# Implementing an **Open64**-based tool for improving the performance of **MPI** programs

### Anthony Danalis Lori Pollock Martin Swany John Cavazos



University of Delaware

## Problem:

### **Communication Overhead in Cluster Computing**



## Problem:

### **Communication Overhead in Cluster Computing**



• Inherent to the Application

• Hard or impossible to change

## **Problem:**

### **Communication Overhead in Cluster Computing**



### **Requirements:**

- Achieve high-performance communication
- Simplify the MPI code developers write

Have your cake + Eat your cake

### Requirements:

- Achieve high-performance communication
- Simplify the MPI code developers write



### Requirements:

- Achieve high-performance communication
- Simplify the MPI code developers write

## **Proposed Solution:**

An **automatic** system that **transforms** simple communication code into efficient code.



## Requirements:

- Achieve high-performance communication
- Simplify the MPI code developers write

## **Proposed Solution:**

An **automatic** system that **transforms** simple communication code into efficient code.



Enables legacy parallel MPI applications to perform better, even if written without any knowledge of this system



# Our Framework : ASPhALT\*



\*Automated System for Parallel AppLication Transformation

# Our Framework : ASPhALT\*



\*Automated System for Parallel AppLication Transformation

# Integration with Open64



### Pros:

- Portability
- Rapid Prototyping

### Cons:

Underutilization of existing:

- Program analysis
- Program transformations

# Challenges & Opportunities

Gravel enables explicit memory registration & rendezvous handshaking at the application layer

- Memory Registration Location
  - Smart placement can reduce registration cost
  - x Choosing location requires inter-procedural analysis
- Rendezvous Protocol Choice
  - Early handshake initiation can overlap the control messages
  - Advanced rendezvous can reduce control messages & increase overlap
  - x Control flow analysis required to map MPI code to Gravel protocol

### Memory Registration: context awareness matters

Frequent registering & unregistering message buffers induces significant cost



Compiler Tool: Fully context aware

```
foo(){
  register_msg_buffer()
  do i=1, N
     bar()
  enddo
  unregister_msg_buffer()
bar(){
  isend()
  wait()
```

# Gravel rendezvous protocols

### Use of wrong rendezvous protocol would lead to deadlock





# Mapping MPI to Gravel rendezvous

Trivial: execution order = textual order

```
mpi_irecv()
```

Overview

do i=1,N sbuf[ i ] = ... enddo

```
mpi_isend( sbuf )
mpi_waitall()
```

More complex: execution order = reverse textual order

```
do i=1,N
  if(i > 1) then
     mpi irecv()
     mpi waitall()
  endif
  do j=1,N
     sbuf[ i ][ j ] = ...
  enddo
  if (i < N) then
     mpi_isend( sbuf )
  endif
enddo
```

Example: NAS SP {x,y,z}\_solve.f

# Further Optimizations (future work)





# **Current & Future Directions**

- > Break MPI collectives into loops of send()/recv()
- Use Loop Nest Optimizations to increase overlap
- Integrate code into Open64 backend (be)

## Questions ?

## Gravel: Performance graphs (1)

GATHER\_MPI\_ASYNCH NP=16 Buff=4MB

GATHER\_GRAVEL\_ONESIDED NP=16 Buff=4MB



# Gravel: Performance graphs (2)

A2A\_MPI\_ASYNCH NP=16 Buff=4MB A2A\_GRAVEL\_ONESIDED NP=16 Buff=4MB 2097152 2097152 
 1048576
 Communication

 524288
 unication

 262144
 tion

 131072
 Overhead

 32768
 usec)

 16384
 Communication

 1048576 Communication

 524288
 unication

 262144
 ation

 131072
 Overhead

 32768
 (usec)

 16384
 c)
 8192 8192 4096 4096 128 128 64 64 32 32 Number of tiles Number of tiles 0 20 40 60 Kernel Exec. Time (msec) Kernel Exec. Time (msec)