# CSE 291: Differentiable Programming
# Final Project: Checkpoint

**Jash Gautam Makhija** (A59025270)
**Animesh Kumar** (A59023817)

---

**Distributed AD: Add OpenMPI support to loma so that it can be distributed to multiple machines.**

## Abstract

The primary goal of our project is to enhance `loma` by integrating support for OpenMPI (Open Message Passing Interface)[1]. This integration will enable `loma` to distribute and parallelize the compilation process across multiple machines, significantly improving compilation speed and efficiency. By splitting a very large dataset across several machines, we can effectively test various inputs for a particular function, leveraging the computational power of distributed systems.

## Loma OpenMPI Extension Architecture[1]

The OpenMPI extension of Loma is supported by the below implementation:

1. The parser codebase has been extended to support another decorator @openMPI. If this decorator is found the parser sets a parameter as is_openMPI as true in the FunctionDef.

2. The compiler codebase has been modified if a loma FunctionDef is found. And generates two different C files one which acts as a parent the other which acts as a child.

3. The parent codebase is compiled as a DLL(shared library) and can be loaded easilt in python. The main purpose of the parent is to spawn $n$(present as a parameter to the parent function) child workers which will actually execute the function. Additionally the parent will send all the required parameter to the the child process and gather and receive the output from the child process. This gathered output is what is sent back to the python interface.

4. The child codebase is compiled as a regular object file. This is necessary as OpenMPI paralliztion can only be invoked in files which have a main function. This child codebase additionally contains a wrapper on top of the concerned function which receives the parameters from the parent and send back the response to the parent.

Figure 1 decribes this implementation briefly. The generated parent and child codebase are present in Listing 3.

---

[1]OpenMPI integration code: https://github.com/animeshk08/loma_public/tree/checkpoint1
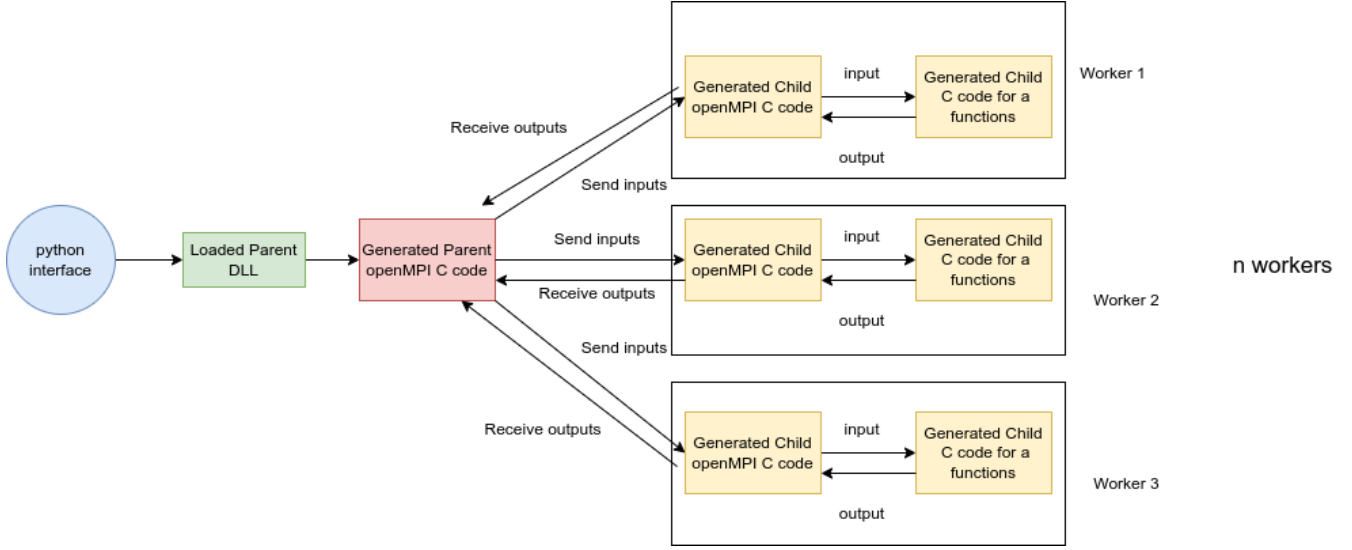
Figure 1: Open MPI implementation in Loma. The loma compiler generates two Open MPI based C code. One DLL which acts as a parent and and spawns multiple child workers. These child workers are the seconds generated files which contains the an OpenMPI wrapper on top of the actual generated C function.

Listing 1: Plus Loma code

```
1  @openMpi
2  def plus(x : In[float], y : In[float]) -> float:
3      return x + y
4
5  d_plus = fwd_diff(plus)
```

Listing 2: Python Testcase

```
1  def test_plus(self):
2      with open('loma_code/plus.py') as f:
3          structs, lib = compiler.compile(f.read(), target = 'openMpi', output_filename
   = '_code/plus')
4      _dfloat = structs['_dfloat']
5      num_worker = 4
6      px_x = [_dfloat(5.0, 0.5), _dfloat(6.0, 0.3), _dfloat(7.0, 0.2), _dfloat(8.0, 0.1)
   ]
7      px_y = [_dfloat(6.0, 1.5), _dfloat(7.0, 2.5), _dfloat(8.0, 3.5), _dfloat(9.0, 4.5)
   ]
8      input1 = (_dfloat * len(px_x))(*px_x)
9      input2 = (_dfloat * len(px_y))(*px_y)
10     py_y = [_dfloat(0, 0)] * num_worker
11     out = (_dfloat * len(py_y))(*py_y)
12     lib.mpi_runner(input1, input2, out, num_worker)
13     result = [(11,2),(13,2.8),(15,3.7),(17,4.6)]
14
15     for worker in range(num_worker):
16         flag = False
17         print(f"Result from Worker {worker+1}: val = {out[worker].val}, dval = {out[
```

2

```
          worker].dval}")
18            for res in result:
19                if abs(out[worker].val - res[0]) < epsilon and abs(out[worker].dval -
      res[1]) < epsilon:
20                    flag = True
21                    break
22            assert flag
```

Listing 3: Generated OpenMpi Code for forward diff

```
1
2 Forward differentiation of function d_plus:
3 def d_plus(x : In[_dfloat], y : In[_dfloat]) -> _dfloat:
4        return make__dfloat(((x).val) + ((y).val),((x).dval) + ((y).dval))
5
6
7 Generated C code for child:
8 #include <mpi.h>
9 #include <math.h>
10 #include <stdlib.h>
11 #include <stdio.h>
12
13 typedef struct {
14     float val;
15     float dval;
16 } _dfloat;
17
18 float plus(float x, float y);
19
20 _dfloat d_plus(_dfloat x, _dfloat y);
21
22 _dfloat make__dfloat(float val, float dval);
23
24 float plus(float x, float y) {
25     return (x) + (y);
26 }
27
28 _dfloat d_plus(_dfloat x, _dfloat y) {
29     return make__dfloat(((x).val) + ((y).val),((x).dval) + ((y).dval));
30 }
31
32 void d_plus_mpi_worker(){
33     MPI_Init(NULL, NULL);
34     MPI_Comm parent_comm;
35     MPI_Comm_get_parent(&parent_comm);
36     if (parent_comm == MPI_COMM_NULL) {
37         MPI_Finalize();
38         return;
39     }
40
41     int world_rank, world_size;
42     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
43     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
44
```

```
45      _dfloat x;
46      MPI_Recv(&x, sizeof(_dfloat), MPI_BYTE, 0, 0, parent_comm,MPI_STATUS_IGNORE);
47
48      _dfloat y;
49      MPI_Recv(&y, sizeof(_dfloat), MPI_BYTE, 0, 0, parent_comm,MPI_STATUS_IGNORE);
50
51      _dfloat out = d_plus(x, y);
52      MPI_Send(&out, sizeof(_dfloat), MPI_BYTE, 0, 0, parent_comm);
53      MPI_Finalize();
54 }
55
56 int main() {
57      d_plus_mpi_worker();
58      return 0;
59 }
60
61 _dfloat make__dfloat(float val, float dval) {
62      _dfloat ret;
63      ret.val = 0;
64      ret.dval = 0;
65      (ret).val = val;
66      (ret).dval = dval;
67      return ret;
68 }
69
70
71 Generated C code for parent:
72 #include <mpi.h>
73 #include <math.h>
74 #include <stdlib.h>
75 #include <stdio.h>
76
77 typedef struct {
78      float val;
79      float dval;
80 } _dfloat;
81
82 void mpi_runner(_dfloat * x, _dfloat * y, _dfloat* output, int __total_work) {
83      MPI_Init(NULL, NULL);
84      int world_rank, world_size;
85      MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
86      MPI_Comm_size(MPI_COMM_WORLD, &world_size);
87
88      if (world_rank == 0) {
89          // This is the parent process
90          MPI_Comm child_comm;
91          MPI_Info info;
92          MPI_Info_create(&info);
93
94          // Spawn NUM_CHILDREN child processes
95          MPI_Comm_spawn("_code/plus.o", MPI_ARGV_NULL, __total_work, info, 0,
     MPI_COMM_SELF, &child_comm, MPI_ERRCODES_IGNORE);
96
97          for (int i = 0; i < __total_work; i++) {
98              // Send input to child process
```

```
99              MPI_Send(&x[i], sizeof(_dfloat), MPI_BYTE, i, 0, child_comm);
100             MPI_Send(&y[i], sizeof(_dfloat), MPI_BYTE, i, 0, child_comm);
101         }
102
103         // Communicate with children and receive responses
104         for (int i = 0; i < __total_work; i++) {
105             _dfloat local_output;
106             MPI_Recv(&local_output, sizeof(_dfloat), MPI_BYTE, MPI_ANY_SOURCE, 0,
        child_comm, MPI_STATUS_IGNORE);
107             output[i] = local_output;
108         }
109
110         MPI_Info_free(&info);
111     }
112     MPI_Finalize();
113 }
```

Listing 4: Result

```
1 Result from Worker 1: val = 13.0, dval = 2.799999952316284
2 Result from Worker 2: val = 15.0, dval = 3.700000047683716
3 Result from Worker 3: val = 11.0, dval = 2.0
4 Result from Worker 4: val = 17.0, dval = 4.599999904632568
```

# Checkpoint 1 Results and Implementation

Till now we have setup the OpenMPI framework integration in Loma and completed the validation of 17 test cases for the forward differentiation module.

To integrate OpenMPI into the current pipeline, the following changes have been implemented:

1. **Decorator Addition**: In the forward differentiation test cases, the decorator `@OpenMpi` was added at the beginning of the function definitions.

2. **MPI Integration in testcase:** The parameters of the function are now each an array with every entry as the value that should be computed in a worker. These paramters along with an output array and total number of workers is sent to the mpi runner function, which resides in the parent code. The parent process then forwards the input parameter at $i^{th}$ index to the $i^{th}$ worker to received the $i^{th}$ output and gathers it back to output array. This aggregated result is then sent to python testcase for validation.

3. **Input Distribution**: The framework allows for distributing different inputs to various child processes, enabling performance checks for different inputs across multiple nodes.

# Scope and Timeline

**Final Submission**

- **Objective**: Implement robust synchronization and communication mechanisms for handling large-scale distributed execution and result aggregation.

- **Future Plan**:
  - We have completed 17 test cases for the forward differentiation module. Our next steps include finalizing the remaining test cases and extending the solution to support reverse differentiation for which we need to incorporate support for sending the output parameter as a pointer in the input parameter.

  - Enable the execution of programs with large input data concurrently across multiple machines. For instance, running computationally intensive algorithms such as gradient descent in parallel, followed by an all-gather operation to consolidate and aggregate outputs from all nodes, thereby significantly reducing overall execution time.

# References

[1] Open MPI. https://github.com/open-mpi/ompi. 2004.