

Cerberus: towards an Executable Semantics for Sequential and Concurrent C11

Kayvan Memarian, Kyndylan Nienhuis, Justus Matthiesen,
James Lingard, Peter Sewell

May 5, 2015

1 - What is C?

Different possible answers:

- ▶ what is described in the **standard**
⇒ ISO/IEC 9899-2011 + various defect reports
- ▶ various “de facto” definitions:
 - ▶ what compilers assume
 - ▶ what programmers believe
 - ▶ what different corpuses of C code actually require to run
- ▶ what analysis and verification tools assume (e.g. CompCertC)

2 - What we want

A formalisation of C11 which:

1. focuses on the various **“de-facto” Cs** (to study existing C codes)
2. is **parametric** (independence from impl. choices/interpretations)
3. supports **C/C++11 concurrency**
4. could be recognisable by those familiar with the C standard

3 - Why C is hard

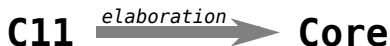
C11 expressions hide a lot their complexity:

- ▶ loose and intricate ordering (**sequence-before relation**)
- ▶ hidden occurrence of memory operations (**boundary of object lifetime**)
- ▶ implicit type conversions (**usual arithmetic conv; integer promotions**)
- ▶ partiality (**undefined behaviour**)
- ▶ parametricity (**implementation-defined choices**)

- ▶ real code uses features outside of ISO C

Direct formalisation of a complex language leads to **heavy and quickly intractable semantic states**.

4 - Semantics by elaboration



- ▶ the elaboration function *explains* each C11 construct
- ▶ Core has *simple constructs* (albeit quite a few)
- ▶ the elaboration produces verbose Core programs but with *explicit behaviour*

4 - Semantics by elaboration

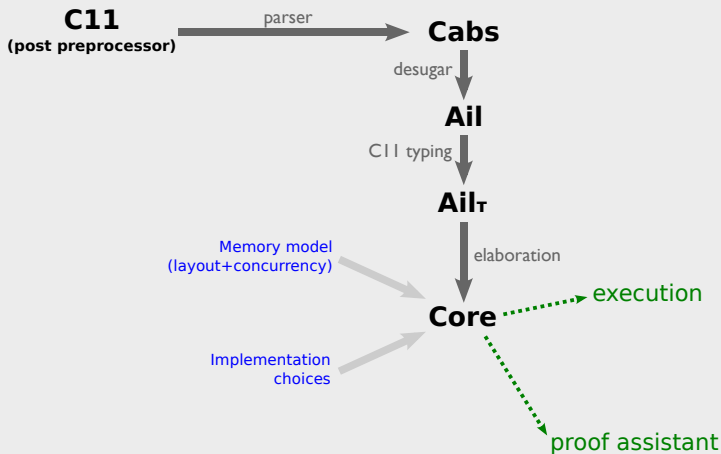
```
int main(void) {  
    int x = 0, y = 1;  
    return x+y;  
}
```

4 - Semantics by elaboration

```
proc main(): eff integer :=
  let strong (x, y) = unseq(create(<alignof>("signed int"), "signed int"),
                           create(<alignof>("signed int"), "signed int")) in
  let strong ()      = store("signed int", x, conv_int("signed int", 0)) in
  let strong ()      = store("signed int", y, conv_int("signed int", 1)) in

  let strong a3 =
    let weak (a1, a2) = unseq(load("signed int", x), load("signed int", y)) in
    overflow("signed int", conv_int("signed int", a1) + conv_int("signed int", a2))
  in
  let strong () = unseq(kill(x), kill(y)) in
  return (conv_int("signed int", a3))
```

5 - Cerberus' structure



Cerberus

6 - The Core language: *some features*

- ▶ functional style with first-order recursive functions
- ▶ simple discrimination between pure and effectful computations
- ▶ explicit memory actions (including object lifetime boundary)
- ▶ a calculus to express the *sequenced-before* relation:

```
let weak (a1, a2) = unseq(E1, E2) in  
store(τ, ptr, a1 + a2)
```

- ▶ C types and implementation-defined constants as values:

```
<Integer.conv_nonrepresentable_signed_integer>("signed int", n)
```

- ▶ explicit undefined behaviours and dynamic-runtime errors:

```
if a2 = 0 then undef(Division.by_zero) else a1 / a2
```

- ▶ others: continuation operators, non-determinism, I/O, ...

6 - The Core language (1/2)

values:

$v ::= \text{unit} \mid \text{null} \mid \text{true} \mid \text{false} \mid [v_1, \dots, v_n] \mid (v_1, \dots, v_n)$
 $\mid \tau \mid \text{unspec}(\tau) \mid n \in \mathbb{Z} \mid \langle \text{ptr} \rangle \mid \langle \text{array-const} \rangle$
 $\mid \langle \text{struct-const} \rangle \mid \langle \text{union-const} \rangle$

pure expressions:

$\text{pe} ::= \text{undef} \mid \text{error} \mid v \mid a \mid \langle \text{impl-name} \rangle \mid \text{cons}(\text{pe}_1, \text{pe}_2)$
 $\mid \text{case_list}(\text{pe}_1, \text{pe}_2, \text{nm}) \mid \text{case_ctype}(\text{pe}_1, \text{pe}_2, \text{nm}_1, \dots, \text{nm}_8)$
 $\mid \text{shift}(\text{pe}, \langle \text{shift-path} \rangle) \mid \text{not}(\text{pe}) \mid \text{pe}_1 \circ \text{pe}_2$
 $\mid \langle \text{mem-op} \rangle(\text{pe}_1 \dots \text{pe}_n) \mid (\text{pe}_1, \dots, \text{pe}_n) \mid \text{nm}(\text{pe}_1 \dots \text{pe}_n)$
 $\mid \text{let } a = \text{pe}_1 \text{ in } \text{pe}_2 \mid \text{if } \text{pe}_1 \text{ then } \text{pe}_2 \text{ else } \text{pe}_3$

$\text{nm} ::= \text{SYMBOL} \mid \text{``impl-def function names''}$

$\langle \text{shift-path} \rangle ::= \{ (\tau_1, \text{pe}_1), \dots, (\tau_n, \text{pe}_n) \}$

6 - The Core language (2/2)

$E ::= pe \mid \mathbf{let} \ a = pe_1 \ \mathbf{in} \ E_2 \mid \mathbf{if} \ pe_1 \ \mathbf{then} \ E_2 \ \mathbf{else} \ E_3 \mid nm\{pe_1, \dots, pe_n\}$
| $\mathbf{nd}(E_1, \dots, E_n) \mid \mathbf{return}(pe)$
| $A \mid \bar{A} \mid \mathbf{skip}$
| $\mathbf{unseq}(E_1, \dots, E_n) \mid \mathbf{let \ weak} \ (a_1, \dots, a_n) = E_1 \ \mathbf{in} \ E_2$
| $\mathbf{let \ strong} \ (a_1, \dots, a_n) = E_1 \ \mathbf{in} \ E_2 \mid \mathbf{let \ atomic} \ a = A_1 \ \mathbf{in} \ P$
| $\mathbf{indet}_i(E) \mid \mathbf{bound}_i(E)$
| $\mathbf{save} \ \delta(a_1 : \tau_1, \dots, a_n : \tau_n) \ \mathbf{in} \ E \mid \mathbf{run} \ \delta(a_1 : E_1, \dots, a_n : E_n)$
| $\mathbf{par}(E_1, \dots, E_n) \mid \mathbf{wait}(\langle tid \rangle)$
| $\mathbf{raise}(event-name) \mid \mathbf{register}(event-name, nm)$

memory actions:

$A ::= \mathbf{create}(pe_\tau, pe_{align}) \mid \mathbf{alloc}(pe_n) \mid \mathbf{kill}(pe_{ptr})$
| $\mathbf{store}(pe_\tau, pe_{ptr}, pe_n) \mid \mathbf{load}(pe_\tau, pe_{ptr}) \mid \dots$

6 - The Core language: *design motivation*

Very different from C11

Designed for the purpose of expressing the behaviour of C11 programs
⇒ each Core construct is motivated by a C11 behaviour

Suitable for usual semantics techniques

- ▶ each construct has only one distinct behaviour
- ▶ functional style
- ▶ strongly typed
- ▶ small-step with continuations stack

6 - The Core language

```
int main(void) {  
    int x = 0, y = 1;  
    return x+y;  
}
```

6 - The Core language

```
proc main(): eff integer :=
  let strong (x, y) = unseq(create(<alignof>("signed int"), "signed int"),
                           create(<alignof>("signed int"), "signed int")) in
  let strong ()      = store("signed int", x, conv_int("signed int", 0)) in
  let strong ()      = store("signed int", y, conv_int("signed int", 1)) in

  let strong a3 =
    let weak (a1, a2) = unseq(load("signed int", x), load("signed int", y)) in
    overflow("signed int", conv_int("signed int", a1) + conv_int("signed int", a2))
  in
  let strong () = unseq(kill(x), kill(y)) in
  return (conv_int("signed int", a3))
```

6 - The Core language

```
int f(int n) {  
    n+1;  
}
```

```
int main(void) {  
    int x = 0, y = 1;  
    return f(x)+y;  
}
```

6 - The Core language

```
proc f(n: pointer): eff integer :=
  let strong a1 =
    let weak (a2, a3) = unseq(load("signed int", n), 1) in
    overflow("signed int", conv_int("signed int", a2) + conv_int("signed int", a3))
  in
  undef("Reached_end_of_function")

proc main(): eff integer :=
  let strong (x, y) = unseq(create(<alignof>("signed int"), "signed int"),
    create(<alignof>("signed int"), "signed int")) in

  let strong () =
    unseq(store("signed int", y, conv_int("signed int", 1)),
      store("signed int", x, conv_int("signed int", 0))) in

  let strong a1 =
    let weak (a2, a3) = unseq(
      let strong ptr1 = create(<alignof>("signed int"), "signed int") in
      let strong tmp1 = load("signed int", x) in
      let strong () = store("signed int", ptr1, tmp1) in
      let strong ret1 = f{ptr1} in
      let strong () = kill(ptr1) in
      ret1
    ,
      load("signed int", y)
    ) in
    overflow("signed int", conv_int("signed int", a2) + conv_int("signed int", a3))
  in

  let strong () = unseq(kill(x), kill(y)) in
  return (conv_int("signed int", a1))
```


7 - Elaboration: the left shift operator ($E_1 \ll E_2$)

© ISO/IEC 2011 – All rights reserved

ISO/IEC 9899:2011 (E)

6.5.7 Bitwise shift operators

Syntax

- 1 *shift-expression*:
additive-expression
shift-expression \ll *additive-expression*
shift-expression \gg *additive-expression*

Constraints

- 2 Each of the operands shall have integer type.

Semantics

- 3 The integer promotions are performed on each of the operands. The type of the result is that of the promoted left operand. If the value of the right operand is negative or is greater than or equal to the width of the promoted left operand, the behavior is undefined.
- 4 The result of $E_1 \ll E_2$ is E_1 left-shifted E_2 bit positions; vacated bits are filled with zeros. If E_1 has an unsigned type, the value of the result is $E_1 \times 2^{E_2}$, reduced modulo one more than the maximum value representable in the result type. If E_1 has a signed type and nonnegative value, and $E_1 \times 2^{E_2}$ is representable in the result type, then that is the resulting value; otherwise, the behavior is undefined.
- 5 The result of $E_1 \gg E_2$ is E_1 right-shifted E_2 bit positions. If E_1 has an unsigned type or if E_1 has a signed type and a nonnegative value, the value of the result is the integral part of the quotient of $E_1 / 2^{E_2}$. If E_1 has a signed type and a negative value, the resulting value is implementation-defined.

7 - Elaboration: the left shift operator ($E_1 \ll E_2$)

$\llbracket (E_1^{ty_1} \ll E_2^{ty_2})^{ty} \rrbracket \rightsquigarrow$

```
let weak (a1, a2) = unseq( $\llbracket E_1 \rrbracket$ ,  $\llbracket E_2 \rrbracket$ ) in
```

```
let a1 = promote(ty1, ty, a1) in
```

```
let a2 = promote(ty2, ty, a2) in
```

```
if a2 < 0 then
```

```
  undef(Negative_shift)
```

```
else if ctype_width(ty) <= a2 then
```

```
  undef(Shift_too_large)
```

```
else
```

```
  let res = (a1 * 2a2) % (<ctype_max>(ty) + 1) in
```

```
  IF AllTypesAux.is_unsigned_integer_type(ty1) THEN
```

```
    res
```

```
  ELSE
```

```
    if is_representable(res, ty) then
```

```
      res
```

```
    else
```

```
      undef(Non_representable_shift)
```

8 - Memory layout model

The elaboration of C into Core is parametric on semantic choices regarding the memory.

(back to the second slide)

- ▶ focuses on the various **“de-facto” Cs** (to study existing C codes)

When it comes to memory related issues, it is common for (system-level) C programs to go outside of ISO C.

⇒ requires an empirical approach:

- ▶ we need to learn the assumptions made by programmers and compilers

8 - Memory layout model (litmus tests)

Can one inspect the value of a pointer to an object whose lifetime has ended?

```
#include <stdio.h>
#include <stdlib.h>
int main(void) {
    int i=0;
    int *pj = (int *) (malloc(sizeof(int)));
    *pj=1;
    printf("&i==pj)=%s\n", (&i==pj)?"true":"false");
    free(pj);
    printf("&i==pj)=%s\n", (&i==pj)?"true":"false");
    // is the == comparison above defined behaviour?
}
```

Undefined in ISO C

8 - Memory layout model (litmus tests)

Can one do relational comparison (with $<$, $>$, ...) of two pointers to separately allocated objects?

```
#include <stdio.h>
int y=2, x=1;
int main(void) {
    int *p=&x *p=&y;
    _Bool b1 = (p < q); // does this have defined behaviour?
    _Bool b2 = (p > q); // does this have defined behaviour?
    printf("(p<q = %s (p>q) = %s\n",
           b1?"true":"false", b2?"true":"false");
}
```

Undefined in ISO C

8 - Memory layout model (web survey)

[1/15] How predictable are reads from padding bytes?

If you zero all bytes of a struct and then write some of its members, do reads of the padding return zero? (e.g. for a bitwise CAS or hash of the struct, or to know that no security-relevant data has leaked into them.)

Will that work in normal C compilers?

- yes
- only sometimes
- no
- don't know
- I don't know what the question is asking

Do you know of real code that relies on it?

- yes
- yes, but it shouldn't
- no, but there might well be
- no, that would be crazy
- don't know

If it won't always work, is that because [check all that apply]:

- you've observed compilers write junk into padding bytes
- you think compilers will assume that padding bytes contain unspecified values and optimise away those reads
- Other:

Comment

9 - C/C++11 concurrency

work of Kyndylan Nienhuis.

⇒ integration with an **operational model** for C/C++11 concurrency.

- ▶ enable the incremental exploration of larger concurrent programs
- ▶ formally proved equivalent to Batty's C11 axiomatic formalisation (in Isabelle/HOL)

parametricity of Core's dynamics on the memory helped:

- ▶ virtually no modification required on the concurrency model
- ▶ but relaxed reads required symbolic evaluation

Conclusion