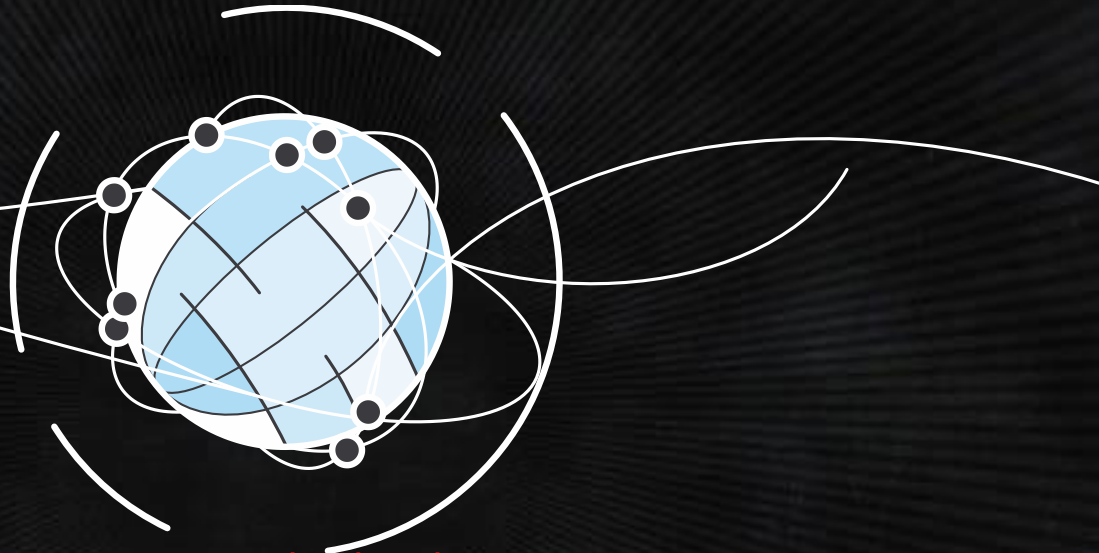# Scaling Tightly Coupled Algorithms on AWS

Dr. Scott Eberhardt

Principle Solutions Architect – HPC, AWS
Visiting Reader, Imperial College

aws

**Research Computing @ AWS**

AWS Worldwide Research & Technical Computing

IT'S ABOUT SCIENCE [NOT SERVERS]

AWS RESEARCH CLOUD PROGRAM
aws.amazon.com/rcp
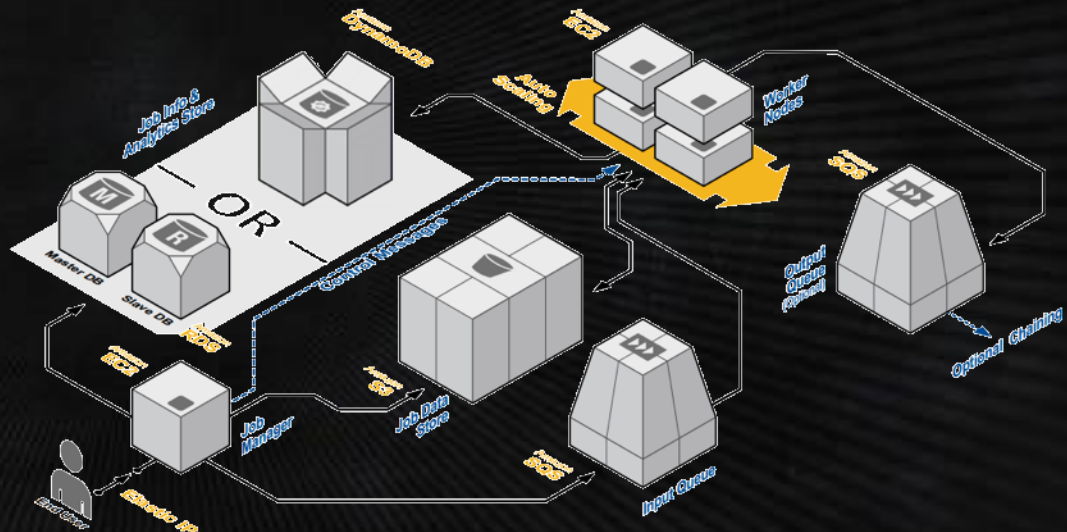
# IT'S ABOUT SCIENCE, NOT SERVERS.

aws.amazon.com/rcp

**#AWSresearchcloud**

aws

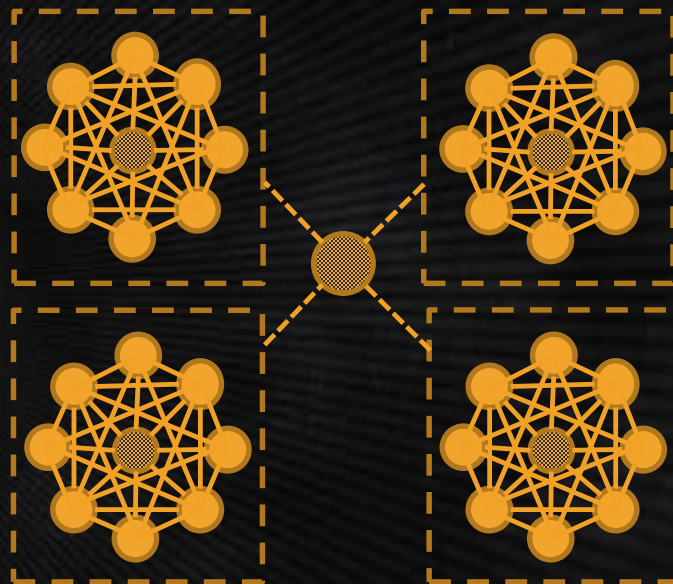# Great Features for HPC Workloads

➤ Experimentation without Fear!

➤ Start and stop instances!

➤ Spot Pricing

➤ Continuous Updates

    ➤ Compute

    ➤ Network

    ➤ Storage

    ➤ Services

# Cloud Improves Workload Throughput

- Run many Jobs in Parallel
  - Eliminate HPC resource contention
  - Eliminate queue wait
  - Use it when you need it
- Right-size clusters and resources
  - Optimize each workload for performance
  - Pay for only what you use

aws

# Cost advantages

## On Premises
Capital Expense Model

- High upfront capital cost
- High cost of ongoing support

## Amazon Web Services
Pay As You Go Model

- Use only what you need
- Multiple pricing models

# Popular HPC workloads on AWS
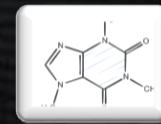
Genomics Processing

Modeling and Simulation

Government and Educational Research

Monte Carlo Simulations

Transcoding and Encoding

Computational Chemistry

... and many more

aws

# Defining HPC – example use cases

**Clustered** (Tightly coupled)

**Data Light**

Minimal requirements for high performance storage

**Data Heavy**

Benefits from access to high performance storage

- Fluid dynamics
- Weather forecasting
- Materials simulations
- Crash simulations

- Seismic processing
- Metagenomics
- Astrophysics
- Deep learning

- Risk modeling
- Molecular modeling
- Contextual search
- Logistics simulations

- Animation and VFX
- Semiconductor verification
- Image processing/GIS
- Genomics

**Distributed / Grid** (Loosely coupled)

aws

Global Infrastructure
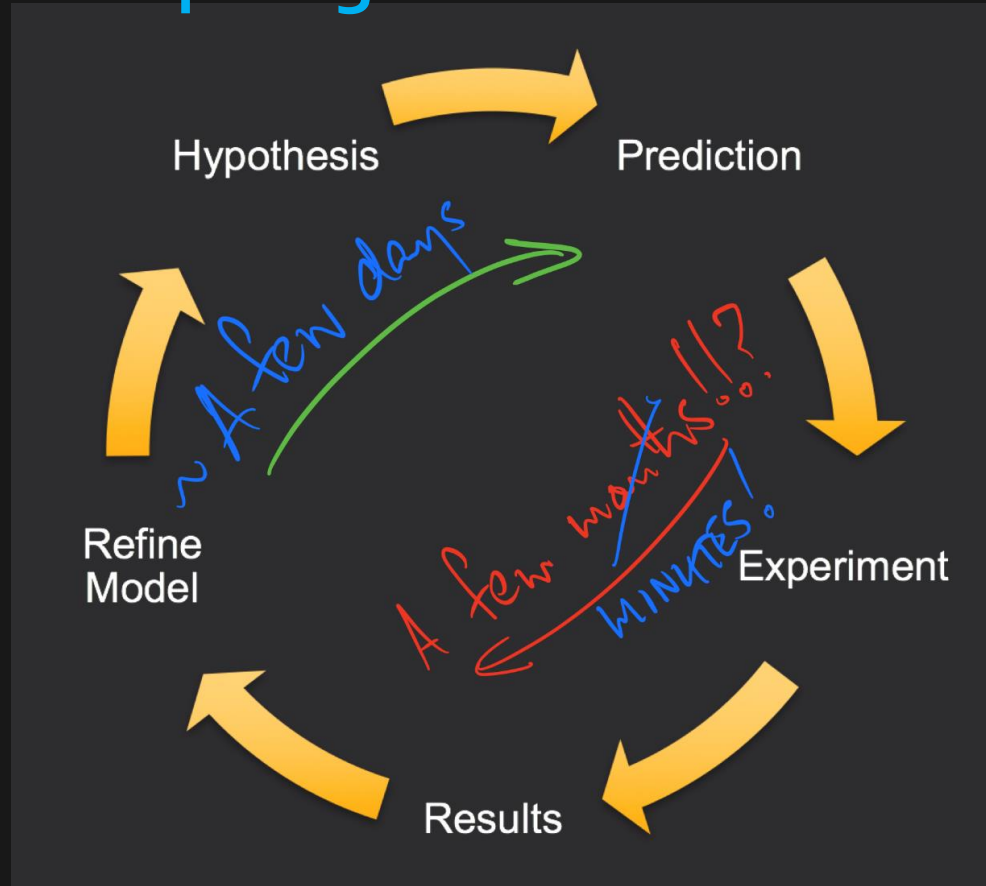
# Important enablers for HPC on the cloud

- Compute performance – CPUs, GPUs, FPGAs

- Memory performance – high RAM requirements in many applications

- Network performance – throughput, latency, and consistency

- Storage performance – including shared filesystems

- Automation and cluster/job management

- Remote graphics for interactive applications

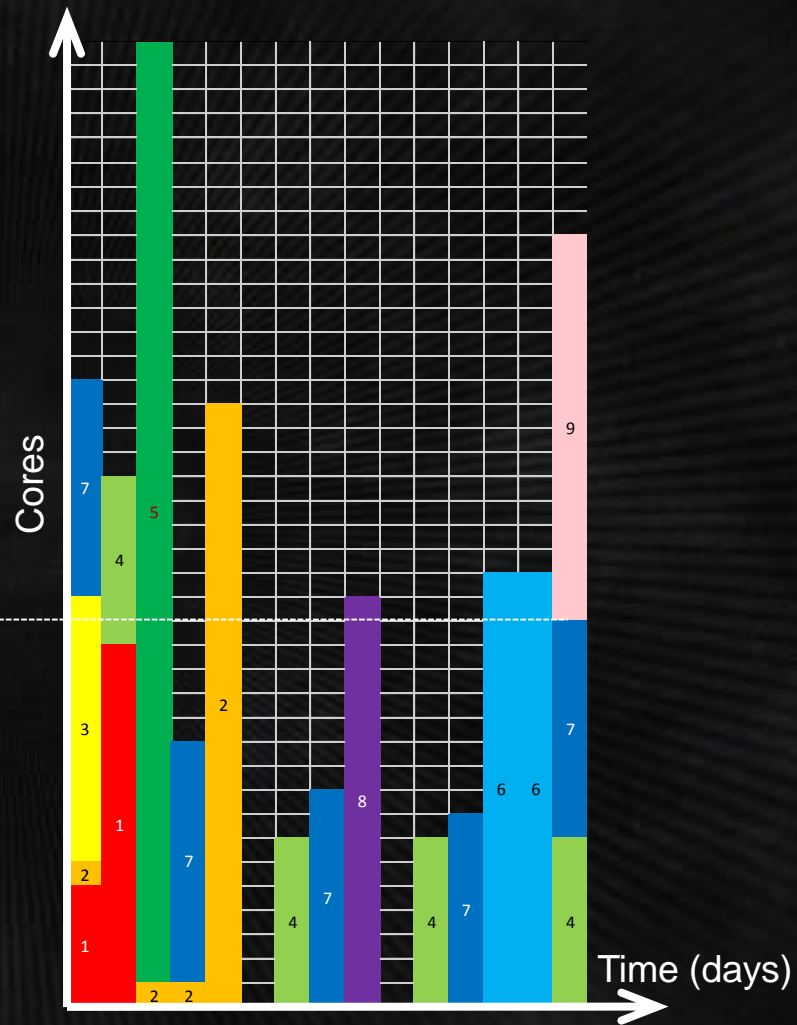- ISV support – including license management
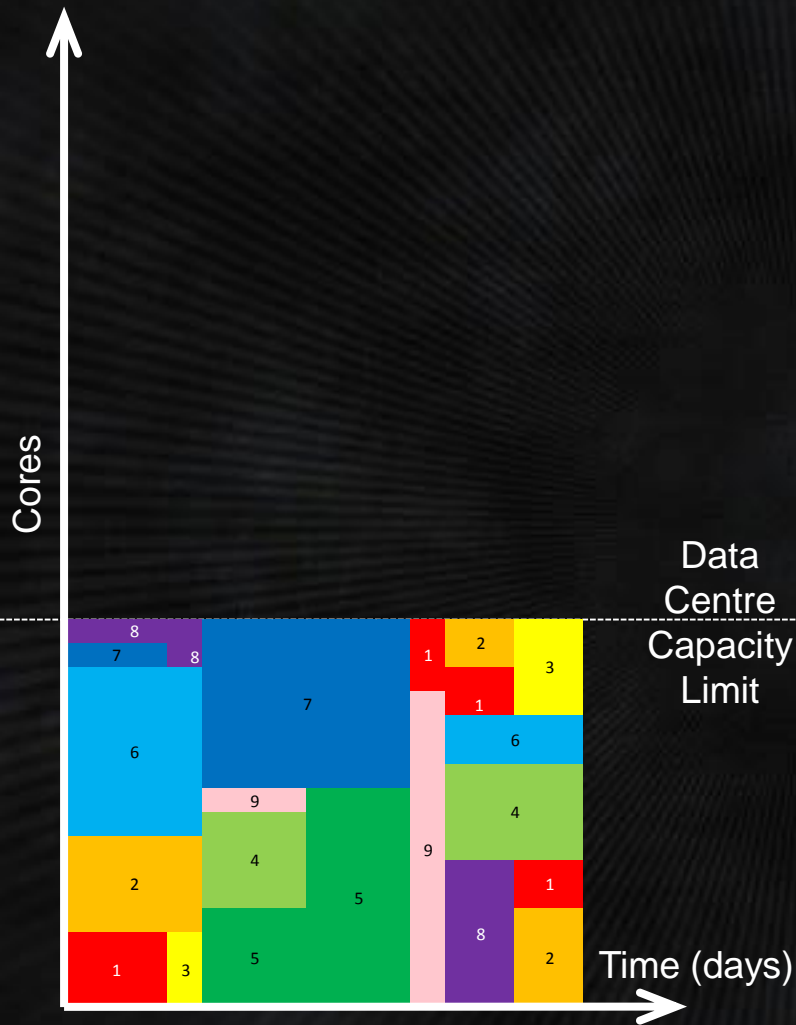
  ...and SCALE

aws

# The Scientific Computing Method

Credit: Aristotle

Cores

Data Centre Capacity Limit
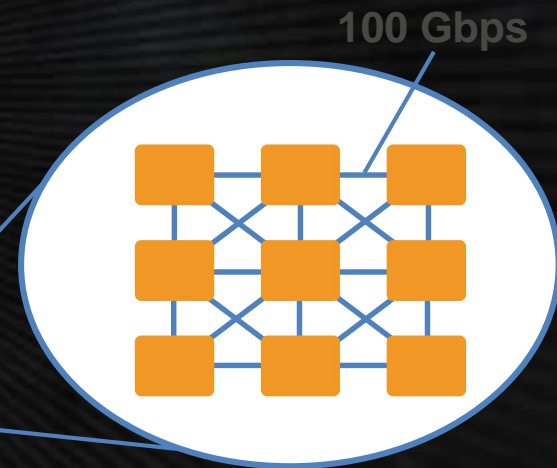
Time (days)

Cores
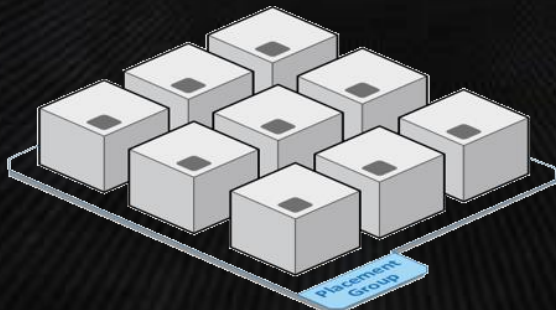
Time (days)

# Deploy multiple HPC clusters

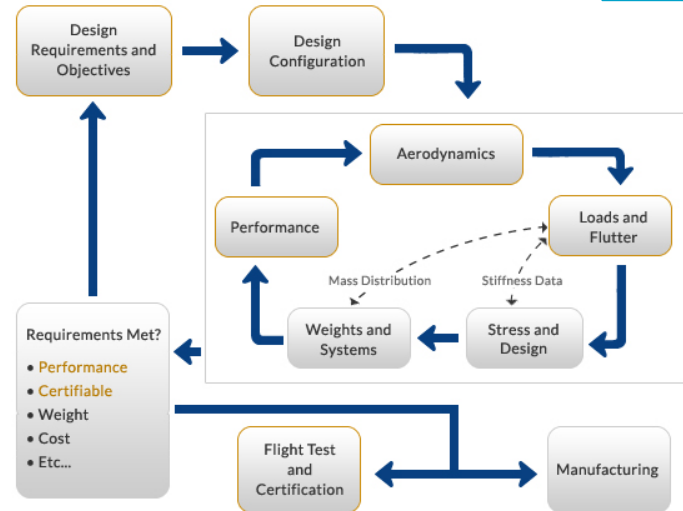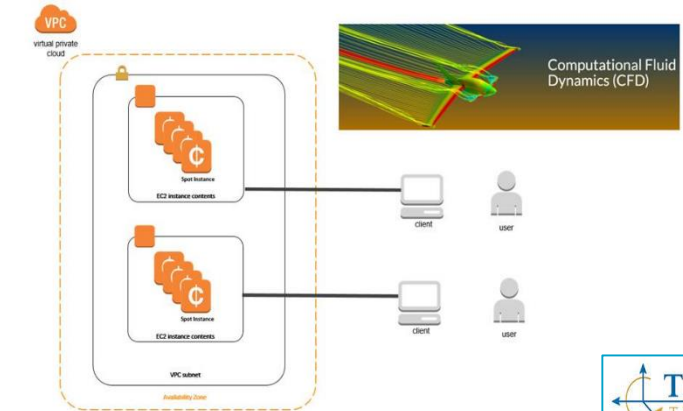Running at the same time, and tuned for each workload

# Which architecture Do I choose?

➤ One size does not fit all

➤ Architectures are an opportunity for optimization

➤ The chosen architecture depends on:
  ➤ The desired user experience
  ➤ The desired deployment method
  ➤ The characteristics of the application

100 Gbps

Placement Group

aws

# Example in aerospace

- Running parallel CFD studies using Siemens STAR-CCM+
  - Goal: shorten the time between Design Requirements and Configuration, and Flight Testing
- 1000+ cores per CFD study, multiple studies required for each workflow iteration
- Job-level optimizations:
  - Enhanced Networking, Placement Groups
  - Amazon Linux, Hyper-threading disabled
- Workflow optimizations:
  - Spot instances, multiple clusters
  - Multiple parallel studies for faster throughput

# Performance considerations

For tightly-coupled cluster workloads



## Test using real-world examples

- Use large cases for testing: do not benchmark scalability using only small examples

## Domain decomposition

- Choose number of cells per core for either pre-core efficiency or for faster results

## Processor states

- Use P-states to reduce processor variability
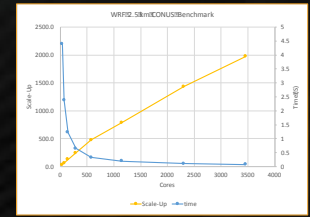
## MPI libraries

- Test with Intel MPI and OpenMPI 4.0, and make use of available tunings

## Network

- Use a placement group

- Enable enhanced networking

## Hyper-threading and affinity

- Test with Hyper-threading (HT) on and off – usually off is best, but not always

- Use CPU affinity to pin threads to CPU cores when HT is off

aws

# Scaling

Amdahl's Scaling

aws

# What is Amdahl's Law?

$$Speedup = \frac{1}{(1-p) + \frac{p}{n}}$$

$p$ – fraction of code that can be run in parallel
$n$ – system resource improvement (number of processors)

$$\text{Ideal:} \lim_{n \to \infty} speedup = \frac{1}{1-p}$$

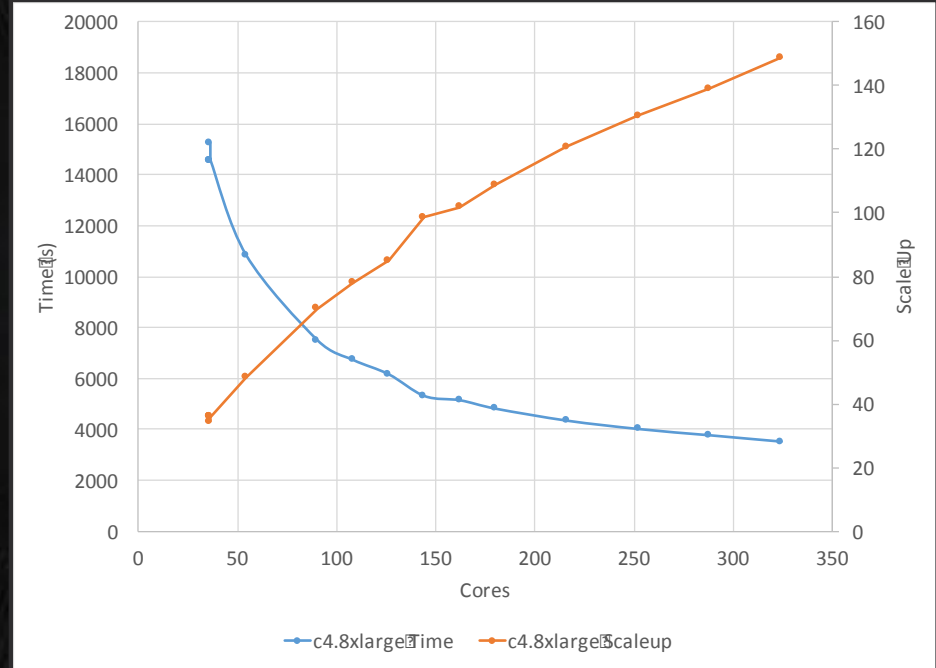$$\text{Or:} \lim_{p \to 1} speedup = n$$

aws

# How is it used?

- Amdahl's Law assumes a fixed problem size.

- Example: a CFD calculation with a fixed number of cells – eg. 40M

- As the number of processors increases, the number of CFD mesh cells per compute core decreases. At some point communication between nodes and cores becomes a bottleneck. This is the driver for low-latency networks.

- Note: The lower the latency the fewer CFD mesh cells per compute node are necessary for good scaling. Lower latency usually means better speedup when considering Amdahl's law.

aws

# Coding Goals

- The following have been routine in code architectures since the first vector computers:

- Choose algorithms where $p$ is maximized
-

- Reprogram codes
- Look for ways to avoid/reduce dependencies
- Consider tradeoffs between recalculating and storing
- Monitor code execution to find bottlenecks

aws

# Structural simulation

# Fluid dynamics – Ansys Fluent

- C4.8xlarge instance type
- 140M cell model
- F1 car CFD benchmark





EXTERNAL FLOW OVER A FORMULA-1 RACE CAR
AWS, C4.8XLARGE

Time/iter (sec)    Speedup    Ideal

aws

# Tightly-coupled HPC – weather

# Resources used in this study

- Archer: Cray XC30 supercomputer
  - two 2.7 GHz, 12-core Intel E5-2697 v2 (Ivy Bridge)
- AWS:
  - z1d: 4.0 GHZ Intel®Xeon®Scalable Processors; 24 core per instance; 16GB Ram per core; 25 Gigabit network bandwidth
  - c5n: 3.0/3.5 GHZ Intel®Xeon®Scalable Processors; 36 core per instance; 5.3 GB Ram per core; 100 Gigabit network bandwidth; New Elastic Fabric Adaptor (EFA) for fast networking

aws

# Methodology

- OpenFOAM v1806 in Double Precision (pimpleFoam)
- Scotch decomposition for solving, hierarchical (i.e constant x/y/z loading) for meshing
- SST-DDES Turbulence Model
- ANSA generated 143/280M cell unstructured mesh
- Time Step=$5e{-}^4$s with 5 inner iterations
- Preconditioned Conjugant Gradient Linear Solver

aws

# Amdahl Scaling: openFoam (pimpleFoam)

**sec/time-step/cell** (y-axis)
- 1.20E-07
- 1.00E-07
- 8.00E-08
- 6.00E-08
- 4.00E-08
- 2.00E-08
- 0.00E+00

**Cores** (x-axis)
- 0.00E+00
- 5.00E+02
- 1.00E+03
- 1.50E+03
- 2.00E+03
- 2.50E+03

Legend:
- z1d (medium mesh)
- z1d (fine Mesh)
- c5n (medium mesh)
- c5n (fine mesh)
- Archer (medium mesh)
- Archer (fine mesh)

**aws**

# Amdahl Scaling, Alternate View

**sec/time-step** (y-axis)

1.80E+01
1.60E+01
1.40E+01
1.20E+01
1.00E+01
8.00E+00
6.00E+00
4.00E+00
2.00E+00
0.00E+00

**Cores** (x-axis)

0.E+00   5.E+02   1.E+03   2.E+03   2.E+03   3.E+03

Legend:
- z1d (medium)
- z1d (fine)
- c5n (medium)
- c5n (fine)
- Archer (medium)
- Archer (fine)

# Amdahl Scaling, 2nd Alternate View

# Scaling

Gustafson-Barsis Scaling (aka Gustafson's Law)
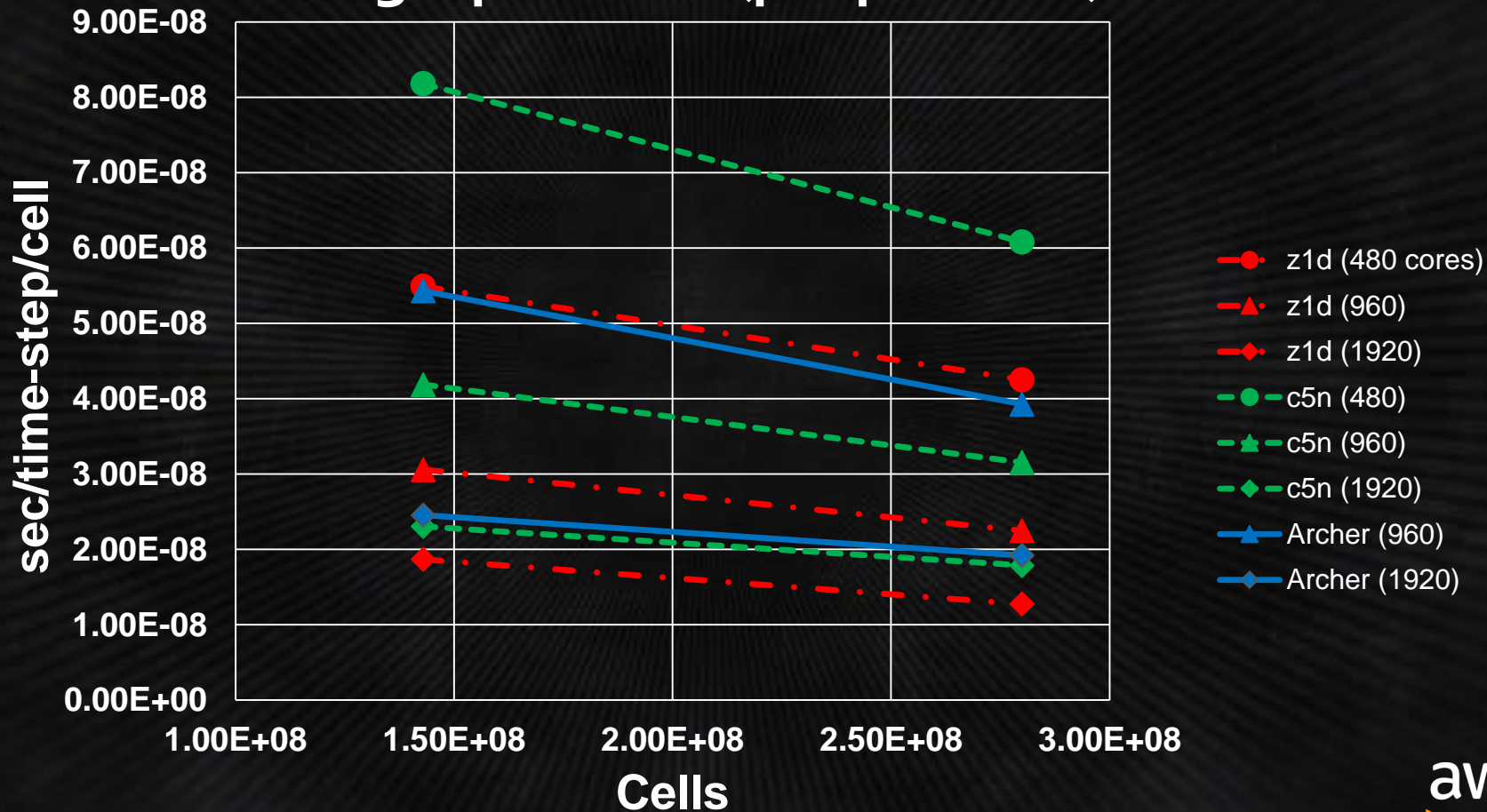
aws

# What is Gustafson's Law?

$$Speedup = (1 - p) + \frac{p}{n}$$

Gustafson's Law assumes a workload will increase in size linearly with the number of processors.  (scalable workload)
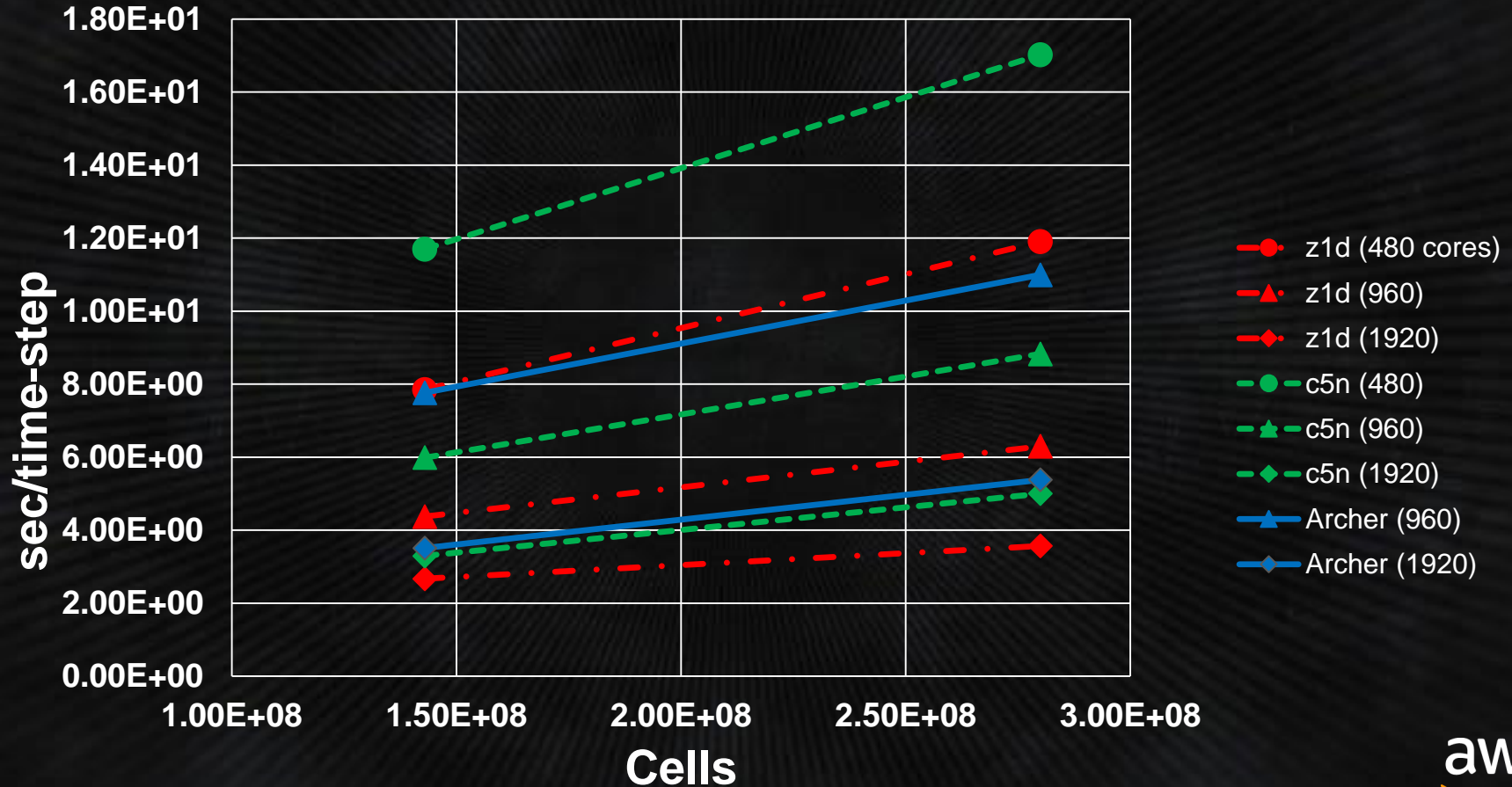
For example, scaling is measured at constant CFD mesh cells per core

$$\text{Ideal:} \lim_{n \to \infty} speedup = 1 - p$$

aws

Gustafson Scaling: openFoam (pimpleFoam)

# Gustafson Scaling: Alternate View

Legend:
- z1d (480 cores)
- z1d (960)
- z1d (1920)
- c5n (480)
- c5n (960)
- c5n (1920)
- Archer (960)
- Archer (1920)

X-axis: Cells
Y-axis: sec/time-step

Gustafson: Speedup ∝ p/n

# Gustafson: workload

Legend:
- z1d (medium) — red circle
- z1d (fine) — red square
- c5n (medium) — green circle
- c5n (fine) — green square
- Archer (medium) — blue circle
- Archer (fine) — blue square

Y-axis: cells/sec/core (workload/sec)
X-axis: Cells/Cores (workload)

aws

# 400 Cell scaling

# Obvious Question:

Why not include 400M cell case with others?

Answer:
They were done by two different people and two different workloads and, therefore, do not follow a consistent process

aws

# Latency

- Amdahl and Gustafson scaling both measure latency.  Latency includes:
  - Non parallel code
  - Memory latency
  - Cache latency
  - Network latency
  - Storage latency

aws

# Comments

- At high cell/core, or low core count, memory or cache more likely to slow processing (possibly see this in Archer)

- At low cell/core, or high core count, network latency more likely to slow processing (see in AWS-z1d and AWS-c5n)

- Total execution scaling based on processor workload (cells/core) is linear in all cases

- Slope of processor workload (cells/core) seems to follow network latency – faster network → steeper slope = better scaling

- Faster network does not mean faster execution time

- Coarser mesh has better scaling than finer mesh in all cases, but AWS-c5n has largest difference (best guess: PCG routines)
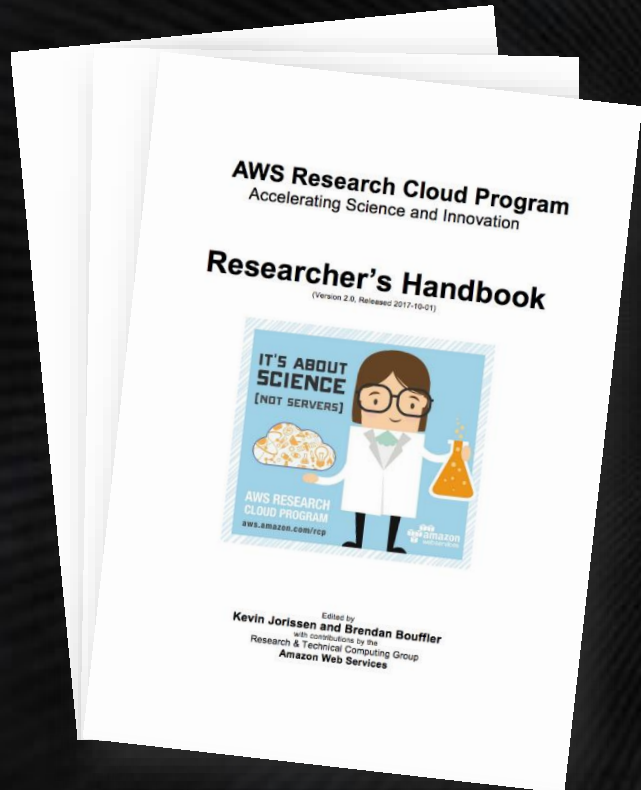
aws

Discussion?

aws

# Acknowledgement

Special thanks to Dr. Neil Ashton of Oxford University for providing his openFoam benchmarking data for the 143M cell and 282M cell cases, and Stephen Sachs, of AWS for the 400M cell cases

**aws**

# AWS Researcher's Handbook

The 200-page "missing manual" for science in the cloud.

Written by Amazon's Research Computing community for scientists.

- Explains foundational concepts about how AWS can accelerate time-to-science in the cloud.

- Step-by-step best practices for securing your environment to ensure your research data is safe and your privacy is protected.

- Tools for budget management that will help you control your spending and limit costs (and preventing any over-runs).

- Catalogue of scientific solutions from partners chosen for their outstanding work with scientists.

aws.amazon.com/rcp

AWS Research Cloud Program
Accelerating Science and Innovation

**Researcher's Handbook**
(Version 2.0, Released 2017-10-01)

IT'S ABOUT SCIENCE [NOT SERVERS]

AWS RESEARCH CLOUD PROGRAM
aws.amazon.com/rcp

Edited by
Kevin Jorissen and Brendan Bouffler
with contributions by the
Research & Technical Computing Group
Amazon Web Services

aws

aws.amazon.com/rcp