

Advanced Algorithmic Differentiation: Memory-reduced Calculation of Adjoint

Andrea Walther and Sri Tadinada
Institut für Mathematik
Humboldt-Universität zu Berlin

CWI Autumn School 2025
Artificial Intelligence in Natural Sciences

Outline

- 1 Motivation
- 2 Basics of Checkpointing
- 3 Summary

Starting Point: Adjoints are Great!

Starting Point: Adjointns are Great!

BUT

$$\text{MEM}(\bar{y}^\top F'(x)) \sim \text{OPS}(F(x))$$

independent of chosen approach!

Starting Point: Adjoint are Great!

BUT

$$\text{MEM}(\bar{y}^\top F'(x)) \sim \text{OPS}(F(x))$$

independent of chosen approach!

Alternatives:

- no assumptions on structure of $F(x)$
⇒ For AD tools: subroutine-oriented checkpointing
see, e.g., Tapenade and OpenAD

Starting Point: Adjoint are Great!

BUT

$$\text{MEM}(\bar{y}^\top F'(x)) \sim \text{OPS}(F(x))$$

independent of chosen approach!

Alternatives:

- no assumptions on structure of $F(x)$
 - ⇒ For AD tools: subroutine-oriented checkpointing
see, e.g., Tapenade and OpenAD
- $F(x)$ = time-stepping procedure
 - transient procedure, i.e., real time-dependent process
 - ⇒ checkpointing
 - pseudo time-stepping, e.g., aerodynamics, spin-up
 - = fixed point iterations
 - ⇒ adapted derivative calculation

Starting Point: Adjointns are Great!

BUT

$$\text{MEM}(\bar{y}^\top F'(x)) \sim \text{OPS}(F(x))$$

independent of chosen approach!

Alternatives:

- no assumptions on structure of $F(x)$
 - ⇒ For AD tools: subroutine-oriented checkpointing
see, e.g., Tapenade and OpenAD
- $F(x)$ = time-stepping procedure
 - transient procedure, i.e., real time-dependent process
 - ⇒ checkpointing
 - pseudo time-stepping, e.g., aerodynamics, spin-up
 - = fixed point iterations
 - ⇒ adapted derivative calculation
- $F(x)$ = iterative solve
 - ⇒ adapted derivative calculation

Starting Point: Adjointns are Great!

BUT

$$\text{MEM}(\bar{y}^\top F'(x)) \sim \text{OPS}(F(x))$$

independent of chosen approach!

Alternatives:

- no assumptions on structure of $F(x)$
 - ⇒ For AD tools: subroutine-oriented checkpointing
see, e.g., Tapenade and OpenAD
- $F(x)$ = time-stepping procedure
 - transient procedure, i.e., real time-dependent process
 - ⇒ **checkpointing**
 - pseudo time-stepping, e.g., aerodynamics, spin-up
 - = fixed point iterations
 - ⇒ adapted derivative calculation
- $F(x)$ = iterative solve
 - ⇒ adapted derivative calculation

Gradient Calculation (Time Integration)

- ① Forward integration: For given x_0

$$x_{i+1} = F_i(x_i, u_i), \quad i = 0, \dots, l-1$$

with x_i state, u_i control

number of time steps l might be not known a priori

- ② Evaluation of target function $J(x, u)$

- ③ Backward integration:

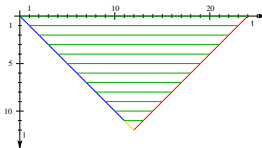
$$\bar{x}_{i-1} = \bar{F}_i(\bar{x}_i, x_{i-1}, u_{i-1}, u_i), \quad i = l, \dots, 1$$

AD, discretization of continuous adjoint, hand coded, ...

- ④ Gradient calculation

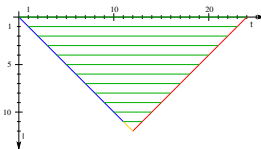
Idea of Checkpointing

Store-Everything Approach:

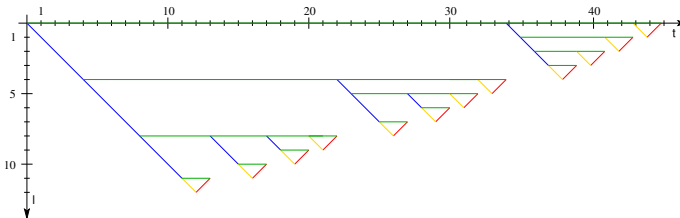


Idea of Checkpointing

Store-Everything Approach:



Checkpointing:



Checkpointing Theory

Goal: Minimal number of recomputations for c checkpoints

Available results:

- / known, constant step costs
(Griewank 1992), (Griewank, Walther 2000), (Kowarz, Walther 2007)
- / known, variable step costs
(Walther 2000), (Hinze, Sternberg 2005)

Checkpointing Theory

Goal: Minimal number of recomputations for c checkpoints

Available results:

- f known, constant step costs
(Griewank 1992), (Griewank, Walther 2000), (Kowarz, Walther 2007)
- f known, variable step costs
(Walther 2000), (Hinze, Sternberg 2005)
- f unknown, constant step costs
(Hinze, Walther, Sternberg 2006), (Stumm, Walther 2010), (Moin, Wang 2010)

Checkpointing Theory

Goal: Minimal number of recomputations for c checkpoints

Available results:

- f known, constant step costs
(Griewank 1992), (Griewank, Walther 2000), (Kowarz, Walther 2007)
- f known, variable step costs
(Walther 2000), (Hinze, Sternberg 2005)
- f unknown, constant step costs
(Hinze, Walther, Sternberg 2006), (Stumm, Walther 2010), (Moin, Wang 2010)
- f known, variable access cost
(Stumm, Walther 2009), (Aupy, Hovland et al. 2015), (Schanen, Marin et al. 2015)

Checkpointing Theory

Goal: Minimal number of recomputations for c checkpoints

Available results:

- f known, constant step costs
(Griewank 1992), (Griewank, Walther 2000), (Kowarz, Walther 2007)
- f known, variable step costs
(Walther 2000), (Hinze, Sternberg 2005)
- f unknown, constant step costs
(Hinze, Walther, Sternberg 2006), (Stumm, Walther 2010), (Moin, Wang 2010)
- f known, variable access cost
(Stumm, Walther 2009), (Aupy, Hovland et al. 2015), (Schanen, Marin et al. 2015)

Also applicable for continuous adjoints!

Checkpointing Theory

Goal: Minimal number of recomputations for c checkpoints

Available results:

- / **known, constant step costs**
(Griewank 1992), (Griewank, Walther 2000), (Kowarz, Walther 2007)
- / known, variable step costs
(Walther 2000), (Hinze, Sternberg 2005)
- / unknown, constant step costs
(Hinze, Walther, Sternberg 2006), (Stumm, Walther 2010), (Moin, Wang 2010)
- / known, variable access cost
(Stumm, Walther 2009), (Aupy, Hovland et al. 2015), (Schanen, Marin et al. 2015)

Also applicable for continuous adjoints!

Theory for Binomial Checkpointing I

Theorem (Complexity Result I)

Given

- l = *total number of time steps to be reversed*
- c = *number of checkpoints that can be used*
- r = *repetition number, i.e., unique integer satisfying*

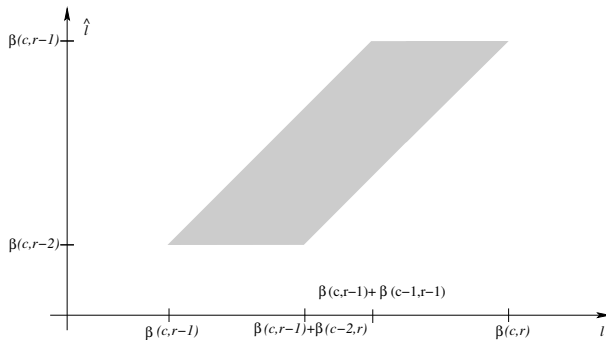
$$\beta(c, r-1) < l \leq \beta(c, r) \equiv \binom{c+r}{c}$$

Then: The minimal number of time step evaluations required equals

$$t(c, l) = rl - \beta(c+1, r-1) + 1$$

Theory for Binomial Checkpointing I

Place of next checkpoint \hat{l} :



Theory for Binomial Checkpointing II

Theorem (Complexity Result II)

Given

- l = total number of time steps to be reversed
- c = number of checkpoints that can be used
- r = repetition number, i.e., unique integer satisfying

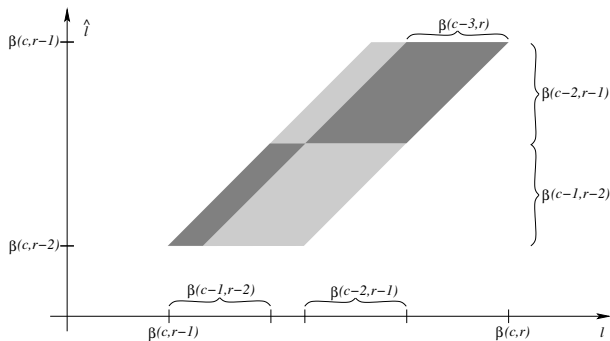
$$\beta(c, r-1) < l \leq \beta(c, r) \equiv \binom{c+r}{c}$$

Then: Among all checkpoint schedules satisfying the previous theorem the minimal number of times a current state is saved onto the stack of checkpoints is given by

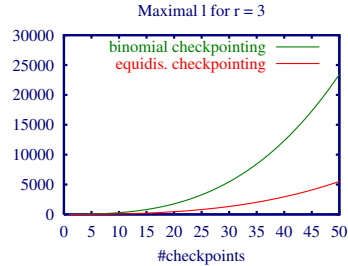
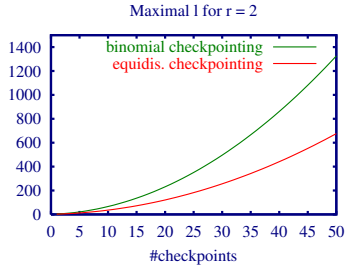
$$q(l, r) = \begin{cases} \beta(c-1, r-1) & \text{if } l \leq \beta(c, r-1) + \beta(c-1, r-1) \\ l - \beta(c, r-1) & \text{if } l \geq \beta(c, r-1) + \beta(c-1, r-1) \end{cases}$$

Theory for Binomial Checkpointing II

Place of next checkpoint \hat{l} :



Repetition Number r vs Time Step Number l



[Walther, Griewank 2004]

The Checkpointing Routine `revolve()`

...

do

 whatodo = `revolve()`

 switch(whatodo)

 case advance: for $\text{old_state} < i \leq \text{next_state}$ `forward(x)`

 case takeshot: `store(x, xstore, current_checkpoint)`

 case firsturn: `init(bar_x, bar_u)`

`reverse(bar_x, bar_u)`

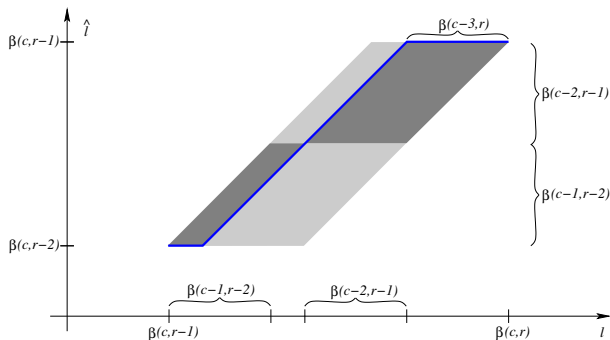
 case youturn: `reverse(bar_x, bar_u)`

 case restore: `restore(x, xstore, current_check)`

while(whatodo \neq terminate)

Current Strategy Chosen by `revolve()`

Place of next checkpoint \hat{l} :



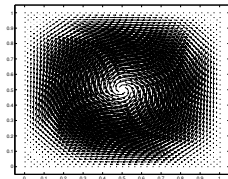
Test Case: Cavity Flow

joint work with Michael Hinze

Cost function of tracking typ, initial state:

$$y_0(x) = e \begin{bmatrix} (\cos 2\pi x_1 - 1) \sin 2\pi x_2 \\ -(\cos 2\pi x_2 - 1) \sin 2\pi x_1 \end{bmatrix},$$

$\Omega := (0, 1) \times (0, 1)$, $T = 1$, $Re = 10$, $e = \text{Euler number}$



Test Case: Cavity Flow

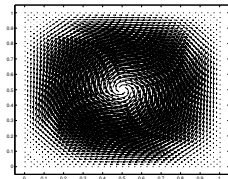
joint work with Michael Hinze

Cost function of tracking typ, initial state:

$$y_0(x) = e \begin{bmatrix} (\cos 2\pi x_1 - 1) \sin 2\pi x_2 \\ -(\cos 2\pi x_2 - 1) \sin 2\pi x_1 \end{bmatrix},$$

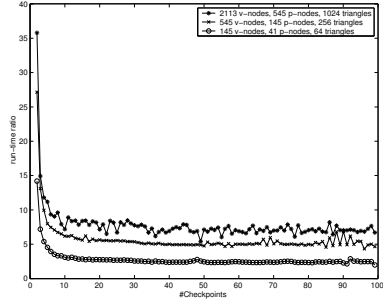
$\Omega := (0, 1) \times (0, 1)$, $T = 1$, $\text{Re} = 10$, $e = \text{Euler number}$

Desired state:

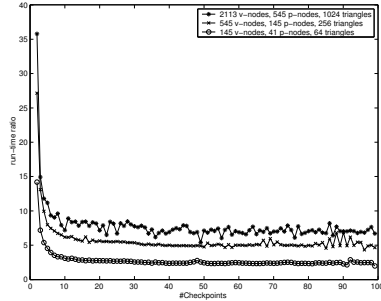


$$\begin{aligned} z(x, t) &= \begin{bmatrix} \varphi_{x_2}(x, t) \\ -\varphi_{x_1}(x, t) \end{bmatrix}, \\ \varphi(x, t) &= \theta(x_1, t)\theta(x_2, t), \\ \theta(y, t) &= (1 - y)^2(1 - \cos kt), \\ y &\in [0, 1]. \end{aligned}$$

Runtime Ratios, 1000 Time Steps



Runtime Ratios, 1000 Time Steps



- Run-time ratio between 2 and 7.
- Only a very small number of checkpoints required!

Summary

- Adjoint computation may result in enormous memory requirement
- Exploit structure!

Summary

- Adjoint computation may result in enormous memory requirement
- Exploit structure!
- Time-dependent problem:
Implicit time-stepping ? Unknown number of time steps?
Apply adapted checkpointing strategy

Summary

- Adjoint computation may result in enormous memory requirement
- Exploit structure!
- Time-dependent problem:
Implicit time-stepping ? Unknown number of time steps?
Apply adapted checkpointing strategy
- Pseudo time-stepping:
Two-phase approach “easy” to apply
One-shot approach more efficient but also more complicated

Summary

- Adjoint computation may result in enormous memory requirement
- Exploit structure!
- Time-dependent problem:
Implicit time-stepping ? Unknown number of time steps?
Apply adapted checkpointing strategy
- Pseudo time-stepping:
Two-phase approach “easy” to apply
One-shot approach more efficient but also more complicated
- Differentiation of Newton-like iterations OK
- Be careful when differentiating solvers like CG, GMRES, ...