

Avoiding Monomorphisation Bottlenecks with Phase-based Splitting

Sophie Kaleba
University of Kent
Canterbury, United Kingdom
S.Kaleba@kent.ac.uk

Stefan Marr
University of Kent
Canterbury, United Kingdom
S.Marr@kent.ac.uk

Richard Jones
University of Kent
Canterbury, United Kingdom
R.E.Jones@kent.ac.uk

Programming languages, such as JavaScript, Ruby or Python, rely on a managed runtime to reach state-of-the-art performance (see respectively V8 [2], JRuby [4], PyPy [7]). Such runtime systems apply aggressive optimisations based on speculative assumptions: one common assumption is that the behaviour of a program remains mostly homogeneous at run time. However, literature [5, 6] on the run-time behaviour of programs shows that programs contain several distinct *phases* (i.e. intervals of time exhibiting a distinct and relatively homogeneous behaviour) and that these phases may also repeat. For instance, the backend of a scientific computation web server can be viewed as such a program: first it analyses the user request, then computes statistics on received data and lastly outputs a plot.

In this talk, we discuss how programming language implementations overlook phases in program behaviour, and we investigate how they could benefit from recognising and adapting to behaviour changes between phases.

We focus on the performance of polymorphic calls and more concretely on two optimisations: lookup caches and call-site method splitting. *Lookup caches* [3] are built upon the assumption that a program has low variability: this kind of cache is stored at call-sites to amortise the cost of method lookup and assumes that only a limited number of call targets will be needed at run time. *Splitting*, first implemented in SELF [3], aims at monomorphising polymorphic call-sites. It copies methods such that their lookup caches are more likely to contain only a single call target. This optimisation helps reduce the dispatch cost, as well as further guide inlining decisions. However, splitting has a memory overhead and may increase compilation time.

Often during the initialisation of a program, lookup caches are filled with call targets that will not be needed in the later phases of the program. This hampers the dispatch performance [1], and possibly prevents inlining. We explored this problem with two micro-benchmarks and test them on TruffleSOM on top of GraalVM and JavaScript on top of NodeJS/V8. In these benchmarks, some call-sites can experience two types of phases: a polymorphic phase or a monomorphic phase. A call-site is considered *phase-sensitive* if its lookup cache is likely to contain non-needed targets during one of these phases, e.g. typically if a polymorphic phase is followed by a monomorphic phase. We measure the impact of phase-sensitivity for both method call-sites and closure application sites.

Preliminary results show that lookup caches are indeed impacted by phased behaviour. On TruffleSOM, a program having a method call-site lookup cache containing non-needed entries is 31% slower than the same program having this call-site being purely monomorphic. Same applies for closure application sites, where having a lookup cache containing several non-needed entries leads to a 14.3%

slowdown. On NodeJS, we obtained respectively a 40% and a 19% slowdown. These results suggest that phase-based monomorphisation based on splitting can increase performance.

We argue that using data on the high-level behaviour of a program, especially phases, could guide splitting heuristics to improve performance. We introduce a proof-of-concept prototype in TruffleSOM that provides a minimal set of source code annotations to identify phase-sensitive call-sites candidates for splitting and to identify phase switches. Concretely, the first time a phase switch occurs, the phase-sensitive call-site is split so that two versions of the method(s) co-exist: one for the polymorphic phase and one for the monomorphic phase.

We apply this prototype to our two micro-benchmarks to respectively split a call-site and a closure application site that are both phase-sensitive. The results are promising: splitting the phase-sensitive call-site leads to an average 18.5% speedup compared to baseline, i.e. a version of the program without annotation. Furthermore, if we focus on the performance at the phase granularity, we reach speedups ranging from 37% to 47.6% during monomorphic phases where the call-site lookup cache used to be filled with non-needed targets. Likewise, splitting the phase-sensitive closure application site leads to a speedup: 10% on average compared to baseline, with speedups ranging from 21% to 23.2% when looked at phase granularity. In addition, this particular use-case highlights an interesting splitting issue: in this benchmark, a function parameter takes a closure, which is later activated. This code structure is common in modern frameworks, causing closure application sites to be polymorphic, without splitting heuristics as in GraalVM being able to resolve it.

These experiments are a starting point to show that phased behaviour is still overlooked, but suggest that a phase-aware lightweight system guiding dynamic optimisations could offer significant performance improvements. For instance, this approach could help reduce the overhead of frameworks and the indirection they introduce by monomorphising performance-critical code. In the last part of this talk, we will describe what could be the future of this project in terms of implementation choices. We will also discuss other dynamic optimisations that could benefit from phase awareness.

REFERENCES

- [1] Guido Chari, Diego Garbervetsky, and Stefan Marr. 2017. A metaobject protocol for optimizing application-specific run-time variability. In *Proceedings of the 12th Workshop on Implementation, Compilation, Optimization of Object-Oriented Languages, Programs and Systems*. 1–5.
- [2] Google. 2008. Google V8 code source repository. <https://github.com/v8/v8>

- [3] Urs Hölzle, Craig Chambers, and David Ungar. 1991. Optimizing dynamically-typed object-oriented languages with polymorphic inline caches. In *European Conference on Object-Oriented Programming*. Springer, 21–38.
- [4] JRuby. 2006. JRuby: The Ruby programming language on the JVM. <http://jruby.org/>
- [5] Priya Nagpurkar. 2007. Analysis, Detection, and Exploitation of Phase Behavior in Java Programs. (2007), 217.
- [6] E. Perelman, M. Polito, J.-Y. Bouguet, J. Sampson, B. Calder, and C. Dulong. 2006. Detecting phases in parallel applications on shared memory architectures. In *Proceedings 20th IEEE International Parallel & Distributed Processing Symposium*. IEEE, Rhodes Island, Greece, 10 pp. <https://doi.org/10.1109/IPDPS.2006.1639325>
- [7] Armin Rigo and Samuele Pedroni. 2006. PyPy’s approach to virtual machine construction. In *Companion to the 21st ACM SIGPLAN symposium on Object-oriented programming systems, languages, and applications (OOPSLA ’06)*. Association for Computing Machinery, Portland, Oregon, USA, 944–953. <https://doi.org/10.1145/1176617.1176753>