

# A Typed C11 Semantics for Interactive Theorem Proving

**Robbert Krebbers**   Freek Wiedijk

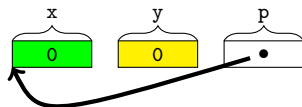
ICIS, Radboud University Nijmegen, The Netherlands

January 13, 2015 @ CPP, Mumbai, India

## What is this C program supposed to do?

```
int x = 0, y = 0, *p = &x;
int f() { p = &y; return 17; }
int main() {
    *p = f();
    printf("x=%d,y=%d\n", x, y);
}
```

Initial state:




Let us try some compilers

- ▶ Clang prints `x=0,y=17`  
`f` is called first, thereafter `p` is evaluated to `&y`
- ▶ GCC prints `x=17,y=0`  
`p` is evaluated to `&x` first, then `f` is called

More subtle: `*p = (p = &y, 17);` has **undefined behavior**

# Contribution




## CH<sub>2</sub>O (Krebbers & Wiedijk)

- ▶ **Compiler independent** C11 semantics in  Coq
- ▶ Operational, executable, and axiomatic semantics

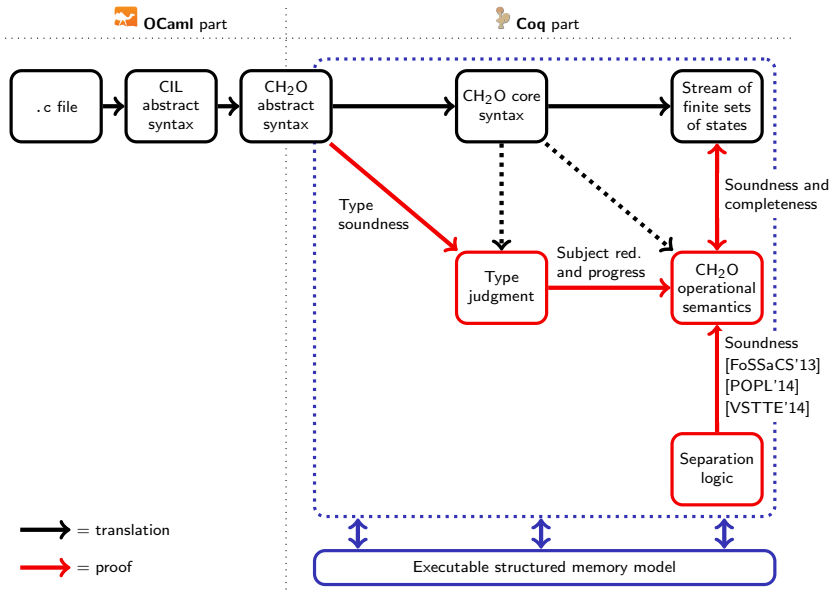
**CPP'15 contribution:** a verified interpreter to explore the **non-deterministic** behaviors of CH<sub>2</sub>O

- ▶ Type system & *weak type safety*
- ▶ Executable semantics & *soundness/completeness*
- ▶ Formal translation from AST & *type soundness*

## Recent related work

	 <b>CompCert</b>	 <b>KCC</b>	 <b>CH<sub>2</sub>O</b>
Compiler indep/close to C11	○	◐	●
Size of C fragment	●	●	◐
Proof assistant support	●	○	●
Type system	○	○	●
Principled core language	●	○	●
Formal translation from AST	○	n/a	●

# Overview of the CH<sub>2</sub>O project





# CH<sub>2</sub>O abstract C

## Formal translation to core C

Conversions include:

- ▶ Named variables to De Bruijn indices
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$
- ▶ Simplification of loops, e.g.

`while(e) s  $\Rightarrow$  catch (loop (if (e) skip else throw 0; catch s))`

- ▶ Expansion of `typedef` and `enum` declarations
- ▶ Translation of constants like `INT_MIN`
- ▶ Translation of compound literals, e.g.  
`(struct S){ .x=1, {4,r}, .y[4+1]=0, q }`

## Theorem (Type soundness)

The translator only produces well-typed CH<sub>2</sub>O core programs

## CH<sub>2</sub>O operational semantics

- ▶ **Zippers** are used to describe non-local control flow
- ▶ **Structured memory model** (as separation algebra) to accurately describe low- versus high-level subtleties of C11
- ▶ **Permissions** (as separation algebra) are used for:
  - ▶ Ruling out expressions like  $(x = 1) + (x = 2)$
  - ▶ Connection with separation logic
- ▶ **Evaluation contexts** for non-deterministic redex selection
- ▶ **Stuck states** for undefined behavior



# CH<sub>2</sub>O operational semantics

## Example of memory state

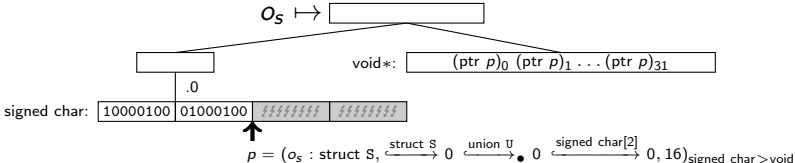
Consider:

```

struct S {
  union U {
    signed char x[2]; int y;
  } u;
  void *p;
} s = { { .x = {33,34} }, s.u.x + 2 }

```

The object in memory may look like:



# Typing of CH<sub>2</sub>O core C

**Expression judgment**  $\Gamma, \Gamma_f, \Delta, \vec{\tau} \vdash e : \tau_r$

- ▶ Struct/union fields:  $\Gamma \in \text{tag} \rightarrow_{\text{fin}} \text{list type}$
- ▶ Functions:  $\Gamma_f \in \text{funname} \rightarrow_{\text{fin}} (\text{list type} \times \text{type})$
- ▶ Memory layout:  $\Delta \in \text{index} \rightarrow_{\text{fin}} (\text{type} \times \text{bool})$
- ▶ De Bruijn variables:  $\vec{\tau} \in \text{list type}$

For example:

$$\frac{\vec{\tau}(i) = \tau}{x_i^{\vec{\tau}} : \tau} \quad \frac{e : \tau_1}{\&e : (\tau*)_r} \quad \frac{\Gamma_f(f) = (\vec{\tau}, \sigma) \quad \vec{e} : \vec{\tau}_r}{f(\vec{e}) : \sigma_r}$$

**Statement judgment**  $\Gamma, \Gamma_f, \Delta, \vec{\tau} \vdash s : (\beta, \tau^?)$

$$\frac{}{\text{skip} : (\text{false}, \perp)} \quad \frac{e : \tau_r}{\text{return } e : (\text{true}, \tau)} \quad \frac{}{\text{goto } l : (\text{true}, \perp)}$$

**State judgment**  $\Gamma, \Gamma_f, \Delta \vdash S : g$  (typically  $g = \text{main}$ )

# Typing of CH<sub>2</sub>O core C

## Type preservation

### Lemma (Type preservation)

*If  $S_1 : g$  and  $S_1 \rightarrow S_2$ , then  $S_2 : g$*

### Theorem (Weak type safety)

If  $S_1$  initial for  $g(\vec{v})$ , then if  $S_1 \rightarrow^* S_2$  we have either:

1. Not finished:  $S_2 \rightarrow S_3$  for some  $S_3$
2. Undefined behavior:  $S_2 = \mathbf{S}(\mathcal{P}, \overline{\text{undef}} \phi_U, m)$
3. Final state:  $S_2 = \mathbf{S}(\epsilon, \overline{\text{return}} g \vec{v}, m)$

# Executable semantics

**Goal:** define  $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$  and extract to OCaml

## Problems:

1. Decomposition  $\mathcal{E}[e_1]$  of expressions is non-deterministic:

$$\mathbf{S}(\mathcal{P}, \mathcal{E}[e_1], m_1) \rightarrow \mathbf{S}(\mathcal{P}, \mathcal{E}[e_2], m_2) \text{ if } (e_1, m_1) \rightarrow_{\text{h}} (e_2, m_2)$$

2. Object identifiers  $o$  for newly allocated memory are arbitrary:

$$\begin{aligned} & \mathbf{S}(\mathcal{P}, (\backslash_{\nu}, \text{local}_{\tau} s), m) \\ & \rightarrow \mathbf{S}((\text{local}_{o:\tau} \square) \mathcal{P}, (\backslash_{\nu}, s), \text{alloc}_{\Gamma} o \tau \text{ false } m) \text{ if } o \notin \text{dom } m \end{aligned}$$

## Solutions:

1. Enumerate all possible decompositions  $\mathcal{E}[e_1]$
2. Pick a canonical object identifier **fresh**  $m$  for  $o$  (makes completeness difficult!)

# Executable semantics

## Soundness and completeness

### Theorem (Soundness)

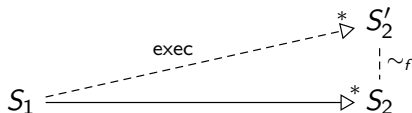
If  $S_2 \in \text{exec } S_1$ , then  $S_1 \rightarrow S_2$

### Definition (Permutation)

We let  $S_1 \sim_f S_2$ , if  $S_2$  is obtained by renaming  $S_1$  with respect to  $f : \text{index} \rightarrow \text{option index}$

### Theorem (Completeness)

If  $S_1 \rightarrow^* S_2$ , then there exists an  $f$  and  $S'_2$  such that:



# Formalization in Coq

Interpreter extracted to 🐨 OCaml from 🦊 Coq

- ▶ **Error monad** for failure of type checking
- ▶ **Set monad** for non-determinism
- ▶ Verified **hash sets** for efficiency

All essential properties proven in Coq:

- ▶ Weak type safety
- ▶ Soundness and completeness of executable semantics
- ▶ Type soundness of translation from AST

Part of ~40.000 LOC constructive and axiom free development

# Conclusion

A programming language semantics should consist of:

- ▶ **Operational semantics**  
Reasoning about program transformations
- ▶ **Axiomatic semantics**  
Correctness proofs of concrete programs
- ▶ **Executable semantics**  
Debugging and testing

Extremely challenging to develop matching versions for C11

**Future work:** still many parts of C11 left to be explored

## Demo and questions

Sources: `http://robertkrebbers.nl/research/ch2o/`