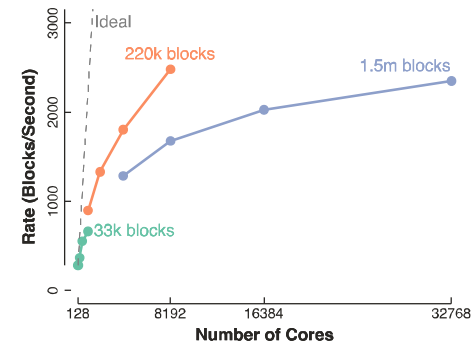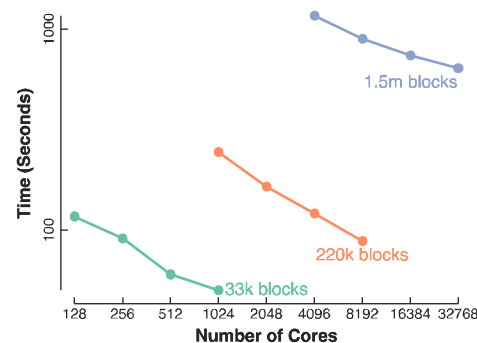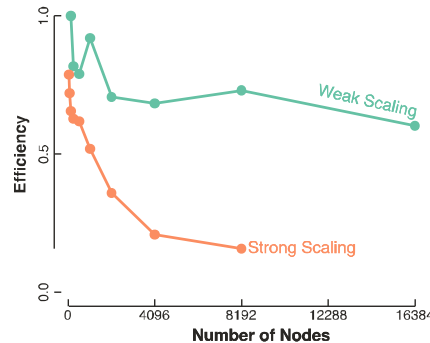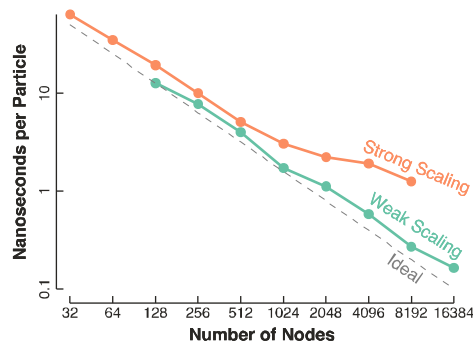*Exceptional service in the national interest*



# Measuring Parallel Software Scalability: You're Doing it Wrong

## SNL Computational Science Seminar Series

October 13, 2015

## Kenneth Moreland

Sandia National Laboratories

# Parallel Algorithm Speedup

$$S(n, p) = \frac{T^*(n)}{T(n, p)}$$

# Parallel Algorithm Speedup

$$S(n,p) = \frac{T^*(n)}{T(n,p)}$$

Serial time for large problem sizes

Cannot be measured in practice

⚠️

# Efficiency

$$E(n, p) = \frac{S(n, p)}{p}$$

# Efficiency

$$E(n, p) = \frac{S(n, p)}{p} = \frac{T^*(n)}{p\, T(n, p)}$$

# Karp-Flatt Metric

$$e(n, p) = \frac{\frac{1}{S(n,p)} - \frac{1}{p}}{1 - \frac{1}{p}}$$
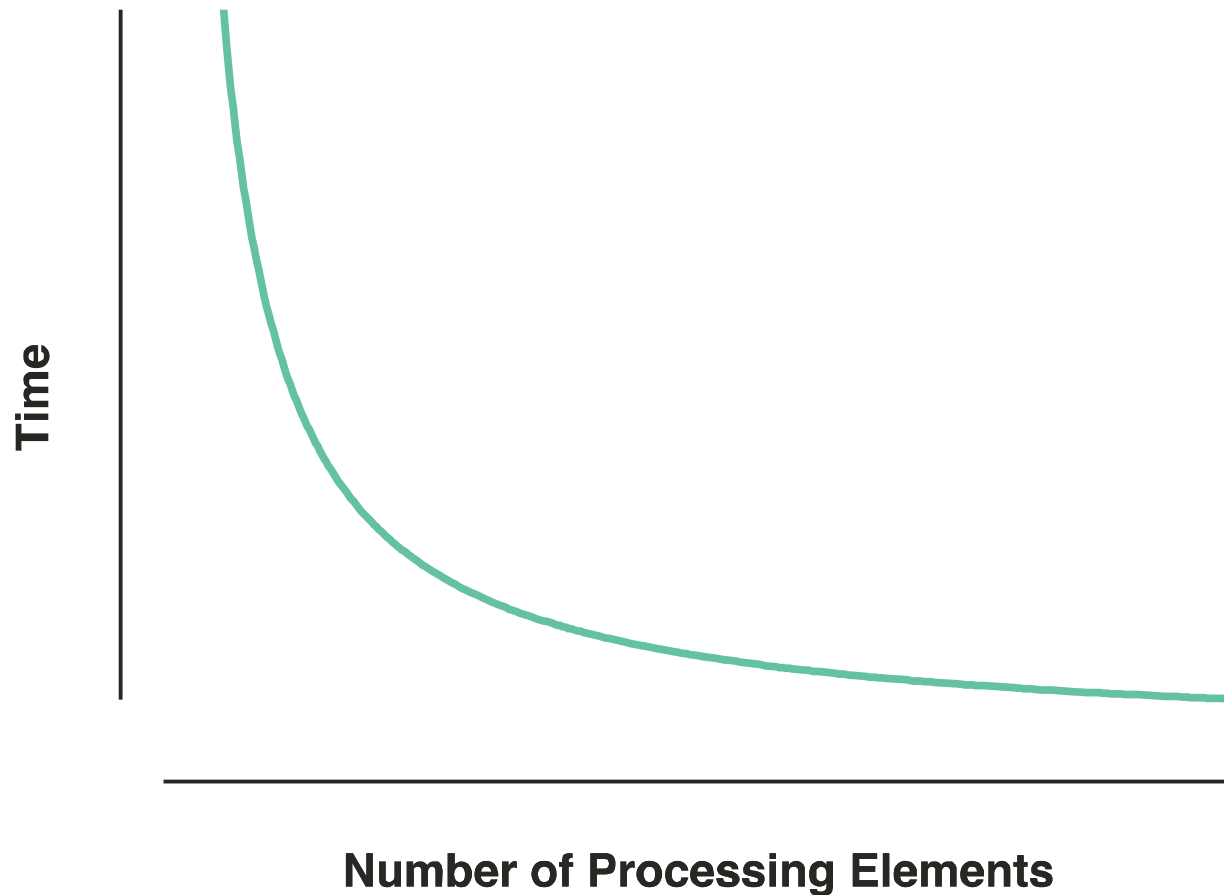
# Isoefficiency Metric
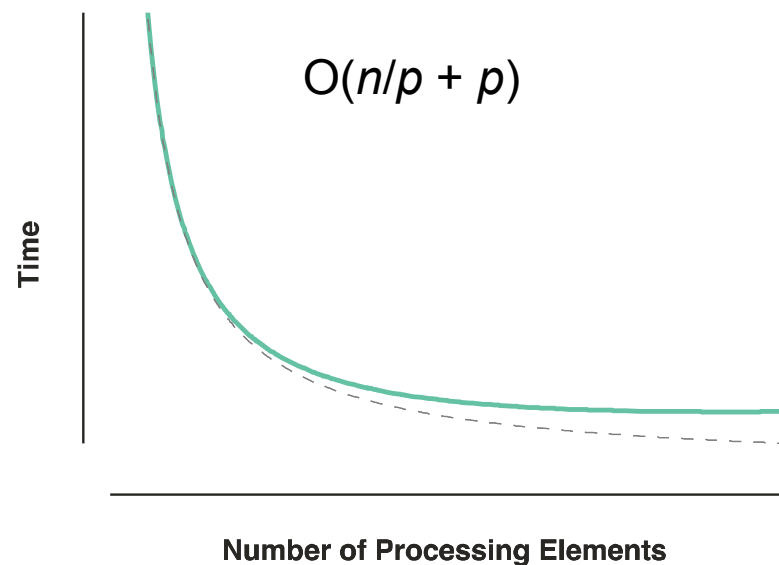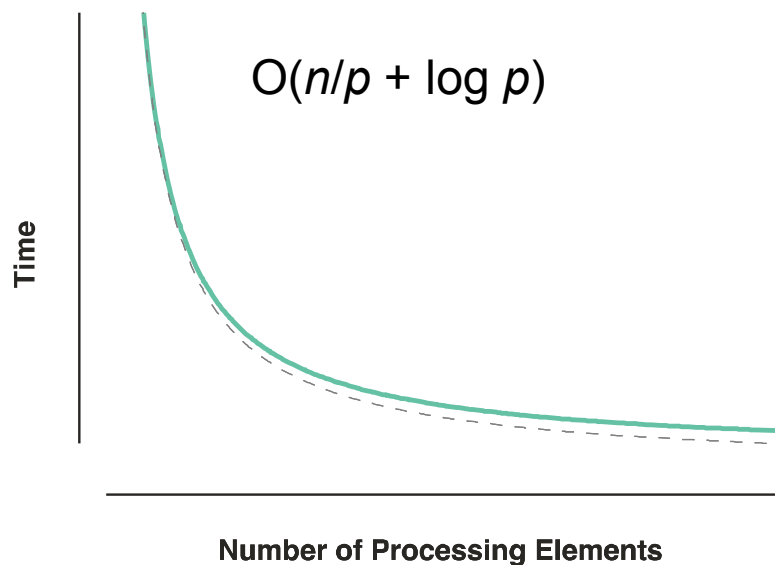
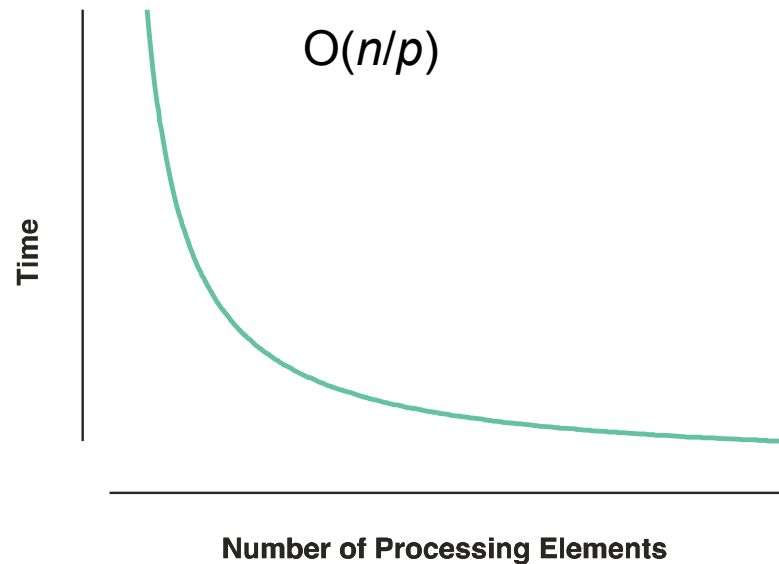$$T(n, 1) \geq \frac{E_d}{1 - E_d} T_o(n, p)$$

# Measuring Scalability in Practice

- **Strong Scaling:** Behavior as processing elements are increased and problem size held constant.
  - Per Amdahl's Law, strong scaling always has its limits.

- **Weak Scaling:** Behavior as processing elements and job size are increased proportionally.
  - Per Gustafson-Barsis Law, weak scaling can possibly be increased indefinitely.

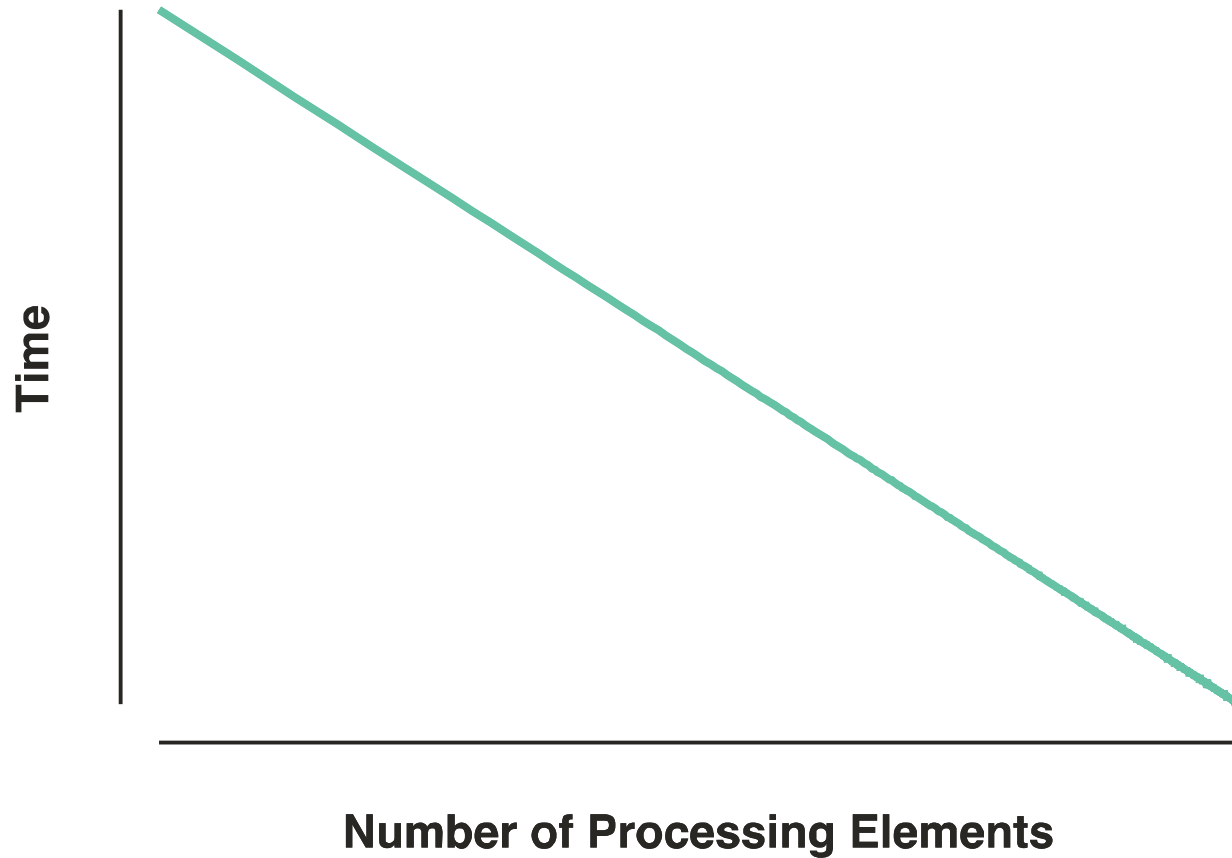- Scaling is often demonstrated with absolute run time over different scales.
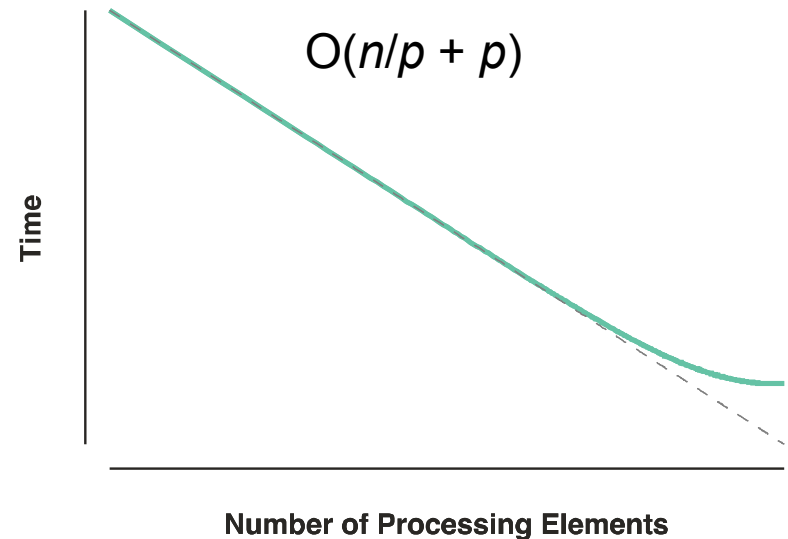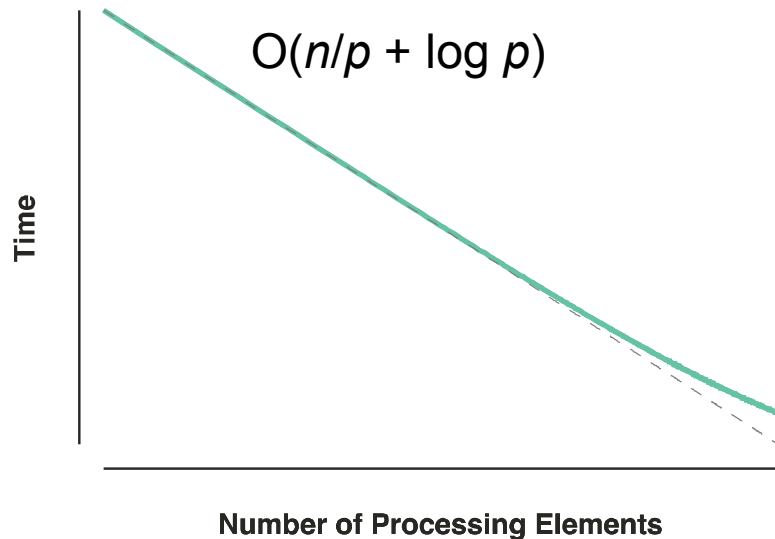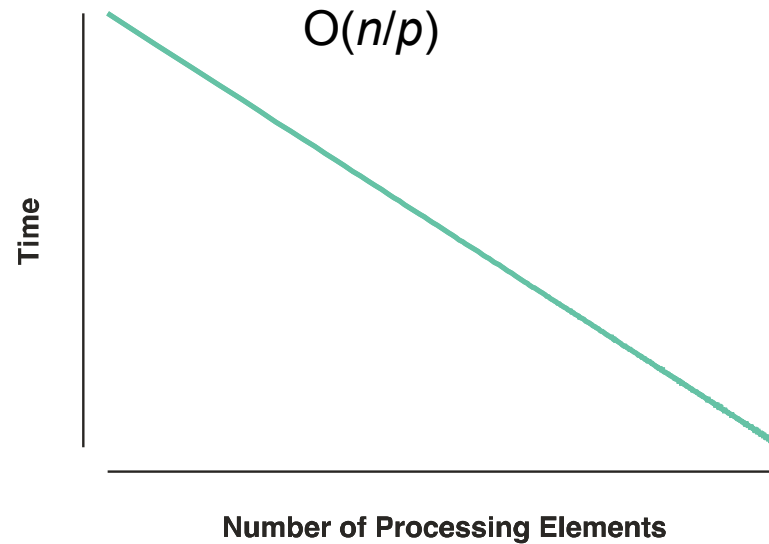
# Demonstrating Strong Scaling

# Measuring Strong Scaling



$O(n/p)$

Time

Number of Processing Elements

$O(n/p + \log p)$

Time

Number of Processing Elements

$O(n/p + p)$

Time

Number of Processing Elements

# Strong Scaling with Log Axes

# Measuring Strong Scaling with Log

O($n/p$)

Time

Number of Processing Elements

O($n/p$ + log $p$)

Time

Number of Processing Elements

O($n/p$ + $p$)

Time

Number of Processing Elements

# Demonstrating Weak Scaling

# Measuring Weak Scaling

O(*n*/*p*)

Time

Number of Processing Elements

O(*n*/*p* + log *p*)

Time

Number of Processing Elements

O(*n*/*p* + *p*)

Time

Number of Processing Elements

# Normalized Weak Scaling

# Scaling with More Visual Precision

- Our position statement: rate and efficiency better represent scaling behavior.

- Although neither rate nor efficiency is a new concept, there is not a lot of consistency in the community.

- Through algebra and examples I will show why rate and efficiency are the "right" metrics to use.

# Rate

$$R(n, p) = \frac{n}{T(n, p)}$$

# Why Use Rate?

$$S(n,p) = \frac{T^*(n)}{T(n,p)}$$

# Why Use Rate?

$$S(n,p) = \frac{T^*(n)}{T(n,p)} = \frac{T^*(n)}{n}R(n,p)$$

# Why Use Rate?

$$S(n,p) = \frac{T^*(n)}{T(n,p)} = \underbrace{\frac{T^*(n)}{n}} R(n,p)$$

Becomes a constant
with $n$ is constant.

# Why Use Rate?

$$S_n(p) \propto R_n(p)$$

# Scaling with Rate

# Measuring Scaling with Rate



O($n/p$)

Rate

Number of Processing Elements

O($n/p + \log p$)

Rate

Number of Processing Elements

O($n/p + p$)

Rate

Number of Processing Elements

# Measuring Scaling with Rate



O($n/p$ + log $p$)

Time — Number of Processing Elements

O($n/p$ + $p$)

Time — Number of Processing Elements

O($n/p$ + log $p$)

Rate — Number of Processing Elements

O($n/p$ + $p$)

Rate — Number of Processing Elements

# Efficiency

$$E(n, p) = \frac{S(n, p)}{p} = \frac{T^*(n)}{p\, T(n, p)}$$

# Measuring Efficiency from Cost

$$C(n,p) = p\,T(n,p)$$

# Measuring Efficiency from Cost

$$E(n, p) = \frac{C^*(n)}{C(n, p)}$$

Minimum (best) cost

# Measuring Efficiency from Cost

$$E(n, p) = \frac{C^*(n)}{C(n, p)}$$

Minimum (best) cost

if $C^*(n) = T^*(n)$

$$E(n, p) = \frac{T^*(n)}{C(n, p)} = \frac{T^*(n)}{p\, T(n, p)}$$

# Scaling with Efficiency

# Measuring Scaling with Efficiency

[This Slide Left Intentionally Blank]

# Unifying Strong and Weak Scaling

$$C_u(n, p) = \frac{C(n, p)}{n} = \frac{p\, T(n, p)}{n}$$

# Unifying Strong and Weak Scaling

$$C_{u,\text{ideal}}(n, p) = \frac{p\, T(n, p)}{n}$$

$$= \frac{p\, O(n/p)}{n}$$

$$= O(1)$$

# Unifying Strong and Weak Scaling

$$E(n, p) = \frac{C_u^*}{C_u(n, p)}$$

# Unifying Rate Across Data Scales

$$E(n, p) = \frac{C_u^*}{C_u(n, p)}$$

# Unifying Rate Across Data Scales

$$E(n, p) = \frac{C_u^*}{C_u(n, p)}$$

$$R(n, p) = \frac{p \, E(n, p)}{C_u^*}$$

# Unifying Rate Across Data Scales

$$E(n, p) = \frac{C_u^*}{C_u(n, p)}$$

$$R(n, p) = \frac{p \, E(n, p)}{C_u^*} \qquad E_{\text{ideal}}(n, p) = 1$$

$$R_{\text{ideal}}(n, p) = \frac{p}{C_u^*}$$

# Rate Across Data Scales

# Use Case 1: Gordon Bell Finalist

- Measurements of HACC code performance

- Excellent Scalability

- Measurements across many scales

- Lots of data provided in paper



HACC: Extreme Scaling and Performance Across Diverse Architectures

# Use Case 1: Gordon Bell Finalist

Use Case 1: Gordon Bell Finalist

# Use Case 1: Gordon Bell Finalist

# Use Case 1: Gordon Bell Finalist

# Use Case 1: Gordon Bell Finalist

# Use Case 2: Imperfect Scaling

- Measures visualization algorithm

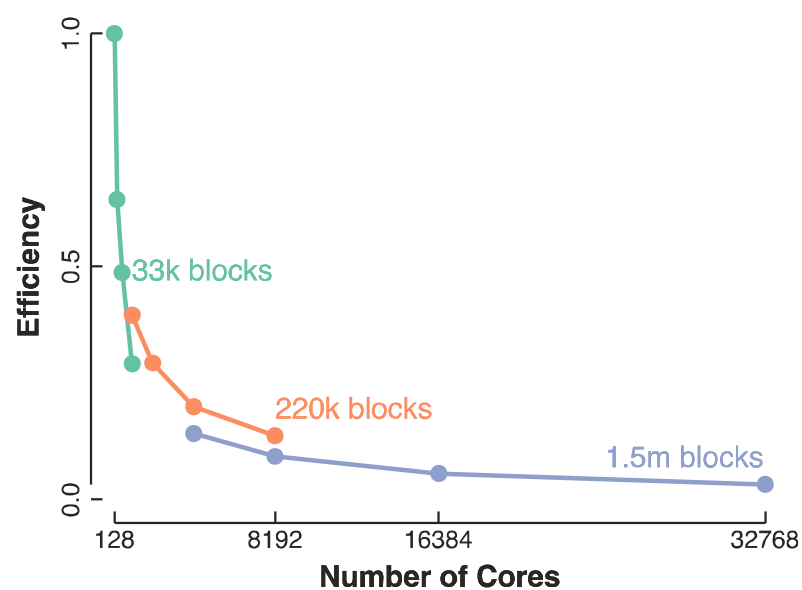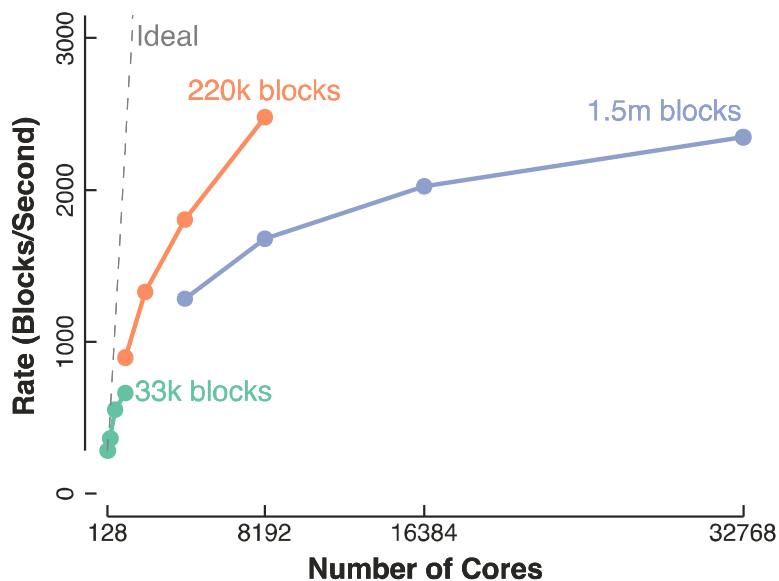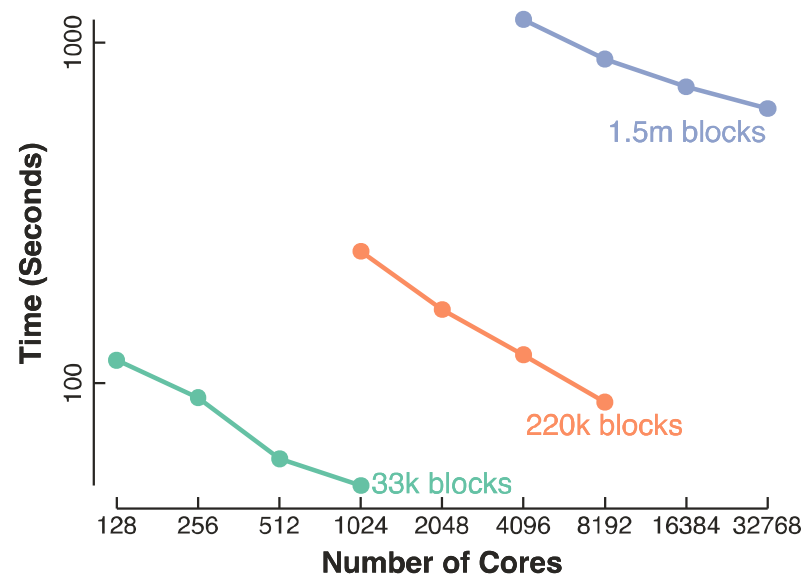- A high communication overhead severely limits scalability
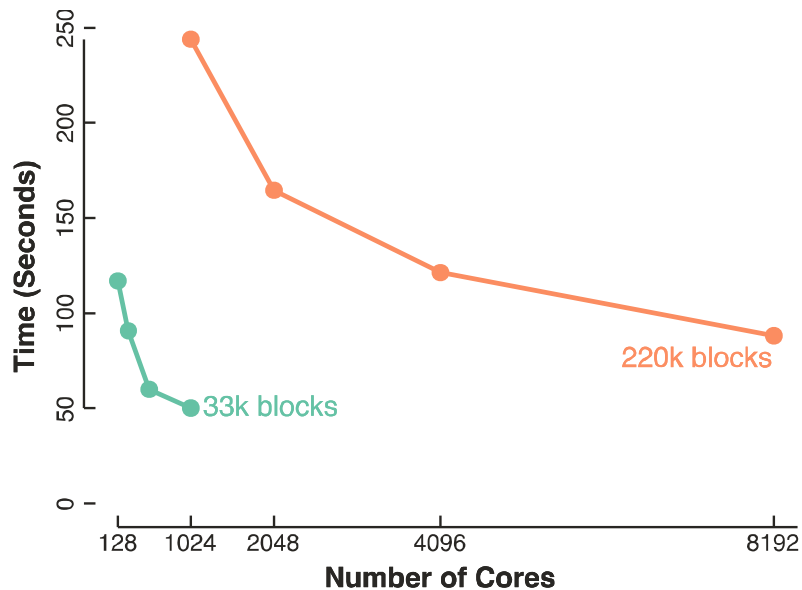
# Use Case 2: Imperfect Scaling

# Use Case 2: Imperfect Scaling

# Use Case 2: Imperfect Scaling

# Final Recommendations

- **Do not rely on running time** for performance analysis. Instead use rate, efficiency, or both.

- **Avoid using log-log scaling** on plot axes, which hides major inefficiencies. If necessary, repeat linear plots at different scales.

- Rather than performing them separately, **incorporate weak and strong scaling studies in one**. Perform several strong scaling studies at different scales of data size. Then find an overall minimal practical cost per unit and plot all the measurements together as demonstrated in the figures in this paper.

# Acknowledgements

- This material is based in part upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Scientific Discovery through Advanced Computing (SciDAC) program under Award Number 12-015215.

- This material is based in part upon work supported by the U.S. Department of Energy, National Nuclear Security Administration, Advanced Simulation and Computing (ASC).

# More Resources