

---

---

Analysis and modeling of  
**Computational Performance**

# Scalability

---

→ Standard parallel performance measures

- speed-up
- efficiency

characterise the so called strong scalability of programs

→ Good strong scalability, being closed to the perfect, linear speed-up, is difficult to obtain, e.g. due to:

- sequential (not possible to be parallelized) parts of the program/algorithm (as in Amdahl law)
- communication

→ It is possible to obtain good parallel performance also for programs with sequential parts and communication

- for many algorithms/programs the speed-up curves become closer to the perfect speed-up with the increasing problem size
  - the most popular way for expressing problem size is the number of (dominant) operations – the operations for which the execution time is proportional to the number of operations



# Scalability

---

- Strong scalability measures the performance for the increasing number of resources used in computations (processors, computational nodes)
- The standard notions of scalability include also the condition of increasing workload
- The increasing workload in case of programs may be the number of dominant operations
  - sometimes the execution time for single processor/core is taken as the workload
  - in the simplified analysis the time per single dominant operation is constant – so the number of dominant operations is proportional to the execution time for single processor/core
- Weak scalability denotes the scalability for the case where the total workload is not constant (as for the strong scalability), but the workload per single processor/core (thread/process) is constant
  - constant workload per thread/process means the total workload proportional to the number of threads/processes

# Weak scalability

---

- For weak scalability study, one can define the execution time as the function of the number of threads/processes and the workload
  - $T_{\parallel} = T_{\parallel}(p, W) = (\text{for weak scalability}) T_{\parallel}(p, pW_0)$
- Then it is possible to define the scaled speed-up,  $S^S(p)$ :
  - $S^S(p) = T_{\parallel}(1, pW_0) / T_{\parallel}(p, pW_0) = p * T_{\parallel}(1, W_0) / T_{\parallel}(p, pW_0)$
- A program/algorithm has linear weak scalability when its scaled speed-up is linear
- Linear weak scalability is equivalent to:
  - **the same execution time for p-times larger problems executed on p cores/processors**
  - the parallel overhead constant per single processor/core for p-times larger problems executed on p cores/processors

# Scalability of computations

---

- Scalability is the key property for obtaining high performance of computations
- Weak scalability (almost linear) can be attributed to many algorithms, while linear strong scalability is extremely rare
- The general relation:
  - $performance = number\_of\_operations / execution\_time$can be transformed to the expression:
  - $performance = speed-up / execution\_time\_per\_single\_operation$ 
    - good speed-up denotes effectively parallelized programs
    - short  $execution\_time\_per\_single\_operation$  characterises well optimized single thread computations
    - hence the total performance has two ingredients:
      - scalability
      - single thread performance

# Execution time modelling

---

- When modelling execution time several simplifying assumptions can be adopted::
- The single thread execution time and the computation time for a single thread in parallel programs are proportional to the number of dominant operations, which for the considered algorithms are the arithmetic operations
    - the notion of arithmetic intensity makes also the execution time of programs with memory limited performance proportional to the number of arithmetic operations
  - The time for performing a single dominant operation is constant and denoted by  $t_c$ 
    - $t_c$  is some amortized time per operation that includes memory accesses, sequential execution system overhead (e.g. memory allocation), etc.
  - Apart from  $t_c$  there are only two other hardware parameters, that characterise the communication time:  $t_s$  and  $t_w$
  - The parallelized computations are perfectly balanced

# Example

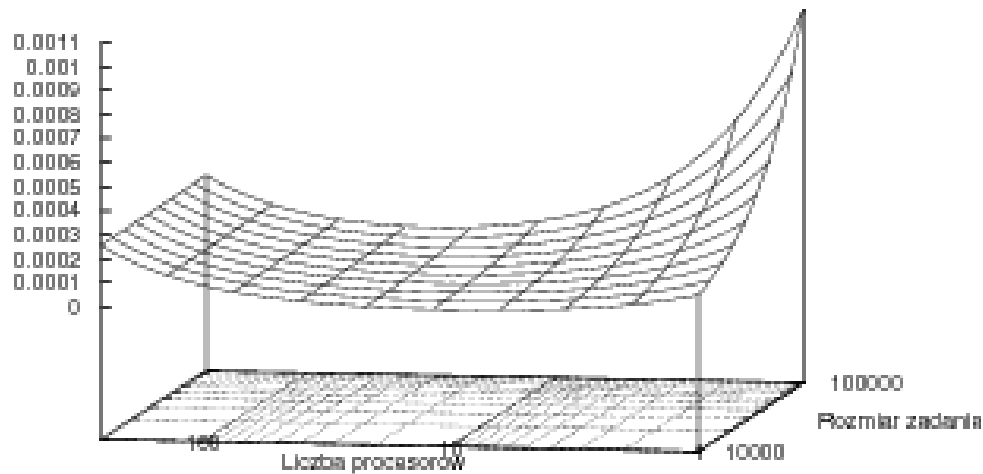
---

- Calculation of the norm of vector with size  $N$
- $N$  additions and multiplications – decomposition to obtain local sums – perfect speed-up possible
- Global sum reduction – communication time dominates arithmetic operations time
- Naive algorithm – all threads/processes send their local sums to the master thread/process
- Execution time modelling:  $T_{\parallel}(p) = 2*N*t_c/p + p*(t_s+8*t_w)$
- Workload (number of operations):  $W = 2*N$
- The same workload per thread/process for weak scalability study:
  - $W(p) = W_1 * p = 2 * N_1 * p$
  - $T_{\parallel}(p) = 2*N*t_c/p + p*(t_s+8*t_w) = 2 * N_1 * t_c + p * (t_s+8*t_w)$
- Simple analysis to obtain: speed-up, efficiency, scaled speed-up, scaled efficiency, isoefficiency function, memory size limited speed-up, etc.

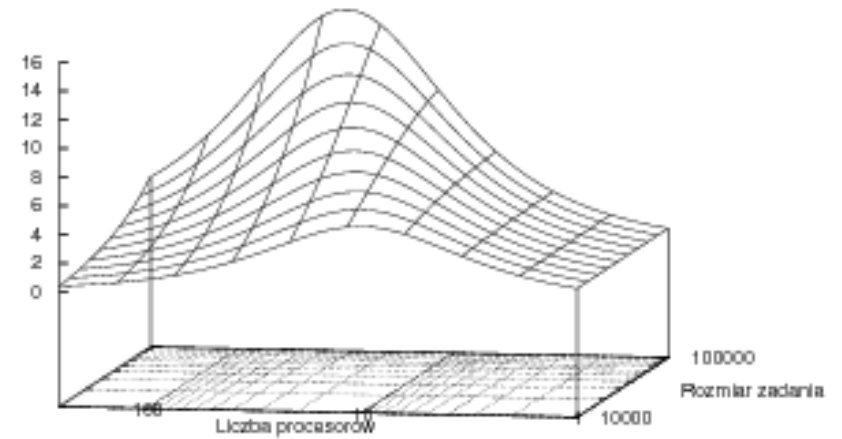


# Example

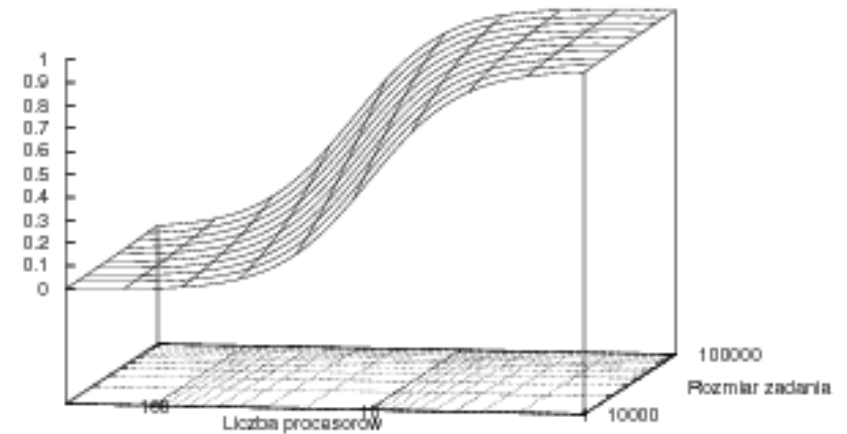
Czas wykonania równoległego



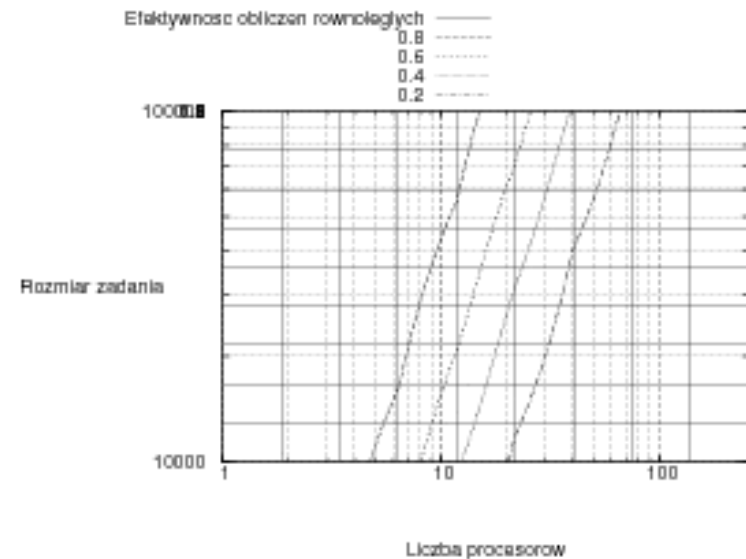
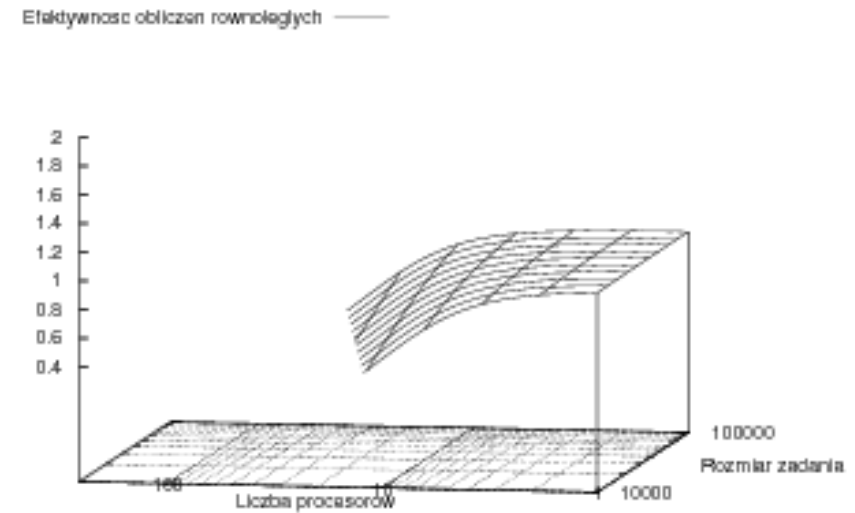
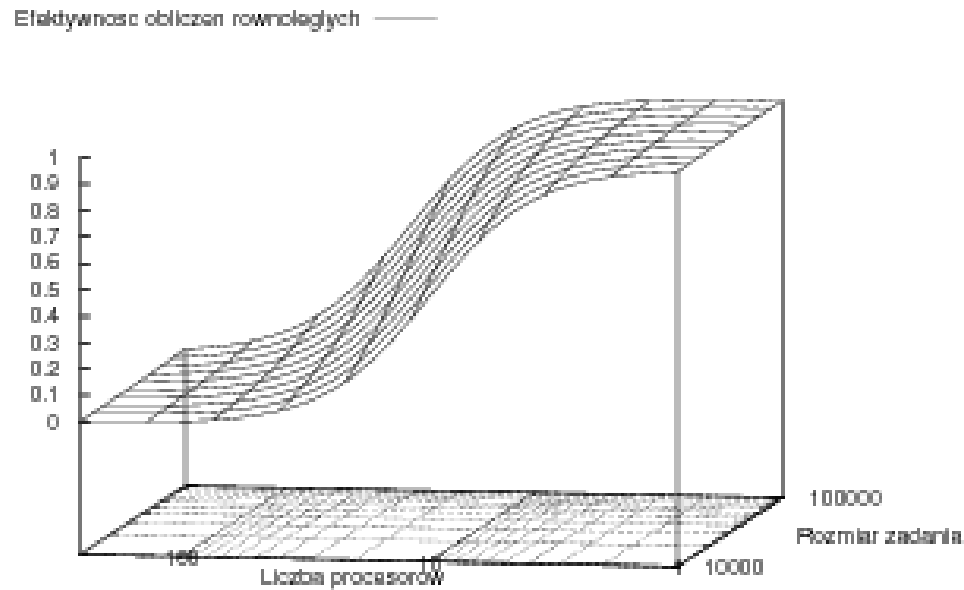
Przyspieszenie obliczeń równoległych



Efektywność obliczeń równoległych



# Example



# Optimization of parallel programs

---

- To minimize execution time for parallel programs, for distributed memory computer architectures, the following steps should be undertaken:
- load balancing
  - minimization of the total size of messages
  - minimization of the number of messages (by increasing the granularity of computations)
  - avoiding network contention
  - reducing the time for additional operations related to parallel computations (e.g. redundant computations – but redundant computations can decrease the communication volume)
  - reducing system overhead (e.g. for synchronization)
  - overlapping computations with communication
  - optimizing single thread execution time
    - including optimizing memory accesses