

Formalizing C in Coq

Robbert Krebbers

ICIS, Radboud University Nijmegen, The Netherlands

October 24, 2014 @ Princeton University, USA

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

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~~ **Challenges** 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
- ▶ Specific choices for unspecified/impl-defined behavior
For example: 32-bits ints
- ▶ Describes **some** undefined behavior
Undefined: dereferencing NULL, ...
But still defined: integer overflow, aliasing violations, ...
- ▶ VST separation logic proofs specific to CompCert

■ ■ **Formalin** (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

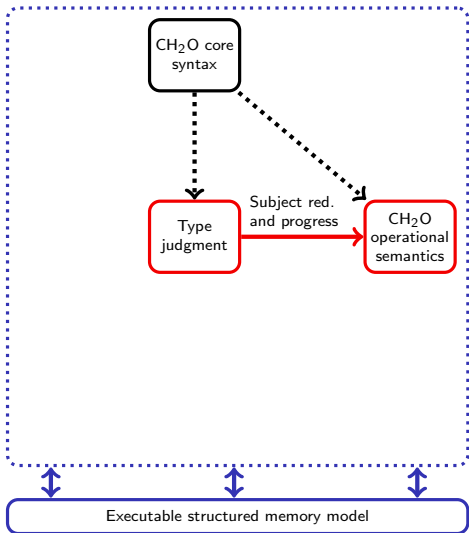
The Formalin project (Krebbbers & Wiedijk,)





OCaml part



Coq part



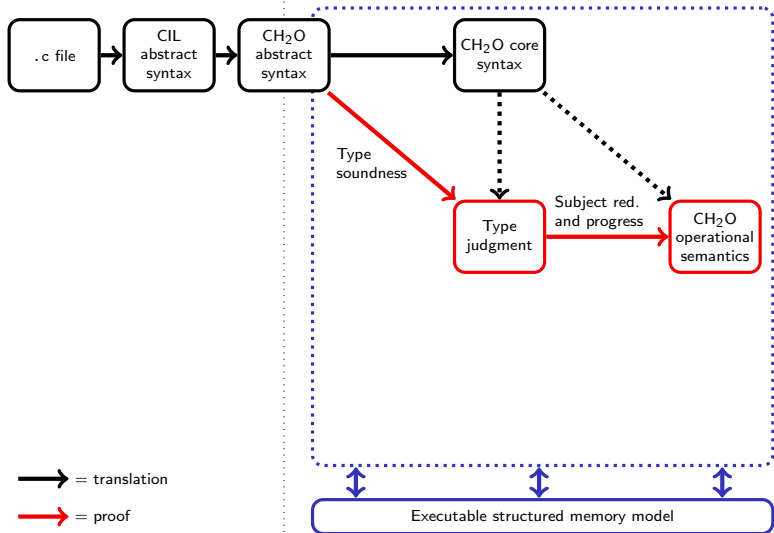
 = translation

 = proof

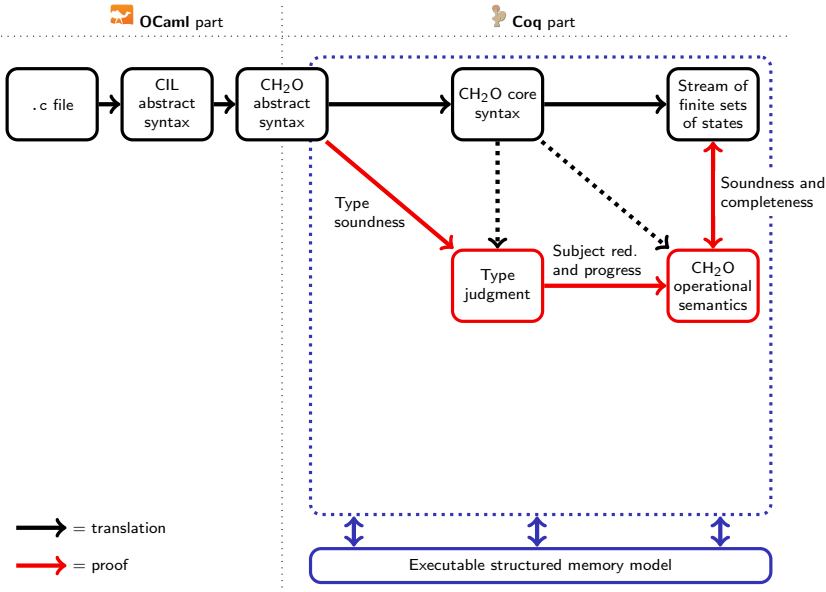
The Formalin project (Krebbbers & Wiedijk,)

 OCaml part

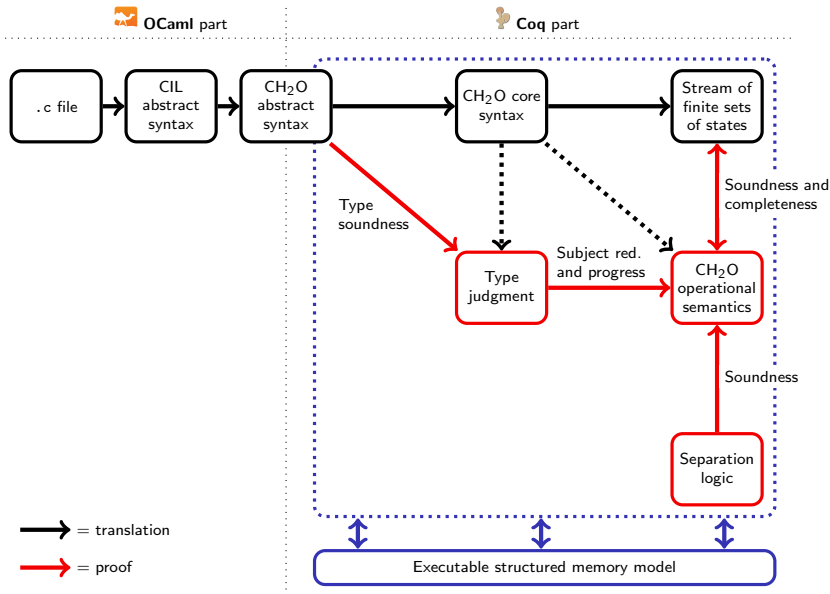
 Coq part



The Formalin project (Krebbbers & Wiedijk,)



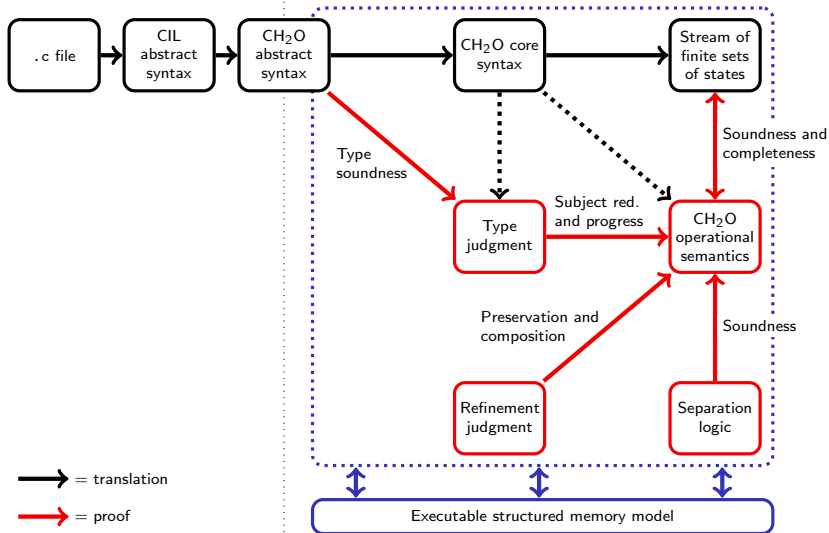
The Formalin project (Krebbbers & Wiedijk,)



The Formalin project (Krebbbers & Wiedijk,)

 OCaml part

 Coq part



Non-local control flow and block scope variables

The C quiz, question 2

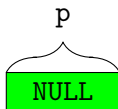
```
int *p = NULL;
l: if (p) {
    return (*p);
} else {
    int j = 17;
    p = &j;
    goto l;
}
```

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:

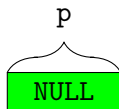


Non-local control flow and block scope variables

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:

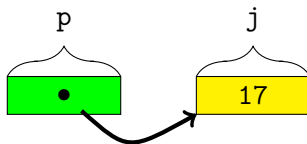


Non-local control flow and block scope variables

The C quiz, question 2

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

memory:

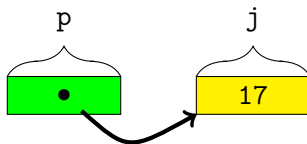


Non-local control flow and block scope variables

The C quiz, question 2

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

memory:

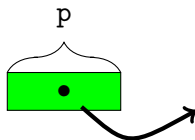


Non-local control flow and block scope variables

The C quiz, question 2

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

memory:

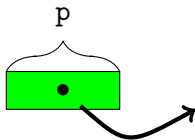


Non-local control flow and block scope variables

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

Operational semantics [Krebbbers/Wiedijk,FoSSaCS'13]

Problem: goto/return/etc. affect memory

Non-local control flow and block scope variables

Operational semantics [Krebbbers/Wiedijk,FoSSaCS'13]

Problem: goto/return/etc. affect memory

Solution:

- ▶ Execute gotos and returns in small steps
 - ▶ Not so much to search for labels, . . .
 - ▶ but to naturally perform required allocations and deallocations

Non-local control flow and block scope variables

Operational semantics [Krebbbers/Wiedijk,FoSSaCS'13]

Problem: goto/return/etc. affect memory

Solution:

- ▶ Execute gotos and returns in small steps
 - ▶ Not so much to search for labels, . . .
 - ▶ but to naturally perform required allocations and deallocations
- ▶ Traversal through the AST in the following directions:
 - ▶ ↘ downwards to the next statement
 - ▶ ↗ upwards to the next statement
 - ▶ ↪ `l` to a label `l`: after a `goto l`
 - ▶ ↵ `v` to the top of the statement after a `return`

Non-local control flow and block scope variables

Operational semantics [Krebbbers/Wiedijk,FoSSaCS'13]

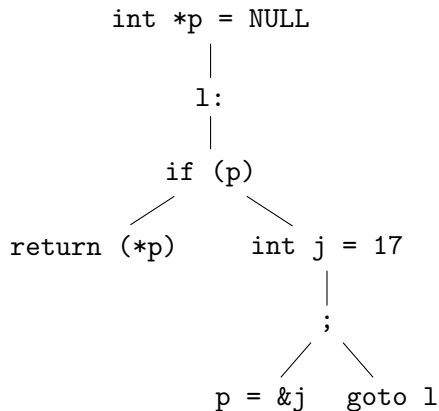
Problem: goto/return/etc. affect memory

Solution:

- ▶ Execute gotos and returns in small steps
 - ▶ Not so much to search for labels, . . .
 - ▶ but to naturally perform required allocations and deallocations
- ▶ Traversal through the AST in the following directions:
 - ▶ ↘ downwards to the next statement
 - ▶ ↗ upwards to the next statement
 - ▶ ↪ `l` to a label `l`: after a `goto l`
 - ▶ ↵ `v` to the top of the statement after a `return`
- ▶ Use a zipper to keep track of the position **and** the stack

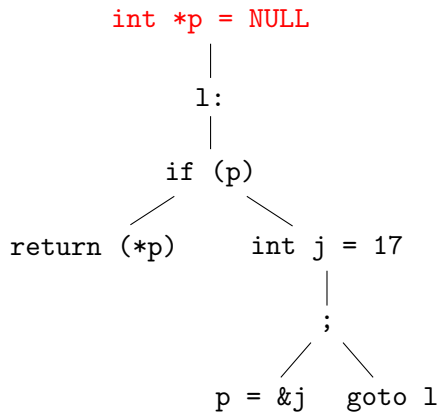
Non-local control flow and block scope variables

Visualization of the operational semantics on an example



Non-local control flow and block scope variables

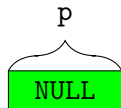
Visualization of the operational semantics on an example



direction:

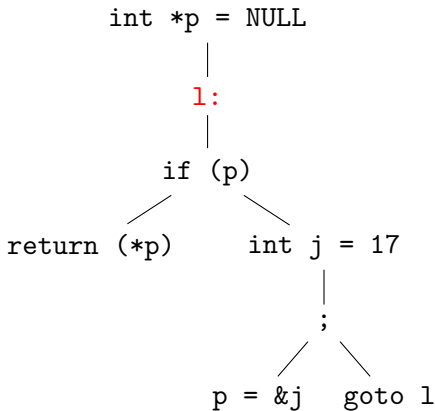


memory:



Non-local control flow and block scope variables

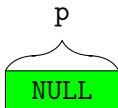
Visualization of the operational semantics on an example



direction:

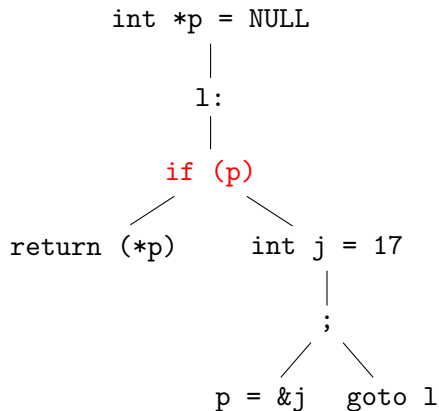


memory:



Non-local control flow and block scope variables

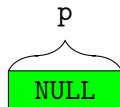
Visualization of the operational semantics on an example



direction:

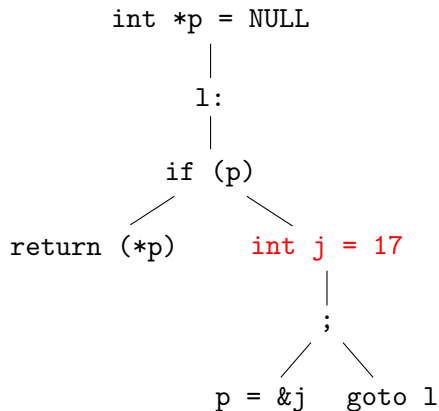


memory:



Non-local control flow and block scope variables

Visualization of the operational semantics on an example



direction:

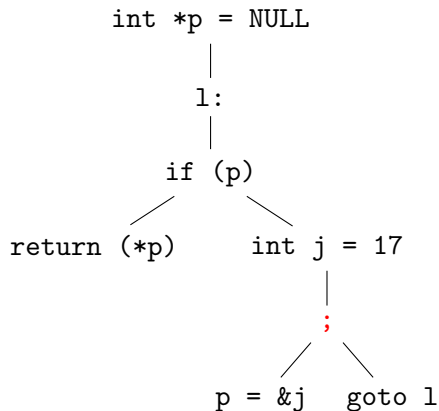


memory:



Non-local control flow and block scope variables

Visualization of the operational semantics on an example



direction:

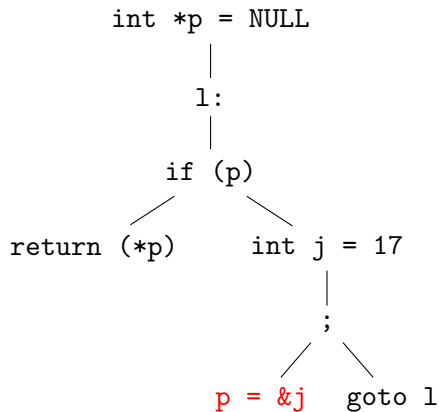


memory:



Non-local control flow and block scope variables

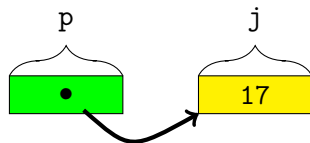
Visualization of the operational semantics on an example



direction:



memory:



Non-local control flow and block scope variables

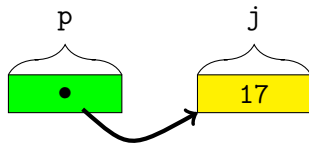
Visualization of the operational semantics on an example



direction:

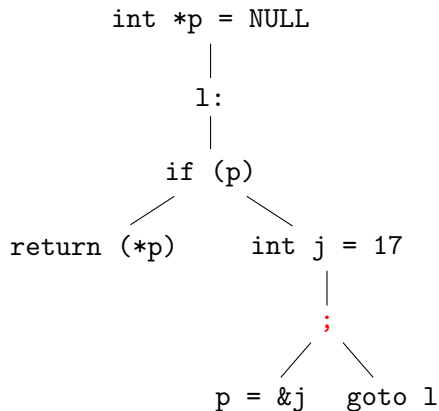


memory:



Non-local control flow and block scope variables

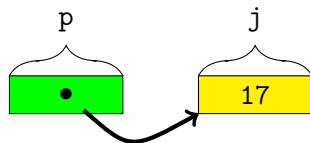
Visualization of the operational semantics on an example



direction:

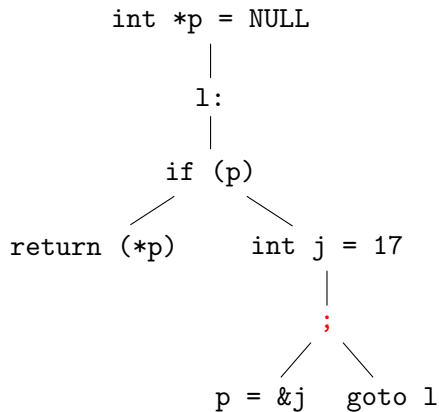


memory:



Non-local control flow and block scope variables

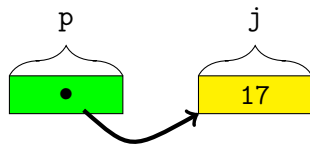
Visualization of the operational semantics on an example



direction:

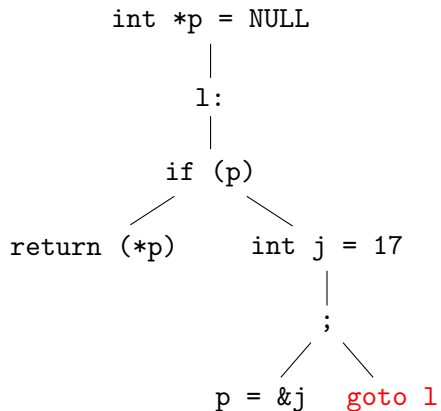


memory:



Non-local control flow and block scope variables

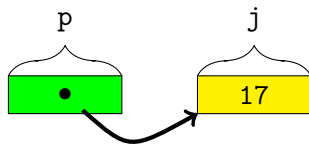
Visualization of the operational semantics on an example



direction:

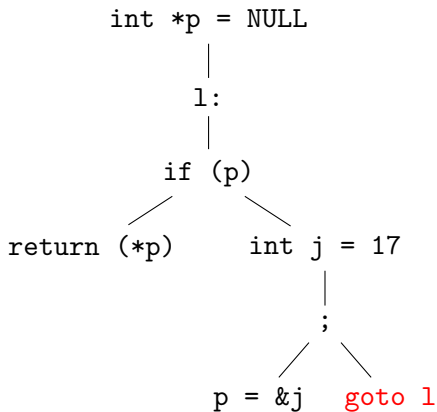


memory:



Non-local control flow and block scope variables

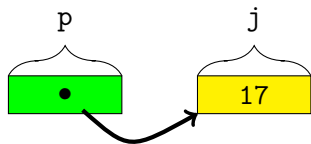
Visualization of the operational semantics on an example



direction:

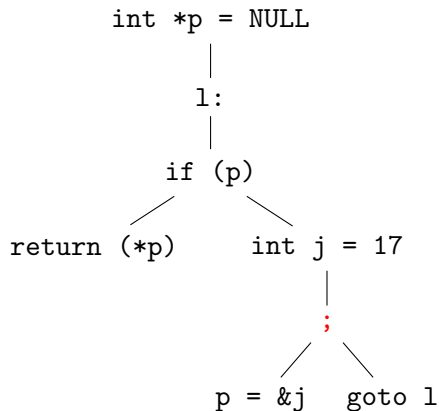


memory:



Non-local control flow and block scope variables

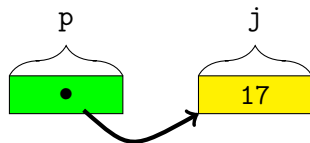
Visualization of the operational semantics on an example



direction:

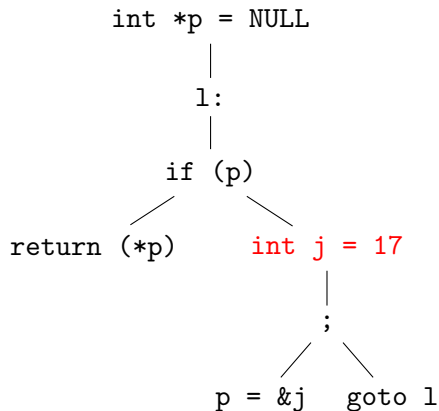


memory:



Non-local control flow and block scope variables

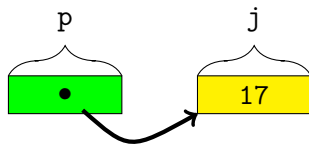
Visualization of the operational semantics on an example



direction:

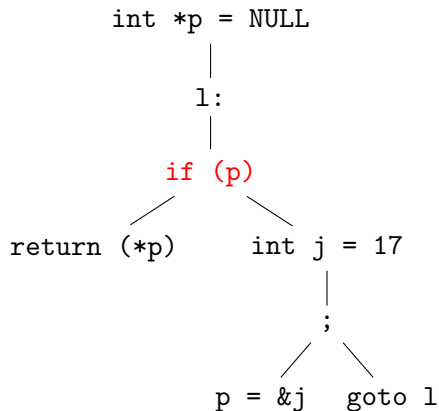


memory:



Non-local control flow and block scope variables

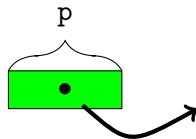
Visualization of the operational semantics on an example



direction:

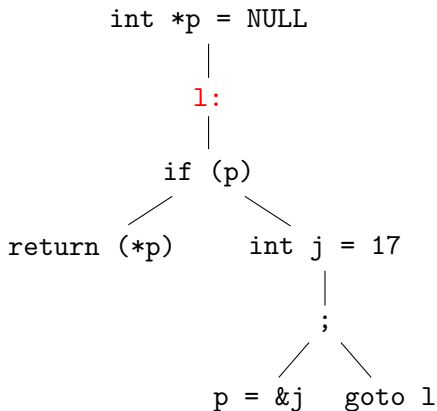


memory:



Non-local control flow and block scope variables

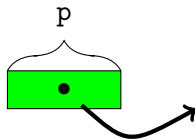
Visualization of the operational semantics on an example



direction:

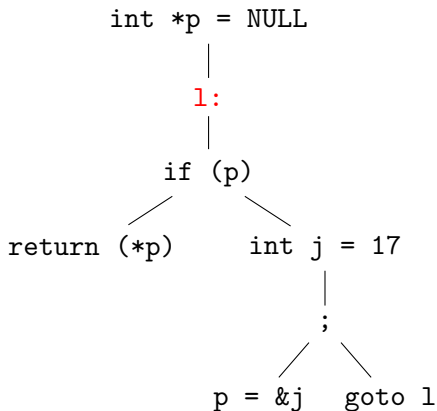


memory:



Non-local control flow and block scope variables

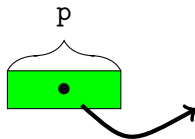
Visualization of the operational semantics on an example



direction:

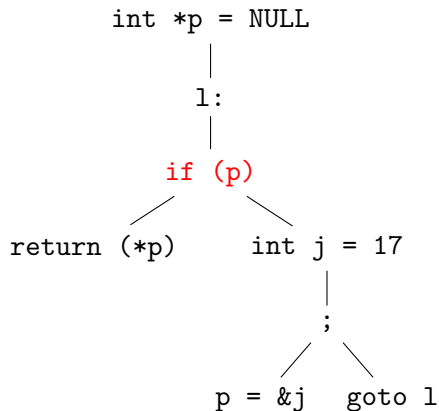


memory:



Non-local control flow and block scope variables

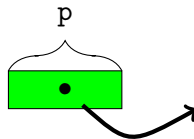
Visualization of the operational semantics on an example



direction:

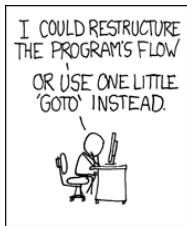


memory:



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

Axiomatic semantics [Krebbbers/Wiedijk, FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

Non-local control flow and block scope variables

Axiomatic semantics [Krebbbers/Wiedijk, FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

where:

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

Non-local control flow and block scope variables

Axiomatic semantics [Krebbbers/Wiedijk, FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

where:

- ▶ $\{P\} s \{Q\}$ is a Hoare triple, as usual
- ▶ Δ maps function names to their pre- and post-conditions

Non-local control flow and block scope variables

Axiomatic semantics [Krebbbers/Wiedijk, FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

where:

- ▶ $\{P\} s \{Q\}$ is a Hoare triple, as usual
 - ▶ Δ maps function names to their pre- and post-conditions
 - ▶ J maps labels to their jumping condition
- When executing a goto l , the assertion $J l$ has to hold

Non-local control flow and block scope variables

Axiomatic semantics [Krebbbers/Wiedijk,FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

where:

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

Non-local control flow and block scope variables

Axiomatic semantics [Krebbbers/Wiedijk,FoSSaCS'13]

CH₂O Hoare sextuples are of the shape

$$\Delta; J; R \vdash \{P\} s \{Q\}$$

where:

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

Remark: the assertions P , Q , J and R correspond to the directions
 \searrow , \nearrow , \curvearrowright and \Uparrow of traversal

Non-local control flow and block scope variables

The block scope variable rule

$$\frac{\Delta; J \uparrow * x_0^T \mapsto -; R \uparrow * x_0^T \mapsto - \vdash \{P \uparrow * x_0^T \mapsto -\} s \{Q \uparrow * x_0^T \mapsto -\}}{\Delta; J; R \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

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

Non-local control flow and block scope variables

The block scope variable rule

$$\frac{\Delta; J \uparrow * x_0^T \mapsto -; R \uparrow * x_0^T \mapsto - \vdash \{P \uparrow * x_0^T \mapsto -\} s \{Q \uparrow * x_0^T \mapsto -\}}{\Delta; J; R \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

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

When leaving a block: the reverse

Non-local control flow and block scope variables

The block scope variable rule

$$\frac{\Delta; J \uparrow * x_0^T \mapsto -; R \uparrow * x_0^T \mapsto - \vdash \{P \uparrow * x_0^T \mapsto -\} s \{Q \uparrow * x_0^T \mapsto -\}}{\Delta; J; R \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

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

When leaving a block: the reverse

Important: symmetry matches gotos going both in and out

Non-local control flow and block scope variables

The block scope variable rule

$$\frac{\Delta; J \uparrow * x_0^T \mapsto -; R \uparrow * x_0^T \mapsto - \vdash \{P \uparrow * x_0^T \mapsto -\} s \{Q \uparrow * x_0^T \mapsto -\}}{\Delta; J; R \vdash \{P\} \text{local}_\tau s \{Q\}}$$

When entering a block:

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

When leaving a block: the reverse

Important: symmetry matches gotos going both in and out

Important: using De Bruijn indexes avoids shadowing

Non-determinism and sequence points

The C quiz, question 3

```
int x = 0, y = 0, *p = &x;

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();
    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();
    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 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);
}
```

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 [Krebbers, POPL'14]

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 [Krebbers, POPL'14]

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 [Krebbers, POPL'14]

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

Operational semantics [Krebbers, POPL'14]

Head reduction for expressions $(e, m) \rightarrow_h (e', m')$

Non-determinism and sequence points

Operational semantics [Krebbers, POPL'14]

Head reduction for expressions $(e, m) \rightarrow_h (e', m')$

- ▶ On assignments: add locked address a to $\Omega_1 \cup \Omega_2$
 $([a]_{\Omega_1} := [v]_{\Omega_2}, m) \rightarrow_h ([v]_{\{a\} \cup \Omega_1 \cup \Omega_2}, \text{lock } a (m[a := v]))$
- ▶ Meanwhile: $\text{lock } a$ makes accesses to a undefined
- ▶ On sequence points: unlock Ω in memory
 $([v]_{\Omega} ? e_2 : e_3, m) \rightarrow_h (e_2, \text{unlock } \Omega m)$ provided ...

Non-determinism and sequence points

Operational semantics [Krebbers, POPL'14]

Head reduction for expressions $(e, m) \rightarrow_h (e', m')$

- ▶ On assignments: add locked address a to $\Omega_1 \cup \Omega_2$
 $([a]_{\Omega_1} := [v]_{\Omega_2}, m) \rightarrow_h ([v]_{\{a\} \cup \Omega_1 \cup \Omega_2}, \text{lock } a (m[a := v]))$
- ▶ Meanwhile: $\text{lock } a$ makes accesses to a undefined
- ▶ On sequence points: unlock Ω in memory
 $([v]_{\Omega} ? e_2 : e_3, m) \rightarrow_h (e_2, \text{unlock } \Omega m)$ provided ...

Gives a local treatment of sequence points

Non-determinism and sequence points

Operational semantics [Krebbers, POPL'14]

Head reduction for expressions $(e, m) \rightarrow_h (e', m')$

- ▶ On assignments: add locked address a to $\Omega_1 \cup \Omega_2$
 $([a]_{\Omega_1} := [v]_{\Omega_2}, m) \rightarrow_h ([v]_{\{a\} \cup \Omega_1 \cup \Omega_2}, \text{lock } a (m[a := v]))$
- ▶ Meanwhile: $\text{lock } a$ makes accesses to a undefined
- ▶ On sequence points: unlock Ω in memory
 $([v]_{\Omega} ? e_2 : e_3, m) \rightarrow_h (e_2, \text{unlock } \Omega m)$ provided ...

Gives a local treatment of sequence points

Small step reduction $\mathbf{S}(\mathcal{P}, \phi, m) \rightarrow \mathbf{S}(\mathcal{P}', \phi', m')$

- ▶ (\mathcal{P}, ϕ) gives the position/direction in the *whole* program
- ▶ Different ϕ s for expressions, statements, function calls
- ▶ Uses evaluation contexts, *i.e.* if $(e_1, m_1) \rightarrow_h (e_2, m_2)$, then

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

Strict aliasing restrictions

What is aliasing?

Aliasing: multiple pointers referring to the same object

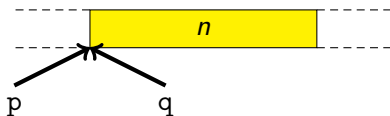
Strict aliasing restrictions

What is aliasing?

Aliasing: multiple pointers referring to the same object

```
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



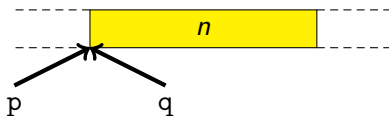
Strict aliasing restrictions

What is aliasing?

Aliasing: multiple pointers referring to the same object

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

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



Optimizing x away is unsound: 17 would be returned

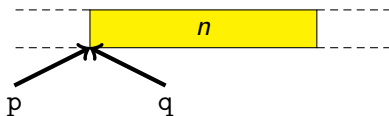
Strict aliasing restrictions

What is aliasing?

Aliasing: multiple pointers referring to the same object

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

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



Optimizing x away is unsound: 17 would be returned

Alias analysis: to determine whether pointers can alias

Strict aliasing restrictions

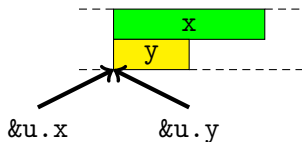
The C quiz, question 4

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 restrictions

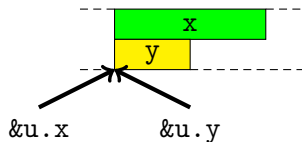
The C quiz, question 4

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);
```



C89 allows `p` and `q` to be aliased, and thus requires `g` to return 3

Strict aliasing restrictions

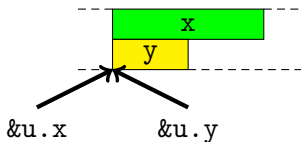
The C quiz, question 4

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);
```



C89 allows p and q to be aliased, and thus requires g to return 3

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 restrictions

How to treat pointers [Krebbbers, CPP'13]

Others (e.g. CompCert)

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

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

Too little information to capture strict aliasing restrictions

Our approach

Strict aliasing restrictions

How to treat pointers [Krebbbers, CPP'13]

Others (e.g. CompCert)

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

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

Too little information to capture strict aliasing restrictions

Our approach

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

Strict aliasing restrictions

How to treat pointers [Krebbbers, CPP'13]

Others (e.g. CompCert)

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

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

Too little information to capture strict aliasing restrictions

Our approach

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

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

Strict aliasing restrictions

How to treat pointers [Krebbbers, CPP'13]

Others (e.g. CompCert)

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

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

Too little information to capture strict aliasing restrictions

Our approach

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

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

A semantics for strict aliasing restrictions

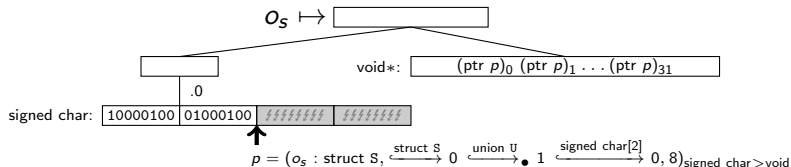
Strict aliasing restrictions

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 restrictions

The C quiz, question 5

```
union U { int x; short y; } u = { .x = 3 };  
short q1() {  
    return u.y;  
}  
short q2() {  
    short *p = &u.y;  
    return *p;  
}
```

Strict aliasing restrictions

The C quiz, question 5

```
union U { int x; short y; } u = { .x = 3 };
short q1() {
    return u.y; // OK
}
short q2() {
    short *p = &u.y;
    return *p; // Undefined
}
```

Type-punning: reading a union using a pointer to another variant

- ▶ C11 vaguely mentions a notion of “visible”
- ▶ GCC only if “the memory is accessed through the union type”

Strict aliasing restrictions

The C quiz, question 5

```
union U { int x; short y; } u = { .x = 3 };
short q1() {
    return u.y; // OK
}
short q2() {
    short *p = &u.y;
    return *p; // Undefined
}
```

Type-punning: reading a union using a pointer to another variant

- ▶ C11 vaguely mentions a notion of “visible”
- ▶ GCC only if “the memory is accessed through the union type”

Formalized by decorating pointers with annotations

Strict aliasing restrictions

Strict-aliasing Theorem

Theorem (Strict-aliasing)

Given:

- ▶ addresses $\Gamma; m \vdash a_1 : \sigma_1$ and $\Gamma; m \vdash a_2 : \sigma_2$
- ▶ with annotations that do not allow type-punning
- ▶ $\sigma_1, \sigma_2 \neq \text{unsigned char}$
- ▶ σ_1 not a subtype of σ_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*) has undefined behavior

Strict aliasing restrictions

Strict-aliasing Theorem

Theorem (Strict-aliasing)

Given:

- ▶ addresses $\Gamma; m \vdash a_1 : \sigma_1$ and $\Gamma; m \vdash a_2 : \sigma_2$
- ▶ with annotations that do not allow type-punning
- ▶ $\sigma_1, \sigma_2 \neq \text{unsigned char}$
- ▶ σ_1 not a subtype of σ_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*) has undefined behavior

Theorem

Compilers can perform type-based alias analysis

CH2O Abstract C

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

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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

- ▶ Named variables to De Bruijn indices

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbbers/Wiedijk, Submitted]

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

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

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

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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

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

- ▶ 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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

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

- ▶ 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`

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

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

- ▶ 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

CH₂O Abstract C

Translation to CH₂O core C in Coq [Krebbers/Wiedijk, Submitted]

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

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

- ▶ 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)

If the translator succeeds, the CH₂O core program is well-typed.

Executable semantics

Goal: define $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$

Executable semantics

Goal: define $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$

Problems:

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

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

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

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

Executable semantics

Goal: define $\text{exec} : \text{state} \rightarrow \mathcal{P}_{\text{fin}}(\text{state})$

Problems:

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

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

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

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

Solutions:

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

Executable semantics

Example of the Coq code

```
Definition cexec ( $\Gamma$  : env Ti) ( $\delta$  : funenv Ti) (S : state Ti) : listset (state Ti) :=
  let 'State k  $\phi$  m := S in
  match  $\phi$  with
  | Stmt  $\searrow$  s =>
    match s with
    | skip => {[ State k (Stmt  $\nearrow$  skip) m ]}
    | goto l => {[ State k (Stmt ( $\curvearrowright$  l) (goto l)) m ]}
    | ...
    | local{ $\tau$ } s =>
      let o := fresh (dom indexset m) in (* use canonical object identifier *)
      {[ State (CLocal o  $\tau$  :: k) (Stmt  $\searrow$  s)
        (mem_alloc  $\Gamma$  o false  $\tau$  m) ]}
    end
  | Expr e =>
    match maybe_EVal e with
    | Some ( $\Omega$ , v) => ...
    | None =>
      '(E, e')  $\leftarrow$  expr_redexes e; (* monadic programming to try each decomposition *)
      match ehexec  $\Gamma$  (get_stack k) e' m with
      | Some (e2, m2) => {[ State k (Expr (subst E e2)) m2 ]}
      | None =>
        match maybe_ECall_redex e' with
        | Some (f,  $\Omega$ s, vs) =>
          {[ State (CFun E :: k) (Call f vs) (mem_unlock ( $\bigcup$   $\Omega$ s) m) ]}
        | _ => {[ State k (Undef (UndefExpr E e')) m ]}
        end
      end
    end
  | ...
```

Executable semantics

Example of the Coq code

```
Definition cexec ( $\Gamma$  : env Ti) ( $\delta$  : funenv Ti) (S : state Ti) : listset (state Ti) :=
  let 'State k  $\phi$  m := S in
  match  $\phi$  with
  | Stmt  $\searrow$  s =>
    match s with
    | skip => {[ State k (Stmt  $\nearrow$  skip) m ]}
    | goto l => {[ State k (Stmt ( $\curvearrowright$  l) (goto l)) m ]}
    | ...
    | local $\{\tau\}$  s =>
      let o := fresh (dom indexset m) in (* use canonical object identifier *)
      {[ State (CLocal o  $\tau$  :: k) (Stmt  $\searrow$  s)
        (mem_alloc  $\Gamma$  o false  $\tau$  m) ]}
    end
  | Expr e =>
    match maybe_EVal e with
    | Some ( $\Omega$ , v) => ...
    | None =>
      '(E,e')  $\leftarrow$  expr_redexes e; (* monadic programming to try each decomposition *)
      match ehexec  $\Gamma$  (get_stack k) e' m with
      | Some (e2,m2) => {[ State k (Expr (subst E e2)) m2 ]}
      | None =>
        match maybe_ECall_redex e' with
        | Some (f,  $\Omega$ s, vs) =>
          {[ State (CFun E :: k) (Call f vs) (mem_unlock ( $\bigcup$   $\Omega$ s) m) ]}
        | _ => {[ State k (Undef (UndefExpr E e')) m ]}
        end
      end
    end
  | ...
```

Executable semantics

Example of the Coq code

```
Definition cexec ( $\Gamma$  : env Ti) ( $\delta$  : funenv Ti) (S : state Ti) : listset (state Ti) :=
  let 'State k  $\phi$  m := S in
  match  $\phi$  with
  | Stmt  $\searrow$  s =>
    match s with
    | skip => {[ State k (Stmt  $\nearrow$  skip) m ]}
    | goto l => {[ State k (Stmt ( $\curvearrowright$  l) (goto l)) m ]}
    | ...
    | local{ $\tau$ } s =>
      let o := fresh (dom indexset m) in (* use canonical object identifier *)
      {[ State (CLocal o  $\tau$  :: k) (Stmt  $\searrow$  s)
        (mem_alloc  $\Gamma$  o false  $\tau$  m) ]}
    end
  | Expr e =>
    match maybe_EVal e with
    | Some ( $\Omega$ , v) => ...
    | None =>
      '(E,e')  $\leftarrow$  expr_redexes e; (* monadic programming to try each decomposition *)
      match ehexec  $\Gamma$  (get_stack k) e' m with
      | Some (e2,m2) => {[ State k (Expr (subst E e2)) m2 ]}
      | None =>
        match maybe_ECall_redex e' with
        | Some (f,  $\Omega$ s, vs) =>
          {[ State (CFun E :: k) (Call f vs) (mem_unlock ( $\bigcup$   $\Omega$ s) m) ]}
        | _ => {[ State k (Undef (UndefExpr E e')) m ]}
        end
      end
    end
  | ...
```

Executable semantics

Soundness and completeness [Krebbers/Wiedijk, Submitted]

Theorem (Soundness)

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

Executable semantics

Soundness and completeness [Krebbbers/Wiedijk, Submitted]

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 $f : \text{index} \rightarrow \text{option index}$.

Executable semantics

Soundness and completeness [Krebbbers/Wiedijk, Submitted]

Theorem (Soundness)

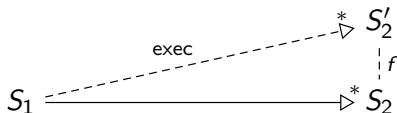
If $S_2 \in \text{exec } S_1$, then $S_1 \rightarrow S_2$.

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:



Executable semantics

Demo

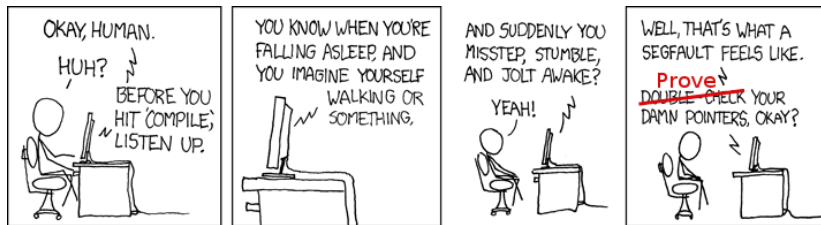
```
$ cat foo.c
int main(void) {
    int x = 10, *p = &x;
    return *p + *p + *p + *p + *p + *p;
}
$ ./ch2o -t foo.c
.....,.....,,.....+;!!?*%$$$$$$$$$%*?!!;:+-,..
.....
"" 60
```

Conclusion

C11 is not a simple language

Questions

Sources: <http://robertkrebbers.nl/research/ch2o/>



<http://xkcd.com/371/>