

# Seq: A High-Performance Language for Bioinformatics

ANONYMOUS AUTHOR(S)\*

The scope and scale of biological data is increasing at an exponential rate, as technologies like next-generation sequencing are becoming radically cheaper and more prevalent. Over the last two decades, the cost of sequencing a genome has dropped from \$100 million to nearly \$100—a factor of over  $10^6$ —and the amount of data to be analyzed has increased proportionally. Yet, as Moore’s Law continues to slow, computational biologists can no longer rely on computing hardware to compensate for the ever-increasing size of biological datasets. In a field where many researchers are primarily focused on biological analysis over computational optimization, the unfortunate solution to this problem is often to simply buy larger and faster machines.

Here, we introduce **Seq**, the first language tailored specifically to bioinformatics, which marries the ease and productivity of Python with C-like performance. Seq is a subset of Python—and in many cases a drop-in replacement—yet also incorporates novel bioinformatics- and computational genomics-oriented data types, language constructs and optimizations. Seq enables users to write high-level, Pythonic code without having to worry about low-level or domain-specific optimizations, and allows for seamless expression of the algorithms, idioms and patterns found in many genomics or bioinformatics applications. On equivalent CPython code, Seq attains a performance improvement of up to two orders of magnitude, and a  $175\times$  improvement once domain-specific language features and optimizations are used. With parallelism, we demonstrate up to a  $650\times$  improvement. Compared to optimized C++ code, which is already difficult for most biologists to produce, Seq frequently attains up to a  $2\times$  improvement, and with shorter, cleaner code. Thus, Seq opens the door to an age of democratization of highly-optimized bioinformatics software.

Additional Key Words and Phrases: computational biology, bioinformatics, programming language, domain-specific language, Python, optimization

## 1 INTRODUCTION

DNA sequencing technologies have revolutionized life sciences and clinical medicine [Mardis 2017]. Today, state-of-the-art cancer treatment involves sequencing a tumor genome for diagnosis and treatment [Kamps et al. 2017]. Improvements in sequencing hardware have drastically reduced the laboratory cost of sequencing; thus, the biggest bottleneck and cost today in the sequence analysis pipeline is computational data analysis [Muir et al. 2016]. Indeed, reductions in sequencing costs over the past two decades have radically outpaced both Moore’s and Kryder’s Laws. Sequencing a genome in 2008 cost more than \$10 million, but over the past several years the cost has fallen below the celebrated \$1,000 mark, and is expected to soon surpass \$100 per genome [Hayden 2014]. Unfortunately, the computational data analysis cost inherent in sequencing pipelines has not improved at the same rate and has already surpassed the cost of sequencing back in 2011 [Sboner et al. 2011]. In fact, it is expected to soon be cheaper to simply re-sequence an individual than to even store their raw sequencing data, let alone analyze it [Weymann et al. 2017]. By the same token, the number and size of sequencing datasets have and continue to grow exponentially, which will require better algorithms and software especially as the gap between Moore’s Law and sequencing data growth continues to widen.

Recent advances in next-generation sequencing (NGS) technologies are continuing this revolution at a rapid rate, providing a means to study various biological processes through a genomic lens. Because of its novel capabilities and vast scale, NGS has even bigger

---

2019. 2475-1421/2019/1-ART1 \$15.00  
<https://doi.org/>

50 computational needs, as terabytes of data need to be processed and analyzed with the  
51 aid of various novel computational methods and tools [Mardis 2017]. These tools (e.g. [Li  
52 et al. 2009a], [Li and Durbin 2009], [Zaharia et al. 2011], [McKenna et al. 2010], [Yorukoglu  
53 et al. 2016], [Shajii et al. 2018]) are used on a daily basis in research laboratories and have  
54 fueled major discoveries such as establishing mutation-disease links [Manolio et al. 2008]  
55 and detection of recent segmental duplications in the genome [Bailey et al. 2001].

56 Despite these advances, however, many contemporary genomic pipelines cannot scale with  
57 the ever-increasing deluge of sequencing data, although some have suggested exploiting the  
58 intrinsic structure of biological data towards this end [Berger et al. 2016]. This inability to  
59 scale has necessitated impractical and expensive *ad hoc* solutions such as frequent hardware  
60 upgrades and constant (re-)implementation of underlying software. Many promising methods  
61 are also too difficult to use and replicate [Peng 2011] because they are often manually tuned  
62 for a single dataset, further fueling the recent replication crisis [Baker 2016]. Finally, many  
63 tools are not maintained due to a lack of personnel, expertise and funds. This fact has led to  
64 many abandoned code repositories that cannot be easily modified to suit researchers' needs,  
65 although being, in theory, vastly superior to the well-maintained alternatives.

66 The root cause of these problems lies in the general-purpose languages that are used for  
67 bioinformatics software development. The most popular programming languages, Python,  
68 R, and occasionally C++, are not designed to efficiently handle and optimize for sequencing  
69 data workflows. Even so, researchers often use high-level languages like Python or R to  
70 analyze NGS data since they allow for quick and easy expression of high-level ideas, despite  
71 a steep performance penalty. Alternatively, a researcher may manually implement low-  
72 level optimizations in a language like C. However, this task requires a considerable time  
73 investment and often results in hard-to-maintain codebases riddled with subtle bugs and  
74 tied to a particular architecture, especially in a field like computational biology where many  
75 researchers are not software engineers by trade. These issues are further exacerbated as  
76 the field shifts towards the use of third-generation portable sequencers that are powered by  
77 resource-limited devices [Lu et al. 2016], which warrant entirely different software designs  
78 and optimizations.

79 In this paper, we introduce *Seq*, a domain-specific language (DSL) and compiler designed  
80 to provide productivity *and* high performance for computational biology. *Seq* is a subset of  
81 Python, and therefore provides Python-level productivity; yet, the compiler can generate  
82 efficient code because the language is statically-typed with compile-time support for Python's  
83 duck typing. *Seq* provides data types tailored to computational genomics and uses domain-  
84 specific information to optimize code. Our DSL allows computational biology experts to  
85 quickly prototype and experiment with new algorithms as they would in Python, without  
86 imposing the burden of learning a new language. Further, *Seq* is designed to hide all low-  
87 level, complex code optimizations from the end user. Unlike libraries, the *Seq* compiler can  
88 perform optimizations such as operator fusion, which are necessary in high-performance  
89 computational biology applications.

90 This paper makes the following contributions:

- 91
- 92 • We introduce *Seq*, the first domain-specific language (DSL) and compiler for computa-  
93 tional biology.
- 94 • We introduce novel genomic-specific data types (e.g. sequence and *k*-mer types) and  
95 operators (e.g. for reverse complementation, *k*-merization, etc.), further augmented  
96 with additional language constructs such as compiler-optimized pipelines and genomic  
97
- 98

99 matching, to both simplify the algorithmic descriptions of complex problems and to  
100 enable domain-specific optimizations.

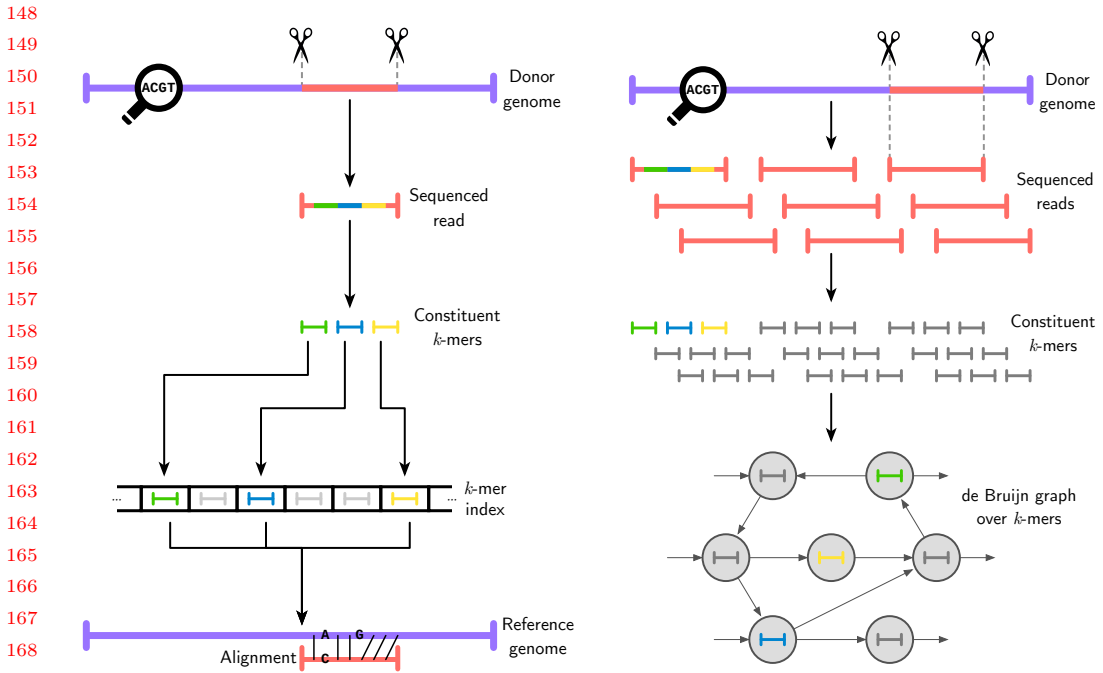
- 101 • We design a statically-typed subset of Python tailored for bioinformatics applications,  
102 which operates entirely without expensive runtime type information and provides  
103 performance comparable to (and in many cases better than) C's.
- 104 • We demonstrate how to use parallelism, prefetching, pipelining and coroutines to  
105 drastically improve the performance of Seq.
- 106 • We provide an implementation of the compiler and standard library.
- 107 • We show that many important and widely-used NGS algorithms can be made 70–  
108 175× faster than their Python counterparts as well as 2× faster than the existing  
109 hand-optimized C++ implementations.

110 The rest of the paper is organized as follows. Section 2 provides a primer on computational  
111 genomics. Section 3 gives an example of the Seq language, followed by a description of  
112 the language design and implementation in Section 4 and domain-specific optimizations in  
113 Section 5. Section 6 evaluates the language, Section 7 describes related work, and the paper  
114 concludes in Section 8.  
115

## 116 2 A PRIMER ON COMPUTATIONAL GENOMICS 117

118 The fundamental data type in computational genomics is the *sequence*, which is conceptually  
119 a string over  $\Sigma = \{A, C, G, T\}$ , representing the four nucleotides (also called *bases*) that  
120 comprise DNA. Sequences come in several different forms, with varying properties such  
121 as length, error profile and metadata. For example, *genome sequencing*—a process that  
122 determines the DNA content of a given biological sample—typically produces *reads*: DNA  
123 sequences roughly 100 bases in length, with a substitution error rate less than 1% and  
124 metadata consisting of a unique identifier and a string of *quality scores*, indicating the  
125 sequencing machine's confidence in each reported base of the read. Reads are often analyzed  
126 in the context of a *reference genome*, a much longer (in the case of human, 3 gigabase-length)  
127 sequence that represents the consensus sequence of an organism's genome in its entirety. A  
128 standard first step in nearly any sequence analysis pipeline is *sequence alignment*, which  
129 is the process of identifying the position in the reference sequence to which a particular  
130 read aligns with the smallest edit distance (although many different formulations of this  
131 problem exist, such as finding *all* alignments under a given edit distance threshold). To this  
132 end, reads are typically first split into fixed length- $k$  contiguous subsequences called *k-mers*,  
133 which are then queried in an index of  $k$ -mers from the reference to guide the alignment  
134 process, as shown in Figure 1a. The index itself is an abstract data type that maps  $k$ -mers to  
135 positions (also called *loci*) in the reference at which they appear, and is often implemented  
136 in practice as a hash table or FM-index [Ferragina and Manzini 2004; Li and Homer 2010].

137 Due to the large memory footprints of these structures (roughly 5 gigabytes for optimized  
138 FM-indices and tens of gigabytes for hash tables) given the size of the genome, coupled with  
139 their poor cache performance, many alignment algorithms spend a significant fraction of their  
140 time time stalled on memory accesses; the fraction of stalled cycles in these applications can  
141 be over 70% depending on the input dataset [Appuswamy et al. 2018]. Once a candidate locus  
142 is found via the index (and possibly after several filtering steps), a full dynamic programming  
143 alignment is performed, usually via a variant of the Smith-Waterman algorithm. Because  
144 dynamic programming alignment is a key kernel in nearly all alignment algorithms, there  
145 has been substantial research into designing hand-optimized implementations that exploit  
146 SIMD vectorization for better performance [Farrar 2006; Suzuki and Kasahara 2018; Šošić  
147



(a) Overview of the alignment process for sequencing data. A sequencing machine produces a *read*: a roughly 100 base pair DNA sequence randomly sampled from the donor’s genome. Most alignment algorithms then split this read into  $k$ -mers—fixed length- $k$  subsequences—and query these  $k$ -mers in an index of  $k$ -mers from the reference genome to determine candidate alignment positions. Finally, full dynamic programming alignment (typically via an adapted Smith-Waterman algorithm) is carried out to produce the final alignment.

(b) Overview of *de novo* genome assembly from sequencing data. Sequenced reads are partitioned into constituent  $k$ -mers, which are then taken to be nodes in a de Bruijn graph whose edges represent  $(k-1)$ -length overlaps. Other formulations use  $(k-1)$ -mers (two for each original  $k$ -mer) as nodes with the original  $k$ -mers represented by the edges. The assembled sequence corresponds to an Eulerian path on this graph.

Fig. 1. Visualizations of two standard computational genomics applications.

and Šikić 2017]. One additional complication in sequence alignment is that, while half of all reads will align in the so-called *forward direction* (i.e. without modification), the other half will only align in the *reverse* direction, meaning the read must be *reverse complemented* before alignment. Reverse complementation of a sequence is an operation where the sequence is reversed, and A-bases are swapped with T-bases while C-bases are swapped with G-bases (and vice versa). The fact that half the reads are reverse complemented with respect to the reference genome is a biproduct of the double-stranded nature of DNA, and ultimately leads to reverse complementation being a very common operation that is done on sequences.

Alongside alignment, another common application in computational genomics is *de novo* assembly, where the reads are used to “reconstruct” the donor genome, in the absence of a predefined reference sequence. While several approaches to this problem exist, perhaps the

197 most common is to again partition the reads into  $k$ -mers, build a de Bruijn graph whose  
198 vertices are these  $k$ -mers with edges indicating that a given  $k$ -mer overlaps with another,  
199 and finally to find an Eulerian path through this graph, which would encode the assembled  
200 sequence [Khan et al. 2018]. An overview of this process can be seen in Figure 1b. As in  
201 alignment, there are several additional steps involved in practice, such as counting and  
202 filtering  $k$ -mers, as well as error correction (as assembly is more sensitive to errors than  
203 alignment) [Simpson and Durbin 2012].

204 Looking further downstream in the genomic analysis pipeline, computational biologists  
205 employ a slew of techniques to handle the problems at hand. However, virtually any  
206 downstream model or algorithm, regardless of its domain (machine learning, graph algorithms,  
207 etc.), is built on top of the sequence manipulation building blocks described above. For  
208 example, structural variation detection (the discovery of novel genomic rearrangements) starts  
209 by analyzing read alignment irregularities to detect potential breakpoints of a rearrangement,  
210 and proceeds by correcting those alignments via more advanced read alignment schemes that  
211 utilize  $k$ -mers and FM-indices. Many other problems, such as mutation calling, gene copy  
212 number variation detection, genome-wide association studies and cancer driver identification,  
213 proceed in a similar fashion. Thus, the common threads between alignment, assembly and  
214 many other applications in the genomics domain are the data types used (i.e. various types of  
215 sequences like reads, reference or fixed-length  $k$ -mers) and the low-level operations performed  
216 on them (i.e. some form of matching, indexing, splitting sequences into subsequences or  
217  $k$ -mers, reverse complementation, etc.). However, these operations are often embedded in  
218 vastly different higher-level algorithms; compare, for example, the dynamic programming  
219 involved in alignment and the de Bruijn graph path finding involved in assembly. For this  
220 reason, we chose to expose these genomics-specific types and operations in a comparatively  
221 lower-level language than many other DSLs, as we discuss in detail below.

222  
223

### 224 3 SEQUENCE $k$ -MERIZATION AND SEEDING – AN EXAMPLE

225 Seq provides built-in language-level facilities for seamlessly expressing many of the types  
226 and design patterns found in genomics applications. As an example, consider reading a set  
227 of sequencing reads from a FASTQ file (a standard format for storing reads) and querying  
228 each read's constituent  $k$ -mers in a genomic index. This process is commonly referred to  
229 as *seeding*, and is the first step in nearly any sequence alignment algorithm [Li and Homer  
230 2010].

231 An implementation of 20-mer seeding in Seq is shown in Figure 2. Seq uses the familiar  
232 syntax of Python, but incorporates several genomics-specific features and optimizations.  
233  $k$ -mer types like `k20` (which represents a  $k$ -mer with 20 bases), for example, allow for easy  
234  $k$ -merization (the process of splitting a sequence into  $k$ -mers, done using `kmers[k20]` in  
235 Seq) and reverse complementation (using `~kmer`). Similarly, pipelining—a natural model  
236 for thinking about processing reads—is easily expressible in Seq, where a user can define  
237 pipelines via the `|>` operator as shown in the figure. Seq can also speed up expensive index  
238 queries via pipeline transformations that allow for effective software prefetching (`prefetch`  
239 keyword). Compare this Seq implementation to the C++ implementation also shown in  
240 Figure 2, which includes extensive boilerplate code for reverse complementation and FASTQ  
241 iteration, and cannot perform the domain-specific pipeline or encoding optimizations made  
242 by the Seq compiler, which in practice we find to attain upwards of 1.5–2× speedups over  
243 optimized C++ implementations.

244  
245

```

246
247
248
249
250 from sys import argv
251 from genomeindex import *
252
253 # index and process 20-mers
254 def process(kmer: k20,
255            index: GenomeIndex[k20]):
256     prefetch index[kmer], index[~kmer]
257     hits_fwd = index[kmer]
258     hits_rev = index[~kmer]
259     ...
260
261 # index over 20-mers
262 index = GenomeIndex[k20](argv[2])
263
264 # stride for k-merization
265 stride = 10
266
267 # sequence-processing pipeline
268 (fastq(argv[1])
269  |> kmers[k20](stride)
270  |> process(index))

```

```

#include <iostream>
#include <fstream>
#include <string>
#include <cstdlib>
#include "GenomeIndex.h"

char revcomp(char base) {
    switch (base) {
        case 'A': return 'T';
        case 'C': return 'G';
        case 'G': return 'C';
        case 'T': return 'A';
        default: return base;
    }
}

void revcomp(char *kmer, int k) {
    for (int i = 0; i < k/2; i++) {
        char a = revcomp(kmer[i]);
        char b = revcomp(kmer[k - i - 1]);
        kmer[i] = b;
        kmer[k - i - 1] = a;
    }
}

void process(char *kmer, int k,
             GenomeIndex &index) {
    auto hits_fwd = index[kmer];
    revcomp(kmer, k);
    auto hits_rev = index[kmer];
    revcomp(kmer, k); // undo
    ...
}

int main(int argc, char *argv[]) {
    const int k = 20;
    const int stride = 10;
    auto *index = GenomeIndex(argv[1], k);
    std::ifstream fin(argv[2]);
    std::string read;
    long line = -1;
    while (std::getline(fin, read)) {
        line++;
        // skip over non-sequences in FASTQ
        if (line % 4 != 1) continue;
        auto *buf = (char *)read.c_str();
        int len = read.size();
        for (int i = 0; i + k <= len; i += stride)
            process(kmer, k, index);
    }
}

```

Fig. 2. Example  $k$ -merization and seeding application in Seq and C++.

#### 4 LANGUAGE DESIGN & IMPLEMENTATION

A critical barrier to any new language’s success in a particular field is its initial adoption, as most potential users already have a set of languages, environments and packages with which they are comfortable. This is particularly true in bioinformatics, where many researchers are biologists first and programmers second. For this reason, the Seq language borrows the syntax and semantics of Python—one of the most widely-used languages in bioinformatics—and adds several genomics-oriented language features and constructs. Indeed, most preexisting Python code will compile and run without modification in Seq, ultimately allowing the user to attain the performance of C/C++ with the programming ease of Python.

To achieve this, we designed a compiler with a static type system. It performs Python-style duck typing and runtime type checking at compile-time, completely eliminating the substantial runtime overhead imposed by the reference Python implementation, CPython, and most other Python implementations alike. Unlike these, we reimplemented all of Python’s language features and built-in facilities from the ground up, completely independent of the CPython runtime. The Seq compiler uses an LLVM [Lattner and Adve 2004] backend, and in general uses LLVM as a framework for performing general-purpose optimizations. Seq programs additionally use a lightweight (<200 LOC) runtime library for I/O and memory allocation; for the latter, CPython’s reference counting is replaced with the Boehm garbage collector [Boehm and Weiser 1988], a widely-used conservative GC that is a drop-in replacement for malloc.

(“Python” is henceforth used as a synonym for “CPython” unless otherwise specified.)

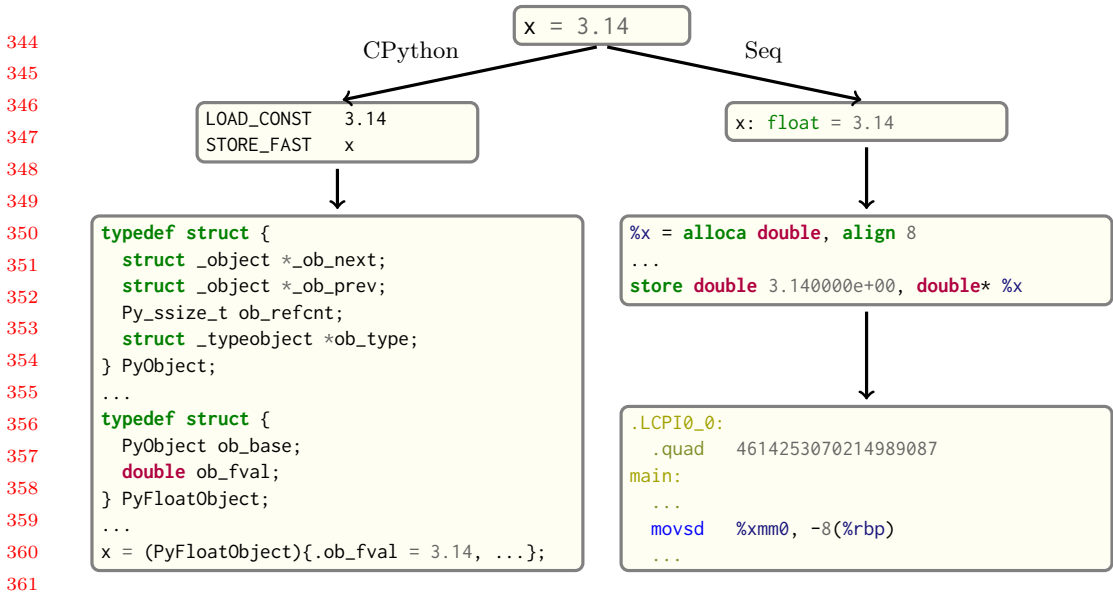


## 4.1 Building a statically-typed Python

Python is an interpreted language, and does not check for type consistency until necessary during runtime. Moreover, Python does not perform strong type checking in the formal sense: for any method call, the Python runtime simply queries the given object's method dictionary (vtable) for the requested method and, if found, dispatches the call there. The same applies to operators, getters, setters and accessors, which are nothing more than syntactic sugar around specially-named *magic methods*. For example, the "+" operator is represented by the `__add__` magic method, while the list lookup `v[i]` is represented as `v.__getitem__(i)`. Importantly, the actual type of `v` is irrelevant (be it a list or anything else), so long as `v` defines `__getitem__`. This approach to typing is commonly referred to as *duck typing*. Further, even variable resolution in Python is realized as simple dictionary lookups (e.g. via the `locals()` dictionary) that are maintained for each scope. On top of all this, these tables are updated during runtime, which allows end users to modify objects' vtables during execution, making it possible to dynamically alter the behavior of a given object in virtually any fashion. This simple and clean design, together with a well-thought-out syntax, enables rapid prototyping and a great deal of flexibility without imposing artificial language design constraints, which is partly what has made Python popular in many different domains, bioinformatics notwithstanding.

However, this dynamism comes with a hefty price in terms of performance, as almost any method invocation or variable reference requires expensive dictionary lookups during runtime. Furthermore, the lack of any type annotations and the dynamic nature of objects necessitates delaying type checks until the given object is actually used, which can sometimes take place *days* after the Python script was initially run in the case of long-running programs (this is a common problem in bioinformatics, where many scripts take a long time to complete due to ever-growing input datasets, whereby rapid initial development is paid for by a slow debugging cycle). Moreover, this lazy approach to typing requires the developer to include numerous manual type checks and large test suites to ensure type soundness during execution. (Python versions 3.6 and later ship with the `mypy` type checker, which supports type annotations with ahead-of-time type checks. However, these annotations are optional and type inconsistencies are treated as warnings, meaning any type annotation is more of a guideline than a strict rule. We did, however, borrow `mypy`'s type syntax.) In most domains, this is a fair price to pay for Python's ease and expressibility as compared to the alternatives. However, the lack of performance is significantly problematic in the context of computational biology, where an average dataset is on the order of hundreds of gigabytes in size, and where even a simple loop construct incurs enough overhead to render Python code hundreds of times slower than its C counterpart. Highly optimized Python implementations, such as PyPy or Numba, do not sufficiently address these problems as they are either bound to the same constraints as the original (CPython) implementation (i.e. dynamic runtime and lazy duck typing), or limited in scope solely to numerical types.

Despite its array of dynamic and runtime-oriented features, the full flexibility provided by Python is not commonly used in many domains. While type flexibility and dynamic object/type modifications are, for example, crucial for rapid web development, they are almost completely absent in high-performance scientific applications, and arguably even slow down the development cycle of such applications. For these reasons, we designed a strongly-typed alternative to Python's runtime, which captures the subset of its dynamic features that is commonly used in the field of computational biology, and that can be resolved



362 Fig. 3. Seq versus CPython during compilation and execution of a simple float assignment. CPython  
363 compiles to bytecode that omits all type information, and instead relies on runtime type information  
364 by virtue of metadata stored alongside the actual float value within the PyFloatObject structure. By  
365 contrast, Seq infers the type of `x` at compile-time and compiles the assignment to LLVM IR, which  
366 encodes type information. LLVM in turn compiles this to assembly or machine code.

367 *at compile time*. In doing so, we trade some largely unneeded dynamism for greatly improved  
368 performance, which by contrast *is* gravely needed in the field.

370 **Basic types.** Python has a relatively simple type system in which all types derive from the  
371 object base type. Some primitive types (such as integers and floats) are, for performance  
372 reasons, implemented directly in C within CPython’s runtime. However, even the C imple-  
373 mentations of these types carry a significant overhead, as they still have to interoperate  
374 nicely with the rest of the Python ecosystem. As can be seen in Figure 3, a simple float  
375 object—arguably the most lightweight type Python has—consists of three pointers, an integer  
376 and finally the float value itself. This “metadata” is necessary for Python’s runtime type  
377 resolution and reference counting (we do note, however, that the `_ob_next` and `_ob_prev`  
378 pointers are compiled into the structure definition conditionally, and can be omitted). For  
379 these reasons, all high-performance Python libraries (such as NumPy) achieve their speed  
380 by dealing primarily with arrays or matrices that can be abstracted away from the Python  
381 runtime to the C level.

382 Seq follows a different design philosophy in terms of types. Primitive types such as `int`, `bool`  
383 and `float` map directly to the equivalent LLVM IR types `i64`, `i8` and `double`, respectively.  
384 As such, they incur no overhead whatsoever. Nevertheless, each of these primitives is still  
385 logically a fully-fledged type with a set of associated methods that can be extended by the  
386 user (e.g. type `int` has a method `__add__` for addition that can be statically patched); there  
387 is no overhead as all method dispatches are resolved by the compiler. Furthermore, Seq  
388 inlines all magic method invocations on primitive types (e.g. an `int.__add__` call is compiled  
389 to a single LLVM `add` instruction).

390 More complex types typically compile to an LLVM aggregate type or a pointer to one.  
391 Aggregate types are used in place of Python’s tuples and named tuples (represented in

392



Seq by a new type construct), and are fully isomorphic to C structs. Pointers to aggregate types—or, more precisely, *reference types*—are used to implement classes (represented in Seq by a class construct). As usual, aggregates are passed by value while reference types are passed by, unsurprisingly, reference. For example, the following Seq expressions have types mapping to the indicated LLVM types:

Seq expression	LLVM IR type	Description
"hello world"	{i64, i8*}	struct of length and character pointer
(1, 0.5, False)	{i64, double, i8}	struct of tuple element types
MyClass()	i8*	pointer to heap-allocated MyClass struct
MyClass().foo	{i8*, void (i8*)*}	struct of self and method function pointer

where the last example assumes `foo` is defined to be a method of `MyClass` that takes no extra arguments and does not return a value (i.e. `def foo(self: MyClass) -> void`).

In order to maintain compatibility with Python, class members can be deduced automatically by lexically analyzing a given class's methods. Python's built-in collection types—`list`, `set` and `dict`—are all modeled as reference types in Seq and bootstrapped as standard library classes implemented in Seq itself.

**Generic functions, methods and types.** Python's lack of static typing allows any function to take objects of any type as an argument. This design philosophy does not translate well to strongly-typed compiled languages that do not use any form of runtime type information, as they typically require each function to explicitly specify input and output types.

Code compatibility with Python is of paramount importance for Seq, as it is unreasonable to expect users to manually annotate (or rewrite) their large codebases. Thus, Seq handles this problem by treating each Python function that does not provide type annotations as a *generic function*, where one or more input or output types cannot be deduced from annotations or a lexical analysis of the function body. In this case, each argument without a type annotation (referred to as an *implicit generic*) is replaced by a concrete type on demand at compile-time. For example, on encountering `f(42)` as in Figure 5, the compiler checks whether there is an instantiation of `f` that accepts an `int` argument, and if so routes the call there. If not, the compiler clones `f`'s AST and creates a new instantiation of the function that accepts specifically an `int` argument. (Note that, as far as LLVM is concerned, instantiations are just simple functions with unique mangled names.) This newly created function would produce an error if, for example, `int` did not contain an appropriate `__mul__` method as required in the function body. Instantiations are created lazily on demand.

Unlike Python, Seq allows users to explicitly mark functions as generic and to specify explicit generic type parameters, allowing more complex type relationships to be expressed. For example, Figure 4 shows a higher-order function that only operates on generic nodes and functions (as `T` is an explicit generic type parameter). The argument types of `item` ensure that the argument function can take the argument node's data as a parameter. Note that it is impossible for Python-style unnamed generics to cover this use-case without explicit `isinstance` checks. As shown in Figure 4, explicit type parameterization is optional even when explicit generics are present, as Seq performs type parameter deduction whenever possible.

Analogous reasoning applies to classes, where class members can be generic. Examples of such classes include `list[T]` and `dict[Key,Value]`. Unlike functions, implicit generics are disallowed in classes as they would impair readability and could lead to ambiguous instantiations during the class member deduction stage. Note that, as far as Seq is concerned, different

```

442 class Node[T]:
443     next: Node[T]
444     data: T
445
446 def item[T,U](n: Node[T], f: function[T,U]) -> list[U]:
447     return [f(n.data)]
448
449 n = Node(None, 5)
450 def foo(x: int) -> str:
451     return str(x)
452 i = item(n, foo) # type parameters deduced as int and str
453 i = item[int,str](n, foo) # explicit specification also OK

```

Fig. 4. Seq’s explicit generic type parameters.

instantiations of functions and classes are treated as different types. Thus,  $f(x: \text{list}[\text{int}])$  and  $f(x: \text{list}[\text{float}])$  are represented internally as two separate functions, which allows Seq to optimize each instantiation according to its concrete argument types (albeit by sacrificing any kind of polymorphism in the current implementation).

*Duck typing.* Seq’s type system is designed to behave like Python’s if one uses Seq as a drop-in Python replacement without specifying explicit types. As long as the methods of every type are known at compile time (an invariant strictly enforced by Seq as it does not allow modifying methods during runtime), the compiler will deduce the argument/return types of all methods and instantiate any generic method as appropriate. Indeed, we find that this static instantiation-on-demand simulates duck typing reasonably well. Explicit type annotations enforce an extra layer of typing discipline on top of duck typing (à la mypy), and as such coexist peacefully with it.

*Type inference.* Any strongly typed language needs a way to infer the type of each variable present in a given program. Languages such as C or Pascal require end users to manually annotate each variable with a type. Other languages, such as C++11 or newer versions of Java, support uni-directional type inference by automatically deducing types of left-hand side terms based on right-hand side types. Initial versions of Seq also used uni-directional type inference, allowing users to say, for instance,  $x = 5$  instead of  $x: \text{int} = 5$ .

However, uni-directional type inference is not able to seamlessly handle a couple of common constructs in the Python language, including empty lists (e.g.  $a = []$ ), nullables (e.g.  $a = \text{None}$ ) and lambda functions (e.g.  $\text{lambda } x: x+1$ ). With uni-directional inference, each of these constructs requires the user to provide manual type annotations (e.g.  $a: \text{list}[\text{int}] = []$ ) even if the type can be inferred later. Because of this, Seq employs a Hindley-Milner-based bi-directional type inference algorithm to automatically annotate such types<sup>1</sup>. We slightly modified the standard Hindley-Milner algorithm to support generic classes, functions and instantiations on demand. We also enforce an invariant where all types within a scope (be it a function scope, class scope or the top-level scope) must be fully deduced by the end of that scope. This implies that a function cannot return a non-instantiated generic type:  $\text{def } f(): \text{return } []$ , for example, will cause a compilation error, but  $\text{def } f[T]() -> \text{list}[T]: \text{return } []$  will compile successfully. Any weakly typed variable or lambda is instantiated as soon as possible (note that Seq treats lambdas as

<sup>1</sup>This is a recent addition to Seq, and is currently still in the testing stage. At the time of writing, Seq’s master branch uses uni-directional type deduction with added support for generics and type-less nullables. Note that this uni-directional version does not support lambdas.

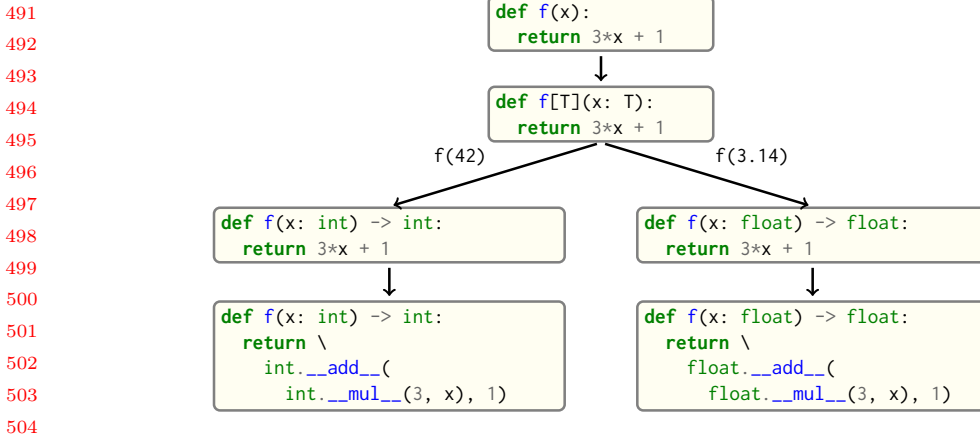


Fig. 5. Seq’s implicit generic type parameters. The function  $f$  is declared to take a parameter  $x$  of unspecified type; the Seq compiler treats the type of  $x$  as generic and clones  $f$  on demand for each new input type, and subsequently deduces return types.

weakly typed constructs and does not generalize them—generalizations are only applied to generic functions defined with `def` and generic classes).

**Limitations.** The strongly-typed nature of Seq does come with some limitations compared to conventional Python. Since all types must be fixed at compile-time, a Seq program cannot (for example) create a collection of elements (e.g. `list`) with varying types. Seq’s tuples are also less versatile than Python’s: they cannot be iterated over if they contain different types, and a list cannot be cast to a tuple easily, as tuple sizes must be known at compile time. Seq also does not support method or class monkey-patching at runtime (but it does support this at compile time), nor indexing into a heterogeneous tuple with a non-constant index (as the type of the resulting expression would be ambiguous). Our type checker and instantiation algorithm also require each function to have a single return type. Finally, while Seq supports class extensions, it does not support subtyping (nor, therefore, fully-fledged polymorphism), meaning that `class A; class B(A)` will copy `A`’s methods to `B` without making `B` a subtype of `A` per se. With these trade-offs, Seq can perform all type-checking at compile time without sacrificing any runtime cycles for type enforcement, and without significantly hindering the expressibility of Python’s syntax. We have found that, especially in bioinformatics software, these language capabilities are seldom required (or at least can almost always be replaced by Seq-conforming alternatives with minimal effort). Consequently, these features are omitted in Seq at the time of writing. A brief list of differences between Seq and Python can be found in Appendix A.

## 4.2 Coroutines and generators

Generators—Python’s answer to streams and lazy data structures—are an integral part of most Python/Seq programs: even simple for-loops are realized as an iteration over a generator. While it would certainly be possible to implement generators as they are in CPython (heap-allocated `generator` objects that expose a `__next__` method for obtaining the next generated value), this would incur a substantial overhead, given how frequently generators are used.

Instead, Seq employs LLVM coroutines (also used by Clang versions 6 and later to implement the C++ Coroutine TS [Nishanov 2017]). The advantage of this approach is that,

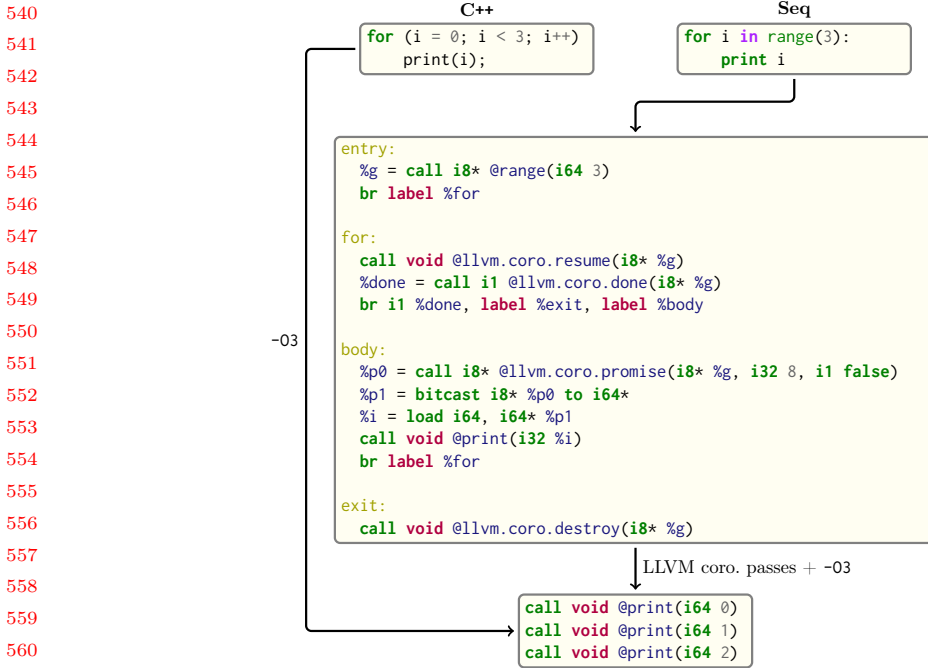


Fig. 6. Compilation of Seq generators. Two semantically identical loops in C++ and Seq are shown in the uppermost boxes. Seq generators are implemented as LLVM coroutines, iteration over which in LLVM IR is shown in the middle box. The LLVM coroutine passes subsequently deduce that the “range” coroutine is created and destroyed in the same function without escaping, and inline/unroll the coroutine to produce code identical to the C++ example’s.

when a generator is created and destroyed in the same function without escaping (by far the most common case in Python, similar to the for-loop example), LLVM’s coroutine passes are able to optimize out all the associated coroutine overhead, like coroutine frame allocation etc. Thereby, a typical Seq for-loop ultimately compiles to identical LLVM IR as the same loop expressed in C or C++, as shown in Figure 6.

### 4.3 Additional and genomics-specific language features

*Sequence and k-mer types.* Seq’s namesake type is indeed the sequence type: `seq`. A `seq` object represents a DNA sequence of any length and—on top of general-purpose string functionality—provides methods for performing common sequence operations such as splitting into subsequences, reverse complementation and *k*-mer extraction. Alongside the `seq` type are *k*-mer types, where e.g. `k1` represents a 1-mer, `k2` a 2-mer and so on, up to `k256` (a reasonable upper bound on *k*-mer length in nearly any genomics application). Evidently, a key difference between the `seq` type

```
dna = s'ACGTACGTACGT' # sequence literal

# (a) split into subsequences of length 3
#       with a stride of 2
for sub in dna.split(3, 2):
    print sub

# (b) split into 5-mers with stride 1
for kmer in dna.kmers[k5](1):
    print kmer
    print ~kmer # reverse complement

# (c) convert entire sequence to 12-mer
kmer = k12(dna)
```

Fig. 7. Example of `seq` and *k*-mer type usage.

589 and  $k$ -mer types is that, for the latter, the length is intrinsic to the type itself, much like  
 590 how common integer types like `int8`, `int32`, etc. have a fixed bit-length.

591 Sequences can be seamlessly converted between these various types, as shown in Figure 7.  
 592 In fact, this pattern is prevalent in many genomics applications, where longer sequences  
 593 (be it a read, reference or anything else) are split into their constituent  $k$ -mers, and each is  
 594 subsequently processed.

595 *Pipelines and partial calls.* Pipelining is a natural  
 596 model for thinking about processing genomic  
 597 data, as sequences are typically processed in  
 598 stages (e.g. read from input file  $\rightarrow$  split into  
 599  $k$ -mers  $\rightarrow$  query  $k$ -mers in index  $\rightarrow$  perform full  
 600 dynamic programming alignment  $\rightarrow$  output re-  
 601 sults to file), and are almost always independent  
 602 of one another as far as this processing is con-  
 603 cerned. Because of this, Seq supports a pipe oper-  
 604 ator: `|>`, similar to F#'s pipe and R's `magrittr`  
 605 (`%>%`) [Bache and Wickham 2014]. Pipeline stages  
 606 in Seq can be regular functions or generators. In  
 607 the case of standard functions, the function is  
 608 simply applied to the input data and the result  
 609 is carried to the remainder of the pipeline, akin  
 610 to F#'s functional piping. If, on the other hand, a stage is a generator, the values yielded by  
 611 the generator are passed lazily to the remainder of the pipeline, which in many ways mirrors  
 612 how piping is implemented in Bash. Note that Seq ensures that generator pipelines do not  
 613 collect any data unless explicitly requested, thus allowing the processing of terabytes of data  
 614 in a streaming fashion with no memory and minimal CPU overhead.

615 An example of pipeline usage is shown in Figure 8, which shows the same two loops from  
 616 Figure 7, but as pipelines. First, note that `split` is a Seq standard library function that takes  
 617 three arguments: the sequence to split, the subsequence length and the stride; `split(..., 3,`  
 618 `2)` is a partial call of `split` that produces a new single-argument function  $f$  where  $f(x) =$   
 619 `split(x, 3, 2)`. The undefined argument(s) in a partial call can be implicit, as in the  
 620 second example: `kmers` (also a standard library function) is a generic function parameterized  
 621 by the target  $k$ -mer type and takes as arguments the sequence to  $k$ -merize and the stride;  
 622 since just one of the two arguments is provided, the first is implicitly replaced by `...` to  
 623 produce a partial call (i.e. the expression is equivalent to `kmers[k5](..., 1)`). Both `split`  
 624 and `kmers` are themselves generators that yield subsequences and  $k$ -mers respectively, which  
 625 are passed sequentially to the last stage of the enclosing pipeline in the two examples.

627 *Pattern matching.* Seq provides the conventional `match`  
 628 construct, which works on integers, lists, strings and tuples.  
 629 An example usage of `match` is shown in Figure 9. A novel  
 630 aspect of Seq's `match` statement is that it also works on  
 631 sequences, and allows for concise recursive representations  
 632 of several sequence operations such as subsequence search,  
 633 reverse complementation tests and base counting, which  
 634 are shown in Figure 10. Sequence patterns consist of literal  
 635 ACGT characters, single-base wildcards (`_`) or "zero or more"  
 636 wildcards (`...`) that match zero or more of any base.

637

```

dna = s'ACGTACGTACGT' # sequence literal

# (a) split into subsequences of length 3
#     with a stride of 2
dna |> split(..., 3, 2) |> echo

# (b) split into 5-mers with stride 1
def f(kmer):
  print kmer
  print ~kmer

dna |> kmers[k5](1) |> f
  
```

Fig. 8. Example of pipeline usage in Seq, where the two loops from Figure 7 are represented as pipelines.

```

def describe(n: int):
  match n:
    case m if m < 0:
      print 'negative'
    case 0:
      print 'zero'
    case m if 0 < m < 10:
      print 'small'
    default:
      print 'large'
  
```

Fig. 9. Example usage of `match`.

```

638 # (a)
639 def has_spaced_acgt(s: seq) -> bool:
640     match s:
641         case s'A_C_G_T...':
642             return True
643         case t if len(t) >= 8:
644             return has_spaced_acgt(s[1:])
645         default:
646             return False
647
648 # (b)
649 def is_own_revcomp(s: seq) -> bool:
650     match s:
651         case s'A...T' or s'T...A' or s'C...G' or s'G...C':
652             return is_own_revcomp(s[1:-1])
653         case s'':
654             return True
655         default:
656             return False
657
658 # (c)
659 type BaseCount(A: int, C: int, G: int, T: int):
660     def __add__(self: BaseCount, other: BaseCount):
661         a1, c1, g1, t1 = self
662         a2, c2, g2, t2 = other
663         return (a1 + a2, c1 + c2, g1 + g2, t1 + t2)
664
665 def count_bases(s: seq) -> BaseCount:
666     match s:
667         case s'A...': return count(s[1:]) + (1,0,0,0)
668         case s'C...': return count(s[1:]) + (0,1,0,0)
669         case s'G...': return count(s[1:]) + (0,0,1,0)
670         case s'T...': return count(s[1:]) + (0,0,0,1)
671         default: return BaseCount(0,0,0,0)

```

Fig. 10. Example usages of match on sequences in Seq. Example (a) checks if a given sequence contains the subsequence A\_C\_G\_T, where \_ is a wildcard base; such an operation may be present in an application that uses *spaced seeds*—non-contiguous  $k$ -mers that are shown to improve accuracy in some settings [Kucherov et al. 2015]. Example (b) checks if the given sequence is its own reverse complement, which is useful in certain sequence hashing schemes [Ondov et al. 2016]. Finally, example (c) counts how many times each base appears in the given sequence, which can e.g. be used to determine GC content (the fraction of bases that are G or C) [Šmarda et al. 2014].

*External functions.* Seq enables seamless interoperability with C and C++ via `cdef` functions, as shown in Figure 11. Primitive types like `int`, `float`, `bool` etc. are directly interoperable with the corresponding types in C/C++, while compound types like tuples are interoperable with the corresponding struct types. Other built-in types like `str` provide methods to convert to C analogs, such as `c_str()` as shown in Figure 11.

```

cdef sqrt(float) -> float
cdef puts(ptr[byte])
print sqrt(100.0)
puts("hello world".c_str())

```

Fig. 11. Example of `cdef` function usage in Seq with C standard library functions `sqrt` and `puts`.



687 *Type extensions.* Seq provides an `extend` key-  
 688 word that allows programmers to add and  
 689 modify methods of various types at compile  
 690 time, including built-in types like `int`  
 691 or `str`. This allows much of the function-  
 692 ality of built-in types to be implemented  
 693 in Seq as type extensions in the standard  
 694 library. Figure 12 shows an example where  
 695 the `int` type is extended to include a `to`  
 696 method that generates integers in a speci-  
 697 fied range, as well as to override the `__mul__`  
 698 magic method to “intercept” integer multi-  
 699 plications. Note that all type extensions are  
 700 performed strictly at compile-time and incur  
 701 no runtime overhead.

702

## 703 5 OPTIMIZATIONS

### 704 5.1 Sequence encoding

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

```

extend int:
    def to(self: int, other: int):
        for i in range(self, other + 1):
            yield i

    def __mul__(self: int, other: int):
        print 'caught int mul!'
        return 42

for i in (5).to(10):
    print i # 5, 6, ..., 10

# prints 'caught int mul!' then '42'
print 2 * 3

```

Fig. 12. Example of type extension in Seq.

The first and most straightforward optimization made by Seq is to 2-bit encode  $k$ -mer objects, as is commonly done in practice in performance-critical applications. In particular, we map  $k$ -mer types to the LLVM IR type `iN` where  $N = 2k$ . This has the advantage of allowing  $k$ -mers up to  $k = 32$  to fit into a single machine word on 64-bit architectures.

In order to support fast  $k$ -mer operations, we also conditionally compile various lookup tables into any Seq program that requires them:

- For reverse complementation, we create a complete 4-mer reverse complement lookup table, implemented as a global length- $4^4$  array indexed by encoded 4-mers, storing the encoded reverse complement of the given 4-mer at each index (hence, this array requires 256 bytes). Then, using the property  $\overline{s_1} \parallel s_2 = \overline{s_2} \parallel \overline{s_1}$  (where  $\overline{s}$  denotes the reverse complement of a sequence  $s$  and  $\parallel$  denotes concatenation), we can construct the reverse complement of an arbitrary  $k$ -mer by partitioning it into 4-mers and concatenating (i.e. shifting and bitwise-ORing) their reverse complements in reverse, each of which is given by the lookup table. For  $k < 4$ , or the remainder of a longer  $k$ -mer whose length is not divisible by 4, we can simply pad with A-bases to obtain a 4-mer then remove the corresponding T-bases in the reverse complemented 4-mer (recall the reverse complement of A is T). Another noteworthy aspect of this scheme is that, if we choose our encoding wisely, we get reversal for free as well. Specifically, if we 2-bit encode base  $b \in \{A, C, G, T\}$  as  $f(b)$  so that  $f(A) = \sim f(T)$  and  $f(C) = \sim f(G)$  (where  $\sim$  is bitwise-NOT), then we can undo the “complementation” component of the reverse complement to obtain the original sequence in reverse by applying a simple bit inversion.
- For converting general sequences into  $k$ -mers, we compile a second lookup table that maps the ASCII characters A, C, G and T to their 2-bit encoded values, implemented also as a length-256 array. The encoding process then iteratively looks up each base in this array to construct the encoded  $k$ -mer.
- For converting  $k$ -mers back to sequences, we use a simple length-4 array that maps the 2-bit encodings back to ASCII. We note, however, that this is a far less common conversion than the previous one.

## 5.2 Parallelism

CPython and many other implementations alike cannot take advantage of parallelism due to the infamous global interpreter lock, a mutex that protects accesses to Python objects, preventing multiple threads from executing Python bytecode at once. Unlike CPython, Seq has no such restriction and supports full multithreading. To this end, Seq supports a *parallel* pipe operator `||>`, which is semantically similar to the standard pipe operator except that it allows the elements sent through it to be processed in parallel by the remainder of the pipeline. Hence, turning a serial program into a parallel one often requires the addition of just a single character in Seq, as shown by Figure 13. Further, a single pipeline can contain multiple parallel pipes, resulting in nested parallelism.

Internally, the Seq compiler uses Tapir [Schardl et al. 2017] with an OpenMP task backend to generate code for parallel pipelines. Logically, parallel pipe operators are similar to parallel-for loops: the portion of the pipeline after the parallel pipe is outlined into a new function that is called by the OpenMP runtime task spawning routines (as in `#pragma omp task` in C++), and a synchronization point (`#pragma omp taskwait`) is added after the outlined segment. Lastly, the entire program is implicitly placed in an OpenMP parallel region (`#pragma omp parallel`) that is guarded by a “single” directive (`#pragma omp single`) so that the serial portions are still executed by one thread (this is required by OpenMP as tasks must be bound to an enclosing parallel region).

## 5.3 Software prefetching for faster genomic index lookups

Large genomic indices—ranging from several to tens or even hundreds of gigabytes—used in many applications in the field result in extremely poor cache performance and, ultimately, a substantial fraction of stalled memory-bound cycles [Appuswamy et al. 2018; Wang et al. 2012; Zhang et al. 2007]. For this reason, Seq performs pipeline optimizations to enable data prefetching and to hide memory latencies, an idea that has also been explored in previous work [Chen et al. 2007; Kiriansky et al. 2018]. The programmer must provide just:

- a `__prefetch__` magic method definition in the index class, which is logically similar to `__getitem__` (indexing construct) but performs a prefetch instead of actually loading the requested value (and can simply delegate to `__prefetch__` methods of built-in types);
- a one-line prefetch hint indicating where a software prefetch should be performed, which can typically be just before the actual load.

An example is shown in Figure 14. First, prefetch statements themselves compile to explicit invocations of the `__prefetch__` method, and functions containing prefetch statements (such as `process` in the figure) are converted by the compiler into coroutines that yield after each prefetch. Then, pipelines containing such functions as stages are transformed into loops that dynamically schedule multiple invocations of the newly created coroutine, where once one invocation yields or terminates, another is resumed or created by the scheduler,

```
dna = s'ACGTACGTACGT' # sequence literal

# (a) split into subsequences of length 3
#       with a stride of 2
dna |> split(..., 3, 2) ||> echo

# (b) split into 5-mers with stride 1
def f(kmer):
    print kmer
    print ~kmer

dna |> kmers[k5](1) ||> f
```

Fig. 13. Example of parallel pipeline usage in Seq, where the two pipelines from Figure 8 are parallelized.

**\_\_prefetch\_\_ magic method**

```

785 class MyIndex: # abstract k-mer index
786     ...
787     def __getitem__(self: MyIndex, kmer: k20):
788         # standard __getitem__
789     def __prefetch__(self: MyIndex, kmer: k20):
790         # similar to __getitem__, but performs prefetch
791
792
793

```

**Function transformations**

```

794 def process(read: seq, index: MyIndex):
795     ...
796     for kmer in read.kmers[k20](step):
797         prefetch index[kmer], index[~kmer]
798         hits = index[kmer]
799         hits_rev = index[~kmer]
800     ...
801     return x
802
803

```

```

804 def process(read: seq, index: MyIndex):
805     ...
806     for kmer in read.kmers[k20](step):
807         index.__prefetch__(kmer)
808         index.__prefetch__(~kmer)
809         yield
810         hits = index[kmer]
811         hits_rev = index[~kmer]
812     ...
813     yield x
814
815

```

**Pipeline transformations**

```

816 FASTQ("reads.fq") # input reads
817 |> process(index) # index lookup
818 |> postprocess # output results
819
820
821
822
823
824
825
826
827

```

```

828 M = ... # num. concurrent tasks
829 N = 0 # next coroutine slot to fill
830 k = 0 # next coroutine to execute
831 states = array(generator[T])(M)
832
833 for read in FASTQ("reads.fq"):
834     if N < M:
835         states[N] = process(read, index)
836         N += 1
837     else:
838         while True:
839             g = states[k]; g.next()
840             if g.done():
841                 postprocess(g.promise())
842                 g.destroy()
843                 states[k] = process(read, index)
844                 break
845             k = (k + 1) % M
846
847 for i in range(N):
848     g = states[i]
849     if not g.done():
850         while not g.done(): g.next()
851         postprocess(g.promise())
852         g.destroy()
853
854
855
856
857
858
859
860
861
862
863

```

Fig. 14. Transformations performed by Seq to enable effective index prefetching. Colored segments under pipeline transformations indicate where the specific stages show up in the resulting code. FASTQ is the standard file format for storing sequencing reads.

834 respectively. The transformed pipeline in Figure 14, for example, has several noteworthy  
 835 components:

- 836 •  $M$  is the number of concurrent coroutines to be processed, which ideally should be  
 837 large enough to saturate the memory bandwidth of the processor as prefetches are  
 838 performed. In practice, we choose  $M$  conservatively to be 16, which also allows for  
 839 software prefetching performed by other parts of the system, such as the garbage  
 840 collector.
- 841 •  $N$  is a variable indicating how many of the  $M$  coroutine slots have been filled, and is  
 842 only used at the start of the loop to actually fill the slots.
- 843 •  $k$  is the next coroutine slot to be resumed by the loop.
- 844 •  $states$  is the array holding the  $M$  coroutine handles/frames (which have type `generator`  
 845 in Seq). In reality this array is stack-allocated in the entry block of the function  
 846 containing the pipeline. ( $T$  is simply the original return type of `process`.)

847 The code generated in the loop body is that of a simple dynamic scheduler where:  
 848

- 849 • The `if N < M` component initially populates the array of pending coroutines `states`.
- 850 • Inside the `else` clause is a loop that iterates cyclically through `states` and resumes  
 851 each coroutine. If a coroutine terminates (i.e. `if g.done()`), then the value returned  
 852 by the coroutine (given by `g.promise()`) is sent through the remainder of the pipeline,  
 853 as it would be in the original untransformed pipeline; then, the coroutine is destroyed  
 854 and a new one is created to take its place.
- 855 • The final loop simply completes any remaining coroutines that have not yet terminated.  
 856 Since the number of such coroutines is at most  $M$ , this loop just executes them  
 857 sequentially.

858 By employing this scheme, the latency of one coroutine’s cache miss can be overlapped  
 859 with useful work from another, increasing memory-level parallelism and overall throughput.  
 860

## 861 6 EVALUATION

862 We evaluated the performance of Seq on the following three benchmark suites, designed to  
 863 mimic both hypothetical and real-world genomics applications:

- 864 (1) *The Computer Language Benchmarks Game* suite [Gouy *n. d.*] restricted to DNA  
 865 benchmarks (3 benchmarks)
- 866 (2) Sequence manipulation suite, developed in-house (3 benchmarks)
- 867 (3) Genomic index queries (2 benchmarks)

869 We compared Seq with C++ (compiled with both GCC and Clang), Julia, Python 2.7 as well  
 870 as Shed Skin-translated [Dufour 2006] and Nuitka-compiled Python binaries. Other “compiled  
 871 Python” implementations, such as Numba, are geared towards numerical rather string or  
 872 DNA processing, and had issues efficiently compiling our benchmarks, or were abandoned.  
 873 All experiments were run on a dual-socket system with Intel Xeon X5690 CPUs (3.46 GHz)  
 874 with 6 cores each (totalling 12 cores and 24 hyper-threads) and 138GB DDR3-1333 RAM  
 875 with 12MB LLC per socket. C++ implementations were compiled with `-O3 -march=native`.  
 876 Julia was run with `--check-bounds=no -O3` parameters. Shed Skin was run with `-l -o`  
 877 optimizations, and Nuitka was run with `noasserts,no_warnings` options. Note that Seq  
 878 binaries (unlike C++ or Julia) do include bounds checks.

879 For the Benchmarks Game suite, we used the FASTA, RevComp and  $k$ -nucleotide mi-  
 880 crobenchmarks. Other benchmarks in this suite are not directly relevant to genomics or  
 881 bioinformatics in general, but we expect Seq’s performance on them to be on par with that  
 882

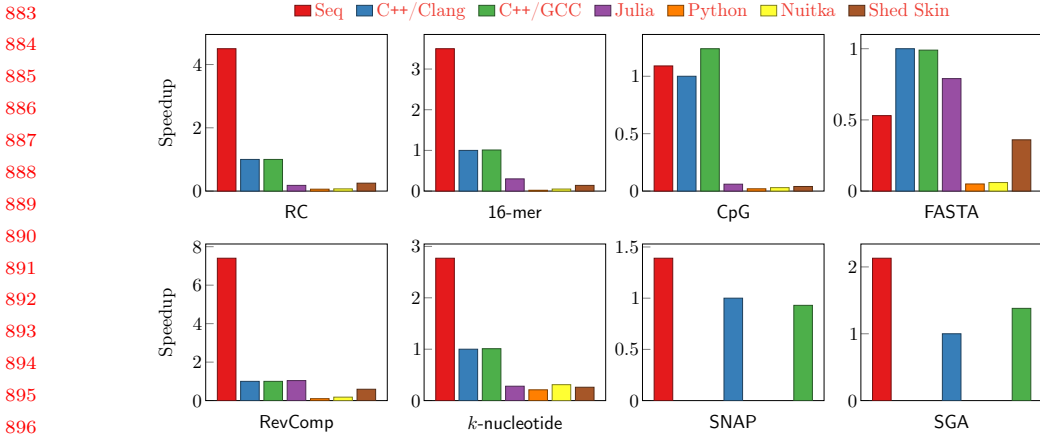


Fig. 15. Seq evaluation results on several genomics benchmarks, showing speedups over Clang. The Seq implementations used in these charts use Seq-specific types and constructs that are not available in Python. Note that only Seq, Clang and GCC were tested on the SNAP and SGA benchmarks. Seq performs at least as good as (and in many cases much better than) the C++ implementations in nearly every benchmark, excluding FASTA which, as we note in the text, is not as common a real-world application as the others.

of other LLVM-backed languages. Briefly, the FASTA benchmark entails generating random sequences in the FASTA format; RevComp entails reverse complementing a set of longer sequences, which is done through Seq’s domain-specific sequence type;  $k$ -nucleotide entails counting  $k$ -mers of various lengths, which we implemented using Seq’s  $k$ -mer types.

For the in-house suite, we designed three microbenchmarks that capture common genomics operations on a large set of reads:

- (1) RC: Output the reverse complement of each read, also implemented using sequence types. Unlike RevComp, this benchmark runs on millions of shorter reads rather than a few long reads.
- (2) 16-mer: Count the number of symmetric 16-mers in all reads (where a 16-mer is symmetric if its first half is identical to the reverse complement of its second half). The Seq implementation for this benchmark uses `match` on sequences, and is based on example (b) in Figure 10.
- (3) CpG: Count the number of CpG regions (i.e. regions that consist of C **and** G characters, such as CGC but *not* CCC or GGG as they lack a G and C, respectively) in all reads, and report the lengths of the shortest and longest CpG regions in the sample.

The final benchmark suite demonstrates the utility of Seq’s domain-specific genomic index query optimizations. Here, we use the genomic indices implemented in the widely-used tools SNAP [Zaharia et al. 2011] (a sequence alignment tool) and SGA [Simpson and Durbin 2012] (a *de novo* assembly tool). SNAP uses a hash table of 20-mers, which we re-implemented from scratch in Seq (see Appendix B). SGA, on the other hand, uses an FM-index, whose C++ implementation we wrapped in Seq. Both indices are used to query  $k$ -mers from our test dataset. For both SNAP and SGA, we compared the performance of our base Seq implementation, a Seq version that performs pipeline optimizations for index prefetching (code difference of one line) and a C++ implementation. Both of these benchmarks consume roughly 30GB of RAM.

	Seq (Py.)	Python	Nuitka	Shed Skin	Julia	Seq (Id.)	Speedup	
932	FASTA	111.3	99.4	15.6	<b>7.1</b>	10.7	0.7–10×	
933	RevComp	97.5	54.3	16.0	9.1	<b>1.3</b>	7–75×	
934	<i>k</i> -nucleotide	255.8	174.1	211.8	196.0	<b>19.7</b>	10–13×	
935	RC	136.1	2,417.3	1,887.2	534.6	764.7	<b>30.2</b>	25–80×
936	CpG	51.3	3,591.0	1,596.4	1,336.9	947.7	<b>51.3</b>	18–70×
937	16-mer	152.6	15,440.2	6,042.3	2,139.9	1,030.9	<b>87.0</b>	22–176×

938 Table 1. Seq runtime compared to Python, Nuitka, Shed Skin and Julia (seconds). “Seq (Py.)” (Pythonic  
 939 Seq) uses the same code as Python, whereas “Seq (Id.)” (idiomatic Seq) uses Seq-specific language  
 940 features and constructs.

941  
 942  
 943 Each benchmark was executed five times for each language/compiler, and the averages are  
 944 reported (with the exception of Julia, Python, Shed Skin and Nuitka, which were run three  
 945 times as they took orders of magnitude longer in some cases). The input dataset consisted  
 946 of 100 million 75bp DNA reads randomly chosen from the HG00123 sample [1000 Genomes  
 947 Project Consortium et al. 2010] (because SGA’s index is an order of magnitude slower than  
 948 SNAP’s, we downsampled our input dataset to 25 million reads for SGA). Results are shown  
 949 in Figure 15, where speedups over Clang-compiled C++ are given.

## 950 6.1 Improvements over Python

951  
 952 Seq programs can be written in one of two ways: in plain Python or using idiomatic Seq  
 953 (as in the implementations described above). The Pythonic implementations embody the  
 954 conventions set by the Python community, and code written in this way can be easily run  
 955 by both Python and Seq without modifications. The alternative style involves the use of  
 956 idiomatic Seq constructs and data types to manipulate genomic data, which are not available  
 957 in plain Python.

958 We compare the performance of Pythonic and idiomatic Seq implementations to that of  
 959 Python, Nuitka, Shed Skin and Julia in Table 1. All of the implementations in the second  
 960 benchmark suite (with the exception of idiomatic Seq) are line-by-line identical in terms  
 961 of the algorithm. For example, the *only* difference between the Pythonic-Seq and Python  
 962 implementations in the RC, CpG and 16-mer benchmarks is the lack of a `str.strip` call  
 963 when reading the input (which Seq’s I/O library does not require). Even by just directly  
 964 running Pythonic code, Seq is able to outperform Python by a factor of 15 to 100.

965 A similar pattern can be seen in the first benchmark suite, where Seq significantly  
 966 outperforms both the Python and Julia implementations. The only exception is the FASTA  
 967 benchmark, where Seq is slightly slower than Julia. While this could be further optimized,  
 968 we chose to keep the version that is most similar to Python. Additionally, we note that  
 969 the FASTA benchmark as specified by the Computer Language Benchmarks Game is not a  
 970 realistic application in genomics, as one would rarely be generating sequences rather than  
 971 reading them from a preexisting dataset.

972 Idiomatic versions further boost the improvement up to 175×, and showcase the impact  
 973 of individual domain-specific optimizations: RC and RevComp utilize Seq’s highly optimized  
 974 reverse complementation constructs; 16-mer showcases the gains—both in terms of readability  
 975 and performance—of sequence-based `match` statements; *k*-nucleotide shows the performance  
 976 improvement gained by using Seq’s native *k*-mer types. A few of the idiomatic versions also  
 977 rely on pipelining to perform further optimizations (described below).

978 Note that runtime becomes prohibitive as the number of reads to be processed increases—  
 979 while the performance of compiled Python (i.e. Nuitka, Shed Skin) and Julia is acceptable if  
 980



	Seq (Py.)	C++/Clang	C++/GCC	Seq (Id.)	Speedup
981	FASTA	<b>5.6</b>	5.7	10.7	0.5×
982	RevComp	9.5	9.5	<b>1.3</b>	7.3×
983	<i>k</i> -nucleotide	54.6	54.3	<b>19.7</b>	2.8×
984	RC	136.1	135.8	136.4	<b>30.2</b> 4.5×
985	CpG	51.3	55.7	<b>44.8</b>	51.3 0.9–1.1×
986	16-mer	152.6	304.5	302.5	<b>87.0</b> 3.5×
987	SNAP	328.1	450.5	327.5	<b>211.9</b> 1.5–2.1×
988	SGA	453.0	569.3	610.1	<b>409.6</b> 1.4–1.5×

989 Table 2. Seq runtime compared to C++ as compiled with Clang and GCC (seconds). “Seq (Py.)” (Pythonic  
 990 Seq) uses the same code as Python, whereas “Seq (Id.)” (idiomatic Seq) uses Seq-specific language  
 991 features and constructs, with the exception of SNAP and SGA, where the difference is just a single  
 992 prefetch statement.

993

994 the number of reads is low (as in the first benchmark suite), it rapidly deteriorates once the  
 995 read count becomes an order of magnitude larger. Even 100 million reads as used here is  
 996 quite minuscule compared to real datasets.

997 Given the results above, the comparisons below focus only on the C++ implementations.

998

## 999 6.2 Improvements over C++

1000 Table 2 compares the Seq and C++ implementations of each benchmark. Again, all of these  
 1001 implementations (with the exception of idiomatic Seq) are line-by-line identical in terms of  
 1002 the algorithm. The performance of Pythonic Seq code is on par with that of C++ code—in  
 1003 most cases, it is the same as or slightly better than Clang’s (we use Clang as a baseline  
 1004 since both Clang and Seq rely on LLVM for general-purpose optimizations). Note that g++  
 1005 is sometimes able to outperform the LLVM-based backends of Seq and Clang. However,  
 1006 once Seq applies domain-specific optimizations, it outperforms even g++ by up to 7×. For  
 1007 example, in the third set of real-world benchmarks (SNAP and SGA), Seq achieves a 50%  
 1008 speedup after adding a one-line domain-specific prefetch statement to the original code,  
 1009 resulting in up to a 2× speedup over C++.

1010

1011 *Prefetch variability:* Index prefetching is useful during genomic index lookups, and is  
 1012 able to speed up both *k*-mer hash tables and FM-indices by 50%. However, we observed  
 1013 the performance of prefetching to be application- and data-dependent: while in almost  
 1014 all evaluated datasets (spanning various technologies such as recent third-generation 10X  
 1015 Genomics linked-reads [Zheng et al. 2016] and “classic” second-generation Illumina short-  
 1016 reads; detailed results omitted for brevity) it produces a steady improvement in the range  
 1017 20–50%, in one dataset it led to a 35% slowdown. However, the fact that it is a one-line  
 1018 change means that any user can easily experiment and judge whether it works well for their  
 1019 use-case.

1020

## 1021 6.3 Effects of parallelization

1022 To evaluate the performance of Seq’s parallel  
 1023 pipelines, we implemented parallel versions  
 1024 of two of our in-house benchmarks, CpG  
 1025 and 16-mer (the third, RC, performs sub-  
 1026 stantially more I/O and hence does benefit  
 1027 much from parallelism), as well as SGA’s  
 1028 FM-index querying (both with and without  
 1029

Threads	1	2	3	4	Speedup
CpG	58.1	29.6	19.9	15.3	3.8×
16-mer	86.7	43.6	29.9	22.8	3.8×
SGA	355.1	184.0	125.4	95.1	3.7×
SGA (pref.)	217.7	128.0	90.3	71.8	3.0×

Table 3. Seq runtimes on multiple threads (seconds).

1030 prefetch optimizations). To this end, we “block” input reads into batches of 100,000, which  
1031 are processed as a whole by each task; tasks themselves are then executed in parallel via  
1032 Seq’s parallel pipe operator.

1033 Results are shown in Table 3, where Seq scales almost linearly up to 4 threads on our  
1034 in-house benchmarks. For these small applications, we find I/O to be a bottleneck beyond  
1035 4 threads. In a real-world setting where reads would take substantially longer to process,  
1036 we would expect I/O to play a less significant role, allowing a greater degree of parallelism.  
1037 Taking these parallel implementations into account, Seq’s largest speedup over Python  
1038 (which cannot be easily parallelized due to the global interpreter lock) is over  $650\times$ . Finally,  
1039 Table 3 also shows that even Seq’s prefetching optimizations benefit from parallelization.

1040

## 1041 7 RELATED WORK

1042 A few methods have been proposed to aid in the development of bioinformatics tools  
1043 and workflows. One approach focuses on building specialized libraries for genomic data  
1044 manipulation. Examples include C++ libraries such as SeqAn [Döring et al. 2008] and htlib [Li  
1045 et al. 2009b], and high-level libraries such as BioPerl and BioPython [Cock et al. 2009].  
1046 However, none of these libraries solve the aforementioned problems, as none can both scale  
1047 with the size of NGS data and allow for sufficiently high-level representations of bioinformatics  
1048 or genomics algorithms. Additionally, the use of libraries prevents optimizations like operator  
1049 fusion, or the pipeline transformations made by Seq. Another line of work focuses on  
1050 integrating various high-level code blocks into a pipeline framework that can be efficiently  
1051 run on large clusters and cloud-based systems—examples include the Broad Institute’s  
1052 HAIL project and Workflow Description Language [Voss et al. 2017]. While these methods  
1053 indeed allow large-scale parallelism and a relatively high-level description of a given problem,  
1054 they are cumbersome to use as they require rather expensive infrastructural setup and  
1055 administration costs, and do not tackle the problem of single-machine optimizations, which  
1056 is still a significant bottleneck in many pipelines.

1057 The DSL proposed in this paper is inspired by many successful DSLs that already exist  
1058 in other fields of computer science [Abadi et al. 2016; Baghdadi et al. 2019; Chafi et al.  
1059 2011; Chiu et al. 2012; Kjolstad et al. 2017, 2016; Ragan-Kelley et al. 2013; Zhang et al.  
1060 2018]. Despite their substantial success in these other areas, computational biology has  
1061 yet to adopt a comparable DSL. SARVAID [Mahadik et al. 2016] is a DSL designed for  
1062 computational genomics applications, which provides a set of high-level genomics kernels and  
1063 exposes them as language constructs. For example, common operations such as `k-merization`,  
1064 `index-generation`, `index-lookup`, `similarity-computation` and `clustering` are provided.  
1065 While such an approach provides efficient implementations of these kernels and combinations  
1066 thereof, it lacks generality, which is gravely needed in the field as new sequencing technologies  
1067 produce new types of data that in turn necessitate novel algorithms. Seq aims to provide  
1068 a more general, lower-level language, with general-purpose constructs that can be used to  
1069 build a variety of kernels efficiently. While there also exist a few general-purpose languages  
1070 optimized for scientific computing such as Julia [Bezanson et al. 2012] and MATLAB, neither  
1071 of these languages is designed for computational biology workflows.

1072 On the Python side, recent Python standards introduced type hints, which allow static  
1073 type checking [van Rossum 2015]. Projects such as Cython [Behnel et al. 2011], PyPy [Bolz  
1074 et al. 2009], Numba [Lam et al. 2015], Shed Skin [Dufour 2006] and Nuitka [Hayen 2012] all  
1075 aim to generate efficient code by relying on ideas such as static type checking and (JIT-)  
1076 compiling rather than interpreting Python. While Seq is similar to these language in that it  
1077 is indeed compiled and uses static type checking, Seq also uses domain-specific information

1078

	Domain	Target	Compilation	Unknown types allowed?	Full Python?	CPython runtime?	Multi-threading?
1079	<b>CPython</b>	General	Bytecode	Interpreted	✓	✓	✓
1080	<b>Seq</b>	Bio.	LLVM IR	AOT	✗	✗	✗
1081	<b>Cython</b>	General	C	AOT	✗	✓	✓
1082	<b>PyPy</b>	General	Bytecode	Interpreted	✓	✓	✓
1083	<b>Numba</b>	Sci.	LLVM IR	JIT	✓	✗	✓
1084	<b>Nuitka</b>	General	C++	AOT	✓	✓	✓
1085	<b>Pythran</b>	Sci.	C++	AOT	✓	✓	✗
1086	<b>Pyston</b>	General	LLVM IR	JIT	✓	✓	✓
1087	<b>HOPE</b>	Astro.	C++	JIT	✗	✗	✗
1088	<b>Shed Skin</b>	General	C++	AOT	✗	✗	✗
1089	<b>Grumpy</b>	General	Go	AOT	✗	✗	✗
1090							
1091							
1092							
1093							

1094 Table 4. Comparison between Seq and other Python implementations. For “Domain”, “Bio.” means  
 1095 computational biology, “Sci.” means scientific computing and “Astro.” means astrophysical computing.  
 1096 “Unknown types” refer to types that cannot be statically determined. Also note that Pyston has several  
 1097 JIT tiers in addition to its LLVM JIT.

1098  
 1099

1100 to apply further code optimizations, and introduces data types that are tailored to the field  
 1101 of computational biology. Furthermore, many of these other implementations still rely on the  
 1102 Python runtime, and are thus bound to its inherent performance overhead. For the sake of  
 1103 completeness, a comprehensive comparison between Seq and other Python implementations  
 1104 is given in Table 4. In this work, we chose to compare to Nuitka and Shed Skin primarily  
 1105 because other similar implementations (e.g. Numba, Pythran, Pyston, Grumpy) are either  
 1106 geared more towards scientific/numerical computing or no longer under active development.

1107  
 1108

## 1108 8 CONCLUSION

1109 We have introduced Seq, a new language for computational biology that offers the productivity  
 1110 of Python and the performance of C. Thereby, Seq bridges the gap between computationalists  
 1111 who seek to write performance-critical code for a particular application, and biologists  
 1112 whose day-to-day workflow involves rapid development and prototyping of new ideas and  
 1113 algorithms. Through Seq, we introduce several novel genomics-specific language constructs  
 1114 and optimizations, which collectively attain a 175× performance improvement over standard  
 1115 Python, and a 2–7× improvement over C++. Future work includes exploring a wider range  
 1116 of domain-specific optimizations that exploit the unique structure of biological data.

1117  
 1118  
 1119  
 1120  
 1121  
 1122  
 1123  
 1124  
 1125  
 1126  
 1127

## REFERENCES

- 1128  
1129 1000 Genomes Project Consortium, Gonçalo R. Abecasis, David Altshuler, Adam Auton, Lisa D. Brooks,  
1130 Richard M. Durbin, Richard A. Gibbs, Matt E. Hurles, and Gil A. McVean. 2010. A Map of Human  
1131 Genome Variation from Population-Scale Sequencing. *Nature* 467, 7319 (Oct. 2010), 1061–1073. <https://doi.org/10.1038/nature09534>  
1132  
1133 Martin Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean, Matthieu Devin,  
1134 Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur, Josh Levenberg, Rajat Monga,  
1135 Sherry Moore, Derek G. Murray, Benoit Steiner, Paul Tucker, Vijay Vasudevan, Pete Warden, Martin  
1136 Wicke, Yuan Yu, and Xiaoqiang Zheng. 2016. TensorFlow: A system for large-scale machine learning.  
1137 In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. 265–283.  
1138 <https://www.usenix.org/system/files/conference/osdi16/osdi16-abadi.pdf>  
1139 R. Appuswamy, J. Fellay, and N. Chaturvedi. 2018. Sequence Alignment Through the Looking Glass. In  
1140 *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. 257–266.  
1141 <https://doi.org/10.1109/IPDPSW.2018.00050>  
1142 Stefan Milton Bache and Hadley Wickham. 2014. magrittr: A forward-pipe operator for R. *R package*  
1143 *version 1, 1* (2014). <https://cran.r-project.org/web/packages/magrittr/vignettes/magrittr.html>  
1144 Riyadh Baghdadi, Jessica Ray, Malek Ben Romdhane, Emanuele Del Sozzo, Abdurrahman Akkas, Yunming  
1145 Zhang, Patricia Suriana, Shoaib Kamil, and Saman Amarasinghe. 2019. Tiramisu: A Polyhedral Compiler  
1146 for Expressing Fast and Portable Code. In *Proceedings of the 2019 IEEE/ACM International Symposium*  
1147 *on Code Generation and Optimization (CGO 2019)*. IEEE Press, Piscataway, NJ, USA, 193–205. <http://dl.acm.org/citation.cfm?id=3314872.3314896>  
1148 J. A. Bailey, A. M. Yavor, H. F. Massa, B. J. Trask, and E. E. Eichler. 2001. Segmental duplications:  
1149 organization and impact within the current human genome project assembly. *Genome Research* 11, 6  
1150 (2001), 1005–1017. <https://doi.org/10.1101/gr.187101>  
1151 Monya Baker. 2016. 1,500 scientists lift the lid on reproducibility. *Nature News* 533, 7604 (2016), 452.  
1152 S. Behnel, R. Bradshaw, C. Citro, L. Dalcin, D.S. Seljebotn, and K. Smith. 2011. Cython: The Best of Both  
1153 Worlds. *Computing in Science Engineering* 13, 2 (2011), 31–39. <https://doi.org/10.1109/MCSE.2010.118>  
1154 Bonnie Berger, Noah M. Daniels, and Y. William Yu. 2016. Computational Biology in the 21st Century:  
1155 Scaling with Compressive Algorithms. *Commun. ACM* 59, 8 (July 2016), 72–80. <https://doi.org/10.1145/2957324>  
1156 Jeff Bezanson, Stefan Karpinski, Viral B Shah, and Alan Edelman. 2012. Julia: A fast dynamic language for  
1157 technical computing. *arXiv* (2012), 1209.5145.  
1158 Hans-Juergen Boehm and Mark Weiser. 1988. Garbage Collection in an Uncooperative Environment. *Softw.*  
1159 *Pract. Exper.* 18, 9 (Sept. 1988), 807–820. <https://doi.org/10.1002/spe.4380180902>  
1160 Carl Friedrich Bolz, Antonio Cuni, Maciej Fijalkowski, and Armin Rigo. 2009. Tracing the Meta-level:  
1161 PyPy’s Tracing JIT Compiler. In *Proceedings of the 4th Workshop on the Implementation, Compilation,*  
1162 *Optimization of Object-Oriented Languages and Programming Systems (ICOOOLPS ’09)*. ACM, New  
1163 York, NY, USA, 18–25. <https://doi.org/10.1145/1565824.1565827>  
1164 Hassan Chafi, Arvind K. Sujeeth, Kevin J. Brown, HyoukJoong Lee, Anand R. Atreya, and Kunle Olukotun.  
1165 2011. A Domain-specific Approach to Heterogeneous Parallelism. *SIGPLAN Not.* 46, 8 (Feb. 2011), 35–46.  
1166 <https://doi.org/10.1145/2038037.1941561>  
1167 Shimin Chen, Anastassia Ailamaki, Phillip B. Gibbons, and Todd C. Mowry. 2007. Improving Hash  
1168 Join Performance Through Prefetching. *ACM Trans. Database Syst.* 32, 3, Article 17 (Aug. 2007).  
1169 <https://doi.org/10.1145/1272743.1272747>  
1170 Charisee Chiw, Gordon Kindlmann, John Reppy, Lamont Samuels, and Nick Seltzer. 2012. Diderot: a  
1171 parallel DSL for image analysis and visualization. In *Acm sigplan notices*, Vol. 47. ACM, 111–120.  
1172 Peter JA Cock, Tiago Antao, Jeffrey T Chang, Brad A Chapman, Cymon J Cox, Andrew Dalke, Iddo  
1173 Friedberg, Thomas Hamelryck, Frank Kauff, Bartek Wilczynski, et al. 2009. Biopython: freely available  
1174 Python tools for computational molecular biology and bioinformatics. *Bioinformatics* 25, 11 (2009),  
1175 1422–1423.  
1176 Andreas Döring, David Weese, Tobias Rausch, and Knut Reinert. 2008. SeqAn an efficient, generic C++  
library for sequence analysis. *BMC Bioinformatics* 9, 1 (2008), 11.  
Mark Dufour. 2006. *Shed skin: An optimizing python-to-c++ compiler*. Master’s thesis. Delft University of  
Technology.  
Michael Farrar. 2006. Striped Smith–Waterman speeds database searches six times over other SIMD  
implementations. *Bioinformatics* 23, 2 (11 2006), 156–161. <https://doi.org/10.1093/bioinformatics/btl582>  
arXiv:<http://oup.prod.sis.lan/bioinformatics/article-pdf/23/2/156/536640/btl582.pdf>

- 1177 Paolo Ferragina and Giovanni Manzini. 2004. Compression Boosting in Optimal Linear Time Using the  
1178 Burrows-Wheeler Transform. In *SODA 2004*. 655–663.
- 1179 Isaac Gouy. [n. d.]. The Computer Language Benchmarks Game. <https://benchmarksgame-team.pages.debian.net/benchmarksgame/>. ([n. d.]).
- 1180 Erika Check Hayden. 2014. Technology: The \$1,000 genome. *Nature* 507, 7492 (mar 2014), 294–295.  
1181 <https://doi.org/10.1038/507294a>
- 1182 K Hayen. 2012. Nuitka. (2012). <http://nuitka.net>
- 1183 Rick Kamps, Rita D Brandão, Bianca J Bosch, Aimee DC Paulussen, Sofia Xanthoulea, Marinus J Blok,  
1184 and Andrea Romano. 2017. Next-generation sequencing in oncology: genetic diagnosis, risk prediction  
1185 and cancer classification. *International Journal of Molecular Sciences* 18, 2 (2017), 308.
- 1186 Abdul Rafay Khan, Muhammad Tariq Pervez, Masroor Ellahi Babar, Nasir Naveed, and Muhammad  
1187 Shoaib. 2018. A Comprehensive Study of De Novo Genome Assemblers: Current Challenges and Future  
1188 Prospective. *Evol Bioinform Online* 14 (20 Feb 2018), 1176934318758650–1176934318758650. <https://doi.org/10.1177/1176934318758650> 29511353[pmid].
- 1189 Vladimir Kiriansky, Haoran Xu, Martin Rinard, and Saman Amarasinghe. 2018. Cimple: Instruction and  
1190 Memory Level Parallelism: A DSL for Uncovering ILP and MLP. In *Proceedings of the 27th International  
1191 Conference on Parallel Architectures and Compilation Techniques (PACT '18)*. ACM, New York, NY,  
1192 USA, Article 30, 16 pages. <https://doi.org/10.1145/3243176.3243185>
- 1193 Fredrik Kjolstad, Stephen Chou, David Lugato, Shoaib Kamil, and Saman Amarasinghe. 2017. Taco: A  
1194 tool to generate tensor algebra kernels. In *Proc. IEEE/ACM Automated Software Engineering*. IEEE,  
1195 943–948.
- 1196 Fredrik Kjolstad, Shoaib Kamil, Jonathan Ragan-Kelley, David IW Levin, Shinjiro Sueda, Desai Chen,  
1197 Etienne Vouga, Danny M Kaufman, Gurtej Kanwar, Wojciech Matusik, et al. 2016. Simit: A language for  
1198 physical simulation. *ACM Transactions on Graphics (TOG)* 35, 2 (2016), 20.
- 1199 Gregory Kucherov, Karel Brinda, and Maciej Sykulski. 2015. Spaced seeds improve k-mer-based metagenomic  
1200 classification. *Bioinformatics* 31, 22 (07 2015), 3584–3592. <https://doi.org/10.1093/bioinformatics/btv419>  
1201 arXiv:<http://oup.prod.sis.lan/bioinformatics/article-pdf/31/22/3584/5027960/btv419.pdf>
- 1202 Siu Kwan Lam, Antoine Pitrou, and Stanley Seibert. 2015. Numba: A LLVM-based Python JIT Compiler. In  
1203 *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC (LLVM '15)*. ACM,  
1204 New York, NY, USA, Article 7, 6 pages. <https://doi.org/10.1145/2833157.2833162>
- 1205 Chris Lattner and Vikram Adve. 2004. LLVM: A compilation framework for lifelong program analysis &  
1206 transformation. In *Proc. Int. Sym. on Code Generation and Optimization*. IEEE Computer Society, 75.
- 1207 Heng Li and Richard Durbin. 2009. Fast and Accurate Short Read Alignment with Burrows-Wheeler  
1208 Transform. *Bioinformatics* 25, 14 (2009), 1754–1760.
- 1209 Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis,  
1210 and Richard Durbin. 2009a. The sequence alignment/map format and SAMtools. *Bioinformatics* 25, 16  
1211 (2009), 2078–2079.
- 1212 Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo  
1213 Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. 2009b. The Sequence  
1214 Alignment/Map Format and SAMtools. *Bioinformatics* 25, 16 (2009), 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>
- 1215 Heng Li and Nils Homer. 2010. A survey of sequence alignment algorithms for next-generation sequencing.  
1216 *Brief Bioinform* 11, 5 (Sep 2010), 473–483. <https://doi.org/10.1093/bib/bbq015> 20460430[pmid].
- 1217 Hengyun Lu, Francesca Giordano, and Zemin Ning. 2016. Oxford Nanopore MinION sequencing and genome  
1218 assembly. *Genomics, Proteomics & Bioinformatics* 14, 5 (2016), 265–279.
- 1219 Kanak Mahadik, Christopher Wright, Jinyi Zhang, Milind Kulkarni, Saurabh Bagchi, and Somali Chaterji.  
1220 2016. SARVAVID: A Domain Specific Language for Developing Scalable Computational Genomics  
1221 Applications. In *Proceedings of the 2016 International Conference on Supercomputing (ICS '16)*. ACM,  
1222 New York, NY, USA, Article 34, 12 pages. <https://doi.org/10.1145/2925426.2926283>
- 1223 Teri A Manolio, Lisa D Brooks, and Francis S Collins. 2008. A HapMap harvest of insights into the genetics  
1224 of common disease. *The Journal of Clinical Investigation* 118, 5 (2008), 1590–1605.
- 1225 ER Mardis. 2017. DNA sequencing technologies: 2006–2016. *Nature Protocols* 12, 2 (2017), 213–218.
- Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernytksy,  
Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, and Mark A. DePristo. 2010. The Genome  
Analysis Toolkit: A MapReduce Framework for Analyzing next-Generation DNA Sequencing Data.  
*Genome Research* 20, 9 (Sept. 2010), 1297–1303. <https://doi.org/10.1101/gr.107524.110>



- 1226 Paul Muir, Shantao Li, Shaoke Lou, Daifeng Wang, Daniel J Spakowicz, Leonidas Salichos, Jing Zhang,  
1227 George M Weinstock, Farren Isaacs, Joel Rozowsky, et al. 2016. The real cost of sequencing: scaling  
1228 computation to keep pace with data generation. *Genome Biology* 17, 1 (2016), 53.
- 1229 Gor Nishanov. 2017. ISO/IEC TS 22277:2017. (Dec 2017). <https://www.iso.org/standard/73008.html>
- 1230 Brian D. Ondov, Todd J. Treangen, Páll Melsted, Adam B. Mallonee, Nicholas H. Bergman, Sergey Koren,  
1231 and Adam M. Phillippy. 2016. Mash: fast genome and metagenome distance estimation using MinHash.  
1232 *Genome Biology* 17, 1 (20 Jun 2016), 132. <https://doi.org/10.1186/s13059-016-0997-x>
- 1233 Roger D Peng. 2011. Reproducible research in computational science. *Science* 334, 6060 (2011), 1226–1227.
- 1234 Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Ama-  
1235 rasinghe. 2013. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in  
1236 image processing pipelines. *ACM SIGPLAN Notices* 48, 6 (2013), 519–530.
- 1237 Andrea Sboner, Ximeng Jasmine Mu, Dov Greenbaum, Raymond K Auerbach, and Mark B Gerstein. 2011.  
1238 The real cost of sequencing: higher than you think! *Genome biology* 12, 8 (2011), 125.
- 1239 Tao B. Schardl, William S. Moses, and Charles E. Leiserson. 2017. Tapir: Embedding Fork-Join Parallelism  
1240 into LLVM’s Intermediate Representation. *SIGPLAN Not.* 52, 8 (Jan. 2017), 249–265. <https://doi.org/10.1145/3155284.3018758>
- 1241 Ariya Shajii, Ibrahim Numanagić, Christopher Whelan, and Bonnie Berger. 2018. Statistical Binning for  
1242 Barcoded Reads Improves Downstream Analyses. *Cell Systems* 7, 2 (2018), 219–226.
- 1243 Jared T. Simpson and Richard Durbin. 2012. Efficient de novo assembly of large genomes using com-  
1244 pressed data structures. *Genome Res* 22, 3 (Mar 2012), 549–556. [https://doi.org/10.1101/gr.126953.11122156294\[pmid\]](https://doi.org/10.1101/gr.126953.11122156294[pmid]).
- 1245 Petr Šmarda, Petr Bureš, Lucie Horová, Ilia J. Leitch, Ladislav Mucina, Ettore Pacini, Lubomír Tichý, Vít  
1246 Grulich, and Olga Rotreklová. 2014. Ecological and evolutionary significance of genomic GC content  
1247 diversity in monocots. *Proceedings of the National Academy of Sciences* 111, 39 (2014), E4096–E4102.  
1248 <https://doi.org/10.1073/pnas.1321152111> arXiv:<https://www.pnas.org/content/111/39/E4096.full.pdf>
- 1249 Hajime Suzuki and Masahiro Kasahara. 2018. Introducing difference recurrence relations for faster semi-global  
1250 alignment of long sequences. *BMC Bioinformatics* 19, 1 (19 Feb 2018), 45. <https://doi.org/10.1186/s12859-018-2014-8>
- 1251 Guido van Rossum. 2015. The Python Library Reference, Release 3.5. Fred L. Drake Jr.
- 1252 K Voss, J Gentry, and G Van der Auwera. 2017. Full-stack genomics pipelining with GATK4 +WDL  
1253 +Cromwell.. In *18th Annual Bioinformatics Open Source Conference*. poster.
- 1254 Martin Šošić and Mile Šikić. 2017. Edlib: a C/C++ library for fast, exact sequence alignment using edit  
1255 distance. *Bioinformatics* 33, 9 (01 2017), 1394–1395. <https://doi.org/10.1093/bioinformatics/btw753>  
1256 arXiv:<http://oup.prod.sis.lan/bioinformatics/article-pdf/33/9/1394/25151249/btw753.pdf>
- 1257 Wendi Wang, Wen Tang, Linchuan Li, Guangming Tan, Peiheng Zhang, and Ninghui Sun. 2012. Investigating  
1258 Memory Optimization of Hash-index for Next Generation Sequencing on Multi-core Architecture. *2012  
1259 IEEE 26th International Parallel and Distributed Processing Symposium Workshops & PhD Forum*  
1260 (2012), 665–674.
- 1261 Deirdre Weymann, Janessa Laskin, Robyn Roscoe, Kasmintan A. Schrader, Stephen Chia, Stephen Yip,  
1262 Winson Y. Cheung, Karen A. Gelmon, Aly Karsan, Daniel J. Renouf, Marco Marra, and Dean A. Regier.  
1263 2017. The cost and cost trajectory of whole-genome analysis guiding treatment of patients with advanced  
1264 cancers. *Molecular Genetics & Genomic Medicine* 5, 3 (2017), 251–260. <https://doi.org/10.1002/mgg3.281>  
1265 arXiv:<https://onlinelibrary.wiley.com/doi/pdf/10.1002/mgg3.281>
- 1266 Deniz Yorukoglu, Yun William Yu, Jian Peng, and Bonnie Berger. 2016. Compressive mapping for next-  
1267 generation sequencing. *Nat Biotech* 34, 4 (2016), 374–376. <http://dx.doi.org/10.1038/nbt.3511> Opinion  
1268 and Comment.
- 1269 Matei Zaharia, William J. Bolosky, Kristal Curtis, Armando Fox, David A. Patterson, Scott Shenker, Ion  
1270 Stoica, Richard M. Karp, and Taylor Sittler. 2011. Faster and More Accurate Sequence Alignment with  
1271 SNAP. *CoRR* abs/1111.5572 (2011). arXiv:[1111.5572](http://arxiv.org/abs/1111.5572) <http://arxiv.org/abs/1111.5572>
- 1272 Di Zhang, Yunquan Zhang, and Jing Chen. 2007. Efficient Construction of FM-index Using Overlapping  
1273 Block Processing for Large Scale Texts. In *Proceedings of the 29th European Conference on IR Research  
1274 (ECIR’07)*. Springer-Verlag, Berlin, Heidelberg, 113–123. <http://dl.acm.org/citation.cfm?id=1763653.1763669>
- 1275 Yunming Zhang, Mengjiao Yang, Riyadh Baghdadi, Shoaib Kamil, Julian Shun, and Saman Amarasinghe.  
1276 2018. GraphIt: A High-performance Graph DSL. *Proc. ACM Program. Lang.* 2, OOPSLA, Article 121  
1277 (Oct. 2018), 30 pages. <https://doi.org/10.1145/3276491>



1275 Grace XY Zheng, Billy T. Lau, Michael Schnall-Levin, Mirna Jarosz, John M. Bell, Christopher M. Hindson,  
1276 Sofia Kyriazopoulou-Panagiotopoulou, Donald A. Masquelier, Landon Merrill, Jessica M. Terry, Patrice A.  
1277 Mudivarti, Paul W. Wyatt, Rajiv Bharadwaj, Anthony J. Makarewicz, Yuan Li, Phillip Belgrader,  
1278 Andrew D. Price, Adam J. Lowe, Patrick Marks, Gerard M. Vurens, Paul Hardenbol, Luz Montesclaros,  
1279 Melissa Luo, Lawrence Greenfield, Alexander Wong, David E. Birch, Steven W. Short, Keith P. Bjornson,  
1280 Pranav Patel, Erik S. Hopmans, Christina Wood, Sukhvinder Kaur, Glenn K. Lockwood, David Stafford,  
1281 Joshua P. Delaney, Indira Wu, Heather S. Ordonez, Susan M. Grimes, Stephanie Greer, Josephine Y. Lee,  
1282 Kamila Belhocine, Kristina M. Giorda, William H. Heaton, Geoffrey P. McDermott, Zachary W. Bent,  
1283 Francesca Meschi, Nikola O. Kondov, Ryan Wilson, Jorge A. Bernate, Shawn Gauby, Alex Kindwall, Clara  
1284 Bermejo, Adrian N. Fehr, Adrian Chan, Serge Saxonov, Kevin D. Ness, Benjamin J. Hindson, and Hanlee P.  
1285 Ji. 2016. Haplotyping germline and cancer genomes using high-throughput linked-read sequencing. *Nat*  
1286 *Biotechnol* 34, 3 (01 Mar 2016), 303–311. <https://doi.org/10.1038/nbt.3432> 26829319[pmid].  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295  
1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323

## 1324 A SEQ VS. PYTHON – A CHEAT SHEET

### 1325 Additional types

- 1326 • seq: Represents a genomic sequence.
- 1327 •  $kN$  ( $1 \leq N \leq 128$ ): Represents a  $k$ -mer of length  $N$ .
- 1328 •  $iN$  ( $1 \leq N \leq 256$ ): Represents a signed  $N$ -bit integer (standard int is an `i64`).
- 1329 •  $uN$  ( $1 \leq N \leq 256$ ): Represents an unsigned  $N$ -bit integer.
- 1330 • `ptr[T]`: Represents a pointer to type  $T$ ; primarily useful for C interoperability.
- 1331 • `array[T]`: Represents an array of type  $T$  (essentially a pointer with a length).

### 1333 Additional keywords

- 1334 • `type`: Declares a named tuple or a type alias.
- 1335 • `match/case/default`: Match statement
- 1336 • `extend`: Adds the given methods to an existing type.
- 1337 • `cdef`: Declares an externally-defined C function.
- 1338 • `prefetch`: Prefetches the given items in the given index objects.

1340 Seq also provides `__ptr__` for obtaining a pointer to a variable, and `__array__` for declaring  
1341 stack-allocated fixed-size arrays.

### 1342 Static types

1344 Because Seq is statically-typed, lists (for example) cannot contain elements of different types  
1345 as they can in Python. Similarly, a variable cannot be assigned to objects of different types,  
1346 nor can a function return objects of different types. Seq currently also does not support  
1347 polymorphism.

### 1348 Tuples

1350 Tuples in Seq are implemented as structures. Consequently, heterogeneous tuples can only  
1351 be indexed by a constant value, since otherwise the type of the index expression would be  
1352 ambiguous. Similarly, iteration over a tuple is only possible if it is homogeneous.

### 1353 Scopes

1355 Seq enforces slightly stricter variable scoping rules than standard Python. In short, a variable  
1356 cannot be first assigned in a block then used afterwards for the first time in the enclosing  
1357 block. This restriction avoids the problem of uninitialized variables.

## 1359 B CODE FROM SNAP BENCHMARK

1360 SNAP uses a hierarchical hash table as its genomic index. To index  $k$ -mers ( $k \geq 16$ ), an array  
1361 of  $4^{k-16}$  quadratic probing hash tables is created, indexed by every possible length- $(k-16)$   
1362 prefix. Each constituent hash table is then a mapping of 16-mer to genomic loci at which  
1363 that 16-mer appears. To handle multiple loci for a single  $k$ -mer, an auxiliary array is used,  
1364 and hash table values can be pointers into this array (determined by whether the value is  
1365 greater than the largest locus). This hierarchical structure exploits the fact that not every  
1366 length- $(k-16)$  prefix appears in the genome with equal frequency, so the size of each internal  
1367 hash table can be chosen based on the corresponding prefix's frequency.

1368 SNAP's index is implemented in C++, and can be found at [https://github.com/amplab/](https://github.com/amplab/snap)  
1369 `snap`. Below, we give the Seq implementation as well as the Seq and C++ code for querying,  
1370 which was used in the SNAP benchmark. Note that loading the precomputed table from  
1371 disk is still done in C++, and wrapped in Seq.

1372

```

1373 # File: hashtable.seq
1374 # Implementation of SNAP aligner's hash table
1375 # https://github.com/amplab/snap/blob/master/SNAPLib/HashTable.{cpp,h}
1376
1377 # Need the following hooks linked to convert C++ SNAPHashTable to Seq object:
1378 # snap_hashtable_ptr(ptr[byte]) -> ptr[tuple[K,V]] -- extract table pointer
1379 # snap_hashtable_len(ptr[byte]) -> int -- extract table length
1380 # snap_hashtable_invalid_val(ptr[byte]) -> V -- extract "invalid" value
1381
1382 QUADRATIC_CHAINING_DEPTH = 5
1383
1384 class SNAPHashTable[K,V]:
1385     table: array[tuple[V,K]]
1386     invalid_val: V
1387
1388     def _hash(k):
1389         key = hash(k)
1390         key ^= int(u64(key) >> u64(33))
1391         key *= 0xff51afd7ed558ccd
1392         key ^= int(u64(key) >> u64(33))
1393         key *= 0xc4ceb9fe1a85ec53
1394         key ^= int(u64(key) >> u64(33))
1395         return key
1396
1397     def __init__(self: SNAPHashTable[K,V], size: int, invalid_val: V):
1398         self.table = array[tuple[V,K]](size)
1399         self.invalid_val = invalid_val
1400
1401         for i in range(size):
1402             self.table[i] = (invalid_val, K())
1403
1404     def __init__(self: SNAPHashTable[K,V], p: ptr[byte]):
1405         cdef snap_hashtable_ptr(ptr[byte]) -> ptr[tuple[V,K]]
1406         cdef snap_hashtable_len(ptr[byte]) -> int
1407         cdef snap_hashtable_invalid_val(ptr[byte]) -> V
1408         self.table = array[tuple[V,K]](snap_hashtable_ptr(p), snap_hashtable_len(p))
1409         self.invalid_val = snap_hashtable_invalid_val(p)
1410
1411     def _get_index(self: SNAPHashTable[K,V], where: int):
1412         return int(u64(where) % u64(len(self.table)))
1413
1414     def get_value_ptr_for_key(self: SNAPHashTable[K,V], k: K):
1415         table = self.table
1416         table_size = table.len
1417         table_index = self._get_index(SNAPHashTable[K,V]._hash(k))
1418         invalid_val = self.invalid_val
1419         entry = table[table_index]
1420
1421         if entry[1] == k and entry[0] != invalid_val:
1422             return ptr[V](table.ptr + table_index)
1423         else:
1424             n_probes = 0
1425             while True:
1426                 n_probes += 1
1427
1428                 if n_probes > table_size + QUADRATIC_CHAINING_DEPTH:
1429                     return ptr[V]()
1430
1431                 diff = (n_probes**2) if n_probes < QUADRATIC_CHAINING_DEPTH else 1
1432                 table_index = (table_index + diff) % table_size
1433
1434                 entry = table[table_index]
1435                 if not (entry[1] != k and entry[0] != invalid_val):
1436                     break
1437
1438             return ptr[V](table.ptr + table_index)
1439
1440     def __prefetch__(self: SNAPHashTable[K,V], k: K):
1441         table = self.table
1442         table_index = self._get_index(SNAPHashTable[K,V]._hash(k))
1443         (self.table.ptr + table_index).__prefetch_r3__()
1444
1445     def __getitem__(self: SNAPHashTable[K,V], k: K):
1446         p = self.get_value_ptr_for_key(k)
1447         return p[0] if p else self.invalid_val
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471

```

```

1422 # File: genomeindex.seq
1423 # Implementation of SNAP aligner's genome index
1424 # https://github.com/amplab/snap/blob/master/SNAPLib/GenomeIndex.{cpp,h}
1425
1426 # Need the following hooks linked to convert C++ GenomeIndex to Seq object:
1427 # snap_index_from_dir(ptr[byte]) -> ptr[byte] -- read object from specified directory
1428 # snap_index_ht_count(ptr[byte]) -> int -- extract hash table count
1429 # snap_index_ht_get(ptr[byte], int) -> ptr[byte] -- extract specified (0-indexed) hash table
1430 # snap_index_overflow_ptr(ptr[byte]) -> ptr[u32] -- extract overflow table pointer
1431 # snap_index_overflow_len(ptr[byte]) -> int -- extract overflow table length
1432 # snap_index_count_of_bases(ptr[byte]) -> int -- extract count of genome bases
1433
1434 from hashtable import SNAPHashTable
1435
1436 class GenomeIndex[K]:
1437     hash_tables: array[SNAPHashTable[k16,u32]]
1438     overflow_table: array[u32]
1439     count_of_bases: int
1440
1441     def _partition(k: K):
1442         n = int(k.as_int())
1443         return (k16(n & ((1 << 32) - 1)), n >> 32)
1444
1445     def __init__(self: GenomeIndex[K], dir: str):
1446         assert k16.len() <= K.len() <= k32.len()
1447         cdef snap_index_from_dir(ptr[byte]) -> ptr[byte]
1448         cdef snap_index_ht_count(ptr[byte]) -> int
1449         cdef snap_index_ht_get(ptr[byte], int) -> ptr[byte]
1450         cdef snap_index_overflow_ptr(ptr[byte]) -> ptr[u32]
1451         cdef snap_index_overflow_len(ptr[byte]) -> int
1452         cdef snap_index_count_of_bases(ptr[byte]) -> int
1453
1454         p = snap_index_from_dir(dir.c_str())
1455         assert p
1456         hash_tables = array[SNAPHashTable[k16,u32]](snap_index_ht_count(p))
1457         for i in range(len(hash_tables)):
1458             hash_tables[i] = SNAPHashTable[k16,u32](snap_index_ht_get(p, i))
1459
1460         self.hash_tables = hash_tables
1461         self.overflow_table = array[u32](snap_index_overflow_ptr(p), snap_index_overflow_len(p))
1462         self.count_of_bases = snap_index_count_of_bases(p)
1463
1464     def __getitem__(self: GenomeIndex[K], seed: K):
1465         kmer, which = GenomeIndex[K]._partition(seed)
1466         table = self.hash_tables[which]
1467         value_ptr = table.get_value_ptr_for_key(kmer)
1468
1469         if not value_ptr or value_ptr[0] == table.invalid_val:
1470             return array[u32](value_ptr, 0)
1471
1472         value = value_ptr[0]
1473
1474         if int(value) < self.count_of_bases:
1475             return array[u32](value_ptr, 1)
1476         else:
1477             overflow_table_offset = int(value) - self.count_of_bases
1478             hit_count = int(self.overflow_table[overflow_table_offset])
1479             return array[u32](self.overflow_table.ptr + overflow_table_offset + 1, hit_count)
1480
1481     def __prefetch__(self: GenomeIndex[K], seed: K):
1482         kmer, which = GenomeIndex[K]._partition(seed)
1483         table = self.hash_tables[which]
1484         table.__prefetch__(kmer)

```

```
1471 # File: snap.seq
1472 # k-mer counting using SNAP's hash table
1473 from genomeindex import *
1474 from sys import argv
1475 type K = k20
1476 good = 0
1477 bad = 0
1478
1479 def process(kmer: K, index: GenomeIndex[K]):
1480     global good, bad
1481     prefetch index[kmer], index[-kmer]
1482     hits = index[kmer]
1483     hits_rc = index[-kmer]
1484
1485     if len(hits) > 0 or len(hits_rc) > 0:
1486         good += 1
1487     else:
1488         bad += 1
1489
1490 assert len(argv) == 3
1491 index = GenomeIndex[K](argv[1])
1492 step = 10
1493 fastq(argv[2]) |> kmers[K](step) |> process(index)
1494 print good, bad
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
```

```

1520 // File: snap.cpp
1521 // k-mer counting using SNAP's hash table
1522 #include <iostream>
1523 #include <fstream>
1524 #include <cstdlib>
1525 #include <stdint>
1526 #include <cassert>
1527 using namespace std;
1528
1529 extern "C" void *snap_index_from_dir(char *dir);
1530 extern "C" void snap_index_lookup(void *idx, char *bases, int64_t *nHits, const unsigned **hits,
1531                                 int64_t *nRCHits, const unsigned **rchHits);
1532
1533 // filter ambiguous bases from sequences
1534 static bool hasN(char *kmer, unsigned len)
1535 {
1536     for (unsigned i = 0; i < len; i++) {
1537         if (kmer[i] == 'N')
1538             return true;
1539     }
1540     return false;
1541 }
1542
1543 int good = 0;
1544 int bad = 0;
1545
1546 int main(int argc, char *argv[])
1547 {
1548     assert(argc == 3);
1549     void *idx = snap_index_from_dir(argv[1]);
1550     const unsigned k = 20;
1551     const unsigned step = 10;
1552     unsigned hit = 0;
1553     unsigned hit_rc = 0;
1554     int64_t n_hits = 0;
1555     int64_t n_hits_rc = 0;
1556     const unsigned *hit_p = &hit;
1557     const unsigned *hit_rc_p = &hit_rc;
1558     ifstream fin(argv[2]);
1559     string read;
1560     long line = -1;
1561
1562     while (getline(fin, read)) {
1563         line++;
1564         if (line % 4 != 1) continue; // skip non-sequences in FASTQ
1565         unsigned max_pos = 0, max_count = 0;
1566         char *buf = (char *)read.c_str();
1567         unsigned len = read.size();
1568
1569         for (unsigned i = 0; i + k <= len; i += step) {
1570             if (hasN(&buf[i], k)) continue;
1571             snap_index_lookup(idx, &buf[i], &n_hits, &hit_p, &n_hits_rc, &hit_rc_p);
1572             ++((n_hits > 0 || n_hits_rc > 0) ? good : bad);
1573         }
1574     }
1575
1576     cout << good << " " << bad << endl;
1577 }

```

1555  
1556  
1557  
1558  
1559  
1560  
1561  
1562  
1563  
1564  
1565  
1566  
1567  
1568