



# **NAVGEM on the Cloud:** Computational Evaluation of Commercial Cloud HPC with a Global Atmospheric Model

- Background
  - Computing Resources
  - Cloud Migration
- Test Workload
  - Navy Global Environmental Model
  - Benchmark Specification
- Low resolution forecast performance
- Elastic Fabric Adapter on AWS EC2
- High resolution forecast performance
- Next Steps

## Navy Supercomputing

- Navy's arm of DoD HPC Modernization Program
- One of five DoD HPC Centers
- Headquartered with Naval Meteorology and Oceanography Command
- Supports various defense computational areas:
  - Climate/Weather/Ocean Modeling and Simulation
  - Computational Structural Mechanics
  - Computational Electromagnetics and Acoustics
  - Space and Astrophysical Science

## Navy DSRC

Stennis Space Center, Mississippi

HPE SGI 8600  
3.05 PFLOPS



Gaffney



Koehr

Cray XC40  
2 PFLOPS



Conrad



Gordon

## DoD Priority: Modernization, Cost-Savings, Redundancy

### Directive History

- 2012 National Defense Authorization Act
- 2012-2014 Navy Approach to Cloud Computing
- 2015 Acquisition and Use of Cloud Services
- 2017 Navy Cloud First Policy
- 2019 Federal Cloud Computing Strategy – Cloud Smart

Performance Plan for Reduction of Resources Required for Data Servers and Centers:

*“Migration... to cloud computing services... that provide a better capability at a lower cost with the same or greater degree of security.”*

- NDAA FY 2012

### Navy’s Emphasis

- Transition public facing websites
- Reduce data centers
- Increasing secure capabilities

*“... reduce investment in traditional, on –premises... data centers ... [including] Special Purpose Processing Nodes...”*

- Navy Cloud First Policy

### Success with Enterprise Applications

- Public-facing Navy websites
  - Fleet and Family has shifted 100 web-based systems to AWS GovCloud.
  - My Navy Portal team awarded for implementing the first Level 4 data enterprise architecture to the cloud.

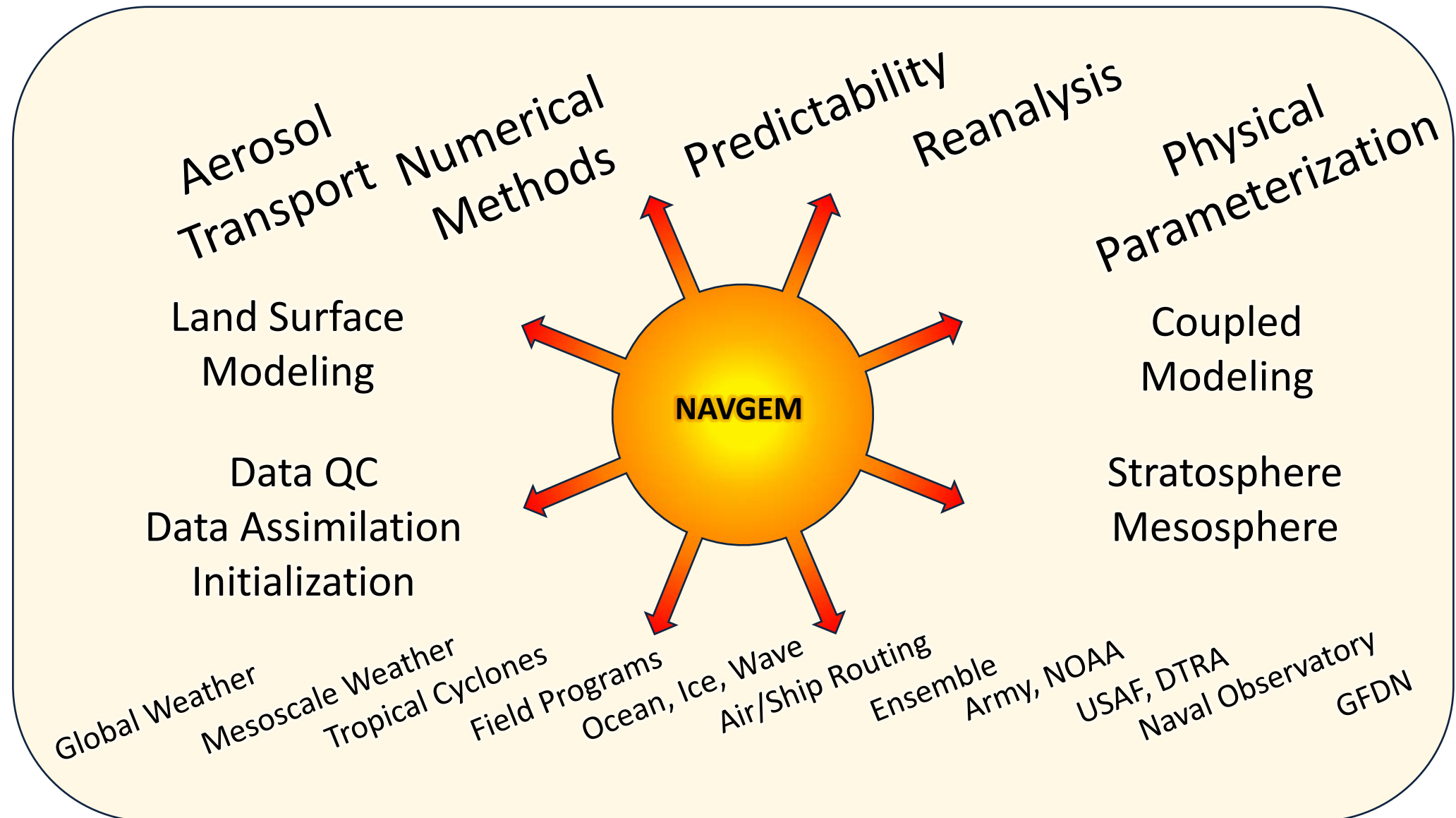




# Test Workload: Navy Global Environmental Model

## NAVGEN

- Global (synoptic) scale weather prediction program written in Fortran/MPI .
- Dynamical core composed of Semi-Lagrangian/Semi-Implicit numerical methods to solve primitive equations on a sphere.
- Output products feed in to numerous external programs and organizations.



# Benchmark Specifications

	Navy DSRC - Conrad (Cray XC40)	AWS c4.8xlarge	Azure H16r	Penguin B30 queue	AWS c5n.18xlarge
CPU	<ul style="list-style-type: none"> <li>- 2.3 GHz Intel Xeon E5-2698 v3 Broadwell</li> <li>- 32 core nodes</li> </ul>	<ul style="list-style-type: none"> <li>- 2.9 GHz Intel Xeon E5-2666 v3 Haswell</li> <li>- 18 core nodes</li> </ul>	<ul style="list-style-type: none"> <li>- 3.2 GHz Intel Xeon E5-2667 v3</li> <li>- 14 core nodes</li> </ul>	<ul style="list-style-type: none"> <li>- 2.4 GHz Intel Xeon E5-2680 v4 Broadwell</li> <li>- 28 core nodes</li> </ul>	<ul style="list-style-type: none"> <li>- 3.0 GHz Intel Xeon Platinum w/ AVX-512</li> <li>- 36 core nodes</li> </ul>
Network	<ul style="list-style-type: none"> <li>- Cray Aries / Dragonfly</li> </ul>	<ul style="list-style-type: none"> <li>- 25 Gbps ethernet with SRIOV</li> </ul>	<ul style="list-style-type: none"> <li>- FDR Infiniband</li> </ul>	<ul style="list-style-type: none"> <li>- Intel OmniPath</li> </ul>	<ul style="list-style-type: none"> <li>- AWS EFA</li> </ul>

## Software Configuration

- Intel Fortran 2018 update 1
- MPI:
  - Intel 2018 update 1
  - EFA: Open MPI 3.1.4
- HDF5 1.8.20

## Platform Parameters

- Hyperthreading disabled
- AWS
  - us-west-2 (OR) region
  - Placement Groups
  - CloudFormation
- Azure
  - US Gov Arizona
  - VMSS
- Penguin on Demand
  - PBS resource manager

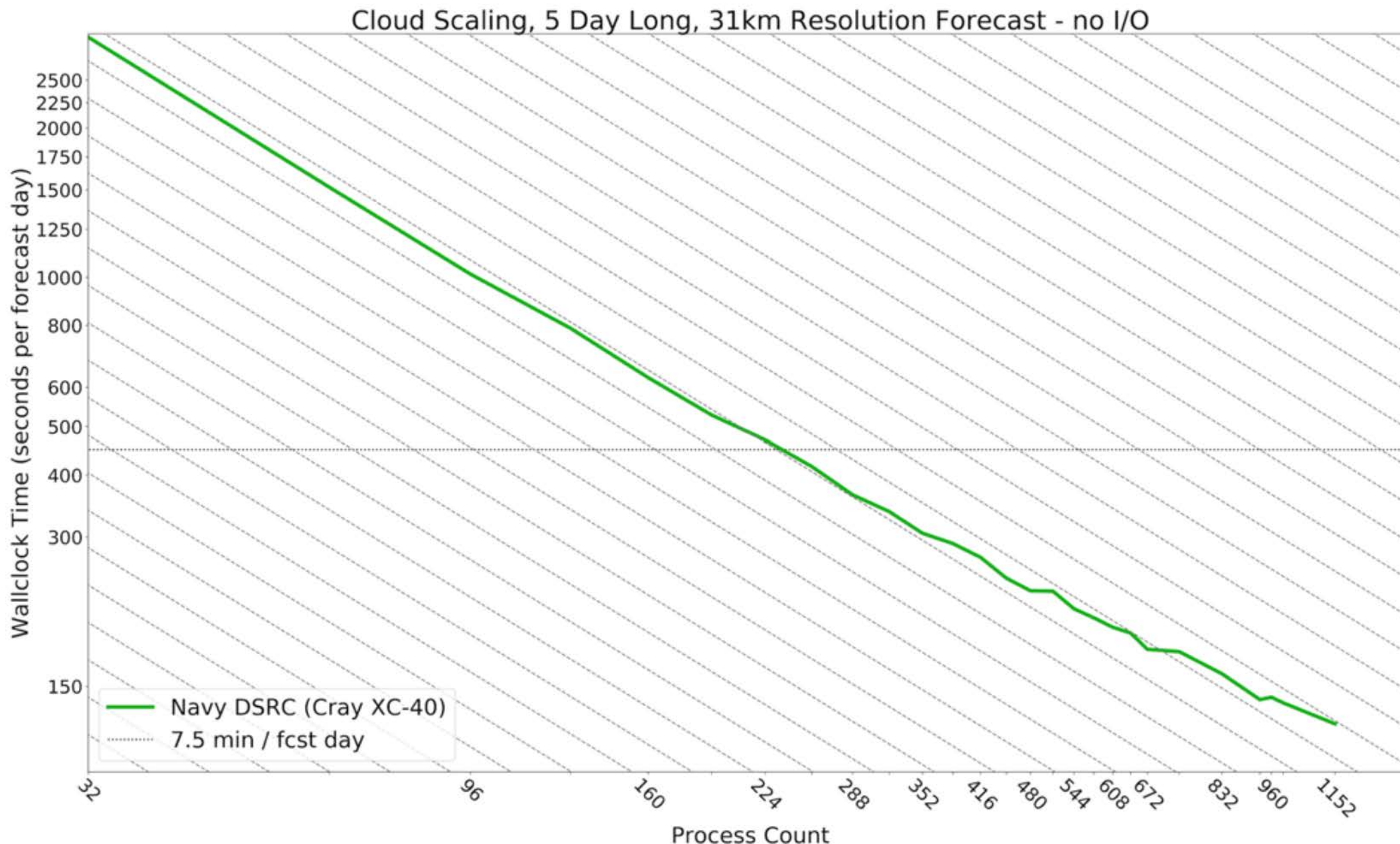
# Low Resolution Forecast: Performance - Navy DSRC

## Platform Specifications:

- 2.3 GHz Intel Xeon E5-2698 v3 Broadwell
- 32 core nodes
- Cray Aries / Dragonfly
- Cray Linux

## Results

- Tested on Conrad
- Good scaling
  - minimal variability
  - minor increase in slope.





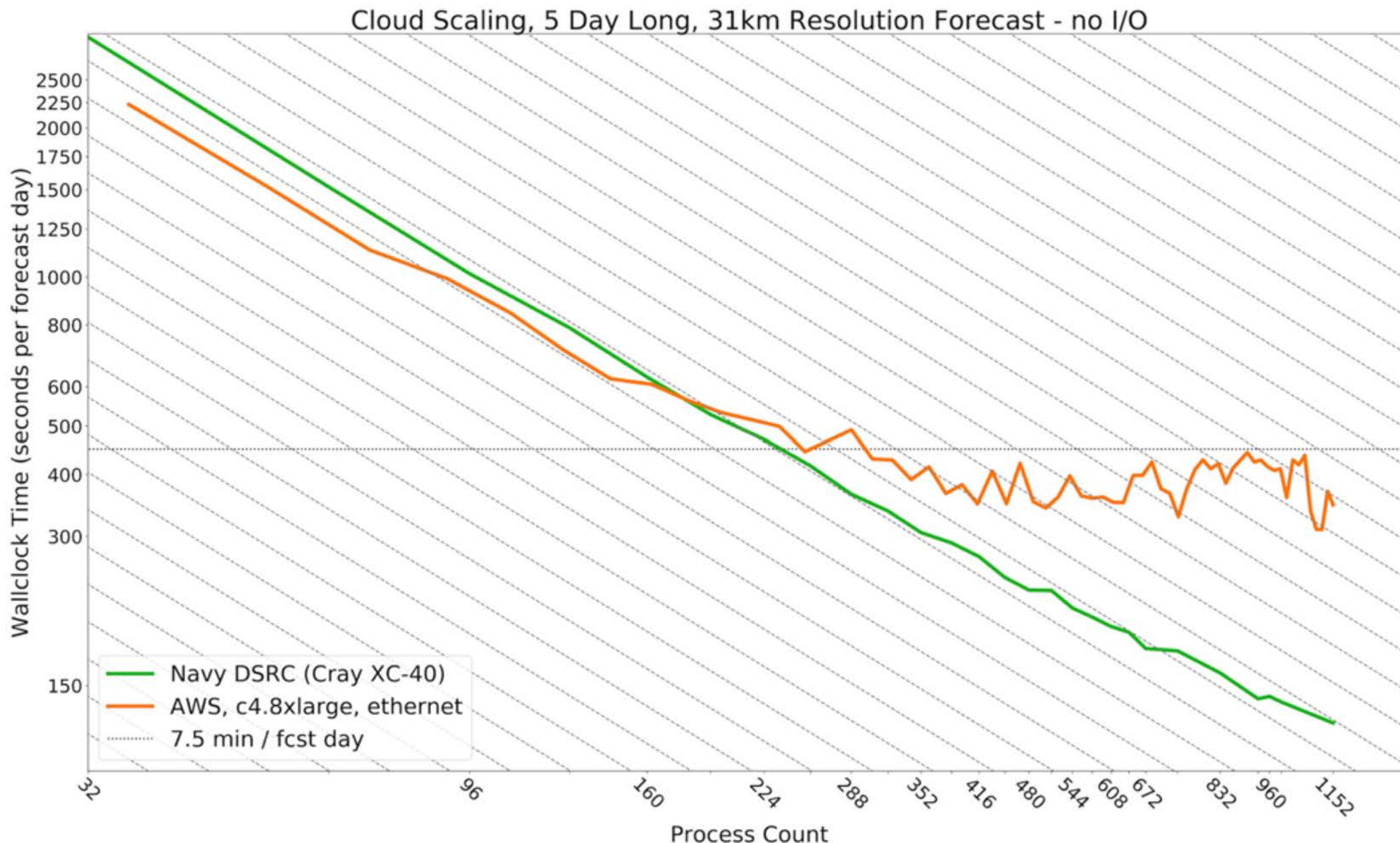
# Low Resolution Forecast: Performance - AWS EC2 1

## Platform Specifications:

- 2.9 GHz Intel Xeon E5-2666 v3 Haswell
- 18 core nodes
- 25 Gbps ethernet with SRIOV
- Amazon Linux

## Results

- AWS's most powerful compute optimized instance at time of testing.
- Adjusted variety of configurations
  - System clocksource
  - Attached storage type
  - Cluster creation tools
  - Processors per node





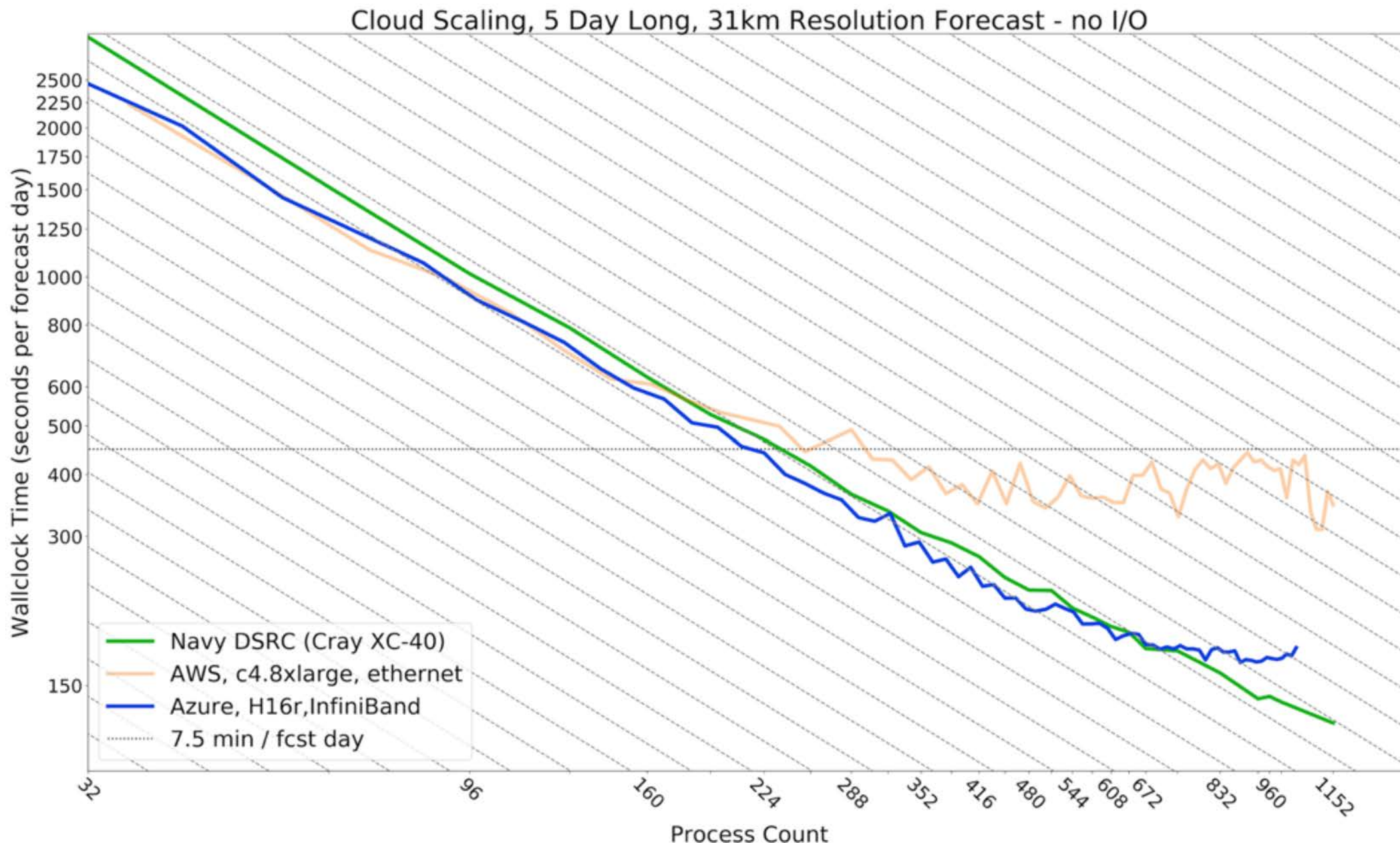
# Low Resolution Forecast: Performance - Azure

## Platform Specifications:

- 3.2 GHz Intel Xeon E5-2667 v3
- 14 core nodes
- FDR Infiniband
- CentOS Linux

## Results

- Infiniband networking likely contributed to improvement.
- Improved performance (plotted) using 14/16 processors per node.
- Eventual flattening of performance.



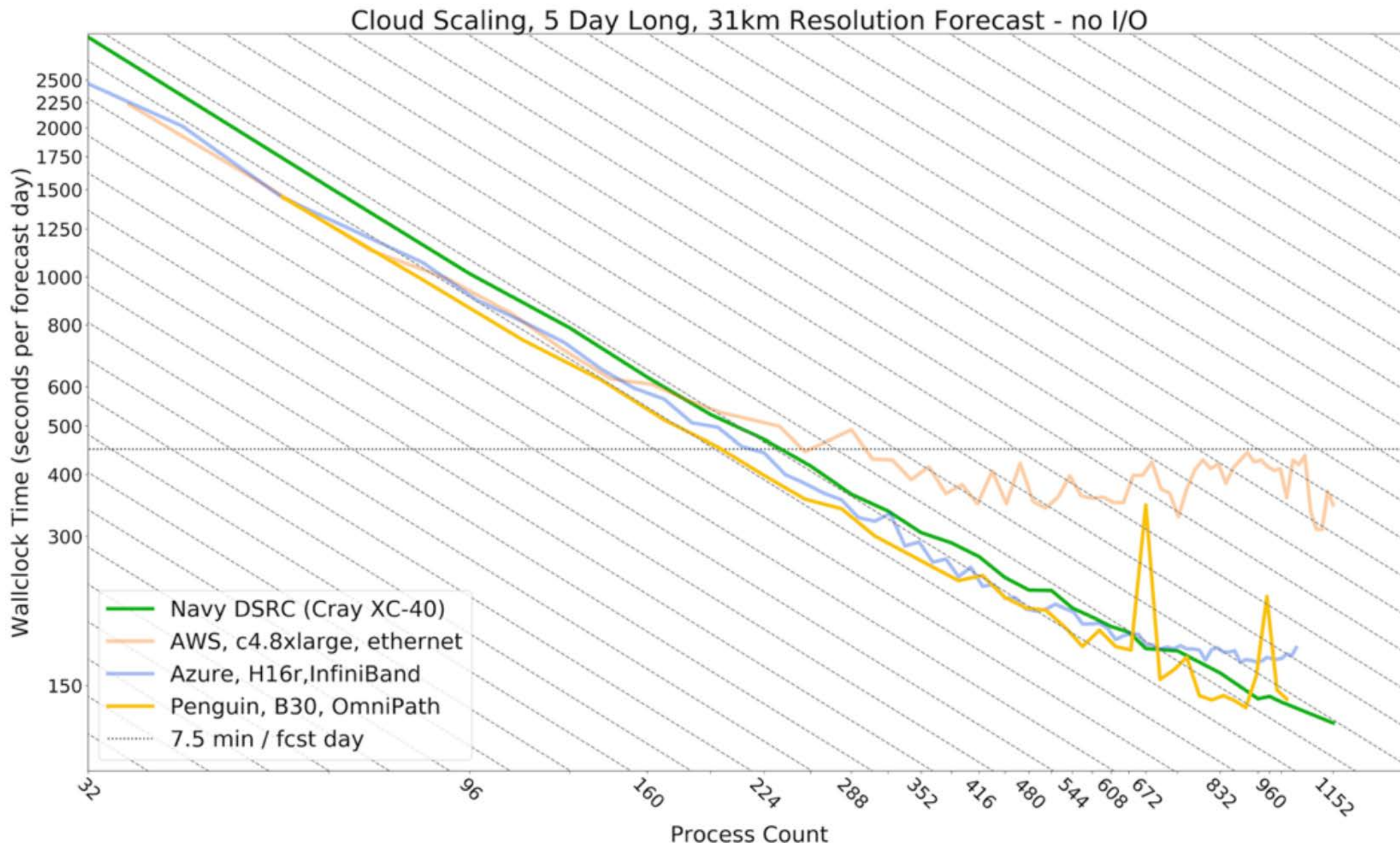
# Low Resolution Forecast: Performance - Penguin

## Platform Specifications:

- 2.4 GHz Intel Xeon E5-2680 v4 Broadwell
- 28 core nodes
- Intel OmniPath

## Results

- Improved scaling expected considering Penguin's "traditional" system design.
- Required use of batch scheduler and waiting for system resources at larger cluster sizes.
- Variability and unexplained spikes at higher core counts.





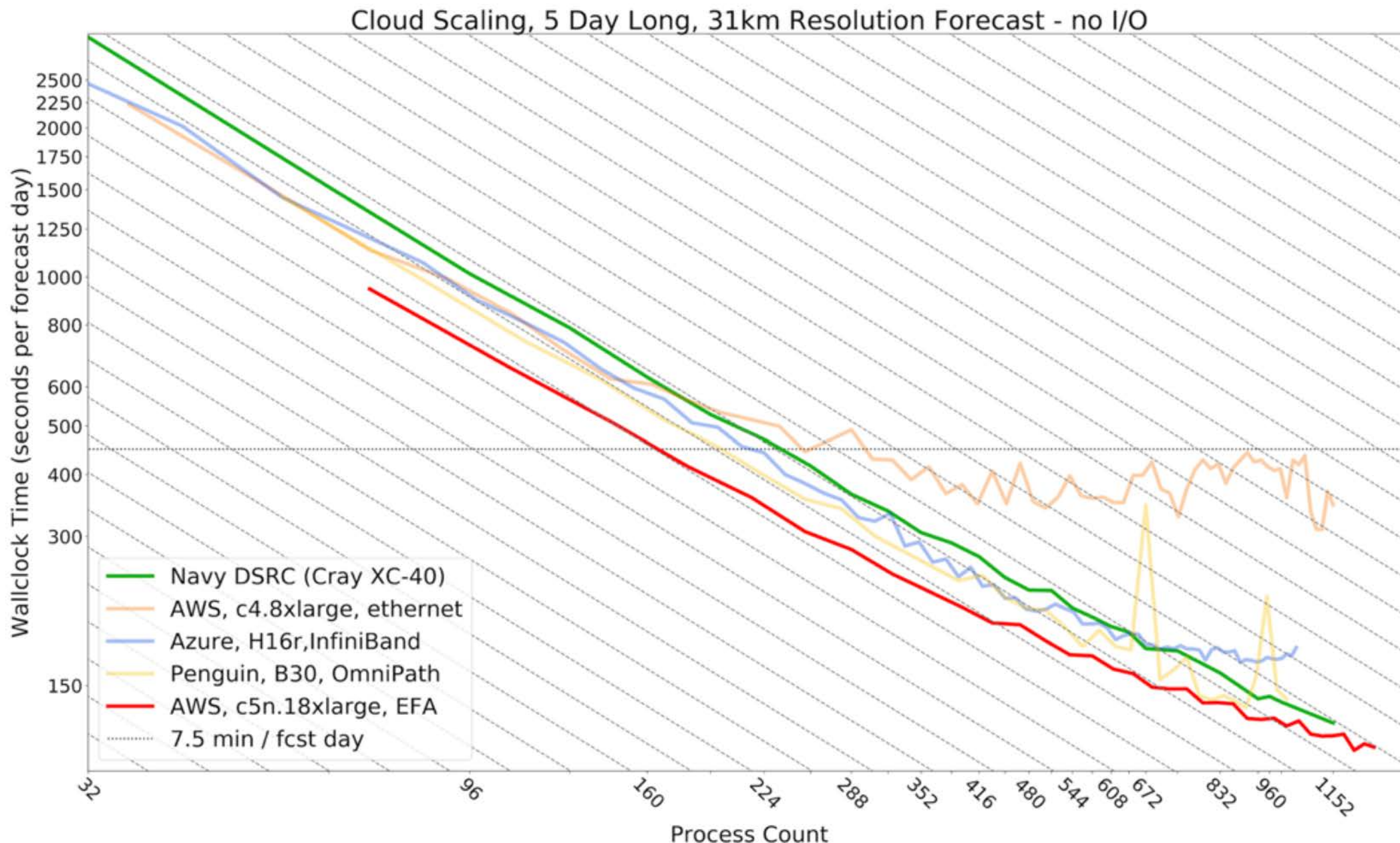
# Low Resolution Forecast: Performance - AWS EC2 EFA

## Platform Specifications:

- 3.0 GHz Intel Xeon Platinum w/ AVX-512
- 36 core nodes
- AWS Elastic Fabric Adapter
- Amazon Linux 2

## Results

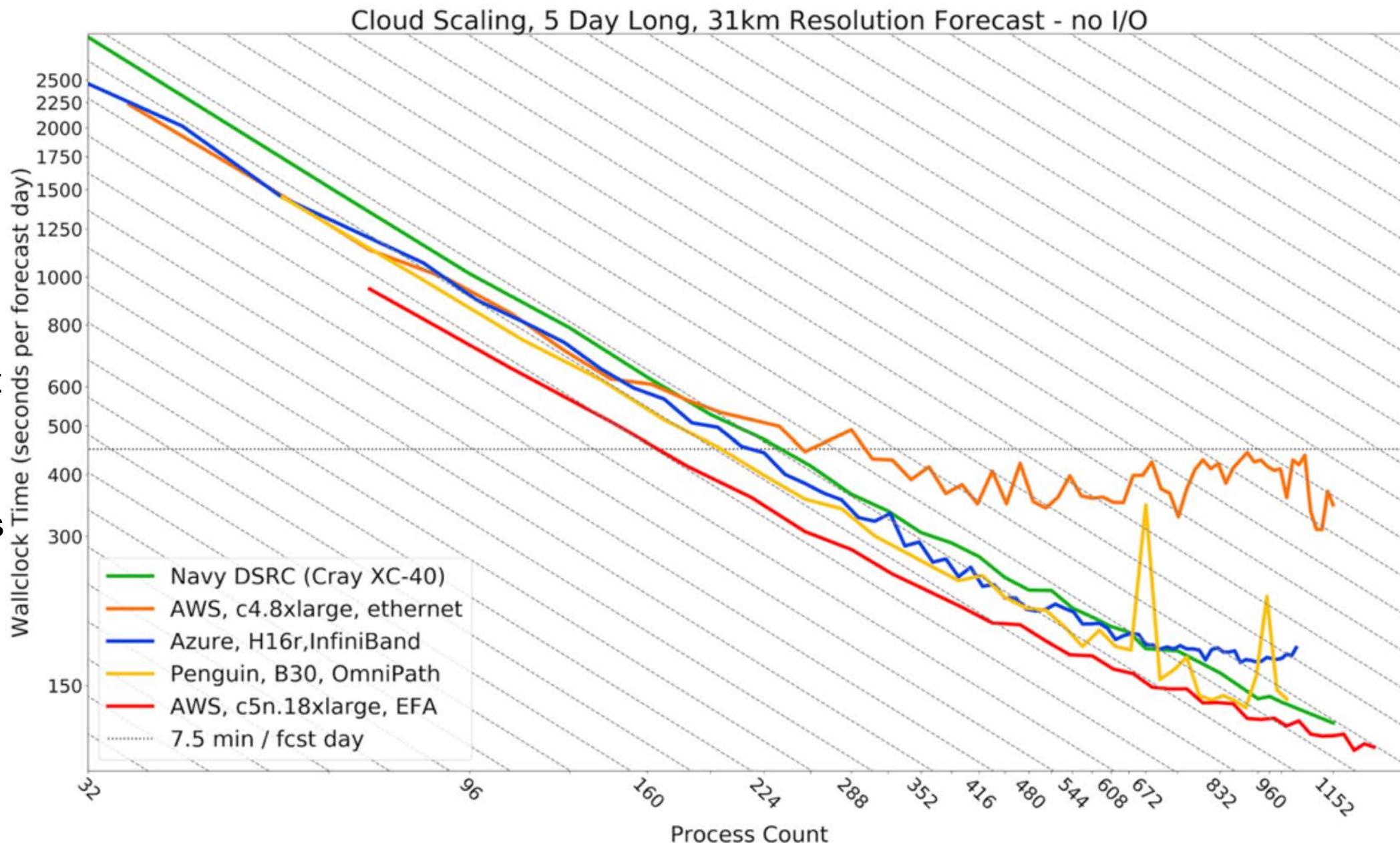
- Smooth scaling up to ~450 cores.
- Slight variability beyond – “stair step” performance beyond.
- Required use of Open MPI; Intel MPI now supported but not yet tested successfully with NAVGEM.



# Low Resolution Forecast: Performance - Comparison

## Performance Improvements: c5n with EFA on AWS EC2

- At the highest core counts:
  - 13% faster than Penguin
  - 39% faster than Azure
  - 192% faster than previous AWS
  - 6% faster than Navy DSRC
- Min size estimated to meet 7.5 min:
  - 6% faster than Penguin
  - 16% faster than Azure
  - 74% faster than previous AWS
  - 29% faster than Navy DSRC
- Min size cost estimate:
  - Penguin: \$12.95
  - Azure: \$21.31
  - Previous AWS: \$18.65
  - C5n with EFA: \$13.76

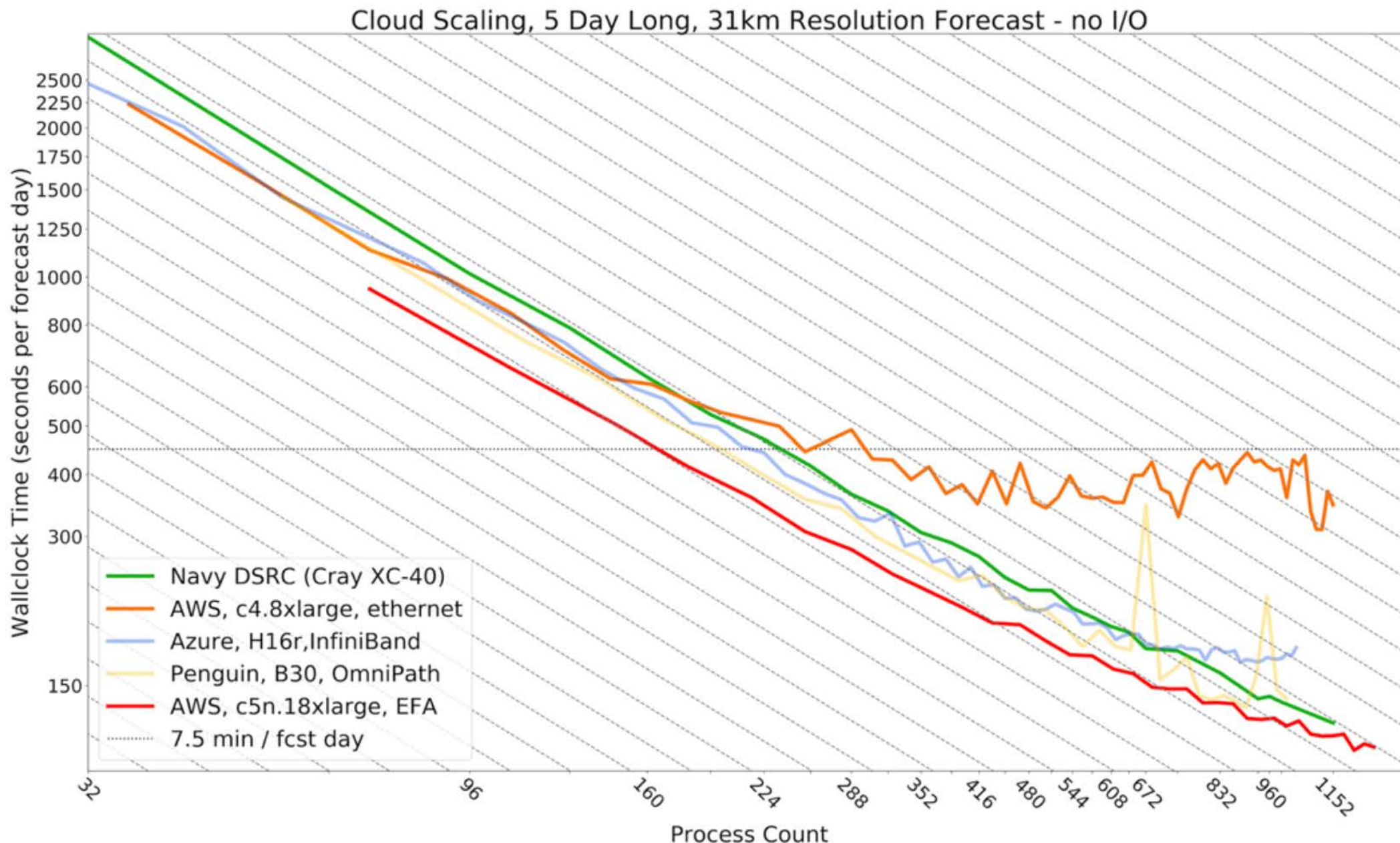




# Elastic Fabric Adapter on AWS EC2

## Elastic Fabric Adapter

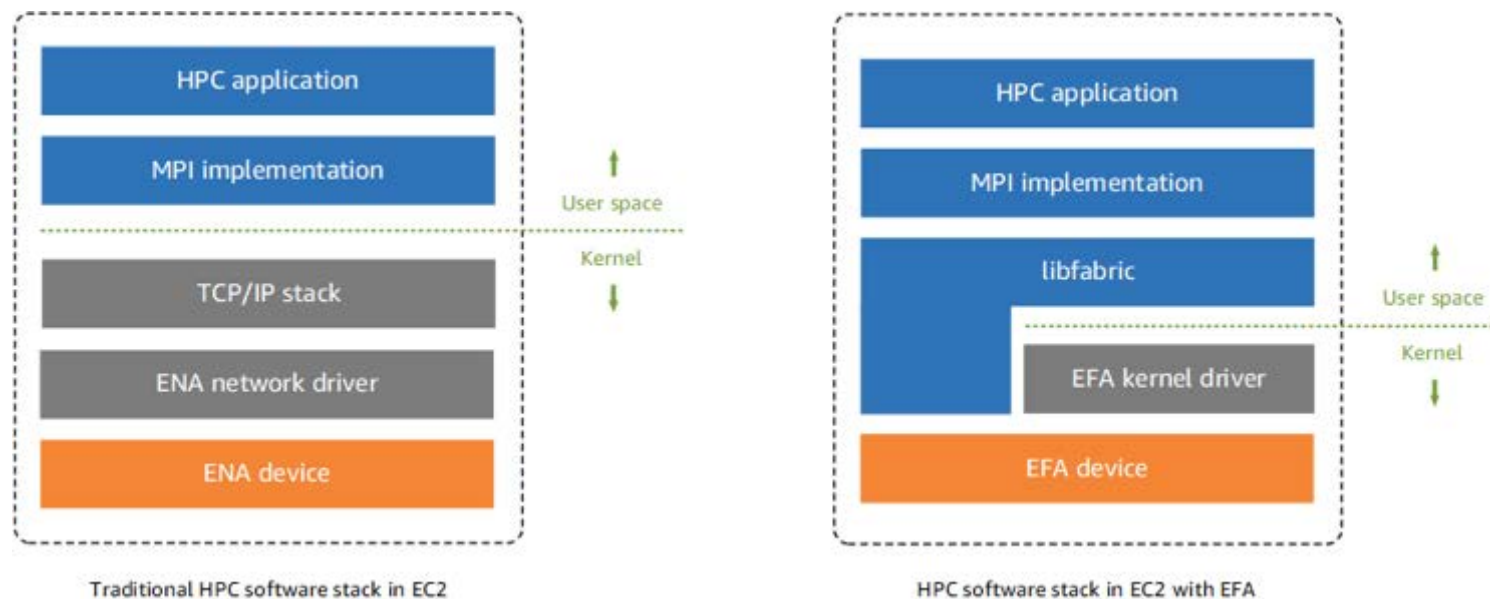
- Updated networking capability launched April 2019
- Hardware:
  - 3<sup>rd</sup> gen Nitro chip
- Software:
  - Scalable Reliable Datagram
- EFA provider has been upstreamed to most recent libfabric release
- Currently available on 4 large instance types.
- Supports Open MPI and Intel MPI.



# Elastic Fabric Adapter on AWS EC2

## Elastic Fabric Adapter

- Updated networking capability launched April 2019.
- Hardware:
  - 3<sup>rd</sup> gen Nitro chip
- Software:
  - Scalable Reliable Datagram
- EFA added as network fabric provider supported by libfabric library.
- Currently available on 4 large instance types.
- Supports Open MPI and Intel MPI.



TCP	InfiniBand	SRD
Stream	Messages	Messages
In-order	In-order	Out-of-order
Single path	Single(ish) path	ECMP spraying with load balancing
High limit on retransmit timeout (>50ms)	Static user-configured timeout (log scale)	Dynamically estimated timeout (usec resolution)
Loss-based congestion control	Semi-static rate limiting (limited set of supported rates)	Dynamic rate limiting
Inefficient software stack	Transport offload with scaling limitations	Scalable transport offload (same number of QPs regardless of cluster size)



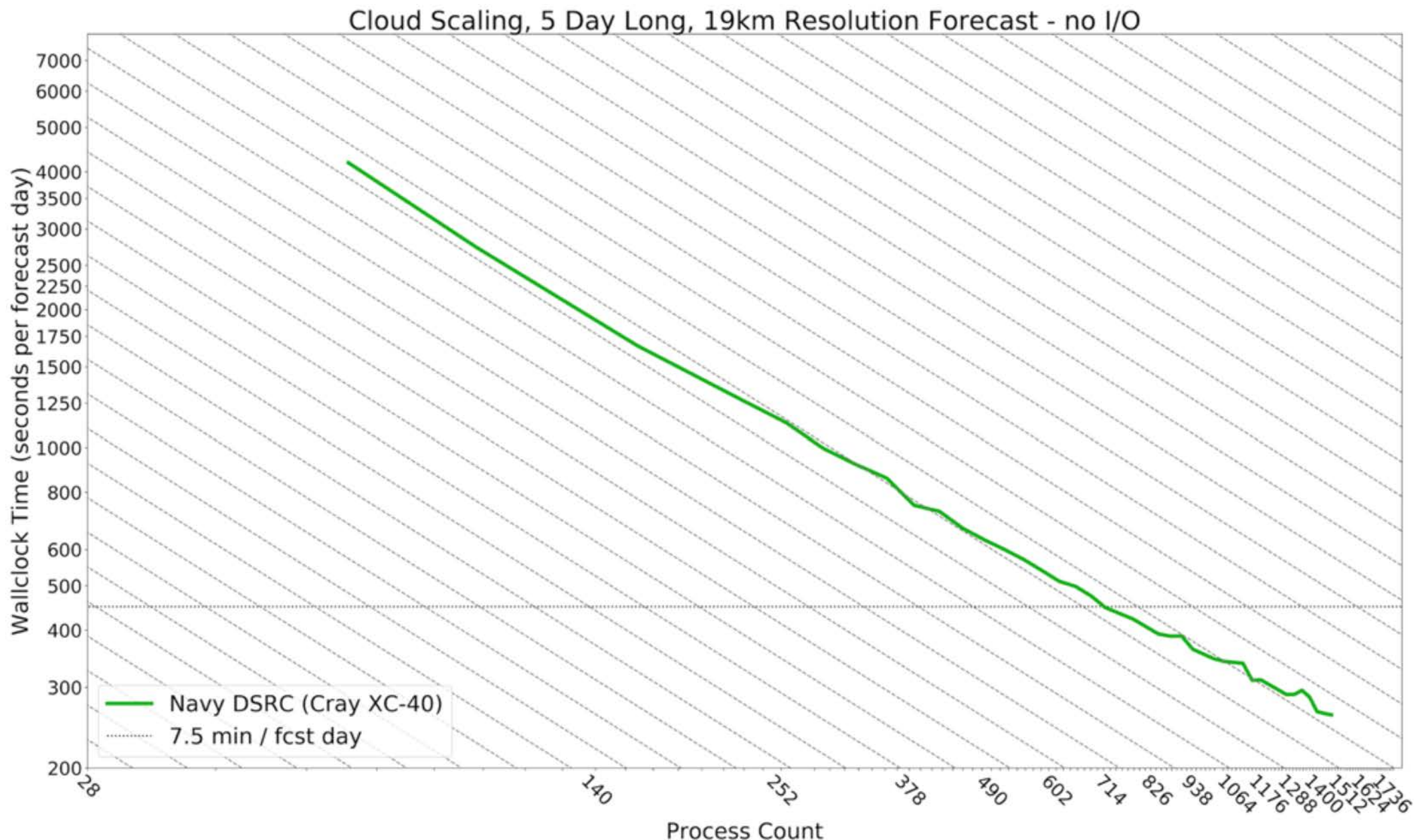
# High Resolution Forecast: Performance - Navy DSRC

## Platform Specifications:

- 2.3 GHz Intel Xeon E5-2698 v3 Broadwell
- 32 core nodes
- Cray Aries / Dragonfly
- Cray Linux

## Results

- Good scaling maintained on Navy resources.
- Larger cluster sizes required to meet 7.5 min/day standard.



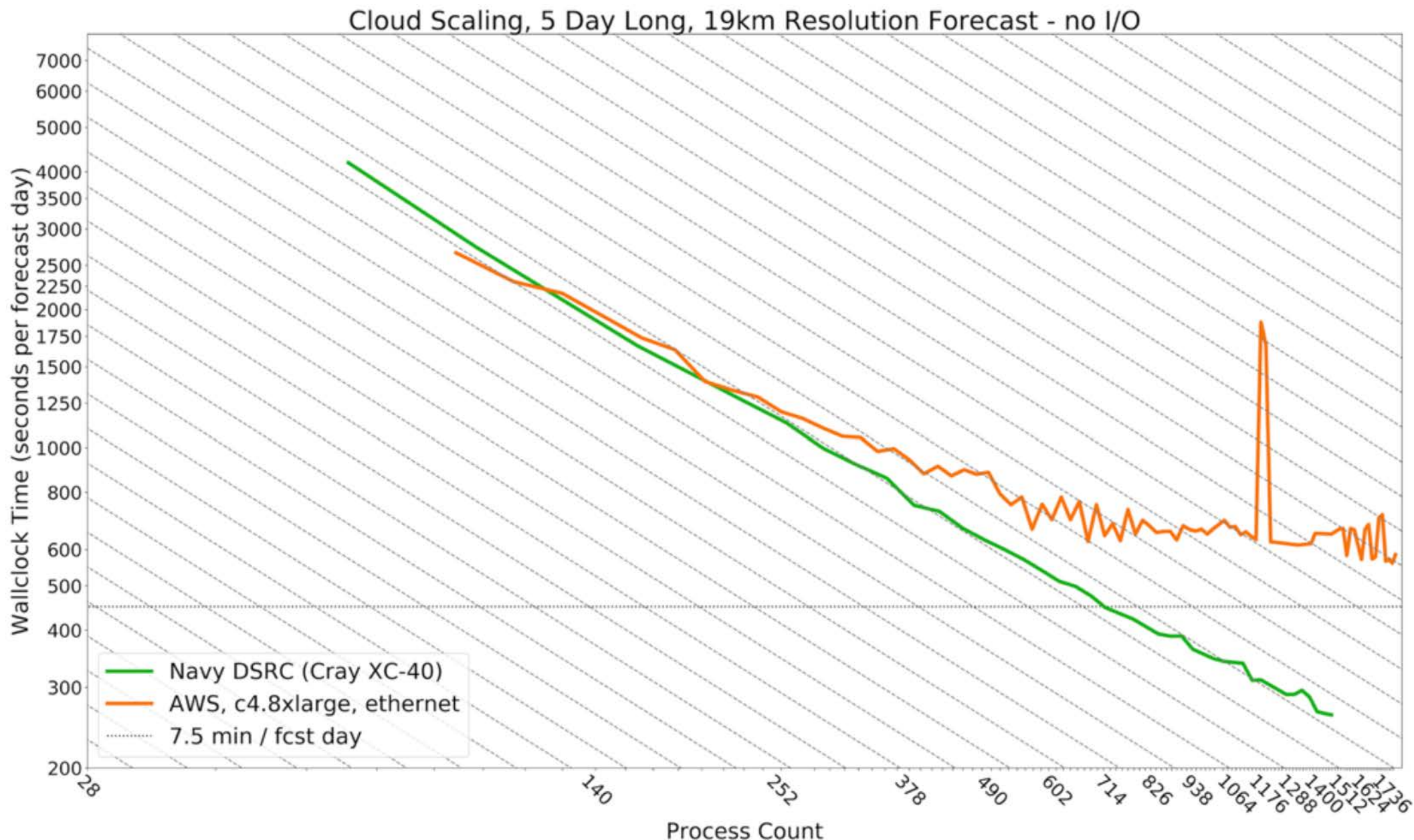
# High Resolution Forecast: Performance - AWS EC2 1

## Platform Specifications:

- 2.9 GHz Intel Xeon E5-2666 v3 Haswell
- 18 core nodes
- 25 Gbps ethernet with SRIOV
- Amazon Linux

## Results

- Variability and performance flattening delayed.
- Cannot meet 7.5 min/day goal.





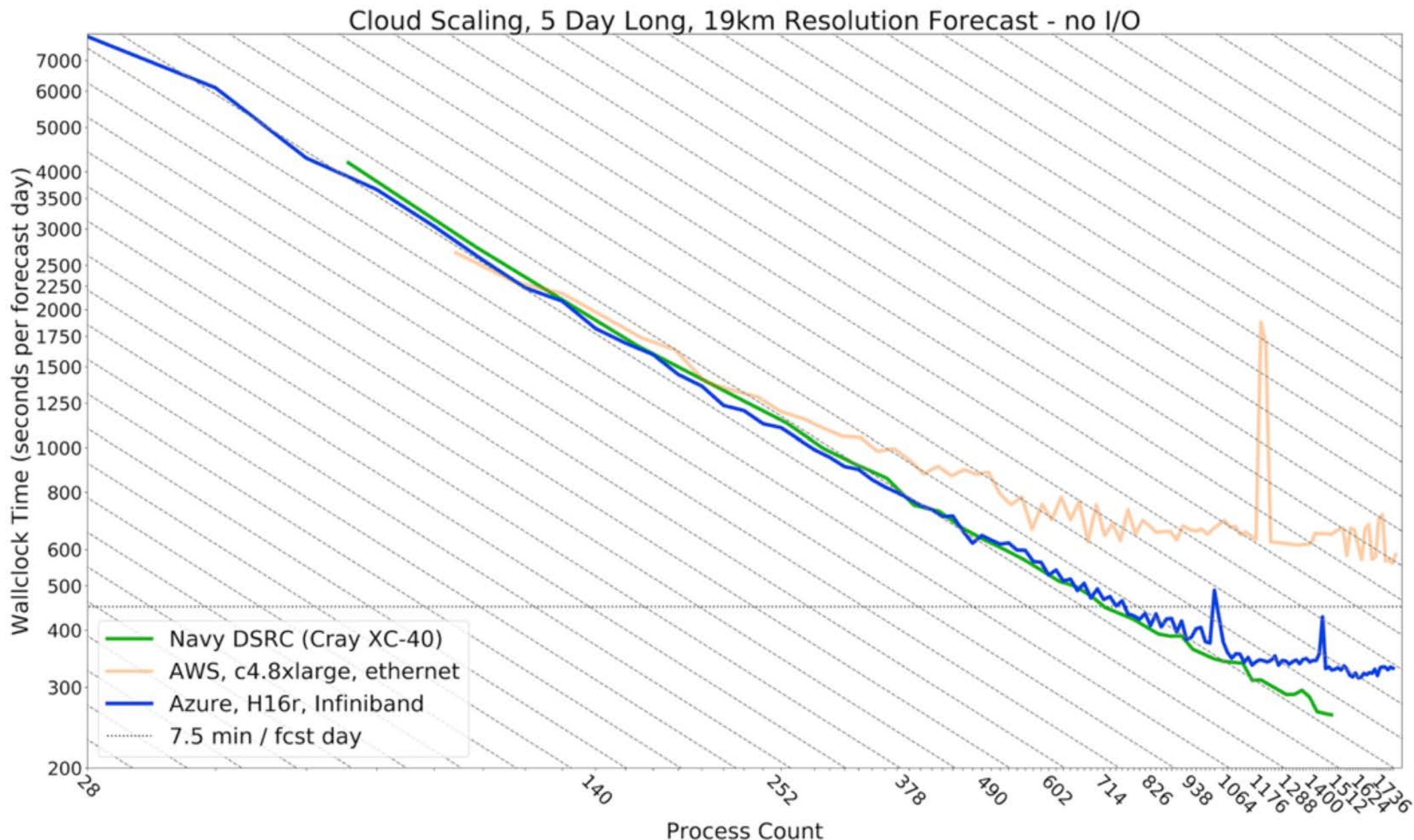
# High Resolution Forecast: Performance - Azure

## Platform Specifications:

- 3.2 GHz Intel Xeon E5-2667 v3
- 14 core nodes
- FDR Infiniband
- CentOS Linux

## Results

- Similar improvement in variability and performance over AWS compared to low resolution forecast.
- Slight performance improvement over Navy DSRC in low resolution forecast has disappeared.



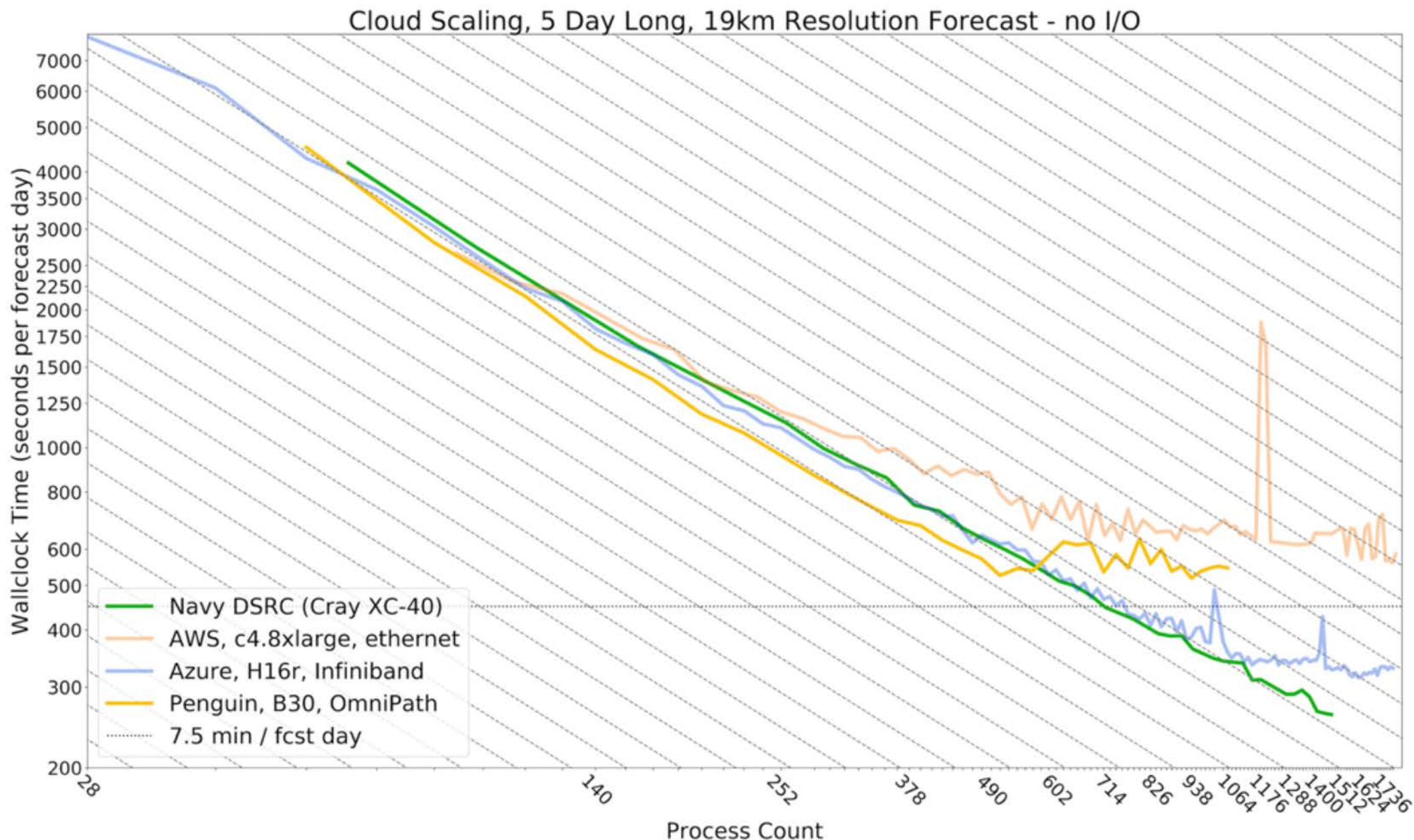
# High Resolution Forecast: Performance - Penguin

## Platform Specifications:

- 2.4 GHz Intel Xeon E5-2680 v4 Broadwell
- 28 core nodes
- Intel OmniPath

## Results

- Good scaling and slight improvement over Navy DSRC until ~500 cores.
- Unable to obtain larger cluster runs to due to system limitations and excessive job queue wait times.





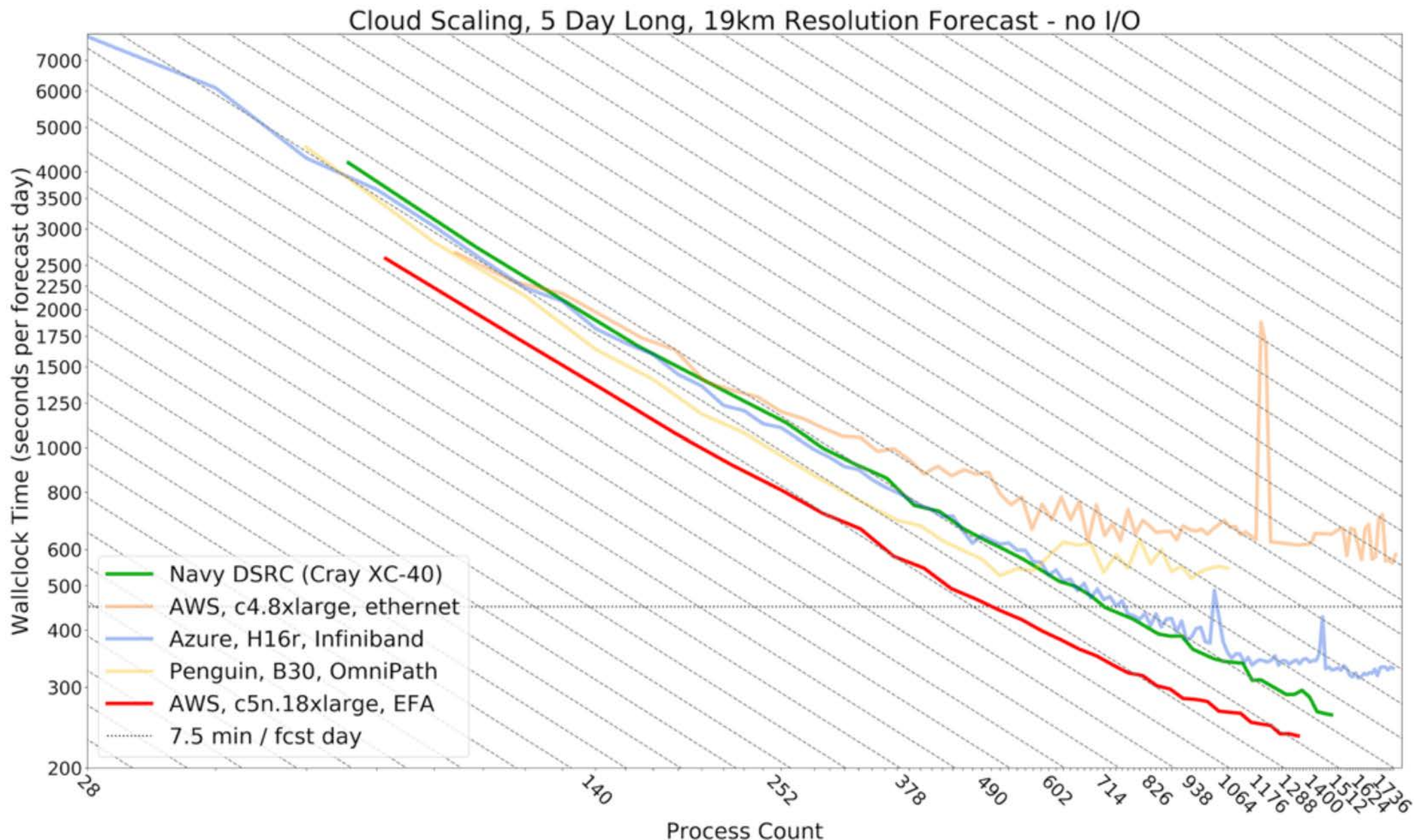
# High Resolution Forecast: Performance - AWS EC2 EFA

## Platform Specifications:

- 3.0 GHz Intel Xeon Platinum w/ AVX-512
- 36 core nodes
- AWS Elastic Fabric Adapter
- Amazon Linux 2

## Results

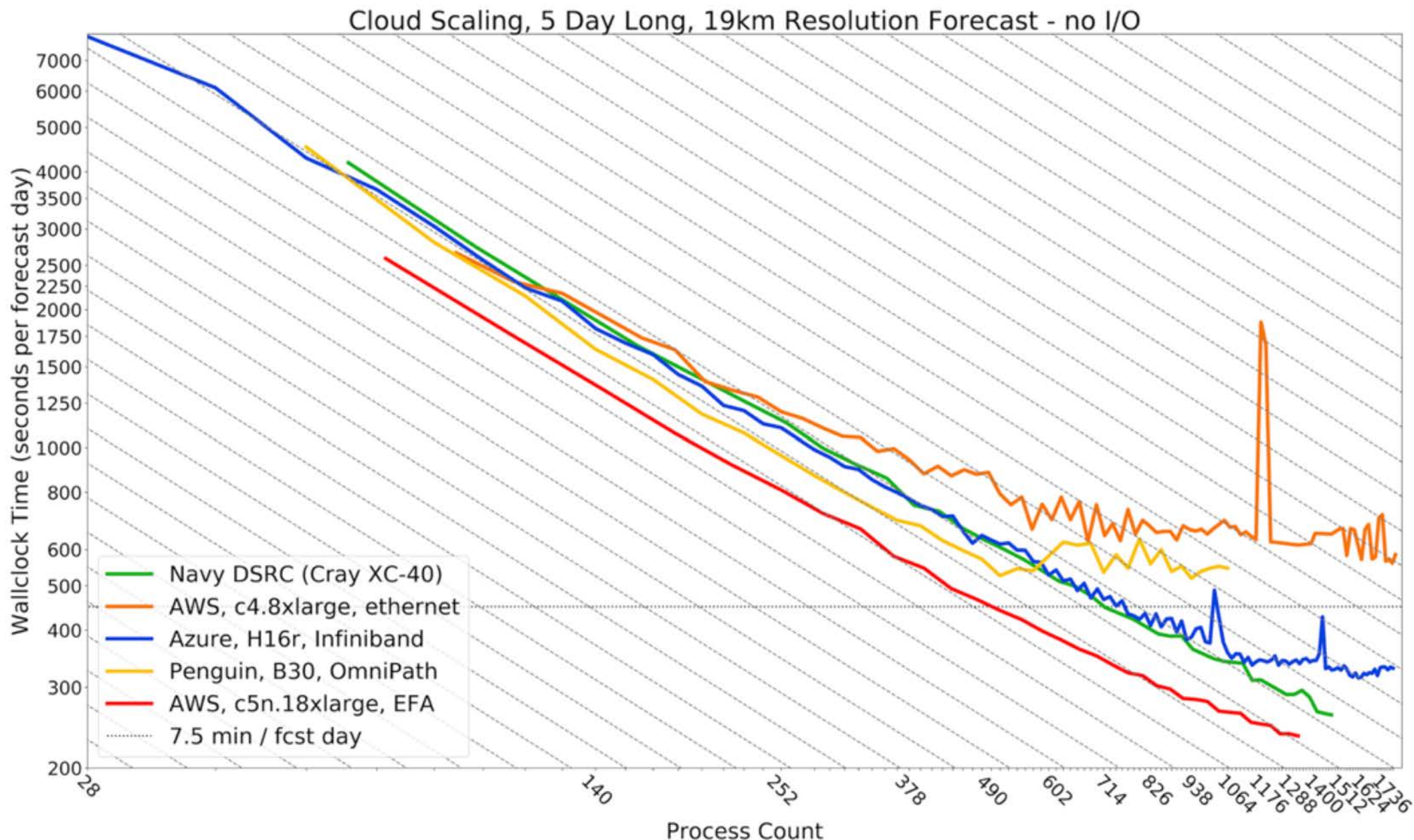
- Similar performance improvement over all platforms as exhibited in low resolution forecast.



# High Resolution Forecast: Performance - Comparison

## Performance Improvements: C5n with EFA on AWS EC2

- At the highest core counts:
  - 107% faster than Penguin
  - 43% faster than Azure
  - 160% faster than previous AWS
  - 25% faster than Navy DSRC
- Min size estimated to meet 7.5 min:
  - 33% faster than Azure
  - 23% faster than Navy DSRC
- Min size forecast cost estimate:
  - Azure: \$82.97
  - C5n with EFA: \$44.02

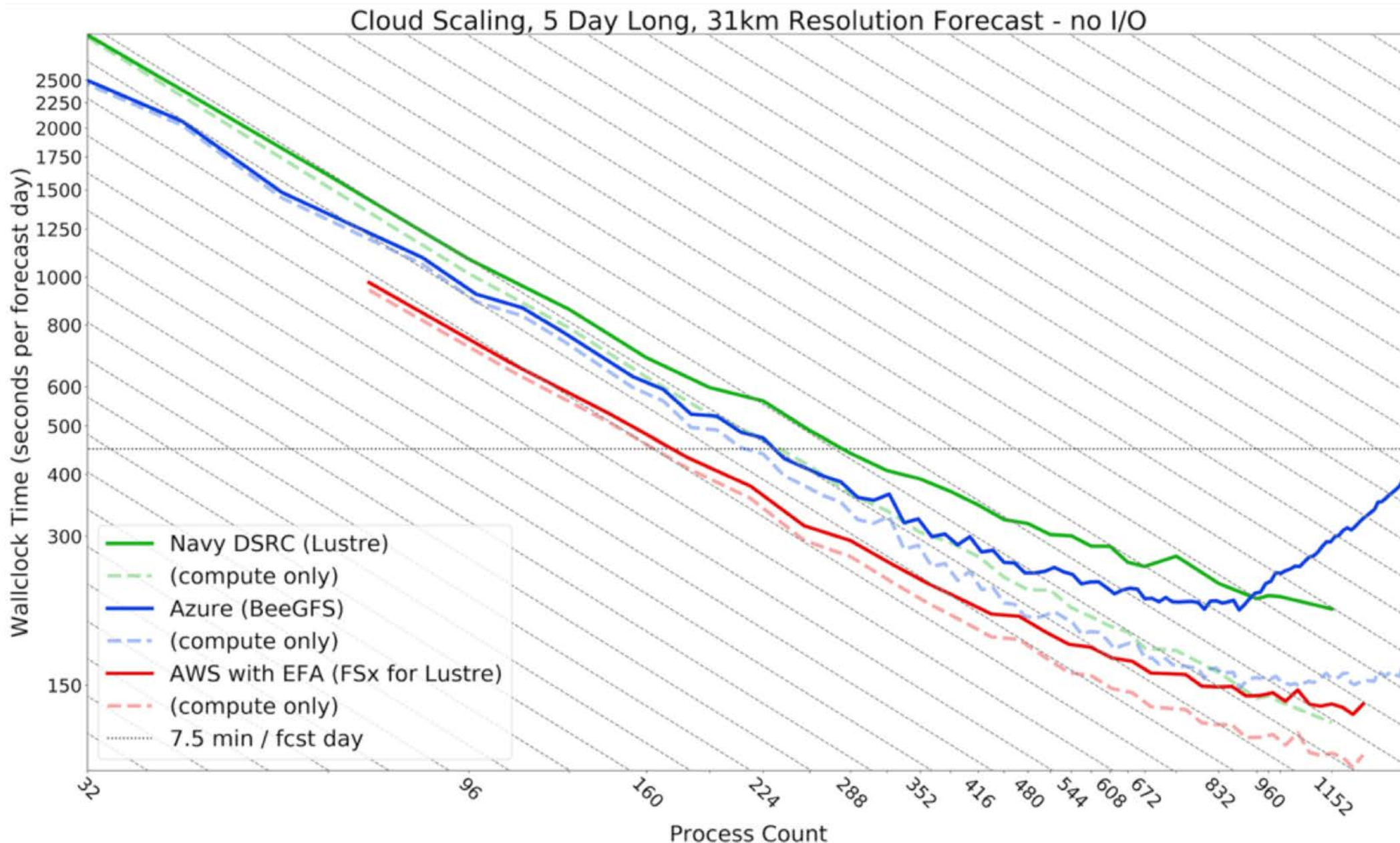




# Low Resolution Forecast with I/O - Comparison

## Incorporating I/O: T425 31 km (low) resolution

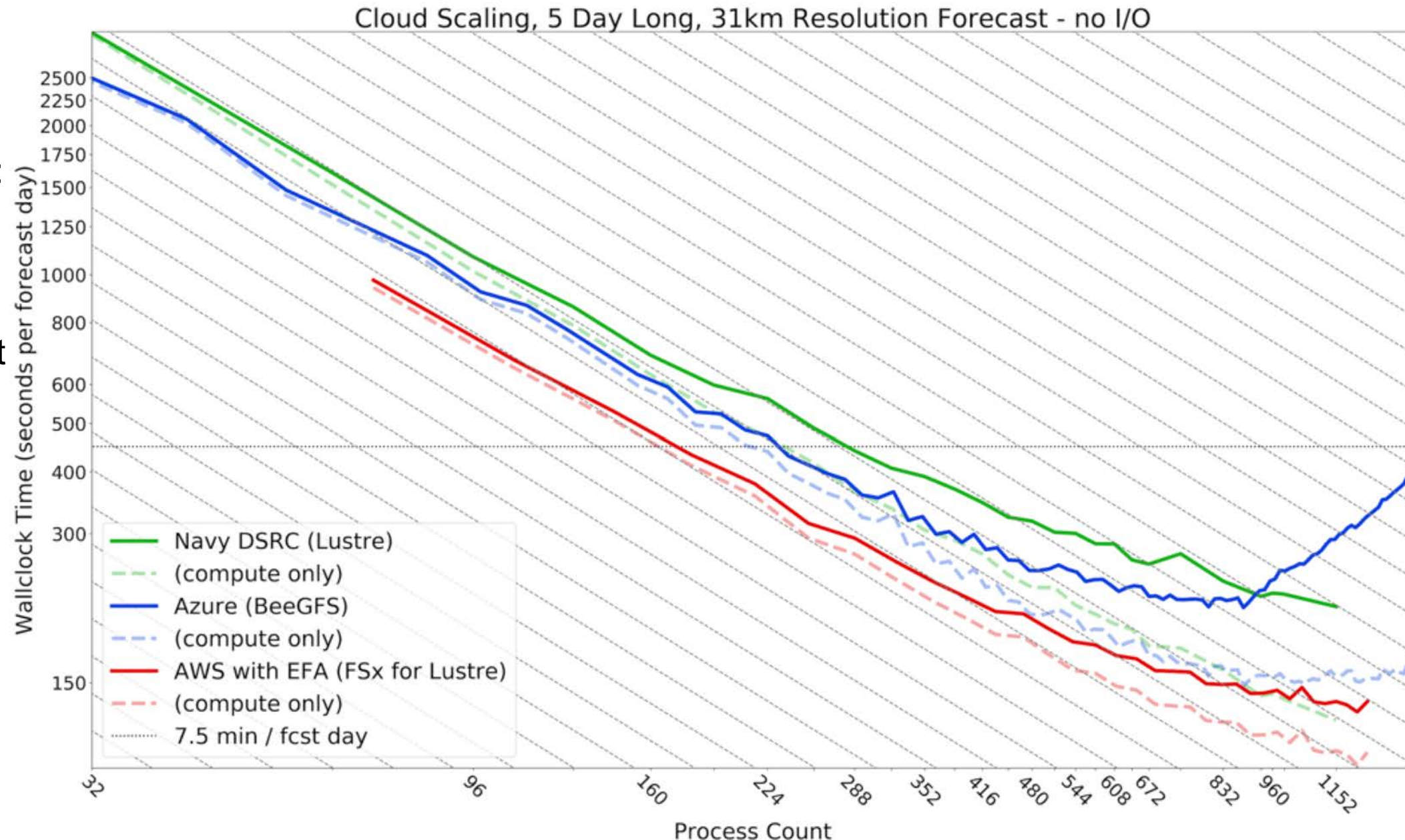
- NAVGEM can output checkpoint data files at specified intervals:
  - ~1.2 GB in size
  - One every sim hour for the first 10 sim hours
  - One every 3 sim hours for the remainder
- All vendors tested offered parallel file systems, we tested on AWS and Azure



# Low Resolution Forecast with I/O - Comparison

## Incorporating I/O: C5n with EFA on AWS EC2

- At the highest core counts:
  - 35% faster than Azure
  - 36% faster than Navy DSRC
- Min size estimated to meet 7.5 min:
  - 18% faster than Azure
  - 33% faster than Navy DSRC
- On AWS, additional cost is minor (<\$1) for small scratch parallel file systems provisioned only for length of model run

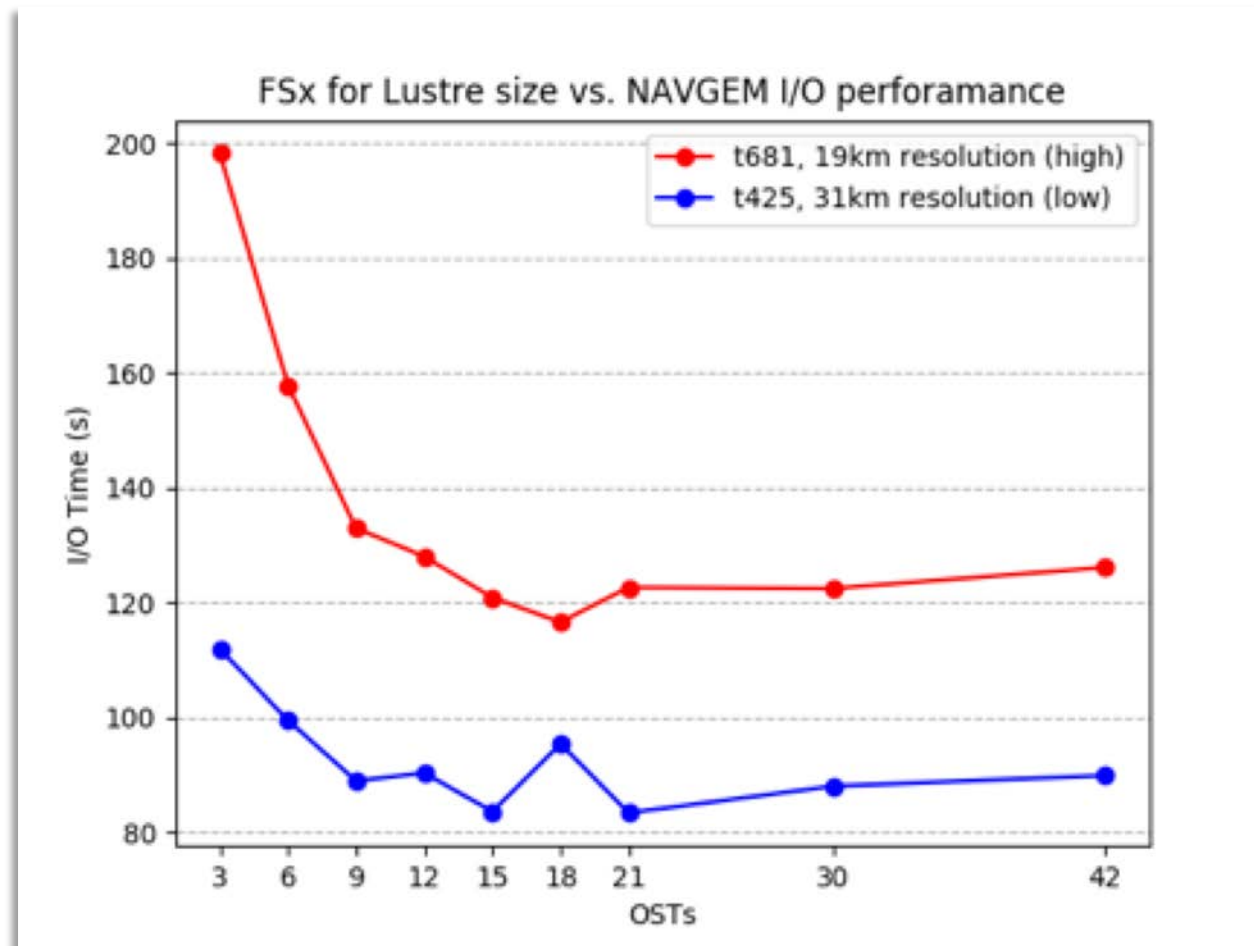




# Low Resolution Forecast with I/O - Comparison

## Incorporating I/O: C5n with EFA on AWS EC2

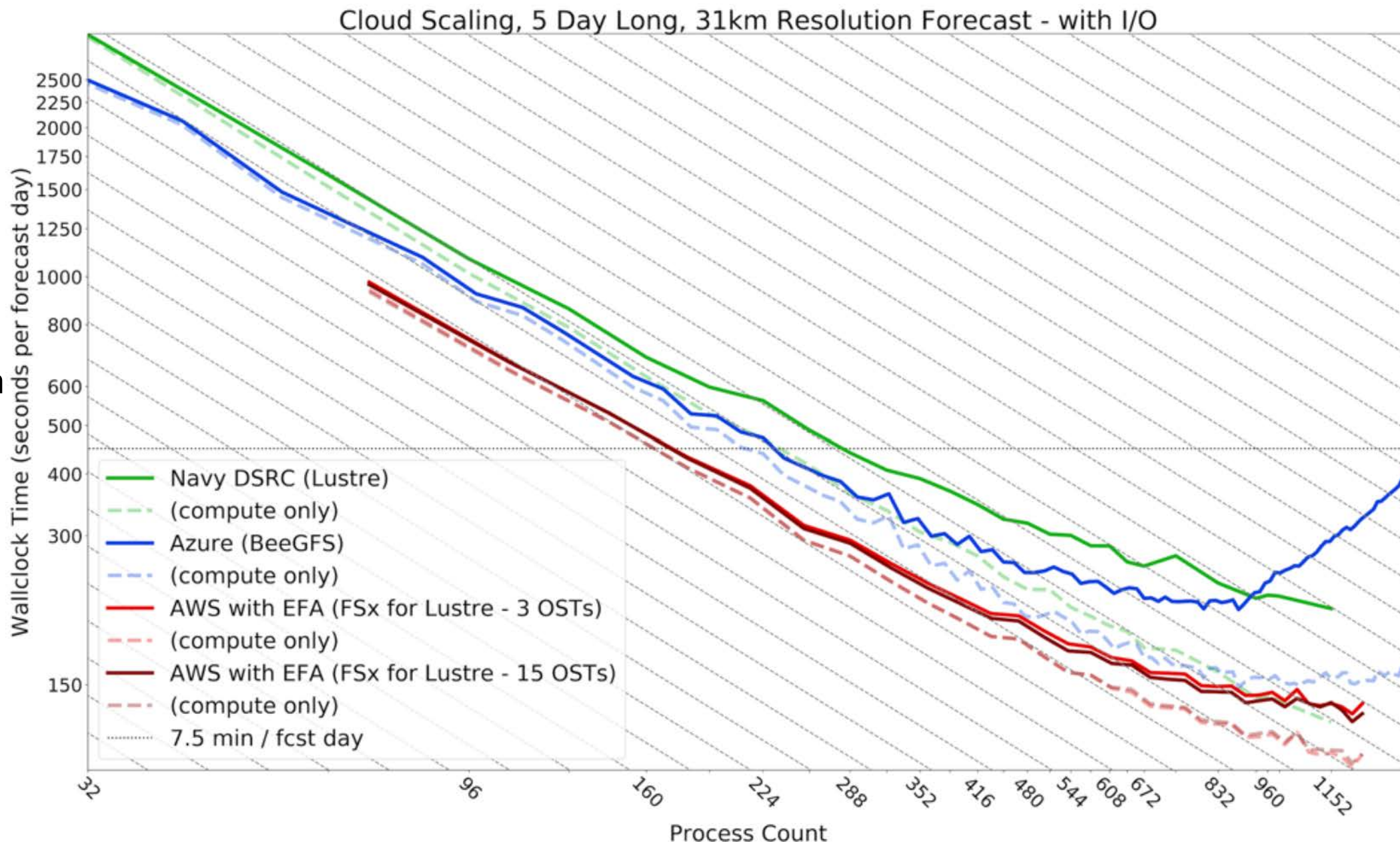
- Cloud allows for “custom” parallel file systems
- FSx for Lustre on AWS: throughput is a function of size provisioned (i.e. larger file system = more Lustre OSTs = higher data throughput)
- Exclusive use offers chance to tune configurations to particular model data



# Low Resolution Forecast with I/O - Comparison

## Incorporating I/O: C5n with EFA on AWS EC2

- Utility of optimizing file system size is questionable.
- Minor improvement:
  - At higher core counts, larger file system only an average of 10 sec faster for a 5-day forecast
- Cost difference can add up:
  - 3 OSTs: + \$0.70/hr
  - 15 OSTs: + \$3.49/hr





## Future Areas of Research

- Test full NAVGEM ensemble
- Test next-generation forecast programs – NEPTUNE
- Incorporate on-going updates to cloud systems to further reduce costs and optimize performance.



U.S. Naval Laboratory Marine Meteorology Division, Monterey, CA

Email contact: [daniel.arevalo@icloud.com](mailto:daniel.arevalo@icloud.com)



U.S. Naval Laboratory Marine Meteorology Division, Monterey, CA