
Faster CPython Documentation

Release 0.0

Victor Stinner

January 29, 2016

1	FAT Python	3
2	Everything in Python is mutable	9
3	Optimizations	13
4	Python bytecode	19
5	Python C API	21
6	AST Optimizers	23
7	Old AST Optimizer	25
8	Register-based Virtual Machine for Python	33
9	Read-only Python	39
10	History of Python optimizations	43
11	Misc	45
12	Kill the GIL?	51
13	Implementations of Python	53
14	Benchmarks	55
15	PEP 509: Add a private version to dict	57
16	PEP 510: Specialized functions with guards	59
17	PEP 511: API for AST transformers	61
18	Random notes about PyPy	63
19	Talks	65
20	Links	67

Contents:

FAT Python



1.1 Intro

The FAT Python project was started by Victor Stinner in October 2015 to try to solve issues of previous attempts of “static optimizers” for Python. The main feature are efficient guards using versionned dictionaries to check if something was modified. Guards are used to decide if the specialized bytecode of a function can be used or not.

Python FAT is expected to be FAT... maybe FAST if we are lucky. FAT because it will use two versions of some functions where one version is specialised to specific argument types, a specific environment, optimized when builtins are not mocked, etc.

See the [fatoptimizer documentation](#) which is the main part of FAT Python.

The FAT Python project is made of multiple parts:

- The [fatoptimizer project](#) is the static optimizer for Python 3.6 using function specialization with guards. It is implemented as an AST optimizer.
- The [fat module](#) is a Python extension module (written in C) implementing fast guards. The [fatoptimizer](#) optimizer uses [fat](#) guards to specialize functions. [fat](#) guards are used to verify assumptions used to specialize the code. If an assumption is no more true, the specialized code is not used. The [fat](#) module is required to run code optimized by [fatoptimizer](#) if at least one function is specialized.
- Python Enhancement Proposals (PEP):
 - PEP 509: [Add a private version to dict](#)
 - PEP 510: [Specialized functions with guards](#)
 - PEP 511: [API for AST transformers](#)
- Patches for Python 3.6:
 - PEP 509: [Add ma_version to PyDictObject](#)
 - PEP 510: [Specialize functions with guards](#)
 - PEP 511: [Add sys.set_code_transformers\(\)](#)
 - Related to the PEP 511:
 - * *DONE*: PEP 511: [Add test.support.optim_args_from_interpreter_flags\(\)](#)
 - * *DONE*: PEP 511: [code.co_lnotab](#): use signed line number delta to support moving instructions in an optimizer
 - * *DONE*: PEP 511: [Add ast.Constant](#) to allow AST optimizer to emit constants
 - * *DONE*: [Lib/test/test_compileall.py](#) fails when run directly
 - * *DONE*: [site](#) ignores [ImportError](#) when running [sitecustomize](#) and [usercustomize](#)
 - * *DONE*: [code_richcompare\(\)](#) don't use constant type when comparing code constants

Announcements and status reports:

- [Status of the FAT Python project](#), January 12, 2016
- ‘FAT’ and fast: [What’s next for Python: Article of InfoWorld](#) by Serdar Yegulalp (January 11, 2016)
- [Python-Dev] [Third milestone of FAT Python](#)
- [Status of the FAT Python project](#), November 26, 2015
- [python-dev] [Second milestone of FAT Python](#) (Nov 2015)
- [python-ideas] [Add specialized bytecode with guards to functions](#) (Oct 2015)

1.2 Getting started

Compile Python 3.6 patched with PEP 509, PEP 510 and PEP 511:

```
hg clone http://hg.python.org/sandbox/fatpython
cd fatpython
./configure --with-pydebug CFLAGS="-O0" && make
```

Install fat:


```
git clone https://github.com/haypo/fat
cd fat
../python setup.py build
cp -v build/lib*/fat.*so ../Lib
cd ..
```

Install fatoptimizer:

```
git clone https://github.com/haypo/fatoptimizer
(cd Lib; ln -s ../fatoptimizer/fatoptimizer .)
```

fatoptimizer is registered by the site module if `-X fat` command line option is used. Extract of `Lib/site.py`:

```
if 'fat' in sys._xoptions:
    import fatoptimizer
    fatoptimizer._register()
```

Check that fatoptimizer is registered with:

```
$ ./python -X fat -c 'import sys; print(sys.implementation.optim_tag)'
fat-opt
```

You must get `fat-opt` (and not `opt`).

1.3 How can you contribute?

The [fatoptimizer project](#) needs the most love. Currently, the optimizer is not really smart. There is a long [TODO list](#). Pick a simple optimization, try to implement it, send a pull request on GitHub. At least, any kind of feedback is useful ;-)

If you know the C API of Python, you may also review the implementation of the PEPs:

- [PEP 509](#): Add `ma_version` to `PyDictObject`
- [PEP 510](#): Specialize functions with guards
- [PEP 511](#): Add `sys.set_code_transformers()`

But these PEPs are still work-in-progress, so the implementation can still change.

1.4 Play with FAT Python

See [Getting started](#) to compile FAT Python.

1.4.1 Disable peephole optimizer

The `-o noopt` command line option disables the Python peephole optimizer:

```
$ ./python -o noopt -c 'import dis; dis.dis(compile("1+1", "test", "exec"))'
1          0 LOAD_CONST           0 (1)
          3 LOAD_CONST           0 (1)
          6 BINARY_ADD
          7 POP_TOP
          8 LOAD_CONST           1 (None)
         11 RETURN_VALUE
```

1.4.2 Specialized code calling builtin function

Test fatoptimizer on builtin function:

```
$ ./python -X fat
>>> def func(): return len("abc")
...

>>> import dis
>>> dis.dis(func)
 1           0 LOAD_GLOBAL             0 (len)
           3 LOAD_CONST               1 ('abc')
           6 CALL_FUNCTION             1 (1 positional, 0 keyword pair)
           9 RETURN_VALUE

>>> import fat
>>> fat.get_specialized(func)
[(<code object func at 0x7f9d3155b1e0, file "<stdin>", line 1>,
[<fat.GuardBuiltins object at 0x7f9d39191198>])]

>>> dis.dis(fat.get_specialized(func)[0][0])
 1           0 LOAD_CONST               1 (3)
           3 RETURN_VALUE
```

The specialized code is removed when the function is called if the builtin function is replaced (here by declaring a `len()` function in the global namespace):

```
>>> len=lambda obj: "mock"
>>> func()
'mock'
>>> fat.func_get_specialized(func)
[]
```

1.4.3 Microbenchmark

Run a microbenchmark on specialized code:

```
$ ./python -m timeit -s 'def f(): return len("abc")' 'f()'
10000000 loops, best of 3: 0.122 usec per loop

$ ./python -X fat -m timeit -s 'def f(): return len("abc")' 'f()'
10000000 loops, best of 3: 0.0932 usec per loop
```

Python must be optimized to run a benchmark: use `./configure && make clean && make` if you previously compiled it in debug mode.

You should compare specialized code to an unpatched Python 3.6 to run a fair benchmark (to also measure the overhead of PEP 509, 510 and 511 patches).

1.5 Run optimized code without registering fatoptimizer

You have to compile optimized .pyc files:

```
# the optimizer is slow, so add -v to enable fatoptimizer logs for more fun
./python -X fat -v -m compileall

# why does compileall not compile encodings/*.py?
./python -X fat -m py_compile Lib/encodings/{__init__,aliases,latin_1,utf_8}.py
```

Finally, enjoy optimized code with no registered optimized:

```
$ ./python -o fat-opt -c 'import sys; print(sys.implementation.optim_tag, sys.get_code_transformers())'
fat-opt []
```

Remember that you cannot import .py files in this case, only .pyc:

```
$ echo 'print("Hello World!")' > hello.py
$ ENV/bin/python -o fat-opt -c 'import hello'
Traceback (most recent call last):
  File "<string>", line 1, in <module>
ImportError: missing AST transformers for 'hello.py': optim_tag='fat-opt', transformers tag='noopt'
```

1.6 Origins of FAT Python

- *Old AST optimizer project*
- *read-only Python*
- Dave Malcolm wrote a patch modifying Python/eval.c to support specialized functions. See the <http://bugs.python.org/issue10399>

Everything in Python is mutable

2.1 Problem

Python Python unittest.mock

Developers like Python because it's possible to modify (almost) everything. This feature is heavily used in unit tests with `unittest.mock` which can override builtin function, override class methods, modify “constants, etc.

Most optimization rely on assumptions. For example, inlining rely on the fact that the inlined function is not modified. Implement optimization in respect of the Python semantics require to implement various assumptions.

2.1.1 Builtin functions

Python Python .

Python provides a lot of builtins functions. All Python applications rely on them, and usually don't expect that these functions are overridden. In practice, it is very easy to override them.

`len()` :

Example overridden the builtin `len()` function:

```
import builtins

def func(obj):
    print("length: %s" % len(obj))

func("abc")
builtins.len = lambda obj: "mock!"
func("abc")
```

Output:

```
length: 3
length: mock!
```

```
len() func LEGB` ``globals len, builtins
```

Technically, the `len()` function is loaded in `func()` with the `LOAD_GLOBAL` instruction which first tries to lookup in frame globals namespace, and then lookup in the frame builtins namespace.

global len builtin:

Example overriding the `len()` builtin function with a `len()` function injected in the global namespace:

```
def func(obj):  
    print("length: %s" % len(obj))  
  
func("abc")  
len = lambda obj: "mock!"  
func("abc")
```

Output:

```
length: 3  
length: mock!
```

Builtins are references in multiple places:

Builtins:

- the `builtins` module
- frames have a `f_builtins` attribute (builtins dictionary)
- the global `PyInterpreterState` structure has a `builtins` attribute (builtins dictionary)
- frame globals have a `__builtins__` variable (builtins dictionary, or builtins module when `__name__` equals `__main__`)

2.1.2 Function code

:

It is possible to modify at runtime the bytecode of a function to modify completely its behaviour. Example:

```
def func(x, y):  
    return x + y  
  
print("1+2 = %s" % func(1, 2))  
  
def mock(x, y):  
    return 'mock'  
  
func.__code__ = mock.__code__  
print("1+2 = %s" % func(1, 2))
```

Output:

```
1+2 = 3  
1+2 = mock
```

2.1.3 Local variables

hack local variable

Technically, it is possible to modify local variable of a function outside the function.

: `hack()` local variable, `x`

Example of a function `hack()` which modifies the `x` local variable of its caller:

```

import sys
import ctypes

def hack():
    # Get the frame object of the caller
    frame = sys._getframe(1)
    frame.f_locals['x'] = "hack!"
    # Force an update of locals array from locals dict
    ctypes.pythonapi.PyFrame_LocalsToFast(ctypes.py_object(frame),
                                           ctypes.c_int(0))

def func():
    x = 1
    hack()
    print(x)

func()

```

Output:

```
hack!
```

2.1.4 Modification made from other modules

Python A B

A Python module A can be modified by a Python module B.

2.1.5 Multithreading

Python B A local variables A B

When two Python threads are running, the thread B can modify shared resources of thread A, or even resources which are supposed to only be access by the thread A like local variables.

B local variables

The thread B can modify function code, override builtin functions, modify local variables, etc.

2.1.6 Python Imports and Python Modules

The Python import path `sys.path` is initialized by multiple environment variables (ex: `PYTHONPATH` and `PYTHONHOME`), modified by the `site` module and can be modified anytime at runtime (by modifying `sys.path` directly).

Moreover, it is possible to modify `sys.modules` which is the “cache” between a module fully qualified name and the module object. For example, `sys.modules['sys']` should be `sys`. It is possible to remove modules from `sys.modules` to force to reload a module. It is possible to replace a module in `sys.modules`.

The eventlet modules injects monkey-patched modules in `sys.modules` to convert I/O blocking operations to asynchronous operations using an event loop.

2.2 Solutions

2.2.1 Make strong assumptions, ignore changes

If the optimizer is an opt-in options, users are aware that the optimizer can make some compromises on the Python semantics to implement more aggressive optimizations.

2.2.2 Static analysis

Analyze the code to ensure that functions don't mutate everything, for example ensure that a function is pure.

Dummy example:

```
def func(x, y):  
    return x + y
```

This function `func()` is pure: it has no side effect. This function will not override builtins, not modify local variables of the caller, etc. It is safe to call this function from anywhere.

It is possible to analyze the code to check that an optimization can be enabled.

2.2.3 Use guards checked at runtime

For some optimizations, a static analysis cannot ensure that all assumptions required by an optimization will respected. Adding guards allows to check assumptions during the execution to use the optimized code or fallback to the original code.

Optimizations

See also [fatoptimizer optimizations](#).

3.1 Inline function calls

Example:

```
def _get_sep(path):
    if isinstance(path, bytes):
        return b'/'
    else:
        return '/'

def isabs(s):
    """Test whether a path is absolute"""
    sep = _get_sep(s)
    return s.startswith(sep)
```

Inline `_get_sep()` into `isabs()` and simplify the code for the `str` type:

```
def isabs(s: str):
    return s.startswith('/')
```

It can be implemented as a simple call to the C function `PyUnicode_Tailmatch()`.

Note: Inlining uses more memory and disk because the original function should be kept. Except if the inlined function is unreachable (ex: “private function?”).

Links:

- [Issue #10399](#): AST Optimization: inlining of function calls

3.2 CALL_METHOD

See [issue #26110](#): Speedup method calls 1.2x

3.3 Move invariants out of the loop

Example:

```
def func(obj, lines):
    for text in lines:
        print(obj.cleanup(text))
```

Become:

```
def func(obj, lines):
    local_print = print
    obj_cleanup = obj.cleanup
    for text in lines:
        local_print(obj_cleanup(text))
```

Local variables are faster than global variables and the attribute lookup is only done once.

3.4 C functions using only C types

Optimizations:

- Avoid reference counting
- Memory allocations on the heap
- Release the GIL

Example:

```
def demo():
    s = 0
    for i in range(10):
        s += i
    return s
```

In specialized code, it may be possible to use basic C types like `char` or `int` instead of Python codes which can be allocated on the stack, instead of allocating objects on the heap. `i` and `s` variables are integers in the range `[0; 45]` and so a simple C type `int` (or even `char`) can be used:

```
PyObject *demo(void)
{
    int s, i;
    Py_BEGIN_ALLOW_THREADS
    s = 0;
    for(i=0; i<10; i++)
        s += i;
    Py_END_ALLOW_THREADS
    return PyLong_FromLong(s);
}
```

Note: if the function is slow, we may need to check sometimes if a signal was received.

3.5 Release the GIL

Many methods of builtin types don't need the *GIL*. Example: `"abc".startswith("def")`.

3.6 Replace calls to pure functions with the result

Examples:

- `len('abc')` becomes `3`
- `"python2.7".startswith("python")` becomes `True`
- `math.log(32) / math.log(2)` becomes `5.0`

Can be implemented in the AST optimizer.

3.7 Constant propagation

Propagate constant values of variables. Example:

Original	Constant propagation
<pre>def func() x = 1 y = x return y</pre>	<pre>def func() x = 1 y = 1 return 1</pre>

Implemented in fatoptimizer.

Read also the [Wikipedia article on copy propagation](#).

3.8 Constant folding

Compute simple operations at the compilation. Usually, at least arithmetic operations (`a+b`, `a-b`, `a*b`, etc.) are computed. Example:

Original	Constant folding
<pre>def func() return 1 + 1</pre>	<pre>def func() return 2</pre>

Implemented in fatoptimizer and the *CPython peephole optimizer*.

See also

- [issue #1346238](#): A constant folding optimization pass for the AST
- [Wikipedia article on constant folding](#).

3.9 Peephole optimizer

See *CPython peephole optimizer*.

3.10 Loop unrolling

Example:

```
for i in range(4):  
    print(i)
```

The loop body can be duplicated (twice in this example) to reduce the cost of a loop:

```
for i in range(0,4,2):  
    print(i)  
    print(i+1)  
i = 3
```

Or the loop can be removed by duplicating the body for all loop iterations:

```
i=0  
print(i)  
i=1  
print(i)  
i=2  
print(i)  
i=3  
print(i)
```

Combined with other optimizations, the code can be simplified to:

```
print('0')  
print('1')  
print('2')  
i = 3  
print('3')
```

Implemented in fatoptimizer

Read also the [Wikipedia article on loop unrolling](#).

3.11 Dead code elimination

- Replace `if 0:` code with `pass`
- `if DEBUG: print("debug")` where `DEBUG` is known to be `False`

Implemented in fatoptimizer and the *CPython peephole optimizer*.

See also [Wikipedia Dead code elimination article](#).

3.12 Load globals and builtins when the module is loaded

Load globals when the module is loaded? Ex: load “print” name when the module is loaded.

Example:

```
def hello():  
    print("Hello World")
```

Become:

```
local_print = print

def hello():
    local_print("Hello World")
```

Useful if `hello()` is compiled to C code.

`fatoptimizer` implements a “copy builtins to constants optimization” optimization.

3.13 Don't create Python frames

Inlining and other optimizations don't create Python frames anymore. It can be a serious issue to debug programs: tracebacks are an important feature of Python.

At least in debug mode, frames should be created.

PyPy supports lazy creation of frames if an exception is raised.

Python bytecode

4.1 CPython peephole optimizer

Implementation: `Python/peephole.c`

Optmizations:

- *Constant folding*
- *Dead code elimination*
- Some other optimizations more specific to the bytecode, like removal of useless jumps and optimizations on conditional jumps

Latest enhancement:

```
changeset: 68375:14205d0fee45
user:      Antoine Pitrou <solipsis@pitrou.net>
date:      Fri Mar 11 17:27:02 2011 +0100
files:     Lib/test/test_peepholer.py Misc/NEWS Python/peephole.c
description:
Issue #11244: The peephole optimizer is now able to constant-fold
arbitrarily complex expressions. This also fixes a 3.2 regression where
operations involving negative numbers were not constant-folded.
```

Should be rewritten as an *AST optimizer*.

4.2 Bytecode

- [byteplay](#)
- [diving-into-byte-code-optimization-in-python](#)
- [BytecodeAssembler](#)
- [ByteplayDoc](#)
- [Hacking Python bytecode](#)

Python C API

5.1 Intro

CPython comes with a C API called the “Python C API”. The most common type is `PyObject*` and functions are prefixed with `Py` (and `_Py` for private functions but you must not use them!).

5.2 Historical design choices

CPython was created in 1991 by Guido van Rossum. Some design choices made sense in 1991 but don’t make sense anymore in 2015. For example, the *GIL* was a simple and safe choice to implement multithreading in CPython. But in 2015, smartphones have 2 or 4 cores, and desktop PC have between 4 and 8 cores. The GIL restricts peak performances on multithreaded applications, even when it’s possible to release the GIL.

5.3 GIL

CPython uses a Global Interpreter Lock called “GIL” to avoid concurrent accesses to CPython internal structures (shared resources like global variables) to ensure that Python internals remain consistent.

See also *Kill the GIL*.

5.4 Reference counting and garbage collector

The C structure of all Python objects inherit from the `PyObject` structure which contains the field `Py_ssize_t ob_refcnt`; . This is a simple counter initialized to 1 when the object is created, increased each time that a variable has a strong reference to the object, and decreased each time that a strong reference is removed. The object is removed when the counter reached 0.

In some cases, two objects are linked together. For example, A has a strong reference to B which has a strong reference to A. Even if A and B are no more referenced outside, these objects are not destroyed because their reference counter is still equal to 1. A garbage collector is responsible to find and break *reference cycles*.

See also the [PEP 442: Safe object finalization](#) implemented in Python 3.4 which helps to break reference cycles.

5.5 Popular projects using the Python C API

- `numpy`
- `PyQt`
- `Mercurial`

AST Optimizers

6.1 Intro

An AST optimizer rewrites the Abstract Syntax Tree (AST) of a Python module to produce a more efficient code.

Currently in CPython 3.5, only basic optimizations are implemented by rewriting the bytecode: *CPython peephole optimizer*.

6.2 Old AST optimizer project

See *old AST optimizer*.

6.3 fatoptimizer

fatoptimizer project: AST optimizer implementing multiple optimizations and can specialize functions using guards of the `fat` module.

6.4 pythran AST

`pythran.analysis.PureFunctions` of pythran project, depend on `ArgumentEffects` and `GlobalEffects` analysis: automatically detect pure functions.

6.5 PyPy AST optimizer

<https://bitbucket.org/pypy/pypy/src/default/pypy/interpreter/astcompiler/optimize.py>

6.6 Cython AST optimizer

<https://mail.python.org/pipermail/python-dev/2012-August/121300.html>

- `Compiler/Optimize.py`
- `Compiler/ParseTreeTransforms.py`

- [Compiler/Builtin.py](#)
- [Compiler/Pipeline.py](#)

6.7 Links

6.7.1 CPython issues

- [Issue #2181](#): optimize out local variables at end of function
- [Issue #2499](#): Fold unary + and not on constants
- [Issue #4264](#): Patch: optimize code to use LIST_APPEND instead of calling list.append
- [Issue #7682](#): Optimisation of if with constant expression
- [Issue #11549](#): Build-out an AST optimizer, moving some functionality out of the peephole optimizer
- [Issue #17068](#): peephole optimization for constant strings
- [Issue #17430](#): missed peephole optimization

6.7.2 AST

- [instrumenting_the_ast.html](#)
- [the-internals-of-python-generator-functions-in-the-ast](#)
- [tlee-ast-optimize](#) branch
- [ast-optimization-branch-elimination-in-generator-functions](#)

Old AST Optimizer

See also *AST optimizers*.

<https://bitbucket.org/haypo/astoptimizer/> was a first attempt to optimize Python. This project was rejected by the Python community because it breaks the Python semantics. For example, it replaces `len("abc")` with `3`. It checks that `len()` was not overridden in the module, but it doesn't check that the builtin `len()` function was not overridden.

Threads on the Python-Dev mailing list:

- [astoptimizer: static optimizer working on the AST](#) (March 2013)
- [Release of astoptimizer 0.3](#) (September 2012): Guido van Rossum, Nick Coghlan and Maciej Fijalkowski complained that too many optimizations broke the Python semantics
- [AST optimizer implemented in Python](#) (August 2012)

The project was created in September 2012. It is now dead and replaced with the `fatoptimizer` project.

7.1 Introduction

`astoptimizer` is an optimizer for Python code working on the Abstract Syntax Tree (AST, high-level representation). It does as much work as possible at compile time.

The compiler is static, it is not a just-in-time (JIT) compiler, and so don't expect better performances than `psyco` or `PyPy` for example. Optimizations depending on the type of functions parameters cannot be done for examples. Only optimizations on immutable types (constants) are done.

Website: <http://pypi.python.org/pypi/astoptimizer>

Source code hosted at: <https://bitbucket.org/haypo/astoptimizer>

7.2 Optimizations

- Call builtin functions if arguments are constants (need “builtin_funcs” feature). Examples:
 - `len("abc") => 3`
 - `ord("A") => 65`
- Call methods of builtin types if the object and arguments are constants. Examples:
 - `u"h\\xe9ho".encode("utf-8") => b"h\\xc3\\xa9ho"`
 - `"python2.7".startswith("python") => True`

- `(32).bit_length() => 6`
- `float.fromhex("0x1.8p+0") => 1.5`
- Call functions of `math` and `string` modules for functions without border effect. Examples:
 - `math.log(32) / math.log(2) => 5.0`
 - `string.atoi("5") => 5`
- Format strings for `str%args` and `print(arg1, arg2, ...)` if arguments are constants and the format string is valid. Examples:
 - `"x=%s" % 5 => "x=5"`
 - `print(1.5) => print("1.5")`
- Simplify expressions. Examples:
 - `not(x in y) => x not in y`
 - `4 and 5 and x and 6 => x and 6`
 - `if a: if b: print("true") => if a and b: print("true")`
- Optimize loops (`range => xrange` needs “`builtin_funcs`” features). Examples:
 - `while True: pass => while 1: pass`
 - `for x in range(3): print(x) => x = 0; print(x); x = 1; print(x); x = 2; print(x)`
 - `for x in range(1000): print(x) => for x in xrange(1000): print(x)` (Python 2)
- Optimize iterators, list, set and dict comprehension, and generators (need “`builtin_funcs`” feature). Examples:
 - `iter(set()) => iter()`
 - `frozenset("") => frozenset()`
 - `(x for x in "abc" if False) => (None for x in ())`
 - `[x*10 for x in range(1, 4)] => [10, 20, 30]`
 - `(x*2 for x in "abc" if True) => (x*2 for x in ("a", "b", "c"))`
 - `list(x for x in iterable) => list(iterable)`
 - `tuple(x for x in "abc") => ("a", "b", "c")`
 - `list(x for x in range(3)) => [0, 1, 2]`
 - `[x for x in ""] => []`
 - `[x for x in iterable] => list(iterable)`
 - `set([x for x in "abc"]) => {"a", "b", "c"} (Python 2.7+) or set(("a", "b", "c"))`
- Replace list with tuple (need “`builtin_funcs`” feature). Examples:
 - `for x in [a, b, c]: print(x) => for x in (a, b, c): print(x)`
 - `x in [1, 2, 3] => x in (1, 2, 3)`
 - `list([x, y, z]) => [x, y, z]`
 - `set([1, 2, 3]) => {1, 2, 3} (Python 2.7+)`

- Evaluate unary and binary operators, subscript and comparison if all arguments are constants. Examples:

```
- 1 + 2 * 3 => 7
- not True => False
- "abc" * 3 => "abccabccabc"
- "abcdef"[:3] => "abc"
- (2, 7, 3)[1] => 7
- frozenset("ab") | frozenset("bc") => frozenset("abc")
- None is None => True
- "2" in "python2.7" => True
- x in [1, 2, 3] => x in {1, 2, 3} (Python 3) or x in (1, 2, 3) (Python 2)
- def f(): return 2 if 4 < 5 else 3 => def f(): return 2
```

- Remove empty loop. Example:

```
- for i in (1, 2, 3): pass => i = 3
```

- Remove dead code. Examples:

```
- def f(): return 1; return 2 => def f(): return 1
- def f(a, b): s = a+b; 3; return s => def f(a, b): s = a+b; return s
- if DEBUG: print("debug") => pass with DEBUG declared as False
- while 0: print("never executed") => pass
```

7.3 Use astoptimizer in your project

To enable astoptimizer globally on your project, add the following lines at the very beginning of your application:

```
import astoptimizer
config = astoptimizer.Config('builtin_funcs', 'pythonbin')
# customize the config here
astoptimizer.patch_compile(config)
```

On Python 3.3, imports will then use the patched compile() function and so all modules will be optimized. With older versions, the compileall module (ex: compileall.compile_dir()) can be used to compile an application with optimizations enabled.

See also the issue [#17515: Add sys.setasthook\(\) to allow to use a custom AST optimizer](#).

7.4 Example

Example with the high-level function optimize_code:

```
from astoptimizer import optimize_code
code = "print(1+1)"
code = optimize_code(code)
exec(code)
```

Example the low-level functions optimize_ast:

```
from astoptimizer import Config, parse_ast, optimize_ast, compile_ast
config = Config('builtin_funcs', 'pythonbin')
code = "print(1+1)"
tree = parse_ast(code)
tree = optimize_ast(tree, config)
code = compile_ast(tree)
exec(code)
```

See also `demo.py` script.

7.5 Configuration

Unsafe optimizations are disabled by default. Use the `Config()` class to enable more optimizations.

Features enabled by default:

- `"builtin_types"`: methods of bytes, str, unicode, tuple, frozenset, int and float types
- `"math"`, `"string"`: constants and functions without border effects of the math / string module

Optional features:

- `"builtin_funcs"`: builtin functions like `abs()`, `str()`, `len()`, etc. Examples:
 - `len("abc") => 3`
 - `ord("A") => 65`
 - `str(123) => "123"`
- `"pythonbin"`: Enable this feature if the optimized code will be executed by the same Python binary: so exactly the same Python version with the same build options. Allow to optimize non-BMP unicode strings on Python < 3.3. Enable the `"platform"` feature. Examples:
 - `u"\\U0010ffff"[0] => u"\\udbff"` or `u"\\U0010ffff"` (depending on build options, narrow or wide Unicode)
 - `sys.version_info.major => 2`
 - `sys.maxunicode => 0x10ffff`
- `"pythonenv"`: Enable this feature if you control the environment variables (like `PYTHONOPTIMIZE`) and Python command line options (like `-Qnew`). On Python 2, allow to optimize int/int. Enable `"platform"` and `"pythonbin"` features. Examples:
 - `__debug__ => True`
 - `sys.flags.optimize => 0`
- `"platform"`: optimizations specific to a platform. Examples:
 - `sys.platform => "linux2"`
 - `sys.byteorder => "little"`
 - `sys.maxint => 2147483647`
 - `os.linesep => "\\n"`
- `"struct"`: struct module, `calcsite()`, `pack()` and `unpack()` functions.
- `"cpython_tests"`: disable some optimizations to workaround issues with the CPython test suite. Only use it for tests.

Use `Config("builtin_funcs", "pythonbin")` to enable most optimizations. You may also enable `"pythonenv"` to enable more optimizations, but then the optimized code will depend on environment variables and Python command line options.

Use `config.enable_all_optimizations()` to enable all optimizations, which may generate invalid code.

7.6 Advices

Advices to help the AST optimizer:

- Declare your constants using `config.add_constant()`
- Declare your pure functions (functions with no border effect) using `config.add_func()`
- Don't use `"from module import *"`. If `"import *"` is used, builtin functions are not optimized anymore for example.

7.7 Limitations

- Operations on mutable values are not optimized, ex: `len([1, 2, 3])`.
- Unsafe optimizations are disabled by default. For example, `len("\U0010ffff")` is not optimized because the result depends on the build options of Python. Enable `"builtin_funcs"` and `"pythonenv"` features to enable more optimizations.
- `len()` is not optimized if the result is bigger than $2^{31}-1$. Enable `"pythonbin"` configuration feature to optimize the call for bigger objects.
- On Python 2, operators taking a bytes string and a unicode string are not optimized if the bytes string has to be decoded from the default encoding or if the unicode string has to be encoded to the default encoding. Exception: pure ASCII strings are optimized. For example, `b"abc" + u"def"` is replaced with `u"abcdef"`, whereas `u"x=%s" % b"\xe9"` is not optimized.
- On Python 3, comparison between bytes and Unicode strings are not optimized because the comparison may emit a warning or raise a `BytesWarning` exception. Bytes string are not converted to Unicode string. For example, `b"abc" < "abc"` and `str(b"abc")` are not optimized. Converting a bytes string to Unicode is never optimized.

7.8 ChangeLog

7.8.1 Version 0.6 (2014-03-05)

- Remove empty loop. Example: `for i in (1, 2, 3): pass => i = 3.`
- Log removal of code
- Fix support of Python 3.4: socket constants are now enum

7.8.2 Version 0.5 (2013-03-26)

- Unroll loops (no support for break/continue yet) and list comprehension. Example: `[x*10 for x in range(1, 4)] => [10, 20, 30].`

- Add `Config.enable_all_optimizations()` method
- Add a more aggressive option to remove dead code (`config.remove_almost_dead_code`), disabled by default
- Remove useless instructions. Example: `"x=1; 'abc'; print(x)" => "x=1; print(x)"`
- Remove empty try/except. Example: `"try: pass except: pass" => "pass"`

7.8.3 Version 0.4 (2012-12-10)

Bugfixes:

- Don't replace `range()` with `xrange()` if arguments cannot be converted to C long
- Disable `float.fromhex()` optimization by default: float may be shadowed. Use `"builtin_funcs"` to enable this optimization.

Changes:

- Add the `"struct"` configuration feature: functions of the `struct` module
- Optimize `print()` on Python 2 with `"from __future__ import print_function"`
- Optimize iterators, list, set and dict comprehension, and generators
- Replace list with tuple
- Optimize `if a: if b: print("true"): if a and b: print("true")`

7.8.4 Version 0.3.1 (2012-09-12)

Bugfixes:

- Disable optimizations on functions and constants if a variable with the same name is set. Example: `"len=ord; print(len('A'))"`, `"sys.version = 'abc'; print(sys.version)"`.
- Don't optimize `print()` function, `frozenset()` nor `range()` functions if `"builtin_funcs"` feature is disabled
- Don't remove code if it contains global or nonlocal. Example: `"def f(): if 0: global x; x = 2"`.

7.8.5 Version 0.3 (2012-09-11)

Major changes:

- Add `astoptimizer.patch_compile(config=None)` function to simply hook the builtin `compile()` function.
- Add `"pythonbin"` configuration feature.
- Disable optimizations on builtin functions by default. Add `"builtin_funcs"` feature to the configuration to optimize builtin functions.
- Remove dead code (optionnal optimization)
- It is now possible to define a callback for warnings of the optimizer
- Drop support of Python 2.5, it is unable to compile an AST tree to bytecode. AST objects of Python 2.5 don't accept arguments in constructors.

Bugfixes:

- Handle `"from math import *"` correctly

- Don't optimize operations if arguments are bytes and unicode strings. Only optimize if string arguments have the same type.
- Disable optimizations on non-BMP unicode strings by default. Optimizations enabled with "pythonbin" feature.

Other changes:

- More functions, methods and constants:
 - bytes, str, unicode: add more methods.
 - math module: add most remaining functions
 - string module: add some functions and all constants
- `not(a in b) => a not in b`, `not(a is b) => a is not b`
- `a if bool else b`
- `for x in range(n) => for x in xrange(n)` (only on Python 2)
- Enable more optimizations if a function is not a generator
- Add `sys.flags.<attr>` and `sys.version_info.<attr>` constants

7.8.6 Version 0.2 (2012-09-02)

Major changes:

- Check input arguments before calling an operator or a function, instead of catching errors.
- New helper functions `optimize_code()` and `optimize_ast()` should be used instead of using directly the `Optimizer` class.
- Support tuple and frozenset types

Changes:

- FIX: add `Config.max_size` to check `len(obj)` result
- FIX: disable non portable optimizations on non-BMP strings
- Support Python 2.5-3.3
- Refactor `Optimizer`: `Optimizer.visit()` now always visit children before calling the optimizer for a node, except for assignments
- Float and complex numbers are no more restricted by the integer range of the configuration
- More builtin functions. Examples: `divmod(int, int)`, `float(str)`, `min(tuple)`, `sum(tuple)`.
- More method of builtin types. Examples: `str.startswith()`, `str.find()`, `tuple.count()`, `float.is_integer()`.
- math module: add `math.ceil()`, `math.floor()` and `math.trunc()`.
- More module constants. Examples: `os.O_RDONLY`, `errno.EINVAL`, `socket.SOCK_STREAM`.
- More operators: `a not in b`, `a is b`, `a is not b`, `+a`.
- Conversion to string: `str()`, `str % args` and `print(arg1, arg2, ...)`.
- Support import aliases. Examples: "import math as M; print(M.floor(1.5))" and "from math import floor as F; print(F(1.5))".
- Experimental support of variables (disabled by default).

7.8.7 Version 0.1 (2012-08-12)

- First public version (to reserve the name on PyPI!)

Register-based Virtual Machine for Python

8.1 Intro

registervm is a fork of CPython 3.3 using register-based bytecode, instead of stack-code bytecode

More information: [REGISTERVM.txt](#)

Thread on the Python-Dev mailing list: [Register-based VM for CPython](#).

The project was created in November 2012.

8.2 Status

- Most instructions using the stack are converted to instructions using registers
- Bytecode using registers with all optimizations enable is usually 10% faster than bytecode using the stack, according to pybench
- registervm generates invalid code, see TODO section below, so it's not possible yet to use it on the Python test suite

8.3 TODO

8.3.1 Bugs

- Register allocator doesn't handle correctly conditional branches: CLEAR_REG is removed on the wrong branch in test_move_instr.
- Fail to track the stack state in if/else. Bug hidden by the register allocator in the following example:
def func(obj): obj.attr = sys.modules['warnings'] if module is None else module
- Don't move globals out of if. Only out of loops? subprocess.py:

```
if mswindows:
    if p2cwrite != -1:
        p2cwrite = msvcrt.open_osfhandle(p2cwrite.Detach(), 0)
```

But do move len() out of loop for:

```
def loop_move_instr():
    length = 0
    for i in range(5):
        length += len("abc") - 1
    return length
```

- Don't remove duplicate LOAD_GLOBAL in "LOAD_GLOBAL ...; CALL_PROCEDURE ...; LOAD_GLOBAL ...": CALL_PROCEDURE has border effect
- Don't remove duplicate LOAD_NAME if a function has a border effect:

```
x=1
def modify():
    global x
    x = 2
print(x)
modify()
print(x)
```

8.3.2 Improvements

- Move LOAD_CONST out of loops: it was done in a previous version, but the optimization was broken by the introduction of CLEAR_REG
- Copy constants to the frame objects so constants can be used as registers and LOAD_CONST instructions can be simplify removed
- Enable move_load_const by default?
- Fix moving LOAD_ATTR_REG: only do that when calling methods. See test_sieve() of test_registervm: primes.append().

```
result = Result()
while 1:
    if result.done:
        break
    func(result)
```

- Reenable merging duplicate LOAD_ATTR
- Register allocation for locale_alias = {...} is very very slow
- "while 1: ... return" generates useless SETUP_LOOP
- Reuse locals?
- implement register version of the following instructions:
 - DELETE_ATTR
 - try/finally
 - yield from
 - CALL_FUNCTION_VAR_KW
 - CALL_FUNCTION_VAR
 - operators: $a \mid b$, $a \& b$, $a \wedge b$, $a \mid= b$, $a \&= b$, $a \wedge= b$
- Deref:
 - add a test using free variables

- Move LOAD_DEREF_REG out of loops
- NAME:
 - test_list_append() of test_registervm.py
 - Move LOAD_NAME_REG out of loop
- Handle JUMP_IF_TRUE_OR_POP: see test_getline() of test_registervm
- Compute the number of used registers in a frame
- Write a new test per configuration option
- Factorize code processing arg_types, ex: disassemblers of dis and registervm modules
- Add tests on class methods
- Fix Inotab

8.4 Changelog

2012-12-21

- Use RegisterTracker to merge duplicated LOAD, STORE_GLOBAL/LOAD_GLOBAL are now also simplified

2012-12-19

- Emit POP_REG to simplify the stack tracker

2012-12-18

- LOAD are now only moved out of loops

2012-12-14

- Duplicated LOAD instructions can be merged without moving them
- Rewrite the stack tracker: PUSH_REG don't need to be moved anymore
- Fix JUMP_IF_TRUE_OR_POP/JUMP_IF_FALSE_OR_POP to not generate invalid code
- Don't move LOAD_ATTR_REG out of try/except block

2012-12-11

- Split instructions into linked-blocks

2012-11-26

- Add a stack tracker

2012-11-20

- Remove useless jumps
- CALL_FUNCTION_REG and CALL_PROCEDURE_REG are fully implemented

2012-10-29

- Remove “if (HAS_ARG(op))” check in PyEval_EvalFrameEx()

2012-10-27

- Duplicated LOAD_CONST and LOAD_GLOBAL are merged (optimization disabled on LOAD_GLOBAL because it is buggy)

2012-10-23

- initial commit, 0f7f49b7083c

8.5 CPython 3.3 bytecode is inefficient

- Useless jump: `JUMP_ABSOLUTE <offset+0>`
- Generate dead code: `RETURN_VALUE; RETURN_VALUE` (the second instruction is unreachable)
- Duplicate constants: see TupleSlicing of pybench
- Constant folding: see astoptimizer project
- `STORE_NAME 'f'; LOAD_NAME 'f'`
- `STORE_GLOBAL 'x'; LOAD_GLOBAL 'x'`

8.6 Rationale

The performance of the loop evaluating bytecode is critical in Python. For Python example, using computed-goto instead of switch to dispatch bytecode improved performances by 20%. Related issues:

- [use computed goto's in ceval loop](#)
- [Faster opcode dispatch on gcc](#)
- [Computed-goto patch for RE engine](#)

Using registers of a stack reduce the number of operations, but increase the size of the code. I expect an significant speedup when all operations will use registers.

8.7 Optimizations

Optimizations:

- Remove useless `LOAD_NAME` and `LOAD_GLOBAL`. For example: `"STORE_NAME var; LOAD_NAME var"`
- Merge duplicate loads (`LOAD_CONST`, `LOAD_GLOBAL_REG`, `LOAD_ATTR`). For example, `"lst.append(1); lst.append(1)"` only gets constant `"1"` and the `"lst.append"` attribute once.

Misc:

- Automatically detect inplace operations. For example, `"x = x + y"` is compiled to `"BINARY_ADD_REG 'x', 'x', 'y'"` which calls `PyNumber_InPlaceAdd()`, instead of `PyNumber_Add()`.
- Move constant, global and attribute loads out of loops (to the beginning)
- Remove useless jumps (ex: `JUMP_FORWARD <relative jump to 103 (+0)>`)

8.8 Algorithm

The current implementation rewrites the stack-based operations to use register-based operations instead. For example, `"LOAD_GLOBAL range"` is replaced with `"LOAD_GLOBAL_REG R0, range; PUSH_REG R0"`. This first step is inefficient because it increases the number of operations.

Then, operations are reordered: PUSH_REG and POP_REG to the end. So we can replace “PUSH_REG R0; PUSH_REG R1; STACK_OPERATION; POP_REG R2” with a single operation: “REGISTER_OPERATION R2, R0, R1”.

Move invariant out of the loop: it is possible to move constants out of the loop. For example, LOAD_CONST_REG are moved to the beginning. We might also move LOAD_GLOBAL_REG and LOAD_ATTR_REG to the beginning.

Later, a new AST to bytecode compiler can be implemented to emit directly operations using registers.

8.9 Example

Simple function computing the factorial of n:

```
def fact_iter(n):
    f = 1
    for i in range(2, n+1):
        f *= i
    return f
```

Stack-based bytecode (20 instructions):

```
0 LOAD_CONST          1 (const#1)
3 STORE_FAST          'f'
6 SETUP_LOOP          <relative jump to 46 (+37)>
9 LOAD_GLOBAL         0 (range)
12 LOAD_CONST         2 (const#2)
15 LOAD_FAST          'n'
18 LOAD_CONST         1 (const#1)
21 BINARY_ADD
22 CALL_FUNCTION       2 (2 positional, 0 keyword pair)
25 GET_ITER
>> 26 FOR_ITER          <relative jump to 45 (+16)>
29 STORE_FAST         'i'
32 LOAD_FAST          'f'
35 LOAD_FAST          'i'
38 INPLACE_MULTIPLY
39 STORE_FAST         'f'
42 JUMP_ABSOLUTE      <jump to 26>
>> 45 POP_BLOCK
>> 46 LOAD_FAST         'f'
49 RETURN_VALUE
```

Register-based bytecode (13 instructions):

```
0 LOAD_CONST_REG      'f', 1 (const#1)
5 LOAD_CONST_REG      R0, 2 (const#2)
10 LOAD_GLOBAL_REG    R1, 'range' (name#0)
15 SETUP_LOOP         <relative jump to 57 (+39)>
18 BINARY_ADD_REG     R2, 'n', 'f'
25 CALL_FUNCTION_REG  4, R1, R1, R0, R2
36 GET_ITER_REG       R1, R1
>> 41 FOR_ITER_REG     'i', R1, <relative jump to 56 (+8)>
48 INPLACE_MULTIPLY_REG 'f', 'i'
53 JUMP_ABSOLUTE      <jump to 41>
>> 56 POP_BLOCK
>> 57 RETURN_VALUE_REG 'f'
```

The body of the main loop of this function is composed of 1 instructions instead of 5.

8.10 Comparative table

Example	S r R	Stack	Register
append(2)	4 1 2	LOAD_FAST 'append' LOAD_CONST 2 (const#2) CALL_FUNCTION (1 positional) POP_TOP	LOAD_CONST_REG R1, 2 (const#2) CALL_PROCEDURE_REG 'append', (1 positional), R1
l[0] = 3	4 1 2	LOAD_CONST 3 (const#1) LOAD_FAST 'l' LOAD_CONST 0 (const#4) STORE_SUBSCR	LOAD_CONST_REG R0, 3 (const#1) LOAD_CONST_REG R3, 0 (const#4) ... STORE_SUBSCR_REG 'l', R3, R0
x = l[0]	4 1 2	LOAD_FAST 'l' LOAD_CONST 0 (const#4) BINARY_SUBSCR STORE_FAST 'x'	LOAD_CONST_REG R3, 0 (const#4) BINARY_SUBSCR_REG 'x', 'l', R3
s.isalnum()	4 1 2	LOAD_FAST 's' LOAD_ATTR 'isalnum' (name#3) CALL_FUNCTION (0 positional) POP_TOP	LOAD_ATTR_REG R5, 's', 'isalnum' (name#3) CALL_PROCEDURE_REG R5, (0 positional)
o.a = 2	3 1 2	LOAD_CONST 2 (const#3) LOAD_FAST 'o' STORE_ATTR 'a' (name#2)	LOAD_CONST_REG R2, 2 (const#3) ... STORE_ATTR_REG 'o', 'a' (name#2), R2
x = o.a	3 1 1	LOAD_FAST 'o' LOAD_ATTR 'a' (name#2) STORE_FAST 'x'	LOAD_ATTR_REG 'x', 'o', 'a' (name#2)

Columns:

- “S”: Number of stack-based instructions
- “r”: Number of stack-based instructions excluding instructions moved out of loops (ex: LOAD_CONST_REG)
- “R”: Total number of stack-based instructions (including instructions moved out of loops)

Read-only Python

9.1 Intro

A first attempt to implement guards was the [readonly PoC](#) (fork of CPython 3.5) which registered callbacks to notify all guards. The problem is that modifying a watched dictionary gets a complexity of $O(n)$ where n is the number of registered guards.

`readonly` adds a `modified` flag to types and a `readonly` property to dictionaries. The guard was notified with the modified key to decide to disable or not the optimization.

More information: [READONLY.txt](#)

Thread on the python-ideas mailing list: [Make Python code read-only](#).

The project was mostly developed in May 2014. The project is now dead, replaced with [FAT Python](#).

9.2 READONLY

This fork on CPython 3.5 adds a machinery to be notified when the Python code is modified. Modules, classes (types) and functions are tracked. At the first modification, a callback is called with the object and the modified attribute.

This machinery should help static optimizers. See this article for more information: http://haypo-notes.readthedocs.org/faster_cpython.html

Examples of such optimizers:

- `astoptimizer` project: replace a function call by its result during the AST compilation
- Learn types of function parameters and local variables, and then compile Python (byte)code to machine code specialized for these types (like Cython)

9.3 Issues with read-only code

- Currently, it's not possible to allow again to modify a module, class or function to keep my implementation simple. With a registry of callbacks, it may be possible to enable again modification and call code to disable optimizations.
- PyPy implements this but thanks to its JIT, it can optimize again the modified code during the execution. Writing a JIT is very complex, I'm trying to find a compromise between the fast PyPy and the slow CPython. Add a JIT to CPython is out of my scope, it requires too much modifications of the code.

- With read-only code, monkey-patching cannot be used anymore. It's annoying to run tests. An obvious solution is to disable read-only mode to run tests, which can be seen as unsafe since tests are usually used to trust the code.
- The `sys` module cannot be made read-only because modifying `sys.stdout` and `sys.ps1` is a common use case.
- The `warnings` module tries to add a `__warningregistry__` global variable in the module where the warning was emitted to not repeat warnings that should only be emitted once. The problem is that the module namespace is made read-only before this variable is added. A workaround would be to maintain these dictionaries in the `warnings` module directly, but it becomes harder to clear the dictionary when a module is unloaded or reloaded. Another workaround is to add `__warningregistry__` before making a module read-only.
- Lazy initialization of module variables does not work anymore. A workaround is to use a mutable type. It can be a dict used as a namespace for module modifiable variables.
- The interactive interpreter sets a `"_"` variable in the builtins namespace. I have no workaround for this. The `"_"` variable is no more created in read-only mode. Don't run the interactive interpreter in read-only mode.
- It is not possible yet to make the namespace of packages read-only. For example, `"import encodings.utf_8"` adds the symbol `"utf_8"` to the `encodings` namespace. A workaround is to load all submodules before making the namespace read-only. This cannot be done for some large modules. For example, the `encodings` has a lot of submodules, only a few are needed.

9.4 STATUS

- Python API:
 - new `function.__modified__` and `type.__modified__` properties: False by default, becomes True when the object is modified
 - new `module.is_modified()` method
 - new `module.set_initialized()` method
- C API:
 - `PyDictObject`: new `"int ma_readonly;"` field
 - `PyTypeObject`: a new `"int tp_modified;"` field
 - `PyFunctionObject`: new `"int func_module;"` and `"int func_initialized;"` fields
 - `PyModuleObject`: new `"int md_initialized;"` field

9.5 Modified modules, classes and functions

- It's common to modify the following attributes of the `sys` module:
 - `sys.ps1`, `sys.ps2`
 - `sys.stdin`, `sys.stdout`, `sys.stderr`
- `"import encodings.latin_1"` sets `"latin_1"` attribute in the namespace of the `"encodings"` module.
- The interactive interpreter sets the `"_"` variable in builtins.
- `warnings`: global variable `__warningregistry__` set in modules
- `functools.wraps()` modifies the wrapper to copy attributes of the wrapped function

9.6 TODO

- builtins modified in `initstdio()`: `builtins.open` modified
- `sys` modified in `initstdio()`: `sys.__stdin__` modified
- `structseq`: types are created modified; same issue with `_ast` types (`Python-ast.c`)
- module, type and function `__dict__`:
 - Drop `dict.setreadonly()`
 - Decide if it's better to use `dict.setreadonly()` or a new subclass (ex: “`dict_maybe_readonly`” or “`namespace`”).
 - Read only dict: add a new `ReadOnlyError` instead of `ValueError`?
 - `sysmodule.c`: `PyDict_DelItemString(FlagsType.tp_dict, “__new__”)` doesn't mark `FlagsType` as modified
 - Getting `func.__dict__` / `module.__dict__` marks the function/module as modified, this is wrong. Use instead a mapping marking the function as modified when the mapping is modified.
 - `module.__dict__` is read-only: similar issue for functions.
- Import submodule. Example: “`import encodings.utf_8`” modifies “`encoding`” to set a new `utf_8` attribute

9.7 TODO: Specialized functions

9.7.1 Environment

- module and type attribute values:
 - (“`module`”, “`os`”, `OS_CHECKSUM`)
 - (“`attribute`”, “`os.path`”)
 - (“`module`”, “`path`”, `PATH_CHECKSUM`)
 - (“`attribute`”, “`path.isabs`”)
 - (“`function`”, “`path.isabs`”)
- function attributes
- set of function parameter types (passed as indexed or keyword arguments)

9.7.2 Read-only state

Scenario:

- 1: `application.py` is compiled. Function A depends on `os.path.isabs`, function B depends on `project.DEBUG`
- 2: application is started, “`import os.path`”
- 3: `os.path.isabs` is modified
- 4: optimized `application.py` is loaded
- 5: `project.DEBUG` is modified

When the function is created, `os.path.isabs` was already modified compared to the `OS_CHECKSUM`.

9.7.3 Example of environments

- The function calls “os.path.isabs”:
 - rely on “os.path” attribute
 - rely on “os.path.isabs” attribute
 - rely on “os.path.isabs” function attributes (except `__doc__`)
- The function “def mysum(x, y):” has two parameters
 - x type is int and y type is int
 - or: x type is str and y type is str
 - (“type is”: check the exact type, not a subclass)
- The function uses “project.DEBUG” constant
 - rely on “project.DEBUG” attribute

9.7.4 Content of a function

- classic attributes: doc, etc.
- multiple versions of the code:
 - required environment of the code
 - bytecode

9.7.5 Create a function

- build the environment
- register on module, type and functions modification

9.7.6 Callback when then environment is modified

xxx

9.7.7 Call a function

xxx

9.8 LINKS

- <http://legacy.python.org/dev/peps/pep-0351/> : Get an immutable copy of arbitrary objects
- <http://legacy.python.org/dev/peps/pep-0416/> : add a new frozendict type => types.MappingProxy added to Python 3.3

History of Python optimizations

- 2002: Creation of the psyco project by Armin Rigo
- 2003-05-05: psyco 1.0 released
- Spring 1997: [Creation of Jython project](#) (initially called JPython) by Jim Hugunin
- xxxx-xx-xx: Creation of IronPython project by xxx
- Creation of PyPy, spin-off of psyco
- mid-2007: PyPy 1.0 released.
- 2009-03: Creation of Unladen Swallow project by xxx
- xxxx-xx-xx: Creation of Pyston project by xxx
- 2012-09: Creation of the [AST optimizer](#) project by Victor Stinner
- 2012-11: Creation of the [registervm](#) project by Victor Stinner
- 2014-05: Creation of [read-only Python](#) PoC by Victor Stinner
- 2015-10: Creation of the [FAT Python](#) project by Victor Stinner

11.1 Ideas

- PyPy CALL_METHOD instructor
- Lazy formatting of Exception message: in most cases, the message is not used. `AttributeError(message) => AttributeError(attr=name)`, lazy formatting for `str(exc)` and `exc.args`.

11.2 Plan

- Modify CPython to be *notified when the Python code is changed*
- *Learn types* of function parameters and variables
- Compile a specialized version of a function using types and platform informations: more efficient bytecode using an *AST optimizer*, or even *emit machine code*. The compilation is done once and not during the execution, it's not a JIT.
- *Choose between bytecode and specialized code* at runtime

Other idea:

- `registervm`: My fork of Python 3.3 using register-based bytecode, instead of stack-code bytecode. Read `REGISTERVM.txt`
- *Kill the GIL?*

11.3 Status

See also the status of individual projects:

- `READONLY.txt`
- `REGISTERVM.txt`
- `astoptimizer` TODO list

11.3.1 Done

- astoptimizer project exists: [astoptimizer](#).
- Fork of CPython 3.5: be notified when the Python code is changed: modules, types and functions are tracked. My fork of CPython 3.5: [readonly](#); read [READONLY.txt](#) documentation.

Note: “readonly” is no more a good name for the project. The name comes from a first implementation using read-only code.

11.3.2 To do

- Learn types
- Enhance astoptimizer to use the type information
- Emit machine code

11.4 Why Python is slow?

11.4.1 Why the CPython implementation is slower than PyPy?

- everything is stored as an object, even simple types like integers or characters. Computing the sum of two numbers requires to “unbox” objects, compute the sum, and “box” the result.
- Python maintains different states: thread state, interpreter state, frames, etc. These informations are available in Python. The common usecase is to display a traceback in case of a bug. PyPy builds frames on demand.
- Cost of maintaining the reference counter: Python programs rely on the garbage collector
- `ceval.c` uses a virtual stack instead of CPU registers

11.4.2 Why the Python language is slower than C?

- modules are mutable, classes are mutable, etc. Because of that, it is not possible to inline code nor replace a function call by its result (ex: `len(“abc”)`).
- The types of function parameters and variables are unknown. Example of missing optimizations:
 - “`obj.attr`” instruction cannot be moved out of a loop: “`obj.attr`” may return a different result at each call, or execute arbitrary Python code
 - `x+0` raises a `TypeError` for “`abc`”, whereas it is a noop for `int` (it can be replaced with just `x`)
 - conditional code becomes dead code when types are known
- `obj.method` creates a temporary bounded method

11.4.3 Why improving CPython instead of writing a new implementation?

- There are already a lot of other Python implementations. Some examples: PyPy, Jython, IronPython, Pyston.
- CPython remains the reference implementation: new features are first implemented in CPython. For example, PyPy doesn’t support Python 3 yet.

- Important third party modules rely heavily on CPython implementation details, especially the *Python C API*. Examples: numpy and PyQt.

11.4.4 Why not a JIT?

- write a JIT is much more complex, it requires deep changes in CPython; CPython code is old (+20 years)
- cost to “warn up” the JIT: Mercurial project is concerned by the Python startup time
- Store generated machine code?

11.5 Learn types

- Add code in the compiler to record types of function calls. Run your program. Use recorded types.
- Range of numbers (predict C int overflow)
- Optional paramters: forceload=0. Dead code with forceload=0.
- Count number of calls to the function to decide if it should be optimized or not.
- Measure time spend in a function. It can be used to decide if it’s useful to release or not the GIL.
- Store type information directly in the source code? Manual type annotation?

11.6 Emit machine code

- Limited to simple types like integers?
- Use LLVM?
- Reuse Cython or numba?
- Replace bytecode with C functions calls. Ex: instead of `PyNumber_Add(a, b)` for `a+b`, emit `PyUnicode_Concat(a, b)`, `long_add(a, b)` or even simpler code without unbox/box
- Calling convention: have two versions of the function? only emit the C version if it is needed?
 - Called from Python: `Python C API, PyObject* func(PyObject *args, PyObject *kwargs)`
 - Called from C (specialized machine code): `C API, int func(char a, double d)`
 - Version which doesn’t need the GIL to be locked?
- Option to compile a whole application into machine code for proprietary software?

11.6.1 Example of (specialized) machine code

Python code:

```
def mysum(a, b):
    return a + b
```

Python bytecode:

0	LOAD_FAST	0	(a)
3	LOAD_FAST	1	(b)
6	BINARY_ADD		
7	RETURN_VALUE		

C code used to executed bytecode (without code to read bytecode and handle signals):

```
/* LOAD_FAST */
{
    PyObject *value = GETLOCAL(0);
    if (value == NULL) {
        format_exc_check_arg(PyExc_UnboundLocalError, ...);
        goto error;
    }
    Py_INCREF(value);
    PUSH(value);
}

/* LOAD_FAST */
{
    PyObject *value = GETLOCAL(1);
    if (value == NULL) {
        format_exc_check_arg(PyExc_UnboundLocalError, ...);
        goto error;
    }
    Py_INCREF(value);
    PUSH(value);
}

/* BINARY_ADD */
{
    PyObject *right = POP();
    PyObject *left = TOP();
    PyObject *sum;
    if (PyUnicode_CheckExact(left) &&
        PyUnicode_CheckExact(right)) {
        sum = unicode_concatenate(left, right, f, next_instr);
        /* unicode_concatenate consumed the ref to v */
    }
    else {
        sum = PyNumber_Add(left, right);
        Py_DECREF(left);
    }
    Py_DECREF(right);
    SET_TOP(sum);
    if (sum == NULL)
        goto error;
}

/* RETURN_VALUE */
{
    retval = POP();
    why = WHY_RETURN;
    goto fast_block_end;
}
```

Specialized and simplified C code if both arguments are Unicode strings:

```

/* LOAD_FAST */
PyObject *left = GETLOCAL(0);
if (left == NULL) {
    format_exc_check_arg(PyExc_UnboundLocalError, ...);
    goto error;
}
Py_INCREF(left);

/* LOAD_FAST */
PyObject *right = GETLOCAL(1);
if (right == NULL) {
    format_exc_check_arg(PyExc_UnboundLocalError, ...);
    goto error;
}
Py_INCREF(right);

/* BINARY_ADD */
PyUnicode_Append(&left, right);
Py_DECREF(right);
if (sum == NULL)
    goto error;

/* RETURN_VALUE */
retval = left;
why = WHY_RETURN;
goto fast_block_end;

```

11.7 Test if the specialized function can be used

Write code to choose between the bytecode evaluation and the machine code.

Preconditions:

- Check if `os.path.isabs()` was modified:
 - current namespace was modified? (os name cannot be replaced)
 - namespace of the `os.path` module was modified?
 - `os.path.isabs` function was modified?
 - compilation: checksum of the `os.py` and `posixpath.py`?
- Check the exact type of arguments
 - x type is str: in C, `PyUnicode_CheckExact(x)`
 - list of int: check the whole array before executing code? fallback in the specialized code to handle non int items?
- Callback to use the slow-path if something is modified?
- Disable optimizations when tracing is enabled
- Online benchmark to decide if preconditions and optimized code is faster than the original code?

Kill the GIL?

See the *Global Interpreter Lock*.

12.1 Why does CPython need a global lock?

Incomplete list:

- Python memory allocation is not thread safe (it should be easy to make it thread safe)
- The reference counter of each object is protected by the GIL.
- CPython has a lot of global C variables. Examples:
 - `interp` is a structure which contains variables of the Python interpreter: modules, list of Python threads, builtins, etc.
 - `int` singletons (-5..255)
 - `str` singletons (Python 3: latin1 characters)
- Some third party C libraries and even functions the C standard library are not thread safe: the GIL works around this limitation.

12.2 Kill the GIL

- Require deep changes of CPython code
- The current Python C API is too specific to CPython implementation details: need a new API. Maybe the stable ABI?
- Modify third party modules to use the stable ABI to avoid relying on CPython implementation details like reference counting
- Replace reference counting with something else? Atomic operations?
- Use finer locks on some specific operations (release the GIL)? like operations on builtin types which don't need to execute arbitrary Python code. Counter example: dict where keys are objects different than int and str.

See also `pyparallel`.

Implementations of Python

13.1 Faster Python implementations

- [PyPy](#)
 - AST optimizer of PyPy: [astcompiler/optimize.py](#)
- [Pyston](#)
- [Hotpy](#) and [Hotpy 2](#), based on [GVM](#) (Glasgow Virtual Machine Toolkit)
- [Numba](#): JIT implemented with LLVM, specialized to numeric types (numpy)
- [pymothoa](#) uses LLVM (“don’t support classes nor exceptions”)
- [WPython](#): 16-bit word-codes instead of byte-codes
- [Cython](#)

13.2 Fully Python compliant

- [PyPy](#)
- [Jython](#) based on the JVM
- [IronPython](#) based on the .NET VM
- [Unladen Swallow](#), fork of CPython 2.6, use LLVM. No more maintained
 - Project [announced](#) in 2009, abandoned in 2011
 - [ProjectPlan](#)
 - [Unladen Swallow Retrospective](#)
 - [PEP 3146](#)
- [Pyjion](#)

13.3 Other

- Replace stack-based bytecode with register-based bytecode: old [registervm](#) project

13.4 Fully Python compliant??

- [psyco](#): JIT. The author of pysco, Armin Rigo, co-created the PyPy project.

13.5 Subset of Python to C++

- [Nuitka](#)
- [Python2C](#)
- [Shedskin](#)
- [pythran](#) (no class, set, dict, exception, file handling, ...)

13.6 Subset of Python

- [pymothoa](#): use LLVM; don't support classes nor exceptions.
- [unpython](#): Python to C
- [Perthon](#): Python to Perl
- [Copperhead](#): Python to GPU (Nvidia)

13.7 Language very close to Python

- [Cython](#): “Cython is a programming language based on Python, with extra syntax allowing for optional static type declarations.”
 - based on [Pyrex](#)

Benchmarks

- speed.pypy.org: compare PyPy to CPython 2.7 (what about Python 3?)
- [Intel Language Performance](#) (CPython, 2.7 and default branches)
- [CPython benchmarks](#), come from Unladen Swallow?

See also:

- [Python benchmark sizes](#)

PEP 509: Add a private version to dict

Moved to python.org: [PEP 509 – Add a private version to dict](#)

PEP 510: Specialized functions with guards

Moved to python.org: [PEP 510](https://python.org) – Specialized functions with guards.

PEP 511: API for AST transformers

Moved to python.org: [PEP 511 – API for AST transformers](#)

Random notes about PyPy

18.1 What is the problem with PyPy?

PyPy is fast, much faster than CPython, but it's still not widely used by users. What is the problem? Or what are the problems?

- Bad support of the *Python C API*: PyPy was written from scratch and uses different memory structures for objects. The cpyext module emulates the Python C API but it's slow.
- New features are first developed in CPython. In January 2016, PyPy only supports Python 2.7 and 3.2, whereas CPython is at the version 3.5. It's hard to have a single code base for Python 2.7 and 3.2, Python 3.3 reintroduced `u' . . . '` syntax for example.
- Not all modules are compatible with PyPy: see [PyPy Compatibility Wiki](#). For example, numpy is not compatible with PyPy, but there is a project under development: [pypy/numy](#). PyGTK, PyQt, PySide and wxPython libraries are not compatible with PyPy; these libraries heavily depend on the Python C API. GUI applications (ex: gajim, meld) using these libraries don't work on PyPy :-(Hopefully, a lot of popular modules are compatible with PyPy (ex: Django, Twisted).
- PyPy is slower than CPython on some use cases. Django: "PyPy runs the templating engine faster than CPython, but so far DB access is slower for some drivers." ([source of the quote](#))

If I understood correctly, Pyjston will have same problems than PyPy since it doesn't support the Python C API neither. Same issue for Pyjion?

Talks

Talks about Python optimizations:

- *PyCon UK 2014 - When performance matters ...* <<http://www.egenix.com/library/presentations/PyCon-UK-2014-When-performance-matters/>> by Marc-Andre Lemburg

Links

20.1 Misc links

- “Need for speed” sprint (2006)
- ceval.c: use registers?
 - Java: [Virtual Machine Showdown: Stack Versus Registers](#) (Yunhe Shi, David Gregg, Andrew Beatty, M. Anton Ertl, 2005)
 - Lua 5: [The Implementation of Lua 5.0](#) (Roberto Ierusalimschy, Luiz Henrique de Figueiredo, Waldemar Celes, 2005)
 - Python-ideas: [Register based interpreter](#)
 - [unladen-swallow](#): [ProjectPlan](#): “Using a JIT will also allow us to move Python from a stack-based machine to a register machine, which has been shown to improve performance in other similar languages (Ierusalimschy et al, 2005; Shi et al, 2005).”
- Use a more efficient VM
- [WPython](#): 16-bit word-codes instead of byte-codes
- [Hotpy](#) and [Hotpy 2](#): built using the [GVMT](#) (The Glasgow Virtual Machine Toolkit)
- Search for Python issues of type performance: <http://bugs.python.org/>
- [Volunteer developed free-threaded cross platform virtual machines?](#)

20.2 Other

- [ASP](#): ASP is a SEJITS (specialized embedded just-in-time compiler) toolkit for Python.
- [PerformanceTips](#)