# Advanced Computing Laboratory: Three Research Threads

Markus Püschel
Computer Science

**ETH** *zürich*

**SPIRAL**
www.spiral.net

Marcela Zuluaga

Georg Ofenbeck

Victoria Caparròs Cabezas

Daniele Spampinato

Alen Stojanov

François Serre

Gagandeep Singh

---

**(1)** **Program synthesis for performance**
*Daniele Spampinato [CGO 2014]*

**(2)** **Performance bottleneck modeling**
*Victoria Caparròs Cabezas [IISWC 2014]*

**(3)** **Predicting Pareto-optimal solutions**
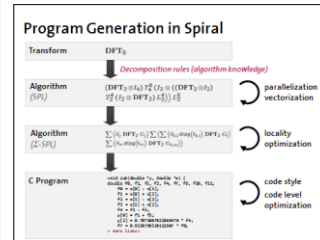*Marcela Zuluaga, Andreas Krause [ICML 2013]*

# Vision:
# Program Synthesis For Performance

Generate highest performance code for mathematical computations directly from a mathematical description

**Approach**

*Mathematical DSLs*
*Rewriting systems for difficult optimizations*
*Compiler*
*Learning and search for fine-tuning*

*Use advanced software platforms*
*for the development of generators*



Example: Linear transforms
www.spiral.net

---

# LGen: Generator for Linear Algebra

$$\gamma = x^T(A + B)y + \delta \quad \longleftarrow \quad \text{A is 2x3, x is 3x1, ...}$$
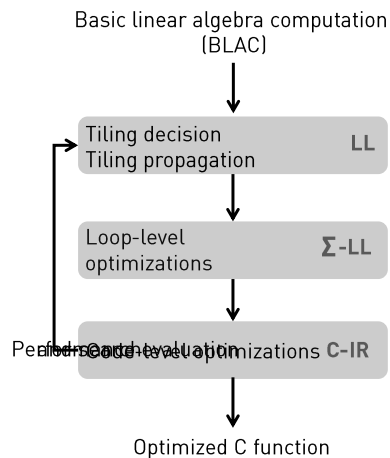
**LGen** *Design similar to Spiral*

```
void f(double const * A, double const * x, double * y) {
  __m128d t0, …;

  t0 = _mm_loadu_pd(A);
  t1 = _mm_load_sd(A + 2);
  ...
  t6 = _mm_hadd_pd(_mm_mul_pd(t0, t4), _mm_mul_pd(t2, t4));
  t7 = _mm_shuffle_pd(t1, t3, 0);
  t8 = _mm_mul_pd(t7, _mm_shuffle_pd(t5, t5, 0));
  t9 = _mm_add_pd(t6, t8);

  _mm_storeu_pd(y, t9);
}
```

# Architecture of LGen

Basic linear algebra computation
(BLAC)

↓

| Tiling decision<br>Tiling propagation | **LL** |

↓

| Loop-level<br>optimizations | **Σ-LL** |

↓

| Code-level optimizations | **C-IR** |

↓

Optimized C function

$$y = Ax$$

$$[y = Ax]_{2,1}$$
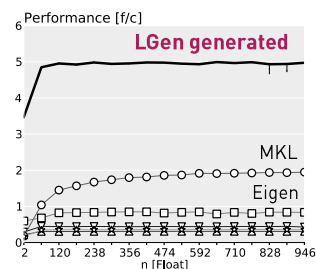
$$y = \sum_{i,j} S_i(G_i \cdots)$$

```
...
Mov (mmMulPs A[0,0], x[0,0]), t[0,0]
...

for(int i = … ) {
  …
  t = _mm_mul_ps(a, x);
  …
}
```

---

# Example Benchmarks

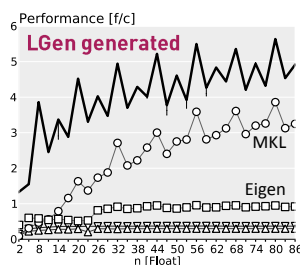**Intel Xeon Westmere**

$$C = \alpha AB + \beta C$$

Performance [f/c]

LGen generated

MKL

Eigen

n [Float]

$$A \in \mathbb{R}^{n \times 4}$$
$$B \in \mathbb{R}^{4 \times 4}$$

**Intel Xeon Westmere**

$$C = \alpha(A_0 + A_1)^T B + \beta C$$

Performance [f/c]

LGen generated

MKL

Eigen

n [Float]

$$A_0 \in \mathbb{R}^{4 \times n}$$
$$B \in \mathbb{R}^{4 \times n}$$

**ARM Cortex-A8**

Performance [f/c]

— LGen
▼ Handwritten fixed (gcc)
▲ Handwritten gen (gcc)
► Handwritten fixed (clang)
◄ Handwritten gen (clang)
■ Eigen-3.2.0
✶ Atlas-3.10.1

LGen generated

n [Float]

## Next Steps

Structured matrices

Higher level algorithms (matrix factorizations etc.)

Multicore

Domain specific extensions
  *Optimization*
  *Machine learning*
  *Communication & Control*

**1** **Program synthesis for performance**
*Daniele Spampinato [CGO 2014]*
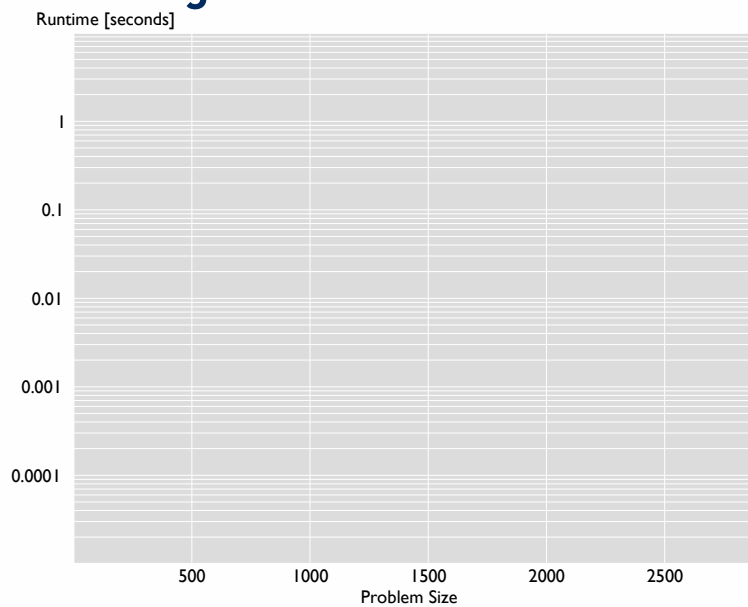
**2** **Performance bottleneck modeling**
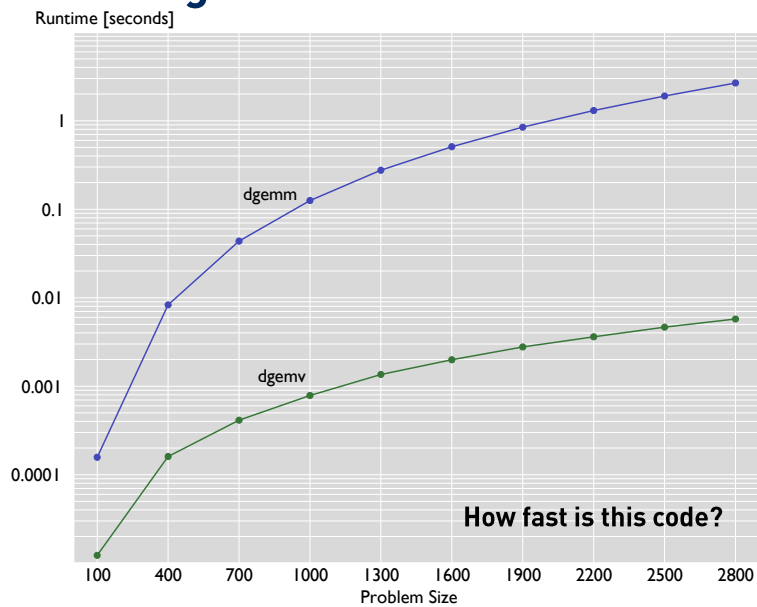*Victoria Caparròs Cabezas [IISWC 2014]*

**3** **Predicting Pareto-optimal solutions**
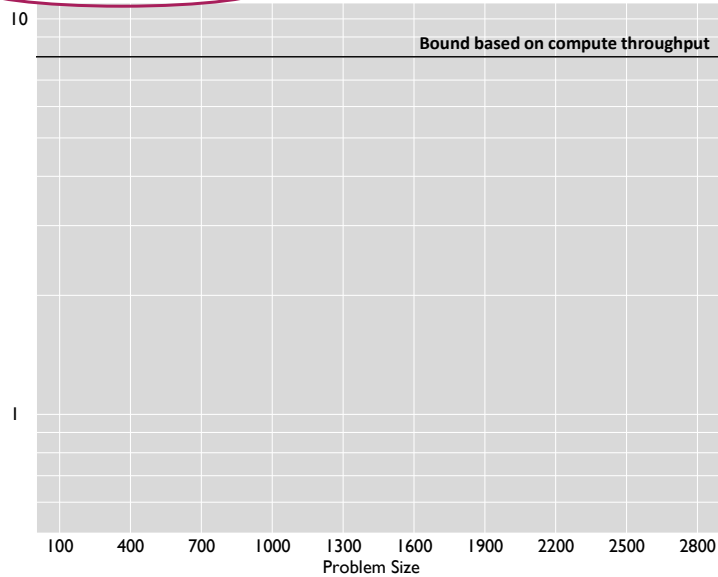*Marcela Zuluaga, Andreas Krause [ICML 2013]*
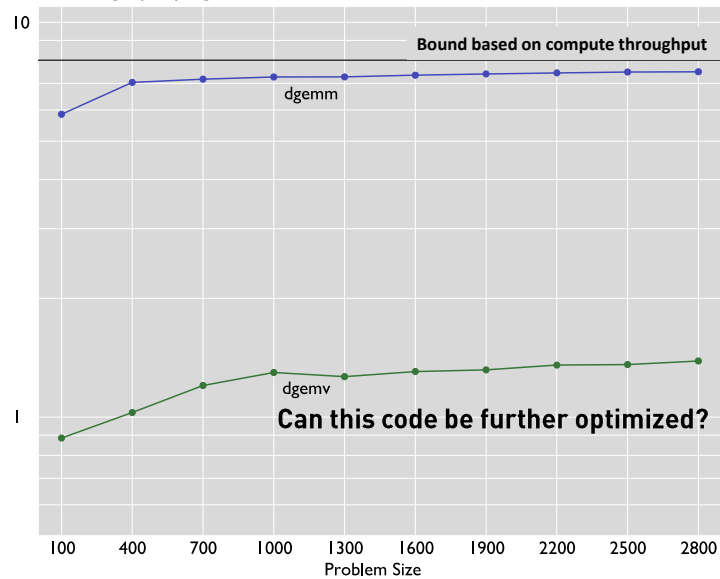
# Measuring Runtime

Runtime [seconds]



Problem Size

# Measuring Runtime

Runtime [seconds]



dgemm

dgemv

**How fast is this code?**

Problem Size

# Measuring Performance

Performance [Flops/Cycle]



Bound based on compute throughput

# Measuring Performance

Performance [Flops/Cycle]



Bound based on compute throughput

dgemm

dgemv

**Can this code be further optimized?**

# Roofline Plot

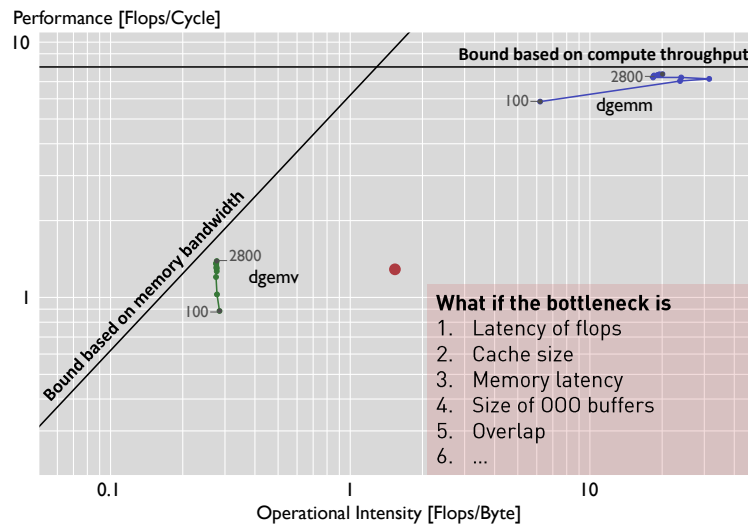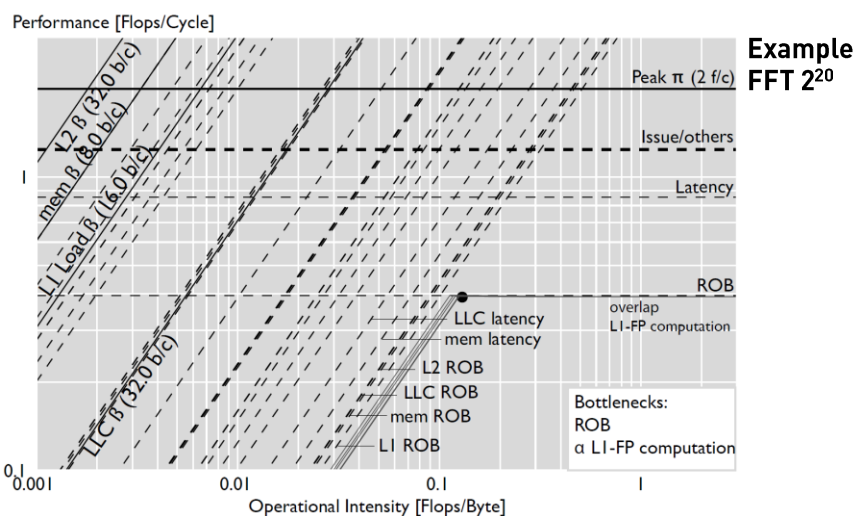Performance [Flops/Cycle]

Operational Intensity [Flops/Byte]

[Williams, 2009] "Roofline: An Insightful Visual Performance Model for Multicore", S. Williams *et al.* Communications of the ACM, 2009

# Roofline Plot

Performance [Flops/Cycle]

**Bound based on compute throughput**

**Bound based on memory bandwidth**

Operational Intensity [Flops/Byte]

[Williams, 2009] "Roofline: An Insightful Visual Performance Model for Multicore", S. Williams *et al.* Communications of the ACM, 2009
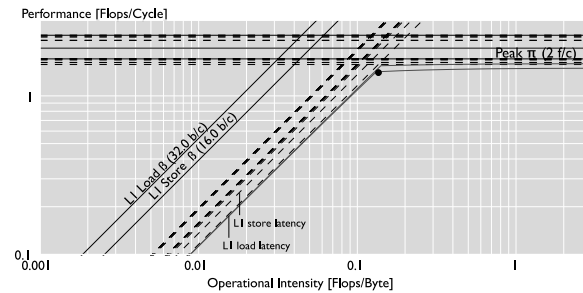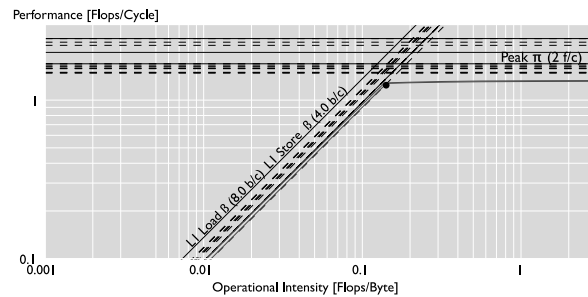
# Roofline Plot

Performance [Flops/Cycle]

Bound based on compute throughput

Bound based on memory bandwidth

2800 · dgemm
100

2800 dgemv
100

What if the bottleneck is
1. Latency of flops
2. Cache size
3. Memory latency
4. Size of OOO buffers
5. Overlap
6. ...

Operational Intensity [Flops/Byte]



# Bottleneck Modelling

Performance [Flops/Cycle]

Example
FFT $2^{20}$

Peak π (2 f/c)

L2 B (32.0 b/c)
mem B (8.0 b/c)
L1 Load B (16.0 b/c)

Issue/others

Latency

ROB

LLC B (32.0 b/c)

LLC latency
mem latency
L2 ROB
LLC ROB
mem ROB
L1 ROB

overlap
L1-FP computation

Bottlenecks:
ROB
α L1-FP computation

Operational Intensity [Flops/Byte]

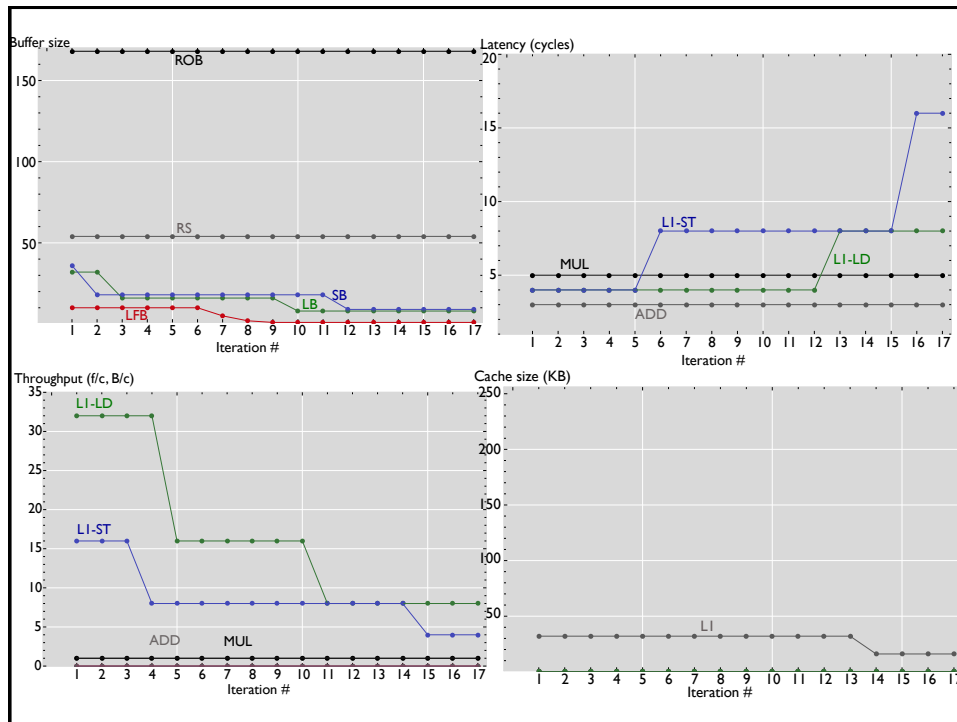# Adapting the Processor

FFT 1024



After adapting



# Evolution of Performance

**1** Program synthesis for performance
*Daniele Spampinato [CGO 2014]*

**2** Performance bottleneck modeling
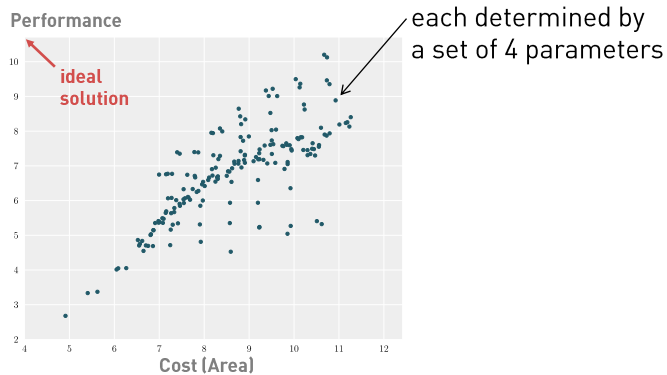*Victoria Caparròs Cabezas [IISWC 2014]*

**3** Predicting Pareto-optimal solutions
*Marcela Zuluaga, Andreas Krause [ICML 2013]*
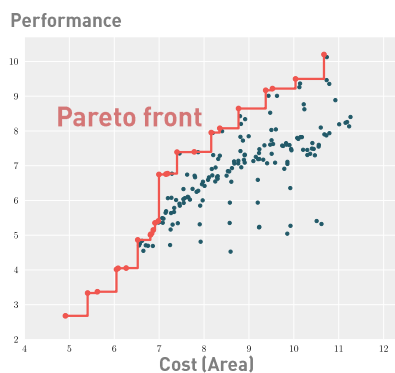
# Multi-Objective Optimization

Case Study: Different Hardware Implementations of a 256-input Sorter

Performance

each determined by
a set of 4 parameters

ideal
solution

10
9
8
7
6
5
4
3
2

4    5    6    7    8    9    10    11    12

Cost (Area)

**Which ones are relevant?**

---

# Multi-Objective Optimization

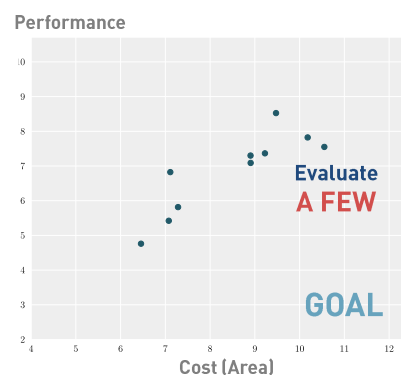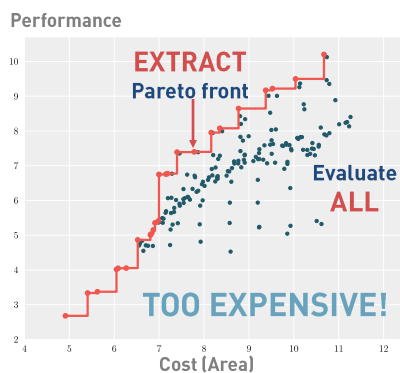Case Study: Different Hardware Implementation of a 256-input Sorter

Performance

10
9
8
7
6
5
4
3
2

Pareto front

4    5    6    7    8    9    10    11    12

Cost (Area)

**How to get it?**

# Multi-Objective Optimization

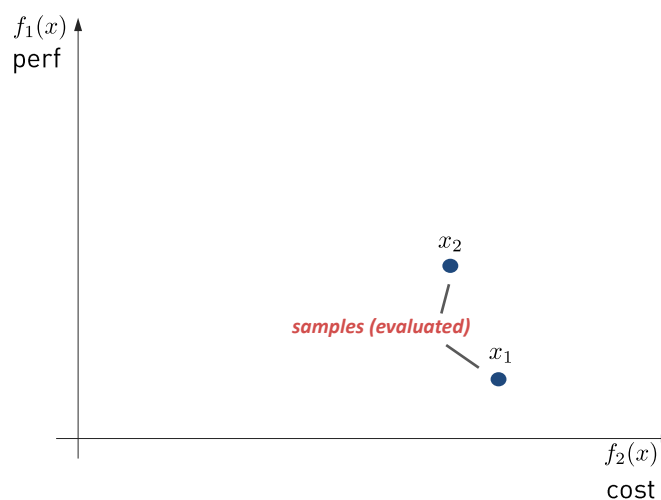Case Study: Different Hardware Implementation of a 256-input Sorter
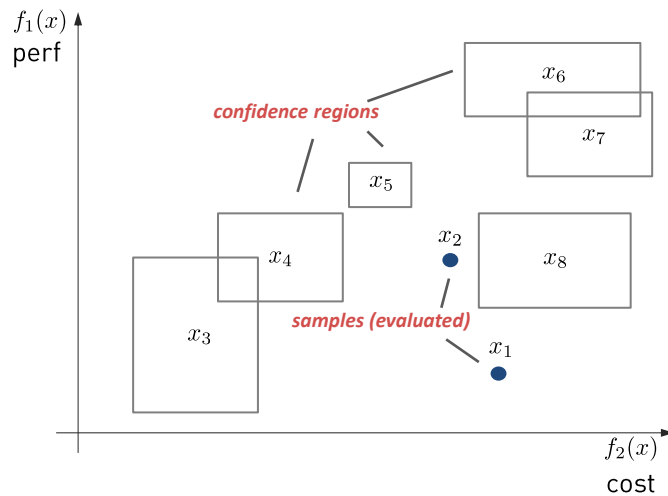


**Our Solution: Pareto Active Learning (PAL)**
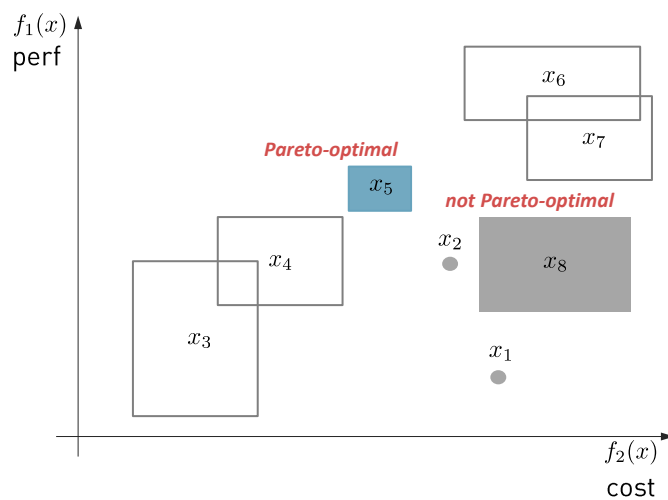
# Running PAL
*Modeling with Gaussian Processes*
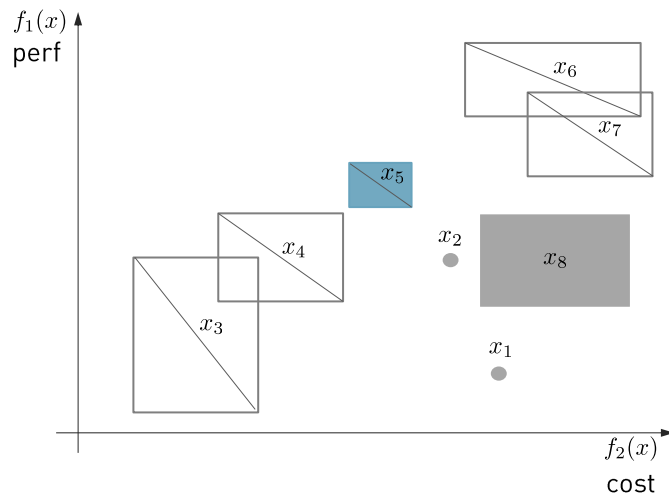
**Running PAL**
*Modeling with Gaussian Processes*

$f_1(x)$ perf — $f_2(x)$ cost

confidence regions

samples (evaluated)

$x_6$, $x_7$, $x_5$, $x_4$, $x_3$, $x_2$, $x_8$, $x_1$



**Running PAL**
*Classification*

$f_1(x)$ perf — $f_2(x)$ cost

Pareto-optimal

not Pareto-optimal

$x_6$, $x_7$, $x_5$, $x_4$, $x_3$, $x_2$, $x_8$, $x_1$

**Running PAL**

*Sampling*



**Running PAL**

*Sampling*

*Modeling, Classification, Sampling*

# Convergence Behaviour

*ParEGO: State-of-the-art evolutionary algorithm*



SNW ($|E| = 206$)    NoC ($|E| = 259$)    SW-LLVM ($|E| = 1023$)

- $\epsilon = 0.001\%$
- $\epsilon = 0.002\%$
- $\epsilon = 0.004\%$
- $\epsilon = 0.008\%$
- $\epsilon = 0.016\%$
- $\epsilon = 0.032\%$
- $\epsilon = 0.064\%$
- $\epsilon = 0.128\%$
- $\epsilon = 0.256\%$
- $\epsilon = 0.512\%$

---

## Program synthesis for performance    Bottleneck modelling    Predicting Pareto fronts



I am happy to discuss more
www.acl.inf.ethz.ch

# Experiments: Data sets

**4 features**      **4 features**      **11 features**



Marcela Zuluaga, Andreas Krause, Peter Milder, Markus Püschel. *Streaming Sorting Networks. DAC 2012*

Oscar Almer, Nigel Topham, Björn Franke. *A Learning- Based Approach to the Automated Design of MP-SoC Networks. ARCS 2011*

Predicting Performance via Automated Feature-Interaction Detection . *N. Siegmund, S. S. Kolesnikov, C. Kastner , S. Apel, D. Batory, M. Rosenmuller, and G. Saake. ICSI 2012*