

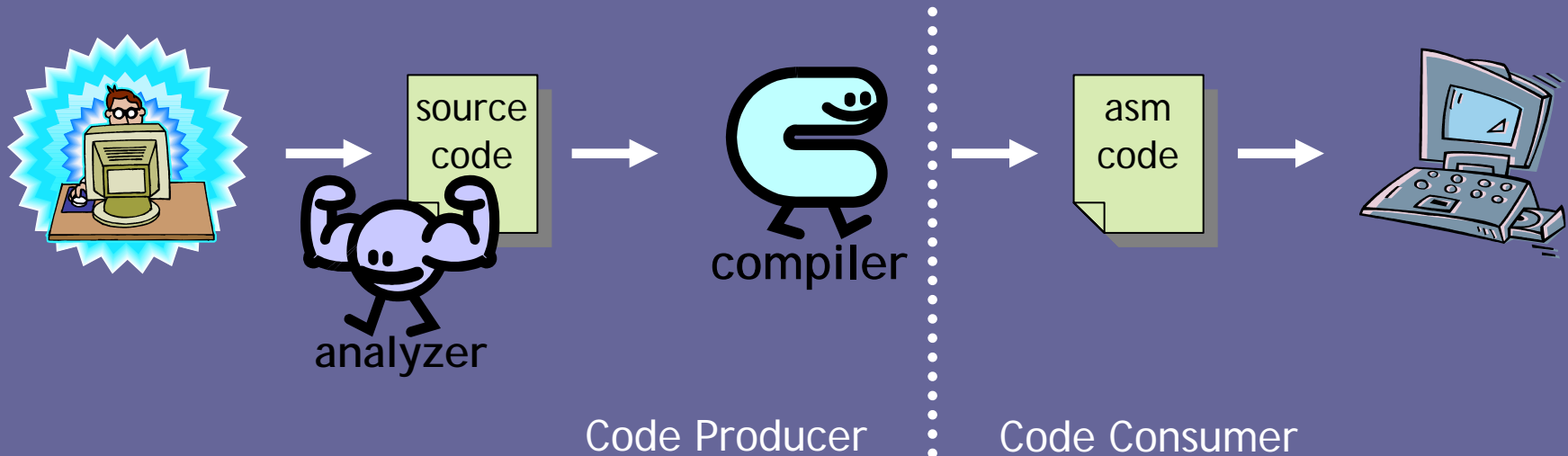
Analysis of Low-Level Code Using Cooperating Decompilers

Bor-Yuh Evan Chang
Matthew Harren
George C. Necula

University of California, Berkeley

SAS 2006

Why Analyze Low-Level Code?



- Analyze what is executed
 - Avoids issues with compiler bugs and underspecified source-language semantics
- Analyze when source is unavailable
 - Applications in mobile code safety assurance

Motivation

Analyzers for low-level code are more difficult and tedious to build

Example: Java Type Analysis

```
class C extends P {  
  void m() {
```

Unsound:

Dependencies must
be carefully

Equal

```
    P p = new P();  
    P c = ... ? new C() : p;  
    ...  
    c.m();
```

```
    rc := m[rsp]
```

```
    if (rc = 0) Lexc
```

```
    r1 := m[r1]
```

```
    r1 :=
```

```
    rsp :=  
    m[rsp] := m[rsp+4]
```

```
    icall [r1]
```

⟨r_c : P, ...⟩

⟨r_c : nonnull P, ...⟩

...⟩

...⟩

Type analysis intermixed
with low-level reasoning

- e.g., args on stack

⟨m[r_{sp}] : nonnull P, ...⟩

Observations

- Handling low-level implementation details is common to many low-level analyzers
 - call stack (provides “local variables”)
 - register allocation
 - dependencies across instructions
- Ad-hoc modularization attempts failed

Goal: Design a modular framework that makes it easy to write high-level analyses

- for different architectures
- for the output of different compilers

Observations

- Intermediate languages abstract varying levels of detail
 - E.g., source language hides compiler details
 - Provides well-specified interface between (de)compiler phases

Proposal: Structure low-level analyses as small, incremental decompilation phases

Basic Idea

```
static void f(C c) { c.m(); }
```

```
f:  
...  
rc := m[rsp+12]  
if (rc = 0) Lexc  
r1 := m[rc]  
r1 := m[r1+28]  
rsp := rsp - 4  
m[rsp] :=  
m[rsp+16]  
icall [r1]  
...
```

```
f(tc):  
rc := tc  
if (rc = 0) Lexc  
r1 := m[rc]  
r1 := m[r1+28]  
t1 := tc  
icall [r1](t1)
```

```
f(c):  
if (c = 0) Lexc  
icall  
[m[m[c]+28]]  
(c)
```

```
f(obj c):  
invokevirtual  
[c, 28]  
()
```

```
f(C c):  
if (c = 0) Lexc  
c.m()
```

Local Variables

Sym Eval

Dynamic Dispatch



Locals



SymEval



OO



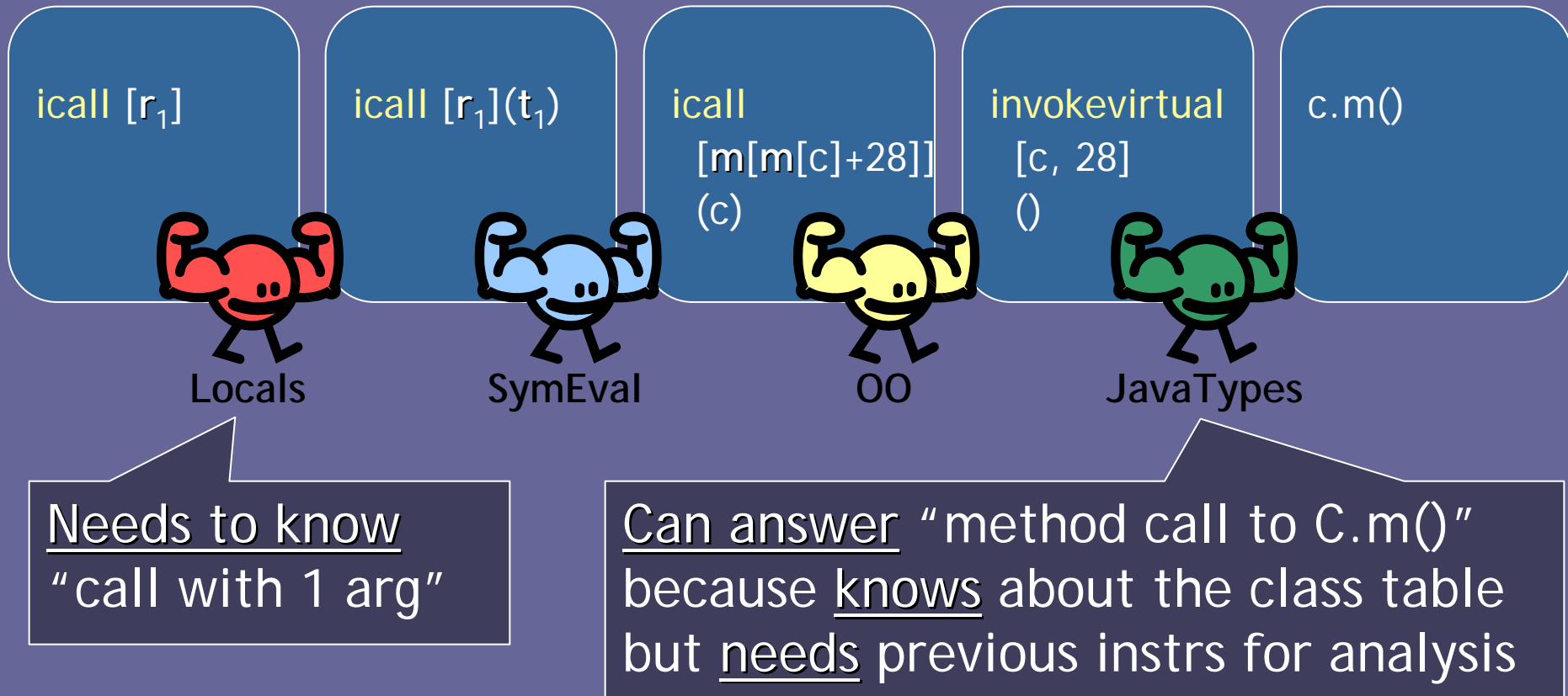
JavaTypes



your analyzer

Difficulties

- Unidirectional communication + analysis is insufficient



Overview of the Framework

- To enable bidirectional communication
 - cannot decompile in stages
 - decompile at all levels simultaneously
- Each decompiler analyzes the preceding instructions before decompiling the next
- ~~“Pipeline”~~ “Reduced product”

Queries

```
static void f(C c) { c.m(); }
```

f:

...

$r_c := m[r_{sp}+12]$

if ($r_c = 0$) L_{exc}

$r_1 := m[r_c]$

$r_1 := m[r_1+28]$

$r_{sp} := r_{sp} - 4$

$m[r_{sp}] :=$

$m[r_{sp}]$

$r_{sp} : sp(-12)$

f(t_c):

$r_c := t_c$

if ($r_c = 0$) L_{exc}

$r_1 := m[r_c]$

$r_1 := m[r_1+28]$

$t_1 := t_c$

$r_1 = m[m[c]+28]$
 $t_1 = c$

f(c):

if ($c = 0$) L_{exc}

$c : \text{nonnull obj}$

f(obj c):

if ($c = 0$) L_{exc}

$c : C$

f(C c):

if ($c = 0$) L_{exc}

icall [r₁]

isFunc(n)

Yes,
1 arg

isMet

Yes,
1 arg

isC?

Yes,
0 args



Locals



SymEval



OO

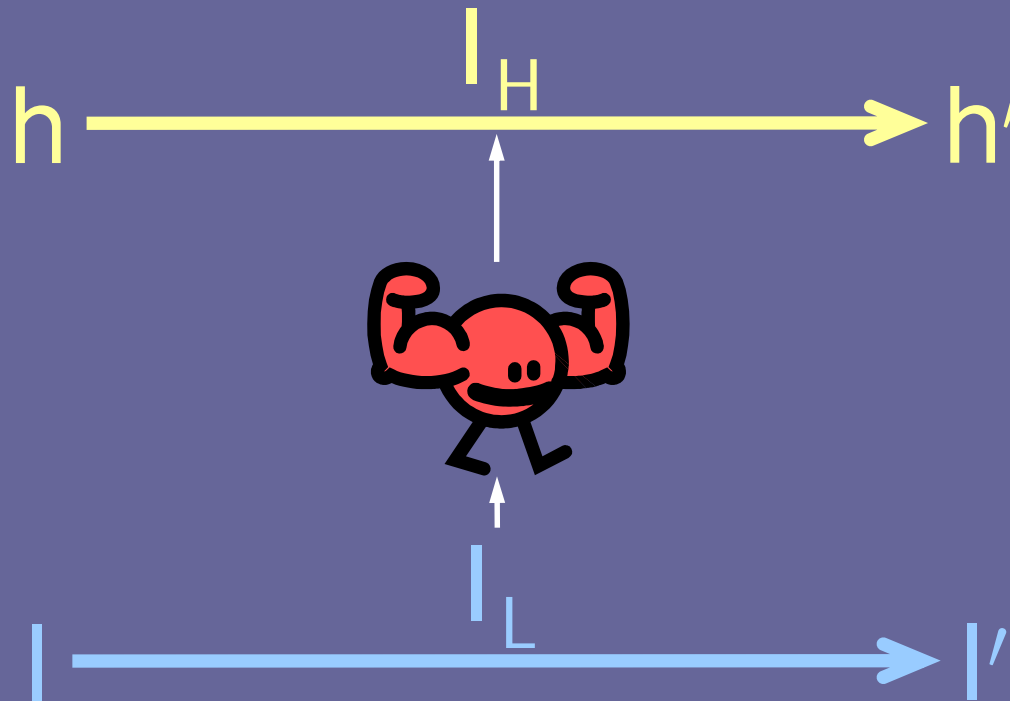


JavaTypes

Communication Summary

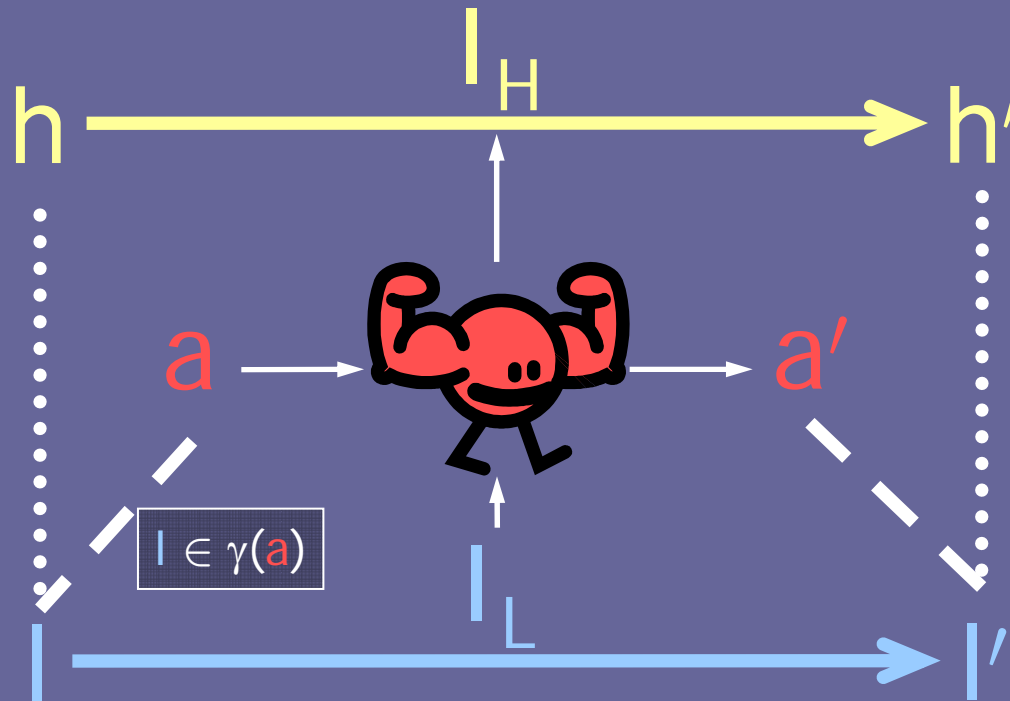
- Low-to-High (Primary)
 - Decompilation Stream
- High-to-Low
 - Queries
 - Initiated by lower-level
 - Questions decompiled
 - Reinterpretations
 - Initiated by higher-level
 - Answers decompiled

Soundness of Decompiler Pipelines



- Operational semantics for each language
- Safety encoded as “not getting stuck”

Soundness of Decompiler Pipelines

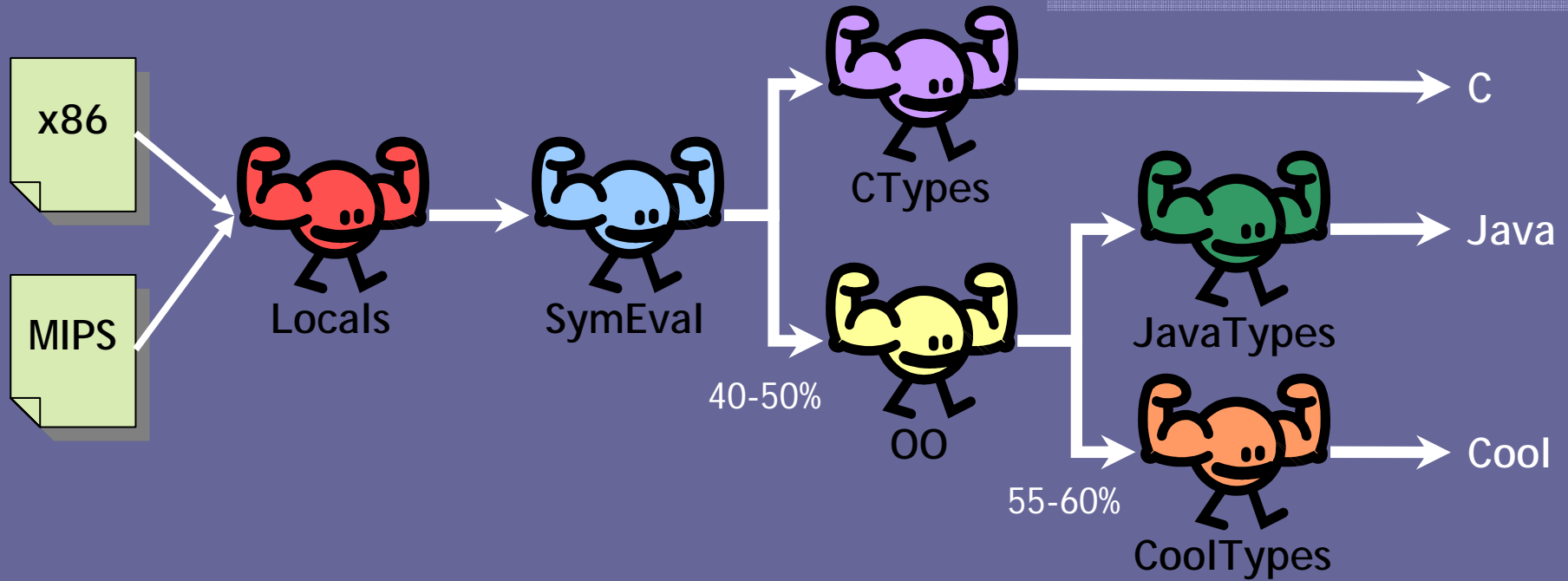


- Safe at high-level implies safe at low-level

Experiments

- Evaluate flexibility
- Evaluate modularity
- Evaluate applicability of existing source-level tools

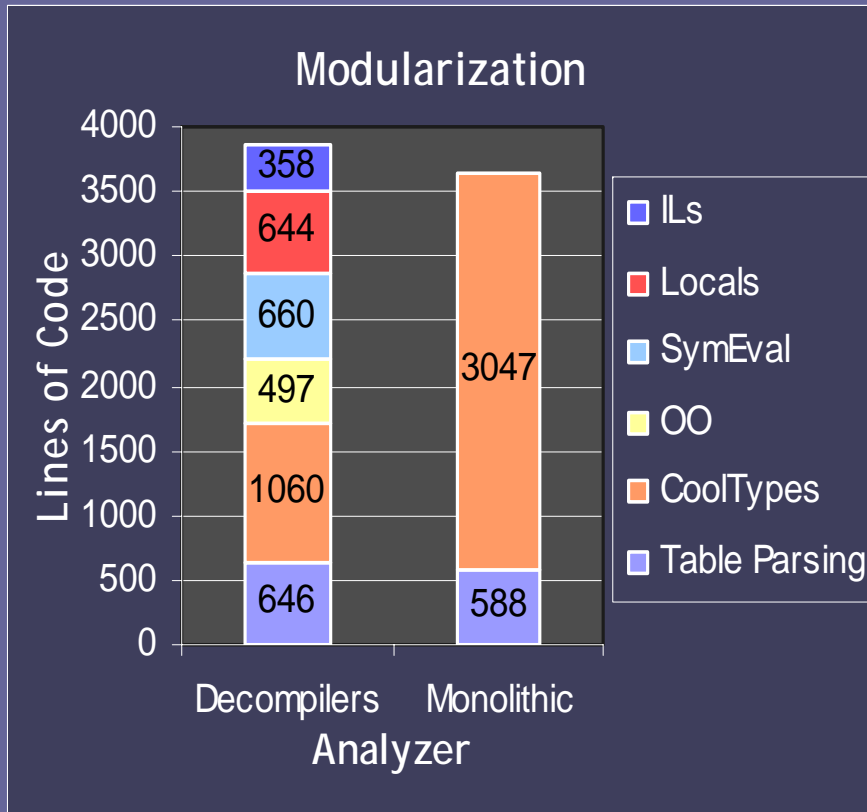
Flexibility



- For the output of gcc, gcj, and compilers for Cool (a “mini-Java”)
- To implement JavaTypes (no exns, interfaces)
 - 3-4 hours, 500 lines

Modularity: Decompilers vs. Monolithic

Type Checking Compiled Cool



- Modules approx. same size

- Compared on 10,339 tests
 - 49 tests, 211 compilers
 - 182 disagreements (pass/fail)
- Decompilers (new)
 - 1 incompletenesses
 - 0 soundness bugs
- Monolithic (used heavily)
 - 5 incompletenesses
 - 3 soundness bugs
 - mishandling of run-time library functions
 - Calls uniform - Locals
 - SSA - SymEval

Applicability of Source-Level Tools

- Experimental Setup

- Compiled 3 previously reported benchmarks for BLAST (B) and Cqual (Q)
- Verified presence (or absence) of bugs as in the original

- Limitations

- Source-level tools needed types
- Recovered types from debugging information
- On optimized code (`qpmouse.c`)
 - gcc -O2 except "merge constants"
 - reads byte in middle of word-sized field

Test Case		Code Size		Decomp	Verification	
		C (kloc)	x86 (kloc)	(s)	Orig (s)	Decomp (s)
<code>qpmouse.c</code>	(B)	8.0	1.9	0.74	0.34	1.26
<code>tlan.c</code>	(B)	10.9	10.7	8.16	41.20	94.30
<code>gamma_dma.c</code>	(Q)	11.2	5.2	2.44	0.97	1.05

Conclusion: Lessons Learned

- Need types for existing source analyses
 - E.g., BLAST on untyped C does not work
- Very useful low-level modules
 - Locals: recover statically-scoped local variables
 - SymEval: recover normalized expr. trees, SSA
- Defining output IL guides analysis impl.
 - IL specifies what analysis should guarantee
 - Leads to small and well-isolated modules

Thank You!

<http://www.cs.berkeley.edu/~bec>

High-to-Low Communication Examples

- Queries
 - “Does e point to a function/method?”
 - “Does e point to an object?”
- Reinterpretations
 - Exceptional successor for try-catch blocks