implementation, on large-scale, distributed architectures, Intel's Task Building Block [26] library, and to some extent the Task Parallel Library [14].

3. Background

3.1 Haskell

The key advantage of Haskell for parallel computation is referential transparency which guarantees that evaluation can happen in any order. This implies that the amount of inherent parallelism in a Haskell program is large such that each sub-expression can be evaluated in parallel. However, this leads to far too fine-grained parallelism and an approach that allows the programmer to specify which computation is worthwhile to be evaluated in parallel is desirable.

The lazy semantics of Haskell has implications on the parallel programming models that can be supported by the language. Unconstrained lazy evaluation is essentially sequential which contradicts how parallel evaluation should proceed. Some degree of eager evaluation is essential in order to arrange computations in parallel. The programmer also needs to specify the evaluation degree of expressions, such that just enough of an expression is evaluated in order to enable other expressions to continue evaluation in parallel.

A number of parallel programming models exist in Haskell and the choice between them depends on the problem at hand and the amount of control needed. They are covered in detail in [32]. For instance, building on top of concurrency primitives and developed entirely as a library, the *Par* monad offers a rather explicit approach requiring manual fork of tasks and communication between threads. *DPH* on the other hand is more implicit and is mainly used for data-parallel problems. In this paper, we focus on *GpH* which is an extension of GHC, a highly optimising, transformation-based compiler and graph-reduction-based runtime system.

3.1.1 GpH

Glasgow parallel Haskell [18, 19, 33] is a minimal, conservative, parallel extension of Haskell, supported by the GHC compiler. It extends standard Haskell by providing two basic primitives for specifying and controlling parallelism.

```
— parallel composition
par :: a -> b -> b
— sequential composition
pseq :: a -> b -> b
```

par allows the programmer to annotate computations that can be usefully be evaluated in parallel. The first argument is sparked and may potentially be executed in parallel with the evaluation of the second argument. *pseq* enforces sequential ordering which is needed to arrange parallel computations.

Naive usage of these primitives can lead to unexpected parallel behaviour, for example generating sparks for already evaluated data. *Evaluation strategies* are abstractions over the primitives to provide an even higher level of control of parallelism. Evaluation strategies provide a clean separation of the coordination aspects from the main computation. For example, *parList* can be used to demand parallel evaluation of each list element, separately from defining the contents of the list. Additionally, evaluation degree and evaluation order can be specified using evaluation strategies.

The following example uses the *parList* strategy to define a *parMap* skeleton:

```
— definition parallel map using strategies
parMap str f xs = map f xs 'using' parList str
— usage
parMap rdeepseq f xs
```

3.2 F#

F# [31] combines the features of a strict, higher-order, impurely functional language of an ML-style, with features of mainstream object-oriented languages. Both paradigms are made available to the programmer, who can make a choice based on the suitability of the paradigm for the application and on his familiarity with the paradigm. The F# implementation compiles to .NET as the VM and can therefore build on highly-optimised VM implementations, and interact with libraries in other languages also targeting this intermediate platform.

F# has good support for concurrency through asynchronous workflows and message passing, and parallelism through the Task Parallel Library. It allows combining these libraries to take advantage of potential parallelism in an application.

3.2.1 Concurrency Features

Asynchronous workflows are intended for writing concurrent and reactive programs that perform asynchronous I/O where avoiding threads to block is necessary. It can however be used for basic parallelisation. It works by wrapping a section of code inside an `async` block which can be run independently without blocking the main thread. For example, a web server can handle requests simultaneously. On a parallel machine, these requests are executed in parallel thus improving the performance.

```
let handleRq rq = async { (* some code here *) }
Async.RunSynchronously (handleRq request)
```

Building on the asynchronous infrastructure, the `MailboxProcessor` type encapsulates an agent-based concurrency implementation similar to that in Erlang [2]. In this model, agents run in parallel and communicate by sending messages to each other. The absence of shared mutable state makes applications scale and simplifies the understanding. Though requiring a quite different application structure, common patterns can be used to implement agent-based applications.

3.2.2 Task Parallel Library

The Task Parallel Library [4, 14] provides an API to simplify the process of adding parallelism to the application. TPL handles many of low-level details such as partitioning of work, scheduling of threads on the threadpool and scaling the degree of concurrency dynamically to exploit all available processors in the most efficient way.

From the programmer's viewpoint, there is no explicit notion of threads. The main construct for task parallelism is built around the concept of a *task* which represents an independent unit of work that is executed in parallel. Tasks are queued to the threadpool which employs work-stealing to balance the load. The number of threads are adjusted during runtime based on the workload and available processors. As several tasks can be mapped to a single thread, this make them relatively lightweight. Tasks can be controlled via a number of built-in methods. The code sample below creates a task to process each item in the sequence concurrently and synchronise at the end.

```
let tasks ts =
  [| for t in ts
    Task.Factory.StartNew(fun () ->
      process t) |]
Task.WaitAll tasks
```

Data parallelism is supported through the `Parallel` class and `PLINQ`. The `Parallel` class implements static methods such as `Parallel.For` and `Parallel.ForEach` for basic loops parallelisation. These methods provide an easy way of parallelising `for` loops. `PLINQ` offers a more declarative model for data parallelism based on Language Integrated Query. It provides a shallow embedding of an SQL-like query language directly in the general purpose, host language (in this case F#) and allows to query XML data, databases or objects from standard data collections. The latter is how we use `PLINQ`. The implementation uses TPL internally for efficient implementation, and is the highest level language mechanism for parallelism in F#.

3.3 Scala

Scala is a statically typed, strict, and multi-paradigm programming language, combining functional and object-oriented features [22]. The language allows the expression of common programming patterns in a concise, elegant and type-safe manner. Scala's compiler targets the Java Virtual Machine (JVM) platform and thus is fully interoperable with Java, empowering the programmer with the full

range of Java libraries and frameworks. The language was also designed with extensibility in mind, meaning that new features can be easily added in the form of new libraries without the need to change the syntax of the language.

A main focus of Scala is to deliver state-of-the-art high-level constructs and abstractions for concurrent programming, emphasising large-scale distribution, scalability and fault-tolerance. Towards this goal it provides a number of programming frameworks, most notably `Scala Actors` [11] for concurrency and `Scala Parallel Collections` [25] for implicit parallelism.

3.3.1 Scala Actors Library

Scala supports concurrency by providing an explicit message passing programming model based on actors. Actors are first-class, light-weight processes that communicate with each other by exchanging asynchronous messages [1]. These messages are gathered in the receiving actor's mailbox. An actor is able to iterate through its mailbox and respond to the various messages it has gathered by using pattern matching, a staple approach in functional programming. Responses to messages include among other actions: creating a new actor; sending a new message to the sender; and changing the underlying behaviour of the receiving actor. This design is motivated by Erlang [2].

```
a ! msg // actor sending a message

// receiving msgs and responding with actions
receive {
  case msg-pattern_1 => action_1
  case msg-pattern_2 => action_2
  case msg-pattern_n => action_n
}
```

3.3.2 Scala Parallel Collections Framework

The `Scala Parallel Collections Framework` provides the `parallel` collections subpackage, which defines parallel implementations for sequences, maps and sets, together with common parallel operations [25]. The programmer can use the method `par` on a sequential collection to invoke its corresponding parallel implementation. With the method `seq` the parallel collection behaves again in a sequential manner. The benefit of this approach is that parallel operations can have the same names as their sequential versions, which means that the programmer can easily introduce parallelism by just providing the method `par` in the right places, as shown below.

```
xs.map( (x: Int) => x + 1 ) // sequential

xs.par.map( (x: Int) => x + 1 ).seq // parallel
```

The implementation of the parallel collections library is built on top of the `Java Fork/Join` framework [13]. This is basically a thread pool implementation that aims to efficiently schedule `fork/join` tasks among available processors. Inspired by `divide-and-conquer` and recursive approaches to parallelism, a `fork/join` task can spawn (`fork`) new tasks and wait for them to finish (`join`) before progressing with the execution. Currently the `fork/join` implementation uses two core techniques, `adaptive work stealing` [8, 13] and `exponential task splitting` [8], in order to efficiently control the task granularity. Although the programmer does not have much control over the level of parallelism provided by the framework, the implementation has been carefully tuned to ensure a high parallel performance.

3.4 Summary

In summary, the following table compares and contrasts the key features of each of the three languages and their parallelism support listed in order of their level of abstraction.

	Key Features	Parallelism Support
Haskell	functional, lazy evaluation, static/inferred typing	<i>Skeleton</i> : parMap <i>Strategies</i> : parList <i>Semi-Exp</i> : par, pseq <i>Explicit</i> : IO threads
F#	functional, imperative, object oriented, strict evaluation, static/inferred typing, .NET interoperability	<i>Skeleton</i> : pmap <i>Semi-Exp</i> : PLINQ, For Async workflows <i>Explicit</i> : tasks
Scala	functional, imperative, object oriented, strict evaluation, static/inferred typing, Java interoperability	<i>Skeleton</i> : pmap <i>Semi-Exp</i> : Par Coll <i>Explicit</i> : Actors

	Haskell	F#	Scala
<i>Skeletons</i>	parMap Strat: parList	pmap	pmap
<i>Semi-explicit</i>	par, pseq	PLINQ, ParFor Async Wkf	Par Coll
<i>Explicit</i>	threads	tasks	Actors

Haskell is designed as a purely functional language and therefore does not include features for object-oriented and imperative programming. However, it does support foreign language integration e.g. C through the FFI library, which can perform unsafe operations inside the IO monad. F# is mostly functional but its design aims at integration with other paradigms from the offset. Scala is mainly influenced by Java and so many of the language concepts are tied to objects. However, it does provide fairly good functional support, though the syntax differs from more traditional functional languages a bit.

All three languages provide advanced type systems with automatic type inference and support high-level approaches for multi-core parallelism. Platform-wise, Haskell has its own graph-reduction-based runtime-system (STGM), F# compiles down to .NET Common Intermediate Language (CIL) and then runs on the Common Language Runtime (CLR), and Scala compiles to Java Byte Code and runs on the Java Virtual Machine (JVM).

In terms of data structures, the most commonly used in all languages are lists, which are lazy by default in Haskell and strict in the other two languages, and mutable arrays. F# also has sequences, in which evaluation is demand-driven, and LazyList from the PowerPack, which is similar to Haskell lists. LazyList uses caching and allows pattern matching unlike sequences.

4. Implementation

In this section, we provide (parallel) implementations for the n-body problem in Haskell, F# and Scala. A detailed discussion of the Haskell implementation is provided in [32]. We produced the F# and Scala versions based on this implementation. We compare the implementations in terms of performance and programmability.

4.1 N-Body Problem

The n-body problem is a common problem in many areas of science. It involves predicting and simulating the motion of a system of N bodies that interact with each other gravitationally. The simulation proceeds over a number of time steps, where the gravitational forces applied by each body to the others are calculated and then

used to update the velocities and positions of the bodies in the system in each iteration.

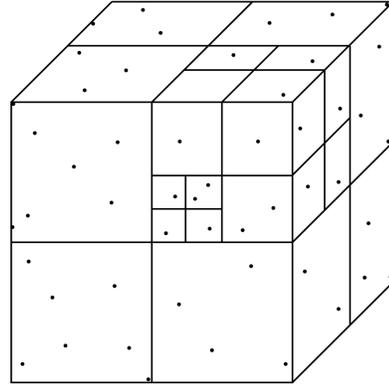


Figure 1. Points in a 3D Barnes-Hut n-body simulation are contained in a region (bounding box) which is sub-divided recursively into smaller regions.

The body-to-body force calculation is a naive method and is not feasible for large number of bodies. Hierarchical force-calculation algorithms, such as the *Barnes-Hut* [3, 23] algorithm provide, an efficient approximation solution. The core of the algorithm consists of two phases: (1) tree construction — an octree is constructed from the list of bodies (see Figure 1); and (2) force calculation — the acceleration due to each body is computed by traversing the tree and approximating bodies that are too far away by using the centre of mass of nearby bodies. The second phase is the most compute-intensive phase. Opportunities for parallelism exist at the tree construction and acceleration calculation stage. The tree construction phase, however, accounts for only a small percentage of the overall time. Therefore we focus on parallelising the force calculation phase, which can be done independently for each body once the tree is constructed. The following is the algorithm for the Barnes-Hut n-body simulation:

```
function doStep (bs,n)
  if (n==0) return bs
  else
    //find the bounding box of the bodies
    bbox = findBounds bs
    //build the BH tree (also calculate centre
      of mass of each region)
    tree = buildTree (bbox,bs)
    //use tree to update velocities and
      positions
    foreach b in bs
      accel = calcAccel b tree
      //deduct acceleration from velocity of
        body b
      b = updateVel b accel
      //move body b
      b = updatePos b
    doStep (bs ,n-1)
```

The *calcAccel* function calculates the acceleration of a given body against the others by traversing the octree. It is the main source of parallelism and can be performed independently for each body.

4.2 Sequential Implementations and Optimisations

The first step in writing parallel code is to come up with an efficient sequential implementation. However, in doing so opportunities for

parallelism should be preserved. Using mutable state to get the best sequential performance destroys most opportunities for parallelism. Instead, keeping the implementation pure eases the parallelisation step. In our case, the sequential Barnes-Hut algorithm is initially implemented in all three languages. The chosen data structure is *list*, as it is the most commonly used in functional languages. Next, we try to improve the sequential implementation by applying a number of generic and language specific optimisations.

4.2.1 Generic Optimisations

A number of general optimisation techniques applies to all three languages. We assess their impact on performance in Section 5. These optimisation techniques are:

Deforestation is used as an instance of (manual) program transformation. For example, elimination of multiple traversals of data structures can improve sequential runtimes by doing fusion, such as merging fold and map, and using function composition with a single map operation instead of two consecutive maps on the same list e.g. `map (f . g) xs` in Haskell; `f << g` in F#; and `f andThen g` in Scala. Some of these transformations are done automatically, e.g. by GHC, with full optimisation enabled.

Tail-call elimination ensures constant stack usage by making sure that recursive functions return an accumulated value in their last call without any further evaluation. This can be often achieved by using the right built-in functions in the language e.g. `foldl` instead of `foldr`. The former is tail recursive and uses an accumulator which is returned in the last call.

Compiler optimisations can be enabled selectively to perform automatic source-to-source transformations, typically on intermediate language code. Notably, the F# and Scala compilers do most optimisations without having to specify any flag (enabled by default), whereas the GHC Haskell compiler allows selective enabling of optimisations and provides several optimisation levels.

4.2.2 Haskell

We start with an initial sequential implementation of the Barnes-Hut algorithm in Haskell, where the expressiveness of the language helps to easily translate the advanced algorithm into functional code. Initially, the program is not well optimised and is unable to execute on a large number of input due to stack overflow. Several iterations of optimisation are required in order to produce an efficient sequential version (reported in [32]). These optimisations are:

Strictness annotation: Often it is not necessary to delay evaluation of values to avoid unnecessary thinking of computations. Forcing evaluation in Haskell is done using the `pseq` primitive, the strict application function (`$(!)`) or simply the exclamation mark (`!`) from `BangPatterns` extension. These annotations are typically placed in data type definitions, to effect every usage of such data. All data types are defined with strict data fields consequently reducing the heap consumption and runtime. Additionally, the `UNPACK` pragma is used to refer to the values directly instead of pointers, thus removing one level of indirection and reducing memory consumption.

foldr/build: Another important optimisation in Haskell is `foldr/build` short cut fusion. This eliminates the intermediate data structures produced by a `build` followed by a `foldr`. The compiler can spot this specific sequence and automatically fuse the code.

4.2.3 F#

The F# version is a direct translation of the Haskell code and the changes are mainly syntactic. In contrast to Haskell's lazy list, we use the default strict list in F#. Other versions using sequence and array are implemented and the results are given in the next section.

F# does not require any optimisations related to laziness as the language is strict by default. After translation from Haskell, the main optimisation involves manually merging fold/map in F#, which is done automatically by using short cut fusion in Haskell. As an example, the operation done in the map is moved into the lambda function of fold.

```
List.fold g acc (List.map (fun x -> f x) xs)
— becomes
List.fold (fun state x -> g state (f x)) acc xs
```

Inlining functions is another way of improving performance. The `inline` annotation indicates that a function definition should be embedded into any code which uses it.

4.2.4 Scala

The Scala version is largely based on the Haskell and the F# implementations and the main differences are in the syntax. This involves wrapping functions inside classes or singleton objects. The two main sequential optimisations are:

Tail recursion optimisation: This generic optimisation technique seems to play a significant role towards achieving good performance in Scala. Tail calls are not natively supported in the JVM, as opposed to .NET, which explains why such a source code transformation is more important than in F#. Although some tail-call optimisation was recently introduced to the Scala compiler, it is still quite basic and only able to convert simple recursive functions into loops, not complicated ones as used in the Barnes-Hut algorithm.

Unnecessary object initialisations removal: Object allocation in Scala is very light-weight, something which is very important as objects are an integral part of the language. Object initialisation, though, causes some additional performance overhead especially when used inside heavy numerical computations such as in the recursive `calcAccel` function, which is the main worker function in the Barnes-Hut code.

4.3 Parallel Implementations and Tuning

The main source of parallelism is the acceleration calculation as confirmed by a high percentage of time spent in this step by time profiling the program. The tree construction stage is insignificant in terms of the overall percentage of time spent in it. Thus, we focus our efforts in parallelising the top level map function that uses the constructed tree and computes the acceleration for each body in the list. The computation of acceleration for each body is independent of the other bodies, which means no synchronisation locks are required.

4.3.1 Haskell

The `strategies` library provides a parallel map implementation `parMap`, which is implemented using `parList`, a higher-order, composable strategy that applies a given evaluation degree to each element in the list in parallel. `parMap` is the starting point to introduce data parallelism in the code. Usually it will give good parallel performance if the list to which it is applied is not too big and the function applied to each list element does enough work to cover the overhead of creating a spark for each list item. With a large input size, `parMap` is inefficient as a spark is generated for each element in the list structure resulting in more overhead than actual benefit of parallelism especially if the work is too fine-grained.

It is usually not a problem to create many sparks in GpH as it amounts to a pointer for each spark created only. However this may lead to too fine-grained parallelism and poor performance.

Parallel tuning: The right balance of spark creation to match the number of cores on the system is important to achieve good parallel performance. If too many sparks are created, they might not end up being taken for execution by the runtime-system; while too few of them may result in under-exploitation of all processing units.

We use strategic chunking as a method to control the number of sparks in Haskell. Two alternative versions of chunking, explicit chunking and clustering, are covered in more detail in [32]. Strategic chunking makes use of an existing higher-order strategy from the library, which performs the chunking implicitly. This involves using `parListChunk` instead of `parList`. The former takes an additional `chunksize` parameter. `parListChunk` breaks the input list into chunks of of the specified chunk size, which is typically calculated depending on the number of processor cores (using `numCapabilities` in GpH) and input size. This ensures that the work is properly balanced among the processors. The two lines of code below are all that is needed to make the parallel implementation scale.

```
chunksize = (length bs) `quot`
            (numCapabilities * 4)
new_bs = map f bs
        `using` parListChunk chunksize rdeepseq
```

4.3.2 F#

We first use asynchronous workflow to implement a parallel map in F#. By marking the function application to each element in the list with the `async` keyword, we have a concurrent map, with each function application not blocking each other. Adding `Async.Parallel` to the pipeline enables the function applications to run in parallel if multiple cores are used. `Async.RunSynchronous` waits to synchronise at the end.

```
let pmap_async f xs =
    seq { for x in xs -> async { return f x } }
    |> Async.Parallel
    |> Async.RunSynchronously
    |> Seq.toList
```

While asynchronous workflow is a fairly easy way to introduce parallelism and get initial speedup, it is also fairly intrusive and changes the code structure. If the main source of parallelism can be identified in one specific higher-order function, or skeleton, one would typically use tasks from the TPL library for independent operations. For instance, similar to the Haskell initial `parMap` implementation, where a spark is created for each list element, we try creating a task for each list element using the task factory pattern from TPL. This surely incurs overhead if the cost of creating a task for an element is higher than the cost of processing (applying a function to) the element.

```
let pmap_tpl_tasks f (xs: list<T>) =
    let createTask x = Task<T>.Factory.StartNew(
        fun () -> f x).Result
    let tasks = xs |> List.map createTask
    tasks
```

```
// chunking
let pmap_tpl_tasks_chunk f (xs: list<T>) =
    let chunks = chunksOf (xs.Length / (numProc *
        2)) xs
    let chunkTask chunk = Task<T>.Factory.
        StartNew(fun () ->
            List.map f chunk).Result
    let tasks = List.map chunkTask (chunks |> Seq.
        toList)
```

```
tasks |> List.concat
```

The code segment above shows explicit task creation for each list element in a naive parallel map implementation (comparable to `parMap` in Haskell). The intention is to compare the overhead of spark versus task creation in Haskell and F# respectively. As we expected this does not give good performance. Thus we use a chunking mechanism to try to limit the number of tasks created. The results are discussed in Section 5.

Using tasks directly is not a good fit for our data-parallel problem. Many higher-level constructs are provided in TPL to achieve a more declarative way of enabling data parallelism. These are typically implemented on top of tasks. PLINQ presents the best choice. It hides the details of task creation and management in its implementation and provides a nice, familiar interface to easily express parallel queries. For example, doing an operation on each element in a list in parallel is enabled by simply marking the container *as parallel* which hints to the underlying system that the latter is to be processed in parallel i.e. converting it into a parallel query. PLINQ uses TPL tasks in the background and handles load balancing across the cores implicitly, though it also offers some limited control. In the PLINQ-style parallel map implementation, the `Select` actually performs a map operation on each element.

```
let pmap_plinq f (xs: list<T>) =
    xs.AsParallel()
        .Select(fun x -> f x) |> Seq.toList
```

Imperative style programming: The other main construct for data parallelism is `Parallel.For/ForEach`. However, this does not prove to be convenient with lists but is mostly useful with mutable arrays, where the action inside the loop is to update elements in the array. Implementing a different version of the algorithm that uses arrays enables us to use `Parallel.For` to introduce parallelism and thus to examine the effect of inplace update.

```
let pmap_tpl_parfor f (xs: array<T>) =
    let new_xs = Array.zeroCreate xs.Length
    Parallel.For(0, xs.Length, (fun i ->
        new_xs.[i] <- f (xs.[i])) |> ignore
    new_xs
```

Alternatively, there already exists a parallel map in the `Array.Parallel` namespace which is a basic implementation and uses `Parallel.For` behind the scene.

```
let res = Array.Parallel.map f arr
```

Tuning: PLINQ and TPL provide some options for tuning, although we find that the default settings are usually sufficient.

Maximum degree of parallelism: A thread pool is used to schedule tasks and the number of threads can be controlled by using maximum degree of parallelism. This configuration specifies the maximum number of concurrently executing tasks.

`Parallel.For` and similar methods take an additional parameter that specifies this parallel option. In PLINQ, the parameter `AsParallel().MaxDegreeOfParallelism` can be set explicitly to the same effect. Being an upper bound this parameter is used to restrict the amount of parallelism at one time in the execution, but does not ensure that the specified number is generated.

Chunking/Partitioning: We implement custom chunking to control the granularity of tasks created explicitly as above — to achieve similar effect as in Haskell.

There is also a `Partitioner` class with static methods to partition collections. It supports partitioning with dynamic allocation, but also range partitioning with static allocation. A partitioner can be passed as argument to `Parallel.For` to specify a custom partitioning. Therefore, it is best used with arrays as it gives the intervals for each partitions.

4.3.3 Scala

Through the use of Parallel Collections in Scala, parallelisation is semi-explicit by using the keyword `par` to call a parallel version of the list, which implements parallel operations.

```
nbody.par.map((b: Body) => new Body(b.mass,
    updatePos(b), updateVel(b))).seq
```

The method `par` is applied on the list of bodies, which calls the default parallel implementation of the list. When we subsequently apply the `map` function on the parallel list, the parallel map function is invoked on the list elements. In this way, it achieves a similar separation of coordination and computation as `parList` in GpH. Finally, we have to convert the result to a sequential list by applying the method `seq` on the results.

Parallel Collections require only small changes to the code to achieve initial speedup. The main disadvantage is that the framework does not currently provide much control over the parallelism. As an example, it is not possible to control how many threads are spawned or to define the size of the underlying thread pool. Instead, these details are handled by the underlying implementation, which uses sophisticated work stealing and chunking techniques.

In a second approach, we implement parallel map skeletons using *futures* from the actors library. A future abstracts over send and receive primitives and represents an object that is created to store a result that has not yet been computed. The result is computed concurrently at a later time and can be collected on demand. A parallel map skeleton is implemented as follows:

```
def pmap[T](f: T => T) (xs: List[T]): List[T] = {
    val tasks = xs.map((x: T) => Futures.future { f
        (x) })
    tasks.map(future => future.apply())
}
```

In `pmap` a `future` is explicitly created for each element in the given list, mapping the given function on the corresponding list element. The results of the parallel map are then returned to the user as the output of the skeleton.

The second skeleton we implement is a parallel map using chunking to explicitly control the granularity of the parallelism, directly corresponding to the initial Haskell implementation:

```
def chunk[T](xs: List[T], size: Int): List[List[T]] = xs.isEmpty match {
    case true => List()
    case false =>
        val split = (xs.take(size), xs.drop(size))
        (split._1) +: chunk(split._2, size)
}
```

```
def pmap_chunk[T](f: T => T, size: Int) (xs: List[T]): List[T] = {
    val chunks = chunk(xs, size)
    val task_chunks = chunks.map((c: List[T]) =>
        Futures.future { c.map((x: T) => f(x)) })
    val tasks = task_chunks.map(future => future.apply())
    tasks.flatten
}
```

5. Results

5.1 Experimental Setup

Platforms: All three languages are supported both on Linux and Windows platforms either natively or through independent (open-source) implementations. This provides ground for comparison of the language implementations across the two platforms and for discussion of the results on each. Haskell’s GHC implementation is

cross-platform. GHC offers the option to compile code down to C and run on a standard C compiler. F#’s official Microsoft implementation is intended to run on the .NET Framework on machines running Windows. The open-source implementation of the runtime, Mono, is available to compile and run F# code under Linux. Scala runs on JVM, which has good support on both platforms. The following are the machines used for the experiments:

	Linux	Windows
Version	CentOS 5.8	XP
Architecture	64-bit	32-bit
CPU	Intel Xeon E5410 2.33GHz	Intel Core i7 860 2.80GHz
# cores	8	4 (8 HyperThreads)
RAM	7986 MB	3520 MB

Language Implementations: The following are the language implementations used with the corresponding version numbers on each platform:

	Linux	Windows
Haskell	GHC 7.4.1	GHC 7.4.1
F#	F# 2.0 / Mono 2.11.1	F# 2.0 / .NET 4.0
Scala	Scala 2.10 / JVM 1.7	Scala 2.10 / JVM 1.7

Input: All measurements are taken using 80,000 bodies as the input and the execution is based on 1 iteration of the n-body simulation. We focus on a single iteration to assess the potential for parallelism in this application core for each language, rather than viewing it as an application tuning exercise.

5.2 Baseline

We use the all-pairs implementations from the Computer Language Shootout website¹ as baseline for comparison. The implementations on the shootout webpage are highly-optimised by experts in each language community. However, the implementations are impure and they make heavy use of inplace updates and other unsafe constructs provided in the languages in order to get the best performance out of the implementation. This approach however seriously hampers parallelisation. Often the whole program would have to change in order to parallelise it. Therefore, we do not take these versions as starting points for our parallelisation. The original shootout runtimes on a Linux x64 Intel Q6600 machine are taken for 5 bodies and 50 million iterations. We take measurements on our machines using 16000 bodies and 1 iteration only, and we exclude the input generation and energy calculation times as we are interested in the main iteration and parallelising it. The runtimes are given in Table 1.

Table 1. Baseline (a) — language shootout results (all-pairs). The numbers in brackets are slow-downs w.r.t. the Haskell version.

	Linux (Original) 5 bodies, 50M iterations	Linux (16k bodies, 1 iteration)	Windows
Haskell	25.23 (1.00)	9.24 (1.00)	7.66 (1.00)
F#	41.36 (1.63)	9.37 (1.01)	4.88 (0.63)
Scala	23.47 (0.93)	5.51 (0.59)	14.25 (1.86)

Both the original runtimes and those taken on our Linux machine show that the Scala implementation is the fastest, followed by Haskell, then F#. However, on the Windows platform, interestingly the F# version, which was the slowest on Linux, performs best.

¹<http://shootout.alioth.debian.org/>

This highlights that the F# .NET implementation on that platform is very well-tuned and Microsoft technologies, for compiler and virtual machine, integrate well together. The Scala version shows the opposite picture: it is the fastest version on Linux but the slowest on Windows. Again, we expect that this is due to Linux being the main target platform for this project. In contrast, Haskell exhibits best sequential performance portability.

We also use our pure all-pairs implementations, which do not use destructive updates, as a baseline. The results are shown in Table 2. Haskell gives the best performance under Linux, and second best under Windows, though closely followed by F# under the same platform. Scala gives decent performance on both platforms. F# on Linux is very slow, by a factor of 8.5 compared to Windows for the pure baseline, as opposed to 1.9 for the impure baseline. We also observe that, as expected, the pure implementations are slower than the impure ones (in Tables 1) between 1.7 to 3.9 times. We intentionally do not use impure features to make parallelisation safe and easy.

Table 2. Baseline (b) — our pure all-pairs implementations (16k bodies, 1 iteration). The numbers in brackets are slow-downs w.r.t. the Haskell version.

	Linux	Windows
Haskell	20.77 (1.00)	15.55 (1.00)
F#	123.22 (5.93)	14.45 (0.92)
Scala	21.88 (1.05)	24.71 (1.58)

Comparing the runtimes from Tables 1 and 2, we can also assess the sequential performance of our all-pairs implementations, giving an estimate of the quality of our sequential code: 49.3% for Haskell, 33.8% for F#, 57.5% in Scala. Avoiding any use of impure features we lose some performance initially, but we gain ample opportunities for parallelisation in exchange. Later in this section we will show that selective usage of impure features *after* parallelisation can further enhance performance. We believe that efficiencies for the Barnes-Hut algorithm are similar, but in the absence of similarly tuned sequential implementations we cannot make a direct comparison for this algorithm.

In the following subsections, we see how these figures relate to the pure Barnes-Hut implementations, where the algorithm is more complex than the all-pairs, on the two platforms. The discussions focus on 3 metrics: performance, programmability and pragmatic aspects of the languages.

5.3 Performance Evaluation

Tables 4 and 5 summarise the runtimes and speedups in each language on Linux and Windows respectively. As ancillary data, Table 3 shows the peak memory usage under the Windows platform, as observed by an external, OS-level task manager.

We use the same data structure, in this case list, across the different implementations to have comparable results. The results highlight a number of interesting points.

5.3.1 Sequential Performance

Tables 4 and 5 show that Haskell gives the best sequential performance on both platforms. This has been made possible due to extensive sequential optimisations, using a range of techniques.

Most notably, the initial, naive Haskell version — without strict data fields — gives a runtime of 479.94s. By using strict data fields and UNPACK pragma, the runtime goes down to 33.02s, amounting to a sequential speedup of 14.5. Enabling foldr/build optimisations by code restructuring, as a GHC specific compiler optimisation, gives a further 23% reduction in runtime to 25.28s. All results are measured under Linux. The Haskell sequential runtime is better

than that of F# under Windows but, using lazy evaluation, it uses a higher memory footprint than F# as shown in Table 3.

Table 3. Peak memory usage on a Windows machine

	Haskell	F#	Scala
Sequential	57 MB	32 MB	58 MB
Parallel (4 cores)	71 MB	36 MB	62 MB

The F# version is a direct translation from Haskell with some F# specific optimisations. Some optimisations native to Haskell e.g. strictness annotations, are not required in F#. Other program transformations, in particular merging fold and map operations, thereby eliminating intermediate data structures, have to be done manually in F# and this improves the runtime from 28.43 to 22.15s. Inlining of functions reduces the runtime by 5% to 21.12s.

Another general observation is that the Mono implementation of the .NET runtime, used under Linux, is not as well optimised as the corresponding Microsoft .NET implementation, used under Windows. Since the Mono project is now focusing on providing a .NET infrastructure on embedded systems, rather than focusing on high-performance computing, this is not surprising. However, this is one of the first systematic comparisons of both platforms for parallel computation. The difference in runtime is particularly remarkable, since the hardware used for running the F#/Mono instance was faster than the hardware for running F#/.NET.

The Scala version is the slowest one under Windows, but significantly better under Linux. This behaviour is consistent with the all-pairs baseline results in Tables 1 and 2. We attribute this difference mainly to the Scala compiler and the memory management overhead imposed by it. It has been reported elsewhere [28] that Scala makes heavy use of boxed types, resulting in a fairly high memory footprint, as shown in Table 3 where the peak memory consumption of the strict Scala implementation (58 MB) is even higher than the lazy Haskell implementation (57 MB). On the Windows machine the smaller amount of main memory will cause this high memory consumption to have a stronger impact on total runtime. Additionally, as a mixed paradigm language, Scala makes heavy use of objects, resulting in initialisation overhead. Under Linux, the initial Scala implementation, without the optimisations described in Section 4.2.4, gives a runtime of 55.44s. By using tail recursion optimisation, the runtime goes down to 45.48s (-18%). Removing unnecessary object instantiations helps to further improve the performance to 39.04s (-14%).

Using Lazy Data Structures: All 3 languages support lazy data structures and here we compare their sequential performance. As reported above, the unoptimised Haskell version, with unconstrained lazy data structures loses a factor of 14.5 in runtime.

The LazyList version in F#, as a direct comparison with this initial Haskell version, exhibits an increase in sequential runtime from 118.12 to 604.19s (Linux), and from 21.12 to 81.68s (Windows), representing a factor of 5 and 4, respectively. The memory usage under Windows peaks at 115 MB (from 32 MB using the default strict list). This suggests that this library is not well-optimised for LazyList and fewer compiler optimisations, aiming at eliminating unnecessary laziness are applied. It is worth noting that this structure, available as part of the PowerPack, is not officially supported and is still under development.

Similarly, in Scala, the use of streams as lazy structures instead of lists results in an increase in sequential runtime from 39.04 to 89.35s (Linux), and from 66.65 to 92.24s (Windows). The memory usage under Windows peaks at 100 MB, which is a 72% increase from the 58 MB using the default strict list.

Table 4. Runtimes (in seconds) on *Linux* (8 cores)

# cores	Haskell		F#				Scala		
	GpH		AsyncWork	TPL			ParColl	Actors	
	parMap	parMapChunk		Tasks	Tasks/chunk	PLINQ		pmap	pmapchunk
seq	25.39	25.28	118.12	118.12	118.12	118.12	39.04	39.04	39.04
1	25.91	26.38	211.78	197.72	209.76	196.14	48.02	45.38	40.01
2	25.77	14.48	129.07	154.09	162.32	120.78	25.72	25.18	22.34
4	22.69	7.41	89.63	128.99	134.54	80.91	16.41	16.42	14.88
8	23.17	4.50	70.41	120.26	122.45	70.67	13.48	14.34	13.26

Table 5. Runtimes (in seconds) on *Windows* (4 cores plus hyperthreading)

# cores	Haskell		F#				Scala		
	GpH		AsyncWork	TPL			ParColl	Actors	
	parMap	parMapChunk		Tasks	Tasks/chunk	PLINQ		pmap	pmapchunk
seq	17.64	17.64	21.12	21.12	21.12	21.12	66.65	66.65	66.65
1	17.77	18.05	21.26	23.31	21.10	21.39	66.96	68.30	67.24
2	17.61	9.41	16.96	26.47	21.36	17.32	57.96	58.63	58.66
4	16.94	6.80	10.18	36.06	21.50	10.56	34.48	33.64	33.84
8 (<i>HT</i>)	17.61	4.77	8.82	36.28	21.05	8.64	26.18	24.74	25.28

5.3.2 Parallel Performance

Figures 2(a) and 2(b) summarise the speedups obtained with the best versions of each language implementation on Linux and Windows, respectively. For F#, we use PLINQ speedups instead of Async Workflow as the differences are small and we want to compare *parallel* construct in the plots. Tables 4 and 5 elaborate on the runtimes of several versions with the best parallel runtimes in each language highlighted. Comparing the performance of the languages, Haskell displays the best speedups, up to 5.6 on 8 cores, and remains scalable. This is achieved through strategic chunking to improve thread granularity.

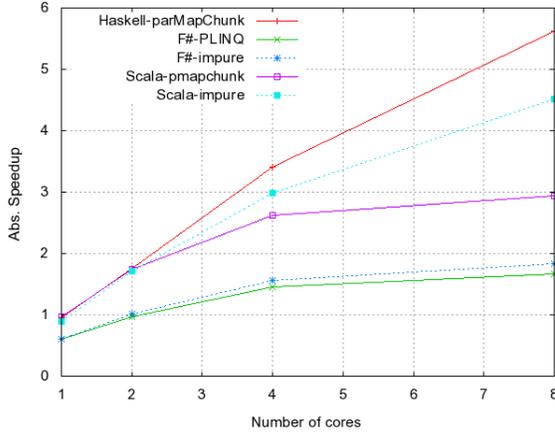
Since this is a data-parallel application, with limited scope for thread-subsumption, such explicit granularity control is crucial, as can be seen from the poor performance of the naive `parMap` implementation, which generates a spark for each list element. In terms of absolute parallel performance, the Haskell version is 3 times faster than the Scala version on an 8-core Linux machine, and 1.5 times faster than F# on an 4-core Windows machine.

Due to the poor performance of Mono, for F# the Windows version is the more interesting one. Here the implementation achieves a respectable speedup of 2 on 4 cores, which increases to 2.4 when using hyperthreading. Notably, the highest-level PLINQ implementation is the best performing on this platform, although speedup in itself is not as good as in the Haskell version. Interestingly, the heap consumption of F# is significantly lower, even with the unoptimised version where it remains the same, than that of the Haskell version (Table 3), but this does not translate into faster runtimes. We conjecture that this is mainly due to GHC performing more aggressive optimisations than F#. We note that the task-based implementations in effect result in a slow-down, mainly due to the high task management overhead, which is reduced when employing chunking. In contrast to the GpH version, F# tasks are mandatory, so this overhead is more pronounced than in GpH. The highest level PLINQ implementation, best suited for data-parallelism gives almost as good a result as the lower level asynchronous workflow implementation (-3.7% on 4 cores). Given the simplicity of the PLINQ code (see Section 4.3) this is a strong argument in favour of this abstraction mechanism for data-parallel code.

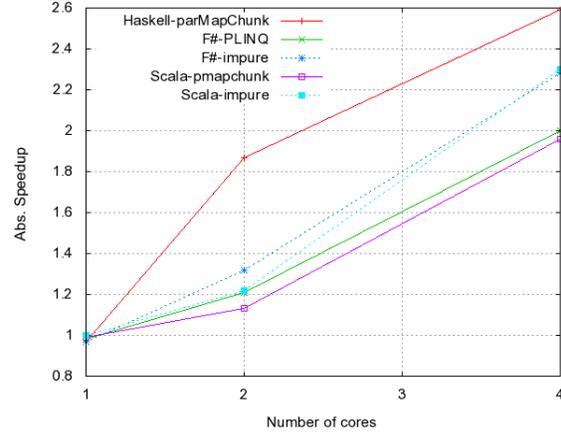
For the Scala version we focus on the better performing Linux version. All three parallel versions exhibit good speedups, although trailing the Haskell results. In this implementation, chunking has a far lower impact on performance compared to the Haskell version. Together with the good 1 processor performance this indicates very efficient task management for actor code in Scala. The highest level parallel collection implementation is within 1.7% of the 8-core performance, almost as good as a tuned actor implementation. The parallel performance of the pure version tails out for higher core counts, with an 8-core speedup of 2.9 using the Parallel Collections, but the impure implementation, discussed below, achieves a speedup of 4.5. This simplest version achieves almost as good a performance as the lower-level actor-based implementations.

Mutable Data Structures: The remaining results use impure language features, in particular mutable data structures, to further improve performance. The use of arrays in F#, which update the body inplace, gives only a small performance gain under Windows, a 1.7% decrease in sequential runtime from 21.12s to 20.76s (Table 6). Interestingly, the memory usage is the same as for the list version. This might be due to the imprecision of using an external, OS-level tool to determine peak usage, as opposed to maximum residency. Under Linux, where the sequential runtime is already slower by a high factor, arrays give a 19% improvement from 118.12 to 95.33s in runtime. The main advantage of using arrays in this context, though, is as mentioned earlier: we can use the built-in parallel map from the `Array.Parallel` module and the `Parallel`. For with the default or custom partitioning. The default partitioner, which creates partitions based on the processor count and input size, exhibits very good parallel performance that does not improve significantly with a custom partitioner (Table 6), leading us to say that the default tuning performed is very efficient.

We also developed an impure functional implementation in Scala using inplace updates. In this case, the maximum residency drops from 114.75 MB to 91.72 MB (20% reduction) on the Linux platform, as obtained by internal, JVM-level monitoring. This directly translates into faster sequential and parallel runtimes than the pure versions under Linux (see Table 7). In Scala this difference is quite remarkable, improving speedup to 4.51 (on 8 cores), as opposed to 2.9 in the pure version (Figure 2(a)). This improve-



(a) Linux (8 cores)



(b) Windows (4 cores)

Figure 2. Absolute speedups on the two platforms for the best versions in all languages and impure versions in F# and Scala.

ment in parallel speedup is most likely due to the reduction in heap contention on this shared memory architecture, in this more memory efficient version. The impure Scala implementation also gains significantly better sequential performance on the Windows platform: 47.86s (28% runtime reduction from the pure Scala implementation) and a 4-cores parallel runtime of 20.78s (39% runtime reduction and 14% speedup improvement from the parallel pure implementation using `pmap_chunk`).

Table 6. F# runtime results (in seconds) using arrays on *Windows* (4 cores plus hyperthreading).

# cores	Best Pure	Array.Par.map	Parallel.For	
			default	partitioner
seq	21.12	20.76	20.76	20.76
1	21.39	21.56	21.2	21.35
2	17.32	15.83	15.66	16.05
4	10.56	9.09	9.09	9.32
8	<i>8.64</i>	7.33	7.33	7.15

Table 7. Runtime results (in seconds) for the impure Scala implementation on *Linux* (8 cores).

# cores	Best Pure	ParColl	pmap	pmapchunk
seq	39.04	32.81	32.81	32.81
1	40.01	36.92	34.86	36.09
2	22.34	17.66	22.05	19.14
4	14.88	11.23	12.55	10.96
8	13.26	8.47	8.57	7.27

5.4 Programmability

Due to the high-level nature of all three languages, introducing parallelism to a pure version is easy and often amounts only to using a suitable data-parallel skeleton, e.g. `parMap` instead of a `map` in the computational core of the application. Implementing the impure versions, we start from a pure version to introduce parallelism and then add inplace updates selectively in such a way that does not require locks as the operations are independent. In the Haskell version we use evaluation strategies, which separate the parallel code from the main computation logic. In another paper [32] we compare

this implementation with alternatives for introducing parallelism in Haskell, in particular with Eden [16] and the `ParMonad` [18], achieving similar performance. In direct comparison, a first parallel version is more easily achieved in Haskell but for optimal performance some parallel performance tuning is needed. In particular, the naive `parMap` does not work well in this example, because its implementation does not automatically introduce chunks, and for competitive performance the `parMapChunk` variant is needed. One strength of Haskell is the modularity in expressing coordination code, which makes it easy to implement customised chunking as discussed in [32]. Haskell provides the least intrusive mechanism for parallelisation. Introducing parallelism with chunking amounts to two additional lines of code. It is also easy to separate issues such as chunking by composing higher-order strategies, for example, using a clustering strategy.

F# supports several paradigms of parallel programming at different levels of abstraction, helping the programmers to make the transition from sequential to parallel programming. The preferred mechanism for data-parallelism is to use the SQL-like `PLINQ` language, which represents the highest level of abstraction, and still gives close to optimal performance. This is an advantage of using F# for data-parallel applications that fit the `PLINQ` abstractions.

Scala offers parallelism support through two libraries. On a data structure level, parallelism is very easily introduced through the keyword `par`, which calls a parallel version of a collection causing most subsequent functions applied on it to be executed in parallel. The lower-level actors message passing model offers explicit messaging between actors, but also allows to build higher-level solutions on top of futures for introducing parallelism.

Optimisations are important for good performance but must be chosen carefully to avoid sequentialisation. The Haskell version requires some modification to the laziness of the code in order to improve sequential performance. In F# and Scala, some optimisations are necessary and easily introduced with the functional style e.g. ensuring recursive functions make tail-calls and avoiding multiple traversals of data structures using function composition in the function argument to a `map`. Several of these (manual) optimisations are performed automatically by GHC for the Haskell version.

Other language features such as purity in each language help to structure the program in a modular way, gaining separation between the main components of the code. For instance, the main n-body code is written in a purely functional way, which makes it easy to reason about the program, offers robustness and modularity, and

efficiently integrates with the world, for example through the main function that generates the input and prints the output.

Both F# and Scala allow integration of object-oriented features with functional programming. Although, we tried to integrate some of the available object-oriented features (e.g. classes) in our F# and Scala implementations, we did not notice any difference in the sequential or parallel performance. That makes intuitive sense because object-orientation is mainly used to help towards large-scale software development, by using features such as polymorphism, inheritance and encapsulation. In this small high performance computing application, though, we found that we did not really require these features to achieve good performance.

One of the main difficulties in using Haskell for high performance computing is to understand the implications of laziness. It is hard to predict when expressions are evaluated, and to estimate how much work is involved. Evaluation strategies provide a powerful mechanism to control evaluation order and degree where needed in order to achieve efficient parallelism. However, laziness allows us to make use of infinite data structures. F# and Scala on the other hand have strict evaluation by default. This makes it easier to predict the performance of the program.

Summarising these programmability aspects, we can provide the following guidance on choosing a language. F# with the PLINQ abstraction is particularly suitable for flat, data-parallelism and it is easy to use for programmers fluent in SQL-style database queries. GpH provides the least intrusive programming model, and is therefore a good choice when aiming to minimise code changes when parallelising an application. GpH and to a lesser extent F# provide rich abstraction hierarchies, that allow the programmer to pick constructs with a suitable degree of control for the application. In direct comparison of these two languages, GpH emphasises a purely declarative model, whereas F# provides more direct control of execution on a lower level. Finally, in terms of the underlying operating system, Haskell is a viable choice on both platforms, exhibiting best sequential performance portability. The same cannot be said for the other two languages: F# works best on Windows, while Scala is best tuned on Linux.

5.5 Pragmatics

Haskell comes with powerful tool support which is helpful in optimising both sequential and parallel algorithms. As an example, time and heap allocation profiling reports, both textual and visual, are useful in identifying the hot-spot of the execution and potential space leaks. Threadscope is a parallel visualisation tool particularly useful to see work distribution across the number of cores. Although F# comes with very good tool support, such as a profiler, unfortunately it is only available on the ultimate version of Microsoft's Visual Studio. Due to the lack of these tools, it was difficult for us to find out why the same code in Haskell performs badly in F#. Free tools are difficult to find, as most third party options are commercial. Scala is based on the JVM and, thus, enjoys a wide range of both monitoring and analysis tools, such as VisualVM, which we used as a JVM-level monitor.

From a pragmatic point of view the rich tool support through VisualStudio and its apparent backing by Microsoft make it a very attractive language in particular for programmers previously unaware of functional programming. There is a tendency in its implementation to hide details of the parallel execution from the programmer, though. While this encourages a high level of abstraction, best supported through data-parallel PLINQ, it also complicates the tuning of the program for the expert parallel programmer.

6. Conclusion

In this paper we compare parallel implementations of the Barnes-Hut algorithm for solving the n-body problem, implemented in

the purely functional language Haskell and in the modern multi-paradigm languages F# and Scala. We assess both programmability and parallel performance, when executing on state-of-the-art multicore machines. We use a number of alternative parallel programming constructs provided in each language, in particular GpH in Haskell, Asynchronous Workflows and TPL in F#, and Parallel Collections and Actors in Scala.

The type of parallelism in this application is data-oriented and to this end, we implemented a number of parallel map skeletons in F# and similarly in Scala. We achieve speedups up to 5.62 in Haskell (8 cores), 2.28 in F# (4 cores) and 4.51 in Scala (8 cores).

With a caveat that our observations are based on parallel variants of only one application, we draw the following conclusions:

- Across all languages, the version using the highest abstraction level also produced the (near) best runtimes.
- Providing first-class parallelism support in the language, through primitives rather than annotations or libraries, is important in the Haskell version in order to explicitly tune thread granularity, e.g. using higher-order functions in a chunking `parMap`.
- Careful (sequential) optimisation of the lazy Haskell version results in sequential performance, surpassing that of the strict F# and Scala versions.
- Aggressive, sequential code optimisations, using impure language features early on, seriously hamper the parallelisation of the code, as can be seen from the (sequential) implementations on the language shootout page, which are a poor starting point for parallelisation due to enforced sequentialisation on mutable data structures. However, selective use of impure features at the end of the parallelisation process can gain notable additional performance, demonstrated, e.g. in the Scala implementation.
- The poor performance of F# on Linux, due to a fairly low-performance .NET implementation provided by Mono, does not make F# a viable choice for parallelism on Linux at the moment.
- The additional expressive power provided by lower-level actor-based code in Scala does not manage to improve performance significantly and therefore the extra programming effort is not justified in this version.

In comparison to other mainstream parallel languages, we found that all languages provide fairly high-level constructs for parallelism but the degree of control provided in them differs. For instance, Haskell allows initial parallelism to be easily specified, as parallel versions of well-known higher-order functions. Since parallelism is provided through primitives, rather than annotations or libraries, the full power of the language is available, facilitating user-customisable chunking to tune parallel performance. On the other hand, F# and Scala confine most control of parallelism inside the runtime-system implementation, aiming for automatic management without any programmer input. Therefore tuning parallelism is more difficult for the expert parallel programmer in these languages.

Finally, we remark that all 3 languages provide easy-to-use, high-level and efficient support for parallelism. Haskell has an edge in its rich libraries and aggressive optimisation. Through laziness, top-level parallelism can be specified in a more modular way, minimising the need to modify components and avoiding any notion of explicit threads. Thus, Haskell is the language of choice for minimally intrusive parallelism. F# and Scala provide several mechanisms for parallelisation, some as high level as Haskell's, some lower level with more potential for performance tuning. Thus, these languages provide a more direct route to controlling low level aspects and can use object-oriented and impure language features.

An online version of this paper, together with the source code of all parallel versions, is available at: <http://www.macs.hw.ac.uk/~dsg/gph/papers/abstracts/fhpc12.html>.

Acknowledgments

We would like to thank our colleagues in the Dependable Systems Group of Heriot-Watt University for their helpful discussions and feedback. We also thank the contributors to the SICSA MultiCore challenge (<http://www.macs.hw.ac.uk/sicsawiki/index.php/MultiCoreChallenge>) for their comparative work on parallel n-body implementations.

References

- [1] G. Agha. *Actors: a Model of Conc. Computation in Distributed Systems*. MIT Press, 1986. ISBN 0262010925.
- [2] J. Armstrong, R. Virding, C. Wikström, and M. Williams. *Conc. Program. in Erlang*. Prentice Hall, second edition, 1995. ISBN 978-0135083017.
- [3] J. Barnes and P. Hut. A Hierarchical O(N log N) Force-calculation Algorithm. *Nature*, 324:446–449, Dec. 1986. DOI: <http://dx.doi.org/10.1038/324446a0>.
- [4] C. Campbell, R. Johnson, A. Miller, and S. Toub. *Parallel Program. with Microsoft .NET — Design Patterns for Decomposition and Coordination on Multicore Architectures*. Microsoft Press, Aug. 2010. URL <http://msdn.microsoft.com/en-us/library/ff963553.aspx>. ISBN 0-7356-5162-0.
- [5] B. L. Chamberlain, D. Callahan, and H. P. Zima. Parallel Programmability and the Chapel Language. *Int'l J. of High Performance Computing Applications*, 21(3):291–312, Aug. 2007. DOI: <http://dx.doi.org/10.1177/1094342007078442>.
- [6] P. Charles, C. Grothoff, V. Saraswat, C. Donawa, A. Kielstra, K. Ebcioglu, C. von Praun, and V. Sarkar. X10: an Object-oriented Approach to Non-uniform Cluster Computing. In *OOPSLA'05 — Int'l Conf. on Object Oriented Program. Systems Languages and Applications*, pages 519–538. ACM Press, 2005. DOI: <http://dx.doi.org/10.1145/1094811.1094852>.
- [7] M. Cole. *Algorithmic Skeletons: Structured Management of Parallel Computation*. MIT Press, Cambridge, MA, 1989. ISBN 0262530864.
- [8] G. Cong, S. Kodali, S. Krishnamoorthy, D. Lea, V. Saraswat, and T. Wen. Solving Large, Irregular Graph Problems Using Adaptive Work-Stealing. In *ICPP'08 — Int'l Conf. on Parallel Processing*, pages 536–545, Portland, OR, Sept. 2008. 10.1109/ICPP.2008.88.
- [9] J. Dean and S. Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. *Commun. ACM*, 51(1):107–113, Jan. 2008. DOI: <http://dx.doi.org/10.1145/1327452.1327492>.
- [10] M. Frigo, C. E. Leiserson, and K. H. Randall. The Implementation of the Cilk-5 Multithreaded Language. In *PLDI'98 — Conf. on Program. Language Design and Impl.*, pages 212–223, Montreal, Quebec, Canada, June 1998. DOI: <http://dx.doi.org/10.1145/277650.277725>.
- [11] P. Haller and M. Odersky. Scala Actors: Unifying Thread-based and Event-based Programming. *Theor. Comput. Sci.*, 410:202–220, 2009. DOI: <http://dx.doi.org/10.1016/j.tcs.2008.09.019>.
- [12] D. Kranz, R. Halstead Jr., and E. Mohr. Mul-T: A High-Performance Parallel Lisp. In *PLDI'89 — Program. Languages Design and Impl.*, pages 81–90, Portland, Oregon, June 21–23, 1989. DOI: <http://dx.doi.org/10.1145/73141.74825>.
- [13] D. Lea. A Java fork/join Framework. In *Java'00 — ACM 2000 Conf. on Java Grande*, pages 36–43. ACM Press, 2000. DOI: <http://dx.doi.org/10.1145/337449.337465>.
- [14] D. Leijen, W. Schulte, and S. Burckhardt. The Design of a Task Parallel Library. In *OOPSLA'09 — Int'l Conf. on Object Oriented Program. Systems Languages and Applications*, pages 227–242. ACM Press, 2009. DOI: <http://dx.doi.org/10.1145/1640089.1640106>.
- [15] X. Leroy, D. Doligez, A. Frisch, J. Garrigue, D. Rémy, and J. Vouillon. *The OCaml System*. Inst. National de Recherche en Informatique et en Automatique (INRIA), July 2011. Documentation and User's Manual.
- [16] R. Loogen, Y. Ortega-Mallén, and R. Peña Marí. Parallel Functional Programming in Eden. *J. Funct. Program.*, 15(3):431–475, May 2005. DOI: <http://dx.doi.org/10.1017/S0956796805005526>.
- [17] D. K. Lowenthal, V. W. Freeh, and G. R. Andrews. Using Fine-grain Threads and Run-time Decision Making in Parallel Computing. *J. Parallel and Distributed Computing*, 37(1):41–54, 1996. DOI: <http://dx.doi.org/10.1006/jpdc.1996.0106>.
- [18] S. Marlow, S. Peyton Jones, and S. Singh. Runtime Support for Multicore Haskell. In *ICFP'09 — Int'l Conf. on Funct. Program.*, pages 65–78, Edinburgh, UK, Aug. 2009. ACM Press. DOI: <http://dx.doi.org/10.1145/1596550.1596563>.
- [19] S. Marlow, P. Maier, H.-W. Loidl, M. K. Aswad, and P. Trinder. Seq no More: Better Strategies for Parallel Haskell. In *Haskell'10 — Haskell Symposium*, Haskell '10, pages 91–102, New York, NY, USA, 2010. ACM. DOI: <http://dx.doi.org/10.1145/1863523.1863535>.
- [20] D. C. J. Matthews and M. Wenzel. Efficient Parallel Programming in Poly/ML and Isabelle/ML. In *DAMP10 — Declarative Aspects of Multicore Program.*, pages 53–62, Madrid, Spain, Nov. 2010. DOI: <http://dx.doi.org/10.1145/1708046.1708058>.
- [21] E. Mohr, D. Kranz, and R. Halstead Jr. Lazy Task Creation: A Technique for Increasing the Granularity of Parallel Programs. *IEEE Trans. on Parallel and Distributed Systems*, 2(3):264–280, July 1991. DOI: <http://dx.doi.org/10.1109/71.86103>.
- [22] M. Odersky et al. An Overview of the Scala Programming Language. Technical Report LAMP-REPORT-2006-001, EPFL Lausanne, Switzerland, 2006. Second edition.
- [23] S. Pfalzner and P. Gibbon. *Many-Body Tree Methods in Physics*, chapter 2: Basic Principles of the Hierarchical Tree Method. Cambridge University Press, 1996. DOI: <http://dx.doi.org/10.1017/CB09780511529368>. ISBN: 9780511529368.
- [24] R. Plasmeijer and M. van Eekelen. *Funct. Program. and Parallel Graph Rewriting*. Addison-Wesley, 1993. ISBN 0-201-41663-8.
- [25] A. Prokopec, P. Bagwell, T. Rompf, and M. Odersky. A Generic Parallel Collection Framework. In *EuroPar'11 — Int'l Conf. on Parallel Processing*, pages 136–147. Springer Verlag, 2011.
- [26] J. Reinders. *Intel Threading Building Blocks: Outfitting C++ for Multi-core Processor Parallelism*. O'Reilly, 2007. ISBN 0596514808.
- [27] J. Reppy, C. Russo, and Y. Xiao. Parallel Concurrent ML. In *ICFP'09 — Int'l Conf. on Funct. Program.*, pages 257–268, Edinburgh, UK, Aug. 2009. ACM Press. DOI: <http://dx.doi.org/10.1145/1596550.1596588>.
- [28] A. Sewe, M. Mezini, A. Sarimbekov, and W. Binder. Da capo con scala: Design and Analysis of a Scala Benchmark Suite for the Java Virtual Machine. In *OOPSLA'11 — Int'l Conf. on Object Oriented Program. Systems Languages and Applications*, pages 657–676. ACM Press, 2011. DOI: <http://dx.doi.org/10.1145/2048066.2048118>.
- [29] Shootout. The Computer Language Benchmarks Game, June 2012. URL <http://shootout.alioth.debian.org>.
- [30] Sun. The Fortress Language. Talks and Posters, June 2012. URL <http://research.sun.com/projects/plrg>.
- [31] D. Syme, A. Granicz, and A. Cisternino. *Expert F#*. Apress Academic, 2007. ISBN 1590598504.
- [32] P. Totoo and H.-W. Loidl. Parallel Haskell Implementations of the n-body Problem. *Conc. and Comp.: Practice and Experience*, 2012. To appear.
- [33] P. Trinder, K. Hammond, H.-W. Loidl, and S. Peyton Jones. Algorithm + Strategy = Parallelism. *J. Funct. Program.*, 8:23–60, January 1998. DOI: <http://dx.doi.org/10.1017/S0956796897002967>.
- [34] P. Trinder, H.-W. Loidl, and R. Pointon. Parallel and Distributed Haskell. *J. Funct. Program.*, 12(5):469–510, July 2002. DOI: <http://dx.doi.org/10.1017/S0956796802004343>.
- [35] P. Trinder, K. Hammond, and H.-W. Loidl. *Encyclopedia of Parallel Computing*, chapter Parallel Functional Languages. Springer Verlag, Sept. 2011. ISBN 978-0387097657.