

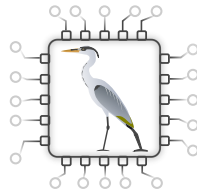
Heron: Modern Hardware Graph Reduction

HAFLANG Project

Craig Ramsay & Rob Stewart

August 2023

Heriot-Watt University

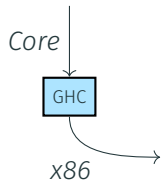


“ I wonder how popular Haskell needs to become for Intel to optimize their processors for my runtime, rather than the other way around. ”

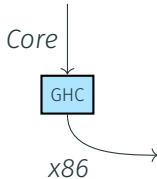
— Simon Marlow, 2009

The Graph Reduction Problem

```
add x y z =  
  x + y + z
```



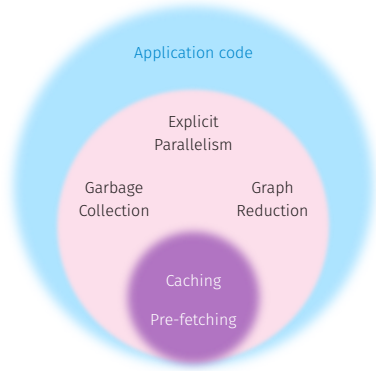
```
add x y z =
  x + y + z
```



```
s1ba_info: ; Code for helper (\a b -> a+b)
.Lc1bo: ; Check for stack space
    leaq -40(%rbp),%rax
    cmpq %r15,%rax
    jb .Lc1bp ; Jump if stack full
.Lc1bq: ; Reduce helper
    movq $stg_upd_frame_info,-16(%rbp)
    movq %rbx,-8(%rbp)
    movq 16(%rbx),%rax ; Load a & b from heap
    movq 24(%rbx),%rbx
    movl $base_GHCziNum_zdfNumInt_closure,%r14d
    ;; Push `a+b` onto stack
    movq $stg_ap_pp_info,-40(%rbp)
    movq %rax,-32(%rbp)
    movq %rbx,-24(%rbp)
    addq $-40,%rbp
    jmp base_GHCziNum_zp_info ; Enter
.Lc1bp: ; Ask RTS for stack space
    jmp *-16(%r13)
```

```
Add_add_info: ; Code for `add`
.Lc1br: ; Check for stack space
    leaq -24(%rbp),%rax
    cmpq %r15,%rax
    jb .Lc1bs ; Jump if stack full
.Lc1bt: ; Check for heap space
    addq $32,%r12
    cmpq 856(%r13),%r12
    ja .Lc1bv ; Jump if heap full
.Lc1bu: ; Reduce `add`
    ;; Build `x+y` thunk on heap
    movq $s1ba_info,-24(%r12)
    movq %r14,-8(%r12)
    movq %rsi,(%r12)
    leaq -24(%r12),%rax
    movl $base_GHCziNum_zdfNumInt_closure,%r14d
    ;; Push `thunk+z` to stack
    movq $stg_ap_pp_info,-24(%rbp)
    movq %rax,-16(%rbp)
    movq %rdi,-8(%rbp)
    addq $-24,%rbp
    jmp base_GHCziNum_zp_info ; Enter
.Lc1bv: ; Ask RTS for heap space
    movq $32,904(%r13)
.Lc1bs: ; Ask RTS for stack space
    movl $Add_add_closure,%ebx
    jmp *-8(%r13)
```

Conventional CPU



= *software*

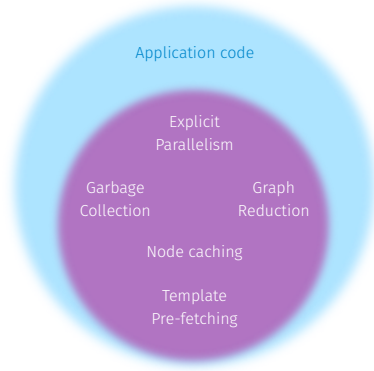


= *RTS software*



= *hardware*

Specialised Graph Reduction Machine

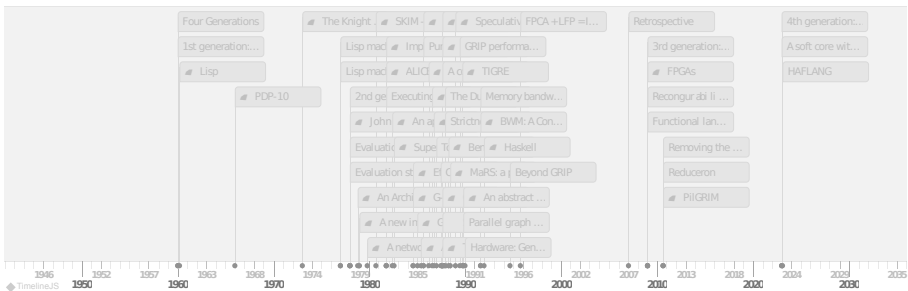


<http://haflang.github.io/history>

FUNCTIONAL HARDWARE 1924 - 2023



A 100 year history of hardware implementations of functional languages.



by [Rob Stewart](#).

Opportunities

Intra-function parallelism

Usually a victim of the von Neumann bottleneck. Worse for lazy, pure languages.

Inter-function parallelism

Exploiting the purity of functional languages.

Hardened, concurrent run-time system

Tasks like garbage collection usually halt reductions in software implementations.

Intra-function parallelism

Usually a victim of the von Neumann bottleneck. Worse for lazy, pure languages.

Inter-function parallelism

Exploiting the purity of functional languages.

Hardened, concurrent run-time system

Tasks like garbage collection usually halt reductions in software implementations.

A Template Instantiation Machine

$p ::= \bar{t}$	(Program)
$t ::=$ $\text{FUN } s \ \alpha \ n =$ $\text{let } \bar{u} \text{ in } v$	(Template)
$u, v ::=$ $\text{APP } \bar{e}$ $ \text{CASE } c \ \bar{e}$ $ \text{PRIM } n \ \bar{e}$	(Applications) (Normal application) (Application with case table) (PRS candidate allocation)
$c ::= \text{TAB } n$	(Case table pointer)
$e ::=$ $\text{CON } \alpha \ n$ $ \text{INT } n$ $ \text{PRI } \alpha \ \otimes$ $ \text{FUN } s \ \alpha \ n$ $ \text{ARG } s \ n$ $ \text{VAR } s \ n$ $ \text{REG } n$	(Atoms) (Constructor tag) (Primitive integers) (Primitive operation) (Function pointer) (Argument pointer) (Application pointer) (Primitive register pointer)

Our “*assembly*” is a non-strict functional language

$p ::= \bar{t}$ (Program)

$t ::=$ (Template)
 $\text{FUN } s \alpha n =$
 $\text{let } \bar{u} \text{ in } v$

$u, v ::=$ (Applications)
 $\text{APP } \bar{e}$ (Normal application)
 $| \text{CASE } c \bar{e}$ (Application with case table)
 $| \text{PRIM } n \bar{e}$ (PRS candidate allocation)

$c ::= \text{TAB } n$ (Case table pointer)

$e ::=$ (Atoms)
 $\text{CON } \alpha n$ (Constructor tag)
 $| \text{INT } n$ (Primitive integers)
 $| \text{PRI } \alpha \otimes$ (Primitive operation)
 $| \text{FUN } s \alpha n$ (Function pointer)
 $| \text{ARG } s n$ (Argument pointer)
 $| \text{VAR } s n$ (Application pointer)
 $| \text{REG } n$ (Primitive register pointer)

Tempates are functions:

- λ -lifted to top-level supercombinators
- In an administrative NF

$p ::= \bar{t}$ (Program)

$t ::=$
 FUN $s \alpha n =$
 let \bar{u} **in** v
(Template)

$u, v ::=$
 APP \bar{e} (Normal application)
 | **CASE** $c \bar{e}$ (Application with case table)
 | **PRIM** $n \bar{e}$ (PRS candidate allocation)

$c ::=$ **TAB** n (Case table pointer)

$e ::=$
 CON αn (Constructor tag)
 | **INT** n (Primitive integers)
 | **PRI** $\alpha \otimes$ (Primitive operation)
 | **FUN** $s \alpha n$ (Function pointer)
 | **ARG** $s n$ (Argument pointer)
 | **VAR** $s n$ (Application pointer)
 | **REG** n (Primitive register pointer)

Applications are **wide** and **tagged**

Case **alternatives** are lifted to **contiguous**
templates

$p ::= \bar{t}$	(Program)
$t ::=$ $\text{FUN } s \alpha n =$ $\text{let } \bar{u} \text{ in } v$	(Template)
$u, v ::=$ $\text{APP } \bar{e}$ $ \text{CASE } c \bar{e}$ $ \text{PRIM } n \bar{e}$	(Applications) (Normal application) (Application with case table) (PRS candidate allocation)
$c ::= \text{TAB } n$	(Case table pointer)
$e ::=$ $\text{CON } \alpha n$ $ \text{INT } n$ $ \text{PRI } \alpha \otimes$ $ \text{FUN } s \alpha n$ $ \text{ARG } s n$ $ \text{VAR } s n$ $ \text{REG } n$	(Atoms) (Constructor tag) (Primitive integers) (Primitive operation) (Function pointer) (Argument pointer) (Application pointer) (Primitive register pointer)

Usual atom suspects...
with four types of pointer¹

¹Foreshadowing is a literary device in which a writer gives an advance hint of what is to come later in the story

$p ::= \bar{t}$	(Program)
$t ::=$ $\text{FUN } s \alpha n =$ $\text{let } \bar{u} \text{ in } v$	(Template)
$u, v ::=$ $\text{APP } \bar{e}$ $ \text{CASE } c \bar{e}$ $ \text{PRIM } n \bar{e}$	(Applications) (Normal application) (Application with case table) (PRS candidate allocation)
$c ::= \text{TAB } n$	(Case table pointer)
$e ::=$ $\text{CON } \alpha n$ $ \text{INT } n$ $ \text{PRI } \alpha \otimes$ $ \text{FUN } s \alpha n$ $ \text{ARG } s n$ $ \text{VAR } s n$ $ \text{REG } n$	(Atoms) (Constructor tag) (Primitive integers) (Primitive operation) (Function pointer) (Argument pointer) (Application pointer) (Primitive register pointer)

Hardware is **fixed size**...
 Beware of \bar{u} and \bar{e}

FUN \top 2 2 =

(fromTo#T)

let APP [ARG \top 0, PRI 2 +, INT 1]
 APP [FUN \top 2 0, VAR \perp 0, ARG \perp 1]

in APP [CON 2 0, ARG \top 0, VAR \perp 1]

FUN T 2 2 =

(fromTo#T)

let APP [ARG T 0, PRI 2 +, INT 1]
 APP [FUN T 2 0, VAR ⊥ 0, ARG ⊥ 1]

in APP [CON 2 0, ARG T 0, VAR ⊥ 1]

SpineLen

Instantiate on stack

FUN T 2 2 =

ApLen

(fromTo#T)

let APP [ARG T 0, PRI 2 +, INT 1]
 APP [FUN T 2 0, VAR ⊥ 0, ARG ⊥ 1]

in APP [CON 2 0, ARG T 0, VAR ⊥ 1]

SpineLen

Instantiate on stack

Instantiate on heap

FUN T 2 2 =

ApLen

(fromTo#T)

let

APP

[ARG T 0, PRI 2 +, INT 1]

APP

[FUN T 2 0, VAR ⊥ 0, ARG ⊥ 1]

MaxAps

in

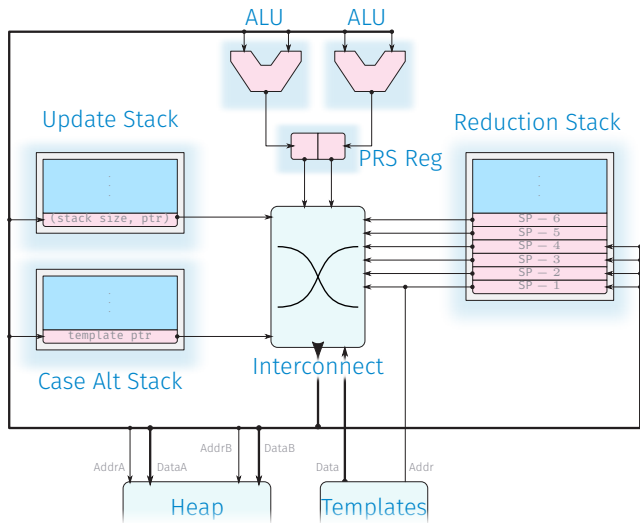
APP

[CON 2 0, ARG T 0, VAR ⊥ 1]

SpineLen

Instantiate on stack

Instantiate on heap



New features in Heron:

Parameterised implementation in Clash, for UltraScale+

Zero-constraint templates

Inline case alternatives

Postfix primitive operations

New features in Heron:

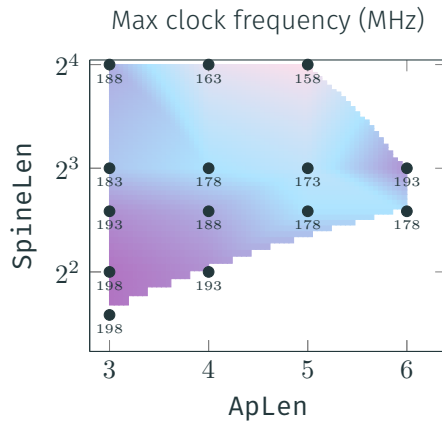
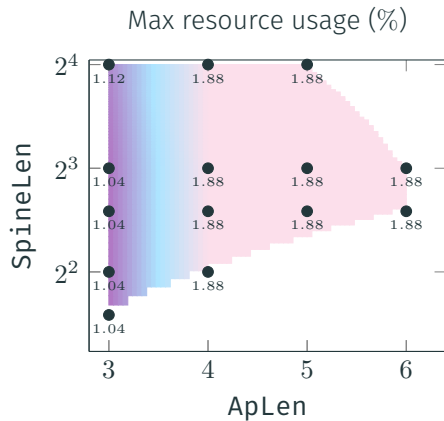
$\times 2 f_{\max}$, $< 2\%$ max usage on Alveo U280

Mean 6% reduction in cycles (max 17%)

Mean 22% reduction in code size (max 34%)

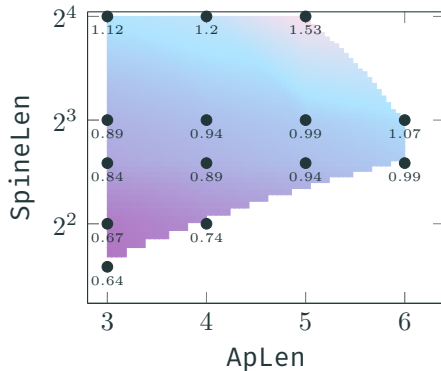
Mean 12% reduction in heap allocs (max 100%)

Circuit Results

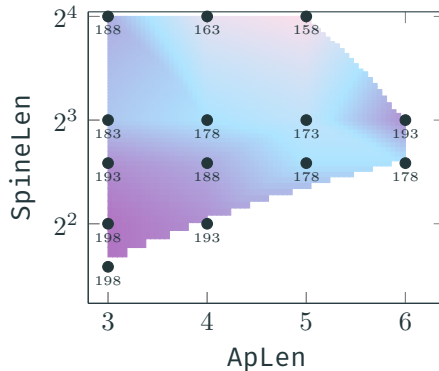


We've now got a **parameter space**.
Which, if any, is best?

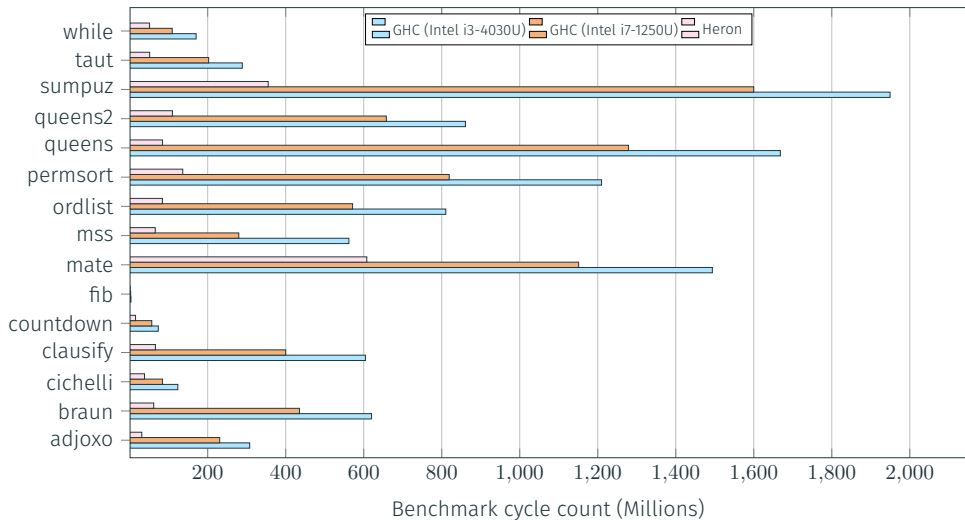
Max resource usage
without UltraRAM (%)



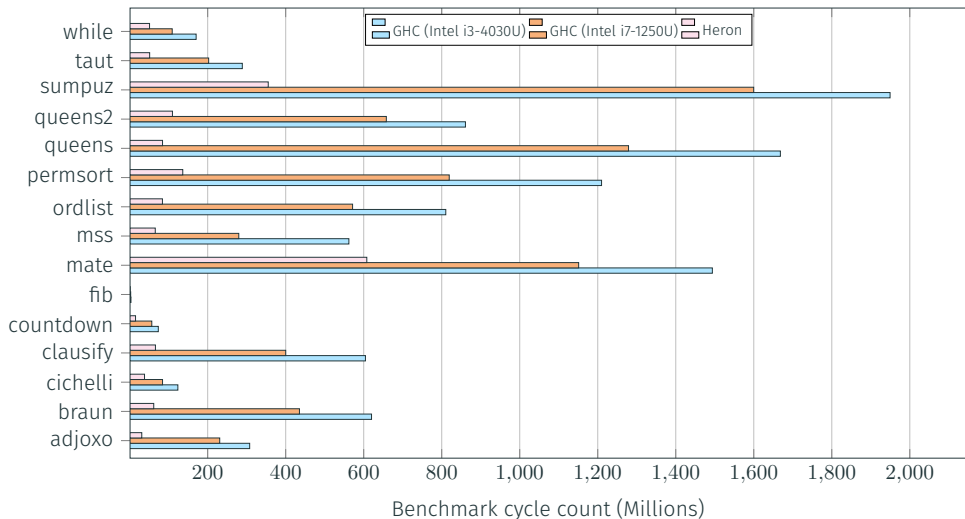
Max clock frequency (MHz)

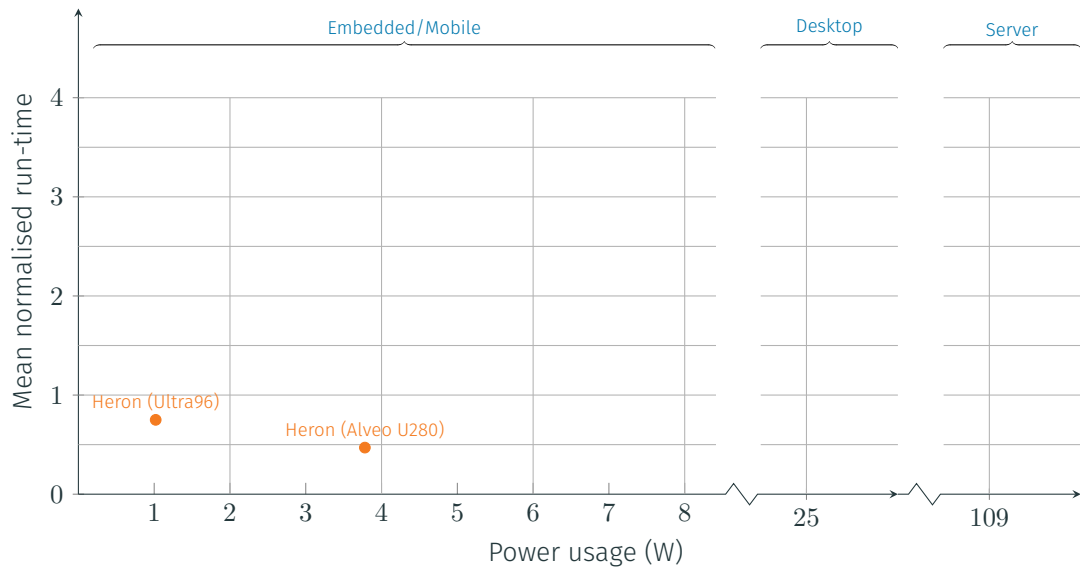


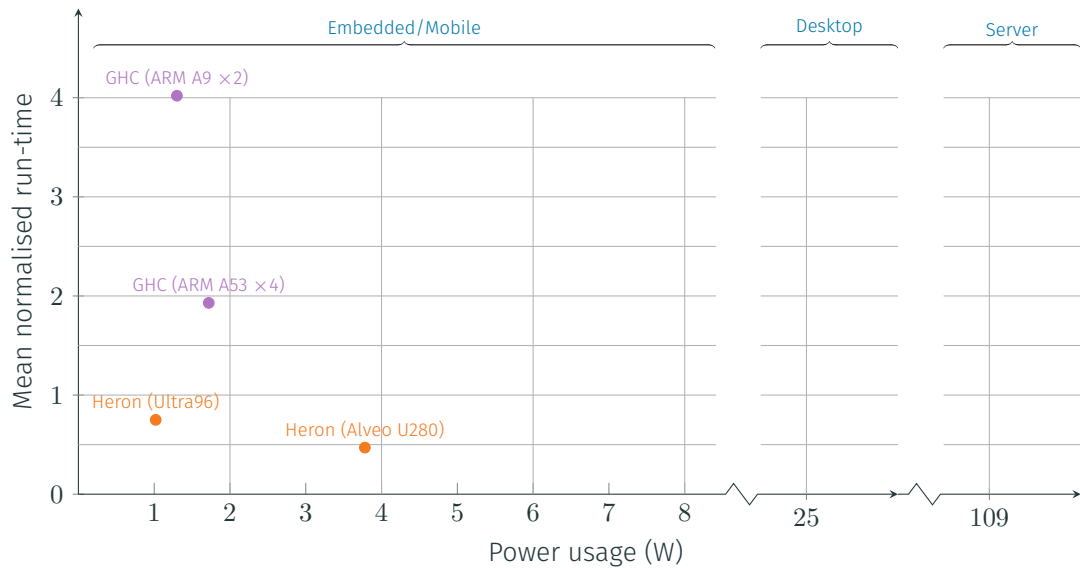
We've now got a **parameter space**.
Which, if any, is best?

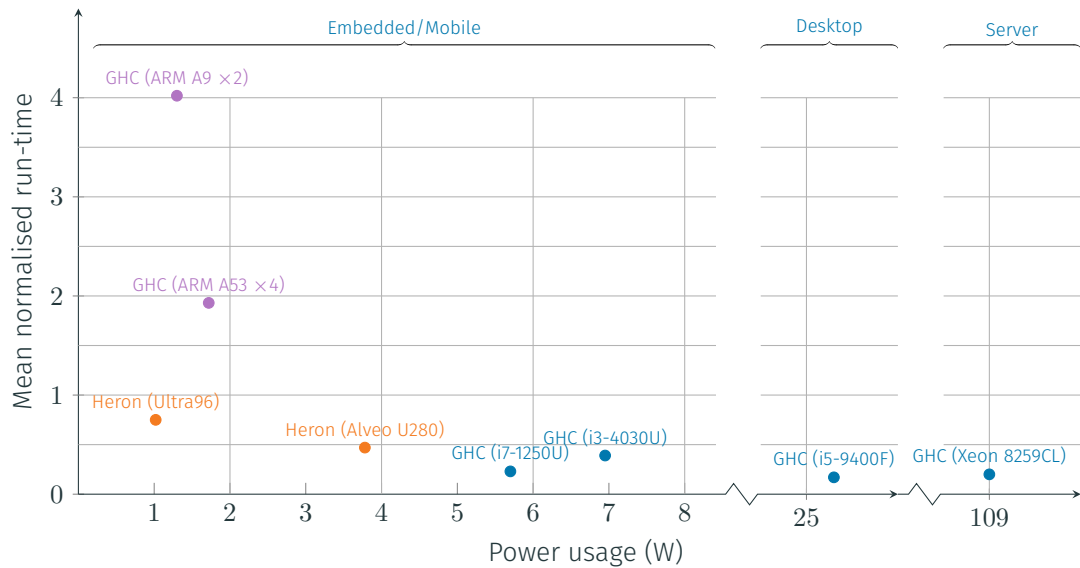


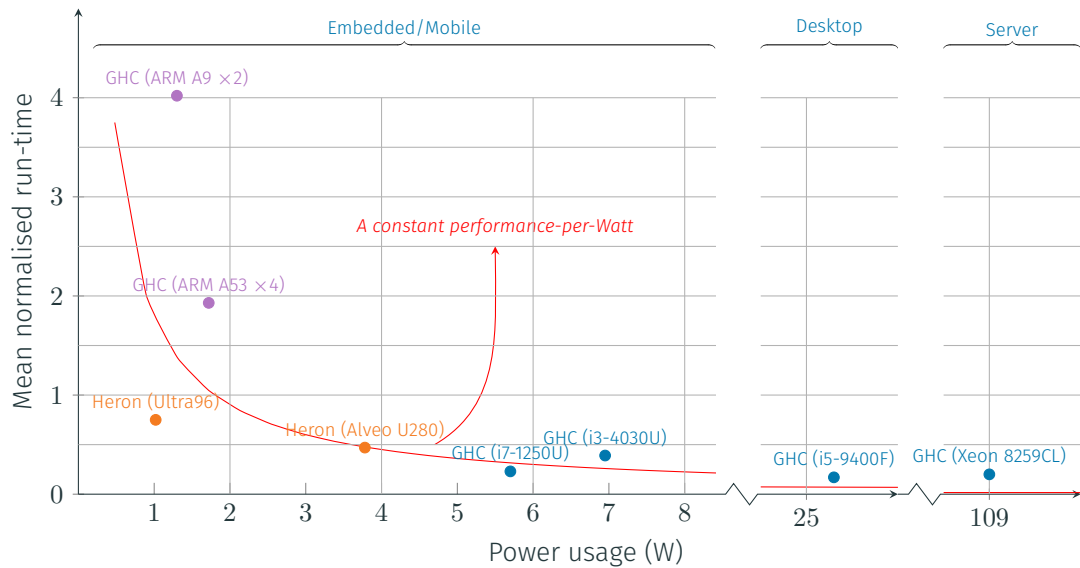
GHC + i3-4030U $\approx \times 9.5$ cycles and GHC + i7-1250U $\approx \times 6.9$ cycles











What's next?

Call To Action

CAFP mailing list

(Computer Architectures for Functional Programming)

<https://groups.google.com/g/cafp>

Appendix

```
add x y z = (x + y) + z;
```



```
add x y z = x + (y + z);
```

add x y z = (x + y) + z;



FUN T 3 0 =

(add)

```
let  APP [  ARG ⊥ 0,    PRI 2 +,    ARG ⊥ 1 ]
in   APP [  VAR ⊥ 0,    PRI 2 +,    ARG ⊥ 2 ]
```

```
fromTo n m = case n <= m of {  
  False -> Nil;  
  True  -> Cons n (fromTo (n + 1) m);  
};
```



```

fromTo n m = case n <= m of {
  False -> Nil;
  True  -> Cons n (fromTo (n + 1) m);
};

```



FUN T 2 0 = (fromTo)

```

let  APP [ ARG T 0, PRI 2 <=, ARG T 1 ]
in   CASE (TAB 1)
      [ VAR ⊥ 0, ARG T 0, ARG T 1 ]

```

FUN T 2 1 = (fromTo#F)

```

let  ∅
in   APP [ CON 0 1 ]

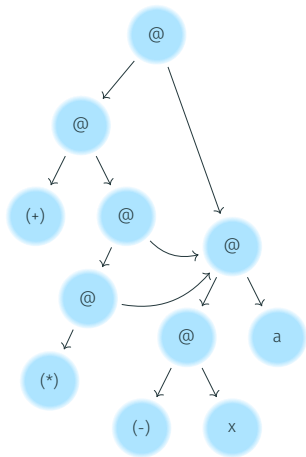
```

FUN T 2 2 = (fromTo#T)

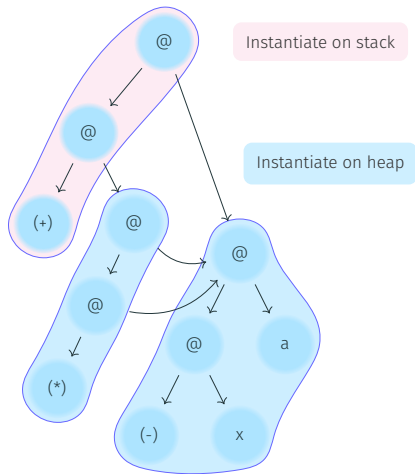
```

let  APP [ ARG T 0, PRI 2 +, INT 1 ]
      APP [ FUN T 2 0, VAR ⊥ 0, ARG ⊥ 1 ]
in   APP [ CON 2 0, ARG T 0, VAR ⊥ 1 ]

```

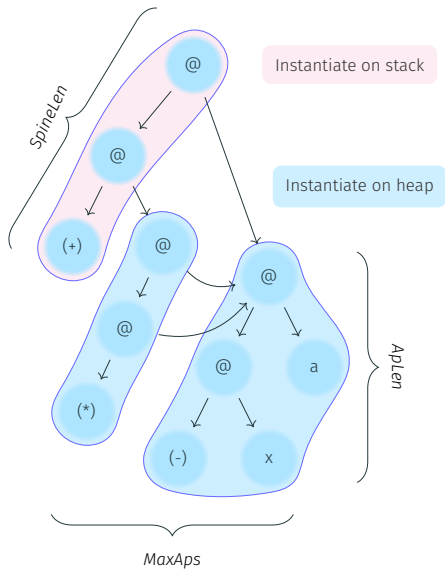


```
y a x = let b = (-) x a  
        c = (*) b b  
        in (+) c b
```

```

y a x = let b = (-) x a
          c = (*) b b
          in (+) c b
  
```



```

y a x = let b = (-) x a
        c = (*) b b
        in (+) c b
  
```