



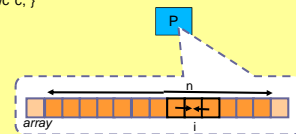
Parallelization Patterns and Algorithmic Skeletons

An Introduction

Example 1: 1D smoothing filter in C

```
...
float filter (float a, b, c) { return wa*a + wb*b + wc*c; }
void main ( int argc, char *argv[] )
{
    float *array = new_FloatArray( n+2 );
    float *tmp = new_FloatArray( n+2 );

    while ( globalerr > 0.1 ) {
        for (i=1; i<=n; i++)
            tmp[i] = filter( array[i-1], array[i], array[i+1] );
        globalerr = 0.0;
        for (i=1; i<=n; i++)
            globalerr = fmax( globalerr, fabs(array[i] - tmp[i]) );
        for (i=1; i<=n; i++)
            array[i] = tmp[i];
    }
}
```



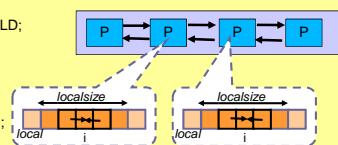
2

Example 1: 1D smoothing filter in C + MPI

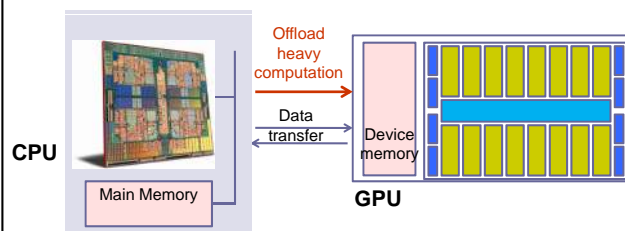
```
...
void main ( int argc, char *argv[] ) {
    MPI_Comm com = MPI_COMM_WORLD;
    MPI_Init ( &argc, &argv );
    MPI_Comm_size ( com, &np );
    MPI_Comm_rank ( com, &me );

    ...
    localsize = (int) ceil ( (float) n / np );
    local = new_FloatArray( localsize + 2 );

    ...
    while ( globalerr > 0.1 ) {
        if (me>0) MPI_Send ( local+1, 1, MPI_FLOAT, left_neighbor, 10, com );
        if (me<np-1) MPI_Send ( last, 1, MPI_FLOAT, right_neighbor, 20, com );
        for (i=1; i<=localsize; i++)
            tmp[i] = filter( local[i-1], local[i], local[i+1] );
        if (me<np-1) MPI_Recv ( tmp, 1, MPI_FLOAT, right_neighbor, 10, com, ... );
        if (me>0) MPI_Recv ( tmp+localsize+1, 1, MPI_FLOAT, left_neighbor, 20, com, ... );
        tmp[1] = filter( local[0], local[1], local[2] );
        tmp[localsize] = filter( local[localsize-1], local[localsize], local[localsize+1] );
        localerr = 0.0;
        for (i=1; i<=localsize; i++) localerr = fmax( localerr, fabs ( local[i]-tmp[i] ) );
        MPI_Allreduce ( &localerr, &globalerr, 1, MPI_FLOAT, MPI_MAX, com );
        for (i=1; i<=localsize; i++)
            local[i] = tmp[i];
    }
}
```

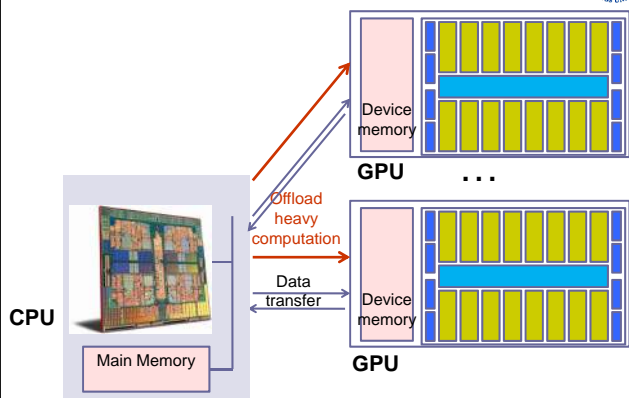


Example 2: GPU-Based Systems



4

Example 2: GPU-Based Systems

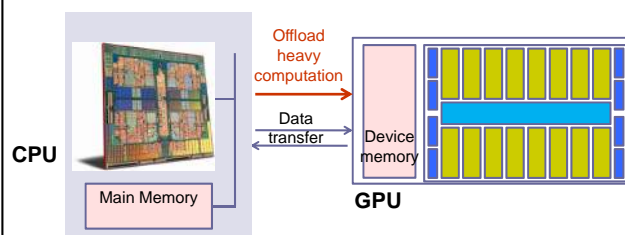


5

Example 2 Programming of GPU-Based Systems

...

- Portability ?
- Programmability ?
- Performance portability ?

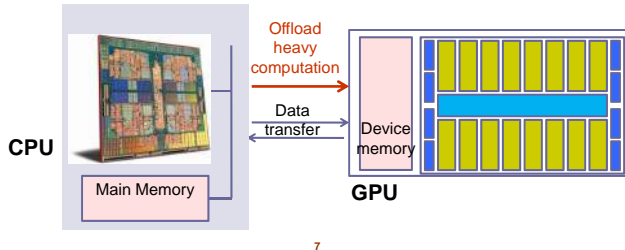


6

Example 2: Programming of GPU-based Systems ... with OpenCL™



- Portability ☺
- Programmability ☹ (low level)
- Performance portability ☹ (requires reoptimization)



7

Complexity of Parallel Algorithms and Programs



- Many different parallel programming models
 - Identify parallelism ("tasks")
 - Synchronization and communication?
 - Memory structure, -consistency model?
 - Resource allocation, mapping, scheduling?
 - Using accelerators, e.g. GPU? Or several GPUs?
- Error prone, hard to debug
- Code portability?
- Performance portability??

Can we make parallel programming as easy as sequential programming?

8

Observation



- Same characteristic form of parallelism, communication, synchronization re-applicable for all occurrences of the same specific structure of computation ((parallel) algorithmic paradigm, building block, pattern, ...)
- Elementwise operations on arrays
- Reductions
- Scan (Prefix-op)
- Divide-and-Conquer
- Farming independent tasks
- Pipelining
- ...
- Most of these have both sequential and parallel implementations

9

Example 1, cont.: 1D smoothing in C + Skeletons

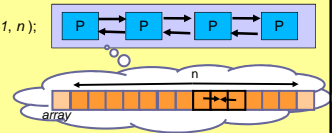


```

float filter (float a, b, c) { return wa*a + wb*b + wc*c; }

float elemError (float a, b) { return fabs (a - b); }

void main (int argc, char *argv[]) {
    ...
    DistrFloatArray *array = new_DistrFloatArray (n + 2);
    DistrFloatArray *tmp = new_DistrFloatArray (n + 2);
    DistrFloatArray *err = new_DistrFloatArray (n + 2);
    ...
    while (globalerr > 0.1) {
        map_with_overlap (filter, 1, tmp, array+1, n);
        map (elemError, err, array+1, tmp, n);
        reduce (fmax, &globalerr, err, n);
        map (copy, array+1, tmp, n);
    }
    ...
}
    
```



10

Data parallelism



Given:

- One or several data containers x with n elements, e.g. array(s) $x = (x_1, \dots, x_n)$, $z = (z_1, \dots, z_n)$, ...
- An operation f on individual elements of x , z , ... (e.g. *incr*, *sqr*, *mult*, ...)

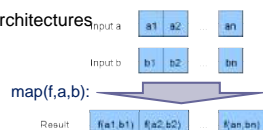
Compute: $y = f(x) = (f(x_1), \dots, f(x_n))$

Parallelizability: Each data element defines a task

- Fine grained parallelism
- Portionable, fits very well on all parallel architectures

Notation with higher-order function:

- $y = \text{map}(f, x)$



Variant: map with overlap: $y_i = f(x_{i-k}, \dots, x_{i+k})$, $i = 0, \dots, n-1$ for some small constant k

11

Data-parallel Reduction

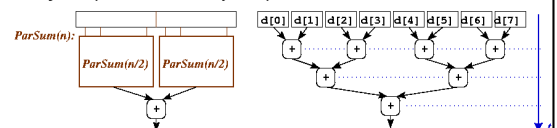


Given:

- A data container x with n elements, e.g. array $x = (x_1, \dots, x_n)$
- A binary associative operation op on individual elements of x (e.g. *add*, *max*, *bitwise-or*, ...)

Compute: $y = OP_{i=1 \dots n} x = x_1 op x_2 op \dots op x_n$

Parallelizability: Exploit associativity of op



Notation with higher-order function:

- $y = \text{reduce}(op, x)$

12

Task farming

Independent computations f_1, f_2, \dots, f_m could be done in parallel and/or in arbitrary order, e.g.

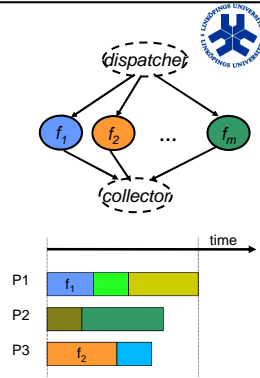
- independent loop iterations
- independent function calls

Scheduling problem

- n tasks onto p processors
- static or dynamic
- Load balancing

Notation with higher-order function:

- $(y_1, \dots, y_m) = \text{farm} (f_1, \dots, f_m) (x_1, \dots, x_n)$



13

Parallel Divide-and-Conquer

(Sequential) Divide-and-conquer:

- **Divide:** Decompose problem instance P in one or several smaller independent instances of the same problem, P_1, \dots, P_k
- For all i : If P_i *trivial*, solve it *directly*.
- Else, solve P_i by recursion.
- **Combine** the solutions of the P_i into an overall solution for P

Parallel Divide-and-Conquer:

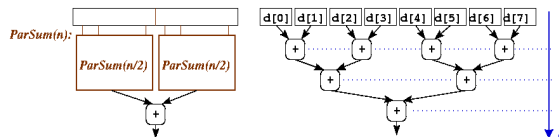
- Recursive calls can be done in parallel.
- Parallelize, if possible, also the divide and combine phase.
- Switch to sequential divide-and-conquer when enough parallel tasks have been created.

Notation with higher-order function:

- $\text{solution} = \text{DC} (\text{divide}, \text{combine}, \text{is trivial}, \text{solved directly}, n, P)$

14

Example: Parallel Divide-and-Conquer



Example: Parallel Sum over integer-array x

Exploit associativity:

$$\text{Sum}(x_1, \dots, x_n) = \text{Sum}(x_1, \dots, x_{n/2}) + \text{Sum}(x_{n/2+1}, \dots, x_n)$$

Divide: trivial, split array x in place

Combine is just an addition.

$$y = \text{DC} (\text{split}, \text{add}, \text{isSmall}, \text{addFewInSeq}, n, x)$$

Data parallel reductions are an important special case of DC.

15

Example: Parallel Divide-and-Conquer (2)

Example: Parallel QuickSort over a float-array x

Divide: Partition the array (elements \leq pivot, elements $>$ pivot)

Combine: trivial, concatenate sorted sub-arrays

$$\text{sorted} = \text{DC} (\text{partition}, \text{concatenate}, \text{isSmall}, \text{qsort}, n, x)$$

16

Pipelining

applies a sequence of dependent computations (f_1, f_2, \dots, f_k) elementwise to data sequence $x = (x_1, \dots, x_n)$

- For fixed x_j , compute $f_i(x_j)$ before $f_{i+1}(x_j)$
- Computations of f_i on different x_j are independent.

Parallelizability: Overlap execution of all f_i for k subsequent x_j

- $\text{time}=1$: compute $f_1(x_1)$
- $\text{time}=2$: compute $f_1(x_2)$ and $f_2(x_1)$
- $\text{time}=3$: compute $f_1(x_3)$ and $f_2(x_2)$ and $f_3(x_1)$
- ...
- Total time: $O((n+k) \max_i(\text{time}(f_i)))$ with k processors

Notation with higher-order function:

- $(y_1, \dots, y_n) = \text{pipe} ((f_1, \dots, f_k), (x_1, \dots, x_n))$

17

Skeletons

Skeletons are reusable, parameterizable components with well defined semantics for which efficient parallel implementations may be available.

Inspired by higher-order functions in functional programming

One or very few skeletons per parallel algorithmic paradigm

- map, farm, DC, reduce, pipe, scan ...

Parameterised in user code

- Customization e.g. by instantiating a skeleton template in a user function

Composition of skeleton instances in program code normally by sequencing+data flow

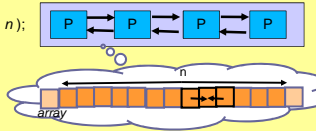
- e.g. $\text{squaresum}(x)$ can be defined by


```
{
  tmp = map( sqr, x );
  return reduce( add, tmp );
}
```

18

Example 1 revisited: 1D smoothing in C + Skeletons

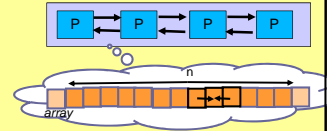
```
...
float filter (float a, b, c) { return wa*a + wb*b + wc*c; }
float elemError (float a, b) { return fabs (a - b); }
...
void main (int argc, char *argv[]) {
    ...
    DistrFloatArray *array = new_DistrFloatArray (n + 2);
    DistrFloatArray *tmp = new_DistrFloatArray (n + 2);
    DistrFloatArray *err = new_DistrFloatArray (n + 2);
    ...
    while (globalerr > 0.1) {
        map_with_overlap( filter, 1, tmp, array+1, n);
        map( elemError, err, array+1, tmp, n);
        reduce( fmax, &globalerr, err, n);
        map( copy, array+1, tmp, n);
    }
    ...
}
```



19

Example 1 revisited: 1D smoothing in C++ + Skeletons

```
...
float filter (float a, b, c) { return wa*a + wb*b + wc*c; }
float elemError (float a, b) { return fabs (a - b); }
...
void main (int argc, char *argv[]) {
    ...
    Vector<float> array = new Vector<float> (n + 2);
    Vector<float> tmp = new Vector<float> (n + 2);
    Vector<float> err = new Vector<float> (n + 2);
    ...
    while (globalerr > 0.1) {
        map_with_overlap( filter, 1, tmp, array);
        map( elemError, err, array, tmp);
        reduce( fmax, &globalerr, err);
        map( copy, array, tmp);
    }
    ...
}
```



20

Generic containers (e.g., Vector<>) provide a cleaner interface and encapsulate metadata (e.g., size) and internal state of operand data (e.g., storage format and location, distribution, modified, ...)

Skeletons (cont.)

Skeletons encapsulate completely all coordination of parallelism and platform-specific issues

Threads/Process creation/termination, communication, synchronization

→ Code portability

→ Reuse of the coordination code across multiple skeleton instances

Skeletons may also have a sequential implementation

Uniform treatment of sequential and parallel programming

21

Skeleton Programming Systems

4 basic approaches for realizing skeletons (esp., parameterisation mechanism):

- Library of higher-order functions (functional or imperative)
- OO class library (subclass and define abstract parameter method(s))
- New language constructs (intrinsic / compiler-known functions)
- Generative programming, Static metaprogramming (Macros / templates)

Many research prototypes, e.g.:

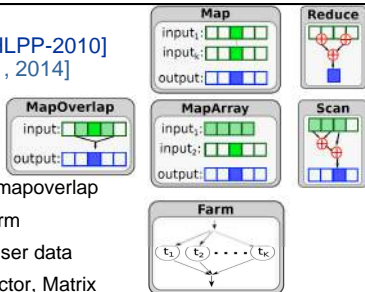
- P3L - C + skeletons
- SCL, Eden, HDC - functional
- eSkel - C + MPI
- Lithium - Java + RMI
- BlockLib - C + macros (generative) + DMA for Cell BE
- muskel, ASSIST - C++, grid computing
- MueSLi, QUAFF - C++ based, MPI
- SkePU, SkelCL - C++ based, for GPU based systems

Domain-specific Skeleton Systems, e.g.

- MallBa (combinatorial optimization: BB, DP, GA, ...)
- MapReduce (distributed data mining, Google)

SkePU [Enmyren, K. HLPP-2010] [Dastgeer 2011, 2014]

- C++ template library
- 6 dataparallel skeletons
 - Map, reduce, scan, mapreduce, maparray, mapoverlap
- 1 task-parallel skeleton: farm
- STL-based containers for user data
 - Smart containers for Vector, Matrix
- Generation of platform-specific variants for user functions
- Multiple back-ends: C, OpenMP, OpenCL, CUDA, StarPU
 - Hybrid CPU-GPU execution (with StarPU backend)
- Multi-GPU support
- Low overhead
- Tunable



23

The SkePU Multi-Backend Skeleton Library

Example: Dot Product

```
#include <iostream>

#include "skepu/vector.h"
#include "skepu/mapreduce.h"

BINARY_FUNC(plus, double, a, b,
return a+b;
)

BINARY_FUNC(mult, double, a, b,
return a*b;
)

....

Generate a skeleton instance (function dotProduct)

Generic container holds operand data

int main()
{
    skepu::MapReduce<mult, plus>
    dotProduct(new mult, new plus);

    skepu::Vector<double> v0(1000,2);
    skepu::Vector<double> v1(1000,2);

    double r = dotProduct(v0,v1);

    std::cout<<"Result: " <<r <<"\n";

    return 0;
}
```

Macro expands to platform-specific code versions for user functions

Invocation at run time

24

SkePU: Generating platform-specific variants of user functions

```
BINARY_FUNC(plus, double, a, b,
return a+b;
}
```



expands to

Macro expands to a struct with platform-specific versions of the user function

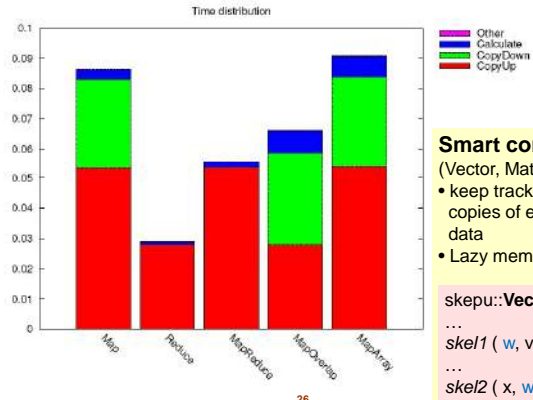
```
struct plus
{
    skepu::FuncType funcType;
    std::string func_CL;
    std::string funcName_CL;
    std::string datatype_CL;
    plus()
    {
        funcType = skepu::BINARY;
        funcName_CL.append("plus");
        datatype_CL.append("double");
        func_CL.append(
            "double plus(double a, double b)\n"
            "{\n"
            "    return a+b;\n"
            "}\n");
    }
    double CPU(double a, double b)
    {
        return a+b;
    }
    _device_ double CU(double a, double b)
    {
        return a+b;
    }
};
```

OpenCL

CPU

CUDA

Importance of Data Transfer Optimization



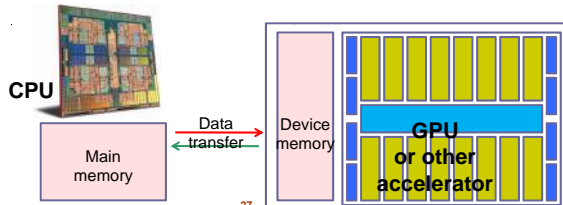
Smart containers
(Vector, Matrix):

- keep track of current copies of element data
- Lazy memory copying

```
skepu::Vector v, w, x;
...
skel1 ( w, v );
...
skel2 ( x, w );
```

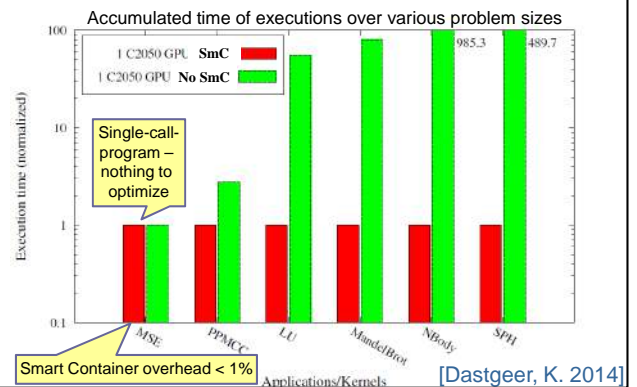
SkePU Containers

- For non-scalar operand data
 - STL-like Vector
 - Matrix
- Handle transparent data access and memory transfers between host and device ("**smart containers**")
 - Lazy memory copying
 - Implements (coarse-grained) memory coherence in software



Evaluation

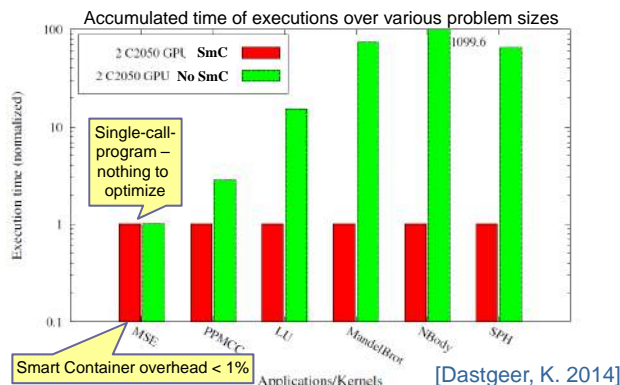
Speedup by (new) Smart Containers – 1 GPU



[Dastgeer, K. 2014]

Evaluation

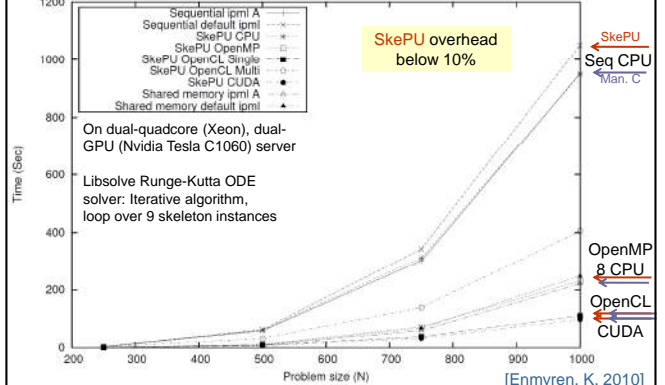
Speedup by (new) Smart Containers – 2 GPUs



[Dastgeer, K. 2014]

Comparison: SkePU vs hand-written code

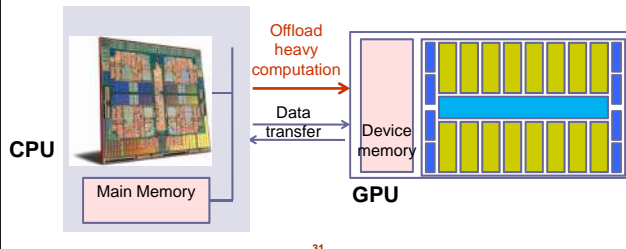
ODE RK-Solver from libsolve



[Enmyren, K. 2010]

Programming of GPU-based Systems ... with SkePU, so far...

- Portability ☺
- Programmability ☺
- Performance portability ☹



31

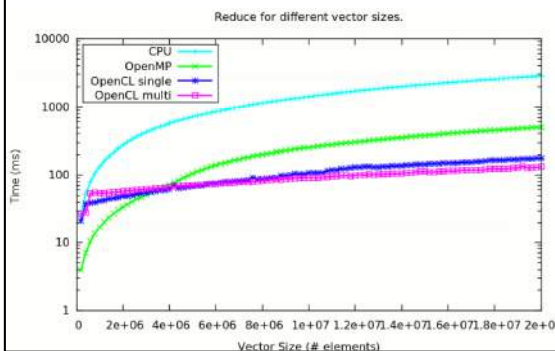
Implementation selection

- For each skeleton call, we need to decide:
 - Which skeleton implementation to use
- The decision depends upon:
 - Skeleton type
 - User function / actual computation
 - Target architecture
 - Problem size
 - etc.

32

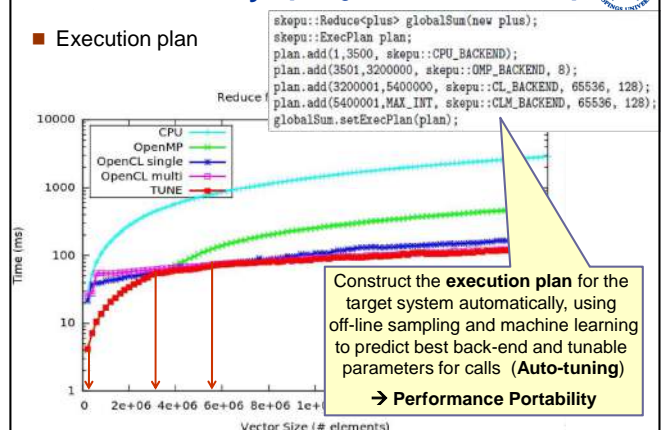
SkePU Tunability [Dastgeer et al. IWMSE-2011]

- Different Back-Ends, each with tunable parameters



SkePU Tunability [Dastgeer et al. IWMSE-2011]

- Execution plan



Auto-Tuning SkePU

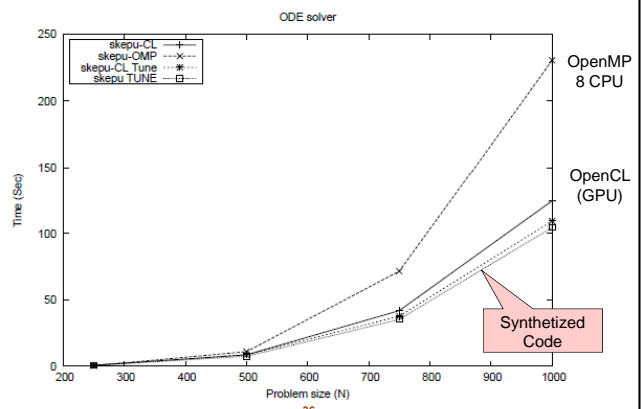
- Use off-line, machine learning approach to predict best back-end and its tunable parameters for given problem size
- Basic tunable parameters for each skeleton:
 - #threads (on CPU), thread block size + grid size (on GPU)

```
1 --- 50000 OMP_BACKEND 8
50001 --- 150000 CL_BACKEND 32 16384
150001 --- 225000 CL_BACKEND 128 2048
225001 --- 1050000 CL_BACKEND 512 2048
...
71500001 --- INFINITY CLM_BACKEND 256 2048
```

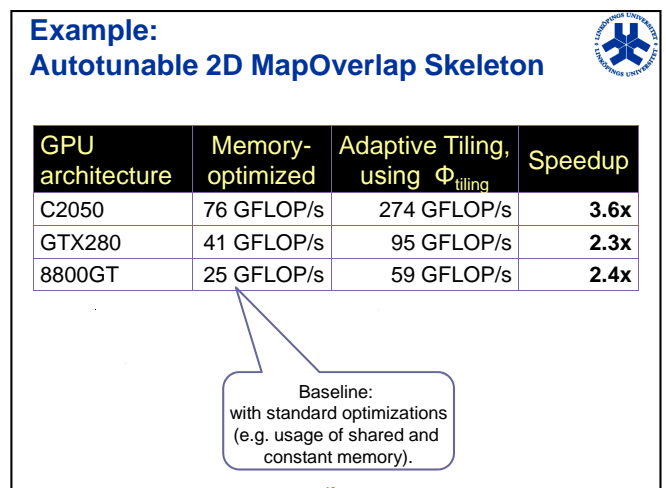
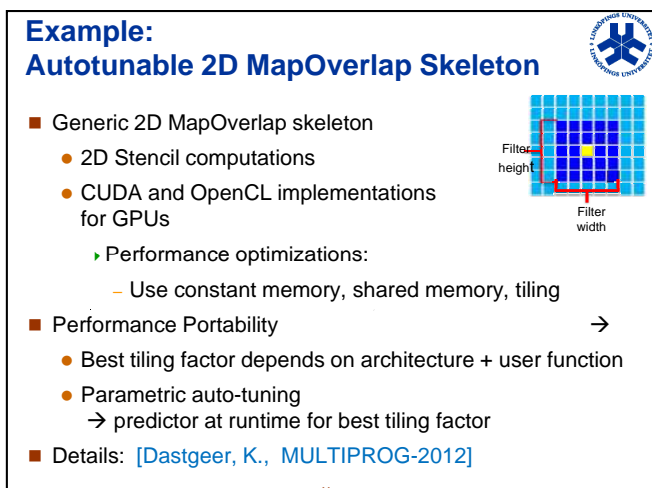
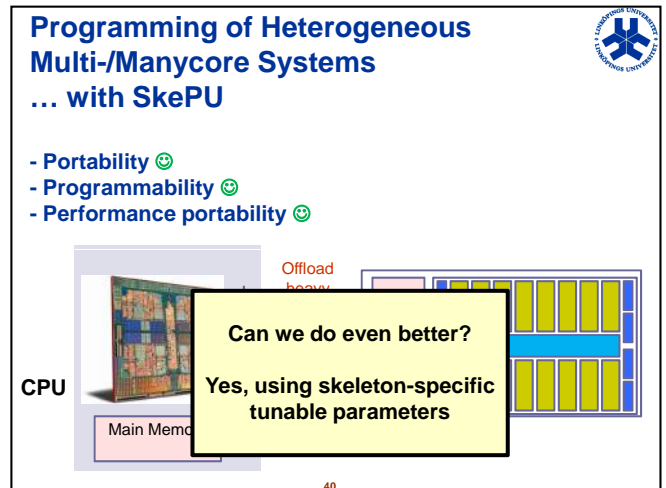
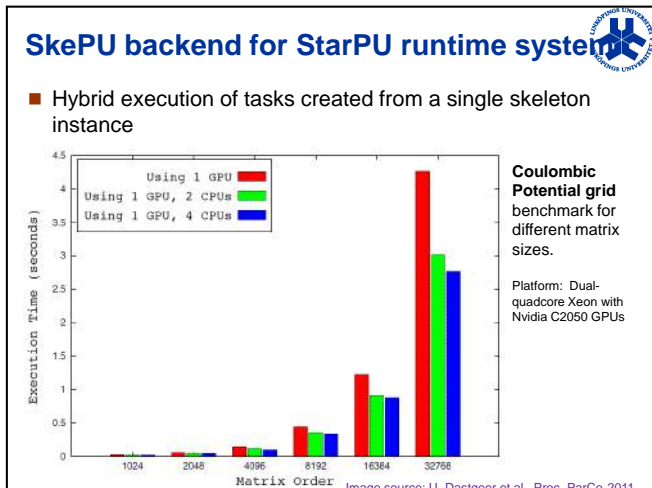
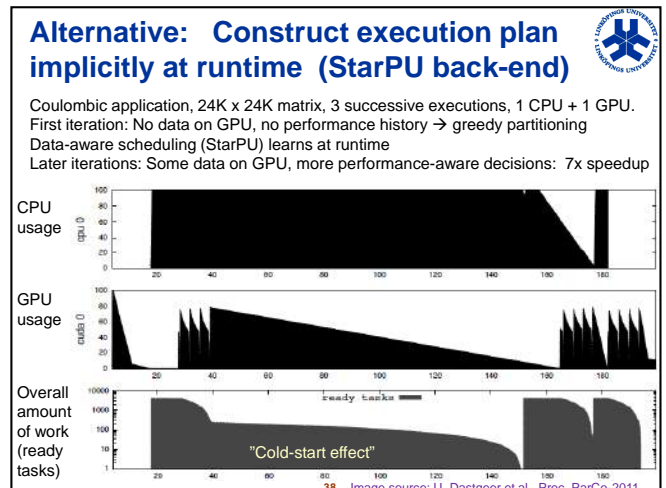
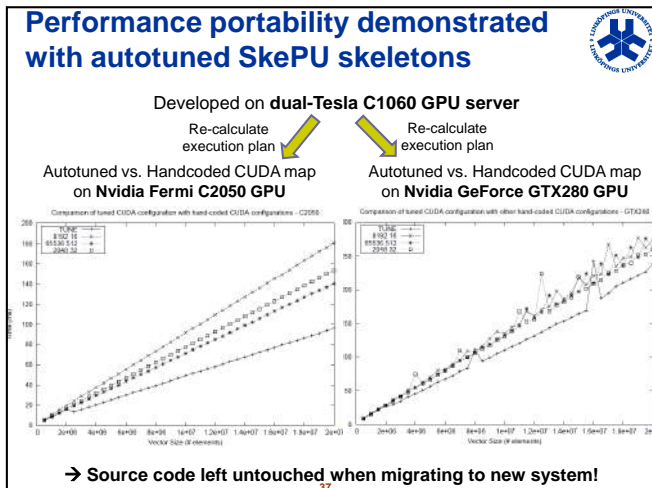
- Performance portability demonstrated across multiple GPU systems [Dastgeer, K. IWMSE'11, APPT'13]

35

RK-ODE-solver, synthesized from tuned SkePU skeletons



36



Conclusion (SkePU)



- SkePU skeletons are pre-defined generic components
 - For frequently occurring algorithmic patterns
 - map, reduce, scan, mapoverlap, farm ...
- Multiple back-ends, multi-GPU support
 - Seq, OpenMP, OpenCL, CUDA
 - StarPU backend for task parallelism and hybrid parallel execution
 - MPI back-end for GPU clusters
- Smart containers to avoid unnecessary data transfers
- Auto-tunable
 - Off-line + on-line tuning of resource allocation for calls
 - Parametric autotuning for specific skeletons



Other Skeleton Programming Frameworks

Example: Intel TBB Algorithm Templates

Intel Threading Building Blocks (TBB)



- Library for programming multicore processors
- extends C++ with a task based parallel programming model including
 - tasks (also fine-grained), no threads
 - high-level parallel *algorithm templates* (functions - nestable),
 - ▶ data-parallel (e.g. parallel for, reduce, scan) and task-parallel (e.g. pipe)
 - ▶ User functions (body) more coarse-grained than individual elements, to better perform on CPU
 - concurrent containers,
 - mutexes, atomic operations, etc.
- sophisticated run-time task scheduling mechanism,
 - At runtime, the TBB run-time system creates tasks and assigns them to threads which the OS schedules to cores
 - Dynamic load balancing by *task stealing*.

45

TBB example

Intel: TBB tutorial, 2010,
www.intel.com and
threadingbuildingblocks.org



- Data-parallel loop in TBB

```
#include "tbb/tbb.h"
void ParallelSquare( float a[], size_t n )
{
    parallel_for( blocked_range<size_t>(0,n),
                  Square(a));
}
```

- Class Square defines a **functor** (= function object, instance of a class containing a member function that overloads the (.) operator → invoking "()" looks like a function call)
 - "element function", "user function":
Works on a contiguous subarray (index subrange) at a time

46

Summary



- **Skeleton programming**
 - Algorithmic paradigms
 - Predefined generic parallel components, parameterized in user code
 - Hiding complexity (parallelism and low-level programming)
 - ☺ Abstraction
 - ☹ Enforces structuring
 - ☺ Parallelization for free
 - ☺ Easier to analyze and transform
 - ☹ Requires complete understanding and rewriting
 - ☹ Available skeleton set does not always fit
 - ☹ May lose some efficiency compared to manual parallelization
- Industry (beyond HPC domain) has adopted skeletons
 - map, reduce, scan in many modern parallel programming APIs
 - ▶ e.g., Intel Threading Building Blocks (TBB): par. for, par. reduce, pipe
 - ▶ NVIDIA Thrust
 - Google MapReduce (for distributed data mining applications)

Thesis projects available!



- Extension of SkePU
 - New skeletons, new containers, new platforms...
- Porting applications to SkePU
 - Medical image visualization
 - Linear equation system solvers
 - ...
- Automated tuning for new optimization objectives (energy)
- Improvements of the autotuning framework
- Static selection for multiple calls to consider inter-call effects
- SkePU is an open-source project
 - Documentation + download: www.ida.liu.se/~chrke/skepu

48



Questions?

Some literature on skeleton programming



- M. Cole: *Algorithmic Skeletons: Structured Management of Parallel Computation*, MIT Press & Pitman, 1989. <http://homepages.inf.ed.ac.uk/mic/Pubs/pubs.html>
- S. Pelagatti: *Structured Development of Parallel Programs*. Taylor and Francis, 1998.
- F. Rabhi and S. Gorlatch (eds.): *Patterns and Skeletons for Parallel and Distributed Computing*. Springer-Verlag, 2003.
- M. Ålind, M. Eriksson, C. Kessler: BlockLib: A Skeleton Library for Cell Broadband Engine. Proc. ACM Int. Worksh. on Multicore Software Engineering, Leipzig, 2008.
- J. Enmyren, C. Kessler: SkePU: A Multi-Backend Skeleton Programming Library for Multi-GPU Systems. Proc. HLPP-2010 Int. Workshop on High-Level Parallel Programming, Sep. 2010, Baltimore, USA. ACM.
- U. Dastgeer: Performance-Aware Component Composition for GPU-Based Systems. PhD thesis, Linköping University, 2014. Chapter 3.
- SkePU Documentation and Download: <http://www.ida.liu.se/~chrke/skepu>



Acknowledgment: SkePU research and development is funded by EU FP7 projects PEPPHER (2010-12), EXCESS (2013-16); and by SeRC (www.e-science.se), 2011-2014

Glossary



■ Performance Portability

... is the ability of a program to automatically adapt to a new execution platform to achieve an automated best-effort optimization of performance on the new target system, without manual rewriting / reoptimization.

■ [Algorithmic] Skeleton

... is a pre-defined, generic software construct for high-level programming that implements a specific *pattern* of control and data flow, that can be *parameterized* by problem-specific code to instantiate a problem-specific function, and whose implementation internally *encapsulates* all platform-specific details such as parallelism, heterogeneity, communication and synchronization.

51