# Performance Prediction based on Inherent Program Similarity

## PACT-2006

*September 19th 2006*

Seattle, Washington

*Kenneth Hoste*[‡], Aashish Phansalkar[†], Lieven Eeckhout[‡],

Andy Georges[‡], Lizy K. John[†] and Koen De Bosschere[‡]

[‡] ELIS, Ghent University, Belgium                    [†] ECE, The University of Texas at Austin

# The performance of a platform is evaluated using benchmarks



**MacBook Performance Benchmarks**
SPEC performance: Up to five times faster than the iBook G4.(2)

| iBook G4 1.42GHz | MacBook Core Duo 2.0GHz | Δ |
|---|---|---|
| 5.7 | 29.1 | 5.1X |

SPECint_rate_base2000
Integer calculation (estimate)

| iBook G4 1.42GHz | MacBook Core Duo 2.0GHz | Δ |
|---|---|---|
| 4.3 | 24.7 | 5.7X |

SPECfp_rate_base2000
Floating-point calculation (estimate)

`www.apple.com/macbook/intelcoreduo.html, May 2006`

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 1/25

UNIVERSITEIT GENT

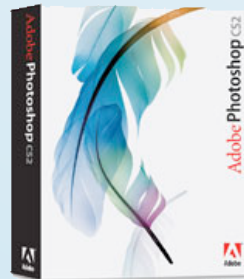# How does that platform perform for my application(s) of interest?

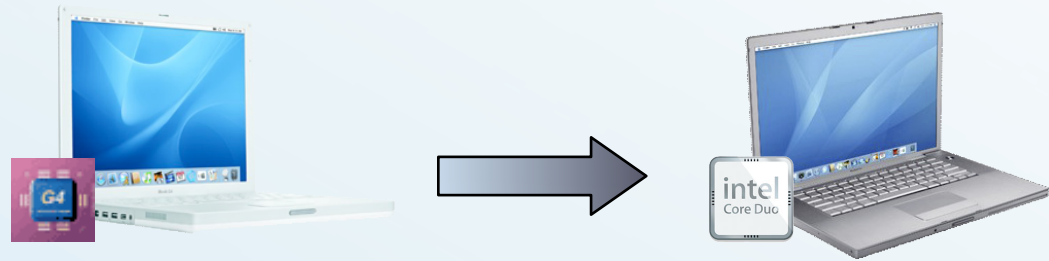## ubiquitous problem in benchmarking



Excel (spreadsheet)

R (statistics)

Photoshop (image processing)

Virtual PC (Windows on Mac)

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 2/25

UNIVERSITEIT
GENT

# Evaluating several platforms for applications of interest is troublesome

✦ Porting

✦ Hardware availability

Playstation 3

✦ Time constraints

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 3/25

UNIVERSITEIT
GENT

# We estimate performance based on program similarity



application of interest

benchmarks

benchmark space

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 4/25

UNIVERSITEIT
GENT

# How do we characterize a program without executing it on each platform?

on

application of interest

**microarchitecture-independent program characteristics**

in our setup: characterization on Alpha using ATOM

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 5/25

UNIVERSITEIT GENT
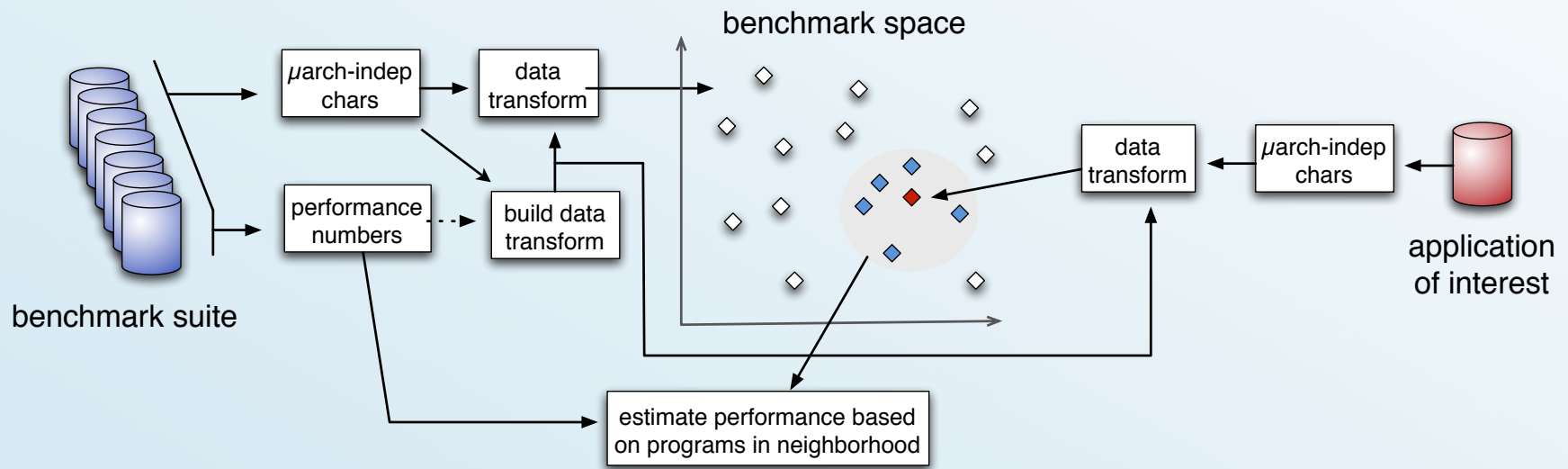
# Microarchitecture-independent program characteristics

6 categories:

‣ instruction level parallelism (ILP)

‣ instruction mix

‣ branch predictability

‣ register traffic

‣ data stream strides

‣ working set size

**totaling 47 program characteristics**
more details are available in the paper

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 6/25

UNIVERSITEIT
GENT

# Estimating performance based on inherent program similarity

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 7/25

UNIVERSITEIT
GENT

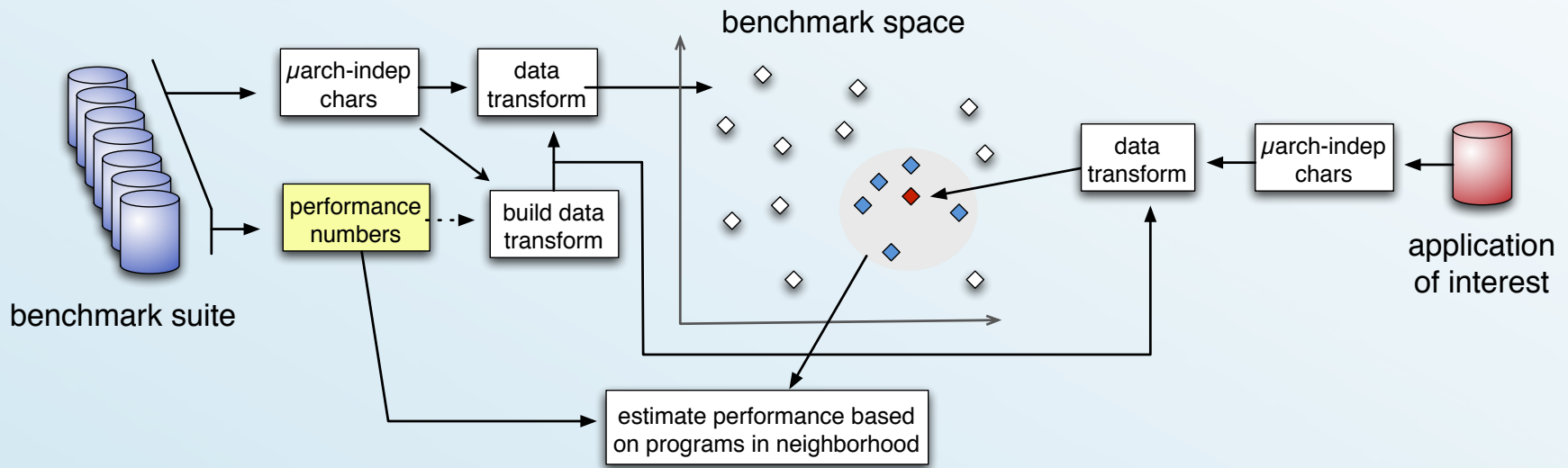# Estimating performance based on inherent program similarity

## *Step 1*



collect program characteristics for the benchmark suite

by instrumenting the benchmarks using ATOM/PIN

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 8/25
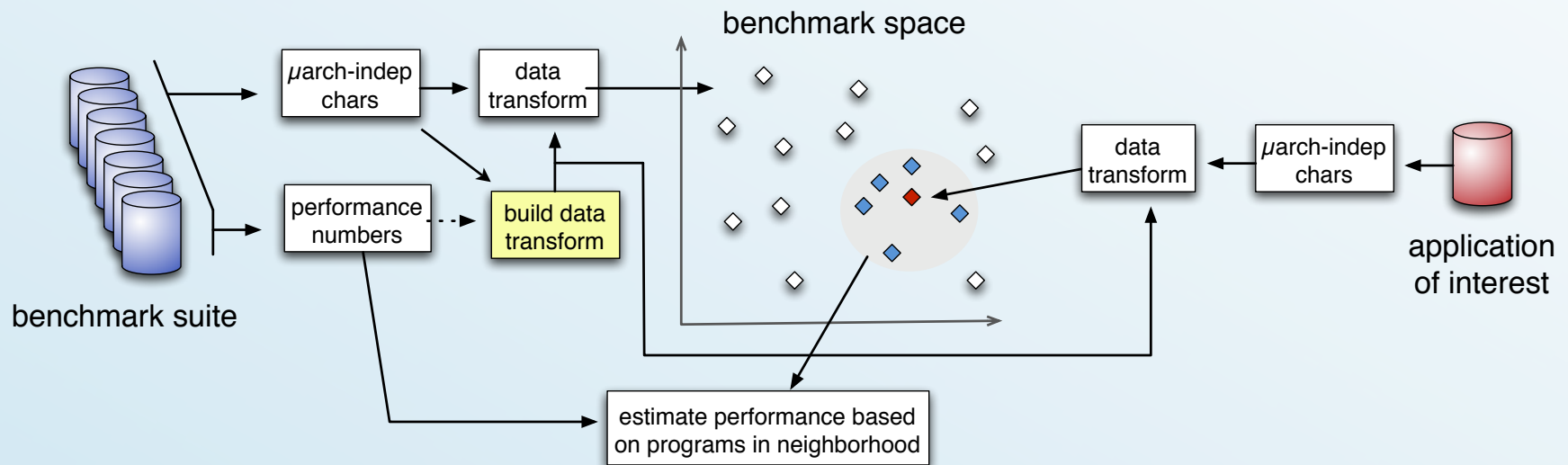
# Estimating performance based on inherent program similarity

## *Step 2*



obtain performance numbers for each benchmark

for example, use the ones published on `spec.org` for SPEC CPU2000

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 9/25

UNIVERSITEIT GENT

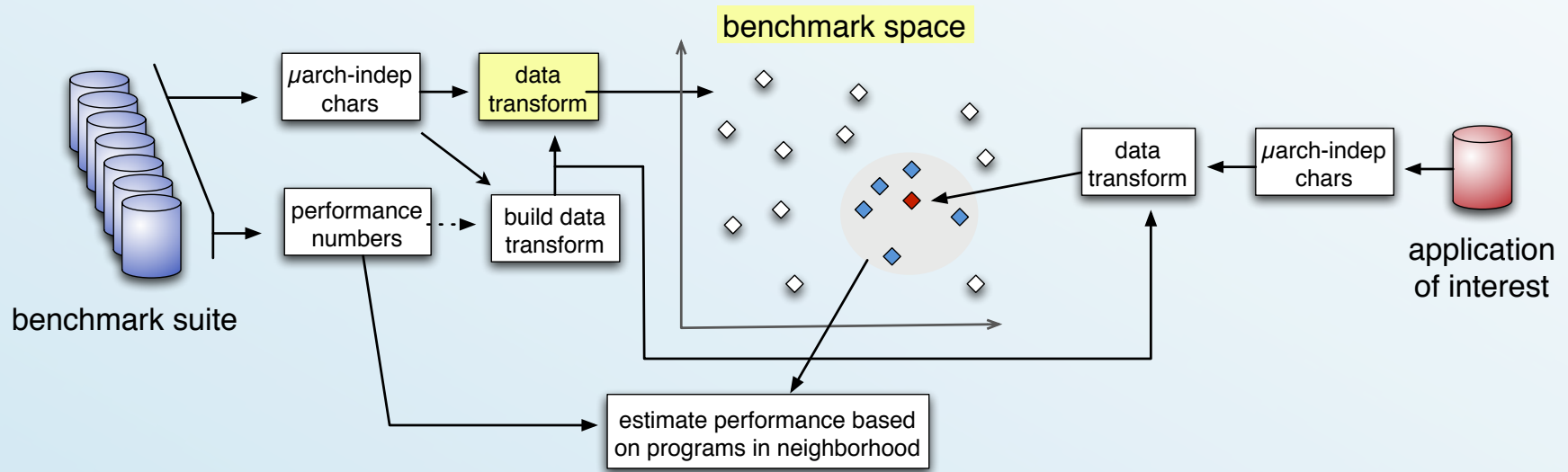# Estimating performance based on inherent program similarity

## *Step 3*



build the data transformation matrix

we evaluate 3 approaches: normalization, PCA and a genetic algorithm

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 10/25

UNIVERSITEIT GENT

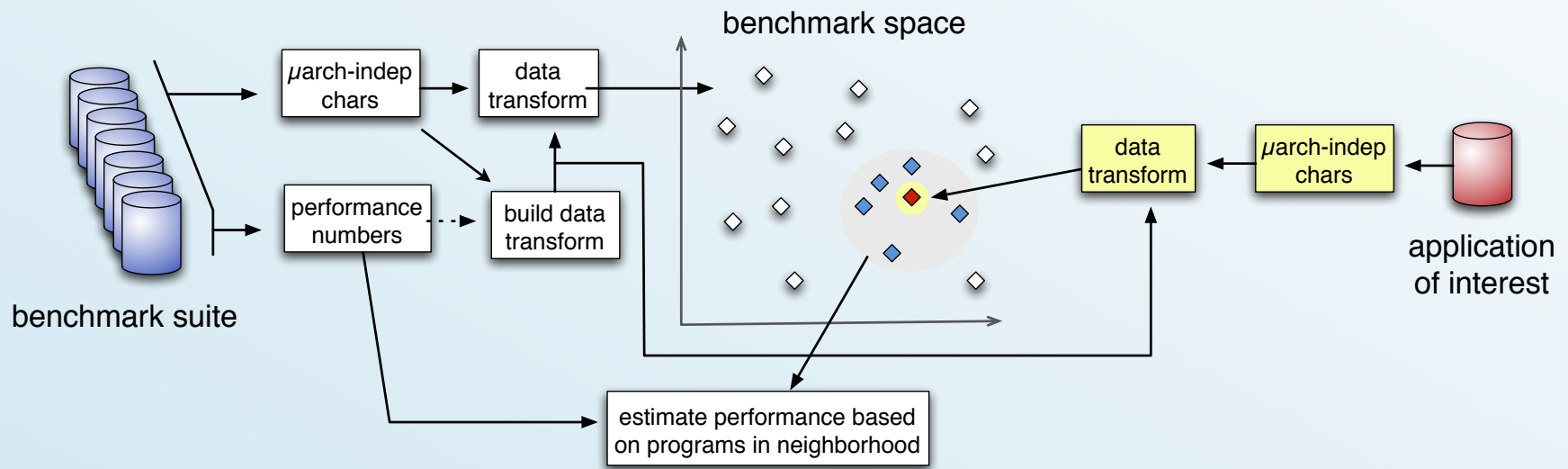# Estimating performance based on inherent program similarity

## *Step 4*



use the program characteristics and data transformation matrix to build the benchmark space

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 11/25

# Estimating performance based on inherent program similarity

*Step 5*



locate the application of interest
in the benchmark space

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 12/25

UNIVERSITEIT
GENT

# Estimating performance based on inherent program similarity

## *Step 6*



estimate performance of the application of interest
using the benchmarks in the neighborhood

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 13/25

UNIVERSITEIT GENT

# Building the data transformation matrix

**problem:**

program characteristics vary in range

*ILP vs instruction mix*

Euclidean distance measure will be biased

**solution:**

normalize characteristics

mean = 0, variance = 1

**data transformation:**

subtract mean and divide by standard deviation

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 14/25

UNIVERSITEIT
GENT

# Building the data transformation matrix

**problem:**

program characteristics are correlated

➡ Euclidean distance gives a higher weight to correlated characteristics

**solution:**

obtain uncorrelated characteristics using Principal Components Analysis (PCA)

**data transformation:**

perform PCA on norm. chars, and normalize PCs

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 15/25

UNIVERSITEIT GENT

# Building the data transformation matrix

**problem:**

some program characteristics are more important for estimating performance than others
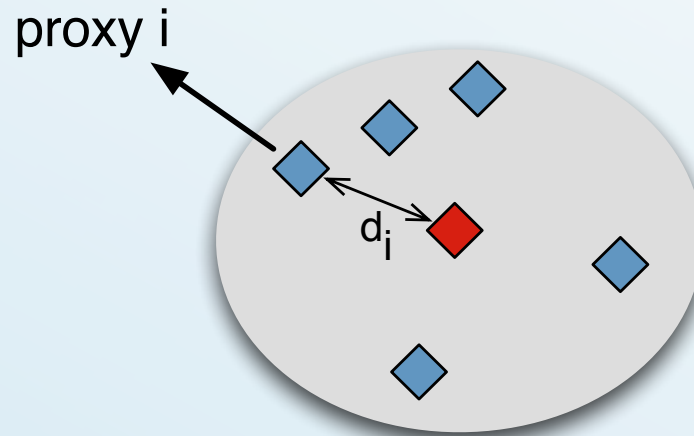
*branch predictability n vs % multiply operations*

**solution:** genetic algorithm

learning how to scale the characteristics so that distance in the benchmark space correlates with difference in performance

**data transformation:**

weigh normalized characteristics

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 16/25

# Estimating performance using application proxies



proxy i

$d_i$

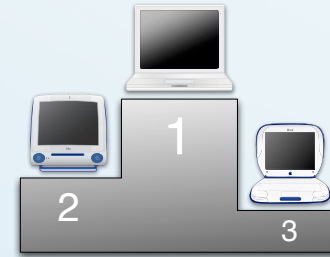$$w_i = \frac{\sum_{i=1}^{n} \frac{1}{d_i}}{d_i}$$

weight for each proxy i

$$S = \frac{1}{\sum_{i=1}^{n} \frac{w_i}{S_i}}$$

weighted harmonic mean

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 17/25

# Determine the hardware platform ranking based on estimated performance

current practice:

rank machines based on average performance



**our approach**:

✦ estimate performance of the application of interest for each machine considered

✦ rank machines based on estimated performance

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 18/25

UNIVERSITEIT
GENT

# Experimental setup

- full SPEC CPU2000 benchmark suite (26 benchmarks)

    methodology is evaluated using *crossvalidation*

    *n-1* benchmarks form benchmark suite,
    *n*th benchmark is application of interest

- speedup numbers for 36 machines

    ✦ different ISAs, processor, configurations and manufacturers

    ✦ taken from SPEC website (`http://www.spec.org`)

UNIVERSITEIT
GENT

# Comparing the estimated ranking with the actual ranking of the machines

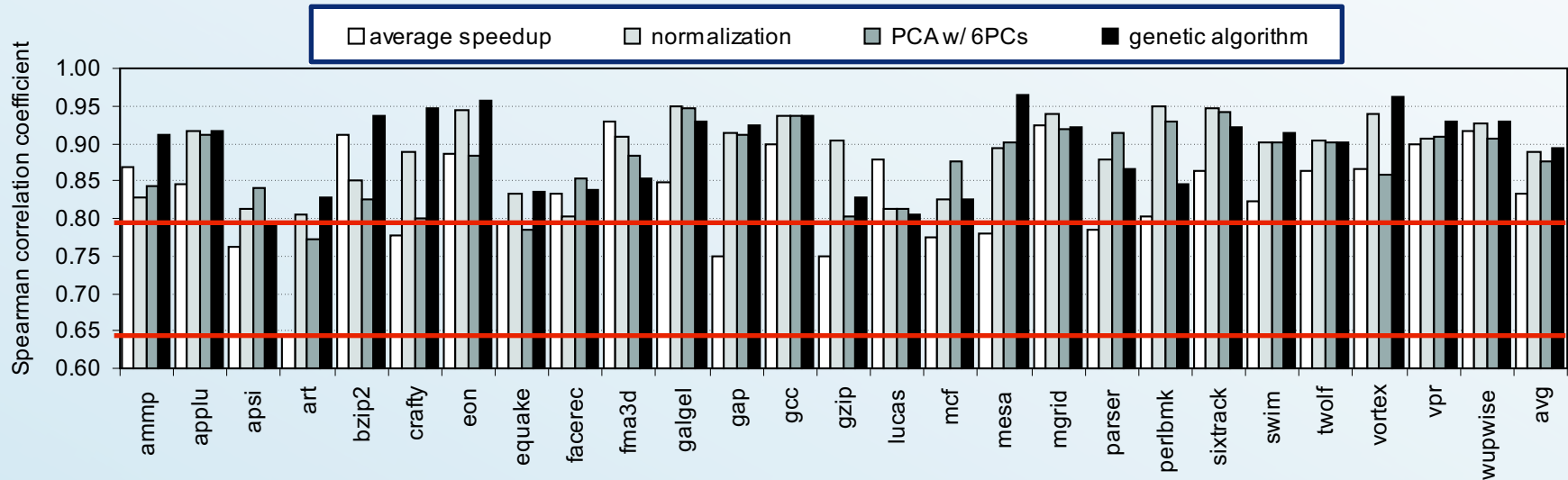Spearman rank correlation coefficient

a value between -1 and 1

quantifies the quality of the estimated ranking
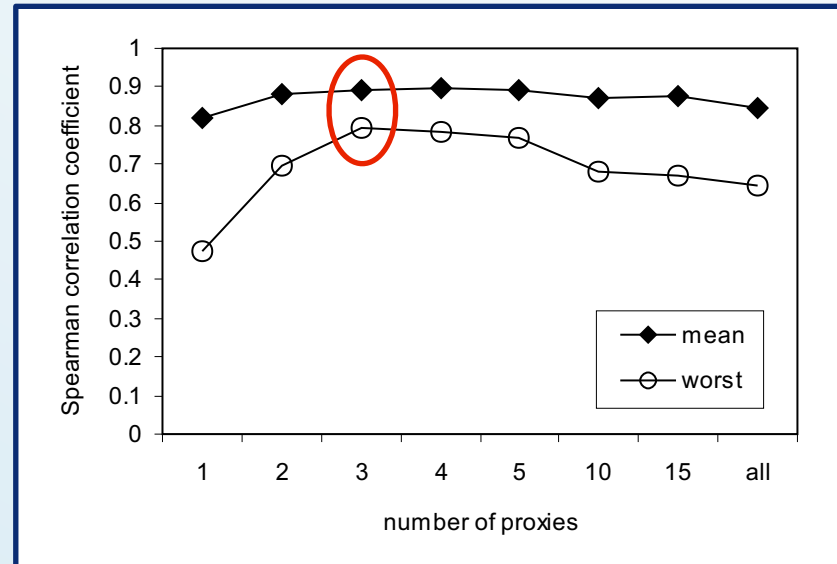
1: perfect

0: random

-1: worst possible (reverse)

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide X/Y

# Our methodology outperforms current practice



| | current practice | norm. | PCA |
|---|---|---|---|
| genetic algorithm | 23/26 | 16/26 | 16/26 |

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 21/25

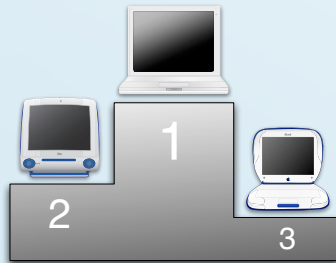UNIVERSITEIT GENT

# How many proxies should we retain?



- a single proxy is not enough

- too many proxies degrades performance

- 3 proxies seems optimal
  (and is used in further evaluation)

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 22/25

UNIVERSITEIT
GENT

# Evaluating the average performance loss

best machine:

current practice: 20%

**our approach: 13.6%**



top 3 machines:

current practice: 19.7%

**our approach: 11.1%**

small differences in rank correlation coefficient can make a big difference in practice

UNIVERSITEIT
GENT

# Some interesting observations

✦ `gcc`, `mcf` and `swim` *never* appear as proxy

✦ `bzip2` is a *frequent* proxy

✦ microarchitecture-independent characteristics seem to be able to capture program behavior across ISAs

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 24/25

UNIVERSITEIT GENT

# Conclusions

✦ identifying the best machine should be done using the characteristics of the application of interest

✦ our methodology using the genetic algorithm is able to scale the benchmark space according to performance differences

✦ some benchmarks have very specific behavior (`gcc`, `mcf, swim`), others exhibit average behavior (`bzip2`)

✦ although program characteristics are measured on one particular ISA, they seem to capture program behavior over different ISAs fairly well

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

slide 25/25

# Questions?

Performance Prediction based on Inherent Program Similarity – Kenneth Hoste – 2006-09-19
Faculty of Engineering – Department of Electronics and Information Systems (ELIS)

UNIVERSITEIT
GENT