# Gradient of MPI-parallel codes

Jean Utke, Paul Hovland — *ANL*

Laurent Hascoët, Valérie Pascual — *INRIA Sophia-Antipolis*

Patrick Heimbach, Chris Hill — *MIT Boston*

Uwe Naumann, Michel Schanen — *RWTH Aachen*

ANL, INRIA, MIT, RWTH

June 14th, 2012

# Outline

# Adjoint Algorithms

Given an **algorithm** P that computes a **function** $f$,
the **"adjoint algorithm"** $\overline{\text{P}}$ computes the **gradient** of $f$.

Adjoint Algorithms

- compute the gradient at a cost
  that is **independent** from the number of inputs of $f$.
- compute the gradient backwards
  $\Rightarrow$ **reversal** of the control-flow and of the data-flow.
- can be built from P by an **Automatic Differentiation** tool.

# Structure of Adjoint Programs

**Algorithm P:**
    **input** a,b,c
    ...
    ...

    u = $g$(a,c)

    ...
    ...

    r = $h$(u,v)

    **output** r

# Structure of Adjoint Programs

**Algorithm P:**

    **input** `a,b,c`

    ...

    ...

    `u = `$g$`(a,c)`

    ...

    ...

    `r = `$h$`(u,v)`

    **output** `r`

**Algorithm $\overline{\text{P}}$:**

        **input** $\overline{\text{r}}$

# Structure of Adjoint Programs

**Algorithm P:**

    **input** `a,b,c`

    ...

    ...

    `u = g(a,c)`

    ...

    ...

    `r = h(u,v)`

    **output** `r`

**Algorithm $\overline{\text{P}}$:**

$$\overline{\text{r}} = 0.0$$
$$\overline{\text{v}} \mathrel{+}= \frac{\partial h}{\partial v} * \overline{\text{r}}$$
$$\overline{\text{u}} \mathrel{+}= \frac{\partial h}{\partial u} * \overline{\text{r}}$$

**input** $\overline{\text{r}}$

# Structure of Adjoint Programs

**Algorithm P:**

    **input** `a,b,c`

    ...

    ...

    `u = `$g$`(a,c)`

    ...

    ...

    `r = `$h$`(u,v)`

    **output** `r`

**Algorithm $\overline{\text{P}}$:**

    $\bar{\text{u}} = 0.0$

    $\bar{\text{c}} \mathrel{+}= \frac{\partial g}{\partial c} * \bar{\text{u}}$

    $\bar{\text{a}} \mathrel{+}= \frac{\partial g}{\partial a} * \bar{\text{u}}$

    ...

    ...

    $\bar{\text{r}} = 0.0$

    $\bar{\text{v}} \mathrel{+}= \frac{\partial h}{\partial v} * \bar{\text{r}}$

    $\bar{\text{u}} \mathrel{+}= \frac{\partial h}{\partial u} * \bar{\text{r}}$

    **input** $\bar{\text{r}}$

# Structure of Adjoint Programs

**Algorithm P:**
   **input** a,b,c
   ...
   ...

   u = $g$(a,c)

   ...
   ...

   r = $h$(u,v)

   **output** r

**Algorithm $\overline{P}$:**
   **output** $\overline{a}$, $\overline{b}$, $\overline{c}$
   ...
   ...
   $\overline{u}$ = 0.0
   $\overline{c}$ += $\frac{\partial g}{\partial c}$ * $\overline{u}$
   $\overline{a}$ += $\frac{\partial g}{\partial a}$ * $\overline{u}$
   ...
   ...
   $\overline{r}$ = 0.0
   $\overline{v}$ += $\frac{\partial h}{\partial v}$ * $\overline{r}$
   $\overline{u}$ += $\frac{\partial h}{\partial u}$ * $\overline{r}$
   **input** $\overline{r}$

# Adjoining Simple Instructions

| Original: | Adjoint: |
|---|---|
| x = a + 2*b | $\overline{a}$ = $\overline{a}$ + $\overline{x}$ ; $\overline{b}$ = $\overline{b}$+2*$\overline{x}$ ; $\overline{x}$ = 0 |
| x = 2*x | $\overline{x}$ = 2*$\overline{x}$ |
| y = sin(x) | $\overline{x}$ = $\overline{x}$ + cos(x)*$\overline{y}$ ; $\overline{y}$ = 0 |
| b = a | $\overline{a}$ = $\overline{a}$ + $\overline{b}$ ; $\overline{b}$ = 0 |
| s = SUM(T(:)) | $\overline{T}$(:) = $\overline{T}$(:) + $\overline{s}$ ; $\overline{s}$ = 0 |
| U(2:9) = U(2:9) + x | $\overline{x}$ = $\overline{x}$ + SUM($\overline{U}$(2:9)) |
| where(T>3) T = T - a | $\overline{a}$ = $\overline{a}$ - SUM($\overline{T}$,T>3) |

All these can be "proved" formally ... but a convenient justification is
**backwards propagation of the influence** on the result.

# Outline

# Adjoining one Global Communication

# Adjoining Global Communications

| Original: | Adjoint: |
|---|---|
| `bcast(x,,P,)` | `reduce(`$\overline{\mathrm{x}}$`,`$\overline{\mathrm{t}}$`,,,SUM,P,)` |
| | $\overline{\mathrm{x}}$`=0.0 ; on P:`$\overline{\mathrm{x}}$`=`$\overline{\mathrm{x}}$`+`$\overline{\mathrm{t}}$ |
| `reduce(x,y,,,SUM,P,)` | `on P:`$\overline{\mathrm{t}}$`=`$\overline{\mathrm{y}}$` ; on P:`$\overline{\mathrm{y}}$`=0.0` |
| | `bcast(`$\overline{\mathrm{t}}$`,,P,) ; `$\overline{\mathrm{x}}$`=`$\overline{\mathrm{x}}$`+`$\overline{\mathrm{t}}$ |
| `allreduce(x,y,,,SUM,)` | `allreduce(`$\overline{\mathrm{y}}$`,`$\overline{\mathrm{t}}$`,,,SUM,)` |
| | $\overline{\mathrm{y}}$`=0.0 ; `$\overline{\mathrm{x}}$`=`$\overline{\mathrm{x}}$`+`$\overline{\mathrm{t}}$ |
| `gather(x,,y,,P,)` | `scatter(`$\overline{\mathrm{y}}$`,,`$\overline{\mathrm{t}}$`,P,)` |
| | `on P:`$\overline{\mathrm{y}}$`=0.0 ; `$\overline{\mathrm{x}}$`=`$\overline{\mathrm{x}}$`+`$\overline{\mathrm{t}}$ |
| `scatter(x,,y,,P,)` | `gather(`$\overline{\mathrm{y}}$`,,`$\overline{\mathrm{t}}$`(:),,P,)` |
| | $\overline{\mathrm{y}}$`=0.0 ; on P:`$\overline{\mathrm{x}}$`=`$\overline{\mathrm{x}}$`+`$\overline{\mathrm{t}}$`(:)` |

# Outline

# Adjoining Point-to-point Communication

- Blocking ⇔ {Nonblocking ; `wait` }
- Analogy `send(a)/receive(b)` with `b = a`



Gives us the adjoints of MPI `isend`, `irecv`, and `wait`.

# Rules for adjoining MPI calls

Adjoining rules depend on the context

- between nonblocking calls and their `wait`
- between the two communicating processes.

| in P | | in $\overline{P}$ |
|------|------|------|
| call | paired with | call |
| `isend(a,r)` | `wait(r)` | `wait(r);`$\overline{a}$`+=t` |
| `wait(r)` | `isend(a,r)` | `irecv(t,r)` |
| `irecv(b,r)` | `wait(r)` | `wait(r);`$\overline{b}$`=0` |
| `wait(r)` | `irecv(b,r)` | `isend(`$\overline{b}$`,r)` |
| `bsend(a)` | `recv(b)` | `recv(t);`$\overline{a}$`+=t` |
| `recv(b)` | `bsend(a)` | `bsend(`$\overline{b}$`);`$\overline{b}$`=0` |
| `ssend(a)` | `recv(b)` | `recv(t);`$\overline{a}$`+=t` |
| `recv(b)` | `ssend(a)` | `ssend(`$\overline{b}$`);`$\overline{b}$`=0` |

# Adjoinable MPI requires more information

Adjoining a `wait` requires some context:

- Each `wait` must know what it waits for (`isend` or `irecv`)
- Each `wait` must know and **see** the travelling variable

Source analysis could find matching `wait` of a non-blocking call, but

- not always, and
- making the travelling variable visible is harder.

Instead, use an adjoinable MPI such as:

- `waitrecv(b,r)`, whose adjoint is `isend(`$\overline{\text{b}}$`, r)`
- `waitsend(a,r)`, whose adjoint is `irecv(t, r)`

# Outline

# Open question: grouped `wait`'s

"`waitall`" groups "`wait`" operations. Improves efficiency.

# Open question: grouped `wait`'s

"`waitall`" groups "`wait`" operations. Improves efficiency.



Adjoining re-introduces separate `wait`'s !
$\Rightarrow$ can we get a `waitall()` back?

# The "anti waitall"

Sometimes, we may introduce a awaitall()

- nonoperational
- placed by the end-user.
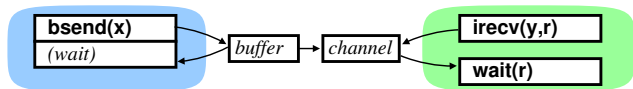


⇒ Allows the adjoint to use a waitall() again.

# Outline

# Data dependence graph for Point-to-Point

- Introduce "*channel*" pseudo variables.
- Nonblocking `isend`/`irecv`



- `send`⇔{`isend`; `wait`}    `recv`⇔{`irecv`; `wait`}
- Buffered `bsend` uses an intermediate copy buffer
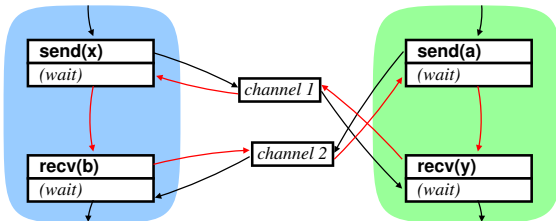  ⇒ immediate return.

# Deadlocks ; Blocking vs Nonblocking

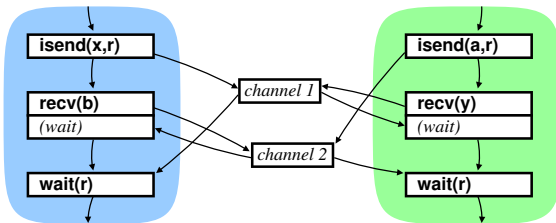Deadlocks are cycles in the data dependence graph:

# Deadlocks ; Blocking vs Nonblocking

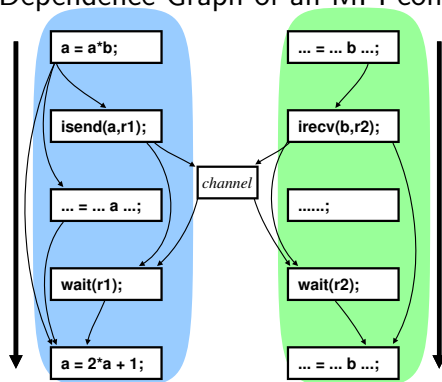Deadlocks are cycles in the data dependence graph:



Splitting the `wait` from the `isend`/`irecv` can solve the problem:
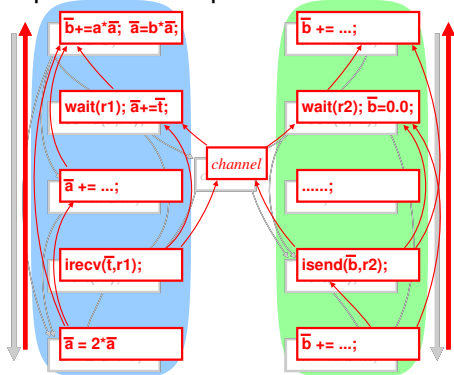


Otherwise, use `bsend`'s.

# Data Dependence Graph of the Adjoint Algorithm

Consider the Data Dependence Graph of an MPI communication.

# Data Dependence Graph of the Adjoint Algorithm

Consider the Data Dependence Graph of an MPI communication.



The Data Dependence Graph of the adjoint communication
seems to just reverse the arrows.
⇒ Adjoining does not introduce deadlocks.

# Outline

# Impact on Data-Flow analyses

AD, like other program transformations, needs preliminary data-flow analyses

Message-passing modifies data-flow $\Rightarrow$ we must adapt analyses.

- If Interprocedural Control-Flow Graph $\rightarrow$ introduce special flow arrows that only convey messages.
- If Call Graph of Flow Graphs $\rightarrow$ introduce special "channel" variables and organize additional fixed-point iterations.

In any case, increases complexity/cost of analyses.

# Outline

# References

- Utke, Hascoët, Heimbach, Hill, Hovland, Naumann **"Towards Adjoinable MPI"** *IEEE IPDPS, 2009*
- Schanen, Naumann, Hascoët, Utke **"Interpretative Adjoints for Numerical Simulation Codes using MPI"** *Procedia Computer Science, 2010*
- Schanen, Naumann, Hascoët, Utke **"Interpretative Adjoints for Numerical Simulation Codes using MPI"** *ICCS 2010*
- Schanen, Förster, Naumann **"Second-Order Algorithmic Differentiation by Source Transformation of MPI Code"** *EuroMPI 2010*
- Naumann, Hascoët, Hill, Hovland, Riehme, Utke **"A Framework for Proving Correctness of Adjoint Message-Passing Programs"** *PVM/MPI Users' Group Meeting, 2008*
- Strout, Kreaseck, Hovland **"Data-Flow analysis for MPI programs"** *ICPP, 2006*
- Kreaseck, Strout, Hovland **"Depth Analysis of MPI programs"** *AMP workshop, PLDI 2010*
- Pascual, Hascoët **"Native handling of Message-Passing communication in Data-Flow analysis"** *AD2012*

# Outline

# Conclusion: an adjoinable MPI ?

- Unlike previous black-box approaches,
  we consider the level of the **individual MPI calls**.
- Our approach requires the **correspondence** between MPI calls
- Static data-flow analysis will not find it in general:
  ⇒ **User input** is necessary:
  - could be pragmas
  - could be an "adjoinable MPI" library

# Outlook

- Ongoing **application to the MITgcm**: validates the adjoining rules of MPI calls. Requires the `awaitall`.

- Look for a **general proof** of correctness, maybe based on PGAS or other abstractions e.g. Data Dependence graphs.

- Develop an **adjoinable MPI library**, to help/induce the user write an adjoinable code..

# Outlook

- Ongoing **application to the MITgcm**: validates the adjoining rules of MPI calls. Requires the `awaitall`.
- Look for a **general proof** of correctness, maybe based on PGAS or other abstractions e.g. Data Dependence graphs.
- Develop an **adjoinable MPI library**, to help/induce the user write an adjoinable code..

## Thank you for your attention !