

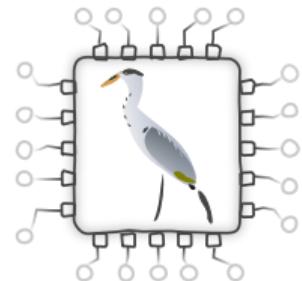
Once bittern, twice shy

Revisiting hardware architectures for lazy functional languages with Heron

Craig Ramsay & Rob Stewart

September 2024

Heriot-Watt University





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



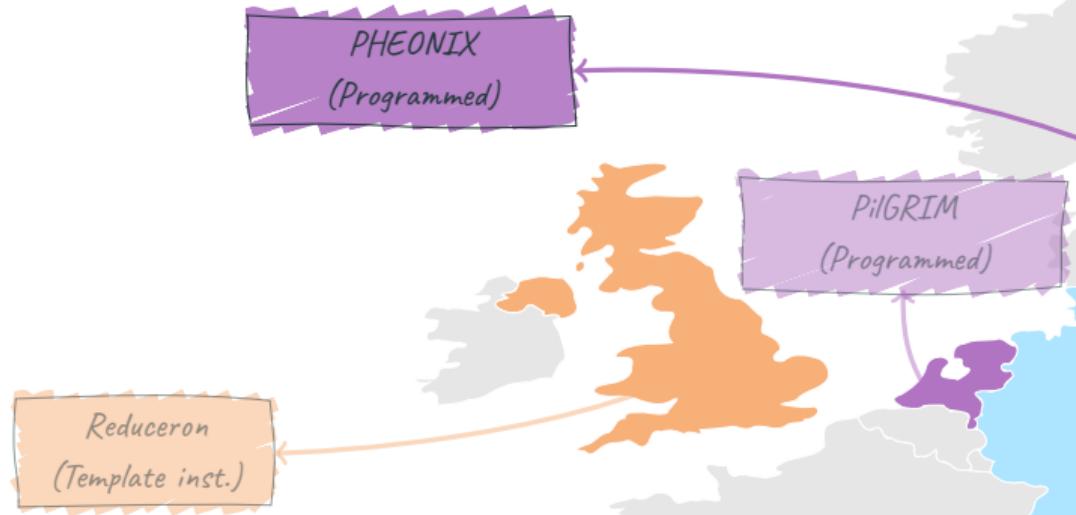
Craig 29.5.99 Self portrait.

the place of
new strat
Universität
The University
registered

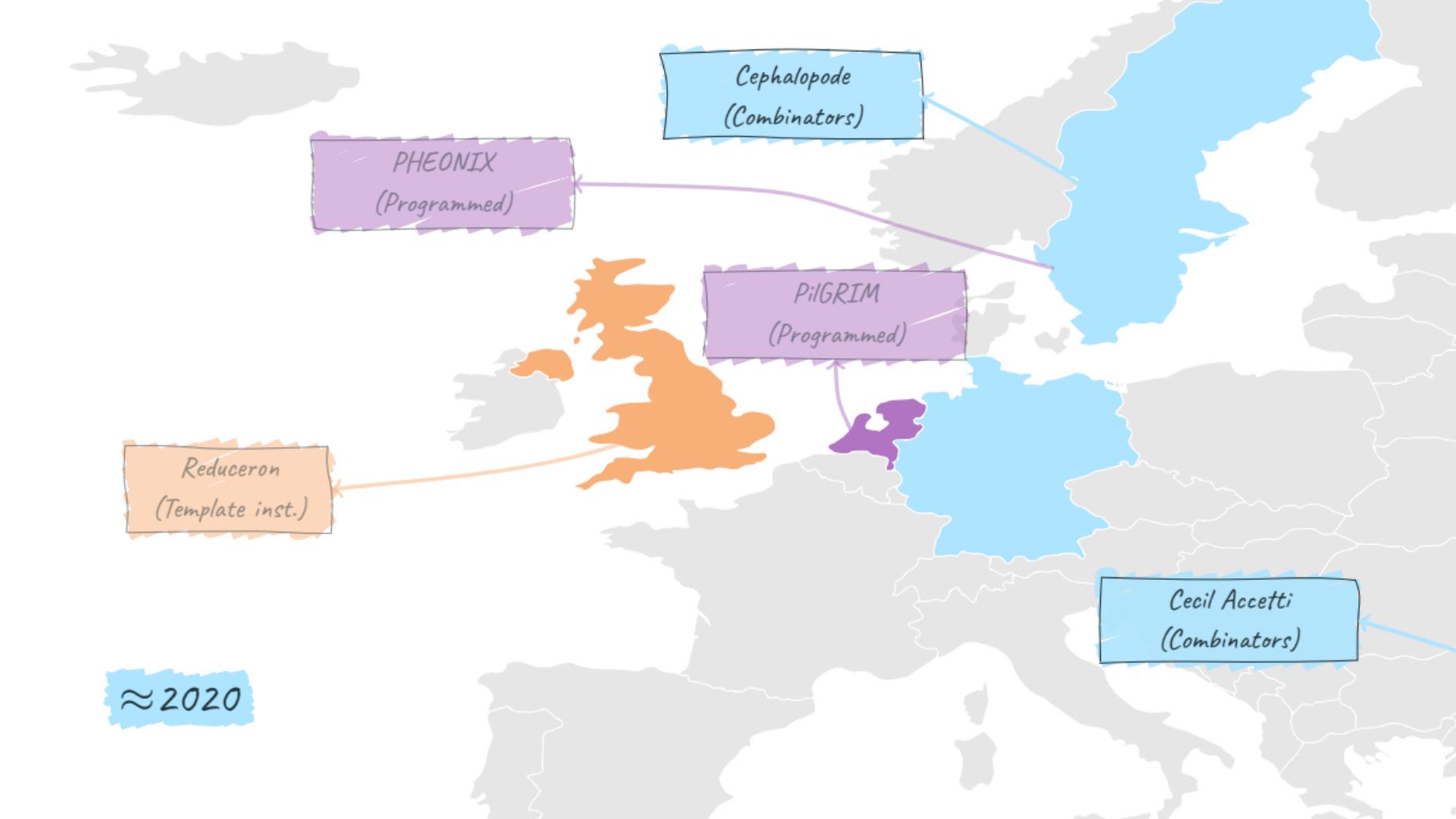
Reduceron
(Template inst.)

PilGRIM
(Programmed)

≈ 2010



≈ 2017



A map of Europe with various cloud patterns in shades of grey, white, orange, purple, and blue. Overlaid on the map are several text boxes connected by arrows.

Reduceron
(Template inst.)

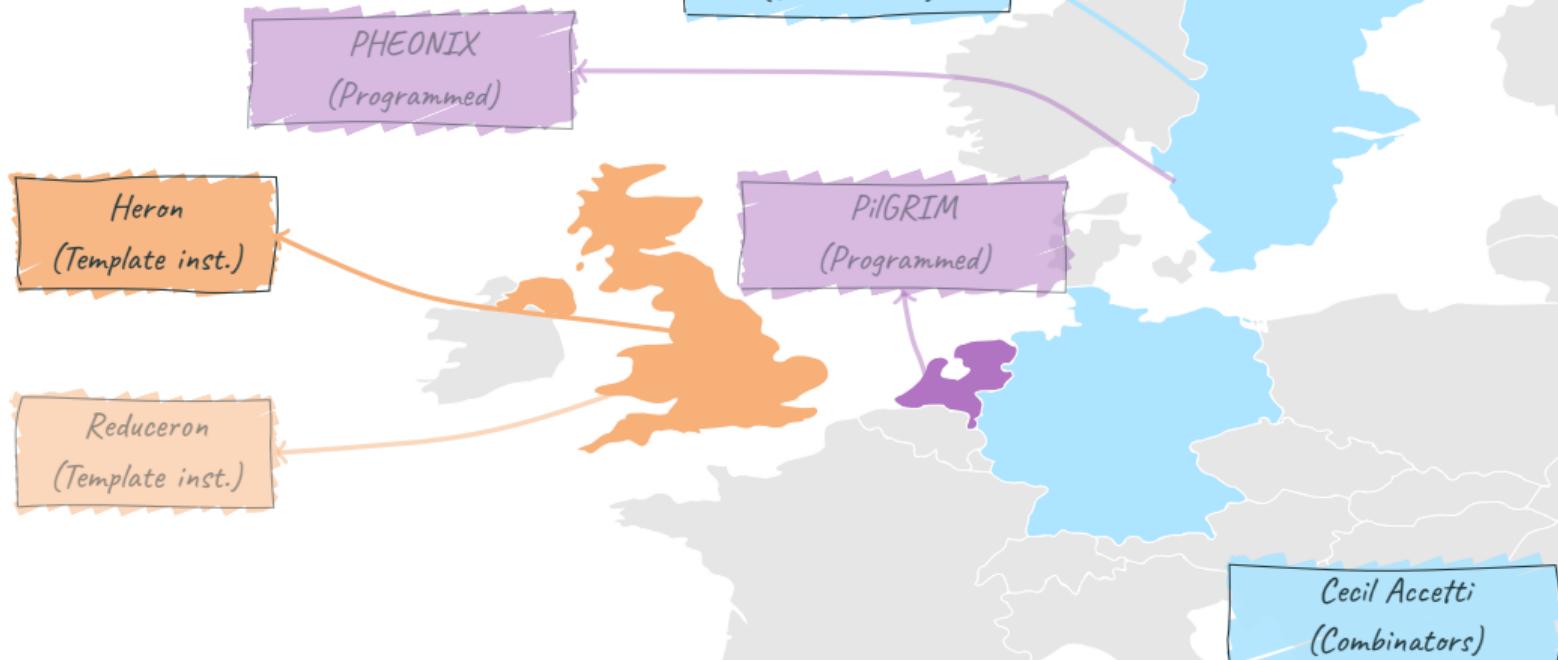
PHEONIX
(Programmed)

Cephalopode
(Combinators)

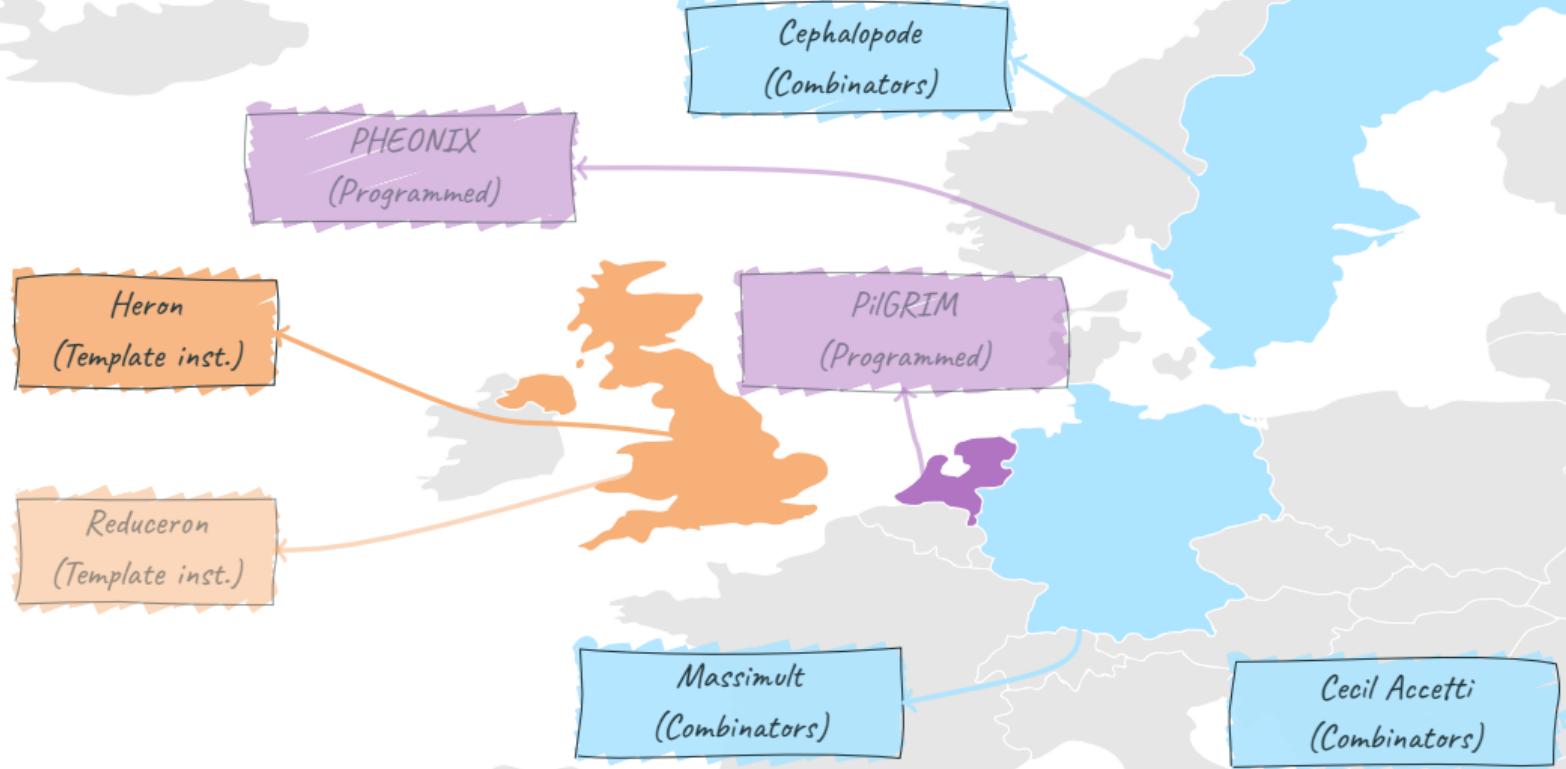
PilGRIM
(Programmed)

Cecil Accetti
(Combinators)

≈ 2020



≈ 2022



≈ 2023

HAFLANG Project

Heron

noun [C]

/'herən/

A graph reduction processor.

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



^oRamsay and Stewart, "Heron: Modern Graph Reduction Hardware".

Translating functions to templates

1) Source

sum acc xs = case xs of

[] -> acc

(y:ys) -> sum (acc + y) ys

Translating functions to templates

1) Source



2) Lifted case alts

$\text{sum acc } xs = \text{case } xs \text{ of}$

$[] \rightarrow acc$

$(y:ys) \rightarrow \text{sum} (acc + y) ys$

$\text{sum acc } xs = \text{case } xs \text{ of}$

$\text{Nil} \rightarrow \text{altNil acc}$

$\text{Cons } y \text{ } ys \rightarrow \text{altCons } y \text{ } ys \text{ acc}$

$\text{altNil acc} = acc$

$\text{altCons } y \text{ } ys \text{ acc} = \text{sum} (acc + y) ys$

Translating functions to templates

2) Lifted case alts

sum acc xs = *case* xs of

Nil → altNil acc

Cons y ys → altCons y ys acc

altNil acc = acc

altCons y ys acc = sum (acc + y) ys

Translating functions to templates

3) Using case tables



2) Lifted case alts

$\text{sum acc } xs = xs \langle \text{altCon}; \text{altNil} \rangle \text{ acc}$

$\text{altNil acc} = acc$

$\text{altCon } y \text{ } ys \text{ acc} = \text{sum } (\text{acc} + y) \text{ } ys$

$\text{sum acc } xs = \text{case } xs \text{ of}$

$\text{Nil} \rightarrow \text{altNil acc}$

$\text{Cons } y \text{ } ys \rightarrow \text{altCon } y \text{ } ys \text{ acc}$

$\text{altNil acc} = acc$

$\text{altCon } y \text{ } ys \text{ acc} = \text{sum } (\text{acc} + y) \text{ } ys$

Translating functions to templates

3) Using case tables

sum acc xs = xs ~~<altCon; altNil>~~ acc

altNil acc = acc

altCon y ys acc = sum (acc + y) ys

Translating functions to templates

3) Using case tables



4) Flatten to ANF

sum acc xs = xs <altCon;altNil> acc

altNil acc = acc

altCon y ys acc = sum (acc + y) ys

sum acc xs = xs <altCon;altNil> acc

altNil acc = acc

altCon y ys acc = let z = acc + y
in sum z ys

Template syntax

$e ::=$

Atoms	
$CON\ a\ n$	(Constructor tag)
$INT\ n$	(Primitive integers)
$PRI\ a\ \otimes$	(Primitive operation)
$FUN\ s\ a\ n$	(Function pointer)
$ARG\ s\ n$	(Argument pointer)
$VAR\ s\ n$	(Application pointer)
$REG\ n$	(Primitive register pointer)

$u, v ::=$

Applications	
$APP\ \bar{e}$	(Normal application)
$CASE\ c\ \bar{e}$	(App with case table)
$PRIM\ n\ \bar{e}$	(PRS candidate)

$c ::= TAB\ n$

Case table pointer

$t ::= let\ \bar{u}\ in\ v$

Template

where $n :: Int$, $a :: Arity$, $s :: IsShared$

Template syntax

$e ::=$ Atoms

$CON\ a\ n$	(Constructor tag)
$INT\ n$	(Primitive integers)
$PRI\ a\ \otimes$	(Primitive operation)
$FUN\ s\ a\ n$	(Function pointer)
$ARG\ s\ n$	(Argument pointer)
$VAR\ s\ n$	(Application pointer)
$REG\ n$	(Primitive register pointer)

$u, v ::=$ Applications

$APP\ \bar{e}$	(Normal application)
$CASE\ c\ \bar{e}$	(App with case table)
$PRIM\ n\ \bar{e}$	(PRS candidate)

$c ::= TAB\ n$ Case table pointer

$t ::= let\ \bar{u}\ in\ v$ Template

where $n :: Int$, $a :: Arity$, $s :: IsShared$

Template syntax

$e ::=$ Atoms

$CON\ a\ n$	(Constructor tag)
$INT\ n$	(Primitive integers)
$PRI\ a\ \otimes$	(Primitive operation)
$FUN\ s\ a\ n$	(Function pointer)
$ARG\ s\ n$	(Argument pointer)
$VAR\ s\ n$	(Application pointer)
$REG\ n$	(Primitive register pointer)

$u, v ::=$ Applications

$APP\ \bar{e}$	(Normal application)
$CASE\ c\ \bar{e}$	(App with case table)
$PRIM\ n\ \bar{e}$	(PRS candidate)

$c ::= TAB\ n$ Case table pointer

$t ::= let\ \bar{u}\ in\ v$ Template

where $n :: Int$, $a :: Arity$, $s :: IsShared$

Template syntax

$e ::=$

Atoms	
$CON\ a\ n$	(Constructor tag)
$INT\ n$	(Primitive integers)
$PRI\ a\ \otimes$	(Primitive operation)
$FUN\ s\ a\ n$	(Function pointer)
$ARG\ s\ n$	(Argument pointer)
$VAR\ s\ n$	(Application pointer)
$REG\ n$	(Primitive register pointer)

$u, v ::=$

Applications	
$APP\ \bar{e}$	(Normal application)
$CASE\ c\ \bar{e}$	(App with case table)
$PRIM\ n\ \bar{e}$	(PRS candidate)
$c ::= TAB\ n$	Case table pointer
$t ::= let\ \bar{u}\ in\ v$	Template

where $n :: Int$, $a :: Arity$, $s :: IsShared$

Template bounds

```
let APP [ ARG True 0,          PRI 2 +,          INT 1 ]
        APP [ FUN True 2 0,          VAR False 0,          ARG False 1 ]
```

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

Template bounds

```
let APP [ ARG True 0,          PRI 2 +,          INT 1 ]  
      APP [ FUN True 2 0,        VAR False 0,       ARG False 1 ]
```

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

SpineLen (6)

Instantiate on stack

Template bounds

ApLen (4)

```
let APP [ ARG True 0, PRI 2 +, INT 1 ]  
      APP [ FUN True 2 0, VAR False 0, ARG False 1 ]
```

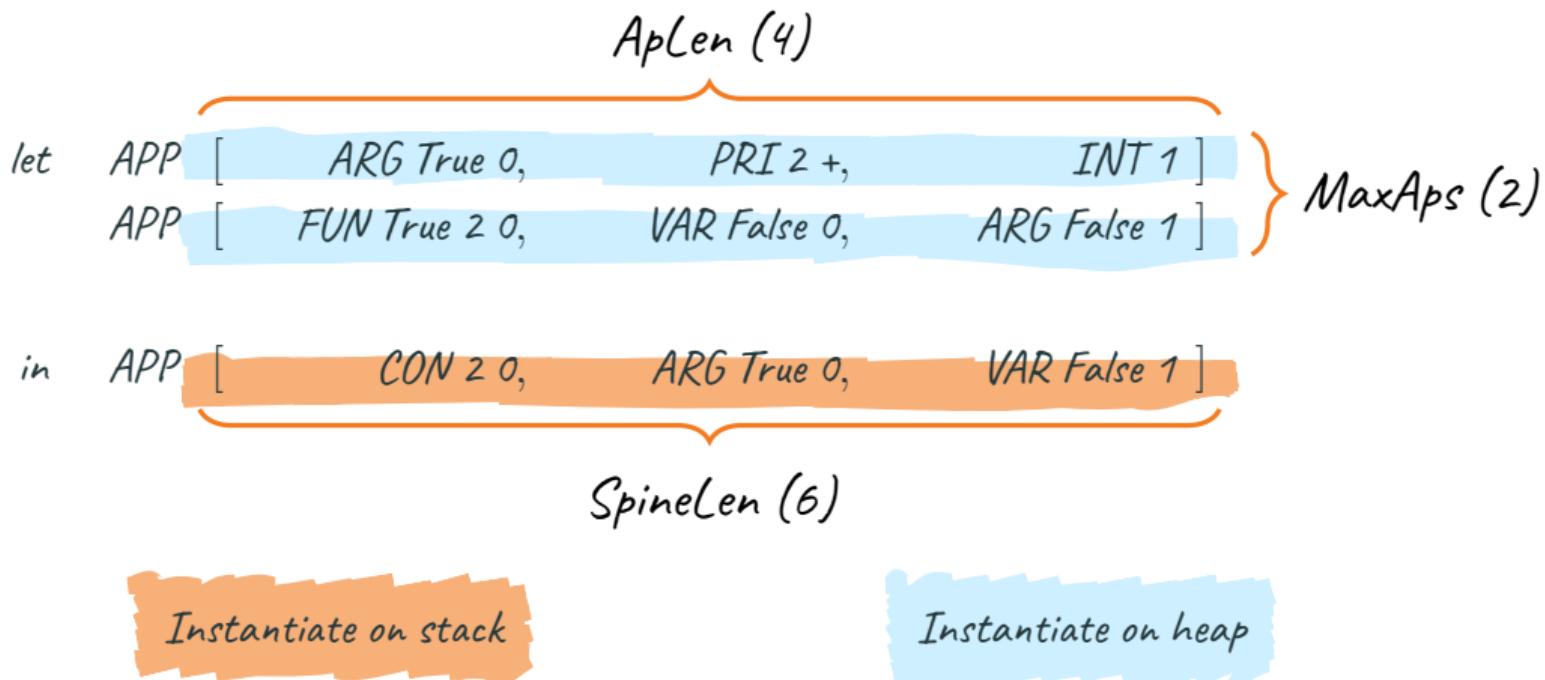
```
in APP [ CON 2 0, ARG True 0, VAR False 1 ]
```

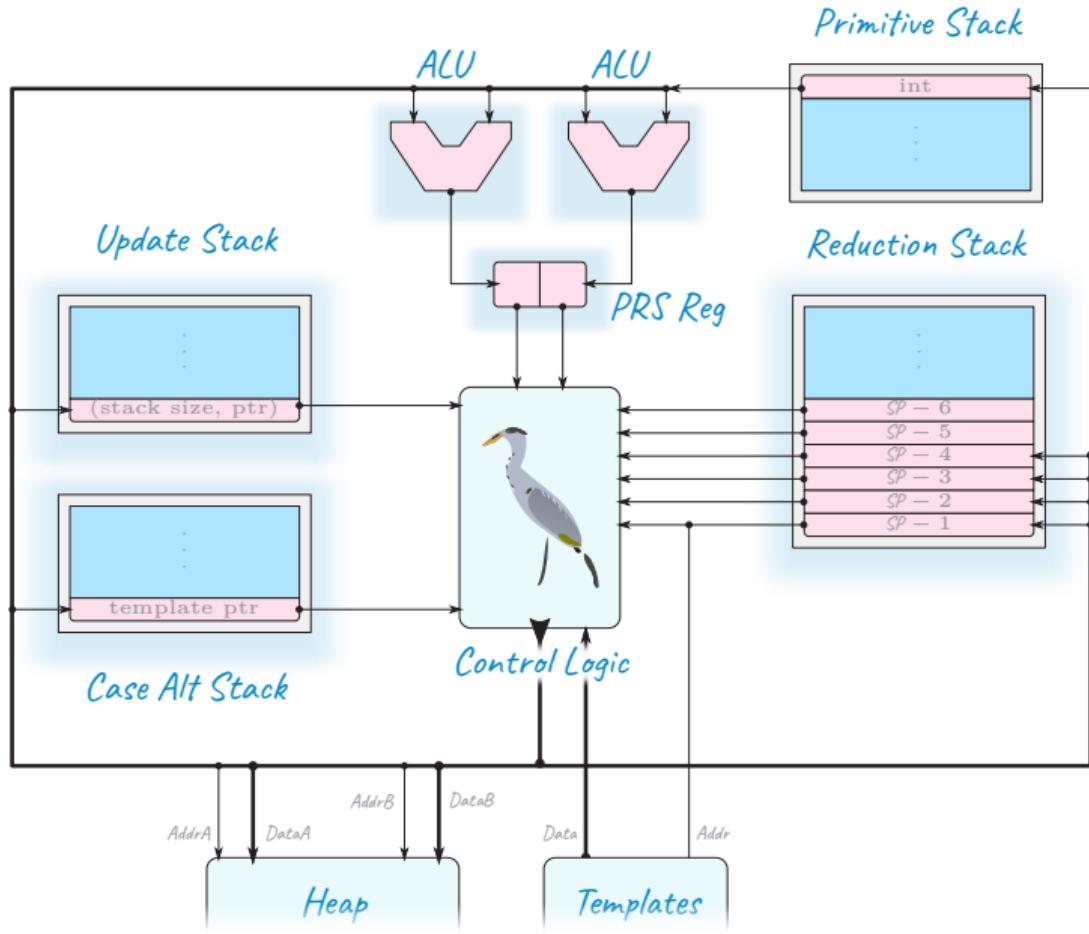
SpineLen (6)

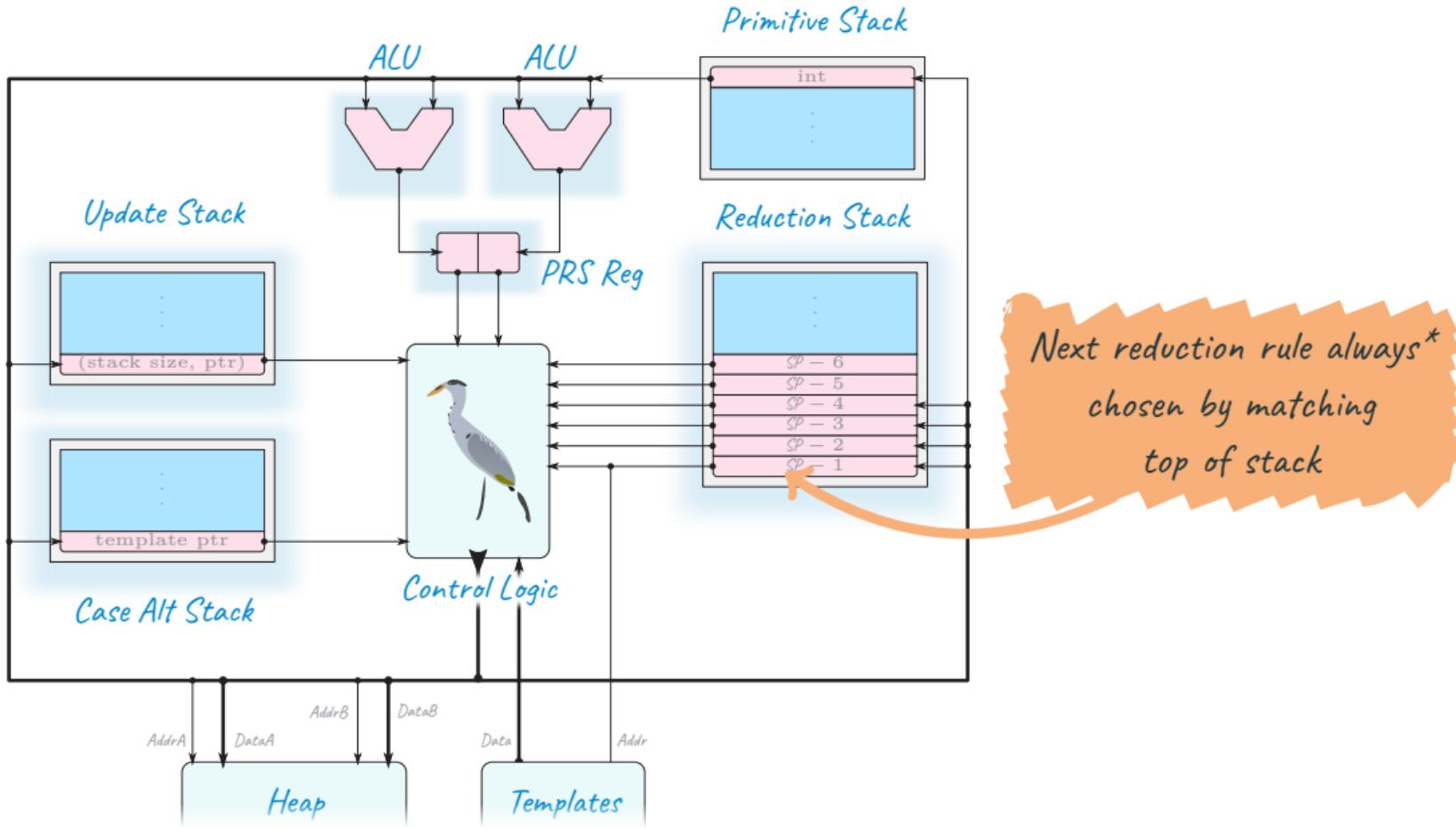
Instantiate on stack

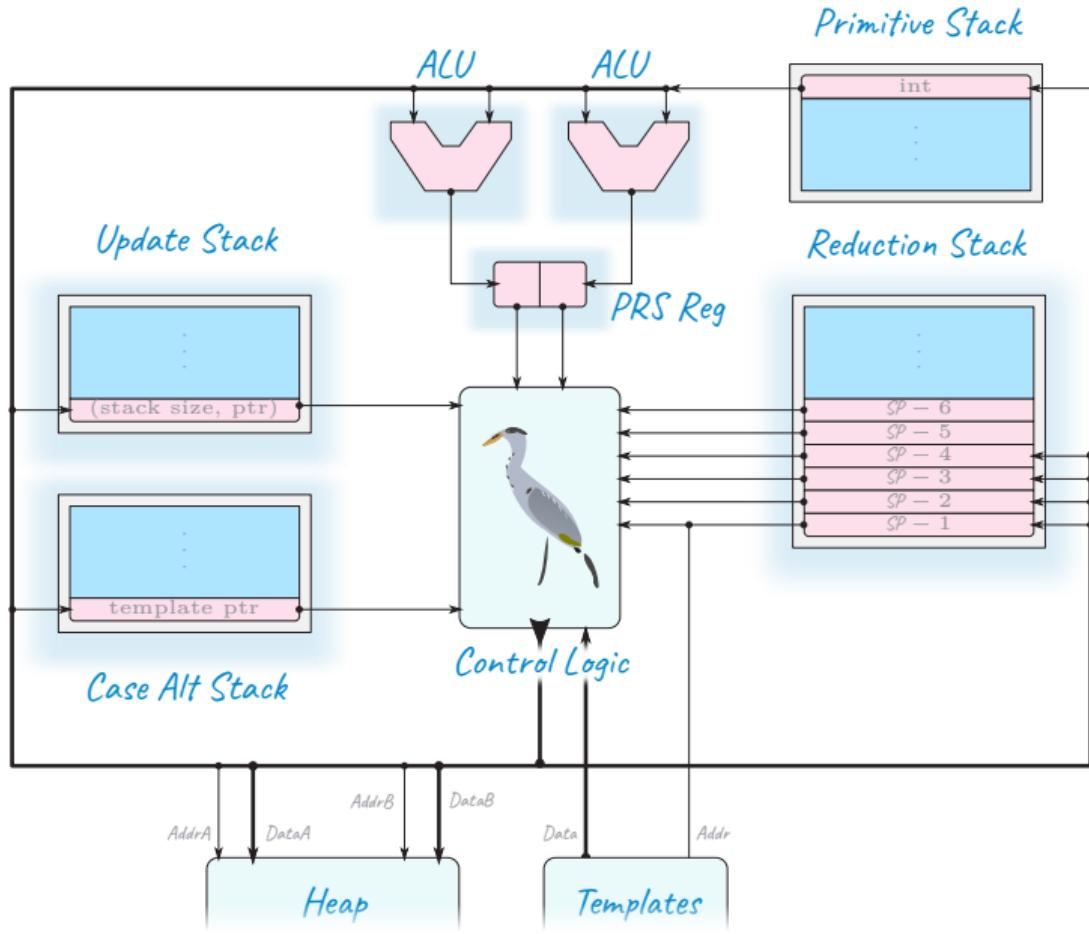
Instantiate on heap

Template bounds



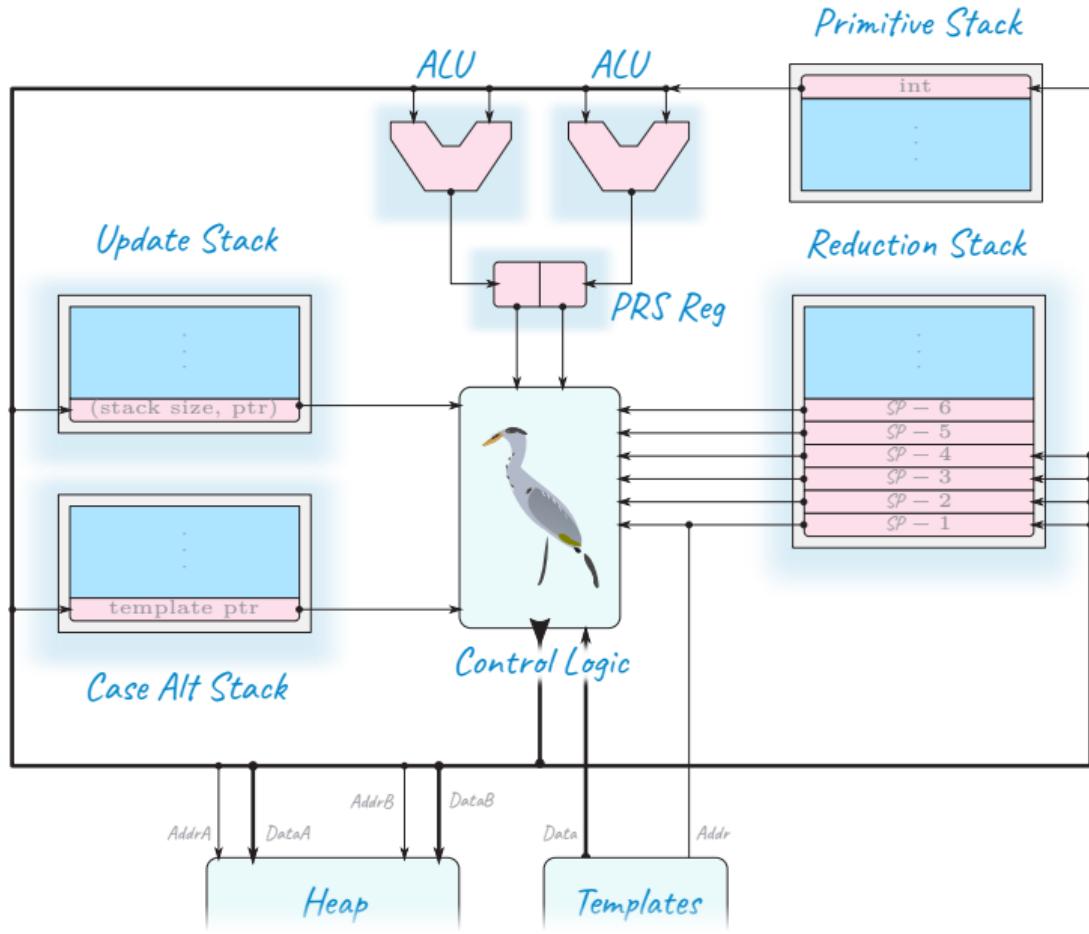






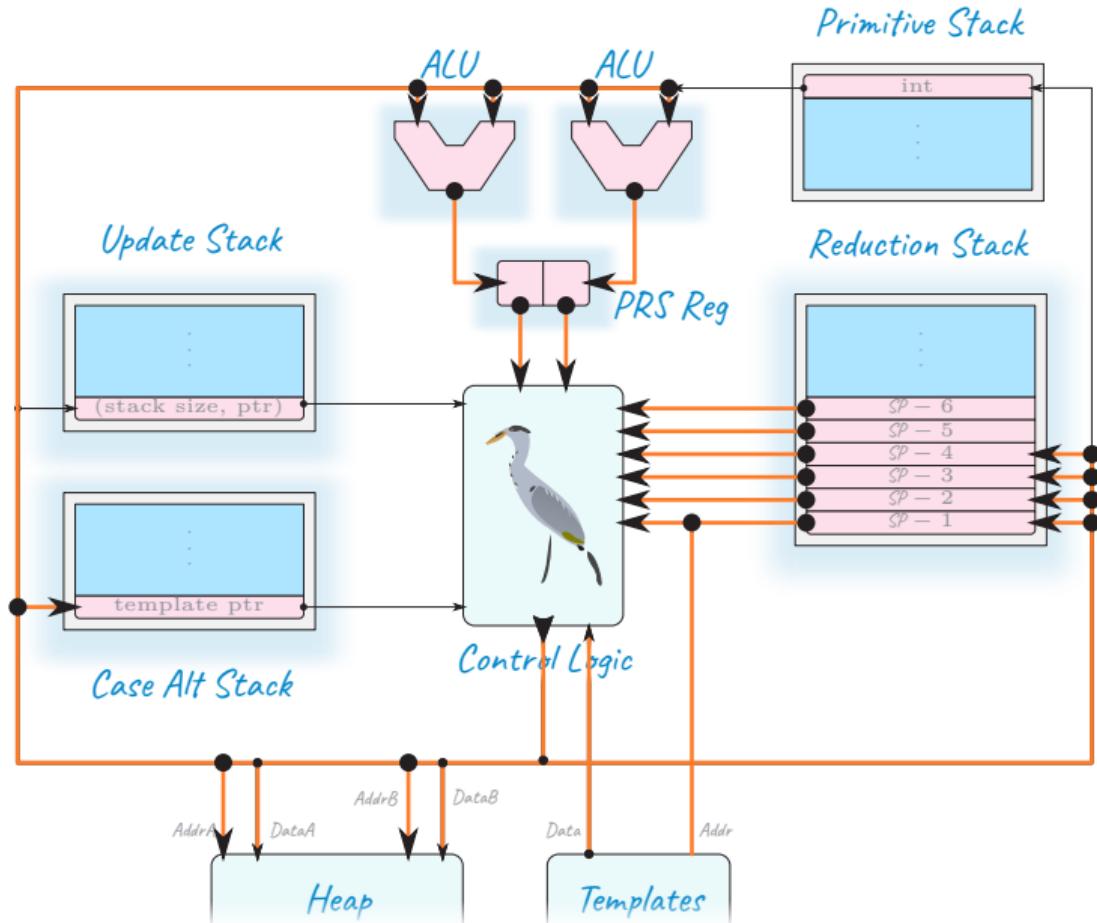
Atoms

- $= \text{FUN } s \text{ } n$
- $| \text{CON } s \text{ } n$
- $| \text{VAR } s \text{ } n$
- $| \text{INT } n$
- $| \text{PRI } a \otimes$
- $| \text{ARG } s \text{ } n$
- $| \text{REG } n$



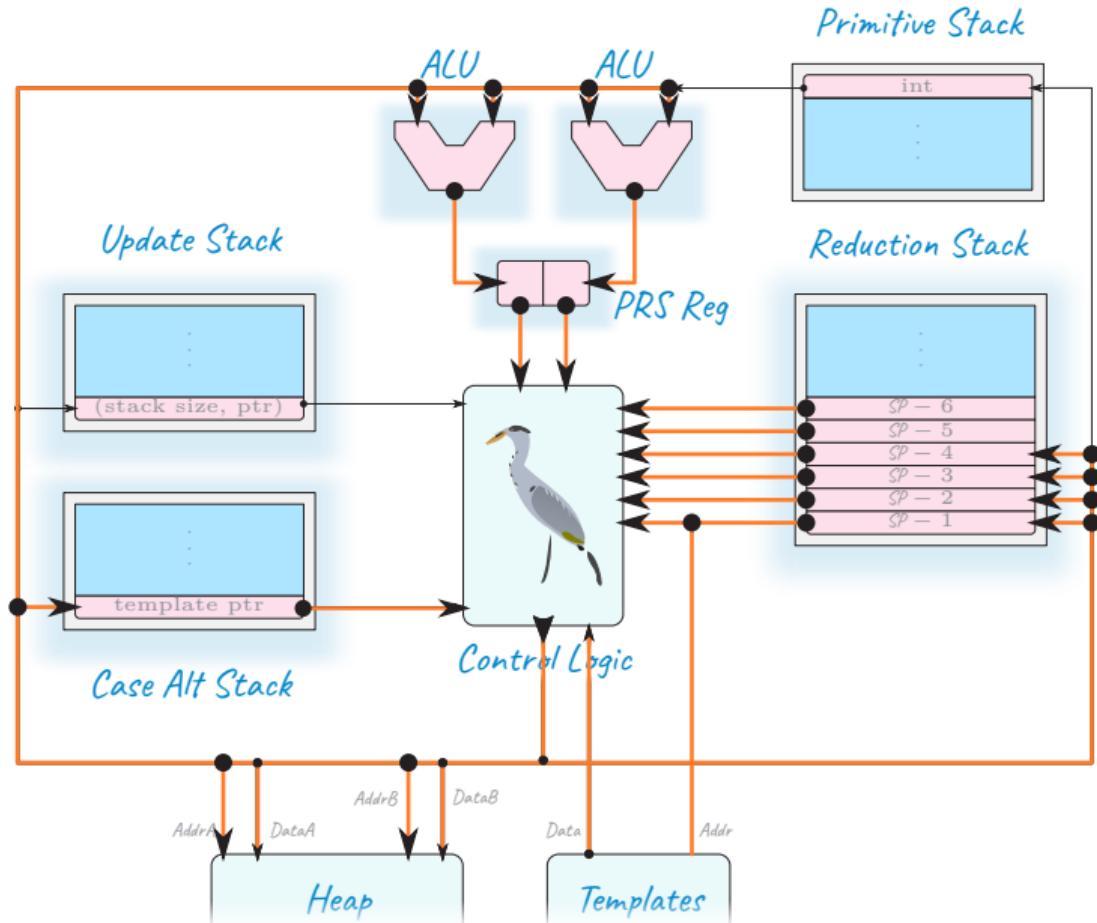
Atoms

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



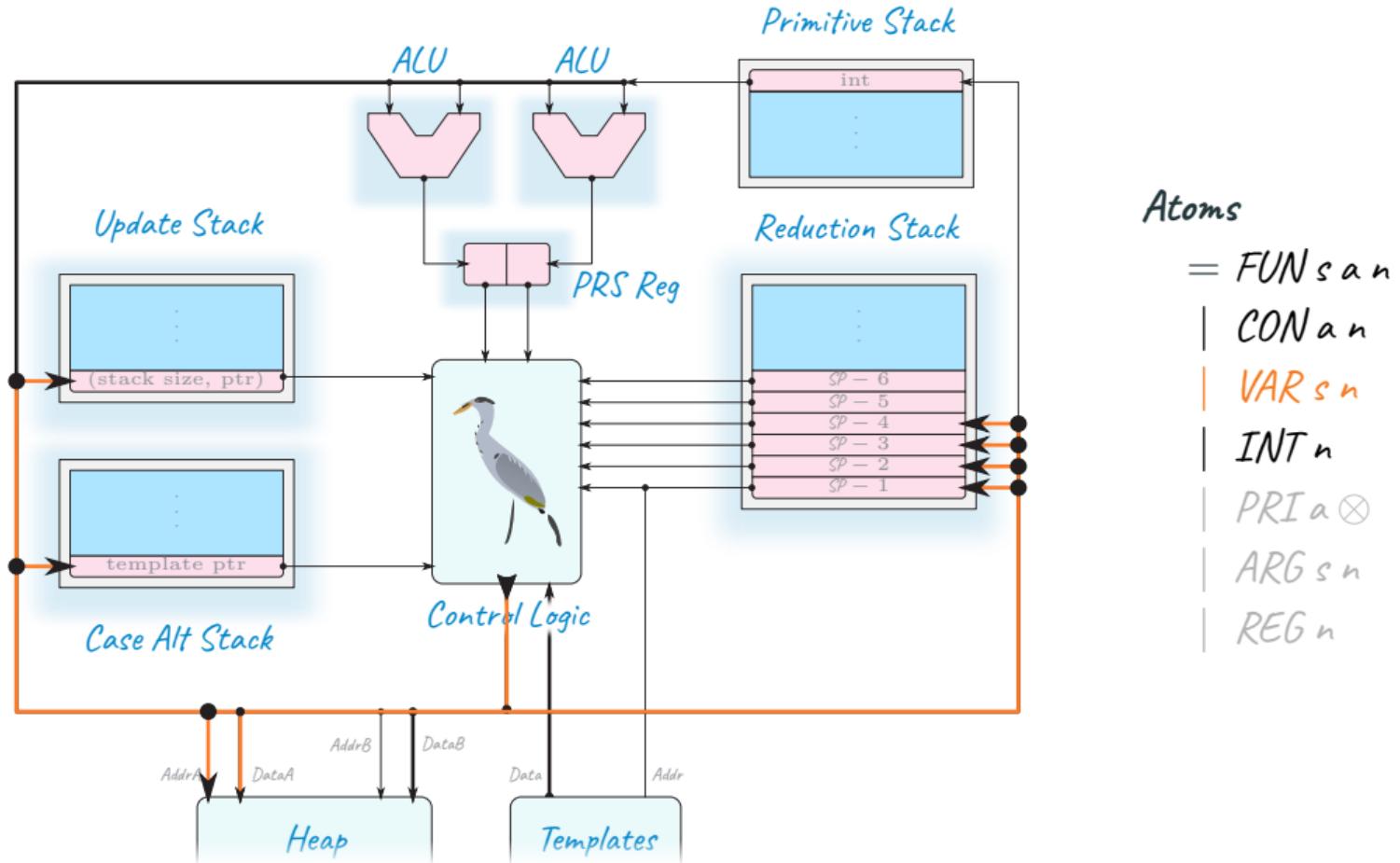
Atoms

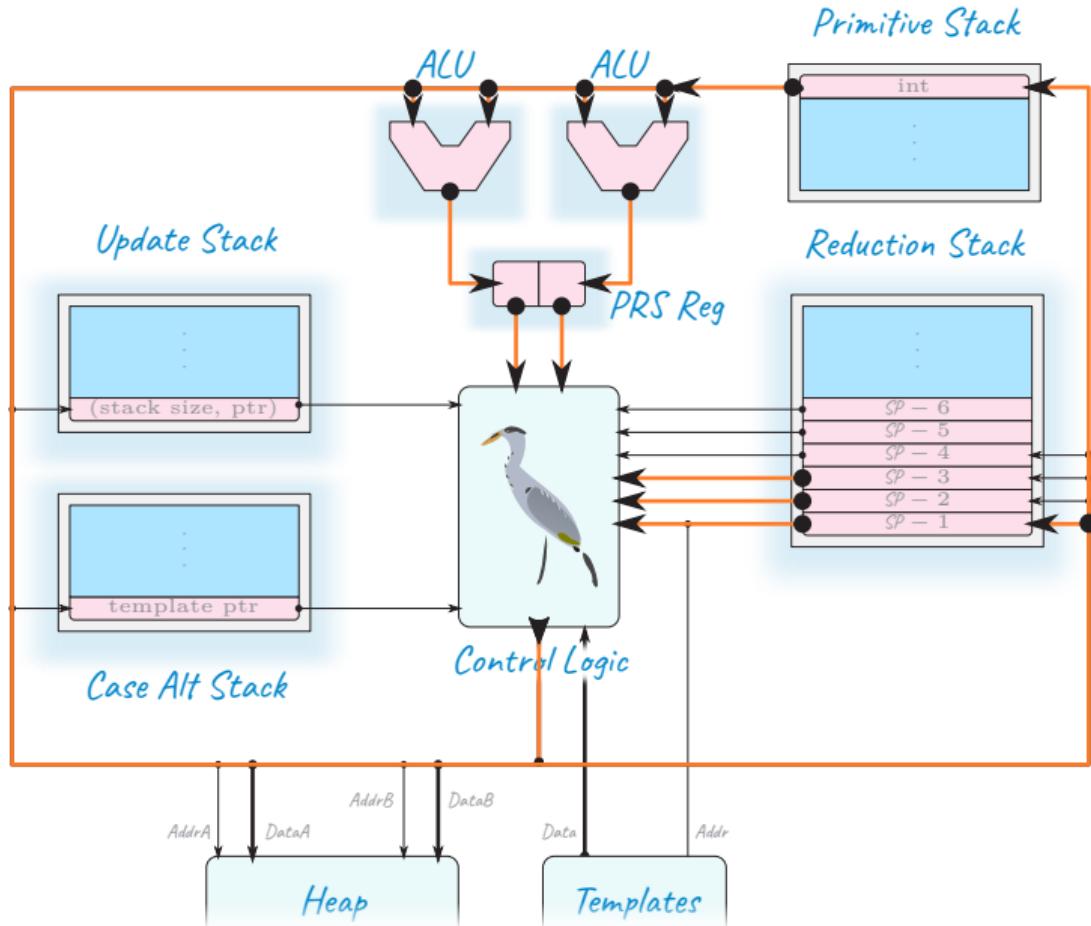
- = `FUN s n`
- | `CON s n`
- | `VAR s n`
- | `INT n`
- | `PRI a ⊗`
- | `ARG s n`
- | `REG n`



Atoms

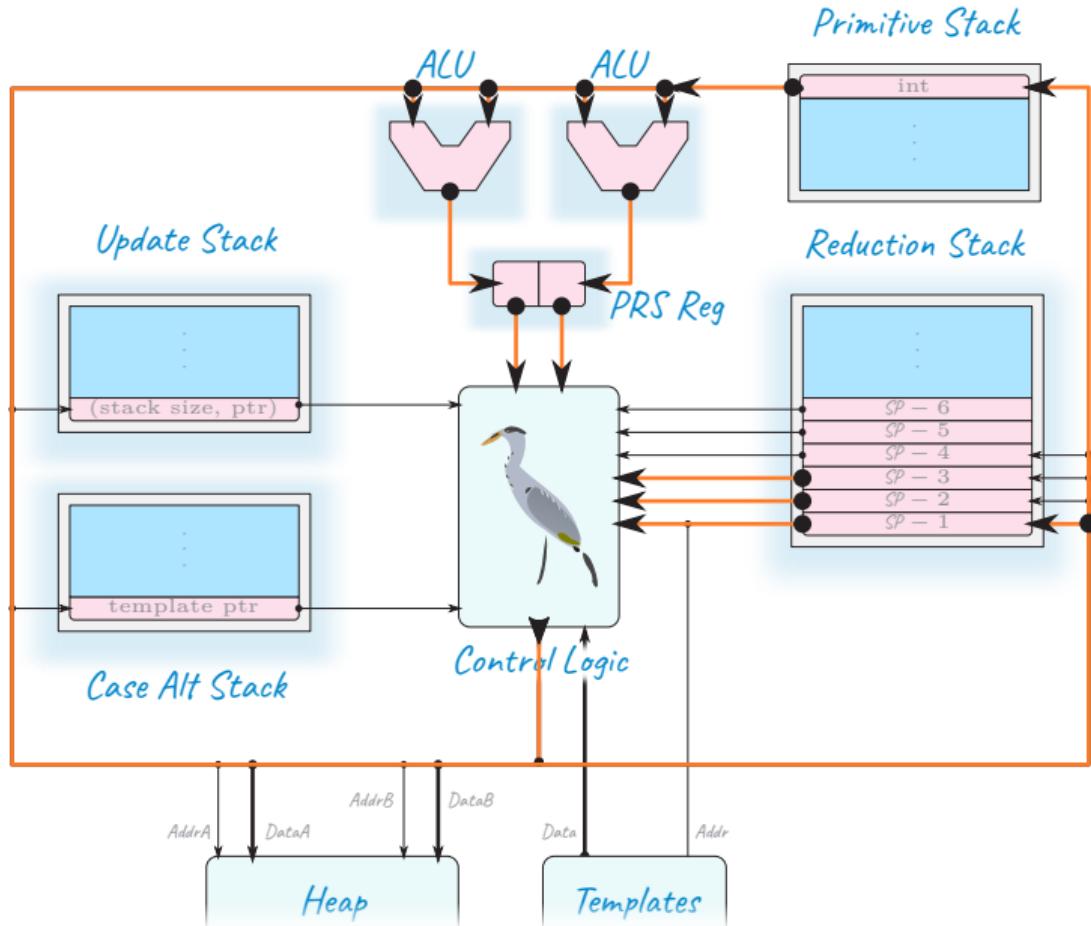
- = **FUN s n**
- | **CON a n**
- | **VAR s n**
- | **INT n**
- | **PRI a ⊗**
- | **ARG s n**
- | **REG n**





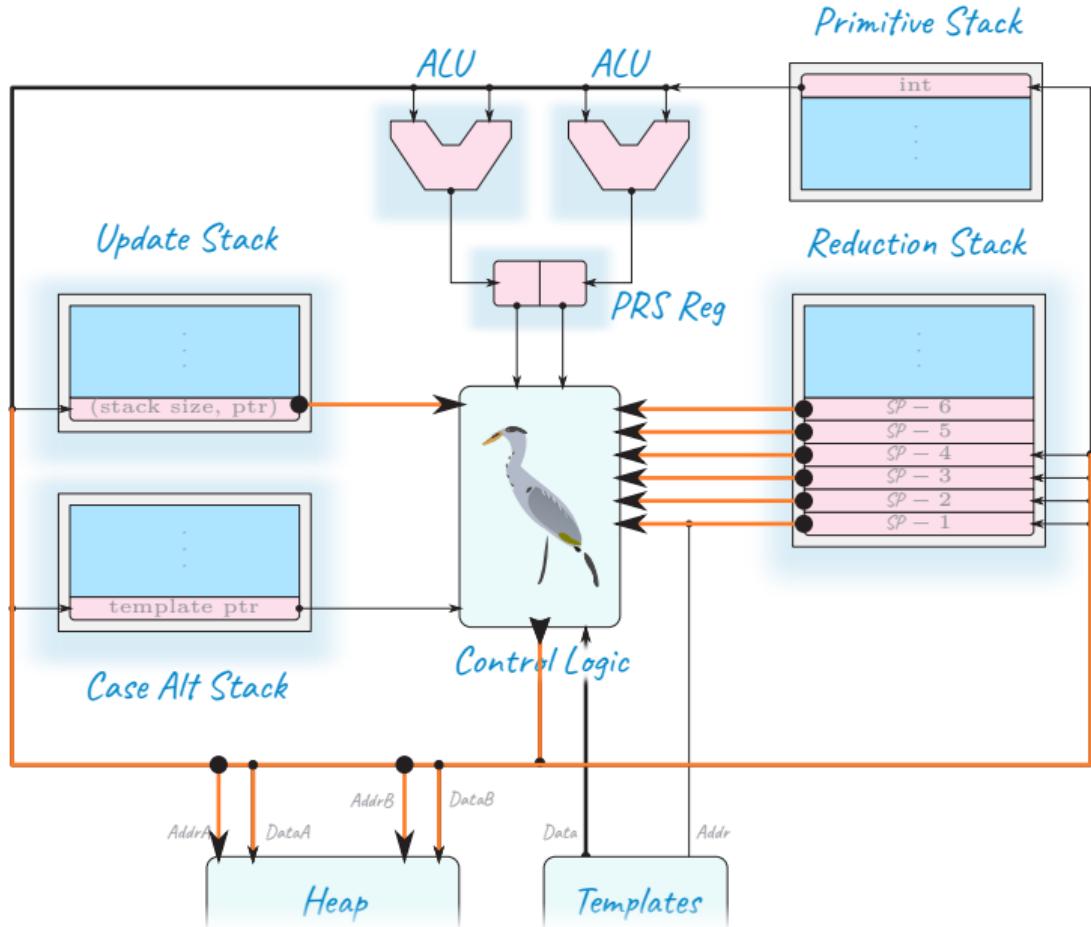
Atoms

- = FUN_{s n}
- | CON_{a n}
- | VAR_{s n}
- | INT_n
- | PRI_{a ⊗}
- | ARG_{s n}
- | REG_n

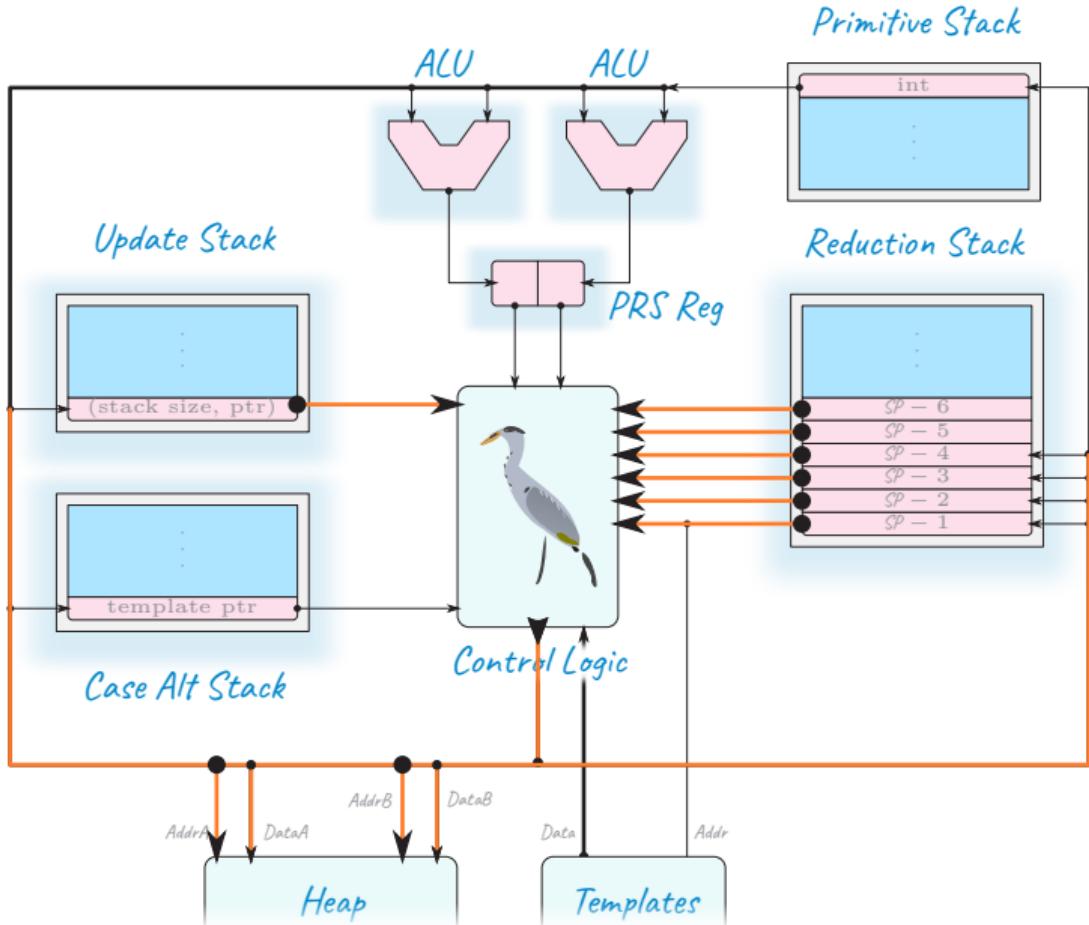


Postfix prims for long spines

$$(fx\ y) + (gz) \\ \Rightarrow f\ x\ y\ g\ z\ +$$



...But what about
heap updates?



Avoid most updates via run-time sharing analysis!

Atoms

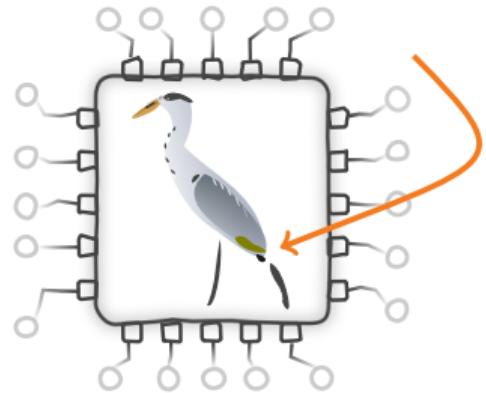
- = **FUN**_{s n}
- | **CON**_n
- | **VAR**_{s n}
- | **INT**_n
- | **PRI**_a ⊗
- | **ARG**_{s n}
- | **REG**_n

Cloaca

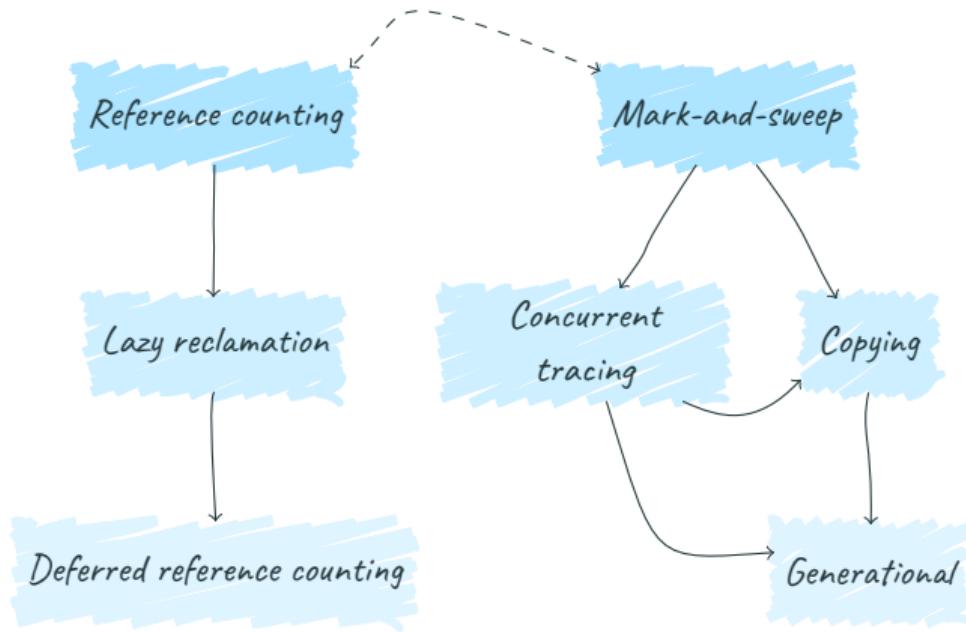
noun [C]

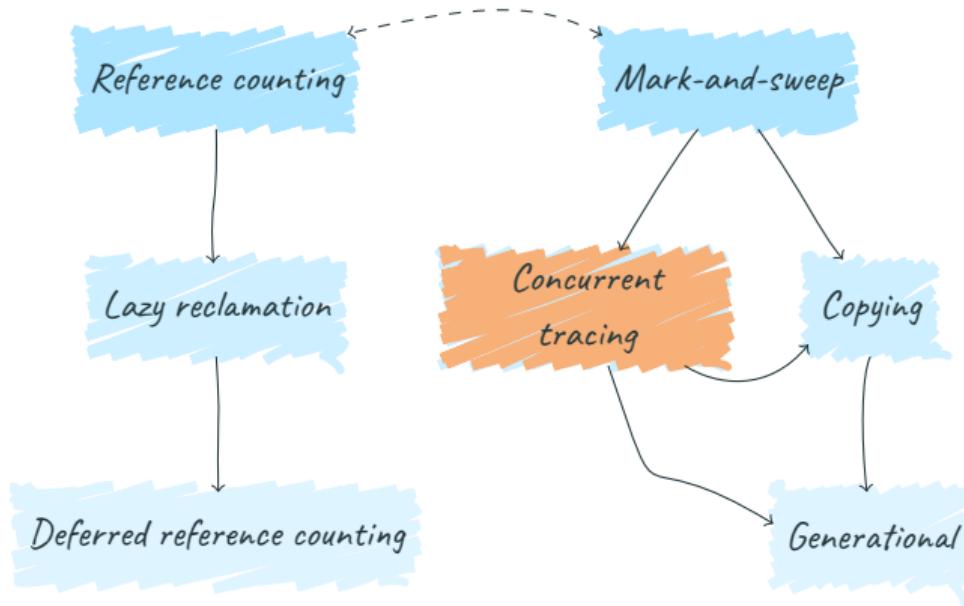
/kloh-ah-kuh/

A concurrent hardware
garbage collector for Heron

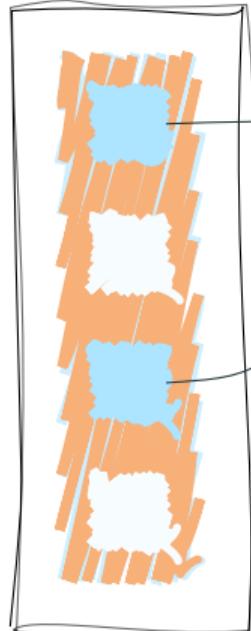


^oRamsay and Stewart, "Cloaca: A Concurrent Hardware Garbage Collector for Non-strict Functional Languages".

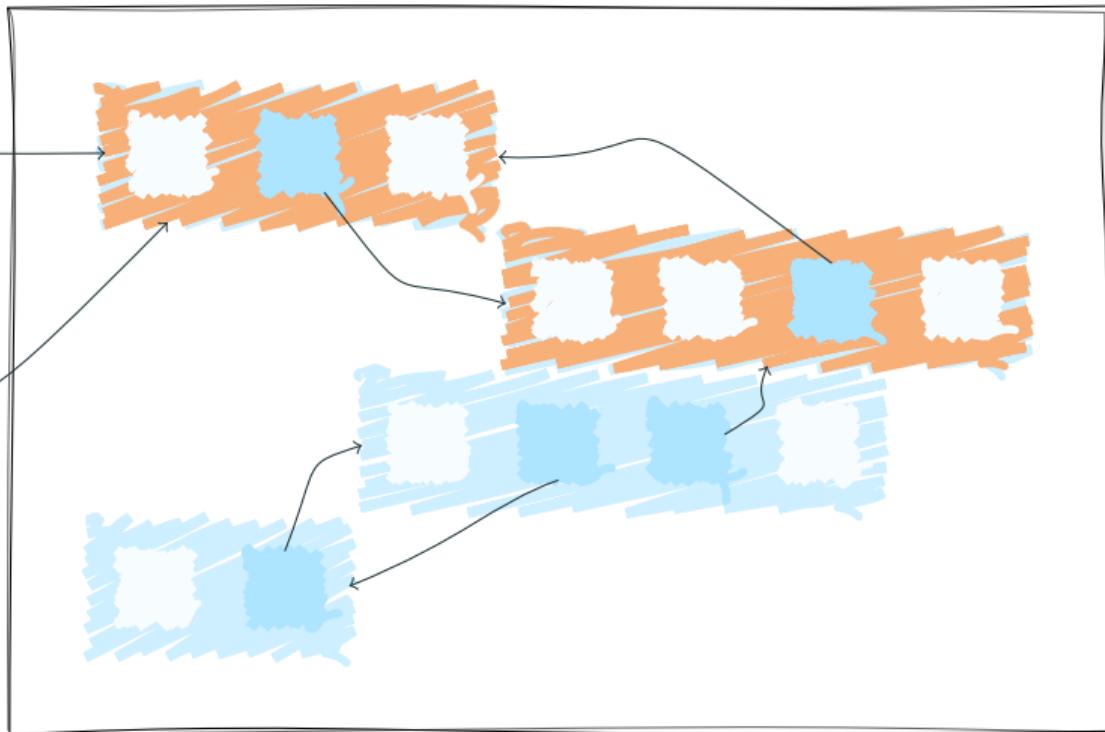




*Stack
(graph roots)*

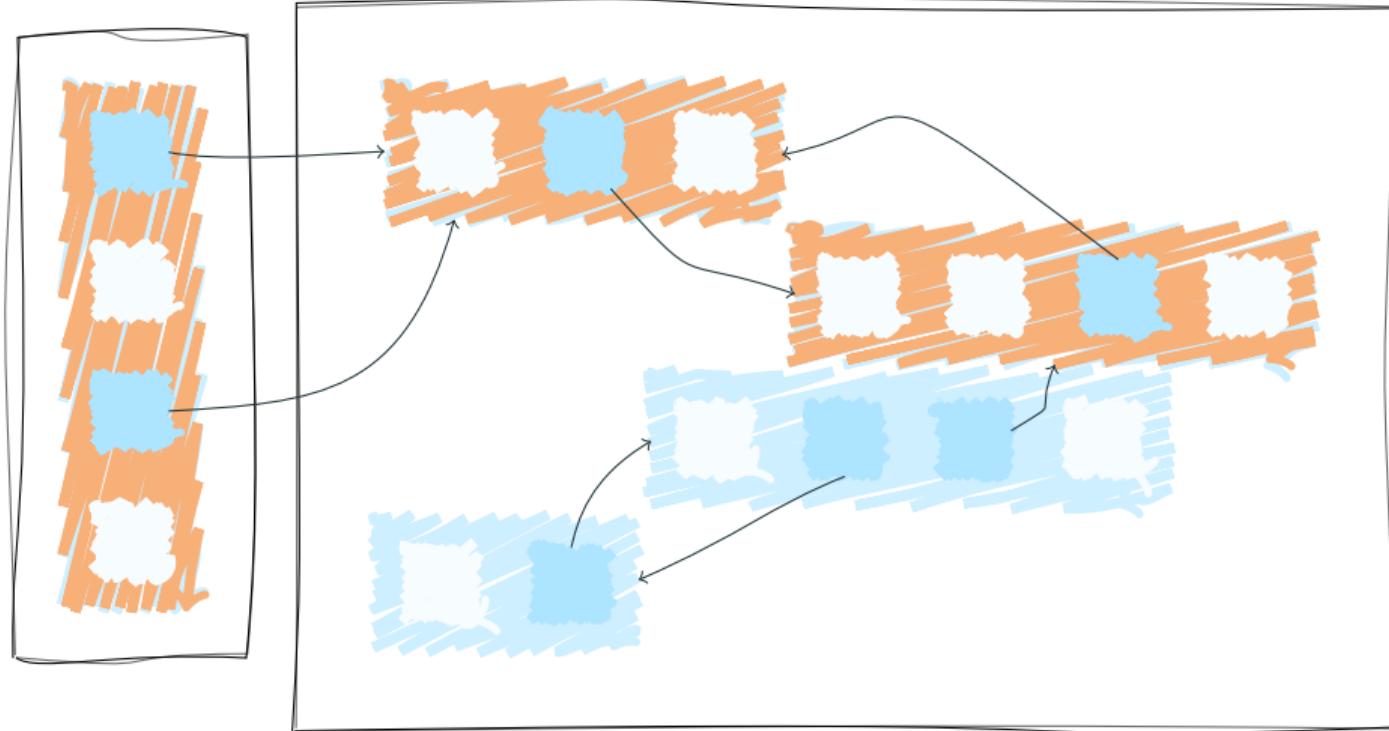


Heap



*Stack
(graph roots)*

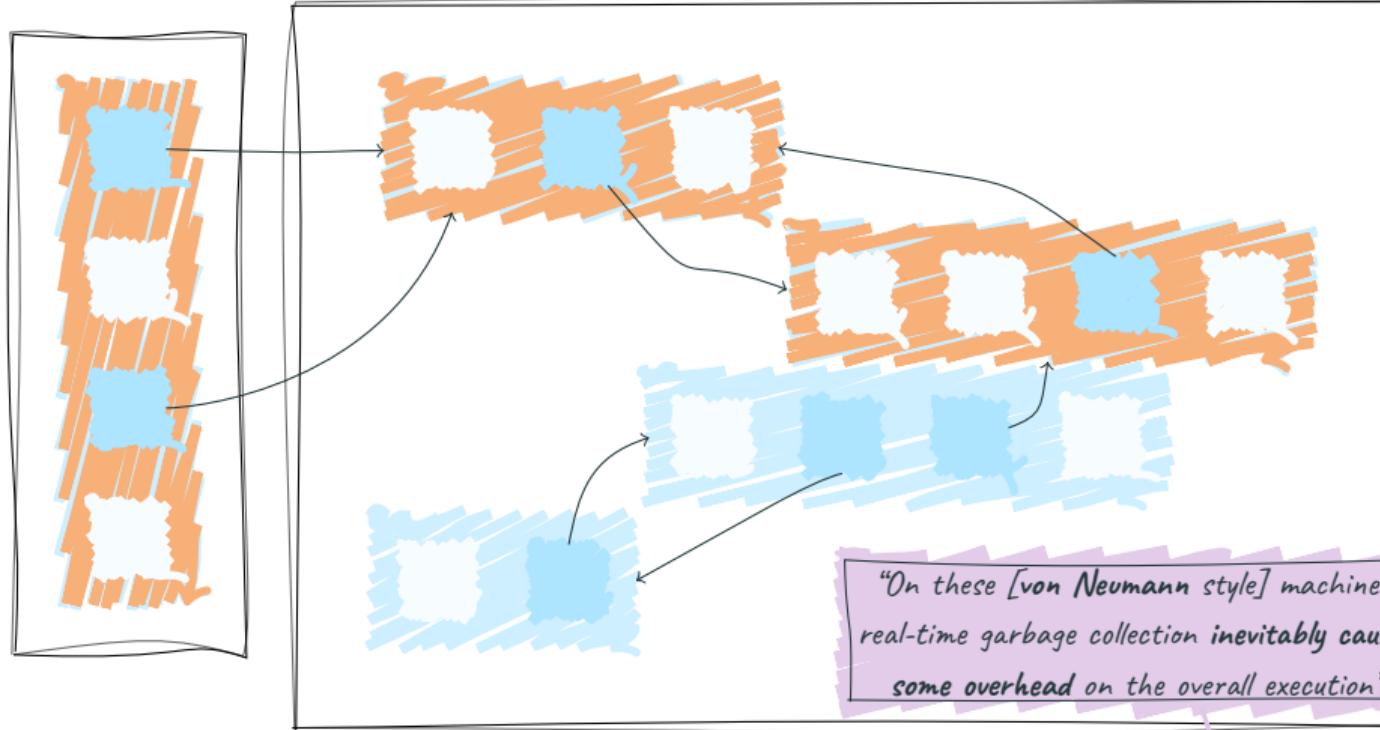
Heap



^o Yuasa, “Real-time garbage collection on general-purpose machines”.

Stack
(graph roots)

Heap



° Yuasa, "Real-time garbage collection on general-purpose machines".

Challenges for concurrent software implementation?

Expect 21% median slowdown for nofib!

1) Allocation depends on GC state

Stop-the-world GC

vs

Concurrent GC

Function alloc (app):

```
    heap[hp] ← app  
    hp++
```

Function alloc (app):

```
if allocBarrier(gcPhase, hp)  
then  
    tag hp as Marked  
else  
    tag hp as Unmarked  
    heap[hp] ← app  
    hp++
```

2) Non-moving GC needs complex allocation

Stop-the-world GC

vs

Concurrent GC

Function alloc (app):

```
    heap[hp] ← app  
    hp++
```

Function alloc (app):

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

3) Prevent graph updates from destroying edges

Stop-the-world GC

vs

Concurrent GC

Function update (nf, a):

```
└ heap[a] ← nf
```

Function update (nf, a):

```
if updateBarrier(gcPhase) then
    x ← heap[a]
    forall y in x's child pointers do
        remember y for marking
    heap[a] ← nf
```

Additional hardware-enabled optimisations

Heron's existing dynamic update avoidance system...

data Atom

= ...

/ Var Shared Int
/ Arg Shared Int

...

Function unwind (a, shared):

...

...

...

if shared and not NF then

 └ push a onto update stack

Heron's existing dynamic update avoidance system...

data Atom

= ...

/ Var Shared Int

/ Arg Shared Int

...

Function unwind (a, shared):

...

...

...

if shared and not NF then

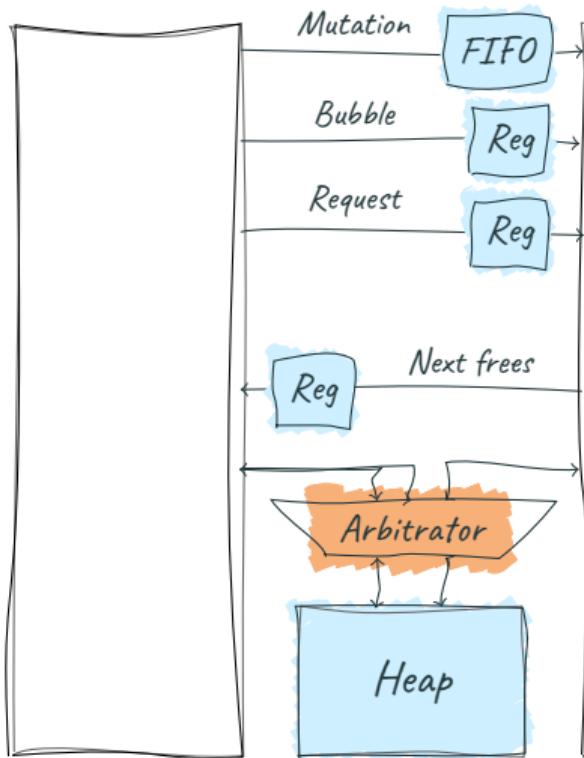
 └ push a onto update stack

if not shared then

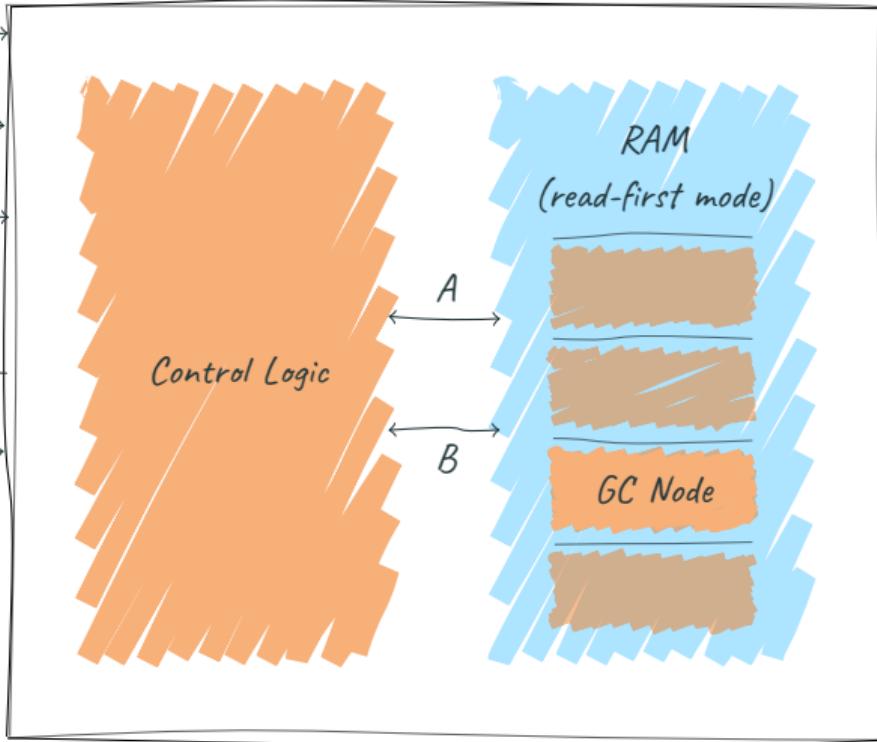
 └ dealloc a

... is just one-bit reference counting with a hat on.

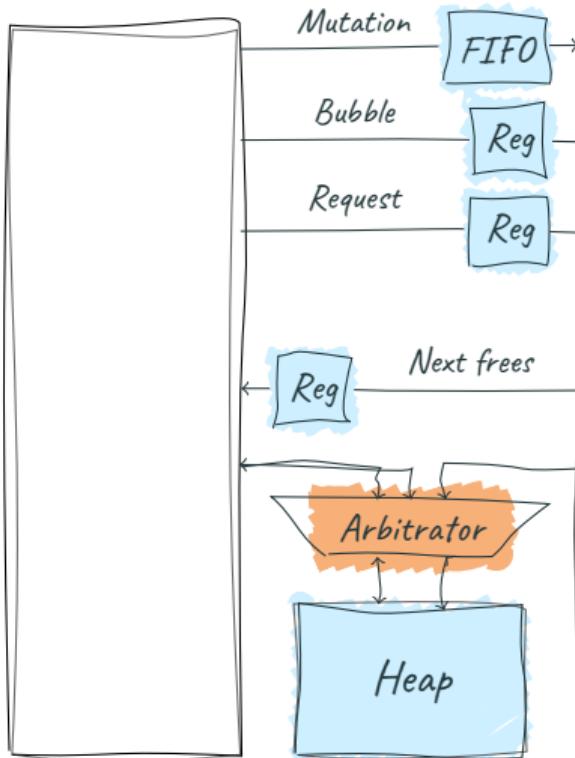
Reduction Core



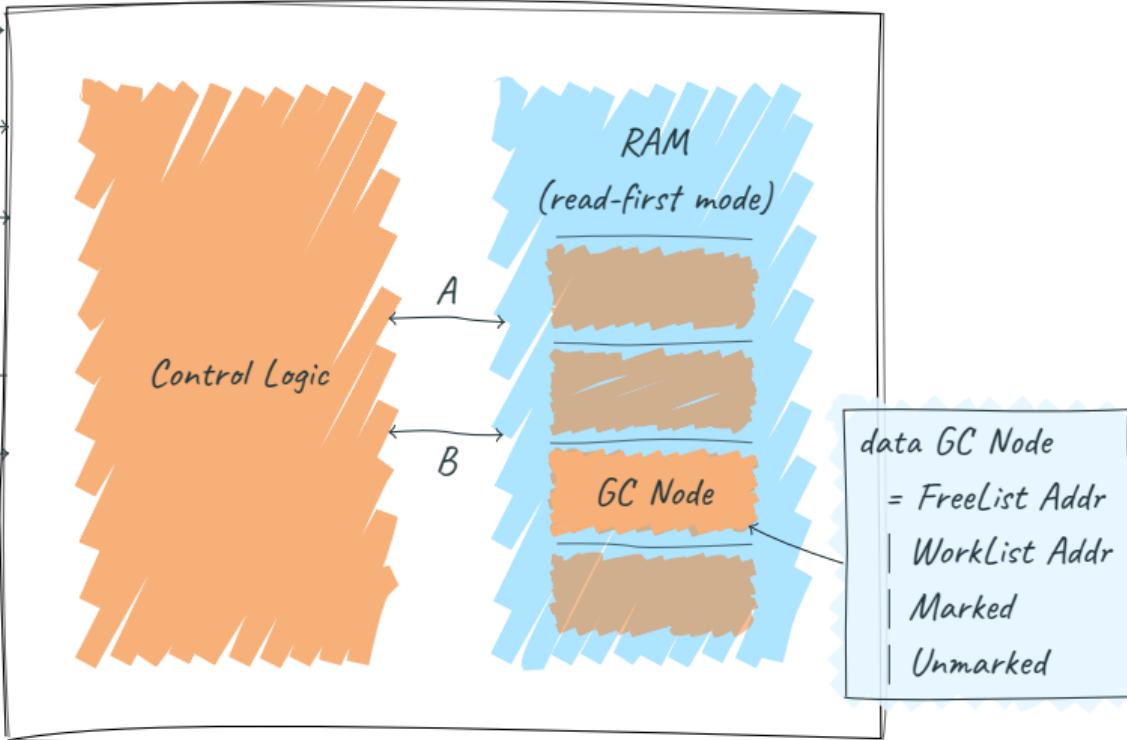
Memory Management



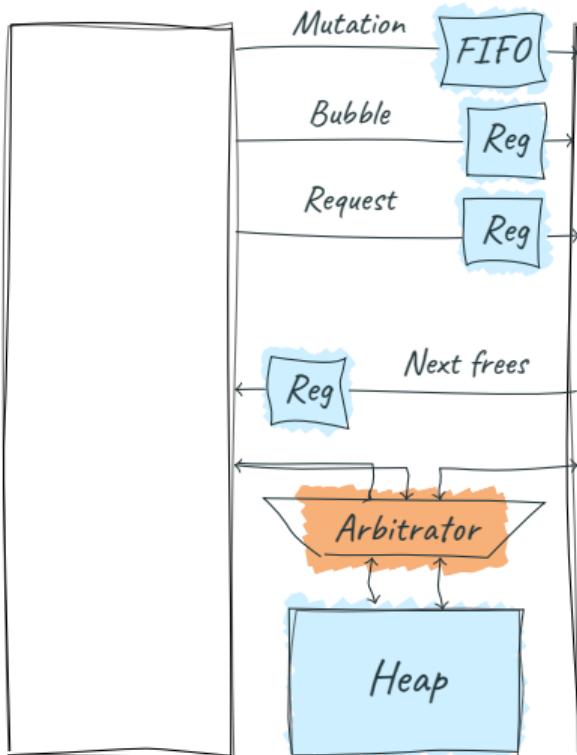
Reduction Core



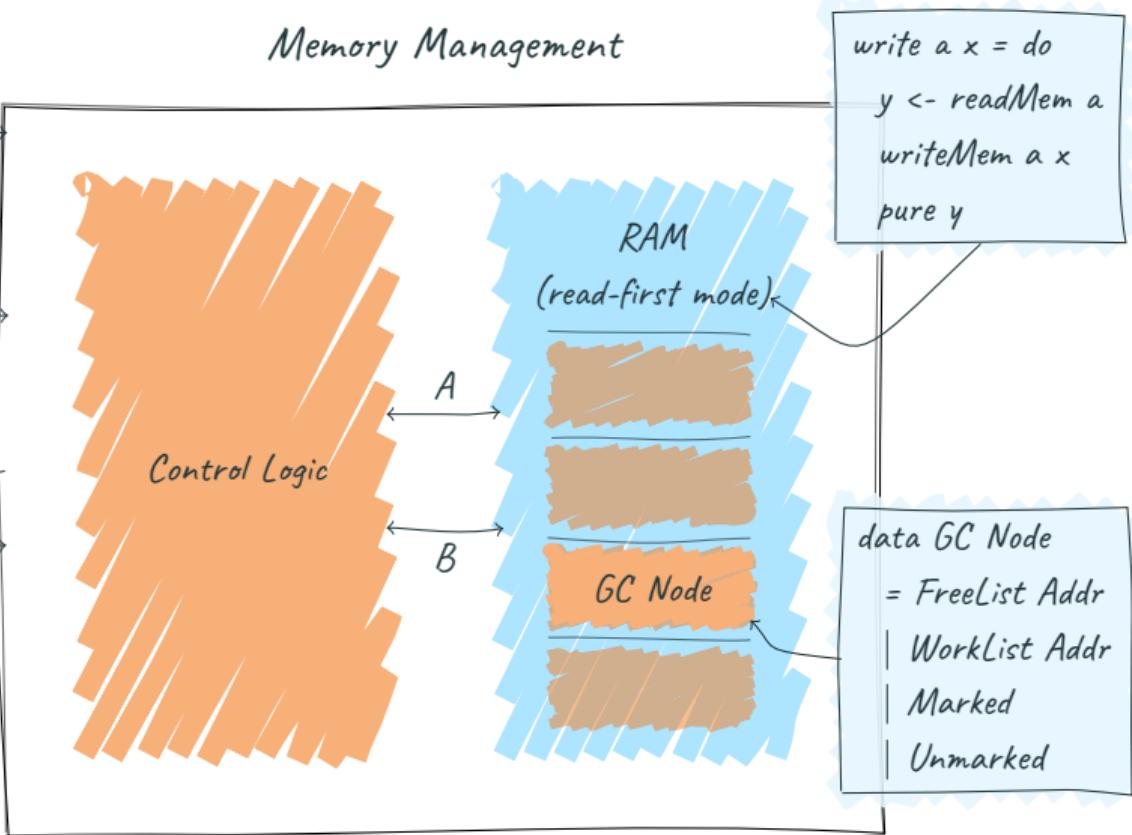
Memory Management



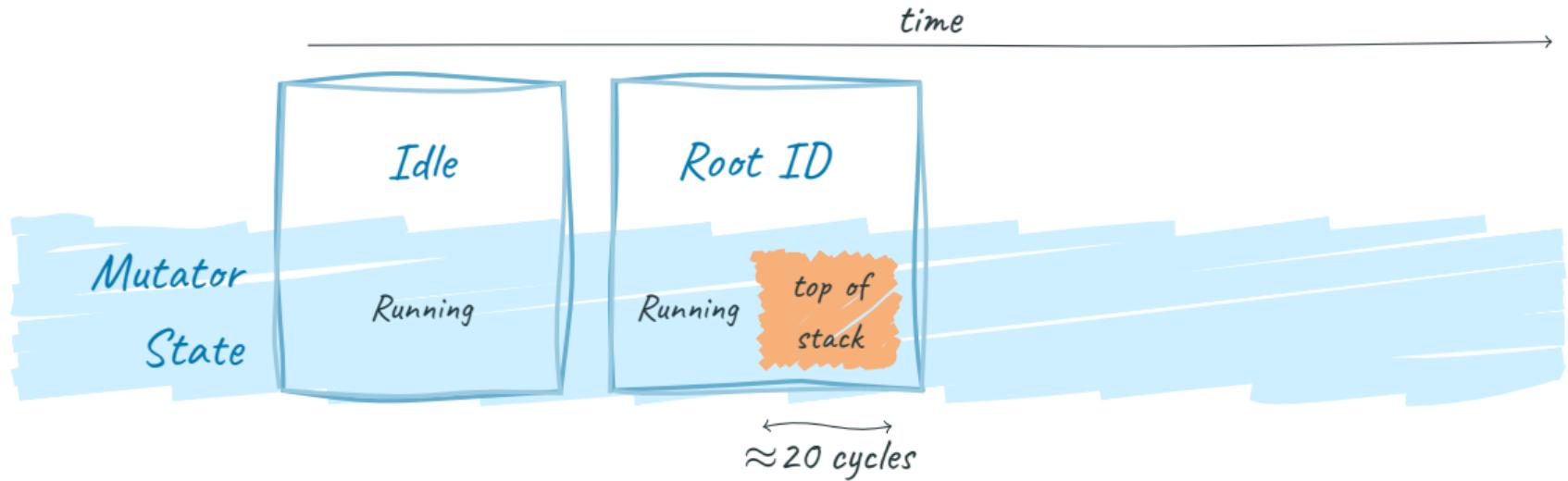
Reduction Core

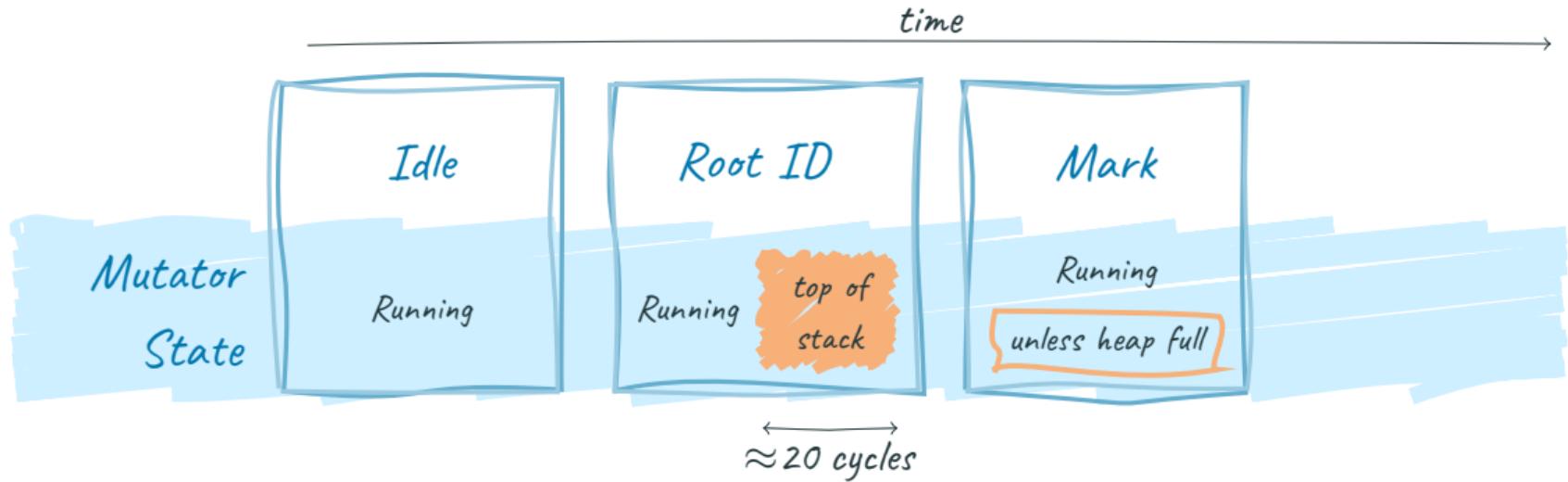


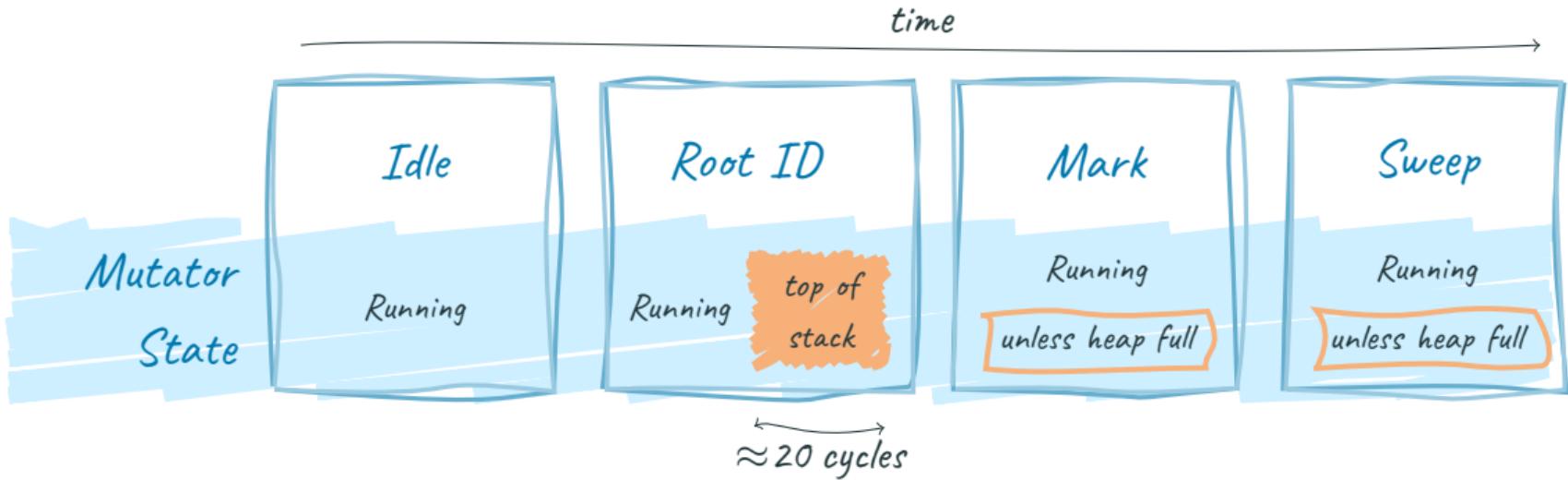
Memory Management











Performance

	Peak GHC working set (KB)	GHC Allocations (MB)	LoC
<i>Adjoxo</i>	47	301	72
<i>Braun</i>	46	0	26
<i>Cichelli</i>	52	41	123
<i>Clausify</i>	77	364	69
<i>Countdown</i>	46	54	62
<i>Knuthbendix</i>	105	54	324
<i>Mate</i>	137	430	293
<i>Mss</i>	46	359	13
<i>Ordlist</i>	46	984	28
<i>Permsort</i>	46	2320	10
<i>Queens</i>	55	1038	25
<i>Queens2</i>	47	1084	20
<i>Sumpuz</i>	48	1293	72
<i>Taut</i>	47	236	37
<i>While</i>	46	264	89

	Peak GHC working set (KB)	GHC Allocations (MB)	LoC
Adjoxo	47	301	72
Braun	46	0	26
Cichelli	52	41	123
Clausify	77	364	69
Countdown	46	54	62
Knuthbendix	105	54	324
Mate	137	430	293
Mss	46	359	13
Ordlist	46	984	28
Permsort	46	2320	10
Queens	55	1038	25
Queens2	47	1084	20
Sumpuz	48	1293	72
Taut	47	236	37
While	46	264	89

	Peak GHC working set (KB)	GHC Allocations (MB)	LoC
Adjoxo	47	301	72
Braun	46	0	26
Cichelli	52	41	123
Clausify	77	364	69
Countdown	46	54	62
Knuthbendix	105	54	324
Mate	137	430	293
Mss	46	359	13
Ordlist	46	984	28
Permsort	46	2320	10
Queens	55	1038	25
Queens2	47	1084	20
Sumpuz	48	1293	72
Taut	47	236	37
While	46	264	89

Which Architectures?

Platform	Heron Xilinx Alveo U280	GHC + Intel i7 1250U Performance	Power-saver
----------	----------------------------	-------------------------------------	-------------

Which Architectures?

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

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 ¹	Cores 15 W or Package 16 W	Cores 2W or Package 6 W

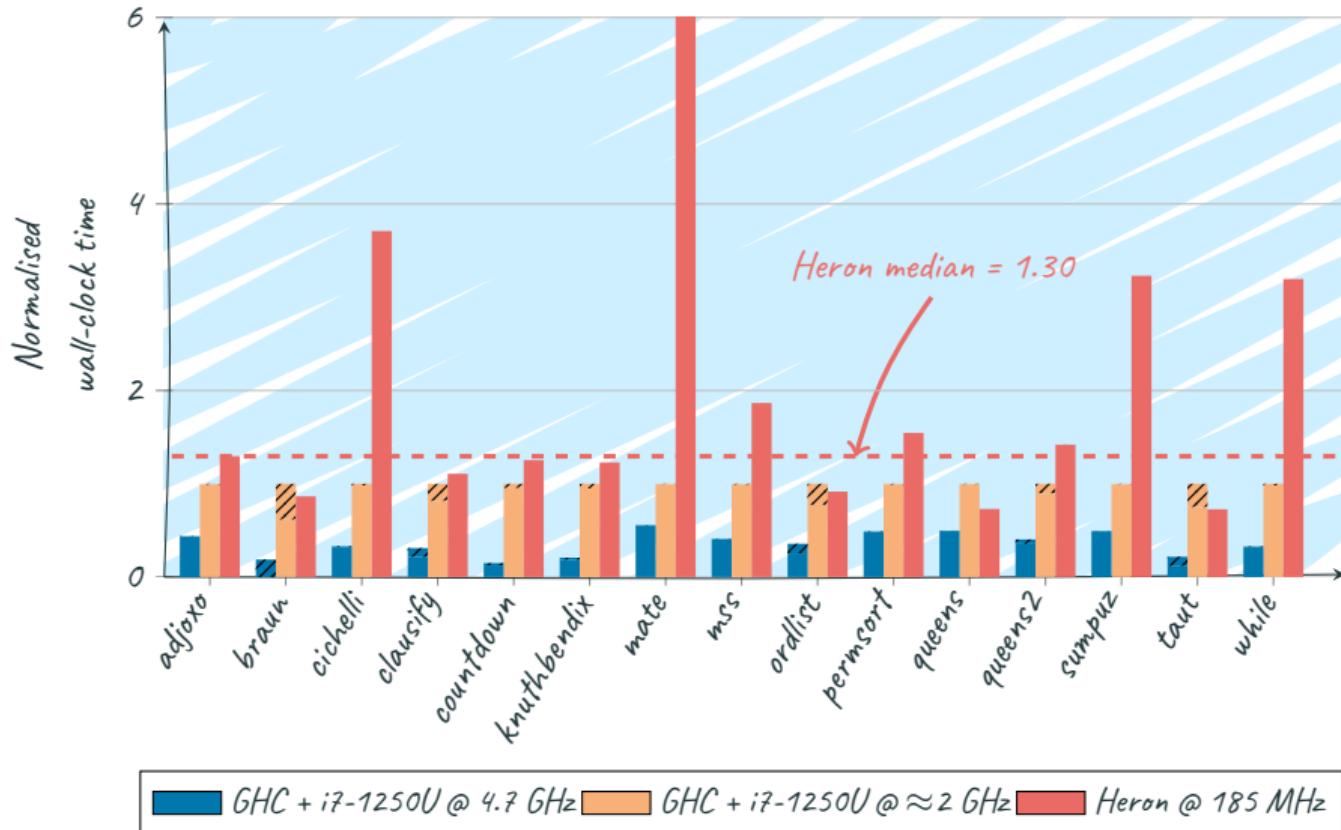
¹ Heron only occupies 1.13% of any resource type though!

Which Architectures?

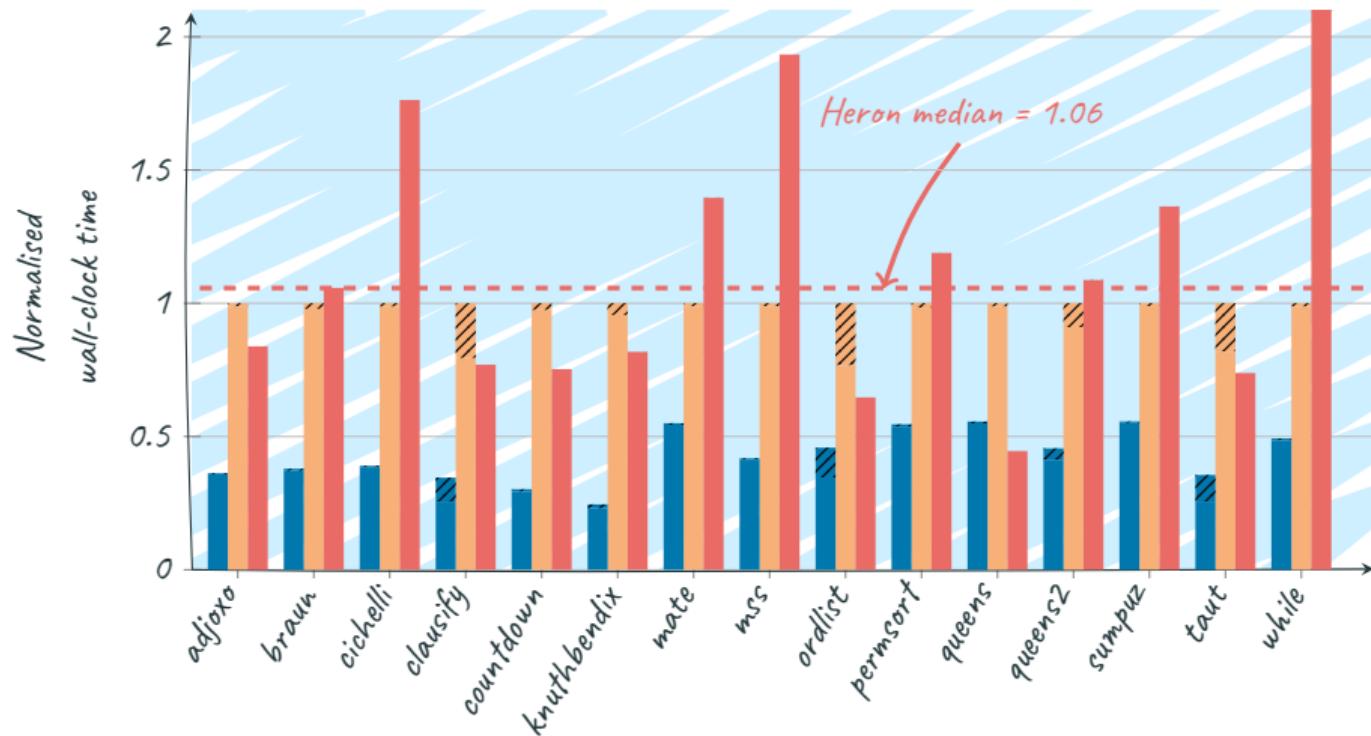
	Heron	GHC + Intel i7 1250U	
Platform	Xilinx Alveo U280	Performance	Power-saver
Clock	185 MHz	4.7 GHz	≈ 2 GHz
Power est.	0.8W dynamic + 3.1 W Static ¹	Cores 15 W or Package 16 W	Cores 2W or Package 6 W
Fabrication	16 nm (FPGA!)	10 nm	10 nm

¹ Heron only occupies 1.13% of any resource type though!

Wall-clock times vs GHC -O2



Wall-clock times vs GHC -00



[█] GHC + i7-1250U @ 4.7 GHz [█] GHC + i7-1250U @ \approx 2 GHz [█] Heron @ 185 MHz

Reflection

Should we widen the local memory bottleneck at all costs?

Should we widen the local memory bottleneck at all costs?

Heron's single-thread performance is competitive with general purpose CPUs!

Should we widen the local memory bottleneck at all costs?

Heron's single-thread performance is competitive with general purpose CPUs!

We may have optimised single-thread
performance to the point of limiting a parallel system²...

²Context switching is expensive when thread state is leaks out into many memories

Should we widen the local memory bottleneck at all costs?

Heron's single-thread performance is competitive with general purpose CPUs!

We may have optimised single-thread
performance to the point of limiting a parallel system²...

Cephalopode and Heron at the two extremes here.

²Context switching is expensive when thread state is leaks out into many memories

How/should we embrace pipelining?

How/should we embrace pipelining?

Non-pipelined core feels natural for template instantiation (or any scheme without an instruction stream)

How/should we embrace pipelining?

Non-pipelined core feels natural for template instantiation (or any scheme without an instruction stream)

Obviously this limits clock frequency...but max clock frequency has been surprisingly fragile and hard to debug

How/should we embrace pipelining?

Non-pipelined core feels natural for template instantiation (or any scheme without an instruction stream)

Obviously this limits clock frequency...but max clock frequency has been surprisingly fragile and hard to debug

Maybe the answer is hardware multi-threading?

Conclusion

A simple, tiny template instantiation core seems to keep up a modern CPU
running at x10 speed!

Conclusion

A simple, tiny template instantiation core seems to keep up a modern CPU running at $\times 10$ speed!

Assisted by a fully concurrent GC (modulo ≈ 20 cycles per pass).

Requires several hardware tricks unavailable to software implementations.

Conclusion

A simple, tiny template instantiation core seems to keep up a modern CPU running at $\times 10$ speed!

Assisted by a fully concurrent GC (modulo ≈ 20 cycles per pass).

Requires several hardware tricks unavailable to software implementations.

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

Conclusion

A simple, tiny template instantiation core seems to keep up a modern CPU running at $\times 10$ speed!

Assisted by a fully concurrent GC (modulo ≈ 20 cycles per pass).

Requires several hardware tricks unavailable to software implementations.

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

Great to see there is some shared interest in hardware architectures for FP once again!

Questions?
