## **Towards Certified** Separate Compilation for Concurrent Programs

**Hanru Jiang**<sup>\*</sup> Hongjin Liang<sup>†</sup> Siyang Xiao\* Junpeng Zha\* Xinyu Feng†

\* University of Science and Technology of China † Nanjing University

# **Compilers are NOT Trustworthy**

### Finding and Understanding Bugs in C Compilers

Eric Eide John Regehr Xuejun Yang Yang Chen

> University of Utah, School of Computing { jxyang, chenyang, eeide, regehr }@cs.utah.edu

- 11 open-source/commercial compilers were tested
- Found **325** bugs, in **EVERY** compiler!

[PLDI 2011]

# **Compilers are NOT Trustworthy**

### Finding and Understanding Bugs in C Compilers

John Regehr Xuejun Yang Yang Chen Eric Eide

> University of Utah, School of Computing {jxyang, chenyang, eeide, regehr}@cs.utah.edu

- 11 open-source/commercial compilers were tested
- Found **325** bugs, in **EVERY** compiler!

Verification of compiler correctness helps:

"The striking thing about our CompCert results is that the middle end bugs we found in all other compilers are absent."

[PLDI 2011]

## **Compilation Correctness**



### **Correct(Compiler) :**

 $\forall S, T . T = Compiler(S) \implies T \subseteq S$ 

Semantic preservation: T has no more observable behaviors (e.g. I/O events by print) than S.

## **Compiler Verification**

# **Compiler Verification**

- Leroy'06: Formal certification of a compiler back-end
- Lochbihler'10: Verifying a compiler for Java threads
- Myreen'10: Verified just-in-time compiler on x86
- Sevcik et al.'11: Relaxed-memory concurrency and verified compilation
- Zhao et al.'13: Formal verification of SSA-based optimizations for LLVM
- Kumar et al.'14: CakeML: A verified implementation of ML
- Stewart et al.'15: Compositional CompCert
- Kang et al.'16: Lightweight Verification of Separate Compilation



# **Compiler Verification**

- Leroy'06: Formal certification of a compiler back-end
- Lochbihler'10: Verifying a compiler for Java threads
- Myreen'10: Verified just-in-time compiler on x86

### Limited support of separate compilation and concurrency!

- Numar et al. 14: CakeIVIL: A vermed implementation of IVIL
- Stewart et al.'15: Compositional CompCert
- Kang et al.'16: Lightweight Verification of Separate Compilation



### Source (e.g. C)



### Real-world programs may consist of multiple components, which will be compiled independently.

## Separate Compilation









## which will be compiled independently.



Real-world programs may consist of multiple components,

### Source (e.g. C)



## Separate Compilation



## Separate Compilation









## Separate Compilation of Concurrent Programs



## Separate Compilation of Concurrent Programs



**Parallel Composition S2 Compiler-2 Parallel Composition T2** 



Can we reuse existing certified compilers (e.g. CompCert) for separate compilation of concurrent programs?



### YES, for data-race-free (DRF) concurrent programs







No race



No race







r1 = 1;	r2 = 2;
r1 = r1 + 1;	r2 = r2 + 1;
<pre>lock();</pre>	<pre>lock();</pre>
x = 1;	x = 2;
v = x + 1:	v = x + 1:

No race

**Non-preemptive:** yield control at

Plausible, but need to address several key challenges







• How to formulate **DRF** in **language independent** manner?





- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?



- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?



- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?



- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?

- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?
- How to support benign-race and relaxed memory models?

- How to formulate DRF in language independent manner?
- How to prove DRF-preservation, compositionally?
- How to support benign-race and relaxed memory models?



"... synchronization primitives are commonly implemented with assembly code that has data races."

— — Hans-J. Boehm, HotPar'11

- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
    - With both external function calls & multi-threaded code

- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
    - With both external function calls & multi-threaded code
- Framework extension:
  - Supports x86-TSO + confined benign-races

- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
    - With both external function calls & multi-threaded code
- Framework extension:
  - Supports x86-TSO + confined benign-races
- **CAS**CompCert:
  - Extends CompCert with Concurrency + Abstraction + Separate compilation
  - Reuses considerable amount of CompCert proofs
  - Racy x86-TSO impl. of locks as synchronization library
# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

### Data-race: read-write / write-write conflicts



Data-race: read-write / write-write conflicts

### **Data-race: read-write / write-write conflicts** Why language-independent? To support cross-language interaction





**Data-race: read-write / write-write conflicts** Why language-independent? To support cross-language interaction abstract away lang. details e.g. interaction semantics



May in different languages





**Data-race: read-write / write-write conflicts** Why language-independent? To support cross-language interaction abstract away lang. details e.g. interaction semantics



May in different languages





### **Data-race: read-write / write-write conflicts** Why language-independent?



- How to formulate DRF
- if we do not even know the concrete reads/writes?





Defined in terms of footprint disjointness

write-set

## **DRF(S1 || S2)**

footprints δ ::= (rs, ws)

read-set

**Defined in terms of** footprint disjointness

- footprints  $\delta ::= (rs, ws)$ write-set read-set
- well-defined language extensional characterization of footprints

**Defined in terms of** footprint disjointness

## **DRF(S1 || S2)**

- footprints  $\delta ::= (rs, ws)$ write-set read-set
- well-defined language extensional characterization of footprints

#### rs

**Defined in terms of** footprint disjointness

## **DRF(S1 || S2)**

- footprints  $\delta ::= (rs, ws)$ write-set read-set
- well-defined language extensional characterization of footprints

#### rs

#### rs

Defined in terms of footprint disjointness

# footprints δ ::= (rs, ws) read-set write-set

**DRF(S1 || S2)** 

 well-defined language extensional characterization of footprints



rs

Defined in terms of footprint disjointness

# footprints δ ::= (rs, ws) read-set write-set

**DRF(S1 || S2)** 

 well-defined language extensional characterization of footprints



#### **Arbitrarily different**

Defined in terms of footprint disjointness

# footprints δ ::= (rs, ws) read-set write-set

**DRF(S1 || S2)** 

 well-defined language extensional characterization of footprints



**Defined in terms of** footprint disjointness

- footprints  $\delta ::= (rs, ws)$ write-set read-set
- well-defined language extensional characterization of footprints

**Defined in terms of** footprint disjointness

### **DRF(S1 || S2)**

• footprints  $\delta ::= (rs, ws)$ read-set write-set

 well-defined language extensional characterization of footprints

**Definition 1** (Well-Defined Languages). wd(tl) iff, for any execution step  $F \vdash (\kappa, \sigma) \xrightarrow{\iota}_{\mathcal{S}} (\kappa', \sigma')$  in this language, all of the following hold (some auxiliary definitions are in Fig. 6):

- (1) forward( $\sigma, \sigma'$ );
- (2) LEffect( $\sigma, \sigma', \delta, F$ );
- (3) For any  $\sigma_1$ , if LEqPre $(\sigma, \sigma_1, \delta, F)$ , then there exists  $\sigma'_1$ such that  $F \vdash (\kappa, \sigma_1) \xrightarrow{\iota}_{\mathcal{S}} (\kappa', \sigma'_1)$  and  $\mathsf{LEqPost}(\sigma', \sigma'_1, \delta, F)$ .
- (4) Let  $\delta_0 = \bigcup \{ \delta \mid \exists \kappa', \sigma' \colon F \vdash (\kappa, \sigma) \xrightarrow{\tau}_{\delta} (\kappa', \sigma') \}$ . For any  $\sigma_1$ , if LEqPre( $\sigma, \sigma_1, \delta_0, F$ ), then for any  $\kappa'_1, \sigma'_1, \iota_1, \delta_1$ ,  $F \vdash (\kappa, \sigma_1) \xrightarrow{\iota_1}_{\delta_1} (\kappa'_1, \sigma'_1) \implies \exists \sigma' . F \vdash (\kappa, \sigma) \xrightarrow{\iota_1}_{\delta_1} (\kappa'_1, \sigma').$



# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

# Compositional CompCert's Argument...

#### DRF(S1 || S2)

T1 = Comp(S1)

T2 = Comp(S2)

## Compositional CompCert's Argument...

#### DRF(S1 || S2)

**7**=**7**

















#### DRF(T1 || T2)





### 2 DRF(T1 || T2)





# DRF(T1 || T2) ? DRF(S1 || S2)





#### **DRF(T1 || T2)**

### How to prove **DRF-preservation?**

# DRF(T1 || T2) 2 DRF(S1 || S2)

# Our Ideas

#### **T1 || T2** $\subseteq$ **S1 || S2** $T1 | T2 \subseteq S1 | S2$

# DRF(T1 || T2) ? **DRF(S1 || S2)**

# Our Ideas

#### $T1 \parallel T2 \subseteq S1 \parallel S2$ U Triv Footprint-preserving $T1 | T2 \subseteq S1 | S2$ simulation T1 | T2 ≤ S1 | S2

# DRF(T1 || T2) ? **DRF(S1 || S2)**

# Our Ideas





### DRF(T1 || T2)



# Our Ideas





### DRF(T1 || T2)



# Our Ideas














































# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

# Outline of this Talk

- Language-independent DRF formulation
- DRF-preservation and key proof structures
- Supporting x86-TSO and confined benign-races in CASCompCert

### **DRF Imposes Strong Restriction on Libraries**



DRF

### **DRF Imposes Strong Restriction on Libraries**







### [spin-lock impl. in Linux 2.6]











- Confined benign-races:



Racy libraries and client code run in separate memory regions

- Confined benign-races:

  - Racy libraries and client code run in separate memory regions Client code be well-synchronized



- Confined benign-races:
  - Racy libraries and client code run in separate memory regions
  - Client code be well-synchronized
  - Racy libraries have race-free abstraction



**Source ₽** 











### race-free abstraction of spin-locks for synchronization



**Source ₽** 




















### Supporting Confined Benign-Race & x86-TSO in Our CASCompCert



### Supporting Confined Benign-Race & x86-TSO in Our CASCompCert



### Supporting Confined Benign-Race & x86-TSO in Our CASCompCert









## Verified 12/20 Passes in CASCompCert







## Verified 12/20 Passes in CASCompCert





#### Including all the translation passes from Clight to x86



Compilation passes	Spec		Proof	
	CompCert	Ours	CompCert	Ours
Cshmgen	515	1021	1071	1503
Cminorgen	753	1556	1152	1251
Selection	336	500	647	783
RTLgen	428	543	821	862
Tailcall	173	328	275	405
Renumber	86	245	117	358
Allocation	704	785	1410	1700
Tunneling	131	339	166	475
Linearize	236	371	349	733
CleanupLabels	126	387	161	388
Stacking	730	1038	1108	2135
Asmgen	208	338	571	1128

Compilation passes	Spec		
	CompCert	Ours	9
Cshmgen	515	1021	
Cminorgen	753	1556	
Selection	336	500	
RTLgen	428	543	
Tailcall	173	328	
Renumber	86	245	
Allocation	704	785	
Tunneling	131	339	
Linearize	236	371	
CleanupLabels	126	387	
Stacking	730	1038	
Asmgen	208	338	



#### **Compiler verification: 100 - 400 more lines of Coq proof for most passes**





#### **Compiler verification: 100 - 400 more lines of Coq proof for most passes**





Framework impl.: > 60k LoC, ~ 1 person year

#### **Compiler verification: 100 - 400 more lines of Coq proof for most passes**



- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
  - Well-defined language for language-independent DRF
  - Footprint-preserving simulation for DRF preservation

- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
  - Well-defined language for language-independent DRF
  - Footprint-preserving simulation for DRF preservation
- Framework extension:
  - Support x86-TSO + confined benign-races

- Language independent verification framework
  - Key semantics components + proof structures
  - Supports separate compilation for race-free concurrent programs
  - Well-defined language for language-independent DRF
  - Footprint-preserving simulation for DRF preservation
- Framework extension:
  - Support x86-TSO + confined benign-races
- **CAS**CompCert:
  - Reused considerable amount of CompCert proofs
  - Racy x86-TSO impl. of locks as synchronization library

