Fixed-work vs fixed-time checkpointing

Anne Benoit

Inria, ENS Lyon & Institut Universitaire de France & IDEaS, GaTech, Atlanta, USA

Quentin Barbut (ENSL), Thomas Herault (UTK, TN & Inria), Lucas Perotin (Vanderbilt University, Nashville, TN), Yves Robert, Frédéric Vivien (Inria, ENSL)

May 14, 2025 – JLESC Workshop

|--|

Fixed-work checkpointing: divide work into chunks, checkpoint each chunk

• Goal: Complete the execution as soon as possible (minimize expected time to execute a fixed amount of work)



- If a fail-stop error strikes, restart the chunk
- Well-understood: periodic strategy, Young-Daly period $W_{YD} = \sqrt{2\mu C}$

Fixed-work checkpointing: divide work into chunks, checkpoint each chunk

• Goal: Complete the execution as soon as possible (minimize expected time to execute a fixed amount of work)



- If a fail-stop error strikes, restart the chunk
- Well-understood: periodic strategy, Young-Daly period $W_{YD} = \sqrt{2\mu C}$

• Goal: Complete as much work as possible (minimize expected time to execute a fixed amount of work)



• If a fail-stop error strikes, what do we do next?

• Goal: Complete as much work as possible (minimize expected time to execute a fixed amount of work)



• If a fail-stop error strikes, what do we do next?



• Goal: Complete as much work as possible (minimize expected time to execute a fixed amount of work)



• If a fail-stop error strikes, what do we do next?



• Goal: Complete as much work as possible (minimize expected time to execute a fixed amount of work)



• If a fail-stop error strikes, what do we do next?



Problem: Fixed-time checkpointing

- Checkpointing strategies for a large-scale parallel application
- Nodes are subject to failures (Exponential distribution of parameter λ , $\mu = \frac{1}{\lambda}$)
- The application executes for a fixed duration, i.e., the length T of its reservation
- Checkpoint cost C, recovery cost R, and downtime D
- The last checkpoint within the reservation plays a particular role:
 ⇒ all the work executed after that checkpoint will be lost
 - \Rightarrow take it close to, or exactly at the very end of the reservation



- Important question, but not yet addressed
- Dual of the fixed-work checkpointing classical problem

JLESC 2025

Problem: Fixed-time checkpointing

- Checkpointing strategies for a large-scale parallel application
- Nodes are subject to failures (Exponential distribution of parameter λ , $\mu = \frac{1}{\lambda}$)
- The application executes for a fixed duration, i.e., the length T of its reservation
- Checkpoint cost C, recovery cost R, and downtime D
- The last checkpoint within the reservation plays a particular role:
 - \Rightarrow all the work executed after that checkpoint will be lost
 - \Rightarrow take it close to, or exactly at the very end of the reservation



- Important question, but not yet addressed
- Dual of the fixed-work checkpointing classical problem

JLESC 2025

Why a fixed-length reservation?

- \bullet Job executes and checkpoints at the end of its reservation of length ${\cal T}$
- Job is killed after T seconds
- Typical scenario:
 - long running HPC application
 - split into multiple smaller reservations
- No failures: Checkpoint at time T C?
 - C not always constant, but may have a bounded support $[C_{min}, C_{max}]$
 - Checkpointing at time $T C_{max}$ is pessimistic and may lead to wasting execution, but no risk taken...
 - Full study in [FTXS'23: When to checkpoint at the end of a fixed-length reservation?]
 - Problem solved if it is possible to checkpoint an any instant, for various probability laws of $C \Rightarrow$ later than at time $T C_{max}!$

• ● ● ● ● ● ●

With failures: Dynamic checkpointing strategy

- Initially, the time left is $t_{left} = T$
- The strategy decides how many checkpoints will be taken, and at which instants, if there is no failure during the whole execution
- Let $t_{end}(i) \leq T$ be the completion time of checkpoint number *i*, with $1 \leq i \leq k$
- If there is no failure up to time $t_{end}(k)$ (when the last checkpoint completes), then the work achieved by the strategy is $W = t_{end}(k) kC$.
- If a (first) failure strikes at time t ≤ T, let ℓ ≤ k be the number of the last checkpoint that completed before the failure:
 - The work done after time $t_{end}(\ell)$ and up to time t is lost
 - The work achieved up to the failure is $W = t_{end}(\ell) \ell C$
- Now after the failure at time t, there is a downtime and a recovery; hence, the time left is $t_{left} = (T t) D R$
- If $t_{left} \ge C$, call recursively the strategy and add the work done then to W
- $\bullet\,$ The objective is to maximize the expected amount of work $\mathbb{E}(\,\mathcal{T})=\mathcal{W}$

With failures: Dynamic checkpointing strategy

• Initially, the time left is $t_{left} = T$

Dynamic approach

• For any value t_{left} , the strategy decides how many checkpoints will be taken, and at which instants, but it is applied until the end only if no failure strikes during the execution.

• The strategy will not be the same for different values of t_{left} : the distances between consecutive checkpoints are recomputed after each failure

Nice but challenging little scheduling problem 🙂

 \sim II $l_{left} \leq c$, call recursively the strategy and add the work done then to m

• The objective is to maximize the expected amount of work $\mathbb{E}(\mathcal{T}) = W$

- No closed-form formula even for a single checkpoint at the end
- If you decide to go for a single checkpoint: do not always place it at the end of the reservation, even if *C* is constant!



• If you decide to go for two checkpoints: almost never use two same-size chunks!



Fixed-time checkpointing

$$\begin{split} \mathbb{E}^{end}(T) &= e^{-\lambda T}(T-C) + (1-e^{-\lambda T}) \int_{t=0}^{T-D-R-C} \lambda t \frac{e^{-\lambda t}}{1-e^{-\lambda T}} \mathbb{E}_{R}^{end}(T-t-D) dt \\ \mathbb{E}_{R}^{end}(t_{left}) &= e^{-\lambda T}(t_{left}-R-C) + \int_{t=0}^{t_{left}-D-R-C} \lambda t \frac{e^{-\lambda t}}{1-e^{-\lambda T}} \mathbb{E}_{R}^{end}(t_{left}-t-D) dt \end{split}$$

$$\bigwedge$$
 No closed-form expression $igodol{arepsilon}$

Fixed-work checkpointing

$$\mathbb{E}(\mathcal{W}) = e^{-\lambda(W+C)}(W+C) + (1 - e^{-\lambda(W+C)})(\mathbb{E}(T_{lost}(W+C)) + \mathbb{E}(T_{rec}) + \mathbb{E}(\mathcal{W}))$$

$$\Rightarrow \mathbb{E}(W) = (\frac{1}{\lambda} + D) e^{\lambda R}(e^{\lambda(W+C)} - 1)$$

- ∢ ≣ ►

Single checkpoint is **not always** at the end of the reservation



Short reservation, T = 6, with D = 0, C = R = 4Strategy STRAT₁: Checkpoint at the end Strategy STRAT₂: Checkpoint before the end

Expected gain GAIN of $STRAT_1$ over $STRAT_2$

• Contrived example: no time left for any work after a failure

• No failure: GAIN =
$$\mathbb{P}_{succ}(6) \times 1$$
:

- First failure strikes after t = 5: GAIN = $\mathbb{P}_{succ}(5)\mathbb{P}_{fail}(1) \times (-1)$
- First failure strikes before t = 5: GAIN = 0

• GAIN =
$$e^{-6\lambda} - e^{-5\lambda}(1 - e^{-\lambda}) = 2e^{-6\lambda} - e^{-5\lambda} = e^{\ln(2)-6\lambda} - e^{-5\lambda}$$

• If $\lambda > \ln(2)$, better to checkpoint before the end of the reservation!

Dynamic threshold-based strategy

• Natural conjecture:

- Easy to prove that expected work increases with T
- Assume same-size segments and last checkpoint at the end
- Conjecture that longer reservations initially require more segments (hence more checkpoints)
- (Sub-optimal) strategy based on this conjecture:
 - Construct sequences of thresholds T_n such that there are exactly n checkpoints for $T_n \leq T \leq T_{n+1}$, use number of checkpoints corresponding to t_{left}



 $GAIN(T, n + 1) = \mathbb{E}(T, n + 1) - \mathbb{E}(T, n)$ Numerically, there is a single solution to GAIN(T, n + 1) = 0 \bigcirc

- YOUNGDALY is the baseline reference, which takes checkpoints periodically following the Young/Daly formula W_{YD} (with mandatory checkpoint at the end)
- NUMERICAL OPTIMUM uses a numerical approximation of the thresholds T_n
- FIRSTORDER uses a first-order approximation of the thresholds:

 $T_{n+1} \approx \sqrt{2n(n+1)\mu C}$

For n = 1, $T_2 \approx \sqrt{4\mu C}$ while $W_{YD} = \sqrt{2\mu C}$ (single checkpoint for longer reservations)

- Optimal strategy: Quantum-based DYNAMICPROGRAMMING algorithm
 - General approach using small time quanta
 - Use dynamic programming to compute the best checkpointing strategy
 - The smaller the quanta, the more accurate the results ...

... but the more costly the algorithm

• • = • • = •

Dynamic programming

- $\mathbb{E}(n, k, \delta)$ optimal work (in expectation)
 - during *n* quanta
 - when taking exactly k checkpoints initially
 - $\delta=0$ if there is no initial recovery with the work, and $\delta=1$ otherwise
- $k_{max} = \lfloor \frac{T^*}{C^*} \rfloor$ maximum number of checkpoints (if we checkpoint all the time!)
- Compute $\mathbb{E}_{opt}(T^*) = \max_{1 \leq k \leq k_{max}} \mathbb{E}(T^*, k, 0)$

Complexity: $O((T^*)^2 k_{max}^2)$

Quantum durationuNumber of quanta $T^* = \frac{T}{u}$ Checkpoint time in numbers of quanta $C^* = \frac{C}{u}$

(4) 国家 (4) 国家

Solution

With one checkpoint

$$\begin{split} \mathbb{E}(n,1,\delta) &= \max_{\delta R^* + C^* + 1 \le i \le n} \mathbb{P}^*_{succ}(i)(i - C^* - \delta R^*) \\ &+ \sum_{f=1}^{i - D^* - R^* - C^*} p^*_f \mathbb{E}(n - f - D^*, 1, 1) \end{split}$$

With *k* checkpoints

$$\mathbb{E}(n,k,\delta) = \max_{\substack{\delta R^* + C^* + 1 \le i \le n - (k-1)C^* \\ f = 1}} \left(\mathbb{P}^*_{succ}(i)(i - C^* - \delta R^* + \mathbb{E}(n-i,k-1,0)) + \sum_{f=1}^i p_f^* \max_{1 \le m \le k} \mathbb{E}(n-f - D^*,m,1) \right)$$

Initialization

$$\mathbb{E}(n, k, 0) = 0$$
 for $n \le kC^*$ (where possibly $n \le 0$)
 $\mathbb{E}(n, k, 1) = 0$ for $n \le R^* + kC^*$ (where possibly $n \le 0$)

• • = • • = •

< A

э

Parameters

- $C \in \{10, 20, 40, 80, 160\}$, and R = C
- *D* ∈ {0,5}
- $\lambda \in \{0.01, 0.001, 0.0001\}$
- $T \in [C, 2000\}$
- Quantum size u = 1 for DYNAMICPROGRAMMING (after experimenting with other sizes)

(arbitrary time-unit - change granularity of the scenarios)

Report proportion of work achieved $= \frac{\text{amount of work achieved}}{T-C}$ (the higher the better)

Results



May 14, 2025

Fixed-work vs fixed-time checkpointing

JLESC 2025

Results



With more failures: Proportion of work when $\lambda = 0.01$, D = 0, and $C \in \{80, 160\}$

May 14, 2025

- YOUNGDALY achieves significantly low performance when only a handful of checkpoints are required, but performs well for long reservations
- NumericalOptimum \geq FirstOrder \geq YoungDaly
- DYNAMICPROGRAMMING only slightly better than NUMERICALOPTIMUM, but much more costly Dynamic threshold strategy very easy to use and close to optimal \bigcirc
- Big differences mainly for high failure rates and short reservations

Details in [FTXS'2024: Checkpointing strategies for a fixed-length execution]

- Fixed-time checkpointing much more difficult than Fixed-work checkpointing
- Optimal solution seems out of reach 😟
- Three good news:
 - Optimal discretized solution as a guarantee of quality
 - Young/Daly performs well for reservations lasting at least, say, a dozen Young/Daly periods $W_{YD} = \sqrt{2\mu C}$
 - Dynamic threshold strategy very easy to use and performs remarkably well

• Future work:

- Extend results to a stochastic framework, where both checkpoint durations and application progress rate are no longer fully deterministic
- Apply to sparse linear algebra problems