

# Rethinking HPC algorithms for Exascale: the case for Adjoint Computations

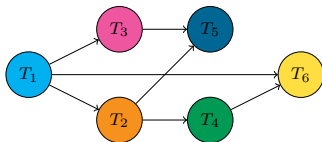
Guillaume Aupy



VANDERBILT  
UNIVERSITY



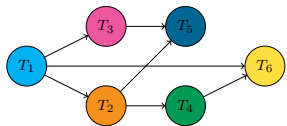
- ▶  $p$  processors (or nodes or cores or computing units).
- ▶ An application represented by a DAG  $\mathcal{G} = (V, E)$ :
  - ▶ Vertices are tasks (or functions to be computed)
  - ▶ Edges represent data dependency (need to be respected)



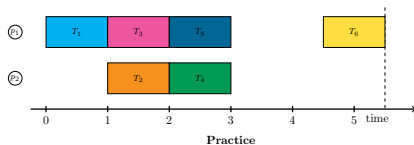
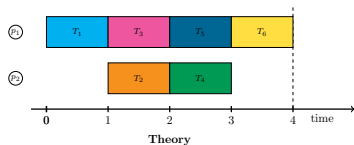




# MOTIVATIONAL EXAMPLE



Assume tasks have unit size and there are two processors.



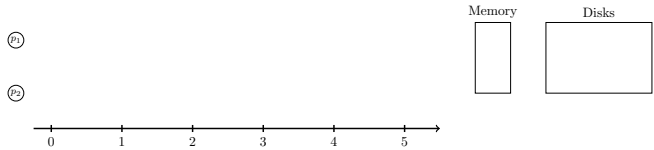
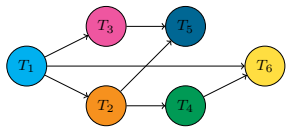
- ▶ In theory, optimal schedule ☺
- ▶ In practice, not what we expected ☹.



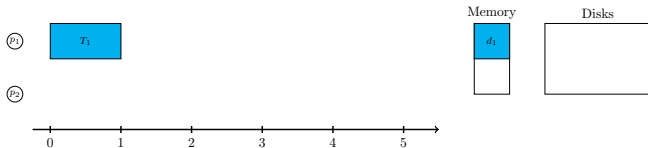
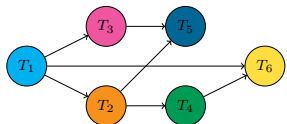




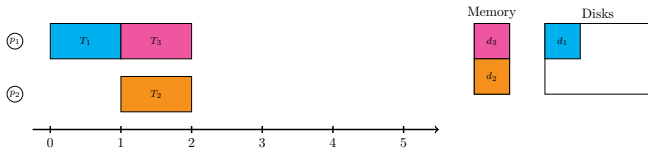
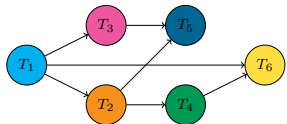
# BACK TO OUR SCHEDULE



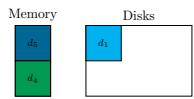
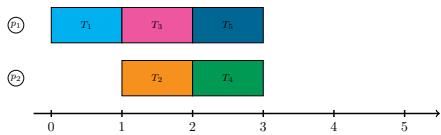
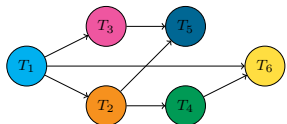
# BACK TO OUR SCHEDULE



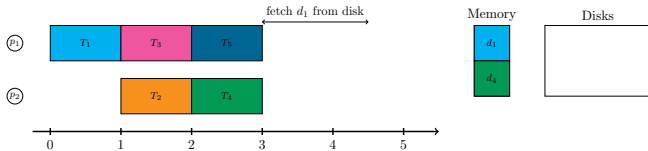
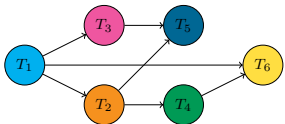
# BACK TO OUR SCHEDULE



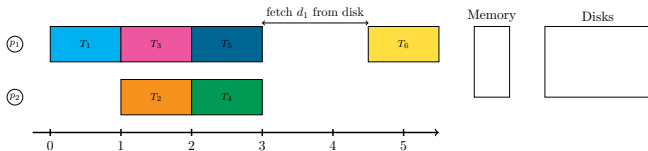
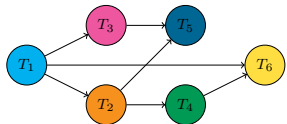
# BACK TO OUR SCHEDULE



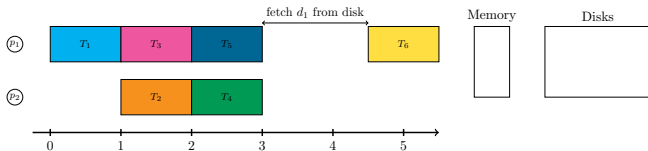
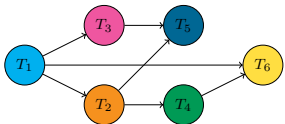
# BACK TO OUR SCHEDULE



# BACK TO OUR SCHEDULE

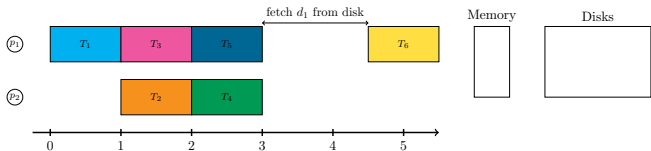
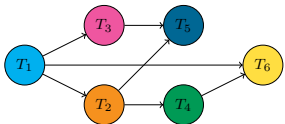


# BACK TO OUR SCHEDULE

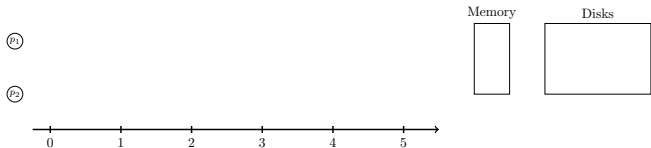


Can we do better (or prove that we cannot)?

# BACK TO OUR SCHEDULE



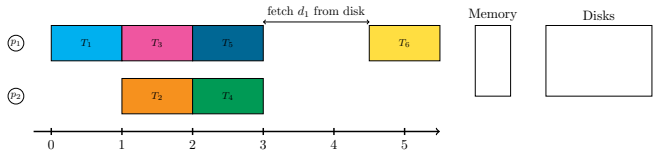
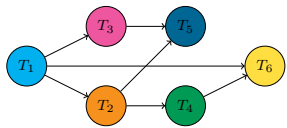
Can we do better (or prove that we cannot)?



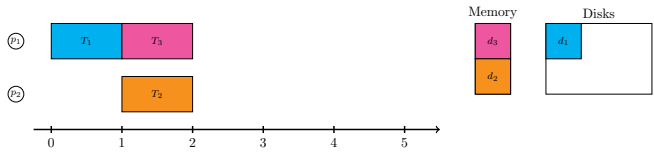




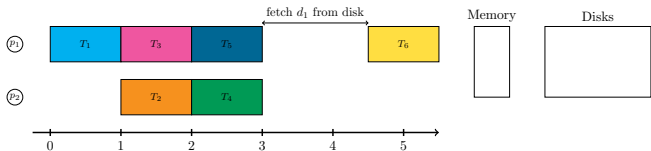
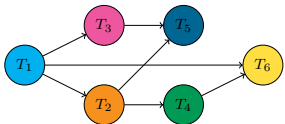
# BACK TO OUR SCHEDULE



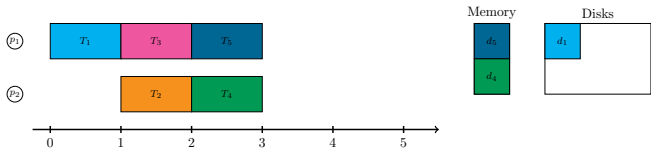
Can we do better (or prove that we cannot)?



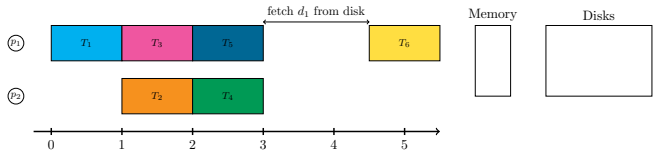
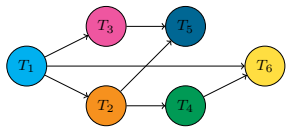
# BACK TO OUR SCHEDULE



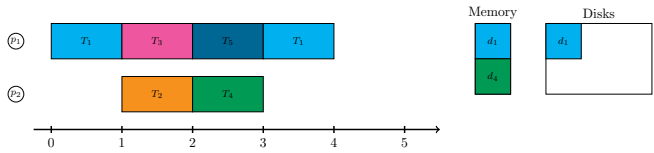
Can we do better (or prove that we cannot)?



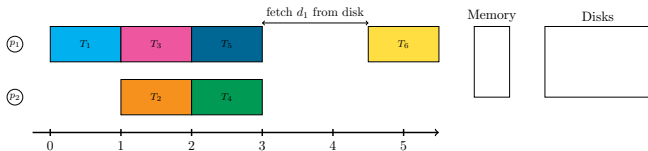
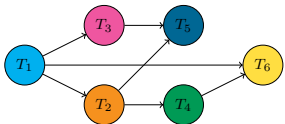
# BACK TO OUR SCHEDULE



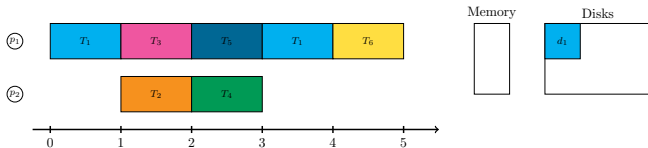
Can we do better (or prove that we cannot)?



# BACK TO OUR SCHEDULE



Can we do better (or prove that we cannot)?



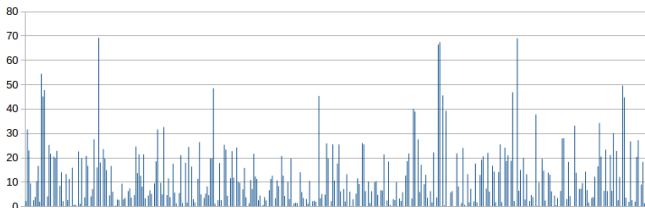




# SOME NUMBERS

(Brief) history of supercomputers at Argonne National Lab:

	Intrepid	Mira
Year	2008-2013	2013-
Peak Perf	0.557 PFlops	10 PFlops
Peak I/O Throughput	88 GB/s	240 GB/s
I/O Tp per Flops	157 GB/PFlops/s	16.8 GB/PFlops/s



Analysis of the Intrepid system @Argonne: I/O throughput decrease (percentage per application, over 400 applications).











## ICE-SHEET MODEL (II)

$$F(u_0) = f_n \circ f_{n-1} \circ \dots \circ f_1 \circ f_0(u_0)$$

$$\nabla F(u_0)\mathbf{y} = Jf_0(u_0)^T \cdot \nabla(f_n \circ f_1)(u_1) \cdot \mathbf{y}$$

$$= Jf_0(u_0)^T \cdot Jf_1(u_1)^T \cdot \dots \cdot Jf_{n-1}(u_{n-1})^T \cdot Jf_n(u_n)^T \cdot \mathbf{y}$$

$Jf^T$  = Transpose Jacobian matrix of  $f$ ;

$u_{i+1}$  =  $f_i(u_i) = f_i(f_{i-1} \circ \dots \circ f_0(u_0))$ .



## A BETTER SOLUTION?

$$\nabla F(u_0)\mathbf{y} = Jf_0(u_0)^T \cdot Jf_1(u_1)^T \cdot \dots \cdot Jf_{n-1}(u_{n-1})^T \cdot Jf_n(u_n)^T \cdot \mathbf{y}$$

```
proc Algo A( $u_0, \mathbf{y}$ )
begin
  Do stuff;
  for  $i = 0$  to  $n$  do
     $u_{i+1} = f_i(u_i)$ ;
    Do stuff;
  end
  Compute  $\nabla F(u_0)\mathbf{y}$ ;
end
```

```
proc Algo B( $u_0, \mathbf{y}$ )
begin
  Do stuff;
  for  $i = 0$  to  $n$  do
     $u_{i+1} = f_i(u_i)$ ;
    Do stuff;
     $v_{i+1} = v_i \cdot Jf_{i+1}(u_{i+1})^T$ ;
  end
end
```





# A BETTER SOLUTION?

$$\nabla F(u_0)\mathbf{y} = Jf_0(u_0)^T \cdot Jf_1(u_1)^T \cdot \dots \cdot Jf_{n-1}(u_{n-1})^T \cdot Jf_n(u_n)^T \cdot \mathbf{y}$$

```

proc Algo A( $u_0, \mathbf{y}$ )
begin
  Do stuff;
  for  $i = 0$  to  $n$  do
     $u_{i+1} = f_i(u_i)$ ;
    Do stuff;
  end
  Compute  $\nabla F(u_0)\mathbf{y}$ ;
end

```

```

proc Algo B( $u_0, \mathbf{y}$ )
begin
  Do stuff;
  for  $i = 0$  to  $n$  do
     $u_{i+1} = f_i(u_i)$ ;
    Do stuff;
     $v_{i+1} = v_i \cdot Jf_{i+1}(u_{i+1})^T$ ;
  end
end

```

What is the problem with Algo B?

$$\nabla F(u_0)\mathbf{y} = \left( \left( \dots \left( Jf_0^T \cdot Jf_1^T \right) \cdot \dots \cdot Jf_{n-1}^T \right) \cdot Jf_n^T \right) \cdot \mathbf{y} \quad n \text{ MatMat ops}$$

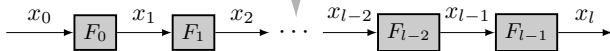
$$\nabla F(u_0)\mathbf{y} = Jf_0^T \cdot \left( Jf_1^T \cdot \dots \cdot \left( Jf_{n-1}^T \cdot \left( Jf_n^T \cdot \mathbf{y} \right) \dots \right) \right) \quad n \text{ MatVec ops}$$



# ADJOINT COMPUTATION

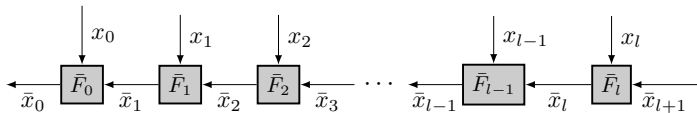
$$F_i(x_i) = x_{i+1} \quad i < l$$

$$\bar{F}_i(x_i, \bar{x}_{i+1}) = \bar{x}_i \quad i \leq l$$



# ADJOINT COMPUTATION

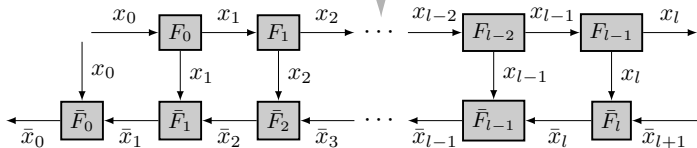
$$F_i(x_i) = x_{i+1} \quad i < l$$
$$\bar{F}_i(x_i, \bar{x}_{i+1}) = \bar{x}_i \quad i \leq l$$



# ADJOINT COMPUTATION

$$F_i(x_i) = x_{i+1} \quad i < l$$

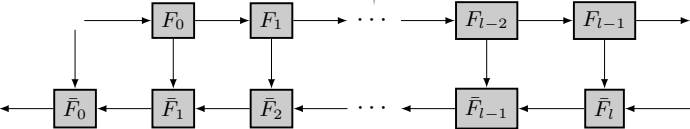
$$\bar{F}_i(x_i, \bar{x}_{i+1}) = \bar{x}_i \quad i \leq l$$



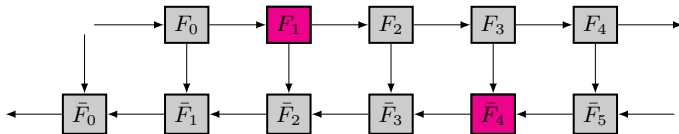
# ADJOINT COMPUTATION

$$F_i(x_i) = x_{i+1} \quad i < l$$

$$\bar{F}_i(x_i, \bar{x}_{i+1}) = \bar{x}_i \quad i \leq l$$

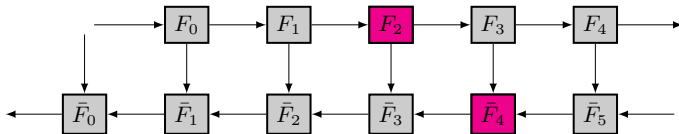


# MODEL OF COMPUTATION



- Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

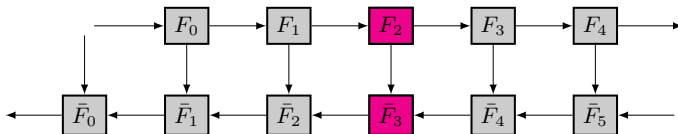
# MODEL OF COMPUTATION



- Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

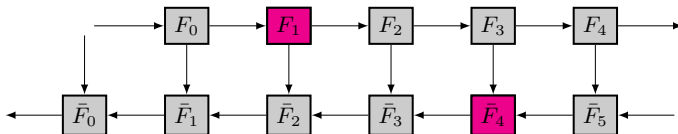


# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

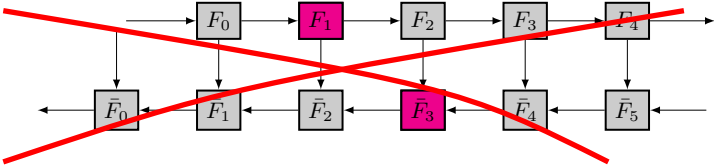
# MODEL OF COMPUTATION



- Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

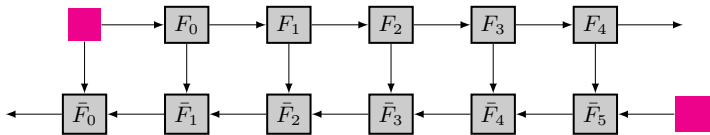


# MODEL OF COMPUTATION



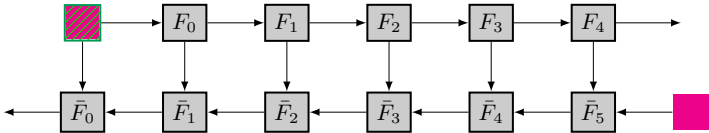
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

# MODEL OF COMPUTATION



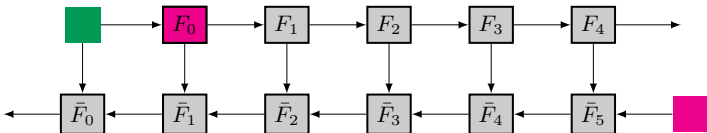
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .

# MODEL OF COMPUTATION



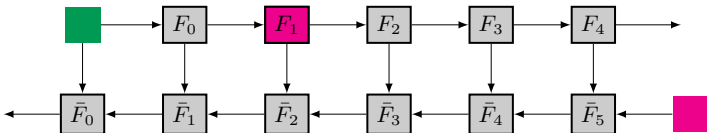
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

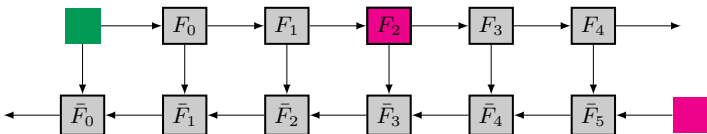
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

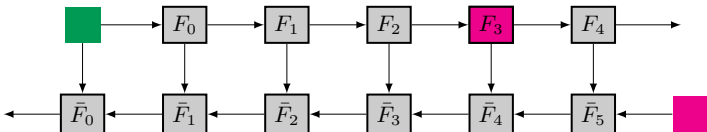


# MODEL OF COMPUTATION



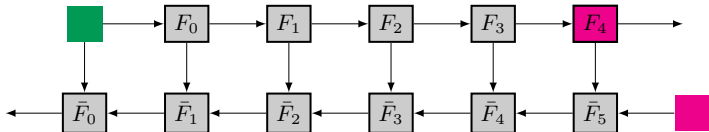
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



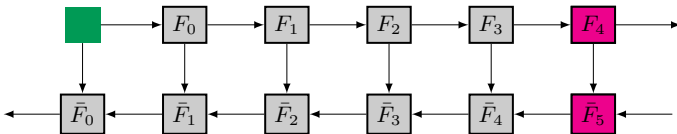
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



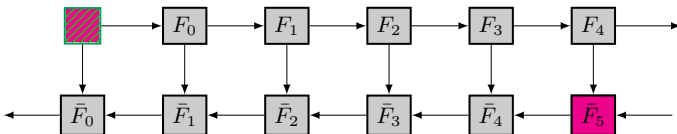
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



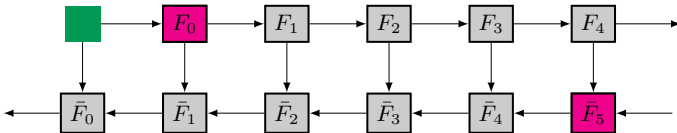
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



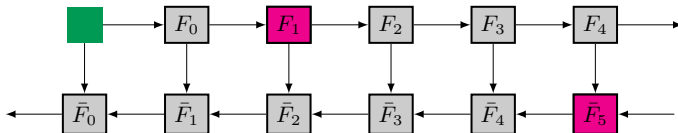
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



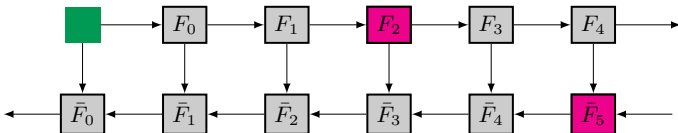
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

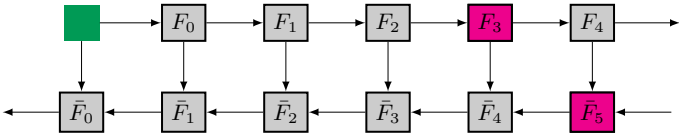
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

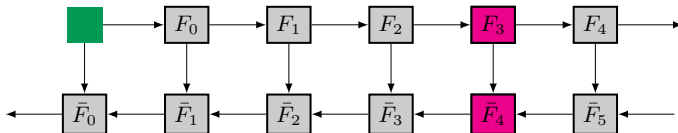


# MODEL OF COMPUTATION



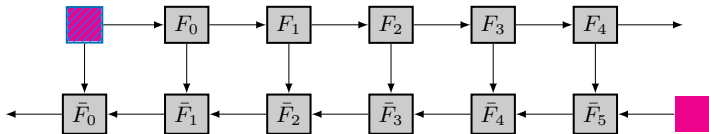
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



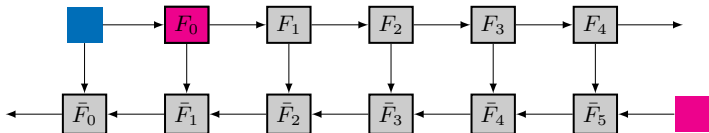
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .

# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

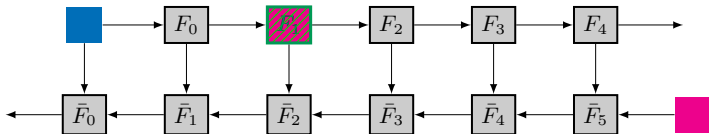
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .



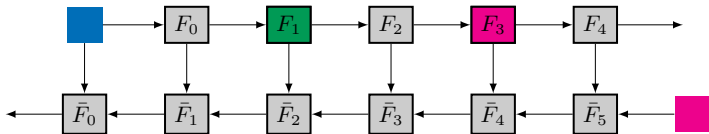
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .



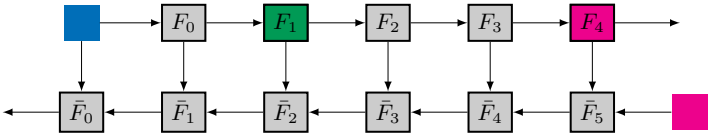
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

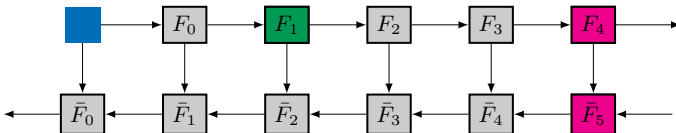


# MODEL OF COMPUTATION



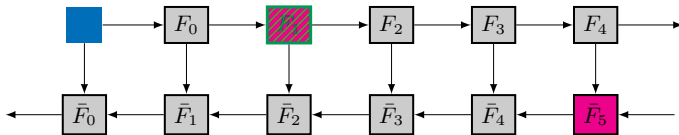
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



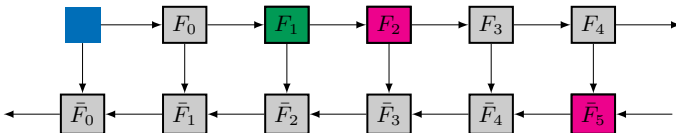
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



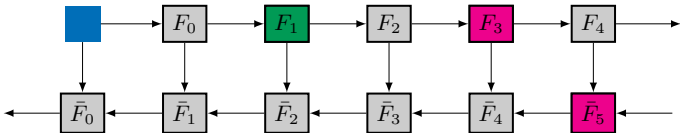
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



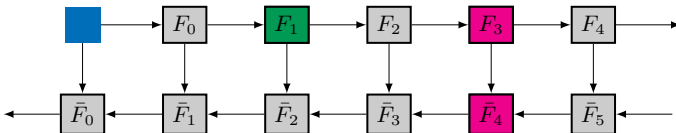
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



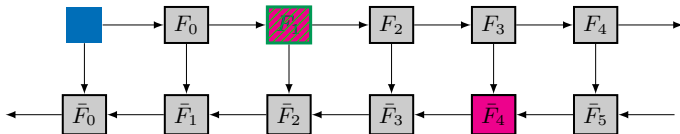
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



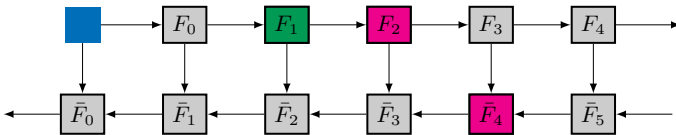
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

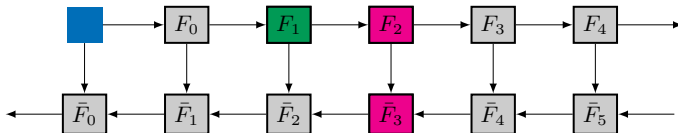
# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

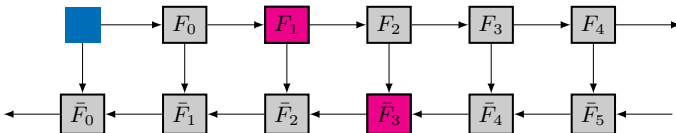


# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# MODEL OF COMPUTATION

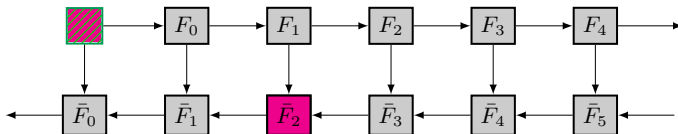


- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .



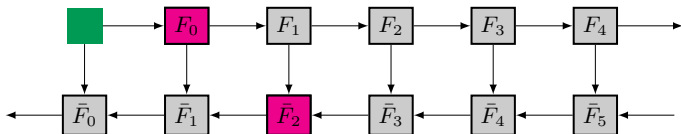


# MODEL OF COMPUTATION



- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

## MODEL OF COMPUTATION



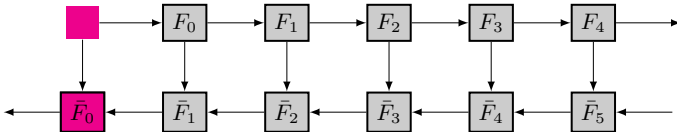
- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .







# MODEL OF COMPUTATION

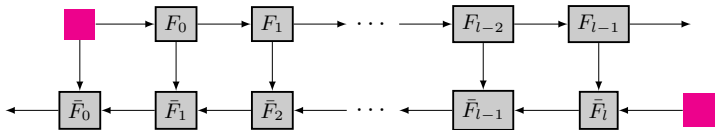


- ▶ Two buffers: store output of computations ( $x_i$  or  $\bar{x}_i$ ). Initial state: contains  $x_0$  and  $\bar{x}_{l+1}$ .
- ▶  $c_m < +\infty$  in-core slots (memory).
  - ▶ Cost to write:  $w_m = 0$ ,
  - ▶ Cost to read:  $r_m = 0$ .
- ▶  $c_d = +\infty$  out-of-core slots (disk).
  - ▶ Cost to write:  $w_d$ ,
  - ▶ Cost to read:  $r_d$ .

# PROBLEM FORMULATION

We want to minimize the makespan of:

		Initial state:
AC graph:	size $l$	
Steps:	$u_f, u_b$	
Memory:	$c_m, w_m = r_m = 0,$	$\mathcal{M}_{\text{ini}} = \emptyset$
Disks:	$c_d = +\infty, w_d, r_d,$	$\mathcal{D}_{\text{ini}} = \emptyset$
Buffers:	$\mathcal{B}^\top, \mathcal{B}^\perp$	$\mathcal{B}_{\text{ini}}^\top = \{x_0\}, \mathcal{B}_{\text{ini}}^\perp = \{\bar{x}_{l+1}\}$



# PREVIOUS WORK

**GW00:** REVOLVE( $l, c_m$ ), optimal algorithm with  $c_m$  memory slots and no disk slots.

**SW09:** SWA\*, an algorithm based on REVOLVE that takes disk storage into account.

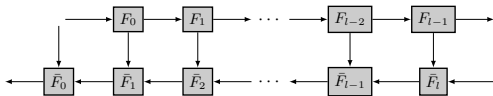
- (i)  $\text{SWA}(l, c_m, c_d, w_d, r_d) \approx \text{REVOLVE}(l, c_d + c_m)$ <sup>1</sup>
- (ii)  $\text{SWA}^*(l, c_m, w_d, r_d) = \min_{c_d=0 \dots l-c_m} \text{SWA}(l, c_m, c_d, w_d, r_d)$

**This work:** optimal algorithm with disk storage.

---

<sup>1</sup>out of the  $c_d + c_m$  slots used by REVOLVE, the  $c_d$  slots the least used are considered disk slots.

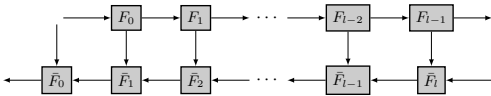
## A GRASP OF THE PROOF



Any algorithm works in two phases:

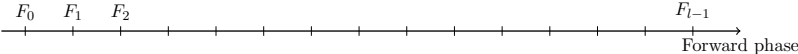
- ▶ The forward phase (before executing  $\bar{F}_l$ );
- ▶ The backward phase (that starts when executing  $\bar{F}_l$ );

# THE FORWARD PHASE

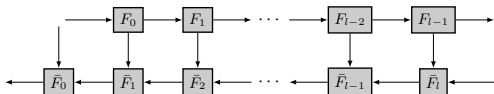


In this phase:

- ▶ We execute all  $F$ -operations;

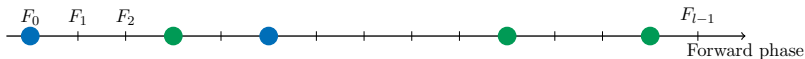


## THE FORWARD PHASE

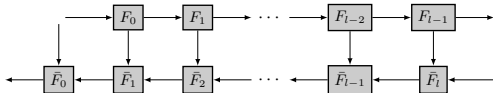


In this phase:

- ▶ We execute all  $F$ -operations;
- ▶ We write some data to disk and/or to memory.



# THE BACKWARD PHASE



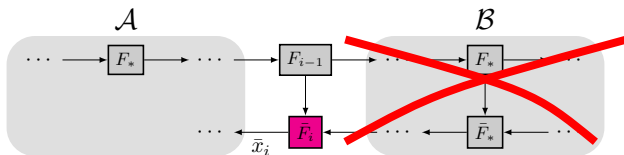
In this phase:

- ▶ We DO NOT write any data to disk (could have been done in the forward phase);
- ▶ All other operations are allowed.

# NO GOING BACK

## Lemma

If  $\bar{x}_i$  is computed, then there are no  $F_j$  for  $i \leq j$  (or operations involving  $\mathcal{B}$ ).

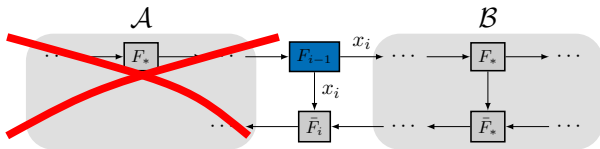




# CHECKPOINT PERSISTENCE (I)

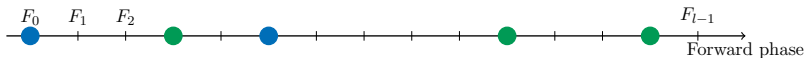
## Lemma

*If  $x_i$  is written (to disk or memory), until we have executed  $\bar{F}_i$  there are no  $F_j$  for  $j < i$  (or operations involving  $\mathcal{A}$ ).*



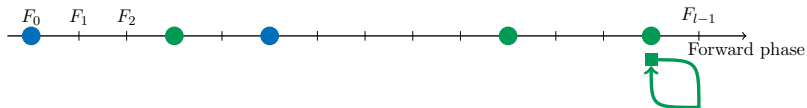
## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:



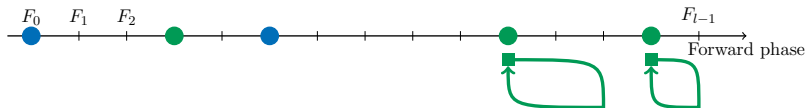
## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:



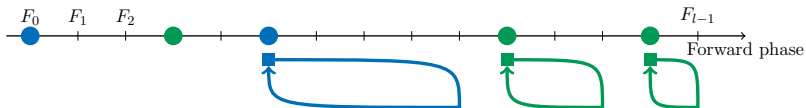
## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:



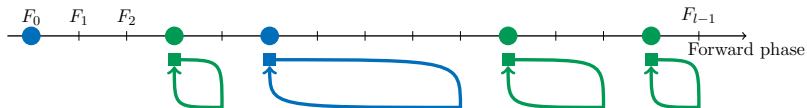
## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:



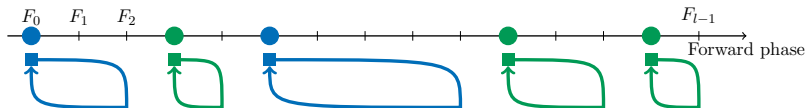
## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:

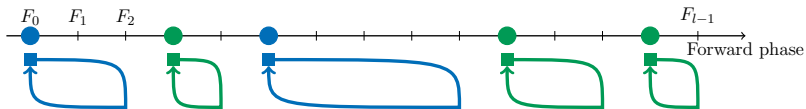


## CHECKPOINT PERSISTENCE (II)

In this case, for a given forward phase, we get a multi-phase backward phase:



# CHARACTERIZING THE BACKWARD PHASE



- ▶  $m$  backward steps to execute;
- ▶ No disk writes or reads;
- ▶  $c$  memory checkpoints available

$$\Rightarrow r_m + \text{REVOLVE}(m, c)$$

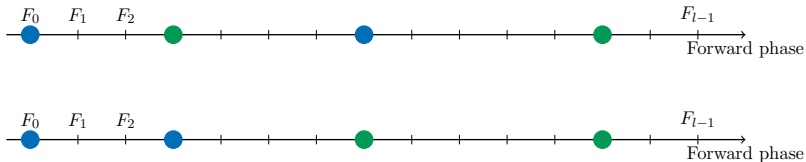




# CHARACTERIZING THE FORWARD PHASE

## Theorem

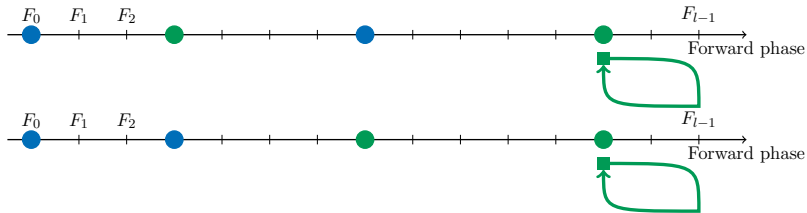
*During the forward phase, first we write to disks, then we write to memory.*



# CHARACTERIZING THE FORWARD PHASE

## Theorem

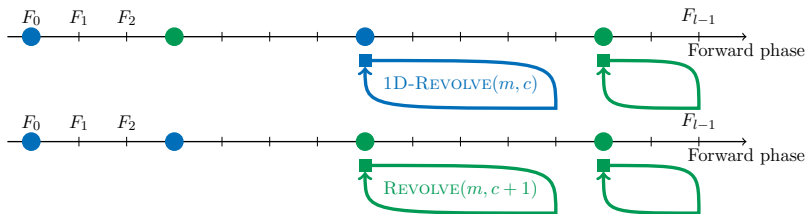
*During the forward phase, first we write to disks, then we write to memory.*



# CHARACTERIZING THE FORWARD PHASE

## Theorem

*During the forward phase, first we write to disks, then we write to memory.*

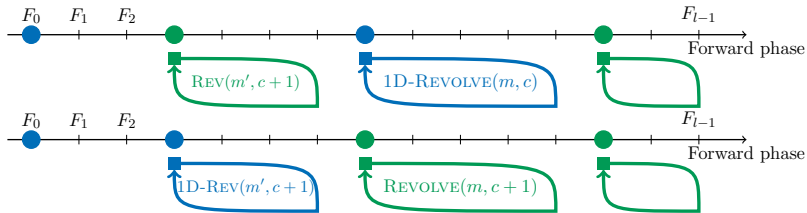


►  $\mathcal{E}xe(\text{REVOLVE}(m, c + 1)) \leq \mathcal{E}xe(\text{1D-REVOLVE}(m, c))$

# CHARACTERIZING THE FORWARD PHASE

## Theorem

*During the forward phase, first we write to disks, then we write to memory.*

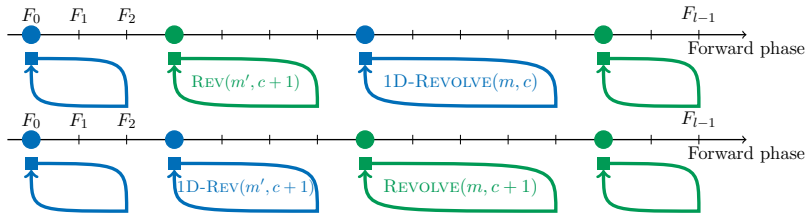


- ▶  $\mathcal{E}xe(\text{REVOLVE}(m, c+1)) \leq \mathcal{E}xe(\text{1D-REVOLVE}(m, c))$
- ▶  $\mathcal{E}xe(\text{1D-REVOLVE}(m', c+1)) \leq \mathcal{E}xe(\text{REVOLVE}(m', c+1))$

# CHARACTERIZING THE FORWARD PHASE

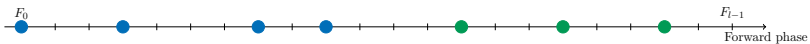
## Theorem

*During the forward phase, first we write to disks, then we write to memory.*

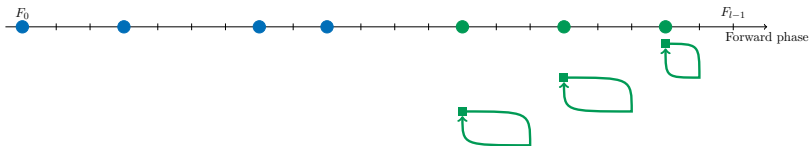


- ▶  $\mathcal{E}xe(\text{REVOLVE}(m, c+1)) \leq \mathcal{E}xe(\text{1D-REVOLVE}(m, c))$
- ▶  $\mathcal{E}xe(\text{1D-REVOLVE}(m', c+1)) \leq \mathcal{E}xe(\text{REVOLVE}(m', c+1))$

# COMPUTING THE OPTIMAL SCHEDULE

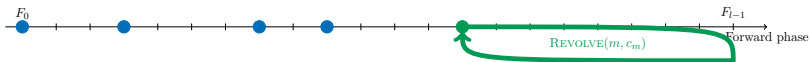


# COMPUTING THE OPTIMAL SCHEDULE

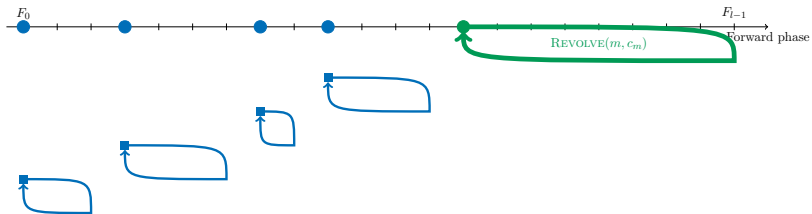




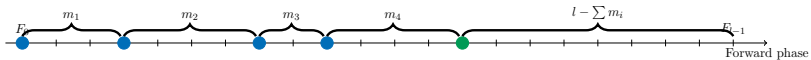
# COMPUTING THE OPTIMAL SCHEDULE



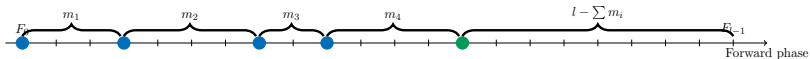
# COMPUTING THE OPTIMAL SCHEDULE



# COMPUTING THE OPTIMAL SCHEDULE



# COMPUTING THE OPTIMAL SCHEDULE

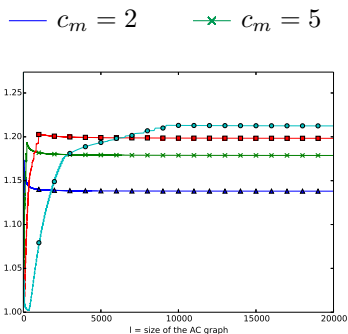


## Theorem

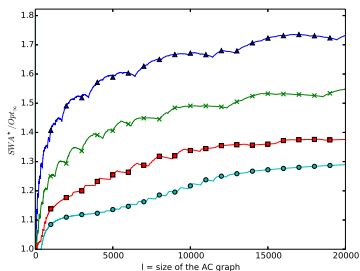
*We can compute the optimal number of disk checkpoints needed and the space between them in  $O(l^2)$  with a dynamic programming algorithm to minimize execution time.*

# IN PRACTICE?

In realistic scenarios we expect to divide the execution time by 2 or 3.



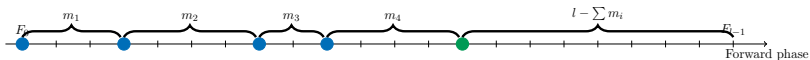
(a)  $w_d = r_d = 1$



(b)  $w_d = r_d = 5$

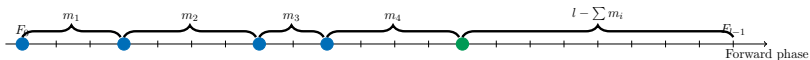
Ratio  $SWA^*(l, c_m, w_d, r_d) / Opt_{\infty}(l, c_m, w_d, r_d)$  as a function of  $l$ .

## GOING FURTHER



We know how to compute the  $m_i$ 's. But the cost of computing them is non-negligible ( $l^2$ ). Can we do better?

## GOING FURTHER

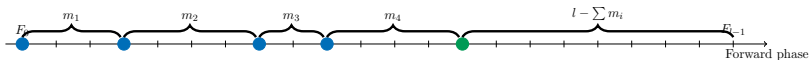


We know how to compute the  $m_i$ 's. But the cost of computing them is non-negligible ( $l^2$ ). Can we do better?

### Theorem (Weak Periodicity)

*Except for a bounded number of them (the bound depends on  $X = (c_m, w_d, r_d)$ ), all the  $m_i$ 's are equal (to  $m_X$ ).*

## GOING FURTHER



We know how to compute the  $m_i$ 's. But the cost of computing them is non-negligible ( $l^2$ ). Can we do better?

### Theorem (Weak Periodicity)

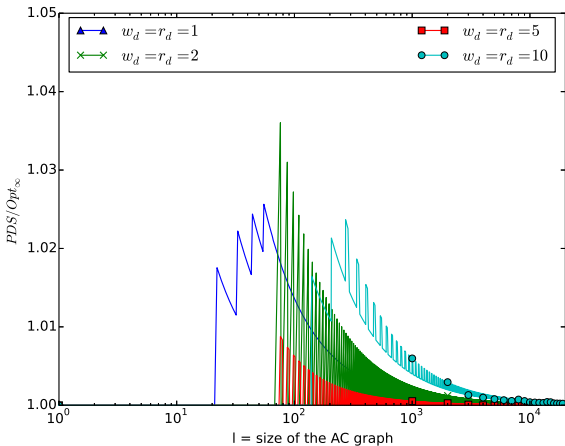
*Except for a bounded number of them (the bound depends on  $X = (c_m, w_d, r_d)$ ), all the  $m_i$ 's are equal (to  $m_X$ ).*

### Corollary

*Writing disk checkpoint every  $m_X$  forward steps is asymptotically optimal.*



# ADDITIONAL DATA



Makespan of periodic algo over optimal: function of  $l$ ,  $c_m = 10$ .

## CONCLUSIONS ON ADJOINT COMPUTATIONS

Some numbers:

- ▶ The adjoint computation in MITgcm runs in  $O(\text{days})$ , the gain induced by our optimal algorithm would be non negligible!
- ▶ Nek5000 runs on 500K cores, two processes/core on Mira. We need to take reliability into account (future work).

References:

- GW00 Griewank and Walther, *Algorithm 799: Revolve: an implementation of checkpointing for the reverse or adjoint mode of computational differentiation*, TOMS, 2000
- SW09 Stumm and Walther, *Multistage approaches for optimal offline checkpointing*, SISC, 2009

# PERSPECTIVES ON I/O MANAGEMENT

Scalable I/O management is a critical problem for Exascale.

Some directions that need to be solved:

- ▶ Models are missing!

*Understanding applications is necessary to design better solutions.*

- ▶ The energy cost of I/O management is barely studied!

*Energy is also one of the limiting factor for the next scale.*

- ▶ Applications need to be redesigned!

*Some data may not be as important as other, can we find new strategies to deal with them?*