# A precise and abstract memory model for C using symbolic values

Frédéric Besson    Sandrine Blazy    **Pierre Wilke**

Rennes, France

# What does this program do?

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  //
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  //
  assert( checksum( (q » 4) « 4 ) == (q & 0xF) );
  int * r = ( q » 4 ) « 4;
  //
  return *r;
}
```

## What does this program do?

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  //
  assert( checksum( (q >> 4) << 4 ) == (q & 0xF) );
  int * r = ( q >> 4 ) << 4;
  //
  return *r;
}
```

# What does this program do?

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  // q = 0xTUVWXYZ5
  assert( checksum( (q >> 4) << 4 ) == (q & 0xF) );
  int * r = ( q >> 4 ) << 4;
  //
  return *r;
}
```

# What does this program do?

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  // q = 0xTUVWXYZ5
  assert( checksum( (q » 4) « 4 ) == (q & 0xF) );
  int * r = ( q » 4 ) « 4;
  // r = 0xTUVWXYZ0 == p
  return *r;
}
```

# What does this program do?

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  // q = 0xTUVWXYZ5
  assert( checksum( (q » 4) « 4 ) == (q & 0xF) );
  int * r = ( q » 4 ) « 4;
  // r = 0xTUVWXYZ0 == p
  return *r;
}
```

"Real life"
Terminates and outputs 42

# What does this program do?

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  // q = 0xTUVWXYZ5
  assert( checksum( (q » 4) « 4 ) == (q & 0xF) );
  int * r = ( q » 4 ) « 4;
  // r = 0xTUVWXYZ0 == p
  return *r;
}
```

ISO C Standard
Undefined behavior

"Real life"
Terminates and outputs 42

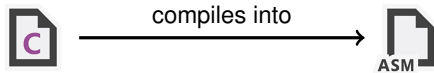# What does this program do?

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  // p = 0xTUVWXYZ0
  *p = 42;
  int * q =  p | (checksum(p) & 0xF);
  // q = 0xTUVWXYZ5
  assert( checksum( (q » 4) « 4 ) == (q & 0xF) );
  int * r = ( q » 4 ) « 4;
  // r = 0xTUVWXYZ0 == p
  return *r;
}
```

ISO C Standard
Undefined behavior

"Real life"
Terminates and outputs 42

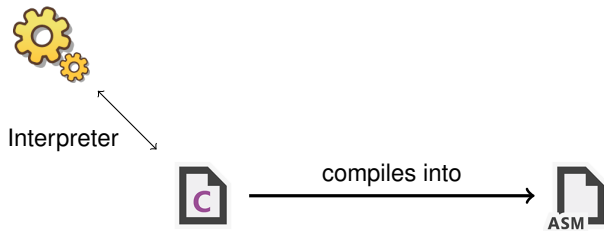This work: an executable semantics for this kind of programs

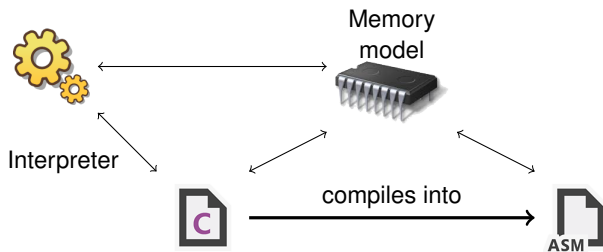# State of the art: formal semantics for C

- Cholera [ESOP'99]: first formal semantics of C
- CompCert [JAR'09]: large subset of C, executable, abstract memory model
- KCC [POPL'12]: large subset of C, executable, abstract memory model
- Krebbers [POPL'14]: closer to the ISO C standard
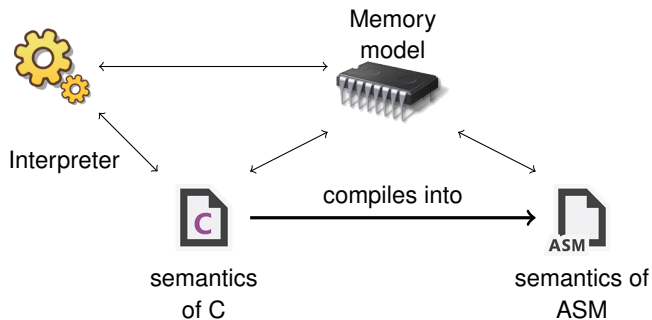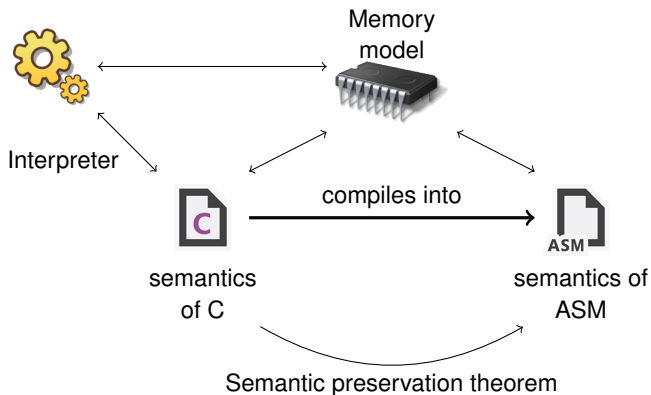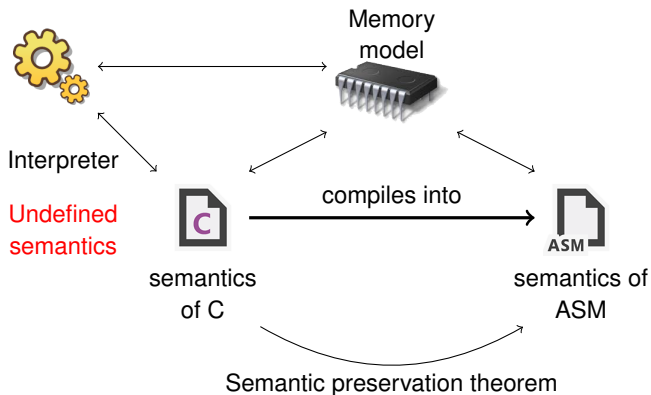- Our work: CompCert + more defined semantics + **low level** memory model

# CompCert



compiles into
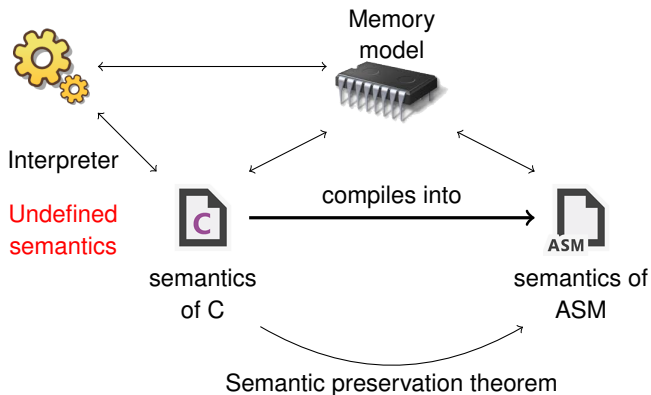
# CompCert



Interpreter

compiles into

# CompCert

# CompCert

# CompCert



Memory model

Interpreter

compiles into

semantics of C

semantics of ASM

Semantic preservation theorem

**If** the program has well-defined semantics
**Then** the compiler does not introduce bugs

COQ

# CompCert



**If** the program has well-defined semantics
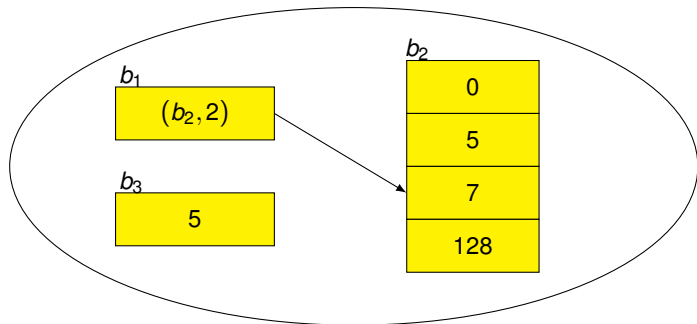**Then** the compiler does not introduce bugs

COQ

# CompCert

# CompCert's memory model

- set of disjoint blocks
- Values:

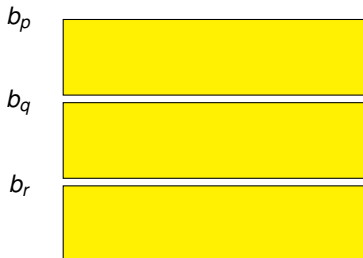$$val ::= i \mid (b, o) \mid \texttt{Vundef}$$

- $load(M, b, o) = \lfloor v \rfloor$
- $store(M, b, o, v) = \lfloor M' \rfloor$
- $alloc(M, lo, hi) = (M', b)$
- $free(M, b) = \lfloor M' \rfloor$



"Good variable" properties:
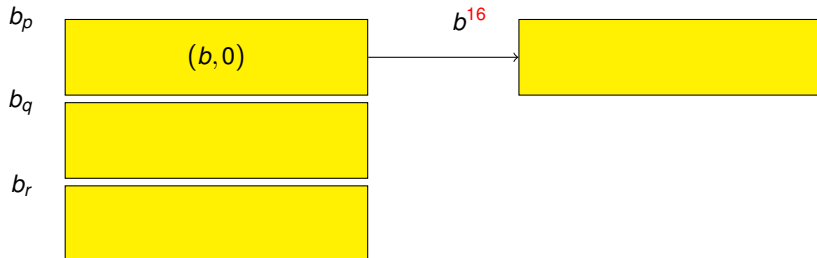$load(store(M, b, o, v), b, o) = v$

# Back to the example

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q » 4) « 4;
  return *r;
}
```

$b_p$

$b_q$
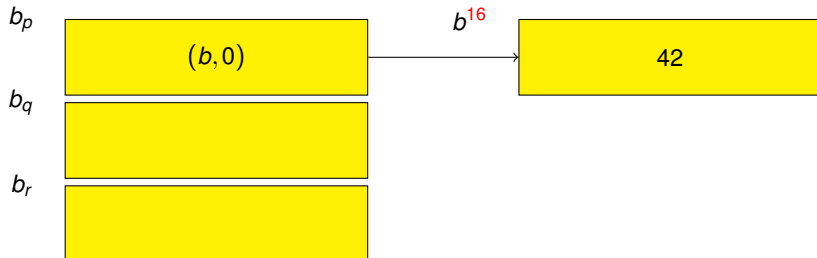
$b_r$

# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```

# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```

# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q » 4) « 4;
  return *r;
}
```



$b_p$
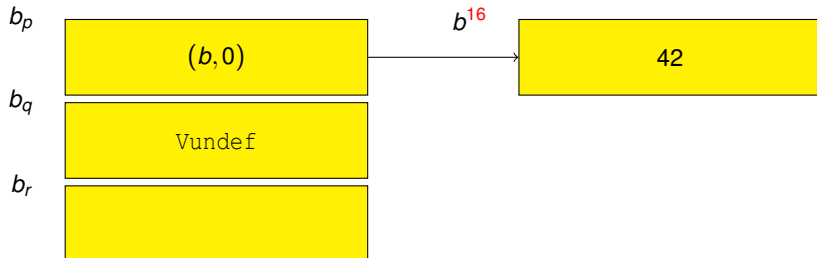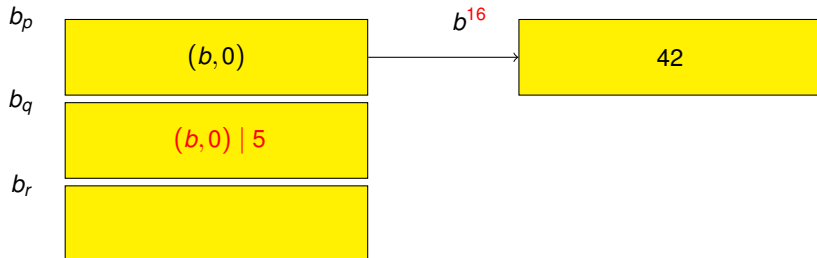
$(b, 0)$

$b^{16}$

42

$b_q$

Vundef

$b_r$

# Back to the example

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q » 4) « 4;
  return *r;
}
```

# Back to the example

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```

# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```
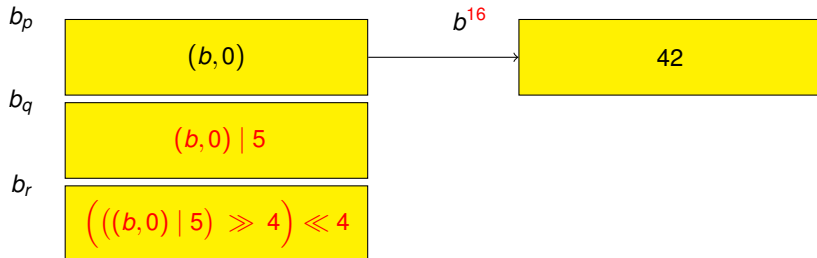
# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```
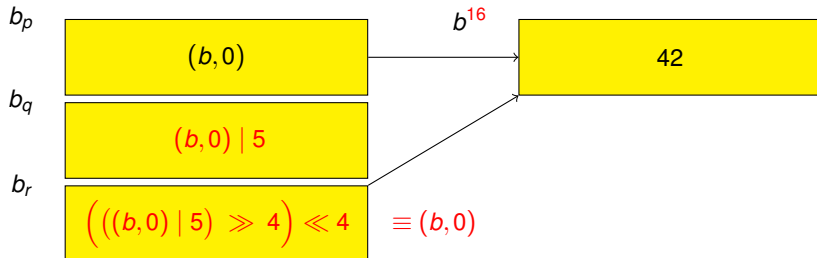


$b_p$

$b_q$

$b_r$

$(b,0)$

$(b,0) \mid 5$

$\left( \left( (b,0) \mid 5 \right) \gg 4 \right) \ll 4$

$\equiv (b,0)$

$b^{16}$ ← Alignment Constraints

42

# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```
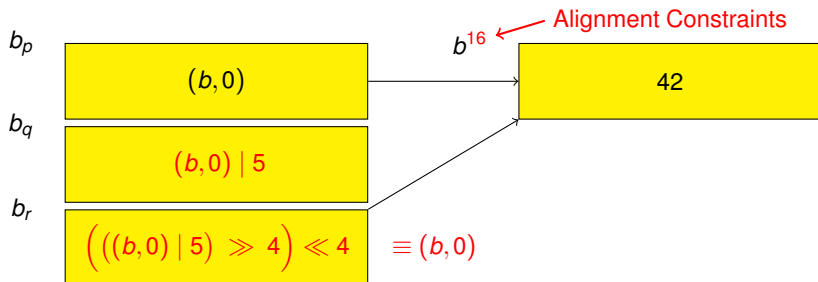
# Back to the example

```c
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q >> 4) << 4;
  return *r;
}
```

# Outline
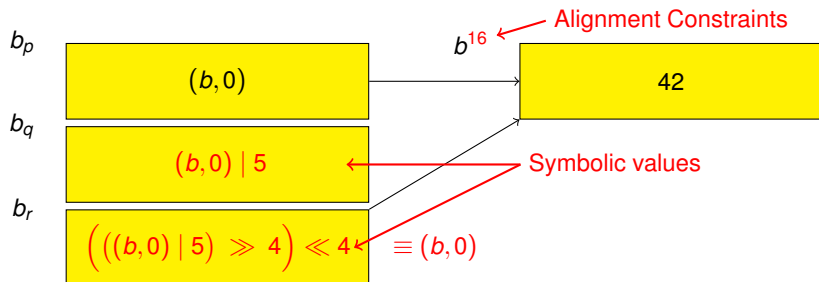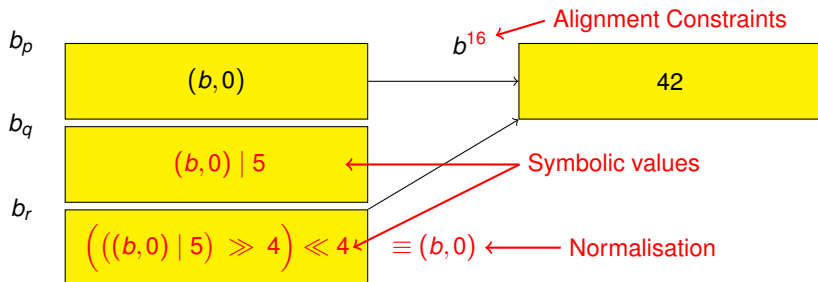
1. A C semantics with symbolic values

2. Normalisation: specification and implementation

3. Experimental evaluation

# Outline

# CompCert's memory model with symbolic values



- Alignment
  - $alloc(M, lo, hi, \textbf{mask}) = (M', b)$

  $$A(b) \ \& \ mask = A(b)$$

- Symbolic values
  - $sv ::= val \ | \ op_1 \ sv \ | \ sv \ op_2 \ sv$
  - $load(M, b, o) = \lfloor \textbf{sv} \rfloor$
  - $store(M, b, o, \textbf{sv}) = \lfloor M' \rfloor$

# CompCert's memory model with symbolic values



$$\text{normalise} : \textit{mem} \rightarrow \textit{sv} \rightarrow \textit{option val}$$

When do we need to normalise symbolic values?
- Memory accesses:
  - **return** *r;
  - *p = 42;
- Control flow:
  - **if** (c) { ... } **else** { ... }

# Adapting the CompCert semantics

Semantic rules

$$\frac{\vdash a, M \to (b, o) \quad \text{load}(M, b, o) = \lfloor v \rfloor}{\vdash *a, M \to v}$$

$$\frac{\vdash a, M \to (b, o) \quad \text{store}(M, b, o, v) = \lfloor M' \rfloor}{\vdash *a = v, M \to M'}$$

$$\frac{\vdash \text{is\_true}(a)}{\vdash \text{if } a \text{ then } s_1 \text{ else } s_2, M \to s_1, M}$$

# Adapting the CompCert semantics

Semantic rules

$$\frac{\vdash a, M \to sv_a \quad \text{normalise}(M, sv_a) = \lfloor (b, o) \rfloor}{\vdash *a, M \to sv}$$

$$\frac{\vdash a, M \to sv_a \quad \text{normalise}(M, sv_a) = \lfloor (b, o) \rfloor}{\text{store}(M, b, o, sv) = \lfloor M' \rfloor}$$
$$\vdash *a = sv, M \to M'$$

$$\frac{\vdash \text{normalise}(M, a) = \lfloor i \rfloor \quad \text{is\_true}(i)}{\vdash \text{if } a \text{ then } s_1 \text{ else } s_2, M \to s_1, M}$$

Interpreter

Symbolic
Memory model

Semantics
of C
with normalisation

# Outline

# Evaluation of symbolic values

- Input:
  - *a memory mapping $A : block \rightarrow int_{32}$*
  - *a symbolic value sv*
- Output: the set of machine integers that *sv* evaluates to.

$$\llbracket \cdot \rrbracket_A : sv \rightarrow \mathcal{P}(int_{32})$$

$$\overline{i \in \llbracket i \rrbracket_A} \qquad \overline{A(b) + o \in \llbracket (b,o) \rrbracket_A} \qquad \overline{n \in \llbracket \text{Vundef} \rrbracket_A}$$

$$\frac{v_1 \in \llbracket e_1 \rrbracket_A \quad \text{eval\_unop}(op_1, v_1) = \lfloor v \rfloor}{v \in \llbracket op_1 \ e_1 \rrbracket_A}$$

$$\frac{v_1 \in \llbracket e_1 \rrbracket_A \quad v_2 \in \llbracket e_2 \rrbracket_A \quad \text{eval\_binop}(op_2, v_1, v_2) = \lfloor v \rfloor}{v \in \llbracket e_1 \ op_2 \ e_2 \rrbracket_A}$$
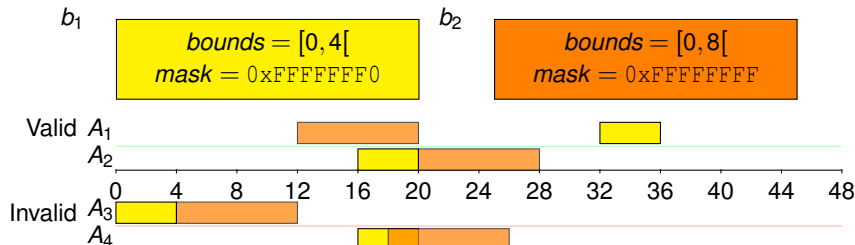
# Valid memory mapping : $A \models M$

A memory mapping $A : block \rightarrow int_{32}$ is valid for memory $M$ iff:

1. addresses from distinct blocks do not overlap,
2. the address of a block satisfies its alignment constraints:

$$A(b) \& \texttt{mask}(M, b) = A(b)$$

3. valid addresses are not null.

**Example:**

# Normalisation: specification

If $\text{normalise}(M, sv) = \lfloor v \rfloor$ then:

- $v \neq \text{Vundef}$
- $\forall A \vDash M, [\![sv]\!]_A = [\![v]\!]_A$

**Example:** Consider block $b$ with bounds $[0, 4[$ and 16-byte aligned.

$sv = \Big( \big( (b, 0) \mid 5 \big) \gg 4 \Big) \ll 4$

$v = (b, 0)$

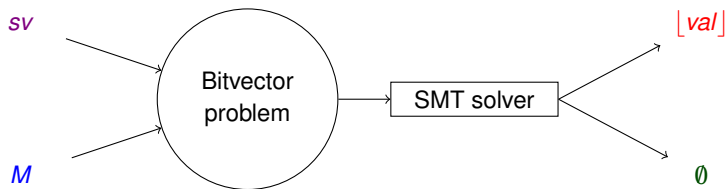| A | $[\![sv]\!]_A$ | $[\![v]\!]_A$ |
|---|---|---|
| $\{b \mapsto 16\}$ | $(((16 + 0) \mid 5) \gg 4) \ll 4 \quad = 16$ | 16 |
| $\{b \mapsto 32\}$ | $(((32 + 0) \mid 5) \gg 4) \ll 4 \quad = 32$ | 32 |
| $\{b \mapsto 16k\}$ | $(((16k + 0) \mid 5) \gg 4) \ll 4 \quad = 16k$ | 16k |

# Normalisation: implementation

```
int main(){
  int * p = (int *) malloc (sizeof (int));
  *p = 42;
  int * q = p | 5;
  int * r = (q » 4) « 4;
  return *r;
}
```



$b_p$

$b^{16}$

$(b, 0)$

42

$b_r$

$\left( \left( (b, 0) \mid 5 \right) \gg 4 \right) \ll 4$

normalise

$(b, 0)$

# How to compute normalisation?

$$\text{normalise}(M, sv) = \left\{ \begin{array}{c} \lfloor val \rfloor \\ 0 \end{array} \right.$$

## Normalisation in our example

Let $A$ be a valid memory layout for memory $M$: $A \vDash M$

$$sv_r = \Big( \big( (b,0) \mid 5 \big) \gg 4 \Big) \ll 4$$

Translation into a bitvector expression

$$bv_r = \Big( \big( (A(b) + 0) \mid 5 \big) \gg 4 \Big) \ll 4$$

**Goal:** find a unique model $(b,i)$ such that: $bv_r = A(b) + i$

Two steps:

- find a model $(b_0, i_0)$ such that $bv_r = A(b_0) + i_0$
  - $(b_0, i_0) = (b, 0)$
- check that this solution is unique: **unsat**$(bv_r = A(b) + i \wedge b \neq b_0)$
  - the solution is indeed unique
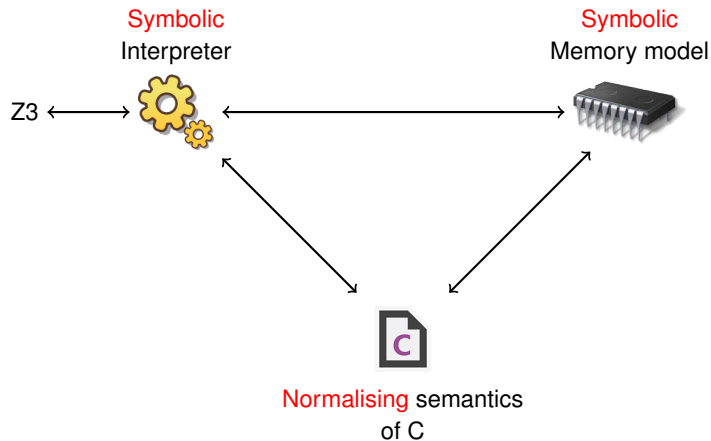  - the normalisation returns $\lfloor (b,0) \rfloor$

# Outline

# Experiments

Real life programs

- a `malloc` implementation (`dlmalloc.c`)
- excerpts from a C library : Public Domain C Library
- excerpts from a cryptographic library: Networking and Cryptographic library
- hand-written C programs

What kind of symbolic values do these programs trigger?

- bitwise operations on pointer
- use of undefined values

# Putting the pieces together



Symbolic
Interpreter

Symbolic
Memory model

Z3

Normalising semantics
of C

Coq development available at `http://www.irisa.fr/celtique/ext/csem/`

# Conclusion

Results:

- non-regression
- more programs have their expected semantics
- limits: some undefined behaviors are not captured by our semantics
  - pointer comparison

Ongoing work:

- Proofs of the memory model
  - *e.g. load*(*store*($M, b, o, $ **sv**)$, b, o$) = **sv**
- Proof of the whole CompCert compiler?
  - memory injections redefined

# Questions?