# Regular Flattening

Troels Henriksen (athas@sigkill.dk)

DIKU
University of Copenhagen

# Agenda

Representation and Fusion

Handling nested parallelism

Basic flattening rules

Incremental flattening

Multi-level parallelism

Final words as time permits

## Representation and Fusion

Handling nested parallelism

Basic flattening rules

Incremental flattening

Multi-level parallelism

Final words as time permits

# "Unzipped" SOACs

## Representation

An array `[](t1, t2, t3...)` is represented in memory as `([]t1, []t2, []t3...)`, i.e. as *multiple arrays*, each containing only primitive values.

Instead of

```
let tmp = map (\(x,y) -> (x-1, y+1))
              (zip xs ys)
let (xs, ys) = unzip xs_ys'
```

we write

```
let (xs, ys) = map (\x y -> (x-1, y+1)) xs ys
```

- In the compiler, **All SOACs accept multiple inputs and produce unzipped results.**
- Arrays of tuples (or records, or sums) do not exist in the core language.
- **Isomorphic to source language**, but this form is simpler in a compiler.

## Loop fusion

```
def increment [n][m] (as:  [n][m]i32) : [n]i32 =
  map (\r -> map (+2) r) a
def sum [n] (a:  [n]i32) : i32 =
  reduce (+) 0 a
def sumrows [n][m] (as: [n][m]i32) : [n]i32 =
  map sum as
```

Suppose we wish to first call increment, then sumrows:

$$sumrows \ (increment \ a)$$

Naively Run increment, then call sumrows.

Problem Manifests intermediate matrix in memory.

Solution *Loop fusion*, which combines loops to avoid intermediate results.

## An example of a fusion rule

The expression

$$\textbf{map } f \ (\textbf{map } g \ a)$$

is *always* equivalent to

$$\textbf{map } (f \circ g) \ a$$

- This is an extremely powerful property that is only true in the absence of side effects.
- Fusion is *the* core optimisation that permits the efficient decomposition of a data parallel program.
- A full fusion engine has much more awkward rules (mostly bookkeeping related to fusing only *some* of several inputs), but safety is guaranteed.

## A fusion example

$$\text{sumrows (increment } a) = \qquad \text{(Initial expression)}$$
$$\textbf{map sum (increment } a) = \qquad \text{(Inline sumrows)}$$
$$\textbf{map sum } (\textbf{map } (\lambda r \rightarrow \textbf{map } (+2)\ r)\ a) = \qquad \text{(Inline increment)}$$
$$\textbf{map } (\text{sum} \circ (\lambda r \rightarrow \textbf{map } (+2)\ r)\ a) = \qquad \text{(Apply \textbf{map}-\textbf{map} fusion)}$$
$$\textbf{map } (\lambda r \rightarrow \text{sum } (\textbf{map } (+2)\ r)\ a) = \qquad \text{(Apply composition)}$$

- We have avoided the temporary matrix, but the composition of sum and the **map** also holds an opportunity for fusion – specifically, **reduce**-**map** fusion.
- Will not cover in detail, but a **reduce** can efficiently apply a function to each input element before engaging in the actual reduction operation.
- Important to remember: a **map** going into a **reduce** is an efficient pattern.

## A shorthand notation for sequences

$$\overline{z}^{(n)} = z_0, \cdots, z_{(n-1)}$$

- The $n$ may be omitted.
- A separator may be implied by context.

$$f\ \overline{v}^{(n)} \equiv f\ v_1\ \cdots\ v_n$$

or a tuple

$$(\overline{v}^{(n)}) \equiv (v_1, \ldots, v_n)$$

or a function type

$$\overline{\tau}^{(n)} \to \tau_{n+1} \equiv \tau_1 \to \cdots \to \tau_n \to \tau_{n+1}.$$

When not all terms under the bar are variant, subscript variant terms with $i$.

$$(\overline{[d]v_i}^{(n)}) = ([d]v_1, \ldots, [d]v_n)$$

and

$$(\overline{[d_i]v_i}^{(n)}) = ([d_1]v_1, \ldots, [d_n]v_n)$$

## Fused constructs

### Convenient shorthands

$$\textbf{redomap} \odot f\,(\overline{d})\,\overline{xs} \equiv \qquad \textbf{reduce} \odot (\overline{d})\,(\textbf{map}\,f\,\overline{xs})$$

$$\textbf{scanomap} \odot f\,(\overline{d})\,\overline{xs} \equiv \qquad \textbf{scan} \odot (\overline{d})\,(\textbf{map}\,f\,\overline{xs})$$

- Emphasises that **reduce**/**scan**-**map** compositions can be considered as a single construct.
- We will see several examples where this is useful.

## Fused constructs

### Convenient shorthands

$$\textbf{redomap} \odot f\ (\overline{d})\ \overline{xs} \equiv \qquad \textbf{reduce} \odot (\overline{d})\ (\textbf{map}\ f\ \overline{xs})$$

$$\textbf{scanomap} \odot f\ (\overline{d})\ \overline{xs} \equiv \qquad \textbf{scan} \odot (\overline{d})\ (\textbf{map}\ f\ \overline{xs})$$

- Emphasises that **reduce**/**scan-map** compositions can be considered as a single construct.
- We will see several examples where this is useful.

**Note:**

$$\textbf{reduce} \odot (\overline{d})\ \overline{xs} \equiv \quad \textbf{reduce} \odot (\overline{d})\ (\textbf{map id}\ \overline{xs}) \equiv \quad \textbf{redomap} \odot \textbf{id}\ (\overline{d})\ \overline{xs}$$

$$\textbf{scan} \odot (\overline{d})\ \overline{xs} \equiv \quad \textbf{scan} \odot (\overline{d})\ (\textbf{map id}\ \overline{xs}) \equiv \quad \textbf{scanomap} \odot \textbf{id}\ (\overline{d})\ \overline{xs}$$
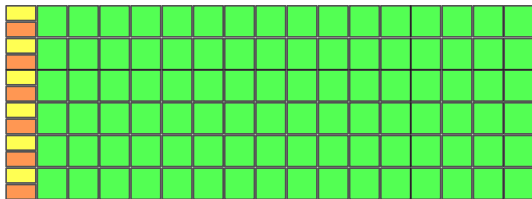
CPU       GPU

- GPUs have *thousands* of simple cores and taking full advantage of their compute power requires *tens of thousands* of threads.
- GPU threads are very *restricted* in what they can do: no stack, no allocation, limited control flow, etc.
- Potential *very high performance* and *lower power usage* compared to CPUs, but programming them is *hard*.

## The SIMT Programming Model



- GPUs are programmed using the SIMT model (*Single Instruction Multiple Thread*).
- Similar to SIMD (*Single Instruction Multiple Data*), but while SIMD has explicit vectors, we provide *sequential scalar per-thread* code in SIMT.

Each thread has its own registers, but they all execute the same instructions at the same time (i.e. they share their instruction pointer).

## SIMT example

For example, to increment every element in an array a, we might use this code:

```
increment(a) {
  tid = get_thread_id();
  x = a[tid];
  a[tid] = x + 1;
}
```

- If a has n elements, we launch n threads, with get_thread_id() returning *i* for thread *i*.
- This is *data-parallel programming*: applying the same operation to different data.
- When we launch a GPU program (*kernel*), we say how many threads should be launched, *all running the same code.*

## Branching

If all threads share an instruction pointer, what about branches?
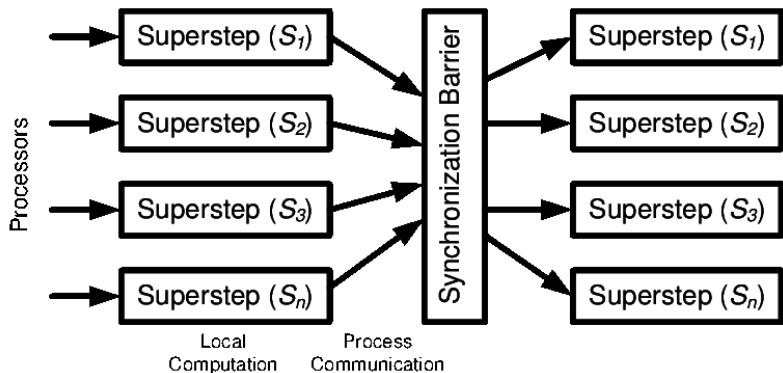
```
mapabs(a) {
  tid = get_thread_id();
  x = a[tid];
  if (x < 0) {
    a[tid] = -x;
  }
}
```

### Masked Execution

Both branches are executed in all threads, but in those threads where the condition is false, a mask bit is set to treat the instructions inside the branch as no-ops.

## Do GPUs exist in theory as well?

GPU programming is a close fit to the *bulk synchronous parallelism* model:



- Supersteps are *threads*, which cannot talk to each other.
- The synchronisation barriers are kernel launches.

[1]Illustration by Aftab A. Chandio.

## A SOAC-kernel correspondence

The compiler *knows*[2] that certain nestings of **map**s correspond to GPU basic blocks.

- **map**s containing scalar code is a kernel with one thread per iteration of the **map**s.
- **map**s containing a single **reduce** is a *segmented reduction*.
- **map**s containing a single **scan** is a *segmented scan*.
- **map**s containing a single **scatter** is a *segmented scatter*.
- ...see the pattern?

**Crucial**: the **map**s must be *perfectly nested* around the operation.

### Perfect nesting of an operation $e_o$

An expression $e$ is a *perfect nesting* of $e_o$ if $e$ has form

$$\textbf{map}\ (\lambda\overline{p} \to e_f)\ \overline{x}$$

where either $e_f = e_o$ or $e_f$ is a perfect nesting of $e_o$.

---

[2] Because it was taught it by Cosmin in PMPH.

## Example

```
map (\xs y ->
        map (\x ->
                x + y)
     xs)
   xss ys
```

- Suppose xss is of shape [n][m].
- This could be compiled to a kernel with $n \times m$ threads, each doing a single $x + y$ operation.

### Problem

Futhark permits *nested* parallelism, but GPUs need *flat* parallel *kernels*.

# Handling nested parallelism

## Problem

Futhark permits *nested* parallelism, but GPUs need *flat* parallel *kernels*.

## Solution

Have the compiler rewrite program to perfectly nested **map**s containing sequential operations, or known parallel patterns such as segmented reduction.

## Handling nested parallelism

### Problem

Futhark permits *nested* parallelism, but GPUs need *flat* parallel *kernels*.

### Solution

Have the compiler rewrite program to perfectly nested **map**s containing sequential operations, or known parallel patterns such as segmented reduction.

```
map (\xs -> let y = reduce (+) 0 xs
            in map (\x -> x + y) xs)
     xss
                   ⇓
let ys = map (\xs -> reduce (+) 0 xs) xss
in map (\xs y -> map (\x -> x + y) xs) xss ys
```

## Flattening via loop fission

The classic map fusion rule:

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

---

[3] *Futhark: Purely Functional GPU-Programming with Nested Parallelism and In-Place Array Updates*, PLDI 2017

## Flattening via loop fission

The classic map fusion rule:

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

We can also apply it backwards to obtain *fission*:

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

This, along with other fission rules (see paper[3]), are applied by the compiler to extract perfect map nests.

---

[3] *Futhark: Purely Functional GPU-Programming with Nested Parallelism and In-Place Array Updates*, PLDI 2017

## Example: (a) Initial program, we inspect the map-nest

```
let (asss, bss) =
  map (\(ps: [m]i32) ->
        let ass = map (\(p: i32): [m]i32 ->
                        let cs = scan (+) 0 (iota p)
                        let r = reduce (+) 0 cs
                        in map (+r) ps) ps
        let bs = loop ws=ps for i < n do
                  map (\as w: i32 ->
                        let d = reduce (+) 0 as
                        let e = d + w
                        in 2 * e) ass ws
        in (ass, bs)) pss
```

We assume the type of pss : [m][m]i32.

## (b) Distribution

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
         let ass = map (\(p: i32): [m]i32 ->
                          let cs = scan (+) 0 (iota p)
                          let r = reduce (+) 0 cs
                          in map (+r) ps) ps
         in ass) pss
let bss: [m][m]i32 =
  map (\ps ass ->
         let bs = loop ws=ps for i < n do
                    map (\as w ->
                           let d = reduce (+) 0 as
                           let e = d + w
                           in 2 * e) ass ws
         in bs) pss asss
```

## (c) Interchanging outermost map inwards

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
        let ass = map (\(p: i32): [m]i32 ->
                        let cs = scan (+) 0 (iota p)
                        let r = reduce (+) 0 cs
                        in map (+r) ps) ps
        in ass) pss
let bss: [m][m]i32 =
  map (\ps ass ->
        let bs = loop ws=ps for i < n do
                    map (\as w ->
                          let d = reduce (+) 0 as
                          let e = d + w
                          in 2 * e) ass ws
        in bs) pss asss
```

## (c) Interchanging outermost map inwards

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
          let ass = map (\(p: i32): [m]i32 ->
                            let cs = scan (+) 0 (iota p)
                            let r = reduce (+) 0 cs
                            in map (+r) ps) ps
          in ass) pss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    map (\ass ws ->
            let ws' = map (\as w ->
                              let d = reduce (+) 0 as
                              let e = d + w
                              in 2 * e) ass ws
            in ws') asss wss
```

## (d) Skipping scalar computation

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
         let ass = map (\(p: i32): [m]i32 ->
                          let cs = scan (+) 0 (iota p)
                          let r = reduce (+) 0 cs
                          in map (+r) ps) ps
         in ass) pss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    map (\ass ws ->
           let ws' = map (\as w ->
                            let d = reduce (+) 0 as
                            let e = d + w
                            in 2 * e) ass ws
           in ws') asss wss
```

## (d) Skipping scalar computation

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
          let ass = map (\(p: i32): [m]i32 ->
                              let cs = scan (+) 0 (iota p)
                              let r = reduce (+) 0 cs
                              in map (+r) ps) ps
          in ass) pss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    map (\ass ws ->
            let ws' = map (\as w ->
                              let d = reduce (+) 0 as
                              let e = d + w
                              in 2 * e) ass ws
            in ws') asss wss
```

## (e) Distributing reduction

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
        let ass = map (\(p: i32): [m]i32 ->
                        let cs = scan (+) 0 (iota p)
                        let r = reduce (+) 0 cs
                        in map (+r) ps) ps
        in ass) pss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    map (\ass ws ->
          let ws' = map (\as w ->
                          let d = reduce (+) 0 as
                          let e = d + w
                          in 2 * e) ass ws
          in ws') asss wss
```

## (e) Distributing reduction

```
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) ->
         let ass = map (\(p: i32): [m]i32 ->
                          let cs = scan (+) 0 (iota p)
                          let r = reduce (+) 0 cs
                          in map (+r) ps) ps
         in ass) pss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    let dss: [m][m]i32 =
      map (\ass ->
             map (\as ->
                    reduce (+) 0 as) ass)
           asss
    in map (\ws ds ->
              let ws' =
                map (\w d -> let e = d + w
                             in 2 * e) ws ds
              in ws') asss dss
```

## (f) Distributing inner map

```
let asss =
  map (\(ps: [m]i32) ->
         let ass = map (\(p: i32): [m]i32 ->
                          let cs = scan (+) 0 (iota p)
                          let r = reduce (+) 0 cs
                          in map (+r) ps) ps
         in ass) pss
let bss: [m][m]i32 = ...
```

## (f) Distributing inner map

```
let rss: [m][m]i32 =
  map (\(ps: [m]i32) ->
        let rs = map (\(p: i32): i32 ->
                        let cs = scan (+) 0 (iota p)
                        let r = reduce (+) 0 cs
                        in r) ps
        in rs) pss
let asss: [m][m][m]i32 =
  map (\(ps: [m]i32) (rs: [m]i32) ->
        map (\(r: i32): [m]i32 ->
              map (+r) ps) rs
      ) pss rss
let bss: [m][m]i32 = ...
```

## (g) Cannot distribute as it would create irregular array

```
let rss: [m][m]i32 =
  map (\(ps: [m]i32) ->
          let rs = map (\(p: i32): i32 ->
                           let cs = scan (+) 0 (iota p)
                           let r = reduce (+) 0 cs
                           in r) ps
          in rs) pss
let asss: [m][m][m]i32 = ...
let bss: [m][m]i32 = ...
```

Array cs has type [p]i32, and p is variant to the innermost map nest.

## (h) These statements are sequentialised

```
let rss: [m][m]i32 =
  map (\(ps: [m]i32) ->
        let rs = map (\(p: i32): i32 ->
                        let cs = scan (+) 0 (iota p)
                        let r = reduce (+) 0 cs
                        in r) ps
        in rs) pss
let asss: [m][m][m]i32 = ...
let bss: [m][m]i32 = ...
```

Array cs has type [p]i32, and p is variant to the innermost map nest.

## Result

```
let rss: [m][m]i32 = map (\ps -> map (...) ps) pss
let asss: [m][m][m]i32 =
  map (\ps rs -> map (\r -> map (...) ps) rs) pss rss
let bss: [m][m]i32 =
  loop wss=pss for i < n do
    let dss: [m][m]i32 = map (\ass -> map (reduce ...) ass)
                             asss
    in map (\ws ds -> map (...) ws ds ) asss dss
```

- From a single kernel with parallelism $m$ to four kernels of parallelism $m^2$, $m^3$, $m^3$, and $m^2$.
- The last two kernels are executed $n$ times each.

## Notation for flat parallelism

Instead of

```
map (\ps rs ->
  map (\r ->
    map (\p -> e)
      ps)
    rs)
  pss rss
```

we write

$$\textbf{segmap} \, (\langle ps, rs \in pss, rss \rangle, \, \langle r \in rs \rangle, \, \langle p \in ps \rangle)$$
$$e$$

## Segmented flat parallel constructs

$$\Sigma = \Sigma', \langle \overline{x} \in \overline{y} \rangle$$

$$
\textbf{segmap } \Sigma \ e \equiv \quad \textbf{map } (\lambda \overline{x_p} \rightarrow \\
\textbf{map } (\lambda \overline{x_{p-1}} \rightarrow \ldots \\
\textbf{map } (\lambda \overline{x_1} \rightarrow e) \ \overline{y_1}) \\
\overline{y_{p-1}}) \\
\overline{y_p}
$$

- Conceptually a perfect nest of **map**s with some operation inside.
- *These* are what trigger GPU code generation.
- Any SOACs left in *e* will be executed sequentially.

## Similarly for reductions and scans

$$\textbf{segred } \Sigma \odot \overline{d} \, e \equiv \ \textbf{map } (\lambda \overline{x_p} \rightarrow$$
$$\textbf{map } (\lambda \overline{x_{p-1}} \rightarrow \ldots$$
$$\textbf{redomap } \odot \ (\lambda \overline{x_1} \rightarrow e) \, (\overline{d}) \ \overline{y_1})$$
$$\overline{y_{p-1}})$$
$$\overline{y_p}$$

$$\textbf{segscan } \Sigma \odot \overline{d} \, e \equiv \ \textbf{map } (\lambda \overline{x_p} \rightarrow$$
$$\textbf{map } (\lambda \overline{x_{p-1}} \rightarrow \ldots$$
$$\textbf{scanomap } \odot \ (\lambda \overline{x_1} \rightarrow e) \, \overline{d} \ \overline{y_1})$$
$$\overline{y_{p-1}})$$
$$\overline{y_p}$$

**Let us look at how one can rewrite SOAC nests to these segmented operations.**

A flattening rewrite

$$\mathcal{G}(\Sigma, e) \Rightarrow e'$$

says that flattening an operation $e$ nested inside a map nest $\Sigma$ produces the expression $e'$.

- *Rewrite rules* define which rewrites are valid.
- Not an algorithm: multiple rewrites are often possible.
- Form in is a bit different and have more details; I have tried to keep it more high level.

## Basic rules

$$\mathcal{G}(\bullet, e) \Rightarrow e \tag{G0}$$

$$\mathcal{G}(\Sigma, e) \Rightarrow \textbf{segmap } \Sigma \ e \tag{G1}$$

$$\mathcal{G}(\Sigma, \textbf{map } (\lambda \overline{x} \to e) \ \overline{xs}) \Rightarrow \mathcal{G}(\Sigma \ \langle \overline{x} \rangle \in \overline{xs}, e) \tag{G2}$$

$$\mathcal{G}(\Sigma, \textbf{redomap } \odot \ (\lambda \overline{x} \to e) \ (\overline{d}) \ \overline{xs}) \Rightarrow \textbf{segred } (\Sigma \ \overline{x} \in \overline{xs}) \ \odot \ \overline{d} \ e \tag{Gr}$$

G0: at top level, choose not to flatten.
G1: inside a nest, emit the corresponding **segmap**.
G2: descend into a **map**.
Gr: for a nested **redomap**, emit a **segred**.

## Distribution rule

$$\mathcal{G}(\Sigma, \textbf{let } \overline{a_0} = e_1 \textbf{ in } e_2) \Rightarrow \textbf{let } \overline{a_p} = \mathcal{G}(\Sigma, e_1) \textbf{ in } \mathcal{G}(\Sigma', e_2) \qquad \text{(G6)}$$

where

$$
\begin{aligned}
\Sigma &= \langle \overline{x_p} \in \overline{y_p} \rangle, \ldots, \langle \overline{x_1} \in \overline{y_1} \rangle \\
\Sigma' &= \langle \overline{x_p} \, \overline{a_{p-1}} \in \overline{y_p} \, \overline{a_p} \rangle, \ldots, \langle \overline{x_1} \, \overline{a_0} \in \overline{y_1} \, \overline{a_1} \rangle
\end{aligned}
$$

and $\overline{a_p}, \ldots, \overline{a_1}$ are fresh names.

**Note:** only applicable when each array in $\overline{a_0}$ is invariant to $\Sigma$.

## Example

We are flattening

$$e = \textbf{map} \ (\lambda \mathsf{xs} \rightarrow \quad \textbf{let} \ \mathsf{y} = \textbf{redomap} \ (+) \ (\lambda \mathsf{x} \rightarrow \mathsf{x}) \ 0 \ \mathit{xs}$$
$$\textbf{in} \ \textbf{map} \ (\lambda \mathsf{x} \rightarrow \mathsf{x} + \mathsf{y}) \ \mathsf{xs})$$
$$\mathsf{xss}$$

## Example

We are flattening

$$e = \textbf{map } (\lambda \text{xs} \rightarrow \ \textbf{let } y = \textbf{redomap } (+) \ (\lambda \text{x} \rightarrow \text{x}) \ 0 \ xs$$
$$\textbf{in map } (\lambda \text{x} \rightarrow \text{x} + \text{y}) \ \text{xs})$$
$$\text{xss}$$

By applying rule G1 we get

$$\mathcal{G}(\bullet, e) = \mathcal{G}(\langle \text{xs} \in \text{xss} \rangle, \ \textbf{let } y = \textbf{redomap } (+) \ (\lambda \text{x} \rightarrow \text{x}) \ 0 \ xs \ )$$
$$\textbf{in map } (\lambda \text{x} \rightarrow \text{x} + \text{y}) \ \text{xs})$$

$$\mathcal{G}(\langle \mathsf{xs} \in \mathsf{xss} \rangle, \ \textbf{let } \mathsf{y} = \textbf{redomap } (+) \ (\lambda \mathsf{x} \to \mathsf{x}) \ 0 \ \mathit{xs})$$
$$\textbf{in map } (\lambda \mathsf{x} \to \mathsf{x} + \mathsf{y}) \ \mathsf{xs})$$

Apply distribution rule by:

$$
\begin{aligned}
\overline{a_0} &= & \mathsf{y} \\
e_1 &= & \textbf{redomap } (+) \ (\lambda \mathsf{x} \to \mathsf{x}) \ 0 \ \mathsf{xs} \\
e_2 &= & \textbf{map } (\lambda \mathsf{x} \to \mathsf{x} + \mathsf{y}) \ \mathsf{xs} \\
\Sigma &= & \langle \mathsf{xs} \in \mathsf{xss} \rangle \\
\Sigma' &= & \langle \mathsf{xs}, \mathsf{y} \in \mathsf{xss}, \mathsf{ys} \rangle \\
\mathcal{G}(\Sigma, e) &= & \textbf{let } \mathsf{ys} = \mathcal{G}(\Sigma, e_1) \ \textbf{in } \mathcal{G}(\Sigma', e_2) \\
&= & \textbf{let } \mathsf{ys} = \textbf{segred } (\langle \mathsf{xs} \in \mathsf{xss} \rangle, \langle \mathsf{x} \in \mathsf{xs} \rangle) \ (\texttt{+}) \ 0 \ \mathsf{x} \\
& & \textbf{in segmap } \langle \mathsf{xs}, \mathsf{y} \in \mathsf{xss}, \mathsf{ys} \rangle \ (\mathsf{x} + \mathsf{y})
\end{aligned}
$$

### Reminder: distribution rule

$$\mathcal{G}(\Sigma, \textbf{let } \overline{a_0} = e_1 \ \textbf{in } e_2) \Rightarrow \textbf{let } \overline{a_p} = \mathcal{G}(\Sigma, e_1) \ \textbf{in } \mathcal{G}(\Sigma', e_2) \qquad \text{(G6)}$$

## Flattening transposition

**rearrange** $(d_1, \cdots, d_n)$ *x* is a generalization of **transpose** in that it rearranges the dimensions of *d*-dimensional array based on a permutation defined by the integer sequence $d_1, \cdots, d_n$. E.g:

$$\textbf{transpose} \equiv \textbf{rearrange}\ (1, 0)$$

## Flattening transposition

**rearrange** $(d_1, \cdots, d_n)$ *x* is a generalization of **transpose** in that it rearranges the dimensions of *d*-dimensional array based on a permutation defined by the integer sequence $d_1, \cdots, d_n$. E.g:

$$\textbf{transpose} \equiv \textbf{rearrange}\ (1, 0)$$

**Flattening rule**

$$\mathcal{G}(\Sigma, \langle x \in y \rangle, \textbf{rearrange}\ (k_1, \ldots, k_n)\ x) \Rightarrow \mathcal{G}(\Sigma, \textbf{rearrange}\ (0, 1+k_1, \ldots, 1+k_n)\ y)$$

Eventually reaches base case where $\Sigma = \bullet$.

## map-loop interchange

$$\mathcal{G}(\Sigma \; \langle \overline{x} \, \overline{y} \in \overline{xs} \; \overline{ys} \rangle, \textbf{loop } \overline{z'} \; \overline{y'} = \overline{z} \; \overline{y} \textbf{ for } i < n \textbf{ do } f \; i \; \overline{q} \; \overline{x} \; \overline{y} \; \overline{y'} \; \overline{z'}) \Rightarrow$$
$$\mathcal{G}(\Sigma, \textbf{loop } \overline{zs'} \; \overline{ys'} = \overline{z^r} \; \overline{ys} \textbf{ for } i < n \textbf{ do map } (f \; i \; \overline{q}) \; \overline{xs} \; \overline{ys} \; \overline{ys'} \; \overline{zs'})$$

where

$$m = \qquad \text{outer size of each of } \overline{xs} \text{ and } \overline{ys}$$
$$\overline{z^r} = \qquad \overline{\textbf{replicate } m \; z_i}$$
$$\{n, \overline{q}, \overline{z}\} \cap \{\overline{x}, \overline{y}\} = \qquad \emptyset$$

and $\overline{zs'}$ and $\overline{ys'}$ are fresh names.

## Informal example of interchange

```
map (\ xs -> loop (xs', j) = (xs, 0) for i < n do
              (map (+j) xs', j + i))
    xss
```

*Becomes after interchange*

```
loop (xss', js) = (xss, replicate m 0) for i < n do
  map (\ xs' j -> (map (+j) xs', j + i))
      xss' js
```

## Validity of interchange

The simple intuition is that

**map** (\x -> **loop** x' = x **for** i < n **do** f x') xs

is equivalent to

**loop** xs' = xs **for** i < n **do map** f xs'

because they both produce

$$[f^n \, xs[0], \ldots, f^n \, xs[m-1]]$$

## Consider Matrix Multiplication

```
for i < n:
  for j < m:
    acc = 0
    for l < p:
      acc += xss[i,l] * yss[l,j]
    res[i,j] = acc
```

```
map (\xs ->
     map (\ys ->
           let zs = map (*) xs ys
           in reduce (+) 0 zs)
          (transpose yss))
    xss
```

## Using `redomap` notation

```
map (\xs ->
     map (\ys ->
          redomap (+) (*) 0 xs ys)
         (transpose yss))
    xss
```

$$\textbf{redomap} \odot f \, 0_\odot \, x \; \equiv \; \textbf{reduce} \odot 0_\odot \, (\textbf{map} \, f \, x)$$

Emphasises that a **map**-**reduce** composition can be turned into a fused tight sequential loop, or into a parallel reduction.

# So how should we parallelise this on GPU?

## So how should we parallelise this on GPU?

*Full flattening*

```
map (\ xs ->
  map (\ ys ->
    redomap (+) (*) 0 xs ys)
    (transpose yss))
  xss
```

- **All parallelism exploited**
- Some communication overhead
- *Best if outer maps don't saturate GPU*

## So how should we parallelise this on GPU?

*Full flattening*

```
map (\ xs ->
  map (\ ys ->
    redomap (+) (*) 0 xs ys)
    (transpose yss))
  xss
```

- **All parallelism exploited**
- Some communication overhead
- *Best if outer maps don't saturate GPU*

*Moderate flattening*

```
map (\ xs ->
  map (\ ys ->
    redomap (+) (*) 0 xs ys)
    (transpose yss))
  xss
```

- **Only outer parallelism**
- The `redomap` can be block tiled
- *Best if outer maps saturate GPU*

- There is no *one size fits all*.
- Both situations may be encountered at program runtime.

## The essence of *incremental flattening*

**From a single source program, for each parallel construct generate multiple
*semantically equivalent* parallelisations, and generate a *single program* that at runtime
picks the *least parallel* that still saturates the hardware.**

- Implemented in the Futhark compiler.
- ...but technique is applicable to any (regular) nested parallelism expressed with
  the common functional array combinators (map, reduce, scan, etc).

## Simple Incremental Flattening

At every level of map-nesting we have two options:

1. Continue flattening inside the map, exploiting the parallelism there.
2. Sequentialise the map body; exploiting only the parallelism on top.

- **Full flattening** in the Blelloch style will do the former, maximising utilised parallelism.
- **Incremental flattening** generates *both* versions and uses a predicate to pick at runtime.
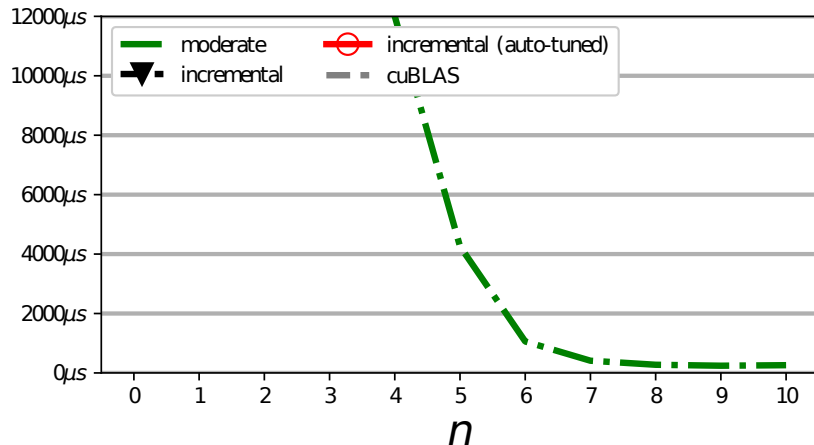
## Multi-versioned matrix multiplication

```
xss : [n][p]i32
yss : [p][m]i32.
if n * m > t0 then
  map (\ xs ->
          map (\ ys ->
                  redomap (+) (*) 0 xs ys)
              (transpose yss))
      xss
else
  map (\ xs ->
          map (\ ys ->
                  redomap (+) (*) 0 xs ys)
              (transpose yss))
      xss
```
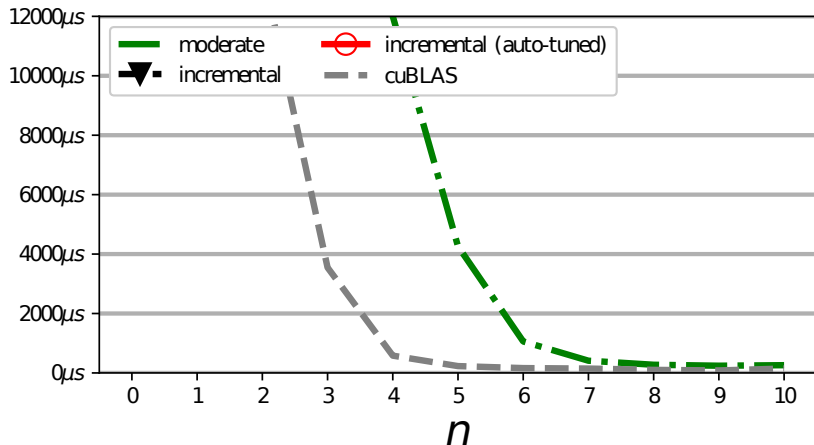
The $t_0$ *threshold parameter* is used to select between the two versions—and should be auto-tuned on the concrete hardware.
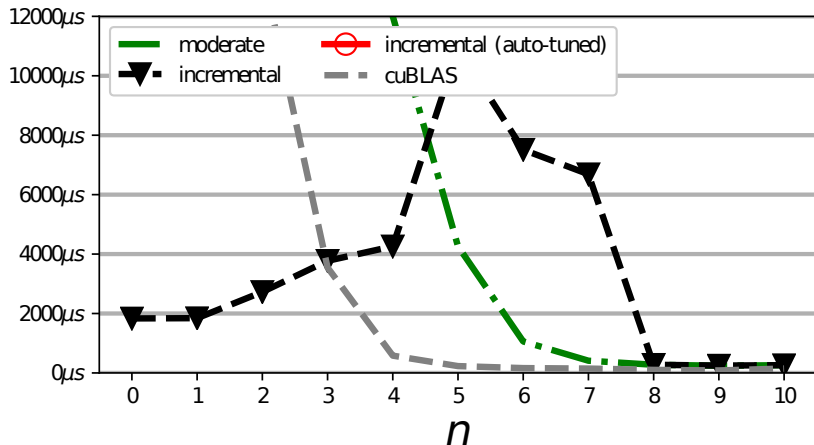
# Matrix multiplication on NVIDIA K40



Multiplying matrices of size $2^n \times 2^m$ and $2^m \times 2^n$, where $m = 25 - 2n$, meaning that work is constant as we vary $n$.

## Matrix multiplication on NVIDIA K40
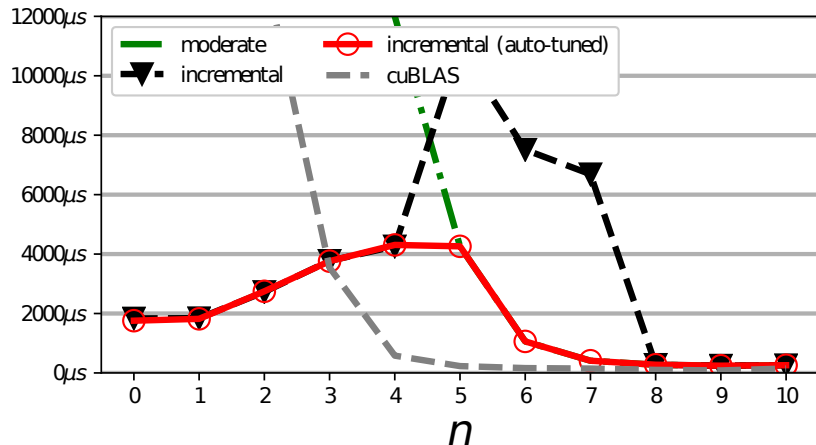
Multiplying matrices of size $2^n \times 2^m$ and $2^m \times 2^n$, where $m = 25 - 2n$, meaning that work is constant as we vary $n$.

## Matrix multiplication on NVIDIA K40

Multiplying matrices of size $2^n \times 2^m$ and $2^m \times 2^n$, where $m = 25 - 2n$, meaning that work is constant as we vary $n$.

## Matrix multiplication on NVIDIA K40



Multiplying matrices of size $2^n \times 2^m$ and $2^m \times 2^n$, where $m = 25 - 2n$, meaning that work is constant as we vary $n$.

## Incremental flattening rule

$$\mathcal{G}(\Sigma, \textbf{map } (\lambda \overline{x} \to e) \, \overline{\textsf{xs}}) \Rightarrow \begin{array}{l} \textbf{if } \text{Parallelism}(\Sigma') \geq t_{\text{top}} \\ \textbf{then segmap } \Sigma' \, e \\ \textbf{else } \mathcal{G}(\Sigma', e) \end{array}$$

where $\Sigma' = \Sigma, \langle \overline{x} \in \overline{\textsf{xs}} \rangle$.

**Example for**

```
map (\xs -> redomap (+) (\x -> x) 0 xs) xss
```

## Incremental flattening rule

$$\mathcal{G}(\Sigma, \textbf{map } (\lambda \overline{x} \to e) \ \overline{\text{xs}}) \Rightarrow \begin{array}{l} \textbf{if } \text{Parallelism}(\Sigma') \geq t_{\text{top}} \\ \textbf{then segmap } \Sigma' \ e \\ \textbf{else } \mathcal{G}(\Sigma', e) \end{array}$$

where $\Sigma' = \Sigma, \langle \overline{x} \in \overline{\text{xs}} \rangle$.

### Example for

```
map (\ xs -> redomap (+) (\x -> x) 0 xs) xss
```

$$\Sigma = \bullet \quad \Sigma' = \langle \text{xs} \in \text{xss} \rangle$$
$$\Sigma' \vdash e \Rightarrow \textbf{segred} \ (\langle \text{xs} \in \text{xss} \rangle, \langle x \in \text{xs} \rangle) \ (+) \ 0 \ x$$
$$\Sigma \vdash ... \Rightarrow \begin{array}{l} \textbf{if } \text{length}(\text{xss}) \geq t_{\text{top}} \\ \textbf{then segmap } \langle \text{xs} \in \text{xss} \rangle \ (\texttt{redomap } (\texttt{+}) \ (\lambda x \to x) \ 0 \ \text{xs}) \\ \textbf{else segred} \ (\langle \text{xs} \in \text{xss} \rangle, \langle x \in \text{xs} \rangle) \ (+) \ 0 \ x \end{array}$$

# Autotuning

- An incrementally flattened program may have dozens of threshold parameters, $t_i$, used to select versions at runtime.
- As we have seen, the default value ($2^{16}$) is often not optimal.
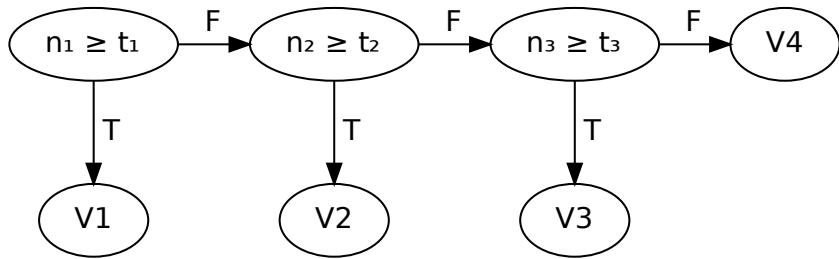
A *configuration P* maps each $t_i$ to an integer $P(t_i)$.

### The search problem

Find the *P* that minimises the cost function $F(P)$, where the the cost function runs the program on a set of user-provided representative datasets and sums the observed runtimes.

- Other cost functions are also possible, e.g. average runtime over datasets.
- **Note:** recompilation is not necessary.

# Briefly on the search procedure[4]



- Suppose we are given training data sets $D_j, j < k$, each of which provide a value $v_{i,j}$ for each threshold parameter $n_i$.
- Starting from the deepest comparison ($t_3$), for each $D_j$ find an $(x_j, y_j)$ that minimises runtime, take the intersection of the intervals, and use that to determine threshold value.
- Tuning time is linear in the number of comparisons.

[4]https://futhark-lang.org/publications/tfp21.pdf

## Using incremental flattening

Compile with a GPU backend (opencl or cuda):

```
$ futhark opencl matmul.fut
```

To autotune:

```
$ futhark autotune -v --backend=opencl matmul.fut
```

Produces matmul.fut.tuning, which is automatically picked up by futhark bench (use --no-tuning to stop this).

Use futhark dev -s --extract-kernels -e matmul.fut to see IR.

## Confession

**I lied when I claimed that GPU threads were completely isolated.**

## Confession

**I lied when I claimed that GPU threads were completely isolated.**

- Most hardware has useful (fixed) levels of parallelism.
- An ideal flattening algorithm maps levels of application parallelism (any number) to hardware parallelism (fixed number) in a way that exploits locality well.

**Example of deep nesting:** a system consists of multiple *datacenters*, that each contain multiple *computers*, that each contain multiple *GPUs*, that each contain multiple *SMs* (next slide), that each run some number of threads.

## Confession

**I lied when I claimed that GPU threads were completely isolated.**
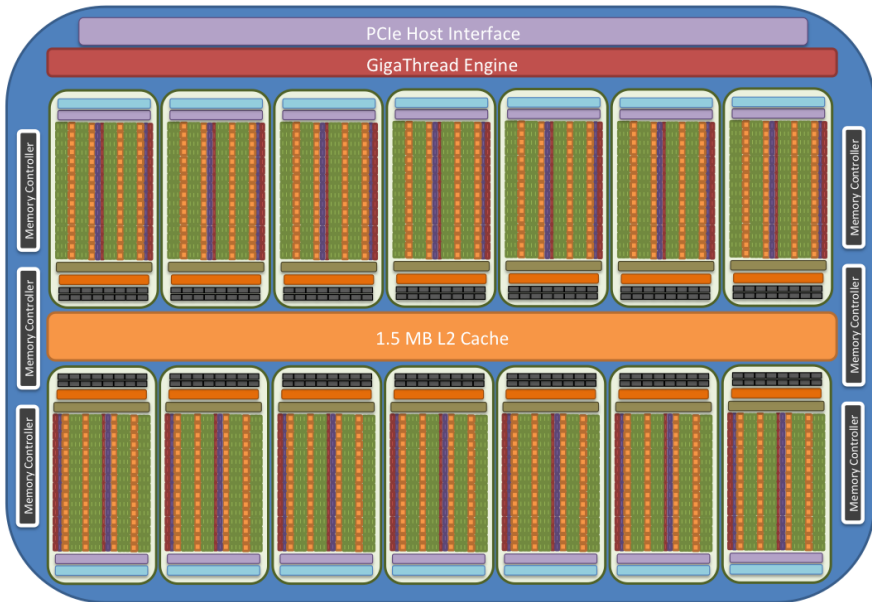
- Most hardware has useful (fixed) levels of parallelism.
- An ideal flattening algorithm maps levels of application parallelism (any number) to hardware parallelism (fixed number) in a way that exploits locality well.

**Example of deep nesting:** a system consists of multiple *datacenters*, that each contain multiple *computers*, that each contain multiple *GPUs*, that each contain multiple *SMs* (next slide), that each run some number of threads.
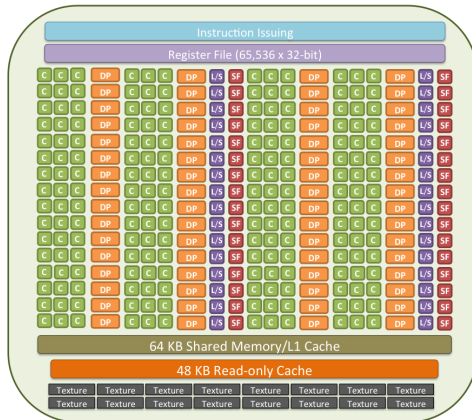
### General principle

"Tasks" at the same hardware level cannot communicate, but can "launch" tasks at a lower level.

# K20 GPU layout

# Streaming Multiprocessor (SM) layout



- C — single precision/integer CUDA core
- DP — double precision FP unit
- L/S — memory load/store unit
- SF — special function unit

## Level-aware segmented operations

$$l \in \{\text{thread}, \text{group}\}$$

- *Group* is the same as a CUDA *thread block*
- Each segmented operation then tagged with the level at which its *body* executes.

$$\mathbf{segmap}_l \ \Sigma \ e$$
$$\mathbf{segscan}_l \ \Sigma \ \odot \ \overline{d} \ e$$
$$\mathbf{segred}_l \ \Sigma \ \odot \ \overline{d} \ e$$

### Restrictions

Both thread and group can occur at top level, but a group construct can contain only thread constructs, and thread cannot contain any segmented constructs.

## Examples

**Each thread transposes part of an array**

$$\textbf{segmap}_{\text{thread}} \; \langle x \in xs \rangle \; (\texttt{transpose x})$$

**Each workgroup transposes part of an array**

$$\textbf{segmap}_{\text{group}} \; \langle x \in xs \rangle \; (\texttt{transpose x})$$

These are both equivalent to map transpose xs.

**Each workgroup sums the row of an array**

$$\textbf{segmap}_{\text{group}} \; \langle xs \in xss \rangle \; (\textbf{segred}_{\text{thread}} \; \langle x \in xs \rangle \; (+) \: 0 \: x)$$

Equivalent to map (reduce (+) 0) xss.

**Tags carry no semantic meaning; used solely for code generation.**

## Example: LocVolCalib

The following is the essential core of the LocVolCalib benchmark from the FinPar suite.

```
map (\xss ->
     map (\xs ->
          let bs = scan ⊕ d_⊕ xs
          let cs = scan ⊗ d_⊗ bs
          in  scan ⊙ d_⊙ cs)
        xss)
   xsss
```

How can we map the application parallelism to hardware parallelism?

## Option I: sequentialise the inner `scans`

$$\textbf{segmap}_{\text{thread}}\ (\langle \text{xss} \in \text{xsss}\rangle, \langle \text{xs} \in \text{xss}\rangle)$$

**let** bs = **scan** $\oplus$ $d_\oplus$ xs
**let** cs = **scan** $\otimes$ $d_\otimes$ bs
**in scan** $\odot$ $d_\odot$ cs

**scan** is relatively expensive in parallel, so this is a good option if the outer dimensions provide enough parallelism.

## Option II: flatten and parallelise inner `scans`

Flattening uses *loop distribution* (or *fission*) to create **map** nests:

```
map (\xss ->
      map (\xs ->
            let bs = scan ⊕ d⊕ xs
            let cs = scan ⊗ d⊗ bs
            in  scan ⊙ d⊙ cs)
          xss)
    xsss
```

## Option II: flatten and parallelise inner `scans`

```
let bsss =
    segscan_thread (⟨xss ∈ xsss⟩, ⟨xs ∈ xss⟩, , ⟨x ∈ xs⟩) ⊕ d_⊕ x
let csss =
    segscan_thread (⟨bss ∈ bsss⟩, ⟨bs ∈ bss⟩, , ⟨b ∈ bs⟩) ⊕ d_⊕ b
in
    segscan_thread (⟨css ∈ csss⟩, ⟨cs ∈ css⟩, , ⟨c ∈ cs⟩) ⊕ d_⊕ c
```

**This is what full flattening will do.**

```
map (\ xss ->
     map (\ xs ->
           let bs = scan ⊕ d⊕ xs
           let cs = scan ⊗ d⊗ bs
           in  scan ⊙ d⊙ cs )
         xss )
    xsss
```

## Option III: Mapping innermost parallelism to the workgroup level

$$\textbf{segmap}_{\text{group}} \ (\langle \text{xss} \in \text{xsss} \rangle, \langle \text{xs} \in \text{xss} \rangle)$$
$$\textbf{let } \text{bs} = \textbf{segscan}_{\text{thread}} \ \langle \text{x} \in \text{xs} \rangle \ \oplus \ d_\oplus \ \text{x}$$
$$\textbf{let } \text{cs} = \textbf{segscan}_{\text{thread}} \ \langle \text{b} \in \text{bs} \rangle \otimes \ d_\otimes \ \text{b}$$
$$\textbf{in } \textbf{segscan}_{\text{thread}} \ \langle \text{c} \in \text{cs} \rangle \otimes \ d_\otimes \ \text{c}$$

- Iterations of outer **segmap**s assigned to GPU workgroups[5].
- Each **segscan**$_{\text{thread}}$ is executed collaboratively by a workgroup and in local memory[6].
- Only works if the innermost parallelism fits in a workgroup.

---

[5] *Thread block* in CUDA
[6] *Shared memory* in CUDA

# LocVolCalib speedup (higher is better)



**NVIDIA K40**

small (baseline: 238ms)
medium (baseline: 342ms)
large (baseline: 5087ms)

**AMD Vega 64**

small (baseline: 110ms)
medium (baseline: 132ms)
large (baseline: 1862ms)

Legend:
- MF
- IF
- AIF
- FinPar (outer parallelism)
- FinPar (all parallelism)

Sequential scans (MF) is the baseline.

## Level-aware incremental flattening

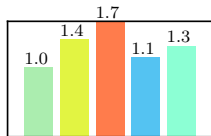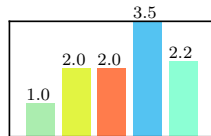$\boxed{\mathcal{G}_l(\Sigma, e) \Rightarrow e'}$ In a map-nest context $\Sigma$, the source expression $e$ can be translated at machine level $l$ into the target expression $e'$.

$$\mathcal{G}_{l+1}(\Sigma, \textbf{map } (\lambda \overline{x} \to e) \, \overline{xs}) \Rightarrow \quad \textbf{if } \mathrm{Par}(\Sigma') \geq t_{\mathrm{top}}$$
$$\textbf{then segmap}_{l+1} \, \Sigma' \, e$$
$$\textbf{else if } \mathrm{Par}(e_{\mathrm{middle}}) \geq t_{\mathrm{intra}}$$
$$\textbf{then segmap}_{l+1} \, \Sigma' \, \mathcal{G}_l(\bullet, e)$$
$$\textbf{else } \mathcal{G}_{l+1}(\Sigma', e)$$

where $\Sigma' = \Sigma, \langle \overline{x} \in \overline{xs} \rangle$ and $t_{\mathrm{top}}, t_{\mathrm{intra}}$ fresh.
In the Futhark compiler, only two levels are handled (thread, group), but we believe the idea generalises well.

# Block tiling

Level-aware constructs can also be used for expressing other powerful optimisations.

Level-aware constructs can also be used for expressing other powerful optimisations.



Threads accessing same memory can cooperatively cache it in on-chip memory.

```
map (\x -> redomap (+) (\y -> y + x) 0 xs) xs
```

After flattening we get this inner-sequential version:

$$\textbf{segmap}_{\text{thread}} \langle x \in xs \rangle \; (\textbf{redomap} \; (+) \; (\lambda y \to y + x) \; 0 \; xs)$$

Operation  One thread for each element of $xs$, and each sequentially traverses $xs$.

Problem  ?

```
map (\x -> redomap (+) (\y -> y + x) 0 xs) xs
```

After flattening we get this inner-sequential version:

$$\textbf{segmap}_{\text{thread}} \ \langle x \in xs \rangle \ (\textbf{redomap} \ (+) \ (\lambda y \rightarrow y + x) \ 0 \ xs)$$

Operation  One thread for each element of xs, and each sequentially traverses xs.

Problem  Poor utilisation of memory bus.

- Many threads simultaneously read same address, which is redundant.
- **Better:** *cooperatively* copy *block* into on-chip memory and iterate from there.

## Strip mining/chunking the outer `segmap`

$$\mathbf{segmap}_{\text{thread}} \; \langle x \in xs \rangle \; (\mathbf{redomap} \; (+) \; (\lambda y \to y + x) \; 0 \; xs)$$

Assuming we can split $xs$ into $m$ equally sized *tiles* each of size $t$, giving $xss$ : $[m][t]f32$, then we can rewrite to

$$\begin{aligned} &\mathbf{segmap}_{\text{group}} \; \langle xs' \in xss \rangle \\ &\quad \mathbf{segmap}_{\text{thread}} \; \langle x \in xs' \rangle \\ &\quad\quad \mathbf{redomap} \; (+) \; (\lambda y \to y + x) \; 0 \; xs \end{aligned}$$

**Question: does this compute the same value as the original?**

## Strip mining/chunking the outer `segmap`

$$\textbf{segmap}_{\text{thread}} \ \langle x \in xs \rangle \ (\textbf{redomap} \ (+) \ (\lambda y \rightarrow y + x) \ 0 \ xs)$$

Assuming we can split $xs$ into $m$ equally sized *tiles* each of size $t$, giving $xss$ : $[m][t]f32$, then we can rewrite to

$$\textbf{segmap}_{\text{group}} \ \langle xs' \in xss \rangle$$
$$\textbf{segmap}_{\text{thread}} \ \langle x \in xs' \rangle$$
$$\textbf{redomap} \ (+) \ (\lambda y \rightarrow y + x) \ 0 \ xs$$

**Question: does this compute the same value as the original?**

- *No*—the original expression had type $[n]f32$, while this has type $[m][t]f32$
- This can be flattened away.

$$\textbf{segmap}_{\text{group}} \; \langle \mathsf{xs'} \in \mathsf{xss} \rangle$$
$$\textbf{segmap}_{\text{thread}} \; \langle \mathsf{x} \in \mathsf{xs'} \rangle$$
$$\textbf{redomap} \; (+) \; (\lambda \mathsf{y} \to \mathsf{y} + \mathsf{x}) \; 0 \; \mathsf{xs}$$

Chunking/strip-mining the **redomap**, we get

$$\textbf{segmap}_{\text{group}} \; \langle \mathsf{xs'} \in \mathsf{xss} \rangle$$
$$\textbf{segmap}_{\text{thread}} \; \langle \mathsf{x} \in \mathsf{xs'} \rangle$$
$$\textbf{loop} \; \mathsf{acc} \; = \; 0 \; \textbf{for} \; \mathsf{ys} \; \textbf{in} \; \mathsf{xss} \; \textbf{do}$$
$$\textbf{redomap} \; (+) \; (\lambda \mathsf{y} \to \mathsf{y} + \mathsf{x}) \; \mathsf{acc} \; \mathsf{ys}$$

$$\textbf{segmap}_{\text{group}} \; \langle \mathsf{xs'} \in \mathsf{xss} \rangle$$
$$\quad \textbf{segmap}_{\text{thread}} \; \langle \mathsf{x} \in \mathsf{xs'} \rangle$$
$$\quad\quad \textbf{redomap} \; (+) \, (\lambda \mathsf{y} \to \mathsf{y} + \mathsf{x}) \; 0 \; \mathsf{xs}$$

Chunking/strip-mining the **redomap**, we get

$$\textbf{segmap}_{\text{group}} \; \langle \mathsf{xs'} \in \mathsf{xss} \rangle$$
$$\quad \textbf{segmap}_{\text{thread}} \; \langle \mathsf{x} \in \mathsf{xs'} \rangle$$
$$\quad\quad \textbf{loop} \; \mathsf{acc} \, = \, 0 \; \textbf{for} \; \mathsf{ys} \; \textbf{in} \; \mathsf{xss} \; \textbf{do}$$
$$\quad\quad\quad \textbf{redomap} \; (+) \, (\lambda \mathsf{y} \to \mathsf{y} + \mathsf{x}) \; \mathsf{acc} \; \mathsf{ys}$$

Distributing and interchanging **segmap**$_{\text{thread}}$ gives

$$\textbf{segmap}_{\text{group}} \; \langle \mathsf{xs'} \in \mathsf{xss} \rangle$$
$$\quad \textbf{loop} \; \mathsf{accs} \, = \, \mathsf{replicate} \; \mathsf{t} \; 0$$
$$\quad \textbf{for} \; \mathsf{ys} \; \textbf{in} \; \mathsf{xss} \; \textbf{do}$$
$$\quad\quad \textbf{segmap}_{\text{thread}} \; \langle \mathsf{x}, \mathsf{acc} \in \mathsf{xs'}, \mathsf{accs} \rangle$$
$$\quad\quad\quad \textbf{redomap} \; (+) \, (\lambda \mathsf{y} \to \mathsf{y} + \mathsf{x}) \; \mathsf{acc} \; \mathsf{ys}$$

$$\textbf{segmap}_{\text{group}} \; \langle \text{xs'} \in \text{xss} \rangle$$
$$\qquad \textbf{loop} \; \text{accs} = \text{replicate t 0}$$
$$\qquad \textbf{for} \; \text{ys} \; \textbf{in} \; \text{xss} \; \textbf{do}$$
$$\qquad\qquad \textbf{segmap}_{\text{thread}} \; \langle x, \text{acc} \in \text{xs'}, \text{accs} \rangle$$
$$\qquad\qquad\qquad \textbf{redomap} \; (+) \; (\lambda y \rightarrow y + x) \; \text{acc} \; \text{ys}$$

Collectively copy ys to shared/local memory

$$\textbf{segmap}_{\text{group}} \; \langle \text{xs'} \in \text{xss} \rangle$$
$$\qquad \textbf{loop} \; \text{accs} = \text{replicate t 0}$$
$$\qquad \textbf{for} \; \text{ys} \; \textbf{in} \; \text{xss} \; \textbf{do}$$
$$\qquad\qquad \textbf{let} \; \text{ys'} = \textbf{copy} \; \text{ys} \; \textbf{in}$$
$$\qquad\qquad\qquad \textbf{segmap}_{\text{thread}} \; \langle x, \text{acc} \in \text{xs'}, \text{accs} \rangle$$
$$\qquad\qquad\qquad\qquad \textbf{redomap} \; (+) \; (\lambda y \rightarrow y + x) \; \text{acc} \; \text{ys'}$$

- Now the many iterations of the **redomap** read from fast on-chip memory rather than slower global memory!
- **copy** done collectively by all threads in group

## The fine print

```
map (\x -> redomap (+) (\y -> y + x) 0 xs) xs
```

to

$$\mathbf{segmap}_{\mathrm{group}} \langle \mathrm{xs'} \in \mathrm{xss} \rangle$$
$$\quad \mathbf{loop} \; \mathrm{accs} = \mathtt{replicate} \; \mathrm{t} \; 0$$
$$\quad \mathbf{for} \; \mathrm{ys} \; \mathbf{in} \; \mathrm{xss} \; \mathbf{do}$$
$$\quad\quad \mathbf{let} \; \mathrm{ys'} = \mathbf{copy} \; \mathrm{ys} \; \mathbf{in}$$
$$\quad\quad\quad \mathbf{segmap}_{\mathrm{thread}} \langle \mathrm{x}, \mathrm{acc} \in \mathrm{xs'}, \mathrm{accs} \rangle$$
$$\quad\quad\quad\quad \mathbf{redomap} \; (+) \; (\lambda \mathrm{y} \rightarrow \mathrm{y} + \mathrm{x}) \; \mathrm{acc} \; \mathrm{ys'}$$

- Very simple case (e.g. xss traversed in both loops)
- 2D tiling much more complex
- The *tile size* t is a sensitive tuning parameter; in this case it should coincide with workgroup size
- Appreciate what a compiler can do for you

## Summary

- There is no *one size fits all:* for optimal performance, we need different amounts of parallelisation for different workloads.
- Incremental flattening generates a *single program* that for varying datasets exploits only as much parallelism as profitable.
- Autotuning for specific hardware and program is needed to select the optimal version at runtime.
- A good IR is as crucial to a compiler as a good language is to a human.