

# ***Bulk:* a Modern C++ BSP Interface**

---

Jan-Willem Buurlage (CWI)

*MasterMath*: Parallel Computing (2018)

- For high-performance computing on distributed-memory systems, BSP is still a (if not *the*) leading model.
- In the last 10 years or so, it has grown again in popularity. It has also found widespread use in industry (MapReduce / Pregel).
- BSP programming usually done using MPI or the various Apache projects (Hama, Giraph, Hadoop).

- *Standard example*: word count. The **map** takes a (file, content) pair, and emits (word, 1) pairs for each word in the content. The **reduce** function sums over all mapped pairs with the same word.
- The **map** and **reduce** are performed in parallel, and are both followed by communication and a bulk synchronization, which means  $\text{MapReduce} \subset \text{BSP}$ !<sup>1</sup>

<sup>1</sup>MapReduce: Simplified Data Processing on Large Clusters, Jeffrey Dean and Sanjay Ghemawat (2004)

BSP for graph processing, used by Google<sup>2</sup> and Facebook<sup>3</sup>.

*The high-level organization of Pregel programs is inspired by Valiant's Bulk Synchronous Parallel model. Pregel computations consist of a sequence of iterations, called supersteps ... It can read messages sent to  $V$  in superstep  $S$  1, send messages to other vertices that will be received at superstep  $S + 1$  ...*

<sup>2</sup>Pregel: A System for Large-Scale Graph Processing – Malewicz et al. (2010)

<sup>3</sup>One Trillion Edges: Graph Processing at Facebook-Scale - Avery Ching et al (2015)

- These frameworks are good for **big data analytics**, not for high-performance scientific computing.
- $\implies$  Most scientific software still built on top of MPI.
- Modern programming languages have **novel features** (safety, abstractions) which can aid parallel programming.

- There are mature implementations of BSPlib for shared and distributed-memory systems<sup>4</sup>.
- Many *Big Data* frameworks are based on (restricted) BSP programming, such as MapReduce (Apache Hadoop), Pregel (Apache Giraph) and so on.
- BSP interfaces that are not based on BSPlib include BSML and Apache Hama.

<sup>4</sup>e.g. Multicore BSP (for C) by Albert Jan Yzelman and BSPonMPI by Wijnand Suijlen

```
#include <bsp.h>

int main() {
    bsp_begin(bsp_nprocs());
    int s = bsp_pid();
    int p = bsp_nprocs();
    printf("Hello World from processor %d / %d", s, p);
    bsp_end();

    return 0;
}
```

## BSPLib: Registering and using variables

```
int x = 0;
bsp_push_reg(&x, sizeof(int));
bsp_sync();

int b = 3;
bsp_put((s + 1) % p, &b, &x, 0, sizeof(int));

int c = 0;
bsp_get((s + 1) % p, &x, 0, &c, sizeof(int));

bsp_pop_reg(&x);
bsp_sync();
```



## BSPlib: Sending messages

```
int tagsize = sizeof(int);  
bsp_set_tagsize(&tagsize);  
bsp_sync();  
  
int tag = 1;  
int payload = 42 + s;  
bsp_send((s + 1) % p, &tag, &payload, sizeof(int));  
bsp_sync();
```

## BSPlib: Receiving messages

```
int packets = 0;
int accum_bytes = 0;
bsp_qsize(&packets, &accum_bytes);

int payload_in = 0;
int payload_size = 0;
int tag_in = 0;
for (int i = 0; i < packets; ++i) {
    bsp_get_tag(&payload_size, &tag_in);
    bsp_move(&payload_in, sizeof(int));
    printf("payload: %i, tag: %i", payload_in, tag_in);
}
```

# A modern BSP interface

- Modern programming languages focus on **safety** and **zero-cost abstractions** to increase programmer productivity, without sacrificing **performance**.
- A modern BSP interface should also have this focus. We want **correct, safe and clear** implementations of BSP programs without taking a performance hit.
- *Modern C++* has a large user base, is widely supported, with a good set of features and (support for) abstractions.

# Bulk: A modern BSP interface

- Bulk is a modern BSPlib replacement.
- Focuses on memory safety, portability, code reuse, and ease of implementation of BSP algorithms.
- Flexible backend architecture. Bulk programs target shared, distributed, or hybrid memory systems.
- Support for various *algorithmic skeletons*, and utility features for logging, benchmarking, and reporting.

## Bulk: Basics

- A BSP computer is captured in an environment (e.g. an MPI cluster, a multi-core processor or a many-core coprocessor).
- In an environment, an SPMD block can be spawned.
- The processors running this block form a parallel world, that can be used to communicate, and for obtaining information about the local process.

```
bulk::backend::environment env;  
env.spawn(env.available_processors(), spmd);
```

```
void spmd(bulk::world& world) {  
    world.log("Hello world from %d / %d\n",  
              world.rank(),  
              world.active_processors());  
}
```

## Bulk: Distributed variables (I)

- Registering and deregistering (`bsp_push_reg`) is replaced by *distributed variables*.

```
auto x = bulk::var<int>(world);  
auto y = x(t).get();  
x(t) = value;
```

- These variables are `var` objects. Their value is generally different on each processor.
- References to remote values are captured in `image` objects, and can be used for reading and writing.

## Bulk: Distributed variables (II)

```
auto x = bulk::var<int>(world);  
auto t = world.next_rank();  
x(t) = 2 * world.rank();  
world.sync();  
// x now equals two times the previous ID  
  
auto b = x(t).get();  
world.sync();  
// b.value() now equals two times the local ID
```

## Bulk: Coarrays (I)

- Distributed variables work well for communicating single values.
- For communication based on (sub)arrays we have coarray objects, loosely inspired by Coarray Fortran.

```
auto xs = bulk::coarray<int>(world, 10);  
xs(t)[5] = 3;  
auto y = xs(t)[5].get();
```

- Images to remote subarrays of a coarray `xs`, are obtained as for variables by `xs(t)`, and can be used to access the remote array.



## Bulk: Coarrays (II)

```
auto xs = bulk::coarray<int>(world, 4);
auto t = world.next_rank();
xs[0] = 1;
xs(t)[1] = 2 + world.rank();
xs(t)[{2, 4}] = {123, 321};
world.sync();
// xs is now [1, 2 + world.prev_rank(), 123, 321]
```

## Bulk: Message passing queues (I)

- One-sided mailbox communication using message passing, which in Bulk is carried out using a queue. Greatly simplified compared to previous BSP interfaces, without losing power or flexibility.

```
// single integer, and zero or more reals  
auto q1 = bulk::queue<int, float[]>(world);  
// sending matrix nonzeros around (i, j, a_ij)  
auto q2 = bulk::queue<int, int, float>(world);
```

- Message structure is defined in the construction of a queue: optionally attach tags, or define your own record structure.

## BSPlib: Sending messages

```
int tagsize = sizeof(int);  
bsp_set_tagsize(&tagsize);  
bsp_sync();  
  
int tag = 1;  
int payload = 42 + s;  
bsp_send((s + 1) % p, &tag, &payload, sizeof(int));  
bsp_sync();
```

## Bulk: Sending messages

```
auto q = bulk::queue<int, int>(world);  
q(world.next_rank()).send(1, 42 + s);  
world.sync();
```

## BSPlib: Receiving messages

```
int packets = 0;
int accum_bytes = 0;
bsp_qsize(&packets, &accum_bytes);

int payload_in = 0;
int payload_size = 0;
int tag_in = 0;
for (int i = 0; i < packets; ++i) {
    bsp_get_tag(&payload_size, &tag_in);
    bsp_move(&payload_in, sizeof(int));
    printf("payload: %i, tag: %i", payload_in, tag_in);
}
```

## Bulk: Receiving messages

```
for (auto [tag, content] : queue) {  
    world.log("payload: %i, tag: %i", content, tag);  
}
```

## Bulk: Beyond tags

- In addition, Bulk supports sending arbitrary data either using custom structs, or by composing messages on the fly. For example, to send a 3D tensor element with indices and its value.

```
auto q = bulk::queue<int, int, int, float>(world);  
q(world.next_rank()).send(1, 2, 3, 4.0f);  
q(world.next_rank()).send(2, 3, 4, 5.0f);  
world.sync();  
  
for (auto [i, j, k, value] : queue) {  
    world.log("element: A(%i, %i, %i) = %f", i, j, k, value);  
}
```

- Multiple queues can be constructed, which eliminates a common use case for tags.

## Bulk: Skeletons

*// dot product*

```
auto xs = bulk::coarray<int>(world, s);
auto ys = bulk::coarray<int>(world, s);
auto result = bulk::var<int>(world);
for (int i = 0; i < s; ++i) {
    result.value() += xs[i] * ys[i];
}
auto alpha = bulk::foldl(result,
    [](int& lhs, int rhs) { lhs += rhs; });
```

*// finding global maximum*

```
auto maxs = bulk::gather_all(world, max);
max = *std::max_element(maxs.begin(), maxs.end());
```



## Bulk: Example application (I)

- In parallel regular sample sort, there are two communication steps.
  1. Broadcasting  $p$  equidistant samples of the sorted local array.
  2. Moving each element to the appropriate remote processor.

*// Broadcast samples*

```
auto samples = bulk::coarray<T>(world, p * p);  
for (int t = 0; t < p; ++t)  
    samples(t)[{s * p, (s + 1) * p}] = local_samples;  
world.sync();
```

*// Contribution from  $P(s)$  to  $P(t)$*

```
auto q = bulk::queue<int, T[]>(world);  
for (int t = 0; t < p; ++t)  
    q(t).send(block_sizes[t], blocks[t]);  
world.sync();
```

## Bulk: Word count

- The *word count* example (MapReduce) can be implemented in Bulk as follows. First the **map** phase:

```
auto words = bulk::queue<std::string>(world);  
if (s == 0) {  
    auto f = std::fstream("examples/data/alice.txt");  
    std::string word;  
    while (f >> word) {  
        words(hash(word) % p).send(word);  
    }  
}  
world.sync();
```

## Word count (II)

- Then the **reduce** phase:

```
auto counts = std::map<std::string, int>{};
for (auto word : words) {
    if (counts.find(word) != counts.end()) {
        counts[word]++;
    } else {
        counts[word] = 1;
    }
}

auto report = bulk::queue<std::string, int>(world);
for (auto [word, count] : counts) {
    report(0).send(word, count);
}

world.sync();
```

## Bulk: Shared-memory results

**Table 1:** Speedups of parallel sort and parallel FFT compared to `std::sort` from `libstdc++`, and the sequential algorithm from FFTW 3.3.7, respectively.

	$n$	$p = 1$	$p = 2$	$p = 4$	$p = 8$	$p = 16$	$p = 32$
Sort	$2^{20}$	0.93	1.95	3.83	6.13	8.10	12.00
	$2^{21}$	1.01	2.08	4.11	7.28	10.15	15.31
	$2^{22}$	0.88	1.82	3.58	5.99	10.27	13.92
	$2^{23}$	0.97	1.90	3.63	6.19	11.99	16.22
	$2^{24}$	0.93	1.79	3.21	6.33	8.47	14.76
FFT	$2^{23}$	0.99	1.07	2.08	2.77	5.60	5.51
	$2^{24}$	1.00	1.26	2.14	3.07	5.68	6.08
	$2^{25}$	1.00	1.23	2.22	3.09	5.80	6.05
	$2^{26}$	0.99	1.24	2.01	3.28	5.48	5.97

## Bulk: Shared-memory benchmarks

**Table 2:** The BSP parameters for MCBSP and the C++ thread backend for Bulk.

Method	$r$ (GFLOP/s)	$g$ (FLOPs/word)	$l$ (FLOPs)
MCBSP (spinlock)	0.44	2.93	326
MCBSP (mutex)	0.44	2.86	10484
Bulk (spinlock) <i>*new*</i>	0.44	5.55	467
Bulk (mutex)	0.44	5.65	11702

# Outlook

- Further performance improvements for the `thread` and the MPI backends.
- Implementing popular BSP algorithms to provide case studies as a learning tool for new Bulk users.
- **Applications:** tomography, imaging science, sparse linear algebra.
- Currently working on syntax/support for **distributions:** partitionings, multi-indexing, 2D/3D computations.

## Bulk: Partitionings

```
auto phi = bulk::cyclic_partitioning<1>({size}, {p});
auto psi = bulk::cyclic_partitioning<2, 2>({n, n}, {M, N});
auto chi = bulk::block_partitioning<2, 2>({n, n}, {M, N});
// And: irregular, cartesian, tree, ...

// In LU decomposition: is a_kk assigned to us?
if (phi.owner({k, k}) == world.rank())
// What is the global index of local element (i, j)
phi.local_to_global({i, j}, {s, t})
// What is the size of my local data
phi.local_size(world.rank())
// What is my 'multi-index'?
auto [s, t] = bulk::unflatten<2>(phi.grid(), world.rank());
// What processor owns global element (i, j)?
phi.grid_owner({i, j})
```

# Conclusion

- Modern interface for writing parallel programs, safer and clearer code
- Works together with other libraries because of generic containers and higher-level functions.
- Works across more (mixed!) platforms than other libraries.
- Open-source, MIT licensed. Documentation at <http://jwbuurlage.github.io/Bulk>. Current version: **v1.1.0**.



MapReduce: Simplified Data Processing on Large Clusters, Jeffrey Dean and Sanjay Ghemawat (2004)



Pregel: A System for Large-Scale Graph Processing – Malewicz et al. (2010)



One Trillion Edges: Graph Processing at Facebook-Scale - Avery Ching et al (2015)



Buurlage JW., Bannink T., Bisseling R.H. (2018) Bulk: A Modern C++ Interface for Bulk-Synchronous Parallel Programs. Euro-Par 2018: Parallel Processing.