# Automatic Differentiation for Optimum Design, Applied to Sonic Boom Reduction

Laurent Hascoët, Mariano Vázquez, Alain Dervieux

Laurent.Hascoet@sophia.inria.fr

TROPICS Project, INRIA Sophia-Antipolis

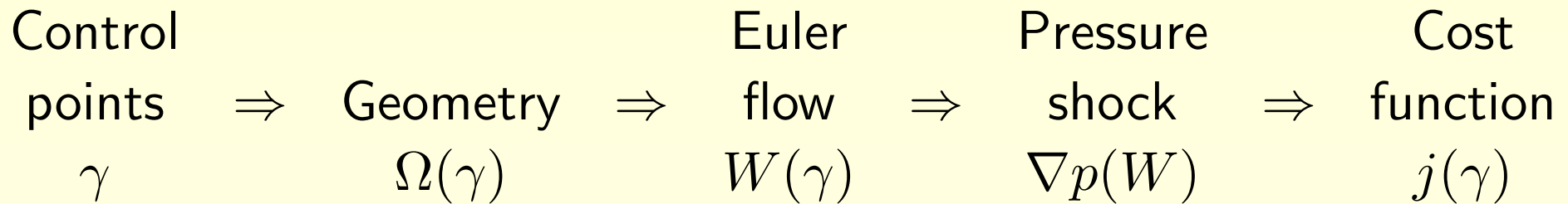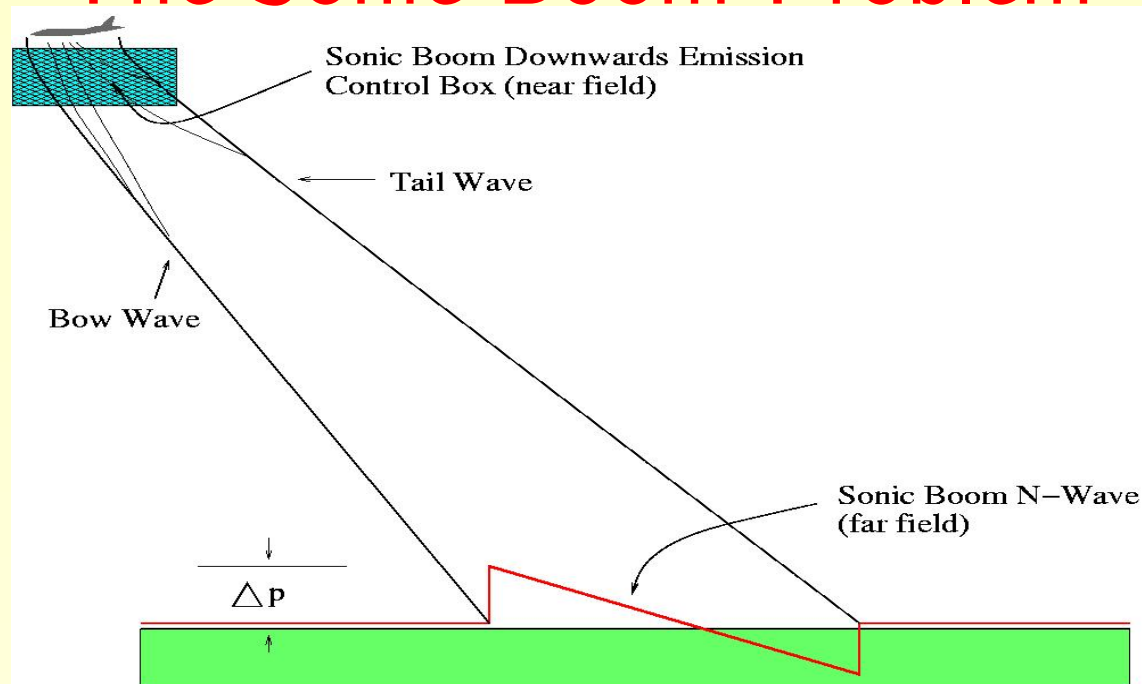AD Workshop, Cranfield, June 5-6, 2003

# PLAN:

- **Part 1:** A gradient-based shape optimization to reduce the Sonic Boom

- **Part 2:** Utilization and improvements of reverse A.D to compute the Adjoint

- **Conclusion**

# PART 1:

# GRADIENT-BASED SONIC BOOM REDUCTION

- The Sonic Boom optimization problem

- A mixed Adjoint/AD strategy

- Resolution of the Adjoint and Gradient

- Numerical results

# The Sonic Boom Problem



Sonic Boom Downwards Emission
Control Box (near field)

Tail Wave

Bow Wave

Sonic Boom N−Wave
(far field)

$\triangle \mathbf{P}$

Control
points $\quad \Rightarrow \quad$ Geometry $\quad \Rightarrow \quad$ Euler
flow $\quad \Rightarrow \quad$ Pressure
shock $\quad \Rightarrow \quad$ Cost
function

$\gamma \qquad\qquad \Omega(\gamma) \qquad\qquad W(\gamma) \qquad\qquad \nabla p(W) \qquad\qquad j(\gamma)$

# Gradient-based Optimization

$j(\gamma)$ models the strength of the downwards Sonic Boom

$\Rightarrow$ Compute the gradient $j'(\gamma)$ and use it in an optimization loop!

$\Rightarrow$ Use reverse-mode AD to compute this gradient

# Differentiate the iterative solver?

$W(\gamma)$ is defined implicitly through the Euler equations

$$\Psi(\gamma, W) = 0$$

$\Rightarrow$ The program uses an iterative solver

$\Rightarrow$ Brute force reverse AD differentiates the whole iterative process

- Does it make sense?

- Is it efficient ?

# A mixed Adjoint/AD strategy

Back to the mathematical adjoint system:

$$
\left\{
\begin{aligned}
\Psi(\gamma, W) &= 0 \quad \textit{(state system)} \\[2em]
\frac{\partial J}{\partial W}(\gamma, W) - (\frac{\partial \Psi}{\partial W}(\gamma, W))^t \cdot \Pi &= 0 \quad \textit{(adjoint system)} \\[2em]
j'(\gamma) = \frac{\partial J}{\partial \gamma}(\gamma, W) - (\frac{\partial \Psi}{\partial \gamma}(\gamma, W))^t \cdot \Pi &= 0 \quad \textit{(optimality condition)}
\end{aligned}
\right.
$$

lower level $\Rightarrow$ reverse AD *cf Part 2*

upper level $\Rightarrow$ hand-coded specific solver for $\Pi$

# Upper level

Solve $\dfrac{\partial \Psi}{\partial W}(\gamma, W))^t \cdot \Pi = \dfrac{\partial J}{\partial W}(\gamma, W)$

$\Rightarrow$ Use a matrix-free solver

Preconditioning: "defect correction"

$\Rightarrow$ use the inverse of $1^{st}$ order $\dfrac{\partial \Psi}{\partial W}$ to precondition $2^{nd}$ order $\dfrac{\partial \Psi}{\partial W}$

# Overall Optimization Loop

Loop:

- compute $\Pi$ with a matrix-free solver
- use $\Pi$ to compute $j'(\gamma)$
- 1-D search in the direction $j'(\gamma)$
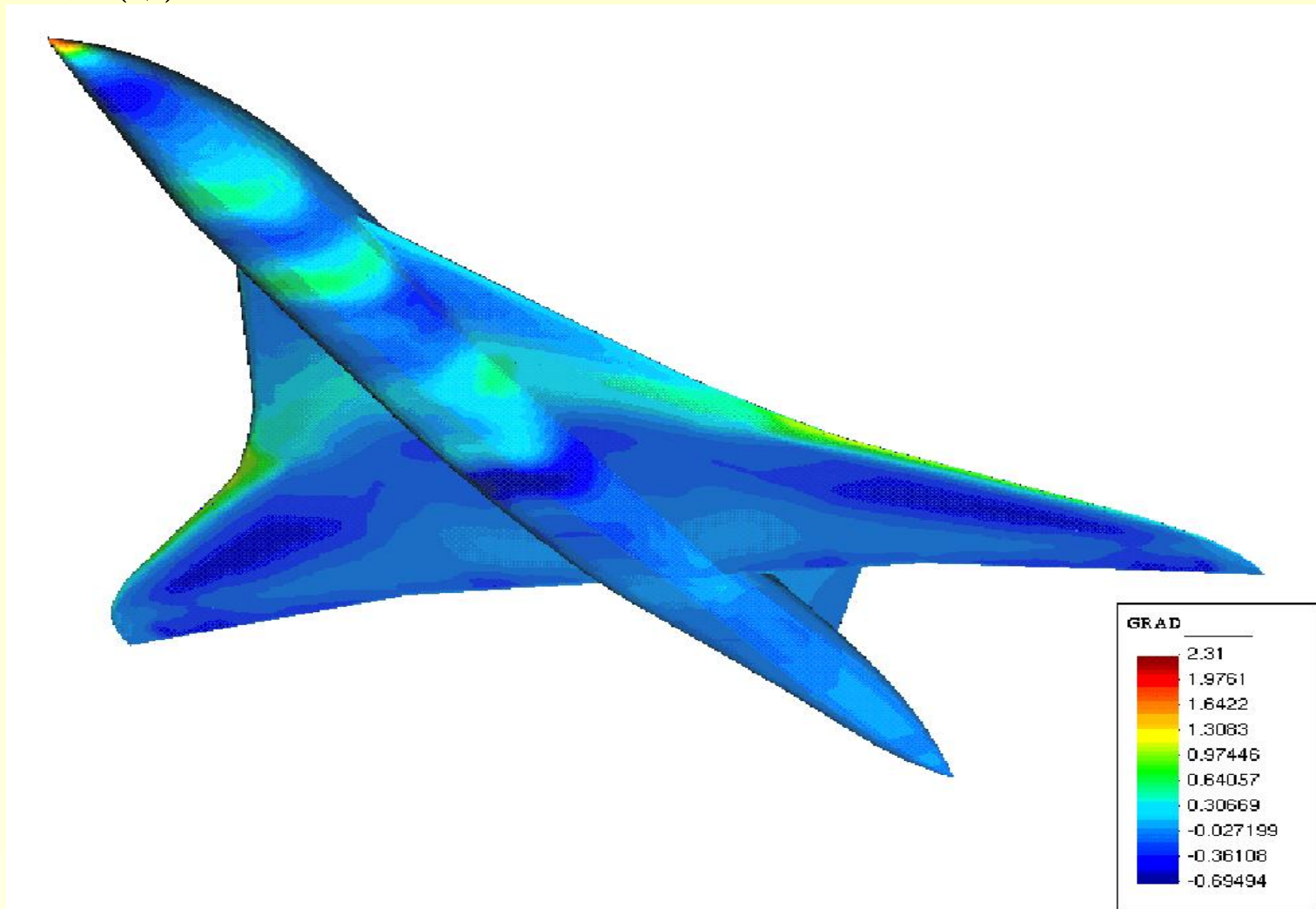- update $\gamma$ (using transpiration conditions)

# Overall Optimization Loop

Loop:

- compute $\Pi$ with a matrix-free solver
- use $\Pi$ to compute $j'(\gamma)$
- 1-D search in the direction $j'(\gamma)$
- update $\gamma$ (using transpiration conditions)

Performances:

Assembling $\dfrac{\partial \Psi}{\partial W}(\gamma, W))^t$ . $\Pi$ takes about 7 times as long as assembing the state residual $\Psi(\gamma, W))$

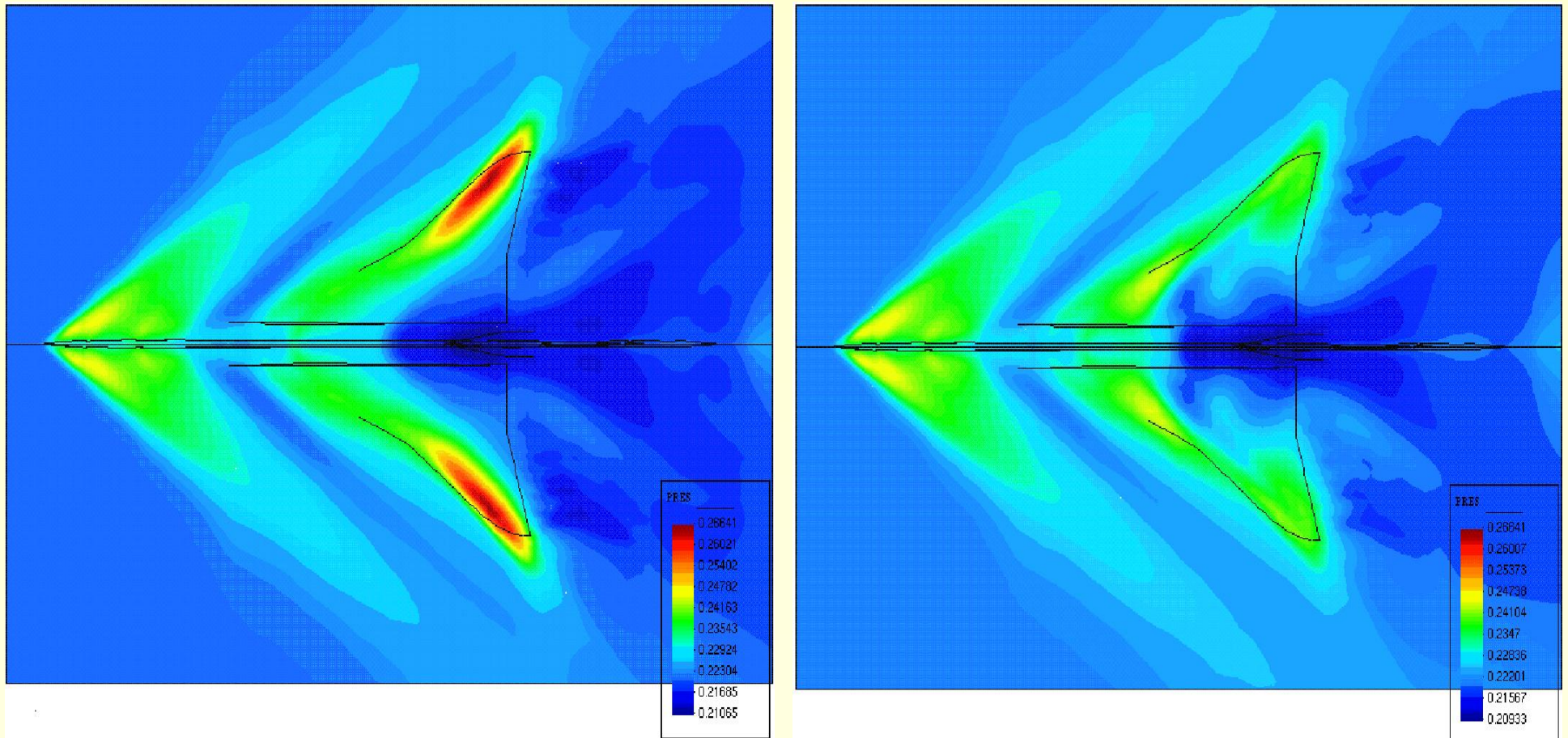$\Rightarrow$ Solving for $\Pi$ takes about 4 times as long as solving for $W$.

# Numerical results 1

Gradient $j'(\gamma)$ on the skin of the plane:

# Numerical results 2

Improvement of the sonic boom after 8 optimization cycles:

# PART 2:

# REVERSE AD TO COMPUTE THE ADJOINT

- Some principles of Reverse AD

- The "Tape" memory problem, the "Checkpointing" tactic

- Impact of Data Dependences Analysis

- Impact of In-Out Data Flow Analysis

# Principles of reverse AD

AD rewrites source programs to make them compute derivatives.

consider: $\qquad P: \ \{I_1; I_2; \ldots I_p; \}$ implementing $f : \mathbf{R}^m \to \mathbf{R}^n$

# Principles of reverse AD

AD rewrites source programs to make them compute derivatives.

consider: $\quad\quad\quad\quad P: \;\; \{I_1; I_2; \ldots I_p; \} \;\;$ implementing $f: \mathbf{IR}^m \to \mathbf{IR}^n$

identify with: $\quad\quad f = \;\; f_p \circ f_{p-1} \circ \cdots \circ f_1$

name: $\quad\quad\quad\quad x_0 = \;\; x \;$ and $\; x_k = f_k(x_{k-1})$

# Principles of reverse AD

AD rewrites source programs to make them compute derivatives.

consider: $\qquad P: \quad \{I_1; I_2; \ldots I_p; \}$ implementing $f : \mathbf{IR}^m \to \mathbf{IR}^n$

identify with: $\qquad f = \quad f_p \circ f_{p-1} \circ \cdots \circ f_1$

name: $\qquad x_0 = \quad x$ and $x_k = f_k(x_{k-1})$

chain rule: $\qquad f'(x) = \quad f'_p(x_{p-1}).f'_{p-1}(x_{p-2}).\ldots.f'_1(x_0)$

# Principles of reverse AD

AD rewrites source programs to make them compute derivatives.

consider: $\quad\quad\quad\quad P: \{I_1; I_2; \ldots I_p;\}$ implementing $f : \mathbf{IR}^m \to \mathbf{IR}^n$

identify with: $\quad\quad\quad f = f_p \circ f_{p-1} \circ \cdots \circ f_1$

name: $\quad\quad\quad\quad x_0 = x$ and $x_k = f_k(x_{k-1})$

chain rule: $\quad\quad f'(x) = f'_p(x_{p-1}).f'_{p-1}(x_{p-2}).\ldots.f'_1(x_0)$

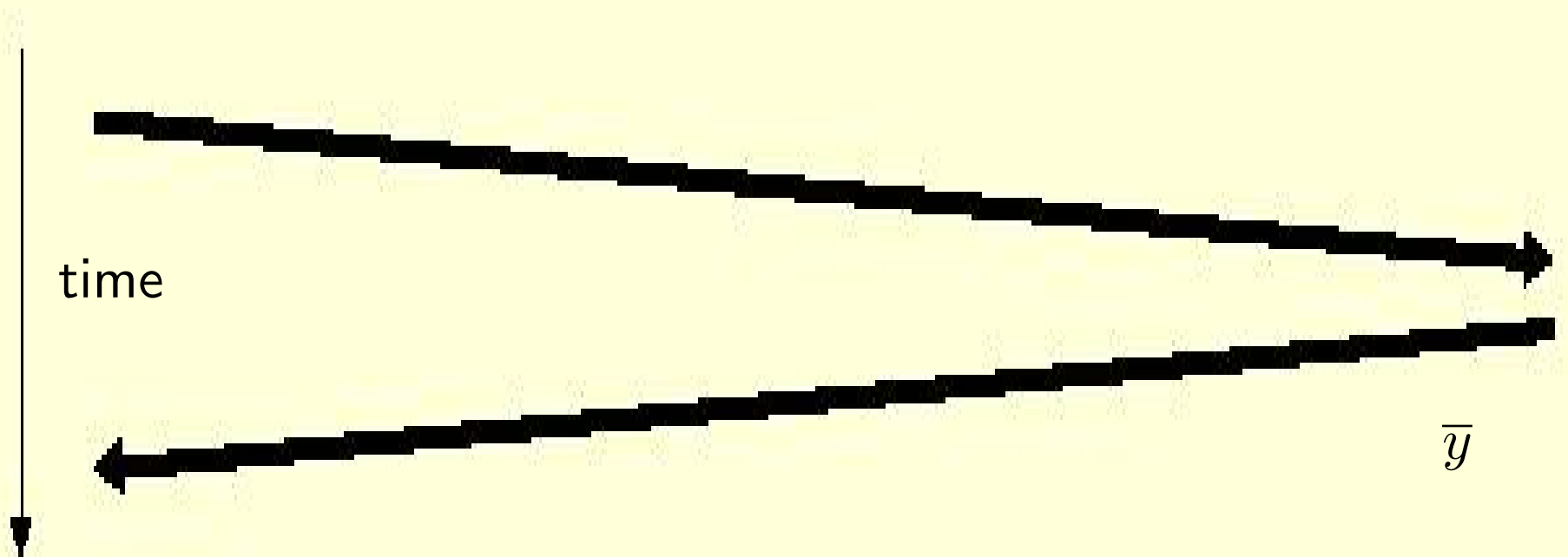$f'(x)$ generally too large and expensive $\Rightarrow$ take useful views!

$$\dot{y} = f'(x).\dot{x} = f'_p(x_{p-1}).f'_{p-1}(x_{p-2}).\ldots.f'_1(x_0).\dot{x} \quad \text{tangent AD}$$
$$\overline{x} = f'^*(x).\overline{y} = f'^*_1(x_0).\ldots.f'^*_{p-1}(x_{p-2}).f'^*_p(x_{p-1}).\overline{y} \quad \text{reverse AD}$$

Evaluate both from right to left !

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0). \ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$



time

$\overline{y}$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0). \ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

time

$$f_p'^*(x_{p-1}). \qquad \overline{y}$$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0)..... f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$



time

$$f_{p-1}'^*(x_{p-2}). \quad f_p'^*(x_{p-1}). \quad \overline{y}$$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0)\ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

time

$$f_{p-1}'^*(x_{p-2}).\quad f_p'^*(x_{p-1}).\quad \overline{y}$$

$$\overline{x} = f_1'^*(x_0).\qquad \cdots$$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0). \ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

time

$$x_{p-1} = f_{p-1}(x_{p-2});$$

$$f_{p-1}'^*(x_{p-2}). \quad f_p'^*(x_{p-1}). \quad \overline{y}$$

$$\overline{x} = f_1'^*(x_0). \qquad \cdots$$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0)\ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

$$x_{p-2} = f_{p-2}(x_{p-3});$$

$$x_{p-1} = f_{p-1}(x_{p-2});$$

time

$$f_{p-1}'^*(x_{p-2}). \quad f_p'^*(x_{p-1}). \quad \overline{y}$$

$$\overline{x} = f_1'^*(x_0). \qquad \cdots$$

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0).\ldots.f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

$x_0;$

$x_1 = f_1(x_0);$     *forward sweep*

$\ldots$     $x_{p-2} = f_{p-2}(x_{p-3});$

$x_{p-1} = f_{p-1}(x_{p-2});$

time

$f_{p-1}'^*(x_{p-2}).$   $f_p'^*(x_{p-1}).$   $\overline{y}$

$\overline{x} = f_1'^*(x_0).$    $\ldots$

*backward sweep*

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0)\ldots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

$x_0;$

$x_1 = f_1(x_0);$     *forward sweep*

$\cdots$    $x_{p-2} = f_{p-2}(x_{p-3});$

$x_{p-1} = f_{p-1}(x_{p-2});$

time

retrieve

$f_{p-1}'^*(x_{p-2}).$   $f_p'^*(x_{p-1}).$   $\overline{y}$

$\overline{x} = f_1'^*(x_0).$    $\cdots$

*backward sweep*

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0).\dots f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$

$x_0;$

$x_1 = f_1(x_0);$

*forward sweep*

$\dots$ $\quad x_{p-2} = f_{p-2}(x_{p-3});$

$x_{p-1} = f_{p-1}(x_{p-2});$

time

retrieve $\qquad$ retrieve

$\overline{y}$

$\overline{x} = f_1'^*(x_0).$ $\quad \dots$ $\qquad f_{p-1}'^*(x_{p-2}).$ $\quad f_p'^*(x_{p-1}).$

*backward sweep*

# Reverse AD is more tricky than Tangent AD

$$\overline{x} = f'^*(x).\overline{y} = f_1'^*(x_0).\ldots.f_{p-1}'^*(x_{p-2}).f_p'^*(x_{p-1}).\overline{y}$$



$x_0;$

$x_1 = f_1(x_0);$

*forward sweep*

$x_{p-2} = f_{p-2}(x_{p-3});$

$x_{p-1} = f_{p-1}(x_{p-2});$

time

retrieve

retrieve

retrieve

$\overline{x} = f_1'^*(x_0).$ $\cdots$ $f_{p-1}'^*(x_{p-2}).$ $f_p'^*(x_{p-1}).$ $\overline{y}$

*backward sweep*

# Memory usage ("Tape") is the bottleneck!

# Example:

Program fragment:

$$\ldots$$
$$v_2 \;=\; 2 * v_1 \;+\; 5$$
$$v_4 \;=\; v_2 \;+\; p_1 * v_3/v_2$$
$$\ldots$$

# Example:

Program fragment:

$$
\begin{aligned}
&\ldots \\
v_2 &= 2 * v_1 + 5 \\
v_4 &= v_2 + p_1 * v_3 / v_2 \\
&\ldots
\end{aligned}
$$

Corresponding transposed Partial Jacobians:

$$
f'^*(x) = \ldots
\begin{pmatrix}
1 & 2 & & \\
 & 0 & & \\
 & & 1 & \\
 & & & 1
\end{pmatrix}
\begin{pmatrix}
1 & & & 0 \\
 & 1 & & 1-\frac{p_1 * v_3}{v_2^2} \\
 & & 1 & \frac{p_1}{v_2} \\
 & & & 0
\end{pmatrix}
\ldots
$$

# Reverse mode on the example

...

$$\bar{v}_2 \ = \ \bar{v}_2 \ + \ \bar{v}_4 * (1 - p_1 * v_3/v_2^2)$$
$$\bar{v}_3 \ = \ \bar{v}_3 \ + \ \bar{v}_4 * p_1/v_2$$
$$\bar{v}_4 \ = \ 0$$

$$\bar{v}_1 \ = \ \bar{v}_1 \ + \ 2 * \bar{v}_2$$
$$\bar{v}_2 \ = \ 0$$

...

# Reverse mode on the example

...

$$v_2 \; = \; 2 * v_1 \; + \; 5$$

$$v_4 \; = \; v_2 \; + \; p_1 * v_3/v_2$$

...

---

...

$$\bar{v}_2 \; = \; \bar{v}_2 \; + \; \bar{v}_4 * (1 - p_1 * v_3/v_2^2)$$
$$\bar{v}_3 \; = \; \bar{v}_3 \; + \; \bar{v}_4 * p_1/v_2$$
$$\bar{v}_4 \; = \; 0$$

$$\bar{v}_1 \; = \; \bar{v}_1 \; + \; 2 * \bar{v}_2$$
$$\bar{v}_2 \; = \; 0$$

...

# Reverse mode on the example

...

Push($v_2$)

$v_2 \;=\; 2 * v_1 \;+\; 5$

Push($v_4$)

$v_4 \;=\; v_2 \;+\; p_1 * v_3 / v_2$

...

---

...

Pop($v_4$)

$\bar{v}_2 \;=\; \bar{v}_2 \;+\; \bar{v}_4 * (1 - p_1 * v_3 / v_2^2)$

$\bar{v}_3 \;=\; \bar{v}_3 \;+\; \bar{v}_4 * p_1 / v_2$

$\bar{v}_4 \;=\; 0$

Pop($v_2$)

$\bar{v}_1 \;=\; \bar{v}_1 \;+\; 2 * \bar{v}_2$

$\bar{v}_2 \;=\; 0$

...

# Reverse mode on our application

From subroutine `Psi` :
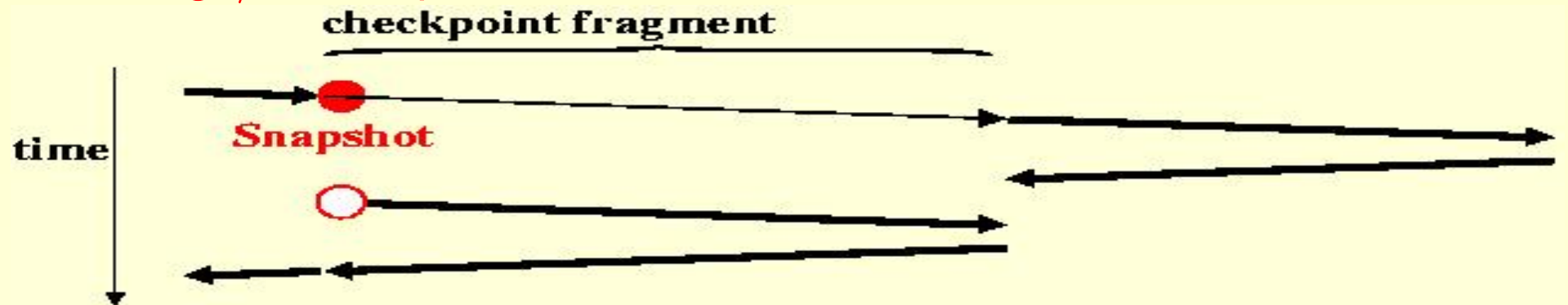
`Psi`: $\gamma, W \mapsto \Psi(\gamma, W),$

Use reverse AD to build subroutine $\overline{\texttt{Psi}}$ :

$\overline{\texttt{Psi}}$: $\gamma, W, \overline{\Psi} \mapsto \dfrac{\partial \Psi}{\partial W}(\gamma, W))^t . \overline{\Psi}$

# Reverse mode on our application

From subroutine `Psi` :

`Psi:` $\gamma, W \mapsto \Psi(\gamma, W),$

Use reverse AD to build subroutine $\overline{\texttt{Psi}}$ :

$\overline{\texttt{Psi}}:$ $\gamma, W, \overline{\Psi} \mapsto \dfrac{\partial \Psi}{\partial W}(\gamma, W))^t . \overline{\Psi}$

But the Tape grows too large on large meshes!

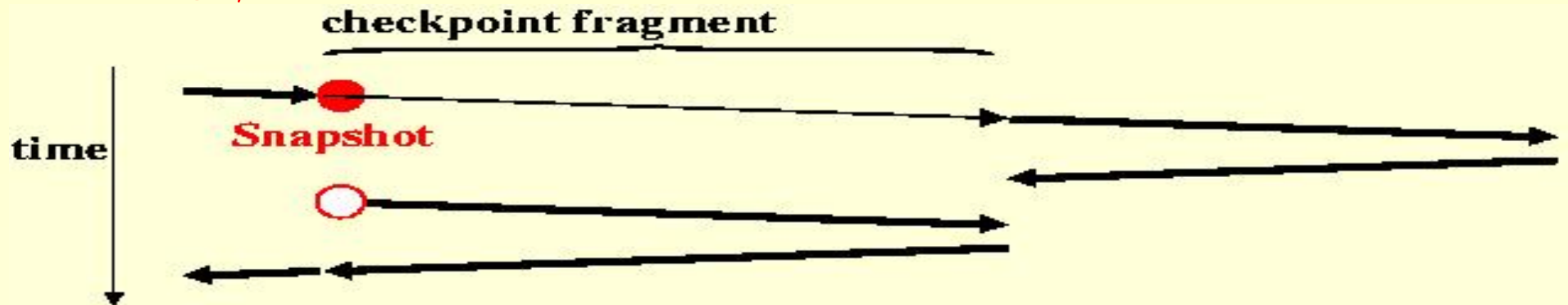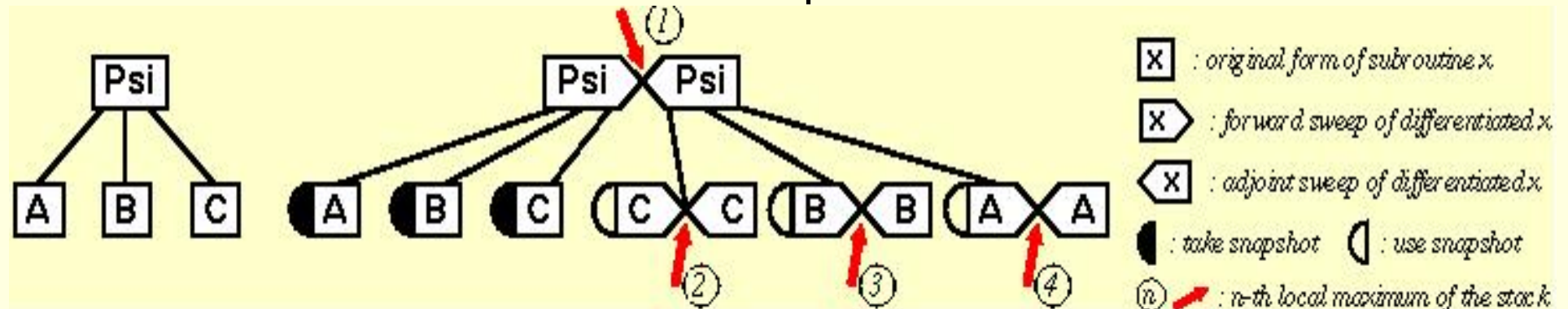# The Checkpointing mechanism

A Storage/Recomputation tradeoff:



checkpoint fragment

Snapshot

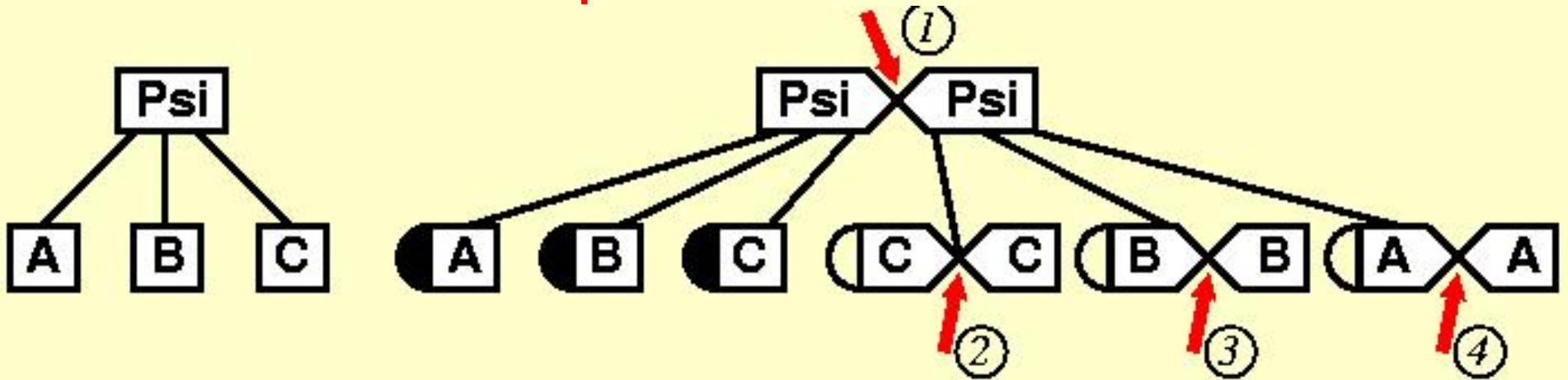time

# The Checkpointing mechanism

A Storage/Recomputation tradeoff:



TAPENADE does it on the Call Graph :



Tape size reaches 4 local maxima.

# Tape size maxima



| Tape local maximum # | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| | 12.40 | 12.37 | 13.60 | 9.66 |

773 R*8/node: (still) too expensive in memory
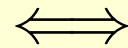
$\Rightarrow$ Use Data Flow analysis

# Data-Flow "to the rescue" (1)

Data Dependence Graphs of P and $\overline{\text{P}}$ are isomorphic, so...
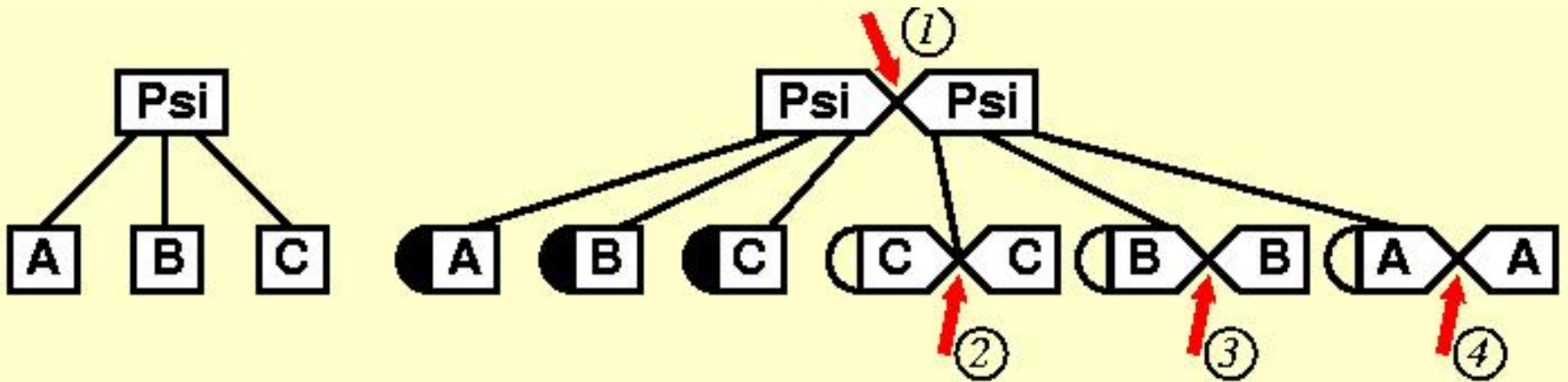improve Independent-Iterations loops ("$II$-loops")

Standard:

```
do // i = 1,N
    body(i)
end
do i = N,1
    body(i)
end
```

$\Longleftrightarrow$

Improved:

```
do i = 1,N
    body(i)
    body(i)
end
```

where in the Standard box the second $body(i)$ has an overline, and in the Improved box the second $body(i)$ has an overline: $\frac{body(i)}{\overline{body(i)}}$

| Tape local maximum # | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| No modification: | 12.40 | 12.37 | 13.60 | 9.66 |
| *II*-loops improvement: | 12.38 | 7.98 | 4.10 | 0.02 |

Improvement on (4), but hidden by (1) !

# Data-Flow "to the rescue" (2)

Analyze the size of Snapshots:



Snapshot = **IN**(**checkpoint**) $\bigcap$ **OUT**(**checkpoint and after**)

| Tape local maximum # | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| No modification: | 12.40 | 12.37 | 13.60 | 9.66 |
| Only snapshot reduction: | 1.02 | 0.85 | 9.70 | 9.33 |
| Only $II$-loops improvement: | 12.38 | 7.98 | 4.10 | 0.02 |
| Both improvements: | 1.02 | 0.61 | 0.22 | 0.02 |

58 R*8/node: quite acceptable !

# CONCLUSION:

- **Part 1:** Brute-force reverse AD, including on the iterative solver, is a hazardous strategy $\Rightarrow$ define a manual strategy *before* AD.

- **Part 1:** ... but the matrix-free solver proves a delicate step.

- **Part 2:** Reverse AD can use reasonable memory space, thanks to data flow analyses.

- **Advertisment!** Tapenade  : an AD tool on the web
  http://tapenade.inria.fr:8080/tapenade
  or alternatively FTP from our web site
  http://www.inria.fr/tropics.