# Performance Prediction through Time Measurements

**Roman Iakymchuk**

AICES Graduate School, RWTH Aachen
iakymchuk@aices.rwth-aachen.de

International Conference on High Performance Computing
October 12-14, 2011
Kyiv, Ukraine

$$Performance = \frac{\#FLOPS}{Execution\_time}$$

*#FLOPS* is known **a priori**

**RWTH**AACHEN
UNIVERSITY

$$Performance = \frac{\#FLOPS}{Execution\_time}$$

*#FLOPS* is known **a priori**

## Modeling Performance

- Target—linear algebra algorithms

$$Performance = \frac{\#FLOPS}{Execution\_time}$$

*#FLOPS* is known **a priori**

## Modeling Performance

- Target—linear algebra algorithms

LAPACK   ScaLAPACK   FLAME   PLAPACK   ATLAS

↘   ↓   ↙

BLAS

- Model time of the kernel operations (BLAS) by conducting only **few** measurements

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**

# Timing Methodologies

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

# Timing Methodologies

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

Wall time

- High resolution (cycle-accurate)

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

Wall time

- High resolution (cycle-accurate)
- Includes other processes

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

Wall time

- High resolution (cycle-accurate)
- Includes other processes

## Improving the accuracy of timings

- Take multiple timing samples

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

Wall time

- High resolution (cycle-accurate)
- Includes other processes

## Improving the accuracy of timings

- Take multiple timing samples
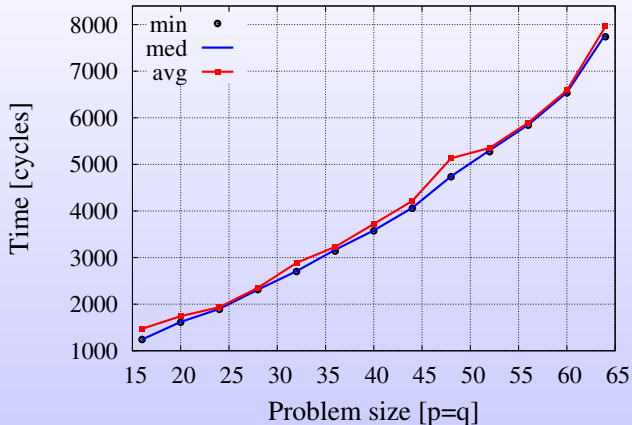- Select **median** timing for CPU time

## Choosing a system timer

CPU time

- if machine heavily loaded and no I/O and parallelism
- **Low resolution**
- Over- and under-reports time

Wall time

- High resolution (cycle-accurate)
- Includes other processes

## Improving the accuracy of timings

- Take multiple timing samples
- Select **median** timing for CPU time
- Select minimum for wall time

Figure: In-cache timing of GER

- GER:
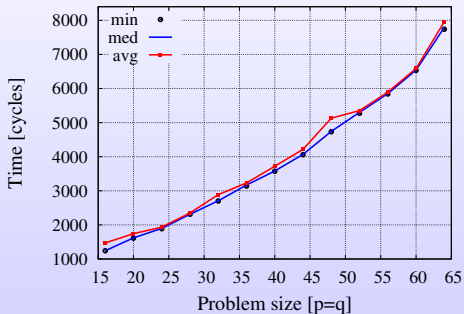  $A := A + \alpha x y^T$
- The cycle-accurate wall timer is used

# Time Measurements

Figure: In-cache



Figure: Out-of-cache

Source: R. Clint Whaley
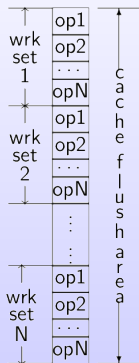(UTSA-CS)

- $size(\textit{flush\_area}) =$
  $\textit{Associativity} \times size(\textit{cache})$

Source: R. Clint Whaley
(UTSA-CS)

- $size(\textit{flush\_area}) =$
  $\textit{Associativity} \times size(\textit{cache})$

- Conduct `n_rep` timing samples of an algorithm

Source: R. Clint Whaley
(UTSA-CS)

- $size(\textit{flush\_area}) =$
  $\quad \textit{Associativity} \times size(\textit{cache})$

- Conduct `n_rep` timing samples of an algorithm

- On each iteration of `n_rep` loop **operands** are out-of-cache

## BLAS subroutines

- Apply polynomial interpolation

$$Execution\_time = a_k n^k + a_{k-1} n^{k-1} + ... + a_1 n + a_0$$

## BLAS subroutines

- Apply polynomial interpolation

  *Execution_time* $= a_k n^k + a_{k-1} n^{k-1} + ... + a_1 n + a_0$

- Complexity of a BLAS subroutine is at most $O(n^3) \rightarrow k \leq 3$

## BLAS subroutines

- Apply polynomial interpolation

  *Execution_time* $= a_k n^k + a_{k-1} n^{k-1} + ... + a_1 n + a_0$

- Complexity of a BLAS subroutine is at most $O(n^3) \rightarrow k \leq 3$
- Perform $4 - 6$ measurements on each memory level

## BLAS subroutines

- Apply polynomial interpolation

  $Execution\_time = a_k n^k + a_{k-1} n^{k-1} + ... + a_1 n + a_0$

- Complexity of a BLAS subroutine is at most $O(n^3) \rightarrow k \leq 3$

- Perform $4 - 6$ measurements on each memory level

- Solve a linear least squares problem

## BLAS subroutines

- Apply polynomial interpolation

  *Execution_time* $= a_k n^k + a_{k-1} n^{k-1} + ... + a_1 n + a_0$

- Complexity of a BLAS subroutine is at most $O(n^3) \rightarrow k \leq 3$
- Perform $4 - 6$ measurements on each memory level
- Solve a linear least squares problem

## Higher level algorithms

$$Execution\_time = \sum_{i=1}^{n-1} \texttt{Model\_subroutines\_time}(i)$$

Figure: $3 \times 3$ partitioning of $A$.

# A Case Study: LU Factorization

**Partition**
$$A \to \left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$$
    **where** $A_{TL}$ is $0 \times 0$

**While** $m(A_{TL}) < m(A)$ **do**
    **Repartition**
$$\left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \to \left( \begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$
        **where** $\alpha_{11}$ is $1 \times 1$

$$a_{21} := a_{21}/\alpha_{11} \qquad \text{SCAL}$$
$$A_{22} := A_{22} - a_{21}a_{12}^T \qquad \text{GER}$$

    **Continue with**
$$\left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left( \begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$$
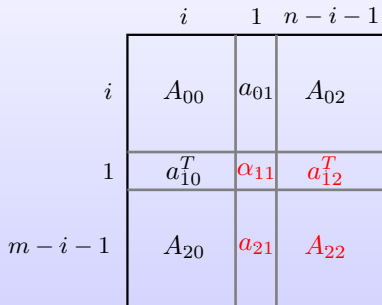
**endwhile**

**Partition**

$A \rightarrow \left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right)$
    **where** $A_{TL}$ is $0 \times 0$

**While** $m(A_{TL}) < m(A)$ **do**
    **Repartition**

$\left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left( \begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$
    **where** $\alpha_{11}$ is $1 \times 1$

$a_{21} := a_{21}/\alpha_{11}$       **SCAL**
$A_{22} := A_{22} - a_{21}a_{12}^T$    **GER**

**Continue with**

$\left( \begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left( \begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right)$
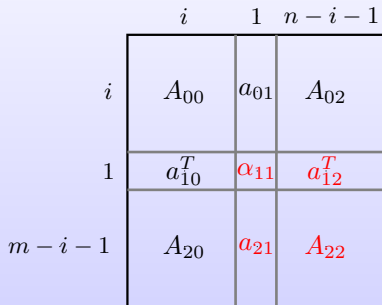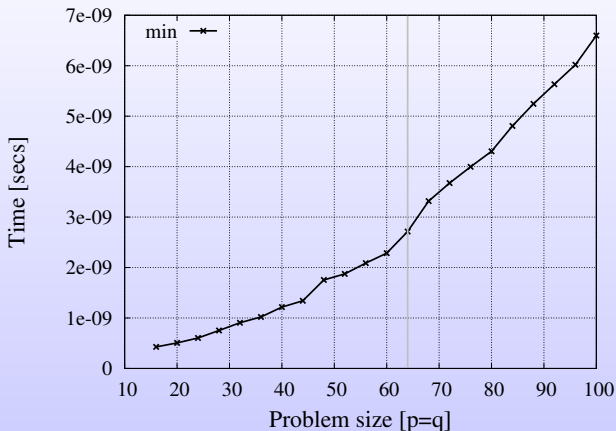
**endwhile**

| | $i$ | $1$ | $n-i-1$ |
|---|---|---|---|
| $i$ | $A_{00}$ | $a_{01}$ | $A_{02}$ |
| $1$ | $a_{10}^T$ | $\alpha_{11}$ | $a_{12}^T$ |
| $m-i-1$ | $A_{20}$ | $a_{21}$ | $A_{22}$ |

Figure: $3 \times 3$ partitioning of $A$.

GER performs more than 96 % of the *#FLOPS* in the LU

**RWTH**AACHEN
**UNIVERSITY**



- Intel Harpertown @3.0 GHz
- L1 (32 KB) and L2 (6 MB) caches
- Apply parabolic interpolation on L1 & L2 caches

Figure: Piecewise-parabolic behavior of GER

### GER

Complexity of GER is $O(n^2) \rightarrow$

$$Execution\_time = a_2 n^2 + a_1 n + a_0$$

## GER

Complexity of GER is $O(n^2) \rightarrow$

$$Execution\_time = a_2 n^2 + a_1 n + a_0$$

## The unblocked LU

- Prediction:

$$Execution\_time = \sum_{i=1}^{n-1} \texttt{Model\_GER\_time}(i)$$

## GER

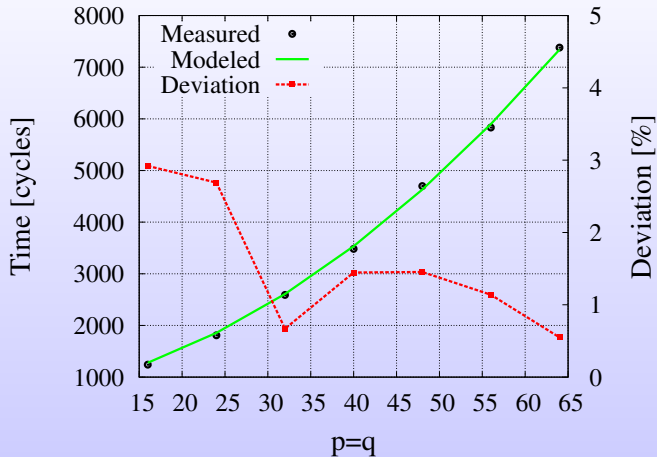Complexity of GER is $O(n^2) \rightarrow$

$$Execution\_time = a_2 n^2 + a_1 n + a_0$$

## The unblocked LU

- Prediction:

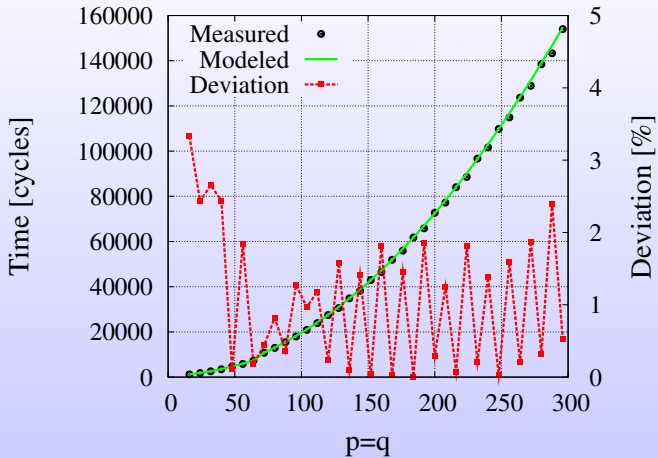$$Execution\_time = \sum_{i=1}^{n-1} \texttt{Model\_GER\_time}(i)$$

- In total, GER is measured only 8-12 times

Figure: Predicting the execution time of GER on Harpertown

- GER from the GotoBLAS library is used
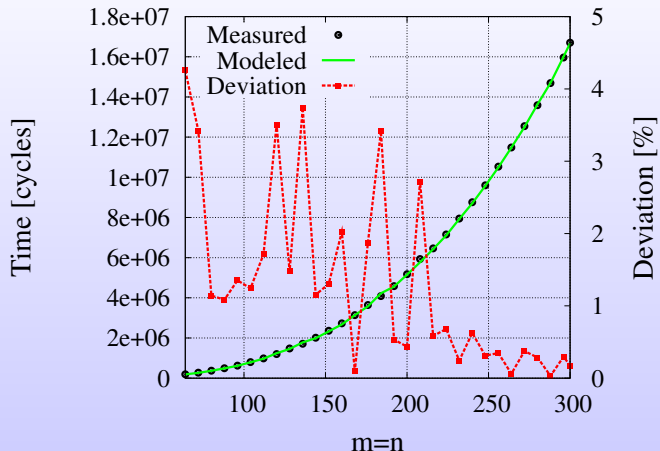- $p \leq 64$ fit in the L1 cache
- The deviation decreases; it is less than 3 %

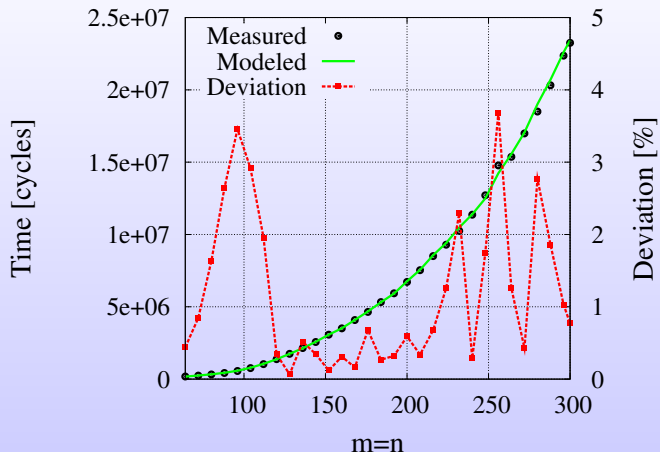Figure: Predicting the execution time of GER on Harpertown

- $p < 300$ fit in the L1 and L2

- The deviation is less than 2 %

**RWTH**AACHEN
UNIVERSITY



- Closer to origin the deviation is higher
- When $m = n$ increases the deviation $\to 0$

Figure: Modeling the execution time of the LU on **Harpertown**

Figure: Modeling the execution time of the LU on **Barcelona**

- Each core has L1(64 KB), L2(512 KB), and L3(2 MB)
- The results have higher variance
- The deviation is less than 3 %

- The **approach** was validated by modeling the execution time of GER and the LU factorization

- The **approach** was validated by modeling the execution time of GER and the LU factorization
- The experiments were conducted on two **different architectures**

**RWTH**AACHEN
UNIVERSITY

- The **approach** was validated by modeling the execution time of GER and the LU factorization

- The experiments were conducted on two **different architectures**

- The **deviation** is mostly less than 2-3 %