

# *Once bittern, twice shy*

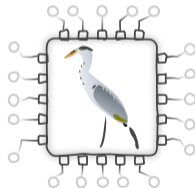
*Revisiting hardware architectures for lazy functional languages with Heron*

---

*Craig Ramsay & Rob Stewart*

*September 2024*

*Heriot-Watt University*

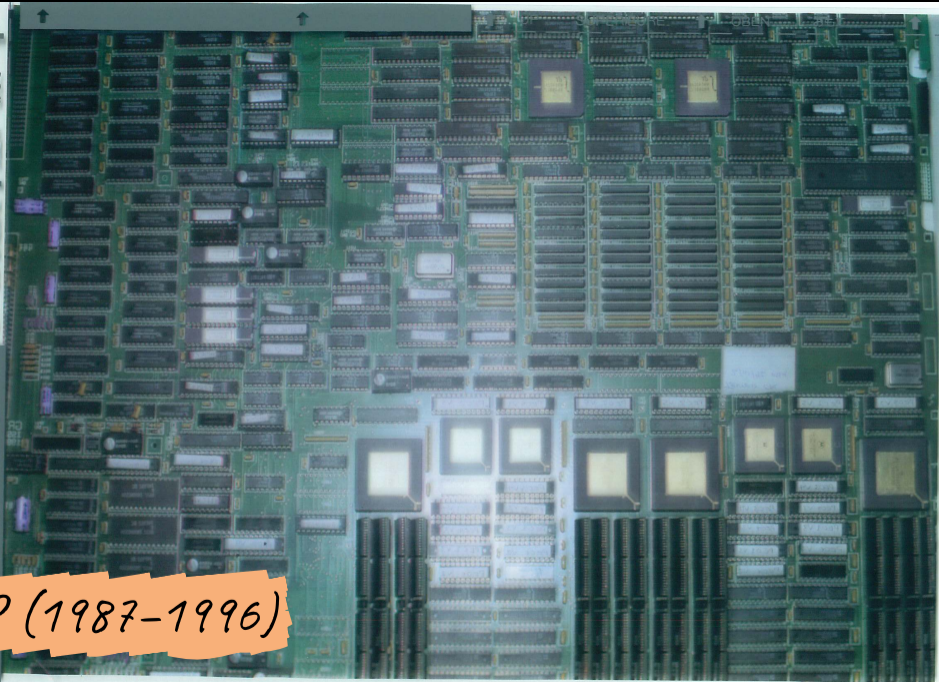




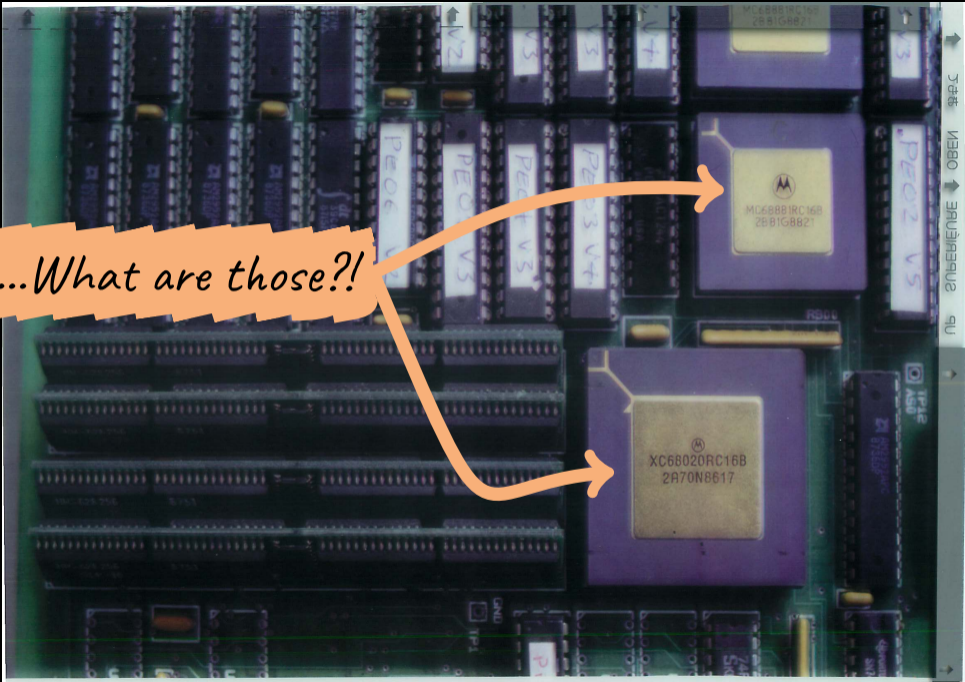
<https://haflang.github.io/history>

UP SUPERIEURE OBEN 上

GRIP (1987-1996)



...What are those?!

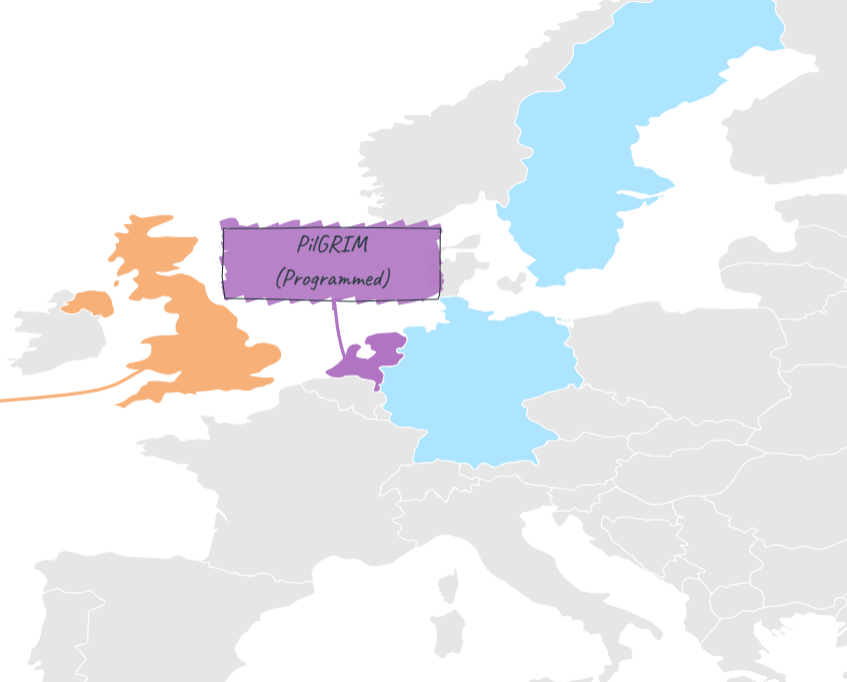




Reduceron  
(Template inst.)

PiIGRIM  
(Programmed)

≈ 2010

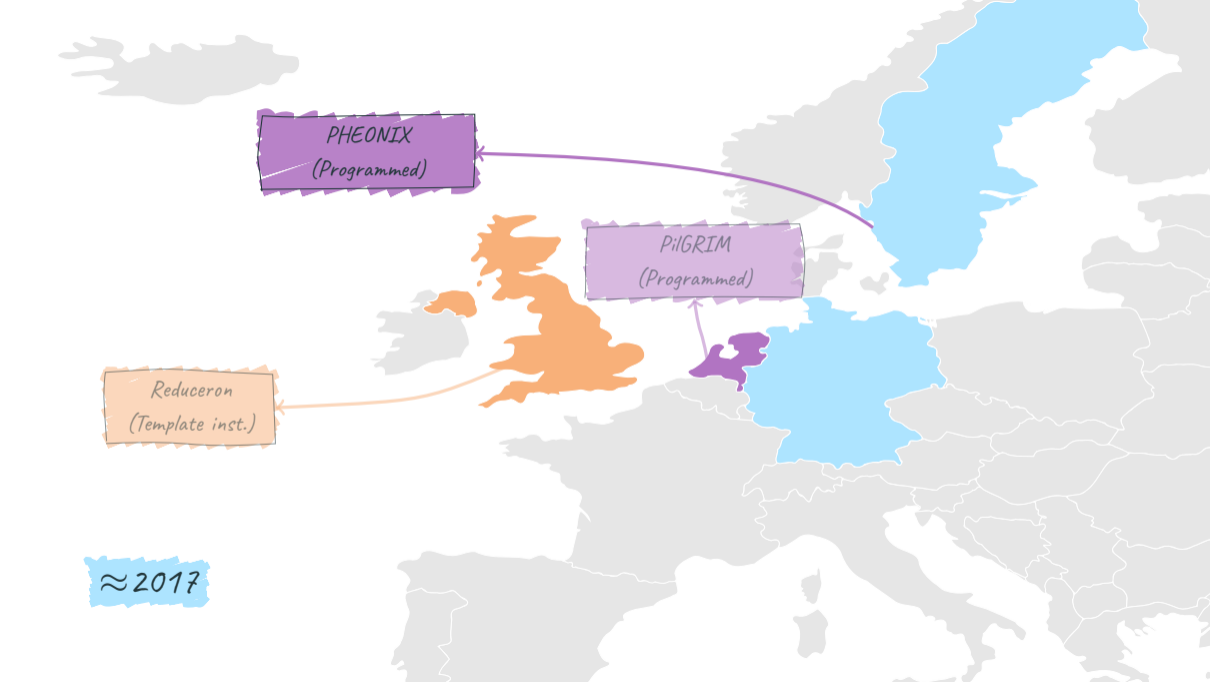


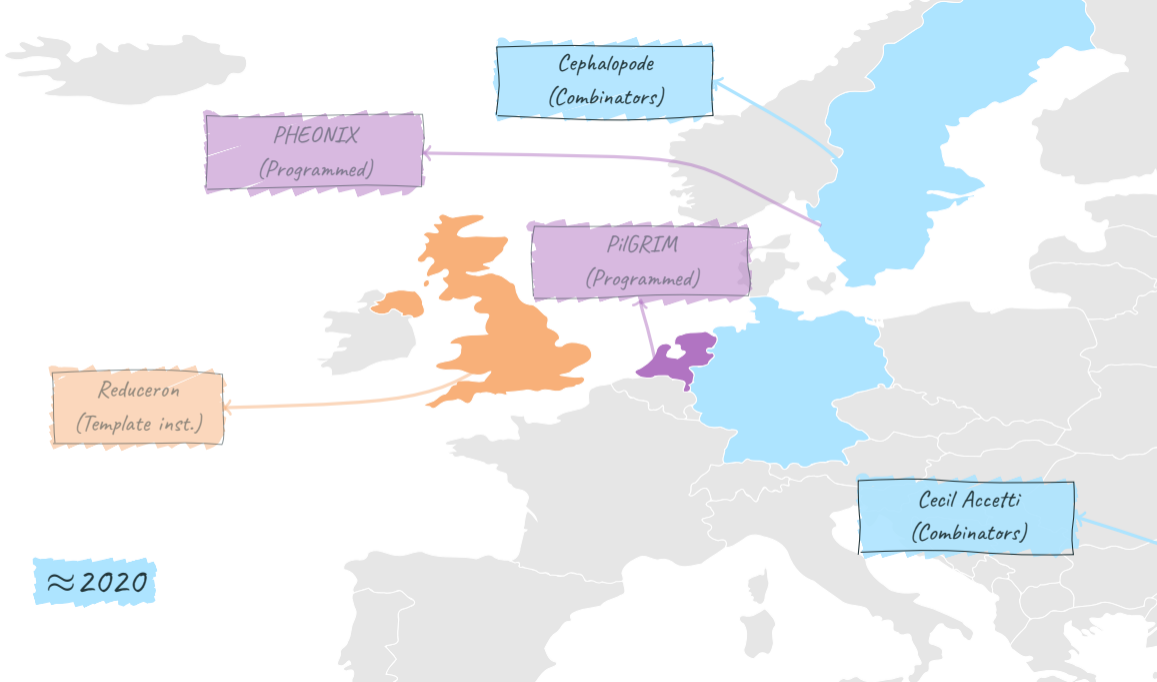
PHEONIX  
(Programmed)

PiIGRIM  
(Programmed)

Reduceron  
(Template inst.)

≈ 2017





Cephalopode  
(Combinators)

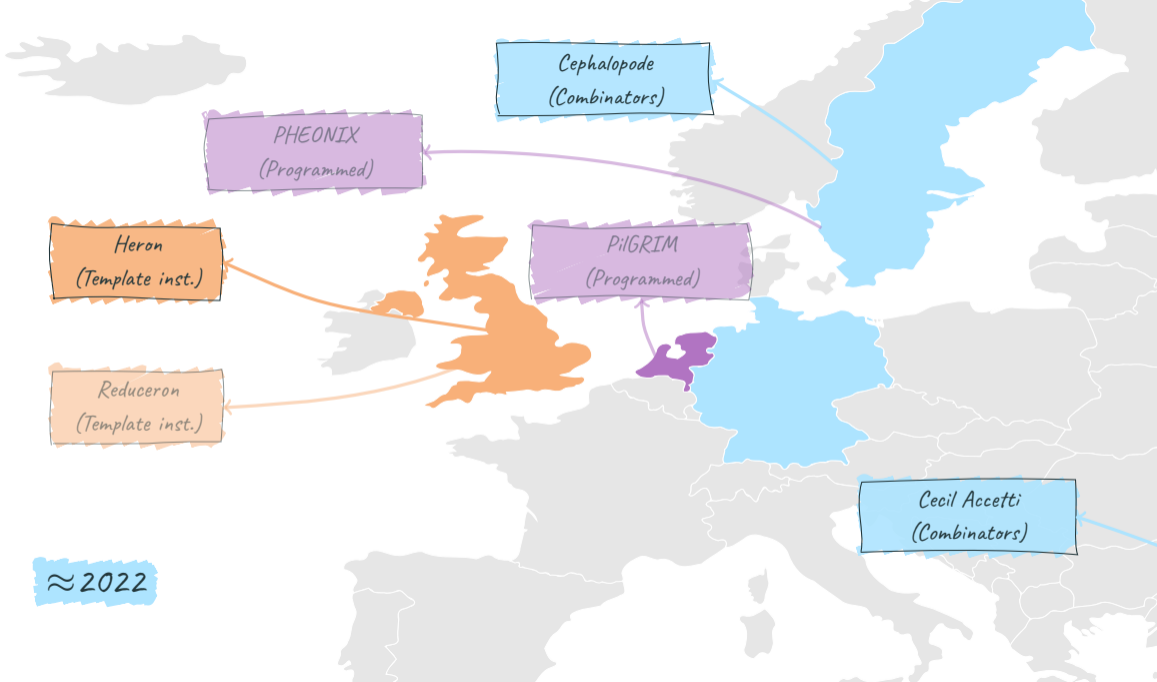
PHEONIX  
(Programmed)

PiIGRIM  
(Programmed)

Reduceron  
(Template inst.)

Cecil Accetti  
(Combinators)

≈ 2020



Cephelopode  
(Combinators)

PHEONIX  
(Programmed)

Heron  
(Template inst.)

PilGRIM  
(Programmed)

Reduceron  
(Template inst.)

Cecil Accetti  
(Combinators)

≈ 2022



Cephalopode

(Combinators)

PHEONIX

(Programmed)

Heron

(Template inst.)

PiIGRIM

(Programmed)

Reduceron

(Template inst.)

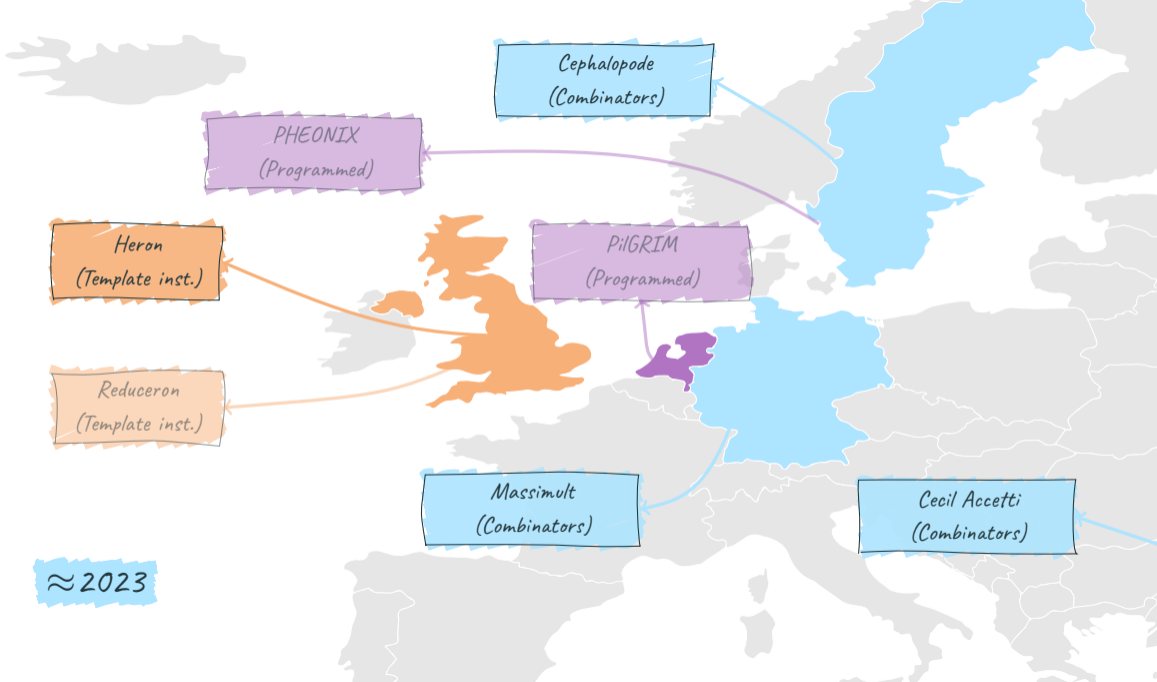
Massimult

(Combinators)

Cecil Accetti

(Combinators)

≈ 2023

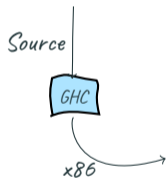


# *HAFLANG Project*

---

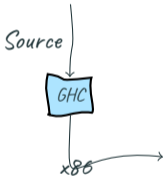
`add x y z =`

`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)

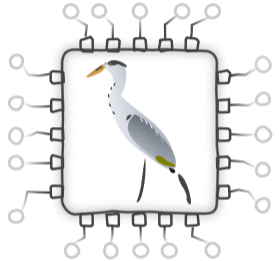
# Heron

noun [C]

/'herən/

*A graph reduction processor.*

*Performs template instantiation in one clock cycle via multiple, wide, multi-ported memories.*



°Ramsay and Stewart, "Heron: Modern Graph Reduction Hardware".

# Heron

noun [C]

/'herən/

*A graph reduction processor.*

*Performs template instantiation in one clock cycle via multiple, wide, multi-ported memories.*



*°Ramsay and Stewart, "Heron: Modern Graph Reduction Hardware".*

## Template example

`enumFrom :: Int -> [Int]`

`enumFrom n = let m = n + 1`

`ns = enumFrom m`

`in Cons n ns`



## Template example

$enumFrom :: Int \rightarrow [Int]$

$enumFrom\ n = let\ m = n + 1$

$ns = enumFrom\ m$

$in\ Cons\ n\ ns$

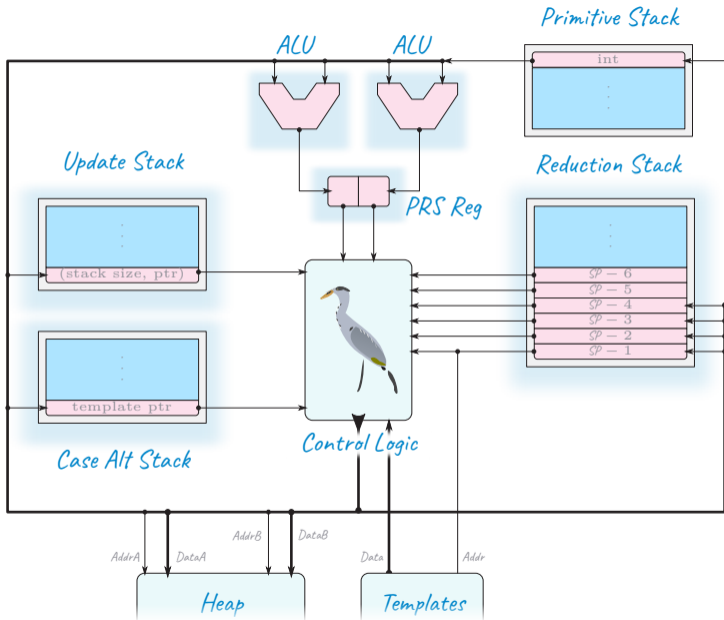


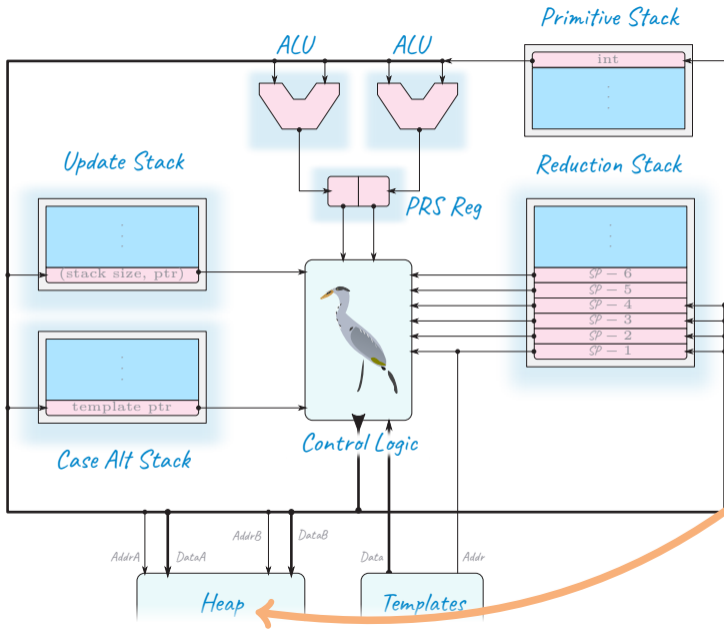
let APP [ ARG True 0, PRI 2 +, INT 1 ]

APP [ FUN True 1 0, VAR False 0, ]

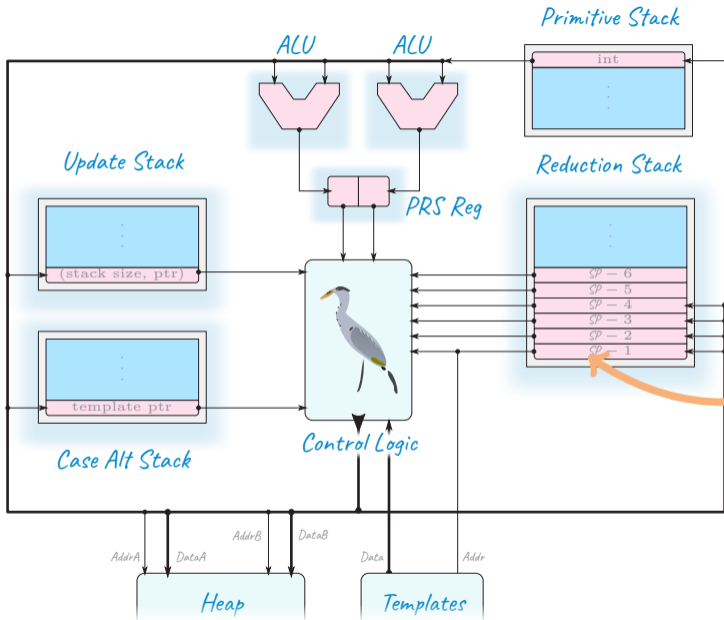
in APP [ CON 2 0, ARG True 0, VAR False 1 ]



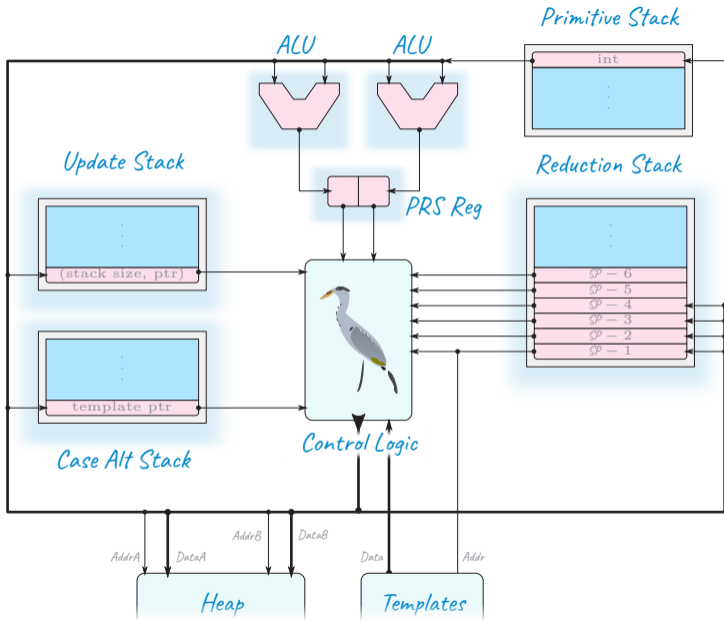




Heap is on-chip & small  
(think L1/L2 cache)

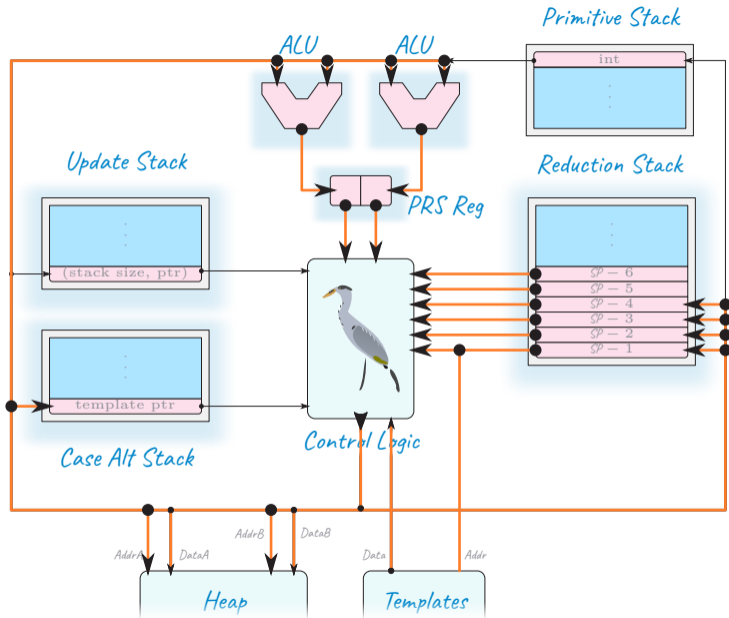


Next reduction rule always\*  
chosen by matching  
top of stack



Atoms

= FUN  $s_n$   
 | CON  $a_n$   
 | VAR  $s_n$   
 | INT  $n$   
 | PRI  $a \otimes$   
 | ARG  $s_n$   
 | REG  $n$



Atoms

= FUN s a n

| CON a n

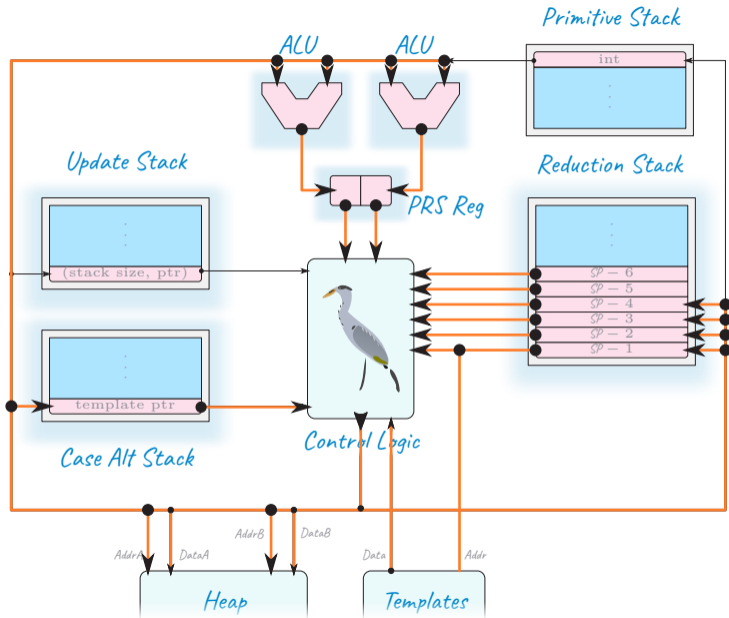
| VAR s n

| INT n

| PRI a ⊗

| ARG s n

| REG n



Atoms

= FUN  $s_n$

| CON  $a_n$

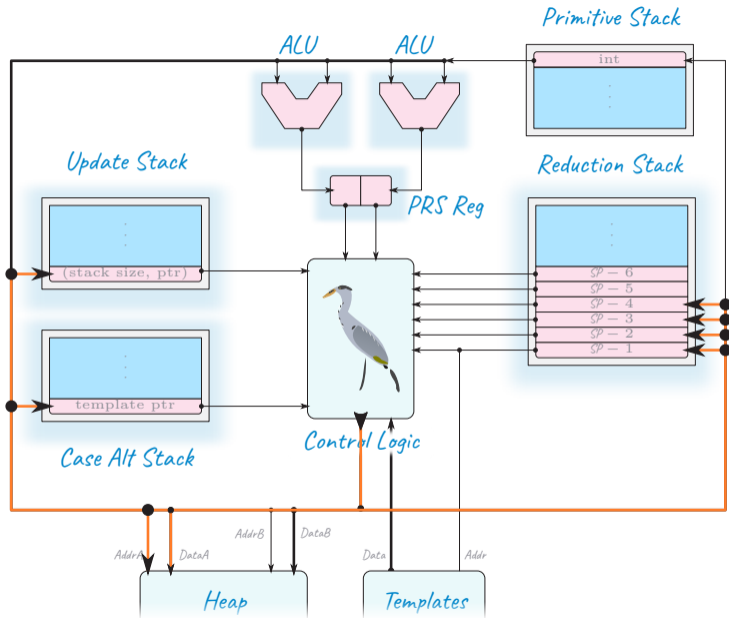
| VAR  $s_n$

| INT  $n$

| PRI  $a \otimes$

| ARG  $s_n$

| REG  $n$



Atoms

=  $FUN s a n$

|  $CON a n$

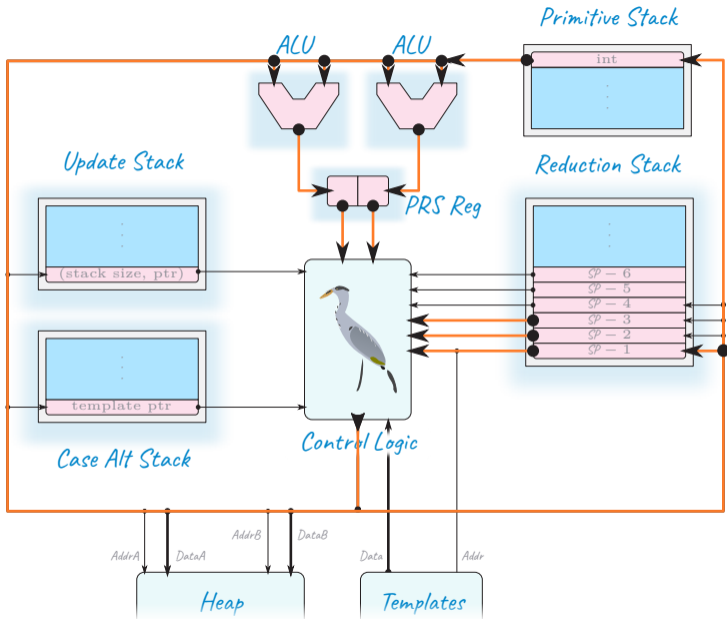
|  $VAR s n$

|  $INT n$

|  $PRI a \otimes$

|  $ARG s n$

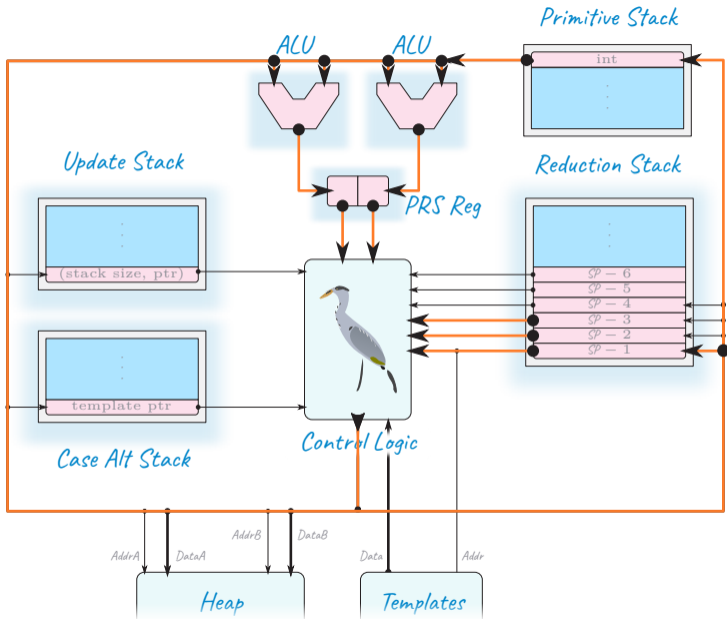
|  $REG n$



Atoms

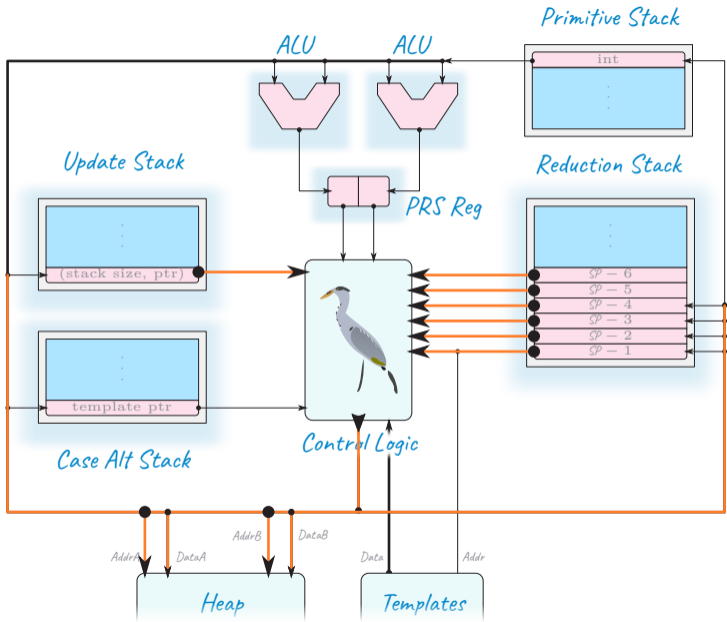
- = FUN s a n
- | CON a n
- | VAR s n
- | INT n
- | PRI a ⊗
- | ARG s n
- | REG n



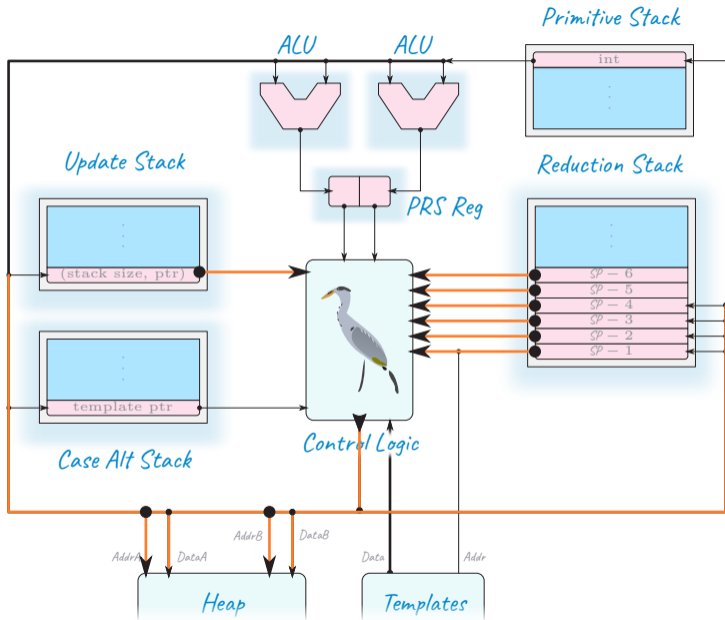


Postfix prims for long spines

$$(f \times y) + (g z) \Rightarrow f \times y \ g \ z +$$



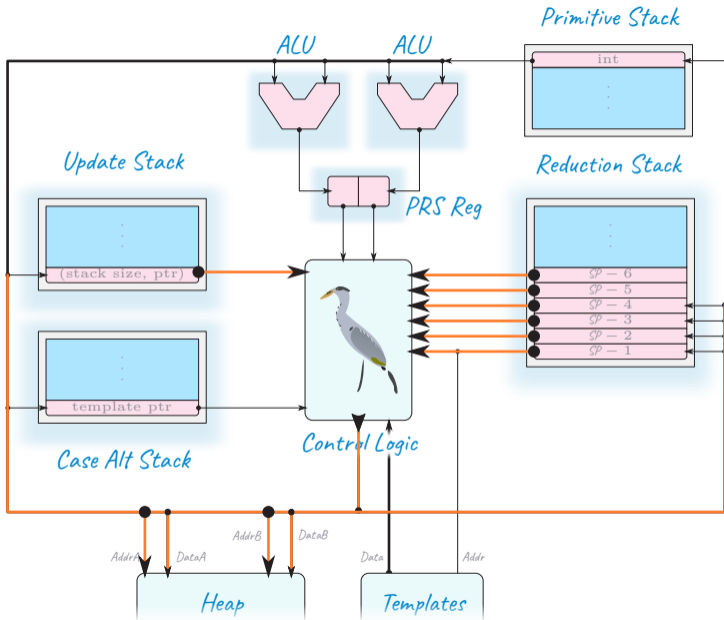
...But what about heap updates?



Avoid most updates via  
run-time sharing analysis!

Atoms

=  $FUN s a n$   
 |  $CON a n$   
 |  $VAR s n$   
 |  $INT n$   
 |  $PRI a \otimes$   
 |  $ARG s n$   
 |  $REG n$



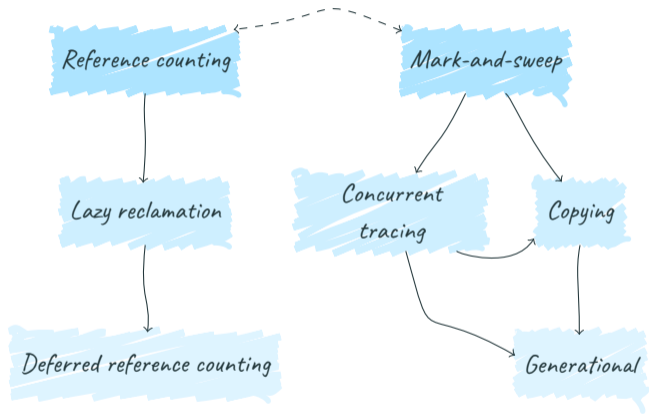
Avoid most updates via  
run-time sharing analysis!

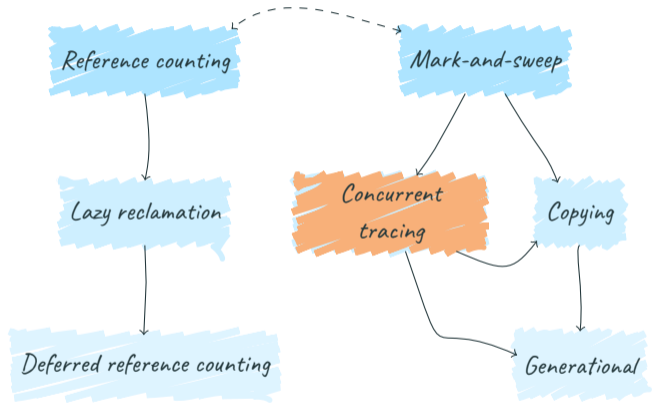
Atoms

$= FUN s a n$   
 $| CON a n$   
 $| VAR s n$   
 $| INT n$   
 $| PRI a \otimes$   
 $| ARG s n$   
 $| REG n$

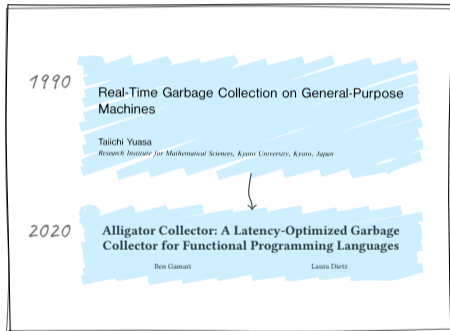
$\approx$  One-bit

reference counting!

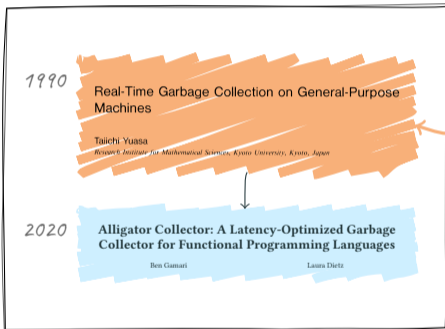




## Software for Concurrent GC



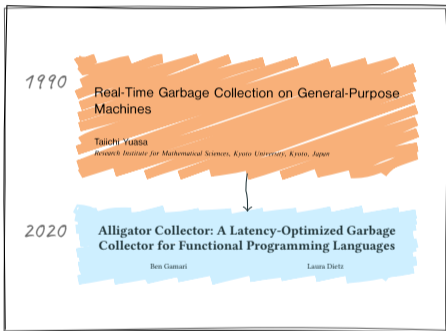
## Software for Concurrent GC



*"On these [von Neumann style] machines, real-time garbage collection inevitably causes some overhead on the overall execution"*



## Software for Concurrent GC

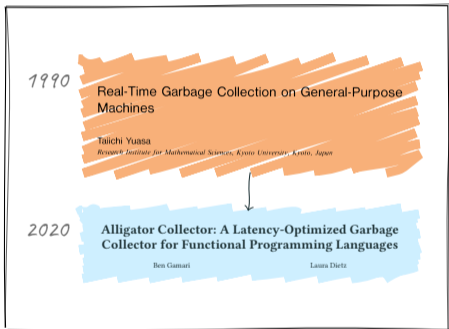


Function alloc (app):

```
heap[hp] ← app
```

```
hp++
```

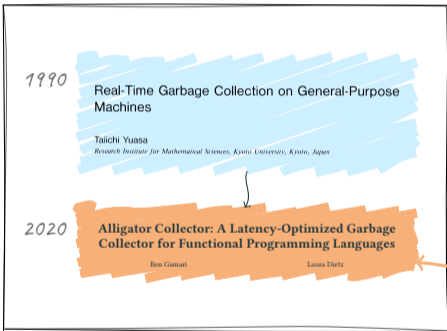
## Software for Concurrent GC



Function alloc (app):

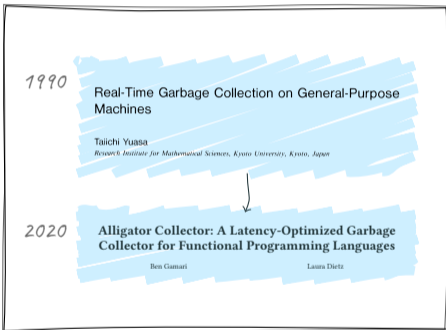
```
a ← pop from freelist
if allocBarrier(gcPhase, a)
then
  | tag a as Marked
else
  | tag a as Unmarked
  | heap[a] ← app
```

## Software for Concurrent GC



*"The nofib cases are quite mixed [...] most tests slow down, with a median of +21%"*

## Software for Concurrent GC



### Key Observation:

Stock CPUs sequentialise write-barriers  
(trades-off GC latency for throughput)

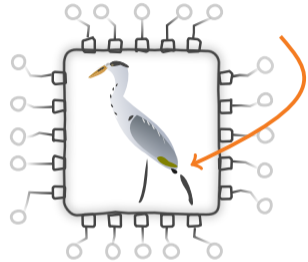
Custom hardware + read-  
first memories grants us both

# Cloaca

noun [C]

/kloh-ah-kuh/

*A concurrent hardware  
garbage collector for Heron*



*°Ramsay and Stewart, "Cloaca: A Concurrent Hardware Garbage Collector for Non-strict Functional Languages".*

# Cloaca

noun [C]

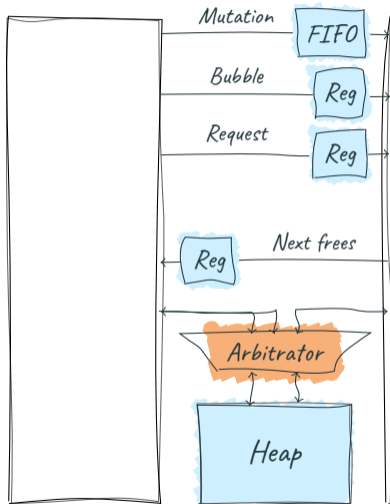
/kloh-ah-kuh/

*A concurrent hardware  
garbage collector for Heron*

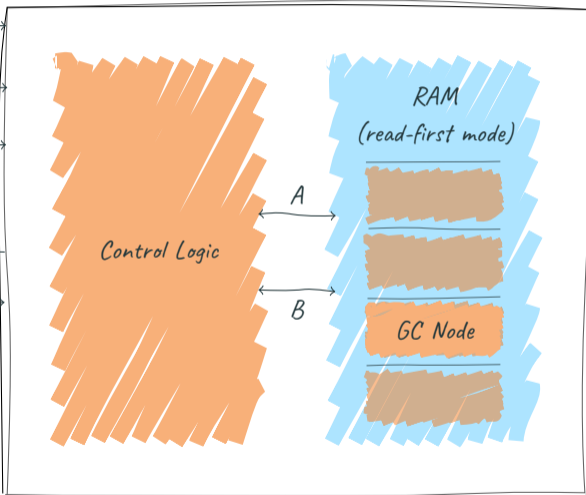


*°Ramsay and Stewart, "Cloaca: A Concurrent Hardware Garbage Collector for Non-strict Functional Languages".*

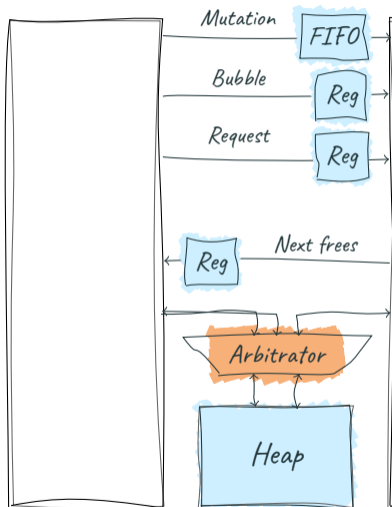
## Reduction Core



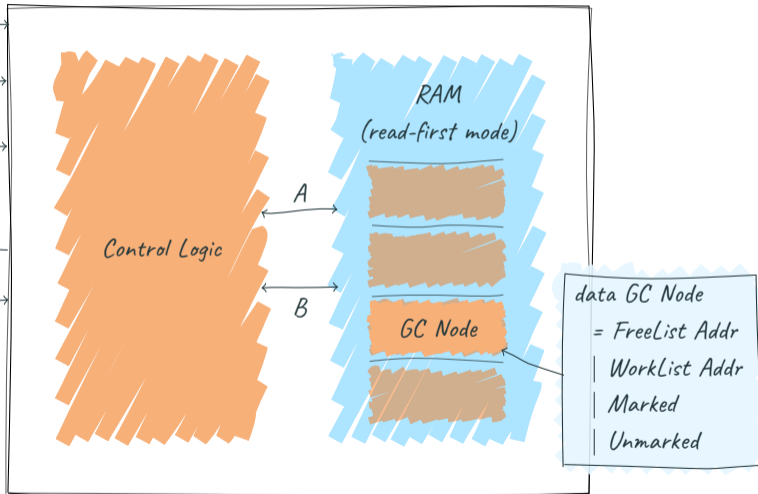
## Memory Management



## Reduction Core

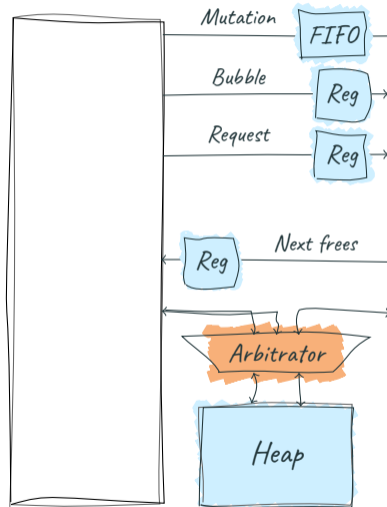


## Memory Management

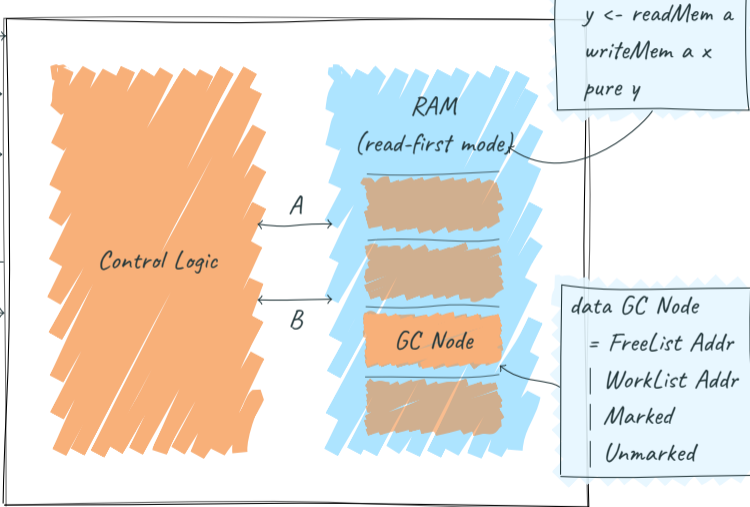




## Reduction Core

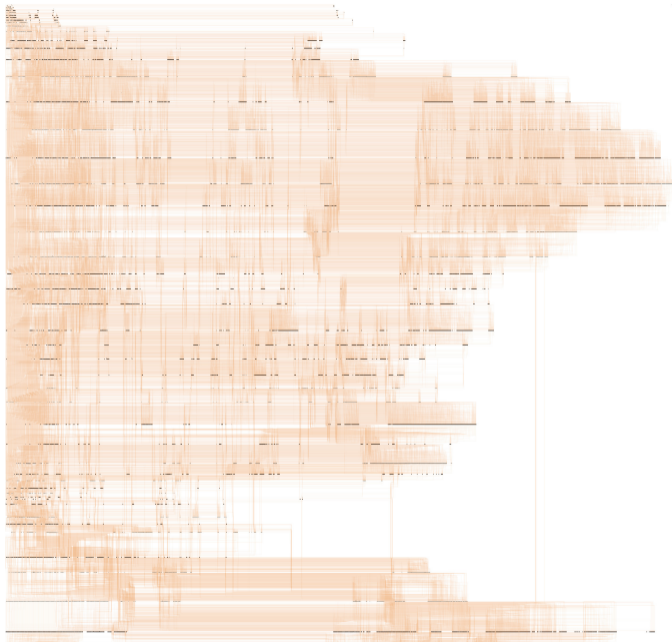


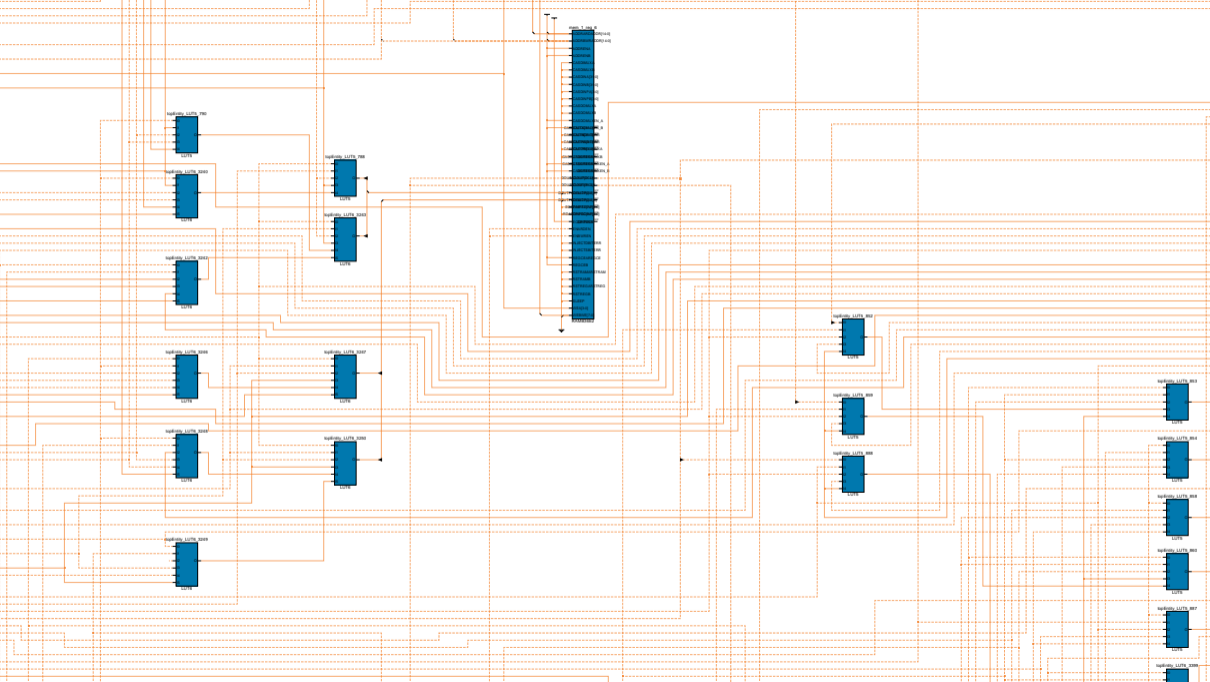
## Memory Management



# *Results*

---





# Which Architectures?

<i>Platform</i>	<i>Heron</i> <i>Xilinx Alveo U280</i>	<i>GHC + Intel i7 1250U</i>	<i>Performance</i>   <i>Power-saver</i>
-----------------	--	-----------------------------	---

# Which Architectures?

	<i>Heron</i>	<i>GHC + Intel i7 1250U</i>	
<i>Platform</i>	<i>Xilinx Alveo U280</i>	<i>Performance</i>	<i>Power-saver</i>
<i>Clock</i>	<i>185 MHz</i>	<i>4.7 GHz</i>	<i>≈ 2 GHz</i>

# Which Architectures?

	Heron	GHC + Intel i7 1250U	
Platform	Xilinx Alveo U280	Performance	Power-saver
Clock	185 MHz	4.7 GHz	≈ 2 GHz
Power est.	0.8W dynamic + 3.1 W Static <sup>1</sup>	Cores 15 W or Package 16 W	Cores 2W or Package 6 W

<sup>1</sup>Heron only occupies 1.13% of any resource type though!

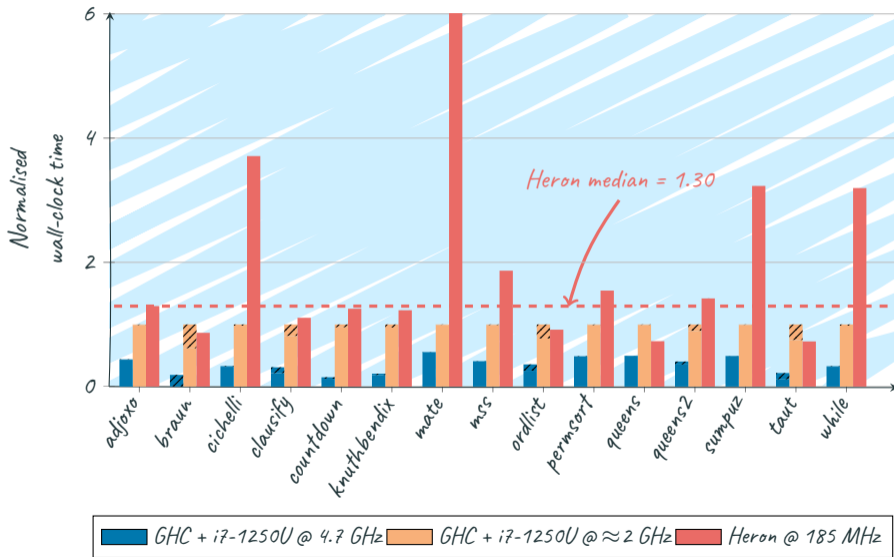
# Which Architectures?

	Heron	GHC + Intel i7 1250U	
Platform	Xilinx Alveo U280	Performance	Power-saver
Clock	185 MHz	4.7 GHz	≈ 2 GHz
Power est.	0.8W dynamic + 3.1 W Static <sup>1</sup>	Cores 15 W or Package 16 W	Cores 2W or Package 6 W
Fabrication	16 nm (FPGA!)	10 nm	10 nm

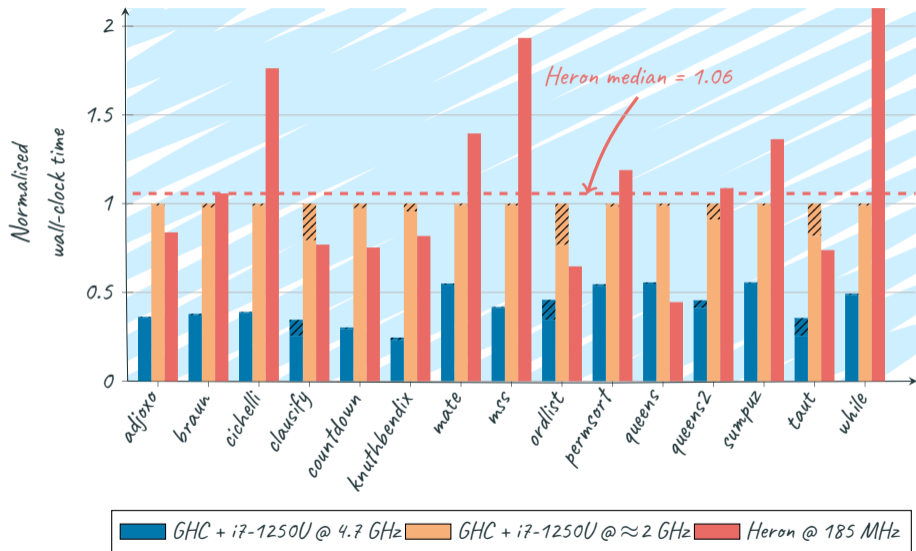
<sup>1</sup>Heron only occupies 1.13% of any resource type though!



## Wall-clock times vs GHC -O2



## Wall-clock times vs GHC -O0



## Conclusion

*A tiny template instantiation core nearly keeps up with a modern CPU  
running at x10 speed!*

## Conclusion

*A tiny template instantiation core nearly keeps up with a modern CPU  
running at x10 speed!*

*Assisted by a fully concurrent GC (modulo  $\approx 20$  cycles per pass).*

## Conclusion

*A tiny template instantiation core nearly keeps up with a modern CPU  
running at x10 speed!*

*Assisted by a fully concurrent GC (modulo  $\approx 20$  cycles per pass).*

*Lays a path towards a single-chip multi-core architecture.*

## Conclusion

*A tiny template instantiation core nearly keeps up with a modern CPU  
running at x10 speed!*

*Assisted by a fully concurrent GC (modulo  $\approx 20$  cycles per pass).*

*Lays a path towards a single-chip multi-core architecture.*

*We want to see more research towards custom FP architectures!*

*Questions?*

---