

Fast, precise dynamic checking of types and bounds “in C”

Stephen Kell

`stephen.kell@cl.cam.ac.uk`



Computer Laboratory
University of Cambridge

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker, (struct commit *)obj))  
        return -1;  
    return 0;  
}
```

Tool wanted

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker, (struct commit *)obj))  
        return -1;  
    return 0;  
}
```

↖ ↗
CHECK this
(at run time)

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker, (struct commit *)obj))  
        return -1;  
    return 0;  
}
```

↖ ↗
CHECK this
(at run time)

But also wanted:

- binary-compatible
- source-compatible
- ... for real, idiomatic code in (say) C
- reasonable performance

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker, (struct commit *)obj))  
        return -1;  
    return 0;  
}
```

↖ ↗
CHECK this
(at run time)

But also wanted:

- binary-compatible
- source-compatible
- ... for real, idiomatic code in (say) C
- reasonable performance

Enter libcrunch, which does the above.

The user's-eye view

- `$ crunchcc -o myprog ... # + other front-ends`

The user's-eye view

- `$ crunchcc -o myprog ...` # + other front-ends
- `$./myprog` # runs normally

The user's-eye view

- `$ crunchcc -o myprog ... # + other front-ends`
- `$./myprog # runs normally`
- `$ LD_PRELOAD=libcrunch.so ./myprog # does checks`

The user's-eye view

- `$ crunchcc -o myprog ... # + other front-ends`
- `$./myprog # runs normally`
- `$ LD_PRELOAD=libcrunch.so ./myprog # does checks`
- `myprog: Failed __is_a_internal(0x5a1220, 0x413560
a.k.a. "uint$32") at 0x40dade, allocation was a
heap block of int$32 originating at 0x40daa1`

Reminiscent of Valgrind (Memcheck), but different...

- not checking memory definedness, in-boundsness, etc..
- ... in fact, *assume correct* w.r.t. these!
- provide & exploit *run-time type information*

Sketch of the instrumentation for C

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker,  
  
        (struct commit *)obj))  
        return -1;  
    return 0;  
}
```

Sketch of the instrumentation for C

```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker,  
        (CHECK(_is_a(obj, "struct_commit")),  
        (struct_commit *)obj)))  
        return -1;  
    return 0;  
}
```

Sketch of the instrumentation for C

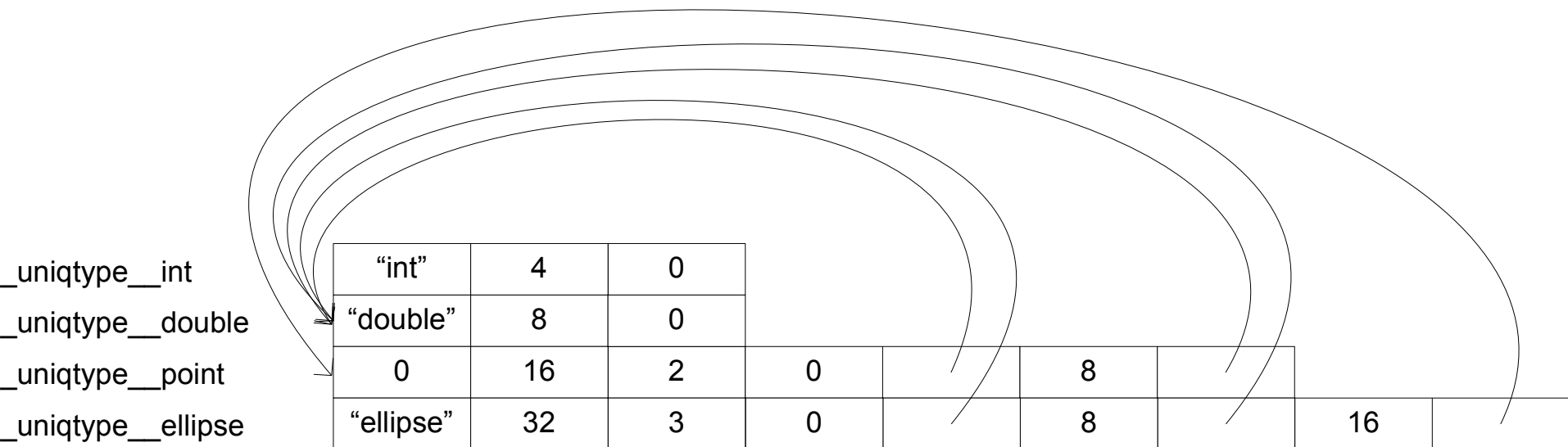
```
if (obj->type == OBJ_COMMIT) {  
    if (process_commit(walker,  
        (CHECK(__is_a(obj, "struct_commit")),  
        (struct_commit *)obj)))  
        return -1;  
    return 0;  
}
```

Need a runtime which

- provides a fast `__is_a()` function
- ... and a few other flavours of check
- by efficiently tracking *allocations*
- ... and attaching reified type info

Reified, unique data types

```
struct ellipse {  
    double maj, min;  
    struct point { double x, y; } ctr;  
};
```



- also model: stack frames, functions, pointers, arrays, ...
- unique → “exact type” test is a pointer comparison
- `__is_a()` is a short search over containment edges

Is it really that simple? What about...?

- untyped malloc() et al.
- opaque pointers, a.k.a. void*
- conversion of pointers to integers and back
- function pointers
- pointers to pointers
- “simulated subtyping”
- {custom, nested} heap allocators
- alloca()
- “sloppy” (non-standard-compliant) code
- unions, varargs, memcpy()

Is it really that simple? What about...?

- **untyped malloc() et al.**
- opaque pointers, a.k.a. void*
- conversion of pointers to integers and back
- function pointers
- **pointers to pointers**
- “simulated subtyping”
- {custom, nested} heap allocators
- alloca()
- “sloppy” (non-standard-compliant) code
- unions, varargs, memcpy()

What data type is being malloc()'d?

Use intraprocedural “sizeofness” analysis

```
size_t sz = sizeof (struct Foo);
```

```
/* ... */
```

```
malloc(sz);
```

Sizeofness propagates, a bit like dimensional analysis.

What data type is being malloc()'d?

Use intraprocedural “sizeofness” analysis

```
size_t sz = sizeof (struct Foo);
```

```
/* ... */
```

```
malloc(sz);
```

Sizeofness propagates, a bit like dimensional analysis.

```
malloc(sizeof (Blah) + n * sizeof (struct Foo))
```

What data type is being malloc()'d?

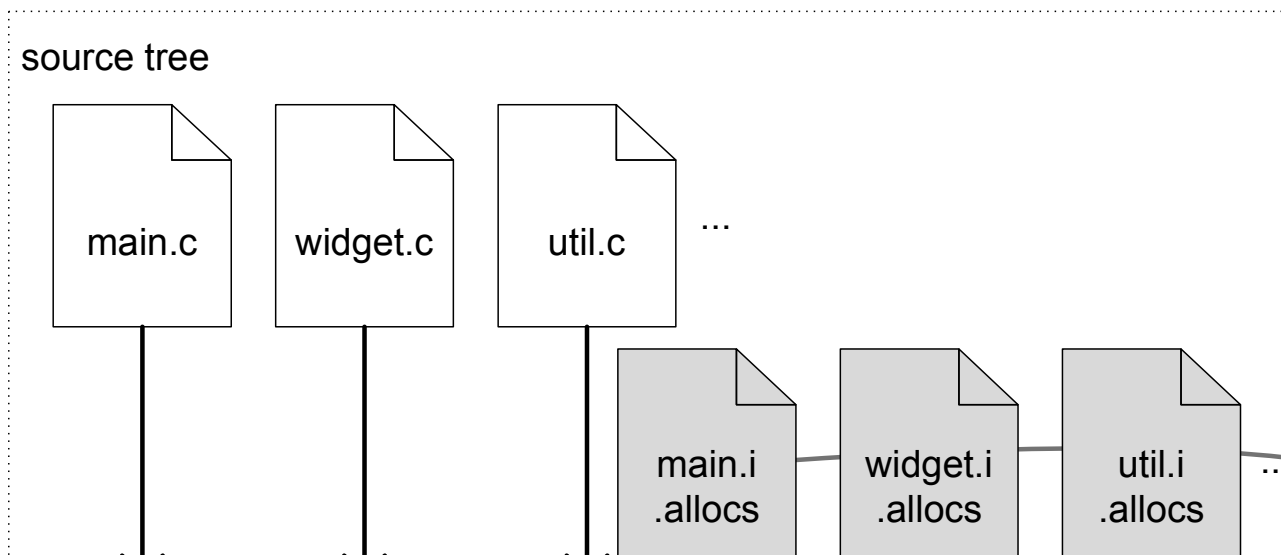
Use intraprocedural “sizeofness” analysis

```
size_t sz = sizeof (struct Foo);  
/* ... */  
malloc(sz);
```

Sizeofness propagates, a bit like dimensional analysis.

```
malloc(sizeof (Blah) + n * sizeof (struct Foo))
```

Dump *typed allocation sites* from compiler, for later pick-up



Polymorphism via multiply-indirected void

```
void sort_eight_special (void **pt){  
    void *tt [8];  
    register int i;  
    for(i=0;i<8;i++)tt[i]=pt[i];  
    for(i=XUP;i<=TUP;i++){pt[i]=tt[2*i]; pt[OPP_DIR(i)]=tt[2*i+1];}  
}  
neighbor = (int **)calloc(NDIRS, sizeof(int *));  
sort_eight_special ((void **) neighbor ); // <-- must allow!
```

- solution: tolerate casts from T^{**} to void^{**} ...
- and check *writes through* void^{**}
- ... *against the underlying object type* (here $\text{int}^{*[]}$)

Performance data: C-language SPEC CPU2006 benchmarks

bench	normal/s	crunch %	nopreload
bzip2	4.95	+6.8%	+1.4%
gcc	0.983	+160 %	- %
gobmk	14.6	+11 %	+2.0%
h264ref	10.1	+3.9%	+2.9%
hmmer	2.16	+8.3%	+3.7%
lbm	3.42	+9.6%	+1.7%
mcf	2.48	+12 %	(-0.5%)
milc	8.78	+38 %	+5.4%
sjeng	3.33	+1.5%	(-1.3%)
sphinx3	1.60	+13 %	+0.0%
perlbench			

Experience on “correct” code

benchmark	compile fixes	run-time false positives		
		instances	unique (of which...)	
			total	unhelpful
bzip2	0	48	3	3
gcc	1	3×10^5	14	3
gobmk	0	0	0	0
h264ref	2	27	2	0
hmmer	0	0	0	0
lbm	0	5×10^7	8	0
mcf	0	0	0	0
milc	0	0	0	0
sjeng	0	0	0	0
sphinx3	0	0	0	0

A “helpful” false positive?

```
typedef double LBM_Grid[SIZE_Z*SIZE_Y*SIZE_X*N_CELL_ENTRIES];  
typedef LBM_Grid* LBM_GridPtr;  
  
#define MAGIC_CAST(v) ((unsigned int*) ((void*) (&(v))))  
#define FLAG_VAR(v) unsigned int* const _aux_ = MAGIC_CAST(v)  
  
// ...  
  
#define TEST_FLAG(g,x,y,z,f) \  
    ((*MAGIC_CAST(GRID_ENTRY(g, x, y, z, FLAGS))) & (f))  
  
#define SET_FLAG(g,x,y,z,f) \  
{FLAG_VAR(GRID_ENTRY(g, x, y, z, FLAGS)); (*_aux_) |= (f);}
```

- check `memcpy()`, `realloc()`, etc..
- check syscalls (e.g. `read()`)
- add a bounds checker (improve on SoftBound)
- add a GC (precise! improve on Boehm)
- check unions and varargs
- always initialize pointers
- check unsafe writes through `char*`
- safely address-takeable union members (!)

Good prospects for all of the above! (ask me)

- check `memcpy()`, `realloc()`, etc..
- check syscalls (e.g. `read()`)
- **add a bounds checker (improve on SoftBound)**
- add a GC (precise! improve on Boehm)
- check unions and varargs
- always initialize pointers
- check unsafe writes through `char*`
- safely address-takeable union members (!)

Good prospects for all of the above! (ask me)

Plenty of existing tools do bounds checking

Memcheck (coarse), ASan (fine-ish), SoftBound (fine) ...

- detect out-of-bounds pointer/array use
- first two also catch some temporal errors

Problems remaining:

- overhead at best 50–100% (ASan & SoftBound)
- problems mixing uninstrumented code (libraries)
- *false positives and negatives*
- abort on false positives

Difficult-to-check bounds errors and non-errors

```
typedef struct {int x[2]; char y[2];} blah;
```

```
blah z = { {0, 0}, "!" };
```

```
*(z.x + 2);           // error: subobject overflow
```

```
((int*) &z)[2];      // error: after bounds-narrowing cast
```

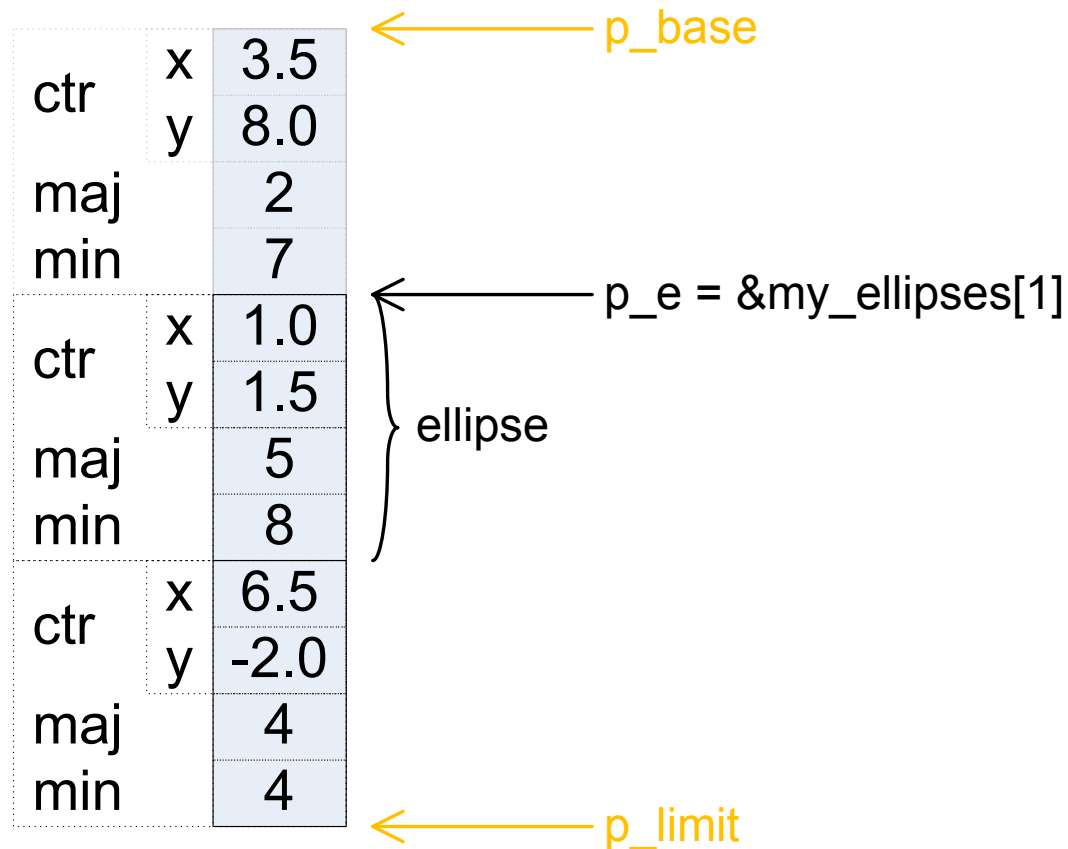
```
*((z.x + 42) - 42);  // non-error: via invalid (OOB) intermediate
```

```
((blah *) z.x) ->y; // non-error: after bounds-widening cast
```

```
*(int*)( intptr_t )z.x; // non-error: via integer
```

```
* strfry (z.y);      // non-error: after uninstrumented code
```

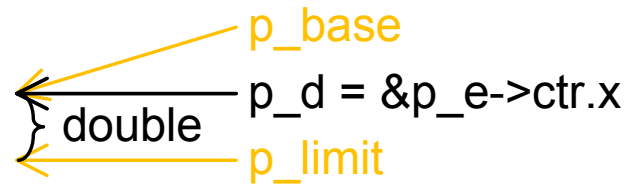
Existing checkers using per-pointer metadata

ctr	x	3.5	
	y	8.0	
maj	2		
min	7		
ctr	x	1.0	
	y	1.5	
maj	5		
min	8		
ctr	x	6.5	
	y	-2.0	
maj	4		
min	4		

```
struct ellipse {
    struct point {
        double x, y;
    } ctr;
    double maj;
    double min;
} my_ellipses[3];
```

Existing checkers using per-pointer metadata

ctr	x	3.5
	y	8.0
maj		2
min		7
ctr	x	1.0
	y	1.5
maj		5
min		8
ctr	x	6.5
	y	-2.0
maj		4
min		4



```
struct ellipse {  
    struct point {  
        double x, y;  
    } ctr;  
    double maj;  
    double min;  
} my_ellipses[3];
```

Without type information, pointer bounds lose precision

ctr	x	3.5
	y	8.0
maj		2
min		7
ctr	x	1.0
	y	1.5
maj		5
min		8
ctr	x	6.5
	y	-2.0
maj		4
min		4

```
struct ellipse {
    struct point {
        double x, y;
    } ctr;
    double maj;
    double min;
} my_ellipses[3];
```

Given allocation type and pointer type, bounds are implicit

ellipse[3]

ctr	x	3.5
	y	8.0
maj		2
min		7
ctr	x	1.0
	y	1.5
maj		5
min		8
ctr	x	6.5
	y	-2.0
maj		4
min		4

p_e = &my_ellipses[1]

ellipse

```
struct ellipse {
    struct point {
        double x, y;
    } ctr;
    double maj;
    double min;
} my_ellipses[3];
```

Given allocation type and pointer type, bounds are implicit

double { ellipse[3]	ctr	x	3.5
		y	8.0
	maj		2
		min	
	ctr		x
		y	1.5
	maj		5
		min	
	ctr		x
y		-2.0	
maj		4	
	min		4

← double p_d = &p_e->ctr.x

```
struct ellipse {  
    struct point {  
        double x, y;  
    } ctr;  
    double maj;  
    double min;  
} my_ellipses[3];
```

Given allocation type and pointer type, bounds are implicit

ellipse[3]

ctr	x	3.5
	y	8.0
maj		2
min		7
ctr	x	1.0
	y	1.5
maj		5
min		8
ctr	x	6.5
	y	-2.0
maj		4
min		4

p_f = (ellipse*) p_d

ellipse

```
struct ellipse {
    struct point {
        double x, y;
    } ctr;
    double maj;
    double min;
} my_ellipses[3];
```


The importance of being type-aware (when bounds-checking)

```
struct driver    { /* ... */ } *d = /* ... */;  
struct i2c_driver { /* ... */ struct driver driver; /* ... */ };  
  
#define container_of(ptr, type, member) \  
    ((type *) ( char *(ptr) - offsetof(type,member) ))  
  
i2c_drv = container_of(d, struct i2c_driver, driver);
```

The importance of being type-aware (when bounds-checking)

```
struct driver    { /* ... */ } *d = /* ... */;  
struct i2c_driver { /* ... */ struct driver driver; /* ... */ };
```

```
#define container_of(ptr, type, member) \  
((type *) ( char *) (ptr) - offsetof(type, member) )
```

```
i2c_drv = container_of(d, struct i2c_driver, driver);
```

SoftBound is oblivious to casts, even though they matter:

- bounds of `d`: just the smaller struct
- bounds of the `char*`: the whole allocation
- bounds of `i2c_drv`: the bigger struct

If only we knew the *type* of the storage!

Write a bounds-checker consuming per-allocation metadata

- avoid these false positives
- avoid `libc` wrappers, ...
- robust to uninstrumented callers/callees
- performance?

Making it fast:

- cache bounds: make pointers “locally fat, globally thin”
- only check *derivation*, not *use*

On x86-64, use noncanonical addresses as trap reps

FFFFFFFF FFFFFFFF

Canonical "higher half"

FFF8000 00000000

Noncanonical
addresses

00007FFF FFFFFFFF

Canonical "lower half"

00000000 00000000

No more deref checking

```
int ret = 0;
for (int i = 0; i < n; ++i)
{ struct list_node *p = malloc(sizeof (struct list_node ));
  p->next = head;
  head = p;
}
for (int i = 0; i < m; ++i)
{ unsigned out = 0;
  for (struct list_node *p = head; p; p = p->next)
  { out += p->x; }
  ret += out;
}
return ret;
```

It “works”!

- still turning the handle...

Early indications

- array-based programs: SoftBound-like perf
 - ◆ mostly +50–100%
- linked structures: much faster!
 - ◆ as expected
- building my own SoftBound-alike to compare...

Run-time type info enables efficient and helpful checking

- source- and binary-compatible
- low overhead type checking
- precise & fast bounds checking
- good prospects for extension

Code is here: <http://github.com/stephenrkell/libcrunch/>

Thanks for your attention. Questions?