# HIGH PERFORMANCE STENCIL CODE GENERATION WITH LIFT

**Bastian Hagedorn** | Larisa Stoltzfus | Michel Steuwer | Sergei Gorlatch | Christophe Dubach

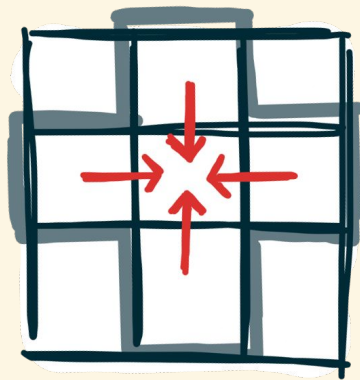WWU MÜNSTER

THE UNIVERSITY of EDINBURGH

University of Glasgow

# WHY STENCIL COMPUTATIONS?



**Stencil computations** are a class of kernels which update *neighboring* array elements according to a fixed pattern, called *stencil*.
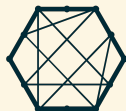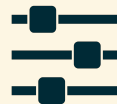
Frequently occur in:

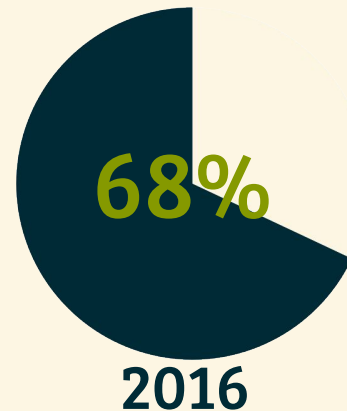♥ Medical Imaging        🌎 Physics Simulations

⬡ Machine Learning        🎚 PDE Solvers

# WHY STENCIL COMPUTATIONS?

Stencil compute time:

HPC Center
München

**lrz**

**49%**

2017

HPC Center
Stuttgart

H L R I S

**68%**

2016

Frequently occur in:

Medical Imaging

Physics Simulations

Machine Learning

PDE Solvers

# YET ANOTHER STENCIL PAPER?



2005     2007     2009     2011     2013     2015     2018

ICS'05     ICS'09     CGO'12     CGO'15     CLUSTER'17

PLDI'07     SC'10     CLUSTER'13     WOLFHPC'16     CGO'18

# DOMAIN SPECIFIC LANGUAGES

PATUS  Pochoir  PARTANS  Halide  ...  *DSL*

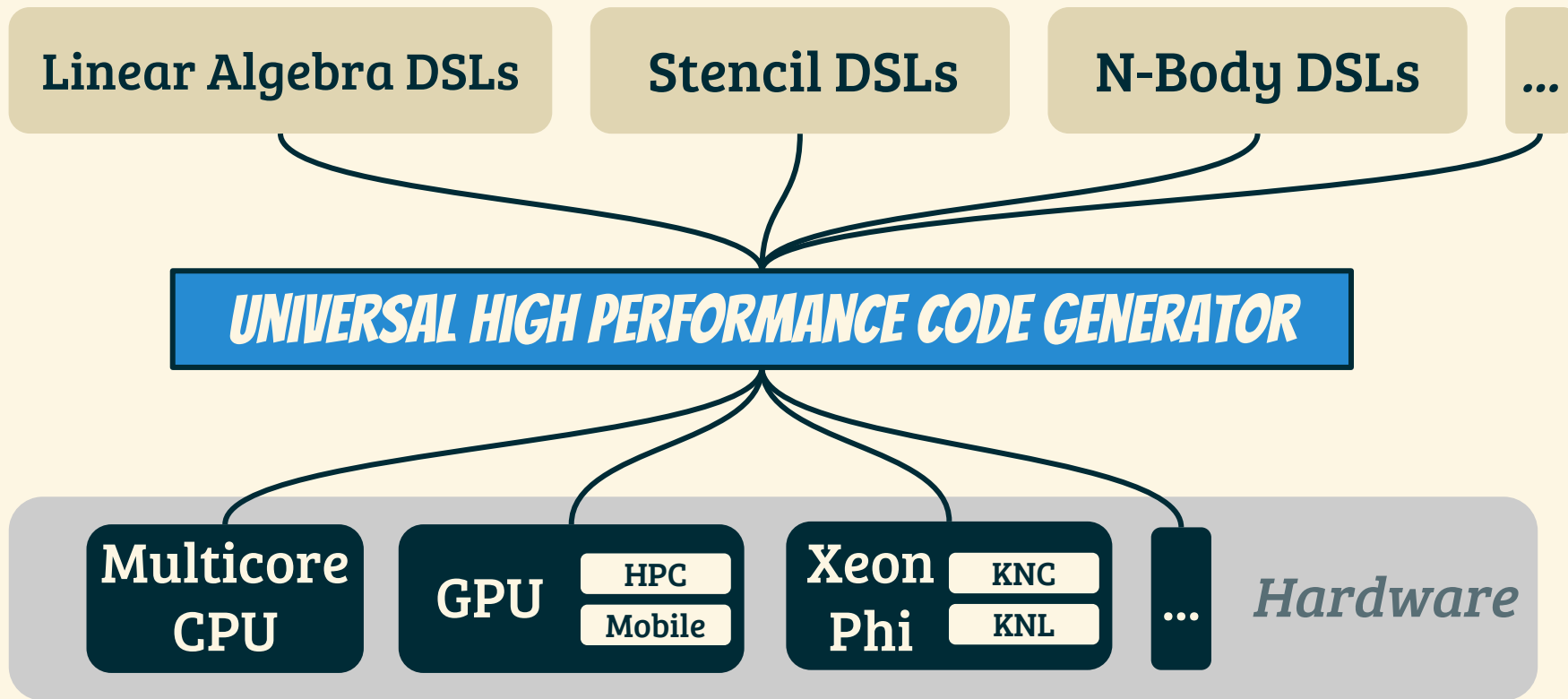# EXPLOITING DOMAIN KNOWLEDGE

**DSL**

PATUS | Pochoir | PARTANS | Halide | ...

**Hardware**

Multicore CPU | GPU [HPC, Mobile] | Xeon Phi [KNC, KNL] | ...

# EXPLOITING DOMAIN KNOWLEDGE

Stencil DSLs

Multicore CPU

GPU — HPC, Mobile

Xeon Phi — KNC, KNL

...

Hardware

# EXPLOITING DOMAIN KNOWLEDGE

# APPROACHING PERFORMANCE PORTABILITY

| Linear Algebra DSLs | Stencil DSLs | N-Body DSLs | ... |

**UNIVERSAL HIGH PERFORMANCE CODE GENERATOR**

Multicore CPU

GPU — HPC / Mobile

Xeon Phi — KNC / KNL

...

*Hardware*

# APPROACHING PERFORMANCE PORTABILITY

Linear Algebra DSLs    Stencil DSLs    N-Body DSLs    ...

**LIFT**

Multicore CPU    GPU [HPC | Mobile]    Xeon Phi [KNC | KNL]    ...    *Hardware*

LIFT

DSL          DSL          DSL

**Multicore CPU**    **GPU** HPC Mobile    **Xeon Phi** KNC KNL    ...    *Hardware*

# LIFT'S HIGH-LEVEL PRIMITIVES

*map(□→□)*

*reduce(⊕)*

*split(n)*

*join*

*zip*

# LIFT'S HIGH-LEVEL PRIMITIVES

map(□→■)

reduce(⊕)

split(n)

join

zip

dotproduct.lift

a   b

# LIFT'S HIGH-LEVEL PRIMITIVES

map(□→□)

reduce(⊕)

split(n)

join

zip

## dotproduct.lift

$zip$(a,b)

# LIFT'S HIGH-LEVEL PRIMITIVES

map(□→□)

reduce(⊕)

split(n)

join

zip

dotproduct.lift

a    b

map(*, zip(a,b))

# LIFT'S HIGH-LEVEL PRIMITIVES



map(□→□)

reduce(⊕)

split(n)

join

zip

dotproduct.lift

$$reduce(+,0,\ map(*,\ zip(a,b)))$$

# LIFT'S HIGH-LEVEL PRIMITIVES

*map(□→▨)*

*reduce(⊕)*

*split(n)*

*join*

*zip*

stencil.lift?
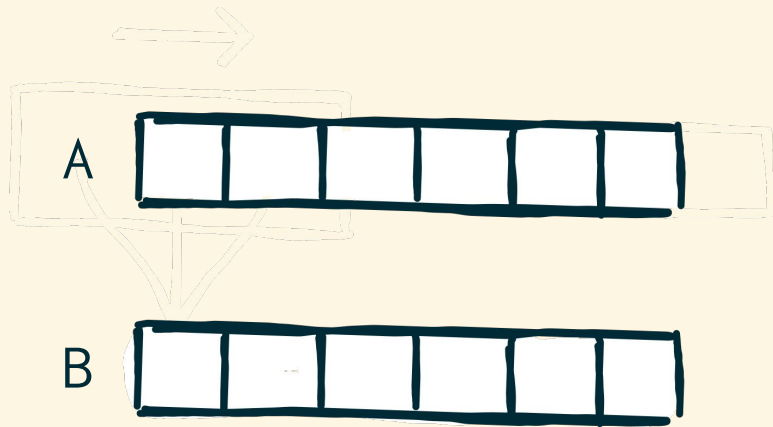
# Can we express stencil computations in Lift?

Let's look at a simple stencil example...

# WHAT ARE STENCIL COMPUTATIONS?

## 3-point-stencil.c

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```

A

B

# WHAT ARE STENCIL COMPUTATIONS?

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```

# WHAT ARE STENCIL COMPUTATIONS?

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```
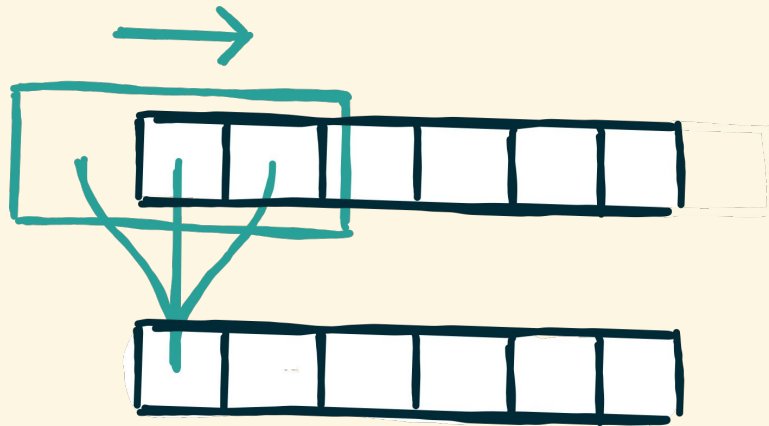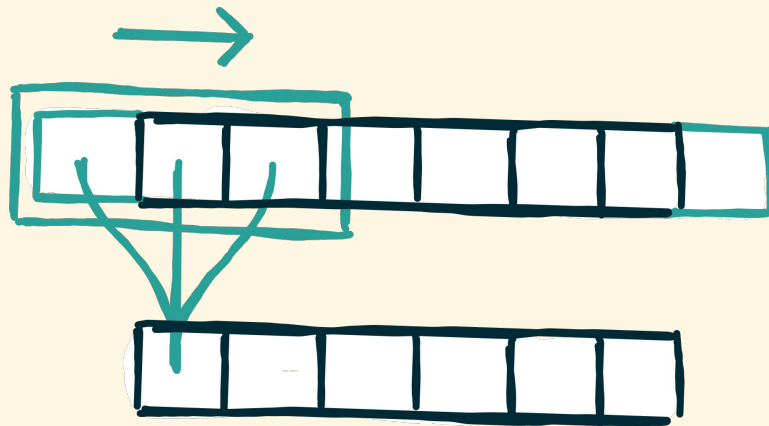
# STENCIL COMPUTATIONS IN LIFT
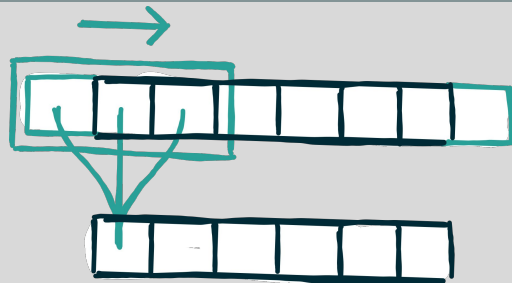
map(□→□)

reduce(⊕)

split(n)

join

zip

3-point-stencil.lift

# STENCIL COMPUTATIONS IN LIFT

$map(\square \rightarrow \blacksquare)$ 
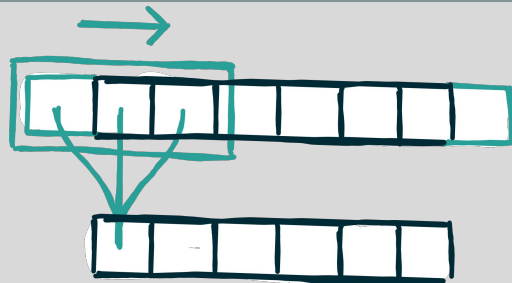
$reduce(\oplus)$ 

$split(n)$ 

$join$ 

$zip$ 

$stencil$ 

**3-point-stencil.lift**



Add specialized primitive:  Job done?

# STENCIL COMPUTATIONS IN LIFT

$map(\square \rightarrow \blacksquare)$

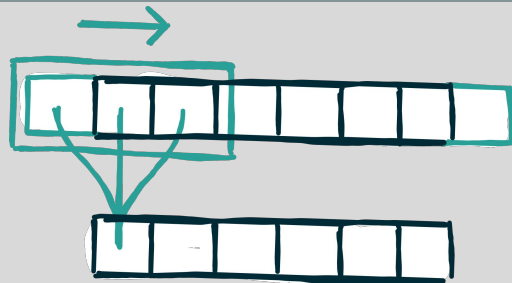$reduce(\oplus)$

$split(n)$

$join$

$zip$

$stencil$



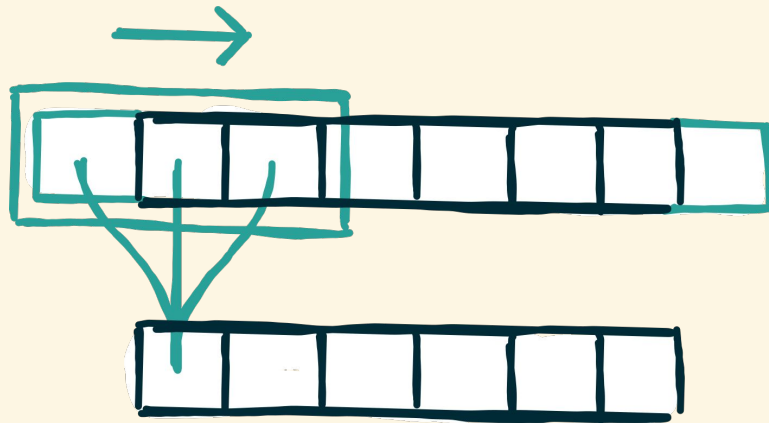3-point-stencil.lift



**Add specialized primitive: Job done?**

⊘ **No Reuse**
  *of existing primitives and optimizations*

⊘ **Domain-specific**
  *rather than generic*

⊘ **Multidimensional?**
  *is it composable?*

# DECOMPOSING STENCIL COMPUTATIONS

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```

# DECOMPOSING STENCIL COMPUTATIONS

### 3-point-stencil.c

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {    // ( a )
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```
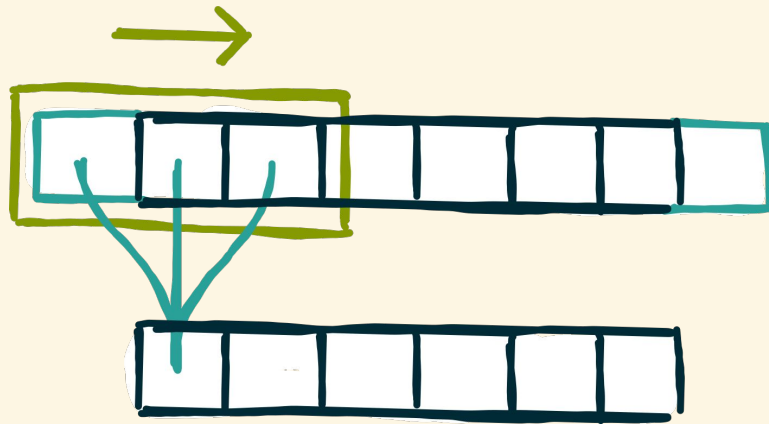
**(a)  access neighborhoods for every element**

# DECOMPOSING STENCIL COMPUTATIONS

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {    // ( a )
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;          // ( b )
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }
 B[ i ] = sum ; }
```
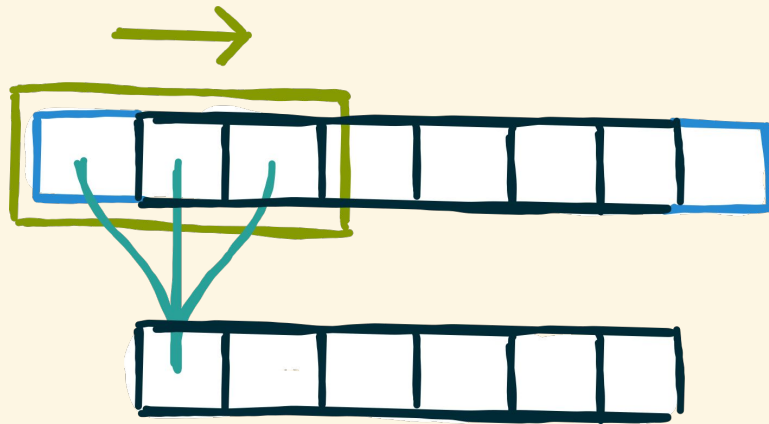


(a)   access **neighborhoods** for every element
(b)   specify **boundary handling**

# DECOMPOSING STENCIL COMPUTATIONS

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {    // ( a )
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;          // ( b )
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }                // ( c )
 B[ i ] = sum ; }
```
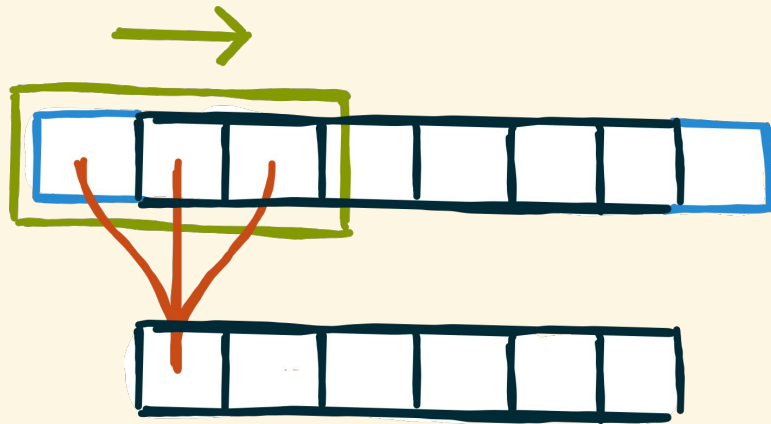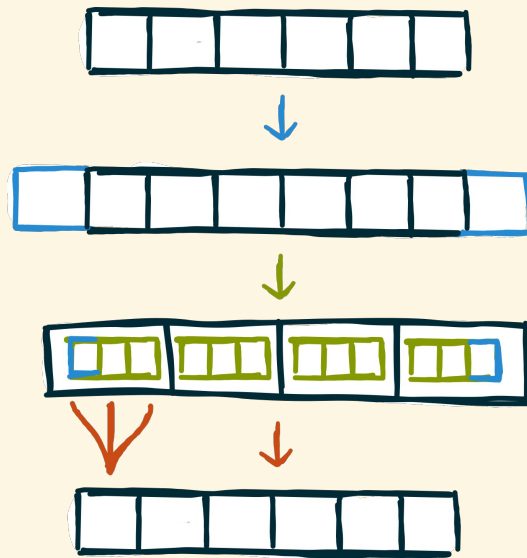
(a)  access **neighborhoods** for every element

(b)  specify **boundary handling**

(c)  apply **stencil function** to neighborhoods

# DECOMPOSING STENCIL COMPUTATIONS

**3-point-stencil.c**

```c
for (int i = 0; i < N ; i ++) {
    int sum = 0;
    for ( int j = -1; j <= 1; j ++) {    // ( a )
        int pos = i + j;
        pos = pos < 0 ? 0 : pos;          // ( b )
        pos = pos > N - 1 ? N - 1 : pos;
        sum += A[ pos ]; }                // ( c )
 B[ i ] = sum ; }
```

(a)  access **neighborhoods** for every element
(b)  specify **boundary handling**
(c)  apply **stencil function** to neighborhoods
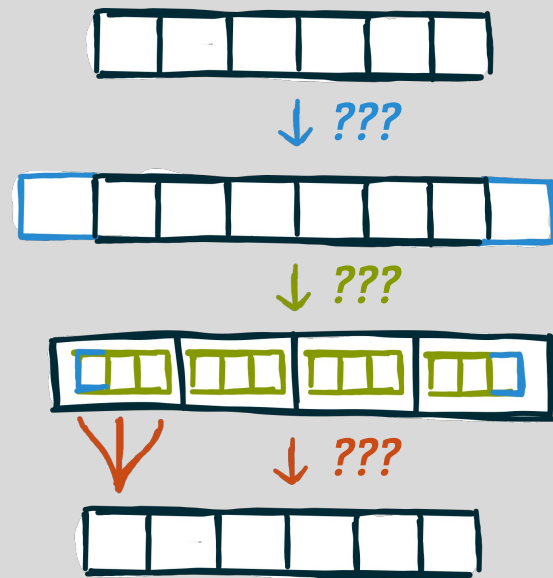
# STENCIL COMPUTATIONS IN LIFT

map(□→▢)  

reduce(⊕)  

split(n)  

join  

zip  

## 3-point-stencil.lift



↓ ???

↓ ???

↓ ???

# STENCIL COMPUTATIONS IN LIFT

$map(\square \rightarrow \blacksquare)$ 

$reduce(\oplus)$ 

$split(n)$ 

$join$ 

$zip$ 

## 3-point-stencil.lift



↓ ???

✔ **Reuse *map***
*allows to reuse existing rewrite rules*

↓ ???

✔ **Simplicity**
*one primitive per task*

✔ **Multidimensional**
*easily composable*

↓ *map*

# BOUNDARY HANDLING USING PAD

## pad ( reindexing )



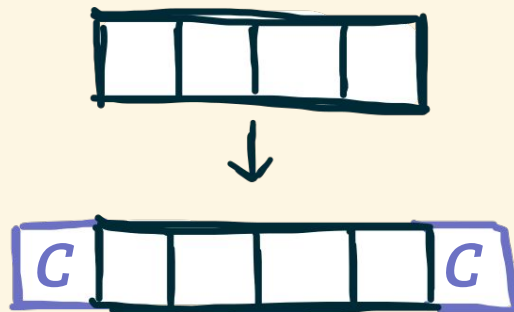### pad-reindexing.lift

```
clamp(i, n) = (i <  0) ? 0  :
               ((i >= n) ? n-1:i)


pad(1,1,clamp, [a,b,c,d]) =
    [a,a,b,c,d,d]
```
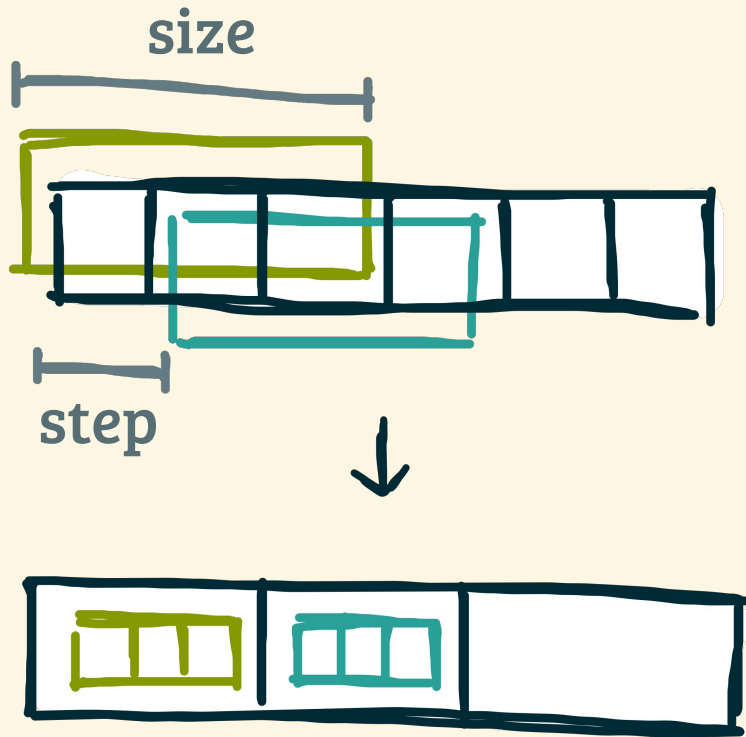
## pad ( constant )



### pad-constant.lift

```
constant(i, n) = C


pad(1,1,constant, [a,b,c,d]) =
    [C,a,b,c,d,C]
```
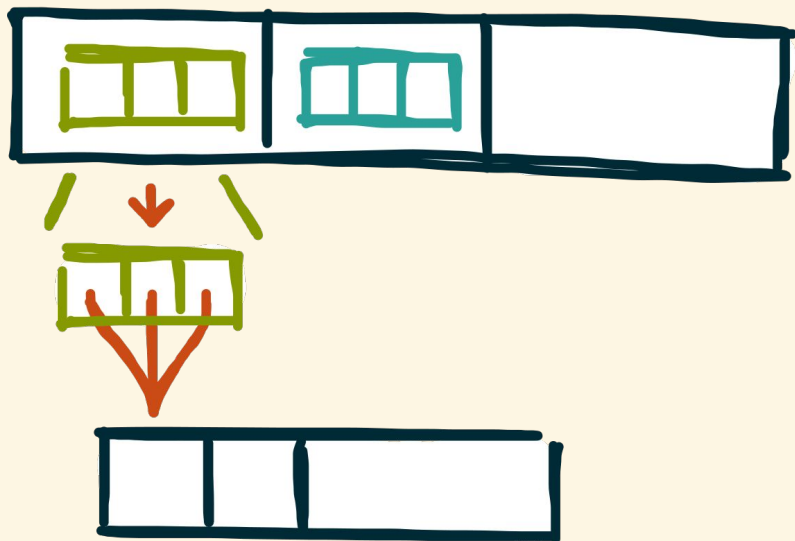
# NEIGHBORHOOD CREATION USING SLIDE



slide-example.lift

slide(3,1,[a,b,c,d,e]) =

[[a,b,c],[b,c,d],[c,d,e]]

# APPLYING STENCIL FUNCTION USING MAP



**sum-neighborhoods.lift**

```
map(nbh =>
    reduce(add, 0.0f, nbh))
```

# PUTTING IT TOGETHER

map(□→▨)  □□□ → ▨▨▨

reduce(⊕)  □□□ → □

split(n)  □□□□ → ▣▣

join  ▣▣ → □□□□

zip  ▣▣ / ▨▨ → ◄▨►◄▨►◄▨►

pad(l,r,b)  □□□ → ▨□□□▨

slide(n,s)  □□□□ → ▣▣

## stencil1D.lift

```
def stencil1D =
  fun(A =>
    map(reduce(add, 0.0f),
      slide(3,1,
        pad(1,1,clamp,A))))
```
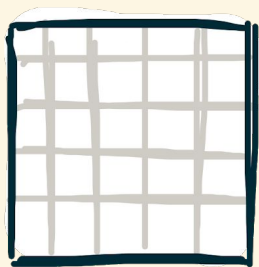
↓ pad

↓ slide

↓ map

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

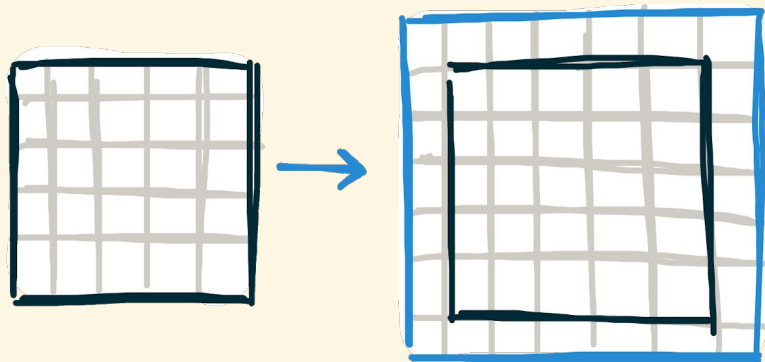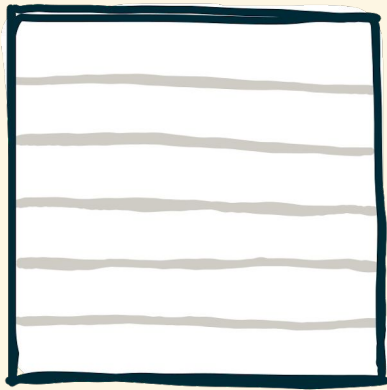# MULTIDIMENSIONAL STENCIL COMPUTATIONS

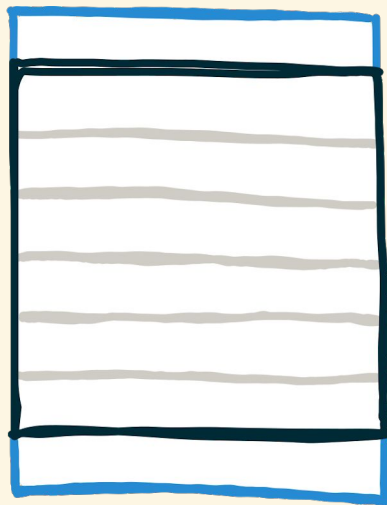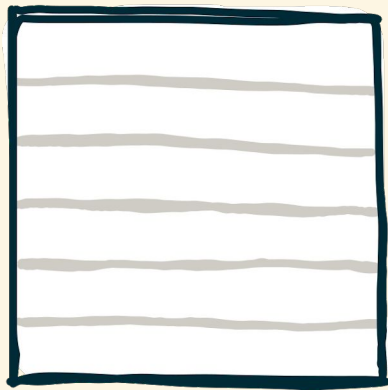are expressed as compositions of intuitive, generic 1D primitives

$pad_2(1,1,clamp,\text{input})$

# MULTIDIMENSIONAL BOUNDARY HANDLING USING PAD $_2$

input

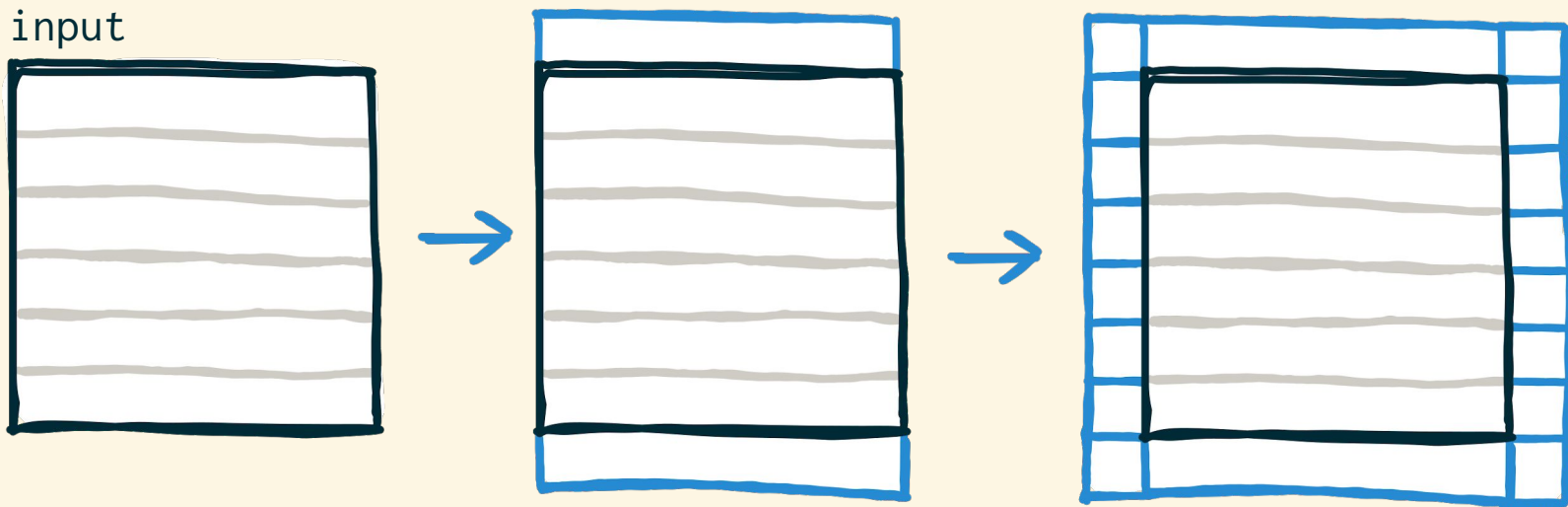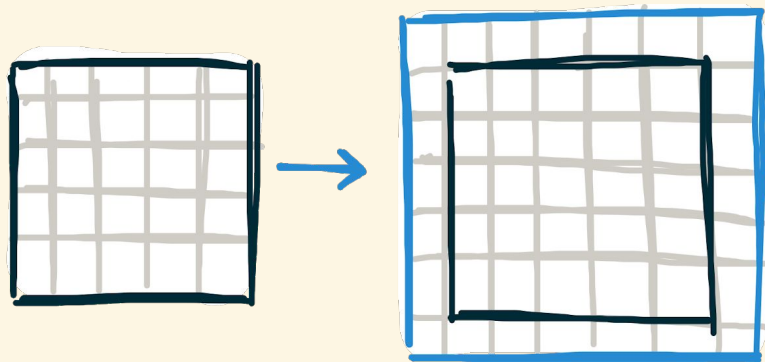$pad_2 =$

# MULTIDIMENSIONAL BOUNDARY HANDLING USING PAD $_2$

input



$$pad_2 = \qquad pad(\texttt{l,r,b,input})$$

# MULTIDIMENSIONAL BOUNDARY HANDLING USING PAD 2

input

$$pad_2 = map(pad(l,r,b,pad(l,r,b,input)))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives



$$pad_2(1,1,clamp,\text{input})$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

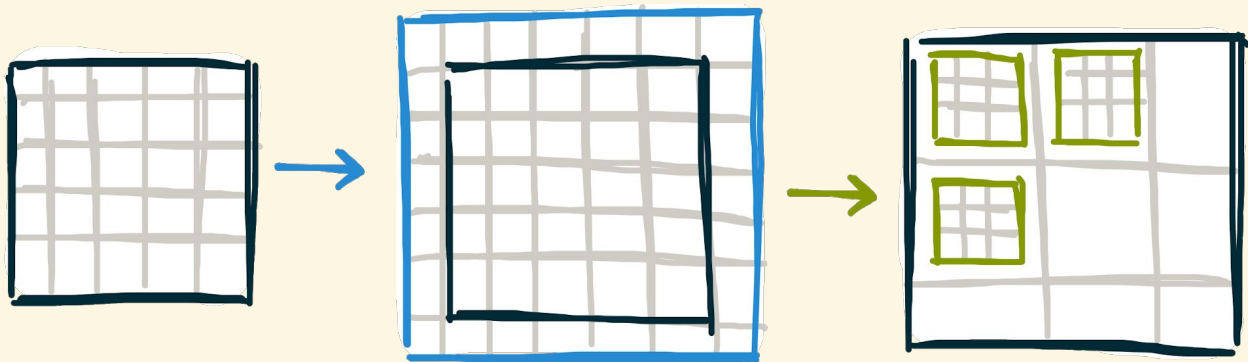are expressed as compositions of intuitive, generic 1D primitives
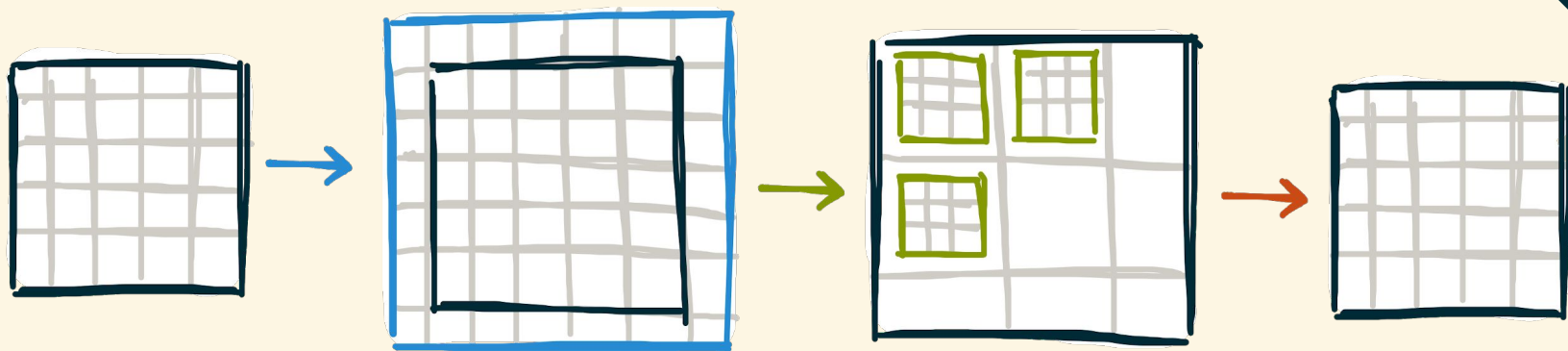
$$slide_2(3,1,\ pad_2(1,1,clamp,\texttt{input}))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

$$map_2(sum, \ slide_2(3,1, \ pad_2(1,1,clamp,\text{input})))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

**are expressed as compositions of intuitive, generic 1D primitives**

$$map_3(sum, \ slide_3(3,1, \ pad_3(1,1,clamp,input)))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

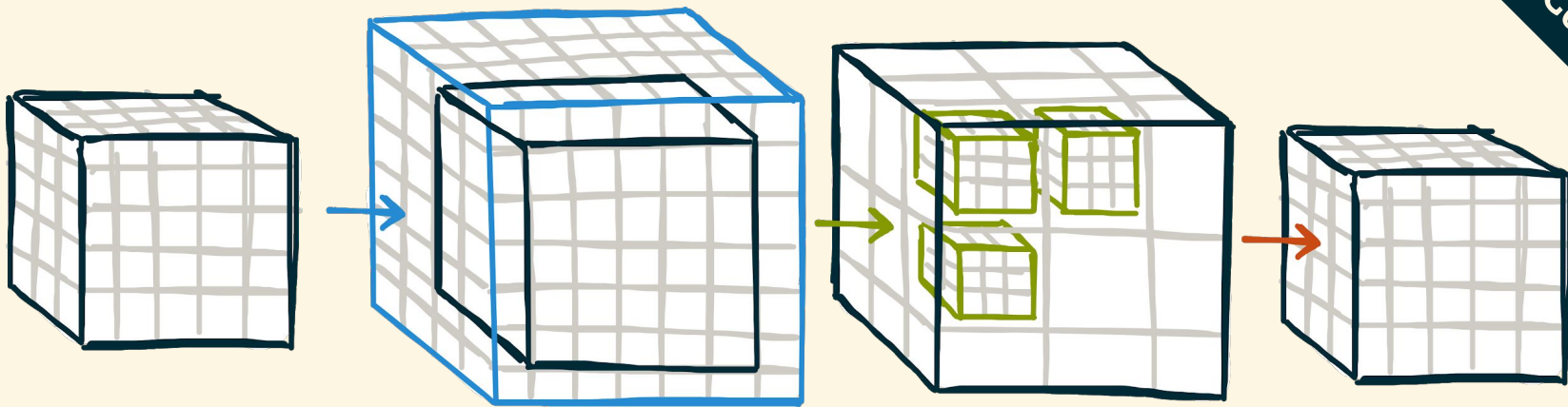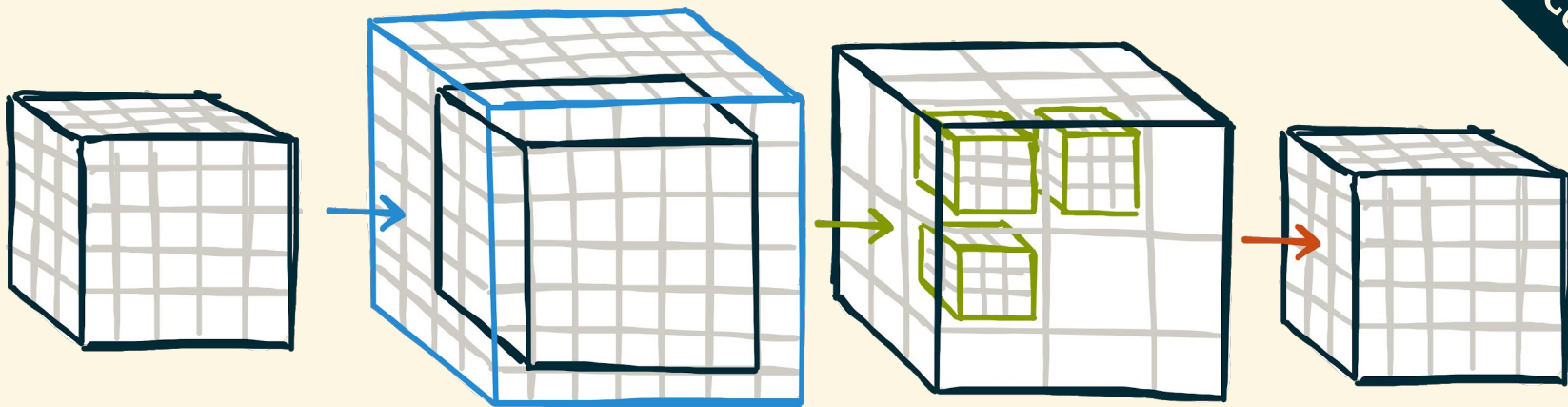$$map_3(sum, \; slide_3(3,1, \; pad_3(1,1,clamp,input)))$$

**Advantages:** ✓ Compact Language    ✓ Reuse Rewrites    ✓ Simple Compilation

# REUSING EXISTING REWRITE RULES

## Divide & Conquer

*map*(f, A)

# REUSING EXISTING REWRITE RULES

**Divide & Conquer**

$map(f, A) \mapsto join(map(map(f), split(n, A)))$

# OPTIMIZATION: OVERLAPPED TILING



✓ **Exploit Locality**

*Close neighborhoods share elements that can be grouped in tiles*

✓ **Shared Memory**

*Fast memory can be used to cache tiles*

# OPTIMIZATION: OVERLAPPED TILING
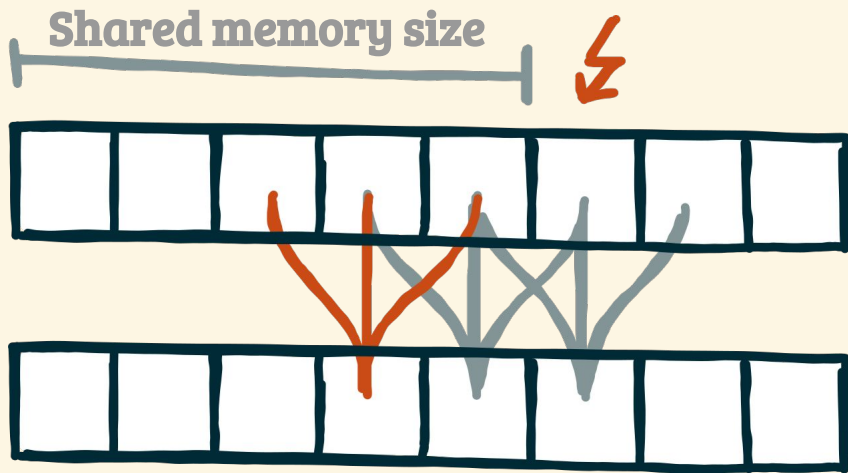


**Exploit Locality**

*Close neighborhoods share elements that can be grouped in tiles*

**Shared Memory**

*Fast memory can be used to cache tiles*

# OPTIMIZATION: OVERLAPPED TILING

overlap

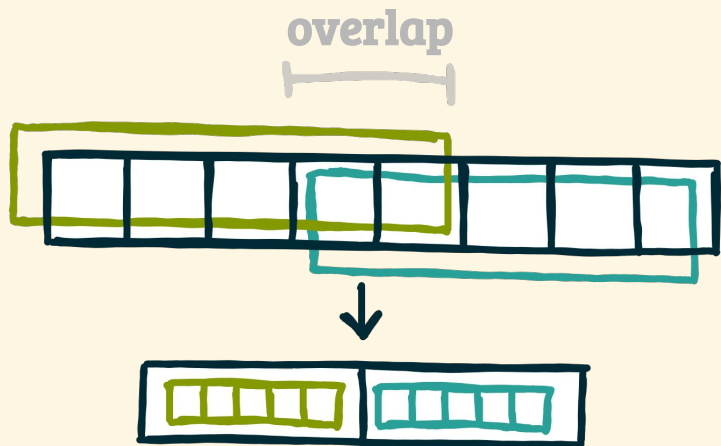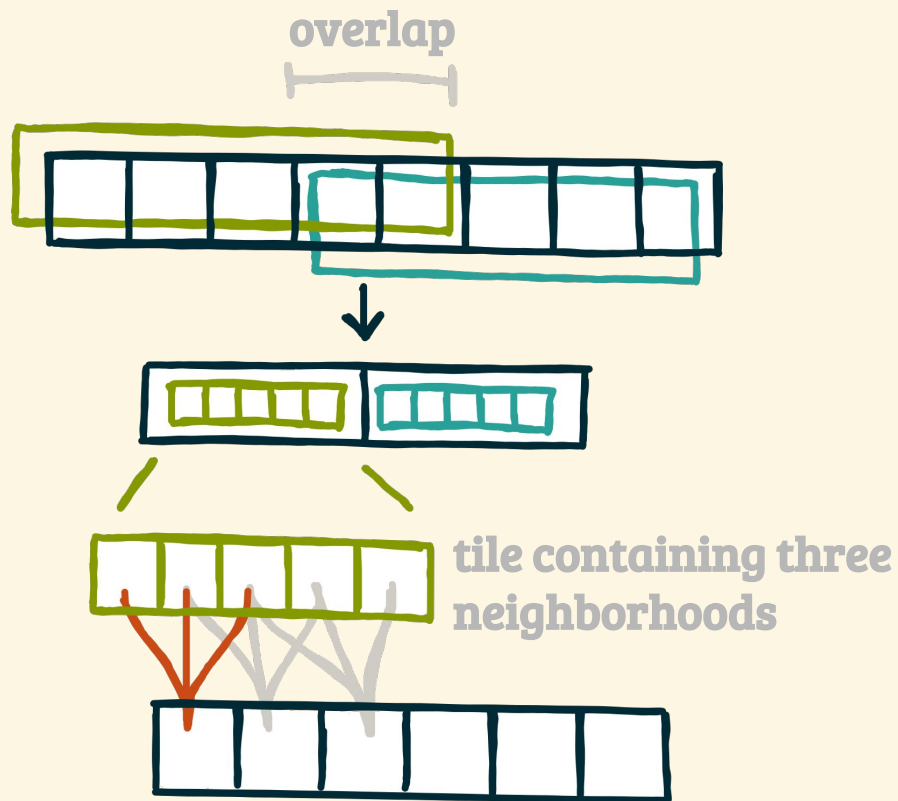## ✓ Exploit Locality

*Close neighborhoods share elements that can be grouped in tiles*

## ✓ Shared Memory

*Fast memory can be used to cache tiles*

# OPTIMIZATION: OVERLAPPED TILING



overlap

tile containing three neighborhoods

✓ **Exploit Locality**

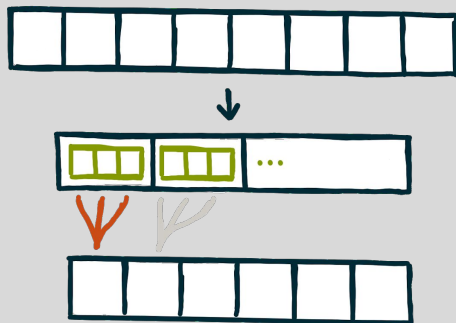*Close neighborhoods share elements that can be grouped in tiles*

✓ **Shared Memory**

*Fast memory can be used to cache tiles*
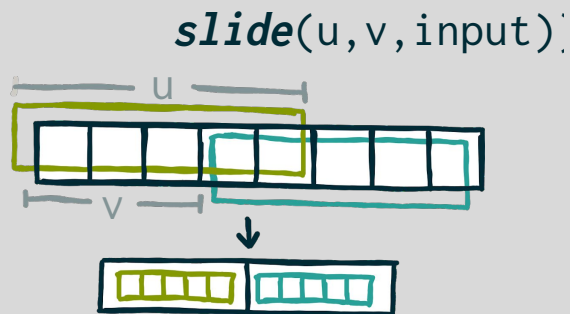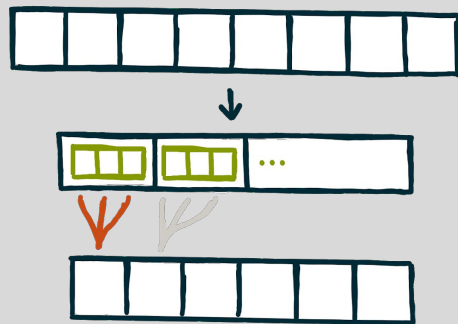
# OVERLAPPED TILING AS A REWRITE RULE

## overlapped tiling rule

`map(f, slide(3,1,input))`

# OVERLAPPED TILING AS A REWRITE RULE
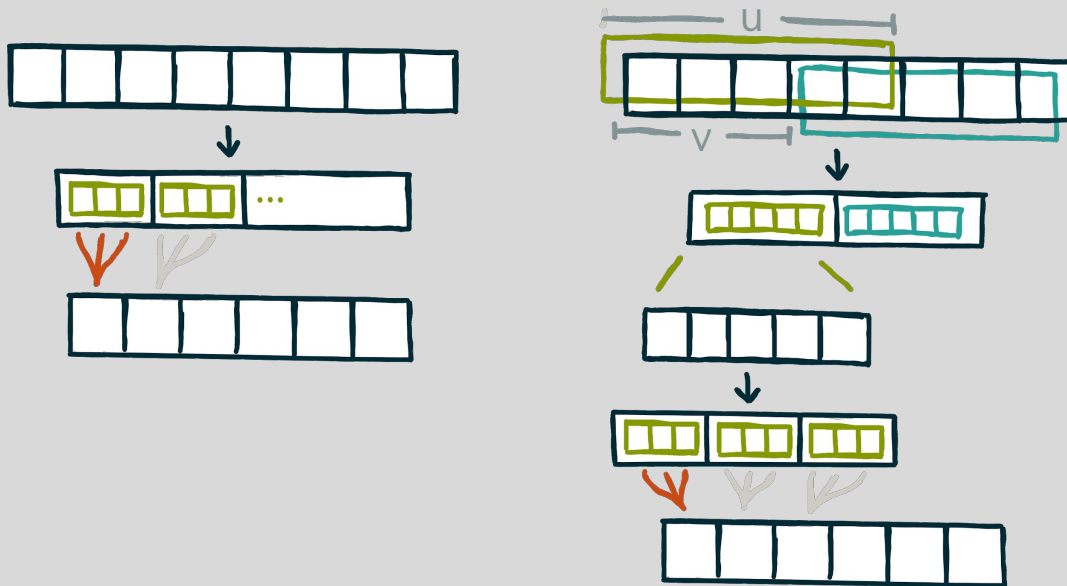


**overlapped tiling rule**

`map(f, slide(3,1,input))` ↦

`slide(u,v,input)`

# OVERLAPPED TILING AS A REWRITE RULE



## overlapped tiling rule

$$map(f, slide(3,1,input))$$

$$\mapsto$$

$$join(map(\text{tile} \Rightarrow$$
$$map(f, slide(3,1,\text{tile})),$$
$$slide(u,v,input)))$$

# EXPERIMENTAL EVALUATION



**14 Benchmarks**
*6 hand-optimized*
*8 polyhedral compilation*

DSL    DSL    DSL

**High-Level IR**

**Explore Optimizations**
**by rewriting**

**Low-Level Program**

**Multicore CPU**

**GPU**    HPC    Mobile

**Xeon Phi**    KNC    KNL

...    *Hardware*

# EXPERIMENTAL EVALUATION

DSL   DSL   DSL

**14 Benchmarks**
*6 hand-optimized*
*8 polyhedral compilation*

**High-Level IR**

**< 3h Exploration**
*per benchmark*

**Explore Optimizations**
**by rewriting**

**Low-Level Program**

**Multicore CPU**

**GPU**   HPC   Mobile

**Xeon Phi**   KNC   KNL

...   *Hardware*

# EXPERIMENTAL EVALUATION

DSL     DSL     DSL

**14 Benchmarks**
*6 hand-optimized*
*8 polyhedral compilation*

**High-Level IR**

**< 3h Exploration**
*per benchmark*

**Explore Optimizations**
**by rewriting**

**up to 20 algorithmically**
**different variants**
*+ auto-tuning of numerical parameters*

**Low-Level Program**

| Multicore CPU | GPU | HPC | Xeon Phi | KNC | ... | *Hardware* |
|---|---|---|---|---|---|---|
| | | Mobile | | KNL | | |

# EXPERIMENTAL EVALUATION

DSL    DSL    DSL

**14 Benchmarks**
*6 hand-optimized*
*8 polyhedral compilation*

**High-Level IR**

**< 3h Exploration**
*per benchmark*

**Explore Optimizations**
**by rewriting**

**up to 20 algorithmically**
**different variants**
*+ auto-tuning of numerical parameters*

**Low-Level Program**

**3 GPU Architectures**
*2 Desktop GPUs*
*1 Mobile GPU*

Multicore CPU

**GPU**
HPC
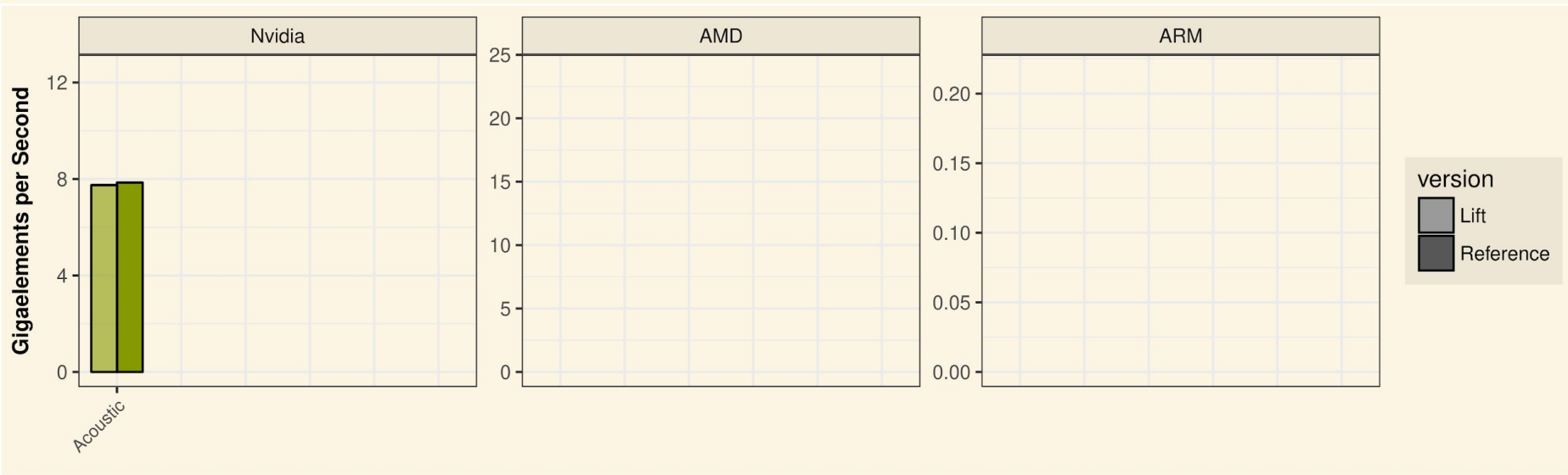**Mobile**

**Xeon Phi**
KNC
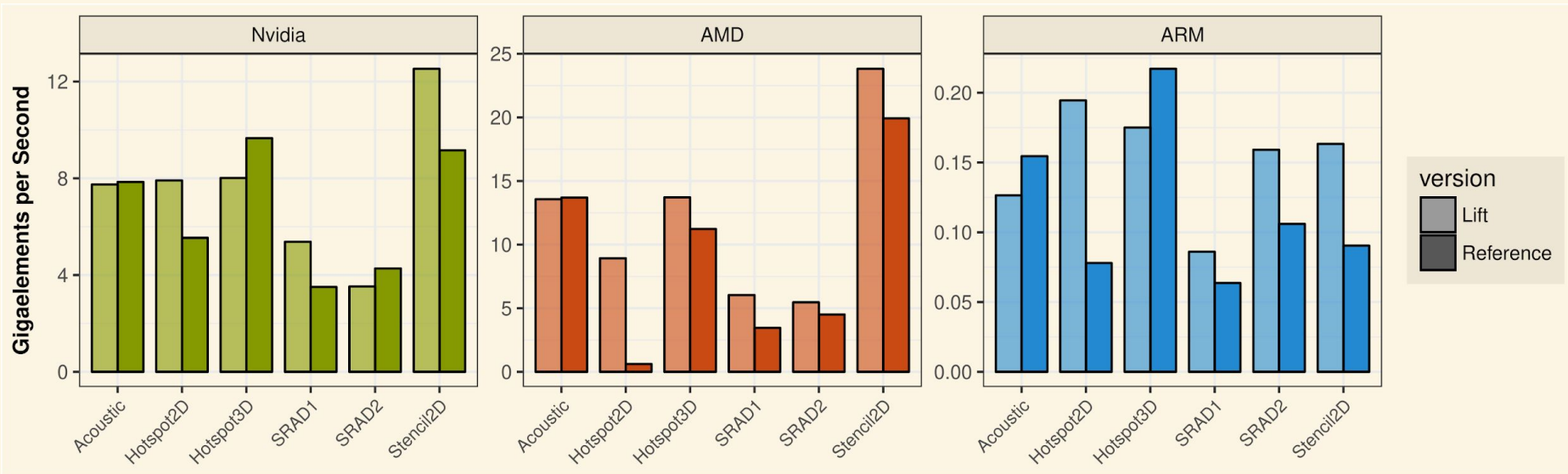KNL

...

*Hardware*

# COMPARISON WITH HAND-OPTIMIZED CODES

## higher is better
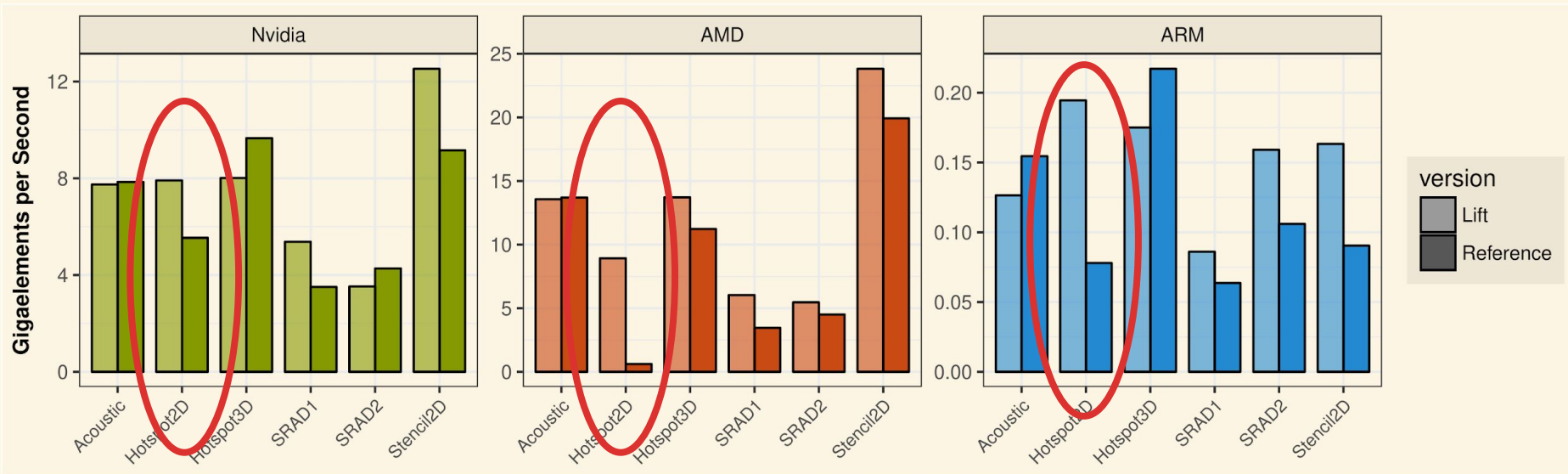
# COMPARISON WITH HAND-OPTIMIZED CODES

## higher is better



**Lift achieves the same performance
as hand optimized code**
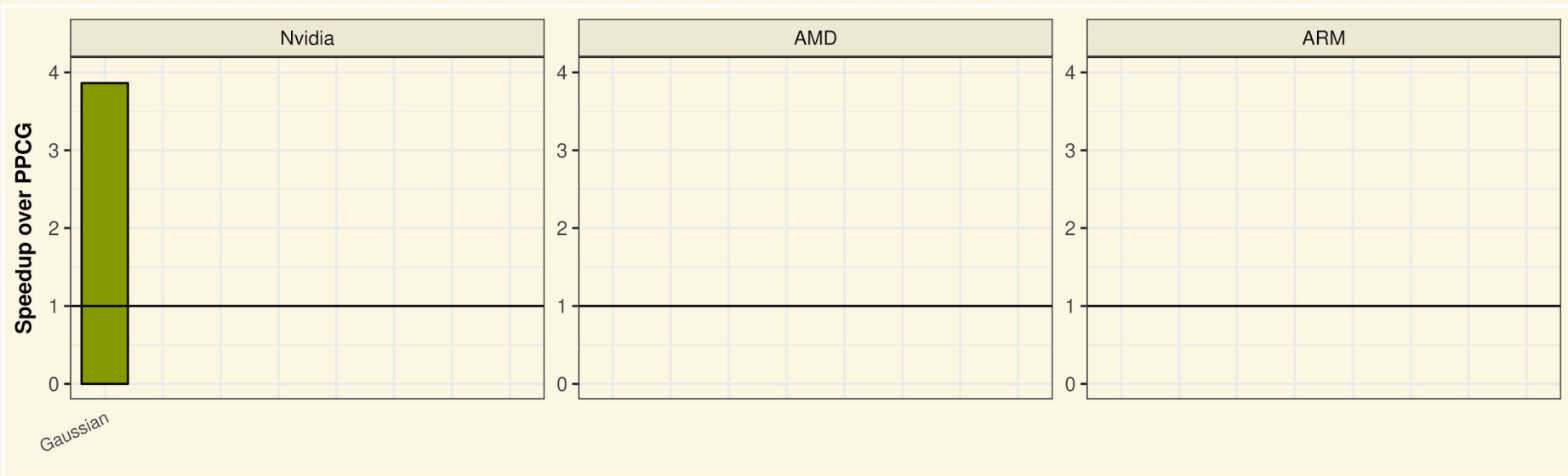
# COMPARISON WITH HAND-OPTIMIZED CODES

higher is better

**Lift achieves the same performance
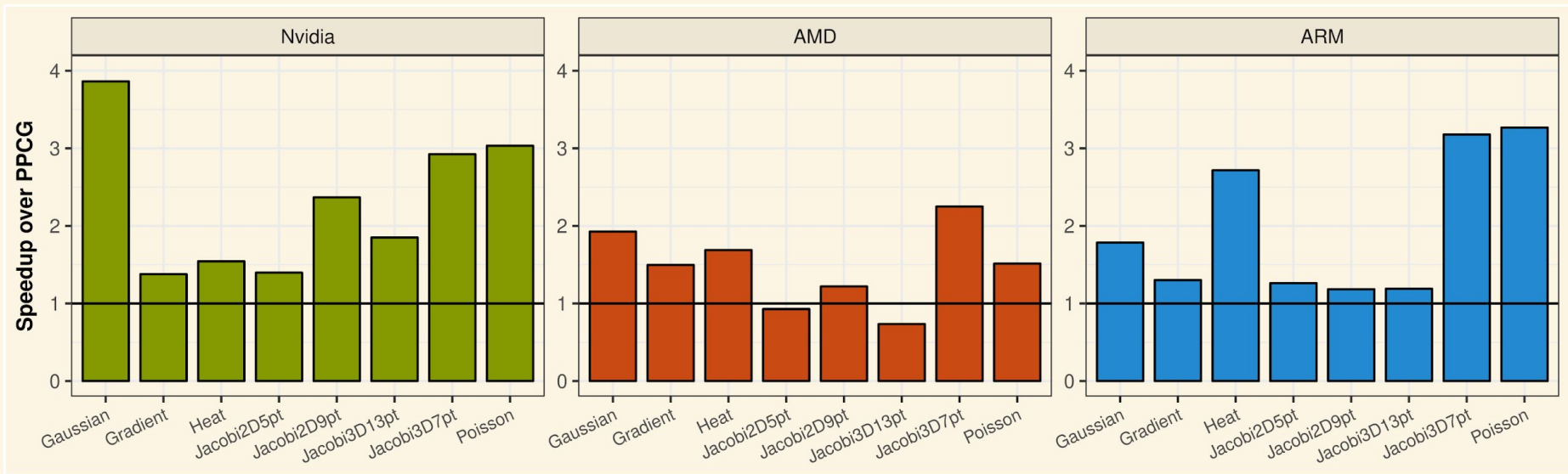as hand optimized code**

# COMPARISON WITH POLYHEDRAL COMPILATION

higher is better

# COMPARISON WITH POLYHEDRAL COMPILATION

higher is better



**Lift outperforms state-of-the-art optimizing compilers**

# A STEP FURTHER
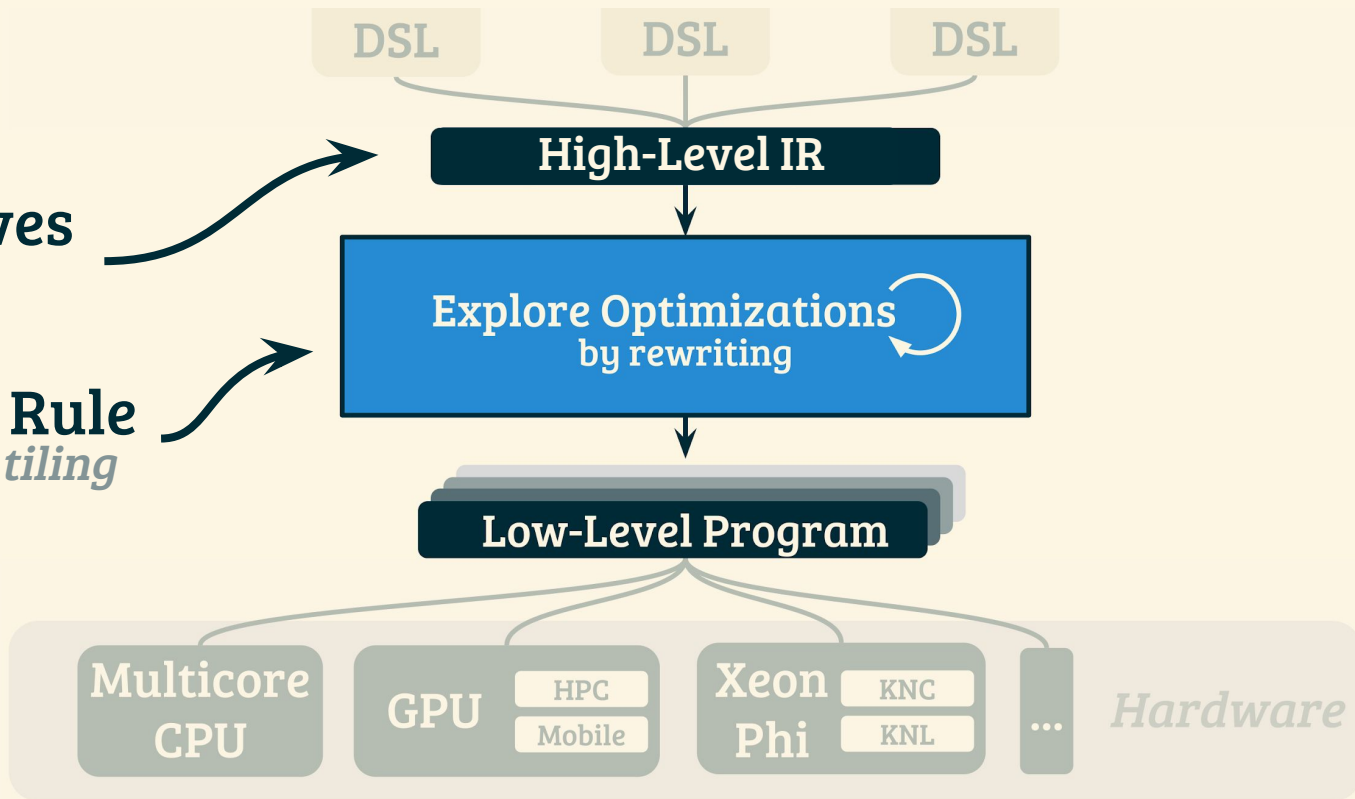
**towards multi-kernel stencil computations**



```
conv

def conv =
 fun(A =>
   map2D(multWeights,
     slide2D(3,1,
       pad2D(1,1,clamp,A)))))
```
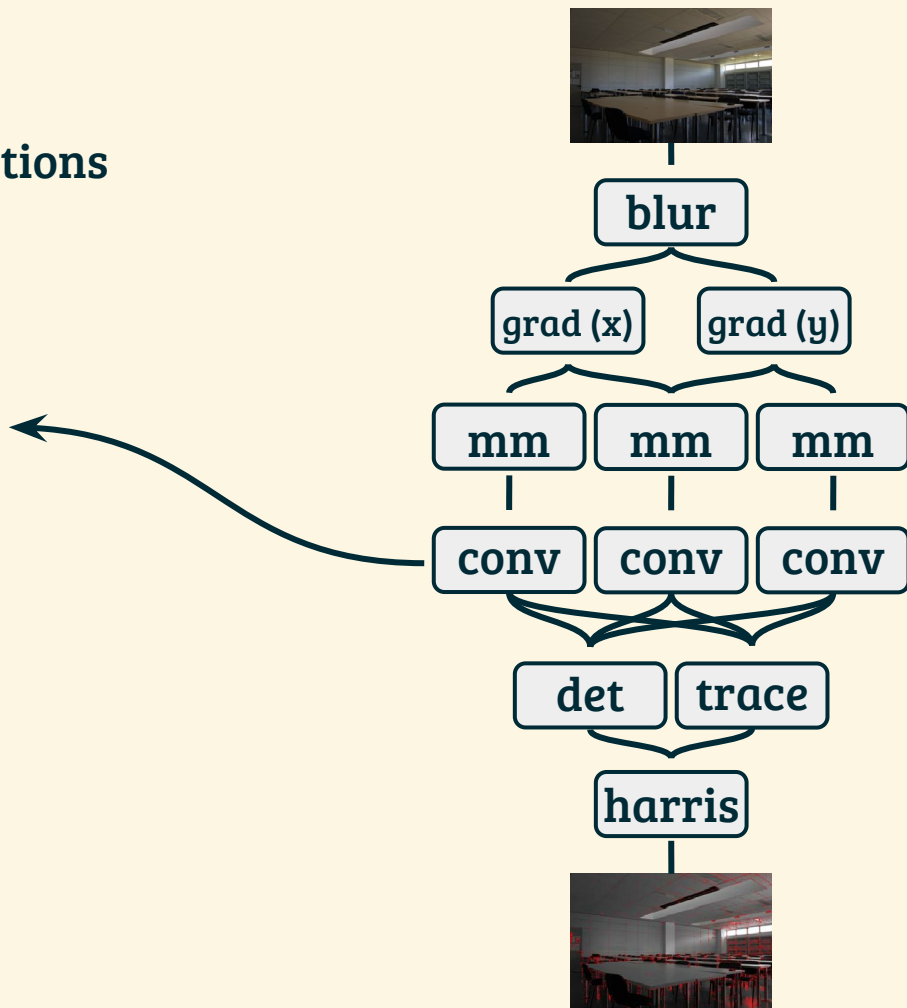
# A STEP FURTHER
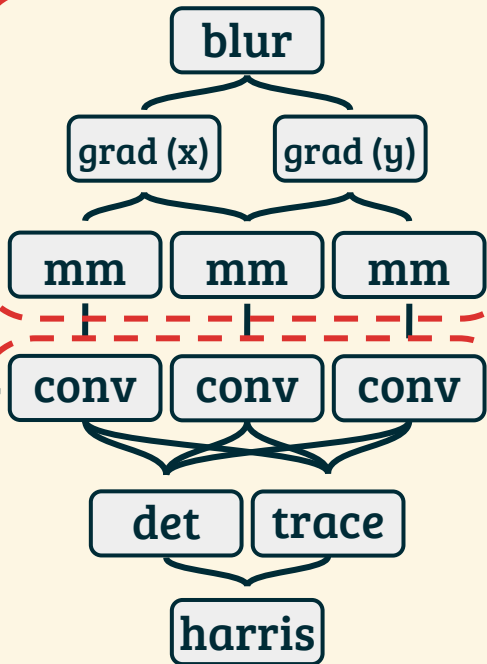## towards multi-kernel stencil computations

**conv**

```
def conv =
 fun(A =>
   map2D(multWeights,
     slide2D(3,1,
       pad2D(1,1,clamp,A)))))
```



stage 1

blur
grad (x)    grad (y)
mm    mm    mm

stage 2

conv    conv    conv
det    trace
harris

# A STEP FURTHER

towards multi-kernel stencil computations



```
conv

 def conv =
  fun(A =>
    map2D(multWeights,
      slide2D(3,1,
        pad2D(1,1,clamp,A)))))
```
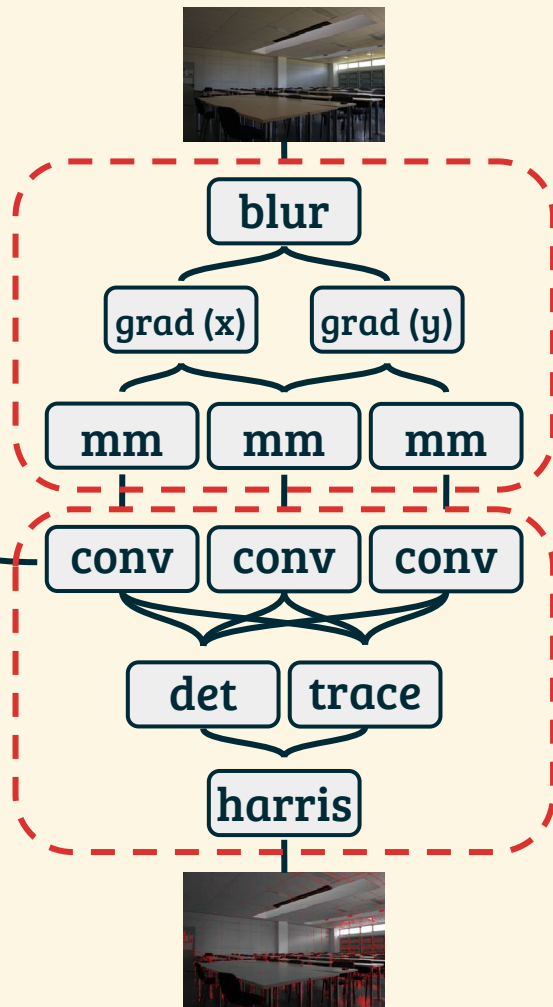
**Kernel Fusion via Rewriting**

**Guided Search Space Exploration**

# LIFT IS OPEN SOURCE!

more info at:

# lift-project.org

❝ Paper    CGO Artifact    Source Code

🏆 Best Paper Award (CGO'18)

Bastian Hagedorn: b.hagedorn@wwu.de

TO BE CONTINUED...