



# Compiling a Subset of APL into Performance Efficient GPU Programs

Martin Elsmann, DIKU, University of Copenhagen

*Joined work with Troels Henriksen, Martin Dybdal,  
Henrik Urms, Anna Sofie Kiehn, and Cosmin Oancea*

**@Dyalog'16**

# Motivation

## Goal:

High-performance at the fingertips of domain experts.

## Why APL:

APL provides a **powerful and concise** notation for array operations.

APL programs are inherently parallel - not just parallel, but **data-parallel**.

There is lots of APL code around - some of which is looking to run faster!



## Challenge:

APL is dynamically typed. To generate efficient code, we need **type inference**:

- Functions are *rank-polymorphic*.
- Built-in operations are overloaded.
- Types are *value-sensitive* (e.g., any integer 0,1 is considered boolean).

Type inference algorithm compiles APL into a *typed array intermediate language* called **TAIL** (ARRAY'14).



# APL Supported Features

Dfns-syntax for functions and operators (incl. trains).

Dyalog APL compatible built-in operators and functions (limitations apply).

Scalar extensions, identity item resolution, overloading resolution.

## Limitations:

- Static scoping and static rank inference
- Limited support for nested arrays
- Whole-program compilation
- No execute!



`else ← { (αα*α) (ωω*(~α) )ω }`

`mean ← +/÷≠`



# TAIL - as an IL



- Type system **expressive enough** for many APL primitives.
- Simplify certain primitives into other constructs...
- Multiple backends...

APL	op(s)	TySc(op)
	addi, ...	int → int → int
	add, ...	double → double → double
~	iota	int → [int] <sup>1</sup>
∇	each	∀αβγ. (α → β) → [α] <sup>γ</sup> → [β] <sup>γ</sup>
/	reduce	∀αγ. (α → α) → α → [α] <sup>γ</sup>
ρ	shape	∀αγ. [α] <sup>γ</sup> → Sh γ
ρ	reshape0	∀αγγ'. Sh γ' → [α] <sup>γ</sup> → [α] <sup>γ'</sup>
ρ	reshape	∀αγγ'. Sh γ' → α → [α] <sup>γ</sup> → [α] <sup>γ'</sup>
φ	reverse	∀αγ. [α] <sup>γ</sup> → [α] <sup>γ</sup>
φ	rotate	∀αγ. int → [α] <sup>γ</sup> → [α] <sup>γ</sup>
\	transp	∀αγ. [α] <sup>γ</sup> → [α] <sup>γ</sup>
\	transp2	∀αγ. Sh γ → [α] <sup>γ</sup> → [α] <sup>γ</sup>
†	take	∀αγ. int → α → [α] <sup>γ</sup> → [α] <sup>γ</sup>
↓	drop	∀αγ. int → [α] <sup>γ</sup> → [α] <sup>γ</sup>
⊔	first	∀αγ. α → [α] <sup>γ</sup> → α
⊔	zipWith	∀α <sub>1</sub> α <sub>2</sub> βγ. (α <sub>1</sub> → α <sub>2</sub> → β) → [α <sub>1</sub> ] <sup>γ</sup> → [α <sub>2</sub> ] <sup>γ</sup> → [β] <sup>γ</sup>
,	cat	∀αγ. [α] <sup>γ+1</sup> → [α] <sup>γ+1</sup> → [α] <sup>γ+1</sup>
,	cons	∀αγ. [α] <sup>γ</sup> → [α] <sup>γ+1</sup> → [α] <sup>γ+1</sup>
,	snoc	∀αγ. [α] <sup>γ+1</sup> → [α] <sup>γ</sup> → [α] <sup>γ+1</sup>

**Value typing** Γ ⊢ v : τ

$$\frac{\vdash \delta : \rho}{\Gamma \vdash [\delta]^\delta : [\text{int}]^\rho} \quad (7) \quad \frac{\vdash \delta : \rho}{\Gamma \vdash [\delta]^\delta : [\text{double}]^\rho} \quad (8)$$

$$\frac{}{\Gamma \vdash [\delta]^{(n)} : \text{Sh } n} \quad (9) \quad \frac{\Gamma, x : \tau \vdash e : \tau'}{\Gamma \vdash \lambda x. e : \tau \rightarrow \tau'} \quad (10)$$

$$\frac{}{\Gamma \vdash [n]^{(1)} : \text{VI } n} \quad (11) \quad \frac{}{\Gamma \vdash n : \text{I } n} \quad (12)$$

**Expression typing** Γ ⊢ e : τ

$$\frac{\Gamma \vdash e : \text{I } n}{\Gamma \vdash [e] : \text{VI } n} \quad (13) \quad \frac{\Gamma(x) \geq \tau}{\Gamma \vdash x : \tau} \quad (14) \quad \frac{\tau \subseteq \tau'}{\Gamma \vdash e : \tau} \quad (15)$$

$$\frac{\Gamma \vdash e_i : \kappa \quad i = [0; n]}{\Gamma \vdash [e^{(n)}] : [\kappa]^1} \quad (16) \quad \frac{\Gamma \vdash e_i : \text{int} \quad i = [0; n]}{\Gamma \vdash [e^{(n)}] : \text{Sh } n} \quad (17)$$

$$\frac{\Gamma \vdash e_1 : \tau' \rightarrow \tau \quad \Gamma \vdash e_2 : \tau'}{\Gamma \vdash e_1 e_2 : \tau} \quad (18) \quad \frac{\Gamma \vdash e_1 : \tau \quad \text{fv}(\vec{\alpha}\vec{\gamma}) \cap \text{fv}(\Gamma, \tau') = \emptyset}{\Gamma, x : \forall \vec{\alpha}\vec{\gamma}. \tau \vdash e_2 : \tau'} \quad (19)$$

$$\frac{\text{TySc}(op) \geq \tau_0 \rightarrow \dots \rightarrow \tau_{(n-1)} \rightarrow \tau \quad \text{arity}(op) = n \quad \Gamma \vdash e_i : \tau_i \quad i = [0; n]}{\Gamma \vdash op(\vec{e}^{(n)}) : \tau} \quad (20)$$

**Small Step Reductions** e ↦ e' / err

$$\frac{e \hookrightarrow e' \quad E \neq [\cdot]}{E[e] \hookrightarrow E[e']} \quad (21) \quad \frac{e \hookrightarrow \text{err} \quad E \neq [\cdot]}{E[e] \hookrightarrow \text{err}} \quad (22)$$

$$\frac{}{\text{let } x = v \text{ in } e \hookrightarrow e[v/x]} \quad (23) \quad \frac{}{(\lambda x. e) v \hookrightarrow e[v/x]} \quad (24)$$

$$\frac{}{[\vec{a}^{(n)}] \hookrightarrow [\vec{a}^{(n)}]^{(n)}} \quad (25)$$

$$\frac{i = i_1 + i_2}{\text{addi}(i_1, i_2) \hookrightarrow i} \quad (26) \quad \frac{d = d_1 + d_2}{\text{add}(d_1, d_2) \hookrightarrow d} \quad (27)$$

$$\frac{n \geq 0}{\text{iota}(n) \hookrightarrow [1, \dots, n]^{(n)}} \quad (28) \quad \frac{n < 0}{\text{iota}(n) \hookrightarrow \text{err}} \quad (29)$$

$$\frac{e = [v_f a_0, \dots, v_f a_{(n-1)}]}{\text{each}(v_f, [\vec{a}^{(n)}]^\delta) \hookrightarrow \text{reshape0}(\delta, e)} \quad (30)$$

$$\frac{\delta = \langle \vec{n}, m \rangle \quad k = \text{product}(\vec{n}) \quad i = [0; m[ \quad e_i = v_f a_{(i+k)} (\dots (v_f a_{(i+k+m-1)}) v \dots )]}{\text{reduce}(v_f, v, [\vec{a}^{(n)}]^\delta) \hookrightarrow \text{reshape0}(\langle \vec{n}, [e^{(k)}] \rangle)} \quad (31)$$

$$\frac{m = \text{product}(\delta') \quad f(i) = i \bmod n \quad n > 0}{\text{reshape}(\delta', a, [\vec{a}^{(n)}]^\delta) \hookrightarrow [a_{f(0)}, \dots, a_{f(m-1)}]^\delta} \quad (32)$$

$$\frac{m = \text{product}(\delta') \quad a_i = a \quad i = [0; m[ \quad \text{reshape}(\delta', a, []^\delta) \hookrightarrow [a_0, \dots, a_{(m-1)}]^\delta} \quad (33)$$

$$\frac{\delta' = \text{rev}(\delta) \quad f = \text{fromSh}_{\delta'} \circ \text{rev} \circ \text{toSh}_{\delta}}{\text{transp}([\vec{a}^{(n)}]^\delta) \hookrightarrow [a_{f(0)}, \dots, a_{f(n-1)}]^\delta} \quad (34)$$


# TAIL Example



## APL:

```
mean ← +/÷≠
var ← mean({ω*2}←-mean)
stddev ← {ω*0.5} var
all ← mean, var, stddev
□ ← all 54 44 47 53 51 48 52 53 52 49 48
```

Type check: Ok

Evaluation:

```
[3](50.0909,8.8099,2.9681)
```

Simple interpreter

## TAIL:

```
let v2:[int]1 = [54,44,47,53,51,48,52,53,52,49,48,52] in
let v1:[int]0 = 11 in
let v15:[double]1 = each(fn v14:[double]0 =>
  subd(v14,divd(i2d(reduce(addi,0,v2)),i2d(v1))),each(i2d,v2)) in
let v17:[double]1 = each(fn v16:[double]0 => powd(v16,2.0),v15) in
let v21:[double]0 = divd(reduce(add,0.0,v17),i2d(v1)) in
let v31:[double]1 = each(fn v30:[double]0 =>
  subd(v30,divd(i2d(reduce(addi,0,v2)),i2d(v1))),each(i2d,v2)) in
let v33:[double]1 = each(fn v32:[double]0 => powd(v32,2.0),v31) in
let v41:[double]1 =
  prArrD(cons(divd(i2d(reduce(addi,0,v2)),i2d(v1)),[divd(reduce(add
d,0.0,v33),i2d(v1)),powd(v21,0.5)])) in 0
```



# Compiling Primitives



APL: Guibas and Wyatt, POPL'78

```
dot ← {  
  WA ← (1↓ρω), ρα  
  KA ← (⊃ρρ) - 1  
  VA ← ι ⊃ ρWA  
  ZA ← (KAϕ-1↓VA), -1↑VA  
  TA ← ZAWAρα  
  WB ← (-1↓ρ)α, ρω  
  KB ← ⊃ ρρ  
  VB ← ι ⊃ ρWB  
  ZB0 ← (-KB) ↓ KB ϕ ι(⊃ρVB)  
  ZB ← (-1↓(ι KB)), ZB0, KB  
  TB ← ZBWBρω  
  αα / TA ωω TB  
}
```

```
A ← 3 2 ρ ι 5  
B ← ϕ A  
R ← A + dot × B  
R2 ← ×/ +/ R
```

TAIL:

```
let v1:[int]2 = reshape([3,2],iotaV(5)) in  
let v2:[int]2 = transp(v1) in  
let v9:[int]3 = transp2([2,1,3],reshape([3,3,2],v1)) in  
let v15:[int]3 = transp2([1,3,2],reshape([3,2,3],v2)) in  
let v20:[int]2 = reduce(addi,0,zipWith(muli,v9,v15)) in  
let v25:[int]0 = reduce(muli,1,reduce(addi,0,v20)) in  
i2d(v25)
```

Evaluating  
Result is `[](65780.0)`

**Notice:** Quite a few simplifications happen at TAIL level..



# Futhark

Pure eager **functional language** with second-order parallel array constructs.

Support for “imperative-like” language constructs for iterative computations (i.e., graph shortest path).

A **sequentialising** compiler...

Close to performance obtained with hand-written OpenCL GPU code.

```
fun [int] addTwo ([int] a) = map(+2, a)
fun int sum ([int] a) = reduce(+, 0, a)
fun [int] sumrows ([[int]] as) = map(sum, a)
fun int main(int n) =
  loop (x=1) = for i<n do x*(i+1)
in x
```



## Performs general optimisations

- *Constant folding*. E.g., remove branch inside code for `take(n,a)` if  $n \leq \lceil \rho a \rceil$ .
- *Loop fusion*. E.g., fuse the many small “vectorised” loops in idiomatic APL code.

## Attempts at flattening nested parallelism

- E.g., reduction (`/`) inside each (`⌈`).

## Allows for indexing and sequential loops

- Needed for indirect indexing and `*`.

## Performs low-level GPU optimisations

- E.g., optimise for coalesced memory accesses.



# An Example



## APL:

```
f ← { 2 ÷ ω + 2 }      ⍝ Function \x. 2 / (x+2)
X ← 1000000           ⍝ Valuation points per unit
domain ← 10 × (ιX) ÷ X ⍝ Integrate from 0 to 10
integral ← +/ (f¨domain)÷X ⍝ Compute integral
```



## TAIL:

```
let domain:<double>1000000 =
  eachV(fn v4:[double]0 => muld(10.0,v4),
    eachV(fn v3:[double]0 => divd(v3,1000000.0),
      eachV(i2d,iotaV(1000000)))) in
let integral:[double]0 =
  reduce(addy,0.0,
    eachV(fn v9:[double]0 => divd(v9,1000000.0),
      eachV(fn v7:[double]0 => divd(2.0,addy(v7,2.0)),
        domain))) in
integral
```



## Futhark - before optimisation:

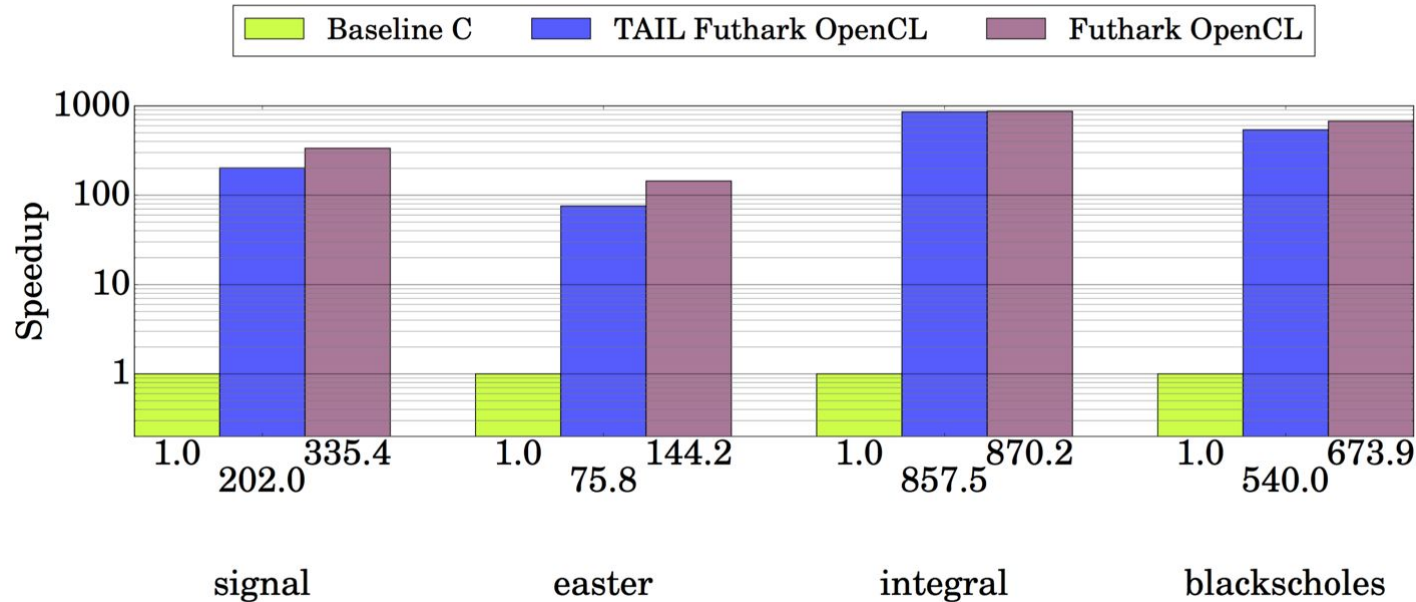
```
let domain =
  map (fn (t_v4: f64): f64 => 10.0f64*t_v4)
    (map (fn (t_v3: f64): f64 => t_v3/1000000.0f64)
      (map i2d (map (fn (x: int): int => x+1)
        (iota 1000000)))) in
let integral =
  reduce (+) 0.0f64
    (map (fn (t_v9: f64): f64 => t_v9/1000000.0f64)
      (map (fn (t_v7: f64): f64 => 2.0f64/(t_v7+2.0f64))
        domain)) in
integral
```

**Notice:** TAIL2Futhark compiler  
is quite straightforward...





# Performance Compute-bound Examples



## Integral benchmark:

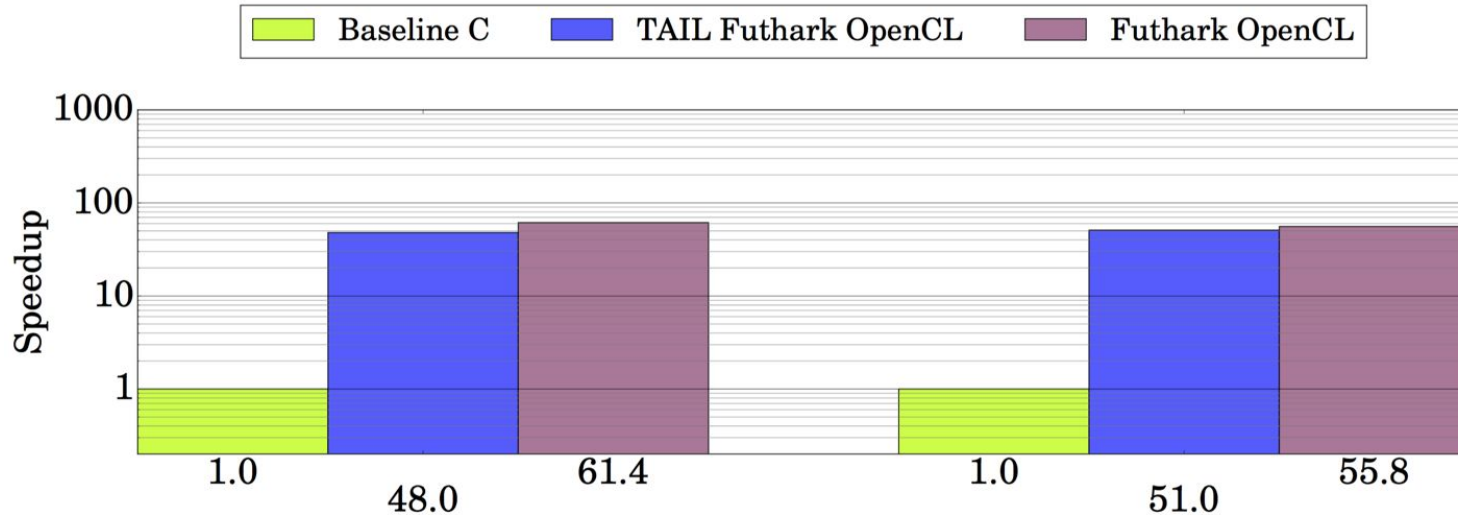
```
f ← { 2 ÷ ω + 2 }
X ← 10000000
domain ← 10 × (ιX) ÷ X
integral ← +/ (f¨domain)÷X
```

```
⌘ Function \x. 2 / (x+2)
⌘ Valuation points per unit
⌘ Integrate from 0 to 10
⌘ Compute integral
```

*OpenCL runtimes from an NVIDIA GTX 780  
CPU runtimes from a Xeon E5-2650 @ 2.6GHz*



# Performance Stencils



Life benchmark:

life

hotspot

```
life ← {  
  rs ← { (-1φω) + ω + 1φω }  
  n ← (rs -1eω) + (rs ω) + rs 1eω  
  (n=3) v (n=4) ∧ ω  
}  
res ← +/+/ (life * 100) board
```



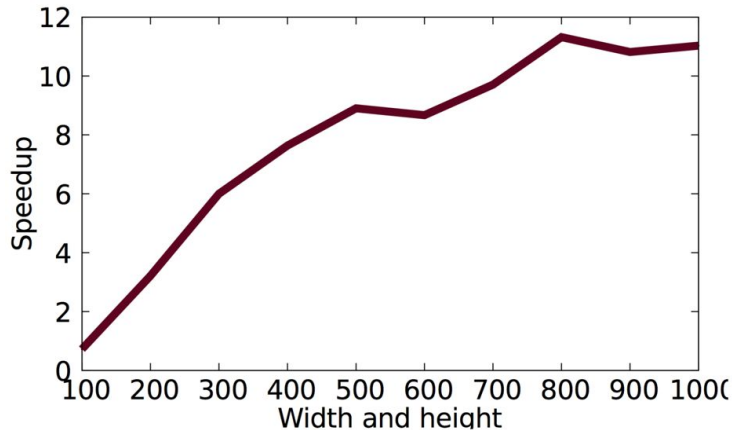
# Different Mandelbrot Implementations

## Parallel inner loop:

*mandelbrot1.apl*

```
seq for i < depth:  
  par for j < n:  
    points[j] = f(points[j])
```

**Memory bound**

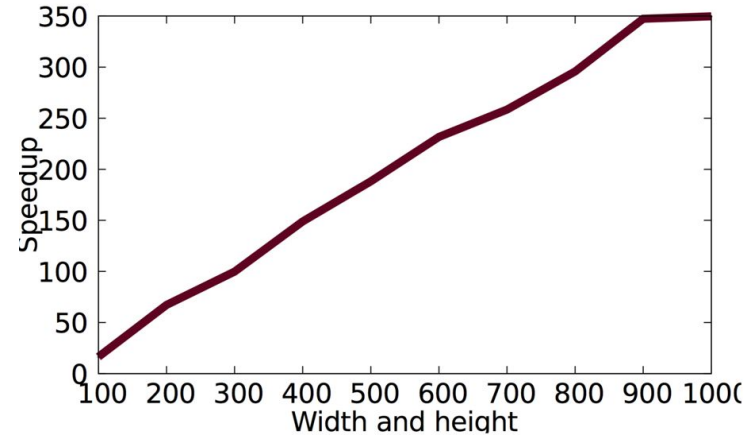


## Parallel outer loop:

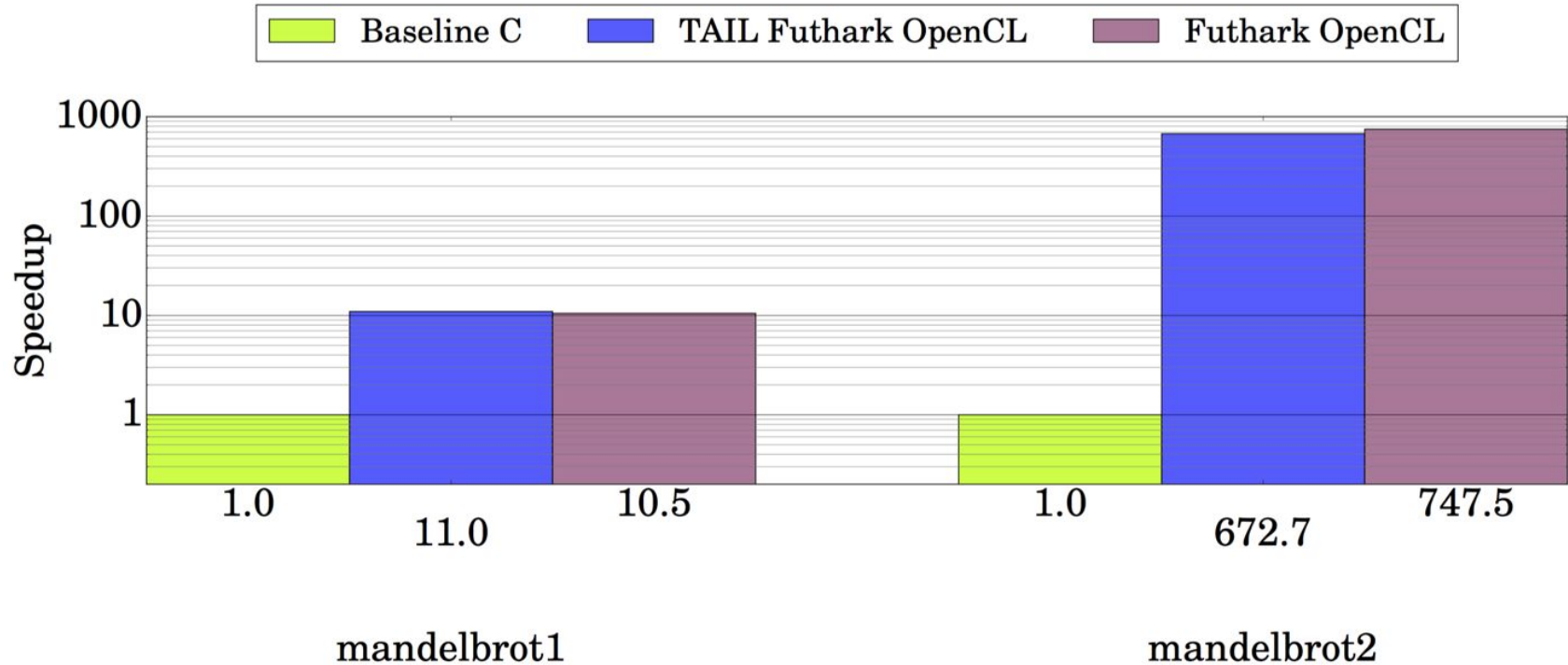
*mandelbrot2.apl*

```
par for j < n:  
  p = points[j]  
  seq for i < depth:  
    p = f(p)  
  points[j] = p
```

**Compute bound**



# Performance Mandelbrot

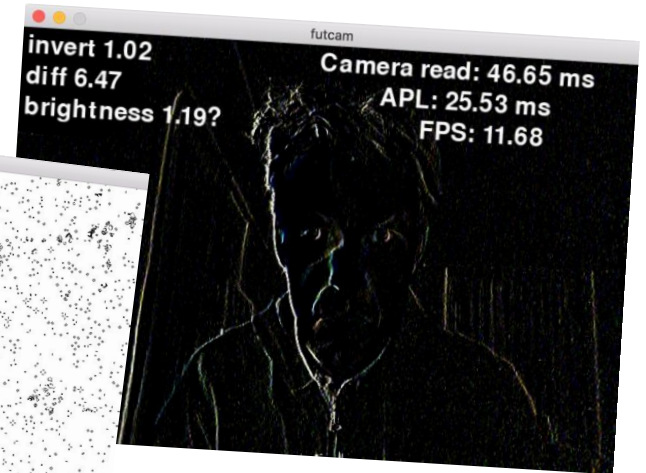
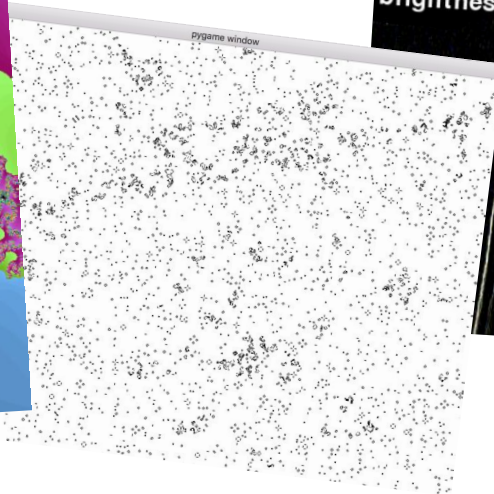
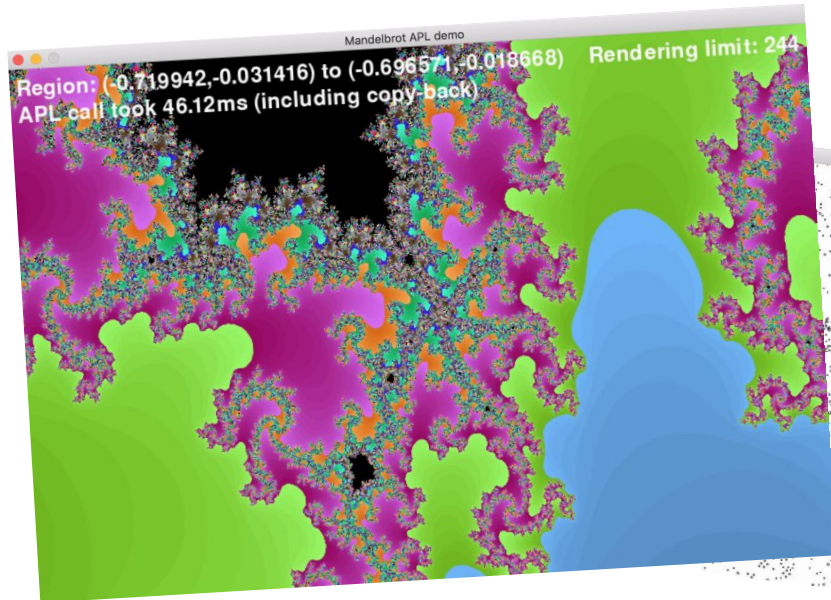


# Interoperability Demos

*Mandelbrot, Life, APlCam*

With Futhark, we can generate reusable *modules* in various languages (e.g, Python) that internally execute on the GPU using OpenCL.

```
onChannels ← {  
  m ← 3 1 2 0 ω  
  m ← (αα h w ρm) ; (αα h w ρ1+m) ; αα h w ρ2+m  
  2 3 1 0 3 h w ρ m  
}  
diff ← {  
  n ← [degree  
  255|n×+ω-1φω  
}  
image ← diff onChannels image
```



# Related Work

## APL Compilers

- Co-dfns compiler by Aaron Hsu. Papers in ARRAY'14 and ARRAY'16.
- C. Grelck and S.B. Scholz. *Accelerating APL programs with SAC*. APL'99.
- R. Bernecky. *APEX: The APL parallel executor*. MSc Thesis. University of Toronto. 1997.
- L.J. Guibas and D.K. Wyatt. *Compilation and delayed evaluation in APL*. POPL'78.

## Type Systems for APL like Languages

- K. Trojahner and C. Grelck. *Dependently typed array programs don't go wrong*. NWPT'07.
- J. Slepak, O. Shivers, and P. Manolios. *An array-oriented language with static rank polymorphism*. ESOP'14.

## Futhark work

- Papers on language and optimisations available from [hiperfit.dk](http://hiperfit.dk).
- Futhark available from [futhark-lang.org](http://futhark-lang.org).

## Other functional languages for GPUs

- Accelerate. Haskell library/embedded DSL.
- Obsidian. Haskell embedded DSL.
- FCL. Low-level functional GPU programming. FHPC'16.

## Libraries for GPU Execution

- Thrust, cuBLAS, cuSPARSE, ...



# Conclusions

- We have managed to get a (small) subset of APL to run efficiently on GPUs.
  - <https://github.com/HIPERFIT/futhark-fhpc16>.
  - <https://github.com/henrikurms/tail2futhark>.
  - <https://github.com/melsman/apltail>.

# Future Work

- More real-world benchmarks.
- Support a wider subset of APL.
- Improve interoperability...
- Add support for APL “type annotations” for specifying programmer intentions...

**HIPERFIT**



# mandelbrot1.apl and mandelbrot2.apl

A grid-size in left argument (e.g., (1024 768))  
 A X-range, Y-range in right argument

```
mandelbrot1 ← {
  X ← ⍳α ⋄ Y ← ⍳1↓α
  xRng ← 2↑ω ⋄ yRng ← 2↓ω
  dx ← ((xRng[2]) - xRng[1]) ÷ X
  dy ← ((yRng[2]) - yRng[1]) ÷ Y
  cx ← Y X ρ (xRng[1]) + dx × 1X      A real plane
  cy ← ⍉ X Y ρ (yRng[1]) + dy × 1Y    A img plane
  mandel1 ← {
    zx ← Y X ρ ω[1] ⋄ zy ← Y X ρ ω[2]
    count ← Y X ρ ω[3]                A count plane
    zzx ← cx + (zx × zx) - zy × zy
    zzy ← cy + (zx × zy) + zx × zy
    conv ← 4 > (zzx × zzx) + zzy × zzy
    count2 ← count + 1 - conv
    (zzx zzy count2)
  }
  pl ← Y X ρ 0                        A zero-plane
  N ← 255                              A iterations
  res ← (mandel1 * N) (pl pl pl)
  res[3] ÷ N                            A count plane
}
```

```
mandelbrot2 ← {
  X ← ⍳α ⋄ Y ← ⍳1↓α
  xRng ← 2↑ω ⋄ yRng ← 2↓ω
  dx ← ((xRng[2]) - xRng[1]) ÷ X
  dy ← ((yRng[2]) - yRng[1]) ÷ Y
  cxA ← Y X ρ (xRng[1]) + dx × 1X    A real plane
  cyA ← ⍉ X Y ρ (yRng[1]) + dy × 1Y  A img plane
  N ← 255                              A iterations
  mandel1 ← {
    cx ← α ⋄ cy ← ω
    f ← {
      arg ← ω
      x ← arg[1] ⋄ y ← arg[2]
      count ← arg[3]
      dummy ← arg[4]
      zx ← cx + (x × x) - (y × y)
      zy ← cy + (x × y) + (x × y)
      conv ← 4 > (zx × zx) + zy × zy
      count2 ← count + 1 - conv
      (zx zy count2 dummy)
    }
    res ← (f * N) (0 0 0 'dummy')     A N iterations
    res[3]
  }
  res ← cxA mandel1" cyA
  res ÷ N
}
```

