TORSTEN HOEFLER

# Automatic Performance Models for the Masses
**Static and dynamic techniques for application performance modeling**

**in collaboration with Alexandru Calotoiu and Felix Wolf @ RWTH Aachen with students Arnamoy Bhattacharyya and Grzegorz Kwasniewski @ SPCL presented at  Co-Design 2015, HPC China, Wuxi, China**

PASC16
Platform for Advanced Scientific Computing Conference
Lausanne Switzerland | 08-10 June 2016

CLIMATE & WEATHER
SOLID EARTH
LIFE SCIENCE
CHEMISTRY & MATERIALS
PHYSICS
COMPUTER SCIENCE & MATHEMATICS
ENGINEERING
EMERGING DOMAINS

sighpc

acm

# Use-cases for performance modeling

1. **Scalability bug prediction**
   - Find latent scalability bugs early on (before machine deployment)

     *SC13: A. Calotoiu, TH, M. Poke, F. Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes*

2. **Automated performance (regression) testing**
   - Performance modeling as part of a software engineering discipline in HPC

     *ICS'15: S. Shudler, A. Calotoiu, T. Hoefler, A. Strube, F. Wolf: Exascaling Your Library: Will Your Implementation Meet Your Expectations?*
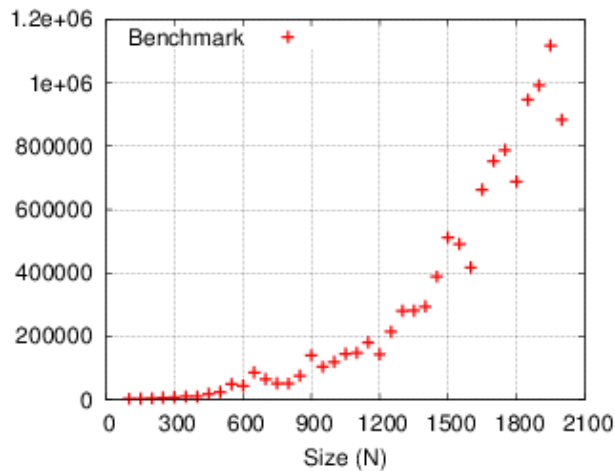
3. **Guided or automated performance optimization**
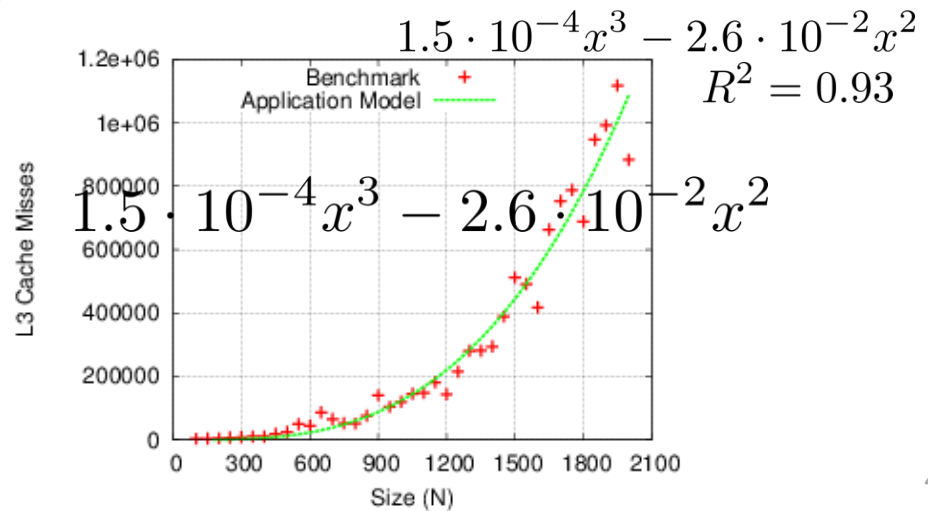   - E.g., near-optimal job scheduling

     *Nan Ding, Wei Xue, et al. (forthcoming)*

4. **Hardware/Software co-design (how to architect systems)**
   - Zhiwei Xu's "efficiency first" design



$$1.5 \cdot 10^{-4} x^3 - 2.6 \cdot 10^{-2} x^2$$
$$R^2 = 0.93$$

**vs.**

$$1.5 \cdot 10^{-4} x^3 - 2.6 \cdot 10^{-2} x^2$$

4

# But how to measure and report performance?

- **We all think we know it but it's harder than I thought!**

    *How many measurements?*
    *How to summarize data?*
    *How to summarize many
       processes?*

    *…*

- **Attempt to establish a rigorous practice**
    - Clarify common problems
    *e.g., Which mean to use when,
    common statistics issues, …*
    - Good start for students
    *12 simple concise rules*

- **My thesis: give up on (performance) reproducibility?**



TH, R. Belli: Scientific Benchmarking of Parallel Computing Systems, SC15 (online already!)

5

# Manual analytical performance modeling

**Identify kernels**
- Parts of the program that dominate its performance at larger scales
- Identified via small-scale tests and intuition

**Create models**
- Laborious process
- Still confined to a small community of skilled experts

- **Disadvantages**
  - Time consuming
  - Error-prone, may overlook unscalable code

# Our first step: scalability bug detector

```
main() {
  foo()
  bar()
  compute()
}
```

Instrumentation

• All functions

Performance measurements (profiles)

$p_1 = 128$      $p_4 = 1,024$
$p_2 = 256$      $p_5 = 2,048$
$p_3 = 512$      $p_6 = 4,096$

Weak scaling

**Input**

Automated modeling

**Output**

Ranking:
1. Asymptotic
2. Target scale $p_t$

1. foo
2. compute
3. main
4. bar
[…]

7

# Primary focus on scaling trend



**Common performance analysis chart in a paper**

## Our ranking

1. $F_1$
2. $F_3$
3. $F_2$

# Primary focus on scaling trend



## Our ranking

1. $F_1$
2. $F_3$
3. $F_2$

**Actual measurement in laboratory conditions**

# Primary focus on scaling trend

## Our ranking



1. $F_1$
2. $F_3$
3. $F_2$

**ETH** *zürich*

# How to mechanize the expert? → Survey!

Computation

LU
$t(p) \sim c$

FFT
$t(p) \sim \log_2(p)$

Naïve N-body
$t(p) \sim p$

…

Samplesort
$t(p) \sim p^2 \log_2^2(p)$

Communication

LU
$t(p) \sim c$

FFT
$t(p) \sim \log_2(p)$

Naïve N-body
$t(p) \sim p$

…

Samplesort
$t(p) \sim p^2$

# Survey result: performance model normal form

$$f(p) = \overset{n}{\underset{k=1}{\mathring{a}}} c_k \times p^{i_k} \times \log_2^{j_k}(p)$$

$$
\begin{aligned}
&n \in \mathbb{N} \\
&i_k \in I \\
&j_k \in J \\
&I, J \in \mathbb{Q}
\end{aligned}
$$

$n = 1$

$I = \{0, 1, 2\}$

$J = \{0, 1\}$

$c_1$        $c_1 \times \log(p)$

$c_1 \times p$        $c_1 \times p \times \log(p)$

$c_1 \times p^2$        $c_1 \times p^2 \times \log(p)$

A. Calotoiu, T. Hoefler, M. Poke, F. Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes, SC13

# Survey result: performance model normal form

$$f(p) = \mathring{\mathop{a}\limits_{k=1}^{n}} c_k \times p^{i_k} \times \log_2^{j_k}(p)$$

$n \in \mathbb{N}$
$i_k \in I$

$n = 2$

$I = \{0, 1, 2\}$

$J = \{0, 1\}$

$c_1 + c_2 \times p$

$c_1 + c_2 \times p^2$

$c_1 + c_2 \times \log(p)$

$c_1 + c_2 \times p \times \log(p)$

$c_1 + c_2 \times p^2 \times \log(p)$

$c_1 \cdot \log(p) + c_2 \cdot p$

$c_1 \cdot \log(p) + c_2 \cdot p \cdot \log(p)$

$c_1 \cdot \log(p) + c_2 \cdot p^2$

$c_1 \cdot \log(p) + c_2 \cdot p^2 \cdot \log(p)$

$c_1 \cdot p + c_2 \cdot p \cdot \log(p)$

$c_1 \cdot p + c_2 \cdot p^2$

$c_1 \cdot p + c_2 \cdot p^2 \cdot \log(p)$

$c_1 \cdot p \cdot \log(p) + c_2 \cdot p^2$

$c_1 \cdot p \cdot \log(p) + c_2 \cdot p^2 \cdot \log(p)$

$c_1 \cdot p^2 + c_2 \cdot p^2 \cdot \log(p)$

# Our automated generation workflow



A. Calotoiu, T. Hoefler, M. Poke, F. Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes, SC13

# Model refinement



Input data

$$\{(p_1, t_1), ..., (p_6, t_6)\}$$

$$n = 1; \overline{R}_0^2 = -\yen$$

Hypothesis generation;
hypothesis size $n$

$$c_1 \qquad\qquad c_1 \times \log(p)$$
$$c_1 \times p \qquad\qquad c_1 \times p \times \log(p)$$
$$c_1 \times p^2 \qquad\qquad c_1 \times p^2 \times \log(p)$$

Hypothesis evaluation
via cross-validation

$$c_1 \times \log(p)$$

Computation of $\overline{R}_n^2$
for best hypothesis

$$R^2 = 1 - \frac{residualSumSquares}{totalSumSquares}$$

$$\overline{R}^2 = 1 - (1 - R^2) \cdot \frac{n}{6 - n - 1}$$

No

$$\overline{R}_{n-1}^2 > \overline{R}_n^2 \ \acute{U}$$

$$n = n_{\max}$$

$$n++$$

Yes

Scaling model

$$I = \{0, 1, 2\}; J = \{0, 1\}; n_{\max} = 2$$

# Evaluation overview



$$I = \left\{ \frac{0}{2}, \frac{1}{2}, \frac{2}{2}, \frac{3}{2}, \frac{4}{2}, \frac{5}{2}, \frac{6}{2} \right\}$$

$$J = \{0, 1, 2\}$$

$$n = 5$$

| Sweep3D | MILC | HOMME | XNS |
|---------|------|-------|-----|

# HOMME

- **Core of the Community Atmospheric Model (CAM)**
- Spectral element dynamical core on a cubed sphere grid



| Kernel<br>[3 of 194] | Model [s]<br>t = f(p) | Predictive error [%]<br>$p_t = 130k$ |
|---|---|---|
| box_rearrange → MPI_Reduce | $0.026 + 2.53{\times}10^{-6}\,p{\times}\sqrt{p} + 1.24{\times}10^{-12}\,p^3$ | 57.02 |
| vlaplace_sphere_vk | 49.53 | 99.32 |
| compute_and_apply_rhs | 48.68 | 1.65 |

$$p_i \leq 15\text{k}$$

# HOMME

- **Core of the Community Atmospheric Model (CAM)**
- Spectral element dynamical core on a cubed sphere grid



| Kernel [3 of 194] | Model [s] t = f(p) | Predictive error [%] $p_t$ = 130k |
|---|---|---|
| box_rearrange → MPI_Reduce | $3.63 \times 10^{-6} p \times \sqrt{p} + 7.21 \times 10^{-13} p^3$ | 30.34 |
| vlaplace_sphere_vk | $24.44 + 2.26 \times 10^{-7} p^2$ | 4.28 |
| compute_and_apply_rhs | 49.09 | 0.83 |

$$p_i \pounds 43k$$

# HOMME

# It works, great! Or not?

- **We face several problems:**
  - Multiple models – when did we collect enough data?
    *if(np < 1.000) a(); else b();*
  - Multiparameter modeling – search space explosion
    *Interesting instance of the curse of dimensionality*
  - Modeling overheads – traces do not scale
    *Cross validation (leave-one-out) is slow and*
    *Our current profiling requires a lot of storage (>TBs)*

# First step: simple compiler analyses

- **Automatic kernel detection in LLVM**
  - Loop call graph – each loop/function as kernel (recursively)
  - Determine relevant input parameters for each kernel
  
  *Massive pruning possible!*

- **Online model generation (using online LASSO)**
  
  *Constant (little) amount of data stored*

- **Automatic profiling rate limiting**
  
  *Profile less as models gain confidence*



Quality: NAS UA and Mantevo MiniFE

Overhead: Mantevo

# Second step: counting loop iterations

```
for (j = 1; j <= n; j = j*2)
        for (k = j; k <= n; k = k++)
                veryComplicatedOperation(j,k);
```

Polyhedral model



$$N = (n+1)\log_2 n - n + 2$$

A.I. Barvinok. A polynomial time algorithm for counting integral points in polyhedra when the dimension is fixed, Math. Oper. Res., 1994

**ETH** *zürich*

# Counting arbitrary affine loop nests

- ## Affine loops

```
x=x0;                    // Initial assignment
while(c^T x < g)         // Loop guard
    x=Ax + b;            // Loop update
```

- ## Perfectly nested affine loops

```
while(c_1^T x < g_1) {
    x = A_1 x + b_1;
    while(c_2^T x < g_2) {
        . . .
        x = A_{k-1} x + b_{k-1};
        while(c_k^T x < g_k) {
            x = A_k x + b_k;
            while(c_{k+1}^T x < g_{k+1})  { . . .  }
            x = U_k x + v_k;  }
        x = U_{k-1} x + v_{k-1};
        . . .}
    x = U_1 x + v_1;}
```

$A_k, U_k \in \mathbb{R}^{m \times m}$, $b_k, v_k, c_k \in \mathbb{R}^m$, $g_k \in \mathbb{R}$ and $k = 1 \ldots r$.

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$
\begin{aligned}
&\texttt{while}(c_1^T x < g_1) \ \{ \\
&\quad x = A_1 x + b_1; \\
&\quad \texttt{while}(c_2^T x < g_2) \ \{ \\
&\qquad \dots \\
&\qquad x = A_{k-1} x + b_{k-1}; \\
&\qquad \texttt{while}(c_k^T x < g_k) \ \{ \\
&\qquad\quad x = A_k x + b_k; \\
&\qquad\quad \texttt{while}(c_{k+1}^T x < g_{k+1}) \ \{ \dots \ \} \\
&\qquad\quad x = U_k x + v_k; \ \} \\
&\qquad x = U_{k-1} x + v_{k-1}; \\
&\qquad \dots \} \\
&\quad x = U_1 x + v_1; \}
\end{aligned}
$$

**ETH**zürich

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}\begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

```
while(c_1^T x < g_1) {
  x = A_1 x + b_1;
  while(c_2^T x < g_2) {
    ...
    x = A_{k-1} x + b_{k-1};
    while(c_k^T x < g_k) {
      x = A_k x + b_k;
      while(c_{k+1}^T x < g_{k+1}) { ... }
      x = U_k x + v_k; }
    x = U_{k-1} x + v_{k-1};
    ...}
  x = U_1 x + v_1;}
```

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

$$while\left( \begin{pmatrix} 1 & 0 \end{pmatrix} \begin{pmatrix} j \\ k \end{pmatrix} < n/p + 1 \right)\{$$

$$\texttt{while}(c_1^T x < g_1) \; \{$$
$$x = A_1 x + b_1;$$
$$\quad \texttt{while}(c_2^T x < g_2) \; \{$$
$$\quad \quad \ldots$$
$$\quad \quad x = A_{k-1} x + b_{k-1};$$
$$\quad \quad \texttt{while}(c_k^T x < g_k) \; \{$$
$$\quad \quad \quad x = A_k x + b_k;$$
$$\quad \quad \quad \texttt{while}(c_{k+1}^T x < g_{k+1}) \; \{ \ldots \; \}$$
$$\quad \quad \quad x = U_k x + v_k; \; \}$$
$$\quad \quad x = U_{k-1} x + v_{k-1};$$
$$\quad \ldots \}$$
$$x = U_1 x + v_1; \}$$

$$\}$$

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

$$while((1 \quad 0) \begin{pmatrix} j \\ k \end{pmatrix} < n/p + 1)\{$$

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$while((0 \quad 1) \begin{pmatrix} j \\ k \end{pmatrix} < m)\{$$

$$while(c_1^T x < g_1) \ \{$$
$$\quad x = A_1 x + b_1;$$
$$\quad while(c_2^T x < g_2) \ \{$$
$$\quad \quad \ldots$$
$$\quad \quad x = A_{k-1} x + b_{k-1};$$
$$\quad \quad while(c_k^T x < g_k) \ \{$$
$$\quad \quad \quad x = A_k x + b_k;$$
$$\quad \quad \quad while(c_{k+1}^T x < g_{k+1}) \ \{ \ldots \ \}$$
$$\quad \quad \quad x = U_k x + v_k; \ \}$$
$$\quad \quad x = U_{k-1} x + v_{k-1};$$
$$\quad \quad \ldots \}$$
$$x = U_1 x + v_1; \}$$

$$\}$$

$$\}$$

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}\begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

$$while((1 \quad 0)\begin{pmatrix} j \\ k \end{pmatrix} < {n}/{p}+1)\{$$

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}\begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$while((0 \quad 1)\begin{pmatrix} j \\ k \end{pmatrix} < m)\{$$

$$\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}\begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$\}\begin{pmatrix} j \\ k \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}\begin{pmatrix} j \\ k \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$\}$$

$$\mathtt{while}(c_1^T x < g_1) \ \{$$
$$x = A_1 x + b_1;$$
$$\quad \mathtt{while}(c_2^T x < g_2) \ \{$$
$$\quad \dots$$
$$\quad x = A_{k-1} x + b_{k-1};$$
$$\quad \mathtt{while}(c_k^T x < g_k) \ \{$$
$$\quad\quad x = A_k x + b_k;$$
$$\quad\quad \mathtt{while}(c_{k+1}^T x < g_{k+1}) \ \{\dots\}$$
$$\quad\quad x = U_k x + v_k; \ \}$$
$$\quad x = U_{k-1} x + v_{k-1};$$
$$\quad \dots\}$$
$$x = U_1 x + v_1;\}$$

# Counting arbitrary affine loop nests

- **Example**

```
for (j=1; j < n/p + 1; j= j*2)
        for (k=j; k < m; k = k + j )
                veryComplicatedOperation(j,k);
```

$$x = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} x + \begin{pmatrix} 1 \\ 0 \end{pmatrix};$$

$$while((1 \quad 0)x < \frac{n}{p} + 1) \{$$

$$x = \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix} x + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$while((0 \quad 1)x < m) \{$$

$$x = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} x + \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\} x = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix} x + \begin{pmatrix} 0 \\ 0 \end{pmatrix};$$

$$\}$$

---

$$\texttt{while}(c_1^T x < g_1) \ \{$$
$$\quad x = A_1 x + b_1;$$
$$\quad \texttt{while}(c_2^T x < g_2) \ \{$$
$$\quad\quad \ldots$$
$$\quad\quad x = A_{k-1} x + b_{k-1};$$
$$\quad\quad \texttt{while}(c_k^T x < g_k) \ \{$$
$$\quad\quad\quad x = A_k x + b_k;$$
$$\quad\quad\quad \texttt{while}(c_{k+1}^T x < g_{k+1}) \ \{\ldots\ \}$$
$$\quad\quad\quad x = U_k x + v_k; \ \}$$
$$\quad\quad x = U_{k-1} x + v_{k-1};$$
$$\quad\quad \ldots \}$$
$$\quad x = U_1 x + v_1; \}$$

where $\quad x = \begin{pmatrix} j \\ k \end{pmatrix}$

T. Hoefler, G. Kwasniewski: Automatic Complexity Analysis of Explicitly Parallel Programs, SPAA'14

# Overview of the whole system



**Parallel program**

```
do i = , procCols
    call mpi_irecv( buff,  , dp_type, reduce_exch_proc(i),
                    i, mpi_comm_world, request, ierr )
    call mpi_send( buff2, , dp_type, reduce_exch_proc(i),
                    i, mpi_comm_world, ierr )
    call mpi_wait( request, status, ierr )
enddo

do i = id *n/p, ( id + )* n/p
    do j = , nSize
        call compute
```

**LLVM**

**Loop extraction**

**Affine loop synthesis**

$$\mathbf{while}\,(c_1^T x < g_1)\;\{$$
$$x = A_1 x + b_1;$$
$$\mathbf{while}\,(c_2^T x < g_2)\;\{$$
$$\dots$$
$$x = A_{k-1} x + b_{k-1};$$
$$\mathbf{while}\,(c_k^T x < g_k)\;\{$$
$$x = A_k x + b_k;$$
$$\mathbf{while}\,(c_{k+1}^T x < g_{k+1})\;\{\dots\}$$
$$x = U_k x + v_k;\;\}$$
$$x = U_{k-1} x + v_{k-1};$$
$$\dots\}$$
$$x = U_1 x + v_1;\}$$

**Closed form representation**

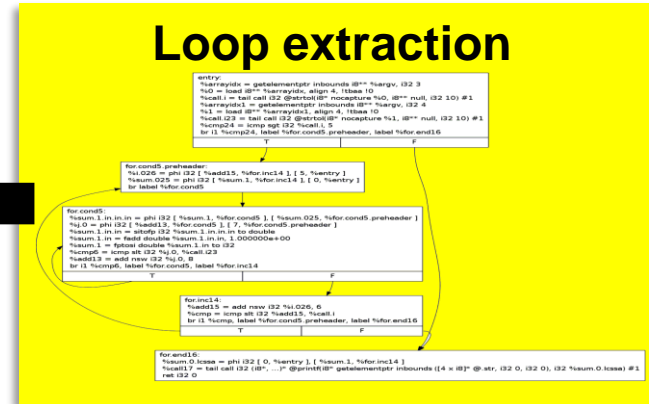$$x(i_1,\dots,i_r) = A_{final}(i_1,\dots,i_r) \cdot x_0 + b_{final}(i_1,\dots,i_r)$$

with

$$i_r = 0..n_k(x_{0,k}), k = 1...r$$

**Number of iterations**

$$N = \sum_{i_1=0}^{n_1(x_{0,1})} \sum_{i_2=0}^{n_2(x_{0,2})} \dots \sum_{i_{r-1}=0}^{n_{r-1}(x_{0,r-1})} n_r(x_{0,r}).$$

**Program analysis**

$$W = N\Big|_{p=1}$$

$$D = N\Big|_{p\to\infty}$$

# Case study: NAS EP

$$N(m,p) = \left\lceil \frac{2^{m-16} \cdot (u + 2^{16})}{p} \right\rceil$$

$$W = T_1 \approx 2^m$$
$$D = T_\infty \approx 1$$

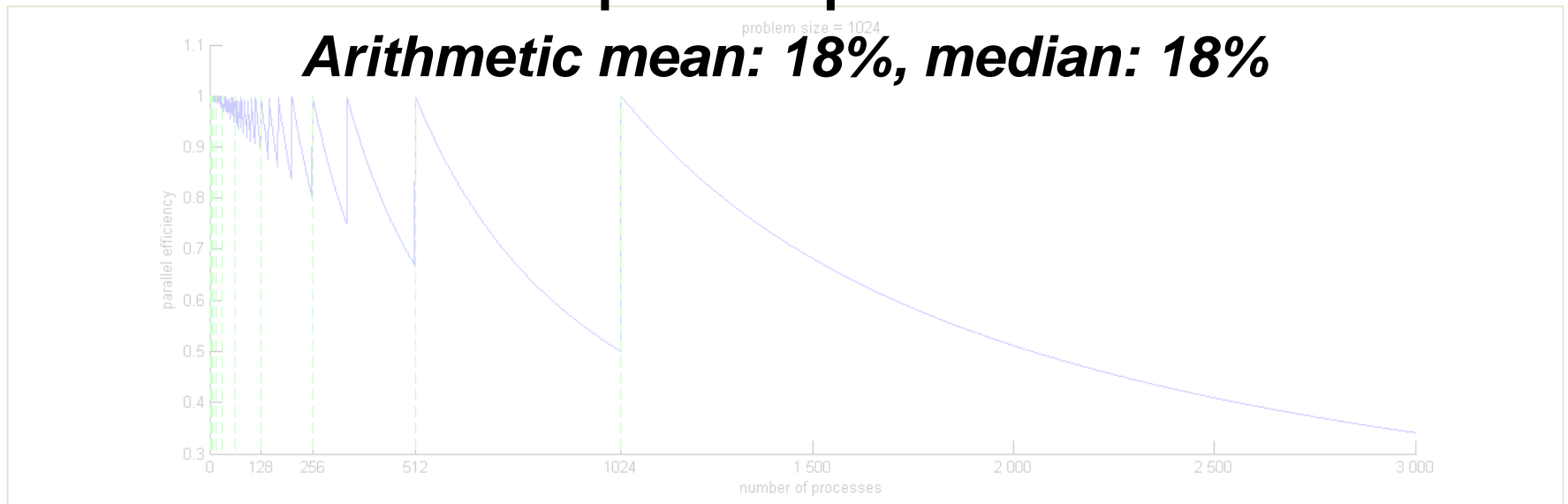$$E_P = \frac{2^m}{p \left\lceil \frac{2^m}{p} \right\rceil}$$

```
u:    do i=1,100
      ik =kk/2
      if (ik .eq. 0) goto 130
      continue
```

**15 applications (NAS/Mantevo/Mibench):**
- **100% of loops were treated (with unknowns)**
- **9-45% of loops were predicted exact**
  - *Arithmetic mean: 18%, median: 18%*

# What problems are remaining?

- **Well, what about non-affine loops?**
  - More general abstract interpretation (next step)
  - Not decidable → will always have undefined terms

$$N = \frac{\mathtt{na} \cdot u}{\mathtt{nprows}}$$

- **Back to PMNF?**
  - Generalize to multiple input parameters
    *a) Bigger search-space* ☹
    *b) Bigger trace files* ☹

$$f(p) = \sum_{k=1}^{n} c_k \times p^{i_k} \times \log_2^{j_k}(p)$$

- **Combine static (loop counting) and dynamic approach (PMNF)**
  - Find number of loop iterations, replace $u_x$ with $\mathrm{PMNF}_x$
  - Find similar kernels (use only one PMNF for similar ones)
  - Remove irrelevant input parameters for each PMNF
  - Other simple optimizations: batch model update, etc.
  - Result: higher accuracy, lower overheads

# Performance Analysis 2.0 – Automatic Models

**A call for action: use performance modeling for rigorous designs**

- **Especially for co-design (systems and applications)!**
  *New architectures, e.g., FPGA assessment*
- **High-performance programming as a science**
  *Learn from natural sciences!*
- **Start with the students: teach rigorous analysis and modeling**

A. Calotoiu, T. Hoefler, M. Poke, F. Wolf: Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes. *Supercomputing (SC13).*

T. Hoefler, G. Kwasniewski: Automatic Complexity Analysis of Explicitly Parallel Programs. *SPAA 2014.*
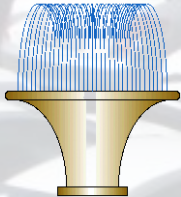
A. Bhattacharyya, T. Hoefler: PEMOGEN: Automatic Adaptive Performance Modeling during Program Runtime, *PACT 2014*

S. Shudler, A. Calotoiu, T. Hoefler, A. Strube, F. Wolf: Exascaling Your Library: Will Your Implementation Meet Your Expectations? *ICS 2015*

A. Bhattacharyya, G. Kwasniewski, T. Hoefler: Using Compiler Techniques to Improve Automatic Performance Modeling, *PACT 2015*