

Introduction

Sorting is ubiquitous in computer science: it is used as a subroutine in many algorithms, and sort functions are included in the standard library of every modern programming language. It is not surprising, then, that sorting algorithms and implementations are an area of active research. However, nearly all sorting research operates under the assumption that comparisons have a unit cost: all that matters is the number of comparisons made. Python’s sort is especially characteristic of this trend: the built-in sort algorithm, Timsort, employs sophisticated optimizations based on assumptions about the distribution of real-world data to try to reduce the number of comparisons necessary to sort a list, taking up nearly half the code in `Objects/listobject.c`.

The unit-cost comparison assumption is reasonable in low-level models of computation, where comparisons only cost a few machine operations, and issues such as cache utilization are much more important than the cost of individual comparisons. However, in dynamic languages like Python, comparisons are extremely expensive, which makes the unit cost assumption tenuous. In this project, it is demonstrated that it is possible to *vastly* increase the speed of Python’s sort routine, not by reducing the number of comparisons made, but by reducing the cost of the individual comparisons, by exploiting the homogeneous structure of real-world data.

The dynamic dispatch problem

Python is a dynamic language. This means that functions must check the types of their arguments each time they are called, and execute different implementations for different argument types. The procedure by which the interpreter selects the correct implementation of a polymorphic function given arbitrary argument types is called *dynamic dispatch* [3].

The comparison operator in Python, like every other operator, is polymorphic, and is thus subject to the cost of dynamic dispatch each time it is applied. In general, this is unavoidable, and part of the inherent cost for which one gains the many benefits a dynamic language offers. However, in the common situation where one can show (at runtime) that the objects being compared will have the same type across function applications, having to pay the needless cost of dynamic dispatch is unfortunate, and one wonders whether it can be avoided altogether.

Inline caching

One way in which dynamic language implementations, specifically Just-in-Time (JIT) compilers, have sought to avoid dynamic dispatch is called *inline caching* [2]. The observation is the following: in practice, the input types at any given polymorphic call site will often remain constant across calls. As long as the homogeneity assumption holds, directly calling a cached implementation is then much faster than looking it up every time. This is a gamble: if the assumption fails, time was spent caching for no benefit, yielding a net loss. However, because the assumption nearly always holds, the average-case savings dwarf the overhead. *The success of inline caching demonstrates that by making a homogeneity assumption, and applying that assumption to run-time observations, it is possible avoid a great deal of the overhead of dynamic dispatch.*

The list homogeneity assumption

It is thus clear that to avoid the overhead of dynamic dispatch, we have to make assumptions about the lists we’re going to be sorting. We will operate under a weak assumption: that lists are homogeneous with respect to type. This assumption is motivated by two considerations:

- Lists are idiomatically homogeneous data structures. To quote the Python documentation (emphasis added) [1]:
“*Lists are mutable, and their elements are usually homogeneous and are accessed by iterating over the list.*”
- It is difficult to define comparison in the heterogeneous context. While some types are comparable to each other, like `longs` and `floats`, and of course subtypes of the same type are often comparable, such pairs are rare in practice: it’s “like comparing apples and oranges”. Truly heterogeneous lists simply don’t have well-defined orderings.

Exploiting homogeneity

The list homogeneity assumption is powerful, because it allows us to justify a pre-sort check: instead of checking types and other properties every time we compare, and then executing the appropriate compare function, we can check *once* for each key, and select the compare function ahead of time! Of course, this is only possible if the keys are homogeneous with respect to whatever properties determine their compare functions, but that is *precisely* our assumption. We can simply replace calls to `PyObject_RichCompareBool()` with calls to a function pointer, which can be set to point to whichever (potentially optimized) compare implementation our pre-sort check indicates is appropriate. This turns $O(n \log n)$ sort-time checks into $O(n)$ pre-sort checks.

There are, of course, many specific properties we could check for, since each built-in type has its own special considerations. For example, strings can have varying character widths, but it’s reasonable to assume lists of strings are homogeneous with respect to character width (since most applications use ASCII), so that check can be moved out of the sort loop. The optimizations used in this project are detailed in the red boxes in the flow chart to the right.

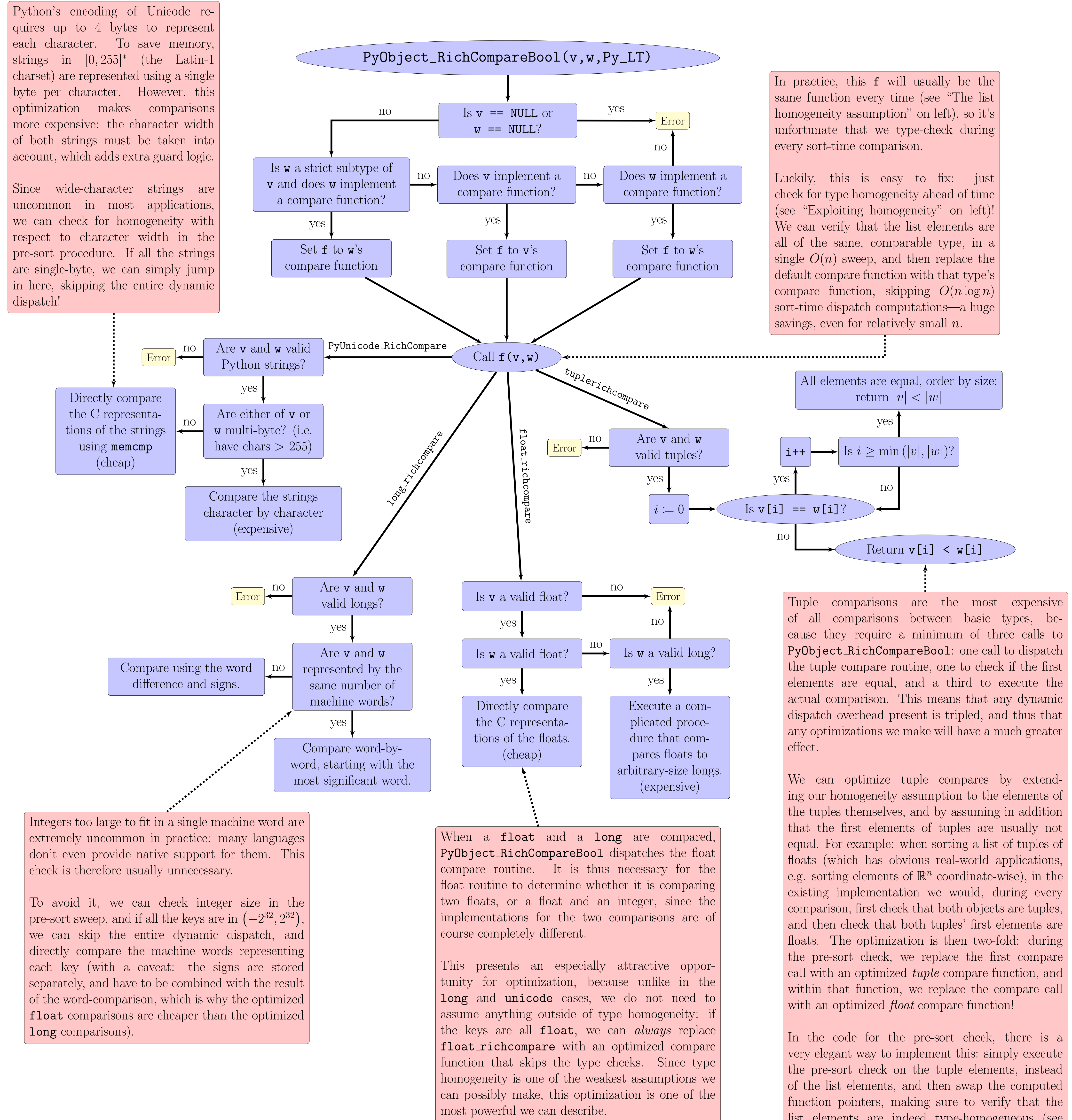
Adding Data-Aware Sort Optimizations to CPython

“Comparisons are extraordinarily expensive in Python, precisely because of the maze of test-and-branch code it requires just to figure out which bottom-level comparison function to invoke each time. That’s why I spent months of my life (overall) devising a sequence of sorting algorithms for Python that reduced the number of comparisons needed.”

— Tim Peters, lead Python developer and author of the Timsort algorithm [5]

Exploring the maze

The flow chart below presents a simplified version of the “maze” of dynamic dispatch logic necessary to compare two Python objects. The blue and yellow boxes correspond to blocks of code in the CPython implementation. The red boxes point out places where type- or other homogeneity found in practice can be exploited to avoid a great deal of that logic during the sort.



Benchmarking methodology

Benchmarks were performed according to recommendations in [7] and [4]:

- Build Python with `./configure --enable-optimizations`
- Isolate CPU core with `isolcpus=cpu.list` and `rcu.nocbs=cpu.list`
- Disable frequency scaling with `python -m perf system tune`
- Serialize with `CPUID` and time with `RDTSCP`

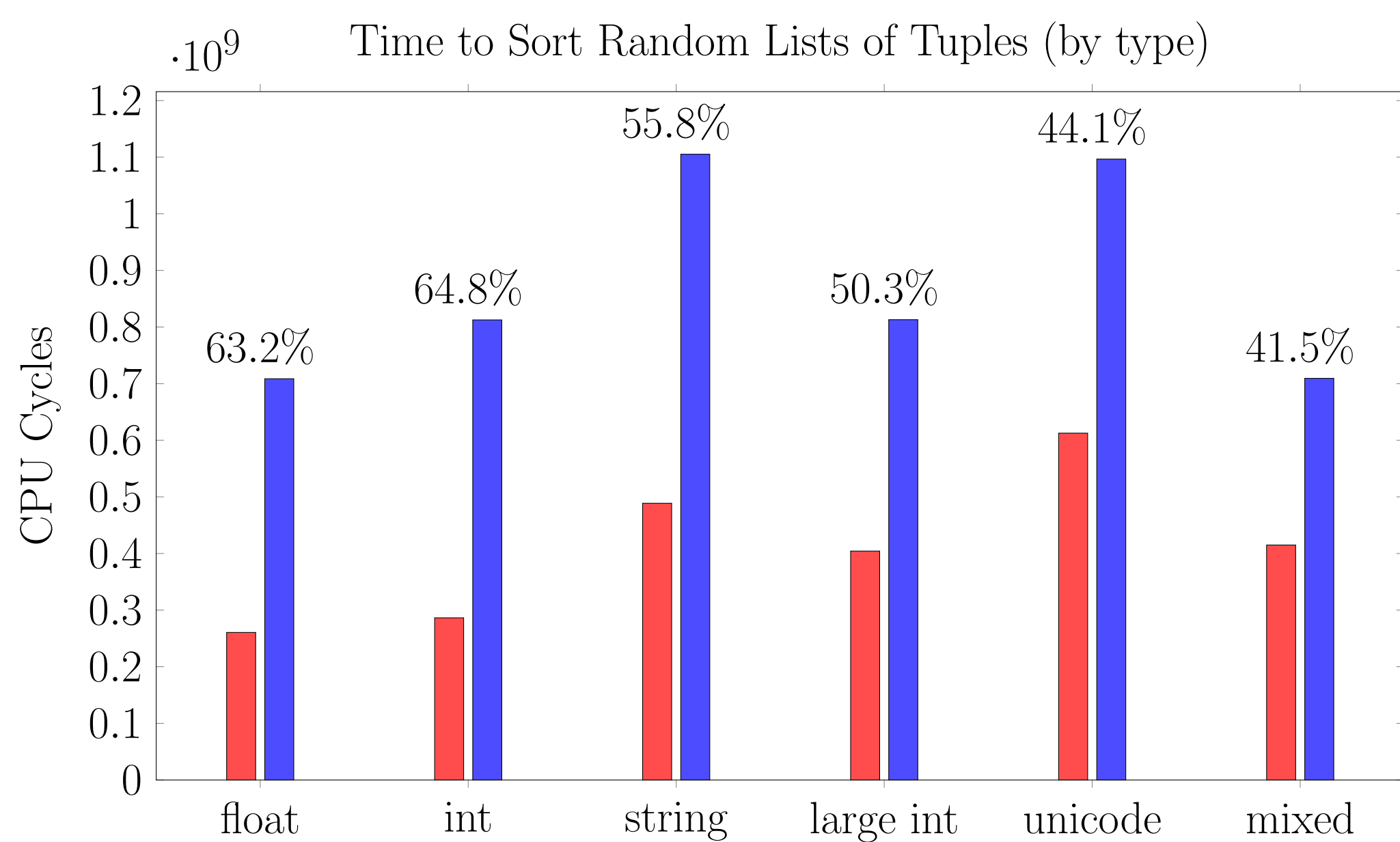
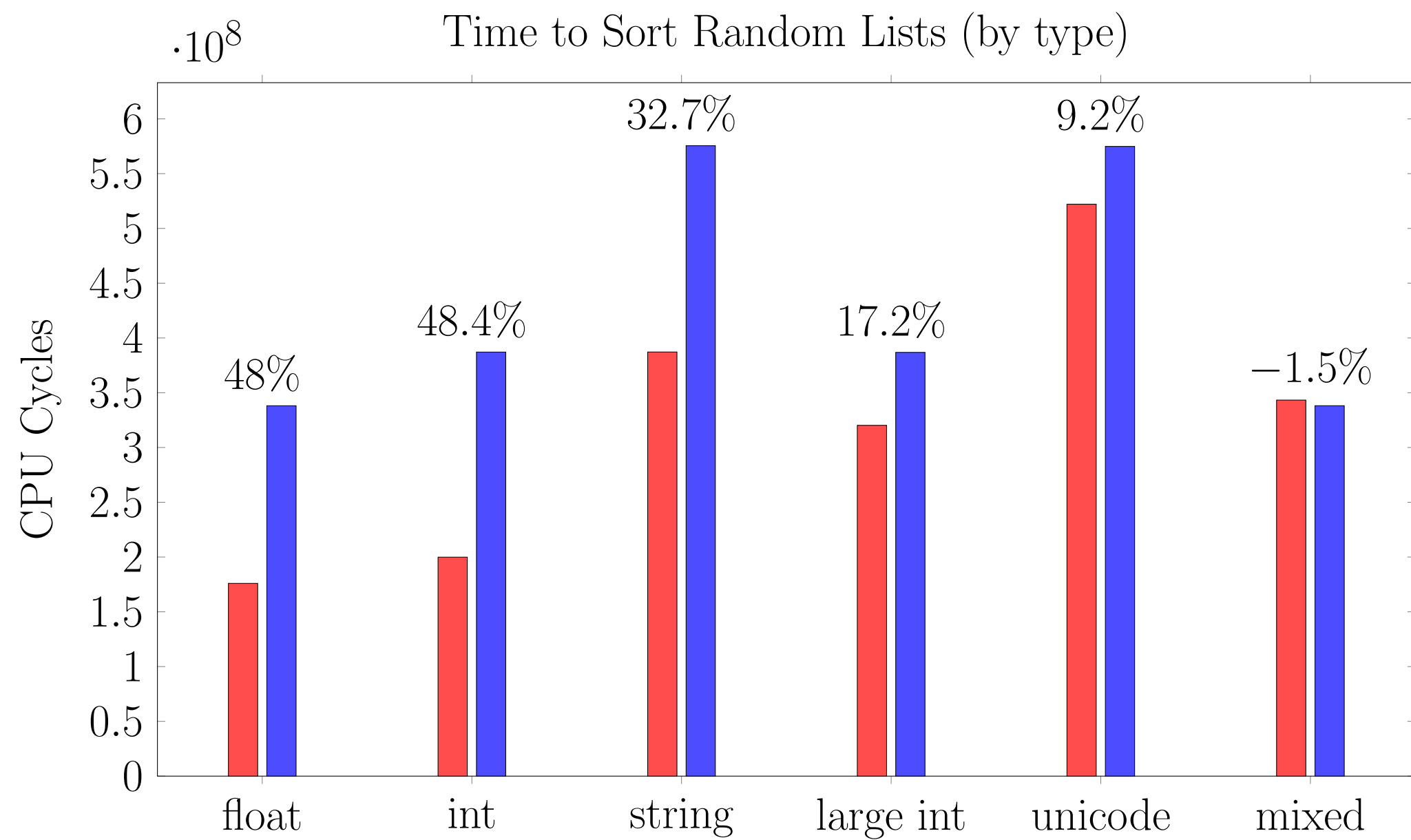
12 different benchmarks were run: sorting random lists of six different types, as well as random lists of tuples of those types. The benchmark lists were constructed using the comprehensions below:

- float: `[random() for _ in range(0,n)]`
- int: `[int((1<<31)*random() - (1<<30)) for _ in range(0,n)]`
- large int: `int_list(n) + [1<<64]`
- string: `[str(random()) for _ in range(0,n)]`
- unicode: `string_list(n) + ["\uffff"]`
- mixed (float + int): `float_list(n) + [0]`
- tuples of `type`: `[(x,) for x in type_list(n)]`

The complete code for the benchmarks can be found in the supplement.

Results

The graphs below give the times it took to run each of the benchmarks on the current CPython codebase (blue) and on the patched codebase (red). Numbers above the bars are percent differences. Confidence intervals could not be shown on the graphs because standard deviation over 1000 iterations was less than 0.3% of the mean for all the benchmarks.



Note that in all the cases commonly encountered in practice (float, int, and string), the improvements are greater than 30%, which compares favorably to numbers commonly reported in the literature, e.g. [6]. Further, note that in the uncommon worst-case (mixed), where the entire pre-sort check is executed without any ability to perform optimization (because the homogeneity assumption is violated), the cost of overhead is less than 1.5%, an order of magnitude smaller than the gain for which we make this near-certain gamble.

Conclusion

In this project, it was shown that by focusing on reducing the cost of the comparison operation, instead of trying to reduce the total number of comparisons, and by making weak assumptions about the distribution of real-world data, it is possible to dramatically reduce the cost of CPython’s sort routine, by bypassing nearly all the overhead due to needless dynamic dispatch.

References

- [1] *Data Structures*. docs.python.org/3.3/tutorial/datastructures.html. Python manual page.
- [2] L. Peter Deutsch and Allan M. Schiffman. “Efficient Implementation of the Smalltalk-80 System”. In: *Proceedings of the 11th ACM SIGACT-SIGPLAN Symposium on Principles of Programming Languages*. POPL ’84. New York, NY, USA: ACM, 1984, pp. 297–302.
- [3] Scott Milton and Heinz W. Schmidt. *Dynamic Dispatch in Object-Oriented Languages*. Tech. rep. CSIRO – Division of Information Technology, 1994.
- [4] Gabriele Paoloni. *How to Benchmark Code Execution Times on Intel IA-32 and IA-64 Instruction Set Architectures*. Tech. rep. Intel, 2010.
- [5] Tim Peters. Post on the python-dev mailing list. URL: mail.python.org/pipermail/python-dev/2016-October/146648.html.
- [6] Tim Peters. *An Adaptive, Stable, Natural Mergesort*. Included with the CPython source at `Objects/listsort.txt`. 2002.
- [7] Victor Stinner. *My journey to a stable benchmark*. hypo-notes.readthedocs.io/microbenchmark.html. 2016.