

# The C standard formalized in Coq, what's next?

Robbert Krebbers

Aarhus University, Denmark

May 13, 2016 @ Cambridge Computer Laboratory, UK

# What is this program supposed to do?

The C quiz, question 1

```
int main() {  
    int x;  
    int y = (x = 3) + (x = 4);  
    printf("x=%d,y=%d\n", x, y);  
}
```

# What is this program supposed to do?

The C quiz, question 1

```
int main() {  
    int x;  
    int y = (x = 3) + (x = 4);  
    printf("x=%d,y=%d\n", x, y);  
}
```

Let us try some compilers

- ▶ Clang prints `x=4,y=7`, seems just left-right

# What is this program supposed to do?

The C quiz, question 1

```
int main() {  
    int x;  
    int y = (x = 3) + (x = 4);  
    printf("x=%d,y=%d\n", x, y);  
}
```

Let us try some compilers

- ▶ Clang prints `x=4,y=7`, seems just left-right
- ▶ GCC prints `x=4,y=8`, does not correspond to any order

# What is this program supposed to do?

The C quiz, question 1

```
int main() {
    int x;
    int y = (x = 3) + (x = 4);
    printf("x=%d,y=%d\n", x, y);
}
```

Let us try some compilers

- ▶ Clang prints `x=4,y=7`, seems just left-right
- ▶ GCC prints `x=4,y=8`, does not correspond to any order

This program violates the **sequence point** restriction

- ▶ due to two unsequenced writes to `x`
- ▶ resulting in **undefined behavior**
- ▶ thus both compilers are right

# Underspecification in C11

- ▶ **Unspecified behavior:** two or more behaviors are allowed  
*For example: order of evaluation in expressions* (+57 more)
- ▶ **Implementation defined behavior:** like unspecified behavior, but the compiler has to document its choice  
*For example: size and endianness of integers* (+118 more)
- ▶ **Undefined behavior:** the standard imposes no requirements at all, the program is even allowed to crash  
*For example: dereferencing a NULL or dangling pointer, signed integer overflow, ...* (+201 more)

# Underspecification in C11

- ▶ **Unspecified behavior:** two or more behaviors are allowed  
*For example: order of evaluation in expressions* (+57 more)  
Non-determinism
- ▶ **Implementation defined behavior:** like unspecified behavior, but the compiler has to document its choice  
*For example: size and endianness of integers* (+118 more)  
Parametrization
- ▶ **Undefined behavior:** the standard imposes no requirements at all, the program is even allowed to crash  
*For example: dereferencing a NULL or dangling pointer, signed integer overflow, ...* (+201 more)  
No semantics/crash state

# Why does C use underspecification that heavily?

**Pros** for optimizing compilers:

- ▶ More optimizations are possible
- ▶ High run-time efficiency
- ▶ Easy to support multiple architectures



# Why does C use underspecification that heavily?

## **Pros** for optimizing compilers:

- ▶ More optimizations are possible
- ▶ High run-time efficiency
- ▶ Easy to support multiple architectures

## **Cons** for programmers/formal methods people:

- ▶ Portability and maintenance problems
- ▶ Hard to capture precisely in a semantics
- ▶ Hard to formally reason about

# Approaches to underspecification

## **CompCert** (Leroy *et al.*) / **VST** (Appel *et al.*)

- ▶ Main goal: verification of/w.r.t. CompCert compiler in Coq
- ▶ Semantics only needs to be correct for CompCert compiler  
*For example: integer overflow and aliasing violations **not** UB*

## **KCC** (Ellison & Rosu, Hathhorn *et al.*)

- ▶ Main goal: compiler independent C11 semantics in  $\mathbb{K}$
- ▶ Describes **most** unspecified and undefined behavior
- ▶ No proof assistant support

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

- ▶ Main goal: compiler independent C11 semantics in Coq
- ▶ Describes **all** unspecified and undefined behavior
- ▶ Describes **some** implementation-defined behavior  
*For example: no legacy architectures with 1s' complement*

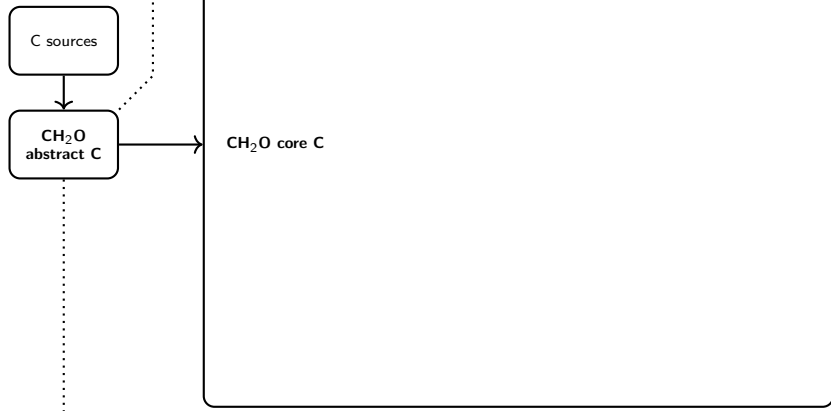
## **Cerberus** (Sewell *et al.*)

- ▶ Main goal: '*defacto*' C11 semantics in LEM
- ▶ Improve standard to match the way C is used in practice

# The CH<sub>2</sub>O project

 OCaml part

 Coq part



# The CH<sub>2</sub>O project

 OCaml part

 Coq part

C sources

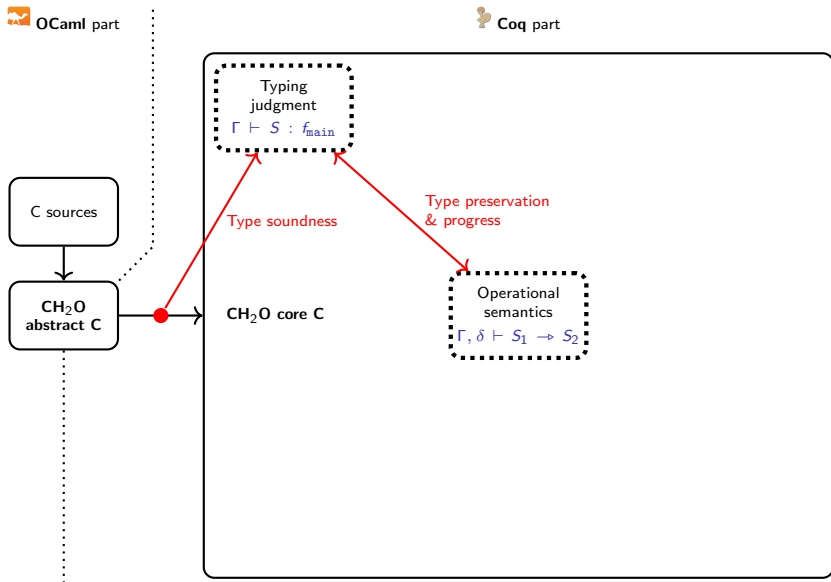
CH<sub>2</sub>O  
abstract C

CH<sub>2</sub>O core C

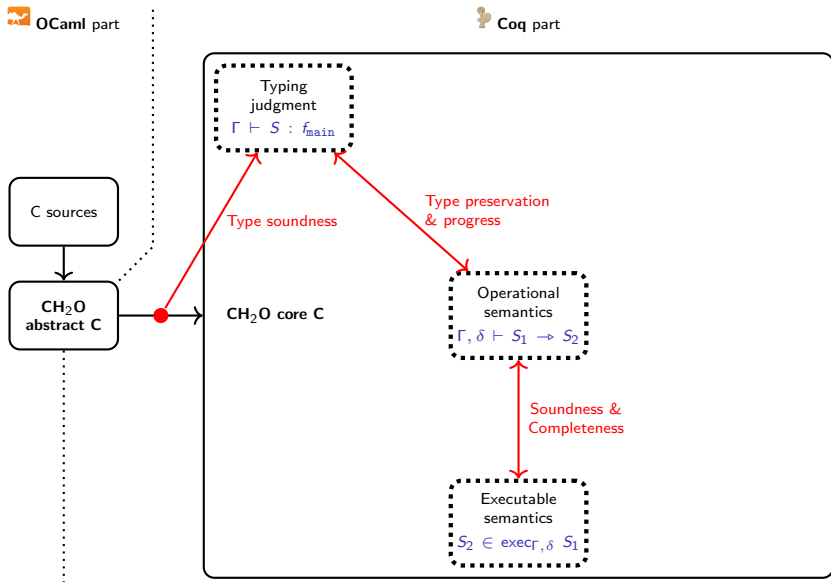
Operational  
semantics

$\Gamma, \delta \vdash S_1 \rightarrow S_2$

# The CH<sub>2</sub>O project



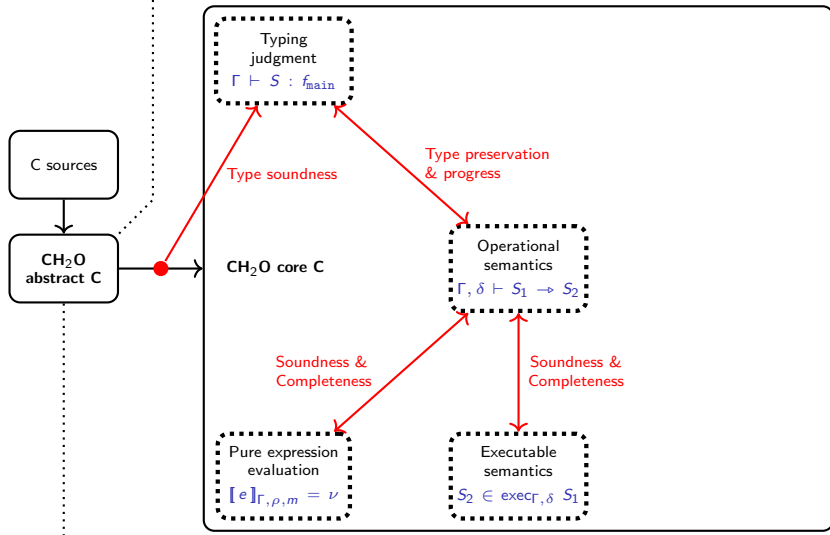
# The CH<sub>2</sub>O project



# The CH<sub>2</sub>O project

 OCaml part

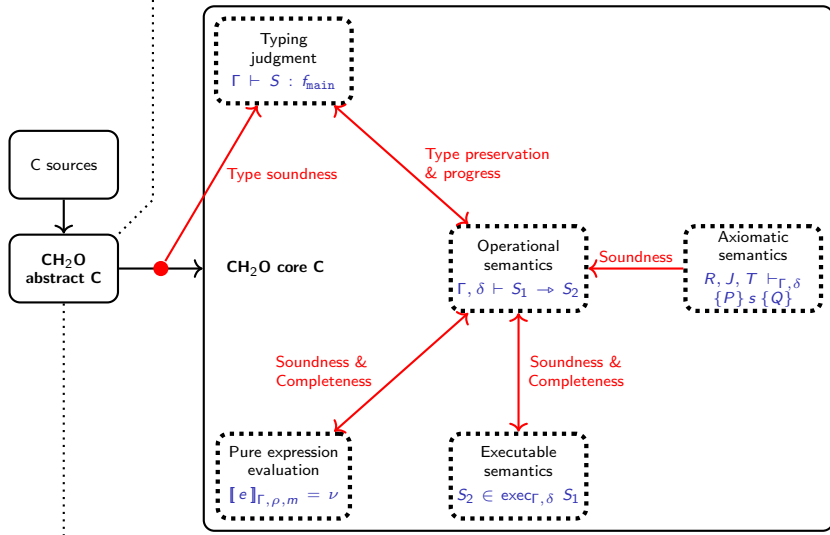
 Coq part



# The CH<sub>2</sub>O project

 OCaml part

 Coq part

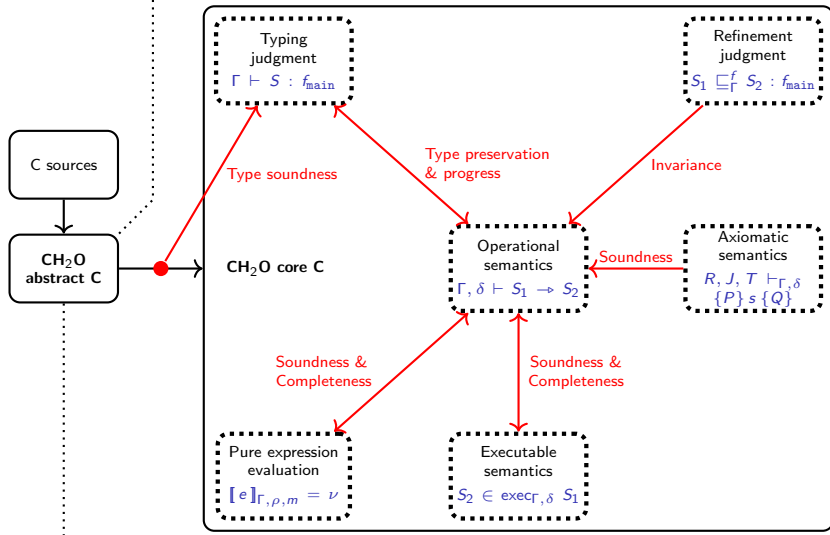




# The CH<sub>2</sub>O project

 OCaml part

 Coq part



# Non-local control flow and block scope variables

## The C quiz, question 2

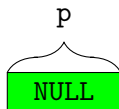
```
int *p = NULL;
l: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto l;
}
```

# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:

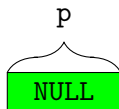


# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:



# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:

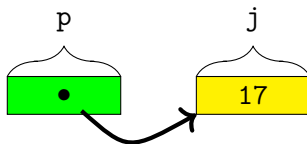


# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:

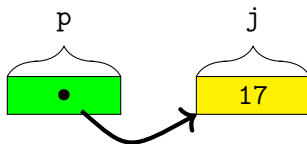


# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:

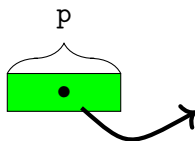


# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

memory:



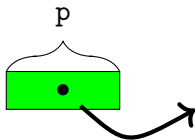


# Non-local control flow and block scope variables

The C quiz, question 2

```
int *p = NULL;
1: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto 1;
}
```

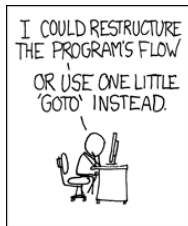
memory:



C11, 6.2.4p2: the value of a pointer becomes indeterminate when the object it points to (or just past) reaches the end of its lifetime.  
⇒ **Undefined behavior**

# Non-local control flow and block scope variables

Goto considered harmful?



<http://xkcd.com/292/>


# Non-local control flow and block scope variables

Goto considered harmful?



<http://xkcd.com/292/>

Not necessarily:

  $\vdash \{P\} \dots \text{goto main\_sub3}; \dots \{Q\}$

# Non-local control flow and block scope variables

Separation logic for non-local control

**Statement judgment:**

$$R, J, T \vdash \{P\} s \{Q\}$$

# Non-local control flow and block scope variables

Separation logic for non-local control

**Statement judgment:**

$$R, J, T \vdash \{P\} s \{Q\}$$

where:

- ▶  $\{P\} s \{Q\}$  is a Hoare triple, as usual

# Non-local control flow and block scope variables

Separation logic for non-local control

**Statement judgment:**

$$R, J, T \vdash \{P\} s \{Q\}$$

where:

- ▶  $\{P\} s \{Q\}$  is a Hoare triple, as usual
- ▶  $R$  has to hold to execute a return

# Non-local control flow and block scope variables

Separation logic for non-local control

## Statement judgment:

$$R, J, T \vdash \{P\} s \{Q\}$$

where:

- ▶  $\{P\} s \{Q\}$  is a Hoare triple, as usual
  - ▶  $R$  has to hold to execute a return
  - ▶  $J$  maps labels to their jumping condition
- When executing a goto  $l$ , the assertion  $Jl$  has to hold

# Non-local control flow and block scope variables

Separation logic for non-local control

## Statement judgment:

$$R, J, T \vdash \{P\} s \{Q\}$$

where:

- ▶  $\{P\} s \{Q\}$  is a Hoare triple, as usual
- ▶  $R$  has to hold to execute a return
- ▶  $J$  maps labels to their jumping condition  
When executing a goto  $l$ , the assertion  $Jl$  has to hold
- ▶  $T$  maps breaks/continues to their jumping condition



# Non-local control flow and block scope variables

Separation logic for non-local control

## Statement judgment:

$$R, J, T \vdash \{P\} s \{Q\}$$

where:

- ▶  $\{P\} s \{Q\}$  is a Hoare triple, as usual
- ▶  $R$  has to hold to execute a return
- ▶  $J$  maps labels to their jumping condition  
When executing a goto  $l$ , the assertion  $Jl$  has to hold
- ▶  $T$  maps breaks/continues to their jumping condition

Example:

$$\frac{}{R, J, T \vdash \{Jl\} \text{goto } l \{Q\}}$$

$$\frac{}{R, J, T \vdash \{Jl\} l : \{Jl\}}$$

# Non-local control flow and block scope variables

The block scope variable rule

$$\frac{R \uparrow * x_0 \mapsto -, J \uparrow * (x_0 \mapsto - : \tau), T \uparrow * (x_0 \mapsto - : \tau) \quad \vdash \{P \uparrow * (x_0 \mapsto - : \tau)\} s \{Q \uparrow * (x_0 \mapsto - : \tau)\}}{R, J, T \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

- ▶ The De Bruijn indexes are lifted:  $(-) \uparrow$
- ▶ The memory is extended:  $(-) * (x_0 \mapsto - : \tau)$

When leaving a block: the reverse

# Non-local control flow and block scope variables

The block scope variable rule

$$\frac{R \uparrow * x_0 \mapsto -, J \uparrow * (x_0 \mapsto - : \tau), T \uparrow * (x_0 \mapsto - : \tau) \quad \vdash \{P \uparrow * (x_0 \mapsto - : \tau)\} s \{Q \uparrow * (x_0 \mapsto - : \tau)\}}{R, J, T \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

- ▶ The De Bruijn indexes are lifted:  $(-) \uparrow$
- ▶ The memory is extended:  $(-) * (x_0 \mapsto - : \tau)$

When leaving a block: the reverse

**Important:**

- ▶ Symmetry matches gotos going both in and out
- ▶ Using De Bruijn indexes avoids shadowing

# Non-determinism and sequence points

The C quiz, question 3

```
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);
}
```

# Non-determinism and sequence points

The C quiz, question 3

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

Let us try some compilers

- ▶ Clang prints `x=0,y=17`
- ▶ GCC prints `x=17,y=0`

Non-determinism appears even in innocently looking code

# Non-determinism and sequence points

## Separation logic for C expressions

**Observation:** non-determinism corresponds to concurrency

**Idea:** use the separation logic rule for parallel composition

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\}}{\{P_1 * P_2\} e_1 \odot e_2 \{Q_1 * Q_2\}}$$

# Non-determinism and sequence points

## Separation logic for C expressions

**Observation:** non-determinism corresponds to concurrency

**Idea:** use the separation logic rule for parallel composition

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\}}{\{P_1 * P_2\} e_1 \odot e_2 \{Q_1 * Q_2\}}$$

What does this mean:

- ▶ Split the memory into two disjoint parts
- ▶ Prove that  $e_1$  and  $e_2$  can be executed safely in their part
- ▶ Now  $e_1 \odot e_2$  can be executed safely in the whole memory

# Non-determinism and sequence points

## Separation logic for C expressions

**Observation:** non-determinism corresponds to concurrency

**Idea:** use the separation logic rule for parallel composition

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\}}{\{P_1 * P_2\} e_1 \odot e_2 \{Q_1 * Q_2\}}$$

What does this mean:

- ▶ Split the memory into two disjoint parts
- ▶ Prove that  $e_1$  and  $e_2$  can be executed safely in their part
- ▶ Now  $e_1 \odot e_2$  can be executed safely in the whole memory

**Disjointness  $\Rightarrow$  no sequence point violation** (accessing the same location twice in one expression)



# Non-determinism and sequence points

Hoare triples

**Expression judgment:**  $\{P\} e \{Q\}$

# Non-determinism and sequence points

## Hoare triples

**Expression judgment:**  $\{P\} e \{Q\}$

$Q : \text{val} \rightarrow \text{assert}$

If  $P$  holds beforehand, then

- ▶  $e$  does not crash
- ▶  $Q \ v$  holds afterwards when terminating with  $v$

# Non-determinism and sequence points

Some actual rules

Binary operators:

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\} \quad \forall v_1 v_2 . (Q_1 v_1 * Q_2 v_2 \models \exists v' . (v_1 \odot v_2) \Downarrow v' \wedge Q' v')}{\{P_1 * P_2\} e_1 \odot e_2 \{Q'\}}$$

# Non-determinism and sequence points

Some actual rules

Binary operators:

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\} \quad \forall v_1 v_2. (Q_1 v_1 * Q_2 v_2 \models \exists v'. (v_1 \odot v_2) \Downarrow v' \wedge Q' v')}{\{P_1 * P_2\} e_1 \odot e_2 \{Q'\}}$$

Simple assignments:

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\} \quad \text{Writable} \subseteq \text{kind } \gamma \quad \forall p v. \left( Q_1 p * Q_2 v \models \exists v'. (\tau) v \Downarrow v' \wedge \left( (p \xrightarrow[\mu]{\gamma} - : \tau) * \left( (p \xrightarrow[\mu]{\text{lock } \gamma} |v'| \circ : \tau) \multimap Q' v' \right) \right) \right)}{\{P_1 * P_2\} e_1 := e_2 \{Q'\}}$$

# Non-determinism and sequence points

Some actual rules

Binary operators:

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\} \quad \forall v_1 v_2. (Q_1 v_1 * Q_2 v_2 \models \exists v'. (v_1 \odot v_2) \Downarrow v' \wedge Q' v')}{\{P_1 * P_2\} e_1 \odot e_2 \{Q'\}}$$

Simple assignments:

$$\frac{\{P_1\} e_1 \{Q_1\} \quad \{P_2\} e_2 \{Q_2\} \quad \text{Writable} \subseteq \text{kind } \gamma \quad \forall p v. \left( Q_1 p * Q_2 v \models \exists v'. (\tau) v \Downarrow v' \wedge \left( (p \xrightarrow[\mu]{\gamma} - : \tau) * ((p \xrightarrow[\mu]{\text{lock } \gamma} |v'|_o : \tau) \multimap Q' v') \right) \right)}{\{P_1 * P_2\} e_1 := e_2 \{Q'\}}$$

Comma:

$$\frac{\{P\} e_1 \{\lambda_. P' \diamond\} \quad \{P'\} e_2 \{Q\}}{\{P\} (e_1, e_2) \{Q\}}$$

# Strict-aliasing

What is aliasing?

**Aliasing:** multiple pointers referring to the same object

```
int x, *p = &x, *q = &x; // p and q are aliased
```

# Strict-aliasing

What is aliasing?

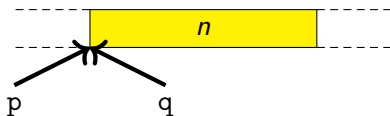
**Aliasing:** multiple pointers referring to the same object

```
int x, *p = &x, *q = &x; // p and q are aliased
```

Tricky with functions:

```
int f(int *p, int *q) {  
    int x = *q; *p = 17; return x;  
}
```

If  $p$  and  $q$  alias, the original value  $n$  of  $*p$  is returned



# Strict-aliasing

What is aliasing?

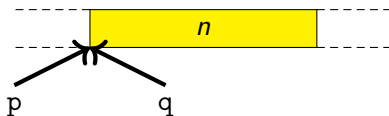
**Aliasing:** multiple pointers referring to the same object

```
int x, *p = &x, *q = &x; // p and q are aliased
```

Tricky with functions:

```
int f(int *p, int *q) {  
    int x = *q; *p = 17; return * *q;  
}
```

If p and q alias, the original value  $n$  of  $*p$  is returned



Eliminating x is unsound: 17 would be returned



# Strict-aliasing

## Alias analysis

**Alias analysis:** to determine whether pointers can alias

# Strict-aliasing

## Alias analysis

**Alias analysis:** to determine whether pointers can alias

Consider a similar function:

```
short g(int *p, short *q) {  
    short x = *q; *p = 17; return x;  
}
```

# Strict-aliasing

## Alias analysis

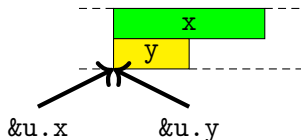
**Alias analysis:** to determine whether pointers can alias

Consider a similar function:

```
short g(int *p, short *q) {  
    short x = *q; *p = 17; return x;  
}
```

And call it with aliased pointers:

```
union { int x; short y; } u;  
u.y = 3;  
g(&u.x, &u.y);
```



# Strict-aliasing

## Alias analysis

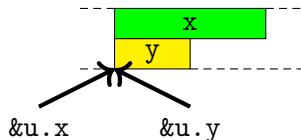
**Alias analysis:** to determine whether pointers can alias

Consider a similar function:

```
short g(int *p, short *q) {  
    short x = *q; *p = 17; return x;  
}
```

And call it with aliased pointers:

```
union { int x; short y; } u;  
u.y = 3;  
g(&u.x, &u.y);
```



C99/C11 allow **type-based alias analysis**: reads/writes with “the wrong type” cause undefined behavior  
⇒ A compiler can **assume** that  $p$  and  $q$  do not alias

# Strict-aliasing

How to treat pointers

## Others (e.g. CompCert)

Memory: a finite map of objects which consist of **arrays** of bytes

Pointers: pairs  $(o, i)$  where  $o$  identifies the object, and  $i$  the **offset** into that object

Too little information to formalize strict-aliasing

## Our approach

# Strict-aliasing

How to treat pointers

## Others (e.g. CompCert)

Memory: a finite map of objects which consist of **arrays** of bytes

Pointers: pairs  $(o, i)$  where  $o$  identifies the object, and  $i$  the **offset** into that object

Too little information to formalize strict-aliasing

## Our approach

A finite map of objects which consist of well-typed **trees** of bits

# Strict-aliasing

## How to treat pointers

### Others (e.g. CompCert)

Memory: a finite map of objects which consist of **arrays** of bytes

Pointers: pairs  $(o, i)$  where  $o$  identifies the object, and  $i$  the **offset** into that object

Too little information to formalize strict-aliasing

### Our approach

A finite map of objects which consist of well-typed **trees** of bits

Pairs  $(o, r)$  where  $o$  identifies the object, and  $r$  the **path through the tree** in that object

# Strict-aliasing

## How to treat pointers

### Others (e.g. CompCert)

Memory: a finite map of objects which consist of **arrays** of bytes

Pointers: pairs  $(o, i)$  where  $o$  identifies the object, and  $i$  the **offset** into that object

Too little information to formalize strict-aliasing

### Our approach

A finite map of objects which consist of well-typed **trees** of bits

Pairs  $(o, r)$  where  $o$  identifies the object, and  $r$  the **path through the tree** in that object

A semantics for strict-aliasing



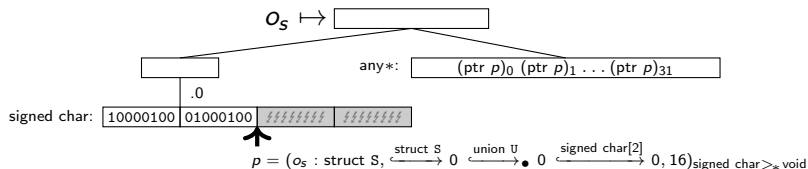
# Strict-aliasing

Example of the memory as a structured forest

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 looks like:



# Strict-aliasing

## Theorem (Strict-aliasing)

*Given:*

- ▶ *addresses  $\Gamma, \Delta \vdash a_1 : \sigma_1$  and  $\Gamma, \Delta \vdash a_2 : \sigma_2$*
- ▶ *with annotations that do not allow type-punning*
- ▶  *$\sigma_1, \sigma_2 \neq \text{unsigned char}$*
- ▶  *$\sigma_1$  not a subtype of  $\sigma_2$  and vice versa*

*Then there are two possibilities:*

1.  *$a_1$  and  $a_2$  do not alias*
2. *accessing  $a_1$  after  $a_2$  (and vice versa) is undefined*

# Strict-aliasing

## Theorem (Strict-aliasing)

*Given:*

- ▶ *addresses  $\Gamma, \Delta \vdash a_1 : \sigma_1$  and  $\Gamma, \Delta \vdash a_2 : \sigma_2$*
- ▶ *with annotations that do not allow type-punning*
- ▶  *$\sigma_1, \sigma_2 \neq \text{unsigned char}$*
- ▶  *$\sigma_1$  not a subtype of  $\sigma_2$  and vice versa*

*Then there are two possibilities:*

1.  *$a_1$  and  $a_2$  do not alias*
2. *accessing  $a_1$  after  $a_2$  (and vice versa) is undefined*

## Corollary

*Compilers can perform type-based alias analysis*

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

$x \in \text{string} ::= \text{Set of strings}$   
 $k \in \text{cintrank} ::= \text{char} \mid \text{short} \mid \text{int}$   
 $\quad \mid \text{long} \mid \text{long long} \mid \text{ptr}$   
 $si \in \text{signedness} ::= \text{signed} \mid \text{unsigned}$   
 $\tau_i \in \text{cinttype} ::= si? k$   
 $\tau \in \text{ctype} ::= \text{void} \mid \text{def } x \mid \tau_i \mid \tau*$   
 $\quad \mid \overrightarrow{\tau x} \rightarrow \tau \mid \tau[e]$   
 $\quad \mid \text{struct } x \mid \text{union } x$   
 $\quad \mid \text{enum } x \mid \text{typeof } e$   
 $e \in \text{cexpr} ::= x \mid \text{const}_{\tau_i} z \mid \text{string } \vec{z}$   
 $\quad \mid \text{sizeof } \tau \mid \text{alignof } \tau$   
 $\quad \mid \text{offsetof } \tau x$   
 $\quad \mid \tau \text{ min} \mid \tau \text{ max} \mid \tau \text{ bits}$   
 $\quad \mid \&e \mid *e \mid e . x$   
 $\quad \mid e_1 \alpha e_2$   
 $\quad \mid e(\vec{e}) \mid \text{alloc}_{\tau} e \mid \text{free } e$   
 $\quad \mid \odot_u e \mid e_1 \odot e_2 \mid (\tau)l$   
 $\quad \mid e_1 \&\& e_2 \mid e_1 \parallel e_2$   
 $\quad \mid (e_1, e_2) \mid e_1 ? e_2 : e_3$   
 $r \in \text{crefseg} ::= [e] \mid .x$

$l \in \text{cinit} ::= e \mid \{\overrightarrow{\tau} := l\}$   
 $\text{sto} \in \text{cstorage} ::= \text{static} \mid \text{extern} \mid \text{auto}$   
 $s \in \text{cstmt} ::= e \mid \text{return } e?$   
 $\quad \mid \text{goto } x \mid x : s$   
 $\quad \mid \text{break} \mid \text{continue}$   
 $\quad \mid \{s\}$   
 $\quad \mid \overrightarrow{\text{sto}} \tau x := l? ; s$   
 $\quad \mid \text{typedef } x := \tau ; s$   
 $\quad \mid \text{skip} \mid s_1 ; s_2$   
 $\quad \mid \text{while}(e) s$   
 $\quad \mid \text{do } s \text{ while}(e)$   
 $\quad \mid \text{for}(e_1 ; e_2 ; e_3) s$   
 $\quad \mid \text{if } (e) s_1 \text{ else } s_2$   
 $d \in \text{decl} ::= \text{struct } \overrightarrow{\tau} \vec{x} \mid \text{union } \overrightarrow{\tau} \vec{x}$   
 $\quad \mid \text{enum } \overrightarrow{x} := e? : \tau_i$   
 $\quad \mid \text{typedef } \tau$   
 $\quad \mid \text{global } l? : \overrightarrow{\text{sto}} \tau$   
 $\quad \mid \text{fun } s : \overrightarrow{\text{sto}} \tau$   
 $\Theta \in \text{decls} ::= \text{list}(\text{string} \times \text{decl})$

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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices

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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values

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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

- ▶ Simplification of loops, e.g. `while(e) s` becomes

`catch (loop (if (e') skip else throw 0 ; catch s'))`



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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

- ▶ Simplification of loops, e.g. `while(e) s` becomes

`catch (loop (if (e') skip else throw 0 ; catch s'))`

- ▶ Expansion of `typedef` and `enum` declarations

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

## Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

- ▶ Simplification of loops, e.g. `while(e) s` becomes

`catch (loop (if (e') skip else throw 0 ; catch s'))`

- ▶ Expansion of `typedef` and `enum` declarations
- ▶ Translation of constants like `INT_MIN`

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

Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

- ▶ Simplification of loops, e.g. `while(e) s` becomes

`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

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

## Translation to CH<sub>2</sub>O core C in Coq

- ▶ Named variables to De Bruijn indices
- ▶ Disambiguate l-values and r-values
- ▶ Sound/complete constant expression evaluation, e.g. in  $\tau[e]$

$$\llbracket e \rrbracket_{\Gamma, \text{locals } \mathcal{P}, m} = v \quad \text{iff} \quad \mathbf{S}(\mathcal{P}, e, m) \rightarrow^* \mathbf{S}(\mathcal{P}, v, m)$$

- ▶ Simplification of loops, e.g. `while(e) s` becomes  
`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

Theorem (Type soundness)

*Translation only produces well-typed CH<sub>2</sub>O core C programs*

# Conclusion

Formal methods can be applied to real programming languages

- ▶ Large part of the C11 standard formalized in Coq
- ▶ Many oddities in the C11 standard text discovered
- ▶ Metatheory is important to establish sanity of specification
- ▶ Executable semantics important to test specification
- ▶ Extensions of separation logic developed

# Future work

STW project "Sovereign" (Radboud University)

*Modular and practical verification of C programs*



- ▶ Develop design patterns to classify critical parts
- ▶ Formalize Misra C as a sublanguage of CH<sub>2</sub>O
- ▶ Develop proof infrastructure
- ▶ Connect CH<sub>2</sub>O to CompCert
- ▶ Methods for proving security properties
- ▶ Case studies at Nuclear Research Group and Rijkswaterstaat

# Future work

## More features

- ▶ Formalized parser and preprocessor
- ▶ Floating point arithmetic
- ▶ Bitfields and `_Bool`
- ▶ Untyped malloc
- ▶ Variadic functions
- ▶ Register storage class
- ▶ Type qualifiers
- ▶ External functions and I/O

# Future work

## Improve executable semantics

- ▶ Better error messages
- ▶ Use more efficient data structures
- ▶ Perform optimizations
- ▶ More desugaring in Coq instead of OCaml
- ▶ Use on large test suites (e.g. CSmith or Cerberus tests)



# Future work

Symbolic execution for separation logic for expressions

**Expression judgment:**  $A \vdash \{P\} e \{Q\}$

Invariant



Symbolic execution:

- ▶ Use static analysis to determine which objects are written to
- ▶ Put read-only objects in invariant:

$$\frac{A_1 * A_2 \vdash \{P\} e \{Q\}}{A_1 \vdash \{A_2 * P\} e \{A_2 * Q\}}$$

- ▶ Invariant can be freely shared, but must be maintained by each atomic expression (in sequential C, function calls are atomic)

# Future work

## Concurrency

- ▶ Concurrency primitives: locks, message passing, ...
  - ▶ Rule out *any* racy concurrency
  - ▶ Well-understood and easy to reason about [Hobor, Appel, ...]
- ▶ Sequentially consistent concurrency
  - ▶ Thread-pool semantics
  - ▶ Difficult to reason about
  - ▶ Works well in separation logic [O'Hearn, Svendsen, Dinsdale-Young, Birkedal, Parkinson, Dreyer, Turon, ...]
  - ▶ Not sound with respect to C11 concurrency
- ▶ Weak memory concurrency
  - ▶ Still open problems w.r.t. semantics [Sewell, Batty, ...]
  - ▶ Very challenging in separation logic [Vafeiadis, ...]

# Questions



PhD thesis & Coq sources:  
<http://robertkrebbers.nl/thesis.html>