

# Deadline Miss Minimization Scheduling for License-Constrained CAE Jobs in Hybrid Cloud Infrastructure

- Mohamed Noaman, **Srishti Dasgupta**, Michael Gerndt -

28th Workshop on Job Scheduling Strategies for Parallel Processing 39th IEEE International Parallel and Distributed Processing System

Milan, Italy

3 - 4 June, 2025





#### Motivation & Problem Statement

#### What is a CAE Simulation



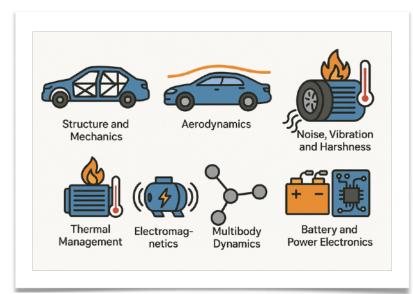


"Computer-Aided Engineering(CAE) software has revolutionized car modeling by enabling virtual testing and optimization of structural, thermal, aerodynamic, and dynamic performance — drastically reducing time, cost, and physical prototypes."



#### Bottlenecks in the CAE Simulation Field







Often computationally heavy requiring hundreds of cores



CAE tools drive safety & innovation—but come with tight deadlines, expensive licenses, and compute bottlenecks.









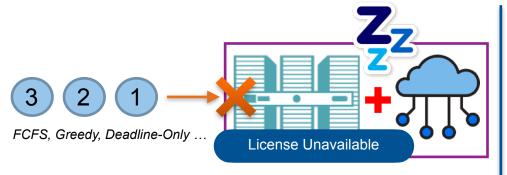






#### Why Advanced Scheduling is Needed?

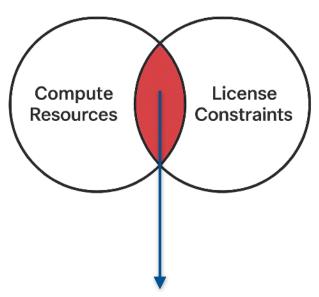




Current schedulers overlook the interdependency of licenses and compute — leading to missed deadlines and wasted cloud costs.



License should also be treated as Resource



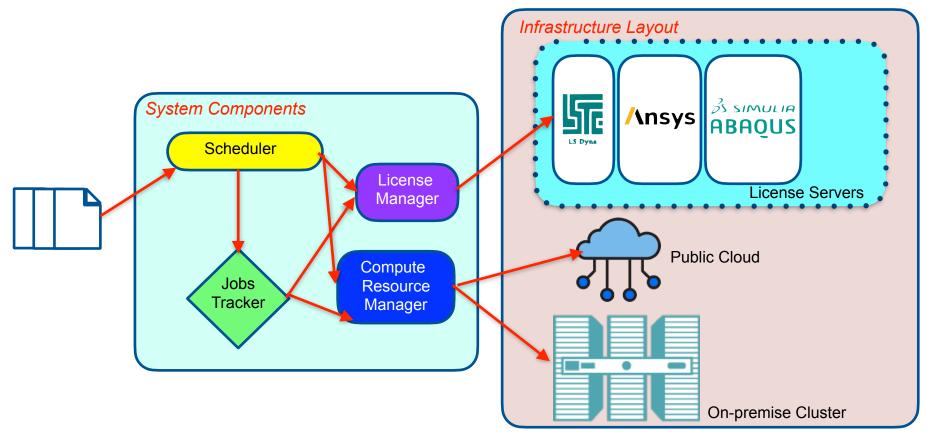
License-Aware Deadline
Miss Minimization Algorithm



Proposed Solution: LADMM\_backfill Algorithm

# System Design





# LADMM\_backfill Algorithm: Core Objectives



Built in principle of EASY backfilling strategy



Aimed at minimizing deadline misses



Integrates real-time availability of floating licenses



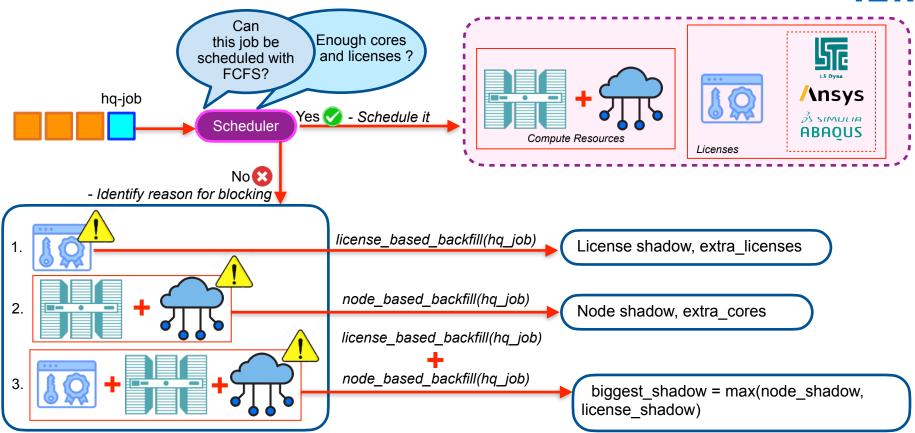
Favors jobs requiring critical resources



Handles both On-Prem HPC clusters and public cloud nodes

# LADMM\_backfill Algorithm: Initial Resource Allocation





## LADMM\_backfill Algorithm: Scoring & Backfilling



• Total Score for Job j:

$$S_j = S_{critical} + S_{deadline}$$

Critical Resources Heuristics(S<sub>critical</sub>)

$$S_{critical} = C_n \cdot t_j + C_T \cdot (R_j/R_{min})$$

 $C_n$  : critical resources licensing factor  $t_i$  : number of licenses required by jobj

 $oldsymbol{\mathcal{C}_{T}}$  : critical resources runtime factor

 $R_{ extit{min}}$  : runtime of the shortest job in the queue

 $oldsymbol{R}_{oldsymbol{j}}$  ; runtime of jobj

• Deadline Heuristics(Sdeadline)

$$S_{deadline} = \begin{cases} (D_B \cdot D_{min}) / (D_j - T_{current} + R_j), \\ If D_j - T_{current} + R_j > 0 \\ 0, otherwise \end{cases}$$

D<sub>B</sub> : deadline urgency weight

 $D_{ exttt{min}}$  ; smallest deadline gap in the queue

 $D_j$  : deadline of job j

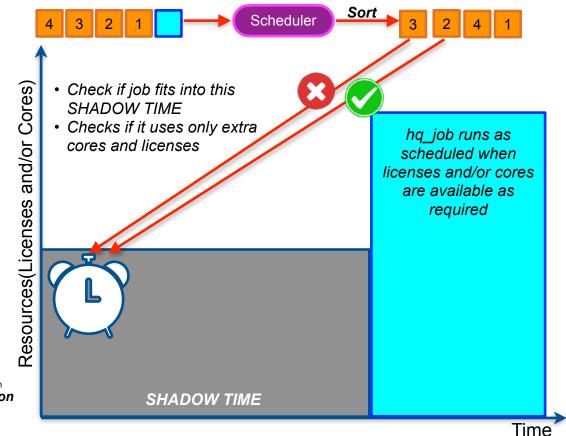
Tcurrent : current simulation time

 $R_i$  : runtime of jobj

#### How scoring was tuned?

• A Pareto front was used 125000 jobs were simulated with varying  $D_B,\,C_T\,\&\,C_n$ 

 Selected weights that maximized License Minimization & minimized Deadline Misses





# **Experimental Setup and Results**

#### **Experimental Setup**



- <u>Licensing Model</u>
  - Floating License Servers
  - 3 CAE Licenses chosen: LS-Dyna, Abaqus & Ansys with their cost models [1]
  - Fixed number of licenses available: Ranging from 50 to 250
- Job Traces
  - From Parallel Workfloads Archive [2]
    - Submission Time, Core Count, runtime
  - Deadline = Runtime \* Factor(2 to 5)
- Each Job assigned a CAE software randomly, with Ansys jobs running specifically on On-premise
- Job Scale:
  - Total jobs: 45000 with core requirements 1 to 50
- Infrastructure :
  - On-premise: considered 528 total cluster cores (48 and 96 cores machines)
  - Cloud: based on AWS m7 instances with 1:1 vcpu to physical core mapping, reflecting license constraints
    - Pay-As-You go model with per-minute billing
  - Used simulus [3] to build the infrastructure

#### **Baselines Used**



- 1. First Come First Served: **FCFS**
- 2. Greedy Technique sorting jobs w.r.t deadlines : **GREEDY\_deadline**
- 3. Heuristics-based Backfilling: **HEURISTICS\_backfill** 
  - · Attempts backfilling via scoring

S<sub>critical</sub> Prioritises jobs requiring more licenses and shorter runtimes, deemed "critical"

S<sub>deadline</sub>" Prioritizes jobs closer to deadline, with respect to the average job queue

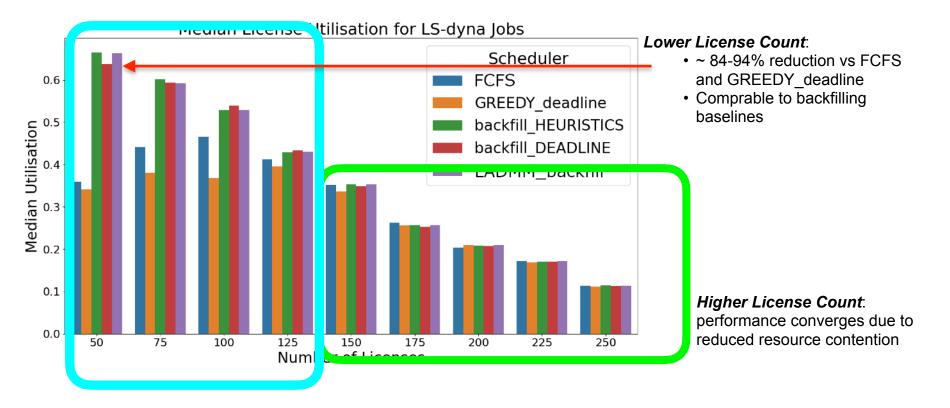
**S**<sub>starvation</sub> Boosts long-waiting jobs to avoid starvation

- 4. Deadline-based Backfilling : **DEADLINE\_backfill** 
  - Checks if a job can meet its deadline, and attempts safe backfilling using shadow time calculations if not
  - If above is not possible, moves job to an "impossible queue"

<sup>\*\*</sup> Also uses a metric "Ansys Ratio" which monitors and prioritized Ansys jobs on on-premise if they are being starved.

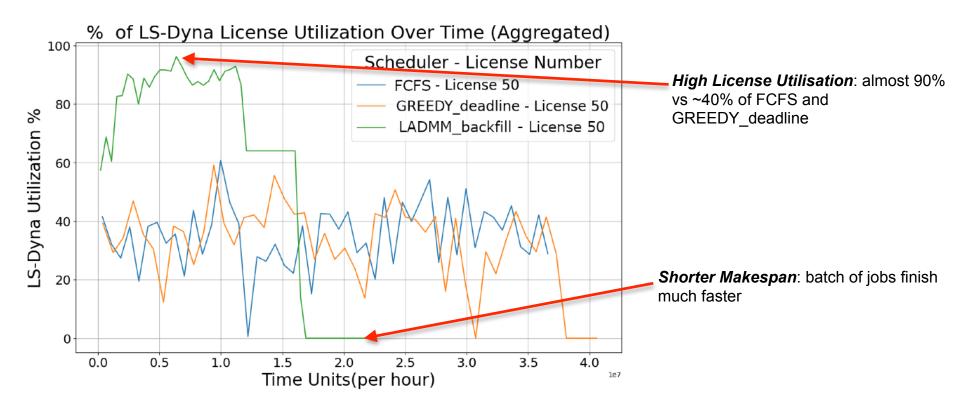
#### Median License Utilization of LS-Dyna Jobs





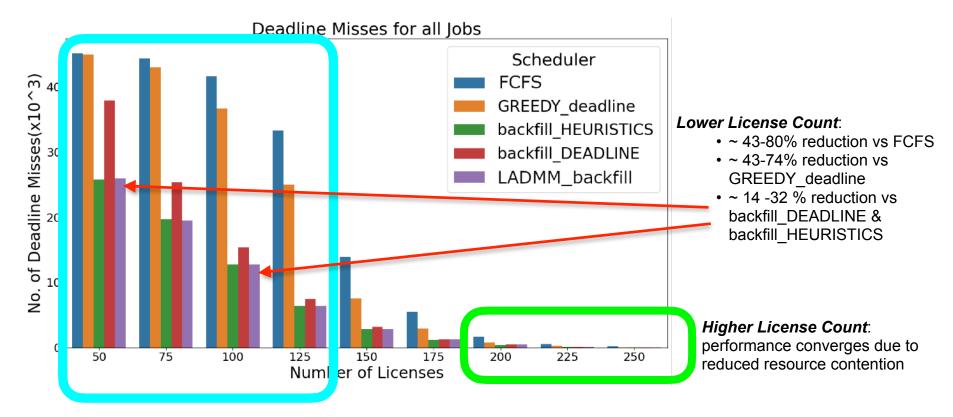
## % of License Utilization of LS-Dyna Jobs





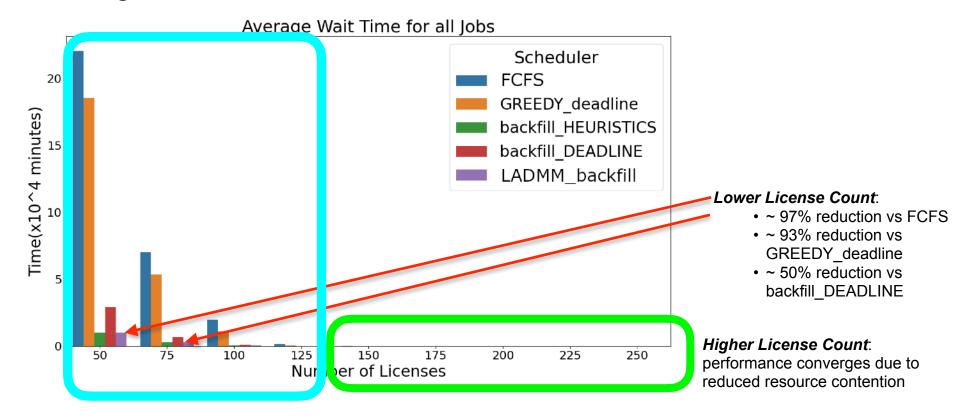
#### **Deadline Miss Ratio**





# Average Wait Time





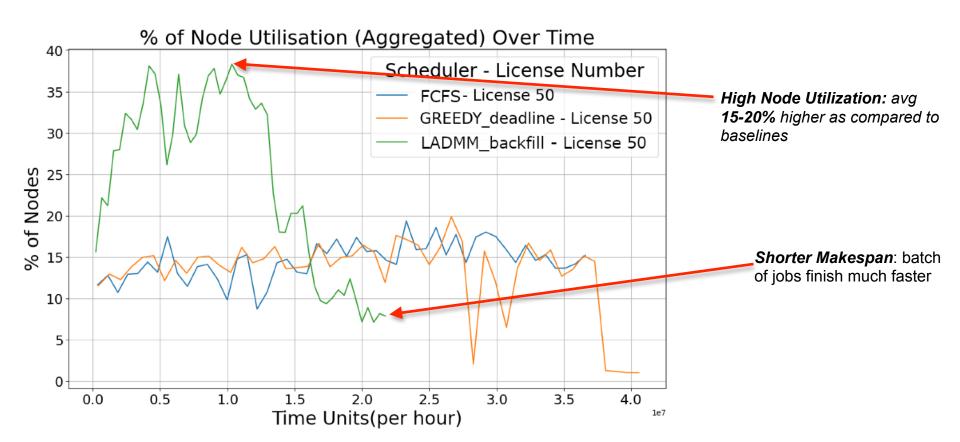
# Public Cloud Costs(x1000 \$)



License Number	FCFS	GREEDY deadline	BACKFILLING heuristics	BACKFILLING deadline	LADMM backfill
50	0	0	3353.485	0	0
75	0	0	229.009	0	0
100	0	0	934.808	0	0
125	1.062	0	907.626	.965	0
150	1.106	324	313.283	4.069	4.130
175	7.261	2.649	125.788	4.788	5.392
200	6.292	6.099	45.241	5.882	7.470
225	7.705	7682	13.616	5.033	8.880
250	10.062	8.186	14.365	5.943	9.899

#### Global Node Utilisation





#### Conclusion & Future Work



- Proposed **LADMM\_backfill**: a license-aware, deadline-driven scheduler for hybrid HPC-cloud environments.
- Integrates license and core availability using optimized scoring heuristics.
- Outperforms FCFS, Greedy, and other backfilling baselines in:

  - Sost: Maintained or reduced, especially under tight license constraints.
- Extend to workflow-based scheduling with inter-task dependencies.
- Support pay-as-you-go licensing models for greater flexibility.
- Integrate dynamic pricing and ML-based demand prediction.
- Validate at larger scales and across diverse CAE workloads for robustness.

# Thank You! Questions?

[1]Henkel, N., Treiber, S.: Cost-effective sizing of your HPC cluster for CAE sim ulations. In: IEEE Int. Conf. High Perform. Comput. (2015). https://api.semanticscholar.orgCorpusID:35686422 [2]https://www.cs.huji.ac.il/labs/parallel/workload/ [3]https://pypi.org/project/simulus/