Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Inline Caching meets Quickening

Stefan Brunthaler

Institut für Computersprachen
Technische Universität Wien

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Smalltalk

Deutsch and Schiffman: "Efficient Implementation of the
Smalltalk-80 System", 1984:

## "dynamic locality of type usage"
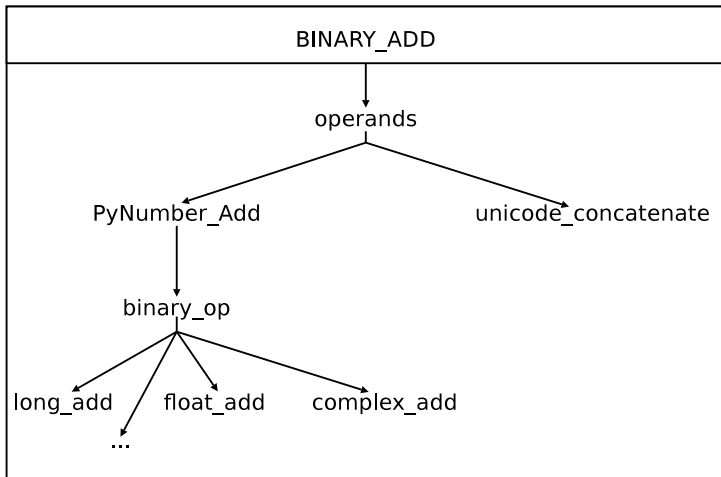
The actual operand types for an instruction within a sequence
tend to remain fixed for about 95% of the time.

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Inline Caching



Ad-hoc polymorphism in Python 3

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# **No** Inline Caching

Many popular interpreters do *not* use inline caching.

- Perl
- Python
- Ruby
- Tcl
- ...

None of the above have dynamic translators, or use known hash-table based look-up caches as used in Smalltalk implementations of the early 80s.
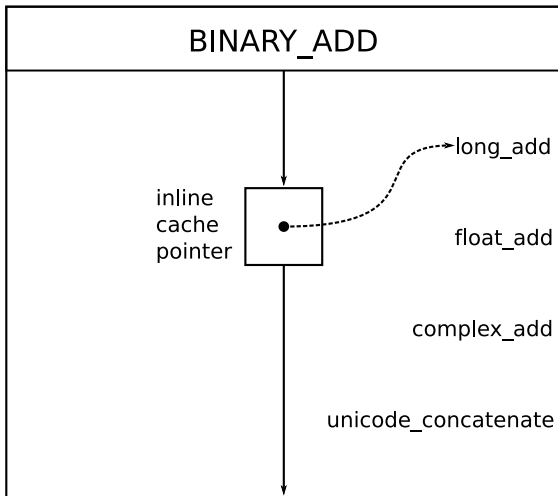
# Why do we care?

$$1000 : 10 : 1$$

Ertl, M. A. and D. Gregg, *The structure and performance of efficient interpreters.*, J. Instruction-Level Parallelism **5** (2003).

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
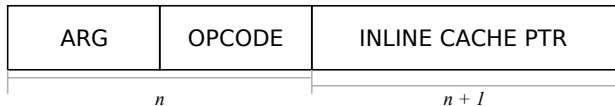**Basic Technique**
Quickening
Evaluation
Conclusion

# Adding Inline Cache Pointer



BINARY_ADD *instance* with long_add cached.

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
**Basic Technique**
Quickening
Evaluation
Conclusion

# Adding Inline Cache Pointer

| ARG | OPCODE | INLINE CACHE PTR |
|:---:|:---:|:---:|
| $n$ | | $n + 1$ |

Efficient basic technique without dynamic translation:

- transform into regular instruction format
- interleave instructions and inline cache pointers
- relocate jump offsets

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Discussion

Pros:

- easy to implement
- more efficient than Smalltalk-style look-up caches (**no** hash-tables)
- improved data locality

Cons:

- memory consumption ➡ profiling
- additional indirect call

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# What is quickening?

The Java virtual machine uses *quick instructions*, i.e., some instructions are replaced by more efficient *quick* instructions after the first execution. The non-quick instruction usually does initialization work.

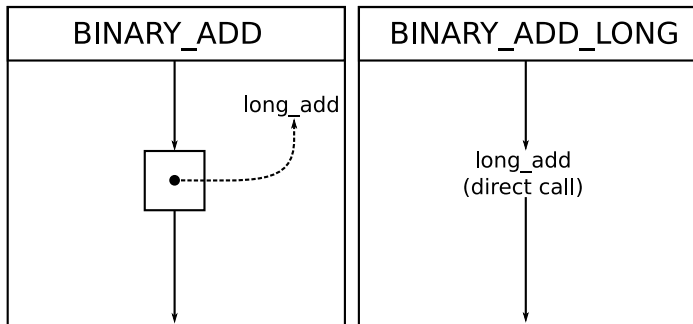Quickening specializes instructions with respect to their operand values.

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Idea

Create optimized derivatives based on result of resolving
ad-hoc polymorphism.

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
Evaluation
Conclusion

# Inline Caching meets Quickening



Replacing the indirect call with a direct call.

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
**Quickening**
Evaluation
Conclusion

# Discussion

Pros:

- easy to implement
- more efficient than previous technique
  ➦ eliminates indirect call
- enables inlining of functions by compiler
- does not require changing instruction format (for many use cases)

Cons:

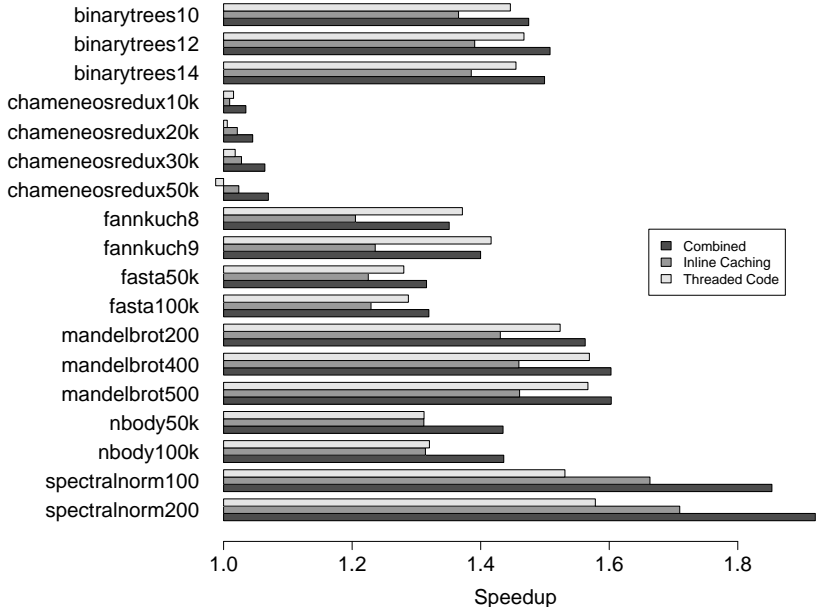- requires instruction set extension ➦ generator
- increased instruction cache requirements
  ➦ desktop/server vs. mobile devices

Inline Caching
meets
Quickening

Stefan
Brunthaler

Motivation
Basic Technique
Quickening
**Evaluation**
Conclusion

# Benchmarks



Speedup

# Take away message

*Efficient* inline caching
is possible *without* dynamic trans-
lation—and *worthwhile*.