

# PLAN Virtual Machine for AMD64 Implementation Manual

Benjamin Summers and Elliot Glaysher

February 19, 2026

## Contents

<b>1</b>	<b>Formatting Conventions</b>	<b>1</b>
1.1	Assembly Code	1
1.2	C Code	1
<b>2</b>	<b>Register Conventions</b>	<b>1</b>
2.1	System V Amd64 Register Conventions	1
2.2	Global Register Allocations	1
2.3	Calling Pins and Laws	1
2.4	Calling Thunks	1
2.5	Calling Asm From C	1
<b>3</b>	<b>Convenience Macros</b>	<b>2</b>
3.1	pdrop	2
3.2	ppush	2
3.3	ppop	2
3.4	Block Saves	2
3.5	sub0	2
3.6	cap	2
3.7	jzero	2
3.8	btw	2
3.9	swapq	2
3.10	swapst	2
<b>4</b>	<b>Pointer Tagging</b>	<b>2</b>
4.1	Essential Properties	2
4.2	Operations on Heap References	2
4.2.1	Extracting the Raw Pointer	2
4.2.2	Extracting an Type Enum	2
4.2.3	Direct Numbers	2
4.2.4	Numbers	2
4.2.5	Thunks, App, Laws	2
4.2.6	Laws	2
4.2.7	Specific Closure Shapes	2
4.3	Asserts and Debugging	2
4.3.1	Assertion Failure	2
4.3.2	Assert Buddy List Pointer	2
4.3.3	No 0xBADEEFBADBEEF Values	3
<b>5</b>	<b>Heap Object Layout</b>	<b>3</b>
5.1	GC Headers	3
5.2	Extracting Data from Heap Objects	3
5.3	Big Numbers	3
5.4	Closures	3
5.5	Thunks	3
5.6	Liquid Pins	3
5.7	Frozen Pins	3
5.8	Mega Pins	3
5.9	Claim and Reserve	3
5.10	Pending Changes	3
<b>6</b>	<b>Allocation</b>	<b>3</b>
6.1	mkbuffer	3
6.2	mkword	3
6.3	mkword64	3
6.4	mkdouble	3
<b>7</b>	<b>Numeric Utilities</b>	<b>3</b>
7.1	openat	3
7.2	unpack	3
<b>8</b>	<b>Pinning</b>	<b>3</b>
8.1	Hash Table	3
8.1.1	ht_create	3
8.1.2	htkey	4
8.1.3	htkpos	4
8.1.4	htval	4
8.1.5	ht_has	4
8.1.6	ht_put	4
8.1.7	ht_resize	4
<b>9</b>	<b>Evaluation</b>	<b>4</b>
9.1	Macros	4
9.1.1	EVAL	4
9.1.2	ERAX	4
9.1.3	JRAX	4
9.1.4	FORCE	4
9.1.5	ENAT	4
9.2	Thunk Executioners	4
9.2.1	xdone	4
9.2.2	xhole	4
9.2.3	xvar	4
9.2.4	xhead	4
9.2.5	xunknow	4
9.2.6	xunknownpupdate	4
<b>10</b>	<b>PLAN vs XPLAN Primops</b>	<b>4</b>
<b>11</b>	<b>Numeric Primops</b>	<b>4</b>
11.1	Nil	4
11.2	ToBit	4
11.3	Nat	4
11.4	Increment	5
11.5	Decrement	5
11.6	Addition	5
11.6.1	opadd	5
11.6.2	fastadd	5
11.6.3	slowadd	5
11.6.4	bufadd	5
11.6.5	bufadd1	5
11.7	Subtract	5
11.7.1	opsub	5
11.7.2	fastsub	5
11.7.3	slowsub	5
11.7.4	bufsub	5
11.7.5	bufsub1	5
11.8	Multiply	5
11.8.1	opmul	5
11.8.2	fastmul	5
11.8.3	slowmul	5
11.8.4	bufmul	5
11.8.5	bufmul1	5
11.9	Div	5
11.9.1	opdiv	5
11.9.2	fastdiv	6
11.9.3	slowdiv	6
11.9.4	divloop	6
11.10	Left Shift	6
11.10.1	oplsh	6
11.10.2	fastlsh	6
11.10.3	slowlsh	6
11.10.4	biglsh	6
11.10.5	buflsh	6
11.11	Read a Bit	6
11.12	Write a Bit	6
11.13	Clear a Bit	6
11.14	Read a Byte	6
11.15	Write a Byte	7
11.16	Write a 64-Bit Word	7
11.17	Write a 64-Bit Word In-Place	7
11.18	Read a Byte Range	7
11.19	Fast Copy of a Byte-Range	7
11.19.1	opsnore	7
11.19.2	fastsnore	7
11.19.3	smolnsnre	7
11.19.4	fumelsnre	7
11.19.5	bufsnore	7
11.19.6	pokesnre	7
<b>12</b>	<b>Closure Primops</b>	<b>7</b>
12.1	Closure Size	7
12.2	Closure Head	7
12.3	Last: Final Element	7
12.4	Init: Drop Final Element	7
12.5	Row: Closure from List	7
<b>13</b>	<b>Effectful Primops</b>	<b>8</b>
13.1	PeekOp	8
13.1.1	fastpeekop	8
13.2	PokeOp	8
13.2.1	fastpokeop	8
13.2.2	pokeraw	8
<b>14</b>	<b>Heaps and Process Layout</b>	<b>8</b>
<b>15</b>	<b>The Per Thread Heap</b>	<b>8</b>
15.1	Cheney Sizes	8
15.1.1	build_heap_table	8
15.2	Cheney State	8
15.3	Cheney Collection	8
15.4	asmgc	8
<b>16</b>	<b>The Buddy Allocator</b>	<b>8</b>
16.1	Buddy Configuration	8
16.2	Buddy Bucket Structure	8
16.2.1	bucket_for_request	8
16.3	GC Headers	8
16.4	Buddy Node Structure	8
16.4.1	nodeget	8
16.4.2	nodeset	8
16.4.3	ptr_for_node	9
16.4.4	node_for_ptr	9
16.5	List Utilities	9
16.5.1	list_init	9
16.5.2	list_push	9
16.5.3	list_remove	9
16.5.4	list_pop	9
16.6	Aligned Bit Trees	9
16.6.1	aligned_bit_set	9
16.6.2	aligned_bit_get	9
16.7	Buddy Collection	9
16.7.1	Data Structures	9
16.7.2	Locks	9
16.7.3	atomic_max	9
16.7.4	project_or	9
16.8	Buddy Sweeping	9
16.8.1	recursively_prune	9
<b>17</b>	<b>The Persistence Heap</b>	<b>9</b>
17.1	Low Level Page Storage	9
17.1.1	Mounting The Current Binary With Mmap	9
17.1.2	Persistence Format	10
17.1.3	Persistence System Initialization	10
17.1.4	Reading a Superblock	10
17.1.5	Writing a Superblock	10
17.1.6	Determining Which Superblock To Use	10
17.1.7	Allocating Raw Pages	10
17.1.8	Msyncing pages	10
17.2	Persisting PLAN values	10
17.3	Persistence Collection	10
17.4	Scrubbing and Validation	10
17.4.1	Calculating CRC32 Checksums	10
17.4.2	Validating Each Item	10
17.4.3	Scrubbing An Entire Pin Tree	10
<b>18</b>	<b>Naive Compiler</b>	<b>11</b>
18.1	Wrappers	11
18.1.1	Recognizing Wrappers	11
18.1.2	match	11
<b>19</b>	<b>Profiling</b>	<b>11</b>
19.1	Profile Format	11
19.2	Time Format	11
19.3	Push Profile Page	11
19.4	Zero Remaining Profile Page	11
19.5	recordevent	11
19.6	oprofile	11
19.7	SetProfEnabledOp	11
19.8	GetProfOp	11
<b>20</b>	<b>Seed</b>	<b>11</b>
20.1	The Boot Seed	11
20.2	The Seed Format	11
20.2.1	Introduction	11
20.2.2	The Header	11
20.2.3	Numbers	11
20.2.4	Scopes and References	11
20.2.5	Fragments	11
20.2.6	Encoding Pins and Laws	11
20.2.7	Important Properties	11
20.3	The Seed Loader	11
20.3.1	mmapfile	12
20.3.2	seedfile	12
20.3.3	seed	12
20.3.4	frags	12
20.3.5	frag	12
<b>21</b>	<b>Parallelism</b>	<b>12</b>
21.1	Mutex	12
21.1.1	mutex_init	12
21.1.2	mutex_lock	12
21.1.3	mutex_unlock	12
21.2	Condition Variables	12
21.2.1	condition_init	12
21.2.2	condition_signal	12
21.2.3	condition_broadcast	12
21.2.4	condition_wait	13
21.3	Read/Write Locks	13
21.3.1	rwlock_init	13
21.3.2	rwlock_read_lock	13
21.3.3	rwlock_read_unlock	13
21.3.4	rwlock_write_lock	13
21.3.5	rwlock_write_unlock	13
21.4	Barrings	13
21.4.1	barrier_set	13
21.4.2	barrier_done	13
21.4.3	barrier_wait	13
21.5	Threads	13
21.5.1	thread_create	13
21.5.2	thread_join	13
21.6	Futex Helpers	13
21.6.1	futex_wait_private	13
21.6.2	futex_wake_private	13
21.6.3	futex_cmp_requeue_private	13
<b>22</b>	<b>Proposals</b>	<b>13</b>
22.1	Seed Format Changes	13
22.1.1	Constant Reference Width	13
22.1.2	Encoding Graphs	13
22.1.3	Smaller Words	13
22.2	Pointed Tagging Optimizations	13
22.2.1	First Alternative	13
22.2.2	Second Alternative	14
22.3	Buffers / FATs	14
22.3.1	xbuf	14
22.3.2	GC1 Collection	14
22.3.3	BigNats always FAT	14
22.4	Actors	14
22.5	Specialized Executioners	14
22.6	Fredlists in Persistence	14
22.7	Register Reform	14
22.8	Killing the Shadow Stack	14
22.9	Faster Evaluation Preludes	14
22.10	Exit Segfaults	14
22.11	GC2	

### 3 Convenience Macros

#### 3.1 pdrop

```
.macro pdrop n
    .set bytes, 8 * \n
    add    r15, bytes
.endm
```

Drops a certain number of items from the shadow stack.

#### 3.2 ppush

Pushes zero or more registers to the shadow stack.

```
.macro ppush regs:vararg
    .ifnc "\regs", ""
        .set n\@, 0
        .irp r, \regs
            .set n\@, n\@ + 1
        .endr

        .set total, 8 * n\@
        sub    r15, total

        .set offset\@, 0
        .irp r, \regs
            mov    qword ptr [r15 + offset\@], \r
            .set offset\@, offset\@ + 8
        .endr
    .endif
.endm
```

#### 3.3 ppop

Restores a sequence of registers from the shadow stack. Note that the order corresponds with the order using for ppush. If you ppush a, b, c, you also ppop a, b, c

```
.macro ppop regs:vararg
    .ifnc "\regs", ""
        .set n\@, 0
        .irp r, \regs
            .set n\@, n\@ + 1
        .endr

        .set total, 8 * n\@

        .set offset\@, 0
        .irp r, \regs
            mov    \r, qword ptr [r15 + offset\@]
            .set offset\@, offset\@ + 8
        .endr
        add    r15, total
    .endif
.endm
```

#### 3.4 Block Saves

These are only used in the division logic, in order to push multiple values to the C stack. This should probably eventually be replaced with a generic macro like ppush and ppop.

```
.macro pop_r8_to_rdi
    pop    r8
    pop    rcf
    pop    rdx
    pop    rsi
    pop    rdi
.endm
```

```
.macro push_rdi_to_r8
    push   rdi
    push   rsi
    push   rdx
    push   rcx
    push   r8
.endm
```

#### 3.5 sub0

sub0 is unsigned subtraction with a floor at 0. If the subtraction would underflow, the result is zero.

```
.macro sub0 a, b
    cmp    \a, \b
    cmovb \a, \b
    sub    \a, \b
.endm
```

#### 3.6 cap

cap(r,max) just computes r=min(r,max) in order to keep r within a certain bounds. Both inputs must be in registers.

```
.macro cap r, max
    cmp    \r, \max
    cmova \r, \max
.endm
```

#### 3.7 jzero

Jump to a label if a register is equal to zero.

```
.macro jzero reg, label
    test   \reg, \reg
    jz     \label
.endm
```

#### 3.8 btw

```
.macro btw r1                                # bits_to_words(x) = (x+63)>>6
    add    \r1, 63
    shr    \r1, 6
.endm
```

#### 3.9 swapq

Using two temporary registers r1, r2, swap the values behind two locations in memory.

```
.macro swapq r1, r2
    mov    \r1, \m1
    mov    \r2, \m2
    mov    \m1, \r2
    mov    \m2, \r1
.endm
```

#### 3.10 swapstk

Using two temporary registers r1, r2, swap the two to values on the stack.

```
.macro swapstk r1, r2
    swapq \r1, \r2, [r15], [r15+8]
.endm
```

### 4 Pointer Tagging

Numbers that fit within 63-bits are stored unboxed, directly in the value. A high bit being set indicates that the value is a heap reference.

Heap references are represented as normal pointer but with tagging data stored in the high 16 bits. The following conventions are used:

u63 0nnnnnnnnnnnn - 0x0000..0x7fff (n=nat data)  
NAT 10000100000000 - 0x8200  
PIN 11000100000000cm - 0xC40x (c=hascrc32, m=ismegapin)  
LAW 1100100000000000 - 0xC800  
CLZ 1101000ttttzzzz - 0xD0tz (t=tag, z=size)  
THK 1110000000000000 - 0xE000

The m flag indicates that a pin is a megapin, which is a pinned pin. Megapins cache information about all of their subpins.

The c flag indicates that the pin has been hashed using crc32 (which is directly supported by the CPU). Pins that have this information are bigger, having the hash data append to the end of the structure.

If the closure tag information is 0b1111, that indicates that the tag is too big and must be loaded from the closure itself.

If the closure size is 0b0000, that indicates that the size is too big and must be loaded from the closure itself. 0 is chosen because 0 is not a valid closure size.

#### 4.1 Essential Properties

This scheme works well because it maintains a number of essential properties:

1. The first property is that direct numbers are encoded directly in the normal way, without any need to encode or decode them. This also makes for very tight fast-paths for numeric primops.
2. The second property is that all of the information needed in order to determine the type is available in the high byte. This makes it possible to perform comparisons on it, after roll with al, r12b, etc. It also makes it very visible in hex-printouts.
3. The third property, because the actual representations are ordered by type, we can recognize ranges of possible types with a comparison. For example, if the high-bit is greater than or equal to the smallest possible closure value, then the result is either a thunk or a closure. Similarly, anything smaller than a pin is a number; either direct or indirect.
4. The fourth property is that each heap reference can be determined by checking a specific bit. This only works if we already know that the reference is not a direct number, but if we do know that we can do the determination with a single 'bt' instruction.
5. The fifth property is that each type of heap reference has a different number of leading zeros, so we can use lzcnt to convert from a heap reference to an enum of all possible heap types.
6. The sixth property is that, the metadata associated with a type also lives in the second highest byte, which has the same advantages as having the type live in the top-most byte.

#### 4.2 Operations on Heap References

##### 4.2.1 Extracting the Raw Pointer

Converting a tagged value into a pointer just requires two shifts to zero out the high bits.

```
.macro ptr r
    shl    \r, 16
    shr    \r, 16
.endm
```

```
.macro ptr_r, from
    mov    \r, 48
    bzh   \r, \from, \r
.endm
```

```
.macro dref_r
    ptr    \r
    mov    \r, [\r]
.endm
```

```
.macro dref_r, from
    mov    \r, \from
    dref   \r
.endm
```

```
.macro refo_r, o
    ptr    \r
    lea    \r, [\r + \o]
.endm
```

```
.macro refo_r, src, o
    mov    \r, \src
    refo   \r, \o
.endm
```

```
.macro drefo_r, o
    ptr    \r
    mov    \r, [\r + \o]
.endm
```

```
.macro drefo_r, src, o
    mov    \r, \src
    drefo  \r, \o
.endm
```

```
.macro clzix_r, src, i
    .set   ix, 8+(\i * 8)
    drefo  \r, \src, ix
.endm
```

```
.macro drefi_r, i
    ptr    \r
    mov    \r, [\r + \i]
.endm
```

```
.macro drefi_r, src, i
    mov    \r, PTR(src)[i]
.endm
```

```
.macro drefo_r, o
    ptr    \r
    mov    \r, [\r + \o]
.endm
```

```
.macro jclz_r, reg, lab
    .set   ix, 8+(\i * 8)
    drefo  \r, \reg, ix
.endm
```

```
.macro drefi_r, reg, lab
    ptr    \r
    mov    \r, PTR(r)[i]
.endm
```

```
.macro drefo_r, reg, lab
    ptr    \r
    mov    \r, [\r + \i]
.endm
```

There are also slightly more efficient variants which can be used if we know that the input is not a direct number:

jdirect jumps to a label if the register is a direct number, and jheap jumps if it is anything else.

```
.macro jdirect reg, label
    test   \reg, \reg
    jns    \label
.endm
```

```
.macro jheap reg, label
    test   \reg, \reg
    js    \label
.endm
```

jdirect jumps to a label if the the register is a direct number, and jheap jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

jdtag jumps to a label if the register is a direct number, and jttag jumps if it is anything else.

```
.macro jdtag r, lab
    .shl   \r, 2
    .lzcnt \r, \r
.endm
```

```
.macro jttag_out, from
    mov    \out, \from
    jttag \out
.endm
```

jdtag jumps to a label if the the register is a direct number, and jttag jumps if it's anything else.

&lt;pre

### 4.3.3 No 0xBADBEEFBADBEEF Values

Validates that the register does not contain the sentinel 0xBADBEEFBADBEEF value.

TODO: An imm64 is kind of heavyweight. Can we just check the lower word and then jump on low word to the rest of the check somehow?

```
.macro nobeef reg
    movabs r12, 0xBADBEEFBADBEEF
    cmp    \reg, r12
    je     assert_fail
.endm
```

## 5 Heap Object Layout

The word immediately preceding each heap object is a GC record which includes its size and type. This is needed by the first-generation GC because Cheney's algorithm needs to do a linear walk of the heap.

At the moment, this size is \*also\* used to encode the size of numbers and closures, but in the future these will be stored separately. (Because of this, the GC size is stored in bits, not words.)

The GC header also has a separate type for NF closures and WHNF closures, which is used to cache `Force` operations to avoid duplicate work.

### 5.1 GC Headers

The GC header is always the word right before the one pointed at. It contains a type and a bit-size, which is always a multiple of 64 for things besides nats.

```
.macro getgc r          # gcheader(r) = PTR(r)[-8]
    ptr   \r
    mov   \r, [r - 8]
.endm

.macro gcsz hdr         # Size of allocation box
    add   \hdr, 16_128
    shr   \hdr, 14
.endm

.macro gcsz_ reg, hdr   # reg = gcsz(hdr)
    mov   \reg, \hdr
    gcsz  \reg
.endm

.macro gcbits hdr       # bits = gcbits(hdr)
    shr   \hdr, 8
.endm
```

We define a number of wrapper macros to formalize the usage patterns for extracting data from this word in different ways.

```
.macro bitsiz r
    getgc  \r
    gcbits \r
.endm

.macro bitsiz_ r, from
    mov    \r, \from
    bitsiz \r
.endm

.macro bytsiz r
    bitsiz \r
    add   \r, 7
    shr   \r, 3
.endm

.macro bytsiz_ r, from
    mov    \r, \from
    bytsiz \r
.endm

.macro dbitsz r
    xchg  \r, r12
    mov   \r, 64
    lzcnt \r12, r12
    sub   \r, r12
.endm

.macro dbitsz_ r, from
    lzcnt \r12, \from
    mov   \r, 64
    sub   \r, r12
.endm

.macro dbytsz_ r, from
    dbitsz_\r, \from
    add   \r, 7
    shr   \r, 3
.endm

.macro wordsz r
    getgc  \r
    gcsz   \r
.endm

.macro wordsz_ r, from
    mov    \r, \from
    wordsz \r
.endm

.macro slots_ r, from
    mov    \r, \from
    slots  \r
.endm

.macro slots r          # wordsize, assuming not a nat
    getgc  \r
    shr   \r, 14
.endm

.macro clzsz r          # Sz(x) = slots(x)-1
    slots  \r
    dec   \r
.endm

.macro clzsz_ r, from
    mov    \r, \from
    clzsz \r
.endm
```

### 5.2 Extracting Data from Heap Objects

Getting the tag of a closure is just a dereference.

```
.macro clzhdr r, from
    dref_ \r, \from
.endm

.macro ix0 r
    drefo_ \r, 8
.endm

.macro ix0_ r, from
    drefo_ \r, \from, 8
.endm

.macro ixi r
    drefo_ \r, 16
.endm

.macro ixi_ r, from
    drefo_ \r, \from, 16
.endm

.macro unpin r
    dref   \r
.endm

.macro unpin_ r, from
    dref_ \r, \from
.endm

.macro name r
    dref   \r
.endm

.macro name_ r, from
    dref_ \r, \from
.endm

.macro arity r
    drefo_ \r, 8
.endm

.macro arity_ r, from
    drefo_ \r, \from, 8
.endm

.macro body r
    drefo_ \r, 16
.endm

.macro body_ r, from
    drefo_ \r, \from, 16
.endm
```

## 5.3 Big Numbers

TODO

### 5.4 Closures

TODO

### 5.5 Thunks

TODO

### 5.6 Liquid Pins

TODO

### 5.7 Frozen Pins

TODO

### 5.8 Mega Pins

TODO

### 5.9 Claim and Reserve

TODO

### 5.10 Pending Changes

Once we add support for non-normalized pins and laws on the 1st-gen heap, we will need to have a NF flag for pins and laws as well.

TODO: explicit nat and closure size to enable in-place shrinking.

## 6 Allocation

### 6.1 mkbuffer

This creates a new bignum directly on GC2 in order to prevent it from moving around. This allows us to get pointers into nats for use during syscalls, and for casting to function pointers in order to generate code for pins and laws.

This takes a single argument: the number of bytes that we will need.

The process here is that we allocate one byte more than needed and set the high byte to 1 (to create a BAR of the right size).

buddy\_malloc does not zero the memory, so we do that. It also does not set the actual bit-size, and we need that, so we overwrite the GC header it creates with our own.

```
.global mkbuffer
mkbuffer:
    mov   r12, rdi          # mkbuffer: (rdi=size)
    mov   rbx, 8             # r12=bytes
    cmp   rdi, rbx           #
    shr   rdi, 3             # size = max(size, 8)
    inc   rdi               #
    mov   rbx, rdi           # words = bytes/8
    xor   rsi, rsi           # words++
    mov   rsi, rdi           # rbx=words
    xor   eax, eax           # tag=0
    call  buddy_malloc        # buddy_malloc() (rax=ptr)
    mov   rsi, rax           # copy pointer to rsi
    mov   rcx, rbx           # count=words
    mov   rdi, rax           # dst=ptr
    sub   rsi, rdi           # value=0
    mov   rdi, rax           # memcopy
    mov   rsi, rdi           # set high byte
    mov   rax, 0x8200000000000000 # result=NATMASK
    or    rax, rsi           # result |= ptr
    mov   rsi, rdi           # bits = (bytes*8)+1
    mov   rax, rsi           # hdr = (bits, tag=0)
    sub   rsi, 8              # replace header w/ correct sz
    mov   [rsi-8], r12        # return result
```

TODO: make sure this is actually collected properly.

TODO: make sure that GC2 correctly handles an actual bit-size, not just a word-size wrapped in a bit-size. If not, we can just provide a fake GC header as the first word within this allocation. We are going to have to do that eventually, once we include the mark and the intrinsic linked list.

### 6.2 mkword

`mkword` takes any 64-bit number, and returns a valid PLAN value. The result will either be a direct number, or a bignum depending on the size.

Because this is often used to finalize a result, this uses the unconventional calling convention of having both the argument and the return value in `rax`. Only `r8` is clobbered.

Note that the GC header is set to 0x4000 which is tag=0, size=64: 0x4000 = 64<<8 + 0. It is always this size, because all valid single-word heap-allocated nats are exactly 64-bits wide.

TODO: explicit nat and closure size to enable in-place shrinking.

## 7 Numeric Utilities

### 7.1 openat

Given any PLAN value, coerces the input into a number (non-numeric values are cast to 0), and returns a pointer to the `u64[]` containing the number along with its bit-width (special-casing 64 if 0).

This routine is used for numeric operations, and it solves two problems at once:

1. Handle the edge-case of a non-numeric input being passed to a operation on natural numbers. (All PLAN jets handle this by casting non-nat value to 0).
2. Converts direct and indirect natural numbers into a uniform representation, which makes it possible to use the same code for all combinations of representations (indirect - indirect), (direct - indirect), (indirect - direct), and (direct - direct).

There is one special cases, which is a bit of a hack: When the input is zero, we actually return 64 for the bit-size. In practice, we will always use the bit-size to compute the word-size of the resulting buffer, which is always one, even if the input is zero.

```
opennat: # rdi=plan_value,rsi=buffer          # openat:
    jdirect rdi, opennat.direct
    test   rax, rax           # if high bit is set
    js    mkword64             # tailcall mkword64(word)
    ret
```

### 7.2 unpack

Unpacks the top two items of the plan stack into the following standardized parts of a stack frame. This function has the unusual calling convention that it reads the top two items from the shadow stack, and then unpacks those items into the caller's stack frame:

Offset	Description
-8	x pointer tag
-16	x qword buffer for direct
-24	x bitsize
-32	x buffer
-40	x buffer word size
-48	y pointer tag
-56	y qword buffer for direct
-64	y bitsize
-72	y buffer
-80	y buffer word size

```
unpack: # unpack:
    mov   rdi, [r15+8]          # arg1
    mov   rdi, [rbp-8], rdi      # rbp[-8] = arg1
    mov   rsi, [r15]              # arg2
    mov   rsi, [rbp-48], rdi      # rbp[-48] = arg2
    lea   rsi, [rbp-16]          #
    call  opennat                #
    mov   rax, rsi                #
    mov   rax, [rbp-24], rax      #
    mov   rax, [rbp-32], rdx      #
    btw  rax, rdx                #
    mov   rax, rsi                #
    mov   rax, [rbp-56]          #
    call  opennat                #
    mov   rax, [rbp-64], rax      #
    mov   rax, [rbp-72], rdx      #
    btw  rax, rdx                #
    mov   rax, [rbp-80]          #
    ret
```

## 8 Pinning

Pinning is our system for taking a tree of values up to the next pin boundary and compacting it into a single contiguous segment, recording metadata about

### 8.1 Hash Table

Because pins can be arbitrarily large, we must define data structures that can grow during pin construction to perform the size calculations, pointer mapping, and pointer deduplication. We implement a simple linear probing hash table that maps a qword to the entire structure can be directly freed at the end.

```
typedef struct { uint64_t key, val; } entry_t;
typedef struct {
    uint64_t lg;           // exponent of the power of 2 size of table
    uint64_t t;            // mask for current size to apply to hashes
    uint64_t count;        // # of items, tracked for resizing
    entry_t *tbl[];        // flexible array
} ht_t;
```

### 8.1.1 ht\_create

We create a hash table by passing in the exponent of the power of two desired capacity of the table.

```
.global ht_create
ht_create:
    mov   rax, rdi          # ht_create(rdi=lg):
    mov   rsi, 1              # move to rax fro sh1
    cmp   rsi, 1              # base 1
    shr   rsi, 1              # size = 1 << lg
    inc   rsi                # save registers
    mov   rsi, 1              # bytes = size
    xor   rsi, 1              # bytes *= 2 (kevys+vals)
    mov   rsi, 1              # tag=Hash Table
    shr   rsi, 1              # buddy_malloc -> rax=buffer
    mov   rsi, 1              # restore registers
    shr   rsi, 1              # h->lg = lg;
    mov   rsi, 1              # mask = pow of two size
    shr   rsi, 1              # mask-- (now a mask of 1s)
    mov   rsi, 1              # h->mask = size - 1
    inc   rsi                # restore size
    mov   rsi, 1              # h->count = 0
    xor   rsi, 1              # save return value
    shr   rsi, 1              # rdi = beginning of kv buffer
    lea   rsi, [rsi+24]        # rsi = size * 2
    xor   rsi, 1              # set value to 0
    mov   rsi, 1              # zero out table memory
    mov   rsi, 1              # restore return value
    ret
```

## 9 Pinning

Pinning is our system for taking a tree of values up to the next pin boundary and compacting it into a single contiguous segment, recording metadata about

### 9.1 ht\_create

Because pins can be arbitrarily large, we must define data structures that can grow during pin construction to perform the size calculations, pointer mapping, and pointer deduplication. We implement a simple linear probing hash table that maps a qword to the entire structure can be directly freed at the end.

```
typedef struct { uint64_t key, val; } entry_t;
typedef struct {
    uint64_t lg;           // exponent of the power of 2 size of table
    uint64_t t;            // mask for current size to apply to hashes
    uint64_t count;        // # of items, tracked for resizing
    entry_t *tbl[];        // flexible array
} ht_t;
```

### 9.1.1 ht\_create

We create a hash table by passing in the exponent of the power of two desired capacity of the table.

```
.global ht_create
ht_create:
    mov   rax, rdi          # ht_create(rdi=lg):
    mov   rsi, 1              # move to rax fro sh1
    cmp   rsi, 1              # base 1
    shr   rsi, 1              # size = 1 << lg
    inc   rsi                # save registers
    mov   rsi, 1              # bytes = size
    xor   rsi, 1              # bytes *= 2 (kevys+vals)
    mov   rsi, 1              # tag=Hash Table
    shr   rsi, 1              # buddy_malloc -> rax=buffer
    mov   rsi, 1              # restore registers
    shr   rsi, 1              # h->lg = lg;
    mov   rsi, 1              # mask = pow of two size
    shr   rsi, 1              # mask-- (now a mask of 1s)
    mov   rsi, 1              # h->mask = size - 1
    inc   rsi                # restore size
    mov   rsi, 1              # h->count = 0
    xor   rsi, 1              # save return value
    shr   rsi, 1              # rdi = beginning of kv buffer
    lea   rsi, [rsi+24]        # rsi = size * 2
    xor   rsi, 1              # set value to 0
    mov   rsi, 1              # zero out table memory
    mov   rsi, 1              # restore return value
    ret
```

### 8.1.2 htkey

Given the base pointer and an index, write the key value at idx. We do things this way because it's cleaner at the call site to operate on masking the index.

```
.macro htkey base, idx, out
    mov    \out, \idx
    shl    \out, 4      # idx * 16; can't rdx*16 inside mov.
    mov    \out, [\base + \out + 24] # base->tbl[idx].key
.endm
```

### 8.1.3 htkpos

Given the base pointer and an index, write the pointer to the key at idx.

```
.macro htkpos base, idx, out
    mov    \out, \idx
    shl    \out, 4      # idx * 16; can't rdx*16 inside mov.
    lea    \out, [\base + \out + 24] # base->tbl[idx].key
.endm
```

### 8.1.4 htval

Like htkey, but points to the value.

```
.macro htval base, idx, out
    mov    \out, \idx
    shl    \out, 4      # idx * 16; can't rdx*16 inside mov.
    mov    \out, [\base + \out + 32] # base->tbl[idx].val
.endm
```

### 8.1.5 ht\_has

Returns the index at which a key exists, or -1 if key isn't found.

```
.global ht_has
ht_has:
    mov    rax, -1          # ht_has(rdi=ht*, rsi=key):
    xor    rax, rsi         #   initialize hash to -1
    crc32 rsi, rsi         #   crc32 on key value as hash
    mov    r8, [rdi + 8]    #   r8 = h->mask
    and    rax, r8          #   idx = hash & h->mask
ht_has.loop:
    htkey rdi, rax, rcx
    test   rcx, rcx
    jz     ht_has.not_found
    cmp    rdx, rsi
    je     ht_has.return
    inc    rax
    and    rax, r8
    jmp    ht_has.loop
ht_has.not_found:
    mov    rax, -1
ht_has.return:
    ret
```

### 8.1.6 ht\_put

Sets key to value, overwriting the current value if one is set.

When we put a value in a hash table, the first thing we have to do is make sure that the hash table is less than 70% full. If it is, we invoke the resizing behaviour, which doubles the size and rehashes everything.

```
ht_put:
    mov    rax, [rdi + 16]          # ht_put(rdi=ht*, rsi=key, rdx=val):
    inc    rax                      #   rax = tbl->count
    cvtsi2sd xmm0, rax             #   convert to double in xmm0
    mov    rax, [rdi + 8]           #   rax = tbl->mask
    inc    rax                      #   rax is now capacity
    cvtsi2sd xmm1, rax             #   convert to double in xmm1
    divsd xmm0, xmm1               #   xmm0 = (>count++) / (h->mask++)
    movsd xmm1, qword ptr [hash_table_load_factor]
    ucomisd xmm0, xmm1             #   if percent < load_factr
    jbe    ht_put.has_valid_hash_table
    ppush rdi, rsi, rdx            #   save registers
    call   ht_resize
    ppop  rdi, rsi, rdx            #   resize the table
    mov    rdi, rax                #   rdi is now resized table
ht_put.has_valid_hash_table:
    mov    rax, -1
    crc32 rsi, rsi
    mov    r8, [rdi + 8]
    and    rax, r8
ht_put.loop:
    mov    r9, rax
    shl    r9, 4
    lea    r9, [rdi + r9 + 24]
    mov    r10, [r9]
    test   r10, r10
    jz     ht_put.increment
    cmp    r10, rsi
    je     ht_put.store
    inc    rax
    and    rax, r8
    jmp    ht_put.loop
ht_put.increment:
    mov    r8, [rdi + 16]
    inc    r8
    mov    [rdi + 16], r8
ht_put.store:
    mov    [r9], rsi
    mov    [r9 + 8], rdx
    ret
```

We also have to store the floating point value for 0.7 in a place in rodata so it's loadable.

```
.section .rodata
.align 8
hash_table_load_factor:
.double 0.7
.section .data
```

### 8.1.7 ht\_resize

Resizing the hash table is incrementing the logarithmic size by 1 and then copying each item over.

```
ht_resize:
    push   rdi
    mov    rdi, [rdi]
    inc    rdi
    call   ht_create
    pop    rdi
    mov    r8, [rax + 8]
    mov    r10, [rdi + 16]
    mov    [rax + 16], r10
    mov    rsi, [rdi + 8]
    inc    rsi
    mov    rdx, 0
ht_resize.outer_loop:
    htkey rdi, rdx, rcx
    test   rcx, rcx
    jz     ht_resize.continue
    mov    r10, -1
    crc32 r10, rdx
    and    r10, r8
    jmp    ht_resize.inner_loop
ht_resize.inner_loop:
    htkey rax, r10, r11
    test   r11, r11
    jz     ht_resize.inner_loop_target
    inc    r10
    and    r10, r8
    jmp    ht_resize.inner_loop
ht_resize.inner_loop_target:
    htkey pos rax, r10, r11
    htkey pos rdi, rdx, rcx
    mov    r10, [rcx]
    mov    [r11], r10
    mov    r10, [rcx + 8]
    mov    [r11 + 8], r10
ht_resize.continue:
    inc    rdx
    cmp    rdx, rsi
    jb     ht_resize.outer_loop
    push   rax
    call   buddy_free
    pop    rax
    ret
```

## 9 Evaluation

### 9.1 Macros

#### 9.1.1 EVAL

EVAL Causes target to be evaluated, result in target, clobbers rax and r12. tosave registers are preserved

```
.macro EVAL target, tosave:vararg
    mov    rax, \target
    jnthk rax, 777f
    ppush \tosave
    dref_ r12, rax
    call_ r12
    mov    \target, rax
    ppop  \tosave
777:
.endm
```

#### 9.1.2 ERAX

This is a specialized version of EVAL, causes rax to be evaluated, clobbers r12. All of the registers in tosave are preserved

```
.macro ERAX tosave:vararg
    jnthk rax, 777f
    ppush \tosave
    dref_ r12, rax
    call_ r12
    ppop  \tosave
777:
.endm
```

#### 9.1.3 JRAX

JRAX evaluates and returns rax (uses a tail-call if rax is a thunk).

```
.macro JRAX
    jnthk rax, 777f
    dref_ r12, rax
    jmp    r12
777:
    ret
```

#### 9.1.4 FORCE

This evaluates the target to normal form.

```
.macro FORCE target, tosave:vararg
    ppush \tosave
    mov    rax, \target
    call   force.newabi
    mov    \target, rax
    ppop  \tosave
.endm
```

#### 9.1.5 ENAT

ENAT evaluates a register and converts it into nat, clobbing r11 and r12. tosave registers are preserved

```
.macro ENAT reg, tosave:vararg
    EVAL \reg, \tosave
    ncast \reg, \tosave
.endm
```

#### 9.2 Thunk Executioners

##### 9.2.1 xdone

xdone is used for thunks which have already been evaluated, it just returns the cached values, which is in the first slot.

Note that the garbage collector automatically recognizes thunks of this form, ignores any extra slots, and shrinks the result.

```
xdone:
    drefo rax, 8
    ret
```

##### 9.2.2 xhole

This is used for thunks which are in the process of being evaluated. This makes it possible to detect cases where an evaluation depends on its own result.

TODO: Should GC also shrink these down to one word, and not treat any of the slots as live?

```
xhole:
    mov    rdi, 1 # stdout
    lea    rsi, [loopstr]
    mov    rdx, 9
    call   syscall_write
    mov    rsi, 1
    jmp    syscall_exit_group
loopstr: .string "<<loop>>\n"
```

##### 9.2.3 xvar

This is an executioner for a dummy thunk which wraps another value with may also be a thunk. This is just used to handle certain edge cases that come up when implementing LETREC.

```
xvar:
    ppush rax
    drefo rax, 8
    ERAX
    ppop  rdi          # rdi=thunk
    ptr   rdi, PTR(thunk)
    mov   [rdi+8], rax
    lea   rsi, [xdone]
    mov   [rdi], rsi
    ret
```

##### 9.2.4 xhead

xhead is a specialized executioner for computing the Init of a closure. The size of the closure must be at least 2.

This is used during the implementation of Open primop, so that pattern matching on the host operating system and for directly poking the process itself.

```
xhead:
    ppush rax
    drefo rax, 8
    clsz_ rax, rax
    ppush rax
    dref_ rdi, rax
    call_ r12
    mov   \target, rax
    ppop  r10, rax
    lea   rsi, [xdone]
    mov   [r10], rdx
    mov   [r10+8], rax
    ret
```

##### 9.2.5 xunknown

This is a wrapper around xunknowndoupdate which performs thunk update.

```
xunknown:
    lea    rdx, [xhole]
    ptr   r10, rax
    mov   [r10], rdx
    ppush rax
    call   xunknowndoupdate
    pop   r10
    ptr   r10, rax
    lea   rsi, [xdone]
    mov   [r10], rdx
    mov   [r10+8], rax
    ret
```

##### 9.2.6 xunknowndoupdate

Some TODOs:

- Branch on the size hand small sizes using registers (not stack).
- Specialized entry-points for small sizes.

```
xunknowndoupdate:
thapple:
    ptr_  r8, rax      # r8 = ptr thunk
    mov   rax, [r8-8]   # fetch RECORD
    shr   rax, 14       # heapsz (in words)
    dec   rax          # size = heapsz-1 (ignoring fp)
    cmp   rax, 2
    je    thapple1
    cmp   rax, 4        # TODO: use a jump table (for arities up to 6)
    je    thapple2
    cmp   rax, 4
    je    thapple3
    mov   r11, rax
    mov   r11, [r11-1]
    jmp   apply1
thapple1:
    mov   rax, [r8+8]
    mov   rdi, [r8+16]
    mov   rsi, [r8+24]
    jmp   apply1
thapple2:
    mov   rax, [r8+8]
    mov   rdi, [r8+16]
    mov   rsi, [r8+24]
    jmp   apply2
thapple3:
    mov   rax, [r8+8]
    mov   rdi, [r8+16]
    mov   rsi, [r8+24]
    mov   rdx, [r8+32]
    jmp   apply3
    apply
```

##### 9.2.7 ToBit

The opposite of Nil is ToBit, which casts a value to 0 or 1.

```
optobit:
    EVAL rdi
    fastnil rdi
fasttobit:
    EVAL rdi
    fasttobit rdi
    # opntbl(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if not 0
    #   return
```

##### 9.3 Nat

Casts the input to a natural number.

```
.macro OPNAT reg, tosave:vararg
    EVAL \reg, \tosave
    fastnil \reg, \tosave
    ncast \reg, \tosave
.endm
```

##### 9.4.1 EVAL

EVAL Causes target to be evaluated, result in target, clobbers rax and r12. tosave registers are preserved

```
.macro EVAL target, tosave:vararg
    mov    rax, \target
    jnthk rax, 777f
    ppush \tosave
    dref_ r12, rax
    call_ r12
    mov    \target, rax
    ppop  \tosave
777:
.endm
```

##### 9.4.2 ERAX

This is a specialized version of EVAL, causes rax to be evaluated, clobbers r12. All of the registers in tosave are preserved

```
.macro ERAX tosave:vararg
    jnthk rax, 777f
    ppush \tosave
    dref_ r12, rax
    call_ r12
    ppop  \tosave
777:
.endm
```

##### 9.4.3 JRAX

JRAX evaluates and returns rax (uses a tail-call if rax is a thunk).

```
.macro JRAX
    jnthk rax, 777f
    dref_ r12, rax
    jmp    r12
777:
    ret
```

##### 9.4.4 FORCE

This evaluates the target to normal form.

```
.macro FORCE target, tosave:vararg
    ppush \tosave
    mov    rax, \target
    call   force.newabi
    mov    \target, rax
    ppop  \tosave
.endm
```

##### 9.4.5 ENAT

ENAT evaluates a register and converts it into nat, clobbing r11 and r12. tosave registers are preserved

```
.macro ENAT reg, tosave:vararg
    EVAL \reg, \tosave
    ncast \reg, \tosave
.endm
```

##### 9.5 Nil

We expose a quick check for zero equality, returning one if the input value is zero.

```
opnil:
    # opnil(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if 0
    #   return
```

##### 9.5.1 fastnil

fastnil is used for thunks which have already been evaluated, it just returns the cached values, which is in the first slot.

```
fastnil:
    # opntbl(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if not 0
    #   return
```

##### 9.5.2 fasttobit

fasttobit is used for thunks which are in the process of being evaluated. This makes it possible to detect cases where an evaluation depends on its own result.

```
fasttobit:
    # opntbl(rdi=val):
    #   evaluate
    #   fasttobit:
    #   set return value to 0
    #   possible return value of 1
    #   is val not zero?
    #   set return value to 1 if not 0
    #   return
```

##### 9.5.3 opnat

Numeric operations need to handle both direct numbers and indirect numbers, and they also need to evaluate their inputs if they are thunks.

```
.macro OPNAT target, tosave:vararg
    EVAL \target, \tosave
    fastnil \target, \tosave
    ncast \target, \tosave
.endm
```

##### 9.5.4 opntbl

However, the runtime itself implements XPLAN, which greatly extends this set of operations, as well as providing a number of mechanisms for interacting the host operating system and for directly poking the process itself.

```
opntbl:
    # opntbl(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if not 0
    #   return
```

##### 9.5.5 opnil

When running PLAN code, access to the lawful subset of these operations is done by doing jet matching in the online compiler, but in XPLAN these are invoked directly as primitives.

```
opnil:
    # opnil(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if 0
    #   return
```

##### 9.5.6 fastnil

fastnil is used for thunks which have already been evaluated, it just returns the cached values, which is in the first slot.

```
fastnil:
    # opntbl(rdi=val):
    #   evaluate
    #   fastnil:
    #   set return value to 0
    #   possible return value of 1
    #   is val?
    #   set return value to 1 if not 0
    #   return
```

##### 9.5.7 fasttobit

fasttobit is used for thunks which are in the process of being evaluated. This makes it possible







```

    mov      [r9 + rcx*8+8], rax          # r9[i+i] = item
    slots_  rax, rdi                   # get gcsz
    cmp     rax, 3                     # if (gcsz < 3)
    jb      oprow.break              # break
    drefo  rdi, 16                    # list = list.ix(1)
    inc    rcx                      # i++
    jmp    oprow.loop               # continue
oprow.break:
    mov    rax, r8                  # break:
oprow.ret:
    ret                           # res = buf
                                # ret:
                                # return res

```

## 13 Effectful Primops

### 13.1 PeekOp

TODO: there is no point in having an offset here, the caller can do the add themselves.

<pre> PeekOp:     mov    rsi, 3     call   unpackop.sized     ENAT  rdi, rsi, rdx     ENAT  rsi, rdi, rdx     ENAT  rdx, rdi, rsi     # fallthrough to fastpeekop </pre>	<pre> # rdi=[bufptr off sz] # # rdi = c dst pointer as nat # rsi = offset inside rdi # rdx = length in bytes to copy </pre>
--	---

#### 13.1.1 fastpeekop

```

ARG rdi = c src pointer as nat
ARG rsi = offset inside rdi
ARG rdx = length in bytes to copy

```

<pre> .global fastpeekop fastpeekop:     add    rdi, rsi     lea    rcx, [rdx+8]     shr    rcx, 3     ppush rdi, rdx, rcx     mov    rdi, rcx     call   reserve     ppop   rsi, rcx, r8     mov    byte ptr [rax+rcx], 1     mov    rdi, rax     rep    movsb     mov    rdi, r8     jmp    claim </pre>	<pre> # fastpeekop: rdi=src rsi=off rdx=n # pointless offset logic # words = (bytes+8)/8 # save dest/bytes/words # rax=buf=reserve(words) # rsi=ptr rcx=bytes r8=words # set high 1 byte # dst=buf # copy bytes (src=ptr rcx=bytes) # rdi=words # return claim(words) </pre>
--	--

### 13.2 PokeOp

And here is the actual primop which unpacks, evaluates, and casts all of the inputs.

TODO: the offset should be from the nat, instead of into the destination. The caller can trivially add to the pointer before they call us, but they can't easily get a pointer into the middle of a nat (we can).

```
    ENAT    rax, r  
    ENAT    rsi, r  
    ENAT    rdx, r  
    ENAT    rcx, r  
# fallthrough
```

13.2.1 fastpoke

```
ppush    rbx, rbp          # save rbx, rbp
add     rdi, rdx          # dst = bufAddr +
mov     rbx, rdi          # rbx=ptr
mov     rdi, rsi          # arg1=src
call    fastbytsz        # rax=bytes (clob
```

```
        mov    rbp, rcx          # rbp-extras
        subo   rbp, rax          # extra=(n-sz) (floor to 0)
        sub    rcx, rbp          # count -= extras
        mov    rdi, rbx          # dst=ptr
        add    rbx, rcx          # rbx=tailptr = ptr+count
        call   pokeraw          # pokeraw()
        mov    rdi, rbx          # dst=tailptr
        mov    rcx, rbp          # cnt=extras
        xor    eax, eax          # val=0
        rep    stosb             # memset(dst,cnt,0)
        ppop   rbx, rbp          # restore rbx/rbp
        ret                  # return 0
```

```
rep  
ppop  
ret
```

## Heaps and Process Layout

At the highest level, the PLAN runtime mmap's two regions:

1. A heap which is managed using a buddy allocator. The per-thread heaps, other memory allocated by the interpreter and all of the unpersisted pins.
2. In binaries created by `boot`, a file backed persistence heap which is stored in the program binary itself after the ELF region of the executable and which is mmap'd on startup.

This memory is garbage collected in different ways:

1. Each per-thread heap is a first generation, and is collected with a thread-local moving GC using Cheney's algorithm. A new per-thread heap is allocated from the buddy heap, live values are copied to it, and the old heap is explicitly freed from the buddy heap.
2. The buddy allocator heap is also garbage collected. On a collection of this second generation, each thread submits a list of roots to the buddy allocator, including interpreter objects like the heap and any pin objects referenced by the stack or the first generation heap.
3. The persistence file is also garbage collected. On a collection of this third generation, live regions are calculated, and then we request the kernel to punch a hole on pages between the live regions.

Right now, this collection only happens offline, but we hope to make this online soon.

## The Per Thread Heap

Each thread has its own heap, where the heap itself is allocated from the main buddy allocation arena.

### Cheney Sizes

The Cheney heap can be in only one of a set of prechosen heap sizes. We maintain a mapping from "heap rank" to sizes. The first 20 entries are the Fibonacci sequence starting with 12 and 38 starting at fib number 1 in that sequence, representing the start of the ramp up phase of the Cheney heap. After the first 20 entries we switch to 20% growth per term up to a heap of 583 megabytes, at which you should just crash. We allocate space for this table:

```
set HEAP_MAX_RANK ,      51

global heap_sizes
align 8
heap_sizes:
    .zero HEAP_MAX_RANK * 4
```

#### build\_heap\_table

When we startup, we must build a table of the appropriate heap bucket sizes. These are the sizes that we allocate during Cheney collection.

We follow the precedent of the Erlang BEAM VM by sizing our heap allocations by a modified Fibonacci sequence, starting in the same place: 12 and 38 words. Also following BEAM, we switch from doing Fibonacci

```
    -  
    mo  
    mo  
    mo  
build_heap  
    mo
```

```
        mov    edi, [heap_sizes + r8*4-8] # edi = heap_s
        add    esi, edi                 # esi += edi
        mov    [heap_sizes + r8*4], esi # hs[i] = hs[i]
        inc    r8                      # i++
        cmp    r8, 21                  # if (i < 21)
```

```

        js      build_heap_table.fibloop    # then keep doing fibonacci
build_heap_table.mulloop:
        mov    esi, [heap_sizes + r8*4-4]  # esi = heap_sizes[i - 1]
        lea    eax, [rsi + rsi*2]          # eax = rsi * 3
        add    eax, eax                  # eax = eax * 2 (or: rsi * 6)
        xor    edx, edx                # clear dividend

```

```
mov    ecx, 5          # div by 5
div    ecx
mov    [heap_sizes + r8*4], eax # eax = (6 * rsi) / 5 = 1.2 * rsi
inc    r8              # heap_sizes[i] = 1.2 * hs[i - 1]
cmp    r8, HEAP_MAX_RANK # if (i < HEAP_MAX_RANK)
js     build_heap_table.mulloop # then loop
ret
# return
```

- If the sum of

- If the sum of the surviving words and the needed words are over 50% of capacity of the current rank, then grow the target heap rank by at least one, increasing up until the sum fits, since the program could ask for much more memory than the current size.
- Otherwise, we're sized correctly between 25% and 50% of capacity and we keep the same target

Where `min_alloc` and `max_alloc` are the minimum and maximum values between the two `log2_size` ranges.

```
.global buddy_mmap_ptr, base_ptr, base_end_ptr  
buddy_mmap_ptr: .quad 0
```

```
base_ptr: .quad 0  
base_end_ptr: .quad 0
```

We'll document the different other structures as we go along. There are two main structures in the buddy allocator: the buckets and the nodes.

The “bucket number” of an allocation is the depth of the buddy tree where this allocation would go, and is used as a general identifier of size, even outside the bucket list.

Given a requested allocation size, round it up to the bucket that can hold that allocation.

TODO: This can actually be done without a loop with bsr.

```
.global bucket_for_request
bucket_for_request:
    shl    rdi, 3          # convert to bytes
    mov    rax, [bucket_count]
    dec    rax
    mov    rsi, [min_alloc]
```

```
bucket_for_request.loop:
    cmp      rsi, rdi          #
    jnb     bucket_for_request.done
    dec      rax
    shl      rsi, 1
    jmp     bucket_for_request.loop
bucket_for_request.done:
    ret
```

**16.3 GC Headers**

TODO: Figure out a better place for this to live.

Every node has a GC header. The tag bit on that header can be:

- 0b0000 - 0x0 - Nat
- 0b0001 - 0x1 - Pin
- 0b0010 - 0x2 - Law
- 0b0011 - 0x3 - ...

```
0b0000 - 0x0 - NIL (all)  
0b0100 - 0x4 - Clz (whnf)  
0b0101 - 0x5 - Thunk  
0b0110 - 0x6 - Megapin  
  
0b0111 - 0x7 - Heap  
0b1000 - 0x8 - Hash Table  
0b1001 - 0x9 - Bump Slab
```

For garbage collection purposes, we need to be able to traverse the buddy tree. We instead must represent three states:

```
.set UNALLOCATED, 0
.set INNER_NODE, 1
.set ALLOC_GC, 2
.set ALLOC_MANUAL, 3
```

There are two ways to refer to a node:

- Directly by pointer to the first byte of a memory range.
- By raw index into the linearized node tree.

You can switch between these two representations as long as you have the “bucket number”, the index which represents the depth of this allocation in the tree.

#### 16.4.1 nodeget

Since nodes in our tree can be in one of three states, our lookup function must parse out those two bits from the byte.

```
.global nodeget
nodeget:                                # nodeget(rdi=index):
    mov    ecx, edi                      # copy index for shift
    and    ecx, 0x03                     # ecx = index % 4
    shl    ecx, 1                        # shift = (index % 4) * 2
    shr    rdi, 2                        # byteIndex = index / 4
```

and  
ret

The core primitives of our buddy tree must be thread safe, so we write them using atomics. Individual values in the tree are owned by different threads, but the packing means that setting one of the four values in a byte needs to not stomp other thread's work, so we use a `cmpxchg` loop to ensure we don't overwrite the value in a different thread set.

```
.global nodeset
nodeset:
    mov    ecx, edi          # nodeset(rdi=index, rsi=tristate):
    and    ecx, 0x03         # copy index for shift
                                # ecx = index % 4
```

```

    shr    edi, 2                                # byteIndex = index / 4
    mov    r9d, 0x03                             # r9d = 2 bit mask
    shl    r9d, cl                               # shift mask into place
    and    esi, 0x03                             # constrain to tristate
    shl    esi, cl                               # shift state into place
    mov    rcx, [node_state_ptr]                 # load current base pointer
nodeset.try:
    mov    al, [rcx + rdi]                      # load current full byte
    mov    r9b, al                               # masking done

```

```
        and    r8b, r9b          # r8b = original & mask
        xor    r8b, al           # r8b = original with cleared bits
        or     r8b, sil          # r8b = new value inserted
lock cmpxchg [rcx + rdi], r8b      # CAS xchange al/[rcx+rdi]/r8b
jnz     nodeset.try          # retry if contended
ret                           # return
```

### 16.4.3 ptr\_for\_node

Given a node index and a bucket number, return the pointer to the actual node.

```
.global ptr_for_node
ptr_for_node:                                # ptr_for_node: (rdi=index, rsi=bucket)
    mov    rax, 1                           # base one
    shl    rax, rsi                         #    sh1 = (1 << rcx)
    inc    rdi                            #    + 1
    sub    rdi, rax                         #    index = index - (1 << bucket) + 1
    mov    rax, [max_alloc_log2]
    sub    rax, rsi
    shr    rdi, rcx                         # offset = index << (MAX_ALLOC_LOG2 - bucket)

    mov    rax, base_ptr
    add    rax, rdi                         # rax = base_ptr + offset
    ret
```

### 16.4.4 node\_for\_ptr

Given a pointer and a bucket number, calculate the node index.

```
.global node_for_ptr
node_for_ptr:                                # node_for_ptr: (rdi=ptr, rsi=bucket)
    mov    rax, rdi
    sub    rax, base_ptr                   # rax = ptr - base_ptr

    mov    rax, [max_alloc_log2]
    sub    rax, rsi
    shr    rax, rcx                         # rax ==> MAX_ALLOC_LOG2 - bucket

    mov    rdx, 1
    mov    rax, rsi
    shr    rdx, rcx                         # rdx = i << bucket

    add    rax, rdx
    dec    rax
    ret
```

## 16.5 List Utilities

Buddy tree free lists are represented inline, where the memory for the free space linked list node lives right at where the allocation will be in the future. The free list is a doubly linked list for fast insertion and removal.

Each list starts with the following struct:

```
typedef struct list_t {
    struct list_t* prev;
    struct list_t* next;
} list_t;
```

### 16.5.1 list\_init

The empty list points to itself.

```
.global list_init
list_init:
    mov    [rdi], rdi                      # list_init(rdi=list*):
    mov    [rdi+8], rdi                    #    list->prev = list;
    ret                                #    list->next = list;
                                         #    return
```

### 16.5.2 list\_push

```
.global list_push
list_push:
    ablptr rdi
    ablptr rsi
    mov    r8, [rdi]
    ablptr r8
    mov    [rsi], r8
    mov    [rsi+8], rdi
    mov    [r8+8], rsi
    mov    [rdi], rsi
    ret                                # list_push(rdi=list*, rsi=entry):
                                         #    input list pointer valid
                                         #    input new entry pointer valid
                                         #    list_t* prev = list->prev;
                                         #    entry->prev = prev;
                                         #    entry->next = list;
                                         #    prev->next = entry;
                                         #    list->prev = entry;
                                         #    return
```

### 16.5.3 list\_remove

Unlinks a list entry from its list.

```
.global list_remove
list_remove:
    ablptr rdi
    mov    r8, [rdi]
    ablptr r8
    mov    r9, [rdi+8]
    ablptr r9
    mov    [r8+8], r9
    mov    [r9], r8
    ret                                # list_remove(rdi=entry*):
                                         #    input entry pointer valid
                                         #    list_t* prev = entry->prev;
                                         #    prev pointer valid
                                         #    list_t* next = entry->next;
                                         #    next pointer valid
                                         #    prev->next = next;
                                         #    next->prev = prev;
                                         #    return
```

### 16.5.4 list\_pop

```
.global list_pop
list_pop:
    mov    r8, [rdi]                      # list_pop(rdi=list*) -> rax:
    cmp    r8, rdi                      #    list_t* back = list->prev;
    je     r80
    ablptr r8
    mov    r9, r8
    call   list_remove
    mov    rax, rdi                      #    assert pointer is in buddy heap
    ret                                #    set back as return value
                                         #    return
```

## 16.6 Aligned Bit Trees

We perform liveness marking against an aligned bit tree, which mirrors the structure of the buddy tree.

We have the logical structure where bit[0] is the head of a binary tree and its left child is  $2*i + 1$  and its right child is  $2*i + 2$ . The naive problem with this structure is that it is never word aligned. Meanwhile, we want to use this structure to be quadword aligned as much as possible so that we can do optimized walks over the data.

The compromise is that we just add one. Logical 0 is physical 1, logical 1 is physical 2, etc. This means all of level 0-5 sits in the first quadword, level 6 sits exactly in the second quadword, level 7 sits in the third and fourth quadword and so forth. All we have to do to achieve alignment is add a single bit.

### 16.6.1 aligned\_bit\_set

Atomically sets a specific bit, returning the previous value.

```
.global aligned_bit_set
aligned_bit_set:
    inc    rdi
    mov    rax, rdi
    shr    rax, 6
    and    rdi, 63
    lock bts QWORD PTR [rsi+rax*8], rdi
    setc   al
    movzx  eax, al
    ret
```

### 16.6.2 aligned\_bit\_get

```
.global aligned_bit_get
aligned_bit_get:
    inc    rdi
    mov    rax, rdi
    shr    rax, 6
    and    rdi, 63
    bt     QWORD PTR [rsi+rax*8], rdi
    setc   al
    movzx  eax, al
    ret
```

## 16.7 Buddy Collection

While our buddy allocator provides explicit `buddy_malloc` and `buddy_free` calls for the cases where a process can control the lifecycle of data, since user pins are placed on the heap, the entire buddy allocator must provide a garbage collector. Our collection is concurrent: the collector runs on its own thread, and will handle the three states as such:

```
.global collector_lock
collector_lock:
    .zero  8*7
```

`bucket_locks_ptr` will be allocated when as part of initialization to the right size of one u32 mutex for every bucket. The bucket locks are a per-bucket level lock needed to be taken before a thread reads or writes to the buckets or the `node_state` for that level.

When acquiring multiple locks, they should be acquired from the smallest allocation size (the highest bucket number) to the largest allocation size in the case of splitting.

```
.global bucket_locks_ptr
bucket_locks_ptr: .quad 0
```

The locking strategy is that you either take a write lock on `collector_lock`, or take a read lock on `collector_lock` and take a fine grained bucket lock. The sweeping process will take a coarse write `collector_lock`, while malloc will take a read `collector_lock` and will take fine grained locks for the buckets it needs to minimize contention inside malloc.

### 16.7.3 atomic\_max

Given a 64-bit word of child mark bits (packed as [child0, child1, child2, ... child63]), where the children for parent  $i$  are at positions  $2*i + 1$  and  $2*i + 2$ , compute for each parent: parent's bit = `child_left` OR `child_right`.

This is done by:

- Shifting a copy of the 64-bit word
- ORing it with the original 64-bit word.
- Using PEXT to extract even-indexed bits (mask 0x5555...5555)

`project_or` expects its output to be pre-zeroed. It takes advantage of this fact by skipping running any bit manipulation code on any zero input since it'll result in a zero output.

```
.equ ODD_MASK, 0x5555555555555555
.global project_or
project_or:
    push   r15
    mov    rsi, ODD_MASK
    xor    rax, rsi
    project_or.loop:
    mov    rax, [rsi+8*rcx]
    test   rax, rax
    jz     project_or.next
    mov    r8, rax
    shr    r8, 1
    or    rax, r8
    pext  r9, rax, r15
    mov    rax, [rdi+4*rcx], r9d
    ret
```

Once the `explore_tree` is filled, it "explores", meaning it looks at completed explore structure to find live items and then recurse into its dependencies. The sweep tree should have every live allocation reachable from the items marked in the `explore_tree`.

```
.global sweep_tree_ptr
sweep_tree_ptr: .quad 0
```

Because we don't want to perform a linear scan of the entire explore tree, which will cause the kernel to back those pages, we instead keep track of the highest index in the final bucket layer so we can bound scanning.

```
.global explore_max_idx
explore_max_idx: .quad 0
```

Once the `explore_tree` is filled, it looks at completed explore structure to find live items and then recurse into its dependencies. The sweep tree should have every live allocation reachable from the items marked in the `explore_tree`.

```
.global sweep_tree_ptr
sweep_tree_ptr: .quad 0
```

These three bitmaps take a lot of space, and we don't want to have to zero them out, paging them into memory. They are very sparse, so we keep record the size of the bitmap in bytes so we can pass it to the kernel for clearing.

```
.global bitmap_bytes
bitmap_bytes: .quad 0
```

Finally, once we've filled the sweep tree, we perform a sweep where we move step by step through the tree. The actual state of the sweeper is as a global variable because sometimes the `buddy_free` process must interact and change the sweeping state if the two are accessing the same location. `sweep_level` will be -1 when not running.

```
.global sweep_level
sweep_level: .quad -1

.global sweep_node_ptr
sweep_node_ptr: .quad 0

.global sweep_action_ptr
sweep_action_ptr: .quad 0
```

### 16.7.2 Locks

Garbage collection of the buddy heap is concurrent, so we need some locks.

First is the `collector_lock`. A reader-biased reader/writer lock where malloc is a reader and the collector is a writer. This is the main lock that has to be taken whenever reading or writing to the `mark_state`.

The design is that the collector should only progress through the tree while there are no mallocs, because a malloc could be writing to the `mark_state` while calculating dependencies, which it is responsible for doing if the collector thread is sweeping.

```
.global collector_lock
collector_lock:
    .zero  8*7
```

`bucket_locks_ptr` will be allocated when as part of initialization to the right size of one u32 mutex for every bucket. The bucket locks are a per-bucket level lock needed to be taken before a thread reads or writes to the buckets or the `node_state` for that level.

When acquiring multiple locks, they should be acquired from the smallest allocation size (the highest bucket number) to the largest allocation size in the case of splitting.

```
.global bucket_locks_ptr
bucket_locks_ptr: .quad 0
```

The locking strategy is that you either take a write lock on `collector_lock`, or take a read lock on `collector_lock` and take a fine grained bucket lock. The sweeping process will take a coarse write `collector_lock`, while malloc will take a read `collector_lock` and will take fine grained locks for the buckets it needs to minimize contention inside malloc.

```
.global sweep_level
sweep_level: .quad -1
```

```
.global sweep_node_ptr
sweep_node_ptr: .quad 0
```

```
.global sweep_action_ptr
sweep_action_ptr: .quad 0
```

At the end of this,  $i$ , bucket and every child node of it will be UNALLOCATED in the `node_state` and will not appear in any freelist for allocation. (The caller is responsible for doing something with  $i$ , bucket.)

TODO: The assembly version of this has been temporarily removed and reverted to the C `recursively_prune_c` while we rework the collector to be concurrent, and add more checking code to the process.

## 17 The Persistence Heap

### 17.1 Low Level Page Storage

#### 17.1.1 Mounting The Current Binary With Mmap

Normally, the Linux kernel will prevent a binary that's being executed from being opened for writing. If you try to open a writeable file descriptor, the kernel will return an ETXTBSY error. This means we have to go through some hoop that have to be performed right during a program's `_start` entry point.

```
.global explore_tree_ptr
explore_tree_ptr: .quad 0
```

Because we don't want to perform a linear scan of the entire explore tree, which will cause the kernel to back those pages, we instead keep track of the highest index in the final bucket layer so we can bound scanning.

```
.global explore_max_idx
explore_max_idx: .quad 0
```

Once the `explore_tree` is filled, it looks at completed explore structure to find live items and then recurse into its dependencies. The sweep tree should have every live allocation reachable from the items marked in the `explore_tree`.

```
.global sweep_tree_ptr
sweep_tree_ptr: .quad 0
```

These three bitmaps take a lot of space, and we don't want to have to zero them out, paging them into memory. They are very sparse, so we keep record the size of the bitmap in bytes so we can pass it to the kernel for clearing.

```
.global bitmap_bytes
bitmap_bytes: .quad 0
```

Finally, once we've filled the sweep tree, we perform a sweep where we move step by step through the tree. The actual state of the sweeper is as a global variable because sometimes the `buddy_free` process must interact and change the sweeping state if the two are accessing the same location. `sweep_level` will be -1 when not running.

```
.global sweep_level
sweep_level: .quad -1
```

```
.global sweep_node_ptr
sweep_node_ptr: .quad 0
```

```
.global sweep_action_ptr
sweep_action_ptr: .quad 0
```

The locking strategy is that you either take a write lock on `collector_lock`, or take a read lock on `collector_lock` and take a fine grained bucket lock. The sweeping process will take a coarse write `collector_lock`, while malloc will take a read `collector_lock` and will take fine grained locks for the buckets it needs to minimize contention inside malloc.

```
.global sweep_level
sweep_level: .quad -1
```

```
.global sweep_node_ptr
sweep_node_ptr: .quad 0
```

```
.global sweep_action_ptr
sweep_action_ptr: .quad 0
```

At the end of this,  $i$ , bucket and every child node of it will be UNALLOCATED in the `node_state` and will not appear in any freelist for allocation. (The caller is responsible for doing something with  $i$ , bucket.)

TODO: The assembly version of this has been temporarily removed and reverted to the C `recursively_prune_c` while we rework the collector to be concurrent, and add more checking code to the process.

### 17.1.2 Data Structures

#### 17.1.2.1 The Persistence Heap

When the concurrent garbage collector wants to collect, it pauses for every thread to mark its roots in the explore tree. Individual threads use `shallow_mark` to mark the bottom of this tree, and when everything is submitted, the sweep thread then uses `project_explore_tree` to make a binary bit tree used for exploration.

```
.global explore_tree_ptr
explore_tree_ptr: .quad 0
```

Because we don't want to perform a linear scan of the entire explore tree, which will cause the kernel to back those pages, we instead keep track of the highest index in the final bucket layer so we can bound scanning.

```
.global explore_max_idx
explore_max_idx: .quad 0
```

Once the `explore_tree` is filled, it looks at completed explore structure to find live items and then recurse into its dependencies. The sweep tree should have every live allocation reachable from the items marked in the `explore_tree`.

```
.global sweep_tree_ptr
sweep_tree_ptr: .quad 0
```

These three bitmaps take a lot of space, and we don't want to have to zero them out, paging them into memory. They are very sparse, so we keep record the size of the bitmap in bytes so we can pass it to the kernel for clearing.

```
.global bitmap_bytes
bitmap_bytes: .quad 0
```

Finally, once we've filled the sweep tree, we perform a sweep where we move step by step through the tree. The actual state of the sweeper is as a global variable because sometimes the `buddy_free` process must interact and change the sweeping state if the two are accessing the same location. `sweep_level` will be -1 when not running.

```
.global sweep_level
sweep_level: .quad -1
```

```
.global sweep_node_ptr
sweep_node_ptr: .quad 0
```

```
.global sweep_action_ptr
sweep_action_ptr: .quad 0
```

At the end of this,  $i$ , bucket and every child node of it will be UNALLOCATED in the `node_state` and will not appear in any freelist for allocation. (The caller is responsible for doing something with  $i$ , bucket.)

TODO: The assembly version of this has been temporarily removed and reverted to the C `recursively_prune_c` while we rework the collector to be concurrent, and add more checking code to the process.

### 17.1.3 The Persistence Heap

#### 17.1.3.1 The Persistence Heap

When the concurrent garbage collector wants to collect, it pauses for every thread to mark its roots in the explore tree. Individual threads use `shallow_mark` to mark the bottom of this tree, and when everything is submitted, the sweep thread then uses `project_explore_tree` to make a binary bit tree used for exploration.

```
.global explore_tree_ptr
explore_tree_ptr: .quad 0
```

Because we don't want to perform a linear scan of the entire explore tree, which will cause the kernel to back those pages, we instead keep track of the highest index in the final bucket layer so we can bound scanning.

Finally, here's the helper function that parses the size of the ELF file. The definition is based off the elf.h headers, which we don't want to bundle into our project as a build dependency, and just calculates the length of the ELF portion of the binary:

```
.global _elf_binary_end
_elf_binary_end:
    movzx   eax, WORD PTR _elf_header[rip+58]
    movzx   edx, WORD PTR _elf_header[rip+60]
    imul    eax, edx
    add     eax, DWORD PTR _elf_header[rip+40]
    ret
```

Finally, we have to store some read only constants used only in initialization:

```
.section .data
## ELF header of the executable loaded first thing on execution.
.align 16
elf_header:
.zero 64

.section .rodata
proc_self_exe:
.string "/proc/self/exe"
proc_fd_path:
.string "/proc/self/fd/3"
plan_interpreter:
.string "plan_interpreter"
empty_str:
.asciz ""
```

This execution trampoline is *only* used in binaries that have the persistence system compiled in. It is not included in rpm and the plan shell because it confuses gdb when you try to restart a binary. Those binaries use a small shim where `_start` just immediately calls `_sharedbegin`.

### 17.1.2 Persistence Format

Our persistence heap starts at the first 4096 aligned page after the end of the ELF data. Our file format is two pages for superblocks, and then a series of data pages. We hardcode a page size of 4096 because that's the most common page size, and we don't want the file format to diverge on systems where that isn't true.

We allocate two pages, each dedicated to one superblock structure. Each superblock is on its own page since filesystems only guarantees that individual pages get committed. We then mmap the entire area of the file past this point.

```
const long page_size = 4096;
const long page_size_words = page_size / 8;

#define MAGIC_NUMBER 0x31764e414c50 // "PLANvi" in ASCII

#pragma pack(push, 1)           // Save alignment and set to 1 byte alignment
struct superblock {
    uint64_t magic;           // Magic number to identify our file type
    uint8_t sequence;          // Sequence number, wraps around 8-bit
    uint32_t val_checksum;     // Expected checksum for val.
    Val val;                  // Tagged pointer into persistence heap.
    size_t write_offset;       // Current write position within the data area
    uint32_t checksum;         // Checksum of the superblock (excluding this field)
};
#pragma pack(pop)              // Restore previous alignment

struct superblock active_sb = {0};
off_t active_offset = 0;        // Binary location of current superblock
off_t inactive_offset = 0;      // Binary location of next superblock
```

We have two superblocks with sequence numbers and CRC32 checksums, so that if we fail during writing of a new superblock, we can detect that case on startup and use the previous superblock.

We compute that by taking the checksum of the rest of the packed struct.

```
uint32_t calculate_superblock_checksum(struct superblock *sb) {
    return calculate_crc32c((const char*)sb,
                           sizeof(struct superblock) - sizeof(uint32_t));
}
```

### 17.1.3 Persistence System Initialization

Each main function has the choice to start the persistence subsystem or not. If a binary does want to initialize the system, they call `init_persistence` with the file descriptor of the file they want to mount. This file descriptor is usually the `self_fd` descriptor that was opened during `_restart`, but can also be a target file when building an image for the first time.

`init_persistence` will calculate the locations of the superblocks, mount the region past the superblocks in an mmap, figure out which is the active superblock (if any), and record the bounds of the current persistence image.

```
void init_persistence(int fd) {
    persistence_fd = fd;

    struct stat st;
    if (syscall_fstat(persistence_fd, &st) < 0) {
        syscall_close(persistence_fd);
        die("Failed to get file size");
    }
    filesize = st.st_size;

    off_t binary_end = _elf_binary_end();
    off_t sb1_offset = round_up_align(binary_end, page_size);
    off_t sb2_offset = sb1_offset + page_size;

    // Data starts at the third page
    size_t data_start = sb1_offset + 2 * page_size;
    headersize = data_start;

    // Mount the data into memory. We do this before processing the superblocks
    // so we can do a top level item validation to make sure it was written
    // correctly.
    //
    // TODO: Lay out memory maps so they use all address space, collaborating
    // with the buddy allocator.
    const u64 rw = 3; // PROT_READ | PROT_WRITE
    const u64 map_flags = MAP_SHARED | MAP_FIXED;
    const u64 base_size = 1ULL << 39;
    persistence_start = syscall_mmap((u64*)base_addr,
                                      base_size, rw,
                                      map_flags, persistence_fd,
                                      data_start);

    // Read and check the validity of the two superblocks.
    struct superblock sb1, sb2;
    bool have_sb1 = false, have_sb2 = false;
    if (filesize >= sb1_offset + sizeof(struct superblock)) {
        have_sb1 = read_superblock(persistence_fd, sb1_offset, &sb1);
    }
    if (filesize >= sb2_offset + sizeof(struct superblock)) {
        have_sb2 = read_superblock(persistence_fd, sb2_offset, &sb2);
    }

    find_active_superblock(sb1, have_sb1, sb1_offset,
                          sb2, have_sb2, sb2_offset);

    persistence_cur = persistence_start + active_sb.write_offset;
    persistence_end = persistence_start + (base_size / 8);
}
```

### 17.1.4 Reading a Superblock

When we read or write a superblock, we do it by reading directly from the file, instead of reading or writing to a memory mapped page. (We are following what LMDB does here.)

We verify that we did read the right amount of data, verify it has the right magic number, verify that the checksum of the superblock is right, and verify that the checksum for the value the superblock points to is correct.

Returns true if everything is valid, false if anything is invalid.

```
bool read_superblock(int fd, off_t offset, struct superblock *sb) {
    ssize_t bytes_read = syscall_pread64(
        fd, (char*)sb, sizeof(struct superblock), offset);

    if (bytes_read != sizeof(struct superblock)) {
        if (bytes_read == 0 || bytes_read == -22 /* EINVAL */ ) {
            // This might be a new file or truncated file
            return false;
        }
        die("Failed to read superblock");
    }

    if (sb->magic != MAGIC_NUMBER) {
        printf("Invalid superblock (wrong magic number) at offset %ld\n",
               (long)offset);
        return false; // Block is not valid
    }

    uint32_t expected_checksum = sb->checksum;
    uint32_t actual_checksum = calculate_superblock_checksum(sb);
    if (expected_checksum != actual_checksum) {
        printf("Superblock corruption detected (checksum mismatch) at offset %ld\n",
               (long)offset);
        return false; // Block is not valid
    }

    expected_checksum = sb->val_checksum;
    actual_checksum = crc32_checksum_for(sb->val);
    if (expected_checksum != actual_checksum) {
        printf("Superblock corruption detected (content mismatch) at offset %ld\n",
               (long)offset);
        return false; // Block is not valid
    }

    return true;
}
```

### 17.1.5 Writing a Superblock

Writing a superblock calculates and sets the checksum of the rest of the block, and then writes and fsyncs the file descriptor to make sure the data is persisted.

```
void write_superblock(int fd, off_t offset, struct superblock *sb) {
    sb->checksum = calculate_superblock_checksum(sb);

    ssize_t bytes_written = syscall_pwrite64(
        fd, (char*)sb, sizeof(struct superblock), offset);
    if (bytes_written != sizeof(struct superblock)) {
        die("Failed to write superblock");
    }

    // Swap offsets.
    off_t tmp = active_offset;
    active_offset = inactive_offset;
    inactive_offset = tmp;
}
```

### 17.1.6 Determining Which Superblock To Use

On startup, in `init_persistence`, we have to initialize the state of which superblock is active and which one isn't. We look at which superblocks exist in the file, and initialize the state to the newest valid superblock, initializing a new superblock if there's no persistence data in the file.

```
void find_active_superblock(struct superblock sb1,
                            bool have_sb1,
                            off_t sb1_offset,
                            struct superblock sb2,
                            bool have_sb2,
                            off_t sb2_offset) {
    if (have_sb1 && have_sb2) {
        // Handle wraparound by checking if the difference (considering
        // unsigned overflow) is less than 128
        if ((uint8_t)sb1.sequence - sb2.sequence) < 128) {
            active_sb = sb1;
            active_offset = sb1_offset;
            inactive_offset = sb2_offset;
        } else {
            active_sb = sb2;
            active_offset = sb2_offset;
            inactive_offset = sb1_offset;
        }
    } else if (have_sb1) {
        active_sb = sb1;
        active_offset = sb1_offset;
        inactive_offset = sb2_offset;
    } else if (have_sb2) {
        active_sb = sb2;
        active_offset = sb2_offset;
        inactive_offset = sb1_offset;
    } else {
        // Write the initial superblock
        active_sb.magic = MAGIC_NUMBER;
        active_sb.sequence = 1;
        active_sb.val_checksum = natural_crc32(0);
        active_sb.val = 0;
        active_sb.write_offset = 0;
        write_superblock(persistence_fd, sb1_offset, &active_sb);

        active_offset = sb1_offset;
        inactive_offset = sb2_offset;
    }
}
```

### 17.1.7 Allocating Raw Pages

The low level primitive in the persistence file is the allocation of pages. Right now, we allocate a number of pages at the end of the file. This function only makes sure the file is large enough to contain the requested allocation (and calling ftruncate if it isn't) and

This function does no syncing. It just hands out a pointer to persistence backed memory.

```
u64* allocate_persistence_space_for(size_t sz) {
    size_t required_size = headersize
                        + ((persistence_cur - persistence_start) * 8)
                        + (sz * 8);

    if (filesize < required_size) {
        size_t new_size = round_up_align(required_size, page_size);
        if (syscall_ftruncate(persistence_fd, new_size) == -1) {
            die("Failed to extend file");
        }
    }

    filesize = new_size;

    // Our current pointer is always page aligned.
    u64* now = persistence_cur;
    persistence_cur += round_up_align(sz, page_size_words);
    return now;
}
```

### 17.1.8 Msyncing pages

The caller to `allocate_persistence_space_for` is responsible for msyncing the pages once they are written to.

```
#define MS_SYNC 4

void msync_region(u64* begin, u64 word_size) {
    // Round word_size turned into bytes to the nearest page.
    u64 rounded_size = round_up_align(word_size * 8, page_size);

    // Sync the region
    if (syscall_msync(begin, rounded_size, MS_SYNC) == -1) {
        die("Failed to msync data");
    }
}
```

## 17.2 Persisting PLAN values

### 17.3 Persistence Collection

### 17.4 Scrubing and Validation

All pins persisted have the crc32 bit set, meaning they have checksums to make sure they point to the right data and haven't been corrupted either by the system or by the garbage collector.

We calculate the checksums with the x86 crc32 instruction, which does Castagnoli CRC.

For natural numbers, we just do a single CRC instruction on the direct natural number.

```
u64 natural_crc32(u64 nat) {
    _int64_t crc = 0xFFFFFFFF;
    crc = _mm_crc32_u64(crc, nat);
    crc = 0xFFFFFFFF;
    return crc;
}
```

For general PLAN values, note that we run the CRC32 loop over the heap header one space before the pointer to the item:

```
u64 crc32_checksum_for(Val item) {
    if (((item > 63) == 0) {
        return natural_crc32(item);
    } else {
        u64* ptr = PTR(item);
        u64 bitsz = get_bitsz(item);
        u64 wordsz = (bitsz + 63) / 64;

        return calculate_crc32c((char*)(ptr - 1), (wordsz + 1) * 8);
    }
}
```

All usages of CRC32 are bound at load time by the following implementation, which is unrolled to operate in quads for efficiency.

```
void calculate_crc32c(struct superblock sb1,
                      bool have_sb1,
                      off_t sb1_offset,
                      struct superblock sb2,
                      bool have_sb2,
                      off_t sb2_offset) {
    if (have_sb1 && have_sb2) {
        // Handle wraparound by checking if the difference (considering
        // unsigned overflow) is less than 128
        if ((uint8_t)sb1.sequence - sb2.sequence) < 128) {
            active_sb = sb1;
            active_offset = sb1_offset;
            inactive_offset = sb2_offset;
        } else {
            active_sb = sb2;
            active_offset = sb2_offset;
            inactive_offset = sb1_offset;
        }
    } else if (have_sb1) {
        active_sb = sb1;
        active_offset = sb1_offset;
        inactive_offset = sb2_offset;
    } else if (have_sb2) {
        active_sb = sb2;
        active_offset = sb2_offset;
        inactive_offset = sb1_offset;
    } else {
        // Write the initial superblock
        active_sb.magic = MAGIC_NUMBER;
        active_sb.sequence = 1;
        active_sb.val_checksum = natural_crc32(0);
        active_sb.val = 0;
        active_sb.write_offset = 0;
        write_superblock(persistence_fd, sb1_offset, &active_sb);

        active_offset = sb1_offset;
        inactive_offset = sb2_offset;
    }
}
```

### 17.4.2 Validating Each Item

For each item, we check that its checksum matches the expected crc32 of the value. So to scrub a whole recursive pin structure, we set up an empty red-black tree with a bump allocator so that we can track the regions of memory we've already checked, and we then start the recursive process.

```
void validate_item(Val item, u64 expected) {
    _int64_t crc = crc32_checksum_for(item);
    if (crc != expected) {
        printf("CRC32 differ for %lx: expected=%lx, actual=%lx\n",
               item, expected, crc);
        die("scrub failed");
    }
}
```

17.4.3 Scrubbing An Entire Pin Tree

Our superblock structure records both a PLAN value and the expected crc32 of the value. So to scrub a whole recursive pin structure, we set up an empty red-black tree with a bump allocator so that we can track the regions of memory we've already checked.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a megapin. This leverages megapin structure so that we don't have to recur through the pin tree and can just linearly walk the megapin index.

```
void scrub_item(Val item, u64 expected, rb_tree* tree, bump_alloc_t* a) {
    if ((item > 63) == 0) {
        // Check a direct reference.
        u64 crc = natural_crc32(item);
        if (crc != expected) {
            printf("Header hash differs from expected for %lx\n", item);
            die("scrub failed");
        }
    }
}
```

scrub\_item is the recursive step, where we check each item in a tree of pins. It uses the same `mark_item_in_tree` tracking system used by the persistence collector, because pin dependencies are massively duplicated in practice.

This structure checks the validity of a single item, by checking that its own crc32 that proceeds it matches the expected checksum, that the item's bytes match the checksum, and then recurring based on whether this is a pin or a

```
    scrub_item(ptr[6 + i], pin_crcs[i], tree, a
}
}
}
```

```
ANG fsi -- The law to in  
RET rax -- The primop me  
RET rdx -- The primop ci  
RET r8 -- The law arity  
  
| match :  
|
```

```
jnkal    rax,  match.none          # if (body !~ 0[a b]) goto none
ix0_     rdx,  rax                # rdx = body.ix(0)
jnpin    rdx,  match.none          # if (rdx !~ <i>) goto none
unpin   rdx                      # rdx = rdx.item
ihasn   rdx,  match.none          # if (ldimset(rdx).note ==
```

```

    ix1      rax                      # x = body.ix(1)
    mov      rcx, r8                  # i = arity
match.loop:
    test     rcx, rcx                # if (i == 0)
    jz      match.head              # head
    jnkal   rax, match.none        # if !(x ~ 0[a b]) goto none
    ix1_    r9, rax                 # r9=b
    cmp      r9, rcx                # if (b != i)
    jne     match.none              # goto none

```

```
        jmp  
match.head:  
        jdir  
        jnqu  
        dref
```

```
    ret                                # return
match.noquote:
    cmp      rax, r8
    ja       match.found               # if (x > arity)
                                                # then goto found
match.none:
    xor     r8d, r8d                # none:
    xor     eax, eax                # arity=0
    xor     edx, edx                # key=0
    xor     edx, edx                # op=0
    ret                                # return
```

The native runtime has built in lightweight profiling events. A third event is used to start recording profiling events, and another to return the

Each profiling record is two values, where the first is a 64-bit integer containing information and 60-bits of unix time at 8ns resolution, and the second is a 64-bit integer representing the actual law value being run. This allows us to store the data directly in memory in a form that is a) fast and b) easy to parse.

The memory is laid out as a linked list of records. Therefore, the format for each profiling record is:

```
5 [prev_page a1 a2 b1 b2 c1 c2]
```

```
.section .data
.global profile_base, p
profile_base:    .quad 0
```

```
profile_end: .quad 0
```

## 19.2 Time Format

We choose a 60-bit 8ns resolution because that's good up until the year 2262 while being close to a resolution that's drowned out by kernel overhead.

We use the following macro to turn a timespec pointer filled by the kernel into our time format:

```
.macro PTIME60 dst, src
    mov    \dst, [\src+0]                      # rax = tv_sec (time_t)
    mov    r12, [\src+8]                       # r12 = tv_nsec (long)
    imul   \dst, \dst, 125000000                # ticks = tv_sec * (1e9 / 8)
    shr    r12, 3                            # tv_nsec >= 3
    add    \dst, r12                          # sum ticks
    mov    r12, 60                           # literal 60
    bzhi  \dst, \dst, r12                    # mask 60 bits.
```

### 19.3 Push Profile Page

ction and what gets called repeatedly to add a new item to the linked list.

```

set PROFILE_SIZE, 4088
push_profile:
    mov rdi, PROFILE_SIZE
    mov rsi, 0x004
    call buddy_malloc
    test rax, rax
    js push_profile.oom
push_profile.prepareclosure:
    mov r8, [profile_base]
    test r8, r8
    jz push_profile.write
    push rax
    call zero_remaining_profile
    pop rax
    mov r9, 0xD050000000000000
    or r8, r9
push_profile.write:
    mov qword ptr [rax], 5
    mov [rax + 8], r8
    mov [profile_base], rax
    mov rcx, rax
    add rcx, PROFILE_SIZE*8
    mov [profile_end], rcx
    add rax, 16
    mov [profile_next], rax
    ret
push_profile.oom:
    ud2

```

# push\_profile:  
# rdi = PROFILE\_SIZE  
# rsi = closure tag  
# buddy\_malloc(PROFILE\_SIZE, 4)  
# if malloc returned null  
# then handle oom  
# prepareclosure:  
# check current profile\_base  
# if profile\_base == NULL  
# then skip cleanup  
# save rax  
# zero out remaining profile page  
# restore rax  
# write correct header  
# or header with profile\_base ptr  
# write  
# Set hd to 5  
# 5[profile\_base a1 a2 b1 b2...]  
# save new profile\_base  
# end = start  
# end = start + size  
# set end  
# increment to first entry  
# set next write location

---

## 4 Zero Remaining Profile Page

don't want to initialize a full profile linked list closure at allocation time because that will destroy cache performance. We only zero the remaining unused structure when we have to either push another profile page in the reverse linked list structure, or when we have to hand this chain of pages to the program.

```

global zero_remaining_profile
zero_remaining_profile:
    mov rdi, [profile_next]
    mov rcx, [profile_end]
    sub rcx, rdi
    shr rcx, 3
    xor eax, eax
    rep stosq
    ret

```

# zero\_remaining\_profile:  
# rdi = profile\_next  
# rcx = profile\_end  
# rcx = remaining space  
# bytes to words  
# store 0  
# memset  
# return

---

## 5 recordevent

is the low level recording implementation for everything else. This takes an input of the tag and the value to write to the log.

```

## TODO: This calls the SYS_clock_gettime call. YOU REALLY DO NOT WANT TO DO
## THIS IN PRODUCTION CODE. This introduces SYSCALL OVERHEAD to a common
## action we're doing all the time. The correct thing to do is to call the
## VDSO version of this, but that's going to require parsing the auxiliary
## vector at startup, then do minimal ELF parsing to find lookup the pointer
## to '__vds(clock_gettime)'. That's all hard. So for getting from 0 to 1, we
## punt and do the wrong thing. But the data is going to be quasi-invalid
## until we do that.

global recordevent
recordevent:
    sub rsp, 32
    mov [rsp + 16], rdi
    mov [rsp + 24], rsi
recordevent.check:
    mov rax, [profile_next]
    add rax, 16
    mov rcx, [profile_end]
    cmp rax, rcx
    jae recordevent.overflow
recordevent.askfortime:
    mov rax, 228
    mov rdi, 1
    mov rsi, rsp
    syscall
recordevent.record:
    PTIME60 rax, rsp
    shl rax, 3
    mov rdi, [rsp + 16]
    or rax, rdi
    mov rcx, [profile_next]
    mov [rcx], rax
    mov rdi, [rsp + 24]
    mov [rcx + 8], rdi
    lea rcx, [rcx + 16]
    mov [profile_next], rcx
    add rsp, 32
    ret
recordevent.overflow:
    ppush rdi, rsi
    call push_profile
    ppop rdi, rsi
    jmp recordevent.askfortime

```

# recordevent(rdi=tag, rsi=2nd val)  
# stack space for timespec\*  
# save tag  
# save 2nd val  
# check free space:  
# get current next pointer  
# increment by 2 64bit words.  
# get end pointer  
# if (next >= profile\_end)  
# then handle overflow  
# ask kernel for time:  
# SYS\_clock\_gettime  
# CLOCK\_MONOTONIC  
# rsi = &timespec  
# call kernel  
# record record:  
# timespec to 60-bit 8ns time.  
# make room for tag  
# rdi = input tag  
# OR time with tag  
# get current pointer  
# low word: time+type tag  
# retrieve 2nd val  
# high word: law pointer  
# advance next pointer by 16  
# save updated pointer  
# restore stack frame  
#  
# overflow:  
# save input during allocation  
# new profile frame  
# restore input  
# continue in before path

---

## 6 oprofile

main generator of profiling events is the oprofile function, which is the jet of Profile function. The file function is defined as just returning its second argument.

the fast path, when profiling is disabled (ie there are no profiling frames to write to), the overhead is just fused compare-and-jump and one register move.

must absolutely prevent references to GC1 from leaking into GC2 because that breaks our garbage collection semantics. So we handle this by ignoring any Profile call that contains data on a thread's Cheney heap. All long term Law objects should be on GC2 or GC3.

```

oprofile:
    cmp qword ptr [profile_next], 0
    jne oprofile.start
oprofile.default:
    mov rax, rsi
    JRAX
oprofile.start:
    jdirect rdi, oprofile.savespace
    ptr_rcx, rdi
    cmp rcx, r14
    jb oprofile.checkincheney
oprofile.savespace:
    sub rsp, 16
    mov [rsp], rdi
    mov [rsp + 8], rsi
oprofile.recordbegin:
    mov rsi, rdi
    mov rdi, 0
    call recordevent
oprofile.eval:
    mov rsi, [rsp + 8]
    EVAL rsi
    mov [rsp + 8], rsi
oprofile.recordend:
    mov rsi, [rsp]
    mov rdi, 1
    call recordevent
oprofile.cleanup:
    mov rax, [rsp + 8]
    add rsp, 16
    JRAX
oprofile.checkincheney:
    cmp rcx, [heap_addr]
    jae oprofile.default
    jmp oprofile.savespace

```

# oprofile(rdi=law rsi=ret val):  
# if profiling is enabled  
# then jump to profile handling  
# default:  
# setup 2nd argument  
# evaluate 2nd argument  
# start:  
# if direct, start immediately  
# rcx = PTR(rdi)  
# if (1st < heap\_end)  
# then check if in cheney heap  
# save space:  
# stack space for two  
# save law pointer  
# save eval thunk  
# record begin:  
# rsi = "2nd val"  
# rdi = type tag 0 (begin)  
# recordevent(0, law ptr)  
# eval:  
# get saved thunk  
# evaluate it  
# put saved evaluation  
# record end:  
# rsi = "2nd val"  
# rdi = type tag 1 (end)  
# recordevent(0, law ptr)  
# cleanup:  
# retrieve result  
# undo stack space  
# evaluate 2nd argument  
# check in cheney:  
# if (PTR(rdi) >= heap\_addr)  
# then can't record value  
# otherwise continue recording

---

## 7 SetProfEnabledOp

profiling system is explicitly enabled and disabled by the using code. This XPLAN operation takes 0 or whether it's enabled.

```
        call      push_profile  
        ret  
SetProfEnabledOp.disable:  
        mov      qword ptr [pr  
        mov      qword ptr [pr
```

```
jne    GetProfUp.enabled      # then jump to profile handling
        mov    rax, 0           # else return 0
        ret                  # return
GetProfOp.enabled:
        call   zero_remaining_profile # zero the remaining profile
        mov    rax, [profile_base]  # get current closure base ptr.
        mov    r8, 0xD050000000000000 # write correct header
        or     rax, r8            # make rax tagged pointer
        mov    qword ptr[profile_base], 0 # reset profile_base stack
```

```
push    rax                      # save result
call    push_profile              # push a new profile buffer
pop    rax                      # restore result
ret
```

The first argument to the `plan` program is a seed to run which contains an XPLAN program. xseed programs have a calling convention similar to C programs: The seed file is loaded, the function is passed the rest of the command-line arguments as an array of nats, the result is evaluated, cast to a nat, and then passed to the `exit(2)` syscall.

This section will document the seed format, and the implementation of the seed loader.

This will attempt to communicate the basic idea of Seed without getting caught up in implementation concerns. The idea is that, if you understand the basic format layout, that should give you pretty strong intuitions about the how encoding and decoding work.

I will avoid talk about /why/ choices were made, since many choices are informed by implementation concerns, and I want to avoid getting caught up in that. Instead, I will only describe the format itself.

A seed file implements a template for a PLAN value, with a special number of parameters. These parameters will need to be passed into the loader routine.

Because TELN values often contain a lot of shared structure, this is represented explicitly in the format. The format contains a number of fragments which can refer to each-other, and the final fragment is the actual resulting value.

- The number of 8-bit numbers.
- The number of fragments.

- And then the data for each single-byte number.

### 20.2.4 Scopes and References

At each point in the seed, everything that came before is "in scope". All of the arguments, all of the numbers, and every preceding tree binding.

For a certain number of bindings in scope, we require a certain number of bits. For example, if there are four things in scope, each reference will require two bits.

parameters. This is encoded as:

- The size of the size in unary.
- The size in binary.
- The actual reference bits
- Zero or more following nodes, one per parameter.

Because the high-bit of the actual size is always 1, so the high bit is omitted in order to save space. The only edge-case is the case where there are zero parameters, but this is easy to handle.

Here are some examples of sizes and how they are encoded. Something tricky here is that these examples are listed LSB first, which means that the actual binary values need to be read in reverse.

```
0 -> 1.  
1 -> 01.  
2 -> 001.0  
3 -> 001.1  
4 -> 00010.00  
5 -> 00010.11  
6 -> 00010.01  
7 -> 00010.11  
8 -> 00011.000
```

The format only includes numbers and applications, so there actually isn't any way to encode pins and laws. However, the resulting template does not need to be in normal form, so we can just take `MkPin` and `MkLaw` as arguments. And we actually do not need to take `MkLaw` as an argument, because it can be constructed using just `MkPin`.

The specific conventions around which argument are passed in depends on the context.

When loading a boot seed, we pass in `MkPin` and we construct `MkLaw` ourselves. Since the calling convention for `MkLaw` is going to get more complicated, we likely should pass in `MkLaw` as well.

When loading a single node of a Pin DAG, we currently pass in `MkPin` but only use it to construct laws, since all of the sub-pins are passed in as references.

### 20.3.1 mmapfile

Given a path and a pointer to a statbuf (structure produced by fstat syscall), this opens a file, gets the size, and then loads the entire file into memory using mmap.

Returns the buffer, uses SysV calling conventions.

```
.global mmapfile
mmapfile:
    push r12          # rdi=path rsi=statbuf
    push rbx
    mov r12, rsi      # r12=statbuf
    xor esi, esi      # flags=0
    xor edx, edx      # mode=0
    call syscall_open_chk
    mov r12, rax
    mov rdi, rax
    mov rsi, r12      # statbuf = &seedstat
    call syscall_fstat_chk
    xor edi, edi
    mov [r12+48], edi # len=file.size
    mov ecx, 1         # prot=1
    mov r9, rbx
    xor r9d, r9d
    call syscall_mmap_chk
    mov r12, rax
    mov rdi, rbx
    call syscall_close_chk
    mov rax, r12
    pop rbx
    pop r12
    ret              # return
```

### 20.3.2 seedfile

Given a filename, load and return the seed, following SysV calling conventions.

```
.global seedfile # RPN testing system calls this.
seedfile:
    sub rsp, 160        # 144+padding
    mov rsi, rsp          # rsi = &stat (stack local)
    call mmapfile
    mov [rsp+152], rax
    mov rdi, rax          # save buf
    mov rsi, rax          # buf
    call seed             # seed(buf)
    mov rdi, [rsp+152]    # ptr = buf
    mov rsi, [rsp+48]     # len = stat.st_size
    mov [rsp+48], rax     # save result
    call syscall_munmap_chk
    mov rax, [rsp+48]     # munmap(ptr, len)
    add rsp, 160          # restore result
    add rsp, 160          # restore stack
    ret                  # return
```

### 20.3.3 seed

Given a pointer to a seed buffer, load the encoded value and return it via rax.

This just loads the header, loads all the sizes, and then calls into frags.

TODO: Document this.

```
seed:
    mov r12, rdi
    push r12
    call seed.inner
    pop r12
    ret

seed.inner:
    mov rcx, [r12]
    lea rsi, [r12 + 40]

seed.holeoop:
    test rcx, rcx
    jz seed.sizes
    xor edi, edi
    push rsi
    push rcx
    push r12
    call mkpin
    pop r12
    pop rcx
    pop rsi
    ppush rax
    dec rcx
    jmp seed.holeoop

seed.sizes:
    mov rcx, [r12 + 8]

seed.sizeoop:
    test rcx, rcx
    jz seed.bigs
    mov r8, [rsi]
    sub r15, 8
    mov [r15], r8
    add rsi, 8
    dec rcx
    jmp seed.sizeoop

seed.bigs:
    mov rcx, [r12 + 8]

seed.bigloop:
    test rcx, rcx
    jz seed.words
    dec rcx
    mov r8, [r15 + rcx*8]
    mov rdi, r8
    push rax
    push r8
    push r12
    call reserve
    pop r12
    pop r8
    mov rdi, r8
    call claim
    pop rsi
    pop rcx
    mov [r15 + rcx*8], rax
    jmp seed.bigloop

seed.words:
    mov rcx, [r12 + 16]

seed.wordoop:
    test rcx, rcx
    jz seed.bytes
    mov rax, [rsi]
    push r12
    call mkword
    pop r12
    sub r15, 8
    mov [r15], rax
    add rsi, 8
    dec rcx
    jmp seed.wordoop

seed.bytes:
    mov rcx, [r12 + 24]

seed.byteoop:
    test rcx, rcx
    jz seed.frags
    movzx r8, byte ptr [rsi]
    sub r15, 8
    mov [r15], r8
    inc rsi
    dec rcx
    jmp seed.byteoop

seed.frags:
    mov r8, rsi
    shr r8, 3
    shl r8, 3
    mov rdx, [r12 + 24]
    mov rdi, rdx
    add rdi, [r12 + 0]
    add rdi, [r12 + 8]
    add rdi, [r12 + 16]
    shl rdx, 61
    shr rdx, 58
    mov rdx, [r12 + 32]
    call frags
    mov rax, [r15]
    mov rcx, [r12]
    add rcx, [r12 + 8]
    add rcx, [r12 + 16]
    add rcx, [r12 + 24]
    add rcx, [r12 + 32]
    lea r15, [r15 + rcx*8]
    ret
```

### 20.3.4 frags

Basically, the just loads a certain number of tree fragments, and appends each one to the environment.

```
rdi - ARG - The total size of the environment table.
rsi - ARG - The environment (an array of plan values).
rdx - ARG - The number of used bits in the current input word.
rcx - ARG - The number of tree fragments to load.
r8 - ARG - A pointer to the current input word.
r10 - TMP - The number of fragments which have been read
r9 - TMP - Misc
```

This basically just calls frag() in a for loop, but the reference width may grow each time, so we need to recalculate that for each fragment.

All this function does is decode these inputs and uses them to construct a lazy function call:

```
((head arg1) arg2) arg3)
```

The logic here is pretty straightforward, and it's not that much code, but it is quite a lot to understand since there is a lot of state.

```
ARG rdi - refSz (bitSz of a reference into the environment)
ARG rsi - table (environment)
ARG rdx - ub (the number of bits that have been used in the current word).
ARG r8 - fp (pointer into the current word of the input buffer).
TMP rcx - temporary, or width-of-width
TMP r9 - temporary, or number of arguments.
TMP r10 - the next word from the input.
```

The arguments are essentially state that is threaded through the entire process. We don't write these to the stack except to avoid clobbering when calling mkapp().

The tricky bits here are just the decoding of the bit-packed input. The best way to understand this code is to first understand the format, and then work through this logic on paper with some small examples.

A couple of finer points to keep in mind:

- Closures sizes up to  $2^{32}$  (or something like that) fit in a single word. By not supporting that, the initial data for each node always fits in 64 bits, which means that all of our bit-decoding logic can work with registers.

- However, we need to read 64 bits of data at a bit-offset. In order to keep things simple and have less state, we just do this every time:

```
word = (p[0]<<n) | (p[1]<<n)
```

This means that we are re-reading the same words multiple times, but this isn't a major concern because of caching, and it keeps the complexity tractable.

- The handling of the zero-parameters case adds little code, but is a little bit tricky.

```
# state: rdi=refsz rsi=env rdx=used r8=ptr
# local: r10=word rcx=tmp r9=params
```

```
frags:
    push rbx
    push rbp
    push r12
    xor r12, r12
    mov rbp, rdi
    mov r12, rdx
    ret
```

```
frags.loop:
    cmp r12, r12
    jae frags.ret
    lea rdi, [rbp-1]
    lzcnt r9, rdi
    mov rdi, 64
    sub rdi, r9
    mov rsi, r15
    call frag
    shr r10, rdx
    mov r12, rdx
    or r10, r9
    jz frags.loop
    jne frags.head
```

```
frag.size:
    xor r9, r9
    tzcnt r10, r10
    shr r10, rdx
    inc r10, 1
    test r10, rdx
    jz frag.size
    sub r10, r9
    mov r15, r8
    add rsi, 8
    dec r10
    dec rdx
    jmp frags.loop
```

```
seed.args:
    dec rdx
    add rdx, rcx
    add rdx, rdx
    dec rdx
    xor r9, r9
    bts r9, r10
    and r9, r10
    bts r9, r10
    shr r10, rdx
    or r10, r9
    jz frag.size
```

```
frag.head:
    xor r10, r10
    bts r10, r10
    dec r10
    and r10, r9
    bts r10, r10
    shr r10, rdx
    or r10, r9
    jz frag.size
```

```
frag.args:
    dec rdx
    add rdx, rdx
    add rdx, rdx
    dec rdx
    xor r9, r9
    bts r9, r10
    and r9, r10
    bts r9, r10
    shr r10, rdx
    or r10, r9
    jz frag.size
```

```
frag.loop:
    test r9, r9
    jz frag.break
    mov [rsp], r9
    call frag.recur
    mov [rsp+8], r8
    call mapp()
    mov r8, [rsp]
    mov r9, [rsp+8]
    mov rdi, [rsp+16]
    mov rsi, [rsp+24]
    dec r9
    add r8, 8
    jne frag.loop
    jne frag.size
```

```
frag.break:
    add rsp, 32
    ret
```

```
# frag:
#   avoid clobbering r12
#   enter recursive logic
#   restore r12
#   return
```

```
# recur:
#   frame=(params, ptr, refs, env)
#   save refs
#   word = ptr[0]
#   if (used == 0)
#     eax = (shr only words with rdx)
#   word << used
```

```
#   remain = 64-used
#   next = ptr[1]
#   word = (word | next<<remain)
```

```
#   word = (word << used)
#   params = 0
#   szsz = ctz(word)
```

```
#   word >= szsz+1 (advance)
#   used++ (for 0 case)
#   if (!szsz)
#     goto head (0 case has no size)
```

```
#   args:
#   used-- (undo inc)
```

```
#   used += szsz*2
#   swid = (szsz-1)
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   head: (r10=word r9=params)
```

```
#   rax = word & ((2**refsz)-1)
#   ref = env[r10]
```

```
#   used += refs
#   word = (word << used)
#   word >= szsz+1 (advance)
#   used++ (for 0 case)
#   if (!szsz)
#     goto head (0 case has no size)
```

```
#   args:
#   used-- (undo inc)
```

```
#   used += szsz*2
#   swid = (szsz-1)
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid
```

```
#   params = word&((2**swid)-1)
#   set high bit
```

```
#   word >> swid</pre
```

### 21.2.4 condition\_wait

```
.global condition_wait
condition_wait:
    push r12          # save register
    mov r12d, [rdi + 8] # save current sequence var
    push rdi          # save input cvar*
    mov rdi, [rdi]
    call mutex_unlock # set up mutex (pointer)
    pop rdi          # restore cvar*
    push rdi          # save cvar*
    add rdi, 8        # rdi = cvar->seq
    mov rdi, r12d     # edx = sequence var expected val
    call futex_wait_private # call futex wait
    pop rdi          # restore cvar*
    mov rdi, [rdi]     # rdi is now mutex for rest of fun
condition_wait.wakeup_loop:
    mov eax, MUTEX_LOCKED_WAITERS # prepare locked+waiters state
    loc xchg DWORD PTR [rdi], eax # grab mutex, mark contended
    test eax, eax # if mutex was unlocked before
    jz condition_wait.done # then we now hold mutex
    push rdi          # save mutex pointer
    mov edx, MUTEX_LOCKED_WAITERS # expect mutex to be locked
    call futex_wait_private # call futex wait
    pop rdi          # restore mutex pointer
    jmp condition_wait.wakeup_loop # continue
condition_wait.done:
    pop r12          # done:
    ret             # restore r12
    # return
```

### 21.3 Read/Write Locks

Our reader writer locks are one 56-byte structure with the reader and writer structs inlined into the structure.

Offset	Size	Description
+0	unsigned int, 4 bytes	mutex
+4	4 bytes	padding (for 8-byte alignment)
+8	pointer, 8 bytes	readers_cv.mutex pointer
+16	unsigned int, 4 bytes	readers_cv.sequence var
+20	4 bytes	padding (for 8-byte alignment)
+24	pointer, 8 bytes	writers_cv.mutex pointer
+32	unsigned int, 4 bytes	writers_cv.sequence var
+36	4 bytes	padding (for 8-byte alignment)
+40	unsigned int, 4 bytes	active_readers
+44	unsigned int, 4 bytes	writer_active
+48	unsigned int, 4 bytes	waiting_writers
+52	4 bytes	padding (for 8-byte alignment)

Our reader/writer lock currently has a strong reader preference; this could lead to writer starvation. This is fine for now, because the only usage of rwlocks is communication between the buddy allocator and the concurrent garbage collector where the allocator allocating should starve forward progress in collection compared to starving allocation completion on work threads.

### 21.3.1 rwlock\_init

```
.global rwlock_init
rwlock_init:
    push r12          # save r12
    mov r12, rdi       # save rwlock pointer
    call mutex_init   # initialize front mutex
    lea rdi, [r12 + 8] # rdi = &rwlock->readers_cv
    lea rsi, [r12]     # rsi = &rwlock->writers_cv
    call condition_init # initialize readers cvar
    lea rdi, [r12 + 24] # rdi = &rwlock->writers_cvar
    call condition_init # initialize writers cvar
    mov DWORD PTR [r12+40], 0 # active_readers = 0
    mov DWORD PTR [r12+44], 0 # writer_active = 0
    mov DWORD PTR [r12+48], 0 # waiting_writers = 0
    mov rdi, r12       # restore rwlock pointer
    pop r12          # restore r12
    ret             # return
```

### 21.3.2 rwlock\_read\_lock

```
.global rwlock_read_lock
rwlock_read_lock:
    push r12          # save r12
    mov r12, rdi       # save rwlock pointer
    call mutex_lock   # lock our mutex
rwlock_read_lock.wait_loop:
    cmp DWORD PTR [r12 + 44], 0 # if writer_active == 0
    jz rwlock_read_lock.no_writers # then advance to next step
    lea rdi, [r12 + 8] # rdi = rwlock->readers_cv
    call condition_wait # wait on readers_cv
    mov rdi, r12       # restore rwlock in rdi
    jmp rwlock_read_lock.wait_loop # try again
rwlock_read_lock.no_writers:
    inc DWORD PTR [r12 + 40] # increment readers atomically
    mov rdi, r12       # rdi = rwlock*
    call mutex_unlock   # mutex_unlock(rwlock*)
    pop r12          # restore r12
    ret             # return
```

### 21.3.3 rwlock\_read\_unlock

```
.global rwlock_read_unlock
rwlock_read_unlock:
    push r12          # save r12
    mov r12, rdi       # save rwlock pointer
    call mutex_lock   # lock the mutex
    dec DWORD PTR [r12 + 40] # decrement active_readers
    jnz rwlock_read_unlock.done # if active_readers, then done
    cmp DWORD PTR [r12 + 48], 0 # then done
    jz rwlock_read_unlock.done # rdi = &mutex->writers_cv
    lea rdi, [r12 + 24] # signal one writer
    call condition_signal # signal one writer
rwlock_read_unlock.done:
    mov rdi, r12       # rdi = rwlock*
    call mutex_unlock   # unlock the mutex
    pop r12          # restore r12
    ret             # return
```

### 21.3.4 rwlock\_write\_lock

```
.global rwlock_write_lock
rwlock_write_lock:
    push r12          # save r12
    mov r12, rdi       # save rwlock pointer
    call mutex_lock   # lock the mutex
    inc DWORD PTR [r12 + 40] # increment waiting_writers
rwlock_write_lock.wait_loop:
    cmp DWORD PTR [r12 + 40], 0 # if there are active readers
    jnz rwlock_write_lock.wait # then wait
    cmp DWORD PTR [r12 + 44], 0 # if there's no writer active
    jz rwlock_write_lock.done # then finish up and take lock
rwlock_write_lock.wait:
    lea rdi, [r12 + 24] # rdi = &rwlock->writers_cv
    call condition_wait # wait on writers_cv
    mov rdi, r12       # restore rdi
    jmp rwlock_write_lock.wait_loop # try again
rwlock_write_lock.done:
    mov DWORD PTR [r12 + 44], 1 # mark writer as active
    dec DWORD PTR [r12 + 48] # decrement waiting_writers
    mov rdi, r12       # unlock mutex
    call mutex_unlock   # mutex_unlock(rdi)
    pop r12          # restore r12
    ret             # return
```

### 21.3.5 rwlock\_write\_unlock

```
.global rwlock_write_unlock
rwlock_write_unlock:
    push r12          # save r12
    mov r12, rdi       # save rwlock pointer
    call mutex_lock   # lock the mutex
    mov DWORD PTR [r12 + 44], 0 # mark writer as inactive
    cmp DWORD PTR [r12 + 48], 0 # if waiting_writers == 0
    jz rwlock_write_unlock.signal # then signal readers
    lea rdi, [r12 + 24] # rdi = &rwlock->writers_cv
    call condition_signal # signal one writer
    jmp rwlock_write_unlock.done # we woke the next writer
rwlock_write_unlock.signal:
    lea rdi, [r12 + 8] # rdi = &rwlock->readers_cv
    call condition_broadcast # broadcast to all readers
rwlock_write_unlock.done:
    mov rdi, r12       # rdi = rwlock*
    call mutex_unlock   # unlock the mutex
    pop r12          # restore r12
    ret             # return
```

### 21.4 Barriers

A barrier waits on a value to become 0. Used for count down tasks, usually where there's one waiter who farms out tasks to threads.

#### 21.4.1 barrier\_set

```
.global barrier_set
barrier_set:
    lock xchg [rdi], rsi # barrier_set(rdi=barrier*, rsi=cnt):
    ret                 # *barrier = count
    # return
```

#### 21.4.2 barrier\_done

```
.global barrier_done
barrier_done:
    mov eax, -1 # decrement by 1
    lock xadd dword ptr [rdi], eax # locked perform subtract
    cmp eax, 1 # if we aren't last thread
    jne barrier_done.done # then don't wake anyone
    mov rdi, FUTEX_WAKE_ALL # otherwise wake everyone
    call futex_wake_private # call futex
barrier_done.done:
    ret
```

#### 21.4.3 barrier\_wait

```
.global barrier_wait
barrier_wait:
    mov edx, [rdi] # current number of waiters
    test edx, edx # if this is zero
    jz barrier_wait.done # then we have don't have to wait
    push rdi # save barrier*
    call futex_wait_private # call futex_wait
    pop rdi # restore rdi
    jmp barrier_wait # try again
barrier_wait.done:
    ret
```

### 21.5 Threads

We use the following structure for a thread:

Offset	Size	Description
+0	unsigned int, 4 bytes	tid
+4	4 bytes	padding
+8	unsigned int, 4 bytes	futex_exit_flag
+12	4 bytes	padding
+16	pointer, 8 bytes	function pointer to run
+20	void*, 8 bytes	argument to pass to function
+24	void*, 8 bytes	return value storage
+32	void*, 8 bytes	stack pointer
+40	void*, 8 bytes	stack pointer
+48	size_t, 8 bytes	stack size

#### 21.5.1 thread\_create

Thread create takes three arguments: A pointer to empty memory for the thread struct in rdi, a function pointer to run on the thread in rsi, and an argument to pass to that function pointer in rdx.

The one thing that's subtle here is the handling of the start path. So clone3 will return 0 to the child thread, and the pid to the parent thread. When we jump into the child thread, we have to immediately set up the stack by creating a base stack frame and make sure we're stack aligned.

```
.equ CLONE_FLGS, CLONE_VM | CLONE_FS | CLONE_FILES | CLONE_SIGHAND | CLONE_THREAD
.equ EIGHTMB, 8 * 1024 * 1024

.global thread_create
thread_create:
    push rbp
    mov rbp, rsp
    push r12
    mov r12, rdi
    call mutex_lock
    mov qword PTR [r12 + 8], 0
    mov qword PTR [r12 + 16], rsi
    mov qword PTR [r12 + 24], rdx
    mov qword PTR [r12 + 32], 0
    mov qword PTR [r12+48], EIGHTMB
    mov rax, SYS_MMAP
    mov rdi, 0
    mov rsi, [r12 + 48]
    mov rdx, PROT_READ|PROT_WRITE
    mov r10, MAP_PRIVATE|MAP_ANONYMOUS
    mov r8, -1
    mov r9, 0
    syscall
    cmp rax, 0
    jne thread_create.error
    mov rdi, [r12 + 40]
    lea rsi, [r12 + 24]
    call condition_wait
    mov rdi, r12
    jmp thread_create.error
thread_create.error:
    leave
    ud2
    thread_create.start:
    xor rbp, rbp
    and rbp, -16
    mov rdi, [r12 + 24]
    push r12
    call [r12 + 16]
    mov r12, [r12 + 32]
    xor r10, r10
    xor r8, r8
    syscall
    cmp rax, 0
    jne thread_create.error
    je thread_create.start
    mov r12, [r12 + 40]
    pop r12
    mov rax, 0
    leave
    ret
thread_create.start:
    xor rbp, rbp
    and rbp, -16
    mov rdi, [r12 + 24]
    push r12
    call [r12 + 16]
    mov r12, [r12 + 32]
    xor r10, r10
    xor r8, r8
    syscall
    cmp rax, 0
    jne thread_create.error
    je thread_create.start
    mov r12, [r12 + 40]
    pop r12
    mov rax, 0
    leave
    ret
thread_create.error:
    leave
    kill
    start:
    clear frame pointer for new thrd
    ensure 16-byte stack alignment
    rdi = arg to function pointer
    save r12
    call function pointer
    restore r12
    set exit mark
    mark thread as exited
    rdi = atomic exit_flag
    wake a single waiter
    futex_wake()
    function return val as exit code
    and call exit
```

#### 21.5.2 thread\_join

Thread joining takes a pointer to a thread structure and optionally a pointer to memory to return a value from the thread.

```
.global thread_join
thread_join:
    cmp dword ptr [rdi + 8], 0 # if futex_exit_flag is not set
    jz thread_join.futex_wait # then wait for thread
    test rsi, rsi # if no result pointer
    jz thread_join.munmap_stack # then skip setting result
    mov rax, [rdi + 32] # get result
    mov rsi, [rdi + 40] # set result
    mov rdi, [rdi + 48] # munmap_stack:
    syscall
    # rdi = stack pointer
    # rsi = stack size
    # munmap(...)
```

```
thread_join.futex_wait:
    push rdi
    push rsi
    call futex_wait_private # futex_wait()
    pop rsi
    pop rdi
    jmp thread_join
```

rdi is a u32 memory location to wake, while edx is the number of waiters to wake up.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

rdi is a u32 memory location to wake, while rsi is the memory location to return a value from the thread.

## 22.2.2 Second Alternative

Another alternative, which continue to allow the high-byte to uniquely define heap reference types, could be:

```
u63 $$$$$$$$$$$$$$  
NAT 10000010ssssss small indirect nat, (s=bitsize)  
BAT 10000010ssssss big indirect nat, (s=bitsize)  
PAT 10000100ssssss pinned nat, (s=bitsize)  
PIN 100010000000cm (c=hascrc32, m=ismegapin)  
LAW 1001000000000000  
CLZ 1010000ttttzzzz (t=tag, z=size)  
THK 1100000000000000
```

## 22.3 Buffers / FATs

Right now, large, unpinned natural numbers live directly on the moving heap, but this has a number of downsides:

- This requires that this data be copied on every GC which is expensive for large buffers.
- We want to be able to store generated code for laws in natural numbers, and this code will need to not be moved by garbage collection.
- We want to be able to allocate temporary buffers in order to work with syscalls, and we need to be sure that those are not moved.

Here is a proposed sequence of proposals for introducing this concept and then integrating it deeper and deeper into the system.

### 22.3.1 xbuf

GC2 already supports top-level nat allocations, and the bucket-marking system used for all GC2 references should already work for marking such nats, so a basic version of this should be possible with very little change.

The heap layout and pointer tagging scheme would be totally conventional, just a normal NAT allocation, just allocated in GC2. Just introduce a few XPLAN primitives:

- `xbuf : Nat -> Nat` constructs a buffer of byte-size(n) as a Nat which is allocated as a top-level allocation at a fixed location.

This should crash if asked to allocate a size smaller than 8, the actual allocation should be one byte bigger than requested and the high byte should be set to 1 (making this a Bar), so that normal in-place NAT operations can be used to manipulate this.

- `xptr : Nat -> Nat` converts from a number in fixed memory (either a fixed nat, or a nat that lives inside of a pin), and converts it into a pointer.

This can be used with system calls to read data into a nat (even at an offset), or to write data from a nat (for example, to serialize data directly from pinned nat without needing to make a copy).

This can also be used if we dynamically generate code, to get a code pointer that can be attached as the judgement for a pin or a law.

- `xtouch: a -> b -> b` is used to guarantee that a temporary buffer is not freed until we are doing using a raw pointer into it. This is equivalent to Seq, but with the guarantee that it will not be optimized away, even if an optimizer can prove that the value has already been evaluated. (right now, there is no difference since we do not optimizer Seq, but we will eventually).

Significant RPN tests should be added which demonstrate that this all works correctly and is correctly collected.

This should be sufficient for basic usage, but it has a number of operational problems which we become a problem under heavier production use.

### 22.3.2 GC1 Collection

Heavy use of simple fixed buffers will result in a much higher allocation rate in GC2, and keeping this allocation rate low (through pin compaction) is a central pillar of our design.

However, fixed nats are not pins and cannot be referenced from pins or shared between actors, so they should in theory be collectable directly by GC1.

Here's an outline of how such a design could work:

- First, give fixed nats (FATs) their own pointer tag. There is enough left space for this in the high byte, and the `jnat/jnnat` checks are already range-checks, so they should be able to support multiple types at once without any change.
- Second, change the layout of FATs so that they have a mark bit and an intrinsic linked list. GC2 still requires a GC header on the outside, so the layout should be something like this:

```
| GC2HDR | mark | next | GC1HDR | word | word | word
```

- Reserve the first stack slot of the shadow (r15) stack to be a pointer to the FAT list. Every FAT allocation prepends itself to this list.

- During GC1 collection, when we encounter a reference to a FAT, set the mark.

- After GC1 collection, traverse the linked list of FATs, and free+delete any node which is not marked.

This scheme should add very little overhead, and the result is that FATs are freed promptly instead of clobbering up the heap between GC2 collections.

It's important to think through the synchronization issues between GC1 and GC2 collections very carefully, but intuitively this seems like it shouldn't be a problem.

- When GC2 is not running there is no sync issue.
- During initial mark collection, this scheme changes nothing.
- After initial mark collection, all old/reachable things are marked.
- During sweep, all allocations are allocated marked.
- The sweep traversal uses a lock that is shared with free(). So, if GC1 frees a FAT while GC2 is sweeping, that should not be a problem.

### 22.3.3 BigNats always FAT

In addition to the explicit `xbuf` operation, all large numbers should be allocated as FATs. This keeps the GC1 heaps small, which improves GC frequency, and avoids copying these large binary objects on every collection.

The conventional cut-off for this is 64 bytes.

This is conceptually trivial, but will be somewhat annoying in practice because it breaks the `reserve/claim` protocol. However, if we wait until we add support for nats that can be truncated in place, this problem should go away, since that change will avoid the need for this special dance altogether.

## 22.4 Actors

Our GC architecture is designed to be thread-safe, but we currently only use it with one thread at a time. Getting things to \*actually\* work with multiple threads at once will likely require quite a lot of effort.

## 22.5 Specialized Executioners

In addition to the unknown executioners or thunks, we should create specialized executioners for saturated (and oversaturated) calls to known operations.

Something like this, for example:

```
{THUNK(5), xknown2, &addop, Add, x, y}  
{THUNK(7), xknown2over, &addop, Add, 0, 1, 2, 3}
```

This would eliminate a lot of overhead, and wouldn't add all that much complexity. However, we will need a compiler before we can make use of this.

## 22.6 Freelists in Persistence

Right now, our persistence file format is append only, with garbage collection of regions only punching holes in the file, zeroing them and freeing the space. But this doesn't allow reuse of this address space, and very long running persistence files will run out.

After a garbage collection pass, we have a map of what regions of the file are free. We should write this map to the file and allocate from this freelist.

## 22.7 Register Reform

First, treating the C stack as GC roots will allow us to get rid of the r15 stack, freeing up another register. It might be worth using rbp as the heap pointer and rbx as the heap end, which would free up all of the numbered registers.

Second, the convention of using r12 as a scratch register conflicts very badly with its usual role as a callee-saved register, and another register should replace it for that purpose: maybe r11?

Third, we gain a lot of value from using custom calling conventions in a lot of places, but having reliable callee-saved registers would make a lot of things easier. These rules must be respected everywhere, or else you can't rely on them anywhere.

Together, this changes would make r12, r13, r14, and r15 available as callee-saved registers.

## 22.8 Killing the Shadow Stack

Killing the shadow stack would free up a register, and it would also mean that each actor has one less region to worry about.

In particular, the Erlang process model has a single allocation area per-actor which contains the stack and the heap, and they grow in opposite directions. This makes allocating a new actor super cheap, and makes it per-actor memory footprint very small for small/helper actors.

However, right now we use the RPN/shadow stack quite heavily. Most uses are just for register spilling, but we also use it as an actual stack in quite a few places. All of these systems would need to be reworked somehow.

- In the actual RPN debugger, used for the whole test suite. For the basic tests, we can probably find some reasonable way to write these tests directly in C/C++.

For more complex tests, we can probably just use seeds / planlisp / Sire.

- The seed loader uses the stack to represent environments, and also as a logical stack for mkapp. The former could just use a single explicit stack-frame, and the latter can be rewritten to use normal register conventions.

- The graph reduction engine uses the shadow stack to handle oversaturation. I'm not sure what the alternative would be, but it should be possible.

- In Judge, the environment and all nested sub-expressions are explicitly laid out on the stack. This will all go away with the new template expansion approach.

- There are likely a few more uses that have been overlooked by the above list.

One other major complication is that scanning the C stack is not at all safe. And this is going to make the interface between C/asn even trickier. We somehow need to avoid scanning C frames. We will either need an explicit roots system or to avoid GC1 allocation from C code at all.

## 22.9 Faster Evaluation Preludes

At present, almost all primops begin with a prelude which evaluates all of the strict arguments. Something like this:

```
# opcopy: rdi=cnt rsi=s0f rdx=d0f rcx=s0f r8=d0f  
opcopy:  
    NAT    rdi, rsi, rdx, rcx, r8      # eval+cast cnt  
    ENAT   rsi, rdi, rdx, rcx, r8      # eval+cast s0f  
    ENAT   rdx, rsi, rdi, rcx, r8      # eval+cast dof  
    EVAL   rcx, rsi, rdx, rdi, r8     # eval src  
    EVAL   r8,   rsi, rdx, rcx, rdi   # eval dst  
fastcopy:  
    ...
```

These evaluation macros EVAL and ENAT are optimized to do as little work as possible in the cases where the routine is not given a thunk. And they accomplish this by only flushing to the stack when \*actually\* given a thunk.

This works by examining each argument one by one. If it is a thunk everything is flushed to the stack, the relevant values are evaluated, and everything is restored from the stack. So, for the above example, if we are given 5 thunks, we will do something like 20 memory writes and 20 memory reads.

Unfortunately, in practice, the routines are almost always given thunks. The only time when we would be given something \*besides\* a thunk is when a law body passes in a hard-coded constant to a function. This happens, but it is very much not the common case.

Once we have an optimizing law compiler, law judgment itself will do a lot of evaluation directly before deferring to the graph reduction engine, and this will greatly increase the set of places where things besides thunks are passed in.

But such an optimizing law compiler should \*also\* be able to just perform the requisite evaluations itself, and then directly call into the fast path (`fastcopy` in this case), avoiding the entire prelude.

Given this combination of factors, it might be significantly more efficient to just flush everything to the stack, evaluate each item one-by-one, and the restore everything. In this case, the memory options would be reduced to 10 writes and 10 reads, which is significantly less work.

1. Save all five arguments to the stack.

2. Evaluate each argument, replacing the stack slot with the result.

3. Restore all of the registers from the stack.

And then, when using a more sophisticated judgment strategy, we would instead do the evaluation in the caller, and then directly invoke the fastpath, skipping this entire prelude.

## 22.10 Exit Segfaults

The code seems to segfault on CTRL-C, but only on the first run. What causes this?

## 22.11 GC2 Under Pressure

What's the best way to handle heavy load of GC2? Since our abstraction hides the memory hierarchy, having any sort of hard-limits on the size of data in each region is not great.

How can we avoid needing to make the buddy allocator resizable?

How can we make sure that heavy use of GC2 swaps, instead of killing the whole process?

## 22.12 GC2/GC3 Mutation

One of the invariants that we rely on in order to have a concurrent garbage collector without any write barriers is the complete lack of mutation in GC2 and GC3.

Because of the design of our system, this invariant is fairly easy to maintain. However, there are a few cases where it would be nice to violate it.

The big question is, can we sneak in these exceptions without losing the properties which make concurrent GC tractile?

Here are the edge-cases that would be nice to introduce:

- Redirections: If a GC2 pin is persisted, we would like to redirect all other references to the on-disk variant.

- Redirections: If we discover via hashing that two pins are identical, in either GC2 or GC3, we would like to be able to redirect one to the other in order to avoid duplication.

- Lazy Hashing: It would be nice to be able to wait to calculate the cryptographic hash of pins in GC2 until they are actually needed.

- Codegen: Since pins in GC2 and GC3 have stable locations, we can avoid needing to have a 'constants' table for generated code, and instead generate machine code which directly references values and code, by raw pointer.

- Moving code from GC2 to GC3. Once we attach generated code to pins and laws, the generated code will somehow need to be updated in order to patch-up the pointers to point to new locations.

- Compaction: It would be nice to be able to re-organize the heap on the fly, to reclaim parts of the address space and reduce fragmentation.

However, this would require some combination of the above things in order to work well.

Unfortunately, in practice, the routines are almost always given thunks. The only time when we would be given something \*besides\* a thunk is when a law body passes in a hard-coded constant to a function. This happens, but it is very much not the common case.

Once we have an optimizing law compiler, law judgment itself will do a lot of evaluation directly before deferring to the graph reduction engine, and this will greatly increase the set of places where things besides thunks are passed in.

But such an optimizing law compiler should \*also\* be able to just perform the requisite evaluations itself, and then directly call into the fast path (`fastcopy` in this case), avoiding the entire prelude.

Given this combination of factors, it might be significantly more efficient to just flush everything to the stack, evaluate each item one-by-one, and the restore everything. In this case, the memory options would be reduced to 10 writes and 10 reads, which is significantly less work.

1. Save all five arguments to the stack.

2. Evaluate each argument, replacing the stack slot with the result.

3. Restore all of the registers from the stack.

And then, when using a more sophisticated judgment strategy, we would instead do the evaluation in the caller, and then directly invoke the fastpath, skipping this entire prelude.

## 22.13 Separate Code Section

We've completely punted on this for now, since it introduces a lot of complexity. And that punting should continue for the foreseeable future.

However, ideally, we would have a separate area in GC2 and in GC3 for code. This way we don't have to map all of our data into executable pages. Also, code can be more compact if it can take advantage of the fact that it is collocated in the GC graph, and avoid relative jumps/references, etc.

However! This may introduce so much complexity that it is simply not worth it, despite the significant upsides.

## 22.14 Lazy Pins

We can't do this yet because this will definitely break compatibility with the Haskell runtime, but one major change which should result in a massive performance improvement is lazy pins and laws.

Right now, whenever we construct a Pin, we normalize the pin and immediately copy everything to GC2.

However, when doing something like inserting a page of JSON into a hutch database in the fulltag demo, we are constructing a huge pile of temporary pins which immediately become garbage.

In the new PLAN standard, pins are \*not\* guaranteed to be normalized. Also, we can separate out pin construction in GC1 from pin freezing (normalization and movement to GC2). The result will be that all of these ephemeral pins never make it to GC2 and never need to be copied.

At the end of each step, when we are ready to merge our final change into the resulting state, \*only then\* do we recursively freeze all of the pins in the result.

And it also means that we don't need to copy all of the data within a pin each time we construct a pin. The new version of a pin can share data with the previous version.

For example, if you do a single insert into the hutchhikers table and a new root.

I'm not sure if this is advisable