

# Process-Oriented Parallel Programming with an Application to Data-Intensive Computing

Edward Givelberg

July 22, 2014

## Abstract

We introduce *process-oriented programming* as a natural extension of object-oriented programming for parallel computing. It is based on the observation that every class of an object-oriented language can be instantiated as a process, accessible via a remote pointer. The introduction of process pointers requires no syntax extension, identifies processes with programming objects, and enables processes to exchange information simply by executing remote methods. Process-oriented programming is a high-level language alternative to multithreading, MPI and many other languages, environments and tools currently used for parallel computations. It implements natural *object-based parallelism* using only minimal syntax extension of existing languages, such as C++ and Python, and has therefore the potential to lead to widespread adoption of parallel programming. We implemented a prototype system for running processes using C++ with MPI and used it to compute a large three-dimensional Fourier transform on a computer cluster built of commodity hardware components. Three-dimensional Fourier transform is a prototype of a data-intensive application with a complex data-access pattern. The process-oriented code is only a few hundred lines long, and attains very high data throughput by achieving massive parallelism and maximizing hardware utilization.

## 1 Introduction

The first commercially available microprocessor CPUs appeared in the early 1970s. These were single-processor devices that operated with a clock rate of less than 1 MHz. Over the course of the following three decades increasingly faster and cheaper CPUs were built. This was achieved in large part by persistently increasing the clock rate. By 2004 the CPU clock rates reached the 3-4 GHz range and heat dissipation became a major problem. In order to continue improving the performance and the cost, microprocessor designers turned to parallel computing. Today, nearly all computing devices (servers, tablets, phones, etc.) are built using processors with multiple computing cores, whose operating frequency is less than 3.5 GHz. The industry move to parallel computing succeeded because practically every application contains tasks that can be executed in parallel. And yet, a decade later the vast majority of computer programs are still being written for execution on a single processor, and parallel computation is being realized primarily using threads, sequential processes that share memory.

Parallel programming is generally recognized as difficult, and has been a subject of extensive research (see [3], [9], [11] and references therein). The problem with threads is eloquently described in [10]. The author paints a bleak scenario:

“If [...] programmers make more intensive use of multithreading, the next generation of computers will become nearly unusable.”

In the scientific computing community parallel programs are typically written in Fortran and C with OpenMP [2] and MPI [1]. Dozens of high-level languages for parallel programming have also been developed, but presently none of them is widely used. Even the so-called *embarrassingly parallel* computations are not embarrassingly easy to implement.

In this paper we develop a new framework for parallel programming, which we call *process-oriented programming*. It is based on the fundamental observation that any class in an object-oriented language can be instantiated as a process, accessible via a remote pointer [6]. Such a process instantiates an object of the class and acts as a server to other processes, remotely executing the class interface methods on this object. The introduction of remote pointers enables a straightforward extension of object-oriented programming languages to process-oriented programming with hardly any syntax additions. We show that process-oriented programming is an efficient framework for parallel programming, and we propose it as a replacement for multithreading, MPI and many other languages, environments and tools currently used for parallel computations. We implemented a prototype system for running processes using C++ with MPI and investigated process-oriented programming in the context of data-intensive computing.

The rapid growth of generated and collected data in business and academia creates demand for increasingly complex and varied computations with very large data sets. The data sets are typically stored on hard drives, so the cost of accessing and moving small portions of the data set is high. Nevertheless, a large number of hard drives can be used in parallel to significantly reduce the amortized cost of data access. In [7], we argued that a *data-intensive computer* can be built, using widely available (“commodity”) hardware components, to solve general computational problems involving very large data sets.

The primary challenge in the construction of the data-intensive computer lies in software engineering. The software framework must balance programmer productivity with efficient code execution, i.e. big data applications with complex data access patterns must be realizable using a small amount of code, and this code must attain high data throughput, using massive parallelism. In this paper we aim to demonstrate that process-oriented programming is the right framework for the realization of the data-intensive computer.

We chose the computation of a large three-dimensional Fourier transform as the subject of our study primarily because it can be considered as a prototype of a difficult data-intensive problem. We show that using processes our application can be realized with only a few hundred lines of code, which are equivalent to approximately 15000 lines of C++ with MPI. We also show that even with the complex data access patterns required for the computation of the 3D Fourier transform, our code attains very high data throughput by achieving massive parallelism and maximizing hardware utilization.

In section 2 we describe a simple model of storage of a data set as a collection of data pages on multiple hard-drives. This model is used in the examples in section 3, where we introduce processes. The process-oriented implementation of the Fourier transform is described in section 4 and the efficiency of the computation is analyzed in section 5. We conclude with a discussion in section 6.

Our presentation uses C++, but can be easily applied to any object-oriented language. `size_t` is a large non-negative integer type used in C++ to represent the size of a data object in bytes.

## 2 The Data Set

### 2.1 Data Pages

We represent an  $N_1 \times N_2 \times N_3$  array of complex double precision numbers as a collection of  $NP_1 \times NP_2 \times NP_3$  pages, where each page is a small complex double precision array of size  $n_1 \times n_2 \times n_3$ . First, we define a `Page` class which stores `n` bytes of unstructured data:

```

class Page
{
public:
    Page(size_t n, unsigned char * data);
    ~Page();
protected:
    size_t n;
    unsigned char * data;
};

```

The `ArrayPage` class is derived from the `Page` class to handle three-dimensional complex double-precision array blocks:

```

class ArrayPage:
public Page
{
public:
    ArrayPage(
        int n1, int n2, int n3,
        double * data
    );
    // constructor that allocates data:
    ArrayPage(int n1, int n2, int n3);
    void transpose13();
    void transpose23();
    :
private:
    int n1, n2, n3;
}

```

The `ArrayPage` class may include functions for local computation with the array page data, such as the transpose functions that are important in the computation of the Fourier transform (see section 4.2). Throughout this paper the arrays have equal dimensions and the distinct variables `n1`, `n2` and `n3` are maintained only for the clarity of exposition.

## 2.2 Storage devices

We store data pages on hard drives, using a single large file for every available hard drive. The following `PageDevice` class controls the hard drive I/O.

```

class PageDevice
{
public:
    PageDevice(
        string filename,
        size_t NumberOfPages,
        size_t PageSize
    );
    ~PageDevice();
    void write(Page * p, size_t PageIndex);
    void read(Page * p, size_t PageIndex);
protected:
    string filename;
    size_t NumberOfPages;
    size_t PageSize;
private:

```

```

    int file_descriptor;
};

```

The implementation of this class creates a file `filename` of `NumberOfPages * PageSize` bytes. Pages of data are stored in the `PageDevice` object using a `PageIndex` address, where `PageIndex` ranges from 0 to `NumberOfPages - 1`. The `write` method copies a data page of size `PageSize` to the location with an offset `PageIndex * PageSize` from the beginning of the file `filename`. Similarly, the `read` method reads a page of data stored at a given integer address in the `PageDevice`. Linux direct I/O functions are used in the class implementation. For example, the Linux `open` function, used with the `O_DIRECT` flag, attempts to minimize cache effects of the I/O to and from the specified file.

For array pages we define the `ArrayPageDevice`, which extends `PageDevice`, as follows:

```

class ArrayPageDevice:
public PageDevice
{
public:
    ArrayPageDevice(
        string filename,
        size_t NumberOfPages,
        int nn1, int nn2, int nn3
    ):
        n1(nn1), n2(nn2), n3(nn3),
        PageDevice(
            filename,
            NumberOfPages,
            2 * n1 * n2 * n3 * sizeof(double)
        )
    {}
    void write_transpose13(Page * p, size_t PageIndex);
    void read_transpose13(Page * p, size_t PageIndex);
    void write_transpose23(Page * p, size_t PageIndex);
    void read_transpose23(Page * p, size_t PageIndex);
    :
private:
    int n1, int n2, int n3;
};

```

The implementation of the transpose methods is very simple. For example:

```

void ArrayPageDevice::
    read_transpose13(
        Page * p, size_t PageIndex
    )
{
    read(p, PageIndex);
    p->transpose13();
}

```

In addition to the transpose methods the `ArrayPageDevice` class may provides various other methods for computing with  $n1 \times n2 \times n3$  blocks of complex double precision data. Furthermore, in section 4.3 we extend `ArrayPageDevice` to include caching.

## 2.3 Page map

Since we use multiple hard drives to store a single large array object, we introduce the **PageMap** class to specify the storage layout of a given array. The **PageMap** translates logical array page indices into storage addresses. A storage address consists of an id of the server, corresponding to the hard drive where the data page is stored, and the page index, indicating the address of the page on that hard drive.

```
typedef
struct
{
    int server_id;
    size_t page_index;
}
address;

struct PageMap
{
    virtual address PageAddress(
        int i1, int i2, int i3
    );
};
```

In most applications the **PageMap** would be a simple, on-the-fly computable function, but it is also possible to implement it using an array of pre-computed values.

## 3 Processes

In this section we introduce processes into object-oriented programming languages and show that a complete framework for parallel programming with processes is obtained using only very small syntax extension.

### 3.1 Process creation, destruction and remote pointers

Programming objects can be naturally interpreted as processes. Upon creation, such a process instantiates the object and proceeds to act as a server to other processes, remotely executing the class interface methods on this object. For example, a program running on the computing node **machine0** can create a new **PageDevice** process on **machine1**, as follows:

```
size_t number_of_pages = 1024;
size_t page_size = 32 * 1024 * 1024;
char * remote_machine = "machine1";

PageDevice * storage
= new(remote_machine)
  PageDevice(
    "pagefile",
    number_of_pages,
    page_size
  );
```

It can then generate a page of data and store it on the remote **machine1** using the remote pointer **storage**:

```
Page * page = GenerateDataPage();
size_t PageAddress = 17;
storage->write(page, PageAddress);
```

The new `PageDevice` process on `machine1` acts as a server which listens on a communications port, accepts commands from other processes, executes them and sends results back to the clients. The client-server protocol is generated by the compiler from the class description. Remote pointer dereferencing triggers a sequence of events, that includes several client-server communications, data transfer and execution of code on both the client and the server machines.

Process semantics and remote pointers extend naturally to simple objects, as shown in the following example:

```
double * data
    = new(remote_machine) double[1024];
data[7] = 3.1415;
double x = data[2];
```

When this code is executed on `machine0`, a new process is created on `remote_machine`. This process allocates a block of 1024 doubles and deploys a server that communicates with the parent client running on `machine0`. The execution of `data[7] = 3.1415;` requires communication between the client and the server, including sending the numbers 7 and 3.1415 from the client to the server. Similarly, the execution of the following command leads to an assignment of the local variable `x` with a copy of the remote double `data[2]` obtained over the network using client-server communications. We emphasize that code execution is sequential: each instruction, and all communications associated with it, is completed before the following instruction is executed.

Finally, we remark that the notion of the class destructor in C++ extends naturally to process objects: destruction of a remote object causes termination of the remote process and completion of the corresponding client-server communications.

## 3.2 Parallel computation

The sequential programming model requires an execution of an instruction to complete before the next instruction is executed. We defined remote method execution to conform to this model, and therefore the calling process is kept idle until it is notified that the remote method has completed. In order to enable parallel computation we introduce the `start` keyword to indicate that the calling process may proceed with the execution of the next statement without waiting for the current statement to complete.

*Example:* A shared memory computation is constructed by providing access to the previously defined `data` block to several computing processes:

```
const int N = 64;
class ComputingProcess;
ComputingProcess * process_group[N];

for (int i = 0; i < N; i++)
    start process_group[i] = new(machine[i])
        ComputingProcess(data);
```

The `start` keyword can be used also for remote method and remote function calls, as is shown below.

An array of remote pointers defines a group of processes. It is easy to assign ids to the processes and to make them aware of the other processes in the group. This enables subsequent inter-process communication by remote method execution. Extending the example above, we assume that `ComputingProcess` is derived from the following `ProcessGroupMember` class.

```
class ProcessGroupMember
{
```

```

public:
    int ID() const
    { return id; }
    int NumberOfProcesses() const
    { return N; }
    ProcessGroupMember ** ProcessGroup() const
    { return group; }
protected:
    virtual void SetProcessGroup(
        int my_id,
        int my_N,
        ProcessGroupMember ** my_group
    );
private:
    int id;
    int N;
    ProcessGroupMember ** group;
};

```

The parameter `my_group` of the `SetProcessGroup` method is a remote pointer to an array of remote processes, so a shallow copy implementation of `SetProcessGroup` will result in redundant future communications. The following deep copy implementation of `SetProcessGroup`, which copies the entire remote array of remote pointers to a local array of remote pointers, is preferable:

```

void ProcessGroupMember::SetProcessGroup(
    int my_id,
    int my_N,
    ProcessGroupMember ** my_group
)
{
    id = my_id;
    N = my_N;
    group = new ProcessGroupMember * [N];

    for (int i = 0; i < N; i++)
        // remote copy:
        group[i] = my_group[i];
}

```

After instantiating the processes in `process_group` the master process can form a process group as follows:

```

for (int id = 0; id < N; id++)
    process_group[id]->SetProcessGroup(
        i, N, process_group
    );

```

Alternatively, the master process can execute the `SetProcessGroup` calls in parallel, but this requires synchronizing the processes at the end. A standard way to do this is using barrier functions. A process executing the barrier function call must wait until all other processes in the group execute a barrier function call before proceeding with the execution of the next statement:

```

for (int id = 0; id < N; id++)
    start process_group[id]->SetProcessGroup(
        i, N, process_group
    );

process_group->barrier();

```

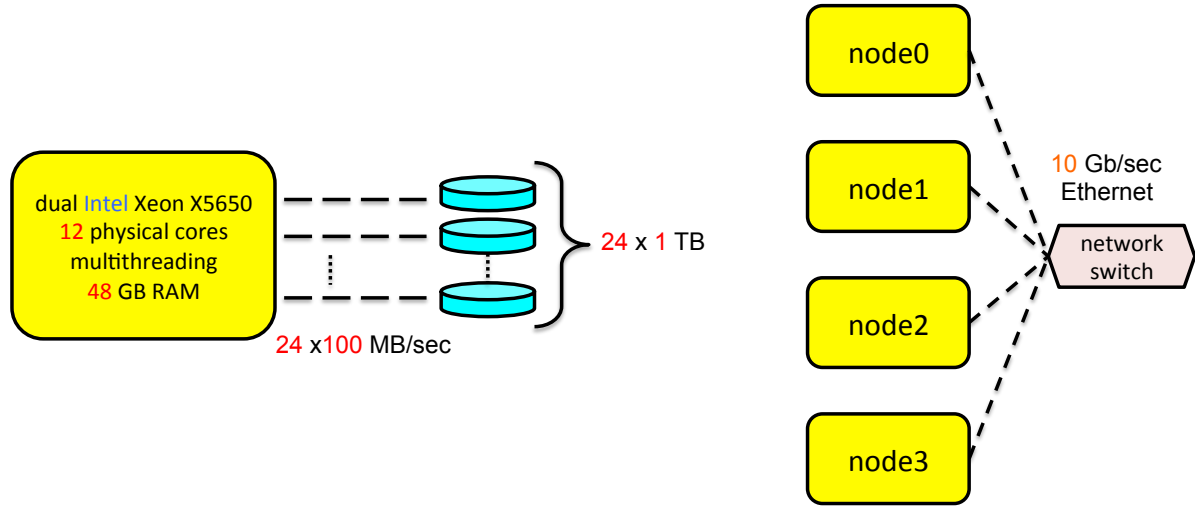


Figure 1: Cluster configuration for data intensive computing: The nodes are interconnected by a high speed network (10 Gb/sec Ethernet). Each node has a 12 core Intel Xeon processor with 48 GB of RAM and 24 attached 1 TB hard-drives. Each hard drive has an I/O throughput in the range of 100-150 MB/sec.

### 3.3 Persistent objects and processes

We view a large data set as a collection of *persistent processes*, which provide access to portions of the data set, as well as methods for computing with it. Persistent processes are objects that can be destroyed only by explicitly calling the destructor. The runtime system is responsible for storing process representation, and activating and de-activating processes, as needed. Processes can be accessed using a symbolic object address, similar to addresses used by the Data Access Protocol (DAP) [5], for example:

```
PageDevice * page_device =
    "http://data/set/PageDevice/34";
```

Because persistence is especially important for large data objects, in this paper we assume that all processes and all objects are persistent.

## 4 A program to compute the Fourier transform

In this section we show how to use processes in a data-intensive application on a cluster, such as the one shown in Figure 1. We chose a large three-dimensional Fourier transform as an example of a data-intensive application and we give an extensive description of the process-oriented code. We consider the situation where the array data set is represented by a group of processes, and the Fourier transform is a separate application whose processes interact with the array processes. In section 4.3 we show how the efficiency of the parallel computation can be improved using server-side caching.

### 4.1 The Array

We described a method of storing a large data set on multiple hard drives in section 2. We now define the `Array` class that is used in the Fourier transform computation.



The `Array` class describes a complex double precision three-dimensional array. It specifies the array domain, its storage layout and data access methods. We first define an auxiliary class to describe rectangular three-dimensional domains:

```
class Domain
{
public:
    Domain(
        int N11, int N12,
        int N21, int N22,
        int N31, int N32
    );
    :
};
```

The collection of hard drives attached to the cluster nodes can be turned into a distributed disk-based random access memory by launching a `PageDevice` process for every hard-drive on the cluster node to which this hard-drive is attached, but for computations with arrays we launch `ArrayPageDevice` processes that enable us to perform some of the array computations “close to the data”:

```
ArrayPageDevice ** page_server =
    new ArrayPageDevice * [NumberOfServers];

for (
    int i = 0; i < NumberOfServers; i ++
)
    page_server[i] = new(machine[i])
        ArrayPageDevice(
            filename[i],
            NumberOfPages,
            N1, N2, N3
        );
```

The `ArrayPageDevice` processes could be launched instead of the `PageDevice` processes, however in most applications it would be advantageous to launch them alongside existing `PageDevice` processes. In this case `ArrayPageDevice` must include a constructor that takes a pointer to an existing `PageDevice` process as a parameter. There would be essentially no communication overhead when the `ArrayPageDevice` process is launched on the same node as the corresponding `PageDevice` process. The `PageDevice::read` method, for example, will copy a page from a hard drive directly into a memory buffer that is accessible by the corresponding `ArrayPageDevice` process.

Array storage layout is determined by the `PageMap` object. The assignment of array pages to servers may affect the degree of parallelism that can be achieved in the computation. We use the circulant map, which assigns array page  $(i1, i2, i3)$  to the server  $(i1 + i2 + i3) \% \text{NumberOfServers}$ . To complete the specification of the `PageMap`, we assign the first available `PageIndex` address within the target `PageServer`.

The `Array` class provides methods for a client process to compute over a small array subdomain. The client may use a small (e.g. 4 GB) memory buffer to assemble the subdomain from array pages that reside within `page_server` processes, as determined by the array `pagemap`.

```
class Array    // complex double
{
public:
    Array(
        Domain * ArrayDomain,
```

```

    Domain * PageDomain,
    int NumberOfServers,
    ArrayPageDevice ** page_server,
    PageMap * pagemap
);
~Array();
void read(Domain * d, double * buffer);
void write(Domain * d, double * buffer);
void read_transpose13(
    Domain * d, double * buffer
);
void write_transpose13(
    Domain * d, double * buffer
);
void read_transpose23(
    Domain * d, double * buffer
);
void write_transpose23(
    Domain * d, double * buffer
);
private:
    Domain * ArrayDomain;
    Domain * PageDomain;
    int NumberOfServers;
    ArrayPageDevice * page_server;
    PageMap * pagemap;
};

```

An `Array` object is constructed by a single process, which can then pass the object pointer to any group of processes. Because an `Array` is a persistent object, it can also be accessed using a symbolic address, as described in section 3.3.

Transpose I/O methods are used to compute with long and narrow array subdomains that are not aligned with the third dimension. In section 4.2 we show an application of the transpose methods to the computation of the Fourier transform. The transpose of array pages can be computed either on the client or on the servers. For example, the implementation of `Array::read_transpose13` can assemble the data by executing `ArrayPageDevice::read_transpose13` on the appropriate `page_server` processes. Alternatively, it can execute `ArrayPageDevice::read` methods on these servers, followed by `ArrayPage::transpose13` on the received pages.

## 4.2 FFT Processes

We compute the three-dimensional Fourier transform using three separate functions, `fft1`, `fft2` and `fft3`, each performing one-dimensional Fourier transforms along the corresponding dimension. The functions `fft1` and `fft2` are similar to `fft3`, except that subdomains are read and written using the I/O transpose operations of the `Array` class. Having read and transposed the data into a local memory buffer, Fourier transforms are computed along the third dimension. The result is then transposed back and written to the array. We therefore restrict our description below to the `fft3` function.

The computation of the `fft3` function is performed in parallel using several `FFT3` client processes:

```

class FFT3:
    public ProcessGroupMember
{
public:

```

```

    FFT3(int sign, Array * a);
    ~FFT3() { delete buffer; }
    void ComputeTransform();
private:
    double * buffer;
    int sign;
    Array * a;
};

```

We divide the array into  $n$  slabs along the first dimension. The master process launches  $n$  FFT3 processes, assigning each process to a slab.

```

FFT3 ** fft = new FFT3 * [n];

for (int i = 0; i < n; i++)
    fft[i] = new(node[i]) FFT3(sign, a);

for (int i = 0; i < n; i++)
    fft[i]->SetProcessGroup(i, n, fft);

for (int i = 0; i < n; i++)
    start fft[i]->ComputeTransform();

```

Each process maintains a buffer that can hold a page line, where a page line is a collection of  $NP_3$  pages with page indices

$$\{(i_1, i_2, i_3) : i_3 = 0, 1, \dots, NP_3 - 1\}.$$

For complex double precision array of  $128^3$  pages, with each page consisting of  $128^3$  points, the page line buffer is 4 GB. Each process computes:

```

void FFT3::ComputeTransform()
{
    Domain * PageLine;
    for every PageLine in the slab
    {
        a->read(PageLine, buffer);
        call FFTW(sign, buffer);
        a->write(PageLine, buffer);
        ProcessGroup()->barrier();
    }
}

```

An `fft` process assembles a page line by reading the pages from appropriate page servers. For each page line one FFTW [4] function call computes a set of  $n_1 \times n_2$  one-dimensional complex double Fourier transforms of size  $N_3$  each. The result pages are then sent back to the page servers and stored on hard drives. Although the `fft` processes do not communicate with each other, they share a common network and a common pool of page servers. The barrier synchronization at the end of each iteration is not strictly necessary.

### 4.3 Parallel execution; caching

Each page line `read` and `write` operation consists of disk I/O and network transfer, with disk I/O being significantly more time consuming. We implement caching on the page server in order to carry out part of the disk I/O in parallel with the FFTW computation.

```

class ArrayDevice:
    public ArrayPageDevice

```

```

{
public:
    ArrayPageDevice(
        string filename,
        int NumberOfPages,
        int n1, int n2, int n3,
        int Nc
    ):
        ArrayPageDevice(
            filename,
            NumberOfPages,
            n1, n2, n3,
        ),
        NumberOfCachePages(Nc)
    {}
    void ReadIntoCache(size_t page_index);
    void read(Page * p, size_t page_index);
    :
    :
private:
    int NumberOfCachePages;
    Page ** cache;
};

```

An ArrayDevice server is configured with a cache that can hold a line of pages. In order to use ArrayDevice instead of ArrayPageDevice we add the following ReadIntoServerCache method to the Array class:

```

void Array::ReadIntoServerCache(
    Domain * domain
)
{
    for every page (i1, i2, i3) in domain
    {
        address a =
            pagemap->PageAddress(i1, i2, i3);
        int id = a.server_id;
        size_t i = a.page_index;
        start server[id]->ReadIntoCache(i);
    }
}

```

The ArrayDevice::read method executes ArrayPageDevice::read if the requested page is not found in its cache; otherwise it returns (sends over the network) the cached page.

The transform computation is now reorganized, so that instructions to read the next line of pages are sent to the page servers before the FFT of the current page line is started.

```

void FFT3::ComputeTransform()
{
    Domain * PageLine;
    for every PageLine in the slab
    {
        a->read(PageLine, buffer);
        Domain * NextPageLine = next page line;
        start a->ReadIntoServerCache(NextPageLine);
        call FFTW(sign, buffer);
        a->write(PageLine, buffer);
        ProcessGroup()->barrier();
    }
}

```

```
}

```

After the completion of the FFT the servers are instructed to write pages to hard drives. An execution of the `write` method will take place only after the server has completed `ArrayDevice::ReadIntoCache`. The efficiency of server-side caching depends on the relative timing of the FFT computation and the page I/O. If necessary, `write` and `WriteFromCache` methods can be implemented in the `ArrayDevice` class analogously.

## 5 Computation of the Fourier Transform

In this section we describe the computation of a large (64 TB) data-intensive Fourier transform on a small (8 nodes, 96 cores) cluster. We describe our prototype implementation of the process framework and analyze the performance of the Fourier transform computation.

### 5.1 Software Implementation of Processes

The Fourier transform application was developed by completing the process-oriented code which is sketched in sections 2 and 4, translating it into C++ with MPI and linking it against an auxiliary library of tools for implementation of basic process functionality. We describe the procedure for `ArrayDevice` processes. In order to instantiate `ArrayDevice` processes, we created an `ArrayDevice_server` class and a C++ program file `ArrayDevice_process.cpp` containing the following code:

```
// start the server and disconnect from parent process
int main(int argc, char **argv)
{
    MPI_Init(&argc, &argv);
    MPI_Comm parent;
    MPI_Comm_get_parent(&parent);
    ArrayDevice_server * s =
        new ArrayDevice_server();
    MPI_Comm_disconnect(&parent);
    MPI_Finalize();
}
```

In addition, we implemented the function

```
void LaunchProcess(
    const string & process_file_name,
    const string & machine_address,
    string & server_address
);
```

`LaunchProcess` uses `MPI_Comm_spawn` to create an MPI process by running the executable `process_file_name` on the remote machine specified by `machine_address`. The launched process starts a server which uses `MPI_Open_port` to open a port and return the port address in the `server_address` output parameter. The server process then disconnects from the launching process.

In order to implement remote method execution we also created an `ArrayDevice_client` class. The `ArrayDevice_client` and the `ArrayDevice_server` classes are automatically constructed from the `ArrayDevice` class. For the purposes of this example we use a simple name mangling scheme to first generate the following `ArrayDevice_interface` class:

```
class ArrayDevice_interface
{

```

```

:
:
public:
    virtual void ArrayDevice_create(    // constructor
        string filename,
        size_t numberofpages,
        size_t pagesize
    );
    virtual void ArrayDevice_destroy(); // destructor
    virtual void ArrayDevice_write(Page * p, size_t page_index);
    virtual void ArrayDevice_ReadIntoCache(
        int n,
        size_t * page_index
    );

    // commands used in client-server protocol
    enum ArrayDevice_command
    {
        create_CMD = 0,
        destroy_CMD = 1,
        write_CMD = 2,
        ReadIntoCache_CMD = 3
    };
    :
};

```

The `ArrayDevice_interface` class contains meta-information obtained from the `ArrayDevice` class. Both the `ArrayDevice_client` and the `ArrayDevice_server` classes are derived from `ArrayDevice_interface`. `ArrayDevice_client` is also derived from a general `ProcessClient` class, and similarly `ArrayDevice_server` is derived from `ProcessServer`. Jointly, the `ArrayDevice_client`, `ArrayDevice_server` pair implement remote procedure calls (RPC) for `ArrayDevice`. The client-server communications protocol for `ArrayDevice` uses the `ArrayDevice_interface` class and the general-purpose functionality implemented in `ProcessClient` and `ProcessServer`. The client-server implementation uses a simple object serialization library to send method parameters and results over the network. This serialization library is also used to implement a rudimentary file-based object persistence mechanism. Information is sent over the network using an MPI-based communications software layer.

The translation procedure we implemented for the processes of the Fourier transform application can be extended and incorporated into a C++ compiler. It converts several hundred lines of process-oriented code into a C++/MPI application for computing Fourier transform, which is approximately 15000 lines long. A very small subset of MPI functions is used. The following is the almost complete list:

- `MPI_Send`, `MPI_Recv` – for inter-process communication,
- `MPI_barrier` – for process synchronization,
- `MPI_Open_port`, `MPI_Close_port`, `MPI_Comm_accept` – to implement client-server functionality,
- `MPI_Comm_spawn`, `MPI_Info_create`, `MPI_Info_set` – to spawn processes on specified target machines,
- `MPI_Comm_connect`, `MPI_Comm_disconnect`, `MPI_Comm_get_parent` – to manage process connections.

Both MPICH and Intel MPI were used in our computation, and with both implementations we encountered problems with some MPI functions. The most significant bugs were found in `MPI_Comm_spawn` and

`MPI_Comm_disconnect.`

## 5.2 Fourier Transform Computation

We used the process prototype implementation to carry out a computation of the Fourier transform of a  $(16K)^3$ -point array of complex double precision numbers. (We use the notation  $1K = 1024$ ,  $16K = 16384 = 128 \times 128$ .) The total size of the array is 64 TB. The computations were carried out on a cluster of 8 nodes, interconnected with a 10 Gb/sec Ethernet, similar to the cluster depicted in Figure 1. Each computational node has an Intel<sup>®</sup> Xeon<sup>®</sup> X5650 12 core CPU with hyperthreading, rated at 2.67GHz, and 24 attached 1 TB hard drives. The hard drives are manufactured by Samsung, model SpinPoint F3 HD103SJ, with manufacturer listed latency of 4.14 ms and average seek time of 8.9 ms. Benchmark hard drive read and write throughput is reported at over 100 MB/sec.

The  $(16K)^3$ -point array was broken up into  $128^3$  `ArrayPage` pages of  $128 \times 128 \times 128$  points each. The resulting page size is 32 MB. The choice of the page size is constrained by the latency and seek time of the system's hard drives: the smaller the page size, the lower the overall disk I/O throughput. We measured a typical page read/write time in the range of 0.25-0.35 sec.

We used 4 of the 8 available nodes to store the `Array` object, creating 24 `ArrayDevice` processes on each node, one process for each available hard drive. Because `ArrayDevice` processes are primarily dedicated to disk I/O it is possible to run 24 processes on a 12-core node. The process framework makes it easy to shift computation closer to data by extending the `ArrayDevice` class, and in such case relatively more powerful CPUs may be needed to run the server processes.

We used the other 4 nodes to run 16 processes of the Fourier transform application, 4 processes per node. Each of the 16 Fourier transform processes was assigned an array slab of  $8 \times 128 \times 128$  pages. The process computes the transform of its slab line by line in  $8 \times 128 = 1024$  iterations. Although the 16 processes are independent of each other, they compete among themselves for service from 96 page servers.

The wall clock time for a single iteration of `ComputeTransform` (see section 4.3) generally ranged from 68 to 78 seconds, with the average of approximately 73 seconds. The total speed of data processing (including reading, computing and writing the data) has been therefore close to 1 GB/sec. In the next section we analyze the performance in detail and indicate a number of ways to substantially improve it.

## 5.3 Performance analysis

The computation of the Fourier transform was completed in the course of several long (over 10 hours) continuous runs. Our implementation of persistence mechanism for processes made it possible to stop and restart the computation multiple times. The primary reason for long runs was to test the stability and robustness of our implementation. We instrumented the code to measure the utilization of the system's components: the network, the hard drives and the processors. Because the results did not vary substantially over the course of the computation, we present a detailed analysis of a typical iteration of `ComputeTransform`.

The synchronization of the processes at the end of each iteration is not strictly necessary, but we found that it did not significantly affect performance, and made the code easier to analyze. On the other hand, we found that introducing additional barriers within the iteration would slow down the computation significantly. We now present detailed measurements of the component phases of the iteration.

At the beginning of every iteration each process reads a page line consisting of 128 pages, which are evenly distributed among the 96 servers by the circulant map (see section 4.1). Except for the first iteration, the required pages have already been read from the hard drive and placed in the memory of the corresponding servers. Each server has either 21 or 22 pages to send to the clients. In total, 64 GB of data is sent from

the servers to the clients. Accordingly, the combined size of the server caches in each of the server nodes is slightly more than 16 GB. We timed the parallel reading of page lines by the 16 client processes in a typical time step:

```
15: Array::read 11520-11647 x 2432-2559 x 0-16383:: 4 GB, 17.3638 sec, 235.894 MB/sec
12: Array::read 9216-9343 x 2432-2559 x 0-16383:: 4 GB, 18.6695 sec, 219.395 MB/sec
0: Array::read 0-127 x 2432-2559 x 0-16383:: 4 GB, 21.6104 sec, 189.538 MB/sec
1: Array::read 768-895 x 2432-2559 x 0-16383:: 4 GB, 22.2821 sec, 183.824 MB/sec
11: Array::read 8448-8575 x 2432-2559 x 0-16383:: 4 GB, 22.4682 sec, 182.302 MB/sec
8: Array::read 6144-6271 x 2432-2559 x 0-16383:: 4 GB, 22.7035 sec, 180.412 MB/sec
14: Array::read 10752-10879 x 2432-2559 x 0-16383:: 4 GB, 22.9213 sec, 178.698 MB/sec
5: Array::read 3840-3967 x 2432-2559 x 0-16383:: 4 GB, 27.1039 sec, 151.122 MB/sec
6: Array::read 4608-4735 x 2432-2559 x 0-16383:: 4 GB, 27.0914 sec, 151.192 MB/sec
2: Array::read 1536-1663 x 2432-2559 x 0-16383:: 4 GB, 27.2604 sec, 150.255 MB/sec
4: Array::read 3072-3199 x 2432-2559 x 0-16383:: 4 GB, 28.1986 sec, 145.256 MB/sec
3: Array::read 2304-2431 x 2432-2559 x 0-16383:: 4 GB, 28.3351 sec, 144.556 MB/sec
7: Array::read 5376-5503 x 2432-2559 x 0-16383:: 4 GB, 28.9623 sec, 141.425 MB/sec
10: Array::read 7680-7807 x 2432-2559 x 0-16383:: 4 GB, 29.027 sec, 141.11 MB/sec
9: Array::read 6912-7039 x 2432-2559 x 0-16383:: 4 GB, 29.5628 sec, 138.553 MB/sec
13: Array::read 9984-10111 x 2432-2559 x 0-16383:: 4 GB, 30.1746 sec, 135.743 MB/sec
```

The first number in each row is the client process id, followed by the description of the domain, domain size, the time it took to read it and the corresponding throughput. The fastest process completes the execution of the `Array::read` method (including sending the command and the parameter to page servers) in approximately 17.5 seconds, the slowest – in about 30 seconds. The faster processes proceed to start the execution of `Array::ReadIntoServerCache` immediately after completion of `Array::read`. The typical aggregate throughput during the parallel execution of `Array::read` is therefore significantly larger than 2 GB/sec. The maximal possible throughput to the four nodes computing the transform is approximately 4 GB/sec.

In the `Array::ReadIntoServerCache` phase of the computation 16 clients send small messages to 96 page servers with instructions to read a total of  $16 \times 128$  pages. These commands are queued for execution in the servers, so that the clients do not wait for the completion of the execution. The time measurements of the parallel execution of `start Array::ReadIntoServerCache` by the 16 client processes in a typical time step were:

```
14: Array::ReadIntoServerCache 10752-10879 x 2560-2687 x 0-16383:: 9.36148 sec
13: Array::ReadIntoServerCache 9984-10111 x 2560-2687 x 0-16383:: 2.10771 sec
10: Array::ReadIntoServerCache 7680-7807 x 2560-2687 x 0-16383:: 3.25717 sec
12: Array::ReadIntoServerCache 9216-9343 x 2560-2687 x 0-16383:: 9.5796 sec
8: Array::ReadIntoServerCache 6144-6271 x 2560-2687 x 0-16383:: 9.81484 sec
11: Array::ReadIntoServerCache 8448-8575 x 2560-2687 x 0-16383:: 13.6131 sec
9: Array::ReadIntoServerCache 6912-7039 x 2560-2687 x 0-16383:: 2.67452 sec
5: Array::ReadIntoServerCache 3840-3967 x 2560-2687 x 0-16383:: 5.16807 sec
2: Array::ReadIntoServerCache 1536-1663 x 2560-2687 x 0-16383:: 5.0223 sec
4: Array::ReadIntoServerCache 3072-3199 x 2560-2687 x 0-16383:: 4.08604 sec
0: Array::ReadIntoServerCache 0-127 x 2560-2687 x 0-16383:: 10.6741 sec
3: Array::ReadIntoServerCache 2304-2431 x 2560-2687 x 0-16383:: 3.94899 sec
1: Array::ReadIntoServerCache 768-895 x 2560-2687 x 0-16383:: 9.98827 sec
6: Array::ReadIntoServerCache 4608-4735 x 2560-2687 x 0-16383:: 5.16787 sec
15: Array::ReadIntoServerCache 11520-11647 x 2560-2687 x 0-16383:: 14.9185 sec
7: Array::ReadIntoServerCache 5376-5503 x 2560-2687 x 0-16383:: 3.27696 sec
```

These results are significantly worse than expected. Our implementation of processes repeatedly establishes and breaks client-server connections. We found the `MPI_Comm_connect` function to be very fast, but with increasing number of client-server connections, it sporadically performed hundreds of times slower than usual. An implementation of caching connections is likely to reduce the total time for this phase of the computation to about 3 seconds.

We timed the execution of the FFTW function call by every processor.



```

14: fftw 10752-10879 x 2432-2559 x 0-16383 5.53611 sec
0: fftw 0-127 x 2432-2559 x 0-16383 5.86336 sec
6: fftw 4608-4735 x 2432-2559 x 0-16383 5.91158 sec
4: fftw 3072-3199 x 2432-2559 x 0-16383 5.99485 sec
8: fftw 6144-6271 x 2432-2559 x 0-16383 6.01364 sec
12: fftw 9216-9343 x 2432-2559 x 0-16383 6.03876 sec
10: fftw 7680-7807 x 2432-2559 x 0-16383 6.06033 sec
2: fftw 1536-1663 x 2432-2559 x 0-16383 6.07081 sec
11: fftw 8448-8575 x 2432-2559 x 0-16383 6.38444 sec
7: fftw 5376-5503 x 2432-2559 x 0-16383 6.38421 sec
15: fftw 11520-11647 x 2432-2559 x 0-16383 6.42032 sec
3: fftw 2304-2431 x 2432-2559 x 0-16383 6.42975 sec
9: fftw 6912-7039 x 2432-2559 x 0-16383 14.001 sec
13: fftw 9984-10111 x 2432-2559 x 0-16383 14.5631 sec
5: fftw 3840-3967 x 2432-2559 x 0-16383 14.5582 sec
1: fftw 768-895 x 2432-2559 x 0-16383 14.878 sec

```

The last 4 processes (9,13,5 and 1) ran on the same node. Throughout our computation these 4 processes executed the FFTW library function call significantly slower than the processes running on other nodes. Additional investigation of the configuration of this node is needed to speed up the computation.

The reading of the pages on the server side is done concurrently with the FFTW computation. We include the measurements for a few of the 96 servers:

```

:
34: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.59864 sec, 125.745 MB/sec
25: ArrayDevice::ReadIntoCache: 21 pages, 672 MB, 5.614 sec, 119.701 MB/sec
57: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.58319 sec, 126.093 MB/sec
52: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.59719 sec, 125.777 MB/sec
40: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.61367 sec, 125.408 MB/sec
73: ArrayDevice::ReadIntoCache: 21 pages, 672 MB, 5.57943 sec, 120.442 MB/sec
21: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.64487 sec, 124.715 MB/sec
33: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.63134 sec, 125.015 MB/sec
89: ArrayDevice::ReadIntoCache: 21 pages, 672 MB, 5.56771 sec, 120.696 MB/sec
46: ArrayDevice::ReadIntoCache: 22 pages, 704 MB, 5.62876 sec, 125.072 MB/sec
:

```

The reading throughput is close to the maximal throughput for this type of hard drives. For every server the page reading commands are scheduled before the page writing commands of the last phase of the iteration. The writing of the pages will therefore start only after the completion of the page read commands. The clients write pages in parallel, with the typical timing as follows:

```

8: Array::write 6144-6271 x 2432-2559 x 0-16383:: 4 GB, 22.3279 sec, 183.448 MB/sec
4: Array::write 3072-3199 x 2432-2559 x 0-16383:: 4 GB, 22.784 sec, 179.776 MB/sec
11: Array::write 8448-8575 x 2432-2559 x 0-16383:: 4 GB, 22.7165 sec, 180.31 MB/sec
12: Array::write 9216-9343 x 2432-2559 x 0-16383:: 4 GB, 23.083 sec, 177.447 MB/sec
15: Array::write 11520-11647 x 2432-2559 x 0-16383:: 4 GB, 23.0278 sec, 177.872 MB/sec
7: Array::write 5376-5503 x 2432-2559 x 0-16383:: 4 GB, 23.2437 sec, 176.22 MB/sec
0: Array::write 0-127 x 2432-2559 x 0-16383:: 4 GB, 23.7938 sec, 172.145 MB/sec
14: Array::write 10752-10879 x 2432-2559 x 0-16383:: 4 GB, 24.5307 sec, 166.975 MB/sec
10: Array::write 7680-7807 x 2432-2559 x 0-16383:: 4 GB, 24.2152 sec, 169.15 MB/sec
6: Array::write 4608-4735 x 2432-2559 x 0-16383:: 4 GB, 24.6648 sec, 166.066 MB/sec
2: Array::write 1536-1663 x 2432-2559 x 0-16383:: 4 GB, 24.7291 sec, 165.635 MB/sec
3: Array::write 2304-2431 x 2432-2559 x 0-16383:: 4 GB, 25.3921 sec, 161.31 MB/sec
9: Array::write 6912-7039 x 2432-2559 x 0-16383:: 4 GB, 21.8906 sec, 187.113 MB/sec
13: Array::write 9984-10111 x 2432-2559 x 0-16383:: 4 GB, 22.0268 sec, 185.956 MB/sec
1: Array::write 768-895 x 2432-2559 x 0-16383:: 4 GB, 21.7936 sec, 187.945 MB/sec
5: Array::write 3840-3967 x 2432-2559 x 0-16383:: 4 GB, 22.2789 sec, 183.851 MB/sec

```

The results for page writing are fairly uniform, with the total time for each process between 22 and 25.5 seconds. The aggregate throughput for this phase is therefore over 2.5 GB/sec. We did not implement explicit caching for writing pages. The implementation of `Array::write` is analogous to the implementation of `Array::ReadIntoServerCache` (see section 4.3):

```

void Array::write(Domain * domain)
{
    for every page in domain
        start write page to the appropriate server
}

```

There is a limited caching effect as a result of the `start` command: having transmitted the page to the server, the client disconnects and proceeds to transmit pages to other servers, while the server starts writing the page to disk only after the client has disconnected. We found that writing pages to hard drive with the `O_DIRECT` flag is about twice as fast as reading, presumably because of buffering. The typical throughput measured was between 230 and 240 MB/sec:

```

:
22: ArrayDevice::write: 1 page, 32 MB, 0.136958 sec, 233.648 MB/sec
78: ArrayDevice::write: 1 page, 32 MB, 0.135615 sec, 235.962 MB/sec
73: ArrayDevice::write: 1 page, 32 MB, 0.13452 sec, 237.883 MB/sec
64: ArrayDevice::write: 1 page, 32 MB, 0.136484 sec, 234.46 MB/sec
8: ArrayDevice::write: 1 page, 32 MB, 0.137042 sec, 233.505 MB/sec
54: ArrayDevice::write: 1 page, 32 MB, 0.136533 sec, 234.376 MB/sec
87: ArrayDevice::write: 1 page, 32 MB, 0.136293 sec, 234.788 MB/sec
:

```

The conclusion of our performance analysis is that the presented computation could be sped up by 25% or more, but greater benefits can be derived from a more balanced hardware configuration. The aggregate throughput of the 24 hard drives of a cluster node is about 3 GB/sec, about 3 times the capacity of the incoming network connection. Furthermore, the FFTW computation takes only 10-20% of the total iteration time. It appears that a 2-4-fold increase in the network capacity of the present cluster is likely to result in a more balanced system with better hardware utilization, and a total runtime of under 20 seconds per iteration.

## 6 Discussion and Conclusions

In this paper we introduced process-oriented programming as a natural extension of object-oriented programming for parallel computing. We implemented a prototype of the process framework and carried out a data-intensive computation. We have shown that a complex and efficient application can be built using only a few hundred lines of process-oriented code, which is equivalent to many thousands of lines of object-oriented code with MPI. The process-oriented code in this paper is an extension of C++, but processes can be introduced into any object-oriented language. The syntax extension is minimal. In C++, for example, it amounts to adding a parameter to the `new` operator and introducing the keyword `start`. Combined with the fact that creating a process can be thought of as simply placing an object on a remote machine, it suggests that a lot of existing code can be easily modified to run in parallel. Potentially the most important impact of the process-oriented extension of languages, such as C++ and Python, is a widespread adoption of parallel programming, as application developers realize the ability to easily create processes instead of using thread libraries and to place different objects on different CPU cores.

The process-oriented programming model is based on a simple hardware abstraction: the computer consists of a collection of processors, interconnected by a network, where each processor is capable of running multiple processes. The run-time system is responsible for mapping the abstract model onto a concrete hardware system, and must provide the programmer with system functions describing the state of the hardware. The hardware abstraction of the process-oriented model makes it possible to create portable parallel applications and applications that run in the cloud.

Processes are accessible by remote pointers. Syntactically, executing a method on a remote process is not different from method execution on an object. Any class of an object-oriented language can be interpreted as a process, but even more importantly, in the process-oriented framework only processes that are class instances are allowed. We argue that *object-based parallelism* is a high level abstraction, which is naturally suitable for reasoning about parallelism. Although shared memory and message passing can be realized in a process-oriented language, these are lower implementation-level models. Process inheritance is a powerful aspect of object-based parallelism, as it enables the definition of new processes in terms of previously defined processes. In combination with process pointers, it gives the programmer the flexibility to adapt the computation to the hardware by adding simple methods to class definitions (see the comments at the end of section 4.1 about the computation of array page transpose).

A process-oriented program, like a typical sequential program starts with a single main process. The main process may launch new processes on remote machines, as easily as it can create objects. In contrast with MPI, processes are explicitly managed by the programmer. Process launching is part of the program, and is not determined by the runtime command line parameters. Using process pointers and language data structures, the programmer can form groups of processes, assign process ids and perform tasks that in MPI would require using communicators.

Processes exchange information by executing remote methods, rather than via shared memory or message passing. We'd like to use the analogy that writing programs with message passing today is like writing programs with GOTOs fifty years ago: it is easy to write intractable code. And just like GOTO statements are used in assembly languages, message passing is a low-level language construct underlying remote method execution in process-oriented programming.

We introduced the `start` keyword to enable parallel execution of remote methods. In order to use the keyword the programmer must decide whether there is a need to wait for the remote task to complete before proceeding with the computation. In general, this decision is easy and intuitive, but keeping track of task dependencies is not. Each process executes only one method at a time and remote method execution requests are queued. The programmer must keep in mind the state of the execution queue for every process. This is possible only for very simple scenarios. We used barriers to synchronize processes. Barrier synchronization helps the programmer to keep track of the execution queues of the processes, but it may reduce the parallelism of the computation.

We view a large data object as a collection of persistent processes. For a large data object a negligible amount of additional storage space is needed to store serialized processes alongside the data. Process persistence is needed to enable stopping and restarting a computation, and to make a data set accessible to several simultaneous applications. It is also needed to develop basic mechanisms for fault tolerance. In this paper we showed that the process-oriented view of a large data object is very powerful: using only a small amount of code the programmer can copy and reformat very large data objects, and even carry out complex operations, such as the Fourier transform. Yet, software users and application developers tend to see a data set as consisting of “just data”, and being independent of a specific programming language. We stop short of suggesting a solution for this problem, but in this context it is worth recalling the CORBA standard [8].

The introduction of process-oriented programming was motivated by our research in data-intensive computing. The data-intensive Fourier transform computation was carried out on a small cluster of 8 nodes (96 cores). Our measurements indicate that Petascale data-intensive computations can be efficiently carried out on a larger cluster, with more nodes and significantly increased network bandwidth. Such a cluster can serve as a general-purpose data-intensive computer, whose operating system and applications are developed as process-oriented programs.

The introduction of process-oriented programming in this paper is far from complete. It is merely the

first step of an extensive research program. We tested a substantial subset of the process framework in a prototype implementation. A full-fledged implementation must include a compiler and a run-time system that substantially expand the basic prototype.

## 7 Acknowledgements

I am grateful to J. J. Bunn for discussions that significantly improved the presentation of the material.

## References

- [1] MPI: A Message-Passing Interface Standard, version 3.0. <http://www.mpi-forum.org/>.
- [2] OpenMP 4.0 Specifications. <http://openmp.org/wp/openmp-specifications/>.
- [3] J. Diaz, C. Munoz-Caro, and A. Nino. A survey of parallel programming models and tools in the multi and many-core era. *Parallel and Distributed Systems, IEEE Transactions on*, 23(8):1369–1386, 2012.
- [4] M. Frigo and S. G. Johnson. The design and implementation of fftw3. *Proceedings of the IEEE*, 93(2):216–231, 2005.
- [5] J. Gallagher, N. Potter, T. Sgouros, S. Hankin, and G. Flierl. The data access protocol-DAP 2.0, 2004.
- [6] E. Givelberg. Object-oriented parallel programming. April 2014. arXiv:1404.4666 [cs.PL].
- [7] E. Givelberg, A. Szalay, K. Kanov, and R. Burns. An architecture for a data-intensive computer. In *Proceedings of the first international workshop on Network-aware data management*, pages 57–64. ACM, 2011.
- [8] M. Henning. The rise and fall of CORBA. *Queue*, 4(5):28–34, June 2006.
- [9] M. Herlihy and N. Shavit. *The Art of Multiprocessor Programming, Revised Reprint*. Elsevier, 2012.
- [10] E. A. Lee. The problem with threads. *Computer*, 39(5):33–42, 2006.
- [11] T. G. Mattson, B. A. Sanders, and B. L. Massingill. *Patterns for parallel programming*. Pearson Education, 2004.