# Cloaca
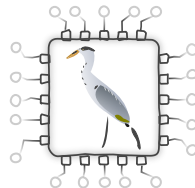
A Concurrent Hardware Garbage Collector for Non-strict Functional Languages

Craig Ramsay & Rob Stewart

September 2024

Heriot-Watt University

## Software for Concurrent GC

1960

Recursive Functions of Symbolic Expressions
and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

1990

Real-Time Garbage Collection on General-Purpose
Machines

Taiichi Yuasa
*Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan*

2020

**Alligator Collector: A Latency-Optimized Garbage
Collector for Functional Programming Languages**

Ben Gamari          Laura Dietz

## Hardware for FP

Can Programming Be Liberated from the von
Neumann Style? A Functional Style and Its
Algebra of Programs

John Backus

1977

BWM
a concrete machine for graph reduction

Lennart Augustsson

1991

*Reduceron reconfigured and re-evaluated*

MATTHEW NAYLOR and COLIN RUNCIMAN

2012

**Heron: Modern Hardware Graph Reduction**

Craig Ramsay          Robert Stewart

2024

## Software for Concurrent GC

**1960** — Recursive Functions of Symbolic Expressions and Their Computation by Machine

John McCarthy, Massachusetts Institute of Technology, Cambridge, Mass.

**1990** — Real-Time Garbage Collection on General-Purpose Machines

Taiichi Yuasa
Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan

**2020** — Alligator Collector: A Latency-Optimized Garbage Collector for Functional Programming Languages

Ben Gamari          Laura Dietz

## Hardware for FP

**1977** — Can Programming Be Liberated from the von Neumann Style? A Functional Style and Its Algebra of Programs

John Backus

**1991** — BWM
a concrete machine for graph reduction

Lennart Augustsson

**2012** — Reduceron reconfigured and re-evaluated

MATTHEW NAYLOR and COLIN RUNCIMAN

**2024** — Heron: Modern Hardware Graph Reduction

Craig Ramsay          Robert Stewart

*Widen von Neumann bottleneck with wide memories and complex stack mutations*

## Software for Concurrent GC

*1960*

Recursive Functions of Symbolic Expressions
and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

*1990*

Real-Time Garbage Collection on General-Purpose
Machines

Taiichi Yuasa
*Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan*

*2020*

Alligator Collector: A Latency-Optimized Garbage
Collector for Functional Programming Languages

Ben Gamari

## Hardware for FP

*1977*

Can Programming Be Liberated from the von
Neumann Style? A Functional Style and Its
Algebra of Programs

John Backus

*1991*

BWM
a concrete machine for graph reduction

Lennart Augustsson

*2012*

*Reduceron reconfigured and re-evaluated*

MATTHEW NAYLOR and COLIN RUNCIMAN

*2024*

**Heron: Modern Hardware Graph Reduction**

Craig Ramsay                    Robert Stewart

*Mean 0.6 hand-reductions per cycle*

## Software for Concurrent GC

## Hardware for FP

*1960*

Recursive Functions of Symbolic Expressions
and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

*1990*

Real-Time Garbage Collection on General-Purpose
Machines

Taiichi Yuasa
Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan

*2020*

Alligator Collector: A Latency-Optimized Garbage
Collector for Functional Programming

Ben Gamari

Can Programming Be Liberated from the von
Neumann Style? A Functional Style and Its
Algebra of Programs

John Backus

*1977*

BWM
a concrete machine for graph reduction

Lennart Augustsson

*1991*

*Reduceron reconfigured and re-evaluated*

MATTHEW NAYLOR and COLIN RUNCIMAN

*2012*

Heron: Modern Hardware Graph Reduction

Craig Ramsay                    Robert Stewart

*2024*

Mean 0.6 hand-reductions per cycle
x5 better than GHC* per cycle

Software for Concurrent GC

1960 — Recursive Functions of Symbolic Expressions and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

1990 — Real-Time Garbage Collection on General-Purpose Machines

Taiichi Yuasa
Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan

2020 — Alligator Collector: A Latency-Optimized Garbage Collector for Functional Programming

Ben Gamari

Hardware for FP

1977 — Can Programming Be Liberated from the von Neumann Style? A Functional Style and Its Algebra of Programs

John Backus

1991 — BWM
a concrete machine for graph reduction

Lennart Augustsson

2012 — Reduceron reconfigured and re-evaluated

MATTHEW NAYLOR and COLIN RUNCIMAN

2024 — Heron: Modern Hardware Graph Reduction

Craig Ramsay          Robert Stewart

Mean 0.6 hand-reductions per cycle
x5 better than GHC* per cycle
...but a tiny Heron is <200 MHz

## Software for Concurrent GC

**1960** — Recursive Functions of Symbolic Expressions and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

**1990** — Real-Time Garbage Collection on General-Purpose Machines

Taiichi Yuasa
*Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan*

**2020** — **Alligator Collector: A Latency-Optimized Garbage Collector for Functional Programming Languages**

Ben Gamari                    Laura Dietz

## Hardware for FP

**1977** — Can Programming Be Liberated from the von Neumann Style? A Functional Style and Its Algebra of Programs

John Backus

**1991** — A concrete machine for graph reduction
Daniel Amagbe... n

**2012** — *Reduceron reconfigured and re-evaluated*

MATTHEW NAYLOR and COLIN RUNCIMAN

**2024** — **Heron: Modern Hardware Graph Reduction**

Craig Ramsay                  Robert Stewart

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

Software for Concurrent GC

1960 — Recursive Functions of Symbolic Expressions and Their Computation by Machine, Part I

John McCarthy, Massachusetts Institute of Technology, Cambridge,

1990 — Real-Time Garbage Collection on General-Purpose Machines

Taiichi Yuasa
Research Institute for Mathematical Sciences, Kyoto University, Kyoto, Japan

2020 — Alligator Collector: A Latency-Optimized Garbage Collector for Functional Programming Languages

Ben Gamari          Laura Dietz

Hardware for FP

1977 — Can Programming Be Liberated from the von Neumann Style? A Functional Style and Its Algebra of Programs

John Backus

1991 — BWM
a concrete machine for graph reduction

Lennart Augustsson

2012 — Reduceron reconfigured and re-evaluated

RUNCIMAN

2024 — Heron: Modern Hardware Graph Reduction

Craig Ramsay          Robert Stewart

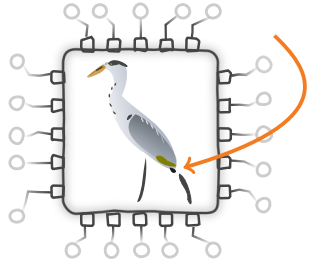"The nofib cases are quite mixed [...] most tests **slow down**, with a median of **+21%**"
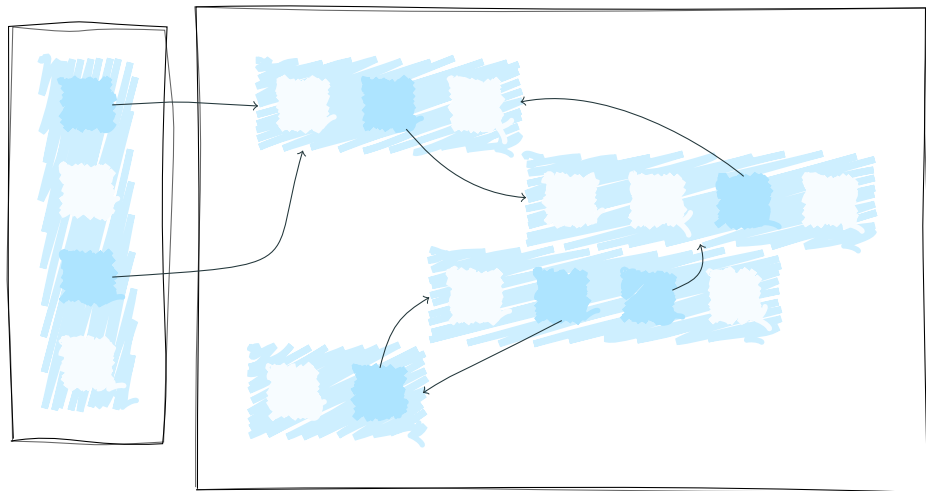
# Cloaca

noun [ C ]

/kloh-ah-kuh/

The system responsible for all the waste generated by a heron
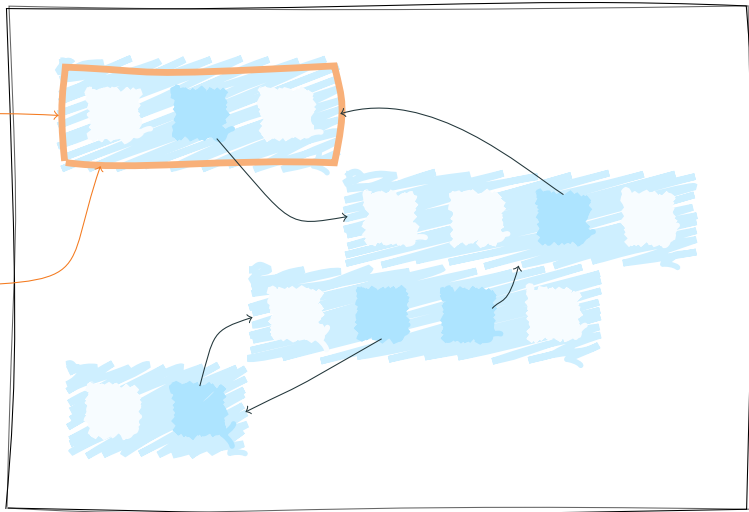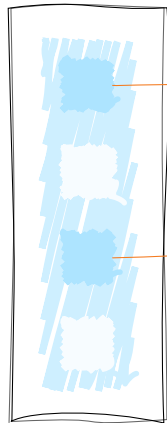
# Tracing example

Stack
(graph roots)

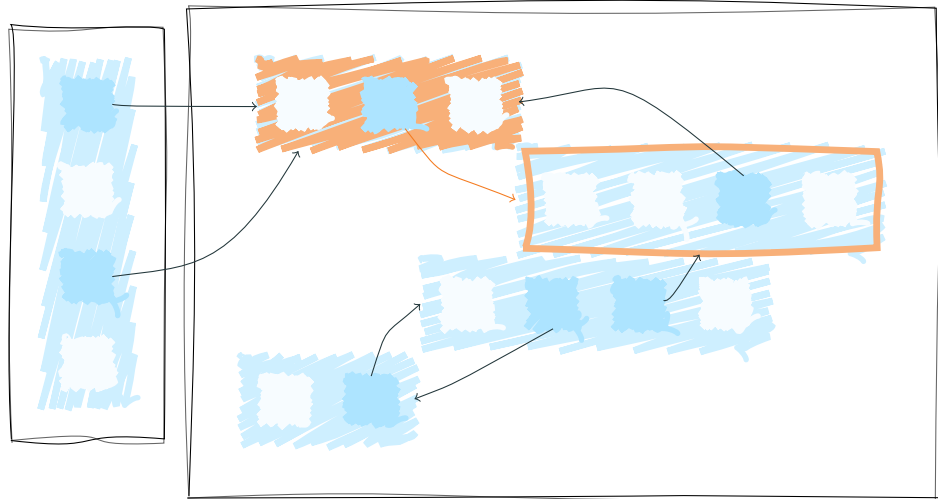Heap

Stack
(graph roots)

Heap
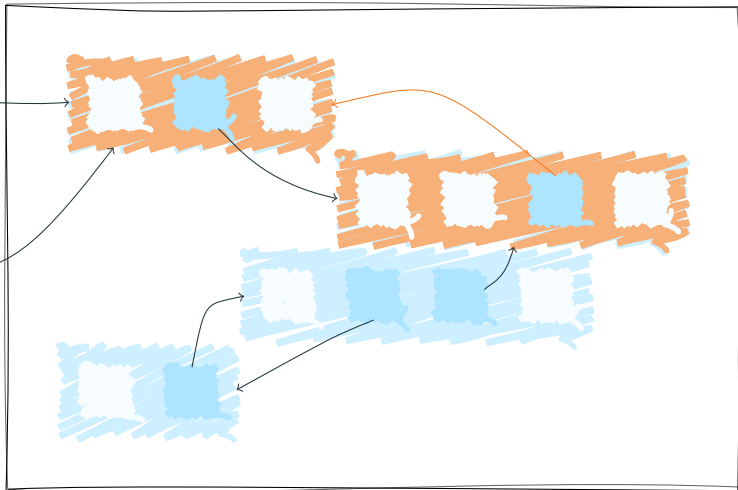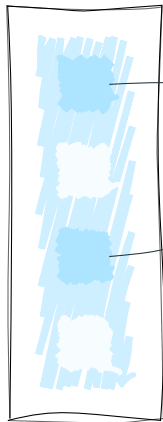
Root ID
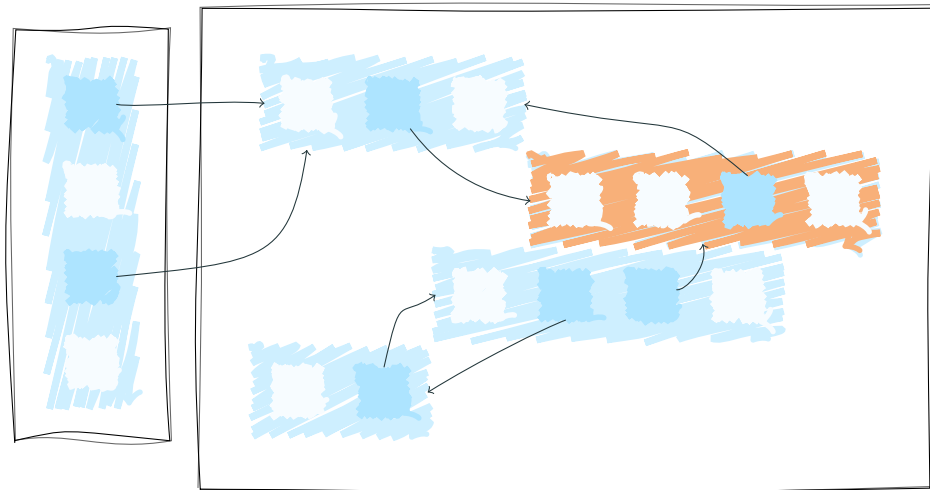
Stack
(graph roots)

Heap

Marking
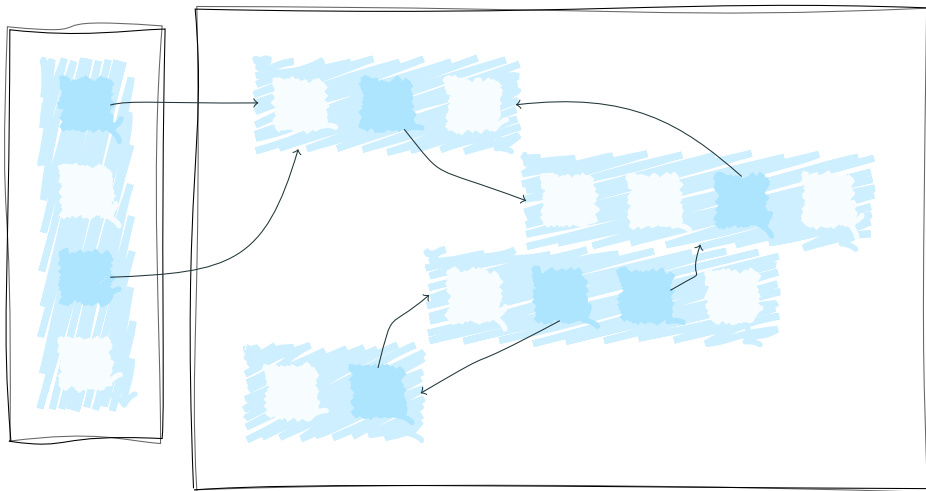
Stack
(graph roots)

Heap

Marking

Stack
(graph roots)

Heap

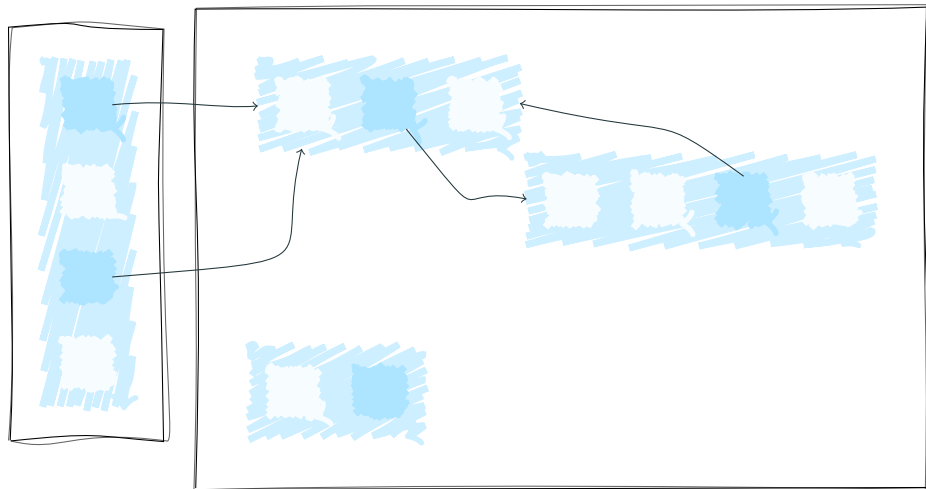Sweeping
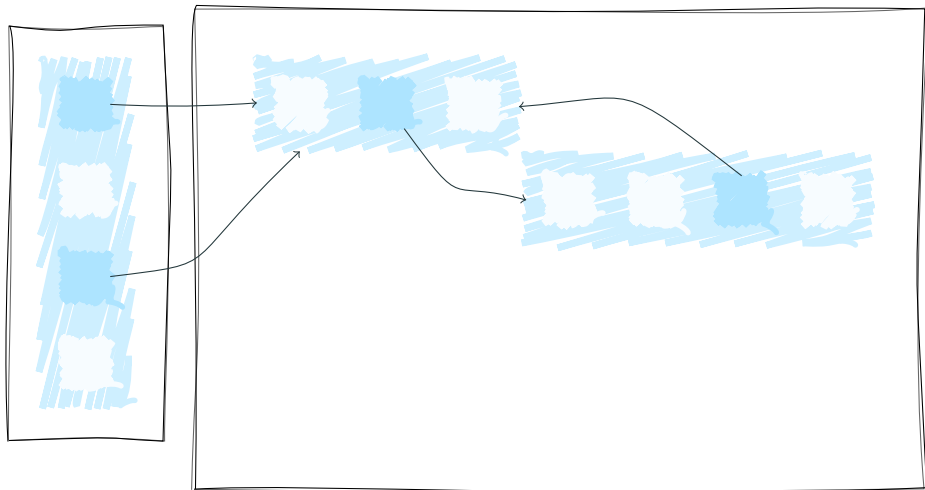
Stack
(graph roots)

Heap

Sweeping

Stack
(graph roots)

Heap

Sweeping

Stack
(graph roots)

Heap

Sweeping

Stack
(graph roots)
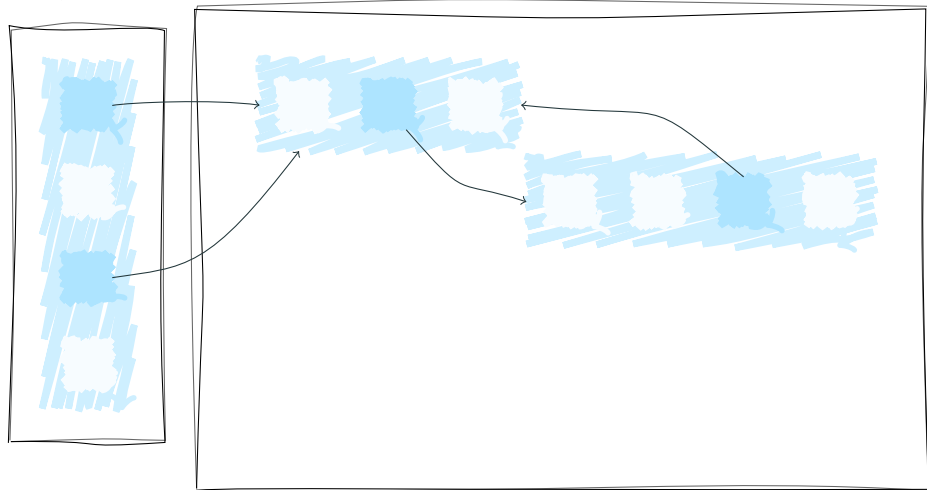
Heap

Sweeping

# 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
    ...
```
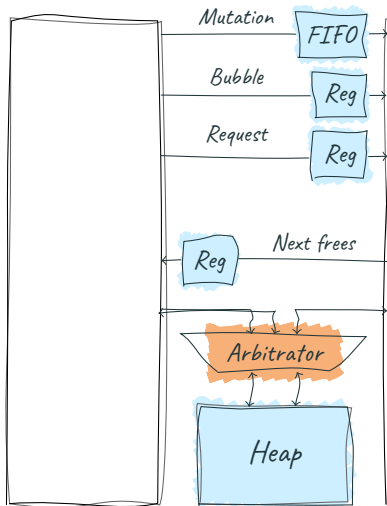
**Function** unwind (a, shared):
- ...
- ...
- ...
- **if** shared and not NF **then**
  └ push a onto update stack
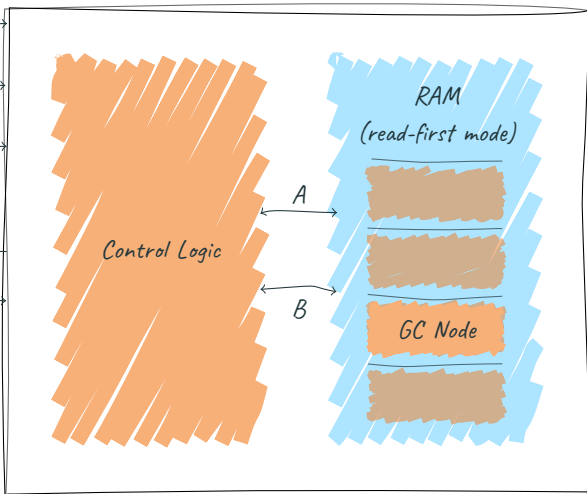- **if** not shared **then**
  └ dealloc a

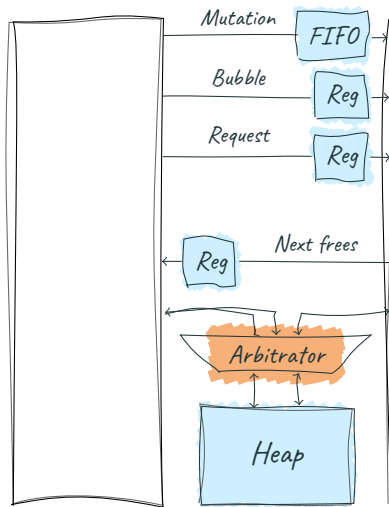*... is just one-bit reference counting with a hat on.*

# Architecture

**Reduction Core**

**Memory Management**

Mutation — FIFO

Bubble — Reg

Request — Reg

Reg — Next frees

Arbitrator

Heap

Control Logic

A

B

RAM
(read-first mode)

GC Node

**Reduction Core**

**Memory Management**

Mutation — FIFO

Bubble — Reg

Request — Reg

Next frees — Reg

Arbitrator

Heap

Control Logic

RAM
(read-first mode)

A

B

GC Node

data GC Node
= FreeList Addr
| WorkList Addr
| Marked
| Unmarked

Reduction Core

Memory Management

Mutation

FIFO

Bubble

Reg

Request

Reg

Next frees

Reg

Arbitrator

Heap

Control Logic

A

B

RAM
(read-first mode)

GC Node

```
write a x = do
  y <- readMem a
  writeMem a x
  pure y
```

```
data GC Node
  = FreeList Addr
  | WorkList Addr
  | Marked
  | Unmarked
```

time

Idle

Running

Mutator State

time

Idle

Root ID

Mark

Mutator State

Running

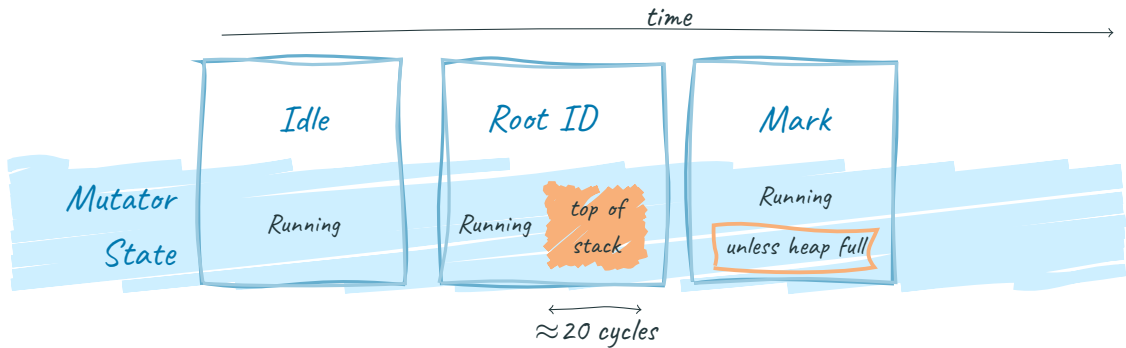Running    top of stack

Running
unless heap full

≈ 20 cycles

# Results

| | 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 |

| | 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 |

GC worst-case pause

GHC

Heron

$X = 3$

Max GC pause (μs)

Heap size / Peak working set

● adjoxo   ■ braun   ▲ cichelli   ▼ clausify   ⬟ countdown   ● knuthbendix   ◆ mate   ⊗ mss
◆ ordlist   ◇ permsort   ◇ queens   ● queens2   ■ sumpuz   ▲ taut   ◇ while

GC wall-clock overhead

GHC — Heron

%GC / Heap size / Peak working set

Legend: adjoxo, braun, cichelli, clausify, countdown, knuthbendix, mate, mss, ordlist, permsort, queens, queens2, sumpuz, taut, while

Concurrent tracing GC $\rightarrow$ high throughput and low latency?

Concurrent tracing GC → high throughput and low latency?

Stock processors + software struggle to maintain throughput.
Write-barriers and friends are hard.

Concurrent tracing GC → high throughput and low latency?

Stock processors + software struggle to maintain throughput.
Write-barriers and friends are hard.

Custom hardware with dual-port read-first memories
can handle them in a single cycle.

Concurrent tracing GC → high throughput and low latency?

Stock processors + software struggle to maintain throughput.
Write-barriers and friends are hard.

Custom hardware with dual-port read-first memories
can handle them in a single cycle.

Cloaca often pauses for only $\approx 20$ cycles for a GC pass,

Concurrent tracing GC → high throughput and low latency?

Stock processors + software struggle to maintain throughput.
Write-barriers and friends are hard.

Custom hardware with dual-port read-first memories
can handle them in a single cycle.

Cloaca often pauses for only ≈ 20 cycles for a GC pass,
and catches ≈ 50% of all our garbage before tracing.

Questions?

# Break glass in case of emergency

```
tails  [ ]       = [ ]
tails  (x : xs) = (x : xs) :  tails  xs

inits  xs =
  case  xs  of
    [ ]        -> [[ ]]
    (y : ys) -> xs :  inits  (init  xs)

segments xs = concatMap tails (inits xs)

mss = maximum . map sum . segments

main = let x = 0 - 50
           y = 150
       in  mss $ enumFromTo x y
```

| $e ::=$ | | Expressions |
| | $\bar{e}$ | (Application) |
| | $\mid$ *case* $e$ *of* $\bar{a}$ | (Case expression) |
| | $\mid$ *let* $\bar{b}$ *in* $e$ | (Let expression) |
| | $\mid n$ | (Integer) |
| | $\mid x$ | (Variable) |
| | $\mid \otimes$ | (Primitive Op) |
| | $\mid f$ | (Function) |
| | $\mid K$ | (Constructor) |

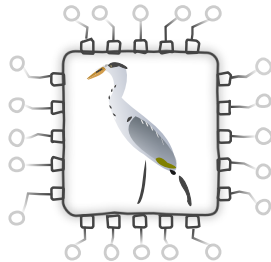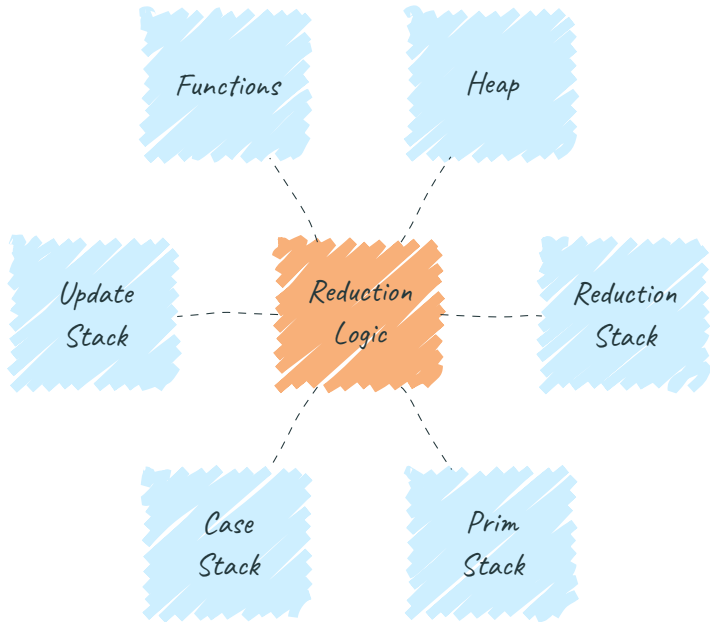| $a ::= K\,\bar{x} \rightarrow e$ | Case alternative |
| $b ::= x \mapsto e$ | Let binding |
| $d ::= f\,\bar{x} = e$ | Function definition |

# Heron

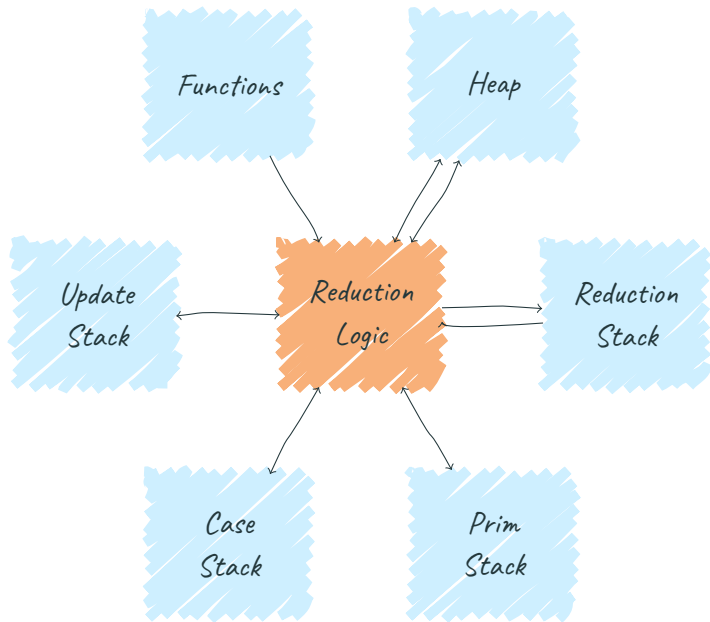noun [C]
/'"herən/

A processor for lazy functional languages.
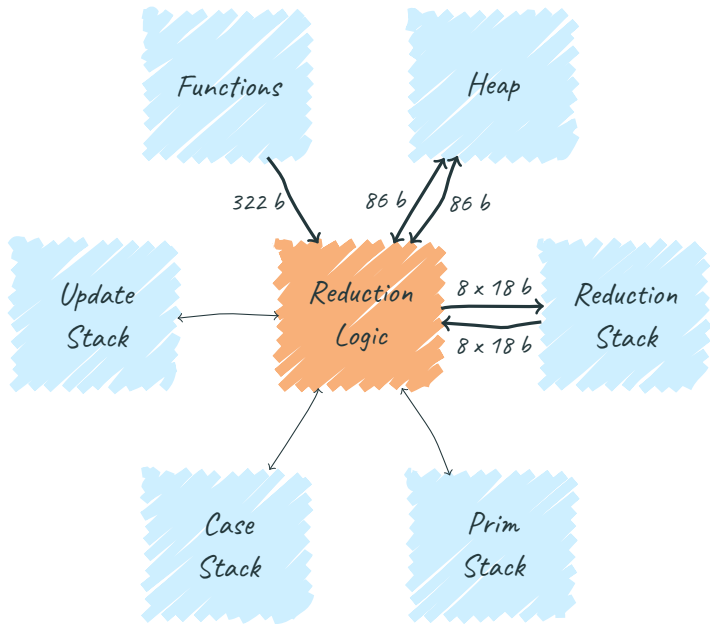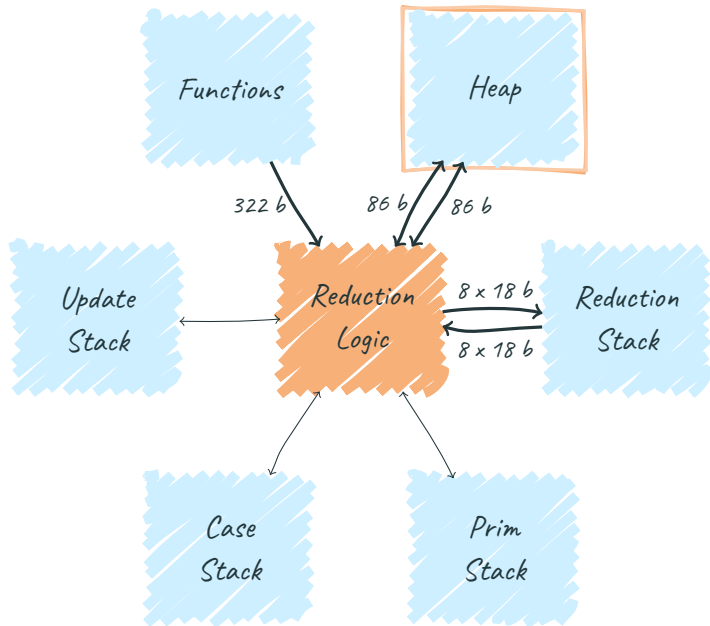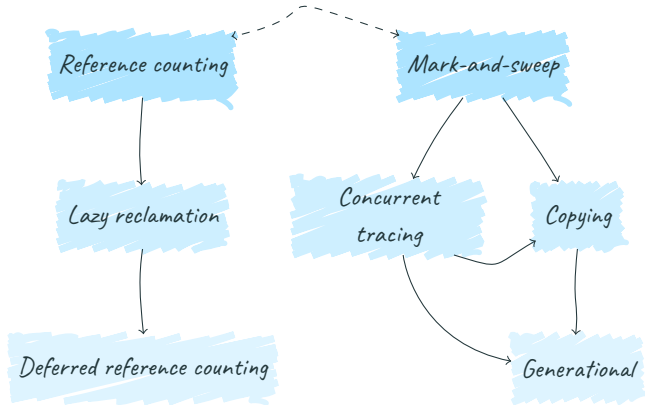Performs beta reduction in one clock cycle via multiple,
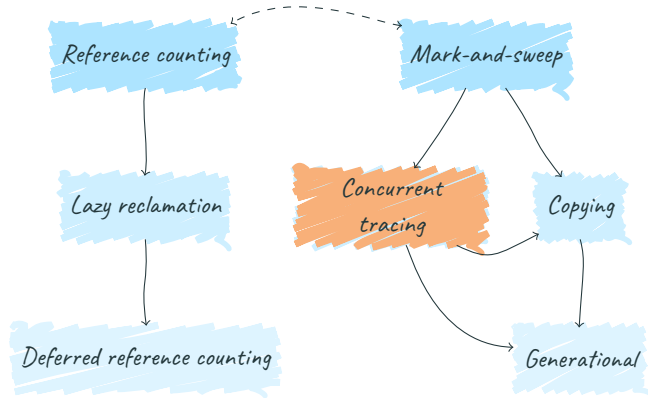wide, multi-ported memories.

Reduction
Logic

Functions

Heap

Update Stack

Reduction Logic

Reduction Stack

Case Stack

Prim Stack

Functions

Heap

Update
Stack

Reduction
Logic

Reduction
Stack

Case
Stack

Prim
Stack

```mermaid
graph TD
    Reference_counting["Reference counting"] -.- Mark_and_sweep["Mark-and-sweep"]
    Reference_counting --> Lazy_reclamation["Lazy reclamation"]
    Lazy_reclamation --> Deferred_reference_counting["Deferred reference counting"]
    Mark_and_sweep --> Concurrent_tracing["Concurrent tracing"]
    Mark_and_sweep --> Copying["Copying"]
    Concurrent_tracing --> Copying
    Concurrent_tracing --> Generational["Generational"]
    Copying --> Generational
```

- Reference counting
- Mark-and-sweep
- Lazy reclamation
- Concurrent tracing
- Copying
- Deferred reference counting
- Generational

GC worst-case pause

GHC · Reduceron · Heron

Max GC pause ($\mu s$) vs. Heap size / Peak working set

$X = 3$

adjoxo · braun · cichelli · clausify · countdown · knuthbendix · mate · mss
ordlist · permsort · queens · queens2 · sumpuz · taut · while

GC wall-clock overhead

**GHC** · **Reduceron** · **Heron**

Axes (all three plots): %GC vs Heap size / Peak working set

Heron plot annotation: $X = 3$

Legend:
- adjoxo
- braun
- cichelli
- clausify
- countdown
- knuthbendix
- mate
- mss
- ordlist
- permsort
- queens
- queens2
- sumpuz
- taut
- while

Total GC wall-clock pause

Normalised GC wall-clock time

adjoxo, braun, cichelli, clausify, countdown, knuthbendix, mate, nss, ordlist, permsort, queens, queens2, sumpuz, taut, while

GHC (Intel i7-1250U – power-saver @ ≈ 2 GHz)   GHC (Intel i7-1250U – performance @ 4.7 GHz)   Heron @ 185 MHz