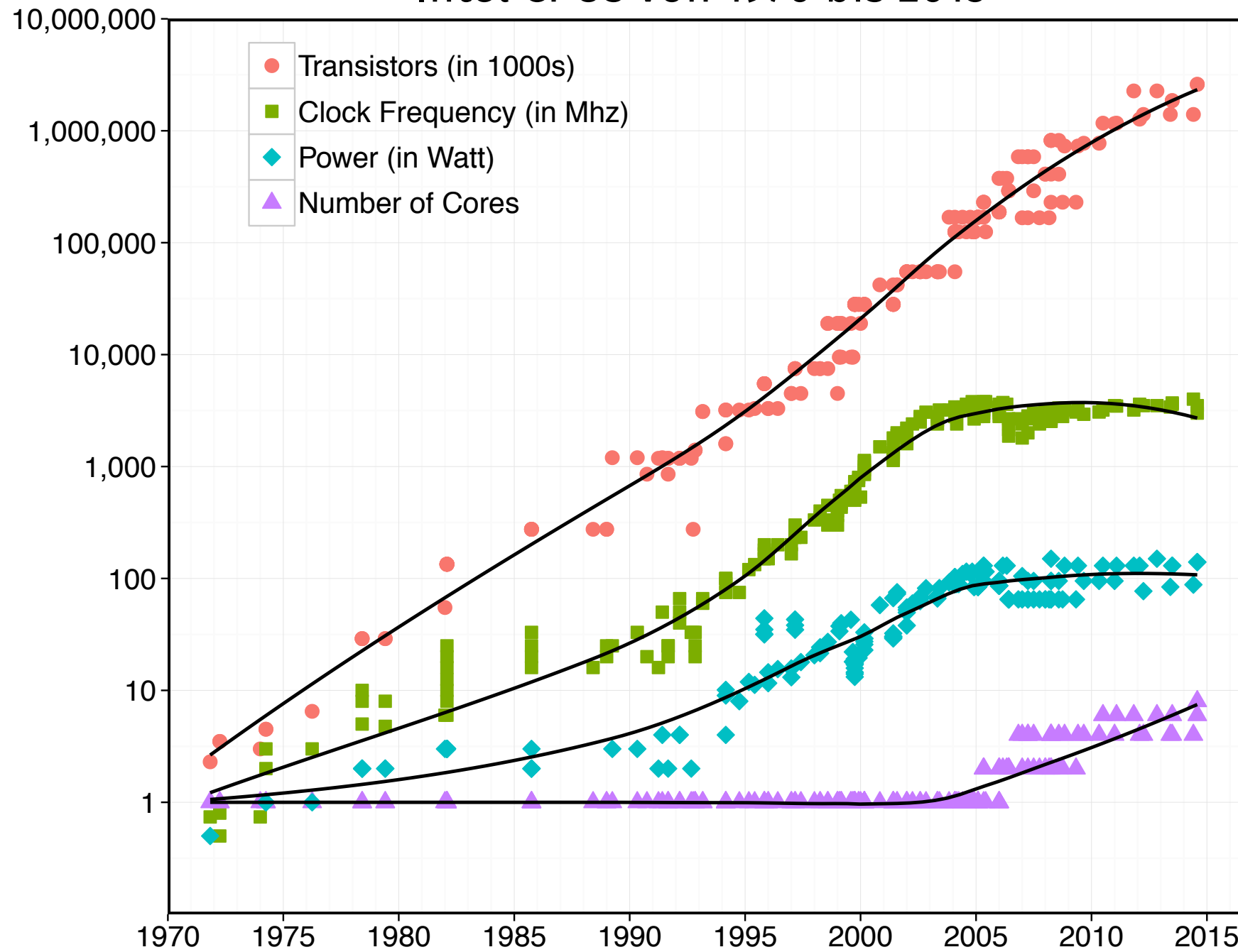


IMPROVING PROGRAMMABILITY
AND PERFORMANCE PORTABILITY
ON MANY-CORE PROCESSORS

MICHEL STEUWER

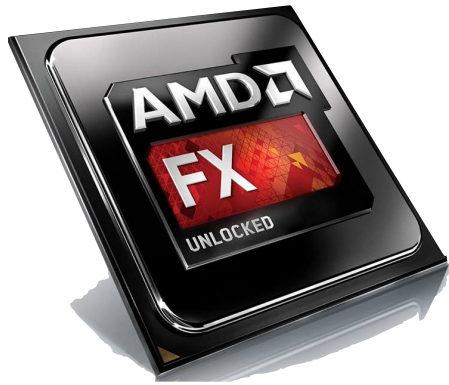
Die Manycore Ära

Intel CPUs von 1970 bis 2015



Inspiriert von Herb Sutter "The Free Lunch is Over:
A Fundamental Turn Towards
Concurrency in Software"

Die Manycore Ära



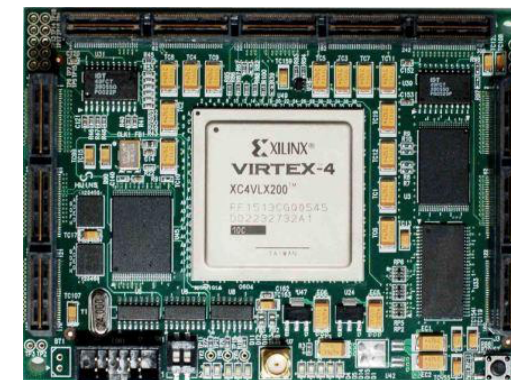
multicore CPUs



GPUs



Beschleuniger



FPGAs

Agenda

Meine Dissertation adressiert zwei zentrale Herausforderungen:

- I. Die *Herausforderung der Programmierbarkeit*
- II. Die *Herausforderung der Performance-Portabilität*

TEIL I

Die Herausforderung der Programmierbarkeit

Programmierung mit OpenCL

- Beispiel: Parallele Summation eines Arrays in OpenCL

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Programmierung mit OpenCL

- Beispiel: Parallele Summation eines Arrays in OpenCL

Kernel Funktion wird parallel von vielen work-items ausgeführt

```
kernel void reduce(global float* g_idata, global float* g_odata,
                   unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n ? g_idata[i] : 0);
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Work-items werden durch eine globale id identifiziert

Programmierung mit OpenCL

- Beispiel: Parallele Summation eines Arrays in OpenCL

Work-items werden zu work-groups zusammengefasst

Lokale id innerhalb einer work-group

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```


Programmierung mit OpenCL

- Beispiel: Parallele Summation eines Arrays in OpenCL

Großer, aber langsamer globaler Speicher

Kleiner, aber schneller lokaler Speicher

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i   = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
            barrier(CLK_LOCAL_MEM_FENCE);
        }
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Barrieren für Speicherkonsistenz

Programmierung mit OpenCL

- Beispiel: Parallele Summation eines Arrays in OpenCL

```
kernel void reduce(global float* g_idata, global float* g_odata,
                  unsigned int n, local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i    = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);
    // do reduction in local memory
    for (unsigned int s=1; s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    // write result for this work-group to global memory
    if (tid == 0) g_odata[get_group_id(0)] = l_data[0];
}
```

Funktional korrekte Implementierungen in OpenCL sind schwierig!

DAS SKELCL PROGRAMMIERMODEL

Das SkelCL Programmiermodell

Drei Abstraktionen zu OpenCL hinzugefügt:

- **Parallele Datencontainer**

für eine einheitliche Speicherverwaltung zwischen CPU und (mehreren) GPUs

- implizite Speichertransfers zwischen CPU und GPU
- *lazy copying* minimiert den Datentransfer

- **Wiederkehrende Muster paralleler Programmierung (Algorithmische Skelette)**

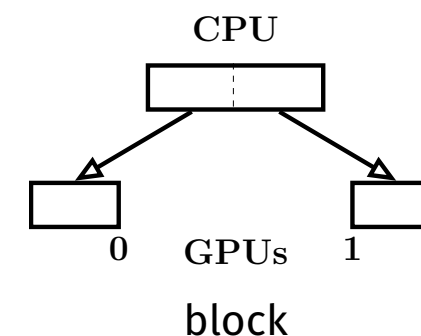
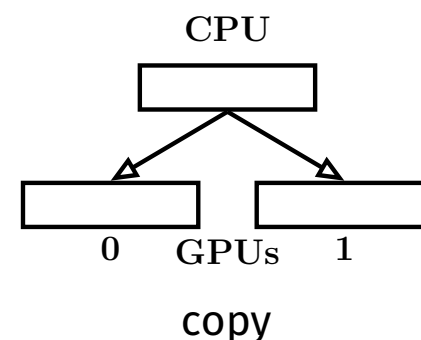
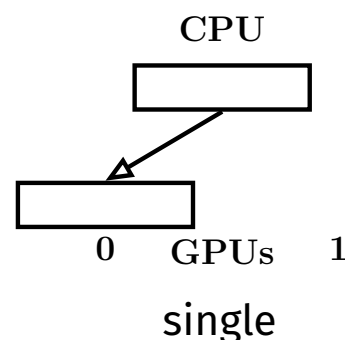
für eine vereinfachte Beschreibung paralleler Berechnungen

$$\text{zip } (\oplus) [x_1, \dots, x_n] [y_1, \dots, y_n] = [x_1 \oplus y_1, \dots, x_n \oplus y_n]$$

$$\text{reduce } (\oplus) \oplus_{\text{id}} [x_1, \dots, x_n] = \oplus_{\text{id}} \oplus x_1 \oplus \dots \oplus x_n$$

- **Daten Verteilungen**

für einen transparenten Datentransfer in Systemen mit mehreren GPUs.



Die SkelCL Softwarebibliothek am Beispiel

$\text{dotProduct } A \ B = \text{reduce } (+) \ 0 \ (\text{zip } (\times) \ A \ B)$

```
#include <SkelCL/SkelCL.h>
#include <SkelCL/Zip.h>
#include <SkelCL/Reduce.h>
#include <SkelCL/Vector.h>

float dotProduct(const float* a, const float* b, int n) {
    using namespace skelcl;
    skelcl::init( 1_device.type(deviceType::ANY) );

    auto mult = zip([](float x, float y) { return x*y; });
    auto sum = reduce([](float x, float y) { return x+y; }, 0);

    Vector<float> A(a, a+n); Vector<float> B(b, b+n);

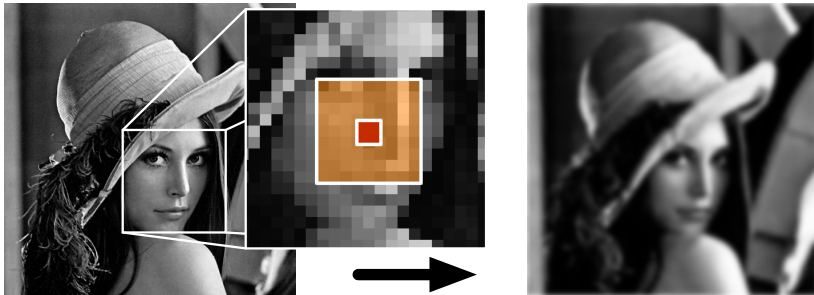
    Vector<float> C = sum( mult(A, B) );

    return C.front();
}
```

Neue Algorithmische Skelette

Stencil Berechnungen

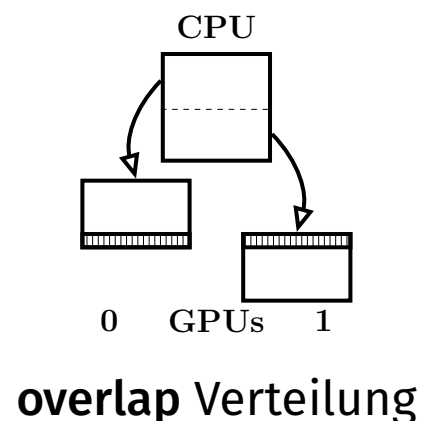
Beispiel: Gaußscher Weichzeichner



$\text{gauss } M = \text{stencil } f \ 1 \ \bar{0} \ M$

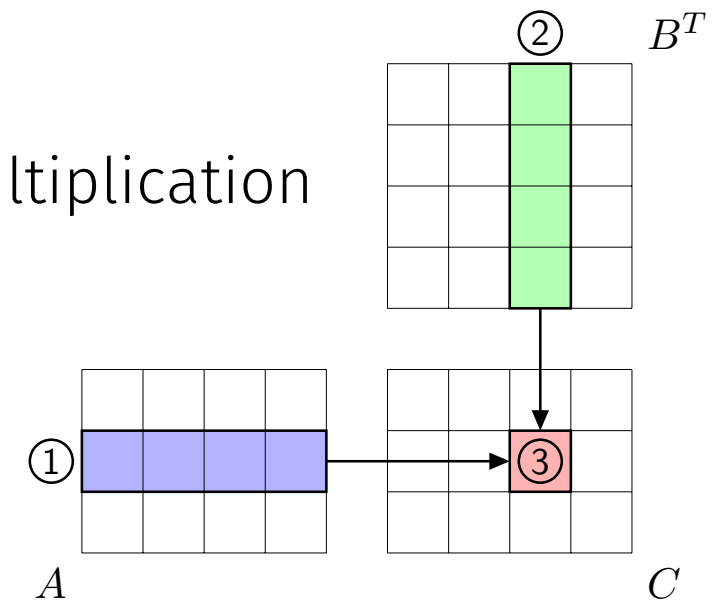
wo f die Funktion ist welche den Gaußschen Weichzeichner beschreibt

Unterstützung für mehrere GPUs:



Allpairs Berechnungen

Example:
Matrix Multiplication



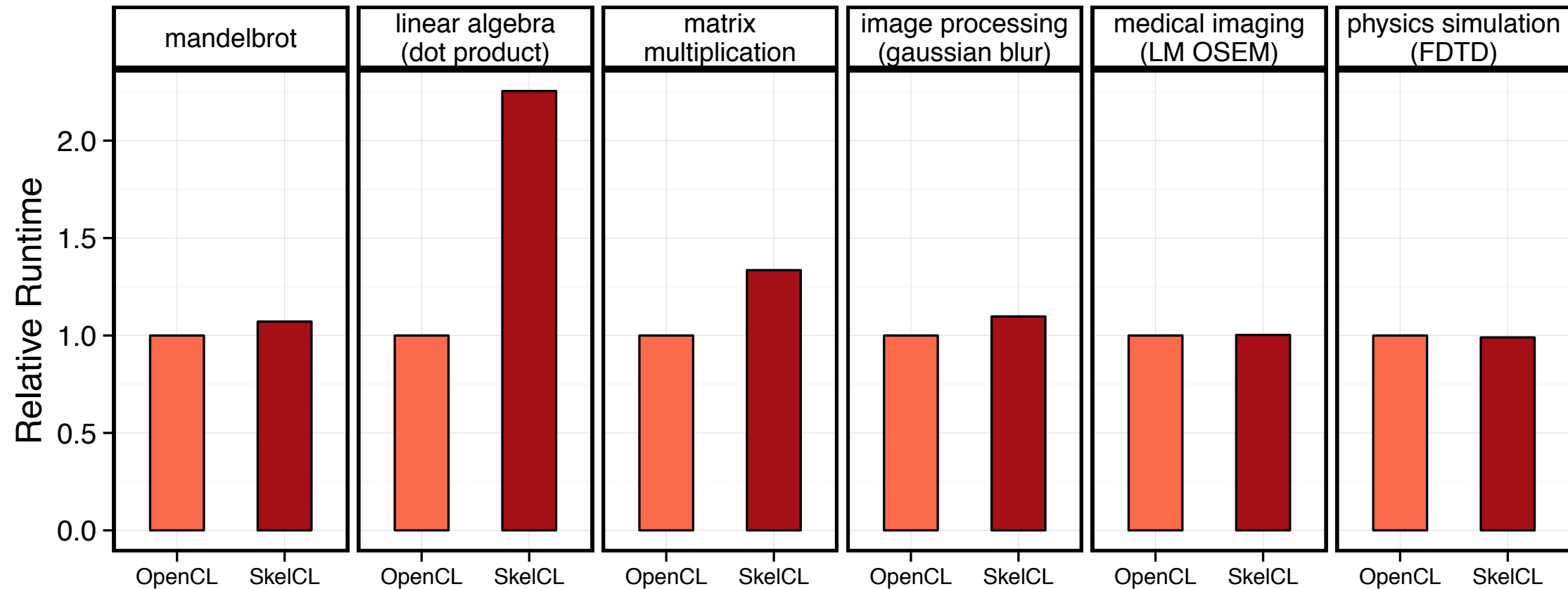
$$A \times B = \text{allpairs dotProduct } A \ B^T$$

Optimierung für zipReduce Muster:

$$\text{dotProduct } a \ b = \text{zipReduce } (+) \ 0 \ (\times) \ a \ b$$

Unterstützung für mehrere GPUs mit **block** und **copy** Verteilung

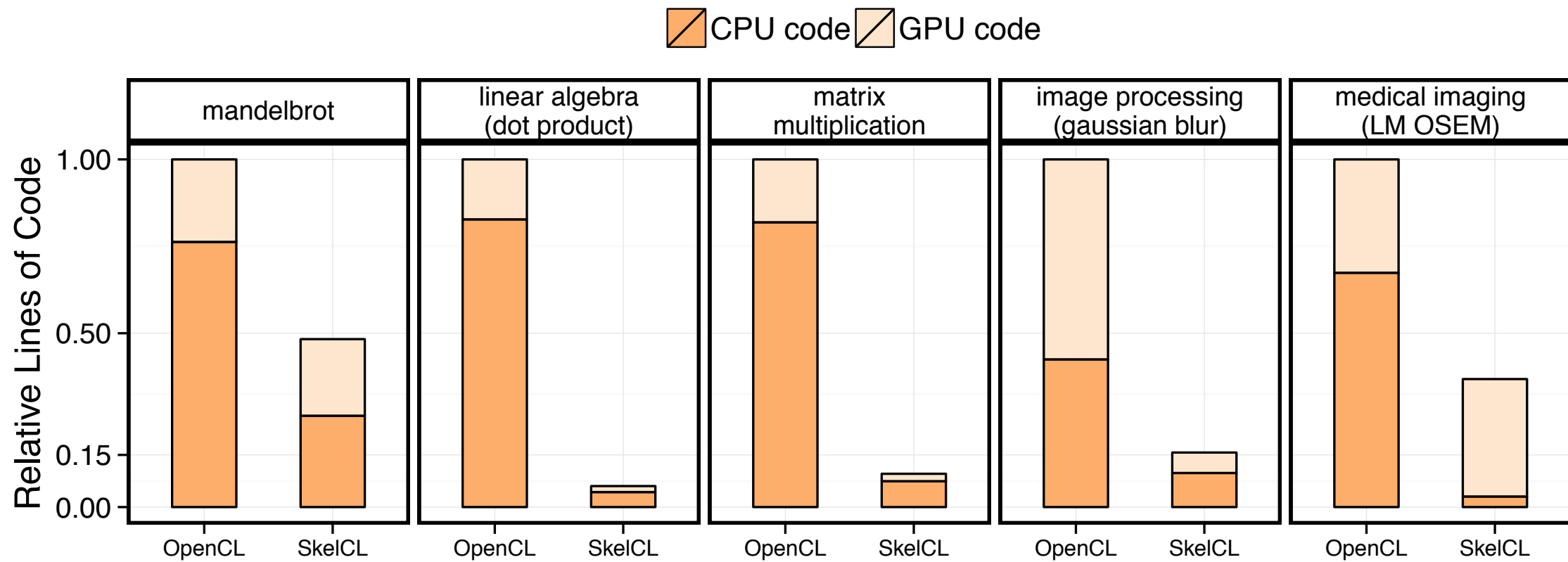
SkelCL Evaluation — Geschwindigkeit



SkelCL nahe an der Geschwindigkeit von OpenCL!

(Ausnahme: dot product ... mehr dazu in Teil II)

SkelCL Evaluation — Produktivität



SkelCL Programme sind signifikant kürzer!

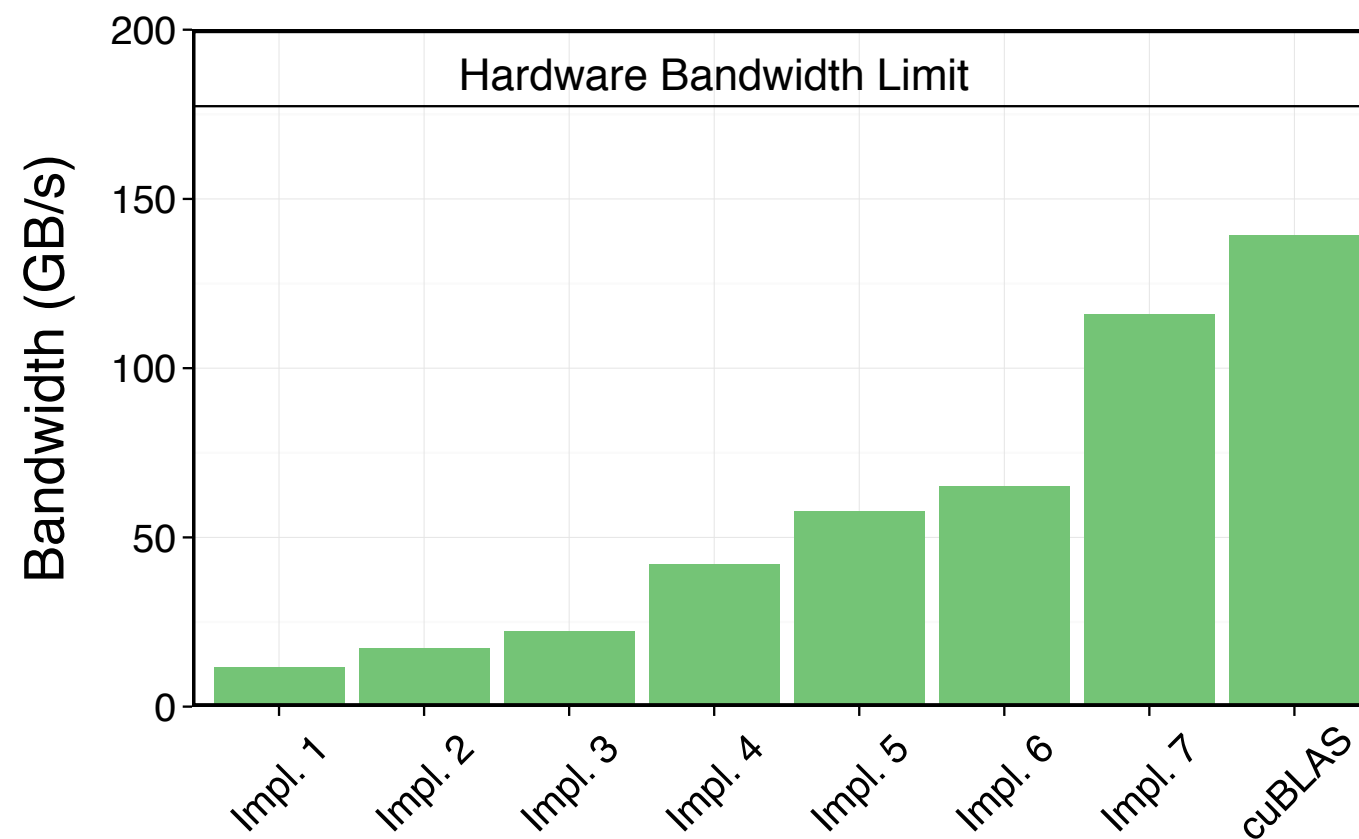
TEIL II

Die Herausforderung der Performance-Portabilität

EIN NEUER ANSATZ ZUR
PERFORMANCE PORTABLEN
CODEGENERIERUNG

OpenCL und Performance-Portabilität

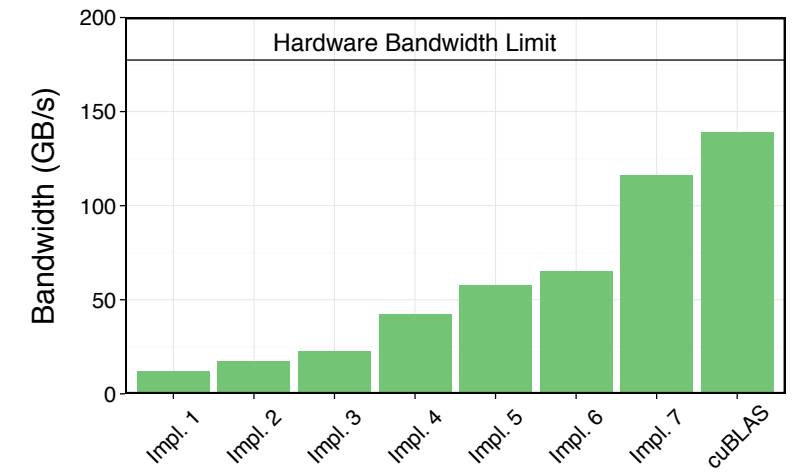
- Beispiel: Parallele Summation eines Arrays in OpenCL
- Vergleich von 7 OpenCL Implementierungen von Nvidia



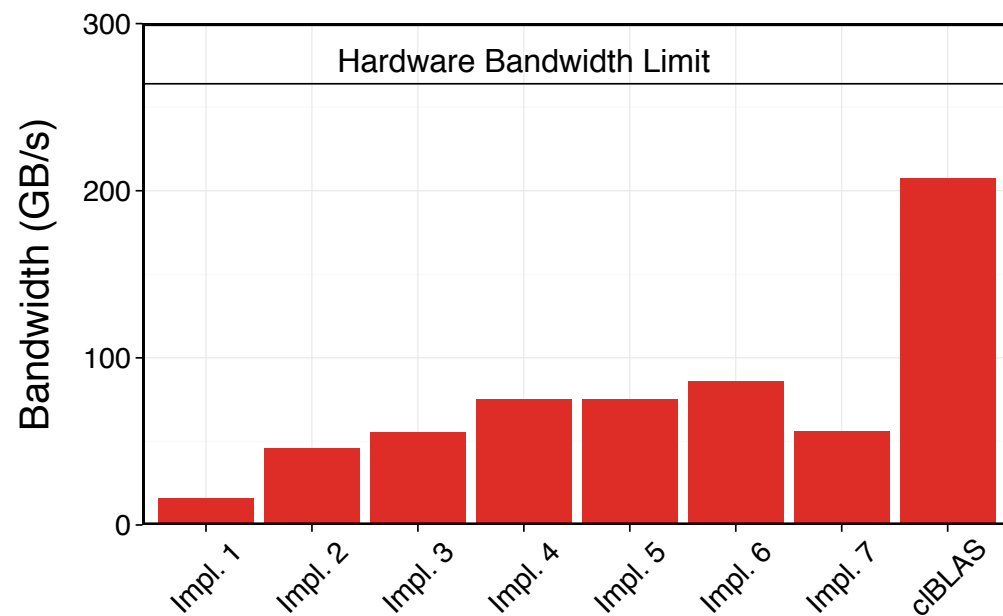
(a) Nvidia's GTX 480 GPU.

OpenCL und Performance-Portabilität

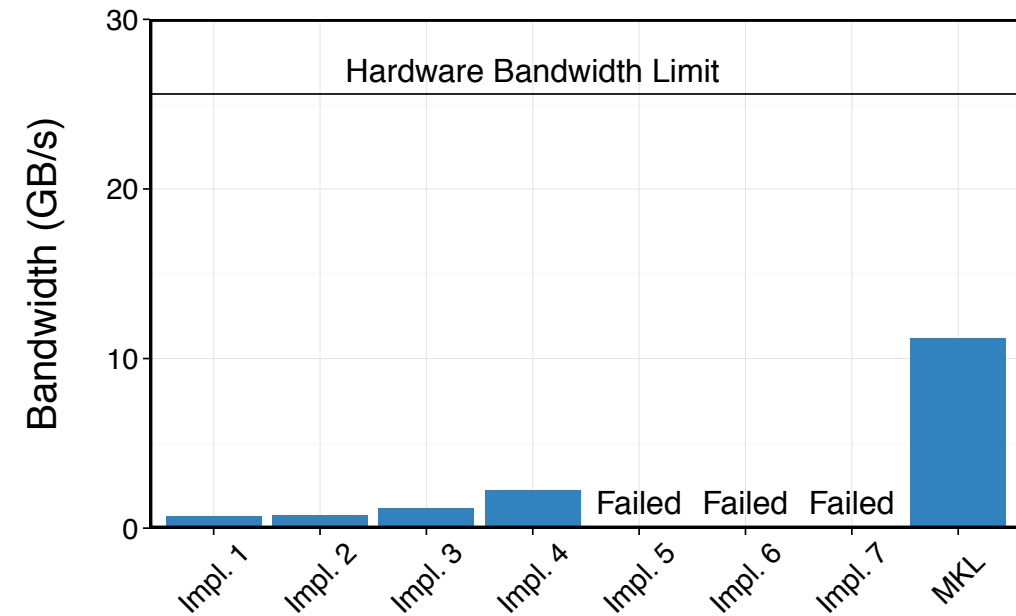
- Beispiel: Parallele Summation eines Arrays in OpenCL
- Vergleich von 7 OpenCL Implementierungen von Nvidia



(a) Nvidia's GTX 480 GPU.



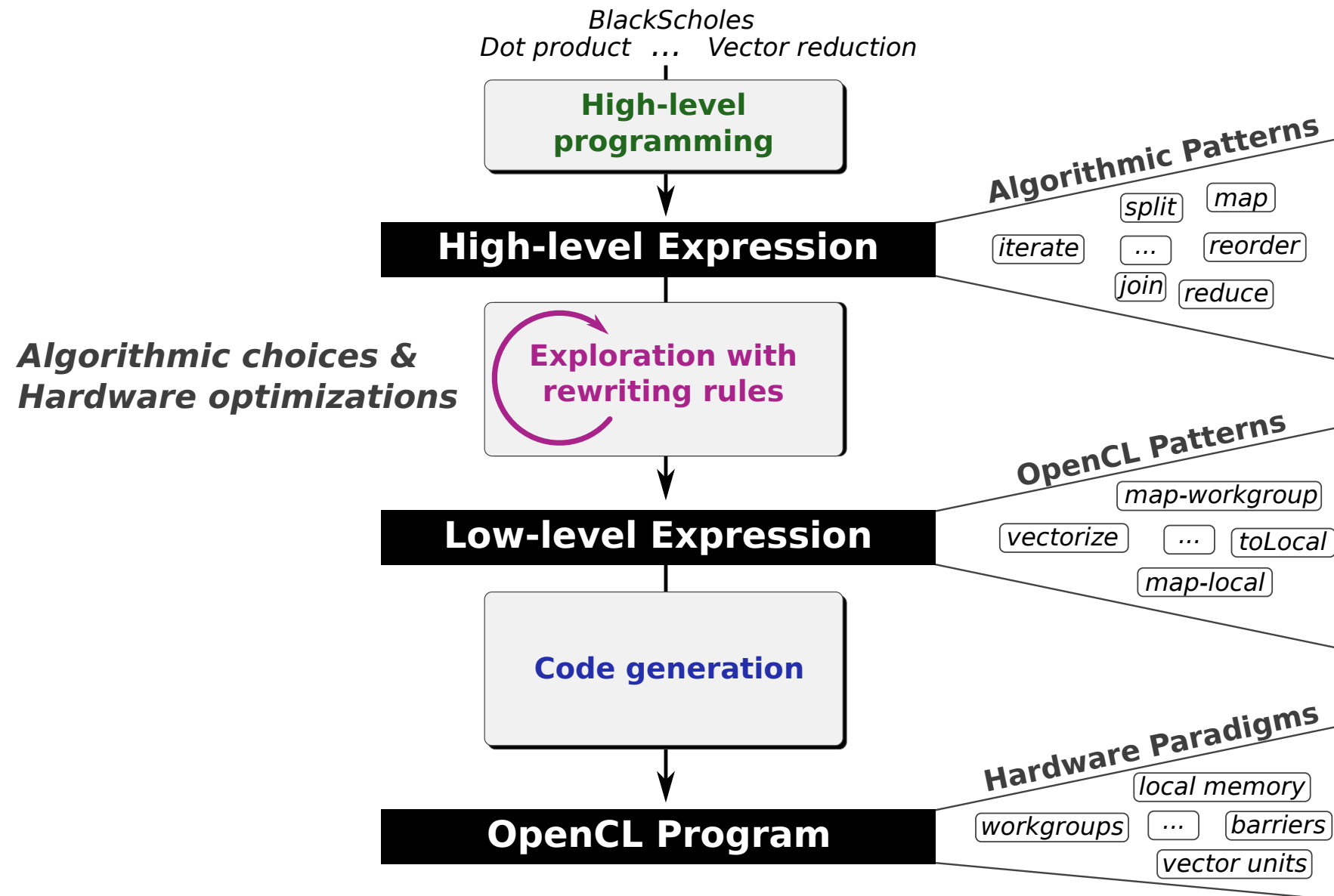
(b) AMD's HD 7970 GPU.



(c) Intel's E5530 dual-socket CPU.

Performance in OpenCL ist nicht portabel!

Performance portable Codegenerierung mit Transformationsregeln



Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

① Algorithmische Primitive

$$\mathit{map}_{A,B,I} : (A \rightarrow B) \rightarrow [A]_I \rightarrow [B]_I$$

$$\mathit{zip}_{A,B,I} : [A]_I \rightarrow [B]_I \rightarrow [A \times B]_I$$

$$\mathit{reduce}_{A,I} : ((A \times A) \rightarrow A) \rightarrow A \rightarrow [A]_I \rightarrow [A]_1$$

$$\mathit{split}_{A,I} : (n : \text{size}) \rightarrow [A]_{n \times I} \rightarrow [[A]_n]_I$$

$$\mathit{join}_{A,I,J} : [[A]_I]_J \rightarrow [A]_{I \times J}$$

$$\begin{aligned} \mathit{iterate}_{A,I,J} : (n : \text{size}) \rightarrow ((m : \text{size}) \rightarrow [A]_{I \times m} \rightarrow [A]_m) \\ \rightarrow [A]_{I^n \times J} \rightarrow [A]_J \end{aligned}$$

① High-Level Programme

$scal = \lambda a. map (*a)$

$asum = reduce (+) 0 \circ map abs$

$dot = \lambda xs\ ys. (reduce (+) 0 \circ map (*)) (zip xs ys)$

$gemv = \lambda mat\ xs\ ys\ \alpha\ \beta. map (+) ($
 $zip (map (scal \alpha \circ dot xs) mat) (scal \beta ys))$

Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >= 8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >= 4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >= 2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦
  join ◦ map-warp (
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32
  ) ◦ split 64 ◦
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦
  split (blockSize/128) ◦ reorder-stride 128
) ◦ split blockSize
```

③

```
kernel
void reduce6(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local volatile float* l_data) {
  unsigned int tid = get_local_id(0);
  unsigned int i =
    get_group_id(0) * (get_local_size(0)*2)
    + get_local_id(0);

  unsigned int gridSize =
    WG_SIZE * get_num_groups(0);
  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_idata[i];
    if (i + WG_SIZE < n)
      l_data[tid] += g_idata[i+WG_SIZE];
    i += gridSize; }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (WG_SIZE >= 256) {
    if (tid < 128) {
      l_data[tid] += l_data[tid+128]; }
    barrier(CLK_LOCAL_MEM_FENCE); }
  if (WG_SIZE >= 128) {
    if (tid < 64) {
      l_data[tid] += l_data[tid+ 64]; }
    barrier(CLK_LOCAL_MEM_FENCE); }
  if (tid < 32) {
    if (WG_SIZE >= 64) {
      l_data[tid] += l_data[tid+32]; }
    if (WG_SIZE >= 32) {
      l_data[tid] += l_data[tid+16]; }
    if (WG_SIZE >= 16) {
      l_data[tid] += l_data[tid+ 8]; }
    if (WG_SIZE >= 8) {
      l_data[tid] += l_data[tid+ 4]; }
    if (WG_SIZE >= 4) {
      l_data[tid] += l_data[tid+ 2]; }
    if (WG_SIZE >= 2) {
      l_data[tid] += l_data[tid+ 1]; } }
  if (tid == 0)
    g_odata[get_group_id(0)] = l_data[0];
}
```

② Algorithmische Transformationsregeln

- Transformationsregeln sind **semantikerhaltend**
- Drücken Auswahl bei der algorithmische Implementierungen aus

Split-Join Zerlegung:

$$\text{map } f \rightarrow \text{join} \circ \text{map } (\text{map } f) \circ \text{split } n$$

Map Zusammenschluss:

$$\text{map } f \circ \text{map } g \rightarrow \text{map } (f \circ g)$$

Reduktionsregeln:

$$\text{reduce } f \ z \rightarrow \text{reduce } f \ z \circ \text{reducePart } f \ z$$

$$\text{reducePart } f \ z \rightarrow \text{reducePart } f \ z \circ \text{reorder}$$

$$\text{reducePart } f \ z \rightarrow \text{join} \circ \text{map } (\text{reducePart } f \ z) \circ \text{split } n$$

$$\text{reducePart } f \ z \rightarrow \text{iterate } n \ (\text{reducePart } f \ z)$$

② OpenCL Primitive

Primitive

OpenCL Konzept

mapGlobal

Work-items

mapWorkgroup / mapLocal

Work-groups

mapSeq / reduceSeq

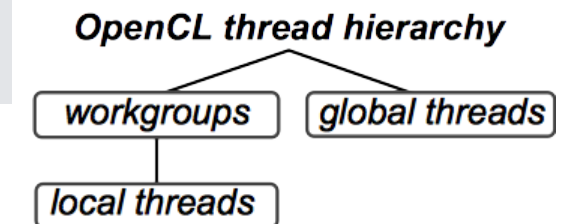
Sequentielle Implementierungen

toLocal / toGlobal

Speicherbereiche

mapVec / splitVec / joinVec

Vektorisierung



② OpenCL Transformationsregeln

- Drücken hardware-spezifische Optimierungen aus

Map:

$$\text{map } f \rightarrow \text{mapWorkgroup } f \mid \text{mapLocal } f \mid \text{mapGlobal } f \mid \text{mapSeq } f$$

Lokaler/ Globaler Speicher:

$$\text{mapLocal } f \rightarrow \text{toLocal } (\text{mapLocal } f) \quad \text{mapLocal } f \rightarrow \text{toGlobal } (\text{mapLocal } f)$$

Vektorisierung:

$$\text{map } f \rightarrow \text{joinVec} \circ \text{map } (\text{mapVec } f) \circ \text{splitVec } n$$

Map-Reduktion Zusammenschluss:

$$\text{reduceSeq } f \ z \circ \text{mapSeq } g \rightarrow \text{reduceSeq } (\lambda (acc, x). f (acc, g x)) \ z$$

Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (  
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦  
  join ◦ map-warp (  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦  
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32  
  ) ◦ split 64 ◦  
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦  
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦  
  split (blockSize/128) ◦ reorder-stride 128  
) ◦ split blockSize
```

③

```
kernel  
void reduce6(global float* g_idata,  
             global float* g_odata,  
             unsigned int n,  
             local volatile float* l_data) {  
  unsigned int tid = get_local_id(0);  
  unsigned int i =  
    get_group_id(0) * (get_local_size(0)*2)  
    + get_local_id(0);  
  
  unsigned int gridSize =  
    WG_SIZE * get_num_groups(0);  
  l_data[tid] = 0;  
  while (i < n) {  
    l_data[tid] += g_idata[i];  
    if (i + WG_SIZE < n)  
      l_data[tid] += g_idata[i+WG_SIZE];  
    i += gridSize; }  
  barrier(CLK_LOCAL_MEM_FENCE);  
  
  if (WG_SIZE >= 256) {  
    if (tid < 128) {  
      l_data[tid] += l_data[tid+128]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (WG_SIZE >= 128) {  
    if (tid < 64) {  
      l_data[tid] += l_data[tid+ 64]; }  
    barrier(CLK_LOCAL_MEM_FENCE); }  
  if (tid < 32) {  
    if (WG_SIZE >= 64) {  
      l_data[tid] += l_data[tid+32]; }  
    if (WG_SIZE >= 32) {  
      l_data[tid] += l_data[tid+16]; }  
    if (WG_SIZE >= 16) {  
      l_data[tid] += l_data[tid+ 8]; }  
    if (WG_SIZE >=  8) {  
      l_data[tid] += l_data[tid+ 4]; }  
    if (WG_SIZE >=  4) {  
      l_data[tid] += l_data[tid+ 2]; }  
    if (WG_SIZE >=  2) {  
      l_data[tid] += l_data[tid+ 1]; } }  
  if (tid == 0)  
    g_odata[get_group_id(0)] = l_data[0];  
}
```

Beispiel: Parallele Summation

① $\text{vecSum} = \text{reduce } (+) 0$

rewrite rules

code generation

②

```
vecSum = reduce ◦ join ◦ map-workgroup (
  join ◦ toGlobal (map-local (map-seq id)) ◦ split 1 ◦
  join ◦ map-warp (
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 1 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 2 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 4 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 8 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 16 ◦
    join ◦ map-lane (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 32
  ) ◦ split 64 ◦
  join ◦ map-local (reduce-seq (+) 0) ◦ split 2 ◦ reorder-stride 64 ◦
  join ◦ toLocal (map-local (reduce-seq (+) 0)) ◦
  split (blockSize/128) ◦ reorder-stride 128
) ◦ split blockSize
```

③

```
kernel
void reduce6(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local volatile float* l_data) {
  unsigned int tid = get_local_id(0);
  unsigned int i =
    get_group_id(0) * (get_local_size(0)*2)
    + get_local_id(0);

  unsigned int gridSize =
    WG_SIZE * get_num_groups(0);
  l_data[tid] = 0;
  while (i < n) {
    l_data[tid] += g_idata[i];
    if (i + WG_SIZE < n)
      l_data[tid] += g_idata[i+WG_SIZE];
    i += gridSize; }
  barrier(CLK_LOCAL_MEM_FENCE);

  if (WG_SIZE >= 256) {
    if (tid < 128) {
      l_data[tid] += l_data[tid+128]; }
    barrier(CLK_LOCAL_MEM_FENCE); }
  if (WG_SIZE >= 128) {
    if (tid < 64) {
      l_data[tid] += l_data[tid+ 64]; }
    barrier(CLK_LOCAL_MEM_FENCE); }
  if (tid < 32) {
    if (WG_SIZE >= 64) {
      l_data[tid] += l_data[tid+32]; }
    if (WG_SIZE >= 32) {
      l_data[tid] += l_data[tid+16]; }
    if (WG_SIZE >= 16) {
      l_data[tid] += l_data[tid+ 8]; }
    if (WG_SIZE >= 8) {
      l_data[tid] += l_data[tid+ 4]; }
    if (WG_SIZE >= 4) {
      l_data[tid] += l_data[tid+ 2]; }
    if (WG_SIZE >= 2) {
      l_data[tid] += l_data[tid+ 1]; } }
  if (tid == 0)
    g_odata[get_group_id(0)] = l_data[0];
}
```


③ Muster basierte OpenCL Codegenerierung

- Generiere OpenCL Code für jedes OpenCL Primitiv

mapGlobal f xs →

```
for (int g_id = get_global_id(0); g_id < n;  
     g_id += get_global_size(0)) {  
    output[g_id] = f(xs[g_id]);  
}
```

reduceSeq f z xs →

```
T acc = z;  
for (int i = 0; i < n; ++i) {  
    acc = f(acc, xs[i]);  
}
```

⋮

⋮

Transformationsregeln definieren
einen Suchraum gültiger Implementierungen

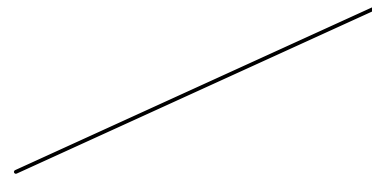
$$\begin{array}{c} \textit{reduce (+) 0} \\ | \\ \textit{reduce (+) 0} \circ \textit{reducePart (+) 0} \end{array}$$

Transformationsregeln definieren einen Suchraum gültiger Implementierungen

$reduce (+) 0$

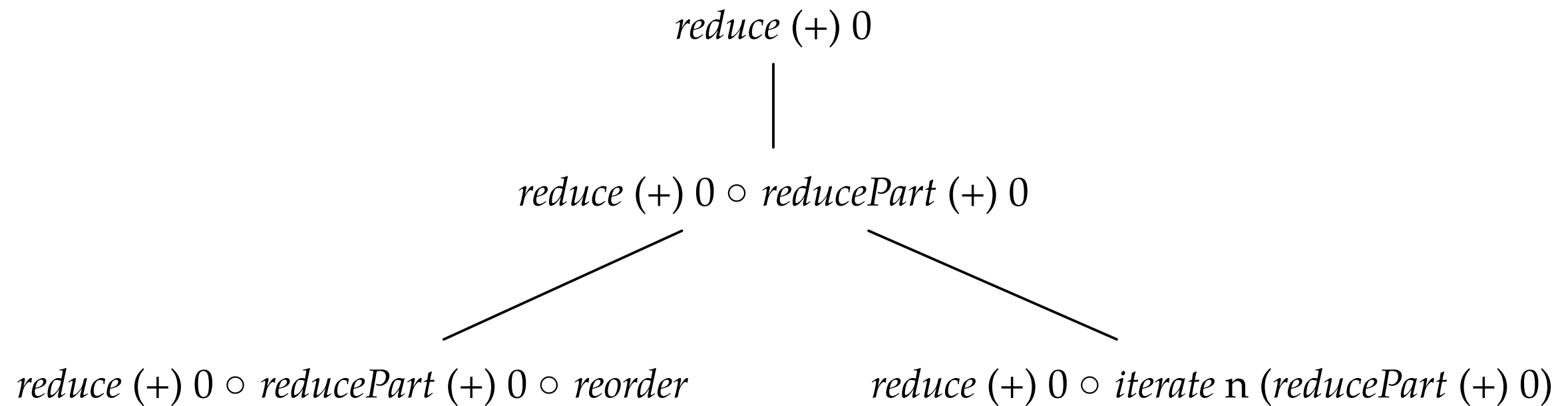


$reduce (+) 0 \circ reducePart (+) 0$

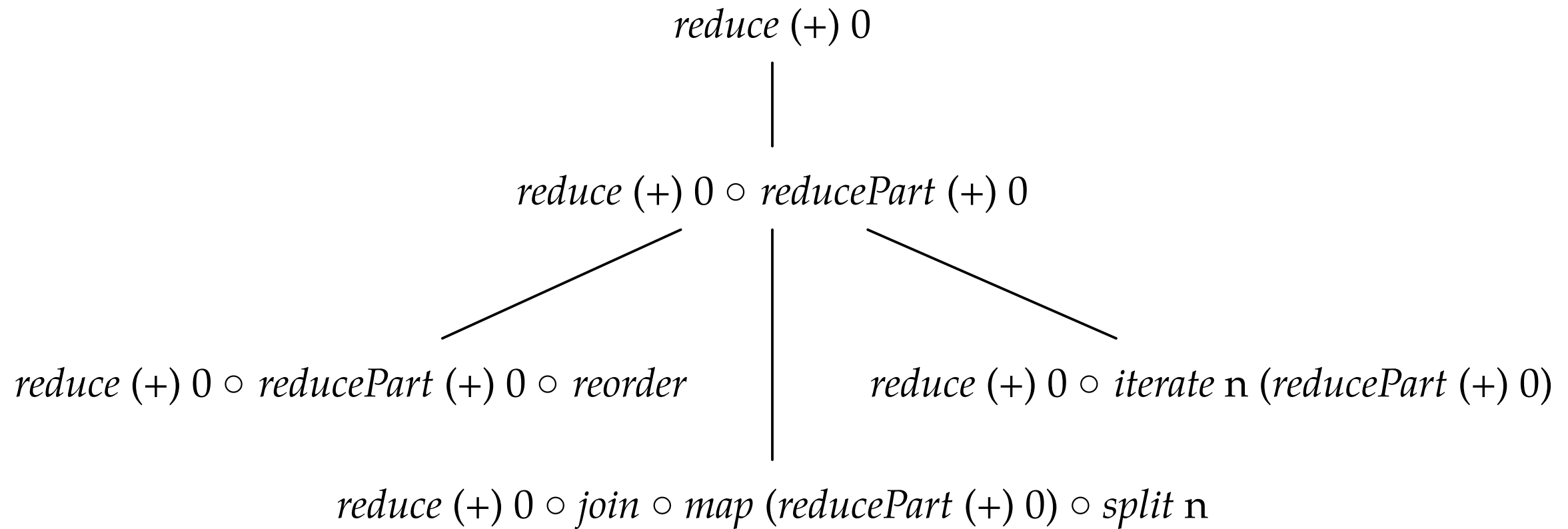


$reduce (+) 0 \circ reducePart (+) 0 \circ reorder$

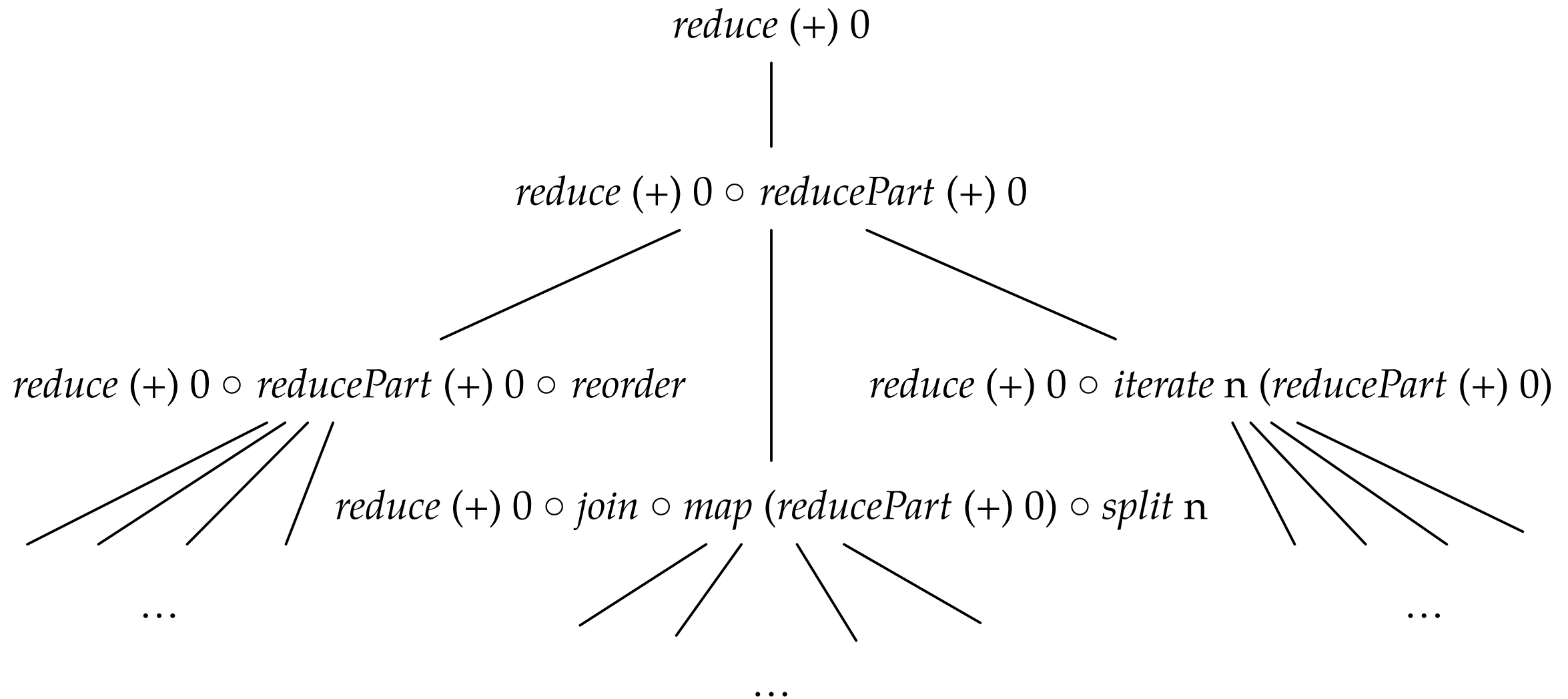
Transformationsregeln definieren einen Suchraum gültiger Implementierungen



Transformationsregeln definieren einen Suchraum gültiger Implementierungen



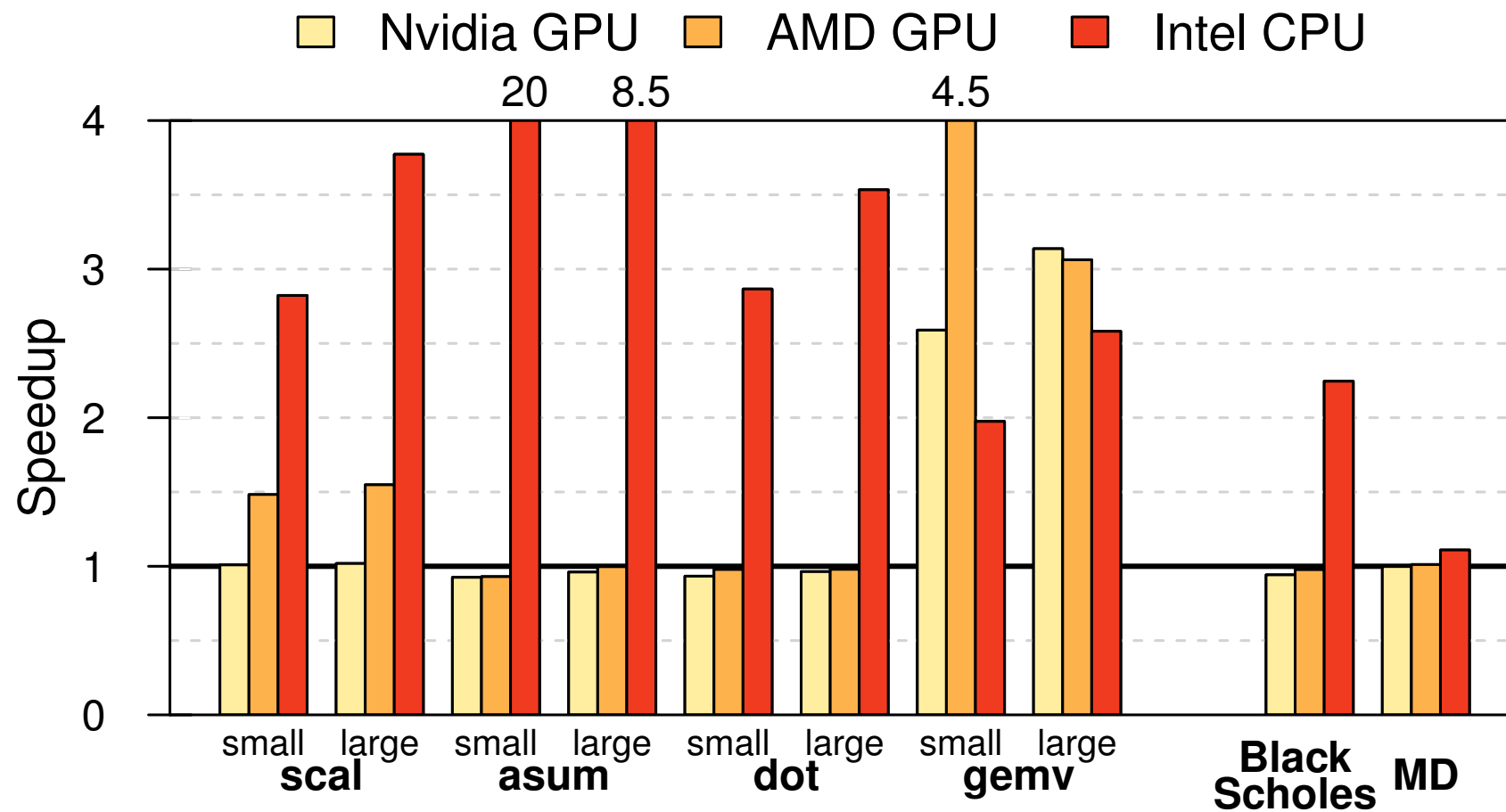
Transformationsregeln definieren einen Suchraum gültiger Implementierungen



- Vollautomatische Suche nach guten Implementierungen möglich!
(Eine einfache Suchstrategie ist in der Dissertation beschrieben)

Evaluation – Geschwindigkeit

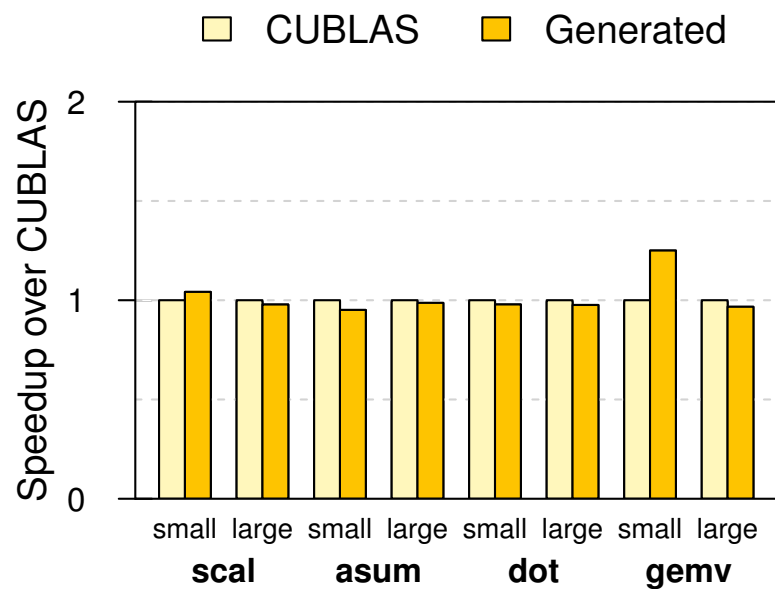
gegenüber einer funktional portablen Implementierung



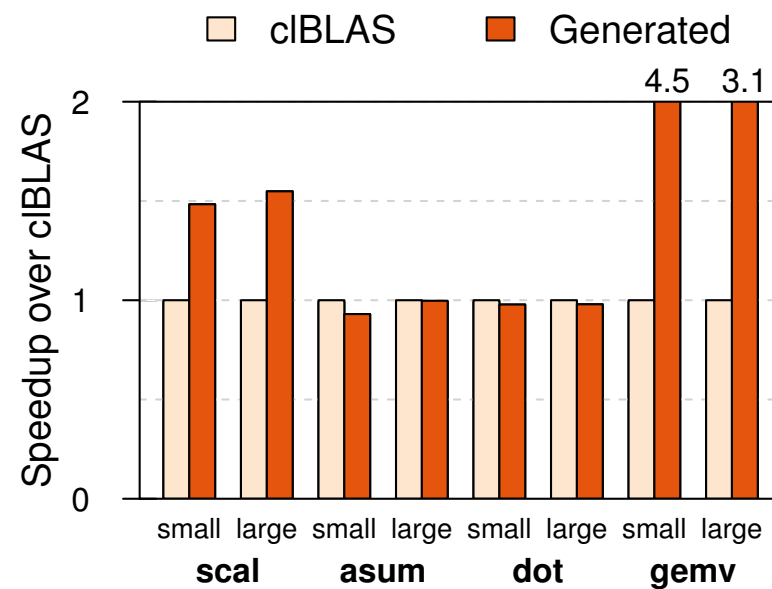
Bis zu **20x** Speedup gegenüber der funktional portablen clBLAS Implementierung

Evaluation – Geschwindigkeit

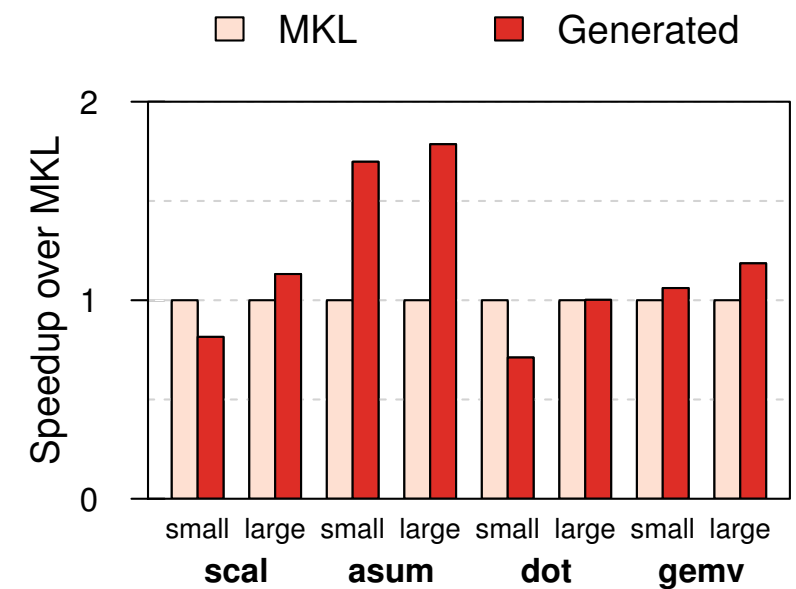
gegenüber Hardware spezifischen Implementierungen



(a) Nvidia GPU



(b) AMD GPU



(c) Intel CPU

- Automatisch generierter Code vs. handoptimierten Code
- Konkurrenzfähige Ergebnisse vs. hochoptimierte Implementierungen
- Bis zu **4.5x** Speedup für gemv auf der AMD GPU

Zusammenfassung

- Um die *Herausforderung der Programmierbarkeit* zu adressieren:
 - Ein neuer Ansatz zur Programmierung von Systemen mit mehreren GPUs
 - Zwei neue formell definierte und implementierte algorithmische Skelette
- Um die *Herausforderung der Performance-Portabilität* zu adressieren:
 - Ein formelles System zur Transformation muster-basierter Programme
 - Ein Codegenerator der Performance-Portabilität erreicht

Ergebnisse der Suche

Automatisch Gefundene Ausdrücke

$$asum = reduce (+) 0 \circ map abs$$


Nvidia
GPU

$$\lambda x. (reduceSeq \circ join \circ join \circ mapWorkgroup (toGlobal (mapLocal (reduceSeq (\lambda(a, b). a + (abs b)) 0)) \circ reorderStride 2048) \circ split 128 \circ split 2048) x$$

AMD
GPU

$$\lambda x. (reduceSeq \circ join \circ joinVec \circ join \circ mapWorkgroup (mapLocal (reduceSeq (mapVec 2 (\lambda(a, b). a + (abs b))) 0 \circ reorderStride 2048) \circ split 128 \circ splitVec 2 \circ split 4096) x$$

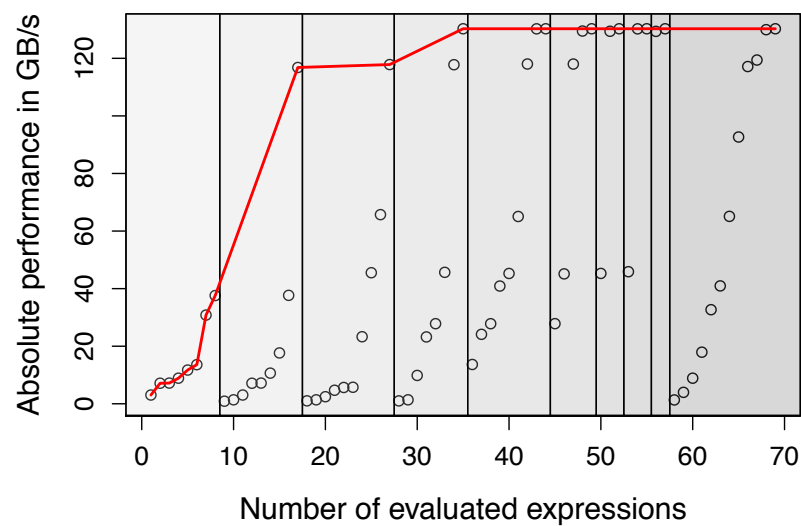
Intel
CPU

$$\lambda x. (reduceSeq \circ join \circ mapWorkgroup (join \circ joinVec \circ mapLocal (reduceSeq (mapVec 4 (\lambda(a, b). a + (abs b))) 0) \circ splitVec 4 \circ split 32768) \circ split 32768) x$$

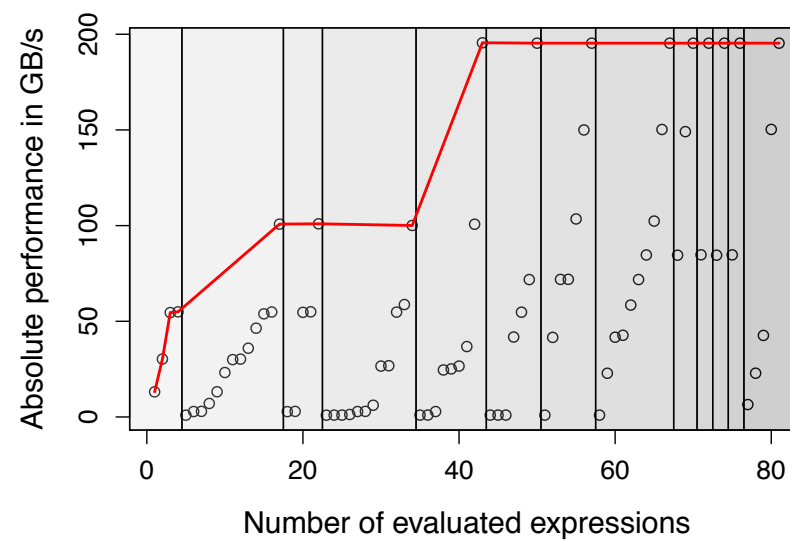
Gesucht für: **Nvidia** GTX 480 GPU, **AMD** Radeon HD 7970 GPU, **Intel** Xeon E5530 CPU

Ergebnisse der Suche

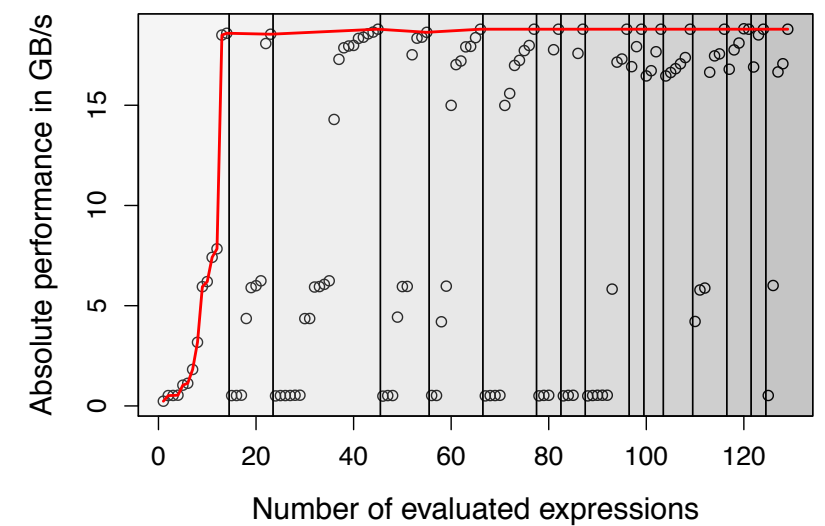
Effizienz der Suche



(a) Nvidia GPU



(b) AMD GPU



(c) Intel CPU

- Die Suche hat auf jeder Plattform weniger als 1 Stunde gedauert
- Durchschnittliche Zeit zur Ausführung eines Kandidaten weniger als 1/2 Sekunde

Fazit des Beispiels

- Optimieren in OpenCL ist kompliziert
 - Verständnis für die Zielarchitektur benötigt
- Veränderungen im Program nicht offensichtlich

```
kernel
void reduce0(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i = get_global_id(0);
    l_data[tid] = (i < n) ? g_idata[i] : 0;
    barrier(CLK_LOCAL_MEM_FENCE);

    for (unsigned int s=1;
         s < get_local_size(0); s*= 2) {
        if ((tid % (2*s)) == 0) {
            l_data[tid] += l_data[tid + s];
        }
        barrier(CLK_LOCAL_MEM_FENCE);
    }
    if (tid == 0)
        g_odata[get_group_id(0)] = l_data[0];
}
```

Nicht Optimierte Implementierung

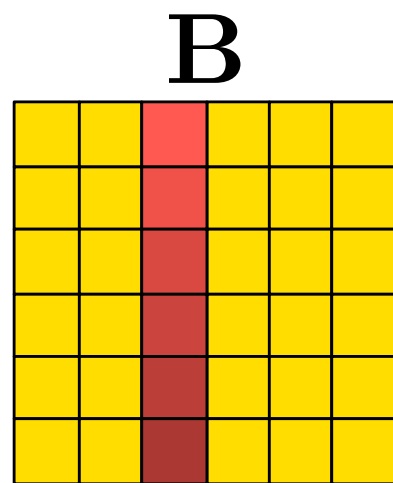
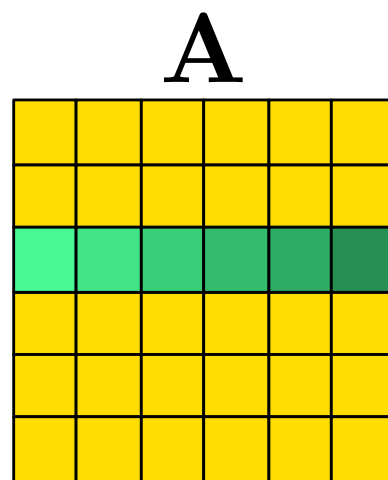
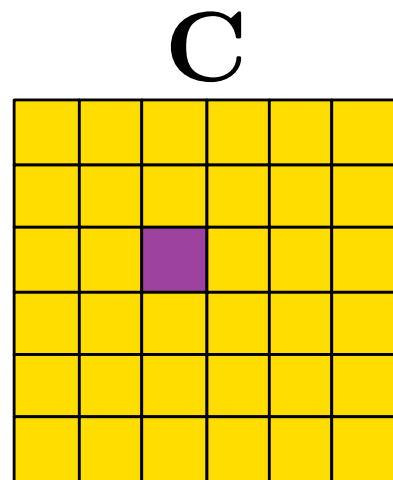
```
kernel
void reduce6(global float* g_idata,
             global float* g_odata,
             unsigned int n,
             local volatile float* l_data) {
    unsigned int tid = get_local_id(0);
    unsigned int i =
        get_group_id(0) * (get_local_size(0)*2)
        + get_local_id(0);

    unsigned int gridSize =
        WG_SIZE * get_num_groups(0);
    l_data[tid] = 0;
    while (i < n) {
        l_data[tid] += g_idata[i];
        if (i + WG_SIZE < n)
            l_data[tid] += g_idata[i+WG_SIZE];
        i += gridSize; }
    barrier(CLK_LOCAL_MEM_FENCE);

    if (WG_SIZE >= 256) {
        if (tid < 128) {
            l_data[tid] += l_data[tid+128]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (WG_SIZE >= 128) {
        if (tid < 64) {
            l_data[tid] += l_data[tid+ 64]; }
        barrier(CLK_LOCAL_MEM_FENCE); }
    if (tid < 32) {
        if (WG_SIZE >= 64) {
            l_data[tid] += l_data[tid+32]; }
        if (WG_SIZE >= 32) {
            l_data[tid] += l_data[tid+16]; }
        if (WG_SIZE >= 16) {
            l_data[tid] += l_data[tid+ 8]; }
        if (WG_SIZE >= 8) {
            l_data[tid] += l_data[tid+ 4]; }
        if (WG_SIZE >= 4) {
            l_data[tid] += l_data[tid+ 2]; }
        if (WG_SIZE >= 2) {
            l_data[tid] += l_data[tid+ 1]; } }
    if (tid == 0)
        g_odata[get_group_id(0)] = l_data[0];
}
```

Voll Optimierte Implementierung

Matrix Multiplikation



$A \times B =$
 $\text{map}(\lambda \text{ rowA} \mapsto$
 $\text{map}(\lambda \text{ colB} \mapsto$
 $\text{dotProduct}(\text{rowA}, \text{colB})$
 , $\text{transpose}(B)$)
 , A)

Suche für Matrix Multiplikation

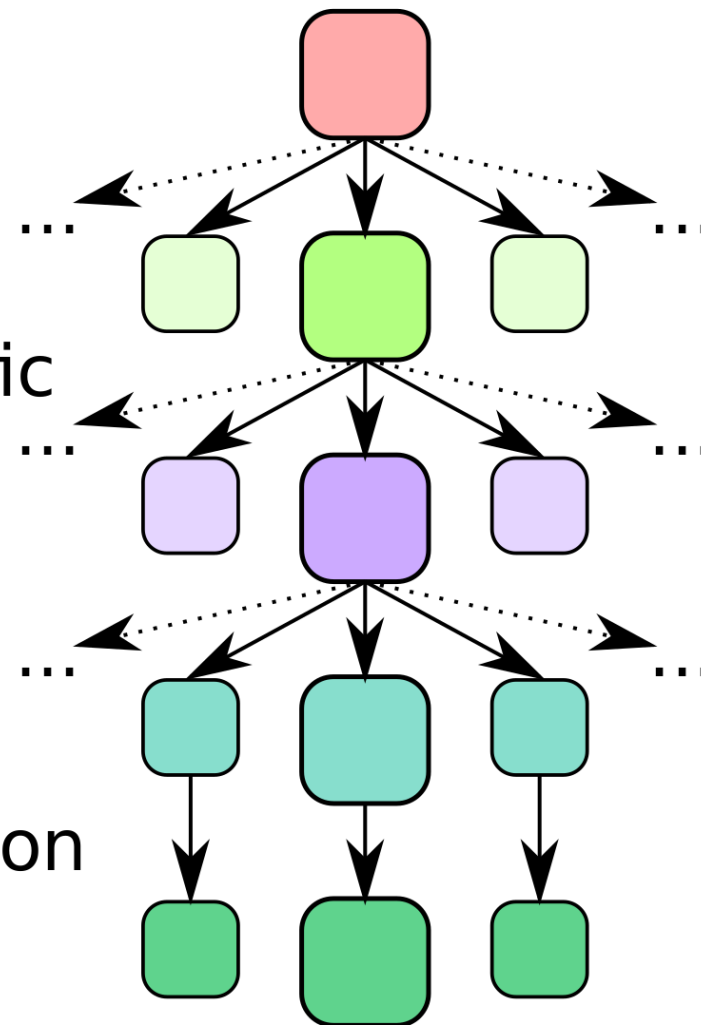
Phases:

Algorithmic
Exploration

OpenCL specific
Exploration

Parameter
Exploration

Code Generation



Program Variants:

High-Level Program 1

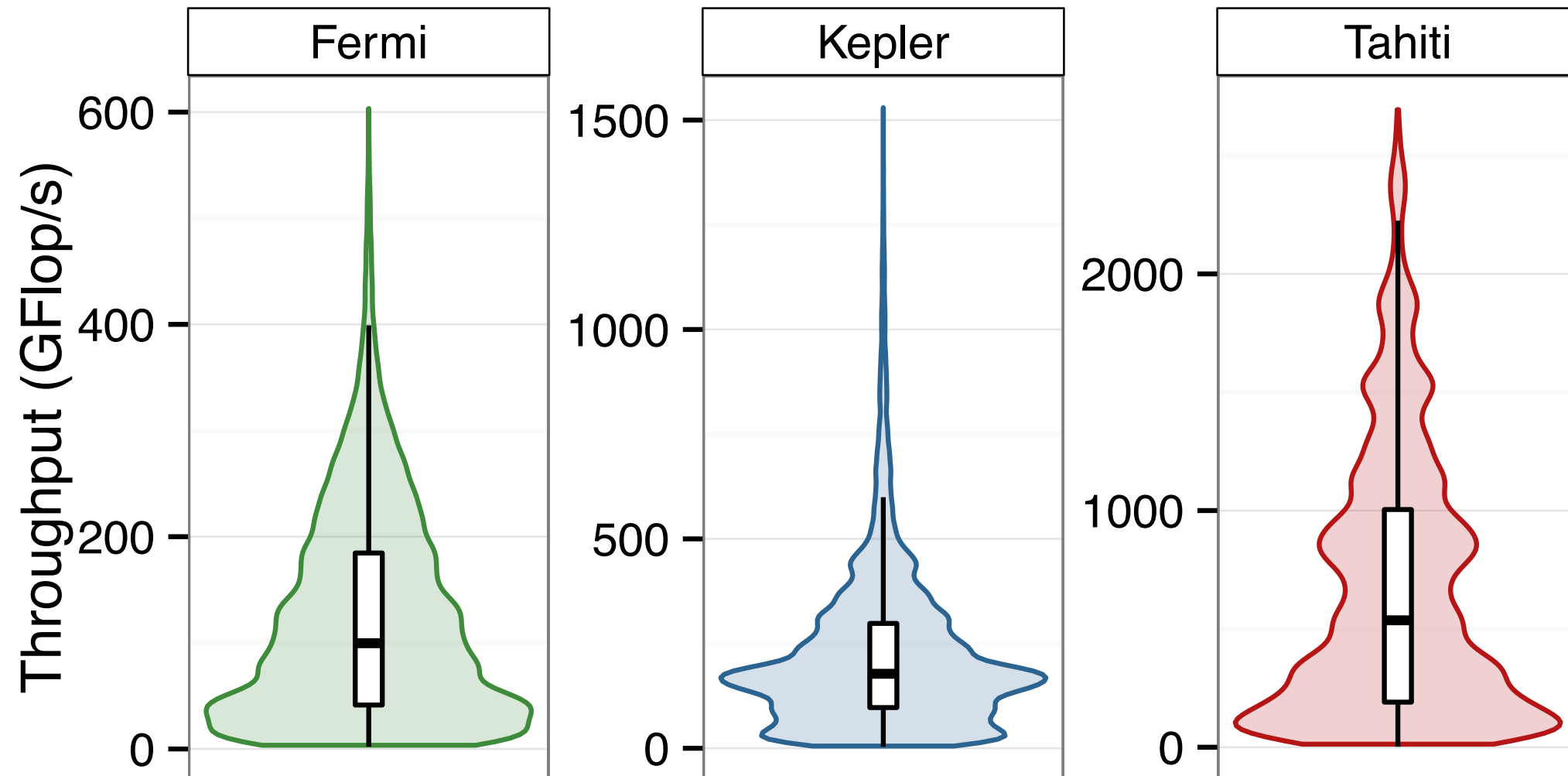
Algorithmic
Rewritten Program 8

OpenCL Specific
Program 760

Fully Specialized
Program 46,000

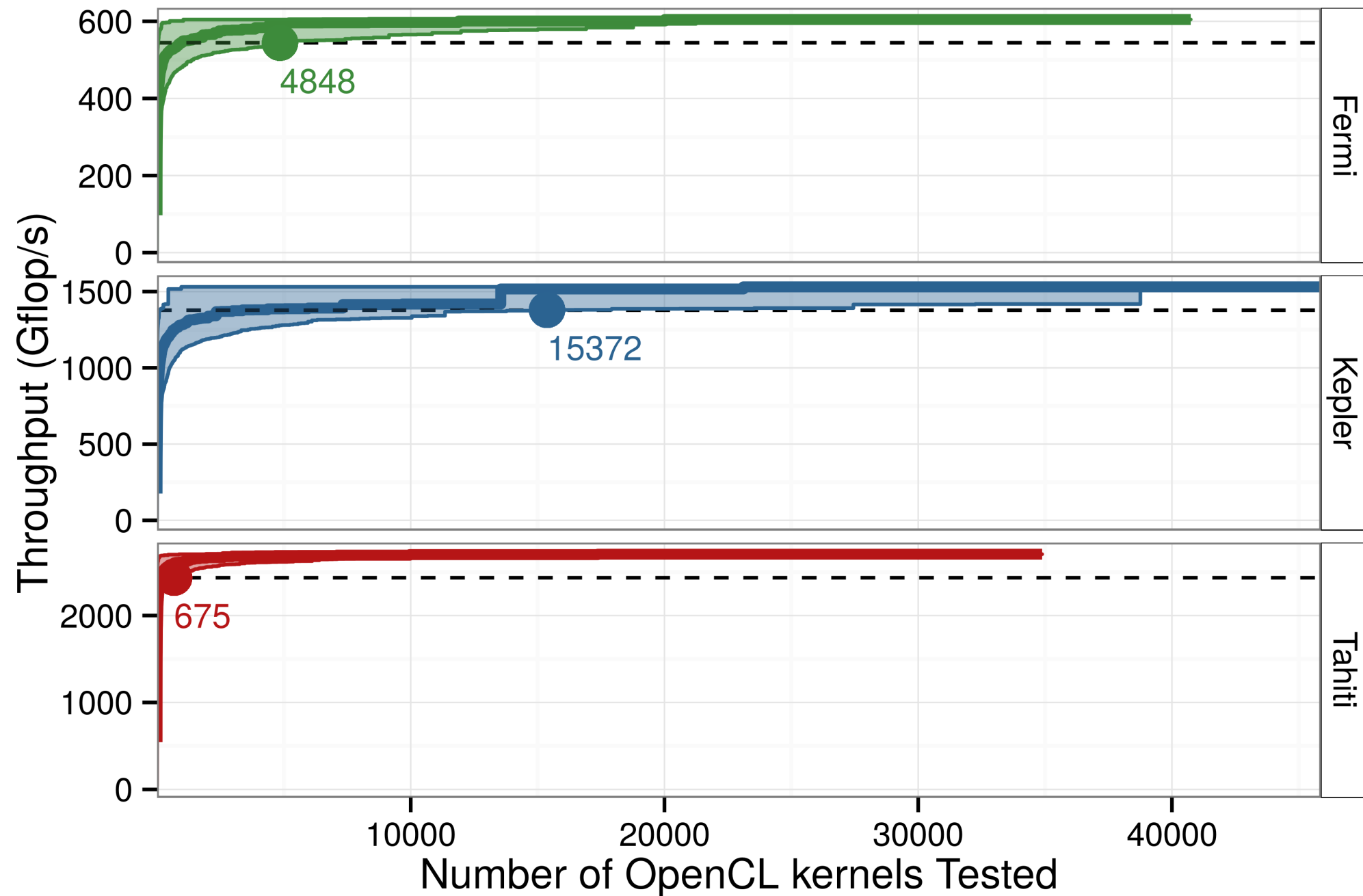
OpenCL Code 46,000

Suchraum für Matrix Multiplikation



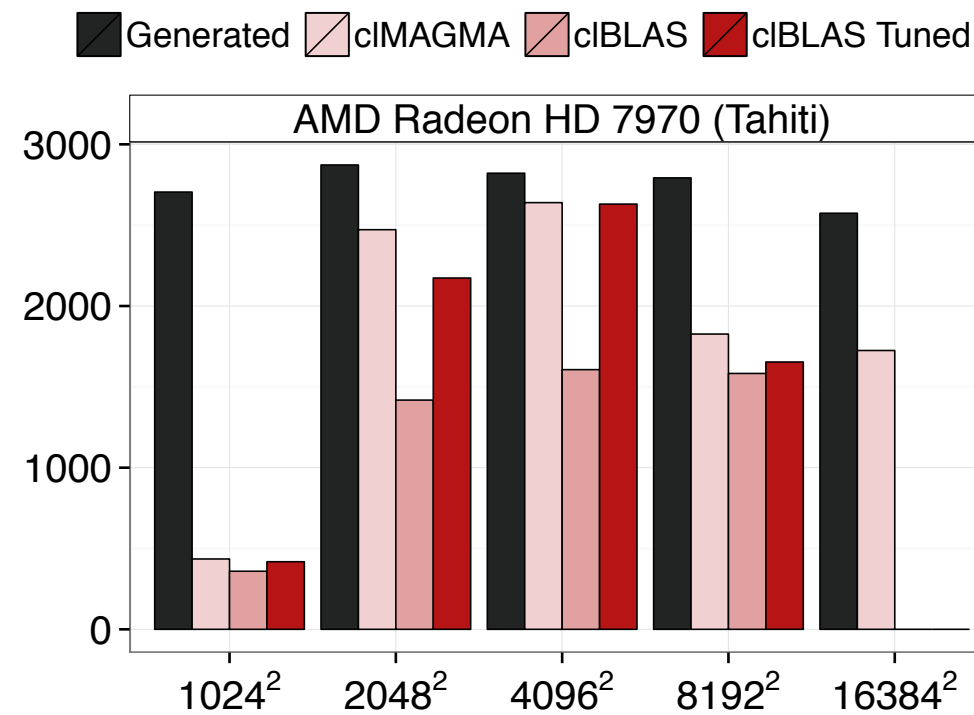
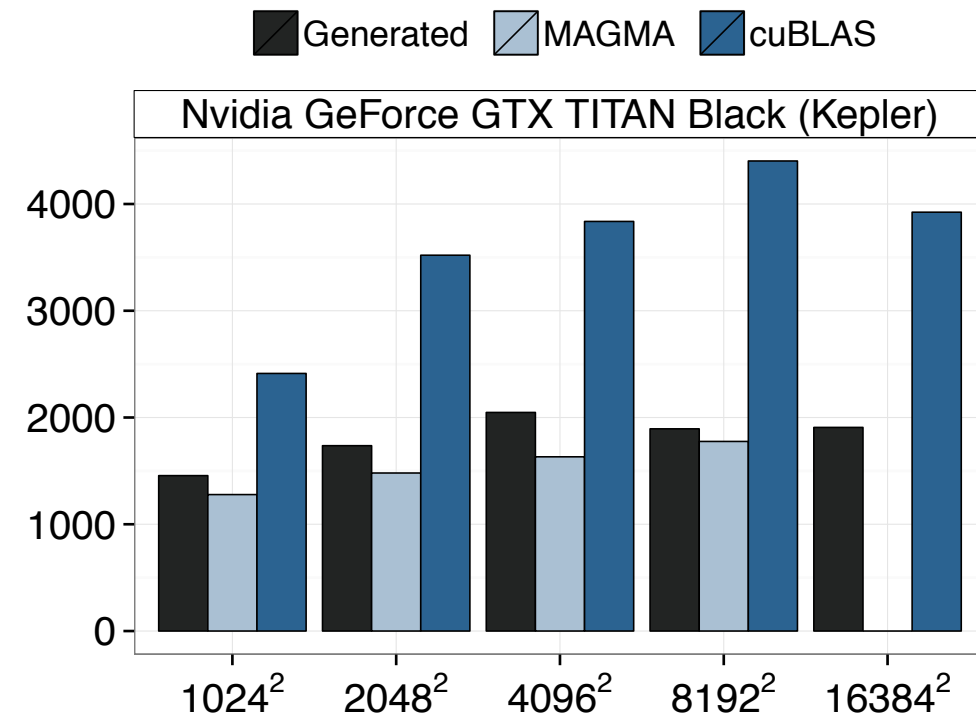
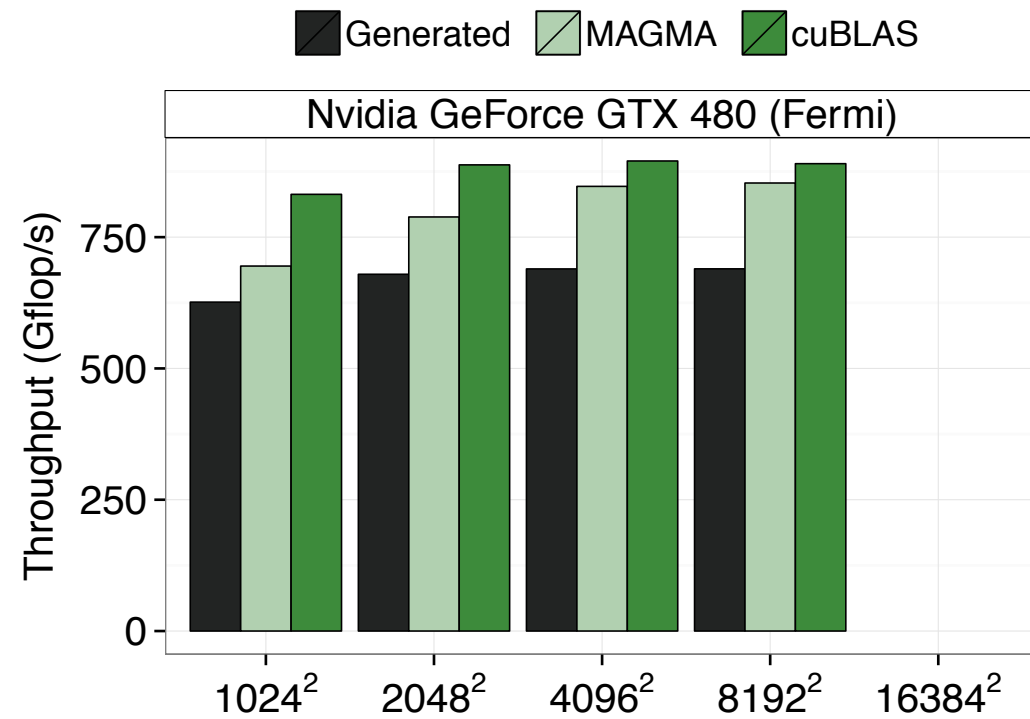
Nur einige generierte OpenCL Programme mit sehr guter Performance

Performance Entwicklung für Matrix Multiplikation



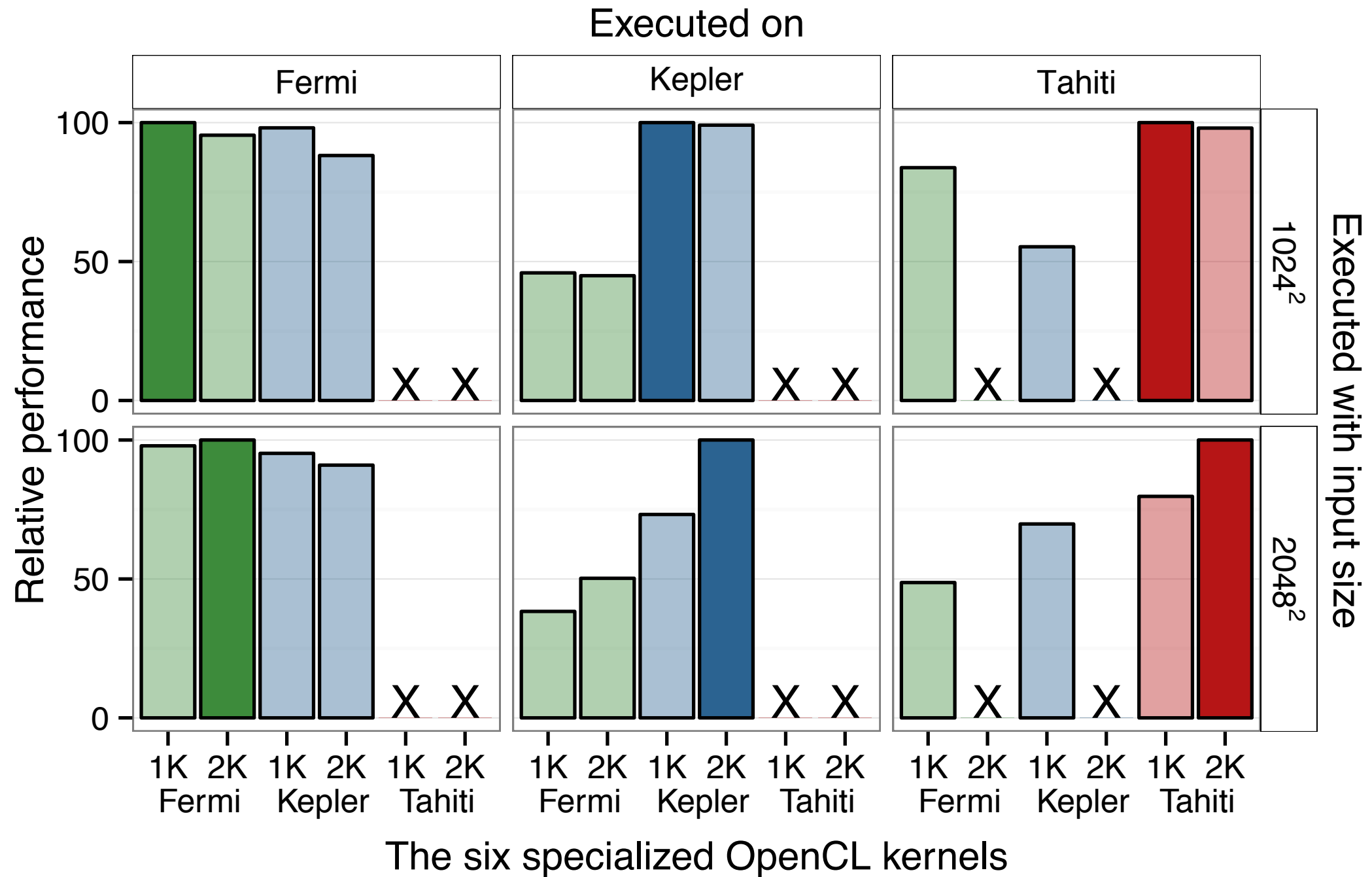
Selbst mit einer einfachen zufälligen Strategie kann man erwarten schnell ein Program mit guter Performance zu finden

Evaluation für Matrix Multiplikation



Performance nahe oder sogar besser als handoptimierte MAGMA Bibliothek

Performance-Portabilität von Matrix Multiplikation



Generierte Programme sind spezialisiert für GPU und Eingabegröße