# 1968 - Go To Statement Considered Harmful



EWD215 - 0

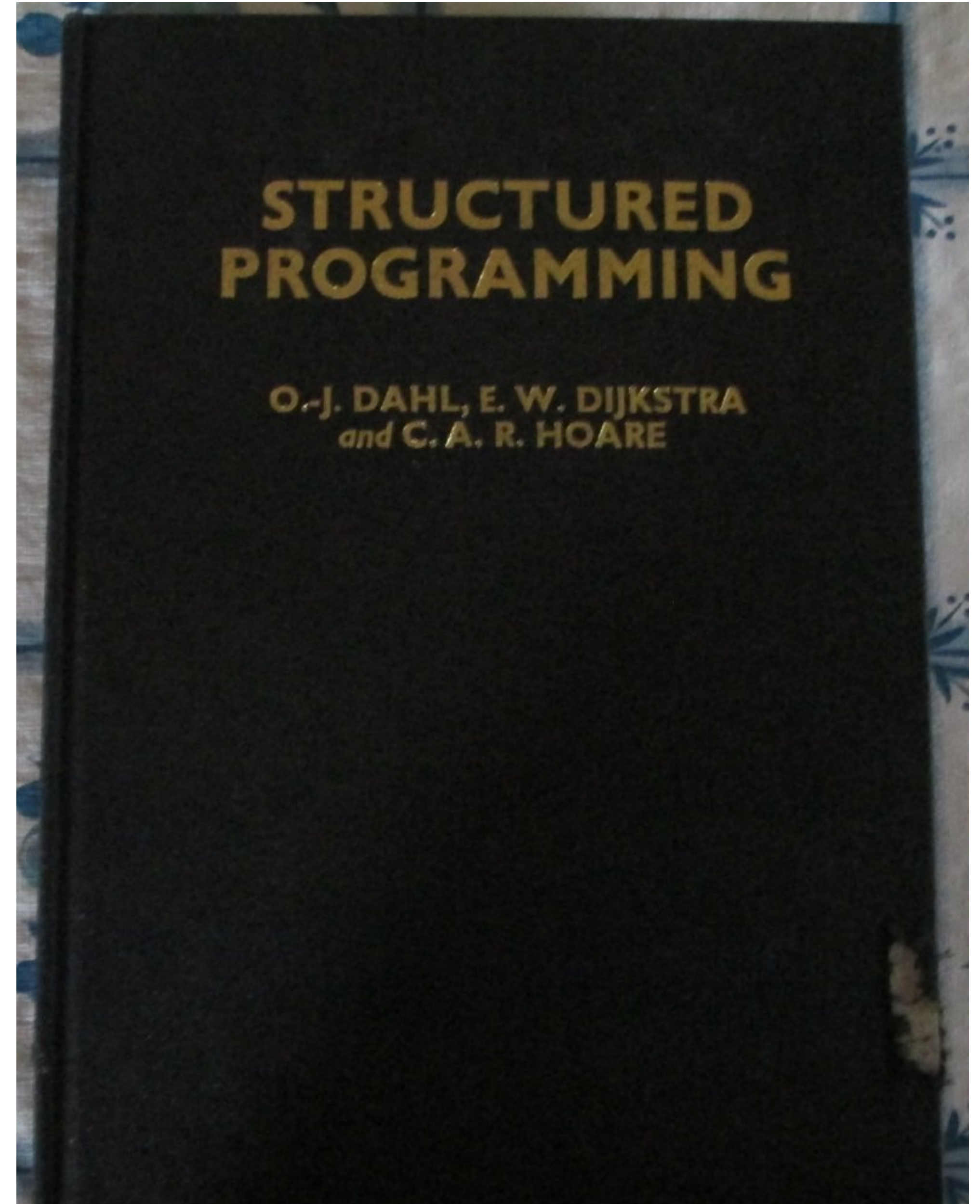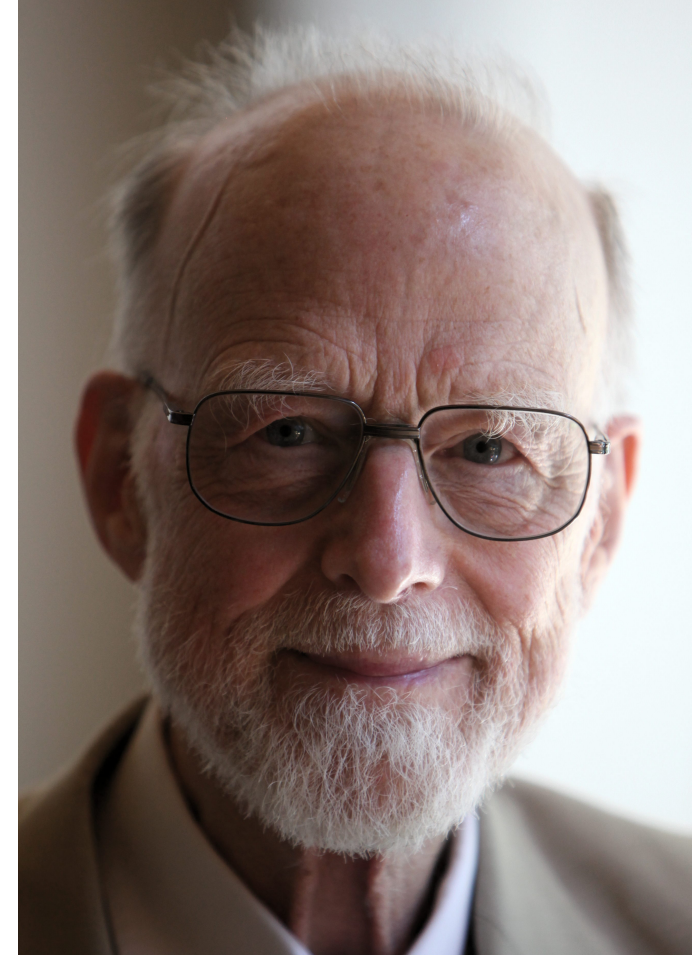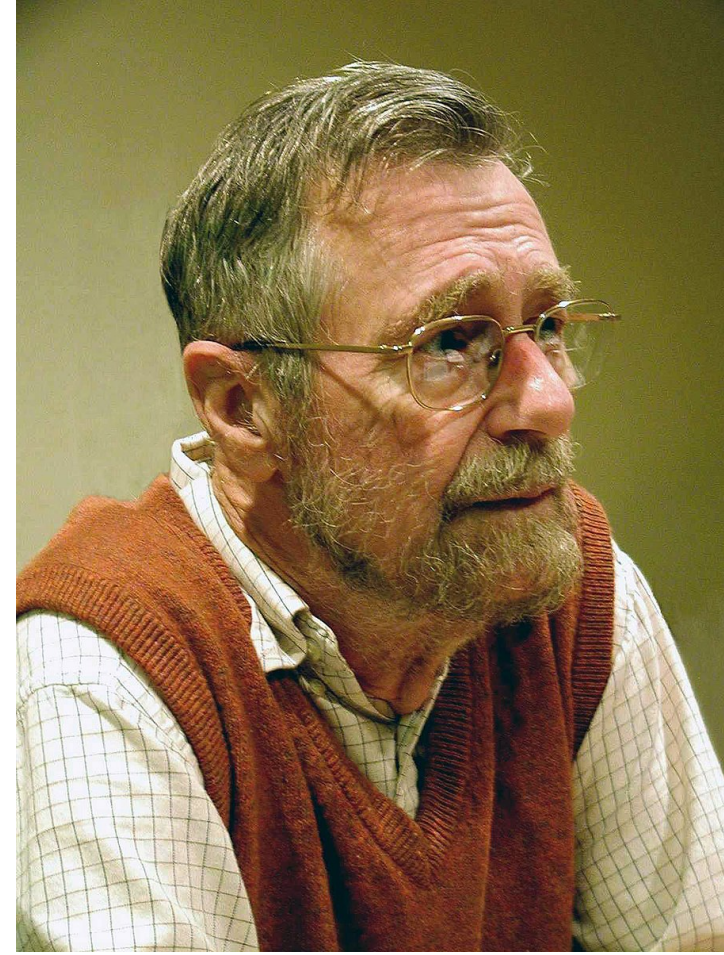A Case against the GO TO Statement.

by Edsger W.Dijkstra
Technological University
Eindhoven, The Netherlands

Since a number of years I am familiar with the observation that the quality of programmers is a decreasing function of the density of go to statements in the programs they produce. Later I discovered why the use of the go to statement has such disastrous effects and did I become convinced that the go to statement should be abolished from all "higher level" programming languages (i.e. everything except -perhaps- plain machine code). At that time I did not attach too much importance to this discovery; I now submit my considerations for publication because in very recent discussions in which the subject turned up, I have been urged to do so.

My first remark is that, although the programmer's activity ends when he has constructed a correct program, the process taking place under control

# 1972 - Structured Programming





STRUCTURED PROGRAMMING

O.-J. DAHL, E. W. DIJKSTRA
and C. A. R. HOARE

```
procedure matmult (A, B, C, m, n, p);
        array A, B, C; integer m, n, p;
begin integer i, j, k;
            for i := 1 step 1 until m do
            for j := 1 step 1 until n do
            begin C[i, j] := 0;
                    for k := 1 step 1 until p do
                    C[i, j] := C[i, j] + A[i, k] × B[k, i]
            end
end;
```

# 1989 - Structured *Parallel* Programming



$$D\_C\ indivisible\ split\ join\ f\ =\ F$$
$$where\ F\ P\ =\quad f\ P,\ \ if\ indivisible\ P$$
$$=\quad join\,(map\ F\ (split\ P))\,,\ otherwise$$

Research Monographs in
Parallel and Distributed Computing

MURRAY COLE

Algorithmic
Skeletons:
Structured
Management
of Parallel
Computation

4

# 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

a     b     →

$zip$(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

$$reduce(+, 0, \textbf{\textit{map}}(*, \textbf{\textit{zip}}(a,b)))$$

# LIFT'S HIGH-LEVEL PRIMITIVES

*map(□→▪)* 

*reduce(⊕)* 

*split(n)* 

*join* 

*zip* 

## matrixMult.lift

```
map(λ rowA ↦
  map(λ colB ↦
    dotProduct(rowA, colB)
  , transpose(B))
, A)
```

# IMPLEMENTATION CHOICES AS REWRITE RULES

## Divide & Conquer

*map*(f, A)
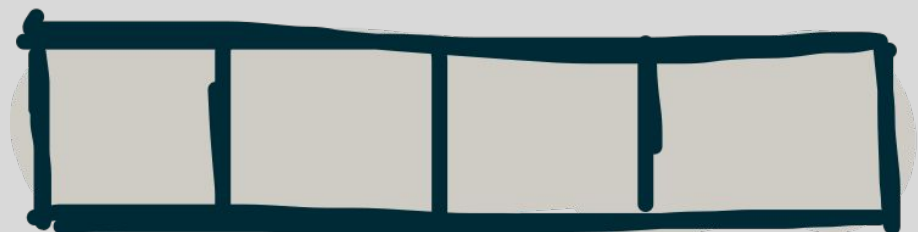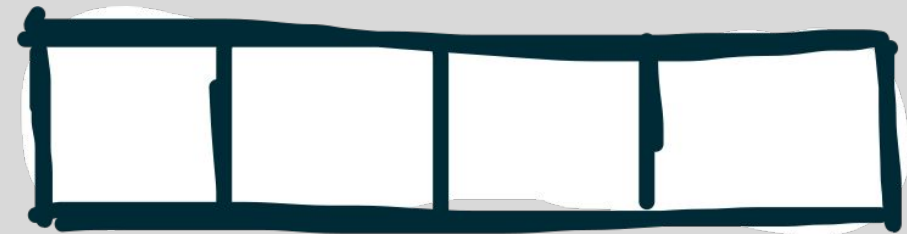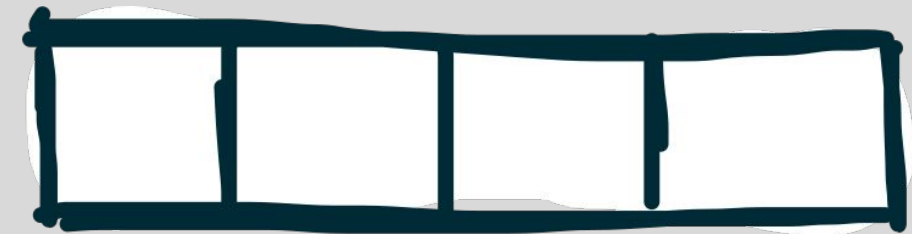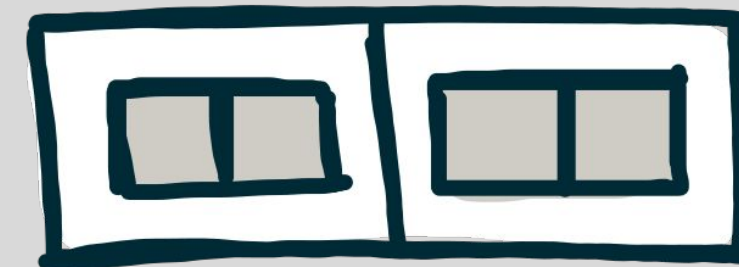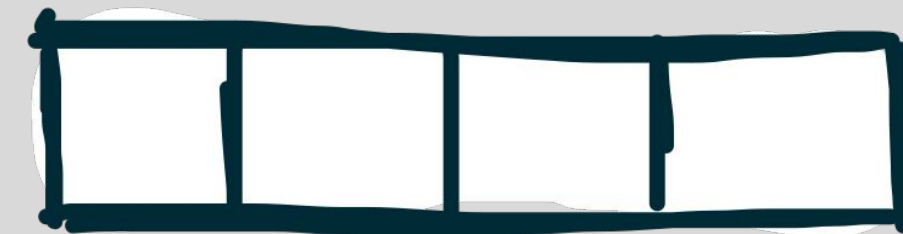
# IMPLEMENTATION CHOICES AS REWRITE RULES

## Divide & Conquer

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

# LIFT'S LOW LEVEL (OPENCL) PRIMITIVES

| Lift primitive | OpenCL concept |
|---|---|
| mapGlobal | Work-items |
| mapWorkgroup mapLocal | Work-groups |
| mapSeq reduceSeq | Sequential implementations |
| toLocal, toGlobal | Memory areas |
| mapVec, splitVec, joinVec | Vectorisation |

# REWRITING INTO OPENCL

Map rules:

map f ↦ mapGlobal f | mapWorkgroup f | mapLocal f | mapSeq f

Local / global memory:

mapLocal f ↦ toLocal (mapLocal f)     mapLocal f ↦ toGlobal (mapLocal f)

Vectorization:

map f ↦ joinVec ∘ map (mapVec f) ∘ splitVec n

# OPTIMIZATIONS AS MACRO RULES
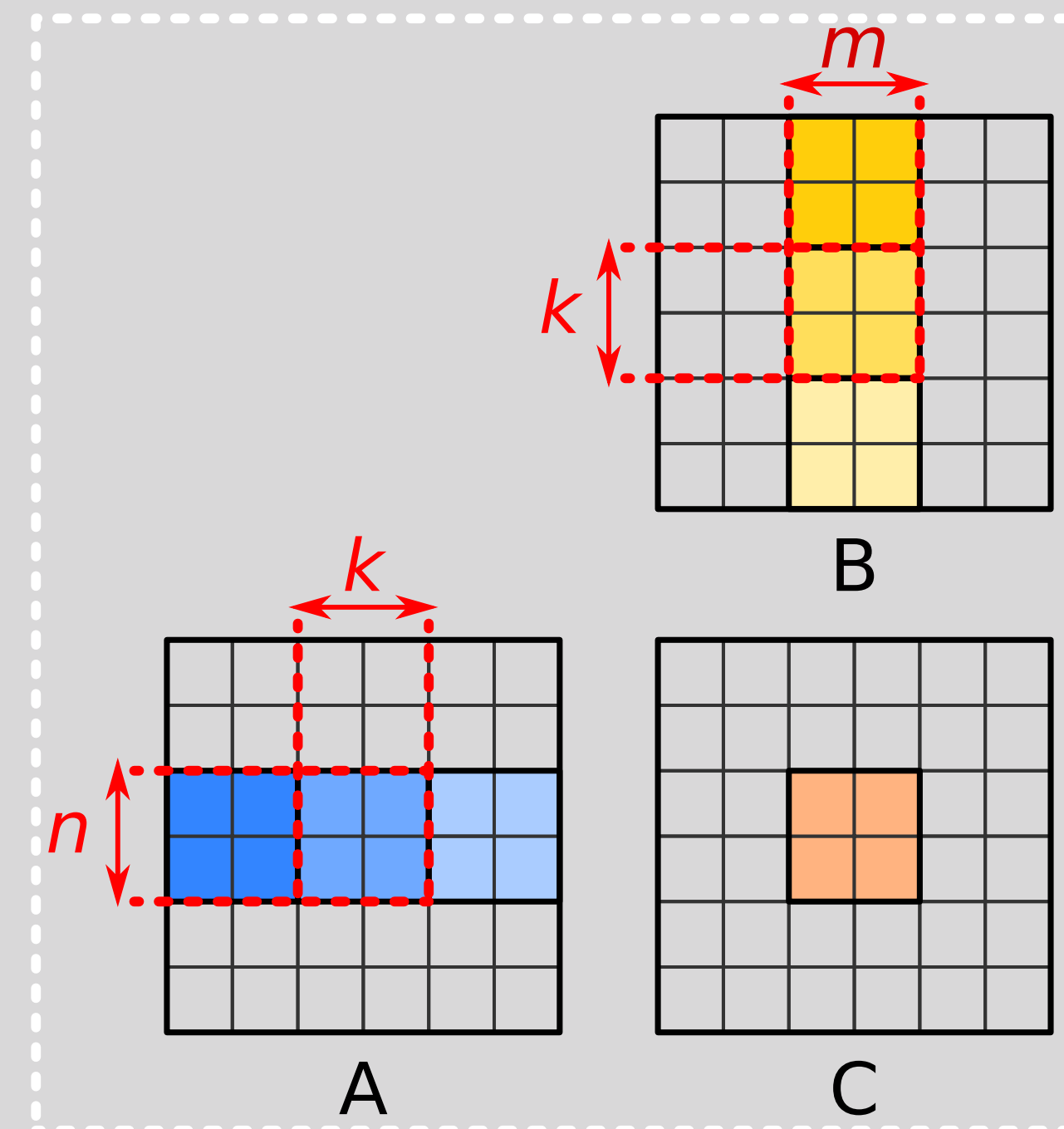
## 2D Tiling

Naïve matrix multiplication

```
1  map(λ arow .
2    map(λ bcol .
3      reduce(+, 0) ∘ map(×) ∘ zip(arow, bcol)
4    , transpose(B))
5  , A)
```

⬇ Apply tiling rules

```
1   untile ∘ map(λ rowOfTilesA .
2     map(λ colOfTilesB .
3       toGlobal(copy2D) ∘
4       reduce(λ (tileAcc, (tileA, tileB)) .
5         map(map(+)) ∘ zip(tileAcc) ∘
6         map(λ as .
7           map(λ bs .
8             reduce(+, 0) ∘ map(×) ∘ zip(as, bs)
9           , toLocal(copy2D(tileB)))
10          , toLocal(copy2D(tileA)))
11        ,0, zip(rowOfTilesA, colOfTilesB))
12    ) ∘ tile(m, k, transpose(B))
13  ) ∘ tile(n, k, A)
```



[GPGPU'16]

# EXPLORATION BY REWRITING
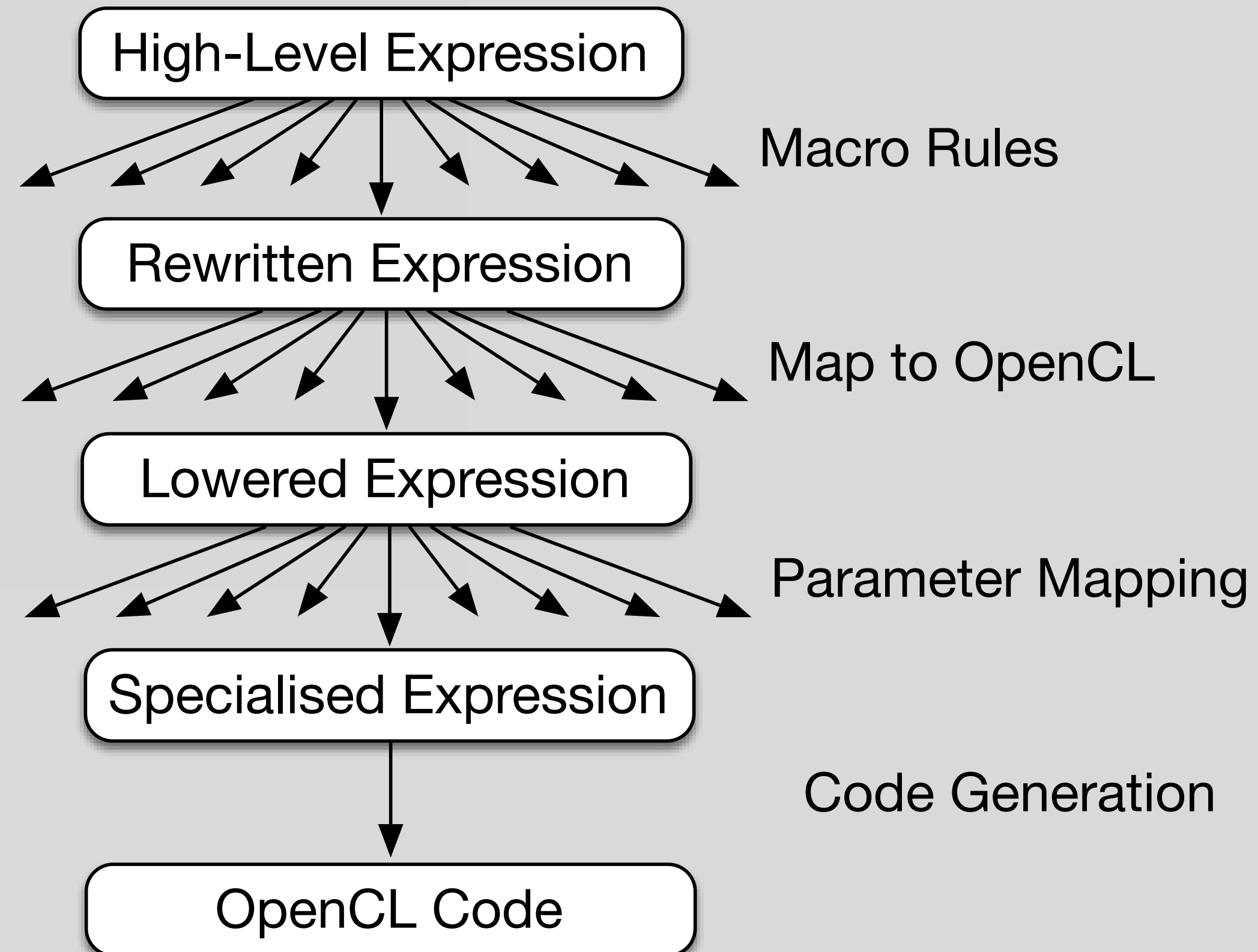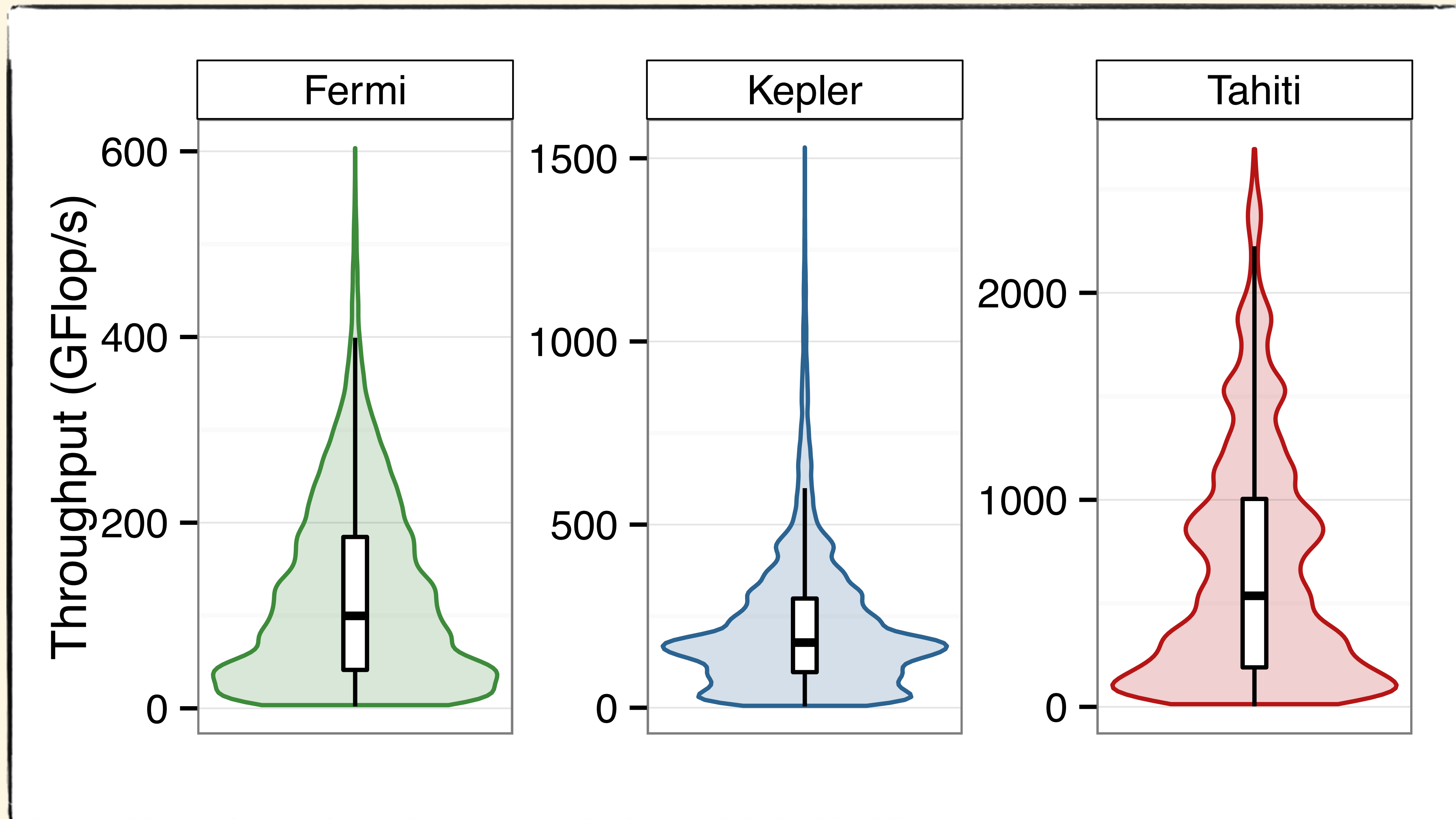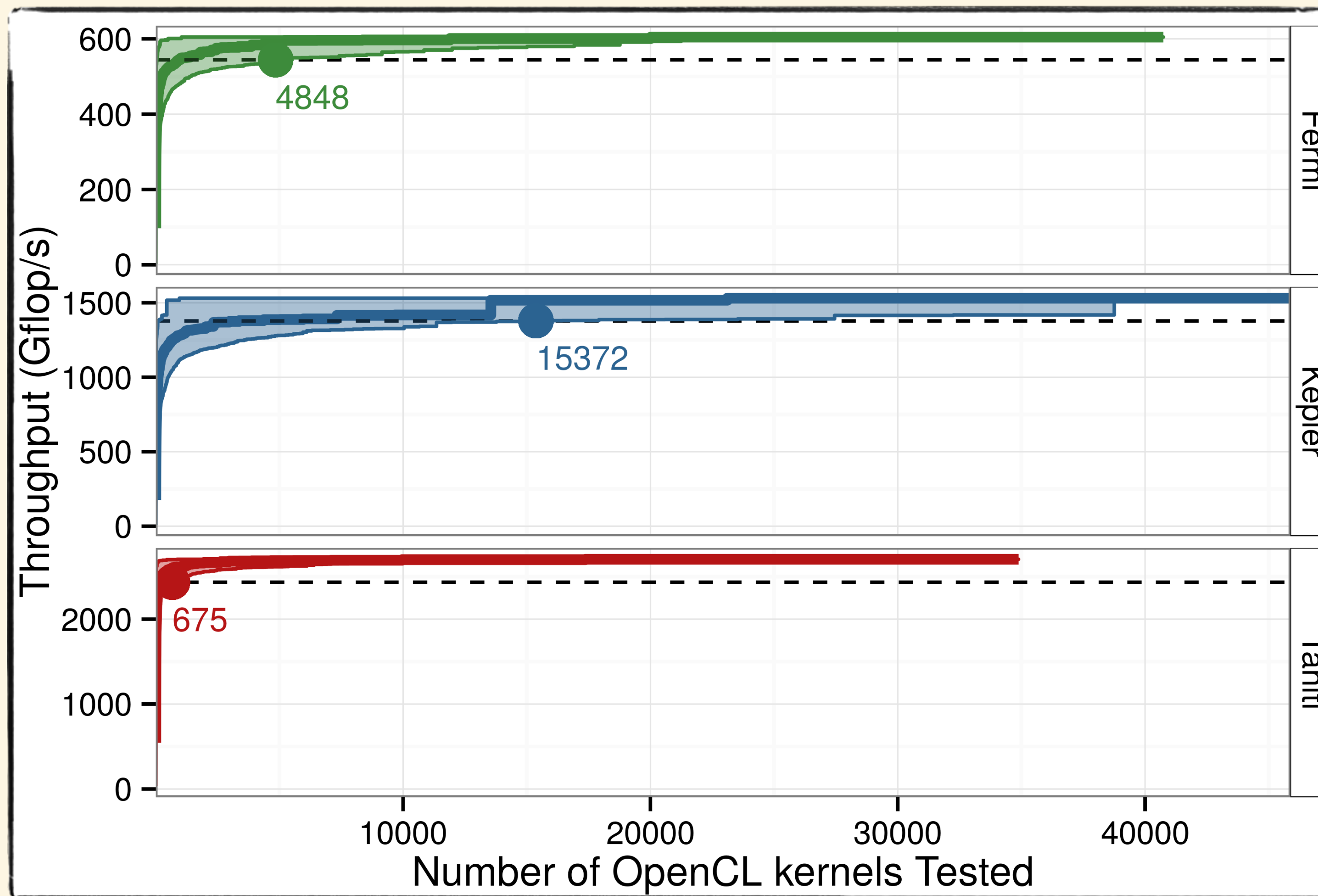
# EXPLORATION SPACE MATRIX MULTIPLICATION



Only few generated code with very good performance
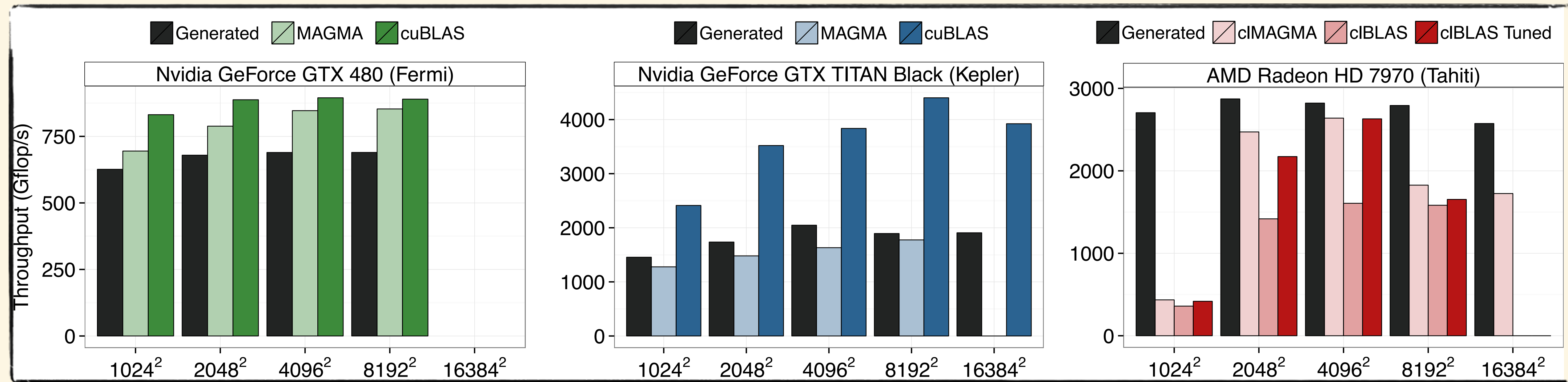
# EVEN RANDOMISED SEARCH WORKS WELL!



Still: One can expect to find a good performing kernel quickly!
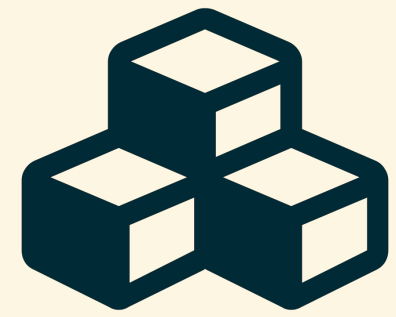
# PERFORMANCE RESULTS MATRIX MULTIPLICATION

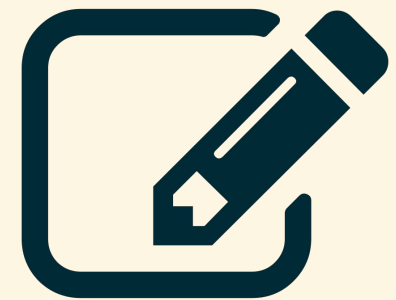

Performance close or better than hand-tuned MAGMA library

# 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 ; }
```



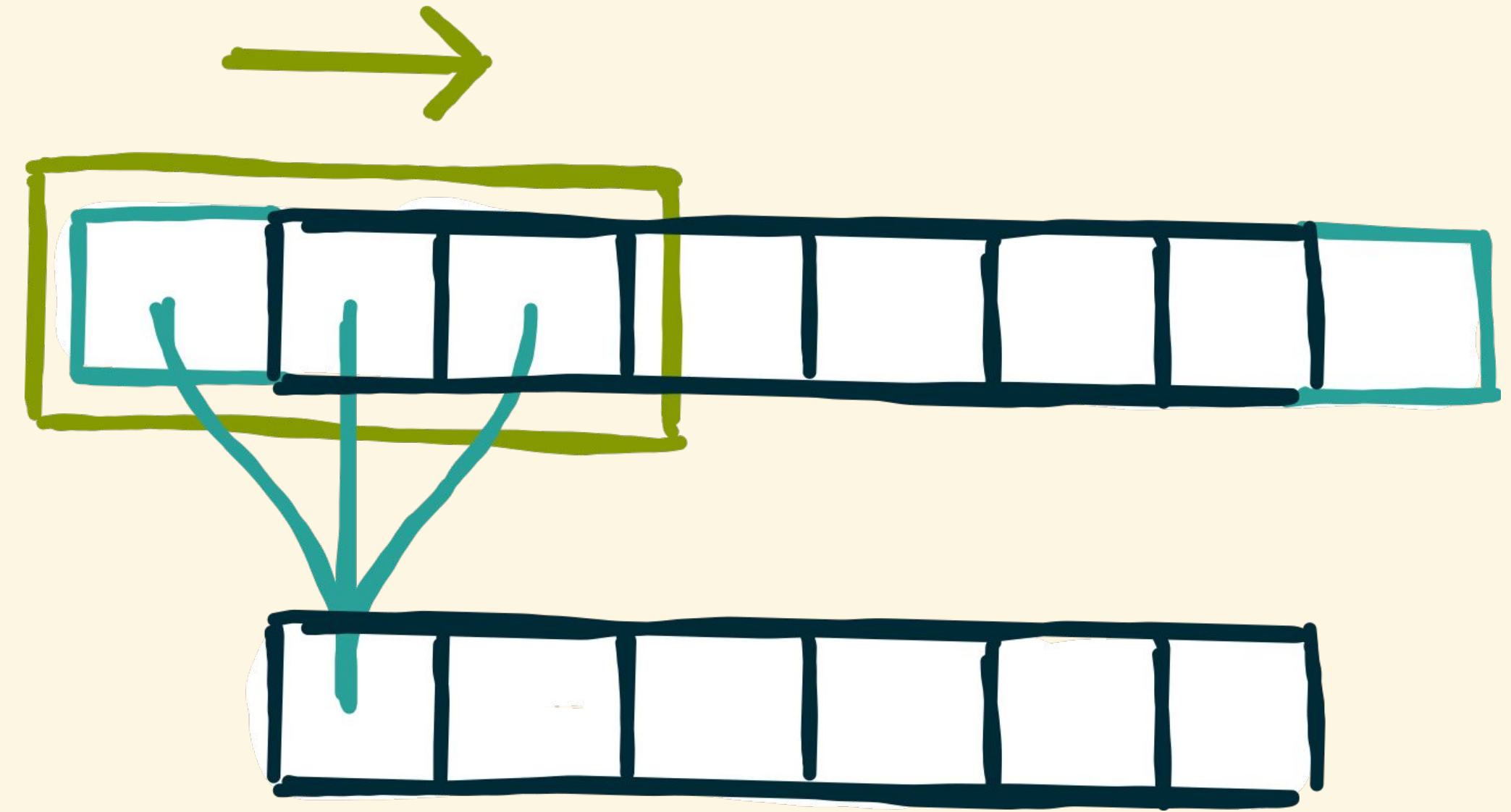(a)   access **neighborhoods** for every element

# 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 ; }
```
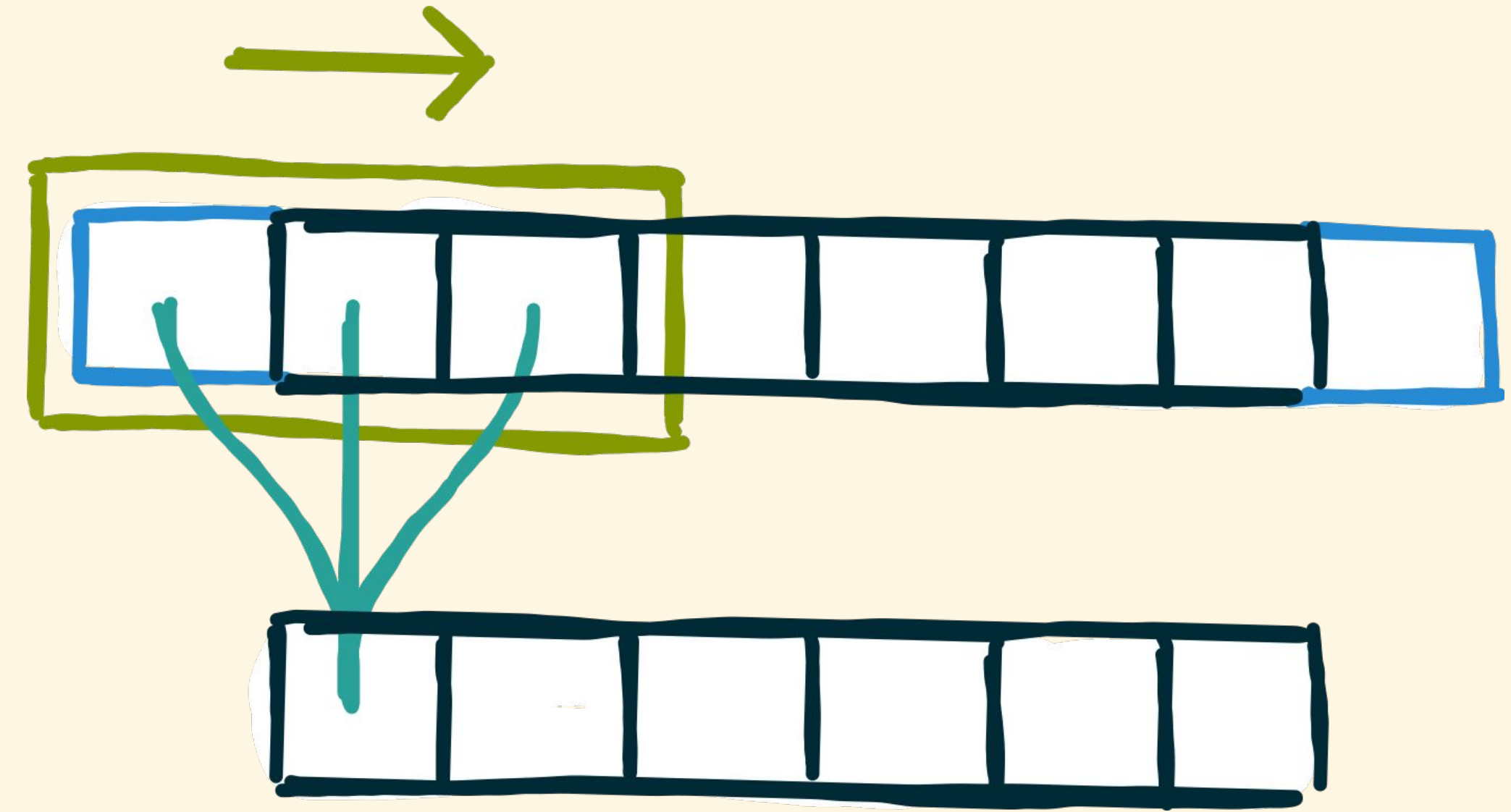


**(a)** access **neighborhoods** for every element

**(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 ; }
```
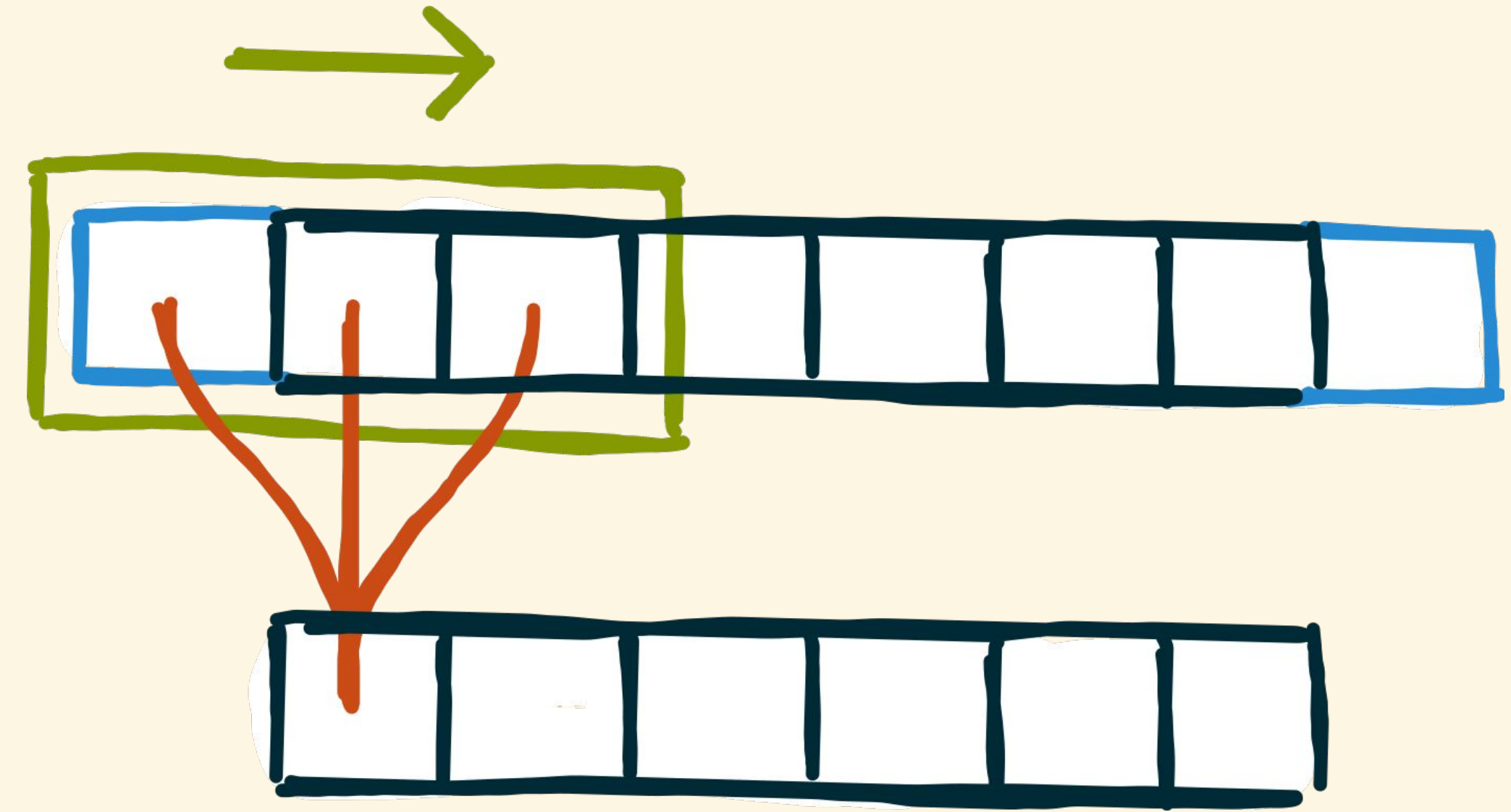


(a)  access **neighborhoods** for every element

(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 ; }
```
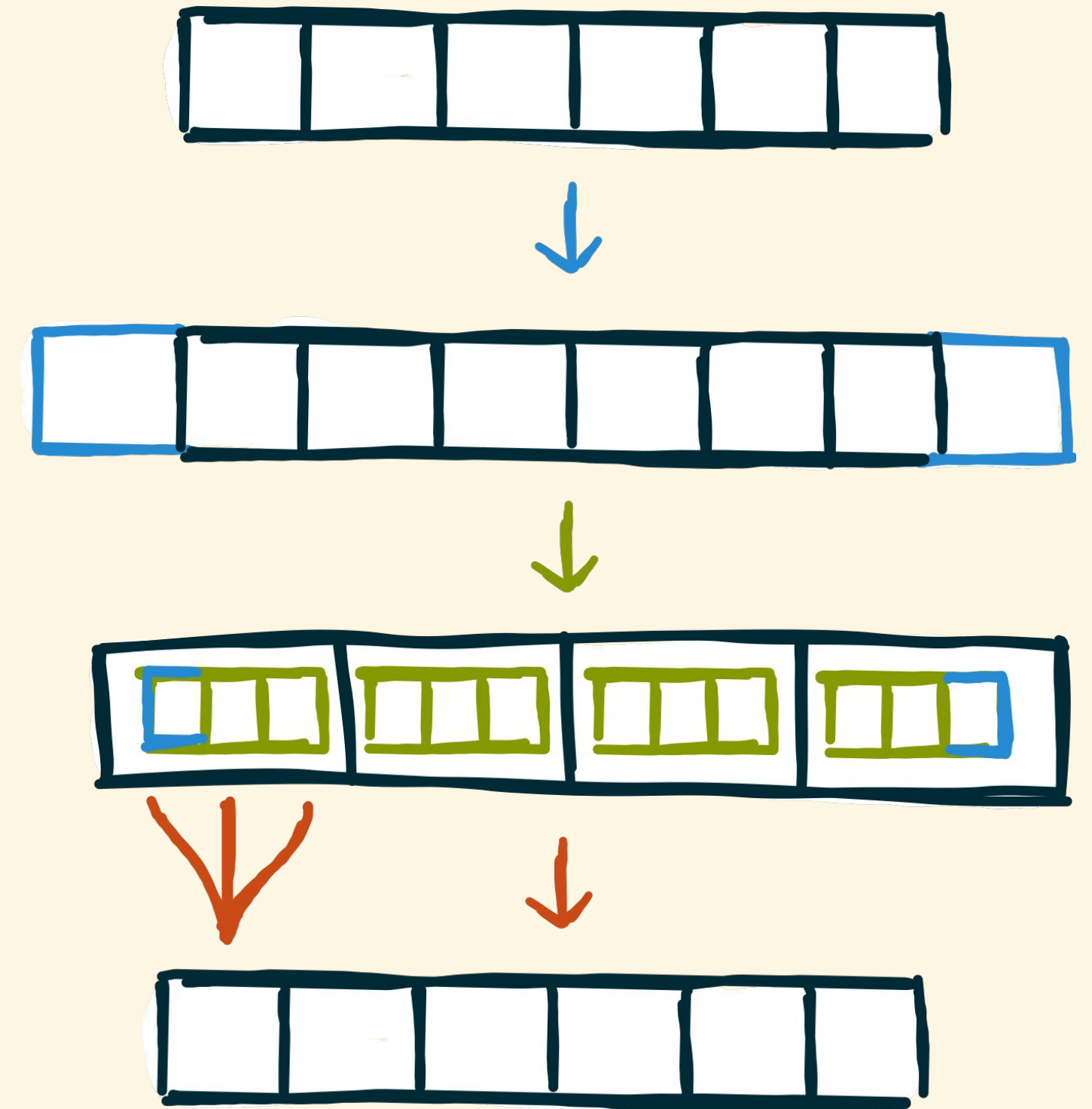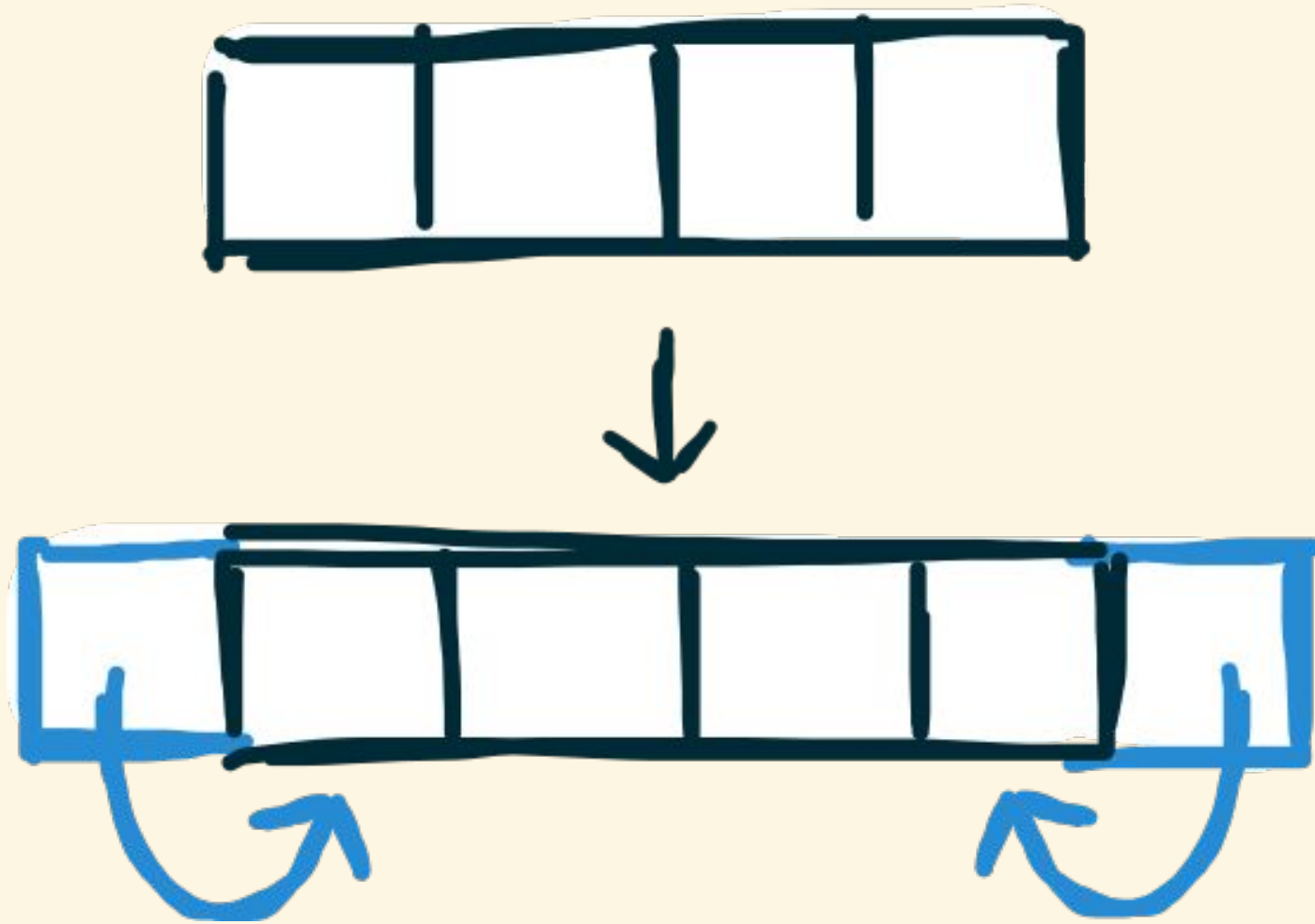
(a)  access **neighborhoods** for every element

(b)  specify **boundary handling**

(c)  apply **stencil function** to neighborhoods
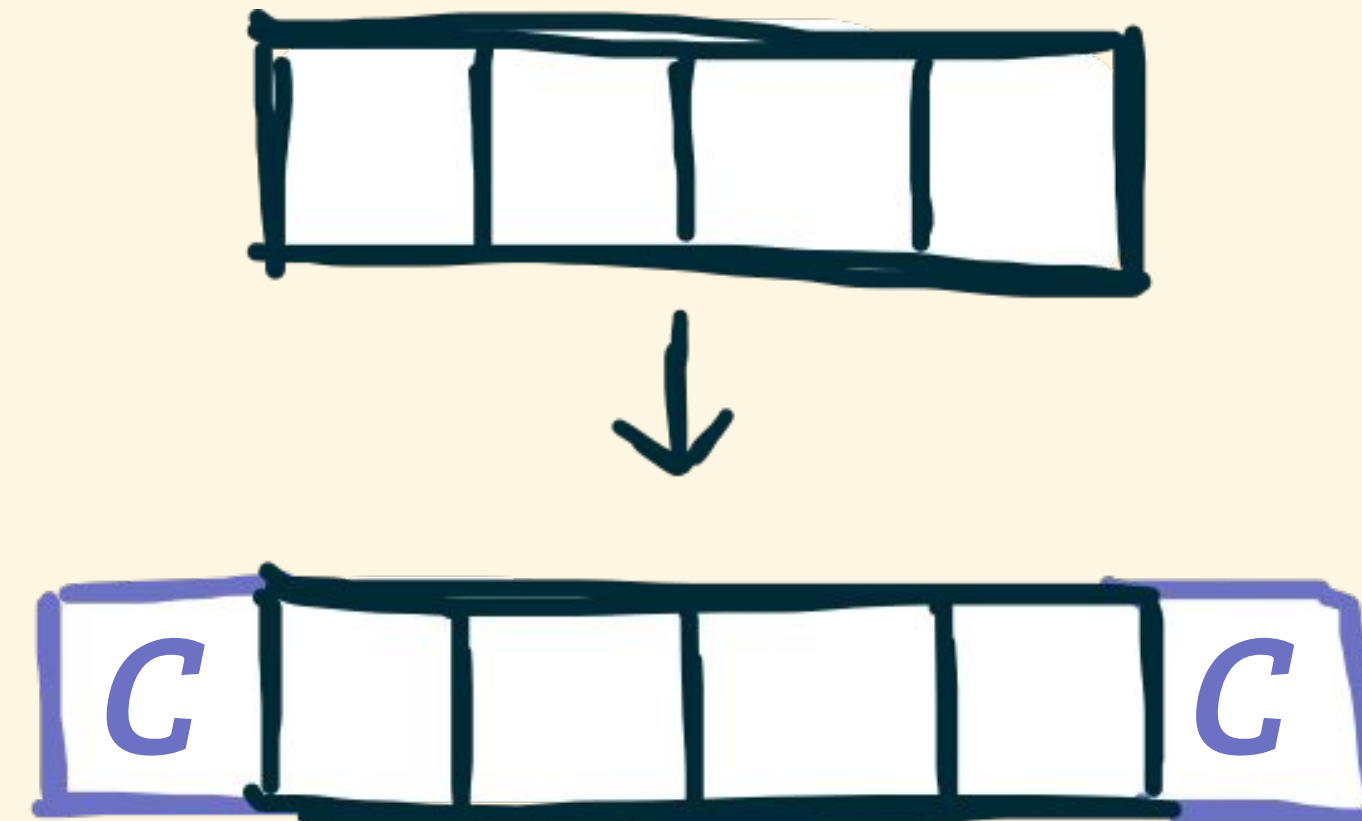
# BOUNDARY HANDLING USING *PAD*

## *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]
```
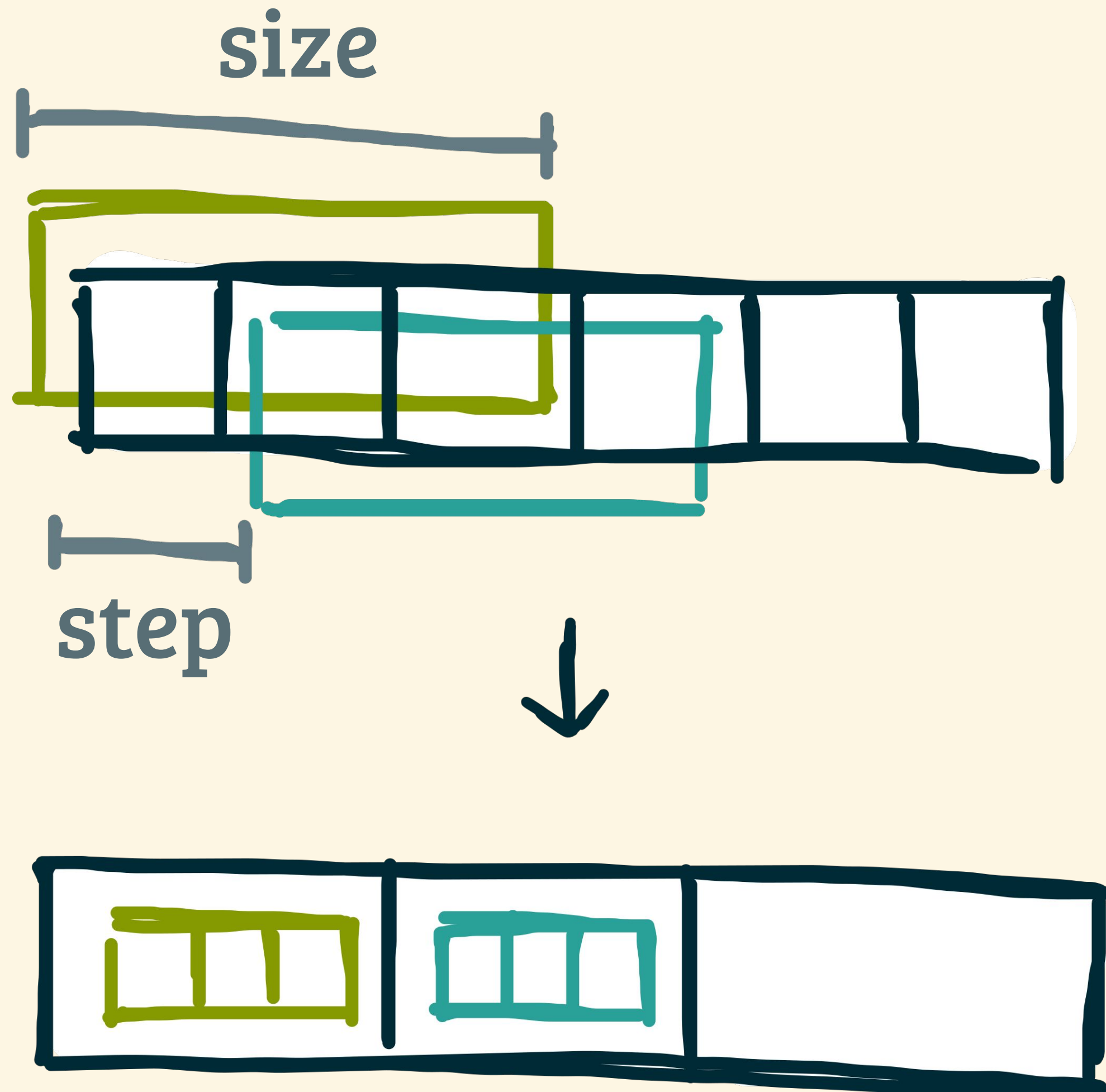
## *pad ( constant )*



### pad-constant.lift

```
constant(i, n) = C


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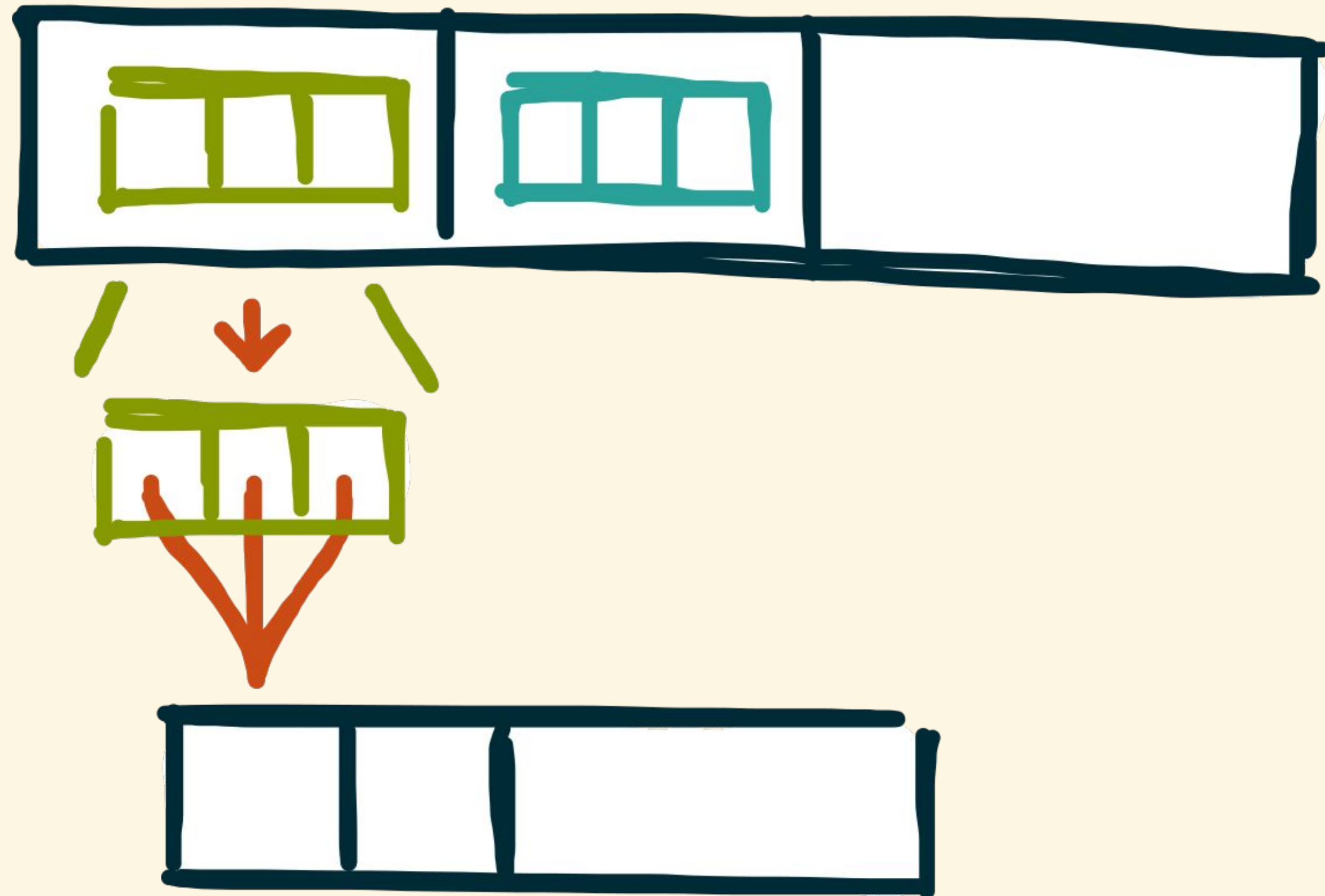
# NEIGHBORHOOD CREATION USING SLIDE

size

step

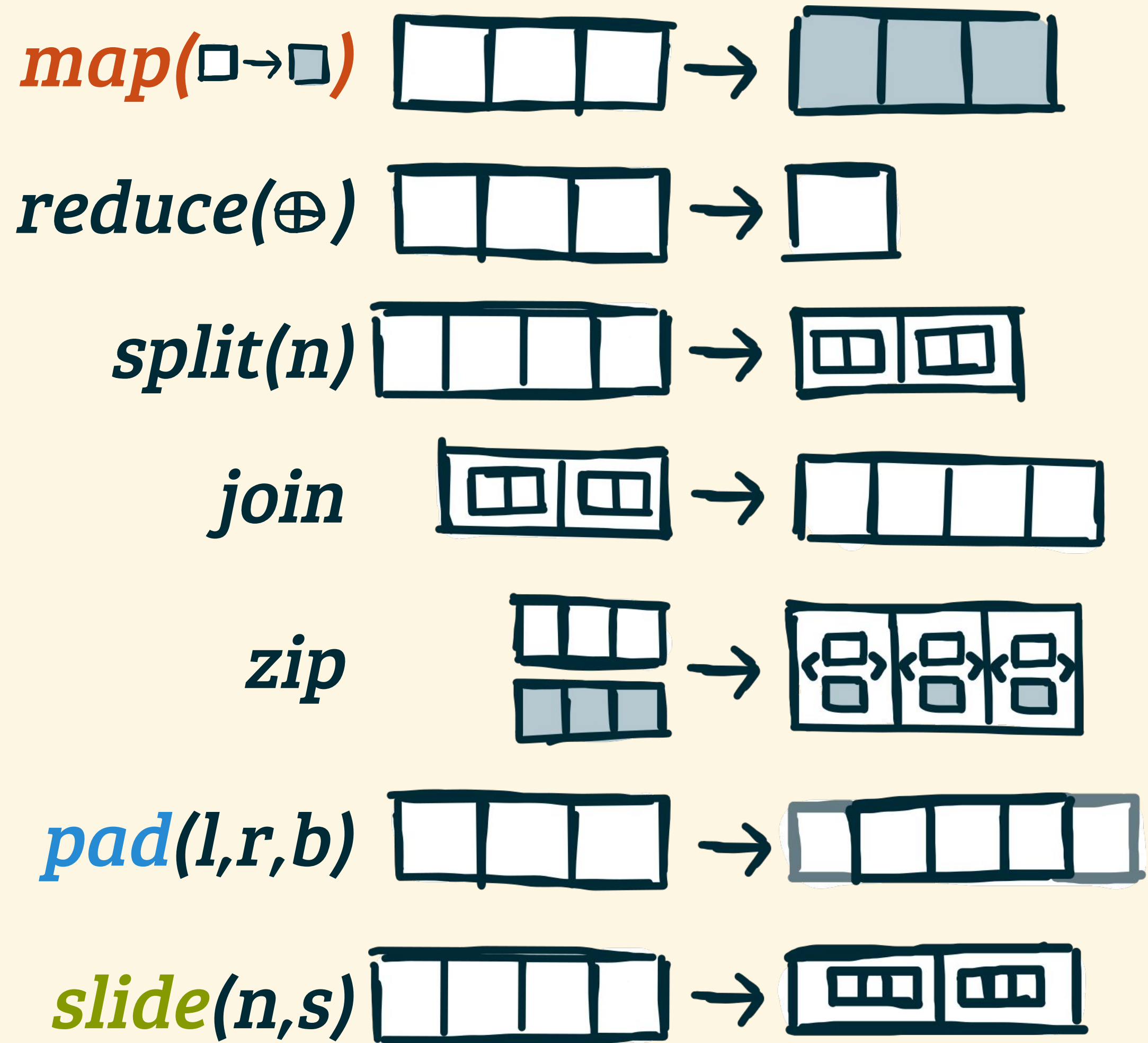slide-example.lift

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

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

# APPLYING STENCIL FUNCTION USING MAP



sum-neighborhoods.lift

```
map(nbh =>
    reduce(add, 0.0f, nbh))
```
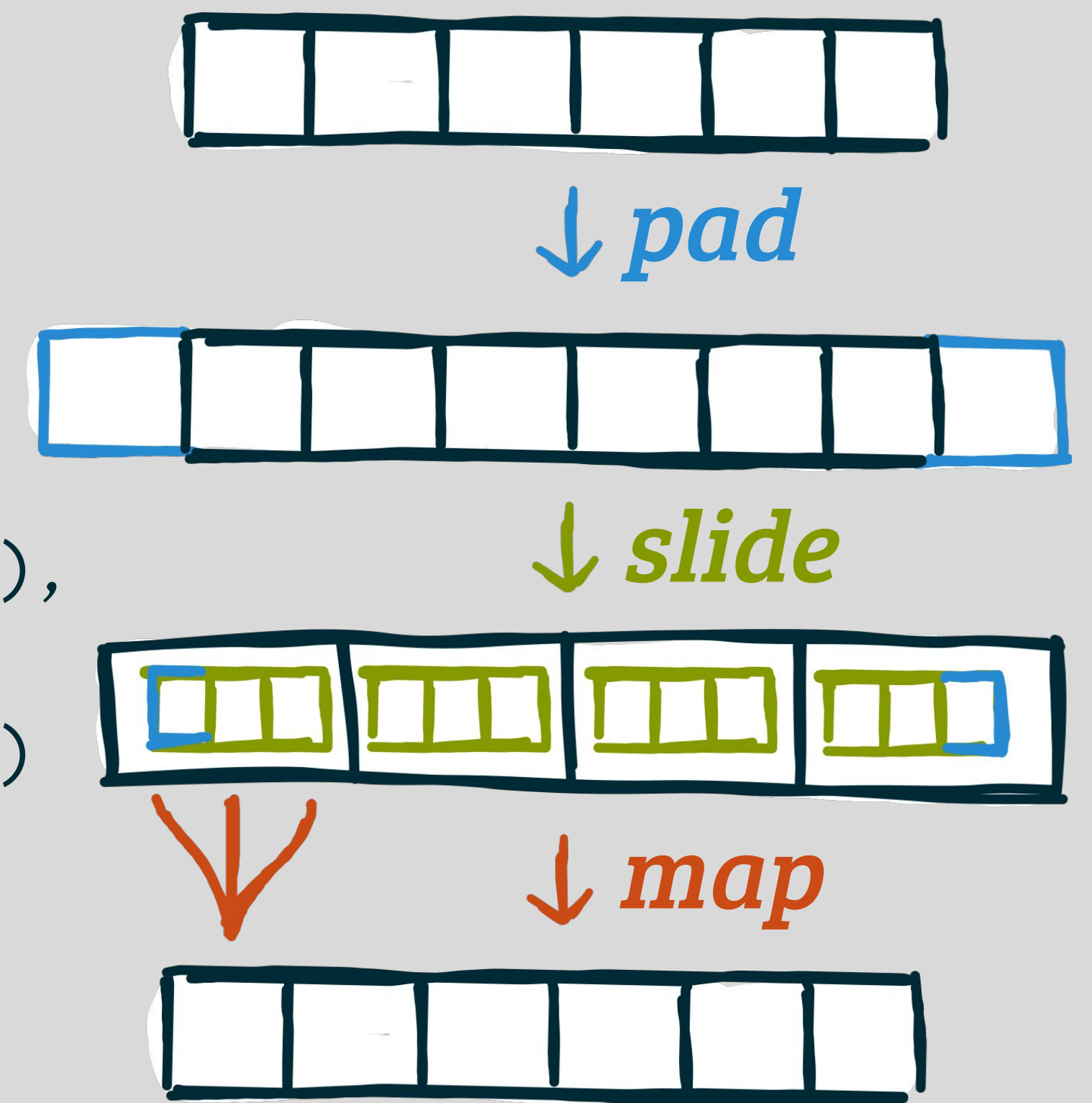
# PUTTING IT TOGETHER

**map(□→▨)**

**reduce(⊕)**

**split(n)**

**join**

**zip**

**pad(l,r,b)**

**slide(n,s)**

## stencil1D.lift

↓ *pad*

↓ *slide*

↓ *map*

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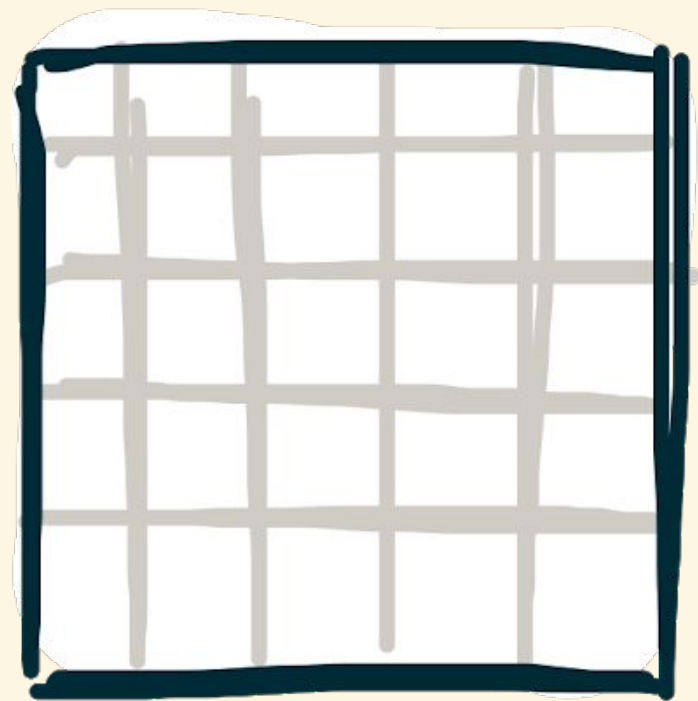
# MULTIDIMENSIONAL STENCIL COMPUTATIONS

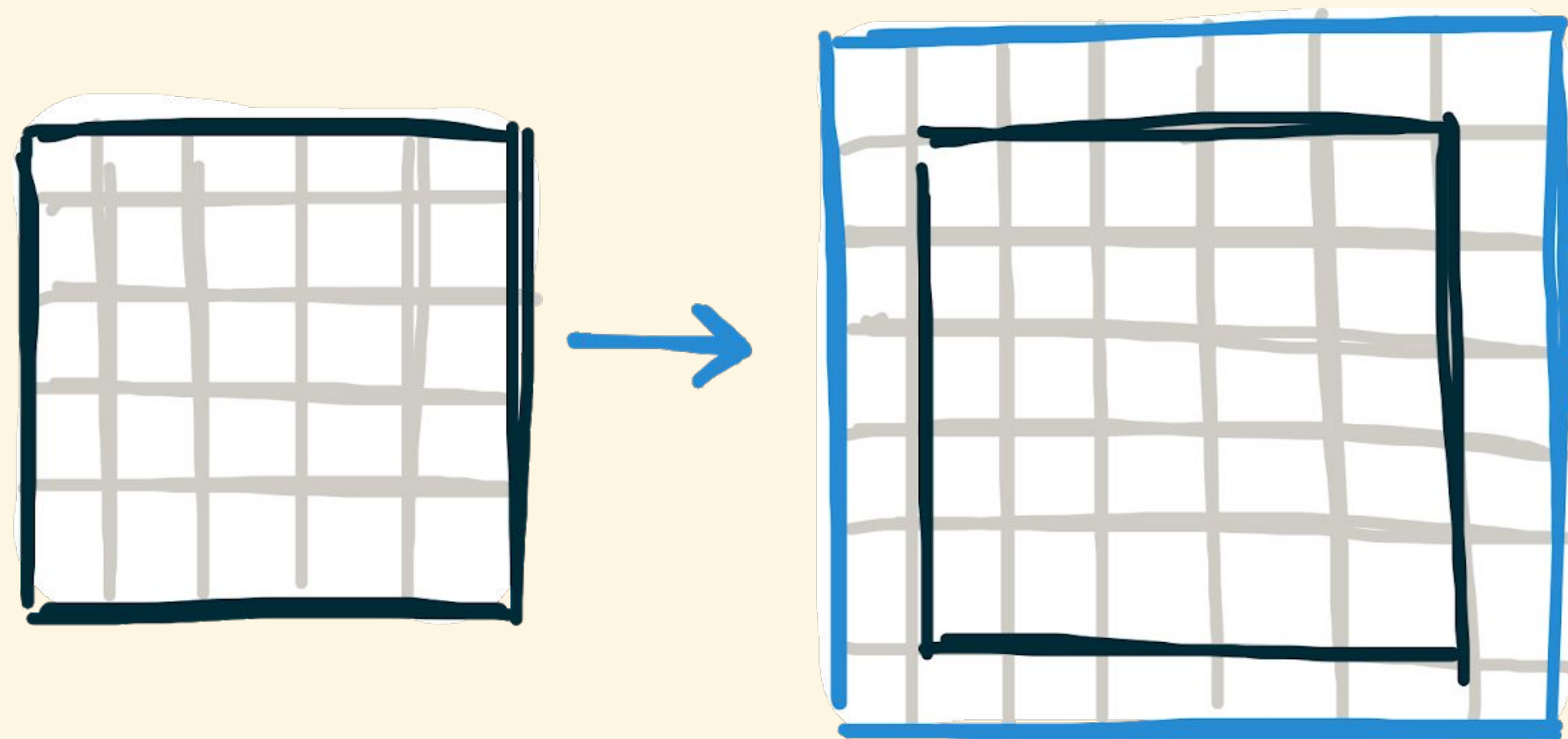are expressed as compositions of intuitive, generic 1D primitives

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

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

$$pad_2(1,1,clamp,input)$$
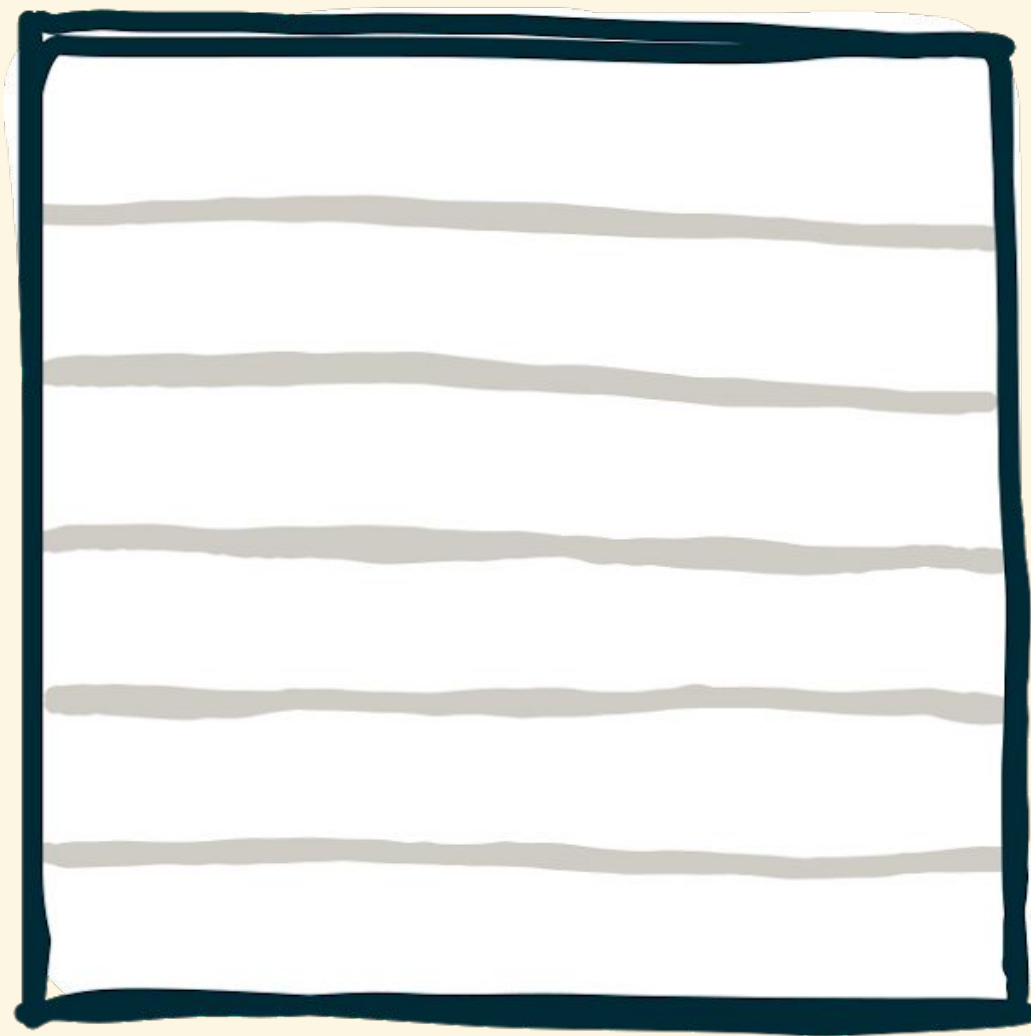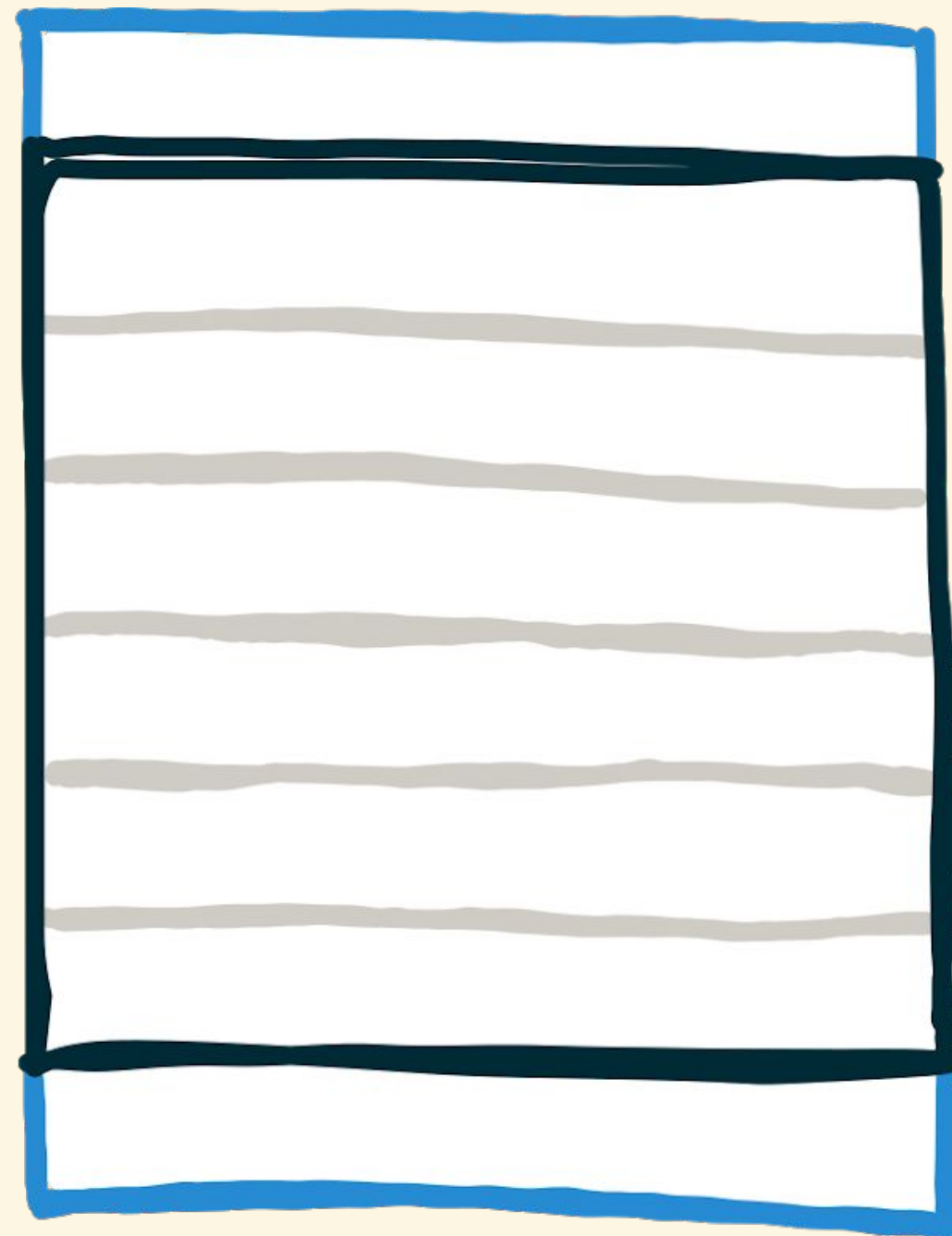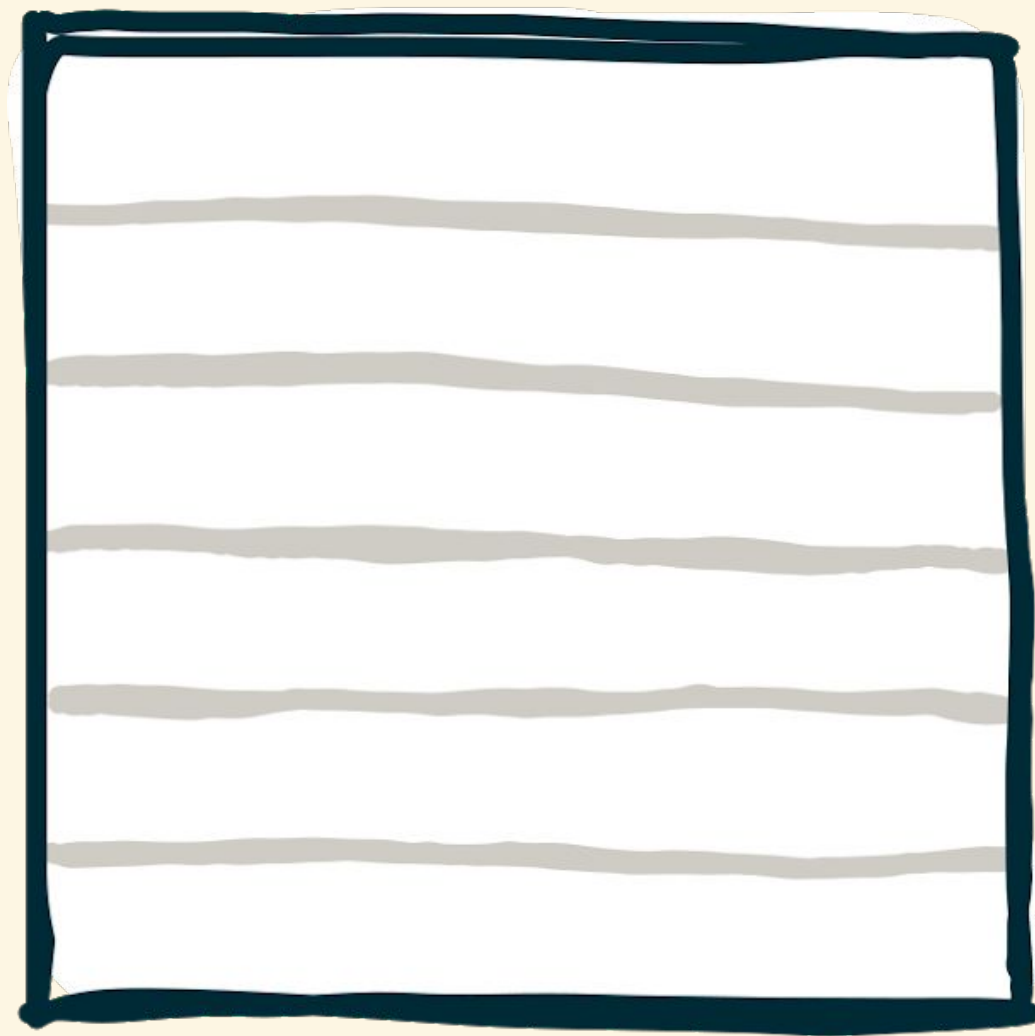
input
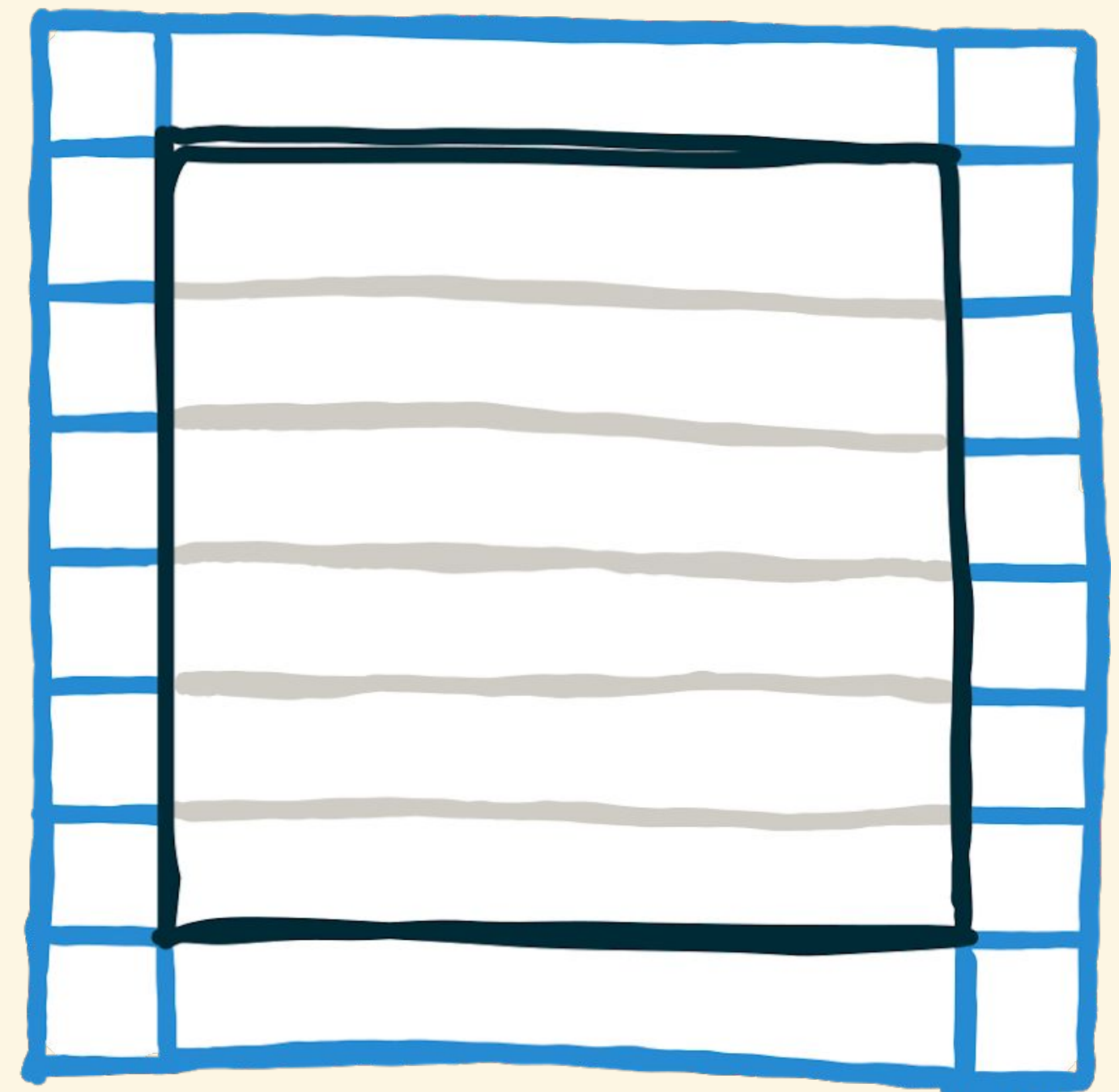
$pad_2 =$
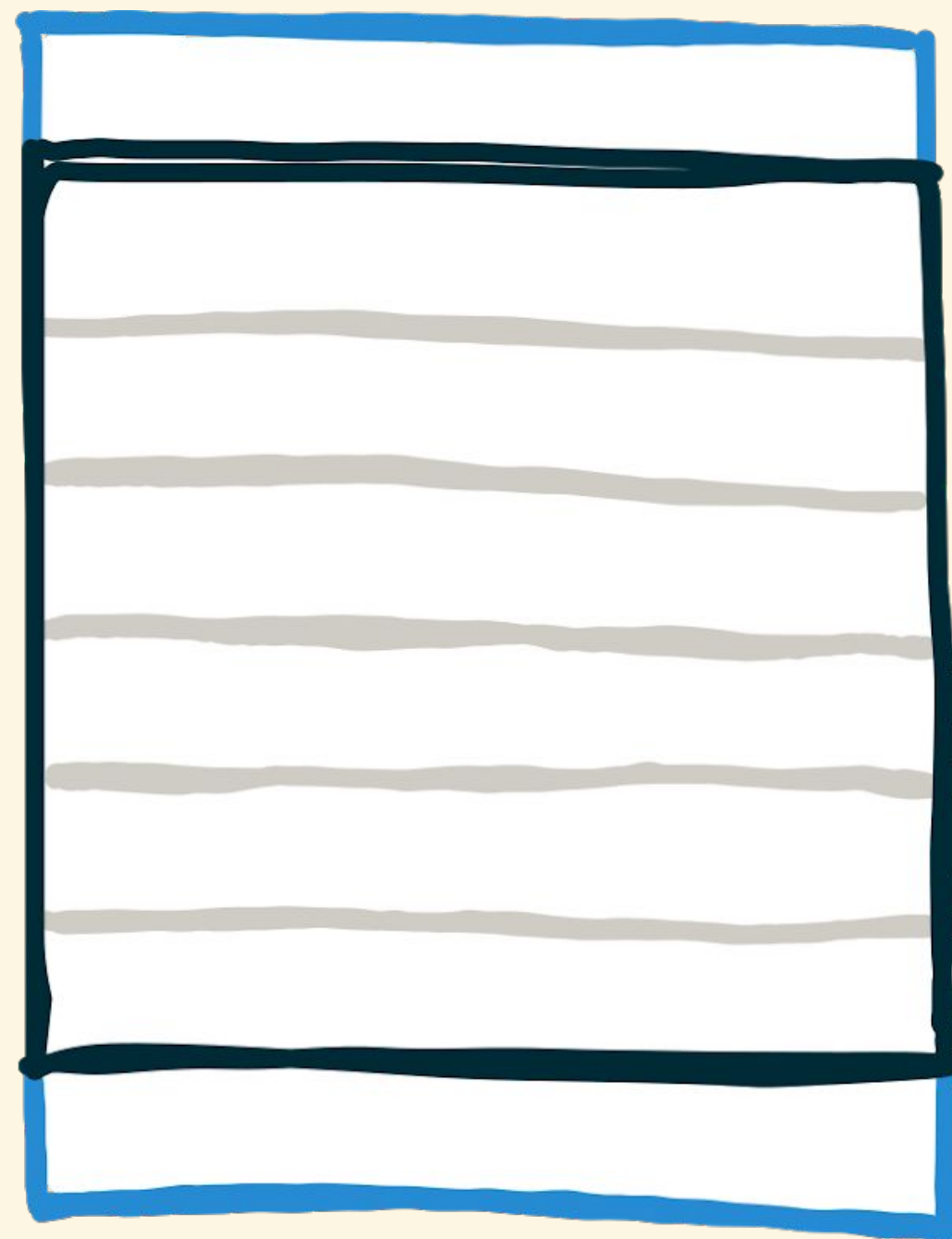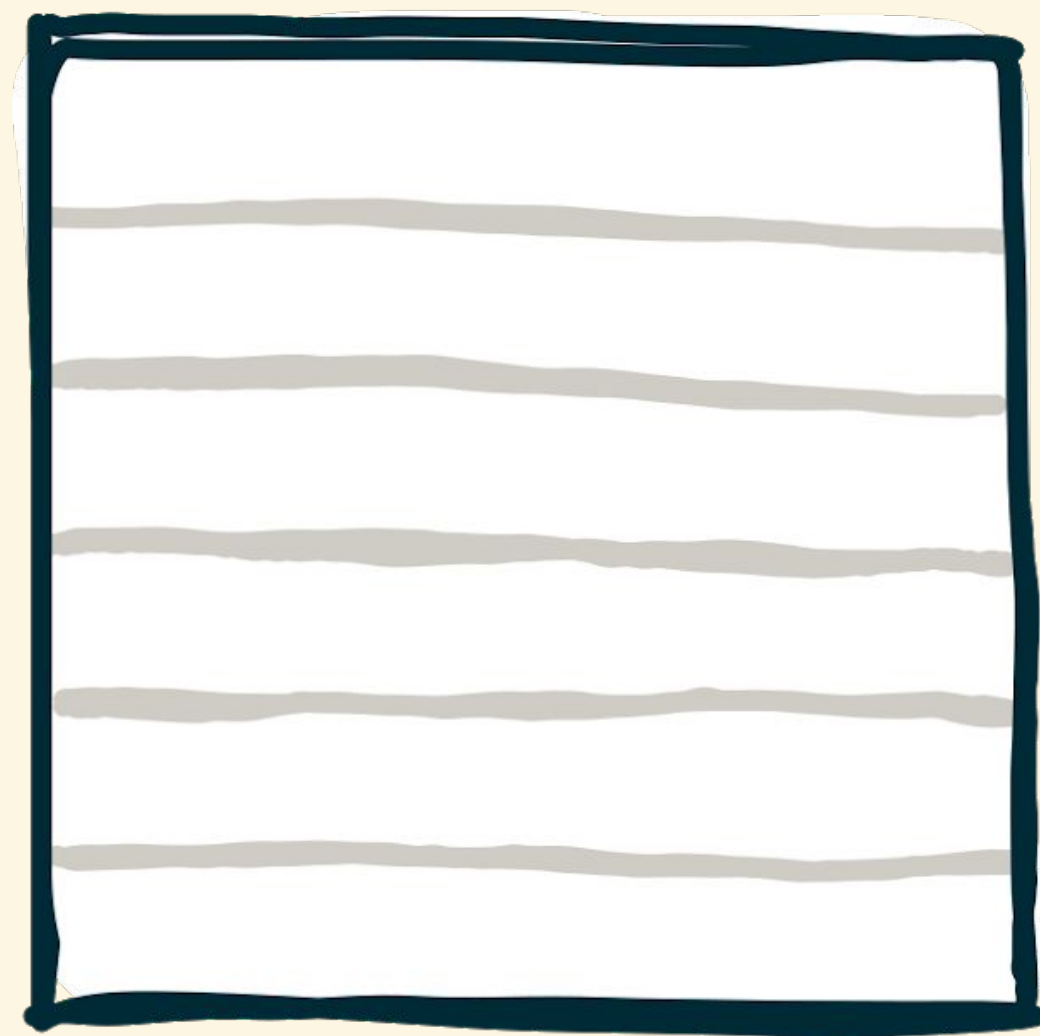
# MULTIDIMENSIONAL BOUNDARY HANDLING USING $PAD_2$

input



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

# MULTIDIMENSIONAL BOUNDARY HANDLING USING $PAD_2$

input



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

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives



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

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

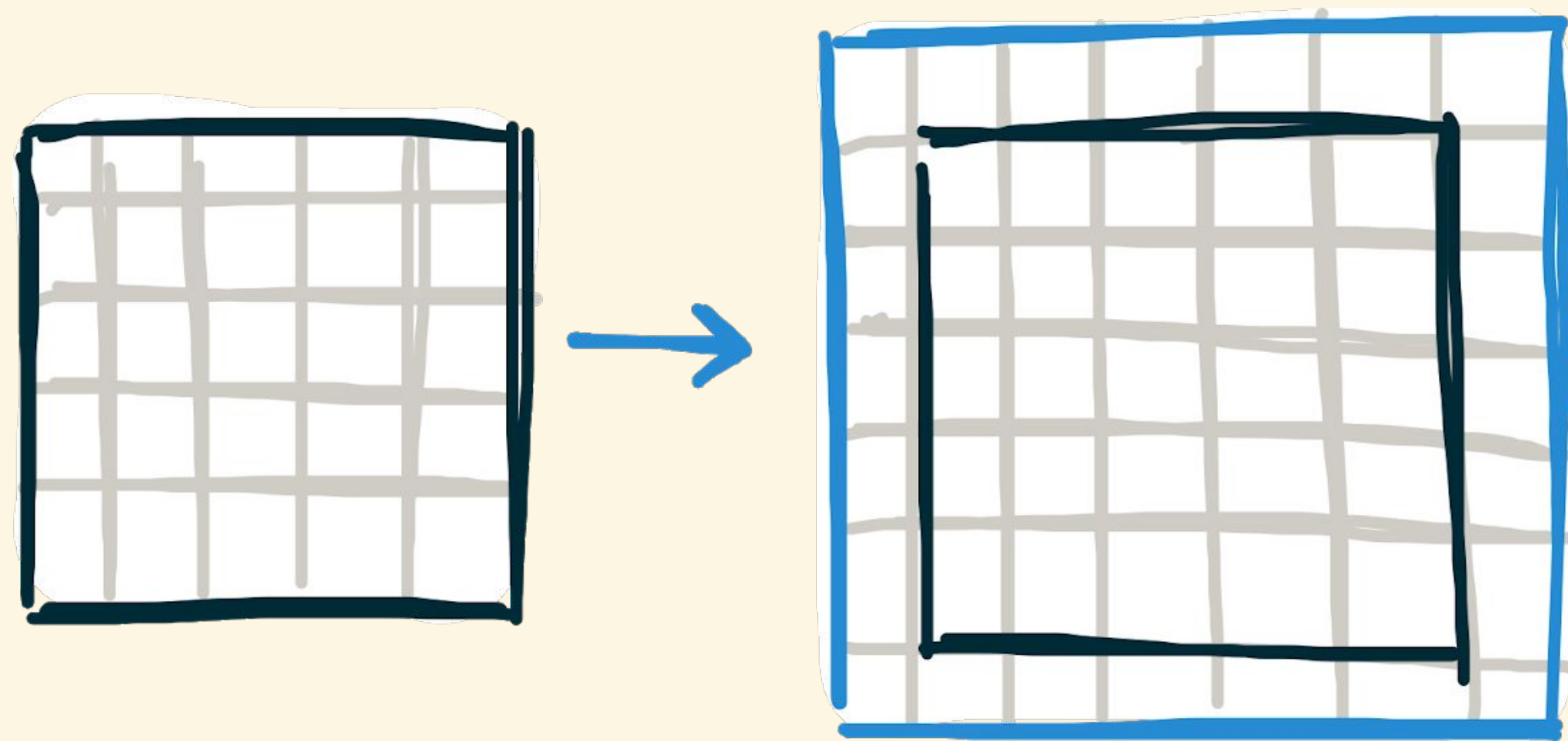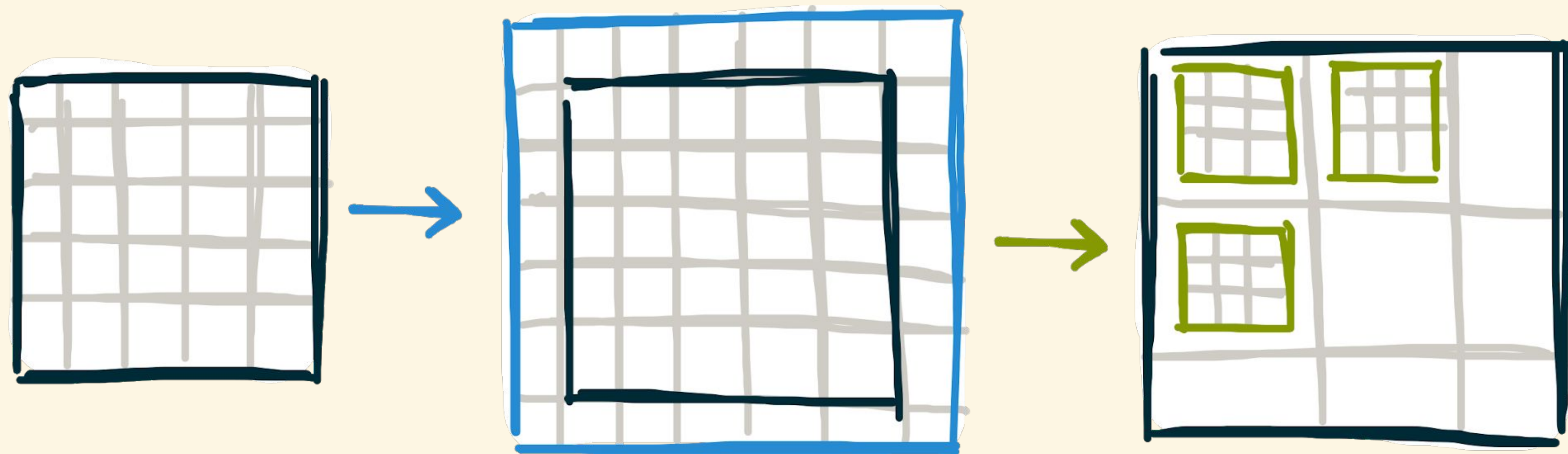$$slide_2(3,1,\ pad_2(1,1,clamp,input))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

are expressed as compositions of intuitive, generic 1D primitives

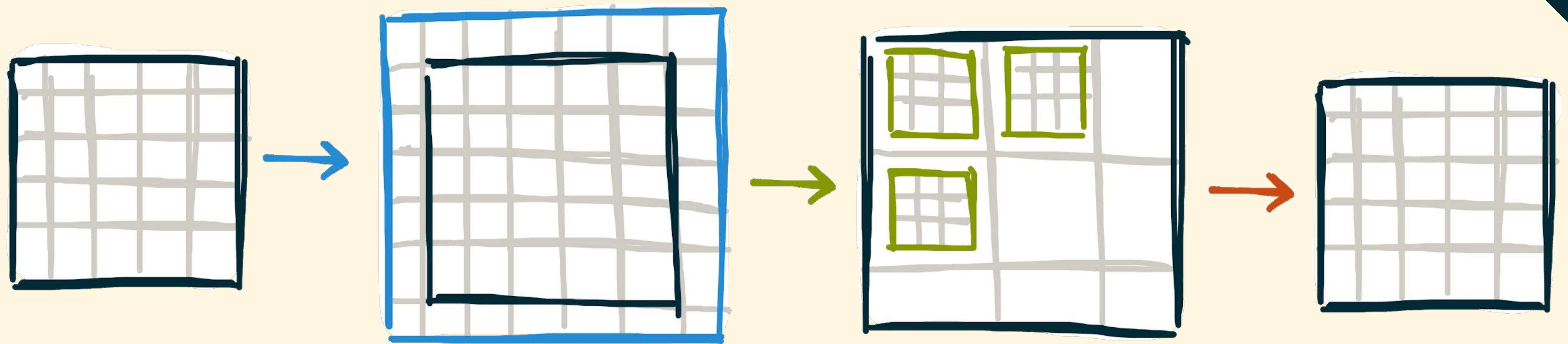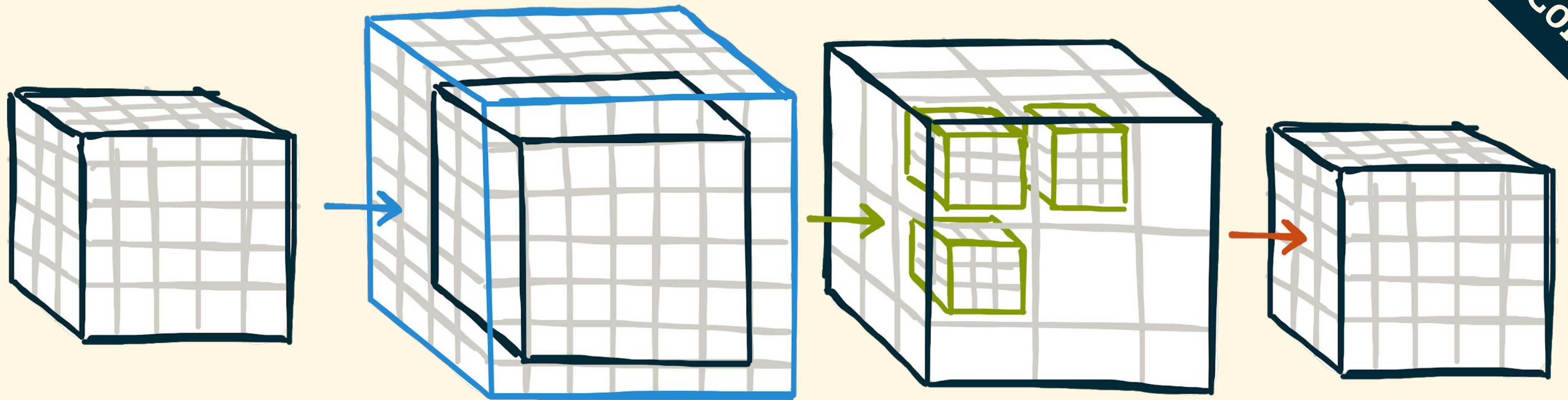$$map_2(sum,\ slide_2(3,1,\ pad_2(1,1,clamp,input)))$$

# MULTIDIMENSIONAL STENCIL COMPUTATIONS

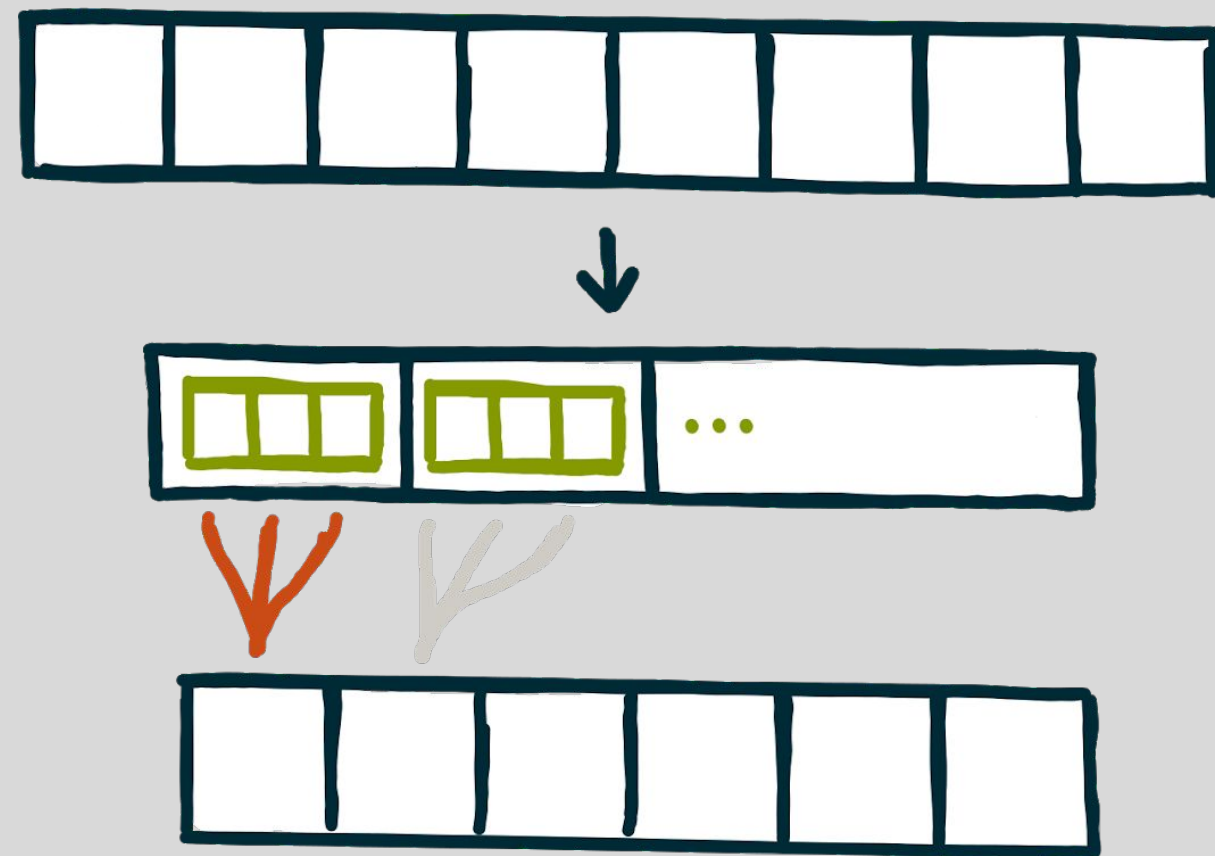are expressed as compositions of intuitive, generic 1D primitives

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

# OVERLAPPED TILING AS A REWRITE RULE
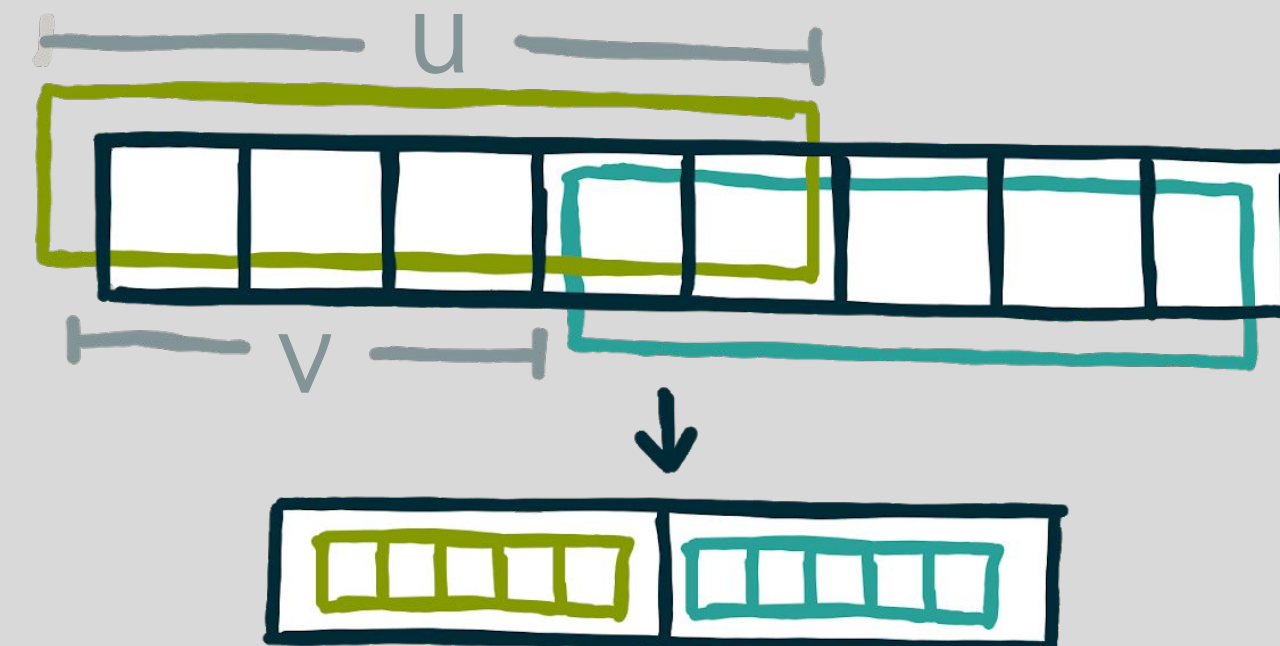
## overlapped tiling rule
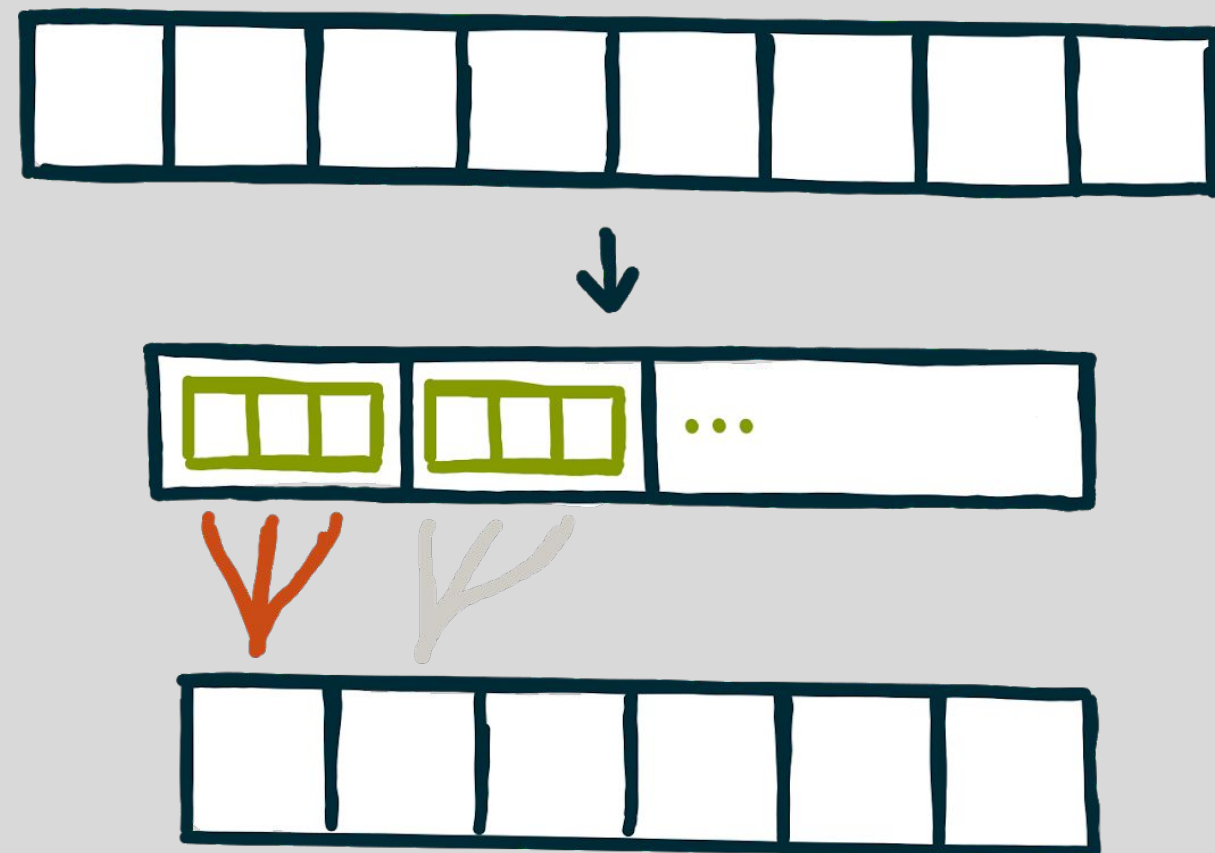
`map`(f, `slide`(3,1,input))

# OVERLAPPED TILING AS A REWRITE RULE

## overlapped tiling rule

$map(f, slide(3,1,input)) \mapsto$
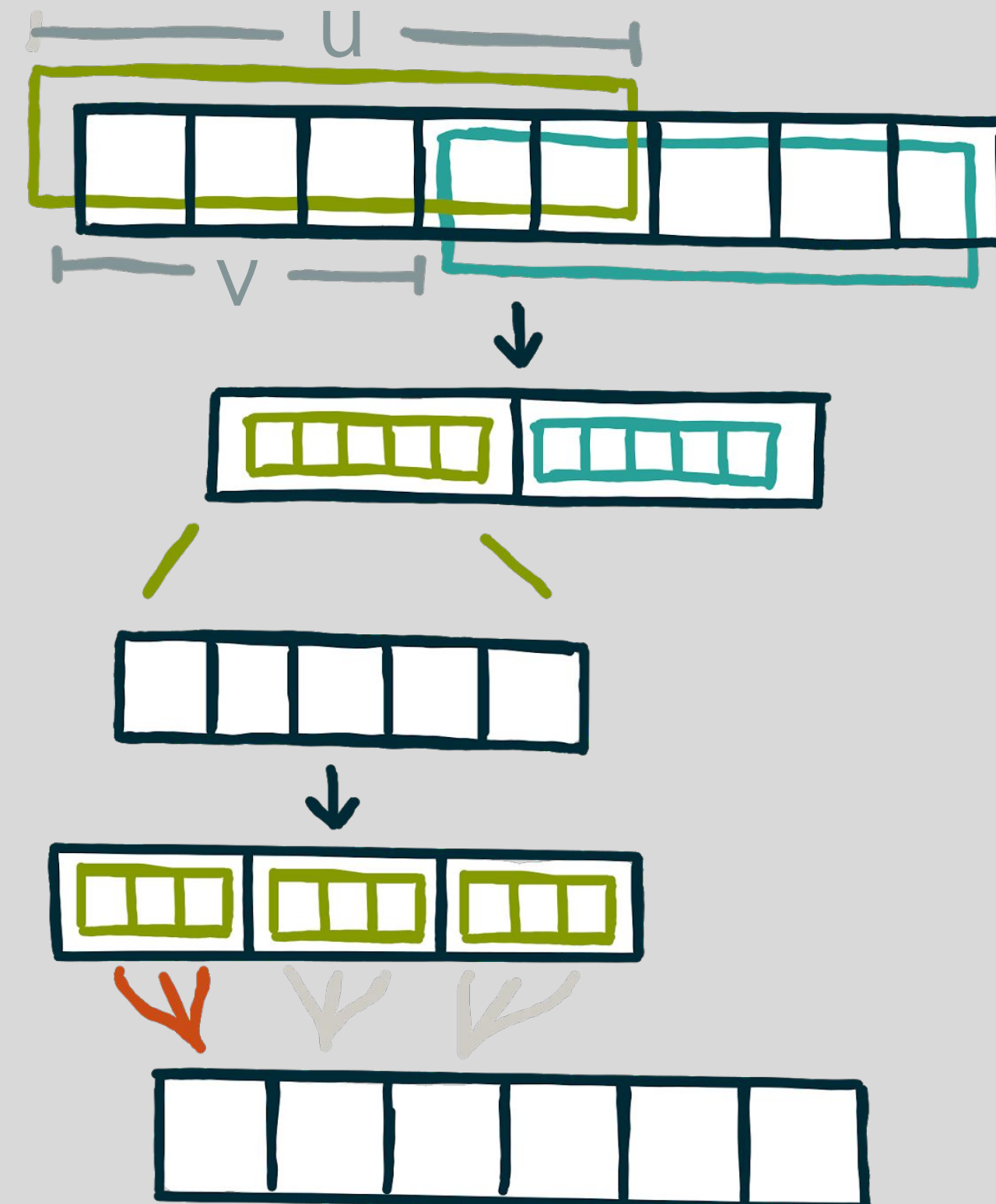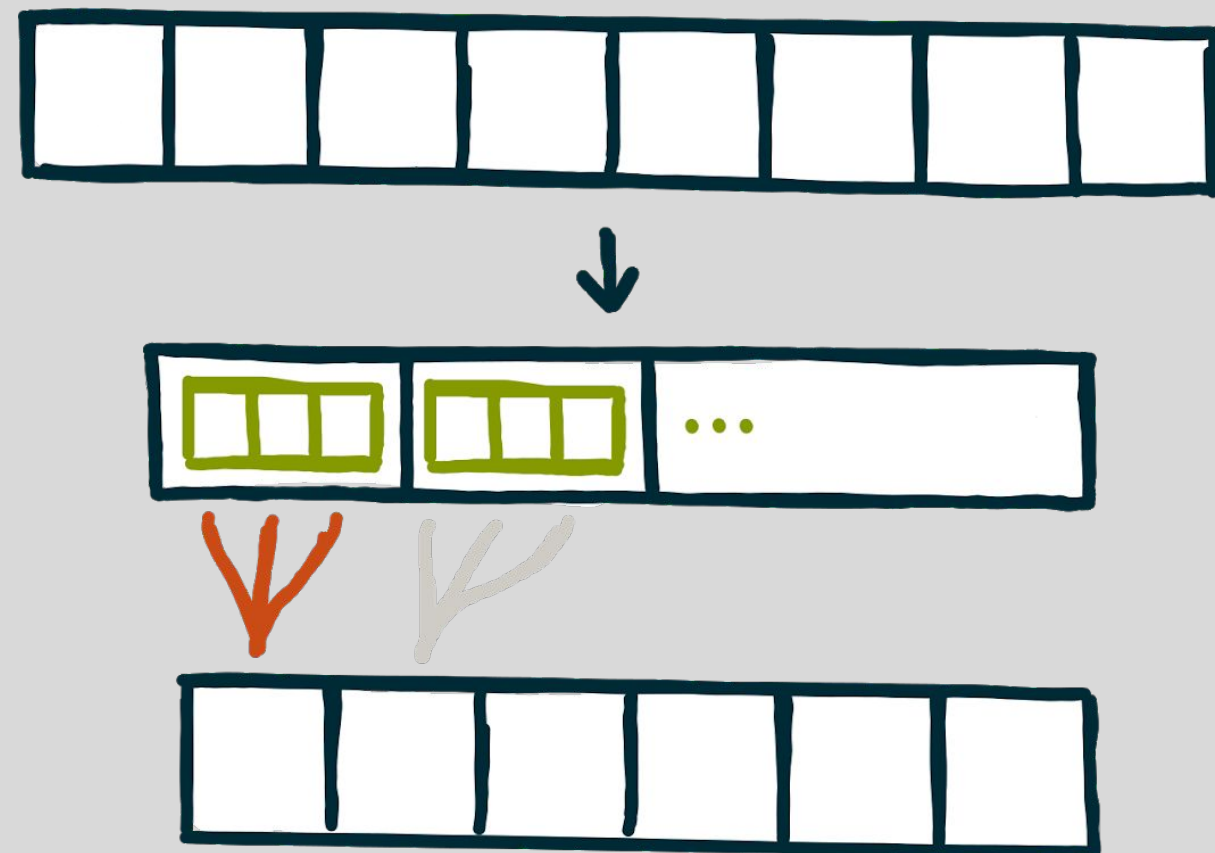
$slide(u,v,input)$

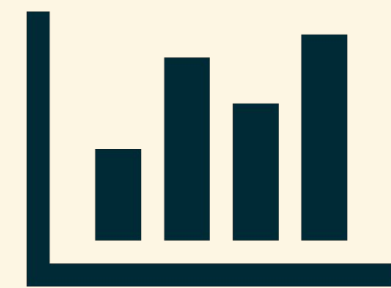# OVERLAPPED TILING AS A REWRITE RULE



overlapped tiling rule

map(f, slide(3,1,input)) ↦ join(map(tile ⇒
                                     map(f, slide(3,1,tile)),
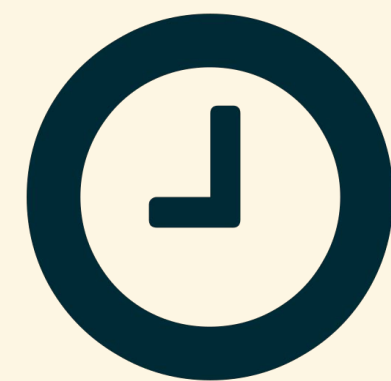                                     slide(u,v,input)))

# EXPERIMENTAL EVALUATION

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

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

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

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

DSL    DSL    DSL

**High-Level IR**

**Explore Optimizations by rewriting**

**Low-Level Program**

Multicore CPU

GPU    HPC    Mobile

Xeon Phi    KNC    KNL

...    Hardware

# COMPARISON WITH HAND-OPTIMIZED CODES

higher is better



Lift achieves the same performance
as hand optimized code

# COMPARISON WITH POLYHEDRAL COMPILATION

## higher is better



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

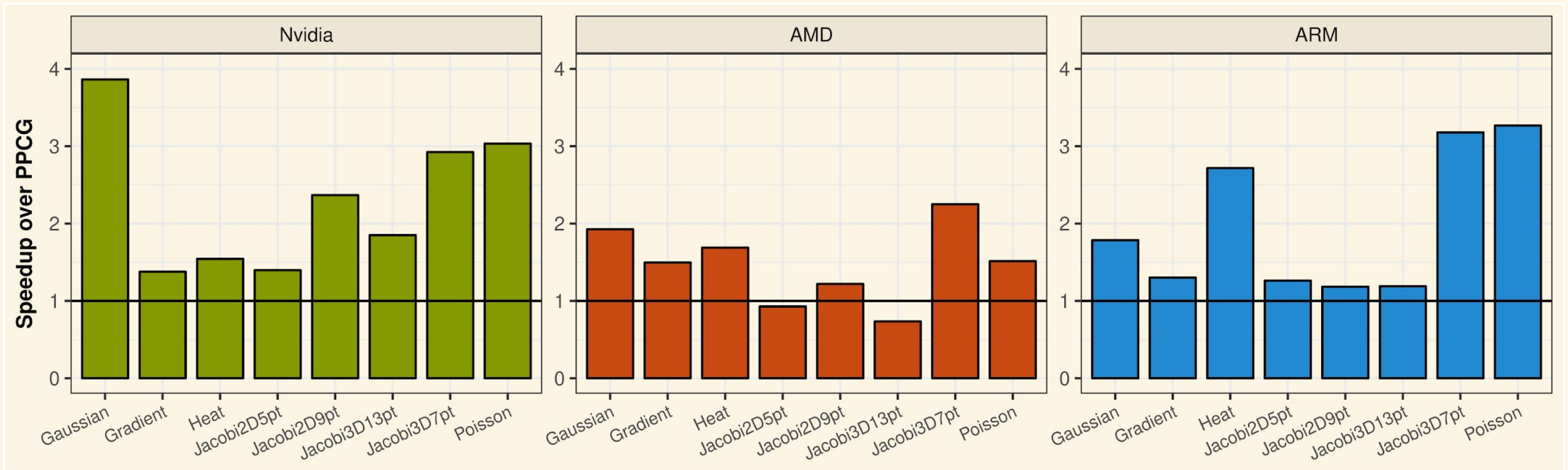# LIFT IS OPEN SOURCE!

more info at:
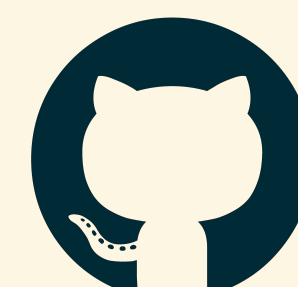
# lift-project.org

" Paper    Artifacts    Source Code



Naums Mogers    Lu Li    Christophe Dubach    Bastian Hagedorn    Toomas Remmelg    Larisa Stoltzfus    Michel Steuwer    Federico Pizzuti    Adam Harries

# Automatic Matching of Legacy Code to Heterogeneous APIs: An Idiomatic Approach

Philip Ginsbach
The University of Edinburgh
philip.ginsbach@ed.ac.uk

Toomas Remmelg
The University of Edinburgh
toomas.remmelg@ed.ac.uk

Michel Steuwer
University of Glasgow
michel.steuwer@glasgow.ac.uk

Bruno Bodin
The University of Edinburgh
bbodin@ed.ac.uk

Christophe Dubach
The University of Edinburgh
christophe.dubach@ed.ac.uk

Michael F. P. O'Boyle
The University of Edinburgh
mob@ed.ac.uk

## Abstract

Heterogeneous accelerators often disappoint. They provide the prospect of great performance, but only deliver it when using vendor specific optimized libraries or domain specific languages. This requires considerable legacy code modifications, hindering the adoption of heterogeneous computing.
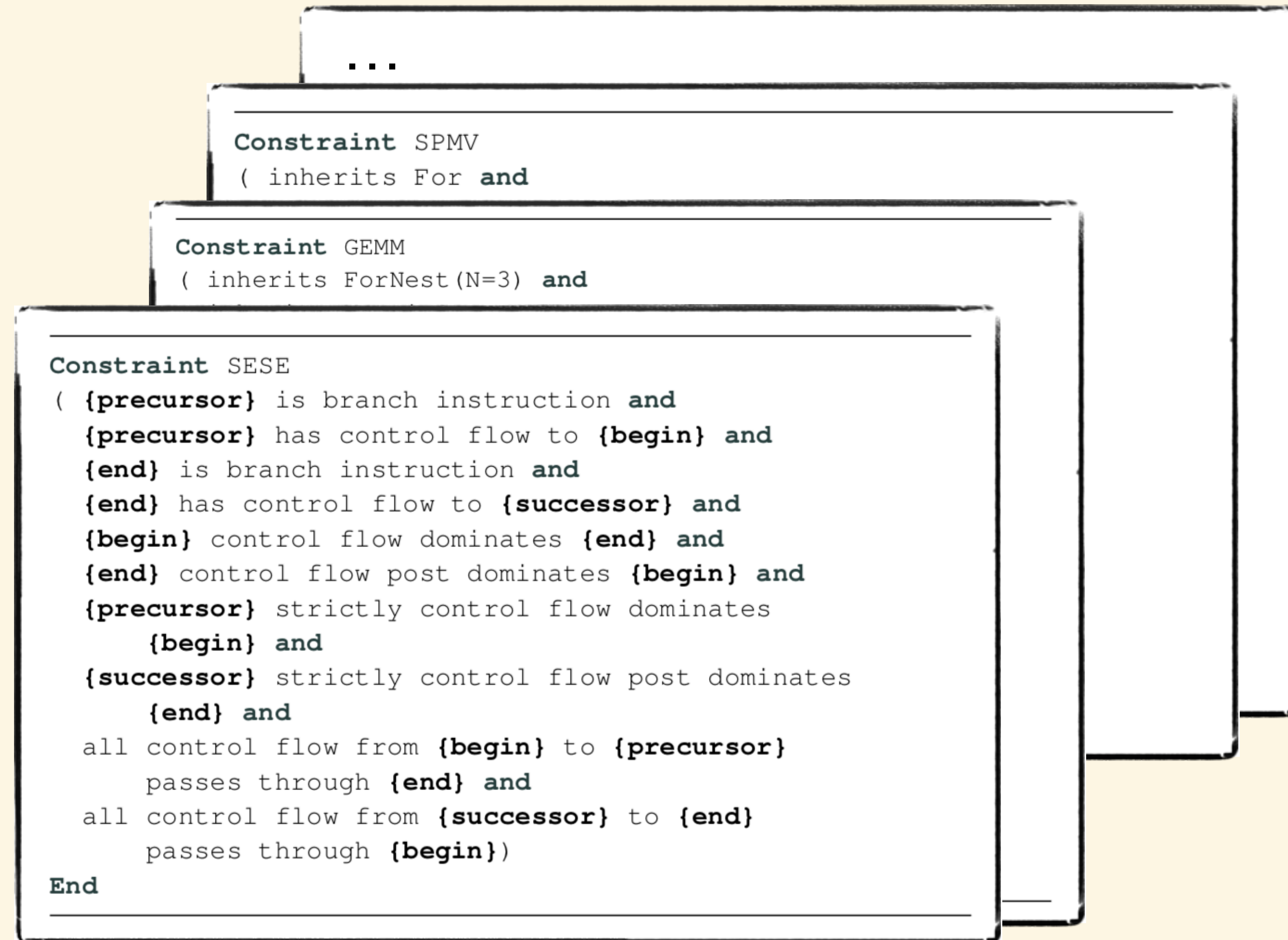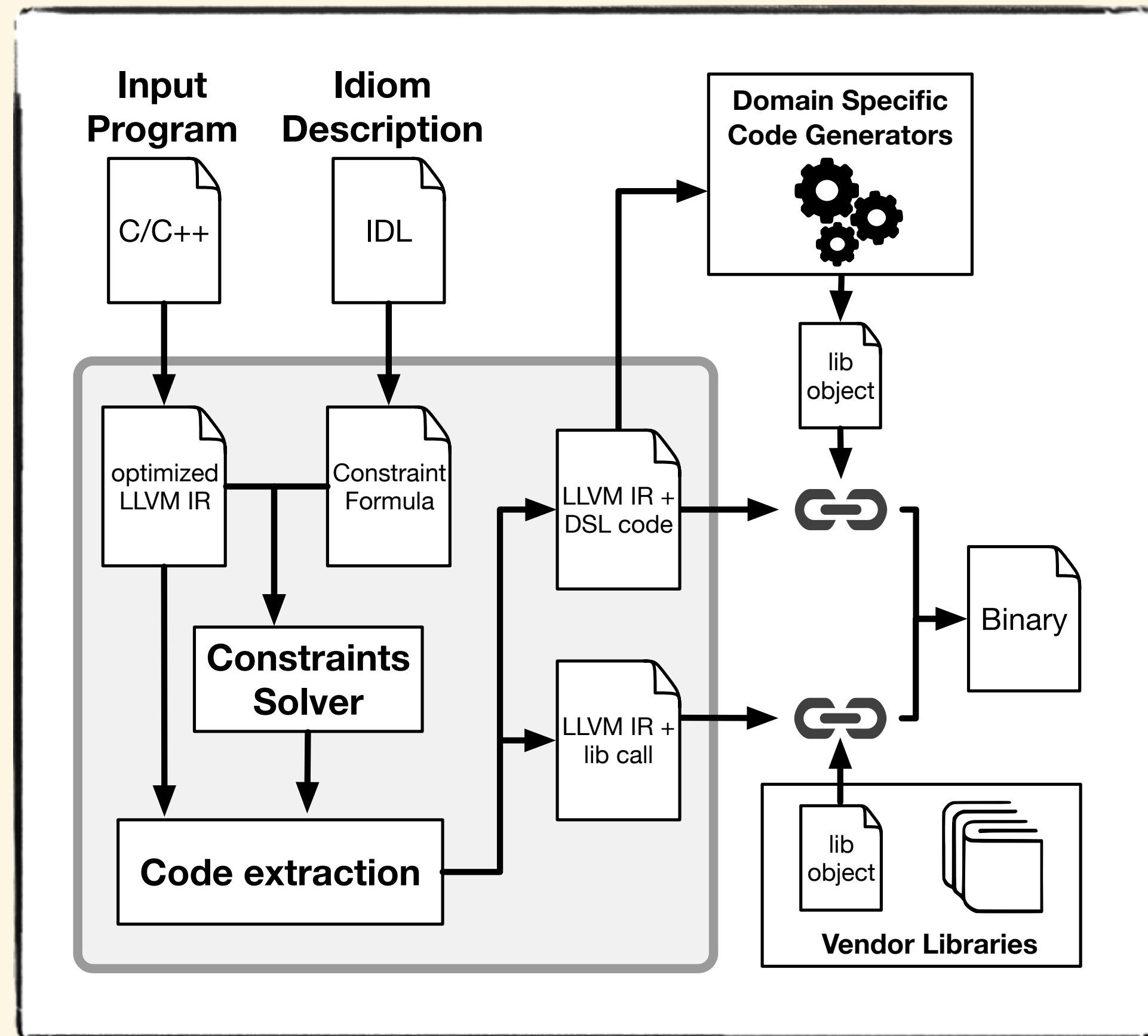
This paper develops a novel approach to automatically

## 1 Introduction

Heterogeneous accelerators provide the potential for great performance. However, achieving that potential is difficult. General purpose languages such as OpenCL [36] provide portability, but the achieved performance often disappoints [29]. This shortfall has led vendors to deliver specialized libraries to bridge the gap [2]. Alternatively, domain specific

# IDIOM DETECTION VIA CONSTRAINT LANGUAGE

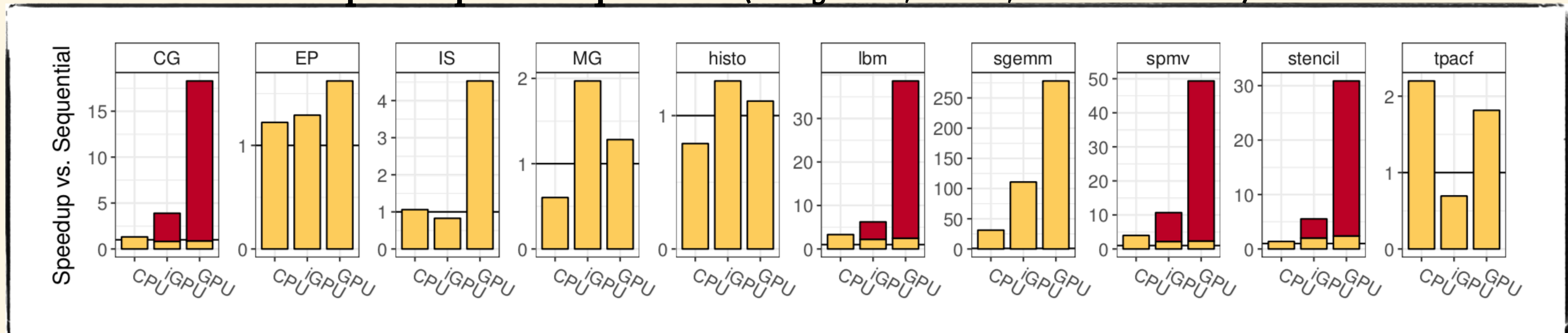# PERFORMANCE RESULTS

## Runtime Coverage of detected Idioms (NAS PB + Parboil)



## Speedup vs. Sequential (using BLAS, Halide, Lift as backends)