Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

# Compiling and Optimizing Scripting Languages

Paul Biggar and David Gregg

Department of Computer Science and Statistics
Trinity College Dublin

LLNL, 17th March, 2009

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Motivation

- User needs web page in 0.5 seconds
  - Execution time
  - DB access
  - Network latency
  - Browser rendering

- Easier maintainance

- What if execution was:
  - 2x as fast?
  - 10x as fast?

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Outline

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Outline


1. **Introduction to phc**

2. Current state of phc
   - Challenges to compilation?
   - phc solution: use the C API
   - Speedup

3. Next for phc - Analysis and Optimization
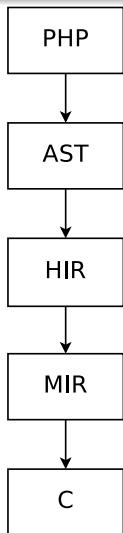   - Simple Optimizations
   - Advanced Optimizations

4. Security

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## **phc**

- http://phpcompiler.org
- Ahead-of-time compiler for PHP
- Edsko de Vries, John Gilbert, Paul Biggar
- BSD license
- Latest release: 0.2.0.3 - compiles non-OO
- svn trunk: compiles most OO

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Structure of **phc**

```
┌─────────┐
│   PHP   │
└─────────┘
     │
     ▼
┌─────────┐
│   AST   │
└─────────┘
     │
     ▼
┌─────────┐
│   HIR   │
└─────────┘
     │
     ▼
┌─────────┐
│   MIR   │
└─────────┘
     │
     ▼
┌─────────┐
│    C    │
└─────────┘
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## PHP

PHP

AST

HIR

MIR

C

```php
<?php

  echo "hello", "world!";

?>
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

# AST

PHP

**AST**

HIR

MIR

C

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## HIR

PHP

AST

**HIR**

MIR

C

```php
<?php
  $x = $a + $b + $c + $d;
?>
```

```php
<?php
  $TLE0 = ($a + $b);
  $TLE1 = ($TLE0 + $c);
  $x = ($TLE1 + $d);
?>
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## MIR

PHP

AST

HIR

**MIR**

C

```php
<?php
  while ($cond)
    echo "hello", "world!";
?>

<?php
L7:
  $TLE0 = !$cond;
  if ($TLE0) goto L3 else goto L6;
L6:
  print('hello');
  print('world!');
  goto L7;
L3:
?>
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Plugins



*http://phpcompiler.org/doc/latest/devmanual.html*

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

# XML



```xml
<?xml version="1.0"?>
<AST:PHP_script xmlns:AST="http://www.phpcompiler.org/phc-1.1">
  <AST:Statement_list>
    <AST:Eval_expr>
      <AST:Method_invocation>
        <AST:Target xsi:nil="true" />
        <AST:METHOD_NAME>
          <value>echo</value>
        </AST:METHOD_NAME>
        <AST:Actual_parameter_list>
          <AST:Actual_parameter>
            <bool><!-- is_ref -->false</bool>
            <AST:STRING>
              <value>hello</value>
            </AST:STRING>
          </AST:Actual_parameter>
          <AST:Actual_parameter>
            <bool><!-- is_ref -->false</bool>
            <AST:STRING>
              <value>world!</value>
            </AST:STRING>
          </AST:Actual_parameter>
        </AST:Actual_parameter_list>
      </AST:Method_invocation>
    </AST:Eval_expr>
    <AST:Nop>
    </AST:Nop>
  </AST:Statement_list>
</AST:PHP script>
```
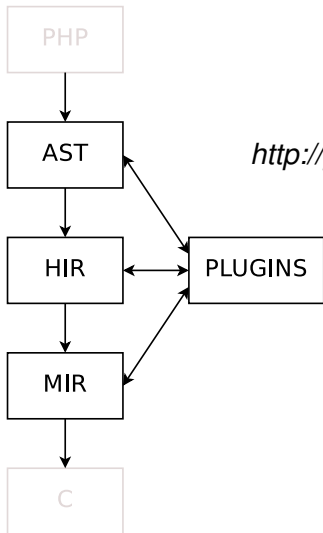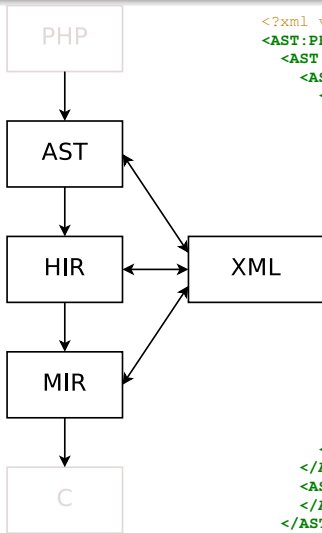
Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Outline

**1** Introduction to phc

**2** Current state of phc
  - Challenges to compilation?
  - phc solution: use the C API
  - Speedup

**3** Next for phc - Analysis and Optimization
  - Simple Optimizations
  - Advanced Optimizations

**4** Security

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## SAC 2009

## A Practical Solution for Scripting Language Compilers

Paul Biggar, Edsko de Vries and David Gregg

Department of Computer Science and Statistics
Trinity College Dublin

ACM Symposium on Applied Computing - PL track
12th March, 2009

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Sneak peak

- Problem: Scripting languages present "unique" problems (in practice)

- Solution: Re-use as much of the ~~Canonical~~ *Reference Implementation* as possible.

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Outline

1. Introduction to phc

2. **Current state of phc**
   - Challenges to compilation?
   - phc solution: use the C API
   - Speedup

3. Next for phc - Analysis and Optimization
   - Simple Optimizations
   - Advanced Optimizations

4. Security

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Undefined

*The PHP group claim that they have the final say in the specification of PHP. This group's specification is an implementation, and there is no prose specification or agreed validation suite. There are alternate implementations [...] that claim to be compatible (they don't say what this means) with some version of PHP.*

D. M. Jones. Forms of language specification: Examples from commonly used computer languages. ISO/IEC JTC1/SC22/OWG/N0121, February 2008.

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Batteries included

```
abs()                               apc_load_constants()                array_intersect()                array_values()
acos()                              apc_sma_info()                      array_intersect_assoc()          array_walk()
acosh()                             apc_store()                         array_intersect_key()            array_walk_recursive()
addcslashes()                       apd_breakpoint()                    array_intersect_uassoc()         ArrayIterator::current()
addslashes()                        apd_callstack()                     array_intersect_ukey()           ArrayIterator::key()
aggregate()                         apd_clunk()                         array_key_exists()               ArrayIterator::next()
aggregate_info()                    apd_continue()                      array_keys()                     ArrayIterator::rewind()
aggregate_methods()                 apd_croak()                         array_map()                      ArrayIterator::seek()
aggregate_methods_by_list()         apd_dump_function_table()           array_merge()                    ArrayIterator::valid()
aggregate_methods_by_regexp()       apd_dump_persistent_resources()     array_merge_recursive()          ArrayObject::__construct()
aggregate_properties()              apd_dump_regular_resources()        array_multisort()                ArrayObject::append()
aggregate_properties_by_list()      apd_echo()                          array_pad()                      ArrayObject::count()
aggregate_properties_by_regexp()    apd_get_active_symbols()            array_pop()                      ArrayObject::getIterator()
aggregation_info()                  apd_set_pprof_trace()               array_product()                  ArrayObject::offsetExists()
apache_child_terminate()            apd_set_session()                   array_push()                     ArrayObject::offsetGet()
apache_get_modules()                apd_set_session_trace()             array_rand()                     ArrayObject::offsetSet()
apache_get_version()                apd_set_socket_session_trace()      array_reduce()                   ArrayObject::offsetUnset()
apache_getenv()                     array()                             array_reverse()                  arsort()
apache_lookup_uri()                 array_change_key_case()             array_search()                   ascii2ebcdic()
apache_note()                       array_chunk()                       array_shift()                    asin()
apache_request_headers()            array_combine()                     array_slice()                    asinh()
apache_reset_timeout()              array_count_values()                array_splice()                   asort()
apache_response_headers()           array_diff()                        array_sum()                      aspell_check()
apache_setenv()                     array_diff_assoc()                  array_udiff()                    aspell_check_raw()
apc_add()                           array_diff_key()                    array_udiff_assoc()              aspell_new()
apc_cache_info()                    array_diff_uassoc()                 array_udiff_uassoc()             aspell_suggest()
apc_clear_cache()                   array_diff_ukey()                   array_uintersect()               assert()
apc_compile_file()                  array_fill()                        array_uintersect_assoc()         assert_options()
apc_define_constants()              array_fill_keys()                   array_uintersect_uassoc()        atan()
apc_delete()                        array_filter()                      array_unique()                   atan2()
apc_fetch()                         array_flip()                        array_unshift()                  atanh()
```

Jeff Atwood, Coding Horror, May 20th, 2008
*http://www.codinghorror.com/blog/archives/001119.html*

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Change between releases

```php
<?php
  var_dump (0x9fa0ff0b);
?>
```

### PHP 5.2.1 (32-bit)

int(2147483647)

### PHP 5.2.3 (32-bit)

float(2678128395)

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Run-time code generation

```php
<?php
  eval ($argv[1]);
?>



<?php
  include ("mylib.php");
  ...
  include ("plugin.php");
  ...
?>
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Outline

1. Introduction to phc

2. **Current state of phc**
   - Challenges to compilation?
   - **phc solution: use the C API**
   - Speedup

3. Next for phc - Analysis and Optimization
   - Simple Optimizations
   - Advanced Optimizations

4. Security

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Use C API

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## More detail

| PHP | zval |
| --- | --- |
| Python | PyObject |
| Ruby | VALUE |
| Lua | TValue |

H. Muhammad and R. Ierusalimschy. C APIs in extension and extensible languages. Journal of Universal Computer Science, 13(6):839–853, 2007.

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Simple listings: $i = 0

```
// $i = 0;
{
  zval* p_i;
  php_hash_find (LOCAL_ST, "i", 5863374, p_i);
  php_destruct (p_i);
  php_allocate (p_i);
  ZVAL_LONG (*p_i, 0);
}
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Example: $i = 0

```
// $i = 0;
{
  if (local_i == NULL)
  {
    local_i = EG (uninitialized_zval_ptr);
    local_i->refcount++;
  }
  zval **p_lhs = &local_i;

  zval *value;
  if ((*p_lhs)->is_ref)
  {
    // Always overwrite the current value
    value = *p_lhs;
    zval_dtor (value);
  }
  else
  {
    ALLOC_INIT_ZVAL (value);
    zval_ptr_dtor (p_lhs);
    *p_lhs = value;
  }

  ZVAL_LONG (value, 0);
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Example: $i = $j

```
// $i = $j;
{
  if (local_i == NULL)
  {
    local_i = EG (uninitialized_zval_ptr);
    local_i->refcount++;
  }
  zval **p_lhs = &local_i;

  zval *rhs;
  if (local_j == NULL)
    rhs = EG (uninitialized_zval_ptr);
  else
    rhs = local_j;

  if (*p_lhs != rhs)
  {
    if ((*p_lhs)->is_ref)
    {
      // First, call the destructor to remove any data structures
      // associated with lhs that will now be overwritten
      zval_dtor (*p_lhs);
      // Overwrite LHS
      (*p_lhs)->value = rhs->value;
      (*p_lhs)->type = rhs->type;
      zval_copy_ctor (*p_lhs);
    }
    else
    {
      zval_ptr_dtor (p_lhs);
      if (rhs->is_ref)
      {
        // Take a copy of RHS for LHS
        *p_lhs = zvp_clone_ex (rhs);
      }
      else
      {
        // Share a copy
        rhs->refcount++;
        *p_lhs = rhs;
      }
    }
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Example: printf ($f)

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Applicability

- Everything
    - Perl
    - PHP
    - Ruby
    - Tcl – *I think*

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Applicability

- Everything
  - Perl
  - PHP
  - Ruby
  - Tcl – *I think*

- Except specification
  - Lua
  - Python

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Applicability

- Everything
  - Perl
  - PHP
  - Ruby
  - Tcl – *I think*

- Except specification
  - Lua
  - Python

- Not at all
  - Javascript

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Outline

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Original Speed-up

# 0.1x
(10 times slower than the PHP interpreter)

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## The problem with copies

```php
<?php
  for ($i = 0; $i < $n; $i++)
    $str = $str . "hello";
?>



<?php
  for ($i = 0; $i < $n; $i++)
  {
    $T = $str . "hello";
    $str = $T;
  }
?>
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Optimization

- Constant folding

```php
<?php
   ...
   $T = "5" + true;
   ...
?>
```

```php
<?php
   ...
   $T = 6;
   ...
?>
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Optimization

- Constant folding
- Constant pooling

```php
<?php
  $sum = 0;
  for ($i = 0; $i < 10; $i=$i+1)
  {
    $sum .= "hello";
  }
?>
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

# Optimization

- Constant folding
- Constant pooling
- Function caching

```c
// printf ($f);
static php_fcall_info printf_info;
{
  php_fcall_info_init ("printf", &printf_info);

  php_hash_find (
    LOCAL_ST, "f", 5863275, &printf_info.params);

  php_call_function (&printf_info);
}
```

Introduction to phc
**Current state of phc**
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Optimization

- Constant folding
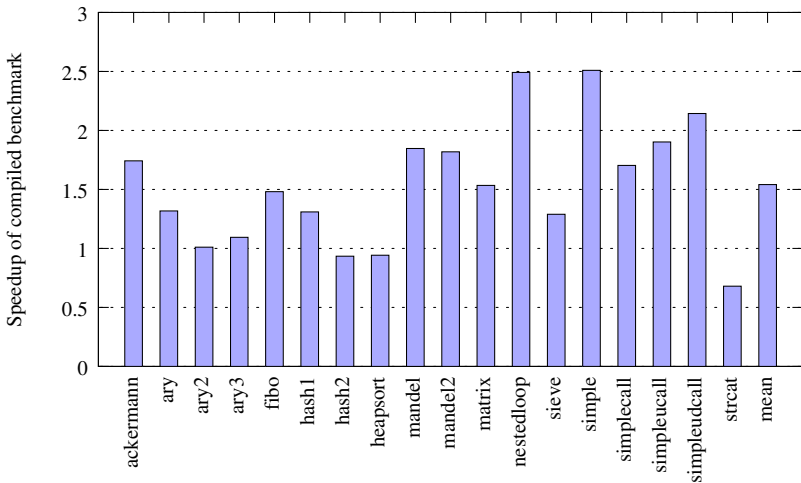- Constant pooling
- Function caching
- Pre-hashing

```
// $i = 0;
{
  zval* p_i;
  php_hash_find (LOCAL_ST, "i", 5863374, p_i);
  php_destruct (p_i);
  php_allocate (p_i);
  ZVAL_LONG (*p_i, 0);
}
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Optimization

- Constant folding
- Constant pooling
- Function caching
- Pre-hashing
- Symbol-table removal

```
// $i = 0;
{
  php_destruct (local_i);
  php_allocate (local_i);
  ZVAL_LONG (*local_i, 0);
}
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Challenges to compilation?
phc solution: use the C API
Speedup

## Current speed-up

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Outline

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Outline

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Intra-procedural optimizations



- Dead-code elimination
- Sparse-conditional **constant propagation**

Introduction to phc
Current state of phc
**Next for phc - Analysis and Optimization**
Security

Simple Optimizations
Advanced Optimizations

## Type-inference

```php
<?php

  function a ($x, $y)
  {
    $str = $x . $y;
    ...
    return $str;
  }

?>
```

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## User-space handlers

- __toString
- __get
- __set
- __isset
- __unset
- __sleep
- __wake
- __call
- __callStatic
- ...

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

# C API handlers

- read_property
- read_dimension
- get
- set
- cast_object
- has_property
- unset_property
- ...

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Unknown types propagate

- local symbol table
- global symbol table
- return values
- reference parameters
- callee parameters

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Outline

1. Introduction to phc

2. Current state of phc
   - Challenges to compilation?
   - phc solution: use the C API
   - Speedup

3. Next for phc - Analysis and Optimization
   - Simple Optimizations
   - Advanced Optimizations

4. Security

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Analysis design

- Must model types precisely
  - (Possibly unnamed) fields, arrays, variables and method calls

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Analysis design

- Must model types precisely
  - (Possibly unnamed) fields, arrays, variables and method calls

- Uses and definitions incomplete
  - Can't use *def-use chains*
  - Can't use *SSA*

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Analysis design

- Must model types precisely
    - (Possibly unnamed) fields, arrays, variables and method calls

- Uses and definitions incomplete
    - Can't use *def-use chains*
    - Can't use *SSA*

- Imprecise callgraph

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Algorithm

- Abstract Execution / Interpretation

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Algorithm

- Abstract Execution / Interpretation

- Points-to analysis
  - *-sensitive

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Algorithm

- Abstract Execution / Interpretation

- Points-to analysis
  - *-sensitive

- Constant-propagation
  - Precision
  - Array-indices/field names
  - Implicit conversions

A. Pioli. Conditional pointer aliasing and constant propagation. Master's thesis, SUNY at New Paltz, 1999.

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Algorithm

- Abstract Execution / Interpretation

- Points-to analysis
    - *-sensitive

- Constant-propagation
    - Precision
    - Array-indices/field names
    - Implicit conversions

- Type-inference
    - Virtual calls
    - Function annotations

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Complex cases

- Hashtables
- Implicit conversions
- Variable-variables
- $GLOBALS
- Static includes
- $SESSION
- Compiler temporaries

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Interesting thoughts

- Strip off first loop iteration

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Interesting thoughts

- Strip off first loop iteration
- JITs or Gal/Franz Tracing?

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

Simple Optimizations
Advanced Optimizations

## Interesting thoughts

- Strip off first loop iteration
- JITs or Gal/Franz Tracing?
- Use string transducer analysis

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

# Outline

1. Introduction to phc

2. Current state of phc
   - Challenges to compilation?
   - phc solution: use the C API
   - Speedup

3. Next for phc - Analysis and Optimization
   - Simple Optimizations
   - Advanced Optimizations

4. Security

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

# Security

- Davis - if we include it, we'll do better

Sound and Precise Analysis of Web Applications
for Injection Vulnerabilities
Gary Wassermann, Zhendong Su, PLDI'07.

Static approximation of dynamically generated Web pages
Yasuhiko Minamide, WWW 2005

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Security

- Davis - if we include it, we'll do better
- Tuwien/Pixy - taint analysis (literal analysis + points to)

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Security

- Davis - if we include it, we'll do better
- Tuwien/Pixy - taint analysis (literal analysis + points to)
- Utrecht/Stanford - dont remember

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Summary

- Re-use existing run-time for language
- Better yet: standardize libraries (and language?), including FFI

- Analysis needs to be precise, and whole-program
- Pessimistic assumptions spread

- Language, implementation and community need to be fixed
  - All related?

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

### Thanks

# **phc** needs contributors

- contribute:
  http://phpcompiler.org/contribute.html
- mailing list: phc-general@phpcompiler.org

- slides: http://www.cs.tcd.ie/~pbiggar/
- contact: paul.biggar@gmail.com

Introduction to phc
Current state of phc
Next for phc - Analysis and Optimization
Security

## Complex cases

- Hashtables
- Implicit conversions
- Variable-variables
- $GLOBALS
- Static includes
- $SESSION
- Compiler temporaries