# Parallel Programming and Heterogeneous Computing

A3 - Performance Metrics

Max Plauth, Sven Köhler, Felix Eberhardt, Lukas Wenzel and Andreas Polze

Operating Systems and Middleware Group

# Performance



- Which car is faster?
  - … for transporting several large boxes
  - … for winning a race





*Performance depends not only on an execution environment but also on the workload it executes!*

# Recap
# Optimization Goals

- Decrease **Latency** – process a single workload faster (= **speedup**)
- Increase **Throughput** – process more workloads in the same time
- ➤ Both are **Performance** metrics

- **Scalability**: make best use of additional resources
    - **Scale Up**: Utilize additional resources on a machine
    - **Scale Out**: Utilize resources on additional machines

- **Cost/Energy Efficiency**:
    - minimize cost/energy requirements for given performance objectives
    - *alternatively: maximize performance for given cost/energy budget*

- **Utilization**: minimize idle time (=waste) of available resources

- **Precision-Tradeoffs**: trade performance for precision of results

# Scaling Behavior

Different responses of performance metrics to scaling (additional resources):

- **Speedup**:
  **More resources ~ less time** executing the **same workload**

  › strong scaling

- **Scaled Speedup**:
  **More resources ~ same time** executing a **larger workload**

  › weak scaling

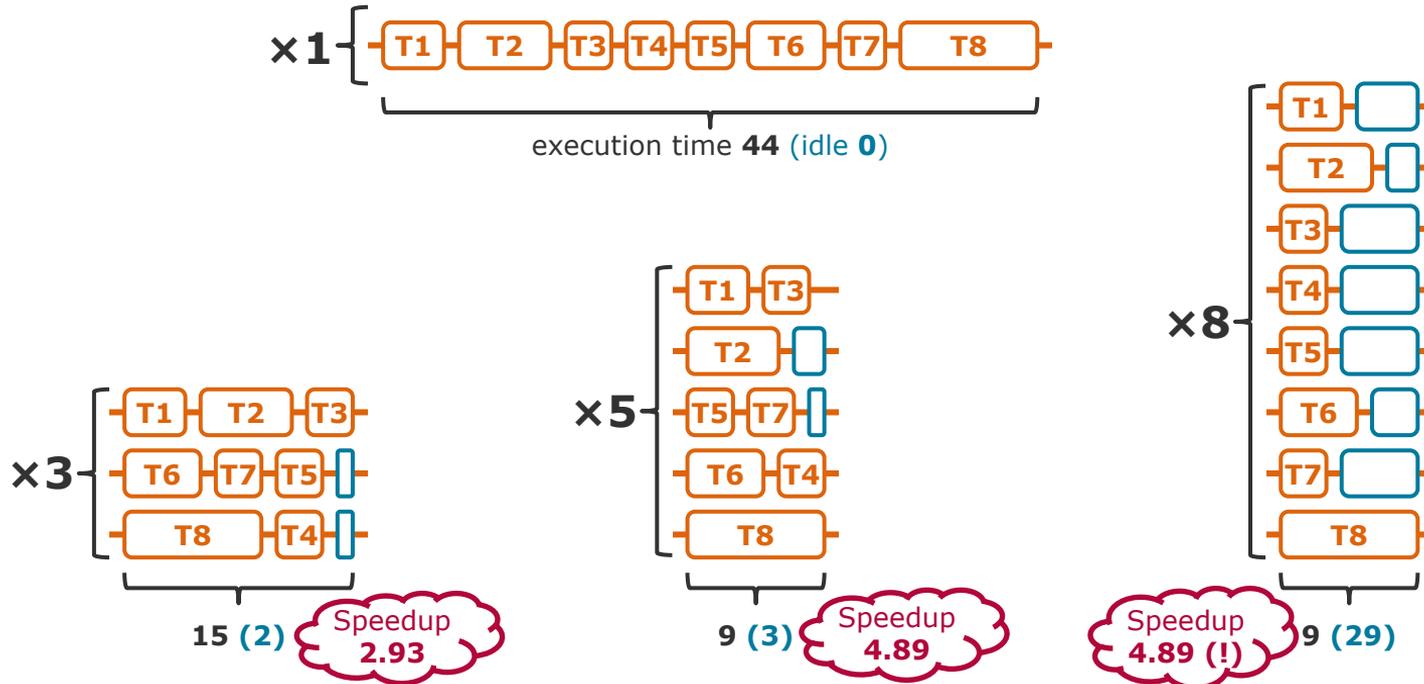- Linear speedup = resources and workload execution scale by same factor

**ParProg 2020 A3
Performance
Metrics**

Lukas Wenzel

Chart **4**

*A workload consists of multiple tasks,*
*containing different amounts of operations each.*

×1

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 |

execution time **44** (idle **0**)

×3

| T1 | T2 | T3 |
| T6 | T7 | T5 |
| T8 | T4 |

**15 (2)**

Speedup
**2.93**

×5

| T1 | T3 |
| T2 |
| T5 | T7 |
| T6 | T4 |
| T8 |

**9 (3)**

Speedup
**4.89**

×8

| T1 |
| T2 |
| T3 |
| T4 |
| T5 |
| T6 |
| T7 |
| T8 |

Speedup
**4.89 (!)**
**9 (29)**

**ParProg 2020 A3
Performance
Metrics**

Lukas Wenzel

Chart **5**

# Anatomy of a Workload

The **longest task** puts a **lower bound on** the shortest **execution time.**

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 |

Modeling discrete tasks is impractical → simplified **continuous model.**

$$T(N) = \frac{T_{par}}{N} + T_{seq}$$

Replace absolute times by **parallelizable fraction** $P$:

$$T_{par} = T_1 \cdot P$$
$$T_{seq} = T_1 \cdot (1 - P)$$

$$T(N) = T_1 \cdot \left( \frac{P}{N} + (1 - P) \right)$$

Chart **6**

# [Amdahl1967]
# Amdahl's Law

Amdahl's Law derives the speedup $s_{Amdahl}(N)$ for a parallelization degree $N$

$$s_{Amdahl}(N) = \frac{T_1}{T(N)} = \frac{T_1}{T_1 \cdot \left(\frac{P}{N} + (1 - P)\right)} = \frac{1}{\frac{P}{N} + (1 - P)}$$

Even for **arbitrarily large** $N$, the speedup converges to a **fixed limit**
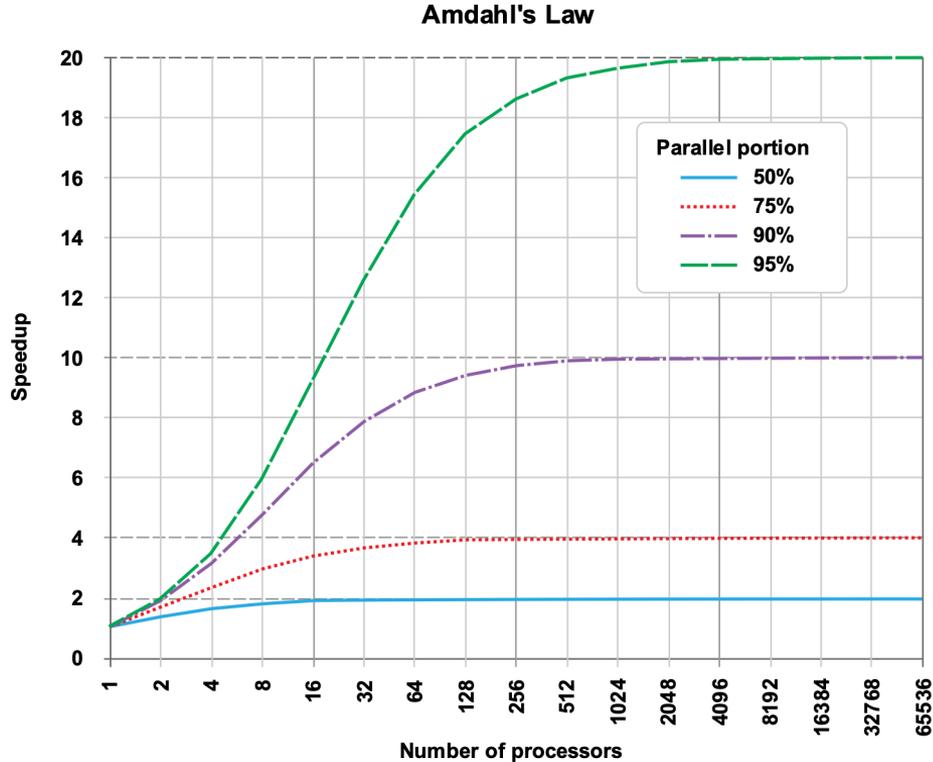
$$\lim_{N \to \infty} s_{Amdahl}(N) = \frac{1}{1 - P}$$

*For getting reasonable speedup out of 1000 processors, the sequential part must be substantially below 0.1%*

# [Amdahl1967]
# Amdahl's Law



Amdahl's Law

By Daniels220 at English Wikipedia, CC BY-SA 3.0,
https://commons.wikimedia.org/w/index.php?curid=6678551

# [Amdahl1967]
# Amdahl's Law

*Regardless of processor count, **90% parallelizable** code allows not more than a **speedup by factor 10**.*

➢ Parallelism requires highly parallelizable workloads to achieve a speedup
▪ What is the sense in large parallel machines?

Amdahl's law assumes a simple speedup scenario!
➢ isolated execution of a **single workload**
➢ **fixed workload size**

# [Gustafson1988]
# Gustafson-Barsis' Law

Consider a **scaled speedup scenario**, allowing a variable workload size **w**.

Amdahl ~ *What is the shortest execution time for a given workload?*

Gustafson-Barsis ~ *What is the largest workload for a given execution time?*

$$w_1 \sim T_{par} + T_{seq}$$

$$w(N) \sim N \cdot T_{par} + T_{seq}$$

**Assumption: The parallelizable part of a workload contributes useful work when replicated.**
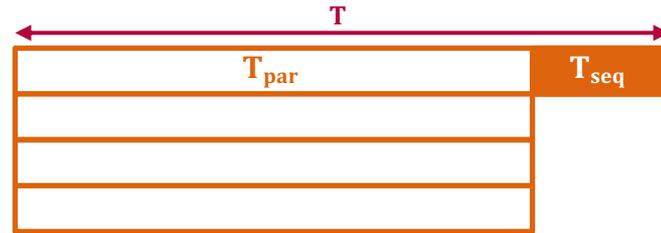
Chart **10**

# [Gustafson1988]
# Gustafson-Barsis' Law

Determine the scaled speedup $s_{Gustavson}(N)$ through
the increase in workload size $w(N)$ over the fixed execution time $T$

$$s_{Gustafson}(N) = \frac{w(N)}{w_1} = \frac{T \cdot (P \cdot N + (1-P))}{T \cdot (P + (1-P))} = P \cdot N + (1-P)$$



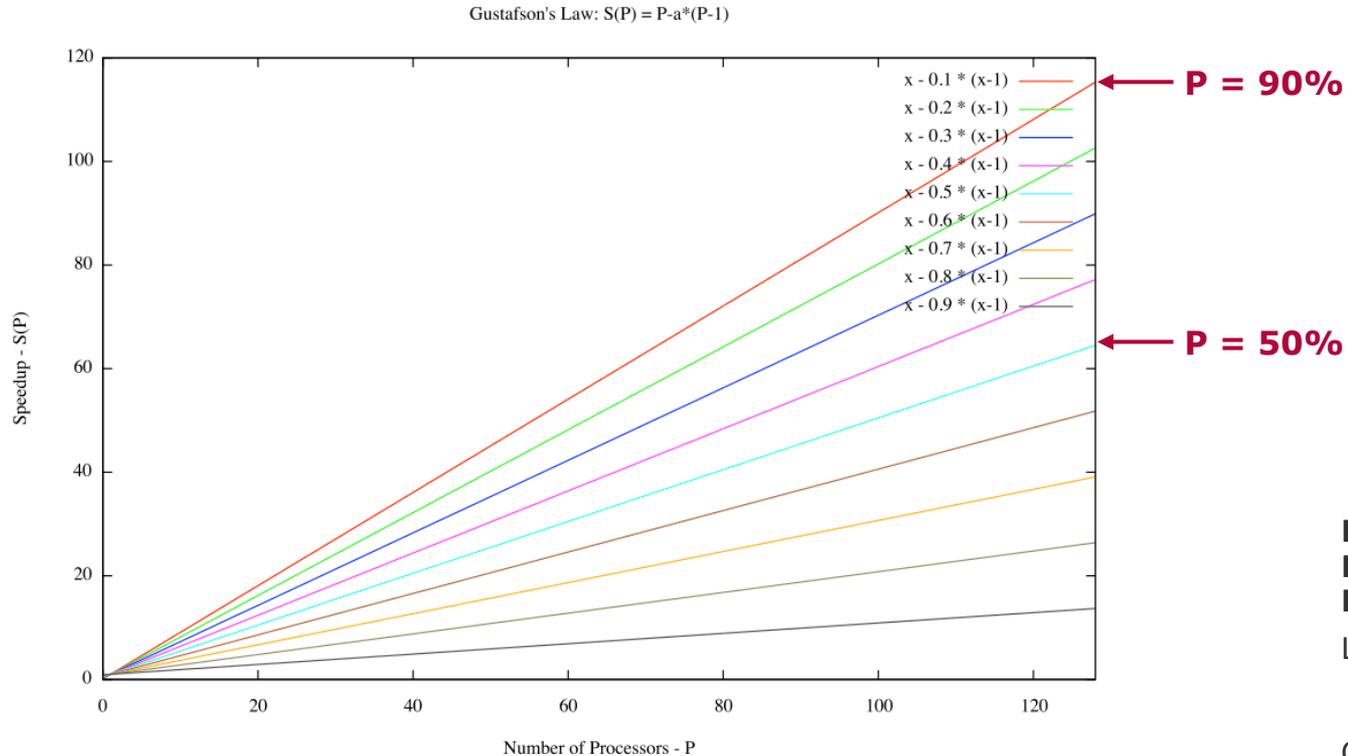$$w_1 \sim T_{par} + T_{seq}$$

$$w(N) \sim N \cdot T_{par} + T_{seq}$$

**Assumption: The parallelizable part of a workload contributes useful work when replicated.**

**ParProg 2020 A3 Performance Metrics**

Lukas Wenzel

Chart **11**

# [Gustafson1988]
# Gustafson-Barsis' Law



Gustafson's Law: S(P) = P-a*(P-1)

P = 90%

P = 50%

x - 0.1 * (x-1)
x - 0.2 * (x-1)
x - 0.3 * (x-1)
x - 0.4 * (x-1)
x - 0.5 * (x-1)
x - 0.6 * (x-1)
x - 0.7 * (x-1)
x - 0.8 * (x-1)
x - 0.9 * (x-1)

Speedup - S(P)

Number of Processors - P

**ParProg 2020 A3
Performance
Metrics**

Lukas Wenzel

Chart **12**

# [Karp1990]
# Karp-Flatt-Metric

Parallel fraction $P$ is a hypothetical parameter and not easily deduced from a given workload.

➢ Karp-Flatt-Metric determines sequential fraction $Q = 1 - P$ empirically

1. Measure baseline execution time $T_1$
   by executing workload on a single execution unit

2. Measure parallelized execution time $T(N)$
   by executing workload on $N$ execution units

3. Determine speedup $s(N) = {T_1}/{T(N)}$

4. Calculate Karp-Flatt-Metric

$$Q(N) = \frac{\frac{1}{s(N)} - \frac{1}{N}}{1 - \frac{1}{N}}$$

**ParProg 2020 A3**
**Performance**
**Metrics**

Lukas Wenzel

Chart **13**

**The Karp-Flatt-Metric is derived by rearranging Amdahl's Law.**

$$\frac{1}{s(N)} = \frac{T(N)}{T_1} \quad ; \quad T(N) = \left(\frac{1-Q}{N} + Q\right) \cdot T_1$$

$$\frac{1}{s(N)} = \frac{\left(\frac{1-Q}{N} + Q\right) \cdot T_1}{T_1}$$

$$\frac{1}{s(N)} = \frac{1-Q}{N} + Q = \frac{1}{N} + \left(1 - \frac{1}{N}\right) \cdot Q$$

$$\frac{1}{s(N)} - \frac{1}{N} = \left(1 - \frac{1}{N}\right) \cdot Q$$

$$\frac{\frac{1}{s(N)} - \frac{1}{N}}{1 - \frac{1}{N}} = Q$$

**ParProg 2020 A3
Performance
Metrics**

Lukas Wenzel

Chart **14**

# [Karp1990]
# Karp-Flatt-Metric

Observing $Q(N)$ for different $N$ gives an indication,
how the workload reacts to different degrees of parallelism:

- $Q(N)$ close to $0$ ~ *high parallel fraction, workload benefits from parallelization*

- $Q(N)$ close to $1$ ~ *low parallel fraction, workload can not use parallel resources*

- $Q(N)$ increases with $N$ ~ *workload suffers from parallelization overhead*

- $Q(N)$ decreases with $N$ ~ *workload scales well*

Observing $Q(N)$ for different implementation variants of the workload can reveal bottlenecks.

# [Leiserson2008]
# A More Detailed View

Directed Acyclic Graph to model a workload:

- Nodes represent operations
- Edges express dependencies between operations

Work **T** - Total workload execution time

$T_1$ - Execution time with a single processor

*~ number of nodes*

$T_P$ - Execution time with P processors

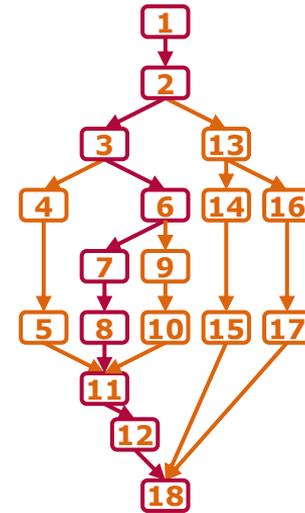$T_\infty$ - Execution time with arbitrary number of processors

*~ graph diameter*

**Work Law** $T_P \geq {}^{T_1}\!/_P$
(processors can not process multiple operations at once)

**Span Law** $T_P \geq T_\infty$
(execution order can not break dependencies)

$T_1 = 18$

$T_\infty = 9$

# Literature

**[Amdahl1967]**

Amdahl, Gene M. "Validity of the single processor approach to achieving large scale computing capabilities." *Proceedings of the AFIPS Spring Joint Computer Conference*. 483-485. 1967.

**[Gustafson1988]**

Gustafson, John L. "Reevaluating Amdahl's law." *Communications of the ACM* 31.5 (1988): 532-533.

**[Karp1990]**

Karp, Alan H. and Flatt, Horace P. "Measuring parallel processor performance." *Communications of the ACM* 33.5 (1990): 539-543.

**[Leiserson2008]**

Leiserson, Charles E. and Mirman, Ilya B. "How to survive the multicore software revolution (or at least survive the hype)." *Cilk Arts* 1 (2008): 11.

**ParProg 2020 A3 Performance Metrics**

Lukas Wenzel

Chart **17**

And now for a break and
a cup of Oolong.

*or beverage of your choice