

Differentiable Programming

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 Architecture¹

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 easily 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 parallization 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 describes this implementation briefly. The generated parent and child codebase are present in Listing 7.

¹OpenMPI integration code: https://github.com/animeshk08/loma_public/tree/final_project

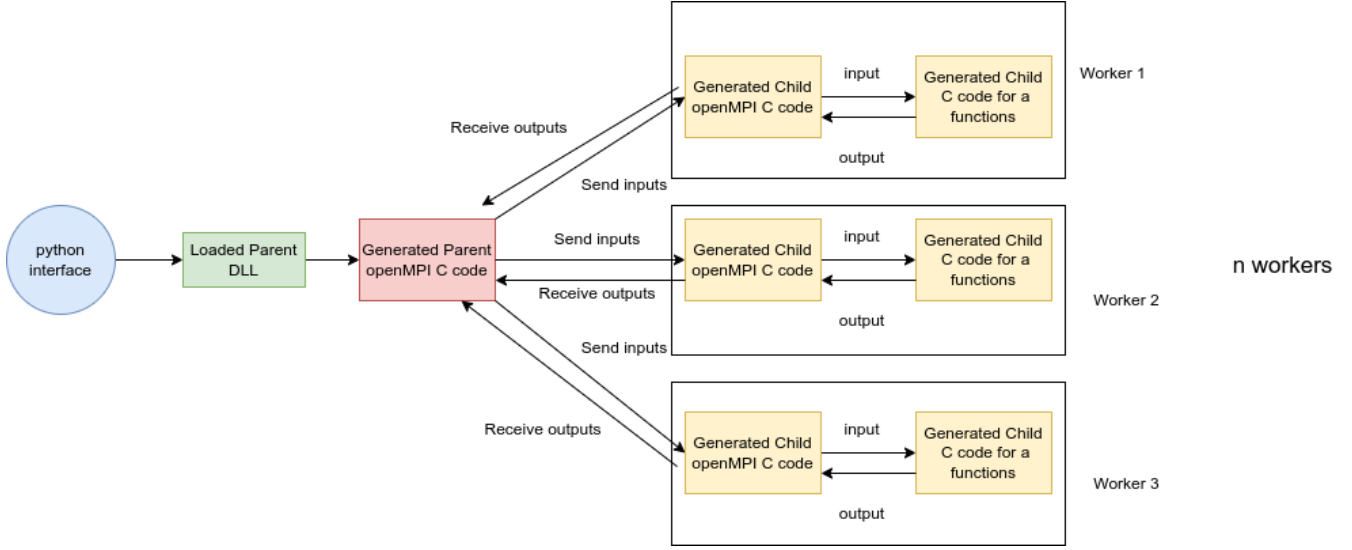


Figure 1: Open MPI implementation in Loma. The loma compiler generates two Open MPI based C code. One DLL which acts as a parent and spawns multiple child workers. These child workers are the second generated files which contain an OpenMPI wrapper on top of the actual generated C function.

Listing 1: Forward Differentiation Array Input Loma code

```

1 @openMpi
2 def array_input(x : In[Array[float]]) -> float:
3     return x[0] + x[1]
4
5 d_array_input = fwd_diff(array_input)

```

Listing 2: Python Testcase

```

1 def test_array_input(self):
2     with open('loma_code/array_input.py') as f:
3         structs, lib = compiler.compile(f.read(), target = 'openMpi',
4         output_filename = '_code/array_input')
5         _dfloat = structs['_dfloat']
6         num_worker = 2
7         px_x = [_dfloat(0.7, 0.8), _dfloat(0.3, 0.5), _dfloat(0.3, 0.4), _dfloat(0.4,
8         0.6)]
9         py_size = [2, 2]
10        input = (_dfloat * len(px_x))(*px_x)
11        input2 = (ctypes.c_int * len(py_size))(*py_size)
12        py_y = [_dfloat(0, 0)] * num_worker
13        out = (_dfloat * len(py_y))(*py_y)
14        lib.mpi_runner(input, input2, out, num_worker)
15        result = []
16        for i in range(0, num_worker+1, 2):
17            result.append((px_x[i].val + px_x[i+1].val, px_x[i].dval + px_x[i+1].dval
18        ))
19        print(result)

```

```

18     for worker in range(num_worker):
19         flag = False
20         print(f"Result from Worker {worker+1}: val = {out[worker].val}, dval = {
out[worker].dval}")
21         for res in result:
22             if abs(out[worker].val - res[0]) < epsilon and abs(out[worker].dval
- res[1]) < epsilon:
23                 flag = True
24                 break
25         assert flag

```

Listing 3: Generated OpenMpi Code for forward diff

```

1
2 Forward differentiation of function d_array_input:
3 def d_array_input(x : In[Array[_dfloat]]) -> _dfloat:
4     return make__dfloat(((x)[(int)(0)].val) + ((x)[(int)(1)].val),(((x)[(int)
)(0)].dval) + ((x)[(int)(1)].dval))
5
6 Generated C code for child:
7
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 float array_input(float* x);
18 _dfloat d_array_input(_dfloat* x);
19 _dfloat make__dfloat(float val, float dval);
20 float array_input(float* x) {
21     return ((x)[(int)(0)]) + ((x)[(int)(1)]);
22 }
23 _dfloat d_array_input(_dfloat* x) {
24     return make__dfloat(((x)[(int)(0)].val) + ((x)[(int)(1)].val),(((x)[(int)
)(0)].dval) + ((x)[(int)(1)].dval));
25 }
26 void d_array_input_mpi_worker() {
27
28     MPI_Init(NULL, NULL);
29
30     MPI_Comm parent_comm;
31     MPI_Comm_get_parent(&parent_comm);
32
33     if (parent_comm == MPI_COMM_NULL) {
34         MPI_Finalize();
35         return;
36     }
37
38     int world_rank, world_size;
39     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);

```

```

40     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
41     int x_size;
42 MPI_Recv(&x_size, 1, MPI_INT, 0, 0, parent_comm, MPI_STATUS_IGNORE);
43 _dfloat* x = (_dfloat*)malloc(x_size * sizeof(_dfloat)); MPI_Recv(x, x_size * sizeof(
44     _dfloat), MPI_BYTE, 0, 0, parent_comm, MPI_STATUS_IGNORE);
45     _dfloat out = d_array_input(x);
46     MPI_Send(&out, sizeof(_dfloat), MPI_BYTE, 0, 0, parent_comm);
47 MPI_Finalize();}
48
49 int main() {
50     d_array_input_mpi_worker();
51     return 0;
52 }
53 _dfloat make__dfloat(float val, float dval) {
54     _dfloat ret;
55     ret.val = 0;
56     ret.dval = 0;
57     (ret).val = val;
58     (ret).dval = dval;
59     return ret;
60 }
61
62 Generated C code for parent:
63
64 #include <mpi.h>
65 #include <math.h>
66 #include <stdlib.h>
67 #include <stdio.h>
68
69 typedef struct {
70     float val;
71     float dval;
72 } _dfloat;
73 void mpi_runner(_dfloat* x, int* x_size, _dfloat* output, int __total_work) {
74
75     MPI_Init(NULL, NULL);
76
77     int world_rank, world_size;
78     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
79     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
80
81     if (world_rank == 0) {
82         // This is the parent process
83         MPI_Comm child_comm;
84         MPI_Info info;
85         MPI_Info_create(&info);
86         // Spawn NUM_CHILDREN child processes
87         MPI_Comm_spawn("_code/array_input.o", MPI_ARGV_NULL, __total_work, info, 0,
88             MPI_COMM_SELF, &child_comm, MPI_ERRCODES_IGNORE);
89         int* array_end = (int*)malloc(1 * sizeof(int));
90         int end = 0;
91         int array_index = 0;
92         for (int i = 0; i < 1; i++){
93             array_end[i] = 0;

```

```

93     }
94     for (int i = 0; i < __total_work; i++) {
95         array_index = 0;
96         // Send input to child process
97         MPI_Send(&x_size[i], 1, MPI_INT, i, 0, child_comm);
98
99         _dfloat* x_i = (_dfloat*)malloc(x_size[i] * sizeof(_dfloat));
100
101         end = array_end[array_index];
102         for (int j = 0; j < x_size[i]; j++) {
103             x_i[j] = x[end++];
104         }
105         array_end[array_index] = end;
106         array_index = array_index + 1;
107         MPI_Send(x_i, sizeof(_dfloat)*x_size[i], MPI_BYTE, i, 0, child_comm)
108     };
109
110     // Communicate with children and receive responses
111     for (int i = 0; i < __total_work; i++) {
112         _dfloat local_output;
113         MPI_Recv(&local_output, sizeof(_dfloat), MPI_BYTE, i, 0, child_comm,
114 MPI_STATUS_IGNORE);
115
116         output[i] = local_output;
117     }
118     MPI_Info_free(&info);
119     MPI_Finalize();
120 }

```

Listing 4: Result

```

1 Result from Worker 1: val = 1.0, dval = 1.29999999523162842
2 Result from Worker 2: val = 0.70000000476837158, dval = 1.0

```

Listing 5: Reverse Differentiation Plus Loma code

```

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

```

Listing 6: Python Testcase

```

1 def test_plus_rev(self):
2     with open('loma_code/plus_rev.py') as f:
3         structs, lib = compiler.compile(f.read(), target = 'openMpi',
4         output_filename = '_code/plus_rev')
5         num_worker = 2
6         py_x = [ctypes.c_float(0.0), ctypes.c_float(0.0)]

```

```

6      py_y = [ ctypes.c_float(0.0) , ctypes.c_float(0.0) ]
7      x = [5.0 ,1.0]
8      y = [6.0 ,4.0]
9      dout = [3.0 ,4.0]
10     arg_dx = ( ctypes.c_float * len(py_x))(*py_x)
11     arg_x = ( ctypes.c_float * len(x))(*x)
12     arg_dy = ( ctypes.c_float * len(py_y))(*py_y)
13     arg_y = ( ctypes.c_float * len(y))(*y)
14     arg_dreturn = ( ctypes.c_float * len(dout))(*dout)
15     lib.mpi_runner(arg_x, arg_dx, arg_y, arg_dy, arg_dreturn , num_worker)
16     for i in range(num_worker):
17         print(f"Result from Worker {i+1}: {arg_dx[i]}, {arg_dy[i]}")
18         assert abs(arg_dx[i] - dout[i]) < epsilon and abs(arg_dy[i] - dout[i])
< epsilon

```

Listing 7: Generated OpenMpi Code for reverse differentiation

```

1      Reverse differentiation of function d_plus:
2  def d_plus(x : In[float], _dx : Out[float], y : In[float], _dy : Out[float],
   _dreturn : In[float]) -> void:
3      tape : Array[float, 1]
4      tape_ptr : int = (int)(0)
5      _adj_0 : float
6      _adj_1 : float
7      _adj_2 : float
8      _adj_3 : float
9      _adj_4 : float
10     _dx = (_dx) + (_dreturn)
11     _dy = (_dy) + (_dreturn)
12
13  Generated C code for child:
14
15  #include <mpi.h>
16  #include <math.h>
17  #include <stdlib.h>
18  #include <stdio.h>
19
20  typedef struct {
21      float val;
22      float dval;
23  } _dfloat;
24  float plus(float x, float y);
25  void d_plus(float x, float* _dx, float y, float* _dy, float _dreturn);
26  _dfloat make__dfloat(float val, float dval);
27  float plus(float x, float y) {
28      return (x) + (y);
29  }
30  void d_plus(float x, float* _dx, float y, float* _dy, float _dreturn) {
31      float tape[1];
32      for (int _i = 0; _i < 1; _i++) {
33          tape[_i] = 0;
34      }
35      int tape_ptr = (int)(0);
36      float _adj_0;

```

```

37     _adj_0 = 0;
38     float _adj_1;
39     _adj_1 = 0;
40     float _adj_2;
41     _adj_2 = 0;
42     float _adj_3;
43     _adj_3 = 0;
44     float _adj_4;
45     _adj_4 = 0;
46     (*_dx) = ((*_dx)) + (_dreturn);
47     (*_dy) = ((*_dy)) + (_dreturn);
48 }
49 void d_plus_mpi_worker() {
50
51     MPI_Init(NULL, NULL);
52
53     MPI_Comm parent_comm;
54     MPI_Comm_get_parent(&parent_comm);
55
56     if (parent_comm == MPI_COMM_NULL) {
57         MPI_Finalize();
58         return;
59     }
60
61     int world_rank, world_size;
62     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
63     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
64     float x;
65     MPI_Recv(&x, 1, MPI_FLOAT, 0, 0, parent_comm, MPI_STATUS_IGNORE);
66     float _dx;
67     float y;
68     MPI_Recv(&y, 1, MPI_FLOAT, 0, 0, parent_comm, MPI_STATUS_IGNORE);
69     float _dy;
70     float _dreturn;
71     MPI_Recv(&_dreturn, 1, MPI_FLOAT, 0, 0, parent_comm, MPI_STATUS_IGNORE);
72     d_plus(x, &_dx, y, &_dy, _dreturn); MPI_Send(&_dx, 1, MPI_FLOAT, 0, 0, parent_comm);
73     MPI_Send(&_dy, 1, MPI_FLOAT, 0, 0, parent_comm);
74
75     MPI_Finalize();}
76
77 int main() {
78     d_plus_mpi_worker();
79     return 0;
80 }
81 __dfloat make__dfloat(float val, float dval) {
82     __dfloat ret;
83     ret.val = 0;
84     ret.dval = 0;
85     (ret).val = val;
86     (ret).dval = dval;
87     return ret;
88 }
89
90 Generated C code for parent:
91

```

```

92 #include <mpi.h>
93 #include <math.h>
94 #include <stdlib.h>
95 #include <stdio.h>
96
97 typedef struct {
98     float val;
99     float dval;
100 } __dfloat;
101 void mpi_runner(float* x, float* _dx, float* y, float* _dy, float* _dreturn, int
    __total_work) {
102
103     MPI_Init(NULL, NULL);
104
105     int world_rank, world_size;
106     MPI_Comm_rank(MPI_COMM_WORLD, &world_rank);
107     MPI_Comm_size(MPI_COMM_WORLD, &world_size);
108
109     if (world_rank == 0) {
110         // This is the parent process
111         MPI_Comm child_comm;
112         MPI_Info info;
113         MPI_Info_create(&info);
114         // Spawn NUM_CHILDREN child processes
115         MPI_Comm_spawn("_code/plus_rev.o", MPI_ARGV_NULL, __total_work, info, 0,
MPI_COMM_SELF, &child_comm, MPI_ERRCODES_IGNORE);
116         int* array_end = (int*)malloc(0 * sizeof(int));
117         int end = 0;
118         int array_index = 0;
119         for (int i = 0; i < 0; i++){
120             array_end[i] = 0;
121         }
122         for (int i = 0; i < __total_work; i++) {
123             array_index = 0;
124             // Send input to child process
125             MPI_Send(&x[i], 1, MPI_FLOAT, i, 0, child_comm);
126             MPI_Send(&y[i], 1, MPI_FLOAT, i, 0, child_comm);
127             MPI_Send(&_dreturn[i], 1, MPI_FLOAT, i, 0, child_comm);
128         }
129
130         // Communicate with children and receive responses
131         int* rec_end = (int*)malloc(0 * sizeof(int));
132         int rend = 0;
133         int rec_array_index = 0;
134         for (int i = 0; i < 0; i++){
135             rec_end[i] = 0;
136         }
137         for (int i = 0; i < __total_work; i++) {
138             rec_array_index = 0;
139             float _temp__dx; MPI_Recv(&_temp__dx, 1, MPI_FLOAT, i, 0, child_comm,
MPI_STATUS_IGNORE);
140             _dx[i] = _temp__dx;
141             float _temp__dy; MPI_Recv(&_temp__dy, 1, MPI_FLOAT, i, 0, child_comm,
MPI_STATUS_IGNORE);
142             _dy[i] = _temp__dy;

```



```

143
144     }
145     MPI_Info_free(&info);} MPI_Finalize();}

```

Listing 8: Result

```

1 Result from Worker 1: 3.0 , 3.0
2 Result from Worker 2: 4.0 , 4.0

```

Implementation

To enable Loma to handle OpenMPI operations effectively, we made several significant modifications to the code base. These enhancements are detailed below:

1. **Codegen_openMPI.py**: This introduces a sophisticated architecture for managing OpenMPI operations across parent and child nodes.

(a) A parent-child structure was developed. The parent function is responsible for the following:

- i. **Spawning Children**: It initiates child processes according to the user specified number of workers.
- ii. **Data Distribution**: Inputs from the test case are broken up and forwarded to the respective child nodes.
- iii. **Result Aggregation**: It aggregates the results received from child nodes and returns them to the test case.

The parent node acts as a conduit between the test case and the child nodes, relaying parameters to the children and gathering results for return to the test case.

(b) The parent node incorporates the main OpenMPI main function, which processes all function parameters. Key functionalities include:

- i. **Parameter Handling**: Each parameter is transmitted as an array containing values for each child node.
- ii. **Array Management**: When an array is transmitted, it is amalgamated into a single large array. An additional array indicating the size of each child's data portion is also sent.
- iii. **Data Distribution**: During the spawning of child nodes, the parent node splits the large array into chunks based on the size array and sends these chunks to the corresponding child nodes.

(c) At the child nodes:

- i. **Data Reception**: The nodes receive their respective data portions.
- ii. **Function Execution**: They invoke the main function with the received input values.

- iii. **Result Transmission:** The results/output of forward differentiation or the necessary variables for reverse differentiation are sent back to the parent node.
- (d) The parent node performs the following:
 - i. **Data Collection:** It collects outputs from all child nodes.
 - ii. **Aggregation:** The results are aggregated into a comprehensive array, which is then sent back to the test case.
- (e) Special attention is given to the output of arrays:
 - i. **Flattening:** Arrays from all child nodes are flattened into a single array.
 - ii. **Reassembly:** At the test case level, the array is split back into its original segments based on the size information and the results are displayed accordingly.
- (f) **Special Consideration:**
 - i. **Data Types:** Careful handling of various data types such as float and integers is implemented.
 - ii. **Multiple Outputs:** For multiple outputs, especially arrays, the system considers size values to split and combine arrays properly.
 - iii. In scenarios with array outputs, the parent function aggregates values sequentially from the children; in other cases, the order of reception is randomized to optimized for speed.
- 2. **Python Test Case's:** The Python test cases are adapted to support multiple value transmissions and parallel processing:
 - (a) **Worker Specification:** For functions like `array_input`, multiple test inputs are processed in parallel. For instance, if the number of workers is set to 2, then two different test input configurations can be tested in parallel.
 - (b) **Array Handling:** Inputs for each worker are combined into a single array, with an accompanying size parameter specifying each child worker's array size.
 - (c) **Output Validation:** If arrays are outputted, the parent's output is split according to size parameters to validate correctness.
 - (d) **Expected Results:** Expected outcomes are computed and stored in a results array. The output from the parent is checked against these expected results for correctness.
 - (e) We tested our program against 30 test cases for forward and reverse differentiation respectively.

Experiment and Observation

To evaluate the performance and benefits of our OpenMPI-based automatic differentiation system, we conducted two key experiments:

1. Generality and Efficacy of Our Implementation:

- (a) **Objective:** Demonstrate the robustness and versatility of our system by optimizing a third-order multivariate polynomial using a gradient descent routine, where the gradient is computed via automatic differentiation.
- (b) **Implementation:**
- i. **Task Distribution:** The task is distributed across 2 worker nodes.
 - **Node Responsibilities:**
 - One worker node computes the gradient in the x -direction.
 - The other worker node computes the gradient in the y -direction.
 - **Gradient Calculation:** Each node performs gradient calculations independently.
 - **Position Update:** The results from both nodes are aggregated to update the position in the optimization space, moving towards the minimum of the loss function.

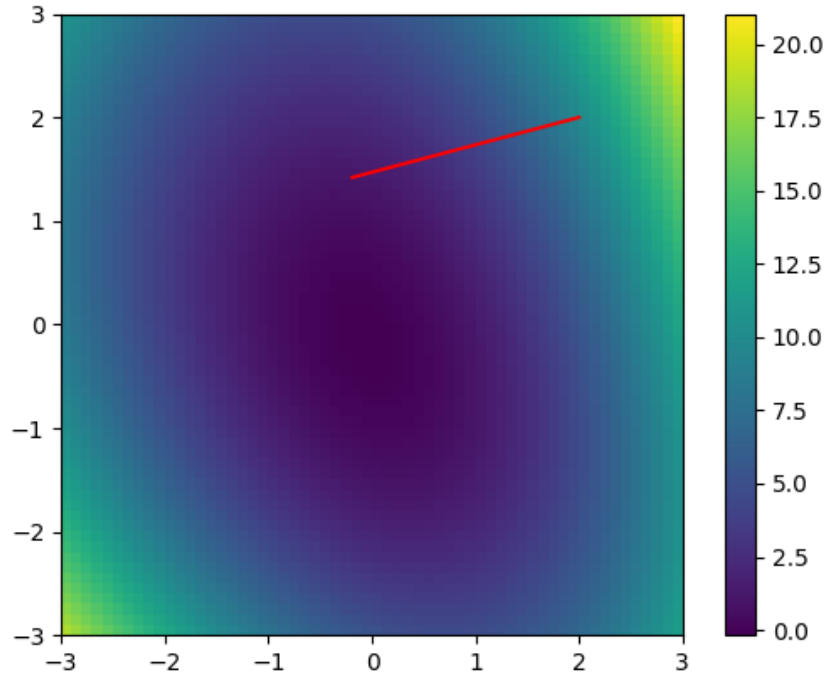


Figure 2: Optimization of a third-order multivariate polynomial using Loma with OpenMPI. The loss function landscape and the optimization trajectory (red curve) show the system starting from (2,2) and progressing towards the center where the loss is minimized.

- (c) **Results:**

- **Figure 2:** Illustrates the optimization trajectory overlaid on the loss function landscape. The red curve indicates the movement from the initial point $(2, 2)$ towards the center, achieving a lower loss.
- (d) **Limitation:** The current implementation limits the parent process to spawn only once, preventing multiple iterations within a gradient descent loop. Consequently, the system computes gradients and updates positions only for a single iteration. This limitation arises from the restriction that the `lib.mpi_runner` function cannot be invoked repeatedly within a loop to dynamically update parameters.

2. Performance Evaluation: OpenMPI vs. Sequential Execution

- (a) **Objective:** Assess the performance gains achieved by our OpenMPI-based implementation compared to a traditional sequential execution approach.
- (b) **Experimental Setup:**
- Test Scenario 1: Varying the number of input values for a single test case to measure execution time.**
 - **Methodology:**
 - **OpenMPI Implementation:** Deployed 100 worker nodes to handle input data in parallel.
 - **Sequential Implementation:** Employed a standard loop to process each input test case individually.
 - **Results:**
 - **Figure 3:** Shows the time comparison between OpenMPI and sequential execution for varying input sizes. Initially, the sequential method is faster due to the overhead associated with initializing MPI processes. However, as input size scales up to 10^8 , OpenMPI outperforms due to its ability to handle large datasets more efficiently by distributing tasks across multiple nodes.
 - Test Scenario 2: Evaluating performance for varying iteration counts in a while loop.**
 - **Methodology:**
 - **OpenMPI Implementation:** Distributed tasks across 100 worker nodes for each iteration.
 - **Sequential Implementation:** Executed iterations in a conventional loop.
 - **Results:**
 - **Figure 4:** Demonstrates that for small iteration counts, both approaches have comparable performance. As the iteration count increases, OpenMPI gains a performance edge due to its ability to parallelize the workload effectively.

The experiments were run on a Apple Macbook Pro (3.2 GHz Octa-Core Apple M1 Pro) with 16 GB RAM. These results underscore the potential of integrating OpenMPI with automatic differentiation

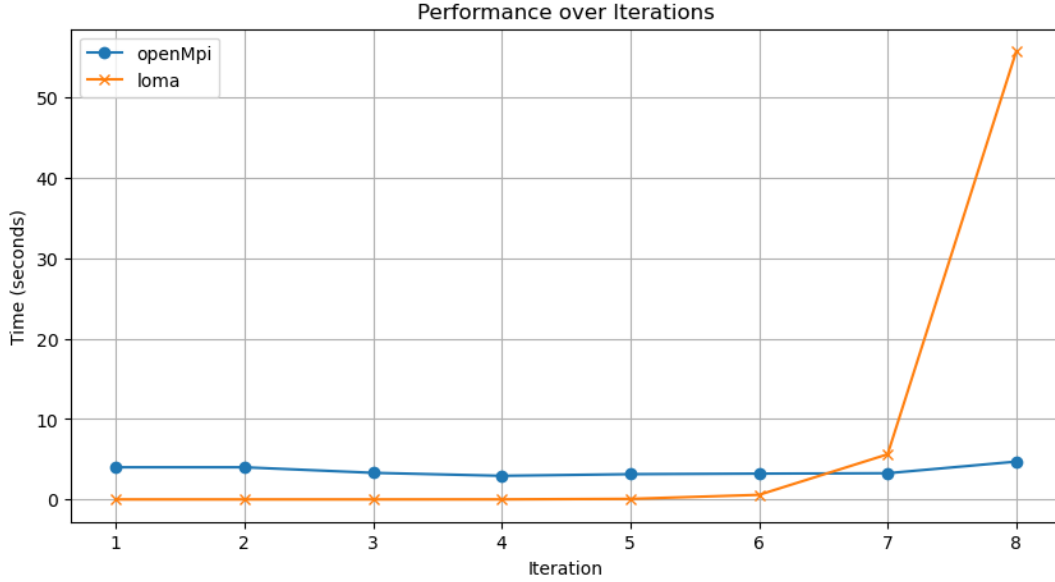


Figure 3: Time comparison for executing varied input sizes between OpenMPI and a traditional sequential approach. The graph indicates that OpenMPI’s time advantage becomes significant as input size increases.

frameworks to enhance computational efficiency and scalability for various applications in scientific computing and machine learning.

Limitation

We have created a simple and straightforward OpenMPI extension to the Loma Compiler. However, there are still a few limitations which are listed below:

1. **Single Function OpenMPI Conversion:** The OpenMPI extension generates two C compiled files to be created, a regular compiler object file for the child which includes the actual function call wrapped in an OpenMPI environment wrapper, and a second, a parent which is compiled as a DLL and which spawns the child code. Thus, the child needs to have a main function and the parent needs to know which child object file to call. Such tight coupling limited our current implementation to convert only a single method to be compiled by the OpenMPI compiler. A possible fix for this can include creating multiple child object files and including the name of the same as a parameter to the parent code base.
2. **Single Execution Environment:** Since the OpenMPI environment can only be "initialized" and "finalized" once in a process, the current OpenMPI extension can run only one single instance of an execution. This means that current multiple OpenMPI test cases cannot run together or multiple calls to the same OpenMPI based parent-child framework are also not supported. A possible fix for this is to invoke the OpenMPI executions as separate processes within the Python-based invocation mechanism.
3. **Limited Support for Arrays:** The current implementation of OpenMPI struggles with arrays as input, since in such cases the actual input to the OpenMPI wrapper will be multiple

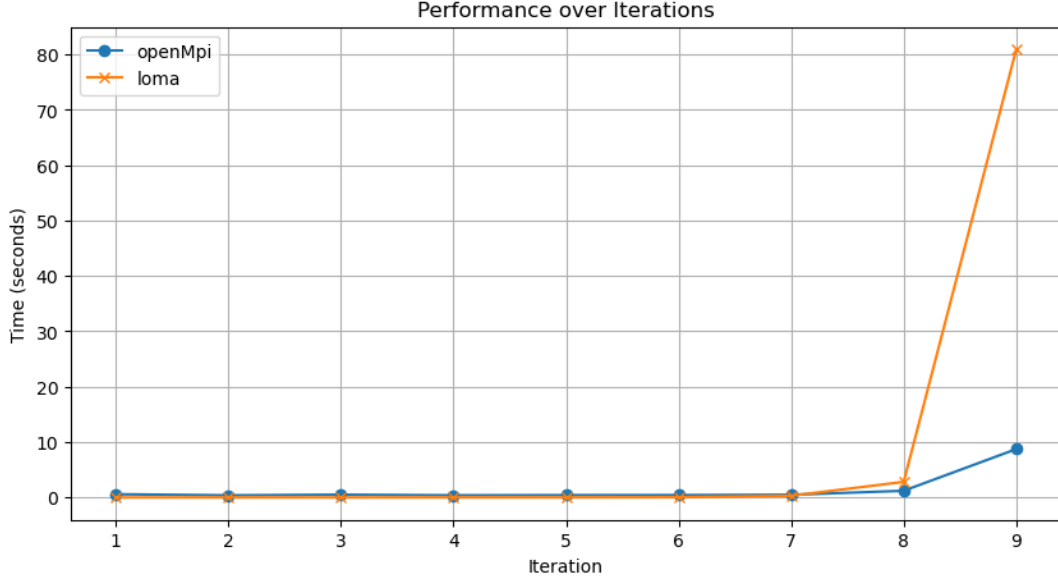


Figure 4: Time comparison for executing varied iteration counts in a while loop between OpenMPI and a traditional sequential approach. OpenMPI exhibits increasing efficiency with higher iteration counts.

such arrays, or a 2-D array. Since 2D arrays are hard to compile using the `ctypes` library, the OpenMPI extension suffers from the same problems. The current implementation flattens the 2D array and passes another list of delimiter indexes to unflatten the array later. Another possible way could be to convert the 2D array to an array of structs. Nevertheless, these approaches are obviously not scalable and hence we claim only limited support for arrays in our implementation.

Result

In this project, we have completed a simple, straightforward, and comprehensible extension to the Loma compiler to support distributed AD using OpenMPI. As part of this project, we have achieved the following milestones:

1. Created an innovative way to compile and spawn parent-child based MPI interface.
2. Supported native message passing between parent and child nodes.
3. Supported backwards compatibility with existing Loma test cases to compile and run in an OpenMPI environment.
4. Demonstrated a simple real-world example (single step of gradient descent) to show applications and advantages of OpenMPI based distributed AD.
5. Conducted a simple benchmark analysis to compare the performance of native Loma compiled C code and OpenMPI extension compiled C code.

References

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