

DF00100 Advanced Compiler Construction
 DF21500 Multicore Computing
 TDDD56 Multicore and GPU Programming

Autotuning
A short introduction

Christoph Kessler, IDA,
 Linköpings universitet



Motivation

- Modern (high-end) computer architectures are too complex
 - Some final machine parameters may not be statically (well-)known
 - Caches (multiple levels, capacity, associativity, replacement policy)
 - Memory latency
 - ILP and pipelining:
Dynamic dispatch, out-of-order execution, speculation, branching
 - Parallelism and contention for shared resources
 - OS scheduler
 - Paging
 - Perf. not well predictable e.g. for manual or compiler optimization
- Some program parameters (problem sizes, data locality etc.) may not be statically known
- Different algorithms / implementation variants may exist for a computation
- Hardcoded manual optimizations lead to non-performance-portable code
- Compiler optimizations are limited and may have unexpected side effects / interferences

2

Motivation (cont.)

→ Thousands of knobs that we could turn to tune performance!

- Which ones and how?
- Avoid hardcoded performance tuning



C. Kessler, IDA, Linköpings universitet.



Performance Portability for User-level code?

Avoid hard-coded adaptations / optimizations such as:

```
if (avail_num_threads() > 1)
    in_parallel {
        sort(a, n/2); // on first half of resources
        sort(&a[n/2], n-n/2); // on the other half
    }
else ... (do it in serial)
```

NO!

```
if (available(GPU))
    gpusort(a,n);
else
    qsort(a,n);
```

NO!

```
if (n < CACHESIZE/4)
    mergesort(a,n);
else
    quicksort(a,n);
```

NO!

C. Ke

Idea: Autotuning – Automatic optimization for unknown target system using Machine Learning

■ Given: Training data and initial program version
 → Observed performance on target
 → Machine learning algorithm
 → Optimization strategy (choice of some parameter(s))
 → Automatic code generation / adaptation for target platform
 and possibly repeat this process

■ for libraries: **autotuning library generators**,
 for compilers: **iterative compilation**
 for dynamic composition: **context-aware composition**

■ Typical examples:

- Find the best blocking factor(s) for loops or loop nests to automatically adapt to target cache behavior
- Find the right sequence and settings of compiler optimizations
- Select among different algorithms for same operation
- How many cores/threads / which processors/accelerators to use?

C. Kessler, IDA, Linköpings universitet.



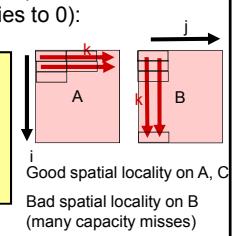
Recall: Tiled Matrix-Matrix Multiplication (1)

- Matrix-Matrix multiplication $C = A \times B$
 here for square ($n \times n$) matrices C, A, B , with n large ($\sim 10^3$):

$$\bullet C_{ij} = \sum_{k=1..n} A_{ik} B_{kj} \quad \text{for all } i, j = 1..n$$

- Standard algorithm for Matrix-Matrix multiplication
 (here without the initialization of C -entries to 0):

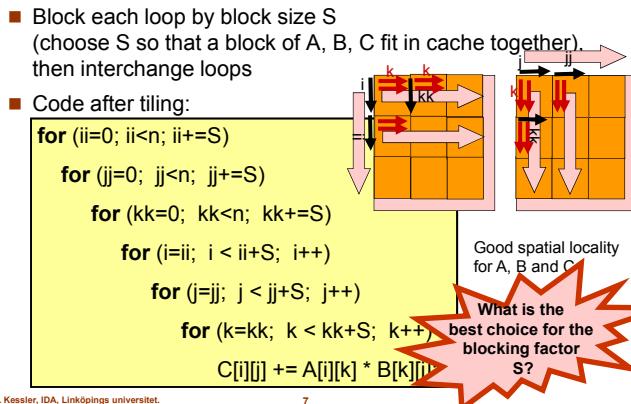
```
for (i=0; i<n; i++)
    for (j=0; j<n; j++)
        for (k=0; k<n; k++)
            C[i][j] += A[i][k] * B[k][j];
```



C. Kessler, IDA, Linköpings universitet.

6

Recall: Tiled Matrix-Matrix Multiplication (2)



Recall: Loop Unroll-And-Jam

unroll the outer loop
and fuse the resulting inner loops:

```

for i from 1 to N do
    for j from 1 to N do
        a[i] ← a[i] + b[j] unroll&jam:
    od
od
    
```

```

for i from 1 to N step 2 do
    for j from 1 to N do
        a[i] ← a[i] + b[j]
        a[i+1] ← a[i+1] + b[j]
    od
od
    
```

The same conditions as for loop interchange (for the two innermost loops after the unrolling step) must hold

(for a formal treatment see [Allen/Kennedy'02, Ch. 2.4.1])

- + increases reuse in inner loop
- + less overhead

What is the best choice for the unroll factor (here, 2)?

C. Kessler, IDA, Linköpings universitet. 8

Auto-tuning linear algebra library ATLAS (1)

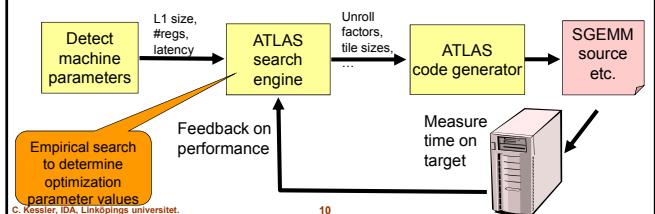
- BLAS = Basic Linear Algebra Subroutines
 - standard numerical library for Fortran, C
 - frequently used in high-performance applications
- Level-1 BLAS: Vector-vector operations e.g. dot product
- Level-2 BLAS: Matrix-vector operations
- Level-3 BLAS: Matrix-matrix operations, esp., generic versions of dense LU decomposition and Matrix mult.
 - SGEMM: $C := \alpha A * B + \beta C$
for matrices A,B,C, scalars α,β
 - is ordinary Matrix-Matrix multiplication for $\alpha=1, \beta=0$

C. Kessler, IDA, Linköpings universitet.

9

Auto-tuning linear algebra library ATLAS (2)

- ATLAS is a generator for optimized BLAS libraries
 - Tiling to address L1 cache
 - Unroll-and-jam / scalar replacement to exploit registers
 - Use multiply-accumulate and SIMD instructions where available
 - Schedule computation and memory accesses
- Outperforms vendor-specific BLAS implementations



Remark

- Off-line sampling and tuning by greedy heuristic search
 - Happens once for each new system at library deployment (generation) time
 - Can be expensive
- Not practical for less static scenarios or costly sampling
 - Fast predictors needed – full execution or even simulation is not feasible
 - Usually constructed by machine learning
 - Shortens the feedback loop
 - Could be adapted dynamically (on-line sampling/tuning)

C. Kessler, IDA, Linköpings universitet.

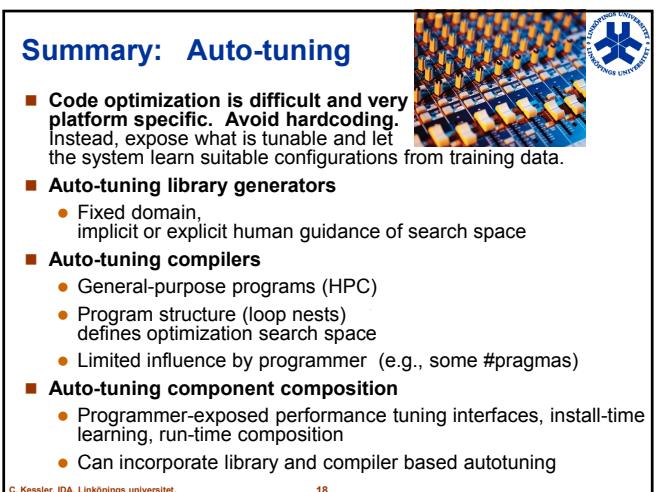
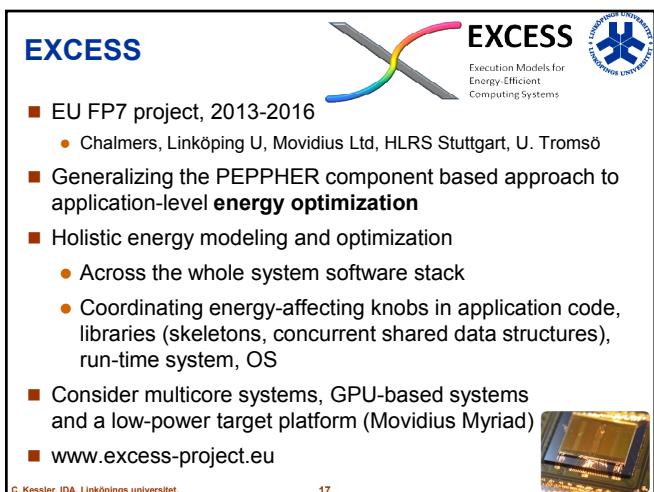
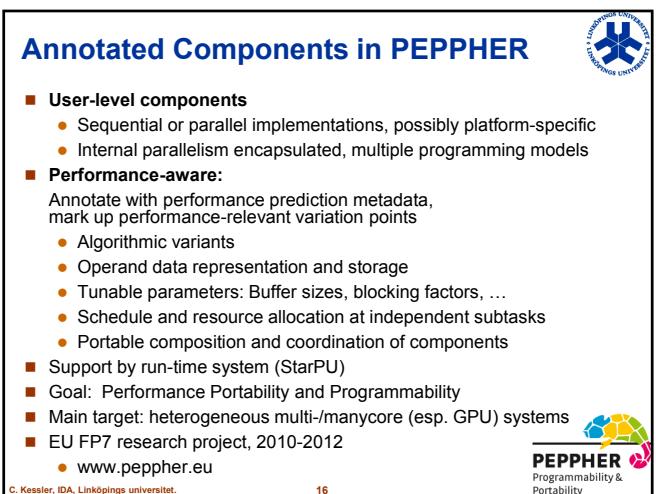
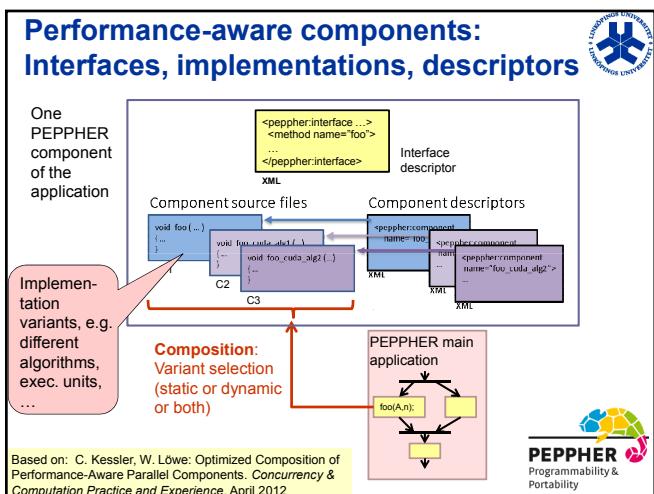
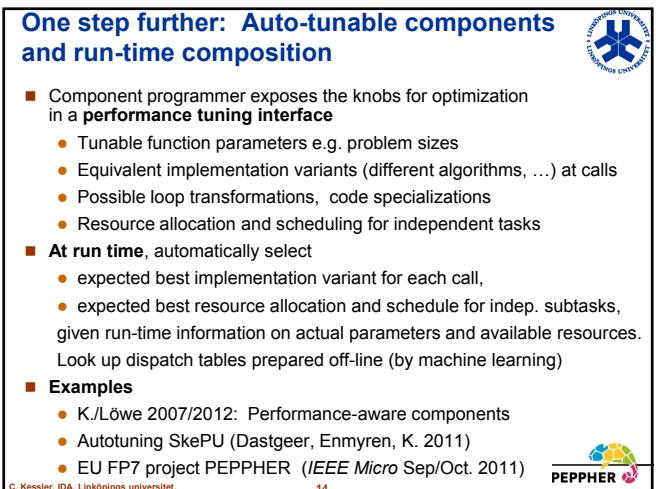
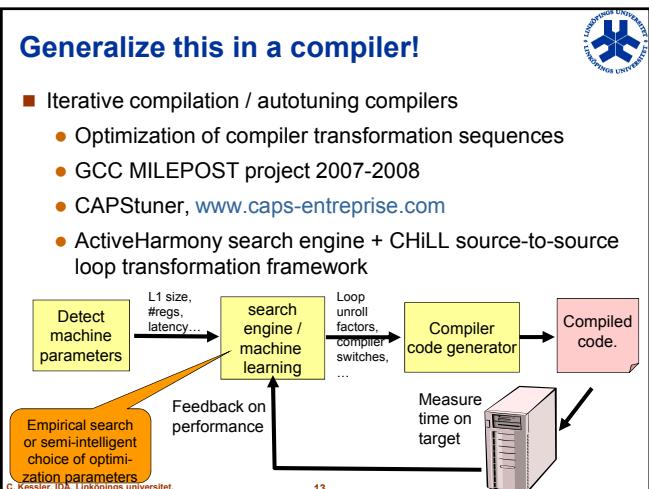
11

Further auto-tuning library generators

- Linear Algebra
 - ATLAS
 - PhiPAC
 - OSKI
- FFT and other signal processing
 - FFTW [Frigo'99]
 - SPIRAL [Püschel et al. 2005]
- Sorting, searching etc.
 - STAPL [Rauchwerger et al.]
 - [Li, Padua, Garzaran CGO'94]
 - [Brewer'95]
 - [Olszewski, Voss PDPTA-2004]

C. Kessler, IDA, Linköpings universitet.

12



Master thesis projects available!

- Infrastructure e.g. extending composition tool, containers, ...
- Case studies: PEPPERing applications, new platforms
- Improved static prediction support
- ...



C. Kessler, IDA, Linköpings universitet.

19

References

- **On ATLAS:**
J. Demmel, J. Dongarra, V. Eijkhout, E. Fuentes, A. Petitet, R. Vuduc, R. C. Whaley, K. Yelick: Self-adapting linear algebra algorithms and software. *Proceedings of the IEEE* 93(2):293-312, Feb. 2005
- **On FFTW:**
M. Frigo: A fast Fourier transform compiler. Proc. PLDI-1999, p.169-180, ACM.
- **On SPIRAL:**
M. Püschel et al.: SPIRAL: Code generation for DSP transforms. *Proceedings of the IEEE* 93(2):232-275, Feb. 2005
- **On iterative compilation:**
 - F. Bodin: Compilers in the Many-Core Era. *HiPEAC-2009*, http://www.caps-entreprise.com/Hipeac09_KeyNoteFBO.pdf
 - **On ActiveHarmony + CHILL:**
A. Tiwari, C. Chen, J. Chame, M. Hall, J. Hollingsworth: A scalable auto-tuning framework for compiler optimization. Proc. IPDPS-2009, pp. 1-12, IEEE.
- **On General Component Autotuning:**
 - C. Kessler, W. Löwe: A Framework for Performance-Aware Composition of Explicitly Parallel Components. Proc. ParCo-2007 conference, pp. 227-234, and later papers e.g.
 - C. Kessler, W. Löwe: Optimized Composition of Performance-Aware Parallel Components. *Concurrency & Computation Practice and Experience*, April 2012.
 - J. Ansel, C. Chan, Y. Wong, M. Olszewski, Q. Zhao, A. Edelman, S. Amarasinghe: PetaBricks: A Language and Compiler for Algorithmic Choice. Proc. PLDI-2009. ACM.
 - U. Dastgeer, L. Li, C. Kessler: The PEPPER Composition Tool. ... MuCoCoS'12, IEEE.

C. Kessler, IDA, Linköpings universitet.

20

